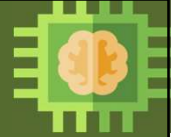


Elective Course

Course Code: CS4103

Autumn 2025-26



Lecture #34

Artificial Intelligence for Data Science

Week-9:

MACHINE LEARNING (Part II)

Classification using K-Nearest Neighbor (KNN) Algorithm

Course Instructor:

Dr. Monidipa Das

Assistant Professor

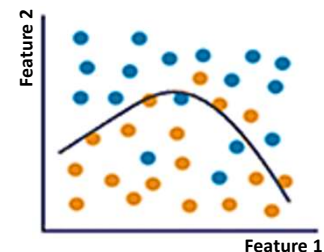
Department of Computational and Data Sciences

Indian Institute of Science Education and Research Kolkata, India 741246

Classification



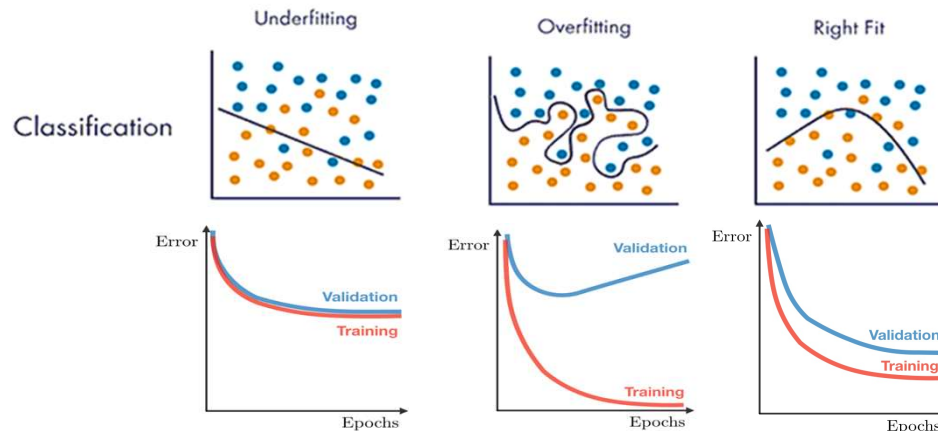
- Computer program asked to specify which of k categories some input belongs to
 - Learning algorithm is asked to produce a function $f: \mathbb{R}^n \rightarrow \{1, \dots, k\}$, where n = no of input variables
 - When $y = f(x)$ model assigns input vector x to a category identified by a numeric code y
- Other variants of classification task:
 - f outputs a probability distribution over classes



Generalization, Underfitting and Overfitting



In ML, **generalization** is the ability to perform well on previously unobserved inputs



Underfitting: Inability to obtain low enough error rate on the training set

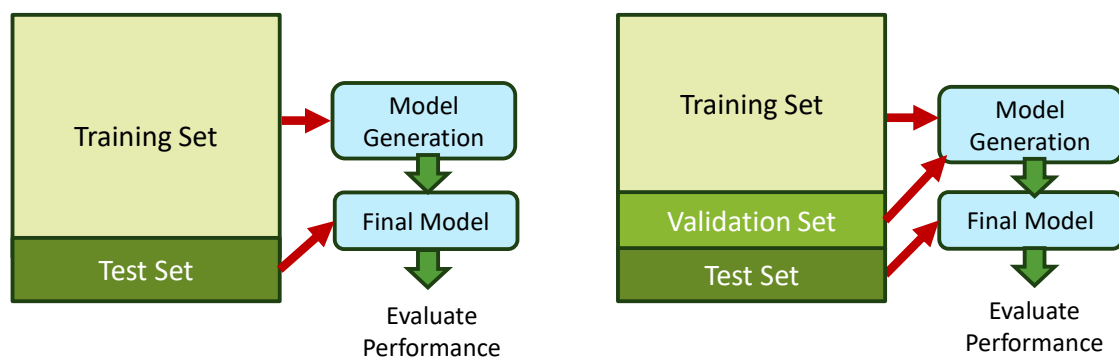
Overfitting: Gap between training error and testing error is too large

Dr. Monidipa Das, Department of CDS, IISER Kolkata

Data Splitting



- Training Set:** For training of the model.
- Validation Set:** For unbiased evaluation of the model.
- Test Set:** For final evaluation of the model.

