

Problem 2

1. (a) 1. Maximize the conditional log-likelihood is easier. Why?  
 0/5 2. Maximize the log of a function and the function itself is the same  
 Not a reason.

5/5 (b) Because log function is a monotonically increasing function

35/40

20/20

$$2. (a) \ell(w) = \log \left[ \prod_{j=1}^n \left( \sum_{c=1}^C I(y_j=c) \frac{e^{w_c^T x_j}}{\sum_{c=1}^C e^{w_c^T x_j}} \right) \right]$$

$$= \sum_{j=1}^n \sum_{c=1}^C \left[ I(y_j=c) w_c^T x_j - \log \sum_{c=1}^C e^{w_c^T x_j} \right] \checkmark$$

20/20

$$(b) \nabla_{w_c} \ell(w) = \sum_{j=1}^n \left( \frac{I(y_j=c)}{\nabla_{w_c} w_c^T x_j} - \nabla_{w_c} \log \sum_{c=1}^C e^{w_c^T x_j} \right)$$

$$= \sum_{j=1}^n \left[ I(y_j=c) x_j - \frac{e^{w_c^T x_j} x_j}{\sum_{c=1}^C e^{w_c^T x_j}} \right]$$

$$= \sum_{j=1}^n x_j [I(y_j=c) - P(y_j=c | x_j, w)] \checkmark$$

20/20

$$(c) \text{ When } c=c, \frac{\partial \ell(w)}{\partial w_c} = \sum_{j=1}^n x_j [I(y_j=c) - P(y_j=c | x_j, w)]$$

$$= \sum_{j=1}^n x_j \left[ I(y_j=c) - \frac{e^{w_c^T x_j}}{\sum_{c=1}^C e^{w_c^T x_j}} \right]$$

$$= \sum_{j=1}^n x_j (x_j)^T \left[ \left( \frac{e^{w_c^T x_j}}{\sum_{c=1}^C e^{w_c^T x_j}} \right)^2 - \left( \frac{e^{w_c^T x_j}}{\sum_{c=1}^C e^{w_c^T x_j}} \right) \right]$$

$$= \sum_{j=1}^n P(y_j=c | x_j, w) [P(y_j=c | x_j, w) - 1] x_j (w_j)^T \checkmark$$

when  $c' \neq c$

$$\begin{aligned}\frac{\partial^2 J(w)}{\partial w_c \partial w_{c'}} &= \left\{ \sum_{j=1}^n x_j \left[ I(y_j = c) - P(y_j = c | x_j, w) \right] \right\}' \\&= \sum_{j=1}^n x_j \left[ I(y_j = c) - \frac{e^{w_c^T x_j}}{\sum_{c'=1}^C e^{w_{c'}^T x_j}} \right]' \\&= \sum_{j=1}^n x_j (x_j)^T \cdot \frac{e^{w_c^T x_j} \cdot e^{w_{c'}^T x_j}}{\left( \sum_{c'=1}^C e^{w_{c'}^T x_j} \right)^2} \\&= \sum_{j=1}^n (P(y_j = c | x_j, w) \cdot P(y_j = c' | x_j, w) \cdot x_j (x_j)^T)\end{aligned}$$
