

# Lecture 8

(1)

Review: Estimating the mean of a Gaussian using Scheffé  
(for  $d=1$ )

Step 1: Guesslist  $\mathcal{L} = \{P_{\hat{\mu}}, \hat{\mu} \in \{x_1, \dots, x_n\}\}$

(Since  $P(\exists i \text{ s.t. } |x_i - \mu| \leq \sqrt{\frac{2\sigma^2 \log \frac{2}{\epsilon}}{n}}) \geq 1 - \epsilon/2$ )

Step 2: Use Scheffé to find a  $\hat{P} \in \mathcal{L}$  s.t.

$$d(\hat{P}, P) \leq 9 \cdot \sqrt{\frac{1}{2n} \log \frac{2}{\epsilon}} + 8\Delta$$

Here Scheffé set  $A(i, j)$  correspond to intervals of the form  $(-\infty, a)$  or  $(a, \infty)$ .

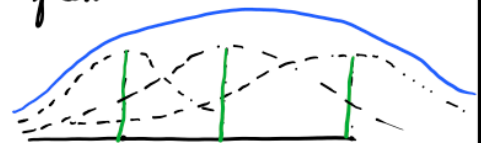
$$\text{Thus, } \Delta \leq d_{KS}(P_{\mu}, \underbrace{M_n}_{\text{empirical measure}}) \leq O\left(\sqrt{\frac{1}{n} \log \frac{1}{\epsilon}}\right).$$

Agenda: \* Learning Gaussian Mixtures

**A** Learning the distribution vs learning parameters

Given samples  $X_1, \dots, X_n$  generated iid from

$$\sum_{j=1}^K w_j N(\mu_j, \sigma_j^2 I_{d \times d}),$$



find:

- the weights  $w$ , means  $\mu$ , variances  $\sigma^2$  (parameter learning)
- a mixture  $\hat{P}$  that is close to  $P$  (proper learning of distribution)
- any distribution  $\hat{P}$  that is close to  $P$  (improper learning)

We will go with (b).

A difficulty in (a) is identifying which samples are associated with which component. This can only be done if we assume sufficient separation between the means.

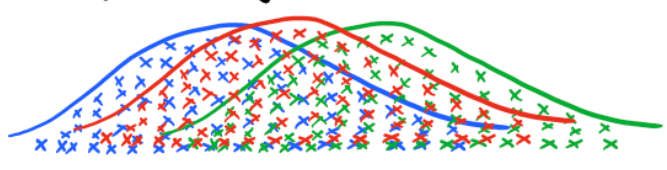
→ This difficulty can be circumvented in (b).

PAC formulation (We assume  $\sigma_j^2 = \sigma^2$  is known)

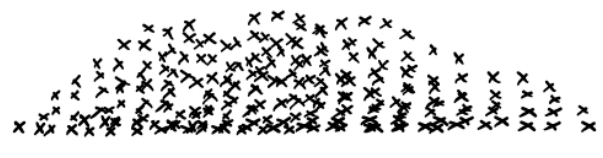
$$\epsilon(n, \delta, d) = \min_{\hat{\theta} = (\hat{\omega}, \hat{\mu})} \max_{\theta = (\omega, \mu)} P_{\theta} (d(P_{\theta}, P_{\hat{\theta}}(x^n)) > \delta)$$

B The 1-dimensional case

→ Difficulty in using empirical mean estimation:



↓ seen as



which samples are coming from blue, red, and green

If we have many, many samples, we can detect subtle differences b/w histograms corresponding to different mixtures.

Large Sample + poly(n) computational complexity  
⇓  
infeasibility of the algorithm

→ A Scheffé based approach

Step (i) Form a guesslist

(3)

Note that

$$\begin{aligned} & d(P_{\omega, \mu}, P_{\omega', \mu'}) \\ & \leq \sum_{i=1}^k |\omega_i - \omega'_i| + \sum_{i=1}^k \omega_i d(P_{\mu_i}, P_{\mu'_i}) \\ & \leq \sum_{i=1}^k |\omega_i - \omega'_i| + \sum_{i=1}^k \frac{\omega_i |\mu_i - \mu'_i|}{2\sigma} \end{aligned}$$

Let  $\Omega = \{ \omega : \omega_i \in \{0, \frac{\delta}{2k}, \frac{2\delta}{2k}, \dots, 1\}, 1 \leq i \leq k, \sum_{i=1}^k \omega_i = 1 \}$ .

Our desired guesslist is then given by

$$\mathcal{L} = \{ P_{\omega, \mu} : \omega \in \Omega, \mu_i \in \{x_1, \dots, x_n\} \forall 1 \leq i \leq k \}.$$

Then,

$$\begin{aligned} & P_{\omega, \mu} \left( \forall \hat{P} \in \mathcal{L}, d(\hat{P}, P_{\omega, \mu}) > \delta \right) \\ & \leq P_{\omega, \mu} \left( \forall (\hat{\mu}_1, \dots, \hat{\mu}_k) \text{ s.t. } \hat{\mu}_i \in \{x_1, \dots, x_n\}, \right. \\ & \quad \left. \sum_{j=1}^k \omega_j d(P_{\mu_j}, P_{\hat{\mu}_j}) > \frac{\delta}{2} \right) \\ & \leq P_{\omega, \mu} \left( \exists 1 \leq j \leq k \text{ s.t. } d(P_{\mu_j}, P_{x_i}) > \frac{\delta}{2} \text{ for} \right. \\ & \quad \left. \text{every } 1 \leq i \leq n \right) \end{aligned}$$

$$\leq \sum_{j=1}^k P_{\omega, \mu} \left( \forall 1 \leq i \leq n, |\mu_j - x_i| > \sigma \delta \right)$$

