





```
In [ ]: from sklearn.model_selection import GridSearchCV

# defining parameter range
param_grid = {'C': [0.1, 1, 10],
              'gamma': [1, 0.1],
              'kernel': ['rbf']}

grid = GridSearchCV(SVC(probability = True), param_grid, refit = True, verbose = 3)

# fitting the model for grid search
grid.fit(train_df, ydata)
```

## Apply Model

```
In [35]: from sklearn.neighbors import KNeighborsClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier
from sklearn.tree import DecisionTreeClassifier
```

## Train Model

```
In [ ]: ydata

In [23]: model2 = RandomForestClassifier(n_estimators = 100)

In [ ]: model2.fit(X, ydata)

In [40]: from sklearn.svm import SVC

In [ ]: svcclassifier = SVC(kernel='rbf', probability = True)
svcclassifier.fit(X, ydata)

In [ ]: from sklearn.calibration import CalibratedClassifierCV

In [ ]: svm = LinearSVC()
clf = CalibratedClassifierCV(svm)
clf.fit(train_df, ydata)

In [ ]: logisticRegr = LogisticRegression()

In [ ]: model = LogisticRegression()
model.fit(train_df, ydata)
```

## Column to predict :

```
In [ ]: Y = vectorizer.transform(test_d['page_description'])
```

## Word2Vec Implementation

```
In [22]: import gensim
from gensim.models import Word2Vec, KeyedVectors
l=0
sentences_to_train = []
sentences_to_test=[]
for sentences in train_d['page_description']:
    sentences_to_train.append(sentences.split())
for sentences in test_d['page_description']:
    sentences_to_test.append(sentences.split())

In [23]: w2v_model = Word2Vec(sentences_to_train, min_count=2, vector_size=200, workers=4)
#w2v_model=KeyedVectors.load_word2vec_format('GoogleNews-vectors-negative300.bin', binary=True)

In [24]: text = train_d['page_description'].apply(gensim.utils.simple_preprocess)
test = test_d['page_description'].apply(gensim.utils.simple_preprocess)
train_d['web'] = text
#test_d['web'] = test

In [25]: make_vector = []
def avg_word2_vec(df):
    sc = np.zeros(200)
    i=0
    for text1 in df:
        if text1 not in w2v_model.wv.key_to_index:
            continue
        sq = w2v_model.wv.get_vector(text1)
        sc += sq
        i+=1
    sc = sc/i
    sc = np.array(sc)
    make_vector.append(sc)
    return sc

In [26]: make_test=[]
def avg_word2_test(df):
    sc = np.zeros(200)
    i=0
    for text1 in df:
        if text1 not in w2v_model.wv.key_to_index:
            continue
        sq = w2v_model.wv.get_vector(text1)
        sc += sq
        i+=1
    sc = sc/i
    sc = np.array(sc)
    make_test.append(sc)
    return sc

In [27]: vectorized = train_d['web'].apply(Lambda x:avg_word2_vec(x))
vectorizedtest = test_d['web'].apply(Lambda x:avg_word2_test(x))

In [28]: make_vector = np.array(make_vector)
make_test = np.array(make_test)
make_vector.shape

Out[28]: (4437, 200)
```

```
In [40]: from sklearn.model_selection import GridSearchCV
from sklearn.svm import SVC, NuSVC
param_grid = {'C': [0.1, 1, 10, 100], 'gamma': [1.0, 1.0, 0.01, 0.001]}
grid = GridSearchCV(SVC(probability=True), param_grid, refit=True, verbose=2)
grid.fit(np.array(make_vector), train_d['label'])

Fitting 5 folds for each of 16 candidates, totalling 80 fits
[CV] END .....C=0.1, gamma=1; total time= 9.6s
[CV] END .....C=0.1, gamma=1; total time= 9.5s
[CV] END .....C=0.1, gamma=1; total time= 9.4s
[CV] END .....C=0.1, gamma=1; total time= 9.3s
[CV] END .....C=0.1, gamma=0.1; total time= 6.5s
[CV] END .....C=0.1, gamma=0.1; total time= 6.6s
[CV] END .....C=0.1, gamma=0.1; total time= 7.1s
[CV] END .....C=0.1, gamma=0.01; total time= 7.3s
[CV] END .....C=0.1, gamma=0.01; total time= 7.5s
[CV] END .....C=0.1, gamma=0.01; total time= 7.4s
[CV] END .....C=0.1, gamma=0.001; total time= 11.2s
[CV] END .....C=0.1, gamma=0.001; total time= 10.7s
[CV] END .....C=0.1, gamma=0.001; total time= 10.6s
[CV] END .....C=1, gamma=1; total time= 10.5s
[CV] END .....C=1, gamma=1; total time= 10.6s
[CV] END .....C=1, gamma=1; total time= 10.5s
[CV] END .....C=1, gamma=0.1; total time= 6.6s
[CV] END .....C=1, gamma=0.1; total time= 6.6s
[CV] END .....C=1, gamma=0.1; total time= 6.5s
[CV] END .....C=1, gamma=0.01; total time= 6.3s
[CV] END .....C=1, gamma=0.01; total time= 6.1s
[CV] END .....C=1, gamma=0.01; total time= 7.1s
[CV] END .....C=1, gamma=0.001; total time= 7.3s
[CV] END .....C=1, gamma=0.001; total time= 7.1s
[CV] END .....C=10, gamma=1; total time= 11.6s
[CV] END .....C=10, gamma=1; total time= 11.4s
[CV] END .....C=10, gamma=0.1; total time= 11.8s
[CV] END .....C=10, gamma=0.1; total time= 7.4s
[CV] END .....C=10, gamma=0.1; total time= 7.3s
[CV] END .....C=10, gamma=0.01; total time= 7.3s
[CV] END .....C=10, gamma=0.01; total time= 6.9s
[CV] END .....C=10, gamma=0.01; total time= 5.9s
[CV] END .....C=10, gamma=0.01; total time= 5.9s
[CV] END .....C=10, gamma=0.001; total time= 6.2s
[CV] END .....C=10, gamma=0.001; total time= 6.3s
[CV] END .....C=10, gamma=0.001; total time= 6.2s
[CV] END .....C=10, gamma=0.001; total time= 6.0s
[CV] END .....C=100, gamma=1; total time= 11.1s
[CV] END .....C=100, gamma=0.1; total time= 11.0s
[CV] END .....C=100, gamma=0.1; total time= 9.7s
[CV] END .....C=100, gamma=0.1; total time= 9.9s
[CV] END .....C=100, gamma=0.1; total time= 10.3s
[CV] END .....C=100, gamma=0.1; total time= 9.7s
[CV] END .....C=100, gamma=0.1; total time= 9.9s
[CV] END .....C=100, gamma=0.01; total time= 7.0s
[CV] END .....C=100, gamma=0.01; total time= 7.2s
[CV] END .....C=100, gamma=0.01; total time= 7.1s
[CV] END .....C=100, gamma=0.01; total time= 7.2s
[CV] END .....C=100, gamma=0.01; total time= 5.8s
[CV] END .....C=100, gamma=0.01; total time= 5.9s
[CV] END .....C=100, gamma=0.01; total time= 5.8s
[CV] END .....C=100, gamma=0.01; total time= 5.9s

Out[40]: GridSearchCV(estimator=SVC(probability=True),
                      param_grid={'C': [0.1, 1, 10, 100],
                                   'gamma': [1, 0.1, 0.01, 0.001]},
                      verbose=2)
```

```
In [41]: print(grid.best_estimator_)
SVC(C=100, gamma=0.001, probability=True)

In [39]: model9=SVC(C=100, gamma=0.001, probability=True)
model9.fit(np.array(make_vector), train_d['label'])

Out[39]: SVC(C=1, gamma=0.1, probability=True)
```

```
In [ ]: ypred = model9.predict_proba(np.array(make_test))
```

