

Lee 7 Math 390 2/21/19

$X_{row} \in \{\text{Red}, \text{Green}\}$

$X \in \{0, 1\}$ binary encoding

$g(x) = \begin{cases} \bar{y}_r & \text{if } x = \text{red} \\ \bar{y}_g & \text{if } x = \text{green} \end{cases}$ this is the OLS sol!

Can we prove this? let $\sum x_i = n_g$ # green $\Rightarrow n - n_g = n_r$ # red

$$\bar{X} = \frac{\sum x_i}{n} = \frac{n_g}{n} = p, \text{ prop. of green}$$

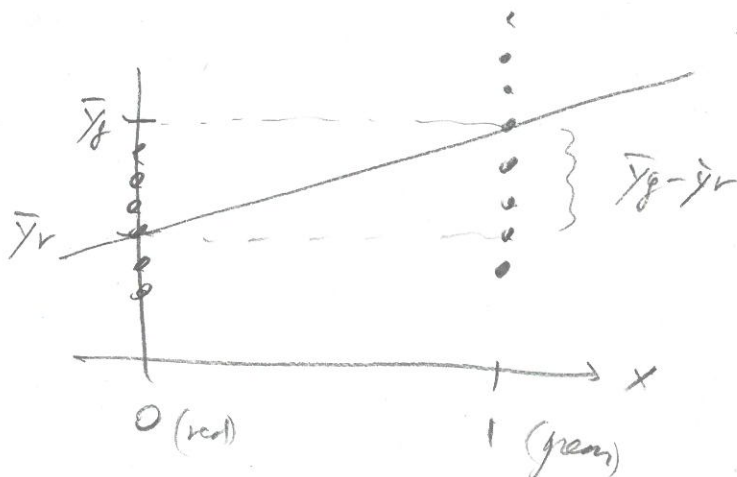
$$\bar{Y} = \frac{\sum y_i}{n} = \frac{\sum_{i: \text{green}} y_i + \sum_{i: \text{red}} y_i}{n} = \frac{\sum_{i: \text{green}} y_i}{n} \cdot \frac{n_g}{n_g} + \frac{\sum_{i: \text{red}} y_i}{n} \cdot \frac{n_r}{n_r} = \bar{Y}_g \cdot p + \bar{Y}_r (1-p)$$

$$b_1 = \frac{\frac{n_g \bar{y}_g}{\sum x_i y_i - n \bar{x} \bar{y}}}{\frac{\sum x_i^2 - n \bar{x}^2}{n_g - np^2}} = \frac{n_g \bar{y}_g - np \bar{y}}{n_g - np^2} \cdot \frac{1}{\frac{1}{n}} = \frac{p \bar{y}_g - p \bar{y}}{p - p^2} = \frac{\bar{y}_g - \bar{y}}{1-p}$$

$$= \frac{\bar{y}_g - (\bar{y}_g p + \bar{y}_r (1-p))}{1-p} = \frac{(1-p) \bar{y}_g - (1-p) \bar{y}_r}{1-p} = \bar{y}_g - \bar{y}_r$$

$$b_0 = \bar{y} - b_1 \bar{x} = p \bar{y}_g + (1-p) \bar{y}_r - (\bar{y}_g - \bar{y}_r) p = (1-p) \bar{y}_r + p \bar{y}_r = \bar{y}_r \quad \checkmark$$

Visual



Note: X_{row} has $L=2$

but only one intercept is necessary. why? Intercept takes place of the red category. This is called a "reference level". Red is the reference level. b_1 is effect of green rel. to red. Not "overall effect of red".

What if $x_{\text{color}} = \{\text{red, green, blue}\}$

We can create two dummies $x_1 = \mathbb{1}_{\text{green}}^{x_{\text{color}}}$, $x_2 = \mathbb{1}_{\text{blue}}^{x_{\text{color}}}$,

If $x_1 = x_2 = 0 \Rightarrow x_{\text{color}} = \text{red}$. Red is the reference category.

$\mathcal{H} = \{w_0 + w_1 x_1 + w_2 x_2 : w_0, w_1, w_2 \in \mathbb{R}\}$ $p=2!$ Even though one variable.
The answer will be the same

$$b_0 = \bar{y}_r$$

$$b_1 = \bar{y}_g - \bar{y}_r$$

$$b_2 = \bar{y}_b - \bar{y}_r$$

Can $x_1=1$ & $x_2=2$ at the same time? No...

So you never add $b_0 + b_1 + b_2$ all together.

