

BAUHAUS UNIVERSITÄT WEIMAR
INTRODUCTION TO MACHINE LEARNING



SUBMITTED BY

Anjana MuraleedharanNair(125512)

Isabel Maria Binu (125514)

Vishal Sanjay Shivam (125353)

Sharat Anand (125404)

Exercise 1: Linear Models

(a) $l(c, y(x))$ - pointwise loss

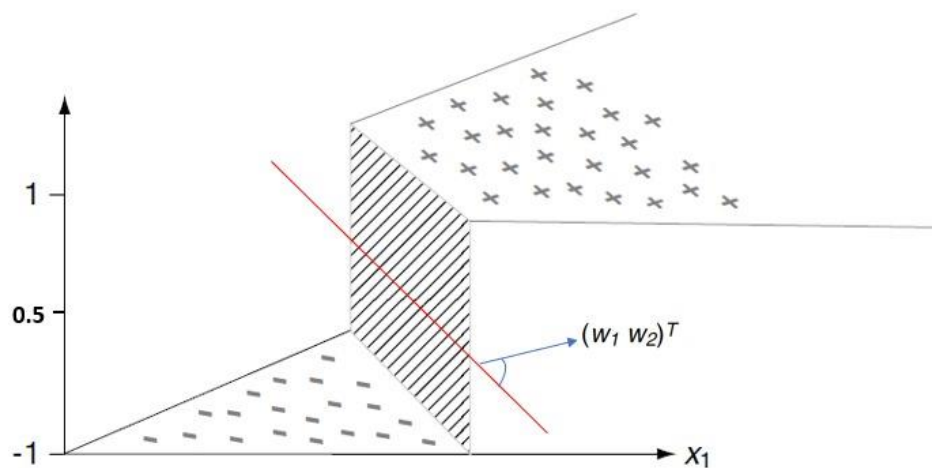
$L(w)$ - Global loss function

$L(w)$ - Error function

w - the entire parameter vector

\vec{w} - direction vector excluding w_0 .

(b)



(c) Both linear regression and logistic regression involve optimising w through a loss minimization problem, where "loss" refers to the interpretation of residuals. Since the residuals in linear regression are squared, outliers are given a high weight. However, logistic regression models an increasing confidence in class membership probability with increasing hyperplane distance. This difference in interpretation typically leads to a different parameter vector w and, consequently, a different decision boundary for logistic regression compared to linear regression.

(d) Lasso regression includes a penalty term based on the L1 norm of the weight vector ($\|w\|_1$). This penalty encourages sparsity, meaning it tends to drive some of the coefficients to exactly zero during optimization. This is due to the geometric shape of the L1 norm, which has corners at the axes, making it more likely for the optimization process to hit exactly zero when minimizing the combined loss and penalty term. In contrast, the L2 penalty used in ridge regression ($\|w\|_2$) tends to shrink coefficients towards but not exactly to zero, as it penalizes large values without promoting true sparsity.

(e) Gradient descent can't be directly applied to $L_0/1$ regularization due to non-differentiability at points where vectors have zero elements.

Exercise 2: Pointwise Loss Functions

(a) $l\sigma(0, y(x)) = \log(1 + e^{w^T x})$

Sigmoid function $\rightarrow \sigma(z) = 1/(1 + e^z)$

$$\sigma(w^T x) = 1/(1 + e^{w^T x})$$

$$\sigma(-a) = 1 - \sigma(a) \text{ (given)}$$

$$\sigma(-w^T x) = 1 - \sigma(w^T x)$$

$$\sigma(w^T x) = 1 - \sigma(-w^T x)$$

$$\sigma(w^T x) = 1 - (1/(1 + e^{w^T x})) = (1 + e^{w^T x} - 1)/(1 + e^{w^T x})$$

