# Project

June 24, 2023

# 1 Predicting Heart Disease using Machine Learning

This notebook will introduce some foundation machine learning and data science concepts by exploring the problem of heart disease **classification**.

It is intended to be an end-to-end example of what a data science and machine learning **proof of concept** might look like.

## 1.1 What is classification?

Classification involves deciding whether a sample is part of one class or another (**single-class classification**). If there are multiple class options, it's referred to as **multi-class classification**.

## 1.2 What we'll end up with

Since we already have a dataset, we'll approach the problem with the following machine learning modelling framework.

More specifically, we'll look at the following topics.

- **Exploratory data analysis (EDA)** - the process of going through a dataset and finding out more about it.
- **Model training** - create model(s) to learn to predict a target variable based on other variables.
- **Model evaluation** - evaluating a models predictions using problem-specific evaluation metrics.
- **Model comparison** - comparing several different models to find the best one.
- **Model fine-tuning** - once we've found a good model, how can we improve it?
- **Feature importance** - since we're predicting the presence of heart disease, are there some things which are more important for prediction?
- **Cross-validation** - if we do build a good model, can we be sure it will work on unseen data?
- **Reporting what we've found** - if we had to present our work, what would we show someone?

## 1.3 1. Problem Definition

In our case, the problem we will be exploring is **binary classification** (a sample can only be one of two things).

This is because we're going to be using a number of differnet **features** (pieces of information) about a person to predict whether they have heart disease or not.

In a statement,

> Given clinical parameters about a patient, can we predict whether or not they have heart disease?

## 1.4 2. Data

What you'll want to do here is dive into the data your problem definition is based on. This may involve, sourcing, defining different parameters, talking to experts about it and finding out what you should expect.

The original data came from the Cleveland database from UCI Machine Learning Repository.

Howevever, we've downloaded it in a formatted way from Kaggle.

The original database contains 76 attributes, but here only 14 attributes will be used. **Attributes** (also called **features**) are the variables what we'll use to predict our **target variable**.

Attributes and features are also referred to as **independent variables** and a target variable can be referred to as a **dependent variable**.

> We use the independent variables to predict our dependent variable.

Or in our case, the independent variables are a patients different medical attributes and the dependent variable is whether or not they have heart disease.

## 1.5 3. Evaluation

The evaluation metric is something you might define at the start of a project.

Since machine learning is very experimental, you might say something like,

> If we can reach 95% accuracy at predicting whether or not a patient has heart disease during the proof of concept, we'll pursure this project.

The reason this is helpful is it provides a rough goal for a machine learning engineer or data scientist to work towards.

However, due to the nature of experimentation, the evaluation metric may change over time.

## 1.6 4. Features

Features are different parts of the data. During this step, you'll want to start finding out what you can about the data.

One of the most common ways to do this, is to create a **data dictionary**.

### 1.6.1 Heart Disease Data Dictionary

A data dictionary describes the data you're dealing with. Not all datasets come with them so this is where you may have to do your research or ask a **subject matter expert** (someone who knows about the data) for more.

The following are the features we'll use to predict our target variable (heart disease or no heart disease).

1. age - age in years
2. sex - (1 = male; 0 = female)
3. 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
4. trestbps - resting blood pressure (in mm Hg on admission to the hospital)
   - anything above 130-140 is typically cause for concern
5. chol - serum cholestoral in mg/dl
   - serum = LDL + HDL + .2 * triglycerides
   - above 200 is cause for concern
6. fbs - (fasting blood sugar > 120 mg/dl) (1 = true; 0 = false)
   - '>126' mg/dL signals diabetes
7. restecg - resting electrocardiographic results
   - 0: Nothing to note
   - 1: ST-T Wave abnormality
     - can range from mild symptoms to severe problems
     - signals non-normal heart beat
   - 2: Possible or definite left ventricular hypertrophy
     - Enlarged heart's main pumping chamber
8. thalach - maximum heart rate achieved
9. exang - exercise induced angina (1 = yes; 0 = no)
10. oldpeak - ST depression induced by exercise relative to rest
    - looks at stress of heart during excercise
    - unhealthy heart will stress more
11. slope - the slope of the peak exercise ST segment
    - 0: Upsloping: better heart rate with excercise (uncommon)
    - 1: Flatsloping: minimal change (typical healthy heart)
    - 2: Downslopins: signs of unhealthy heart
12. ca - number of major vessels (0-3) colored by flourosopy
    - colored vessel means the doctor can see the blood passing through
    - the more blood movement the better (no clots)
13. thal - thalium stress result
    - 1,3: normal
    - 6: fixed defect: used to be defect but ok now
    - 7: reversable defect: no proper blood movement when excercising
14. target - have disease or not (1=yes, 0=no) (= the predicted attribute)

**Note:** No personal identifiable information (PPI) can be found in the dataset.

It's a good idea to save these to a Python dictionary or in an external file, so we can look at them later without coming back here.

```
[3]:  # Regular EDA and plotting libraries
      import numpy as np # np is short for numpy
      import pandas as pd # pandas is so commonly used, it's shortened to pd
      import matplotlib.pyplot as plt
```

```python
import seaborn as sns # seaborn gets shortened to sns

# We want our plots to appear in the notebook
%matplotlib inline

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

## Model evaluators
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
```

[4]:
```python
df = pd.read_csv("../data/heart-disease.csv") # 'DataFrame' shortened to 'df'
df.shape # (rows, columns)
```

[4]: (303, 14)

[5]:
```python
# Let's check the top 5 rows of our dataframe
df.head()
```

[5]:
```
    age  sex  cp  trestbps  chol  fbs  restecg  thalach  exang  oldpeak  slope  \
0    63    1   3       145   233    1        0      150      0      2.3      0
1    37    1   2       130   250    0        1      187      0      3.5      0
2    41    0   1       130   204    0        0      172      0      1.4      2
3    56    1   1       120   236    0        1      178      0      0.8      2
4    57    0   0       120   354    0        1      163      1      0.6      2

   ca  thal  target
0   0     1       1
1   0     2       1
2   0     2       1
3   0     2       1
4   0     2       1
```
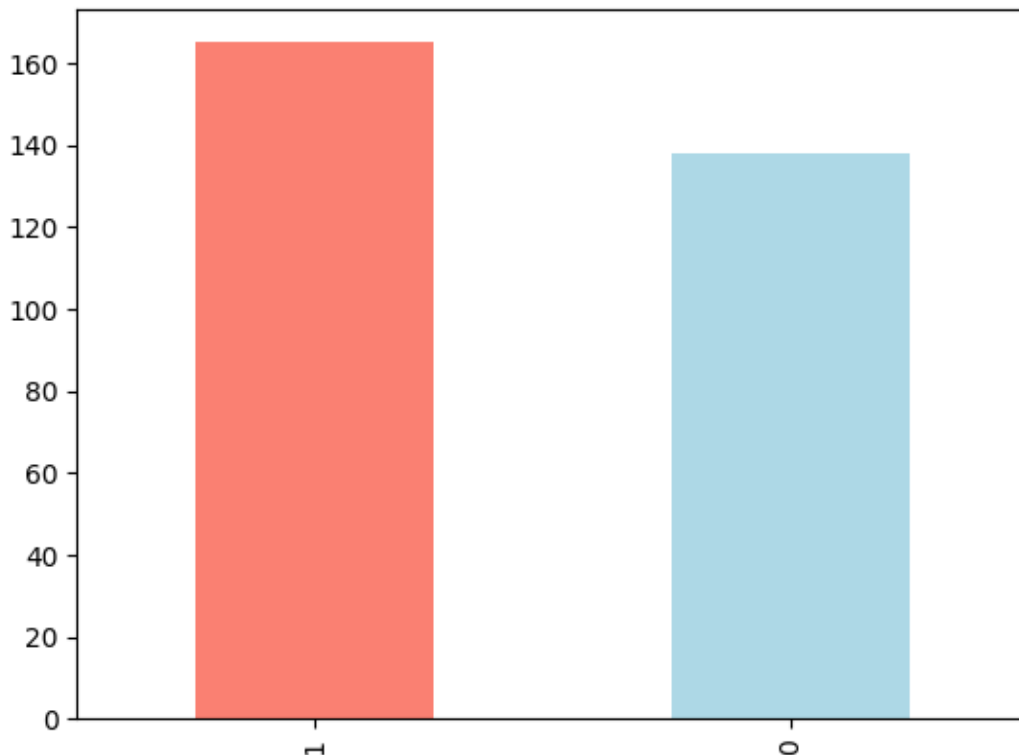
[6]:
```python
# Let's see how many positive (1) and negative (0) samples we have in our
 ↪dataframe
df.target.value_counts()
```

