# Text Processing using Jupyter notebook on Stampede

Abhishek Manoj Kumar, Chetan Kumar and Deepak Kumar

University of Massachusetts Dartmouth

December 20, 2016

**Abstract**

Text processing is used to extract knowledge and insights from any unstructured or raw data. We are going to classify text using python and NLTK (Natural language tool kit) library. We will be making use of Jupyter notebook which is a web application that allows to produce interactive code and output representations which will make the output reproducible and easy to share. We also configured Jupyter on stampede which is called as Jupyter Hub, so the browser and kernel(s) will be on separate machine and Jupyter Hub is a multi-user server for Jupyter Notebooks.

## 1    Introduction

Text processing also known as text analysis is the process of extracting useful information from any text document. The obtained knowledge can be useful for different purposes such as doing sentiment analysis, finding out users and products relation to enhance marketing strategies, security applications, academic applications etc. Most of the text processing work is already in existence and our key contribution is implementing text processing on Jupyter Notebook. Jupyter notebook makes it easy to code in upto 42 languages and its easy to share data, code and visualization with annotations. Also notebooks can be shared as as package by using email, Dropbox, GitHub, etc. We also installed Jupyter Notebook on Stampede and remotely accessed it with by local machine. This facilitates us to upload files or computational jobs from local machine directly to stampede using Jupyter notebook without using any code repository to upload the code by local machine and pull that on stampede. We can also use the Jupyter terminal to access stampede terminal which provides direct access to stampede. Further we used Natural Language Toolkit (NLTK) which makes use of python programs to process human language data. We are making use of text processing libraries provided by NLTK for sentence segmentation, word tokenization, parts of speech (POS) tagging and particular grammar junking. We will define each text processing term along with the implementation below. You can process text obtained from any blogs or online sources and find the information like nouns including names of people, places, relationship between them. This kind of information does have applications in security, for example by doing text processing on the social media data which have discussions about security i.e. by analyzing what people are talking about the current issues on the internet, can give us insights into the future and prevent

any threats to the national security. Another possible application would be marketing, for example by finding the relationship between people and products, they are talking about over the internet and consumer targeted marketing can be achieved.

# 2 Background

In Text processing, usually documents are processed in a certain order. The documents containing texts are processed by the natural language processing techiques in order to extract the linguistic structure of the text. Futher information is extracted using the information extraction engine. The information extracted includes sentence boundaries, part-of-speech, sentiments etc[1]. Implementation of text mining on high performance computing to process exponentially growing text data is shown in [2]. In [2] research papers were mined to extract information from the entire paper which was previously limited to the abstract, this lead to better understanding of the research as only 20 percent of the information can be found from the abstract.

# 3 Setting up Jupyter Notebook on Stampede

Anaconda consists of packages such as Python, Jupyter Notebook and commonly used scientific and data science packages. Stampede is a super computer which has Linux based operating system. The steps to install Anaconda and setting its packages are given below:

- On terminal of linux type the following command.
  $>>> wget https://3230d63b5fc54e62148e-$
  $c95ac804525aac4b6dba79b00b39d1d3.ssl.cf1.$
  $rackcdn.com/Anaconda2 - 4.0.0 - Linux - x86_64.sh$
  It will download Anaconda*.sh file which is in 64 bit format.

- Install Anaconda by using the following command
  $>>> bash Anaconda2 - 4.0.0 - Linux - x86_64.sh$
  You need to enter many times $>> yes >>$ and also can specify the directory or it will get the default directory.

- Install Python3 because dependency of Jupyterhub is on Python3
  $>>> sudo apt - get install python3 - pip$

- Install nodejs/npm
  $>>> sudo apt - get install npm nodejs - legacy$

- Install proxy with npm
  $>>> npm install - g configurable - http - proxy$

- Install Jupyterhub
  $>>> pip3 install jupyterhub$

- When we have installed the required dependencies, now we install Jupyter notebook
  $>>> pip3 install - -upgrade notebook$

Now you need to setup remote Jupyter server on stampede in order to access the Jupyter notebook by your local machine using that server. To do that we need to specify the login node because by default Stampede assigns random login node whenever someone logs in into stampede.

- $>>>$ *ssh username@login1.stampede.tacc.u texas.edu*

Here at the place of 'username' you need to write the user name of TACC as "tg837687" and we have specified login node that is 'login1', likewise you need to define a login node. You need to enter the TACC password and TACC token. TACC token can be obtained by using "TACC token" application on a smart phone. To start getting the tokens, first you need to log into your TACC account on the web and then scan the QR code using smart phone application.

Afterwards we need to specify the port on which Jupyter server will be accessed.

- $>>>$ *ipythonnotebook − −no − browser − −port = 7000*

We have specified port 7000 on which our Jupyter server is setup but we can assign any other port if the given port is not available.
The server will be running on the stampede as you can see on the terminal when you have setup the port. Now open a new terminal and create a ssh tunnel directly to stampede and connect to port 7000.

- $>>>$ *ssh−N−f−Llocalhost : 8888 : localhost : 7000username@login1.stampede.tacc.utexas.edu*

Here we specified that '8888' port on local machine will be connected to the remote server on port '7000'.
Now open up a browser and point to

- http://localhost:8888

You will see screen of ipython notebook which is running on Stampede login node. You can also open terminal in the following interface which will show the directories of stampede.

In the following figure **??** you can see the interface of Jupyter notebook running on stampede and remotely accessed by local machine. We can also upload the files using the Jupyter notebook interface directly to stampede as shown in figure **??**. We can access stampede terminal remotely and you can see the uploaded file 'prob2a' in figure **??**.

## 3.1   Installing NLTK on stampede

First on terminal of jupyter notebook run 'import NLTK' and if it is not available then follow the below steps to install NLTK library on jupyter notebook.
To install NLTK

- $>>>$ pip install -U nltk

To test if installation was successful go to jupyter terminal and write 'python import nltk'.

Figure 1: Jupyter notebook interface on local machine

# 4 Experimental Work and Results

## 4.1 HTML Cleaning

We get the HTML code which contains the tags with the text data and we remove the HTML tags in order to get the cleaned text data for processing. The result after cleaning HTML tags can be seen in figure **??**.

## 4.2 Segmentation

Sentence segmentation is the process of dividing a block of text into individual sentences and the strategy applied is based on punctuation mostly in English and some other languages. The output of the sentence segmentation can be seen in figure **??**.

## 4.3 Word Tokenization

Word tokenization refers to the process of dividing the sentence into separate words and which acts as a first step towards text processing. The output of the word tokenization is given in figure **??**.

## 4.4 Parts of Speech (POS) Tagging

In POS tagging, each token is read and part of speech(noun, pronoun, verb, adjective, adverb, preposition, conjunction and interjection) is assigned to it. The output of the POS tagging is given in figure **??**.

Figure 2: Jupyter notebook interface to upload files

## 4.5 Chunking

We define some expression such as finding out nouns phrase chunks by using output of POS tagging to find particular words. We can further make use of this information to find the relationship between entities. The result obtained after chunking can be seen in figure 1.In the figures obtained, we can see chunking done to find the simple chunking, proper noun and adjective chunking can be seen in 3 and noun phrase chunking can be seen in figure 2.

Chunking allows us to define expressions to obtain meaningful relationships. In noun phrase chunking we have defined the expression to obtain the relationship between noun and phrase. Similarly in proper noun and adjective chunking we get the relationship between noun and the adjective.

## 4.6 Extracting entities

From a given block of text we have extracted entities as seen in figure **??**.

