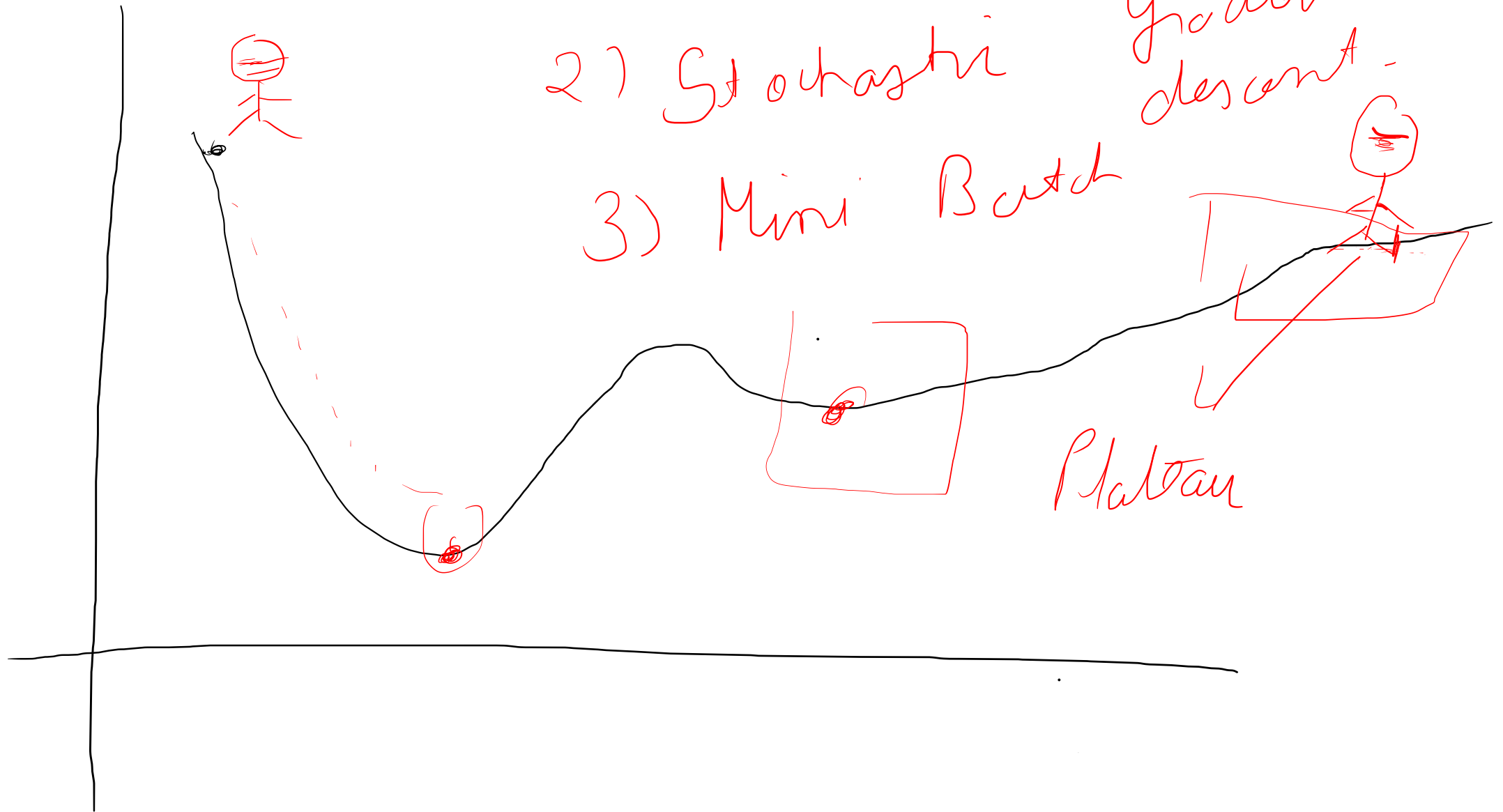# Starting of Gradient descent :→

1) Batch
2) Stochastic
3) Mini Batch

Gradient descent.
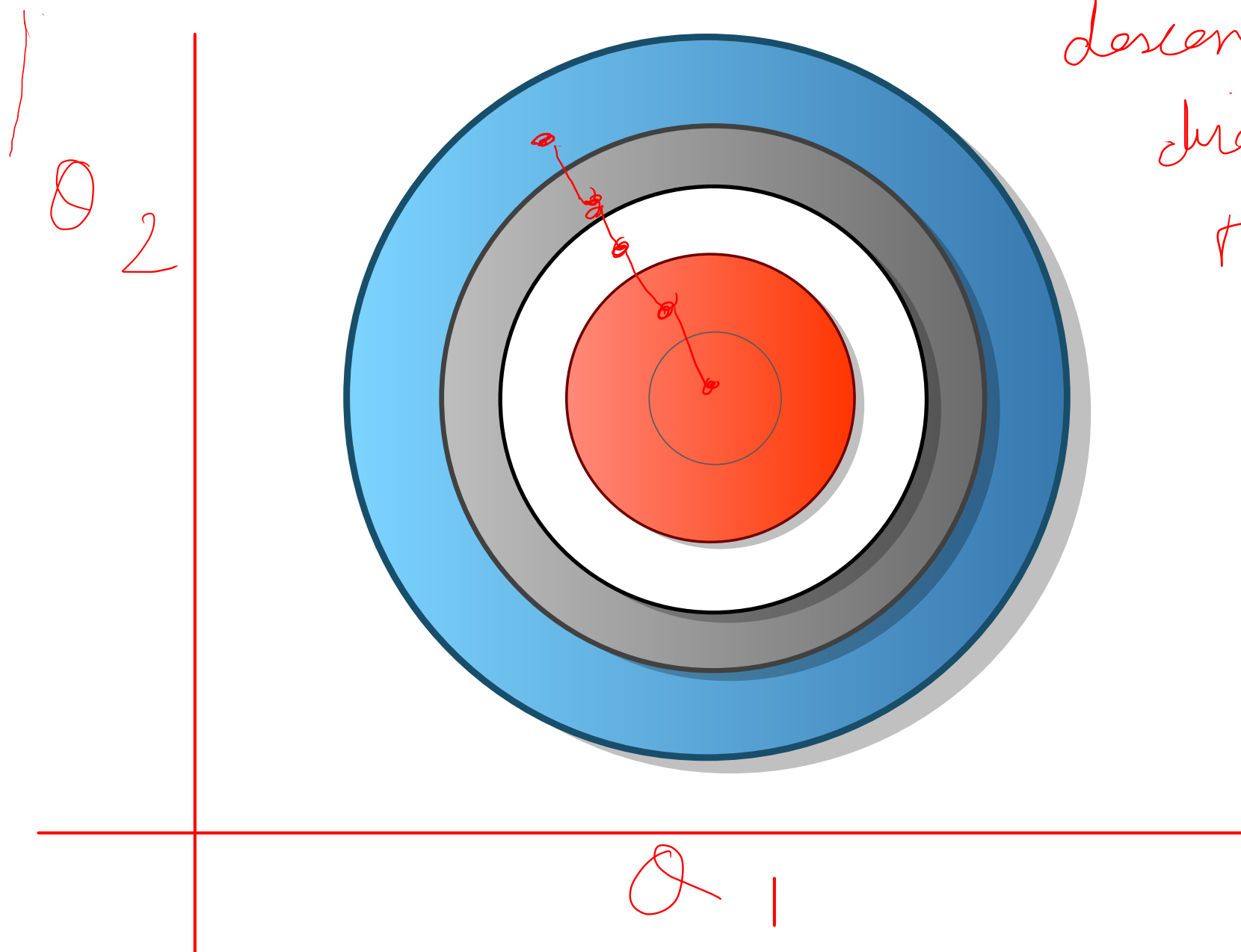
Plateau

Why we use MSE as our cost function
in Linear regression.

