CS280 Fall 2018 Assignment 1 Part A

Name: 于莘

Student ID: 2018233131

① $KL(P_{emp} \| q) = \int P_{emp}(x)(\log P_{emp}(x) - \log q(x;\hat{\theta})) dx$

$$= \int P_{emp}(x) \log P_{emp}(x) dx - \int P_{emp}(x) \log q(x;\hat{\theta}) dx \quad -- ①$$

$\int P_{emp}(x) \log q(x;\hat{\theta}) dx = \int \frac{1}{n} \sum_{i=1}^{n} \delta(x, x_i) \log \ddagger q(x;\hat{\theta}) dx = \frac{1}{n} \sum_{i=1}^{n} \ddagger \log q(x;\hat{\theta}) \int \delta(x, x_i) dx$

$$= \frac{1}{n} \sum_{i=1}^{n} \log(q(x;\theta)) \overset{\triangle}{=} \hat{\theta}$$

So when $\hat{\theta}$ is the maximum likelihood estimator, we get the minimum $KL(P_{emp} \| q)$, which means $\arg\min_q KL(P_{emp} \| q)$ ‡ is obtained by $q(x)$.

② $(J(w))''$ to slove

$J(w) = -\frac{1}{|D|} \sum_{i \in D} \log \sigma(\pm 1 \cdot x_i^T w) + \lambda \|w\|_2^2$, label $\# \log \sigma(y_i x_i^T w)$ as ①

$\dfrac{\partial(\log \sigma(x_i^T w))}{\partial w_j} = \dfrac{-\partial(\log(1+e^{-x_i^T w}))}{\partial w_j} = \left(\dfrac{e^{-x_i^T w}}{1+e^{-x_i^T w}}\right)\left(\dfrac{x_i^T}{x_{ij}}\right) = (1-\sigma(x_i^T w)) x_{ij}$

$\dfrac{\partial^2(\log \sigma(x_i^T w))}{\partial w_j \partial w_k} = \sigma(x_i^T w)(1-\sigma(x_i^T w)) x_{ij} \cdot x_{ik} \quad \dfrac{= \sigma(x_i^T w)(1-\sigma(x_i^T w)) x_i^T \cdot x_i^T}{1}$

∵ $\lambda \|w\|^2$ is convex, $\dfrac{\partial(\log \sigma(y_i x_i^T w))}{\partial w} = \dfrac{\sigma(y_i x_i^T w)(1-\sigma(y_i x_i^T w)) y_i x_i^T}{\sigma(y_i x_i^T w)}$

$$= (1-\sigma(y_i x_i^T w)) y_i x_i^T$$

∴ $\dfrac{\partial((1-\sigma(y_i x_i^T w)) y_i x_i^T)}{\partial w} = -\sigma(y_i x_i^T w)(1-\sigma(y_i x_i^T w)) \cdot (y_i x_i^T)^2$

because $\sigma(x) \in (0,1)$, $\text{①}'' \downarrow 1-\sigma(\frac{1}{2}x_i^T w) \in (0,1)$ <span>PAGE 2</span>

$(\frac{1}{2}x_i^T)^2 > 0$, so $\text{①}''$ always $< 0$

when it add the coefficient $-\frac{1}{|D|}$, it will become always $> 0$

so the first subject is also convex [the hessian matrix $\geq 0$]

convex + convex $\Rightarrow$ convex $\Longrightarrow$ local optimal = global optimal

$\Rightarrow$ False

False    L1 norm prefer a sparse $\hat{w}$, but L2 norm prefer a average value $\hat{w}$

if we consider $J(w)$ as a loss function, when we optimize and get the $\text{argmin}_w J(w) = \hat{w}$, it's not sparse.

