

Lecture 17

Topics:

I) Bivariate Normal prediction given/predicting z-scores, raw scores or percentiles.

II) Start estimation.

Since raw scores, percentiles & z-scores are all essentially equivalent ways of specifying the same thing prediction for BV (X, Y) can apply to all three.

e.g. Midterm 1 scores: X

$$E(X) = 10.62, \text{sd}(X) = 1.46$$

Midterm 2 scores: Y

$$E(Y) = 8.54, \text{sd}(Y) = 2.08$$

Correlation b/w X & Y

$$\rho = 0.155$$

Question: Suppose you got $X=8$ on midterm 1. what is your predicted percentile for Y among

(a) the whole class

(b) the subgroup of students who also got $X=8$ on Midterm 1.

Answer: First compute the predicted z-score for Y (among the whole class).

$$(\text{predicted z-score}) = p(\text{observed z-score})$$

$$= 0.155 \left(\frac{8 - E(X)}{\text{sd}(X)} \right)$$

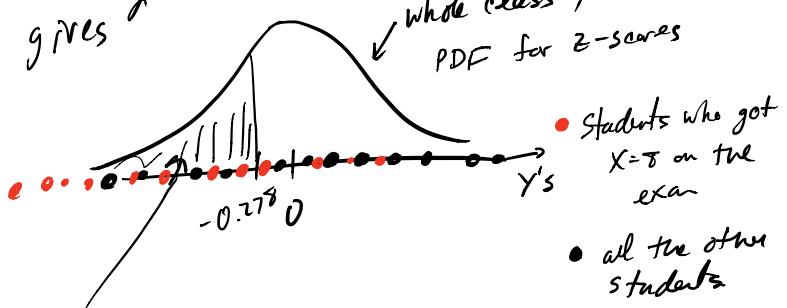
$$= 0.155 \left(\frac{8 - 10.62}{1.46} \right)$$

$$= -0.278$$

$$\text{since predicted z-score} = \frac{E(Y|X=8) - E(Y)}{\text{sd}(Y)}$$

-0.278 describes the z-score prediction among the whole class.

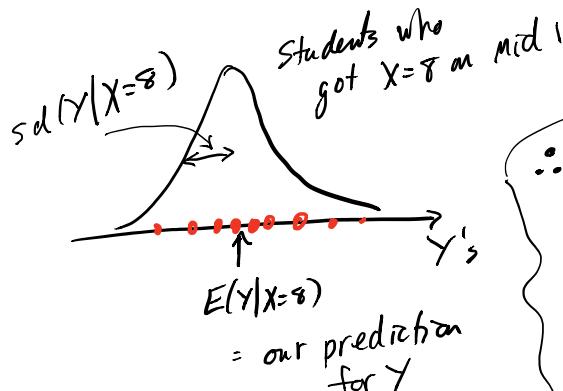
Relating z-scores to percentiles gives



This area = 0.39

∴ for (a) the predicted percentile on Midterm 2 (Y) is 39th percentile.

What about (b)?



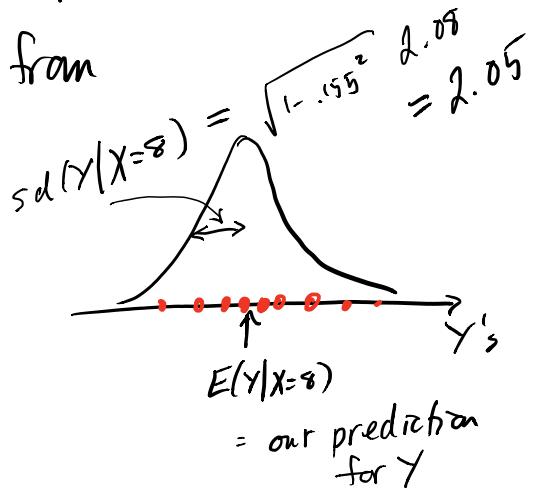
∴ Among these students the predicted z-score is 50th percentile

e.g. Same problem from (3)
last example.

