3/23/2023

Categorical (independent) variables in regression:

Dummy variables, also known as indicator variables, are used in regression analysis to represent categorical variables in a quantitative form. In regression analysis, we typically use continuous variables as predictors to estimate the relationship between the predictors and the outcome variable. However, when we have categorical variables as predictors, we need to encode them as dummy variables to include them in the regression model.

The reason for encoding categorical variables as dummy variables is to represent the categorical variable as a set of binary variables. For example, suppose we have a categorical variable "color" with three possible values: red, blue, and green. We can create three dummy variables: "red," "blue," and "green," each of which takes on the value of 1 when the corresponding value of the categorical variable is present and 0 otherwise. These dummy variables can then be included as predictors in a regression model.

Encoding categorical variables as dummy variables allows us to estimate the effect of each category on the outcome variable, while controlling for other variables in the model. It also allows us to model non-linear relationships between the categories and the outcome variable.

There are several methods for encoding dummy variables in regression analysis, including:

**One-hot encoding**: This method involves creating a separate binary variable for each category of the categorical variable. For example, if the categorical variable is "color" with three possible values (red, blue, green), we would create three binary variables (color_red, color_blue, color_green). Each binary variable takes the value of 1 if the corresponding value of the categorical variable is present, and 0 otherwise.

**Effect coding**: In this method, we assign a value of -1 to one category of the categorical variable, and a value of 1/(k-1) to the remaining categories, where k is the number of categories. This allows us to estimate the effect of each category relative to the reference category.

**Helmert coding**: This method involves assigning a value of -1 to the first category of the categorical variable, and the negative of the mean of the previous categories to the subsequent categories. This allows us to estimate the difference between each category and the mean of the previous categories.

**Backward difference coding**: In this method, we assign a value of -1 to the last category of the categorical variable, and the difference between the previous category and the last category to the previous category. This allows us to estimate the difference between each category and the previous category.

The choice of encoding method depends on the research question, the number of categories, and the properties of the data. One-hot encoding is the most common method and is often used when the number of categories is small. Effect coding is useful when we are interested in comparing the effect of each category to a reference category. Helmert coding and backward difference coding are useful when we want to estimate the differences between each category and the previous or subsequent categories.

Ordinal variables are categorical variables that have a natural order or ranking. Examples of ordinal variables include education level (e.g., high school, college, graduate degree) and income level (e.g., low, medium, high).

Ordinal variables can be used in regression analysis as predictors, but it is important to consider the assumptions of the regression model. In particular, when using ordinal variables, the regression model assumes that the distance between adjacent categories is equal.

For example, if we have an ordinal variable for education level with three categories (high school, college, graduate degree), the model assumes that the difference between high school and college is the same as the difference between college and graduate degree. However, this assumption may not be valid if the distance between categories is not equal.

In practice, it may be difficult to determine whether the assumption of equal distance between categories holds. If there is uncertainty about the validity of this assumption, it may be better to treat the ordinal variable as a categorical variable and use one of the methods for encoding dummy variables.

In some cases, it may be useful to encode quantitative variables in categories or bins when doing regression analysis. This is known as binning or discretization.

One reason for binning a quantitative variable is to handle nonlinearity in the relationship between the predictor and the response variable. Binning can help to capture nonlinear relationships that may not be apparent when using the original quantitative variable. For example, if we have a predictor variable of age, binning it into age groups (e.g. 20-29, 30-39, 40-49, etc.) may help to capture nonlinearities in the relationship between age and the response variable.

Another reason for binning a quantitative variable is to handle data sparsity or data imbalance. In some cases, there may be too few data points in certain regions of the variable's range, making it difficult to model the relationship between the predictor and response. Binning can help to aggregate the data into more evenly distributed categories, which may improve the model's performance.

However, there are also potential drawbacks to binning. Binning can lead to information loss, as the original continuous variable is transformed into discrete categories. This can reduce the model's predictive power and increase the risk of overfitting. Additionally, the choice of bin size or number of bins can have a significant impact on the results, and there is no universally optimal choice.

Using *discrete variables* or dummy variables in a regression model can impact the interpretation of the model's coefficients and the overall model fit.

