

Final exam 12/28, close book

12/21 review, 羽题课?

Lecture 12.

EM Algorithm.

- Maximal likelihood estimation with unknown parameters

Example: Let w be Bernoulli r.v., $w = \begin{cases} 1 & \text{w.p. } \delta \\ 0 & \text{w.p. } 1-\delta \end{cases}$, y be a noise observation of w .
⊗ $P_{Y|W}(y|w) = \begin{cases} \varepsilon, & y \neq w \\ 1-\varepsilon, & y=w \end{cases}$

Let $\underline{w} = [w_1, \dots, w_n]^T$ be n iid samples of w .

$$\underline{y} = [y_1, \dots, y_n]^T$$

We do not observe \underline{w} , but observe \underline{y} .

Our goal is to estimate δ and ε .

$$\text{ML: } P_Y(\underline{y}; \varepsilon, \delta) = \prod_{i=1}^n P_{Y_i}(y_i; \varepsilon, \delta) = \prod_{i=1}^n \left[P_{Y_i|W_i}(y_i | w_i=0; \varepsilon, \delta) \cdot P_{W_i}(0; \varepsilon, \delta) + P_{Y_i|W_i}(y_i | w_i=1; \varepsilon, \delta) \cdot P_{W_i}(1; \varepsilon, \delta) \right]$$

hidden/latent r.v.
Bayes rule.

$$(\varepsilon^*, \delta^*) = \arg \max_{\varepsilon, \delta} P_Y(\underline{y}; \varepsilon, \delta)$$

- No close form solution for the optimal ε, δ .
- Algorithm to compute the optimal parameters with hidden variables \Rightarrow EM Alg.

General setup:

- Assume complete data Z , generated by $P_Z(z; x)$
- Can only observe $y = g(z)$, observed data, g : deterministic function.
- In our example: $Z = [\underline{\omega}, \underline{y}]$, $y = \underline{y}$, $x = [\varepsilon, \delta]$

The goal is to find $x^* = \arg \max_x P_Y(y; x) = \arg \max_x \log P_Y(y; x)$

$$\text{Note that } P_Z(z; x) = \underbrace{\sum_y P_{Z|Y}(z|y; x) \cdot P_Y(y; x)}_{P_{Z|Y}(z|g(z); x) \cdot P_Y(g(z); x)} \Rightarrow \text{let } y = g(z) : \log P_Y(y; x) = \log P_Z(z; x) - \log P_{Z|Y}(z|y; x)$$

Do expectation of both sides over $P_{Z|Y}(z|y; x')$:

$$\text{LHS} = \sum_{z'} P_{Z|Y}(z'|y; x') \log P_Y(y; x) = \left(\sum_z P_{Z|Y}(z|y; x') \right) \cdot \log P_Y(y; x) = \log P_Y(y; x)$$

$$\begin{aligned} \text{RHS} &= \mathbb{E}[\log P_Z(z; x) | Y=y, x=x'] - \mathbb{E}[\log P_{Z|Y}(z|y; x) | Y=y, x=x'] \\ &= \sum_{z'} P_{Z|Y}(z|y; x') \log P_Z(z'; x) \stackrel{\text{IIA}}{=} - \sum_z P_{Z|Y}(z|y; x') \log P_{Z|Y}(z|y; x) \stackrel{\text{IIA}}{=} V(x, x') \end{aligned}$$

$$\Rightarrow \log P_Y(y; x) = U(x, x') + V(x, x'), \forall x'$$

Lemma: $V(x, x') \geq V(x', x')$

K-L divergence

$$\text{Proof}: V(x, x') - V(x', x') = \sum_z P_{Z|Y}(z|y; x') \log \frac{P_{Z|Y}(z|y; x')}{P_{Z|Y}(z|y; x)} = D(P_{Z|Y; x'} \parallel P_{Z|Y; x}) \geq 0 \quad \square$$

