

Analysis Of People Relationship Using Word2Vec on Wiki Data

ABHISHEK GUPTA^{1,*} AND AVADHOOT AGASTI^{1,**}

¹ School of Informatics and Computing, Bloomington, IN 47408, U.S.A.

* Corresponding authors: abhigupt@iu.edu

** Corresponding authors: aagasti@iu.edu

project-1: Data mining for a wiki url , April 19, 2017

Given a wiki URL of a person, find out his details like School, Spouse, Coaches, language, alma-mater etc Typically, the wiki page has all this information available but in the free form text. We need to converting it into structured data format so that it can help us analyze the people, from the networks etc We can create a network by navigating the people mentioned in the wiki page. © 2017 <https://creativecommons.org/licenses/>. The authors verify that the text is not plagiarized.

Keywords: Cloud, I524

<https://github.com/cloudmesh/sp17-i524/blob/master/project/S17-IR-P005/report/report.pdf>

CONTENTS

1	Introduction	1
2	Plan	1
3	Design	2
3.1	Wiki crawler	2
3.2	News crawler	2
3.3	Word2Vec model creation	2
3.4	Using the Word2Vec model to find synonyms and relations	3
4	Deployment	3
4.1	Stage1	3
4.2	Stage2	4
4.3	Stage3	4
4.4	Stage4	4
5	Benchmarking	4
6	Discussion	4
6.1	Word2Vec app - key insights	4
6.2	Deployment of Word2Vec app on Chameleon and JetStream cloud - key insights	4
7	Conclusion	4
8	Acknowledgement	5
9	Appendices	5

1. INTRODUCTION

Use spark [1] to load the wiki data and create word vectors. Train it using spark ML [2] and then use the model for analytics and prediction. The training set will use Word2Vec. Word2vec [3] is a group of related models that are used to produce word embeddings. Word2Vec is used to analyze the linguistic context of the words. In this project, we created Word2vec model using Wikipedia data. Our focus is people and organization names occurring in the Wikipedia data and to see if Word2vec can be used to understand relationship between people. Typically Wikipedia page for people and celebrities contain the entire family and friends, colleagues information. Our idea is to use Word2vec to see if using a smaller training set of known relationships whether we can derive similar relationship for anyone who has presence on Wikipedia. This mechanism can be then used to convert the data hidden in textual format to more structured data.

Technology Name	Purpose
spark [1]	data analysis
sparkML [2]	machine learning
python [1]	development
ansible [4]	automated deployment
collectd [5]	statistics collection for benchmarking

2. PLAN

Following table gives a breakdown of tasks in order to complete the project. Assuming week1 starts after submission of the pro-

positional. These work items are high level breakdown on the tasks and may change if needed.

Week	Work Item	Status
week1	Basic POC of Word2Vec using Python	planned
week2	Scripts to download Wiki data	planned
week3	Word2Vec Spark program	planned
week4	Training and measuring accuracy	planned
week5	Ansible Deployment script for Spark	planned
week6	Deployment and test on 2 clouds	planned
week7	Performance measurement	planned
week8	Report Creation(parallel)	planned

3. DESIGN

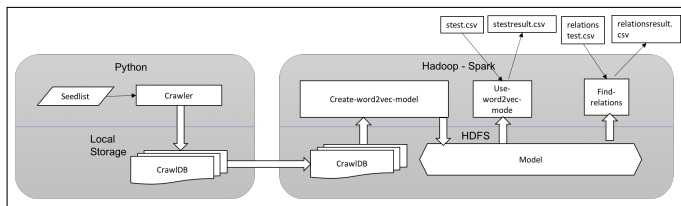


Fig. 1. Data Pipeline.

Figure 1 shows the overall data pipeline for the project. The data pipeline has three important stages:

- **Wiki crawler:** Wiki crawler runs in batch mode on a standalone machine. It can download wikipedia data as explained in section 3.1. Crawler creates Crawldb which is a collection of text files. This crawler can be replaced or augmented with any web-crawler which can download or create the text files.
- **News crawler:** News crawler is responsible for downloading the news articles.
- **CreateWord2VecModel:** This component is responsible for creating the Word2Vec model for the text files in the Crawldb. This model runs on spark and stores the model on HDFS. Section 3.3 describes this component in detail.
- **UseWord2VecModel and FindRelations:** These two components use the pre-created Word2Vec model to find synonym of a word or find the relationships. Section 3.4 describes these components in detail.

