

Machine Learning crash course

Ahmed Hammad

High Energy Accelerator Research Organization (KEK), Japan.

4th summer school at CTP, BUE

References

- 1- Hands-on machine learning with Sikit-Learn and Tensorflow. By Aurélien Géron
- 2- Probabilistic machine learning. By Kevin P.Murphy
- 3- Machine learning with Pytorch and scikit-Learn. By Sebastian Raschka, etal

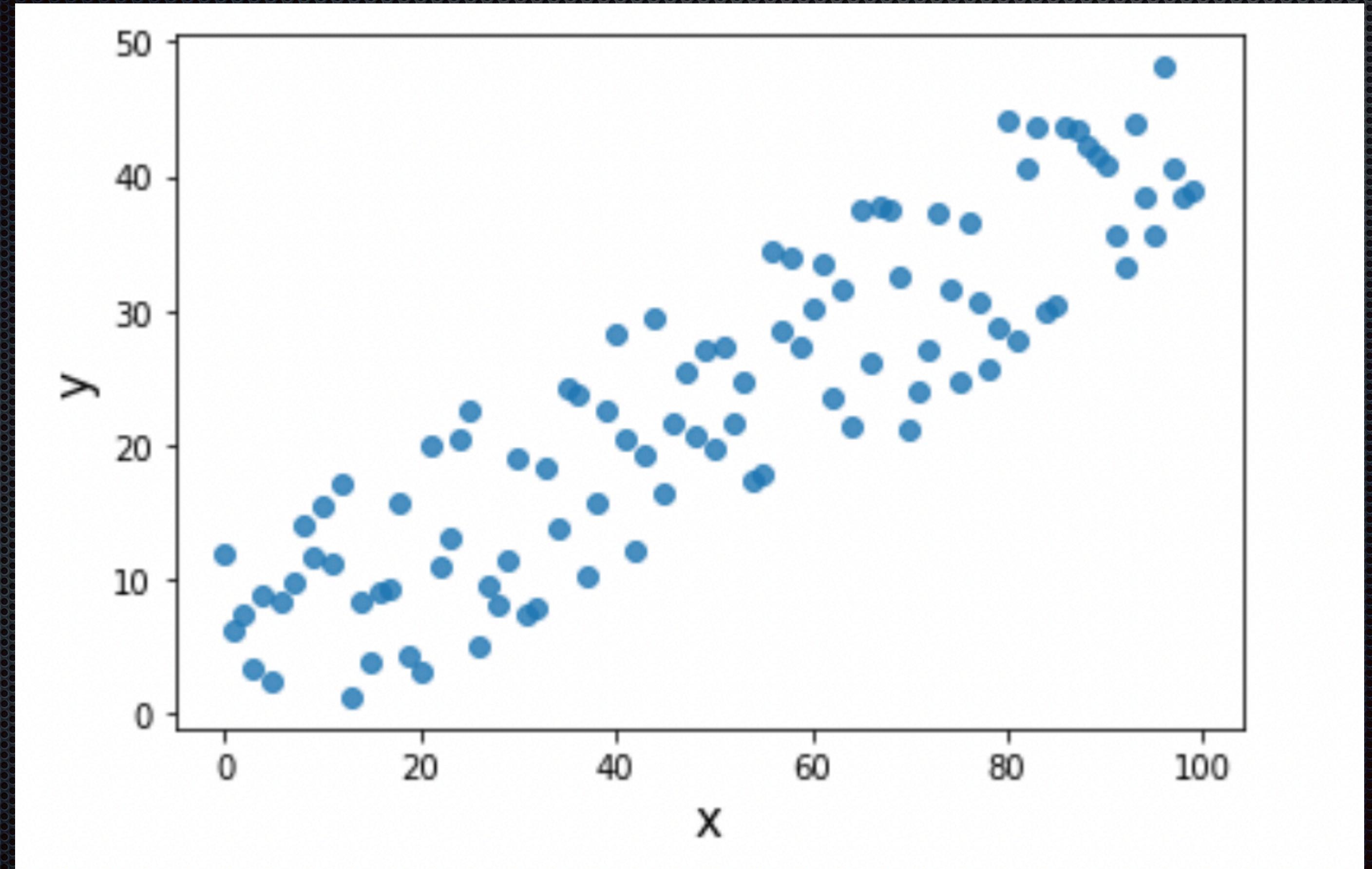
So, what is machine learning ?

$$Y = F(X)$$

Machine learning models approximate the function F for any type of data X

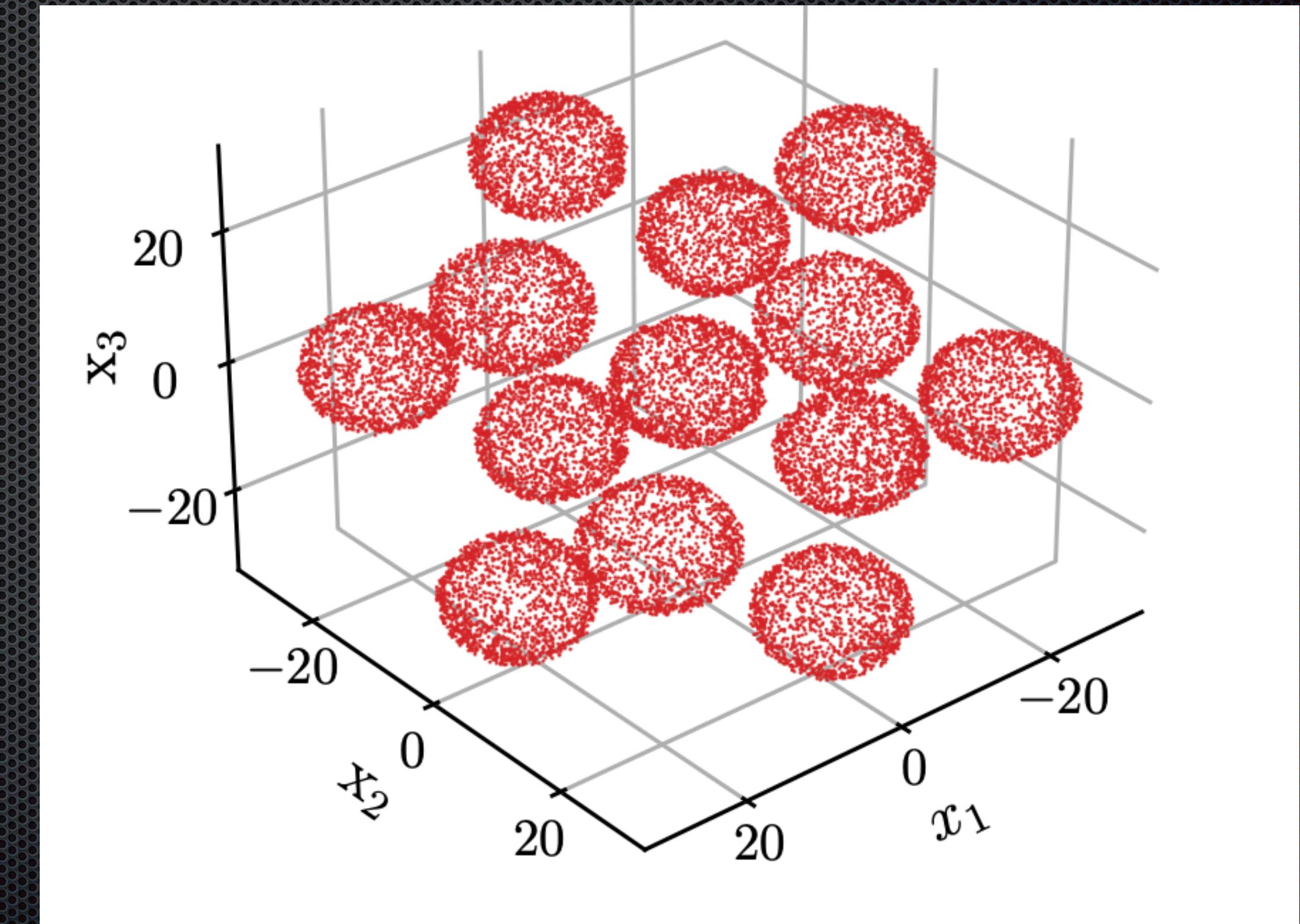
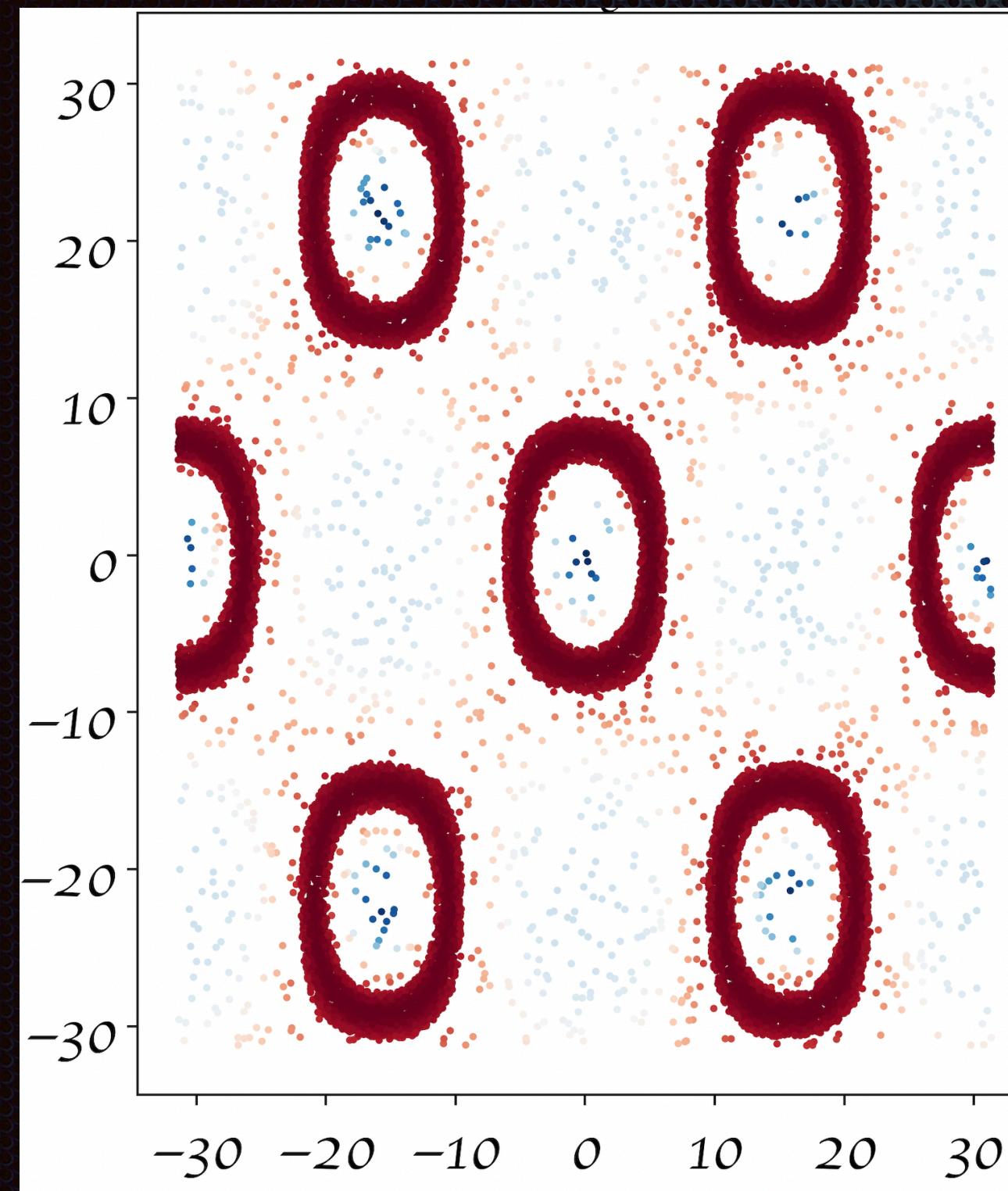
Introduction