As it is a Convex function (which
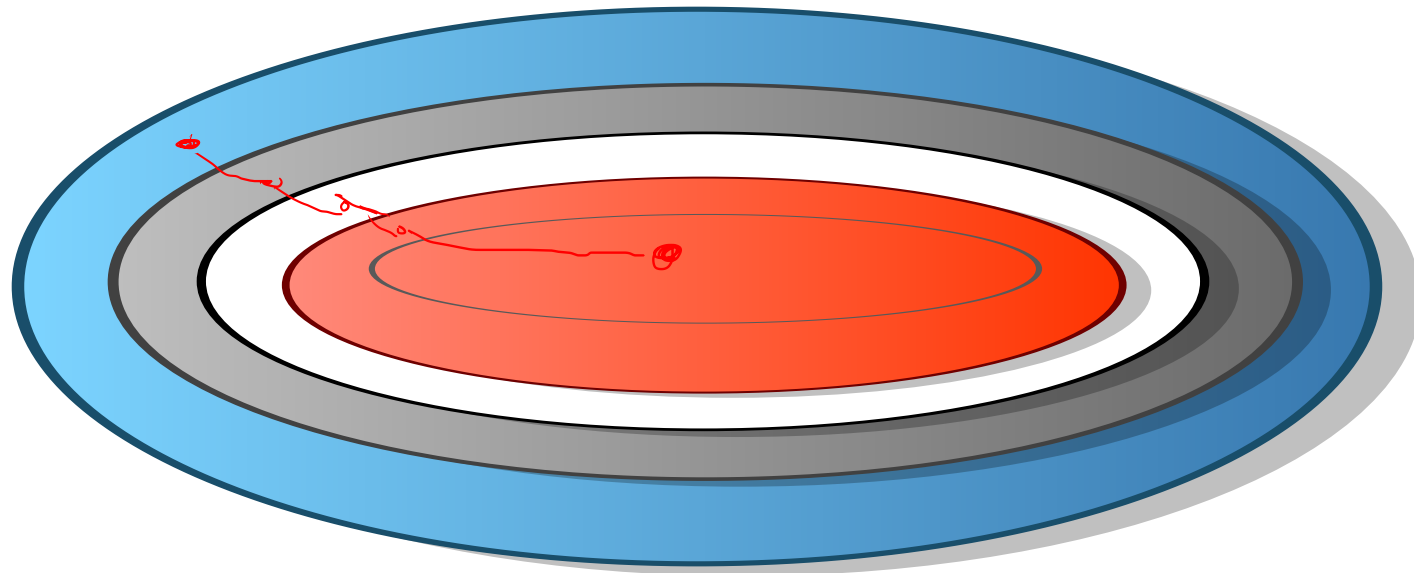means that it has only 1 minima)

V. V. Imb:

While using Gradient descent
our data must be scaled.

If The **x** features are having same scale. The gradient descent will directly travel toward minima

If we are not scaling the data. Then we will travel longer path to get to minima

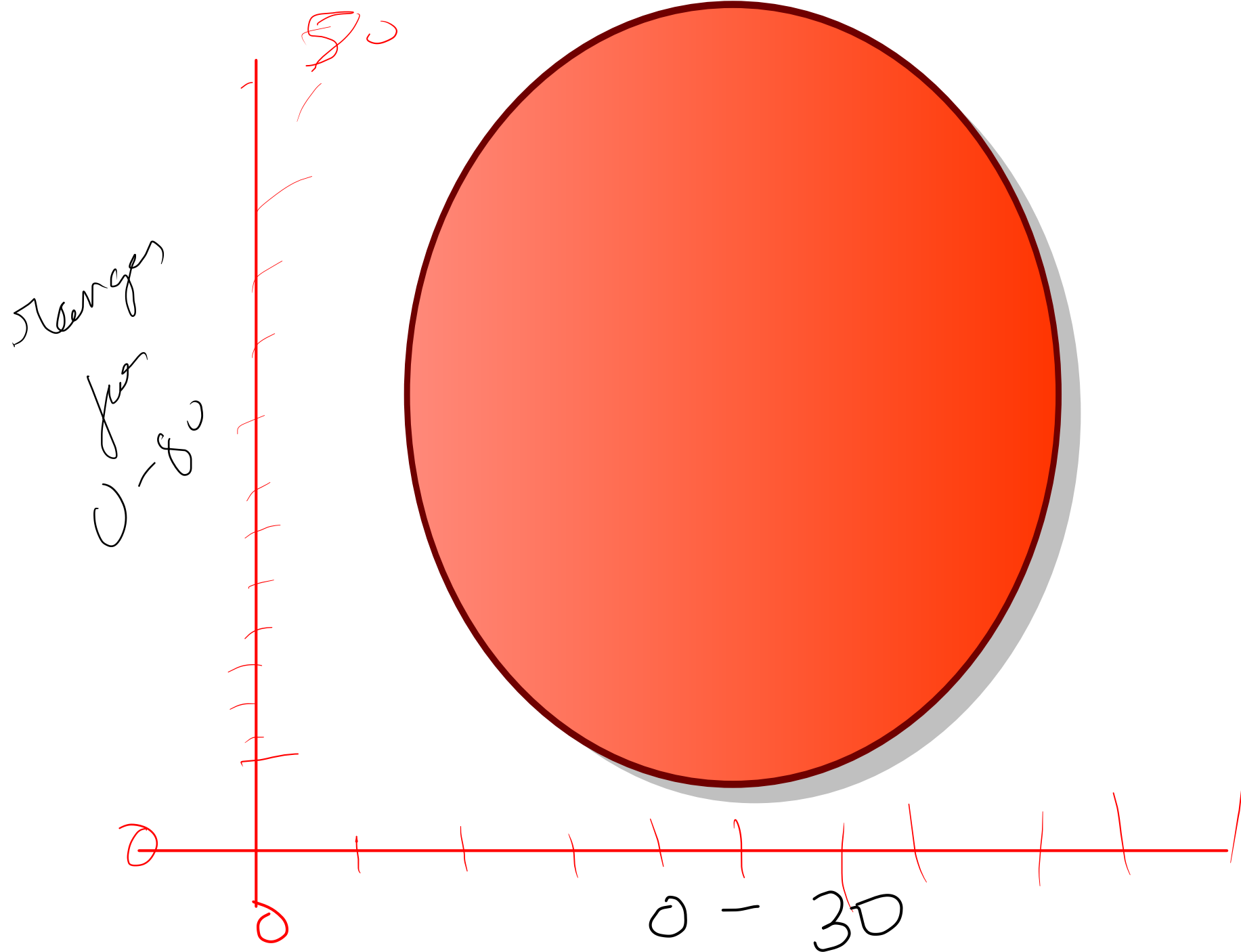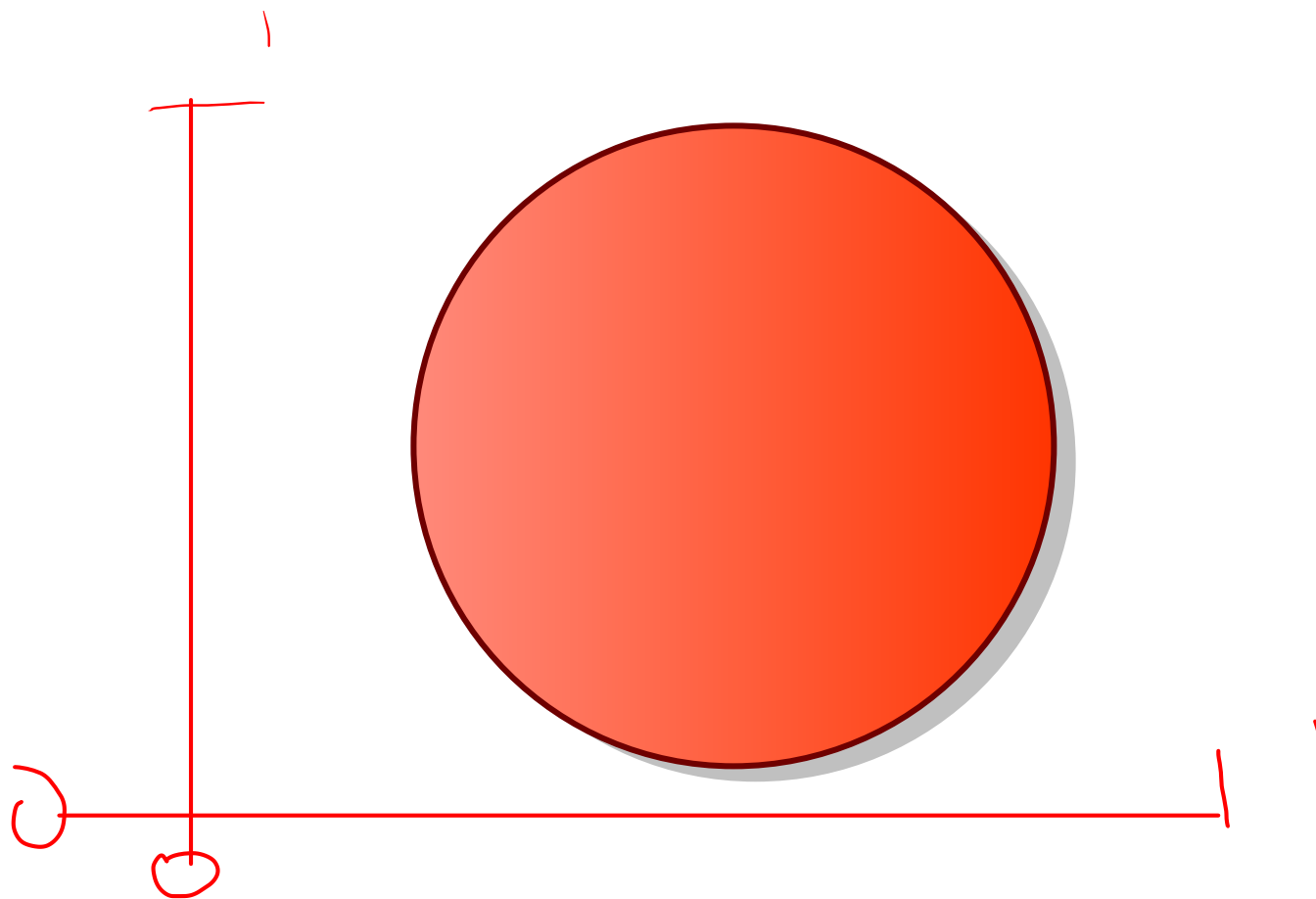| Age of Student | Class of Student |
|---|---|
| 6 | 1 |
| 16 | 10 |
| 30 | 18 |
| 50 | 20 |
| 70. 80. | PhD |

