

5 fold or 10 fold is very common.

- Use $k-1$ fold for training, 1 fold for testing

- $ER = \frac{1}{k} \sum_{i=1}^k ER_i$

Important: all classes w_1, \dots, w_c well represented in all folds.

(B) Training, Validation, Test set
 ↓
 hyperparameter optimization

4-110

4-111

5 Unsupervised Learning

5-1

5-2

5-3

5-4

5-5

5.1 Clustering for a known number of clusters

5.1.1 Problem formulation

Given:

- N samples $S = \{x_n \in \mathbb{R}^d, 1 \leq n \leq N\}$
- number of clusters C

Wanted:

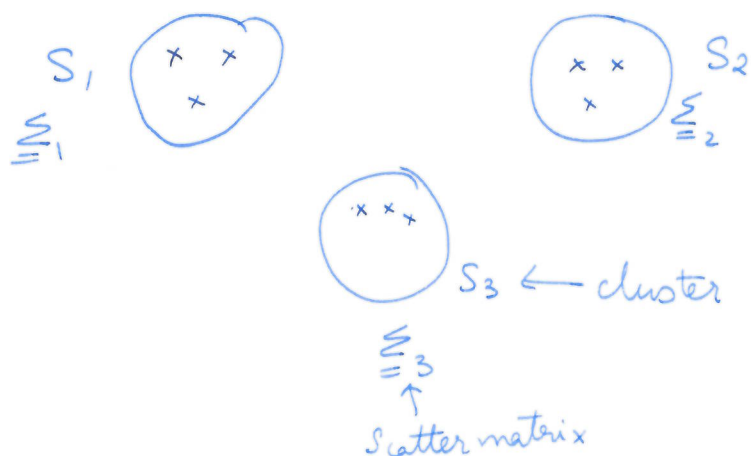
- C clusters $S_i \subset S; 1 \leq i \leq C$ $N_i = \frac{|S_i|}{\text{samples in } S_i}$
- each sample exactly in one cluster
i.e. $\sum N_i = N$

Wishes:

- Samples of one cluster stay "together"
- Samples of different clusters stay "far away"

(3)

Measure for "closeness": scatter matrix \sim covariance matrix



(A) For each cluster S_i :

$$\text{center } \underline{\mu}_i = \frac{1}{N_i} \sum_{x_i \in S_i} x_i$$

$$\text{scatter matrix } \underline{\Sigma}_i = N_i \underbrace{\frac{1}{N_i} \sum_{x_n \in S_i} (x_n - \underline{\mu}_i)(x_n - \underline{\mu}_i)^T}_{\text{sample covariance matrix}}$$

within-cluster scatter matrix

$$\underline{\Sigma}_W = \sum_{i=1}^C \underline{\Sigma}_i$$

(B) Between clusters:

$$\text{total ~~cluster~~ center: } \underline{\mu} = \frac{1}{N} \sum_{n=1}^N x_n$$

$$\text{between-cluster scatter matrix } \underline{\Sigma}_B = \sum_{i=1}^C N_i (\underline{\mu}_i - \underline{\mu})(\underline{\mu}_i - \underline{\mu})^T$$

(C) For all samples:

$$\text{total scatter matrix } \underline{\Sigma}_T = \sum_{n=1}^N (x_n - \underline{\mu})(x_n - \underline{\mu})^T$$

Properties

- $\underline{\Sigma}_T$ fixed ; $\underline{\Sigma}_W, \underline{\Sigma}_B$ depend on S_i
- $\underline{\Sigma}_T = \underline{\Sigma}_W + \underline{\Sigma}_B$

5-6

Goal of clustering

$$\min \text{tr}(\underline{\Sigma}_W) \hat{=} \max \text{tr}(\underline{\Sigma}_B)$$

$$\text{trace} : \text{tr}(\underline{A}) = \sum_i a_{ii}$$

$$\text{tr}(\underline{A}\underline{B}) = \text{tr}(\underline{B}\underline{A})$$

Problem

$$\min_{S_1, \dots, S_{c-1}} J(S_1, \dots, S_{c-1}) = \text{tr}(\underline{\Sigma}_W)$$

$$= \sum_{i=1}^C \sum_{\underline{x}_n \in S_i} \left\| \underline{x}_n - \frac{1}{|S_i|} \sum_{\underline{x}_n \in S_i} \underline{x}_n \right\|^2$$

Discrete valued optimization problem

Q How many possibilities

5-7

5.1.2

K-means algorithm

5

- (+) much faster
- (-) don't find the best sol.ⁿ

Idea:

- If μ_i are well known: classify all x_n to their nearest mean (4.2.1) \Rightarrow cluster S_i

- If S_i known: recompute (improved) centers.

$$\mu_i = \frac{1}{|S_i|} \sum_{x_n \in S_i} x_n$$

5-8

5-9

5-10

\vdots

5-13

5.1.3 GMM

$\{\underline{x}(n)\}$, $M = \# \text{ modes} \hat{=} \# \text{ clusters} = C$

$\downarrow \text{EM}$

$\hat{\underline{\mu}}_m$, $\hat{\underline{C}}_m$, $q_{mn} = P(\underline{x}(n) \text{ belongs to mode/cluster } m)$
 $1 \leq m \leq M$

5.2 Clustering for an unknown number of clusters

5.1 : known number of clusters C ,

k -means, GMM \rightarrow sensitive to wrong C

In practice : C often unknown

\Rightarrow need other clustering algorithms