$$\sigma(w^T x) = e^{w^T x} / (1 + e^{w^T x})$$

Taking natural logarithm on both sides

$$\log(\sigma(w^T x)) = \log(e^{w^T x} / (1 + e^{w^T x}))$$

$$= \log(e^{w^T x}) - \log(1 + e^{w^T x})$$

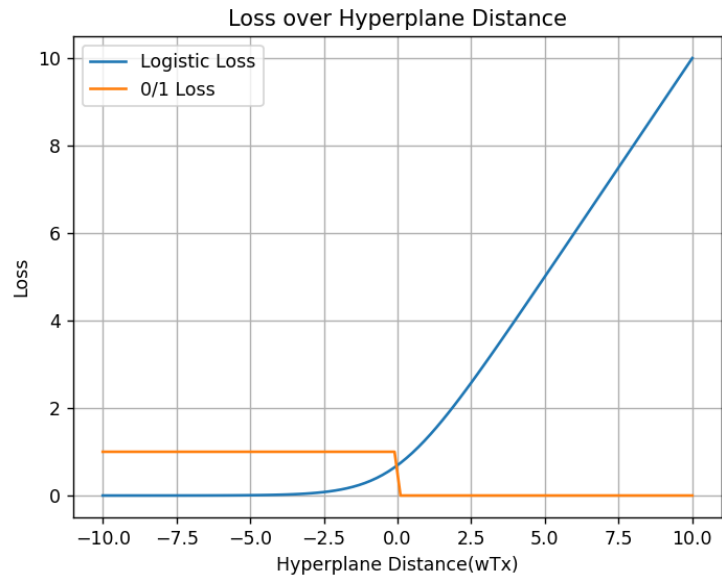
$$\log(a/b) = \log(a) - \log(b)$$

$$= w^T x - \log(1 + e^{w^T x}) \quad (\log(e^a) = a)$$

Substituting $w^T x = 0$

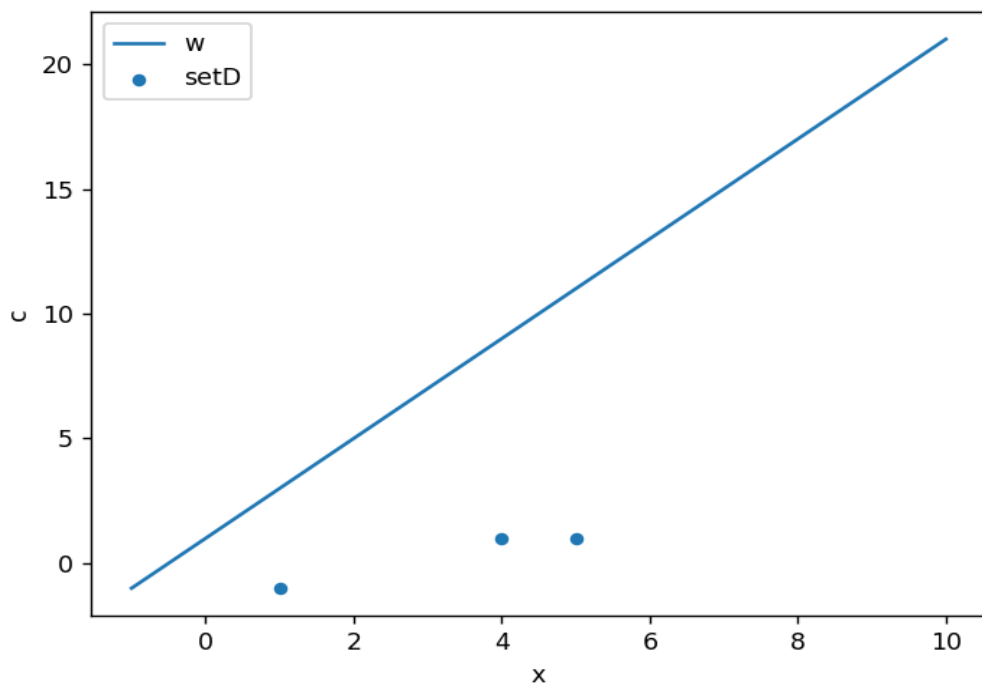
$$l\sigma(0, y(x)) = \log(1 + e^{w^T x})$$

(b)



Exercise 3 : Gradient Descent

(a)



(b) For $x_1=(4,1)$, $c=1$

$$y(x_1) = w^T x_1 = 9$$

$$l_2(c, y(x_1)) = \frac{1}{2} * (c - y(x_1))^2 = \frac{1}{2} * (1-9)^2 = 32$$

(c)

$$l_2 = \frac{1}{2} * (\mathbf{c} - \mathbf{w}^T \mathbf{x})^2$$

$$\partial l_2 / \partial \mathbf{w} = -(\mathbf{c} - \mathbf{w}^T \mathbf{x}) * \partial / \partial \mathbf{w} (\mathbf{w}^T \mathbf{x}) \quad (\mathbf{w}^T \mathbf{x} = \mathbf{w}_0 \mathbf{x}_0 + \mathbf{w}_1 \mathbf{x}_1)$$

$$= -(\mathbf{c} - \mathbf{w}^T \mathbf{x}) \mathbf{x}_0$$

$$\partial l_2 / \partial \mathbf{w} = -(\mathbf{c} - \mathbf{w}^T \mathbf{x}) * \partial / \partial \mathbf{w} (\mathbf{w}^T \mathbf{x}) \quad (\mathbf{w}^T \mathbf{x} = \mathbf{w}_0 \mathbf{x}_0 + \mathbf{w}_1 \mathbf{x}_1)$$

$$= -(\mathbf{c} - \mathbf{w}^T \mathbf{x}) \mathbf{x}_1$$

$$(\partial l_2 / \partial \mathbf{w}_0, \partial l_2 / \partial \mathbf{w}_1)^T = \begin{bmatrix} -(\mathbf{c} - \mathbf{w}^T \mathbf{x}) \mathbf{x}_0 \\ -(\mathbf{c} - \mathbf{w}^T \mathbf{x}) \mathbf{x}_1 \end{bmatrix}$$

$$= -\delta * \mathbf{x}$$

where $\delta = -(\mathbf{c} - \mathbf{w}^T \mathbf{x})$ and $\mathbf{x} =$

$$\begin{bmatrix} \mathbf{x}_0 \\ \mathbf{x}_1 \end{bmatrix}$$

(d) $\partial l_2 / \partial \mathbf{w} = -(\mathbf{c} - \mathbf{w}^T \mathbf{x}) * \partial / \partial \mathbf{w} (\mathbf{w}^T \mathbf{x})$

$$= -(\mathbf{c} - \mathbf{w}^T \mathbf{x}) \mathbf{x}$$

$$\partial l_2 / \partial \mathbf{x} = -(\mathbf{c} - \mathbf{w}^T \mathbf{x}) * \partial / \partial \mathbf{x} (\mathbf{w}^T \mathbf{x})$$

$$= -(\mathbf{c} - \mathbf{w}^T \mathbf{x}) \mathbf{w}$$

Loss gradient for $\mathbf{x}_1 = -(1-9).1 = 8$

Loss gradient for $\mathbf{w} = -(1-9).4 = 32$

(e)

$$\Delta \mathbf{w} = -\eta * \nabla L_2(\mathbf{w})$$

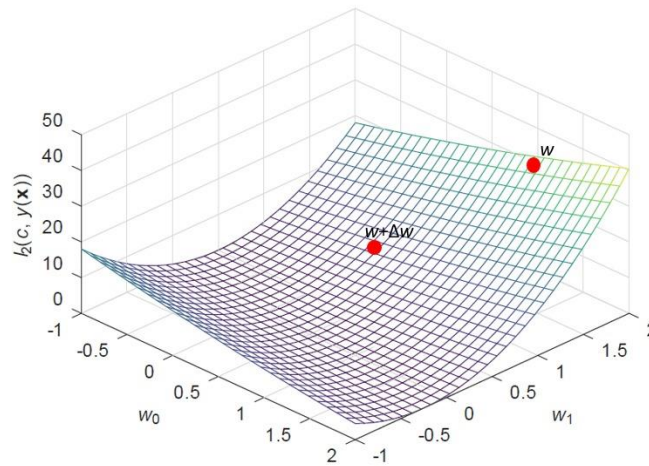
$$\nabla L_2 = -\Sigma(\mathbf{c} - \mathbf{w}^T \mathbf{x}) * \mathbf{x}$$

$$\Delta \mathbf{w} = -\eta * (\partial l_2 / \partial \mathbf{w}_0, \partial l_2 / \partial \mathbf{w}_1) = -0.03 * \begin{bmatrix} 8 \\ 32 \end{bmatrix}$$

$$= \begin{bmatrix} -0.24 \\ -0.96 \end{bmatrix}$$

$$w + \Delta w = \begin{matrix} 0.76 \\ 1.04 \end{matrix}$$

(f)



Exercise 4 : Regularization

(a) Option A - steadily increase

$$L(w) = \text{RSS}_{\text{tr}}(w) + \lambda w^T w$$

$$\text{RSS}_{\text{tr}}(\mathbf{w}) = \sum_{(x_i, y_i) \in D_{\text{tr}}} (y_i - \mathbf{w}^T x_i)^2$$

Here λ is increased from 0.

If $\lambda=0$, then $L(w)=\text{RSS}_{\text{tr}}(w)$

λ controls the impact of $w^T \cdot w$

So, when λ increases, the $\text{RSS}_{\text{tr}}(w)$ value also steadily increases.

(b) Option E - decrease initially, then eventually start increasing in a U shape.

When λ equals zero, the model effectively captures the training data. Therefore, we adjust the λ value to prevent overfitting. However, the λ parameter's variation can cause the model to fit the D_{test} only up to a certain threshold. Consequently, an initial increase in λ leads to reduced error on D_{test} . Nevertheless, beyond a certain threshold, the model ceases to fit the D_{test} , resulting in elevated $\text{RSS}_{\text{test}}(w)$

So, $\text{RSS}_{\text{test}}(w)$ will decrease initially, then eventually start increasing in a U shape.