# Supervised Machine Learning

## Classification, Regression, Time Series

# Recap – Machine Learning Framework
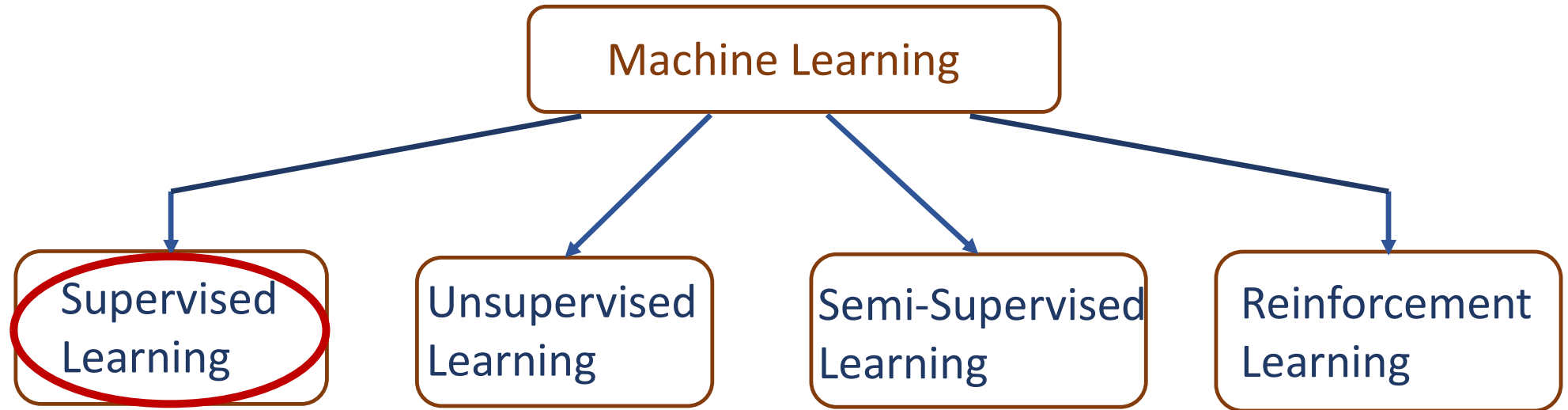
$$y = f(x)$$

output          prediction          feature or
                function            representation

- The input is converted to a vector **x**

- The output is a value indicated by **y**

- Depending on the nature of **x** and **y**, we define different types of learning

# Categorization of ML Based on Learning

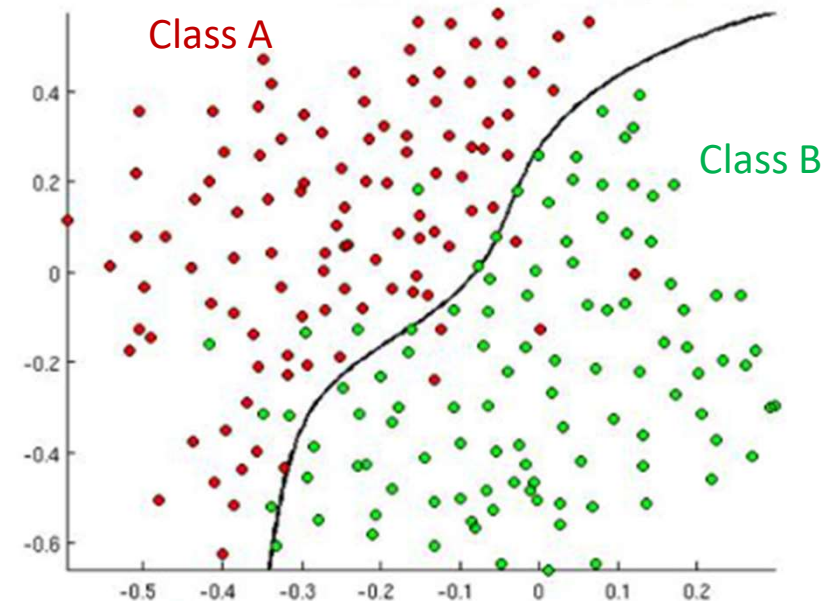# Supervised Learning

- Machine learning that are designed to learn by examples

- It is trained with labelled data
  - Feature vectors: $x_{ij}, \; i = 1..N, \; j = 1..M$
  - Output values: $y_i, \; i = 1..N$

