

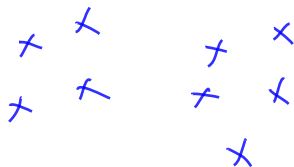
SMAI-M20-L34: Beyond Supervised Learning

C. V. Jawahar

IIIT Hyderabad

November 9, 2020

"Spectral clustering"



- 1 Consider a hierarchical classifier
 - Create an MST
 - Successively remove the longest or largest edges.

What objective does it optimize? maximizes? minimizes? Is the answer unique? Does the performance depend on initialization?

- 2 Is clustering unique? Where do we use clustering?

$$\begin{bmatrix} d(x_i, x_j) \end{bmatrix}_{N \times N}$$

N - vertices
all pairs

Recap:

- **Matrix Factorization and Applications:** (i) SVD, (ii) Eigen Decomposition (iii) Matrix Completion (iv) LSI (v) Recommendations
- **Dimensionality Reduction and Applications:** (i) Feature Selection and Extraction (ii) PCA (iii) LDA (iv) Eigen face
- **Supervised Learning:** Formulation, Conceptual Issues, Concerns etc. (i) Loss Functions and Optimization (ii) Probabilistic View, Bayesian View, MLE (iii) Eigen Vector based optimization (iv) Gradient Descent: Stochastic and Batch GD (v) Classification and Regression
- **Classifiers:** (i) Nearest Neighbour, (ii) Notion of a Linear Classifier (iii) Perceptrons (iv) Bayesian Optimal Classifier (v) Logistic Regression (vi) Multiclass classification architectures (v) Decision Trees (vi) SVMs (hard margin, soft margin, kernel) (vii) Kernel trick and kernelized algorithms
- **Neural Network Architectures and Learning** (i) Neuron model, Single Layer Perceptrons (ii) SLP (iii) MLP (iv) Backpropagation (v) Chain rule (vi) Activations (vii) challenges in optimization (viii) Momentum (ix) Convolutional Layer (x) Recurrent/Feedback networks (xi) Auto-encoder and unsupervised learning
- **Beyond Simple Supervised Learning** (i) Paradigms of Learning (ii)

This Lecture:

$\{L, U\}$

$\{Text\}$

① K-Means Clustering

- ① Initialize Clusters randomly
- ② Compute Means for each cluster
- ③ Re-Assign samples based on the distances to the means
- ④ Repeat until no-change

if dist
are all
small
a

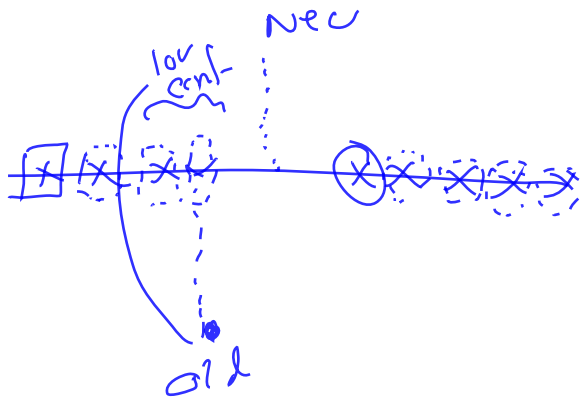
② Self-Training - *Semi supervised*

- ① Small number of Labelled and Large number of unlabelled examples
- ② Build a classifier with labelled examples.
- ③ Predict on unlabelled and use these "pseudo-labels" as "labels" and train the new classifier.
- ④ Take only selected (say most confident samples for re-training).

Questions? Comments?

$\|v^T x\|$

$\text{sign}(w^T x)$



"zero shot"
"few shot"

Discussions Point - I

Consider K-Means (with K=2)

- ① Six 1D samples

$$\{-3, -2, -1, +1, +2, +3\}$$

with $C_1 = \{-3, -2, +1\}$ and $C_2 = \{+3, +2, -1\}$

- ② Six 1D samples

$$\{-3, -2, -1, +1, +2, +3\}$$

with $C_1 = \{-3, +2, +1\}$ and $C_2 = \{+3, -2, -1\}$

- ③ Four 2D samples:

$$(-1, -1), (-1, +1), (+1, +1), (+1, -1)$$

with initialization as:

- $C_1 = \{(-1, -1), (-1, +1)\}$ and $C_2 = \{(+1, -1), (+1, +1)\}$
- $C_1 = \{(+1, +1), (-1, +1)\}$ and $C_2 = \{(+1, -1), (-1, -1)\}$

- ④ Comment on the optimization problem that K-Means solves (i) sensitivity to initialization (ii) convergence (iii) Local and global minima etc.

$$I = \sum_k \sum_{x_i \in C_k} \|x_i - \mu_k\|^2$$

Handwritten notes: "zero for simi." with an arrow pointing to the distance formula.

