```python
import pandas as pd
import numpy as np
from sklearn.svm import SVC
import matplotlib.pyplot as plt
from sklearn.preprocessing import StandardScaler


#Importing the Data from the give files
train_set = pd.read_csv("tweets-train-data.csv", sep= ",", header = None, names =
["Tweet", "Date", "Retweets", "Likes", "Location"])
test_set = pd.read_csv("tweets-test-data.csv", sep= ",", header = None, names =
["Tweet", "Date", "Retweets", "Likes", "Location"])
train_targets = pd.read_csv("tweets-train-targets.csv", sep= ",", header = None,
names = ["Author"])
test_targets = pd.read_csv("tweets-test-targets.csv", sep= ",", header = None, names
= ["Author"])

#Formating the data
training = pd.concat([train_set, train_targets], axis = 1)
training = training[training["Date"].astype(str).str.startswith("2016")]
training["Likes"] = training["Likes"].astype(float)
training["Date"], training["Time"] = training["Date"].str.split("T").str
training["Date"].astype("datetime64")
print(training)

#Likes information
Trump = training[training["Author"] == "DT"]
Clinton = training[training["Author"] == "HC"]
Clinton_likes = Clinton["Likes"]
Trump_likes = Trump["Likes"]
print(Trump_likes)
print(Clinton_likes)

#Ploting the retweets
Clinton_rt = Clinton["Retweets"]
Trump_rt = Trump["Retweets"]

plt.plot(Clinton_rt, "c")
plt.plot(Trump_rt, "r")
plt.xlabel("Tweet id")
plt.ylabel("Number of rts")
plt.legend(("HC", "DT"))
plt.grid(True)

plt.show()
plt.close()

#Clinton most popular tweet (rts)
print(max(Clinton_rt))
top_Clinton_rt = training.loc[training["Retweets"]==max(Clinton_rt)]
print(top_Clinton_rt)

#Trump most popular tweet (rts)
print(max(Trump_rt))
top_Trump_rt = training.loc[training["Retweets"]==max(Trump_rt)]
print(top_Trump_rt)

#Clinton most liked tweet
print(max(Clinton_likes))
top_Clinton_likes = training.loc[training["Likes"]==max(Clinton_likes)]
print(top_Clinton_likes)

#Trump most liked tweet
print(max(Trump_likes))
top_Trump_likes = training.loc[training["Likes"]==max(Trump_likes)]
print(top_Trump_likes)

#Formating the testing data
testing = pd.concat([test_set, test_targets], axis = 1)
testing = testing[testing["Date"].astype(str).str.startswith("2016")]
testing["Likes"] = testing["Likes"].astype(float)
testing["Date"], testing["Time"] = testing["Date"].str.split("T").str
```

```python
69    testing["Date"].astype("datetime64")
70    print(testing)
71
72    #Extracting the features
73    train_features = training["Author"]
74    train_rts = training["Retweets"].values.reshape(-1,1)
75    test_features = testing["Author"]
76    test_rts = testing["Retweets"].values.reshape(-1,1)
77
78
79    #Normalize the tweets length
80    a = np.mean(training["Tweet"].apply(len))
81    b = np.mean(testing["Tweet"].apply(len))
82
83    training["word_length"] = training["Tweet"].apply(len)-a
84    testing["word_length"] = testing["Tweet"].apply(len)-b
85
86    #Re-assign
87    Trump = training[training["Author"] == "DT"]
88    Clinton = training[training["Author"] == "HC"]
89
90    Trump_average_length = Trump["word_length"]
91    Clinton_average_length = Clinton["word_length"]
92
93    #Spaces count
94    training["space_count"] = training["Tweet"].str.count(" ")
95    training["a_count"] = training["Tweet"].str.count("a")
96    testing["space_count"] = testing["Tweet"].str.count(" ")
97    testing["a_count"] = testing["Tweet"].str.count("a")
98
99    #Taking logarithms (monotonic function)
100   training["log_likes"] = np.log(training["Likes"])
101   training["log_rts"] = np.log(training["Retweets"])
102   testing["log_likes"] = np.log(testing["Likes"])
103   testing["log_rts"] = np.log(testing["Retweets"])
104
105   #Dot count
106   training["dot_count"] = training["Tweet"].str.count(".")
