$$Y = X\theta^0 + \varepsilon \qquad\qquad (y_i, x_i)$$

$$\varepsilon \sim \xi \qquad \mathbb{E}[\varepsilon] = 0, \qquad Var(\varepsilon) = \sigma^2$$

$$\hat{\theta} \in \arg\min_{\theta \in \mathbb{R}^{p+1}} \| Y - X\theta \|^2 \qquad X \in \mathbb{R}^{n \times (p+1)} \quad Y \in \mathbb{R}^n$$

$$\downarrow \qquad\qquad\qquad \theta^0, \hat{\theta} \in \mathbb{R}^{p+1}$$

$$(X^T X) \, \hat{\theta} = X^T y \qquad \text{Normal equation}$$

▷ The solution to OLS is unique iff $(X^T X)^{-1} \; \exists$

exer 1 $\qquad X = 1_n = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \qquad x, y, \qquad$ Which is the OLS?

$$\hat{\theta} = (X^T X)^{-1} X^T y = (1_n^T 1_n)^{-1} 1_n^T y = \left[ (1 \cdots 1) \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \right]^{-1} 1_n y =$$

$$= n^{-1} 1_n^T y = \frac{1}{n} \sum_{i=1}^n y_i = \bar{y}$$

exer 2

Show $\qquad \|Y - \hat{Y}\|^2 \le \| Y - \bar{y} 1_n \|^2$

Let $X \in (1_n, \tilde{X})$ for $\tilde{X} \in \mathbb{R}^{n \times p}$, $X \in \mathbb{R}^{n \times (p+1)}$, $1_n = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \in \mathbb{R}^n$

Let $\theta = (\theta_0, \theta_1, \ldots, \theta_r)^\top = (\theta_0, \widetilde{\theta})^\top$, $\theta_0 \in \mathbb{R}$, $\widetilde{\theta} \in \mathbb{R}^s$

$$\min_{\theta \in \mathbb{R}^{p+1}} \| Y - X\theta \|^2 = \min_{\substack{\theta_0 \in \mathbb{R}, \\ \widetilde{\theta} \in \mathbb{R}^p}} \| Y - (1_n, \widehat{x})\begin{pmatrix} \theta_0 \\ \theta \end{pmatrix} \|^2$$

$$\leq \min_{\theta_0} \| Y - 1_n \theta_0 \|^2 = \| Y - \bar{Y} 1_n \|^2$$

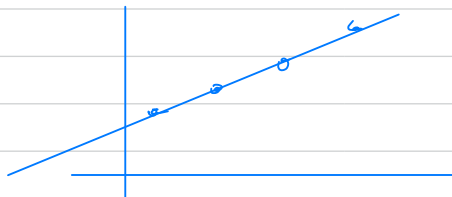thus $\quad \| Y - \widehat{Y} \|^2 \leq \| Y - \bar{Y} 1_n \|^2 \qquad \qquad \Uparrow$

<u>Proposition</u> $\qquad 0 \overset{①}{\leq} R^2 \overset{②}{\leq} 1$

Recall
$$R^2 = 1 - \frac{\| Y - \widehat{Y} \|^2}{\| Y - \bar{Y} 1_n \|^2} \qquad Y = \widehat{Y}$$

② for $\widehat{Y} = Y \Rightarrow R^2 = 1$

① by exer2



---

Recall : $\quad X \to$ random matrix
$\qquad$ $A, b$ deterministic $\quad$ matrix / vector
$\qquad$ $= \text{cov}(AX+b) = A \, \text{cov}(X) \, A^\top$

<u>Properties</u> $\qquad$ (we assume $\ker(X) = 0$, i.e $(X^\top X)^{-1}$ $\exists$s

