

# Feature Engineering

## Agenda

In this session, we will cover the following concepts with the help of a business use case:

- Features
- Feature engineering and its methods
  - Applying domain expertise
  - Combining features
  - Encoding of categorical variables

## Features

- In most cases, machine learning algorithms are fed a collection of numeric examples that are stacked on top of each other to form a two-dimensional feature matrix.
- This matrix has one **example** per row and one **function** per column.
- Understanding the task of these features is critical when doing some kind of learning on the dataset.
- Since the algorithms create their own internal transformations, deep learning features are normally simple.
- This method necessitates a vast volume of data and sacrifices interpretability; but, in image processing and natural language processing applications, these trade offs are often worthwhile.
- Feature engineering is needed for the majority of other use cases in order to translate data into a machine learning-ready format.
- The features chosen are critical for both the data's interpretability and the model's performance.
- A precise machine learning framework cannot be built without feature engineering.

## What Is Feature Engineering?

- Feature Engineering repurposes already modified features to build new ones, making it possible for any Machine Learning algorithm to comprehend and master any pattern.
- It is the process of **creating new features** from raw data using domain knowledge to make the machine learning algorithms work.
- It is one of the fundamental concepts of machine learning and it is both difficult and expensive.
- Every new resource must enhance the performance in some way; otherwise, it would have the opposite effect, worsening the final outcome. When this happens, we must use feature selection instead of feature engineering.
- According to Andrew Ng, "Coming up with features is difficult, time-consuming, requires expert knowledge."
- According to Pedro Domingos, feature engineering is the art of coming up with new features with more predictive power using:
  - Experience
  - Domain expertise
  - Empirical processes
- The following are some of the most important components of feature engineering:
  1. Imputation
  2. Handling Outliers
  3. Binning
  4. One-Hot Encoding
  5. Grouping Operations
  6. Scaling

## Most Important Techniques in Feature Engineering

- Applying domain expertise
- Combining features
- Encoding of categorical variables

### Applying Domain Expertise

Domain expertise or domain knowledge is nothing but expertise in a particular field, such as Education, Healthcare, Consumer Goods, and Retail. A domain expert is someone who is not related to the technology aspect but has in-depth knowledge about the particular industry, how it is shaping up, the trends, and the things that might impact the industry.

For example, you are called in to develop a particular application for a consumer goods company, specifically an apparel and footwear one. The application that you build has to align to the industry and its various facets, and you being a technology expert wouldn't know much about it. This is where a domain expert will come in and explain how that industry works and what would be the best way to have the application built.

### Combining Features

This is one of the most crucial part of feature engineering that involves analyzing the data set and understanding the variable(s) to combine one or more features to create new ones. This process involves intuitive decision-making and a small amount of research about the domain of the data set.

Combining features helps to increase the performance of the model which is fitted on the data.

### Encoding of Categorical Variables

Let's first see what encoding is.

#### What Is Encoding?

Encoding is used to transform the categorical variable into numerical features.

For example, we have an attribute called gender in a data set where values are male and female. In this case, the numerical encoded version of the values will be 1 for male, 0 for female, or vice versa.

Since different kinds of categorical variables capture different amount of information, we need different techniques to encode them.

### Label Encoding

Label encoding is a handy technique to encode categorical variables. However, such encoded nominal variables might end up being misinterpreted as ordinal.

#### Encoding Can Be Performed on the Following Types of Variables:

##### Nominal Variables

Consider the following three categorical variables and their values:

**Color:** Blue, Green, Red, Yellow

**Educational Qualification:** Primary School, Secondary School, Graduate, Post-Graduate, PhD

**Salary Bracket:** 0-50,000, 50,001-100,000, 100,001-150,000, 150,001-200,000

Although all three of them are categorical variables, they are different in the amount of information they convey. Let's look at them one by one.

Color conveys blue is different from red. That's all. The value of the variable is not meant to capture any relative difference among the values. Such variables are called **Nominal Variables**.

### Ordinal Variables

Now consider educational qualifications.

The value "graduate" does not only convey that it is different from the value say "Primary school", it also implies that it is more in terms of qualification than "Primary School".

Such variables are called **Ordinal Variables** because they convey a sense of order.

In our example, Primary School < Secondary School < Graduate < Post-Graduate < PhD. (in terms of qualification).

### Interval Variables

The third variable, salary bracket is similar to educational qualification by conveying order (a person earning 50,001-100,000 earns more salary than 0-50,000).

However, here, apart from knowing the order, we also know the interval between the values.

Here we can say that the averages of each of the values are separated by 50,000. Such variables are called **Interval variables**.

## Techniques Used for Encoding Variables

There are two types of broadly used algorithms which perform the task of encoding of variables.

