# Feature Transformation:

Data → Numerical
↳ Categorical data



| Nominal | Ordinal |
|---|---|
| → No relation or No order between Variables | There is a order or relation b/w the variables. |
| eg: parle-g, monaco, krack-jack. | eg: good, better, best eg: |

# Ordinal Encoding
↳ ~~Nominal~~ Ordinal Data

{ for output colomn Target column
eg: Yes or No, Cwrn prediction, we use label encoder

eg: Education

High School)
UG
PG
PG
UG
HS
HS

} Ordinal data.
PG > UG > HS
HS => 0
UG => 1
PG => 2
} Transformed based on order.

Label Encoding ———→ output column

Ordinal Encoding ———→ input column

# One hot encoding !

⇒ why? ↳ Nominal chartegorical

↳ it doesnot have any orde

color

$\left.\begin{array}{l} \text{Yellow} \\ \text{Green} \\ \text{Red} \end{array}\right\}$ No order

if we encode with ordinal or Label encode. ML model will prioritize more on having greater numbers.

No matter how much unique categories you have in a column.

⇒ Dummy Variable Trap.

we remove one column (generally first column)

|  | Y | r | G |
|---|---|---|---|
| Yellow | 1 | 0 | 0 |
| red | 0 | 1 | 0 |
| Green | 0 | 0 | 1 |

remove ←

⇒ ## Multicollinearity:

In ML we should not have any rel$^n$
btw Independent columns.

by doing OHE :

|       | Y | r | G |
|-------|---|---|---|
| Yell  | 1 | 0 | 0 | $\Sigma = 1$
| reel  | 0 | 1 | 0 | $\Sigma = 1$
| Green | 0 | 0 | 1 | $\Sigma = 1$

Summation = 1

$\Sigma = 1$ of each row

this is not correct for linear models

linear models → linear regression

→ logistic regression.

then How would I be able to represent
Yellow column.

|   | r | G |
|---|---|---|
| ~~Y~~ | 0 | 0 | → If both are
| ~~0~~ | 1 | 0 | zero then it's yellow,
| ~~0~~ | 0 | 1 | problem solved.

# # OHE using most frequent variables

Brand nominal $\rightarrow$ 40 diff
brand cars

$\downarrow$

"OHE"

$\downarrow$

Increase dimensions

$\downarrow$

creat most freq categories
to "others".

# # Column Transformer:

Imputing

| Age | City | grades | review |
|-----|------|--------|--------|
| | | OHE | OE |

so missing values

Now, here we have to transform all the features with some different techniques

we will be getting many numpy arrays. and it may cause trouble while making Machine learning pipelinee.

So, we use column transformer