

sms-spam-v2-stop

April 23, 2019

V2: + Delete the stop words + All models uses the CountVectorizer to do the the preprocessing
Goal of this notebook to test several classifiers on the data set with different features

0.0.1 Let's begin

First of all necessary imports

```
In [1]: import numpy as np
import pandas as pd
import matplotlib as mpl
import matplotlib.pyplot as plt
import seaborn as sns
import string
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.model_selection import train_test_split
from nltk.corpus import stopwords
from sklearn.preprocessing import LabelEncoder
%matplotlib inline
```

Let's read the data from csv file

```
In [2]: sms = pd.read_csv('../data/sms-spam.csv', delimiter=',', encoding='latin-1')

sms.head()
```

```
Out[2]:
```

	v1	v2	Unnamed: 2	\
0	ham	Go until jurong point, crazy.. Available only ...		NaN
1	ham	Ok lar... Joking wif u oni...		NaN
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...		NaN
3	ham	U dun say so early hor... U c already then say...		NaN
4	ham	Nah I don't think he goes to usf, he lives aro...		NaN

	Unnamed: 3	Unnamed: 4
0	NaN	NaN
1	NaN	NaN
2	NaN	NaN
3	NaN	NaN
4	NaN	NaN

Now drop “unnamed” columns and rename v1 and v2 to “label” and “message”

```
In [3]: sms = sms.drop(['Unnamed: 2', 'Unnamed: 3', 'Unnamed: 4'], axis=1)
        sms = sms.rename(columns = {'v1': 'label', 'v2': 'message'})
```

Let’s look into our data

```
In [4]: sms.groupby('label').describe()
```

```
Out[4]:
```

	message	count	unique	top freq
label				
ham		4825	4516	Sorry, I'll call later 30
spam		747	653	Please call our customer service representativ... 4

Intresting that “Sorry, I’ll call later” appears only 30 times here =)

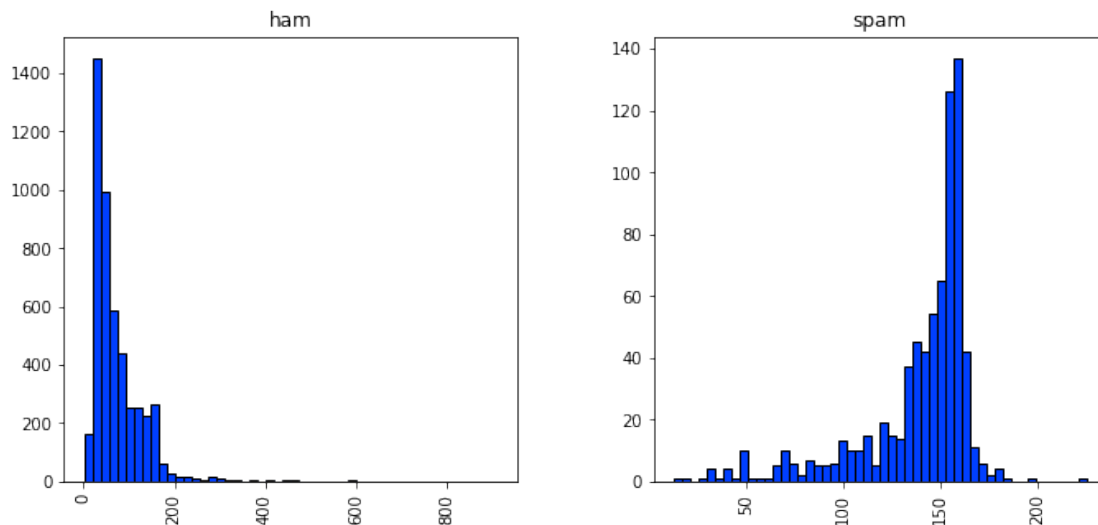
Now let’s create new feature “message length” and plot it to see if it’s of any interest

```
In [5]: sms['length'] = sms['message'].apply(len)
        sms.head()
```

```
Out[5]:
```

	label	message	length
0	ham	Go until jurong point, crazy.. Available only ...	111
1	ham	Ok lar... Joking wif u oni...	29
2	spam	Free entry in 2 a wkly comp to win FA Cup fina...	155
3	ham	U dun say so early hor... U c already then say...	49
4	ham	Nah I don't think he goes to usf, he lives aro...	61

```
In [6]: mpl.rcParams['patch.force_edgecolor'] = True
        plt.style.use('seaborn-bright')
        sms.hist(column='length', by='label', bins=50, figsize=(11,5))
        plt.savefig("../img/sms-length-distribution.eps")
        plt.show()
```