### There are few libraries required to perform encoding variables:

- Pandas - It helps to retrieve datasets, handle missing data and perform data wrangling.
- NumPy - It helps to perform numerical operations in the dataset.
- sklearn.preprocessing - It helps in data transformation.

```
In [ ]: 1 #Select the cell and click on run icon to import libraries.  
2  
3 import pandas as pd  
4 import numpy as np  
5  
6 # Import Label encoder  
7 from sklearn import preprocessing
```

#### Instruction:

Download the **Iris.csv** dataset file from Course Resources and upload it in the lab using the Up arrow as shown below on the View Tab

#### Note 1:

- The **iris\_df** is a dataframe that stores data imported from a CSV file in rows and columns format.
- The **head()** function helps to view the first few data present in the **iris\_df** dataframe.

```
In [ ]: 1 #Select the cell and click on run icon to retrieve and view the dataset.  
2 # Import dataset  
3  
4 iris_df = pd.read_csv('Iris.csv')  
5 iris_df.head()
```

Out[4]:

	sepal_length	sepal_width	petal_length	petal_width	species
0	5.1	3.5	1.4	0.2	setosa
1	4.9	3.0	1.4	0.2	setosa
2	4.7	3.2	1.3	0.2	setosa
3	4.6	3.1	1.5	0.2	setosa
4	5.0	3.6	1.4	0.2	setosa

**Observations from the above output:**

View few records of Iris dataset

**Note 2:**

The **info()** function helps to understand the dataset, the column name, total null values, and data type.

```
In [ ]: 1 #Select the cell and click on run icon  
2 iris_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 150 entries, 0 to 149  
Data columns (total 5 columns):  
 #   Column      Non-Null Count  Dtype     
 ---  -----      -----          -----  
 0   sepal_length 150 non-null    float64  
 1   sepal_width  150 non-null    float64  
 2   petal_length 150 non-null    float64  
 3   petal_width  150 non-null    float64  
 4   species      150 non-null    object    
dtypes: float64(4), object(1)  
memory usage: 6.0+ KB
```

**Observations from the above output:**

The variable **species** is a categorical Dtype '**object**'.

**Note 3:**

The **LabelEncoder()** function is used to convert categorical variables into numerical. According to our dataset, variable **species** is categorical, so convert **species** into numerical type.

```
In [ ]: 1 #Select the cell and click on run icon to define LabelEncoder  
2 label_encoder = preprocessing.LabelEncoder()
```

**Note 4:**

The **unique()** function is used to identify distinct rows present in the **iris\_df** dataframe.

```
In [ ]: 1 #Select the cell and click on run icon  
2 iris_df['species'].unique()
```

```
Out[9]: array(['setosa', 'versicolor', 'virginica'], dtype=object)
```

#### Observations from the above output:

There are three categories of variable `species` such as 'setosa', 'versicolor', and 'virginica' to be encoded.

#### Note 5:

- The `fit_transform()` method calculates the mean and variance of each feature and transforms all the features using the respective mean and variance.
- The `head()` function helps to view the first few data present in the `iris_df` dataframe.

```
In [ ]: 1 #Select the cell and click on run icon  
2 iris_df['species']= label_encoder.fit_transform(iris_df['species'])  
3 iris_df.head()
```

```
Out[10]:   sepal_length  sepal_width  petal_length  petal_width  species
```

0	5.1	3.5	1.4	0.2	0
1	4.9	3.0	1.4	0.2	0
2	4.7	3.2	1.3	0.2	0
3	4.6	3.1	1.5	0.2	0
4	5.0	3.6	1.4	0.2	0

#### Observations from the above output:

The variable `species` is converted into numerical variable as 0's and 1's.

```
In [ ]: 1 #Select the cell and click on run icon  
2  
3 iris_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 150 entries, 0 to 149  
Data columns (total 5 columns):  
 #   Column      Non-Null Count  Dtype     
---  --          --          --  
 0   sepal_length  150 non-null    float64  
 1   sepal_width   150 non-null    float64  
 2   petal_length  150 non-null    float64  
 3   petal_width   150 non-null    float64  
 4   species       150 non-null    int64  
dtypes: float64(4), int64(1)  
memory usage: 6.0 KB
```

#### Observations from the above output:

The variable `species` Dtype is converted to 'int64'.

Using LabelEncoder, we can convert a categorical variable into a numerical variable.

## One Hot Encoding

A data set with more dimensions requires more parameters for the model to understand, and that means more rows to reliably learn those parameters.

The effect of using One Hot Encoder is the addition of a number of columns (dimensions).

If the number of rows in the data set is fixed, the addition of extra dimensions without adding more information for the models to learn from can have a detrimental effect on the eventual model accuracy.

If the number of rows in the dataset is fixed and the addition of extra dimension without adding more information in the model would affect the model accuracy.

In the below example, you can see that category of Variable X is encoded into separate columns such as Variable X\_Blue, Variable X\_Yellow, Variable X\_Red.

