SC310005 Artificial Intelligence

Lecture 5: Supervised Learning (Part I)

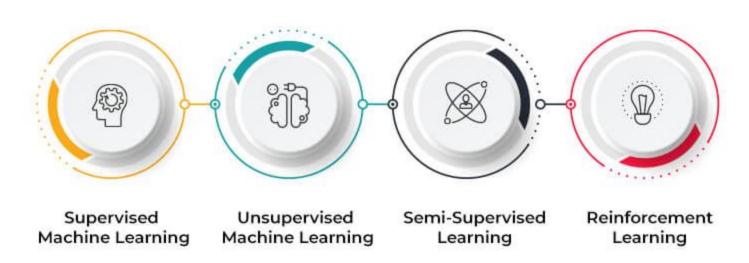
teerapong.pa@chula.ac.th

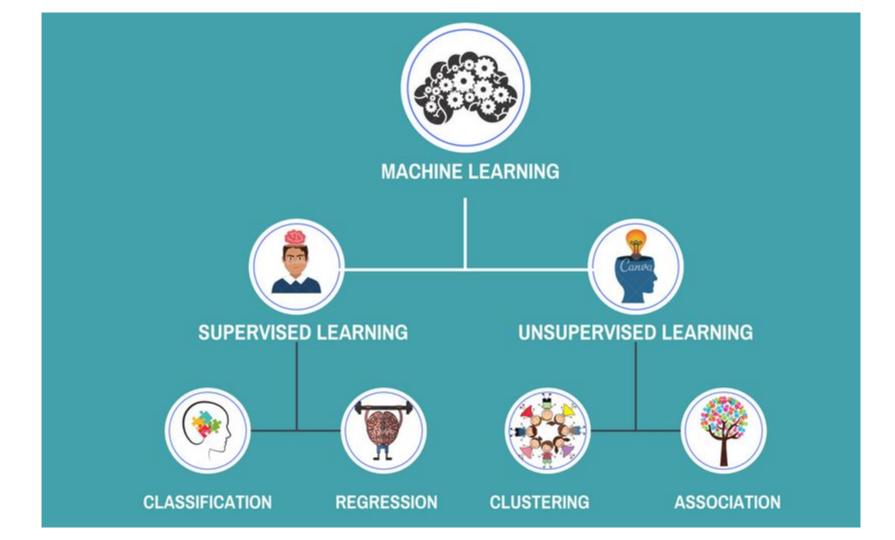
Reference

- https://www.labellerr.com/blog/supervised-vs-unsupervised-learning-whats-the-difference/
- https://www.turing.com/kb/importance-of-decision-trees-in-machine-learning
- https://towardsdatascience.com/decision-trees-in-machine-learning-641b9c4e
 8052
- https://www.section.io/engineering-education/hyperparmeter-tuning/
- https://medium.com/titansoft-engineering/decision-tree-easily-explained-c775
 c80abecb