\Rightarrow If we can find $x \neq x'$ such that $U(x, x') \geq U(x', x')$, then

$$\log P_Y(y; x) = U(x, x') + V(x, x') \geq U(x', x') + V(x', x') = \log P_Y(y; x')$$

EM Algorithm:

(i) Initialization : choose a $\hat{x}^{(0)}$

(ii) Repeat until convergence :

E-step : given the previous estimation $\hat{x}^{(n)}$, compute.
 \uparrow
 expectation. $U(x, \hat{x}^{(n)}) = \mathbb{E}[\log P_{\epsilon}(z; x) | Y=y; x = \hat{x}^{(n)}]$

M-step : Find $\hat{x}^{(n+1)}$ maximizing $U(\cdot; \hat{x}^{(n)})$

$$\begin{aligned} \uparrow \\ \text{maximization} \end{aligned} \quad \hat{x}^{(n+1)} &= \arg \max_x U(x; \hat{x}^{(n)}) \\ \Rightarrow U(\hat{x}^{(n+1)}, \hat{x}^{(n)}) &\geq U(\hat{x}^{(n)}, \hat{x}^{(n)}) \end{aligned}$$

- We can get a sequence of $\hat{x}^{(0)}, \hat{x}^{(1)}, \dots$ such that

$$P_Y(y; \hat{x}^{(0)}) \leq P_Y(y; \hat{x}^{(1)}) \leq \dots \dots$$

since $P_Y(y; x) \leq 1$, it is a non-decreasing bound sequence
 \Rightarrow it must converges.

- EM Algorithm converges to a stationary point of the likelihood function, i.e.

Let x^* be the convergent point, then $\frac{\partial}{\partial x} P_Y(y; x)|_{x=x^*} = 0$

EM Algorithm for mixture model.

Mixture model :

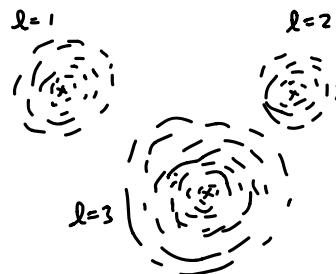
Assume data generated by the following process:

1. sample $l_i \in \{1, \dots, k\}$, $i=1, \dots, m$ and $l_i \stackrel{iid}{\sim} \text{multinomial}(\underline{\phi})$, $\underline{\phi} = [\phi_1, \dots, \phi_k]$
 $P(l_i=j) = \phi_j, \forall j=1 \dots k, i=1, \dots, m$ $\sum_{i=1}^k \phi_i = 1$

2. Sample observation y_i from some distribution $P(l_i, y_i)$

$$P(l_i, y_i) = P(l_i) \cdot \underbrace{P(y_i | l_i)}_{\text{generating the } y_i \text{ from } l_i}$$

- Mixture Gaussian model: $P(y_i | l_i=j) \sim N(\mu_j, \Sigma_j)$



- In our notation, $x = [\underline{\phi}, \underline{\mu}, \underline{\Sigma}] : \begin{cases} \underline{\phi} = [\phi_1, \dots, \phi_k] \\ \underline{\mu} = [\mu_1, \dots, \mu_k] \\ \underline{\Sigma} = [\Sigma_1, \dots, \Sigma_k] \end{cases}$
 $z = [l_1, \dots, l_m, y_1, \dots, y_m]$
 $y = [y_1, \dots, y_m]$

The EM Algorithm: $(\hat{\underline{\phi}}^{(n)}, \hat{\underline{\mu}}^{(n)}, \hat{\underline{\Sigma}}^{(n)})$

$$\begin{aligned} \text{E-step} : U(x, \hat{x}^{(n)}) &= \mathbb{E}_{P_z|Y}(\cdot | y; \hat{x}^{(n)}) [\log P_z(z; x) | Y=y; \hat{x}^{(n)}] \\ &= \sum_{i=1}^m \sum_{l_i=1}^k P(l_i | y_i; \hat{x}^{(n)}) \log P(l_i, y_i; x) \end{aligned}$$

$$\text{M-step} : \hat{x}^{(n+1)} = \underset{x}{\operatorname{argmax}} \quad U(x; \hat{x}^{(n)})$$

