

# Ray Jennings

## CS 677, Fall 2022

### Project #3

```
In [1]: import numpy as np
import pandas as pd

import seaborn as sns
import matplotlib.pyplot as plt
from ipynondisplay import display, HTML

from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
from sklearn.metrics import accuracy_score
from sklearn.metrics import roc_auc_score

from sklearn.model_selection import GridSearchCV
from sklearn.linear_model import LogisticRegression
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier

from xgboost import XGBClassifier

import eli5
from eli5.sklearn import PermutationImportance

import shap
import time
import time, line_tabular
```

#### Part 1 - EDA

```
In [2]: # Read Heart Failure dataset
#
df = pd.read_csv('./heart_failure_clinical_records_dataset.csv', index_col=False)
df.reset_index(drop=True, inplace=True)

df.head()
```

	age	anaemia	creatinine_phosphokinase	diabetes	ejection_fraction	high_blood_pressure	platelets	serum_creatinine
0	75.0	0	582	0	20		1	265000.00
1	65.0	0	7861	0	38		0	263358.03
2	65.0	0	146	0	20		0	162000.00
3	60.0	1	111	0	20		0	210000.00
4	65.0	1	160	1	20		0	327000.00

#### Check for missing values

```
In [3]: df.isna().any()
```

	age	anaemia	creatinine_phosphokinase	diabetes	ejection_fraction	high_blood_pressure	platelets	serum_creatinine	sex	smoking	time	DEATH_EVENT
Out[3]:	age	anaemia	creatinine_phosphokinase	diabetes	ejection_fraction	high_blood_pressure	platelets	serum_creatinine	sex	smoking	time	DEATH_EVENT
	dtype: bool	dtype: bool	dtype: bool	dtype: bool	dtype: bool	dtype: bool	dtype: bool	dtype: bool	dtype: bool	dtype: bool	dtype: bool	dtype: bool

```
In [4]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 299 entries, 0 to 298
Data columns (total 13 columns):
 #   Column              Non-Null Count  Dtype
---  -
 0   age                  299 non-null   float64
 1   anaemia              299 non-null   int64
 2   creatinine_phosphokinase  299 non-null   int64
 3   diabetes             299 non-null   int64
 4   ejection_fraction    299 non-null   int64
 5   high_blood_pressure   299 non-null   int64
 6   platelets            299 non-null   float64
 7   serum_creatinine      299 non-null   float64
 8   serum_sodium          299 non-null   int64
 9   sex                  299 non-null   int64
10  smoking              299 non-null   int64
11  time                 299 non-null   int64
12  DEATH_EVENT          299 non-null   int64
dtypes: float64(3), int64(10)
memory usage: 38.5 KB
```

```
In [5]: df.describe()
```

	age	anaemia	creatinine_phosphokinase	diabetes	ejection_fraction	high_blood_pressure	platelets	serum_creatinine	sex	smoking	time	DEATH_EVENT
count	299.000000	299.000000	299.000000	299.000000	299.000000	299.000000	299.000000	299.000000	299.000000	299.000000	299.000000	299.000000
mean	60.833893	0.431438	581.839465	0.418060	38.038362	0.351171	263358.03	119.94809	0.496107	0.702881	11834841	0.478136
std	11.894809	0.496107	970.287881	0.494067	14.000000	0.000000	0.000000	25100.000000	0.500000	0.500000	212500.000000	0.500000
min	40.000000	0.000000	23.000000	0.000000	0.000000	0.000000	0.000000	25100.000000	0.000000	0.000000	212500.000000	0.000000
25%	51.000000	0.000000	116.500000	0.000000	30.000000	0.000000	0.000000	212500.000000	0.000000	0.000000	212500.000000	0.000000
50%	60.833893	0.431438	581.839465	0.418060	38.038362	0.351171	263358.03	119.94809	0.496107	0.702881	11834841	0.478136
75%	70.000000	1.000000	970.287881	0.494067	14.000000	0.000000	0.000000	25100.000000	0.500000	0.500000	212500.000000	0.500000
max	95.000000	1.000000	7861.000000	1.000000	80.000000	1.000000	850000.000000	850000.000000	1.000000	1.000000	850000.000000	1.000000

#### Apply Standard Scaler to entire dataframe - except target feature

Target Feature of 1 == Died

Target Feature of 0 == Lived

```
In [6]: X_data = df.iloc[:,12:]
print("Target Feature:")
print(X_data.columns)
print("\n")

ss = StandardScaler()
df = pd.DataFrame(ss.fit_transform(df), columns=df.columns)
df.head()
```

