

# Probabilistic Programming and Bayesian Methods for Hackers

Cameron DAVIDSON-PILON

June 8, 2013



# CONTENTS

## 0.1 Probabilistic Programming

## 0.2 and Bayesian Methods for Hackers

**Version 0.1** Welcome to *Bayesian Methods for Hackers*. The full Github repository, and additional chapters, is available at [github/Probabilistic-Programming-and-Bayesian-Methods-for-Hackers](https://github.com/Probabilistic-Programming-and-Bayesian-Methods-for-Hackers). We hope you enjoy the book, and we encourage any contributions!

## 0.3 Chapter 1

---

### 0.3.1 The Philosophy of Bayesian Inference

You are a skilled programmer, but bugs still slip into your code. After a particularly difficult implementation of an algorithm, you decide to test your code on a trivial example. It passes. You test the code on a harder problem. It passes once again. And it passes the next, *even more difficult*, test too! You are starting to believe that there may be no bugs in this code...

If you think this way, then congratulations, you already are a Bayesian practitioner! Bayesian inference is simply updating your beliefs after considering new evidence. A Bayesian can rarely be certain about a result, but he or she can be very confident. Just like in the example above, we can never be 100% sure that our code is bug-free unless we test it on every possible problem; something rarely possible in practice. Instead, we can test it on a large number of problems, and if it succeeds we can feel more *confident* about our code. Bayesian inference works identically: we update our beliefs about an outcome; rarely can we be absolutely sure unless we rule out all other alternatives.

#### The Bayesian state of mind

Bayesian inference differs from more traditional statistical inference by preserving *uncertainty* about our beliefs. At first, this sounds like a bad statistical technique. Isn't statistics all about deriving *certainty* from randomness? To reconcile this, we need to start thinking like Bayesians.

The Bayesian world-view interprets probability as measure of *believability in an event*, that is, how confident we are in an event occurring. In fact, we will see in a moment that this is the natural interpretation of probability.

For this to be clearer, we consider an alternative interpretation of probability: *Frequentist* methods assume that probability is the long-run frequency of events (hence the bestowed title). For example, the *probability of plane accidents* under a frequentist philosophy is interpreted as the *long-term frequency of plane accidents*. This makes logical sense for many probabilities of events, but becomes more difficult to understand when events have no long-term frequency of occurrences. Consider: we often assign probabilities to outcomes of presidential elections, but the election itself only happens once! Frequentists get around this by invoking alternative realities and saying across all these universes, the frequency of occurrences defines the probability.

Bayesians, on the other hand, have a more intuitive approach. Bayesians interpret a probability as measure of *belief*, or confidence, of an event occurring. Simply, a probability is a summary of an opinion. An individual who assigns a belief of 0 to an event has no confidence that the event will occur; conversely, assigning a belief of 1 implies that the individual is absolutely certain of an event occurring. Beliefs between 0 and 1 allow for weightings of other outcomes. This definition agrees with the probability of a plane accident example, for having observed the frequency of plane accidents, an individual's belief should be equal to that frequency, excluding any outside information. Similarly, under this definition of probability being equal to beliefs, it is clear how we can speak about probabilities (beliefs) of presidential election outcomes: how confident are you candidate *A* will win?

Notice in the paragraph above, I assigned the belief (probability) measure to an *individual*, not to Nature. This is very interesting, as this definition leaves room for conflicting beliefs between individuals. Again, this is appropriate for what naturally occurs: different individuals have different beliefs of events occurring, because they possess different *information* about the world. The existence of different beliefs does not imply that anyone is wrong. Consider the following examples demonstrating the relationship between individual beliefs and probabilities:

- I flip a coin, and we both guess the result. We would both agree, assuming the coin is fair, that the probability of heads is  $1/2$ . Assume, then, that I peek at the coin. Now I know for certain what the result is: I assign probability 1.0 to either heads or tails. Now what is *your* belief that the coin is heads? My knowledge of the outcome has not changed the coin's results. Thus we assign different probabilities to the result.
- Your code either has a bug in it or not, but we do not know for certain which is true, though we have a belief about the presence or absence of a bug.
- A medical patient is exhibiting symptoms  $x$ ,  $y$  and  $z$ . There are a number of diseases that could be causing all of them, but only a single disease is present. A doctor has beliefs about which disease.

This philosophy of treating beliefs as probability is natural to humans. We employ it constantly as we interact with the world and only see partial evidence. Alternatively, you have to be *trained* to think like a frequentist.

To align ourselves with traditional probability notation, we denote our belief about event  $A$  as  $P(A)$ . We call this quantity the *prior probability*.

John Maynard Keynes, a great economist and thinker, said "When the facts change, I change my mind. What do you do, sir?" This quote reflects the way a Bayesian updates his or her beliefs after seeing evidence. Even — especially — if the evidence is counter to what was initially believed, the evidence cannot be ignored. We denote our updated belief as  $P(A|X)$ , interpreted as the probability of  $A$  given the evidence  $X$ . We call the updated belief the *posterior probability* so as to contrast it with the prior probability. For example, consider the posterior probabilities (read: posterior beliefs) of the above examples, after observing some evidence  $X$ :

1.  $P(A)$  : the coin has a 50 percent chance of being heads.  $P(A|X)$  : You look at the coin, observe a heads has landed, denote this information  $X$ , and trivially assign probability 1.0 to heads and 0.0 to tails.
2.  $P(A)$  : This big, complex code likely has a bug in it.  $P(A|X)$  : The code passed all  $X$  tests; there still might be a bug, but its presence is less likely now.
3.  $P(A)$  : The patient could have any number of diseases.  $P(A|X)$  : Performing a blood test generated evidence  $X$ , ruling out some of the possible diseases from consideration.

It's clear that in each example we did not completely discard the prior belief after seeing new evidence  $X$ , but we *re-weighted the prior* to incorporate the new evidence (i.e. we put more weight, or confidence, on some beliefs versus others).

By introducing prior uncertainty about events, we are already admitting that any guess we make is potentially very wrong. After observing data, evidence, or other information, we update our beliefs, and our guess becomes *less wrong*. This is the alternative side of the prediction coin, where typically we try to be *more right*.

## Bayesian Inference in Practice

If frequentist and Bayesian inference were programming functions, with inputs being statistical problems, then the two would be different in what they return to the user. The frequentist inference function would return a number, whereas the Bayesian function would return *probabilities*.

For example, in our debugging problem above, calling the frequentist function with the argument “My code passed all  $X$  tests; is my code bug-free?” would return a *YES*. On the other hand, asking our Bayesian function “Often my code has bugs. My code passed all  $X$  tests; is my code bug-free?” would return something very different: a probabilities of *YES* and *NO*. The function might return:

*YES*, with probability 0.8; *NO*, with probability 0.2

This is very different from the answer the frequentist function returned. Notice that the Bayesian function accepted an additional argument: “*Often my code has bugs*”. This parameter is the *prior*. By including the prior parameter, we are telling the Bayesian function to include our belief about the situation. Technically this parameter in the Bayesian function is optional, but we will see excluding it has its own consequences.

## Incorporating evidence

As we acquire more and more instances of evidence, our prior belief is *washed out* by the new evidence. This is to be expected. For example, if your prior belief is something ridiculous, like “I expect the sun to explode today”, and each day you are proved wrong, you would hope that any inference would correct you, or at least align your beliefs better. Bayesian inference will correct this belief.

Denote  $N$  as the number of instances of evidence we possess. As we gather an *infinite* amount of evidence, say as  $N \rightarrow \infty$ , our Bayesian results align with frequentist results. Hence for large  $N$ , statistical inference is more or less objective. On the other hand, for small  $N$ , inference is much more *unstable*: frequentist estimates have more variance and larger confidence intervals. This is where Bayesian analysis excels. By introducing a prior, and returning probabilities (instead of a scalar estimate), we *preserve the uncertainty* that reflects the instability of statistical inference of a small  $N$  dataset.

One may think that for large  $N$ , one can be indifferent between the two techniques since they offer similar inference, and might lean towards the computational-simpler, frequentist methods. An individual in this position should consider the following quote by Andrew Gelman (2005)[1], before making such a decision:

Sample sizes are never large. If  $N$  is too small to get a sufficiently-precise estimate, you need to get more data (or make more assumptions). But once  $N$  is “large enough,” you can start subdividing the data to learn more (for example, in a public opinion poll, once you have a good estimate for the entire country, you can estimate among men and women, northerners and southerners, different age groups, etc.).  $N$  is never enough because if it were “enough” you’d already be on to the next problem for which you need more data.

## Are frequentist methods incorrect then?

**No.**

Frequentist methods are still useful or state-of-the-art in many areas. Tools like Least Squares linear regression, LASSO regression, EM algorithm etc. are all very powerful and incredibly fast. Bayesian methods are a compliment to solve the problems these solutions cannot or to gain further insight into the underlying system by offering more flexibility in modeling.

## A note on *Big Data*

Paradoxically, big data's predictive analytic problems are actually solved by relatively simple algorithms [2][4]. Thus we can argue that big data's prediction difficulty does not lie in the algorithm used, but instead on the computational difficulties of storage and execution on big data. (One should also consider Gelman's quote from above and ask "Do I really have big data?" )

The much more difficult analytic problems involve *medium data* and, especially troublesome, *really small data*. Using a similar argument as Gelman's above, if big data problems are *big enough* to be readily solved, then we should be more interested in the *not-quite-big enough* datasets.

## Our Bayesian framework

We are interested in beliefs, which can be interpreted as probabilities by thinking Bayesian. We have a *prior* belief in event  $A$ , beliefs formed by previous information, e.g., our prior belief about bugs being in our code before performing tests.

Secondly, we observe our evidence. To continue our buggy-code example: if our code passes  $X$  tests, we want to update our belief to incorporate this. We call this new belief the *posterior* probability. Updating our belief is done via the following equation, known as Bayes' Theorem, after its discoverer Thomas Bayes:

$$P(A|X) = \frac{P(X|A)P(A)}{P(X)} \quad (1)$$

(2)

$$\propto P(X|A)P(A) \quad (\propto \text{ is proportional to } ) \quad (3)$$

The above formula is not unique to Bayesian inference: it is a mathematical fact with uses outside Bayesian inference. Bayesian inference merely uses it to connect prior probabilities  $P(A)$  with an updated posterior probabilities  $P(A|X)$ .

**Example: Mandatory coin-flip example** Every statistics text must contain a coin-flipping example, I'll use it here to get it out of the way. Suppose, naively, that you are unsure about the probability of heads in a coin flip (spoiler alert: it's 50%). You believe there is some true underlying ratio, call it  $p$ , but have no prior opinion on what  $p$  might be.

We begin to flip a coin, and record the observations: either  $H$  or  $T$ . This is our observed data. An interesting question to ask is how our inference changes as we observe more and more data? More specifically, what do our posterior probabilities look like when we have little data, versus when we have lots of data.

Below we plot a sequence of updating posterior probabilities as we observe increasing amounts of data (coin flips).

```
In [4]: """
The book uses a custom matplotlibrc file, which provides the unique styles.
If executing this book, and you wish to use the book's styling, provided:
1. Overwrite your own matplotlibrc file with the rc-file provided in the book.
See http://matplotlib.org/users/customizing.html
2. Also in the styles is bmh_matplotlibrc.json file. This can be used in
only this notebook. Try running the following code:

import json
s = json.load( open("../styles/bmh_matplotlibrc.json") )
matplotlib.rcParams.update(s)

"""

#the code below can be passed over, as it is currently not important.
%pylab inline
```

```

figsize( 11, 9)

import scipy.stats as stats

dist = stats.beta
n_trials = [0,1,2,3,4,5,8,15, 50, 500]
data = stats.bernoulli.rvs(0.5, size = n_trials[-1] )
x = np.linspace(0,1,100)

for k, N in enumerate(n_trials):
    sx = subplot( len(n_trials)/2, 2, k+1)
    plt.xlabel("$p$, probability of heads") if k in [0,len(n_trials)-1] else
    plt.setp(sx.get_yticklabels(), visible=False)
    heads = data[:N].sum()
    y = dist.pdf(x, 1 + heads, 1 + N - heads )
    plt.plot( x, y, label= "observe %d tosses,\n %d heads"%(N,heads) )
    plt.fill_between( x, 0, y, color="#348ABD", alpha = 0.4 )
    plt.vlines( 0.5, 0, 4, color = "k", linestyles = "--", lw=1 )

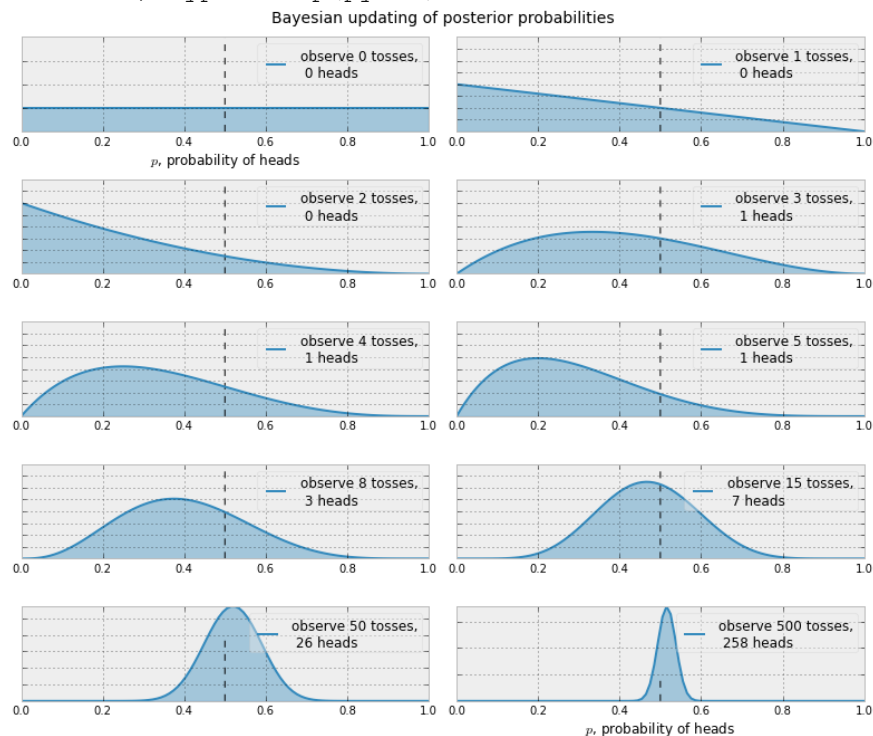
    leg = plt.legend()
    leg.get_frame().set_alpha(0.4)
    plt.autoscale(tight = True)

plt.suptitle( "Bayesian updating of posterior probabilities",
              y = 1.02,
              fontsize = 14);

plt.tight_layout()

```

Welcome to pylab, a matplotlib-based Python environment [backend: module://IPython.zmq.pylab]. For more information, type 'help(pylab)'.



The posterior probabilities are represented by the curves, and our confidence is proportional to the height of the curve. As the plot above shows, as we start to observe data our posterior probabilities start to shift and move around. Eventually, as we observe more and more data (coin-flips), our probabilities will lump closer and closer around the true value of  $p = 0.5$  (marked by a dashed line).

Notice that the plots are not always *peaked* at 0.5. There is no reason it should be: recall we assumed we did not have a prior opinion of what  $p$  is. In fact, if we observe quite extreme data, say 8 flips and only 1 observed heads, our distribution would look very biased *away* from lumping around 0.5. As more data accumulates, we would see more and more probability being assigned at  $p = 0.5$ .

The next example is a simple demonstration of the mathematics of Bayesian inference.

**Example: Bug, or just sweet, unintended feature?** Let  $A$  denote the event that our code has **no bugs** in it. Let  $X$  denote the event that the code passes all debugging tests. For now, we will leave the prior probability of no bugs as a variable, i.e.  $P(A) = p$ .

We are interested in  $P(A|X)$ , i.e. the probability of no bugs, given our debugging tests  $X$ . To use the formula above, we need to compute some quantities.

What is  $P(X|A)$ , i.e., the probability that the code passes  $X$  tests *given* there are no bugs? Well, it is equal to 1, for a code with no bugs will pass all tests.

$P(X)$  is a little bit trickier: The event  $X$  can be divided into two possibilities, event  $X$  occurring even though our code *indeed has* bugs (denoted  $\sim A$ , spoken *not A*), or event  $X$  without bugs ( $A$ ).  $P(X)$  can be represented as:

$$P(X) = P(X \text{ and } A) + P(X \text{ and } \sim A) \quad (4)$$

(5)

$$= P(X|A)P(A) + P(X|\sim A)P(\sim A) \quad (6)$$

(7)

$$= P(X|A)p + P(X|\sim A)(1 - p) \quad (8)$$

We have already computed  $P(X|A)$  above. On the other hand,  $P(X|\sim A)$  is subjective: our code can pass tests but still have a bug in it, though the probability there is a bug present is reduced. Note this is dependent on the number of tests performed, the degree of complication in the tests, etc. Let's be conservative and assign  $P(X|\sim A) = 0.5$ . Then

$$P(A|X) = \frac{1 \cdot p}{1 \cdot p + 0.5(1 - p)} \quad (9)$$

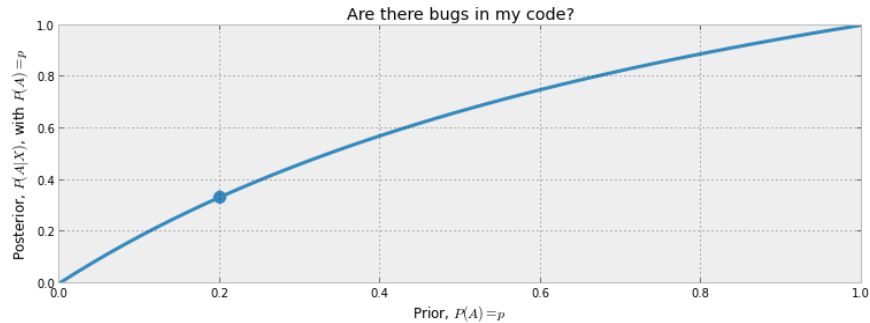
(10)

$$= \frac{2p}{1 + p} \quad (11)$$

This is the posterior probability. What does it look like as a function of our prior,  $p \in [0, 1]$ ?

```
In [19]: figsize(12.5,4)
p = np.linspace( 0,1, 50)
plt.plot( p, 2*p/(1+p), color = "#348ABD", lw = 3 )
#plt.fill_between( p, 2*p/(1+p), alpha = .5, facecolor = ["#A60628"])
plt.scatter( 0.2, 2*(0.2)/1.2, s = 140, c = "#348ABD" )
plt.xlim( 0, 1)
plt.ylim( 0, 1)
plt.xlabel( "Prior, $P(A) = p$" )
plt.ylabel( "Posterior, $P(A|X)$, with $P(A) = p$" )
plt.title( "Are there bugs in my code?" );
```





We can see the biggest gains if we observe the  $X$  tests passed when the prior probability,  $p$ , is low. Let's settle on a specific value for the prior. I'm a strong programmer (I think), so I'm going to give myself a realistic prior of 0.20, that is, there is a 20% chance that I write code bug-free. To be more realistic, this prior should be a function of how complicated and large the code is, but let's pin it at 0.20. Then my updated belief that my code is bug-free is 0.33.

Recall that the prior is a probability:  $p$  is the prior probability that there *are no bugs*, so  $1 - p$  is the prior probability that there *are bugs*.

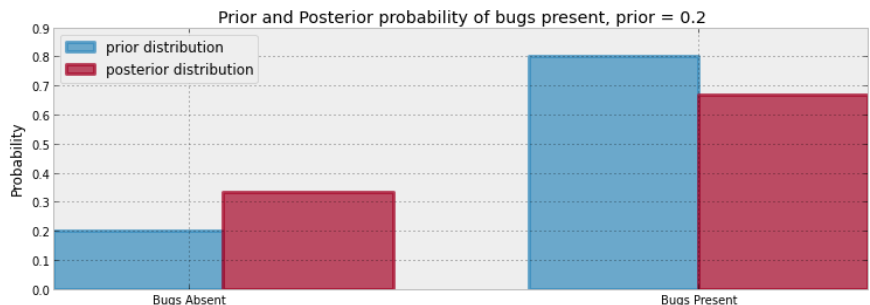
Similarly, our posterior is also a probability, with  $P(A|X)$  the probability there is no bug *given we saw all tests pass*, hence  $1 - P(A|X)$  is the probability there is a bug *given all tests passed*. What does our posterior probability look like? Below is a graph of both the prior and the posterior probabilities.

```
In [1]: figsize( 12.5, 4 )
        colours = ["#348ABD", "#A60628"]

        prior = [0.20, 0.80]
        posterior = [1./3, 2./3]
        plt.bar( [0, .7], prior ,alpha = 0.70, width = 0.25, \
                color = colours[0], label = "prior distribution",
                lw = "3", edgecolor = colours[0])

        plt.bar( [0+0.25, .7+0.25], posterior ,alpha = 0.7, \
                width = 0.25, color = colours[1],
                label = "posterior distribution",
                lw = "3", edgecolor = colours[1])

        plt.xticks( [0.20, .95], ["Bugs Absent", "Bugs Present"] )
        plt.title("Prior and Posterior probability of bugs present, prior = 0.2")
        plt.ylabel("Probability")
        plt.legend(loc="upper left");
```



Notice that after we observed  $X$  occur, the probability of bugs being absent increased. By increasing the number of tests, we can approach confidence (probability 1) that there are no bugs present.

This was a very simple example of Bayesian inference and Bayes rule. Unfortunately, the mathematics necessary to perform more complicated Bayesian inference only becomes more difficult, except for artificially constructed cases. We will later see that this type of mathematical analysis is actually unnecessary. First we must broaden our modeling

tools. The next section deals with *probability distributions*. If you are already familiar, feel free to skip (or at least skim), but for the less familiar the next section is essential.

---

## 0.3.2 Probability Distributions

**Let's quickly recall what a probability distribution is:** Let  $Z$  be some random variable. Then associated with  $Z$  is a *probability distribution function* that assigns probabilities to the different outcomes  $Z$  can take. Graphically, a probability distribution is a curve where the probability of an outcome is proportional to the height of the curve. You can see examples in the first figure of this chapter.

We can divide random variables into three classifications:

- **$Z$  is discrete:** Discrete random variables may only assume values on a specified list. Things like populations, movie ratings, and number of votes are all discrete random variables. Discrete random variables become more clear when we contrast them with...
- **$Z$  is continuous:** Continuous random variable can take on arbitrarily exact values. For example, temperature, speed, time, color are all modeled as continuous variables because you can progressively make the values more and more precise.
- **$Z$  is mixed:** Mixed random variables assign probabilities to both discrete and continuous random variables, i.e. it is a combination of the above two categories.

### Discrete Case

If  $Z$  is discrete, then its distribution is called a *probability mass function*, which measures the probability  $Z$  takes on the value  $k$ , denoted  $P(Z = k)$ . Note that the probability mass function completely describes the random variable  $Z$ , that is, if we know the mass function, we know how  $Z$  should behave. There are popular probability mass functions that consistently appear: we will introduce them as needed, but let's introduce the first very useful probability mass function. We say  $Z$  is *Poisson*-distributed if:

$$P(Z = k) = \frac{\lambda^k e^{-\lambda}}{k!}, \quad k = 0, 1, 2, \dots$$

What is  $\lambda$ ? It is called the parameter, and it describes the shape of the distribution. For the Poisson random variable,  $\lambda$  can be any positive number. By increasing  $\lambda$ , we add more probability to larger values, and conversely by decreasing  $\lambda$  we add more probability to smaller values. One can describe  $\lambda$  as the *intensity* of the Poisson distribution.

Unlike  $\lambda$ , which can be any positive number, the value  $k$  in the above formula must be a non-negative integer, i.e.,  $k$  must take on values 0,1,2, and so on. This is very important, because if you wanted to model a population you could not make sense of populations with 4.25 or 5.612 members.

If a random variable  $Z$  has a Poisson mass distribution, we denote this by writing

$$Z \sim \text{Poi}(\lambda)$$

One very useful property of the Poisson random variable, given we know  $\lambda$ , is that its expected value is equal to the parameter, i.e.:

$$E[Z \mid \lambda] = \lambda$$

We will use this property often, so it's something useful to remember. Below we plot the probability mass distribution for different  $\lambda$  values. The first thing to notice is that by increasing  $\lambda$  we add more probability to larger values occurring. Secondly, notice that although the graph ends at 15, the distributions do not. They assign positive probability to every non-negative integer.

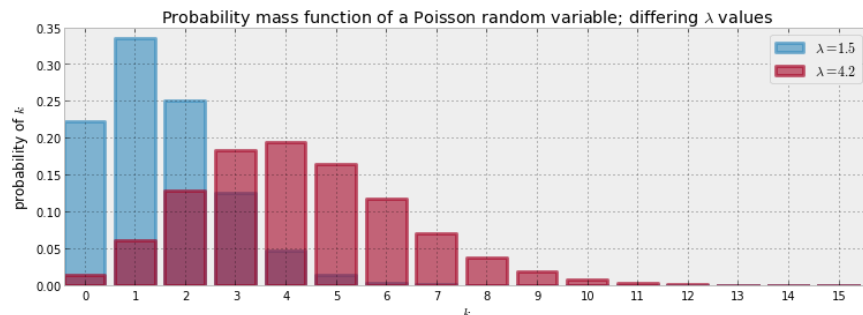
```
In [48]: figsize( 12.5, 4)

import scipy.stats as stats
a = np.arange( 16 )
poi = stats.poisson
lambda_ = [1.5, 4.25 ]

plt.bar( a, poi.pmf( a, lambda_[0]), color=colours[0],
        label = "\lambda = %.1f"%lambda_[0], alpha = 0.60,
        edgecolor = colours[0], lw = "3")

plt.bar( a, poi.pmf( a, lambda_[1]), color=colours[1],
        label = "\lambda = %.1f"%lambda_[1], alpha = 0.60,
        edgecolor = colours[1], lw = "3")

plt.xticks( a + 0.4, a )
plt.legend()
plt.ylabel("probability of $k$")
plt.xlabel("$k$")
plt.title("Probability mass function of a Poisson random variable; differing
\lambda$ values");
```



## Continuous Case

Instead of a probability mass function, a continuous random variable has a *probability density function*. This might seem like unnecessary nomenclature, but the density function and the mass function are very different creatures. An example of continuous random variable is a random variable with a *exponential density*. The density function for an exponential random variable looks like:

$$f_Z(z|\lambda) = \lambda e^{-\lambda z}, \quad z \geq 0$$

Like the Poisson random variable, an exponential random variable can only take on non-negative values. But unlike a Poisson random variable, the exponential can take on *any* non-negative values, like 4.25 or 5.612401. This makes it a poor choice for count data, which must be integers, but a great choice for time data, or temperature data (measured in Kelvins, of course), or any other precise *and positive* variable. Below are two probability density functions with different  $\lambda$  value.

When a random variable  $Z$  has an exponential distribution with parameter  $\lambda$ , we say  $Z$  is *exponential* and write

$$Z \sim \text{Exp}(\lambda)$$

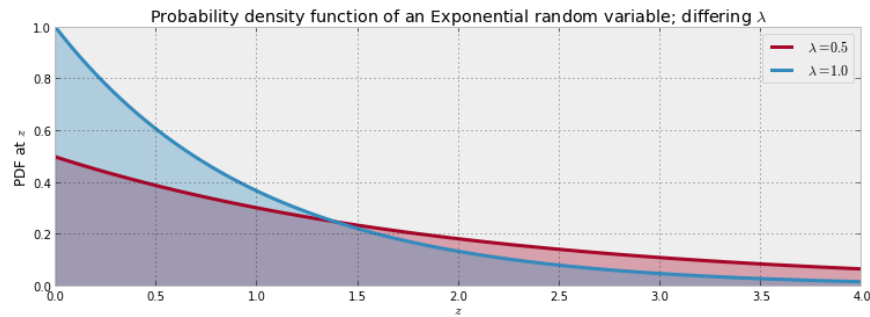
Given a specific  $\lambda$ , the expected value of an exponential random variable is equal to the inverse of  $\lambda$ , that is:

$$E[Z | \lambda] = \frac{1}{\lambda}$$

```
In [6]: a = np.linspace(0,4, 100)
expo = stats.expon
lambda_ = [0.5, 1]

for l,c in zip(lambda_, colours):
    plt.plot( a, expo.pdf( a, scale=1./l), lw=3,
              color=c, label = "$\lambda = %.1f$" % l)
    plt.fill_between( a, expo.pdf( a, scale=1./l), color=c, alpha = .33)

plt.legend()
plt.ylabel("PDF at $z$")
plt.xlabel("$z$")
plt.title("Probability density function of an Exponential random variable;
differing $\lambda$");
```



## But what is $\lambda$ ?

**This question is what motivates statistics.** In the real world,  $\lambda$  is hidden from us. We only see  $Z$ , and must go backwards to try and determine  $\lambda$ . The problem is so difficult because there is not a one-to-one mapping from  $Z$  to  $\lambda$ . Many different methods have been created to solve the problem of estimating  $\lambda$ , but since  $\lambda$  is never actually observed, no one can say for certain which method is better!

Bayesian inference is concerned with *beliefs* about what  $\lambda$  is. Rather than try to guess  $\lambda$  exactly, we can only talk about what  $\lambda$  is likely to be by assigning a probability distribution to  $\lambda$ .

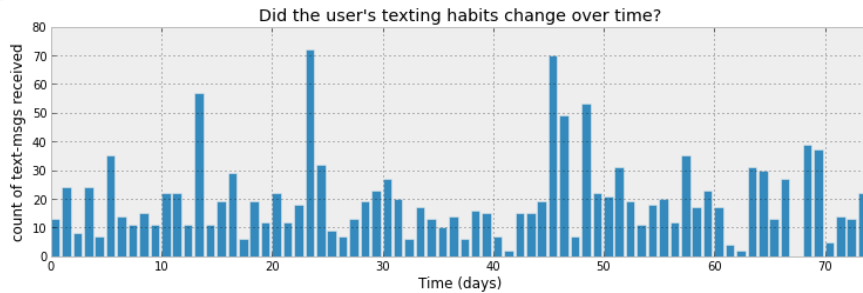
This might seem odd at first: after all,  $\lambda$  is fixed, it is not (necessarily) random! How can we assign probabilities to a non-random event. Ah, we have fallen for the frequentist interpretation. Recall, under our Bayesian philosophy, we *can* assign probabilities if we interpret them as beliefs. And it is entirely acceptable to have *beliefs* about the parameter  $\lambda$ .

**Example: Inferring behaviour from text-message data** Let's try to model a more interesting example, concerning text-message rates:

You are given a series of text-message counts from a user of your system. The data, plotted over time, appears in the graph below. You are curious if the user's text-messaging habits changed over time, either gradually or suddenly. How can you model this? (This is in fact my own text-message data. Judge my popularity as you wish.)

```
In [5]: figsize( 12.5, 3.5 )
count_data = np.loadtxt("data/txtdata.csv")
n_count_data = len(count_data)
```

```
plt.bar( np.arange( n_count_data ), count_data, color = "#348ABD" )
plt.xlabel( "Time (days)" )
plt.ylabel("count of text-msgs received")
plt.title("Did the user's texting habits change over time?")
plt.xlim( 0, n_count_data );
```



Before we begin, with respect to the plot above, would you say there was a change in behaviour during the time period? How can we start to model this? Well, as I conveniently already introduced, a Poisson random variable would be a very appropriate model for this *count* data. Denoting day  $i$ 's text-message count by  $C_i$ ,

$$C_i \sim \text{Poisson}(\lambda)$$

We are not sure about what the  $\lambda$  parameter is though. Looking at the chart above, it appears that the rate might become higher at some later date, which is equivalently saying the parameter  $\lambda$  increases at some later date (recall a higher  $\lambda$  means more probability on larger outcomes, that is, higher probability of many texts.).

How can we mathematically represent this? We can think, that at some later date (call it  $\tau$ ), the parameter  $\lambda$  suddenly jumps to a higher value. So we create two  $\lambda$  parameters, one for behaviour before the  $\tau$ , and one for behaviour after. In literature, a sudden transition like this would be called a *switchpoint*:

$$\lambda = \begin{cases} \lambda_1 & \text{if } t < \tau \\ \lambda_2 & \text{if } t \geq \tau \end{cases}$$

If, in reality, no sudden change occurred and indeed  $\lambda_1 = \lambda_2$ , the  $\lambda$ 's posterior distributions should look about equal.

We are interested in inferring the unknown  $\lambda$ s. To use Bayesian inference, we need to assign prior probabilities to the different possible values of  $\lambda$ . What would be good prior probability distributions for  $\lambda_1$  and  $\lambda_2$ ? Recall that  $\lambda_i$ ,  $i = 1, 2$ , can be any positive number. The *exponential* random variable has a density function for any positive number. This would be a good choice to model  $\lambda_i$ . But, we need a parameter for this exponential distribution: call it  $\alpha$ .

$$\lambda_1 \sim \text{Exp}(\alpha) \tag{12}$$

$$\lambda_2 \sim \text{Exp}(\alpha) \tag{13}$$

$\alpha$  is called a *hyper-parameter*, or a *parent-variable*, literally a parameter that influences other parameters. The influence is not too strong, so we can choose  $\alpha$  liberally. A good rule of thumb is to set the exponential parameter equal to the inverse of the average of the count data, since we're modeling *lambda* using an Exponential distribution we can use the expected value identity shown earlier to get:

$$\frac{1}{N} \sum_{i=0}^N C_i \approx E[\lambda | \alpha] = \frac{1}{\alpha}$$

Alternatively, and something I encourage the reader to try, is to have two priors: one for each  $\lambda_i$ ; creating two exponential distributions with different  $\alpha$  values reflects a prior belief that the rate changed after some period.

What about  $\tau$ ? Well, due to the randomness, it is too difficult to pick out when  $\tau$  might have occurred. Instead, we can assign an *uniform prior belief* to every possible day. This is equivalent to saying

$$\tau \sim \text{DiscreteUniform}(1,70) \tag{14}$$

$$\tag{15}$$

$$\Rightarrow P(\tau = k) = \frac{1}{70} \tag{16}$$

So after all this, what does our overall prior for the unknown variables look like? Frankly, *it doesn't matter*. What we should understand is that it would be an ugly, complicated, mess involving symbols only a mathematician would love. And things would only get uglier the more complicated our models become. Regardless, all we really care about is the posterior distribution. We next turn to PyMC, a Python library for performing Bayesian analysis, that is agnostic to the mathematical monster we have created.

### 0.3.3 Introducing our first hammer: PyMC

PyMC is a Python library for programming Bayesian analysis [3]. It is a fast, well-maintained library. The only unfortunate part is that documentation can be lacking in areas, especially the bridge between beginner to hacker. One of this book's main goals is to solve that problem, and also to demonstrate why PyMC is so cool.

We will model the above problem using the PyMC library. This type of programming is called *probabilistic programming*, an unfortunate misnomer that invokes ideas of randomly-generated code and has likely confused and frightened users away from this field. The code is not random. The title is given because we create probability models using programming variables as the model's components, that is, model components are first-class primitives in this framework.

B. Cronin [5] has a very motivating description of probabilistic programming:

Another way of thinking about this: unlike a traditional program, which only runs in the forward directions, a probabilistic program is run in both the forward and backward direction. It runs forward to compute the consequences of the assumptions it contains about the world (i.e., the model space it represents), but it also runs backward from the data to constrain the possible explanations. In practice, many probabilistic programming systems will cleverly interleave these forward and backward operations to efficiently home in on the best explanations.

Due to its poorly understood title, I'll refrain from using the name *probabilistic programming*. Instead, I'll simply use *programming*, as that is what it really is.

The PyMC code is easy to follow along: the only novel thing should be the syntax, and I will interrupt the code to explain sections. Simply remember we are representing the model's components ( $\tau, \lambda_1, \lambda_2$ ) as variables:

```
In [6]: import pymc as mc

alpha = 1.0/count_data.mean() #recall count_data is
                               #the variable that holds our txt counts

lambda_1 = mc.Exponential( "lambda_1", alpha )
lambda_2 = mc.Exponential( "lambda_2", alpha )

tau = mc.DiscreteUniform( "tau", lower = 0, upper = n_count_data )
```

In the above code, we create the PyMC variables corresponding to  $\lambda_1$ ,  $\lambda_2$ . We assign them to PyMC's *stochastic variables*, called stochastic variables because they are treated by the backend as random number generators. We can test this by calling their built-in `random()` method.

```
In [7]: print "Random output:", tau.random(),tau.random(), tau.random()
```

Random output: 62 61 53

```
In [8]: @mc.deterministic
def lambda_( tau = tau, lambda_1 = lambda_1, lambda_2 = lambda_2 ):
    out = np.zeros( n_count_data )
    out[:tau] = lambda_1 #lambda before tau is lambda1
    out[tau:] = lambda_2 #lambda after tau is lambda2
    return out
```

This code is creating a new function `lambda_`, but really we think of it as a random variable: the random variable  $\lambda$  from above. Note that because `lambda_1`, `lambda_2` and `tau` are random, `lambda_` will be random. We are **not** fixing any variables yet. The `@mc.deterministic` is a decorator to tell PyMC that this is a deterministic function, i.e., if the arguments were deterministic (which they are not), the output would be deterministic as well.

```
In [9]: observation = mc.Poisson( "obs", lambda_, value = count_data, observed = True
model = mc.Model( [observation, lambda_1, lambda_2, tau] )
```

The variable `observation` combines our data, `count_data`, with our proposed data-generation scheme, given by the variable `lambda_`, through the `value` keyword. We also set `observed = True` to tell PyMC that this should stay fixed in our analysis. Finally, PyMC wants us to collect all the variables of interest and create a `Model` instance out of them. This makes our life easier when we try to retrieve the results.

The below code will be explained in the Chapter 3, but this is where our results come from. One can think of it as a *learning* step. The machinery being employed is called *Markov Chain Monte Carlo* (which I delay explaining until Chapter 3). It returns thousands of random variables from the posterior distributions of  $\lambda_1$ ,  $\lambda_2$  and  $\tau$ . We can plot a histogram of the random variables to see what the posterior distribution looks like. Below, we collect the samples (called *traces* in MCMC literature) in histograms.

```
In [14]: ### Mysterious code to be explained in Chapter 3.
mcmc = mc.MCMC(model)
mcmc.sample( 40000, 10000, 1 )
```

[\*\*\*\*\*100%\*\*\*\*\*] 40000 of 40000 complete

```
In [15]: lambda_1_samples = mcmc.trace( 'lambda_1' )[:]
lambda_2_samples = mcmc.trace( 'lambda_2' )[:]
tau_samples = mcmc.trace( 'tau' )[:]
```

```
In [33]: figsize(12.5, 10)
#histogram of the samples:

ax = plt.subplot(311)
ax.set_autoscaley_on(False)

plt.hist( lambda_1_samples, histtype='stepfilled', bins = 30, alpha = 0.85,
          label = "posterior of $\lambda_1$", color = "#A60628",normed = True)
plt.legend(loc = "upper left")
plt.title(r"Posterior distributions of the variables $\lambda_1, \lambda_2$")
plt.xlim([15,30])
plt.xlabel("$\lambda_2$ value")
plt.ylabel("probability")
```

```

ax = plt.subplot(312)
ax.set_autoscaley_on(False)

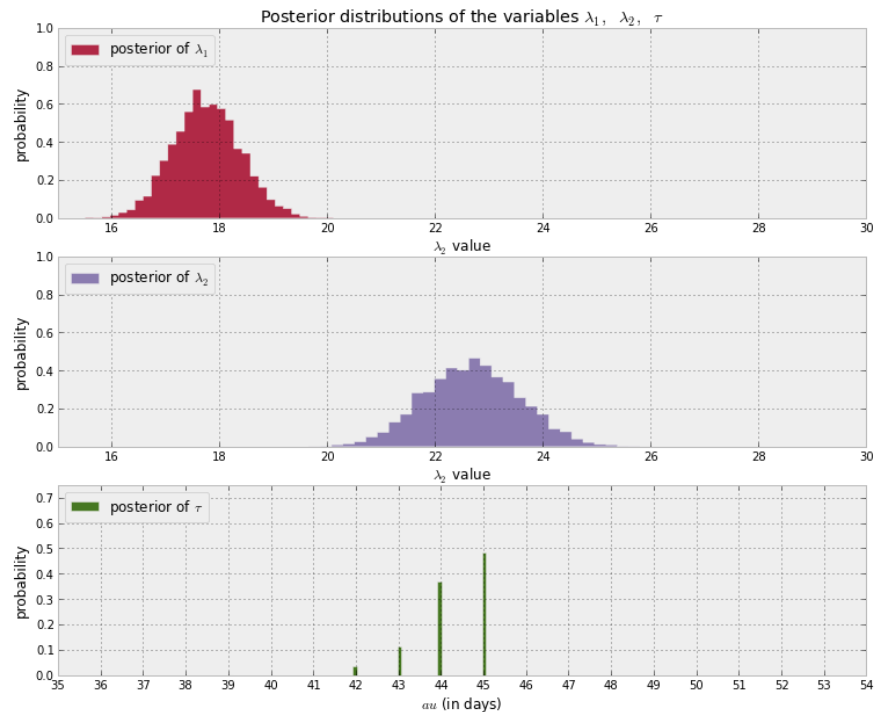
plt.hist( lambda_2_samples, histtype='stepfilled', bins = 30, alpha = 0.85,
          label = "posterior of  $\lambda_2$ ", color="#7A68A6", normed = True)
plt.legend(loc = "upper left")
plt.xlim([15,30])
plt.xlabel(" $\lambda_2$  value")
plt.ylabel("probability")

plt.subplot(313)

w = 1.0/ tau_samples.shape[0] * np.ones_like( tau_samples )
plt.hist( tau_samples, bins = n_count_data, alpha = 1,
          label = r"posterior of  $\tau$ ",
          color="#467821", weights=w, rwidth =2. )
plt.xticks( np.arange( n_count_data ) )

plt.legend(loc = "upper left");
plt.ylim([0, .75])
plt.xlim([35, len(count_data)-20])
plt.xlabel(" $\tau$  (in days)")
plt.ylabel("probability");

```



## Interpretation

Recall that the Bayesian methodology returns a *distribution*, hence we now have distributions to describe the unknown  $\lambda$ 's and  $\tau$ . What have we gained? Immediately we can see the uncertainty in our estimates: the more variance in the distribution, the less certain our posterior belief should be. We can also say what a plausible value for the parameters might be:  $\lambda_1$  is around 18 and  $\lambda_2$  is around 23. What other observations can you make? Look at the data again, do these seem reasonable? The distributions of the two *lambdas* are positioned very differently, indicating that it's likely there was a change in the user's text-message behaviour.



Also notice that the posterior distributions for the  $\lambda$ 's do not look like any exponential distributions, though we originally started modeling with exponential random variables. They are really not anything we recognize. But this is OK. This is one of the benefits of taking a computational point-of-view. If we had instead done this mathematically, we would have been stuck with a very analytically intractable (and messy) distribution. Via computations, we are agnostic to the tractability.

Our analysis also returned a distribution for what  $\tau$  might be. Its posterior distribution looks a little different from the other two because it is a discrete random variable, hence it doesn't assign probabilities to intervals. We can see that near day 45, there was a 50% chance the users behaviour changed. Had no change occurred, or the change been gradual over time, the posterior distribution of  $\tau$  would have been more spread out, reflecting that many values are likely candidates for  $\tau$ . On the contrary, it is very peaked.

## Why would I want samples from the posterior, anyways?

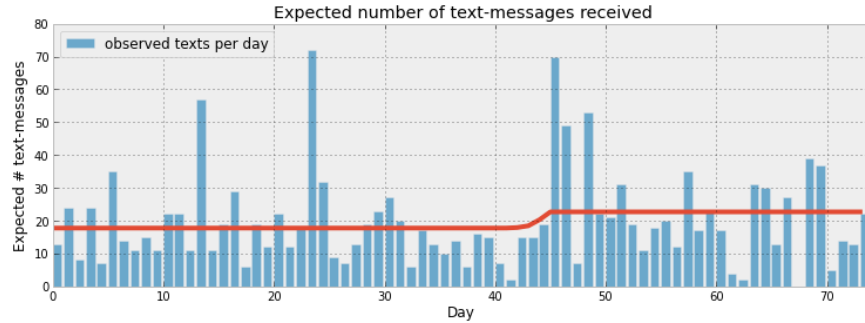
We will deal with this question for the remainder of the book, and it is an understatement to say we can perform amazingly useful things. For now, let's end this chapter with one more example. We'll use the posterior samples to answer the following question: what is the expected number of texts at day  $t$ ,  $0 \leq t \leq 70$ ? Recall that the expected value of a Poisson is equal to its parameter  $\lambda$ , then the question is equivalent to *what is the expected value of  $\lambda$  at time  $t$ ?*

In the code below, we are calculating the following: Let  $i$  index samples from the posterior distributions. Given a day  $t$ , we average over all possible  $\lambda_i$  for that day  $t$ , using  $\lambda_i = \lambda_{1,i}$  if  $t < \tau_i$  (that is, if the behaviour change hadn't occurred yet), else we use  $\lambda_i = \lambda_{2,i}$ .

```
In [25]: figsize( 12.5, 4)
# tau_samples, lambda_1_samples, lambda_2_samples contain
# N samples from the corresponding posterior distribution
N = tau_samples.shape[0]
expected_texts_per_day = np.zeros(n_count_data)
for day in range(0, n_count_data):
    # ix is a bool index of all tau samples corresponding to
    # the switchpoint occurring prior to value of 'day'
    ix = day < tau_samples
    # Each posterior sample corresponds to a value for tau.
    # for each day, that value of tau indicates whether we're "before"
    # (in the lambda1 "regime") or
    # "after" (in the lambda2 "regime") the switchpoint.
    # by taking the posterior sample of lambda1/2 accordingly, we can average
    # over all samples to get an expected value for lambda on that day.
    # As explained, the "message count" random variable is Poisson distributed
    # and therefore lambda (the poisson parameter) is the expected value of the count.
    expected_texts_per_day[day] = (lambda_1_samples[ix].sum()
                                   + lambda_2_samples[~ix].sum() ) / N

plt.plot( range( n_count_data), expected_texts_per_day, lw =4, color = "#1f77b4")
plt.xlim( 0, n_count_data )
plt.xlabel( "Day" )
plt.ylabel( "Expected # text-messages" )
plt.title( "Expected number of text-messages received")
#plt.ylim( 0, 35 )
plt.bar( np.arange( len(count_data) ), count_data, color = "#348abd", alpha=0.5,
         label="observed texts per day")

plt.legend(loc="upper left");
```



Our analysis shows strong support for believing the user’s behavior did change ( $\lambda_1$  would have been close in value to  $\lambda_2$  had this not been true), and the change was sudden rather than gradual (demonstrated by  $\tau$ ’s strongly peaked posterior distribution). We can speculate what might have caused this: a cheaper text-message rate, a recent weather-2-text subscription, or a new relationship. (The 45th day corresponds to Christmas, and I moved away to Toronto the next month leaving a girlfriend behind.)