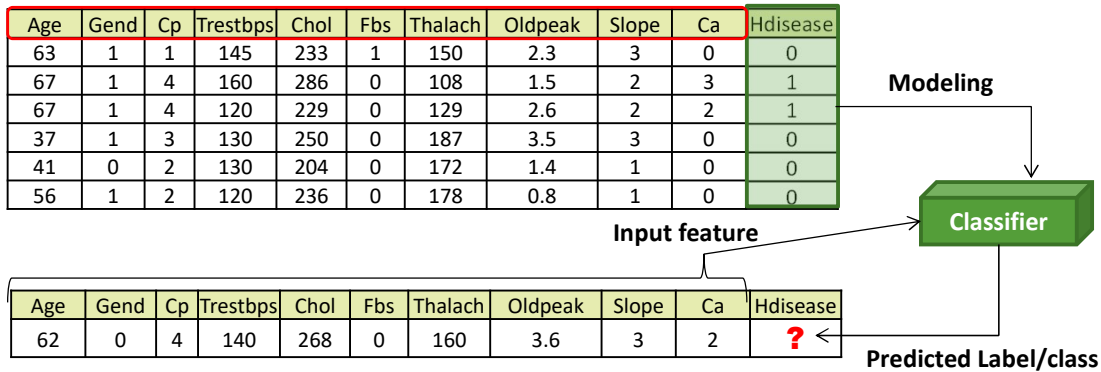


Dr. Monidipa Das, Department of CDS, IISER Kolkata

How does classification work?

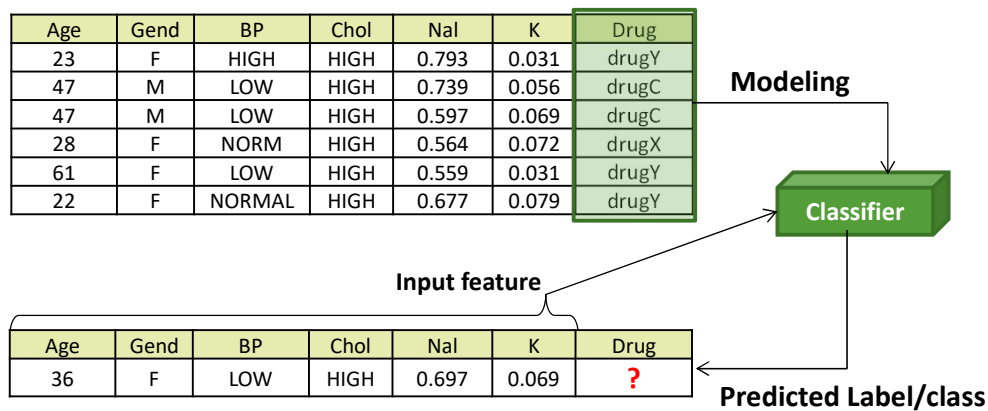


- Classification determines the class label for an unlabeled data point (test case)



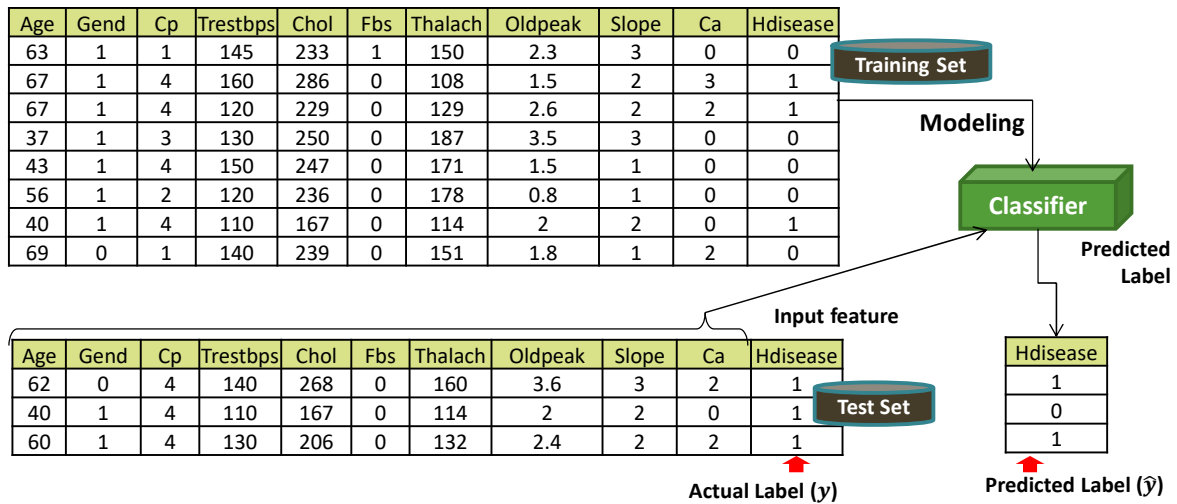
Dr. Monidipa Das, Department of CDS, IISER Kolkata

Multi-class Classification: Example



Dr. Monidipa Das, Department of CDS, IISER Kolkata

Evaluation Metrics for Classification



Dr. Monidipa Das, Department of CDS, IISER Kolkata

Evaluation Metrics for Classification: Jaccard Index



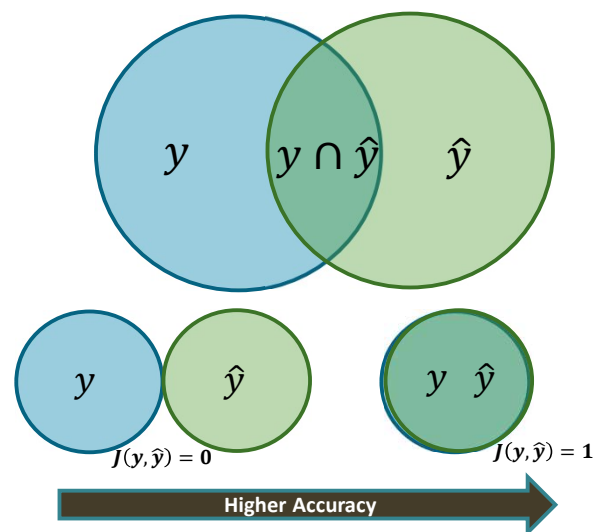
- y : Actual labels
- \hat{y} : Predicted labels

$$J(y, \hat{y}) = \frac{|y \cap \hat{y}|}{|y| + |\hat{y}| - |y \cap \hat{y}|}$$

$$y: [1, 0, 0, 0, 0, 1, 1, 0, 1, 0, 1]$$

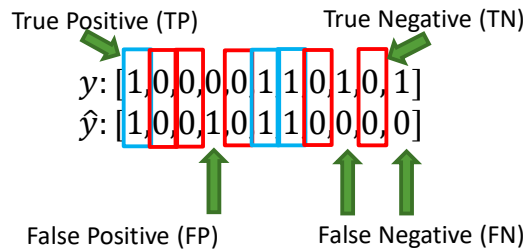
$$\hat{y}: [1, 0, 0, 1, 1, 0, 1, 1, 0, 0, 0]$$

$$J(y, \hat{y}) = \frac{8}{11 + 11 - 8} = 0.57$$



Dr. Monidipa Das, Department of CDS, IISER Kolkata

Evaluation Metrics for Classification: F1-Score



$$F1\text{-score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{(\text{Precision} + \text{Recall})}$$

$$\text{Precision} = \frac{TP}{(TP + FP)}$$

$$\text{Recall} = \frac{TP}{(TP + FN)}$$

	Precision	Recall	F1-score
Hdisease=1 (Positive)	0.75	0.60	0.67
Hdisease=0 (Negative)	0.71	0.83	0.77

