

EM Algorithm

EE698V - Machine Learning for Signal Processing

Vipul Arora



Announcements

- Coding Test
- Coding Tutorial

How to compare two pdf's?

Entropy

X

$x \in \{1, 2, \dots, N\}$

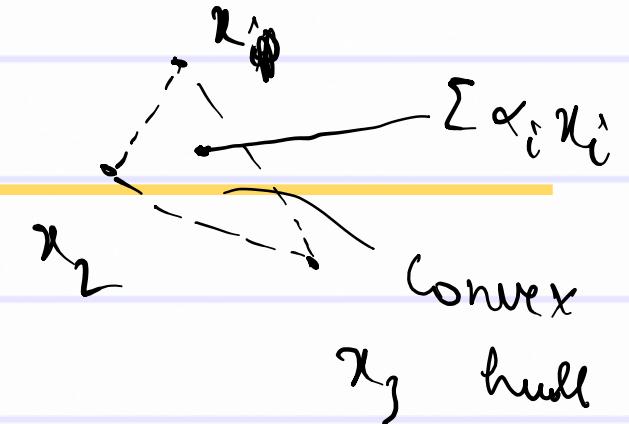
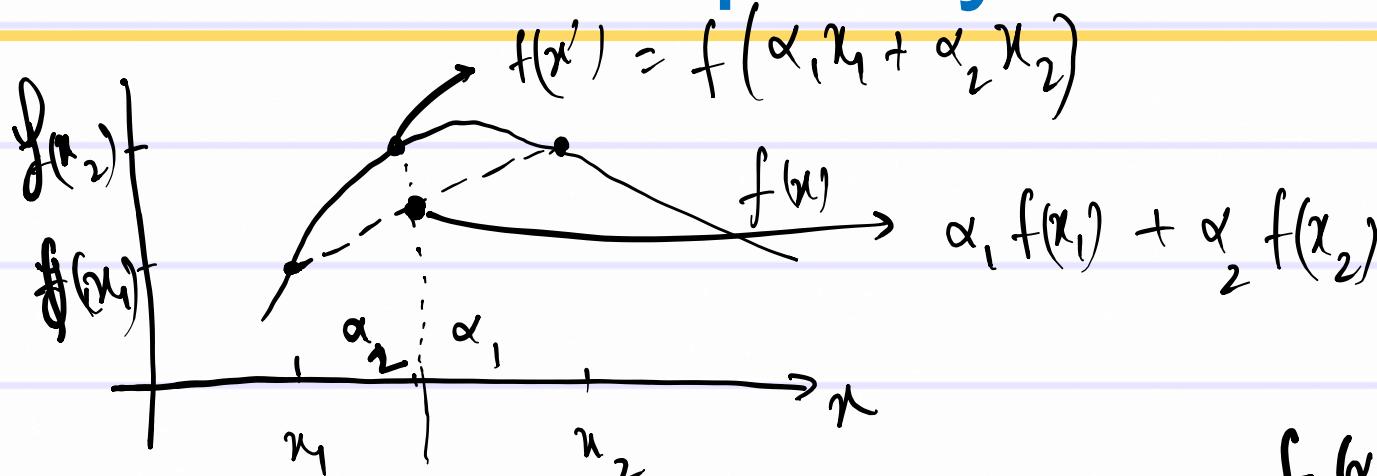
- $H(p(x)) = -\sum_x p(x) \log(p(x))$ for discrete r.v. (Shannon)
- It is always ≥ 0
- When is it minimum?
- Is it upper bounded? When?
- $H(p(x)) = -\int p(x) \log(p(x))dx$ for continuous r.v.
 - Called as continuous entropy
 - It is not lower bounded
 - We will discuss discrete case (Shannon entropy) only

$p(x) \in [0, 1]$ discr.
r.v.