**Exercises** 1. Using `lambda_1_samples` and `lambda_2_samples`, what is the mean of the posterior distributions of  $\lambda_1$  and  $\lambda_2$ ?

In [17]: `#type your code here.`

2. What is the expected percentage increase in text-message rates? hint: compute the mean of `lambda_1_samples/lambda_2_samples`. Note that this quantity is very different from `lambda_1_samples.mean()/lambda_2_samples.mean()`.

In [16]: `#type your code here.`

3. What is the mean of  $\lambda_1$  **given** we know  $\tau$  is less than 45. That is, suppose we have new information as we know for certain that the change in behaviour occurred before day 45. What is the expected value of  $\lambda_1$  now? (You do not need to redo the PyMC part, just consider all instances where `tau_samples<45`.)

In []: `#type your code here.`

## References

- [1] Gelman, Andrew. N.p.. Web. 22 Jan 2013. [http://andrewgelman.com/2005/07/n\\_is\\_never\\_larg/](http://andrewgelman.com/2005/07/n_is_never_larg/).
- [2] Norvig, Peter. 2009. *The Unreasonable Effectiveness of Data*.
- [3] Patil, A., D. Huard and C.J. Fonnesebeck. 2010. PyMC: Bayesian Stochastic Modelling in Python. *Journal of Statistical Software*, 35(4), pp. 1-81.
- [4] Jimmy Lin and Alek Kolcz. Large-Scale Machine Learning at Twitter. *Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data (SIGMOD 2012)*, pages 793-804, May 2012, Scottsdale, Arizona.
- [5] Cronin, Beau. “Why Probabilistic Programming Matters.” 24 Mar 2013. Google, Online Posting to Google . Web. 24 Mar. 2013. <https://plus.google.com/u/0/107971134877020469960/posts/KpeRdJKR6Z1>.

```
In [3]: from IPython.core.display import HTML
def css_styling():
    styles = open("../styles/custom.css", "r").read()
    return HTML(styles)
css_styling()
```

Out[3]:

## 0.4 Chapter 2

---

This chapter introduces more PyMC syntax and design patterns, and ways to think about how to model a system from a Bayesian perspective. It also contains tips and data visualization techniques for assessing goodness-of-fit for your Bayesian model.

### 0.4.1 A little more on PyMC

#### Parent and Child relationships

To assist with describing Bayesian relationships, and to be consistent with PyMC's documentation, we introduce *parent and child* variables.

- *parent variables* are variables that influence another variable.
- *child variable* are variables that are affected by other variables, i.e. are the subject of parent variables.

A variable can be both a parent and child. For example, consider the PyMC code below.

```
In [8]: import pymc as mc
parameter = mc.Exponential( "poisson_param", 1 )
data_generator = mc.Poisson("data_generator", parameter )
data_plus_one = data_generator + 1
```

`parameter` controls the parameter of `data_generator`, hence influences its values. The former is a parent of the latter. By symmetry, `data_generator` is a child of `parameter`.

Likewise, `data_generator` is a parent to the variable `data_plus_one` (hence making `data_generator` both a parent and child variable). Although it does not look like one, `data_plus_one` should be treated as a PyMC variable as it is a *function* of another PyMC variable, hence is a child variable to `data_generator`.

This nomenclature is introduced to help us describe relationships in PyMC modeling. You can access a variables children and parent variables using the `children` and `parents` attributes attached to variables.

```
In [9]: print "Children of `parameter`: "
print parameter.children
print "\nParents of `data_generator`: "
print data_generator.parents
print "\nChildren of `data_generator`: "
print data_generator.children
```

```

Children of `parameter`:
set([<pymc.distributions.Poisson 'data_generator' at 0x0000000008DF5198>])

Parents of `data_generator`:
'mu': <pymc.distributions.Exponential 'poisson_param' at 0x0000000008DF57B8>

Children of `data_generator`:
set([<pymc.PyMCObjects.Deterministic '(data_generator_add_1)' at 0x0000000008DE3CC0>])

```

Of course a child can have more than one parent, and a parent can have many children.

## PyMC Variables

All PyMC variables also expose a `value` attribute. This method produces the *current* (possibly random) internal value of the variable. If the variable is a child variable, its value changes given the variable's parents' values. Using the same variables from before:

```
In [10]: print "parameter.value =", parameter.value
         print "data_generator.value =", data_generator.value
         print "data_plus_one.value =", data_plus_one.value
```

```
parameter.value = 2.13659208293
data_generator.value = 3
data_plus_one.value = 4
```

PyMC is concerned with two types of programming variables: *stochastic* and *deterministic*.

- *stochastic variables* are variables that are not deterministic, i.e., even if you knew all the values of the variables' parents (if it even has any parents), it would still be random. Included in this category are instances of classes `Poisson`, `DiscreteUniform`, and `Exponential`.
- *deterministic variables* are variables that are not random if the variables' parents were known. This might be confusing at first: a quick mental check is *if I knew all of variable foo's parent variables, I could determine what foo's value is*.

We will detail each below.

### Initializing Stochastic variables

Initializing a stochastic variable requires a `name` argument, plus additional parameters that are class specific. For example:

```
some_variable = mc.DiscreteUniform( "discrete_uni_var", 0, 4 )
```

where 0,4 are the `DiscreteUniform`-specific upper and lower bound on the random variable. The [PyMC docs](#) contain the specific parameters for stochastic variables. (Or use `??` if you are using IPython!)

The `name` attribute is used to retrieve the posterior distribution later in the analysis, so it is best to use a descriptive name. Typically, I use the Python variable's name as the name.

For multivariable problems, rather than creating a Python array of stochastic variables, addressing the `size` keyword in the call to a *Stochastic* variable creates multivariate array of (independent) stochastic variables. The array behaves like a Numpy array when used like one, and references to its `value` attribute return Numpy arrays.

The `size` argument also solves the annoying case where you may have many variables  $\beta_i$ ,  $i = 1, \dots, N$  you wish to model. Instead of creating arbitrary names and variables for each one, like:

```
beta_1 = mc.Uniform( "beta_1", 0, 1)
beta_2 = mc.Uniform( "beta_2", 0, 1)
...
```

we can instead wrap them into a single variable:

```
betas = mc.Uniform( "betas", 0, 1, size = N )
```

### Calling `random()`

We can also call on a stochastic variable's `random()` method, which (given the parent values) will generate a new, random value. Below we demonstrate this using the texting example from the previous chapter.

```
In [11]: lambda_1 = mc.Exponential( "lambda_1", 1 ) #prior on first behaviour
lambda_2 = mc.Exponential( "lambda_2", 1 ) #prior on second behaviour
tau = mc.DiscreteUniform( "tau", lower = 0, upper = 10 ) #prior on behaviour

print "lambda_1.value = %.3f"%lambda_1.value
print "lambda_2.value = %.3f"%lambda_2.value
print "tau.value = %.3f"%tau.value
print

lambda_1.random(), lambda_2.random(), tau.random()

print "After calling random() on the variables..."
print "lambda_1.value = %.3f"%lambda_1.value
print "lambda_2.value = %.3f"%lambda_2.value
print "tau.value = %.3f"%tau.value
```

```
lambda_1.value = 0.450
lambda_2.value = 0.807
tau.value = 6.000
```

After calling `random()` on the variables...

```
lambda_1.value = 0.847
lambda_2.value = 3.247
tau.value = 3.000
```

The call to `random` stores a new value into the variable's `value` attribute. In fact, this new value is stored in the computer's cache for faster recall and efficiency.

### Warning: *Don't update stochastic variables' values in-place.*

Straight from the PyMC docs, we quote [4]:

Stochastic objects' values should not be updated in-place. This confuses PyMC's caching scheme... The only way a stochastic variable's value should be updated is using statements of the following form:

```
A.value = new_value
```

The following are in-place updates and should **never** be used:

```
A.value += 3
A.value[2,1] = 5
A.value.attribute = new_attribute_value
```

## Deterministic variables

Since most variables you will be modeling are stochastic, we distinguish deterministic variables with a `pymc.deterministic` wrapper. (If you are unfamiliar with Python wrappers (also called decorators), that's no problem. Just prepend the `pymc.deterministic` decorator before the variable declaration and you're good to go. No need to know more. ) The declaration of a deterministic variable uses a Python function:

```
@mc.deterministic
def some_deterministic_var(v1=v1,):
    #jelly goes here.
```

For all purposes, we can treat the object `some_deterministic_var` as a variable and not a Python function.

Prepending with the wrapper is the easiest way, but not the only way, to create deterministic variables. This is not completely true: elementary operations, like addition, exponentials etc. implicitly create deterministic variables. For example, the following returns a deterministic variable:

```
In [12]: type( lambda_1 + lambda_2 )
```

```
pymc.PyMCObjects.Deterministic
Out [12]:
```

The use of the `deterministic` wrapper was seen in the previous chapter's text-message example. Recall the model for  $\lambda$  looked like:

$$\lambda = \lambda_1 \text{ if } t < \tau \lambda_2 \text{ if } t \geq \tau$$

And in PyMC code:

```
In [13]: n_data_points = 5 # in CH1 we had ~70 data points

@mc.deterministic
def lambda_( tau = tau, lambda_1 = lambda_1, lambda_2 = lambda_2 ):
    out = np.zeros(n_data_points)
    out[:tau] = lambda_1 #lambda before tau is lambda1
    out[tau:] = lambda_2 #lambda after tau is lambda1
    return out
```

Clearly, if  $\tau$ ,  $\lambda_1$  and  $\lambda_2$  are known, then  $\lambda$  is known completely, hence it is a deterministic variable.

Inside the deterministic decorator, the Stochastic variables passed in behave like scalars or Numpy arrays ( if multivariable), and *not* like Stochastic variables. For example, running the following:

```
@mc.deterministic
def some_deterministic( stoch = some_stochastic_var ):
    return stoch.value**2
```

will return an `AttributeError` detailing that `stoch` does not have a `value` attribute. It simply needs to be `stoch**2`. During the learning phase, it the variables `value` that is repeatedly passed in, not the actual variable.

Notice in the creation of the deterministic function we added defaults to each variable used in the function. This is a necessary step, and all variables *must* have default values.

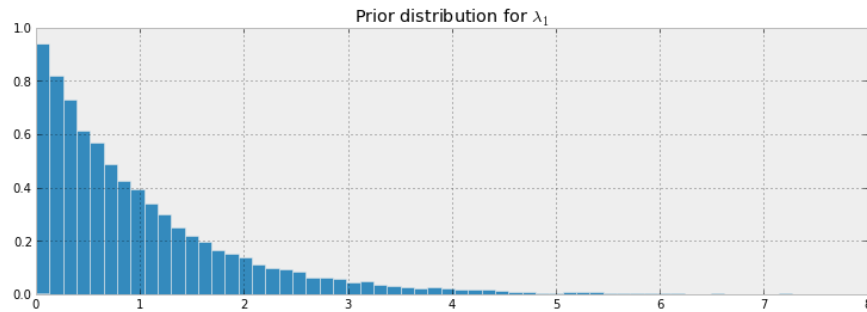
## Including observations in the Model

At this point, it may not look like it, but we have fully specified our priors. For example, we can ask and answer questions like “What does my prior distribution of  $\lambda_1$  look like?”

```
In [14]: %pylab inline
figsize(12.5, 4)

samples = [ lambda_1.random() for i in range( 20000) ]
hist( samples, bins = 70, normed=True )
plt.title( "Prior distribution for  $\lambda_1$ " )
plt.xlim( 0, 8);
```

(0, 8)  
Out [14]:



To frame this in the notation of the first chapter, though this is a slight abuse of notation, we have specified  $P(A)$ . Our next goal is to include data/evidence/observations  $X$  into our model.

PyMC stochastic variables have a keyword argument `observed` which accepts a boolean (`False` by default). The keyword `observed` has a very simple role: fix the variable's current value, i.e. make `value` immutable. We have to specify an initial `value` in the variable's creation, equal to the observations we wish to include, typically an array (and it should be a Numpy array for speed). For example:

```
In [17]: data = np.array( [10, 5] )
fixed_variable = mc.Poisson( "fxd", 1, value = data, observed = True )
print "value: ", fixed_variable.value
print "calling .random()"
fixed_variable.random()
print "value: ", fixed_variable.value
```

```
value: [10 5]
calling .random()
value: [10 5]
```

This is how we include data into our models: initializing a stochastic variable to have a *fixed value*.

To complete our text message example, we fix the PyMC variable `observations` to the observed dataset.

```
In [18]: #we're using some fake data here
data = np.array( [ 10, 25, 15, 20, 35] )
obs = mc.Poisson( "obs", lambda_, value = data, observed = True )
print obs.value
```

```
[10 25 15 20 35]
```

## Finally...

We wrap all the created variables into a `mc.Model` class. With this `Model` class, we can analyze the variables as a single unit.

```
In [19]: model = mc.Model( [obs, lambda_, lambda_1, lambda_2, tau] )
```

## 0.4.2 Modeling approaches

A good starting thought to Bayesian modeling is to think about *how your data might have been generated*. Position yourself in an omniscient position, and try to imagine how *you* would recreate the dataset.

In the last chapter we investigated text message data. We begin by asking how our observations may have been generated:

1. We started by thinking “what is the best random variable to describe this count data?” A Poisson random variable is a good candidate because it can represent count data. So we model the number of sms’s received as sampled from a Poisson distribution.
2. Next, we think, “Ok, assuming sms’s are Poisson-distributed, what do I need for the Poisson distribution?” Well, the Poisson distribution has a parameters  $\lambda$ .
3. Do we know  $\lambda$ ? No. In fact, we have a suspicion that there are *two*  $\lambda$  values, one for the earlier behaviour and one for the latter behaviour. We don’t know when the behaviour switches though, but call the switchpoint  $\tau$ .
4. What is a good distribution for the two  $\lambda$ s? The exponential is good, as it assigns probabilities to positive real numbers. Well the exponential distribution has a parameter too, call it  $\alpha$ .
5. Do we know what the parameter  $\alpha$  might be? No. At this point, we could continue and assign a distribution to  $\alpha$ , but it’s better to stop once we reach a set level of ignorance: whereas we have a prior belief about  $\lambda$ , (“it probably changes over time”, “it’s likely between 10 and 30”, etc.), we don’t really have any strong beliefs about  $\alpha$ . So it’s best to stop here.

What is a good value for  $\alpha$  then? We think that the  $\lambda$ s are between 10-30, so if we set  $\alpha$  really low (which corresponds to larger probability on high values) we are not reflecting our prior well. Similar, a too-high alpha misses our prior belief as well. A good idea for  $\alpha$  as to reflect our belief is to set the value so that the mean of  $\lambda$ , given  $\alpha$ , is equal to our observed mean. This was shown in the last chapter.

6. We have no expert opinion of when  $\tau$  might have occurred. So we will suppose  $\tau$  is from a discrete uniform distribution over the entire timespan.

Below we give a graphical visualization of this, where arrows denote parent-child relationships. (provided by the [Daft Python library](#) )

PyMC, and other probabilistic programming languages, have been designed to tell these data-generation *stories*. More generally, B. Cronin writes [5]:

Probabilistic programming will unlock narrative explanations of data, one of the holy grails of business analytics and the unsung hero of scientific persuasion. People think in terms of stories - thus the unreasonable power of the anecdote to drive decision-making, well-founded or not. But existing analytics largely fails to provide this kind of story; instead, numbers seemingly appear out of thin air, with little of the causal context that humans prefer when weighing their options.

### Same story; different ending.

Interestingly, we can create *new datasets* by retelling the story. For example, if we reverse the above steps, we can simulate a possible realization of the dataset.

1. Specify when the user’s behaviour switches by sampling from `DiscreteUniform(0, 80)`:

```
In [20]: tau = mc.rdiscrete_uniform(0, 80)
         print tau
```

63

2. Draw  $\lambda_1$  and  $\lambda_2$  from an  $\text{Exp}(\alpha)$  distribution:



```
In [21]: alpha = 1./20.
lambda_1, lambda_2 = mc.rexponential( alpha, 2 )
print lambda_1, lambda_2
```

16.7630558157 33.1568046968

3. For days before  $\tau$ , represent the user's received SMS count by sampling from  $\text{Poi}(\lambda_1)$ , and sample from  $\text{Poi}(\lambda_2)$  for days after  $\tau$ . For example:

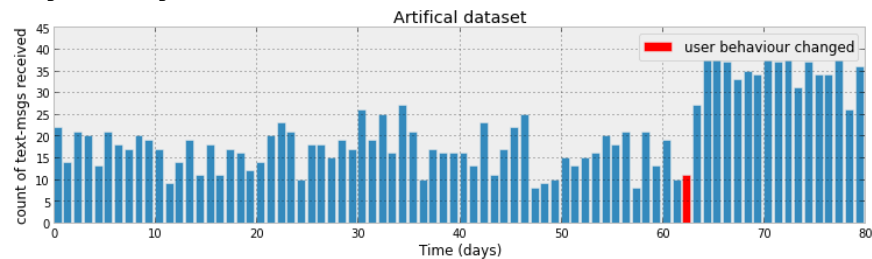
```
In [22]: data = np.r_[ mc.rpoisson( lambda_1, tau ), mc.rpoisson( lambda_2, 80 - tau )
```

4. Plot the artificial dataset:

```
In [23]: plt.bar( np.arange( 80 ), data, color = "#348ABD" )
plt.bar( tau-1, data[tau-1], color = "r", label = "user behaviour changed" )
plt.xlabel( "Time (days)" )
plt.ylabel( "count of text-msgs received" )
plt.title( "Artificial dataset" )
plt.xlim( 0, 80 );
plt.legend();
```

<matplotlib.legend.Legend at 0x12f67710>

Out [23]:

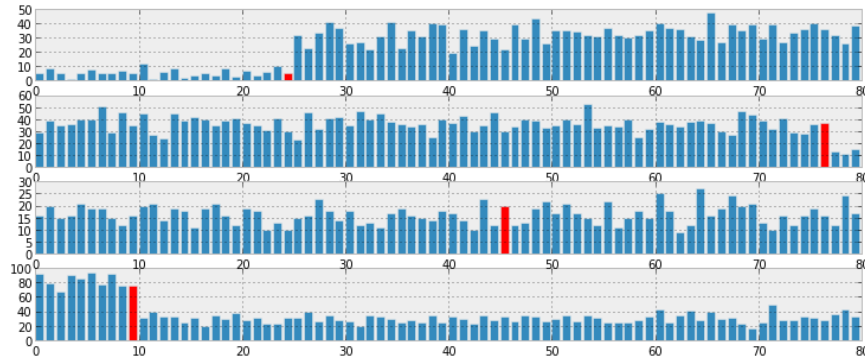


It is okay that our fictional dataset does not look like our observed dataset: the probability is incredibly small it indeed would. PyMC's engine is designed to find good parameters,  $\lambda_i, \tau$ , that maximize this probability.

The ability to generate artificial dataset is an interesting side effect of our modeling, and we will see that this ability is a very important method of Bayesian inference. We produce a few more datasets below:

```
In [29]: def plot_artificial_sms_dataset():
tau = mc.rdiscrete_uniform(0, 80)
alpha = 1./20.
lambda_1, lambda_2 = mc.rexponential( alpha, 2 )
data = np.r_[ mc.rpoisson( lambda_1, tau ), mc.rpoisson( lambda_2, 80-
plt.bar( np.arange( 80 ), data, color = "#348ABD" )
plt.bar( tau-1, data[tau-1], color = "r", label = "user behaviour char
plt.xlim( 0, 80 );

figsize( 12.5, 5)
plt.title( "More example of artificial datasets" )
for i in range(4):
plt.subplot( 4, 1, i)
plot_artificial_sms_dataset()
```



Later we will see how we use this to make predictions and test the appropriateness of our models.

### 0.4.3 An algorithm for human deceit

Likely the most common statistical task is estimating the frequency of events. However, there is a difference between the *observed frequency* and the *true frequency* of an event. The true frequency can be interpreted as the probability of an event occurring. For example, the true frequency of rolling a 1 on a 6-sided die is 0.166. Knowing the frequency of events like baseball home runs, frequency of social attributes, fraction of internet users with cats etc. are common requests we ask of Nature. Unfortunately, in general Nature hides the true frequency from us and we must *infer* it from observed data.

The *observed frequency* is then the frequency we observe: say rolling the die 100 times you may observe 20 rolls of 1. The observed frequency, 0.2, differs from the true frequency, 0.166. We can use Bayesian statistics to infer probable values of the true frequency using an appropriate prior and observed data.

Social data is really interesting as people are not always honest with responses, which adds a further complication into inference. For example, simply asking individuals “Have you ever cheated on a test?” will surely contain some rate of dishonesty. What you can say for certain is that the true rate is less than your observed rate (assuming individuals lie *only about not cheating*; I cannot imagine one who would admit “Yes” to cheating when in fact they hadn’t cheated).

To present an elegant solution to circumventing this dishonesty problem, and to demonstrate Bayesian modeling, we first need to introduce the binomial distribution.

#### The Binomial Distribution

The binomial distribution is one of the most popular distributions, mostly because of its simplicity and usefulness. Unlike the other distributions we have encountered thus far in the book, the binomial distribution has 2 parameters:  $N$ , a positive integer representing  $N$  trials or number of instances of potential events, and  $p$ , the probability of an event occurring in a single trial. Like the Poisson distribution, it is a discrete distribution, but unlike the Poisson distribution, it only weighs integers from 0 to  $N$ . The mass distribution looks like:

$$P(X = k) = \binom{N}{k} p^k (1 - p)^{N-k}$$

If  $X$  is a binomial random variable with parameters  $p$  and  $N$ , denoted  $X \sim \text{Bin}(N, p)$ , then  $X$  is the number of events that occurred in the  $N$  trials (obviously  $0 \leq X \leq N$ ). The larger  $p$  is (while still remaining between 0 and 1), the more events are likely to occur. The expected value of a binomial is equal to  $Np$ . Below we plot the mass probability distribution for varying parameters.

```
In [22]: figsize( 12.5, 4)
import scipy.stats as stats
```

```

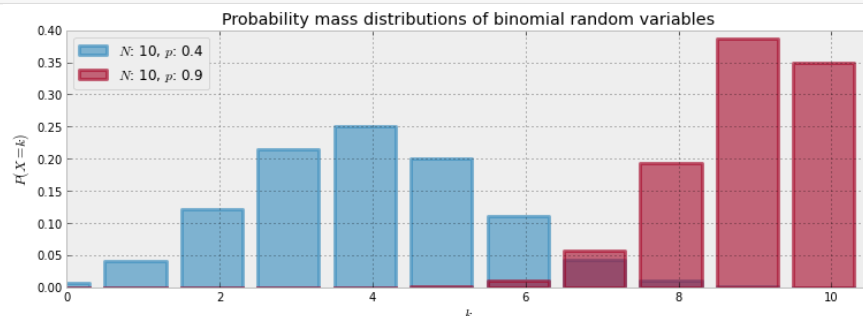
binomial = stats.binom

parameters = [ (10, .4) , (10, .9) ]
colors = ["#348ABD", "#A60628"]

for i in range(2):
    N, p = parameters[i]
    _x = np.arange( N+1 )
    plt.bar( _x - 0.5, binomial.pmf( _x, N, p ), color = colors[i],
            edgecolor=colors[i],
            alpha = 0.6,
            label = "$N$: %d, $p$: %.1f"%(N,p),
            linewidth=3)

plt.legend(loc="upper left")
plt.xlim(0, 10.5)
plt.xlabel("$k$")
plt.ylabel("$P(X = k)$")
plt.title("Probability mass distributions of binomial random variables");

```



The special case when  $N = 1$  corresponds to the Bernoulli distribution. If  $X \sim \text{Ber}(p)$ , then  $X$  is 1 with probability  $p$  and 0 with probability  $1 - p$ . The Bernoulli distribution is useful for indicators, e.g.  $Y = X\alpha + (1 - X)\beta$  is  $\alpha$  with probability  $p$  and  $\beta$  with probability  $1 - p$ .

There is another connection between Bernoulli and Binomial random variables. If we have  $X_1, X_2, \dots, X_N$  Bernoulli random variables with the same  $p$ , then  $Z = X_1 + X_2 + \dots + X_N \sim \text{Binomial}(N, p)$ .

The expected value of a Bernoulli random variable is  $p$ . This can be seen by noting the more general Binomial random variable has expected value  $Np$  and setting  $N = 1$ .

**Example: Cheating among students** We will use the binomial distribution to determine the frequency of students cheating during an exam. If we let  $N$  be the total number of students who took the exam, and assuming each student is interviewed post-exam (answering without consequence), we will receive integer  $X$  “Yes I did cheat” answers. We then find the posterior distribution of  $p$ , given  $N$ , some specified prior on  $p$ , and observed data  $X$ .

This is a completely absurd model. No student, even with a free-pass against punishment, would admit to cheating. What we need is a better *algorithm* to ask students if they had cheated. Ideally the algorithm should encourage individuals to be honest while preserving privacy. The following proposed algorithm is a solution I greatly admire for its ingenuity and effectiveness:

In the interview process for each student, the student flips a coin, hidden from the interviewer. The student agrees to answer honestly if the coin comes up heads. Otherwise, if the coin comes up tails, the student (secretly) flips the coin again, and answers “Yes, I did cheat” if the coin flip lands heads, and “No, I did not cheat”, if the coin flip lands tails. This way, the interviewer does not know if a “Yes” was the result of a guilty plea, or a Heads on a second coin toss. Thus privacy is preserved and the researchers receive honest answers.

I call this the Privacy Algorithm. One could of course argue that the interviewers are still receiving false data since some “Yes”s are not confessions but instead randomness, but an alternative perspective is that the researchers are

discarding approximately half of their original dataset since half of the responses will be noise. But they have gained a systematic data generation process that can be modeled. Furthermore, they do not have to incorporate (perhaps somewhat naively) the possibility of deceitful answers. We can use PyMC to dig through this noisy model, and find a posterior distribution for the true frequency of liars. Suppose 100 students are being surveyed for cheating, and we wish to find  $p$ , the proportion of cheaters. There are a few ways we can model this in PyMC. I'll demonstrate the most explicit way, and later show a simplified version. Both versions arrive at the same inference. In our data-generation model, we sample  $p$ , the true proportion of cheaters, from a prior. Since we are quite ignorant about  $p$ , we will assign it a  $\text{Uniform}(0, 1)$  prior.

```
In [5]: import pymc as mc
        N = 100
        p = mc.Uniform( "freq_cheating", 0, 1)
```

Again, thinking of our data-generation model, we assign Bernoulli random variables to the 100 students: 1 implies they cheated and 0 implies they did not.

```
In [6]: true_answers = mc.Bernoulli( "truths", p, size = N)
```

If we carry out the algorithm, the next step that occurs is the first coin-flip each student makes. This can be modeled again by sampling 100 Bernoulli random variables with  $p = 1/2$ : denote a 1 as a *Heads* and 0 as a *Tails*.

```
In [7]: first_coin_flips = mc.Bernoulli( "first_flips", 0.5, size = N)
        print first_coin_flips.value
```

```
[False False False  True False  True False  True  True False  True  True
  True False False  True  True  True  True False False  True  True False
  True  True  True False False  True  True  True  True False  True  True
  False False False False False False  True  True False  True  True False
  True False  True  True  True  True  True False False  True  True False
  True False  True False  True False False  True  True False  True False
  True  True False  True  True  True False False  True  True False  True
  False  True False  True  True False False False  True False  True False
  False False  True False]
```

Although *not everyone* flips a second time, we can still model the possible realization of second coin-flips:

```
In [8]: second_coin_flips = mc.Bernoulli("second_flips", 0.5, size = N)
```

Using these variables, we can return a possible realization of the *observed proportion* of “Yes” responses. We do this using a PyMC deterministic variable:

```
In [9]: @mc.deterministic
        def observed_proportion( t_a = true_answers,
                                fc = first_coin_flips,
                                sc = second_coin_flips ):
            observed = fc*t_a + (1-fc)*sc
            return observed.sum()/float(N)
```

The line  $fc*t_a + (1-fc)*sc$  contains the heart of the Privacy algorithm. Elements in this array are 1 *if and only if* i) the first toss is heads and the student cheated or ii) the first toss is tails, and the second is heads, and are 0 else. Summing this vector and dividing by  $\text{float}(N)$  produces a proportion.

```
In [10]: observed_proportion.value
```

```
0.54000000000000004
```

```
Out [10]:
```

Next we need a dataset. After performing our coin-flipped interviews the researchers received 35 “Yes” responses. To put this into a relative perspective, if there truly were no cheaters, we should expect to see on average 1/4 of all responses being a “Yes” (half chance of having first coin land Tails, and another half chance of having second coin land Heads), so about 25 responses in a cheat-free world. On the other hand, if *all students cheated*, we should expect to see on approximately 3/4 of all response be “Yes”.

The researchers observe a Binomial random variable, with  $N = 100$  and  $p = \text{observed\_proportion}$  with value = 35:

```
In [13]: X = 35
```

```
observations = mc.Binomial("obs", N, observed_proportion, observed = True,
```

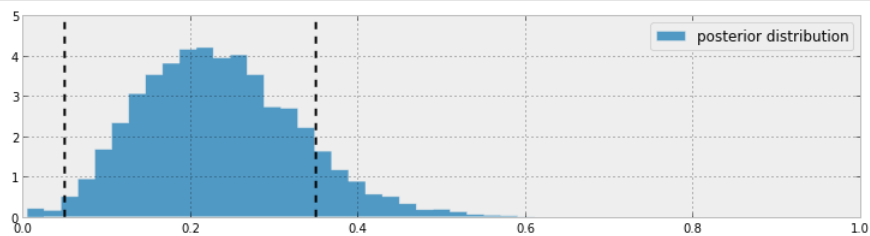
Below we add all the variables of interest to a Model container and run our black-box algorithm over the model.

```
In [14]: model = mc.Model( [p, true_answers, first_coin_flips,
                           second_coin_flips, observed_proportion, observations] )

### To be explained in Chapter 3!
mcmc = mc.MCMC( model )
mcmc.sample( 120000, 80000, 4 )
```

```
[*****100%*****] 120000 of 120000 complete
```

```
In [17]: figsize(12.5, 3 )
p_trace = mcmc.trace("freq_cheating")[:]
plt.hist( p_trace, histtype="stepfilled" , normed = True,
         alpha = 0.85, bins = 30, label = "posterior distribution",
         color = "#348ABD" )
plt.vlines( [.05, .35], [0,0], [5,5], linestyles = "--" )
plt.xlim(0,1)
plt.legend();
```



With regards to the above plot, we are still pretty uncertain about what the true frequency of cheaters might be, but we have narrowed it down to a range between 0.05 to 0.35 (marked by the dashed lines). This is pretty good, as *a priori* we had no idea how many students might have cheated (hence the uniform distribution for our prior). On the other hand, it is also pretty bad since there is a .3 length window the true value most likely lives in. Have we even gained anything, or are we still too uncertain about the true frequency?

I would argue, yes, we have discovered something. It is implausible, according to our posterior, that there are *no cheaters*, i.e. the posterior assigns low probability to  $p = 0$ . Since we started with an uniform prior, treating all values of  $p$  as equally plausible, but the data ruled out  $p = 0$  as a possibility, we can be confident that there were cheaters.

This kind of algorithm can be used to gather private information from users and be *reasonably* confident that the data, though noisy, is truthful.

## Alternative PyMC Model

Given a value for  $p$  (which from our god-like position we know), we can find the probability the student will answer yes:

$$P(\text{"Yes"}) = P(\text{Heads on first coin})P(\text{cheater}) + P(\text{Tails on first coin})P(\text{Heads on second coin}) \quad (17)$$

$$= \frac{1}{2}p + \frac{1}{2}\frac{1}{2} \quad (18)$$

$$= \frac{p}{2} + \frac{1}{4} \quad (20)$$

$$= \frac{p}{2} + \frac{1}{4} \quad (21)$$

Thus, knowing  $p$  we know the probability a student will respond “Yes”. In PyMC, we can create a deterministic function to evaluate the probability of responding “Yes”, given  $p$ :

```
In [18]: p = mc.Uniform( "freq_cheating", 0, 1)

@mc.deterministic
def p_skewed( p = p ):
    return 0.5*p + 0.25
```

I could have typed `p_skewed = 0.5*p + 0.25` instead for a one-liner, as the elementary operations of addition and scalar multiplication will implicitly create a deterministic variable, but I wanted to make the deterministic boilerplate explicit for clarity’s sake.

If we know the probability of respondents saying “Yes”, which is `p_skewed`, and we have  $N = 100$  students, the number of “Yes” responses is a binomial random variable with parameters  $N$  and `p_skewed`.

This is where we include our observed 35 “Yes” responses. In the declaration of the `mc.Binomial`, we include `value = 35` and `observed = True`.

```
In [19]: yes_responses = mc.Binomial( "number_cheaters", 100, p_skewed,
                                     value = 35, observed = True )
```

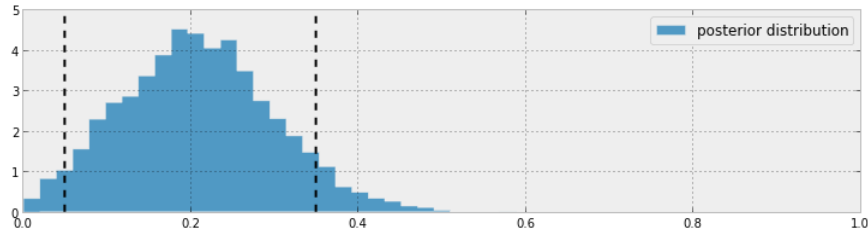
Below we add all the variables of interest to a `Model` container and run our black-box algorithm over the model.

```
In [22]: model = mc.Model( [yes_responses, p_skewed, p ] )

### To Be Explained in Chapter 3!
mcmc = mc.MCMC(model)
mcmc.sample( 12500, 2500 )
```

```
[*****100%*****] 12500 of 12500 complete
```

```
In [23]: figsize(12.5, 3 )
p_trace = mcmc.trace("freq_cheating")[:]
plt.hist( p_trace, histtype="stepfilled" , normed = True,
          alpha = 0.85, bins = 30, label = "posterior distribution",
          color = "#348ABD" )
plt.vlines( [.05, .35], [0,0], [5,5], linestyle = "--" )
plt.xlim(0,1)
plt.legend();
```



## More PyMC Tricks

### Protip: *Lighter* deterministic variables with `Lambda` class

Sometimes writing a deterministic function using the `@mc.deterministic` decorator can seem like a chore, especially for a small function. I have already mentioned that elementary math operations *can* produce deterministic variables implicitly, but what about operations like indexing or slicing? Built-in `Lambda` functions can handle this with the elegance and simplicity required. For example,

```
beta = mc.Normal( "coefficients", 0, size=(N,1) )
x = np.random.randn( (N,1) )
linear_combination = mc.Lambda( lambda x=x, beta = beta: np.dot( x.T, beta ) )
```

### Protip: Arrays of PyMC variables

There is no reason why we cannot store multiple heterogeneous PyMC variables in a Numpy array. Just remember to set the `dtype` of the array to `object` upon initialization. For example:

```
In [16]: N = 10
x = np.empty( N , dtype=object )
for i in range(0, N):
    x[i] = mc.Exponential('x_%i' % i, (i+1)**2)
```

The remainder of this chapter examines some practical examples of PyMC and PyMC modeling:

**Example: Challenger Space Shuttle Disaster** On January 28, 1986, the twenty-fifth flight of the U.S. space shuttle program ended in disaster when one of the rocket boosters of the Shuttle Challenger exploded shortly after lift-off, killing all seven crew members. The presidential commission on the accident concluded that it was caused by the failure of an O-ring in a field joint on the rocket booster, and that this failure was due to a faulty design that made the O-ring unacceptably sensitive to a number of factors including outside temperature. Of the previous 24 flights, data were available on failures of O-rings on 23, (one was lost at sea), and these data were discussed on the evening preceding the Challenger launch, but unfortunately only the data corresponding to the 7 flights on which there was a damage incident were considered important and these were thought to show no obvious trend. The data are shown below (see [1]):

```
In [24]: figsize( 12.5, 3.5 )
np.set_printoptions(precision=3, suppress=True)
challenger_data = np.genfromtxt("data/challenger_data.csv", skip_header =
                               usecols=[1,2], missing_values="NA", delimit

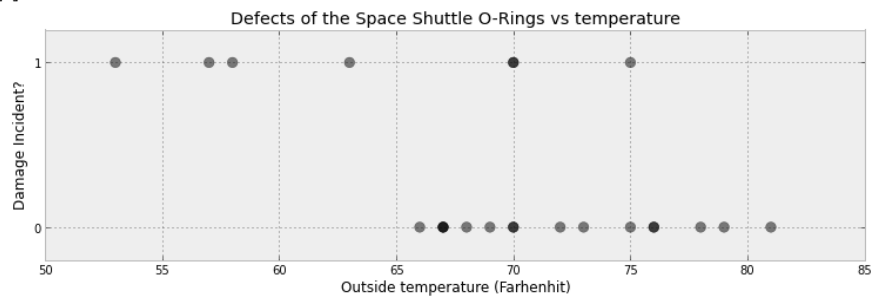
#drop the NA values
challenger_data = challenger_data[ ~np.isnan(challenger_data[:,1]) ]

#plot it, as a function of tempature (the first column)
print "Temp (F), O-Ring failure?"
print challenger_data
```

```
plt.scatter( challenger_data[:,0], challenger_data[:,1], s = 75, color="k",
            alpha = 0.5)
plt.yticks([0,1])
plt.ylabel("Damage Incident?")
plt.xlabel("Outside temperature (Fahrenheit)" )
plt.title("Defects of the Space Shuttle O-Rings vs temperature");
```

Temp (F), O-Ring failure?

```
[[ 66.  0.]
 [ 70.  1.]
 [ 69.  0.]
 [ 68.  0.]
 [ 67.  0.]
 [ 72.  0.]
 [ 73.  0.]
 [ 70.  0.]
 [ 57.  1.]
 [ 63.  1.]
 [ 70.  1.]
 [ 78.  0.]
 [ 67.  0.]
 [ 53.  1.]
 [ 67.  0.]
 [ 75.  0.]
 [ 70.  0.]
 [ 81.  0.]
 [ 76.  0.]
 [ 79.  0.]
 [ 75.  1.]
 [ 76.  0.]
 [ 58.  1.]]
```



It looks clear that *the probability* of damage incidents occurring increases as the outside temperature decreases. We are interested in modeling the probability here because it does not look like there is a strict cutoff point between temperature and a damage incident occurring. The best we can do is ask “At temperature  $t$ , what is the probability of a damage incident?”. The goal of this example is to answer that question.

We need a function of temperature, call it  $p(t)$ , that is bounded between 0 and 1 (so as to model a probability) and changes from 1 to 0 as we increase temperature. There are actually many such functions, but the most popular choice is the *logistic function*.

$$p(t) = \frac{1}{1 + e^{\beta t}}$$

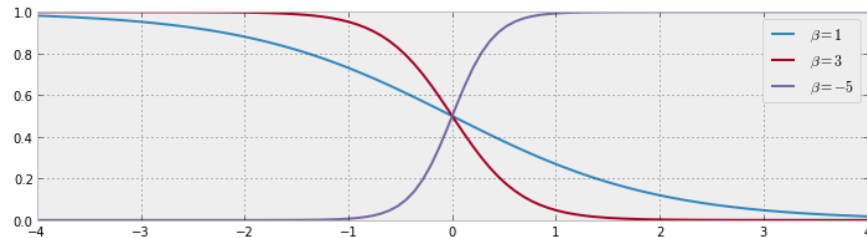
In this model,  $\beta$  is the variable we are uncertain about. Below is the function plotted for  $\beta = 1, 3, -5$ .



```
In [25]: figsize(12,3)

def logistic( x, beta):
    return 1.0/( 1.0 + np.exp( beta*x) )

x = np.linspace( -4, 4, 100 )
plt.plot(x, logistic( x, 1), label = r"$\beta = 1$")
plt.plot(x, logistic( x, 3), label = r"$\beta = 3$")
plt.plot(x, logistic( x, -5), label = r"$\beta = -5$")
plt.legend();
```



But something is missing. In the plot of the logistic function, the probability changes only near zero, but in our data above the probability changes around 65 to 70. We need to add a *bias* term to our logistic function:

$$p(t) = \frac{1}{1 + e^{\beta t + \alpha}}$$

Some plots are below, with differing  $\alpha$ .

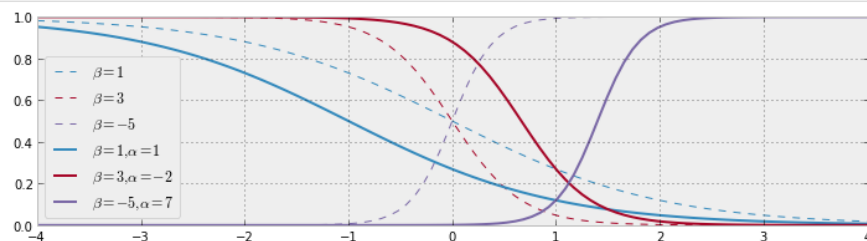
```
In [27]: def logistic( x, beta, alpha=0):
    return 1.0/( 1.0 + np.exp( np.dot( beta, x) + alpha) )

x = np.linspace( -4, 4, 100 )

plt.plot(x, logistic( x, 1), label = r"$\beta = 1$", ls= "--", lw =1 )
plt.plot(x, logistic( x, 3), label = r"$\beta = 3$", ls= "--", lw =1)
plt.plot(x, logistic( x, -5), label = r"$\beta = -5$", ls= "--", lw =1)

plt.plot(x, logistic( x, 1,1), label = r"$\beta = 1, \alpha = 1$",
    color = "#348ABD")
plt.plot(x, logistic( x, 3, -2), label = r"$\beta = 3, \alpha = -2$",
    color="#A60628")
plt.plot(x, logistic( x, -5, 7), label = r"$\beta = -5, \alpha = 7$",
    color = "#7A68A6")

plt.legend(loc = "lower left");
```



Adding a constant term  $\alpha$  amounts to shifting the curve left or right (hence why it is called a *bias*.)

Let's start modeling this in PyMC. The  $\beta, \alpha$  parameters have no reason to be positive, bounded or relatively large, so they are best modeled by a *Normal random variable*, introduced next.

## Normal distributions

A Normal random variable, denoted  $X \sim N(\mu, 1/\tau)$ , has a distribution with two parameters: the mean,  $\mu$ , and the *precision*,  $\tau$ . Those familiar with the Normal distribution already have probably seen  $\sigma^2$  instead of  $\tau$ . They are in fact reciprocals of each other. The change was motivated by simpler mathematical analysis and is an artifact of older Bayesian methods. Just remember: The smaller  $\tau$ , the larger the spread of the distribution (i.e. we are more uncertain); the larger  $\tau$ , the tighter the distribution (i.e. we are more certain). Regardless,  $\tau$  is always positive.

The probability density function of a  $N(\mu, 1/\tau)$  random variable is:

$$f(x|\mu, \tau) = \sqrt{\frac{\tau}{2\pi}} \exp\left(-\frac{\tau}{2}(x - \mu)^2\right)$$

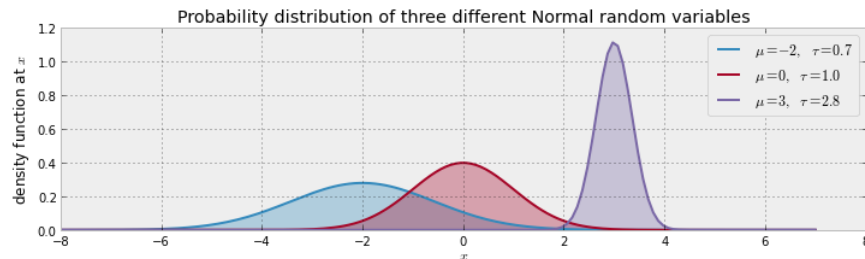
We plot some different density functions below.

```
In [28]: import scipy.stats as stats

nor = stats.norm
x = np.linspace( -8, 7, 150 )
mu = (-2, 0, 3)
tau = (.7, 1, 2.8 )
colors = ["#348ABD", "#A60628", "#7A68A6"]
parameters = zip( mu, tau, colors )

for _mu, _tau, _color in parameters:
    plt.plot( x, nor.pdf( x, _mu , scale = 1./_tau ), \
              label = "$\mu = %d, \tau = %.1f$" % (_mu, _tau), color = _color )
    plt.fill_between( x, nor.pdf( x, _mu, scale = 1./_tau ), color = _color,
                    alpha = .33)

plt.legend(loc = "upper right")
plt.xlabel("$x$")
plt.ylabel("density function at $x$")
plt.title("Probability distribution of three different Normal random variables")
```



A Normal random variable can take on any real number, but the variable is very likely to be relatively close to  $\mu$ . In fact, the expected value of a Normal is equal to its  $\mu$  parameter:

$$E[X|\mu, \tau] = \mu$$

and its variance is equal to the inverse of  $\tau$ :

$$Var(X|\mu, \tau) = \frac{1}{\tau}$$

Below we continue our modeling of the Challenger space craft:

```
In [29]: import pymc as mc

temperature = challenger_data[:,0]
D = challenger_data[:,1] #defect or not?

#notice the 'value' here. We explain why below.
beta = mc.Normal( "beta", 0, 0.001, value = 0 )
alpha = mc.Normal( "alpha", 0, 0.001, value = 0 )

@mc.deterministic
def p( t = temperature, alpha = alpha, beta = beta):
    return 1.0/( 1. + np.exp( beta*t + alpha) )
```

We have our probabilities, but how do we connect them to our observed data? A *Bernoulli* random variable with parameter  $p$ , denoted  $\text{Ber}(p)$ , is a random variable that takes value 1 with probability  $p$ , and 0 else. Thus, our model can look like:

$$\text{Defect Incident, } D_i \sim \text{Ber}( p(t_i) ), \quad i = 1..N$$

where  $p(t)$  is our logistic function and  $t_i$  are the temperatures we have observations about. Notice in the above code we had to set the values of  $\beta$  and  $\alpha$  to 0. The reason for this is that if  $\beta$  and  $\alpha$  are very large, they make  $p$  equal to 1 or 0. Unfortunately, `mc.Bernoulli` does not like probabilities of exactly 0 or 1, though they are mathematically well-defined probabilities. So by setting the coefficient values to 0, we set the variable  $p$  to be a reasonable starting value. This has no effect on our results, nor does it mean we are including any additional information in our prior. It is simply a computational caveat in PyMC.

```
In [30]: p.value

array([[ 0.5,  0.5,  0.5,  0.5,  0.5,  0.5,  0.5,  0.5,  0.5,  0.5,  0.5,
         0.5,  0.5,  0.5,  0.5,  0.5,  0.5,  0.5,  0.5,  0.5,  0.5,
         0.5]])
```

```
In [48]: # connect the probabilities in 'p' with our observations through a
# Bernoulli random variable.
observed = mc.Bernoulli( "bernoulli_obs", p, value = D, observed=True)

model = mc.Model( [observed, beta, alpha] )

#mysterious code to be explained in Chapter 3
map_ = mc.MAP(model)
map_.fit()
mcmc = mc.MCMC( model )
mcmc.sample( 120000, 100000, 2 )
```

```
[*****100%*****] 120000 of 120000 complete
```

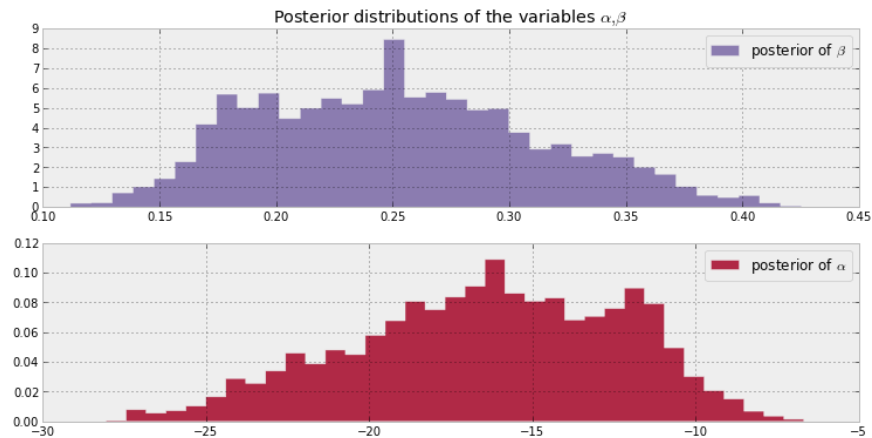
We have trained our model on the observed data, now we can sample values from the posterior. Let's look at the posterior distributions for  $\alpha$  and  $\beta$ :

```
In [51]: alpha_samples = mcmc.trace( 'alpha' )[:, None] #best to make them 1d
beta_samples = mcmc.trace( 'beta' )[:, None]

figsize(12.5, 6)

#histogram of the samples:
plt.subplot(211)
plt.title(r"Posterior distributions of the variables $\alpha, \beta$")
plt.hist( beta_samples, histtype='stepfilled', bins = 35, alpha = 0.85, \
         label = r"posterior of $\beta$", color = "#7A68A6",normed = True )
plt.legend()
```

```
plt.subplot(212)
plt.hist( alpha_samples, histtype='stepfilled', bins = 35, alpha = 0.85, \
        label = r"posterior of  $\alpha$ ", color = "#A60628",normed = True)
plt.legend();
```



Regarding the spread of the data, we are very uncertain about what the true parameters might be (though considering the low sample size and the large overlap of defects-to-nondefects this behaviour is perhaps expected).

