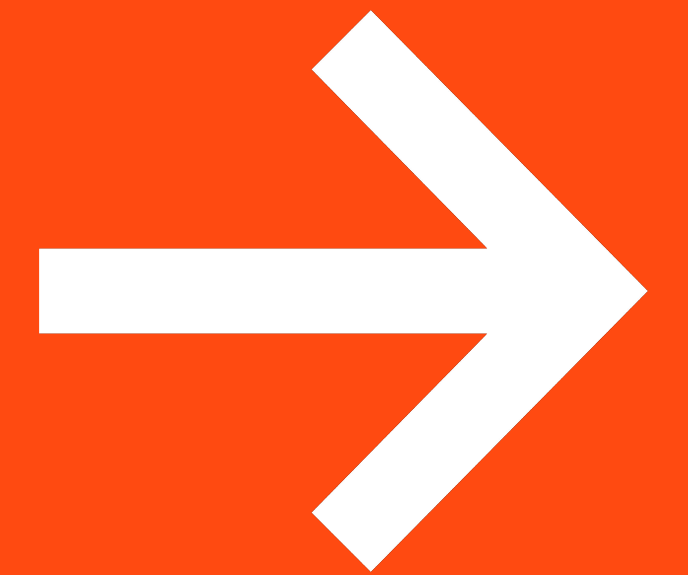


**>>> neue fische**

School and Pool for Digital Talent

# Gradient Descent





**Motivation**

**Intro Gradient descent**

**Learning rate**

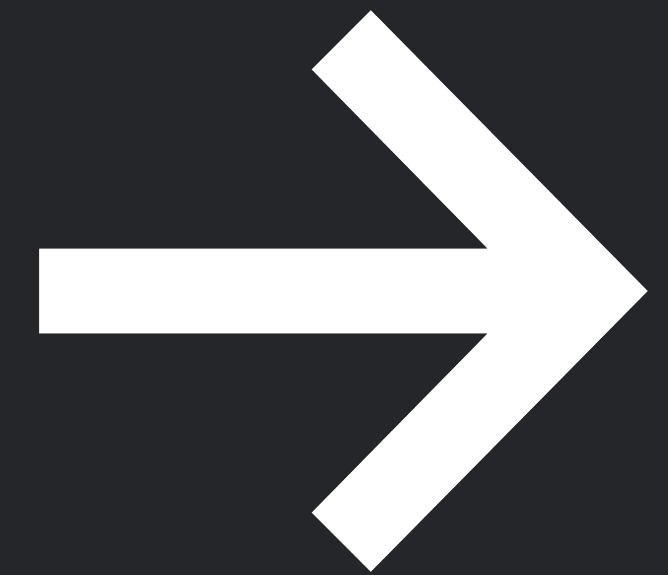
**Summary**



Gradient Descent

# Part 1

# Motivation

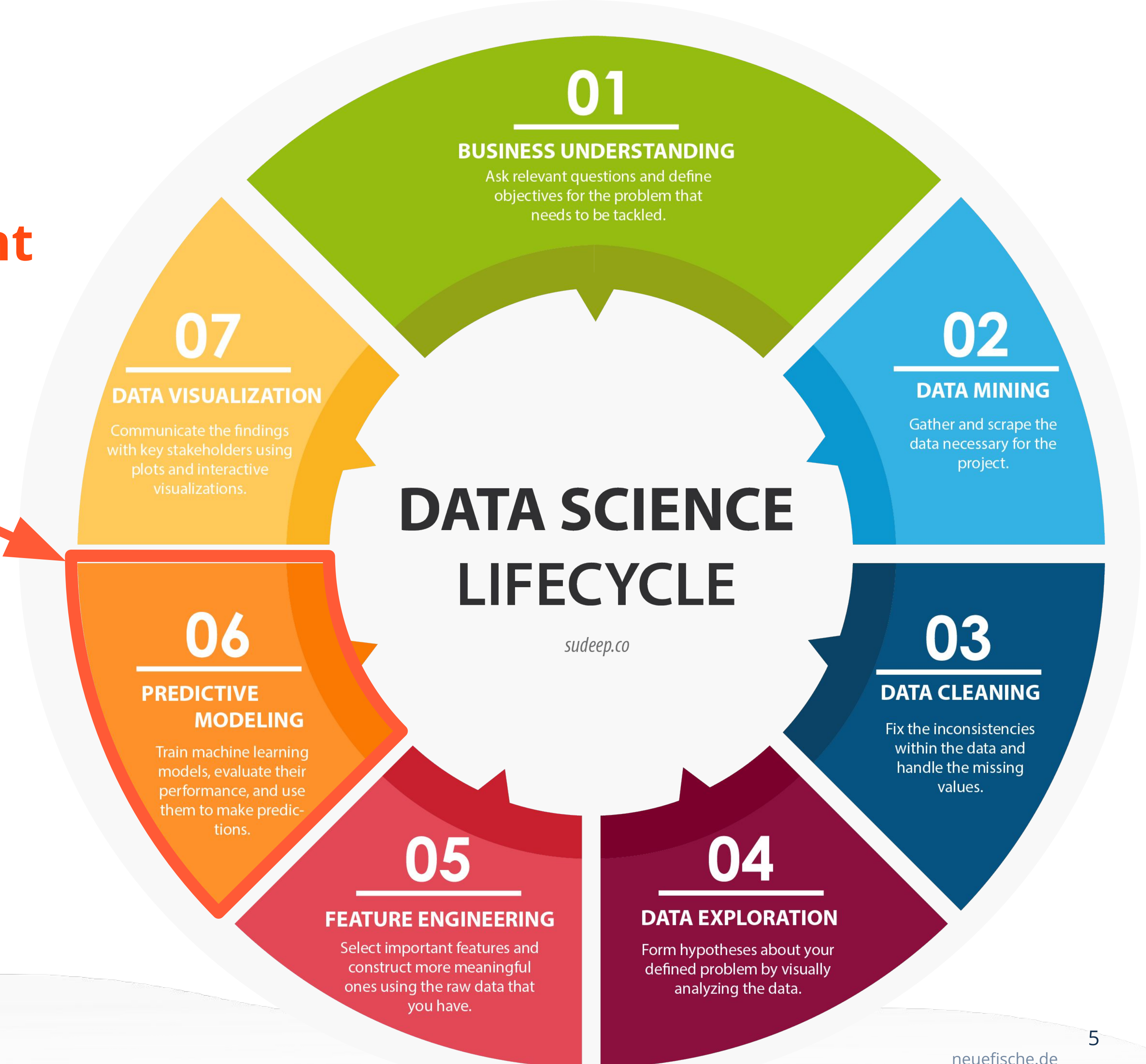


## Orientation

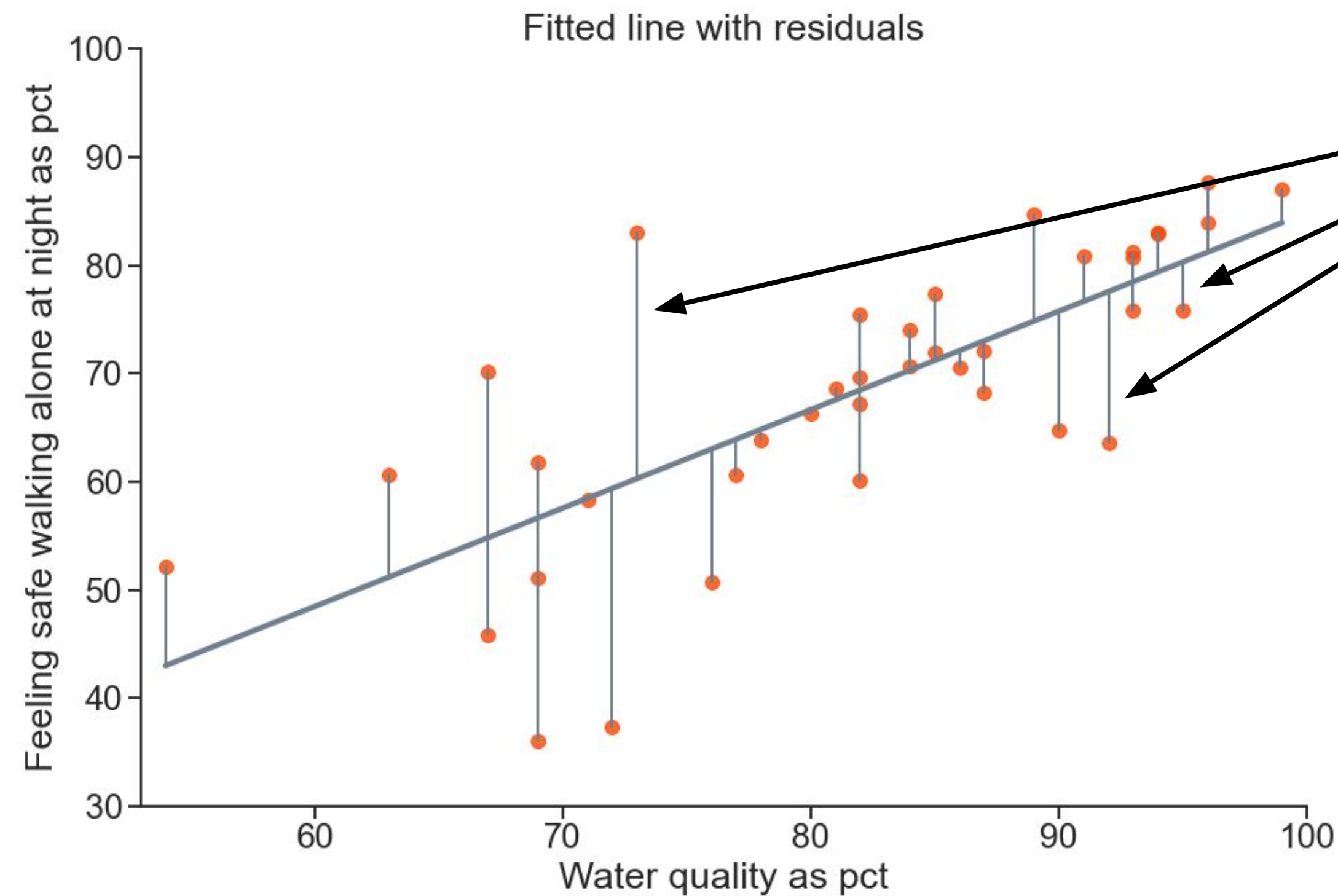
# Where is **Gradient Descent** used in the Data Science Lifecycle?

## 06 PREDICTIVE MODELING:

- select a ML algorithm
- train the ML model
- evaluate the performance
- make predictions



# Recap of Optimization of Linear Regression with OLS



Model

$$y = b_0 + b_1 \cdot x + e$$

Residuals

$$e_i = y_i - \hat{y}_i$$

OLS:

minimize the sum of squared residuals

$$\sum_{i=1}^n e_i^2$$

results in Normal Equation:

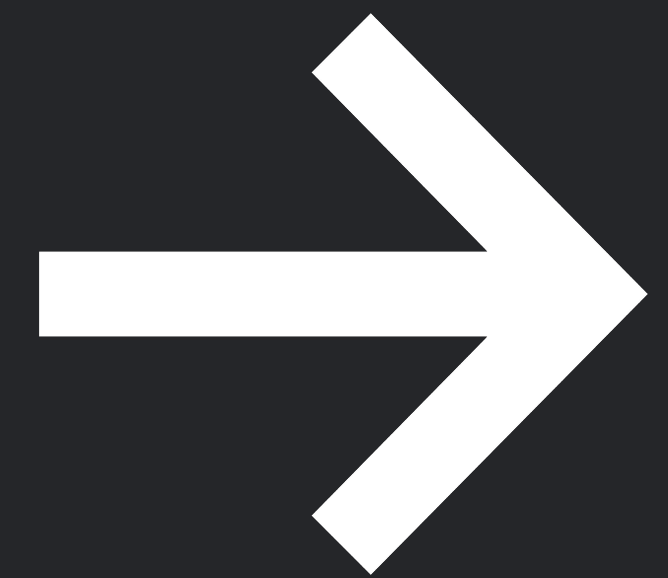
$$b = (X^T X)^{-1} X^T y$$

→ closed form solution

Gradient Descent

Part 2

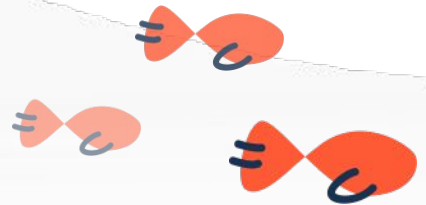
Intro Gradient descent



## Definition

Gradient Descent is an **iterative optimization algorithm** to find the minimum of a function.

→ e.g. of a cost function

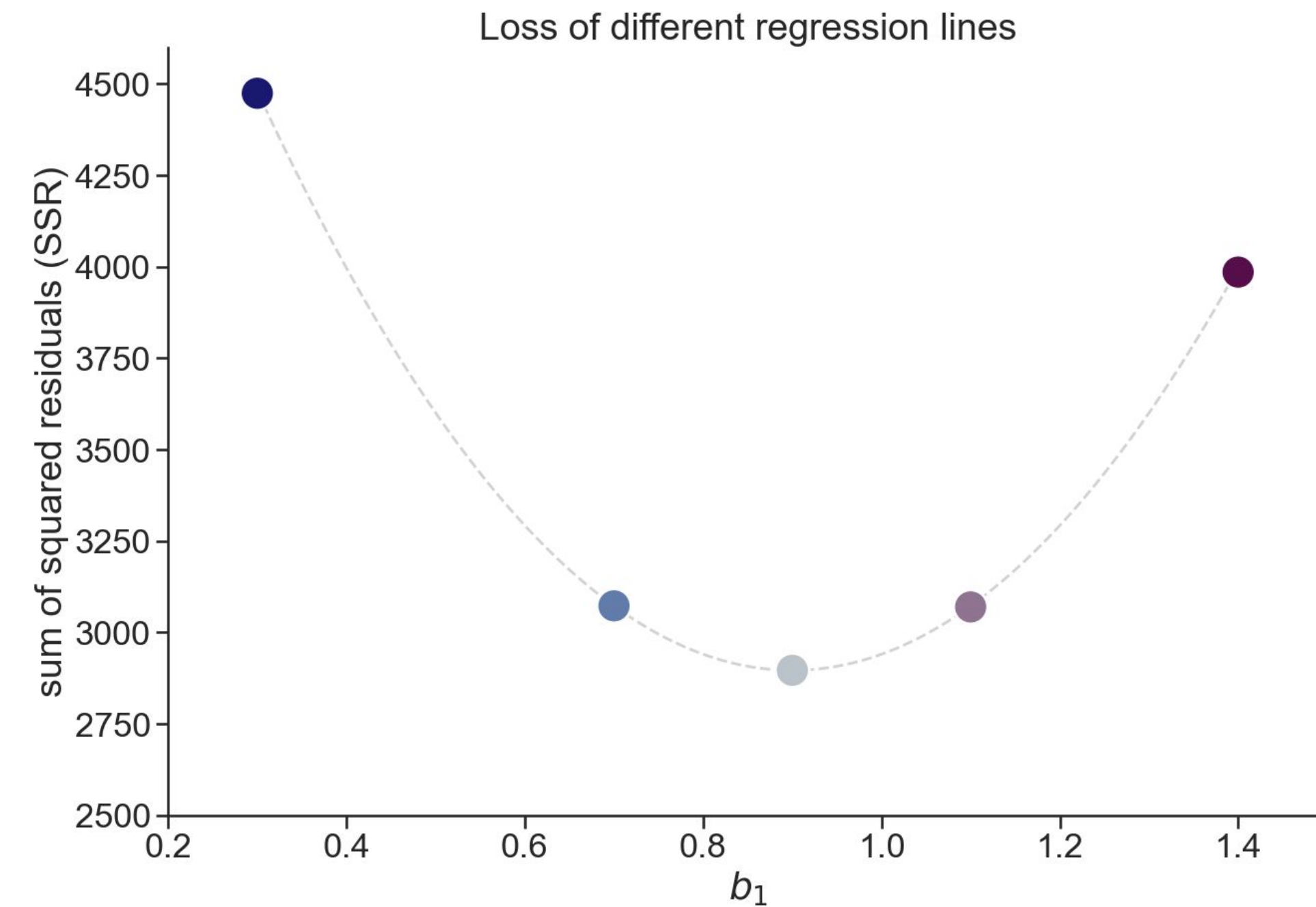
