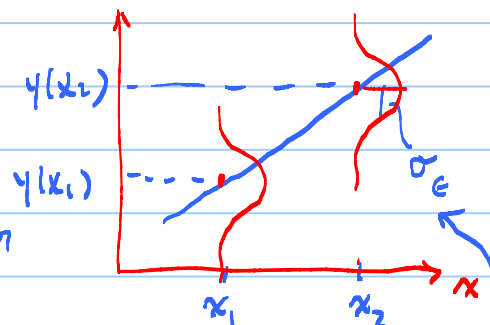<u>Lecture 25 (Ch. 11)</u>

Last time: In preparation for doing inference in regression, we introduced The <u>prob model</u> for regression:

At a given $x$, $y \sim N(\ \mu\ ,\ \sigma\ )$
$$\underset{y(x)}{=} \qquad \underset{\sigma_E}{=}$$

The $\mu$ param. (ie. center) of The Normal dist. of $y$'s is allowed to vary with $x$.



The $\sigma$ param. (denoted $\sigma_E$) is <u>not</u> a function of $x$, and is estimated/approximated with $S_e = \sqrt{\dfrac{SSE}{n-(k+1)}}$ $\quad k = \#$ of $\beta$'s.

Then,

1) $Y(x) = \alpha + \beta x + \cdots \ = $ true mean of $y$, at a given $x$.

2) $\hat{y}(x) = \hat{\alpha} + \hat{\beta} x + \cdots = $ estimated mean of $y$, given $x$

3) 95% of $y$'s, at given $x$ are within $y(x) \pm 1.96\ \sigma_E$

4) $\text{prob}(a < y < b) = pr\left( \dfrac{a - Y(x)}{\sigma_E} < \underset{z}{\underbrace{\dfrac{Y - Y(x)}{\sigma_E}}} < \dfrac{b - Y(x)}{\sigma_E} \right) = \cdots$

$\qquad$ (Table I)

5) more (below)

$\quad$ e.g. $\hat{\beta}$ (and $\hat{\alpha}$) is now a random variable!  $\quad$ like $\bar{x}$

$\qquad$ It has a distribution!  $\qquad\qquad\qquad$ like $\bar{x} \sim N(\mu, \sigma_{\bar{x}})$

$\qquad$ It has a prob!  $\qquad\qquad\qquad\qquad$ like $pr(\bar{x} > \bar{x}_{obs})$

$\qquad$ ⋮

$\quad$ We can build a CI for $\beta$ (and $\alpha$)  $\quad$ like CI for $\mu_x$

```
n = 10
n.trial = 64

x = c(1:n)
y_true = 10 + 2*x
sigma_eps = 15

par(mfrow=c(8,8),mar=c(0,0,0,0))
set.seed(123)
for(trial in 1:n.trial){
y_obs = y_true + rnorm(n,0,sigma_eps)
lm.1 = lm(y_obs ~ x)
plot(x, y_obs)
abline(10,2, col=2)
abline(lm.1, col=4)
}
```

⟵ Note That The x-values are The same across trials.
( in The kind of regression we are doing, x has no uncertainty;
only y does.)

$\sigma_{\hat{\alpha}}$

$\sigma_{\hat{\beta}}$

$\alpha$

$\hat{\alpha}$

$\beta$

$\hat{\beta}$

Let's build a CI (and hyp. test) for ONE $\beta$ : $Y_i = \alpha + \beta x_i + \epsilon_i$

**Theorem:** If $\epsilon \sim N(0, \sigma_\epsilon^2)$, Then $\hat{\beta}$ is normal with params:

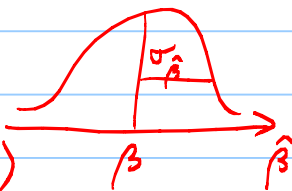Expected value (or mean) of The sampling dist. of $\hat{\beta}$

$$E[\hat{\beta}] \equiv \mu_{\hat{\beta}} = \beta \quad \leftarrow \text{pop. slope}$$

$$\sqrt{V[\hat{\beta}]} \equiv \sigma_{\hat{\beta}} = \frac{\sigma_\epsilon}{\sqrt{S_{xx}}} \quad (\text{not obvious})$$

$$\hookrightarrow S_{xx} = \sum_i^n (x_i - \bar{x})^2 = (n-1)S_x^2$$

Defn. of sample var.

