

COMP307/AIML420 Week 3 (Tutorial)

1. Announcements

- Assignment 1 (**15%**)
- Helpdesk sessions

2. Sets

- Training and Test sets
- Validation set

3. Datasets

- Instances
- Features and feature vectors
- Class label

4. 3-K Algorithms

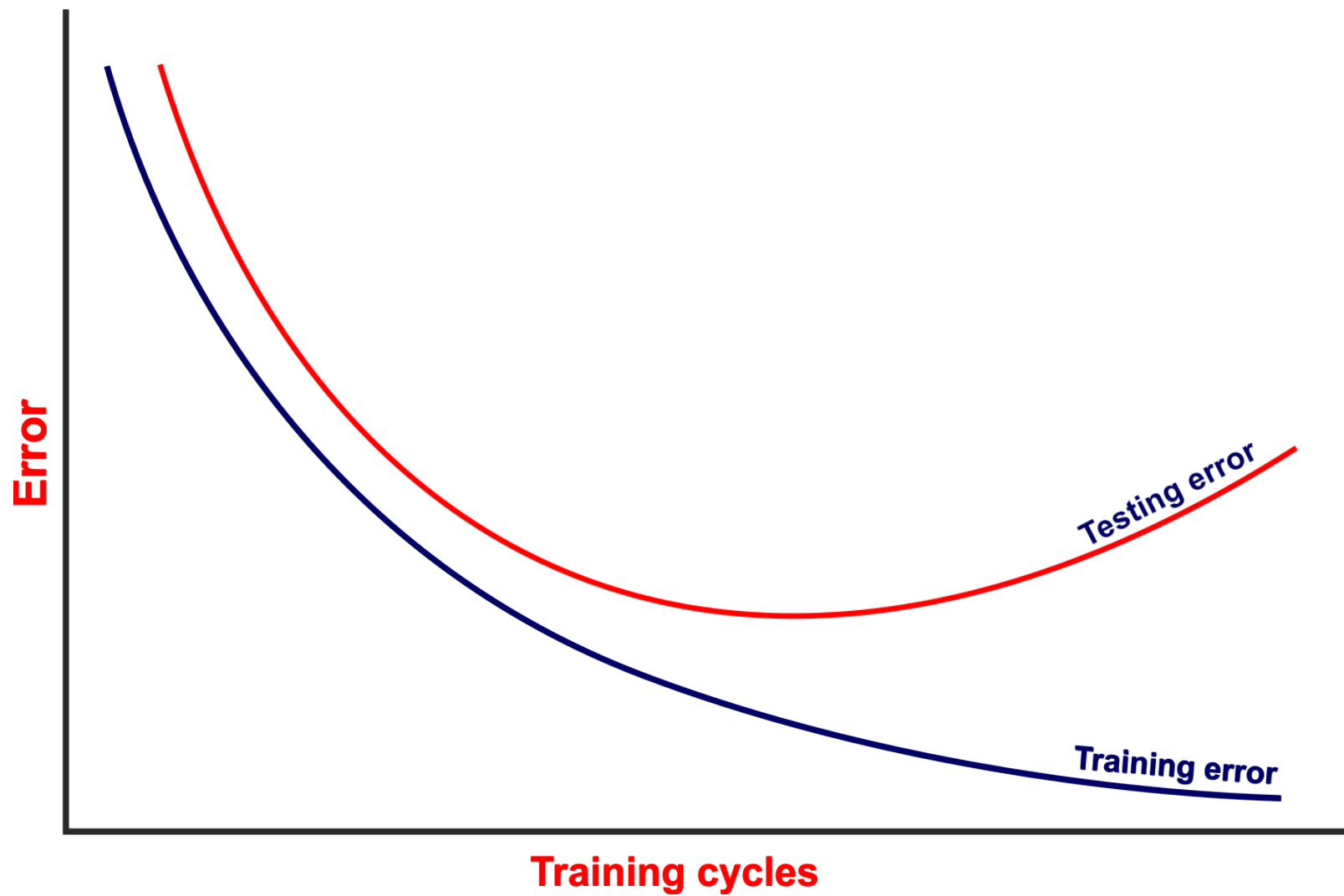
- k-Nearest Neighbour
- k-Means Clustering
- k-fold Cross Validation

