

Summary

1. Regression model / function

⇒ Linear model of conditional mean

2. Estimation of linear model

2.1 Simple linear regression (SLR)

[SLR1] Model.

[SLR2] Definition of least-squares loss function.

[SLR3] Ordinary least-squares estimates and alternate forms.
(Derivation through differentiation = 0)

[SLR4] Show that when OLS estimate is the
unique minimizer of least-squares loss.

① Two key properties of convex functions

(do not require proof, only need to use these conclusions)

(1.1) A convex function $f: \mathbb{R}^d \rightarrow \mathbb{R}$ has a unique
global minimizer that is the stationary
point of f (i.e. \bar{x} such that $\nabla f(\bar{x}) = 0$)

(1.2) Twice-differentiable f

is strictly convex if its Hessian is positive definite;
is convex if its Hessian is semipositive definite.

② Quadratic function $v^T A v$ (A is symmetric w.l.g.)

is strictly convex if A is positive definite;
is convex if A is semipositive definite.

- ③ Show when the least-squares loss function is convex or strictly convex by deriving the Hessian matrix.

[SLR5] Properties of OLS estimates

- ① Identities on sample residuals $\sum_{i=1}^n \hat{R}_i = \sum_{i=1}^n \hat{R}_i X_i = 0$
- ② Data center
- ③ Relationship with Pearson correlation
- ④ Non-symmetry in X & Y

2.2 Multiple linear regression (MLR)

[MLR1] Model.

[MLR2] Ordinary least-squares estimates

- (1) Least-squares loss function ✓
- (2) Derivation through multivariate differentiation ✓
- (3) Numerical OLS solution via QR decomposition ✓

Fall 2020 Qualifying Exam - Option B

3. This problem investigates fixed design linear regressions in high-dimensional settings. Suppose that we observe a sample of n observations, $(\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$ for $i = 1, \dots, n$. Let $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n]^T \in \mathbb{R}^{n \times d}$ denote the design matrix and $\mathbf{y} = (y_1, \dots, y_n)^T \in \mathbb{R}^n$ be the response vector. Consider a linear model with i.i.d. mean-zero Gaussian noise

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{n \times n}), \quad (1)$$

where $\boldsymbol{\beta} \in \mathbb{R}^d$ is the unknown coefficient vector, $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^T \in \mathbb{R}^n$ is the noise term, and $\mathbf{I}_{n \times n}$ is an n -by- n identity matrix. Many modern applications of this model are high-dimensional, in that the number of features d is comparable to, or even larger than, the sample size n . Assume that $d = n$, and \mathbf{X} has orthonormal columns such that $\mathbf{X}^T \mathbf{X} = \mathbf{I}_{d \times d}$. Consider the following regularized estimator for $\boldsymbol{\beta}$,

$$\hat{\boldsymbol{\beta}}_\lambda = \arg \min_{\boldsymbol{\beta}} \{ \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda \|\boldsymbol{\beta}\|^2 \}, \quad (2)$$

where λ is an unknown tuning parameter, $\|\cdot\|$ denotes the vector 2-norm; i.e., $\|\mathbf{a}\| = \left(\sum_{j=1}^d |a_j|^2 \right)^{1/2}$ for a vector $\mathbf{a} = (a_1, \dots, a_d)^T \in \mathbb{R}^d$.

(a) Let $\lambda = 0$. (2) $\Rightarrow \hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$ ✓ when $\lambda = 0$

- i. Derive the expression for $\hat{\boldsymbol{\beta}}_0$ by solving the optimization problem in (2) and find its distribution.
- ii. Consider the prediction error for a new observation of the form $y_{\text{new}} = \mathbf{x}_{\text{new}}^T \boldsymbol{\beta} + \varepsilon_{\text{new}}$, where $\mathbf{x}_{\text{new}} \in \mathbb{R}^d$ is a fixed covariate vector and $\varepsilon_{\text{new}} \sim \mathcal{N}(0, 1)$ is a noise term independent of $\{\varepsilon_i\}_{i=1}^n$. Find the expected squared prediction error, $\mathbb{E}(y_{\text{new}} - \mathbf{x}_{\text{new}}^T \hat{\boldsymbol{\beta}})^2$, for this new observation.