Independent Features:

```
Index(['age', 'anaemia', 'creatinine_phosphokinase', 'diabetes',
       'ejection_fraction', 'high_blood_pressure', 'platelets',
       'serum_creatinine', 'serum_sodium', 'sex', 'smoking', 'time'],
      dtype='object')
```

#### Correlation Heatmap

```
In [7]: mask = np.triu(np.ones_like(df.corr(), dtype=bool))
plt.figure(figsize=(16, 10))
heatmap = sns.heatmap(df.corr(), mask=mask, vmin=-1, vmax=1, annot=True)
heatmap.set_title('Heart Failure Correlation Heatmap', fontdict={'fontsize':18, 'pad':12})
```

	age	anaemia	creatinine_phosphokinase	diabetes	ejection_fraction	high_blood_pressure	platelets	serum_creatinine	serum_sodium	sex	smoking	time	DEATH_EVENT
age	1.000	0.088	0.002	-0.13	0.032	-0.044	-0.049	0.006	-0.046	-0.095	0.08	-0.15	-0.75
anaemia	0.088	1.000	0.002	-0.13	0.032	-0.044	-0.049	0.006	-0.046	-0.095	0.08	-0.15	-0.75
creatinine_phosphokinase	0.002	0.002	1.000	0.002	0.032	-0.044	-0.049	0.006	-0.046	-0.095	0.08	-0.15	-0.75
diabetes	-0.13	-0.13	0.002	1.000	0.032	-0.044	-0.049	0.006	-0.046	-0.095	0.08	-0.15	-0.75
ejection_fraction	0.032	0.032	0.032	0.032	1.000	-0.044	-0.049	0.006	-0.046	-0.095	0.08	-0.15	-0.75
high_blood_pressure	-0.044	-0.044	-0.044	-0.044	-0.044	1.000	-0.049	0.006	-0.046	-0.095	0.08	-0.15	-0.75
platelets	-0.049	-0.049	-0.049	-0.049	-0.049	-0.049	1.000	0.006	-0.046	-0.095	0.08	-0.15	-0.75
serum_creatinine	0.006	0.006	0.006	0.006	0.006	0.006	0.006	1.000	-0.046	-0.095	0.08	-0.15	-0.75
serum_sodium	-0.046	-0.046	-0.046	-0.046	-0.046	-0.046	-0.046	-0.046	1.000	-0.095	0.08	-0.15	-0.75
sex	-0.095	-0.095	-0.095	-0.095	-0.095	-0.095	-0.095	-0.095	-0.095	1.000	0.08	-0.15	-0.75
smoking	0.08	0.08	0.08	0.08	0.08	0.08	0.08	0.08	0.08	0.08	1.000	-0.15	-0.75
time	-0.15	-0.15	-0.15	-0.15	-0.15	-0.15	-0.15	-0.15	-0.15	-0.15	-0.15	1.000	-0.75
DEATH_EVENT	-0.75	-0.75	-0.75	-0.75	-0.75	-0.75	-0.75	-0.75	-0.75	-0.75	-0.75	-0.75	1.000

#### Correlation with Target (DEATH\_EVENT)

```
In [8]: df.corrwith(df['DEATH_EVENT'])
```

	age	anaemia	creatinine_phosphokinase	diabetes	ejection_fraction	high_blood_pressure	platelets	serum_creatinine	serum_sodium	sex	smoking	time	DEATH_EVENT
age	0.253729	0.066270	0.062728	-0.400143	-0.268603	-0.049329	-0.079351	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329
anaemia	0.066270	0.062728	-0.400143	-0.268603	-0.049329	-0.079351	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329
creatinine_phosphokinase	0.062728	-0.400143	-0.268603	-0.049329	-0.079351	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329
diabetes	-0.400143	-0.268603	-0.049329	-0.079351	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329
ejection_fraction	-0.268603	-0.049329	-0.079351	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329
high_blood_pressure	-0.049329	-0.079351	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329
platelets	-0.079351	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329
serum_creatinine	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329
serum_sodium	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329
sex	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329
smoking	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329
time	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329
DEATH_EVENT	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329	-0.049329

```
In [9]: plt.figure(figsize=(8, 10))
heatmap = sns.heatmap(df.corr(), mask=mask, vmin=-1, vmax=1, annot=True)
heatmap.set_title('Features Correlating with (DEATH_EVENT)', fontdict={'fontsize':16, 'pad':8})
```

