Lecture 16  Math 390.03-02  4/6/16

n draws from bivariate distr $\langle X, Y \rangle$ — (response)/outcome/dep. var.

$\langle X_1, Y_1 \rangle, \langle X_2, Y_2 \rangle, \ldots, \langle X_n, Y_n \rangle$ — (covariate)/regressor/indep. var.

too easy
Control w/t
dep/indep in
prob.

$$X \xrightarrow[\text{through}]{\text{affects}\atop\text{change}} Y$$

a mechanism
$f$ and error $\varepsilon$

$$Y = f(X) + \varepsilon$$

we limit $\quad f \in \mathcal{F}_{lin} := \{ \beta_0 + \beta_1 x \text{ s.t. } \beta_0, \beta_1 \in \mathbb{R} \} \subset \mathcal{F}$

all possible
functions
of $X$ is
too big of a
space to work
with

and $\quad \varepsilon \neq h(X)$

not a function of $X$

Forget everything about r.v.'s and prob for

a moment...

$$\left. \begin{array}{l} y_1 = \beta_0 + \beta_1 x_1 \\ y_2 = \beta_0 + \beta_1 x_2 \\ \quad \vdots \\ y_n = \beta_0 + \beta_1 x_n \end{array} \right\}$$

overdetermined
system
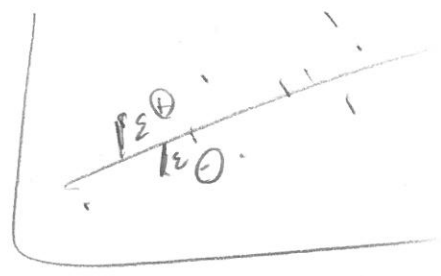need $\varepsilon$ so
make these
into eq's.

Gauss (1795)

Legendre (1805)

Let's try to figure out $\beta_0$ and $\beta_1$



Let's try to minimize the error terms

$$\hat{\beta_0}, \hat{\beta_1} = \arg\min_{\beta_0, \beta_1} \sum_{i=1}^{n} l(\varepsilon_i)$$

loss function

Laplace's idea...
non-unique

$L_1 : l_1(\varepsilon) = |\varepsilon|$   absolute loss

$L_2 : l_2(\varepsilon) = \varepsilon^2$   sqd err / quadrature loss

$L_{0-1} : l_{0-1}(\varepsilon) = \begin{cases} 0 & \text{if } |\varepsilon| \le c \\ 1 & \text{if } |\varepsilon| > c \end{cases}$

lets go with $L_2$ loss

error is diff between lin. mod. and actual $y$

$$\hat{\beta_0}, \hat{\beta_1} = \arg\min_{\beta_0, \beta_1} \sum_{i=1}^{n} \varepsilon_i^2 = \sum_{i=1}^{n} \left( y_i - (\beta_0 + \beta_1 x_i) \right)^2$$

$-\beta_0 - \beta_1 x_i$

$y_i^2 - 2y_i \bar{y} + \bar{y}^2$

$(n-1) s_y^2 = \sum (y_i - \bar{y})^2$
$= \sum y_i^2 - 2\bar{y} \sum y_i + \sum \bar{y}^2$
$= \sum y_i^2 - 2n\bar{y}^2 + n\bar{y}^2$
$= \sum y_i^2 - n\bar{y}^2$

$= \sum y_i^2 - 2y_i\beta_0 - 2y_i x_i \beta_1 + \beta_0^2 2\beta_0\beta_1 x_i + \beta_1^2 x_i^2$

$= \sum y_i^2 + \beta_1^2 \sum x_i^2 + n\beta_0^2 - 2\bar{y}n\beta_0 - 2\beta_1 \sum y_i x_i + 2\beta_0\beta_1 \bar{x} n$

$\frac{\partial}{\partial \beta_0} = -2\bar{y}n - 2\sum y_i x_i + 2n\beta_0 + \beta_1 \bar{x} n \stackrel{=0}{\Rightarrow} b_0 = \bar{y}n + \sum y_i x_i - b_1 \bar{x} n$

$\frac{\partial}{\partial \beta_1} = 2\beta_1 \sum x_i^2 - 2\sum y_i x_i + 2\beta_0 \bar{x} n \stackrel{=0}{\Rightarrow} \beta_1 = \sum y_i x_i - b_0 \bar{x} n$

$$r := \frac{s_{xy}}{s_x s_y} \Rightarrow r s_x s_y = s_{xy} := \sum (y_i - \bar{y})(x_i - \bar{x}) = \sum x_i y_i - x_i \bar{y} - y_i \bar{x} + \bar{y}\bar{x}$$

