

sms-spam-v3-stop

April 23, 2019

V3: + Delete the stop words + LSTM uses the Tokenizer.fit_on_texts(data) then Tokenizer.texts_to_sequences(data) to do the preprocessing + Other models uses the TfidfVectorizer to do 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 TfidfVectorizer
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()

```

```

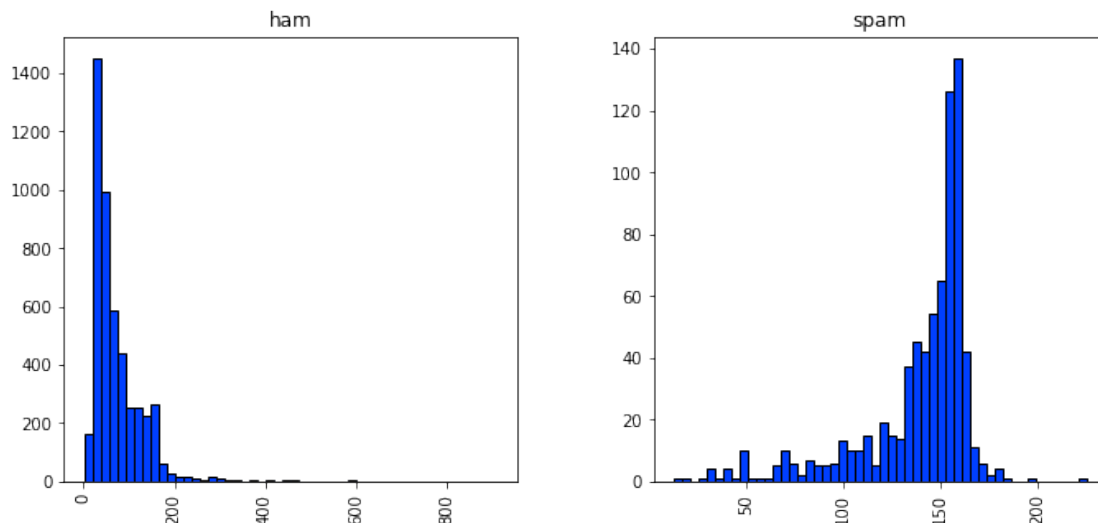
Out[5]:  label      message  length
0   ham  Go until jurong point, crazy.. Available only ...    111
1   ham                      Ok lar... Joking wif u oni...    29
2  spam  Free entry in 2 a wkly comp to win FA Cup fina...   155
3   ham  U dun say so early hor... U c already then say...    49
4   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 = TfidfVectorizer("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:2
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
```

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,
```

```
/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
/Users/alex/anaconda/envs/gc/lib/python3.6/site-packages/ipykernel/__main__.py:2: DataConversi
    from ipykernel import kernelapp as app
/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/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
```

