# ILS - Z 604: Assignment #3

Due on Tuesday, April 19, 2016

*Prof. Xiazhong Liu*

**vpatani (Vivek Patani)**

# Contents

vpatani (Vivek Patani)     ILS - Z 604 (Prof. Xiazhong Liu ): Assignment #3

Page 2 of 10

# Task 1

Listing 1 Shows the IPynb Script.

Listing 1: Shows the IPynb Script for cleaning

```perl
#!/usr/bin/perl

#Part 1
import os
import graphlab
graphlab.product_key.set_product_key('Your API Key')
location = 'http://www.vivekpatani.tk/resources/reuters.csv'
sf = graphlab.SFrame.read_csv(location, header=False)
```

- How to configure the graphLab for AWS:

  - Register as a student and Dato will send you a key.

  - Register the key as $graphlab.product-key$ function as shown above.

  - Execute in order to register



Listing 2 Shows the IPynb Script.

Listing 2: Shows the IPynb Script for counting

```perl
#Part 2
word_list = graphlab.text_analytics.count_words(sf['X2'])
docs = sf['word_list'].dict_trim_by_values(2)
docs = docs.dict_trim_by_keys(graphlab.text_analytics.stopwords(),exclude=True)
```

- Now we count the words and eliminate the stop words and words not crossing threshold

  - Use the $SFrame.read-csv$ function giving input of the location of your file.

  - View sf to confirm the number of rows.

– Once we have the data ready, we use the libraries to count the words that exist in the dataset. We use an Sarray.SArray to store the word and their count.

– The next step is to see what words have appeared lesser than 2 times and we pass a parameter as 2 to the function $trim - by - value$.

– Next step is to eliminate the stop word. We pass *doc* as the input and the the output is over ridden on *doc* itself by passing it through *stopwords* function.

– Just to make sure we've got it right, just check *doc*[0].



• This completes the first part of **cleaning data** and **generating features**.
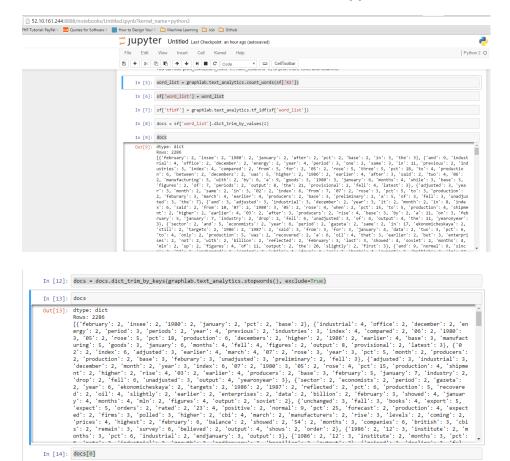
## Task 2

Listing 3 Shows the IPynb Script.

Listing 3: Shows the IPynb Script for generating Features

```
#Part 3
word_list = graphlab.text_analytics.count_words(sf['X2'])
sf['word_list'] = word_list
sf['tfidf'] = graphlab.text_analytics.tf_idf(sf['word_list'])
```

- Here the we simply create a list(bag) of words for each document.

- This is straightforward since we just generate a word list by the $count_words$ function by giving input as the the feature column over which we need to do this.

- Next step is to add it the frame generated by us.

- Next run the TF-IDF from the analytics library of graph lab over the word list generated.

- We generated two things

  1. Word List
  2. TF-IDF Values

- This is how the Frame will look like



- Now, we have attached the **features** that will help us analyse the data.

# Task 3

Listing 4 Shows the IPynb Script.

Listing 4: Shows the IPynb Script for generating Model & Evaluating it

```
#Part 4
train_data, test_data = sf.random_split(0.8)
test_data
len(sf)
len(test_data)
len(train_data)
```

- In this we have cleaned the data and build the elements required to build a model.

- We first divide the set into **Train** & **Test** in order to build and evaluate the model. We split the data in the ratio of 80:20, where 80 is train and 20 is test.

- This is how the Train & Test Frame will look like

In [17]: `train_data, test_data = sf.random_split(0.8)`

In [19]: `test_data`

Out[19]:

| X1 | X2 | word_list | tfidf |
|---|---|---|---|
| ipi | Swiss industrial output rose nine pct in the ... | {'and': 4, 'industrial': 1, 'stood': 1, 'office': ... | {'and': 0.5586978285529463, ... |
| ipi | Swedish industrial production rose 26 pc ... | {'and': 2, 'spell': 1, 'all': 1, 'industrial': ... | {'and': 0.2793489142764731, ... |
| ipi | US industrial production rose 05 pct in February ... | {'six': 1, 'month': 3, 'paper': 1, 'still': 1, ... | {'six': 2.5034502275001693, ... |
| ipi | The Bank of France expects a continued ... | {'sector': 4, 'all': 2, 'unemployment': 1, ... | {'sector': 12.360671780853533, ... |
| ipi | Canadas industrial product price index rose ... | {'yearly': 1, 'and': 3, 'industrial': 1, ... | {'yearly': 5.169609486893219, 'a ... |
| ipi | British manufacturers expect output to grow ... | {'rating': 1, 'all': 1, 'consider': 1, 'polled': ... | {'rating': 6.348264483234865, 'a ... |
| ipi | Industrial production rose 48 pct on a ... | {'and': 1, 'industrial': 3, 'show': 1, 'decemb ... | {'and': 0.13967445713823656, ... |
| ipi | Japans preliminary industrial production ... | {'and': 1, 'adjusted': 3, 'industrial': 1, 'goo ... | {'and': 0.13967445713823656, ... |
| ipi | Swedish industrial production rose 15 pc ... | {'sector': 1, 'and': 2, 'industrial': 1, ... | {'sector': 3.0901679452133832, ... |
| ipi | A leading Soviet economist said the ... | {'years': 1, 'inflating': 1, 'economist': 1, ... | {'years': 1.8624410548793402, ... |

[? rows x 4 columns]
Note: Only the head of the SFrame is printed. This SFrame is lazily evaluated.
You can use len(sf) to force materialization.