I mean (n-1)sxy = here

$$= \sum y_i x_i - n\bar{x}\bar{y} + n\bar{y}\bar{x} + n\bar{y}\bar{x}$$

$$\sum \varepsilon_i^2 = s_y^2 + n\bar{y}^2 + \beta_1^2 s_x^2 + \beta_1^2 n\bar{x}^2 + n\beta_0^2 - 2\bar{y}n\beta_0 - 2\beta_1 s_{xy} + 2\beta_1 n\bar{y}\bar{x} + 2\beta_0\beta_1 \bar{x}n$$

This should be (n-1)s^2_y instead of just s^2_y

$$b_1 := \frac{\partial}{\partial \beta_1}( \ ) \overset{set}{=} 0$$

$$\Rightarrow 2\beta_1 (s_x^2 + n\bar{x}^2) - 2 s_{xy} - 2n\bar{y}\bar{x} + 2\beta_0 \bar{x}n = 0$$

$$\Rightarrow b_1 = \frac{s_{xy} + n\bar{y}\bar{x} - \beta_0 \bar{x}n}{s_y^2 + n\bar{x}^2} = \frac{s_{xy} + n\bar{y}\bar{x} - (\bar{y} - b_1\bar{x})\bar{x}n}{s_y^2 + n\bar{x}^2}$$

$$b_0 := \frac{\partial}{\partial \beta_0}( \ ) \overset{set}{=} 0$$

$$\Rightarrow 2n\beta_0 - 2\bar{y}n + 2\beta_1 \bar{x}n = 0$$

$$b_0 = \frac{\bar{y}n - \beta_1 \bar{x}n}{n} = \bar{y} - b_1 \bar{x}$$

$$\Rightarrow s_{xy} + n\bar{y}\bar{x} - n\bar{x}\bar{y} + b_1 \bar{x}^2 n = b_1 s_y^2 + b_1 n\bar{x}^2$$

$$\Rightarrow b_1 = \frac{s_{xy}}{s_y^2} = \frac{r s_x s_y}{s_y^2} = r\frac{s_x}{s_y}$$

The (n-1)'s cancel out between numerator and denominator here

"L.S. estimates"

Gauss's method of L.S.

— line of least fit

$$\hat{Y}_i = b_0 + b_1 x_i$$

"$\beta_1$"

best est of $\beta_0$

$$b_0 \overset{?}{=} \beta_0 \ , \ b_1 \overset{?}{=} \beta_1 \quad No...$$

Now... r.v. $\langle X, Y \rangle$ joint density $\#$.

$$Y = \beta_0 + \beta_1 X + \varepsilon \quad \text{where } E[\varepsilon] = 0$$

$$E(Y) = \underset{X}{E}[\underset{Y}{E}[Y|X]] = E_X[\beta_0 + \beta_1 X + 0] = \beta_0 + \beta_1 E(X)$$

What if $\varepsilon_1, \ldots, \varepsilon_n \overset{iid}{\sim} N(0, \sigma^2)$ $\left(\overset{Gore}{OLS}\ \text{Assumption}\right)$

$\left(\begin{array}{c}\text{ord}\\\text{lest}\\\text{sqs}\end{array}\right)$ $\left\{\begin{array}{l}\text{①} \text{Normality of errors}\\\text{②} \text{Indep of errors}\\\text{③} \text{Homoskedasticity} \text{ all var} = \sigma^2\end{array}\right.$

Condition on $X$ to keep things simple (real order terms — PhD in Stats $\text{req'd}$)

independent of $X$ ~~$\longrightarrow$~~

$$Y \mid X \sim N\left(\beta_0 + \beta_1 X, \ \sigma^2\right)$$

Why is normality reasonable?   CLT

Consider $Y = \alpha + \beta_1 X + \beta_2 X_2 + \ldots + \beta_p X_p$ (no $\varepsilon$ here)

when $X_2, \ldots, X_p$ are unobserved but i.i.d.

