

Elements of Statistical Learning.

- For several cases the goal is to use the inputs to predict the values of the outputs. (Supervised Learning)

Qualitative variables : Categorical or discrete variables.

- Distinction in output type $\left\{ \begin{array}{ll} \text{Regression} & \text{Quantitative outputs} \\ \text{Classification} & \text{Qualitative outputs.} \end{array} \right.$

- Qualitative values: "Success" or "Failure"
"Survived" or "died"

Input variable : X

Quantitative output Y

Qualitative output G

Observed values are written in lowercase; i th of X is written as x_i

Linear models :

Given a vector of inputs $X^T = (x_1, x_2, \dots, x_p)$, we predict the output Y via the model

$$\hat{Y} = \hat{\beta}_0 + \sum_{j=1}^n x_j \hat{\beta}_j \quad (2.1)$$

Include $\hat{\beta}_0$ in the vector of coefficients $\hat{\beta}$ and then write the linear model in vector form:

$$\hat{Y} = X^T \hat{\beta} \quad (2.2)$$

How do we fit the linear model to a set of training data?

By far the most popular is the method of least squares. In this approach, we pick the coefficients β to minimize the residual sum of squares

$$RSS(\beta) = \sum_{i=1}^n (y_i - x_i^T \beta)^2 \quad (2.3)$$

y_i = Real values
 $\hat{y}_i = x_i^T \beta$ Predicted values.

Quadratic function of the parameters \Rightarrow Minimum always exists, but may not be unique. We can write:

$$RSS(\beta) = (y - X\beta)^T (y - X\beta) \quad (2.4)$$

$$X \in \mathbb{R}^{n \times p}$$

(Each row is an input vector)

$$y \in \mathbb{R}^n$$

(Vector of the outputs)

$$\alpha = x^T A x$$

$$\frac{\partial \alpha}{\partial x} = x^T (A + A^T)$$

$$\alpha = \sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j \quad \text{Taking } \frac{\partial \alpha}{\partial x_k} = \sum_{ij} a_{ij} x_i \delta_{jk} + \sum_{ij} a_{ij} \delta_{ik} x_j$$

$$\frac{\partial \alpha}{\partial x_k} = \sum_i a_{ik} x_i + \sum_j a_{kj} x_j = x^T A^T + x^T A$$

$$\frac{\partial \alpha}{\partial x} = x^T (A^T + A) \quad \text{if the matrix is symmetric:}$$

$$\frac{\partial \alpha}{\partial x} = 2 x^T A.$$

$$\frac{\partial R}{\partial \beta} = -2 x^T (y - X\beta) \Rightarrow x^T (y - X\beta) = 0$$

Organizing: $x^T y - x^T X \beta = 0$

$$\hat{\beta} = (X^T X)^{-1} X^T y = 0 \quad (2.6)$$

The fitted value at i th input is $\hat{y}_i = \hat{y}_i(x_i) = x_i^T \hat{\beta}$

• Probabilistic approach:

Regression problem with $p(y|x) = \mathcal{N}(y | f(x), \sigma^2)$

$x \in \mathbb{R}^D$ inputs
 $y \in \mathbb{R}$ targets. $\hat{y} = f(x) + \epsilon$.

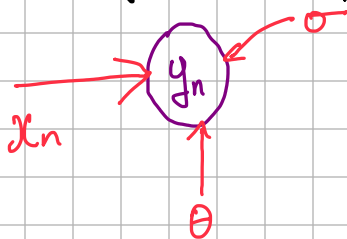
and $\epsilon = \mathcal{N}(0, \sigma^2)$ Gaussian distribution with 0 mean and variance σ^2 .

Model parameters θ (Learning this)

$$p(y|x, \theta) = \mathcal{N}(y | x^T \theta, \sigma^2) \Leftrightarrow y = x^T \theta + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2)$$

→ Parameter estimation:

We are given a training set $\mathcal{D} := \{(x_1, y_1), \dots, (x_n, y_n)\}$
Consisting of N inputs $x \in \mathbb{R}^D$



y_i and y_j are conditionally independent given their respective inputs.

$$p(y|x, \theta) = p(y_1, \dots, y_n | x_1, x_2, \dots, x_n; \theta)$$

$$p(y|x, \theta) = \prod_{i=1}^n p(y_i | x_i, \theta) = \prod_{i=1}^n \mathcal{N}(y_i | x_i^T \theta, \sigma^2)$$

Maximum likelihood estimation:

θ_{ML} : Maximizing likelihood (Maximize the predictive distribution of the training data given the model parameters:

$$\theta_{ML} = \underset{\theta}{\operatorname{argmax}} p(y|x, \theta)$$

Using the properties of logarithms on products, we minimize the negative log-likelihood:

$$-\log p(y|x;\theta) = -\log \prod_{i=1}^n p(y_i|x_i;\theta) = -\sum_{i=1}^n \log(p(y_i|x_i;\theta))$$

Factorization due to independence.

$$\log p(y_n|x_n;\theta) = -\frac{1}{2\sigma^2} (y_n - x_n^T \theta)^2 + c.$$

$$L(\theta) := \frac{1}{2\sigma^2} \sum_{n=1}^n (y_n - x_n^T \theta)^2$$

$$L(\theta) = \frac{1}{2\sigma^2} (y - X^T \theta)^T (y - X \theta) = \frac{1}{2\sigma^2} \|y - X \theta\|^2.$$

Computing the Gradient:

$$\begin{aligned} \frac{\partial L}{\partial \theta} &= \frac{\partial L}{\partial \theta} \left(\frac{1}{2\sigma^2} (y - X \theta)^T (y - X \theta) \right) \\ &= \frac{1}{2\sigma^2} \frac{d}{d\theta} (y^T y - 2y^T X \theta + \theta^T X^T X \theta) \\ &= \frac{1}{\sigma^2} (-y^T X + \theta^T X^T X) \end{aligned}$$

Again $\hat{\theta}_{HL} = (X^T X)^{-1} X^T y.$

Now we can use monomial to express different features.
For example for a second order polynomial:

$$\Phi = \begin{bmatrix} 1 & x_1 & x_1^2 \\ 1 & x_2 & x_2^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 \end{bmatrix}$$

$$x_n \in \mathbb{R}^D$$

$$y_n \in \mathbb{R}, \quad n = 1, 2, \dots, N$$

↓
Number of inputs.

$$\hat{\theta}_{HL} = (X^T X)^{-1} X^T y$$