Encoding of categorical variables

One Hot Encoding:

Index	Variable X
1	Blue
2	Red
3	Yellow
4	Red
5	Red
6	Blue
7	Yellow

Index	Variable X_Blue	Variable X_Red	Variable X_Yellow
1	1	0	0
2	0	1	0
3	0	0	1
4	0	1	0
5	0	1	0
6	1	0	0
7	0	0	1

There are few libraries required to perform **One Hot Encoding**:

- **datasets**: It helps to load the datasets from sklearn library.
- **OneHotEncoder**: It helps to encode categorical variable.

In [ ]:

```
1 #Select the cell and click on run icon
2 from sklearn import datasets
3 from sklearn.preprocessing import OneHotEncoder
```

### Note 6:

The given code helps to load the **Iris dataset** and create a dataframe **iris\_data**, and **y** variable consists of target values.

```
In [ ]: 1 #Select the cell and click on run icon
        2 iris_data = datasets.load_iris()
        3
        4 iris_data = pd.DataFrame(data=np.c_[iris_data["data"], iris_data["target"]], columns=
        5 y = iris_data.target.values
```

### Note 7:

The `head()` function is used to view the first few data present in the `iris_data` dataframe.

```
In [ ]: 1 #Select the cell and click on run icon  
2 iris.data.head()
```

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)	target
0	5.1	3.5	1.4	0.2	0.0
1	4.9	3.0	1.4	0.2	0.0
2	4.7	3.2	1.3	0.2	0.0
3	4.6	3.1	1.5	0.2	0.0
4	5.0	3.6	1.4	0.2	0.0

#### **Observations from the above output:**

There are four dependent variable and one independent variable present in the `iris_data` dataframe.

```
In [ ]: 1 # Select the cell and click on run icon to view the target values  
2 y
```

#### **Observations from the above output:**

The variable `target` is converted into an array.

### Note 8:

- In OneHotEncoder, the categories set as 'auto' determines the categories automatically.
  - The **reshape()** function reshapes an array without modifying its data.
  - The **toarray()** function helps to represent the data in array format.

```
In [ ]: 1 #Select the cell and click on run icon
2 onehotencoder = OneHotEncoder(categories='auto')
3 y = onehotencoder.fit_transform(y.reshape(-1,1))
4 print(y.toarray())
[[1. 0. 0.]
 [1. 0. 0.]
 [1. 0. 0.]
 [1. 0. 0.]
 [1. 0. 0.]
 [1. 0. 0.]
 [1. 0. 0.]
 [1. 0. 0.]
 [1. 0. 0.]
 [1. 0. 0.]
 [1. 0. 0.]
 [1. 0. 0.]
 [1. 0. 0.]
 [1. 0. 0.]
 [1. 0. 0.]
 [1. 0. 0.]
 [1. 0. 0.]
 [1. 0. 0.]
 [1. 0. 0.]
 [1. 0. 0.]
```

**Observations from the above output:**

With the help of above output, you can see that variable `target` is converted into numerical variable as 0's and 1's into three separate categories .

**Note 9:**

The `get_dummies()` function is used for data manipulation and converts categorical variable to numerical variables .

```
In [ ]: 1 #Select the cell and click on run icon
2 pd.get_dummies(iris_data.target).head()
```

```
Out[8]: 0.0 1.0 2.0
0   1   0   0
1   1   0   0
2   1   0   0
3   1   0   0
4   1   0   0
```

**Observations from the above output:**

The target values are converted into dummy values in three different columns.

## Hashing

Hashing is a technique that combines more than one category of a categorical variable into one single category.

There are three advantages in one-hot encoding or set-of word model:

- It's much simpler to code

- It doesn't require any dictionaries to prepare
- It's friendly to online learning where you can train on a data set that doesn't fit in memory because you need to see each example only once. One-hot encoding will not work well with online learning because to prepare dictionaries you need to see the whole data set first

## Feature Hashing

Feature hashing is a important technique for handling sparse and high-dimensional features in machine learning.

- It is fast, simple, memory-efficient, and well-suited to online learning sceanrios.
- It converts unique tokens into integers.
- It operates on the exact strings that you provide as input and does not perform any linguistic analysis or preprocessing.

## Techniques to merge categories:

### Business Logic

Consider a column in the dataset corresponds to "zip codes". There are 182 zip codes in New York state and it is impractical to use each zip code as a separate category. So, to tackle this situation we can merge the zip codes according to localities. This helps to reduce the number of categories and results in meaningful aggregation of zip code.

### Frequency

- It is not possible to apply business logic every time. In such cases, perform hashing using the frequency of occurrence.
- To combine levels using their frequency, we first look at the frequency distribution of each level and combine levels having frequency say less than 5% of total observation (can be changed based on distribution).
- This is an effective method to deal with rare levels.
- We can also combine levels by considering the response rate of each level. We can simply combine levels having similar response rates into the same group.

Now we'll see all the concepts explained with the help of a business use case.