- It maps the input to an output based on previous input-output pairs, through a mapping function, $Y = f(X)$

- Depending on the nature of $y$, we define:
  1. Classification
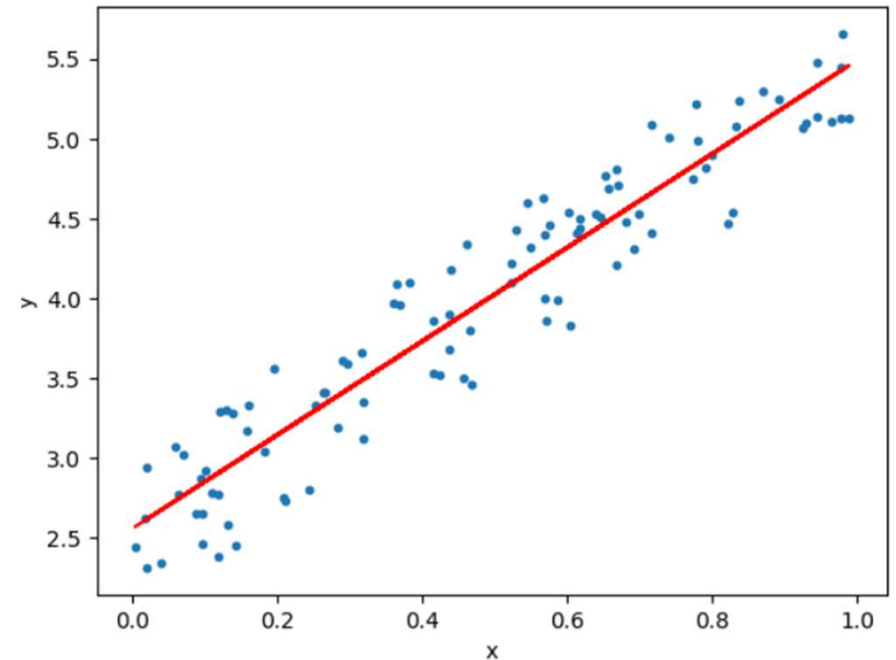  2. Regression

# Classification

- Classification predicts a discrete value/class label

- Some common classification algorithms include:
    - Decision Trees
    - Support Vector Machines
    - Naïve Bayes Classifier

- Classification is often used for tasks such as:
    - Spam filtering
    - Image classification
    - Text classification

# Regression

- Regression is a type of machine learning that predicts a continuous value.

- Some common regression algorithms include:
    - Linear regression
    - Polynomial regression

- Regression is often used for tasks such as:
    - Predicting the price of a house
    - Predicting the number of sales
    - Predicting the risk of a disease

e.g., a house's [Area, Age] ($\mathbf{x}$) vs. its Price($y$)

$y = f(x)$, interpolating (approximating) a function from examples

# Time Series Model

- Time series is a sequence of observations often ordered in time

- Popular Problem: Given a sequence, predict future samples

- Applications:
  - Meteorology,
  - Finance,
  - Marketing, etc

| Year | Sales (in Million) |
|------|--------------------|
| 1921 | 251 |
| 1931 | 279 |
| 1941 | 319 |
| 1951 | 261 |
| 1961 | 439 |
| 1971 | 348 |
| 1981 | 585 |

- We want a machine learning model to understand sequences, not samples
- Assume we have a sequence of measurements, and we want to take N sequential measurements and predict the next one

# Time Series Prediction

Model a sequence; Predict next

# Nearest Neighbour Classifier

Classification with Association by Similarity

# Nearest Neighbour Classifier

**X$_{test}$: [42.5, 39]**

$\rightarrow$ *Is the person diabetic or not?*

- How do we compare a test sample to a classification?

