# LINEAR REGRESSION - 1

**Cars24**

**Optimising Cost**

Buy ← Cars24 ← Sell

$1000

$800  evaluator
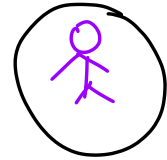
Car

① Develop system online
↳ data about the car

Analysis

Computer Vision
↳ Photographs
{ ↳ Police Record
{ ↳ Service Record.

1000 requests →

↳ Insurance data }
↳ pollution record. }

Output → Continuous → Selling Price ⟩ Regression.

Do we know the target variable → YES ⟩ Supervised

---

Cat
Elephant
dog

🔵 0
🩷 0
🔵 0
🔴 2
🔵 0
🩷 0
🔵 0
🔴 0
🔵 0
🩷 1
🔴 2
🔵 0
🔴 2

(Order)

1000

0

500

1000

Label Encoding

2 Categories

# Target Encoding

| Target Encoding | Target | (HIT/FLOP) (1) (0) |
|---|---|---|

$P[T = 1 \mid Category]$

$0 \rightarrow$ Rom    H   — 0.5

1   drama    H   — 0.66

2   Thriller    H   — 1

0   Rom    F   — 0.5

1   drama    H   — 0.66

1   drama    F   — 0.66

Rom   H
Rom   F     $1/2 = 0.5$

drama   H
drama   H    $2/3 = 0.66$
drama   F

Thriller   H    $1/1 = 1$

1000 classes $\rightarrow$ (TE)

**How do you think we should handle the large number of categories in make and model column?**

38 users have participated

| A | One Hot Encoding | 18% |
|---|---|---|
| B | Label Encoding | 13% |
| C | Target Variable Encoding | 68% |

1000's car → C1

C2

2 column

(1000

C1    C2    C2    - . (1000)

1     0     0    0   0   0

0     1     2    2   2

1000

Input $\longrightarrow$ Models $\rightarrow$ Output

high

① Scale

Standardisation    Normalisation

$99.7\%$

$-3 \underline{\quad\quad} 3$
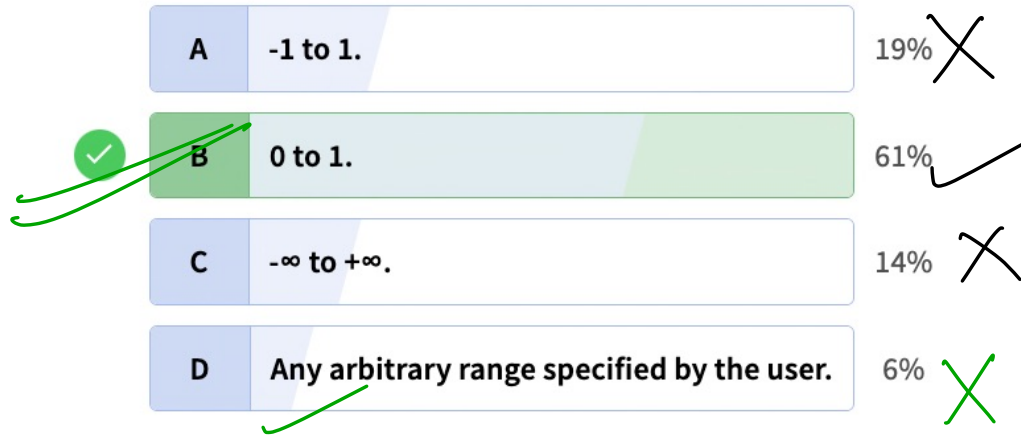
$\boxed{\dfrac{x - \mu}{\sigma}}$    $[0, 1]$  $\dfrac{x - min}{max - min}$