TYPES OF MACHINE LEARNING

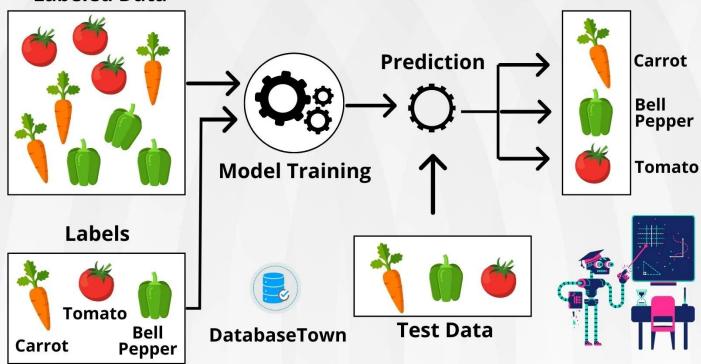


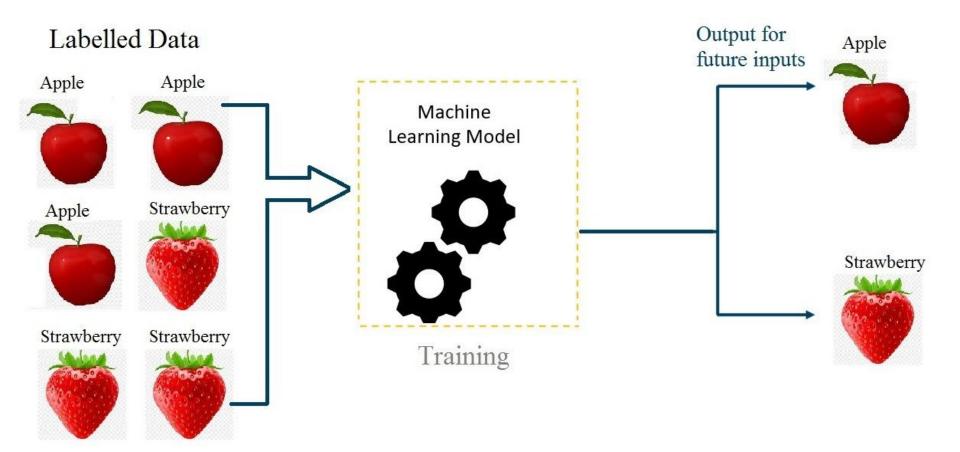


SUPERVISED LEARNING

Supervised machine learning is a branch of artificial intelligence that focuses on training models to make predictions or decisions based on labeled training data.

Labeled Data





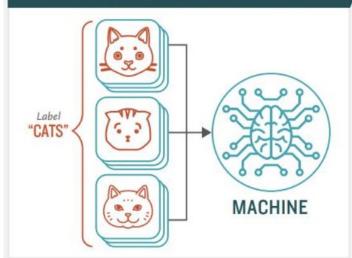
How Supervised Machine Learning Works

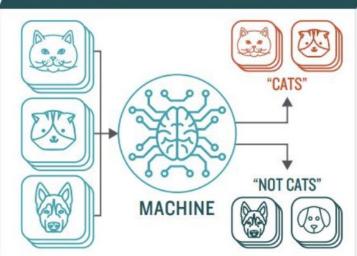
STEP I

Provide the machine learning algorithm categorized or "labeled" input and output data from to learn

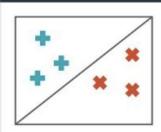
STEP 2

Feed the machine new, unlabeled information to see if it tags new data appropriately. If not, continue refining the algorithm