Approximated function can be very simple: One dimension-Linear function



Introduction

Low dimensions - Non-linear function



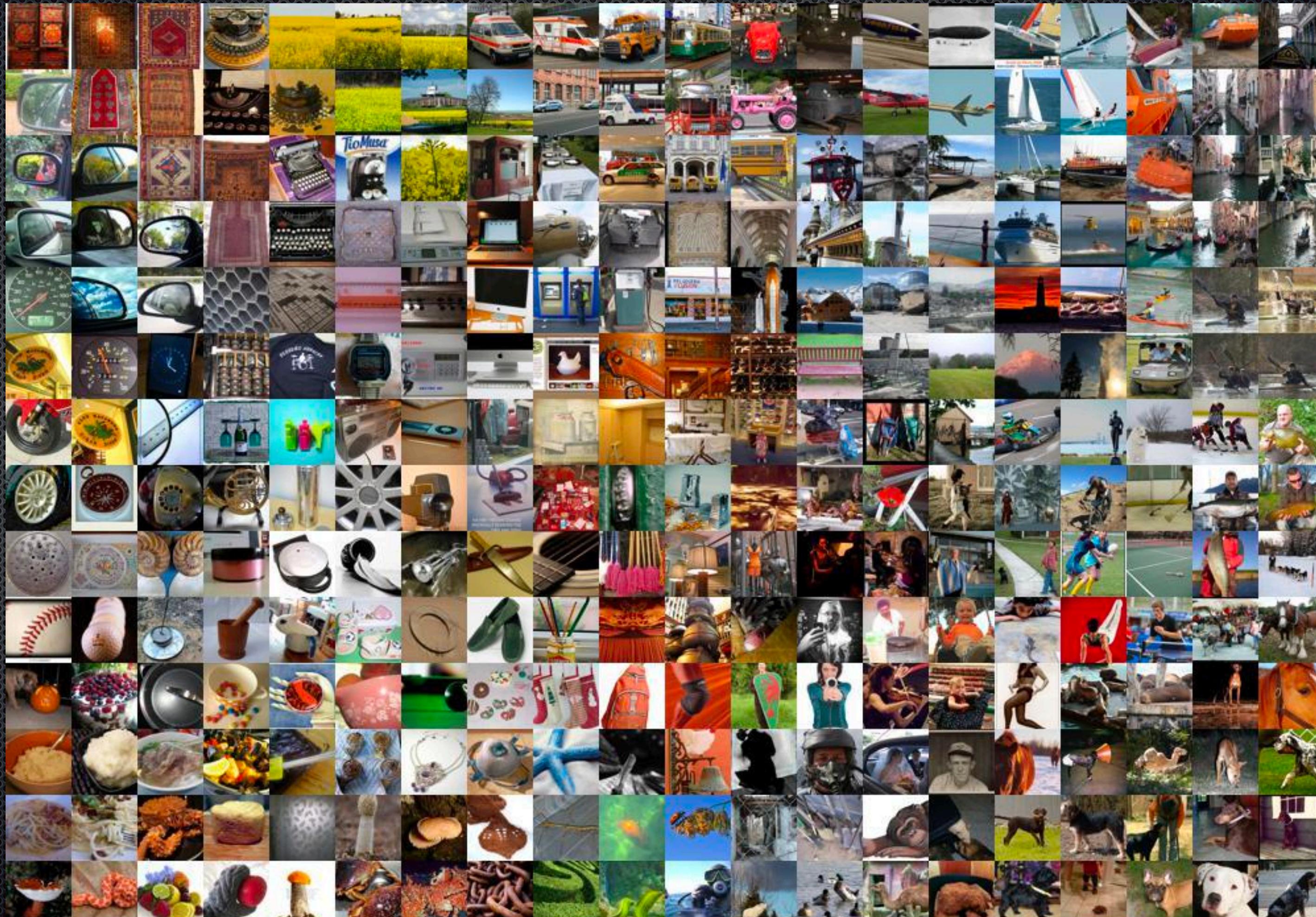
Introduction

ImageNet dataset

Size: 167.6 GB

High dimensions - Non-linear function

$$d = (100000, 225, 225, 3)$$



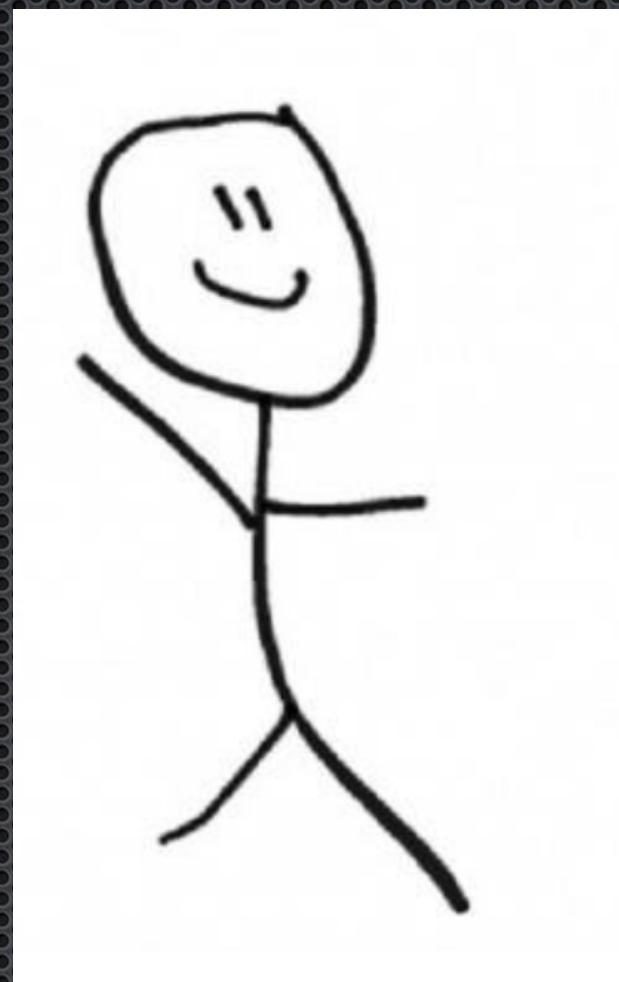
Introduction

How does machine learning approximate F ?

Inputs



Data scientist



0
10
20
30

Bunch of numbers

3	1	4	1
5	9	2	6
5	3	5	8
9	7	9	3
2	3	8	4
6	2	6	4

2	1	8	8	1
2	8	5	0	5
2	4	9	4	5
2	3	6	2	8
7	7	1	5	2

Find a pattern



New data

Prediction

Multi-Layers perceptron

	Epoch 000,313	Learning rate 0.0001	Activation ReLU	Regularization None	Regularization rate 0	Problem type Classification
--	------------------	-------------------------	--------------------	------------------------	--------------------------	--------------------------------

DATA

Which dataset do you want to use?



Ratio of training to test data: 50%

Noise: 20

Batch size: 6

REGENERATE

FEATURES

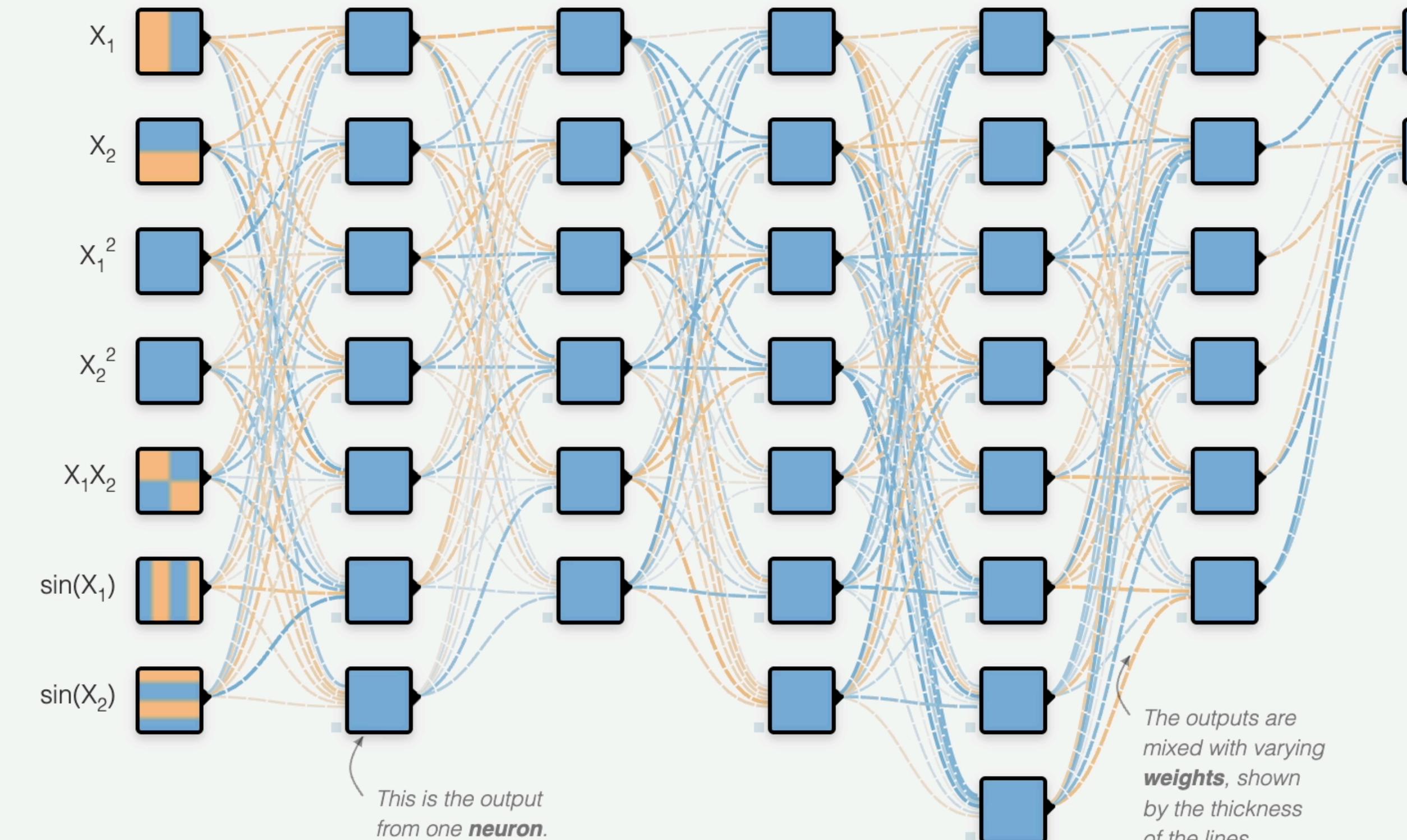
Which properties do you want to feed in?