The Mixture Gaussian :

$$P(l_i=j | y_i; \hat{x}^{(n)}) = \frac{P(y_i | l=j; \hat{x}^{(n)}) \cdot P(l_i=j; \hat{x}^{(n)})}{P(y_i; \hat{x}^{(n)})} \triangleq w_{ij}$$

$$\underbrace{\sum_{j=1}^k \underbrace{P_{e,r}(y_i, l_i=j; \hat{x}^{(n)})}_{P(y_i | l=j; \hat{x}^{(n)}) \cdot P(l_i=j; \hat{x}^{(n)})}}$$

\Rightarrow Compute $U(x; \hat{x}^{(n)}) \quad \forall x$

$$\text{find } x^* = \arg \max U(x; \hat{x}^{(n)})$$

M-step : updating parameters

$$\hat{\phi}_j^{(n+1)} = \frac{1}{m} \sum_{i=1}^m w_{ij}$$

$$\hat{\mu}_j^{(n+1)} = \frac{\sum_{i=1}^m w_{ij} \cdot y_i}{\sum_{i=1}^m w_{ij}}$$

$$\hat{\Sigma}_j^{(n+1)} = \frac{\sum_{i=1}^m w_{ij} \cdot (y_i - \hat{\mu}_j^{(n)}) (y_i - \hat{\mu}_j^{(n)})^\top}{\sum_{i=1}^m w_{ij}}$$

Minimal Mean square Error Estimator

X : prior distribution $P_X(x)$, $x \rightarrow P_{Y|X} \rightarrow Y$

$$\hat{X}_{MMSE}(Y) \triangleq \underset{\hat{x}}{\operatorname{argmin}} \mathbb{E}[(x - \hat{x}(Y))^2]$$

↑ mean square error.

$$\text{Thm: } \hat{X}_{MMSE}(Y) = \mathbb{E}[X|Y]$$

$$\begin{aligned} \text{Proof: for any } f(Y), \mathbb{E}[(x - f(Y))^2] &= \mathbb{E}\left[\left((x - \mathbb{E}[X|Y]) + (\mathbb{E}[X|Y] - f(Y))\right)^2\right] \\ &= \mathbb{E}[(x - \mathbb{E}[X|Y])^2] + 2 \mathbb{E}[(x - \mathbb{E}[X|Y])(\mathbb{E}[X|Y] - f(Y))] + \mathbb{E}[(\mathbb{E}[X|Y] - f(Y))^2] \\ &\quad \text{!} \\ \mathbb{E}[(x - \mathbb{E}[X|Y])(\mathbb{E}[X|Y] - f(Y))] &= \mathbb{E}[x \cdot \mathbb{E}[X|Y]] - \mathbb{E}[(\mathbb{E}[X|Y])^2] - \mathbb{E}[x \cdot f(Y)] + \mathbb{E}[\mathbb{E}[X|Y] \cdot f(Y)] \\ &\quad \text{!} \\ &= \sum_{x,y} x' P_{X|Y}(x|y) \cdot P_Y(y) \cdot \left(\sum_x x \cdot P_{X|Y}(x|y) \right) \\ &\quad \text{!} \\ &= \sum_y P_Y(y) \cdot \left(\sum_x x \cdot P_{X|Y}(x|y) \right)^2 \\ &\quad \text{!} \\ &= \mathbb{E}[x \cdot f(Y)] \\ &= 0 \end{aligned}$$

$$\begin{aligned} \mathbb{E}[(x - f(Y))^2] &= \mathbb{E}[(x - \mathbb{E}[X|Y])^2] + \mathbb{E}[(\mathbb{E}[X|Y] - f(Y))^2] \\ &\geq \mathbb{E}[(x - \mathbb{E}[X|Y])^2] \end{aligned}$$

$$\text{Example: } Y = X + Z \quad X \sim N(0, \sigma_x^2), \quad Z \sim N(0, \sigma_z^2)$$