107   testing["dot_count"] = testing["Tweet"].str.count(".")
108
109   #Hours
110   training["hour"] = training["Time"].str[0:2]
111   pd.to_numeric(training["hour"])
112   testing["hour"] = testing["Time"].str[0:2]
113   pd.to_numeric(testing["hour"])
114
115   #Months
116   training["month"] = training["Date"].str[5:7]
117   pd.to_numeric(training["month"])
118   testing["month"] = testing["Date"].str[5:7]
119   pd.to_numeric(testing["month"])
120
121   scale = StandardScaler()
122
123   log_training_matrix = training[["log_likes", "log_rts", "word_length", "a_count",
      "space_count", "dot_count", "month", "hour", "Retweets"]]
124   log_testing_matrix = testing[["log_likes", "log_rts", "word_length", "a_count",
      "space_count", "dot_count", "month", "hour", "Retweets"]]
125
126   transformed_x = scale.fit_transform(log_training_matrix, train_features)
127   transformed_y = scale.fit_transform(log_testing_matrix, test_features)
128
129   #Training using SVC
130
131   clf = SVC(C=10, kernel = "rbf", gamma = 0.02)
132   clf.fit(transformed_x, train_features)
133
134   y_predict = clf.predict(transformed_y)
135
136   from sklearn import metrics
137
138   report = metrics.classification_report(test_features, y_predict)
```

```python
139     metrics.accuracy_score(test_features, y_predict)
140
141     #K-fold cross Validation
142     try:
143         from sklearn.model_selection import KFold, cross_val_score
144         legacy = False
145     except ImportError:
146         from sklearn.model_selection import KFold, cross_val_score
147         legacy = True
148
149     #K=3
150     if legacy:
151         kf = KFold(len(train_features), n_folds = 3, shuffle = True, random_state = 42)
152     else:
153         kf = KFold(n_splits = 3, shuffle = True, random_state = 42)
154
155     gamma_values = [0.01, 0.001, 0.0001, 0.00001]
156     accuracy_scores = []
157
158     #K-Fold cross validation algorithm:
159     #-Train predictor
160     #-Compute score
161     for gamma in gamma_values:
162
163         #Train classifier
164         clf = SVC(C=10, kernel = "rbf", gamma = gamma)
165
166         #Score
167         if legacy:
168             scores = cross_val_score(clf, transformed_x, train_features, cv = kf,
                    scoring="accuracy")
169         else:
170             scores = cross_val_score(clf, transformed_x, train_features, cv =
                    kf.split(transformed_x), scoring = "accuracy")
171
172         #Compute Average score
173         accuracy_score = scores.mean()
174         accuracy_scores.append(accuracy_score)
175
176     #Best value of gamma
177     best_index = np.array(accuracy_scores).argmax()
178     best_gamma = gamma_values[best_index]
179
180     #Train with the selected gamma
181     clf = SVC(C=10, kernel = "rbf", gamma = best_gamma)
182     clf.fit(transformed_x, train_features)
183
184     #Evaluate on the test set
185     y_predict = clf.predict(transformed_y)
186     accuracy = metrics.accuracy_score(test_features, y_predict)
187
188     print(accuracy)
189
190     print(best_gamma)
191
192     from sklearn.model_selection import learning_curve
193
194     plt.figure(2)
195     plt.title("Learning curve")
196     plt.xlabel("Training examples")
197     plt.ylabel("Score")
198     plt.grid(True)
199
200
201     clf = SVC(C=10, kernel = "rbf", gamma = best_gamma)
202
203     train_sizes, train_scores, val_scores = learning_curve(clf, transformed_x,
            train_features, scoring = "accuracy")
204
205     train_scores_mean = np.mean(train_scores, axis = 1)
206     train_scores_std = np.std(train_scores, axis = 1)
207     val_scores_mean = np.mean(val_scores, axis = 1)
```

```
208    val_scores_std = np.std(val_scores, axis = 1)
209