- We find distance to feature vectors of known classes - Association by Similarity

- We assign label of that sample which is nearest to the test sample

| BMI | Age | Diabetic |
|-----|-----|----------|
| 32.6 | 49 | 1 |
| 34.2 | 23 | 1 |
| 22.4 | 31 | 0 |
| 25.7 | 43 | 0 |
| 29.8 | 15 | 0 |
| 31 | 58 | 1 |
| 43.2 | 65 | 0 |
| 37.6 | 52 | 1 |

# Nearest Neighbour Classifier

| BMI | Age | Diabetic |
|------|-----|----------|
| 32.6 | 49 | 1 |
| 34.2 | 23 | 1 |
| 22.4 | 31 | 0 |
| 25.7 | 43 | 0 |
| 29.8 | 15 | 0 |
| 31 | 58 | 1 |
| 43.2 | 65 | 0 |
| 37.6 | 52 | 1 |

| Feature Vector | Label |
|----------------|-------|
| $X_1$ [32.6, 49] | 1 |
| $X_2$ [34.2, 23] | 1 |
| $X_3$ [22.4, 31] | 0 |
| $X_4$ [25.7, 43] | 0 |
| $X_5$ [29.8, 15] | 0 |
| $X_6$ [31.0, 58] | 1 |
| $X_7$ [43.2, 65] | 0 |
| $X_8$ [37.6, 52] | 1 |

iHub-Data-FMML 2023

# Nearest Neighbour Classifier

| Feature Vector | Label |
|---|---|
| $X_1$ [32.6, 49] | 1 |
| $X_2$ [34.2, 23] | 1 |
| $X_3$ [22.4, 31] | 0 |
| $X_4$ [25.7, 43] | 0 |
| $X_5$ [29.8, 15] | 0 |
| $X_6$ [31.0, 58] | 1 |
| $X_7$ [43.2, 65] | 0 |
| $X_8$ [37.6, 52] | 1 |

| Distance |
|---|
| 14.07 |
| 18.02 |
| 21.63 |
| 17.27 |
| 27.15 |
| 22.21 |
| 26.01 |
| 13.89 |

$X_{test}$: [42.5, 39]

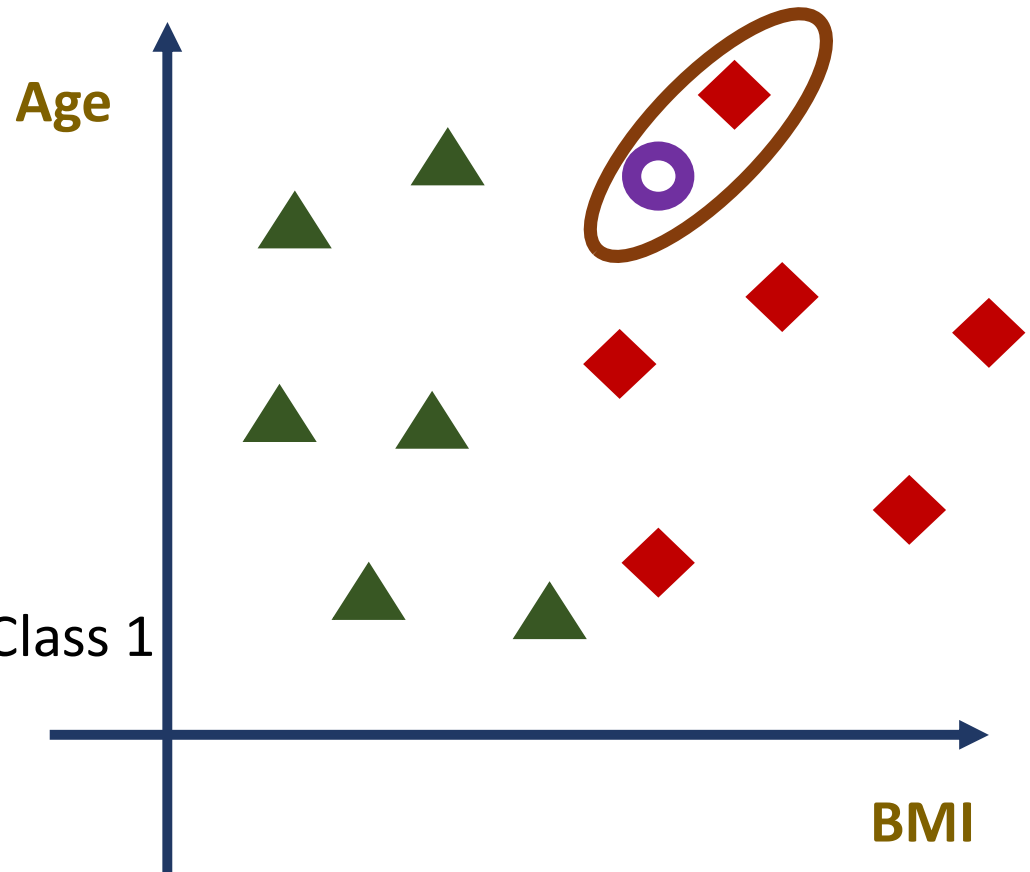$$\sqrt{(42.5 - 32.6)^2 + (39 - 49)^2}$$
$$= 14.07$$