5. Decision Trees

- DT learning vs learned DT
- Impurity measure Conditions
- Pruning

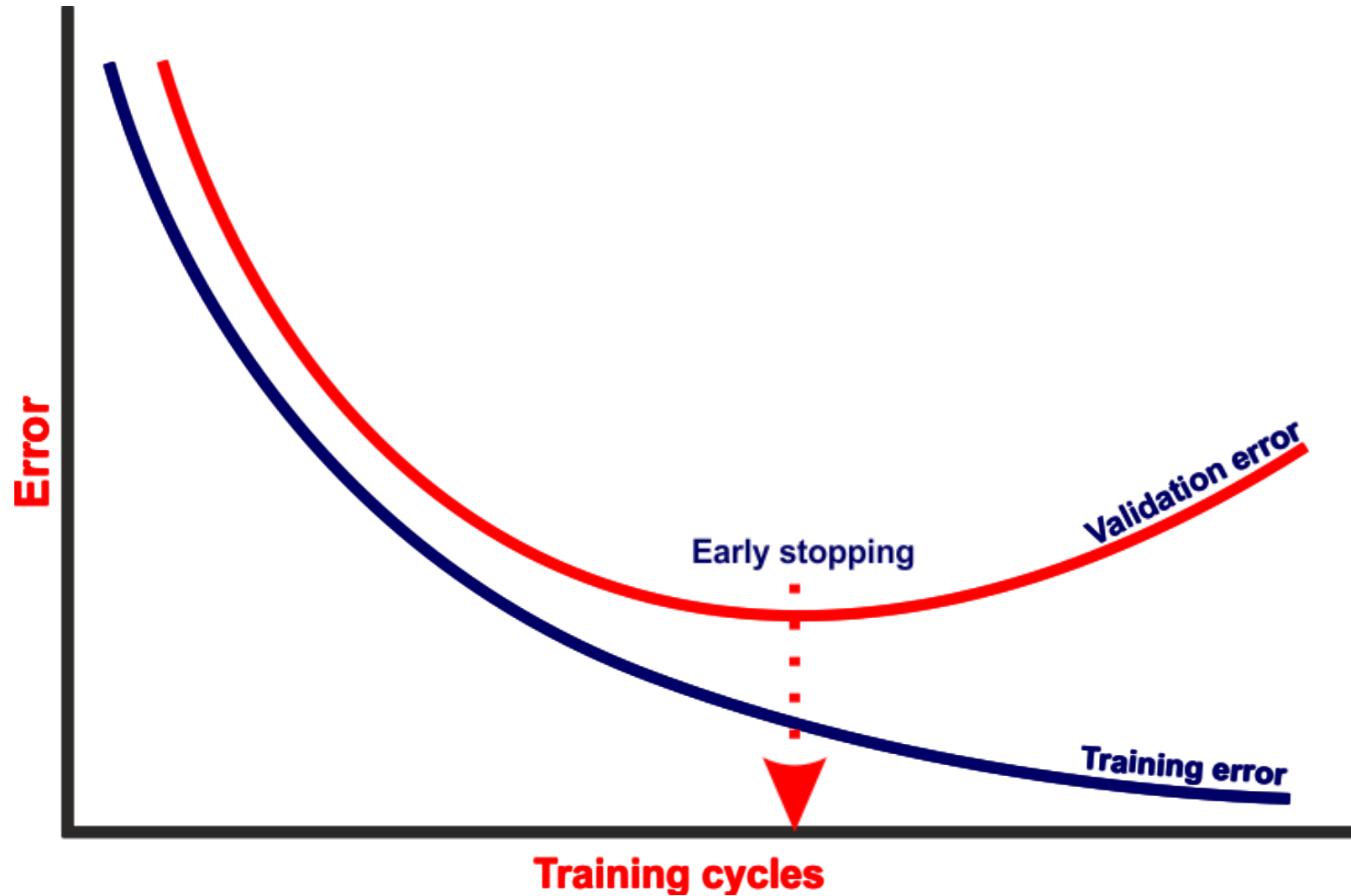
