

Feature Engineering

One Hot Encoding:

Feature Engineering - It refers to manipulation such as addition, deletion, Combination, mutation of your dataset to improve Machine learning model training, leading to better performance and greater Accuracy.

One-Hot Encoding:

One hot encoding is a technique that we use to represent Categorical variables as numerical values in Machine learning model.

Advantages:

- Straightforward to implement.
- Does not require hours of variable exploration.
- Does not expand massively the feature Space.

[No of columns in dataset]

Disadvantages:

- Does not add any information that may make the Variable more Predictive.
- Does not keep the information of the ignored labels.

Suppose, we have the "Gender" feature which includes male and female. We need to change this category feature to number or numerical values. It can be represented as 0 (or) 1. 0 for male and 1 for female.

Suppose we have the dataset,

Employee-ID	Gender	Remarks
0	Male	Nice
1	Female	Good
2	Female	Great
3	Male	Great
4	Female	Nice

After applying One hot encoding,

It looks,

Employee-ID	Remarks_Good	Remarks_Great	Remarks_Nic
0	0	0	1
1	-	-	0
2	-	-	0

Gender-female Gender-male

	0	1
0	1	0
1	-	-

One Hot Encoding - Variables with many changes.

```
import pandas as pd
```

```
import numpy as np
```

```
data = pd.read_csv('mercedesbenz-econ', usecols = ['x1', 'x2', 'x3',  
          'x4', 'x5', 'x6'])
```

```
data.head()
```

To find out Unique Categories in each and every column

```
for col in data.columns:
```

```
print (col, ':', len(data[col].unique()), 'labels')
```

Performing One-hot encoding

```
pd.get_dummies(data, drop_first = True).shape
```

```
Out: (4209, 117)
```

↓
117 Columns [It leads to Curse of

dimensionality]

According to this, we use to take top 10 features from the each feature.

(It's Based on KDD Cup in kaggle)

Selecting top 10

One-hot encoding perform top 10 features and remaining are zero.

Step. # Let's find the top 10 most frequent categories

for the variable X_2

data[' X_2 '].value_counts().sort_values(ascending=False).head(10)

top-10 = [x for x in data[' X_2 '].value_counts()]

sort_values(ascending=False).head(10).index]

Now make binary variable (0 and 1) for encoding.

for label in top-10:

data[label] = np.where(data[' X_2 '] == label, 1,

dat[[' X_2 '] + top-10].head(40)

Perform one hot encoding for top 10

```
def One-hot-top-x (df, variable, top_x_labels):
```

```
    for label in top_x_labels:
```

```
        df [variable + '-' + label] = np.where
```

```
(data [variable] == label, 1, 0)
```

```
# Read data against fresh list of elements of label
```

```
data = pd.read_csv ('mercedesbenz.csv', usecols=
```

```
[ 'x1', 'x2', 'x3', 'x4', 'x5', 'x6' ] )
```

```
One-hot-top-x (data, 'x2', top_10)
```

```
data.head ()
```

↑

Similar applying for x_1, x_2, x_3 upto x_6 . Finally

removes all the x_1, x_2, x_3 upto x_6 columns.

One-hot Encoding well for ~~noise~~ data because
take most frequent elements and make better
prediction.

Types of Encoding Techniques:

Whenever speak about the encoding techniques,

firstly comes to the mind is Categorical Variable

(Qualitative Variable)

Because, we need to change Categorical Value to Numerical like, Gender [Male, Female] like this.

Types,

1. Nominal encoding.

[Order does not matter]

Gender

Male

State

Keywork

Female

spain

Male

Keywork

This type of data is known as the Nominal encoding.

2. Ordinal encoding.

[Order does matter]

Education dataset,

X BE PPE

y BCom .

Z Phd

A BE

Order does, The rearrangement preference based on Order.

Phd - 1

BE - 2

Bcom - 3

Nominal Encoding

Ordinal Encoding

↳ One hot Encoding.

↳ Label Encoding.

↳ One hot Encoding with Many Categorical.

↳ Target guided Ordinal Encoding.

↳ Mean Encoding.

One hot Encoding:

we already discussed the one hot Encoding,

Suppose,

State Germany France Spain

Germany

France

Spain

Normally, If applied one hot encoding. (Any of the One Column) may be deleted.

After apply one hot encoding,

State	Germany	France	Spain
Germany	1	0	0
France	0	1	0
Spain	0	0	1

The Spain column removed, after that

- 0 → Represent Germany
- 1 → Represent France
- 0 0 → Represent Spain

Suppose, Consider another feature [pincode]

Pincode

600001

100 pincodes,

600002

The one hot encoding creates
99 features.

