

Problem 1

Thursday, October 26, 2017 6:25 PM

$$\begin{aligned}
 (a) \quad \nabla_w l(w) &= \nabla_w \left[\sum_{i=1}^n -y_i \log h(x_i) - (1-y_i) \log (1-h(x_i)) \right] \\
 &= \sum_{i=1}^n -y_i \cdot \nabla (\log h(x_i)) - (1-y_i) \nabla \log (1-h(x_i)) \\
 &= \sum_{i=1}^n -y_i \cdot \frac{\nabla h(x_i)}{h(x_i)} - (1-y_i) \frac{\nabla (1-h(x_i))}{1-h(x_i)} \quad (1) \\
 \nabla h(x_i) &= \nabla_w \frac{1}{1 + \exp(-w^T x_i)} = \frac{-1 \cdot (-x_i) \exp(-w^T x_i)}{(1 + \exp(-w^T x_i))^2} \\
 &= \frac{x_i \cdot \exp(-w^T x_i)}{(1 + \exp(-w^T x_i))^2} = h(x_i) \cdot \exp(-w^T x_i) \cdot x_i
 \end{aligned}$$

$$\begin{aligned}
 (1) &= \sum_{i=1}^n -y_i \frac{\nabla h(x_i)}{h(x_i)} + (1-y_i) \frac{\nabla h(x_i)}{1-h(x_i)} \\
 &= \sum_{i=1}^n -y_i \frac{x_i \exp(-w^T x_i)}{1 + \exp(-w^T x_i)} + (1-y_i) \frac{x_i}{1 + \exp(-w^T x_i)} \\
 &= \sum_{i=1}^n -y_i x_i + \frac{x_i}{1 + \exp(-w^T x_i)} \\
 &= \sum_{i=1}^n -x_i (y_i - h(x_i)) \\
 &= -x \cdot (\underline{y} - h(x))
 \end{aligned}$$

where $h(x)$ is a vector of $h(x_i)$, x is $[x_1 \dots x_n] \in \mathbb{R}^{m \times n}$

$$\begin{aligned}
 (b) \quad \nabla^2 l(w) &= \nabla_w [-x^T \cdot (\underline{y} - h(x))] \\
 &= \nabla_w x^T \cdot h(x)
 \end{aligned}$$

$$= X S X^T$$

where S is a diagonal matrix with
diagonal element

$$S_{ii} = h^2(x_i) \exp(-w^T x_i).$$

$$X \in \mathbb{R}^{m \times n}$$

Since $h(x_i) > 0$ and $\exp(-w^T x_i) > 0$,
it is positive semi-definite and thus
 $h(w)$ is convex & has no local minima other
than the global one.

Problem 2

Sunday, October 29, 2017 9:14 PM

$$(a) \quad \log L(\alpha; X) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \alpha)^2$$

$$\frac{\partial \log L(\alpha; \underline{x})}{\partial \alpha} = \frac{1}{\delta^2} \sum_{i=1}^n (x_i - \alpha) = 0$$

$$\Rightarrow \sum_{i=1}^n (x_i - \alpha) = 0$$

$$\hat{\alpha} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\log L(\alpha, \sigma^2; \underline{x}) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \alpha)^2$$

$$\frac{\partial \log L(\alpha, \delta^2; X)}{\partial (\delta^2)} = -\frac{n}{2\delta^2} + \frac{1}{2\delta^4} \sum_{i=1}^n (X_i - \alpha)^2$$

According to previous one for α , the MLE of $[\hat{\alpha}, \hat{\beta}]^\top$ solves the system of equations

$$\frac{1}{\hat{\sigma}^2} \sum_{i=1}^n (X_i - \hat{X}) = 0 \quad (1)$$

$$-\frac{n}{\sum \hat{\beta}^2} + \frac{1}{2(\hat{\beta}^2)^2} \sum_{i=1}^n (x_i - \hat{\alpha})^2 = 0 \quad (2)$$

$$(1) \Rightarrow \hat{\alpha} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2.$$