## Slacking :

```
In [ ]: ypred1 = model1.predict_proba(Y).T[1]
ypred2 = svcclassifier.predict_proba(Y).T[1]
ypred3 = model2.predict_proba(Y).T[1]
ypred = (ypred1+ypred2+ypred3)/3
```

## Vooting :

```
In [36]: from sklearn.ensemble import VotingClassifier
from sklearn.tree import DecisionTreeClassifier
```

## create a voting classifier with hard voting

```
In [41]: voting_classifier_hard = VotingClassifier(
estimators = [('dtc', DecisionTreeClassifier(random_state=42)),
              ('lr', LogisticRegression()),
              ('gmb', SVC())],
          voting='hard')
```

## create a voting classifier with soft voting

```
In [42]: voting_classifier_soft = VotingClassifier(
estimators = [('dtc', DecisionTreeClassifier(random_state=42)),
              ('lr', LogisticRegression()),
              ('gmb', SVC())],
          voting='soft')
```

## Make prediction using hard and soft vooting :

```
In [ ]: voting_classifier_hard.fit(train_df, ydata)
y_pred_vch = voting_classifier_hard.predict(Y)

voting_classifier_soft.fit(train_df, ydata)
y_pred_vcs = voting_classifier_soft.predict(Y)

In [ ]:

In [ ]: ypred.shape

In [ ]: ypred = ypred.T[1]
ypred = ypred.T

In [ ]: ypred.shape

In [ ]: ypred

In [ ]: sub = pd.read_csv("../input/aid-escalating-internet-coverage/sample_submission.csv")
sub.head()
```

"if model after suffling dataframe based on alchemy category is used uncomment line 2"

```
In [ ]: sub['label'] = ypred
#sub['link_id'] = testdf['link_id'] #if model after suffling dataframe based on alchemy category is used
sub.head()
```

## Expoting the output file:

```
In [ ]: sub.to_csv('submission14.csv', index = False)
```

```
In [44]: from sklearn.model_selection import StratifiedKFold
from sklearn.metrics import roc_auc_score as rac
```

## Checking Accuracy using ROC Curve and AUC:

```
In [ ]: skf = StratifiedKFold(n_splits = 4, shuffle = True)
ydata = train_d['label']
print(skf)
model = KNeighborsClassifier(n_neighbors=50, leaf_size = 20)
#StratifiedKFold(n_splits=4, random_state=None, shuffle=True)
racs = []
for train_index, test_index in skf.split(train_df, ydata):
    #print("TRAIN:", train_index, "TEST:", test_index)
    X_train, X_test = train_df[train_index], train_df[test_index]
    y_train, y_test = ydata[train_index], ydata[test_index]
    model.fit(X_train, y_train)
    ypred = model.predict_proba(X_test)
    ypred = ypred.T[1]
    racs.append(rac(y_test, ypred))
print(np.mean(racs))

In [ ]:

In [ ]:
```