

### 1) Sigmoid function derivatives $\sigma(\eta)$

The sigmoid function is written as  $\sigma(\eta) = \frac{1}{1+e^{-\eta}} = \frac{e^\eta}{1+e^\eta}$ , where  $0 < \sigma(\eta) < 1$ .  
Show that  $\frac{d\sigma(\eta)}{d\eta} = \sigma(\eta)[1 - \sigma(\eta)]$  and  $\frac{d\log\sigma(\eta)}{d\eta} = 1 - \sigma(\eta)$ .

Solution

$$\begin{aligned}\sigma(\eta) &= \frac{1}{1+e^{-\eta}} = \frac{e^\eta}{1+e^\eta}, \quad 0 < \sigma(\eta) < 1 \\ \frac{d\sigma(\eta)}{d\eta} &= -\frac{-e^{-\eta}}{(1+e^{-\eta})^2} = \frac{e^{-\eta}}{(1+e^{-\eta})^2} = \frac{1}{1+e^{-\eta}} \left( \frac{e^{-\eta}}{1+e^{-\eta}} \right) = \frac{1}{1+e^{-\eta}} \left( 1 - \frac{1}{1+e^{-\eta}} \right) = \sigma(\eta)[1 - \sigma(\eta)] \\ \frac{d\log\sigma(\eta)}{d\eta} &= \frac{1}{\sigma(\eta)} \cdot \frac{d\sigma(\eta)}{d\eta} = 1 - \sigma(\eta)\end{aligned}$$

### 2) Logistic Regression Likelihood & Cross-Entropy

Let  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ , where  $\mathbf{x}_n \in \mathbb{R}^D$  and  $y_n \in \mathbb{R}$ , be the training data of a binary logistic regression model with weights  $\mathbf{w} \in \mathbb{R}^D$ . The probability of sample  $(\mathbf{x}_n, y_n)$  belonging to class 1 is  $p(y=1|\mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x})$ , while the probability of belonging to class 0 is  $p(y=0|\mathbf{x}, \mathbf{w}) = 1 - \sigma(\mathbf{w}^T \mathbf{x})$ . Compute the likelihood  $\mathcal{L}(\mathcal{D}|\mathbf{w})$  of data  $\mathcal{D}$  given the model parameters  $\mathbf{w}$ , as well as the cross-entropy error  $\mathcal{E}(\mathbf{w}) = -\log\mathcal{L}(\mathcal{D}|\mathbf{w})$ .

Solution

Input:  $\mathbf{x} \in \mathbb{R}^D$

Output:  $y \in \{0, 1\}$

Training data:  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$

Model:

$$\begin{aligned}p(y=1|\mathbf{x}, \mathbf{w}) &= \sigma(\mathbf{w}^T \mathbf{x}) \\ p(y=0|\mathbf{x}, \mathbf{w}) &= 1 - \sigma(\mathbf{w}^T \mathbf{x}), \quad \sigma(\eta) = \frac{1}{1+e^{-\eta}} \\ f(\mathbf{x}) : \mathbf{x} \rightarrow y, \quad f(\mathbf{x}) &= \begin{cases} 1, & p(y=1|\mathbf{x}, \mathbf{w}) > 0.5 \\ 0, & \text{otherwise} \end{cases} = \begin{cases} 1, & \sigma(\mathbf{w}^T \mathbf{x}) > 0.5 \\ 0, & \text{otherwise} \end{cases}\end{aligned}$$

Model parameters: Weights  $\mathbf{w} \in \mathbb{R}^D$  (to be learned)

Data likelihood for 1 training sample:

$$p(y_n|\mathbf{x}_n, \mathbf{w}) = \begin{cases} \sigma(\mathbf{w}^T \mathbf{x}_n), & y_n = 1 \\ 1 - \sigma(\mathbf{w}^T \mathbf{x}_n), & y_n = 0 \end{cases} = [\sigma(\mathbf{w}^T \mathbf{x}_n)]^{y_n} [1 - \sigma(\mathbf{w}^T \mathbf{x}_n)]^{1-y_n}$$

Data likelihood for all training data:

$$L(\mathcal{D}|\mathbf{w}) = \prod_{n=1}^N p(y_n|\mathbf{x}_n, \mathbf{w}) = \prod_{n=1}^N [\sigma(\mathbf{w}^T \mathbf{x}_n)]^{y_n} [1 - \sigma(\mathbf{w}^T \mathbf{x}_n)]^{1-y_n}$$

Log-likelihood for all training data:

$$l(\mathcal{D}|\mathbf{w}) = \sum_{n=1}^N \left\{ y_n \log [\sigma(\mathbf{w}^T \mathbf{x}_n)] + (1 - y_n) \log [1 - \sigma(\mathbf{w}^T \mathbf{x}_n)] \right\}$$

Cross-entropy error (negative log-likelihood):

$$\mathcal{E}(\mathbf{w}) = - \sum_{n=1}^N \left\{ y_n \log [\sigma(\mathbf{w}^T \mathbf{x}_n)] + (1 - y_n) \log [1 - \sigma(\mathbf{w}^T \mathbf{x}_n)] \right\}$$