$$(b) \quad l(\underline{\mu}; \underline{X}_1, \dots, \underline{X}_n) = \frac{1}{(2\pi)^{nP/2} |\Sigma|^{n/2}} e^{-\frac{1}{2} \sum_{k=1}^n (\underline{X}_k - \underline{\mu})^\top \Sigma^{-1} (\underline{X}_k - \underline{\mu})}$$

$$= \frac{1}{(2\pi)^{nP/2} |\Sigma|^{n/2}} e^{-\frac{1}{2} \text{tr}(\Sigma^{-1} S)} e^{-\frac{1}{2} n(\bar{\underline{X}} - \underline{\mu})^\top \Sigma (\bar{\underline{X}} - \underline{\mu})}$$

where $\underline{X}_k \quad k \in \{1, 2, \dots, n\} \in \mathbb{R}^p$, and where

$$S = \sum_{k=1}^n (\underline{X}_k - \bar{\underline{X}}) (\underline{X}_k - \bar{\underline{X}})^\top \quad \text{where } \bar{\underline{X}} \text{ is the sample mean of } \underline{X}_k.$$

$$\text{Then } \log l(\underline{\mu}; \underline{X}_1, \dots, \underline{X}_n) = C - \frac{n}{2} \log |\Sigma| - \frac{1}{2} \text{tr}(\Sigma^{-1} S) - \frac{1}{2} n \cdot (\bar{\underline{X}} - \underline{\mu})^\top \Sigma^{-1} (\bar{\underline{X}} - \underline{\mu})$$

where C does not depend on the parameter $\underline{\mu}$.

$$\frac{\partial \log l(\underline{\mu}; \underline{X}_1, \dots, \underline{X}_n)}{\partial \underline{\mu}} = 0$$

$$\Rightarrow \bar{\underline{X}} - \underline{\mu} = 0$$

$$\therefore \hat{\underline{\mu}} = \bar{\underline{X}} = \frac{1}{n} \sum_{i=1}^n \underline{X}_i$$

Problem 3

Sunday, October 29, 2017 9:41 PM

$$\begin{aligned}
 (a) H(X) - H(X|Y) &= - \int p(x) \ln p(x) dx + \int p(x,y) \ln p(x|y) dx dy \\
 &= - \int p(x) \ln p(x) dx + \int p(x,y) \ln p(x|y) dy dx \\
 &= \int (\int p(x,y) \ln p(x|y) dy) - p(x) \ln p(x) dx. \quad (1)
 \end{aligned}$$

$$\ln p(x|y) = \ln \frac{p(x,y)}{p(y)} = \ln \frac{p(y|x)p(x)}{p(y)}$$

$$= \ln p(y|x) + \ln p(x) - \ln p(y)$$

$$\int p(x,y) \cdot \ln p(x|y) dy$$

$$= \int p(x,y) \ln p(y|x) dy + \int p(x,y) \ln p(x) dy - \int p(x,y) \ln p(y) dy. \quad (2)$$

Plug in (2) into (1) we get

$$\begin{aligned}
 &\int (\int p(x,y) \ln p(y|x)) dy + \int p(x,y) dy \cdot \ln p(x) - \int p(x,y) \ln p(y) dy \\
 &- p(x) \ln p(x) dx \\
 &= -H(Y|X) + \int \int p(x,y) \ln \frac{p(x)}{p(y)} - p(x) \ln p(x) dy dx \quad (3)
 \end{aligned}$$

$$\text{Since } P(X) = \int p(x,y) dy$$

Then the second term of (3) becomes

$$\begin{aligned}
 &\int \int p(x,y) \ln \frac{p(x)}{p(y)} - p(x,y) \ln p(x) dy dx \\
 &= - \int \int p(x,y) \ln p(y) dy dx
 \end{aligned}$$

