

Unsupervised Learning



Unsupervised Learning

- Given feature vectors x_1, \dots, x_n
No labels

Clustering: group "similar" objects together

Main Ingredients

- 1) feature vectors $\underline{x}_j \in X$
- 2) # of clusters k
- 3) similarity or dissimilarity measure

$$\text{Similarity} \rightarrow s(x, x') \quad \text{Common conversion: } s(x, x') = \frac{-d(x, x')}{2\pi^2}$$

$$\text{dissimilarity} \rightarrow d(x, x') \quad \|x - x'\|^2 = \|x\|^2 + \|x'\|^2 - 2\langle x, x' \rangle$$

If featurevecs normalized to unit length

$$\|x_j\|_2 = 1 \quad \forall j, \quad s(x, x') := \langle x, x' \rangle \quad \text{then}$$

$$d(x, x') = \|x - x'\|^2 = 1 + 1 - 2 \langle x, x' \rangle \\ = 2 - 2 s(x, x')$$

$$d(x, x') = 2(1 - s(x, x'))$$

- Feature based Clustering

↳ Works with features x_1, \dots, x_n

- Similarity-based Clustering

↳ Works with the pairwise similarity matrix ($n \times n$)

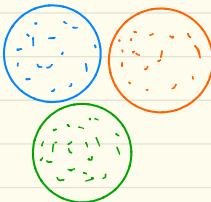
$$[s(x_i, x_j)]$$

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

Goal: Infer labels $y_1, \dots, y_n \in \{1, k\}$

Generative model for clustering

(not given colors)



$$P(x, y | \theta) = \underbrace{\prod_{j=1}^k P(y_j)}_{P(y)} \mathcal{N}(\mu_{y_j}, \sigma^2 I_d)(x)$$

$$P(x | y, \theta) = \mathcal{N}(\mu_y, \sigma^2 I_d)(x)$$

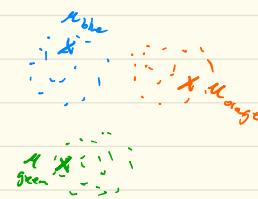
$$\theta = (\mu_1, \dots, \mu_k, \sigma^2)$$

ML estimates of labels y_1, \dots, y_n & parameters (μ_1, \dots, μ_k)

$$\left(\hat{y}_1, \dots, \hat{y}_n \right) = \arg \max_{\substack{y_1, \dots, y_n \in \{1, k\} \\ \mu_1, \dots, \mu_k \in \mathbb{R}^d}} \prod_{j=1}^n \underbrace{P(x_j | y_j, \theta)}_{\frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{\|x_j - \mu_{y_j}\|^2}{2\sigma^2}}}$$

$$= \arg \min_{\substack{y_1, \dots, y_n \\ \mu_1, \dots, \mu_k}} \sum_{j=1}^n \|x_j - \mu_{y_j}\|_2^2$$

Within-cluster sum of squares (WCSS)



K-means clustering (Lloyd's algorithm - 1957 Bell Labs)

Initialize : μ_1, \dots, μ_K

For $t = 1$ until stopping condition

Assign Labels : $\forall j, y_j = \underset{\ell \in \{1, \dots, K\}}{\operatorname{argmin}} \|x_j - \mu_\ell\|_2^2$
 \uparrow
(For all j)

Update Centers : $\forall \ell, \mu_\ell = \frac{1}{|\{j : y_j = \ell\}|} \sum_{j: y_j = \ell} x_j$

end for

Alternative way

For $t = 1$ until stopping condition

$\forall j, y_j^{(t)} = \underset{\ell \in \{1, \dots, K\}}{\operatorname{argmin}} \|x_j - \frac{1}{|\{j : y_j^{(t-1)} = \ell\}|} \sum_{j: y_j^{(t-1)} = \ell} x_j\|_2^2$

end for

$$\|x_j - \frac{1}{n_\ell} \sum_{j: y_j^{(t-1)} = \ell} x_j\|^2 = \langle x_j - \frac{1}{n_\ell} \sum x, x_j - \frac{1}{n_\ell} \sum x \rangle$$

$$= \langle x_j, x_j \rangle - \frac{2}{n_\ell} \sum_j \langle x_j, x_j \rangle + \sum \sum \langle x_j, x_j \rangle$$

Replacing $\langle x, x \rangle$ with $k(x, x)$ gives us the
Kernel k-means Algorithm

Initialization

- 1) Randomly (uniformly, without replacement) choose K points among training set
- 2) Randomly assign labels (uniformly over $1, \dots, K$) to each training point.
- 3) K-means ++ (farthest point clustering)
 - 1st point picked uniformly at random
 - sequentially pick from remaining points with the probability which is proportional to squared distance to closest cluster point among cluster points chosen thus far.