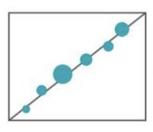


TYPES OF PROBLEMS TO WHICH IT'S SUITED



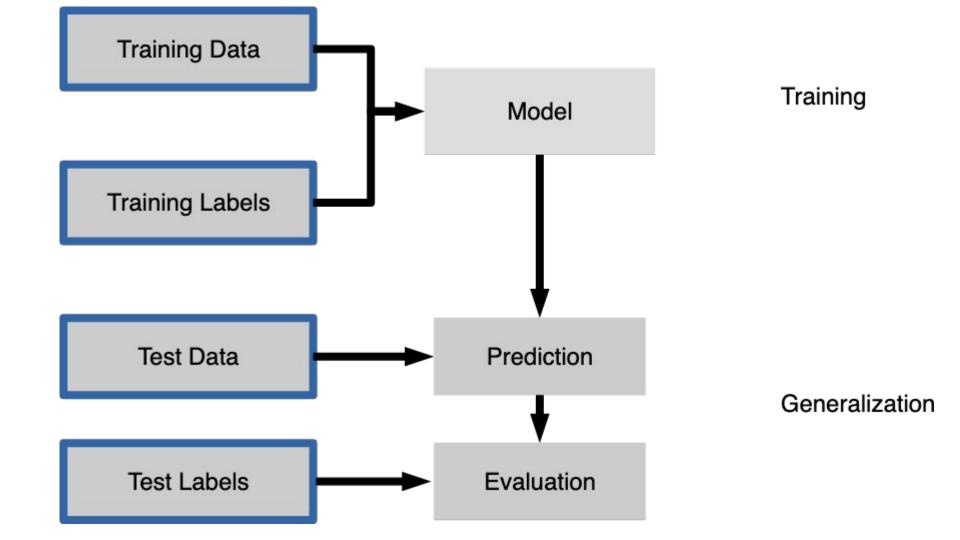
CLASSIFICATION

Sorting items into categories

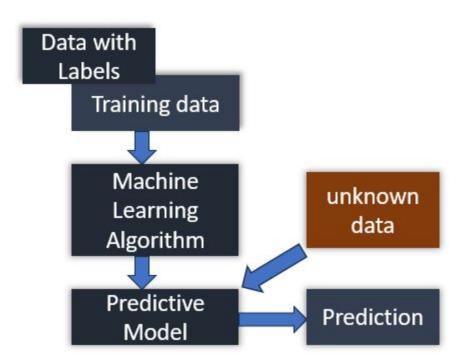


REGRESSION

Identifying real values (dollars, weight, etc.)



Supervised Learning







Visit KDnuggets.com for more cheatsheets and additional learning resources

Scikit-learn CheatSheet



Scikit-learn is an open-source Python library for all kinds of predictive data analysis. You can perform classification, regression, clustering, dimensionality reduction, model tuning, and data preprocessing tasks.

Loading the Data

Classification

from sklearn import datasets
X, y = datasets.load_wine(return_X_y=True)

Regression

diabetes = datasets.load_diabetes() X, y = diabetes.data, diabetes.target

Training And Test Data

from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(
X, y, random_state=0

Preprocessing the Data

Standardization

from sklearn.preprocessing import StandardScaler scaler = StandardScaler() scaled_X_train = scaler.fit_transform(X_train) scaled_X_test = scaler.transform(X_test)

Normalization

from sklearn.preprocessing import Normalizer norm = Normalizer() norm_X_train = norm.fit_transform(X_train) norm_X_test = norm.transform(X_test)

Binarization

from sklearn.preprocessing import Binarizer binary = Binarizer(threshold=0.0) binary_X = binary.fit_transform(X)

Encoding Categorical Features

from sklearn.preprocessing import LabelEncoder lab_enc = LabelEncoder()
y = lab_enc.fit_transform(y)

Imputer

from sklearn.impute import SimpleImputer imp_mean = SimpleImputer(missing_values=0, strategy='mean') imp_mean.fit_transform(X_train)

Supervised Learning Model

Linear Regression

from sklearn.linear_model import LinearRegression Ir = LinearRegression()

Support Vector Machines

from sklearn.svm import SVC svm_svc = SVC(kernel='linear')

Naive Bayes

from sklearn.naive_bayes import GaussianNB gnb = GaussianNB()

Unsupervised Learning Model

Principal Component Analysis

from sklearn.decomposition import PCA pca = PCA(n_components=2)

K Means

from sklearn.cluster import KMeans kmeans = KMeans(n_clusters=5, random_state=0)

Model Fitting

Supervised Learning

Ir.fit(X_train, y_train)
svm_svc.fit(X_train, y_train)

Unsupervised Learning

model = pca.fit_transform(X_train) kmeans.fit(X_train)

Prediction

Supervised Learning

y_pred = Ir.predict_proba(X_test) y_pred = svm_svc.predict(X_test)

Unsupervised Learning

y_pred = kmeans.predict(X_test)

Evaluation

Accuracy Score

Ir.score(X_test, y_test)

from sklearn.metrics import accuracy_score accuracy_score(y_test, y_pred)

Classification Report

from sklearn.metrics import classification_report print(classification_report(y_test, y_pred))

Mean Squared Error

from sklearn.metrics import mean_squared_error mean_squared_error(y_test, y_pred)

R2 Score

from sklearn.metrics import r2_score r2_score(y_test, y_pred)

Adjusted Rand Index

from sklearn.metrics import adjusted_rand_score adjusted_rand_score(y_test, y_pred)

