

Encoding categorical Data

Categorical data

Ordinal

Nominal

There is a relationship between the categories present in your data.

For Ex - Un, Pl, HS

PS > UC > HS

Excellent > Good > Bad

We have an order.

categories doesn't have any relation between them

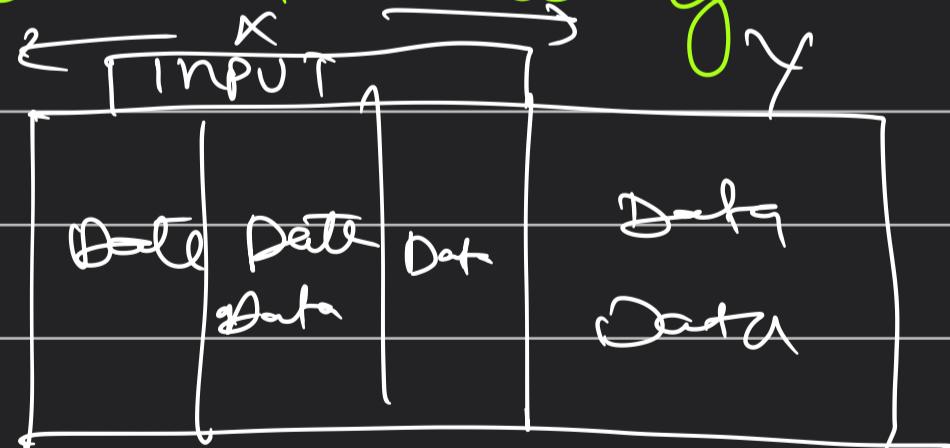
For Ex - State, Equity Branches.

We use One Hot Encoding Here

Why Encoding

Now our categorical data is present in the form of strings & ML Algo require numerical values.

Ordinal Encoding



Whenever we have input columns, collectively called X , if there is some ordinal data present there, we use ordinal encoder there.

But if our Y (Output Variable) is categorical in nature. Then we don't use ordinal encoder, we use something called "Label Encoder".

So Input - X \Rightarrow Ordinal Encoder
Output - Y \Rightarrow Label Encoder

Ex	Education	But we know
	High School 0	PG > VH > HS
	PG 2	
	VH 1	PH \rightarrow 2
	VH 1	VH \rightarrow 1
	PG 2	HS \rightarrow 0
	High School 0	
	2	+ ordinal
		Done with Encoding

Nominal Encoding

is a categorical encoding where we don't have order between our categorical data.

Ex Colors, Brands of Car they all have equal importance

\Rightarrow Nominal Encoding

Color	Target	Color_A	Color_B	Color_R	Target
Yellow	0	1	0	0	0
Blue	1	0	1	0	1
Red	1	0	0	1	1

Yellow	0	1	0	0	0
Blue	1	0	1	0	1
Red	0	0	0	1	0

of color

Noe our string in categorical data
got converted into vectors.

$$\text{Yellow} = [1, 0, 0]$$

$$\text{Blue} = [0, 1, 0]$$

$$\text{Red} = [0, 0, 1]$$

• Dummy Variable Trap

Here we remove one column or variable
from our Input variables.

Say we N columns as input variable
but we shall only use $N-1$ columns.

But why do we remove this?

This is due to Multicollinearity)

Input columns may have some
mathematical relationship between.

If this exist then those input variables
are not suitable for ML.

Because in ML we require Independent
variable ^{as input} & only our Target Variable
should be dependant on "Input".

So to break this multicollinearity we

remove one column.

For ex we remove color-Y

? what would be color of yellow now

→ previously,

$$\text{Yellow} = [1, 0, 0]$$

$$\text{Blue} = [0, 1, 0]$$

$$\text{Red} = [0, 0, 1]$$

Now after removing color-Y

$$\text{Yellow} = [0, 0]$$

$$\text{Blue} = [1, 0]$$

$$\text{Red} = [0, 1]$$

This removing of column helped us
reduce dimension.

- One hot encoding using most frequent variables

Let say we have list of brands
say count - 40.

We need to encode them there are few
brands that comes most frequent.

So what we can do is put the least
frequent brands in their own category
called "others"

This way we reduce the dimensionality

