

Ch 12 Experimental Design

12/1

Ex 12.2 Increasing accuracy

In parameter estimation, more accuracy means smaller standard error.

As we once saw in Ch 9, usually the only thing we have control over is n .

But in practice, n can't be infinite and maybe we can't even control it either.

What can we do?

ex: Let $\hat{\theta} = \bar{X}_1 - \bar{X}_2$ difference in means

For large sample $V(\hat{\theta}) = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$

and standard error is $\sqrt{V(\hat{\theta})} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$

Assume you can sample n times.

How do we choose n_1, n_2 ?

Use $n_1 + n_2 = n$ or $n_2 = n - n_1$.

The variance as a function of n_1 is

$$V(\hat{\theta}) = \frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n - n_1}$$

We want to minimize V .

$$\begin{aligned}\frac{dV}{dn_1} = V'(n_1) &= \frac{-\sigma_1^2}{n_1^2} - \frac{\sigma_2^2}{(n-n_1)^2} \quad (-) \\ &= -\frac{\sigma_1^2}{n_1^2} + \frac{\sigma_2^2}{(n-n_1)^2}.\end{aligned}$$

1200.

$$\text{Note } V''(n_1) = \frac{2\sigma_1^2}{n_1^3} + \frac{2\sigma_2^2}{(n-n_1)^3} > 0.$$

So any critical number is a local min.

$$V'(n_1) = 0 \Rightarrow -\sigma_1^2(n-n_1)^2 + \sigma_2^2 n_1^2 = 0$$

$$-\sigma_1^2 n^2 + 2\sigma_1^2 n n_1 - \sigma_1^2 n_1^2 + \sigma_2^2 n_1^2 = 0$$

$$(\sigma_2^2 - \sigma_1^2)n_1^2 + 2n\sigma_1^2 n_1 - n^2\sigma_1^2 = 0.$$

$$\begin{aligned}\text{discriminant } (2n\sigma_1^2)^2 - 4(\sigma_2^2 - \sigma_1^2)(-n^2\sigma_1^2) \\ &= 4n^2\sigma_1^4 + 4n^2\sigma_1^2\sigma_2^2 - 4n^2\sigma_1^2 \\ &= 4n^2\sigma_1^2\sigma_2^2.\end{aligned}$$

$$\text{and } n_1 = \frac{-2n\sigma_1^2 \pm 2n\sigma_1\sigma_2}{2(\sigma_2^2 - \sigma_1^2)} = \frac{n\sigma_1(\sigma_1 \pm \sigma_2)}{(\sigma_2^2 - \sigma_1^2)}$$

WLOG, let $\sigma_2^2 > \sigma_1^2$. Then we require $\sigma_1 \neq \sigma_2$ term in top to have $n_1 > 0$.

$$n_1 = \frac{+n\sigma_1(\sigma_2 - \sigma_1)}{\sigma_2^2 - \sigma_1^2} = \left(\frac{\sigma_1}{\sigma_2 + \sigma_1} \right) n.$$

10p3

Note: This implies that to minimize V for difference in means we ~~pick~~ ^{sample} more of the type that has inherently large variance

$$n_1 \approx \frac{\sigma_1}{\sigma_1 + \sigma_2} n, \quad n_2 \approx \frac{\sigma_2}{\sigma_1 + \sigma_2} n$$

ex: best-fit line $y = \beta_0 + \beta_1 x + \epsilon$.

Recall a C.I for β_i is $\hat{\beta}_i \pm t^* \sqrt{V(\hat{\beta}_i)}$
 $\Rightarrow \hat{\beta}_i \pm t^* c_{ii} \sigma^2$

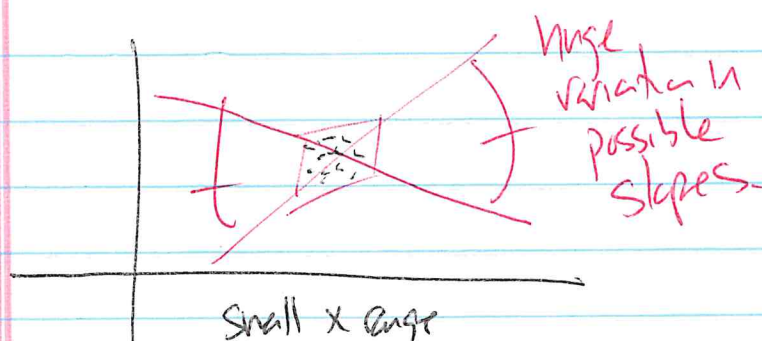
where c_{ii} is from $(X^T X)^{-1} = \begin{bmatrix} \sum x_i^2 / n S_{xx} & -\bar{x} / S_{xx} \\ -\bar{x} / S_{xx} & 1 / S_{xx} \end{bmatrix}$

