

Ensembling and Stacking

Import Libraries

```
In [1]: 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
```

Read datasets

```
In [2]: train = pd.read_csv('../input/train.csv')
test = pd.read_csv('../input/test.csv')

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

train.head(3)
```

```
Out[2]:   PassengerId  Survived  Pclass          Name     Sex   Age  SibSp  Parch     Ticket  Fare Cabin Embarked
0            1         0      3  Braund, Mr. Owen Harris   male  22.0      1     0  A/5 21171  7.2500   NaN      S
1            2         1      1    Cumings, Mrs. John Bradley
                           (Florence Briggs Th... female  38.0      1     0       PC 17599  71.2833  C85      C
2            3         1      3       Heikkinen, Miss. Laina female  26.0      0     0  STON/O2.
                           3101282  7.9250   NaN      S
```

Feature Engineering

```
In [3]: full_data = [train, test]

# 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)

# 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-Z][a-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', 'Mme', 'Ms'], 'Miss')
    dataset['Title'] = dataset['Title'].replace('Mlle', '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 ;

```

Feature selection

```
In [4]: 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)
```

Visualisations

```
In [5]: train.head(3)
```

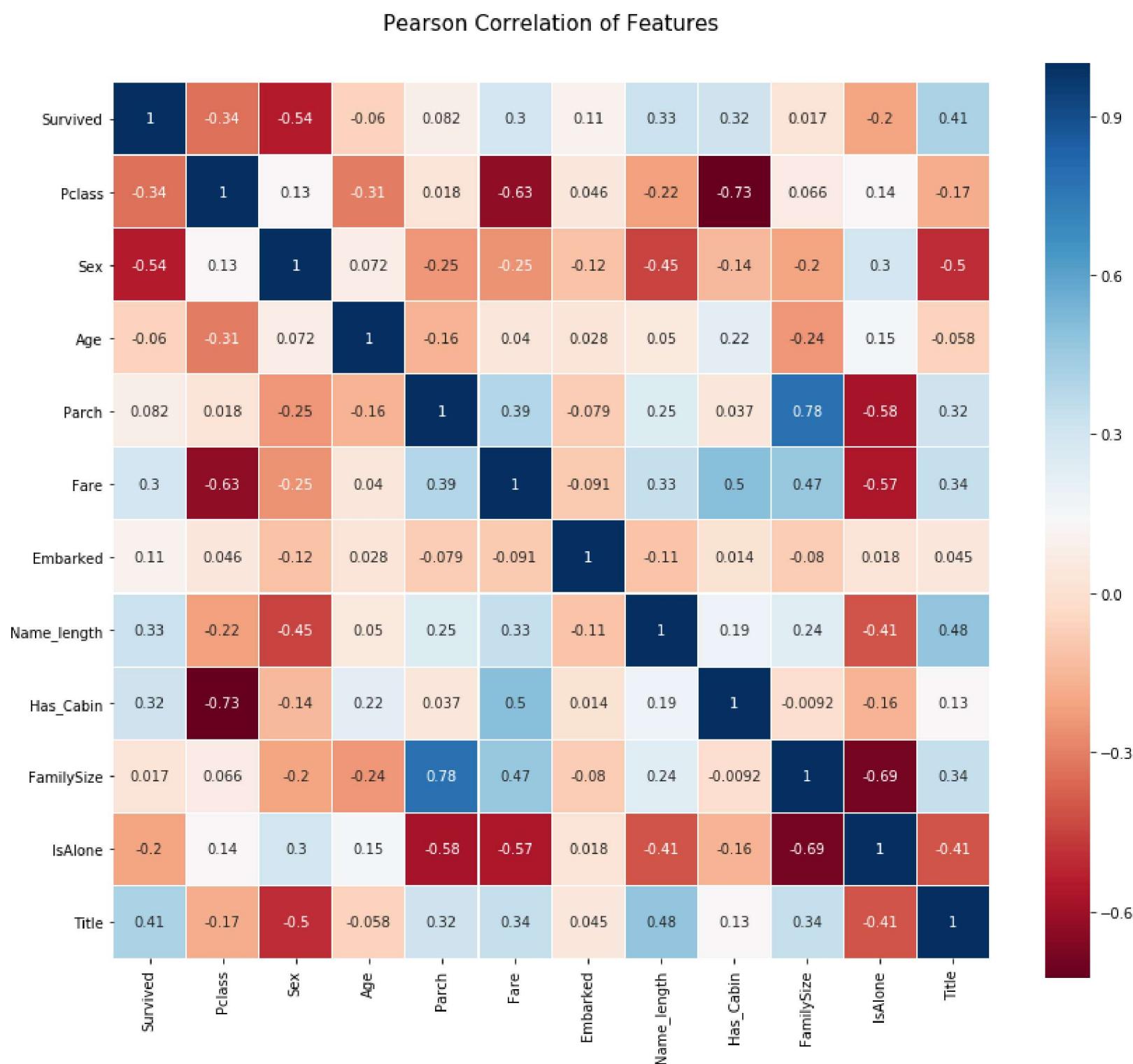
	Survived	Pclass	Sex	Age	Parch	Fare	Embarked	Name_length	Has_Cabin	FamilySize	IsAlone	Title
0	0	3	1	1	0	0	0	23	0	2	0	1
1	1	1	0	2	0	3	1	51	1	2	0	3
2	1	3	0	1	0	1	0	22	0	1	1	2

