


```
[In 190]: #writing lemmatize function
def lemmatizer_function(our_list):
    main_list=[]
    for i in our_list:
        main_list.append(lemmatizer.lemmatize(i, pos='n'))
    return main_list

#performing lemmatization on our word list
word_string=lemmatizer_function(word_string)
t_pos1=lemmatizer_function(t_pos1)
t_pos2=lemmatizer_function(t_pos2)
t_pos3=lemmatizer_function(t_pos3)
t_pos4=lemmatizer_function(t_pos4)
t_pos5=lemmatizer_function(t_pos5)
t_pos6=lemmatizer_function(t_pos6)

In [191]: #marking data whether positive or negative
pos_docs=[(t_pos1,'pos'),(t_pos2,'pos'),(t_pos3,'pos'),(t_pos4,'pos'),(t_pos5,'pos'),(t_pos6,'pos')]
neg_docs=[(word_string[0:40],'neg'),(word_string[40:60],'neg'),(word_string[60:70],'neg'),
          (word_string[70:90],'neg'),(word_string[90:138],'neg'),(word_string[138:158],'neg')]

In [192]: # Create a training and test set
train_pos_docs = pos_docs[0:4]
test_pos_docs = pos_docs[4:6]
train_neg_docs = neg_docs[0:4]
test_neg_docs = neg_docs[4:6]

# Merge the training sets and the test sets
training_docs = train_pos_docs+train_neg_docs
testing_docs = test_pos_docs+test_neg_docs

In [193]: # building object of Sentiment Analyzer
sentim_analyzer = SentimentAnalyzer()

In [194]: # Mark up the negative sections of the doc
all_words_neg = sentim_analyzer.all_words(mark_negation(doc) for doc in training_docs)

In [195]: # Return all words that occur at least 2 times
unigram_feats = sentim_analyzer.unigram_word_feats(all_words_neg, min_freq=2)

#printing number words which have minimum frequency of 2.
print(len(unigram_feats))
3

In [196]: #adding feature extractor to our model
sentim_analyzer.add_feat_extractor(extract_unigram_feats, unigrams=unigram_feats)

In [197]: # Apply the feature extractor to the docs
training_set = sentim_analyzer.apply_features(training_docs)
test_set = sentim_analyzer.apply_features(testing_docs)

In [198]: # using NaiveBayesClassifier modeling to train on our data
trainer = NaiveBayesClassifier.train
classifier = sentim_analyzer.train(trainer, training_set)

Training classifier

In [199]: # Print out the evaluation metrics
sentim_analyzer.evaluate(test_set)

Evaluating NaiveBayesClassifier results...

Out[199]: {'Accuracy': 0.5,
'Precision [neg]': None,
'Recall [neg]': 0.0,
'F-measure [neg]': None,
'Precision [pos]': 0.5,
'Recall [pos]': 1.0,
'F-measure [pos]': 0.6666666666666666}
```

Here we can notice the accuracy of our model is just 50% because we didn't gave much data to train our model

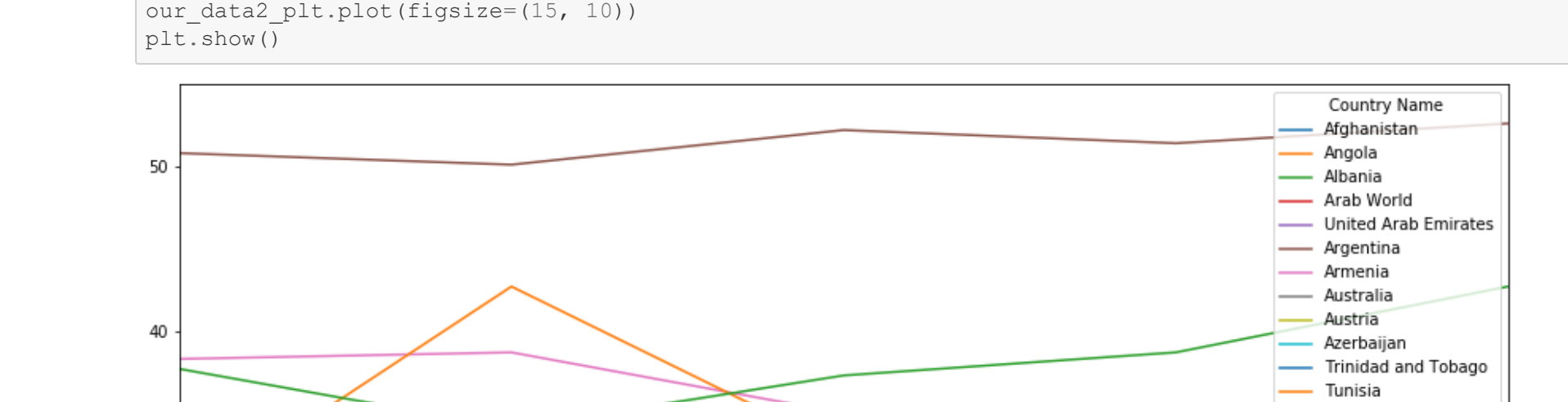
```
In [100]: # Print out most informative features
classifier.show_most_informative_features(15)

Most Informative Features
contains(inflation) = False      pos : neg      = 3.0 : 1.0
contains(recession) = False     pos : neg      = 3.0 : 1.0
contains(time) = False          pos : neg      = 1.8 : 1.0

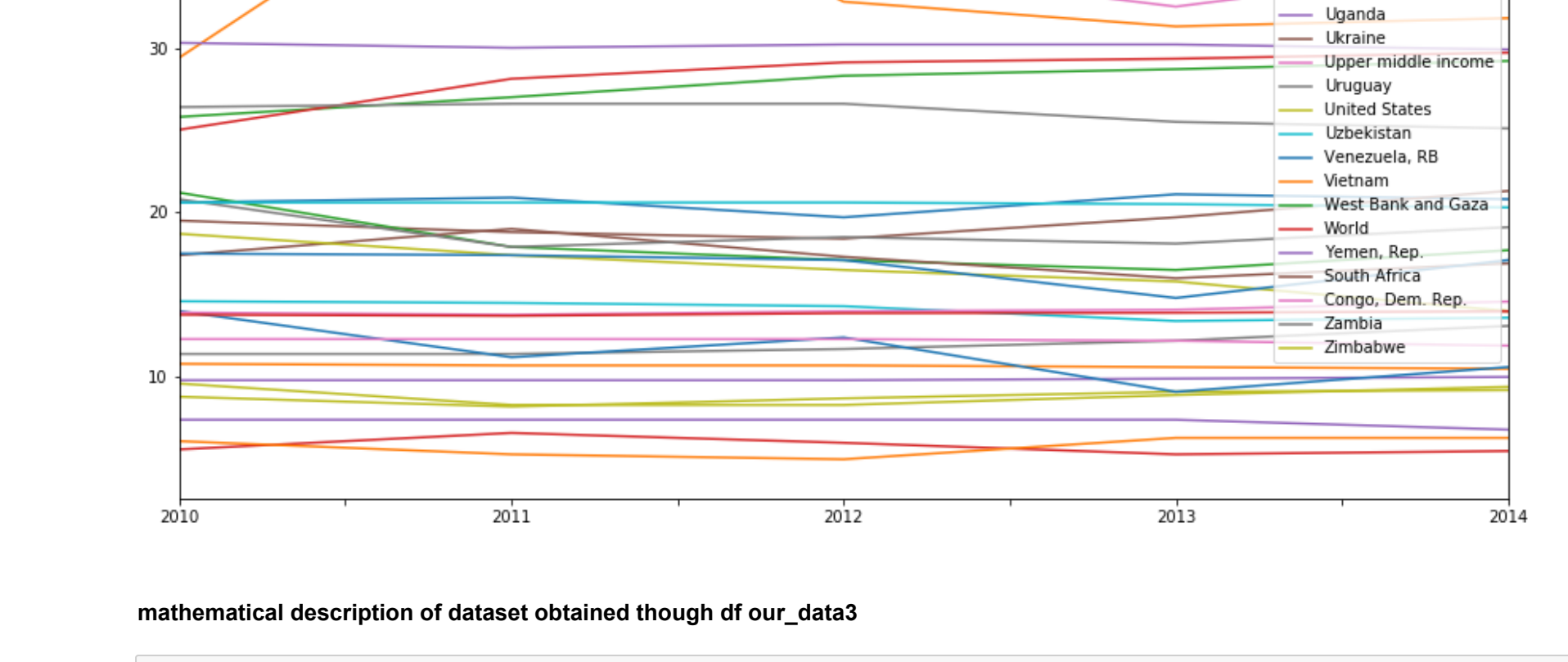
In [101]: #making prediction through our model
for doc in test_set:
    print(classifier.classify(doc[0]))

pos
pos
pos
pos
```

Unemployment range depicted through a boxplot



Unemployment range depicted through a line graph



mathematical description of dataset obtained though off our_data3

```
In [104]: # describing statistics parameters (like mean, mode, median etc.) of our 3rd DataFrame
our_data3.describe()

Out[104]:
```

	Estimated Unemployment Rate (%)	Estimated Employed	Estimated Labour Participation Rate (%)
count	740.000000	7.400000e+02	740.000000
mean	11.787946	7.204460e+06	42.630122
std	10.721298	8.067989e+06	8.111094
min	0.000000	4.942004e+04	13.330000
25%	4.657500	1.190404e+06	36.062500
50%	8.350000	4.744178e+06	41.160000
75%	15.887500	1.127549e+07	45.505000
max	76.740000	4.577751e+07	72.570000

Heatmap of off our_data3 and we can see there is no relation between the columns

```
In [105]: #using seaborn library to plot heat map
import matplotlib inline
plt.figure(figsize = (12,8))
ax = sns.heatmap(our_data3.corr(), cmap = "YlGnBu", linewidths = 0.30, annot = True)
sns.set(font_size=10)
ax.set_yticklabels(['Estimated Unemployment Rate (%)', 'Estimated Employed', 'Estimated Labour Participation Rate (%)'])
plt.xticks(rotation=degrees)
plt.show()
```



Dropout rates among Men and Women depicted through a diagonal graph

```
In [106]: #dividing our data to test and train our model
x=our_data1.drop(['Year','Month','Date','White','Black','Asian','Hispanic','Men','Women'],axis=1).dropna()
y=our_data1['Dropout_Rate']
x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.25,random_state=42)

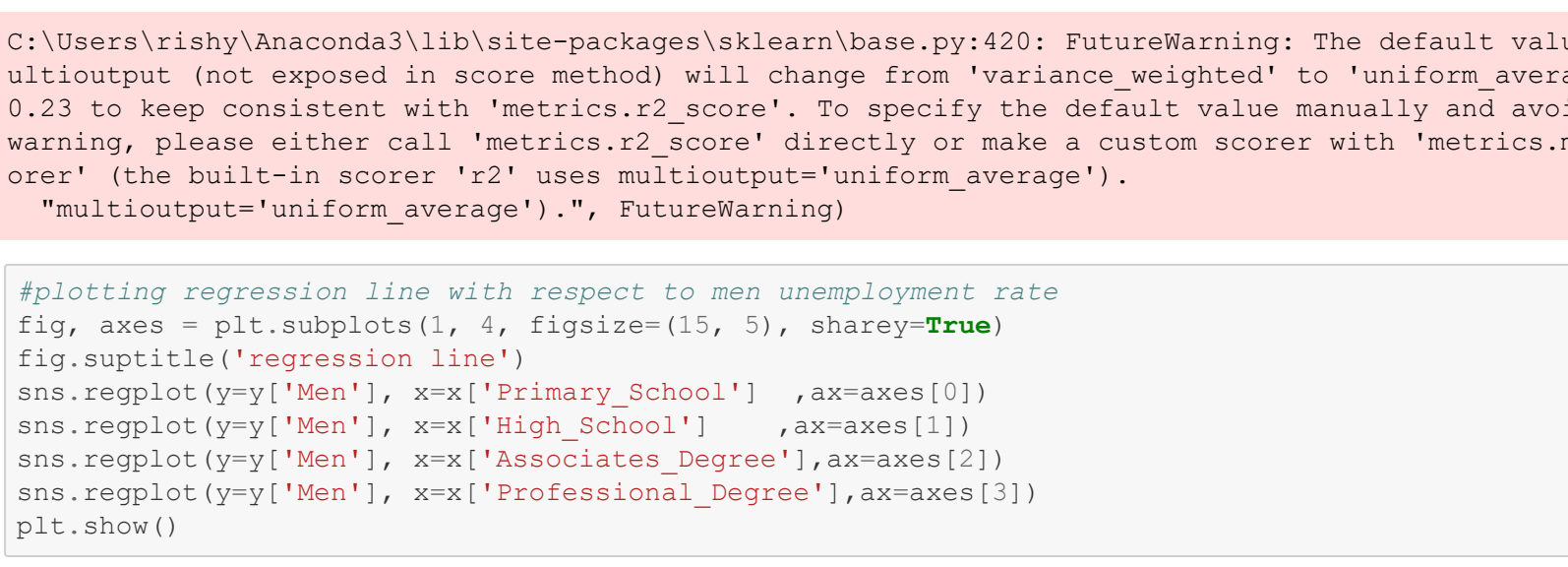
In [107]: #building pipeline which first scale our data then perform linear regression on it
pipe = Pipeline([('scaler', StandardScaler()), ('linearregression', LinearRegression())])
pipe.fit(x_train,y_train)
y_pred=pipe.predict(x_test)

In [108]: #Model score, mean squared error & root mean squared error
print('Model score = ',pipe.score(x_test,y_test))
print('mean_sqd_error = ',mean_squared_error(y_test,y_pred))
print('root_mean_squared_error of is = ',np.sqrt(mean_squared_error(y_test,y_pred)))

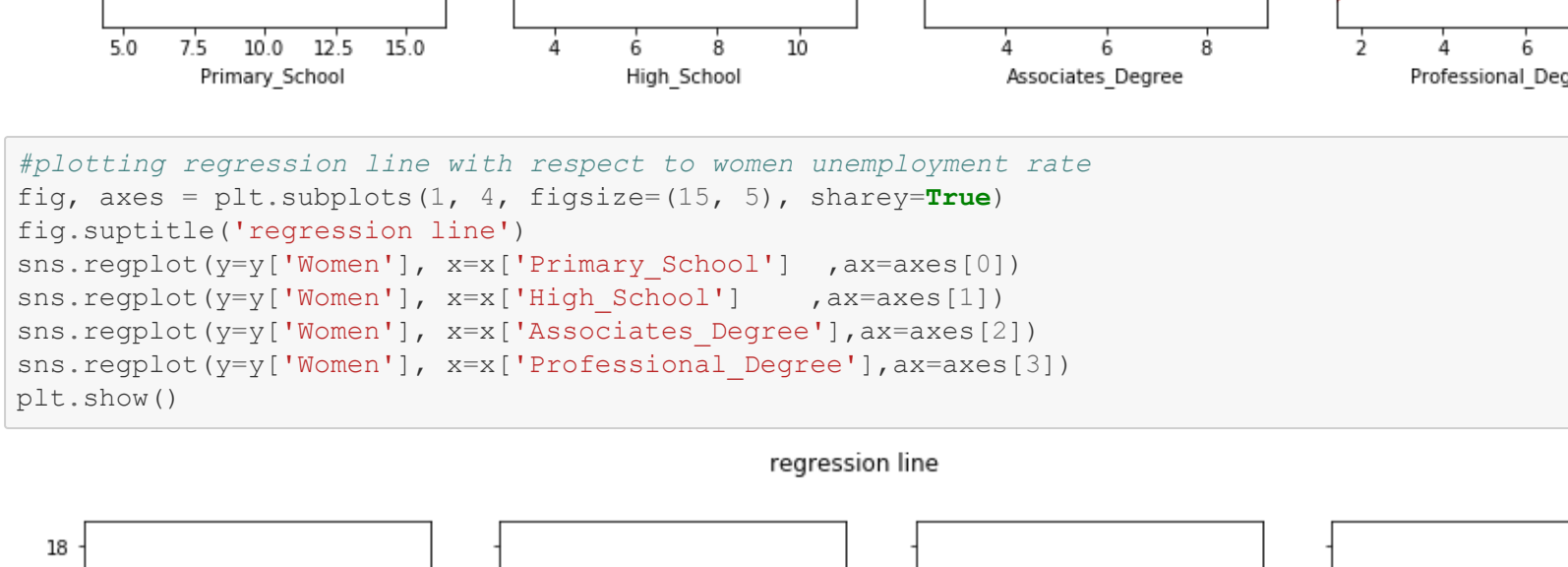
Model score = 0.9860170750978912
mean_sqd_error is = 0.054932518363035
root_mean_squared_error of is = 0.23437843723054036
```

C:\Users\irishy\Anaconda3\lib\site-packages\sklearn\base.py:420: FutureWarning: The default value of n_jobs is 1 which will be changed to -1 in version 0.23 to keep consistent with 'metrics.r2_score'. To specify the default value manually and avoid the warning, please either call 'metrics.r2_score' directly or make a custom scorer with 'metrics.make_scorer' (the built-in scorer 'r2' uses multioutput='uniform_average').

```
In [109]: #plotting regression line with respect to men unemployment rate
fig, axes = plt.subplots(4, figsize=(15, 5), sharey=True)
fig.suptitle('regression line')
sns.regplot(y=y['Men'], x=x['Primary_School'], ax=axes[0])
sns.regplot(y=y['Men'], x=x['High_School'], ax=axes[1])
sns.regplot(y=y['Men'], x=x['Associates_Degree'], ax=axes[2])
sns.regplot(y=y['Men'], x=x['Professional_Degree'], ax=axes[3])
plt.show()
```



```
In [110]: #plotting regression line with respect to women unemployment rate
fig, axes = plt.subplots(4, figsize=(15, 5), sharey=True)
fig.suptitle('regression line')
sns.regplot(y=y['Women'], x=x['Primary_School'], ax=axes[0])
sns.regplot(y=y['Women'], x=x['High_School'], ax=axes[1])
sns.regplot(y=y['Women'], x=x['Associates_Degree'], ax=axes[2])
sns.regplot(y=y['Women'], x=x['Professional_Degree'], ax=axes[3])
plt.show()
```



CONCLUSION

We have extracted data from multiple websites to compare and understand the effects of education on unemployment numbers around the world. It was clear through our dataset that unemployment in developed countries was much less as compared to unemployment in underdeveloped countries. This theory was also confirmed through the double bar graphs that we plotted for urban and rural India unemployment. A similar trend was seen here, with urban unemployment numbers being lower than that in rural India. We were also able to identify a correlation between dropout rates and unemployment in men and women through sklearn library. There was a clear correlation between the education dropout rates and unemployment. This stayed true without significant difference for both genders. From our data it is clear that education/skills plays an important role in decreasing unemployment.

Here I took help from documentation of these major libraries:

1. pandas for building data frames https://pandas.pydata.org/docs/getting_started/overview.html
2. numpy to perform maths operations <https://numpy.org/doc/>
3. matplotlib to plot graphs from our databases <https://matplotlib.org/stable/users/index.html>
4. seaborn to plot heat map and pair plot graph <https://seaborn.pydata.org/generated/seaborn.heatmap.html>
5. beautifulsoup to perform web scraping <https://beautiful-soup-4.readthedocs.io/en/latest/>
6. wordcloud to draw unemployment word cloud https://amueller.github.io/word_cloud/
7. sklearn to build regression model https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LinearRegression.html
8. nltk to perform sentiment analysis <https://www.nltk.org/api/nltk.sentiment.html>

If formatting is not proper then please check out pdf version of it.