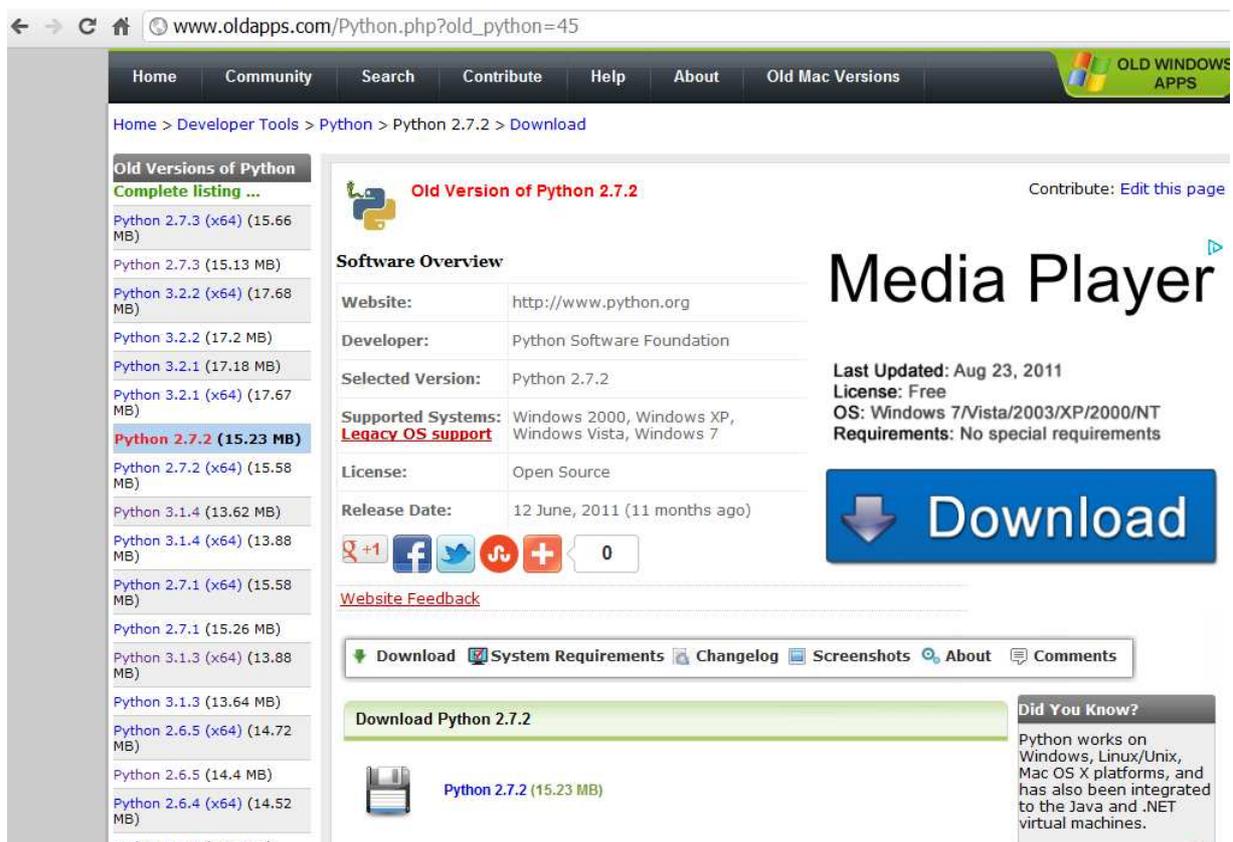


A Python program to calculate the Good-Turing frequencies

Python is a computer programming language. A complicating factor with Python is that there are many different versions, which are not 100% compatible. In particular, Python 3 is different from Python 2. Because our algorithm has been written in Python 2, we must install this version. Because this is not the most recent release, it is better not to go to the official Python site (where installing old releases is rather cumbersome), but to another site such as <http://www.oldapps.com/> (use a google search **download old releases python 2** if you can't easily find what you are looking for). Let's install Python 2.7. If you're working under Windows, you need the 32 bit version.