Denoting by  $N_j$  the number of samples from  $P_{\mu_j}$ ,

$$\begin{aligned} & P_{\omega, \mu} \left( \forall 1 \leq i \leq n, |\mu_j - x_i| > \sigma \delta \right) \\ & \leq \sum_{l=0}^n P_{\omega, \mu} (N_j = l) e^{-\frac{l\delta^2}{2}} \leq P_{\omega, \mu} (N_j \leq t) + e^{-t\delta^2/8}, \end{aligned}$$

for every  $t$ . We choose  $t = n w_j \theta$ ,  $\theta < 1$ . (4)

The right-side is bounded by (assuming  $w_j \leq 1/2$ )

$$e^{-c \cdot n \frac{w_j (1-\theta)^2}{w_j}} + e^{-n w_j \theta \delta^2 / 8}$$

which is less than  $\epsilon$  if  $n = O\left(\frac{1}{w_j \delta^2} \log \frac{1}{\epsilon}\right)$ .

Thus,  $n = O\left(\frac{1}{\min_j w_j} \cdot \frac{1}{\delta^2} \cdot \log \frac{k}{\epsilon}\right)$

suffice. We can improve by ignoring weights  $w_i \leq \frac{\delta}{4k}$ .

Then, our required prob. is bounded above by

$$\sum_{j \in [k], w_j > \frac{\delta}{4k}} P_{\omega, \mu} \left( \exists i \quad d(P_{\mu_j}, P_{x_i}) > \frac{\delta}{4} \right)$$

which can be bounded as above to get

$$n = O\left(\frac{k}{\delta^3} \log \frac{k}{\epsilon}\right) \text{ suffice.}$$

Step (iii) Scheffé Selector

Our list now contains  $\left(\frac{k}{\delta}\right)^k \cdot \binom{n}{k} = O\left(\frac{k^{2k}}{\delta^{4k}} \log^k \frac{k}{\epsilon}\right)$

Also, note that the Scheffé sets  $A(i, j)$  still

appear to be intervals. Thus, we can bound  $\Delta$

using  $c \cdot d_{KS}(P_\mu, \mu_n)$ .

[C] Can we use ML for selection?

The Yes Case We can do this if we assume bounds for

the probabilities in the guesslist over the support of the unknown  $P$ . (5)

Theorem Consider a set  $\mathcal{L}$  of distributions on  $\mathcal{X}$  with densities w.r.t.  $\nu$  s.t.  $\exists \hat{P} \in \mathcal{L}$  satisfying

$$D(P \parallel \hat{P}) \leq \delta.$$

Furthermore, assume  $\alpha \leq Q(x)$  for all  $x \in \text{supp}(P)$ , for every  $\hat{P} \in \mathcal{L}$ . Denote by  $\hat{P}_{ML}$  the maximizer of  $\sum_{i=1}^n \ln \frac{1}{Q(x_i)}$  over the set  $\mathcal{L}$ . Then,

$$P(D(P \parallel \hat{Q}_{ML}) > 4\delta) \leq (|\mathcal{L}| + 1) \exp\left(-\frac{2n \cdot \delta^2}{\log^2(1/\alpha)}\right).$$

Proof. The function  $\Lambda(Q) = -\sum_{i=1}^n \ln Q(x_i)$  is called the score function. Note that

$$\mathbb{E}_P[\Lambda(Q)] = n(H(P) + D(P \parallel Q)).$$

Note that  $\hat{P}_{ML}$  will satisfy the required condition if

$$\Lambda(\hat{P}) < \min \{ \Lambda(Q) : D(P \parallel Q) > 4\delta \},$$

which in turn will hold if

$$\Lambda(\hat{P}) \leq n(H(P) + 2\delta) \text{ and}$$

$$\Lambda(Q) \geq n(H(P) + 3\delta) \text{ for every } Q \text{ with } D(P \parallel Q) > 4\delta.$$

We bound the complement of this latter event.

⑥

$$\begin{aligned}
& P(\Delta(\hat{P}) > n(H(P) + 2\delta)) \\
&= P(\Delta(\hat{P}) > n(H(P) + \delta) + n\delta) \\
&\leq P(\Delta(\hat{P}) > n(H(P) + D(P \parallel \hat{P})) + n\delta) \\
&= P(\Delta(\hat{P}) > \mathbb{E}_P[\Delta(\hat{P})] + n\delta) \\
&\leq \exp\left(-\frac{n\delta^2}{\log^2(1/\alpha)}\right) \quad (\text{using Hoeffding's inequality})
\end{aligned}$$

Similarly, for a  $Q$  s.t.  $D(P \parallel Q) > 4\delta$ ,

$$\begin{aligned}
& P(\Delta(Q) \leq n(H(P) + 3\delta)) \\
&\leq P(\Delta(Q) \leq \mathbb{E}_P[\Delta(Q)] - n\delta) \\
&\leq \exp\left(-\frac{n\delta^2}{\log^2(1/\alpha)}\right).
\end{aligned}$$

The claimed bound follows from the union bound. ■

### The No example

$$P_1 = \text{unif}[-1, 1], \quad P_2 = \text{unif}[\delta, 1 + \delta], \quad P = \text{unif}[0, 1].$$

We see  $n$  samples from  $P$ .

Let  $N \equiv \#$  of samples in  $[0, \delta)$ .

ML chooses  $P_1$  if  $N > 0$ , which happens with prob.

$$P(N > 0) = 1 - P(N = 0) = 1 - (\delta)^n$$

But  $d(P_1, P) = 1/2$  and  $d(P_2, P) = \delta \Rightarrow$  ML doesn't

choose the closer  $P_i$  with prob.  $\rightarrow 1$  as  $n \rightarrow \infty$ .