

Source code

Project work Implementation code (Source code of your project)

Preparing the tools

We're going to use pandas, Matplotlib and NumPy for data analysis and manipulation.

In [32]:

```
# Import all the tools we need

# Regular EDA (exploratory data analysis) and plotting libraries
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

# we want our plots to appear inside the notebook
%matplotlib inline

# Models from Scikit-Learn
from sklearn.linear_model import LogisticRegression
from sklearn.neighbors import KNeighborsClassifier
from sklearn.ensemble import RandomForestClassifier

# Model Evaluations
from sklearn.model_selection import train_test_split, cross_val_score
from sklearn.model_selection import RandomizedSearchCV, GridSearchCV
from sklearn.metrics import confusion_matrix, classification_report
from sklearn.metrics import precision_score, recall_score, f1_score
from sklearn.metrics import plot_roc_curve
```

Load data

In [4]:

```
df = pd.read_csv("heart-disease.csv")
df.shape # (rows, columns)
```

Data Exploration (exploratory data analysis or EDA)

The goal here is to find out more about the data and become a subject matter expert on the dataset you're working with.

1. What question(s) are you trying to solve?
2. What kind of data do we have and how do we treat different types?
3. What's missing from the data and how do you deal with it?
4. Where are the outliers and why should you care about them?
5. How can you add, change or remove features to get more out of your data?

In [5]:

```
df.head()
```

In [6]:

```
df.tail()
```

In [7]:

```
# Let's find out how many of each class there  
df["target"].value_counts()
```

In [8]:

```
df["target"].value_counts().plot(kind="bar", color=["salmon", "lightblue"]);
```

In [10]:

```
df.info()
```

In [11]:

```
# Are there any missing values?  
df.isna().sum()
```

In [12]:

```
df.describe()
```

Heart Disease Frequency according to Sex

In [14]:

```
df.sex.value_counts()
```

In [15]:

```
# Compare target column with sex column  
pd.crosstab(df.target, df.sex)
```

In [19]:

```
# Create a plot of crosstab  
pd.crosstab(df.target, df.sex).plot(kind="bar",  
                                     figsize=(10, 6),  
                                     color=["salmon", "lightblue"])  
  
plt.title("Heart Disease Frequency for Sex")  
plt.xlabel("0 = No Disease, 1 = Disease")  
plt.ylabel("Amount")  
plt.legend(["Female", "Male"]);  
plt.xticks(rotation=0);
```

Age vs. Max Heart Rate for Heart Disease

In [26]:

```

# Create another figure
plt.figure(figsize=(10, 6))

# Scatter with positive examples
plt.scatter(df.age[df.target==1],
            df.thalach[df.target==1],
            c="salmon")

# Scatter with negative examples
plt.scatter(df.age[df.target==0],
            df.thalach[df.target==0],
            c="lightblue")

# Add some helpful info
plt.title("Heart Disease in function of Age and Max Heart Rate")
plt.xlabel("Age")
plt.ylabel("Max Heart Rate")
plt.legend(["Disease", "No Disease"]);

```

In [27]:

```

# Check the distribution of the age column with a histogram
df.age.plot.hist();

```

Heart Disease Frequency per Chest Pain Type

1. cp - chest pain type
 - 0: Typical angina: chest pain related decrease blood supply to the heart
 - 1: Atypical angina: chest pain not related to heart
 - 2: Non-anginal pain: typically esophageal spasms (non heart related)
 - 3: Asymptomatic: chest pain not showing signs of disease

In [28]:

```
pd.crosstab(df.cp, df.target)
```

In [30]:

```

# Make the crosstab more visual
pd.crosstab(df.cp, df.target).plot(kind="bar",
                                   figsize=(10, 6),
                                   color=["salmon", "lightblue"])

# Add some communication
plt.title("Heart Disease Frequency Per Chest Pain Type")
plt.xlabel("Chest Pain Type")
plt.ylabel("Amount")
plt.legend(["No Disease", "Disease"])
plt.xticks(rotation=0);

```

```
df.head()
```

In [33]:

```
# Make a correlation matrix  
df.corr()
```

In [46]:

```
# Let's make our correlation matrix a little prettier  
corr_matrix = df.corr()  
fig, ax = plt.subplots(figsize=(15, 10))  
ax = sns.heatmap(corr_matrix,  
                  annot=True,  
                  linewidths=0.5,  
                  fmt=".2f",  
                  cmap="YlGnBu");  
bottom, top = ax.get_ylim()  
ax.set_ylim(bottom + 0.5, top - 0.5)
```

