





```
X = df.iloc[:,:]
y = df.iloc[:,:]

df.shape

Treat Imbalance Data

In [ ]:
y.value_counts()

In [ ]:
ros = RandomOverSampler(sampling_strategy='all',random_state=0)

In [ ]:
new_X, new_y = ros.fit_resample(X, y)

In [ ]:
new_y().value_counts()

In [ ]:
new_X

Train Test Split Cont'd

In [ ]:
X.values, y.values

In [ ]:
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)

In [ ]:
X_train.shape, X_test.shape, y_train.shape, y_test.shape

Train Test Split to create Train, Validation and Test Set

In [ ]:
#set test set size
X_df, X_test, y_df, y_test = train_test_split(X.values, y.values, test_size=0.2, random_state=0)

In [ ]:
X_df.shape, X_test.shape, y_df.shape, y_test.shape

In [ ]:
X_train, X_val = train_test_split(X_df, test_size=0.2, random_state=0)

In [ ]:
X_train.shape, X_val.shape

Feature Selection

Using SelectKBest

In [ ]:
X_new = SelectKBest(f_regression, k=10).fit_transform(X_train,y_train)

In [ ]:
X_new[0:5]

In [ ]:

Univariate Feature Selection

In [ ]:
select_feature = SelectKBest(chi2, k=10).fit(X_train,y_train)

In [ ]:
select_feature.scores_

Recursive Feature Elimination

In [ ]:
rfe = RFE(estimator=XGBRegressor(),n_features_to_select=10,verbose=1, step=1)

In [ ]:
rfe.fit(X_train,y_train)

In [ ]:
selected_rfe_features = pd.DataFrame({'Feature':list(X_train.columns),'Ranking':rfe.ranking_})

In [ ]:
selected_rfe_features

Recursive Feature Elimination with Cross Validation

In [ ]:
rfecv = RFECV(estimator=XGBRegressor(), cv=5,scoring="neg_mean_squared_error",verbose=1, step=1)

In [ ]:
rfecv.fit(X_train,y_train)

In [ ]:
print("Optimal no of features:", rfecv.n_features_)

In [ ]:
print("Best features:", rfecv.support_ )

Feature Scaling

In [ ]:
X_train

In [ ]:
encoder = LabelEncoder()

In [ ]:
scaler = StandardScaler()

In [ ]:
minmax = MinMaxScaler()

In [ ]:
one = OneHotEncoder()

In [ ]:
X_train_scaled = minmax.fit_transform(X_train)

In [ ]:
X_test_scaled = minmax.transform(X_test)

In [ ]:
X_train_scaled

In [ ]:
X_test_scaled

In [ ]:

In [ ]:

Model Training

Using PyCaret

In [ ]:
exp_reg = setup(data = df, target = '', session_id=0, normalize=True)

In [ ]:
compare_models(exclude=['catboost','lightgbm','lda','qda','mlp','ada','nb','ridge','rbfsvm','svm'],fold=5) #For

In [ ]:
compare_models(exclude=['omp','br','ad','pac','rannac','ts','huber','kr','svm','knn','dt','rf','et','ada','gb','mlp','xgboost','lightgbm','catboost'],fold=5) # For Regressor

In [ ]:
model_selected = create_model('catboost')

In [ ]:
print(model_selected)

In [ ]:
tuned_model = tune_model(catboost, optimizer='mse')

In [ ]:
print(tuned_model)

In [ ]:
plot_model(tuned_model)

In [ ]:
plot_model(tuned_model, plot = 'error')

In [ ]:
plot_model(tuned_model, plot='feature')

In [ ]:
interpret_model(tuned_model)

In [ ]:
evaluate_model(tuned_model)

In [ ]:
predict_model(tuned_model)

In [ ]:
final_model = finalize_model(tuned_model)

In [ ]:
unseen_predictions = predict_model(final_model, data=data_unseen)
unseen_predictions.head()

Using Regression or Classification Models

In [ ]:
reg_model

In [ ]:
classf_model

K-Fold Cross-Validation (Generalization Performance)

In [ ]:
lasso = Lasso(random_state=0)

In [ ]:
kf = KFold(n_splits=5, shuffle=True, random_state=0)

In [ ]:
lasso_cv = cross_validate(estimator=lasso, X=X_train_scaled, y=y_train, scoring="neg_root_mean_squared_error",
                           cv=kf, n_jobs=-1,return_train_score=True)

In [ ]:
lasso_cv

In [ ]:
np.mean(lasso_cv["train_score"]), np.std(lasso_cv["train_score"])

