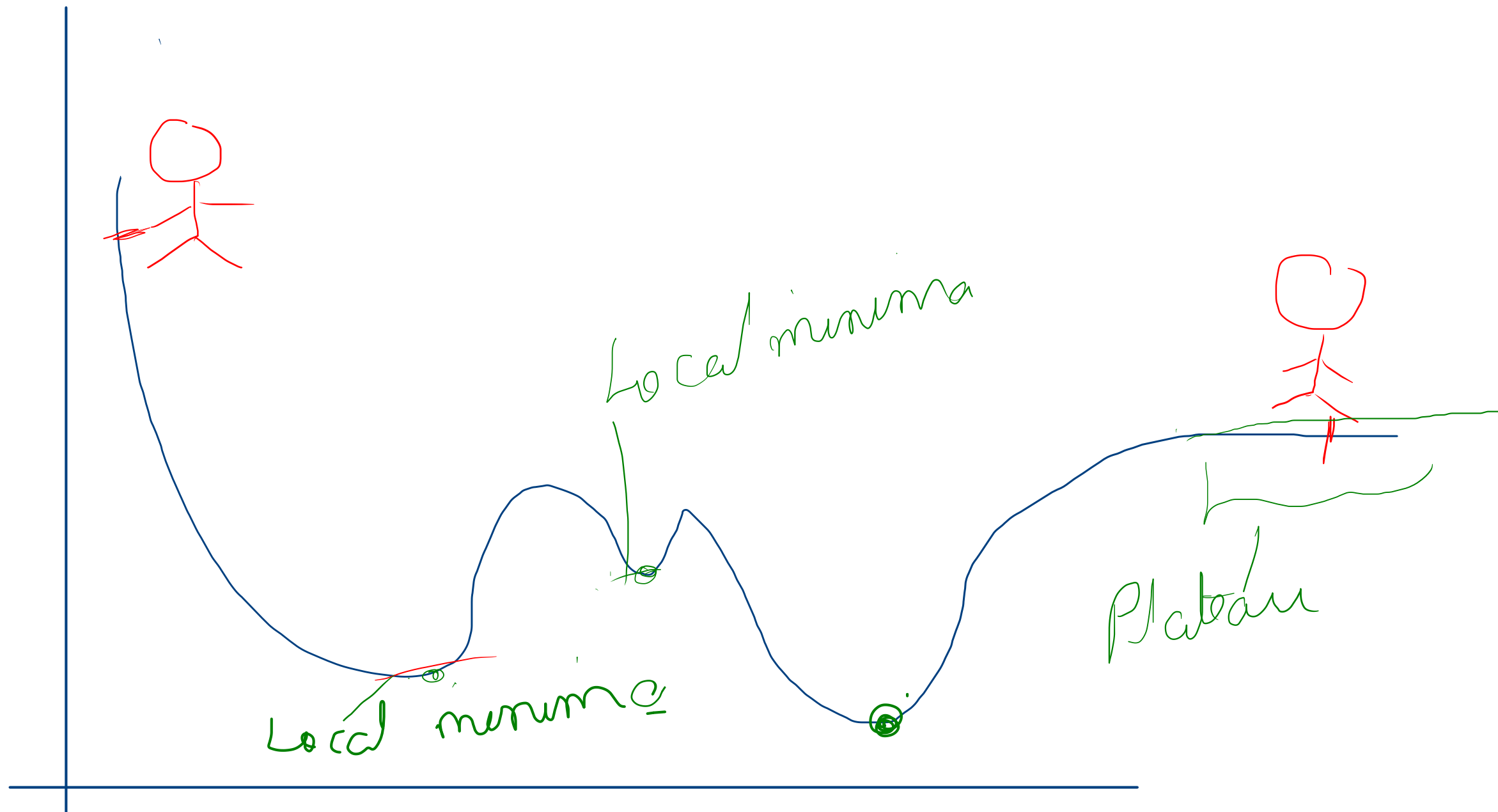


# Starting understanding gradient descent.



# Types of gradient descent $\rightarrow$

1) Batch

2) Stochastic

3) Mini-Batch.