**Ch. 7**

If $x \sim N(\mu_x, \sigma_x)$, Then $\bar{x}$ is Normal with params

$$E[\bar{x}] = \mu_{\bar{x}} = \mu_x$$

$$\sqrt{V[\bar{x}]} = \sigma_{\bar{x}} = \sigma_x/\sqrt{n}$$

---

Since $\hat{\beta} \sim N(\beta, \sigma_{\hat{\beta}})$, Then

$$z = \frac{\hat{\beta} - \beta}{\sigma_{\hat{\beta}}} = \frac{\hat{\beta} - \beta}{\sigma_\epsilon/\sqrt{S_{xx}}} \sim N(0,1)$$

$$\longrightarrow t = \frac{\hat{\beta} - \beta}{se/\sqrt{S_{xx}}} \sim t\text{-dist.}$$

$$df = n - ② \quad k+1$$

$$\bar{x} \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$

$$z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} \sim N(0,1)$$

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}} \sim t\text{-dist.}$$

$$df = n - 1$$

---

Then, from self-evident fact

$$pr(-t^* < t < t^*) = \text{Conf. level}$$

$$df = n - 2 \quad (\text{Table VI})$$

C.I. for $\beta$ : $\hat{\beta} \pm t^* \dfrac{se}{\sqrt{S_{xx}}}$

$H_0 : \beta = \beta_0$

$H_1 : \beta \;\square\; \beta_0$

$$t_{obs} = \frac{\hat{\beta}_{obs} - \beta_0}{se/\sqrt{S_{xx}}}$$

p-value = $(1,2) \cdot pr\left(\hat{\beta} \;\square\; \hat{\beta}_{obs}\right) =$

$\underset{1 \text{ or } 2 - \text{sided.}}{\Big\uparrow}$

$= (1,2) \; pr\left(t \;\square\; t_{obs}\right)$

$= \text{Table VI}, \; df = n-2$

$$pr(-t^* < t < t^*) = \text{Conf. level}$$

$$df = n - 1$$

C.I for $\mu$ : $\bar{x} \pm t^* \dfrac{s}{\sqrt{n}}$

$H_0 : \mu = \mu_0$

$H_1 : \mu \;\square\; \mu_0$

$$t_{obs} = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$$

p-value $= (1,2) \; pr\left(\bar{x} \;\square\; \bar{x}_{obs}\right)$

$= (1,2) \; pr\left(t \;\square\; t_{obs}\right)$

$= \text{Table VI}, \; df = n-1$

| problem 11.17 [Revised] |
| --- |

$n = 13$   $x =$ nickel content,   $y =$ percentage austentite.

Data:   $\sum(x_i - \bar{x})^2 = 1.183$ $\qquad = S_{xx}$

$\qquad\qquad \sum(y_i - \bar{y})^2 = 0.0508$ $\qquad = S_{yy} \quad = SST$

$\qquad\qquad \sum(x_i - \bar{x})(y_i - \bar{y}) = 0.2073$ $\quad = S_{xy}$

Question: Is There a statistically significant ($\alpha = 0.05$) relationship between $x$ and $y$? Hint: $SS_{exp} = \hat{\beta} S_{xy}$

1) C.I. $\beta$:   $\hat{\beta} \pm t^* \dfrac{S_e}{\sqrt{S_{xx}}}$

$\hat{\beta} = \dfrac{S_{xy}}{S_{xx}} = \dfrac{.2073}{1.183} = .1752 \longrightarrow$ SSE $= SST - SS_{exp}$

$\qquad\qquad\qquad\qquad\qquad\qquad = .0508 - (.1752)(.2073) = .014$

$S_e = \sqrt{\dfrac{SSE}{n-2}} = \sqrt{\dfrac{.014}{13-2}} = 0.0357$

$\therefore$ 95% CI for $\beta$:   $.1752 \pm 2.201 \left( \dfrac{.0357}{\sqrt{1.183}} \right) = 0.0328$ $\qquad = (0.10, 0.24)$

$\qquad\qquad\qquad\qquad$ $df = 13-2$

1) We are 95% Confident That The pop. $\beta$ is in here.