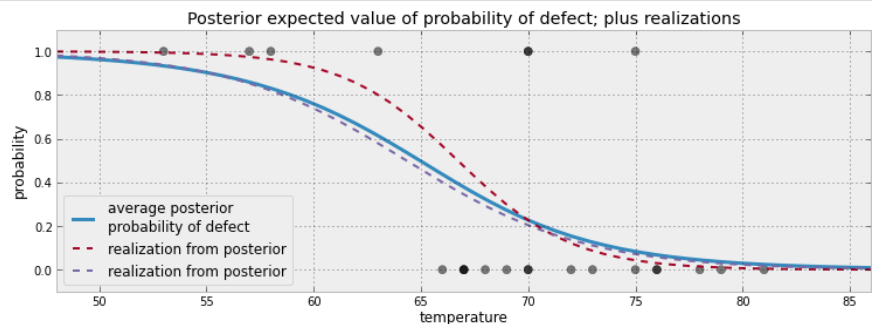
Next, let's look at the *expected probability* for a specific value of the temperature. That is, we average over all samples from the posterior to get a likely value for  $p(t_i)$ .

```
In [52]: t = np.linspace( temperature.min() - 5, temperature.max()+5, 50 )[:,None]
p_t= logistic( t.T, beta_samples, alpha_samples )

mean_prob_t = p_t.mean(axis=0)
```

```
In [54]: figsize( 12.5, 4)

plt.plot( t, mean_prob_t, lw = 3, label = "average posterior \nprobability
of defect")
plt.plot( t, p_t[0, :], ls="--",label="realization from posterior" )
plt.plot( t, p_t[-2, :], ls="--", label="realization from posterior" )
plt.scatter( temperature, D, color = "k", s = 50, alpha = 0.5 )
plt.title("Posterior expected value of probability of defect; plus realiz
plt.legend(loc= "lower left")
plt.ylim( -0.1, 1.1 )
plt.xlim( t.min(), t.max() )
plt.ylabel("probability")
plt.xlabel("temperature");
```



Above we also plotted two possible realizations of what the actual underlying system might be. Both are equally likely as any other draw. The blue line is what occurs when we average all the 20000 possible dotted lines together.

An interesting question to ask is for what temperatures are we most uncertain about the defect-probability? Below we plot the expected value line **and** the associated 95% intervals for each temperature.

```
In [55]: from scipy.stats.mstats import mquantiles

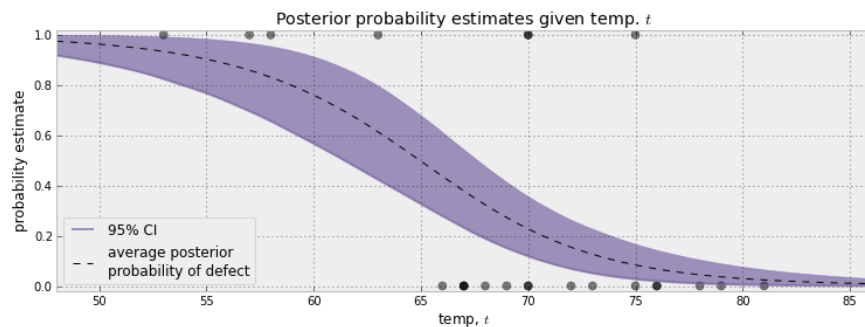
# vectorized bottom and top 5% quantiles for "confidence interval"
qs = mquantiles(p_t, [0.05, 0.95], axis=0)
plt.fill_between(t[:,0], *qs, alpha = 0.7,
                 color = "#7A68A6")

plt.plot(t[:,0], qs[0], label="95% CI", color = "#7A68A6", alpha = 0.7)

plt.plot( t, mean_prob_t, lw = 1, ls= "--", color = "k",
          label = "average posterior \nprobability of defect")

plt.xlim( t.min(), t.max() )
plt.ylim( -0.02, 1.02 )
plt.legend(loc="lower left")
plt.scatter( temperature, D, color = "k", s = 50, alpha = 0.5 )
plt.xlabel("temp, $t$")

plt.ylabel("probability estimate")
plt.title("Posterior probability estimates given temp. $t$");
```



The *95% credible interval*, or *95% CI*, painted in purple, represents the interval, for each temperature, that contains 95% of the distribution. For example, at 65 degrees, we can be 95% sure that the probability of defect lies between 0.25 and 0.75.

More generally, we can see that as the temperature nears 60 degrees, the CI's spread out over [0,1] quickly. As we pass 70 degrees, the CI's tighten again. This can give us insight about how to proceed next: we should probably test more O-rings around 60-65 temperature to get a better estimate of probabilities in that range. Similarly, when reporting to scientists your estimates, you should be very cautious about simply telling them the expected probability, as we can see this does not reflect how *wide* the posterior distribution is.

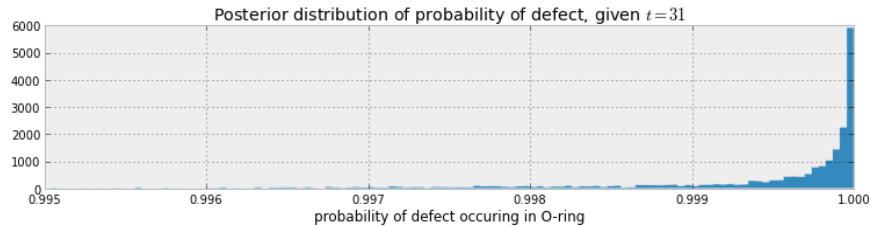
### What about the day of the Challenger disaster?

On the day of the Challenger disaster, the outside temperature was 31 degrees Fahrenheit. What is the posterior distribution of a defect occurring, given this temperature? The distribution is plotted below. It looks almost guaranteed that the Challenger was going to be subject to defective O-rings.

```
In [71]: figsize(12.5, 2.5)

prob_31 = logistic( 31, beta_samples, alpha_samples )

plt.xlim( 0.995, 1)
plt.hist( prob_31, bins = 1000, normed = True, histtype='stepfilled' )
plt.title( "Posterior distribution of probability of defect, given $t = 31$")
plt.xlabel( "probability of defect occurring in O-ring" );
```



## Is our model appropriate?

The skeptical reader will say “You deliberately chose the logistic function for  $p(t)$  and the specific priors. Perhaps other functions or priors will give different results. How do I know I have chosen a good model?” This is absolutely true. To consider an extreme situation, what if I had chosen the function  $p(t) = 1, \forall t$ , which guarantees a defect always occurring: I would have again predicted disaster on January 28th. Yet this is clearly a poorly chosen model. On the other hand, if I did choose the logistic function for  $p(t)$ , but specified all my priors to be very tight around 0, likely we would have very different posterior distributions. How do we know our model is an expression of the data? This encourages us to measure the model’s **goodness of fit**.

We can think: *how can we test whether our model is a bad fit?* An idea is to compare observed data (which if we recall is a *fixed* stochastic variable) with artificial dataset which we can simulate. The rational is that if the simulated dataset does not appear similar, statistically, to the observed dataset, then likely our model is not accurately represented the observed data.

Previously in this Chapter, we simulated artificial dataset for the SMS example. To do this, we sampled values from the priors. We saw how varied the resulting datasets looked like, and rarely did they mimic our observed dataset. In the current example, we should sample from the *posterior* distributions to create *very plausible datasets*. Luckily, our Bayesian framework makes this very easy. We only need to create a new `Stochastic` variable, that is exactly the same as our variable that stored the observations, but minus the observations themselves. If you recall, our `Stochastic` variable that stored our observed data was:

```
observed = mc.Bernoulli( "bernoulli_obs", p, value = D, observed=True)
```

Hence we create:

```
simulated_data = mc.Bernoulli("simulation_data", p )
```

Let’s simulate 10 000:

```
In [57]: simulated = mc.Bernoulli( "bernoulli_sim", p)
         N = 10000

         mcmc = mc.MCMC( [simulated, alpha, beta, observed ] )
         mcmc.sample( N )
```

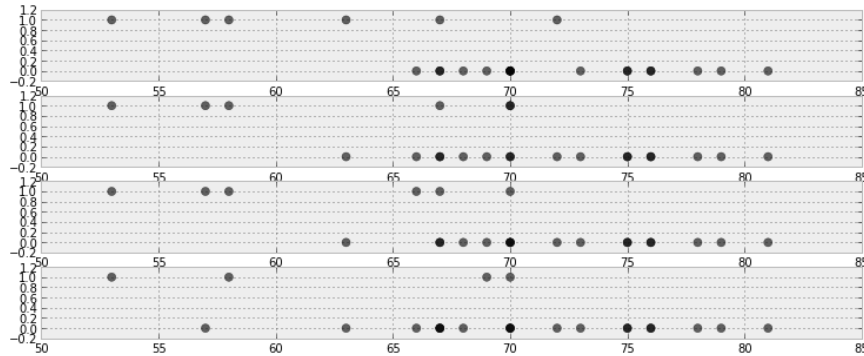
```
[*****100%*****] 10000 of 10000 complete
```

```
In [58]: figsize(12.5, 5)

         simulations = mcmc.trace("bernoulli_sim")[:]
         print simulations.shape

         plt.title( "Simulated dataset using posterior parameters")
         figsize( 12.5, 6)
         for i in range(4):
             ax = subplot( 4, 1, i+1)
             plt.scatter( temperature, simulations[1000*i,:], color = "k",
                         s = 50, alpha = 0.6 );
```

(10000L, 23L)



Note that the above plots are different (if you can think of a cleaner way to present this, please send a pull request and answer [here!](#)).

We wish to assess how good our model is. “Good” is a subjective term of course, so results must be relative to other models.

We will be doing this graphically as well, which may seem like an even less objective method. The alternative is to use *Bayesian p-values*. These are still subjective, as the proper cutoff between good and bad is arbitrary. Gelman emphasises that the graphical tests are more illuminating [7] than p-value tests. We agree.

The following graphical test is a novel data-viz approach to logistic regression. The plots are called *separation plots*[8]. For a suite of models we wish to compare, each model is plotted on an individual separation plot. I leave most of the technical details about separation plots to the very accessible [original paper](#), but I’ll summarize their use here.

For each model, we calculate the proportion of times the posterior simulation proposed a value of 1 for a particular temperature, i.e. compute  $P(\text{Defect} = 1 | t, \alpha, \beta)$  by averaging. This gives us the posterior probability of a defect at each data point in our dataset. For example, for the model we used above:

```
In [65]: posterior_probability = simulations.mean(axis=0)
print "posterior prob of defect | realized defect "
for i in range( len(D) ):
    print "%.2f                | %d"%(posterior_probability[i], D[i])
```

posterior prob of defect		realized defect
0.44		0
0.26		1
0.29		0
0.35		0
0.39		0
0.19		0
0.16		0
0.26		0
0.82		1
0.59		1
0.26		1
0.07		0
0.39		0
0.91		1
0.39		0
0.12		0
0.26		0
0.04		0
0.10		0
0.06		0

0.12		1
0.10		0
0.79		1

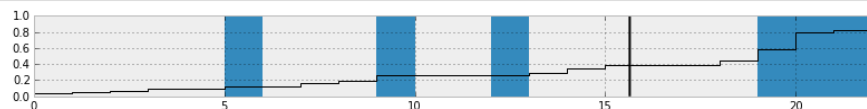
Next we sort each column by the posterior probabilities:

```
In [66]: ix = np.argsort( posterior_probability )
print "probb | defect "
for i in range( len(D) ):
    print "%.2f | %d"%(posterior_probability[ix[i]], D[ix[i]])
```

probb		defect
0.04		0
0.06		0
0.07		0
0.10		0
0.10		0
0.12		1
0.12		0
0.16		0
0.19		0
0.26		1
0.26		0
0.26		0
0.26		1
0.29		0
0.35		0
0.39		0
0.39		0
0.39		0
0.44		0
0.59		1
0.79		1
0.82		1
0.91		1

We can present the above data better in a figure: I've wrapped this up into a `separation_plot` function.

```
In [80]: from separation_plot import separation_plot
figsize( 11., 1.5 )
separation_plot( posterior_probability, D )
```



The snaking-line is the sorted probabilities, blue bars denote defects, and empty space (or grey bars for the optimistic readers) denote non-defects. As the probability rises, we see more and more defects occur. On the right hand side, the plot suggests that as the posterior probability is large (line close to 1), then more defects are realized. This is good behaviour. Ideally, all the blue bars *should* be close to the right-hand side, and deviations from this reflect missed predictions.

The black vertical line is the expected number of defects we should observe, given this model. This allows the user to see how the total number of events predicted by the model compares to the actual number of events in the data.



It is much more informative to compare this to separation plots for other models. Below we compare our model (top) versus three others:

1. the perfect model, which predicts the posterior probability to be equal 1 if a defect did occur.
2. a completely random model, which predicts random probabilities regardless of temperature.
3. a constant model: where  $P(D = 1 | t) = c, \forall t$ . The best choice for  $c$  is the observed frequency of defects, in this case 7/23.

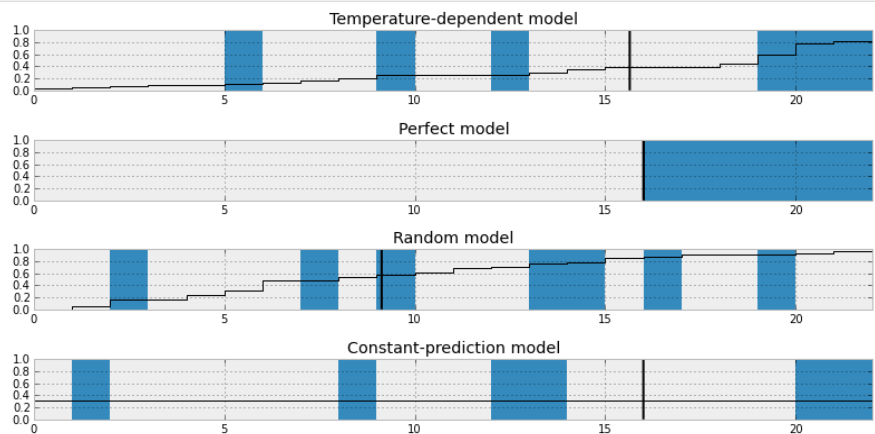
```
In [91]: figsize( 11., 1.25 )

# our temperature-dependent model
separation_plot( posterior_probability, D )
plt.title("Temperature-dependent model")

# perfect model
# i.e. the probability of defect is equal to if a defect occurred or not.
p = D
separation_plot( p, D )
plt.title("Perfect model")

# random predictions
p = np.random.rand( 23 )
separation_plot( p, D )
plt.title("Random model")

# constant model
constant_prob = 7./23*np.ones(23)
separation_plot( constant_prob, D )
plt.title("Constant-prediction model");
```



In the random model, we can see that as the probability increases there is no clustering of defects to the right-hand side. Similarly for the constant model.

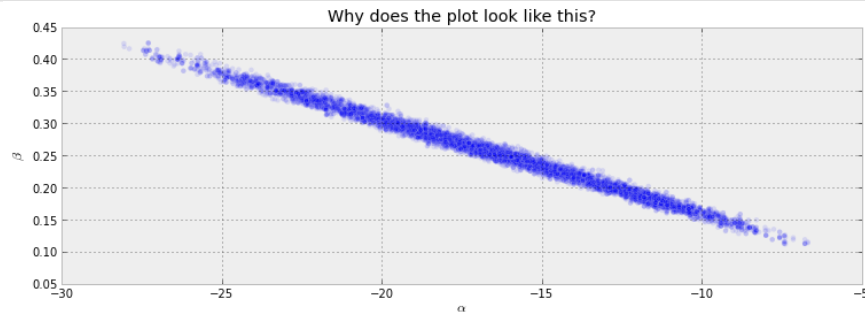
The perfect model, the probability line is not well shown, as it is stuck to the bottom and top of the figure. Of course the perfect model is only for demonstration, and we cannot infer any scientific inference from it.

**Exercises** 1. Try putting in extreme values for our observations in the cheating example. What happens if we observe 25 affirmative responses? 10? 50? 2. Try plotting  $\alpha$  samples versus  $\beta$  samples. Why might the resulting plot look like this?

```
In [90]: #type your code here.
figsize(12.5, 4 )

plt.scatter( alpha_samples, beta_samples, alpha = 0.1 )
```

```
plt.title( "Why does the plot look like this?" )
plt.xlabel( r"$\alpha$" )
plt.ylabel( r"$\beta$" );
```



## References

- [1] Dalal, Fowlkes and Hoadley (1989), JASA, 84, 945-957.
- [2] German Rodriguez. Datasets. In WWS509. Retrieved 30/01/2013, from <http://data.princeton.edu/wws509/datasets/#smoking>.
- [3] McLeish, Don, and Cynthia Struthers. STATISTICS 450/850 Estimation and Hypothesis Testing. Winter 2012. Waterloo, Ontario: 2012. Print.
- [4] Fonnesbeck, Christopher. "Building Models." PyMC-Devs. N.p., n.d. Web. 26 Feb 2013. <http://pymc-devs.github.com/pymc/modelbuilding.html>.
- [5] Cronin, Beau. "Why Probabilistic Programming Matters." 24 Mar 2013. Google, Online Posting to Google+. Web. 24 Mar. 2013. <https://plus.google.com/u/0/107971134877020469960/posts/KpeRdJKR6Z1>.
- [6] S.P. Brooks, E.A. Catchpole, and B.J.T. Morgan. Bayesian animal survival estimation. Statistical Science, 15: 357-376, 2000
- [7] Gelman, Andrew. "Philosophy and the practice of Bayesian statistics." British Journal of Mathematical and Statistical Psychology. (2012): n. page. Web. 2 Apr. 2013.
- [8] Greenhill, Brian, Michael D. Ward, and Audrey Sacks. "The Separation Plot: A New Visual Method for Evaluating the Fit of Binary Models." American Journal of Political Science. 55.No.4 (2011): n. page. Web. 2 Apr. 2013.

```
In [1]: from IPython.core.display import HTML
def css_styling():
    styles = open("../styles/custom.css", "r").read()
    return HTML(styles)
css_styling()
```

Out[1]:

## 0.5 Chapter 3

---

## 0.5.1 Opening the black box of MCMC

The previous two chapters hid the inner-mechanics of PyMC, and more generally Markov Chain Monte Carlo (MCMC), from the reader. The reason for including this chapter is three-fold. The first is that any book on Bayesian inference must discuss MCMC. I cannot fight this. Blame the statisticians. Secondly, knowing the process of MCMC gives you insight into whether your algorithm has converged. (Converged to what? We will get to that) Thirdly, we'll understand *why* we are returned thousands of samples from the posterior as a solution, which at first thought can be odd.

### The Bayesian landscape

When we setup a Bayesian inference problem with  $N$  unknowns, we are implicitly creating an  $N$  dimensional space for the prior distributions to exist in. Associated with the space is an additional dimension, which we can describe as the *surface*, or *curve*, that sits on top of the space, that reflects the *prior probability* of a particular point. The surface on the space is defined by our prior distributions. For example, if we have two unknowns  $p_1$  and  $p_2$ , and priors for both are Uniform(0, 5), the space created is a square of length 5 and the surface is a flat plane that sits on top of the square (representing that every point is equally likely).

```
In [2]: %pylab inline
import scipy.stats as stats
figsize(12.5, 4)

import matplotlib.pyplot as plt
from mpl_toolkits.mplot3d import Axes3D

fig = plt.figure()
x = y = np.linspace(0, 5, 100 )
X, Y = np.meshgrid(x, y)

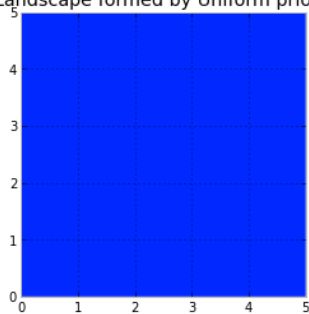
subplot(121)
uni_x = stats.uniform.pdf( x, loc=0, scale = 5)
uni_y = stats.uniform.pdf(x, loc=0, scale = 5)
M = np.dot( uni_x[:,None], uni_y[None,:] )
im = plt.imshow(M, interpolation='none', origin='lower',
               cmap=cm.jet, vmax=1, vmin = -.15, extent=(0, 5, 0, 5))

plt.xlim( 0, 5)
plt.ylim( 0, 5)
plt.title( "Landscape formed by Uniform priors." )

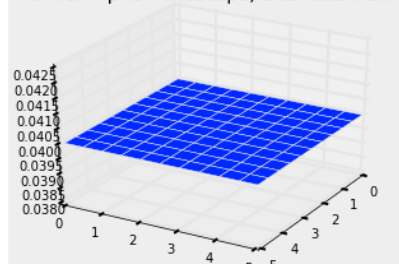
ax = fig.add_subplot(122, projection='3d')
ax.plot_surface(X, Y, M, cmap=cm.jet, vmax=1, vmin = -.15 )
ax.view_init( azimuth = 390)
plt.title( "Uniform prior landscape; alternate view");
```

Welcome to pylab, a matplotlib-based Python environment [backend: module://IPython.zmq.pylab]. For more information, type 'help(pylab)'.

Landscape formed by Uniform priors.



Uniform prior landscape; alternate view



Alternatively, if the two priors are  $\text{Exp}(3)$  and  $\text{Exp}(10)$ , then the space is all positive numbers on the 2-D plane, and the surface induced by the priors looks like a water fall that starts at the point (0,0) and flows over the positive numbers.

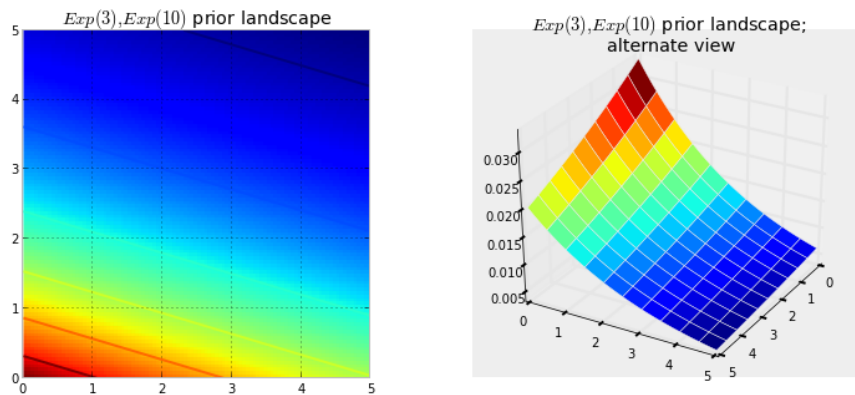
The plots below visualize this. The more dark red the color, the more prior probability is assigned to that location. Conversely, areas with darker blue represent that our priors assign very low probability to that location.

```
In [22]: figsize( 12.5,5 )

fig = plt.figure()
subplot(121)

exp_x = stats.expon.pdf( x, scale = 3)
exp_y = stats.expon.pdf(x, scale = 10)
M = np.dot( exp_x[:,None], exp_y[None,:] )
CS = plt.contour( X, Y, M )
im = plt.imshow(M, interpolation='none', origin='lower',
               cmap=cm.jet, extent=(0,5,0,5))
#plt.xlabel( "prior on $p_1$" )
#plt.ylabel( "prior on $p_2$" )
plt.title( "$\text{Exp}(3), \text{Exp}(10)$ prior landscape" )

ax = fig.add_subplot(122, projection='3d')
ax.plot_surface(X, Y, M, cmap=cm.jet )
ax.view_init( azimuth = 390)
plt.title( "$\text{Exp}(3), \text{Exp}(10)$ prior landscape; \nalternate view")
```



These are simple examples in 2D space, where our brains can understand surfaces well. In practice, spaces and surfaces generated by our priors can be much higher dimensional.

If these surfaces describe our *prior distributions* on the unknowns, what happens to our space after we incorporate our observed data  $X$ ? The data  $X$  does not change the space, but it changes the surface of the space by *pulling and stretching the fabric of the prior surface* to reflect where the true parameters likely live. More data means more pulling and stretching, and our original shape becomes mangled or insignificant compared to the newly formed shape. Less data, and our original shape is more present. Regardless, the resulting surface describes the *posterior distribution*.

Again I must stress that it is, unfortunately, impossible to visualize this in large dimensions. For two dimensions, the data essentially *pushes up* the original surface to make *tall mountains*. The tendency of the observed data to *push up* the posterior probability in certain areas is checked by the prior probability distribution, so that less prior probability means more resistance. Thus in the double-exponential prior case above, a mountain (or multiple mountains) that might erupt near the (0,0) corner would be much higher than mountains that erupt closer to (5,5), since there is more resistance (low prior probability) near (5,5). The peak reflects the posterior probability of where the true parameters are likely to be found. Importantly, if the prior has assigned a probability of 0, then no posterior probability will be assigned there.

Suppose the priors mentioned above represent different parameters  $\lambda$  of two Poisson distributions. We observe a few data points and visualize the new landscape:

```

In [20]: ### create the observed data

#sample size of data we observe
N = 1

#the true parameters, but of course we do not see these values...
lambda_1_true = 1
lambda_2_true = 3

#...we see the data generated, dependent on the above two values.
data = np.concatenate( [
    stats.poisson.rvs( lambda_1_true, size = (N,1)),
    stats.poisson.rvs( lambda_2_true, size = (N,1))
], axis=1 )
print "observed (2-dimensional, sample size = %d):"%N, data

#plotting details.
x = y = np.linspace(.01,5, 100 )
likelihood_x = np.array( [stats.poisson.pmf( data[:,0], _x) for _x in x])
likelihood_y = np.array( [stats.poisson.pmf( data[:,1], _y) for _y in y])
L = np.dot( likelihood_x[:,None], likelihood_y[None,:] )

```

```
observed (2-dimensional, sample size = 1): [[2 3]]
```

```

In [21]: figsize( 12.5,12 )

subplot(221)
uni_x = stats.uniform.pdf( x,loc=0, scale = 5)
uni_y = stats.uniform.pdf(x, loc=0, scale = 5)
M = np.dot( uni_x[:,None], uni_y[None,:] )
im = plt.imshow(M, interpolation='none', origin='lower',
    cmap=cm.jet, vmax=1, vmin = -.15, extent=(0,5,0,5))
plt.scatter( lambda_1_true, lambda_2_true, c = "k", s = 50 )
plt.xlim( 0, 5)
plt.ylim( 0, 5)
plt.title( "Landscape formed by Uniform priors on $p_1, p_2$." )

subplot(223)
plt.contour( X, Y, M*L )
im = plt.imshow(M*L, interpolation='none', origin='lower',
    cmap=cm.jet, extent=(0,5,0,5))
plt.title( "Landscape warped by data observations;\n Uniform priors on $p_1, p_2$." )
plt.scatter( lambda_1_true, lambda_2_true, c = "k", s = 50 )
plt.xlim( 0, 5)
plt.ylim( 0, 5)

subplot(222)
exp_x = stats.expon.pdf( x,loc=0, scale = 3)
exp_y = stats.expon.pdf(x, loc=0, scale = 10)
M = np.dot( exp_x[:,None], exp_y[None,:] )

plt.contour( X, Y, M )
im = plt.imshow(M, interpolation='none', origin='lower',
    cmap=cm.jet, extent=(0,5,0,5))
plt.scatter( lambda_1_true, lambda_2_true, c = "k", s = 50 )
plt.xlim( 0, 5)
plt.ylim( 0, 5)
plt.title( "Landscape formed by Exponential priors on $p_1, p_2$." )

subplot(224)
plt.contour( X, Y, M*L )
im = plt.imshow(M*L, interpolation='none', origin='lower',
    cmap=cm.jet, extent=(0,5,0,5))

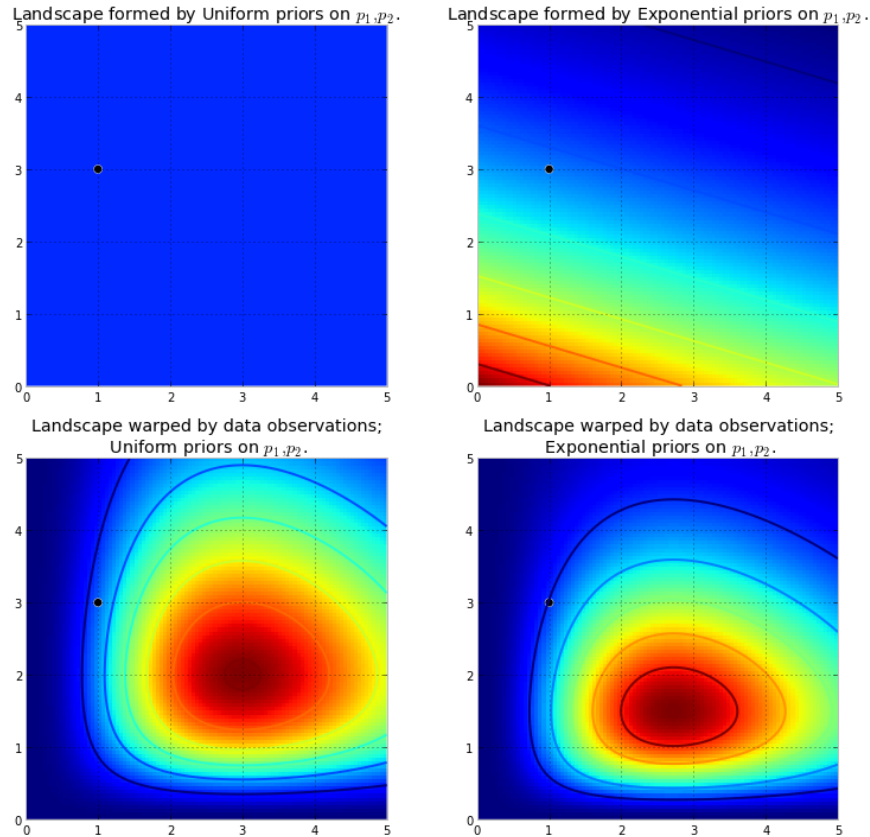
plt.scatter( lambda_1_true, lambda_2_true, c = "k", s = 50 )
plt.title( "Landscape warped by data observations;\n Exponential priors on $p_1, p_2$." )

```

```

$p_1, p_2$.")
plt.xlim( 0, 5)
plt.ylim( 0, 5);

```



The plot on the left is the deformed landscape with the Uniform(0, 5) priors, and the plot on the right is the deformed landscape with the exponential priors. Notice that the posterior landscapes look different from one another, though the data observed is identical in both cases. The reason is as follows. The exponential-prior landscape, on the right, puts very little posterior weight on values in the upper right corner: this is because *the prior does not put much weight there*, whereas the uniform-prior landscape is happy to put posterior weight there. Also, the highest-point, corresponding to the darkest red, is biased towards (0,0) in the exponential case, which is the result from the exponential prior putting more prior weight in the (0,0) corner.

The black dot represents the true parameters. Even with 1 sample point, the mountain attempts to contain the true parameter. Of course, inference with a sample size of 1 is incredibly naive, and choosing such a small sample size was only illustrative. Try observing how the mountain changes as you vary the sample size.

## Exploring the landscape using the MCMC

We should explore the deformed posterior space generated by our prior surface and observed data to find the posterior mountain. However, we cannot naively search the space: any computer scientist will tell you that traversing  $N$ -dimensional space is exponentially difficult in  $N$ : the size of the space quickly blows-up as we increase  $N$  (see [the curse of dimensionality](#)). What hope do we have to find these hidden mountains? The idea behind MCMC is to perform an intelligent search of the space. To say "search" implies we are looking for a particular object, which is perhaps not an accurate description of what MCMC is doing. Recall: MCMC returns *samples* from the posterior distribution, not the distribution itself. Stretching our mountainous analogy to its limit, MCMC performs a task similar to repeatedly asking "How likely is this pebble I found to be from the mountain I am searching for?", and completes its task by returning thousands of accepted pebbles in hopes of reconstructing the original mountain. In MCMC and

PyMC lingo, the returned sequence of “pebbles” are the samples, more often called the *traces*.

When I say MCMC intelligently searches, I mean MCMC will *hopefully* converge towards the areas of high posterior probability. MCMC does this by exploring nearby positions and moving into areas with higher probability. Again, perhaps “converge” is not an accurate term to describe MCMC’s progression. Converging usually implies moving towards a point in space, but MCMC moves towards a *broader area* in the space and randomly walks in that area, picking up samples from that area.

At first, returning thousands of samples to the user might sound like being an inefficient way to describe the posterior distributions. I would argue that this is extremely efficient. Consider the alternative possibilities::

1. Returning a mathematical formula for the “mountain ranges” would involve describing a N-dimensional surface with arbitrary peaks and valleys.
2. Returning the “peak” of the landscape, while mathematically possible and a sensible thing to do as the highest point corresponds to most probable estimate of the unknowns, ignores the shape of the landscape, which we have previously argued is very important in determining posterior confidence in unknowns.

Besides computational reasons, likely the strongest reason for returning samples is that we can easily use *The Law of Large Numbers* to solve otherwise intractable problems. I postpone this discussion for the next chapter. With the thousands of samples, we can reconstruct the posterior surface by organizing them in a histogram.

## Algorithms to perform MCMC

There is a large family of algorithms that perform MCMC. Most of these algorithms can be expressed at a high level as follows: (Mathematical details can be found in the appendix. )

1. Start at current position.
2. Propose moving to a new position (investigate a pebble near you).
3. Accept/Reject the new position based on the position’s adherence to the data and prior distributions (ask if the pebble likely came from the mountain).
4. - If you accept: Move to the new position. Return to Step 1.  
- Else: Do not move to new position. Return to Step 1.
5. After a large number of iterations, return the positions.

This way we move in the general direction towards the regions where the posterior distributions exist, and collect samples sparingly on the journey. Once we reach the posterior distribution, we can easily collect samples as they likely all belong to the posterior distribution.

If the current position of the MCMC algorithm is in an area of extremely low probability, which is often the case when the algorithm begins (typically at a random location in the space), the algorithm will move in positions *that are likely not from the posterior* but better than everything else nearby. Thus the first moves of the algorithm are not reflective of the posterior.

In the above algorithm’s pseudocode, notice that only the current position matters (new positions are investigated only near the current position). We can describe this property as *memorylessness*, i.e. the algorithm does not care *how* it arrived at it’s current position, only that it is there.

## Other approximation solutions to the posterior

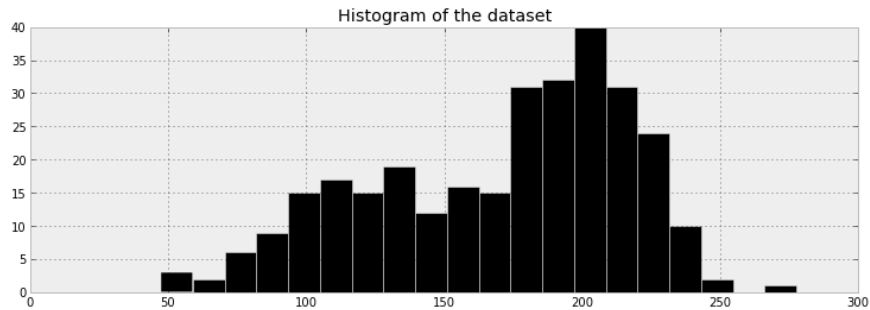
Besides MCMC, there are other procedures available for determining the posterior distributions. A [Laplace approximation](#) is an approximation of the posterior using simple functions. A more advanced method is [Variational Bayes](#).

**Example: Unsupervised Clustering using a Mixture Model** Suppose we are given the following dataset:

```
In [112]: figsize( 12.5, 4)
data = np.loadtxt( "data/mixture_data.csv", delimiter="," )

hist( data, bins = 20, color = "k" )
plt.title("Histogram of the dataset" )
print data[ :10 ], "..."
```

```
[ 115.85679142  152.26153716  178.87449059  162.93500815  107.02820697
 105.19141146  118.38288501  125.3769803  102.88054011  206.71326136] ...
```



What does the data suggest? It appears the data has a bimodal form, that is, it appears to have two peaks, one near 120 and the other near 200. Perhaps there are *two clusters* within this dataset.

This dataset is a good example of the data-generation modeling technique from last chapter. We can propose *how* the data might have been created. I suggest the following data generation algorithm:

1. For each data point, choose cluster 1 with probability  $p$ , else choose cluster 2.
2. Draw a random variate from a Normal distribution with parameters  $\mu_i$  and  $\sigma_i$  where  $i$  was chosen in step 1.
3. Repeat.

This algorithm would create a similar effect as the observed dataset, so we choose this as our model. Of course, we do not know  $p$  or the parameters of the Normal distributions. Hence we must infer, or *learn*, these unknowns.

Denote the Normal distributions  $Nor_0$  and  $Nor_1$  (having variables' index start at 0 is just Pythonic). Both currently have unknown mean and standard deviation, denoted  $\mu_i$  and  $\sigma_i$ ,  $i = 0, 1$  respectively. A specific data point can be from either  $Nor_0$  or  $Nor_1$ , and we assume that the data point is assigned to  $Nor_0$  with probability  $p$ .

An appropriate way to assign data points to clusters is to use a PyMC `Categorical` stochastic variable. Its parameter is a  $k$ -length array of probabilities that must sum to one and its `value` attribute is a integer between 0 and  $k - 1$  randomly chosen according to the crafted array of probabilities. (In our case  $k = 2$ ) *A priori*, we do not know what the probability of assignment to cluster 1 is, so we create a uniform variable over 0,1 to model this. Call this  $p$ . Thus the probability array we enter into the `Categorical` variable is  $[p, 1-p]$ .

```
In [93]: import pymc as mc

p = mc.Uniform( "p", 0, 1)

assignment = mc.Categorical("assignment", [p, 1-p], size = data.shape[0] )
print "prior assignment, with p = %.2f:"%p.value
print assignment.value[ :10 ], "..."
```

```
prior assignment, with p = 0.45:
[1 1 0 1 1 1 0 1 0 1] ...
```

Looking at the above dataset, I would guess that the standard deviations of the two Normals are different. To maintain ignorance of what the standard deviations might be, we will initially model them as uniform on 0 to 100. Really we are talking about  $\tau$ , the *precision* of the Normal distribution, but it is easier to think in terms of standard deviation. Our PyMC code will need to transform our standard deviation into precision by the relation:



$$\tau = \frac{1}{\sigma^2}$$

In PyMC, we can do this in one step by writing:

```
taus = 1.0/mc.Uniform( "stds", 0, 100, size= 2)**2
```

Notice that we specified `size=2`: we are modeling both  $\tau$ s as a single PyMC variable. Note that it does not induce a necessary relationship between the two  $\tau$ s, it is simply for succinctness.

We also need to specify priors on the centers of the clusters. The centers are really the  $\mu$  parameters in this Normal distributions. Their priors can be modeled by a Normal distribution. Looking at the data, I have an idea where the two centers might be — I would guess somewhere around 120 and 190 respectively, though I am not very confident in these eyeballed estimates. Hence I will set  $\mu_0 = 120$ ,  $\mu_1 = 190$  and  $\sigma_{0,1} = 10$  (recall we enter the  $\tau$  parameter, so enter  $1/\sigma^2 = 0.01$  in the PyMC variable.)

```
In [94]: taus = 1.0/mc.Uniform( "stds", 0, 100, size= 2)**2
        centers = mc.Normal( "centers", [120, 190], [0.01, 0.01], size= 2 )

        """
        The below deterministic functions map a assignment, in this case 0 or 1,
        to a set of parameters, located in the (1,2) arrays 'taus' and 'centers.'
        """

        @mc.deterministic
        def center_i( assignment = assignment, centers = centers ):
            return centers[ assignment]

        @mc.deterministic
        def tau_i( assignment = assignment, taus = taus ):
            return taus[ assignment]

        print "Random assignments: ", assignment.value[ :4 ], "..."
        print "Assigned center: ", center_i.value[ :4], "..."
        print "Assigned precision: ", tau_i.value[ :4], "..."
```

```
Random assignments: [1 1 0 1] ...
Assigned center: [ 174.87922993  174.87922993  129.56485407  174.87922993] ...
Assigned precision: [ 0.00012929  0.00012929  0.00050314  0.00012929] ...
```

```
In [95]: #and to combine it with the observations:
        observations = mc.Normal( "obs", center_i, tau_i, value = data, observed =
        #below we create a model class
        model = mc.Model( [p, assignment, taus, centers ] )
```

PyMC has an MCMC class, `MCMC` in the main namespace of PyMC, that implements the MCMC exploring algorithm. We initialize it by passing in a `Model` instance:

```
mcmc = mc.MCMC( model )
```

The method for asking the MCMC to explore the space is `sample( iterations )`, where `iterations` is the number of steps you wish the algorithm to perform. We try 50000 steps below:

```
In [96]: mcmc = mc.MCMC( model )
        mcmc.sample( 50000 )
```

```
[*****100%*****] 50000 of 50000 complete
```

Below I plot the paths, or “traces”, the unknown parameters (centers, precisions, and  $p$ ) have taken thus far. The traces can be retrieved using the `trace` method in the MCMC object created, which accepts the assigned PyMC variable name. For example, `mcmc.trace("centers")` will retrieve a Trace object that can be indexed (using `[:]` or `.gettrace()` to retrieve all traces, or fancy-indexing like `[1000:]`).

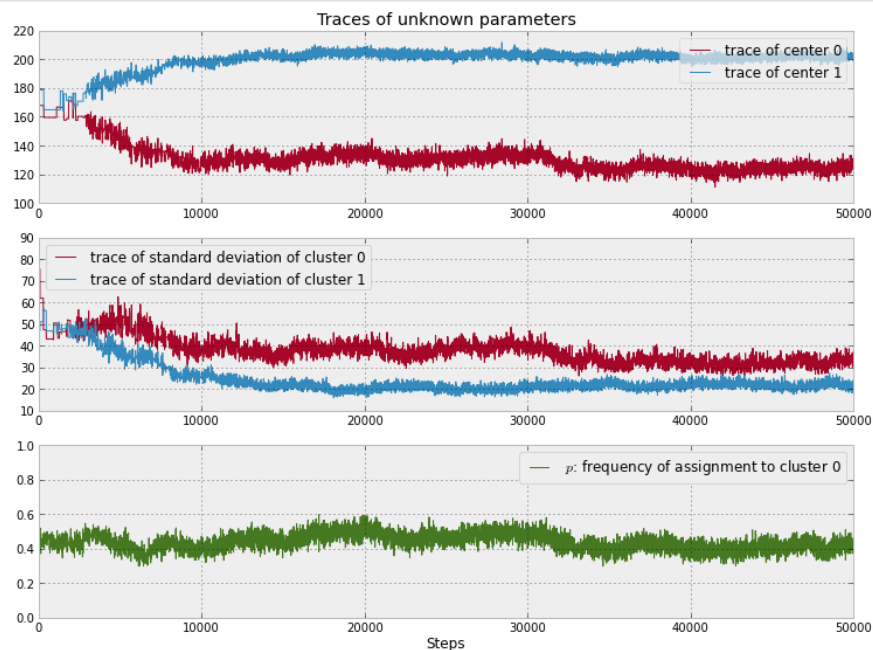
```
In [99]: figsize( 12.5, 9 )
subplot(311)
lw = 1
center_trace = mcmc.trace("centers")[:]

#for pretty colors later in the book.
colors = ["#348ABD", "#A60628"] if center_trace[-1,0] > center_trace[-1,1]
else [ "#A60628", "#348ABD" ]

plot( center_trace[:,0], label = "trace of center 0", c = colors[0], lw =
plot( center_trace[:,1], label = "trace of center 1", c = colors[1], lw =
plt.title( "Traces of unknown parameters" )
leg = plt.legend(loc = "upper right")
leg.get_frame().set_alpha(0.7)

subplot(312)
std_trace = mcmc.trace("stds")[:]
plot( std_trace[:,0], label = "trace of standard deviation of cluster 0",
      c = colors[0], lw = lw )
plot( std_trace[:,1], label = "trace of standard deviation of cluster 1",
      c = colors[1], lw = lw )
plt.legend(loc = "upper left")

subplot(313)
p_trace = mcmc.trace("p")[:]
plot( p_trace, label = "$p$: frequency of assignment to cluster 0",
      color = "#467821", lw = 1 )
plt.xlabel( "Steps" )
plt.ylim(0,1)
plt.legend();
```



Notice the following characteristics:

1. The traces converges, not to a single point, but to a *distribution* of possible points. This is *convergence* in an MCMC algorithm.
2. Inference using the first few thousand points is a bad idea, as they are unrelated to the final distribution we are interested in. Thus is it a good idea to discard those samples before using the samples for inference. We call this period before converge the *burn-in period*.
3. The traces appear as a random “walk” around the space, that is, the paths exhibit correlation with previous positions. This is both good and bad. We will always have correlation between current positions and the previous positions, but too much of it means we are not exploring the space well. This will be detailed in the Diagnostics section later in this chapter.