This can be generalised to any categorical variable.

Does it matter if variable is nominal or ordinal?

No... except if you force the relationship in y to be monotonic too...

if $x_{\text{size}} = \{\text{low, medium, high}\}$

↑
set reference category

$$\hat{y}(\text{low}) = \bar{y}_L$$

$$\hat{y}(\text{medium}) = \bar{y}_M$$

$$\hat{y}(\text{high}) = \bar{y}_H$$

if I force $\hat{y}(\text{low}) \leq \hat{y}(\text{medium}) \leq \hat{y}(\text{high})$,

this OLS alg. does not always work.

Hard problems!

Consider two r.v.s X, Y . They are dep. if

$$\exists x_1, x_2 \text{ s.t. } P(Y/X=x_1) \neq P(Y/X=x_2)$$

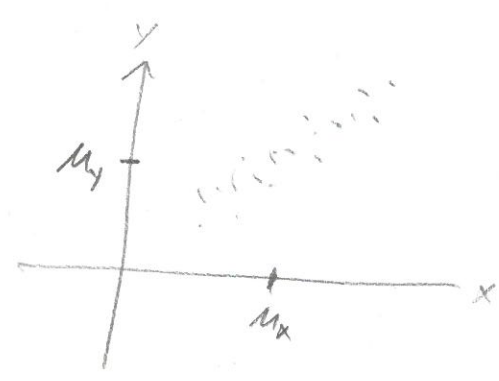
AKA "independent"

Remember: $\rho, r, \sigma_{xy}, s_{xy}$

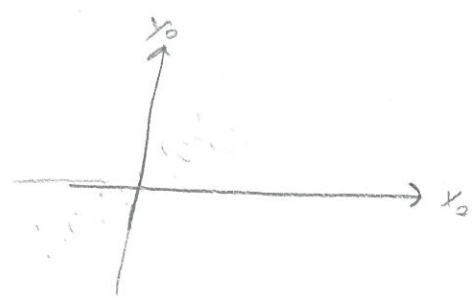
Conclusion

$$\rho = \text{Corr}[X, Y] := \frac{\sigma_{xy}}{\sigma_x \sigma_y} \quad \text{estimated by } r :=$$

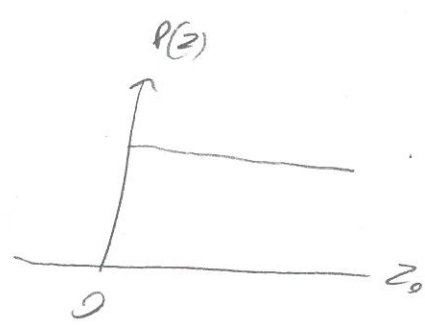
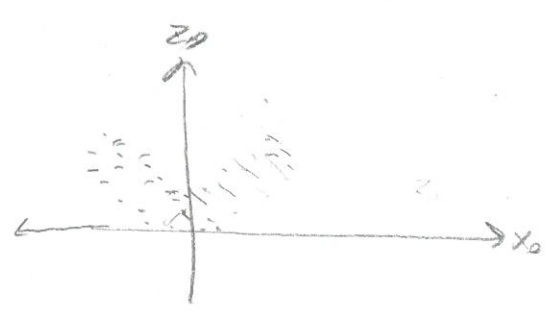
$$\sigma_{xy} := E[(X - \mu_x)(Y - \mu_y)] \quad \text{estimated by } s_{xy}$$



\Rightarrow



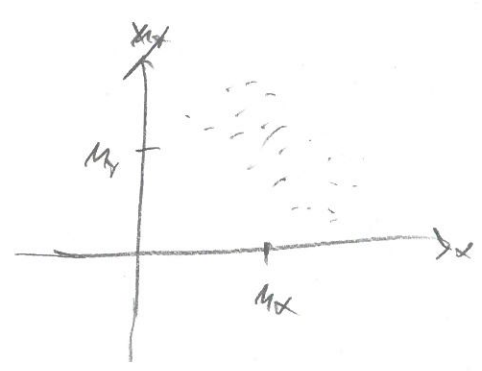
$$\text{let } Z_0 = X_0 Y_0$$



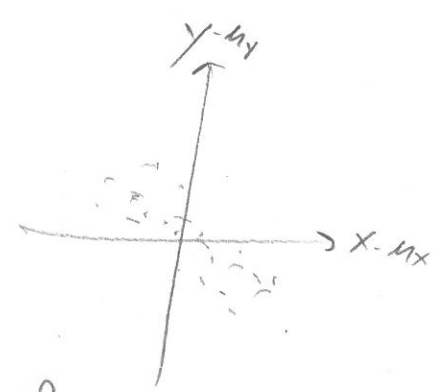
$$\sigma_{xy} := E[(X - \mu_x)(Y - \mu_y)] > 0$$

$$E[Z] > 0$$

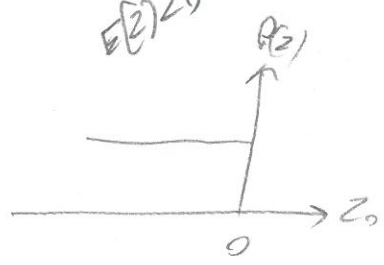
if lin. ass & + slope



\Rightarrow

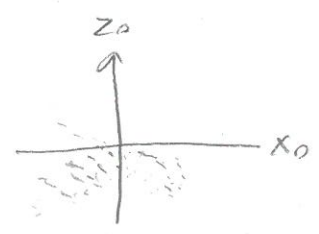


$$E[Z] < 0$$



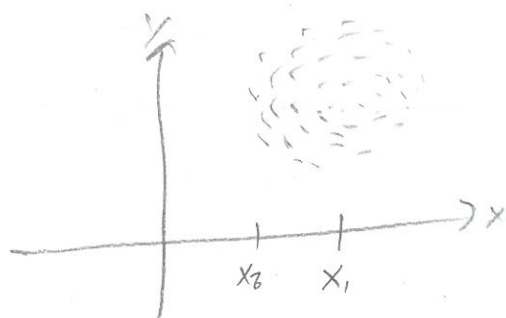
$$\sigma_{xy} < 0$$

if lin. ass & - slope

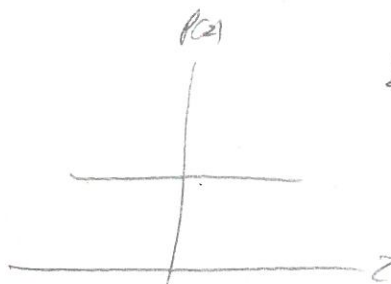
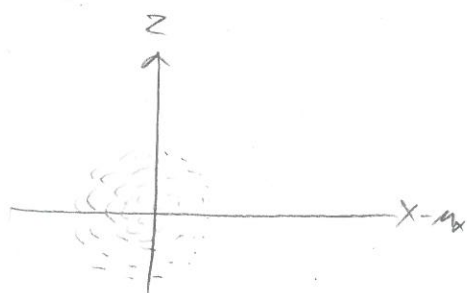
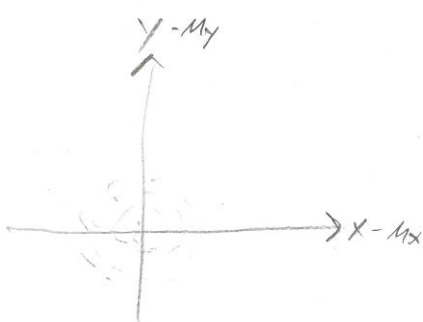
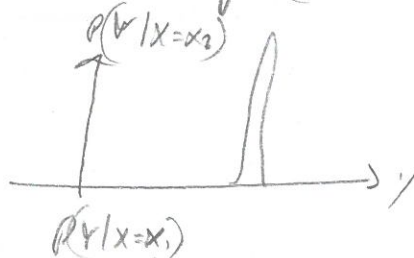


\Rightarrow

What if



Clearly dep. (& assoc.)



$$E(Z) \approx 0 \Rightarrow \sigma_{XY} \approx 0$$

assoc. \nRightarrow correlation (linear assoc.)
 better name for "corr."

corr \in {association type}

$$\rho \in [-1, 1], \quad r \in [-1, 1]$$

For $p=1$

$$R^2 = r^2 \quad (\text{that's where it gets its name from})$$

\uparrow
 on HW.

minimum \uparrow
 maximum \downarrow

$$y \in \mathbb{R}$$

(5)

OLS with $p=2 \Rightarrow \mathcal{H} = \{w_0 + w_1 x_1 + w_2 x_2 : w_0, w_1, w_2 \in \mathbb{R}\}$

$$SSE = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - w_0 - w_1 x_{i1} - w_2 x_{i2})^2 = \dots \text{long} \dots$$

Then $\argmin_{w_0} \{SSE\}$, $\argmin_{w_1} \{SSE\}$, $\argmin_{w_2} \{SSE\}$ lots of derivatives

Easier way?