Cross-Validation

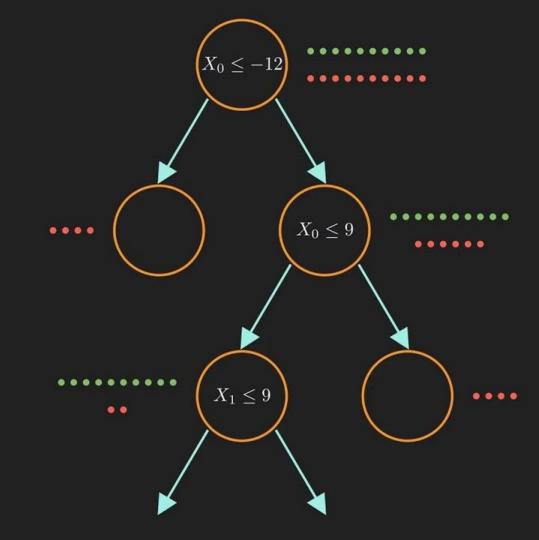
from sklearn.model_selection import cross_val_score cross_val_score(Ir, X, y, cv=5, scoring='f1_macro')

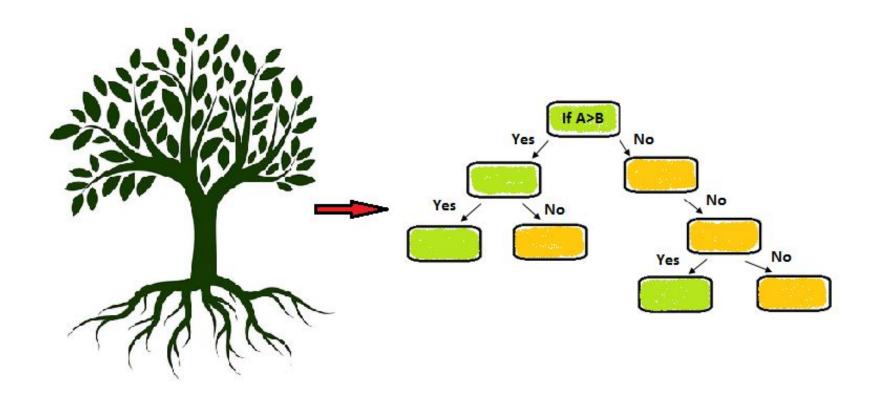
Model Tuning

from sklearn.model_selection import GridSearchCV parameters = {kerne!',(linear', 'rbf'), 'C:{f, 10}} model = GridSearchCV(svm_svc, parameters) model.fit(X_train, y_train) print(model.best_score_) print(model.best_stimator_)

Subscribe to KDnuggets News

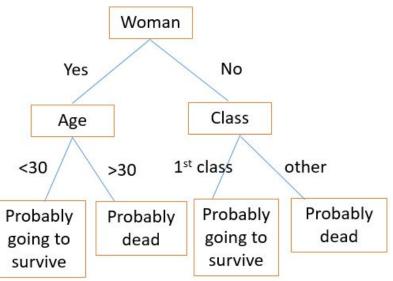
Decision Tree Classifier

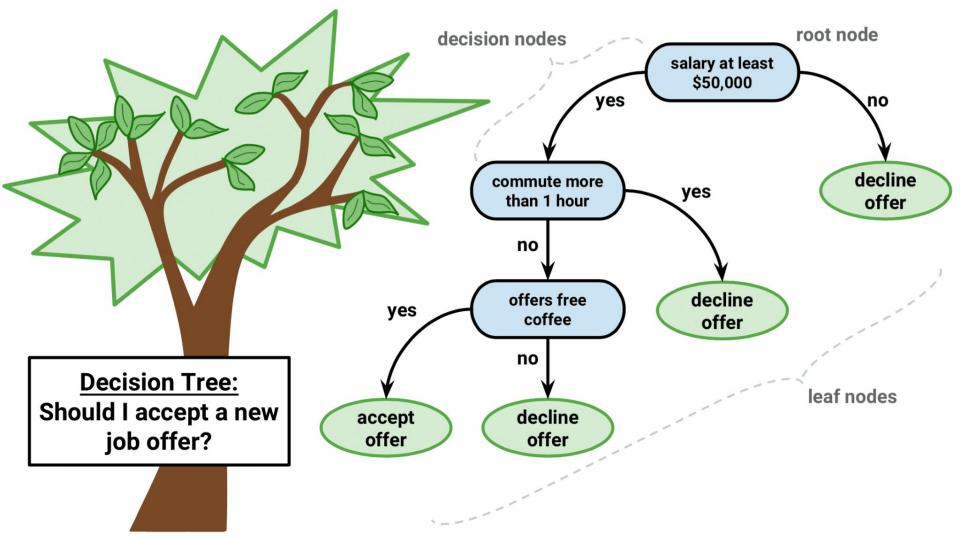




A normal tree A decision tree







Algorithms used in Decision Trees

ID3
Iterative Dichotomiser 3, an extension of D3. Uses top-down approach to build a decision tree.

