

Setup

```
In [128]:  
import seaborn as sns  
from rfpimp import *  
import numpy as np  
import scipy  
from rfpimp import plot_corr_heatmap  
import pandas as pd  
import PipelineProfiler  
import AutoSklearn_regressor  
from auto_sklearn import permutation_importance  
from auto_sklearn.model_selection import train_test_split  
from scikit_learn import RandomForestRegressor  
from sklearn.cluster import hierarchy as hc  
from matplotlib.pyplot import figure  
from matplotlib.pyplot import plt  
import rfpimp.  
from rfpimp import LogisticRegression, RANSACRegressor, ARDRegression  
from sklearn.linear_model import LinearRegression, SVR  
from sklearn.neighbors import KNeighborsClassifier  
from sklearn.svm import SVC  
from sklearn.ensemble import RandomForestClassifier, ExtraTreesRegressor, AdaBoostRegressor, GradientBoostingRegressor  
from sklearn.neural_network import MLPRegressor,MLPClassifier  
# Data processing, modeling, and model evaluation  
from sklearn.preprocessing import StandardScaler  
from sklearn.metrics import f1_score, classification_report  
from sklearn.metrics import confusion_matrix  
from sklearn.neural_network import MLPRegressor,MLPClassifier  
# Data processing, modeling, and model evaluation  
from sklearn.preprocessing import StandardScaler  
from sklearn.metrics import f1_score, classification_report, plot_confusion_matrix  
import matplotlib.pyplot as plt # for data visualization  
from pandas.io.json import json_normalize  
import statsmodels.formula.api as sm  
import statsmodels.linear_model as lm  
import LogisticRegression, Ridge  
from sklearn.feature_selection import RFE  
from sklearn.svm import SVR  
from sklearn import preprocessing  
from sklearn.neighbors import KNeighborsClassifier, KNeighborsRegressor  
from sklearn.ensemble import RandomForestRegressor, RandomForestClassifier  
warnings.filterwarnings("ignore")
```

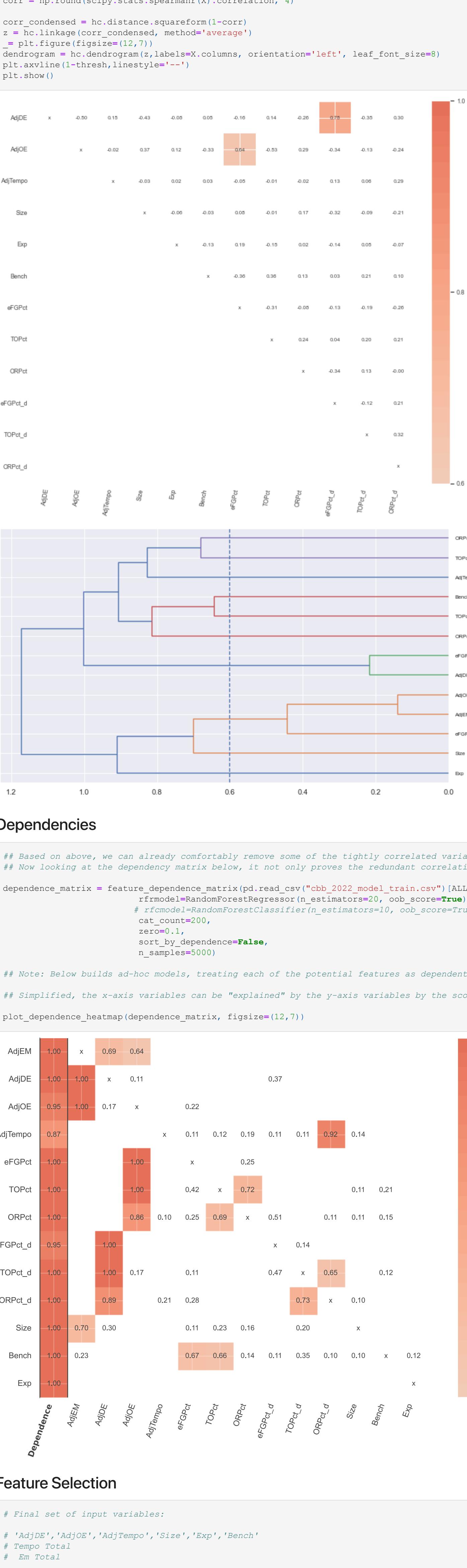
Data

```
In [129]:  
## KenPom (http://www.kenpom.com)  
df = pd.read_csv("cbb_2022_model_train.csv")
```

Distribution

```
In [130]:  
## Distribution of standard and kenPom stats, denoted by "BASE_FEATURES"  
BASE_FEATURES = ['AdjEM', 'AdjDE', 'AdjOE', 'AdjTempo']  
g = sns.PairGrid(df[BASE_FEATURES])  
g.map_diag(sns.histplot)  
g.map_offdiag(sns.scatterplot)
```

```
Out[130]: <seaborn.axisgrid.PairGrid at 0x149c7ce80>
```



```
In [131]:  
## Distribution of decomposed stats  
BASE_FEATURES = ['AdjEM', 'AdjDE', 'AdjOE', 'AdjTempo']  
DECOMPOSED_O_FEATURES = ['eFGPct_d', 'TOPct_d', 'ORPct_d']  
OTHER = ['Size', 'Bench']  
ALL_FEATURES = BASE_FEATURES + DECOMPOSED_O_FEATURES + DECOMPOSED_D_FEATURES + OTHER  
g = sns.PairGrid(df[DECOMPOSED_O_FEATURES])  
g.map_diag(sns.histplot)  
g.map_offdiag(sns.scatterplot)
```

```
##OFFENSE##
```

```
Out[131]: <seaborn.axisgrid.PairGrid at 0x14e53bb38>
```



Correlations

```
In [116]:  
## PLOT A: Correlation Heatmap ###  
## Choosing Spearman's coefficient because  
## a) It's better for identifying continuous, monotonic relationships, as we expect between these features  
## b) Doesn't assume linearity  
viz = plot_corr_heatmap(pd.read_csv("cbb_2022_model_train.csv")[ALL_FEATURES], figsize=(10,10))  
viz.view()
```

```
## PLOT B: Hierarchical analysis ##
```

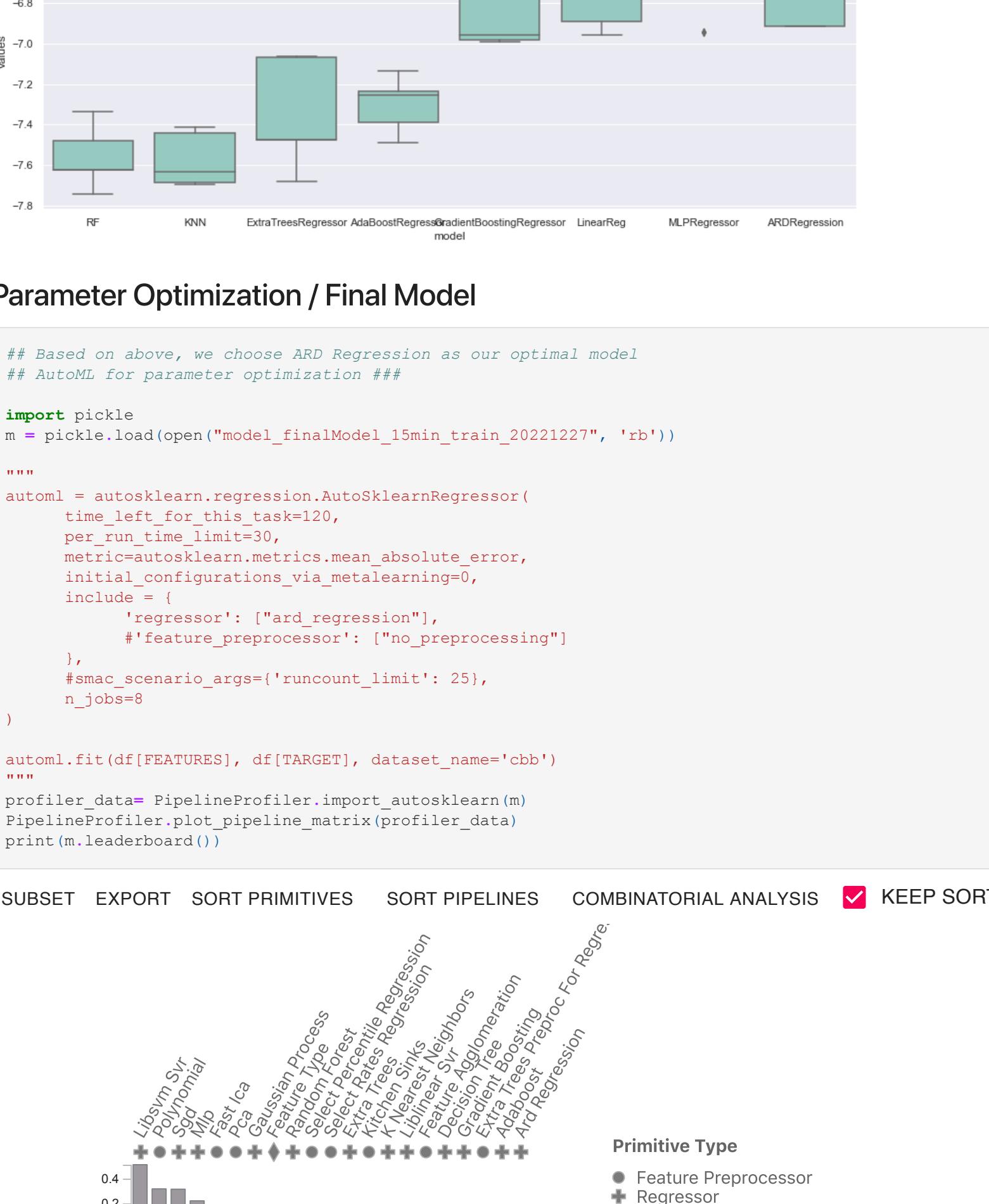
```
X = pd.read_csv("cbb_2022_model_train.csv")[ALL_FEATURES]  
# Arbitrary threshold  
thresh=0.4  
corr = np.round(scipy.stats.spearmanr(X).correlation, 4)
```

```
corr_condensed = hc.distance.squareform(1-corr)  
z = hc.linkage(corr_condensed, method="average")  
plt.figure(figsize=(12, 7))
```

```
dendrogram = hc.dendrogram(z, labels=X.columns, orientation="left", leaf_font_size=8)
```

```
plt.title(1-thresh,linestyle="--")
```

```
plt.show()
```



