

Lab for week 3: Working with probability distributions

Author: Sharon Goldwater, Ida Szubert, Henry S. Thompson
Date: 2014-09-01, updated 2015-10-01, 2016-09-29, 2017-09-25, 2018-08-20, 2019-09-20
Copyright: This work is licensed under a [Creative Commons Attribution-NonCommercial 4.0 International License](https://creativecommons.org/licenses/by-nc/4.0/)¹: You may re-use, redistribute, or modify this work for non-commercial purposes provided you retain attribution to any previous author(s).

This lab is available as a [web page](#)² or [pdf document](#)³.

Goals and motivation of this lab

This lab aims to build your intuitions about estimating probability distributions from data. It also provides example code, e.g., for generating bar plots and for sampling from a discrete probability distribution.

Some students are still Python beginners, so we have written most of the code already and have added a lot of explanatory comments. You'll have a chance to modify some of the code to help you understand it better.

修改

If you are new to Python programming

This lab can be done by itself. However you may also want to do 2014's [Lab 2](#)⁴ for some additional practice with Python and further heavily-commented example code. There are also instructions on how to access the Python help documentation. The [Lab 2 solutions](#)⁵ are also available, but you should work through as much as you can before looking at them.

If you are experienced with Python

We have kept the coding style here very simple, and have avoided some of the more advanced language features of Python to make the lab accessible to newcomers.

Please help those around you who may be new to Python. Learning to explain technical concepts to others is an important transferrable skill, and solidifies your own knowledge as well.

If you finish early, the 'Going Further' section suggests some more advanced topics and efficiency issues to explore.

¹<https://creativecommons.org/licenses/by-nc/4.0/>.

²<http://www.inf.ed.ac.uk/teaching/courses/anlp/labs/lab3.html>

³<http://www.inf.ed.ac.uk/teaching/courses/anlp/labs/lab3.pdf>

⁴<http://www.inf.ed.ac.uk/teaching/courses/anlp/labs/lab2.html>

⁵http://www.inf.ed.ac.uk/teaching/courses/anlp/labs/lab_solutions.html

⁶<http://www.inf.ed.ac.uk/teaching/courses/anlp/labs/lab3.py>

⁷<http://stackoverflow.com/questions/11373192/generating-discrete-random-variables-with-specified-weights-using-scipy-or-numpy>

⁸<https://wiki.python.org/moin/PythonSpeed/PerformanceTips>

Preliminaries

First, create a directory for this lab inside your `labs` directory:

```
cd ~/anlp/labs
mkdir lab3
cd lab3
```

Download the file [lab3.py](#)⁶ into your `lab3` directory: From the [Lab 3 web page](#)², right-click on the link and select *Save link as...*, then navigate to your `lab3` directory to save.

Running code in iPython

As you may already know, there are different ways to run Python code. In this class, we will use the Spyder IDE, based on the iPython interpreter. To start the IDE, just type the following into a terminal:

```
spyder3 &
```

You'll get a new full-screen window, with a text editor on the left and two smaller sub-windows on the right: A **Usage** note on top and an iPython console below.

Using the tabs at the bottom edge of the top one, change it to the **Variable explorer** view.

In the iPython console you should see a `In [1]:` prompt, and can now do things like this (where the `Out` line is the interpreter's response):

```
In [1]: 3+4
Out[1]: 7
```

Now open `lab3.py` in the editor (File -> Open...) and run the code by clicking the green arrow button on the main menu bar. If a dialogue window pops up the first time you click the arrow, please leave all options as default and click **Run** in that window.

What happened? Can you see which lines of code in `lab3.py` generated the text output you got?

Alternatively, you can run the code by typing:

```
%run lab3.py
```

in the iPython console window. If you get a `file not found` error, you may be in the wrong directory: try typing `cd ~/anlp/labs/lab3` in the interpreter and then try again.

(Note: commands that start with '%' are specific to iPython, and may not work in other Python interpreters. Also, UNIX commands like `cd` and `ls` work in iPython, but not all interpreters.)