The screenshot shows a web browser window with the URL www.oldapps.com/Python.php?old_python=45. The page title is "Old Version of Python 2.7.2". The main content area features a "Software Overview" table with the following information:

Website:	http://www.python.org
Developer:	Python Software Foundation
Selected Version:	Python 2.7.2
Supported Systems:	Windows 2000, Windows XP, Windows Vista, Windows 7
License:	Open Source
Release Date:	12 June, 2011 (11 months ago)

Below the table, there are social media sharing icons for Google+, Facebook, Twitter, Dribbble, and a plus sign, with a counter showing 0. A large blue "Download" button is prominently displayed. To the right of the download button, the text reads: "Last Updated: Aug 23, 2011", "License: Free", "OS: Windows 7/Vista/2003/XP/2000/NT", and "Requirements: No special requirements".

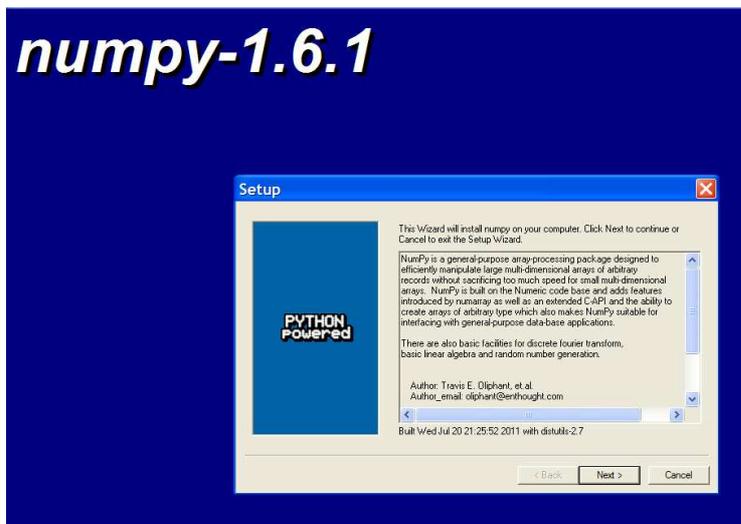
On the left side of the page, there is a sidebar titled "Old Versions of Python" with a "Complete listing ..." link. The list includes various versions of Python 2.7 and 3.2, with "Python 2.7.2 (15.23 MB)" highlighted in red. Below the sidebar, there are navigation links: "Download", "System Requirements", "Changelog", "Screenshots", "About", and "Comments".

At the bottom of the main content area, there is a "Download Python 2.7.2" section with a file icon and the text "Python 2.7.2 (15.23 MB)". To the right of this section, there is a "Did You Know?" box stating: "Python works on Windows, Linux/Unix, Mac OS X platforms, and has also been integrated to the Java and .NET virtual machines."