$\log p(x) \in (-\infty, 0]$

- $\log p(x) \in [0, \infty)$

Jensen's Inequality



$$f(\alpha_1 x_1 + \alpha_2 x_2) \geq \alpha_1 f(x_1) + \alpha_2 f(x_2)$$

$$f\left(\sum_i \alpha_i x_i\right) \geq \sum_i \alpha_i f(x_i) \quad \text{if } f'' \leq 0$$

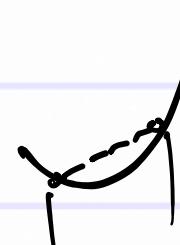
and

$$f\left(\sum_i \alpha_i x_i\right) \leq \sum_i \alpha_i f(x_i) \quad \text{if } f'' \geq 0$$

$$\alpha_i \geq 0 \forall i$$

$$x = \frac{\alpha_1 x_1 + \alpha_2 x_2}{(\alpha_1 + \alpha_2)}$$

let $\sum_i \alpha_i = 1$



$$y' =$$

Lower Bound

\log is $f'' \leq 0$

\therefore

$$\log\left(\sum_i \alpha_i x_i\right) \geq \sum_i \alpha_i \log(x_i)$$

$$H(p(x)) = -\sum_x p(x) \log p(x)$$

$$\underbrace{\alpha_i}_{\alpha_i}$$

$$\alpha_i \geq 0$$

$$\sum_i \alpha_i = 1$$

$$\leq \log \sum_i p_i^2 \cancel{p(x)}$$

$$\geq -\log \sum_i p_i^2$$

$$\boxed{\sum_i p^2(x_i) \leq 1}$$

$$\log(\cancel{\downarrow}) \leq 0$$

$$-\log\left(\sum_i p(x_i)\right) \geq 0$$

$$\Rightarrow H(p(x)) \geq 0$$

$$H(p(x)) = 0, \text{ when } p(x_i) = \sum_i \alpha_i$$

Upper Bound

$$H(p(x)) = -\sum_x p(x) \log p(x)$$
$$= \sum_x p(x) \log \frac{1}{p(x)}$$

Using Jensen's Ineq. $\log \sum_i \alpha_i x_i \geq \sum_i \alpha_i \log x_i$

$$H(p(x)) \leq \log \left(\sum_x p(x) \frac{1}{p(x)} \right)$$
$$= \log \sum_x 1$$
$$= \log N \quad \begin{matrix} N \text{ is the no. of possible} \\ \text{values that r.v. } x \text{ can take} \end{matrix}$$

Cross Entropy between 2 pdf's

- $H(p(x), q(x)) = -\sum_x p(x) \log(q(x))$ for discrete r.v.
- It is always ≥ 0
- When is it minimum?
- Is it upper bounded?
- $H(p(x), q(x)) = -\int p(x) \log(q(x))dx$ for continuous r.v.
- We will study the discrete case only

$$H(p(x), q(x)) \geq 0$$

$= 0$ when $p(x) = q(x)$

No Upper Bound

$$H(p^N, q^N) = - \sum_{i=1}^N p(x_i) \log q(x_i)$$

unif, $p(x_i) = \frac{1}{N} \forall i$

$$= - \sum_i \frac{1}{N} \log q(x_i)$$

for i^* $q(x_{i^*}) \rightarrow 0$

then $H \rightarrow \infty$

KL divergence between 2 pdf's

$$= - \sum p(x) \log q(x) + \sum p(x) \log p(x)$$
$$= H(p(x), q(x)) - H(p(x))$$

- $\text{KL}(p(x)||q(x)) = - \sum_x p(x) \log \frac{q(x)}{p(x)}$ for discrete r.v.
- $\text{KL}(p(x)||q(x)) = - \int p(x) \log \frac{q(x)}{p(x)} dx$ for continuous r.v.
- It is always ≥ 0
- When is it minimum?
- Is it upper bounded? No
- Can you express it in terms of entropy?

if $p(x) = q(x)$, then
 $\text{KL} = 0$

- Can we use $H(p||q)$ as distance metric?
 - only if one is δ
- Can we use KL div. as distance metric?
 - Yes, but one should not be 0 when the other is non-zero.