x_1 x_2 x_1^2 x_2^2 x_1x_2 $\sin(x_1)$ $\sin(x_2)$

6 HIDDEN LAYERS

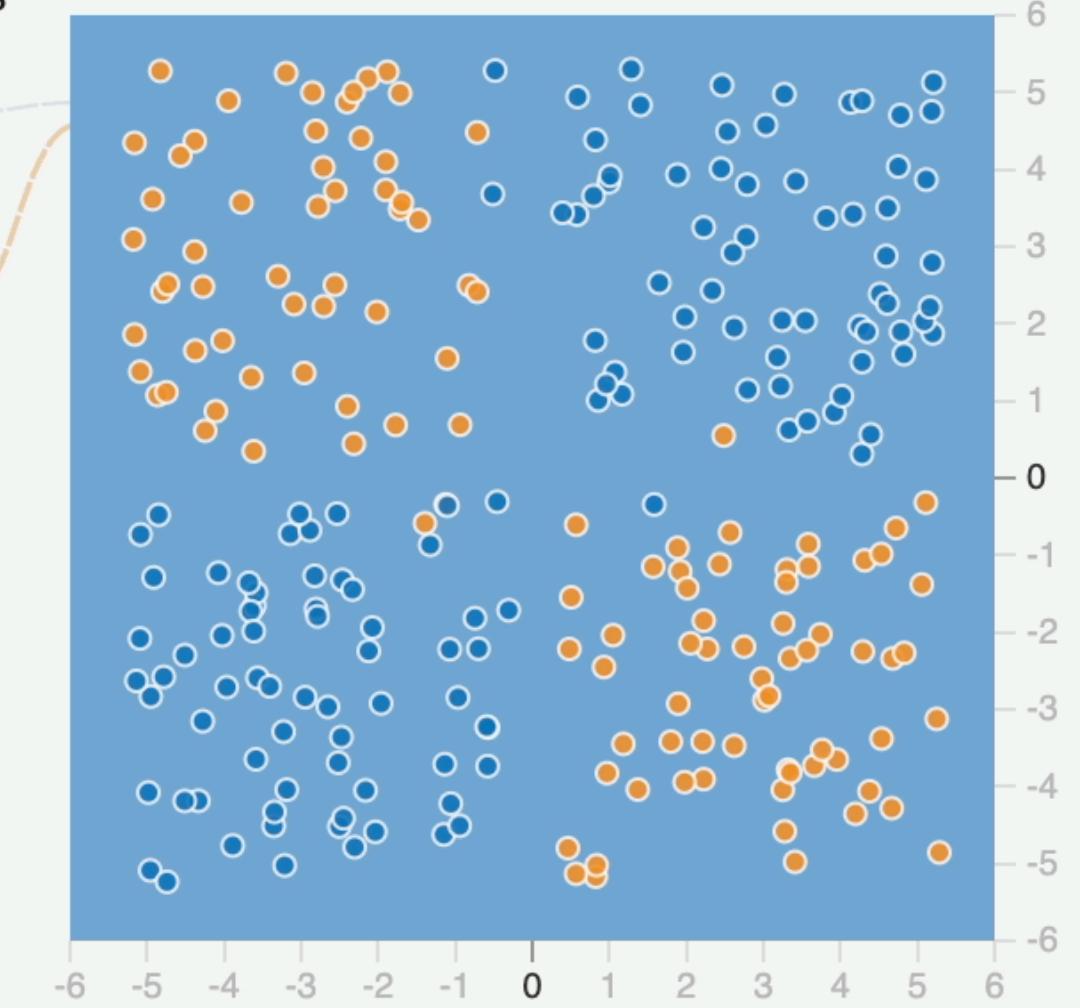
+ - + - + - + - + - + -

7 neurons 6 neurons 7 neurons 8 neurons 6 neurons 2 neurons



OUTPUT

Test loss 0.495
Training loss 0.495



Show test data Discretize output

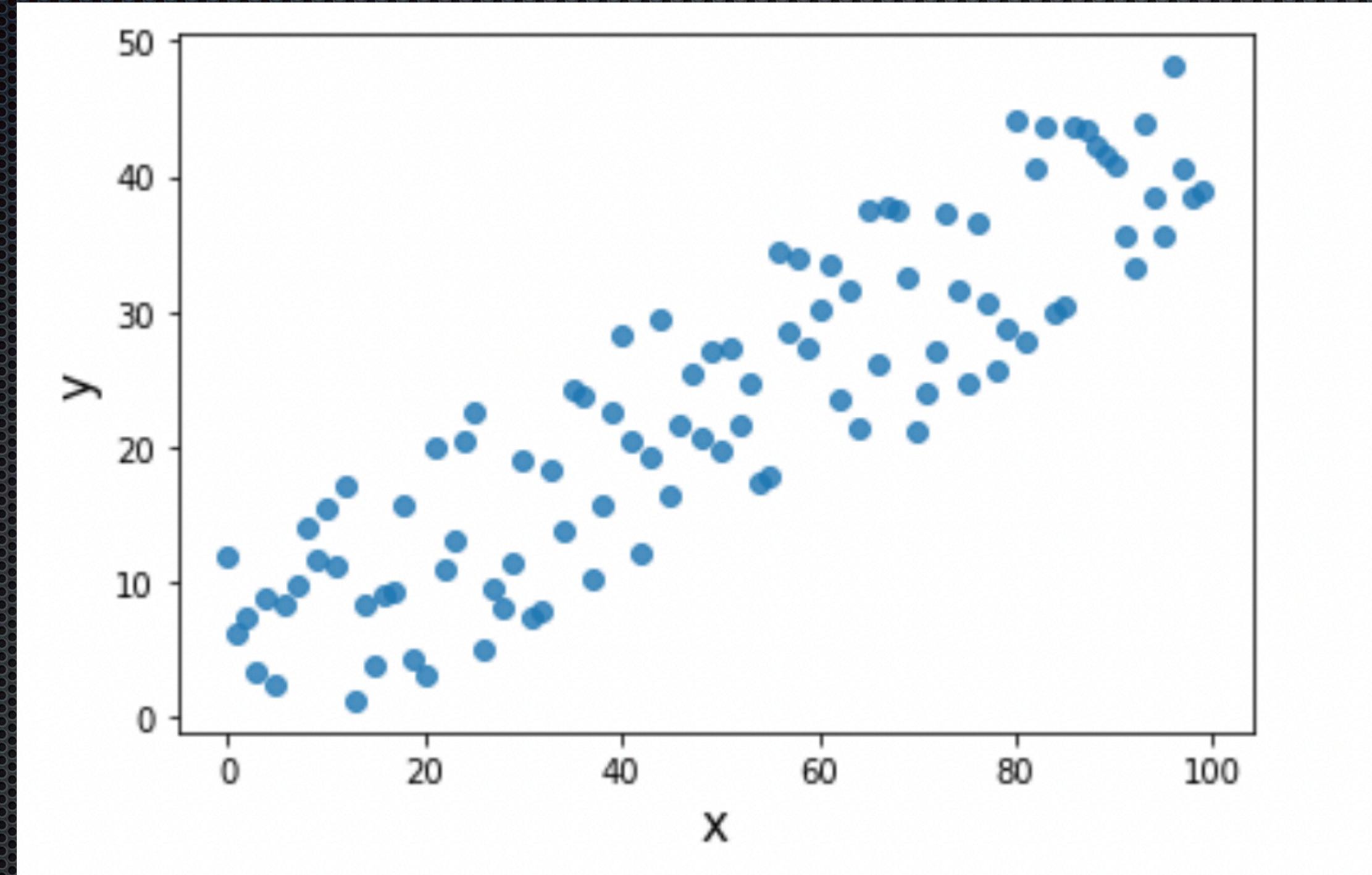
The outputs are mixed with varying **weights**, shown by the thickness of the lines.

Multi-Layers perceptron

Well, seems good! But how it works ?