## There is a libraries required to perform encoding variables:

- FeatureHasher - Feature hashing can be employed in document classification and helps to extract important features from large text corpus.

```
In [ ]: 1 #Select the cell and click on run icon
          2 from sklearn.feature_extraction import FeatureHasher
```

### Instruction:

Download the `vgsales.csv` dataset file from Course Resources and upload it in the lab using the Up arrow as shown below on the View Tab

### Note 10:

- The `game_df` is dataframe to store the data imported from the csv as rows and columns table format.
- The encoding is specified as 'UTF-8' as the file type converion of the csv file.

- The **head()** function helps to view first few data present in the dataframe.

```
In [ ]: 1 #Select the cell and click on run icon
2 game_df = pd.read_csv("vgsales.csv", encoding="utf-8")
3 game_df.head()
```

Out[20]:

	Rank	Name	Platform	Year	Genre	Publisher	NA_Sales	EU_Sales	JP_Sales	Other_Sales	GI
0	1	Wii Sports	Wii	2006.0	Sports	Nintendo	41.49	29.02	3.77	8.46	
1	2	Super Mario Bros.	NES	1985.0	Platform	Nintendo	29.08	3.58	6.81	0.77	
2	3	Mario Kart Wii	Wii	2008.0	Racing	Nintendo	15.85	12.88	3.79	3.31	
3	4	Wii Sports Resort	Wii	2009.0	Sports	Nintendo	15.75	11.01	3.28	2.96	
4	5	Pokemon Red/Pokemon Blue	GB	1996.0	Role-Playing	Nintendo	11.27	8.89	10.22	1.00	

### Note 11:

The **columns** provides the name of the columns present in the `game_df` dataframe.

```
In [ ]: 1 #Select the cell and click on run icon
2 game_df.columns
```

Out[21]:

```
Index(['Rank', 'Name', 'Platform', 'Year', 'Genre', 'Publisher', 'NA_Sales',
       'EU_Sales', 'JP_Sales', 'Other_Sales', 'Global_Sales'],
      dtype='object')
```

### Observations from the above output:

The column names such as Rank, Name, Platform, Year, Genre, Publisher, NA\_Sales, EU\_Sales, JP\_Sales, Other\_Sales, and Global\_Sales present in the `game_df` dataframe.

### Note 12:

The **shape** represents the number of elements in each dimension and returns a tuple with each index having the number of corresponding elements.

```
In [ ]: 1 #Select the cell and click on run icon
2 game_df.shape
```

Out[22]:

```
(16598, 11)
```

### Observations we can draw from the above output:

There are 16598 rows and 11 columns present in the `game_df` dataframe.

### Note 13:

- The **iloc** is integer-index based that helps to extract the required number of rows or columns by specifying the integer index value present in the dataframe.
- Example: `iloc[1:7]` represents six rows, as it subtracts one from seven.

```
In [ ]: 1 #Select the cell and click on run icon  
2 game_df[['Name', 'Platform', 'Year', 'Genre', 'Publisher']].iloc[1:7]
```

Out[23]:

	Name	Platform	Year	Genre	Publisher
1	Super Mario Bros.	NES	1985.0	Platform	Nintendo
2	Mario Kart Wii	Wii	2008.0	Racing	Nintendo
3	Wii Sports Resort	Wii	2009.0	Sports	Nintendo
4	Pokemon Red/Pokemon Blue	GB	1996.0	Role-Playing	Nintendo
5	Tetris	GB	1989.0	Puzzle	Nintendo
6	New Super Mario Bros.	DS	2006.0	Platform	Nintendo

**Observations we can draw from the above output:**

There are six rows present in the `game_df` dataframe as specified in the `iloc`.

**Note 14:**

- The `unique()` function extracts unique values present in a particular variable or dataframe.
- The `len()` function calculates the total number of elements present of a variable or dataframe.
- The `print()` function is used to display the result.

```
In [ ]: 1 #Select the cell and click on run icon  
2 u_genres = np.unique(game_df["Genre"])  
3 print("Total game genres:", len(u_genres))  
4 print(u_genres)
```

Total game genres: 12  
['Action' 'Adventure' 'Fighting' 'Misc' 'Platform' 'Puzzle' 'Racing'  
'Role-Playing' 'Shooter' 'Simulation' 'Sports' 'Strategy']

**Observations we can draw from the above output:**

The '`u_genres`' variable stores unique features such as Action, Adventure, Fighting, Misc, Platform, Puzzle, Racing, Role-Playing, Shooter, Simulation, Sports, and Strategy and the total count is 12 present in the '`genre`' column of `game_df` dataframe.

