

Lecture 4

19

Linear Algebra of Linear Regression

- Residual Sum of Squares
- Normal Equation
- Ordinary Least Squares

In the next few lectures we will discuss linear methods for solving regression problems.

Let $X = (X_1, \dots, X_p)^T$ be an input vector.

We want to predict output $Y \in \mathbb{R}$, $\hat{Y} = f(X) \approx Y$.

Let's revisit the linear regression method.

In linear regression: $f(X) = \beta_0 + \sum_{j=1}^p \beta_j X_j = X^T \beta$, where $X = \begin{bmatrix} 1 \\ X_1 \\ \vdots \\ X_p \end{bmatrix}$ $\beta = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_p \end{bmatrix}$

unknown parameters or coefficients

Remark: Note that inputs X_1, \dots, X_p can come from different sources:

- Quantitative inputs, $X_j \in \mathbb{R}$
- Transformation of quantitative inputs, $X_j = \log \tilde{X}_j$, $X_j = \tilde{X}_j^2$, etc.
- Basis expansions, such as $X_2 = X_1^2$, $X_3 = X_1^3$, ..., $X_p = X_1^p$. \Rightarrow fitting a polynomial.
- Interactions between inputs, such as $X_1, X_2, X_3 = X_1 \cdot X_2$
- Dummy coding of qualitative inputs: if \tilde{X}_j has three levels (say, red, blue, pink)

then we can create X_1, X_2, X_3 to represent \tilde{X}_j .

$$X_1 = \begin{cases} 1 & \text{if } \tilde{X}_j = \text{red} \\ 0 & \text{otherwise} \end{cases} \quad X_2 = \begin{cases} 1 & \text{if } \tilde{X}_j = \text{blue} \\ 0 & \text{otherwise} \end{cases}$$

$$X_3 = \begin{cases} 1 & \text{if } \tilde{X}_j = \text{pink} \\ 0 & \text{otherwise} \end{cases}$$

Remark: The linear regression model is called "linear" because it is linear in parameters (not in inputs).

In Lecture 2, we found an estimate $\hat{\beta}$ of the vector of regression parameters β by first finding its optimal value with respect to the mean squared error and then estimating this optimal value using the training data.

$$MSE(\beta) = E[(Y - \hat{Y})^2] = E[(Y - X^T \beta)^2] \rightarrow \min$$

Lecture 2
page 9

It is instructive to see how β can be estimated directly from the data using the residual sum of squares:

$$RSS(\beta) = \sum_{i=1}^N (y_i - x_i^T \beta)^2, \quad x_i^T = (1, x_{i1}, \dots, x_{ip})$$

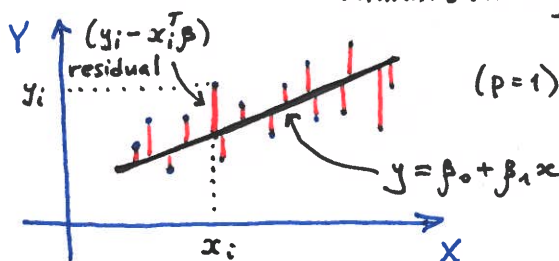
Remark: RSS is the finite sample version of MSE.

(up to a multiplicative factor $\frac{1}{N}$ which is irrelevant for minimization)

Geometrically, $RSS(\beta)$ is an intuitive measure of the quality of the linear fit to the data:

$$RSS(\beta_1) < RSS(\beta_2)$$

\Rightarrow linear fit corresponding to β_1 is better.



$$\beta = E[X X^T]^{-1} E[X Y]$$

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

where

$$X = \begin{bmatrix} 1 & x_{11} & \dots & x_{1p} \\ \vdots & \vdots & & \vdots \\ 1 & x_{i1} & \dots & x_{ip} \\ \vdots & \vdots & & \vdots \\ 1 & x_{N1} & \dots & x_{Np} \end{bmatrix}$$

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$$

training data

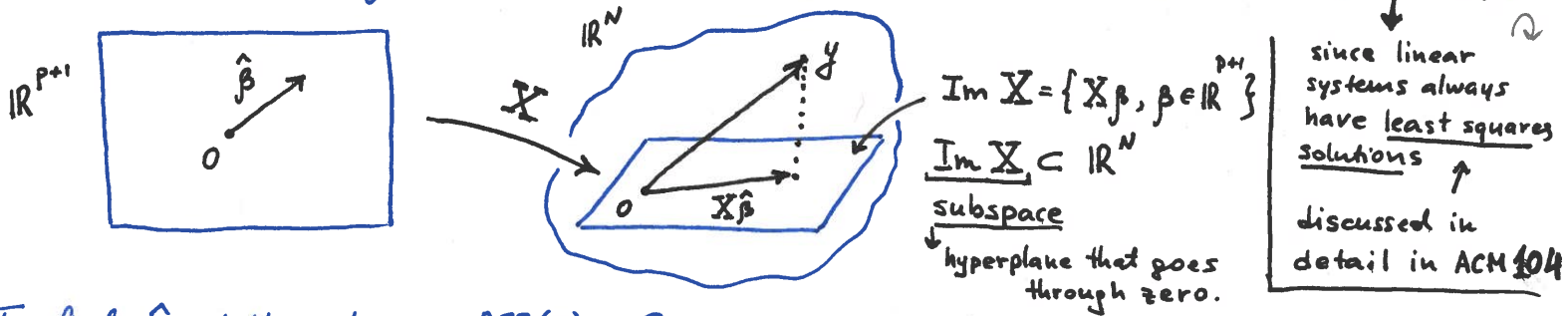
The value $\hat{\beta}$ of β that minimizes the RSS, $\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^{p+1}} \text{RSS}(\beta)$ has an important linear algebraic interpretation.

