

# 10/3 Clustering scribe (cont.)

recall

k-means

Input:  $\{x_i\}_{i=1}^n$

$k$

Search for  $C = \{c_1, \dots, c_k\}$

$\mu = \{\mu_1, \dots, \mu_k\}$

$$\min_{C, \mu} \sum_{l=1}^k \sum_{i \in c_l} \|x_i - \mu_l\|_2^2 \quad (\text{kmeans})$$

$$\min_{C, \mu} \sum \sum \|x_i - \mu_l\|_1 \quad (\text{k medians})$$

$$\min_{C, \mu} \sum \sum \|x_i - \mu_l\|_2 \quad (L_2-L_1 \text{ norm})$$

Initial seeding

① random

② by singular value decomposition

Consistency - can random seeding lead to of kmeans - the true  $\mu$  centroids??

- If  $C^*, \mu^*$  are the true minimizers.

Strong Law of Large of Numbers (1960-80)

$$L_n(\mu) = \int \|x - \mu\|^2 dP_n(x)$$

$$dP_n(x) = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}(x) dx \quad \# \text{empirical measure of data}$$

def  $A = \{a_1, \dots, a_n\}, x \in \mathbb{R}^d$

$$\|x - A\| = \min_{1 \leq i \leq n} \|x - a_i\|$$

As  $L_n(\mu) \rightarrow \infty, L(\mu)$  and

$$f(\mu) = \int \|x - \mu\|^2 dP(x), \text{ but does the minimizer converge?}$$

Consistency of Lloyd's + New analysis

(KK '01, McSherry '01)

1A

Requirement of "separation"

→ bound for misclassification error.

assuming same true centroidal partition

$\|\mu_k^* - \mu_l^*\|$  if  $k \neq l$  needs to be bigger

than a gap  $\Delta_{kl}$

$\|\mu_k^* - \mu_l^*\| > \Delta_{kl}$  where  $\Delta_{kl} = \frac{\text{number of points in set } C_l^*}{n}$

$$\Delta_{kl} > c \times \frac{\sigma_k}{\sqrt{w_{\min}}} \quad \text{and } w_{\min} = \min_l \frac{|C_l^*|}{n}$$

$\sigma_k^2$  = variance of data in  $k^{\text{th}}$  cluster.

and if we can assume  $\sigma_k \approx \sigma_l$

(The) Lloyd's and SVD misclassification  $\leq \epsilon \cdot n$  points w/ high probability

Practical considerations of k-means

- what is  $k$ ?
- how to initialize
- failing cases of k-means

diameter  $\rightarrow \sigma_k$  might be too large

Compared to separation  $\|\mu_k - \mu_l\|$

cluster size  $\rightarrow$  too small cluster. If cluster sizes are not evenly balanced

Voronoi cells  $\rightarrow$  cluster must be convex and piecewise. can't do concave

not the same as tolerance.

$$\min_{\mu} L_n(\mu) \rightarrow \min_{\mu} f(\mu)$$

but this doesn't give us the misclassification rate

## Spectral Clustering

Given  $\{x_i\}_{i=1}^n$  and  $k$ ,

① Build affinity matrix  $W_{n \times n}$  by  $knn$   $\epsilon$ -neighborhood or  $W_{ij} = e^{-\frac{\|x_i - x_j\|^2}{\epsilon}}$

$W$  has to be positive and symmetric

② Consider the eig decomp of  $L$  graph Laplacian

$$L_{un} = D - W \quad (\text{unnormalized})$$

$$L_{rw} = D^{-1}(D - W) = I - P \quad (\text{Sym-Markov '00})$$

$L_{sym}$  = ~~compute the first  $k$~~  (NJW '02)  
eigenval and eigvec of the

③ Apply  $k$ -means to  $\Psi = [\psi_1, \dots, \psi_k]_{n \times k}$

$\psi_i$  is the  $i$ -th row of  $\Psi$

Remark If  $k=2$ , you can use truncation of truncation of  $\text{sign}(\psi_2)$ , b.c. first eigenval is constant

Can use  $\text{sign}(\psi_2)$  to indicate clustering

## Definition (Connected components)

On  $G = (V, E)$  on each edge  $(i, j) \in E$ ,

$W_{ij} > 0$ . If node  $i$  is connected to  $j$ , there is a path from  $i$  to  $j$ . So then

Set  $A$  is a connected component, if every pair of  $i$  and  $j$  is connected and

$A$  is the maximum set that satisfies this condition to preserve connectivity

Prop (Eigenspace of  $\lambda=0$  of  $L$  and connected components in  $G$ )

Suppose the graph has  $k$  conn. comp  $A_1, \dots, A_k$ . Then the eigenspace of  $\lambda=0$  of dim  $k$  is spanned by  $\{\mathbb{1}_{A_1}, \dots, \mathbb{1}_{A_k}\}$  where  $\mathbb{1}_{A_i}(i) = \begin{cases} 1 & \text{if } i \in A_i \\ 0 & \text{otherwise} \end{cases}$

Proof

Suppose  $f$  is eigenvector of  $\lambda=0$ . So

$$L f = 0 \quad \text{and} \quad f^T L f = 0. \quad \text{We also}$$

$$\text{know } f^T L f = \frac{1}{2} \sum_{(i,j) \in E} W_{ij} (f_i - f_j)^2$$

$$\text{and } f^T L f = 0 \Leftrightarrow f_i = f_j \text{ whenever}$$

$$W_{ij} > 0$$

Thus,  $f$  is piecewise constant in each of conn. comp. Meanwhile,

for each  $v \in \text{span}\{\mathbb{1}_{A_1}, \dots, \mathbb{1}_{A_k}\}$

$$v^T L v = 0$$

Exercise:

- Does this generalize to  $L_{rw}$  and  $L_{sym}$ ?
- How leads to consistency of spectral clustering?