**Note 15:**

- The `FeatureHasher()` function extracts only important features from the dataset.

```
In [ ]: 1 #Select the cell and click on run icon
2 fh = FeatureHasher(n_features=6, input_type='string')
3 hashed_features = fh.fit_transform(game_df[ "Genre"])
4 hashed_features = hashed_features.toarray()
5 new_game_df = pd.concat([game_df[['Name', 'Genre']], pd.DataFrame(hashed_features)], axis=1)
6
7 new_game_df.head()
```

Out[25]:

	Name	Genre	0	1	2	3	4	5
0	Wii Sports	Sports	-2.0	2.0	0.0	-2.0	0.0	0.0
1	Super Mario Bros.	Platform	0.0	2.0	2.0	-1.0	1.0	0.0
2	Mario Kart Wii	Racing	-1.0	0.0	0.0	0.0	0.0	-1.0
3	Wii Sports Resort	Sports	-2.0	2.0	0.0	-2.0	0.0	0.0
4	Pokemon Red/Pokemon Blue	Role-Playing	-1.0	1.0	2.0	0.0	1.0	-1.0

## Grouping Operations

There are two types of grouping operations:

- 1. Categorical Grouping:** Grouping based on aggregate function using lambda or grouping using a pivot table.
- 2. Numeric Grouping:** Grouping numerical columns using mean and sum functions.

## Problem Statements:

You work in HR analytics. You have been given a task to predict the employees who are going to leave the organization, so that their replacement process can be started within the available time frame.

The dataset is created by IBM data scientists to uncover the factors that lead to employee attrition and explore important questions such as 'show me a breakdown of distance from home by job role and attrition' or 'compare average monthly income by education and attrition'.

The name of the dataset attributes are self explanatory. Some required attributes are explained below.

- Education
  - Below College
  - College
  - Bachelor
  - Master
  - Doctor
- EnvironmentSatisfaction
  - Low
  - Medium
  - High
  - Very High
- JobInvolvement
  - Low
  - Medium
  - High
  - Very High
- JobSatisfaction
  - Low
  - Medium
  - High

- Very High
- PerformanceRating
  - Low
  - Good
  - Excellent
  - Outstanding
- RelationshipSatisfaction
  - Low
  - Medium
  - High
  - Very High
- WorkLifeBalance
  - Bad
  - Good
  - Better
  - Best

## There are few libraries required to perform encoding variables:

- pandas - It helps to retrieve datasets, handle missing data and in data wrangling.
- Numpy - It helps to perform numerical operations in the dataset.
- seaborn and matplotlib- It helps to data visualization.
- warnings - It helps to segregate the warnings.

## Import libraries

```
In [ ]: 1 #Select the cell and click on run icon
2 import numpy as np
3 import pandas as pd
4 import seaborn as sns
5
6 import matplotlib.pyplot as plt
7 %matplotlib inline
8
9 import warnings
10 warnings.filterwarnings('ignore')
```

### Instruction:

Download the `dataset.csv` dataset file from Course Resources and upload it in the lab using Up arrow shown below View Tab

## Read the dataset

### Note 16:

- The `df` is dataframe to store the data imported from the csv as rows and columns table format.
- The `head()` function helps to view first few data present in the `df` dataframe.

```
In [ ]: 1 #Select the cell and click on run icon  
2 df = pd.read_csv("dataset.csv")  
3 df.head()
```

Out[5]:

Index	Age	Attrition	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EducationField	EmployeeCount	EnvironmentSatisfaction	HourlyRate	JobInvolvement	JobLevel	JobRole	JobSatisfaction	MaritalStatus	OverTime	PercentSalaryHike	RelationshipSatisfaction	StandardHours	TotalWorkingYears	WorkLifeBalance	YearsAtCompany	YearsInCurrentRole	YearsOnSite	YearsWithCurrManager
0	41	Yes	Travel_Rarely	1102	Sales	1	2	Life Sciences	10000	Low	40	High	Manager	Manager	Low	Married	No	0.5	Low	40	10	Low	10	10	10	
1	49	No	Travel_Frequently	279	Research & Development	8	1	Life Sciences	10000	Low	40	High	Manager	Manager	Low	Married	No	0.5	Low	40	10	Low	10	10	10	
2	37	Yes	Travel_Rarely	1373	Research & Development	2	2	Other	10000	Low	40	High	Manager	Manager	Low	Married	No	0.5	Low	40	10	Low	10	10	10	
3	33	No	Travel_Frequently	1392	Research & Development	3	4	Life Sciences	10000	Low	40	High	Manager	Manager	Low	Married	No	0.5	Low	40	10	Low	10	10	10	
4	27	No	Travel_Rarely	591	Research & Development	2	1	Medical	10000	Low	40	High	Manager	Manager	Low	Married	No	0.5	Low	40	10	Low	10	10	10	

5 rows × 35 columns

## Note 17:

The **columns** provides the name of the columns present in the `df` dataframe.

```
In [ ]: 1 #Select the cell and click on run icon  
2 df.columns
```

```
Out[6]: Index(['Age', 'Attrition', 'BusinessTravel', 'DailyRate', 'Department',
   'DistanceFromHome', 'Education', 'EducationField', 'EmployeeCount',
   'EmployeeNumber', 'EnvironmentSatisfaction', 'Gender', 'HourlyRate',
   'JobInvolvement', 'JobLevel', 'JobRole', 'JobSatisfaction',
   'MaritalStatus', 'MonthlyIncome', 'MonthlyRate', 'NumCompaniesWorked',
   'Over18', 'OverTime', 'PercentSalaryHike', 'PerformanceRating',
   'RelationshipSatisfaction', 'StandardHours', 'StockOptionLevel',
   'TotalWorkingYears', 'TrainingTimesLastYear', 'WorkLifeBalance',
   'YearsAtCompany', 'YearsInCurrentRole', 'YearsSinceLastPromotion',
   'YearsWithCurrManager'],
  dtype='object')
```

#### **Observations we can draw from the above output:**

It denotes the name of the columns present in the `df` dataframe.

## Note 18:

The **shape** represents the number of elements in each dimension and returns a tuple with each index having the number of corresponding elements.

```
In [ ]: 1 #Select the cell and click on run icon  
2 df.shape
```

Out[7]: (1470, 35)

### **Observations we can draw from the above output:**

There are 1470 rows and 35 columns present in the `df` dataframe.

## Check if there are any null values

```
In [ ]: 1 #Select the cell and click on run icon  
2 df.isna().sum()
```

```
Out[8]: Age          0  
Attrition      0  
BusinessTravel  0  
DailyRate       0  
Department     0  
DistanceFromHome 0  
Education       0  
EducationField   0  
EmployeeCount    0  
EmployeeNumber   0  
EnvironmentSatisfaction 0  
Gender          0  
HourlyRate      0  
JobInvolvement   0  
JobLevel         0  
JobRole          0  
JobSatisfaction  0  
MaritalStatus    0  
MonthlyIncome    0  
NumCompaniesLastYear 0  
OverTime         0  
PercentHrWorked  0  
StockOptionLevel 0  
TotalWorkingYears 0  
WorkLifeBalance  0
```

### Observations:

```
There are no null values present in the df dataframe.
```

## Now, print the count of values for each column

```
In [ ]: 1 for i in df.columns:  
2     print (i , ":" , df[i].value_counts())  
3     print ("_ *40)  
4     print ("_ *40)
```

```
Age : 35    78  
34    77  
31    69  
36    69  
29    68  
32    61  
30    60  
33    58  
38    58  
40    57  
37    50  
27    48  
28    48  
42    46  
39    42  
45    41  
41    40  
26    39  
46    33  
..    ..
```

### Observations we can draw from the above output:

It is an imbalanced data as value count in attrition column is imbalanced.

Over18 is Y (Yes) across all the employees and it is not beneficial for further use.

`StandardHours` is 80 which is common for all and it is not beneficial for further use.

`EmployeeCount` , `EmployeeNumber` are unique which are also not going to help us predict our end result.

**Now, print the unique values in each column to understand the above output more clearly**

```
In [ ]: 1 for i in df.columns:  
2     print (i , ":" , df[i].unique())  
3     print (" _ _ _ _ _")  
4     print (" _ _ _ _ _")  
  
Age : [41 49 37 33 27 32 59 30 38 36 35 29 31 34 28 22 53 24 21 42 44 46 39 43  
50 26 48 55 45 56 23 51 40 54 58 20 25 19 57 52 47 18 60]  
- - - - -  
- - - - -  
- - - - -  
- - - - -  
Attrition : ['Yes' 'No']  
- - - - -  
- - - - -  
- - - - -  
- - - - -  
BusinessTravel : ['Travel_Rarely' 'Travel_Frequently' 'Non-Travel']  
- - - - -  
- - - - -  
- - - - -  
- - - - -  
DailyRate : [1102 279 1373 1392 591 1005 1324 1358 216 1299 809 153 670 1  
103 1389 334 1123 1219 371 673 1218 419 391 699 1282 1125 691  
477 705 924 1459 125 895 813 1273 869 890 852 1141 464 1240  
1257 661 721 1260 1265 422 1211 1222 626 1224 1422 1227 1142 515
```

### Note 19:

The **drop()** function helps to drop the columns which are not important features in the dataset.

```
In [ ]: 1 df=df.drop(['Over18','EmployeeNumber','EmployeeCount','StandardHours'],axis=1)
```

## Note 20:

After dropping, check the number of columns and rows present in the dataset.

```
In [ ]: 1 df.shape
```