- The length of each is

In [20]: `len(sf)`
Out[20]: 2286

In [21]: `len(test_data)`
Out[21]: 450

In [22]: `len(train_data)`
Out[22]: 1836

- We use train to learn the model and data.

---

- Test will be used to evaluate features of the model such as accuracy.

- Now we start building models with different approaches:

  1. Boosted Tree Classifier: Listing 5 Shows the IPynb Script.

     Listing 5: Shows the IPynb Script for generating Model & Evaluating it

```
model = graphlab.boosted_trees_classifier.create(train_data, target='X1',
                                                 max_iterations=2,
                                                 max_depth = 3)
predictions = model.classify(test_data)
results = model.evaluate(test_data)
```

– After training the data with the given model.
– Output:

```
Boosted trees classifier:

-------------------------------------------------------

Number of examples        : 1756
Number of classes         : 34
Number of feature columns    : 1
Number of unpacked features : 17269
Create disk column page 1/1

+-----------+--------------+-------------------+-------------------+---------------------+---------------------+
| Iteration | Elapsed Time | Training-accuracy | Training-log_loss | Validation-accuracy | Validation-log_loss |
+-----------+--------------+-------------------+-------------------+---------------------+---------------------+
| 1         | 10.136764    | 0.876993          | 1.235948          | 0.575000            | 2.121445            |
| 2         | 19.866346    | 0.890091          | 0.934086          | 0.525000            | 1.927115            |
| 3         | 29.492107    | 0.899772          | 0.741076          | 0.537500            | 1.795780            |
| 4         | 39.086573    | 0.910592          | 0.603383          | 0.512500            | 1.716426            |
| 5         | 48.789179    | 0.917426          | 0.499902          | 0.512500            | 1.668092            |
| 6         | 58.514369    | 0.919134          | 0.421551          | 0.525000            | 1.626862            |
| 7         | 68.342350    | 0.921982          | 0.361099          | 0.525000            | 1.573787            |
| 8         | 78.620567    | 0.922551          | 0.314064          | 0.537500            | 1.547404            |
| 9         | 89.087766    | 0.923690          | 0.275237          | 0.537500            | 1.539655            |
| 10        | 99.270443    | 0.924260          | 0.245745          | 0.537500            | 1.526494            |
+-----------+--------------+-------------------+-------------------+---------------------+---------------------+

PROGRESS: Creating a validation set from 5 percent of training data. This may take a while.
         You can set ``validation_set=None`` to disable validation tracking.
```

– After running the test, the accuracy is 71%:

– Output:

```
In [12]: results = model.evaluate(test_data)

         External memory mode: 1 batches


In [13]: results

Out[13]: {'accuracy': 0.7103004291845494,
          'auc': 0.9698714070751959,
          'confusion_matrix': Columns:
                  target_label     str
                  predicted_label  str
                  count    int

          Rows: 88

          Data:
          +----------------+------------------+-------+
          |  target_label  | predicted_label  | count |
          +----------------+------------------+-------+
          |    pet-chem    |     pet-chem     |   3   |
          |    veg-oil     |  oilseedsoybean  |   4   |
          |    graincorn   |     graincorn    |   8   |
          |      jobs      |       gnp        |   2   |
          |      jobs      |       jobs       |   5   |
          |  dlrmoney-fx   | interestmoney-fx |   1   |
          |    tradebop    |     tradebop     |   2   |
          |      gnp       |       cpi        |   2   |
          |    reserves    |       gold       |   2   |
          | money-fxdlryen |   money-fxdlr    |   6   |
          +----------------+------------------+-------+
          [88 rows x 3 columns]
          Note: Only the head of the SFrame is printed.
          You can use print_rows(num_rows=m, num_columns=n) to print more rows and columns.,
          'f1_score': 0.6662522123798048,
          'log_loss': 1.4751669587448086,
          'precision': 0.7308917606173222,
          'recall': 0.6861378490790255,
          'roc_curve': Columns:
                  threshold    float
                  fpr      float
                  tpr      float
                  p        int
                  n        int
                  class    int

          Rows: 3400034
```

– The prediction excerpt
– Output:

In [16]: predictions

Out[16]:

| class | probability |
|---|---|
| ipi | 0.463541865349 |
| ipi | 0.495206207037 |
| ipi | 0.49030277133 |
| ipi | 0.490942507982 |
| ipi | 0.494604974985 |
| ipi | 0.397208571434 |
| ipi | 0.0526305325329 |
| alum | 0.338211506605 |
| ipi | 0.490942507982 |
| reserves | 0.563541531563 |

[466 rows x 2 columns]
Note: Only the head of the SFrame is p
You can use print_rows(num_rows=m,