[6]:
```
1    165
0    138
Name: target, dtype: int64
```

```
[7]:  # Normalized value counts
      df.target.value_counts(normalize=True)
```

```
[7]:  1    0.544554
      0    0.455446
      Name: target, dtype: float64
```

```
[8]:  # Plot the value counts with a bar graph
      df.target.value_counts().plot(kind="bar", color=["salmon", "lightblue"]);
```



```
[9]:  df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 303 entries, 0 to 302
Data columns (total 14 columns):
 #   Column    Non-Null Count  Dtype
---  ------    --------------  -----
 0   age       303 non-null    int64
 1   sex       303 non-null    int64
 2   cp        303 non-null    int64
 3   trestbps  303 non-null    int64
 4   chol      303 non-null    int64
 5   fbs       303 non-null    int64
```

5

```
6    restecg    303 non-null    int64
7    thalach    303 non-null    int64
8    exang      303 non-null    int64
9    oldpeak    303 non-null    float64
10   slope      303 non-null    int64
11   ca         303 non-null    int64
12   thal       303 non-null    int64
13   target     303 non-null    int64
dtypes: float64(1), int64(13)
memory usage: 33.3 KB
```

[10]: `df.describe()`

[10]:
|       | age        | sex        | cp         | trestbps   | chol       | fbs        |
|-------|------------|------------|------------|------------|------------|------------|
| count | 303.000000 | 303.000000 | 303.000000 | 303.000000 | 303.000000 | 303.000000 |
| mean  | 54.366337  | 0.683168   | 0.966997   | 131.623762 | 246.264026 | 0.148515   |
| std   | 9.082101   | 0.466011   | 1.032052   | 17.538143  | 51.830751  | 0.356198   |
| min   | 29.000000  | 0.000000   | 0.000000   | 94.000000  | 126.000000 | 0.000000   |
| 25%   | 47.500000  | 0.000000   | 0.000000   | 120.000000 | 211.000000 | 0.000000   |
| 50%   | 55.000000  | 1.000000   | 1.000000   | 130.000000 | 240.000000 | 0.000000   |
| 75%   | 61.000000  | 1.000000   | 2.000000   | 140.000000 | 274.500000 | 0.000000   |
| max   | 77.000000  | 1.000000   | 3.000000   | 200.000000 | 564.000000 | 1.000000   |

|       | restecg    | thalach    | exang      | oldpeak    | slope      | ca         |
|-------|------------|------------|------------|------------|------------|------------|
| count | 303.000000 | 303.000000 | 303.000000 | 303.000000 | 303.000000 | 303.000000 |
| mean  | 0.528053   | 149.646865 | 0.326733   | 1.039604   | 1.399340   | 0.729373   |
| std   | 0.525860   | 22.905161  | 0.469794   | 1.161075   | 0.616226   | 1.022606   |
| min   | 0.000000   | 71.000000  | 0.000000   | 0.000000   | 0.000000   | 0.000000   |
| 25%   | 0.000000   | 133.500000 | 0.000000   | 0.000000   | 1.000000   | 0.000000   |
| 50%   | 1.000000   | 153.000000 | 0.000000   | 0.800000   | 1.000000   | 0.000000   |
| 75%   | 1.000000   | 166.000000 | 1.000000   | 1.600000   | 2.000000   | 1.000000   |
| max   | 2.000000   | 202.000000 | 1.000000   | 6.200000   | 2.000000   | 4.000000   |

|       | thal       | target     |
|-------|------------|------------|
| count | 303.000000 | 303.000000 |
| mean  | 2.313531   | 0.544554   |
| std   | 0.612277   | 0.498835   |
| min   | 0.000000   | 0.000000   |
| 25%   | 2.000000   | 0.000000   |
| 50%   | 2.000000   | 1.000000   |
| 75%   | 3.000000   | 1.000000   |
| max   | 3.000000   | 1.000000   |

[11]: `df.sex.value_counts()`
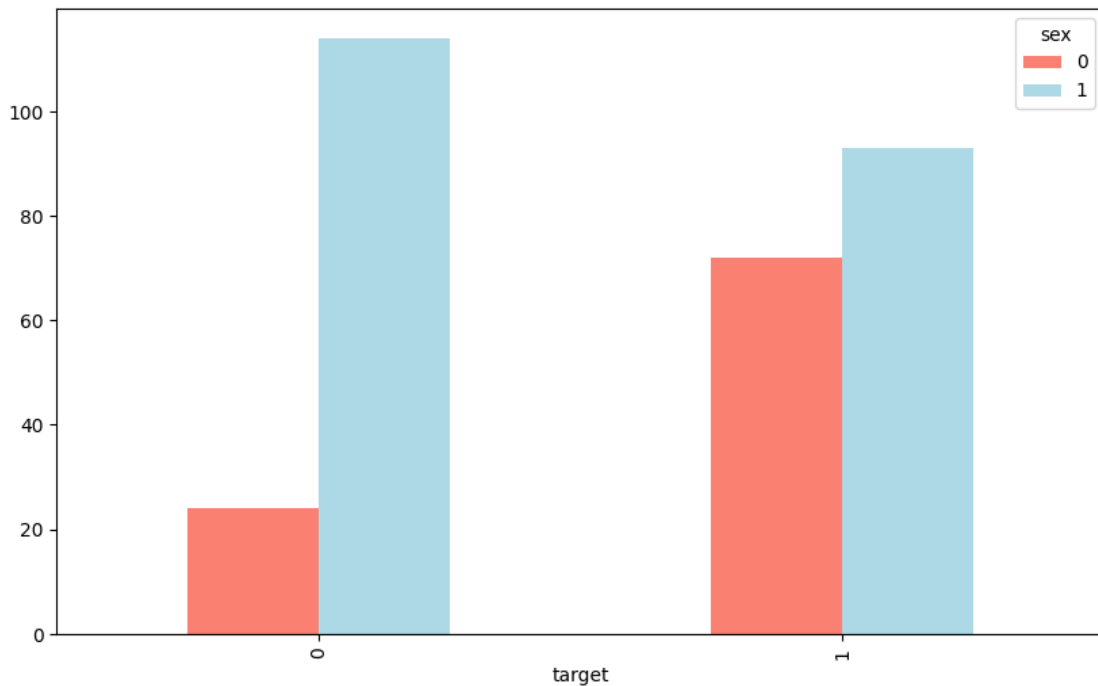
[11]:
```
1    207
0     96
Name: sex, dtype: int64
```

