

Linear Model

Fish

2019 年 6 月 22 日

Content

1	What	2
2	one sample has one property	2
2.1	background	2
2.2	analysis	2
2.3	processing	3
3	one sample has multiply property	4
3.1	analysis	4
3.2	processing	5
3.3	deep analysis	7
3.4	penalty factor	7
3.5	$ w $ and $ w ^2$	8
4	machine learning and data analysis	8

1 What

define the example x described by d property, $x = (x_1; x_2; \dots; x_d)$, among the variable x , x_i is decided by i th property of x .

linear model showed as follow:

$$f(x) = w_1x_1 + w_2x_2 + \dots + w_dx_d + b \quad (1)$$

whether x has d property

vector as well:

$$f(x) = w^T x + b \quad (2)$$

2 one sample has one property

to simplify the situation, we assume that the input x property has only one.

2.1 background

we have dataset: $D = (x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$, to study a model predict tab of output more correctly.

2.2 analysis

to study model: $f(x) = w^T x + b$, make the predict result $f(x_i) \cong y_i$, the w, b need to be known.

mean square error is corresponding to Euclidean distance. The way solve model based on mean square error named least square method. according to linear regression, least square method try to find a line that make the Euclidean distance sum of sample to line be least.

2.3 processing

find least value of equation 3 and 4, get the w^*, b^*

$$(w^*, b^*) = \arg \min_{(w, b)} \sum_{i=1}^m (f(x_i) - y_i)^2 \quad (3)$$

$$= \arg \min_{(w, b)} \sum_{i=1}^m (wx_i + b - y_i)^2 \quad (4)$$

the progress solving the above equation named least square parameter estimation. there differenetiate w and b, could get following part:

$$\frac{\partial E(w, b)}{\partial w} = \sum_{i=1}^m 2(-y_i + b + wx_i)x_i = 2 \sum_{i=1}^m [wx_i^2 - (y_i - b)x_i] \quad (5)$$

$$\frac{\partial E(w, b)}{\partial b} = \sum_{i=1}^m 2(-y_i + b + wx_i) = 2[mb - \sum_{i=1}^m (y_i - wx_i)] \quad (6)$$

then optimizer solution fo closed-form is shown:

$$w = \frac{\sum_{i=1}^m y_i(x_i - \bar{x})}{\sum_{i=1}^m x_i^2 - \frac{1}{m}(\sum_{i=1}^m x_i)^2} \quad (7)$$

$$b = \frac{1}{m} \sum_{i=1}^m (y_i - wx_i) \quad (8)$$

3 one sample has multiply property

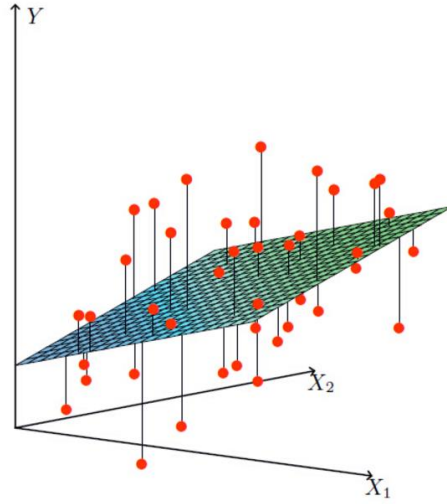


Fig 1: fitting plane in sample space

the common condition is as section one show, dataset D has d property.
at present, it should try to study

$$f(\mathbf{x}_i) = \mathbf{w}^T \mathbf{x}_i + b \quad (9)$$

let

$$f(\mathbf{x}_i) \cong y_i \quad (10)$$

this condition named multivariate linear regression.

3.1 analysis

similarly, estimate w and b by least square method. conveniently,

$$\hat{\mathbf{w}} = (\mathbf{w}; b) \quad (11)$$

dataset D show as a matrix \mathbf{X} , shape $m \times (d + 1)$

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1d} & 1 \\ x_{21} & x_{22} & \cdots & x_{2d} & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{m1} & x_{m2} & \cdots & x_{md} & 1 \end{pmatrix} = \begin{pmatrix} x_1^T & 1 \\ x_2^T & 1 \\ \vdots & \vdots \\ x_m^T & 1 \end{pmatrix} \quad (12)$$

$$\text{tag } \mathbf{y} = (y_1; y_2; \cdots; y_m) \quad (13)$$

so fitting plane

$$h_w(x) = w_0 + w_1x + w_2x + \dots + w_dx, w_0 \text{ is bias} \quad (14)$$

in summary,

$$h_w(x) = \sum_{i=0}^n w_i x_i = w^T x \quad (15)$$

3.2 processing

ε error, represent the difference between true and predict value.

for every sample,

$$y^i = w^T x^i + \varepsilon^i \quad (16)$$

whether ε obey the independent and identical distribution, and observe mean 0 and var θ^2 distribution.

ε obey gauge distribution,

$$p(\varepsilon^i) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\varepsilon^i)^2}{2\sigma^2}\right) \quad (17)$$

substitution equation 16 into equation 17, as following:

$$p(y^i | x^i; w) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^i - w^T x^i)^2}{2\sigma^2}\right) \quad (18)$$

likelihood function:

$$L(w) = \prod_{i=1}^m p(y^i | x^i; w) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^i - w^T x^i)^2}{2\sigma^2}\right) \quad (19)$$

logarithmic likelihood function:

$$\log L(w) = \log \prod_{i=1}^m p(y^i | x^i; w) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^i - w^T x^i)^2}{2\sigma^2}\right) \quad (20)$$

expand equation 20th:

$$\begin{aligned} L(w) &= \sum_{i=1}^m \log \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^i - w^T x^i)^2}{2\sigma^2}\right) \\ &= m \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{\sigma^2} \cdot \frac{1}{2} \sum_{i=1}^m (y^i - w^T x^i)^2 \end{aligned} \quad (21)$$

target: $L(w)$ more big and more better

$$J(w) = \frac{1}{2} \sum_{i=1}^m (y^i - w^T x^i)^2 \quad (\text{least square method}) \quad (22)$$

target function:

$$\begin{aligned} J(w) &= \frac{1}{2} \sum_{i=1}^m (h_w(x^i) - y^i)^2 \\ &= \frac{1}{2} (\mathbf{X}\hat{\mathbf{w}} - \mathbf{y})^T (\mathbf{X}\hat{\mathbf{w}} - \mathbf{y}) \end{aligned} \quad (23)$$

solve partial derivation:

$$\begin{aligned} \nabla_w J(w) &= \nabla_w \left(\frac{1}{2} (\mathbf{X}\hat{\mathbf{w}} - \mathbf{y})^T (\mathbf{X}\hat{\mathbf{w}} - \mathbf{y}) \right) \\ &= \nabla_w \left(\frac{1}{2} (\hat{\mathbf{w}}^T \mathbf{X}^T - \mathbf{y}^T) (\mathbf{X}\hat{\mathbf{w}} - \mathbf{y}) \right) \\ &= \nabla_w \left(\left(\frac{1}{2} (\hat{\mathbf{w}}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - \hat{\mathbf{w}}^T \mathbf{X}^T \mathbf{y} - \mathbf{y} \mathbf{X} \hat{\mathbf{w}} + \mathbf{y}^T \mathbf{y}) \right) \right) \\ &= \frac{1}{2} (2\mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} - \mathbf{X}^T \mathbf{y} - \mathbf{X}^T \mathbf{y}) \\ &= \mathbf{X}^T \mathbf{X} \hat{\mathbf{w}} - \mathbf{X}^T \mathbf{y} \end{aligned} \quad (24)$$

let partial derivation be equal to 0:

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \quad (25)$$

estimate method: original estimate part

$$R^2 = 1 - \frac{\sum_{i=1}^m (\hat{y}_i - y_i)^2}{\sum_{i=1}^m (y_i - \bar{y})^2} \quad (26)$$

when R^2 value is very close approximation to 1, the model fits dataset D well.

3.3 deep analysis

if $X^T X$ is irreversible or avoid over fitting, increase λ destabilization

$$\theta = (X^T X + \lambda I)^{-1} X^T y \quad (27)$$

to remember conclusion simply,

$$X\theta = y \Rightarrow X^T X\theta = X^T y \Rightarrow \theta = (X^T X)^{-1} X^T y \quad (28)$$

after add λ destabilization:

- $X^T X$ is positive semidefinite matrix: for any non-zero vector u

$$\theta X^T X \theta = (X\theta)^T X\theta \xrightarrow{\text{define } v=X\theta} v^T v \geq 0 \quad (29)$$

- for any real $\lambda > 0$, $X^T X + \lambda I$ is positive definite matrix, and is reversible. be sure that regerssion equation must be significative.

$$\theta = (X^T X + \lambda I)^{-1} X^T y \quad (30)$$

3.4 penalty factor

- the target function of linear regerssion:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2 \quad (31)$$

- increase square sum loss:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2 + \lambda \sum_{j=1}^m \theta_j^2 \quad (32)$$

- actually, assume that parameter θ obey Gaussian distribution

3.5 $|w|$ and $|w|^2$

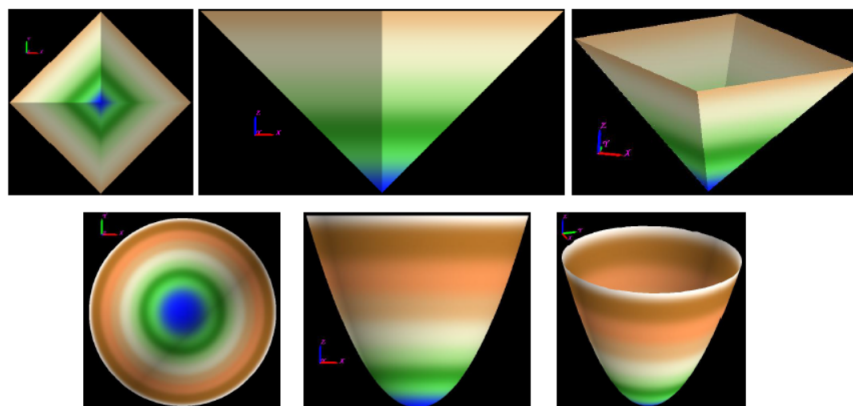


Fig 2: $|w|$ and $|w|^2$

4 machine learning and data analysis



Fig 3: machine learning and data analysis

- cross validation:
example: 10 fold cross validation