2) There is a 95% prob That a random CI will cover $\beta$.

3) Corollary: Relationship is statistically significant (zero not in CI).

2) $H_0 : \beta = 0$ $\qquad t_{obs} = \dfrac{.1752 - 0}{.0328} = 5.31$,

$H_1 : \beta \neq 0$

$\qquad$ p-value $= 2\,pr(\hat{\beta} > \hat{\beta}_{obs}) = 2\,pr(t > t_{obs})$

$\qquad\qquad\qquad\qquad = 2\,pr(t > 5.31) < 0.001$

p-value $< \alpha$ $\qquad\qquad\qquad\qquad\qquad$ $\uparrow$

$\therefore$ Evidence That $\beta \neq 0$. (same conclusion $\qquad$ Table VI

$\qquad\qquad\qquad\qquad$ as above). $\qquad\qquad$ $df = 13-2$

In summary: We have 2 ways of testing whether there is a relationship between 2 continuous variables.

Note that the test of $\beta = 0$ is equivalent to testing if there is a linear relationship between $x$ and $y$. But if a linear relationship is all that you are testing, then we can test the population correlation coeff

$$H_0: \rho = 0$$

$$H_1: \rho \neq 0$$

The test statistic for this test is a bit weird:

$$\Rightarrow t = \frac{r - 0}{\sqrt{\frac{1-r^2}{n-2}}} \quad \text{has a } t \text{ distr, with } df = n-2.$$

Recall $r = S_{xy}/\sqrt{S_{xx}S_{yy}}$

This way, you take your data $(x_i, y_i)$, compute the sample correl. coeff $(r)$, then $t_{obs}$, and then p-value, all without any fitting.

3) For the above example:

$$H_0: \rho = 0$$
$$H_1: \rho \neq 0$$

$$r = \frac{S_{xy}}{\sqrt{S_{xx}S_{yy}}} = \cdots = .8456$$

$$t_{obs} = \frac{r-0}{\sqrt{\frac{1-r^2}{n-2}}} = \cdots = 5.3 \quad \leftarrow \text{Same value as } t_{obs} \text{ we got above when testing } \beta.$$

p-value $= 2 \text{ prob}(t > t_{obs}) = $ Same as above.

$\therefore$ Same conclusion.

book problem

We have now done inference on $\beta$ (and $\alpha$),

What about <u>multiple</u> regression (ie. multiple $x$'s and $\beta$'s)?

In going from $y = \alpha + \beta x$  <span style="color:red">(1+1 parm)</span>  <span style="color:red">#of $\beta$'s.</span>

to $Y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_k x_k$  <span style="color:red">(k+1 params)</span>

things generalize in a straight forward way.

Basically, all that happens is <span style="color:red">$df = n-2 \longrightarrow df = n - (k+1)$</span>

This happens every where, e.g.

1) The estimate of $\sigma_\epsilon^2$ is $s_e^2 = \dfrac{SSE}{n - (k+1)}$

2) The df associated with t-test <span style="color:red">changes: $n-2 \longrightarrow n-(k+1)$</span>

Finally, don't forget that the issues of collinearity, interaction, non-linearity, ..., and overfitting all return when doing multiple reg.

But, the presence of <u>multiple</u> $\beta$'s allows for 1.5 more tests:

<span style="color:red">0.5)</span> $\begin{cases} H_0: & \beta_i : \square \beta_0 \\ H_1: & \beta_i : \square \beta_0 \end{cases}$  <span style="color:red">Is the $i^{th}$ predictor useful?</span>

<span style="color:blue">I'm saying 1.5 tests because this one is a straight forward generalization of the t-test for a single $\beta$ (see next page).</span>

<span style="color:red">1)</span> $\begin{cases} H_0: & \beta_1 = \beta_2 = \cdots = \beta_k = 0 \\ H_1: & \text{At least 1 } \beta_i \text{ is} \neq 0 \end{cases}$  <span style="color:red">Are any of the predictors useful? (Test of "model utility".)</span>

<span style="color:blue">In $Y = \alpha + \beta_1 x_1 + \beta_2 x_2 + \cdots \beta_k x_k$, if all $\beta_i = 0$, then <u>none</u> of the predictors are useful for predicting $y$.</span>