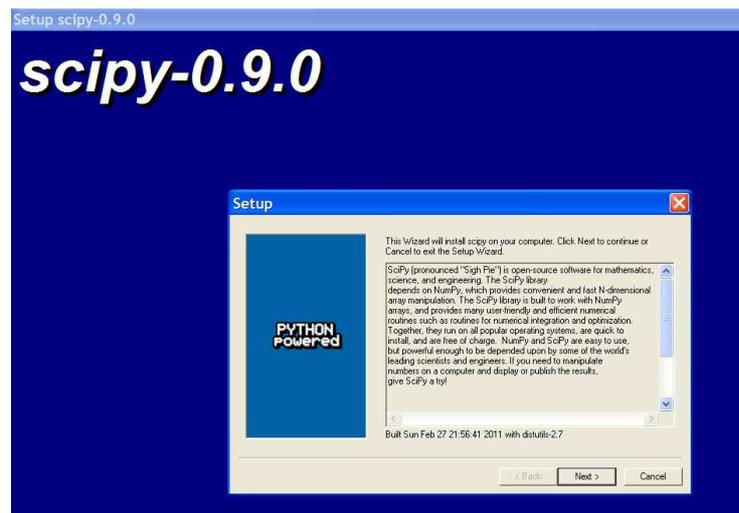
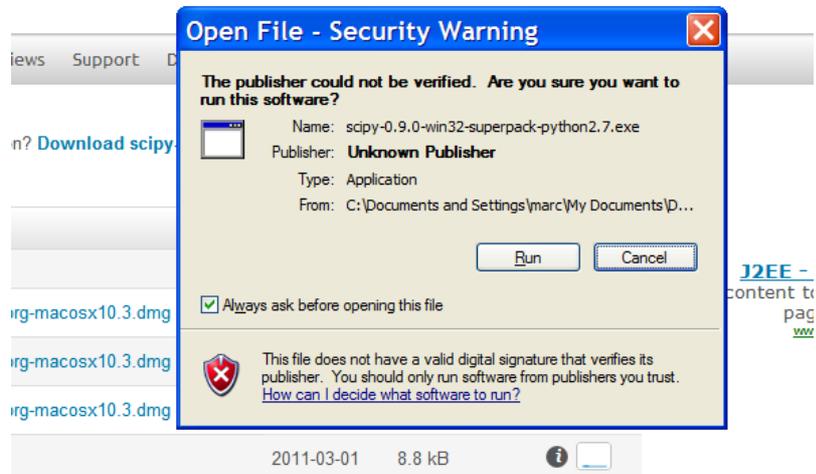
Follow the full installation instructions until you have Python 2.7 installed.

Because the Good-Turing algorithm works with mathematical libraries, you need to install two additional packages Numpy and Scipy. Be careful because these libraries are version dependent. So, make sure you use the correct version. For Python 2.7, this is currently at <http://sourceforge.net/projects/numpy/files/NumPy/1.6.1/>. The name of the file is *numpy-1.6.1-win32-superpack-python2.7.exe*



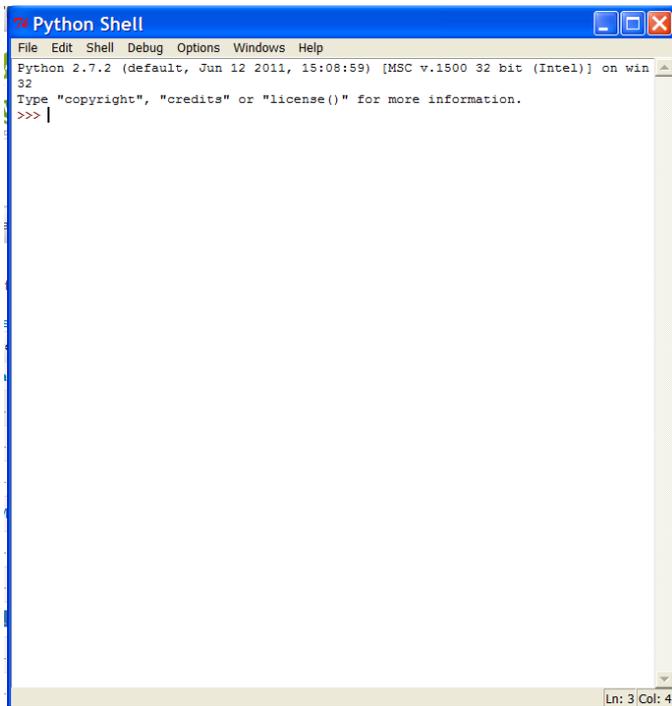
Finally, we need to install Scipy from (at the time of writing) <http://sourceforge.net/projects/scipy/files/scipy/0.9.0/> (name of the file *scipy-0.9.0-win32-superpack-python2.7.exe*)

fic Library for Python .a. edschofield, ericjones, jarrodmillman, teoliphant, tvaugt



You can now open the Python shell using Programs under the Windows Start button.