Mutual Information between two r.v.'s

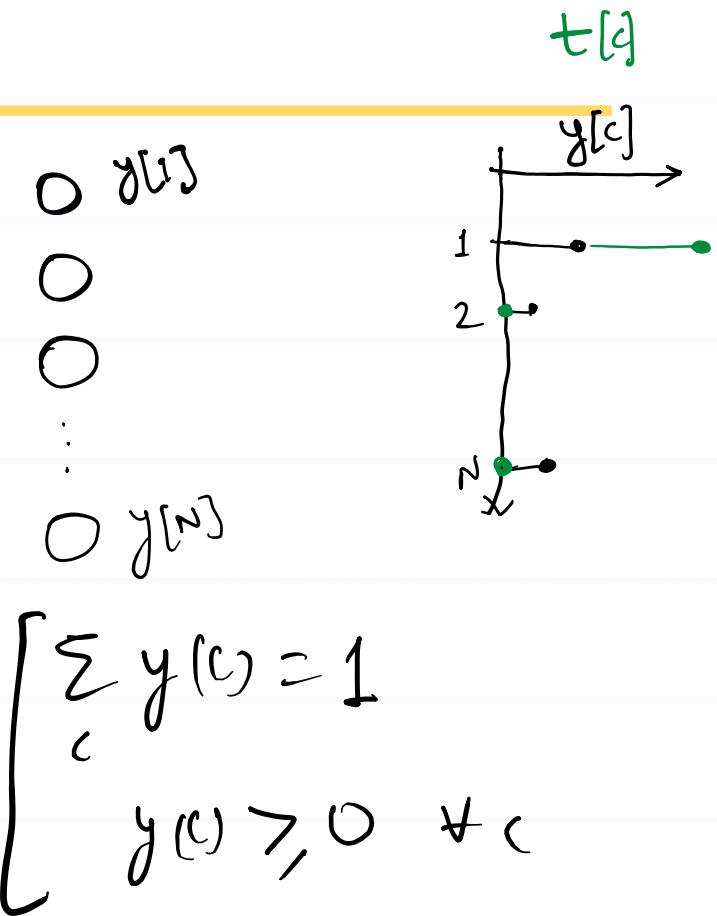
- $I(X; Y) = H(X) - H(X|Y) = H(X) + H(Y) - H(X, Y)$
- Express it in terms of KL divergence
- It is always ≥ 0
- When is it minimum?
- Is it upper bounded?

Relation with KL divergence

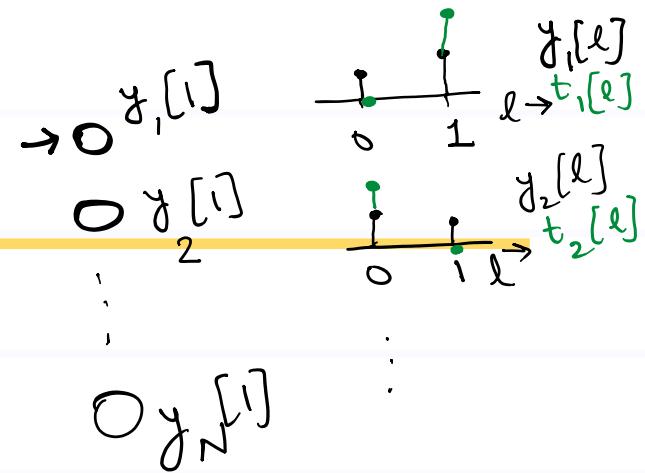
$$\begin{aligned} I(x; y) &= - \sum_x p(x) \log p(x) + \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(x)p(y)} \\ &= - \sum_{x,y} p(x,y) \log p(x) + \sum_{x,y} p(x,y) \log \frac{p(x,y)}{p(y)} \\ &= - \sum_{x,y} p(x,y) \log \left(\frac{p(x)p(y)}{p(x,y)} \right) \\ &= KL[p(x,y) || p(x)p(y)] \end{aligned}$$