To achieve further convergence, we will perform more MCMC steps. Starting the MCMC again after it has already been called does not mean starting the entire algorithm over. In the pseudo-code algorithm of MCMC above, the only position that matters is the current position (new positions are investigated near the current position), implicitly stored in PyMC variables’ `value` attribute. Thus it is fine to halt an MCMC algorithm and inspect its progress, with the intention of starting it up again later. The ‘value’ attributes are not overwritten.

We will sample the MCMC one hundred thousand more times and visualize the progress below:

```
In [100]: mcmc.sample( 100000 )
```

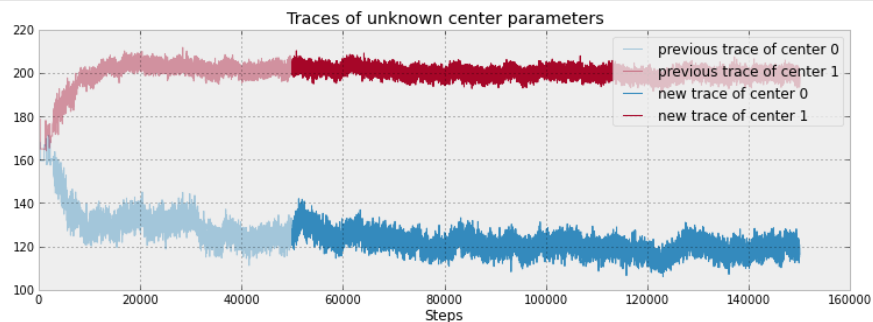
```
[*****100%*****] 100000 of 100000 complete
```

```
In [101]: figsize(12.5, 4)
center_trace = mcmc.trace("centers", chain = 1)[: ]
prev_center_trace = mcmc.trace("centers", chain = 0)[: ]

x = np.arange(50000)
plot(x, prev_center_trace[:,0], label = "previous trace of center 0",
     lw = lw, alpha = 0.4 , c= colors[1] )
plot(x, prev_center_trace[:,1], label = "previous trace of center 1",
     lw = lw, alpha = 0.4, c = colors[0] )

x = np.arange(50000, 150000)
plot(x, center_trace[:,0], label = "new trace of center 0", lw = lw, c = colors[1] )
plot(x, center_trace[:,1], label = "new trace of center 1", lw = lw, c = colors[0] )

plt.title( "Traces of unknown center parameters" )
leg = plt.legend(loc = "upper right")
leg.get_frame().set_alpha(0.8)
plt.xlabel( "Steps" );
```



The `trace` method in the MCMC instance has a keyword argument `chain`, that indexes which call to `sample` you would like to be returned. (Often we need to call `sample` multiple times, and the ability to retrieve past samples is a useful procedure). The default for `chain` is -1, which will return the samples from the latest call to `sample`.

## Cluster Investigation

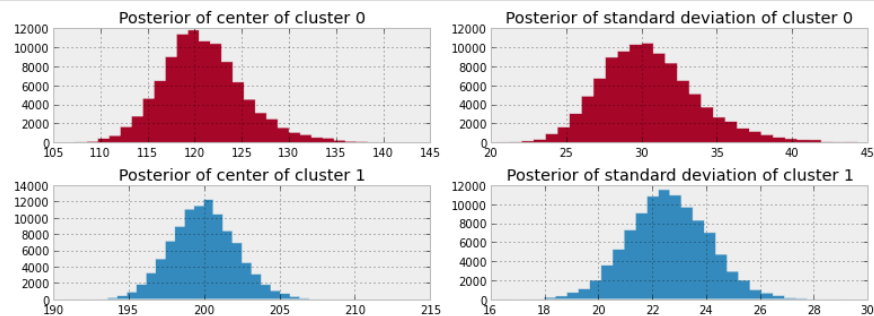
We have not forgotten our main challenge: identify the clusters. We have determined posterior distributions for our unknowns. We plot the posterior distributions of the center and standard deviation variables below:

```
In [102]: figsize( 11.0, 4 )
std_trace = mcmc.trace("stds")[:]

_i = [ 1,2,3,0]
for i in range(2):
    subplot(2,2, _i[ 2*i ] )
    plt.title("Posterior of center of cluster %d"%i)
    plt.hist( center_trace[:, i], color = colors[i],bins = 30,
             histtype="stepfilled" )

    subplot(2,2, _i[ 2*i + 1 ] )
    plt.title( "Posterior of standard deviation of cluster %d"%i )
    plt.hist( std_trace[:, i], color = colors[i], bins = 30,
             histtype="stepfilled" )
    #plt.autoscale(tight=True)

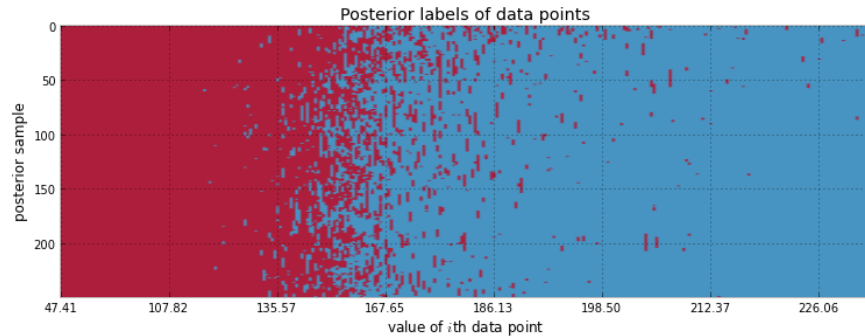
plt.tight_layout()
```



The MCMC algorithm has proposed that the most likely centers of the two clusters are near 120 and 200 respectively. Similar inference can be applied to the standard deviation.

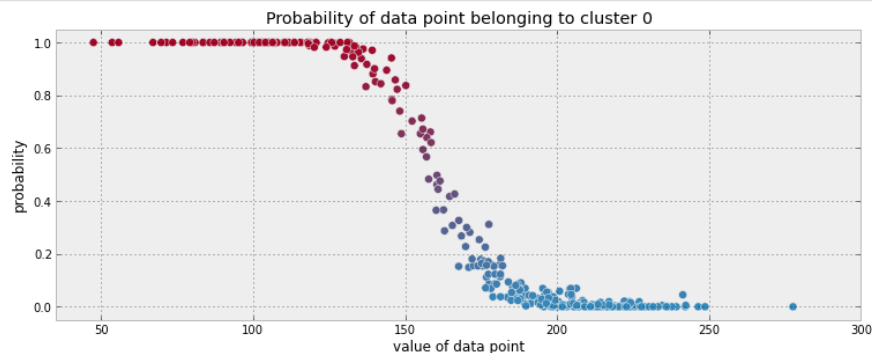
We are also given the posterior distributions for the labels of the data point, which is present in `mcmc.trace("assignment")`. Below is a visualization of this. The y-axis represents a subsample of the posterior labels for each data point. The x-axis are the sorted values of the data points. A red square is an assignment to cluster 1, and a blue square is an assignment to cluster 0.

```
In [103]: figsize( 12.5, 4.5 )
cmap = mpl.colors.ListedColormap(colors)
imshow( mcmc.trace("assignment")[:,400,np.argsort(data) ],
        cmap = cmap ,aspect=.4, alpha = .9 )
xticks(np.arange(0, data.shape[0], 40 ), [ "%.2f"%s for s in np.sort( data
plt.ylabel("posterior sample")
plt.xlabel("value of $i$th data point")
plt.title("Posterior labels of data points" );
```



Looking at the above plot, it appears that the most uncertainty is between 150 and 170. The above plot slightly misrepresents things, as the x-axis is not a true scale (it displays the value of the  $i$ th sorted data point.) A more clear diagram is below, where we have estimated the *frequency* of each data point belonging to the labels 0 and 1.

```
In [104]: cmap = mpl.colors.LinearSegmentedColormap.from_list("BMH", colors )
assign_trace = mcmc.trace("assignment")[:]
scatter( data, 1-assign_trace.mean(axis=0), cmap=cmap,
        c=assign_trace.mean(axis=0), s = 50)
plt.ylim( -0.05, 1.05 )
plt.xlim( 35, 300)
plt.title( "Probability of data point belonging to cluster 0" )
plt.ylabel("probability")
plt.xlabel("value of data point" );
```



Even though we modeled the clusters using Normal distributions, we didn't get just a single Normal distribution that *best* fits the data (whatever our definition of best is), but a distribution of values for the Normal's parameters. How can we choose just a single pair of values for the mean and variance and determine a *sorta-best-fit* gaussian?

One quick and dirty way (which has nice theoretical properties we will see in Chapter 5), is to use the *mean* of the posterior distributions. Below we overlay the Normal density functions, using the mean of the posterior distributions as the chosen parameters, with our observed data:

```
In [106]: norm = stats.norm
x = np.linspace( 20, 300, 500 )
posterior_center_means = center_trace.mean(axis=0)
posterior_std_means = std_trace.mean(axis=0)
posterior_p_mean = mcmc.trace("p")[:].mean()

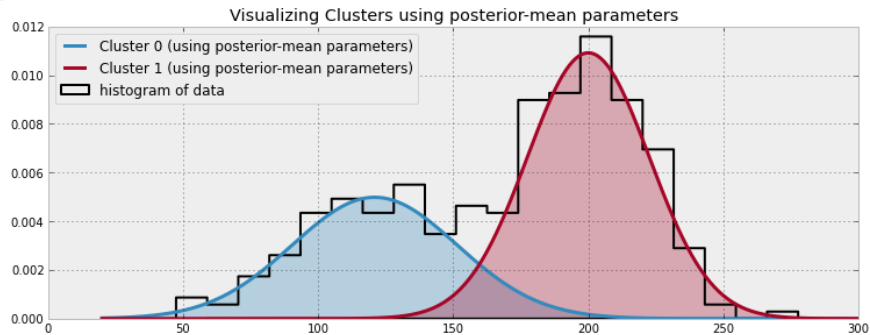
hist( data, bins = 20, histtype="step", normed = True, color = "k",
      lw = 2, label = "histogram of data" )
y = posterior_p_mean*norm.pdf(x, loc=posterior_center_means[0],
                             scale=posterior_std_means[0] )
plt.plot(x, y, label = "Cluster 0 (using posterior-mean parameters)", lw=3)
plt.fill_between( x, y, color = colors[1], alpha = 0.3 )
```

```

x = (1-posterior_p_mean)*norm.pdf(x, loc=posterior_center_means[1],
                                scale=posterior_std_means[1] )
plt.plot(x, y, label = "Cluster 1 (using posterior-mean parameters)", lw = 2)
plt.fill_between( x,y, color = colors[0], alpha = 0.3, )

plt.legend(loc = "upper left")
plt.title( "Visualizing Clusters using posterior-mean parameters" );

```



### Important: Don't mix posterior samples

In the above example, a possible (though less likely) scenario is that cluster 0 has a very large standard deviation, and cluster 1 has a small standard deviation. This would still satisfy the evidence, albeit less so than our original inference. Alternatively, it would be incredibly unlikely for *both* distributions to have a small standard deviation, as the data does not support this hypothesis at all. Thus the two standard deviations are *dependent* on each other: if one is small, the other must be large. In fact, *all* the unknowns are related in a similar manner. For example, if a standard deviation is large, the mean has a wider possible space of realizations. Conversely, a small standard deviation restricts the mean to a small area.

During MCMC, we are returned vectors representing samples from the unknown posteriors. Elements of different vectors cannot be used together, as this would break the above logic: perhaps a sample has returned that cluster 1 has a small standard deviation, hence all the other variables in that sample would incorporate that and be adjusted accordingly. It is easy to avoid this problem though, just make sure you are indexing traces correctly.

Another small example to illustrate the point. Suppose two variables,  $x$  and  $y$ , are related by  $x + y = 10$ . We model  $x$  as a Normal random variable with mean 4 and explore 300 samples.

```

In [2]: import pymc as mc

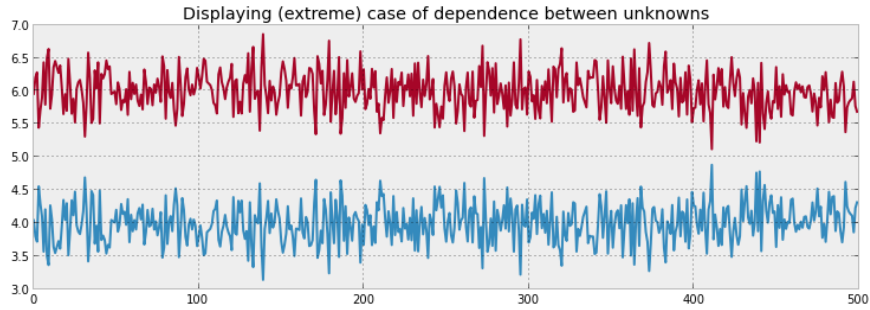
x = mc.Normal("x", 4, 10 )
y = mc.Lambda( "y", lambda x=x: 10-x, trace = True )

ex_mcmc = mc.MCMC( mc.Model( [x,y] ) )
ex_mcmc.sample( 500 )

plt.plot( ex_mcmc.trace("x")[:] )
plt.plot( ex_mcmc.trace("y")[:] )
plt.title( "Displaying (extreme) case of dependence between unknowns" );

```

[\*\*\*\*\*100%\*\*\*\*\*] 500 of 500 complete



As you can see, the two variables are not unrelated, and it would be wrong to add the  $i$ th sample of  $x$  to the  $j$ th sample of  $y$ , unless  $i = j$ .

### Returning to Clustering: Prediction

The above clustering can be generalized to  $k$  clusters. Choosing  $k = 2$  allowed us to visualize the MCMC better, and examine some very interesting plots.

What about prediction? Suppose we observe a new data point, say  $x = 175$ , and we wish to label it to a cluster. It is foolish to simply assign it to the *closer* cluster center, as this ignores the standard deviation of the clusters, and we have seen from the plots above that this consideration is very important. More formally: we are interested in the *probability* (as we cannot be certain about labels) of assigning  $x = 175$  to cluster 1. Denote the assignment of  $x$  as  $L_x$ , which is equal to 0 or 1, and we are interested in  $P(L_x = 1 | x = 175)$ .

A naive method to compute this is to re-run the above MCMC with the additional data point appended. The disadvantage with this method is that it will be slow to infer for each novel data point. Alternatively, we can try a *less precise*, but much quicker method.

We will use Bayes Theorem for this. If you recall, Bayes Theorem looks like:

$$P(A|X) = \frac{P(X|A)P(A)}{P(X)}$$

In our case,  $A$  represents  $L_x = 1$  and  $X$  is the evidence we have: we observe that  $x = 175$ . For a particular sample set of parameters for our posterior distribution,  $(\mu_0, \sigma_0, \mu_1, \sigma_1, p)$ , we are interested in asking “Is the probability that  $x$  is in cluster 1 *greater* than the probability it is in cluster 0?”, where the probability is dependent on the chosen parameters.

$$P(L_x = 1|x = 175) < P(L_x = 0|x = 175) \tag{22}$$

$$\tag{23}$$

$$\frac{P(x = 175|L_x = 1)P(L_x = 1)}{P(x = 175)} < \frac{P(x = 175|L_x = 0)P(L_x = 0)}{P(x = 175)} \tag{24}$$

As the denominators are equal, they can be ignored (and good riddance, because computing the quantity  $P(x = 175)$  can be difficult).

$$P(x = 175|L_x = 1)P(L_x = 1) < P(x = 175|L_x = 0)P(L_x = 0)$$

```
In [385]: norm_pdf = stats.norm.pdf
p_trace = mcmc.trace("p")[:]
x = 175

v = p_trace*norm_pdf(x, loc = center_trace[:,0], scale = std_trace[:,0] )
    (1-p_trace)*norm_pdf(x, loc = center_trace[:,1], scale = std_trace[:,1] )

print "Probability of belonging to cluster 1:", v.mean()
```

Probability of belonging to cluster 1: 0.9406

Giving us a probability instead of a label is a very useful thing. Instead of the naive

$L = 1$  if  $\text{prob} > 0.5$  else 0

we can optimize our guesses using *loss function*, of which the entire fifth chapter is devoted to.

### Using MAP to improve convergence

If you ran the above example yourself, you may have noticed that our results were not consistent: perhaps your cluster division was more scattered, or perhaps less scattered. The problem is that our traces are a function of the *starting values* of the MCMC algorithm.

It can be mathematically shown that letting the MCMC run long enough, by performing many steps, the algorithm *should forget its initial position*. In fact, this is what it means to say the MCMC converged (in practice though we can never achieve total convergence). Hence if we observe different posterior analysis, it is likely because our MCMC has not fully converged yet, and we should not use samples from it yet (we should use a larger burn-in period).

In fact, poor starting values can prevent any convergence, or significantly slow it down. Ideally, we would like to have the chain start at the *peak* of our landscape, as this is exactly where the posterior distributions exist. Hence, if we started at the “peak”, we could avoid a lengthy burn-in period and incorrect inference. Generally, we call this “peak” the *maximum a posterior* or, more simply, the *MAP*.

Of course, we do not know where the MAP is. PyMC provides an object that will approximate, if not find, the MAP location. In the PyMC main namespace is the `MAP` object that accepts a PyMC `Model` instance. Calling `.fit()` from the `MAP` instance sets the variables in the model to their MAP values.

```
map_ = mc.MAP( model )
map_.fit()
```

The `MAP.fit()` methods has the flexibility of allowing the user to choose which optimization algorithm to use (after all, this is an optimization problem: we are looking for the values that maximize our landscape), as not all optimization algorithms are created equal. The default optimization algorithm in the call to `fit` is `scipy`'s `fmin` algorithm (which attempts to minimize the *negative of the landscape*). An alternative algorithm that is available is Powell's Method, a favourite of PyMC blogger [Abraham Flaxman](#) [1], by calling `fit(method='fmin_powell')`. From my experience, I use the default, but if my convergence is slow or not guaranteed, I experiment with Powell's method.

The MAP can also be used as a solution to the inference problem, as mathematically it is the *most likely* value for the unknowns. But as mentioned earlier in this chapter, this location ignores the uncertainty and doesn't return a distribution.

Typically, it is always a good idea, and rarely a bad idea, to prepend your call to `mcmc` with a call to `MAP(model).fit()`. The intermediate call to `fit` is hardly computationally intensive, and will save you time later due to a shorter burn-in period.



## Speaking of the burn-in period

It is still a good idea to provide a burn-in period, even if we are using MAP prior to calling `MCMC.sample`, just to be safe. We can have PyMC automatically discard the first  $n$  samples by specifying the `burn` parameter in the call to `sample`. As one does not know when the chain has fully converged, I like to assign the first *half* of my samples to be discarded, sometimes up to 90% of my samples for longer runs. To continue the clustering example from above, my new code would look something like:

```
model = mc.Model( [p, assignment, taus, centers ] )

map_ = mc.MAP( model )
map_.fit() #stores the fitted variables' values in foo.value

mcmc = mc.MCMC( model )
mcmc.sample( 100000, 50000 )
```

## 0.5.2 Diagnosing Convergence

### Autocorrelation

Autocorrelation is a measure of how related a series of numbers is with itself. A measurement of 1.0 is perfect positive autocorrelation, 0 no autocorrelation, and -1 is perfect negative correlation. If you are familiar with standard *correlation*, then autocorrelation is just how correlated a series,  $x_t$ , at time  $t$  is with the series at time  $t - k$ :

$$R(k) = \text{Corr}(x_t, x_{t-k})$$

For example, consider the two series:

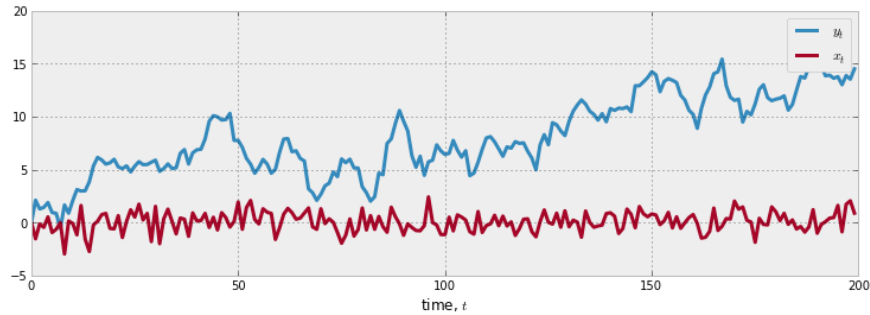
$$x_t \sim \text{Normal}(0, 1), \quad x_0 = 0$$
$$y_t \sim \text{Normal}(y_{t-1}, 1), \quad y_0 = 0$$

which has an example paths like:

```
In [40]: figsize(12.5, 4)

import pymc as mc
x_t = mc.rnormal( 0, 1, 200 )
x_t[0] = 0
y_t = np.zeros( 200 )
for i in range(1,200):
    y_t[i] = mc.rnormal( y_t[i-1], 1 )

plt.plot( y_t, label= "$y_t$", lw = 3 )
plt.plot( x_t, label= "$x_t$", lw = 3 )
plt.xlabel("time, $t$")
plt.legend();
```



One way to think of autocorrelation is “If I know the position of the series at time  $s$ , can it help me know where I am at time  $t$ ?” In the series  $x_t$ , the answer is No. By construction,  $x_t$  are random variables. If I told you that  $x_2 = 0.5$ , could you give me a better guess about  $x_3$ ? No.

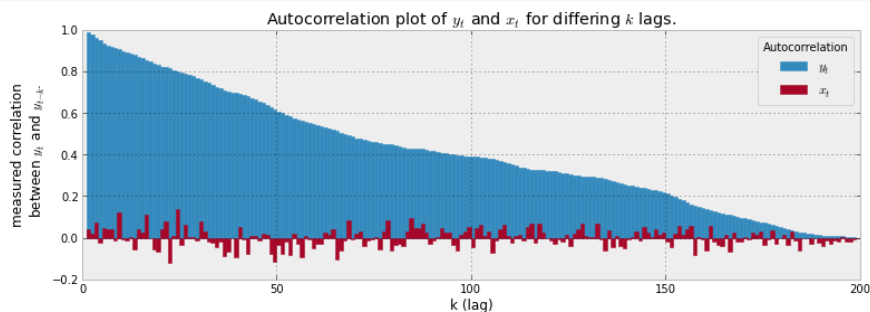
On the other hand,  $y_t$  is autocorrelated. By construction, if I knew that  $y_2 = 10$ , I can be very confident that  $y_3$  will not be very far from 10. Similarly, I can even make a (less confident guess) about  $y_4$ : it will probably not be near 0 or 20, but a value of 5 is not too unlikely. I can make a similar argument about  $y_5$ , but again, I am less confident. Taking this to its logical conclusion, we must concede that as  $k$ , the lag between time points, increases the autocorrelation decreases. We can visualize this:

```
In [41]: def autocorr(x):
#from http://tinyurl.com/afz57c4
result = np.correlate(x, x, mode = 'full')
result = result / np.max(result)
return result[result.size/2:]

colors = [ "#348ABD", "#A60628", "#7A68A6" ]

x = np.arange(1, 200)
plt.bar(x, autocorr( y_t ) [1:], width =1, label="$y_{t}$",
        edgcolor = colors[0], color = colors[0] )
plt.bar(x, autocorr( x_t ) [1:], width =1, label = "$x_{t}$",
        color = colors[1], edgcolor = colors[1] )

plt.legend( title="Autocorrelation" )
plt.ylabel("measured correlation \n between $y_{t}$ and $y_{t-k}$.")
plt.xlabel("k (lag)")
plt.title("Autocorrelation plot of $y_{t}$ and $x_{t}$ for differing $k$ lags.")
```



Notice that as  $k$  increases, the autocorrelation of  $y_t$  decreases from a very high point. Compare with the autocorrelation of  $x_t$  which looks like noise (which it really is), hence we can conclude no autocorrelation exists in this series.

### How does this relate to MCMC convergence?

By the nature of the MCMC algorithm, we will always be returned samples that exhibit autocorrelation (this is because of the step from your current position, move to a position near you).

A chain that is [Isn't meandering exploring?] exploring the space well will exhibit very high autocorrelation. Visually, if the trace seems to meander like a river, and not settle down, the chain will have high autocorrelation.

This does not imply that a converged MCMC has low autocorrelation. Hence low autocorrelation is not necessary for convergence, but it is sufficient. PyMC has a built-in autocorrelation plotting function in the `Matplot` module.

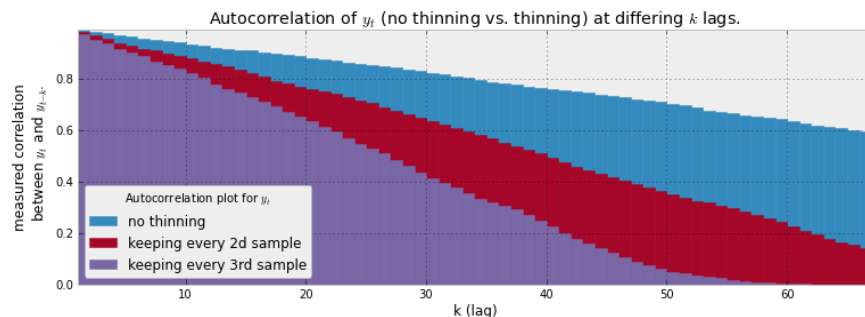
## Thinning

Another issue can arise if there is high-autocorrelation between posterior samples. Many post-processing algorithms require samples to be *independent* of each other. This can be solved, or at least reduced, by only returning to the user every  $n$ th sample, thus removing some autocorrelation. Below we perform an autocorrelation plot for  $y_t$  with differing levels of thinning:

```
In [25]: max_x = 200/3+1
x = np.arange(1,max_x )

plt.bar(x, autocorr( y_t ) [1:max_x], edgecolor=colors[0],
        label="no thinning", color = colors[0], width =1 )
plt.bar(x, autocorr( y_t [::2] ) [1:max_x], edgecolor=colors[1],
        label="keeping every 2d sample", color = colors[1], width=1 )
plt.bar(x, autocorr( y_t [::3] ) [1:max_x], width =1, edgecolor = colors[2],
        label="keeping every 3rd sample", color = colors[2] )

plt.autoscale(tight=True)
plt.legend( title="Autocorrelation plot for $y_t$",loc="lower left" )
plt.ylabel("measured correlation \nbetween $y_t$ and $y_{t-k}$.")
plt.xlabel("k (lag)")
plt.title("Autocorrelation of $y_t$ (no thinning vs. thinning) \
at differing $k$ lags.");
```



With more thinning, the autocorrelation drops quicker. There is a tradeoff though: higher thinning requires more MCMC iterations to achieve the same number of returned samples. For example, 10 000 samples unthinned is 100 000 with a thinning of 10 (though the latter has less autocorrelation).

What is a good amount of thinning. The returned samples will always exhibit some autocorrelation, regardless of how much thinning is done. So long as the autocorrelation tends to zero, you are probably ok. Typically thinning of more than 10 is not necessary.

PyMC exposes a thinning parameter in the call the `sample`, for example: `sample( 10000, burn = 5000, thinning = 5)`.

## `pymc.Matplot.plot()`

It seems silly to have to manually create histograms, autocorrelation plots and trace plots each time we perform MCMC. The authors of PyMC have included a visualization tool for just this purpose.

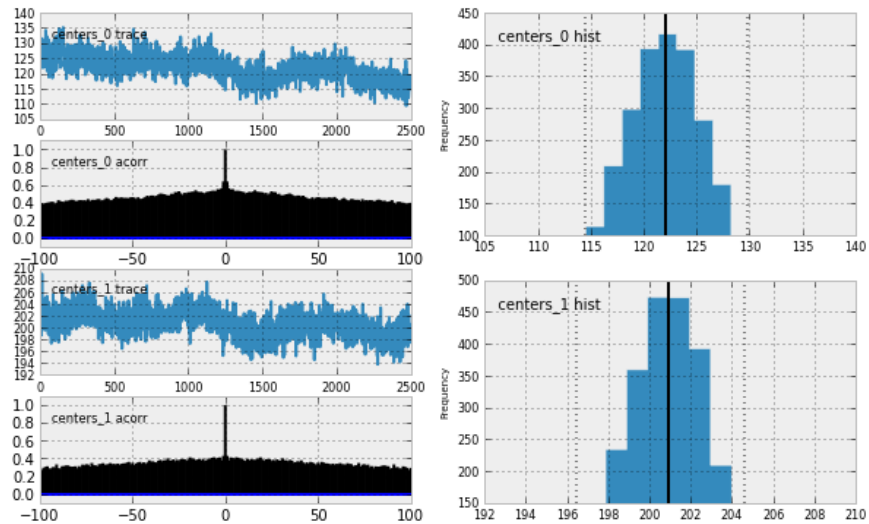
As the title suggests, the `pymc.Matplot` module contains a poorly named function `plot`, which I prefer to import as `mcplot` so there is no conflict with other namespaces. `plot`, or `mcplot` as I suggest, accepts an MCMC object and will return posterior distributions, traces and auto-correlations for each variable (up to 10 variables).

Below we use the tool to plot the centers of the clusters, after sampling 25 000 more times and `thinning = 10`.

```
In [110]: from pymc.Matplot import plot as mcplot
```

```
mcmc.sample( 25000, 0, 10)
mcplot( mcmc.trace("centers",2), common_scale = False )
```

```
[*****100%*****] 25000 of 25000 completePlotting centers_0
Plotting centers_1
```



There are really two figures here, one for each unknown in the `centers` variable. In each figure, the subfigure in the top left corner is the trace of the variable. This is useful for inspecting that possible “meandering” property that is a result of non-convergence. The largest plot on the right-hand side is the histograms of the samples, plus a few extra features. The thickest vertical line represents the posterior mean, which is a good summary of posterior distribution. The interval between the two dashed vertical lines in each the posterior distributions represent the *95% credible interval*, not to be confused with a *95% confidence interval*. I won’t get into the latter, but the former can be interpreted as “there is a 95% chance the parameter of interested lies in this interval”. (Changing default parameters in the call to `mcplot` provides alternatives to 95%.) When communicating your results to others, it is incredibly important to state this interval. One of our purposes for studying Bayesian methods is to have a clear understanding of our uncertainty in unknowns. Combined with the posterior mean, the 95% credible interval provides a reliable interval to communicate the likely location of the unknown (provided by the mean) *and* the uncertainty (represented by the width of the interval). The plots titled `center_0_acorr` and `center_1_acorr` are the generated autocorrelation plots. They look different than the ones I have displayed above, but the only difference is that 0-lag is centered in the middle of the figure, whereas I have 0 centered to the left.

### 0.5.3 Useful tips for MCMC

Bayesian inference would be the *de facto* method if it weren’t for MCMC’s computational difficulties. In fact, MCMC is what turns most users off practical Bayesian inference. Below I present some good heuristics to help convergence and speed up the MCMC engine:

## Intelligent starting values

It would be great to start the MCMC algorithm off near the posterior distribution, so that it will take little time to start sampling correctly. We can aid the algorithm by telling where we *think* the posterior distribution will be by specifying the `value` parameter in the `Stochastic` variable creation. Often we possess a guess about this anyways. For example, if we have data from a Normal distribution, and we wish to estimate the  $\mu$  parameter, then a good starting value would be the *mean* of the data.

```
mu = mc.Uniform( "mu", 0, 100, value = data.mean() )
```

For most parameters in models, there is a frequentist estimate of it. These estimates are a good starting value for our MCMC algorithms. Of course, this is not always possible for some variables, but including as many appropriate initial values is always a good idea. Even if your guesses are wrong, the MCMC will still converge to the proper distribution, so there is little to lose.

This is what using MAP tries to do, by giving good initial values to the MCMC. So why bother specifying user-defined values? Well, even giving MAP good values will help it find the maximum a-posterior.

Also important, *bad initial values* are a source of major bugs in PyMC and can hurt convergence.

## Priors

## Covariance matrices and eliminating parameters

## The Folk Theorem of Statistical Computing

*If you are having computational problems, probably your model is wrong.*

## 0.5.4 Conclusion

PyMC provides a very strong backend to performing Bayesian inference, mostly because it has abstracted the inner mechanics of MCMC from the user. Despite this, some care must be applied to ensure your inference is not being biased by the iterative nature of MCMC.

## References

1. Flaxman, Abraham. "Powell's Methods for Maximization in PyMC." *Healthy Algorithms*. N.p., 9 Oct 2012. Web. 28 Feb 2013. <http://healthyalgorithms.com/2012/02/09/powells-method-for-maximization-in-pymc/>.

```
In [1]: from IPython.core.display import HTML
def css_styling():
    styles = open("../styles/custom.css", "r").read()
    return HTML(styles)
css_styling()
```

Out[1]:

## 0.6 Chapter 4

---

## 0.6.1 The greatest theorem never told

This chapter focuses on an idea that is always bouncing around our minds, but is rarely made explicit outside books devoted to statistics. In fact, we've been using this simple idea in every example thus far.

### The Law of Large Numbers

Let  $Z_i$  be  $N$  independent samples from some probability distribution. According to *the Law of Large numbers*, so long as the expected value  $E[Z]$  is finite, the following holds,

$$\frac{1}{N} \sum_{i=1}^N Z_i \rightarrow E[Z], \quad N \rightarrow \infty.$$

In words:

The average of a sequence of random variables from the same distribution converges to the expected value of that distribution.

This may seem like a boring result, but it will be the most useful tool you use.

### Intuition

If the above Law is somewhat surprising, it can be made more clear by examining a simple example.

Consider a random variable  $Z$  that can take only two values,  $c_1$  and  $c_2$ . Suppose we have a large number of samples of  $Z$ , denoting a specific sample  $Z_i$ . The Law says that we can approximate the expected value of  $Z$  by averaging over all samples. Consider the average:

$$\frac{1}{N} \sum_{i=1}^N Z_i$$

By construction,  $Z_i$  can only take on  $c_1$  or  $c_2$ , hence we can partition the sum over these two values:

$$\frac{1}{N} \sum_{i=1}^N Z_i = \frac{1}{N} \left( \sum_{Z_i=c_1} c_1 + \sum_{Z_i=c_2} c_2 \right) \tag{25}$$

(26)

$$= c_1 \sum_{Z_i=c_1} \frac{1}{N} + c_2 \sum_{Z_i=c_2} \frac{1}{N} \tag{27}$$

(28)

$$= c_1 \times (\text{approximate frequency of } c_1) \tag{29}$$

(30)

$$+ c_2 \times (\text{approximate frequency of } c_2) \tag{31}$$

(32)

$$\approx c_1 \times P(Z = c_1) + c_2 \times P(Z = c_2) \tag{33}$$

(34)

$$= E[Z] \tag{35}$$

Equality holds in the limit, but we can get closer and closer by using more and more samples in the average. This Law holds for *any distribution*, minus some pathological examples that only mathematicians have fun with.

## Example

Below is a diagram of the Law of Large numbers in action for three different sequences of Poisson random variables.

We sample `sample_size= 100000` Poisson random variables with parameter  $\lambda = 4.5$ . (Recall the expected value of a Poisson random variable is equal to its parameter.) We calculate the average for the first  $n$  samples, for  $n = 1$  to `sample_size`.

```
In [3]: %pylab inline

figsize( 12.5, 5 )
import pymc as mc

sample_size = 100000
expected_value = lambda_ = 4.5
poi = mc.rpoisson
N_samples = range(1, sample_size, 100)

for k in range(3):

    samples = poi( lambda_, size = sample_size )

    partial_average = [ samples[:i].mean() for i in N_samples ]

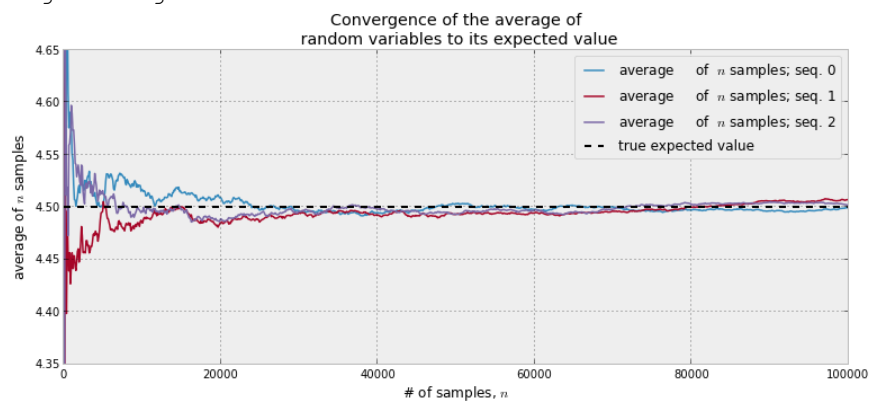
    plt.plot( N_samples, partial_average, lw=1.5, label="average \
of $n$ samples; seq. %d"%k)

plt.plot( N_samples, expected_value*np.ones_like( partial_average), \
ls = "--", label = "true expected value", c = "k" )

plt.ylim( 4.35, 4.65)
plt.title( "Convergence of the average of \n random variables to its \
expected value" )
plt.ylabel( "average of $n$ samples" )
plt.xlabel( "# of samples, $n$" )
plt.legend()
```

<matplotlib.legend.Legend at 0x615b6f0>

Out [3]:



Looking at the above plot, it is clear that when the sample size is small, there is greater variation in the average (compare how *jagged and jumpy* the average is initially, then *smooths* out). All three paths *approach* the value 4.5, but just flirt with it as  $N$  gets large. Mathematicians and statistician have another name for *flirting*: convergence.

Another very relevant question we can ask is *how quickly am I converging to the expected value?* Let's plot something new. For a specific  $N$ , let's do the above trials thousands of times and compute how far away we are from the true expected value, on average. But wait — *compute on average?* This is simply the law of large numbers again! For example, we are interested in, for a specific  $N$ , the quantity:

$$D(N) = \sqrt{E \left[ \left( \frac{1}{N} \sum_{i=1}^N Z_i - 4.5 \right)^2 \right]}$$

(We take the square root so the dimensions of the above quantity and our random variables are the same). As the above is an expected value, it can be approximated using the law of large numbers: instead of averaging  $Z_i$ , we calculate the following multiple times and average them:

$$Y_k = \left( \frac{1}{N} \sum_{i=1}^N Z_i - 4.5 \right)^2$$

i.e., we consider the average

$$\frac{1}{N_Y} \sum_{k=1}^{N_Y} Y_k \approx D(N)$$

where  $N_Y$  is some suitably large number.

```
In [4]: figsize( 12.5, 4)
N_Y = 250
D_N_results = []
N_array = np.arange( 0, 50000, 2500 )
lambda_ = 4.5
expected_value = 4.5

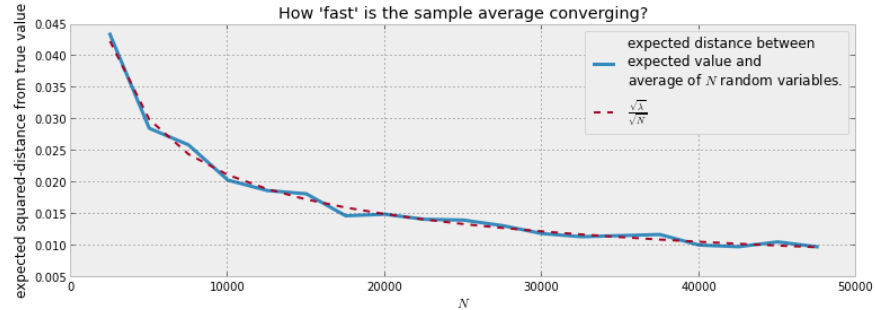
def D_N( n ):
    Z = poi( lambda_, size = (n, N_Y) )
    average_Z = Z.mean(axis=0)
    return np.sqrt( ( (average_Z - expected_value)**2 ).mean() )

for n in N_array:
    D_N_results.append( D_N(n) )

plt.xlabel( "$N$" )
plt.ylabel( "expected squared-distance from true value" )
plt.plot(N_array, D_N_results, lw = 1,
         label="expected distance between\n\
expected value and \naverage of $N$ random variables.")
plt.plot( N_array, np.sqrt(expected_value)/np.sqrt(N_array), lw = 3, ls =
         label = r"$\frac{\sqrt{\lambda}}{\sqrt{N}}$" )
plt.legend()
plt.title( "How 'fast' is the sample average converging? " )
```



<matplotlib.text.Text at 0x12098e10>  
Out [4]:



As expected, the expected distance between our sample average and the actual expected value shrinks as  $N$  grows large. But also notice that the *rate* of convergence decreases, that is, we need only 10 000 additional samples to move from 0.020 to 0.015, a difference of 0.005, but 20 000 more samples to again decrease from 0.015 to 0.010, again only a 0.005 decrease.

It turns out we can measure this rate of convergence. Above I have plotted a second line, the function  $\sqrt{\lambda}/\sqrt{N}$ . This was not chosen arbitrarily. In most cases, given a sequence of random variable distributed like  $Z$ , the rate of converge to  $E[Z]$  of the Law of Large Numbers is

$$\frac{\sqrt{\text{Var}(Z)}}{\sqrt{N}}$$

This is useful to know: for a given large  $N$ , we know (on average) how far away we are from the estimate. On the other hand, in a Bayesian setting, this can seem like a useless result: Bayesian analysis is OK with uncertainty so what's the *statistical* point of adding extra precise digits? Though drawing samples can be so computationally cheap that having a *larger*  $N$  is fine too.

### How do we compute $\text{Var}(Z)$ though?

The variance is simply another expected value that can be approximated! Consider the following, once we have the expected value (by using the Law of Large Numbers to estimate it, denote it  $\mu$ ), we can estimate the variance:

$$\frac{1}{N} \sum_{i=1}^N (Z_i - \mu)^2 \rightarrow E[(Z - \mu)^2] = \text{Var}(Z)$$

### Expected values and probabilities

There is an even less explicit relationship between expected value and estimating probabilities. Define the *indicator function*

$$\mathbb{1}_A(x) = \begin{cases} 1 & x \in A \\ 0 & \text{else} \end{cases}$$

Then, by the law of large numbers, if we have many samples  $X_i$ , we can estimate the probability of an event  $A$ , denoted  $P(A)$ , by:

$$\frac{1}{N} \sum_{i=1}^N \mathbb{1}_A(X_i) \rightarrow E[\mathbb{1}_A(X)] = P(A)$$

Again, this is fairly obvious after a moments thought: the indicator function is only 1 if the event occurs, so we are summing only the times the event occurs and dividing by the total number of trials (consider how we usually approximate probabilities using frequencies). For example, suppose we wish to estimate the probability that a  $Z \sim \text{Exp}(.5)$  is greater than 10, and we have many samples from a  $\text{Exp}(.5)$  distribution.

$$P(Z > 10) = \sum_{i=1}^N \mathbb{1}_{z>10}(Z_i)$$

```
In [4]: import pymc as mc
        N = 10000
        print np.mean( [ mc.rexponential( 0.5 )>10 for i in range(N) ] )
```

0.0058

## What does this all have to do with Bayesian statistics?

*Point estimates*, to be introduced in the next chapter, in Bayesian inference are computed using expected values. In more analytical Bayesian inference, we would have been required to evaluate complicated expected values represented as multi-dimensional integrals. No longer. If we can sample from the posterior distribution directly, we simply need to evaluate averages. Much easier. If accuracy is a priority, plots like the ones above show how fast you are converging. And if further accuracy is desired, just take more samples from the posterior.

When is enough enough? When can you stop drawing samples from the posterior? That is the practitioners decision, and also dependent on the variance of the samples (recall from above a high variance means the average will converge slower).

We also should understand when the Law of Large Numbers fails. As the name implies, and comparing the graphs above for small  $N$ , the Law is only true for large sample sizes. Without this, the asymptotic result is not reliable. Knowing in what situations the Law fails can give use *confidence in how unconfident we should be*. The next section deals with this issue.

## 0.6.2 The Disorder of Small Numbers

The Law of Large Numbers is only valid as  $N$  gets *infinitely* large: never truly attainable. While the law is a powerful tool, it is foolhardy to apply it liberally. Our next example illustrates this.

**Example: Aggregated geographic data** Often data comes in aggregated form. For instance, data may be grouped by state, county, or city level. Of course, the population numbers vary per geographic area. If the data is an average of some characteristic of each the geographic areas, we must be conscious of the Law of Large Numbers and how it can *fail* for areas with small populations.

We will observe this on a toy dataset. Suppose there are five thousand counties in our dataset. Furthermore, population number in each state are uniformly distributed between 100 and 1500. The way the population numbers are generated is irrelevant to the discussion, so we do not justify this. We are interested in measuring the average height of individuals per county. Unbeknownst to the us, height does **not** vary across county, and each individual, regardless of the county he or she is currently living in, has the same distribution of what their height may be:

$$\text{height} \sim \text{Normal}(150, 15)$$

We aggregate the individuals at the county level, so we only have data for the *average in the county*. What might our dataset look like?

```

In [20]: figsize( 12.5, 4)
std_height = 15
mean_height = 150

n_counties = 5000
pop_generator = mc.rdiscrete_uniform
norm = mc.rnormal

#generate some artificial population numbers
population = pop_generator(100, 1500, size = n_counties )

average_across_county = np.zeros( n_counties )
for i in range( n_counties ):
    #generate some individuals and take the mean
    average_across_county[i] = norm(mean_height, 1./std_height**2,
                                   size=population[i] ).mean()

#where are the extreme populations?
i_min = np.argmin( average_across_county )
i_max = np.argmax( average_across_county )

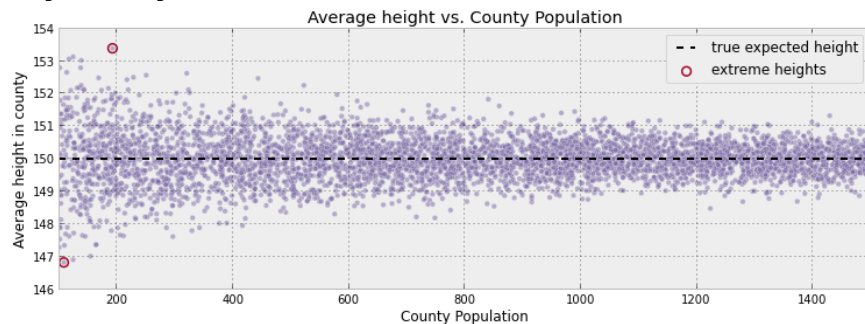
#plot population vs. average
plt.scatter( population, average_across_county, alpha = 0.5, c="#7A68A6" )
plt.scatter( [ population[i_min], population[i_max] ],
            [ average_across_county[i_min], average_across_county[i_max] ],
            s = 60, marker = "o", facecolors = "none",
            edgecolors = "#A60628", linewidths = 1.5,
            label="extreme heights" )

plt.xlim( 100, 1500 )
plt.title( "Average height vs. County Population" )
plt.xlabel( "County Population" )
plt.ylabel( "Average height in county" )
plt.plot( [100, 1500], [150, 150], color = "k", label = "true expected \
height", ls="--" )
plt.legend(scatterpoints = 1)

```

<matplotlib.legend.Legend at 0xc5cb030>

Out [20]:



What do we observe? *Without accounting for population sizes* we run the risk of making an enormous inference error: if we ignored population size, we would say that the county with the shortest and tallest individuals have been correctly circled. But this inference is wrong for the following reason. These two counties do *not* necessarily have the most extreme heights. The error is that the calculated average of the small population is not a good reflection of the true expected value of the population (which should be  $\mu = 150$ ). The sample size/population size/ $N$ , whatever you wish to call it, is simply too small to invoke the Law of Large Numbers effectively.

