

# Understanding SHAP's Approach to Binary-Encoded Categorical Features

SHAP (SHapley Additive exPlanations) has become a cornerstone methodology in explainable AI, providing crucial insights into how machine learning models make predictions. When dealing with categorical features—especially those transformed through binary encoding—understanding how SHAP interprets and processes these features requires special consideration. This report explores the mechanisms, challenges, and best practices for handling binary-encoded categorical features within the SHAP framework.

## Fundamentals of SHAP Values for Categorical Features

SHAP values, derived from cooperative game theory, quantify each feature's contribution to a model's prediction. These values distribute the difference between a prediction and the average prediction among all features. For numerical features, this process is relatively straightforward. However, categorical features require encoding before most machine learning models can process them, creating additional complexity in interpretation.

When categorical variables undergo binary encoding, each category is represented as a series of binary digits (0s and 1s), creating multiple columns from a single original feature. Unlike one-hot encoding which creates a separate column for each category, binary encoding reduces dimensionality by using  $\log_2(n)$  binary features to represent  $n$  categories. This efficiency comes with interpretability challenges for SHAP analysis, as the individual binary features lack the immediate semantic meaning that one-hot encoded features retain.

The feature importance in SHAP is typically defined as the mean absolute value of the SHAP values for a given feature<sup>[1]</sup>. This definition applies to both numerical and encoded categorical features, though aggregation becomes necessary for the latter to recover the importance of the original categorical variable.

## SHAP Calculation Process for Encoded Features

SHAP does not inherently distinguish between original numerical features and transformed categorical features. Instead, it calculates contribution values for each input feature as presented to the model. When a categorical feature is encoded into multiple binary columns, SHAP treats each of these columns as separate features during its computation process.

For binary classification problems, SHAP values are typically interpreted in terms of log odds, while multiclass targets use softmax transformations to appropriately scale the contributions<sup>[2]</sup>. This mathematical foundation applies regardless of whether the input features are numerical or encoded categorical variables.

## Aggregation Methods for Binary-Encoded Categorical Features

A key challenge with binary-encoded categorical features is that the SHAP values are calculated for each binary column independently, but understanding the overall importance of the original categorical feature requires aggregating these values appropriately.

According to guidance from SHAP's creator, to compute the feature importance of a categorical variable encoded with one-hot encoding, you should sum the SHAP values of all components for each observation first, and then calculate the mean of the absolute values across observations<sup>[1]</sup>. This same principle applies to binary encoding:

1. For each observation, sum the SHAP values across all binary columns derived from a single categorical feature
2. Take the absolute value of each sum (to capture both positive and negative contributions)
3. Calculate the mean across all observations to get the global importance of the categorical feature

This aggregation method ensures that the importance attributed to the categorical feature reflects its total impact on predictions, rather than diluting it across multiple encoded columns<sup>[1]</sup>. Without proper aggregation, the importance of categorical features may be underestimated compared to numerical features.

## Interpretation Challenges for Binary Encoding

The interpretation of SHAP values for binary-encoded features presents unique challenges. When examining force plots or summary plots, you may observe patterns like "Gender\_female=0" with a contribution margin X and "Gender\_male=1" with a contribution margin Y<sup>[3]</sup>. These patterns can be confusing because the binary digits themselves don't have inherent meaning—they're just mathematical representations of categories.

For example, if a categorical feature "Color" with values "Red," "Green," and "Blue" is binary encoded, it might create two columns where:

- Red =
- Green = <sup>[1]</sup>
- Blue = <sup>[1]</sup>

The SHAP values for "Color\_bit1=0" and "Color\_bit2=1" have no intuitive interpretation in isolation. Their meaning only emerges when mapped back to the original categories and aggregated appropriately.

## Implementation Approaches and Practical Considerations

Several practical approaches exist for handling binary-encoded categorical features in SHAP analysis:

## Direct Aggregation in Post-Processing

After obtaining SHAP values for all features, you can programmatically aggregate values for binary columns belonging to the same categorical feature. This requires maintaining metadata about which binary columns correspond to which original features.

```
# Example aggregation approach (conceptual)
shap_values_aggregated = {}
for feature in original_categorical_features:
    binary_columns = get_binary_columns_for_feature(feature)
    feature_shap_values = shap_values[:, binary_columns].sum(axis=1)
    shap_values_aggregated[feature] = np.mean(np.abs(feature_shap_values))
```

## Alternative Encoding Strategies

Some researchers have explored alternative encoding approaches specifically designed to work better with SHAP. For instance, research has investigated hybridizing target-encoded and SHAP-encoded features to improve algorithm selection for optimization problems<sup>[4]</sup>. While these approaches don't directly address binary encoding, they highlight the importance of considering encoding strategies in conjunction with SHAP analysis.

