

F-tests / ANOVA in MLR Recap

1, 2, or 3

We can use `sm.stats.anova_lm(model, typ = 1)` to generate ANOVA tables in MLR.

The `typ` argument is super important — it determines if we get sequential or partial F-test results.

Typ = 1 \Rightarrow Uses sequential sum of squares

reduction in SSE from H_0 to H_1

source	SS	df	F	p	H_0	H_1
X_1	<u>$SS(X_1)$</u>	1	$\frac{SSX_1/df}{SSE_F/n-p}$		$n-1$ $Y \sim 1$	$n-2$ $Y \sim X_1$
X_2	$SS(X_2 X_1)$	1	$\frac{SS/1}{SSE_F/n-p}$	<u>1.02</u>	$n-2$ $Y \sim X_1$	$n-3$ <u>$Y \sim X_1 + X_2$</u>
X_3	$SS(X_3 X_1, X_2)$	1			$Y \sim X_1 + X_2$	$Y \sim X_1 + X_2 + X_3$
resid	<u>SSE_F</u>	<u>$n-p$</u>				

⊛ order matters for `typ=1`

For ex, if X_2 's row has $p \approx .02 \Rightarrow$

X_2 is a sig. predictor GIVEN that

X_1 is included in the model.

④

$$\frac{(SSE_{H_0} - SSE_{H_1}) / (df_{H_0} - df_{H_1})}{SSE_F / (df_F)}$$

For typ=2 \Rightarrow uses partial sum of squares

			reduction in SSE from H_0 to H_1				
source	SS	df	F	p		H_0	H_1
X_1	$SS(X_1 X_2, X_3)$	1	$\frac{SS}{SSE_F/n-p}$	1.4		$Y \sim X_2 + X_3$	$Y \sim X_1 + X_2 + X_3$
X_2	$SS(X_2 X_1, X_3)$	1		1.5		$Y \sim X_1 + X_3$	$Y \sim X_1 + X_2 + X_3$
X_3	$SS(X_3 X_1, X_2)$	1		.0001		$Y \sim X_1 + X_2$	$Y \sim X_1 + X_2 + X_3$
resid	SSE_F	$n-p$					

⊛ order doesn't matter for typ=2

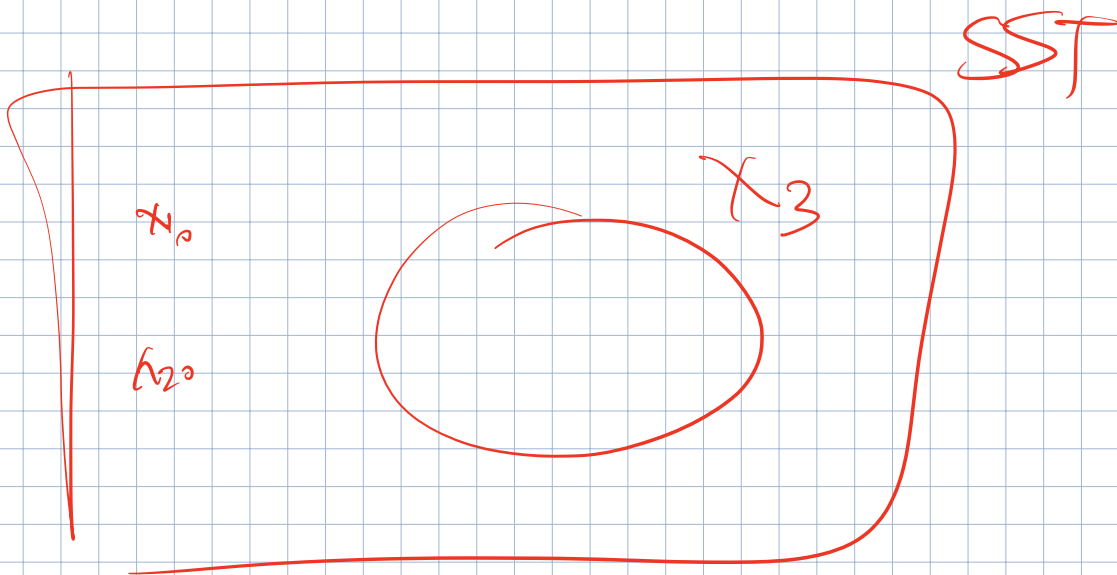
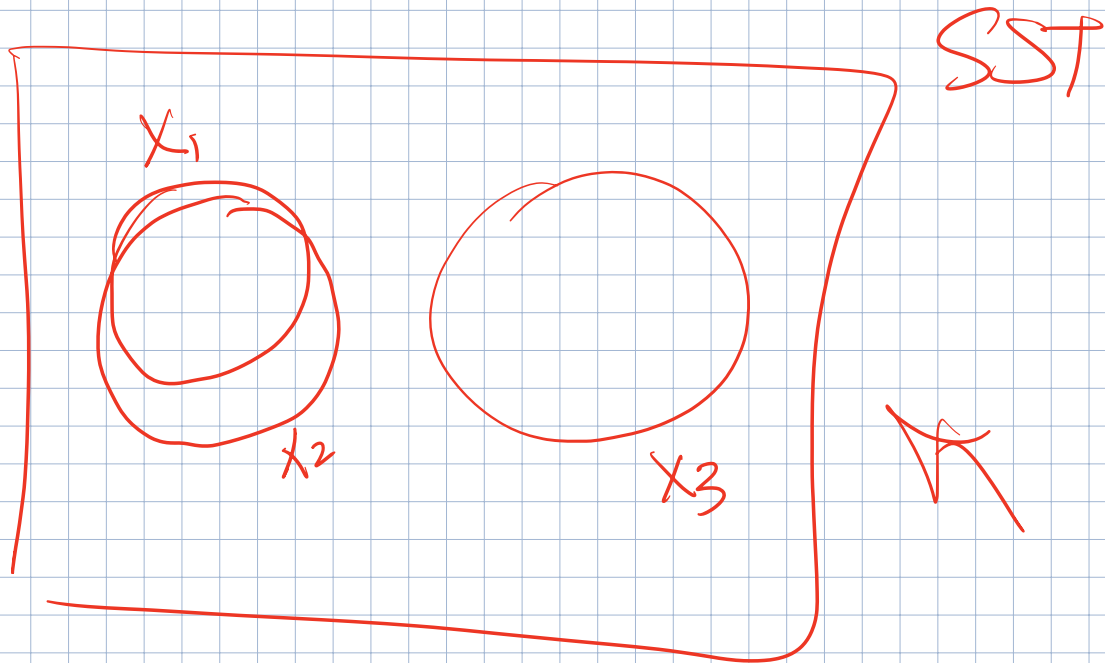
For ex: if $p < .0001$ for X_3 's row:

