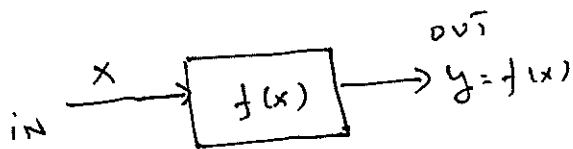
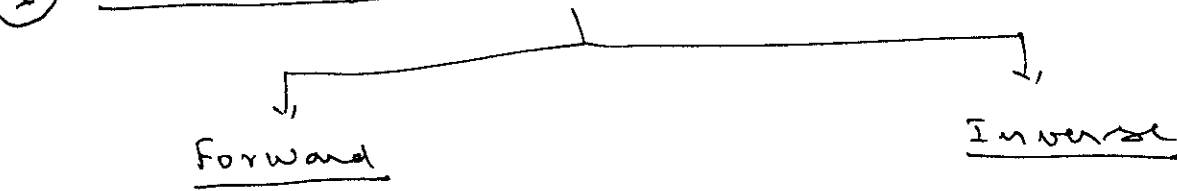
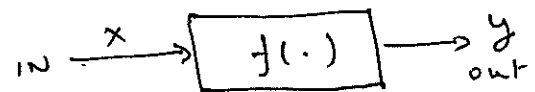


STATISTICAL LEAST SQUARESS. Lakshmi VarahanApril 30, 2022① A classification of problems:

- $f(x)$ is known
- x is given
- compute y



- x and y are given
- Identify $f(\cdot)$
- x - is x-ray, y - x-ray picture
- Identify $f(\cdot)$

② Simple pendulum: $T = 2\pi \sqrt{l/g}$

- If l and g are given, compute T - Forward
- If $(l_i, T_i), 1 \leq i \leq m$ given, estimate g - Inverse

③ Economics:

change in
unemployment
↑
~~observe it.~~
observe it.

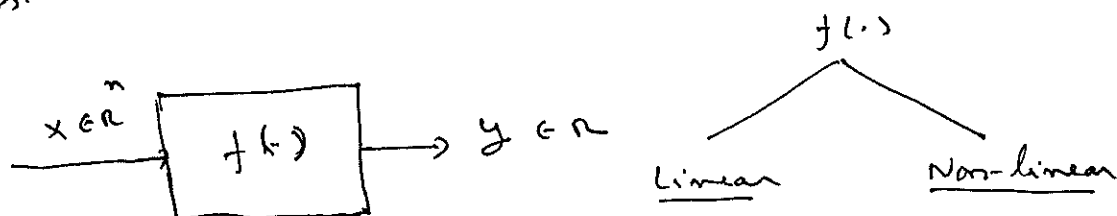
$= W_0 + W_1 * \text{change in GDP} + \text{noise}$
 want to discover the relation
 = output
 ↑
 we can observe it

(2)

① Satellite problem: Thermal energy radiated is proportional to the 4th power of the temperature of the radiating surface: $E = \alpha T^4$

- knowing E as measured by the satellite and α , estimate T .

② statistical
Least Squares method is one of the methods for modelling that is used in solving inverse problems.



W-parameter ① Example 1: $y = x^T w + \mathcal{N}$ $\mathcal{N} \sim \text{Noise} \sim \mathcal{N}(0, \sigma^2)$ \uparrow Known

- Here y and x and properties of noise \mathcal{N} are known. Find / estimate w .

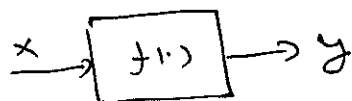
- Here $f(x) = x^T w$ is linear in x and w .

② Example 2: $y = a e^{bx} + \mathcal{N}$



- Here $f(x) = a e^{bx}$. It is nonlinear in the unknown (a, b) and also nonlinear in x .

③ Example 3: $y = a_0 + a_1 x + a_2 x^2 + \mathcal{N}$



- $f(x) = a_0 + a_1 x + a_2 x^2$. Nonlinear in x and linear in parameters (a_0, a_1, a_2) .