Looks like the lengthy is the message, more likely it is a spam. Let’s not forget this

0.0.2 Text processing and vectorizing our meddages

Let's create new data frame. We'll need a copy later on

```
In [7]: text_feat = sms['message'].copy()
```

Now define our tex precessing function. It will remove any punctuation and stopwords aswell.

```
In [8]: def text_process(text):  
  
    text = text.translate(str.maketrans('', '', string.punctuation))  
    text = [word for word in text.split() if word.lower() not in stopwords.words('engl.  
  
    return " ".join(text)
```

```
In [9]: text_feat = text_feat.apply(text_process)
```

```
In [10]: vectorizer = CountVectorizer("english")
```

```
In [11]: features = vectorizer.fit_transform(text_feat)
```

```
In [12]: labels = LabelEncoder().fit_transform(sms['label'])  
labels = labels.reshape(-1,1)
```

```
In [13]: text_feat.shape
```

```
Out[13]: (5572,)
```

```
In [14]: features.shape
```

```
Out[14]: (5572, 9403)
```

0.0.3 Classifiers and predictions

First of all let's split our features to test and train set

```
In [15]: features_train, features_test, labels_train, labels_test = train_test_split(features,
```

Now let's import bunch of classifiers, initialize them and make a dictionary to itereate through

```
In [16]: from sklearn.linear_model import LogisticRegression  
from sklearn.svm import SVC  
from sklearn.naive_bayes import MultinomialNB  
from sklearn.tree import DecisionTreeClassifier  
from sklearn.neighbors import KNeighborsClassifier  
from sklearn.ensemble import RandomForestClassifier  
from sklearn.ensemble import AdaBoostClassifier  
from sklearn.ensemble import BaggingClassifier  
from sklearn.ensemble import ExtraTreesClassifier  
from sklearn.metrics import precision_score, recall_score, accuracy_score, f1_score
```

```
/Users/alex/anaconda/envs/gc/lib/python3.6/site-packages/sklearn/ensemble/weight_boosting.py:20
from numpy.core.umath_tests import inner1d
```

```
In [17]: svc = SVC(kernel='sigmoid', gamma=1.0)
         knc = KNeighborsClassifier(n_neighbors=49)
         mnb = MultinomialNB(alpha=0.2)
         dtc = DecisionTreeClassifier(min_samples_split=7, random_state=111)
         lrc = LogisticRegression(solver='liblinear', penalty='l1')
         rfc = RandomForestClassifier(n_estimators=31, random_state=111)
         abc = AdaBoostClassifier(n_estimators=62, random_state=111)
         bc = BaggingClassifier(n_estimators=9, random_state=111)
         etc = ExtraTreesClassifier(n_estimators=9, random_state=111)
```

```
In [18]: clfs = {'SVC' : svc, 'KN' : knc, 'NB': mnb, 'DT': dtc, 'LR': lrc, 'RF': rfc, 'AdaBoost' : abc, 'ET': etc}
```

Let's make functions to fit our classifiers and make predictions

```
In [19]: def train_classifier(clf, feature_train, labels_train):
         clf.fit(feature_train, labels_train)
```

```
In [20]: def predict_labels(clf, features):
         return (clf.predict(features))
```

