

$$③ \text{ Given (4.5): } l(\theta) = \prod_{i: y_i=1} p(x_i) \prod_{i': y_i'=0} (1-p(x_i'))$$

where:

(i)  $\theta$  = vector of parameters on which  $p$  depends, which is  $(\beta_0, \beta_1)$  in above equation in the text

(ii)  $p(x) = (e^{\beta_0 + \beta_1 x}) / (1 + e^{\beta_0 + \beta_1 x})$  for the binomial likelihood in the text, although we could generalize it to have  $x_1, \dots, x_K$  instead of just  $x$  easily

(iii) Per above setup, we have  $M=2$  classes, thus

$y_i$  takes value in  $\{0, 1\}$  in this setup (i.e.

this is dummy encoding). Let  $m$  denote class and thus

(iv) we have  $n$  samples

$m \in \{0, 1\}$

Taking the negative of the log of (4.5) yields:

$$\begin{aligned} -\log(l(\theta)) &= -\log \left( \prod_{i: y_i=1} p(x_i) \cdot \prod_{i': y_i'=0} (1-p(x_i')) \right) \\ &= - \left[ \log \left( \prod_{i: y_i=1} p(x_i) \right) + \log \left( \prod_{i': y_i'=0} (1-p(x_i')) \right) \right] \quad \text{property of logs} \\ \text{of the target} &= - \left[ \sum_{i: y_i=1} \log(p(x_i)) + \sum_{i': y_i'=0} \log(1-p(x_i')) \right] \quad \text{property of logs} \end{aligned}$$

We note the above encoding is "binary encoding", and (10.14) and thus the below will assume "one-hot" encoding. Thus we now index each actual target as  $y_{im} = \underset{\text{value of target}}{\text{target for } i^{\text{th}} \text{ sample in the column}}$  for class  $m$ . Thus, given the way we've split the sums, each  $y_{im}=1$ , and thus we can write:

↓ BACK ↓

$$-\log(\ell(\theta)) = - \left[ \sum_{i:y_i=1} \log(p(x_i))^{y_{ii}} + \sum_{i':y_{i'}=0} \log(1-p(x_{i'}))^{y_{i'0}} \right]$$

Further, to align our notation with our new one-hot encoding, we write:

$$f_1(x_i) = p(x_i)$$

$$f_0(x_i) = 1 - p(x_i)$$

Thus,

$$\begin{aligned} -\log(\ell(\theta)) &= - \left[ \sum_{i:y_i=1} \log(f_1(x_i))^{y_{ii}} + \sum_{i':y_{i'}=0} \log(f_0(x_{i'}))^{y_{i'0}} \right] \\ &= - \left[ \sum_{i:y_i=1} y_{ii} \cdot \log(f_1(x_i)) + \sum_{i':y_{i'}=0} y_{i'0} \cdot \log(f_0(x_{i'})) \right] \end{aligned}$$

property  
of  
 $\log$

Lastly, we can now group these two sums into one sum by realizing that:

(i) There are  $n$  total samples

(ii) Given our one-hot encoding,  $y_{ii}$  will equal 1 when  $y_i=1$

and 0 when  $y_i=0$ , and  $y_{i'0}$  will equal 1 when  $y_{i'}=0$  and 0 otherwise, (for the  $i^{\text{th}}$  sample)

Thus,

$$-\log(\ell(\theta)) = - \sum_{i=1}^n \sum_{m=0}^1 y_{im} \cdot \log(f_m(x_i))$$

(10.14) when  $M=2$

To clarify, since exactly half of the  $y_{im}$ 's are 0 (since each sample will have one 1 and one 0), half of these terms will be zero.

which demonstrates that the negative log of the likelihood expression (4.5) is equivalent to the negative multinomial log-likelihood (10.14) when there are  $M=2$  classes.