✓ *The test sample has diabetes*

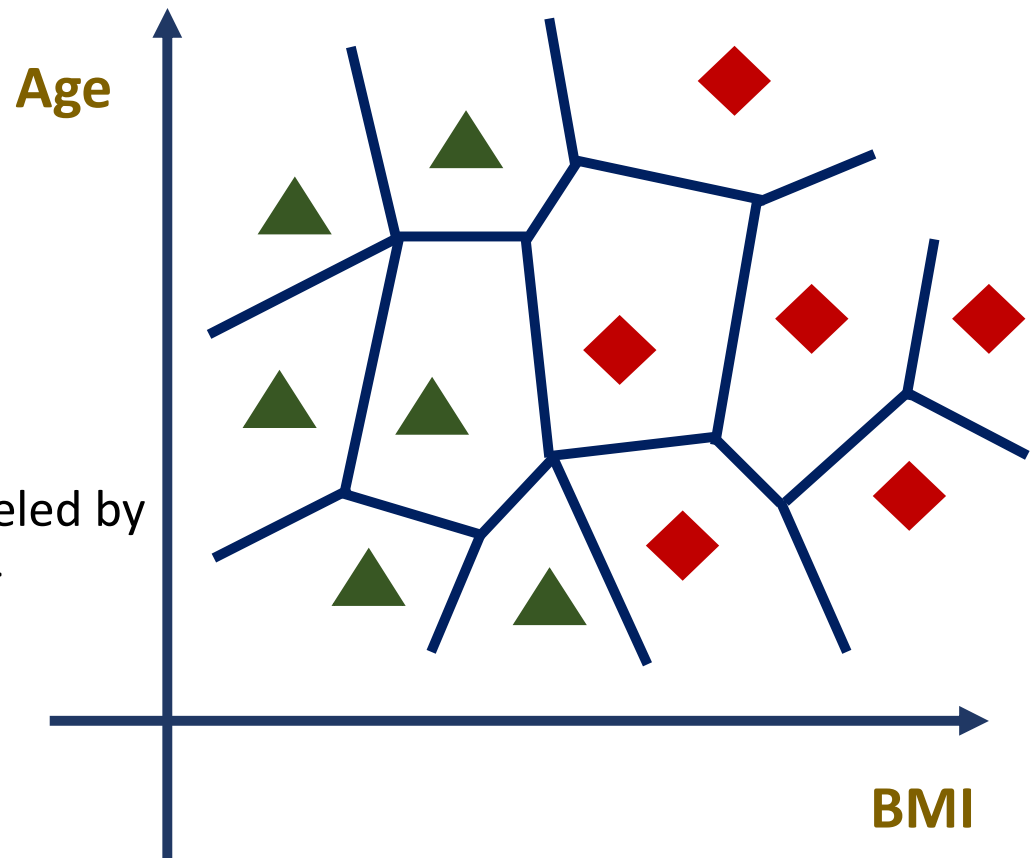# Visualization in Feature Space



Class 1

Class 0

- The nearest sample is a ◆

- We assign the test sample to Class 1
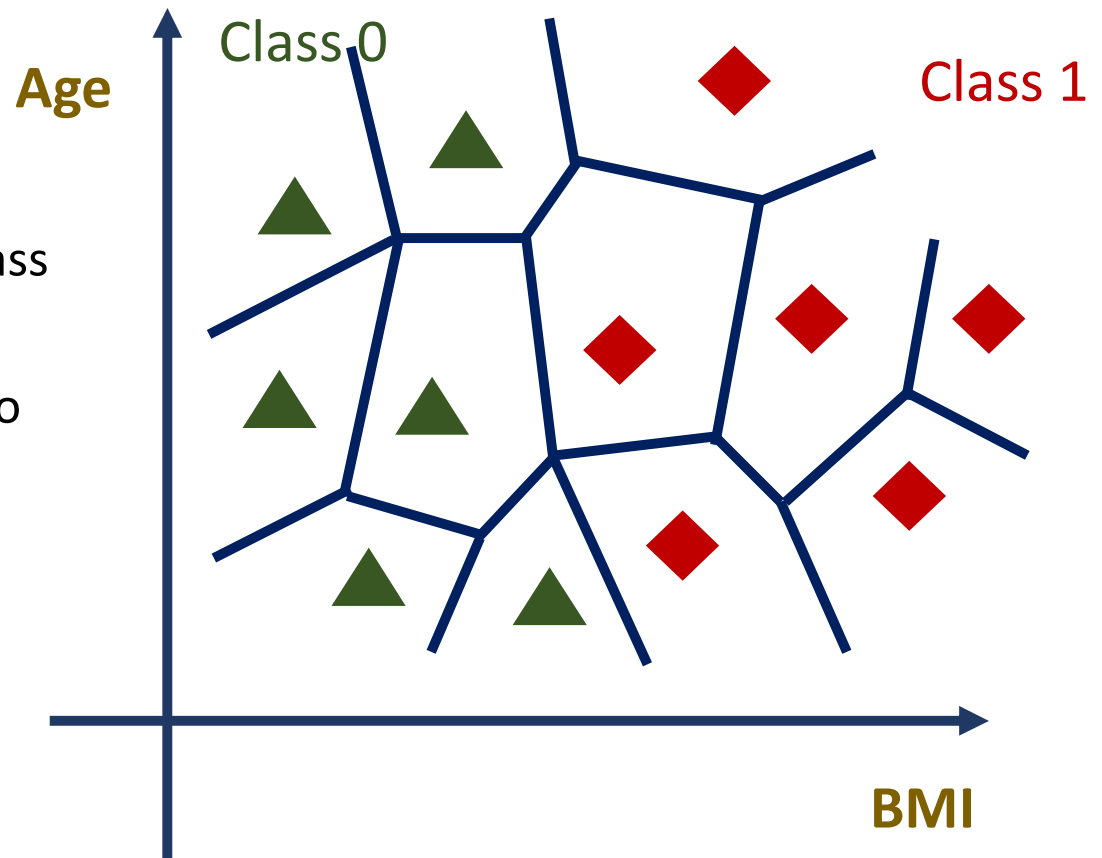
**Age**

**BMI**

iHub-Data-FMML 2023

# Class Boundaries

- The classifier effectively partitions the feature space into cells consisting of all points closer to a given training point $(x_1, x_2)$ than to any other training points.

- All points in such a cell are thus labeled by the category of the training point – **Voronoi tessellation** of the space

