

kaggle

Search

Competitions

Datasets

Kernels

Discussion

Learn

...

🏆

Introduction to Ensembling/Stacking in Python

Python notebook using data from [Titanic: Machine Learning from Disaster](#) · 332,704 views · tutorial, xgboost, ensembling

^

3079

Fork8707

...

- Version 95
- [🔗 95 commits](#)
- Notebook
- Data
- Output
- Log
- Comments

## Introduction

This notebook is a very basic and simple introductory primer to the method of ensembling (combining) base learning models, in particular the variant of ensembling known as Stacking. In a nutshell stacking uses as a first-level (base), the predictions of a few basic classifiers and then uses another model at the second-level to predict the output from the earlier first-level predictions.

The Titanic dataset is a prime candidate for introducing this concept as many newcomers to Kaggle start out here. Furthermore even though stacking has been responsible for many a team winning Kaggle competitions there seems to be a dearth of kernels on this topic so I hope this notebook can fill somewhat of that void.

I myself am quite a newcomer to the Kaggle scene as well and the first proper ensembling/stacking script that I managed to chance upon and study was one written in the AllState Severity Claims competition by the great Faron. The material in this notebook borrows heavily from Faron's script although ported to factor in ensembles of classifiers whilst his was ensembles of regressors. Anyway please check out his script here:

Stacking Starter (<https://www.kaggle.com/mmueller/allstate-claims-severity/stacking-starter/run/390867>) : by Faron

Now onto the notebook at hand and I hope that it manages to do justice and convey the concept of ensembling in an intuitive and concise manner. My other standalone Kaggle script (<https://www.kaggle.com/arthurtok/titanic/simple-stacking-with-xgboost-0-808>) which implements exactly the same ensembling steps (albeit with different parameters) discussed below gives a Public LB score of 0.808 which is good enough to get to the top 9% and runs just under 4 minutes. Therefore I am pretty sure there is a lot of room to improve and add on to that script. Anyways please feel free to leave me any comments with regards to how I can improve

In [1]:

```
# Load in our libraries
import pandas as pd
import numpy as np
import re
import sklearn
import xgboost as xgb
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline

import plotly.offline as py
py.init_notebook_mode(connected=True)
import plotly.graph_objs as go
import plotly.tools as tls

import warnings
warnings.filterwarnings('ignore')

# Going to use these 5 base models for the stacking
from sklearn.ensemble import (RandomForestClassifier, AdaBoostClassifier,
                               GradientBoostingClassifier, ExtraTreesClassifier)
from sklearn.svm import SVC
from sklearn.cross_validation import KFold
```

Now we will proceed much like how most kernels in general are structured, and that is to first explore the data on hand, identify possible feature engineering opportunities as well as numerically encode any categorical features.

In [2]:

```
# Load in the train and test datasets
train = pd.read_csv('../input/train.csv')
test = pd.read_csv('../input/test.csv')

# Store our passenger ID for easy access
PassengerId = test['PassengerId']

train.head(3)
```

Out[2]:

|   | PassengerId | Survived | Pclass | Name  | Sex    | Age  | SibSp | Parch | Ticket           |
|---|-------------|----------|--------|---|--------|------|-------|-------|------------------|
| 0 | 1           | 0        | 3      | Braund, Mr. Owen Harris                           | male   | 22.0 | 1     | 0     | A/5 21171        |
| 1 | 2           | 1        | 1      | Cumings, Mrs. John Bradley (Florence Briggs Th... | female | 38.0 | 1     | 0     | PC 17599         |
| 2 | 3           | 1        | 3      | Heikkinen, Miss. Laina                            | female | 26.0 | 0     | 0     | STON/O2. 3101282 |

Well it is no surprise that our task is to somehow extract the information out of the categorical variables

### Feature Engineering

Here, credit must be extended to Sina's very comprehensive and well-thought out notebook for the feature engineering ideas so please check out his work

Titanic Best Working Classifier (<https://www.kaggle.com/sinakhorami/titanic/titanic-best-working-classifier>) : by Sina

In [3]:

```
full_data = [train, test]

# Some features of my own that I have added in
# Gives the length of the name
train['Name_length'] = train['Name'].apply(len)
test['Name_length'] = test['Name'].apply(len)
# Feature that tells whether a passenger had a cabin on the Titanic
train['Has_Cabin'] = train["Cabin"].apply(lambda x: 0 if type(x) == float else 1)
test['Has_Cabin'] = test["Cabin"].apply(lambda x: 0 if type(x) == float else 1)

# Feature engineering steps taken from Sina
# Create new feature FamilySize as a combination of SibSp and Parch
for dataset in full_data:
    dataset['FamilySize'] = dataset['SibSp'] + dataset['Parch'] + 1
# Create new feature IsAlone from FamilySize
```

```

for dataset in full_data:
    dataset['IsAlone'] = 0
    dataset.loc[dataset['FamilySize'] == 1, 'IsAlone'] = 1
# Remove all NULLS in the Embarked column
for dataset in full_data:
    dataset['Embarked'] = dataset['Embarked'].fillna('S')
# Remove all NULLS in the Fare column and create a new feature CategoricalFare
for dataset in full_data:
    dataset['Fare'] = dataset['Fare'].fillna(train['Fare'].median())
train['CategoricalFare'] = pd.qcut(train['Fare'], 4)
# Create a New feature CategoricalAge
for dataset in full_data:
    age_avg = dataset['Age'].mean()
    age_std = dataset['Age'].std()
    age_null_count = dataset['Age'].isnull().sum()
    age_null_random_list = np.random.randint(age_avg - age_std, age_avg + age_std, size=age_null_count)
    dataset['Age'][np.isnan(dataset['Age'])] = age_null_random_list
    dataset['Age'] = dataset['Age'].astype(int)
train['CategoricalAge'] = pd.cut(train['Age'], 5)
# Define function to extract titles from passenger names
def get_title(name):
    title_search = re.search('([A-Za-z]+)\.', name)
    # If the title exists, extract and return it.
    if title_search:
        return title_search.group(1)
    return ""
# Create a new feature Title, containing the titles of passenger names
for dataset in full_data:
    dataset['Title'] = dataset['Name'].apply(get_title)
# Group all non-common titles into one single grouping "Rare"
for dataset in full_data:
    dataset['Title'] = dataset['Title'].replace(['Lady', 'Countess', 'Capt', 'Col', 'Don', 'Dr', 'Major', 'Rev', 'Sir', 'Jonkheer', 'Dona'], 'Rare')

    dataset['Title'] = dataset['Title'].replace('Mlle', 'Miss')
    dataset['Title'] = dataset['Title'].replace('Ms', 'Miss')
    dataset['Title'] = dataset['Title'].replace('Mme', 'Mrs')

for dataset in full_data:
    # Mapping Sex
    dataset['Sex'] = dataset['Sex'].map( {'female': 0, 'male': 1} ).astype(int)

    # Mapping titles
    title_mapping = {"Mr": 1, "Miss": 2, "Mrs": 3, "Master": 4, "Rare": 5}
    dataset['Title'] = dataset['Title'].map(title_mapping)
    dataset['Title'] = dataset['Title'].fillna(0)

    # Mapping Embarked
    dataset['Embarked'] = dataset['Embarked'].map( {'S': 0, 'C': 1, 'Q': 2} ).astype(int)

    # Mapping Fare
    dataset.loc[ dataset['Fare'] <= 7.91, 'Fare'] = 0
    dataset.loc[(dataset['Fare'] > 7.91) & (dataset['Fare'] <= 14.454), 'Fare'] = 1

