# (X) Feature Engineering

## (1) Feature Selection

selecting attributes which fit best with independent var & target var

There are certain features which are more imp than others.

→ Chi-squared test
→ Correlation coeff scores
→ LASSO
→ Ridge Regression

## (2) Feature Transformation

Transforming original feature to func of orig features.

→ Scaling
→ Discretization
→ binning, filling
→ Missing data values

To reduce right skewness of data we use [log]

* ③ **Feature Extraction**
* When data is large ⇒ redunda
* For tabular data ~~use~~ PCA
* For Image use line, edge detecti

## SIMPLE LINEAR REGRESSION :-

Relationship b/w dependent &
independent var can be expressed i
st line



Linear Regression

Simple linear                    Multiple linear

Notes

Simple linear = When X & Y have
linear relationship $y = mx + c$
To chk this, XY Scatterplot = linea

SYEDA DARAQSHAN

& residual plot shows random pattern

## Multiple Linear Regression

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 \cdots$$

1 - Independent var (cont)

> 1 Independent var $(x_1, x_2)$

For Regression

① Convert Categorical var to Cont var

1·1 = Label Encoder

↳ For Dichotomous Var

  Yes = 1      No = 0

↳ For Nominal Var

Red = 0      Grey = 2
Blue = 1

↳ For Ordinal Var

SYEDA DARAQSHAN

08  $y =$ that sinking feeling

09  Ranking with labels.

10  | Ordinal x | Encoded x |
    |-----------|-----------|
    | Bad       | 0         |
11  | Shit      | 1         |
    | Ah, tak   | 2         |

12  ✳ When facing nominal var that
13  should not have diff weightag

14  Red = 0 ⟹ Grey has more weightage
15  Blue = 1      than Red, blue
    Grey = 2          & so on.
16  But All are same.

17  So ONE HOT ENCODING used
18  (Creating Dummy Var)

Notes If var is 3 colours then you
should only be using 2 dummy
var.

- Run model again with chosen vary trying one of rem var at each time & sticking with best.
- Repeat until adding does not improve model

$$X_1, X_2, X_3 \cdots$$
$$\downarrow best$$

$$X_2, X_1 \mid X_2, X_3 \cdots$$
$$\downarrow best$$

② Backward elimination

- Start with all var
- Try model out multiple times, excluding one var at a time.
- Remove var that causes model to improve the most when it is left out.
- Repeat

$$X_1 + X_2 + X_3 + X_4$$
$$\uparrow$$
$$\text{best when removed}$$

SYEDA DARAQSHAN

$$X_1 + X_3 + X_4$$
$$\uparrow$$

$$X_1 + X_3 \implies \underline{\underline{X_3}}$$
$$\uparrow$$

③ Bidirectional Elimination

forward + backward

Like added $X_1$ $X_2$ $X_5$ $X_8$

del $X_2$.

COLLINEARITY , CORRELATION

↳ helps to get rid of var that
are skewing data

Correlation — describes relation
b/w 2 var.

If extremely correlated then
collinear

SYEDA DARAQSHAN