Pearson Correlation Heatmap

generate some correlation plots of the features to observe the relationship between each feature. To accomplish this, we'll utilize the Seaborn plotting package, which conveniently allows us to create heatmaps 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]: <matplotlib.axes._subplots.AxesSubplot at 0x79433cb22cf8>
```



Learn from the Plots

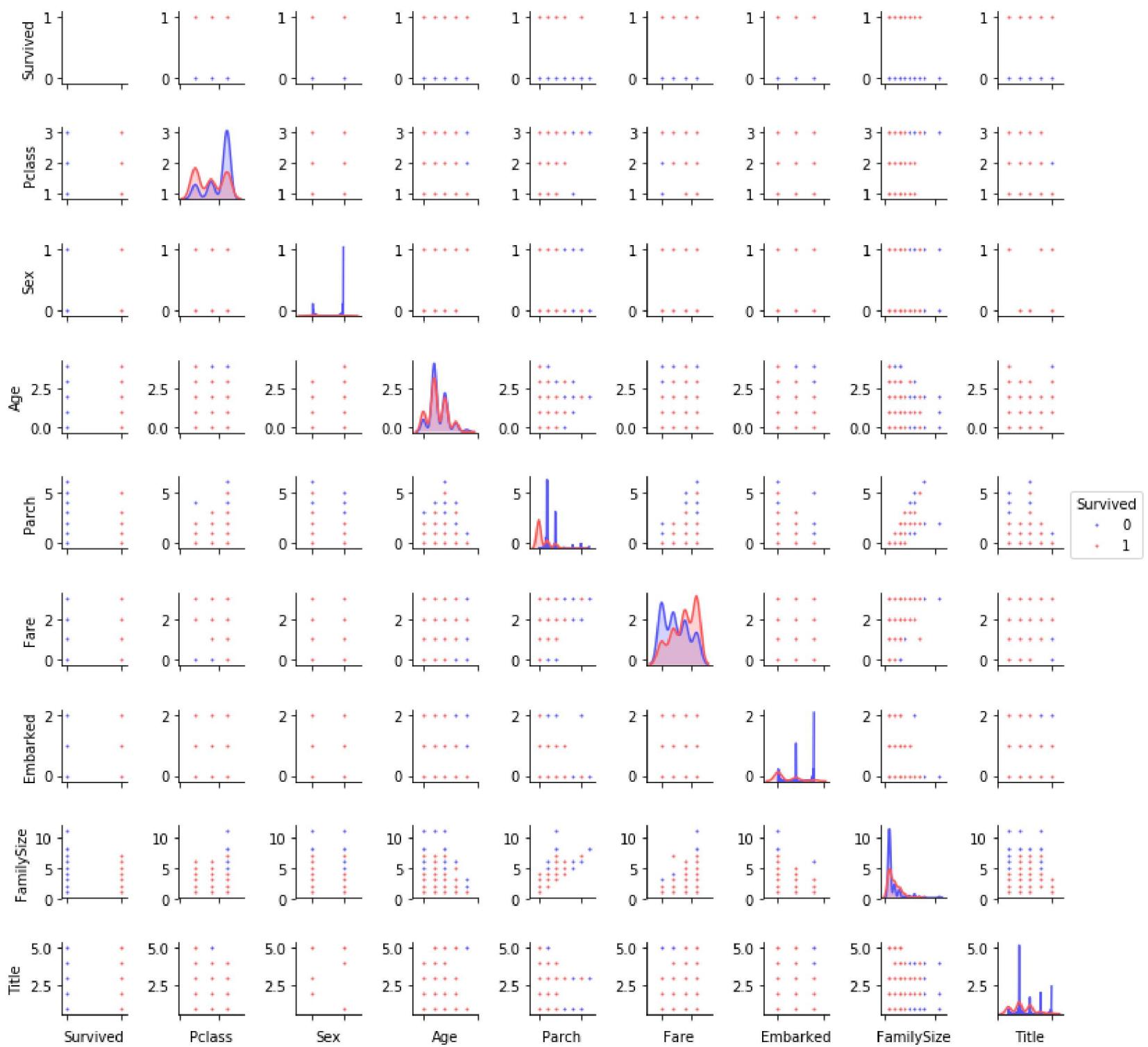
One insight from the Pearson Correlation plot is that there aren't too many strongly correlated features. This is advantageous when feeding these features into your learning model because it indicates minimal redundant or superfluous data in our training set, ensuring each feature carries unique information. The two most correlated features are Family size and Parch (Parents and Children). For the purposes of this exercise, I'll retain both features.

Pairplots:

Lastly, let's generate pairplots to visualize the data distribution across different features.

```
In [7]: g = sns.pairplot(train[[u'Survived', u'Pclass', u'Sex', u'Age', u'Parch', u'Fare', u'Embarked', u'FamilySize', u'Title']], hue='Survived', palette = 'seismic',size=1.2,diag_kind = 'kde',diag_kws=dict(shade=True),set(xticklabels=[]))

Out[7]: <seaborn.axisgrid.PairGrid at 0x794337757fd0>
```



Ensembling & Stacking models

```
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_)
```

Out-of-Fold 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 5 Machine Learning Models

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

Parameters

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.

```
In [10]: # 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
}
```

```
In [11]: # Create 5 objects that represent the 5 models
rf = SklearnHelper(clf=RandomForestClassifier, seed=SEED, params=rf_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, params=gb_params)
svc = SklearnHelper(clf=SVC, seed=SEED, params=svc_params)
```

Creating NumPy arrays out of train and test sets

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

Output of the First level Predictions

```
In [13]: 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_test) # 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

```
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)

[0.12409257 0.20142268 0.03274976 0.02172969 0.07164986 0.02373582
 0.10891636 0.06480882 0.06816018 0.01341809 0.26931617]
[0.12091429 0.38084424 0.03075754 0.01673822 0.05615037 0.02829344
 0.04750693 0.08287509 0.04331143 0.02191649 0.17069195]
[0.032 0.014 0.018 0.066 0.038 0.008 0.692 0.012 0.052 0. 0.068]
[0.07273928 0.0318522 0.09614415 0.03090687 0.10686161 0.05763081
 0.4035061 0.02355397 0.07132947 0.0238005 0.08167503]
```

```
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, 0.01066667]
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

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
)
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'].values,
    #     #color = np.random.randn(500), #set color equal to a variable
    #     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(

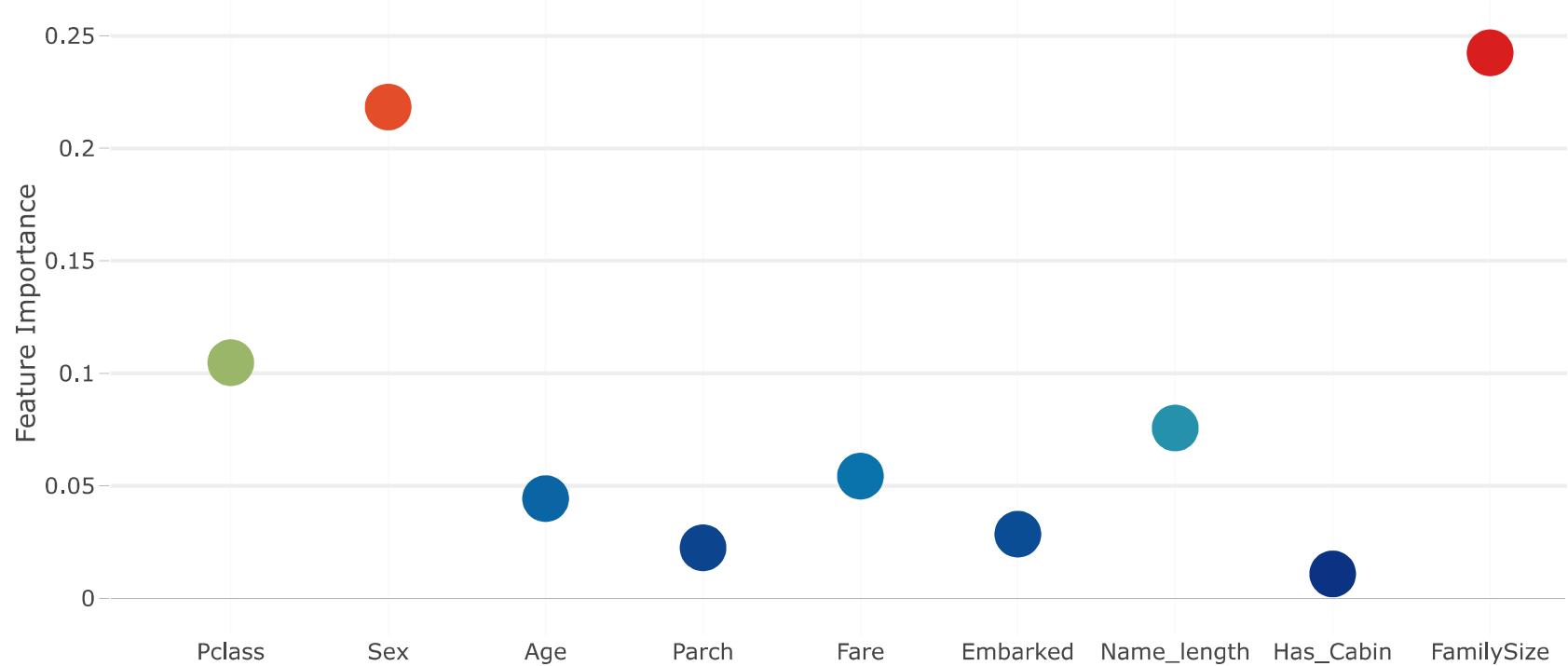
```

```

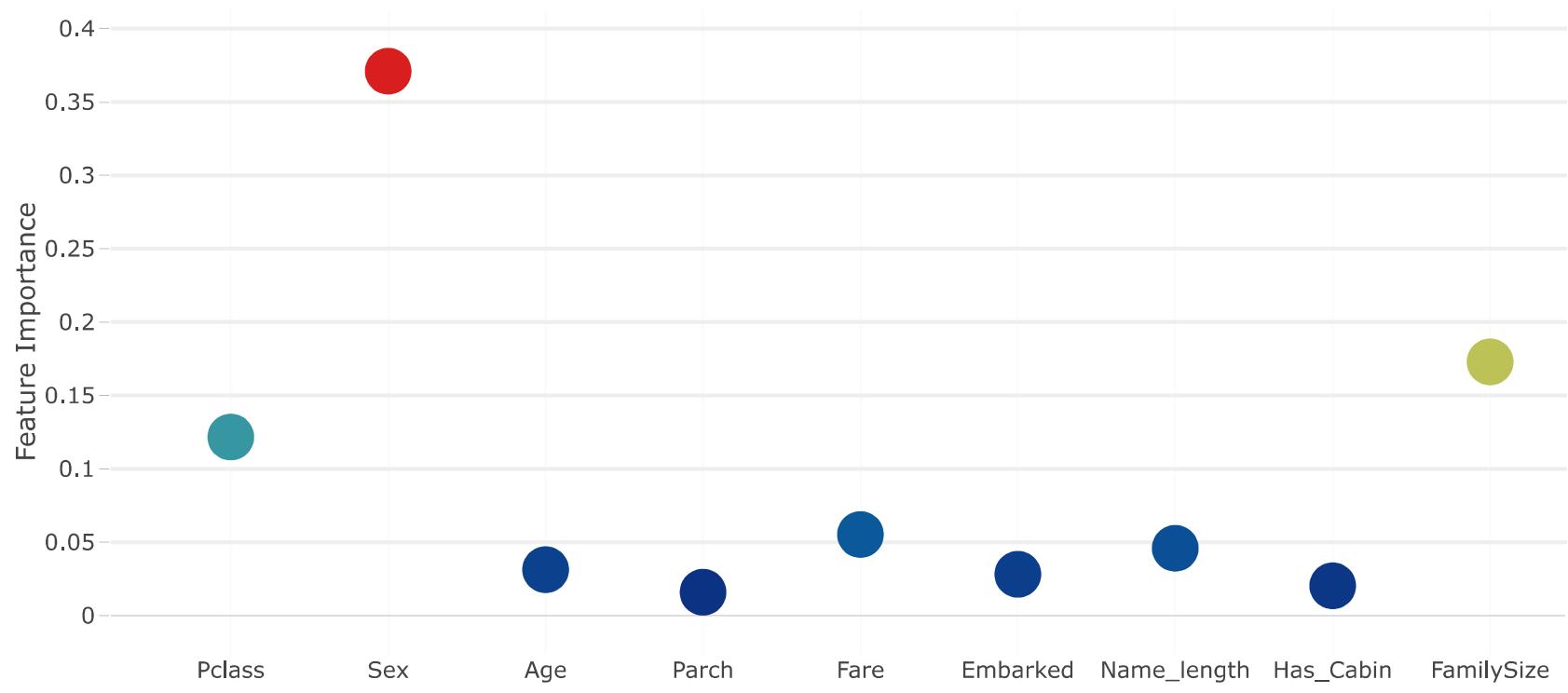
#         title= 'Pop',
#         ticklen= 5,
#         zeroline= False,
#         gridwidth= 2,
#     ),
yaxis=dict(
    title= 'Feature Importance',
    ticklen= 5,
    gridwidth= 2
),
showlegend= False
)
fig = go.Figure(data=data, layout=layout)
py.iplot(fig,filename='scatter2010')

```

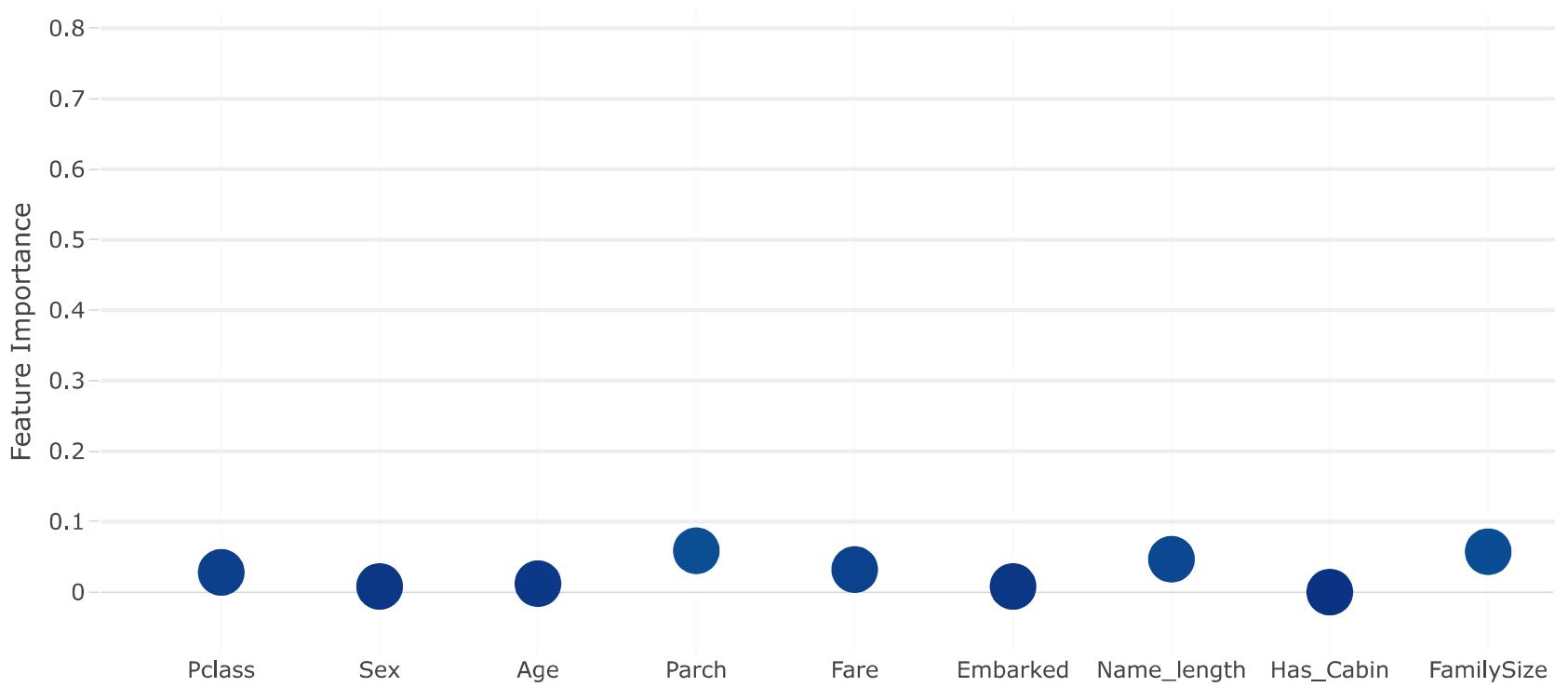
Random Forest Feature Importance



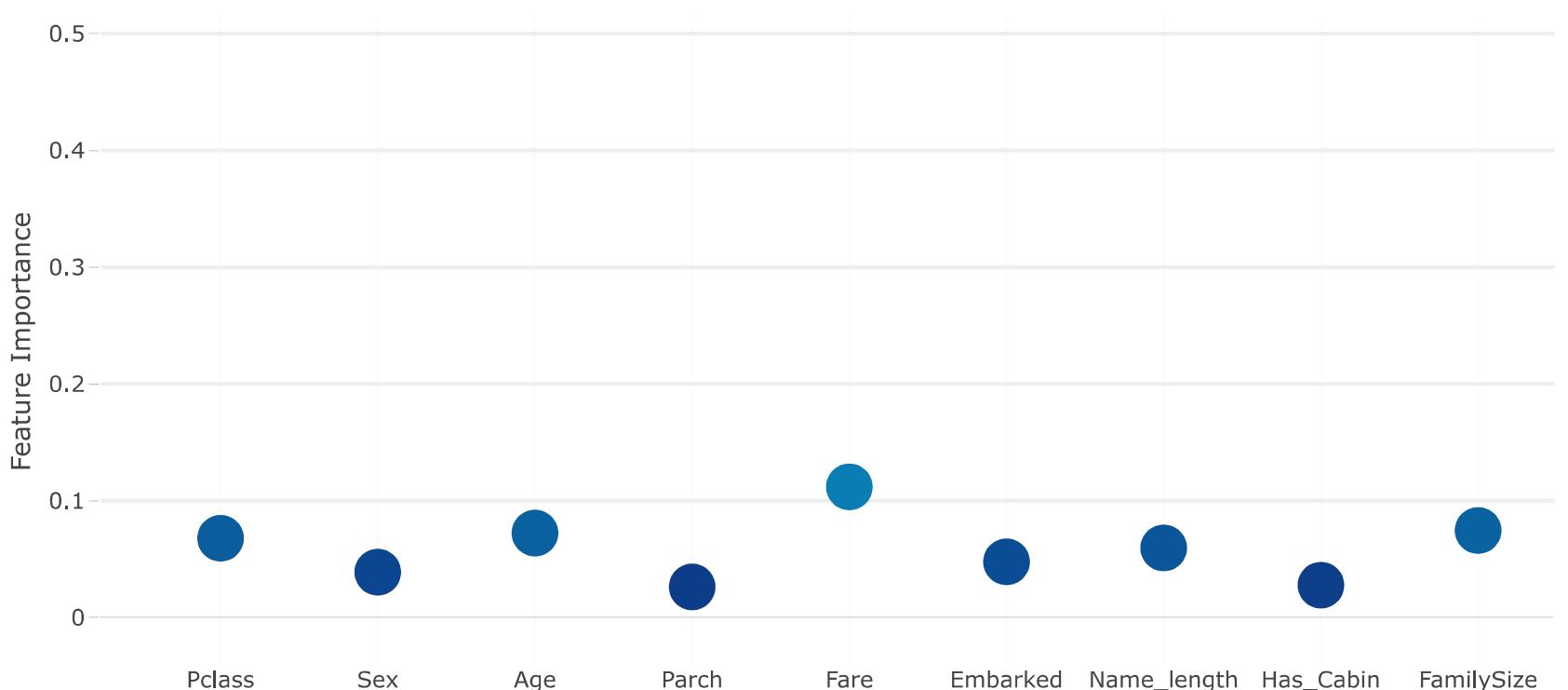
Extra Trees Feature Importance



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]:	AdaBoost feature importances	