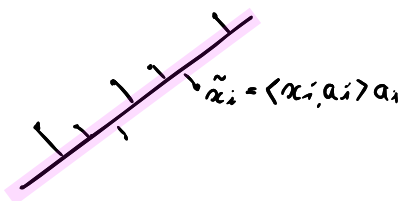


Unsupervised

I. Principal Component Analysis

- The 1st principal axis of projections.



axis is a direction maximizing the variance

- 1st principal axis is solution to:

$$\max_{a \in \mathbb{R}^p} a^T V a \quad \text{s.t. } \|a\| = 1$$

$$\mathcal{L}(a, \lambda) = a^T V a - \lambda (\|a\|^2 - 1)$$

$$\hookrightarrow a^T V a = a^T \lambda a = \lambda \|a\|^2 = \lambda \quad \text{largest eigen-value.}$$

Coordinates of x_i on projected principal axis a_1 are the 1st principal component.

$$c_{1,i} = \langle x_i, a_1 \rangle$$

- 2nd principal axis: $\max_a a^T V a \quad \text{s.t. } \|a\| = 1, a \perp a_1$

$$\hookrightarrow V a_2 = \lambda_2 a_2 \quad \text{In large dimension } p \gg n, \quad \frac{1}{n} K c = \lambda c, \quad K_{ij} = x_i^T x_j$$

...

\Rightarrow Reconstruction in the new base: $x_i = \sum_{k=1}^r c_{i,k} a_k$

$$\text{Simplified representation in dimension } s \ll r: \quad \tilde{X} = \tilde{C} \tilde{A}^T = (c_1 \dots c_s) \begin{pmatrix} a_1 \\ \vdots \\ a_s^T \end{pmatrix}$$

$$\text{Explained variance by } c_k = \frac{\lambda_k}{\sum_{j=1}^p \lambda_j}$$

Limitations: • scale variant \rightarrow Need standardization

• Complexity $O(p^3)$

• Might be impossible if p is too large.

- Kernel PCA

$$\Phi(x)^T \Phi(x') = k(x, x')$$

$$\text{Solve } C a = \lambda a \quad \text{where } C = \frac{1}{n} \sum \Phi(x_i) \Phi(x_i)^T$$

$$\text{Problem: } a_k = \frac{1}{n \lambda} (x^T C e_k) = \frac{1}{n \lambda} \sum_{i=1}^n c_{i,k} \Phi(x_i)$$

In the standard case, c_1 is the linear combination of variables that maximize the sum of squared correlations of variables:

$$c_1 = \operatorname{argmax}_{c \in \operatorname{Vect}(x^1 \dots x^p)} \sum_{j=1}^p R^2(x^j, c_j)$$

↳ Project on 2D correlation circle: $\cos \theta_{j,k} = \frac{\langle x^j, c^k \rangle}{\sqrt{x_j^T D x_j \cdot c_k^T D c_k^T}}$

angle between var. j and component k .

And length of line contribution of component c : $CTR(i, c) = \frac{p_i c_i^2}{\lambda}$

• Modified metric PCA: $V_a = \lambda_a$, $V = \langle x, x^T \rangle M \langle x, x^T \rangle$

• Modified weight PCA: $V = X^T D X - g g^T$, g barycenter, D matrix of weights.

II. Dependence Analysis

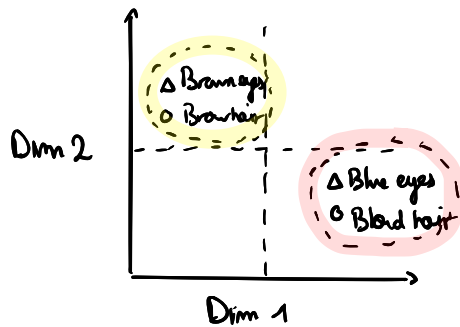
Used for discrete variables.

	Blue Eyes	Brown Eyes	Blond Hair	Brown Hair
Person 1	1		1	
Person 2	1			1
Person 3		1		1
...				

⇒

	Blond Hair	Brown Hair
Blue Eyes	16	8
Brown Eyes	4	29

Idea: We have 2 profiles (line, column) and we want to assess which categories of these profiles are correlated. It does however not say anything about intensity of relation.