K-means ++ is guaranteed to obtain a within-cluster sum of squares which is not more than $O(\log k)$ of minimum possible WCSS.

Mixture of Gaussians (GM)

$$p(x, y | \theta) = \underbrace{\lambda_y}_{p(y|\theta)} \underbrace{N(\mu_y, \Sigma_y)}_{p(x|y, \theta)}(x)$$

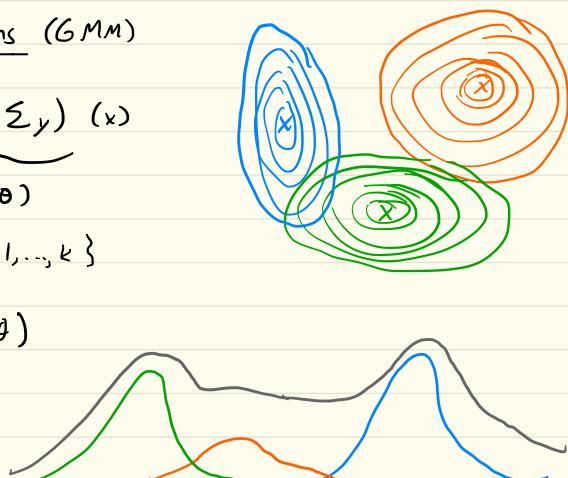
$$\Theta = \{(\lambda_y, \mu_y, \Sigma_y), y=1, \dots, k\}$$

$$\Rightarrow p(x | \theta) = \sum_{y=1}^k p(x, y | \theta)$$

$$= \lambda_1 N(\mu_1, \Sigma_1)(x)$$

$$+ \lambda_2 N(\mu_2, \Sigma_2)(x)$$

$$+ \dots + \lambda_k N(\mu_k, \Sigma_k)(x)$$

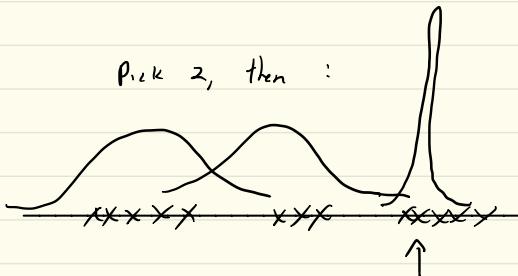


Idea of soft membership to clusters

$$\{ P(Y_j = l \mid x_j, \theta), \quad l=1, \dots, k \} \quad \text{Prob. of membership to cluster } l$$

Given points, — ~~XXXXXX~~ ~~XXX~~ ~~XXXX~~

Pick 2, then :



Finding ML estimates of θ is ill-posed:

Can make likelihood go to ∞ by choosing a S function for one of the components centered at one of the data points.

\Rightarrow need a prior for θ , specifically for $\Sigma_1, \dots, \Sigma_k$

E - M algorithm (Expectation Maximization)

$$\underline{\text{Idea}} \quad (x, y) \sim p(x, y | \theta) \quad , \quad \theta \sim \pi(\theta) \rightarrow \text{prior}$$

Complete data: (x, y)

$$\text{Ideally: } \hat{\theta} = \underset{\theta}{\operatorname{argmax}} \left[\ln p(x, y | \theta) + \ln \pi(\theta) \right]$$

Incomplete data: X

but, Y is not given, only X is
(latent, hidden)

$$\int \ln(\rho(\mathcal{X}, \mathcal{Y} | \theta) \pi(\theta)) \rho(\mathcal{Y} | \mathcal{X}, \pi) d\theta$$

Solution :

- "average-out" \mathcal{T} in cost, then maximize over θ
- But "arg-out" requires $P(\mathcal{T} | \mathcal{X}, \theta)$
- Use previous θ to do avg'ing (expectation step) then
maximize over θ (M-step)

Expectation step (E-step)

$$\text{Replace } \ln p(x, y | \theta) + \ln \pi(\theta) \text{ with } E[\ln p(x, y | \theta) + \ln \pi(\theta) \mid X = x]_{\theta = \theta}$$

M maximization step (M-step)

$$\theta^{(t)} = \underset{\theta}{\operatorname{argmax}} \quad Q(\theta, \theta^{(t-1)})$$

$$Q(\theta, \theta^{(t-1)})$$

Geometric

$$\underbrace{p(x | \theta^{(t)}) \pi(\theta^{(t)})}_{\text{unconditioned likelihood}} \leq p(x | \theta^+) \pi(\theta^+)$$

unconditioned likelihood

E-M algorithm

Initialize $\hat{\theta}_0$

For $t = 1, \dots, t_{\max}$ or stopping cond.

$$\underline{\text{E-step}}: Q(\theta, \hat{\theta}_{t-1}) = E_{\bar{Y}} \left[\ln(p(\bar{X}, \bar{Y} | \theta) \pi(\theta)) \mid \bar{X}, \hat{\theta}_{t-1} \right]$$

$$\underline{\text{M-step}}: \hat{\theta}_t = \underset{\theta}{\operatorname{arg\,max}} \quad Q(\theta, \hat{\theta}_{t-1})$$

end for

$$\cdot E_{\bar{Y}} \left[\ln(p(\bar{X}, \bar{Y} | \theta) \pi(\theta)) \mid \bar{X}, \hat{\theta}_{t-1} \right] = \int \ln(p(\bar{X}, \bar{Y} | \theta) \pi(\theta)) \cdot p(\bar{Y} | \bar{X}, \hat{\theta}_{t-1}) d\bar{Y}$$

• key property

$$p(\bar{X} | \hat{\theta}_{t+1}) \pi(\hat{\theta}_{t+1}) \geq p(\bar{X} | \hat{\theta}_t) \pi(\hat{\theta}_t)$$

6MM

$$p(x, y | \theta) = \underbrace{\lambda_y N(\mu_x, \sigma^2_y)}_{p(y|v|\theta)} \underbrace{\rho(x|v, \theta)}_{p(x|v, \theta)}$$

$$\Theta = \{(\lambda_y, \mu_x, \sigma^2_y), y=1, \dots, k\}$$

$$\bar{X} = (\underline{x}_1, \dots, \underline{x}_n)$$

$$\bar{Y} = (y_1, \dots, y_n) \rightarrow \text{cluster labels (hidden)}$$

$$P(\mathbf{X}, \mathbf{Y} | \theta) = \prod_{j=1}^n p(x_j, y_j | \theta)$$

($\mathcal{M}(\theta) \sim \text{Unif.}$)

$$(\text{no prior}) = \prod_{j=1}^n \lambda_{y_j} N(\mu_{y_j}, \sigma_{y_j}) (x_j)$$

$$= \prod_{j=1}^n \prod_{\ell=1}^k (\lambda_\ell N(\mu_\ell, \sigma_\ell)(x_j) I(y_j = \ell))$$

$$\ln P(\mathbf{X}, \mathbf{Y} | \theta) = \sum_{j=1}^n \sum_{\ell=1}^k I(y_j = \ell) \ln (\lambda_\ell N(\mu_\ell, \sigma_\ell)(x_j))$$

$$E\text{-step} \quad E_{\bar{\theta}} [\ln P(\mathbf{X}, \mathbf{Y} | \theta) | \mathbf{X}, \bar{\theta}_{+1}]$$

$$= \sum_{j=1}^n \sum_{\ell=1}^k E[I(y_j = \ell) | X_j = x_j, \bar{\theta}_{+1}] \cdot \ln (\lambda_\ell N(\mu_\ell, \sigma_\ell)(x_j))$$

$$E[I(y_j = \ell) | X_j = x_j, \bar{\theta}_{+1}]$$

$$=: \gamma_{x_j}^{(+1)}$$

$$\sum_{j=1}^n \sum_{\ell=1}^k \gamma_{x_j}^{(+1)} \ln (\lambda_\ell N(\mu_\ell, \sigma_\ell)(x_j))$$

