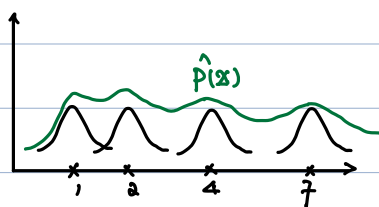


Observation: larger  $K$ , smoother output.

### 3) Poisson Window Estimation with Gaussian Kernels



→ Each sample has a Gaussian Distribution.

→ Add all the probability

→ Scale to have CDF = 1

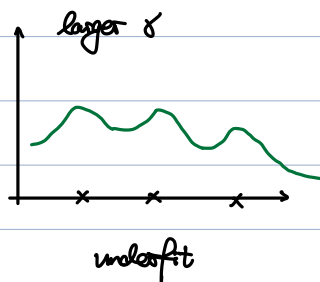
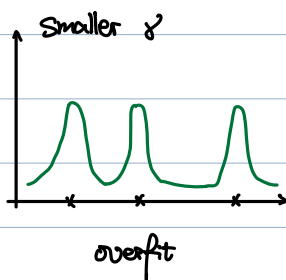
Details: Standard normal distribution (PDF)

$$\phi(z) = \frac{1}{\sqrt{2\pi}} e^{-\frac{z^2}{2}}$$

$$\hat{p}(x) = \frac{1}{K} \sum_{i=1}^N \phi\left(\frac{\|x - x_i\|}{\gamma}\right)$$

$\gamma$ : kernel width, width of each sample affected neighbor.

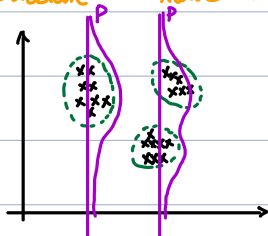
$K$ : Normalizing constant



Note: If more points → complex computation.

## Sec 22.

### 4) Gaussian Mixture Model



Each cluster with its own Gaussian distribution

It suffices to use only  $k$  Gaussian distributions, one for each cluster.

- mean:  $\mu_1, \mu_2, \dots, \mu_k \in \mathbb{R}^d$

covariance:  $\Sigma_1, \Sigma_2, \dots, \Sigma_k \in \mathbb{R}^d \times \mathbb{R}^d$

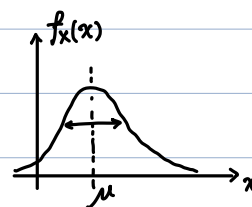
Review:

$d=1$ : PDF of Gaussian R.V.  $X$  with mean  $\mu$ , and variance  $\sigma^2$ .

$$f_X(x) \equiv \mathcal{N}(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

$d=2$ : PDF of jointed Gaussian R.V.  $X_1$  and  $X_2$ , with mean  $\mu_1$  &  $\mu_2$ ,

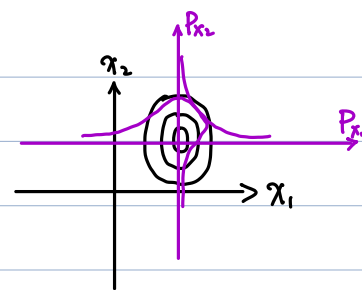
variances  $\sigma_1^2$  &  $\sigma_2^2$ , and correlation coefficient  $\rho$ .



$$\rho = \frac{E[(X_1 - \mu_1)(X_2 - \mu_2)]}{\sigma_1 \cdot \sigma_2}$$

$$f_{x_1, x_2}(x_1, x_2) \equiv \mathcal{N}(x_1, x_2; \mu_1, \mu_2, \sigma_1^2, \sigma_2^2, \rho)$$

$$= \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \cdot e^{-\frac{1}{2(1-\rho^2)}\left[\left(\frac{x_1-\mu_1}{\sigma_1}\right)^2 - \frac{2\rho(x_1-\mu_1)(x_2-\mu_2)}{\sigma_1\sigma_2} + \left(\frac{x_2-\mu_2}{\sigma_2}\right)^2\right]}$$



General  $d \gg 1$ :