In [ ]:
np.mean(lasso_cv["test_score"]), np.std(lasso_cv["test_score"])

In [ ]:

In [ ]:

In [ ]:

In [ ]:

Using TPOT

In [ ]:
# tpot = TPOTClassifier(generations=3,population_size=10,scoring='accuracy', cv=5, verbosity=2, random_state=0,
#                       early_stop=1)

In [ ]:
# tpot.fit(X_train,y_train)

In [ ]:
# tpot.score(X_test, y_test)

Using XGBoost (Scikit-Learn)

Using RandomSearchCV

In [ ]:
model = XGBRegressor(random_state=0, n_estimators=100, objective='reg:squarederror')

In [ ]:
model = XGBClassifier(random_state=0, n_estimators=100, objective='softmax:multi')

In [ ]:
parameters = {'max_depth': np.arange(3,10,1),
              'eta': np.arange(0.05,0.3,0.05),
              'n_estimators': np.arange(100,1000,100),
              'min_child_weight': np.arange(1,4,1),
              'gamma': np.arange(0,10,2),
              'subsample': np.arange(0.5,0.9,0.1),
              'colsample_bytree': np.arange(0.5,0.9,0.1),
              'reg_alpha': np.arange(0,1,0.1),
              'reg_lambda': np.arange(0,1,0.1)
              }

In [ ]:
randm = RandomizedSearchCV(estimator=model, param_distributions = parameters, cv = 5, n_iter = 50,
                           n_jobs=-1, scoring='')

In [ ]:
randm.fit(X, y)

In [ ]:
randm.best_estimator_

In [ ]:
randm.best_score_

In [ ]:
randm.best_params_

In [ ]:

In [ ]:

Final Model

In [ ]:
xgbmodel = XGBRegressor(random_state=0, n_estimators=100, objective='reg:squarederror')

In [ ]:
xgbmodel = XGBClassifier(random_state=0, n_estimators=100, objective='binary:logistic')

In [ ]:
xgbmodel = XGBClassifier(random_state=0, n_estimators=100, objective='softmax:multi')

In [ ]:
xgbmodel.fit(X_train_scaled,y_train,eval_set=[(X_test_scaled,y_test)],eval_metric='rmse',early_stopping_rounds=10)

In [ ]:
xgbmodel.fit(X_train_scaled,y_train,eval_set=[(X_test_scaled,y_test)],eval_metric='error',early_stopping_rounds=10)

In [ ]:
xgbmodel.fit(X_train_scaled,y_train,eval_set=[(X_test_scaled,y_test)],eval_metric='mlogloss',early_stopping_rounds=10)

In [ ]:
y_pred = xgbmodel.predict(X_test_scaled)

In [ ]:
y_pred

Model Evaluation

In [ ]:
cm = confusion_matrix(y_test,y_pred)
cm

In [ ]:
fig, ax = plt.subplots(figsize=(10,5))
sns.heatmap(cm, annot=True,fmt='.4g',linewidths=2, cmap='viridis')
plt.xlabel('True label')
plt.ylabel('Predicted label')
plt.show()

In [ ]:
fig, ax = plt.subplots(figsize=(10,5))
sns.heatmap(cm, annot=True,fmt='.4g',linewidths=2, cmap='viridis')
plt.xlabel('True label')
plt.ylabel('Predicted label')
plt.show()

In [ ]:
print(classification_report(y_test,y_pred))

In [ ]:
plot_roc_curve(xgbmodel,X_test,y_test)
plt.show()

In [ ]:
mse = mean_squared_error(y_test,y_pred)

In [ ]:
rmse = np.sqrt(mse)

In [ ]:
rmse

In [ ]:
r2score = r2_score(y_test,y_pred)
r2score

In [ ]:
fig, ax = plt.subplots(figsize=(10,8))
sns.regplot(x=y_test, y=y_pred, ax=ax)
plt.title("Plot to compare actual vs predicted", fontsize=20)
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.show()

Plot Feature Importances

In [ ]:
rf.feature_importances_

In [ ]:
feat_importances = pd.Series(rf.feature_importances_, index=X.columns)

In [ ]:
feat_importances

In [ ]:
feat_importances.nlargest(10).plot(kind='barh', figsize=(10,10))
plt.title('Feature Importances')
plt.show()

The permutation based importance

In [ ]:
perm_importance = permutation_importance(rf,X_test,y_test, random_state=0, scoring='neg_mean_squared_error')

In [ ]:
sorted_idx = perm_importance.importances_mean.argsort()
plt.figure(figsize=(10,10))
plt.title("Permutation-based Importance")
plt.barh(X.columns[sorted_idx], perm_importance.importances_mean[sorted_idx])
plt.xlabel("Permutation Importance")
plt.show()

Compute Importance from SHAP Values

In [ ]:
explainer = shap.TreeExplainer(rf)

In [ ]:
shap_values = explainer.shap_values(X_test)

In [ ]:
shap.summary_plot(shap_values, X_test, plot_type="bar")

In [ ]:
shap.summary_plot(shap_values, X_test)

Available importance_types = ['weight', 'gain', 'cover', 'total_gain', 'total_cover']

In [ ]:
X.columns

In [ ]:
xgbmodel.get_booster().feature_names = X.columns

In [ ]:
fig, ax = plt.subplots(figsize=(20,10))
xgb.plot_importance(xgbmodel.get_booster(),ax=ax)
plt.show()

In [ ]:
xgb.to_graphviz(xgbmodel,num_trees=100)

Example:
f='gain'

XGBClassifier.get_booster().get_score(importance_type= f)

Plot Tree

In [ ]:
X.columns

In [ ]:
plt.figure(figsize=(40,25))
plot_tree(treeclf, feature_names=X.columns,class_names=['0','1'], fontsize=14, filled=True)
plt.show()

In [ ]:

Cross-Validation

In [ ]:
cv = cross_val_score(xgbmodel,X,y,cv=5,verbose=1,scoring='')

In [ ]:
cv.mean()

Using XGBoost (API)

In [ ]:

In [ ]:
dtrain = xgb.DMatrix(data=X_train,label=y_train)
dtest = xgb.DMatrix(data=X_test,label=y_test)

In [ ]:
params = {'n_estimators':,
          'learning_rate':,
          'max_depth':,
          'objective': '',
          'num_class':,
          'seed': 0,
          'eval_metric': ''}

In [ ]:
xgbmodel = xgb.train(params=params,dtrain=dtrain,num_boost_round=100,evals=[(dtest,"Test")],
                    early_stopping_rounds=10)

In [ ]:
y_pred = xgbmodel.predict(dtest)

In [ ]:
y_pred

Cross-Validation (API)

In [ ]:
cv = xgb.cv(params=params,
            dtrain=dtrain,
            num_boost_round=100,
            nfold=5,
            stratified=False,
            folds=None,
            metrics='merror',
            obj=None,
            feval=None,
            maximize=False,
            early_stopping_rounds=10,
            fpreproc=None,
            as_pandas=True,
            verbose_eval=None,
            show_stdv=True,
            seed=0,
            callbacks=None,
            shuffle=True, )

In [ ]:
cv

In [ ]:
cv['test-merror-mean'].min()

Model Evaluation (Classification)

In [ ]:
cm = confusion_matrix(y_test,y_pred)
cm

In [ ]:
fig, ax = plt.subplots(figsize=(10,5))
sns.heatmap(cm, annot=True,fmt='.4g',linewidths=2, cmap='viridis')
plt.xlabel('True label')
plt.ylabel('Predicted label')
plt.show()

In [ ]:
fig, ax = plt.subplots(figsize=(10,5))
plot_confusion_matrix(xgbmodel,X_test,y_test,values_format='.4g',ax=ax)
plt.show()

In [ ]:
print(classification_report(y_test,y_pred))

In [ ]:
plot_roc_curve(xgbmodel,X_test,y_test)
plt.show()

Model Evaluation (Regression)

In [ ]:
mse = mean_squared_error(y_test,y_pred)

In [ ]:
mse

In [ ]:
rmse = np.sqrt(mse)

In [ ]:
rmse

In [ ]:
r2score = r2_score(y_test,y_pred)
r2score

In [ ]:
fig, ax = plt.subplots(figsize=(10,8))
sns.regplot(x=y_test, y=y_pred, ax=ax)
plt.title("Plot to compare actual vs predicted")
plt.xlabel("Predicted")
plt.ylabel("Actual")
plt.show()

Table Formatted View

In [ ]:
table = X_test.copy()

In [ ]:
table["True Value"] = y_test.copy()

In [ ]:
table["Predicted"] = np.round(lr_pred,2)

In [ ]:
table

Cross-Validation

In [ ]:
cv = cross_val_score(xgbmodel,X,y,cv=5,verbose=1,scoring='accuracy')

In [ ]:
cv.mean()

Feature Selection

In [ ]:
df.columns

In [ ]:
df2 = df[['']]

In [ ]:
df2

In [ ]:
X = df2.iloc[:,0:7]
y = df2.iloc[:,7]

In [ ]:
X.values, y.values

In [ ]:
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)

In [ ]:
xgbmodel1 = XGBClassifier(random_state=0, n_estimators=100, objective='binary:logistic')

In [ ]:
xgbmodel1 = XGBRegressor(random_state=0, n_estimators=100, objective='reg:squarederror')

In [ ]:
xgbmodel1.fit(X_train,y_train,eval_set=(X_test,y_test),eval_metric='error',early_stopping_rounds=10)

In [ ]:
xgbmodel1.fit(X_train,y_train,eval_set=(X_test,y_test),eval_metric='rmse',early_stopping_rounds=10)

In [ ]:
y_pred = xgbmodel1.predict(X_test)

In [ ]:
y_pred

In [ ]:
fig, ax = plt.subplots(figsize=(5,5))
plot_confusion_matrix(xgbmodel1,X_test,y_test,values_format='4g',ax=ax)
plt.show()

In [ ]:
print(classification_report(y_test,y_pred))

In [ ]:
plot_roc_curve(xgbmodel1,X_test,y_test)
plt.show()

Model Prediction

In [ ]:
testdata = pd.read_csv()

In [ ]:
answer = xgbmodel.predict(testdata)

In [ ]:
answer

Model Tuning

Using RandomSearchCV

In [ ]:
model = XGBClassifier(random_state=0, n_estimators=100, objective='binary:logistic')

In [ ]:
model = XGBRegressor(random_state=0, n_estimators=100, objective='reg:squarederror')

In [ ]:
parameters = {'max_depth': np.arange(3,10,1),
              'learning_rate': np.arange(0.05,0.3,0.03),
              'n_estimators': np.arange(100,1000,100),
              'min_child_weight': np.arange(1,4,1),
              'gamma': np.arange(0,50,2),
              'subsample': np.arange(0.5,0.9,0.1),
              'colsample_bytree': np.arange(0.5,0.9,0.1)
              }

In [ ]:
randm = RandomizedSearchCV(estimator=model, param_distributions = parameters, cv = 5, n_iter = 50,
                           n_jobs=-1, scoring='')

In [ ]:
randm.fit(X, y)

In [ ]:
randm.best_estimator_

In [ ]:
randm.best_score_

In [ ]:
randm.best_params_

Using GridSearchCV

In [ ]:
model = XGBClassifier(random_state=0, n_estimators=100, objective='binary:logistic')

In [ ]:
parameters = {'max_depth': np.arange(3,10,1),
              'learning_rate': np.arange(0.05,0.3,0.03),
              'n_estimators': np.arange(100,1000,100),
              'min_child_weight': np.arange(1,4,1),
              'gamma': np.arange(0,50,2),
              'subsample': np.arange(0.5,0.9,0.1),
              'colsample_bytree': np.arange(0.5,0.9,0.1)
              }

In [ ]:
grids = GridSearchCV(estimator=model,param_grid=parameters,scoring='accuracy',
                     n_jobs=-1,cv=5,verbose=1,return_train_score=True)

In [ ]:
grids.fit(X,y)

In [ ]:
grids.best_estimator_

In [ ]:

Final Model

In [ ]:
xgbnew = XGBClassifier(random_state=, n_estimators=, objective='binary:logistic',max_depth=,
                      gamma=, min_child_weight=,learning_rate=,subsample=,colsample_bytree=)

In [ ]:
xgbnew = XGBRegressor(random_state=0, n_estimators=, objective='binary:logistic',max_depth=,
                      gamma=, min_child_weight=,learning_rate=,subsample=,colsample_bytree=)

In [ ]:
xgbnew.fit(X_train,y_train,eval_set=(X_test,y_test),eval_metric='error',early_stopping_rounds=10)

In [ ]:
y_pred = xgbnew.predict(X_test)

In [ ]:
y_pred

In [ ]:
fig, ax = plt.subplots(figsize=(5,5))
plot_confusion_matrix(xgbnew,X_test,y_test,values_format='4g',ax=ax)
plt.show()

In [ ]:
print(classification_report(y_test,y_pred))

In [ ]:
plot_roc_curve(xgbnew,X_test,y_test)
plt.show()

Save the Model

In [ ]:
filename = 'model.sav'
dump(xgbnew,open(filename,'wb'))

Load the Model

In [ ]:
loaded_model = load(open(filename,'rb'))

In [ ]:
loaded_model

Python code done by Dennis Lam

In [ ]:
```