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

```
[12]: sex       0    1
      target
      0        24  114
      1        72   93
```
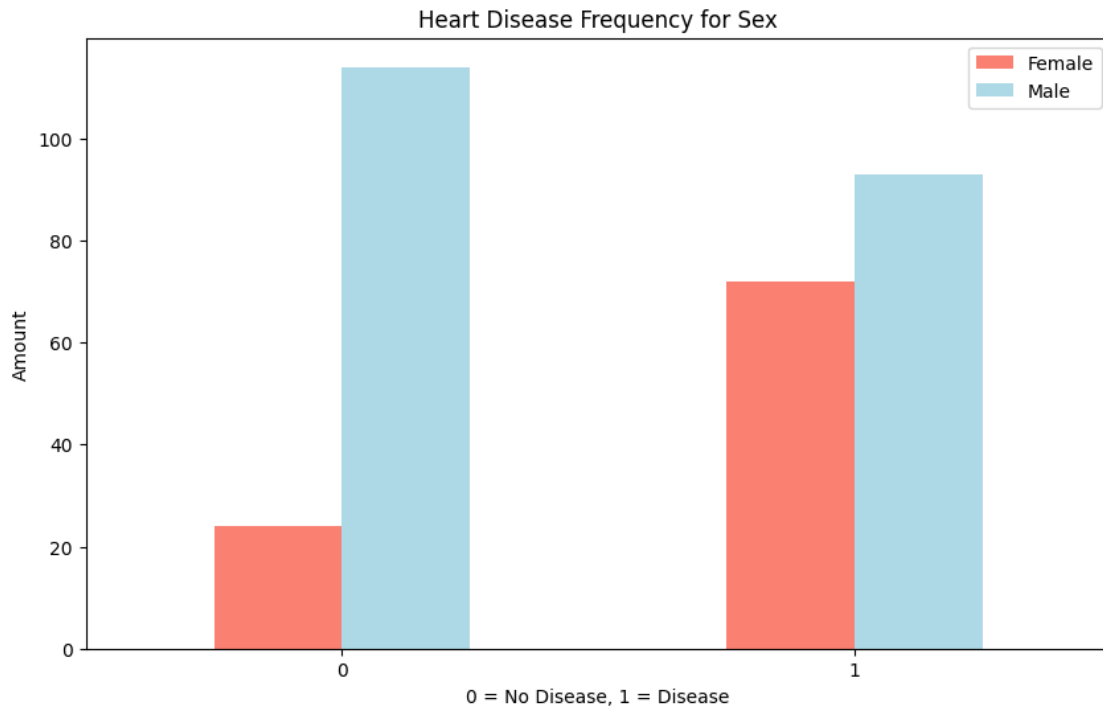
```
[13]: # Create a plot
      pd.crosstab(df.target, df.sex).plot(kind="bar",
                                          figsize=(10,6),
                                          color=["salmon", "lightblue"]);
```



```
[14]: # Create a plot
      pd.crosstab(df.target, df.sex).plot(kind="bar", figsize=(10,6),␣
        ↪color=["salmon", "lightblue"])

      # Add some attributes to it
      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); # keep the labels on the x-axis vertical
```
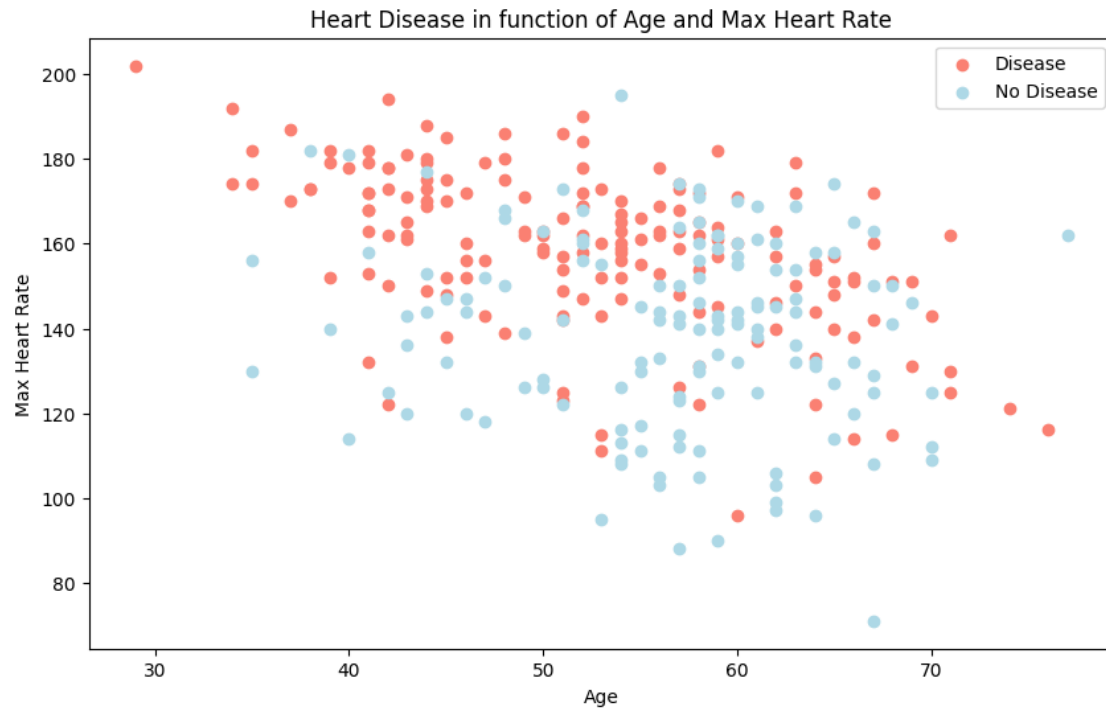
Heart Disease Frequency for Sex

```
[15]: # Create another figure
      plt.figure(figsize=(10,6))

      # Start with positve examples
      plt.scatter(df.age[df.target==1],
                  df.thalach[df.target==1],
                  c="salmon") # define it as a scatter figure

      # Now for negative examples, we want them on the same plot, so we call plt again
      plt.scatter(df.age[df.target==0],
                  df.thalach[df.target==0],
                  c="lightblue") # axis always come as (x, y)

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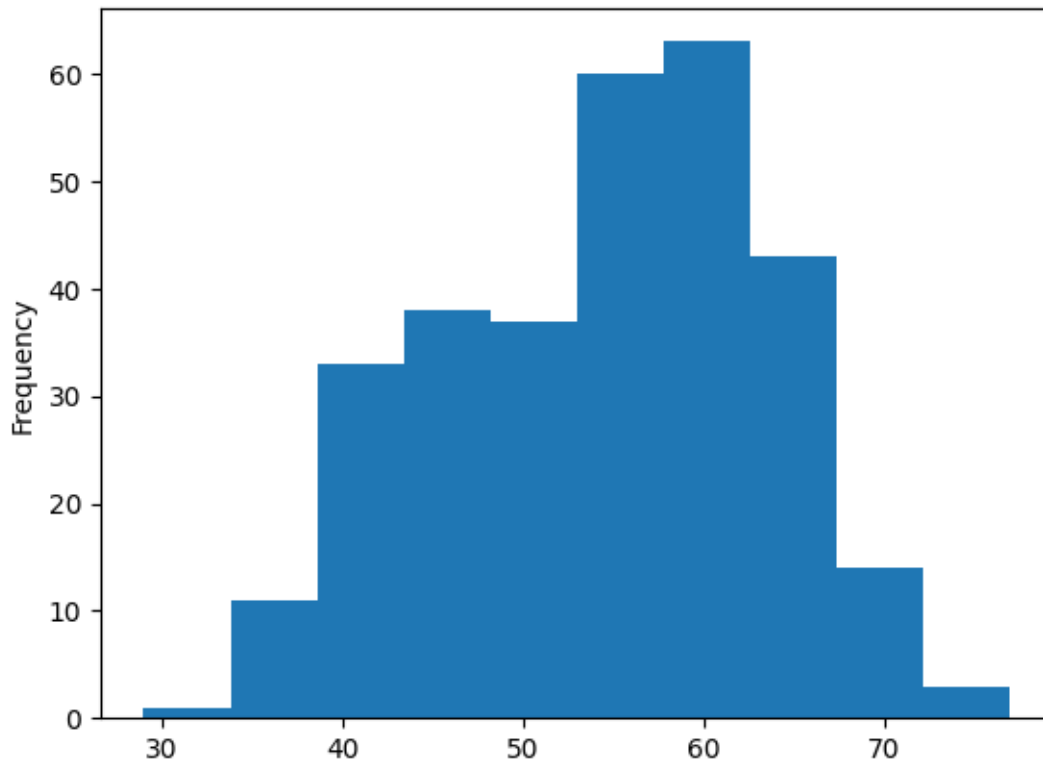
Heart Disease in function of Age and Max Heart Rate

[16]: ```python
# Histograms are a great way to check the distribution of a variable
df.age.plot.hist();
```

```
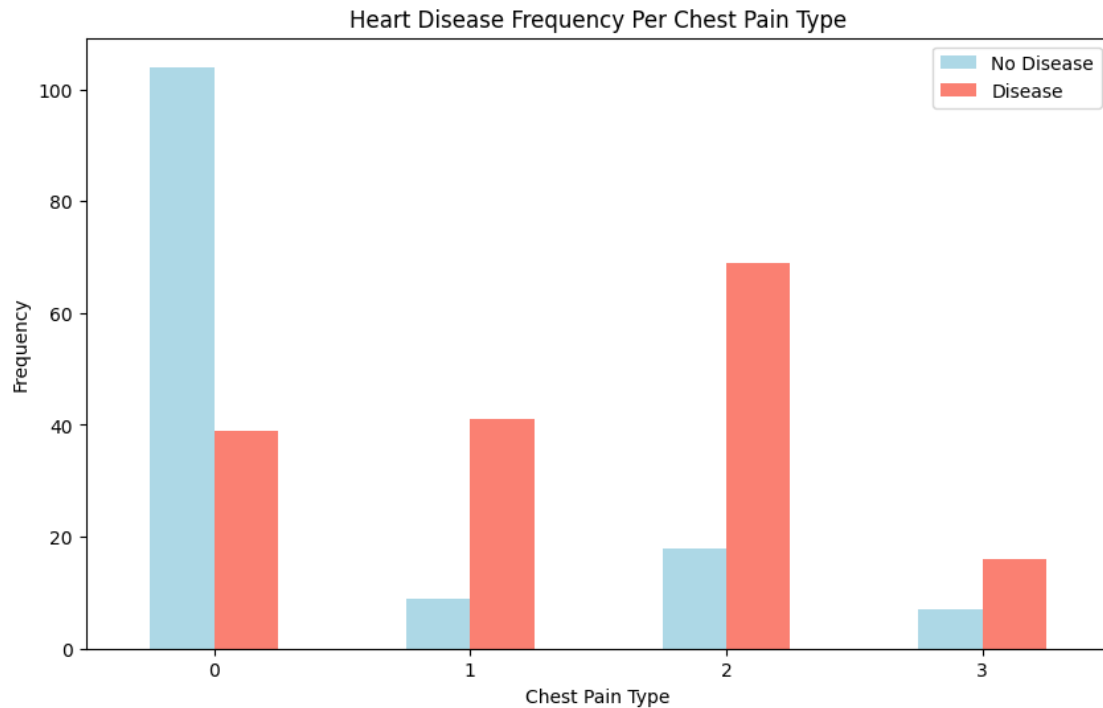[17]: pd.crosstab(df.cp, df.target)
```

```
[17]: target     0    1
      cp
      0         104   39
      1           9   41
      2          18   69
      3           7   16
```

```
[18]: # Create a new crosstab and base plot
      pd.crosstab(df.cp, df.target).plot(kind="bar",
                                         figsize=(10,6),
                                         color=["lightblue", "salmon"])

      # Add attributes to the plot to make it more readable
      plt.title("Heart Disease Frequency Per Chest Pain Type")
      plt.xlabel("Chest Pain Type")
      plt.ylabel("Frequency")
      plt.legend(["No Disease", "Disease"])
      plt.xticks(rotation = 0);
```