	age	anaemia	creatinine_phosphokinase	diabetes	ejection_fraction	high_blood_pressure	platelets	serum_creatinine	serum_sodium	sex	smoking	time	DEATH_EVENT
age	1.000	0.088	0.002	-0.13	0.032	-0.044	-0.049	0.006	-0.046	-0.095	0.08	-0.15	-0.75
anaemia	0.088	1.000	0.002	-0.13	0.032	-0.044	-0.049	0.006	-0.046	-0.095	0.08	-0.15	-0.75
creatinine_phosphokinase	0.002	0.002	1.000	0.002	0.032	-0.044	-0.049	0.006	-0.046	-0.095	0.08	-0.15	-0.75
diabetes	-0.13	-0.13	0.002	1.000	0.032	-0.044	-0.049	0.006	-0.046	-0.095	0.08	-0.15	-0.75
ejection_fraction	0.032	0.032	0.032	0.032	1.000	-0.044	-0.049	0.006	-0.046	-0.095	0.08	-0.15	-0.75
high_blood_pressure	-0.044	-0.044	-0.044	-0.044	-0.044	1.000	-0.049	0.006	-0.046	-0.095	0.08	-0.15	-0.75
platelets	-0.049	-0.049	-0.049	-0.049	-0.049	-0.049	1.000	0.006	-0.046	-0.095	0.08	-0.15	-0.75
serum_creatinine	0.006	0.006	0.006	0.006	0.006	0.006	0.006	1.000	-0.046	-0.095	0.08	-0.15	-0.75
serum_sodium	-0.046	-0.046	-0.046	-0.046	-0.046	-0.046	-0.046	-0.046	1.000	-0.095	0.08	-0.15	-0.75
sex	-0.095	-0.095	-0.095	-0.095	-0.095	-0.095	-0.095	-0.095	-0.095	1.000	0.08	-0.15	-0.75
smoking	0.08	0.08	0.08	0.08	0.08	0.08	0.08	0.08	0.08	0.08	1.000	-0.15	-0.75
time	-0.15	-0.15	-0.15	-0.15	-0.15	-0.15	-0.15	-0.15	-0.15	-0.15	-0.15	1.000	-0.75
DEATH_EVENT	-0.75	-0.75	-0.75	-0.75	-0.75	-0.75	-0.75	-0.75	-0.75	-0.75	-0.75	-0.75	1.000

#### Part 2 - Machine Learning Models

```
In [10]: # Train/Test Split - 30% for test
#
X_train, X_test, y_train, y_test = train_test_split(X_data, y_data, test_size=0.3, random_state=1)

y_train = y_train.values
y_test = y_test.values

y_train = y_train.reshape(1, -1)
y_test = y_test.reshape(1, -1)

y_train = y_train.flatten()
y_test = y_test.flatten()
```

#### Logistic Regression - Grid Search

```
In [11]: parameters = [
    {'solver': ['lbfgs'],
     'penalty': ['l2', 'none']},
    {'solver': ['liblinear'],
     'penalty': ['l1', 'l2']},
    {'solver': ['newton-cg'],
     'penalty': ['l2', 'none']},
    {'solver': ['sag'],
     'penalty': ['l2', 'none']},
    {'solver': ['saga'],
     'penalty': ['l1', 'l2', 'none']},
    {'solver': ['elasticnet'],
     'l1_ratio': [1]}
]

lr_classifier = GridSearchCV(estimator=LogisticRegression(random_state=0), param_grid=parameters,
                             cv=5, scoring='roc_auc', return_train_score=True, n_jobs=-1)

lr_classifier.fit(X_train, y_train)

# Show best parameters
print("\nBest Parameters:")
print(lr_classifier.best_params_)
print(lr_classifier.best_estimator_)

predict = lr_classifier.predict(X_test)
accuracy_lr_classifier = np.mean(predict == y_test)

print("Logistic Regression Accuracy: " + str(accuracy_lr_classifier))
print("ROC-AUC Score: " + str(roc_auc_score(y_test, predict, average=None)))

Fitting 5 folds for each of 12 candidates, totalling 60 fits

Best Parameters:
{'penalty': 'l2', 'solver': 'lbfgs'}

Logistic Regression Accuracy: 0.8333333333333334
ROC-AUC Score: 0.768629807923077
```

