

## Lecture 9

①

Review \* Scheffé selector for learning 1-dimensional mixtures of Gaussian:

→  $\Omega \equiv \epsilon/2k$  grid of the prob. simplex  $\mathcal{P}_k$

$\mathcal{L} = \{P_{\omega, \mu} : \omega \in \Omega, \mu_i \in \{X_1, \dots, X_n\} \text{ for all } 1 \leq i \leq k\}$

Use Scheffé selector to find the best match

\* ML selector can also work if our guess-list satisfies more conditions, but fails in general.

Agenda \* Learning d-dimensional Gaussian mixtures

[A] Distance based clustering for learning Gaussian mixtures in higher dimensions

Consider  $X_1, \dots, X_n$  generated iid from  $\sum_{j=1}^k \omega_j \mathcal{N}(\mu_j, \sigma^2 I_{d \times d})$ .

How does the distance b/w two samples generated from the same mean differ from that b/w two samples from a different mean?

If  $X, Y \sim \mathcal{N}(\mu, \sigma^2 I_{d \times d})$ ,

$$E \|X - Y\|_2^2 = \sum_{i=1}^d E[(X_i - Y_i)^2] = 2d\sigma^2$$

If  $X \sim N(\mu_1, \sigma^2 I_{d \times d})$  and  $Y \sim N(\mu_2, \sigma^2 I_{d \times d})$ , ②

$$\begin{aligned} \mathbb{E} \|X - Y\|_2^2 &= \sum_{i=1}^d \mathbb{E} [X_i^2 + Y_i^2 - 2X_i Y_i] \\ &= \sum_{i=1}^d (\sigma^2 + \mu_{1i}^2 + \sigma^2 + \mu_{2i}^2 - 2\mu_{1i}\mu_{2i}) \\ &= 2d\sigma^2 + \|\mu_1 - \mu_2\|^2 \end{aligned}$$

Also, in the first case, using Chebyshev's inequality

$$\begin{aligned} P\left( \left| \|X - Y\|_2^2 - \mathbb{E}[\|X - Y\|_2^2] \right| > t \right) &\leq \frac{\text{Var}(\|X - Y\|_2^2)}{t^2} \\ &= \frac{d}{t^2} \text{Var}((X_i - Y_i)^2) \\ &= O\left(\frac{d}{t^2} \cdot \sigma^4\right). \end{aligned}$$

Thus, in this case, with significant probability

$$\|X - Y\|_2^2 = 2d\sigma^2 \pm O(\sqrt{d}\sigma^2)$$

On the other hand, in the second case

$$\begin{aligned} &P\left( \left| \|X - Y\|_2^2 - \mathbb{E}[\|X - Y\|_2^2] \right| > t \right) \\ &\leq \frac{1}{t^2} \sum_{i=1}^d \text{Var}((X_i - Y_i)^2) \\ &= \frac{1}{t^2} \sum_{i=1}^d \left[ (\mu_{1i} - \mu_{2i})^4 + 12(\mu_{1i} - \mu_{2i})^2 \sigma^2 + 12\sigma^4 \right] - (2\sigma^2 + (\mu_{1i} - \mu_{2i})^2)^2 \\ &= \frac{1}{t^2} \sum_{i=1}^d 8(\mu_{1i} - \mu_{2i})^2 \sigma^2 + 8\sigma^4 = \frac{8\sigma^2}{t^2} (\|\mu_1 - \mu_2\|^2 + d\sigma^2) \end{aligned}$$

Thus, in the second case, with large prob.

$$\|x - \gamma\|_2^2 = 2d\sigma^2 + \|\mu_1 - \mu_2\|_2^2 \pm O(\sigma(\|\mu_1 - \mu_2\|_2 + \sqrt{d}\sigma)) \quad (3)$$

Therefore, we can reliably distinguish two points from the same cluster from those coming from two different clusters iff

$$\|\mu_1 - \mu_2\|_2^2 - O(\sigma \|\mu_1 - \mu_2\|_2) \gg O(\sqrt{d}\sigma^2)$$

which holds if

$$\|\mu_1 - \mu_2\|_2 > cd^{1/4}\sigma.$$

\* The effort in this line of research is to seek algorithms that circumvent the requirement of  $\|\mu_1 - \mu_2\|_2$  to grow  $d$ .