# Class Boundaries – Partition of Feature Space

- We now ignore boundaries between samples of the same class

- The decision boundary is found to be piece-wise linear

Class 0
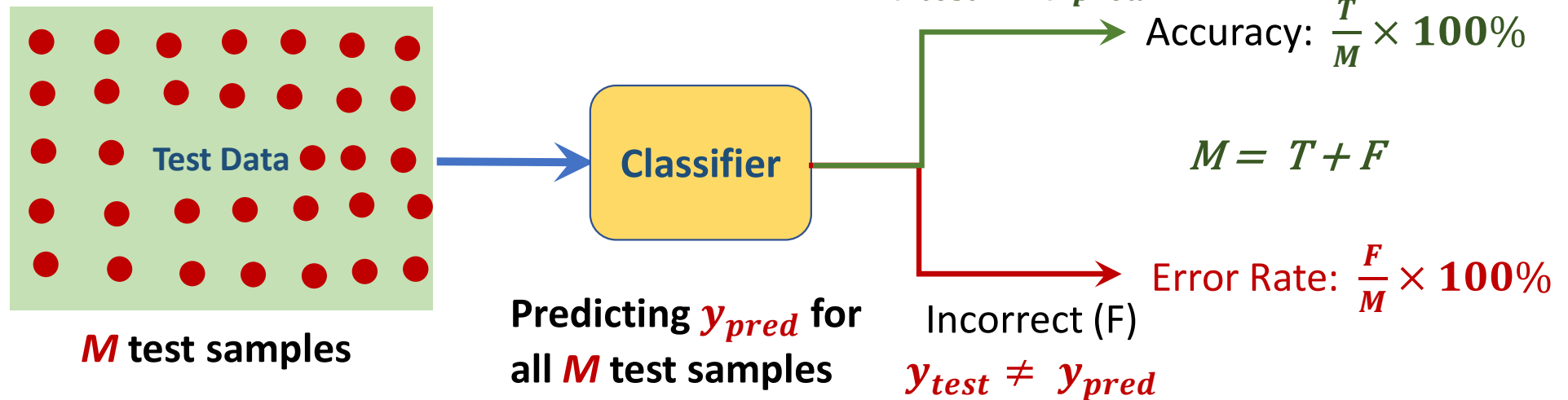
**Age**

Class 1

**BMI**

iHub-Data-FMML 2023

# The Classification Algorithm

**Problem**

- Given:
  - A set of training $n$ training samples: $(x_i, y_i)$
  - A set of $m$ test samples: $(x_{test}, y_{test}), m \ll n$
- Find:
  - Label $(x_{test})$ using Similarity Measure and return $y_{pred}$
- Find accuracy of the prediction
  - Evaluation of a classifier

# Evaluating a Classifier

- Several Metrics used to evaluate a classifier



Correct (T), $y_{test} = y_{pred}$

Accuracy: $\frac{T}{M} \times \mathbf{100\%}$

$M = T + F$

Error Rate: $\frac{F}{M} \times \mathbf{100\%}$

Incorrect (F) $y_{test} \neq y_{pred}$

**Test Data**

**$M$ test samples**

**Classifier**

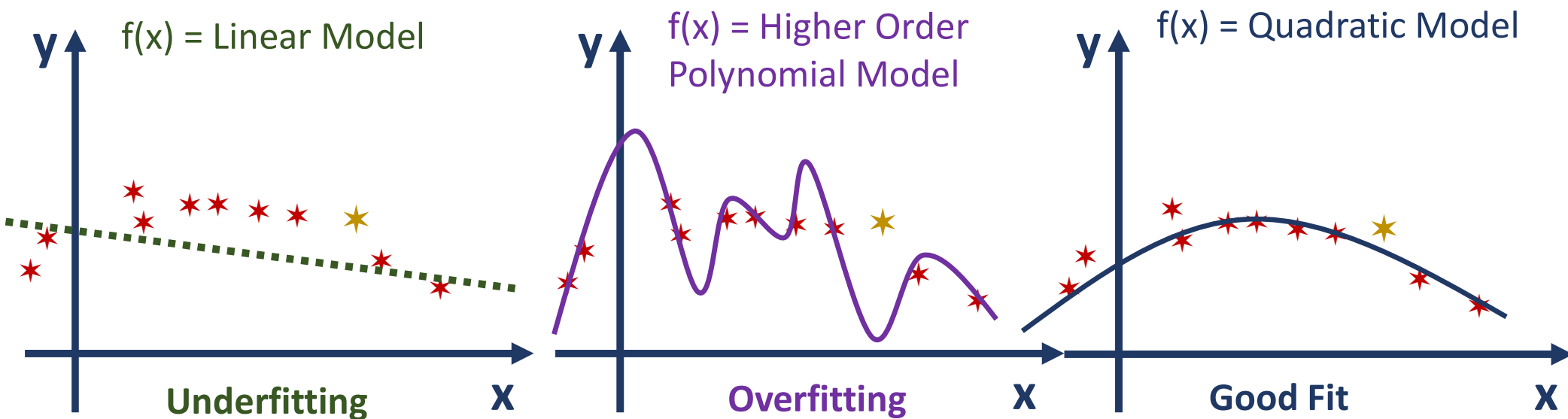**Predicting $y_{pred}$ for all $M$ test samples**

- In a 2-class classifier, what does an accuracy of 65% tell us about the classifier?