Consider  $d$  jointly Gaussian RVs,  $X_1, X_2, \dots, X_d$ , with  $\mu_1, \mu_2, \dots, \mu_d$  and covariance matrix.

$$\Sigma \triangleq \begin{bmatrix} \text{cov}(X_1, X_1) & \text{cov}(X_1, X_2) & \dots & \text{cov}(X_1, X_d) \\ \vdots & \ddots & \ddots & \vdots \\ \text{cov}(X_d, X_1) & \dots & \dots & \text{cov}(X_d, X_d) \end{bmatrix}$$

$$\text{where } \text{cov}(X_i, X_j) = E[(X_i - \mu_i)(X_j - \mu_j)]$$

$$\Leftrightarrow \Sigma = E[(X - \mu)(X - \mu)^T]$$

$$f_X(x) = \mathcal{N}(x, \mu, \Sigma)$$

$$= \frac{1}{(2\pi)^{d/2} \cdot \det(\Sigma)^{1/2}} \cdot e^{-\frac{1}{2}(x - \mu)^T \Sigma^{-1} (x - \mu)}$$

$$e^{-\frac{1}{2} \frac{(x - \mu)^2}{\sigma^2}}$$

$$\sqrt{2\pi\sigma^2}$$

Data Generalization Model (Hypothesis)

(for each data point in  $\mathcal{D}$ )

First, pick a bump (cluster)

$j \in \{1, 2, \dots, k\}$  according to probability distribution  $\{w_1, w_2, \dots, w_k\}$

$$\begin{cases} w_i > 0 \\ \sum_{i=1}^k w_i = 1 \end{cases}$$

Then, generate  $x$  according to probability density  $P(x|j) = \mathcal{N}(x; \mu_j, \Sigma_j)$

$$= \frac{1}{(2\pi)^{d/2} \cdot \det(\Sigma_j)^{1/2}} \cdot e^{-\frac{1}{2}(x - \mu_j)^T \Sigma_j^{-1} (x - \mu_j)}$$

$\Rightarrow$  the overall PDF for  $x$  is

$$P(x) = \sum_{j=1}^k w_j P(x|j) = \sum_{j=1}^k w_j \cdot \mathcal{N}(x; \mu_j, \Sigma_j)$$

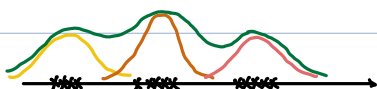
Example:  $d=1, k=3$

$$\mathcal{D} = \{x_1, x_2, \dots, x_N\}$$

For each index  $n$ ,

random switch  $\rightarrow$

$$\begin{aligned} w_1 & \xrightarrow{\text{random switch}} \mathcal{N}(x; \mu_1, \sigma_1^2) \rightarrow \\ w_2 & \xrightarrow{\text{random switch}} \mathcal{N}(x; \mu_2, \sigma_2^2) \rightarrow x_n \\ w_3 & \xrightarrow{\text{random switch}} \mathcal{N}(x; \mu_3, \sigma_3^2) \rightarrow \end{aligned}$$



Given  $\mathcal{D}$ , need to estimate  $\mu_1, \mu_2, \mu_3$  &  $\sigma_1, \sigma_2, \sigma_3$ .

Given  $\mathcal{D} = \{x_1, \dots, x_N\}$ , and  $k$ . find a GMM  $\Omega = \{w_j, \mu_j, \Sigma_j\}_{j=1}^k$  that is the best fit for  $\mathcal{D}$ .

Best-fit: MLE

$$\hat{P}_{\Omega}(\mathcal{D}) = \prod_{n=1}^N \hat{P}(x_n)$$

is maximized.

$$\text{Minimize } E_{\Omega}(\mathcal{D}) \triangleq -\log \hat{P}_{\Omega}(\mathcal{D})$$

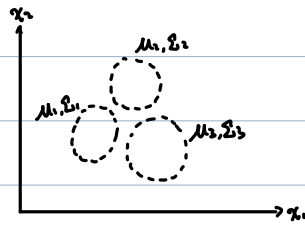
$$= -\sum_{n=1}^N \log \hat{P}(x_n)$$

$$E_n(\Omega) = -\log \hat{P}(x_n)$$

How to minimize  $E_n(\Omega)$ ?