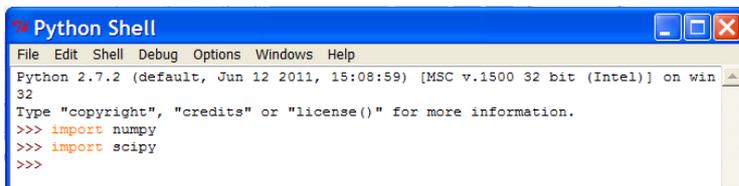




```
Python Shell
File Edit Shell Debug Options Windows Help
Python 2.7.2 (default, Jun 12 2011, 15:08:59) [MSC v.1500 32 bit (Intel)] on win
32
Type "copyright", "credits" or "license()" for more information.
>>> |
```

Ln: 3 Col: 4

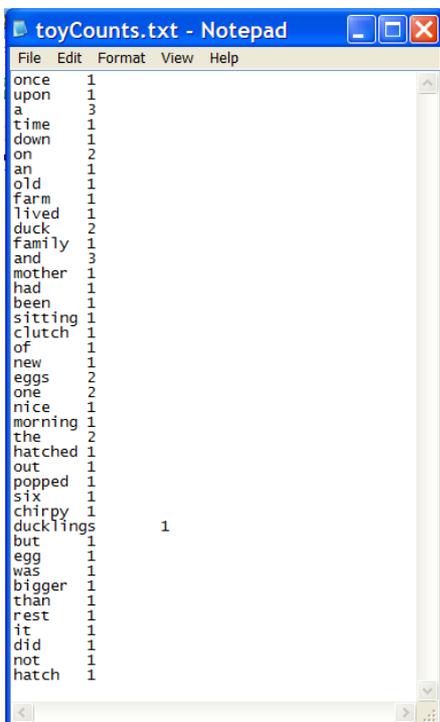
Write *import numpy* and enter; write *import scipy* and enter



```
Python Shell
File Edit Shell Debug Options Windows Help
Python 2.7.2 (default, Jun 12 2011, 15:08:59) [MSC v.1500 32 bit (Intel)] on win
32
Type "copyright", "credits" or "license()" for more information.
>>> import numpy
>>> import scipy
>>>
```

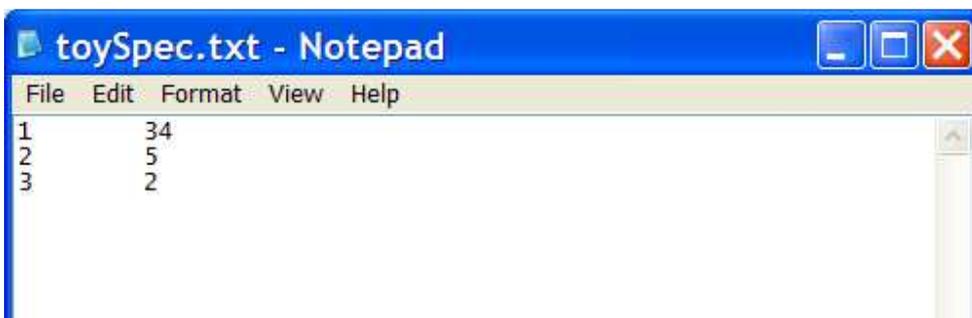
Now that you have python and the associated libraries installed, you can run the Good-Turing algorithm. To install the algorithm, save the file sgt.tar.gz or sgt.zip in the directory you want to. These are compressed files. The former is in a much-used format in the programming world, but which MS Windows does not recognize. If you have no special decompressor installed, you'll better work with the latter.

