

## L6. Deal with Categorical Predictors.

### — Dummy variables in MLR

So far we have only dealt with Numeric variables: just

through in the model directly, then  $\vec{b} = \underbrace{(\vec{X}^T \vec{X})^{-1} \vec{X}^T \vec{y}}$

can be plug in directly  
with the predictor values.

What if we have categorical predictors?

GENDER : male female

RACE : white, asian, black, ...

COLOR : red, green, blue.

We are collecting this kind of information from observations.

Where the values can be  $K$  different levels, which are

not numeric. We need to change it to be "NUMBERS"

so we can perform the analysis we have used.

• Dummy Variables. — MOST COMMON.

The most common way to "recode" categorical variable is to define dummy variables.:

$$X_i = \begin{cases} L_1 & \text{when obs } i \text{ is in category } L_1 \\ L_2 & \text{when obs } i \text{ is in category } L_2 \\ L_3 & \\ \vdots & \\ L_K & \end{cases}$$

$$X_{i-L_1} = \begin{cases} 1 & \text{if } i^{\text{th}} \text{ obs is in } L_1 \\ 0 & \text{o.w. } (L_2, L_3, \dots, L_K) \end{cases}$$

$\Leftrightarrow$

$$X_{i-L_2} = \begin{cases} 1 & \text{if } i^{\text{th}} \text{ obs is in } L_2 \\ 0 & \text{o.w. } (L_1, L_3, \dots, L_K) \end{cases}$$

$\vdots$

$$X_{i-L_{K-1}} = \begin{cases} 1 & \text{if } i^{\text{th}} \text{ obs is in } L_{K-1} \\ 0 & \text{o.w. } (L_1, \dots, L_{K-2}, L_K) \end{cases}$$

$X_{i-L_K}$  doesn't need to define: if  $X_{i-L_1} = \dots = X_{i-L_{K-1}} = 0$   
then  $i^{\text{th}}$  obs is in  $L_K$ .

For ex : eye color of 5 obs.

Obs	Color	Color-Brown	Color-Blue
1	Brown	1	0
2	Blue	0	1
3	Blue	0	1
4	Green	0	0
5	Brown	1	0

Dummy coding →

Green.

One categorical variable with  $k$  levels will be recoded as  $(k-1)$  dummy variables. then the regression fitted is

$$y_i \sim \underbrace{X_{i1} \beta_1 + X_{i2} \beta_2 + \dots + X_{i,k-1} \beta_{k-1}}_{\text{original } X_1} + X_{i2} + X_{i3} + \dots + X_{ip-1}$$

1. How to read the estimated coefficient?

Again, use the eye-color example. the fitted line would

look like :

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 * \text{Color-brown} + \hat{\beta}_2 * \text{Color-blue}$$

$\hat{\beta}_1$  : when Color-Brown = 1.

$$\textcircled{1} \hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 * \text{Color-Brown.} \quad (\text{Color-blue} = 0 \\ \text{when Color-Brown} = 1)$$

compared to

$$\textcircled{2} \hat{y}_i = \hat{\beta}_0 \quad (\text{Color-blue} = 0 \text{ and} \\ \text{Color-Brown} = 0, \\ \text{AKA. Color = Blue})$$

$\textcircled{1} - \textcircled{2} = \hat{\beta}_1$  : measuring the change in  $y$  when color changes from Blue (Baseline) to Brown (Color-Brown)

Summary: the coef. of  $X_{Lj}$  measures the change caused in  $y$  when the obs is in the level  $L_j$  compared to the baseline level.

Python: coef for fuel-type [gas] = -5.6297.

means the cars with gas fuel-type has, on average, 5.6297 less city mpg than diesel type.

## 2. How to read t-test?

t test on  $H_0: \beta_{L_j} = 0$  v.s.  $H_1: \beta_{L_j} \neq 0$

Rejection indicates there is a significant difference in  $Y|L_j$  compared to  $Y|L_{\text{baseline}}$ .

As long as one test for the ~~a~~ dummy variables is significant, we can conclude the original categorical variable is significant to  $Y$ .

Ex: Python: `citympg ~ drive-wheels`

	p-value of t test
<code>drivewheels_fwd :</code>	0.006
<code>drivewheels_rwd :</code>	0.185

Indicates there's significant (positive) change in city-mpg from baseline (4wd) to fwd.

No sig. change from baseline (4wd) to rwd.

In general, drivewheels has significant impact on citympg.



3. How to read ANOVA? Both typ=1 and 2 shows the original categorical variable, instead of dummies.  
 type=1.

	df	SS
Drivewheels	2	<del>SSD</del>
fuel-type	1	
enginesize	1	
Residuals	<u>n-5.</u>	

Reduced:  $y_i = \beta_0 + \epsilon_i$

Full:  $y_i = \beta_0 + \beta_1 * \text{Drivewheels\_fuel} + \beta_2 * \text{Drivewheels\_rwd} + \epsilon_i$

In Full model:  $\beta_0$ , two

effects with drivewheels, 1 coefficient with fuel-type

1 coeff with enginesize

Total of 5 parameters estimated

type=2.

	df	SS
Drivewheels	2	
fuel-type	1	
enginesize	1	
Residuals	n-5	

Reduced:  $y_i = \beta_0 + \beta_3 * \text{fuel-type\_gas} + \beta_4 * \text{enginesize} + \epsilon_i$

Full:  $y_i = \beta_0 + \beta_1 * \text{Drivewheel\_fuel} + \beta_2 * \text{Drivewheel\_rwd} + \beta_3 * \text{fuel-type\_gas} + \beta_4 * \text{enginesize} + \epsilon_i$

Note: There are cases that you would want to treat the predictor as categorical, but the data was either recorded with numbers:

①

color	
1	where 1 = "red"
2	2 = "green"
2	3 = "blue"
3	..
3	
4	
:	

Then we need to do `smf.ols('y ~ C(x)')`  
change x into categorical

otherwise it would be analysed as numeric variable

② Or for an ordinal variable, change it into dummies  
can give more information.

Age-group

1	<18
2	18 ≤ <30
3	30 ≤ <40
4	40 ≤ <50
5	≥50

the ~~the~~ levels comes with  
numerical order. change it into  
dummies can see changes between  
levels.

Note : Potential problem with dummies.

The information was collected in one categorical variable, but in the model we need to estimate  $(k-1)$  parameters. AKA, we are adding  $(k-1)$  predictors in the model, that could cause two problems :

- ① Multicollinearity ( Later )
- ② Not enough data to support a precise estimation: we all know that more data can have a more accurate estimate.

In Regression, the  $\frac{\text{Number of parameters}}{\text{Number of observation}} \geq 1$

(More details can be decided in "power analysis, not covered!")