Let's rewrite RSS in the matrix form: $\text{RSS}(\beta) = \|y - X\beta\|^2$.

Therefore, $\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^{p+1}} \text{RSS}(\beta) = \arg \min_{\beta \in \mathbb{R}^{p+1}} \|X\beta - y\| \leftarrow \text{Euclidean norm}$

$\Rightarrow \hat{\beta}$ is a solution of linear system $X\beta = y$ in the least squares sense.

This interpretation of $\hat{\beta}$ allows to conclude that a global minimizer of RSS exists.



To find $\hat{\beta}$, let's solve $\nabla_{\beta} \text{RSS}(\beta) = 0$.

$$\text{RSS}(\beta) = (y - X\beta)^T (y - X\beta) = y^T y - \beta^T X^T y - y^T X \beta + \beta^T X^T X \beta$$

$$\nabla_{\beta} \text{RSS}(\beta) = -\nabla_{\beta} \beta^T X^T y - \nabla_{\beta} y^T X \beta + \nabla_{\beta} \beta^T X^T X \beta$$

$$\bullet \beta^T X^T y = \sum_{j=0}^p \beta_j (X^T y)_j \Rightarrow \nabla_{\beta} \beta^T X^T y = X^T y$$

$$\bullet y^T X \beta = \sum_{j=0}^p (y^T X)_j \beta_j \Rightarrow \nabla_{\beta} y^T X \beta = (y^T X)^T = X^T y$$

$$\bullet \text{Let } X^T X = A \leftarrow \begin{matrix} (p+1) \times (p+1) \\ \text{matrix} \end{matrix} \Rightarrow (\nabla_{\beta} \beta^T X^T X \beta)_k = \frac{\partial}{\partial \beta_k} \left(\sum_{i,j=0}^p a_{ij} \beta_i \beta_j \right)$$

$$= \sum_{i,j=0}^p a_{ij} \frac{\partial \beta_i}{\partial \beta_k} \beta_j + \sum_{i,j=0}^p a_{ij} \beta_i \frac{\partial \beta_j}{\partial \beta_k} = \sum_{j=0}^p a_{kj} \beta_j + \sum_{i=0}^p a_{ik} \beta_i$$

$= \delta_{ik} = \begin{cases} 1 & i=k \\ 0 & i \neq k \end{cases}$

$= \delta_{jk} = \begin{cases} 1 & j=k \\ 0 & j \neq k \end{cases}$

a_{ki} since A is symmetric

$$= 2 \sum_{j=0}^p a_{kj} \beta_j = 2 (A\beta)_k = 2 (X^T X \beta)_k \Rightarrow \nabla_{\beta} \beta^T X^T X \beta = 2 X^T X \beta$$

$$\text{So, } \nabla_{\beta} \text{RSS}(\beta) = -2 X^T y + 2 X^T X \beta.$$

Therefore, $\hat{\beta}$ is a solution of $X^T X \beta = X^T y$ \leftarrow normal equation (linear system of equations on β_0, \dots, β_p)

The normal equation always has a solution.

The $(p+1) \times (p+1)$ matrix $X^T X$ is the Gram matrix associated with $x^{(0)}, \dots, x^{(p)}$. (21)

It is always positive semidefinite, $X^T X \geq 0$.

① Suppose $X^T X$ is nonsingular $\Leftrightarrow X^T X > 0$ is positive definite

Then the normal equation has the unique solution

$\Leftrightarrow x^{(0)}, \dots, x^{(p)}$ are linearly indepen.

$\Leftrightarrow \ker X = \{\beta : X\beta = 0\} = \{0\}$

$\Leftrightarrow \text{rank } X = p+1$

columns of X

$$x^{(j)} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

$x^{(j)}$ consists of N observations of input X_j .

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

typically happens in applications.

To show that $\hat{\beta}$ is the global minimum, we need to check that the Hessian

$$H(RSS) = \left(\frac{\partial^2 RSS}{\partial \beta_s \partial \beta_k} \right)$$

matrix of second order partial derivatives.

of RSS at $\hat{\beta}$ is positive definite, $H(RSS) > 0$.

$$\bullet \frac{\partial RSS}{\partial \beta_k} = (\nabla_{\beta} RSS)_k = (2X^T X \beta - 2X^T y)_k = 2 \sum_{j=0}^p (X^T X)_{kj} \beta_j - (2X^T y)_k$$

$$\bullet \frac{\partial^2 RSS}{\partial \beta_s \partial \beta_k} = \frac{\partial}{\partial \beta_s} \left(2 \sum_{j=0}^p (X^T X)_{kj} \beta_j \right) = 2 \sum_{j=0}^p (X^T X)_{kj} \delta_{js} = 2 (X^T X)_{ks}$$

Therefore $H(RSS) = 2X^T X > 0 \Rightarrow$

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

is indeed the global minimum of RSS .

② Suppose $X^T X$ is singular.

This can happen if does not happen often in applications

$$\bullet N < p+1 \Rightarrow \text{rank } X < p+1$$

$X = \boxed{}$ is "flat"
#observations \leq #predictors

It is called the ordinary least squares (OLS) estimate of the regression parameter $\beta = (\beta_0 \dots \beta_p)^T$

deserves to be boxed twice on the same page

There are redundant predictors: one X_j is a linear combination of other predictors.

In this case, the normal equation has infinitely many solutions, all delivering the smallest possible value of RSS .

Example: Let $p=3$, $X_3 = X_1 + X_2 \Rightarrow y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 = \beta_0 + (\beta_1 + a)X_1 + (\beta_2 + a)X_2 + (\beta_3 - a)X_3$.

Then if $(\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2, \hat{\beta}_3)$ is a global minimizer of RSS ,

then so is $(\hat{\beta}_0, \hat{\beta}_1 + a, \hat{\beta}_2 + a, \hat{\beta}_3 - a)$ for any $a \in \mathbb{R}$.

From linear algebra we know that the general solution of a linear system $Ax = b$ is $x^* + \ker A$, where x^* is a particular solution and $\ker A = \{x : Ax = 0\}$.

The general solution of the normal equation is then $\hat{\beta} = \hat{\beta}^* + \ker(X^T X)$

It can be shown that $\ker(X^T X) = \ker X$:

A particular solution $\hat{\beta}^*$ can be found using the notion of pseudoinverse matrix, which is based on the king of matrix decompositions: the singular value decomposition (SVD)

- 1) $\beta \in \ker X \Rightarrow X\beta = 0 \Rightarrow X^T X\beta = 0 \Rightarrow \beta \in \ker(X^T X)$
- 2) $\beta \in \ker(X^T X) \Rightarrow X^T X\beta = 0 \Rightarrow \beta^T X^T X\beta = 0 \Rightarrow (X\beta)^T \cdot X\beta = 0 \Rightarrow \|X\beta\|^2 = 0 \Rightarrow X\beta = 0 \Rightarrow \beta \in \ker X$

Thm Let A be an $m \times n$ matrix of rank $r > 0$. Then it can be decomposed as follows :

$$A = P \cdot \Sigma \cdot Q^T \leftarrow \text{SVD of } A$$

- P is $m \times r$ with orthonormal columns, $P^T P = I_r$
- Σ is $r \times r$ diagonal, $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_r)$ singular values of A
- Q is $n \times r$ with orthonormal columns, $Q^T Q = I_r$

Remark Since the 1st column of X is $\begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$, its rank ≥ 1 . \Rightarrow SVD of X exists.

The SVD allows to generalize the notion of inverse matrix.

Def Let A be an $m \times n$ matrix and $A = P \Sigma Q^T$ be its SVD.

The pseudoinverse (aka Moore - Penrose inverse) is $A^+ = Q \Sigma^{-1} P^T$

It can be shown that

- If A is nonsingular $\Rightarrow A^+ = A^{-1}$
 - If $\ker A = \{0\} \Rightarrow A^+ = (A^T A)^{-1} A^T$
- Proved in ACM 104 (direct check)

Let $X = P \Sigma Q^T$ be the SVD of the training inputs matrix.

Consider $\hat{\beta}^* = X^+ y$. Then $\hat{\beta}^*$ is a solution of the normal equation $X^T X \beta = X^T y$

Indeed :

- $X^T X \hat{\beta}^* = Q \Sigma P^T P \Sigma Q^T Q \Sigma^{-1} P^T y = Q \Sigma P^T y$
- $X^T y = Q \Sigma P^T y$

So, if $X^T X$ is singular, then the general solution of the normal equation is

$$\hat{\beta} = X^+ y + \ker X$$

Combining the two cases together :

The general solution of the normal equation is

$$\hat{\beta} = \begin{cases} X^+ y = (X^T X)^{-1} X^T y & \text{if } X^T X \text{ is nonsing} \\ X^+ y + \ker X & \text{if } X^T X \text{ is sing} \Rightarrow \ker X \neq \{0\} \end{cases}$$

only if $\ker X = \{0\}$ $\ker X = \{0\}$

Proved in ACM 104