When using a discrete variable in a regression model, the model fits a separate intercept and slope for each level of the variable. For example, if we have a categorical variable for gender with two levels (male and female), the model will fit separate coefficients for male and female. This allows us to estimate the effect of each level of the variable on the response variable, relative to the reference level.

When using dummy variables in a regression model, the interpretation of the coefficients can be slightly different. Dummy variables are binary variables that take on a value of 1 or 0 to indicate the presence or absence of a particular category. In this case, the reference level is represented by a 0 for all dummy variables. The coefficient for each dummy variable represents the change in the response variable

associated with a change from the reference level to the level represented by the dummy variable, while holding all other variables constant.

One potential issue with using dummy variables in a regression model is that it can lead to collinearity, where one or more dummy variables are perfectly correlated with each other. This can make it difficult to estimate the individual effects of each level of the variable and can lead to unstable estimates of the coefficients.

We can create dummy variables manually using ifelse() functions, or there are packages that can perform the transformations for us. It's always necessary to inspect the outcome of the transformations after they are performed to make sure we don't include redundant information in the dataframe.

**Original Data**

| Team | Points |
|------|--------|
| A | 25 |
| A | 12 |
| B | 15 |
| B | 14 |
| B | 19 |
| B | 23 |
| C | 25 |
| C | 29 |

**One-Hot Encoded Data**

| Team_A | Team_B | Team_C | Points |
|--------|--------|--------|--------|
| 1 | 0 | 0 | 25 |
| 1 | 0 | 0 | 12 |
| 0 | 1 | 0 | 15 |
| 0 | 1 | 0 | 14 |
| 0 | 1 | 0 | 19 |
| 0 | 1 | 0 | 23 |
| 0 | 0 | 1 | 25 |
| 0 | 0 | 1 | 29 |

In the example above, the three columns are not needed to encode three levels here. At least one of the Team dummy variables should be removed from the resulting dataset.

**Factor analysis** is a statistical method used to identify patterns in a large dataset by reducing the number of variables and summarizing the underlying relationships between the original variables. It is often used in the social sciences, psychology, and other fields to identify underlying constructs or dimensions that are not directly measured.

The basic idea behind factor analysis is to identify a smaller number of factors, which are linear combinations of the original variables, that capture the majority of the information in the data. These factors are derived by looking for patterns of correlations between the original variables. Variables that are highly correlated with each other are assumed to be measuring the same underlying construct or factor.

There are several methods for extracting factors in factor analysis, including principal component analysis (PCA) and maximum likelihood estimation (MLE). Once the factors have been extracted, they can be rotated to simplify their interpretation and make them more meaningful. The most common rotation methods are Varimax, Promax, and Oblimin.

Factor analysis can be used for several purposes, including:

*Dimension reduction*: Factor analysis can be used to reduce the number of variables in a dataset while retaining as much of the original information as possible.

*Identifying underlying constructs*: Factor analysis can be used to identify the underlying constructs or dimensions that are not directly measured in a dataset.

*Validating a measurement instrument*: Factor analysis can be used to determine whether the items in a measurement instrument are measuring the same underlying construct.

Overall, factor analysis is a powerful tool for identifying patterns and relationships in a large dataset, and can be useful for exploring complex datasets in a variety of fields.

Factor analysis involves several steps. Here is a general overview of the process:

Determine the research question: The first step is to determine the research question or problem that you want to investigate. This will guide your selection of variables and the type of factor analysis to use. Select the variables: Identify the variables that are relevant to your research question. These variables should be measured on an interval or ratio scale.

Check for data suitability: Check the suitability of the data for factor analysis. You can do this by examining the correlation matrix to ensure that the variables are highly correlated with each other. Choose a factor analysis method: There are different factor analysis methods to choose from, including principal component analysis (PCA), maximum likelihood estimation (MLE), and others. The choice of method depends on the research question and the characteristics of the data.

Determine the number of factors: Decide on the number of factors to extract. This can be based on a scree plot, eigenvalues, or other statistical methods.

