

Lecture 7 : 28th Jan 2026:

Density Estimation

Recall : $h^*(x) = \underset{f}{\operatorname{argmin}} \left[\mathbb{E}_{P_{xy}} [L(h(x), y)] \right]$

for 0-1 loss,

$$h_B^*(y) = \begin{cases} 1 & \text{if } p_{y=1|x} > p_{y=0|x} \\ 0 & \text{otherwise} \end{cases}$$

* Implementing the optimal Bayes classifier demands estimating densities from data

Eg: Given D , Estimate $p_x, p_y, p_{x|y}$

Maximum Likelihood (Minimal KL) Estimation:

Given $V = \{\omega_1, \omega_2, \dots, \omega_n\} \sim \text{iid } P_V$,

Estimate p_y

Start with p_θ &

$$\Theta = \underset{\theta}{\operatorname{argmin}} D_{KL}(P_V || P_\theta)$$

$$\hat{\theta} = \operatorname{argmax}_{\theta} \left[\frac{1}{n} \sum_{i=1}^n \log p_{\theta}(\mathbf{v}_i) \right]$$

$$\mathbf{v}_i \sim \text{iid } P_v$$

Examples of ML Estimation:

a) $p_{\theta}(\mathbf{v}) \sim N(\mathbf{v}; \theta, \mathbf{I})$ model

$$\mathbf{v} \in \mathbb{R}^d, \theta \in \mathbb{R}^d$$

we have $\mathbf{v}_1, \mathbf{v}_2 \dots \mathbf{v}_n \sim \text{iid } P_v$

$$p_{\theta}(\mathbf{v}) = \frac{1}{(2\pi)^{d/2} |\mathbf{I}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{v}-\theta)^T \mathbf{I} (\mathbf{v}-\theta) \right\}$$

Consider

$$l(\theta) = \frac{1}{n} \sum_{i=1}^n \log p_{\theta}(\mathbf{v}_i)$$

$$\propto -\frac{1}{n} \sum_{i=1}^n \|\mathbf{v}_i - \theta\|_2^2$$

$$\theta^* = \operatorname{argmax}_{\theta} l(\theta)$$

$$\theta^* = \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \quad \text{MLE -}$$

2) Example : Generalized Discrete RV

$$D = \{ \omega_1, \omega_2, \dots, \omega_n \} \sim \text{iid } P_V$$

$\omega_i \in \{a_1, a_2, \dots, a_m\}$ with Prob

$$\{p_1, p_2, \dots, p_m\}$$

parameters $\Theta = \{p_1, p_2, \dots, p_m\}$

To express the mass fun^c for this RV, we need
One-hot representation.

One-hot Representation :

For each ω_i , define an auxiliary RV Z_i as follows :

$$Z_i = [z_i^1, z_i^2, \dots, z_i^m], \quad z_i^j \in \{0, 1\}$$

$$z_i^j = \begin{cases} 1 & \text{for } j: \omega_i = a_j \\ 0 & \text{otherwise} \end{cases}$$

$$\text{Eg : } D = \{2, 3, 4, 1, 6, \dots\}$$

$$\omega_i \in \{1, 2, 3, 4, 5, 6\}$$

$$z_1 = [0 \ 1 \ 0 \ 0 \ 0 \ 0]$$

$$z_2 = [0 \ 0 \ 1 \ 0 \ 0 \ 0]$$

New data $D = \{z_1, z_2, \dots, z_n\}$

$$z_i \in \{0,1\}^m$$

$$p_\theta(u_i) = \prod_{j=1}^m p_j^{z_{ij}}$$

$$\Theta = \{p_1, p_2, \dots, p_m\}$$

$$\theta^* = \underset{\Theta}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n \log p_\theta(u_i)$$

$$= \underset{\Theta}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n \log \frac{m}{\prod_{j=1}^m p_j^{z_{ij}}}$$

$$= \underset{\Theta}{\operatorname{argmax}} \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m z_{ij} \log p_j$$

