

# DEALING WITH CATEGORICAL COLUMNS

By:

M. Ashish Reddy

DSWP -> Batch-5

# Encoding Categorical Data

- There are three common approaches for converting ordinal and categorical variables to numerical values. They are:
  1. Ordinal Encoding
  2. One-Hot Encoding
  3. Dummy Variable Encoding

# 1. Ordinal Encoding

- In ordinal encoding, each unique category value is assigned an integer value.
- **For example**, “red” is 1, “green” is 2, and “blue” is 3.
- This is called an ordinal encoding or an integer encoding and is easily reversible. Often, integer values starting at zero are used.

## Input

```
1 # example of a ordinal encoding
2 from numpy import asarray
3 from sklearn.preprocessing import OrdinalEncoder
4 # define data
5 data = asarray(['red'], ['green'], ['blue']))
6 print(data)
7 # define ordinal encoding
8 encoder = OrdinalEncoder()
9 # transform data
10 result = encoder.fit_transform(data)
11 print(result)
```

## Output

```
1 [['red']]
2 [['green']]
3 [['blue']]
4 [[2.]]
5 [[1.]]
6 [[0.]]
```

## 2. One-Hot Encoding

- For categorical variables where no ordinal relationship exists, the integer encoding may not be enough, at best, or misleading to the model at worst.

- Input

```
1 # example of a one hot encoding
2 from numpy import asarray
3 from sklearn.preprocessing import OneHotEncoder
4 # define data
5 data = asarray([[ 'red' ], [ 'green' ], [ 'blue' ]])
6 print(data)
7 # define one hot encoding
8 encoder = OneHotEncoder(sparse=False)
9 # transform data
10 onehot = encoder.fit_transform(data)
11 print(onehot)
```

Output

```
1 [ 'red' ]
2 [ 'green' ]
3 [ 'blue' ]
4 [0. 0. 1.]
5 [0. 1. 0.]
6 [1. 0. 0.]
```

# 3. Dummy Variable Encoding

- The one-hot encoding creates one binary variable for each category.

- Input

```
1 # example of a dummy variable encoding
2 from numpy import asarray
3 from sklearn.preprocessing import OneHotEncoder
4 # define data
5 data = asarray(['red'], ['green'], ['blue']))
6 print(data)
7 # define one hot encoding
8 encoder = OneHotEncoder(drop='first', sparse=False)
9 # transform data
10 onehot = encoder.fit_transform(data)
11 print(onehot)
```

## Output

```
1 [['red']]
2 [['green']]
3 [['blue']]
4 [[0. 1.]
5  [1. 0.]
6  [0. 0.]
```

# When to use a Label Encoding vs. One Hot Encoding

- We apply One-Hot Encoding when:
  - 1.The categorical feature is **not ordinal** (like the countries above)
  - 2.The number of categorical features is less so one-hot encoding can be effectively applied
- We apply Label Encoding when:
  - 1.The categorical feature is **ordinal** (like Jr. kg, Sr. kg, Primary school, high school)
  2. The number of categories is quite large as one-hot encoding can lead to high memory consumption

Thank You!