## Logistic Regression - Optimization

**3a)** Show that the first order derivative (i.e., gradient vector) of the cross-entropy function is

$$\nabla \mathcal{E}(\mathbf{w}) = \frac{\partial \mathcal{E}(\mathbf{w})}{\partial \mathbf{w}} = \sum_{n=1}^N \underbrace{(\sigma(\mathbf{w}^T \mathbf{x}_n) - y_n)}_{\text{error}} \mathbf{x}_n$$

Solution

We apply the chain rule for each of the terms of the  $\nabla \mathcal{E}(\mathbf{w})$  sum.

$$\begin{aligned} \nabla \mathcal{E}(\mathbf{w}) &= - \sum_{n=1}^N \left\{ y_n \frac{1}{\sigma(\mathbf{w}^T \mathbf{x}_n)} \sigma'(\mathbf{w}^T \mathbf{x}_n) [1 - \sigma(\mathbf{w}^T \mathbf{x}_n)] \mathbf{x}_n + (1 - y_n) \frac{1}{1 - \sigma(\mathbf{w}^T \mathbf{x}_n)} [1 - \sigma(\mathbf{w}^T \mathbf{x}_n)] [1 - (1 - \sigma(\mathbf{w}^T \mathbf{x}_n))] (-1) \mathbf{x}_n \right\} \\ &= - \sum_{n=1}^N \left\{ y_n [1 - \sigma(\mathbf{w}^T \mathbf{x}_n)] \mathbf{x}_n - (1 - y_n) [1 - (1 - \sigma(\mathbf{w}^T \mathbf{x}_n))] \mathbf{x}_n \right\} \\ &= - \sum_{n=1}^N \left\{ y_n [1 - \sigma(\mathbf{w}^T \mathbf{x}_n)] \mathbf{x}_n - (1 - y_n) \sigma(\mathbf{w}^T \mathbf{x}_n) \mathbf{x}_n \right\} \\ &= - \sum_{n=1}^N [y_n - y_n \sigma(\mathbf{w}^T \mathbf{x}_n) - \sigma(\mathbf{w}^T \mathbf{x}_n) + y_n \sigma(\mathbf{w}^T \mathbf{x}_n)] \mathbf{x}_n \\ &= \sum_{n=1}^N \underbrace{(\sigma(\mathbf{w}^T \mathbf{x}_n) - y_n)}_{\text{error}} \mathbf{x}_n \end{aligned}$$

No closed-form solution that minimizes the cross-entropy function.

We use an approximate method, e.g. gradient descent, so we need to compute  $\nabla \mathcal{E}(\mathbf{w})$ . Gradient descent update:  $\mathbf{w}_{k+1} := \mathbf{w}_k - \alpha(k) \nabla \mathcal{E}(\mathbf{w})$

**3b)** Show that the Hessian of the cross-entropy function is  $\mathbf{H} = \frac{\partial^2 \mathcal{E}(\mathbf{w})}{\partial^2 \mathbf{w}} = \nabla \left( (\nabla \mathcal{E}(\mathbf{w}))^T \right) = \sum_{n=1}^N \sigma(\mathbf{w}^T \mathbf{x}_n) \cdot (1 - \sigma(\mathbf{w}^T \mathbf{x}_n)) \cdot (\mathbf{x}_n \cdot \mathbf{x}_n^T)$  and show that it is positive semi-definite.

Solution

$$\begin{aligned} \mathbf{H} &= \frac{\partial^2 \mathcal{E}(\mathbf{w})}{\partial^2 \mathbf{w}} = \nabla \left( (\nabla \mathcal{E}(\mathbf{w}))^T \right) = \nabla \left( \sum_{n=1}^N (\sigma(\mathbf{w}^T \mathbf{x}_n) - y_n) \mathbf{x}_n^T \right) \\ \mathbf{H} &= \frac{\partial}{\partial \mathbf{w}} \left[ \sum_{n=1}^N (\sigma(\mathbf{w}^T \mathbf{x}_n) \cdot \mathbf{x}_n^T - y_n \mathbf{x}_n^T) \right] \\ &= \sum_{n=1}^N \frac{\partial}{\partial \mathbf{w}} [\sigma(\mathbf{w}^T \mathbf{x}_n)] \cdot \mathbf{x}_n^T \quad (\text{chain rule}) \\ &= \sum_{n=1}^N \underbrace{\sigma(\mathbf{w}^T \mathbf{x}_n)}_{\in [0,1]} \cdot \underbrace{(1 - \sigma(\mathbf{w}^T \mathbf{x}_n))}_{\in [0,1]} \cdot \underbrace{(\mathbf{x}_n \cdot \mathbf{x}_n^T)}_{\in \mathbb{R}^{D \times D}} \end{aligned}$$

For all  $\mathbf{v} \in \mathbb{R}^D$ , substituting  $\mu_n = \sigma(\mathbf{w}^T \mathbf{x}_n) (1 - \sigma(\mathbf{w}^T \mathbf{x}_n)) \geq 0$ , we have:

$$\mathbf{v}^T \mathbf{H} \mathbf{v} = \mathbf{v}^T \left( \sum_{n=1}^N \mu_n \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{v} = \sum_{n=1}^N (\mu_n \mathbf{x}_n^T \mathbf{v})^T (\mathbf{x}_n^T \mathbf{v}) = \sum_{n=1}^N \mu_n \|\mathbf{x}_n^T \mathbf{v}\|_2^2 \geq 0$$