

A6(a)

We have regularized negative log-likelihood function:

$$J(w, b) = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i(b + x_i^T w))) + \lambda \|w\|_2^2.$$

Gradient with respect to w

For a single example (x_i, y_i) , define:

$$\ell_i(w, b) = \log(1 + \exp(-y_i(b + x_i^T w))).$$

Let

$$z_i = -y_i(b + x_i^T w).$$

Then

$$\ell_i(w, b) = \log(1 + \exp(z_i)).$$

Step 1: Differentiate $\log(1 + \exp(z_i))$ w.r.t. w Using the chain rule:

$$\frac{d}{dw} \log(1 + \exp(z_i)) = \frac{1}{1 + \exp(z_i)} \cdot \frac{d}{dw} \exp(z_i).$$

Since

$$z_i = -y_i(b + x_i^T w),$$

we have

$$\frac{dz_i}{dw} = -y_i \frac{d}{dw} (b + x_i^T w) = -y_i x_i.$$

Thus,

$$\frac{d}{dw} \exp(z_i) = \exp(z_i) (-y_i x_i).$$

Substituting back:

$$\frac{d}{dw} \log(1 + \exp(z_i)) = \frac{1}{1 + \exp(z_i)} \exp(z_i) (-y_i x_i).$$

Rewriting in terms of the original variable,

$$\frac{d}{dw} \log(1 + \exp(-y_i(b + x_i^T w))) = -y_i x_i \frac{\exp(-y_i(b + x_i^T w))}{1 + \exp(-y_i(b + x_i^T w))}.$$

Step 2: We observe that

$$\frac{\exp(-y_i(b + x_i^T w))}{1 + \exp(-y_i(b + x_i^T w))} = \frac{(1 + \exp(-y_i(b + x_i^T w))) - 1}{1 + \exp(-y_i(b + x_i^T w))} = 1 - \frac{1}{1 + \exp(-y_i(b + x_i^T w))}.$$

Since $\mu_i(w, b)$ is defined as

$$\mu_i(w, b) = \frac{1}{1 + \exp(-y_i(b + x_i^T w))},$$

it follows that

$$\frac{\exp(-y_i(b + x_i^T w))}{1 + \exp(-y_i(b + x_i^T w))} = 1 - \mu_i(w, b).$$

Hence,

$$-y_i x_i \frac{\exp(-y_i(b + x_i^T w))}{1 + \exp(-y_i(b + x_i^T w))} = -y_i x_i [1 - \mu_i(w, b)].$$

Step 3: Sum over all i and include regularization The total cost is

$$J(w, b) = \frac{1}{n} \sum_{i=1}^n \ell_i(w, b) + \lambda \|w\|_2^2.$$

Therefore,

$$\frac{\partial}{\partial w} \left[\frac{1}{n} \sum_{i=1}^n \ell_i(w, b) \right] = \frac{1}{n} \sum_{i=1}^n \left[-y_i x_i (1 - \mu_i(w, b)) \right].$$

The derivative of $\lambda \|w\|_2^2$ with respect to w is $2\lambda w$. Combining these gives

$$\boxed{\nabla_w J(w, b) = \frac{1}{n} \sum_{i=1}^n \left[-y_i x_i (1 - \mu_i(w, b)) \right] + 2\lambda w.}$$

■

Gradient with respect to b

Step 1: Differentiate w.r.t. b Again let

$$z_i = -y_i(b + x_i^T w).$$

Then

$$\frac{dz_i}{db} = -y_i,$$

which implies

$$\frac{d}{db} \exp(z_i) = \exp(z_i) (-y_i).$$

Thus,

$$\frac{d}{db} \log(1 + \exp(z_i)) = \frac{1}{1 + \exp(z_i)} \exp(z_i) (-y_i).$$

Rewriting in original variables,

$$\frac{d}{db} \log\left(1 + \exp(-y_i(b + x_i^T w))\right) = -y_i \frac{\exp(-y_i(b + x_i^T w))}{1 + \exp(-y_i(b + x_i^T w))}.$$

Using the same observation as before,

$$= -y_i [1 - \mu_i(w, b)].$$

Step 2: Sum over i and simplify Hence,

$$\frac{\partial}{\partial b} \left[\frac{1}{n} \sum_{i=1}^n \ell_i(w, b) \right] = -\frac{1}{n} \sum_{i=1}^n y_i [1 - \mu_i(w, b)].$$

Since $\lambda \|w\|_2^2$ does not involve b ,

$$\nabla_b J(w, b) = -\frac{1}{n} \sum_{i=1}^n y_i [1 - \mu_i(w, b)].$$

Moving the negative sign inside the summation,

$$\boxed{\nabla_b J(w, b) = \frac{1}{n} \sum_{i=1}^n (\mu_i(w, b) - 1) y_i.}$$

■