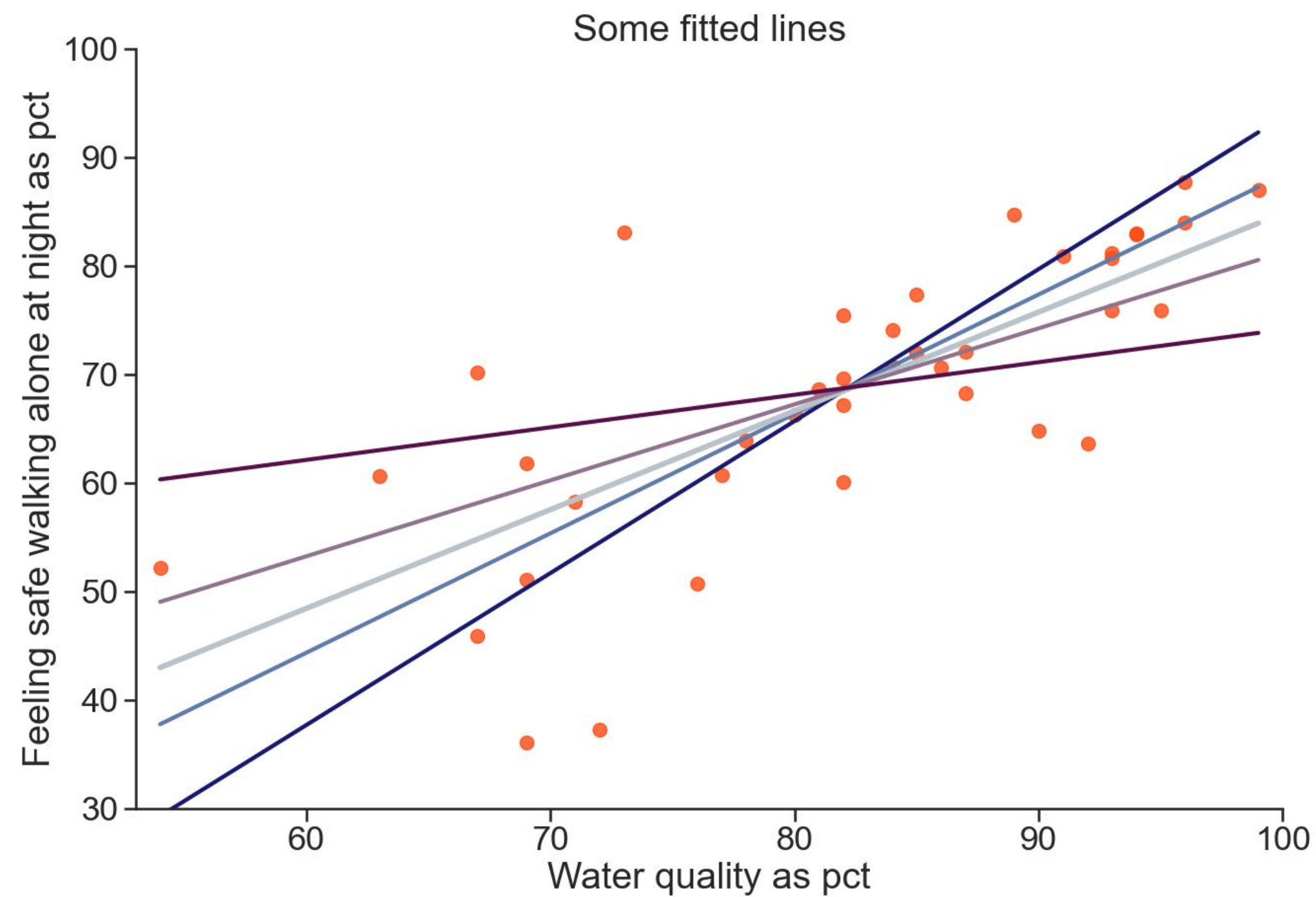




