

Linear Regression

Thursday, January 16, 2020

13:34

Linear Regression

Real-value target function

(instead of yes/no, we estimate a real value)

Approximate some function $y_n = f(x_n)$ by a linear function

MODEL SETUP

$$\mathcal{D} = \{(x_1, y_1), \dots\} \text{ same as before}$$

$$x_i \in \mathbb{R}^{d+1} \quad y_i \in \mathbb{R} \quad \underline{w} \in \mathbb{R}^{d+1}$$

Regression: $\hat{y}_n = \underline{w}^T \underline{x}_n$
(predict by $\underline{w}^T \cdot \underline{x}_n$)

COMPACT REPRESENTATION

Data matrix

$$X = \begin{bmatrix} x_{10} = 1, & x_{11}, & \dots & x_{1d} \\ x_{20} = 1, & x_{21}, & \dots & x_{2d} \\ x_{30} = 1, & x_{31}, & \dots & x_{3d} \\ x_{N0} = 1, & x_{N1}, & \dots & x_{Nd} \end{bmatrix} \in \mathbb{R}^{N \times (d+1)}$$

$$X = \begin{bmatrix} \underline{x}_1^T \\ \underline{x}_2^T \\ \vdots \\ \underline{x}_N^T \end{bmatrix}$$

we represent input data as a matrix for compact computation

Observation Vector

$$\underline{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} \in \mathbb{R}^N$$

$N \times (d+1)$
size
1

$$L[y_N]$$

$$\Rightarrow \hat{\underline{y}} = \begin{bmatrix} \underline{x}_1^T \underline{w} \\ \underline{x}_2^T \underline{w} \\ \vdots \\ \underline{x}_N^T \underline{w} \end{bmatrix} = \underline{X} \underline{w}$$

N size \nearrow $N \times (d+1)$ size \downarrow
 $(d+1)$ size \uparrow

ideally $\hat{\underline{y}} = \underline{y}$ after training.

Case ① $N < d+1$:

Many Solutions (# of unknowns > # of equations)

Case ② $N = d+1$: (overfitting)

Single/Unique solution: $\hat{\underline{y}} = \underline{y}$

Case ③ $N > d+1$:

No solution for $\hat{\underline{y}} = \underline{y}$

$\Rightarrow \hat{\underline{y}} \neq \underline{y}$ (but approximation is good enough)

\Rightarrow Typically we want $N \gg d+1$ so we do not suffer from overfitting

Loss Function

intuition: we can measure distance:

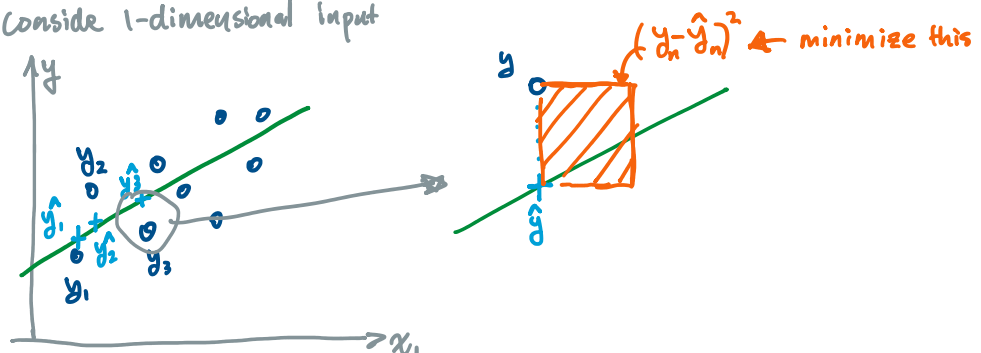
$$E_{in}(\underline{w}) = \frac{1}{N} \|\underline{y} - \hat{\underline{y}}\|_2^2 = \frac{1}{N} \sum_{n=1}^N (y_n - \hat{y}_n)^2$$

$$= \frac{1}{N} \|\underline{y} - \underline{X}\underline{w}\|_2^2 = \frac{1}{N} \sum_{n=1}^N (y_n - \underline{x}_n^T \underline{w})^2$$

error of each data: e_n

VISUAL EXAMPLE:

Consider 1-dimensional input



we want optimised weights

$$\underline{w}^* = \underset{\underline{w} \in \mathbb{R}^{d+1}}{\operatorname{argmin}} (E_{in}(\underline{w}))$$

$$= \underset{\underline{w}}{\operatorname{argmin}} \left(\frac{1}{N} \|\underline{y} - X\underline{w}\|^2 \right)$$

How do we minimize this?

Recall $\|\underline{a}\|^2 = \underline{a}^T \underline{a}$

$$E_{in}(\underline{w}) = \frac{1}{N} (\underline{y} - X\underline{w})^T (\underline{y} - X\underline{w})$$

$$= \frac{1}{N} \left(\underline{y}^T \underline{y} - \underbrace{\underline{w}^T X \underline{y}}_{-2\underline{y}^T X \underline{w}} + \underline{w}^T X^T X \underline{w} \right)$$

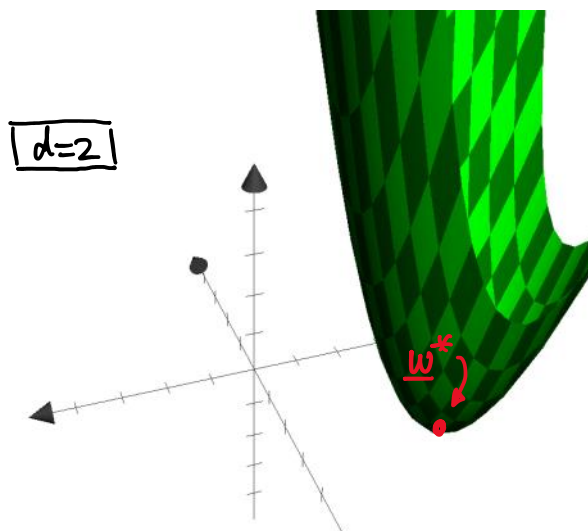
$$= \frac{1}{N} \left(\underline{y}^T \underline{y} - 2\underline{y}^T X \underline{w} + \underline{w}^T X^T X \underline{w} \right)$$

convexity determined by $X^T X$
 \hookrightarrow if all eigenvalues are \oplus / \ominus / mix of \oplus and \ominus

this looks like quadratic. 

Quadratic form in \underline{w} is convex

\Rightarrow there exists a local minimum.



To get to w^* , find global minimum.

$$\nabla(E_{in}(\underline{w})) = 0$$

LINEAR ALG RULES

$$\blacktriangleright f(\underline{w}) = \underline{w}^T \underline{v} = \underline{v}^T \underline{w}$$

$$\blacktriangleright f(\underline{w}) = \underline{w}^T \underline{v} = \underline{v}^T \underline{w}$$

$$\nabla f(\underline{w}) = \underline{v}$$

$$\blacktriangleright \nabla(\underline{w}^T A \underline{w}) = (A + A^T) \underline{w} = 2A \underline{w}$$

$$\nabla(E_{in}(\underline{w})) = \frac{1}{N} (0 - 2\underline{y}^T \underline{x} + 2\underline{x}^T \underline{x} \underline{w}) = 0$$

$$\Rightarrow \cancel{2} \underline{y}^T \underline{x} = \cancel{2} \underline{x}^T \underline{x} \underline{w}^*$$

$$\Rightarrow \underline{w}^* = (\underline{x}^T \underline{x})^{-1} \underline{x}^T \underline{y}$$

Least square solution.

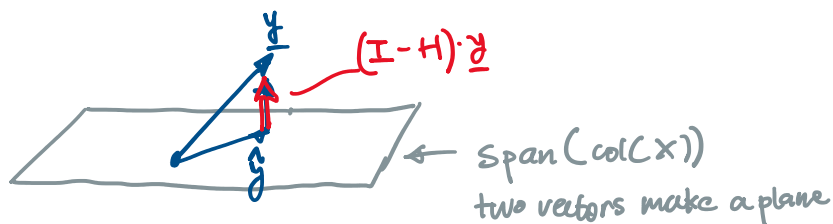
Finally, we can use \underline{w}^* to predict:

$$\hat{\underline{y}} = \underline{x} \underline{w}^* = \underbrace{\underline{x} (\underline{x}^T \underline{x})^{-1} \underline{x}^T}_{\text{H matrix}} \cdot \underline{y}$$

Special: This is a closed form solution.

Geometric Interpretation

$$\underline{x} \underline{w} = \sum_{i=0}^d w_i \cdot \underline{q}_i \quad \text{where } \underline{q}_i \text{ are columns of } \underline{x}$$



Regularization

Helps with overfitting

Regularized Least Square (Regression)

$$\operatorname{argmin} \left(\frac{1}{N} \|\underline{y} - X\underline{w}\|^2 + \lambda \|\underline{w}\|^2 \right)$$

$\lambda > 0$

We add a term to penalize large weights.

$$f(\underline{w}) = \underbrace{\frac{1}{N} \|\underline{y} - X\underline{w}\|^2}_{E_{in}(\underline{w})} + \lambda \underbrace{\|\underline{w}\|^2}_{\underline{w}^T \underline{w}}$$

$$\begin{aligned} \nabla f(\underline{w}) &= \frac{1}{N} (2X^T X \underline{w} + 2X^T \underline{y}) + 2\underline{w}\lambda = 0 \\ \Rightarrow \underline{w}^* &= (X^T X + \lambda I)^{-1} X^T \underline{y} \end{aligned}$$

How to choose " λ "?

\Rightarrow selected using validation

collected data broken into three types:

- \rightarrow training data : $\lambda_1, \lambda_2 \dots \lambda_M \Rightarrow \underline{w}_{\lambda_1}^*, \underline{w}_{\lambda_2}^* \dots \underline{w}_{\lambda_M}^*$
 \uparrow gives
- \rightarrow validation data : $\underline{E}_{valid_1}, \underline{E}_{valid_2} \dots \underline{E}_{valid_M}$
 \uparrow try different λ \searrow also gives
 \rightarrow pick $\underline{w}_{\lambda_m}^*$
- \rightarrow test data : $E_{test}(\underline{w}_{\lambda_m}^*)$
(not involved in training)