

**Solution 1:**

- (a) For a binary classification problem the model can be written as:

$$\hat{y} = 1 \Leftrightarrow \pi(\mathbf{x}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \mathbf{x})} \geq a$$

This can be reformulated, s.t. for  $a \in (0, 1)$

$$\begin{aligned} \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \mathbf{x})} &\geq a \\ \Leftrightarrow 1 + \exp(-\boldsymbol{\theta}^T \mathbf{x}) &\leq a^{-1} \\ \Leftrightarrow \exp(-\boldsymbol{\theta}^T \mathbf{x}) &\leq a^{-1} - 1 \\ \Leftrightarrow -\boldsymbol{\theta}^T \mathbf{x} &\leq \log(a^{-1} - 1) \\ \Leftrightarrow \boldsymbol{\theta}^T \mathbf{x} &\geq -\log(a^{-1} - 1) \end{aligned}$$

For  $a = 0.5$  we get:

$$\hat{y} = 1 \Leftrightarrow \boldsymbol{\theta}^T \mathbf{x} \geq -\log(0.5^{-1} - 1) = -\log(2 - 1) = -\log(1) = 0$$

This means the linear decision boundary is defined by a hyperplane equation, i.e.,  $\boldsymbol{\theta}^T \mathbf{x} = 0$  and it divides the input space in the "1"-space ( $\boldsymbol{\theta}^T \mathbf{x} \geq 0$ ) and in the "0"-space ( $\boldsymbol{\theta}^T \mathbf{x} < 0$ ).

- (b) When the threshold  $a = 0.5$  is chosen, the losses of missclassified observations, i.e.,  $L(\hat{y} = 0|y = 1)$  and  $L(\hat{y} = 1|y = 0)$ , are weighted equally. This means  $a = 0.5$  is a sensible threshold if one does not need to avoid one type of misclassification more than the other.

Intuitively it makes sense to cut off at 0.5 because, if the probability for 1 is closer to 1 than to 0, we would intuitively choose 1 rather than 0.

**Solution 2:**

$$\text{a) } \pi_1(x) = \frac{\exp(\theta_1^T x)}{\exp(\theta_1^T x) + \exp(\theta_2^T x)}$$

$$\pi_2(x) = \frac{\exp(\theta_2^T x)}{\exp(\theta_1^T x) + \exp(\theta_2^T x)}$$

$$\pi_1(x) = \frac{1}{(\exp(\theta_1^T x) + \exp(\theta_2^T x)) / \exp(\theta_1^T x)} = \frac{1}{1 + \exp(-\theta^T x)} \text{ where } \theta = \theta_1 - \theta_2 \text{ and } \pi_2(x) = 1 - \pi_1(x)$$

- b) When using softmax regression the posterior class probability for the class  $k$  is modeled, s.t.

$$\pi_k(x) = \frac{\exp(\theta_k^T x)}{\sum_{j=1}^g \exp(\theta_j^T x)}.$$

A single observation is multinomially distributed, i.e.,

$$\mathcal{L}_i = \mathbb{P}(Y^{(i)} = y^{(i)} | x^{(i)}, \theta_1, \dots, \theta_g) = \prod_{j=1}^g \pi_j(x^{(i)})^{\mathbb{1}_{\{y^{(i)}=j\}}}.$$

Under the assumption that the observations are conditionally independent the likelihood of the data can be expressed, s.t.

$$\mathcal{L} = \mathbb{P}(Y^{(1)} = y^{(1)}, \dots, Y^{(n)} = y^{(n)} | x^{(1)}, \dots, x^{(n)}, \theta_1, \dots, \theta_g) = \prod_{i=1}^n \prod_{j=1}^g \pi_j(x^{(i)})^{\mathbb{1}_{\{y^{(i)}=j\}}}.$$

(By following the maximum likelihood principle, we should look for parameters  $\theta_1, \dots, \theta_g$ , which maximize the expression above.)

Now we want the empirical risk to be a *sum* of loss function values, not a *product* recall:

$$\mathcal{R}_{\text{emp}} = \sum_{i=1}^n L(y^{(i)}, f(\mathbf{x}^{(i)})).$$

We can turn the product into a sum by taking its log since the same parameters maximize the expressions. Finally, we convert the maximization into minimization by multiplying with -1. So we end up with the so-called cross-entropy loss function:

$$L(y, f(\mathbf{x})) = - \sum_{j=1}^g \mathbb{1}_{\{y=j\}} \log[\pi_j(x)].$$

We see that for the softmax regression the loss function is equal to the negative log-likelihood of one observation. Thus the associated empirical risk is the negative log-likelihood of the complete data set.

- c) Since the subtraction of any fixed vector from all  $\theta_k$  does not change the prediction, one set of parameters is "redundant". Thus we set  $\theta_g = (0, \dots, 0)$ . Hence for  $g$  classes we get  $g - 1$  discriminant functions from the softmax  $\hat{\pi}_1(x), \dots, \hat{\pi}_{g-1}(x)$  which can be interpreted as probability. The probability for class  $g$  can be calculated by using  $\hat{\pi}_g = 1 - \sum_{k=1}^{g-1} \hat{\pi}_k(x)$ . To estimate the class we are using majority vote:

$$\hat{y} = \arg \max_k \hat{\pi}_k(x)$$

The parameter of the softmax regression is defined as parameter matrix where each class has its own parameter vector  $\theta_k$ ,  $k \in \{1, \dots, g - 1\}$ :

$$\theta = [\theta_1, \dots, \theta_{g-1}]$$