$\boxed{\text{Good News:}}$ The test for each $\beta_i$ is the same as the
t-test for a single $\beta$, except for $\boxed{df = n-(k+1)}$

E.g. suppose we want to test $\beta_3$ :

$H_0: \beta_3 \,\square\, \beta_0 \quad (e.g.\ 0)$
$H_1: \beta_3 \,\square\, \beta_0$

$t_{obs} = \dfrac{\hat{\beta}_3 - \beta_0}{S_e/\sqrt{S_{x_3 x_3}}}$

p-value $= (1,2)\ pr(t\,\square\, t_{obs})$

$\phantom{p-value} = \cdots$  $df = n-(k+1)$

$\phantom{p-value == } \text{Table } VI$

C.I. for $\beta_3$ :

$\hat{\beta}_3 \pm t^* \dfrac{\overbrace{S_e}}{\sqrt{S_{x_3 x_3}}} \qquad \overline{\sqrt{\dfrac{SSE}{n-(k+1)}}}$

$df = n-(k+1)$

Technically, in multiple regression $S_{\hat{\beta}}$
is NOT $S_e/\sqrt{S_{xx}}$. The denominator ends-up
being a more complicated function of $x$'s.
But when the predictors are completely
uncorrelated, then this formula is OK.

$\boxed{\text{Note}}$: even though we are testing $\underline{\underline{ONE}}$ $\beta_i$, the df is $n-(k+1)$

$\boxed{\text{Bad News:}}$ If you test each of the $\beta_i$ separately, it's almost
guaranteed that some of the $\beta$'s will pass the test, ie. give small p-value,
ie. are found to be useful, when in fact. They are $\underline{not}$.

Here is the proof, but it's only $\boxed{FYI}$

$\boxed{\text{Bad News}}$ : If you test each of the $\beta_i$ separately, you will
make many more Type I errors than $\alpha\%$ of the time!

Consider 3 $\beta$'s : $\beta_1, \beta_2, \beta_3$

Type I errors : $\underbrace{(\beta_1 \neq 0 \mid \beta_1 = 0)}_{e_1} \quad \underbrace{(\beta_2 \neq 0 \mid \beta_2 = 0)}_{e_2} \quad \underbrace{(\beta_3 \neq 0 \mid \beta_3 = 0)}_{e_3}$

You may commit the errors $e_1$ OR $e_2$ OR $e_3$
OR $(e_1$ and $e_2)$ OR $(e_1$ and $e_3)$ OR $(e_2$ and $e_3)$ OR $(e_1$ and $e_2$ and $e_3)$.

It can be shown that the prob of making at least 1 Type I error
approaches 1 as the number of tests increases.

**Good News** : Enter the test of model utility!

Thm. $F = \dfrac{R^2/(k)}{(1-R^2)/(n-(k+1))}$ "numerator df" $\sim$ F-distribution with $df = (k,\ n-(k+1))$ "denominator df"

$\therefore$ p-value $= pr(F \geq F_{obs})$ — Just like in 1-way ANOVA where $H_1$: Atleast ...

Then, if p-value $< \alpha$, we can reject $H_0$ $(\beta_1 = \beta_2 = \cdots = \beta_k = 0)$ in favor of $H_1$ (atleast 1 $\beta_i$ is not zero).

This F-test allows you to do ONE test to find out if any of the predictors are useful for predicting y. This is very useful if k is large, because it tells you if __any__ of the predictors are useful. I.e. it tells you if there is a "needle in the haystack," to begin with!

[IF] you get a significant result (ie. p-value $< \alpha$) from the test of model utility, then there is evidence that at least one of the predictors is useful. [THEN] you can do separate tests on each of the $\beta$'s to see __which__ predictors are useful. (see next page).

But [IF] the F-test comes back as non-significant, then there is no evidence that any of the predictors are useful. [THEN], you don't have to test each predictor, separately. This will not only save time, but more importantly, it will save you from the danger of making multiple Type I errors (ie. declaring some predictor as useful, when infact, it is not).

Recall That

- "bad" things happen if you keep adding terms to a regression model. Specifically, overfitting happens.
- overfitting is <u>not</u> a black and white thing — it happens gradually, and in degrees, as you add more terms.
- even a complete "garbage" term can lead to overfitting.

what happens to F (and its p-value)?

more terms $\longrightarrow$ higher $R^2$ $\longrightarrow$ higher F $\longrightarrow$ lower p-value.

I.e. If you keep throwing enough predictors into a model (regression or otherwise), the F-test of model utility will find at least 1 useful predictor, regardless of whether or not the predictors are actually useful.

So, you MUST <u>be Thoughtful</u> about adding terms to regression
$\hookrightarrow$ ANY model!

The k-dependence of the formula for F does complicate things a bit but you can ignore it, because the real problem arises from $R^2$ approaching 1, as the # of predictors increases. Still, we <u>can</u> pay attention to the k-dependence:

$$F = \frac{R^2/k}{(1-R^2)/[n-(k+1)]} = \frac{R^2}{1-R^2}\left(\frac{n-(k+1)}{k}\right) = \frac{R^2}{1-R^2}\left(\frac{n-1}{k}-1\right)$$

Now, technically k must be less than $(n-1)$, otherwise $F<0$, which it cannot be. So, $k < n-1$, in which case the largest allowed value of k is $n-2$, and so $\frac{n-1}{k}$ is at most $\frac{n-1}{n-2}$, i.e. a constant! Then, we're back to looking at how $R^2$ grows.

FYI

(11.66)

```
The regression equation is
durpr = -0.912 + 0.161 formconc + 0.220 catratio + 0.0112 temp + 0.102 time

Predictor       Coef      StDev        T        p
Constant      -0.9122    0.8755     -1.04    0.307
formconc       0.16073   0.06617     2.43    0.023
catratio       0.21978   0.03406     6.45    0.000
temp           0.011226  0.004973    2.26    0.033
time           0.10197   0.05874     1.74    0.095
S = 0.8365 R-Sq = 69.2% R-Sq(adj) = 64.3%

Analysis of Variance
Source        DF       SS      MS        F       P
Regression     4   39.3769  9.8442   14.07   0.000
Error         25   17.4951  0.6998
Total         29   56.8720
```

You have already learned what all These numbers are, from The prelab. But now, we are going to do everything by hand.

Also $F = \dfrac{MS_{expl}}{MS_{Err}} = \dfrac{9.8442}{0.6998}$

$n - (k+1)$

a) Is The model useful?   $= n-1$

F-test:    $F_{obs} = \dfrac{R^2/k}{(1-R^2)/(n-(k+1))} = \dfrac{\left(\dfrac{.692}{4}\right)}{\left(\dfrac{1-.692}{30-(4+1)}\right)} = 14.04$

p-value $= prob(F > F_{obs}) = prob(F > 14.04) < .001$

According To Table VIII, $df = (4, 25)$

∴ At any reasonable $\alpha$, we can reject $H_0$ (That all $\beta_i = 0$) in favor of $H_1$ (That at least 1 of The $\beta_i \neq 0$). I.e. The model is useful.

b) Estimate, in a way that conveys info about precision & reliability, the average change in durability press rating associated with a 1-degree increase in curing temperature, when all other predictors remain fixed. (if there is NO collinearity)

I.e. what's the C.I for $\beta_{temp}$. $t^*$ at $df = n - (k+c) = 25$

95% CI: $\hat{\beta} \pm t^* \dfrac{Se}{\sqrt{S_{xx}}} = .0112 \pm 2.060 \dfrac{0.8365}{\sqrt{?}}$ ← not given!

2-sided ↗   ↳ temp.

$\underbrace{\phantom{xxxxxxxxx}}_{\text{std. err in } \hat{\beta}}$ } from printout.

∴ $.0112 \pm 2.060 \,(.004973) \Rightarrow (.001, .021)$

This is the interval estimate of $\beta_{temp}$. It's useful as it is, but we we can also see that $\beta_{temp} \neq 0$

We can build the CI for all the other $\beta_i$ :

| C.I. for | | | | |
|---|---|---|---|---|
| conc | .1607 | $\pm$ 2.060 | (.06617) = | (0.02, 0.30) |
| cat ratio | .2198 | $\pm$ " | .03406 = | (0.15, 0.29) |
| temp | .0112 | $\pm$ " | .00497 = | (0.001, 0.02) |
| time | .10197 | $\pm$ " | .0587 = | (-0.02, 0.22) ← |

Note that 3 $\beta$'s are non-zero.
   ↳ At least 1  ✓

Given that there is no evidence that "time" is a useful predictor, you may remove it from the regression model so that the "smaller" model will be less likely to overfit data.