### Note 21:

The `columns` provides the name of the columns present in the `df` dataframe.

```
In [ ]: 1 df.columns
```

```
Out[12]: Index(['Age', 'Attrition', 'BusinessTravel', 'DailyRate', 'Department',  
       'DistanceFromHome', 'Education', 'EducationField',  
       'EnvironmentSatisfaction', 'Gender', 'HourlyRate', 'JobInvolvement',  
       'JobLevel', 'JobRole', 'JobSatisfaction', 'MaritalStatus',  
       'MonthlyIncome', 'MonthlyRate', 'NumCompaniesWorked', 'OverTime',  
       'PercentSalaryHike', 'PerformanceRating', 'RelationshipSatisfaction',  
       'StockOptionLevel', 'TotalWorkingYears', 'TrainingTimesLastYear',  
       'WorkLifeBalance', 'YearsAtCompany', 'YearsInCurrentRole',  
       'YearsSinceLastPromotion', 'YearsWithCurrManager'],  
      dtype='object')
```

**Observations from the above output:**

We can see a lot columns which deal with satisfaction, let's observe them closely.

```
In [ ]: 1 df[['RelationshipSatisfaction', 'JobSatisfaction', 'EnvironmentSatisfaction', 'JobInvolvement']]
```

```
Out[13]:
```

	RelationshipSatisfaction	JobSatisfaction	EnvironmentSatisfaction	JobInvolvement
0	1	4	2	3
1	4	2	3	2
2	2	3	4	2
3	3	3	4	3
4	4	2	1	3
...	...	...	...	...
1465	3	4	3	4
1466	1	1	4	2
1467	2	2	2	4
1468	4	2	4	2
1469	1	3	2	4

1470 rows × 4 columns

**Observations from the above output:**

All the satisfactions are measured from 1 to 4 in incremental order, along with WorkLifeBalance.

Higher the value, higher the satisfaction.

So, we can combine all the columns that convey satisfaction detail and make them one.

**Note 22:**

Let us now calculate the mean of all the types of satisfaction. Based on a condition if the mean value is greater than 2.35 then it returns one else zero

```
In [ ]: 1 df['TotalSatisfaction_mean'] = (df['RelationshipSatisfaction'] + df['EnvironmentSat']
2                                + df['JobSatisfaction'] + df['JobInvolvement']) / 2
3
4 def Satif(df) :
5     if df['TotalSatisfaction_mean'] > 2.35 :
6         return 1
7     else :
8         return 0
9
10
11 df['Satif'] = df.apply(lambda df:Satif(df) ,axis = 1)
12 df['Satif']
```

```
Out[18]: 0      0
1      1
2      1
3      1
4      1
..
1465    1
1466    0
1467    1
1468    1
1469    1
Name: Satif, Length: 1470, dtype: int64
```

```
In [ ]: 1 df.shape
```

```
Out[29]: (1470, 33)
```

Let's keep all the columns for now.

### Note 23:

Let us now create a separate column for job satisfaction.

```
In [ ]: 1 df['JobSatisf_mean'] = (df['JobSatisfaction'] + df['JobInvolvement']) / 2
```

```
In [ ]: 1 df.shape
```

```
Out[31]: (1470, 34)
```

### Note 24:

We can understand people who switch companies frequently have more tendency to leave a company. So let's create a new column called MovingPeople .

```
In [ ]: 1 def MovingPeople(df) :
2     if df['NumCompaniesWorked'] > 4:
3         return 1
4     else:
5         return 0
6 df['MovingPeople'] = df.apply(lambda df:MovingPeople(df), axis = 1)
7 df['MovingPeople']
```

```
Out[35]: 0      1
1      0
2      1
3      0
4      1
..
1465    0
1466    0
1467    0
1468    0
1469    0
Name: MovingPeople, Length: 1470, dtype: int64
```

```
In [ ]: 1 df.shape
```

```
Out[21]: (1470, 34)
```

### Note 25:

Create a column using `DistanceFromHome` column

```
In [ ]: 1 def LongDis(df) :
2     if df['DistanceFromHome'] > 11:
3         return 1
4     else :
5         return 0
6 df['LongDis'] = df.apply(lambda df:LongDis(df) ,axis = 1)
7 df['LongDis']
```

```
Out[28]: 0      0
1      0
2      0
3      0
4      0
..
1465    1
1466    0
1467    0
1468    0
1469    0
Name: LongDis, Length: 1470, dtype: int64
```

### Note 26:

Create a column using `TrainingTimesLastYear` column

```
In [ ]: 1 def MiddleTraining(df) :
2     if df['TrainingTimesLastYear'] >= 3 and df['TrainingTimesLastYear'] <= 6:
3         return 1
4     else:
5         return 0
6 df['MiddleTraining'] = df.apply(lambda df:MiddleTraining(df) ,axis = 1)
```

### Note 27:

Create a column to view number of years worked in each company

```
In [ ]: 1 df['Time_in_each_comp'] = (df['Age'] - 20) / ((df['NumCompaniesWorked'] + 1)  
2 df['Time_in_each_comp']
```

```
Out[34]: 0      2.333333  
1      14.500000  
2      2.428571  
3      6.500000  
4      0.700000  
...  
1465    3.200000  
1466    3.800000  
1467    3.500000  
1468    9.666667  
1469    4.666667  
Name: Time_in_each_comp, Length: 1470, dtype: float64
```

```
In [ ]: 1 df.shape
```

```
Out[31]: (1470, 37)
```