<u>property1</u>: $\widehat{\theta}$ is unbiased $\theta^*$, i.e $\mathbb{E}[\widehat{\theta}] = \theta^*$ $\qquad$ only reads

$$\mathbb{E}[\widehat{\theta}] = \mathbb{E}[(X^\top X)^{-1} X^\top Y] = \mathbb{E}[(X^\top X)^{-1} X^\top (X\theta^* + \varepsilon)] =$$

$$= (X^\top X)^{-1} \underbrace{X^\top X} \theta^* + (X^\top X)^{-1} X^\top \underbrace{\mathbb{E}[\varepsilon]}_{0} = \theta^*$$

**P2:** $\operatorname{cov}(\hat{\theta}) = \sigma^2 (X^TX)^{-1}$

$\operatorname{cov}(\hat{\theta}) = \operatorname{cov}(\underbrace{(X^TX)^{-1}X^TY}_{A}) = (X^TX)^{-1}X^T\underbrace{\operatorname{cov}(Y)}_{\sigma^2 I}((X^TX)^{-1}X^T)^T$

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

**property 3**      $\hat{\theta}$ is the $\underset{\text{min var}}{\underline{\text{Best}}}$ $\underset{AY}{\underline{\text{Linear}}}$ $\underset{\mathbb{E}[\hat{\theta}]=\theta^*}{\underline{\text{Unbiased}}}$ Estimator (BLUE) $\theta^*$

By linear we mean that $\hat{\theta} = AY$. Note that

$\hat{\theta}_{OLS} = \underline{(X^TX)^{-1}}X^TY$. Thus, in OLS $A = (X^TX)^{-1}X^T$

$\mathbb{E}[AY] = \theta^*$   (its unbiased)

$\Longleftrightarrow$   $A\mathbb{E}[Y] = \theta^*$ $\Longleftrightarrow$ $A\mathbb{E}[X\cdot\theta^* + \varepsilon] = \theta^*$

$\Rightarrow$ $\underline{AX\theta^*} = \theta^*$        $\Longleftrightarrow$      $AX = I$ .

$\operatorname{cov}(AY) = A\operatorname{cov}(Y)A^T = \sigma^2 AA^T$

$\Rightarrow$ We choose $A$ so that the diagonal of $AA^T$ is min
subject to those that satisfy $AX = I$

$$AA^T = \left( \left( A - (X^TX)^{-1}X^T \right) + (X^TX)^{-1}X^T \right) \left( A - (X^TX)^{-1}X^T + (X^TX)^{-1}X^T \right)^T$$

$$\Rightarrow \left( A - (X^TX)^{-1}X^T \right) \left( (X^TX)^{-1}X^T \right)^T = (A - (X^TX)^{-1}X^T)(X(X^TX)^{-1}) =$$

$$\Rightarrow = \underbrace{AX}_{I}(X^TX)^{-1} = (X^TX)^{-1}X^TX(X^TX)^{-1} = 0$$

$$AA^T = \underbrace{\left( A - (X^TX)^{-1}X^T \right)}_{B} \underbrace{\left( A - (X^TX)^{-1}X^T \right)^T}_{B^T} + (X^TX)^{-1}$$

$$BB^T \qquad [BB^T]_{ii} = \sum_k B_{ik}^2 \qquad \geq 0$$

The diagonal in $AA^T$ (ie, the variance of the linear

estimator $AY$) is $\geq 0$ always     When is it "$= 0$"?

We choose $A$. $AA^T = 0$ ∴ diag$(AA^T) = 0$

$\Rightarrow \quad A = (X^TX)^{-1}X^+ \quad \rightarrow$ Note. this choice satisfies $AX = I$

$\Rightarrow \quad A = (X^TX)^{-1}X^T \quad$ is the BLUE ▢

exer 3    Show that the predicted value $\hat{Y}$ is invariant to
          linear changes on $X$

$$X = [x_0, \ldots , x_p] \qquad Z = [c_0 x_0, c_0 x_{11} - \ldots , c_p x_p]$$

How to write the transformed prob?

$$D = \text{diag}(c_0, c_1, \ldots, c_p)$$

$$Z = XD$$

$$\hat{\theta}_x = (X^TX)^{-1} X^T y$$

$$\hat{\theta}_z = (Z^TZ)^{-1} Z^T y = ((XD)^T(XD))^{-1}(XD)^T y =$$

$$= (D^T X^T XD)^{-1} D^T X^T y = D^{-1}(X^TX)^{-1} \underbrace{D^{-1}D}_{I} X^T y =$$

$$= D^{-1} \hat{\theta}_x \qquad \color{red}{\text{Note! typo}}$$

$$x_0 \in \mathbb{R}^{p+1} \quad \rightarrow \text{point prediction} \qquad \hat{\theta}_x^T x_0$$

let $z_0$ the point with linear changes $\qquad z_0 = D x_0$

point prediction for $z_0$?

$$\hat{\theta}_z^T z_0 = (D^{-1}\hat{\theta}_x)^T \cdot (Dx_0) = \hat{\theta}_x^T D^{-1} D x_0 = \hat{\theta}_x^T x_0 \quad \square$$

## Cochran's lemma

Rem     X is fixed     Gaussian noise

$$\varepsilon_i \sim N(0, \sigma^2) \qquad\qquad y = X\theta^* + \varepsilon \quad, \quad X \text{ full rank}$$

$$\hat{\sigma}^2 = \frac{1}{n-p-1} \sum_{i=1}^{n} \hat{\varepsilon}_i^2 \qquad\qquad \boxed{\begin{array}{l}\text{slide 59}\\ \text{in session 1}\\ H^2 = H = H^T\end{array}}$$

Hat matrix.     $H = X(X^TX)^{-1}X^T$

$$\hat{y} = X\hat{\theta} = \underbrace{X(X^TX)^{-1}X^T}_{H} y = HY$$

· Note that $(I - H) X = 0$

(1) $\hat{\theta} \sim N(\theta^*, \sigma^2 (X^T X)^{-1})$

(2) $(n-p-1) \dfrac{\hat{\sigma}^2}{\sigma^2} \sim \chi^2_{n-p-1}$

(3) $\hat{\theta}, \hat{\sigma}^2$ are independent

(4) $\mathbb{E}[\hat{\sigma}^2] = \sigma^2$ is unbiased

(5) - Relation to T-student distn : next week

for (1)

$\hat{\theta} = (X^T X)^{-1} X^T Y = (X^T X)^{-1} X^T (X \theta^* + \varepsilon) =$

$= (X^T X)^{-1} X^T X \theta^* + (X^T X)^{-1} X^T \varepsilon$

$\varepsilon \sim N(0, \sigma^2)$

$\hat{\theta}$ is gaussian, characterized $\mathbb{E}[\hat{\theta}]$, $Var(\hat{\theta})$

(2) $V = (V_1, V_2)$

$V_1$ is a basis for span $(x)$
$V$ is orthogonal $\in \mathbb{R}^{n \times n}$

$$\rightarrow \begin{cases} V_1^T(1-H) = 0 \\ V_2^T(1-H) = V_2^T \end{cases} \qquad \begin{vmatrix} H \cdot X \\ (1-H)X \end{vmatrix} \quad \text{hat matrix } X$$

$$(n-p-1)\,\hat{\sigma}^2 = \frac{1}{\#} \sum_{i=1}^{n} \hat{\varepsilon}_i^2 = \| Y - \widehat{Y} \|^2 =$$

$$\| Y - HY \|^2 = \| (1-H)\cdot Y \|^2 = \| (1-H)(X\theta^* + \varepsilon) \|^2$$

$$= \| (1-H)\varepsilon \|^2 = \| V^T (1-H)\varepsilon \|^2 =$$

$$= \| V_2^T \hat{\varepsilon} \|^2$$

$$\text{Note} \quad \text{let} \quad \tilde{\varepsilon} = \frac{V_2^T \varepsilon}{\varepsilon} \quad \text{then}$$

$$(n-p-1)\frac{\hat{\sigma}^2}{\sigma^2} = \sum_{i=1}^{n-p-1} \tilde{\varepsilon}_i^2 \qquad \varepsilon \sim N$$

$$\sim \chi^2_{n-p-1}$$

$$(3) \qquad \text{Note} \quad X^T(1-H) = 0$$

$$\frac{1}{n-p-1}\hat{\sigma}^2 = \| \hat{\varepsilon} \|^2 = \| (1-H)\varepsilon \|^2 \qquad \Big\{$$

$$\hat{\theta} - \theta^* = (X^TX)^{-1} X^T \varepsilon$$

$$(4) \qquad \mathbb{E}[\hat{\sigma}^2] = \sigma^2$$

$$\mathbb{E}\left[ \frac{1}{n-p-1} \frac{\hat{\sigma}^2}{\sigma^2} \right] = n-p-1$$

$$\underbrace{\quad}_{\text{random var.}} \qquad \chi^2_{n-p-1} \text{ based on } (2)$$

$$\frac{1}{n-p-1} \cdot \frac{\mathbb{E}[\hat{\sigma}^2]}{\sigma^2} = n - p -$$

$$\mathbb{E}[\hat{\sigma}^2] = \sigma^2$$

---

# MAXIMUM LIKELIHOOD ESTIMATION

let $X_i \sim N(\mu, \sigma^2)$. $S = \{X_1, \ldots, X_n\}$ is a

sample. Give the MLE for the params, $\hat{\mu}, \hat{\sigma}^2$.

Recall, the density is

$$p(x_i, \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu_i)^2}{2\sigma^2}\right)$$

1) give the likelihood of the sample $x_1 \ldots x_n$

2) give the log-likelihood

3) derivate w.r.t. $\mu$ and $\sigma^2$

4) solve the equation in (3)

$S = \{x_1, \ldots, x_n\}$  $\quad x_i \sim N(\_, \_)$

1) $\ell(\mu, \sigma^2; S) = \prod_{i=1}^{\hat{n}} p(x_i; \mu, \sigma^2) =$

$= \prod_{i=1}^{n} (2\pi\sigma)^{-\frac{1}{2}} \exp\left(\frac{-(x_i - \mu)^2}{2\sigma^2}\right)$

$= \cdot (2\pi\sigma)^{-\frac{n}{2}} \exp\left(-\sum \cdot \frac{(x_i - \mu)^2}{2\sigma^2}\right)$

2) $\mathcal{L}(\mu, \sigma^2, S) = -\frac{n}{2}\lg(2\pi) - \frac{n}{2}\lg\sigma^2 - \frac{1}{2\sigma^2}\sum(x_i - \mu)^2$

3) $\underset{\mu, \sigma^2}{\text{argmax}} \ \mathcal{L}(\mu, \sigma^2; S) \iff \nabla\mathcal{L}(\mu, \sigma^2, S) = 0$

$\dfrac{\partial \mathcal{L}}{\partial \mu} = \dfrac{1}{\sigma^2} \sum(x_i - \hat{\mu})^2 = 0$

$\sum(x_i - \hat{\mu}) = 0 \quad \rightarrow \quad \boxed{\hat{\mu} = \frac{1}{n}\sum x_i = \bar{x}}$

$\dfrac{\partial\mathcal{L}}{\partial\sigma^2} = \dfrac{-n}{2\hat{\sigma}^2} - \left(\dfrac{1}{2}\sum(x_i - \hat{\mu})^2\right) \underbrace{\dfrac{d}{\hat{\sigma}^2}\left(\dfrac{1}{\sigma^2}\right)}_{} \overset{=0}{\longrightarrow} \dfrac{1}{(\hat{\sigma}^2)^2}$

$\dfrac{1}{\hat{\sigma}^2}\left(\dfrac{1}{\hat{\sigma}^2}\sum(x_i - \mu_i)^2 - n\right) = 0$

