

- ① K-means clustering
- ② K-nearest neighbors
- ③ neural network

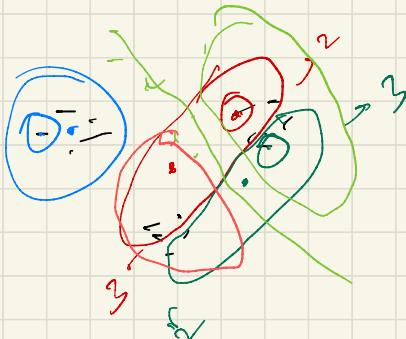
• Selection (Initialization)

Iteration:

Random

① Assignment

② Re-center



$$x_1, x_2, \dots, x_n \in \mathbb{R}^d$$

$$x \rightarrow p_1, p_2, \dots, p_k$$

$$p_1 \cup p_2 \cup \dots \cup p_k = X$$

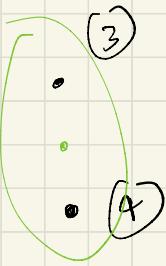
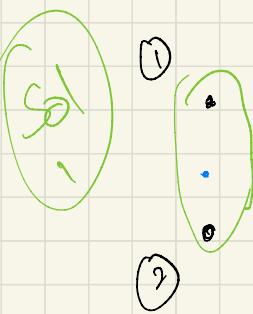
$$p_i \cap p_j = \emptyset$$

$$\left\{ \text{cost}(X, C) = \sum_{i=1}^n \|x_i - c(i)\|_2^2 \mid |C|=k \right. \\ \left. c(i) \in C \right\}$$

• Lloyd's algorithm

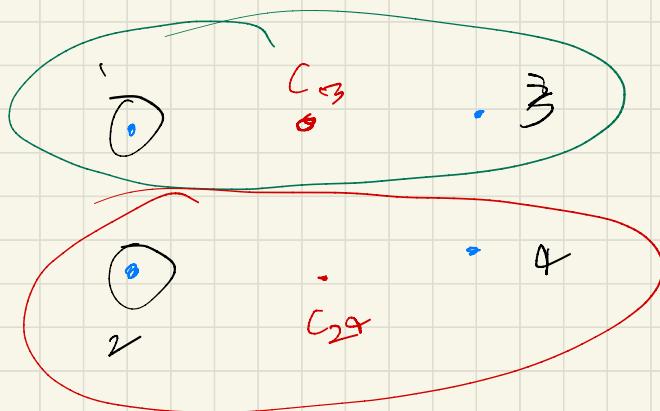
→ K-means algorithm

$$c(i) = \min_{c \in C} \|x - c\|_2$$



$$\sum_{i=1}^4 (x_i - c_i)^2$$

Random seeding $C_i = \{1, 2\}$



① Assignment

$$1 \rightarrow 1$$

② Re-center

$$c_{13}$$

$$2 \rightarrow 2$$

$$c_{24}$$

$$3 \rightarrow 1$$

$$4 \rightarrow 2$$

K-means ++: ^{"2008"}

- Seeding

- K-means Algorithm

C_1, C_2, \dots, C_k , $C_i = \{c_1, c_2, \dots, c_{i-1}\}$

① First center is picked v.a.r.

② For $i = 2$ to k :

$$P_{2n}[C_i = x_j] \propto \text{cost}(x_j, C_{i-1})$$

$$\text{cost}(x_j, C_{i-1}) = d(x_j, C_{i-1})$$

$O(nkd)$

$$\begin{array}{ccccccc} & & 10 & & & 100 & \\ 1 & \xrightarrow{\quad} & & \xrightarrow{\quad} & 3 & \xrightarrow{\quad} & \frac{100}{202} \\ & & & & & & \\ & & & & & 4 & \xrightarrow{\quad} \frac{101}{202} = \frac{1}{2} \\ & & & & & & \\ & & & & & 101 & \end{array}$$

$$\frac{1}{202} \checkmark$$

$$P_{2n}[C_2 = x_2] \propto 1$$

$$P_{2n}[C_2 = x_3] \propto 100$$

$$P_{2n}[C_2 = x_4] \propto 100$$

How to choose k :