(b) Let $\lambda > 0$.

- i. Derive the expression for ridge regression estimator $\hat{\boldsymbol{\beta}}^{\text{ridge}}$ by solving the optimization problem in (2).
 - ii. Consider the prediction error for a new observation of the form $y_{\text{new}} = \mathbf{x}_{\text{new}}^T \boldsymbol{\beta} + \varepsilon_{\text{new}}$, where $\mathbf{x}_{\text{new}} \in \mathbb{R}^d$ is a fixed covariate vector and $\varepsilon_{\text{new}} \sim \mathcal{N}(0, 1)$ is a noise term independent of $\{\varepsilon_i\}_{i=1}^n$. Find the expected squared prediction error, $\mathbb{E}(y_{\text{new}} - \mathbf{x}_{\text{new}}^T \hat{\boldsymbol{\beta}}^{\text{ridge}})^2$, for this new observation.
 - iii. For this part of the question only, assume that $\|\boldsymbol{\beta}\| = 1$. Derive the optimal λ that minimizes the mean squared error for the ridge estimator, $\mathbb{E}\|\hat{\boldsymbol{\beta}}^{\text{ridge}} - \boldsymbol{\beta}\|^2$. Discuss how you would find λ in practice when $\|\boldsymbol{\beta}\|$ is unknown.
- (c) Suppose that a prior distribution $\boldsymbol{\beta}^{\text{prior}} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \Phi)$ is imposed to the model (1), where σ^2 is an unknown variance parameter, and Φ is a known positive definite matrix. Furthermore, assume that $\boldsymbol{\beta}^{\text{prior}}$ and $\boldsymbol{\varepsilon}$ are independent.
- i. Find the marginal distribution of \mathbf{y} .
 - ii. Derive the method-of-moments estimator for σ^2 based on the marginal distribution of \mathbf{y} .

When X has full column rank, \Rightarrow [Assumed all below.]

① OLS algebraic solution $\hat{\beta} = (X^T X)^{-1} X^T Y$

② OLS numerical solution through $X = QR$

$$\Rightarrow \hat{\beta} = R^{-1} Q^T Y$$

Exercise 1:

(a) Express hat matrix $X(X^T X)^{-1} X^T$ through $X = QR$

(b) $Q = \begin{pmatrix} 0 & \frac{1}{\sqrt{2}} \\ 0 & \frac{1}{\sqrt{2}} \\ 1 & 0 \end{pmatrix}$ Derive QQ^T and $Q^T Q$

(c) Solve $\begin{pmatrix} 2 & 1 & 1 \\ 0 & 1 & 2 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix} = \begin{pmatrix} 2 \\ 4 \\ 1 \end{pmatrix} \Rightarrow \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \hat{\beta}_3 \end{pmatrix}$

Solutions

(a) $X = QR$ ① $Q^T Q = I_p$ ② R upper triangular matrix

$$\text{Hat matrix} = X(X^T X)^{-1} X^T$$

$$= QR \left((QR)^T QR \right)^{-1} (QR)^T \quad (AB)^T = B^T A^T$$

$$= QR \left(R^T \underbrace{Q^T Q}_{I_p} R \right)^{-1} R^T Q^T$$

$$= QR (R^T R)^{-1} R^T Q^T$$

$$= \underbrace{QR R^{-1}}_I \times \underbrace{(R^T)^{-1} R^T}_I Q^T$$

$$= QQ^T \quad (\text{NOT NECESSARILY IDENTITY})$$

$$\underbrace{\begin{matrix} n \times p & p \times n \\ & n \times n \end{matrix}}$$

$$(Q^T Q = I_p)$$

X has full column rank
 R invertible
 R^T invertible
 $(R^T R)^{-1} = R^{-1} \cdot (R^T)^{-1}$

$$(b) \quad Q = \begin{pmatrix} 0 & \frac{1}{\sqrt{2}} \\ 0 & \frac{1}{\sqrt{2}} \\ -1 & 0 \end{pmatrix}$$