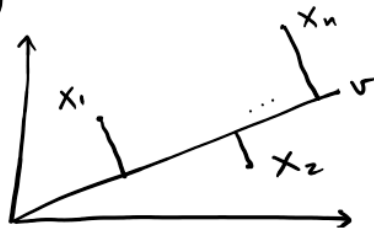
### [B] Projection to a lower dimensional space

\* Idea (Vempala-Wang '02) Do distance based clustering after projecting on a lower dimensional space.

### Review of Singular Value Decomposition (SVD)

Chapter 3 of "Foundations of data science," Blum, Hopcroft and Kannan

#### (a) Best fitting line



Denote by  $A$  the  $(n \times d)$  data matrix obtained by stacking the data vectors  $X_i \in \mathbb{R}^d$  as rows.

We need to minimize  $\sum_{i=1}^n \|X_i - (X_i \cdot v)v\|_2^2$  over  $v$ ,

which by Pythagorean Theorem is the same as obtaining

$$v_i = \operatorname{argmax}_{v \text{ s.t. } \|v\|_2=1} \sum_{i=1}^n (X_i \cdot v)^2 = \operatorname{argmax}_{v \text{ s.t. } \|v\|_2=1} \|Av\|_2^2 \quad (4)$$

Denote  $\sigma_i(A) = \|Av_i\|_2 \rightarrow$  this corresponds to the

$$\text{spectral norm } \|A\|_2 = \max_{v \text{ s.t. } \|v\|_2 \leq 1} \|Av\|_2$$

### (b) Best fitting k-rank space

We now move to the more general problem of finding a k-dimensional vector space  $V$  that minimizes

$$\sum_{i=1}^d \|X_i - \operatorname{proj}_V(X_i)\|_2^2.$$

Such a space can be found by the following greedy procedure:

$$1) \quad v_1 = \operatorname{argmax}_{v \text{ s.t. } \|v\|_2=1} \|Av\|_2$$

2) Repeat for  $i=2, \dots, k$

$$v_i = \operatorname{argmax}_{v \text{ s.t. } v \perp \{v_1, \dots, v_{i-1}\}} \|Av\|_2$$

$$\|v\|_2 = 1$$

Suppose that we continue the process above till we cannot find a  $v$  in step(2). Suppose  $\{v_1, \dots, v_{r_2}\}$  denote the obtained vectors. Then,

$\rightarrow$  row-rank of  $A$  is  $r_2$

$\rightarrow \{v_1, \dots, v_{r_2}\}$  constitutes an orthonormal basis of the row-space of  $A$  ( $\equiv \operatorname{span}\{x_1, \dots, x_n\}$ )

Note that  $\|X_i\|_2^2 = \sum_{j=1}^r (X_{i,j} v_j)^2$  and summing over  $i$

(5)

we get

$$\begin{aligned} \sum_{i=1}^n \|X_i\|_2^2 &= \sum_{i=1}^n \sum_{j=1}^r (X_{i,j} v_j)^2 = \sum_{j=1}^r \sum_{i=1}^n (X_{i,j} v_j)^2 \\ &= \sum_{j=1}^r \|A v_j\|_2^2 = \sum_{j=1}^r \sigma_j^2(A), \end{aligned}$$

where  $\sigma_j(A) = \|A v_j\|_2 \equiv j^{\text{th}}$  singular value.

Thus,  $\sum_{j=1}^r \sigma_j^2(A) = \|A\|_F^2 \equiv$  Frobenius norm of  $A$

→ The vectors  $v_1, \dots, v_r$  are called the right-singular vectors.

Let  $u_i = \frac{1}{\sigma_i(A)} A v_i$ ,  $1 \leq i \leq r$ .

The vectors  $u_1, \dots, u_r$  are the left-singular vectors.

(c) Singular Value Decomposition (SVD)

Consider  $\sum_{i=1}^r \sigma_i(A) u_i v_i^T \equiv \tilde{A}$ .

Then,  $\tilde{A} v_i = \sigma_i(A) u_i = A v_i$  for  $1 \leq i \leq r$ . Since  $v_1, \dots, v_r$

constitutes a basis for the row-space of  $A$ ,  $A = \tilde{A}$ .

(d) Best rank- $k$  approximation

Let  $A_k = \text{span}\{v_1, \dots, v_k\}$ .

Theorem. For any matrix  $B$  of rank at most  $k$ ,

(i)  $\|A - A_k\|_F \leq \|A - B\|_F$ , (ii)  $\|A - A_k\|_2 \leq \|A - B\|_2$ .

$\sigma_{k+1}$