, Elbow method

, $k=1 \rightarrow 50$



k -nearest neighbours:

- Linear regression, neural networks,

parametric models.

$$\hookrightarrow \mathbf{x} \in \mathbb{R}^d, \quad \mathbf{w} \in \mathbb{R}^d, \quad b \in \mathbb{R}.$$

$$f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$$

$$\min \sum_{i=1}^n \| \mathbf{w}^\top \mathbf{x}_i + b - y_i \|^2$$

- Non-parametric

\mathbf{x}_i, y_i

\mathbf{Wx}

- Table look-up-

- k -nearest neighbours

"Metric"

→

1-nearest neighbour

K-nn:

- ① Find the "k" nearest neighbours to the query point
- ② Output the mode of all the labels of the k-nn.

$$k \rightarrow \text{"odd"}$$

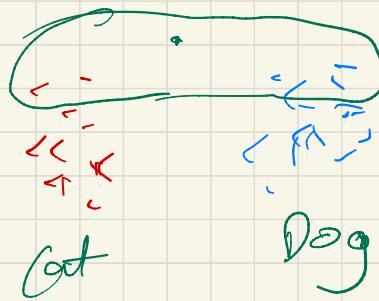
(F)

Loss fn:

Train

Test

0-1 error



Cat

Dog

2

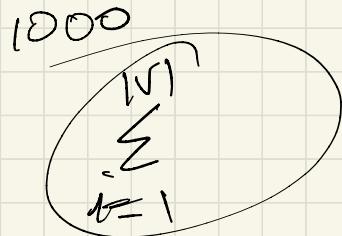
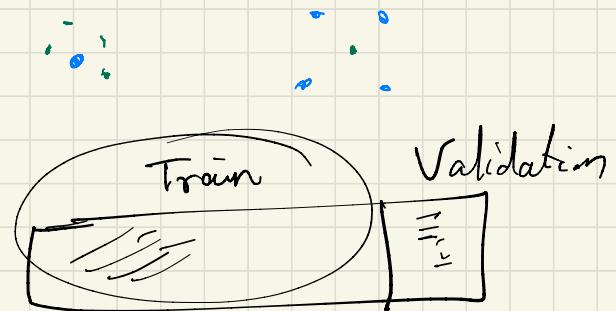
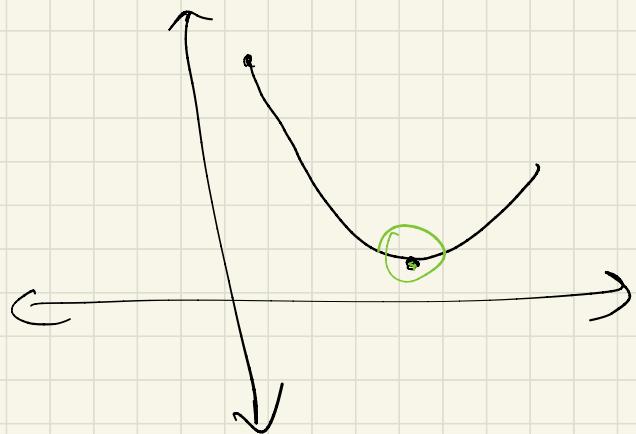
3

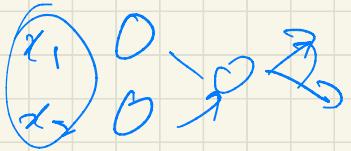
Over-fitting
Under-fitting

How to Pick "k":

Cross-validation

$k=1, 3, 5, 7, 9, \dots, 99$

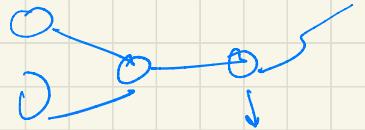




③

w

x_1
 x_2
1



②

3×100