We provide more damning evidence against this inference. Recall the population numbers were uniformly distributed over 100 to 1500. Our intuition should tell us that the counties with the most extreme population heights should also be uniformly spread over 100 to 1500, and certainly independent of the county's population. Not so. Below are the population sizes of the counties with the most extreme heights.

```
In [11]: print "Population sizes of 10 'shortest' counties: "
print population[ np.argsort( average_across_county )[:10] ]
print
print "Population sizes of 10 'tallest' counties: "
print population[ np.argsort( -average_across_county )[:10] ]
```

Population sizes of 10 'shortest' counties:  
 [181 168 229 110 156 123 222 154 498 375]

Population sizes of 10 'tallest' counties:  
 [105 114 111 236 373 244 183 278 234 268]

Not at all uniform over 100 to 1500. This is an absolute failure of the Law of Large Numbers.

**Example: Kaggle's U.S. Census Return Rate Challenge** Below is data from the 2010 US census, which partitions populations beyond counties to the level of block groups (which are aggregates of city blocks or equivalents). The dataset is from a Kaggle machine learning competition some colleagues and I participated in. The objective was to predict the census letter mail-back rate of a group block, measured between 0 and 100, using census variables (median income, number of females in the block-group, number of trailer parks, average number of children etc.). Below we plot the census mail-back rate versus block group population:

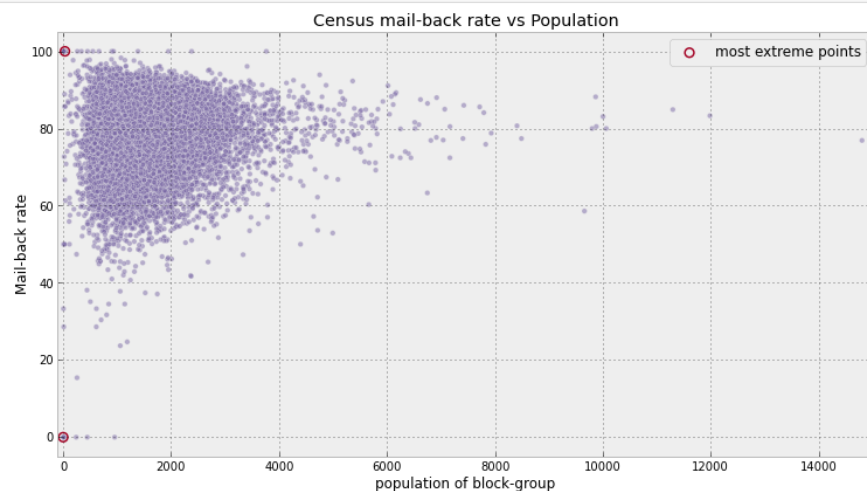
```
In [89]: figsize( 12.5, 6.5 )
data = np.genfromtxt( "./data/census_data.csv", skip_header=1,
                    delimiter= "," )

plt.scatter( data[:,1], data[:,0], alpha = 0.5, c="#7A68A6" )
plt.title("Census mail-back rate vs Population")
plt.ylabel("Mail-back rate")
plt.xlabel("population of block-group")
plt.xlim(-100, 15e3 )
plt.ylim( -5, 105)

i_min = np.argmin( data[:,0] )
i_max = np.argmax( data[:,0] )

plt.scatter( [ data[i_min,1], data[i_max, 1] ],
            [ data[i_min,0], data[i_max,0] ],
            s = 60, marker = "o", facecolors = "none",
            edgecolors = "#A60628", linewidths = 1.5,
            label="most extreme points")

plt.legend(scatterpoints = 1);
```



The above is a classic phenomenon in statistics. I say *classic* referring to the “shape” of the scatter plot above. It follows a classic triangular form, that tightens as we increase the sample size (as the Law of Large Numbers becomes more exact).

I am perhaps overstressing the point and maybe I should have titled the book “*You don’t have big data problems!*”, but here again is an example of the trouble with *small datasets*, not big ones. Simply, small datasets cannot be processed using the Law of Large Numbers. Compare with applying the Law without hassle to big datasets (ex. big data). I mentioned earlier that paradoxically big data prediction problems are solved by relatively simple algorithms. The paradox is partially resolved by understanding that the Law of Large Numbers creates solutions that are *stable*, i.e. adding or subtracting a few data points will not affect the solution much. On the other hand, adding or removing data points to a small dataset can create very different results.

For further reading on the hidden dangers of the Law of Large Numbers, I would highly recommend the excellent manuscript [The Most Dangerous Equation](#).

**Example: How Reddits ranks comments** You may have disagreed with the original statement that the Law of Large numbers is known to everyone, but only implicitly in our subconscious decision making. Consider ratings on online products: how often do you trust an average 5-star rating if there is only 1 reviewer? 2 reviewers? 3 reviewers? We implicitly understand that with such few reviewers that the average rating is **not** a good reflection of the true value of the product.

This has created flaws in how we sort items, and more generally, how we compare items. Many people have realized that sorting online search results by their rating, whether the objects be books, videos, or online comments, return poor results. Often the seemingly top videos or comments have perfect ratings only from a few enthusiastic fans, and truly more quality videos or comments are hidden in later pages with *falsely-substandard* ratings of around 4.8. How can we correct this?

Consider the popular site Reddit (I purposefully did not link to the website as you would never come back). The site hosts links to stories or images, and a very popular part of the site are the comments associated with each link. Redditors can vote up or down on each comment (called upvotes and downvotes). Reddit, by default, will sort comments by Top, that is, the best comments.

How would you determine which comments are the best? There are a number of ways to achieve this:

1. *Popularity*: A comment is considered good if it has many upvotes. A problem with this model is that a comment with hundreds of upvotes, but thousands of downvotes. While being very *popular*, the comment is likely more controversial than best.
2. *Difference*: Using the *difference* of upvotes and downvotes. This solves the above problem, but fails when we consider the temporal nature of comments. Comments can be posted many hours after the original link submission. The difference method will bias the *Top* comments to be the oldest comments, which have accumulated more upvotes than newer comments, but are not necessarily the best.
3. *Time adjusted*: Consider using Difference divided by the age of the comment. This creates a *rate*, something like *difference per second*, or *per minute*. An immediate counter example is, if we use per second, a 1 second old comment with 1 upvote would be better than a 100 second old comment with 99 upvotes. One can avoid this by only considering at least  $t$  second old comments. But what is a good  $t$  value? Does this mean no comment younger than  $t$  is good? We end up comparing unstable quantities with stable quantities (young vs. old comments).
4. *Ratio*: Rank comments by the ratio of upvotes to total number of votes (upvotes plus downvotes). This solves the temporal issue, such that new comments who score well can be considered Top just as likely as older comments, provided they have many upvotes to total votes. The problem here is that a comment with a single upvote (ratio = 1.0) will beat a comment with 999 upvotes and 1 downvote (ratio = 0.999), but clearly the latter comment is *more likely* to be better.

I used the phrase *more likely* for good reason. It is possible that the former comment, with a single upvote, is in fact a better comment than the later with 999 upvotes. The hesitation to agree with this is because we have not seen the other

999 potential votes the former comment might get. Perhaps it will achieve an additional 999 upvotes and 0 downvotes and be considered better than the latter, though not likely.

What we really want is an estimate of the *true upvote ratio*. Note that the true upvote ratio is not the same as the observed upvote ratio: the true upvote ratio is hidden, and we only observe upvotes vs. downvotes (one can think of the true upvote ratio as “what is the underlying probability someone gives this comment a upvote, versus a downvote”). So the 999 upvote/1 downvote comment probably has a true upvote ratio close to 1, which we can assert with confidence thanks to the Law of Large Numbers, but on the other hand we are much less certain about the true upvote ratio of the comment with only a single upvote. Sounds like a Bayesian problem to me. One way to determine a prior on the upvote ratio is that look at the historical distribution of upvote ratios. This can be accomplished by scrapping Reddit’s comments and determining a distribution. There are a few problems with this technique though:

1. Skewed data: The vast majority of comments have very few votes, hence there will be many comments with ratios near the extremes (see the “triangular plot” in the above Kaggle dataset), effectively skewing our distribution to the extremes. One could try to only use comments with votes greater than some threshold. Again, problems are encountered. There is a tradeoff between number of comments available to use and a higher threshold with associated ratio precision.
2. Biased data: Reddit is composed of different subpages, called subreddits. Two examples are *r/aww*, which posts pics of cute animals, and *r/politics*. It is very likely that the user behaviour towards comments of these two subreddits are very different: visitors are likely friend and affectionate in the former, and would therefore upvote comments more, compared to the latter, where comments are likely to be controversial and disagreed upon. Therefore not all comments are the same.

In light of these, I think it is better to use a Uniform prior.

With our prior in place, we can find the posterior of the true upvote ratio. The Python script `comments_for_top_reddit_pic.py` will scrap the comments from the current top picture on Reddit. Below is the picture, and some comments:

```
In [1]: from IPython.core.display import Image
        #adding a number to the end of the %run call with get the ith top photo.
        %run top_pic_comments.py 2

        Image(top_post_url)
```

Version 2.0.14 of praw is outdated. Version 2.0.15 was released Saturday April 06, 2013.

Title of submission:

3 Wolves in Quebec, Canada

<http://i.imgur.com/aHes4aF.jpg>

```
Out[43]: """
          contents: an array of the text from all comments on the pic
          votes: a 2d numpy array of upvotes, downvotes for each comment.
          """
          n_comments = len(contents )
          comments = np.random.randint( n_comments, size=4)
          print "Some Comments (out of %d total) \n-----"%n_comments
          for i in comments:
              print ''' + contents[i] + '''
              print "upvotes/downvotes: ",votes[i,:]
              print
```

Some Comments (out of 49 total)

```
-----
"Shaggydog, Nymeria, and Ghost?"
upvotes/downvotes: [134 17]
```

```
"I call bullshit. Not one of them is howling, nor is there a moon present."
upvotes/downvotes: [80 13]
```

```
"That would make a badass t shirt."  
upvotes/downvotes: [10 0]
```

```
"That is the alpha male and he is alerting his pack that he has found a Tim Hortons."  
upvotes/downvotes: [2 0]
```

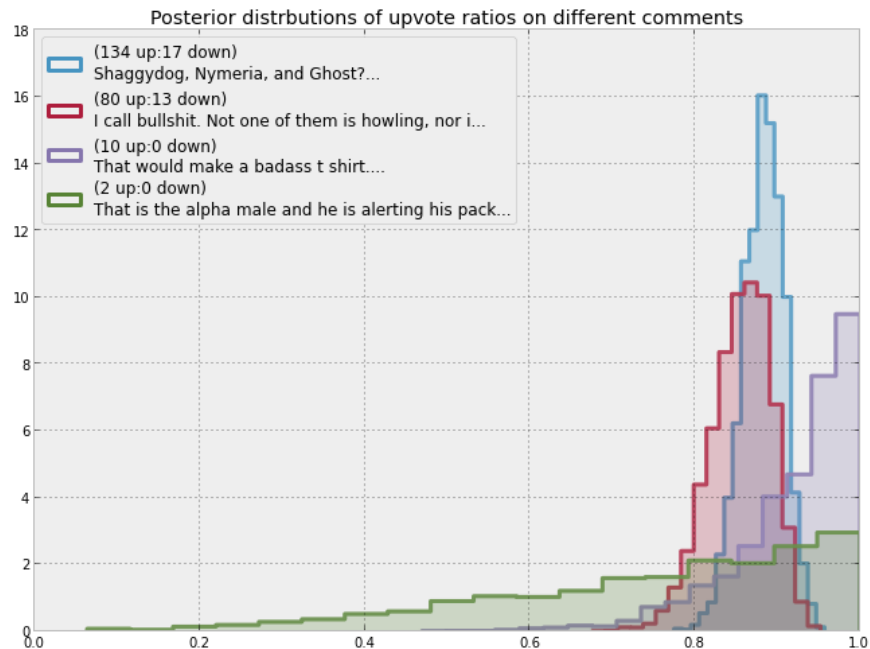
For a given true upvote ratio  $p$  and  $N$  votes, the number of upvotes will look like a Binomial random variable with parameters  $p$  and  $N$ . (This is because of the equivalence between upvote ratio and probability of upvoting versus downvoting, out of  $N$  possible votes/trials). We create a function that performs Bayesian inference on  $p$ , for a particular comment's upvote/downvote pair.

```
In [44]: import pymc as mc  
  
def posterior_upvote_ratio( upvotes, downvotes, samples = 20000):  
    N = upvotes + downvotes  
    upvote_ratio = mc.Uniform( "upvote_ratio", 0, 1 )  
    observations = mc.Binomial( "obs", N, upvote_ratio, value = upvotes,  
  
    model = mc.Model( [upvote_ratio, observations ] )  
    map_ = mc.MAP(model).fit()  
    mcmc = mc.MCMC(model)  
  
    mcmc.sample(samples, samples/4)  
  
    return mcmc.trace("upvote_ratio")[:]
```

Below are the resulting posterior distributions.

```
In [45]: figsize( 11., 8)  
posteriors = []  
colours = ["#348ABD", "#A60628", "#7A68A6", "#467821", "#CF4457"]  
for i in range(len(comments)):  
    j = comments[i]  
    posteriors.append( posterior_upvote_ratio( votes[j, 0], votes[j,1] ) )  
plt.hist( posteriors[i], bins = 18, normed = True, alpha = .9,  
          histtype="step", color = colours[i%5], lw = 3,  
          label = ' (%d up:%d down)\n%s...'%(votes[j, 0], votes[j,1], comments[i]) )  
plt.hist( posteriors[i], bins = 18, normed = True, alpha = .2,  
          histtype="stepfilled", color = colours[i], lw = 3, )  
#plt.title( 'upvotes:downvotes: %d:%d'%(votes[j,0], votes[j,1]) )  
  
plt.legend(loc="upper left")  
plt.xlim( 0, 1)  
plt.title("Posterior distributions of upvote ratios on different comments")
```

[\*\*\*\*\*100%\*\*\*\*\*] 20000 of 20000 complete



Some distributions are very tight, others have very long tails (relatively speaking), expressing our uncertainty with what the true upvote ratio might be.

## Sorting!

We have been ignoring the goal of this exercise: how do we sort the comments from *best to worst*? Of course, we cannot sort distributions, we must sort scalar numbers. There are many ways to distill a distribution down to a scalar: expressing the distribution through its expected value, or mean, is one way. Choosing the mean bad choice though. This is because the mean does not take into account the uncertainty of distributions.

I suggest using the *95% least plausible value*, defined as the value such that there is only a 5% chance the true parameter is lower (think of the lower bound on the 95% credible region). Below are the posterior distributions with the 95% least-plausible value plotted:

```
In [46]: N = posteriors[0].shape[0]
lower_limits = []

for i in range(len(comments)):
    j = comments[i]
    plt.hist( posteriors[i], bins = 20, normed = True, alpha = .9,
              histtype="step", color = colours[i], lw = 3,
              label = ' (%d up:%d down)\n%s...' % (votes[j, 0], votes[j,1], comments[i]) )
    plt.hist( posteriors[i], bins = 20, normed = True, alpha = .2,
              histtype="stepfilled", color = colours[i], lw = 3, )
    v = np.sort( posteriors[i] ) [ int(0.05*N) ]
    #plt.vlines( v, 0, 15 , color = "k", alpha = 1, linewidths=3 )
    plt.vlines( v, 0, 10 , color = colours[i], linestyle = "--", linewidth=3 )
    lower_limits.append(v)
    plt.legend(loc="upper left")

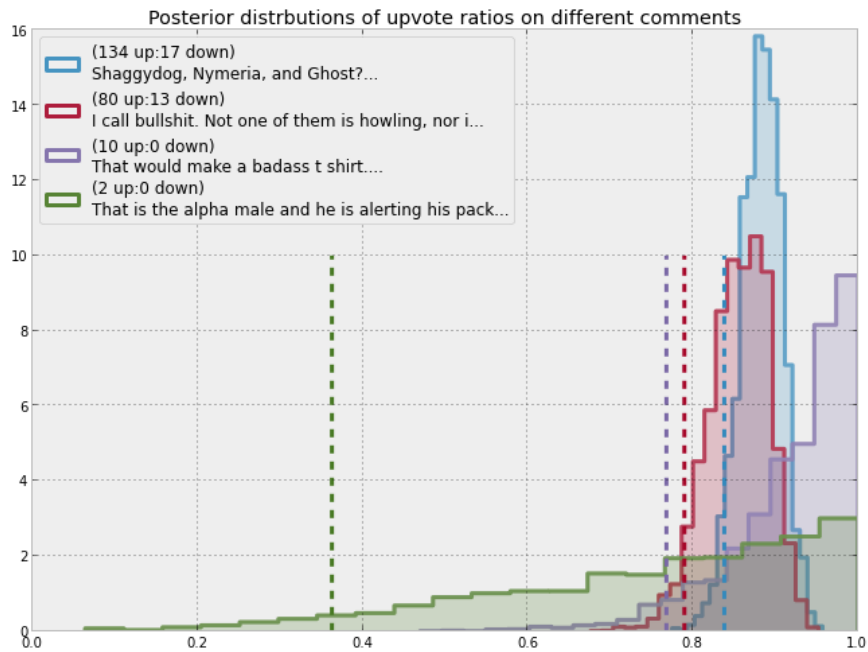
plt.legend(loc="upper left")

plt.title("Posterior distributions of upvote ratios on different comments")
order = argsort( -np.array( lower_limits ) )
```



```
print order, lower_limits
```

```
[0 1 2 3] [0.83778530918250105, 0.78945340111675133, 0.7676689367397137, 0.362408958352504]
```



The best comments, according to our procedure, are the comments that are *most-likely* to score a high percentage of upvotes. Visually those are the comments with the 95% least plausible value close to 1.

Why is sorting based on this quantity a good idea? By ordering by the 95% least plausible value, we are being the most conservative with what we think is best. That is, even in the worst case scenario, when we have severely overestimated the upvote ratio, we can be sure the best comments are still on top. Under this ordering, we impose the following very natural properties:

1. given two comments with the same observed upvote ratio, we will assign the comment with more votes as better (since we are more confident it has a higher ratio).
2. given two comments with the same number of votes, we still assign the comment with more upvotes as *better*.

### But this is too slow for real-time!

I agree, computing the posterior of every comment takes a long time, and by the time you have computed it, likely the data has changed. I delay the mathematics to the appendix, but I suggest using the following formula to compute the lower bound very fast.

$$\frac{a}{a+b} - 1.65 \sqrt{\frac{ab}{(a+b)^2(a+b+1)}}$$

where

$$a = 1 + u \tag{36}$$

$$\tag{37}$$

$$b = 1 + d \tag{38}$$

$$\tag{39}$$

$$\tag{40}$$

$u$  is the number of upvotes, and  $d$  is the number of downvotes. The formula is a shortcut in Bayesian inference, which will be further explained in Chapter 6 when we discuss priors in more detail.

```
In [52]: def intervals(u,d):
          a = 1. + u
          b = 1. + d
          mu = a/(a+b)
          std_err = 1.65*np.sqrt( (a*b)/( (a+b)**2*(a+b+1.) ) )
          return ( mu, std_err )

          print "Approximate lower bounds:"
          posterior_mean, std_err = intervals(votes[:,0],votes[:,1])
          lb = posterior_mean - std_err
          print lb
          print

          print "Top 40 Sorted according to approximate lower bounds:"
          print
          order = np.argsort( -lb )
          ordered_contents = []
          for i in order[:40]:
              ordered_contents.append( contents[i] )
              print votes[i,0], votes[i,1], contents[i]
              print "-----"
```

Approximate lower bounds:

```
[ 0.85943711  0.83951434  0.8458571  0.84536604  0.86388258  0.83651439
  0.89239065  0.7929375  0.85447643  0.70806405  0.79018506  0.82545948
  0.70408053  0.79198214  0.60100251  0.6931046  0.6931046  0.6530083
  0.7137931  0.60091591  0.60091591  0.60091591  0.60091591  0.60091591
  0.52182581  0.60100251  0.53055613  0.53055613  0.53055613  0.56085485
  0.51184301  0.4722123  0.43047887  0.43047887  0.43047887  0.53055613
  0.43047887  0.43047887  0.43047887  0.53055613  0.43047887  0.43047887
  0.43047887  0.43047887  0.43047887  0.43047887  0.43047887  0.43047887
  0.43047887]
```

Top 40 Sorted according to approximate lower bounds:

22 0 Being as an ocean?

-----

354 43 Ok, three-wolf-joke-in-my-second-language time:

Un jours trois loups chassent ensemble dans la fort: un Qubcois, un Franco-Ontarien et un P

"Oh ben \*tabarnak\*" dit le Franco-Ontarien. "L on fait quoi, crise? Les chasseurs vont nou

"J'ai une ide!" exclame le Qubcois. "Si on s'mange une patte on pourrait nous librer. C'est

"Osti d'sapp," rpond le Qubcois, "l on va avoir le sauver!" Les deux retournent sur le lie

"Bordel, vous avez fait comment les mecs? Je m'suis dj bouff trois pattes et je suis encor

-----  
362 46 This was done by photographer Daniel Parent. [[Source] (<http://500px.com/photo/75056>)

Here are some of my favorites from the collection:

- OP's one, called ["Three Weary Wolves"] (<http://i.imgur.com/dW2bsks.jpg>) [707 x 900 Versi

- ["The Crossing" (with 12 wolves)] (<http://i.imgur.com/J5bKJsu.jpg>)

- ["Above and Beyond"] (<http://i.imgur.com/bJHbrNa.jpg>)

- ["On Patrol"] (<http://i.imgur.com/TKyST9h.jpg>)

-----  
44 3 Are they also in a Christian hard rock band?

-----  
33 2 Princess Mononoke?

-----  
805 126

```
      /      -'      /      /      /      /      /      /      /
     _/|    '- _ /    /-_/    ~' |    /      /
    /      ". |    /|      /      /      /      /
   |      ;    T / |      ' ,    |      /
  -      -      -      -      -      -      -      -      -
   /      /      /      /      /      /      /      /
  /      /      /      /      /      /      /      /
 /      /      /      /      /      /      /      /
-      -      -      -      -      -      -      -      -
      /      /      /      /      /      /      /      /
     /      /      /      /      /      /      /      /
    /      /      /      /      /      /      /      /
   /      /      /      /      /      /      /      /
  /      /      /      /      /      /      /      /
 /      /      /      /      /      /      /      /
/      /      /      /      /      /      /      /
(____.'_____)_|_|
```

-----  
134 17 Shaggydog, Nymeria, and Ghost?

-----  
106 13 [Was this picture taken in Rivire-du-Loup?] (<http://www.badumtss.net/>)

-----  
21 1 Dire Wolves. Fucking sexy ass Dire Wolves.

-----  
80 13 I call bullshit. Not one of them is howling, nor is there a moon present.

-----  
17 1 Yayy Quebec on the front page :)

-----  
10 0 That would make a badass t shirt.

-----  
23 4 where's moon moon?

-----  
11 1 Good Ol' Quebec, Canada. Land of wolves, maple syrup and Labatt 50.

-----  
22 4 No retarded comment from me. This is a nice picture!

-----

6 0 I recognize this from a Being As An Ocean band t-shirt.

Such a cool picture.

6 0 Looks like a boy band posing for their new album cover.

5 0 Is that Wolf+1's new album cover?

7 1 On a tu pris ces photos Rivire-du-Loup ?? :-)

7 1 Mon tabarnac c'est dont ben des beaux loups esti!

4 0 This image has actually been reprinted into a REALLY nice shirt by the band Being As An

<http://merchnow.com/products/149014/alpine-wolf-white>

4 0 where, specifically in Quebec was that photo taken?

4 0 I like this one, one wolf is facing one way another facing the other way and the wolf is

4 0 reminds me of princess mononoke

4 0 [As a Magic: The Gathering card] (<http://i.imgur.com/KI74xCc.jpg>)

6 1 Winter is coming.

3 0 I suppose this is the time to post [this] (<http://www.reddit.com/r/pics/comments/1dlqs4>)

Me, in Quebec, joking around with some wolves.

3 0 Looks like the Christ + 1 album cover.

3 0 Where exactly was it taken?

3 0 J'adore!

3 0 Looks like Cartman tried to make an christian rock album cover

9 3 Did Eric Cartman stage this for an album cover?

5 1 [3 Wolf River] (<http://i.imgur.com/W01N0G4.png>) saved my marriage! I used to wear boring

Cue 3 Wolf River: this magnificent fucking shirt has not one, NOT TWO, BUT \*THREE\* stunning

The amount of sheer manliness contained within these threads was not meant for mortal men,

9 4 I hope those Stark kids pick up after them.

2 0 Dude i love how quebec looks, lived there nearly my entire life

2 0 3 wolf canada...

2 0 MAJESTIC AS FUCK!

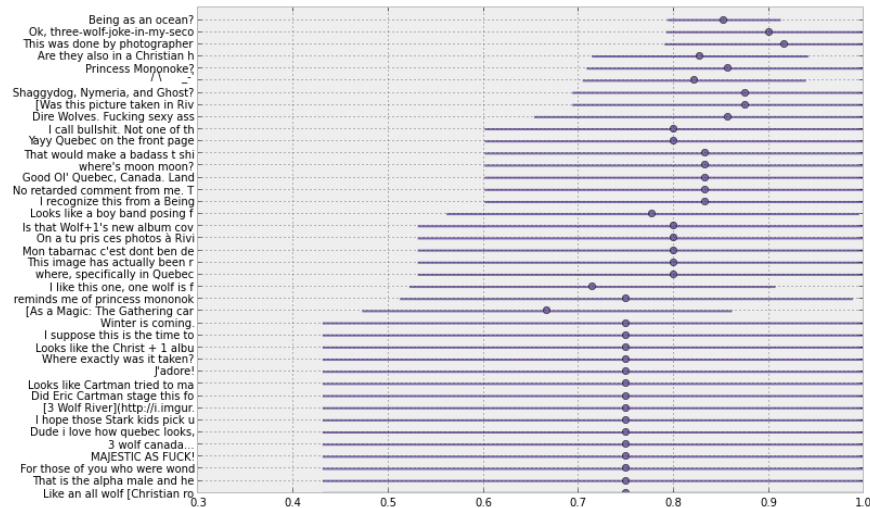
-----  
2 0 For those of you who were wondering where those pictures came from, a lot of them were  
-----

2 0 That is the alpha male and he is alerting his pack that he has found a Tim Hortons.  
-----

2 0 Like an all wolf [Christian rock band] (<http://imgur.com/9yR3V6R.jpg>)  
-----

We can view the ordering visually by plotting the posterior mean and bounds, and sorting by the lower bound. In the plot below, notice that the left error-bar is sorted (as we suggested this is the best way to determine an ordering), so the means, indicated by dots, do not follow any strong pattern.

```
In [53]: r_order = order[::-1][:40]
plt.errorbar( posterior_mean[r_order], np.arange( len(r_order) ),
             xerr=std_err[r_order], xuplims=True, capsize=0, fmt="o",
             color = "#7A68A6" )
plt.xlim( 0.3, 1)
plt.yticks( np.arange( len(r_order)-1, -1, -1 ), map( lambda x: x[:30], orde
```



In the graphic above, you can see why sorting by mean would be sub-optimal.

## Conclusion

While the Law of Large Numbers is cool, it is only true so much as its name implies: with large sample sizes only. We have seen how our inference can be affected by not considering *how the data is shaped*.

1. By (cheaply) drawing many samples from the posterior distributions, we can ensure that the Law of Large Number applies as we approximate expected values (which we will do in the next chapter).
2. Bayesian inference understands that with small sample sizes, we can observe wild randomness. Our posterior distribution will reflect this by being more spread rather than tightly concentrated. Thus, our inference should be correctable.
3. There are major implications of not considering the sample size, and trying to sort objects that are unstable leads to pathological orderings. The method provided above solves this problem.

## Appendix

**Derivation of sorting comments formula** Basically what we are doing is using a Beta prior (with parameters  $a = 1, b = 1$ , which is a uniform distribution), and using a Binomial likelihood with observations  $u, N = u + d$ . This means our posterior is a Beta distribution with parameters  $a' = 1 + u, b' = 1 + (N - u) = 1 + d$ . We then need to find the value,  $x$ , such that 0.05 probability is less than  $x$ . This is usually done by inverting the CDF (Cumulative Distribution Function), but the CDF of the beta, for integer parameters, is known but is a large sum [3].

We instead using a Normal approximation. The mean of the Beta is  $\mu = a'/(a' + b')$  and the variance is

$$\sigma^2 = \frac{a'b'}{(a' + b')^2(a' + b' + 1)}$$

Hence we solve the following equation for  $x$  and have an approximate lower bound.

$$0.05 = \Phi\left(\frac{(x - \mu)}{\sigma}\right)$$

$\Phi$  being the cumulative distribution for the normal distribution

**Exercises** 1. How would you estimate the quantity  $E[\cos X]$ , where  $X \sim \text{Exp}(4)$ ? What about  $E[\cos X|X < 1]$ , i.e. the expected value *given* we know  $X$  is less than 1? Would you need more samples than the original samples size to be equally as accurate?

```
In [14]: ## Enter code here
import scipy.stats as stats
exp = stats.expon( scale=4 )
N = 1e5
X = exp.rvs( N )
## ...
```

2. The following table was located in the paper “Going for Three: Predicting the Likelihood of Field Goal Success with Logistic Regression” [2]. The table ranks football field-goal kickers by there percent of non-misses. What mistake have the researchers made?

---

### Kicker Careers Ranked by Make Percentage

Rank
Kicker
Make %
Number of Kicks
1
Garrett Hartley
87.7
57
2
Matt Stover

86.8

335

3

Robbie Gould

86.2

224

4

Rob Bironas

86.1

223

5

Shayne Graham

85.4

254

...

...

...

51

Dave Rayner

72.2

90

52

Nick Novak

71.9

64

53

Tim Seder

71.0

62

54

Jose Cortez

70.7

75

55

Wade Richey

66.1

## References

1. Wainer, Howard. *The Most Dangerous Equation*. American Scientist, Volume 95.
2. Clark, Torin K., Aaron W. Johnson, and Alexander J. Stimpson. "Going for Three: Predicting the Likelihood of Field Goal Success with Logistic Regression." (2013): n. page. Web. 20 Feb. 2013.
3. [http://en.wikipedia.org/wiki/Beta\\_function#Incomplete\\_beta\\_function](http://en.wikipedia.org/wiki/Beta_function#Incomplete_beta_function)

```
In [1]: from IPython.core.display import HTML
def css_styling():
    styles = open("../styles/custom.css", "r").read()
    return HTML(styles)
css_styling()
```

Out [1]:

## 0.7 Chapter 5

---

### 0.7.1 Would you rather lose an arm or a leg?

Statisticians can be a sour bunch. Instead of considering their winnings, they only measure how much they have lost. In fact, they consider their wins as *negative losses*. But what's interesting is *how they measure their losses*.

For example, consider the following example:

A meteorologist is predicting the probability of a possible hurricane striking his city. He estimates, with 95% confidence, that the probability of it *not* striking is between 99% - 100%. He is very happy with his precision and advises the city that a major evacuation is unnecessary. Unfortunately, the hurricane does strike and the city is flooded.

This stylized example shows the flaw in using a pure accuracy metric to measure outcomes. Using a measure that emphasizes estimation accuracy, while an appealing and *objective* thing to do, misses the point of why you are even performing the statistical inference in the first place: results of inference. The author Nassim Taleb of *The Black Swan* and *Antifragility* stresses the importance of the *payoffs* of decisions, *not the accuracy*. Taleb distills this quite succinctly: "I would rather be vaguely right than very wrong."

### Loss Functions

We introduce what statisticians and decision theorists call *loss functions*. A loss function is a function of the true parameter, and an estimate of that parameter

$$L(\theta, \hat{\theta}) = f(\theta, \hat{\theta})$$

The important point of loss functions is that it measures how *bad* our current estimate is: the larger the loss, the worse the estimate is according to the loss function. A simple, and very common, example of a loss function is the *squared-error loss*:

$$L(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$$



The squared-error loss function is used in estimators like linear regression, UMVUEs and many areas of machine learning. We can also consider an asymmetric squared-error loss function, something like:

$$L(\theta, \hat{\theta}) = \begin{cases} (\theta - \hat{\theta})^2 & \hat{\theta} < \theta \\ c(\theta - \hat{\theta})^2 & \hat{\theta} \geq \theta, \quad 0 < c < 1 \end{cases}$$

which represents that estimating a value larger than the true estimate is preferable to estimating a value below. A situation where this might be useful is in estimating web traffic for the next month, where an over-estimated outlook is preferred so as to avoid an underallocation of server resources.

A negative property about the squared-error loss is that it puts a disproportionate emphasis on large outliers. This is because the loss increases quadratically, and not linearly, as the estimate moves away. That is, the penalty of being three units away is much less than being five units away, but the penalty is not much greater than being one unit away, though in both cases the magnitude of difference is the same:

$$\frac{1^2}{3^2} < \frac{3^2}{5^2}, \quad \text{although } 3 - 1 = 5 - 3$$

This loss function imposes that large errors are *very* bad. A more *robust* loss function that increases linearly with the difference is the *absolute-loss*

$$L(\theta, \hat{\theta}) = |\theta - \hat{\theta}|$$

Other popular loss functions include:

- $L(\theta, \hat{\theta}) = \mathbb{1}_{\hat{\theta} \neq \theta}$  is the zero-one loss often used in machine learning classification algorithms.
- $L(\theta, \hat{\theta}) = -\hat{\theta} \log(\theta) - (1 - \hat{\theta}) \log(1 - \theta)$ , ; ;  $\hat{\theta} \in [0, 1]$ , ;  $\theta \in [0, 1]$ , called the *log-loss*, also used in machine learning.

Historically, loss functions have been motivated from 1) mathematical convenience, and 2) they are robust to application, i.e., they are objective measures of loss. The first reason has really held back the full breadth of loss functions. With computers being agnostic to mathematical convenience, we are free to design our own loss functions, which we take full advantage of later in this Chapter.

With respect to the second point, the above loss functions are indeed objective, in that they are most often a function of the difference between estimate and true parameter, independent of signage or payoff of choosing that estimate. This last point, its independence of payoff, causes quite pathological results though. Consider our hurricane example above: the statistician equivalently predicted that the probability of the hurricane striking was between 0% to 1%. But if he had ignored being precise and instead focused on outcomes (99% chance of no flood, 1% chance of flood), he might have advised differently.

By shifting our focus from trying to be incredibly precise about parameter estimation to focusing on the outcomes of our parameter estimation, we can customize our estimates to be optimized for our application. This requires us to design new loss functions that reflect our goals and outcomes. Some examples of more interesting loss functions:

- $L(\theta, \hat{\theta}) = \frac{|\theta - \hat{\theta}|}{\theta(1 - \theta)}$ , ; ;  $\hat{\theta}, \theta \in [0, 1]$  emphasizes an estimate closer to 0 or 1 since if the true value  $\theta$  is near 0 or 1, the loss will be *very* large unless  $\hat{\theta}$  is similarly close to 0 or 1. This loss function might be used by a political pundit who's job requires him or her to give confident "Yes/No" answers. This loss reflects that if the true parameter is close to 1 (for example, if a political outcome is very likely to occur), he or she would want to strongly agree as to not look like a skeptic.
- $L(\theta, \hat{\theta}) = 1 - \exp(-(\theta - \hat{\theta})^2)$  is bounded between 0 and 1 and reflects that the user is indifferent to sufficiently-far-away estimates. It is similar to the zero-one loss above, but not quite as penalizing to estimates that are close to the true parameter.

- Complicated non-linear loss functions can be programmed:

```
def loss(true_value, estimate):
    if estimate*true_value > 0:
        return abs(estimate - true_value)
    else:
        return abs(estimate)*(estimate - true_value)**2
```

- Another example is from the book *The Signal and The Noise*. Weather forecasters have an interesting loss function for their predictions. [2]

People notice one type of mistake — the failure to predict rain — more than other, false alarms. If it rains when it isn't supposed to, they curse the weatherman for ruining their picnic, whereas an unexpectedly sunny day is taken as a serendipitous bonus.

[The Weather Channel's bias] is limited to slightly exaggerating the probability of rain when it is unlikely to occur — saying there is a 20 percent chance when they know it is really a 5 or 10 percent chance — covering their butts in the case of an unexpected sprinkle.

As you can see, loss functions can be used for good and evil: with great power, comes great — well you know.

## 0.7.2 Loss functions in the real world

So far we have been under the unrealistic assumption that we know the true parameter. Of course if we knew the true parameter, bothering to guess an estimate is pointless. Hence a loss function is really only practical when the true parameter is unknown.

In Bayesian inference, we have a mindset that the unknown parameters are really random variables with prior and posterior distributions. Concerning the posterior distribution, a value drawn from it is a *possible* realization of what the true parameter could be. Given that realization, we can compute a loss associated with an estimate. As we have a whole distribution of what the unknown parameter could be (the posterior), we should be more interested in computing the *expected loss* given an estimate. This expected loss is a better estimate of the true loss than comparing the given loss from only a single sample from the posterior.

First it will be useful to explain a *Bayesian point estimate*. The systems and machinery present in the modern world are not built to accept posterior distributions as input. It is also rude to hand someone over a distribution when all that is asked for is an estimate. In the course of an individual's day, when faced with uncertainty we still act by distilling our uncertainty down to a single action. Similarly, we need to distill our posterior distribution down to a single value (or vector in the multivariate case). If the value is chosen intelligently, we can avoid the flaw of frequentist methodologies that mask the uncertainty and provide a more informative result. The value chosen, if from a Bayesian posterior, is a Bayesian point estimate.

Suppose  $P(\theta|X)$  is the posterior distribution of  $\theta$  after observing data  $X$ , then the following function is understandable as the *expected loss of choosing estimate  $\hat{\theta}$  to estimate  $\theta$* :

$$l(\hat{\theta}) = E_{\theta} \left[ L(\theta, \hat{\theta}) \right]$$

This is also known as the *risk* of estimate  $\hat{\theta}$ . The subscript  $\theta$  under the expectation symbol is used to denote that  $\theta$  is the unknown (random) variable in the expectation, something that at first can be difficult to consider.

We spent all of last chapter discussing how to approximate expected values. Given  $N$  samples  $\theta_i$ ,  $i = 1, \dots, N$  from the posterior distribution, and a loss function  $L$ , we can approximate the expected loss of using estimate  $\hat{\theta}$  by the Law of Large Numbers:

$$\frac{1}{N} \sum_{i=1}^N L(\theta_i, \hat{\theta}) \approx E_{\theta} \left[ L(\theta, \hat{\theta}) \right] = l(\hat{\theta})$$

Notice that measuring your loss via an *expected value* uses more information from the distribution than the MAP estimate which, if you recall, will only find the maximum value of the distribution and ignore the shape of the distribution. Ignoring information can over-expose yourself to tail risks, like the unlikely hurricane, and leaves your estimate ignorant of how ignorant you really are about the parameter.

Similarly, compare this with frequentist methods, that traditionally only aim to minimize the error, and not considering the *loss associated with the result of that error*. Compound this with the fact that frequentist methods are almost guaranteed to never be absolutely accurate. Bayesian point estimates fix this by planning ahead: your estimate is going to be wrong, you might as well err on the right side of wrong.

**Example: Optimizing for the Showcase on The Price is Right** Bless you if you are ever chosen as a contestant on the Price is Right, for here we will show you how to optimize your final price on the *Showcase*. For those who forget the rules:

1. Two contestants compete in *The Showcase*.
2. Each contestant is shown a unique suite of prizes.
3. After the viewing, the contestants are asked to bid on the price for their unique suite of prizes.
4. If a bid price is over the actual price, the bid's owner is disqualified from winning.
5. If a bid price is under the true price by less than \$250, the winner is awarded both prizes.

The difficulty in the game is balancing your uncertainty in the prices, keeping your bid low enough so as to not bid over, and trying to bid close to the price.

Suppose we have recorded the *Showcases* from previous *The Price is Right* episodes and have *prior* beliefs about what distribution the true price follows. For simplicity, suppose it follows a Normal:

$$\text{True Price} \sim \text{Normal}(\mu_p, \sigma_p)$$

In a later chapter, we will actually use *real Price is Right Showcase data* to form the historical prior, but this requires some advanced PyMC use so we will not use it here. For now, we will assume  $\mu_p = 35000$  and  $\sigma_p = 7500$ .

We need a model of how we should be playing the *Showcase*. For each prize in the prize suite, we have an idea of what it might cost, but this guess could differ significantly from the true price. (Couple this with increased pressure being onstage and you can see why some bids are so wildly off). Let's suppose your beliefs about the prices of prizes also follow Normal distributions:

$$\text{Prize}_i \sim \text{Normal}(\mu_i, \sigma_i), \quad i = 1, 2$$

This is really why Bayesian analysis is great: we can specify what we think a fair price is through the  $\mu_i$  parameter, and express uncertainty of our guess in the  $\sigma_i$  parameter.

We'll assume two prizes per suite for brevity, but this can be extended to any number. The true price of the prize suite is then given by  $\text{Prize}_1 + \text{Prize}_2 + \epsilon$ , where  $\epsilon$  is some error term.

We are interested in the updated True Price given we have observed both prizes and have belief distributions about them. We can perform this using PyMC.

Lets make some values concrete. Suppose there are two prizes in the observed prize suite:

1. A trip to wonderful Toronto, Canada!
2. A lovely new snowblower!

We have some guesses about the true prices of these objects, but we are also pretty uncertain about them. I can express this uncertainty through the parameters of the Normals:

snowblower  $\sim$  Normal(3000, 500) (41)

(42)

Toronto  $\sim$  Normal(12000, 3000) (43)

(44)

(45)

For example, I believe that the true price of the trip to Toronto is 12 000 dollars, and that there is a 68.2% chance the price falls 1 standard deviation away from this, i.e. my confidence is that there is a 68.2% chance the snowblower is in [9 000, 15 000].

We can create some PyMC code to perform inference on the true price of the suite.

```
In [2]: import scipy.stats as stats

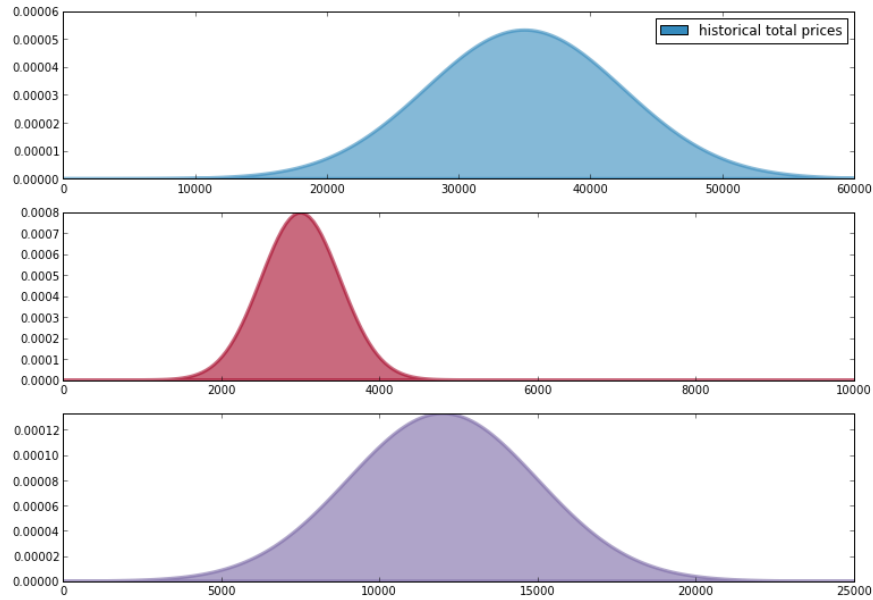
figsize( 12.5, 9 )

norm_pdf = stats.norm.pdf

subplot(311)
x = np.linspace(0, 60000, 200)
sp1 = plt.fill_between(x, 0, norm_pdf(x, 35000, 7500),
                      color = "#348ABD", lw = 3, alpha = 0.6,
                      label = "historical total prices")
p1 = Rectangle((0, 0), 1, 1, fc=sp1.get_facecolor()[0])
legend([p1], [sp1.get_label()])

subplot(312)
x = np.linspace(0, 10000, 200)
sp2 = plt.fill_between(x, 0, norm_pdf(x, 3000, 500),
                      color = "#A60628", lw = 3, alpha = 0.6,
                      label="snowblower price guess")
p2 = Rectangle((0, 0), 1, 1, fc=sp2.get_facecolor()[0])
legend([p2], [sp2.get_label()])

subplot(313)
x = np.linspace(0, 25000, 200)
sp3 = plt.fill_between(x, 0, norm_pdf(x, 12000, 3000),
                      color = "#7A68A6", lw = 3, alpha = 0.6,
                      label = "Trip price guess")
plt.autoscale(tight=True)
p3 = Rectangle((0, 0), 1, 1, fc=sp3.get_facecolor()[0])
legend([p3], [sp3.get_label()])
```



```
In [1]: %pylab inline
import pymc as mc

data_mu = [ 3e3, 12e3]
data_std = [ 5e2, 3e3 ]

mu_prior = 35e3
std_prior = 75e2

true_price = mc.Normal( "true_price", mu_prior, 1.0/std_prior**2 )

prize_1 = mc.Normal( "first_prize", data_mu[0], 1.0/data_std[0]**2 )
prize_2 = mc.Normal( "second_prize", data_mu[1], 1.0/data_std[1]**2 )
price_estimate = prize_1 + prize_2

@mc.potential
def error( true_price = true_price, price_estimate = price_estimate ):
    return mc.normal_like( true_price, price_estimate, 1/(3e3)**2 )

model = mc.Model([true_price, prize_1, prize_2, price_estimate, error ])
mcmc = mc.MCMC(model)
mcmc.sample( 50000, 10000)

price_trace = mcmc.trace( "true_price" )[:]
```

[\*\*\*\*\*100%\*\*\*\*\*] 50000 of 50000 complete

```
In [28]: figsize( 12.5, 4)

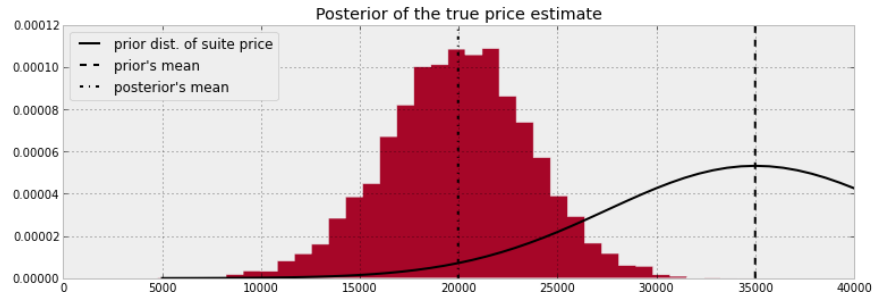
import scipy.stats as stats

x = np.linspace( 5000, 40000 )
plt.plot( x, stats.norm.pdf(x, 35000, 7500), c = "k", lw = 2,
         label = "prior dist. of suite price")

_hist = plt.hist( price_trace, bins = 35, normed= True, histtype= "stepfilled")
plt.title( "Posterior of the true price estimate" )
```

```
plt.vlines( mu_prior, 0, 1.1*np.max(_hist[0] ), label = "prior's mean",
           linestyle="--" )
plt.vlines( price_trace.mean(), 0, 1.1*np.max(_hist[0] ), \
           label = "posterior's mean", linestyle="-.")
plt.legend(loc = "upper left")
```

<matplotlib.legend.Legend at 0x13508978>  
Out [28]:



Notice that because of our two observed prizes and subsequent guesses (including uncertainty about those guesses), we shifted our mean price estimate down about \$15 000 dollars from the previous mean price.