$$\underbrace{Q^T Q}_{\substack{2 \times 3 \quad 3 \times 2 \\ 2 \times 2}} = \begin{pmatrix} 0 & 0 & -1 \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & 0 \end{pmatrix} \begin{pmatrix} 0 & \frac{1}{\sqrt{2}} \\ 0 & \frac{1}{\sqrt{2}} \\ -1 & 0 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = I_{2 \times 2}$$

$$\underbrace{Q Q^T}_{\substack{3 \times 2 \quad 2 \times 3 \\ 3 \times 3}} = \begin{pmatrix} 0 & \frac{1}{\sqrt{2}} \\ 0 & \frac{1}{\sqrt{2}} \\ -1 & 0 \end{pmatrix} \begin{pmatrix} 0 & 0 & -1 \\ \frac{1}{\sqrt{2}} & \frac{1}{\sqrt{2}} & 0 \end{pmatrix} \\ = \begin{pmatrix} \frac{1}{2} & \frac{1}{2} & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \\ 0 & 0 & 1 \end{pmatrix} \neq \text{Identity matrix}$$

$$(c) \quad \begin{pmatrix} 2 & 1 & 1 \\ 0 & 1 & 2 \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix} = \begin{pmatrix} 2 \\ 4 \\ 1 \end{pmatrix}$$

Upper triangular 3×3

$$\Leftrightarrow \begin{cases} 2\beta_1 + \beta_2 + \beta_3 = 2 \\ \beta_2 + 2\beta_3 = 4 \\ \beta_3 = 1 \end{cases} \Rightarrow \begin{cases} 2\beta_1 + 2 + 1 = 2 \Rightarrow \beta_1 = -\frac{1}{2} \\ \beta_2 + 2 \times 1 = 4 \Rightarrow \beta_2 = 2 \\ \beta_3 = 1 \end{cases}$$

From bottom \Leftarrow

Numerical OLS Normal equation \Rightarrow Upper triangular system

Exercise 2

$$X^T X$$

Two ways to think about this matrix product

(i) By row (observations)

$$X_{n \times p} = \begin{pmatrix} \text{---} x_{(1)}^T \text{---} \\ \text{---} x_{(2)}^T \text{---} \\ \vdots \\ \text{---} x_{(n)}^T \text{---} \end{pmatrix} \quad X^T = \begin{pmatrix} | & | & \dots & | \\ x_{(1)} & x_{(2)} & \dots & x_{(n)} \\ | & | & & | \end{pmatrix}$$

$(x_{(j)} \in \mathbb{R}^p, \text{ by default } x_{ij} \text{ are column vectors})$
 $j = 1 \dots n$

(2) By columns (covariates)

$$X_{n \times p} = \begin{pmatrix} | & | & \dots & | \\ x^{(1)} & x^{(2)} & \dots & x^{(p)} \\ | & | & & | \end{pmatrix} \quad X^T_{p \times n} = \begin{pmatrix} \text{---} x^{(1)T} \text{---} \\ \text{---} x^{(2)T} \text{---} \\ \vdots \\ \text{---} x^{(p)T} \text{---} \end{pmatrix}$$

$x^{(j)} \in \mathbb{R}^n \quad j = 1, \dots, p$

Exercise: Express $X^T X$ under two views

Solutions:

$$\textcircled{1} \quad \mathbf{X}^T \mathbf{X} = \begin{pmatrix} | & & | \\ \mathbf{x}_{(1)} & \dots & \mathbf{x}_{(n)} \\ | & & | \end{pmatrix} \begin{pmatrix} - \mathbf{x}_{(1)}^T - \\ \vdots \\ - \mathbf{x}_{(n)}^T - \end{pmatrix}$$

$$= \sum_{i=1}^n \mathbf{x}_{(i)} \mathbf{x}_{(i)}^T \Rightarrow \text{sum of } n \text{ } p \times p \text{ matrices}$$

