

Python Libraries for Data Science

Python for Data Science

Gopi Subramanian

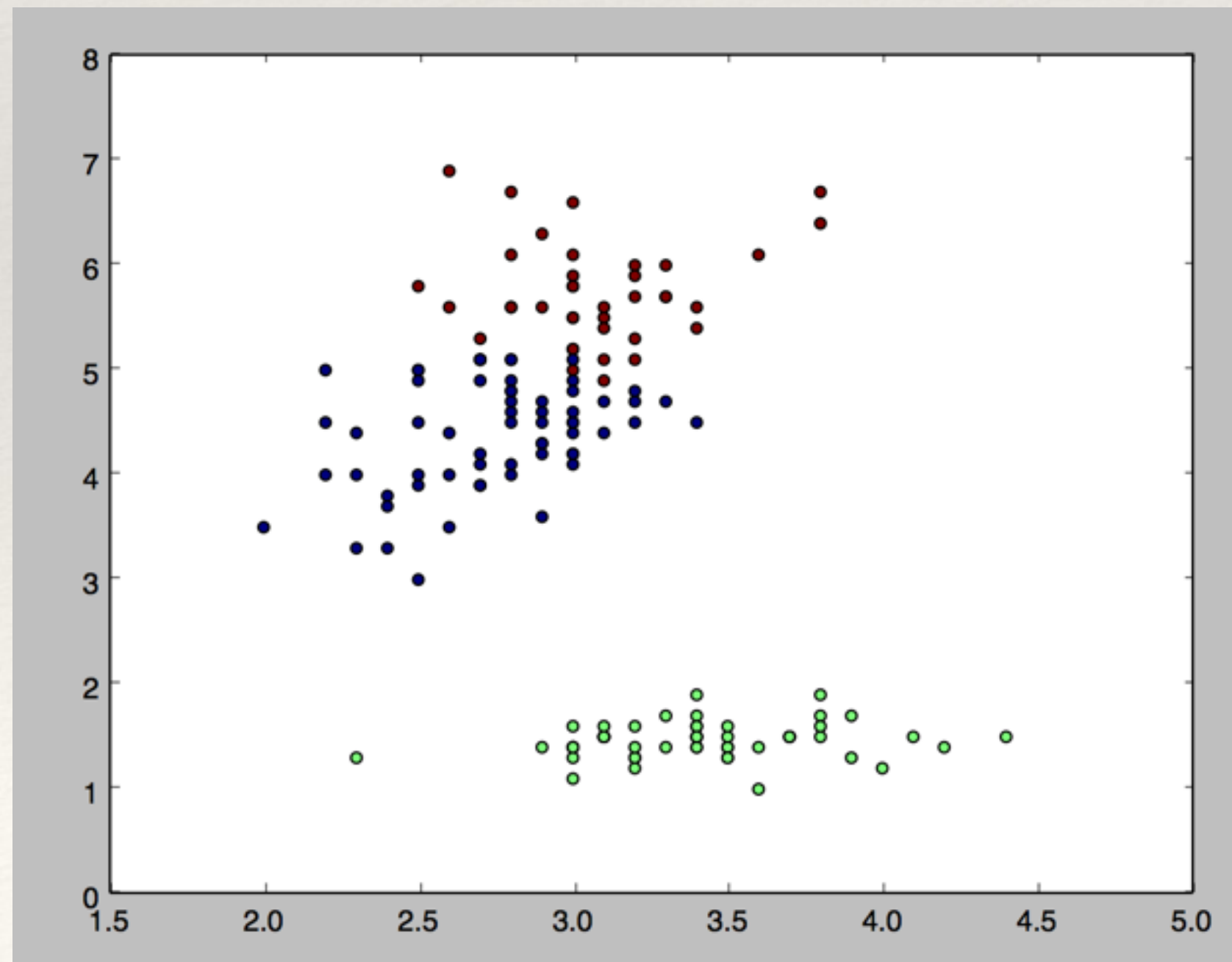
Warm up

Nearest Neighbour

Compare Sentence
Jaccard's distance

```
1  """
2  A Pure Python Nearest Neighbour
3  Algorithm
4
5  Mar-07-2016
6  Gopi Subramanian
7  """
8
9  # Small input
10 data = 'A dog chased a cat.\n \
11         The cat ran away from the dog.\n \
12         Cat are from tiger family.\n \
13         Dog is loyal. \n \
14         Cat and dog is pet.'
15
16 # Process to get words
17 sentences = data.split('\n')
18 words = [ sent.strip().lower().split(' ') for sent in sentences ]
19
20 def jacc_similarity(l1,l2):
21     s1 = set(l1)
22     s2 = set(l2)
23     return len(s1.intersection(s2)) / (1.0 * len(s1.union(s2)))
24
25 # Get all possible pairs and their
26 # Similarity
27 for i in range(len(words)):
28     for j in range(i+1, len(words)):
29         print i,j,jacc_similarity(words[i], words[j])
30
31 # Find Nearset Neighbours
32 test = "some cat pet".split(' ')
33 sim = [ jacc_similarity(test,word) for word in words]
34
35 max_sim = max(sim)
36
37 # Cluster assignment
38 for i in range(len(sim)):
39     if sim[i] == max_sim:
40         print words[i], sim[i]
```