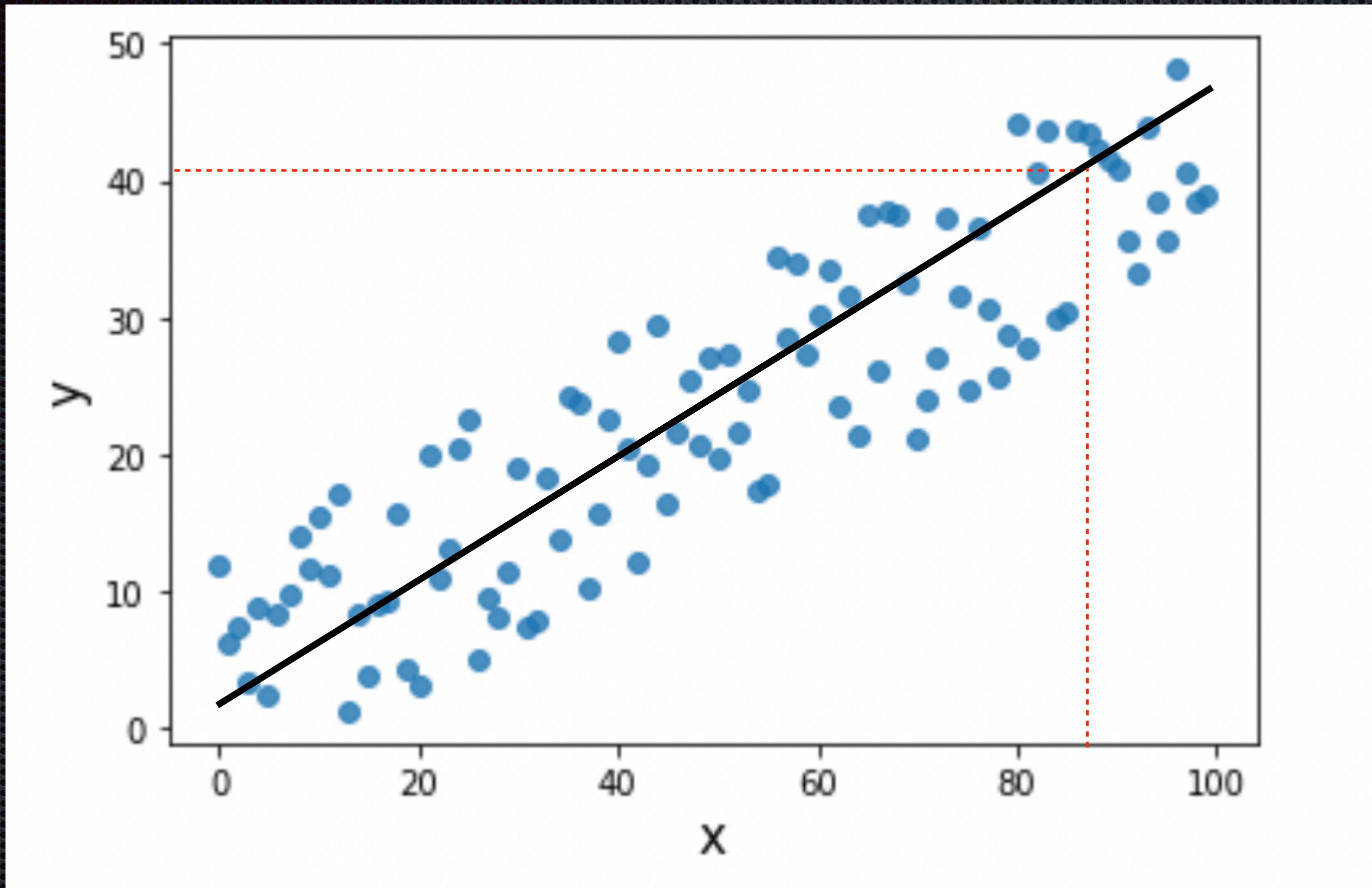
Linear regression

For a continuous data set X



Given new (unseen) x point, can you predict the corresponding Y ?

Linear regression

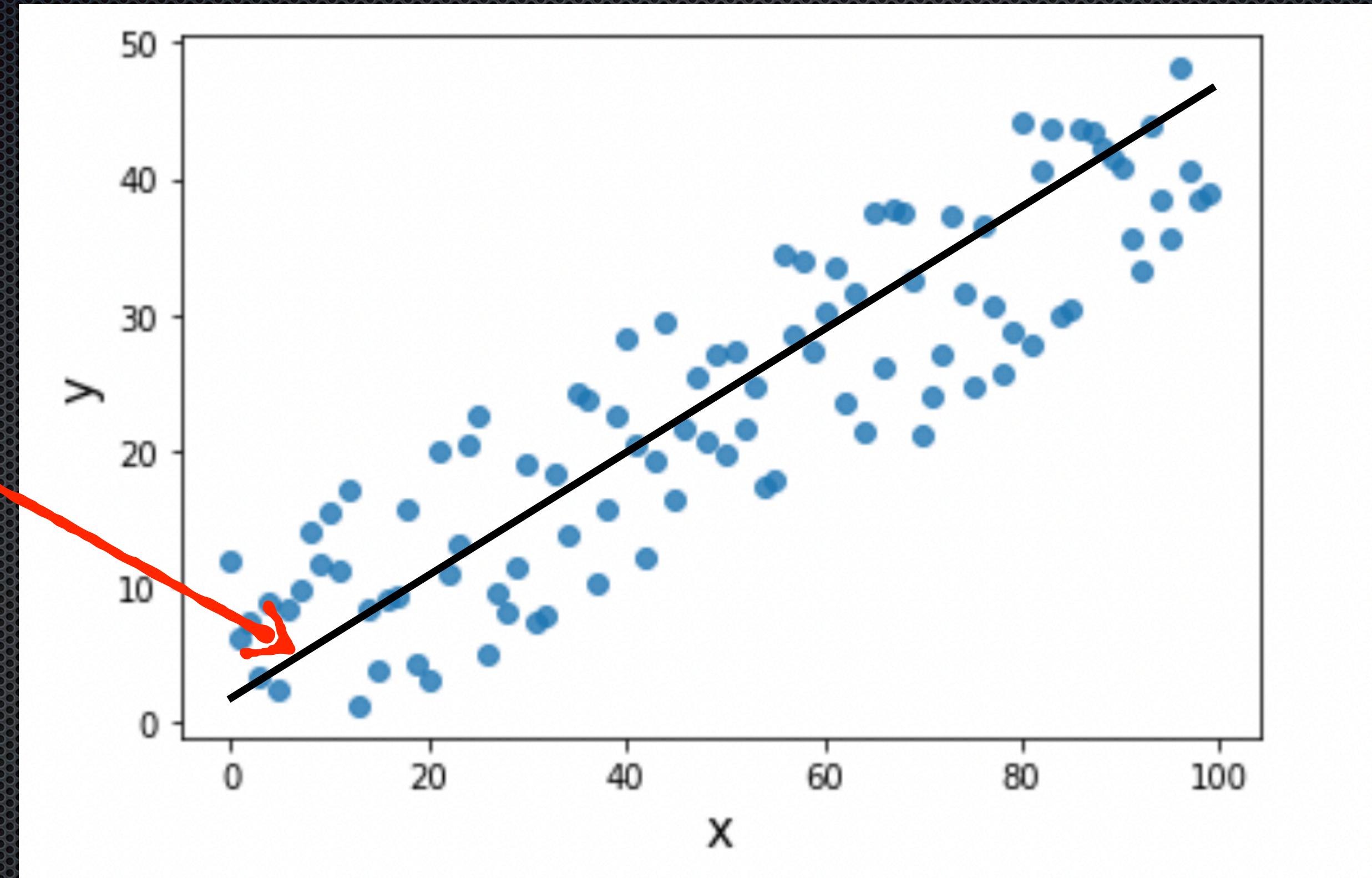


Once we found the line that best fit the data we can interpret the **Y** for new **x** points

Linear regression

How to find the line that fit the given data ?

$$\hat{Y} = B_0 + B_1 x$$



With given \mathbf{Y} and \mathbf{x} we try to find B_0 and B_1 that fit the line to the data

Linear regression

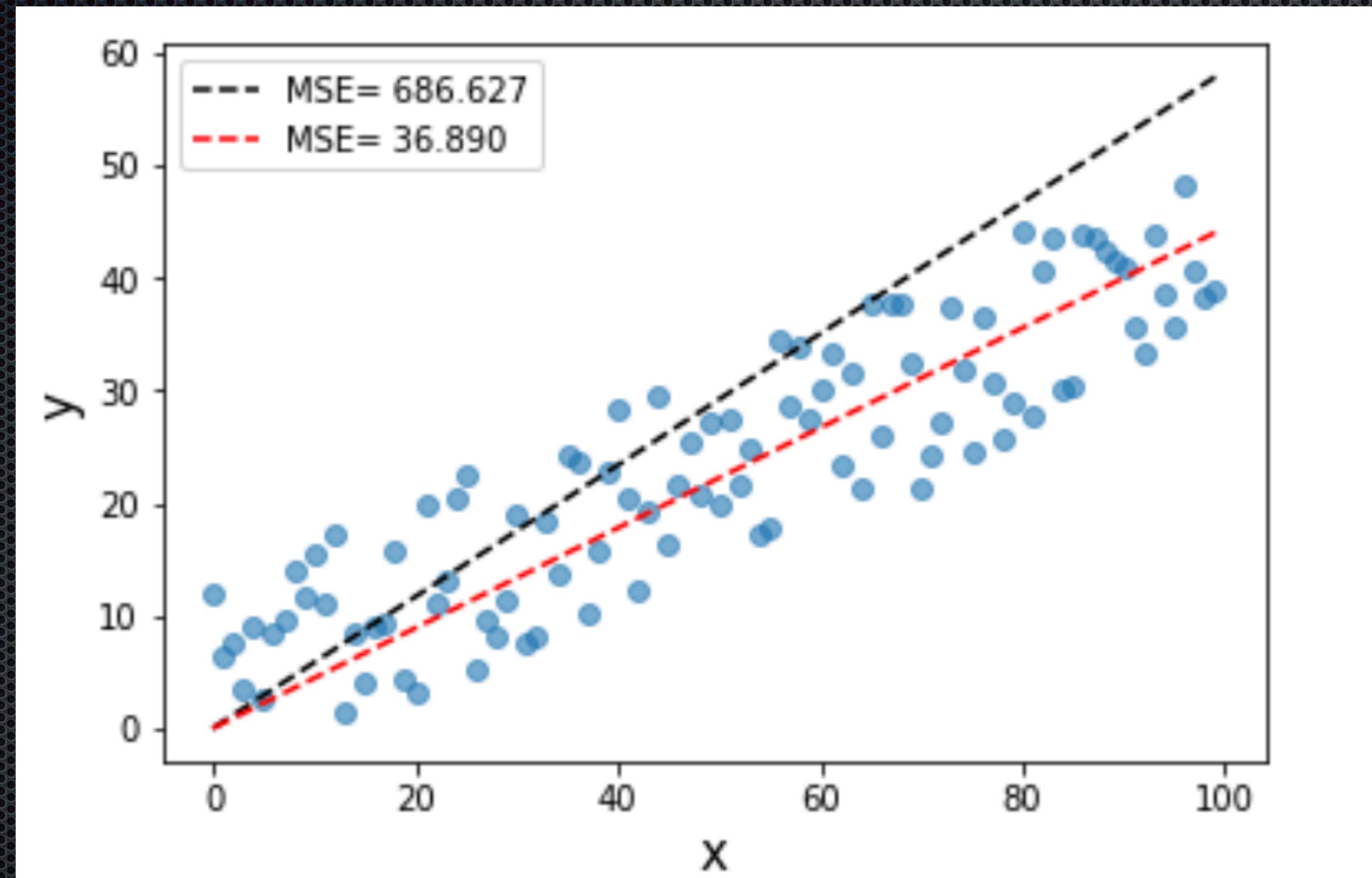
Question:

We can draw infinite number of lines that can fit the data!!
Which line shall we consider ?

Linear regression

Question:

We can draw infinite number of lines that can fit the data!!
Which line shall we consider ?

