

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

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

Evaluation through cross-entropy error (no regularization):

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

Evaluation through cross-entropy error ( $l_2$ -norm regularization):

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

$$\nabla \mathcal{E}(\mathbf{w}) = \sum_{n=1}^N (\sigma(\mathbf{w}^T \mathbf{x}_n) - y_n) \mathbf{x}_n + 2\lambda \mathbf{w}$$

Approximate solution through gradient descent:

$$\mathbf{w}(k+1) := \mathbf{w} - \alpha \left( \sum_{n=1}^N (\sigma(\mathbf{w}^T \mathbf{x}_n) - y_n) \mathbf{x}_n + 2\lambda \mathbf{w} \right)$$

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

The above Hessian is positive semi-definite, since both matrices  $\mathbf{x}_n \mathbf{x}_n^T$  and  $\mathbf{I}_{D \times D}$  are positive semi-definite.