Gradient Descent.

But before this we have  
more topics to discuss

What would happen if our data is on different scale - or features have different scales?

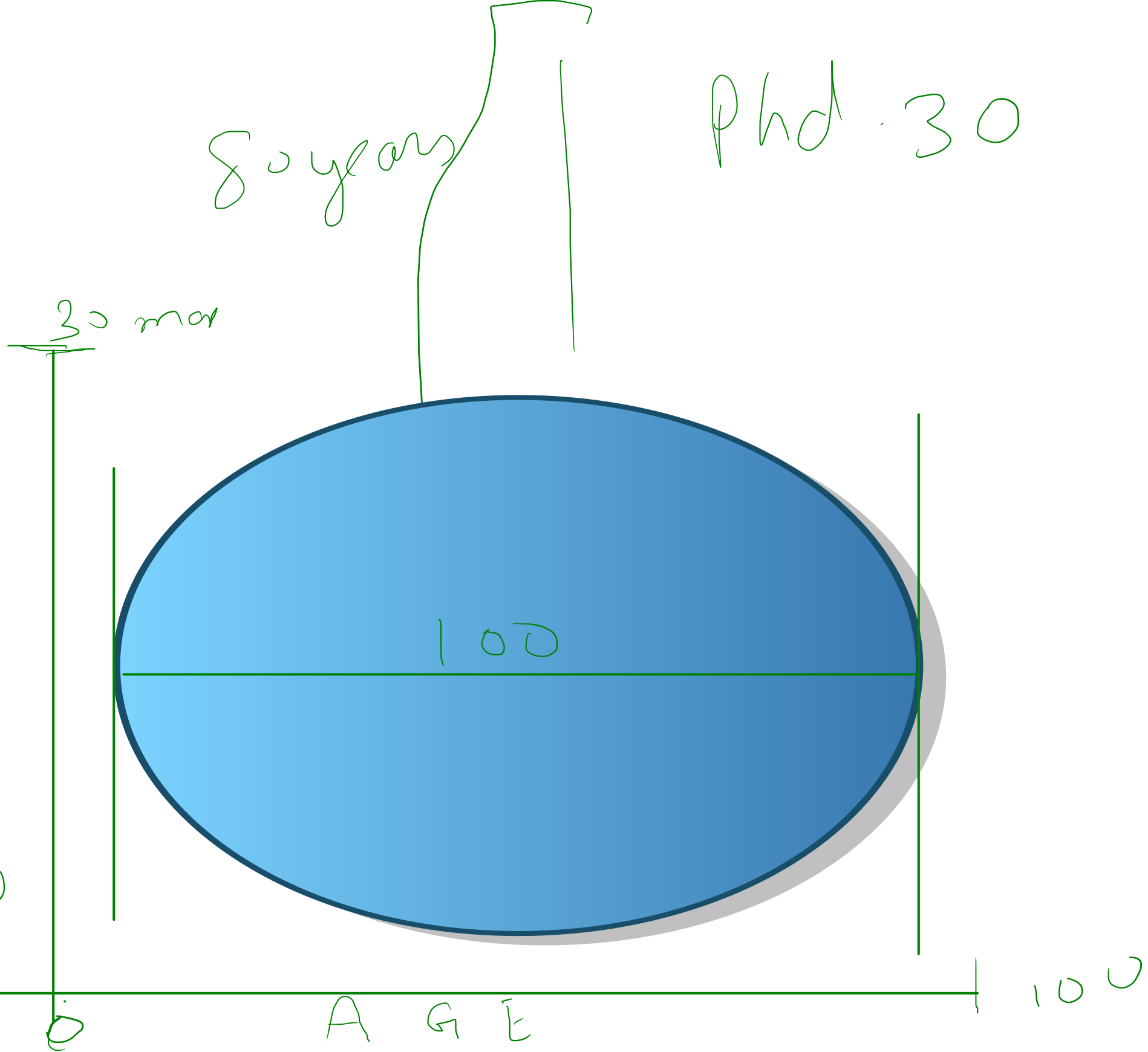
Student Age		Student class	
50	6	10	10
	16	15	10
	20	30	10

80 years

Phd. 30

30 mar

C  
L  
A  
S  
S  
D



2 ways by which we can  
scale our data.

