

Optimizing Search Engine Relevance

By
Samaneh Karimi

LinkedIn: <https://www.linkedin.com/in/samane-karimi-52178459/>

Outline

❖ Summary

- Preprocessing
- Techniques and models used
- Evaluation results

❖ Conclusion and final recommendations

❖ Python implementation and visual walk-through of the outputs

Summary of my approach

1. Analyzing raw data
2. Preprocessing
3. Model selection
4. Implementing the models and evaluation

Step1) Analyzing raw data

Step 1) Analyzing raw data

- **Goal: Knowing the characteristics of the data**

- Balanced or imbalanced?

Class	1 (Relevant)	0 (Non-relevant)
#Instance (ratio)	34987 (44%)	45060 (56%)

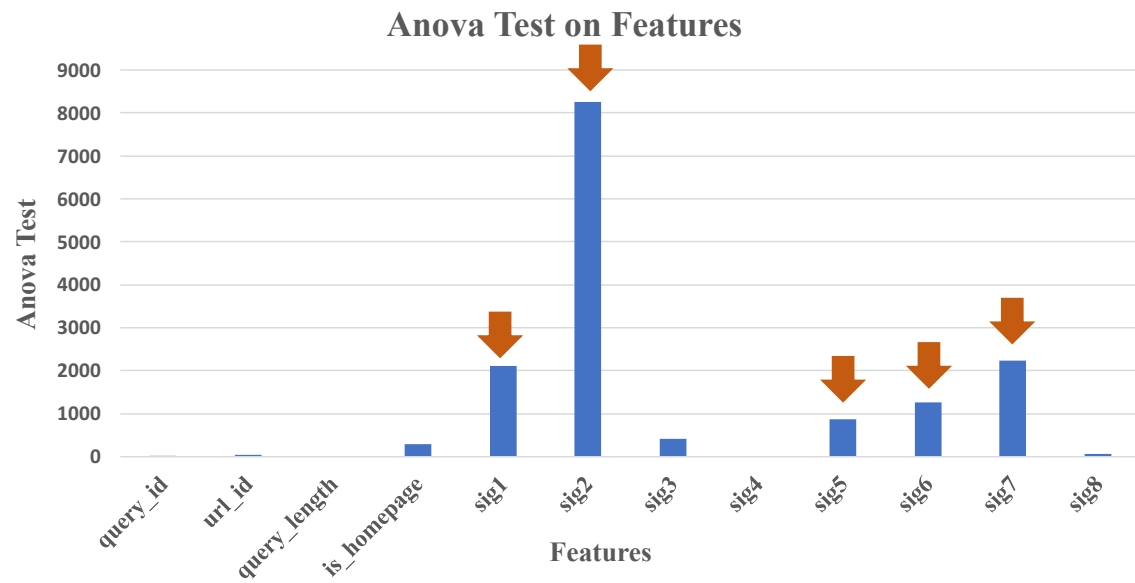
=> balanced (no need for using imbalanced data handling techniques)

- Analyzing feature importance:
 - Using ANOVA Test

Step 1) Analyzing the raw data (Cont.)

- Top 5 features based on ANOVA test:

1. sig2
2. sig7
3. sig1
4. sig6
5. sig5



Step 2) Preprocessing

- Handling missing values
- Duplicate removal
- Feature scaling
- Standardization
- Feature selection (dimensionality reduction)

Step 2) Preprocessing

- **Handling missing values**
 - No missing value was found in the dataset.
- **Duplicate removal**
 - The redundant data may appear both in the training and testing sets and cause inaccurate learning of the model.
 - No duplicate was found in the dataset.
- **Feature scaling**
 - The majority of machine learning algorithms perform much better when dealing with features that are on the same scale.
- **Standardization:**
 - Centering the feature columns at mean 0 with standard deviation 1 so that the feature columns have the same parameters as a standard normal distribution (zero mean and unit variance).

Step 2) Preprocessing (Cont.)

- **Feature Selection (Dimensionality Reduction)**
 - The process of reducing the number of input variables when developing a predictive model to both reduce the computational cost of modeling and to improve the performance of the model.
- **Two approaches tested in this project:**
 1. Statistics for filter-based feature selection methods:
 - Since we have numerical input and binary annotations in our dataset => ANOVA correlation coefficient is used.
 2. Principal Component Analysis (PCA)
 - finds directions of maximal variance of data

