# 20th Mar Assignment

March 23, 2023

## 1 Assignment 44

**Q1. What is data encoding? How is it useful in data science?**

**Ans. Data encoding is the process of converting data from one format or representation to another format. In data science, data encoding is useful for several reasons:**

- Compression: Encoding can reduce the size of the data, making it easier to store and transmit. This is especially important when dealing with large datasets that would otherwise require a lot of storage space or network bandwidth.

- Standardization: Encoding can be used to convert data into a standard format, making it easier to work with and analyze. For example, different datasets may use different date formats, such as "MM/DD/YYYY" or "YYYY/MM/DD". Encoding can be used to convert all the dates into a common format for easier analysis.

- Security: Encoding can be used to encrypt data, making it more difficult for unauthorized users to access or read. This is important when dealing with sensitive or confidential data.

- Machine Learning: Many machine learning algorithms require data to be encoded in a specific way. For example, categorical data (such as "red", "green", and "blue") may need to be encoded as numerical values (such as 1, 2, and 3) before it can be used in a machine learning algorithm.

**Overall, data encoding is an important tool for data scientists to use in order to make data more manageable, standardized, secure, and usable for analysis and machine learning.**

**Q2. What is nominal encoding? Provide an example of how you would use it in a real-world scenario.**

**Ans. Nominal encoding is a type of data encoding that is used to convert categorical variables into numerical variables. In nominal encoding, each unique category is assigned a unique integer value. Unlike ordinal encoding, the values assigned to each category have no inherent order or meaning.**

**For example, suppose we have a dataset of customer information for an online store, and one of the categorical variables is "country of residence". The possible values for this variable are "USA", "Canada", "Mexico", "Germany", "France", and "Japan".**

To use this variable in a machine learning algorithm, we could use nominal encoding to convert it into numerical values. One possible encoding scheme would be:

- USA: 1
- Canada: 2
- Mexico: 3
- Germany: 4
- France: 5
- Japan: 6 #### In this encoding scheme, each unique country is assigned a unique integer value. These values have no inherent order or meaning, but they can be used in a machine learning algorithm that requires numerical inputs.

A real-world scenario where nominal encoding might be used is in predicting customer churn for a subscription-based service. One of the categorical variables in the dataset might be "plan type", which could have values like "basic", "standard", and "premium". Nominal encoding could be used to convert this variable into numerical values, which could then be used in a machine learning algorithm to predict which customers are most likely to churn based on their plan type.

**Q3. In what situations is nominal encoding preferred over one-hot encoding? Provide a practical example.**

**Ans.** Nominal encoding and one-hot encoding are both useful techniques for converting categorical variables into numerical variables. However, there are some situations where nominal encoding may be preferred over one-hot encoding.

One situation where nominal encoding may be preferred is when the number of unique categories is large. One-hot encoding requires creating a new binary column for each unique category, which can lead to a very large number of columns if there are many unique categories. Nominal encoding, on the other hand, requires only a single column to represent all the categories, which can be more efficient in terms of memory and computational resources.

Another situation where nominal encoding may be preferred is when there is no clear order or hierarchy among the categories. One-hot encoding assumes that there is a clear order or hierarchy among the categories, but this may not always be the case. In such situations, using nominal encoding may be more appropriate.

For example, consider a dataset of restaurant reviews that includes a categorical variable for the type of cuisine, which includes categories such as "Chinese", "Italian", "Mexican", "Indian", "Japanese", and "Thai". One-hot encoding would require creating a separate binary column for each of these categories, which could result in a large number of columns. Nominal encoding, on the other hand, would only require a single column to represent all the categories. Since there is no clear order or hierarchy among these cuisine types, nominal encoding may be the more appropriate choice.

**Q4.** Suppose you have a dataset containing categorical data with 5 unique values. Which encoding technique would you use to transform this data into a format suitable for machine learning algorithms? Explain why you made this choice.

**Ans.** The choice of encoding technique depends on the nature of the categorical data and the requirements of the machine learning algorithm being used. However, assuming that the categorical data has no inherent order or hierarchy, and that the machine learning algorithm can handle numerical inputs, one possible choice of encoding technique is nominal encoding.

Nominal encoding would assign a unique integer value to each of the 5 unique categories in the dataset. This would allow the categorical data to be represented as numerical data, which can be used as input to a machine learning algorithm. Since there is no inherent order or hierarchy among the categories, and the number of unique categories is relatively small, nominal encoding would be a simple and efficient choice.

It's worth noting that other encoding techniques, such as one-hot encoding or ordinal encoding, could also be used depending on the specific requirements of the machine learning algorithm and the nature of the categorical data. One-hot encoding, for example, would be more appropriate if there are many unique categories, while ordinal encoding would be more appropriate if the categories have a clear order or hierarchy.

**Q5.** In a machine learning project, you have a dataset with 1000 rows and 5 columns. Two of the columns are categorical, and the remaining three columns are numerical. If you were to use nominal encoding to transform the categorical data, how many new columns would be created? Show your calculations.

**Ans.** If we use nominal encoding to transform the two categorical columns in the dataset, we would create two new columns, one for each categorical column. Each new column would have a unique integer value assigned to each category in the original categorical column.

Assuming that each categorical column has k unique categories, the number of new columns created using nominal encoding would be k. In this case, we have two categorical columns, so we would create two new columns.

Therefore, in the given machine learning project, if we use nominal encoding to transform the categorical data, we would create **2 new columns.**

**Q6.** You are working with a dataset containing information about different types of animals, including their species, habitat, and diet. Which encoding technique would you use to transform the categorical data into a format suitable for machine learning algorithms? Justify your answer.

Ans. The choice of encoding technique depends on the nature of the categorical data and the requirements of the machine learning algorithm being used. However, assuming that the categorical variables in the animal dataset have no inherent order or hierarchy, and that the machine learning algorithm can handle numerical inputs, one possible choice of encoding technique is nominal encoding.

Nominal encoding would assign a unique integer value to each unique category in each of the categorical variables in the animal dataset. For example, the "species" variable might have categories like "lion", "tiger", and "bear", while the "habitat" variable might have categories like "forest", "desert", and "ocean". Nominal encoding would assign a unique integer value to each of these categories, allowing them to be represented as numerical data that can be used as input to a machine learning algorithm.

One reason to choose nominal encoding is that it is a simple and efficient technique that can work well with small to medium-sized datasets like the one described. Another reason is that there is no inherent order or hierarchy among the categories in each variable, so nominal encoding would be appropriate.

Alternatively, one-hot encoding could also be used if the number of unique categories is small, or if the machine learning algorithm being used can benefit from having each category represented as a separate binary column. However, one-hot encoding can result in a large number of new columns if there are many unique categories, which can be a problem for larger datasets. Therefore, nominal encoding is a good choice for this animal dataset.

**Q7.** You are working on a project that involves predicting customer churn for a telecommunications company. You have a dataset with 5 features, including the customer's gender, age, contract type, monthly charges, and tenure. Which encoding technique(s) would you use to transform the categorical data into numerical data? Provide a step-by-step explanation of how you would implement the encoding.

**Ans.** For the customer churn prediction project, we have a dataset with one categorical feature "contract type" and four numerical features "gender", "age", "monthly charges", and "tenure". To transform the categorical data into numerical data, we can use the following encoding techniques:

1. Nominal Encoding: Since there is no inherent order or hierarchy among the categories of "contract type", we can use nominal encoding to assign a unique integer value to each category. We can create a new column in the dataset called "contract_type_encoded" and assign a unique integer value to each category, such as 1 for "Month-to-month", 2 for "One year", and 3 for "Two year".

2. No encoding: Since "gender" is a binary categorical feature with only two possible values, we can represent it as numerical data without any encoding. We can assign a value of 0 to "Male" and 1 to "Female".

3. No encoding: Since "age", "monthly charges", and "tenure" are numerical features, they do not require any encoding.

**Therefore, the final encoded dataset would have 5 columns: "gender", "age", "monthly charges", "tenure", and "contract_type_encoded".**

**Here are the steps to implement the encoding:**

1. Load the dataset into a pandas DataFrame.
2. Use nominal encoding to transform the "contract type" feature into a new column called "contract_type_encoded" using the replace() function in pandas.
3. Convert the "gender" feature into numerical data using the replace() function in pandas.
4. Check the data types of all features using the dtypes attribute in pandas. Ensure that all features are represented as numerical data types.
5. If necessary, scale or normalize the numerical features to ensure that they have similar ranges and magnitudes. Use the encoded dataset to train and test machine learning models to predict customer churn.