- 1) Standard Scaler / Standardization
- 2) Min-max Scaling / Normalization

## Star normalization:

$$x_{\text{new}} = \frac{x - \mu}{\sigma}$$

where  $\mu = \text{mean}$

$\sigma = \text{standard deviation}$

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

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

Standardization will scale our data such that mean of data would be zero. Standard deviation will be 1 for all the data.

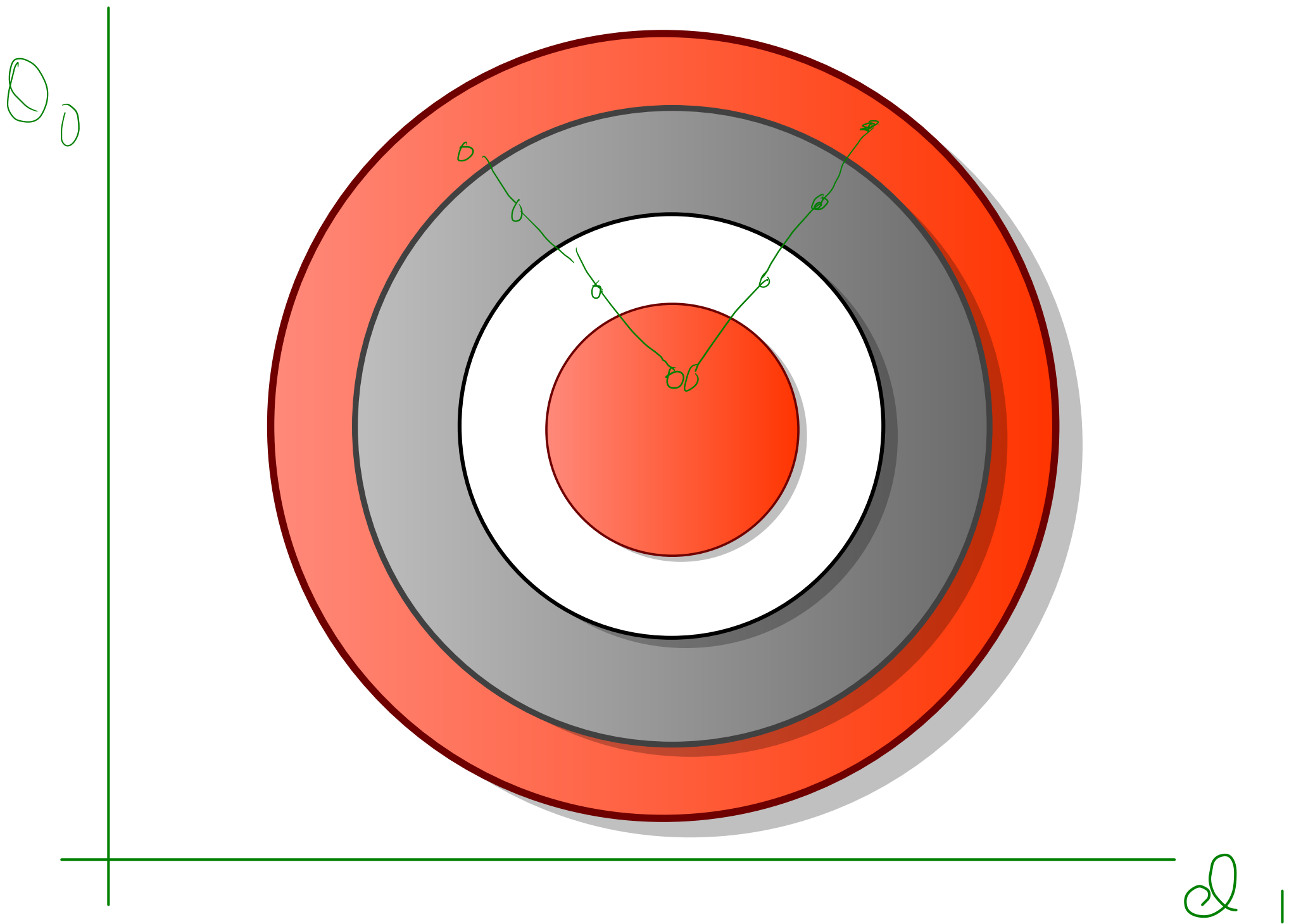
Normalization: - or Min max scaling

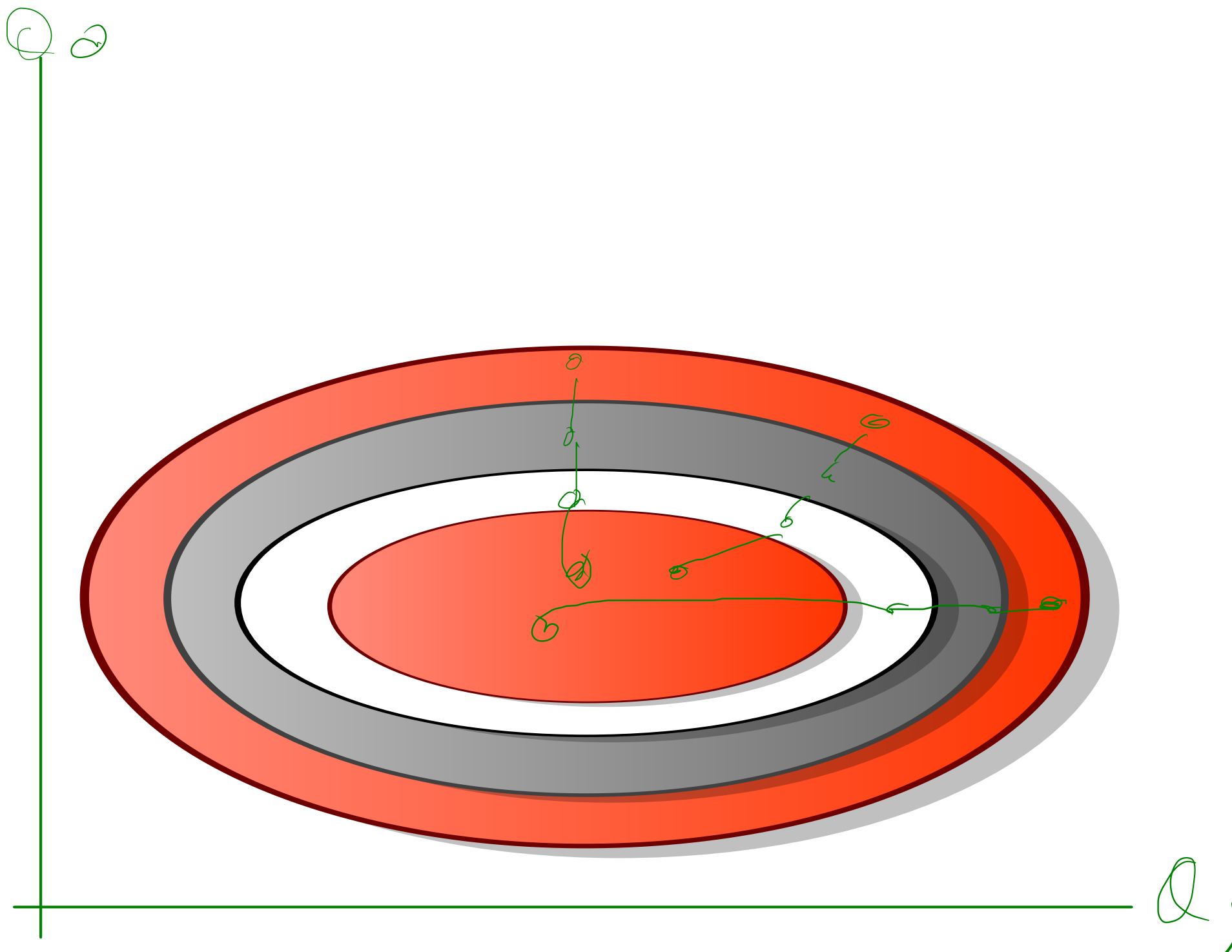
$$x_{i \text{ new}} = \frac{x_i - \min}{\max - \min}$$

