

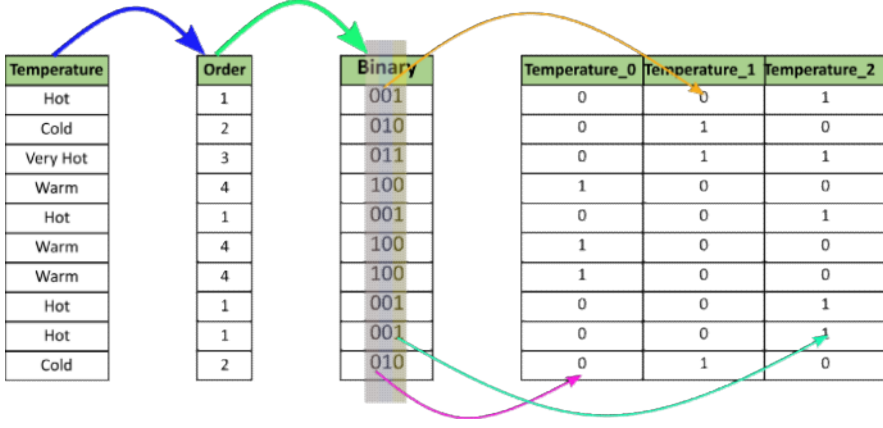
# Encoding methods

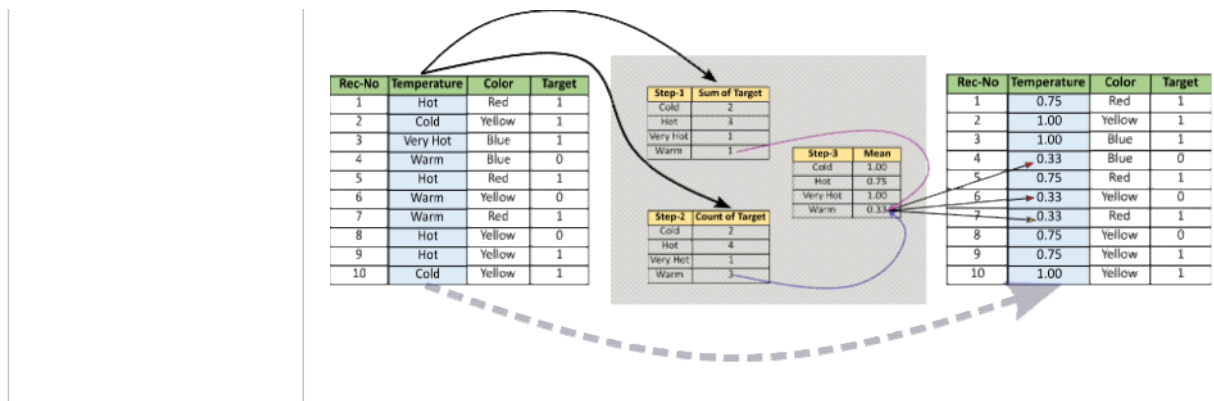
09 October 2022 10:56 PM

<https://analyticsindiamag.com/a-complete-guide-to-categorical-data-encoding/#:~:text=Encoding%20categorical%20data%20is%20a,provided%20to%20the%20different%20models.&text=In%20the%20field%20of%20data,preparation%20is%20a%20mandatory%20task.>

<https://towardsdatascience.com/all-about-categorical-variable-encoding-305f3361fd02>

<https://www.analyticsvidhya.com/blog/2020/08/types-of-categorical-data-encoding/>

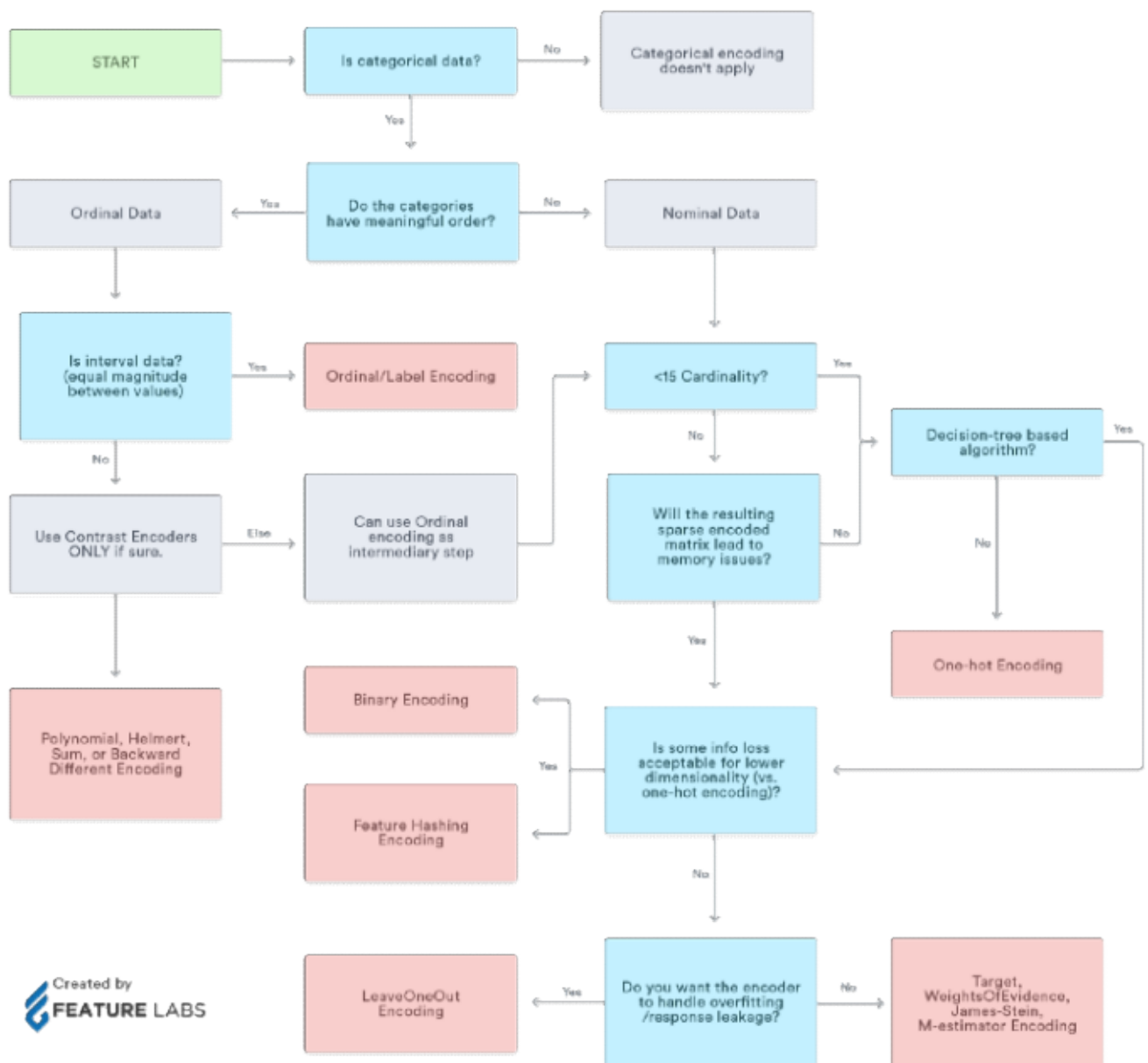
Encoding method	Description
Label/Ordinal encoding	<ul style="list-style-type: none"> <li>• <u>Converts each label into integer, with value representing sequence</u></li> <li>• Used for ordinal data</li> </ul>
One-Hot Encoding	<ul style="list-style-type: none"> <li>• <u>Each class gets becomes a new variable, called dummy variables - maps each class into binary</u></li> <li>• Used for nominal data</li> </ul> <p>Note: Dummy encoding is the same but you remove one level</p>
Effect/Deviation/Sum Encoding	<ul style="list-style-type: none"> <li>• <u>Classes are put into -1, 0, 1 format (has an additional -1 compared to one-hot)</u></li> <li>• <a href="https://jamesmccaffrey.wordpress.com/2019/10/17/effect-coding-vs-one-hot-encoding-for-neural-networks/">https://jamesmccaffrey.wordpress.com/2019/10/17/effect-coding-vs-one-hot-encoding-for-neural-networks/</a></li> <li>• Usually used in multiple linear regression, when two categorical variables statistically interact</li> <li>• People tend to use one-hot over effect because it looks nicer and it is symmetrical</li> <li>• However, both works the same - no free lunch theorem</li> </ul>
Binary Encoding	<ul style="list-style-type: none"> <li>• <u>Converts a category into binary digits</u></li> <li>• Requires less columns than one-hot</li> </ul>  <p>Note: Base-N Encoding is an extension of binary encoding, where the base is a hyperparameter</p>
Mean/Target Encoding	<ul style="list-style-type: none"> <li>• <u>Converting a class into the mean of the target variable for its class (type of Bayesian encoding)</u></li> <li>• Brings out relation between similar categories</li> <li>• Does not affect volume of data</li> <li>• Helps in faster learning</li> <li>• BUT, tends to overfit</li> </ul>



Note: There are two types of categorical data - **Nominal** (no order) and **Ordinal** (order matters)

<https://towardsdatascience.com/all-about-categorical-variable-encoding-305f3361fd02>

### Categorical Encoding Methods Cheat-Sheet



Dataset

- 33 variables; 32 features & 1 target variables