

Dawney, A. (2012), Think Bayes: Bayesian Statistics Made Simple, version 1.0.1, Green Tea Press
<http://www.greenteapress.com/thinkbayes/thinkbayes.pdf>

Chapter 2

Computational Statistics

2.1 Distributions

In statistics a **distribution** is a set of values and their corresponding probabilities.

For example, if you roll a six-sided die, the set of possible values is the numbers 1 to 6, and the probability associated with each value is $1/6$.

As another example, you might be interested in how many times each word appears in common English usage. You could build a distribution that includes each word and how many times it appears.

To represent a distribution in Python, you could use a dictionary that maps from each value to its probability. I have written a class called `Pmf` that uses a Python dictionary in exactly that way, and provides a number of useful methods. I called the class `Pmf` in reference to a **probability mass function**, which is a way to represent a distribution mathematically.

`Pmf` is defined in a Python module I wrote to accompany this book, `thinkbayes.py`. You can download it from <http://thinkbayes.com/thinkbayes.py>. For more information see Section 0.3.

To use `Pmf` you can import it like this:

```
from thinkbayes import Pmf
```

The following code builds a `Pmf` to represent the distribution of outcomes for a six-sided die:

```
pmf = Pmf()  
for x in [1,2,3,4,5,6]:  
    pmf.Set(x, 1/6.0)
```

Pmf creates an empty Pmf with no values. The Set method sets the probability associated with each value to 1/6.

Here's another example that counts the number of times each word appears in a sequence:

```
pmf = Pmf()
for word in word_list:
    pmf.Incr(word, 1)
```

Incr increases the “probability” associated with each word by 1. If a word is not already in the Pmf, it is added.

I put “probability” in quotes because in this example, the probabilities are not normalized; that is, they do not add up to 1. So they are not true probabilities.

But in this example the word counts are proportional to the probabilities. So after we count all the words, we can compute probabilities by dividing through by the total number of words. Pmf provides a method, Normalize, that does exactly that:

```
pmf.Normalize()
```

Once you have a Pmf object, you can ask for the probability associated with any value:

```
print pmf.Prob('the')
```

And that would print the frequency of the word “the” as a fraction of the words in the list.

Pmf uses a Python dictionary to store the values and their probabilities, so the values in the Pmf can be any hashable type. The probabilities can be any numerical type, but they are usually floating-point numbers (type float).

2.2 The cookie problem

In the context of Bayes's theorem, it is natural to use a Pmf to map from each hypothesis to its probability. In the cookie problem, the hypotheses are B_1 and B_2 . In Python, I represent them with strings:

```
pmf = Pmf()
pmf.Set('Bowl 1', 0.5)
pmf.Set('Bowl 2', 0.5)
```

This distribution, which contains the priors for each hypothesis, is called (wait for it) the **prior distribution**.

To update the distribution based on new data (the vanilla cookie), we multiply each prior by the corresponding likelihood. The likelihood of drawing a vanilla cookie from Bowl 1 is $3/4$. The likelihood for Bowl 2 is $1/2$.

```
pmf.Mult('Bowl 1', 0.75)
pmf.Mult('Bowl 2', 0.5)
```

Mult does what you would expect. It gets the probability for the given hypothesis and multiplies by the given likelihood.

After this update, the distribution is no longer normalized, but because these hypotheses are mutually exclusive and collectively exhaustive, we can **renormalize**:

```
pmf.Normalize()
```

The result is a distribution that contains the posterior probability for each hypothesis, which is called (wait now) the **posterior distribution**.

Finally, we can get the posterior probability for Bowl 1:

```
print pmf.Prob('Bowl 1')
```

And the answer is 0.6. You can download this example from <http://thinkbayes.com/cookie.py>. For more information see Section 0.3.

2.3 The Bayesian framework

Before we go on to other problems, I want to rewrite the code from the previous section to make it more general. First I'll define a class to encapsulate the code related to this problem:

```
class Cookie(Pmf):

    def __init__(self, hypos):
        Pmf.__init__(self)
        for hypo in hypos:
            self.Set(hypo, 1)
        self.Normalize()
```

A Cookie object is a Pmf that maps from hypotheses to their probabilities. The `__init__` method gives each hypothesis the same prior probability. As in the previous section, there are two hypotheses:

```
hypos = ['Bowl 1', 'Bowl 2']
pmf = Cookie(hypos)
```

Cookie provides an Update method that takes data as a parameter and updates the probabilities:

```
def Update(self, data):
    for hypo in self.Values():
        like = self.Likelihood(data, hypo)
        self.Mult(hypo, like)
    self.Normalize()
```

Update loops through each hypothesis in the suite and multiplies its probability by the likelihood of the data under the hypothesis, which is computed by Likelihood:

```
mixes = {
    'Bowl 1':dict(vanilla=0.75, chocolate=0.25),
    'Bowl 2':dict(vanilla=0.5, chocolate=0.5),
}

def Likelihood(self, data, hypo):
    mix = self.mixes[hypo]
    like = mix[data]
    return like
```

Likelihood uses mixes, which is a dictionary that maps from the name of a bowl to the mix of cookies in the bowl.

Here's what the update looks like:

```
pmf.Update('vanilla')
```

And then we can print the posterior probability of each hypothesis:

```
for hypo, prob in pmf.Items():
    print hypo, prob
```

The result is

```
Bowl 1 0.6
Bowl 2 0.4
```

which is the same as what we got before. This code is more complicated than what we saw in the previous section. One advantage is that it generalizes to the case where we draw more than one cookie from the same bowl (with replacement):

```
dataset = ['vanilla', 'chocolate', 'vanilla']
for data in dataset:
    pmf.Update(data)
```

The other advantage is that it provides a framework for solving many similar problems. In the next section we'll solve the Monty Hall problem computationally and then see what parts of the framework are the same.

The code in this section is available from <http://thinkbayes.com/cookie2.py>. For more information see Section 0.3.

2.4 The Monty Hall problem

To solve the Monty Hall problem, I'll define a new class:

```
class Monty(Pmf):  
  
    def __init__(self, hypos):  
        Pmf.__init__(self)  
        for hypo in hypos:  
            self.Set(hypo, 1)  
        self.Normalize()
```

So far Monty and Cookie are exactly the same. And the code that creates the Pmf is the same, too, except for the names of the hypotheses:

```
hypos = 'ABC'  
pmf = Monty(hypos)
```

Calling Update is pretty much the same:

```
data = 'B'  
pmf.Update(data)
```

And the implementation of Update is exactly the same:

```
def Update(self, data):  
    for hypo in self.Values():  
        like = self.Likelihood(data, hypo)  
        self.Mult(hypo, like)  
    self.Normalize()
```

The only part that requires some work is Likelihood:

```
def Likelihood(self, data, hypo):  
    if hypo == data:  
        return 0  
    elif hypo == 'A':  
        return 0.5  
    else:  
        return 1
```