In part a, we found out that at least one of the $\beta_i \neq 0$.
To see which one(s), we test each of them!

$$H_0: \beta_i = 0 \qquad vs. \qquad H_1: \beta_i \neq 0 \qquad \text{for each } i.$$

c) E.g.
$$H_0: \beta_{formald.} = 0$$
$$H_1: \beta_{formald.} \neq 0$$

$\frac{Se}{\sqrt{S_{xx}}}$
↑↑
formalda.
(from Table)

$$t_{obs} = \frac{.16073 - 0}{.06617} = 2.43 \qquad (\text{check the output!})$$

even though testing 1 $\beta$.

$$p\text{-value} = 2 \text{ prob}(t > t_{obs}) = 2(.012) \qquad df = n - (k+1) = 25$$

$$= .024 \qquad (\text{check output!})$$

So, $p$-value $< \alpha \implies$ formaldehyde provides useful info.

In fact, look at all the $p$-values:

look at last col. of printout.

$$.023, \quad .000, \quad .033, \quad .095 \leftarrow$$

At $\alpha = .05$ 
$\beta \neq 0$    $\beta \neq 0$    $\beta \neq 0$    $\beta =$ cannot tell
formald.    cat.      temp      time

consistent with the conclusions in part b.

(FYI)
Note these $p$-values are <u>different</u> from what you would get if you did $y = \alpha + \beta_1 x_1$, $y = \alpha + \beta_2 x_2$, .... etc. and tested if each of these $\beta_i$ are zero. The multiple regression model is more correct because it does take into account the correlations between predictors. See ch. 3 lects.

We have learned that if p-value < alpha, then there's evidence to reject H0 in favor of H1. For the test of model utility, p-value = pr(F > F_obs). So, for that p-value to be less than alpha, F_obs must be larger than some critical value.

a) At alpha=0.05, find the critical value of F_obs for a multiple regression problem involving four betas, and 30 cases.

b) Find the critical value of R^2 (above which p-value < alpha). Hint: The F-ratio appearing in the test of model utility depends on R^2 of the model. So, if you know the critical value of F (as in part a), then you know the critical value of R^2.

Moral: Like all other tests we have studied, the reject/no-reject decision can be based on the critical value of some statistic, i.e. without a p-value. For the test of model utility, the decision can be made by comparing F_obs with some critical value (e.g. found in part a), or even by comparing R^2_obs with its critical value (e.g. found in part b).

We have seen that adding useless predictors to a regression model will increase R2. Here, let's examine what our inference methods say if the predictors are, in fact, useless. Suppose the true/pop fit is y = 1,(i.e., no x at all), and so a possible sample from the population could be the following:

```
set.seed(123)    # Use this line to make sure we all get the same answes.
n = 20
y = 1 + rnorm(n,0,1)
```

a) Write code to make data on 10 useless predictors (and no useful predictors) each from unif(-1,+1), fit the model y = alpha + beta1 x1 + ... + beta10 x10, perform the test of model utility, and perform t-tests on each of the 10 coefficients to see if they are zero. Show/turn-in your R code.

b) According to the F-test of model utility, are any of the predictors useful at alpha = 0.1?

c) According to the t-tests, are any of the predictors useful at alpha = 0.1? See the solns to make sure you understand the moral of this exercise.

Consider a multiple regression problem with k betas on the right-hand side. Suppose all of the k predictors are completely useless. But, of course, we don't know that, so we test each of the betas individually. Our hyp. testing formalism assures that each test has prob. alpha of finding the predictor useful (when in fact it's useless).
a) what's the prob. of finding j useful predictors out of k predictors? Hint: Here you should recognize a familiar string of words here!
b) What's the prob. that at least 1 of the k predictors will be found to be useful (when it's not)?
c) Plot that prob. vs. k = 1:100, for alpha=0.05, and for alpha=0.01
(Make sure you check the soln, later)