Now iterate through classifiers and save the results

```
In [21]: import time
```

```
In [22]: pred_scores = []
         for k,v in clfs.items():
             since = time.time()

             train_classifier(v, features_train, labels_train)
             time_elapsed = time.time() - since

             pred = predict_labels(v, features_test)
             pred_scores.append((k, [precision_score(labels_test, pred), recall_score(labels_test, pred)]))
```

```
/Users/alex/anaconda/envs/gc/lib/python3.6/site-packages/sklearn/utils/validation.py:578: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using y = column_or_1d(y, warn=True)
/Users/alex/anaconda/envs/gc/lib/python3.6/site-packages/ipykernel/__main__.py:2: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using y = column_or_1d(y, warn=True)
from ipykernel import kernelapp as app
/Users/alex/anaconda/envs/gc/lib/python3.6/site-packages/sklearn/metrics/classification.py:113: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using y = column_or_1d(y, warn=True)
'precision', 'predicted', average, warn_for)
/Users/alex/anaconda/envs/gc/lib/python3.6/site-packages/sklearn/metrics/classification.py:113: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using y = column_or_1d(y, warn=True)
'precision', 'predicted', average, warn_for)
/Users/alex/anaconda/envs/gc/lib/python3.6/site-packages/ipykernel/__main__.py:2: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using y = column_or_1d(y, warn=True)
from ipykernel import kernelapp as app
/Users/alex/anaconda/envs/gc/lib/python3.6/site-packages/sklearn/utils/validation.py:578: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using y = column_or_1d(y, warn=True)
```

```

y = column_or_1d(y, warn=True)
/Users/alex/anaconda/envs/gc/lib/python3.6/site-packages/sklearn/utils/validation.py:578: Data
y = column_or_1d(y, warn=True)
/Users/alex/anaconda/envs/gc/lib/python3.6/site-packages/ipykernel/__main__.py:2: DataConversi
from ipykernel import kernelapp as app

```

```

In [23]: df = pd.DataFrame.from_items(pred_scores,orient='index', columns=['Precision', 'Recall', 'F1', 'Training Time (s)'])
df