Loss Functions

- Multi class classification
 - The random variable is $c \in \{1, \dots, N\}$
 - The ground truth pdf is $t[c]$
 - The estimated pdf is $y[c]$
 - The metric to minimize is $H(t[c], y[c])$



Loss Functions



- Multi-label classification
 - The random variables are $l_i \in \{0,1\}$, for $i \in \{1, \dots, N\}$
 - The ground truth pdf is $t_i[l_i]$
 - The estimated pdf is $y_i[l_i]$
 - The metric to minimize is $H(t_i[l_i], y_i[l_i])$
 - The output neuron may estimate $y_i[l_i = 1]$, the metric becomes
 - $H(t_i[l_i], y_i[l_i]) = -t_i[0] \log y_i[0] - t_i[1] \log y_i[1]$
 - Since, l_i can take only two values, $y_i[l_i = 0] = 1 - y_i[l_i = 1]$

EM Algorithm Optimizing with Auxiliary Functions

Reference: PRML Section 9.4

Optimizing latent variable models

$$p(x) = \sum_k p(z_k) N(x | \mu_k, \sigma_k^2)$$

$$\log p(z|\theta) p(x|z, \theta)$$

- We saw $p(X, Z|\theta)$ has a nice form, i.e., without summation.
Hence, log likelihood is easy to maximize
- But Z is a latent variable, i.e., unobserved. So mostly we get $p(X|\theta)$ which we decompose as $\sum_Z p(X, Z|\theta)$. It is difficult to maximize log likelihood because of summation.
- Can we convert $\log p(X|\theta)$ into $\log p(X, Z|\theta)$?

$$\log \sum$$



Getting rid of summation inside log

$$\bullet \log p(X|\theta) = \log \sum_Z \underbrace{\frac{p(X, Z|\theta)}{q(Z)}}_{\text{v}_i} \times \underbrace{q(Z)}_{\alpha_i}$$

or

$$\bullet \sum_i \log p(x_i|\theta) = \sum_i \log \sum_{z_i} \underbrace{\frac{p(x_i, z_i|\theta)}{q(z_i)}}_{\text{v}_i} \times \underbrace{q(z_i)}_{\alpha_i} \geq \sum_i \sum_{z_i} q(z_i) \log \frac{p(x_i, z_i|\theta)}{q(z_i)}$$

$\log x$

$\log \sum_i \alpha_i v_i \geq \sum_i \alpha_i \log v_i$

Jensen's Ineq.

$L(q(z), \theta)$

$$\sum_i \alpha_i = 1; \alpha_i > 0 \forall i$$

Auxiliary Loss Function

- $\log p(X|\theta) \geq \sum_Z q(Z) \log \frac{p(X, Z|\theta)}{q(Z)} = \mathcal{L}(q(Z), \theta)$

Gap

*q assigned
from data*

- $\log p(X|\theta) - \mathcal{L}(q(Z), \theta) = \text{KL} [q(Z)||p(Z|X, \theta)]$

$$\begin{aligned} \text{LHS} &= \underbrace{\sum_Z q(z) \cdot \log p(x|\theta)}_{=1} - \sum_Z q(z) \log \frac{p(x,z|\theta)}{q(z)} \\ &= \sum_Z q(z) \log \underbrace{\frac{p(x|\theta) \cdot q(z)}{p(x,z|\theta)}}_{\frac{q(z)}{p(z|x,\theta)}} \\ &= \text{RHS} \end{aligned}$$