$$\begin{matrix} p \times 1 & 1 \times p \\ \hline p \times p \end{matrix}$$

$$\textcircled{2} \quad \mathbf{X}^T \mathbf{X} = \begin{pmatrix} - \mathbf{x}^{(1)T} - \\ - \mathbf{x}^{(2)T} - \\ \vdots \\ - \mathbf{x}^{(p)T} - \end{pmatrix} \begin{pmatrix} | & | & & | \\ \mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \dots & \mathbf{x}^{(p)} \\ | & | & & | \end{pmatrix}$$

$$= \begin{pmatrix} \underbrace{\mathbf{x}^{(1)T} \mathbf{x}^{(1)}}_{1 \times 1} & \underbrace{\mathbf{x}^{(1)T} \mathbf{x}^{(2)}}_{1 \times 1} & \dots & \underbrace{\mathbf{x}^{(1)T} \mathbf{x}^{(p)}}_{1 \times 1} \\ \mathbf{x}^{(2)T} \mathbf{x}^{(1)} & \mathbf{x}^{(2)T} \mathbf{x}^{(2)} & \dots & \mathbf{x}^{(2)T} \mathbf{x}^{(p)} \\ \vdots & \vdots & \ddots & \vdots \\ \underbrace{\mathbf{x}^{(p)T} \mathbf{x}^{(1)}}_{1 \times 1} & \underbrace{\mathbf{x}^{(p)T} \mathbf{x}^{(2)}}_{1 \times 1} & \dots & \underbrace{\mathbf{x}^{(p)T} \mathbf{x}^{(p)}}_{1 \times 1} \end{pmatrix}_{p \times p}$$

$$= \left(\mathbf{x}^{(i)T} \mathbf{x}^{(j)} \right)_{1 \leq i, j \leq p} \quad (p \times p \text{ matrix})$$

(Use for HW1-Q3)

Exercise 3

Example: Suppose Y_1, \dots, Y_n have common mean β

The least square estimate $\min_{\beta} \underbrace{\sum_{i=1}^n (Y_i - \beta)^2}_{L(\beta)}$

① Specific case derivative $\frac{\partial L(\beta)}{\partial \beta} = -2 \sum_{i=1}^n (Y_i - \beta) = 0 \Rightarrow \hat{\beta} = \frac{\sum_{i=1}^n Y_i}{n} = \bar{Y}$
(Quadratic, unique minimizer)

② General Alternatively, we can linear regression model as

Plug-in design matrix $Y = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}_{n \times 1}$ $E(Y) = \begin{pmatrix} \beta \\ \beta \\ \vdots \\ \beta \end{pmatrix}_{n \times 1} = \mathbf{1}_n \times \beta$ $\mathbf{1}_n = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}_{n \times 1}$

$X = \mathbf{1}_n$
 $n \times 1$

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

$= (\mathbf{1}_n^T \mathbf{1}_n)^{-1} \mathbf{1}_n^T Y$
 $1 \times n \quad n \times 1$

$= n^{-1} \times \mathbf{1}_n^T Y = \frac{1}{n} \sum_{i=1}^n Y_i = \bar{Y} \quad \hat{\beta} = \bar{Y}$

These two views should be equivalent.

Exercise

$$E(Y|x) = \alpha + \beta x \quad \text{in SLR}$$

$$(\hat{\alpha}, \hat{\beta})$$

$$\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

↗ intercept

equivalent

Hint: (1) $\mathbf{X} = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}$

$n \times 2$ ↘ 1-dim covariate

(2) $\mathbf{X}^T \mathbf{X}$ 2×2 matrix

$2 \times n$ $n \times 2$

(3) Inverse: a general invertible 2×2 matrix

$$\begin{pmatrix} a & b \\ c & d \end{pmatrix}^{-1} = \frac{1}{ad-bc} \begin{pmatrix} d & -b \\ -c & a \end{pmatrix} \quad (\text{if } ad \neq bc)$$

2×2

Solution

(1) $\mathbf{X}^T \mathbf{X}$

$$= \begin{pmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_n \end{pmatrix} \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}$$

$$= \begin{pmatrix} n & \sum_{i=1}^n x_i \\ \sum_{i=1}^n x_i & \sum_{i=1}^n x_i^2 \end{pmatrix} = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

$$(2) \quad (X^T X)^{-1}$$

$$= \frac{1}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2} \begin{pmatrix} \sum_{i=1}^n x_i^2 & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & n \end{pmatrix}$$

$$(3) \quad \underbrace{X^T}_{2 \times n} \underbrace{Y}_{n \times 1}$$

$$= \begin{pmatrix} 1 & \dots & 1 \\ x_1 & \dots & x_n \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_i y_i \end{pmatrix}$$

$$(4) \quad \underbrace{\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix}}_{2 \times 1} = \underbrace{(X^T X)^{-1}}_{2 \times 2} \underbrace{X^T Y}_{2 \times 1}$$

$$= \frac{1}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2} \begin{pmatrix} \sum_{i=1}^n x_i^2 & -\sum_{i=1}^n x_i \\ -\sum_{i=1}^n x_i & n \end{pmatrix} \times \begin{pmatrix} \sum_{i=1}^n y_i \\ \sum_{i=1}^n x_i y_i \end{pmatrix}$$

$$= \frac{1}{n \sum_{i=1}^n x_i^2 - \left(\sum_{i=1}^n x_i \right)^2} \begin{pmatrix} \sum_{i=1}^n x_i^2 \times \sum_{i=1}^n y_i - \sum_{i=1}^n x_i \sum_{i=1}^n x_i y_i \\ -\sum_{i=1}^n x_i \sum_{i=1}^n y_i + n \sum_{i=1}^n x_i y_i \end{pmatrix}$$

Use $\sum_{i=1}^n Y_i = n \bar{Y}$ $\sum_{i=1}^n X_i = n \bar{X}$

$$= \frac{1}{n \sum_{i=1}^n X_i^2 - \left(\sum_{i=1}^n X_i \right)^2} \begin{pmatrix} \sum_{i=1}^n X_i^2 \times n \bar{Y} - n \bar{X} \cdot \sum_{i=1}^n X_i Y_i \\ - n^2 \bar{X} \bar{Y} + n \sum_{i=1}^n X_i Y_i \end{pmatrix}$$

Equivalent to

$$= \begin{pmatrix} \bar{Y} - \hat{\beta} \bar{X} \\ \frac{\sum_{i=1}^n X_i Y_i - n \bar{X} \bar{Y}}{\sum_{i=1}^n X_i^2 - n \bar{X}^2} \end{pmatrix}$$

OLS estimates $\begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix}$

by differentiation

2.4 Unique minimizer?

The LS solution of MLR is unique as long as the columns of X are linearly independent.

(has full column rank)

Proof: (1) Hessian of $L(\beta) = \|Y - X\beta\|^2$

$$\frac{\partial L(\beta)}{\partial \beta} = -2 X^T (Y - X^T \beta)$$

$$p \times 1 \quad = \quad 2 X^T X \beta - \underbrace{2 X^T Y}$$

$$\text{Hessian} \quad \frac{\partial^2 L(\beta)}{\partial \beta \partial \beta^T} = 2 X^T X - 0 \quad (\text{Quadratic})$$

$p \times p$

(2) Hessian is positive definite $\Rightarrow L(\beta)$ is strictly convex and has a unique global minimizer

(Remains to show)

\Rightarrow Hessian matrix is p.d. if columns of X are linearly independent

For $v \neq 0 \in \mathbb{R}^p$,

$$\underbrace{v^T}_{1 \times p} \underbrace{(X^T X)}_{p \times n \quad n \times p} \underbrace{v}_{p \times 1} = \underbrace{(Xv)^T}_{1 \times 1} (Xv) \quad (Xv)^T = v^T X^T$$
$$= \|Xv\|^2$$

Since X 's columns are linearly independent,

$$\Rightarrow Xv \neq 0 \Rightarrow \|Xv\|^2 > 0$$

Therefore, the Hessian of $L(\beta): 2 X^T X$ is positive definite.