

Encoding

30 May 2023 20:39

How to convert categorical data into numerical?

- * Label Encoder
- * One Hot Encoding / Get Dummies

Label Encoder

Label Encoder is one of the preprocessing method in machine learning which is used to convert categorical data into numerical data.

Gender	Gender
Male	1
Female	0
Male	1
Female	0

Height	Height
Short	1
Medium	0
Tall	2

Notes

- * Label Encoder assigns numbers for alphabetical order

Color	Color
Red	2
Green	1
Blue	0

Limitation

Gender	Gender
Male	1
Female	0

- * Label Encoder leads to priority issue.
- * A label with high value will be considered as to have priority than the label with low priority.

$$1 > 0$$

Mark	Review
Fail	bad
Pas	Good

For these data label encoding.

One Hot Encoding / Get Dummies.

It is one of the preprocessing method used to convert categorical data into numerical data.

- * It is used for nominal data [Data has no order].
- * It creates number of columns equal to the number of categories.

Gender	Gender_male	Gender_femal
Male	1	0

$$\sim \rightarrow 1 \rightarrow \text{Male}$$

* It creates number of columns - true

<u>Gender</u>	Gender-male	Gender-female
Male	1	0
Female	0	1
Male	1	0
Female	0	1

1 1

[1, 0] → Male
[0, 1] → Female

<u>Color</u>	Color-red	Color-green	Color-blue	
Red	1	0	0	[1, 0, 0]-red
green	0	1	0	[0, 1, 0]-green
blue	0	0	1	[0, 0, 1]-blue
Red green	1	1	0	

Limitation

- * # of columns increases and results in overfitting.
- Overfitting → Training score is good and testing score is less.

Breaks 10 - 10-10 Pm

How to overcome?

- * Use one hot encoding when categories are less than 3.
- * Drop one of the existing column which you have created.

<u>Color</u>	Color-red	X	Color-green	Color-blue	Color-orange	
red	1	X	0	0	0	[0, 0, 0]-red
green	0		1	0	0	[1, 0, 0]-green
blue	0		0	1	0	[0, 1, 0]-blue
orange	0		0	0	1	[0, 0, 1]-orange
red	1		0	0	0	
green	0		1	0	0	
blue	0		0	1	0	