```

```

dataset.loc[(dataset['Fare'] > 14.454) & (dataset['Fare'] <= 31), 'Fare'] = 2
dataset.loc[ dataset['Fare'] > 31, 'Fare']
= 3
dataset['Fare'] = dataset['Fare'].astype(int)

# Mapping Age
dataset.loc[ dataset['Age'] <= 16, 'Age']
= 0
dataset.loc[(dataset['Age'] > 16) & (dataset['Age'] <= 32), 'Age'] = 1
dataset.loc[(dataset['Age'] > 32) & (dataset['Age'] <= 48), 'Age'] = 2
dataset.loc[(dataset['Age'] > 48) & (dataset['Age'] <= 64), 'Age'] = 3
dataset.loc[ dataset['Age'] > 64, 'Age'] = 4 ;

```

In [4]:

```

# Feature selection
drop_elements = ['PassengerId', 'Name', 'Ticket', 'Cabin', 'SibSp']
train = train.drop(drop_elements, axis = 1)
train = train.drop(['CategoricalAge', 'CategoricalFare'], axis = 1)
test = test.drop(drop_elements, axis = 1)

```

All right so now having cleaned the features and extracted relevant information and dropped the categorical columns our features should now all be numeric, a format suitable to feed into our Machine Learning models. However before we proceed let us generate some simple correlation and distribution plots of our transformed dataset to observe how

## Visualisations

In [5]:

```
train.head(3)
```

Out[5]:

|   | Survived | Pclass | Sex | Age | Parch | Fare | Embarked | Name_length | Has_Cabin | Fare |
|---|----------|--------|-----|-----|-------|------|----------|-------------|-----------|------|
| 0 | 0        | 3      | 1   | 1   | 0     | 0    | 0        | 23          | 0         | 2    |
| 1 | 1        | 1      | 0   | 2   | 0     | 3    | 1        | 51          | 1         | 2    |
| 2 | 1        | 3      | 0   | 1   | 0     | 1    | 0        | 22          | 0         | 1    |

## Pearson Correlation Heatmap

let us generate some correlation plots of the features to see how related one feature is to the next. To do so, we will utilise the Seaborn plotting package which allows us to plot heatmaps very conveniently as follows

In [6]:

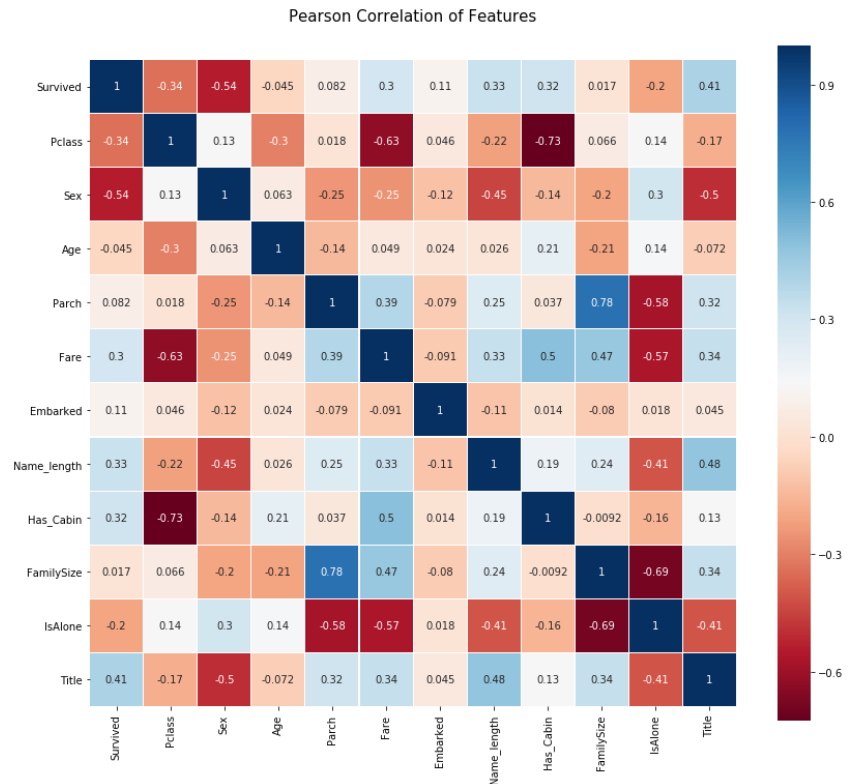
```

colormap = plt.cm.RdBu
plt.figure(figsize=(14,12))
plt.title('Pearson Correlation of Features', y=1.05, size=15)
sns.heatmap(train.astype(float).corr(),linewidths=0.1,vmax=1.0,
             square=True, cmap=colormap, linecolor='white', annot=True)