Now uncomment the line `plot_distributions(distribution, str_probs)` near the end of the file. If you run the code using `%run lab3.py` you'll see no difference, because your change was not saved. But if you use the green arrow button, the file will be automatically saved and the new version will be used, so you will now see a plot in the output. Remember to save your changes when using `%run`.

You can also run individual commands in the interpreter. For example, to produce a histogram of some counts we generated from a uniform distribution, you could type in:

```
values = list(range(1, 10))
counts = [16, 8, 12, 12, 11, 10, 9, 12, 10]
plot_histogram(values, counts)
```

After running a code file in Spyder, you have access to all the variables and functions defined in it. In fact we just used one of them: the `plot_histogram()` function.

To see that you can also access the variables, type `distribution` in the interpreter. What is the value of this variable?

Variables are also displayed in the variable explorer window. Try double clicking on the `distribution` variable there to see what is displayed.

In the remainder of the lab we provide commands which you can run in the interpreter window of Spyder.

Getting help in iPython (and in general)

Before moving on, let's make sure you know how to get help. There are several ways to do this:

1. Simply type `help()` at the iPython interpreter prompt; you can then type a module name (e.g., `string`) and you will get some documentation. Actually a *lot* of documentation, but maybe not the most helpful kind. So, it may be better to try options 2 or 3. (Hit `q` to exit help.)
2. To get a list of all the methods available for a particular variable, you can use tab completion in iPython. This can also work to tell you all the methods available for particular datatype, like a string. For example, type the following two lines into the interpreter, but don't hit `<enter>` after the second line:

```
ss = "a string"  
ss.
```

Now type `<tab>`. You should see a list of all of the methods you can call on `ss`, which is to say, all the methods for strings.

3. You can use the `help` tab in the upper right window of spyder. Notice that you can even get help for functions defined in your own code, if you have written documentation for them. For example, click on the call to `plot_histogram` further back up in the console window, then type `Ctrl+i` and see what appears in the Help window. Or just type `help(plot_histogram)` in the interpreter. Notice where the documentation comes from in `lab3.py`.
4. Often it's most appropriate to just search the Internet, if you want more general information, don't know the name of the function you're looking for, want examples of how to use particular functions, or whatever.

Examining the initial output

When you type `distribution` into the interpreter, the output is the value of the `distribution` variable. It is a dictionary that represents a probability distribution over different characters.

Look at the top of the main body of code where the probabilities of each outcome in `distribution` are defined, and compare them to the values output by the interpreter. You might see some tiny differences. For example, the probability of `a` is defined as 0.2, but might be printed as 0.20000000000000001. Other results later in the lab might also be slightly different than you expect. This is because tiny rounding errors can happen when the computer converts numbers from base 10 (which we use) to base 2 (which the computer uses internally) and back again. Most programming languages only show numbers to five or six decimal places, so you won't see the error, but Python shows more, so you do see it.

Now, type `%run lab3.py` again and look at the plot. The plot has two sets of bars, "True" and "Est". "True" plots the probabilities of the different outcomes defined by the `distribution` variable. "Est" will eventually plot the probabilities as estimated from some data. Right now the values of the bars are incorrect.

Look now at what else was printed out in the interpreter window. By looking at the printout, the variable explorer window, and the code, can you see what `str_list`, `str_counts`, and `str_probs` are intended to be, and what datatypes they are? Which of these three variables has an incorrect value?

str prob没有 归一化

Normalizing a distribution

Look at the function `normalize_counts`. What is this function supposed to do? What is it actually doing? Fix the function so that it does what it is supposed to do. (You may want to use the `sum` function to help you.)

What is a simple error check you could perform on the return value of `normalize_counts` to make sure it is a probability distribution?

Comparing estimated and true probabilities

After you fix the `normalize_counts` function, rerun the whole file to see a comparison between two distributions. One is the *true* distribution over characters. The other is an *estimate* of that distribution which is based on the observed counts of each character in the randomly generated sequence (which in turn was generated from the true distribution). 反过来

What is the name for the kind of estimate produced here?