Ranges from 0-80

50

0

0 — 30

What if we put any b/w 0-1 & does also b/w 0 & 1. Then

Q) How do we scale the data?

Ans) There are 2 scaling methods.

1) Standard scaling

2) Min - max scaling.

# Standard Scaling :->

Properties :-> It scales the data such that **mean is zero**

2) All the data lies b/w 1 standard deviation. ( b/w $1\sigma$ )

$$x_{i_{new}} = \frac{x_i - \mu}{\sigma} \begin{cases} \mu = \text{mean} \\ \sigma = \text{standard deviation} \end{cases}$$

$$x_{i,new} = \frac{x_i - \mu}{\sigma}$$

When $\mu = \frac{1}{n} \sum_{i=1}^{n} x_i$

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (x_i - \mu)^2}$$

Also called :- Standardization

Min - max Scaling $\Longrightarrow$

Also called Normalization.

$$x_{i, new} = \frac{x_i - min}{max - min}$$

The $x_{i, new}$ will always be
b/w 0 & 1

Training a model means
to find the best possible
parameters. To minimize lost
function like MSE.
If we have more features. Then
algorithm will have to find
more parameters. and the search
becomes complex.

1) Batch Gradient Descent.

2) Stochastic Gradient descent.

3) Mini - Batch " "

Partial Derivative -

1) We initialize all parameters random-ly.

2) Than we change 1 weights while keeping all others fixed & calculate change in cost function

$$\frac{\partial}{\partial \theta_j}(MSE) = \frac{2}{m}(\theta^T \cdot \chi^{(i)} - y^{(i)}) * \chi_j^{(i)}$$

$$\nabla_\theta MSE(\theta) = \begin{pmatrix} \frac{\partial}{\partial \theta_0} MSE(\theta) \\ \frac{\partial}{\partial \theta_1} MSE(\theta) \\ \vdots \\ \frac{\partial}{\partial \theta_n} \end{pmatrix}$$

$$MSE(\theta) = \frac{2}{m} X^T (X \cdot \theta - y)$$

$$\nabla_\theta MSE(\theta) = \frac{2}{m} X^T (X \cdot \theta - y)$$

1) Batch gradient descent

will load all the data in every

step.

So Batch gradient descent is
very slow on large datasets.

But in

$$\theta^{(\text{new step})} = \theta - \eta \nabla_\theta MSE(\theta)$$

$$\llcorner \text{Learning rate}.$$