Extract the factors: Extract the factors using the chosen method. This will produce a matrix of factor loadings, which represents the strength of the relationship between each variable and each factor.

Interpret the factors: Interpret the factors based on the factor loadings. This involves examining the variables with high loadings on each factor and determining what they have in common.

Rotate the factors: Rotate the factors to make them easier to interpret. This involves rotating the factor loadings to simplify the pattern of relationships between the variables and the factors.
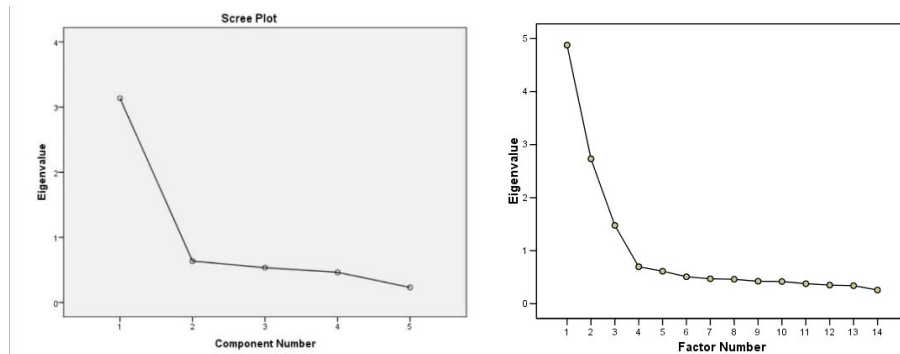
Evaluate the model: Evaluate the model to ensure that it is a good fit for the data. This can be done by examining the goodness-of-fit statistics and comparing the model to alternative models.

Overall, factor analysis is a complex statistical method that requires careful consideration of the research question, data, and interpretation of the results. It is important to have a clear understanding of the underlying assumptions and limitations of the method before applying it to a dataset.

A scree plot is a graphical representation of the eigenvalues associated with the factors extracted in a factor analysis. It is used to determine the appropriate number of factors to retain in the analysis. The scree plot is a line graph where the x-axis represents the number of factors and the y-axis represents the corresponding eigenvalues. The plot typically shows a sharp drop in eigenvalues for the first few factors, followed by a more gradual decline in the subsequent factors. The point at which the eigenvalues start to level off is often used as a cut-off for retaining factors, as factors beyond this point are likely to

explain relatively little variance in the data. However, the decision of how many factors to retain ultimately depends on the research question, the characteristics of the data, and other statistical methods. Scree plots are a useful tool for visualizing the eigenvalues and providing a simple way to assess the dimensionality of the data in factor analysis.

The scree plot on the left suggests we should retain only one factor. The one on the right indicates three should be kept.

Resources:
1. https://www.statisticssolutions.com/dummy-coding-the-how-and-why/
2. https://www.analyticsvidhya.com/blog/2020/08/types-of-categorical-data-encoding/
3. https://towardsdatascience.com/all-about-categorical-variable-encoding-305f3361fd02
4. https://timeseriesreasoning.com/contents/dummy-variables-in-a-regression-model/
5. https://stats.oarc.ucla.edu/spss/faq/coding-systems-for-categorical-variables-in-regression-analysis/
6. https://www.learndatasci.com/glossary/dummy-variable-trap/
7. https://machinelearningmastery.com/one-hot-encoding-for-categorical-data/
8. https://www.iosrjournals.org/iosr-jm/papers/vol2-issue5/A0250107.pdf
9. https://towardsdatascience.com/an-introduction-to-discretization-in-data-science-55ef8c9775a2
10. https://machinelearningmastery.com/discretization-transforms-for-machine-learning/
11. https://stats.oarc.ucla.edu/spss/seminars/introduction-to-factor-analysis/a-practical-introduction-to-factor-analysis/
12. https://www.qualtrics.com/experience-management/research/factor-analysis/
13. https://www.theanalysisfactor.com/factor-analysis-1-introduction/
14. https://www.researchgate.net/figure/Limitations-of-scatterplots-for-managing-discrete-variables-a-The-scatterplot-matrix_fig2_319328127
15. https://www.statology.org/one-hot-encoding-in-r/