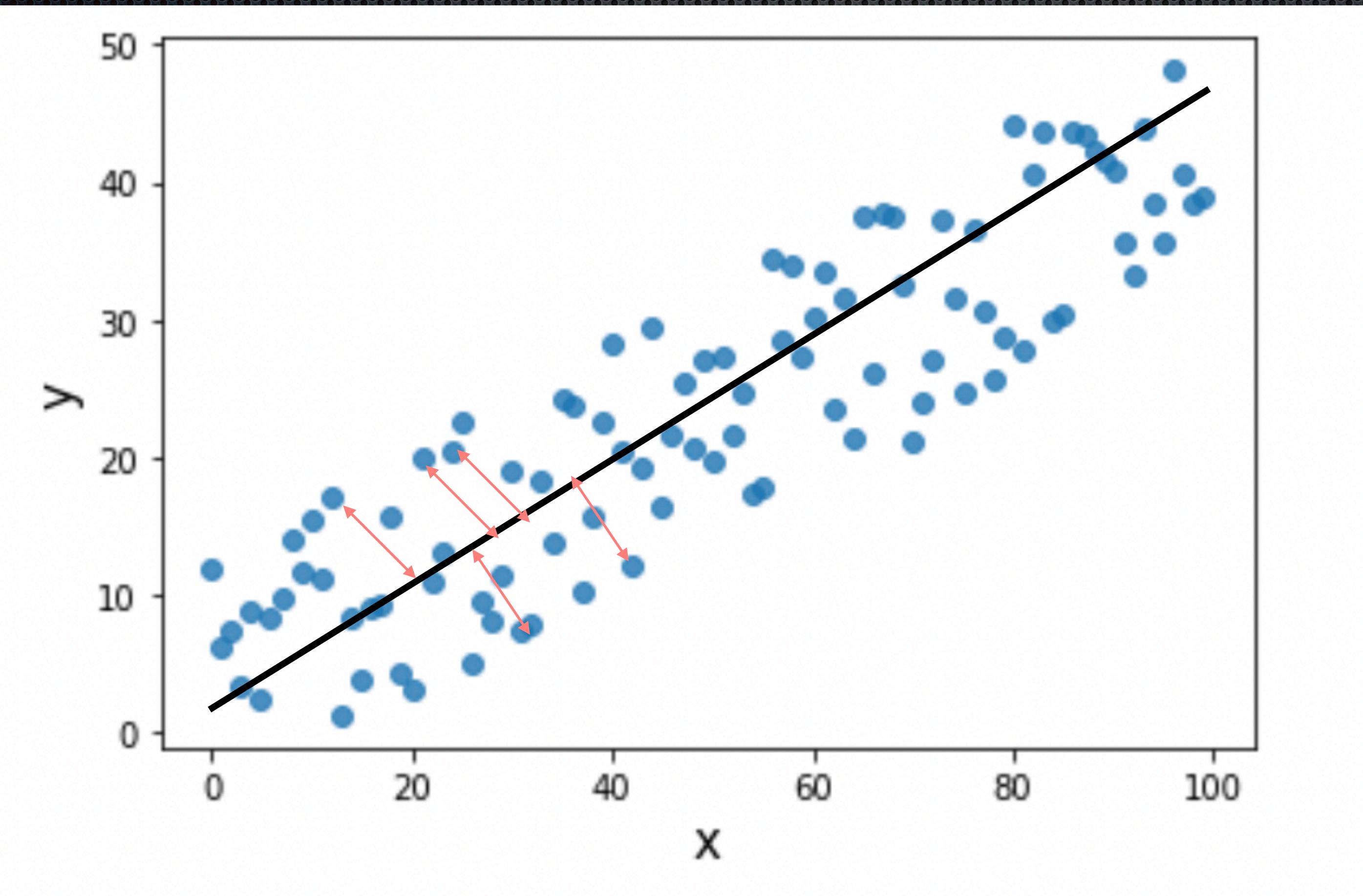


We have to consider the line that minimizes the **error function**

Linear regression-loss function

Loss function:

It mostly quantifies the difference between the model predictions and the true values



Linear regression-loss function

Types of linear regression loss functions

Absolute mean error:

$$\text{AME} = \frac{1}{m} \sum_{i=1}^m |Y_i - \hat{Y}_i| = \frac{1}{m} \sum_{i=1}^m |Y_i - (B_0 + B_1 x_i)|$$

Mean square error:

$$\text{MSE} = \frac{1}{m} \sum_{i=1}^m |Y_i - \hat{Y}_i|^2 = \frac{1}{m} \sum_{i=1}^m |Y_i - (B_0 + B_1 x_i)|^2$$

Root mean square error:

$$\text{RMSE} = \sqrt{\frac{1}{m} \sum_{i=1}^m |Y_i - \hat{Y}_i|^2} = \sqrt{\frac{1}{m} \sum_{i=1}^m |Y_i - (B_0 + B_1 x_i)|^2}$$

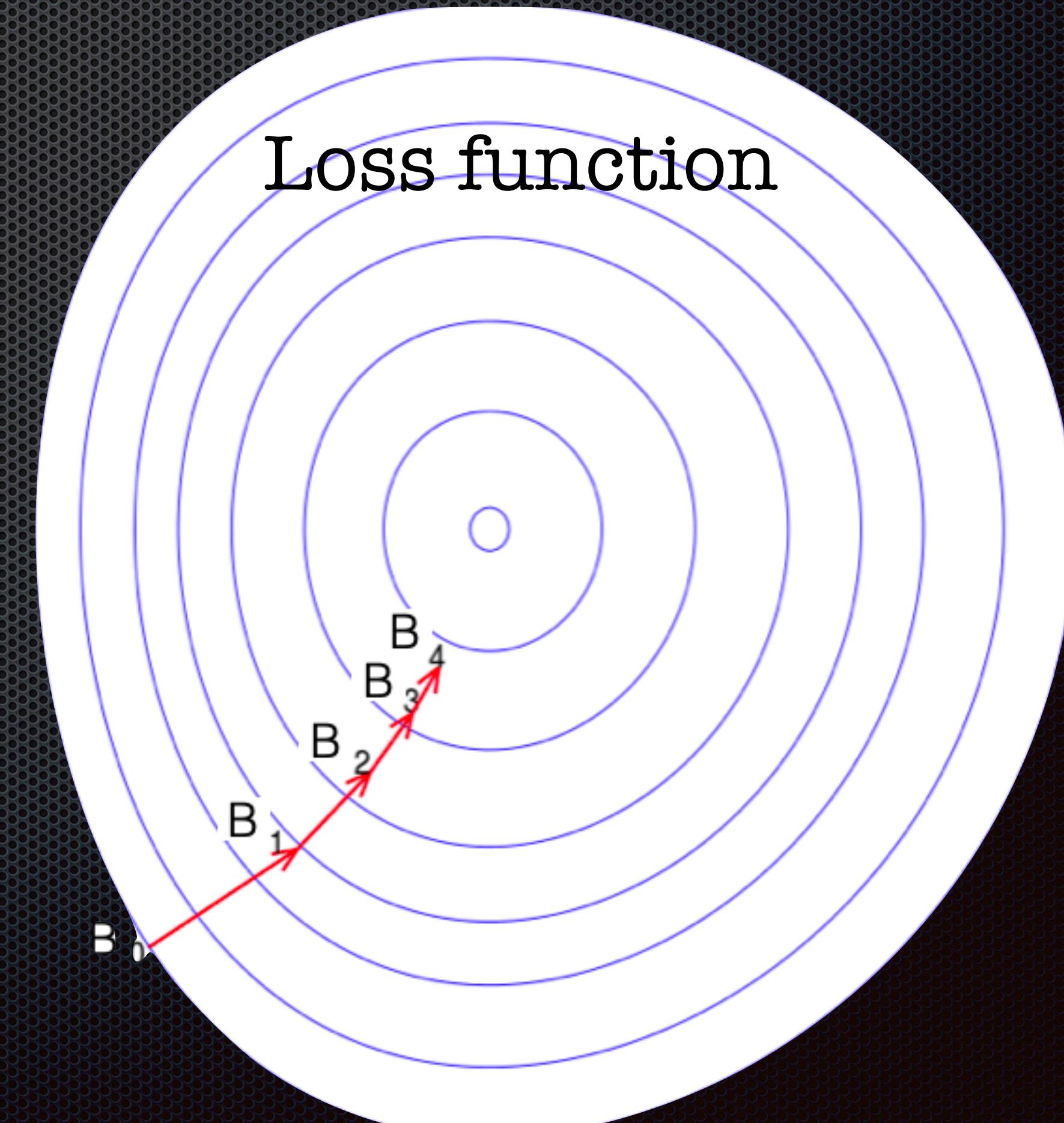
Gradient Descent method

Question:

Well, now we know that to find the best fit line we need to minimize the loss function, but how ??

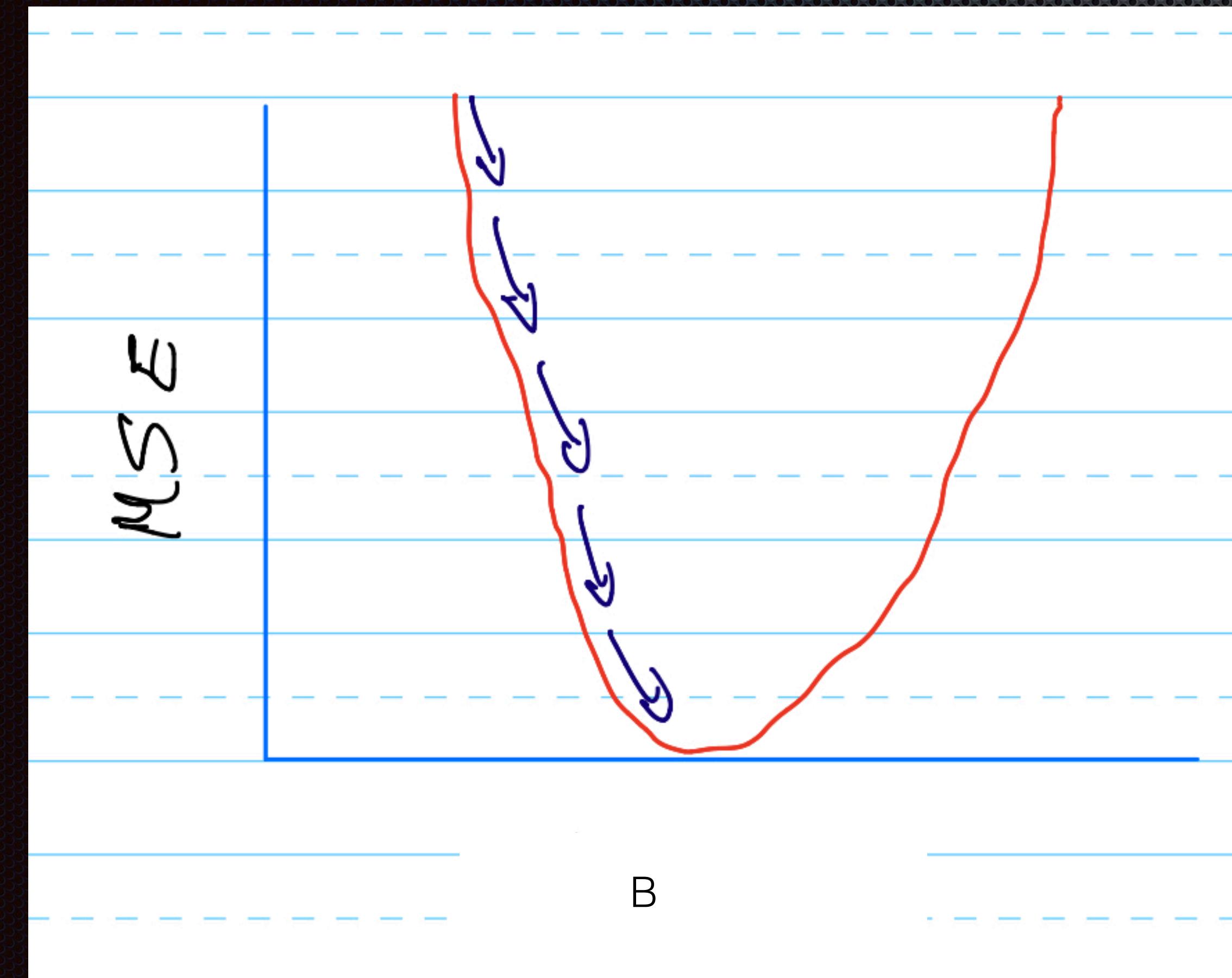
Find the best fit parameters using the gradient descent method

$$B_{\text{new}}^i = B_{\text{old}}^i - \eta \nabla \text{Loss}(B_{\text{old}}^i)$$