```

Out[6]:

&lt;matplotlib.axes.\_subplots.AxesSubplot at 0x7f191b84cef0&gt;



### Takeaway from the Plots

One thing that the Pearson Correlation plot can tell us is that there are not too many features strongly correlated with one another. This is good from a point of view of feeding these features into your learning model because this means that there isn't much redundant or superfluous data in our training set and we are happy that each feature carries with it some unique information. Here are two most correlated features are that of Family size and Parch (Parents and Children). I'll still leave both features in for the purposes of this exercise.

### Pairplots

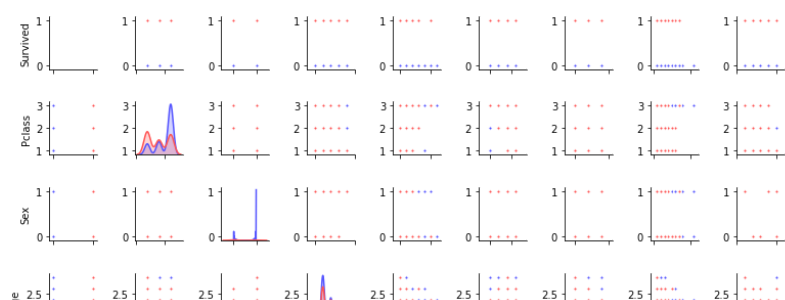
Finally let us generate some pairplots to observe the distribution of data from one feature to the other. Once again we use Seaborn to help us.

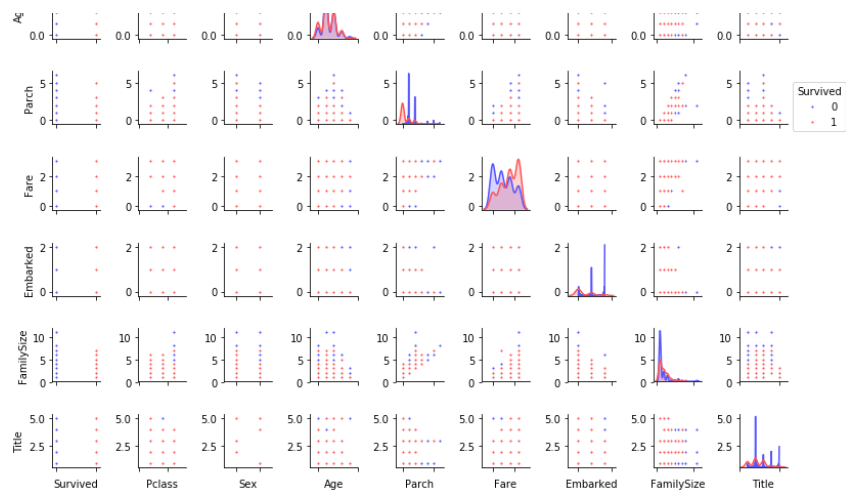
In [7]:

```
g = sns.pairplot(train[['Survived', 'Pclass', 'Sex', 'Age', 'Parch', 'Fare', 'Embarked', 'Name_length', 'Has_Cabin', 'FamilySize', 'IsAlone', 'Title']], hue='Survived', palette='seismic', size=1.2, diag_kind='kde', diag_kws=dict(shade=True), plot_kws=dict(s=10))
g.set(xticklabels=[])
```

Out[7]:

&lt;seaborn.axisgrid.PairGrid at 0x7f191376cba8&gt;





## Ensembling & Stacking models

Finally after that brief whirlwind detour with regards to feature engineering and formatting, we finally arrive at the meat and gist of this notebook.

Creating a Stacking ensemble!

### Helpers via Python Classes

Here we invoke the use of Python's classes to help make it more convenient for us. For any newcomers to programming, one normally hears Classes being used in conjunction with Object-Oriented Programming (OOP). In short, a class helps to extend some code/program for creating objects (variables for old-school peeps) as well as to implement functions and methods specific to that class.

In the section of code below, we essentially write a class *SklearnHelper* that allows one to extend the inbuilt methods (such as train, predict and fit) common to all the Sklearn classifiers. Therefore this cuts out redundancy as won't need to write the same methods five times if we wanted to invoke five different classifiers.

```
In [8]:
# Some useful parameters which will come in handy later on
ntrain = train.shape[0]
ntest = test.shape[0]
SEED = 0 # for reproducibility
NFOLDS = 5 # set folds for out-of-fold prediction
kf = KFold(ntrain, n_folds= NFOLDS, random_state=SEED)

# Class to extend the Sklearn classifier
class SklearnHelper(object):
    def __init__(self, clf, seed=0, params=None):
        params['random_state'] = seed
        self.clf = clf(**params)

    def train(self, x_train, y_train):
        self.clf.fit(x_train, y_train)

    def predict(self, x):
        return self.clf.predict(x)

    def fit(self, x, y):
        return self.clf.fit(x, y)

    def feature_importances(self, x, y):
        print(self.clf.fit(x, y).feature_importances_)
```

```
# Class to extend XGboost classifier
```

Bear with me for those who already know this but for people who have not created classes or objects in Python before, let me explain what the code given above does. In creating my base classifiers, I will only use the models already present in the Sklearn library and therefore only extend the class for that.

**def init** : Python standard for invoking the default constructor for the class. This means that when you want to create an object (classifier), you have to give it the parameters of clf (what sklearn classifier you want), seed (random seed) and params (parameters for the classifiers).

The rest of the code are simply methods of the class which simply call the corresponding methods already existing within the sklearn classifiers. Essentially, we have created a wrapper class to extend the various Sklearn classifiers so that this should help us reduce having to write the same code over and over when we implement multiple learners to our stacker.

### Out-of-Fold Predictions

Now as alluded to above in the introductory section, stacking uses predictions of base classifiers as input for training to a second-level model. However one cannot simply train the base models on the full training data, generate predictions on the full test set and then output these for the second-level training. This runs the risk of your base model predictions already having "seen" the test set and therefore overfitting when feeding these predictions.

```
In [9]:
def get_oof(clf, x_train, y_train, x_test):
    oof_train = np.zeros((ntrain,))
    oof_test = np.zeros((ntest,))
    oof_test_skf = np.empty((NFOLDS, ntest))

    for i, (train_index, test_index) in enumerate(kf):
        x_tr = x_train[train_index]
        y_tr = y_train[train_index]
        x_te = x_train[test_index]

        clf.train(x_tr, y_tr)

        oof_train[test_index] = clf.predict(x_te)
        oof_test_skf[i, :] = clf.predict(x_test)

    oof_test[:] = oof_test_skf.mean(axis=0)
    return oof_train.reshape(-1, 1), oof_test.reshape(-1, 1)
```

## Generating our Base First-Level Models

So now let us prepare five learning models as our first level classification. These models can all be conveniently invoked via the Sklearn library and are listed as follows:

1. Random Forest classifier
2. Extra Trees classifier
3. AdaBoost classifier
4. Gradient Boosting classifier
5. Support Vector Machine

### Parameters

Just a quick summary of the parameters that we will be listing here for completeness,



**n\_jobs** : Number of cores used for the training process. If set to -1, all cores are used.

**n\_estimators** : Number of classification trees in your learning model ( set to 10 per default)

**max\_depth** : Maximum depth of tree, or how much a node should be expanded. Beware if set to too high a number would run the risk of overfitting as one would be growing the tree too deep

**verbose** : Controls whether you want to output any text during the learning process. A value of 0 suppresses all text while a value of 3 outputs the tree learning process at every iteration.

Please check out the full description via the official Sklearn website. There you will find that there are a whole host of other useful parameters that you can play around with.

```
In [10]:
# Put in our parameters for said classifiers
# Random Forest parameters
rf_params = {
    'n_jobs': -1,
    'n_estimators': 500,
    'warm_start': True,
    #'max_features': 0.2,
    'max_depth': 6,
    'min_samples_leaf': 2,
    'max_features' : 'sqrt',
    'verbose': 0
}

# Extra Trees Parameters
et_params = {
    'n_jobs': -1,
    'n_estimators':500,
    #'max_features': 0.5,
    'max_depth': 8,
    'min_samples_leaf': 2,
    'verbose': 0
}

# AdaBoost parameters
ada_params = {
    'n_estimators': 500,
    'learning_rate' : 0.75
}

# Gradient Boosting parameters
gb_params = {
    'n_estimators': 500,
    #'max_features': 0.2,
    'max_depth': 5,
    'min_samples_leaf': 2,
    'verbose': 0
}

# Support Vector Classifier parameters
svc_params = {
    'kernel' : 'linear',
    'C' : 0.025
}
```

Furthermore, since having mentioned about Objects and classes within the OOP framework, let us now create 5 objects that represent our 5 learning models via our Helper Sklearn Class we defined earlier.

```
In [11]:
# Create 5 objects that represent our 4 models
rf = SklearnHelper(clf=RandomForestClassifier, seed=SEED, params=r
f_params)
et = SklearnHelper(clf=ExtraTreesClassifier, seed=SEED, params=et_
params)
ada = SklearnHelper(clf=AdaBoostClassifier, seed=SEED, params=ada_
params)
gb = SklearnHelper(clf=GradientBoostingClassifier, seed=SEED, para
ms=gb_params)
svc = SklearnHelper(clf=SVC, seed=SEED, params=svc_params)
```

### Creating NumPy arrays out of our train and test sets

Great. Having prepared our first layer base models as such, we can now ready the training and test test data for input into our classifiers by generating NumPy arrays out of their original dataframes as follows:

```
In [12]:
# Create Numpy arrays of train, test and target ( Survived) datafram
es to feed into our models
y_train = train['Survived'].ravel()
train = train.drop(['Survived'], axis=1)
x_train = train.values # Creates an array of the train data
x_test = test.values # Creates an array of the test data
```

### Output of the First level Predictions

We now feed the training and test data into our 5 base classifiers and use the Out-of-Fold prediction function we defined earlier to generate our first level predictions. Allow a handful of minutes for the chunk of code below to run.

```
In [13]:
# Create our OOF train and test predictions. These base results will
be used as new features
et_oof_train, et_oof_test = get_oof(et, x_train, y_train, x_test)
# Extra Trees
rf_oof_train, rf_oof_test = get_oof(rf,x_train, y_train, x_test) #
Random Forest
ada_oof_train, ada_oof_test = get_oof(ada, x_train, y_train, x_tes
t) # AdaBoost
gb_oof_train, gb_oof_test = get_oof(gb,x_train, y_train, x_test) #
Gradient Boost
svc_oof_train, svc_oof_test = get_oof(svc,x_train, y_train, x_test
) # Support Vector Classifier

print("Training is complete")
```

Training is complete

### Feature importances generated from the different classifiers

Now having learned our the first-level classifiers, we can utilise a very nifty feature of the Sklearn models and that is to output the importances of the various features in the training and test sets with one very simple line of code.