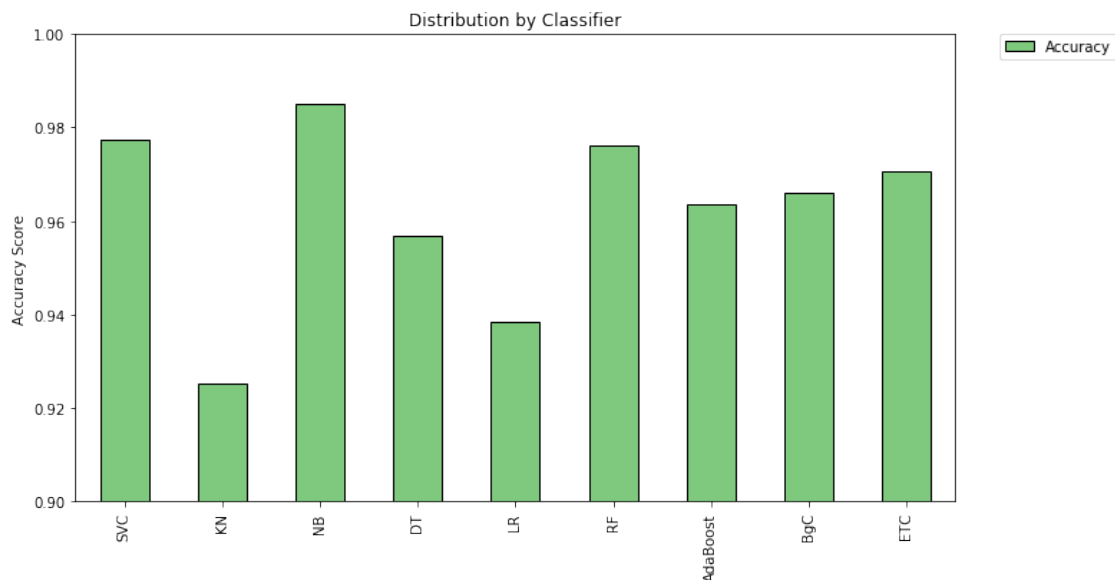
```
In [23]: # pred_scores
```

```
In [24]: df = pd.DataFrame.from_items(pred_scores,orient='index', columns=['Precision', 'Recall', 'Accuracy', 'F1', 'Training Time (s)'])
df
```

```
Out[24]:
```

	Precision	Recall	Accuracy	F1	Training Time (s)
SVC	0.980198	0.853448	0.977273	0.912442	0m 0.3706s
KN	0.990826	0.465517	0.925239	0.633431	0m 0.0014s
NB	0.951965	0.939655	0.985048	0.945770	0m 0.0017s
DT	0.866972	0.814655	0.956938	0.840000	0m 0.1950s
LR	0.900621	0.625000	0.938397	0.737913	0m 0.0085s
RF	1.000000	0.827586	0.976077	0.905660	0m 1.2973s
AdaBoost	0.967213	0.762931	0.963517	0.853012	0m 2.5404s
BgC	0.910798	0.836207	0.965909	0.871910	0m 1.0032s
ETC	1.000000	0.788793	0.970694	0.881928	0m 0.9153s

```
In [25]: df.plot(kind='bar', y="Accuracy", ylim=(0.9,1.0), 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-v3-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 [26]: from sklearn.ensemble import VotingClassifier

In [27]: eclf = VotingClassifier(estimators=[('BgC', bc), ('ETC', etc), ('RF', rfc), ('Ada', adaboost)])

In [28]: eclf.fit(features_train, labels_train)

/Users/alex/anaconda/envs/gc/lib/python3.6/site-packages/sklearn/preprocessing/label.py:95: DataConversionWarning: Data has dtype object, but could be converted to dtype float, which could improve performance.
  y = column_or_1d(y, warn=True)
/Users/alex/anaconda/envs/gc/lib/python3.6/site-packages/sklearn/preprocessing/label.py:128: DataConversionWarning: Data has dtype object, but could be converted to dtype float, which could improve performance.
  y = column_or_1d(y, warn=True)

Out[28]: 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 [29]: pred = eclf.predict(features_test)

/Users/alex/anaconda/envs/gc/lib/python3.6/site-packages/sklearn/preprocessing/label.py:151: DataConversionWarning: Data has dtype object, but could be converted to dtype float, which could improve performance.
  if diff:

In [30]: print(precision_score(labels_test, pred), recall_score(labels_test, pred), accuracy_score(labels_test, pred))

0.9897959183673469 0.8362068965517241 0.9760765550239234 0.9065420560747662

```

Better but nope.