Average $\rightarrow 0.72$

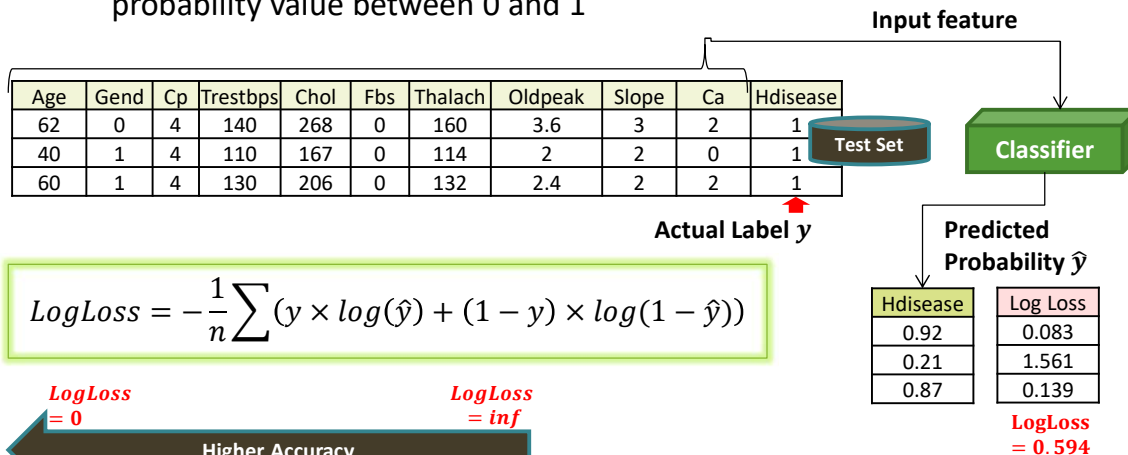
Confusion Matrix	
Actual Label	Hdisease=1 (Positive)
	3 (TP)
Actual Label	Hdisease=0 (Negative)
	1 (FP)
Actual Label	Predicted Label
	2 (FN)
Actual Label	Predicted Label
	5 (TN)

Dr. Monidipa Das, Department of CDS, IISER Kolkata

Evaluation Metrics for Classification: Log Loss



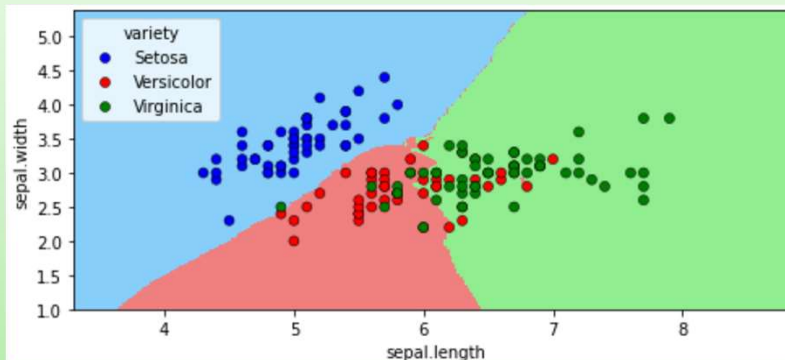
- Performance of a classifier where the predicted output is a probability value between 0 and 1



Dr. Monidipa Das, Department of CDS, IISER Kolkata



K-Nearest Neighbors (KNN)



Dr. Monidipa Das, Department of CDS, IISER Kolkata

Determining class using KNN



Age	Gend	Cp	Trestbps	Chol	Fbs	Thalach	Oldpeak	Slope	Ca	Hdisease
63	1	1	145	233	1	150	2.3	3	0	0
67	1	4	160	286	0	108	1.5	2	3	1
67	1	4	120	229	0	129	2.6	2	2	1
37	1	3	130	250	0	187	3.5	3	0	0
41	0	2	130	204	0	172	1.4	1	0	0
56	1	2	120	236	0	178	0.8	1	0	0