Step 3) Model selection

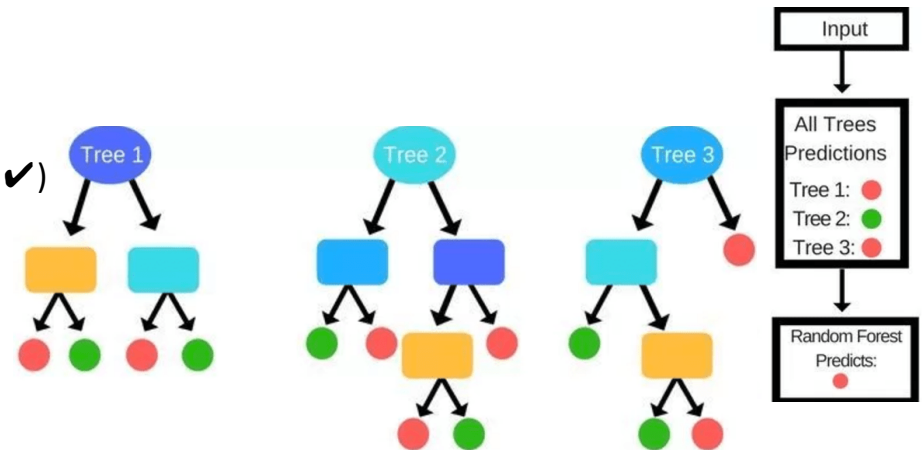
- Random Forest
- Support Vector Machines (SVM)
- Deep Learning (MLP)

Step 3) Model selection

Random Forest

- The **random forest** approach is an ensemble (bagging) method where **deep trees**, fitted on bootstrap samples, are combined to produce an output with lower variance. For each test data, the output class would be the mode of the classes of the individual trees.

- My reasons for choosing it:
 - It runs efficiently on large databases. (✓)
 - Since it takes the **advantage** of ensemble learning, reduces variance and thus helps us avoid overfitting.



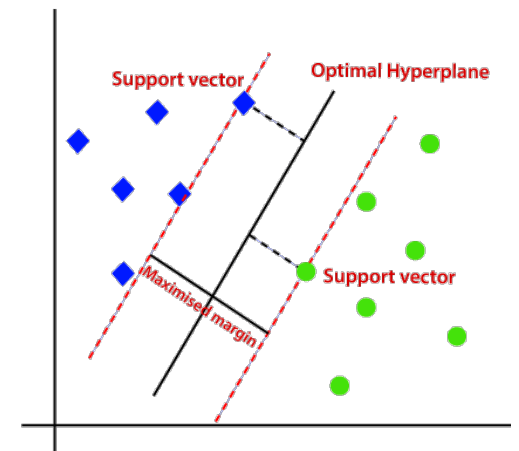
11

Image Ref: <https://syncdreview.com/2017/10/24/how-random-forest-algorithm-works-in-machine-learning/>

Step 3) Model selection

Support Vector Machines (SVM)

- SVM is a supervised machine learning algorithm that uses hyper planes to separate out different classes.
- Two types:
 - Linear
 - Non-linear (Kernel-based)
- My reasons for choosing it:
 - SVMs are suitable for numerical features (✓)
 - SVMs can handle non-linearly separable data (kernel trick)(✓)

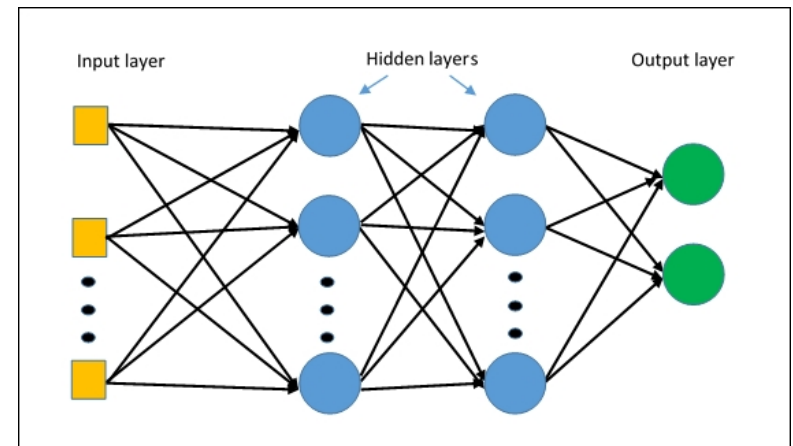


Step 3) Model selection

Deep Learning (MLP)

- Multilayer Perceptrons (MLPs) is a feedforward neural network that consists of multiple fully-connected layers.