0.0.5 RNN

Define the RNN structure.

```

In [31]: 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.

0.0.6 Process the data

- Tokenize the data and convert the text to sequences.
- Add padding to ensure that all the sequences have the same shape.
- There are many ways of taking the *max_len* and here an arbitrary length of 500 is chosen. (From the Fig, almost all the sentences have the length < 200)

```
In [32]: features_lstm = text_feat
        labels_lstm = labels
```

```
In [33]: max_words = 1000
        max_len = 200 # n_features
        tok = Tokenizer(num_words=max_words)
        tok.fit_on_texts(features_lstm)
        sequences = tok.texts_to_sequences(features_lstm)
        features_lstm = sequence.pad_sequences(sequences,maxlen=max_len)
```

```
In [34]: features_lstm.shape
```

```
Out[34]: (5572, 200)
```

```
In [35]: labels_lstm.shape
```

```
Out[35]: (5572, 1)
```

```
In [36]: features_lstm_train, features_lstm_test, labels_lstm_train, labels_lstm_test = train_test_split(features_lstm, labels_lstm, test_size=0.2, random_state=42)
```

```
In [37]: 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 [38]: model = RNN()
        model.summary()
        model.compile(loss='binary_crossentropy',optimizer=RMSprop(),metrics=['accuracy'])
```

Layer (type)	Output Shape	Param #
inputs (InputLayer)	(None, 200)	0

embedding_1 (Embedding)	(None, 200, 50)	50000

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: 96,337		
Trainable params: 96,337		
Non-trainable params: 0		

```
In [39]: since = time.time()
```

```
model.fit(features_lstm_train, labels_lstm_train, epochs=10, batch_size=128, validation_data=(features_lstm_test, labels_lstm_test),
          callbacks=[EarlyStopping(monitor='val_loss', min_delta=0.0001)])
```

```
time_elapsed = time.time() - since
```

```
Train on 3120 samples, validate on 780 samples
```

```
Epoch 1/10
```

```
3120/3120 [=====] - 7s 2ms/step - loss: 0.3713 - acc: 0.8462 - val_loss: 0.1536
```

```
Epoch 2/10
```

```
3120/3120 [=====] - 7s 2ms/step - loss: 0.1536 - acc: 0.9529 - val_loss: 0.0720
```

```
Epoch 3/10
```

```
3120/3120 [=====] - 8s 2ms/step - loss: 0.0720 - acc: 0.9804 - val_loss: 0.0720
```

```
In [40]: print('Training complete in {:.0f}m {:.4f}s'.format(
          time_elapsed // 60, time_elapsed % 60))
```

```
Training complete in 0m 22.1877s
```

```
In [41]: pred = (np.asarray(model.predict(features_lstm_test, batch_size=128))).round()
```

```
In [42]: pred_scores.append(("LSTM", [precision_score(labels_lstm_test, pred), recall_score(labels_lstm_test, pred)]))
```

0.0.7 gcForest

```
In [43]: import sys
          sys.path.append("../")
```



```

from gcforest.gcforest import GCForest
from gcforest.utils.config_utils import load_json

```

```

In [44]: 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_estimators": 100})
    ca_config["estimators"].append({"n_folds": 5, "type": "MultinomialNB", "alpha": 0.01})

    config["cascade"] = ca_config
    return config

```

```

In [45]: 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

# gc.fit_transform(features_train, labels_train, features_test, labels_test)

```

```

[ 2019-04-23 22:05:32,580][cascade_classifier.fit_transform] X_groups_train.shape=[(3900, 9403)]
[ 2019-04-23 22:05:32,824][cascade_classifier.fit_transform] group_dims=[9403]
[ 2019-04-23 22:05:32,825][cascade_classifier.fit_transform] group_starts=[0]
[ 2019-04-23 22:05:32,826][cascade_classifier.fit_transform] group_ends=[9403]
[ 2019-04-23 22:05:32,827][cascade_classifier.fit_transform] X_train.shape=(3900, 9403),X_test.shape=(1000, 9403)
[ 2019-04-23 22:05:33,012][cascade_classifier.fit_transform] [layer=0] look_indexs=[0], X_cur_train.shape=(3900, 9403)
[ 2019-04-23 22:05:37,035][kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_0 - 5-fold) = 0.999
[ 2019-04-23 22:05:41,197][kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_0 - 5-fold) = 0.999
[ 2019-04-23 22:05:45,674][kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_0 - 5-fold) = 0.999
[ 2019-04-23 22:05:50,323][kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_0 - 5-fold) = 0.999
[ 2019-04-23 22:05:54,284][kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_0 - 5-fold) = 0.999
[ 2019-04-23 22:05:54,285][kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_0 - 5-fold) = 0.999
[ 2019-04-23 22:05:54,575][kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_1 - 5-fold) = 0.999
[ 2019-04-23 22:05:54,885][kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_1 - 5-fold) = 0.999
[ 2019-04-23 22:05:55,182][kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_1 - 5-fold) = 0.999
[ 2019-04-23 22:05:55,417][kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_1 - 5-fold) = 0.999
[ 2019-04-23 22:05:55,648][kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_1 - 5-fold) = 0.999

```