If everything went well, you should have four files in your directory: sgt.py, sgtInp.py (be careful: the middle letter is the capital I), tmpCounts.txt, and tmpSpec.txt. To these, we can add two more files related to the toy corpus we are working with in the article: toyCounts.txt and toySpec.txt. To create these, open Notepad or some other basic text processor (not MS Word!) and enter your data. Make sure there is a tab between the entries of the two columns. This is how the toyCounts.txt file looks like:



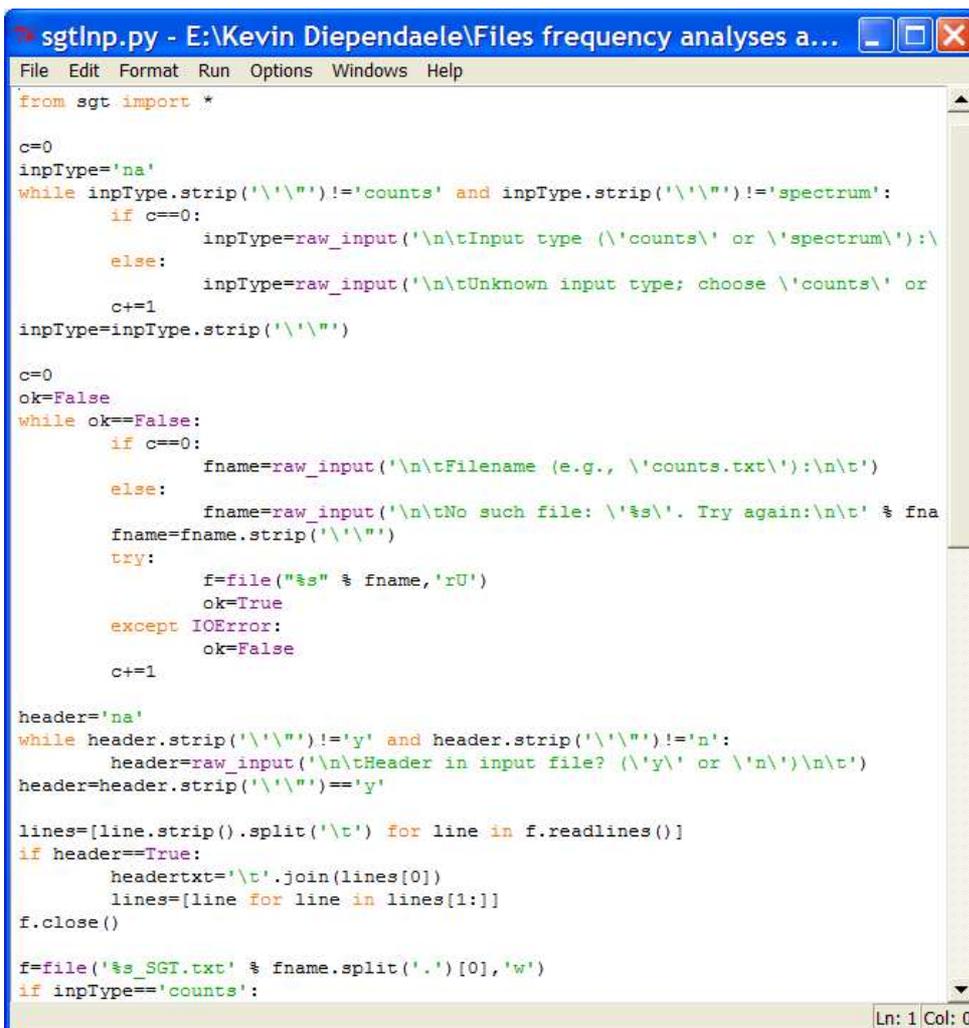
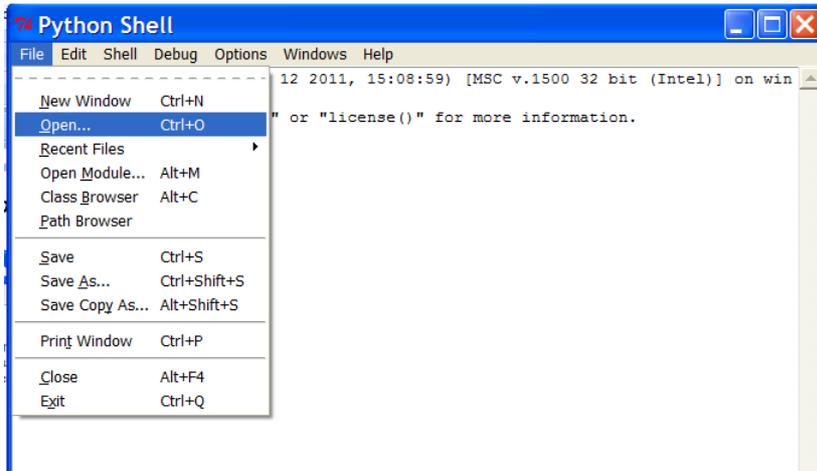
```
once 1
upon 1
a 3
time 1
down 1
on 2
an 1
old 1
farm 1
lived 1
duck 2
family 1
and 3
mother 1
had 1
been 1
sitting 1
clutch 1
of 1
new 1
eggs 2
one 2
nice 1
morning 1
the 2
hatched 1
out 1
popped 1
six 1
chirpy 1
ducklings 1
but 1
egg 1
was 1
bigger 1
than 1
rest 1
it 1
did 1
not 1
hatch 1
```