SGD?

## Lecture 23



$$D = \{x_1, x_2, \dots, x_N\}$$

$$P(x) = \sum_{j=1}^K w_j \cdot \mathcal{N}(x_j; \mu_j, \Sigma_j)$$

Given set  $D$  and  $K$ , find the GMM  $\Omega = \{w_j, \mu_j, \Sigma_j\}_{j=1}^K$  that is the "best fit" for  $D$ .

Maximize the likelihood

$$\hat{P}_\Omega(D) = \prod_{n=1}^N \hat{P}(x_n)$$

$$\equiv \text{Minimize } E_n(\Omega) = -\log \hat{P}_\Omega(D)$$

$$= \sum_{n=1}^N -\log \hat{P}(x_n)$$

$$= \sum_{n=1}^N -\log \underbrace{\sum_{j=1}^K w_j \cdot \mathcal{N}(x_n; \mu_j, \Sigma_j)}_{E_n(\Omega)}$$

Can use SGD? No

< Complicated  $E_n$   
 $\sum_j w_j = 1$  . constrained optimization

## The Expectation Maximization (EM) Algorithm

An alternating optimization approach with 2 subproblems.

### Subproblem 1: (M step)

Given  $D = \{x_1, \dots, x_N\}$ , suppose for each  $x_i \in D$ , we know the bump (Gaussian distribution) it belongs to.

Let  $B_j \in D$  denotes the  $j$ th bump.

i.e. all points in  $B_j$  are sampled from

$$P(x|j) = \mathcal{N}(x; \mu_j, \Sigma_j)$$

Estimate  $w_j, \mu_j, \Sigma_j$  for  $j=1, 2, \dots, K$

$$w_j = \frac{N_j}{N} \quad N_j: \text{\# of points in } B_j$$

$$\mu_j = \frac{1}{N_j} \sum_{x_n \in B_j} x_n \quad (\text{sample mean})$$

$$\Sigma_j = \frac{1}{N_j} \sum_{x_n \in B_j} (x_n - \mu_j)(x_n - \mu_j)^T$$

(sample covariance)

### Subproblem 2: (E step)

Given  $\Omega = \{\omega_j, \mu_j, \Sigma_j\}_{j=1}^K$ , estimate bump membership, i.e. for each  $x_n$ , find bump  $j$  that is most likely to produce  $x_n$ .

$$x_n \in B_{j^*} \text{ if } j^* = \underset{j \in \{1, 2, \dots, K\}}{\operatorname{argmax}} P(j|x_n)$$

Recall: Bayes' Theorem

$$P(j|x_n) = \frac{P(x_n|j) \cdot P(j)}{P(x_n)} = \frac{P(x_n|j) \cdot P(j)}{\sum_{i=1}^K P(x_n|i) \cdot P(i)} = \frac{\mathcal{N}(x_n; \mu_j, \Sigma_j) \cdot \omega_j}{\sum_{i=1}^K \mathcal{N}(x_n; \mu_i, \Sigma_i) \cdot \omega_i}$$

$$\Rightarrow j^* = \underset{j}{\operatorname{argmax}} \mathcal{N}(x_n; \mu_j, \Sigma_j) \cdot \omega_j$$

EM algo. Summary: (Hard decisions)

1. Initialization: start with arbitrary bump membership for each  $x_n$ .
  2. Estimate  $\Omega = \{\omega_j, \mu_j, \Sigma_j\}_{j=1}^K$  given bump membership  $\{B_1, \dots, B_K\}$  (subproblem 1)
  3. Estimate bump membership  $\{B_1, \dots, B_K\}$  given  $\Omega$  (subproblem 2)
- Repeat step 2 & 3 until convergence.

Note: ① Convergence is guaranteed.

②  $B_1, \dots, B_K$  are auxiliary ("hidden" variables).