- My reasons for choosing it is that MLPs are suitable for*:
 - Classification prediction problems (✓)
 - Tabular datasets (✓)



Ref*: <https://machinelearningmastery.com/when-to-use-mlp-cnn-and-rnn-neural-networks/>

Image Ref: <https://deepai.org/machine-learning-glossary-and-terms/multilayer-perceptron>

Step 4) Implementing the models and evaluation

- Evaluation Protocol
- Experiments

The Evaluation Protocol

➤ K-fold cross validation

- To avoid overfitting
- To avoid wasting data as the validation set

➤ Tuning hyper-parameters

- I implemented Grid Search using a hyper-parameter optimization tool (`sklearn.model_selection.GridSearchCV`)

➤ Evaluation measures

- Precision, Recall, F1 and Accuracy

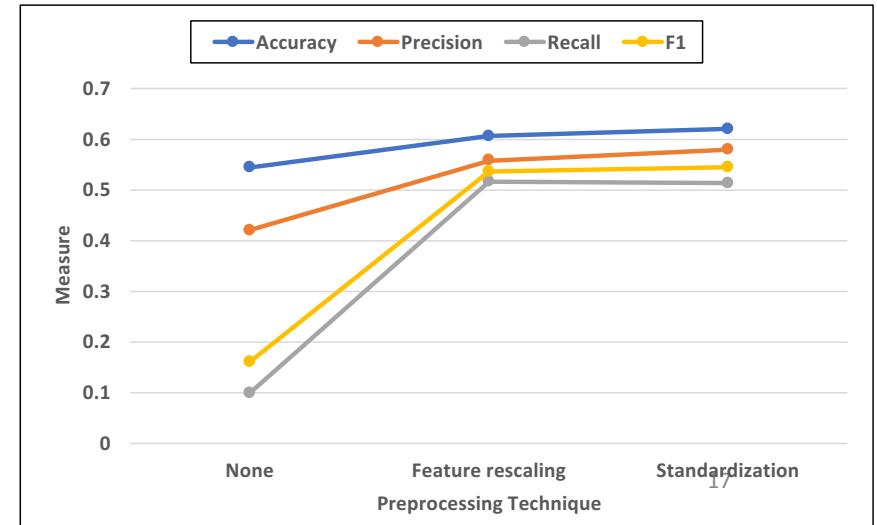
Experiments

- The impact of feature scaling and standardization
- The impact of dimension reduction techniques
- Hyper-parameters tuning for the classification methods
- Comparing the performance of the classification methods on the test set

The impact of feature scaling and standardization

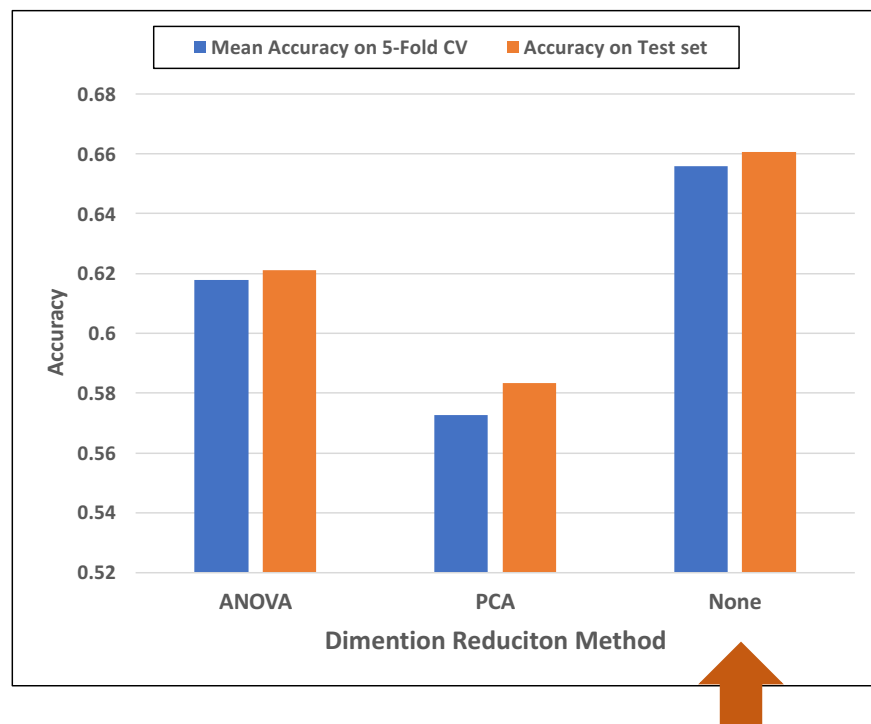
- Classification Method:
 - SVM
- Feature rescaling using min-max normalization
- Standardization: Centering the feature columns at mean 0 with standard deviation 1

- ❖ **Standardization has the highest impact on SVM's performance**
- ❖ **Standardization also decreased SVM's training time significantly.**



The impact of dimension reduction techniques

- Experimental setup:
 - Classification Method:
 - Random Forest
 - 5-fold Cross Validation
 - Evaluation Measure:
 - Training over 5-olds: Mean accuracy
 - Test set: Accuracy using the optimal classifier



- ❖ Results: The best performance on both training and testing is when no dimension reduction technique is used.