Extra Trees feature importances	Gradient Boost feature importances	Random Forest feature importances	features	mean
0	0.028	0.121657	0.067961	0.104741	Pclass	0.080590
1	0.008	0.370983	0.038893	0.218370	Sex	0.159062
2	0.012	0.031296	0.072378	0.044327	Age	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,
```

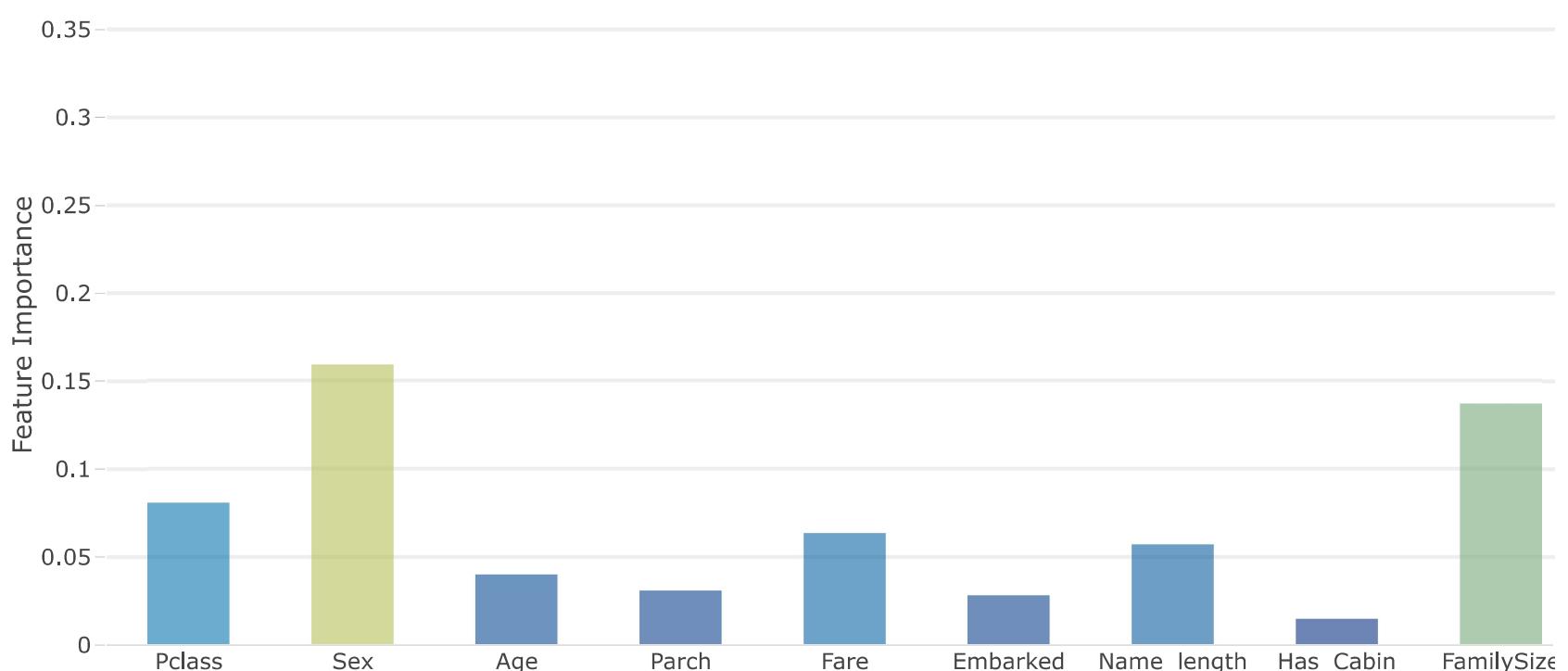
```

        colorscale='Portland',
        showscale=True,
        reversescale = False
    ),
    opacity=0.6
)

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

Now that we've acquired our first-level predictions, consider this process as essentially constructing a fresh set of features to serve as training data for the subsequent classifier. Following the code below, our new columns consist of the first-level predictions generated by our previous classifiers, which we use to train the next classifier.

```
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]: AdaBoost  ExtraTrees  GradientBoost  RandomForest
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          0.0
3          1.0         1.0          1.0          1.0
4          0.0         0.0          0.0          0.0
```

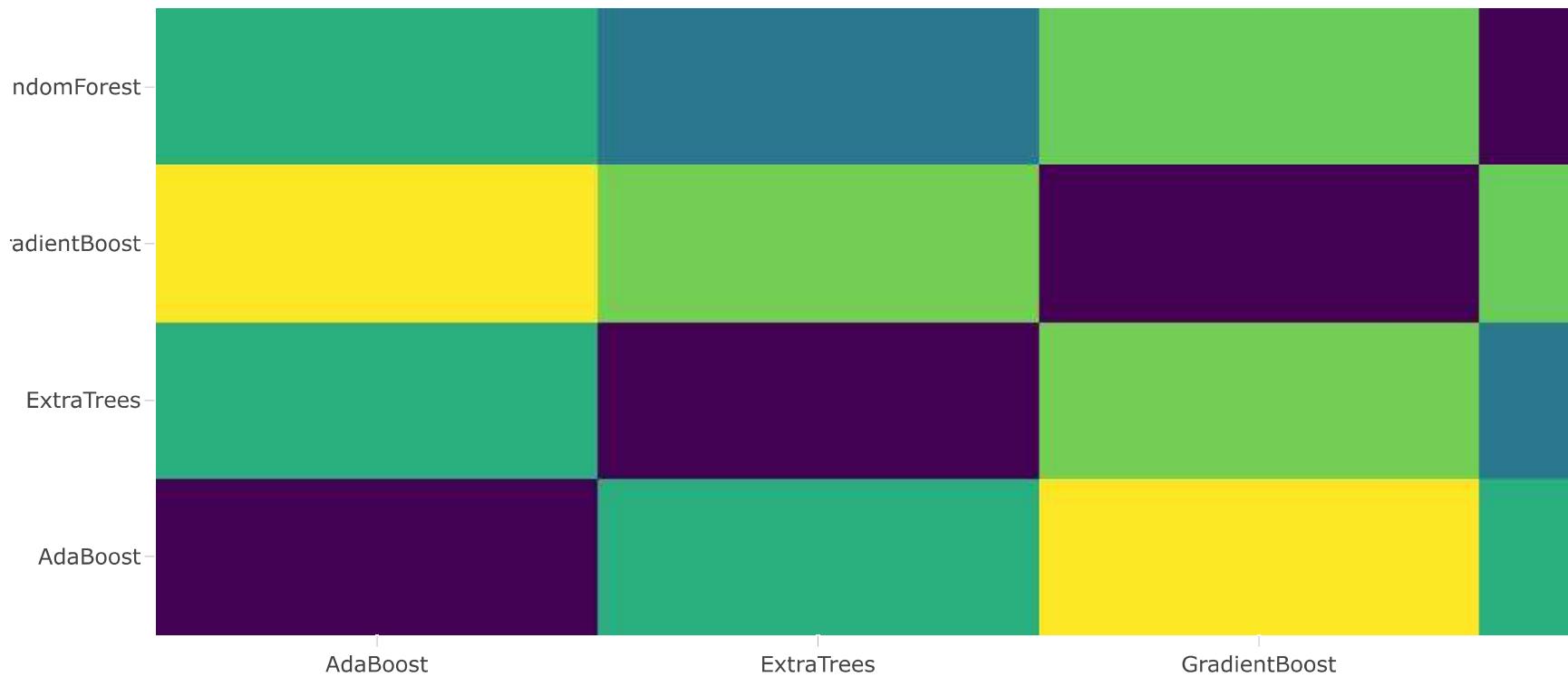
Correlation Heatmap of the Second Level Training set

```
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')

```



```
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)
```

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

```
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)
```

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

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