In [49]:

Modelling

In [50]:

```
df.head()
```

In [51]:

```
# Split data into X and y  
X = df.drop("target", axis=1)
```

```
y = df["target"]
```

In [52]:

```
X
```

In [53]:

```
y
```

In [54]:

```
# Split data into train and test sets  
np.random.seed(42)
```

```
# Split into train & test set  
X_train, X_test, y_train, y_test = train_test_split(X,  
                                                    y,  
                                                    test_size=0.2)
```

In [55]:

```
X_train
```

In [56]:

```
y_train, len(y_train)
```

In [57]:

```
# Put models in a dictionary
models = {"Logistic Regression": LogisticRegression(),
          "KNN": KNeighborsClassifier(),
          "Random Forest": RandomForestClassifier()}

# Create a function to fit and score models
def fit_and_score(models, X_train, X_test, y_train, y_test):
    """
    Fits and evaluates given machine learning models.
    models : a dict of differetn Scikit-Learn machine learning models
    X_train : training data (no labels)
    X_test : testing data (no labels)
    y_train : training labels
    y_test : test labels
    """
    # Set random seed
    np.random.seed(42)
    # Make a dictionary to keep model scores
    model_scores = {}
    # Loop through models
    for name, model in models.items():
        # Fit the model to the data
        model.fit(X_train, y_train)
        # Evaluate the model and append its score to model_scores
        model_scores[name] = model.score(X_test, y_test)
    return model_scores
```

In [58]:

```
model_scores = fit_and_score(models=models,
                              X_train=X_train,
                              X_test=X_test,
                              y_train=y_train,
                              y_test=y_test)
```

```
model_scores
```

Model Comparison

In [61]:

```
model_compare = pd.DataFrame(model_scores, index=["accuracy"])
model_compare.T.plot.bar();
```

Hyperparameter tuning (by hand)

In [63]:

```
# Let's tune KNN

train_scores = []
test_scores = []

# Create a list of different values for n_neighbors
neighbors = range(1, 21)

# Setup KNN instance
knn = KNeighborsClassifier()

# Loop through different n_neighbors
for i in neighbors:
    knn.set_params(n_neighbors=i)

    # Fit the algorithm
    knn.fit(X_train, y_train)

    # Update the training scores list
    train_scores.append(knn.score(X_train, y_train))

    # Update the test scores list
    test_scores.append(knn.score(X_test, y_test))
```

In [64]:

```
train_scores
```

In [66]:

```
test_scores
```

In [68]:

```
plt.plot(neighbors, train_scores, label="Train score")
plt.plot(neighbors, test_scores, label="Test score")
plt.xticks(np.arange(1, 21, 1))
plt.xlabel("Number of neighbors")
plt.ylabel("Model score")
plt.legend()

print(f"Maximum KNN score on the test data: {max(test_scores)*100:.2f}%")
```

In [78]:

```
# Create a hyperparameter grid for LogisticRegression
log_reg_grid = {"C": np.logspace(-4, 4, 20),
```

```

        "solver": ["liblinear"]}

# Create a hyperparameter grid for RandomForestClassifier
rf_grid = {"n_estimators": np.arange(10, 1000, 50),
          "max_depth": [None, 3, 5, 10],
          "min_samples_split": np.arange(2, 20, 2),
          "min_samples_leaf": np.arange(1, 20, 2)}

```

Now we've got hyperparameter grids setup for each of our models, let's tune them using RandomizedSearchCV...

In [74]:

```

# Tune LogisticRegression

np.random.seed(42)

# Setup random hyperparameter search for LogisticRegression
rs_log_reg = RandomizedSearchCV(LogisticRegression(),
                                param_distributions=log_reg_grid,
                                cv=5,
                                n_iter=20,
                                verbose=True)

# Fit random hyperparameter search model for LogisticRegression
rs_log_reg.fit(X_train, y_train)

```

In [75]:

```
rs_log_reg.best_params_
```

In [76]:

```
rs_log_reg.score(X_test, y_test)
```

Now we've tuned LogisticRegression(), let's do the same for RandomForestClassifier()...

In [79]:

```

# Setup random seed
np.random.seed(42)

# Setup random hyperparameter search for RandomForestClassifier
rs_rf = RandomizedSearchCV(RandomForestClassifier(),
                            param_distributions=rf_grid,
                            cv=5,
                            n_iter=20,
                            verbose=True)

# Fit random hyperparameter search model for RandomForestClassifier()
rs_rf.fit(X_train, y_train)

```

In [80]:

```

# Find the best hyperparameters
rs_rf.best_params_

```

In [81]:

```
# Evaluate the randomized search RandomForestClassifier model
rs_rf.score(X_test, y_test)
```

Hyperparameter Tuning with GridSearchCV

Since our LogisticRegression model provides the best scores so far, we'll try and improve them again using GridSearchCV...

In [83]:

```
# Different hyperparameters for our LogisticRegression model
log_reg_grid = {"C": np.logspace(-4, 4, 30),
                "solver": ["liblinear"]}
```

```
# Setup grid hyperparameter search for LogisticRegression
gs_log_reg = GridSearchCV(LogisticRegression(),
                           param_grid=log_reg_grid,
                           cv=5,
                           verbose=True)
```

```
# Fit grid hyperparameter search model
gs_log_reg.fit(X_train, y_train);
```

In [84]:

```
# Check the best hyperparameters
gs_log_reg.best_params_
```

In [85]:

```
# Evaluate the grid search LogisticRegression model
gs_log_reg.score(X_test, y_test)
```

Evaluating our tuned machine learning classifier, beyond accuracy

- ROC curve and AUC score
- Confusion matrix
- Classification report
- Precision
- Recall
- F1-score

... and it would be great if cross-validation was used where possible.

To make comparisons and evaluate our trained model, first we need to make predictions.

In [87]:

```
# Make predictions with tuned model
y_preds = gs_log_reg.predict(X_test)
```

In [88]:

```
y_preds
```

In [89]:

```
y_test
```

In [90]:

In [91]:

In [96]:

```
plot_conf_mat(y_test, y_preds)
```

In [97]:

Calculate evaluation metrics using cross-validation

In [98]:

In [99]:

In [100]:

cv acc

In [102]:

```
cv_acc = np.mean(cv_acc)
cv_acc
```

In [103]:

```
# Cross-validated precision
cv_precision = cross_val_score(clf,
                                X,
                                Y,
                                cv=5,
                                scoring="precision")
cv_precision=np.mean(cv_precision)
cv_precision
```

In [104]:

```
# Cross-validated recall
cv_recall = cross_val_score(clf,
                             X,
                             Y,
                             cv=5,
                             scoring="recall")
cv_recall = np.mean(cv_recall)
cv_recall
```

In [105]:

```
# Cross-validated f1-score
cv_f1 = cross_val_score(clf,
                         X,
                         Y,
                         cv=5,
                         scoring="f1")
cv_f1 = np.mean(cv_f1)
cv_f1
```

In [107]:

```
# Visualize cross-validated metrics
cv_metrics = pd.DataFrame({"Accuracy": cv_acc,
                           "Precision": cv_precision,
                           "Recall": cv_recall,
                           "F1": cv_f1},
                           index=[0])

cv_metrics.T.plot.bar(title="Cross-validated classification metrics",
                      legend=False);
```

Feature Importance

Feature importance is another as asking, "which features contributed most to the outcomes of the model and how did they contribute?"

Finding feature importance is different for each machine learning model. One way to find feature importance is to search for "(MODEL NAME) feature importance".

Let's find the feature importance for our LogisticRegression model...

In [110]:

```
# Fit an instance of LogisticRegression
clf = LogisticRegression(C=0.20433597178569418,
                        solver="liblinear")
```

```
clf.fit(X_train, y_train);
```

In [111]:

```
# Check coef_
clf.coef_
```

In [114]:

```
df.head()
```

In [113]:

```
# Match coef's of features to columns
feature_dict = dict(zip(df.columns, list(clf.coef_[0])))
feature_dict
```

In [115]:

```
# Visualize feature importance
feature_df = pd.DataFrame(feature_dict, index=[0])
feature_df.T.plot.bar(title="Feature Importance", legend=False);
```

In [117]:

```
pd.crosstab(df["sex"], df["target"])
```

In [119]:

```
pd.crosstab(df["slope"], df["target"])
```