210    #Training scores: mean
211    plt.plot(train_sizes, train_scores_mean, '*' ,color = "r", label = "Training Score")
212    #Training scores: std
213    plt.fill_between(train_sizes, train_scores_mean - train_scores_std,
       train_scores_mean + train_scores_std, alpha = 0.1, color = "r")
214    #Validation scores: mean
215    plt.plot(train_sizes, val_scores_mean,'*', color = "c", label = "Cross Validation
       Scores")
216    #Validation scores: std
217    plt.fill_between(train_sizes, val_scores_mean - val_scores_std, val_scores_mean +
       val_scores_std, alpha = 0.1, color = "c")
218
219    plt.ylim(0.05, 1.3)
220    plt.legend()
221    plt.show()
222    plt.close()
223
224    try:
225        from sklearn.model_selection import GridSearchCV
226    except ImportError:
227        from sklearn.grid_search import GridSearchCV
228    possible_parameters = {
229            'C' : [10,4,3,2],
230            'gamma' : [0.01,0.001,0.0001,0.00001]
231            }
232
233    svc = SVC(kernel = "rbf")
234
235    #The GridSearch is a classifier. Hence, we try to fit it with the training data,
236    #and then use it for prediction
237    clf = GridSearchCV(svc, possible_parameters, n_jobs = 4, cv = 3)
238    clf.fit(transformed_x, train_features)
239
240    y_predidct = clf.predict(transformed_y)
241    accuracy = metrics.accuracy_score(test_features, y_predict)
242
243    #Training phase
244    clf = SVC(C=10, kernel = "rbf", gamma = 0.01)
245    clf.fit(transformed_x, train_features)
246    #Prediction phase
247    y_predict = clf.predict(transformed_y)
248
249    report = metrics.classification_report(test_features, y_predict)
250    metrics.accuracy_score(test_features, y_predict)
251
252
253    print(report)
254
255    #Performance parameters
256    from sklearn.model_selection import cross_val_score
257
258    accuracy = cross_val_score(clf, transformed_x, train_features, cv = 10, scoring =
       "accuracy")
259
260    #Labels
261    labels = train_features.map(lambda x: 1 if x == "HC" else 0).values
262
263    precision = cross_val_score(clf, transformed_x, labels, cv = 10, scoring =
       "precision")
264
265    recall = cross_val_score(clf, transformed_x, labels, cv = 10, scoring = "recall")
266
267    f1 = cross_val_score(clf, transformed_x, labels, cv=10, scoring = "f1")
268
269    print("Avg Precision: {}".format(round(precision.mean(), 3)))
270    print("Avg Recall: {}".format(round(recall.mean(), 3)))
271    print("Avg Accuracy: {}".format(round(accuracy.mean(), 3)))
272    print("Avg f1: {}".format(round(f1.mean(), 3)))
273
274
```

```python
275    #Retweets plot
276    plt.hist(training["Retweets"][training["Author"]=="HC"], color = "c")
277    plt.hist(training["Retweets"][training["Author"]=="DT"], color = "r")
278    plt.title("Retweets")
279    plt.legend()
280    plt.grid(True)
281    plt.show()
282    plt.close()
283
284
285    #Log Retweets plot
286    plt.hist(training["log_rts"][training["Author"]=="HC"], color = "c")
287    plt.hist(training["log_rts"][training["Author"]=="DT"], color = "r")
288    plt.title("Log Retweets")
289    plt.legend()
290    plt.grid(True)
291    plt.show()
292    plt.close()
293
294    #Likes plot
295    plt.hist(training["Likes"][training["Author"]=="HC"], color = "c")
296    plt.hist(training["Likes"][training["Author"]=="DT"], color = "r")
297    plt.title("Likes")
298    plt.legend()
299    plt.grid(True)
300    plt.show()
301    plt.close()
302
303    #Log likes plot
304    plt.hist(training["log_likes"][training["Author"]=="HC"], color = "c")
305    plt.hist(training["log_likes"][training["Author"]=="DT"], color = "r")
306    plt.title("Log Likes")
307    plt.legend()
308    plt.grid(True)
309    plt.show()
310    plt.close()
311
```