Save the file in the same directory as the file sgtInp.py is located in. Now make the toySpec.txt file. Again make sure that there is a tab between the entries on a line.

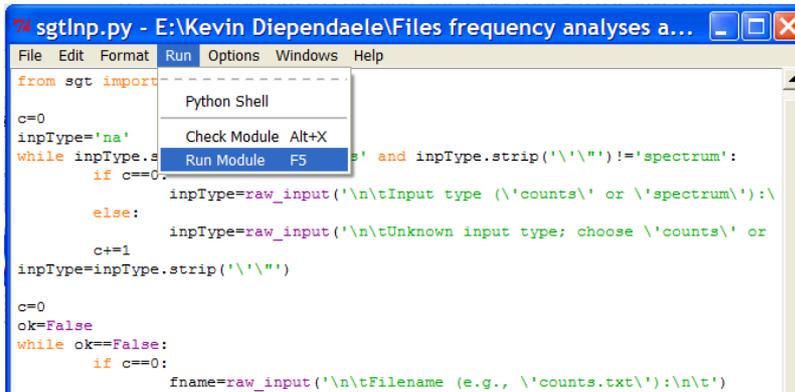


```
1 34
2 5
3 2
```

Go to your Python Shell and open the file *sgtInp.py* from the directory in which you saved it.