K-Means Clustering



```
1  """
2  A simple python script to introduce
3  numerical libraries, |
4
5      * scikit-learn,
6      * numpy
7      * matplotlib
8
9  1. Load iris dataset
10 2. Perform KMeans
11 3. Plot output
12
13 Mar-07-2016
14 Gopi Subramanian
15 """
16
17 # Load Libraries
18 import numpy as np
19 import matplotlib.pyplot as plt
20 from sklearn.cluster import KMeans
21 from sklearn import datasets
22
23 # Let us use Iris dataset
24 iris = datasets.load_iris()
25 X = iris.data
26 Y = iris.target
27
28 # Fit KMeans
29 model = KMeans(n_clusters = 3)
30 model.fit(X)
31
32 # Plot output
33 fig = plt.figure(1)
34 plt.clf()
35 plt.scatter(X[:,1],X[:,2],c = model.labels_)
36 plt.show()
37
```


Decision tree classifier

A Decision Tree Classifier

Accuracy Metrics to evaluate the classifier

```
1  """
2  A simple python script to introduce
3  numerical libraries,
4
5      * scikit-learn,
6      * numpy
7      * matplotlib
8
9  1. Load iris dataset
10 2. Perform Classification
11
12 Mar-07-2016
13 Gopi Subramanian
14 """
15
16 # Load Libraries
17 import numpy as np
18 from sklearn import datasets
19 from sklearn.tree import DecisionTreeClassifier
20 from sklearn.metrics import accuracy_score
21
22
23 # Let us use Iris dataset
24 iris = datasets.load_iris()
25 x = iris.data
26 y = iris.target
27
28 # Build a classifier
29 estimator = DecisionTreeClassifier()
30 estimator.fit(x,y)
31 predicted_y = estimator.predict(x)
32
33 # Find model accuracy
34 print "Model accuracy = %.2f"%(accuracy_score(y,predicted_y) * 100) + "%\n"
35
```

Libraries

Library	Features
Numpy	N Dimensional Array Objects
Matplotlib	Graph Library
Scikit-Learn	Machine Learning Algorithms
Pandas	R Like Data Frame for Data Analysis
NetworkX	Graph Mining
NLTK	Text Processing
Simulation	SimPy
Disco	Light weight Distributed Computing (Map Reduce)
Theano	Neural Networks / Deep Learning

Numpy

- ❖ Numerical Python
- ❖ Enriches Python Data Structures with multi-dim arrays and matrices
- ❖ Scipy` - A companion library with function minimisation and fourier transforms functions.
- ❖ Python when combined with Numpy, Scipy and Matplotlib is a good alternative to MATLAB or OCTAVE.
- ❖ <http://www.python-course.eu/numpy.php>, A good tutorial on numpy

Function Minimisation

Scipy' Minimize scalar function

An exponential function

$$-\exp^{-(x-.7)^2}$$

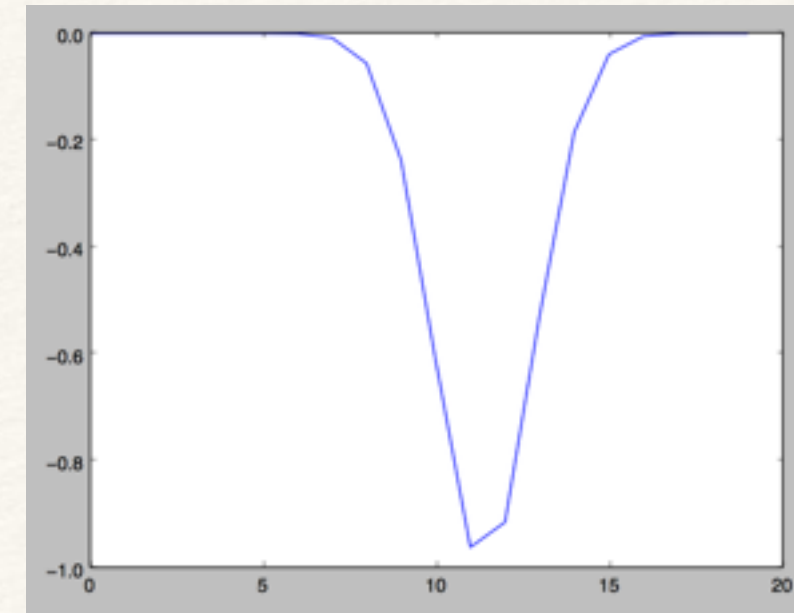
```
0.6999999999784
```

```
fun: -1.0
```

```
nfev: 10
```

```
nit: 9
```

```
x: 0.69999999997839409
```

