Homework-2
- Theory Part.

# Q.1)

Acc. to bayes theorem we know, the probability of the hidden variables over the observations are

$$P(Z/X) = \frac{P(X/Z) \; P(Z)}{P(X)}$$

$$= \frac{P(X/Z) \; P(Z)}{\sum_{Z} P(X,Z)}$$

we can Perive the variational lower bound $L'$ by applying log of the observational probability on the denominator of the $P(Z/X)$

$$\Rightarrow \log \sum_{Z} P(X,Z)$$

As we approximate the $P(Z/X)$ using $q(Z)$ we add it to the equation of the ELBO.

$$\Rightarrow \log \int_{Z} \frac{P(X,Z) \; q(x)}{q(x)}$$

$$\Rightarrow \log \left( E_q \left[ \frac{P(x,z)}{\cancel{E}\; q(z)} \right] \right)$$

The above equation can also be written as:

$$E_q \left[ \log \frac{P(x,z)}{q(z)} \right] \leq \log P(x)$$

$$\Rightarrow E_q \log P(x,z) - E_q \log q(z)$$

which is the variational lower bound.

→ Now we try to prove,

$$KL\left( q(z) \| P(z/x) \right) = - E_q \left[ \log P(x,z) \right]$$

$$+ E_q \log q(z)$$

$$+ \log P(x)$$

We know KL divergence is used to define the similarity between two different distributions in our code ity $P$ & $q$

when we talk about the joint probability $P(z, x)$ to & factorize it.

$$P(z/x) = P(z, x) \, P(z) - \text{(A)}$$

• Next if we apply log and expectations

$E_q \log (P(x/z) P(z))$

$$= E_{qv} [\log P(x/z) + E_{qv} [\log P(z)]$$

- then we apply integral & write it as

$$\Rightarrow \int q(z) \, (\log P(z) - \log q(z)) \, dz$$

$$\Rightarrow \int q(z) \log \frac{P(z)}{q(z)} \, dz.$$

by the help of the above equation we can write the KL divergence as

$$KL (q(z) || P(z/x)) \Rightarrow \int q(z) \log \frac{q(z)}{P(z/x)} \, dz$$

As here our distributions are $q(z)$ & $P(z/x)$

$$= \int q(z)\left(\log q(z) - \log P(z/x)\right) dz$$

$$= \int q(z) \log q(z) - q(z) \log P(z/x) \, dz$$

$$= \int q(z) \log q(z) \, dz - \int q(z) \log P(z/x) \, dz$$

$$= E_q\left(\log q(z)\right) dz - E_q\left(\log P(z/x)\right)$$

$$= E_q\left(\log q(z)\right) - E_q\left(\log \frac{P(x,z)}{P(x)}\right)$$

$\Rightarrow$ we write $P(z/x) = P(z,x) \cdot P(z)$

~~so changed it as PL·~~

$$= E_q\left(\log q(z)\right) - E_q\left[\log P(x,z) - \log P(x)\right]$$

$$= E_q\left(\log q(z) - \log P(x,z)\right) + E_q\left[\log P(x)\right]$$

$$= E_q \log q(z) - \log E_q \log P(x,z) + E_q \log P(x)$$

$\Rightarrow -E_q \log P(x,z) + E_q \log q(z) + E_q \log P(x)$

# Q.2)

The given $\delta(x)$ is the binary classifier with the classification:

$$\delta(x) = \begin{cases} +1 & \text{if } P(y=+1|x=x) \geq P(y=-1|x=x) \\ -1 & \text{otherwise} \end{cases}$$

$t(x)$ is another classifier with an error rate of $R(t)$ which can be defined as

$$R(t) = P(y \neq t(x))$$

which says when the classification $t(x)$ gives is not the same as the expected Label in output. we can write that as

$$R(t) = \cancel{P(x)} E_x\left[P(y \neq t(x)/X=x)\right]$$

Now we need to show that the error rate of the binary classifier $R(\delta)$ is greater than or equal to $R(t)$ which is any binary classifier.

All the features $x \in \mathcal{X}$ y labled $y \in y$ $\mathcal{Y}$
$x, y$ are distributed according to $\rho$.

for an binary classifier
$f: x \to y$ its $0-1$ loss $l(yf(x))$ is

$$l(y, f(x)) = \mathbb{P}(y \neq f(x)) \text{ or}$$

we can say in other words,

$$l(y, f(x)) = \begin{cases} 1 & \text{if } y \neq f(x) \\ 0 & \text{otherwise} \end{cases}$$

the error rate $R(f)$ is defined as

$$R(f) = E\left[l(y, f(x))\right] = \mathbb{P}(y \neq f(x))$$

if $R(\delta) \geq R(f)$ then,

$$\eta(x) \geq 1/2 \iff \frac{\mathbb{P}(y = +1 \mid X = x)}{\mathbb{P}(y = -1 \mid X = x)} \geq 1$$

for any $\delta$

$$\mathbb{P}(y = \delta(x) \mid X = x) = 1 - \left\{ \mathbb{P}(y = 1, \delta(x) = 1 \mid X = 1) \right.$$

$$+ \mathbb{P}(y = -1, g(x) = -1 \mid X = x)$$

$$= 1 - \left\{ E\left[\mathbb{P}(y = 1)\right] \eta(x) + E\left[\delta(x = -1) \right. \right.$$
$$\left. \left. \eta(x)\right]\right\}$$