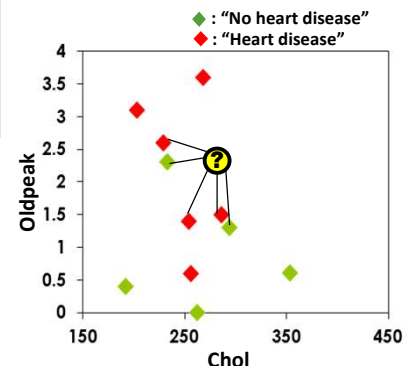
Age	Gend	Cp	Trestbps	Chol	Fbs	Thalach	Oldpeak	Slope	Ca	Hdisease
54	1	4	124	266	0	109	2.2	2	1	?

0 → "No heart disease"

← 1-NN

1 → "Heart disease"

← 5-NN

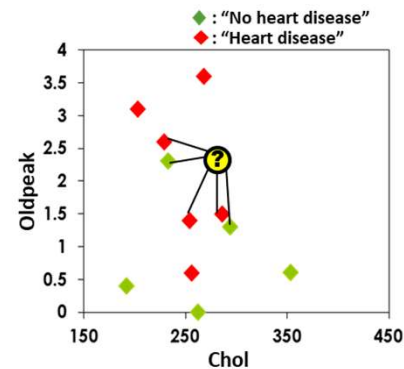


Dr. Monidipa Das, Department of CDS, IISER Kolkata

K-Nearest Neighbors (KNN)



- A method for **classifying** cases/data-points based on their similarity to others cases/data-points
- Data-points that are near to each other are said to be "**neighbors**"
- Assume data-points with the same class labels are near to each other.



Dr. Monidipa Das, Department of CDS, IISER Kolkata

KNN Algorithm (basic steps)



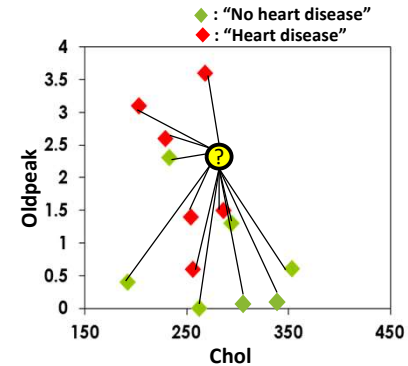
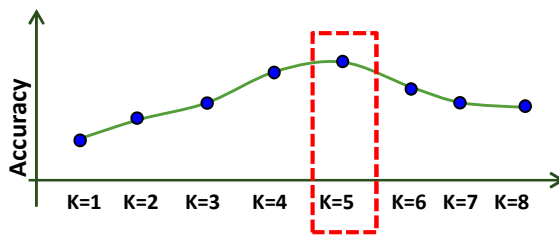
- Pick a value for K
- Calculate the distance of unknown data-points from all the known data-points
- Select K known data-points that are "nearest" to the unknown data-points
- Get the labels of the K selected data-points
- Return the mode of the K labels

Dr. Monidipa Das, Department of CDS, IISER Kolkata

Finding the best value for K

15

- K=1 Class label 0
- K=7 Class label 1
- K=13 ?



Better to avoid even number as a value for K in the KNN algorithm to prevent ties

Dr. Monidipa Das, Department of CDS, IISER Kolkata

Calculating Distance

16

- Most common way--- **Euclidean Distance**

$$x_1 = (x_{11}, x_{12}, \dots, x_{1m})$$

$$x_2 = (x_{21}, x_{22}, \dots, x_{2m})$$

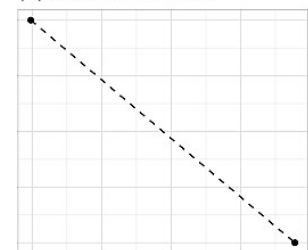
$$EDist(x_1, x_2) = \sqrt{(x_{11} - x_{21})^2 + (x_{12} - x_{22})^2 + \dots + (x_{1m} - x_{2m})^2}$$