Heart Disease Frequency Per Chest Pain Type



```
[19]:  # Find the correlation between our independent variables
       corr_matrix = df.corr()
       corr_matrix
```

```
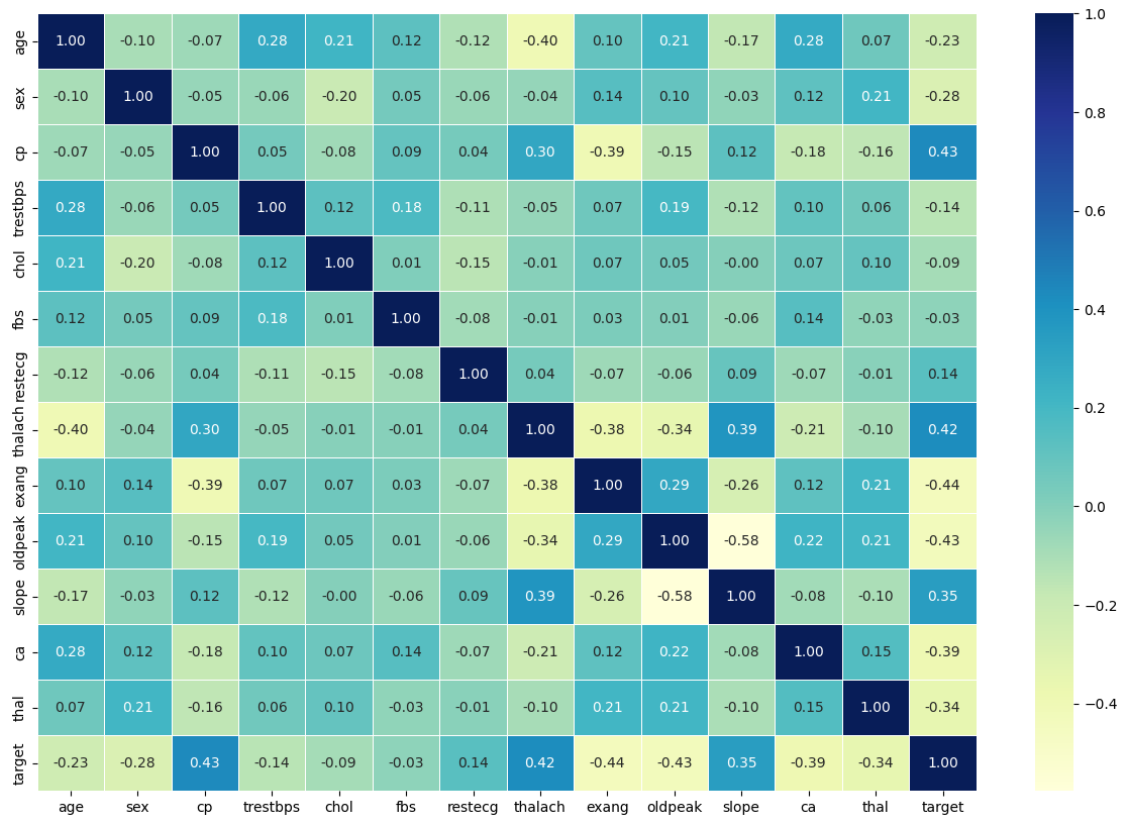[19]:                age       sex        cp  trestbps      chol       fbs  \
       age       1.000000 -0.098447 -0.068653  0.279351  0.213678  0.121308
       sex      -0.098447  1.000000 -0.049353 -0.056769 -0.197912  0.045032
       cp       -0.068653 -0.049353  1.000000  0.047608 -0.076904  0.094444
       trestbps  0.279351 -0.056769  0.047608  1.000000  0.123174  0.177531
       chol      0.213678 -0.197912 -0.076904  0.123174  1.000000  0.013294
       fbs       0.121308  0.045032  0.094444  0.177531  0.013294  1.000000
       restecg  -0.116211 -0.058196  0.044421 -0.114103 -0.151040 -0.084189
       thalach  -0.398522 -0.044020  0.295762 -0.046698 -0.009940 -0.008567
       exang     0.096801  0.141664 -0.394280  0.067616  0.067023  0.025665
       oldpeak   0.210013  0.096093 -0.149230  0.193216  0.053952  0.005747
       slope    -0.168814 -0.030711  0.119717 -0.121475 -0.004038 -0.059894
       ca        0.276326  0.118261 -0.181053  0.101389  0.070511  0.137979
       thal      0.068001  0.210041 -0.161736  0.062210  0.098803 -0.032019
       target   -0.225439 -0.280937  0.433798 -0.144931 -0.085239 -0.028046

                  restecg   thalach     exang   oldpeak     slope        ca  \
       age      -0.116211 -0.398522  0.096801  0.210013 -0.168814  0.276326
       sex      -0.058196 -0.044020  0.141664  0.096093 -0.030711  0.118261
       cp        0.044421  0.295762 -0.394280 -0.149230  0.119717 -0.181053
```

```
trestbps  -0.114103 -0.046698   0.067616   0.193216 -0.121475   0.101389
chol      -0.151040 -0.009940   0.067023   0.053952 -0.004038   0.070511
fbs       -0.084189 -0.008567   0.025665   0.005747 -0.059894   0.137979
restecg    1.000000  0.044123  -0.070733  -0.058770  0.093045  -0.072042
thalach    0.044123  1.000000  -0.378812  -0.344187  0.386784  -0.213177
exang     -0.070733 -0.378812   1.000000   0.288223 -0.257748   0.115739
oldpeak   -0.058770 -0.344187   0.288223   1.000000 -0.577537   0.222682
slope      0.093045  0.386784  -0.257748  -0.577537  1.000000  -0.080155
ca        -0.072042 -0.213177   0.115739   0.222682 -0.080155   1.000000
thal      -0.011981 -0.096439   0.206754   0.210244 -0.104764   0.151832
target     0.137230  0.421741  -0.436757  -0.430696  0.345877  -0.391724

              thal     target
age        0.068001 -0.225439
sex        0.210041 -0.280937
cp        -0.161736  0.433798
trestbps   0.062210 -0.144931
chol       0.098803 -0.085239
fbs       -0.032019 -0.028046
restecg   -0.011981  0.137230
thalach   -0.096439  0.421741
exang      0.206754 -0.436757
oldpeak    0.210244 -0.430696
slope     -0.104764  0.345877
ca         0.151832 -0.391724
thal       1.000000 -0.344029
target    -0.344029  1.000000
```

[20]: 
```python
# Let's make it look a little prettier
corr_matrix = df.corr()
plt.figure(figsize=(15, 10))
sns.heatmap(corr_matrix,
            annot=True,
            linewidths=0.5,
            fmt= ".2f",
            cmap="YlGnBu");
```

## 1.7 5. Modeling

We've explored the data, now we'll try to use machine learning to predict our target variable based on the 13 independent variables.