③ $L(\theta) = \sum\limits_{i=1}^{m} \log P(x_i;\theta) = \sum\limits_{i=1}^{m} \log \sum\limits_{z} P(x, z; \theta)$

$\qquad = \sum\limits_{i} \log \sum\limits_{z^{(i)}} Q_i(z^{(i)}) \frac{P(x^{(i)}, z^{(i)};\theta)}{Q_i(z^{(i)})}$

$\qquad \geq \sum\limits_{i} \sum\limits_{z^{(i)}} Q_i(z^{(i)}) \log \frac{P(x^{(i)}, z^{(i)};\theta)}{Q_i(z^{(i)})}$

$\qquad = L(\theta^{(t)})$

$w_j^{(i)} = Q_i(z^{(i)} = j) = P(z^{(i)} = j | x^{(i)}; \phi, \mu, \Sigma)$

$\sum\limits_{i=1}^{m} \sum\limits_{z^{(i)}} Q_i(z_*^{(i)}) \log \frac{P(x^{(i)}, z^{(i)}; \phi, \mu, \Sigma)}{Q_i(z^{(i)})}$ ----- I

$= \sum\limits_{i=1}^{m} \sum\limits_{j=1}^{m} w_j^{(i)} \log \frac{\frac{1}{(2\pi)^{\frac{n}{2}}|\Sigma_j|^{\frac{1}{2}}} \exp(-\frac{1}{2}(x^{(i)}-\mu_j)^T \Sigma_j^{-1}(x^{(i)}-\mu_j)) \cdot \phi_j}{w_j^{(i)}}$

$\nabla_{\mu_1} I = -\nabla_{\mu_1} \sum\limits_{i=1}^{m} \sum\limits_{j=1}^{K} w_j^{(i)}(x^{(i)}-\mu_j)^T \cdot \frac{1}{2} \Sigma_j^{-1}(x^{(i)}-\mu_j)$

$\qquad = \frac{1}{2} \sum\limits_{i=1}^{m} w_1^{(i)} \nabla_{\mu_1} 2\mu_1^T \Sigma_1^{-1} x^{(i)} - \mu_1^T \Sigma_1^{-1} \mu_1 = \sum\limits_{i=1}^{m} w_1^{(i)}(\Sigma_1^{-1} x^{(i)} - \Sigma_1^{-1}\mu_1)$

令其为0 得 $\mu_l = \dfrac{\sum_1^m w_l^{(i)} x^{(i)}}{\sum_1^m w_l^{(i)}}$ , $\Rightarrow \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \log \phi_j$

又： $\sum_{j=1}^k \phi_j = 1$

∴ let $L(\phi) = \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} \log \phi_j + \beta(\sum_{j=1}^k \phi_j - 1)$

$\dfrac{\partial L(\phi)}{\partial \phi_j} = \sum_{i=1}^m \dfrac{w_j^{(i)}}{\phi_j} + \beta \overset{\wedge}{=} 0$

$\Rightarrow \phi_j = \dfrac{\sum_{i=1}^m w_j^{(i)}}{-\beta}$ $\Rightarrow -\beta = \sum_{i=1}^m \sum_{j=1}^k w_j^{(i)} = \sum_{i=1}^m 1 = m$

∴ $\phi_j = \dfrac{1}{m} \sum_{i=1}^m w_j^{(i)}$

$\nabla_{\mu_k} \sum_{n=1}^N \sum_{k=1}^k \gamma_{nk} \log \dfrac{\frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma_k|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(x_n - \mu_k)^T \Sigma_k^{-1}(x_n - \mu_k)\right) \pi_k}{\gamma_{nk}}$

$= \left(\sum_{n=1}^N \gamma_{nk}\right)(x_n - \mu_k)\Sigma_k^{-1}$

$= \sum_1^N \gamma_{nk} \Sigma_k^{-1}(x_n - \mu_k)$

❶ $\sum_{k=1}^k \pi_k = 1$ , $\pi$ is the above $\phi$ , we have already proved that

$\pi_k = \dfrac{1}{N}\sum_{n=1}^N \gamma_{nk}$

❷ $\Sigma_k$ can be the $\mu$ above, so we get $\bar\Sigma_k = \dfrac{\sum_{i=1}^m \gamma_{nk} x_i}{\sum_{m=1}^m \gamma_n}$

which is $\dfrac{\partial \ln \ell(X|\mu,\Sigma)}{\partial \mu} = -\dfrac{N}{2}\Sigma^{-1} + \dfrac{1}{2}\Sigma^{-1}\left[\sum_{n=1}^N (x_n - \mu)(x_n - \mu)^T\right]\Sigma^{-1}$

$\overset{\beta}{=} 0$

$\Rightarrow \bar\Sigma_k = \dfrac{1}{N}\sum_{n=1}^N (x_n - \mu)(x_n - \mu)^T$