A frequentist, seeing the two prizes and having the same beliefs about their prices, would bid  $\mu_1 + \mu_2 = 35000$ , regardless of any uncertainty. Meanwhile, the *naive Bayesian* would simply pick the mean of the posterior distribution. But we have more information about our eventual outcomes; we should incorporate this into our bid. We will use the loss function above to find the *best bid* (*best* according to our loss).

What might a contestant's loss function look like? I would think it would look something like:

```
def showcase_loss( guess, true_price, risk = 80000):
    if true_price < guess:
        return risk
    elif abs( true_price - guess ) <= 250:
        return -2*np.abs( true_price )
    else:
        return np.abs( true_price - guess - 250 )
```

where `risk` is a parameter that defines of how bad it is if your guess is over the true price. A lower `risk` means that you are more comfortable with the idea of going over. If we do bid under and the difference is less than \$250, we receive both prizes (modeled here as receiving twice the original prize). Otherwise, when we bid under the `true_price` we want to be as close as possible, hence the `else` loss is a increasing function of the distance between the guess and true price. For every possible bid, we calculate the *expected loss* associated with that bid. We vary the `risk` parameter to see how it affects our loss:

```
In [26]: figsize( 12.5, 7 )
         #numpy friendly showdown_loss
         def showdown_loss( guess, true_price, risk = 80000):
             loss = np.zeros_like( true_price )
             ix = true_price < guess
             loss[~ix] = np.abs( guess - true_price[~ix] )
             close_mask = [ abs( true_price - guess ) <= 250 ]
             loss[close_mask] = -2*true_price[close_mask]
             loss[ix] = risk
             return loss

         guesses = np.linspace( 5000, 50000, 70 )
         risks = np.linspace( 30000, 150000, 6 )
         expected_loss = lambda guess, risk: \
             showdown_loss( guess, price_trace, risk ).mean()

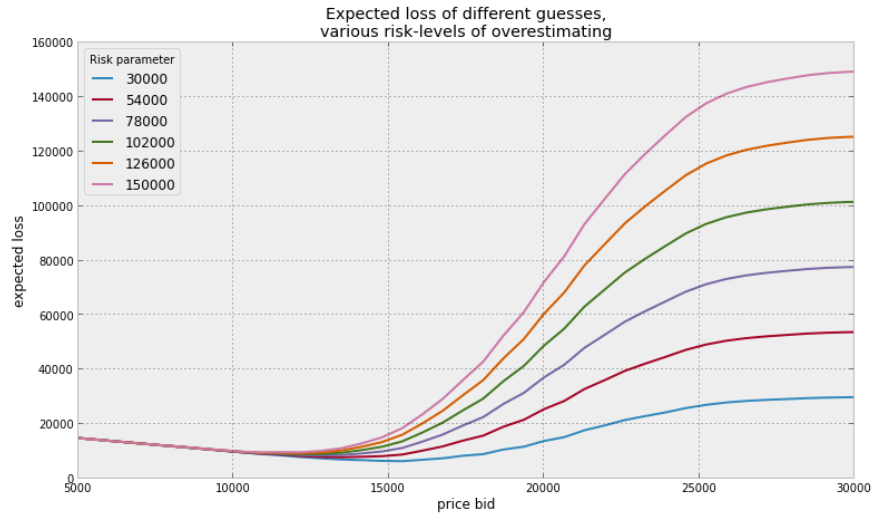
         for _p in risks:
             results = [ expected_loss( _g, _p ) for _g in guesses ]
```

```

plt.plot( guesses, results, label = "%d"%_p )

plt.title("Expected loss of different guesses, \nvarious risk-levels of \
overestimating")
plt.legend( loc="upper left", title="Risk parameter")
plt.xlabel("price bid")
plt.ylabel("expected loss")
plt.xlim( 5000, 30000 );

```



## Minimizing our losses

It would be wise to choose the estimate that minimizes our expected loss. This corresponds to the minimum point on each of the curves above. More formally, we would like to minimize our expected loss by finding the solution to

$$\arg \min_{\hat{\theta}} E_{\theta} [ L(\theta, \hat{\theta}) ]$$

The minimum of the expected loss is called the *Bayes action*. We can solve for the Bayes action using Scipy's optimization routines. The function `fmin` in `scipy.optimize` module uses an intelligent search to find a minimum (not necessarily a *global* minimum) of any uni- or multivariate function. For most purposes, `fmin` will provide you with a good answer.

We'll compute the minimum loss for the *Showcase* example above:

```

In [28]: import scipy.optimize as sop

ax = subplot(111)

for _p in risks:
    _color = ax._get_lines.color_cycle.next()
    _min_results = sop.fmin( expected_loss, 15000, args=( _p, ), disp = False )
    _results = [ expected_loss( _g, _p ) for _g in guesses ]
    plt.plot( guesses, _results, color = _color )
    plt.scatter( _min_results, 0, s = 60, \
                color= _color, label = "%d"%_p )
    plt.vlines( _min_results, 0, 120000, color = _color, linestyle="--" )
    print "minimum at risk %d: %.2f"% ( _p, _min_results )

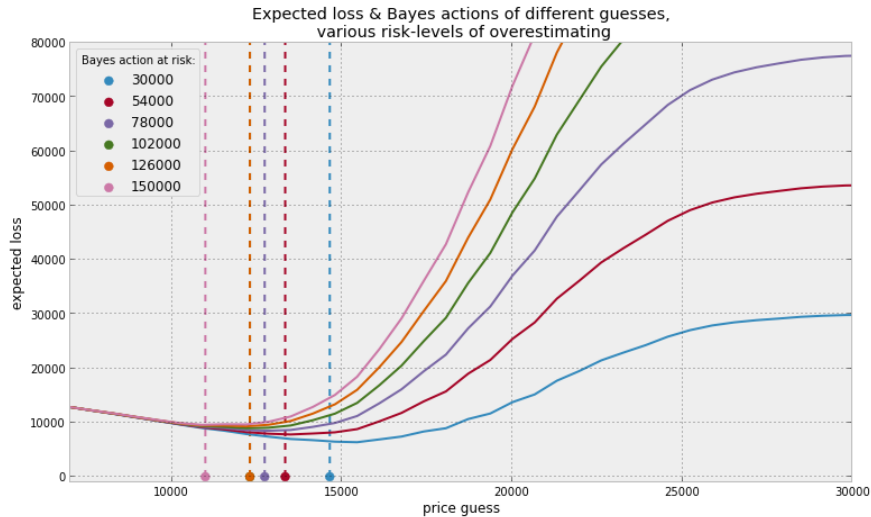
plt.title("Expected loss & Bayes actions of different guesses, \n \
various risk-levels of overestimating")
plt.legend( loc="upper left", scatterpoints = 1, title = "Bayes action at

```

```
plt.xlabel("price guess")
plt.ylabel("expected loss")
plt.xlim( 7000, 30000 )
plt.ylim( -1000, 80000)
```

```
minimum at risk 30000: 14629.68
minimum at risk 54000: 13319.39
minimum at risk 78000: 12716.49
minimum at risk 102000: 12281.71
minimum at risk 126000: 12281.71
minimum at risk 150000: 10972.57
(-1000, 80000)
```

Out [28]:



As intuition suggests, as we decrease the risk threshold (care about overbidding less), we increase our bid, willing to edge closer to the true price. It is interesting how far away our optimized loss is from the posterior mean, which was about 20 000.

Suffice to say, in higher dimensions being able to eyeball the minimum expected loss is impossible. Hence why we require use of Scipy's `fmin` function.

## Shortcuts

For some loss functions, the Bayes action is known in closed form. We list some of them below:

- If using the mean-squared loss, the Bayes action is the mean the posterior distribution, i.e. the value

$$E_{\theta} [\theta]$$

minimizes  $E_{\theta} [(\theta - \hat{\theta})^2]$ . Computationally this requires us to calculate the average of the posterior samples [See chapter 4 on The Law of Large Numbers]

- Whereas the *median* of the posterior distribution minimizes the expected absolute-loss. The sample median of the posterior samples is an appropriate and very accurate approximation to the true median.
- In fact, it is possible to show that the MAP estimate is the solution to using a loss function that shrinks to the zero-one loss.

Maybe it is clear now why the first-introduced loss functions are used most often in the mathematics of Bayesian inference: no complicated optimizations are necessary. Luckily, we have machines to do the complications for us.



## Machine Learning via Bayesian Methods

Whereas frequentist methods strive to achieve the best precision about all possible parameters, machine learning cares to achieve the best *prediction* among all possible parameters. Of course, one way to achieve accurate predictions is to aim for accurate predictions, but often your prediction measure and what frequentist methods are optimizing for are very different.

For example, least-squares linear regression is the most simple active machine learning algorithm. I say active as it engages in some learning, whereas predicting the sample mean is technically *simpler*, but is learning very little if anything. The loss that determines the coefficients of the regressors is a squared-error loss. On the other hand, if your prediction loss function (or score function, which is the negative loss) is not a squared-error, like AUC, ROC, precision, etc., your least-squares line will not be optimal for the prediction loss function. This can lead to prediction results that are suboptimal.

Finding Bayes actions is equivalent to finding parameters that optimize *not parameter accuracy* but an arbitrary performance measure, however we wish to define performance (loss functions, AUC, ROC, precision/recall etc.).

The next two examples demonstrate these ideas. The first example is a linear model where we can choose to predict using the least-squares loss or a novel, outcome-sensitive loss.

The second example is adapted from a Kaggle data science project. The loss function associated with our predictions is incredibly complicated.

**Example: Financial prediction** Suppose the future return of a stock price is very small, say 0.01 (or 1%). We have a model that predicts the stock's future price, and our profit and loss is directly tied to us acting on the prediction. How should we measure the loss associated with the model's predictions, and subsequent future predictions? A squared-error loss is agnostic to the signage and would penalize a prediction of -0.01 equally as bad a prediction of 0.03:

$$0.01 - (-0.01)^2 = (0.01 - 0.03)^2 = 0.004$$

If you had made a bet based on your model's prediction, you would have earned money with a prediction of 0.03, and lost money with a prediction of -0.01, yet our loss did not capture this. We need a better loss that takes into account the *sign* of the prediction and true value. We design a new loss that is better for financial applications below:

```
In [32]: figsize(12.5, 4)
def stock_loss( true_return, yhat, alpha = 100. ):
    if true_return*yhat < 0:
        #opposite signs, not good
        return alpha*yhat**2 - sign( true_return )*yhat \
            + abs( true_return )
    else:
        return abs( true_return - yhat )

true_value = .05
pred = np.linspace( -.04, .12, 75 )

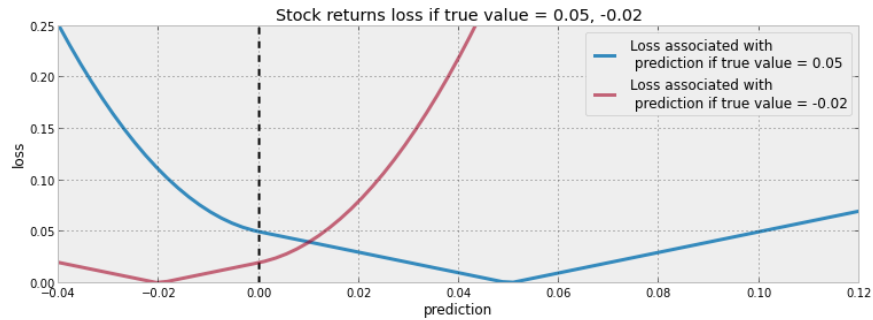
plt.plot( pred, [stock_loss( true_value, _p) for _p in pred], \
    label = "Loss associated with\n prediction if true value = 0.05", lw = 2,
    plt.vlines( 0, 0, .25, linestyles="--" )

plt.xlabel( "prediction" )
plt.ylabel( "loss" )
plt.xlim( -0.04, .12 )
plt.ylim( 0, 0.25)

true_value = -.02
plt.plot( pred, [stock_loss( true_value, _p) for _p in pred], alpha = 0.6,
    label = "Loss associated with\n prediction if true value = -0.02", lw = 2,
    plt.legend()
```

```
plt.title("Stock returns loss if true value = 0.05, -0.02" )
```

<matplotlib.text.Text at 0x176fbc18>  
Out [32]:



Note the change in the shape of the loss as the prediction crosses zero. This loss reflects that the user really does not want to guess the wrong sign, especially be wrong *and* a large magnitude.

Why would the user care about the magnitude? Why is the loss not 0 for predicting the correct sign? Surely, if the return is 0.01 and we bet millions we will still be (very) happy.

Financial institutions treat downside risk, as in predicting a lot on the wrong side, and upside risk, as in predicting a lot on the right side, similarly. Both are seen as risky behaviour and discouraged. Hence why we have an increasing loss as we move further away from the true price. (With less extreme loss in the direction of the correct sign.)

We will perform a regression on a trading signal that we believe predicts future returns well. Our dataset is artificial, as most financial data is not even close to linear. Below, we plot the data along with the least-squares line.

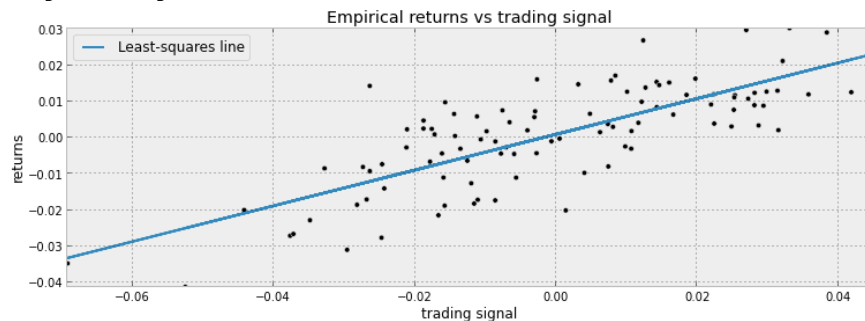
```
In [56]: ## Code to create artificial data
N = 100
X = 0.025*np.random.randn(N )
Y = 0.5*X + 0.01*np.random.randn( N)

ls_coef_ = np.cov( X, Y ) [0,1]/np.var(X)
ls_intercept = Y.mean() - ls_coef_*X.mean()

plt.scatter( X, Y, c="k")
plt.xlabel("trading signal")
plt.ylabel("returns")
plt.title( "Empirical returns vs trading signal" )
plt.plot( X, ls_coef_*X + ls_intercept, label = "Least-squares line")
plt.xlim( X.min(), X.max())
plt.ylim( Y.min(), Y.max() )
plt.legend( loc="upper left" )
```

<matplotlib.legend.Legend at 0xdbff070>

Out [56]:



We perform a simple Bayesian linear regression on this dataset. We look for a model like:

$$R = \alpha + \beta x + \epsilon$$

where  $\alpha, \beta$  are our unknown parameters and  $\epsilon \sim \text{Normal}(0, 1/\tau)$ . The most common priors on  $\beta$  and  $\alpha$  are Normal priors. We will also assign a prior on  $\tau$ , so that  $\sigma = 1/\sqrt{\tau}$  is uniform over 0 to 100 (equivalently then  $\tau = 1/\text{Uniform}(0, 100)^2$ ).

```
In [11]: import pymc as mc
from pymc.Matplot import plot as mcplot

std = mc.Uniform( "std", 0, 100, trace = False ) #this needs to be explained

@mc.deterministic
def prec( U = std ):
    return 1.0/( U )**2

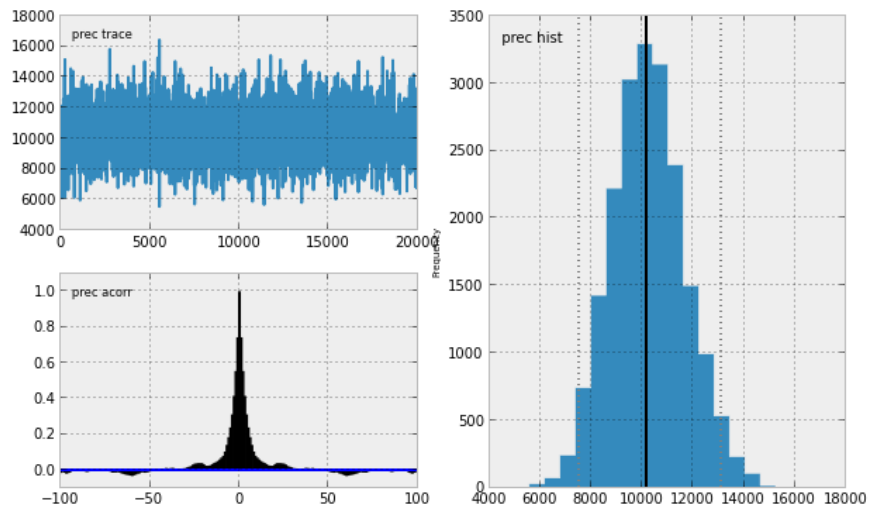
beta = mc.Normal( "beta", 0, 0.0001 )
alpha = mc.Normal( "alpha", 0, 0.0001 )

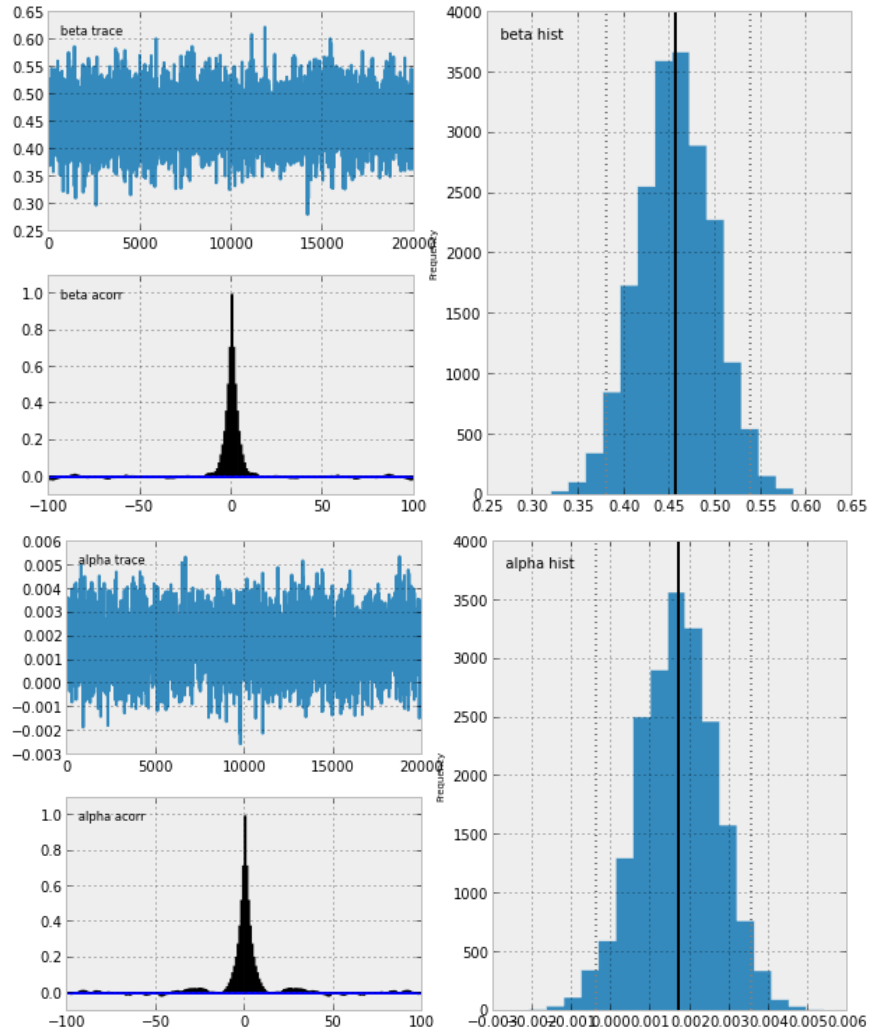
@mc.deterministic
def mean( X = X, alpha = alpha, beta = beta ):
    return alpha + beta*X

obs = mc.Normal( "obs", mean, prec, value = Y, observed = True)
model = mc.Model( {"obs":obs, "beta_0":beta, "alpha_0":alpha, "prec":prec} )
mcmc = mc.MCMC( model )

mcmc.sample( 100000, 80000 )
mcplot( mcmc )
```

[\*\*\*\*\*100%\*\*\*\*\*] 100000 of 100000 complete  
 Plotting prec  
 Plotting beta  
 Plotting alpha





It appears the MCMC has converged so we may continue.

For a specific trading signal, call it  $x$ , the distribution of possible returns has the form:

$$R_i(x) = \alpha_i + \beta_i x + \epsilon$$

where  $\epsilon \sim \text{Normal}(0, 1/\tau_i)$  and  $i$  indexes our posterior samples. We wish to find the solution to

$$\arg \min_r E_{R(x)} [ L(R(x), r) ]$$

according to the loss given above. This  $r$  is our Bayes action for trading signal  $x$ . Below we plot the Bayes action over different trading signals. What do you notice?

```
In [13]: figsize( 12.5, 6 )
from scipy.optimize import fmin

def stock_loss( price, pred, coef = 500):
    """vectorized for numpy"""
    sol = np.zeros_like(price)
    ix = price*pred < 0
```

```

sol[ix] = coef*pred**2 - sign(price[ix])*pred + abs(price[ix])
sol[~ix] = abs( price[~ix] - pred )
return sol

tau_samples = mcmc.trace( "prec" )[:]
alpha_samples = mcmc.trace( "alpha" )[:]
beta_samples = mcmc.trace( "beta" )[:]

N = tau_samples.shape[0]

noise = 1./np.sqrt(tau_samples)*np.random.randn(N)

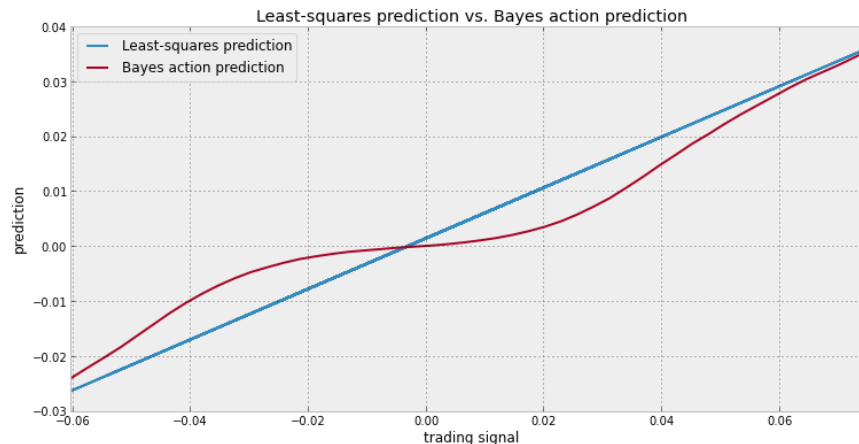
possible_outcomes = lambda signal: alpha_samples + beta_samples*signal \
    + noise

opt_predictions = np.zeros(50)
trading_signals = np.linspace( X.min(), X.max(), 50 )
for i, _signal in enumerate( trading_signals ):
    _possible_outcomes = possible_outcomes( _signal )
    tomin = lambda pred: stock_loss( _possible_outcomes, pred).mean()
    opt_predictions[i] = fmin( tomin, 0, disp = False )

plt.xlabel("trading signal")
plt.ylabel("prediction")
plt.title( "Least-squares prediction vs. Bayes action prediction" )
plt.plot( X, ls_coef*X + ls_intercept, label="Least-squares prediction")
plt.xlim( X.min(), X.max())
plot( trading_signals, opt_predictions, label="Bayes action prediction")
plt.legend( loc="upper left" )

```

<matplotlib.legend.Legend at 0x15f0a160>  
Out [13]:



What is interesting about the above graph is that when the signal is near 0, and many of the possible returns outcomes are possibly both positive and negative, our best (with respect to our loss) prediction is to predict close to 0, hence *take on no position*. Only when we are very confident do we enter into a position. I call this style of model a *sparse prediction*, where we feel uncomfortable with our uncertainty so choose not to act. (Compare with the least-squares prediction which will rarely, if ever, predict zero).

A good sanity check that our model is still reasonable: as the signal becomes more and more extreme, and we feel more and more confident about the positive/negativeness of returns, our position converges with that of the least-squares line.

The sparse-prediction model is not trying to *fit* the data the best (according to a *squared-error loss* definition of *fit*). That honor would go to the least-squares model. The sparse-prediction model is trying to find the best prediction *with respect to our stock\_loss-defined loss*. We can turn this reasoning around: the least-squares model is not try

to *predict* the best (according to a *stock-loss* definition of *predict*). That honor would go the *sparse prediction* model. The least-squares model is trying to find the best fit of the data *with respect to the squared-error loss*.

**Example: Kaggle contest on *Observing Dark World*** A personal motivation for learning Bayesian methods was trying to piece together the winning solution to Kaggle's *Observing Dark Worlds* contest. From the contest's website:

There is more to the Universe than meets the eye. Out in the cosmos exists a form of matter that outnumbers the stuff we can see by almost 7 to 1, and we don't know what it is. What we do know is that it does not emit or absorb light, so we call it Dark Matter. Such a vast amount of aggregated matter does not go unnoticed. In fact we observe that this stuff aggregates and forms massive structures called Dark Matter Halos. Although dark, it warps and bends spacetime such that any light from a background galaxy which passes close to the Dark Matter will have its path altered and changed. This bending causes the galaxy to appear as an ellipse in the sky.

The contest required predictions about where dark matter was likely to be. The winner, [Tim Salimans](#), used Bayesian inference to find the best locations for the halos (interestingly, the second-place winner also used Bayesian inference). With Tim's permission, we provided his solution [1] here:

1. Construct a prior distribution for the halo positions  $p(x)$ , i.e. formulate our expectations about the halo positions before looking at the data.
2. Construct a probabilistic model for the data (observed ellipticities of the galaxies) given the positions of the dark matter halos:  $p(e|x)$ .
3. Use Bayes' rule to get the posterior distribution of the halo positions, i.e. use to the data to guess where the dark matter halos might be.
4. Minimize the expected loss with respect to the posterior distribution over the predictions for the halo positions:  $\$ = \arg \min \{prediction\} E\{p(x|e)\}[L(prediction, x)]$ , i.e. *tune our prediction to be as good as possible for the given error metric*. The loss function in this problem is very complicated. For the very determined, the loss function is contained in the file `DarkWorldsMetric.py` in parent folder. Though I suggest not reading it all, suffice to say the loss function is about 160 lines of code — not something that can be written down in a single mathematical line. The loss function attempts to measure the accuracy of prediction, in a Euclidean distance sense, and that no shift-bias is present. More details can be found on the metric's [main page](#).

We will attempt to implement Tim's winning solution using PyMC and our knowledge of loss functions.

## The Data

The dataset is actually 300 separate files, each representing a sky. In each file, or sky, are between 300 and 720 galaxies. Each galaxy has an  $x$  and  $y$  position associated with it, ranging from 0 to 4200, and measures of ellipticity:  $e_1$  and  $e_2$ . Information about what these measures mean can be found [here](#), but for our purposes it does not matter besides for visualization purposes. Thus a typical sky might look like the following:

```
In [1]: from draw_sky2 import draw_sky

n_sky = 3 #choose a file/sky to examine.
data = np.genfromtxt( "data/Train_Skies/Train_Skies/\
Training_Sky%d.csv"%n_sky,
                    dtype = None,
                    skip_header = 1,
                    delimiter = ",",
                    usecols = [1,2,3,4])

print "Data on galaxies in sky %d."%n_sky
print "position_x, position_y, e_1, e_2 "
print data[:3]

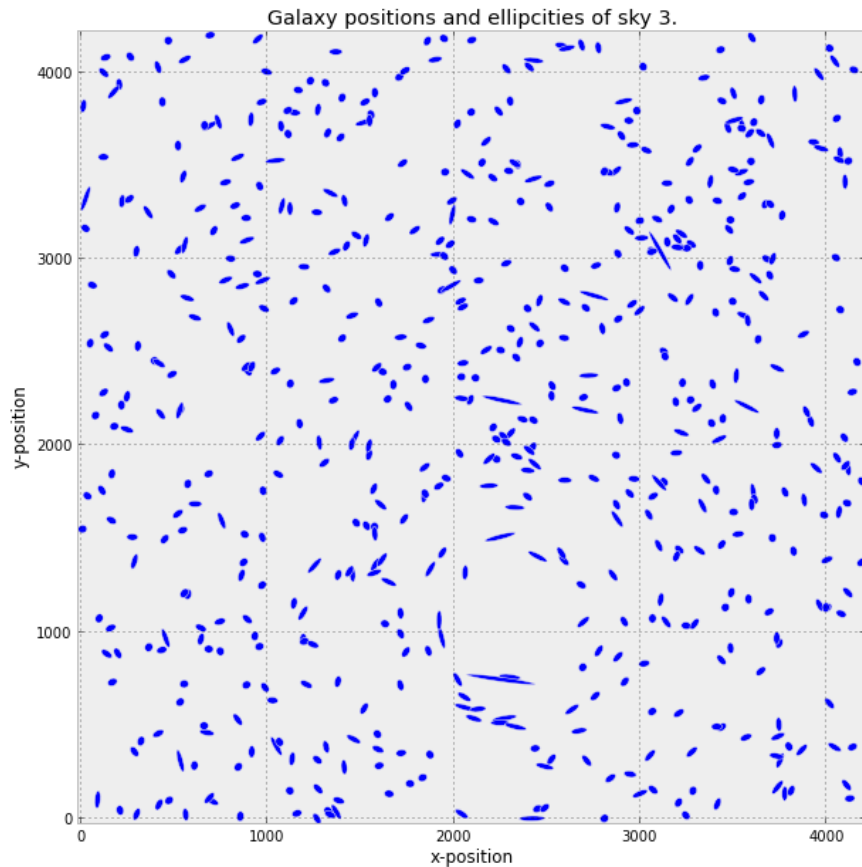
fig = draw_sky( data )
plt.title("Galaxy positions and ellipticities of sky %d."%n_sky)
```

```
plt.xlabel( "x-position")
plt.ylabel( "y-position" );
```

Data on galaxies in sky 3.

```
position_x, position_y, e_1, e_2
```

```
[[ 1.62690000e+02  1.60006000e+03  1.14664000e-01 -1.90326000e-01]
 [ 2.27228000e+03  5.40040000e+02  6.23555000e-01  2.14979000e-01]
 [ 3.55364000e+03  2.69771000e+03  2.83527000e-01 -3.01870000e-01]]
```



## Priors

Each sky has one, two or three dark matter halos in it. Tim's solution details that his prior distribution of halo positions was uniform, i.e.

$$x_i \sim \text{Uniform}(0, 4200)$$

$$y_i \sim \text{Uniform}(0, 4200), \quad i = 1, 2, 3$$

Tim and other competitors noted that most skies had one large halo and other halos, if present, were much smaller. Larger halos, having more mass, will influence the surrounding galaxies more. He decided that the large halos would have a mass distributed as a *log*-uniform random variable between 40 and 180 i.e.

$$m_{\text{large}} = \log \text{Uniform}(40, 180)$$

and in PyMC,

```
exp_mass_large = mc.Uniform( "exp_mass_large", 40, 180)
@mc.deterministic
def mass_large(u = exp_mass_large):
    return np.log( u )
```

(This is what we mean when we say *log-uniform*.) For smaller galaxies, Tim set the mass to be the logarithm of 20. Why did Tim not create a prior for the smaller mass, nor treat it as a unknown? I believe this decision was made to speed up convergence of the algorithm. This is not too restrictive, as by construction the smaller halos have less influence on the galaxies.

Tim logically assumed that the ellipticity of each galaxy is dependent on the position of the halos, the distance between the galaxy and halo, and the mass of the halos. Thus the vector of ellipticity of each galaxy,  $e_i$ , are *children* variables of the vector of halo positions ( $\mathbf{x}$ ,  $\mathbf{y}$ ), distance (which we will formalize), and halo masses.

Tim conceived a relationship to connect positions and ellipticity by reading literature and forum posts. He supposed the following was a reasonable relationship:

$$e_i | (\mathbf{x}, \mathbf{y}) \sim \text{Normal} \left( \sum_{j=\text{halo positions}} d_{i,j} m_j f(r_{i,j}), \sigma^2 \right)$$

where  $d_{i,j}$  is the *tangential direction* (the direction in which halo  $j$  bends the light of galaxy  $i$ ),  $m_j$  is the mass of halo  $j$ ,  $f(r_{i,j})$  is a *decreasing function* of the Euclidean distance between halo  $j$  and galaxy  $i$ .

Tim's function  $f$  was defined:

$$f(r_{i,j}) = \frac{1}{\min(r_{i,j}, 240)}$$

for large halos, and for small halos

$$f(r_{i,j}) = \frac{1}{\min(r_{i,j}, 70)}$$

This fully bridges our observations and unknown. This model is incredibly simple, and Tim mentions this simplicity was purposefully designed: it prevents the model from overfitting.

## Training & PyMC implementation

For each sky, we run our Bayesian model to find the posteriors for the halo positions — we ignore the (known) halo position. This is slightly different than perhaps traditional approaches to Kaggle competitions, where this model uses no data from other skies nor the known halo location. That does not mean other data are not necessary — in fact, the model was created by comparing different skies.

```
In [7]: def euclidean_distance( x, y):
        return np.sqrt( ( ( x - y )**2 ).sum(axis=1) )

def f_distance( gxy_pos, halo_pos, c):
    # foo_position should be a 2-d numpy array
    return np.maximum(euclidean_distance( gxy_pos, halo_pos), c)[:,None]

def tangential_distance( glxy_position, halo_position):
    # foo_position should be a 2-d numpy array
```



```

delta = glxy_position - halo_position
t = (2*np.arctan( delta[:,1]/delta[:,0] ))[:,None]
return np.concatenate( [-np.cos( t), -np.sin(t) ], axis=1)

import pymc as mc

#set the size of the halo's mass
mass_large = mc.Uniform( "mass_large", 40, 180, trace = False)

#set the initial prior position of the halos, it's a 2-d Uniform dist.
halo_position = mc.Uniform( "halo_position", 0, 4200, size = (1,2) )

@mc.deterministic
def mean(mass=mass_large, h_pos = halo_position, glx_pos=data[:,2]):
    return mass/f_distance( glx_pos, h_pos, 240)*\
           tangential_distance( glx_pos, h_pos )

```

```

In [18]: ellpty = mc.Normal( "ellipticity", mean, 1./0.05, observed = True,
                           value = data[:,2:] )
model = mc.Model( [ellpty, mean, halo_position, mass_large] )
mcmc = mc.MCMC( model )
map_ = mc.MAP( model )
map_.fit()
mcmc.sample(200000, 140000, 3)

```

[\*\*\*\*\*100%\*\*\*\*\*] 200000 of 200000 complete

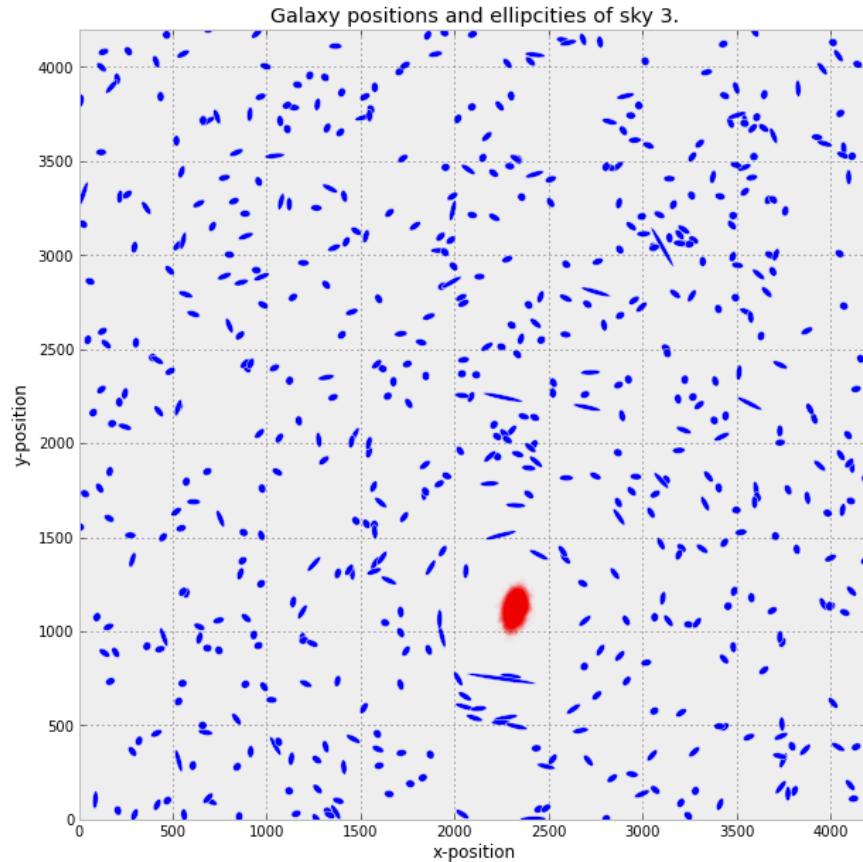
Below we plot a “heatmap” of the posterior distribution. (Which is just a scatter plot of the posterior, but we can visualize it as a heatmap.)

```

In [20]: t = mcmc.trace("halo_position")[:].reshape( 20000,2)

fig = draw_sky( data )
plt.title("Galaxy positions and ellipticities of sky %d."%n_sky)
plt.xlabel( "x-position" )
plt.ylabel( "y-position" )
scatter( t[:,0], t[:,1], alpha = 0.015, c = "r" )
plt.xlim( 0, 4200 )
plt.ylim(0, 4200 );

```



The most probable position reveals itself like a lethal wound.

Associated with each sky is another data point, located in `./data/Training_halos.csv` that holds the locations of up to three dark matter halos contained in the sky. For example, the night sky we trained on has halo locations:

```
In [13]: halo_data = np.genfromtxt( "data/Training_halos.csv",
                                   delimiter = ",",
                                   usecols = [1, 2, 3, 4, 5, 6, 7, 8, 9],
                                   skip_header = 1)
print halo_data[ n_sky]
```

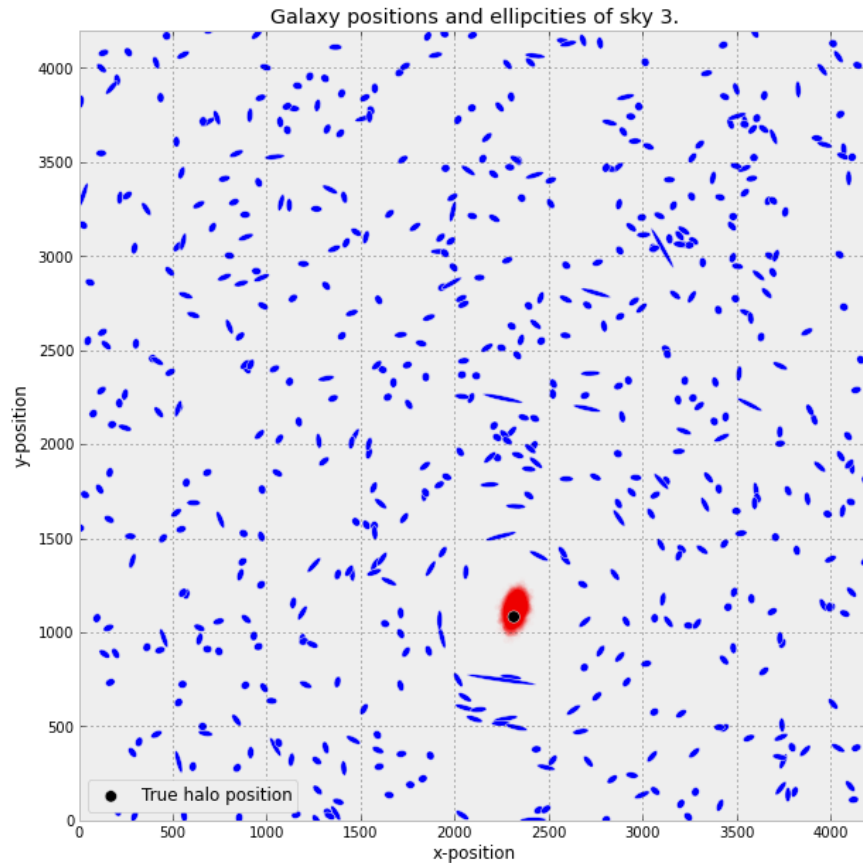
```
[ 3.00000000e+00  2.78145000e+03  1.40691000e+03  3.08163000e+03
 1.15611000e+03  2.28474000e+03  3.19597000e+03  1.80916000e+03
 8.45180000e+02]
```

The third and fourth column represent the true  $x$  and  $y$  position of the halo. It appears that the Bayesian method has located the halo within a tight vicinity.

```
In [24]: fig = draw_sky( data )
plt.title("Galaxy positions and ellipticities of sky %d."%n_sky)
plt.xlabel( "x-position" )
plt.ylabel( "y-position" )
plt.scatter( t[:,0], t[:,1], alpha = 0.015, c = "r" )
plt.scatter( halo_data[n_sky-1][3], halo_data[n_sky-1][4],
            label = "True halo position",
            c = "k", s = 70)
plt.legend(scatterpoints = 1, loc = "lower left")
plt.xlim( 0, 4200 )
plt.ylim(0, 4200 );
```

```
print "True halo location:", halo_data[n_sky][3], halo_data[n_sky][4]
```

True halo location: 1408.61 1685.86



Perfect. Our next step is to use the loss function to optimize our location. A naive strategy would be to simply choose the mean:

```
In [25]: mean_posterior = t.mean(axis=0).reshape(1,2)
print mean_posterior
```

```
[[ 2324.07677813  1122.47097816]]
```

```
In [25]: from DarkWorldsMetric import main_score

_halo_data = halo_data[n_sky-1]

nhalo_all = _halo_data[0].reshape(1,1)
x_true_all = _halo_data[3].reshape(1,1)
y_true_all = _halo_data[4].reshape(1,1)
x_ref_all = _halo_data[1].reshape(1,1)
y_ref_all = _halo_data[2].reshape(1,1)
sky_prediction = mean_posterior

print "Using the mean:"
main_score( nhalo_all, x_true_all, y_true_all, \
            x_ref_all, y_ref_all, sky_prediction)

#what's a bad score?
print
```

```

random_guess = np.random.randint( 0, 4200, size=(1,2) )
print "Using a random location:", random_guess
main_score( nhalo_all, x_true_all, y_true_all, \
            x_ref_all, y_ref_all, random_guess )
print

```

Using the mean:

Your average distance in pixels you are away from the true halo is 31.1499201664  
Your average angular vector is 1.0  
Your score for the training data is 1.03114992017

Using a random location: [[2755 53]]

Your average distance in pixels you are away from the true halo is 1773.42717812  
Your average angular vector is 1.0  
Your score for the training data is 2.77342717812

This is a good guess, it is not very far from the true location, but it ignores the loss function that was provided to us. We also need to extend our code to allow for up to two additional, *smaller* halos: Let's create a function for automatizing our PyMC.

```

In [9]: from pymc.Matplot import plot as mcplot

def halo_posteriors( n_halos_in_sky, galaxy_data,
                    samples = 5e5, burn_in = 34e4, thin = 4):

    #set the size of the halo's mass
    """
    exp_mass_large = mc.Uniform( "exp_mass_large", 40, 180)
    @mc.deterministic
    def mass_large( exp_mass_large = exp_mass_large ):
        return np.log( exp_mass_large )
    """

    mass_large = mc.Uniform( "mass_large", 40, 180 )

    mass_small_1 = 20
    mass_small_2 = 20

    masses = np.array( [ mass_large, mass_small_1, mass_small_2], dtype=ob

    #set the initial prior positions of the halos, it's a 2-d Uniform dist
    halo_positions = mc.Uniform( "halo_positions", 0, 4200,
                                size = (n_halos_in_sky,2)) #notice this size

    fdist_constants = np.array( [240, 70, 70] )

    @mc.deterministic
    def mean(mass=masses, h_pos = halo_positions, glx_pos=data[:,2],
            n_halos_in_sky = n_halos_in_sky):

        _sum = 0
        for i in range( n_halos_in_sky ):
            _sum += mass[i]/f_distance( glx_pos, h_pos[i, :], fdist_constan
            tangential_distance( glx_pos, h_pos[i, :] )

        return _sum

    ellpty = mc.Normal( "ellipcity", mean, 1./0.05, observed = True,
                        value = data[:,2:] )

    model = mc.Model( [ellpty, mean, halo_positions, mass_large] )

```

```

map_ = mc.MAP( model)
map_.fit(method="fmin_powell)

mcmc = mc.MCMC( model)
mcmc.sample(samples, burn_in, thin)
return mcmc.trace( "halo_positions" )[:]

```

```

In [10]: n_sky =215
data = np.genfromtxt( "data/Train_Skies/Train_Skies/\
Training_Sky%d.csv"%(n_sky),
                    dtype = None,
                    skip_header = 1,
                    delimiter = ",",
                    usecols = [1,2,3,4])

```

```

In [11]: #there are 3 halos in this file.
samples = 10.5e5
traces = halo_posteriors( 3, data, samples = samples,
                          burn_in = 9.5e5,
                          thin= 10 )

```

[\*\*\*\*\*100%\*\*\*\*\*] 1050000 of 1050000 complete

```

In [14]: fig = draw_sky( data )
plt.title("Galaxy positions and ellipcities of sky %d."%n_sky)
plt.xlabel( "x-position")
plt.ylabel( "y-position" )

colors = ["#467821", "#A60628", "#7A68A6"]

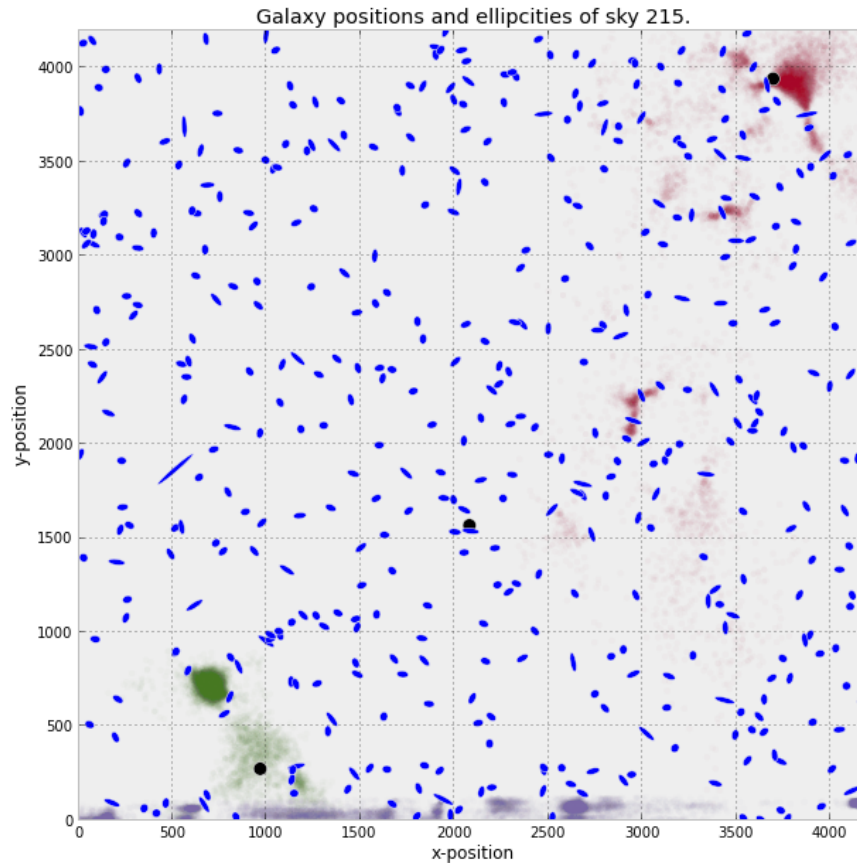
for i in range( traces.shape[1] ):
    plt.scatter( traces[:, i, 0], traces[:, i, 1], c = colors[i], alpha =

for i in range( traces.shape[1] ):
    plt.scatter( halo_data[n_sky-1][3 + 2*i], halo_data[n_sky-1][4 + 2*i],
                label = "True halo position",
                c = "k", s = 90)

plt.legend(scatterpoints = 1)
plt.xlim( 0, 4200 )
plt.ylim( 0, 4200 );

```

(0, 4200)  
Out [14]:



In []:

This looks pretty good, though it took a long time for the system to (sort of) converge. Our optimization step would look something like this:

```
In [35]: _halo_data = halo_data[n_sky-1]
print traces.shape

mean_posterior = traces.mean(axis=0).reshape( 1,4 )
print mean_posterior

nhalo_all = _halo_data[0].reshape(1,1)
x_true_all = _halo_data[3].reshape(1,1)
y_true_all = _halo_data[4].reshape(1,1)
x_ref_all = _halo_data[1].reshape(1,1)
y_ref_all = _halo_data[2].reshape(1,1)
sky_prediction = mean_posterior

print "Using the mean:"
main_score( [1], x_true_all, y_true_all, \
            x_ref_all, y_ref_all, sky_prediction)

#what's a bad score?
print
random_guess = np.random.randint( 0, 4200, size=(1,2) )
print "Using a random location:", random_guess
main_score( [1], x_true_all, y_true_all, \
```

```

print x_ref_all, y_ref_all, random_guess )
(10000L, 2L, 2L)
[[ 48.55499317 1675.79569424 1876.46951857 3265.85341193]]
Using the mean:
Your average distance in pixels you are away from the true halo is 37.3993004245
Your average angular vector is 1.0
Your score for the training data is 1.03739930042

Using a random location: [[2930 4138]]
Your average distance in pixels you are away from the true halo is 3756.54446887
Your average angular vector is 1.0
Your score for the training data is 4.75654446887

```

### 0.7.3 References

1. Antifragile: Things That Gain from Disorder. New York: Random House. 2012. ISBN 978-1-4000-6782-4.
2. [Tim Saliman's solution to the Dark World's Contest](#)
3. Silver, Nate. The Signal and the Noise: Why So Many Predictions Fail — but Some Don't. 1. Penguin Press HC, The, 2012. Print.

```

In [3]: from IPython.core.display import HTML
def css_styling():
    styles = open("../styles/custom.css", "r").read()
    return HTML(styles)
css_styling()

```

Out[3]:

## 0.8 Chapter 6

This chapter of [Bayesian Methods for Hackers](#) focuses on the most debated and discussed part of Bayesian methodologies: how to choose an appropriate prior distribution. We also present how the prior's influence changes as our dataset increases, and an interesting relationship between priors and penalties on linear regression.