Remember our problem?

> Given clinical parameters about a patient, can we predict whether or not they have heart disease?

That's what we'll be trying to answer.

And remember our evaluation metric?

> If we can reach 95% accuracy at predicting whether or not a patient has heart disease during the proof of concept, we'll pursure this project.

That's what we'll be aiming for.

But before we build a model, we have to get our dataset ready.

Let's look at it again.

```
[21]: df.head()
```

```
[21]:    age  sex  cp  trestbps  chol  fbs  restecg  thalach  exang  oldpeak  slope  \
     0   63    1   3       145   233    1        0      150      0      2.3      0
     1   37    1   2       130   250    0        1      187      0      3.5      0
     2   41    0   1       130   204    0        0      172      0      1.4      2
     3   56    1   1       120   236    0        1      178      0      0.8      2
     4   57    0   0       120   354    0        1      163      1      0.6      2

         ca  thal  target
     0    0     1       1
     1    0     2       1
     2    0     2       1
     3    0     2       1
     4    0     2       1
```

```python
[23]: # Random seed for reproducibility
      np.random.seed(42)

      # Split into train & test set
      X_train, X_test, y_train, y_test = train_test_split(X, # independent variables
                                                          y, # dependent variable
                                                          test_size = 0.2) #␣
      ↪percentage of data to use for test set
```

### 1.7.1  Model choices

Now we've got our data prepared, we can start to fit models. We'll be using the following and comparing their results.

1. Logistic Regression - `LogisticRegression()`
2. K-Nearest Neighbors - `KNeighboursClassifier()`
3. RandomForest - `RandomForestClassifier()`

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

      # Create 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 different Scikit-Learn machine learning models
          X_train : training data
          X_test : testing data
          y_train : labels assosciated with training data
          y_test : labels assosciated with test data
          """
          # Random seed for reproducible results
```

```
        np.random.seed(42)
        # Make a list 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
```

```
[25]: 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
```

```
C:\Users\felip\anaconda3\lib\site-
packages\sklearn\linear_model\_logistic.py:458: ConvergenceWarning: lbfgs failed
to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
    https://scikit-learn.org/stable/modules/linear_model.html#logistic-
regression
  n_iter_i = _check_optimize_result(
```

```
[25]: {'KNN': 0.6885245901639344,
       'Logistic Regression': 0.8852459016393442,
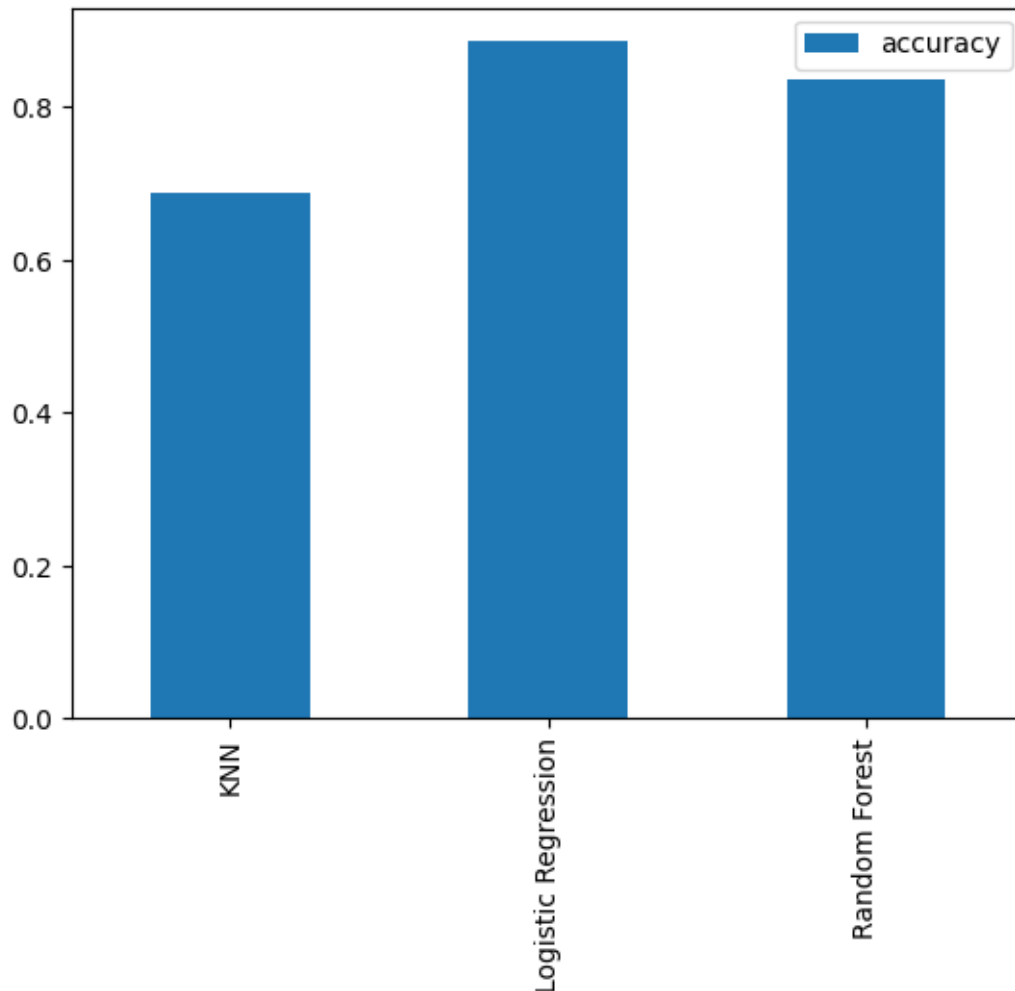       'Random Forest': 0.8360655737704918}
```

## 1.8 Model Comparison

Since we've saved our models scores to a dictionary, we can plot them by first converting them to a DataFrame.

```
[26]: model_compare = pd.DataFrame(model_scores, index=['accuracy'])
      model_compare.T.plot.bar();
```

- **Hyperparameter tuning** - Each model you use has a series of dials you can turn to dictate how they perform. Changing these values may increase or decrease model performance.
- **Feature importance** - If there are a large amount of features we're using to make predictions, do some have more importance than others? For example, for predicting heart disease, which is more important, sex or age?
- **Confusion matrix** - Compares the predicted values with the true values in a tabular way, if 100% correct, all values in the matrix will be top left to bottom right (diagnol line).
- **Cross-validation** - Splits your dataset into multiple parts and train and tests your model on each part and evaluates performance as an average.
- **Precision** - Proportion of true positives over total number of samples. Higher precision leads to less false positives.
- **Recall** - Proportion of true positives over total number of true positives and false negatives. Higher recall leads to less false negatives.
- **F1 score** - Combines precision and recall into one metric. 1 is best, 0 is worst.
- **Classification report** - Sklearn has a built-in function called `classification_report()` which returns some of the main classification metrics such as precision, recall and f1-score.

- **ROC Curve** - Receiver Operating Characterisitc is a plot of true positive rate versus false positive rate.
- **Area Under Curve (AUC)** - The area underneath the ROC

### 1.8.1 Tune KNeighborsClassifier (K-Nearest Neighbors or KNN) by hand

There's one main hyperparameter we can tune for the K-Nearest Neighbors (KNN) algorithm, and that is number of neighbours. The default is 5 (`n_neigbors=5`).

What are neighbours?

Imagine all our different samples on one graph like the scatter graph we have above. KNN works by assuming dots which are closer together belong to the same class. If `n_neighbors=5` then it assume a dot with the 5 closest dots around it are in the same class.

We've left out some details here like what defines close or how distance is calculated but I encourage you to research them.

For now, let's try a few different values of `n_neighbors`.