$$= - \int p(Y) \ln p(Y) dY$$

$$= H(Y)$$

$$\text{Finally, } ① = H(Y) - H(Y|X)$$

$$\text{Therefore, we have } H(X) - H(X|Y) = H(Y) - H(Y|X) = I(X,Y)$$

$$(b) I(X,Y) = H(X) - H(X|Y)$$

$$= - \int p(x) \ln p(x) dx + \int p(x,y) \ln (x|y) dx dy$$

Plug in $X = f(Y)$, we have

$$= - \int p(f(Y)) \ln p(f(Y)) df(Y) + \int p(f(Y), Y) \ln p(f(Y)|Y) df(Y) dY$$

For second term, since $\ln p(f(Y)|Y) = 1$ then it becomes

$$= - \int p(f(Y)) \ln p(f(Y)) df(Y)$$

$$= - \int p(Y) \ln p(Y) dY$$

$$= H(Y)$$

$$\text{Also, since } - \int p(f(Y)) \ln p(f(Y)) df(Y)$$

$$= H(f(Y))$$

$$= H(X)$$

We have $I(X, Y) = H(X) - H(Y)$

$$\begin{aligned}
 (c) \quad \min_{\theta} D_{KL}(\hat{P} \parallel q) &\triangleq \min_{\theta} - \int \hat{P}(x) \ln \frac{q(x|\theta)}{\hat{P}(x)} dx \\
 &= \min_{\theta} - \int \hat{P}(x) \ln q(x|\theta) dx + \int \hat{P}(x) \ln \hat{P}(x) dx \\
 &\propto \min_{\theta} - \int \hat{P}(x) \ln q(x|\theta) dx \\
 &= \min_{\theta} - \int \frac{1}{N} \sum_{i=1}^N I[X=x_i] \ln q(x_i|\theta) dx \\
 &= \min_{\theta} - \frac{1}{N} \sum_{i=1}^N \int \delta(x - x_i) \ln q(x|\theta) dx \\
 &= \min_{\theta} - \frac{1}{N} \sum_{i=1}^N \ln q(x_i|\theta) \\
 &\propto \max_{\theta} \sum_{i=1}^N \ln q(x_i|\theta) \\
 &= \max_{\theta} q(x|\theta) \quad \text{where it is exactly the maximum likelihood} \\
 &\text{estimate } \theta_{ML} \text{ given data } D.
 \end{aligned}$$

(d) To maximize entropy for a continuous variable, we need to constrain the first and second moments of $p(x)$ as well as preserve the normalization constraint in order for this maximum well-defined.

Then the three constraint are the following:

$$\int_{-\infty}^{\infty} p(x) dx = 1$$

$$\int_{-\infty}^{\infty} x p(x) dx = \mu$$

$$\int_{-\infty}^{\infty} (x - \mu)^2 p(x) dx = \sigma^2$$

The objective function is given by

$$\max - \int_{-\infty}^{\infty} p(x) \ln p(x) dx = H(x)$$

Then using Lagrange multipliers to perform the constrained maximization

$$\begin{aligned} - \int_{-\infty}^{\infty} p(x) \ln p(x) dx + \lambda_1 (\int_{-\infty}^{\infty} p(x) dx - 1) \\ + \lambda_2 (\int_{-\infty}^{\infty} x p(x) dx - \mu) \\ + \lambda_3 (\int_{-\infty}^{\infty} (x - \mu)^2 p(x) dx - \sigma^2) = F(p(x)) \end{aligned}$$

$$\text{Let } \frac{\partial F(p(x))}{p(x)} = 0$$

$$\text{we have } p(x) = \exp \{-1 + \lambda_1 + \lambda_2 x + \lambda_3 (x - \mu)^2\}.$$

$$\text{Since } \left\{ \int p(x) dx = 1 \right.$$

$$\left. \int x p(x) dx = \mu \right.$$

$$\left. \int (x - \mu)^2 p(x) dx = \sigma^2 \right.$$

we solve λ_1, λ_2 & λ_3 and plug into $p(x)$ we get

$$\hat{p}(x) = \frac{1}{(2\pi\sigma^2)^{\frac{1}{2}}} \exp \left\{ -\frac{(x - \mu)^2}{2\sigma^2} \right\}$$

where $\hat{p}(x)$ is exactly Gaussian distributed.

Therefore Gaussian distribution has maximum entropy among all distribution

Therefore Gaussian distribution has maximum entropy among all distribution of the same variance.