```

```

Out[23]:

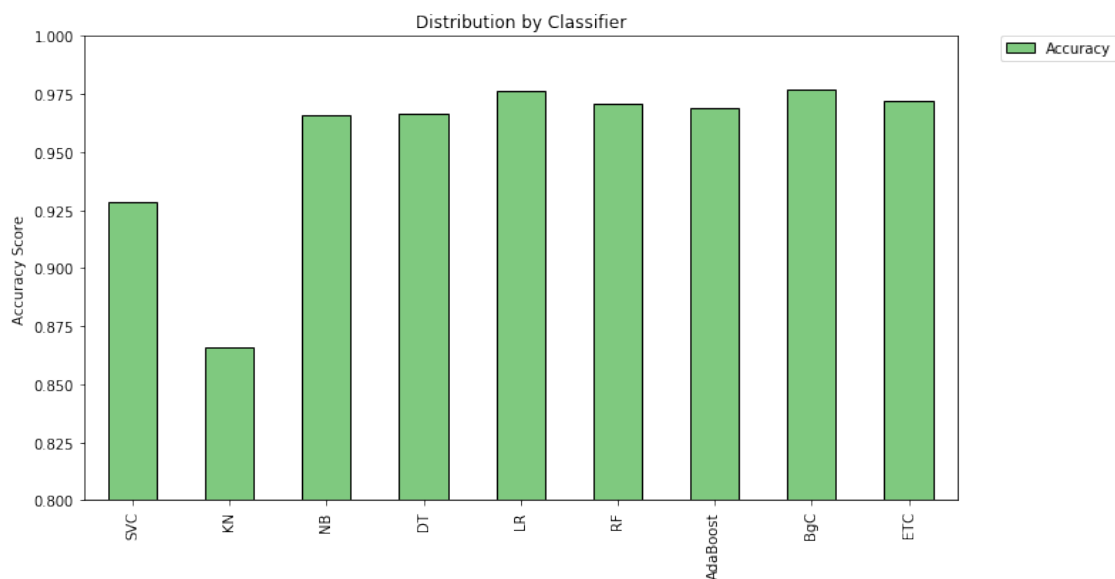
```

	Precision	Recall	Accuracy	F1	Training Time (s)
SVC	0.747619	0.700893	0.928230	0.723502	0m 0.1977s
KN	0.000000	0.000000	0.866029	0.000000	0m 0.0009s
NB	0.819923	0.955357	0.965909	0.882474	0m 0.0019s
DT	0.900000	0.843750	0.966507	0.870968	0m 0.1379s
LR	0.989362	0.830357	0.976077	0.902913	0m 0.0108s
RF	0.988827	0.790179	0.970694	0.878412	0m 1.2109s
AdaBoost	0.934343	0.825893	0.968900	0.876777	0m 2.3704s
BgC	0.969543	0.852679	0.976675	0.907363	0m 0.8051s
ETC	0.963351	0.821429	0.971890	0.886747	0m 0.7646s

```

In [24]: df.plot(kind='bar', y="Accuracy", ylim=(0.8,1.0), figsize=(11,6), align='center', color='green')
plt.xticks(np.arange(9), df.index)
plt.ylabel('Accuracy Score')
plt.title('Distribution by Classifier')
plt.legend(bbox_to_anchor=(1.05, 1), loc=2, borderaxespad=0.)
plt.savefig("../img/sms-acc-basemodel-v2-stop.eps")
plt.show()

```



Looks like ensemble classifiers are not doing as good as expected.

0.0.4 Voting classifier

We are using ensemble algorithms here, but what about ensemble of ensembles? Will it beat NB?

```
In [25]: from sklearn.ensemble import VotingClassifier
```

```
In [26]: eclf = VotingClassifier(estimators=[('BgC', bc), ('ETC', etc), ('RF', rfc), ('Ada', adaboost)])
```

```
In [27]: eclf.fit(features_train, labels_train)
```

```
/Users/alex/anaconda/envs/gc/lib/python3.6/site-packages/sklearn/preprocessing/label.py:95: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using y = column_or_1d(y, warn=True)
/Users/alex/anaconda/envs/gc/lib/python3.6/site-packages/sklearn/preprocessing/label.py:128: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using y = column_or_1d(y, warn=True)
```

```
Out[27]: VotingClassifier(estimators=[('BgC', BaggingClassifier(base_estimator=None, bootstrap=True, bootstrap_features=False, max_features=1.0, max_samples=1.0, n_estimators=9, n_jobs=1, oob_score=False, random_state=111, verbose=0, warm_start=False)), ('ETC', ExtraTreesClassifier(bootstrap=False, learning_rate=1.0, n_estimators=62, random_state=111))], flatten_transform=None, n_jobs=1, voting='soft', weights=None)
```

```
In [28]: pred = eclf.predict(features_test)
```

```
/Users/alex/anaconda/envs/gc/lib/python3.6/site-packages/sklearn/preprocessing/label.py:151: DataConversionWarning: A column-vector y was passed when a 1d array was expected. Please change the shape of y to (n_samples, ), for example using y = column_or_1d(y, warn=True)
if diff:
```

```
In [29]: print(precision_score(labels_test, pred), recall_score(labels_test, pred), accuracy_score(labels_test, pred))
```

```
0.9893617021276596 0.8303571428571429 0.9760765550239234 0.9029126213592233
```

Better but nope.

0.0.5 RNN

Define the RNN structure.

```
In [30]: from keras.models import Model
         from keras.layers import LSTM, Activation, Dense, Dropout, Input, Embedding
         from keras.optimizers import RMSprop
         from keras.preprocessing.text import Tokenizer
         from keras.preprocessing import sequence
         from keras.utils import to_categorical
         from keras.callbacks import EarlyStopping
         from keras.callbacks import Callback
```

Using TensorFlow backend.

```
In [31]: max_words = features_train.shape[0]
        max_len = features_train.shape[1]
```

```
In [32]: def RNN():
        inputs = Input(name='inputs', shape=[max_len])
        layer = Embedding(max_words, 50, input_length=max_len)(inputs)
        layer = LSTM(64)(layer)
        layer = Dense(256, name='FC1')(layer)
        layer = Activation('relu')(layer)
        layer = Dropout(0.5)(layer)
        layer = Dense(1, name='out_layer')(layer)
        layer = Activation('sigmoid')(layer)
        model = Model(inputs=inputs, outputs=layer)
        return model
```

Call the function and compile the model.

```
In [33]: model = RNN()
        model.summary()
        model.compile(loss='binary_crossentropy', optimizer=RMSprop(), metrics=['accuracy'])
```