#### Decision Tree Classifier - Grid Search

```
In [12]: parameters = [
    {'criterion': ['gini', 'entropy', 'log_loss'],
     'splitter': ['best', 'random'],
     'class_weight': [(0, 1, 2), 'balanced']},
    {'criterion': ['gini', 'entropy', 'log_loss'],
     'splitter': ['best', 'random'],
     'class_weight': [(0, 1, 2), 'balanced', 'balanced_subsample']}
]

dtree_classifier = GridSearchCV(estimator=DecisionTreeClassifier(random_state=0), param_grid=parameters,
                                cv=5, scoring='roc_auc', return_train_score=True, n_jobs=-1)

dtree_classifier.fit(X_train, y_train)

# Show best parameters
print("\nBest Parameters:")
print(dtree_classifier.best_params_)
print(dtree_classifier.best_estimator_)

# Run prediction with best parameters
predict = dtree_classifier.predict(X_test)
accuracy_dtree_classifier = np.mean(predict == y_test)

print("Decision Tree Accuracy: " + str(accuracy_dtree_classifier))
print("ROC-AUC Score: " + str(roc_auc_score(y_test, predict, average=None)))

Fitting 5 folds for each of 12 candidates, totalling 60 fits

Best Parameters:
{'class_weight': (0, 1, 2), 'criterion': 'entropy', 'splitter': 'best'}

Decision Tree Accuracy: 0.8333333333333334
ROC-AUC Score: 0.7914663461538461
```

#### Random Forest - Grid Search

```
In [13]: parameters = [
    {'n_estimators': [10, 100, 200, 500],
     'criterion': ['gini', 'entropy', 'log_loss'],
     'bootstrap': [True, False],
     'class_weight': [(0, 1, 2), 'balanced', 'balanced_subsample']}
]

rf_classifier = GridSearchCV(estimator=RandomForestClassifier(random_state=0), param_grid=parameters,
                             cv=5, scoring='roc_auc', return_train_score=True, n_jobs=-1)

rf_classifier.fit(X_train, y_train)

# Show best parameters
print("\nBest Parameters:")
print(rf_classifier.best_params_)
print(rf_classifier.best_estimator_)

# Run prediction with best parameters
predict = rf_classifier.predict(X_test)
accuracy_rf_classifier = np.mean(predict == y_test)

print("Random Forest Accuracy: " + str(accuracy_rf_classifier))
print("ROC-AUC Score: " + str(roc_auc_score(y_test, predict, average=None)))

Fitting 5 folds for each of 72 candidates, totalling 360 fits

Best Parameters:
{'bootstrap': True, 'class_weight': 'balanced_subsample', 'criterion': 'gini', 'n_estimators': 100}

Decision Tree Accuracy: 0.8888888888888888
ROC-AUC Score: 0.8419471153846154
```

#### XGBoost - Grid Search

```
In [14]: parameters = [
    {'n_estimators': [5, 10, 15, 20, 50, 100, 200, 500],
     'booster': ['gbtree', 'gblinear', 'dart'],
     'eval_metric': ['logloss', 'auc', 'rmse'],
     'learning_rate': [0.5, 0.1, 0.05, 0.01]},
    {'n_estimators': [5, 10, 15, 20, 50, 100, 200, 500],
     'booster': ['gbtree', 'gblinear', 'dart'],
     'eval_metric': ['logloss', 'auc', 'rmse'],
     'learning_rate': [0.5, 0.1, 0.05, 0.01]}
]

xgb_classifier = GridSearchCV(estimator=XGBClassifier(random_state=0), param_grid=parameters,
                             cv=5, scoring='roc_auc', return_train_score=True, n_jobs=-1)

xgb_classifier.fit(X_train, y_train)

# Show best parameters
print("\nBest Parameters:")
print(xgb_classifier.best_params_)
print(xgb_classifier.best_estimator_)

# Run prediction with best parameters
predict = xgb_classifier.predict(X_test)
accuracy_xgb_classifier = np.mean(predict == y_test)

print("XGBoost Accuracy: " + str(accuracy_xgb_classifier))
print("ROC-AUC Score: " + str(roc_auc_score(y_test, predict, average=None)))

Fitting 5 folds for each of 288 candidates, totalling 1440 fits

Best Parameters:
{'booster': 'gbtree', 'eval_metric': 'logloss', 'learning_rate': 0.05, 'n_estimators': 200}

Decision Tree Accuracy: 0.8888888888888888
ROC-AUC Score: 0.8647836538461539
```

#### Part 3A

```
In [15]: # Fine one positive row (death) and one negative row (lived) in X_test data
#
y_1 = 0
y_0 = 0
X_sample_1 = None
X_sample_0 = None

for i in range(y_test.size):
    if y_test[i] == 1:
        X_sample_1 = X_test.iloc[i, :]
        y_1 = 1
        break

for i in range(y_test.size):
    if y_test[i] == 0:
        X_sample_0 = X_test.iloc
```