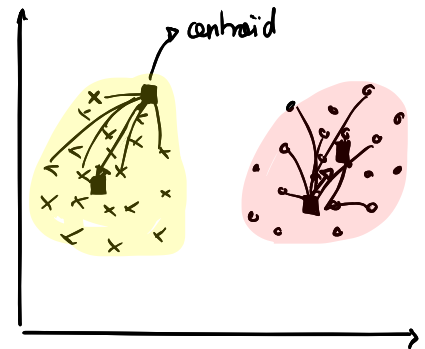
Correspondence analysis can be generalized.
↳ Multiple components > 2.

III. K-means

$$\min_{c_1, \dots, c_k} \sum_{j=1}^k \sum_{i \in G_j} \|x_i - \mu_j\|^2, \quad \mu_j = \frac{1}{|G_j|} \sum_{i \in G_j} x_i \quad \text{center of the cluster.}$$

Start from arbitrary location of clusters

- Find partition induced by nearest cluster
- Update cluster center
- Iterate until cluster convergence



Can be expressed as:

$$\min \sum_{k=1}^K \sum_{x_i, x_j \in C_k} \|x_i - x_j\|^2 \rightarrow \text{minimize distance of points belonging to a cluster}$$

$$\max \sum_{1 \leq k \neq l \leq K} \sum_{(x_i, x_j) \in C_k \times C_l} \|x_i - x_j\|^2 \rightarrow \text{maximize distance between points belonging to different clusters.}$$

• k-means ++

Select initial clusters far away from each other (squared distance criteria)

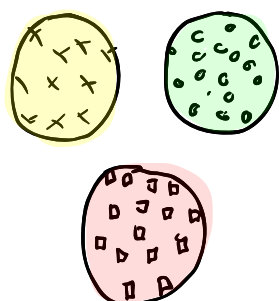
• Minibatch k-means

Update centroid using mini-batch of b samples chosen randomly
↳ Complexity $O(bk)$ instead of $O(nk)$

⇒ Always a finite number of configurations → these algorithms converge.

• Soft clustering

We suppose a gaussian distribution around a cluster. Each point has a probability $p_{ij} \propto \exp(-\|x_i - \mu_j\|^2 / 2\sigma^2)$ to belong to cluster j



Start from arbitrary location of cluster center

- Compute p_{ij}

- Update cluster center: $\mu_j = \frac{\sum_{i=1}^n p_{ij} x_i}{\sum_{i=1}^n p_{ij}}$

- Iterate until convergence

• Fuzzy k-means

Pareto (heavy tail) version of soft clustering: $p_{ij} \propto \frac{1}{\|x_i - \mu_j\|^d}$, $d > 0$

• Weighted k-means

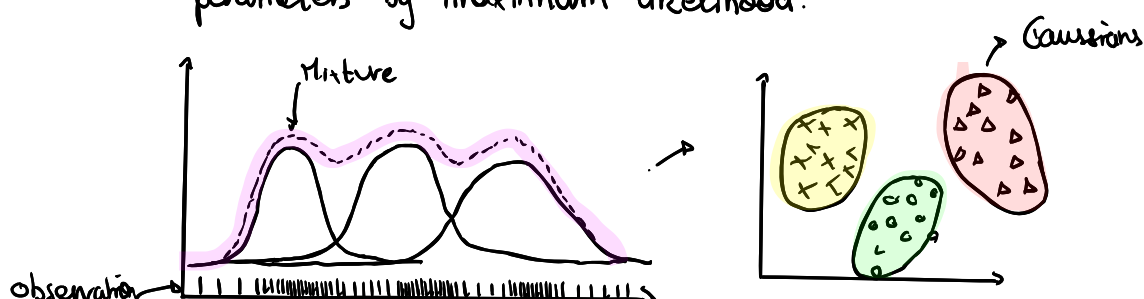
Assign relative importance to each sample w_i .

$$\min \sum_{j=1}^k \sum_{i \in G_j} w_i \|x_i - \mu_j\|^2 \quad \text{where} \quad \mu_j = \frac{\sum_{i \in G_j} w_i x_i}{\sum_{i \in G_j} w_i}$$

IV Gaussian Mixture model

Recall: $f(x) = \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$, $|\Sigma|$ determinant of Σ

Idea: Data may come from one of the gaussian distribution among the mixture. We need to set the number of clusters and identify best parameters by maximum likelihood.