C4.5

Uses C programming language. Succession of ID3. 2

CART

3

5

Classification and Regression Tree. Generates future predictions based on already available values.

CHAID

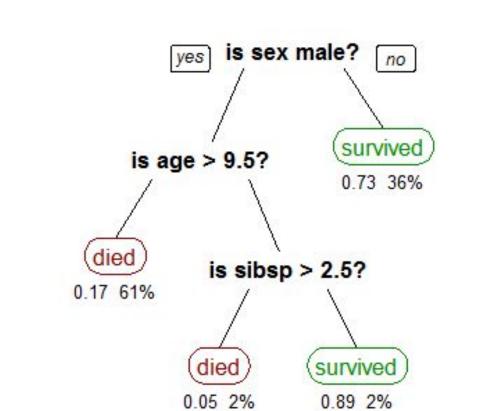
Chi-square automatic interaction detection. Based on adjusted significance testing.

4

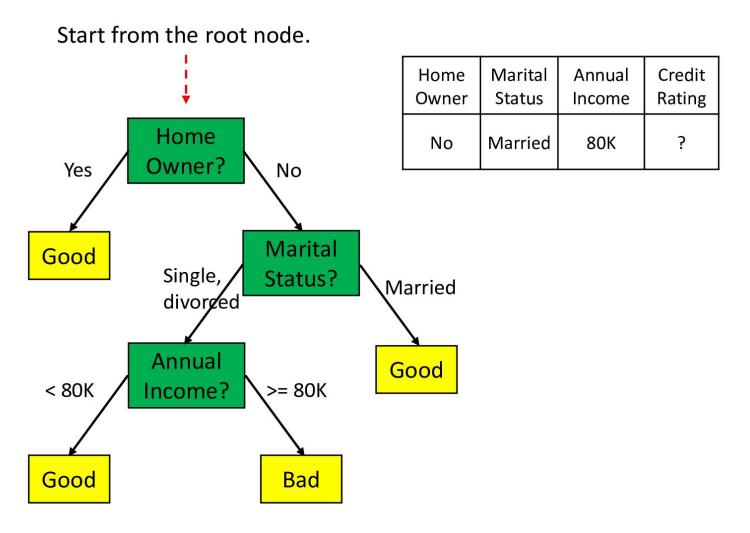
MARS

Multivariate adaptive regression splines. Helps solve non-linear regression problems.





Start from the root node. Marital Credit Home Annual Rating Owner Status Income Home No Married 80K Owner? Yes No Marital Good Single, Status? Married divorged **Annual** Good < 80K Income? >= 80K Good Bad



Install the required libraries

```
[1] # Remember to install the required libraries (pandas, sklearn, matplotlib) if you haven't already done so:

!pip install pandas scikit-learn matplotlib --quiet
```

Libraries

```
[2] import pandas as pd
    import matplotlib.pyplot as plt
    import seaborn as sns
    from sklearn.model_selection import train_test_split
    from sklearn.tree import DecisionTreeClassifier
    from sklearn.metrics import confusion_matrix, classification_report
    from sklearn import tree
[3] # Load the Titanic dataset
    # Replace 'path_to_dataset' with the actual path of your dataset
    data = pd.read_csv('titanic_dataset.csv')
```

Load the Titanic dataset

```
[3] # Load the Titanic dataset
    # Replace 'path_to_dataset' with the actual path of your dataset

data = pd.read_csv('titanic_dataset.csv')
```