given $X=8$ compute

$$P(Y \geq 10 | X=8).$$

Given $X=8$, Y behaves like a random draw from



$$\begin{aligned} &= (-0.278)(2.08) + 8.54 \\ &\quad \uparrow \quad \uparrow \quad \uparrow \\ &\quad \text{predicted} \quad \text{sd}(Y) \quad E(Y) \\ &\quad z\text{-score} \\ &= 7.962 \end{aligned}$$

i.e. given $X=8$, Y behaves like a random draw from $N(7.962, 2.05^2)$.

$$\begin{aligned} \therefore P(Y \geq 10 | X=8) &= P\left(Z \geq \frac{10 - 7.962}{2.05}\right) \\ &\quad \uparrow \\ &\quad z \sim N(0,1) \\ &= P(Z \geq 0.994) \\ &= 0.16. \end{aligned}$$

Basis of estimation

(4)

A typical estimation problem goes like this ... I'm interested in the average house square footage of households in California (μ) I sampled 45 homes at random, measured their sq footage (denoted X_1, \dots, X_{45}). Since $\bar{X} = 2857 \text{ ft}^2$ I'm estimating μ to be 2857 ft^2 .

Now, clearly μ isn't exactly 2857 ft^2 how far off could it be?

trying to quantify this is the main goal of this section.

Note $\mu - \bar{X}$ = estimation error.

There are two possible sources of est error.

(1) Sampling error (due to the randomness in \bar{X}).

Quantified by $\text{var}(\mu - \bar{X})$

(2) systematic error (also called bias)

Quantified by $E(\mu - \bar{X})$.

These combine to give an overall estimate of the typical size of the estimator error:

typical size of the estimation error

(5)

$$= \sqrt{E[(\mu - \bar{x})^2]} \rightarrow \text{called Mean squared error}$$

$$\text{fact} = \sqrt{\underbrace{\text{var}(\mu - \bar{x})}_{\text{quantifies Sampling variability}} + \underbrace{(E(\mu - \bar{x}))^2}_{\text{bias}}}$$

$$= \sqrt{\text{var}(\bar{x}) + (E(\mu - \bar{x}))^2}$$

Since $E(\bar{x}) = \mu$ & $\text{var}(\bar{x}) = \frac{\sigma^2}{n}$ variance of each x_i

for this example we have:

the typical estimation error

$$= \sqrt{E[(\mu - \bar{x})^2]}$$

$$= \sqrt{\text{var}(\bar{x}) + \sigma^2}$$

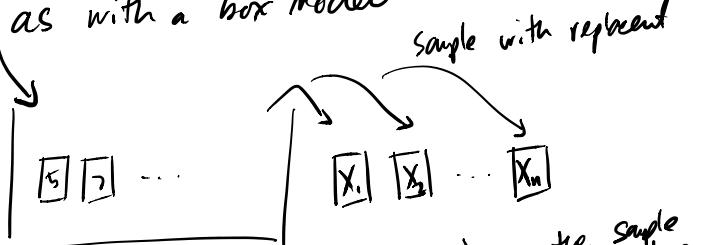
$$= \frac{\sigma}{\sqrt{n}}$$

so if $\sigma = 1000 \text{ ft}^2$ then $\mu - \bar{x}$ is typically off by $\pm \frac{1000}{\sqrt{45}} = 149$.
 & if $\sigma = 200 \text{ ft}^2$ then $\mu - \bar{x}$ is typically off by $\pm \frac{200}{\sqrt{45}} = 29.8$

The above example was special... let's now outline the general case

(6)

A population can be represented as with a box model.



Let θ denote a population parameter, which is computed from these tickets, that you want to estimate.

$\theta - \hat{\theta} = \text{estimation error}$

$\text{bias}(\hat{\theta}) = E(\theta - \hat{\theta})$

= systematic error

$\text{var}(\hat{\theta}) = \text{sampling uncertainty}$

$MSE(\hat{\theta}) = E[(\theta - \hat{\theta})^2]$

= mean squared estimation error

typical estimation error

$$= \sqrt{MSE(\hat{\theta})}$$

Fact:

$$MSE(\hat{\theta}) = \text{var}(\hat{\theta}) + (\text{bias}(\hat{\theta}))^2$$

(7)

combines the
two possible
sources of estimation
error to compute
"typical estimation
error".

(8)