6. Other Questions

Over-fitting

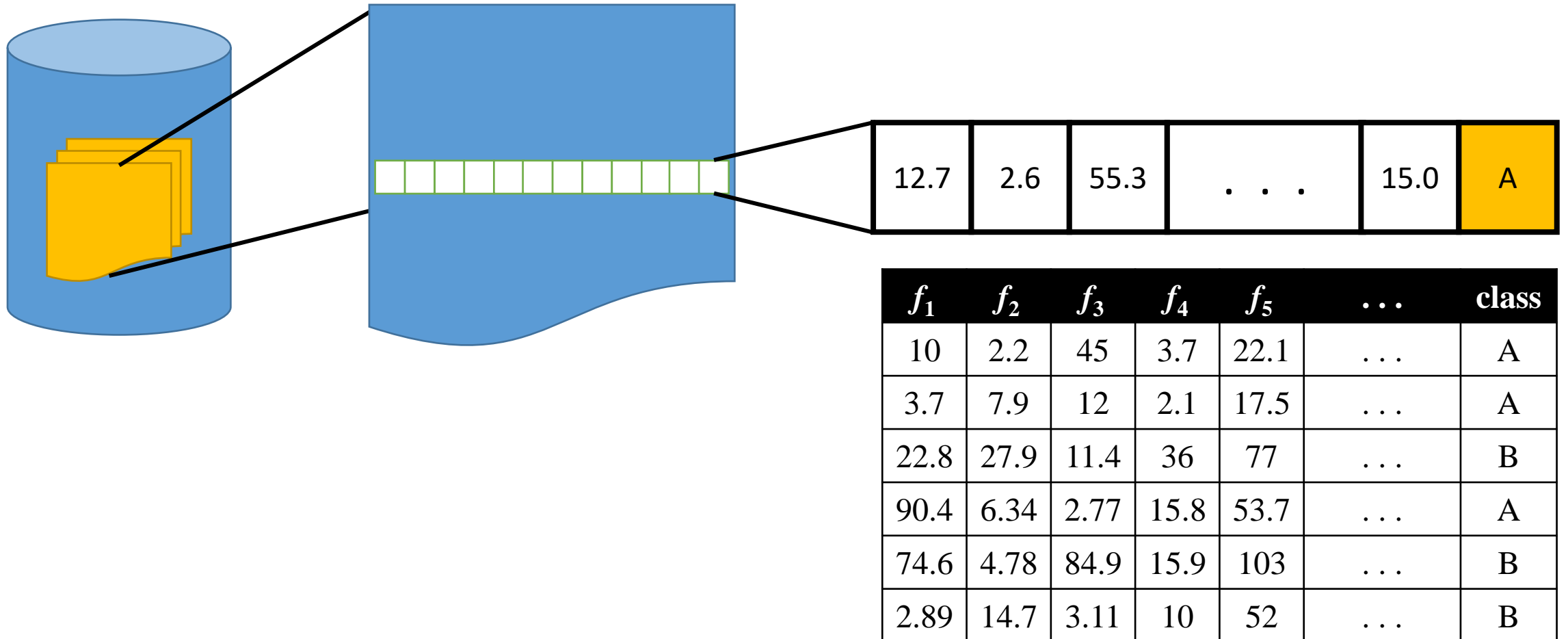


Validation Set

- What?
- Why?
- How?