Display the first few rows of the dataset

| | Display the f ta.head() | irst few r | ows of | the dataset to understand its structure | | | | | | | | |
|---|----------------------------|------------|--------|------------------------------------------------|--------|------|-------|-------|------------------|---------|-------|----------|
| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin | Embarked |
| 0 | 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN | s |
| 1 | 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence Briggs Th | female | 38.0 | 1 | 0 | PC 17599 | 71.2833 | C85 | С |
| 2 | 3 | 1 | 3 | Heikkinen, Miss. Laina | female | 26.0 | 0 | 0 | STON/O2. 3101282 | 7.9250 | NaN | s |
| 3 | 4 | 1 | 1 | Futrelle, Mrs. Jacques Heath (Lily May Peel) | female | 35.0 | 1 | 0 | 113803 | 53.1000 | C123 | s |
| 4 | 5 | 0 | 3 | Allen, Mr. William Henry | male | 35.0 | 0 | 0 | 373450 | 8.0500 | NaN | S |

Data preprocessing

```
# Data preprocessing - handle missing values, feature selection, encoding categorical variables, etc.

# For simplicity, you might want to drop some columns or perform imputation for missing values

# Replace 'selected_features' with the features you want to use for prediction

selected_features = ['Pclass', 'Sex', 'Age', 'SibSp', 'Parch', 'Fare', 'Embarked', 'Survived']

data = data[selected_features].dropna()
```

Data preprocessing (cont.)

| Po | lass | Sex | Age | SibSp | Parch | Fare | Embarked | Survived |
|----|------|--------|------|-------|-------|---------|----------|----------|
| 0 | 3 | male | 22.0 | 1 | 0 | 7.2500 | S | 0 |
| 1 | 1 | female | 38.0 | 1 | 0 | 71.2833 | С | 1 |
| 2 | 3 | female | 26.0 | 0 | 0 | 7.9250 | S | 1 |
| 3 | 1 | female | 35.0 | 1 | 0 | 53.1000 | S | 1 |
| 4 | 3 | male | 35.0 | 0 | 0 | 8.0500 | S | 0 |

Convert categorical variables to numerical using one-hot encoding

```
# Convert categorical variables to numerical using one-hot encoding
    data = pd.get_dummies(data, columns=['Sex', 'Embarked'])
[8]
   data.head()
       Pclass Age SibSp Parch
                                   Fare Survived Sex_female Sex_male Embarked_C Embarked_Q Embarked_S
            3 22.0
                                  7.2500
            1 38.0
                                 71.2833
                                                                      0
                                                                                              0
                                                                                                          0
            3 26.0
                                 7.9250
                                                                                              0
            1 35.0
                               0 53.1000
                                                                      0
                                                                                  0
                                                                                              0
            3 35.0
                                  8.0500
                                                                                              0
                                                                                  0
```

Split the data into features and target variable

```
[9] # Split the data into features and target variable
X = data.drop('Survived', axis=1)
y = data['Survived']
```

Split the data into training and testing sets

```
[10] # Split the data into training and testing sets
    X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=2023)
```

Create and train the decision tree classifier

Reference:

https://scikit-learn.org/stable/modules/generated/sklearn.tree.DecisionTreeClassifier.html

```
[11] # Create and train the decision tree classifier
  clf = DecisionTreeClassifier(random_state=42, max_depth=5)
  clf.fit(X_train, y_train)
```

```
▼ DecisionTreeClassifier
DecisionTreeClassifier(max_depth=5, random_state=42)
```

Make predictions

```
[12] # Make predictions
predictions = clf.predict(X_test)
```

Evaluate the model

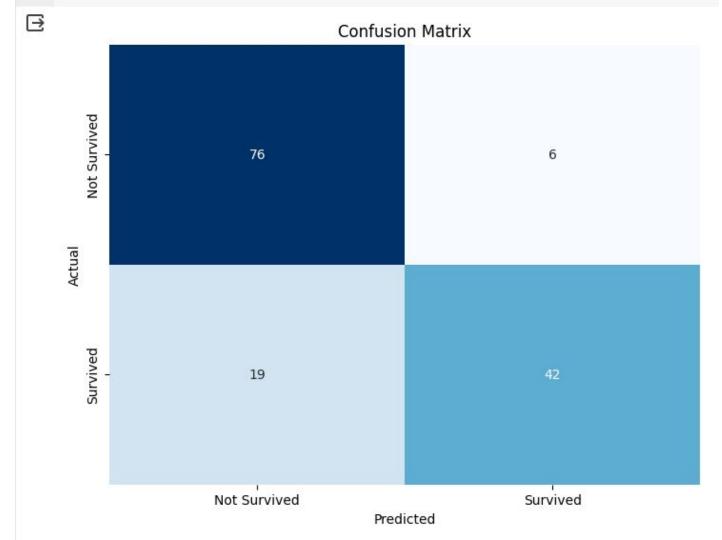
[[76 6]

