

07-09-2025

Agenda: —

FE-1V

Data encoding

label encoding, Ordinal encoding

one hot encoding

(1) Red 400 - 490 Thz
 (2) Orange 490 - 510 Thz
 (3) Yellow 510 - 530 Thz
 (4) Green 530 - 600 Thz
 (5) Blue 600 - 670 Thz
 (6) Indigo 670 - 700 Thz
 (7) Violet 700 - 750 Thz

10
20
30
40
50
60
70

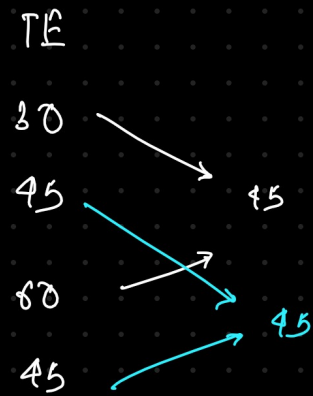
color-freq	category of color	LE (sort alpha beta)	LE(Ordinal) sort by order	B-OHE	R-OHE	I-OHE	V-OHE	O-OHE
625	Blue	0	2	1	0	0	0	0
420	Red	3	0	0	1	0	0	0
680	Indigo	1	3	0	0	1	0	0
690	Indigo	1	3	0	0	1	0	0
710	violet	4	4	0	0	0	1	0
490	orange	2	1	0	0	0	0	1

color-freq	category of color	Target-encoding
625	Blue	value-3
420	Red	value-2
680	Indigo	value-1
690	Indigo	value-1
710	violet	value-4
490	orange	value-5

Indigo → mean
 680, 690 → value-1
 orange → mean
 490 → value-2
 ...
 ...



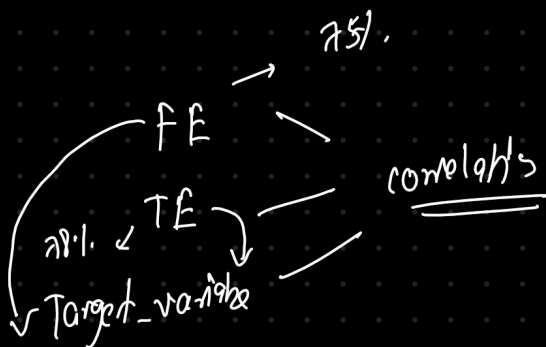
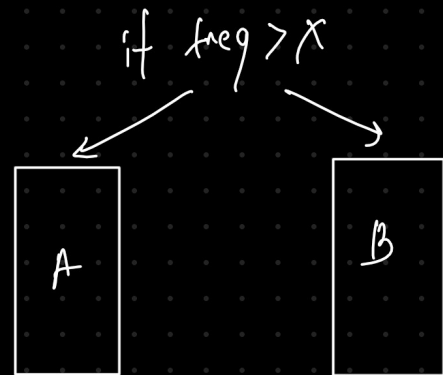
node	Marks
A	-
B	-
A	-
B	-
B	-
C	-



color-freq category of color

frequency - encoding

625	Blue	1
420	Red	1
680	Indigo	2
690	Indigo	2
710	violet	1
490	orange	1



OHE \rightarrow sklearn \rightarrow sparse matrix

coords values

0, 0 1.0

1, 1 1.0

2, 2 1.0

3, 3 1.0

4, 4 1.0

B_OHE, R_OHE, T_OHE, V_OHE, O_OHE

0 [1. 0 0 0 0]

1 [0 1 0 0 0]

2 [0 0 1 0 0]
[0 0 1 0 0]

3 [0 0 0 1 0]

4 [0 0 0 0 1]

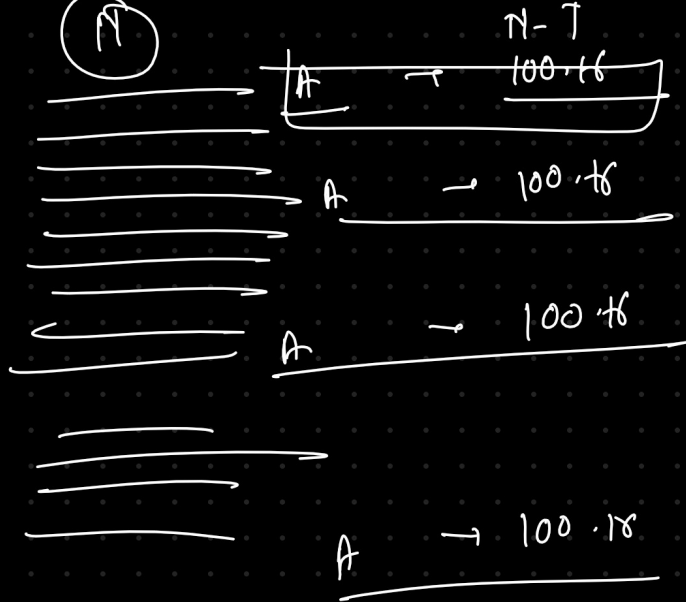
1	0	0	0	0
0	1	0	0	0
0	0	1	0	0
0	0	0	1	0
0	0	0	0	1

← sparse matrix
→ OHE

A B C

→ unique Dep all rows

1420



feature selection :—

df \rightarrow 80 columns
 \downarrow
 100 column

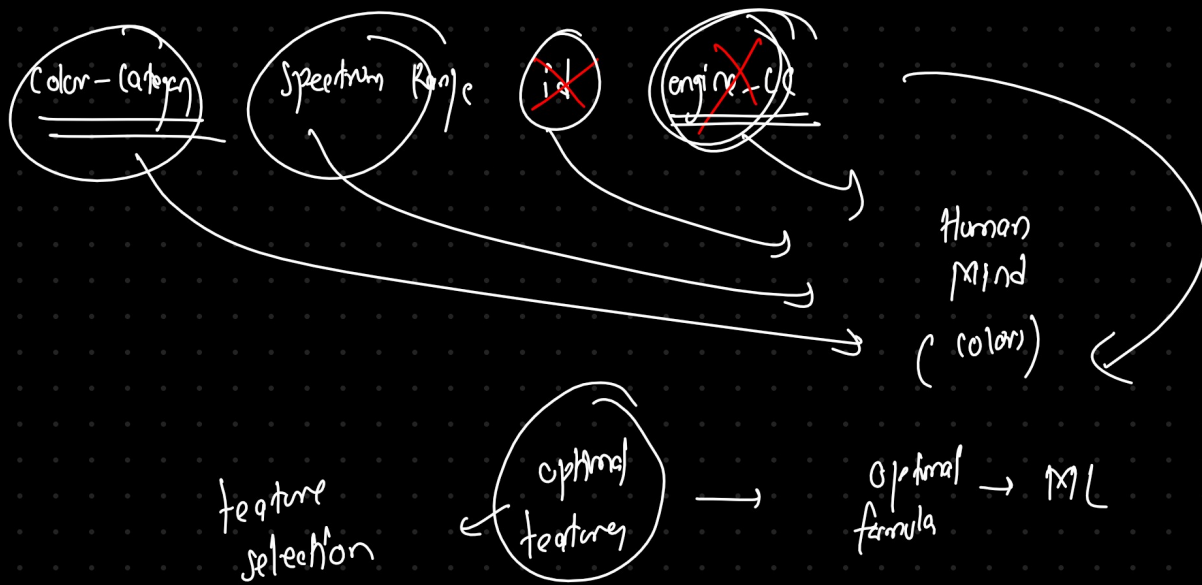
ML

$m_0 + m_1 x_1$ \rightarrow linear regression
 \swarrow
 1 feature

$$\boxed{m_0 + m_1 x_1 + m_2 x_2 \dots + m_{100} x_{100}}$$

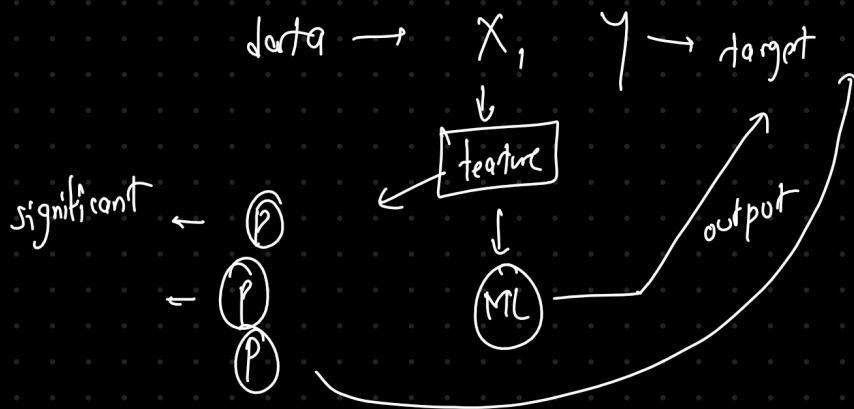
no feature column
 \downarrow
 constant
 \downarrow
 Intercept

$$\boxed{y = mx + c}$$



feature selection :—

(1) Backward Elimination (based p-values in regression)



(1) data



(2) data → X_1 y

(3) $\overset{\text{add constant}}{\hat{X}} \leftarrow ML \leftarrow y$ (learn)

(4) ML → p-value for all feature-column of X

(5) check highest p-value column (least significant feature)

(6) remove that column from X

(7) repeat step 3-6 until all high p-value column are removed.

A	1.78
B	1.98
C	2.56
D	0.03
E	2.30

Target ←

A	1.3
B	1.98
D	0.01
E	1.56

Target ←

A	0.04
D	0.01
E	1.16

Target ←

A	0.04
D	0.03

$A, D < 0.05$

(2) Forward Selection

(1) data → data, columns (collect column names)

(2) X, Y ($X \rightarrow 1000$ row, 12 column)

(3) (only constant)
(1000 row, 1 column) \xrightarrow{X} $ML \leftarrow Y$

(4) $ML \rightarrow R^2$ of all columns
(bigger is better)

(5) get the name of a column having highest R^2 value

(6) add any 1 column from step 1 to X

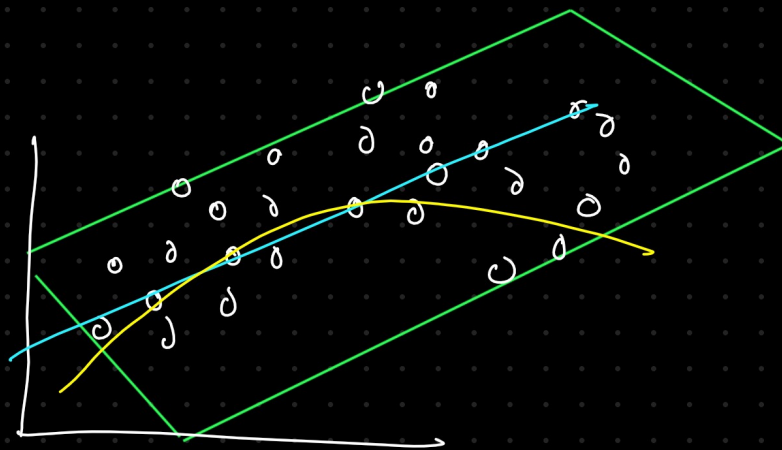
(7) repeat step 3 & 6 until all columns are done

x_0 (constant)	x			y
	A	B	C	D
1	Blue	128	11.12	23
1	Blue	182	11.34	22
1	Red	3000	36.86	21

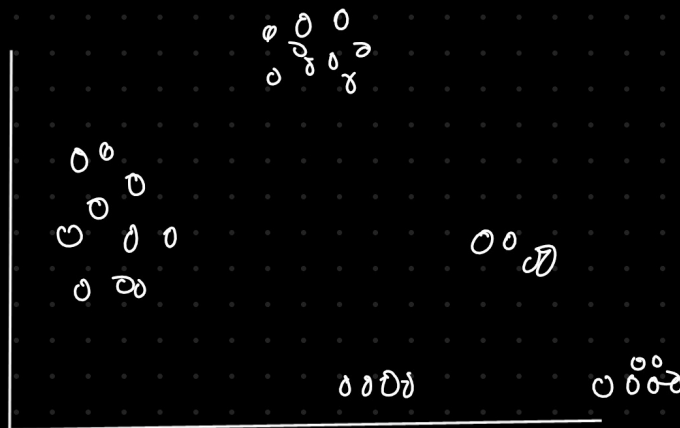
← data →

$x \rightarrow ML$

$$m_0 + m_1 x_1 + m_2 x_2 + m_3 x_3$$



linear
data



Non linear
data

Tree based model

Random forest \rightarrow importance

ML