```
[28]: # Create a list of train scores
      train_scores = []

      # Create a list of test scores
      test_scores = []

      # Create a list of different values for n_neighbors
      neighbors = range(1, 21) # 1 to 20

      # Setup algorithm
      knn = KNeighborsClassifier()

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

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

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

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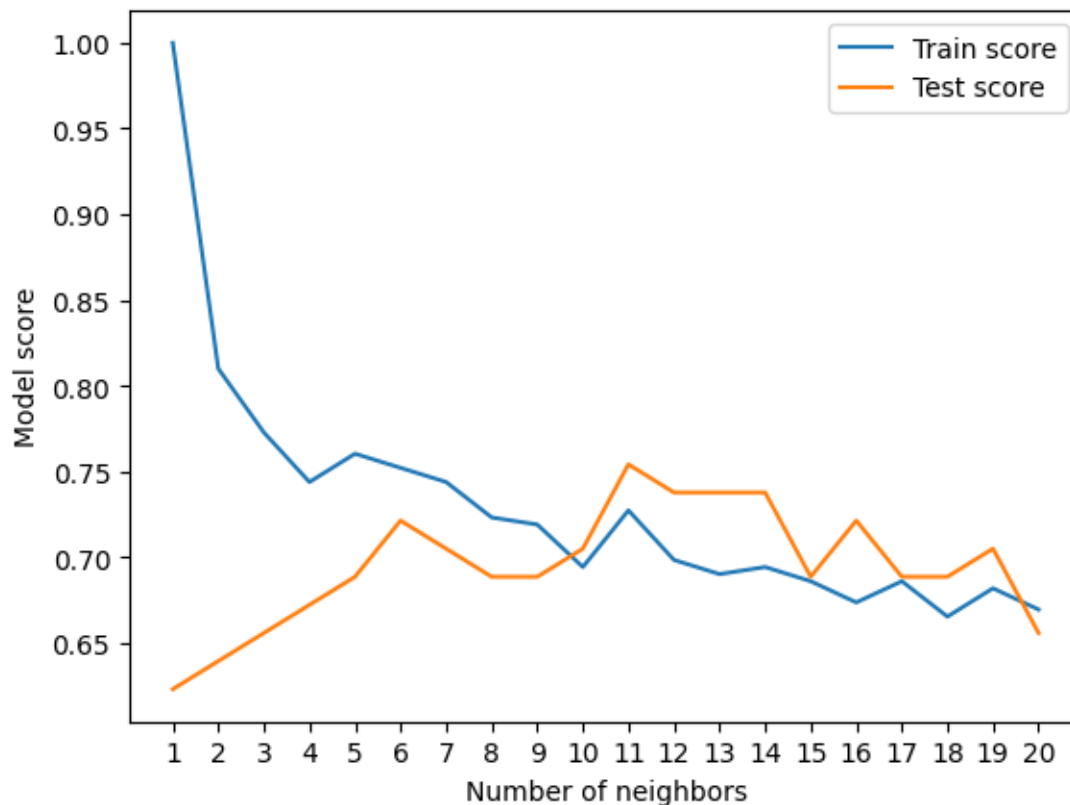          test_scores.append(knn.score(X_test, y_test))
```

```
[29]: train_scores
```

```
[29]: [1.0,
       0.8099173553719008,
       0.7727272727272727,
```

```
   0.743801652892562,
   0.7603305785123967,
   0.7520661157024794,
   0.743801652892562,
   0.7231404958677686,
   0.71900826446281,
   0.6942148760330579,
   0.7272727272727273,
   0.6983471074380165,
   0.6900826446280992,
   0.6942148760330579,
   0.6859504132231405,
   0.6735537190082644,
   0.6859504132231405,
   0.6652892561983471,
   0.6818181818181818,
   0.6694214876033058]
```

[30]:
```python
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}%")
```

```
Maximum KNN score on the test data: 75.41%
```

Looking at the graph, `n_neighbors = 11` seems best.

Even knowing this, the KNN's model performance didn't get near what `LogisticRegression` or the `RandomForestClassifier` did.

Because of this, we'll discard KNN and focus on the other two.

We've tuned KNN by hand but let's see how we can `LogisticsRegression` and `RandomForestClassifier` using `RandomizedSearchCV`.

Instead of us having to manually try different hyperparameters by hand, `RandomizedSearchCV` tries a number of different combinations, evaluates them and saves the best.

### 1.8.2 Tuning models with with `RandomizedSearchCV`

Reading the Scikit-Learn documentation for `LogisticRegression`, we find there's a number of different hyperparameters we can tune.

The same for `RandomForestClassifier`.

Let's create a hyperparameter grid (a dictionary of different hyperparameters) for each and then test them out.

```
[31]:  # Different LogisticRegression hyperparameters
       log_reg_grid = {"C": np.logspace(-4, 4, 20),
                       "solver": ["liblinear"]}

       # Different RandomForestClassifier hyperparameters
       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 let's use `RandomizedSearchCV` to try and tune our `LogisticRegression` model.

We'll pass it the different hyperparameters from `log_reg_grid` as well as set `n_iter = 20`. This means, `RandomizedSearchCV` will try 20 different combinations of hyperparameters from `log_reg_grid` and save the best ones.

```
[33]:  # Setup random seed
       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
       rs_log_reg.fit(X_train, y_train);
```

```
Fitting 5 folds for each of 20 candidates, totalling 100 fits
```

```
[34]:  rs_log_reg.best_params_
```

```
[34]:  {'solver': 'liblinear', 'C': 0.23357214690901212}
```

```
[35]:  rs_log_reg.score(X_test, y_test)
```

```
[35]:  0.8852459016393442
```

Now we've tuned `LogisticRegression` using `RandomizedSearchCV`, we'll do the same for `RandomForestClassifier`.

```
[36]:  # 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
rs_rf.fit(X_train, y_train);
```

Fitting 5 folds for each of 20 candidates, totalling 100 fits

```
[37]:  # Find the best parameters
       rs_rf.best_params_
```

```
[37]:  {'n_estimators': 210,
        'min_samples_split': 4,
        'min_samples_leaf': 19,
        'max_depth': 3}
```

```
[38]:  # Evaluate the randomized search random forest model
       rs_rf.score(X_test, y_test)
```

[38]: 0.8688524590163934

Tuning the hyperparameters for each model saw a slight performance boost in both the `RandomForestClassifier` and `LogisticRegression`.

This is akin to tuning the settings on your oven and getting it to cook your favourite dish just right.

But since `LogisticRegression` is pulling out in front, we'll try tuning it further with `GridSearchCV`.

### 1.8.3 Tuning a model with `GridSearchCV`

The difference between `RandomizedSearchCV` and `GridSearchCV` is where `RandomizedSearchCV` searches over a grid of hyperparameters performing `n_iter` combinations, `GridSearchCV` will test every single possible combination.

In short: * `RandomizedSearchCV` - tries `n_iter` combinations of hyperparameters and saves the best. * `GridSearchCV` - tries every single combination of hyperparameters and saves the best.

```
[39]:  # Different LogisticRegression hyperparameters
       log_reg_grid = {"C": np.logspace(-4, 4, 20),
                       "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);
```

```
Fitting 5 folds for each of 20 candidates, totalling 100 fits
```

[41]:
```
# Check the best parameters
gs_log_reg.best_params_
```

[41]: {'C': 0.23357214690901212, 'solver': 'liblinear'}

[42]:
```
# Evaluate the model
gs_log_reg.score(X_test, y_test)
```

[42]: 0.8852459016393442

[43]:
```
# Make preidctions on test data
y_preds = gs_log_reg.predict(X_test)
```

[44]:
```
y_preds
```

[44]: array([0, 1, 1, 0, 1, 1, 1, 0, 0, 1, 1, 0, 1, 0, 1, 1, 1, 0, 0, 0, 1, 0,
       0, 1, 1, 1, 1, 1, 0, 1, 0, 0, 0, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1,
       1, 0, 1, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0], dtype=int64)

[45]:
```
y_test
```

[45]: array([0, 0, 1, 0, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 1, 1, 1, 0, 0, 0, 1, 0,
       0, 1, 1, 0, 0, 1, 0, 1, 1, 0, 0, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1,
       1, 0, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0], dtype=int64)

### 1.8.4 Confusion matrix

A confusion matrix is a visual way to show where your model made the right predictions and where it made the wrong predictions (or in other words, got confused).

