

sms-spam-v2

April 23, 2019

V2: + 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 =)

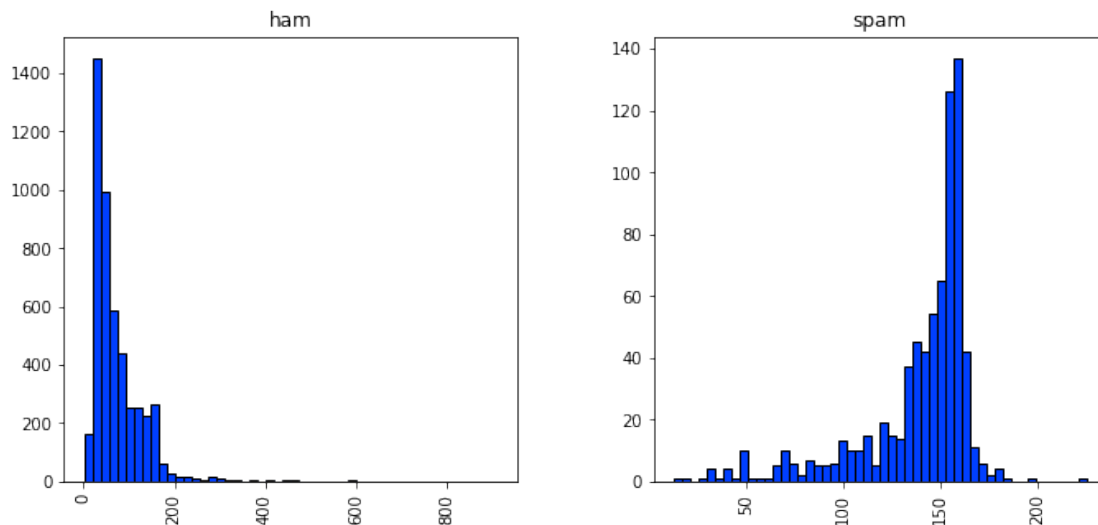
Now let’s create new feature “message length” and plot it to see if it’s of any interest

```
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('en

    #     return " ".join(text)
```

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

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

```
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, 8710)
```

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

```

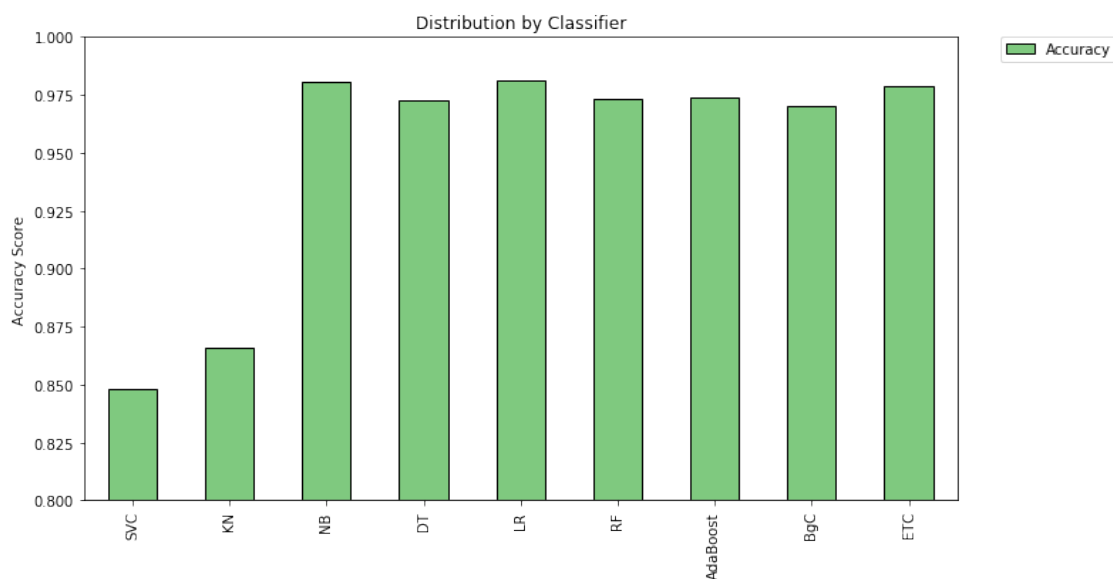
```
In [23]: # pred_scores
```

```
In [24]: df = pd.DataFrame.from_items(pred_scores,orient='index', columns=['Precision', 'Recal
df
```

```
Out[24]:
```

	Precision	Recall	Accuracy	F1	Training Time (s)
SVC	0.427885	0.397321	0.848086	0.412037	0m 0.4724s
KN	0.000000	0.000000	0.866029	0.000000	0m 0.0009s
NB	0.896266	0.964286	0.980263	0.929032	0m 0.0016s
DT	0.927885	0.861607	0.972488	0.893519	0m 0.1232s
LR	0.970588	0.883929	0.980861	0.925234	0m 0.0151s
RF	0.989071	0.808036	0.973086	0.889435	0m 0.8418s
AdaBoost	0.928571	0.870536	0.973684	0.898618	0m 2.2122s
BgC	0.906542	0.866071	0.970096	0.885845	0m 0.7597s
ETC	0.984536	0.852679	0.978469	0.913876	0m 0.5474s

```
In [25]: df.plot(kind='bar', y="Accuracy", ylim=(0.8,1.0), figsize=(11,6), align='center', col
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.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:
  y = column_or_1d(y, warn=True)
/Users/alex/anaconda/envs/gc/lib/python3.6/site-packages/sklearn/preprocessing/label.py:128: DataConversionWarning:
  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:
  if diff:
```

```
In [30]: print(precision_score(labels_test, pred), recall_score(labels_test, pred), accuracy_score(labels_test, pred))
0.9897435897435898 0.8616071428571429 0.9802631578947368 0.9212410501193319
```

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.

```
In [32]: max_words = features_train.shape[0]
        max_len = features_train.shape[1]
```

```
In [33]: 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 [34]: model = RNN()
        model.summary()
        model.compile(loss='binary_crossentropy', optimizer=RMSprop(), metrics=['accuracy'])
```

Layer (type)	Output Shape	Param #
inputs (InputLayer)	(None, 8710)	0
embedding_1 (Embedding)	(None, 8710, 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 [35]: since = time.time()

         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 [=====] - 256s 82ms/step - loss: 0.4340 - acc: 0.8587 - val_
Epoch 2/10
3120/3120 [=====] - 240s 77ms/step - loss: 0.4046 - acc: 0.8654 - val_
Epoch 3/10
3120/3120 [=====] - 242s 78ms/step - loss: 0.4006 - acc: 0.8654 - val_
Epoch 4/10
3120/3120 [=====] - 239s 76ms/step - loss: 0.4020 - acc: 0.8654 - val_
Epoch 5/10
3120/3120 [=====] - 241s 77ms/step - loss: 0.4026 - acc: 0.8654 - val_

In [36]: print('Training complete in {:.0f}m {:.4f}s'.format(
          time_elapsed // 60, time_elapsed % 60))

```

Training complete in 20m 19.1406s

```

In [37]: pred = (np.asarray(model.predict(features_test, batch_size=128))).round()

In [38]: 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 [39]: import sys
         sys.path.append("..")
         from gcforest.gcforest import GCForest
         from gcforest.utils.config_utils import load_json

In [40]: 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": "RandomForestClassifier", "n_estimators": 100})
ca_config["estimators"].append({"n_folds": 5, "type": "MultinomialNB", "alpha": 0.01})
ca_config["estimators"].append({"n_folds": 5, "type": "MultinomialNB", "alpha": 0.01})

config["cascade"] = ca_config
return config

```

```

In [41]: 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 21:26:48,528] [cascade_classifier.fit_transform] X_groups_train.shape=[(3900, 8710), (3900, 8710)]
[ 2019-04-23 21:26:48,733] [cascade_classifier.fit_transform] group_dims=[8710]
[ 2019-04-23 21:26:48,734] [cascade_classifier.fit_transform] group_starts=[0]
[ 2019-04-23 21:26:48,735] [cascade_classifier.fit_transform] group_ends=[8710]
[ 2019-04-23 21:26:48,736] [cascade_classifier.fit_transform] X_train.shape=(3900, 8710),X_test.shape=(3900, 8710)
[ 2019-04-23 21:26:48,905] [cascade_classifier.fit_transform] [layer=0] look_indexes=[0], X_cur_train.shape=(3900, 8710)
[ 2019-04-23 21:26:51,566] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_0 - 5-fold) = 0.999
[ 2019-04-23 21:26:54,048] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_0 - 5-fold) = 0.999
[ 2019-04-23 21:26:56,570] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_0 - 5-fold) = 0.999
[ 2019-04-23 21:26:59,117] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_0 - 5-fold) = 0.999
[ 2019-04-23 21:27:01,705] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_0 - 5-fold) = 0.999
[ 2019-04-23 21:27:01,708] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_0 - 5-fold) = 0.999
[ 2019-04-23 21:27:04,200] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_1 - 5-fold) = 0.999
[ 2019-04-23 21:27:06,701] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_1 - 5-fold) = 0.999
[ 2019-04-23 21:27:09,178] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_1 - 5-fold) = 0.999
[ 2019-04-23 21:27:11,556] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_1 - 5-fold) = 0.999
[ 2019-04-23 21:27:14,051] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_1 - 5-fold) = 0.999
[ 2019-04-23 21:27:14,052] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_1 - 5-fold) = 0.999
[ 2019-04-23 21:27:14,255] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_2 - 5-fold) = 0.999
[ 2019-04-23 21:27:14,458] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_2 - 5-fold) = 0.999
[ 2019-04-23 21:27:14,662] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_2 - 5-fold) = 0.999
[ 2019-04-23 21:27:14,871] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_2 - 5-fold) = 0.999
[ 2019-04-23 21:27:15,074] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_2 - 5-fold) = 0.999
[ 2019-04-23 21:27:15,075] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_2 - 5-fold) = 0.999
[ 2019-04-23 21:27:15,286] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_3 - 5-fold) = 0.999
[ 2019-04-23 21:27:15,497] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_3 - 5-fold) = 0.999

```

```

[ 2019-04-23 21:27:15,704] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_3 - 5_1
[ 2019-04-23 21:27:15,907] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_3 - 5_1
[ 2019-04-23 21:27:16,120] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_3 - 5_1
[ 2019-04-23 21:27:16,121] [kfold_wrapper.log_eval_metrics] Accuracy(layer_0 - estimator_3 - 5_1
[ 2019-04-23 21:27:16,122] [cascade_classifier.calc_accuracy] Accuracy(layer_0 - train.classifi
[ 2019-04-23 21:27:16,299] [cascade_classifier.fit_transform] [layer=1] look_indexs=[0], X_cur
[ 2019-04-23 21:27:17,900] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_0 - 5_1
[ 2019-04-23 21:27:19,368] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_0 - 5_1
[ 2019-04-23 21:27:20,955] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_0 - 5_1
[ 2019-04-23 21:27:22,473] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_0 - 5_1
[ 2019-04-23 21:27:24,048] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_0 - 5_1
[ 2019-04-23 21:27:24,050] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_0 - 5_1
[ 2019-04-23 21:27:25,393] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_1 - 5_1
[ 2019-04-23 21:27:26,958] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_1 - 5_1
[ 2019-04-23 21:27:28,517] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_1 - 5_1
[ 2019-04-23 21:27:29,976] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_1 - 5_1
[ 2019-04-23 21:27:31,435] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_1 - 5_1
[ 2019-04-23 21:27:31,437] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_1 - 5_1
[ 2019-04-23 21:27:31,637] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_2 - 5_1
[ 2019-04-23 21:27:31,842] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_2 - 5_1
[ 2019-04-23 21:27:32,059] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_2 - 5_1
[ 2019-04-23 21:27:32,263] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_2 - 5_1
[ 2019-04-23 21:27:32,465] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_2 - 5_1
[ 2019-04-23 21:27:32,466] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_2 - 5_1
[ 2019-04-23 21:27:32,675] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_3 - 5_1
[ 2019-04-23 21:27:32,878] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_3 - 5_1
[ 2019-04-23 21:27:33,087] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_3 - 5_1
[ 2019-04-23 21:27:33,288] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_3 - 5_1
[ 2019-04-23 21:27:33,490] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_3 - 5_1
[ 2019-04-23 21:27:33,491] [kfold_wrapper.log_eval_metrics] Accuracy(layer_1 - estimator_3 - 5_1
[ 2019-04-23 21:27:33,492] [cascade_classifier.calc_accuracy] Accuracy(layer_1 - train.classifi
[ 2019-04-23 21:27:33,675] [cascade_classifier.fit_transform] [layer=2] look_indexs=[0], X_cur
[ 2019-04-23 21:27:35,241] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_0 - 5_1
[ 2019-04-23 21:27:36,608] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_0 - 5_1
[ 2019-04-23 21:27:37,853] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_0 - 5_1
[ 2019-04-23 21:27:39,119] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_0 - 5_1
[ 2019-04-23 21:27:40,464] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_0 - 5_1
[ 2019-04-23 21:27:40,467] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_0 - 5_1
[ 2019-04-23 21:27:41,826] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_1 - 5_1
[ 2019-04-23 21:27:43,272] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_1 - 5_1
[ 2019-04-23 21:27:44,879] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_1 - 5_1
[ 2019-04-23 21:27:46,158] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_1 - 5_1
[ 2019-04-23 21:27:47,307] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_1 - 5_1
[ 2019-04-23 21:27:47,308] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_1 - 5_1
[ 2019-04-23 21:27:47,510] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_2 - 5_1
[ 2019-04-23 21:27:47,713] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_2 - 5_1
[ 2019-04-23 21:27:47,926] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_2 - 5_1
[ 2019-04-23 21:27:48,140] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_2 - 5_1

```

```

[ 2019-04-23 21:27:48,342] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_2 - 5_1
[ 2019-04-23 21:27:48,344] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_2 - 5_1
[ 2019-04-23 21:27:48,557] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_3 - 5_1
[ 2019-04-23 21:27:48,762] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_3 - 5_1
[ 2019-04-23 21:27:48,971] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_3 - 5_1
[ 2019-04-23 21:27:49,169] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_3 - 5_1
[ 2019-04-23 21:27:49,377] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_3 - 5_1
[ 2019-04-23 21:27:49,378] [kfold_wrapper.log_eval_metrics] Accuracy(layer_2 - estimator_3 - 5_1
[ 2019-04-23 21:27:49,379] [cascade_classifier.calc_accuracy] Accuracy(layer_2 - train.classifi
[ 2019-04-23 21:27:49,557] [cascade_classifier.fit_transform] [layer=3] look_indexs=[0], X_cur
[ 2019-04-23 21:27:50,837] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_0 - 5_1
[ 2019-04-23 21:27:52,095] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_0 - 5_1
[ 2019-04-23 21:27:53,442] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_0 - 5_1
[ 2019-04-23 21:27:54,803] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_0 - 5_1
[ 2019-04-23 21:27:56,057] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_0 - 5_1
[ 2019-04-23 21:27:56,058] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_0 - 5_1
[ 2019-04-23 21:27:57,309] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_1 - 5_1
[ 2019-04-23 21:27:58,655] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_1 - 5_1
[ 2019-04-23 21:28:00,009] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_1 - 5_1
[ 2019-04-23 21:28:01,259] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_1 - 5_1
[ 2019-04-23 21:28:02,600] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_1 - 5_1
[ 2019-04-23 21:28:02,604] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_1 - 5_1
[ 2019-04-23 21:28:02,816] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_2 - 5_1
[ 2019-04-23 21:28:03,020] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_2 - 5_1
[ 2019-04-23 21:28:03,228] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_2 - 5_1
[ 2019-04-23 21:28:03,443] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_2 - 5_1
[ 2019-04-23 21:28:03,645] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_2 - 5_1
[ 2019-04-23 21:28:03,646] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_2 - 5_1
[ 2019-04-23 21:28:03,859] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_3 - 5_1
[ 2019-04-23 21:28:04,063] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_3 - 5_1
[ 2019-04-23 21:28:04,266] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_3 - 5_1
[ 2019-04-23 21:28:04,474] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_3 - 5_1
[ 2019-04-23 21:28:04,677] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_3 - 5_1
[ 2019-04-23 21:28:04,679] [kfold_wrapper.log_eval_metrics] Accuracy(layer_3 - estimator_3 - 5_1
[ 2019-04-23 21:28:04,680] [cascade_classifier.calc_accuracy] Accuracy(layer_3 - train.classifi
[ 2019-04-23 21:28:04,682] [cascade_classifier.fit_transform] [Result] [Optimal Level Detected]

```

```

In [42]: print('Training complete in {:.0f}m {:.4f}s'.format(
           time_elapsed // 60, time_elapsed % 60))

```

Training complete in 1m 16.2139s

```

In [43]: pred = predict_labels(gc,features_test.toarray())
           pred_scores.append(("DCF", [precision_score(labels_test,pred), recall_score(labels_test,

```

```

[ 2019-04-23 21:28:04,785] [cascade_classifier.transform] X_groups_test.shape=[[1672, 8710]]
[ 2019-04-23 21:28:04,884] [cascade_classifier.transform] group_dims=[8710]

```

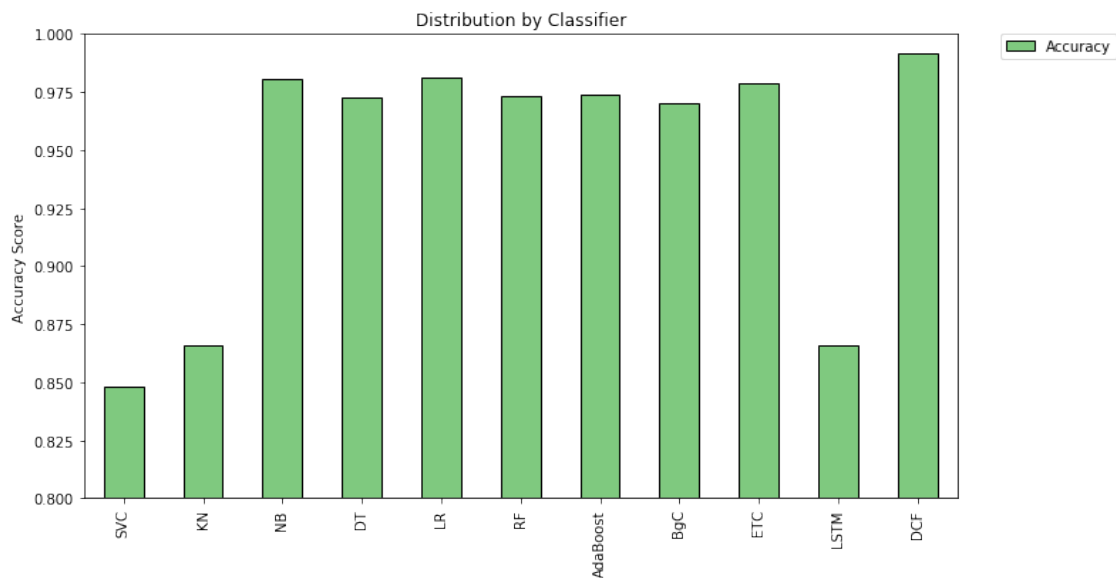
```
[ 2019-04-23 21:28:04,885][cascade_classifier.transform] X_test.shape=(1672, 8710)
[ 2019-04-23 21:28:04,958][cascade_classifier.transform] [layer=0] look_indexes=[0], X_cur_test
```

```
In [44]: df = pd.DataFrame.from_items(pred_scores,orient='index', columns=['Precision', 'Recall', 'Accuracy', 'F1', 'Training Time (s)'])
df
```

```
Out [44]:
```

	Precision	Recall	Accuracy	F1	Training Time (s)
SVC	0.427885	0.397321	0.848086	0.412037	0m 0.4724s
KN	0.000000	0.000000	0.866029	0.000000	0m 0.0009s
NB	0.896266	0.964286	0.980263	0.929032	0m 0.0016s
DT	0.927885	0.861607	0.972488	0.893519	0m 0.1232s
LR	0.970588	0.883929	0.980861	0.925234	0m 0.0151s
RF	0.989071	0.808036	0.973086	0.889435	0m 0.8418s
AdaBoost	0.928571	0.870536	0.973684	0.898618	0m 2.2122s
BgC	0.906542	0.866071	0.970096	0.885845	0m 0.7597s
ETC	0.984536	0.852679	0.978469	0.913876	0m 0.5474s
LSTM	0.000000	0.000000	0.866029	0.000000	20m 19.1406s
DCF	0.981651	0.955357	0.991627	0.968326	1m 16.2139s

```
In [45]: 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.eps")
plt.show()
```



```
In [46]: import pickle
         # dump
         with open("../pkl/sms-gc-v2.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 [ ]:
```

```
In [ ]:
```