Hyper-parameters tuning Random Forest

- Experimental setup:
 - Grid Search over all combinations of the Random Forest parameters as follows
 - 'bootstrap': [True, False]
 - 'max_depth': [10, 20, 30, 40, 50, 60, 70, 80, 90, 100, None]
 - 'max_features': ['auto', 'sqrt']
 - 'min_samples_leaf': [1, 2, 4]
 - 'min_samples_split': [2, 5, 10]
 - 'n_estimators': [200, 400, 600, 800, 1000, 1200, 1400, 1600, 1800, 2000]
 - 5-fold Cross Validation
 - Tuned for both precision and recall measures
- Best parameters set found on development set:
 - For Precision :{'bootstrap': False, 'max_depth': 10, 'max_features': 'sqrt', 'min_samples_leaf': 1, 'min_samples_split': 5, 'n_estimators': 800}
 - For Recall: {'bootstrap': False, 'max_depth': 20, 'max_features': 'sqrt', 'min_samples_leaf': 1, 'min_samples_split': 5, 'n_estimators': 800}

Random Forest performance on the test set tuned based on precision and recall

- Experimental setup:
 - The Random Forest classifier is trained on the training set in two ways:
 - using optimal parameter values tuned for precision (Prec-based)
 - using optimal parameter values tuned for recall (Recall-based).
 - The training and test sets are obtained using Sklearn **train_test_split** module that makes random partitions for the two subsets.
 - `sklearn.metrics.classification_report` is used for reporting the model's performance on the test set.

Random Forest performance on the test set tuned based on precision and recall (Cont.)

- Performance on the test set using precision-based optimal parameter

Measure Class	Precision	Recall	F1-measure	Support
Class 0 (Non-Relevant)	0.67	0.81	0.73	13503
Class 1 (Relevant)	0.66	0.48	0.55	10511

- Performance on the test set using recall-based optimal parameter

Measure Class	Precision	Recall	F1-measure	Support
Class 0 (Non-Relevant)	0.67	0.79	0.72	13503
Class 1 (Relevant)	0.66	0.50	0.57	10511

Hyper-parameters tuning MLP

- Experimental setup:
 - Grid Search over all combinations of the MLP parameters as follow
 - 'hidden_layer_sizes': [(50,50,50), (50,100,50), (100,)]
 - 'activation': ['tanh', 'relu']
 - 'solver': ['sgd', 'adam']
 - 'alpha': [0.0001, 0.05]
 - 'learning_rate': ['constant', 'adaptive']
 - 5-fold Cross Validation
 - Tuned for both precision and recall measures
- Best parameters set found on development set:
 - For Precision: {'activation': 'relu', 'alpha': 0.05, 'hidden_layer_sizes': (100,), 'learning_rate': 'constant', 'solver': 'adam'}
 - For Recall: {'activation': 'relu', 'alpha': 0.0001, 'hidden_layer_sizes': (100,), 'learning_rate': 'constant', 'solver': 'adam'}

MLP performance on the test set tuned based on precision and recall

- Performance on the test set using precision-based optimal parameter

Measure Class	Precision	Recall	F1-measure	Support
Class 0 (Non-Relevant)	0.66	0.80	0.72	13591
Class 1 (Relevant)	0.64	0.46	0.53	10423

- Performance on the test set using recall-based optimal parameter

Measure Class	Precision	Recall	F1-measure	Support
Class 0 (Non-Relevant)	0.57	1.00	0.73	13591
Class 1 (Relevant)	0.82	0.02	0.05	10423

Hyper-parameters tuning SVM

- Experimental setup:
 - Grid Search over all combinations of the SVM parameters as follow
 - {'kernel': ['rbf'], 'gamma': [1e-3, 1e-4]}
 - {'kernel': ['linear'], 'C': [1, 10, 100, 1000]}
 - 5-fold Cross Validation
 - Tuned for both precision and recall measures
- Best parameters set found on development set:
 - For Precision: {'gamma': 0.001, 'kernel': 'rbf'}
 - For Recall: {'gamma': 0.001, 'kernel': 'rbf'}
 - ❖ **The optimal parameter values for precision and recall are the same**