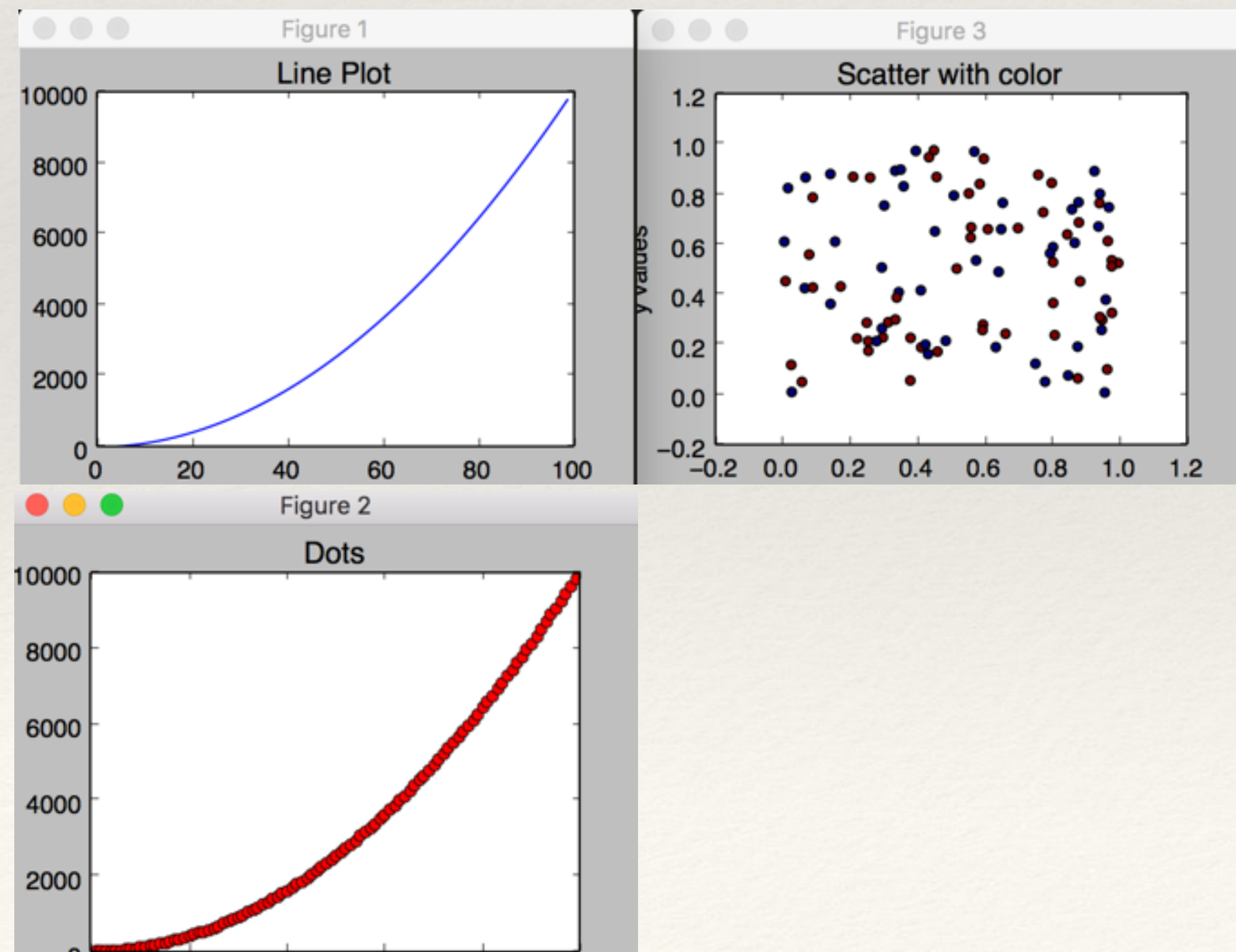


```
1  """
2  Function minimiztion using
3  scipy
4
5  Mar-07-2016
6  Gopi Subramanian
7  """
8
9  from scipy import optimize
10 import numpy as np
11 import matplotlib.pyplot as plt
12
13 # Define a function
14 def afunction(x):
15     return -np.exp(-(x-.7)**2)
16
17 # Generate some input
18 vals = np.arange(-5,5,0.5)
19 fvals =[afunction(x) for x in vals]
20
21 # Plot the function
22 plt.figure(2)
23 plt.plot(fvals)
24
25 plt.show()
26
27 print optimize.brent(afunction)
28
29 #Alternatively
30 print optimize.minimize_scalar(afunction)
```

