

SOLUTION

(a) Firstly, we ^{prove} assume $\theta^* = \max_{\theta} L(D, \theta)$ if $\theta^* = \max_{\theta} \ell(D, \theta)$.

- We assume $\theta' = \max_{\theta} \ell(D, \theta)$ and $\theta'' = \max_{\theta} L(D, \theta)$, where $\theta' \neq \theta''$. Because of the assumption, we have

$$\ell(D, \theta') > \ell(D, \theta'')$$

Since exponential function is strictly increasing, we have

$$e^{\ell(D, \theta')} > e^{\ell(D, \theta')}$$

which is equivalent to

$$L(D, \theta') > L(D, \theta'')$$

However, $\theta' \neq \theta''$ and $\theta'' = \max_{\theta} L(D, \theta)$, which implies a contradiction. We can conclude that θ^* solves $\max_{\theta} L(D, \theta)$ if it solves $\max_{\theta} \ell(D, \theta)$.

Second Similarly, we ^{prove} assume $\theta^* = \max_{\theta} \ell(D, \theta)$ if $\theta^* = \max_{\theta} L(D, \theta)$.

- We assume the contrary: $\theta' = \max_{\theta} L(D, \theta)$, $\theta'' = \max_{\theta} \ell(D, \theta)$, where $\theta' \neq \theta''$. Because of the assumption, we have

$$L(D, \theta') > L(D, \theta'')$$

Since the logarithm function is strictly increasing, we have

$$\log L(D, \theta') > \log L(D, \theta'')$$

which is equivalent to

$$\ell(D, \theta') > \ell(D, \theta'')$$

However, $\theta' \neq \theta''$, $\theta'' = \max_{\theta} \ell(D, \theta)$, which implies a contradiction. We can conclude the θ^* solves $\max_{\theta} \ell(D, \theta)$ if $\theta^* = \max_{\theta} L(D, \theta)$.

q.e.d.

(b) Maximising log-likelihood is better because

i) Summing is less expensive than multiplication, for computers and by hand

ii) Likelihood can be very small, which can result in float number underflow. (Sometimes even float 64 is not enough)

$$\begin{aligned}
 (c) \Pr(y=0|x, \theta) &= 1 - \Pr(y=1|x, \theta) = 1 - \frac{1}{1 + \exp(x\theta)} \\
 &= \frac{1 + \exp(x\theta) - 1}{1 + \exp(x\theta)} = \frac{\exp(x\theta)}{1 + \exp(x\theta)}
 \end{aligned}$$

~~$$(d) \log \Pr(y, x, \theta) = y \log \Pr(\theta)$$~~

$$\begin{aligned}
 (d) \log \Pr(y, x, \theta) &= \log p(x, \theta)^y + \log (1 - p(x, \theta))^{1-y} \\
 &= \log p(x, \theta)^y \cdot (1 - p(x, \theta))^{1-y}
 \end{aligned}$$

y is a binary variable, When $y=0$:

~~$$\log \Pr(y, x, \theta) = \log (1 - p(x, \theta))$$~~

When $y=1$

$$\log \Pr(y, x, \theta) = \log p(x, \theta)^y$$

So we write

$$\log \Pr(y, x, \theta) = y \log p(x, \theta)^y + (1-y) \log (1 - p(x, \theta))$$

q.e.d.

~~Natural logarithm.~~

(e) Assume the logarithm function is natural logarithm.

$$\begin{aligned}\nabla_{\theta} \ell(D, \theta) &= \sum_{i=1}^N \left[\nabla_{\theta} y_i \log p(x_i, \theta) + (1-y_i) \nabla_{\theta} \log (1-p(x_i, \theta)) \right] \\ &= \sum_{i=1}^N \left[y_i \nabla_{\theta} \log p(x_i, \theta) + (1-y_i) \nabla_{\theta} \log (1-p(x_i, \theta)) \right]\end{aligned}$$



Since $P(x_i, \theta) = \frac{1}{1 + \exp(x_i \theta)}$, we have $\nabla_{\theta} \log p(x_i, \theta) = \cancel{P(x_i, \theta)} (1 - P(x_i, \theta)) \cdot X_i$
and $\nabla_{\theta} \log (1 - P(x_i, \theta)) = -P(x_i, \theta) \cdot X_i$

~~$$\nabla_{\theta} \ell(D, \theta) = \sum_{i=1}^N \left[y_i \frac{P(x_i, \theta)(1-P(x_i, \theta))}{P(x_i, \theta)} X_i + (1-y_i) \frac{-P(x_i, \theta)}{1-P(x_i, \theta)} X_i \right]$$~~

$$\begin{aligned}\nabla_{\theta} \ell(D, \theta) &= \sum_{i=1}^N \left[y_i (1 - P(x_i, \theta)) X_i + (1-y_i) \cdot (-P(x_i, \theta) X_i) \right] \\ &= \sum_{i=1}^N \left[y_i X_i - y_i P(x_i, \theta) X_i - P(x_i, \theta) X_i + y_i P(x_i, \theta) X_i \right] \\ &= \sum_{i=1}^N \left[y_i X_i - P(x_i, \theta) X_i \right] = \sum_{i=1}^N X_i (y_i - P(x_i, \theta))\end{aligned}$$

Its Hessian matrix is

$$H_{\theta}(x) = \begin{bmatrix} \sum_{i=1}^N X_i (y_i - P(x_i, \theta)) \\ \sum_{i=1}^N X_i (y_i - P(x_i, \theta)) \\ \vdots \\ \sum_{i=1}^N X_i (y_i - P(x_i, \theta)) \end{bmatrix}$$

, where N is the dimension of the parameter vector θ .