3

③ linear least squares problem : (statistical)

let $y_i = w_0 + x_{i1}w_1 + x_{i2}w_2 \dots + \overbrace{x_{i,n-1}w_{n-1}}^{\text{crossed out}} + v_i$
 \rightarrow ①

for $1 \leq i \leq m$, where assume $m > n$.

Define $y = (y_1, y_2, \dots, y_m)^T \in \mathbb{R}^m$

$w = (w_0, w_1, \dots, w_{n-1})^T \in \mathbb{R}^n$

$$X = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1,n-1} \\ 1 & x_{21} & x_{22} & \dots & x_{2,n-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m1} & x_{m2} & \dots & x_{m,n-1} \end{bmatrix} \in \mathbb{R}^{m \times n}$$

$v = [v_1, v_2, \dots, v_m]^T \in \mathbb{R}^m$

① becomes in matrix-vector notation:

$y = Xw + v$

where $v \sim N(0, \overbrace{R}^{\text{crossed out}})$, σ^2 is known, $R = \sigma^2 I_m$
 \rightarrow ②

Define residual

$r(w) = y - Xw$

④ Cost function : $J : \mathbb{R}^n \rightarrow \mathbb{R}$

$$J(w) = \frac{1}{2} (y - Xw)^T R^{-1} (y - Xw)$$

is the weighted sum of squared residuals.

Goal is to minimize $J(w)$ w.r. to w .

This is a multivariate minimization problem

- Rewrite $J(w)$ by multiplying out:

$$J(w) = \frac{1}{2} [y^T R^{-1} y - y^T R^{-1} x w - (xw)^T R^{-1} y + (xw)^T R^{-1} (xw)] \rightarrow (3)$$

• Recall: $(xw)^T = w^T x^T$

$$\therefore (xw)^T R^{-1} y = w^T x^T R^{-1} y \rightarrow (4)$$

• Recall, if a is a scalar: $a^T = a$.

• Since $y^T R^{-1} x w$ is a scalar:

$$y^T R^{-1} x w = (y^T R^{-1} x w)^T = w^T x^T R^{-1} y \rightarrow (5)$$

When since R is symmetric, so is R^{-1} and $(R^{-1})^T = R^{-1}$.

• Using (4) and (5) in (3):

$$J(w) = \frac{1}{2} [\underbrace{y^T R^{-1} y}_{\text{Independent of } w} - \underbrace{2 y^T R^{-1} x w}_{\text{Linear in } w} + \underbrace{w^T (x^T R^{-1} x) w}_{\text{Quadratic in } w}] \rightarrow (6)$$

• Recall: $f(w) = a^T w \Rightarrow \nabla f(w) = a \quad (a \in \mathbb{R}^n)$

$\boxed{A = A^T}$ $f(w) = w^T A w \Rightarrow \nabla f(w) = 2 A w$ (7)

• Consider the linear term: Set $a^T = y^T R^{-1} x$

$$\left. \begin{aligned} (y^T R^{-1} x) w &= a^T w \\ \nabla [\quad] &= a \end{aligned} \right\} \Rightarrow \nabla [(y^T R^{-1} x) w] = x^T R^{-1} y \rightarrow (8)$$

• Consider the quadratic term: $A = x^T R^{-1} x$

Verify $A^T = A$

(5)

$$\nabla [x^T A x] = 2 A x = 2 (x^T R^{-1} x) w \rightarrow (9)$$

Combining (7) to (9) with (6) and simplify:

$$\nabla J(w) = \frac{1}{2} [-2 x^T R^{-1} y + 2 x^T R^{-1} x w]$$

$$= -x^T R^{-1} [y - x w] \longrightarrow (10)$$

Hessian of $J(w)$:-

$$\nabla^2 J(w) = (x^T R^{-1} x) \in \mathbb{R}^{n \times n} - \text{SPD}$$

Equate the gradient to zero:

$$\nabla J(w) = 0 \Rightarrow (x^T R^{-1} x) w = x^T R^{-1} y \rightarrow (11)$$

If $R = \sigma^2 I_m$, $R^{-1} = \frac{1}{\sigma^2} I_m$ and (11)

becomes

$$(x^T x) w = x^T y \longrightarrow (12)$$

(12) is called Normal equation. Assume $X \in \mathbb{R}^{m \times n}$ is a full rank matrix: $\text{Rank}(X) = n$. This guarantees $x^T x$ is SPD:

$a^T (x^T x) x = (x a)^T (x a) = \|x a\|_2^2 \geq 0$ with equality only if $a = 0$ since the columns of x are linearly independent.