We have evidence that X_3 is a sig. predictor

GIVEN all other predictors are in the model.

Note: ① $\text{typ}=2 \Rightarrow$ order doesn't matter

② $\text{typ}=2 \Rightarrow F = t^2$



Dealing w/ Categorical Predictors

What happens if I have categorical preds:

UX - GENDER: M, F, NB,

- RACE: ...

- EYE COLOR: ...

The problem is in their current form,
this data isn't compatible w/ the

OLS solution $\hat{\beta} = (X^T X)^{-1} X^T y$

X need numbers!

Indicators / Dummy Vars

The most common way to "recode" this data is to define the indicators:

Ex: Eye color has 4 levels:

- blue
- green
- hazel
- brown

$$X_i = \begin{cases} X_{i-\text{blue}} \\ X_{i-\text{green}} \\ X_{i-\text{hazel}} \end{cases}$$

$$X_{i-\text{blue}} = \begin{cases} 1 & \text{if } X_i = \text{"blue"} \\ 0 & \text{if } X_i \neq \text{"blue"} \end{cases}$$

$$X_{i-\text{green}} = \begin{cases} 1 & \text{if } X_i = \text{"green"} \\ 0 & \text{if } X_i \neq \text{"green"} \end{cases}$$

$$X_{i_hazel} = \begin{cases} 1 & \text{if } X_i = \text{'hazel'} \\ 0 & \text{if } X_i \neq \text{'hazel'} \end{cases}$$

X_i	X_{i_blue}	X_{i_green}	X_{i_hazel}
1 green	0	1	0
2 brown	0	0	0
3 brown	0	0	0
4 brown	0	0	0
5 purple (11)	0	0	0
6 blue	1	0	0

we call this indicator variable:

$$X_{i_green} = \mathbb{1}(X_i = \text{green}) = \begin{cases} 1 & \text{if } X_i = \text{green} \\ 0 & \text{if } X_i \neq \text{green} \end{cases}$$

Ex: Predict $Y = \text{dexterity}$

$X_1 = \text{R/L-handed}$

$X_2 = \text{brown/not}$

$$Y = \beta_0 + \beta_1 \mathbb{I}(\text{LH}) + \beta_2 \mathbb{I}(\text{not brown}) + \epsilon_i$$

↓ fit model

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 \mathbb{I}(\text{LH}) + \hat{\beta}_2 \mathbb{I}(\text{not brown})$$

int	6
LH	1
Not brown	0

Interpreting the coefs:

Baseline:

We predict an avg dexterity score of

$\hat{\beta}_0 = 6$ for right handed brown-eyed people.

We predict an avg. dexterity score of

$\hat{\beta}_0 + \hat{\beta}_1 = 6 + 1 = 7$ for left handed brown-eyed
indivs.

Example Design Matrix

$$\begin{pmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 0 & 1 \end{pmatrix}$$