3.1. Wiki crawler

The Wiki Crawler component is useful to download the data from web. We implemented a simple crawler using Python which can deep traverse the wikipedia pages and download the text from it. In our crawler implementation, a user can specify the seed pages from wikipedia. User can also specify the maximum number of pages that are required to be downloaded. The crawler first downloads all the pages specified in the seedlist. It then extract the links from each wikipedia page and puts it in a queue which is internally maintained by the crawler. The crawler then downloads the the linked pages. Since this logic is

implemented in recursive manner, the crawler can potentially download all the wikipedia pages which can be reached from the pages in the seedlist.

We followed the seedlist based crawler approach so that we can retrieve domain specific web pages. A well chosen seedlist can fetch large number of relevant web pages.

Figure 2 is the flowchart of the crawler implementation.

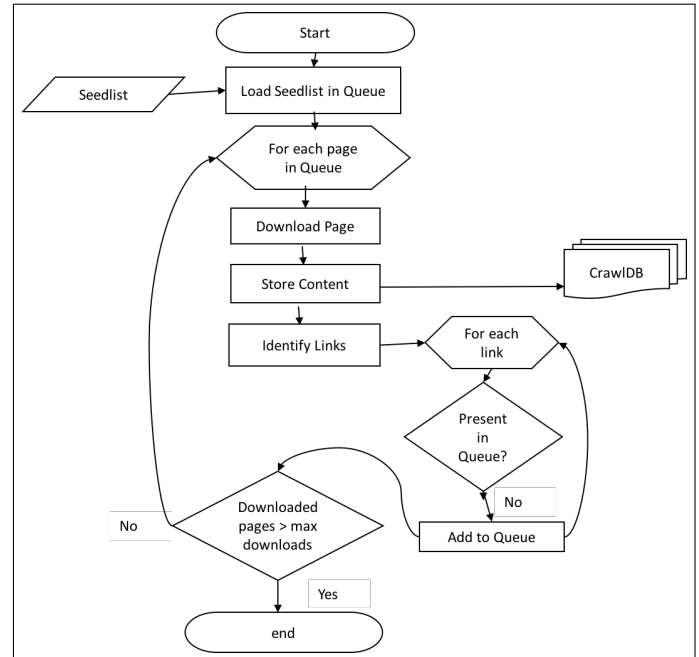


Fig. 2. Flowchart of crawler.

3.2. News crawler

News crawler is crawler implemented in Python. It executes in batch mode and download latest news article related to the topics configured in its seedlist. The news crawler uses Google APIs [6] to search the topics configured in the seedlist. Then it iterates over the result of each search, and downloads the original HTML page contents. The textual portion of the HTML is extracted using goose python library [7]

3.3. Word2Vec model creation

CreateWord2VecModel is Spark application implemented in Python. This application is responsible for creating the Word2Vec model and storing it for later use. Figure 3 explains the steps involved in the Word2Vec model creation. We used Spark Feature Extraction [8] for implementing the steps involved in the Word2Vec model creation. These steps are explained below:

- Read crawled documents from HDFS
- For each crawled document, remove special characters from the text
- Tokenize text to create list of words
- Remove the stop words
- Create Word2Vec model
- Store the model on HDFS

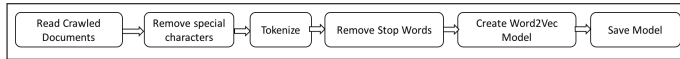


Fig. 3. Steps involved in Word2Vec model creation.

3.4. Using the Word2Vec model to find synonyms and relations

The pre-created Word2Vec model can be queried to find the synonyms or relations between the words. In the context of Word2Vec, synonym of the word is the word that co-occur in similar context [9]. The *UseWord2VecModel* finds the synonyms for the words provided in the *test.csv* file. The results are stored in *sresults.csv* file.

The Word2Vec model is also used to find the relationships. [10] explains how the vector operations can be performed on word vectors to derive relationships. The *FindRelations* spark application performs the vector operations to find relationships. The *relationtest.csv* is input to the *FindRelations* application. The application predicts the fourth word which has same relation to the third word as the second word related to first word. In the example below

Sachin, Anjali, Sourav

If *Anjali* is spouse of *Sachin*, then *FindRelations* application is expected to predict the first name of the person who is spouse of *Sourav*. The result of *FindRelations* is saved in *relationsresult.csv*.

4. DEPLOYMENT

The deployment on the cluster can be accomplished using 2 steps, assuming the cluster is up and running

- Step1: update hosts file *ansible-word2vec/hosts* with the IP of the master. The first node on the cluster becomes the master node.
- Step2: run the script *ansible-word2vec/run.sh*. This script will run the ansible playbooks to accomplish stage1 through stage4 of the deployment process.

Figure 4 shows the deployment stages. The two steps accomplish deployment in multiple stages as discussed in the sections below.

4.1. Stage1

As pre-requisite, we need to create a cluster with 1 or more nodes. We created a 3 node cluster using Cloudmesh [11] command line interface(CLI). Cloudmesh[11] CLI allows you to orchestrate virtual machines(VM) in a cloud environment. For this project, we have used Chameleon and Jetstream cloud providers to orchestrate the VMs using cloudmesh[11] CLI. We can orchestrate a 3 node cluster using following CLI:

```

cm reset
pip uninstall cloudmesh_client
pip install -U cloudmesh_client
cm key add --ssh
cm refresh on
cm cluster define --count 3 --image CC-Ubuntu14.04 --flavor m1.medium
cm hadoop define spark pig
cm hadoop sync
cm hadoop deploy
cm cluster cross_ssh
  
```

We are using Ubuntu14.04 image with m1.medium which comes with 2 CPU, 4GB memory. Also, the nodes created are

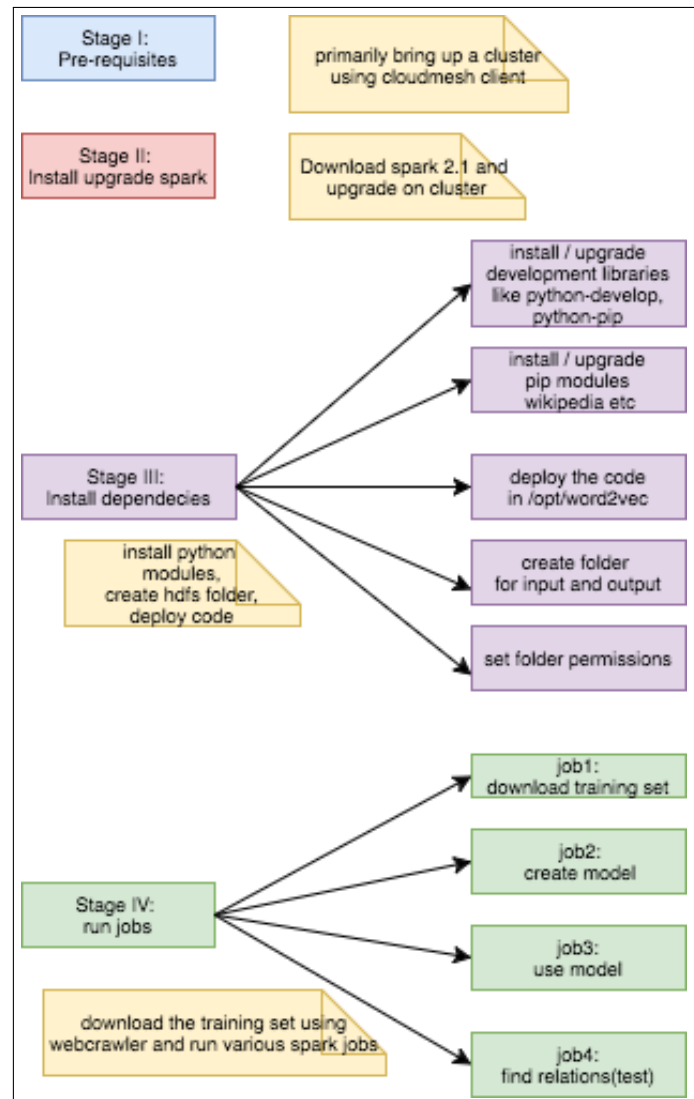


Fig. 4. Deployment stages.