\uparrow signal \uparrow noise

$$\text{PDF } f_{X|Y}(x|y) \sim N\left(\frac{\sigma_x^2}{\sigma_x^2 + \sigma_z^2} y, \frac{\sigma_x^2 \sigma_z^2}{\sigma_x^2 + \sigma_z^2}\right)$$

$$\hat{X}_{MMSE}(Y) = \frac{\sigma_x^2}{\sigma_x^2 + \sigma_z^2} \cdot y, \quad \text{MSE} = \frac{\sigma_x^2 \cdot \sigma_z^2}{\sigma_x^2 + \sigma_z^2} = \left(\frac{\sigma_x^2}{\sigma_x^2 + \sigma_z^2}\right) \cdot \sigma_z^2$$

Remark: If cannot observe y , random guess $\Rightarrow \text{MSE} = \sigma_z^2$

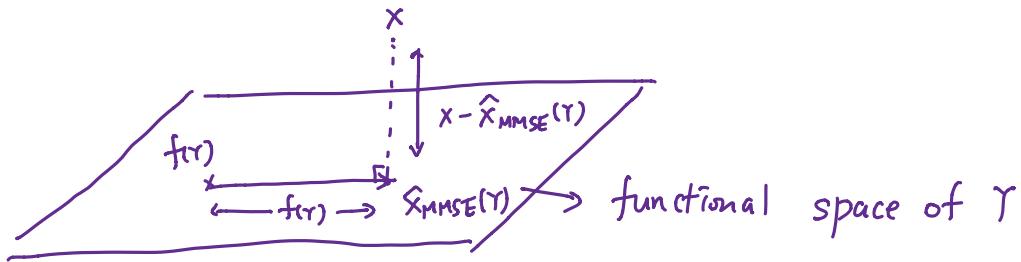
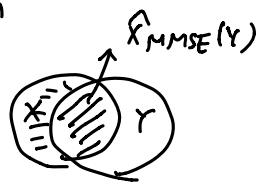
Orthogonality principle:

Thm: suppose \hat{X}_{MMSE} is the MMSE estimator of X given Y , then

$$\mathbb{E}[(X - \hat{X}_{MMSE}(Y)) \cdot f(Y)] = 0, \forall f(Y)$$

\Rightarrow Estimation error and observation are uncorrelated.

Proof: $\mathbb{E}[(X - \mathbb{E}[X|Y]) \cdot f(Y)] = \mathbb{E}[X \cdot f(Y)] - \mathbb{E}[\mathbb{E}[X|Y] \cdot f(Y)] = 0 \quad \square$



Corollary: If $\hat{X}_{MMSE}(Y) = a \cdot Y$, then $a = \frac{\mathbb{E}[XY]}{\mathbb{E}[Y^2]}$

Proof: $\mathbb{E}[(X - aY) \cdot Y] = 0 \Rightarrow a \cdot \mathbb{E}[Y^2] = \mathbb{E}[XY] \quad \square$

Linear regression: $X \in \mathbb{R}$, $Y \in \mathbb{R}^m$, P_{XY}

$$\text{Take } \hat{X}(Y) = A \cdot Y$$

$$\text{Goal: } A^* = \underset{A}{\operatorname{argmin}} \underset{P_{XY}}{\mathbb{E}} [(X - AY)^2]$$

$$\Rightarrow A^* = K_{XY} \cdot K_{YY}^{-1}, K_{XY} = \underset{P_{XY}}{\mathbb{E}} [XY^T], K_{YY} = \underset{P_{XY}}{\mathbb{E}} [YY^T]$$

- If X, Y are jointly Gaussian $\Rightarrow \hat{X}_{MMSE}(Y) = A^* \cdot Y$