## Native Categorical Support

Some model implementations, like certain versions of XGBoost, support categorical features directly without requiring explicit encoding<sup>[5]</sup>. However, this native support may create compatibility challenges with SHAP implementation, requiring special handling to ensure proper explanation generation.

## Visualization Techniques for Binary-Encoded Categorical Features

Standard SHAP visualization methods can be adapted to better represent binary-encoded categorical features:

### Summary Plots

When creating summary plots for models with binary-encoded categorical features, consider grouping binary columns from the same original feature together or pre-aggregating their SHAP values before visualization. This provides a clearer picture of the overall importance of categorical features<sup>[6]</sup>.

### Dependence Plots

For dependence plots, mapping encoded values back to their original categories can enhance interpretability. Instead of showing how SHAP values vary with binary digit values (which isn't semantically meaningful), show how they vary with the original categorical values<sup>[7]</sup>.

## Force Plots

Force plots for individual predictions can be overwhelming when many binary columns are present. Consider creating custom force plots that group and aggregate the impact of all binary columns from the same original feature, presenting them as a single unit.

## Advanced Techniques and Research Directions

Several advanced techniques are emerging to better handle categorical features in SHAP:

### Class-Specific SHAP Analysis

For classification problems, analyzing SHAP values on a per-class basis can provide deeper insights. By subsetting SHAP values according to true class labels, you can identify features that are particularly important for specific classes<sup>[6]</sup>. This technique applies to both numerical and categorical features but can be especially valuable for understanding how categorical features influence different class predictions.

### Hybrid Encoding Approaches

Research suggests that hybridizing different encoding schemes can improve performance in downstream tasks. While some encoding methods might not individually lead to better results, combining approaches like target-encoding and SHAP-based encoding can outperform individual methods<sup>[4]</sup>. This indicates potential for developing specialized encoding strategies that optimize both model performance and explainability.

## Conclusion: Best Practices for Binary-Encoded Categorical Features in SHAP

Based on the evidence from practical implementations and research, several best practices emerge for handling binary-encoded categorical features in SHAP analysis:

1. Consistently aggregate SHAP values across all binary columns derived from the same original categorical feature to accurately assess feature importance.
2. Maintain clear documentation of the encoding scheme used and which binary columns correspond to which original categorical features.
3. Consider feature importance at both the original categorical feature level (through aggregation) and at the individual binary column level to gain comprehensive insights.
4. When visualizing SHAP values, adapt standard plots to better represent categorical features by grouping or aggregating related binary columns.
5. For classification problems, explore class-specific SHAP analyses to understand how categorical features influence predictions for particular classes.

By following these practices, data scientists and machine learning practitioners can more effectively leverage SHAP for explainable AI applications involving categorical features, ensuring that explanations accurately reflect the role of these features in model predictions regardless of the encoding method used.

1. [https://www.reddit.com/r/datascience/comments/s2epy0/computing\\_categorical\\_feature\\_importance\\_using/](https://www.reddit.com/r/datascience/comments/s2epy0/computing_categorical_feature_importance_using/)
2. <https://www.youtube.com/watch?v=2xlgOu22YgE>
3. <https://github.com/slundberg/shap/issues/2270>
4. <https://arxiv.org/html/2407.07439v1>
5. <https://stackoverflow.com/questions/75899158/shap-summary-plots-for-xgboost-with-categorical-data-inputs>
6. <https://stackoverflow.com/questions/65110798/feature-importance-in-a-binary-classification-and-extracting-shap-values-for-one>
7. [https://docs.seldon.io/projects/alibi/en/latest/examples/kernel\\_shap\\_adult\\_lr.html](https://docs.seldon.io/projects/alibi/en/latest/examples/kernel_shap_adult_lr.html)