Matplotlib

- ❖ 2D Plotting Library
- ❖ Plots, Histograms, Scatterplots, bar graphs, error charts
- ❖ <http://matplotlib.org>
- ❖ Pyplot' collection of command style functions
- ❖ Pylab' combines numpy' and pyplot'

Matplotlib



```
1  """
2  Introducing Matplotlib
3  Mar-07-2016
4  Gopi Subramanian
5  """
6
7  import numpy as np
8  import matplotlib.pyplot as plt
9
10 plt.close('all')
11 # Sample x y data for line and simple dot plots
12 x = np.arange(1,100,dtype=float)
13 y = np.array([np.power(xx,2) for xx in x])
14
15 # Line Plot
16 plt.figure(1)
17 plt.plot(x,y)
18 plt.xlabel('x values')
19 plt.ylabel('y values')
20 plt.title('Line Plot')
21
22 # Dot plot
23 plt.figure(2)
24 plt.plot(x,y,'or')
25 plt.xlabel('x values')
26 plt.ylabel('y values')
27 plt.title('Dots')
28
29 # Sample x,y data for scatter plot
30 x = np.random.uniform(size=100)
31 y = np.random.uniform(size=100)
32 labels = np.random.randint(2,size=100)
33
34 # Scatter plot
35 plt.figure(3)
36 plt.scatter(x,y,c=labels)
37 plt.xlabel('x values')
38 plt.ylabel('y values')
39 plt.title('Scatter with color')
40
41 plt.show()
```

Scikit-Learn

- ❖ Python Machine Learning Library
- ❖ <http://scikit-learn.org/stable/>
- ❖ Leverages Numpy, Scipy and matplotlib
- ❖ Range of Machine Learning Algorithms.
- ❖ Easily extendable. Very Clean API

Scikit-Learn API

- ❖ Every estimator object should have the following functions
 - ❖ `fit (X, Y)` for classifiers, `fit(X)` for clustering
 - ❖ `predict(X)`
 - ❖ `fit_predict(X,Y)`


```
model = LinearRegression(normalize=True,fit_intercept=True)
model.fit(x,y)
y_p = model.predict(x)
```

```
model = Ridge(normalize=True,alpha=0.015)
model.fit(x,y)
y_p = model.predict(x)
```


Train / Test Split

- ❖ Split data into train and test
- ❖ Stratified split of data by Y variable

```
1  """
2  Split data as train and test
3
4  Mar-07-2016
5  Gopi Subramanian
6  """
7
8  # Load Libraries
9  import numpy as np
10 from sklearn import datasets
11 from sklearn.tree import DecisionTreeClassifier
12 from sklearn.metrics import accuracy_score
13 from sklearn.crossvalidation import train_test_split
14
15
16 # Let us use Iris dataset
17 iris = datasets.load_iris()
18 x = iris.data
19 y = iris.target
20
21 # Stack x and y together
22 input_dataset = np.colstack([x,y])
23 train,test = train_test_split(input_dataset,test_size = 0.2)
24
25 # Perform Stratified Split data by class variable
26 stratified_split = StratifiedShuffleSplit(y,test_size = 0.2, n_iter = 1)
27
28 for train_indx,test_indx in stratified_split:
29     train = input_dataset[train_indx]
30     test = input_dataset[test_indx]
31
```

More.....

- ❖ Cross Validation
- ❖ Parameter Search - Grid Search and Random Search
- ❖ Data Sets - Iris, Boston Housing, Create Data Sets for classification and clustering
- ❖ Evaluation Metrics - Accuracy, Root Mean Square....
- ❖ Stochastic Gradient Descent, Random Projections....