Dependencies

```
In [143]:  
## Based on above, we can already comfortably remove some of the tightly correlated variables (e.g. remove eFGPct_d, TOPct_d, ORPct_d, etc.)  
## Now looking at the dependency matrix below, it not only proves the redundant correlations, but gives some insight.
```

```
dependence_matrix = feature_dependence_matrix(pd.read_csv("cbb_2022_model_train.csv")[ALL_FEATURES],  
                                              rfmodel=RandomForestRegressor(n_estimators=20, oob_score=True),  
                                              cat_count=200,  
                                              zero_mean=False,  
                                              sort_by=dependence=False,  
                                              n_samples=5000)
```

```
## Note: Below builds ad-hoc models, treating each of the potential features as dependent variables, and measures the error against the real value.
```

```
## Simplified, the x-axis variables can be "explained" by the y-axis variables by the score denoted in D(x,y)
```

```
plot_dependence_heatmap(dependence_matrix, figsize=(12,7))
```

```
Out[143]: <seaborn.axisgrid.PairGrid at 0x159c7f780>
```


Feature Selection

```
In [1]:  
# Final set of input variables:  
# 'AdjDE', 'AdjOE', 'AdjTempo', 'Size', 'Exp', 'Bench'  
# Temp Total  
# Em Total
```

```
# Evaluates various models benchmarked against  
# Mean Absolute Error against real score  
# Error against opening Vegas spread
```

```
# Features  
BASE_FEATURES = ['AdjDE', 'AdjOE', 'AdjTempo',  
                  'Size', 'Exp', 'Bench']
```

```
FEATURES = []  
for f in BASE_FEATURES:  
    FEATURES += [f + '_opp']
```

```
FEATURES += FEATURES + ['home'] + ['tempTotal', 'emTotal']
```

```
TARGET = ['margin']
```

```
#####
```

```
#####
```

```
# PCA transformations (dimensionality reduction)
```

```
# NOT BEING USED, but could be useful
```

```
PCA_TRANSFORMATIONS=2
```

```
bclf = 'scf'.StandardScaler()  
bclf = scf.dropna(subset=FEATURES)
```

```
scaler = StandardScaler()  
X = bclf.query(TRAIN_QUERY)[FEATURES]
```

```
X = scaler.fit_transform(Xt)
```

```
pca = PCA(n_components=PCA_TRANSFORMATIONS)
```

```
pca.fit(Xt)
```

```
X = bclf.fit_transform(X)
```

```
X = pca.fit_transform(X)
```

```
#X = scaler.fit_transform(X)
```

```
#FEATURES=[]
```

```
for i in range(0,PCA_TRANSFORMATIONS):  
    bclf['PCA'+str(i)] = X[:,i,:]
```

```
#FEATURES = FEATURES + ['PCA'+str(i)]
```

```
#####
```

```
#Normalize the data, split data set into train / test
```

```
#bcf = bclf[FEATURES+TARGET+['spread']].dropna()
```

```
#ARBITRARY_SPLIT_BETWEEN_TRAIN_AND_TEST ####
```

```
TRAIN_QUERY="conference='MWC'"
```

```
TEST_QUERY="conference='NWC'"
```

```
X_train = bclf.query(TRAIN_QUERY)[FEATURES+TARGET+['abs_vegas_error', 'spread']].dropna()
```

```
y_train = bclf.query(TRAIN_QUERY)[TARGET].dropna()
```

```
print("Train size: " + str(len(y_train)))
```

```
X_test = bclf.query(TEST_QUERY)[FEATURES+TARGET+['abs_vegas_error', 'spread']].dropna()
```

```
y_test = bclf.query(TEST_QUERY)[TARGET].dropna()
```

```
print("Test size: " + str(len(X_test)))
```

```
#####
```

```
# Evaluate models
```

```
dfs = []
```

```
models = [
```

```
    ('LinearReg', LinearRegression()),
```

```
    ('KNN', KNeighborsRegressor(max_depth=1, max_leaf_nodes=10)),
```

```
    ('GNB', GaussianNB()),
```

```
    ('DTreeReg', DecisionTreeRegressor(max_depth=1, max_leaf_nodes=10)),
```

```
    ('ETreeReg', ExtraTreesRegressor(max_depth=1, max_leaf_nodes=10)),
```

```
    ('ABReg', AdaBoostRegressor(max_depth=1, max_leaf_nodes=10)),
```

```
    ('GRReg', GradientBoostingRegressor(max_depth=1, max_leaf_nodes=10)),
```

```
    ('ARDReg', ARDRegression(max_depth=1, max_leaf_nodes=10)),
```

```
    ('RANSACReg', RANSACRegressor(max_depth=1, max_leaf_nodes=10)),
```

```
    ('MLPReg', MLPRegressor(random_state=1, max_iter=250, hidden_layer_sizes=100)),
```

```
    ('RF', RF)
```

```
]
```

```
results = []
```

```
names = []
```

```
scoring = ['r2', 'neg_median_absolute_error']
```

```
for name, model in models:  
    kfold = model_selection.KFold(n_splits=5, shuffle=True, random_state=0)
```

```
    results = model_selection.cross_val_score(model, X_train[FEATURES].append(X_test[FEATURES]), X_train[TARGET], cv=kfold)
```

```
    print("Mean score for " + name + " is " + str(np.mean(results)))
```

```
    results.append(name)
```

```
    this_df = pd.DataFrame([{'model': name, 'mean': np.mean(results)}], index=[0])
```

```
    this_df.append(this_df, ignore_index=True)
```

```
final = pd.concat(dfs, ignore_index=True)
```

```
#####
```

```
# Performance metrics
```

```
boot_size = 1000
```

```
for model in list(set(final.model.values)):
```

```
    model_df = final.loc[final.model == model]
```

```
    bootstraps = model_df.sample(n=boot_size, replace=True)
```

```
    bootstraps.append(bootstraps)
```

```
    time_metrics = pd.concat([pd.DataFrame(bootstraps['model'].values, columns=['model']),
```

```
                             pd.DataFrame(bootstraps['mean'].values, columns=['mean'])], axis=1)
```

```
    time_metrics['model'] = model
```

```
    time_metrics['mean'] = np.mean(bootstraps['mean'].values)
```

```
    time_metrics['std'] = np.std(bootstraps['mean'].values)
```

```
    time_metrics['lower'] = time_metrics['mean'] - time_metrics['std']
```

```
    time_metrics['upper'] = time_metrics['mean'] + time_metrics['std']
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics.set_index('model', inplace=True)
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', 'mean', 'std', 'lower', 'upper']]
```

```
    time_metrics = time_metrics[['model', '
```