

MULTIVARIATE LOSS FUNCTIONS

Advanced Machine Learning

Multi-Target Prediction: Loss Functions

- In MTP: For a feature vector \mathbf{x} , predict a tuple of scores $f(\mathbf{x}) = (f(x)_1, f(x)_2, \dots, f(x)_l)^T$ for l targets with a function (hypothesis) $f: \mathcal{X} \rightarrow \mathbb{R}^{g_1} \times \dots \times \mathbb{R}^{g_l}$.
- Following loss minimization in machine learning, we need a *multivariate loss function*



$$L: (\mathcal{Y}_1 \times \dots \times \mathcal{Y}_l) \times (\mathbb{R}^{g_1} \times \dots \times \mathbb{R}^{g_l}) \rightarrow \mathbb{R}.$$

- In multi-target regression: $\mathcal{Y}_1 = \dots = \mathcal{Y}_l = \mathbb{R}$, and $g_1 = \dots = g_l = 1$.
- In multi-label binary classification: $\mathcal{Y}_1 = \dots = \mathcal{Y}_l = \{0, 1\}$, and $g_1 = \dots = g_l = 1$.

Learning goals

- Get to know loss functions for multi-target prediction problems
- Understand the difference between instance-wise and decomposable losses
- Know risk minimizer for Hamming and subset 0/1 loss

Instance-wise

Decomposable

MULTIVARIATE LOSS FUNCTIONS

- In MTP: For a feature vector \mathbf{x} , predict a tuple of scores $f(\mathbf{x}) = (f(\mathbf{x})_1, f(\mathbf{x})_2, \dots, f(\mathbf{x})_l)^T$ for l targets with a function (hypothesis) $f: \mathcal{X} \rightarrow \mathbb{R}^{g_1} \times \dots \times \mathbb{R}^{g_l}$.
- We treat two categories: Decomposable and instance-wise
- Following loss minimization in machine learning, we need a *multivariate loss function*



- L is decomposable over targets if $L: (\mathcal{Y}_1 \times \dots \times \mathcal{Y}_l) \times (\mathbb{R}^{g_1} \times \dots \times \mathbb{R}^{g_l}) \rightarrow \mathbb{R}$.

- In multi-target regression: $\mathcal{Y}_1 = \dots = \mathcal{Y}_l = \mathbb{R}$, and $g_1 = \dots = g_l = 1$.

- with single-target binary classification: $\mathcal{Y}_1 = \dots = \mathcal{Y}_l = \{0, 1\}$, and

- Example: Squared error loss (in multivariate regression):

$$L_{\text{MSE}}(\mathbf{y}, f) = \frac{1}{l} \sum_{m=1}^l (y_m - f(\mathbf{x})_m)^2.$$

- Can also be used for cases with missing entries.

INSTANCE-WISE LOSSES

- Hamming loss averages over mistakes in single targets:

- We treat two categories: Decomposable and instance-wise

$$L_H(\mathbf{y}, \mathbf{h}) = \frac{1}{I} \sum_{m=1}^I \mathbb{1}_{[y_m \neq h_m(\mathbf{x})]}$$



- where $h_m(\mathbf{x}) := \mathbb{1}_{[f(\mathbf{x})_m \geq c_m]}$ is the threshold function for target m with threshold c_m .

- Hamming loss is identical to the average 0/1 loss and is decomposable.

$$L(\mathbf{y}, f) = \frac{1}{I} \sum_{m=1}^I L_m(y_m, f(\mathbf{x})_m)$$

- The subset 0/1 loss checks for entire correctness and is not decomposable.

- Example: Squared error loss (in multivariate regression):

$$L_{0/1}(\mathbf{y}, \mathbf{h}) = \mathbb{1}_{[\mathbf{y} \neq \mathbf{h}]} = \max_m \mathbb{1}_{[y_m \neq h_m(\mathbf{x})]}$$

$$L_{\text{MSE}}(\mathbf{y}, f) = \frac{1}{I} \sum_{m=1}^I (y_m - f(\mathbf{x})_m)^2.$$

- Can also be used for cases with missing entries.

HAMMING VS. SUBSET 0/1 LOSS

- The risk minimizer for the Hamming loss is the *marginal mode*:

$$f^*(\mathbf{x})_m = \arg \max_{y_m \in \{0,1\}} \Pr(y_m | \mathbf{x}), \quad m = 1, \dots, l,$$

$$L_H(\mathbf{y}, \mathbf{h}) = \frac{1}{l} \sum_{m=1}^l \mathbb{1}_{[y_m \neq h_m(\mathbf{x})]},$$

while for the subset 0/1 loss it is the *joint mode*:

where $h_m(\mathbf{x}) := [f(\mathbf{x})_m \geq c_m]$ is the threshold function for target m with threshold c_m . $f^*(\mathbf{x}) = \arg \max_{\mathbf{y}} \Pr(\mathbf{y} | \mathbf{x})$.

- Hamming loss is identical to the average 0/1 loss and is decomposable.
- Marginal mode vs. joint mode:

- The *subset 0/1 loss* $\Pr(\mathbf{y})$ for entire correctness and is not decomposable.

0 0 0 0	0.30
0 1 1 1	0.17
1 0 1 1	0.18
1 1 0 1	0.17
1 1 1 0	0.18

Marginal mode: 1 1 1 1

Joint mode: 0 0 0 0

$$L_{0/1}(\mathbf{y}, \mathbf{h}) = \mathbb{1}_{[\mathbf{y} \neq \mathbf{h}]} = \max_m \mathbb{1}_{[y_m \neq h_m]}$$



HAMMING VS. SUBSET 0/1 LOSS

- The risk minimizer for the Hamming loss is the *marginal mode*:

$$f^*(\mathbf{x})_m = \arg \max_{y_m \in \{0,1\}} \Pr(y_m | \mathbf{x}), \quad m = 1, \dots, l,$$

while for the subset 0/1 loss it is the *joint mode*:

$$f^*(\mathbf{x}) = \arg \max_{\mathbf{y}} \Pr(\mathbf{y} | \mathbf{x}).$$

- Marginal mode vs. joint mode:

\mathbf{y}	$\Pr(\mathbf{y})$		
0 0 0 0	0.30		
0 1 1 1	0.17	Marginal mode:	1 1 1 1
1 0 1 1	0.18	Joint mode:	0 0 0 0
1 1 0 1	0.17		
1 1 1 0	0.18		