Problem 4

Wednesday, November 1, 2017 4:53 PM

$$(a) \quad \underline{w} = \arg \min_{\beta, b} \sum_{i=1}^n c_i (y_i - \beta^T \underline{x}_i - b)^2 \\ = \arg \min_{\underline{w}} \|C^{\frac{1}{2}}(\underline{y} - \underline{X}\underline{w})\|_2^2$$

where \underline{C} is diagonal matrix with $C_{ii} = c_i$, and $C = C^{\frac{1}{2}} \cdot C^{\frac{1}{2}}$

$$\frac{\partial L(\underline{w})}{\partial \underline{w}} = -2(C\underline{x})^T C^{\frac{1}{2}}(\underline{y} - \underline{X}\underline{w}) = 0 \\ \Rightarrow (C\underline{x})^T(\underline{y} - \underline{X}\underline{w}) = 0$$

$$\widehat{\underline{w}} = ((C\underline{x})^T \underline{x})^{-1} \cdot (C\underline{x})^T \underline{y} \\ = (\underline{x}^T C \cdot \underline{x})^{-1} \underline{x}^T C \cdot \underline{y}$$

When $c_i = 1 \forall i$, C is identity matrix.

$$\widehat{\underline{w}} = (\underline{x}^T \underline{x})^{-1} \underline{x}^T \underline{y}, \quad \widehat{\underline{w}} = \begin{bmatrix} b \\ \beta \end{bmatrix}$$

For $y_i = \beta^T \underline{x}_i + b + \varepsilon_i$ it can be written as

$$\underline{y} = \underline{X}\underline{w} + \underline{\varepsilon} \quad \text{where } \underline{\varepsilon} | \underline{X} \sim N(0, \sigma^2 I)$$

$$\widehat{\underline{w}}_{ML} = \max_{\underline{w}} p(\underline{y} | \underline{w}, \underline{X})$$

$$\therefore \underline{y} | \underline{w}, \underline{X} \sim N(\underline{X}\underline{w}, \sigma^2 I)$$

$$\begin{aligned}
 & \because \max_{\underline{w}} p(\underline{y} | \underline{x}, \underline{w}) = \max_{\underline{w}} \log p(\underline{y} | \underline{w}) \\
 &= \max_{\underline{w}} -(\underline{y} - \underline{x}\underline{w})^T (\delta^2 I)^{-1} (\underline{y} - \underline{x}\underline{w}) \\
 &= \min_{\underline{w}} \frac{1}{\delta^2} \|\underline{y} - \underline{x}\underline{w}\|_2^2 \\
 &= \min_{\underline{w}} \|\underline{y} - \underline{x}\underline{w}\|_2^2
 \end{aligned}$$

which is a least square regression problem, and therefore

$$\hat{\underline{w}}_{ML} = (\underline{x}^T \underline{x})^{-1} \underline{x}^T \underline{y} = \hat{\underline{w}}_{LS}.$$

\therefore When $c_i = 1 \forall i$, the weighted least squares regression problem is the same as finding the MLE of $\underline{w} = [b, \beta^T]^T$ with y_i being modeled as $y_i = \beta^T x_i + b + \varepsilon_i \forall i$.

(b) Now $\underline{y} | \underline{x}, \underline{w} \sim N(\underline{x}\underline{w}, \Sigma)$, where Σ is diagonal and $\Sigma_{ii} = \delta_i^2$

$$\begin{aligned}
 & \max_{\underline{w}} p(\underline{y} | \underline{x}, \underline{w}) \\
 &= \min_{\underline{w}} (\underline{y} - \underline{x}\underline{w})^T \Sigma^{-1} (\underline{y} - \underline{x}\underline{w}) \\
 &= \min_{\underline{w}} \|\Sigma^{-\frac{1}{2}} (\underline{y} - \underline{x}\underline{w})\|_2^2, \text{ which is exactly weighted LS problem.}
 \end{aligned}$$

According to part a, we have

$$\hat{\underline{w}}_{WLS} = (\underline{x}^T \Sigma^{-1} \underline{x})^{-1} \underline{x}^T \Sigma^{-1} \underline{y}$$

where we change C in part a to Σ^{-1}

where we change $\underline{\underline{C}}$ in part a to $\underline{\underline{\Sigma}}^{-1}$.

∴ The MLE of $\underline{\underline{W}}$ with different noise variance for each i is equivalent to weighted LS solution with weight matrix $\underline{\underline{C}} = \underline{\underline{\Sigma}}^{-1}$.