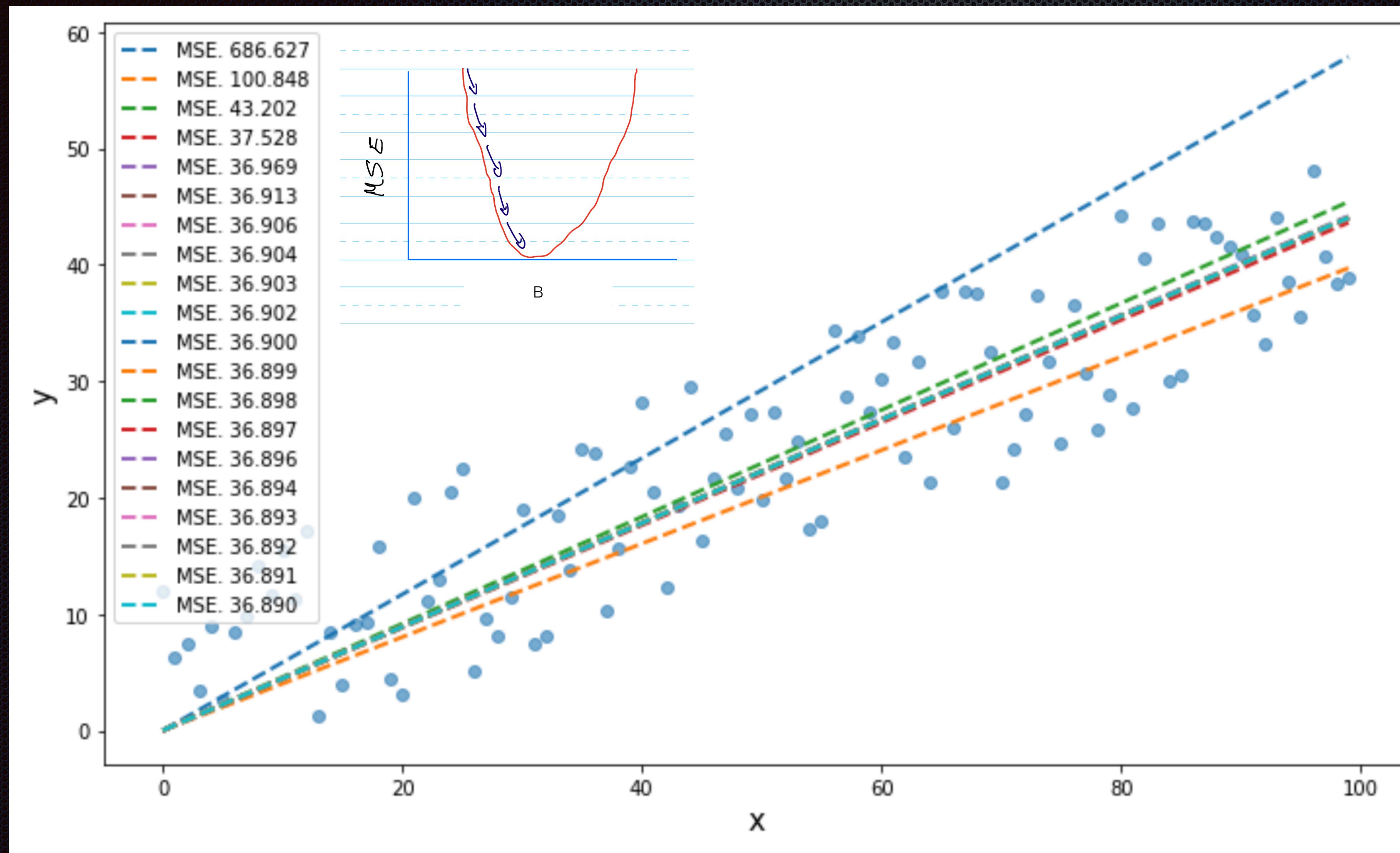
Gradient Descent method

Given a number of iteration, every iteration we update the parameters and calculate the error function until we hit the global minimum



$$B_{\text{new}}^i = B_{\text{old}}^i - \eta \nabla \text{MSE}(B_{\text{old}}^i)$$

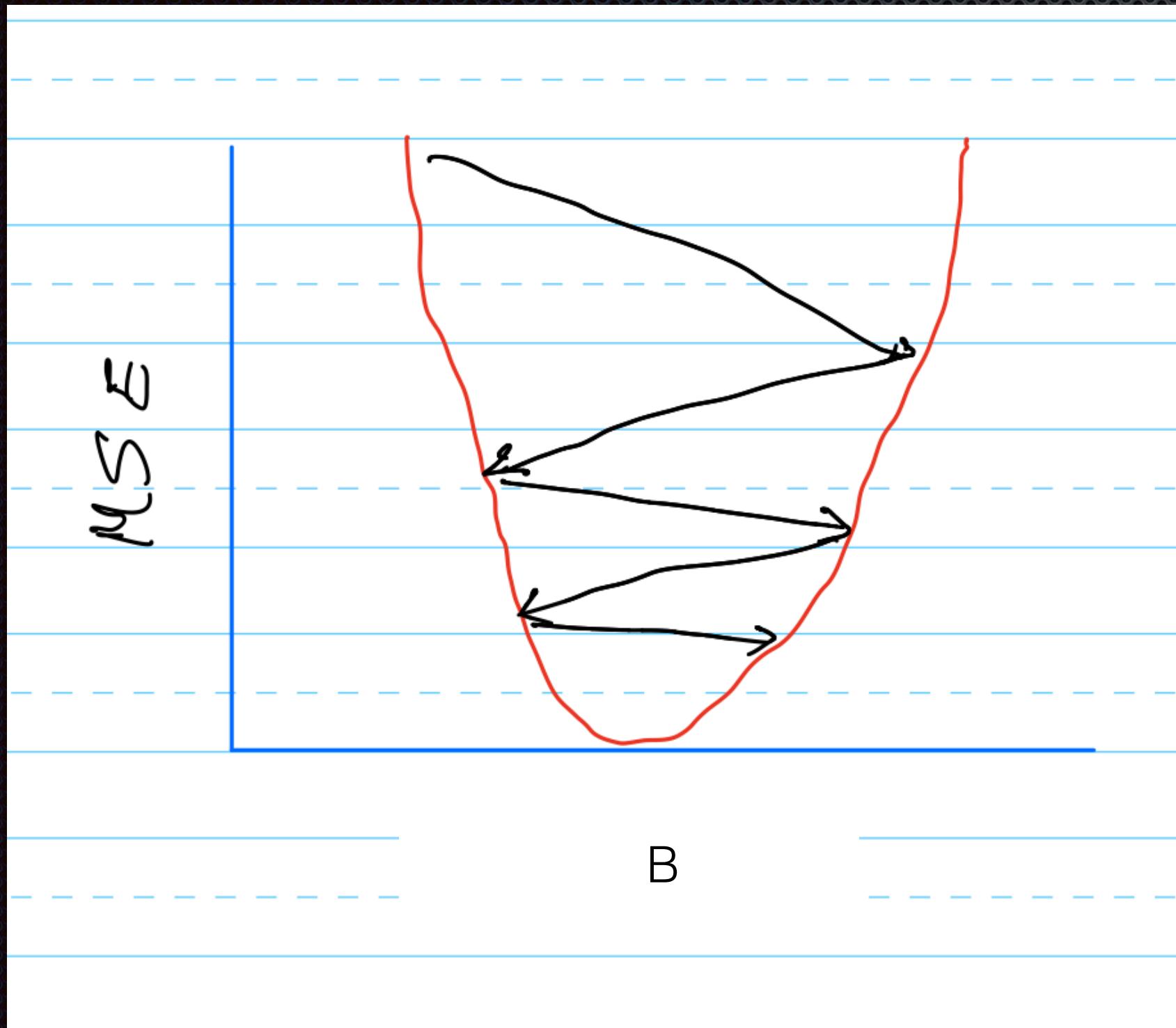
Gradient Descent method



Learning rate

$$B_{\text{new}}^i = B_{\text{old}}^i - \eta \nabla \text{Loss}(B_{\text{old}}^i)$$