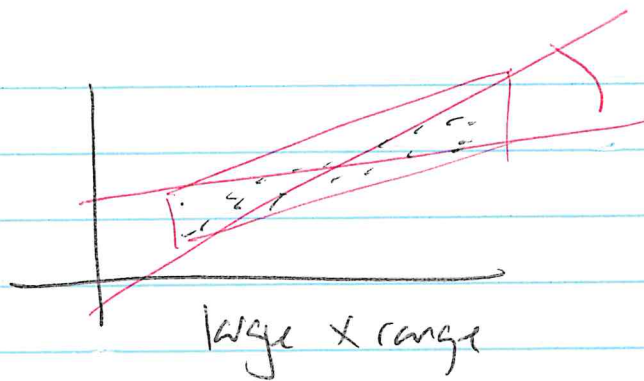
For estimating β_1 , $V(\hat{\beta}_1) = c_{11} \sigma^2 = \frac{\sigma^2}{S_{xx}}$

where $S_{xx} = \sum (x_i - \bar{x})^2$.

* The more spread in x , the better our estimator for slope is.

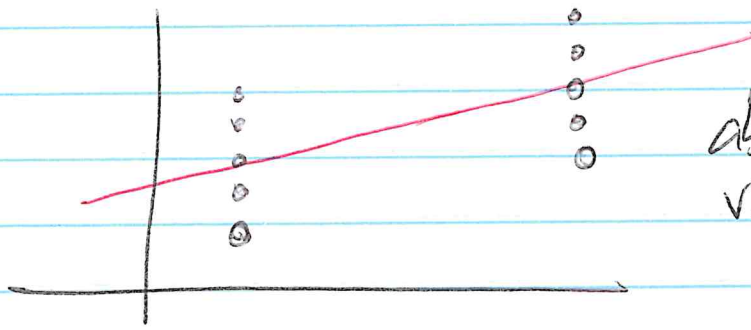
Which kind makes sense...





10p4
smaller possible
range of slope.

But also feels weird... something like



absolutely minimizes
variance.

§ 10.3 Matched-Pairs Experiment

10p5

Real-world experiments often measure a r.v. X
pre- and post-treatment
 (X_1) (X_2)

As these measurements are on the same item or individual, clearly don't expect independence.

So, measuring effectiveness, we can't simply look at $\bar{X}_1 - \bar{X}_2$.

We use X_{1i} and X_{2i} to indicate the pre- and post- data point of the i th individual.

We use the standard assumption that X 's are normally distributed: $X_1 \sim N(\mu_1, \sigma_1^2)$, $X_2 \sim N(\mu_2, \sigma_2^2)$

No reason to expect either $\mu_1 = \mu_2$ or $\sigma_1^2 = \sigma_2^2$.

$$\text{So } E(X_{1i}) = \mu_1, V(X_{1i}) = \sigma_1^2 \\ E(X_{2i}) = \mu_2, V(X_{2i}) = \sigma_2^2$$

$$\text{but } \text{Cov}(X_{1i}, X_{2i}) = \rho \sigma_1 \sigma_2 \\ \text{where } \rho = \frac{\text{Cov}(X_{1i}, X_{2i})}{\sigma_1 \sigma_2}, \text{ the correlation coef.} \\ \text{See Ch 5.}$$

b.p.6

The stat we measure is the difference pre-to post-
 $D_i = X_{1i} - X_{0i}$

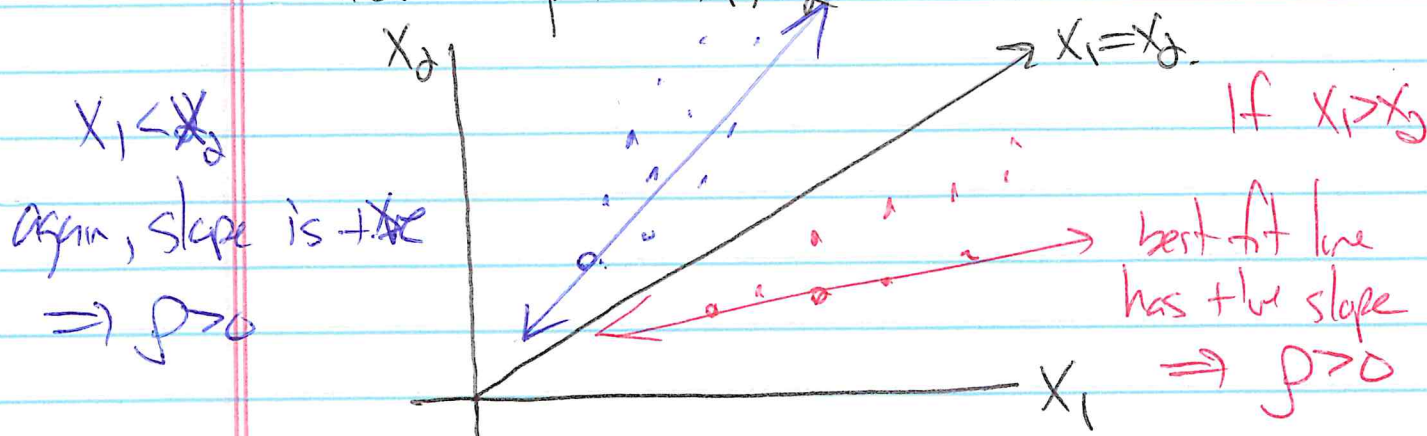
Then $E(D_i) = \mu_1 - \mu_0$
 $V(D_i) = V(X_{1i}) + V(X_{0i}) - 2\text{Cov}(X_{1i}, X_{0i})$
 $= \sigma_1^2 + \sigma_0^2 - 2\rho\sigma_1\sigma_0$

Then $D_i \sim N(\mu_1 - \mu_0, V(D_i)/n)$, the
 sampling dist'n with $n-1$ d.f.

disc: why this is better.

