

9/7/2021

AI 1104

- Analyse two SL models:

- Linear regression

- Artificial neural network

Supervised Learning (SL): $D = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$ (Entire dataset)

$T = \{(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)\}$ (Training dataset)

$V = \{(x_{m+1}, y_{m+1}), \dots, (x_N, y_N)\}$ (Test/Val. dataset)



Predicted value: $\hat{y} = f(x; \theta)$

Loss: $R(\theta) = \sum_{i=1}^n R_i(\theta)$

Cumulative loss
over training
samples

$$= \sum_{i=1}^n d(y_i, \hat{y}_i)$$

$$= \sum_{i=1}^n d(y_i, f(x_i; \theta))$$

$d: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}_+$ have $d(\cdot)$ be a valid distance
Set of non-negative reals
Set of values that our loss can take

Goal: Find $\hat{\theta}^* = \underset{\theta \in \Theta}{\operatorname{arg\,min}} R(\theta)$

(i) Linear regression: Assume $x, y \in \mathbb{R}$, $\hat{y} = \beta_0 + \beta_1 x = f(x; \theta)$ where $\theta = \{\beta_0, \beta_1\}$

$$R_i(\theta) = (\hat{y}_i - y_i)^2 = (y_i - (\beta_0 + \beta_1 x_i))^2$$

i is the sample index

$$\Rightarrow R(\theta) = \sum_{i=1}^n (\hat{y}_i - y_i)^2$$

$$R(\theta) = \|X\theta - y\|_2^2$$

$$\theta = \{\beta_0, \beta_1\}$$

vector,

where $\beta = \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$, $\underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix}$; $X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$

Note: $\underline{y} = X \cdot \beta$

$$\theta^* = \beta^* = \underset{\beta \in \mathbb{H}}{\operatorname{argmin}} \|X \cdot \beta - \underline{y}\|_2^2$$

$$\|X\|_2 = \left(\sum_{i=1}^m x_i^2 \right)^{1/2}$$

$\Rightarrow l_2$ norm of vector

If $e = \underline{y} - \hat{y}$, then our goal is to min. the l_2 norm of e .

$$R(\theta) = \|X \cdot \beta - \underline{y}\|^2 = e^T e = (X \cdot \beta - \underline{y})^T (X \cdot \beta - \underline{y})$$

$$\Rightarrow R(\theta) = \beta^T X^T X \beta - \underbrace{2 \beta^T X^T \underline{y}}_{(2)} + \underline{y}^T \underline{y}$$

Q: let $a, b \in \mathbb{R}^n$, is $a^T b = b^T a$?

$$\nabla_\theta R(\theta) = 0 \Leftrightarrow \nabla_\beta (\beta^T X^T X \beta - 2 \beta^T X^T \underline{y} + \underline{y}^T \underline{y}) = 0$$

$$2 \cdot \underline{x}^T \underline{x} \beta - 2 \underline{x}^T \underline{y} = 0 \quad (\text{multiply by } 1)$$

$$\underline{x}^T \underline{x} \beta = \underline{x}^T \underline{y} \quad \text{--- (1)}$$

If $\underline{x}^T \underline{x}$ is a positive definite matrix ($\Rightarrow \underline{x}^T \underline{x}^{-1}$ exists),

$$\text{then } \beta^* = (\underline{x}^T \underline{x})^{-1} \cdot \underline{x}^T \underline{y} \quad \text{--- (2)}$$

Moore-Penrose inverse of X ; Pseudo-inverse of X .

We can arrive at (1) using scalar or regular differentiation operation.

Specifically $\frac{\partial R(\theta)}{\partial \beta_0}$; $\frac{\partial R(\theta)}{\partial \beta_1}$. Try this as an exercise.

Given: $\hat{y} = \beta_0 + \beta_1 x_1 + \beta_2 x_2^2 + \dots + \beta_m x_m^m = f(x; \theta) : \underline{x}, \underline{y} \in \mathbb{R}$

Q: Find $\theta^* = \hat{\beta}^*$; what is $d(y, \hat{y})$?

A: $\hat{\beta}^* = (X^T X)^{-1} \cdot X^T \underline{y}$ where $X = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots & x_1^m \end{bmatrix}$ $\hat{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_m \end{bmatrix}_{m \times 1}$

Identical solution as before.

$$\underline{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}_{n \times m}$$

$n \times 1$

(ii) Multilayer Perception or Artificial Neural Network (Following Elements of Statistical Learning)

Hastie & Tibshirani