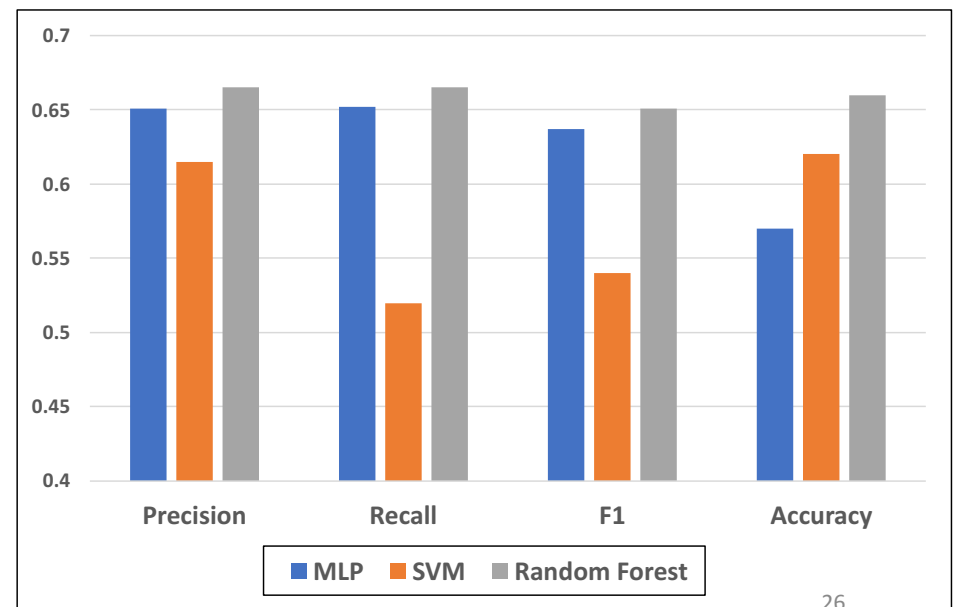
SVM performance on the test set tuned based on precision and recall

- Performance on the test set using precision/recall-based optimal parameter

Measure				
Class	Precision	Recall	F1-measure	Support
Class 0 (Non-Relevant)	0.65	0.70	0.68	13562
Class 1 (Relevant)	0.57	0.52	0.54	10452

Comparing the performance of the classification methods on test set

- All models' performances on the test set using the optimal parameters tuned for precision
- Random Forest outperforms other two models in terms of all measures.
- In terms of running time, Random Forest trained faster than MLP and SVM.



Conclusion and final recommendations

- Standardization decreased the training time of the models.
- The top 5 important features in the dataset according to ANOVA test are sig2, sig7, sig1, sig6 and sig5 respectively.
- PCA and ANOVA feature reduction methods do not improve the classification performance on this dataset.
 - Due to the small number of features (13) compared to the number of instances (80047), dimension reduction increases bias and can cause under-fitting.
- The performance of Random Forest and MLP on the test set is better than SVM in terms of Precision, Recall, F1 and accuracy.
- The Random Forest model trained faster than MLP and SVM.

Conclusion and final recommendations

- The per-class evaluation of all three methods show that the relevant document detection (Class 1) is a harder problem compared to identifying non-relevant documents.
- The relevance optimization problem is similar to the problem of document retrieval in response to query. Therefore, the ranking of documents would be a useful information in our dataset.
- The forth method used in this project is a learning to rank algorithm called SVMRank. The idea of learning the ranking of the documents for each query is an effective approach for the search engine optimization problem.

Python implementation and visual walk-through of the outputs

- The python script attached to this report can be executed to regenerate the results reported on the slides.
- In the first step for running the program, you may read the ReadMe file or enter the following command to see the instructions:
 - `python3 RelevancePrediction.py -h`
- You can pass the following arguments with your command when running the code:
 - -d: The dimension reduction method. (available options are "PCA" and "AnovaTest")
 - -m: The classification model. (available options are "SVM", "MLP" and "RandomForest")
- Example Command:
 - `python3 RelevancePrediction.py -d "AnovaTest" -m "RandomForest"`

Python implementation and visual walk-through of the outputs (Cont.)

- The first function called after entering the command is setup()
- It takes the arguments from input and set the corresponding variables.

```
257 ##### Setup Function #####
258 def setup(argv):
259
260     if len(argv) == 0:
261         print('You must pass some parameters. Use \'-h\' to help.')
262         return
263     if len(argv) == 1 and argv[0] == '-h':
264         f = open('ReadMe.txt', 'r')
265         print(f.read())
266         f.close()
267         return
268     inputFile = "Dataset.csv"#default value
269     outputFolderName = "Results"#default value
270     DimRedMethod = "None"#default value
271     ClassificationMethod = "RandomForest"#default value
272
273     opts, args = getopt.getopt(sys.argv[1:], "d:m:")
274     for opt, arg in opts:
275         if opt == '-d':
276             DimRedMethod = arg
277         elif opt == '-m':
278             ClassificationMethod = arg
279         else:
280             print("Usage: %s -d DimRedMethod -m ClassificationMethod" % sys.argv[0])
281
282     ##### Step One: Preprocessing #####
283     print(".....Preprocessings.....")
284     preprocessedInputDir = preprocessing(inputFile)
285     print("Done!")
286
287     ##### SVM Classification #####
288     if "SVM" in ClassificationMethod:
289         print(".....SVM Classification.....")
290         SVMClassification(preprocessedInputDir, outputFolderName, DimRedMethod)
291         print("Done!")
292
293     ##### RandomForest Classification #####
294     elif "RandomForest" in ClassificationMethod:
295         print(".....Random Forest Classification.....")
296         RandomForestClassification(preprocessedInputDir, outputFolderName, DimRedMethod)
297         print("Done!")
298
```

Python implementation and visual walk-through of the outputs (Cont.)

- Based on the arguments, the next functions are called.

- The next called function is preprocessing.

```
29 ##### Preprocessings Function #####
30 def preprocessing(inputFileDir):
31
32     ##### Duplicate Removal #####
33     print("start preprocessing!")
34     print("inputFileDir: "+str(inputFileDir))
35     dataframe = pd.read_csv(inputFileDir, sep=",")
36     dataframe.drop_duplicates(subset=['query_id','url_id'], keep = 'first', inplace = True)
37
38     ##### Feature Scaling #####
39     array = dataframe.values
40     # separate array into input and output components
41     X = array[:,[2,4,5,6,7,8,9,10,11]]
42     Y = array[:,1]
43     scaler = MinMaxScaler(feature_range=(0, 1))
44     rescaledX = scaler.fit_transform(X)
45
46     ##### Standardization #####
47     # Centering the feature columns at mean 0 with standard deviation
48     scaler = StandardScaler().fit(rescaledX)
49     StandardizedX = scaler.transform(rescaledX)
50
51     preprocX = np.concatenate((array[:,0:2],StandardizedX[:,[0]],array[:,[3]],StandardizedX[:,1:9],array[:,[12]]), axis=1)
52     preprocessedFileDir = ''.join("Preprocessed_"+inputFileDir)
53     with open(preprocessedFileDir,"w") as my_csv:
54         csvWriter = csv.writer(my_csv,delimiter=',')
55         csvWriter.writerow(['query_id','url_id','query_length','is_homepage','sig1','sig2','sig3','sig4','sig5','sig6','sig7','s
56         csvWriter.writerows(preprocX)
57     return preprocessedFileDir
58
```

Python implementation and visual walk-through of the outputs (Cont.)

- Based on the arguments, the next functions are called.

- The next called function is preprocessing.
- In preprocessing function,
 - The duplicate instances are removed.
 - Features are rescaled using min-max normalization.
 - Then, features are standardized.

```
31
32 ##### Preprocessings Function #####
33 def preprocessing(inputFileDir):
34
35     ##### Duplicate Removal #####
36     print("start preprocessing!")
37     print("inputFileDir: "+str(inputFileDir))
38     dataframe = pd.read_csv(inputFileDir, sep=",")
39     dataframe.drop_duplicates(subset=['query_id','url_id'], keep = 'first', inplace = True)
40
41     ##### Feature Scaling #####
42     array = dataframe.values
43     # separate array into input and output components
44     X = array[:,2,4,5,6,7,8,9,10,11]
45     Y = array[:,1]
46     scaler = MinMaxScaler(feature_range=(0, 1))
47     rescaledX = scaler.fit_transform(X)
48
49     ##### Standardization #####
50     # Centering the feature columns at mean 0 with standard deviation
51     scaler = StandardScaler().fit(rescaledX)
52     StandardizedX = scaler.transform(rescaledX)
53
54     preprocX = np.concatenate((array[:,0:2],StandardizedX[:,0]],array[:,3]],StandardizedX[:,1:9],array[:,12])), axis=1)
55     preprocessedFileDir = ''.join("Preprocessed_"+inputFileDir)
56     with open(preprocessedFileDir,"w") as my_csv:
57         csvWriter = csv.writer(my_csv,delimiter=',')
58         csvWriter.writerow(['query_id','url_id','query_length','is_homepage','sig1','sig2','sig3','sig4','sig5','sig6','sig7','s
59         csvWriter.writerow(preprocX)
60     return preprocessedFileDir
```


Python implementation and visual walk-through of the outputs (Cont.)

- After preprocessing, according to the desired classification method set in the arguments, one of the three classification functions are called.

- SVM
- Random Forest
- MLP

```
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309

opts, args = getopt.getopt(sys.argv[1:], "d:m:")
for opt, arg in opts:
    if opt == '-d':
        DimRedMethod = arg
    elif opt == '-m':
        ClassificationMethod = arg
    else:
        print("Usage: %s -d DimRedMethod -m ClassificationMethod" % sys.argv[0])

##### Step One: Preprocessing #####
print(".....Preprocessings.....")
preprocessedInputDir = preprocessing(inputFileName)
print("Done!")

##### SVM Classification #####
if "SVM" in ClassificationMethod:
    print(".....SVM Classification.....")
    SVMClassification(preprocessedInputDir, outputFolderName, DimRedMethod)
    print("Done!")

##### RandomForest Classification #####
elif "RandomForest" in ClassificationMethod:
    print(".....Random Forest Classification.....")
    RandomForestClassification(preprocessedInputDir, outputFolderName, DimRedMethod)
    print("Done!")

##### MLP Classification #####
elif "MLP" in ClassificationMethod:
    print(".....MLP Classification.....")
    MLPClassification(preprocessedInputDir, outputFolderName, DimRedMethod)
    print("Done!")

if __name__ == "__main__":
    setups(sys.argv[1:])
```

SVM Classification

- First, dimension reduction is applied (if set in command options)
- Then, hyper-parameter tuning starts by defining the values of SVM parameters to be searched using Sklearn GridSearchCV
- The GridSearch is performed using 5-fold cross validation over training set for both Precision and Recall
- After finding the optimal parameter values, the model is trained on the training set using the tuned parameter values and tested on the testing set.

```
61 ##### SVM Classification #####
62 def SVMClassification(preprocessedInputDir,outputFolderName,DimRedMethod):
63
64     dataframe = pd.read_csv(preprocessedInputDir, sep=",")
65     array = dataframe.values
66     X = array[:,0:12]
67     Y = array[:,13]
68
69     #writing to outPutFile
70     outFile = open(outputFolderName+"/SVMResults.txt","w")
71
72     #splitting into train and test data
73     X_train,X_test,y_train,y_test = train_test_split(X,Y,test_size=0.3)
74
75     #### Dimension Reduction using PCA ####
76     if(DimRedMethod == "PCA"):
77         pca = PCA(.95)
78         pca.fit(X_train)
79         X_train = pca.transform(X_train)
80         X_test = pca.transform(X_test)
81
82     #### Dimension Reduction using Anova test ####
83     if (DimRedMethod == "AnovaTest"):
84         test = SelectKBest(score_func=f_classif, k=5)
85         fit = test.fit(X_train, y_train)
86         X_train = fit.transform(X_train)
87         X_test = fit.transform(X_test)
88
89     ##### SVM Parameter Tuning #####
90     # Set the parameters by cross-validation
91     tuned_parameters = [{'kernel': ['rbf'], 'gamma': [1e-3, 1e-4], 'C': [1, 10, 100, 1000]},
92                         {'kernel': ['linear'], 'C': [1, 10, 100, 1000]}]
93     scores = {'precision', 'recall'}
94
95     for score in scores:
96         print("# Tuning hyper-parameters for %s" % score)
97         #clf = GridSearchCV(svm.NuSVC(), tuned_parameters, scoring='%s_macro' % score)
98         #clf = GridSearchCV(svm.LinearSVC(), tuned_parameters, scoring='%s_macro' % score)
99         clf = GridSearchCV(SVC(), tuned_parameters, scoring='%s_macro' % score,n_jobs=3)
100         clf.fit(X_train, y_train)
101         print()
102         print("Best parameters set found on development set:")
103         print()
104         print(clf.best_params_)
105         print()
106         print("Grid scores on development set:")
107         print()
108         means = clf.cv_results_['mean_test_score']
109         stds = clf.cv_results_['std_test_score']
110         for mean, std, params in zip(means, stds, clf.cv_results_['params']):
111             print("%0.3f (+/-%0.03f) for %r"
112                   % (mean, std * 2, params))
113         print()
114
115     print("Detailed classification report:")
116     print()
117     print("The model is trained on the full development set.")
118     print("The scores are computed on the full evaluation set.")
119     print()
120     y_true, y_pred = y_test, clf.predict(X_test)
121     print(classification_report(y_true, y_pred))
122     print()
123     outFile.write("the model using optimal parameter value on the test set"+"\\n")
124     outFile.write(classification_report(y_true, y_pred))
125     outFile.close()
```

MLP Classification

- The whole process is the same as SVM, unless the MLP classification module and also MLP parameters for tuning are different.
- Three different network architectures has been tested with different number of layers and neurons.
- Two activation functions are tested including tanh and ReLU
- The result are written in a file in a directory called 'Results'

```
194 ##### Fully Connected MLP Classification #####
195 def MLPClassification(preprocessedInputDir, outputFolderName, DimRedMethod):
196     dataframe = pd.read_csv(preprocessedInputDir, sep=",")
197     array = dataframe.values
198     X = array[:,0:12]
199     Y = array[:,13]
200     X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.3)
201     #writing to outPutFile
202     outFile = open(outputFolderName+"MLPResults.txt", "w")
203
204     #### Dimension Reduction using PCA ####
205     if(DimRedMethod == "PCA"):
206         pca = PCA(.95)
207         pca.fit(X_train)
208         X_train = pca.transform(X_train)
209         X_test = pca.transform(X_test)
210         print("After PCA")
211     #### Dimension Reduction using Anova test ####
212     elif(DimRedMethod == "AnovaTest"):
213         test = SelectKBest(score_func=f_classif, k=5)
214         fit = test.fit(X_train, y_train)
215         X_train = fit.transform(X_train)
216         X_test = fit.transform(X_test)
217
218     #### MLP Parameter Tuning ####
219     tuned_parameters = {
220         'hidden_layer_sizes': [(50,50,50), (50,100,50), (100,)],
221         'activation': ['tanh', 'relu'],
222         'solver': ['sgd', 'adam'],
223         'alpha': [0.0001, 0.05],
224         'learning_rate': ['constant', 'adaptive'],
225     }
226     scores = {'precision', 'recall'}
227     for score in scores:
228         print("# Tuning hyper-parameters for %s" % score)
229         clf = GridSearchCV(MLPClassifier(), tuned_parameters, scoring='%s_macro' % score, n_jobs=3)
230         clf.fit(X_train, y_train)
231         print()
232         print("Best parameters set found on development set:")
233         print()
234         print(clf.best_params_)
235         print()
236         print("Grid scores on development set:")
237         print()
238         means = clf.cv_results_['mean_test_score']
239         stds = clf.cv_results_['std_test_score']
240         for mean, std, params in zip(means, stds, clf.cv_results_['params']):
241             print("%0.3f (+/-0.03f) for %r"
242                   % (mean, std * 2, params))
243         print()
244         print("Detailed classification report:")
245         print()
246         print("The model is trained on the full development set.")
247         print("The scores are computed on the full evaluation set.")
248         print()
249         #Run the model using optimal parameter value on the test set
250         y_true, y_pred = y_test, clf.predict(X_test)
251         print(classification_report(y_true, y_pred))
252         print()
253         outFile.write("the model using optimal parameter value on the test set"+"\\n")
254         outFile.write(classification_report(y_true, y_pred))
255     outFile.close()
```

Random Forest Classification

- Due to the large number of hyper-parameters in Random Forest, in order to reduce the running time, I utilized parallel processing by increasing the value of `n_jobs` parameter in Sklearn.

```
127 ##### Random Forest Classification #####
128 def RandomForestClassification(preprocessedInputDir, outputFolderName, DimRedMethod):
129     dataframe = pd.read_csv(preprocessedInputDir, sep=",")
130     array = dataframe.values
131     X = array[:,0:12]
132     Y = array[:, -1]
133
134     #writing to outPutFile
135     outFile = open(outputFolderName+"/RandomForestResults.txt", "w")
136
137     #splitting into train and test data
138     X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.3)
139
140     #### Dimension Reduction using PCA ####
141     if (DimRedMethod == "PCA"):
142         pca = PCA(.95)
143         pca.fit(X_train)
144         X_train = pca.transform(X_train)
145         X_test = pca.transform(X_test)
146
147     #### Dimension Reduction using Anova test ####
148     if (DimRedMethod == "AnovaTest"):
149         test = SelectKBest(score_func=f_classif, k=5)
150         fit = test.fit(X_train, y_train)
151         X_train = fit.transform(X_train)
152         X_test = fit.transform(X_test)
153
154     ##### RF Parameter Tuning #####
155     #hyper-parameter tuning- Selecting best n_estimators
156     tuned_parameters = {'bootstrap': [True, False],
157                         'max_depth': [10, 20, 30, 40, 50, 60, 70, 80, 90, 100, None],
158                         'max_features': ['auto', 'sqrt'],
159                         'min_samples_leaf': [1, 2, 4],
160                         'min_samples_split': [2, 5, 10],
161                         'n_estimators': [200, 400, 600, 800, 1000, 1200, 1400, 1600, 1800, 2000]}
162     scores = ['precision', 'recall']
163     for score in scores:
164         print("# Tuning hyper-parameters for %s" % score)
165         clf = GridSearchCV(RandomForestClassifier(), tuned_parameters, scoring='%s_macro' % score, n_jobs=3)
166         clf.fit(X_train, y_train)
167         print()
168         print("Best parameters set found on development set:")
169         print()
170         print(clf.best_params_)
171         print()
172         print("Grid scores on development set:")
173         print()
174         means = clf.cv_results_['mean_test_score']
175         stds = clf.cv_results_['std_test_score']
176         for mean, std, params in zip(means, stds, clf.cv_results_['params']):
177             print("%0.3f (+/-%0.03f) for %r"
178                   % (mean, std * 2, params))
179         print()
180         print("Detailed classification report:")
181         print()
182         print("The model is trained on the full development set.")
183         print("The scores are computed on the full evaluation set.")
184         print()
185         #Run the model using optimal parameter value on the test set
186         y_true, y_pred = y_test, clf.predict(X_test)
187         print(classification_report(y_true, y_pred))
188         print()
189         outFile.write("the model using optimal parameter value on the test set"+"\\n")
190         outFile.write(classification_report(y_true, y_pred))
191     outFile.close()
```

Thanks for your attention