

II. Least Squares Estimation

Data $(Y_i, X_i^T)^T, i=1, \dots, n$ not necessarily iid

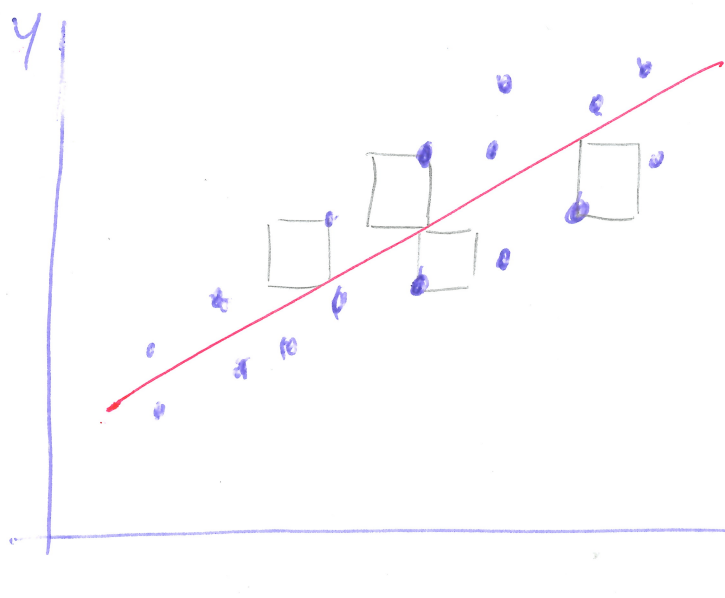
$$Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} \quad X = \begin{pmatrix} X_1^T \\ \vdots \\ X_n^T \end{pmatrix} \quad n = \begin{pmatrix} X^0, \dots, X^{k-1} \end{pmatrix}$$

regressors

$$Y|X \sim (X\beta, \sigma^2 I) \quad , \text{rank}(X) = \underline{r=k} < n$$

unknown parameters (almost surely)

2.1 Sum of squares, LSE & normal equations



MODEL: $E(Y|z=z) = \beta_0 + \beta_1 z = X^T \beta$
 $X = (1, z)^T$

$\hat{\beta} = ?$

TAKE $\hat{\beta} = \underset{\beta}{\operatorname{argmin}} \underbrace{\sum_{i=1}^n (Y_i - X_i^T \beta)^2}_{SS(\beta)}$

Def. 2.1 Sum of squares

Consider a linear model $Y|X \sim (X\beta, \sigma^2 I_n)$.

The function $SS: \mathbb{R}^k \rightarrow \mathbb{R}$ given as

$$SS(\beta) = \sum_{i=1}^n (Y_i - X_i^T \beta)^2 = \|Y - X\beta\|^2 = (Y - X\beta)^T (Y - X\beta),$$

$\beta \in \mathbb{R}^k$ will be called the sum of squares of the model.

Lemma 2.1 Least squares estimator | Let $r=k$.

There exist a unique minimizer to $SS(\beta)$ given as $\hat{\beta} := (X^T X)^{-1} (X^T Y)$.

Proof: $SS(\beta) = \|Y - X\beta\|^2 = (Y - X\beta)^T (Y - X\beta) =$
 $= Y^T Y - 2Y^T X\beta + \beta^T X^T X \beta.$

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

$$\frac{\partial SS}{\partial \beta}(\beta) = 0_k \Leftrightarrow \underbrace{X^T X}_{k \times k, \text{ rank}(X^T X) = k} \beta = X^T Y$$

$$\Leftrightarrow \beta = \underbrace{(X^T X)^{-1}}_{\text{always exists, it is unique}} X^T Y$$

Do we have a minimum?

$$\frac{\partial^2 SS}{\partial \beta \partial \beta^T}(\beta) = \underline{2X^T X}$$

\hookrightarrow for any $\beta \in \mathbb{R}^k$ this is positive definite matrix,

hence $\hat{\beta} := (X^T X)^{-1} X^T Y$ is a unique minimizer to $SS(\beta)$ □

Note: $\hat{\beta}$ solves a linear system

$$X^T X \beta = X^T Y.$$

Def 2.2 Least squares estimator, normal equations

Consider a linear model $\forall X \sim (X\beta, \sigma^2 I_n)$,
 $\text{rank}(X_{n \times k}) = k$. The quantity $\hat{\beta} := (X^T X)^{-1} X^T Y$
will be called the least squares estimator
(LSE) of the vector of regression coeffs. β .
The linear system $X^T X \beta = X^T Y$ will
be called the system of normal equations.