Commonly expect that a treatment will
 have an effect that is relatively uniform.
 either $X_1 > X_0$ or $X_1 < X_0$.

Consider the plot (X_1, X_0)



10p7

Thus $V(D_i) = \sigma_1^2 + \sigma_0^2 - 2\rho\sigma_1\sigma_0$

$< \sigma_1^2 + \sigma_0^2 \leftarrow$ the variance if X_1, X_0 were independent.

The matched pairs sbt yields a smaller confidence interval!

Summary: $D_i = X_{1i} - X_{0i} \quad 1 \leq i \leq n$

Use $t/2$ methods to construct C.I. or hypo. tests.

Recall $n > 30$, using z
 $n \leq 30$, using t .

Here $\bar{D} = \frac{\sum D_i}{n}$, $S_D^2 = \frac{1}{n-1} \sum_1^n (D_i - \bar{D})^2$
 $= \frac{1}{n-1} [\sum D_i^2 - n(\bar{D})^2]$

ex: testing effects of a supplement on cholesterol levels.

"randomly" select 6 people.

X_1 : C-level at beginning

X_0 : C-level at 6 weeks

(milligrams/deciliter)

12p8

Really common assumption for biological things:
 $X \sim N(\mu, \sigma^2)$

Person	1	2	3	4	5	6
X_1	210	235	208	190	170	244
X_d	190	170	210	188	173	228

$$D: \quad 20 \quad 65 \quad -2 \quad 2 \quad -1 \quad 16$$

$$\bar{D} = \frac{100}{6} = 16.7$$

$$D^2: \quad 400 \quad 4225 \quad 4 \quad 4 \quad 1 \quad 256$$

$$S_{DD} = \sum D_i^2 = 4890$$

$$\sum S_D^2 = 4890 - 6(16.7)^2 = 643.332$$

$$S_D \approx 25.36$$

Testing if change occurred $H_0: D=0$
 $H_a: D \neq 0$

$$t = \frac{\bar{D} - \mu_0}{S_D / \sqrt{n}} = \frac{16.7 - 0}{25.4 / \sqrt{6}} = 1.610$$

$$\begin{aligned} p\text{-value} &= P(|t| > 1.610 \mid df=5) \\ &= 2P(t > 1.610 \mid df=5) > 0.10 \end{aligned}$$

bpf.

Not significant

Can not reject H_0

No indication that the supplement caused a change in c-level.

(review)

Find a 90% C.I. for μ_D .

$$\bar{D} \pm t_{0.05} (n-1) S_D / \sqrt{n}$$

$$16.7 \pm 2.015 \cdot \frac{25.36}{\sqrt{6}}$$

$$16.7 \pm 2.015 \cdot 10.353 \text{ (standard error)}$$

$$16.7 \pm 20.821 \text{ or } -4.9 \leq \mu_D \leq 37.6^*$$

Of course, we see $\mu_D = 0$ a possible value in the interval. This is exactly why we could not reject $H_0: \mu_D = 0$ in the prior hypothesis test.

12p 10

Is this really better?

The claim is, this C.I. is tighter than if we assumed independence.

lets do $\mu_1 - \mu_2$ 90% C.I

Back to the data, we get

$$X_1: n_1=6, \bar{X}_1=209.8\bar{3}, S_1^2=725.7\bar{6}$$

$$X_2: n_2=6, \bar{X}_2=193.1\bar{6}, S_2^2=495.3\bar{6}$$

Using old d.f. in means methods,

note $df = n_1 + n_2 - 2 = 10$ is small.

we use

$$(\bar{X}_1 - \bar{X}_2) \pm t_{0.05}(10) S_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$

$$\text{where } S_p = \sqrt{\frac{(n_1-1)S_1^2 + (n_2-1)S_2^2}{n_1+n_2-2}}, \text{ pooled estimator.}$$

$$\text{Here } S_p^2 = \sqrt{\frac{S_1^2 + S_2^2}{2}} = \sqrt{\frac{610.56}{2}} = 24.71$$

$$\text{and standard error } S_p \sqrt{\frac{1}{6} + \frac{1}{6}} = 14.266$$

$$S_0 \quad 16.67 \pm 1.812(14.266)$$

$$\Rightarrow 16.67 \pm 25.85 \leftarrow \text{much bigger error term}$$

WORSE!

with basically the same center!