$$= \sqrt{\sum_{i=1}^m (x_{1i} - x_{2i})^2}$$

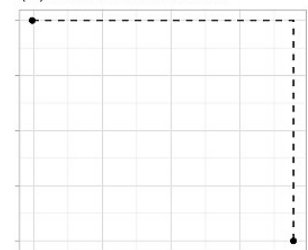
$$MDist(x_1, x_2) = \sum_{i=1}^m |x_{1i} - x_{2i}|$$

$$\text{Minkowski distance: } Dist(x_1, x_2) = \left[\sum_{i=1}^m |x_{1i} - x_{2i}|^p \right]^{1/p}$$

(A) Euclidean distance



(B) Manhattan distance



Dr. Monidipa Das, Department of CDS, IISER Kolkata

Calculating Distance

Patient-1 (x_1)Patient-2 (x_2)

Age	Chol	Oldpeak
54	236	0.8
67	284	1.5

$$EDist(x_1, x_2) = \sqrt{(54 - 67)^2 + (236 - 284)^2 + (0.8 - 1.5)^2} = 49.7$$

$$MDist(x_1, x_2) = |54 - 67| + |236 - 284| + |0.8 - 1.5| = 61.7$$

Dr. Monidipa Das, Department of CDS, IISER Kolkata

Measuring Distance for Categorical Features



- Hamming Distance**

Data point 1

0	0	0	1	1	0	0
1	1	0	1	1	1	0

Data point 2

Hamming distance is: 3

When all the attributes/features are categorical:

	Color	Shape	Size
Data point 1	Red	Circle	Small
Data point 2	Green	Rectangle	Small

mismatch mismatch match

Hamming distance is: 1 + 1 + 0 = 2

Using One-Hot Encoding for Categorical Attributes

	Color	Weight	Length
Data point 1	Red	40.5	5
Data point 2	Green	53	3
Data point 3	Blue	78.5	15

[0,0,1] Red
 [0,1,0] Green
 [1,0,0] Blue

Hamming distance is: 2

Dr. Monidipa Das, Department of CDS, IISER Kolkata

Feature Scaling



Dataset

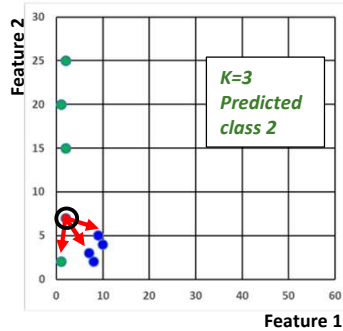
Feature1	Feature2	Class
1	2	1
2	15	1
1	20	1
2	25	1
10	4	2
9	5	2
7	3	2
8	2	2
2	7	?

● class 1
● class 2

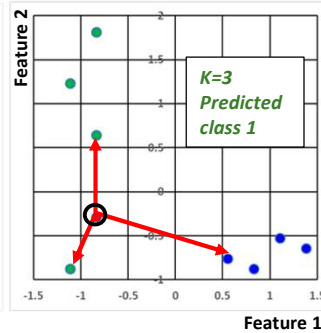
$$x' = \frac{x - \bar{x}}{\sigma}$$

$$x' = \frac{x - \min(x)}{\max(x) - \min(x)}$$

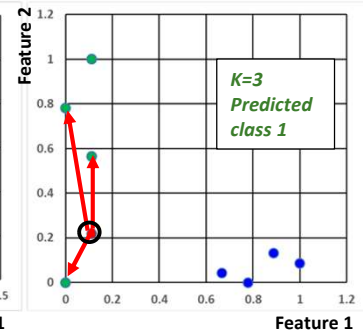
Original Features



Features after Standardization
(Z-score Normalization)



Features after Min-Max
Normalization



Dr. Monidipa Das, Department of CDS, IISER Kolkata



Questions?

Dr. Monidipa Das, Department of CDS, IISER Kolkata