Scikit-Learn allows us to create a confusion matrix using `confusion_matrix()` and passing it the true labels and predicted labels.

[47]:
```
print(confusion_matrix(y_test, y_preds))
```

```
[[25  4]
 [ 3 29]]
```

[49]:
```
# Import Seaborn
import seaborn as sns
sns.set(font_scale=1.5) # Increase font size
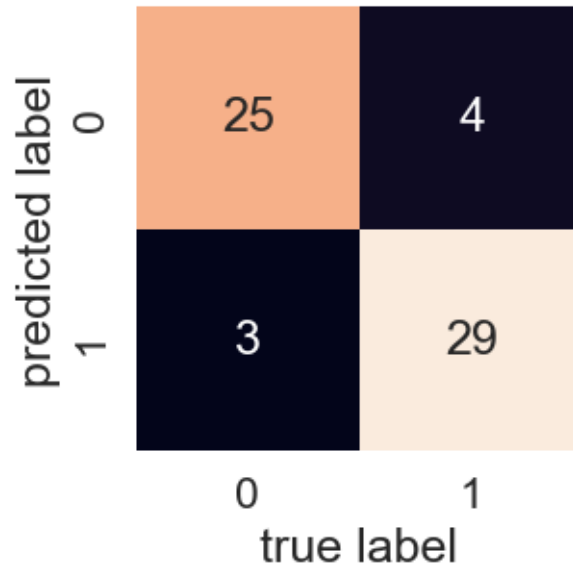
def plot_conf_mat(y_test, y_preds):
    """
    Plots a confusion matrix using Seaborn's heatmap().
    """
```

```
    fig, ax = plt.subplots(figsize=(3, 3))
    ax = sns.heatmap(confusion_matrix(y_test, y_preds),
                     annot=True, # Annotate the boxes
                     cbar=False)
    plt.xlabel("true label")
    plt.ylabel("predicted label")

plot_conf_mat(y_test, y_preds)
```



### 1.8.5  Classification report

We can make a classification report using `classification_report()` and passing it the true labels as well as our models predicted labels.

A classification report will also give us information of the precision and recall of our model for each class.

```
[50]: # Show classification report
      print(classification_report(y_test, y_preds))
```

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.89      | 0.86   | 0.88     | 29      |
| 1            | 0.88      | 0.91   | 0.89     | 32      |
|              |           |        |          |         |
| accuracy     |           |        | 0.89     | 61      |
| macro avg    | 0.89      | 0.88   | 0.88     | 61      |
| weighted avg | 0.89      | 0.89   | 0.89     | 61      |

```
[51]:   # Check best hyperparameters
        gs_log_reg.best_params_

[51]:   {'C': 0.23357214690901212, 'solver': 'liblinear'}

[52]:   # Import cross_val_score
        from sklearn.model_selection import cross_val_score

        # Instantiate best model with best hyperparameters (found with GridSearchCV)
        clf = LogisticRegression(C=0.23357214690901212,
                                 solver="liblinear")

[53]:   # Cross-validated accuracy score
        cv_acc = cross_val_score(clf,
                                 X,
                                 y,
                                 cv=5, # 5-fold cross-validation
                                 scoring="accuracy") # accuracy as scoring
        cv_acc

[53]:   array([0.81967213, 0.90163934, 0.8852459 , 0.88333333, 0.75       ])

[54]:   cv_acc = np.mean(cv_acc)
        cv_acc

[54]:   0.8479781420765027

[55]:   # Cross-validated precision score
        cv_precision = np.mean(cross_val_score(clf,
                                               X,
                                               y,
                                               cv=5, # 5-fold cross-validation
                                               scoring="precision")) # precision as␣
         ↪scoring
        cv_precision

[55]:   0.8215873015873015

[56]:   # Cross-validated recall score
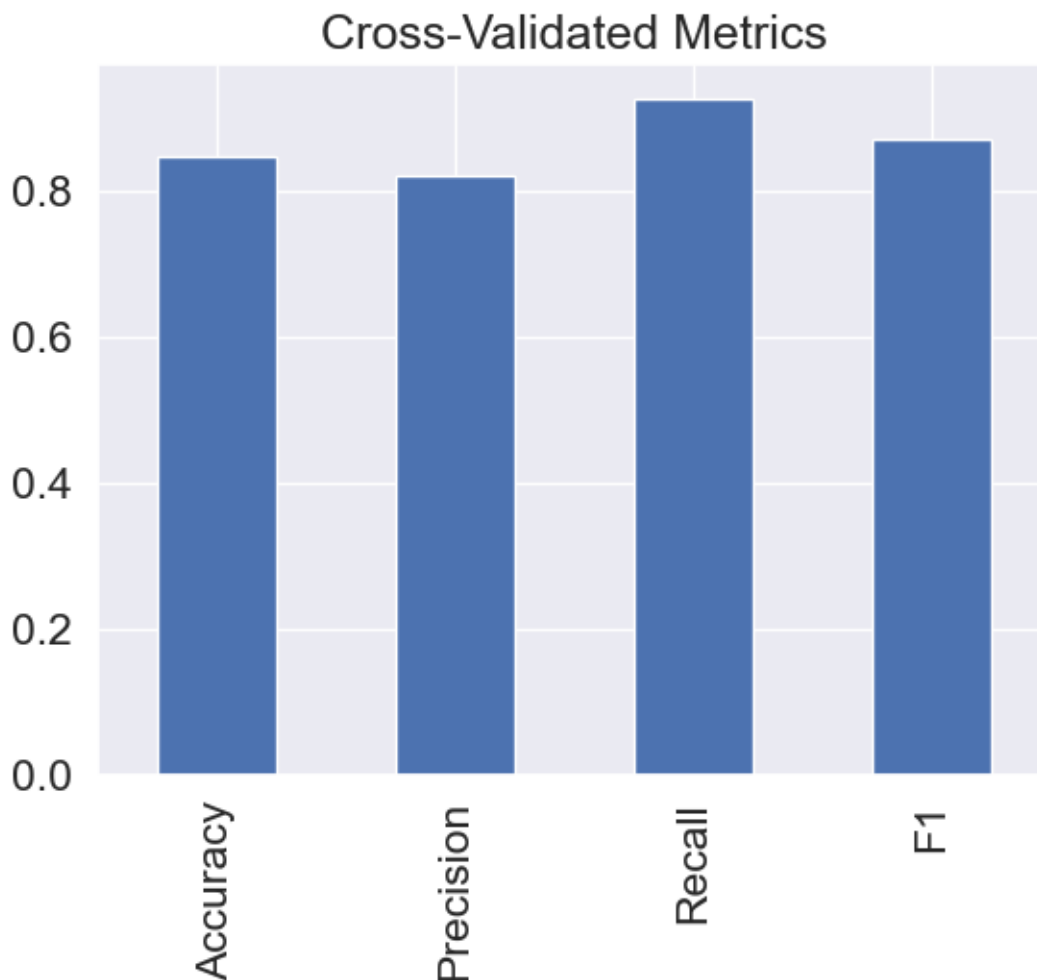        cv_recall = np.mean(cross_val_score(clf,
                                            X,
                                            y,
                                            cv=5, # 5-fold cross-validation
                                            scoring="recall")) # recall as scoring
        cv_recall

[56]:   0.9272727272727274
```

```
[57]:  # Cross-validated F1 score
       cv_f1 = np.mean(cross_val_score(clf,
                                       X,
                                       y,
                                       cv=5, # 5-fold cross-validation
                                       scoring="f1")) # f1 as scoring
       cv_f1
```

[57]:  0.8705403543192143

```
[58]:  # Visualizing 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 Metrics", legend=False);
```

## 1.9 Feature importance

Feature importance is another way of asking, "which features contributing most to the outcomes of the model?"

Or for our problem, trying to predict heart disease using a patient's medical characterisitcs, which charateristics contribute most to a model predicting whether someone has heart disease or not?

Unlike some of the other functions we've seen, because how each model finds patterns in data is slightly different, how a model judges how important those patterns are is different as well. This means for each model, there's a slightly different way of finding which features were most important.

