

Lec 10 3/2/18 Prob 3924

11

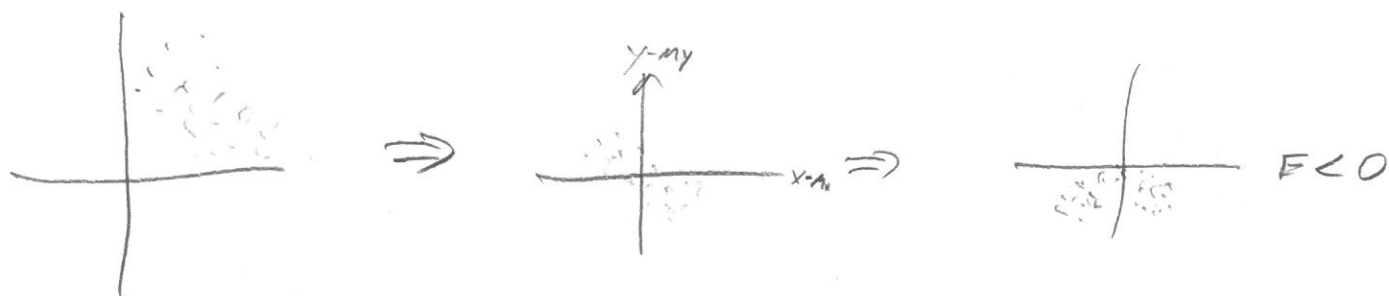
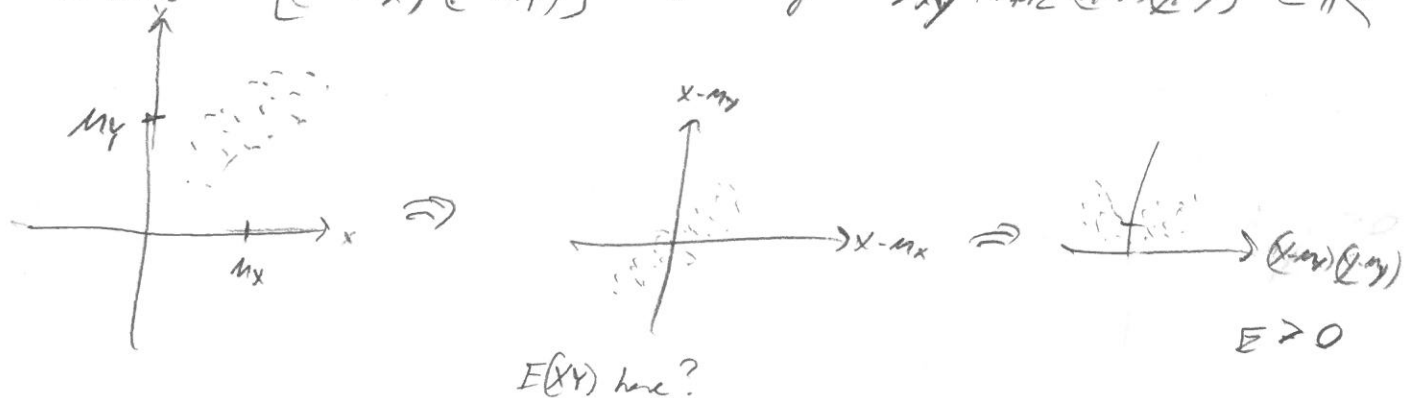
In Prob 241 r.v.'s X, Y were said to be dependent if knowing the value of one affects the dist of the other:

$$P(Y|X=x) \neq P(Y)$$

In data science if knowing a predictor x allows you to know something about y we say X, Y are "associated".

Recall covariance

$$\text{Cov}(X, Y) := E[(X - \mu_X)(Y - \mu_Y)] \text{ estm by } s_{xy} := \frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y}) \in \mathbb{R}$$



The sign of the cov. is important. It indicates if \oplus if $x \uparrow \Rightarrow y \uparrow$

If \ominus if $x \uparrow \Rightarrow y \downarrow$

$$\text{Pearson } r := \text{Corr}(X, Y) := \frac{\text{Cov}(X, Y)}{SE(X)SE(Y)} \in [-1, 1] \text{ standardised! Convenient!}$$

est by

$$r = \frac{s_{xy}}{s_x s_y} \in [-1, 1]$$

We say X, Y are pos. corr if $r > 0$ which means $x \uparrow \Rightarrow y \uparrow$

We say ... neg. if $r < 0$ which means $x \uparrow \Rightarrow y \downarrow$

... not ... if $r = 0$ which means $x \uparrow \Rightarrow y$ unchanged

Covariance is "linear correlation" as it appears in LS! It is a type of association

Previously $y \in \mathbb{R}$, $p=1$, $\mathcal{H} = \{ \vec{w} \cdot \vec{x} = w_0 + w_1 x : w_0 \in \mathbb{R}, w_1 \in \mathbb{R} \}$

Now let $p=2 \Rightarrow \mathcal{H} = \{ w_0 + w_1 x_1 + w_2 x_2 : \vec{w} \in \mathbb{R}^3 \}$

for any \vec{w} , $(y_i - \hat{y}_i)^2$

$$SSE = \sum_{(\vec{x}_i, y_i) \in \mathcal{D}} (y_i - (w_0 + w_1 x_{i1} + w_2 x_{i2}))^2 = \dots \text{ long}$$

Take $\frac{\partial}{\partial w_0} [SSE] \stackrel{set}{=} 0$, $\frac{\partial}{\partial w_1} [SSE] \stackrel{set}{=} 0$, $\frac{\partial}{\partial w_2} [SSE] \stackrel{set}{=} 0$

Can we solve this any better?

$\mathcal{D} = \langle X, \vec{y} \rangle$ where $X = \begin{bmatrix} 1 & x_{11} & x_{12} \\ 1 & x_{21} & x_{22} \\ \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} \end{bmatrix} \in \mathbb{R}^{n \times \overset{p+1}{3}}$

where is $X\vec{w} = \begin{bmatrix} w_0 + w_1 x_{11} + w_2 x_{12} \\ w_0 + w_1 x_{21} + w_2 x_{22} \\ \vdots \\ w_0 + w_1 x_{n1} + w_2 x_{n2} \end{bmatrix} = \begin{bmatrix} \hat{y}_1 \\ \hat{y}_2 \\ \vdots \\ \hat{y}_n \end{bmatrix} \Rightarrow \vec{\hat{y}} = X\vec{w}$

Recall $(\vec{a} + \vec{b})^T = \vec{a}^T + \vec{b}^T$

$$\Rightarrow SSE = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = (\vec{y} - \vec{\hat{y}})^T (\vec{y} - \vec{\hat{y}}) = (\vec{y}^T - \vec{\hat{y}}^T) (\vec{y} - \vec{\hat{y}})$$

$$= \vec{y}^T \vec{y} - \vec{\hat{y}}^T \vec{y} - \vec{y}^T \vec{\hat{y}} + \vec{\hat{y}}^T \vec{\hat{y}}$$

Recall $\vec{v} \in \mathbb{R}^d$

$$\vec{v} \cdot \vec{v} = \sum_{j=1}^d v_j^2$$

Recall $\vec{a} \cdot \vec{b} = \vec{b} \cdot \vec{a}$

$$= \vec{y}^T \vec{y} - 2\vec{\hat{y}}^T \vec{y} + \vec{\hat{y}}^T \vec{\hat{y}}$$

$$= \vec{y}^T \vec{y} - 2(X\vec{w})^T \vec{y} + (X\vec{w})^T (X\vec{w})$$

Recall $(AB)^T = B^T A^T \Rightarrow \vec{y}^T \vec{y} - 2\vec{w}^T X^T \vec{y} + \vec{w}^T X^T X \vec{w}$

Now we need $\frac{\partial}{\partial w_0} [\dots], \frac{\partial}{\partial w_1} [\dots], \dots, \frac{\partial}{\partial w_p} [\dots]$

$$\frac{\partial}{\partial \vec{z}} [a] = \begin{pmatrix} \frac{\partial}{\partial z_1} [a] \\ \vdots \\ \frac{\partial}{\partial z_n} [a] \end{pmatrix} = \vec{0} \quad \vec{z} \in \mathbb{R}^n$$

we can imagine taking the derivative wrt to the whole vector

$$\frac{\partial SSE}{\partial \vec{w}} := \begin{bmatrix} \frac{\partial}{\partial w_0} [SSE] \\ \vdots \\ \frac{\partial}{\partial w_p} [SSE] \end{bmatrix} \quad \text{i.e. the vector derivative}$$

$$\vec{z} [a \vec{z}] = [a_1, c_1]$$

$$\begin{aligned} \text{e.g. } \frac{\partial}{\partial \vec{z}} [g(f(\vec{z})) + g(\vec{z})] &= \begin{bmatrix} \frac{\partial}{\partial c_1} [g(f(\vec{z})) + g(\vec{z})] \\ \vdots \\ \frac{\partial}{\partial c_n} [g(f(\vec{z})) + g(\vec{z})] \end{bmatrix} \\ &= \begin{bmatrix} g \frac{\partial}{\partial c_1} [f(\vec{z})] + \frac{\partial}{\partial c_1} [g(\vec{z})] \\ \vdots \\ g \frac{\partial}{\partial c_n} [f(\vec{z})] + \frac{\partial}{\partial c_n} [g(\vec{z})] \end{bmatrix} \\ &= g \frac{\partial}{\partial \vec{z}} [f(\vec{z})] + \frac{\partial}{\partial \vec{z}} [g(\vec{z})] \end{aligned}$$

$$\text{e.g. } \frac{\partial}{\partial \vec{z}} [\vec{z}^T \vec{b}] = \frac{\partial}{\partial \vec{z}} [c_1 b_1 + c_2 b_2 + \dots + c_n b_n] = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} = \vec{b}$$

where $\vec{z} \in \mathbb{R}^n$, $\vec{b} \in \mathbb{R}^n$

e.g. $\frac{\partial}{\partial \vec{z}} [\vec{z}^T A \vec{z}]$ where $A \in \mathbb{R}^{n \times n}$ and $\vec{z} \in \mathbb{R}^n$ and A symmetric

$$\text{Note } A \vec{z} = \begin{bmatrix} a_{11}c_1 + a_{12}c_2 + \dots + a_{1n}c_n \\ a_{21}c_1 + a_{22}c_2 + \dots + a_{2n}c_n \\ \vdots \\ a_{n1}c_1 + a_{n2}c_2 + \dots + a_{nn}c_n \end{bmatrix}$$

$$\begin{aligned} \vec{z}^T (A \vec{z}) &= c_1 (a_{11}c_1 + a_{12}c_2 + \dots + a_{1n}c_n) \\ &\quad + c_2 (a_{21}c_1 + a_{22}c_2 + \dots + a_{2n}c_n) \\ &\quad \vdots \\ &\quad + c_n (a_{n1}c_1 + a_{n2}c_2 + \dots + a_{nn}c_n) = \sum_{j=1}^n \sum_{i=1}^n a_{ij} c_i c_j \\ &= c_1 (\vec{z}^T \vec{a}_{1\cdot}) + c_2 (\vec{z}^T \vec{a}_{2\cdot}) + \dots + c_n (\vec{z}^T \vec{a}_{n\cdot}) \end{aligned}$$

$$\begin{aligned} \frac{\partial}{\partial \vec{z}} [\vec{z}^T A \vec{z}] &= \begin{bmatrix} \frac{\partial}{\partial c_1} [\vec{z}^T A \vec{z}] \\ \vdots \\ \frac{\partial}{\partial c_n} [\vec{z}^T A \vec{z}] \end{bmatrix} = \begin{bmatrix} 2c_1 a_{11} + 2c_2 a_{12} + \dots + 2c_n a_{1n} \\ \vdots \\ 2c_1 a_{n1} + 2c_2 a_{n2} + \dots + 2c_n a_{nn} \end{bmatrix} \\ &= 2 \vec{a}_{1\cdot}^T \vec{z} \\ &\quad \vdots \\ &= 2 \vec{a}_{n\cdot}^T \vec{z} \end{aligned}$$