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

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

$$\vec{e} = \vec{y} - \hat{\vec{y}}$$

$$(a+b)^T = a^T + b^T$$

Result $\sum e_i^2 = \vec{e}^T \vec{e} = (\vec{y} - \hat{\vec{y}})^T (\vec{y} - \hat{\vec{y}}) = (\vec{y}^T - \hat{\vec{y}}^T) (\vec{y} - \hat{\vec{y}})$

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

$$= \vec{y}^T \vec{y} - 2 \hat{\vec{y}}^T \vec{y} + \hat{\vec{y}}^T \hat{\vec{y}} \quad (a^T b)^T = b^T a$$

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

Now we need to do argmin which means we need $\frac{\partial}{\partial w_0} [SSE] \dots$

Imagine taking deriv. of whole vector

16

$$\frac{\partial \text{SSE}}{\partial \vec{w}} \stackrel{\text{scalar}}{=} \begin{bmatrix} \frac{\partial}{\partial w_0} [\text{SSE}] \\ \frac{\partial}{\partial w_1} [\text{SSE}] \\ \frac{\partial}{\partial w_2} [\text{SSE}] \end{bmatrix} \quad \text{"the vector derivative"}$$

Let's get some rules for this..

let \vec{x} have length n

$$\frac{\partial}{\partial \vec{x}} [a] \quad \text{where } a \text{ is a constant wrt } \vec{x}$$

$$= \begin{bmatrix} \frac{\partial}{\partial x_1} [a] \\ \vdots \\ \frac{\partial}{\partial x_n} [a] \end{bmatrix} = \vec{0}_n$$

let $\vec{a} \in \mathbb{R}^n$ const wrt \vec{x}

$$\frac{\partial}{\partial \vec{x}} [\underbrace{\vec{a}^T \vec{x}}_{\vec{x}^T \vec{a}}] = \begin{bmatrix} \frac{\partial}{\partial x_1} [a_1 x_1 + a_2 x_2 + \dots + a_n x_n] \\ \frac{\partial}{\partial x_2} [] \\ \vdots \\ \frac{\partial}{\partial x_n} [] \end{bmatrix} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} = \vec{a} \neq \vec{a}^T$$

NOTE: different dimension

let b constant wrt to x

$$\frac{\partial}{\partial \vec{x}} [a f(\vec{x}) + b g(\vec{x})] = \begin{bmatrix} \frac{\partial}{\partial x_1} [a f(\vec{x}) + b g(\vec{x})] \\ \vdots \\ \frac{\partial}{\partial x_n} [a f(\vec{x}) + b g(\vec{x})] \end{bmatrix} = \begin{bmatrix} a \frac{\partial}{\partial x_1} [f(\vec{x})] + b \frac{\partial}{\partial x_1} [g(\vec{x})] \\ \vdots \\ a \frac{\partial}{\partial x_n} [f(\vec{x})] + b \frac{\partial}{\partial x_n} [g(\vec{x})] \end{bmatrix} = a \frac{\partial}{\partial \vec{x}} [f(\vec{x})] + b \frac{\partial}{\partial \vec{x}} [g(\vec{x})]$$

quadratic form

let $A \in \mathbb{R}^{n \times n}$ constant wrt \vec{x}

$$\frac{\partial}{\partial \vec{x}} [\vec{x}^T A \vec{x}] \text{ \& symmetric}$$

$$A \vec{x} = \begin{bmatrix} \vec{a}_1 \rightarrow \\ \vec{a}_2 \rightarrow \\ \vdots \\ \vec{a}_n \rightarrow \end{bmatrix} \begin{bmatrix} 1 \\ \vec{x} \\ 1 \end{bmatrix} = \begin{bmatrix} \vec{a}_1 \cdot \vec{x} \\ \vec{a}_2 \cdot \vec{x} \\ \vdots \\ \vec{a}_n \cdot \vec{x} \end{bmatrix}$$

$$A \vec{x} = \begin{bmatrix} a_{11} x_1 + a_{12} x_2 + \dots + a_{1n} x_n \\ a_{21} x_1 + a_{22} x_2 + \dots + a_{2n} x_n \\ \vdots \\ a_{n1} x_1 + a_{n2} x_2 + \dots + a_{nn} x_n \end{bmatrix}$$