adding Constraints to the above problem

$$0 \leq p_j \leq 1 \quad \& \quad \sum_{j=1}^m p_j = 1$$

New objective

$$\underset{\Theta}{\operatorname{argmax}} \lambda(\theta)$$

$$\text{s.t. } 0 \leq p_j \leq 1 \quad \& \quad \sum_{j=1}^m p_j = 1$$

writing the lagrangian

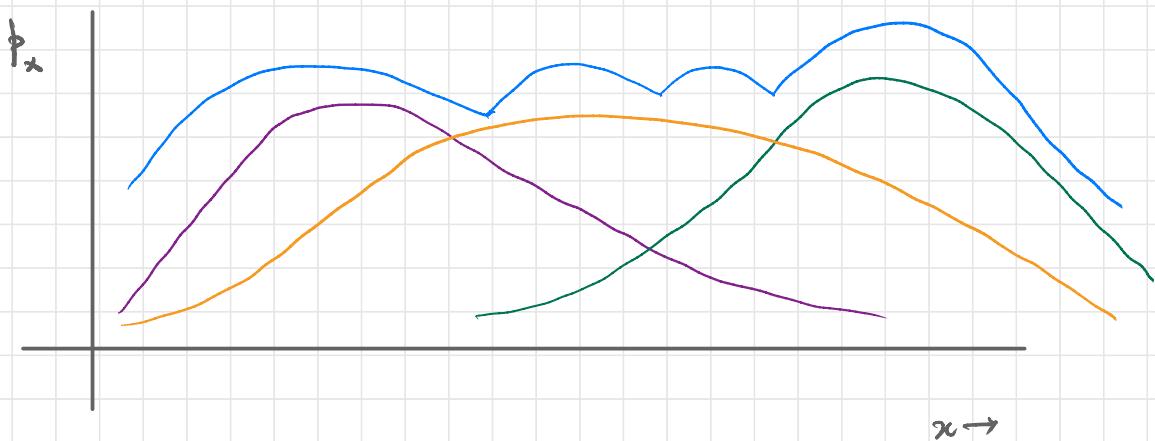
$$\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^m \beta_i^j \log p_j + \lambda \left(\sum_{j=1}^m p_j - 1 \right) + \alpha (p_j - 1) + \beta (-p_j)$$

after Solving the above problem

$$p_j^* = \frac{\# a_j \text{ as occurred}}{n}$$

$$= \frac{1}{n} \sum_{i=1}^n \beta_i^j$$

Density Estimation for Mixture Distributions:



Multi-model densities are not "well-estimated" with uni-model models.

Soluⁿ: Try a multi-modal model.

Mixture density fun^c:

Suppose $\omega \in \mathbb{R}^d$

Define a mixture density

$$\rho_{\theta}(\omega) = \sum_{j=1}^m \alpha_j \rho_{\theta_j}(\omega)$$

$\alpha_j \in [0, 1]$ ρ_{θ_j} : density fun^c

Gaussian Mixture models : GMM.

$$\rho_{\theta_j} \sim NC(\cdot; \theta_j)$$

GMMs are Universal density approximations for CRVs.

Lecture 8 : 2nd Feb 2026

Mixtures Densities & EM

Recall

A mixture density model

$$p_{\theta}(v) = \sum_{j=1}^m \alpha_j p_{\theta_j}(v)$$

Where $p_{\theta_j}(v)$ is a density func

$$0 \leq \alpha_j \leq 1 \quad \text{and} \quad \sum_{j=1}^m \alpha_j = 1$$

Gaussian Mixture Model

$$p_{\theta_j}(v) \sim N(\cdot; \mu_j, \Sigma_j) \quad j=1, \dots, m$$

MLE for mixture density

$$\Theta^* = \underset{\Theta}{\operatorname{argmax}} \sum_{i=1}^n \log p_{\theta}(v_i)$$

in The Case of mixture densities.

$$\Theta = \left\{ \alpha_j, \mu_j, \Sigma_j \right\}_{j=1}^m$$

$$\sum_{i=1}^n \log b_\theta(x_i) = \sum_{i=1}^n \log \left(\sum_{j=1}^m \alpha_j b_{\theta_j}(x_i) \right)$$

Log of Sums.
exp will not get cancelled with log

∴ we need an iterative algorithm.

Latent Variable Models:

Given $D = \{x_1, x_2, \dots, x_n\} \sim \text{iid } IP_x$
 $x_i \in \mathbb{R}^d$

Introduce into the system another RV into the system which is un-observed.

→ this is known as hidden / latent RV. : Z

Mathematically

$\forall x_i \in \mathbb{R}^d$, define $z_i \in \mathbb{R}^k$, typically $k \ll d$

→ x & Z are assumed to be correlated.