Layer (type)	Output Shape	Param #
inputs (InputLayer)	(None, 9403)	0
embedding_1 (Embedding)	(None, 9403, 50)	195000
lstm_1 (LSTM)	(None, 64)	29440
FC1 (Dense)	(None, 256)	16640
activation_1 (Activation)	(None, 256)	0
dropout_1 (Dropout)	(None, 256)	0
out_layer (Dense)	(None, 1)	257
activation_2 (Activation)	(None, 1)	0
Total params: 241,337		
Trainable params: 241,337		
Non-trainable params: 0		

```
In [34]: since = time.time()
```

```
# model.fit(features_train, labels_train, epochs=10, batch_size=128,
```

```

#                                     validation_split=0.2,
#                                     callbacks=[metrics, EarlyStopping(monitor='val_loss',min_delta=

model.fit(features_train,labels_train,batch_size=128,epochs=10,
          validation_split=0.2,callbacks=[EarlyStopping(monitor='val_loss',min_delta=

time_elapsed = time.time() - since

```

Train on 3120 samples, validate on 780 samples

Epoch 1/10

3120/3120 [=====] - 295s 95ms/step - loss: 0.4455 - acc: 0.8465 - val

Epoch 2/10

3120/3120 [=====] - 301s 96ms/step - loss: 0.4044 - acc: 0.8654 - val

Epoch 3/10

3120/3120 [=====] - 357s 114ms/step - loss: 0.4055 - acc: 0.8654 - val

```

In [35]: print('Training complete in {:.0f}m {:.4f}s'.format(
          time_elapsed // 60, time_elapsed % 60))

```

Training complete in 15m 53.9231s

```

In [36]: pred = (np.asarray(model.predict(features_test, batch_size=128))).round()

```

```

In [37]: pred_scores.append(("LSTM", [precision_score(labels_test,pred), recall_score(labels_t

```

/Users/alex/anaconda/envs/gc/lib/python3.6/site-packages/sklearn/metrics/classification.py:113

'precision', 'predicted', average, warn_for)

/Users/alex/anaconda/envs/gc/lib/python3.6/site-packages/sklearn/metrics/classification.py:113

'precision', 'predicted', average, warn_for)

0.0.6 gcForest

```

In [38]: import sys
          sys.path.append("..")
          from gcforest.gcforest import GCForest
          from gcforest.utils.config_utils import load_json

```

```

In [39]: def get_toy_config():
          config = {}
          ca_config = {}
          ca_config["random_state"] = 111
          ca_config["max_layers"] = 10
          ca_config["early_stopping_rounds"] = 3
          ca_config["n_classes"] = 2
          ca_config["estimators"] = []
          ca_config["estimators"].append({"n_folds": 5, "type": "RandomForestClassifier", "n

```



```

ca_config["estimators"].append({"n_folds": 5, "type": "RandomForestClassifier", "1
ca_config["estimators"].append({"n_folds": 5, "type": "MultinomialNB", "alpha": 0
ca_config["estimators"].append({"n_folds": 5, "type": "MultinomialNB", "alpha": 0

config["cascade"] = ca_config
return config

In [40]: config = get_toy_config()
gc = GCForest(config)

# features_train ndarraylabels_train (n_samples, )(n_samples, 1)
features_train = features_train.toarray()
labels_train = labels_train.reshape(-1)

since = time.time()
gc.fit_transform(features_train, labels_train)

time_elapsed = time.time() - since

[ 2019-04-23 21:56:28,200] [cascade_classifier.fit_transform] X_groups_train.shape=[(3900, 9403)
[ 2019-04-23 21:56:28,434] [cascade_classifier.fit_transform] group_dims=[9403]
[ 2019-04-23 21:56:28,435] [cascade_classifier.fit_transform] group_starts=[0]
[ 2019-04-23 21:56:28,436] [cascade_classifier.fit_transform] group_ends=[9403]
[ 2019-04-23 21:56:28,437] [cascade_classifier.fit_transform] X_train.shape=(3900, 9403),X_test
[ 2019-04-23 21:56:28,611] [cascade_classifier.fit_transform] [layer=0] look_indexes=[0], X_cur_1
[ 2019-04-23 21:56:32,297] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_0 - 5_1
[ 2019-04-23 21:56:36,039] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_0 - 5_1
[ 2019-04-23 21:56:39,512] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_0 - 5_1
[ 2019-04-23 21:56:43,081] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_0 - 5_1
[ 2019-04-23 21:56:46,488] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_0 - 5_1
[ 2019-04-23 21:56:46,489] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_0 - 5_1
[ 2019-04-23 21:56:49,923] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_1 - 5_1
[ 2019-04-23 21:56:53,249] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_1 - 5_1
[ 2019-04-23 21:56:56,568] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_1 - 5_1
[ 2019-04-23 21:56:59,982] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_1 - 5_1
[ 2019-04-23 21:57:03,521] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_1 - 5_1
[ 2019-04-23 21:57:03,524] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_1 - 5_1
[ 2019-04-23 21:57:03,752] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_2 - 5_1
[ 2019-04-23 21:57:03,981] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_2 - 5_1
[ 2019-04-23 21:57:04,199] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_2 - 5_1
[ 2019-04-23 21:57:04,421] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_2 - 5_1
[ 2019-04-23 21:57:04,638] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_2 - 5_1
[ 2019-04-23 21:57:04,639] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_2 - 5_1
[ 2019-04-23 21:57:04,856] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_3 - 5_1
[ 2019-04-23 21:57:05,076] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_3 - 5_1
[ 2019-04-23 21:57:05,294] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_3 - 5_1
[ 2019-04-23 21:57:05,506] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_3 - 5_1
[ 2019-04-23 21:57:05,718] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_3 - 5_1