Lemma 2.2 Moments of LSE

Let $Y|X \sim (X\beta, \sigma^2 I_n)$, $\text{rank}(X_{n \times k}) = k < n$.
Then $E(\hat{\beta}|X) = \beta$, $E(\hat{\beta}) = \beta$,
 $\text{var}(\hat{\beta}|X) = \sigma^2 (X^T X)^{-1}$.

Proof: $E(\hat{\beta}|X) = E((X^T X)^{-1} X^T Y | X) =$
 $= (X^T X)^{-1} X^T \underbrace{E(Y|X)}_{X\beta} = (X^T X)^{-1} X^T X \beta = \beta$.

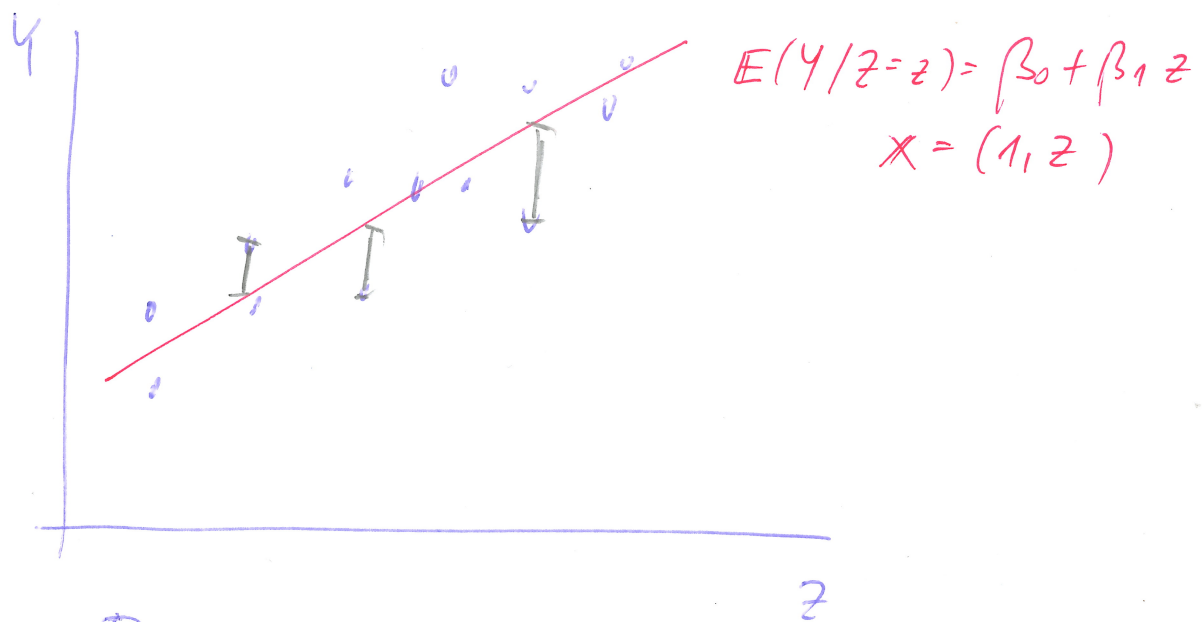
$$E(\hat{\beta}) = E(E(\hat{\beta}|X)) = E(\beta) = \beta$$

$$\text{var}(\hat{\beta}|X) = \text{var}((X^T X)^{-1} X^T Y | X) =$$
$$= (X^T X)^{-1} X^T \underbrace{\text{var}(Y|X)}_{\sigma^2 I_n} X (X^T X)^{-1} =$$

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

Remark: LSE $\hat{\beta}$ is unbiased estimator of β ,
both conditionally (given X) and
unconditionally. \rightarrow what does this
mean?

Remark: Other methods exist on how
to estimate β .



For example:

$$\tilde{\beta} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n |y_i - X_i^T \beta|$$

\equiv L_1 regression

Will not be covered by this lecture.

2.2 Fitted values, residuals, projections

Notation: For a matrix $X = (\underbrace{X^0, \dots, X^{k-1}}_k)_{n \times k}$

$$\mathcal{M}(X) := \left\{ \underset{\sim}{v} : \underset{\sim}{v} = \underbrace{\sum_{j=0}^{k-1} b_j X_j^0}_{= Xb}, b = (b_0, \dots, b_{k-1})^T \in \mathbb{R}^k \right\}$$

\equiv linear span of columns of X

$$\mathcal{M}(X)^\perp := \left\{ \underset{\sim}{u} : u \in \mathbb{R}^n, v^T u = 0 \ \forall v \in \mathcal{M}(X) \right\}$$

\equiv orthogonal complement to $\mathcal{M}(X)$

$\&$ $X_{n \times k}$: $\text{rank}(X) = r \leq k < n$
(column)

$$\begin{array}{l} \mathcal{M}(X) \subset \mathbb{R}^n \\ \mathcal{M}(X)^\perp \subset \mathbb{R}^n \end{array} \quad \left. \begin{array}{l} \text{vector dimension } r \\ n-r \end{array} \right\} \text{a.s.}$$

$$\begin{array}{l} \mathcal{M}(X) \cap \mathcal{M}(X)^\perp = \{0_n\} \\ \cup \\ = \mathbb{R}^n \end{array}$$

Def. 2.3 Regression and residual space of a linear model

Consider a linear model $Y|X \sim (X\beta, \sigma^2 I_n)$,
 $\text{rank}(X_{n \times k}) = r \leq k < n$. The regression space
of the model is a vector space $\mathcal{M}(X)$.