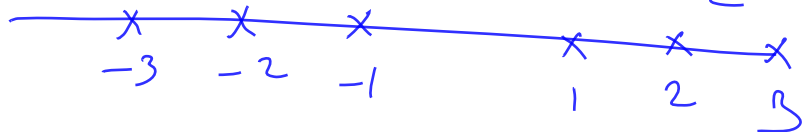
1C me + f

$$\mu_1 = -2$$

$$\mu_1 = 1.77$$

$$\mu_1 = +1$$

$$\mu_2 = -1.77$$

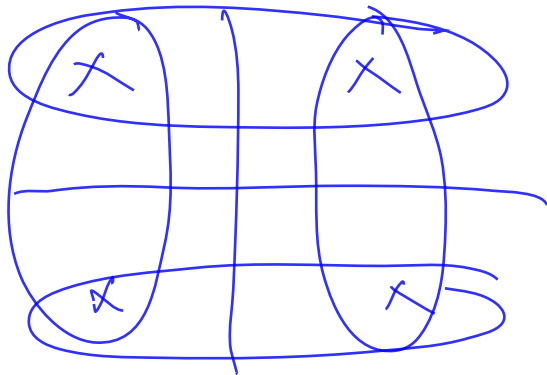


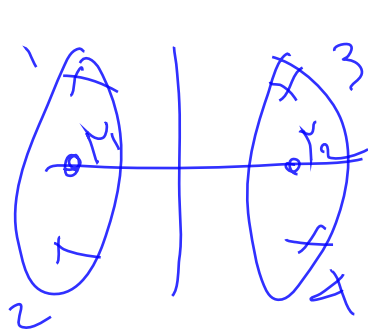
$$C_1 = -2, -1, +1 \rightarrow \mu_1 = -1.33$$

$$C_2 = 3, 2, -1 \rightarrow \mu_2 = +1.33$$

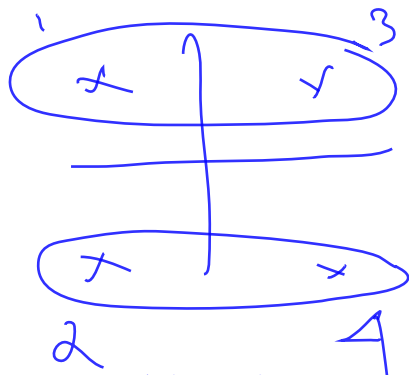
$$C_1 = -3, -2, -1 \rightarrow \mu_1 = -2$$

$$C_2 = 1, 2, 3 \rightarrow \mu_2 = +2$$



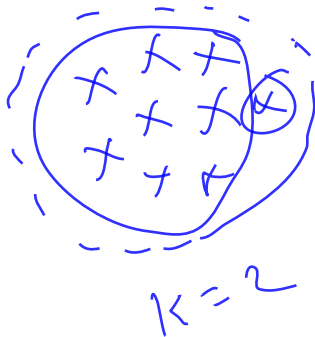


$K=2$



$K=2$

Is same?



Evaluation of Clustering is More Challenging than Classification. Objective is to sensitize the practical issues than exploring fully here.

How do we do performance evaluation in the Clustering problems?

- ① How do we evaluate/quantify the quality of a specific cluster?
- ② Does the objective (value) of the K-Means Obj. fn act as a good measure? How does it change with K ? What should be desirable?
- ③ How do we evaluate/quantify the computational complexity of clustering?
- ④ How do we evaluate the final clusters (results)
 - ① When there is no ground truth (GT)?
 - ② Compare two different outputs, and choose the better one?
 - ③ When K is same (in both) and different?
 - ④ When there is full GT (as if from a manually labelled data as in classification; may be only in some course assignments!!)
 - ⑤ When there is some partial GT (some samples were manually labelled)
 - ⑥ When there is some partial GT (some pairs were labelled as “similar” and some “dissimilar”. (or “must link”, and “can not link” in graph terminologies)

Suggest situations where some of these will be needed.

What Next:?

- ① Programming for Deep Learning.
- ② NN Architectures and NN Learning (winding up)
- ③ Beyond Simple Supervised Learning