

Figure 2: Code for cleaning the Data





### Spam

OMG skinny coffee made me lose seven pounds in seven days! skinny coffee revolution!

Thank GOD skinny coffee ships worldwide! Don't know what I'd do without my skinny coffee!!

steveeeee Get skinny coffee challenge and Burn Excess fat without skipping your fav meals! twenty dollars!

I lost seven pounds of fat in seven days! unreal!! skinny coffee is a miracle!

Get skinny coffee challenge everywhere! ship free worldwide!

seven day challenge, lose those seven pounds, girl!

Aury9forever Get skinny coffee challenge and Burn Excess fat without skipping your fav meals! twenty dollars!

xxellttill Get skinny coffee challenge and Burn Excess fat without skipping your fav meals! twenty dollars!

Never skip your favourite meals just to lose weight, take a challenge on us.

struggling with being over weight? Lose seven pounds in seven days for twenty dollars! Free shipping worldwide!

### Not Spam

Feeling really blue today :( can't wait til my bae gets home!

TFW you wake up and the day already seems overwhelming

There is seriously nothing like coffee, the sunrise, and a good book

Some of y'all be so shallow! Don't hate me cause you ain't me!

oh my god, my boy Lamar be killin it!

If you live to be 100 you should just make up some fake reason why just to mess with peoples heads

Been on hold so long I can't remember who I even called smh

waiter, uh, theres a reflection of a sad lonely man in my soup?

Relationships are mostly you apologizing for saying something hilarious

What do you mean I didn't win I ate more wet t-shirts than anyone else

The results were spam with an entropy score of 4.2801307214326645 and not spam (ham!) with a score of 4.049611539276234 and a combined entropy of 4.2085382753968075. This doesn't make much sense to me because I was under the impression that the less alike or the greater the change in a list of strings, the greater the entropy, and inversely, the more similar the words, the less the entropy. As you can see by my text samples, the spam tweets are very similar, nigh identical, whereas the ham tweets are not. The function calculating entropy is correct because it is using the literal math of entropy and it passes a test that I designed so I am unsure as to what is going on. At first, I considered that it could possibly be the way I was reading in the txt file, and indeed at first it was problematic because I was reading in the full tweet as an item, not the individual words. But once that was ameliorated, I'm not sure what is wrong. Should I have more time at a later date, I would love to play around more with this. I will definitely be fielding this question at the next practical.

```
from __future__ import division
import math

def inputLength(spam):
    return float(len(str(spam)))

def dictCreate(spam):
    freq_dict = {}
    for key in str(spam):
        freq_dict.setdefault(key, 0)
        freq_dict[key] = freq_dict[key] + (1 / inputLength(spam))
    return freq_dict.values()

def main(spam):
    entropy = 0.0
    for v in dictCreate(spam):
        entropy = entropy + (v * math.log2(v))
    return entropy * -1
```

Fig.11: Code to find entropy of a string

## Works Cited

Amazon Web Services. (n.d.). Basic Text Mining with R. Retrieved October 12, 2017, from [https://rstudio-pubs-static.s3.amazonaws.com/132792\\_864e3813b0ec47cb95c7e1e2e2ad83e7.html](https://rstudio-pubs-static.s3.amazonaws.com/132792_864e3813b0ec47cb95c7e1e2e2ad83e7.html)

Liu, E. (2015, November 17). TF-IDF, Term Frequency-Inverse Document Frequency. Retrieved October 12, 2017, from [http://ethen8181.github.io/machine-learning/clustering\\_old/tf\\_idf/tf\\_idf.html](http://ethen8181.github.io/machine-learning/clustering_old/tf_idf/tf_idf.html)

Marmotter. (2016, April 18). Challenge #263 [Easy] Calculating Shannon Entropy of a String • r/dailyprogrammer. Retrieved October 12, 2017, from [https://www.reddit.com/r/dailyprogrammer/comments/4fc896/20160418\\_challenge\\_263\\_easy\\_calculating\\_shannon/](https://www.reddit.com/r/dailyprogrammer/comments/4fc896/20160418_challenge_263_easy_calculating_shannon/)

Turney, P. D., & Pantel, P. (2010). From frequency to meaning: Vector space models of semantics. *Journal of artificial intelligence research*, 37, 141-188.