Datasets and Instances



K-Nearest Neighbour

Training Set

f_1	f_2	f_3	class
10.3	45.7	2.7	A
7.1	80.5	1.1	A
22.3	20.4	9.6	B
30.5	21.2	17.9	B
5.2	67.1	7.7	A
15.6	18.6	11.4	B
11.9	53.4	6.3	A

$$d(\vec{u}, \vec{v}) = \sqrt{\sum_{i=1}^N (\vec{u}_i - \vec{v}_i)^2}$$

$$d(\cdot, \cdot) = 14.84$$

$$d(\cdot, \cdot) = 47.40$$

$$d(\cdot, \cdot) = 24.57$$

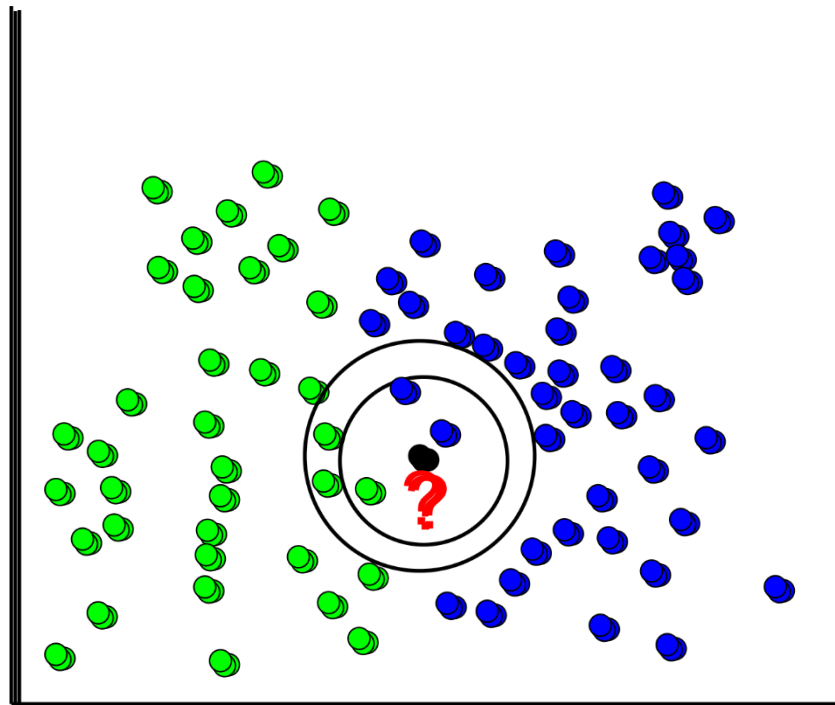
$$d(\cdot, \cdot) = 33.65$$

$$d(\cdot, \cdot) = 33.88$$

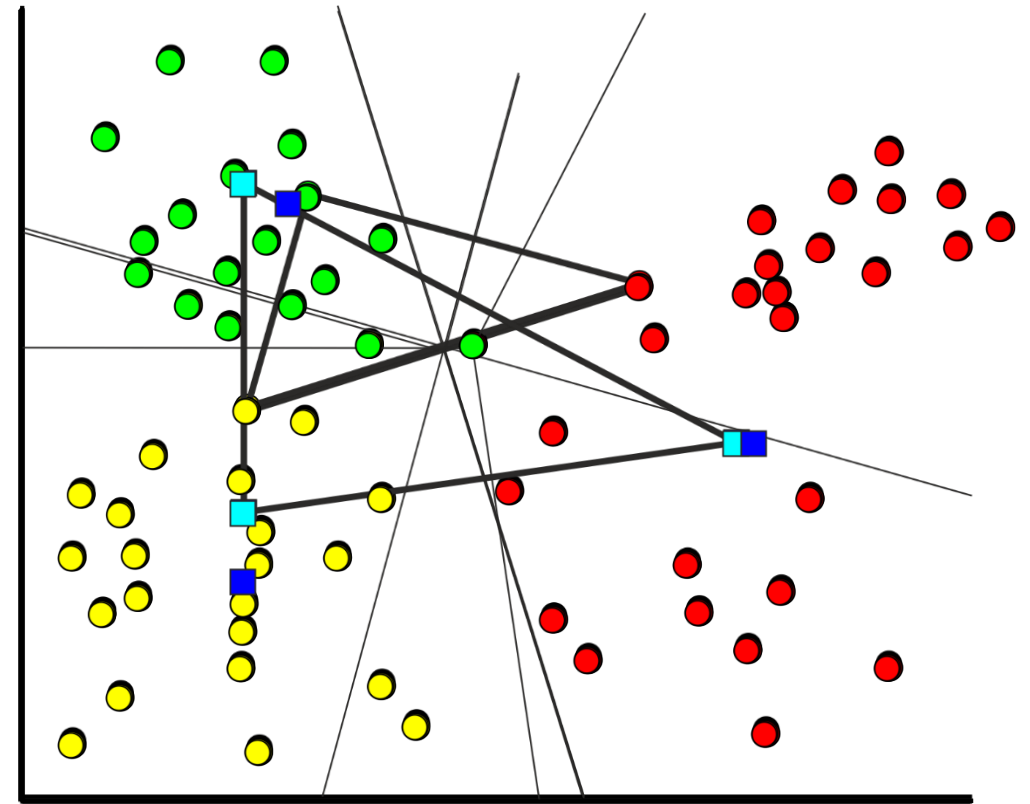
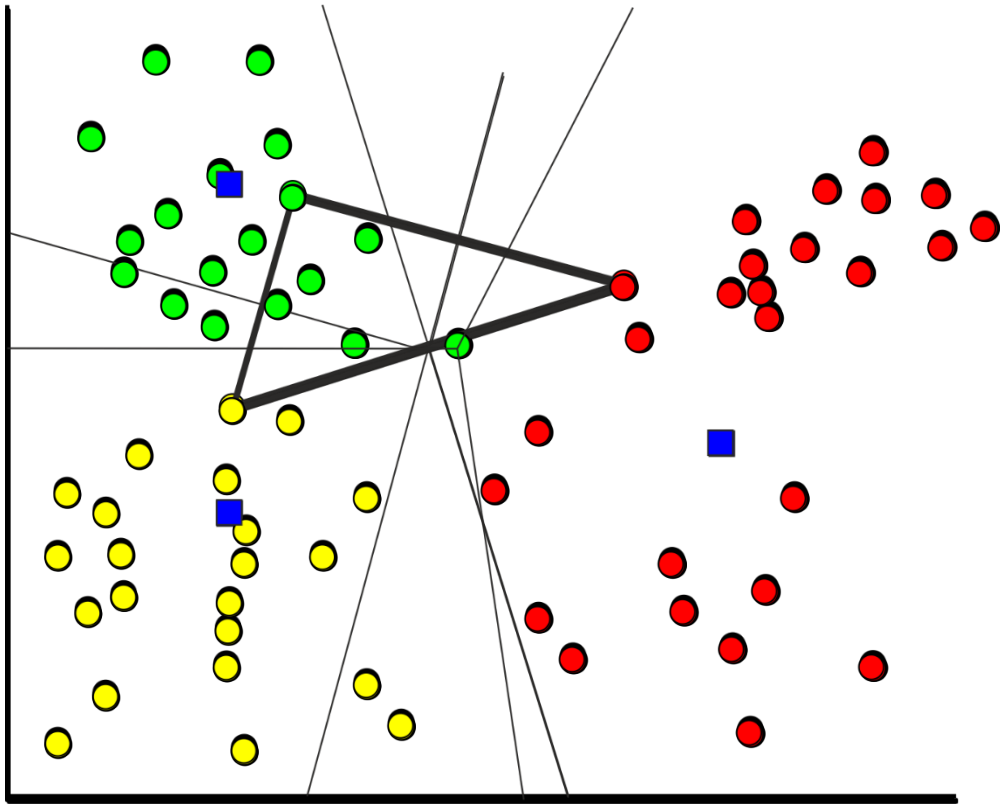
$$d(\cdot, \cdot) = 21.19$$