The data will range b/w - 0 & 1.

Now we will see the effect of scaling on gradient descent.

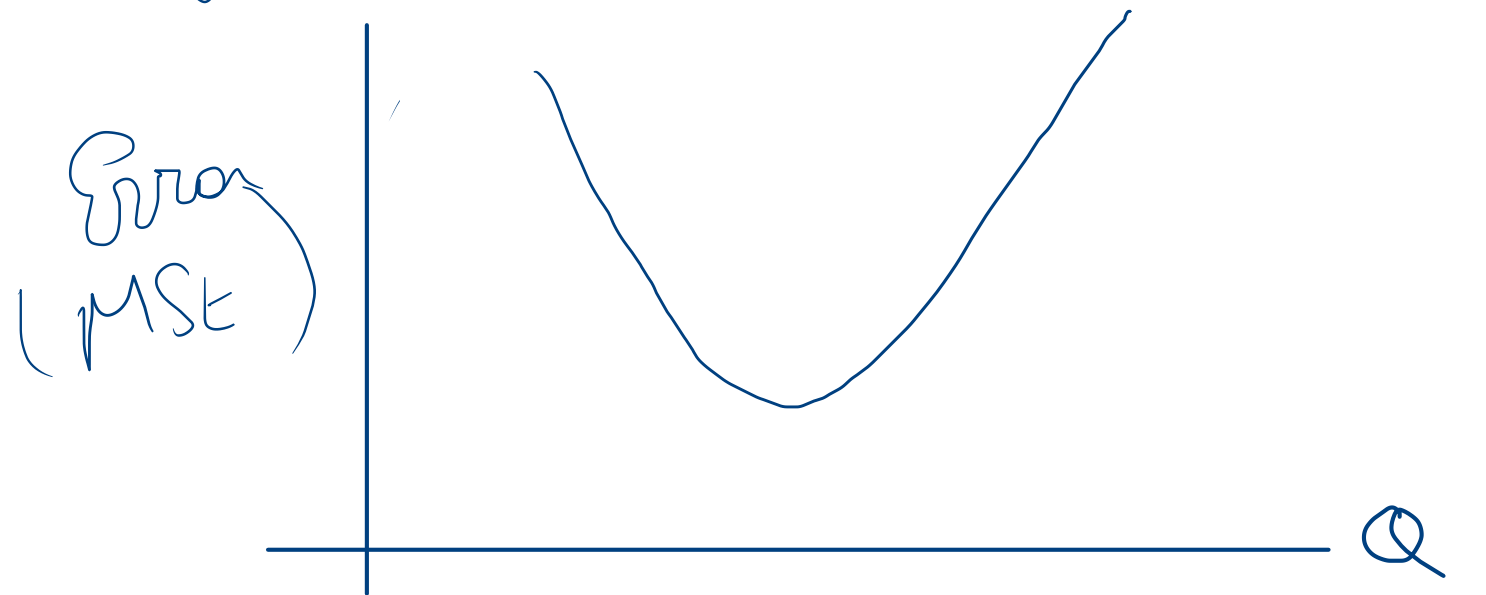






What do we mean by training a  
model?