### 0.8.1 Getting our priorities straight

Up until now, we have mostly ignored our choice of priors. This is unfortunate as we can be very expressive with our priors, but we also must be careful about choosing them. This is especially true if we want to be objective, that is, not express any personal beliefs in the priors.

#### Subjective vs Objective priors

Bayesian priors can be classified into two classes: *objective* priors, which aim to allow the data to influence the posterior the most, and *subjective* priors, which allow the practitioner to express his or her views into the prior.

What is an example of a objective prior? We have seen some already, including the *flat* prior (which is a uniform distribution over the entire possible range of the unknown). Using a flat prior implies we give each possible value an equal weighting. Choosing this type of prior is invoking what is called “The Principle of Indifference”, literally

we have no a priori reason to favor one value over another. Though similar, but not correct, is calling a flat prior over a restricted space an objective prior. For example, if we know  $p$  in a Binomial model is greater than 0.5, then  $\text{Uniform}(0.5, 1)$  is not an objective prior (since we have used prior knowledge) even though it is “flat” over  $[0.5, 1]$ . The flat prior must be flat along the *entire* range of possibilities.

Aside from the flat prior, other examples of objective priors are less obvious, but they contain important characteristics that reflect objectivity. For now, it should be said that *rarely* is a objective prior *truly* objective. We will see this later.

## Subjective Priors

On the other hand, if we added more probability mass to certain areas of the prior, and less elsewhere, we are biasing our inference towards the unknowns existing in the former area. This is known as a subjective, or *informative* prior. In the figure below, the subjective prior reflects a belief that the unknown likely lives around 0.5, and not around the extremes. The objective prior is insensitive to this.

```
In [2]: %pylab inline
import scipy.stats as stats
figsize(12.5,3)
colors = ["#348ABD", "#A60628", "#7A68A6", "#467821"]

x = np.linspace(0,1)
y1, y2 = stats.beta.pdf(x, 1,1), stats.beta.pdf(x, 10,10)

p = plt.plot( x, y1,
             label='An objective prior \n(uninformative, \n"Principle of Indifference" )
plt.fill_between( x, 0, y1, color = p[0].get_color(), alpha = 0.3 )

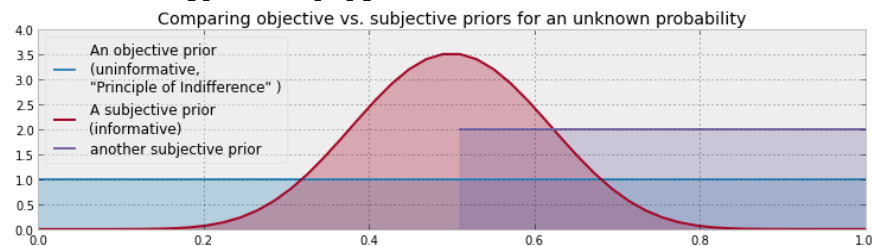
p = plt.plot( x, y2 ,
             label = "A subjective prior \n(informative)" )
plt.fill_between( x, 0, y2, color = p[0].get_color(), alpha = 0.3 )

p = plt.plot( x[25:], 2*np.ones(25), label = "another subjective prior" )
plt.fill_between( x[25:], 0, 2, color = p[0].get_color(), alpha = 0.3 )

plt.ylim(0,4)

plt.ylim(0, 4)
leg = plt.legend(loc = "upper left")
leg.get_frame().set_alpha(0.4)
plt.title("Comparing objective vs. subjective priors for an unknown probab")
```

Welcome to pylab, a matplotlib-based Python environment [backend: module://IPython.zmq.pylab]. For more information, type 'help(pylab)'.



The choice of a subjective prior does not always imply that we are using the practitioner’s subjective opinion: more often the subjective prior was once a posterior to a previous problem, and now the practitioner is updating this posterior with new data. A subjective prior can also be used to inject *domain knowledge* of the problem into the model. We will see examples of these two situations later.



## Decision, decisions...

The choice, either *objective* or *subjective* mostly depend on the problem being solved, but there are a few cases where one is preferred over the other. In instances of scientific research, the choice of an objective prior is obvious. This eliminates any biases in the results, and two researchers who might have differing prior opinions would feel an objective prior is fair. Consider a more extreme situation:

A tobacco company publishes an report with a Bayesian methodology that retreated 60 years of medical research on tobacco use. Would you believe the results? Unlikely. The researchers probably choose a subjective prior that too strongly biased results in their favor.

Unfortunately, choosing an objective prior is not as simple as selecting a flat prior, and even today the problem is still not completely solved. The problem with naively choosing the uniform prior is that pathological issues can arise. Some of these issues are pedantic, but we delay more serious issues to the Appendix of this Chapter (TODO). We must remember that choosing a prior, whether subjective or objective, is still part of the modeling process. To quote Gelman [5]:

... after the model has been fit, one should look at the posterior distribution and see if it makes sense. If the posterior distribution does not make sense, this implies that additional prior knowledge is available that has not been included in the model, and that contradicts the assumptions of the prior distribution that has been used. It is then appropriate to go back and alter the prior distribution to be more consistent with this external knowledge.

If the posterior does not make sense, then clearly one had an idea what the posterior *should* look like (not what one *hopes* it looks like), implying that the current prior does not contain all the prior information and should be updated. At this point, we can discard the current prior and choose a more reflective one.

Gelman [4] suggests that using a uniform distribution, with a large bounds, is often a good choice for objective priors. Although, one should be wary about using Uniform objective priors with large bounds, as they can assign too large of a prior probability to non-intuitive points. Ask: do you really think the unknown could be incredibly large? Often quantities are naturally biased towards 0. A Normal random variable with large variance (small precision) might be a better choice, or an Exponential with a fat tail in the strictly positive (or negative) case.

If using a particularly subjective prior, it is your responsibility to be able to explain the choice of that prior, else you are no better than the tobacco company's guilty parties.

## Empirical Bayes

While not a true Bayesian method, *empirical Bayes* is a trick that combines frequentist and Bayesian inference. As mentioned previously, for (almost) every inference problem there is a Bayesian method and a frequentist method. The significant difference between the two is that Bayesian methods have a prior distribution, with hyperparameters  $\alpha$ , while empirical methods do not have any notion of a prior. Empirical Bayes combines the two methods by using frequentist methods to select  $\alpha$ , and then proceeding with Bayesian methods on the original problem.

A very simple example follows: suppose we wish to estimate the parameter  $\mu$  of a Normal distribution, with  $\sigma = 5$ . Since  $\mu$  could range over the whole real line, we can use a Normal distribution as a prior for  $\mu$ . How to select the prior's hyperparameters, denoted  $(\mu_p, \sigma_p^2)$ ? The  $\sigma_p^2$  parameter can be chosen to reflect the uncertainty we have. For  $\mu_p$ , we have two options:

1. Empirical Bayes suggests using the empirical sample mean, which will center the prior around the observed empirical mean:

$$\mu_p = \frac{1}{N} \sum_{i=0}^N X_i$$

2. Traditional Bayesian inference suggests using prior knowledge, or a more objective prior (zero mean and fat standard deviation).

Empirical Bayes can be argued as being semi-objective, since while the choice of prior model is ours (hence subjective), the parameters are solely determined by the data.

Personally, I feel that Empirical Bayes is *double-counting* the data. That is, we are using the data twice: once in the prior, which will influence our results towards the observed data, and again in the inferential engine of MCMC. This double-counting will understate our true uncertainty. To minimize this double-counting, I would only suggest using Empirical Bayes when you have *lots* of observations, else the prior will have too strong of an influence. I would also recommend, if possible, to maintain high uncertainty (either by setting a large  $\sigma_p^2$  or equivalent.)

Empirical Bayes also violates a theoretical axiom in Bayesian inference. The textbook Bayesian algorithm of:

$$\text{prior} \Rightarrow \text{observed data} \Rightarrow \text{posterior}$$

is violated by Empirical Bayes, which instead uses

$$\text{observed data} \Rightarrow \text{prior} \Rightarrow \text{observed data} \Rightarrow \text{posterior}$$

Ideally, all prior should be specified *before* we observe the data, so that the data does not influence our prior opinions (see the volumes of research by Daniel Kahnem *et. al* about [anchoring](#) ).

## 0.8.2 Useful priors to know about

### The Gamma distribution

A Gamma random variable, denoted  $X \sim \text{Gamma}(\alpha, \beta)$ , is a random variable over the positive real numbers. It is in fact a generalization of the Exponential random variable, that is:

$$\text{Exp}(\beta) \sim \text{Gamma}(1, \beta)$$

This additional parameter allows the probability density function to have more flexibility, hence allows the practitioner to express his or her subjective priors more accurately. The density function for a  $\text{Gamma}(\alpha, \beta)$  random variable is:

$$f(x | \alpha, \beta) = \frac{\beta^\alpha x^{\alpha-1} e^{-\beta x}}{\Gamma(\alpha)}$$

where  $\Gamma(\alpha)$  is the [Gamma function](#), and for differing values of  $(\alpha, \beta)$  looks like:

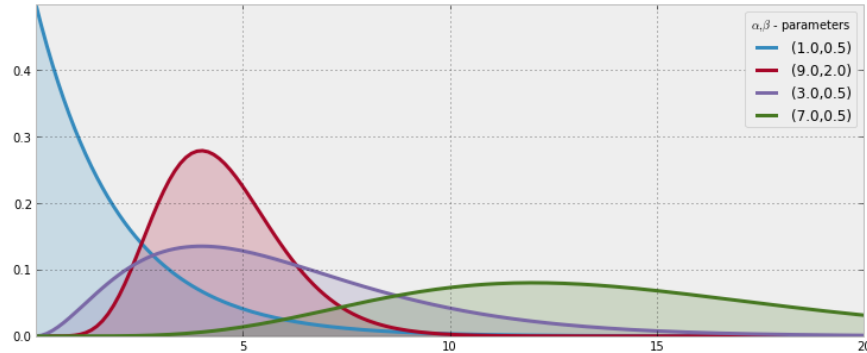
```
In [13]: figsize( 12.5, 5 )
gamma = stats.gamma

parameters = [ (1, 0.5), (9, 2), (3, 0.5), (7, 0.5) ]

x = linspace( 0.001, 20, 150 )
for alpha, beta in parameters:
    y = gamma.pdf(x, alpha, scale=1./beta)
    lines = plt.plot( x, y, label = "(%.1f, %.1f) "%(alpha, beta), lw = 3 )
    plt.fill_between( x, 0, y, alpha = 0.2, color = lines[0].get_color() )
    plt.autoscale(tight=True)

plt.legend(title=r"$\alpha, \beta$ - parameters")
```

<matplotlib.legend.Legend at 0x117f27f0>  
Out [13]:



## The Wishart distribution

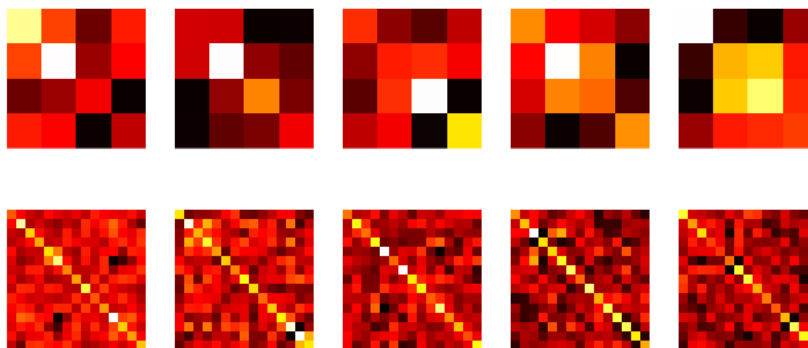
Until now, we have only seen random variables that are scalars. Of course, we can also have *random matrices*! Specifically, the Wishart distribution is a distribution over all **positive semi-definite matrices**. Why is this useful to have in our arsenal? (Proper) covariance matrices are positive-definite, hence the Wishart is an appropriate prior for covariance matrices. We can't really visualize a distribution of matrices, so I'll plot some realizations from the  $5 \times 5$  (above) and  $20 \times 20$  (below) Wishart distribution:

```
In [28]: import pymc as mc

n = 4
for i in range( 10 ):
    ax = plt.subplot( 2, 5, i+1)
    if i >= 5:
        n = 15
    plt.imshow( mc.rwishart( n+1, np.eye(n) ), interpolation="none",
                cmap = plt.cm.hot )

    ax.axis("off")
plt.suptitle("Random matrices from a Wishart Distribution" );
```

Random matrices from a Wishart Distribution



One thing to notice is that the symmetry of these matrices. The Wishart distribution can be a little troubling to deal with, but we will use it in an example later.

## The Beta distribution

You may have seen the term `beta` in previous code in this book. Often, I was implementing a Beta distribution. The Beta distribution is very useful in Bayesian statistics. A random variable  $X$  has a Beta distribution, with parameters

$(\alpha, \beta)$ , if its density function is:

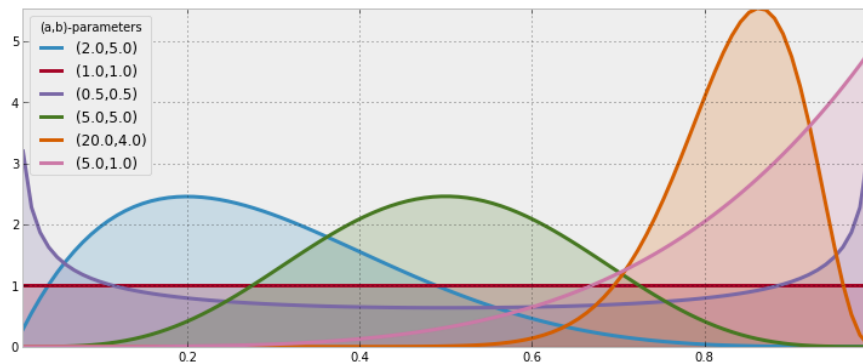
$$f_X(x | \alpha, \beta) = \frac{x^{(\alpha-1)}(1-x)^{(\beta-1)}}{B(\alpha, \beta)}$$

where  $B$  is the **Beta function** (hence the name). The random variable  $X$  is only allowed in  $[0,1]$ , making the Beta distribution a popular distribution for decimal values, probabilities and proportions. The values of  $\alpha$  and  $\beta$ , both positive values, provide great flexibility in the shape of the distribution. Below we plot some distributions:

```
In [10]: figsize( 12.5, 5 )
import scipy.stats as stats

params = [ (2,5), (1,1), (0.5, 0.5), ( 5, 5), (20, 4), (5, 1)]

x = np.linspace( 0.01, .99, 100 )
beta = stats.beta
for a, b in params:
    y = beta.pdf( x, a, b )
    lines = plt.plot( x, y, label = "(%.1f,%.1f)"%(a,b), lw = 3 )
    plt.fill_between( x, 0, y, alpha = 0.2, color = lines[0].get_color() )
    plt.autoscale(tight=True)
plt.ylim(0)
plt.legend(loc = 'upper left', title="(a,b)-parameters");
```



One thing I'd like the reader to notice is the presence of the flat distribution above, specified by parameters  $(1, 1)$ . This is the Uniform distribution. Hence the Beta distribution is a generalization of the Uniform distribution, something we will revisit many times.

There is an interesting connection between the Beta distribution and the Binomial distribution. Suppose we are interested in some unknown proportion or probability  $p$ . We assign a  $\text{Beta}(\alpha, \beta)$  prior to  $p$ . We observe some data generated by a Binomial process, say  $X \sim \text{Binomial}(N, p)$ , with  $p$  still unknown. Then our posterior is *again a Beta distribution*, i.e.  $p|X \sim \text{Beta}(\alpha + X, \beta + N - X)$ . Succinctly, one can relate the two by "a Beta prior with Binomial observations creates a Beta posterior". This is a very useful property, both computationally and heuristically.

In light of the above two paragraphs, if we start with a  $\text{Beta}(1, 1)$  prior on  $p$  (which is a Uniform), observe data  $X \sim \text{Binomial}(N, p)$ , then our posterior is  $\text{Beta}(1 + X, 1 + N - X)$ .

### Example: Bayesian Multi-Armed Bandits *Adapted from an example by Ted Dunning of MapR Technologies*

Suppose you are faced with  $N$  slot machines (colourfully called multi-armed bandits). Each bandit has an unknown probability of distributing a prize (assume for now the prizes are the same for each bandit, only the probabilities differ). Some bandits are very generous, others not so much. Of course, you don't know what these probabilities are. By only choosing one bandit per round, our task is devise a strategy to maximize our winnings.

Of course, if we knew the bandit with the largest probability, then always picking this bandit would yield the maximum winnings. So our task can be phrased as “Find the best bandit, and as quickly as possible”.

The task is complicated by the stochastic nature of the bandits. A suboptimal bandit can return many winnings, purely by chance, which would make us believe that it is a very profitable bandit. Similarly, the best bandit can return many duds. Should we keep trying losers then, or give up?

A more troublesome problem is, if we have found a bandit that returns *pretty good* results, do we keep drawing from it to maintain our *pretty good score*, or do we try other bandits in hopes of finding an *even-better* bandit? This is the exploration vs. exploitation dilemma.

## Applications

The Multi-Armed Bandit problem at first seems very artificial, something only a mathematician would love, but that is only before we address some applications:

- Internet display advertising: companies have a suite of potential ads they can display to visitors, but the company is not sure which ad strategy to follow to maximize sales. This is similar to A/B testing, but has the added advantage of naturally minimizing strategies that do not work (and generalizes to A/B/C/D... strategies)
- Ecology: animals have a finite amount of energy to expend, and following certain behaviours has uncertain rewards. How does the animal maximize its fitness?
- Finance: which stock option gives the highest return, under time-varying return profiles.
- Clinical trials: a researcher would like to find the best treatment, out of many possible treatments, while minimizing losses.
- Psychology: how does punishment and reward effect our behaviour? How do humans' learn?

Many of these questions above are fundamental to the application's field.

It turns out the *optimal solution* is incredibly difficult, and it took decades for an overall solution to develop. There are also many approximately-optimal solutions which are quite good. The one I wish to discuss is one of the few solutions that can scale incredibly well. The solution is known as *Bayesian Bandits*.

## A Proposed Solution

Any proposed strategy is called an *online algorithm* (not in the internet sense, but in the continuously-being-updated sense), and more specifically a reinforcement learning algorithm. The algorithm starts in an ignorant state, where it knows nothing, and begins to acquire data by testing the system. As it acquires data and results, it learns what the best and worst behaviours are (in this case, it learns which bandit is the best). With this in mind, perhaps we can add an additional application of the Multi-Armed Bandit problem:

- Psychology: how does punishment and reward effect our behaviour? How do humans' learn?

The Bayesian solution begins by assuming priors on the probability of winning for each bandit. In our vignette we assumed complete ignorance of these probabilities. So a very natural prior is the flat prior over 0 to 1. The algorithm proceeds as follows:

For each round:

1. Sample a random variable  $X_b$  from the prior of bandit  $b$ , for all  $b$ .
2. Select the bandit with largest sample, i.e. select  $B = \operatorname{argmax}_b X_b$ .
3. Observe the result of pulling bandit  $B$ , and update your prior on bandit  $B$ .
4. Return to 1.

That's it. Computationally, the algorithm involves sampling from  $N$  distributions. Since the initial priors are  $\text{Beta}(\alpha = 1, \beta = 1)$  (a uniform distribution), and the observed result  $X$  (a win or loss, encoded 1 and 0 respectively) is Binomial, the posterior is a  $\text{Beta}(\alpha = 1 + X, \beta = 1 + 1X)$ .

To answer our question from before, this algorithm suggests that we should not discard losers, but we should pick them at a decreasing rate as we gather confidence that there exist *better* bandits. This follows because there is always a non-zero chance that a loser will achieve the status of  $B$ , but the probability of this event decreases as we play more rounds (see figure below).

Below we implement Bayesian Bandits using two classes, `Bandits` that defines the slot machines, and `BayesianStrategy` which implements the above learning strategy.

```
In [1]: from pymc import rbeta

rand = np.random.rand

class Bandits(object):
    """
    This class represents N bandits machines.

    parameters:
        p_array: a (n,) Numpy array of probabilities >0, <1.

    methods:
        pull( i ): return the results, 0 or 1, of pulling
                   the ith bandit.
    """
    def __init__(self, p_array):
        self.p = p_array
        self.optimal = np.argmax(p_array)

    def pull( self, i ):
        #i is which arm to pull
        return rand() < self.p[i]

    def __len__(self):
        return len(self.p)

class BayesianStrategy( object ):
    """
    Implements a online, learning strategy to solve
    the Multi-Armed Bandit problem.

    parameters:
        bandits: a Bandit class with .pull method

    methods:
        sample_bandits(n): sample and train on n pulls.

    attributes:
        N: the cumulative number of samples
        choices: the historical choices as a (N,) array
        bb_score: the historical score as a (N,) array
    """
    def __init__(self, bandits):

        self.bandits = bandits
        n_bandits = len( self.bandits )
        self.wins = np.zeros( n_bandits )
        self.trials = np.zeros(n_bandits )
        self.N = 0
        self.choices = []
        self.bb_score = []
```

```

def sample_bandits( self, n=1 ):

    bb_score = np.zeros( n )
    choices = np.zeros( n )

    for k in range(n):
        #sample from the bandits's priors, and select the largest sample
        choice = np.argmax( rbeta( 1 + self.wins, 1 + self.trials - self.wins ) )

        #sample the chosen bandit
        result = self.bandits.pull( choice )

        #update priors and score
        self.wins[ choice ] += result
        self.trials[ choice ] += 1
        bb_score[ k ] = result
        self.N += 1
        choices[ k ] = choice

    self.bb_score = np.r_[ self.bb_score, bb_score ]
    self.choices = np.r_[ self.choices, choices ]
    return

```

Below we visualize the learning of the Bayesian Bandit solution.

```

In [20]: import scipy.stats as stats
figsize( 11.0, 10)

beta = stats.beta
x = np.linspace(0.001, .999, 200)

def plot_priors(bayesian_strategy, prob, lw = 3, alpha = 0.2, plt_vlines = True):
    ## plotting function
    wins = bayesian_strategy.wins
    trials = bayesian_strategy.trials
    for i in range( prob.shape[0] ):
        y = beta( 1+wins[i], 1 + trials[i] - wins[i] )
        p = plt.plot( x, y.pdf(x), lw = lw )
        c = p[0].get_markeredgcolor()
        plt.fill_between(x,y.pdf(x),0, color = c, alpha = alpha,
                        label="underlying probability: %.2f"%prob[i])
    if plt_vlines:
        plt.vlines( prob[i], 0, y.pdf(prob[i]) ,
                   colors = c, linestyles = "--", lw = 2 )
    plt.autoscale(tight = "True")
    plt.title("Posteriors After %d pull"%bayesian_strategy.N + \
              "s"*(bayesian_strategy.N>1) )
    plt.autoscale(tight=True)
    return

```

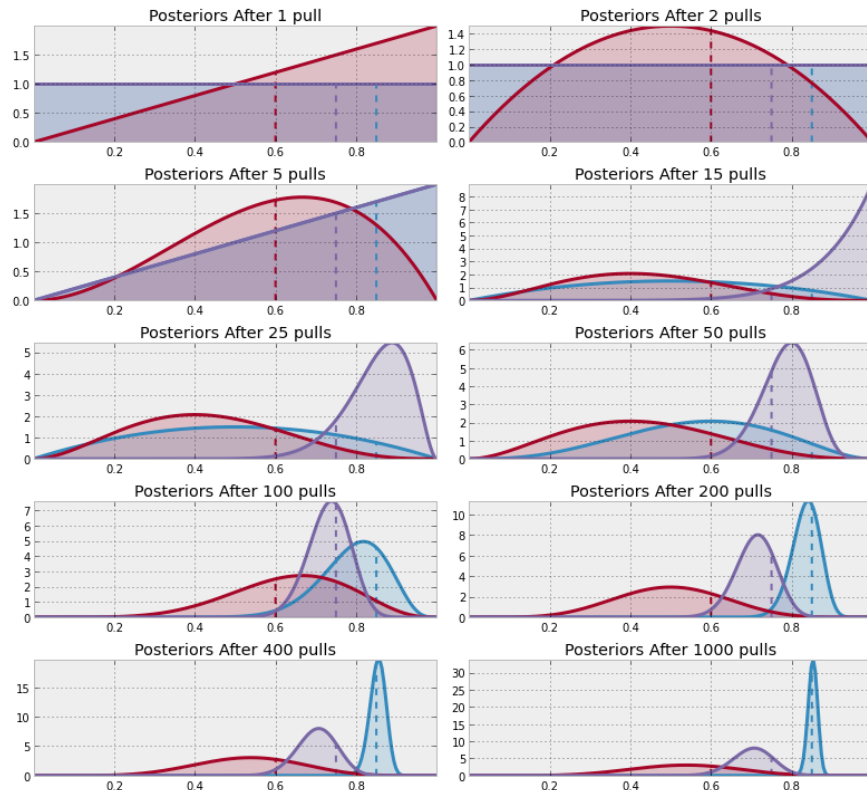
```

In [5]: hidden_prob = np.array([0.85, 0.60, 0.75] )
bandits = Bandits( hidden_prob )
bayesian_strat = BayesianStrategy( bandits )

draw_samples = [1, 1, 3, 10, 10, 25, 50, 100, 200, 600 ]

for j,i in enumerate(draw_samples):
    subplot( 5, 2, j+1)
    bayesian_strat.sample_bandits(i)
    plot_priors( bayesian_strat, hidden_prob )
    #plt.legend()
    plt.autoscale( tight = True )
plt.tight_layout()

```



Note that we don't really care how accurate we become about inference of the hidden probabilities — for this problem we are more interested in choosing the best bandit (or more accurately, becoming *more confident* in choosing the best bandit). For this reason, the distribution of the red bandit is very wide (representing ignorance about what that hidden probability might be) but we are reasonably confident that it is not the best, so the algorithm chooses to ignore it.

From the above, we can see that after 1000 pulls, the majority of the “blue” function leads the pack, hence we will almost always choose this arm. This is good, as this arm is indeed the best.

Below is a D3 app that demonstrates our algorithm updating/learning three bandits. The first figure are the raw counts of pulls and wins, and the second figure is a dynamically updating plot. I encourage you to try to guess which bandit is optimal, prior to revealing the true probabilities, by selecting the arm buttons.

```
In [4]: from IPython.core.display import HTML

        #try executing the below command twice if the first time doesn't work
        HTML(filename = "BanditsD3.html" )
```

Deviations of the observed ratio from the highest probability is a measure of performance. For example, in the long run, optimally we can attain the reward/pull ratio of the maximum bandit probability. Long-term realized ratios less than the maximum represent inefficiencies. (Realized ratios larger than the maximum probability is due to randomness, and will eventually fall below).

## A Measure of Good

We need a metric to calculate how well we are doing. Recall the absolute *best* we can do is to always pick the bandit with the largest probability of winning. Denote this best bandit's probability of  $w_{opt}$ . Our score should be relative to how well we would have done had we chosen the best bandit from the beginning. This motivates the *total regret* of a strategy, defined:



$$R_T = \sum_{i=1}^T (w_{opt} - w_{B(i)}) \quad (46)$$

(47)

$$= Tw^* - \sum_{i=1}^T w_{B(i)} \quad (48)$$

where  $w_{B(i)}$  is the probability of a prize of the chosen bandit in the  $i$  round. A total regret of 0 means the strategy is matching the best possible score. This is likely not possible, as initially our algorithm will often make the wrong choice. Ideally, a strategy's total regret should flatten as it learns the best bandit. (Mathematically we achieve  $w_{B(i)} = w_{opt}$  often)

Below we plot the total regret of this simulation, including the scores of some other strategies:

1. Random: randomly choose a bandit to pull. If you can't beat this, just stop.
2. largest Bayesian credible bound: pick the bandit with the largest upper bound in its 95% credible region of the underlying probability.
3. Bayes-UCB algorithm: pick the bandit with the largest *score*, where score is a dynamic quantile of the posterior (see [4])
4. Mean of posterior: choose the bandit with the largest posterior mean. This is what a human player (sans computer) would likely do.
5. Largest proportion: pick the bandit with the current largest observed proportion of winning.

The code for these are in the `other_strats.py`, where you can implement your own very easily.

```
In [2]: figsize( 12.5, 5 )
        from other_strats import *

        #define a harder problem
        hidden_prob = np.array([0.15, 0.2, 0.1, 0.05] )
        bandits = Bandits( hidden_prob )

        #define regret
        def regret( probabilities, choices ):
            w_opt = probabilities.max()
            return ( w_opt - probabilities[choices.astype(int)] ).cumsum()

        #create new strategies

        strategies= [upper_credible_choice,
                    bayesian_bandit_choice,
                    ucb_bayes ,
                    max_mean,
                    random_choice ]

        algos = []
        for strat in strategies:
            algos.append( GeneralBanditStrat( bandits, strat ))
```

```
In [3]: #train 10000 times
        for strat in algos:
            strat.sample_bandits( 10000)

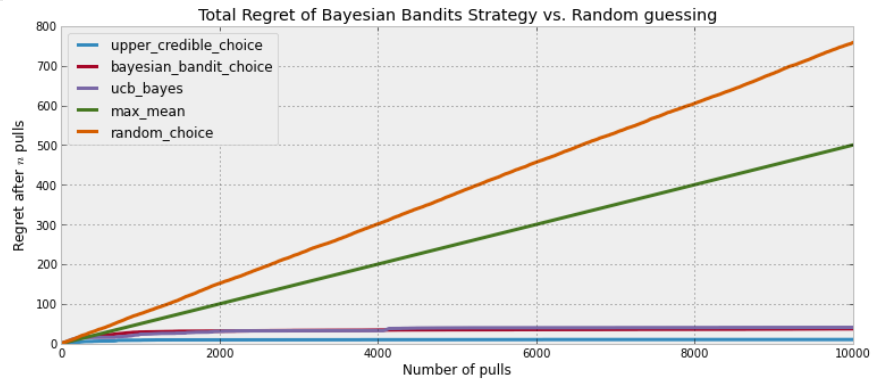
        #test and plot
```

```

for i, strat in enumerate(algos):
    _regret = regret( hidden_prob, strat.choices )
    plt.plot( _regret, label = strategies[i].__name__, lw = 3)

plt.title("Total Regret of Bayesian Bandits Strategy vs. Random guessing")
plt.xlabel("Number of pulls")
plt.ylabel("Regret after $n$ pulls");
plt.legend(loc = "upper left");

```



Like we wanted, Bayesian bandits and other strategies have decreasing rates of regret, representing we are achieving optimal choices. To be more scientific so as to remove any possible luck in the above simulation, we should instead look at the *expected total regret*:

$$\bar{R}_T = E[R_T]$$

It can be shown that any *sub-optimal* strategy's expected total regret is bounded below logarithmically. Formally,

$$E[R_T] = \Omega( \log(T) )$$

Thus, any strategy that matches logarithmic-growing regret is said to “solve” the Multi-Armed Bandit problem [3].

Using the Law of Large Numbers, we can approximate Bayesian Bandit's expected total regret by performing the same experiment many times (500 times, to be fair):

```

In [24]: #this can be slow, so I recommend NOT running it.

trials = 500
expected_total_regret = np.zeros( ( 10000, 3 ) )

for i_strat, strat in enumerate( strategies[:-2] ):
    for i in range(trials):
        general_strat = GeneralBanditStrat( bandits, strat )
        general_strat.sample_bandits(10000)
        _regret = regret( hidden_prob, general_strat.choices )
        expected_total_regret[:,i_strat] += _regret

    plot(expected_total_regret[:,i_strat]/trials, lw =3, label = strat.__name__)

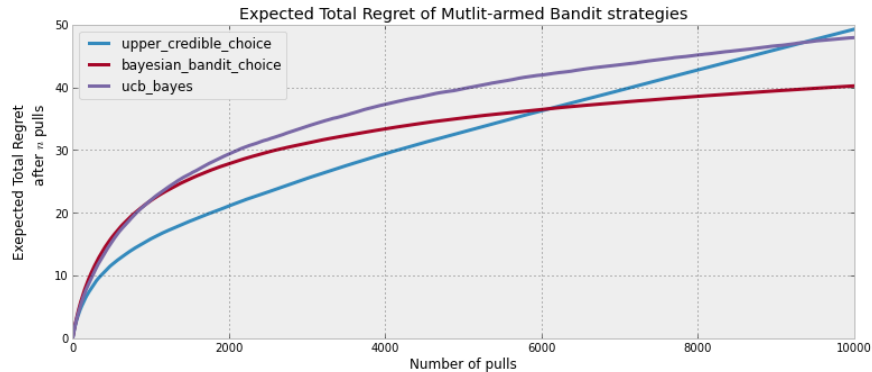
plt.title("Expected Total Regret of Multi-armed Bandit strategies" )
plt.xlabel("Number of pulls")
plt.ylabel("Expected Total Regret \n after $n$ pulls");
plt.legend(loc = "upper left");

```

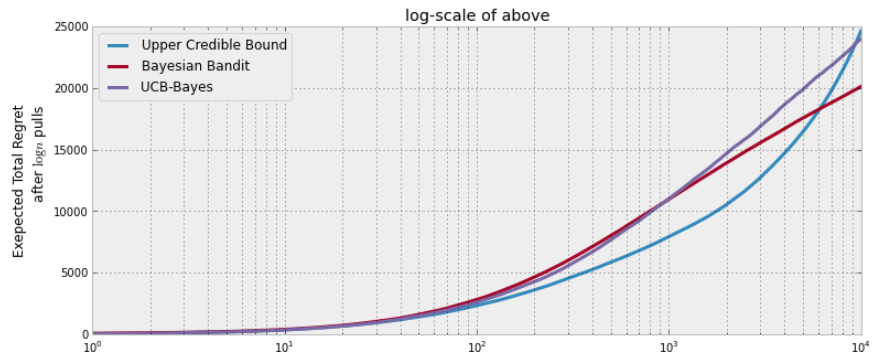
0  
1  
2

<matplotlib.legend.Legend at 0x137a3860>

Out [24]:



```
In [25]: plt.figure()  
[p11, p12, p13] = plt.plot( expected_total_regret[:, [0,1,2]], lw = 3 )  
plt.xscale("log")  
plt.legend([ p11, p12, p13],  
           ["Upper Credible Bound", "Bayesian Bandit", "UCB-Bayes"],  
           loc="upper left")  
plt.ylabel("Exepected Total Regret \n after $\log\{n\}$ pulls");  
plt.title( "log-scale of above" );  
plt.ylabel("Exepected Total Regret \n after $\log\{n\}$ pulls");
```



## Extending the algorithm

Because of the Bayesian Bandits algorithm's simplicity, it is easy to extend. Some possibilities:

- If interested in the *minimum* probability (eg: where prizes are are bad thing), simply choose  $B = \operatorname{argmin} X_b$  and proceed.
- Adding learning rates: Suppose the underlying environment may change over time. Technically the standard Bayesian Bandit algorithm would self-update itself (awesome) by noting that what it thought was the best is starting to fail more often, we can motivate the algorithm to learn changing environments quicker. We simply need to add a *rate* term upon updating:

```
self.wins[ choice ] = rate*self.wins[ choice ] + result  
self.trials[ choice ] = rate*self.trials[ choice ] + 1
```

If  $\text{rate} < 1$ , the algorithm will *forget* its previous wins quicker and there will be a downward pressure towards ignorance. Conversely, setting  $\text{rate} > 1$  implies your algorithm will act more risky, and bet on earlier winners more often and be more resistant to changing environments.

- Hierarchical algorithms: We can setup a Bayesian Bandit algorithm on top of smaller bandit algorithms. Suppose we have  $N$  Bayesian Bandit models, each varying in some behavior (for example different `rate` parameters, representing varying sensitivity to changing environments). On top of these  $N$  models is another Bayesian Bandit learner that will select a sub-Bayesian Bandit. This chosen Bayesian Bandit will then make an internal choice as to which machine to pull. The super-Bayesian Bandit updates itself depending on whether the sub-Bayesian Bandit was correct or not.
- Extending the rewards, denoted  $y_a$  for bandit  $a$ , to random variables from a distribution  $f_{y_a}(y)$  is straightforward. More generally, this problem can be rephrased as “Find the bandit with the largest expected value”, as playing the bandit with the largest expected value is optimal. In the case above,  $f_{y_a}$  was Bernoulli with probability  $p_a$ , hence the expected value for an bandit is equal to  $p_a$ , which is why it looks like we are aiming to maximize the probability of winning. If  $f$  is not Bernoulli, and it is non-negative, which can be accomplished a priori by shifting the distribution (we assume we know  $f$ ), then the algorithm behaves as before:

For each round,

1. Sample a random variable  $X_b$  from the prior of bandit  $b$ , for all  $b$ .
2. Select the bandit with largest sample, i.e. select bandit  $B = \operatorname{argmax}_b X_b$ .
3. Observe the result,  $R \sim f_{y_a}$ , of pulling bandit  $B$ , and update your prior on bandit  $B$ .
4. Return to A

The issue is in the sampling of  $X_b$  drawing phase. With Beta priors and Bernoulli observations, we have a Beta posterior — this is easy to sample from. But now, with arbitrary distributions  $f$ , we have a non-trivial posterior. Sampling from these can be difficult.

- There has been some interest in extending the Bayesian Bandit algorithm to commenting systems. Recall in Chapter 4, we developed a ranking algorithm based on the Bayesian lower-bound of the proportion of upvotes to total total votes. One problem with this approach is that it will bias the top rankings towards older comments, since older comments naturally have more votes (and hence the lower-bound is tighter to the true proportion). This creates a positive feedback cycle where older comments gain more votes, hence are displayed more often, hence gain more votes, etc. This pushes any new, potentially better comments, towards the bottom. J. Neufeld proposes a system to remedy this that uses a Bayesian Bandit solution.

His proposal is to consider each comment as a Bandit, with a the number of pulls equal to the number of votes cast, and number of rewards as the number of upvotes, hence creating a  $\text{Beta}(1 + U, 1 + D)$  posterior. As visitors visit the page, samples are drawn from each bandit/comment, but instead of displaying the comment with the max sample, the comments are ranked according the the ranking of their respective samples. From J. Neufeld’s blog [7]:

[The] resulting ranking algorithm is quite straightforward, each new time the comments page is loaded, the score for each comment is sampled from a  $\text{Beta}(1 + U, 1 + D)$ , comments are then ranked by this score in descending order. . . This randomization has a unique benefit in that even untouched comments ( $U = 1, D = 0$ ) have some chance of being seen even in threads with 5000+ comments (something that is not happening now), but, at the same time, the user will is not likely to be inundated with rating these new comments.

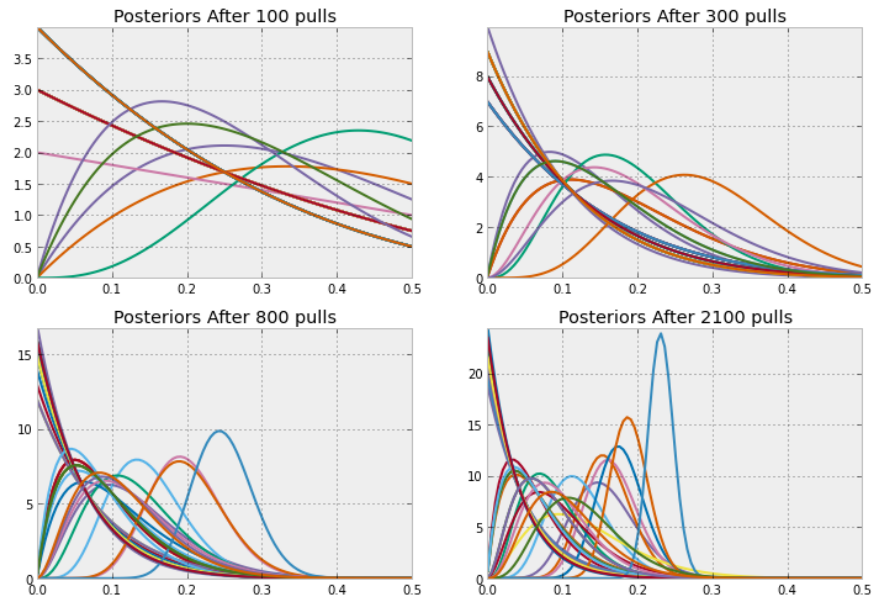
Just for fun, though the colors explode, we watch the Bayesian Bandit algorithm learn 15 different options.

```
In [31]: figsize( 12.0, 8)
         beta = stats.beta
         hidden_prob = beta.rvs(1,13, size = 35 )
         print hidden_prob
         bandits = Bandits( hidden_prob )
         bayesian_strat = BayesianStrategy( bandits )

         for j,i in enumerate([100, 200, 500, 1300 ]):
             subplot( 2, 2, j+1)
             bayesian_strat.sample_bandits(i)
             plot_priors( bayesian_strat, hidden_prob, lw = 2, alpha = 0.0, plt_vl:
```

```
#plt.legend()
plt.xlim(0, 0.5 )
```

```
[ 1.05872577e-02  2.44329406e-02  2.17921542e-01  2.84275845e-02
 1.79636306e-01  1.76434301e-01  1.03803619e-01  3.41231782e-02
 1.51276511e-01  2.02988173e-05  8.25837083e-02  3.00819469e-02
 3.87676920e-02  9.95699292e-02  5.74074773e-02  8.88739680e-02
 3.04542676e-02  7.08154796e-02  8.18108782e-03  1.45369509e-01
 6.87009381e-03  8.52942080e-02  1.28626668e-02  7.84557611e-02
 1.84449164e-01  6.52984488e-02  4.34863358e-02  4.39868639e-02
 4.82796284e-02  7.85578667e-02  2.20346874e-01  4.95440403e-02
 4.70066750e-02  8.37623221e-02  5.95142751e-02]
```



### 0.8.3 Eliciting expert prior

Specifying a subjective prior is how practitioners incorporate domain knowledge about the problem into our mathematical framework. Allowing domain knowledge is useful for many reasons:

- Aids speeds of MCMC convergence. For example, if we know the unknown parameter is strictly positive, then we can restrict our attention there, hence saving time that would otherwise be spent exploring negative values.
- More accurate inference. By weighing prior values near the true unknown value higher, we are narrowing our eventual inference (by making the posterior tighter around the unknown)
- Express our uncertainty better. See the *Price is Right* problem in Chapter 3.

plus many other reasons. Of course, practitioners of Bayesian methods are not experts in every field, so we must turn to domain experts to craft our priors. We must be careful with how we elicit these priors though. Some things to consider:

1. From experience, I would avoid introducing Betas, Gammas, etc. to non-Bayesian practitioners. Furthermore, non-statisticians can get tripped up by how a continuous probability function can have a value exceeding one.
2. Individuals often neglect the rare *tail-events* and put too much weight around the mean of distribution.
3. Related to above is that almost always individuals will under-emphasize the uncertainty in their guesses.

