

Solution to small exercise (Lecture 2, Day 2)

Helene Charlotte Wiese Rytgaard

September 24, 2021

Define:

(i) Log-likelihood loss function:

$$\mathcal{L}(f)(O) = -(Y \log(f(A, X)) + (1 - Y) \log(1 - f(A, X))).$$

(ii) Logistic regression model:

$$f_\varepsilon(A, X) = \text{expit}(\text{logit}(f(A, X)) + \varepsilon H(A, X)),$$

with the so-called "clever covariate",

$$H(A, X) := \frac{2A - 1}{\pi(A | X)},$$

We verify that (i)–(ii) fulfill

$$(1) \quad f_{\varepsilon=0} = f \quad (2) \quad \left. \frac{d}{d\varepsilon} \right|_{\varepsilon=0} \mathcal{L}(f_\varepsilon)(O) = \phi^*(f, \pi)(O).$$

(1) follows immediately. For (2), first recall that

$$\text{expit}(x) = \frac{e^x}{1 + e^x} = \frac{1}{1 + e^{-x}}, \quad \text{logit}(x) = \log\left(\frac{x}{1-x}\right),$$

such that

$$\text{expit}(\text{logit}(x)) = x,$$

and

$$\begin{aligned} \text{expit}(-\text{logit}(x)) &= \text{expit}(-\log(x) + \log(1-x)) \\ &= \text{expit}(\log(1-x) - \log(x)) \\ &= \text{expit}(\text{logit}(1-x)) = 1-x. \end{aligned}$$

Furthermore, it is easily seen that

$$\begin{aligned} \log(\text{expit}(x)) &= -\log(1 + e^{-x}), \\ \log(1 - \text{expit}(x)) &= -\log(1 + e^x), \end{aligned}$$

so that

$$\begin{aligned}\frac{d}{dx} \log(\text{expit}(x)) &= \frac{e^{-x}}{1 + e^{-x}} = \text{expit}(-x), \\ \frac{d}{dx} \log(1 - \text{expit}(x)) &= -\frac{e^x}{1 + e^x} = -\text{expit}(x).\end{aligned}$$

Now we can differentiate the composite functions

$$\begin{aligned}\frac{d}{d\varepsilon} \log(\text{expit}(\text{logit}(x) + \varepsilon h)) &= h \text{expit}(-(\text{logit}(x) + \varepsilon h)), \\ \frac{d}{d\varepsilon} \log(1 - \text{expit}(\text{logit}(x) + \varepsilon h)) &= -h \text{expit}(\text{logit}(x) + \varepsilon h),\end{aligned}$$

where setting $\varepsilon = 0$ gives

$$\begin{aligned}\left. \frac{d}{d\varepsilon} \log(\text{expit}(\text{logit}(x) + \varepsilon h)) \right|_{\varepsilon=0} &= h \text{expit}(-(\text{logit}(x))) = h(1 - x), \\ \left. \frac{d}{d\varepsilon} \log(1 - \text{expit}(\text{logit}(x) + \varepsilon h)) \right|_{\varepsilon=0} &= -h \text{expit}(\text{logit}(x)) = hx.\end{aligned}$$

Applying these steps to $\mathcal{L}(f_\varepsilon)$ now gives:

$$\begin{aligned}\left. \frac{d}{d\varepsilon} \right|_{\varepsilon=0} \mathcal{L}(f_\varepsilon)(O) &= YH(A, X)(1 - f(A, X)) - (1 - Y)H(A, X)f(A, X) \\ &= YH(A, X) - H(A, X)f(A, X) = H(A, X)(Y - f(A, X)).\end{aligned}$$