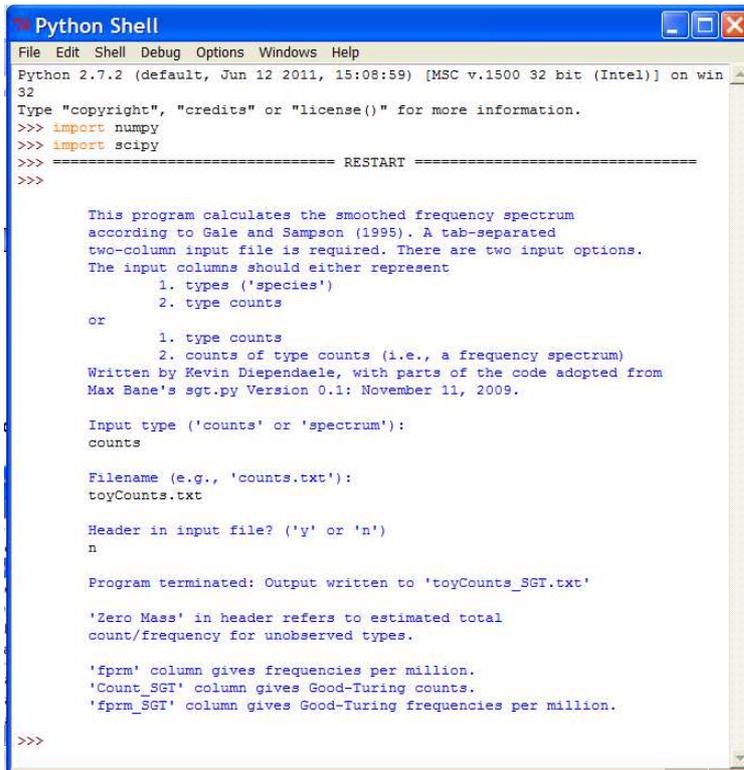


Click on Run:



```
sgtInp.py - E:\Kevin Diependaele\Files frequency analyses a...
File Edit Format Run Options Windows Help
from sgt import
c=0
inpType='na'
while inpType!='s':
    if c==0:
        inpType=raw_input('\n\tInput type (\t'counts\t' or \t'spectrum\t'):\n\t')
    else:
        inpType=raw_input('\n\tUnknown input type; choose \t'counts\t' or \t'spectrum\t'):\n\t')
    c+=1
inpType=inpType.strip('\t')
c=0
ok=False
while ok==False:
    if c==0:
        fname=raw_input('\n\tFilename (e.g., \t'counts.txt\t'):\n\t')
```

Input the required data to get the outputfile:



```
Python Shell
File Edit Shell Debug Options Windows Help
Python 2.7.2 (default, Jun 12 2011, 15:08:59) [MSC v.1500 32 bit (Intel)] on win32
Type "copyright", "credits" or "license()" for more information.
>>> import numpy
>>> import scipy
>>>
===== RESTART =====
>>>
This program calculates the smoothed frequency spectrum
according to Gale and Sampson (1995). A tab-separated
two-column input file is required. There are two input options.
The input columns should either represent
    1. types ('species')
    2. type counts
or
    1. type counts
    2. counts of type counts (i.e., a frequency spectrum)
Written by Kevin Diependaele, with parts of the code adopted from
Max Bane's sgt.py Version 0.1: November 11, 2009.
Input type ('counts' or 'spectrum'):
counts
Filename (e.g., 'counts.txt'):
toyCounts.txt
Header in input file? ('y' or 'n')
n
Program terminated: Output written to 'toyCounts_SGT.txt'
'Zero Mass' in header refers to estimated total
count/frequency for unobserved types.
'fprm' column gives frequencies per million.
'Count_SGT' column gives Good-Turing counts.
'fprm_SGT' column gives Good-Turing frequencies per million.
>>>
```

Open the file `toyCounts_SGT.txt` to see the output. In this file you get the output both as raw counts and as frequencies per million. So, the zero mass is 34/50 or 680,000 pm. The words observed once (or 20,000 pm) get a recalculated frequency of .26/50 or 5,219 pm. The words observed twice have a new frequency of .83/50 or 16,534 pm. Finally, the words observed three times have a new frequency of 1.50/50 or 29,942 pm.