Eliciting priors from non-technical experts is especially difficult. Rather than introduce the notion of probability distributions, priors, etc. that may scare an expert, there is a much simpler solution.

## Trial roulette method

The *trial roulette method* [8] focuses on building a prior distribution by placing counters (think casino chips) on what the expert thinks are possible outcomes. The expert is given  $N$  counters (say  $N = 20$ ) and is asked to place them on a pre-printed grid, with bins representing intervals. Each column would represent their belief of the probability of getting the corresponding bin result. Each chip would represent an  $\frac{1}{N} = 0.05$  increase in the probability of the outcome being in that interval. For example [9]:

A student is asked to predict the mark in a future exam. The figure below shows a completed grid for the elicitation of a subjective probability distribution. The horizontal axis of the grid shows the possible bins (or mark intervals) that the student was asked to consider. The numbers in top row record the number of chips per bin. The completed grid (using a total of 20 chips) shows that the student believes there is a 30% chance that the mark will be between 60 and 64.9.

From this, we can fit a distribution that captures the expert's choice. Some reasons in favor of using this technique are:

1. Many questions about the shape of the expert's subjective probability distribution can be answered without the need to pose a long series of questions to the expert - the statistician can simply read off density above or below any given point, or that between any two points.
2. During the elicitation process, the experts can move around the chips if unsatisfied with the way they placed them initially - thus they can be sure of the final result to be submitted.
3. It forces the expert to be coherent in the set of probabilities that are provided. If all the chips are used, the probabilities must sum to one.
4. Graphical methods seem to provide more accurate results, especially for participants with modest levels of statistical sophistication.

**Example: Stock Returns** Take note stock brokers: you're doing it wrong. When choosing which stocks to pick, an analyst will often look at the *daily return* of the stock. Suppose  $S_t$  is the price of the stock on day  $t$ , then the daily return on day  $t$  is :

$$r_t = \frac{S_t - S_{t-1}}{S_{t-1}}$$

The *expected daily return* of a stock is denoted  $\mu = E[r_t]$ . Obviously, stocks with high expected returns are desirable. Unfortunately, stock returns are so filled with noise that it is very hard to estimate this parameter. Furthermore, the parameter might change over time (consider the rises and falls of AAPL stock), hence it is unwise to use a large historical dataset.

Historically, the expected return has been estimated by using the sample mean. This is a bad idea. As mentioned, the sample mean of a small dataset size has enormous potential to be very wrong (again, see Chapter 4 for full details). Thus Bayesian inference is the correct procedure here, since we are able to see our uncertainty along with probable values.

For this exercise, we will be examining the daily returns of the AAPL, GOOG, MSFT and AMZN. Before we pull in the data, suppose we ask our a stock fund manager (an expert in finance, but see [10] ),

What do you think the return profile looks like for each of these companies?

Our stock broker, without needing to know the language of Normal distributions, or priors, or variances, etc. creates four distributions using the trial roulette method above. Suppose they look enough like Normals, so we fit Normals to them. They may look like:

```
In [2]: import scipy.stats as stats
        figsize(11., 5)
        colors = ["#348ABD", "#A60628", "#7A68A6", "#467821"]
```

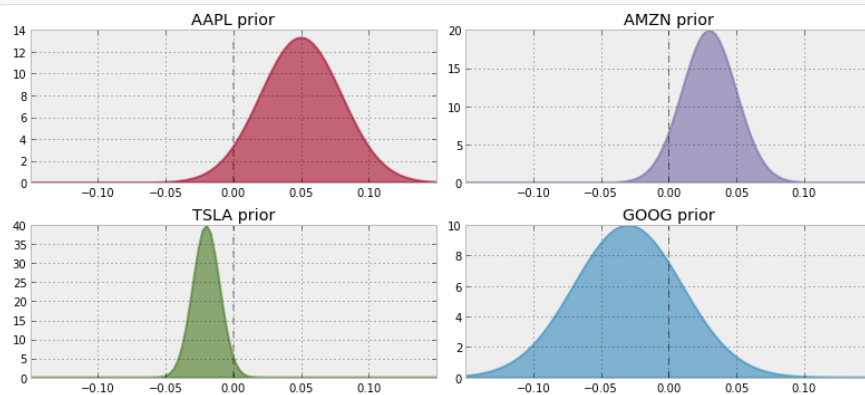
```

normal = stats.norm
x = np.linspace( -0.15, 0.15, 100 )

expert_prior_params = { "AAPL":(0.05, 0.03),
                        "GOOG":(-0.03, 0.04 ),
                        "TSLA": (-0.02, 0.01),
                        "AMZN": (0.03, 0.02 ),
                        }

for i, (name, params) in enumerate(expert_prior_params.iteritems() ):
    subplot(2,2,i)
    y = normal.pdf( x, params[0], scale = params[1] )
    #plt.plot( x, y, c = colors[i] )
    plt.fill_between(x, 0, y, color = colors[i], linewidth=2,
                    edgecolor = colors[i], alpha = 0.6)
    plt.title(name + " prior" )
    plt.vlines(0, 0, y.max(), "k", "--", linewidth = 0.5 )
    plt.xlim(-0.15, 0.15 )
plt.tight_layout()

```



Note that these are subjective priors: the expert has a personal opinion on the stock returns of each of these companies, and is expressing them in a distribution. He's not wishful thinking – he's introducing domain knowledge.

In order to better model these returns, we should investigate the *covariance matrix* of the returns. For example, it would be unwise to invest in two stocks that are highly correlated, since they are likely to tank together (hence why fund managers suggest a diversification strategy). We will use the *Wishart distribution* for this, introduced earlier.

```

In [3]: import pymc as mc

n_observations = 100 #we will truncate the the most recent 100 days.

prior_mu = np.array( [ x[0] for x in expert_prior_params.values() ] )
prior_std = np.array( [ x[1] for x in expert_prior_params.values() ] )

inv_cov_matrix = mc.Wishart( "inv_cov_matrix", n_observations, diag(prior_mu
mu = mc.Normal( "returns", prior_mu, 1, size = 4 )

```

Next we pull historical data for these stocks:

```

In [69]: # I wish I could have used Pandas as a prereq for this book, but oh well.

import datetime
import ystockquote as ysq

stocks = [ "AAPL", "GOOG", "TSLA", "AMZN" ]

enddate = datetime.datetime.now().strftime("%Y-%m-%d") #today's date.
startdate = "2012-09-01"

```

```

stock_closes = {}
stock_returns = {}
CLOSE = 6

for stock in stocks:
    x = np.array( ysq.get_historical_prices( stock, startdate, enddate ) )
    stock_closes[stock] = x[1:,:CLOSE].astype(float)

#create returns:

for stock in stocks:
    _previous_day = np.roll(stock_closes[stock],-1)
    stock_returns[stock] = ((stock_closes[stock] - _previous_day)/_previous

dates = map( lambda x: datetime.datetime.strptime(x, "%Y-%m-%d" ), x[1:n_c

```

In [100]: figsize(12.5, 4)

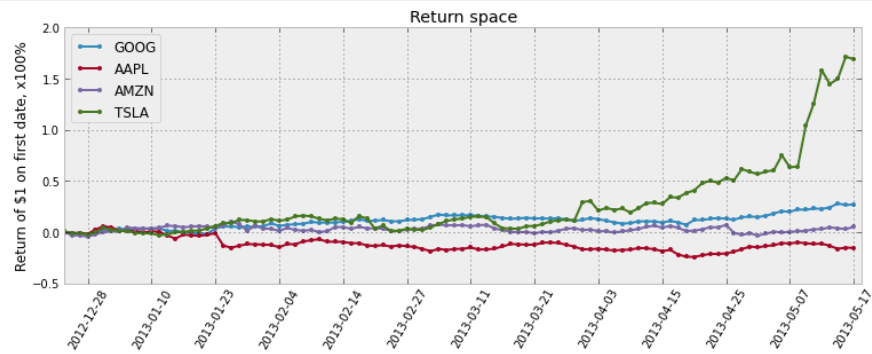
```

for _stock, _returns in stock_returns.iteritems():
    p = plot( (1+_returns)[::-1].cumprod()-1, '-o', label = "%s"%_stock,
              markersize=4, markeredgecolor="none" )

plt.xticks( arange(100)[::-8],
            map(lambda x: datetime.datetime.strptime(x, "%Y-%m-%d"), dates),
            rotation=60);

plt.legend(loc = "upper left")
plt.title("Return space")
plt.ylabel("Return of $1 on first date, x100%");

```



In [121]: figsize(11., 5)

```

returns = np.zeros( (n_observations,4 ) )

for i, (_stock,_returns) in enumerate(stock_returns.iteritems() ):

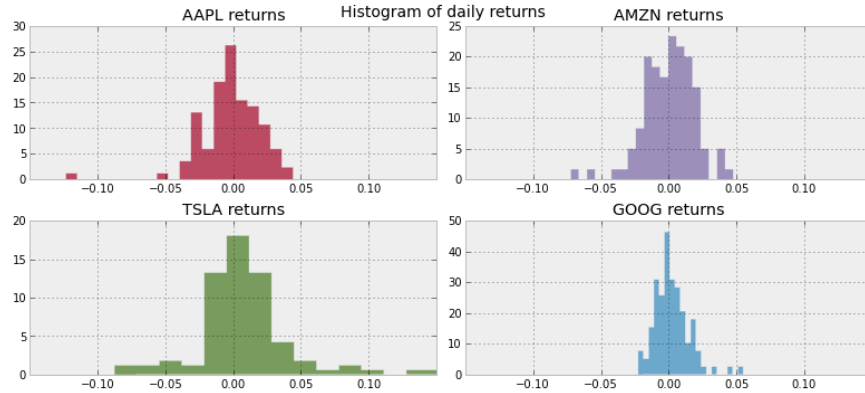
    returns[:,i] = _returns

    subplot(2,2,i)
    plt.hist( _returns, bins = 20,
              normed = True, histtype="stepfilled",
              color = colors[i], alpha = 0.7 )
    plt.title(_stock + " returns" )
    plt.xlim(-0.15, 0.15 )

plt.tight_layout()
plt.suptitle("Histogram of daily returns", size =14);

```



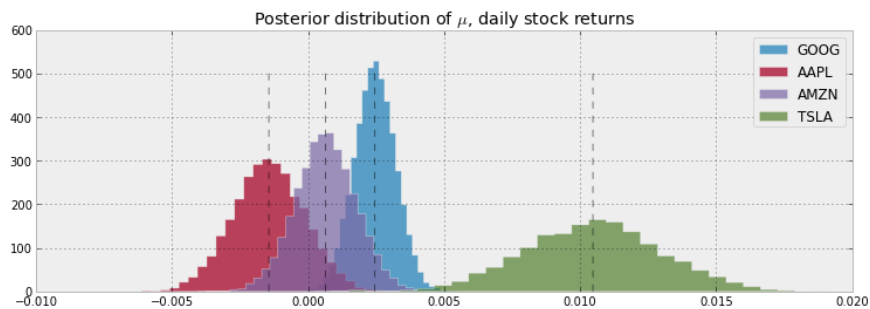


Below we perform the inference on the posterior mean return and posterior covariance matrix.

```
In [104]: obs = mc.MvNormal( "observed returns", mu, inv_cov_matrix, observed = True
model = mc.Model( [obs, mu, inv_cov_matrix] )
mcmc = mc.MCMC( )
mcmc.sample( 150000, 100000, 3 )
```

[\*\*\*\*\*100%\*\*\*\*\*] 120000 of 120000 complete

```
In [122]: figsize(12.5,4)
#examine the mean return first.
mu_samples = mcmc.trace("returns")[:]
for i in range(4):
    plt.hist(mu_samples[:,i], alpha = 0.8 - 0.05*i, bins = 30,
             histtype="stepfilled", normed=True,
             label = "%s"%stock_returns.keys()[i])
plt.vlines( mu_samples.mean(axis=0), 0, 500, linestyle="--", linewidth =
plt.title("Posterior distribution of  $\mu$ , daily stock returns" )
plt.legend();
```



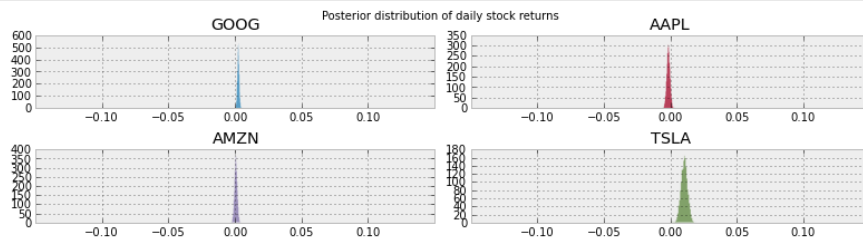
(plots like these are what inspired the book's cover btw)

What can we say about the results above? Clearly TSLA has been a strong performer, and our analysis suggests that it has an almost 1% daily return! Similarly, most of the distribution of AAPL is negative, suggesting that its *true daily return* is negative.

You may not have immediately noticed, but these variables are a whole order of magnitude *less* than our priors on them. For example, to put these one the same scale as the above prior distributions:

```
In [114]: figsize(11.0,3)
for i in range(4):
    subplot(2,2,i+1)
    plt.hist(mu_samples[:,i], alpha = 0.8 - 0.05*i, bins = 30,
            histtype="stepfilled", normed=True, color = colors[i],
            label = "%s"%stock_returns.keys()[i])
    plt.title( "%s"%stock_returns.keys()[i] )
    plt.xlim(-0.15, 0.15 )

plt.suptitle("Posterior distribution of daily stock returns" )
plt.tight_layout()
```



Why did this occur? Recall how I mentioned that finance has a very very low signal to noise ratio. This implies an environment where inference is much more difficult. One should be careful about over interpreting these results: notice (in the first figure) that each distribution is positive at 0, implying that the stock may return nothing. Furthermore, the subjective priors influenced the results. From the fund managers point of view, this is good as it reflects his updated beliefs about the stocks, whereas from a neutral viewpoint this can be too subjective of a result.

Below we show the posterior correlation matrix, and posterior standard deviations. An important caveat to know is that the Wishart distribution models the *inverse covariance matrix*, so we must invert it to get the covariance matrix. We also normalize the matrix to acquire the *correlation matrix*. Since we cannot plot hundreds of matrices effectively, we settle my summarizing the posterior distribution of correlation matrices of showing the *mean posterior correlation matrix* (defined on line 2).

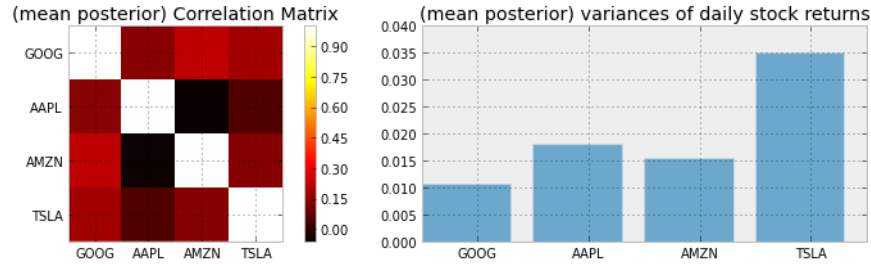
```
In [118]: inv_cov_samples = mcmc.trace("inv_cov_matrix")[:]
mean_covariance_matrix = np.linalg.inv( inv_cov_samples.mean(axis=0) )

def cov2corr( A ):
    """
    covariance matrix to correlation matrix.
    """
    d = np.sqrt(A.diagonal())
    A = ((A.T/d).T)/d
    #A[ np.diag_indices(A.shape[0]) ] = np.ones( A.shape[0] )
    return A

subplot(1,2,1)
plt.imshow( cov2corr(mean_covariance_matrix) , interpolation="none",
            cmap = plt.cm.hot )
plt.xticks( arange(4), stock_returns.keys() )
plt.yticks( arange(4), stock_returns.keys() )
plt.colorbar(orientation="vertical")
plt.title("(mean posterior) Correlation Matrix" )

subplot(1,2,2)
plt.bar(np.arange(4), np.sqrt(np.diag(mean_covariance_matrix)),
        color = "#348ABD", alpha = 0.7 )
plt.xticks( np.arange(4) + 0.5, stock_returns.keys() );
plt.title("(mean posterior) variances of daily stock returns" )

plt.tight_layout();
```



Looking at the above figures, we can say that likely TSLA has an above-average volatility (looking at the return graph this is quite clear). The correlation matrix shows that there are not strong correlations present, but perhaps GOOG and AMZN express a higher correlation (about 0.30).

With this Bayesian analysis of the stock market, we can throw it into a Mean-Variance optimizer (which I cannot stress enough, do not use with frequentist point estimates) and find the minimum. This optimizer balances the tradeoff between a high return and high variance.

$$w_{opt} = \min_w \frac{1}{N} \left( \sum_{i=0}^N \mu_i^T w - \frac{\lambda}{2} w^T \Sigma_i w \right)$$

where  $\mu_i$  and  $\Sigma_i$  are the  $i$ th posterior estimate of the mean returns and the covariance matrix. This is another example of loss function optimization.

### Protips for the Wishart distribution

If you plan to be using the Wishart distribution, read on. Else, feel free to skip this.

In the problem above, the Wishart distribution behaves pretty nicely. Unfortunately, this is rarely the case. The problem is that estimating an  $N \times N$  covariance matrix involves estimating  $\frac{1}{2}N(N - 1)$  unknowns. This is a large number even for modest  $N$ . Personally, I've tried performing a similar simulation as above with  $N = 23$  stocks, and ended up giving considering that I was requesting my MCMC simulation to estimate at least  $\frac{1}{2}23 * 22 = 253$  additional unknowns (plus the other interesting unknowns in the problem). This is not easy for MCMC. Essentially, you are asking you MCMC to traverse 250+ dimensional space. And the problem seemed so innocent initially! Below are some tips, in order of supremacy:

1. Use conjugacy if it applies. See section below.
2. Use a good starting value. What might be a good starting value? Why, the data's sample covariance matrix is! Note that this is not empirical Bayes: we are not touching the prior's parameters, we are modifying the starting value of the MCMC. Due to numerical instability, it is best to truncate the floats in the sample covariance matrix down a few degrees of precision (e.g. instability can cause unsymmetrical matrices, which can cause PyMC to cry.).
3. Provide as much domain knowledge in the form of priors, if possible. I stress *if possible*. It is likely impossible to have an estimate about each  $\frac{1}{2}N(N - 1)$  unknown. In this case, see number 4.
4. Use empirical Bayes, i.e. use the sample covariance matrix as the prior's parameter.
5. For problems where  $N$  is very large, nothing is going to help. Instead, ask, do I really care about *every* correlation? Probably not. Further ask yourself, do I really really care about correlations? Possibly not. In finance, we can set an informal hierarchy of what we might be interested in the most: first a good estimate of  $\mu$ , the variances along the diagonal of the covariance matrix are secondly important, and finally the correlations are least important. So, it might be better to ignore the  $\frac{1}{2}(N - 1)(N - 2)$  correlations and instead focus on the more important unknowns.

## 0.8.4 Conjugate Priors

Recall that a Beta prior with Binomial data implies a Beta posterior. Graphically:

$$\underbrace{\text{Beta}}_{\text{prior}} \cdot \overbrace{\text{Binomial}}^{\text{data}} = \underbrace{\text{Beta}}_{\text{posterior}}$$

Notice the Beta on both sides of this equation (no, you cannot cancel them, this is not a *real* equation). This is a really useful property. It allows us avoid using MCMC, since the posterior is known in closed form. Hence inference and analytics are easy to derive. This shortcut was the heart of the Bayesian Bandit algorithm above. Fortunately, there is an entire family of distributions that have similar behaviour.

Suppose  $X$  comes from, or is believed to come from, a well-known distribution, call it  $f_\alpha$ , where  $\alpha$  are possibly unknown parameters of  $f$ .  $f$  could be a Normal distribution, or Binomial distribution, etc. For particular distributions  $f_\alpha$ , there may exist a prior distribution  $p_\beta$ , such that:

$$\underbrace{p_\beta}_{\text{prior}} \cdot \overbrace{f_\alpha(X)}^{\text{data}} = \underbrace{p_{\beta'}}_{\text{posterior}}$$

where  $\beta'$  is a different set of parameters *but  $p$  is the same distribution as the prior*. A prior  $p$  that satisfies this relationship is called a *conjugate prior*. As I mentioned, they are useful computationally, as we can avoid approximate inference using MCMC and go directly to the posterior. This sounds great, right?

Unfortunately, not quite. There are a few issues with conjugate priors.

1. The conjugate prior is not objective. Hence only useful when a subjective prior is required. It is not guaranteed that the conjugate prior can accommodate the practitioner's subjective opinion.
2. There typically exist conjugate priors for simple, one dimensional problems. For larger problems, involving more complicated structures, hope is lost to find a conjugate prior. For smaller models, Wikipedia has a nice [table of conjugate priors](#).

Really, conjugate priors are only useful for their mathematical convenience: it is simple to go from prior to posterior. I personally see conjugate priors as only a neat mathematical trick, and offer little insight into the problem at hand.

## 0.8.5 Jefferys Priors

TODO.

## 0.8.6 Effect of the prior as $N$ increases

In the first Chapter, I proposed that as the amount of observations, or data, that we possess, the less the prior matters. This is intuitive. After all, our prior is based on previous information, and eventually enough new information will shadow our previous information's value. The smothering of the prior by enough data is also helpful: if our prior is significantly wrong, then the self-correcting nature of the data will present to us a *less wrong*, and eventually *correct*, posterior.

We can see this mathematically. First, recall Bayes Theorem from Chapter 1 that relates the prior to the posterior. The following is a sample from [What is the relationship between sample size and the influence of prior on posterior?](#)[1] on CrossValidated.

The posterior distribution for a parameter  $\theta$ , given a data set  $\mathbf{X}$  can be written as

$$p(\theta|\mathbf{X}) \propto \underbrace{p(\mathbf{X}|\theta)}_{\text{likelihood}} \cdot \overbrace{p(\theta)}^{\text{prior}}$$

or, as is more commonly displayed on the log scale,

$$\log(p(\theta|\mathbf{X})) = c + L(\theta; \mathbf{X}) + \log(p(\theta))$$

The log-likelihood,  $L(\theta; \mathbf{X}) = \log(p(\mathbf{X}|\theta))$ , **scales with the sample size**, since it is a function of the data, while the prior density does not. Therefore, as the sample size increases, the absolute value of  $L(\theta; \mathbf{X})$  is getting larger while  $\log(p(\theta))$  stays fixed (for a fixed value of  $\theta$ ), thus the sum  $L(\theta; \mathbf{X}) + \log(p(\theta))$  becomes more heavily influenced by  $L(\theta; \mathbf{X})$  as the sample size increases.

There is an interesting consequence not immediately apparent. As the sample size increases, the chosen prior has less influence. Hence inference converges regardless of chosen prior, so long as the areas of non-zero probabilities are the same.

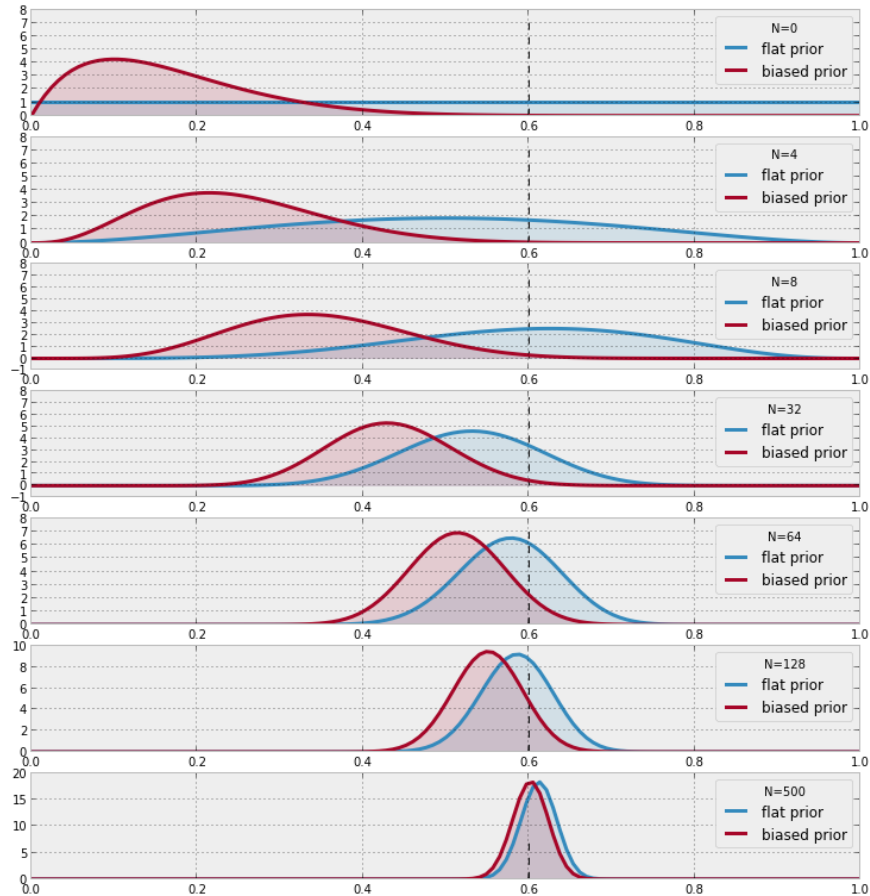
Below we visualize this. We examine the convergence of two posteriors of a Binomial's parameter  $\theta$ , one with a flat prior and the other with a biased prior towards 0. As the sample size increases, the posteriors, and hence the inference, converge.

```
In [138]: figsize( 12.5, 15)

p = 0.6
beta1_params = np.array( [1.,1.] )
beta2_params = np.array( [2,10] )
beta = stats.beta

x = np.linspace(0.00, 1, 125)
data = mc.rbernoulli(p, size=500)

figure()
for i,N in enumerate([0,4,8, 32,64, 128, 500]):
    s = data[:N].sum()
    subplot(8,1,i+1)
    params1 = beta1_params + np.array( [s, N-s] )
    params2 = beta2_params + np.array( [s, N-s] )
    y1,y2 = beta.pdf( x, *params1), beta.pdf( x, *params2)
    plt.plot( x,y1, label = r"flat prior", lw=3 )
    plt.plot( x, y2, label = "biased prior", lw= 3 )
    plt.fill_between( x, 0, y1, color = "#348ABD", alpha = 0.15)
    plt.fill_between( x, 0, y2, color = "#A60628", alpha = 0.15)
    plt.legend(title = "N=%d"%N)
    plt.vlines( p, 0.0, 7.5, linestyles = "--", linewidth=1)
    #plt.ylim( 0, 10)#
```



Keep in mind, not all posteriors will “forget” the prior this quickly. This example was just to show that *eventually* the prior is forgotten. The “forgetfulness” of the prior as we become awash in more and more data is the reason why Bayesian and Frequentist inference eventually converge as well.

### Bayesian perspective of Penalized Linear Regressions

There is a very interesting relationship between a penalized least-squares regression and Bayesian priors. A penalized linear regression is an optimization problem of the form:

$$\operatorname{argmin}_{\beta} (Y - X\beta)^T(Y - X\beta) + f(\beta)$$

for some function  $f$  (typically a norm like  $\|\cdot\|_p$ ).

We will first describe the probabilistic interpretation of least-squares linear regression. Denote our response variable  $Y$ , and features are contained in the data matrix  $X$ . The standard linear model is:

$$Y = X\beta + \epsilon \tag{49}$$

where  $\epsilon \sim \text{Normal}(\mathbf{0}, \sigma\mathbf{I})$ . Simply, the observed  $Y$  is a linear function of  $X$  (with coefficients  $\beta$ ) plus some noise term. Our unknown to be determined is  $\beta$ . We use the following property of Normal random variables:

$$\mu' + \text{Normal}(\mu, \sigma) \sim \text{Normal}(\mu' + \mu, \sigma)$$

to rewrite the above linear model as:

$$Y = X\beta + \text{Normal}(\mathbf{0}, \sigma\mathbf{I}) \quad (50)$$

$$(51)$$

$$Y = \text{Normal}(X\beta, \sigma\mathbf{I}) \quad (52)$$

$$(53)$$

$$(54)$$

In probabilistic notation, denote  $f_Y(y | \beta)$  the probability distribution of  $Y$ , and recalling the density function for a Normal random variable (see [here](#)):

$$f_Y(Y | \beta, X) = L(\beta | X, Y) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{1}{2\sigma^2}(Y - X\beta)^T(Y - X\beta)\right)$$

This is the likelihood function for  $\beta$ . Taking the log:

$$\ell(\beta) = K - c(Y - X\beta)^T(Y - X\beta)$$

where  $K$  and  $c > 0$  are constants. Maximum likelihood techniques wish to maximize this for  $\beta$ ,

$$\hat{\beta} = \text{argmax}_{\beta} -(Y - X\beta)^T(Y - X\beta)$$

Equivalently we can *minimize the negative* of the above:

$$\hat{\beta} = \text{argmin}_{\beta} (Y - X\beta)^T(Y - X\beta)$$

This is the familiar least-squares linear regression equation. Therefore we showed that the solution to a linear least-squares is the same as the maximum likelihood assuming Normal noise. Next we extend this to show how we can arrive at penalized linear regression by a suitable choice of prior on  $\beta$ .

## Penalized least-squares

In the above, once we have the likelihood, we can include a prior distribution on  $\beta$  to derive to the equation for the posterior distribution:

$$P(\beta|Y, X) = L(\beta|X, Y)p(\beta)$$

where  $p(\beta)$  is a prior on the elements of  $\beta$ . What are some interesting priors?

1. If we include *no explicit* prior term, we are actually including an uninformative prior,  $P(\beta) \propto 1$ , think of it as uniform over all numbers.
2. If we have reason to believe the elements of  $\beta$  are not too large, we can suppose that *a priori*:

$$\beta \sim \text{Normal}(\mathbf{0}, \lambda\mathbf{I})$$

The resulting posterior density function for  $\beta$  is *proportional to*:

$$\exp\left(-\frac{1}{2\sigma^2}(Y - X\beta)^T(Y - X\beta)\right) \exp\left(-\frac{1}{2\lambda^2}\beta^T\beta\right)$$

and taking the log of this, and combining and redefining constants, we arrive at:

$$\ell(\beta) \propto K - (Y - X\beta)^T(Y - X\beta) - \alpha\beta^T\beta$$

we arrive at the function we wish to maximize (recall the point that maximizes the posterior distribution is the MAP, or *maximum a posterior*):

$$\hat{\beta} = \operatorname{argmax}_{\beta} - (Y - X\beta)^T(Y - X\beta) - \alpha\beta^T\beta$$

Equivalently, we can minimize the negative of the above, and rewriting  $\beta^T\beta = \|\beta\|_2^2$ :

$$\hat{\beta} = \operatorname{argmin}_{\beta} (Y - X\beta)^T(Y - X\beta) + \alpha\|\beta\|_2^2$$

This above term is exactly Ridge Regression. Thus we can see that ridge regression corresponds to the MAP of a linear model with Normal errors and a Normal prior on  $\beta$ .

3. Similarly, if we assume a *Laplace* prior on  $\beta$ , ie.

$$f_{\beta}(\beta) \propto \exp(-\lambda\|\beta\|_1)$$

and following the same steps as above, we recover:

$$\hat{\beta} = \operatorname{argmin}_{\beta} (Y - X\beta)^T(Y - X\beta) + \alpha\|\beta\|_1$$

which is LASSO regression. Some important notes about this equivalence. The sparsity that is a result of using a LASSO regularization is not a result of the prior assigning high probability to sparsity. Quite the opposite actually. It is the combination of the  $\|\cdot\|_1$  function and using the MAP that creates sparsity on  $\beta$ : **purely a geometric argument**. The prior does contribute to an overall shrinking of the coefficients towards 0 though. An interesting discussion of this can be found in [2].

For an example of Bayesian linear regression, see Chapter 4's example on financial losses.

## References

1. Macro, . "What is the relationship between sample size and the influence of prior on posterior?." 13 Jun 2013. StackOverflow, Online Posting to Cross-Validated. Web. 25 Apr. 2013.
2. Starck, J.-L., , et al. "Sparsity and the Bayesian Perspective." *Astronomy & Astrophysics*. (2013): n. page. Print.
3. Kuleshov, Volodymyr, and Doina Precup. "Algorithms for the multi-armed bandit problem." *Journal of Machine Learning Research*. (2000): 1-49. Print.
4. Gelman, Andrew. "Prior distributions for variance parameters in hierarchical models." *Bayesian Analysis*. 1.3 (2006): 515-533. Print.
5. Gelman, Andrew, and Cosma R. Shalizi. "Philosophy and the practice of Bayesian statistics." *British Journal of Mathematical and Statistical Psychology*. (2012): n. page. Web. 17 Apr. 2013.
6. <http://jmlr.csail.mit.edu/proceedings/papers/v22/kaufmann12/kaufmann12.pdf>
7. James, Neufeld. "Reddit's"best" comment scoring algorithm as a multi-armed bandit task." *Simple ML Hacks*. Blogger, 09 Apr 2013. Web. 25 Apr. 2013.
8. Oakley, J. E., Daneshkhah, A. and O'Hagan, A. Nonparametric elicitation using the roulette method. Submitted to *Bayesian Analysis*.



9. "Eliciting priors from experts." 19 Jul 2010. StackOverflow, Online Posting to Cross-Validated. Web. 1 May. 2013. <http://stats.stackexchange.com/questions/1/eliciting-priors-from-experts>.
10. Taleb, Nassim Nicholas (2007), The Black Swan: The Impact of the Highly Improbable, Random House, ISBN 978-1400063512

```
In [1]: from IPython.core.display import HTML
def css_styling():
    styles = open("../styles/custom.css", "r").read()
    return HTML(styles)
css_styling()
```

Out[1]:

## 0.8.7 Implementation of Salisman's Don't Overfit submission

From Kaggle > In order to achieve this we have created a simulated data set with 200 variables and 20,000 cases. An 'equation' based on this data was created in order to generate a Target to be predicted. Given the all 20,000 cases, the problem is very easy to solve – but you only get given the Target value of 250 cases – the task is to build a model that gives the best predictions on the remaining 19,750 cases.

```
In [30]: import gzip
import requests
import zipfile

url = "https://dl.dropbox.com/s/lnly9gw8pb1xhir/overfitting.zip"

results = requests.get(url)
```

```
In [31]: import StringIO
z = zipfile.ZipFile(StringIO.StringIO(results.content))
#z.extractall()
```

```
In [32]: z.extractall()
```

```
In [38]: z.namelist()
```

```
['overfitting.csv']
Out [38]:
```

```
In [69]: d = z.open('overfitting.csv')
d.readline()
```

```
'case_id,train,Target_Practice,Target_Leaderboard,Target_Evaluate,var_1,var_2,var_3,var_4,
Out [69]:
```

```
In [70]: import numpy as np
```

```
In [63]: M = np.fromstring(d.read(), sep=",")
```

```
In [71]: len(d.read())
```

23919756  
Out [71]:

```
In [75]: np.fromstring?
```

```
In []:
```

```
In [1]: data = np.loadtxt("overfitting.csv", delimiter=",", skiprows=1)
```

```
In [2]: print """
There are also 5 other fields,
case_id - 1 to 20,000, a unique identifier for each row
train - 1/0, this is a flag for the first 250 rows which are the training
Target_Practice - we have provided all 20,000 Targets for this model, so y
Target_Leaderboard - only 250 Targets are provided. You submit your predic
Target_Evaluate - again only 250 Targets are provided. Those competitors v
"""
data.shape
```

There are also 5 other fields,

case\_id - 1 to 20,000, a unique identifier for each row

train - 1/0, this is a flag for the first 250 rows which are the training dataset

Target\_Practice - we have provided all 20,000 Targets for this model, so you can develop yo

Target\_Leaderboard - only 250 Targets are provided. You submit your predictions for the re

Target\_Evaluate - again only 250 Targets are provided. Those competitors who beat the 'ben

(20000L, 205L)  
Out [2]:

```
In [4]: ix_training = data[:,1] == 1
ix_testing = data[:,1] == 0

training_data = data[ ix_training, 5: ]
testing_data = data[ ix_testing, 5: ]

training_labels = data[ ix_training, 2]
testing_labels = data[ ix_testing, 2]

print "training:", training_data.shape, training_labels.shape
print "testing: ", testing_data.shape, testing_labels.shape
```

```
training: (250L, 200L) (250L,)
testing: (19750L, 200L) (19750L,)
```

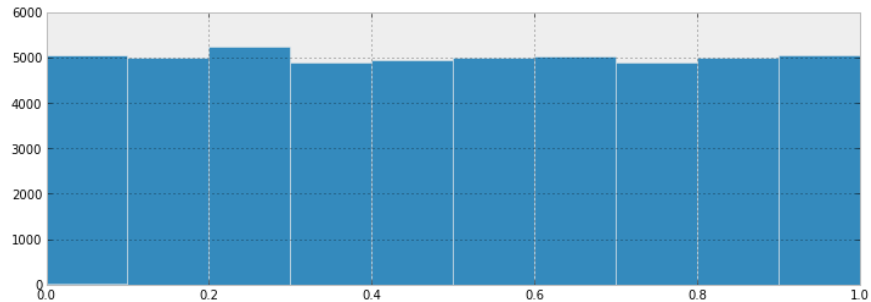
## 0.8.8 Develop Tim's model

He mentions that the X variables are from a Uniform distribution. Let's investigate this:

```
In [5]: figsize( 12, 4 )
```

```
In [6]: hist( training_data.flatten() )  
print training_data.shape[0]*training_data.shape[1]
```

50000



looks pretty right

```
In [7]: import pymc as mc  
to_include = mc.Bernoulli( "to_include", 0.5, size= 200 )
```

```
In [8]: coef = mc.Uniform( "coefs", 0, 1, size = 200 )
```

```
In [9]: @mc.deterministic  
def Z( coef = coef, to_include = to_include, data = training_data ):  
    ym = np.dot( to_include*training_data, coef )  
    return ym - ym.mean()
```

```
In [10]: @mc.deterministic  
def T( z = Z ):  
    return 0.49*(np.sign(z) + 1.02)
```

```
In [11]: obs = mc.Bernoulli( "obs", T, value = training_labels, observed = True)  
model = mc.Model( [to_include, coef, Z, T, obs] )  
map_ = mc.MAP( model )  
map_.fit()
```

Warning: Stochastic to\_include's value is neither numerical nor array with floating-point

```
In [12]: mcmc = mc.MCMC( model )
```

```
In [30]: mcmc.sample(400000, 320000, 8)
```

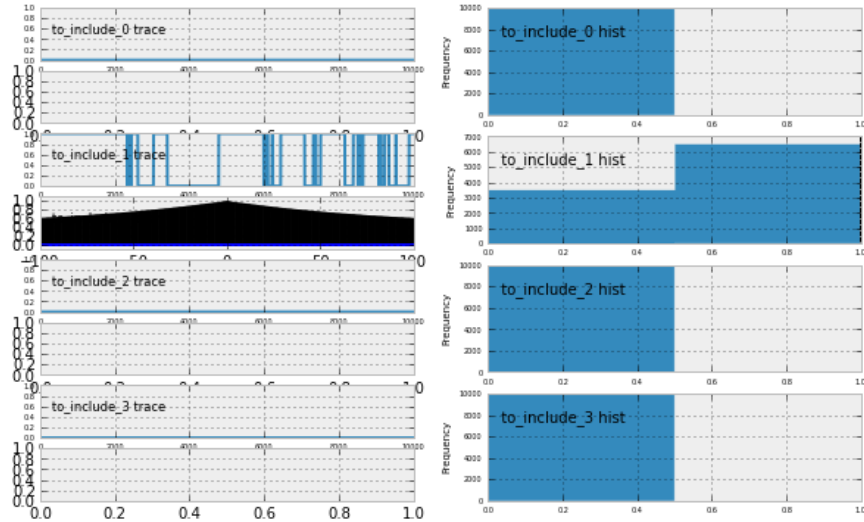
[\*\*\*\*\*100%\*\*\*\*\*] 400000 of 400000 complete

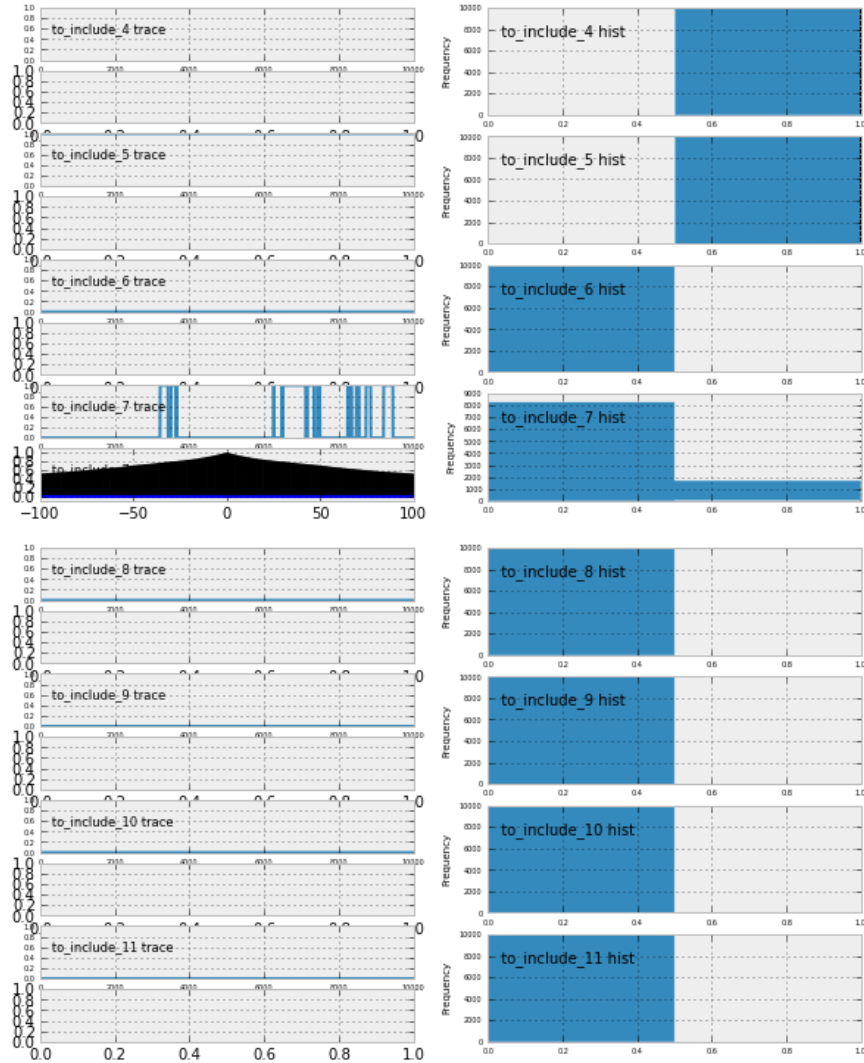
```
In [18]: from pymc.Matplot import plot as mcplot
```

```
In [31]: mcplot( to_include )
```

Plotting to\_include\_0  
 Cannot plot autocorrelation for to\_include\_0  
 Plotting to\_include\_1  
 Plotting to\_include\_2  
 Cannot plot autocorrelation for to\_include\_2  
 Plotting to\_include\_3  
 Cannot plot autocorrelation for to\_include\_3  
 Plotting to\_include\_4  
 Cannot plot autocorrelation for to\_include\_4  
 Plotting to\_include\_5  
 Cannot plot autocorrelation for to\_include\_5  
 Plotting to\_include\_6  
 Cannot plot autocorrelation for to\_include\_6  
 Plotting to\_include\_7  
 Plotting to\_include\_8  
 Cannot plot autocorrelation for to\_include\_8  
 Plotting to\_include\_9  
 Cannot plot autocorrelation for to\_include\_9  
 Plotting to\_include\_10  
 Cannot plot autocorrelation for to\_include\_10  
 Plotting to\_include\_11  
 Cannot plot autocorrelation for to\_include\_11

KeyboardInterrupt:





```
In [21]: model.T
```

-----  
 AttributeError Traceback (most recent call last)

```
<ipython-input-21-6f0fab58ada2> in <module>()
----> 1 model.T
```

```
AttributeError: 'Model' object has no attribute 'T'
```

```
In [32]: (np.round(T.value) == training_labels ).mean()
```

```
0.79600000000000004
```

```
Out [32]:
```

```
In []:
```

List of topics to cover:

- Bayesian solution to overfitting

- Salisman's solution to the Don't Overfit
- Predictive distributions; "how do I evaluate testing data?"
- model fitting, BIC + visualization tools
- Gaussian Processes

Would be nice/cool to cover:

- classification models (using the books text)
- Bayesian networks?

In [6]:

In [6]:

In [6]:

In [6]:

In [10]:

```
from IPython.core.display import HTML
def css_styling():
    styles = open("../styles/custom.css", "r").read()
    return HTML(styles)
css_styling()
```

Out[10]: