

Handling Categorical Data

(pandas) map

(sklearn) ordinal enc

3
2
1
0

owner type \propto Car price

Ordinal
owner_Type

first > 0
second > 1
third > 2
fourth > 3

owner type \propto $\frac{1}{\text{Car price}}$

Nominal

Binary

Transmission

Auto - 0
Manual - 1

Trans. \rightarrow Auto Manual
1 0
0 1
:
:

Location

C1
C2
C3
C4
:
C10

Per Row
↓
① hot



Location	Location → Nominal			
	C1	C2	C3	C4
C1	1	0	0	0
C2	0	1	0	0
C3	0	0	1	0
C4	0	0	0	1

One
Hot
Encoding

Unique values ↑↑ → New Columns ↑↑
Sparsity ↑

unique values ↑↑ (Brand, Model)

↓
Binary Encoder

Brand

Step ①
Decimal

Step ②
Binary

Step ③
Columns

B1 → 0

B2 → 1

B3 → 7

⋮

B50 → 49

000000

000001

000010

⋮

110001

Brand 0 ... Brand 5

0

0

0

0

0

1

0

1

1

...

1

50 unique values

6 Columns

$$2^4 = 16$$

$$2^5 = 32$$

$$2^6 = 64$$

$$2^7 = 128$$

$$2^8 = 256$$

Scaling

f_1	f_2	f_3
1	1000	0.7
↓	↓	↓
10	5000	0.8

different scales

Standardization

Standard Scale

$$\frac{x - \mu}{\sigma}$$

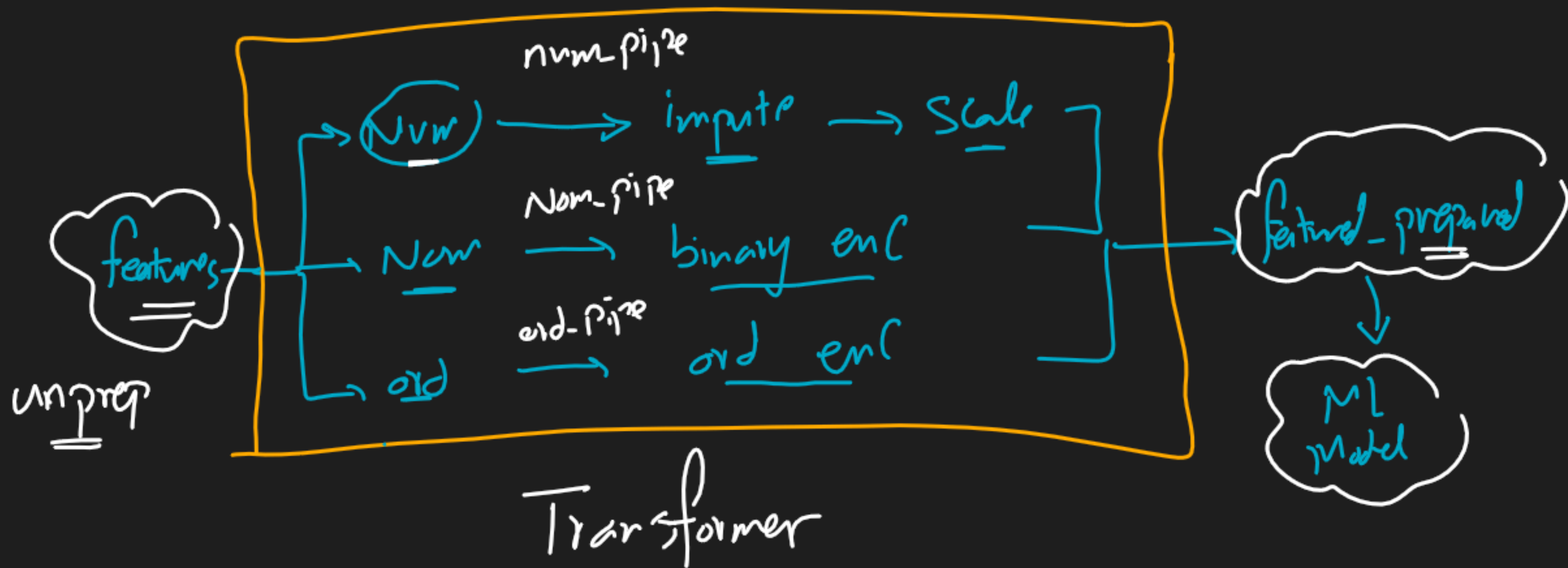
Mean 0
Std 1

Normalization

Min Max Scale

$$\frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

Min 0
Max 1



preprocessing

use

- ① Calculation
- ② Conversion



fit



transform

train

train

test

fillna("median")

Mileage

train
80%

test
20%

NaN

NaN

fit

Calculate(Median) → train

use

Conversion

transform