```

```

[ 2019-04-23 21:57:05,719] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_3 - 5_
[ 2019-04-23 21:57:05,720] [cascade_classifier.calc_accuracy] Accuracy(layer_0 - train.classifi
[ 2019-04-23 21:57:05,907] [cascade_classifier.fit_transform] [layer=1] look_indexes=[0], X_cur_1
[ 2019-04-23 21:57:08,068] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_0 - 5_
[ 2019-04-23 21:57:09,964] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_0 - 5_
[ 2019-04-23 21:57:11,994] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_0 - 5_
[ 2019-04-23 21:57:14,049] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_0 - 5_
[ 2019-04-23 21:57:15,923] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_0 - 5_
[ 2019-04-23 21:57:15,925] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_0 - 5_
[ 2019-04-23 21:57:17,921] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_1 - 5_
[ 2019-04-23 21:57:19,903] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_1 - 5_
[ 2019-04-23 21:57:21,985] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_1 - 5_
[ 2019-04-23 21:57:23,975] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_1 - 5_
[ 2019-04-23 21:57:26,152] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_1 - 5_
[ 2019-04-23 21:57:26,153] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_1 - 5_
[ 2019-04-23 21:57:26,391] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_2 - 5_
[ 2019-04-23 21:57:26,626] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_2 - 5_
[ 2019-04-23 21:57:26,865] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_2 - 5_
[ 2019-04-23 21:57:27,093] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_2 - 5_
[ 2019-04-23 21:57:27,323] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_2 - 5_
[ 2019-04-23 21:57:27,325] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_2 - 5_
[ 2019-04-23 21:57:27,574] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_3 - 5_
[ 2019-04-23 21:57:27,942] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_3 - 5_
[ 2019-04-23 21:57:28,278] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_3 - 5_
[ 2019-04-23 21:57:28,636] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_3 - 5_
[ 2019-04-23 21:57:29,042] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_3 - 5_
[ 2019-04-23 21:57:29,045] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_3 - 5_
[ 2019-04-23 21:57:29,047] [cascade_classifier.calc_accuracy] Accuracy(layer_1 - train.classifi
[ 2019-04-23 21:57:29,350] [cascade_classifier.fit_transform] [layer=2] look_indexes=[0], X_cur_1
[ 2019-04-23 21:57:31,499] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_0 - 5_
[ 2019-04-23 21:57:33,600] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_0 - 5_
[ 2019-04-23 21:57:35,678] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_0 - 5_
[ 2019-04-23 21:57:37,568] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_0 - 5_
[ 2019-04-23 21:57:39,449] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_0 - 5_
[ 2019-04-23 21:57:39,450] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_0 - 5_
[ 2019-04-23 21:57:41,338] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_1 - 5_
[ 2019-04-23 21:57:43,218] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_1 - 5_
[ 2019-04-23 21:57:45,210] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_1 - 5_
[ 2019-04-23 21:57:47,093] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_1 - 5_
[ 2019-04-23 21:57:49,077] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_1 - 5_
[ 2019-04-23 21:57:49,078] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_1 - 5_
[ 2019-04-23 21:57:49,296] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_2 - 5_
[ 2019-04-23 21:57:49,516] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_2 - 5_
[ 2019-04-23 21:57:49,736] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_2 - 5_
[ 2019-04-23 21:57:49,953] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_2 - 5_
[ 2019-04-23 21:57:50,172] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_2 - 5_
[ 2019-04-23 21:57:50,173] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_2 - 5_
[ 2019-04-23 21:57:50,395] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_3 - 5_

