# Foundation of Data Science: Project Report

Rana Al-rubaye

# Data balance

- **Up-sample:** is the process of randomly duplicating observations from the minority class in order to reinforce its signal.
- **Down-Sample:** randomly removing observations from the majority class to prevent its signal from dominating the learning algorithm.
- **Note:** Up and Down sampling give same accuracy while down sampling faster since we have less samples. We offer both options in our solution but we set default to down-sample.

**Example of balancing the dataset using Down-sample**

| 0 | 6787 |
|---|------|
| 1 | 365  |

**We random remove samples from majority that match the number of samples in minority samples.**

| 1 | 365 |
|---|-----|
| 0 | 365 |

```python
df_majority = data[data.fraudulent == 0]
df_minority = data[data.fraudulent == 1]
# Downsample majority class
df_majority_downsampled = resample(df_majority,
            replace=False,  # sample without replacement
            n_samples=df_minority.shape[0],  # to match minority class
            random_state=123)  # reproducible results
# Combine minority class with downsampled majority class
df_sampled = pd.concat([df_majority_downsampled, df_minority])
```

# Preprocessing- TF-IDF ( top 45 words)

- We have text data in columns (title, location, description, and requirements), we need to make all dataset as numeric data.
- Apply TF-IDF on the text data and pick top 45 words. I tried different top from top10 - top >200 when I go over 45 noise increase.

# Preprocessing (1-2)

- Split the columns that has text from columns that has number

```python
# Take only columns that has number
X_clean = data.drop(self._pp_colums, axis=1)
self._shiftCols = X_clean.shape[1]
# Take only columns that has text
X_pre = data[self._pp_colums]
textData = []
for index, row in X_pre.iterrows():
    text = "{} {} {} {}".format(row['title'],
                                row['location'],
                                row['description'],
                                row['requirements'])
    textData.append(text)
documents = pd.DataFrame(textData, columns=['headline_text'])
```

# Preprocessing (2-2)

- Apply TF-IDF on columns that has text. Then combine generated numeric columns with existing numeric columns in dataset

```python
# create the TD-IDF transform
if self.vectorizer == None:
  self.vectorizer = TfidfVectorizer(stop_words='english', norm='l2', use_idf=False, smooth_idf=False)
  # apply TF-IDF on training
  vectors = self.vectorizer.fit_transform(documents["headline_text"])
else:
  # apply TF-IDF on testing
  vectors = self.vectorizer.transform(documents["headline_text"])
# get the numberic data from TD-IDF as 2D array
data_pro = pd.DataFrame(vectors.toarray(), columns=self.vectorizer.get_feature_names())
# create new N columns in datase, where _numberOfTopics is number of  words that we want to pick
default set to 45
for col in range(data_pro.shape[1]):
  X_clean["Topic_{}".format(col)] = np.nan
  if col > self._numberOfTopics:
     break
# merge the data of N words of TD-IDF into dataset
for row in range(data_pro.shape[0]):
  for col in range(data_pro.shape[1]):
     X_clean.iloc[row, self._shiftCols + col] =  data_pro.iloc[row, col]
     if col > self._numberOfTopics:
        Break
```

# Normalize the data

- Put all data in same scale using min/max scale

```
# normalize the data
min_max_scaler = preprocessing.MinMaxScaler()
X_norm = min_max_scaler.fit_transform(X_train)
```

# Model Training

- I use stochastic gradient descent (SGD) learning. Also I notice similar results using different algorithms such as decision tree.
- SGD is the a good way to optimization in Machine Learning

```python
def fit(self, X, y):
    print( "==Training the model ( takes up to 5 minutes )=======")
    [X_train_balance, y_balance] = textNLP.balanceData(X,y)   # balance data
    X_train = textNLP.pre_pro_Cols(X_train_balance) # pre-processing
    X_norm = min_max_scaler.fit_transform(X_train)# normalize data
    # Create stochastic gradient descent
    self.clf = SGDClassifier()
    self.clf.fit(X_norm, y_balance)
```

# Feature selection

- Base on filter-base feature selection results:
  - I notice that starting Top_45 feature column and next features that generated by TF-IDF are not contributing to the model accuracy rather than they lead to noise.
  - I pick only top 45 features and ignore rest.

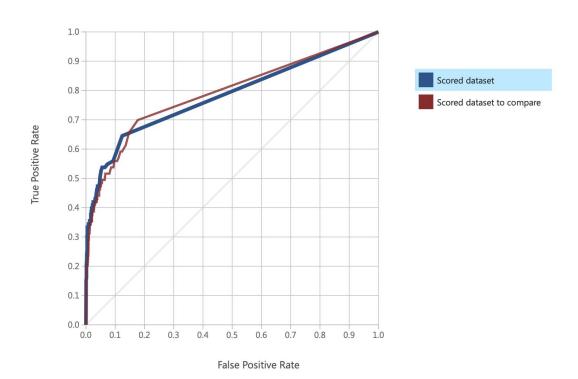Rana Project ❯ Filter Based Feature Selection ❯ Features

| rows | columns |
|------|---------|
| 1 | 9 |

| | fraudulent | has_company_logo | has_questions | telecommuting | Topic_0 | Topic_3 | Topic_1 | Topic_2 | Topic_4 |
|---|---|---|---|---|---|---|---|---|---|
| view as | 1 | 0.256308 | 0.099393 | 0.043413 | 0.042841 | 0.029965 | 0.022064 | 0.00722 | 0.005666 |

# Model Selection and Tuning

Not noticeable impact

# Prediction

```python
def predict(self, X):
    print("====  Testing the model =====")
    # remember to apply the same preprocessing in fit() on test data before making predictions
    X_test = textNLP.pre_pro_Cols(X)
    X_norm = min_max_scaler.fit_transform(X_test)
    return self.clf.predict(X_norm)
```

# Model Evaluation

- To evaluate model, I split the dataset into 80% training and 20% testing, I find these results.

```
=============== Testing the model ====================
(Step 2 of 2) Pre-proces coloums that has text (token,  remove stop
               precision    recall  f1-score   support

           0        0.98      0.86      0.92      1685
           1        0.24      0.75      0.37       103

    accuracy                            0.85      1788
   macro avg        0.61      0.80      0.64      1788
weighted avg        0.94      0.85      0.88      1788
```

# Questions?