The least square solution: \hat{w}_{LS} :

Solving (12):

use cholsky } To solve (13)
• QR decomp.

$$\hat{W}_{LS} = (X^T X)^{-1} X^T y \rightarrow (13)$$

④ unbiasedness: $y = XW + v$

$$\begin{aligned} \therefore \hat{W}_{LS} &= (X^T X)^{-1} X^T [XW + v] \\ &= (X^T X)^{-1} (X^T X) W + (X^T X)^{-1} X^T v \end{aligned}$$

$$\therefore \cancel{E[\hat{W}_{LS}]} = \cancel{W} + (X^T X)^{-1} X^T v \rightarrow (14)$$

$$\therefore E[\hat{W}_{LS}] = W + (X^T X)^{-1} X^T E(v) = W \rightarrow (15)$$

$\therefore \hat{W}_{LS}$ is unbiased.

⑤ Covariance of \hat{W}_{LS} : $\text{Cov}(\hat{W}_{LS}) = \sigma^2 (X^T X)^{-1}$

$$\text{Cov}(\hat{W}_{LS}) = E\left[(\hat{W}_{LS} - E(\hat{W}_{LS}))(\hat{W}_{LS} - E(\hat{W}_{LS}))^T\right] \rightarrow (16)$$

$$\text{From (14): } \hat{W}_{LS} - E(\hat{W}_{LS})$$

$$= W + (X^T X)^{-1} X^T v - W = (X^T X)^{-1} X^T v \rightarrow (17)$$

Substituting (17) in (16)

$$\text{Cov}(\hat{W}_{LS}) = E\left[(X^T X)^{-1} X^T v (X^T X)^{-1} X^T v^T\right]$$

$$= E\left[(X^T X)^{-1} X^T (v v^T) X (X^T X)^{-1}\right]$$

$$= (X^T X)^{-1} X^T \underbrace{[E(v v^T)]}_{= R = \sigma^2 I_m} X (X^T X)^{-1}$$

$$= \sigma^2 (X^T X)^{-1} (X^T X) (X^T X)^{-1} = \underline{\sigma^2 (X^T X)^{-1}} \rightarrow (18)$$

(7)

Recall $X \in \mathbb{R}^{m \times n}$. Do an SVD of X

$$\left. \begin{aligned} X &= U \Lambda^{1/2} V^T, \\ (X^T X) V &= V \Lambda \\ (X X^T) U &= U \Lambda \end{aligned} \right\} \rightarrow (19)$$

What is the total sum of the variance of the components of \hat{W}_{LS} ? This is the $\text{tr}[\text{Cov}(\hat{W}_{LS})]$

$$\text{tr}[\text{Cov}(\hat{W}_{LS})] = \text{tr}[\sigma^2 (X^T X)^{-1}]$$

$$= \sigma^2 \text{tr}[(X^T X)^{-1}]$$

$$= \sigma^2 \text{tr}[V \Lambda^{-1} V^T]$$

$$= \sigma^2 \text{tr}[V \Lambda^{-1} V^T]$$

$$= \sigma^2 \text{tr}[V^T V \Lambda^{-1}]$$

$$= \sigma^2 \text{tr}[\Lambda^{-1}]$$

$$= \sigma^2 \sum_{i=1}^n \frac{1}{\lambda_i} \rightarrow (20)$$

$$\begin{aligned} \text{tr}(ABC) &= \text{tr}(CAB) \\ &= \text{tr}(BCA) \end{aligned}$$

$$V^T V = I$$

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n > 0$$

max min

Consequences:- The smallest eigenvalue of $(X^T X)$ determine the behavior of the total variance of the components of \hat{W}_{LS} . Thus, if the columns of X inherit collinearity property, then λ_n may be positive but very small and $\frac{1}{\lambda_n}$ can be very large.

Question:- How large the estimate \hat{W}_{LS} in this case?

⑥ Estimation of σ^2 :

8

We assumed σ^2 is known. If not, can be estimate σ^2 ?

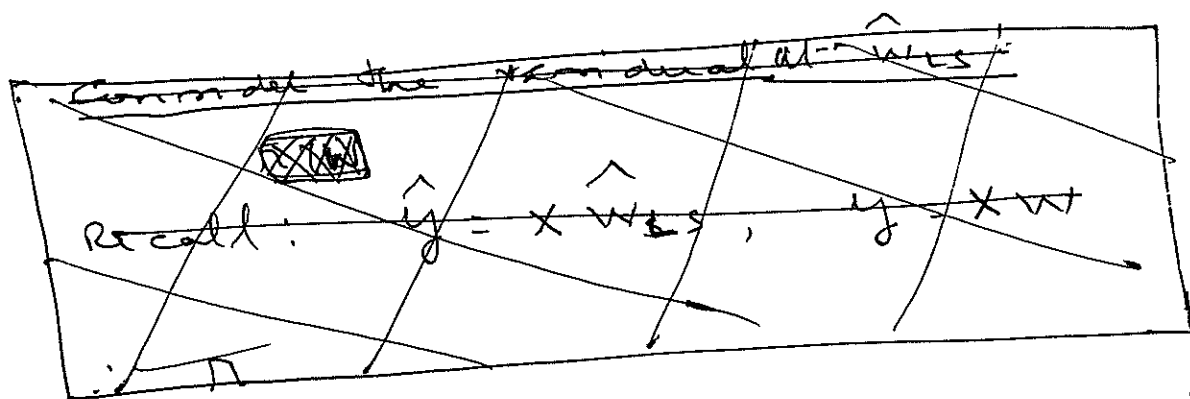
Recall $\hat{w}_{LS} = (X^T X)^{-1} X^T y = X^+ y \rightarrow (21)$

where $X^+ = (X^T X)^{-1} X^T$ is called generalized inverse of X .

Define a new matrix

$$P = X X^+ = X (X^T X)^{-1} X^T$$

Verify $P^T = P$ and $P^2 = P. \rightarrow (22)$



Define fitted value of y given by \hat{y} :

$$\hat{y} = X \hat{w}_{LS} = \underline{X (X^T X)^{-1} X^T} y = P y \rightarrow (23)$$

~~Define~~ Define error e :

$$e = y - \hat{y} = y - P y = (I - P) y \rightarrow (24)$$

Verify $(I - P) X = X - P X = X - X (X^T X)^{-1} (X^T X) = X - X = 0$

$$\begin{aligned} E(e) &= E[(I - P)y] \\ &= E[(I - P)(Xw + v)] \end{aligned}$$

(7)

$$\therefore e = (I - P)y = (I - P)(\cancel{X}w + v)$$

$$= \underbrace{[(I - P)X]}_{=0} w + (I - P)v$$

$$= (I - P)v.$$

Mean of e is zero

$$\therefore E(e) = (I - P)E(v) = 0 \Rightarrow$$

→ (25)

\therefore Total Variance in the components of e :

$$E(e^T e) = E\left[\{(I - P)v\}^T \{(I - P)v\}\right]$$

$$= E[v^T (I - P)^T (I - P)v]$$

$$= E\left[\underbrace{v^T (I - P)v}_{\text{scalar}}\right]$$

$$= E[\text{tr}[v^T (I - P)v]]$$

$$= E[\text{tr}[v v^T (I - P)]]$$

$$= \text{tr}[E(v v^T) (I - P)]$$

$$= \sigma^2 \text{tr}[I_m - P]$$

$$= \sigma^2 \left[\underbrace{\text{tr}(I_m)}_{=m} - \underbrace{\text{tr}(P)}_{=n} \right]$$