As per the Sklearn documentation, most of the classifiers are built in with an attribute which returns feature importances by simply typing in **.featureimportances**. Therefore we will invoke this very useful attribute via our function earliand plot the feature importances as such

In [14]:

```

rf_feature = rf.feature_importances(x_train,y_train)
et_feature = et.feature_importances(x_train, y_train)
ada_feature = ada.feature_importances(x_train, y_train)
gb_feature = gb.feature_importances(x_train,y_train)

[ 0.12512537  0.20195675  0.03187994  0.02117736  0.0720603  0.02
351993
    0.10877493  0.06461801  0.06974652  0.01355575  0.26758514]
[ 0.12082488  0.37460384  0.02701211  0.01713043  0.05593931  0.02
854512
    0.04782111  0.08303969  0.04500919  0.02174382  0.1783305 ]
[ 0.028  0.012  0.02  0.062  0.04  0.01  0.69  0.014  0.05
0.004
    0.07 ]
[ 0.07626914  0.03373915  0.10207353  0.03738547  0.10223908  0.04
940839
    0.39897827  0.01836958  0.07035654  0.02056481  0.09061604]

```

So I have not yet figured out how to assign and store the feature importances outright. Therefore I'll print out the values from the code above and then simply copy and paste into Python lists as below (sorry for the lousy hack)

In [15]:

```

rf_features = [0.10474135, 0.21837029, 0.04432652, 0.02249159,
0.05432591, 0.02854371
    ,0.07570305, 0.01088129 , 0.24247496, 0.13685733 , 0.06128402]
et_features = [ 0.12165657, 0.37098307 ,0.03129623 , 0.01591611
    , 0.05525811 , 0.028157
    ,0.04589793 , 0.02030357 , 0.17289562 , 0.04853517, 0.08910063]
ada_features = [0.028 , 0.008 , 0.012 , 0.05866667,
0.032 , 0.008
    ,0.04666667 , 0. , 0.05733333, 0.73866667, 0.01066
667]
gb_features = [ 0.06796144 , 0.03889349 , 0.07237845 , 0.02628645
    , 0.11194395, 0.04778854
    ,0.05965792 , 0.02774745, 0.07462718, 0.4593142 , 0.01340093]

```

Create a dataframe from the lists containing the feature importance data for easy plotting via the Plotly package.

In [16]:

```

cols = train.columns.values
# Create a dataframe with features
feature_dataframe = pd.DataFrame( {'features': cols,
    'Random Forest feature importances': rf_features,
    'Extra Trees feature importances': et_features,
    'AdaBoost feature importances': ada_features,
    'Gradient Boost feature importances': gb_features
    })

```

### Interactive feature importances via Plotly scatterplots

I'll use the interactive Plotly package at this juncture to visualise the feature importances values of the different classifiers via a plotly scatter plot by calling "Scatter" as follows:

In [17]:

```

# Scatter plot
trace = go.Scatter(
    y = feature_dataframe['Random Forest feature importances'].values,
    x = feature_dataframe['features'].values,
    mode='markers',
    marker=dict(
        sizemode = 'diameter',
        sizeref = 1,
        size = 25,
        # size= feature_dataframe['AdaBoost feature importances'].values,
        #color = np.random.randn(500), #set color equal to a variable
        color = feature_dataframe['Random Forest feature importances'].values,
        colorscale='Portland',
        showscale=True
    ),
    text = feature_dataframe['features'].values
)
data = [trace]

layout= go.Layout(
    autosize= True,
    title= 'Random Forest Feature Importance',
    hovermode= 'closest',
    # xaxis= dict(
    #     title= 'Pop',
    #     ticklen= 5,
    #     zeroline= False,
    #     gridwidth= 2,
    # ),
    yaxis=dict(
        title= 'Feature Importance',
        ticklen= 5,
        gridwidth= 2
    ),
    showlegend= False
)
fig = go.Figure(data=data, layout=layout)
py.iplot(fig, filename='scatter2010')

# Scatter plot
trace = go.Scatter(
    y = feature_dataframe['Extra Trees feature importances'].values,
    x = feature_dataframe['features'].values,
    mode='markers',
    marker=dict(
        sizemode = 'diameter',
        sizeref = 1,
        size = 25,
        # size= feature_dataframe['AdaBoost feature importances'].values,
        #color = np.random.randn(500), #set color equal to a variable
        color = feature_dataframe['Extra Trees feature importances'].values,
        colorscale='Portland',
        showscale=True
    ),
    text = feature_dataframe['features'].values
)

```

```

text = feature_dataframe[ 'features' ].values
)
data = [trace]

layout= go.Layout(
    autosize= True,
    title= 'Extra Trees Feature Importance',
    hovermode= 'closest',
    #    xaxis= dict(
    #        title= 'Pop',
    #        ticklen= 5,
    #        zeroline= False,
    #        gridwidth= 2,
    #    ),
    yaxis=dict(
        title= 'Feature Importance',
        ticklen= 5,
        gridwidth= 2
    ),
    showlegend= False
)
fig = go.Figure(data=data, layout=layout)
py.iplot(fig, filename='scatter2010')

# Scatter plot
trace = go.Scatter(
    y = feature_dataframe['AdaBoost feature importances'].values,
    x = feature_dataframe['features'].values,
    mode='markers',
    marker=dict(
        sizemode = 'diameter',
        sizeref = 1,
        size = 25,
    #    size= feature_dataframe['AdaBoost feature importances'].valu
    es,
        #color = np.random.randn(500), #set color equal to a variabl
    e
        color = feature_dataframe['AdaBoost feature importances'].
    values,
        colorscale='Portland',
        showscale=True
    ),
    text = feature_dataframe['features'].values
)
data = [trace]

layout= go.Layout(
    autosize= True,
    title= 'AdaBoost Feature Importance',
    hovermode= 'closest',
    #    xaxis= dict(
    #        title= 'Pop',
    #        ticklen= 5,
    #        zeroline= False,
    #        gridwidth= 2,
    #    ),
    yaxis=dict(
        title= 'Feature Importance',
        ticklen= 5,
        gridwidth= 2
    ),
    showlegend= False
)
fig = go.Figure(data=data, layout=layout)

```

```

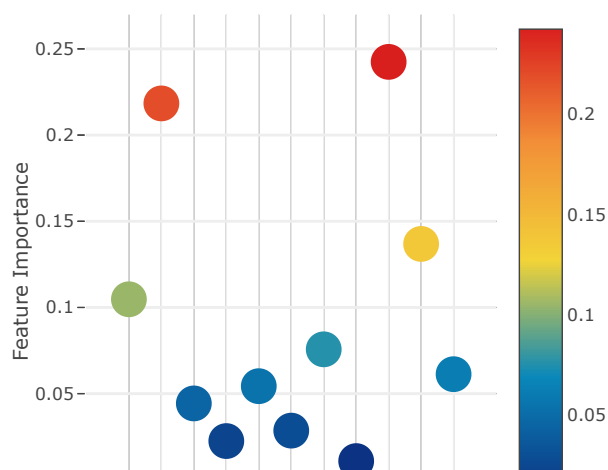
py.iplot(fig, filename='scatter2010')

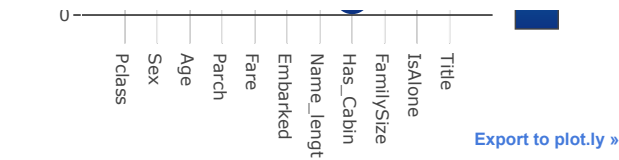
# Scatter plot
trace = go.Scatter(
    y = feature_dataframe['Gradient Boost feature importances'].values,
    x = feature_dataframe['features'].values,
    mode='markers',
    marker=dict(
        sizemode = 'diameter',
        sizeref = 1,
        size = 25,
        # size= feature_dataframe['AdaBoost feature importances'].values,
        #color = np.random.randn(500), #set color equal to a variable
        color = feature_dataframe['Gradient Boost feature importances'].values,
        colorscale='Portland',
        showscale=True
    ),
    text = feature_dataframe['features'].values
)
data = [trace]

layout= go.Layout(
    autosize= True,
    title= 'Gradient Boosting Feature Importance',
    hovermode= 'closest',
    # xaxis= dict(
    #     title= 'Pop',
    #     ticklen= 5,
    #     zeroline= False,
    #     gridwidth= 2,
    # ),
    yaxis=dict(
        title= 'Feature Importance',
        ticklen= 5,
        gridwidth= 2
    ),
    showlegend= False
)
fig = go.Figure(data=data, layout=layout)
py.iplot(fig, filename='scatter2010')

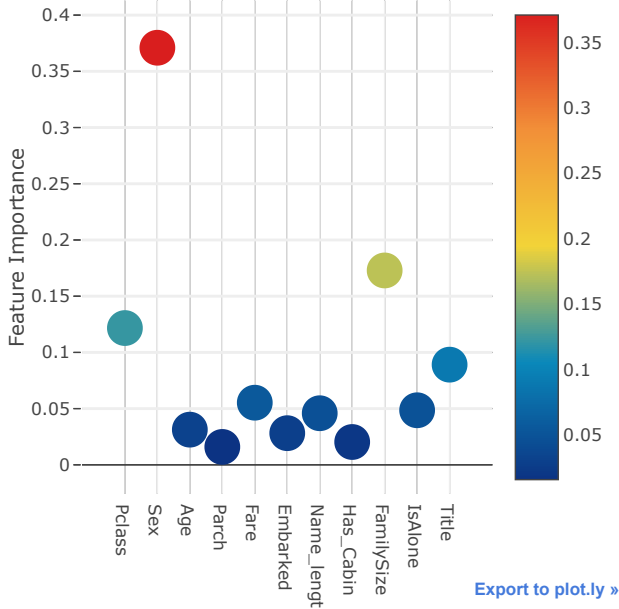
```

Random Forest Feature Importance

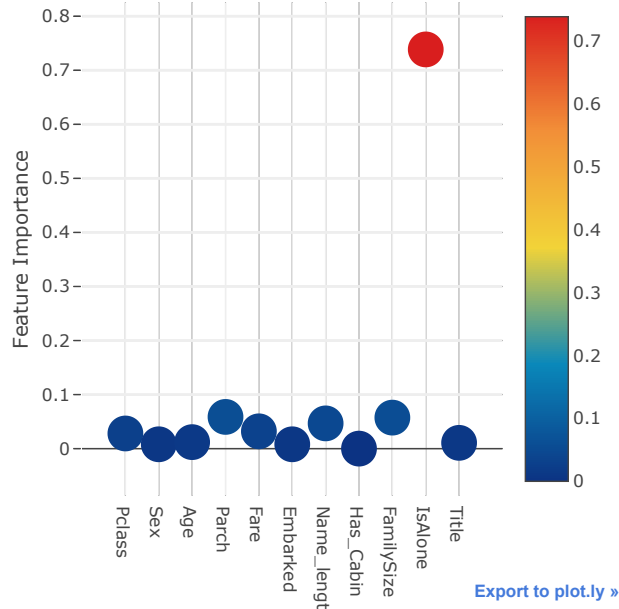




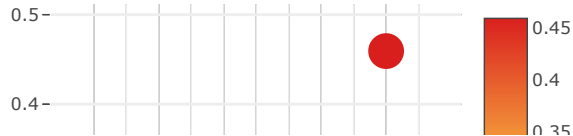
Extra Trees Feature Importance

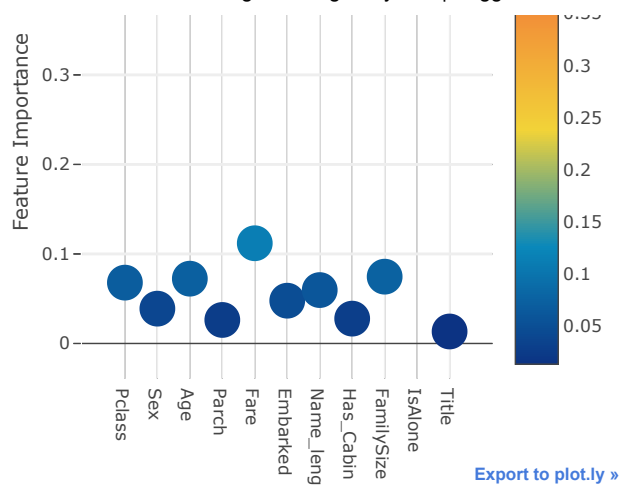


AdaBoost Feature Importance



Gradient Boosting Feature Importance





Now let us calculate the mean of all the feature importances and store it as a new column in the feature importance dataframe.

In [18]:

```
# Create the new column containing the average of values

feature_dataframe['mean'] = feature_dataframe.mean(axis= 1) # axis
= 1 computes the mean row-wise
feature_dataframe.head(3)
```

Out[18]:

|   | features | Random Forest<br>feature<br>importances | Extra Trees<br>feature<br>importances | AdaBoost<br>feature<br>importances | Gradient Boost<br>feature<br>importances | mean     |
|---|----------|---|---------------------------------------|------------------------------------|--|----------|
| 0 | Pclass   | 0.104741                                | 0.121657                              | 0.028                              | 0.067961                                 | 0.080590 |
| 1 | Sex      | 0.218370                                | 0.370983                              | 0.008                              | 0.038893                                 | 0.159062 |
| 2 | Age      | 0.044327                                | 0.031296                              | 0.012                              | 0.072378                                 | 0.040000 |

### Plotly Barplot of Average Feature Importances

Having obtained the mean feature importance across all our classifiers, we can plot them into a Plotly bar plot as follows:

In [19]:

```
y = feature_dataframe['mean'].values
x = feature_dataframe['features'].values
data = [go.Bar(
    x= x,
    y= y,
    width = 0.5,
    marker=dict(
        color = feature_dataframe['mean'].values,
        colorscale='Portland',
        showscale=True,
        reversescale = False
    ),
    opacity=0.6
)]

layout= go.Layout(
    autosize= True,
    title= 'Barplots of Mean Feature Importance'
```

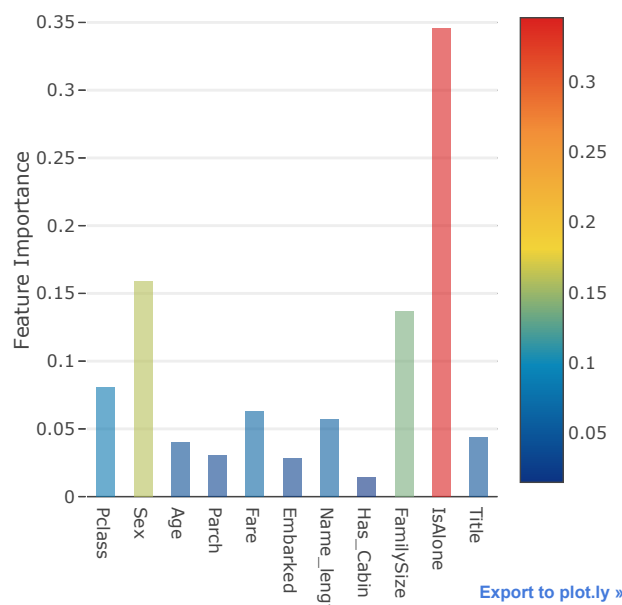


```

# Barplots of Mean Feature Importance
hovermode= 'closest',
# xaxis= dict(
#     title= 'Pop',
#     ticklen= 5,
#     zeroline= False,
#     gridwidth= 2,
# ),
yaxis=dict(
    title= 'Feature Importance',
    ticklen= 5,
    gridwidth= 2
),
showlegend= False
)
fig = go.Figure(data=data, layout=layout)
py.iplot(fig, filename='bar-direct-labels')

```

Barplots of Mean Feature Importance



## Second-Level Predictions from the First-level Output

### First-level output as new features

Having now obtained our first-level predictions, one can think of it as essentially building a new set of features to be used as training data for the next classifier. As per the code below, we are therefore having as our new columns the first-level predictions from our earlier classifiers and we train the next classifier on this.

```

In [20]:
base_predictions_train = pd.DataFrame( {'RandomForest': rf_oof_train.ravel(),
    'ExtraTrees': et_oof_train.ravel(),
    'AdaBoost': ada_oof_train.ravel(),
    'GradientBoost': gb_oof_train.ravel()
    })
base_predictions_train.head()

```

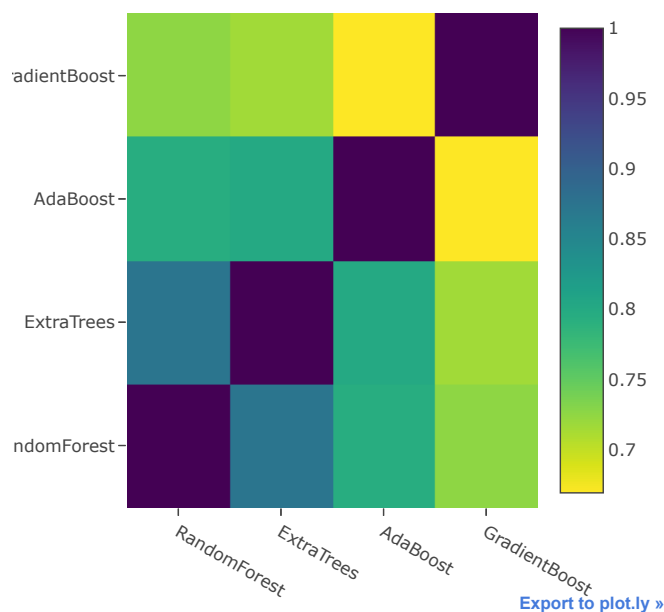
Out[20]:

|   | RandomForest | ExtraTrees | AdaBoost | GradientBoost |
|---|--------------|------------|----------|---------------|
| 0 | 0.0          | 0.0        | 0.0      | 0.0           |
| 1 | 1.0          | 1.0        | 1.0      | 1.0           |
| 2 | 1.0          | 0.0        | 1.0      | 1.0           |
| 3 | 1.0          | 1.0        | 1.0      | 1.0           |
| 4 | 0.0          | 0.0        | 0.0      | 0.0           |

