

COS 314

NEURAL NETWORK

Project 3: Report

Author:

u11013878 - Jaco BEZUIDENHOUT

May 27, 2015

Contents

- 1 Dataset 1**
 - 1.1 Number of chunks used 1
 - 1.1.1 Afrikaans: 1
 - 1.1.2 English: 1
 - 1.2 How frequencies were calculated 1
 - 1.3 How special characters were handled 2
 - 1.4 Store format 2
 - 1.5 How the data was pre-processed 2

Chapter 1

Dataset

I have used a web scraper to scrape data from 2 websites:

- woes.co.za - Creative writing website for Afrikaans
- mibba.com - Creative writing in English

I iterated through the pages and followed all the links to stories posted on the sites. From there I scraped the body of all the stories and used that to asynchronously calculate the frequencies of each of the 26 characters in the alphabet.

On the Afrikaans website I have found that the url 'http://www.woes.co.za/bydrae/kortverhale/5' changes by 5 for every page. I then used

```
for (var page = 0; page < 600; page+=5)
```

to loop through all the pages to gather as many as possible data-sets to evaluate. The Afrikaans stories only counted to about 171 but it was enough to work with.

On the English website the paging was easier on the page 'http://www.mibba.com/Stories/?page=1' and I could just use:

```
for (var page = 0; page < 40; page++)
```

to page through all the pages. English stories was much easier to gather and here I ran through 535 stories.

1.1 Number of chunks used

I used different size chunks. All were greater than 300 characters.

1.1.1 Afrikaans:

- Training Dataset size: 81
- Generalization Dataset size: 90

1.1.2 English:

- Training Dataset size: 425
- Generalization Dataset size: 110

1.2 How frequencies were calculated

When calculating the frequencies I used the callback function of the scraper to run through all the characters in the body of the stories. For every story I made an array of 26 values. With 0 as the default value. I first had to make the entire story lowercase and then I iterated through all the characters. If a (character's ordinal value (ascii value) minus 97) were between 0 and 26 (0 included and 26 not included), then I would increment the number stored in the array at the position of the (character's ordinal; value - 97). This array was appended to the language output file (afrikaans.txt or english.txt). Only the frequencies were stored in the file. This then acted as input data to the neural network.

1.3 How special characters were handled

After some research I decided to just simply ignore the special characters for it is not a clear trademark of the English or Afrikaans language. (If I looked at for instance Russian, then special characters would definitely not be ignored.)

1.4 Store format

Example of a string in english.txt:

```
[0.07302904564315353, 0.025726141078838173, 0.02987551867219917, 0.04232365145228216, 0.1045643153526971,
0.017427385892116183, 0.02074688796680498, 0.04730290456431535, 0.0979253112033195, 0, 0.014937759336099586,
0.054771784232365145, 0.023236514522821577, 0.05975103734439834, 0.08132780082987552, 0.015767634854771784,
0, 0.030705394190871368, 0.06970954356846473, 0.07883817427385892, 0.038174273858921165, 0.01991701244813278,
0.027385892116182572, 0.0008298755186721991, 0.025726141078838173, 0],
```

Each character's frequency are represented as the percentage of the character's occurrence in the entire text. The values were converted to the percentage values before stored into the file. This then allow the scraper to use various lengths of text to summarize the frequency of certain characters per language.

1.5 How the data was pre-processed

The generalization set were made first. I tried to split the data as equally as possible (90 sets for Afrikaans and 110 sets for English).

The training set was not very equally split with 425 sets for English and 81 sets for Afrikaans.