$$d(\cdot, \cdot) = 22.24$$

2.1	33.5	4.7	A
-----	------	-----	---

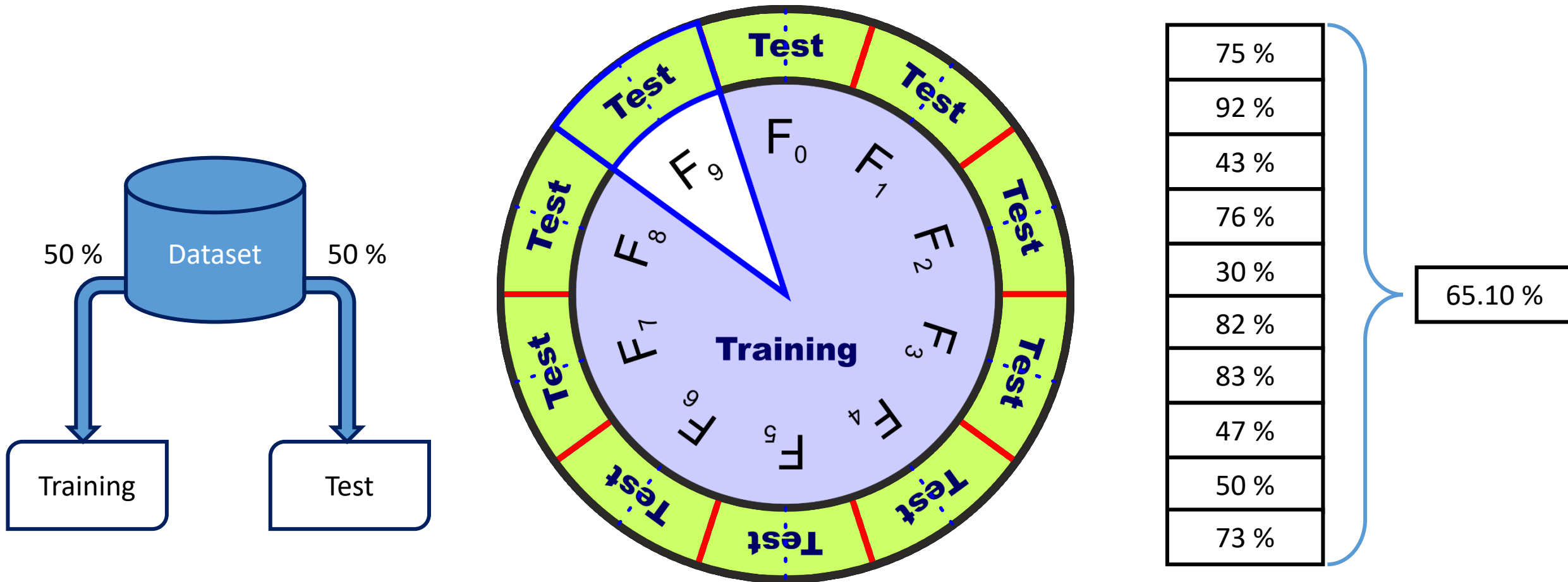


k-Means Clustering



Demo

k-fold Cross Validation



- How do we specify the number of instances in each fold?
- How do we select those instances?

Decision Trees (DT)

- DT learning \neq learned DT
- Impurity measure
 1. 0, if all instances belong to one class
 2. Max, equal number of instances for both classes
 3. Continuous –smooth
- Large vs. small trees?

Other Questions

- Can k-NN cope with categorical data/features?
- Does k-Means always provide the same clusters?
- How do we select/initialise the seeds in k-Means?