having hadoop and spark add-ons. We can test the deployment by checking hdfs and spark-submit CLI work fine.

```
ssh cc@<cluster-ip>
sudo su - hadoop
hdfs
spark-submit
```

At this stage our cluster is ready for further deployments.

4.2. Stage2

By default, cloudmesh installs spark 1.6 but our word2vec solution requires spark 2.1. We need to upgrade spark on the cluster. In order to do so we can run *install_upgrade_spark.yaml* ansible[4] playbook. This will download and unpack spark2.1 tar ball and further update the softlink to point to spark 2.1 folder.

4.3. Stage3

In this step, we upgrade the development libraries for python, and pip, install python modules like wikipedia, request etc, download the code from git repo and install it in */opt/word2vec* folder, set the folder permissions for the */opt/word2vec* folder so that it can be executed by hadoop user. These steps can be achieved using *word2vec_setup.yaml* playbook. After completing this stage, we are ready for running our word2vec solution on the cluster.

4.4. Stage4

This stage primarily deals with submitting the jobs for various purpose. Before we submit the jobs, we need to make sure input folder are created on hdfs. First, we run the crawler to download the training set and upload the data on hdfs. Further we run various jobs to create model and find relations. Along with these jobs we also run some monitoring jobs. The monitoring job queries spark metrics using

```
http://${spark_master}:4040
```

Stage4 steps can be accomplished using *word2vec_execute.yaml* playbook.

At the end of stage4 we also fetch the execution results from the cluster along with the metrics of execution times at various stages. The output files are fetched into */tmp/word2vec_results*

```
ls -lrt /tmp/word2vec_results
```

```
total 56
```

```
-rw-rw-r-- 1 abhigup4 wheel 571 Apr 18 17:36 jobs.csv
-rw-rw-r-- 1 abhigup4 wheel 418 Apr 18 17:36 executors.csv
-rw-rw-r-- 1 abhigup4 wheel 89 Apr 18 17:36 app.csv
-rwxrwxr-x 1 abhigup4 wheel 43 Apr 18 23:43 stest.csv
-rwxrwxr-x 1 abhigup4 wheel 145 Apr 18 23:43 relationsresult.csv
-rw-r--r-- 1 abhigup4 wheel 1833 Apr 18 23:46 stestresult.csv
-rw-r--r-- 1 abhigup4 wheel 844 Apr 18 23:47 relationsresult.csv
```

Files *jobs.csv*, *executors.csv*, and *app.csv* collect the execution time for various jobs. File *relationsresult.csv* file collects the results for sample relations corresponding to *relationstest.csv*. Similarly *stestresult.csv* collects the results corresponding to *stest.csv*.

5. BENCHMARKING

Solution will use collectd [5] to collect statistics. Once the solution is deployed to the cluster. We should benchmark parameters like

- cpu

- memory
- throughput reads/writes

Benchmarking will be done for one or more cloud providers. The deployment scripts should be agnostic to the cloud provider.

6. DISCUSSION

The analysis of the results provided interesting insights. Section 6.1 provides key insights on our experiments with Word2Vec on Wikipedia data. Section 6.2 provide key insights on our experience of deployment and execution of Word2Vec application on Chameleon and Jetstream cloud.

6.1. Word2Vec app - key insights

We trained the Word2Vec model with the wikipedia data of Indian Cricket team members. We used *List of India ODI cricketers* as seedlist and downloaded 500 pages from Wikipedia. With this dataset, we observed that the *synonym* or correlated words were pretty accurate. Some interesting examples and its explanation:

```
sachin -> tendulkar
```

```
world -> cup
```

Sachin Tendulkar [12] is famous Indian cricket player. Naturally, the words *sachin* and *tendulkar* are highly correlated. The Word2Vec model was able to find this correlation.