→ Z can be CRV @ DRV.

with this latent variable model is defined as

$$p_{\theta}(x) = \sum_z p_{\theta}(x, z) \quad \text{if } Z \text{ is DRV}$$

$$p_{\theta}(x) = \int_z p_{\theta}(x, z) dz \quad Z \text{ is CRV}$$

$$\left. \begin{array}{l} p_{\theta}(x|z) \\ p_{\theta}(z|x) \\ p_{\theta}(z) \end{array} \right\} \text{All these exists}$$

MLE for Latent Variable Models:

- * For a latent variable models, both the model parameters & the density over the latent variable needs to be estimated.
- * we have the Likelihood func $l(\theta)$

Note: we have considered one data point

$$\begin{aligned} l(\theta) &= \log p_{\theta}(x) \\ &= \log \sum_z p_{\theta}(x, z) \end{aligned}$$

Assuming Z is DRV.

Suppose $g(z)$ is some density over Z .

$$l(\theta) = \log \sum_z \left[p_\theta(x, z) \frac{g(z)}{\bar{g}(z)} \right]$$

$$= \log \sum_z g(z) \underbrace{\left(\frac{p_\theta(x, z)}{g(z)} \right)}_{\text{Func of } Z}$$

$$= \log \mathbb{E}_{\frac{p_\theta(x, z)}{g(z)}} \left[\frac{p_\theta(x, z)}{g(z)} \right]$$

Using Jensen's inequality & noting that

\log is a concave func

$$\log(\mathbb{E}_x) \leq \mathbb{E} \log(x) : \log(x) \text{ is convex.}$$

\therefore we have

$$\log \mathbb{E}_{\frac{p_\theta(x, z)}{g(z)}} \left(\frac{p_\theta(x, z)}{g(z)} \right) \geq \mathbb{E}_{\frac{p_\theta(x, z)}{g(z)}} \log \left[\frac{p_\theta(x, z)}{g(z)} \right]$$

$$l(\theta) \geq F_\theta(g)$$

* $F_\theta(g)$ is now computable " ; it has
"Sum of log" type terms.

* $F_\theta(q)$ is a lower bound on $\ell(\theta)$

$\ell(\theta)$: log Likelihood @ Evidence.

∴ Evidence lower bound (ELBo)

q : Variational distribution

New optimization problem:

$$\theta^* = \underset{\theta}{\operatorname{argmax}} \ell(\theta)$$

$$\hat{\theta}^*, q = \underset{\theta, q}{\operatorname{argmax}} F_\theta(q) \quad \therefore \text{modified optimization}$$

Question: What choice of $q(z)$ would make ELBo tight?

Consider

$$\ell_\theta - F_\theta(q) = \log p_\theta(x) - \sum_z q(z) \log \frac{p_\theta(x, z)}{q(z)}$$

$$= \log p_\theta(x) - \sum_z q(z) \log \frac{p_\theta(x) p_\theta(z|x)}{q(z)}$$

$$= \log p_\theta(x) - \log p_\theta(x) - \sum_z q(z) \log \frac{p_\theta(z|x)}{q(z)}$$

$$= \sum_z q(z) \log \frac{q(z)}{p_\theta(z|x)}$$

$$= D_{KL}(q(z) || p_\theta(z|x))$$

More $F_\theta(q) = \mathcal{L}(\theta)$ by construction.

$$\mathcal{L}(\theta) - F_\theta(q) = 0 \quad \text{iff} \quad D_{KL}(q(z) || p_\theta(z|x)) = 0$$

$$\therefore q^*(z) = p_\theta(z|x)$$

makes the ELBO tight (exactly same)

* if we know $p_\theta(z|x)$ then we need to optimize the modified optimization only on θ

An iterative Algorithm for optimizing $\mathcal{L}(\theta)$:

1) Initialize θ'

2) For $t=1$ to Convergence

Compute $q^{t+1}(z) = p_{\theta^t}(z|x)$

Expectation Step

$$F_{\theta^t}(q^{t*}) = \mathbb{E}_{q_{\theta^t}^{t*}(z)} \log \frac{p_{\theta}(x, z)}{q_{\theta}^{t*}(z)}$$

Estimate

$$\theta^{t+1} \leftarrow \text{argmax} [F_{\theta^t}(q^{t*})]$$

Maximization Step

To show EM guarantees

$$\ell(\theta^{t+1}) \geq \ell(\theta^t)$$