```

[ 2019-04-23 22:05:55,649] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_1 - 5_)
[ 2019-04-23 22:05:55,650] [cascade_classifier.calc_accuracy] Accuracy(layer_0 - train.classifi
[ 2019-04-23 22:05:55,865] [cascade_classifier.fit_transform] [layer=1] look_indexes=[0], X_cur_1
[ 2019-04-23 22:05:59,040] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_0 - 5_)
[ 2019-04-23 22:06:01,692] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_0 - 5_)
[ 2019-04-23 22:06:04,537] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_0 - 5_)
[ 2019-04-23 22:06:07,564] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_0 - 5_)
[ 2019-04-23 22:06:11,201] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_0 - 5_)
[ 2019-04-23 22:06:11,202] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_0 - 5_)
[ 2019-04-23 22:06:11,441] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_1 - 5_)
[ 2019-04-23 22:06:11,682] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_1 - 5_)
[ 2019-04-23 22:06:11,930] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_1 - 5_)
[ 2019-04-23 22:06:12,168] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_1 - 5_)
[ 2019-04-23 22:06:12,430] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_1 - 5_)
[ 2019-04-23 22:06:12,432] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_1 - 5_)
[ 2019-04-23 22:06:12,433] [cascade_classifier.calc_accuracy] Accuracy(layer_1 - train.classifi
[ 2019-04-23 22:06:12,648] [cascade_classifier.fit_transform] [layer=2] look_indexes=[0], X_cur_1
[ 2019-04-23 22:06:14,951] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_0 - 5_)
[ 2019-04-23 22:06:17,663] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_0 - 5_)
[ 2019-04-23 22:06:20,807] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_0 - 5_)
[ 2019-04-23 22:06:23,452] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_0 - 5_)
[ 2019-04-23 22:06:26,517] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_0 - 5_)
[ 2019-04-23 22:06:26,519] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_0 - 5_)
[ 2019-04-23 22:06:26,818] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_1 - 5_)
[ 2019-04-23 22:06:27,091] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_1 - 5_)
[ 2019-04-23 22:06:27,350] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_1 - 5_)
[ 2019-04-23 22:06:27,602] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_1 - 5_)
[ 2019-04-23 22:06:27,926] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_1 - 5_)
[ 2019-04-23 22:06:27,928] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_1 - 5_)
[ 2019-04-23 22:06:27,929] [cascade_classifier.calc_accuracy] Accuracy(layer_2 - train.classifi
[ 2019-04-23 22:06:28,120] [cascade_classifier.fit_transform] [layer=3] look_indexes=[0], X_cur_1
[ 2019-04-23 22:06:30,550] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_0 - 5_)
[ 2019-04-23 22:06:33,377] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_0 - 5_)
[ 2019-04-23 22:06:36,437] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_0 - 5_)
[ 2019-04-23 22:06:40,795] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_0 - 5_)
[ 2019-04-23 22:06:43,870] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_0 - 5_)
[ 2019-04-23 22:06:43,872] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_0 - 5_)
[ 2019-04-23 22:06:44,297] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_1 - 5_)
[ 2019-04-23 22:06:44,625] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_1 - 5_)
[ 2019-04-23 22:06:45,250] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_1 - 5_)
[ 2019-04-23 22:06:45,607] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_1 - 5_)
[ 2019-04-23 22:06:45,968] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_1 - 5_)
[ 2019-04-23 22:06:45,972] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_1 - 5_)
[ 2019-04-23 22:06:45,973] [cascade_classifier.calc_accuracy] Accuracy(layer_3 - train.classifi
[ 2019-04-23 22:06:46,231] [cascade_classifier.fit_transform] [layer=4] look_indexes=[0], X_cur_1
[ 2019-04-23 22:06:49,152] [kfold_wrapper.log_eval_metrics] Accuracy(layer_4 - estimator_0 - 5_)
[ 2019-04-23 22:06:52,179] [kfold_wrapper.log_eval_metrics] Accuracy(layer_4 - estimator_0 - 5_)
[ 2019-04-23 22:06:55,443] [kfold_wrapper.log_eval_metrics] Accuracy(layer_4 - estimator_0 - 5_)

```

```

[ 2019-04-23 22:06:58,325] [kfold_wrapper.log_eval_metrics] Accuracy(layer_4 - estimator_0 - 5_)
[ 2019-04-23 22:07:01,003] [kfold_wrapper.log_eval_metrics] Accuracy(layer_4 - estimator_0 - 5_)
[ 2019-04-23 22:07:01,004] [kfold_wrapper.log_eval_metrics] Accuracy(layer_4 - estimator_0 - 5_)
[ 2019-04-23 22:07:01,228] [kfold_wrapper.log_eval_metrics] Accuracy(layer_4 - estimator_1 - 5_)
[ 2019-04-23 22:07:01,452] [kfold_wrapper.log_eval_metrics] Accuracy(layer_4 - estimator_1 - 5_)
[ 2019-04-23 22:07:01,774] [kfold_wrapper.log_eval_metrics] Accuracy(layer_4 - estimator_1 - 5_)
[ 2019-04-23 22:07:02,006] [kfold_wrapper.log_eval_metrics] Accuracy(layer_4 - estimator_1 - 5_)
[ 2019-04-23 22:07:02,235] [kfold_wrapper.log_eval_metrics] Accuracy(layer_4 - estimator_1 - 5_)
[ 2019-04-23 22:07:02,236] [kfold_wrapper.log_eval_metrics] Accuracy(layer_4 - estimator_1 - 5_)
[ 2019-04-23 22:07:02,237] [cascade_classifier.calc_accuracy] Accuracy(layer_4 - train.classifi
[ 2019-04-23 22:07:02,428] [cascade_classifier.fit_transform] [layer=5] look_indexs=[0], X_cur_1
[ 2019-04-23 22:07:05,651] [kfold_wrapper.log_eval_metrics] Accuracy(layer_5 - estimator_0 - 5_)
[ 2019-04-23 22:07:08,054] [kfold_wrapper.log_eval_metrics] Accuracy(layer_5 - estimator_0 - 5_)
[ 2019-04-23 22:07:10,155] [kfold_wrapper.log_eval_metrics] Accuracy(layer_5 - estimator_0 - 5_)
[ 2019-04-23 22:07:12,466] [kfold_wrapper.log_eval_metrics] Accuracy(layer_5 - estimator_0 - 5_)
[ 2019-04-23 22:07:14,474] [kfold_wrapper.log_eval_metrics] Accuracy(layer_5 - estimator_0 - 5_)
[ 2019-04-23 22:07:14,475] [kfold_wrapper.log_eval_metrics] Accuracy(layer_5 - estimator_0 - 5_)
[ 2019-04-23 22:07:14,686] [kfold_wrapper.log_eval_metrics] Accuracy(layer_5 - estimator_1 - 5_)
[ 2019-04-23 22:07:14,902] [kfold_wrapper.log_eval_metrics] Accuracy(layer_5 - estimator_1 - 5_)
[ 2019-04-23 22:07:15,121] [kfold_wrapper.log_eval_metrics] Accuracy(layer_5 - estimator_1 - 5_)
[ 2019-04-23 22:07:15,341] [kfold_wrapper.log_eval_metrics] Accuracy(layer_5 - estimator_1 - 5_)
[ 2019-04-23 22:07:15,559] [kfold_wrapper.log_eval_metrics] Accuracy(layer_5 - estimator_1 - 5_)
[ 2019-04-23 22:07:15,560] [kfold_wrapper.log_eval_metrics] Accuracy(layer_5 - estimator_1 - 5_)
[ 2019-04-23 22:07:15,562] [cascade_classifier.calc_accuracy] Accuracy(layer_5 - train.classifi
[ 2019-04-23 22:07:15,730] [cascade_classifier.fit_transform] [layer=6] look_indexs=[0], X_cur_1
[ 2019-04-23 22:07:17,632] [kfold_wrapper.log_eval_metrics] Accuracy(layer_6 - estimator_0 - 5_)
[ 2019-04-23 22:07:19,627] [kfold_wrapper.log_eval_metrics] Accuracy(layer_6 - estimator_0 - 5_)
[ 2019-04-23 22:07:21,705] [kfold_wrapper.log_eval_metrics] Accuracy(layer_6 - estimator_0 - 5_)
[ 2019-04-23 22:07:24,003] [kfold_wrapper.log_eval_metrics] Accuracy(layer_6 - estimator_0 - 5_)
[ 2019-04-23 22:07:26,307] [kfold_wrapper.log_eval_metrics] Accuracy(layer_6 - estimator_0 - 5_)
[ 2019-04-23 22:07:26,309] [kfold_wrapper.log_eval_metrics] Accuracy(layer_6 - estimator_0 - 5_)
[ 2019-04-23 22:07:26,529] [kfold_wrapper.log_eval_metrics] Accuracy(layer_6 - estimator_1 - 5_)
[ 2019-04-23 22:07:26,743] [kfold_wrapper.log_eval_metrics] Accuracy(layer_6 - estimator_1 - 5_)
[ 2019-04-23 22:07:26,960] [kfold_wrapper.log_eval_metrics] Accuracy(layer_6 - estimator_1 - 5_)
[ 2019-04-23 22:07:27,176] [kfold_wrapper.log_eval_metrics] Accuracy(layer_6 - estimator_1 - 5_)
[ 2019-04-23 22:07:27,393] [kfold_wrapper.log_eval_metrics] Accuracy(layer_6 - estimator_1 - 5_)
[ 2019-04-23 22:07:27,394] [kfold_wrapper.log_eval_metrics] Accuracy(layer_6 - estimator_1 - 5_)
[ 2019-04-23 22:07:27,395] [cascade_classifier.calc_accuracy] Accuracy(layer_6 - train.classifi
[ 2019-04-23 22:07:27,396] [cascade_classifier.fit_transform] [Result][Optimal Level Detected]

```

```

In [46]: print('Training complete in {:.0f}m {:.4f}s'.format(
           time_elapsed // 60, time_elapsed % 60))

```

Training complete in 1m 54.8653s

```

In [47]: pred = predict_labels(gc,features_test.toarray())
           pred_scores.append(("DCF", [precision_score(labels_test,pred), recall_score(labels_test,

```

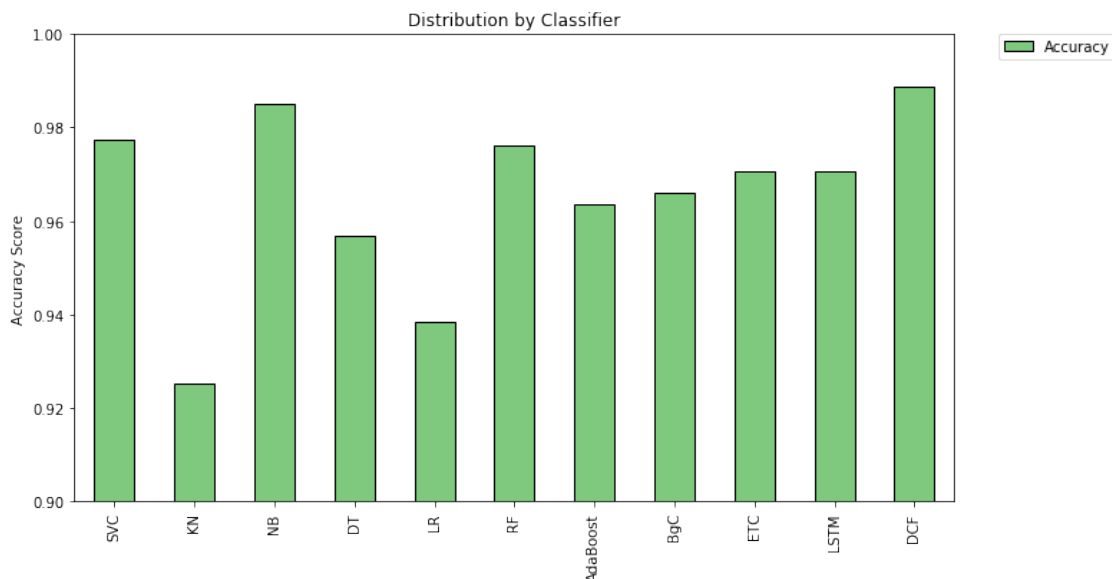
```
[ 2019-04-23 22:07:27,479][cascade_classifier.transform] X_groups_test.shape=[[1672, 9403]]
[ 2019-04-23 22:07:27,578][cascade_classifier.transform] group_dims=[9403]
[ 2019-04-23 22:07:27,579][cascade_classifier.transform] X_test.shape=(1672, 9403)
[ 2019-04-23 22:07:27,642][cascade_classifier.transform] [layer=0] look_indexes=[0], X_cur_test
[ 2019-04-23 22:07:28,498][cascade_classifier.transform] [layer=1] look_indexes=[0], X_cur_test
[ 2019-04-23 22:07:29,283][cascade_classifier.transform] [layer=2] look_indexes=[0], X_cur_test
[ 2019-04-23 22:07:30,073][cascade_classifier.transform] [layer=3] look_indexes=[0], X_cur_test
```

```
In [48]: df = pd.DataFrame.from_items(pred_scores,orient='index', columns=['Precision', 'Recall', 'Accuracy', 'F1', 'Training Time (s)'])
df
```

```
Out[48]:
```

	Precision	Recall	Accuracy	F1	Training Time (s)
SVC	0.980198	0.853448	0.977273	0.912442	0m 0.3706s
KN	0.990826	0.465517	0.925239	0.633431	0m 0.0014s
NB	0.951965	0.939655	0.985048	0.945770	0m 0.0017s
DT	0.866972	0.814655	0.956938	0.840000	0m 0.1950s
LR	0.900621	0.625000	0.938397	0.737913	0m 0.0085s
RF	1.000000	0.827586	0.976077	0.905660	0m 1.2973s
AdaBoost	0.967213	0.762931	0.963517	0.853012	0m 2.5404s
BgC	0.910798	0.836207	0.965909	0.871910	0m 1.0032s
ETC	1.000000	0.788793	0.970694	0.881928	0m 0.9153s
LSTM	1.000000	0.781250	0.970694	0.877193	0m 22.1877s
DCF	0.986301	0.931034	0.988636	0.957871	1m 54.8653s

```
In [49]: df.plot(kind='bar', y="Accuracy", ylim=(0.9,1.0), 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-v3-stop.eps")
plt.show()
```



```
In [50]: import pickle
         # dump
         with open("../pkl/sms-gc-v3-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)
```

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

```
In [ ]:
```