You can usually find an example via the Scikit-Learn documentation or via searching for something like "[MODEL TYPE] feature importance", such as, "random forest feature importance".

Since we're using `LogisticRegression`, we'll look at one way we can calculate feature importance for it.

To do so, we'll use the `coef_` attribute. Looking at the Scikit-Learn documentation for `LogisticRegression`, the `coef_` attribute is the coefficient of the features in the decision function.

We can access the `coef_` attribute after we've fit an instance of `LogisticRegression`.

```
[59]: # Fit an instance of LogisticRegression (taken from above)
      clf.fit(X_train, y_train);
```

```
[60]: # Check coef_
      clf.coef_
```

```
[60]: array([[ 0.00369922, -0.90424087,  0.67472827, -0.0116134 , -0.00170364,
               0.04787689,  0.33490191,  0.02472938, -0.63120403, -0.57590925,
               0.47095125, -0.6516535 , -0.69984203]])
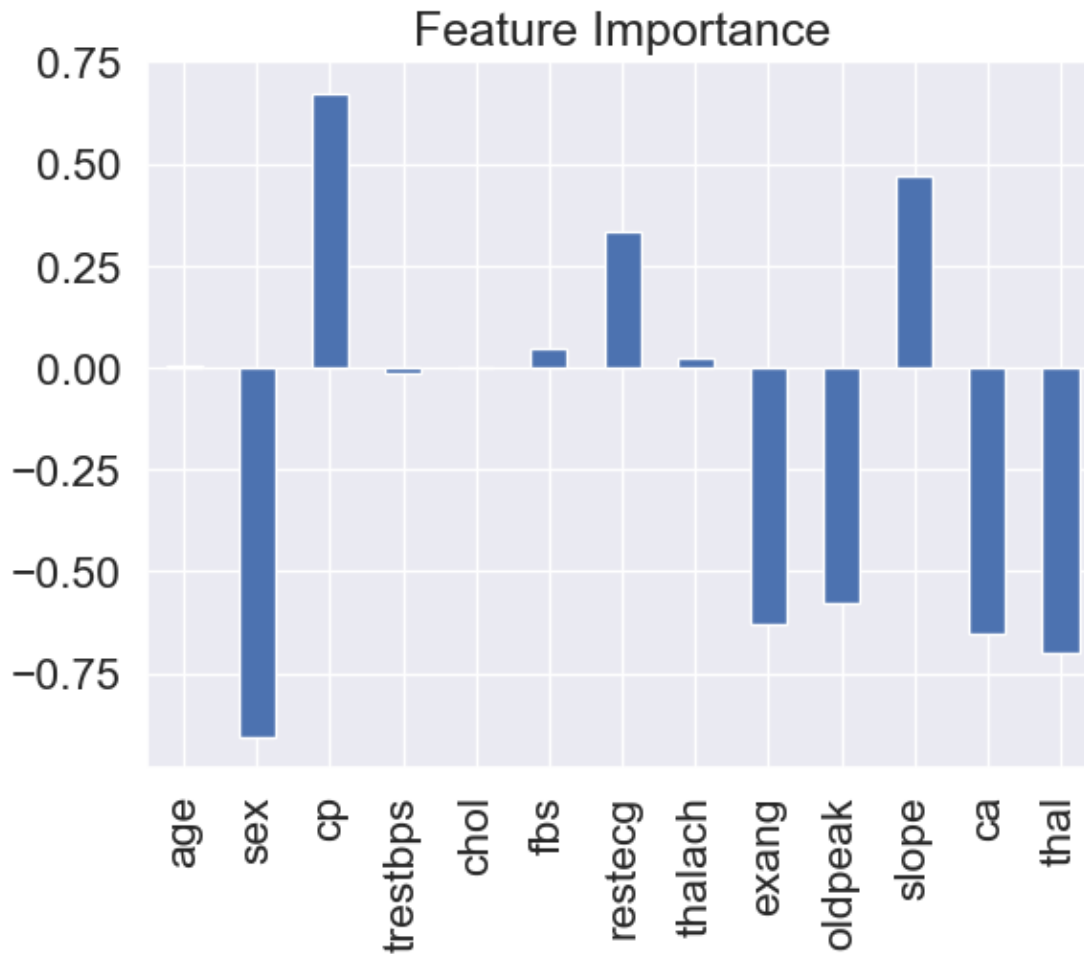```

```
[61]: # Match features to columns
      features_dict = dict(zip(df.columns, list(clf.coef_[0])))
      features_dict
```

```
[61]: {'age': 0.003699219393946938,
       'sex': -0.9042408694997455,
       'cp': 0.6747282718496431,
       'trestbps': -0.01161340293294992,
       'chol': -0.0017036445187558994,
       'fbs': 0.0478768869355857,
       'restecg': 0.33490190640612627,
       'thalach': 0.024729384262128506,
       'exang': -0.6312040272395671,
       'oldpeak': -0.5759092502488238,
       'slope': 0.4709512530750857,
       'ca': -0.6516534979207133,
```

```
'thal': -0.699842030571726}
```

Now we've match the feature coefficients to different features, let's visualize them.

```python
[62]:  # Visualize feature importance
       features_df = pd.DataFrame(features_dict, index=[0])
       features_df.T.plot.bar(title="Feature Importance", legend=False);
```



You'll notice some are negative and some are positive.

The larger the value (bigger bar), the more the feature contributes to the models decision.

If the value is negative, it means there's a negative correlation. And vice versa for positive values.

For example, the `sex` attribute has a negative value of -0.904, which means as the value for `sex` increases, the `target` value decreases.

We can see this by comparing the `sex` column to the `target` column.

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

```
[63]: target    0    1
      sex
      0         24   72
      1        114   93
```

You can see, when `sex` is 0 (female), there are almost 3 times as many (72 vs. 24) people with heart disease (`target = 1`) than without.

And then as `sex` increases to 1 (male), the ratio goes down to almost 1 to 1 (114 vs. 93) of people who have heart disease and who don't.

What does this mean?

It means the model has found a pattern which reflects the data. Looking at these figures and this specific dataset, it seems if the patient is female, they're more likely to have heart disease.

How about a positive correlation?

```
[65]: # Contrast slope (positive coefficient) with target
      pd.crosstab(df["slope"], df["target"])
```

```
[65]: target    0     1
      slope
      0         12    9
      1         91   49
      2         35  107
```

Looking back the data dictionary, we see `slope` is the "slope of the peak exercise ST segment" where: * 0: Upsloping: better heart rate with excercise (uncommon) * 1: Flatsloping: minimal change (typical healthy heart) * 2: Downslopins: signs of unhealthy heart

According to the model, there's a positive correlation of 0.470, not as strong as `sex` and `target` but still more than 0.

This positive correlation means our model is picking up the pattern that as `slope` increases, so does the `target` value.

Is this true?

When you look at the contrast (`pd.crosstab(df["slope"], df["target"])`) it is. As `slope` goes up, so does `target`.

What can you do with this information?

This is something you might want to talk to a subject matter expert about. They may be interested in seeing where machine learning model is finding the most patterns (highest correlation) as well as where it's not (lowest correlation).

Doing this has a few benefits: 1. **Finding out more** - If some of the correlations and feature importances are confusing, a subject matter expert may be able to shed some light on the situation and help you figure out more. 2. **Redirecting efforts** - If some features offer far more value than others, this may change how you collect data for different problems. See point 3. 3. **Less but better** - Similar to above, if some features are offering far more value than others, you could reduce the number of features your model tries to find patterns in as well as improve the ones which offer

the most. This could potentially lead to saving on computation, by having a model find patterns across less features, whilst still achieving the same performance levels.

## 1.10  6. Experimentation

In this case, we know the current model we're using (a tuned version of `LogisticRegression`) along with our specific data set doesn't hit the target we set ourselves.

This is where step 6 comes into its own.

A good next step would be to discuss with your team or research on your own different options of going forward.

- Could you collect more data?

- Could you try a better model? If you're working with structured data, you might want to look into CatBoost or XGBoost.

- We can improve the current models (beyond what we've done so far)?

- If your model is good enough, how would you export it and share it with others? (Hint: check out Scikit-Learn's documentation on model persistance)

The key here is to remember, your biggest restriction will be time. Hence, why it's paramount to minimise your times between experiments.

[ ]: