

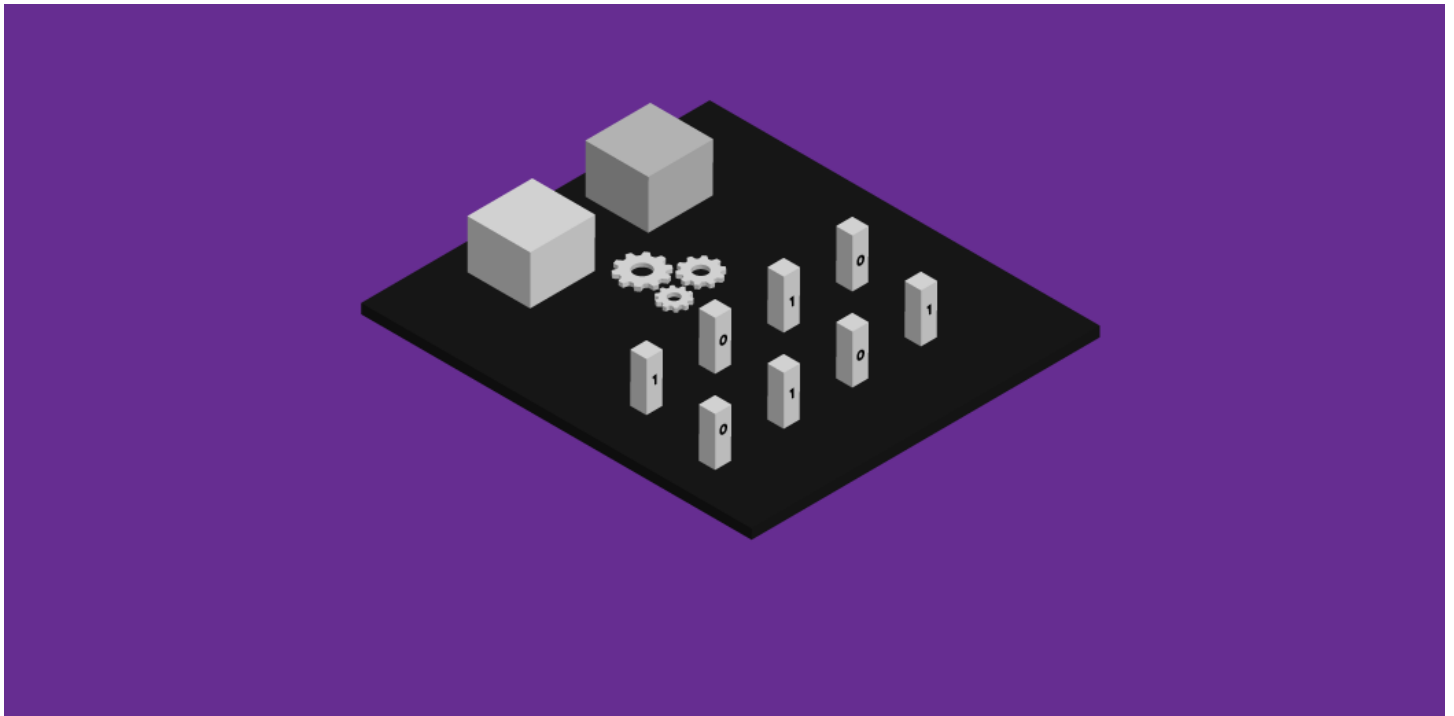
# One-Hot Encoding vs. Label Encoding using Scikit-Learn

[INTERMEDIATE](#)[PYTHON](#)[STRUCTURED DATA](#)[TECHNIQUE](#)

## What is One-Hot Encoding? When should you use One-Hot Encoding over Label Encoding?

These are typical data science interview questions every aspiring data scientist needs to know the answer to. After all, you'll often find yourself having to make a choice between the two in a data science project!

Machines understand numbers, not text. We need to convert each text category to numbers in order for the machine to process them using mathematical equations. Ever wondered how we can do that? What are the different ways?



This is where Label Encoding and One-Hot Encoding come into the picture. We'll discuss both in this article and understand the difference between them.

*Note: Starting your machine learning journey? I recommend taking our comprehensive and popular [Applied Machine Learning course](#)!*

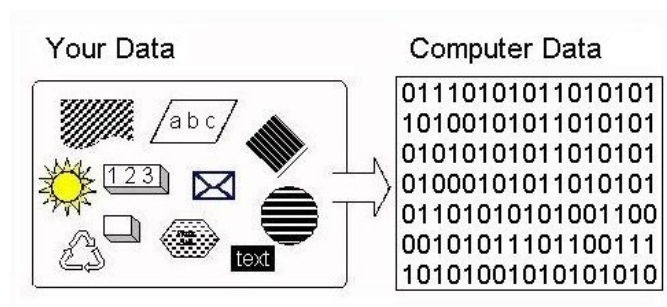
## Table of Contents

1. What is Categorical Encoding?
2. Different Approaches to Categorical Encoding

1. Label Encoding
2. One-Hot Encoding
3. When to use Label Encoding vs. One-Hot Encoding?

## What is Categorical Encoding?

Typically, any structured dataset includes multiple columns – a combination of numerical as well as categorical variables. A machine can only understand the numbers. It cannot understand the text. That's essentially the case with [Machine Learning algorithms](#) too.



That's primarily the reason we need to convert categorical columns to numerical columns so that a machine learning algorithm understands it. This process is called **categorical encoding**.

Categorical encoding is a process of converting categories to numbers.

In the next section, I will touch upon different ways of handling categorical variables.

## Different Approaches to Categorical Encoding

So, how should we handle categorical variables? As it turns out, there are multiple ways of handling Categorical variables. In this article, I will discuss the two most widely used techniques:

- Label Encoding
- One-Hot Encoding

Now, let us see them in detail.

### Label Encoding

**Label Encoding** is a popular encoding technique for handling categorical variables. In this technique, each label is assigned a unique integer based on alphabetical ordering.

Let's see how to implement label encoding in Python using the [scikit-learn library](#) and also understand the challenges with label encoding.

Let's first import the required libraries and dataset:

```
1 #importing the libraries
2 import pandas as pd
3 import numpy as np
4
5 #reading the dataset
6 df=pd.read_csv("Salary.csv")
```

view raw

One8.py hosted with ❤ by GitHub

Output:

Country	Age	Salary
India	44	72000
US	34	65000
Japan	46	98000
US	35	45000
Japan	23	34000

Understanding the datatypes of features:

```
1 print df.info
```

view raw

one9.py hosted with ❤ by GitHub

Output:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 14 entries, 0 to 13
Data columns (total 3 columns):
Country    14 non-null object
Age        14 non-null int64
Salary     14 non-null int64
dtypes: int64(2), object(1)
memory usage: 416.0+ bytes
```

As you can see here, the first column, Country, is the categorical feature as it is represented by the **object data type** and the rest of them are numerical features as they are represented by *int64*.

Now, let us implement label encoding in Python:

```
1 # Import label encoder
2 from sklearn import preprocessing
3 # label_encoder object knows how to understand word labels.
4 label_encoder = preprocessing.LabelEncoder()
5 # Encode labels in column 'Country'.
6 data['Country']= label_encoder.fit_transform(data['Country'])
7 print(data.head())
```

view raw

One3.py hosted with ❤ by GitHub

Output:

Country	Age	Salary
0	44	72000
2	34	65000
1	46	98000
2	35	45000
1	23	34000

As you can see here, label encoding uses alphabetical ordering. Hence, India has been encoded with 0, the US with 2, and Japan with 1.

## Challenges with Label Encoding

In the above scenario, the Country names do not have an order or rank. But, when label encoding is performed, the country names are ranked based on the alphabets. Due to this, there is a very high probability that the model captures the relationship between countries such as India < Japan < the US.

This is something that we do not want! So how can we overcome this obstacle? Here comes the concept of **One-Hot Encoding**.

## One-Hot Encoding

One-Hot Encoding is another popular technique for treating categorical variables. It simply creates additional features based on the number of unique values in the categorical feature. Every unique value in the category will be added as a feature.

One-Hot Encoding is the process of creating dummy variables.

In this encoding technique, each category is represented as a one-hot vector. Let's see how to implement one-hot encoding in Python:

```
1 # importing one hot encoder
2 from sklearn.preprocessing import OneHotEncoder
3 # creating one hot encoder object
4 onehotencoder = OneHotEncoder()
5 #reshape the 1-D country array to 2-D as fit_transform expects 2-D and finally fit the object
6 X = onehotencoder.fit_transform(data.Country.values.reshape(-1,1)).toarray()
7 #To add this back into the original dataframe
8 dfOneHot = pd.DataFrame(X, columns = ["Country_"+str(int(i)) for i in range(data.shape[1])])
9 df = pd.concat([data, dfOneHot], axis=1)
10 #dropping the country column
11 df= df.drop(['Country'], axis=1)
12 #printing to verify
13 print(df.head())
```

[view raw](#)

One4.py hosted with ❤ by GitHub

**Output:**

0	1	2	Age	Salary
1	0	0	44	72000
0	0	1	34	65000
0	1	0	46	98000
0	0	1	35	45000
0	1	0	23	34000

As you can see here, 3 new features are added as the country contains 3 unique values – India, Japan, and the US. In this technique, we solved the problem of ranking as each category is represented by a binary vector.

Can you see any drawbacks with this approach? Think about it before reading on.

## Challenges of One-Hot Encoding: Dummy Variable Trap

One-Hot Encoding results in a Dummy Variable Trap as the outcome of one variable can easily be predicted with the help of the remaining variables.

Dummy Variable Trap is a scenario in which variables are highly correlated to each other.

The Dummy Variable Trap leads to the problem known as **multicollinearity**. Multicollinearity occurs where there is a dependency between the independent features. Multicollinearity is a serious issue in machine learning models like [Linear Regression](#) and [Logistic Regression](#).

So, in order to overcome the problem of multicollinearity, one of the dummy variables has to be dropped. Here, I will practically demonstrate how the problem of multicollinearity is introduced after carrying out the one-hot encoding.

One of the common ways to check for multicollinearity is the Variance Inflation Factor (VIF):

- VIF=1, Very Less Multicollinearity
- VIF<5, Moderate Multicollinearity
- VIF>5, Extreme Multicollinearity (This is what we have to avoid)

Compute the VIF scores:

```
1 # Function to calculate VIF
2 def calculate_vif(data):
3     vif_df = pd.DataFrame(columns = ['Var', 'Vif'])
4     x_var_names = data.columns
5     for i in range(0, x_var_names.shape[0]):
6         y = data[x_var_names[i]]
7         x = data[x_var_names.drop([x_var_names[i]])]
8         r_squared = sm.OLS(y,x).fit().rsquared
9         vif = round(1/(1-r_squared),2)
10        vif_df.loc[i] = [x_var_names[i], vif]
11    return vif_df.sort_values(by = 'Vif', axis = 0, ascending=False, inplace=False)
12
13 X=df.drop(['Salary'],axis=1)
14 calculate_vif(X)
```

[view raw](#)

One5.py hosted with ❤ by GitHub

**Output:**

	Var	Vif
2	2.0	9.94
0	0.0	9.84
1	1.0	7.80
3	3.0	1.23

From the output, we can see that the dummy variables which are created using one-hot encoding have VIF above 5. We have a multicollinearity problem.

Now, let us drop one of the dummy variables to solve the multicollinearity issue:

```
1 df = df.drop(df.columns[[0]], axis=1)
2 calculate_vif(df)
```

[view raw](#)

one6.py hosted with ❤ by GitHub

**Output:**

	Var	Vif
2	3.0	2.60
1	2.0	1.91
0	1.0	1.69

Wow! VIF has decreased. We solved the problem of multicollinearity. Now, the dataset is ready for building the model.

I would recommend you to go through [Going Deeper into Regression Analysis with Assumptions, Plots & Solutions](#) for understanding the assumptions of linear regression.

We have seen two different techniques – Label and One-Hot Encoding for handling categorical variables. In the next section, I will touch upon when to prefer label encoding vs. One-Hot Encoding.

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

This question generally depends on your dataset and the model which you wish to apply. But still, a few points to note before choosing the right encoding technique for your model:

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

## End Notes

As quoted by Jeff Hawkins:

"The key to artificial intelligence has always been the representation."

Representation has been the key for developers and new techniques are emerging now and then to better represent the data and improve the accuracy and learning of our model.

I encourage you to go through the below course to become a machine learning expert:

- [Applied Machine Learning – Beginner to Professional](#)

---

Article Url - <https://www.analyticsvidhya.com/blog/2020/03/one-hot-encoding-vs-label-encoding-using-scikit-learn/>



### [Alakh Sethi](#)

Aspiring Data Scientist with a passion to play and wrangle with data and get insights from it to help the community know the upcoming trends and products for their better future. With an ambition to develop product used by millions which makes their life easier and better.