The residual space of the model is the
orthogonal complement of the regression space,
i.e., a vector space $\mathcal{M}(X)^\perp$.

Def 2.4 Fitted values, residuals

Consider a linear model $Y|X \sim (X\beta, \sigma^2 I_n)$,
 $\text{rank}(X_{n \times k}) = k < n$. The vector

$\hat{Y} := X\hat{\beta} = X(X^T X)^{-1} X^T Y$ will be called

the vector of fitted values of the model,

the vector $U := Y - \hat{Y}$ will be called

the vector of residuals of the model.

Notes:

$$\hat{Y} = X\hat{\beta} \in \mathcal{K}(X)$$

$$\left(= \sum_{j=0}^{k-1} \hat{\beta}_j X^j \right)$$

at the same time $\hat{\beta} = \underset{\beta}{\text{argmin}} \|Y - X\beta\|^2$

That is, $\hat{Y} = X\hat{\beta}$ is the closest point (in L_2 norm)

to Y in the regression space $\mathcal{K}(X)$

Linear algebra: \hat{Y} is projection of Y to $\mathcal{K}(X)$
(with L_2 norm)

Since $\hat{Y} = X(X^T X)^{-1} X^T Y$, the matrix

$H := X(X^T X)^{-1} X^T$ must be the

respective projection matrix,

(which is unique, see LA).

Residuals: $U = Y - \hat{Y} = Y - H \cdot Y = (I_n - H)Y$

$$M = I_n - H$$

Basic properties of matrices H and M:

• $H \cdot H = X(X^T X)^{-1} X^T X (X^T X)^{-1} X^T = X(X^T X)^{-1} X^T = H$

$$M \cdot M = (I - H)(I - H) = I - H - H + \underbrace{H \cdot H}_H = I - H = M$$

Both H and M are idempotent.

• $H^T = (X(X^T X)^{-1} X^T)^T = X(X^T X)^{-1} X^T = H$

$$M^T = (I - H)^T = I - H^T = I - H = M$$

Both H and M are symmetric.

• $H \cdot M = H \cdot (I - H) = H - H \cdot H = H - H = \mathbb{O}_{n \times n}$

$$M \cdot H = \dots = \mathbb{O}_{n \times n}$$

\Rightarrow Fitted values and residuals are mutually perpendicular:

$$U^T \hat{Y} = (MY)^T H Y = Y^T \underbrace{M \cdot H}_{\mathbb{O}_{n \times n}} Y = 0$$

That is, $U \perp \hat{Y}$, where \hat{Y} is projection of Y into $\mathcal{N}(X)$.

LA: $U \in \mathcal{N}(X)^\perp$, $U = MY$ is projection of Y into $\mathcal{N}(X)^\perp$.

$$U = MY \quad (= (I - H)Y = Y - \hat{Y})$$

LA: projection matrix is unique, hence M must be the respective projection matrix into $\mathcal{N}(X)^\perp$.

$$\bullet \quad HX = \underline{X(X^T X)^{-1} X^T X} = X$$

$$MX = (I - H)X = X - HX = X - X = \mathbf{0}_{n \times k}$$

Lemma 2.3 Algebraic properties of fitted values and residuals and related projection matrices

- (i) $\hat{Y} = HY$ and $U = MY$ are projections of Y into $\mathcal{R}(X)$ and $\mathcal{N}(X)^\perp$ respectively.
- (ii) $\hat{Y} \perp U$ (almost surely).
- (iii) H and M are projection matrices into $\mathcal{R}(X)$ and $\mathcal{N}(X)^\perp$, respectively.
- (iv) $H^T = H, \quad M^T = M.$
- (v) $H \cdot H = H, \quad M \cdot M = M.$
- (vi) $H \cdot X = X, \quad MX = \mathbf{0}_{n \times k}.$

Proof: Previous derivations. □

Alternative expression of projection matrices

Now: not necessarily full-rank, i.e.