$X \sim N(\mu_z, \Sigma_z)$, $Z \sim \pi$ latent variables

The parameters are the following: $\Theta = (\pi, \mu, \Sigma)$

• $\pi = (\pi_1 \dots \pi_k)$ mixing distribution (weights)

• $\mu = (\mu_1 \dots \mu_k)$

• $\Sigma = (\Sigma_1 \dots \Sigma_k)$

$$\hookrightarrow p_\Theta(x|Z) = \prod_{i=1}^n \pi_{z_i} f_{z_i}(x_i)$$

Maximum log likelihood: $\hat{\Theta} = \arg \max_{\Theta} l(\Theta; Z)$, $l(\Theta, Z) = \sum_{i=1}^n \log \pi_{z_i} + \sum_{i=1}^n \log f_{z_i}(x_i)$

\hookrightarrow Find solutions $\hat{\pi}_j$, $\hat{\mu}_j$, $\hat{\Sigma}_j$

\hookrightarrow Estimate latent variables: $p_\Theta(Z|x) = \frac{p_\Theta(x, Z)}{p_\Theta(x)}$ & $p_\Theta(x, Z) = \prod_{i=1}^n \pi_{z_i} f_{z_i}(x_i)$

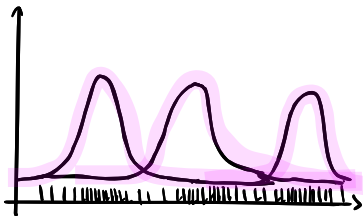
V. Expectation-maximization

Expected likelihood: $\sum_{j=1}^k \sum_{i=1}^n p_{ij} (\log \pi_j + \log f_j(x_i))$, $p_{ij} \propto \pi_j f_j(x_i)$ is the probability that observation i belongs to cluster j .

What is different with GMM? We now consider for each observation its probability to belong to each cluster, and EM allows a clear MLE solution.

$$\hookrightarrow \hat{\pi}_j = \frac{n_j}{n}, \quad \hat{\mu}_j = \frac{1}{n} \sum_{i=1}^n p_{ij} x_i, \quad \hat{\Sigma}_j = \frac{1}{n_j} \sum_{i=1}^n p_{ij} (x_i - \hat{\mu}_j)(x_i - \hat{\mu}_j)^T$$

\Rightarrow This is soft k-means. When $\sigma \rightarrow 0$, this is k-means.



Distance metric between 2 probab. distributions:

- KL divergence: $D(p||q) = \sum_{j=1}^k p_j \log \left(\frac{p_j}{q_j} \right)$

- Entropy: $H(p) = - \sum_{j=1}^k p_j \log p_j$

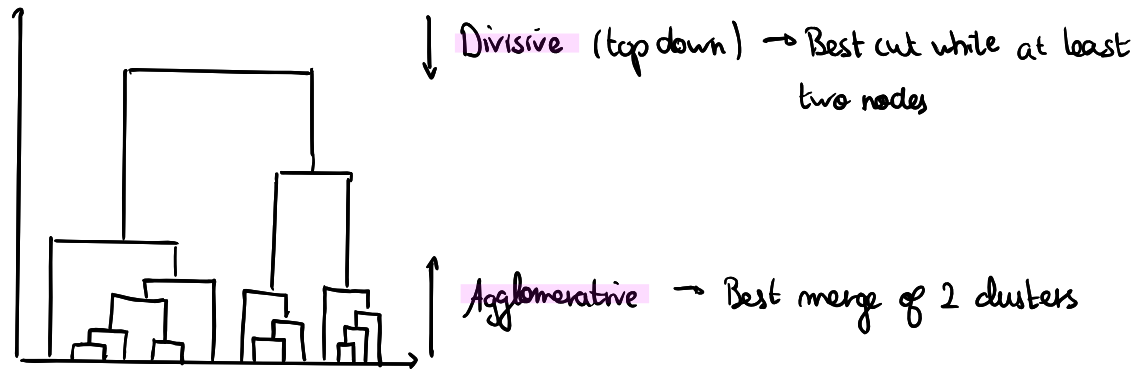
• Symmetric Gaussian mixture model.

\hookrightarrow EM with the set of covariance matrices $\Sigma = (\sigma^2 I, \dots, \sigma^2 I)$

VI. Hierarchical clustering

Aim: Reflect multiscale nature of data.

2 approaches to build a dendrogram



But how do we decide to split 2 far clusters or merge 2 close ones?

- **Minimum distance (single linkage):**

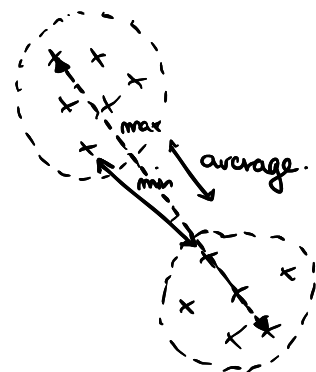
$$d(a, b) = \min_{i \in a, j \in b} \|x_i - x_j\|$$

- **Maximum distance (complete linkage):**

$$d(a, b) = \max_{i \in a, j \in b} \|x_i - x_j\|$$

- **Average (average linkage):**

$$d(a, b) = \frac{1}{|a||b|} \sum_{i \in a, j \in b} \|x_i - x_j\|$$



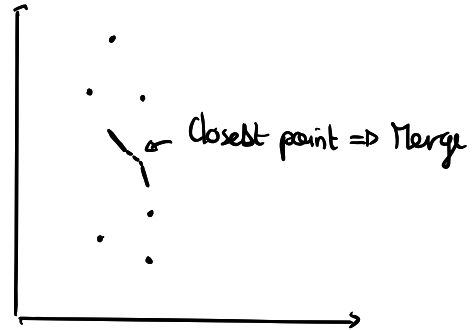
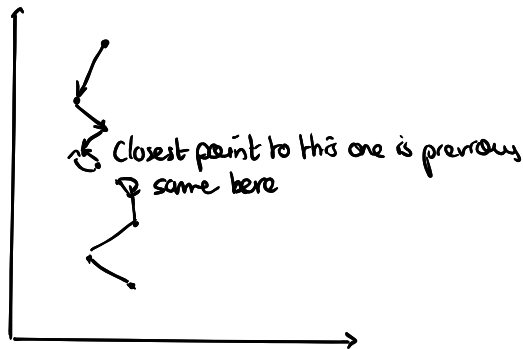
- **Ward's method:** Minimise sum of squared errors

$$S = \sum_{c \in C} \sum_{i \in c} \|x_i - g(c)\|^2 \quad \text{where } g(c) = \frac{1}{|c|} \sum_{i \in c} x_i$$

- **Nearest neighbor chain**

- 1) Start from any cluster
- 2) Build discretised chain of nearest neighbors until 2 clusters are jointly nearest neighbors

- 3) Merge these 2 clusters and proceed with rest of chain until empty
- 4) Go to step 1 if there are at least 2 clusters left.

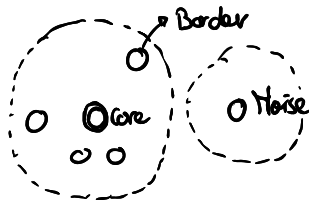


How do we measure performance for clustering?

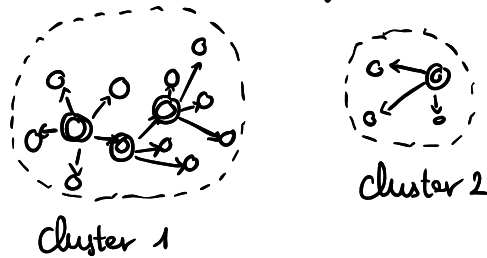
- square distance → Optimized for k-Means
- log-likelihood → Optimized for gaussian mixture
- mutual information (normalized)
- precision, recall, F1-Score

VIII. Density based spatial clustering of applications with noise

- A **core** is a point that has at least k points in its neighborhood.
- A **border** is a non-core point that has at least 1 core in its neighborhood.
- A **noise** is an outlier



- Pick a 1st core point and perform core-search
Assign all points recursively in neighborhood to the same cluster as X



- Advantage** :
- Can fit any shape and size of cluster
 - Can ignore outliers

Limitations : Sensitive to choice of neighborhood.

IX. Non-negative matrix factorisation model (NMF)

$$V_{F \times N} \approx W_{F \times K} \times H_{K \times N}$$

Data \approx Expl. Var. \times Regressors

We want to decompose V into positive matrices : $\begin{cases} W = [w_{fk}] \text{ s.t. } w_{fk} \geq 0 \\ H = [h_{kn}] \text{ s.t. } h_{kn} \geq 0 \end{cases}$

$$\Leftrightarrow V_n \approx \sum_{k=1}^K h_{kn} w_k$$

Used for :

- topic modeling (LSA)
- clustering (K-Means)
- temporal segmentation (videos)
- Source separation (ICA)

Limitations :

- Choice of K
- Non-unique solution

$V \approx WH$ usually obtained through $\min_{W, H \geq 0} D(V | WH)$ reconstruction error.

$$\text{where } D(V | \hat{V}) = \sum_{f=1}^F \sum_{n=1}^N d(v_{fn} | \hat{v}_{fn})$$

What divergence metrics can we use?

• **Euclidean distance**: $d_{\text{Euc}}(x, y) = (x - y)^2$

$d_{\text{Euc}}(\lambda x | \lambda y) = \lambda^2 d_{\text{Euc}}(x | y)$ not scale invariant.

• **Kullback-Leibler**: $d_{\text{KL}}(x, y) = x \log \frac{x}{y} - x + y$

$d_{\text{KL}}(\lambda x | \lambda y) = \lambda d_{\text{KL}}(x | y)$

• **Itahara-Saito (IS)**: $d_{\text{IS}}(x, y) = \frac{x}{y} - \log \frac{x}{y} - 1$

$d_{\text{IS}}(\lambda x | \lambda y) = d_{\text{IS}}(x | y)$ is scale invariant.

We solve optimal H and W values by multiplicative update (MU) and Majorization-minimization.

Limitations:

- monotonicity not guaranteed
- convergence not so good.

X. Independent Component Analysis (ICA)

$$X_{n \times K} \approx A_{F \times K} \cdot S_{n \times K}$$

data expl. var. regressors

Used for source separation:

Mixed signal = Mixing matrix \times Sources

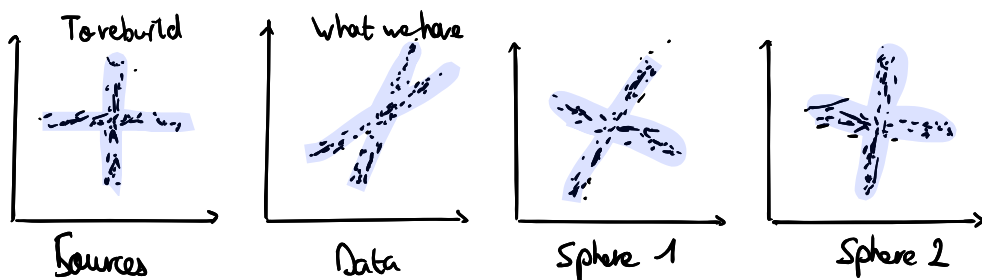
• 1) **Sphering**

Variance of entries of S must be normalized to 1

$$S_{\text{sph}} = A_{\text{sph}}^T X \text{ where } A_{\text{sph}} = E_{\text{IK}} D_K^{-\frac{1}{2}}$$

$$= Z$$

$$\hookrightarrow E(Z Z^T) = I \Rightarrow S_{\text{ISA}} = U_{\text{ISA}} \cdot Z \Rightarrow \text{Find } U \text{ rotation}$$



Once data is sphered:

- Construct numerical criteria $C(Y)$ measuring independence of entries
- Solve $\max_{\mu} C(\mu Z)$

↳ Signals that add up: CLT. We want to be as non-gaussian as possible

How to measure non-gaussianity?

- Kurtosis $\text{kurt}\{Y\} = \mathbb{E}(Y^4) - 3\mathbb{E}(Y^2)^2$
- Negentropy: $J(Y) = H\{Y_{\text{gaussian}}\} - H\{Y\}$

- Fast ICA algorithm

$$\max_{\mu} C(\mu) = |\text{kurt}\{\mu^T Z\}| \text{ s.t. } \mu^T \mu = 1$$

$$\Leftrightarrow \mathcal{L}(\mu, \lambda) = C(\mu) + \lambda(1 - \|\mu\|^2)$$

There must be at most 1 gaussian component inside, and we need to know how many sources there are.