type	count	fprm	count_SGT(Zero Mass=34.0)	fprm_SGT(Zero Mass=680000.0)
once	1	20000.0	0.260947758614	5218.95517228
upon	1	20000.0	0.260947758614	5218.95517228
a	3	60000.0	1.49707845262	29941.5690524
time	1	20000.0	0.260947758614	5218.95517228
down	1	20000.0	0.260947758614	5218.95517228
on	2	40000.0	0.826723860375	16534.4772075
an	1	20000.0	0.260947758614	5218.95517228
old	1	20000.0	0.260947758614	5218.95517228
farm	1	20000.0	0.260947758614	5218.95517228
lived	1	20000.0	0.260947758614	5218.95517228
duck	2	40000.0	0.826723860375	16534.4772075
family	1	20000.0	0.260947758614	5218.95517228
and	3	60000.0	1.49707845262	29941.5690524
mother	1	20000.0	0.260947758614	5218.95517228
had	1	20000.0	0.260947758614	5218.95517228
been	1	20000.0	0.260947758614	5218.95517228
sitting	1	20000.0	0.260947758614	5218.95517228
clutch	1	20000.0	0.260947758614	5218.95517228
of	1	20000.0	0.260947758614	5218.95517228
new	1	20000.0	0.260947758614	5218.95517228
eggs	2	40000.0	0.826723860375	16534.4772075
one	2	40000.0	0.826723860375	16534.4772075
nice	1	20000.0	0.260947758614	5218.95517228
morning	1	20000.0	0.260947758614	5218.95517228
the	2	40000.0	0.826723860375	16534.4772075
hatched	1	20000.0	0.260947758614	5218.95517228
out	1	20000.0	0.260947758614	5218.95517228
popped	1	20000.0	0.260947758614	5218.95517228
six	1	20000.0	0.260947758614	5218.95517228
chirpy	1	20000.0	0.260947758614	5218.95517228
ducklings	1	20000.0	0.260947758614	5218.95517228
but	1	20000.0	0.260947758614	5218.95517228
egg	1	20000.0	0.260947758614	5218.95517228

You get the same information if you run the `toySpec.txt` file:

```

Python Shell
File Edit Shell Debug Options Windows Help
count/frequency for unobserved types.

'fprm' column gives frequencies per million.
'Count_SGT' column gives Good-Turing counts.
'fprm_SGT' column gives Good-Turing frequencies per million.

>>> ===== RESTART =====
>>>

This program calculates the smoothed frequency spectrum
according to Gale and Sampson (1995). A tab-separated
two-column input file is required. There are two input options.
The input columns should either represent
  1. types ('species')
  2. type counts
or
  1. type counts
  2. counts of type counts (i.e., a frequency spectrum)
Written by Kevin Diependaele, with parts of the code adopted from
Max Bane's sgt.py Version 0.1: November 11, 2009.

Input type ('counts' or 'spectrum'):
spectrum

Filename (e.g., 'counts.txt'):
toySpec.txt

Header in input file? ('y' or 'n')
n

Program terminated: Output written to 'toySpec_SGT.txt'

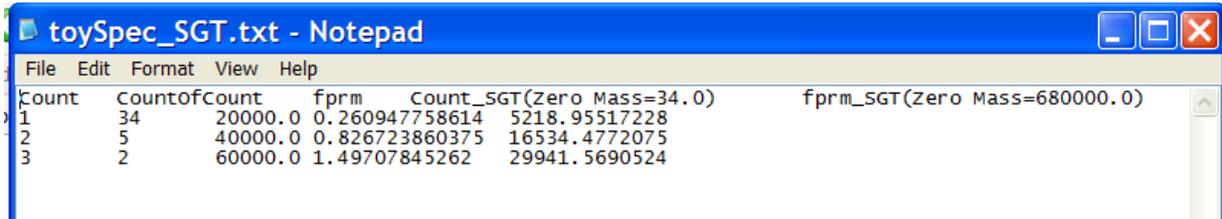
'Zero Mass' in header refers to estimated total
count/frequency for unobserved types.

'fprm' column gives frequencies per million.
'Count_SGT' column gives Good-Turing counts.
'fprm_SGT' column gives Good-Turing frequencies per million.

>>> |
Ln: 71, Col: 4
headers='na'

```

You find the information in the toySpec_SGT.txt file now arranged as a frequency spectrum (a distribution of the different frequency values).



count	CountofCount	fprm	Count_SGT(Zero Mass=34.0)	fprm_SGT(Zero Mass=680000.0)
1	34	20000.0	0.260947758614	5218.95517228
2	5	40000.0	0.826723860375	16534.4772075
3	2	60000.0	1.49707845262	29941.5690524