```

```

[ 2019-04-23 21:57:50,614] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_3 - 5_1
[ 2019-04-23 21:57:50,833] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_3 - 5_1
[ 2019-04-23 21:57:51,055] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_3 - 5_1
[ 2019-04-23 21:57:51,277] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_3 - 5_1
[ 2019-04-23 21:57:51,278] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_3 - 5_1
[ 2019-04-23 21:57:51,279] [cascade_classifier.calc_accuracy] Accuracy(layer_2 - train.classifi
[ 2019-04-23 21:57:51,473] [cascade_classifier.fit_transform] [layer=3] look_indexs=[0], X_cur_1
[ 2019-04-23 21:57:53,451] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_0 - 5_1
[ 2019-04-23 21:57:55,351] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_0 - 5_1
[ 2019-04-23 21:57:57,239] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_0 - 5_1
[ 2019-04-23 21:57:59,318] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_0 - 5_1
[ 2019-04-23 21:58:01,303] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_0 - 5_1
[ 2019-04-23 21:58:01,305] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_0 - 5_1
[ 2019-04-23 21:58:03,398] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_1 - 5_1
[ 2019-04-23 21:58:05,294] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_1 - 5_1
[ 2019-04-23 21:58:07,279] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_1 - 5_1
[ 2019-04-23 21:58:09,261] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_1 - 5_1
[ 2019-04-23 21:58:11,250] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_1 - 5_1
[ 2019-04-23 21:58:11,251] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_1 - 5_1
[ 2019-04-23 21:58:11,479] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_2 - 5_1
[ 2019-04-23 21:58:11,696] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_2 - 5_1
[ 2019-04-23 21:58:11,918] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_2 - 5_1
[ 2019-04-23 21:58:12,134] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_2 - 5_1
[ 2019-04-23 21:58:12,349] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_2 - 5_1
[ 2019-04-23 21:58:12,350] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_2 - 5_1
[ 2019-04-23 21:58:12,569] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_3 - 5_1
[ 2019-04-23 21:58:12,782] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_3 - 5_1
[ 2019-04-23 21:58:12,999] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_3 - 5_1
[ 2019-04-23 21:58:13,213] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_3 - 5_1
[ 2019-04-23 21:58:13,425] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_3 - 5_1
[ 2019-04-23 21:58:13,426] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_3 - 5_1
[ 2019-04-23 21:58:13,427] [cascade_classifier.calc_accuracy] Accuracy(layer_3 - train.classifi
[ 2019-04-23 21:58:13,429] [cascade_classifier.fit_transform] [Result][Optimal Level Detected]

```

```

In [41]: print('Training complete in {:.0f}m {:.4f}s'.format(
           time_elapsed // 60, time_elapsed % 60))

```

Training complete in 1m 45.2865s

```

In [42]: pred = predict_labels(gc,features_test.toarray())
         pred_scores.append(("DCF", [precision_score(labels_test,pred), recall_score(labels_test,