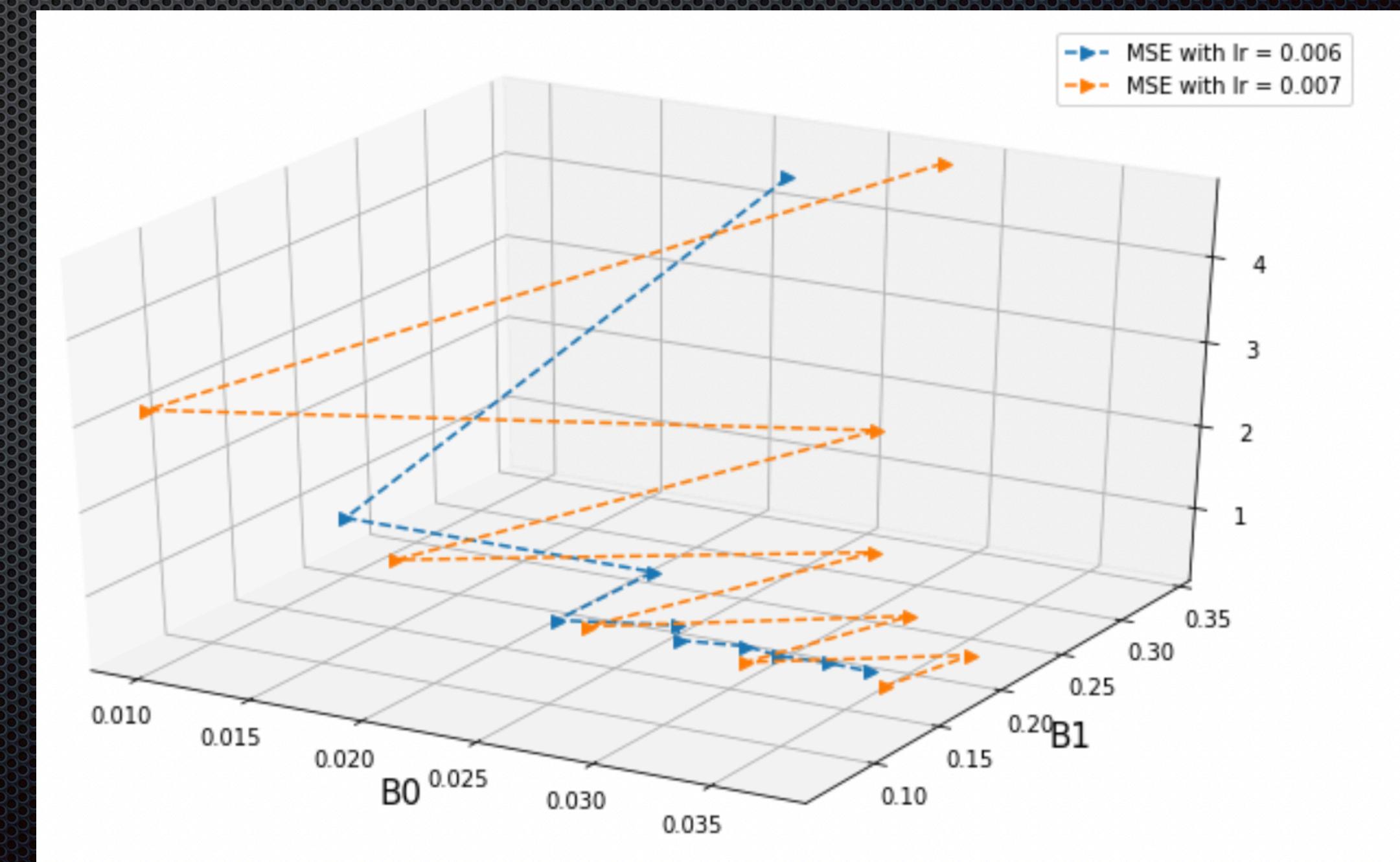
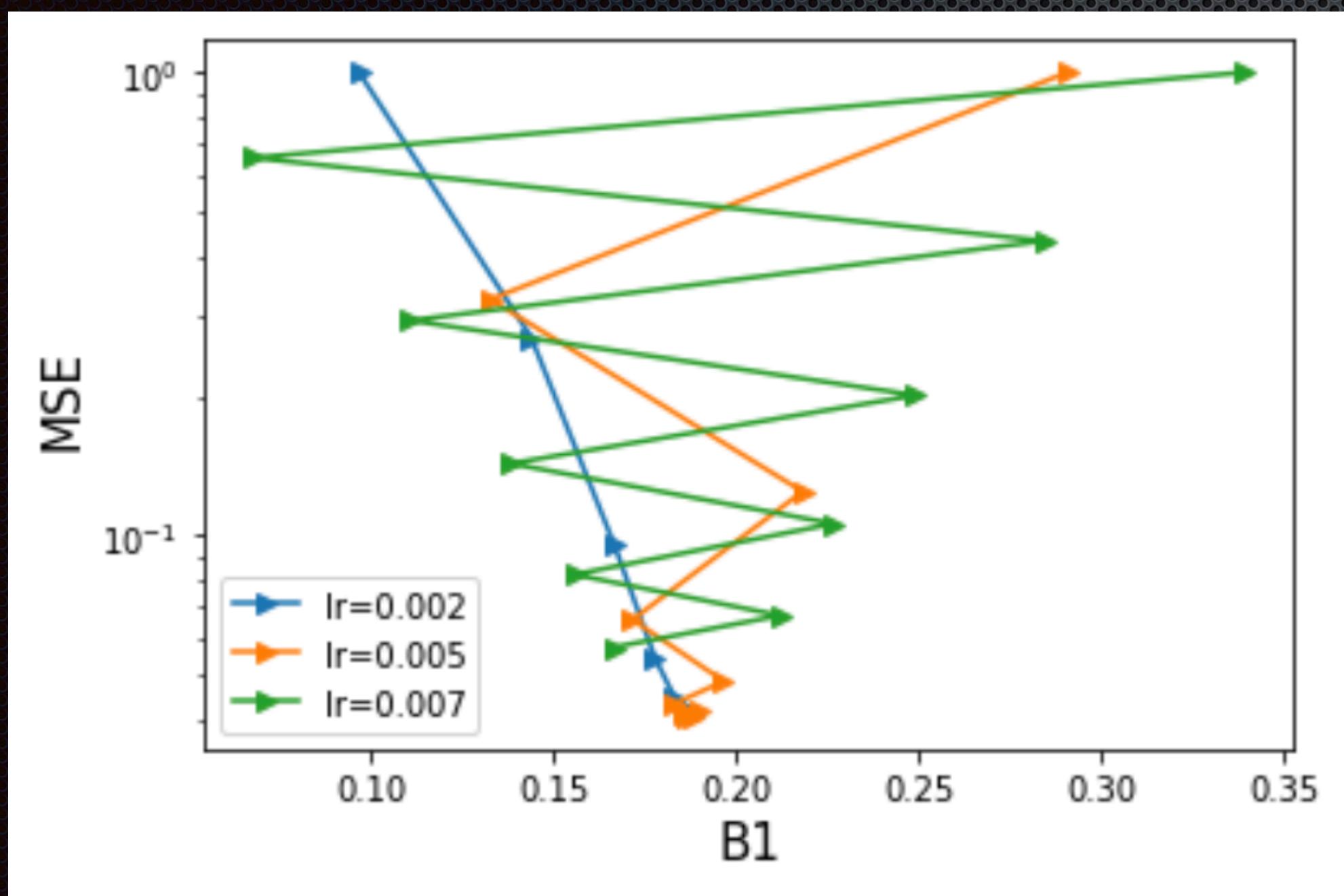
η Is called the learning rate which controls the descent rate of the loss function