$$\Rightarrow E[Y \mid X] = \beta_1 X + \underbrace{\left(\alpha + \sum_{i=2}^{p} \beta_i E[X_i]\right)}_{\beta_0}$$

$$Var[Y \mid X] = \underbrace{\sum_{i=2}^{p} \beta_i^2 Var[X_i]}_{\sigma^2}$$

$$\alpha + \beta_2 X_2 + \ldots + \beta_p X_p \overset{d}{\approx} N(\beta_0, \sigma^2)$$

let $\varepsilon = \sum_{i=2}^{p} \beta_i X_i - \beta_0$   Noise is the result of not seeing a bunch of stuff

(philosophical point)

What is $\hat{\beta}_0$ MLE, $\hat{\beta}_1$ MLE ?

$$\mathcal{L}(\beta_i | \overset{x_i}{x}, y) = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}\left(y_i - (\beta_0 + \beta_1 x_i)\right)^2}$$

$$= \left(\frac{1}{\sqrt{2\pi\sigma^2}}\right)^n e^{-\frac{1}{2\sigma^2} \sum (y_i - b_0 - \beta_1 x_i)^2}$$

$$\ell(\beta_i | x, y) = \ln\left(\overset{b_0}{\mathcal{L}}\right) - \frac{1}{2\sigma^2}$$

$$\ell'\left(\overset{b_0}{\beta_1} | x, y\right) = -\frac{1}{2\sigma^2}\left(2\beta_1 \left(s_x^2 + n\bar{x}^2\right) - 2s_{xy} - 2n\bar{y}\bar{x} + 2\beta_0 n\bar{x}\right) = 0$$

$$\ell'\left(\overset{\beta_1}{b_0} | x, y\right) = -\frac{1}{2\sigma^2}\left(2n b_0 - 2\bar{y}n + 2\beta_1 \bar{x}n\right) = 0$$

$$\hat{\beta}_0^{MLE} = b_0 \quad \text{and} \quad \hat{\beta}_1^{MLE} = b_1$$

the MLE's = the L.S. estimates !

$L_2$ loss is. $\ell(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2$ seems to be "natural" once again...

let $B_0$ be the estimator (r.v.) for $b_0$      $B_0 \sim$ ?   normal()

let $B_1$        "   "   "   "   "   "   "   "        $B_1$        $B_1 \sim$ ?

$E(B_1) = \ldots b_1$      $E(B_0) = \ldots \beta_0$  (unbiased)

We will see soon why this is...

Same thing with the (n-1) terms when I did this derivation in class

This works if $Y \in \mathbb{R}$ ... What about $Y \in \{0,1\}$ i.e. binary?

$P(Y=1|x) \in (0,1)$   by def of the param. space of the Bern. model   Does it matter if $P(Y=0|x)$ ?

$Y \sim Bern\left(P(Y=1|x)\right)$

No, $0/1$ arbitrary ... labels can be flipped. Usually you know what you want to predict

model

linear models $\in \mathbb{R}$

prob $\in [0,1]$

What to do? Introduce "link function", the most famous being:

$$\ell := logit\left(P(Y=1|x)\right) := \ln\left(\frac{P(Y=1|x)}{1-P(Y=1|x)}\right) = \beta_0 + \beta_1 x$$

$$P(Y=1|x) = logit^{-1}(\ell) = \frac{e^\ell}{1-e^\ell} \quad \Rightarrow \quad 1-P(Y=1|x) = 1-\frac{e^\ell}{1-e^\ell} = \frac{1}{1-e^\ell}$$

The denominator should be 1+e^ell ... carry that through everywhere...

Recall $Y_1, \ldots, Y_n \overset{iid}{\sim} Bern(\theta)$

$$\mathcal{L}(\theta; y) = \prod_{i=1}^n \theta^{\Sigma y_i}(1-\theta)^{1-\Sigma y_i} = \theta^{\Sigma y_i}(1-\theta)^{n-\Sigma y_i} \propto \left(\frac{\theta}{1-\theta}\right)^{\Sigma y_i}$$

$$\Rightarrow$$

$$\mathcal{L}(\beta_0,\beta_1; y, x) = \prod_{i=1}^n \left(\frac{e^{\ell_i}}{1-e^{\ell_i}}\right)^{y_i}\left(\frac{1}{1-e^{\ell_i}}\right)^{1-y_i} = \prod_{i=1}^n (1-e^{\ell_i})^{-1} \prod_{i=1}^n e^{\ell_i y_i} = \prod_{i=1}^n (1-e^{\beta_0+\beta_1 x_i})^{-1} e^{\sum_{i=1}^n y_i(\beta_0+\beta_1 x_i)}$$