Now look at the plot comparing the true and estimated distributions. Notice that there are several letters that have very low probability under the true distribution. Due to random chance, some of them will occur in the randomly generated sequence and some will not. 因为随机生成的句子是依赖于真实分布生成的，因为真实分布下有的字母概率很小，所以在随机分布中，可能这样的字母就不会出现

For the low-probability letters that *do* occur in the sequence, are their estimated probabilities generally higher or lower than the true probabilities? What about for the low-probability letters that *do not* occur in the sequence?

低频出现的字母，估计概率会比真实概率高；没有出现的低频字母会显示0估计概率

Effects of sample size

Change the code so that the length of the generated sequence is much smaller or much larger than 50 and rerun the code. Are the estimated probabilities more or less accurate with larger amounts of data? Can you get estimates that no longer include zero probabilities? more accurate yes

Suppose we're estimating unigram probabilities of words (not characters) from a natural language corpus. Would it be possible to adjust the sample size to avoid zero probabilities while still using the same kind of probability estimation you have here? (That is, could you make the corpus big enough to avoid zeros?) Why or why not? 一元概率

no. 因为 有百分之50的词，只会出现一次，所有总是会有概率很少的词，生成时包含不到

Computing the likelihood

In class, we discussed the *likelihood*, defined as the probability of the observed data given a particular model. Here, our true distribution and estimated distribution are different possible models. Let's find the likelihood of each model.

First, change the code back so that you are generating sequences of length 50 again. You can also comment out the line that produces the plot, since we won't be using it again.

Next, fill in the correct code in the `compute_likelihood` function so that it returns the probability of a sequence of data (first argument) given the model provided as the second argument.

Hint: the model just tells you the unigram probability of each character, and the likelihood is $P(\text{data} \mid \text{model})$: that is, the probability of the full sequence of characters using this unigram model.

The function you just defined is being used to compute two different likelihoods using the generated sequence of 50 characters: first, the likelihood of the true distribution, and then the likelihood of the estimated distribution. Look at the values being printed out. Which model assigns higher probability to the data? By how much? (You may want to run the code a few times to see how these values change depending on the particular random sequence that is generated each time.) 假设的高 100、1000倍

Note: make sure you are familiar with the *floating-point notation* used by Python (and other computer programs): A number like $1.2\text{e}4$ means 1.2×10^4 or 12000, similarly $1.2\text{e}-4$ means 1.2×10^{-4} or .00012

Log likelihood

Increase the length of the random sequence of data to 1000 characters. You should find that your program now says both likelihoods are 0. Is that correct? Do you know why this happened? 概率太好

The problem you just saw is one practical reason for using *logs* in so many of the computations we do with probabilities. So, instead of computing the likelihood, we typically compute the *log likelihood* (or negative log likelihood).

One way to try to compute the log likelihood would be to just call the likelihood function you already wrote, and then take the log of the result. However, we don't do things that way. Why not? $\log(0)$

To correctly compute the log likelihood, fill in the body of the `compute_log_likelihood` function with code that is specific to the purpose. Please use log base 10 (see note below). *Hint: remember that $\log xy = \log x + \log y$.*

Note: For this lab, we ask you to use log base 10 because it has an intuitive interpretation. For any x that is a power of 10, $\log_{10} x$ is the number of places to shift the decimal point from 1 to get x . For example $\log_{10} 100 = 2$, and $\log_{10} .01$ is -2. Numbers falling in between powers of 10 will have non-integer logs, but rounding the log value to the nearest integer will still tell you how many decimal places the number has. Consider: what is the relationship between $\log_{10} x$ and the floating-point representation of x ? 帮助表达很小的float

Now, what is the log likelihood of your random sequence of 1000 characters?