[19 42]]

```
[13] # Evaluate the model
    conf_matrix = confusion_matrix(y_test, predictions)
     print("Confusion Matrix:")
     print(conf_matrix)
    Confusion Matrix:
```

Plot confusion matrix as a heatmap

```
# Plot confusion matrix as a heatmap
 plt.figure(figsize=(8, 6))
 sns.heatmap(conf_matrix, annot=True, fmt='d', cmap='Blues', cbar=False,
             xticklabels=['Not Survived', 'Survived'],
             yticklabels=['Not Survived', 'Survived'])
plt.xlabel('Predicted')
plt.ylabel('Actual')
 plt.title('Confusion Matrix')
 plt.show()
```



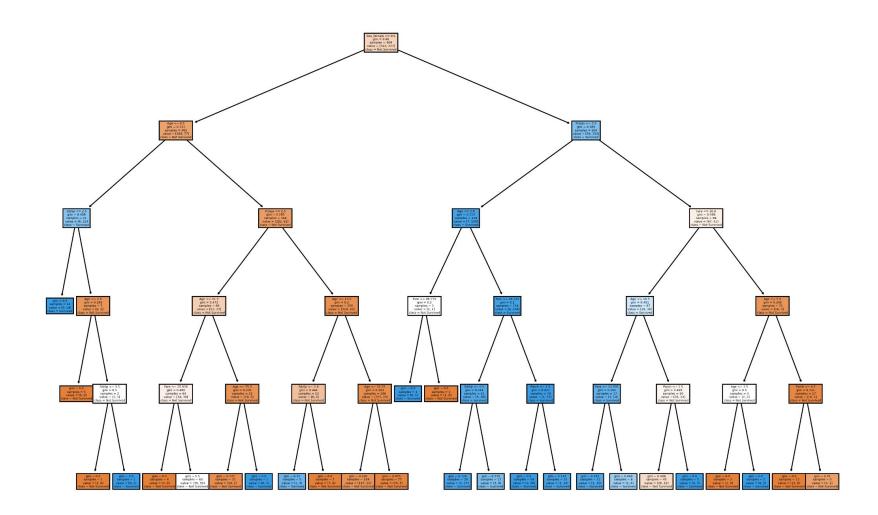
Classification Report

```
[15] print("\nClassification Report:")
    print(classification_report(y_test, predictions, digits=4))
    Classification Report:
                   precision
                                recall f1-score
                                                    support
                      0.8000
                                0.9268
                                           0.8588
                                                         82
                                           0.7706
                                                         61
                      0.8750
                                0.6885
                                           0.8252
                                                        143
        accuracy
                      0.8375
                                0.8077
                                           0.8147
                                                        143
       macro avg
    weighted avg
                                           0.8212
                      0.8320
                                0.8252
                                                        143
```

Plot the Decision Tree

```
# Set high resolution for tree plot plt.figure(figsize=(15, 10), dpi=300) # Adjust figsize and dpi for the desired resolution

# Plot the decision tree (optional) tree.plot_tree(clf, feature_names=X.columns, class_names=['Not Survived', 'Survived'], filled=True) plt.savefig('my_titanic_decision_tree_model.png', format='png', bbox_inches='tight') plt.show()
```



Week 5: Assignment

Predicting Heart Attack Risk Using Decision Tree

You are given a dataset containing various attributes of patients. Your task is to create a machine learning model using a Decision Tree Classifier to predict the likelihood of a heart attack for a patient based on their medical attributes.