600004

600008

:

600124

(Suppose 100 pincodes)

It leads to Curse of

dimensionality.

So, Some other techniques

Used.

Label Encoding:

Label encoding technique used for Ordinal Categories.

Example,

Education	Rank	
BE	1	We are putting values,
ME	3	Based on the Rank.
Phd	4	[It represents the
STATS	2	higher most important].

This is known as Label Encoding.

One hot Encoding with multiple Categories:

Suppose the feature, f_1 have 50 categories.

The Kaggle Competition, KDD Cup Orange proves the [Select top 10 categories and apply One hot encoding] is efficient.

This also proves with Researcher.

So, top 10 \Rightarrow Applying one hot encoding \Rightarrow Then it creates the 9 columns.

Target Guided Ordinal Categories: (Encoding)

functions info into ordinal position label

Suppose we have,

f_1	O/P	Mean	RANK
A	1	$\frac{1}{2} = 0.5$	2
B	1	$\frac{1}{3} = 0.33$	3
C	0	$\frac{1}{2} = 0.5$	2
D	1	$\frac{1}{1} = 1$	1
B	0	$\frac{1}{3} = 0.33$	3
A	0	$\frac{1}{2} = 0.5$	2
C	1	$\frac{1}{2} = 0.5$	2
B	0	$\frac{1}{3} = 0.33$	3

The weighted mean also calculate with various methods in target encoding.

Mean Encoding:

f_1	O/P	Mean (Suppose)
A	1	0.73
B	0	0.6
C	1	0.5
D	1	0.4
A	0	-
B	1	-

this value
is fed into
machine learning
training.

Why use these techniques?

Suppose, we have pincode feature. It contains 1600 categories. The One hot encoding is very large dimensionality. It leads to Curse of dimensionality.

So rather used, Mean Encoding.

Based on the dataset, we use different different encoding techniques.

Pincode:

60001
60002
60003 } Mean \Rightarrow to machine.

FEATURE SCALING:

Why feature Scaling?

Suppose we have the,

height weight BMI

180 78

170 84

Every features have different different units

and magnitude.

Without feature Scaling,

Some algorithm does not work

Some algorithm does not work

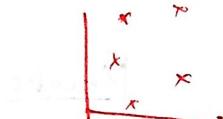
Not working properly.

Example:

k-nearest neighbour based method

Suppose plot the points in this algorithm, there is a huge difference between the datapoints and

Very long distance points covered.



Some Scaling techniques are,

→ Standard Scaler.

→ Minmax Scaler etc.

When apply feature Scaling?

Some algorithm like,

→ Linear Regression for [Gradient descent]

→ Euclidean distance related algorithm like

kNN, k-means [Distance problem avoid]

→ CNN [unit scaling for images] [RGB 0-255]

When does not require feature Scaling?

→ Decision Tree algorithms.

→ Random Forest.

→ Xgboost

Because, it is tree based decision algorithm.

Whenever apply the Scaling or not. Both the scenes are well-performed.

HANDLE MISSING VALUES IN CATEGORICAL VARIABLE:

There are four steps (or) four techniques are available,

→ Delete the Rows.

→ Replace the Most frequent values.

→ Apply Classification classifier algorithm to predict.

→ Apply UnSupervised Machine Learning algorithm.

Consider,

	f_1	f_2	f_3	Op
MALE	23	24	Yes	
	24	25	No	
FEMALE	25	26	Yes	
MALE	26	27	Yes	

Delete the Rows,

→ The very less no. of rows^{Missing} in the huge dataset.

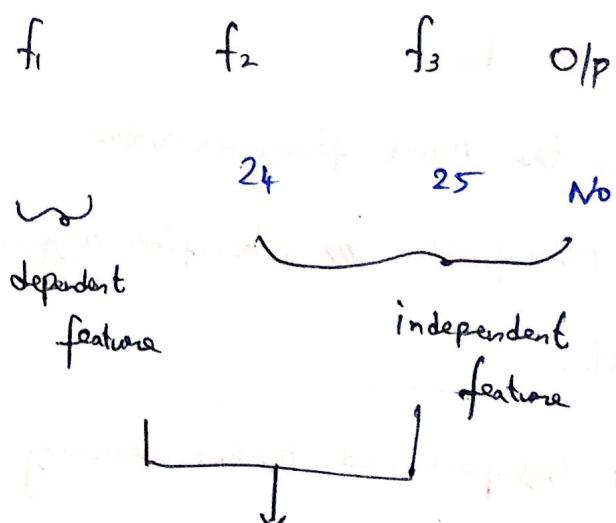
→ Delete the entire row.

Replace With most frequent value,

To fill with most frequent value Using "Mode" method in Statistics.

Apply Classifier algorithm,

Take the row,



Using test data

Remaining things are to be considered as the Training data to predict the outcome.

Apply Unsupervised ML:

Using Clustering technique to group the features.

$$f_2 \quad f_3$$

24 25 → which groups are related
and assign them to predict.

Using Unsupervised ML algorithm. kmeans, kNN etc.

HANDLING CATEGORICAL FEATURES FOR MANY CATEGORIES:

COUNT/FREQUENCY ENCODING:

Higher No of labels (or) Higher No of Categories
is known as High Cardinality.

Definition:

The frequency encoding is replace each label of the Categorical Variable by the Count, this is the amount of times each label appears in the dataset.

Example:

import Pandas as pd

import numpy as np

df = pd.read_csv('mercedesbenz.csv', usecols=[x1, x2])

df.head()

df.shape

Out: (4209, 2)

Applying one hot encoding

pd.get_dummies(df).shape

Out: (4209, 71) # more columns created.

If leads to maybe curse of dimensionality.

Applying frequency label encoding

len(df['x1'].unique())

len(df['x2'].unique())

look how many labels

for col in df.columns[0:]:

print(col, ':', len(df[[col]].unique()), 'labels')

Out:

X1: 27 labels

X2: 444 labels

Take X2 because it is high Cardinality (or) take any.

Count X2 Each Unique Category Value Count

df.X2.value_counts().to_dict()

To dictionary purpose is the need to mapping

Variable -

df_frequency_map = df.X2.value_counts().to_dict()

Now replace the X2 values in the dataset

df.X2 = df.X2.map(df_frequency_map)

df.head()

Before applying:

After Applying:

X1 X2

v at
t av

w n

t h

v h

X1 X2

v 6.000000
t 4

w 137

t 137

v 137

Advantages:

→ It is Very Simple to implement.

→ Does not increase the feature dimensional space.

[Averts leads to Curse of dimensionality]

Disadvantages:

→ If Some of labels have the Same Count, then they are replaced with Some Count and they will loose Some Valuable information.

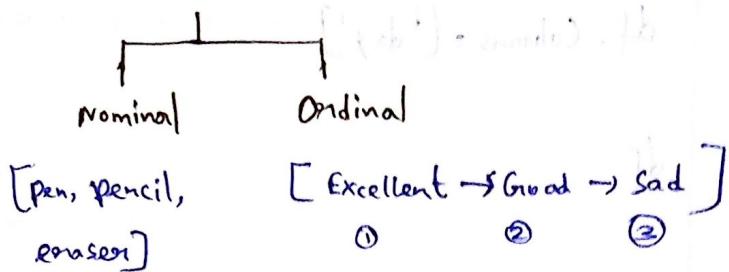
→ Adds Somehow arbitrary numbers, and therefore weights to the different labels, that may not be related to their Predictive Power.

HANDLING ORDINAL CATEGORIES [ORDINAL ENCODING]

Ordinal Categorical Variable:

Ordinal data is a Categorical, Statistical, data type where the Variables have Natural, Ordered Categories and the distances between the categories is not known.

Categorical



nominal & ordinal

Example: [ORDINAL]

→ Student's grade in Exam. [A, B, C or fail]

→ Educational level [Elementary school, high school,

College, graduate, Phd] from rank
1 to 4.

IMPLEMENTATION:-

```
import pandas as pd
```

```
import datetime
```

Create a variable with dates, and from that extract
the weekday.

I create a list of dates with 20 days difference
from today.

Then transform into dataframe -

```
df_base = datetime.datetime.today()
```

```
df_date_list = [df_base + datetime.timedelta(days=x)  
for x in range(0, 20)]
```

df = pd.DataFrame (df_data_list)

(using dict)

df. columns = ['day']

df

df = df.set_index('date')

④

⑤

⑥

⑦

⑧

⑨

⑩

⑪

⑫

⑬

⑭

⑮

⑯

⑰

⑱

⑲

⑳

㉑

㉒

㉓

㉔

㉕

㉖

㉗

㉘

㉙

㉚

㉛

㉜

㉝

㉞

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

㉟

Advantages:

Ordinary least squares is normally used for fitting and generalization.

With respect to machine learning,

- keeps the geometrical information of the variables.
- Straightforward approach.

Disadvantages:

With respect to machine learning,

- Does not add machine learning valuable information.