$$= \sigma^2 (m - n) \rightarrow (26)$$

$$E(v v^T) = \sigma^2 I_m$$

$$\text{tr}(I_m) = m$$

$$P = X(X^T X)^{-1} X^T$$

$$\text{tr}(P) = \text{tr}[X(X^T X)^{-1} X^T]$$

$$= \text{tr}[(X^T X)^{-1} X^T X]$$

$$= \text{tr}[I_n]$$

$$= n$$

This suggests an estimator for σ^2 ;

$$\hat{\sigma}^2 = \frac{e^T e}{(m-n)} \longrightarrow (27)$$

This is an unbiased estimator of σ^2 .

(7) Now to the case where X is such that $(X^T X)$ is nearly singular, or multicollinearity in X

Let $\lambda_1 > \lambda_2 > \dots > \lambda_n > 0 \longrightarrow (28)$

be the eigenvalues of $(X^T X) \in \mathbb{R}^{n \times n}$ - Smaller eigenvalues of X . Then

$$\lambda_1^{1/2} > \lambda_2^{1/2} > \dots > \lambda_n^{1/2} > 0 \longrightarrow (29)$$

are the singular values of X .

$\kappa(X) = \frac{\lambda_1^{1/2}}{\lambda_n^{1/2}}$ is a measure of the collinearity: larger $\kappa(X)$ more likely $(X^T X)$ is near singular.

In this case: $(X^T X) V = \Lambda V$, $\boxed{V V^T = V^T V = I, V^{-1} = V^T}$

$$\Rightarrow (X^T X) = V \Lambda V^T$$

$$\begin{aligned} (X^T X)^{-1} &= (V \Lambda V^T)^{-1} = (V^T)^{-1} \Lambda^{-1} V^{-1} \\ &= V \Lambda^{-1} V^T = \sum_{i=1}^n \frac{1}{\lambda_i} v_i v_i^T \rightarrow (30) \end{aligned}$$

\therefore If λ_n is very small, small changes in λ_n will reflect as big changes in $(X^T X)^{-1}$ and hence $\hat{WLS} = (X^T X)^{-1} X^T y \rightarrow (31)$

⑧ one solution to reduce this sensitivity is ⑪

Ridge regression:

Reformulate by changing the cost function:

$$J_R(w) = \frac{1}{2\sigma^2} (y - Xw)^T (y - Xw) + \frac{1}{2} \alpha \|w\|_2^2 \rightarrow (32)$$

Where α is the penalty parameter.

This penalty term is called a regularization term.

$$\nabla J_R(w) = \frac{1}{2\sigma^2} [-2X^T y + 2(X^T X)w + \alpha 2w]$$

$$= 0$$

$$\Rightarrow (X^T X + \alpha I) = X^T y$$



$$\hat{w}_{LS}^R = (X^T X + \alpha I)^{-1} X^T y \rightarrow (33)$$

$\therefore X_\lambda^+ = (X^T X + \alpha I)^{-1} X^T$

⑨ To make sense of (33) express X in SVD.

$$(X^T X) V = V \Lambda \quad (X X^T) U = U \Lambda, \quad X = U \Lambda^{1/2} V^T$$

$$\downarrow$$

Form (33): $(X^T X) = V \Lambda V^T, \quad V^T V = V V^T = I \rightarrow (34)$

$$\begin{aligned} \rightarrow X_\lambda^+ &= [V \Lambda V^T + \alpha V V^T]^{-1} [U \Lambda^{1/2} V^T]^T \\ &= [V (\Lambda + \alpha I) V^T]^{-1} [U \Lambda^{1/2} V^T]^T \\ &= V (\Lambda + \alpha I)^{-1} \underbrace{V^T V}_{= I_n} \Lambda^{1/2} U^T \end{aligned}$$

$$\begin{aligned} V^{-1} &= V^T \\ (AB)^{-1} &= B^{-1} A^{-1} \end{aligned}$$

$$= V \underbrace{(\Lambda + \alpha I)^{-1} \Lambda^{1/2}}_{\text{SVD of } X^T_\lambda} U^T = \text{SVD of } X^T_\lambda \quad (12)$$

$$= V D U^T \text{ where } D = \text{diag}(d_1, d_2, \dots, d_n)$$

$$d_i = \frac{\lambda_i^{1/2}}{\lambda_i + \alpha} \quad (1 \leq i \leq n)$$

$$\therefore X^T_\lambda = \sum_{i=1}^n \left(\frac{\lambda_i^{1/2}}{\lambda_i + \alpha} \right) V_i U_i^T \rightarrow (36)$$

The addition of α to λ_n stabilizes the estimate \hat{w}_{LS}

(10) For comparison: Consider a special case orthogonal, when columns of X are ~~orthogonal~~: that is: $X^T X = n I$. Then

$$\hat{w}_{LS}^R = (X^T X + \alpha I)^{-1} X^T \underbrace{[XW + N]}_{=Y} = (X^T X + \alpha I)^{-1} [X^T X W + X^T N]$$

$$= (nI + \alpha I)^{-1} n I W + (X^T X + \alpha I)^{-1} X^T N$$

$$= \left(\frac{n}{\alpha + n} \right) W + \frac{X^T N}{\alpha + n}$$

$$\therefore E[\hat{w}_{LS}^R] = \left(\frac{n}{n + \alpha} \right) W + \left(\frac{X^T}{n + \alpha} \right) E(N) = 0$$

$$= \left(\frac{n}{n + \alpha} \right) W \neq W \Rightarrow \underline{\underline{\text{biased}}}$$

$$\text{cov}(\hat{w}_{LS}) = \text{cov}\left(\frac{n}{n+\alpha} w + \frac{x^T u}{n+\alpha}\right)$$

$$= \text{cov}\left(\frac{x^T u}{n+\alpha}\right)$$

$$= \frac{1}{(n+\alpha)^2} E[(x^T u)(x^T u)^T]$$

$$= \frac{1}{(n+\alpha)^2} x^T \underbrace{E(u u^T)}_{= \sigma^2 I_m} x$$

$$= \frac{\sigma^2 I_m}{(n+\alpha)^2} \cancel{x^T x} = \frac{n \sigma^2}{(n+\alpha)^2} I_m$$

$$= \frac{n}{(n+\alpha)^2} \sigma^2 \boxed{n (x^T x)^{-1}}$$

$$= \frac{n^2 \sigma^2}{(n+\alpha)^2} (x^T x)^{-1}$$

Recall $\text{cov}(\hat{w}_{OLS}) = \sigma^2 (x^T x)^{-1}$

$$\therefore \text{cov}(\hat{w}_{LS}^R) < \text{cov}(\hat{x}_{LS})$$

since $\frac{n^2}{(n+\alpha)^2} \sigma^2 (x^T x)^{-1} < \sigma^2 (x^T x)^{-1}$

$$\therefore \text{cov}(\hat{w}_{LS}) - \text{cov}(\hat{w}_{LS}^R)$$

$$= \sigma^2 (x^T x)^{-1} - \frac{n^2}{(n+\alpha)^2} \sigma^2 (x^T x)^{-1}$$

$$\begin{aligned} (x^T x) &= n I_m \\ \therefore (x^T x)^{-1} &= \frac{1}{n} I_m \\ \therefore I_m &= n (x^T x)^{-1} \end{aligned}$$

$$= \sigma^2 (x^T x)^{-1} \left[1 - \frac{n^2}{(n+\alpha)^2} \right]$$

$$= \sigma^2 (x^T x)^{-1} \left[\frac{\alpha (2n + \alpha)}{(n+\alpha)^2} \right] > 0$$

(ii) Bias in $\hat{w}_{LS}^R = \frac{n}{n+\alpha}$

$$\text{Difference in Variance} = \left[1 - \frac{n^2}{(n+\alpha)^2} \right] > 0$$

(ii) \hat{w}_{LS}^R has lesser variance but biased.

This is called bias - Covariance trade-off.