Cricket World Cup [13] is one of the most viewed sporting event and considered as the flagship event of the international cricket. Hence the words *world* and *cup* are highly correlated in the context of cricket. The Word2Vec model was able to find this correlation.

However, when we tried to find the people relationships using the wikipedia dataset, we could not get good accuracy. For example, *Anjali Tendulkar* is wife of *Sachin Tendulkar* while *Dona Ganguly* is wife of *Sourav Ganguly* who is another famous Indian Cricket player [14]. When we provided the test record of

```
sachin,anjali,sourav
```

we did not get good results.

After observing the wikipedia data, we concluded that there is lot of literature in wikipedia about the cricketers. However, there were very few mentions of their family members, coaches, schools etc. Due to this, we were not getting good results on the people relationships.

To augment the Wikipedia data, we decided to crawl news data which provide large amount of articles containing people and their relationships. We implemented the news crawler as explained in the section 3.2. After downloading about 200 news articles of 20 cricketers, our relationships discovery improved a lot. Some interesting examples are given below:

```
sachin,anjali,sourav,dona
```

```
sachin,anjali,dhoni,sakshi
```

```
sachin,cricket,amitabh,hero
```

In first two examples, the Word2Vec model was able to identify the people-people relationships, while in third example, the Word2Vec model was able to identify the people-profession relationships.

6.2. Deployment of Word2Vec app on Chameleon and Jet-Stream cloud - key insights

7. CONCLUSION

Using this wiki analysis we should be able to build a network based on wiki data.

8. ACKNOWLEDGEMENT

We acknowledge our professor Gregor von Laszewski and all associate instructors for helping us and guiding us throughout this project.

9. APPENDICES

TBD

REFERENCES

- [1] "Spark Python API (PySpark)," Web Page, accessed: 2017-02-26. [Online]. Available: <https://spark.apache.org/docs/0.9.1/python-programming-guide.html>
- [2] "Spark ml programming guide," Web Page, accessed: 2017-02-26. [Online]. Available: <https://spark.apache.org/docs/1.2.2/ml-guide.html>
- [3] "Word2Vec, learning vector representation of words," Web Page, accessed: 2017-02-26. [Online]. Available: <https://en.wikipedia.org/wiki/Word2vec>
- [4] Red Hat, Inc., "Ansible documentation," Web Page, Jan. 2015, accessed 2017-01-13. [Online]. Available: <https://docs.ansible.com/ansible/index.html>
- [5] "collectd - the system statistics collection daemon," Web Page, accessed: 2017-02-26. [Online]. Available: <https://collectd.org/>
- [6] "Google custom search," Web Page. [Online]. Available: <https://developers.google.com/custom-search/>
- [7] "Python-goose - article extractor," Code Repo. [Online]. Available: <https://github.com/grangier/python-goose>
- [8] "Extracting, transforming and selecting features," Web Page, accessed: 2017-04-26. [Online]. Available: <https://spark.apache.org/docs/latest/ml-features.html>
- [9] Y. Goldberg and O. Levy, "word2vec explained: deriving mikolov et al.'s negative-sampling word-embedding method," *CoRR*, vol. abs/1402.3722, 2014.
- [10] "Vector representations of words," Web Page, accessed: 2017-04-26. [Online]. Available: <https://www.tensorflow.org/tutorials/word2vec>
- [11] "Cloudmesh client," Code Repo. [Online]. Available: <https://github.com/cloudmesh/client>
- [12] "Sachin tendulkar," Web Page, accessed: 2017-04-15. [Online]. Available: https://en.wikipedia.org/wiki/Sachin_Tendulkar
- [13] "Cricket world cup," Web Page, accessed: 2017-04-15. [Online]. Available: https://en.wikipedia.org/wiki/Cricket_World_Cup
- [14] "Sourav ganguly," Web Page, accessed: 2017-04-15. [Online]. Available: https://en.wikipedia.org/wiki/Sourav_Ganguly