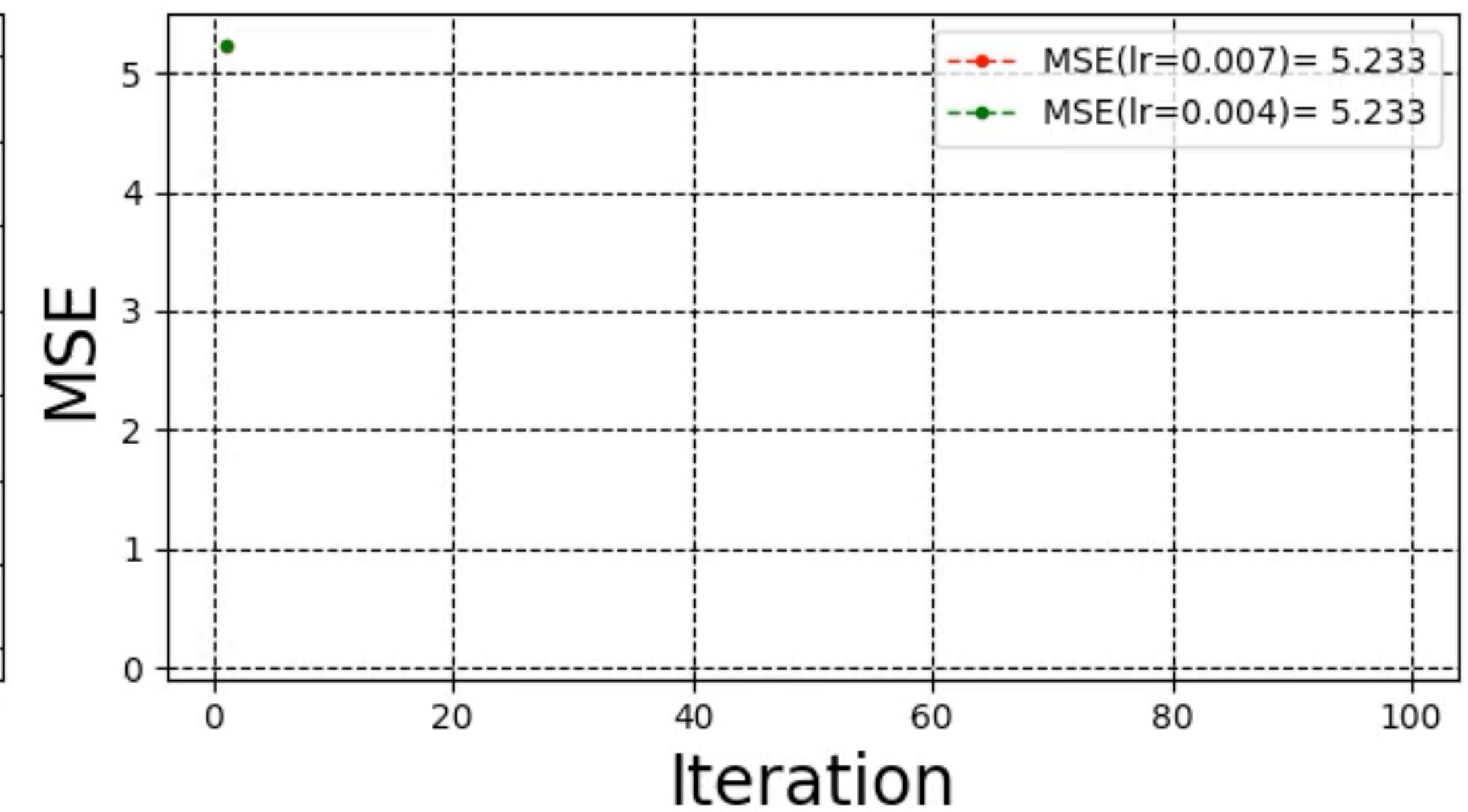
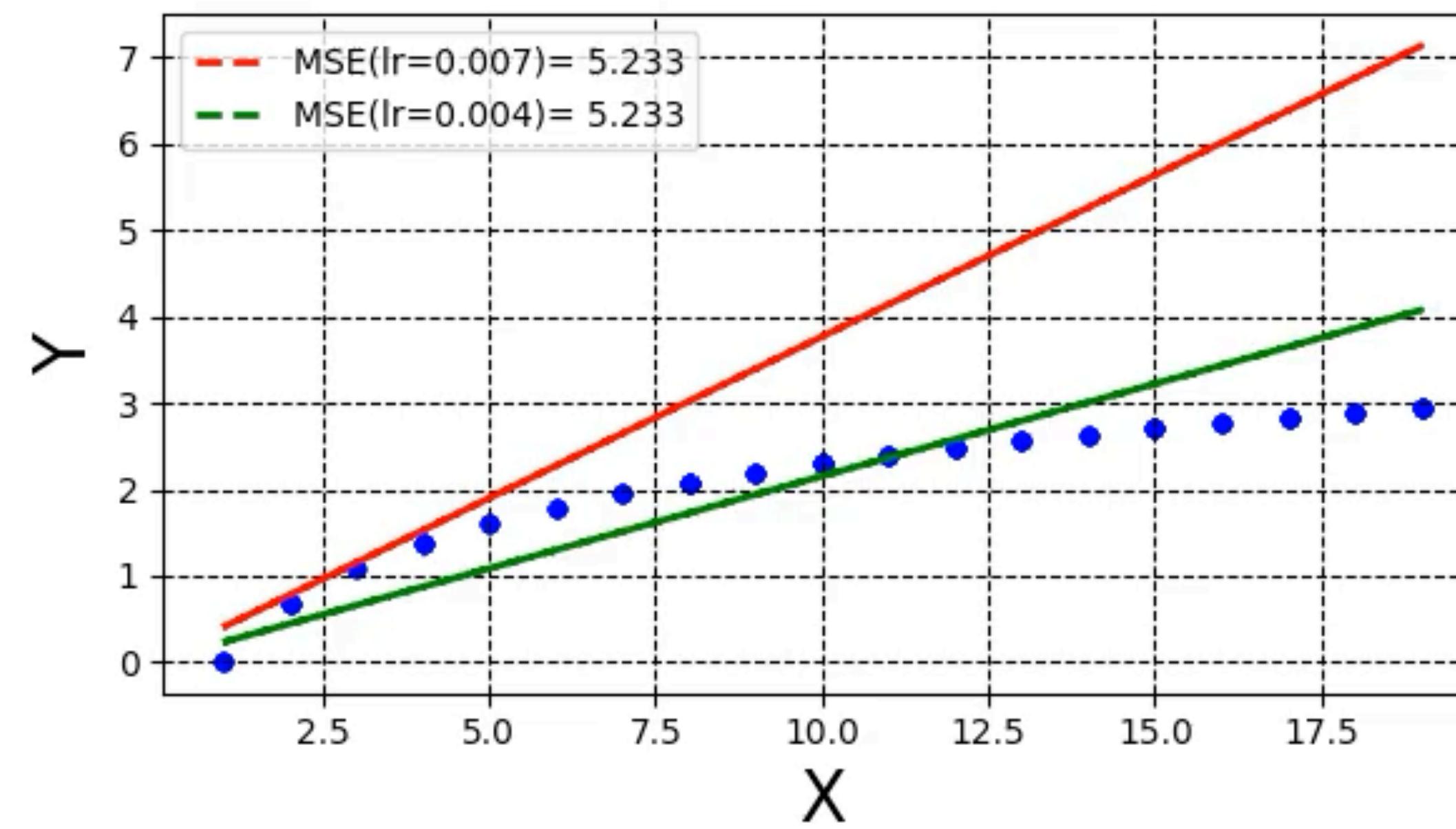
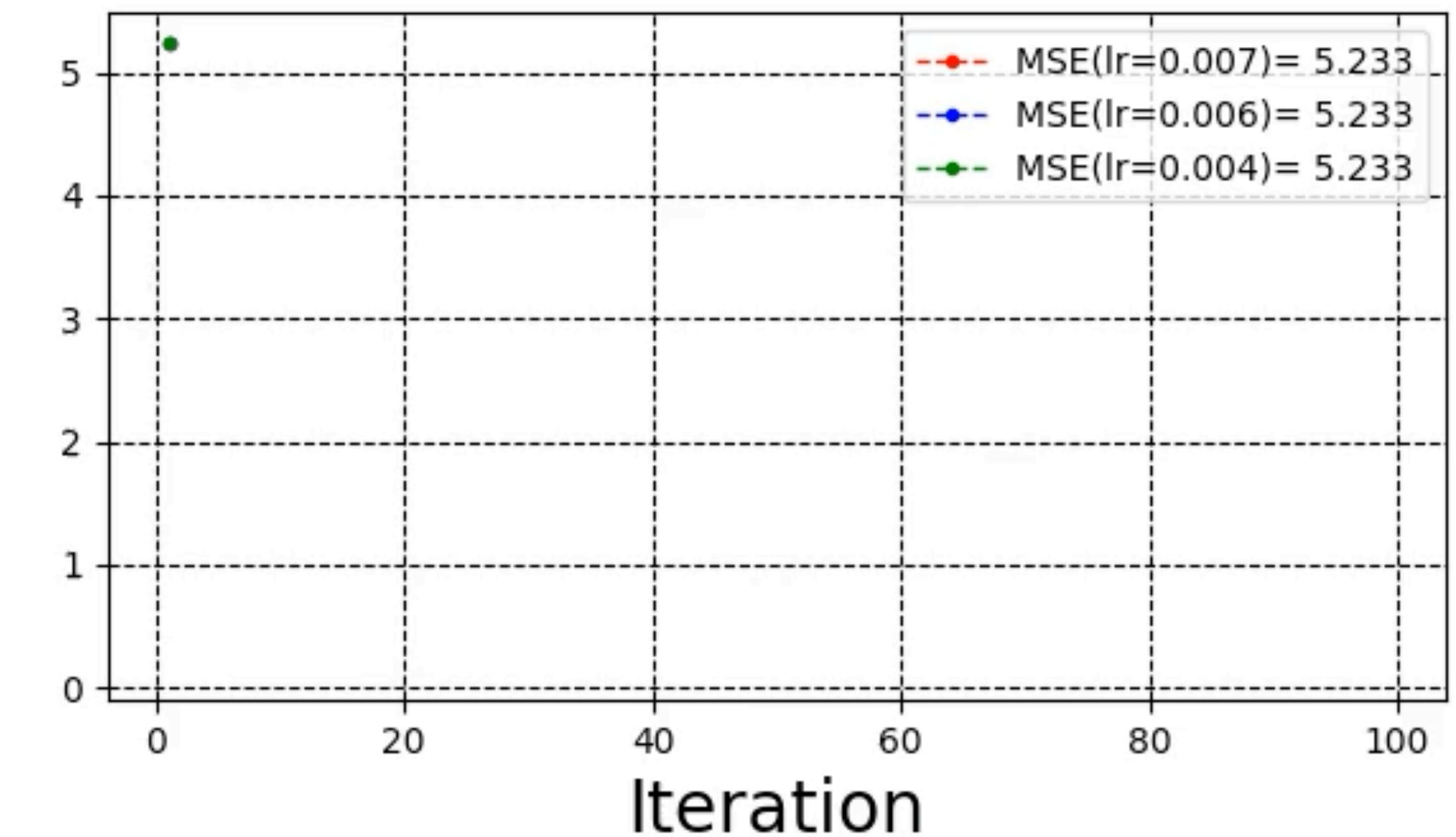
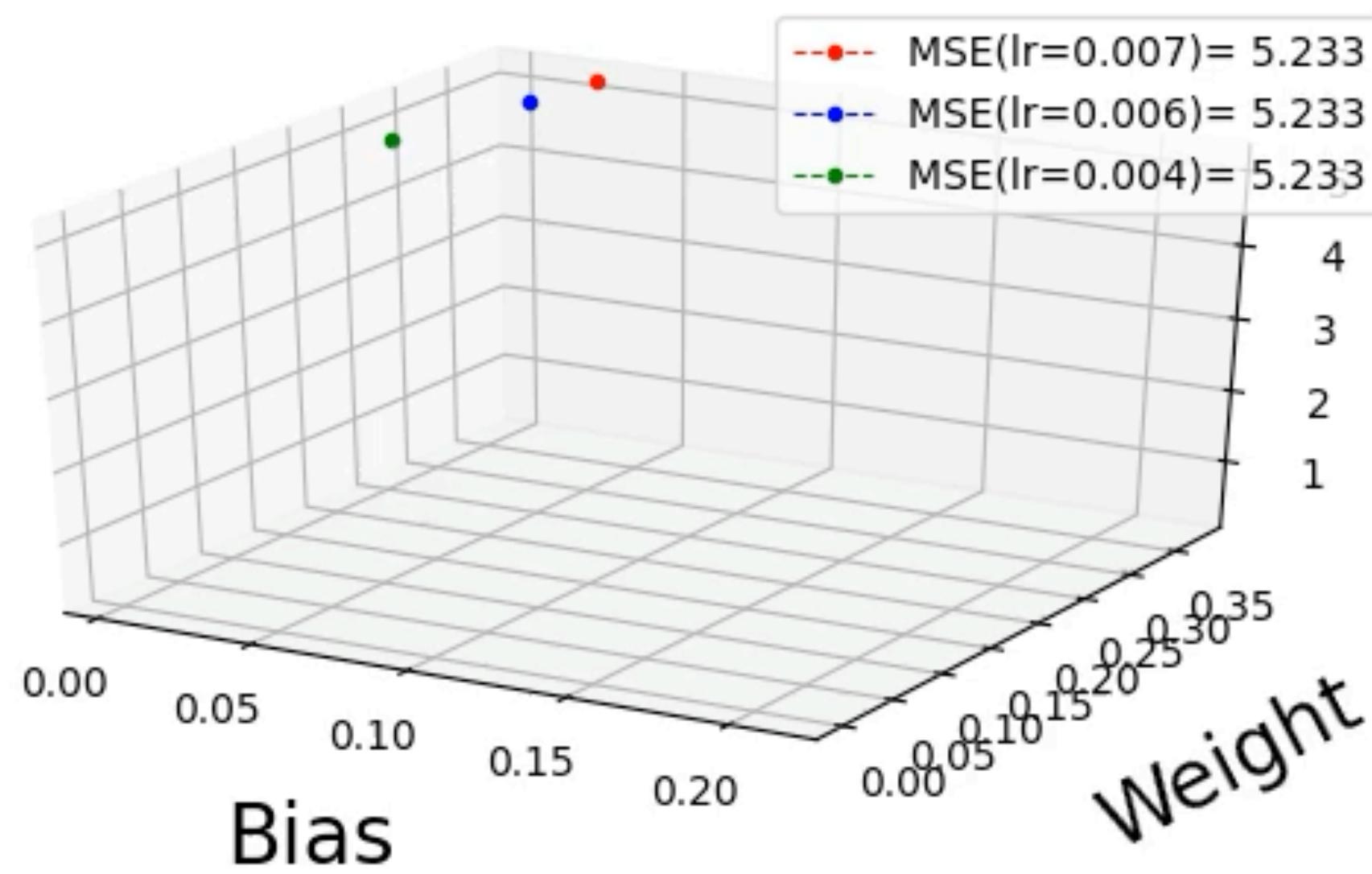
For large learning rate the function oscillates and we will not able to reach the minimum

Learning rate

Effect of different learning rates



Learning rate



To be continued...