### Correlation Heatmap of the Second Level Training set

In [21]:

```
data = [
    go.Heatmap(
        z= base_predictions_train.astype(float).corr().values ,
        x=base_predictions_train.columns.values,
        y= base_predictions_train.columns.values,
        colorscale='Viridis',
        showscale=True,
        reversescale = True
    )
]
py.iplot(data, filename='labelled-heatmap')
```



There have been quite a few articles and Kaggle competition winner stories about the merits of having trained models that are more uncorrelated with one another producing better scores.

In [22]:

```
x_train = np.concatenate(( et_oof_train, rf_oof_train, ada_oof_train,
                             gb_oof_train, svc_oof_train), axis=1)
x_test = np.concatenate(( et_oof_test, rf_oof_test, ada_oof_test,
                             gb_oof_test, svc_oof_test), axis=1)
```

```
g_train_test, y_train_test, y_test =
```

Having now concatenated and joined both the first-level train and test predictions as `x_train` and `x_test`, we can now fit a second-level learning model.

### Second level learning model via XGBoost

Here we choose the eXtremely famous library for boosted tree learning model, XGBoost. It was built to optimize large-scale boosted tree algorithms. For further information about the algorithm, check out the official documentation (<https://xgboost.readthedocs.io/en/latest/>).

Anyways, we call an `XGBClassifier` and fit it to the first-level train and target data and use the learned model to predict the test data as follows:

```
In [23]: gbm = xgb.XGBClassifier(
          #learning_rate = 0.02,
          n_estimators= 2000,
          max_depth= 4,
          min_child_weight= 2,
          #gamma=1,
          gamma=0.9,
          subsample=0.8,
          colsample_bytree=0.8,
          objective= 'binary:logistic',
          nthread= -1,
          scale_pos_weight=1).fit(x_train, y_train)
          predictions = gbm.predict(x_test)
```

Just a quick run down of the XGBoost parameters used in the model:

**max\_depth** : How deep you want to grow your tree. Beware if set to too high a number might run the risk of overfitting.

**gamma** : minimum loss reduction required to make a further partition on a leaf node of the tree. The larger, the more conservative the algorithm will be.

**eta** : step size shrinkage used in each boosting step to prevent overfitting

### Producing the Submission file

Finally having trained and fit all our first-level and second-level models, we can now output the predictions into the proper format for submission to the Titanic competition as follows:

```
In [24]: # Generate Submission File
          StackingSubmission = pd.DataFrame({ 'PassengerId': PassengerId,
                                              'Survived': predictions })
          StackingSubmission.to_csv("StackingSubmission.csv", index=False)
```

### Steps for Further Improvement

As a closing remark it must be noted that the steps taken above just show a very simple way of producing an ensemble stacker. You hear of ensembles created at the highest level of Kaggle competitions which involves monstrous combinations of stacked classifiers as well as levels of stacking which go to more than 2 levels.

Some additional steps that may be taken to improve one's score could be:

1. Implementing a good cross-validation strategy in training the models to find optimal parameter values
2. Introduce a greater variety of base models for learning. The more uncorrelated the results, the better the final score.

## Conclusion

I have this notebook has been helpful somewhat in introducing a working script for stacking learning models. Again credit must be extended to Faron and Sina.

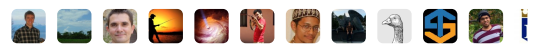
For other excellent material on stacking or ensembling in general, refer to the de-facto Must read article on the website MLWave: Kaggle Ensembling Guide (<http://mlwave.com/kaggle-ensembling-guide/>).

Till next time, Peace Out

This kernel has been released under the [Apache 2.0](#) open source license.

Did you find this Kernel useful?  
Show your appreciation with an upvote


3079



## Data

### Data Sources

- ▼ 🏆 Titanic: Machine Lear...
  - gender\_... 418 x 2
  - test.csv 418 x 11
  - train.csv 891 x 12



### Titanic: Machine Learning from Disaster

Start here! Predict survival on the Titanic and get familiar with ML basics

Last Updated: 7 years ago

#### About this Competition

#### Overview

The data has been split into two groups:

- training set (train.csv)
- test set (test.csv)

**The training set** should be used to build your machine learning models. For the training set, we provide the outcome (also known as the "ground truth") for each passenger. Your model will be based on "features" like passengers' gender and class. You can also use feature engineering to create new features.

**The test set** should be used to see how well your model performs on unseen data. For the test set, we do not provide the ground truth for each passenger. It is your job to predict these outcomes. For each passenger in the test set, use the model you trained to predict whether or not they survived the sinking of the Titanic.

We also include **gender\_submission.csv**, a set of predictions that assume all and only female passengers survive, as an example of what a submission file should look like.

#### Data Dictionary

**VariableDefinitionKey** survival Survival 0 = No, 1 = Yes pclass Ticket class 1 = 1st, 2 = 2nd, 3 = 3rd sex Sex Age Age in years sibsp # of siblings / spouses aboard the Titanic parch # of parents / children aboard the Titanic ticket Ticket number fare Passenger fare cabin Cabin number embarked Port of