

How to handle categorical features in Correlation based feature selection

Let's delve into how to handle categorical features within the framework of Correlation-Based Feature Selection. As we discussed earlier, standard correlation measures like Pearson's work directly with numerical data. Therefore, to include categorical features, we need to transform them into a numerical representation first.

Here's a step-by-step explanation with an example:

Steps:

1. **Encode the Categorical Feature(s):** The crucial first step is to convert your categorical features into a numerical format. The most common and generally recommended technique for nominal (unordered) categorical features in this context is **One-Hot Encoding**.
2. **Calculate Correlation:** Once the categorical features are numerically encoded (as one or more binary columns), you can calculate the correlation between these new numerical columns and your numerical target variable using a standard correlation method like Pearson's.
3. **Interpret the Correlations:** The correlation coefficient will indicate the strength and direction of the linear relationship between the presence of a specific category (represented by the binary column) and the target variable.

Example:

Let's go back to our house price prediction scenario. Assume all other features (size, bedrooms, bathrooms, age) are numerical, but we also have a categorical feature: **Neighborhood**.

The "Neighborhood" feature might have categories like: "Downtown", "Suburb A", "Suburb B", "Rural".

1. Encode the Categorical Feature:

We'll use One-Hot Encoding to transform the "Neighborhood" feature:

Original Neighborhood	Is_Downtown	Is_Suburb_A	Is_Suburb_B	Is_Rural
Downtown	1	0	0	0
Suburb A	0	1	0	0
Suburb B	0	0	1	0
Rural	0	0	0	1
Downtown	1	0	0	0
...

Now, the categorical "Neighborhood" feature is represented by four new numerical (binary) features.

2. Calculate Correlation:

Next, we calculate the Pearson correlation coefficient between each of these new binary "Neighborhood" columns and our numerical target variable, "Price". Let's assume we get the following results:

Encoded Feature	Correlation with Price
Is_Downtown	0.65
Is_Suburb_A	0.40
Is_Suburb_B	0.30
Is_Rural	-0.50

3. Interpret the Correlations:

- **Is_Downtown (Correlation: 0.65):** There is a moderate positive correlation between a house being in the "Downtown" neighborhood and its price. This suggests that houses in the downtown area tend to have higher prices on average.
- **Is_Suburb_A (Correlation: 0.40):** There is a weak positive correlation between a house being in "Suburb A" and its price. Houses in this suburb tend to be slightly more expensive on average compared to the baseline (the category not explicitly represented by a '1' in these columns).
- **Is_Suburb_B (Correlation: 0.30):** Similar to Suburb A, there's a weak positive correlation with price.
- **Is_Rural (Correlation: -0.50):** There is a moderate negative correlation between a house being in a "Rural" area and its price. This suggests that rural houses tend to have lower prices on average compared to other neighborhoods.

Using these Correlations for Feature Selection:

Now, when you are performing correlation-based feature selection, you would consider these correlation values alongside the correlations of your other numerical features (size, bedrooms, etc.). You could then rank all features (both originally numerical and the encoded categorical ones) based on the absolute value of their correlation with the "Price" and select the top-k features or those above a certain threshold.

Important Considerations:

- **Baseline Category:** When interpreting the correlations of one-hot encoded features, one category acts as the implicit baseline (the one where all the binary columns for that feature are 0). The correlations of the other categories are relative to this baseline.
- **Number of Categories:** If a categorical feature has a very large number of unique categories, one-hot encoding can lead to a high-dimensional dataset¹ (many new binary features). In such cases, you might consider other encoding techniques or feature reduction methods after one-hot encoding.
- **Alternative Encoding for Ordinal Features:** If your categorical feature has an inherent order (e.g., "Quality": Low, Medium, High), you could consider Label Encoding (assigning numerical values like 1, 2, 3) and then calculate the correlation. However, be mindful that Pearson correlation will assume a linear relationship between the ordered categories and the target.

In summary, to handle categorical features in Correlation-Based Feature Selection:

1. **One-Hot Encode** nominal categorical features into multiple binary numerical columns.
2. **Calculate the Pearson correlation** between each of these binary columns and the numerical target variable.
3. **Use the absolute values of these correlation coefficients** along with the correlations of your other numerical features to rank and select the most relevant features.

This approach allows you to quantify the linear relationship between the presence of each category of a nominal feature and the target variable within the correlation-based feature selection framework.