### Note 28:

Let us now understand the columns that are categorical and numerical in the `df` dataframe.

```
In [ ]: 1 numeric_df=df.select_dtypes(include=[np.number])  
2  
3 categoric_df=df.select_dtypes(exclude=[np.number])
```

```
In [ ]: 1 numericcol=numeric_df.columns.tolist()  
2 categorycol=categoric_df.columns.tolist()  
3  
4 print ("Category : ",categorycol)  
5 print ("\n Numeric : ",numericcol)
```

```
Category : ['Attrition', 'BusinessTravel', 'Department', 'EducationField', 'Gender', 'JobRole', 'MaritalStatus', 'OverTime']
```

```
Numeric : ['Age', 'DailyRate', 'DistanceFromHome', 'Education', 'EnvironmentSatisfaction', 'HourlyRate', 'JobInvolvement', 'JobLevel', 'JobSatisfaction', 'MonthlyIncome', 'MonthRate', 'NumCompaniesWorked', 'PercentSalaryHike', 'PerformanceRating', 'RelationshipSatisfaction', 'StockOptionLevel', 'TotalWorkingYears', 'TrainingTimesLastYear', 'WorkLifeBalance', 'YearsAtCompany', 'YearsInCurrentRole', 'YearsSinceLastPromotion', 'YearsWithCurrManager', 'TotalSatisfaction_mean', 'Satif', 'MovingPeople', 'LongDis', 'MiddleTraining', 'Time_in_each_comp']
```

### Observations from the above output:

'DailyRate', 'HourlyRate' can be dropped as we have monthly rate

Attributes 'DistanceFromHome', 'EnvironmentSatisfaction', 'HourlyRate', 'JobInvolvement', 'JobSatisfaction', 'NumCompaniesWorked', 'RelationshipSatisfaction', and 'TrainingTimesLastYear' can be removed as these are considered in newly created column attributes

### Note 29:

Let us now drop unnecessary columns present in the `df` dataframe.

```
In [ ]: 1 df=df.drop(['DailyRate', 'DistanceFromHome', 'EnvironmentSatisfaction',
2 'HourlyRate', 'JobInvolvement',
3 'JobSatisfaction', 'NumCompaniesWorked',
4 'RelationshipSatisfaction', 'TrainingTimesLastYear'],axis=1)
```

```
In [ ]: 1 df.shape
```

Out[37]: (1470, 28)

### Note 30:

After extracting the important features from the dataframe. Now, we can convert the categorical features into numerical using **one hot encoding technique**.

```
In [ ]: 1 data = pd.get_dummies(df, columns=categorycol, drop_first=True)
2 print(data.columns)
3 print(data.shape)
```

```
Index(['Age', 'Education', 'JobLevel', 'MonthlyIncome', 'MonthlyRate',
       'PercentSalaryHike', 'PerformanceRating', 'StockOptionLevel',
       'TotalWorkingYears', 'WorkLifeBalance', 'YearsAtCompany',
       'YearsInCurrentRole', 'YearsSinceLastPromotion', 'YearsWithCurrManager',
       'TotalSatisfaction_mean', 'Satif', 'MovingPeople', 'LongDis',
       'MiddleTraining', 'Time_in_each_comp', 'Attrition_Yes',
       'BusinessTravel_Travel_Frequently', 'BusinessTravel_Travel_Rarely',
       'Department_Research & Development', 'Department_Sales',
       'EducationField_Life Sciences', 'EducationField_Marketing',
       'EducationField_Medical', 'EducationField_Other',
       'EducationField_Technical Degree', 'Gender_Male',
       'JobRole_Human Resources', 'JobRole_Laboratory Technician',
       'JobRole_Manager', 'JobRole_Manufacturing Director',
       'JobRole_Research Director', 'JobRole_Research Scientist',
       'JobRole_Sales Executive', 'JobRole_Sales Representative',
       'MaritalStatus_Married', 'MaritalStatus_Single', 'OverTime_Yes'],
      dtype='object')
(1470, 42)
```

```
In [ ]: 1 data.head()
```

Out[39]:

	Age	Education	JobLevel	MonthlyIncome	MonthlyRate	PercentSalaryHike	PerformanceRating	StockOptio
0	41	2	2	5993	19479	11		3
1	49	1	2	5130	24907	23		4
2	37	2	1	2090	2396	15		3
3	33	4	1	2909	23159	11		3
4	27	1	1	3468	16632	12		3

5 rows × 42 columns

### Observations from the above output:

We can see that all the variables present in the dataframe are in the numerical format.

Note: The rest of the encoding methods will be used in further lessons

## The further operations come under feature Selection

Note: In this lesson, we saw the use of the feature engineering methods, but in the next lesson we are going to use one of these methods as a sub component for "Exploratory Data Analysis".

Powered by  simplilearn