$$X^T(A\vec{x}) = [x_1 \ x_2 \ \dots \ x_n] \begin{bmatrix} \vec{q}_1 \cdot \vec{x} \\ \vec{q}_2 \cdot \vec{x} \\ \vdots \\ \vec{q}_n \cdot \vec{x} \end{bmatrix} = x_1 \vec{q}_1 \cdot \vec{x} + x_2 \vec{q}_2 \cdot \vec{x} + \dots + x_n \vec{q}_n \cdot \vec{x}$$

$$x_1(q_{11}x_1 + q_{12}x_2 + \dots + q_{1n}x_n) + x_2(q_{21}x_1 + q_{22}x_2 + \dots) + \dots +$$

$$\frac{\partial}{\partial x_1} [\] = \frac{\partial}{\partial x_1} [\] = (2q_{11}x_1 + q_{12}x_2 + \dots + q_{1n}x_n) + q_{21}x_2 + \dots + q_{n1}x_n$$

Same due to symmetry

$$= 2(q_{11}x_1 + q_{12}x_2 + \dots + q_{1n}x_n) = 2\vec{q}_1 \cdot \vec{x}$$

$$\frac{\partial}{\partial x_2} [\] = q_{12}x_1 + (q_{21}x_1 + 2q_{22}x_2 + q_{23}x_3 + \dots + q_{2n}x_n) + q_{32}x_3 + \dots + q_{n2}x_n$$

$$= 2(q_{21}x_1 + q_{22}x_2 + \dots + q_{2n}x_n)$$

$$= 2\vec{q}_2 \cdot \vec{x}$$

$$\frac{\partial}{\partial \vec{x}} [x^T A x] = 2 \begin{bmatrix} \vec{q}_1 \cdot \vec{x} \\ \vec{q}_2 \cdot \vec{x} \\ \vdots \\ \vec{q}_n \cdot \vec{x} \end{bmatrix} = 2A\vec{x}$$

Math lesson over...

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

Symmetric? $(X^T X)^T = X^T (X^T)^T = X^T X \checkmark$

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

const 3x3

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

$$\Rightarrow (X^T X) \vec{w} = X^T \vec{y}$$

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

OLS estimate

Proof valid for all p .

We assume $X^T X$ is invertible. When? only when

$\text{rank}(X) = p+1$ Note: $\text{rank}(X) \neq p+1$ can't be greater!

" $\dim[\text{colsp}(X)] \neq p+1$ only $p+1$ cols!"



We want to prove $\mathbf{X}^T \mathbf{X}$ is invertible only when $\text{rank}[\mathbf{X}] = p + 1$ i.e. the design matrix is “full rank”. This is equivalent to proving that

$$\text{rank}[\mathbf{X}] = p + 1 \implies \text{rank}[\mathbf{X}^T \mathbf{X}] = p + 1$$

Logically equivalent is the contrapositive:

$$\text{rank}[\mathbf{X}^T \mathbf{X}] \neq p + 1 \implies \text{rank}[\mathbf{X}] \neq p + 1$$

Not equal in this case is equivalent to less than because a matrix cannot have a rank that exceeds its number of columns since $\text{rank}[\mathbf{X}] := \dim \text{colsp}[\mathbf{X}]$ so the above is equivalent to:

$$\text{rank}[\mathbf{X}^T \mathbf{X}] < p + 1 \implies \text{rank}[\mathbf{X}] < p + 1$$

Beginning with the premise on the left hand side, a rank-deficient matrix has at least one non-trivial (i.e. non-zero) vector $\mathbf{v} \in \mathbb{R}^{p+1}$ that maps to the zero vector, i.e. there is at least one direction in the nullspace:

$$\mathbf{X}^T \mathbf{X} \mathbf{v} = \mathbf{0}_{p+1}$$

Saqib noticed we can multiply both sides on the right by \mathbf{v}^T to arrive at:

$$\mathbf{v}^T \mathbf{X}^T \mathbf{X} \mathbf{v} = \mathbf{v}^T \mathbf{0}_{p+1} = 0$$

We can write this alternatively as

$$(\mathbf{X} \mathbf{v})^T \mathbf{X} \mathbf{v} = 0 \implies \sum_{i=1}^n (\mathbf{X} \mathbf{v})_i^2 = 0$$

If all elements of $\mathbf{X} \mathbf{v}$ squared and summed yield zero, every single element must be zero and thus,

$$\mathbf{X} \mathbf{v} = \mathbf{0}_n$$

indicating that the vector \mathbf{v} (which was assumed to be nontrivial above) is in the nullspace of \mathbf{X} indicating that \mathbf{X} is rank deficient and thus $\text{rank}[\mathbf{X}] < p + 1$. ■