$$X_{n \times k}, \text{rank}(X) = r \leq k < n$$

• $\mathcal{U}(X) \equiv$ vector space, vec-dim = r

• $\mathcal{U}(X)^\perp \equiv \quad \quad \quad = n-r$

LA: There exist an orthonormal vector basis of \mathbb{R}^n :

$$n \left(\underbrace{q_1, \dots, q_r}_r \mid \underbrace{n_1, \dots, n_{n-r}}_{n-r} \right) = P$$

basis of $\mathcal{U}(X)$ $\quad \quad \quad \mathcal{U}(X)^\perp$

properties of orth. basis:

$$Q^T Q = I_r$$

$$N^T N = I_{n-r}$$

$$Q^T N = O_{r \times (n-r)}$$

$$N^T Q = O_{(n-r) \times r}$$

$$P^T P = I_n, \quad P^T \text{ is inverse to } P,$$

$$\Rightarrow P \cdot P^T = I_n$$

$$\begin{aligned} \left(\begin{array}{c} Q \\ N \end{array} \right)^T \cdot \begin{pmatrix} Q^T \\ N^T \end{pmatrix} &= Q \cdot Q^T + N N^T \\ &=: \tilde{H} \quad \quad \quad =: \tilde{M} \end{aligned}$$

$$\text{i.e. } I_n = \tilde{H} + \tilde{M}$$

Properties of \tilde{H} (and \tilde{M}):

$$\tilde{H}^T = (Q \cdot Q^T)^T = Q \cdot Q^T = \tilde{H} \quad (\text{symmetric})$$
$$\tilde{H} \cdot \tilde{H} = \underbrace{Q \cdot Q^T \cdot Q \cdot Q^T}_I = Q \cdot Q^T = \tilde{H} \quad (\text{idempotent})$$

(the same for \tilde{M})

$$\tilde{H} \cdot \tilde{M} = \underbrace{Q \cdot Q^T}_0 \cdot N \cdot N^T = 0$$

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

$$\tilde{H}y = Q \cdot \underbrace{Q^T y}_b = Q \cdot b = \sum_{j=1}^k b_j \cdot q_j \in \mathcal{U}(X)$$

$b = (b_1, \dots, b_k)^T$

$$\tilde{M}y = N \cdot \underbrace{N^T y}_c = N \cdot c \in \mathcal{U}(X)^\perp$$

$c = (c_1, \dots, c_{n-k})^T$

$$y = I \cdot y = (\tilde{H} + \tilde{M})y = \tilde{H}y + \tilde{M}y$$

$\in \mathcal{U}(X) \quad \in \mathcal{U}(X)^\perp$

LA: • projections are unique, i.e.

$\tilde{H}y$ must be the projection of y
into $\mathcal{U}(X)$

$\tilde{M}y$ — " — into $\mathcal{U}(X)^\perp$.

• also projection matrices are unique
and hence (in case $k=k$):

$$\boxed{Q \cdot Q^T = \tilde{H} = H = X(X^T X)^{-1} X^T}$$
$$\boxed{N \cdot N^T = \tilde{M} = M = I_n - X(X^T X)^{-1} X^T}$$

↑ They do not depend on choices
of the orthonormal bases.

Summary:

Fitted values \hat{Y} = projection of Y into $\mathcal{M}(X)$

i.e. $\hat{Y} = \underset{\tilde{Y} \in \mathcal{M}(X)}{\text{argmin}} \|Y - \tilde{Y}\|^2$

$$\hat{Y} = H \cdot Y$$

$$H = X(X^T X)^{-1} X^T \quad (\text{if } \text{rank}(X_{n \times k}) = k)$$

$$= Q \cdot Q^T$$

Q = orthonormal basis
of $\mathcal{M}(X)$

Residuals

$$U = Y - \hat{Y} =$$

= projection of Y into $\mathcal{M}(X)^\perp$

$$U = M \cdot Y$$

$$M = I_n - H =$$

$$= I_n - X(X^T X)^{-1} X^T \quad (\text{if } \text{rank}(X) = k)$$

$$= N \cdot N^T$$

N = orthonormal basis
of $\mathcal{M}(X)^\perp$

Terminology:

- H : hat matrix (projection matrix to the ~~the~~ regression space)
- M : residual projection matrix

Remarks: if $\text{rank}(X_{n \times k}) = k < k < n$

- $SS(\beta) = \|Y - X\beta\|^2$ is minimized for any solution to normal equations $X^T X \beta = X^T Y$.
- All solutions to normal equations are $b = (X^T X)^{-} X^T Y$
- Projection matrices can be expressed as $H = X(X^T X)^{-} X^T$, $M = I_n - X(X^T X)^{-} X^T$ where a pseudoinverse $(X^T X)^{-}$ is not unique, but matrices H and M are unique.
- $\hat{Y} = HY = X \underbrace{(X^T X)^{-} X^T Y}_{b} = Xb$ is unique even if b is not unique.
- $U = Y - \hat{Y}$ is also unique.