

# CATEGORICAL VARIABLES AND INTERACTION TERMS

Categorical Variables → Binary Category

Example, Pink Slip = 1 if car has roadworthy  
= 0 otherwise.

$$Price_i = \beta_0 + \beta_1 (\text{Pink Slip}) + \epsilon_i$$

output →

DV: Price	Coef	SE	t	P-value
Intercept	3978.3	1056.2	3.7	0.0002
Pink slip	1625.6	1203.7	1.3	0.179

Inference → A car with a pink slip would command a sale price \$1,626 more than a car without pink slip, on average.

But p value > 0.1 which means p value / pink slip is not that important.

So consider other values,

$$\ln(\text{Price}) = \beta_0 + \beta_1 (\text{Age}_i) + \beta_2 (\text{Age}_i)^2 + \beta_3 \ln(\text{Odometer}_i) + \beta_4 (\text{Pink Slip}) + \epsilon_i$$

Output →

Ln(Price)	Coef	SE	t	P-value
Intercept	9.237	0.27	33.51	0.000
Age	-0.052	0.019	-3.78	0.003
Age <sup>2</sup>	0.001	0.000	4.72	0.000
Ln(Odometer)	-0.198	0.061	-3.24	0.0016
Pink slip	0.158	0.178	0.87	0.3846

$R^2 = 0.367$

Still pink slip have high P-value

$$\ln(\text{Price}) = 9.237 - 0.052(\text{Age}) + 0.001(\text{Age})^2 - 0.198(\text{Odometer}) + 0.158(\text{Pink Slip})$$

Inference → A car with pink slip would command a sale price 15.8% higher than a car without a pink slip on average holding all other variables constant. Because we log have log in dependent variable, so can express variable in percentage. by car.

## Multi level Category

Age Category = 1 if age ≤ 5, = 2 if age > 5 & age ≤ 15, 3 if 15 < age ≤ 35  
= 4 if age > 35, to get this cut check Age vs Price (See below)

$$\ln(\text{Price}) = \beta_0 + \beta_1 (\text{Age Cat}) + \beta_2 \ln(\text{Odometer}) + \beta_3 (\text{Pink Slip}) + \epsilon_i$$

But the problem is it is ordinal, 4 > 3 > 2 > 1 but does every time it will not mean older the car older the price or older the car higher the price, it cannot be normal

So use dummy variable, 8

AgeCat1 = 1 if age ≤ 5  
= 0 otherwise

Avoid AgeCat4, Dummy Variable trap

AgeCat2 = 1 if 5 < Age ≤ 15  
= 0 if otherwise

AgeCat3 = 1 if 15 < age ≤ 35  
= 0 otherwise

$$\ln(\text{Price}) = \beta_0 + \beta_1 (\text{Age Cat2}) + \beta_2 (\text{Age Cat3}) + \beta_3 (\text{Age Cat4}) + \beta_4 \ln(\text{Odometer}) + \beta_5 (\text{Pink Slip}) + \epsilon_i$$

Inference → On average, holding all other variable constant a car in age category 2 will command a price 12.9% lower than age category 1.

Ln(Price)	Coef	SE	t	P-value
Intercept	8.94	0.27	32.09	0.00
AgeCat2	-0.129	0.01	-0.52	0.60
AgeCat3	-0.735	0.26	-2.80	0.006
AgeCat4	0.479	0.32	1.45	0.150
Ln(Odometer)	-0.225	0.06	-3.69	0.004
Pink Slip	0.359	0.17	1.98	0.052

Age category plot Price

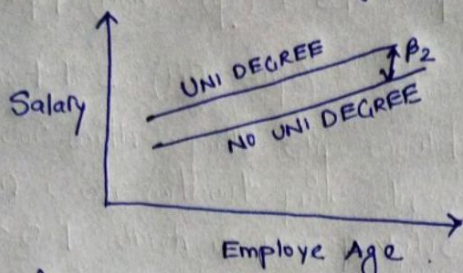


# INTERACTION TERMS

Example - Build a model to explain the salary of all Google's employees.  
 Dependent Variable - Salary Independent Var - 1) Employee Age  
 2) University degree - 1 Yes 0 No.

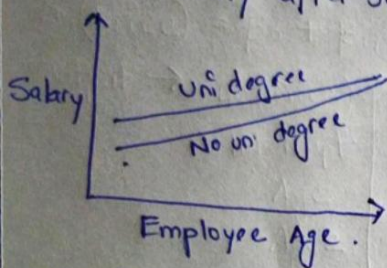
$$\text{Salary}_i = \beta_0 + \beta_1 (\text{Employee Age}) + \beta_2 (\text{Uni degree}) + \epsilon$$

If we plot,



$\beta_2$  will define how much a salary difference is there if we have a uni degree with no uni degree.

But, in reality after some experience, university degree don't matter, it's about experience.



So get the experience term, we need another variable like, Employee Age  $\times$  Uni degree.

$$\text{Salary} = \beta_0 + \beta_1 (\text{Employee Age}) + \beta_2 (\text{Uni degree}) + \beta_3 (\text{Employee Age}) \times (\text{Uni degree}) + \epsilon_i$$

The added term is called Interaction term.

Required when,  $X_1$  affects the relationship between  $X_2$  and  $Y$ .

Example, Age of employee affects the relationship between Uni degree & Salary.

Common misconception  $\rightarrow$  An interaction term is required when  $X_1$  &  $X_2$  are correlated.

Model output  $\rightarrow$

Ln (Price)	Coeff	SE	t	P-value
Intercept	9.125	0.271	33.28	0.00
AgeCat2	-0.181	0.23	-0.76	0.44
AgeCat3	-0.800	0.25	-3.18	0.00
AgeCat4	-0.390	0.42	-0.92	0.35
Ln (Odometer)	-0.290	0.059	-3.53	0.00
Pink slip	0.123	0.18	0.68	0.50
Pink slip $\times$ AgeCat4	1.371	0.45	3.02	0.003

Why we took Age Cat 4?

$\rightarrow$  Here come business/domain knowledge

In this case, we thought vintage cars might cost more.

But if vintage cars and also road worthy then it can add more. (we can drive that car now also).

$\rightarrow$  P value is also showing it is important variable.

Interaction effect  $\rightarrow$  Happens when one independent var interact with other var affect dependent variable.

Example  $\rightarrow$  Drinking diet juice

	Yes	No
Eating Pill	5 lbs 2 lbs	2 lbs 1 lb

Weight loss after 1 week

Highest weight loss is happening when we are taking both pill and drink together.

$\rightarrow$  So we can find this through factor analysis.

$\rightarrow$  Factor analysis extracts maximum common variance from all variable and put them into common score.

$\rightarrow$  Methods which can be used as Factor analysis are -

- 1) Principal component analysis (PCA)
- 2) Maximum likelihood method.
- 3) Common factor analysis.