Text Mining Using Python

- ❖ NLTK - Natural Language Toolkit
- ❖ <http://www.nltk.org>
- ❖ NLTK plugins for Scikit-Learn

RAKE - Rapid Keyphrase extraction

- RAKE is an extremely efficient keyword extraction algorithm and operates on individual documents.
 - Its language and domain independent.
- Given a document, stop word list and a list of phrase delimiters, RAKE extracts candidate phrases.
- The next step is to find out the frequency of the individual words in these phrases.
 - Frequency is the count of occurrence of the word.
- Find Degree of a word is sum of length of all the phrases where the word occurs.
- Finally scoring for each word is done by, $\text{degree}(\text{word}) / \text{freq}(\text{word})$
- Phrase scores are calculated by sum of individual word scores in that phrase.
 - Phrase with very large scores are considered to be keyphrases for the document.
- <https://github.com/subramgo/RAKE>

Text Search Engine

- ❖ <https://github.com/subramgo/Vritti>
- ❖ Word Associations
- ❖ Document Clustering - Non Negative Matrix Factorisation
- ❖ Collocations
- ❖ Uses PyLucene - Document Index

Sentiment Mining

<http://sentiment.christopherpotts.net/index.html#data>

A Nice Tutorial and Several Python implementations by Christopher Potts from Stanford

Simulation Using Python

- ❖ SimPy
- ❖ <http://simpy.readthedocs.org/en/latest/index.html>
- ❖ Let us see a simple example
- ❖ Say a new store is opened, we want to simulate how many checkout counters should be open at any point in time. Since we don't have historical data we simulate.


```

1  """
2  Simulating a Super Market Aisle
3
4  Let us say we want to decide how many ailes should be
5  functional in a new super market. Since it is a
6  new shop, we dont have historic information. We simulate
7  and check out.
8
9
10 Mar-09-2016
11 Gopi Subramanian
12 """
13
14
15 import simpy
16 import random
17 import time
18
19
20 time_spent_at_counter = []
21
22 class Counter(object):
23     def __init__(self,env,no_counters,scan_time):
24         self.env = env
25         self.counters = simpy.Resource(env, no_counters)
26         self.item_scan_time = scan_time
27
28     def checkout(self, customer, no_items):
29         yield self.env.timeout(no_items * self.item_scan_time)
30
31 def customer(env, name, counter, no_items):
32     #total_customers+=1
33     with counter.counters.request() as request:
34         yield request
35         enter = env.now
36         yield env.process(counter.checkout(name, no_items))
37         leaves = env.now
38         #print '%s Enters at %.2f leaves counter at %.2f'%(name,enter,leaves)
39         elapsed = leaves - enter
40         time_spent_at_counter.append(elapsed)
41

```

```

43 def setup(env,no_counters, scan_time):
44     counters = Counter(env, no_counters, scan_time)
45     # Initial customers
46     for i in range(10):
47         no_items = random.randint(5,35)
48         env.process(customer(env, 'Customer %d' % i, counters, no_items))
49     hour = 0
50     while True:
51         # Every hour
52         yield env.timeout(random.randint(48,62))
53         hour = hour + 1
54         no_customers = random.randint(30,200)
55         i += 1
56         start = i
57         for i in range(start,start+no_customers):
58             no_items = random.randint(5,35)
59             env.process(customer(env, 'Customer %d' % i, counters,no_items))
60         print 'Hour %d, Number of Customers %d, %.3f Average minutes spent by a customer at counter ' \
61               % (hour, no_customers, (sum(time_spent_at_counter) / (1.0*len(time_spent_at_counter))))

```

```

64 No_counters = 2
65 # Time it takes to scan an item
66 item_scan_time = 0.3
67
68 env = simpy.Environment()
69 env.process(setup(env, No_counters, item_scan_time))
70
71 # Simulation in minutes
72 # 8 hours
73 env.run(until=480)
74 print '%.3f Average minutes spent by a customer at counter ' \
75       % (sum(time_spent_at_counter) / (1.0*len(time_spent_at_counter)))
76

```

```

Hour 1, Number of Customers 190, 6.960 Average minutes spent by a customer at counter
Hour 2, Number of Customers 50, 7.471 Average minutes spent by a customer at counter
Hour 3, Number of Customers 67, 7.425 Average minutes spent by a customer at counter
Hour 4, Number of Customers 117, 6.750 Average minutes spent by a customer at counter
Hour 5, Number of Customers 133, 6.829 Average minutes spent by a customer at counter
Hour 6, Number of Customers 110, 6.910 Average minutes spent by a customer at counter
Hour 7, Number of Customers 98, 6.519 Average minutes spent by a customer at counter
Hour 8, Number of Customers 63, 6.426 Average minutes spent by a customer at counter
6.391 Average minutes spent by a customer at counter

```

Thank You....

- ❖ Code and Presentation
- ❖ <https://github.com/subramgo/Py4DS>