5

Figure 3: Stampede terminal running remotely



Figure 4: HTML Cleaning

```
Out[4]: [u'Close share panel\nMedia captionMoaza Al Matrooshi: "It is a miracle that I have my son in my hand"\nA woman has given birth
        in London after doctors restored her fertility using frozen ovarian tissue removed when she was a young child.The 24-year-old i
        s thought to be the first in the world to have a baby after having an ovary frozen before the onset of puberty.Moaza Al Matroosh
        i, whose son was delivered at the privately-run Portland Hospital yesterday, told the BBC:\n"It\'s like a miracle.',
        u'"We\'ve been waiting so long for this result - a healthy baby.',
        u'"Her doctor, Sara Matthews, a consultant in gynaecology and fertility, said she was overjoyed for the family - and delighted
         by the hope it offered to others too.',
        u'"This is a huge step forward.',
        u'"We know that ovarian tissue transplantation works for older women, but we've never known if we could take tissue from a chil
        d, freeze it and make it work again.",
        u'"Doctors say it will give hope to many other girls and young women who risk losing the chance of motherhood as a result of tr
        eatment for cancer, blood or immune disorders.Frozen for the futureMoaza Al Matrooshi, who is from Dubai, was born with beta tha
        lassaemia, an inherited blood disorder that is fatal if untreated.She needed chemotherapy, which damages the ovaries, before rec
        eiving a bone marrow transplant from her brother at Great Ormond Street Hospital in London.',
        u'So, prior to treatment, when she was nine years old, she had her right ovary removed in an operation in Leeds, where the tiss
        ue was frozen.Fragments of her ovarian tissue were mixed with cryo-protective agents and slowly reduced in temperature to minus
        196C, before being stored under liquid nitrogen.Last year, surgeons in Denmark transplanted five slivers of the ovarian tissue
        back into her body - four were stitched on to her failed left ovary and one on to the side of her uterus.',
        u"Image copyright\nFERGUS WALSH\nImage caption\nMoaza's baby boy was born in a London hospital\nMoaza had been going through th
        e menopause.",
        u'But after the transplant, her hormone levels began returning to normal, she began ovulating and her fertility was restored.In
         order to maximise the chances of having a child, Moaza and her husband Ahmed underwent IVF treatment.From the eight eggs that w
        ere collected, three embryos were produced, two of which were implanted earlier this year.']
```

Figure 5: Sentence Segmentation

```
Out[5]: [u'Close',
        u'share',
        u'panel',
        u'Media',
        u'captionMoaza',
        u'Al',
        u'Matrooshi',
        u':',
        u'``',
        u'It',
        u'is',
        u'a',
        u'miracle',
        u'that',
        u'I',
        u'have',
        u'my',
        u'son',
        u'in',
        u'my',
        u'hand',
        u"''",
        u'A',
        u'woman',
        u'has',
        u'given',
        u'birth',
        u'in',
        u'London',
        u'after',
        u'doctors',
        u'restored',
```

Figure 6: Word Tokenization

```
Out[6]: [[(u'Close', 'NNP'),
        (u'share', 'NN'),
        (u'panel', 'NN'),
        (u'Media', 'NNP'),
        (u'captionMoaza', 'NN'),
        (u'Al', 'NNP'),
        (u'Matrooshi', 'NNP'),
        (u':', ':'),
        (u'``', '``'),
        (u'It', 'PRP'),
        (u'is', 'VBZ'),
        (u'a', 'DT'),
        (u'miracle', 'NN'),
        (u'that', 'IN'),
        (u'I', 'PRP'),
        (u'have', 'VBP'),
        (u'my', 'PRP$'),
        (u'son', 'NN'),
        (u'in', 'IN'),
        (u'my', 'PRP$'),
```

Figure 7: POS Tagging

7

Figure 8: Chunking



Figure 9: Noun phrase chunking



Figure 10: Proper noun and Adjective Chunking

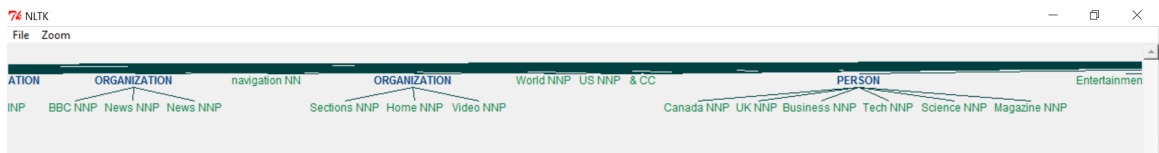

Figure 11: Extracting Entities

Figure 12: Entity Extraction Tree

# 5    Conclusion and Future Work

We have managed to successfully process text on Jupyter notebook running on Stampede, but all the computations were carried out on the login node. Looking ahead we will be trying to connect to the compute node of stampede and process even larger text documents.

# 6    Bitbucket Link

The code has been uploaded on the bitbucket, you can access the code on the given link

https://bitbucket.org/dkumar2/dsc520hpc.git

# References

[1] Miyao, Y., T. Ohta, et al. *Semantic Retrieval for the Accurate Identification of Relational Concepts in Massive Textbases. Coling/ACL, Sydney, Australia, Association for Computational Linguistics..* 2006.

[2] Firat Tekiner, Yoshimasa Tsuruoka, Junichi Tsujii *Highly scalable Text Mining - parallel tagging application.* 2006.

[3] Matthew A.Russell *Mining the Social Web.* 2006.

[4] Steven Bird, Ewan Klein, Edward Loper *Natural Language Processing with Python.* 2006.