Introduction to NumPy (optional)

You don't need to understand how most of the functions in this lab are implemented in order to do the previous parts of the lab. However, you may want to reuse/modify some of them in the future (including on your homework assignment), so it is good to understand how they work. Our code uses functions and datatypes from the NumPy (`numpy`) library. NumPy contains many things that are useful for doing numerical computing. One of its most basic parts, which we used here in several places, is the NumPy `array` datatype. An array stores a sequence of items, like a list, except that all the items must have the same type (e.g., all strings or all integers). If the items are numbers, then the array can be treated as a vector. To see how it works, let's create some arrays and lists we can manipulate. (Note that we imported `numpy` as `np` at the top of our file.)

```
l=range(4)
m=[1,0,2,3]
a=np.arange(4)
b=np.array([2,0,1,1])
c=np.array([2,3])
```

Now, try to predict what each of the following lines will do, then check to see if you are right. Note that a few of these statements will give you errors.

```
l
a
l+[1]
a+[1]
l+1
a+1
l+m
a+b
a+c
2*l
2*a
np.array(l)
a[b]
a[c]
a[m]
```

$a = \text{array}([0,1,2,3]) + [1] = \text{array}([1,2,3,4])$ 广播

```
l[m]
sum(l)
sum(a)
np.product(c)
np.cumsum(a)
np.digitize([.1, .7, 2.3, 1.2, 3.1, .3], a) 默认不包含右边区间, a是bin
```

We included the last two statements because they are used in the `generate_random_sequence` function. If you haven't already, see if you can now understand how that function works.

This section is just a tiny taste of NumPy, and even of arrays (e.g., we can create multi-dimensional arrays and use them as matrices with functions available for standard matrix operations). If you likely to do a lot of numeric programming in Python, we recommend familiarizing yourself with more of NumPy (and SciPy, another package for scientific computing in Python, which we will refer to in the Going Further section).

Going Further

1. Using either the distribution we gave you or one of your own choosing, generate a sequence of outcomes that is small enough to ensure you get some zero-count outcomes. Write a function to compute the Good-Turing estimate of the probability that the next outcome would be (any) previously unseen event. (Good-Turing is described in JM2, or you can look it up on the internet.) Do you run into any problems? Compare the GT estimate to the actual probability of getting one of the unseen events on the next draw from the distribution. Does the quality of the estimate vary depending on things like the proportion of unseen events, the sample size, or the actual probabilities of the unseen events?
2. As of Python 3.7, `dict` and all classes derived from it are guaranteed to preserve insertion order. Try removing the calls to `sorted` in the `print` calls at the end of the lab and see if everything still works as expected. What about the calls to `sorted` in `plot_distributions`? Try removing *them* and see if things are still OK. (Which version of Python is default on DICE?)

Even if you're using Python 3.7, removing these calls has a downside: it renders the code no longer backwards-compatible. How might you ensure that a version without the calls to `sorted` doesn't fail silently when run on earlier versions of Python?

3. Some algorithms in NLP and machine learning require huge numbers of random outcomes, so the efficiency of the sampling function is very important. Usually library functions are written to be efficient, but according to the [StackOverflow page](#)⁷ where I got the code that I modified to make `generate_random_sequence`, SciPy's built-in function for generating discrete random variables is much slower than (the original version of) the hand-built function here. Take a look at fraxel's code (second answer on the page; similar to mine) and EOL's code (fourth answer, using `scipy.stats.rv_discrete()`) and the comments below it. Use the `timeit` function to see if the comments are correct: is the library function really slower, and by how much? (For a fair comparison, you may need to modify code slightly. Make sure you are not timing anything *other* than random sampling: you will need to pre-compute the list of values and probabilities rather than recomputing them each time you call for a sample.)

Now consider `numpy.random.multinomial()`. It does not compute a sequence explicitly, instead returns only the total number of outcomes of each type. But, if that is all you care about, is this function more efficient than the other two?

If you pursue this question carefully, please send your code and results! It will be especially interesting to compare if we get results from multiple people.

In general, if you are interested program efficiency in Python, you may want to look at [this page](#)⁸.