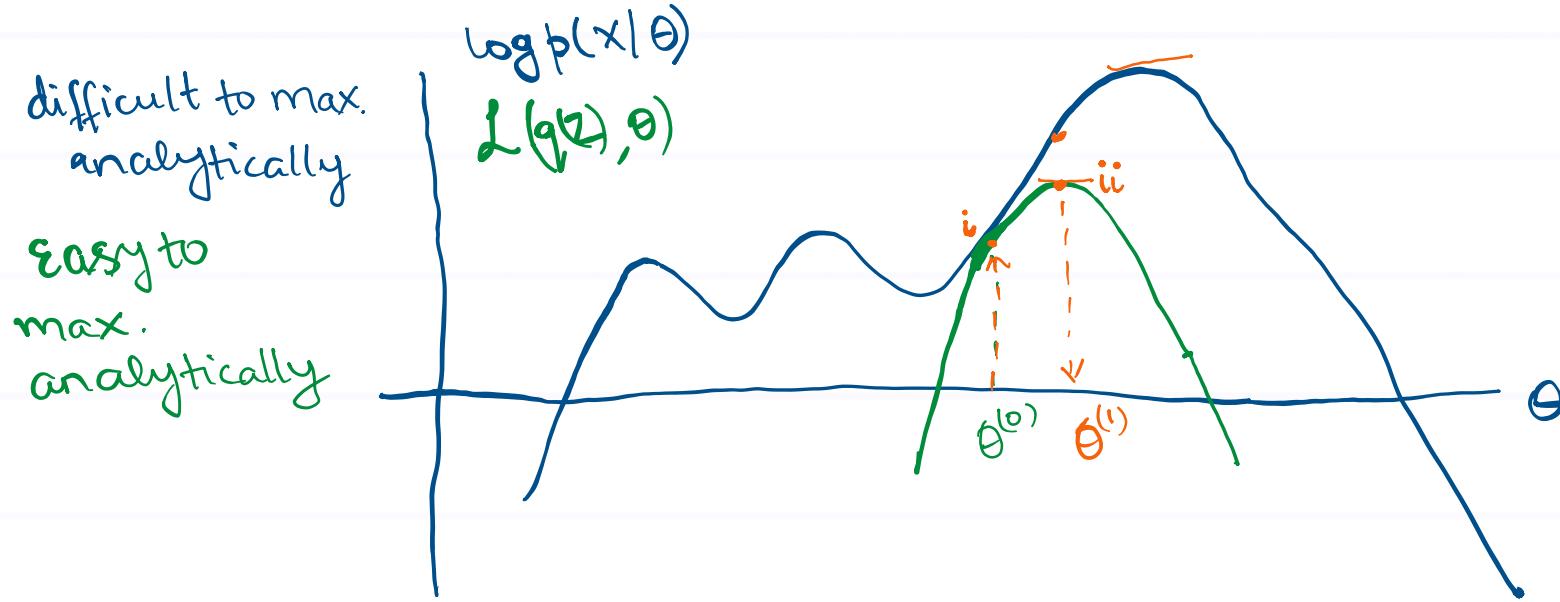
Auxiliary Loss Function

- $\log p(X|\theta) \geq \mathcal{L}(q(Z), \theta)$

max ↗

- $\log p(X|\theta) - \mathcal{L}(q(Z), \theta) = \text{KL} [q(Z)||p(Z|X, \theta)]$

Equal when
 $q(z) = p(z|x, \theta)$



- i. For given $\theta^{(0)}$, choose $q(z) = p(z|x, \theta)$
This gives us \mathcal{L}
- ii. For this \mathcal{L} , find $\max_{\theta} \mathcal{L}$ to get $\theta^{(1)}$.

EM Algorithm

- E-step

- $q \leftarrow \operatorname{argmin}_q \text{KL} [q(Z) || p(Z|X, \theta)]$

$$q(z) = p(z|x, \theta)$$

- M-step

- $\theta \leftarrow \operatorname{argmax}_{\theta} \mathbb{E}_q [\log p(X, Z | \theta)]$

analytically solvable