M-step maximize over

$$\Theta = \{(\lambda_\ell, \mu_\ell, \sigma_\ell)\} \quad \text{Similar to GDA with weights } \gamma_{x_j}^{(+1)}$$

$$\hat{\lambda}_\ell^{(t)} = \frac{1}{n} \sum_{j=1}^n \delta_{\ell_j}^{(t-1)}$$

$$\hat{\mu}_\ell^{(t)} = \frac{\sum_{j=1}^n \delta_{\ell_j}^{(t-1)} x_j}{\sum_{j=1}^n \delta_{\ell_j}^{(t-1)}}$$

$$\hat{\Sigma}_\ell^{(t)} = \frac{\sum_{j=1}^n \delta_{\ell_j}^{(t-1)} (x_j - \hat{\mu}_\ell^{(t)}) (x_j - \hat{\mu}_\ell^{(t)})^T}{\sum_{j=1}^n \delta_{\ell_j}^{(t-1)}}$$

$$\begin{aligned} \delta_{\ell_j}^{(t)} &= P(Y_j = \ell \mid X_j = x_j, \theta_t) \\ &= \frac{\hat{\lambda}_\ell^{(t)}}{\sum_{\ell=1}^k \hat{\lambda}_\ell^{(t)}} N(\hat{\mu}_\ell^{(t)}, \hat{\Sigma}_\ell^{(t)})(x_j) \end{aligned}$$

* Spectral Clustering Tutorial on class webpage
 (+ pseudocode)

Spectral Clustering

- Graph-cut based clustering algorithms

Similarity score : high \Rightarrow similar

Can use to form a weighted or unweighted Graph

Gaussian kernel:

$$s(x_i, x_j) = \exp(-\|x_i - x_j\|^2 / 2\sigma^2)$$

another is: $s(x_i, x_j) = \begin{cases} 1 & \text{if } \|x_i - x_j\| \leq \xi \\ 0 & \text{else} \end{cases}$

another! $s(x_i, x_j) = \begin{cases} 1 & \text{if } x_i \in NN \text{ of } x_j \text{ OR } x_j \in NN \text{ of } x_i \\ 0 & \text{else} \end{cases}$

another: " " " " AND "

How to cut the graph

Graph terminology

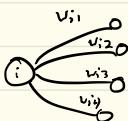
- similarity graph on N nodes [node \equiv feature/training vector]
- Weighted adjacency matrix

$$W_{N \times N}, w_{i,j} = s(x_i, x_j)$$

$$i \left[\begin{array}{c} j \\ \vdots \\ \hline - - - \\ i_{ij} = s(x_i, x_j) \end{array} \right]_n^n$$

Assume symmetric $s(\cdot, \cdot)$
 $s(x, x') = s(x', x)$
 $\Rightarrow W = W^T$

Node degree



degree of node 'i': $d_i := \sum_{j=1}^n w_{i,j}$

Degree Matrix

$$D_{N \times N} := \text{diag}(d_1, \dots, d_n) = \begin{bmatrix} d_1 & & & 0 \\ & \ddots & & \\ 0 & & \ddots & d_n \end{bmatrix}$$

Graph Laplacian $L_{n \times n}$

$$L := D - W$$

Properties of L

1) If $\underline{f} \in \mathbb{R}^n$, $\underline{f} = \begin{bmatrix} f_1 \\ f_2 \\ \vdots \\ f_n \end{bmatrix}_{n \times 1}$

$$\underbrace{\underline{f}^\top L \underline{f}}_{\|f\|} = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n w_{ij} (f_i - f_j)^2$$

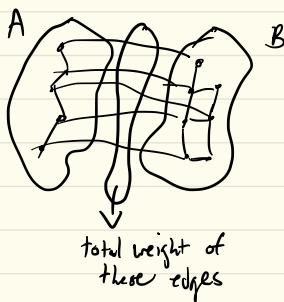
2) $L = L^\top$, 3) (1) $\Rightarrow L \geq 0$ positive semidefinite
symmetric

$$\text{Cut}(A, B)$$

$$\text{CUT}(A, B) := \sum_{i \in A} \sum_{j \in B} w_{ij}$$

$$A \cap B = \emptyset$$

$$A, B = \{1, \dots, n\}$$



K clusters A_1, A_2, \dots, A_k = position of nodes for k adjacent graphs

if $i \neq j$, $A_i \cap A_j = \emptyset$,

$$\bigcup_{i=1}^k A_i = \{1, \dots, k\}$$

$$A^c = \{j : j \in A\} = \{1, \dots, n\} \setminus A$$

$$\text{Ratio Cut : } \sum_{l=1}^k \frac{\text{Cut}(A_l, A^c)}{|A_l|}$$

$$|A_l| = |\{j : j \in A_l\}|$$

= # nodes in cluster l

*

$$F^{\text{Ratio Cut}} = [f_1, \dots, f_k]_{n \times k}$$

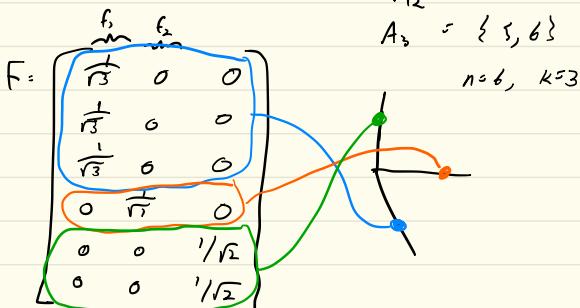
$$f_{il} = \begin{cases} \frac{1}{|A_l|} & \text{if } j \in A_l \\ 0 & \text{else} \end{cases}$$

$$F = \begin{bmatrix} \frac{1}{\sqrt{|A_1|}} & \dots & \frac{1}{\sqrt{|A_k|}} \end{bmatrix}$$

e.g. suppose $A_1 = \{1, 2, 3\}$

$$A_2 = \{4\}$$

$$A_3 = \{5, 6\}$$



$$\underline{\text{Results}} : \quad f^T f = I_k$$

$k \times n \quad n \times k$

$$\text{Ratio Cut } (A_1, A_2, \dots, A_k) = \sum_{\ell=1}^k \underbrace{f_\ell^T L f_\ell}_{1 \times 2}$$

$$\begin{aligned} \text{Recall} \quad abc &= \text{Trace}(abc) \\ (\text{scalar}) &= \text{Tr}(bc) \\ &= \text{Tr}(cab) \end{aligned}$$

$$f_\ell^T L f_\ell = \text{Tr}(f_\ell^T L f_\ell) = \text{Tr}(L f_\ell f_\ell^T)$$

$$\Rightarrow \sum_\ell f_\ell^T L f_\ell = \sum_\ell \text{Tr}(L f_\ell f_\ell^T) = \text{Tr}\left(L \sum_\ell f_\ell f_\ell^T\right)$$

$$= \text{Tr}(L F F^T) = \text{Tr}(F^T L F)$$

Ratio Cut minimization for clustering

$$= \underset{A_1, \dots, A_k}{\text{argmin}} \quad \text{Ratio Cut}(A_1, \dots, A_k)$$

$$= \underset{F \text{ of form } *} {\text{argmin}} \quad \text{Tr}(F^T L F)$$

$$\underset{\underline{F^T F = I}}{\text{argmin}} \quad \text{Tr}(F^T L F)$$

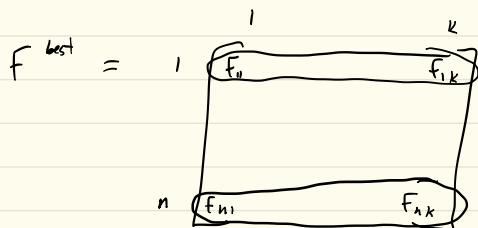
↑ relax requirement of F

Solution to Relaxed Ratio Cut Problem

→ Rayleigh-Ritz Theorem

cols of best F

\equiv orthonormal eigenvectors of the k SMALLEST eigenvalues of L



→ cluster the rows of F^{best} into k groups using k -means

Intuition

D = degree matrix

SC-2

$V = F$

RW: random walk on graph pagerank

NCut

$$\text{NCut}(A_1, \dots, A_k) = \sum_{\ell=1}^k \frac{\text{Cut}(A_\ell, A_\ell^c)}{\text{Vol}(A_\ell)}$$

$$\text{Vol}(A_\ell) = \sum_{j \in A_\ell} \delta_j$$

$$WCSS(y_1, \dots, y_n) = \sum_{j=1}^n \|x_i - M_{y_j}\|_2^2$$

- K-means algorithm
- K-means ++
- GMMs \rightarrow EM algorithm
- Spectral methods:

$W_{n \times n} \rightarrow$ weighted adjacency matrix
 $w_{ij} \equiv$ similarity between x_i, x_j

$$\text{Graph Laplacian : } L = D - W$$

$$D = \text{diag}(W \mathbf{1}_n)$$

$A_1, A_2, \dots, A_k \rightarrow$ k clusters (set of nodes)

Ratio Cut (A_1, \dots, A_k)

$$= \sum_{l=1}^k \frac{\text{cut}(A_l, A_l^c)}{|A_l|}$$

$$* F = \begin{bmatrix} \frac{1}{\sqrt{|A_1|}}, & \dots, & \frac{1}{\sqrt{|A_k|}} \end{bmatrix} \quad F^T F = I_k$$

$$\text{RatioCut}(A_1, \dots, A_k) = \text{Tr}(F^T L F)$$

RatioCut minimization

$$\equiv \underset{A_1, A_k}{\operatorname{argmin}} \text{RatioCut}(A_1, \dots, A_k)$$

$$\equiv \underset{F \text{ orthonormal}}{\operatorname{argmin}} \text{Tr}(F^T L F)$$

Relaxed Ratio Cut minimization

$$\underset{\substack{\text{argmin} \\ F: F^T F = I_k}}{\text{Tr}(F^T L F)}$$

Solution: Rayleigh - Ritz Theorem

cols of F_{best} = k o.n. eigenvalues of L correspond to the
 (relaxed) k smallest eigenvalues of L

$$L = U \Lambda U^T = [u_1 \ u_n] \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_n \end{bmatrix} \begin{bmatrix} u_1^T \\ u_n^T \end{bmatrix}$$

$$F_{\text{best}} = \begin{bmatrix} | & | \\ u_1 & \dots & u_k \\ | & | \end{bmatrix}$$

WCSS ($y_1, \dots, y_n, m_1, \dots, m_k$)

$$= \underset{\substack{\text{argmin} \\ F \text{ of form } *}}{\text{Tr}(F^T (I - k) F)}$$

$$\downarrow \quad \downarrow$$

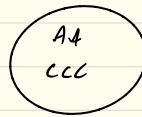
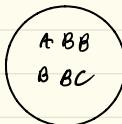
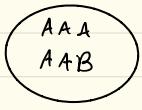
$\overset{n \times n}{\text{id matrix}}$

$$k = \mathbb{X}^T \mathbb{X}, \quad \mathbb{X} = \begin{bmatrix} | & | \\ x_1 & \dots & x_n \\ | & | \end{bmatrix}$$

Metrics

(for clustering when given a reference clustering)

1) Purity



$$\begin{array}{r} \text{dominant} \\ \text{label counts} \end{array} \quad \frac{5 + 4 + 3}{6 + 6 + 5}$$

Formally, $n_{ij} = \# \text{ class } j \text{ points in cluster } \# i$

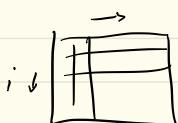
$$n_i = \max_j n_{ij}, \quad \text{purity} = \frac{\sum_{i=1}^k n_i}{n}$$

2) Normalized mutual information (NMI)

$$P(i, j) = \frac{n_{ij}}{n}$$

$$P_u(i) = \sum_j P_{uv}(i, j)$$

$$P_v(j) = \sum_i P_{uv}(i, j)$$



Mutual Information

$$I(u, v) := D(P_{uv} || P_u \cdot P_v)$$

$$NAT = \frac{I(u; v)}{(H(u) + H(v))/2}$$

H: entropy

If K is not known

- 1) penalize # clusters: $\underset{\mu_1, \dots, \mu_K}{\operatorname{argmin}} WCSS(\underset{x_1, \dots, x_n}{y_1, \dots, y_n}, \mu_1, \dots, \mu_K) + \lambda k^p$
regularization parameter
- 2) DP-means ($DP = \text{Dirichlet Process means}$)
 - like k-means
 - Idea: if there exists a point whose distance to closest μ_1, \dots, μ_K is greater than some distance γ , then $K \rightarrow K+1$ & begin new cluster.
- 3) Sum of Norms (SON) clustering

$$\underset{\mu_1, \dots, \mu_n}{\operatorname{argmin}} \left[\sum_{j=1}^n \|x_j - \mu_j\|_2^2 + \lambda \sum_{j=2}^n \sum_{i < j} \|\mu_i - \mu_j\|_p \right]$$

- 4) Chinese Restaurant Process