Now let's try to apply this to our problem at hand

$$= 2 \begin{bmatrix} \vec{a}_{1\cdot} \\ \vec{a}_{2\cdot} \\ \vdots \\ \vec{a}_{n\cdot} \end{bmatrix} \vec{z} = 2 A \vec{z}$$

$$\frac{\partial}{\partial \vec{w}} [SSE] = \frac{\partial}{\partial \vec{w}} [\vec{y}^T \vec{y} - 2 \vec{w}^T X^T \vec{y} + \vec{w}^T X^T X \vec{w}]$$

$$= \frac{\partial}{\partial \vec{w}} [\vec{y}^T \vec{y}] - 2 \frac{\partial}{\partial \vec{w}} [\vec{w}^T X^T \vec{y}] + \frac{\partial}{\partial \vec{w}} [\vec{w}^T \underbrace{X^T X}_{\text{symmetric?}} \vec{w}]$$

$$= \vec{0} - 2 X^T \vec{y} + 2 X^T X \vec{w} \stackrel{\text{set}}{=} \vec{0}$$

$$(X^T X)^T = (X^T X)^T = X^T X \quad \checkmark$$

$$\Rightarrow X^T \vec{y} = X^T X \vec{w} \Rightarrow \boxed{\vec{b} = (X^T X)^{-1} X^T \vec{y}}$$

LS solution!

Note: $(X^T X)$ must be invertible!
 $(p+1 \times n)(n \times p+1)$
 $\in \mathbb{R}^{p+1 \times p+1}$

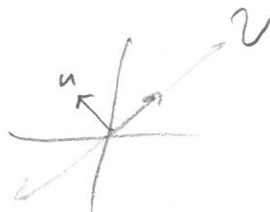
When is $X^T X$ invertible?

Only when $\text{rank}(X) = p+1$ Proof:

Assume $\text{rank}(X) < p+1 \Rightarrow$ there's a non-trivial nullspace i.e. a vector $\vec{u} \neq \vec{0} \in \mathbb{R}^{p+1}$ that can be mapped to $\vec{0}$.

i.e. $X\vec{u} = \vec{0}_n$

$$X = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$



$$V = \text{colspace}(X)$$

is just the line

$$\vec{u} = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

But then if we use $X^T X$ to map vector \vec{u} ,

$$(X^T X) \vec{u} = X^T (X \vec{u}) = X^T \vec{0}_n = \vec{0}_{p+1} \quad X \vec{u} = \vec{0}$$

this means $\vec{u} \in \text{Nullspace}(X^T X)$ and thus $\dim(\text{Nullspace}(X^T X)) > 0 \Rightarrow X^T X$ is not invertible.

What does $\text{rank}(X) = p+1$ mean? Each col is not lin. dep. on other cols.
 \Rightarrow No predictor information is duplicate

$$X = \begin{bmatrix} \text{Salary} - 2000 & \text{Salary} - 2001 & \text{tot_salary} \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \end{bmatrix}$$

not full rank \Rightarrow
 delete tot_salary!
 (or one of the yrs)

$$X = \begin{bmatrix} \text{height-in-feet} & \text{height-in-meters} \\ \vdots & \vdots \\ \vdots & \vdots \\ \vdots & \vdots \end{bmatrix}$$

" \Rightarrow
 "
 "

What is y ?

$$g(\vec{x}^*) = \vec{x}^{*T} \vec{b} = \begin{bmatrix} 1 & x_1^* & \dots & x_p^* \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \\ \vdots \\ b_p \end{bmatrix} = b_0 + b_1 x_1^* + \dots + b_p x_p^*$$