*  −ve values

&  $\mu = 0$

*  I know Range

*  +ve values

**What is the range of values after applying min-max scaling?**
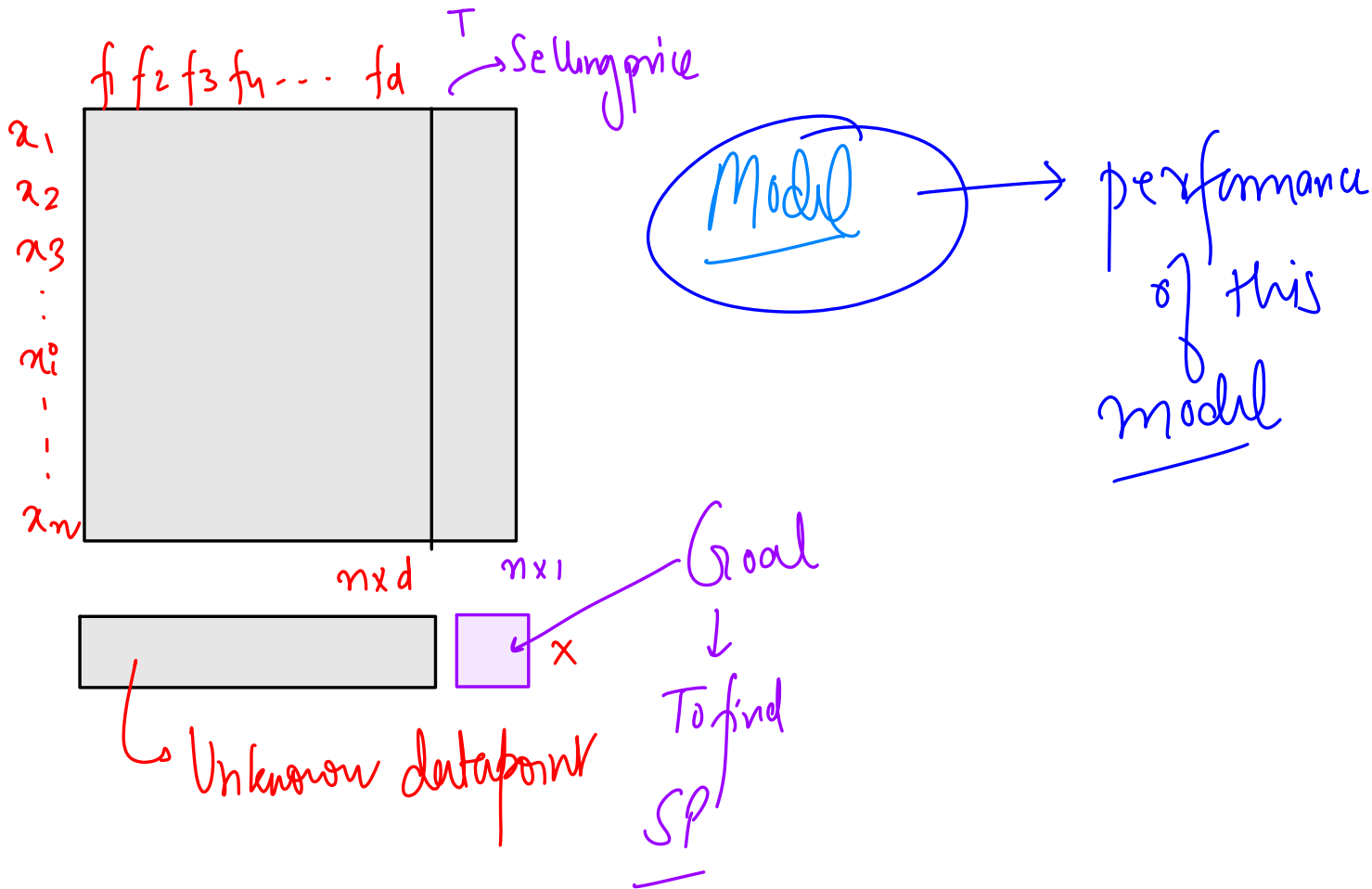
36 users have participated

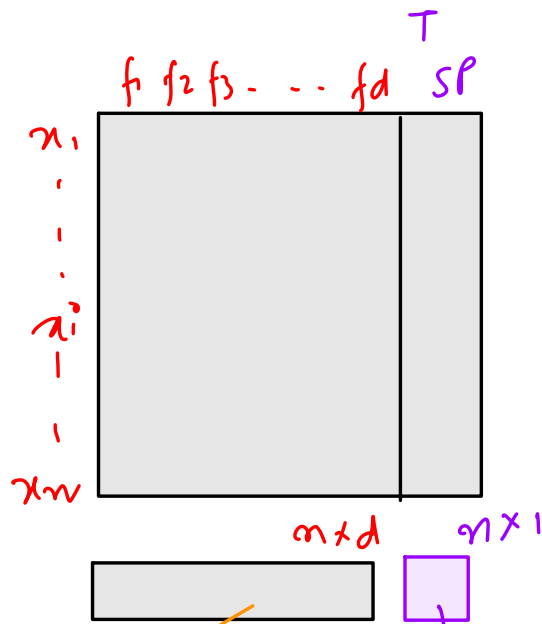| | | | |
|---|---|---|---|
| A | -1 to 1. | 19% | ✗ |
| ✓ B | 0 to 1. | 61% | ✓ |
| C | -∞ to +∞. | 14% | ✗ |
| D | Any arbitrary range specified by the user. | 6% | ✗ |

**Structured data**

Table

**Unstructured data**

Audio, Video, Text, Images

$f_1\ f_2\ f_3\ f_4\ \cdots\ f_d$

T → Selling price

$x_1$
$x_2$
$x_3$
$\vdots$
$x_i$
$\vdots$
$x_n$

$n \times d$   $n \times 1$

Model → performance of this model

↳ Unknown datapoint

X

Goal
↓
To find
SP

T
f₁ f₂ f₃ . . . f_d    SP

$x_1$
.
.
$x_i$
.
.
$x_n$

$n \times d$        $n \times 1$

Unknown      find this out

Dataset → Model → predict
                              unknown

TRAIN { $X_{TRAIN}$    $Y_{TRAIN}$

TEST { $X_{TEST}$      $Y_{TEST}$
                              $Y_{PRED}$

Unknown      find this out

$X_{TRAN}$ → Train Model → prediction

prediction $\rightarrow$ $\underline{X_{Test}}$ $\rightarrow$ $Y_{PRED}$

\# Model perfromanca good

$\Rightarrow$ $Y_{PRED} \backsim Y_{Test}$

\# Model performing badly

$\Rightarrow$ $Y_{PRED} \not\backsim Y_{PEST}$

Retrain Model

$f_1\ f_2\ f_3\ f_4\ f_5\ \ldots\ f_d$  $y$

$x_1$
$x_2$
$x_3$
$x_i$
$\vdots$
$\vdots$
$x_n$

$x_i$

$f_j$

$y_i$

$x_{ij}$

$n = no.\ of\ datapoints$

$d = no.\ of\ features$

$True = y_i$

$Predicted = \hat{y}_i$

$i_{th}\ Sample \rightarrow$

# When should you split your data into training and testing sets?

38 users have participated

| | | |
|---|---|---|
| ✓ A | Before preprocessing the data. | 97% |
| B | After training the model. | 3% |
| C | After evaluating the model's performance. | 0% |

Target

module child mango label apple

$f_1$ $f_2$ $f_3$ $f_n$ $f_2$

min max

scaling

$W_1 f_1 + W_2 f_2 + W_3 f_3 - -$

$W_3 f_n$

KNN

$$f_1 \quad f_2 \quad f_3 \quad f_4 \quad f_5 \quad f_6 \qquad y$$

$$\mu_1 \quad \mu_2 \quad \mu_3 \; \mu_4 \qquad M_y$$
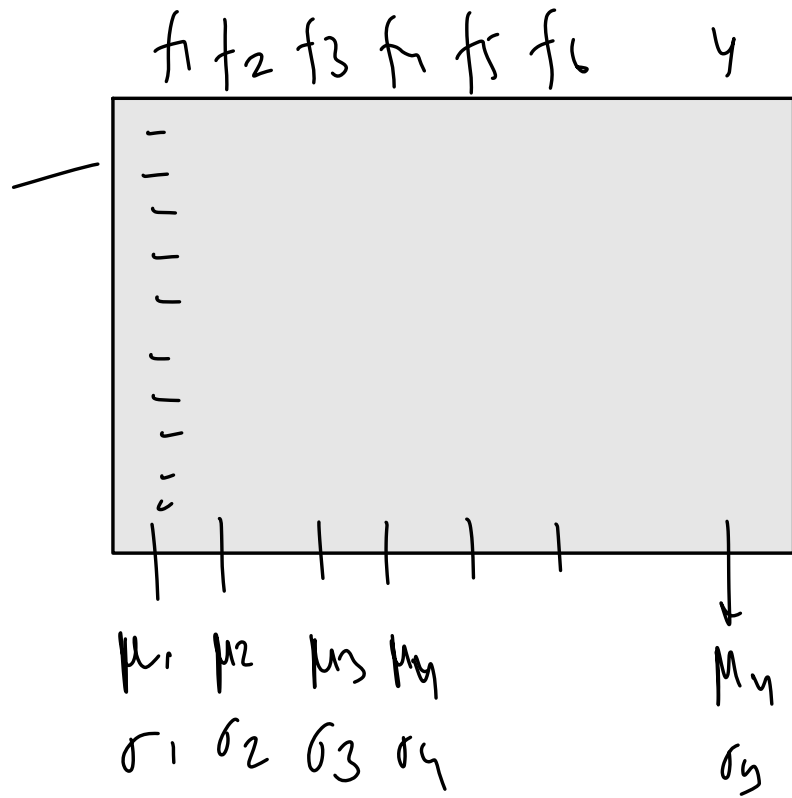
$$\sigma_1 \quad \sigma_2 \quad \sigma_3 \; \sigma_4 \qquad \sigma_y$$

$$\frac{x^0_{11} - \mu_1}{\sigma_1}$$