# Supervised Machine Learning - Methodology

- **Step 1**: A set of training $n$ training samples: $(x_i, y_i), i = 1, .., n$

- **Step 2**: We need to correctly predict labels of unseen $m$ test sample: $(x_{test}, y_{test}), test = 1, .., m$. Predicted value = $y_{pred}$

- **Step 3**: We need to maximize the accuracy on unseen $m$ test sample: $(x_{test}, y_{test}), test = 1, .., m$. $y_{pred} = y_{test}$

- **Assumption**: Test samples come from the same distribution as training samples

- Can we have situations where we do well on training samples but perform badly on test samples? Rote Learning / Memorization

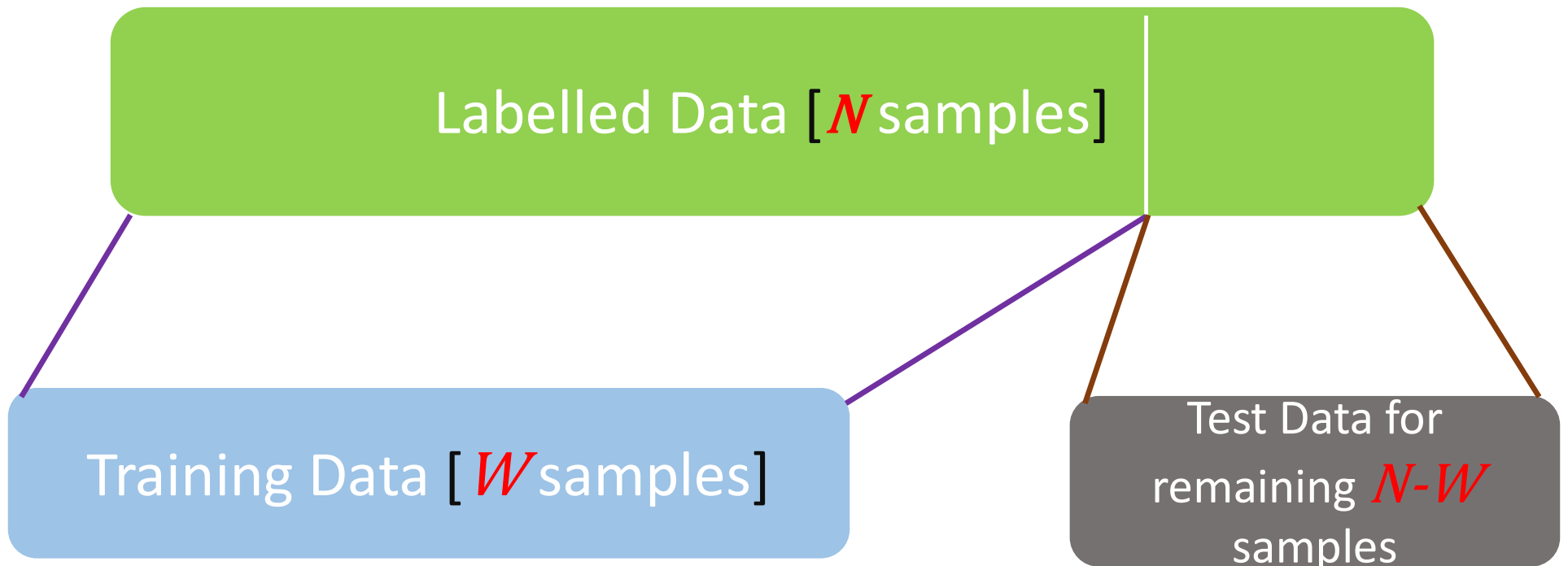# Overfitting vs Generalization

- We try to fit the data with different models in various orders of complexity



Which of the above model will perform best on the unseen test data?

# ML Based on Training-Testing Data

Labelled Data [$N$ samples]

Training Data [$W$ samples]

Test Data for remaining $N$-$W$ samples

- Take care to not leak information from Test Data into the Model with repeated testing

# ML Based on Training - Validation - Testing Data

Labelled Data [$N$ samples]

Training Data [$W$ samples]

Validation Data [$X$ samples]

Test Data for remaining $N$-$W$-$X$ samples

- The validation model is repeatedly used during development
- The test data is used once for the final prediction

iHub-Data-FMML 2023

Training Data / Validation Data → Feature extraction → Choose a model → HyperParams of model: $H_j$ Goal: to predict $f()$ depending on $\Theta_i$

Test Data → Feature extraction

Select $H_j$ and for $\Theta_i$, learn about $f()$ from training data

Evaluate the model using Validation Data and $\Theta_i$

Are there any other ML models

Choose the best model, Classifier/Predictor

Compute accuracy for the test data using $\Theta_i$ of the best classifier

Yes

No

iHub-Data-FMML 2023