We are searching for best parameters  
for  $\theta_0, \theta_1, \dots, \theta_n$  To minimize  
Cost function (like MSE)



Model becomes complex as the number of features increases -

Types of Gradient descent

1) Batch gradient descent.

2) Stochastic

3) Mini batch

"

"

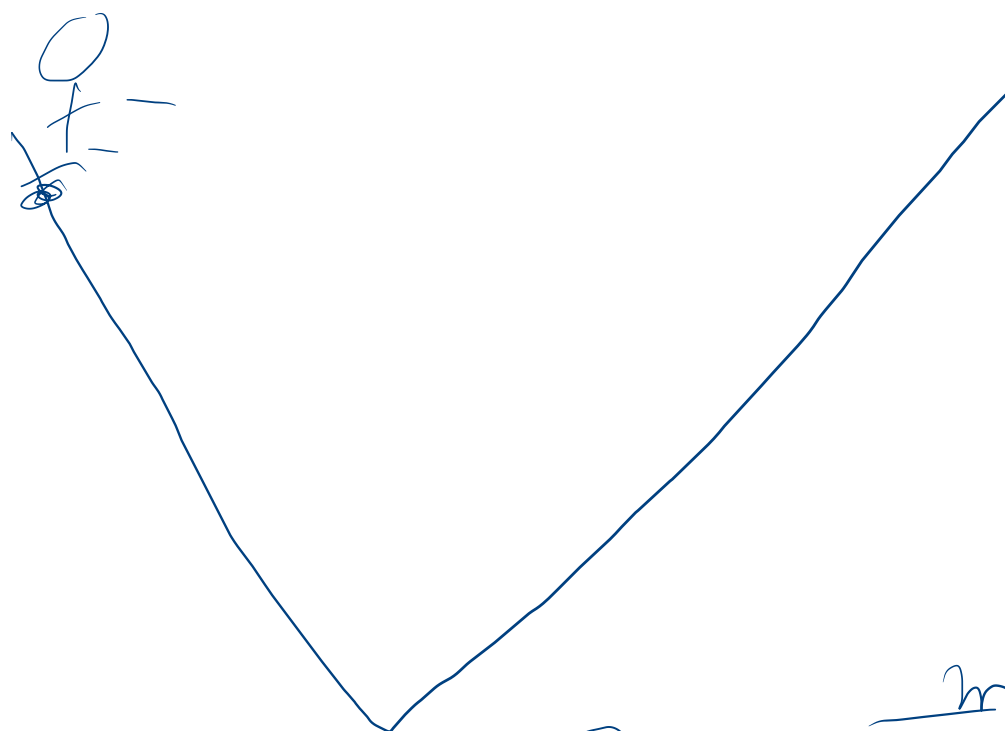
"

"

# Partial Derivatives:

First we initialize model parameter  
 $(\theta_1, \theta_2, \dots, \theta_n)$  randomly.

For each weight we change one  
weight ~~wt~~ & find all other. Then we  
check change in Cost function.  
This is called partial derivatives.



$$\frac{\partial}{\partial \theta_j} (MSE(\theta)) = \frac{2}{m} \sum_{i=1}^m (\theta^T \cdot x^{(i)} - y^{(i)})$$

$x_j^{(i)}$

gradient vector  $\nabla_{\theta} \text{MSE}(\theta)$

$$= \begin{pmatrix} \frac{\partial}{\partial \theta_0} \text{MSE}(\theta) \\ \frac{\partial}{\partial \theta_1} \text{MSE}(\theta_1) \dots \\ \frac{\partial}{\partial \theta_m} \text{MSE}(\theta_m) \end{pmatrix}$$

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

---

Since Batch gradient descent is  
using complete dataset to every  
step. It is really slow on  
training. But still faster than

linear regression.

Gradient descent (step)

$$\theta_{(new)} = \theta - \eta \nabla_{\theta} \text{MSG}(\theta)$$

$\theta$  Learning rate.