

L2 Regularisation for Perceptron

L2 regularisation

- Let us denote by $L(b, \bar{W}, \mathcal{D})$ the Loss of classifying the dataset \mathcal{D} using the model represented by the weight vector \bar{W} and the bias term b
- We would like to impose L2 regularisation on \bar{W} .
- The overall objective to minimise can then be written as follows

$$J(b, \bar{W}, \mathcal{D}) = L(b, \bar{W}, \mathcal{D}) + \lambda ||\bar{W}||_2^2 = L(b, \bar{W}, \mathcal{D}) + \lambda \sum_{i=1}^d w_i^2.$$

Here λ is called the **regularisation coefficient** and is usually set via **cross-validation**.

- The gradient of the overall objective simply becomes the sum of the loss-gradient and the scaled weight vector \bar{W} .

$$\nabla_{b, w_1, \dots, w_d} J(b, \bar{W}, \mathcal{D}) = \nabla_{b, w_1, \dots, w_d} L(b, \bar{W}, \mathcal{D}) + 2\lambda \cdot (0, w_1, \dots, w_d)^T$$

L2 regularisation (SGD version)

$$\nabla_{b, w_1, \dots, w_d} J(b, \bar{W}, \bar{X}, y) = \nabla_{b, w_1, \dots, w_d} L(b, \bar{W}, \bar{X}, y) + 2\lambda \cdot (0, w_1, \dots, w_d)^T$$

Loss function $L(b, \bar{W}, \mathcal{D})$ for Perceptron

$$L(b, \bar{W}, \mathcal{D}) = \sum_{k=1}^n L(b, \bar{W}, \bar{X}_k, y_k) = \sum_{k=1}^n h(-y_k \cdot a_k)$$

where $h(t) = \max(0, t)$

Use the gradient descent method:

$$(b, w_1, \dots, w_d)^T \leftarrow (b, w_1, \dots, w_d)^T - \mu \nabla_{b, w_1, \dots, w_d} L(b, \bar{W}, \mathcal{D})$$

$$\nabla_{b, w_1, \dots, w_d} L(b, \bar{W}, \mathcal{D}) = \sum_{k=1}^n \nabla_{b, w_1, \dots, w_d} L(b, \bar{W}, \bar{X}_k, y_k) = \sum_{k=1}^n \nabla_{b, w_1, \dots, w_d} h(-y_k \cdot a_k)$$

Computation of $\nabla_{b,w_1,\dots,w_d} h(-y_k \cdot a_k)$

If (\bar{X}_k, y_k) is misclassified, then

$$\nabla_{b,w_1,\dots,w_d} h(-y_k \cdot a_k) = \left(\frac{\partial h(-y_k \cdot a_k)}{\partial b}, \frac{\partial h(-y_k \cdot a_k)}{\partial w_1}, \dots, \frac{\partial h(-y_k \cdot a_k)}{\partial w_d} \right)^T = -y_k \cdot \left(1, x_k^{(1)}, x_k^{(2)}, \dots, x_k^{(d)} \right)^T$$

Otherwise

$$\nabla_{b,w_1,\dots,w_d} h(-y_k \cdot a_k) = (0, 0, 0, \dots, 0)^T$$

L2 regularisation for the Perceptron (SGD version)

$$\nabla_{b,w_1,\dots,w_d} J(b, \bar{W}, \bar{X}, y) = \nabla_{b,w_1,\dots,w_d} L(b, \bar{W}, \bar{X}, y) + 2\lambda \cdot (0, w_1, \dots, w_d)^T$$

If (\bar{X}, y) is misclassified, then

$$\nabla_{b,w_1,\dots,w_d} L(b, \bar{W}, \bar{X}, y) = \nabla_{b,w_1,\dots,w_d} h(-y \cdot a) = -y \cdot (1, x_1, \dots, x_d)^T$$

Otherwise

$$\nabla_{b,w_1,\dots,w_d} h(-y \cdot a) = (0, 0, 0, \dots, 0)^T$$

L2 regularisation for the Perceptron (SGD version)

$$\nabla_{b,w_1,\dots,w_d} J(b, \bar{W}, \bar{X}, y) = \nabla_{b,w_1,\dots,w_d} L(b, \bar{W}, \bar{X}, y) + 2\lambda \cdot (0, w_1, \dots, w_d)^T$$

If (\bar{X}, y) is misclassified, then

$$\nabla_{b,w_1,\dots,w_d} J(b, \bar{W}, \bar{X}, y) = -y \cdot (1, x_1, \dots, x_d)^T + 2\lambda \cdot (0, w_1, \dots, w_d)^T$$

Otherwise

$$\nabla_{b,w_1,\dots,w_d} J(b, \bar{W}, \bar{X}, y) = (0, 0, \dots, 0)^T + 2\lambda \cdot (0, w_1, \dots, w_d)^T$$

L2 regularisation for the Perceptron (SGD version)

$$\nabla_{b,w_1,\dots,w_d} J(b, \bar{W}, \bar{X}, y) = \nabla_{b,w_1,\dots,w_d} L(b, \bar{W}, \bar{X}, y) + 2\lambda \cdot (0, w_1, \dots, w_d)^T$$

If (\bar{X}, y) is misclassified, then

$$w_i = w_i \cdot (1 - 2\lambda) + y \cdot x_i$$

$$b = b + y$$

$$\nabla_{b,w_1,\dots,w_d} J(b, \bar{W}, \bar{X}, y) = -y \cdot (1, x_1, \dots, x_d)^T + 2\lambda \cdot (0, w_1, \dots, w_d)^T$$

Otherwise

$$w_i = w_i \cdot (1 - 2\lambda)$$

$$b = b$$

$$\nabla_{b,w_1,\dots,w_d} J(b, \bar{W}, \bar{X}, y) = (0, 0, \dots, 0)^T + 2\lambda \cdot (0, w_1, \dots, w_d)^T$$