```

```

[ 2019-04-23 21:58:13,530] [cascade_classifier.transform] X_groups_test.shape=[[1672, 9403]]
[ 2019-04-23 21:58:13,650] [cascade_classifier.transform] group_dims=[9403]
[ 2019-04-23 21:58:13,651] [cascade_classifier.transform] X_test.shape=(1672, 9403)
[ 2019-04-23 21:58:13,730] [cascade_classifier.transform] [layer=0] look_indexs=[0], X_cur_test

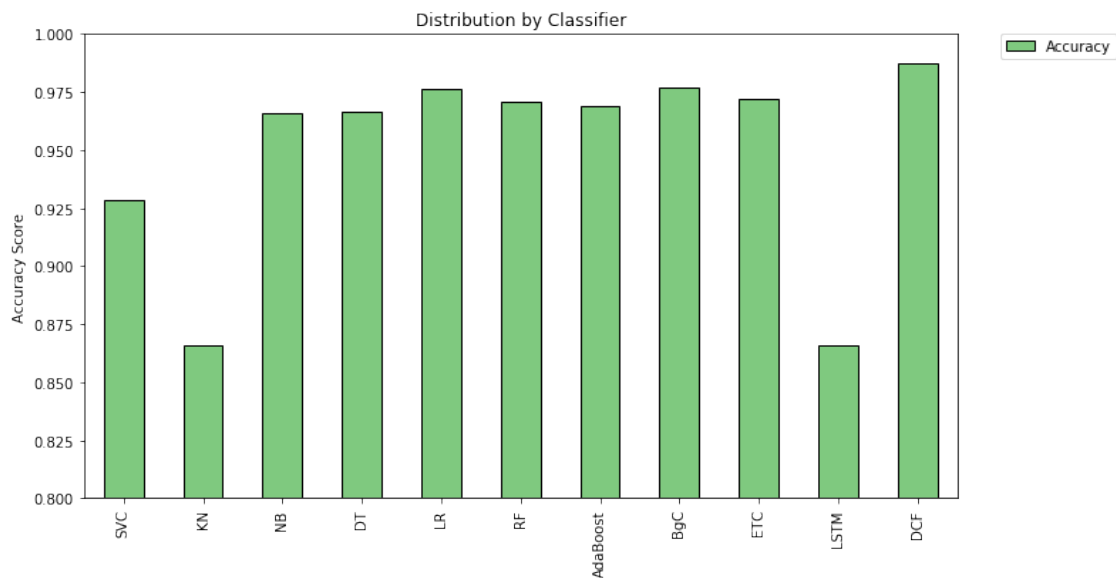
```

```
In [43]: df = pd.DataFrame.from_items(pred_scores,orient='index', columns=['Precision', 'Recall', 'F1', 'Training Time (s)'])
df
```

```
Out[43]:
```

	Precision	Recall	Accuracy	F1	Training Time (s)
SVC	0.747619	0.700893	0.928230	0.723502	0m 0.1977s
KN	0.000000	0.000000	0.866029	0.000000	0m 0.0009s
NB	0.819923	0.955357	0.965909	0.882474	0m 0.0019s
DT	0.900000	0.843750	0.966507	0.870968	0m 0.1379s
LR	0.989362	0.830357	0.976077	0.902913	0m 0.0108s
RF	0.988827	0.790179	0.970694	0.878412	0m 1.2109s
AdaBoost	0.934343	0.825893	0.968900	0.876777	0m 2.3704s
BgC	0.969543	0.852679	0.976675	0.907363	0m 0.8051s
ETC	0.963351	0.821429	0.971890	0.886747	0m 0.7646s
LSTM	0.000000	0.000000	0.866029	0.000000	15m 53.9231s
DCF	0.963470	0.941964	0.987440	0.952596	1m 45.2865s

```
In [44]: df.plot(kind='bar', y="Accuracy", ylim=(0.8,1.0), figsize=(11,6), align='center', color='green')
plt.xticks(np.arange(11), df.index)
plt.ylabel('Accuracy Score')
plt.title('Distribution by Classifier')
plt.legend(bbox_to_anchor=(1.05, 1), loc=2, borderaxespad=0.)
plt.savefig("../img/sms-acc-v2-stop.eps")
plt.show()
```



```
In [45]: import pickle
# dump
with open("../pkl/sms-gc-v2-stop.pkl", "wb") as f:
    pickle.dump(gc, f, pickle.HIGHEST_PROTOCOL)
```

```
# # load  
# with open("../pkl/2018_gc.pkl", "rb") as f:  
#     gc = pickle.load(f)
```

0.0.7 Final verdict - gcForest is your friend in spam detection.

In []: