# 1. Loading and Preprocessing

## Load the breast cancer dataset from sklearn.

In [1]:
```python
from sklearn.datasets import load_breast_cancer
data = load_breast_cancer()
x = data.data
y = data.target
x,y
```

Out[1]:  (array([[1.799e+01, 1.038e+01, 1.228e+02, ..., 2.654e-01, 4.601e-01,
             1.189e-01],
            [2.057e+01, 1.777e+01, 1.329e+02, ..., 1.860e-01, 2.750e-01,
             8.902e-02],
            [1.969e+01, 2.125e+01, 1.300e+02, ..., 2.430e-01, 3.613e-01,
             8.758e-02],
            ...,
            [1.660e+01, 2.808e+01, 1.083e+02, ..., 1.418e-01, 2.218e-01,
             7.820e-02],
            [2.060e+01, 2.933e+01, 1.401e+02, ..., 2.650e-01, 4.087e-01,
             1.240e-01],
            [7.760e+00, 2.454e+01, 4.792e+01, ..., 0.000e+00, 2.871e-01,
             7.039e-02]]),
       array([0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 1, 1,
            0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0,
            0, 0, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 0, 0, 1, 1, 1, 1, 0, 1, 0, 0,
            1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 0, 0,
            1, 1, 1, 0, 1, 1, 0, 0, 1, 1, 1, 0, 0, 1, 1, 1, 1, 0, 1, 1, 0, 1,
            1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 1, 0, 0, 1, 1, 1, 0, 0, 1, 0, 1, 0,
            0, 1, 0, 0, 1, 1, 0, 1, 1, 0, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1,
            1, 1, 0, 1, 1, 1, 1, 0, 0, 1, 0, 1, 1, 0, 0, 1, 1, 0, 0, 1, 1, 1,
            1, 0, 1, 1, 0, 0, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1, 0, 0,
            0, 0, 1, 0, 0, 0, 1, 0, 1, 0, 1, 1, 0, 1, 0, 0, 0, 0, 1, 1, 0, 0,
            1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 1, 0, 1, 1,
            1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0,
            0, 0, 1, 1, 1, 1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 0, 1, 0, 0, 1, 1,
            1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1,
            1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 0, 0,
            0, 1, 1, 1, 1, 0, 1, 0, 1, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0,
            0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 1, 0, 0, 0, 1, 0, 0,
            1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 0, 1, 1, 0, 0, 1, 1,
            1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 0,
            1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 0, 1, 0, 1, 1, 1, 1,
            1, 0, 1, 1, 0, 1, 0, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0,
            1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1,
            1, 1, 1, 0, 1, 0, 1, 1, 0, 1, 1, 1, 1, 1, 0, 0, 1, 0, 1, 0, 1, 1,
            1, 1, 1, 0, 1, 1, 0, 1, 0, 1, 0, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1,
            1, 1, 1, 1, 1, 0, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
            1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 1]))

## Preprocess the data to handle any missing values and perform necessary feature scaling.

```
In [2]:  import numpy as np
         print(np.isnan(x).sum())
```

0

## Explain the preprocessing steps you performed and justify why they are necessary for this dataset.

```
In [4]:  from sklearn.preprocessing import StandardScaler
         scaler = StandardScaler()
         x_scaled = scaler.fit_transform(x)
```

```
In [5]:  x_scaled
```

```
Out[5]:  array([[ 1.09706398, -2.07333501,  1.26993369, ...,  2.29607613,
                   2.75062224,  1.93701461],
                 [ 1.82982061, -0.35363241,  1.68595471, ...,  1.0870843 ,
                  -0.24388967,  0.28118999],
                 [ 1.57988811,  0.45618695,  1.56650313, ...,  1.95500035,
                   1.152255  ,  0.20139121],
                 ...,
                 [ 0.70228425,  2.0455738 ,  0.67267578, ...,  0.41406869,
                  -1.10454895, -0.31840916],
                 [ 1.83834103,  2.33645719,  1.98252415, ...,  2.28998549,
                   1.91908301,  2.21963528],
                 [-1.80840125,  1.22179204, -1.81438851, ..., -1.74506282,
                  -0.04813821, -0.75120669]])
```

# 2. Classification Algorithm Implementation

## 1. Logistic Regression

```
In [7]:  from sklearn.linear_model import LogisticRegression
         log_reg = LogisticRegression(max_iter=10000)
         log_reg.fit(x_scaled, y)
```

```
Out[7]:        ▾        LogisticRegression
         LogisticRegression(max_iter=10000)
```

## 2. Decision Tree Classifier

```python
In [9]:  from sklearn.tree import DecisionTreeClassifier
         tree_clf = DecisionTreeClassifier()
         tree_clf.fit(x_scaled, y)
```

Out[9]:
```
▾ DecisionTreeClassifier

DecisionTreeClassifier()
```

## 3. Random Forest Classifier

```python
In [12]:  from sklearn.ensemble import RandomForestClassifier
          rf_clf = RandomForestClassifier()
          rf_clf.fit(x_scaled,y)
```

Out[12]:
```
▾ RandomForestClassifier

RandomForestClassifier()
```

## 4. Support Vector Machine (SVM)

```python
In [13]:  from sklearn.svm import SVC
          svm_clf = SVC()
          svm_clf.fit(x_scaled, y)
```

Out[13]:
```
▾ SVC

SVC()
```

## 5. k-Nearest Neighbors (k-NN)

```python
In [14]:  from sklearn.neighbors import KNeighborsClassifier
          knn_clf = KNeighborsClassifier()
          knn_clf.fit(x_scaled, y)
```

Out[14]:
```
▾ KNeighborsClassifier

KNeighborsClassifier()
```

# 3. Model Comparison

In [15]:
```python
from sklearn.model_selection import cross_val_score
models = [log_reg, tree_clf, rf_clf, svm_clf, knn_clf]
model_names = ['Logistic Regression', 'Decision Tree', 'Random Forest', 'SV
for model, name in zip(models, model_names):
    scores = cross_val_score(model, x_scaled, y, cv=5, scoring='accuracy')
    print(f"{name}: {scores.mean():.4f}")
```

```
Logistic Regression: 0.9807
Decision Tree: 0.9138
Random Forest: 0.9631
SVM: 0.9736
k-NN: 0.9649
```

# Which algorithm performed the best and which one performed the worst?

Best Algorithm: The one with the highest average accuracy. Worst Algorithm: The one with the lowest average accuracy.

In [ ]: