

# Regression With Dummies

# Binary Variables

When dealing with a binary outcome like “Yes/No” or “Male/Female”, we often convert them into numeric binary variables using a base level. These are sometimes called dummy variables. . For a “Yes/No” scenario, the typical encoding is:

$$x_{ij} = \begin{cases} 1 & \text{if Yes} \\ 0 & \text{if No} \end{cases}$$

# Interpreting Binary Variables in Linear Regression

In linear regression, the coefficient for a binary variable shows the expected change in the dependent variable for the encoded value compared to its base level, all else held equal.

- The **base level** is the category against which all other categories are compared. It serves as a reference point. Typically, it's encoded as 0 in a binary variable setting.
- E.g., Comparing Group A vs. Group B with Group B as the base (encoded as 0), the coefficient for Group A (encoded as 1) suggests the expected difference in the dependent variable for Group A vs. Group B.

# Setting up the Regression Model

Consider a simple linear regression model where we predict a dependent variable  $Y$  using a continuous variable  $X_1$  (e.g., age) and a binary variable  $X_2$  (e.g., gender, encoded as 0 for Male and 1 for Female).

Our regression equation is:

$$E(Y) = \beta_0 + \beta_1 X_1 + \beta_2 X_2$$

Where:

- $\beta_0$  is the intercept.
- $\beta_1$  is the coefficient for the continuous variable  $X_1$ .
- $\beta_2$  is the coefficient for the binary variable  $X_2$ .

# Interpreting the Coefficients

- $\beta_1$  shows the expected change in  $Y$  for a one-unit increase in  $X_1$ , keeping  $X_2$  constant.
- $\beta_2$  shows the expected change in  $Y$  when  $X_2$  is 1 (Female) compared to 0 (Male), keeping  $X_1$  constant.

For instance, if  $\beta_2$  is positive, it suggests that, on average, being Female is associated with a higher value of  $Y$  than being Male, holding  $X_1$  constant.

# Impact on Regression Estimates

Given:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2$$

For a Male (with  $X_2 = 0$ ):

$$Y = \beta_0 + \beta_1 X_1$$

For a Female (with  $X_2 = 1$ ):

$$Y = \beta_0 + \beta_1 X_1 + \beta_2$$

The difference in estimates between the two genders, for a given value of  $X_1$ , is exactly  $\beta_2$ .

# Regression with Dummies

The inclusion of both continuous and binary variables allows for nuanced interpretations:

- 1 The continuous variable's effect is captured by its coefficient, illustrating its linear relationship with the outcome.
- 2 The binary variable's coefficient showcases its difference from the base category, adjusting for the continuous variable.

In practice, the choice of variables and their encoding is crucial for meaningful model interpretations.

# Interpreting Binary Variables in Regression Trees

In regression trees, binary variables help decide data splits:

- Presence (or absence) of a feature (e.g., Group A vs. Group B) drives a branching.
- Decisions are made, leading to target predictions based on terminal leaf mean values.

Unlike linear regression, regression trees lack a singular “coefficient”. They provide a segmented relationship representation.



# Multiple Categories and One-hot Encoding

When a variable has multiple categories, we create several dummy variables. This is also called one-hot encoding. Consider “Region” with levels: North, East, and West:

- ‘Region:North’: 1 if North, 0 otherwise.
- ‘Region:East’: 1 if East, 0 otherwise.

The West category (omitted) becomes our base.

# Interpreting Multiple Categories in Linear Regression

For one-hot encoding in linear regression:

- ① Each dummy's coefficient indicates the expected dependent variable difference compared to the base.
- ② E.g., For 'Region:North', the coefficient shows the expected dependent variable difference for North vs. West, holding other variables constant.

# Regression Model with Multiple Categories

Let's include a continuous variable  $X_1$  (e.g., income) along with our one-hot encoded region variables.

Our regression equation becomes:

$$E(Y) = \beta_0 + \beta_1 X_1 + \beta_2 \text{Region:North} + \beta_3 \text{Region:East}$$

Where:

- $\beta_0$  is the intercept.
- $\beta_1$  is the coefficient for the continuous variable  $X_1$ .
- $\beta_2$  and  $\beta_3$  are coefficients for the dummy variables.

# Interpreting Coefficients with Multiple Categories

- 1  $\beta_1$  captures the expected change in  $Y$  for a one-unit increase in  $X_1$ , holding other variables constant.
- 2  $\beta_2$  shows the expected change in  $Y$  when the region is North, compared to West (base), while keeping  $X_1$  and other factors constant.
- 3  $\beta_3$  provides the expected change in  $Y$  when the region is East, relative to West (base), adjusting for  $X_1$  and other variables.

For instance, a positive  $\beta_2$  suggests that the North region is associated with a higher value of  $Y$  than the West region, given the same income level.

# Impact on Regression Estimates

Given:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 \text{Region:North} + \beta_3 \text{Region:East}$$

For the West region (base):

$$Y = \beta_0 + \beta_1 X_1$$

For the North region:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2$$

For the East region:

$$Y = \beta_0 + \beta_1 X_1 + \beta_3$$

The differences between each region's estimates and the base (West), for a given  $X_1$  value, are  $\beta_2$  and  $\beta_3$ , respectively.

# Regression with Dummies

Incorporating multiple categories in regression:

- ① Enables a nuanced understanding of the effects of each category, adjusting for continuous variables.
- ② Provides insights into differences from the base category.

Effective model interpretation hinges on a thoughtful choice of base category and careful coefficient examination.

# Interpreting Binary Variables in Linear Regression

In linear regression, the coefficient for a binary variable shows the expected change in the dependent variable for the encoded value compared to its base level, all else held equal.

- The choice of the base level can be strategic. Often, it's either a natural reference point, the most frequent category, or a category of special interest.
- By comparing other levels to the base, we isolate the effect of each level, making coefficients interpretable as differences from the base level.
- The absence of a specific category coefficient means the predicted value is the intercept for the base level, and any other coefficients are adjustments from this baseline.

# Choosing a Base Level: An Example

Suppose you're analyzing data on sales of a new tech gadget across various regions: North, South, East, and West.

- ① **Natural Reference Point:** If the company's headquarters are in the North, the North might be chosen as the base level. It can serve as a region for standard comparison since it might be where main operations occur.
- ② **Most Frequent Category:** Upon data exploration, if most sales occur in the East, it might be chosen as the base. This way, other regions' coefficients represent deviations from the major sales region.
- ③ **Category of Special Interest:** If launching a targeted campaign in the West, choosing the West as the base can help gauge its efficacy relative to other regions.

The choice of base level often aligns with the research question or business problem at hand.



# Interpreting Multiple Categories in Regression Trees

In regression trees:

- Each category can be a split point. A split at 'Region:North' divides data based on North membership.
- Trees provide segmented insights without a “coefficient”, with predictions based on terminal leaf means.

# Key Differences in Interpretation

- ① **Linear Regression:** Coefficients offer insights on the dependent variable's expected change based on predictor levels.
- ② **Regression Trees:** Without coefficients, trees use categorical variables for branching, leading to segmented predictions based on terminal leaf values.

Both models offer distinct but valuable insights.