$\boxed{\dfrac{1}{n}\sum(x_i - \mu_i)^2 = \hat{\sigma}^2}$

$$\varepsilon_i \sim N\left(0, \sigma^2\right)$$

Back to regression!

$$y_i = \theta_0^* + \theta_1^* x_i + \varepsilon_i$$

We observe $(x_i, y_i)_{i=1}^n$

$$\mathbb{E}[y_i] = \theta_0^* + \theta_1^* x$$

We want to estimate $\theta_0^*, \theta_1^*, \sigma^2$

1) $\ell(\theta_0, \theta_1, \sigma^2 ; S) = \prod_{i=1}^n P((x_i, y_i), \theta_0, \theta_1, \sigma^2)$

$$= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(\frac{-(y_i - (\theta_0 + \theta_1 x))^2}{2\sigma^2}\right)$$

$$= \left(2\pi\sigma^2\right)^{-\frac{n}{2}} \exp\left(\frac{-1}{2\sigma^2} \sum_{i=1}^n (y_i - \theta_0 - \theta_1 x_i)^2\right)$$

**2)** log likelihood $\mathcal{L}(\theta_0, \theta_1, \sigma^2) = \lg L(\ldots)$

$$= -\frac{n}{2} \lg(2\pi) - \frac{n}{2} \lg \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^{n} (y_i - \theta_0 - \theta_1 x_n)^2$$

**3)** $\underset{\theta_0, \theta_1, \sigma^2}{\arg\max} \mathcal{L}(\theta_0, \theta_1, \sigma_z^2, S)$

Find the partial derivatives w.r.t $\theta_0, \theta_1, \sigma^2$

$$\frac{\partial \mathcal{L}'}{\partial \theta_0} = \frac{1}{\sigma^2} \sum (y_i - \hat{\theta}_0 - \hat{\theta}_1 x_i) = 0$$

$$\boxed{\hat{\theta}_0^{MLE} = \bar{y} + \hat{\theta}_1 \bar{x}}$$

$$\frac{\partial \mathcal{L}}{\partial \theta_1} = \frac{1}{\sigma^2} \sum (y_i - \hat{\theta}_0 - \hat{\theta}_1 x_i) x_i = 0$$

$$\boxed{\hat{\theta}_1^{MLE} = \frac{cov(x, y)}{var(x)}}$$

$$\frac{\partial \mathcal{L}}{\partial \sigma^2} = \frac{-n}{2\hat{\sigma}^2} + \frac{1}{2(\hat{\sigma}^2)^2} \cdot \sum (y_i - \hat{\theta}_0 - \hat{\theta}_1 x_i)^2 = 0$$

$$\frac{+n}{2\hat{\sigma}^2} = \frac{1}{2(\hat{\sigma}^2)^2} \cdot \sum (y_i - \hat{\theta}_0 - \theta_1 x_i)^2$$

$$\hat{\sigma}^2_{MLE} = \frac{1}{n} \sum (y_i - \hat{\theta}_0 - \hat{\theta}_1 x_i)^2$$

$$\hat{\sigma}^2_{OLS} = \frac{1}{n-p-1} \sum_{i=1}^{n} \hat{\varepsilon}_i^2 \qquad X \in \mathbb{R}^{u \, \text{oper}}$$

Simple regression (i.e. 1D) $X, y \in \mathbb{R}^n$

in this case $p = 1$ $\Rightarrow$

$$\widehat{\sigma}^2_{OLS} = \frac{1}{n-2} \sum_{i=1}^{n} \widehat{\varepsilon}_i^2 = \frac{1}{n-2} \sum_{i=1}^{n} (y_i - \widehat{\theta}_0 - \widehat{\theta}_1 x_i)^2$$

$\Rightarrow$

$\widehat{\sigma}^2_{MLE}$ is biased

but $\widehat{\theta}_0^*, \widehat{\theta}_1$ they are the same

$\varepsilon_i \sim \sigma^2$

$\theta_0$