So,

$$P(Y \neq g^\delta(x) \mid x = x) - P(Y \neq f(x) \mid x = x)$$

$$= \eta(x)\left\{E[\![g f(x) = 1]\!] - E(\![g f(x) = 1]\!]\right\} +$$

$$(1 - \eta(x))\left\{E(\![g f(x) = -1]\!] + E(g(x) = -1)\right\}$$

this after solving

$$\Rightarrow (2\eta(x) - 1)\left\{[E(\![g f(x) = 1]\!]]\right.$$

$$\left. - E[\![g(x) = 1]\!]\right\} \leq 0$$

taking the expectation with respect to $x$
gives

$$R(\delta) - R(f) \leq 0$$
$$R(\delta) \leq R(f)$$

Q.3) Given a set of data with $n$-Paired samples in the form of $\{(x_i, y_i)\}_{i=1}^{n}$ where

$x_i$ is the $d$-dimentional vector of $i^{th}$ sample and $y_i$ is the label of the same $i^{th}$ sample.

- The log likelihood function of the logistic regression.

$$\lambda(\beta) = \sum_{i=1}^{n} y_i \beta^T x_i - \log\left(1 + e^{\beta^T x_i}\right)$$

We need to compute $\dfrac{\partial \lambda(\beta)}{\partial (\beta)}$ &

$$\frac{\partial \lambda(\beta)}{\partial \beta \, \partial \beta^T}$$

First computing $\dfrac{\partial \lambda(\beta)}{\partial (\beta)}$

$$\Rightarrow \sum_{i=1}^{N} y_i x_{iN} - \sum_{i=1}^{N} \frac{x_{ij} e^{\beta^T y_i}}{1 + e^{\beta^T x_i}}$$

$$= \sum_{i=1}^{N} \left( y_i - \frac{e^{\beta^T x_i}}{1 + e^{\beta^T x_i}} \right) x_i$$

$$= \sum_{i=1}^{n} \left( y_i x_{ij} - \frac{x_{ij} e^{\beta^T x_i}}{1 + e^{\beta^T x_i}} \right)$$

$$= \sum_{i=1}^{N} y_i x_i - \sum_{i=1}^{N} \frac{x_{ij} e^{\beta^T x_i}}{1 + e^{\beta^T x_i}}$$

$$= \sum_{i=1}^{N} y_i x_i - \sum_{i=1}^{N} P(x_i; \beta) x_{ij}$$

$$= \sum_{i=1}^{N} x_{ij} \left( y_i - P(x_i; \beta) \right) = 0$$

$$\text{for} \quad \frac{\partial l(\beta)}{\partial \beta_j} = 0, \quad i = 0, 1, \dots, P.$$

For the $2^{nd}$ order derivative

$$\frac{\partial l(\beta)}{\partial \beta \, \partial \beta^T}$$

$$= - \sum_{i=1}^{N} \frac{(1 + e^{\beta^T x_i}) e^{\beta^T x_i} x_i^T - (e^{\beta^T x_i})^2 x_{ij} x_i^T}{(1 + e^{\beta^T x_i})^2}$$

$$= -\sum_{i=1}^{N} x_{i\beta}\, x_{i\beta}^{T}\, P(x_i;\beta) - x_{i\beta}\, x_i^{T} P(x_{i};\beta)^{2}$$

$$= -\sum_{i=1}^{N} x_{i\beta}\, x_i^{T}\, P(x_i;\beta)\left(1 - P(x_i;\beta)\right)$$

$$\Rightarrow \frac{\partial\, l(\beta)}{\partial\beta\,\partial\beta^{T}} = -\sum_{i=1}^{N} x_{i\beta}\, x_i^{T}\, P(x_i;\beta)\left(1 - P(x_i;\beta)\right)$$

we solve thes equations for the value of $\beta$, as with the $1^{st}$ derivative the $\beta$ value is not easy or possible ~~for computing~~ using gradient decent. so we compute $\beta^{1}$

$$\beta^{1} = \beta - \frac{f(\beta)}{f'(\beta)}$$

$\Rightarrow f(\beta)$ is the $1^{st}$ derivative $\dfrac{\partial l(\beta)}{\partial(\beta)}$

and the $f'(\beta)$ is the $2^{nd}$ derivative

$$\frac{\partial l(\beta)}{\partial(\beta)\partial(\beta^{T})}$$