## Training = Finding the best fitting line for our data

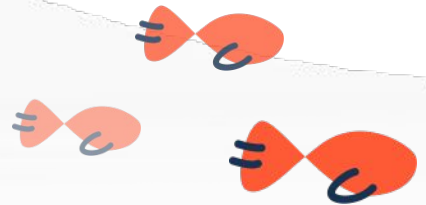
We get to the best fitting line by minimizing the cost function:

$$J(b_0, b_1) = \sum e_i^2 = \sum (y_i - \hat{y}_i)^2 = \sum (y_i - b_0 - b_1 x_i)^2$$



## When to use Gradient Descent?

- GD is a simple **optimization procedure** that can be used for many machine learning algorithms (eg. linear regression, logistic regression, SVM)
- used for “**Backpropagation**” in Neural Nets (variations of GD)
- can be used for **online learning** (needed when Data cannot be stored on only one computer)
- gives **faster results** for problems with many features (computationally simpler)  
(double nr. of features: Normal Equation will take 5-8 times longer)



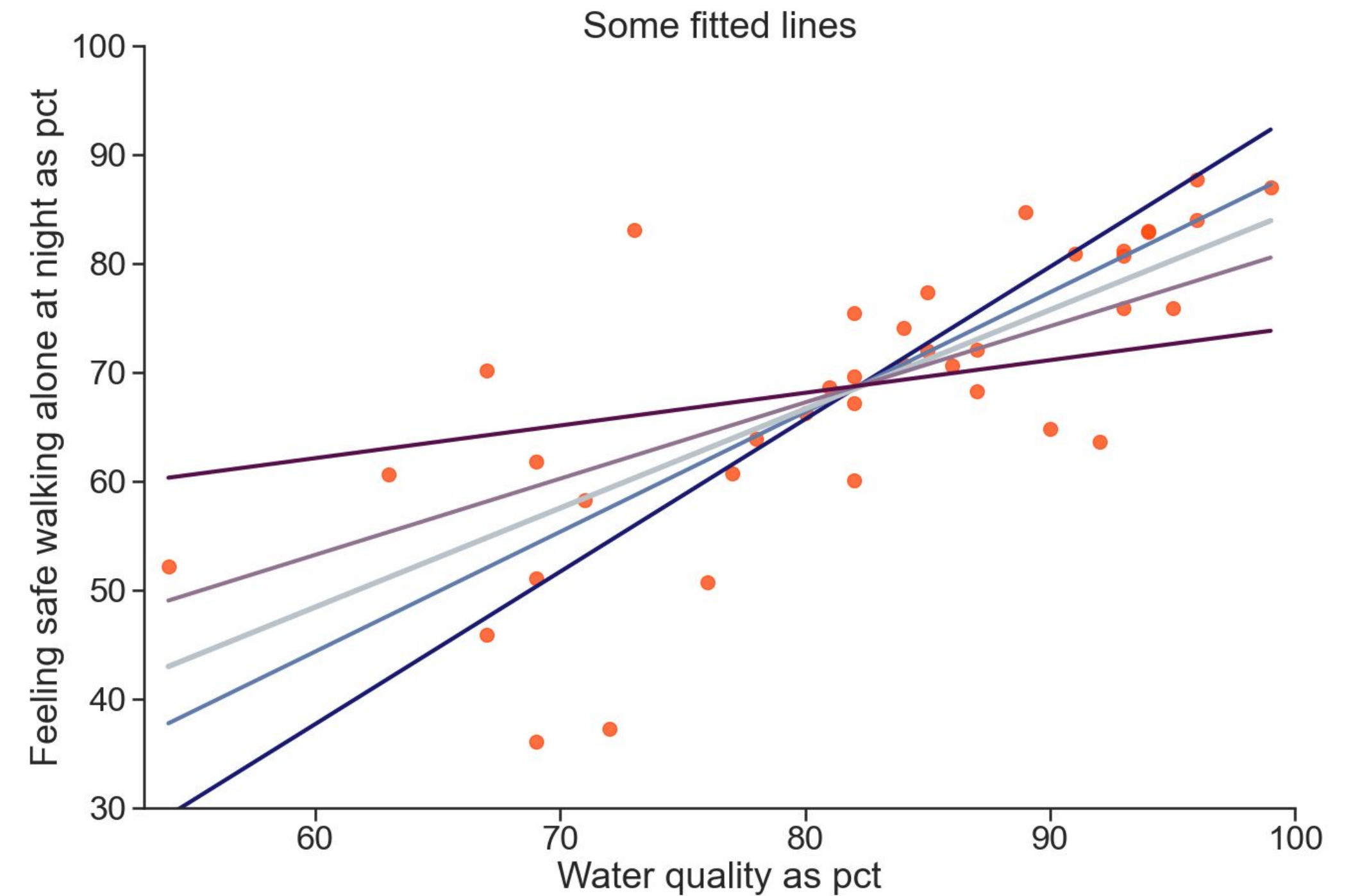
# How to use Gradient Descent?

## First Step

## Define Model

Example:  
least squares with one feature

$$\hat{y} = b_0 + b_1 \cdot x$$



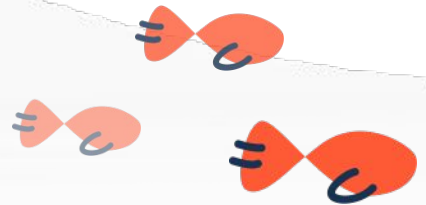
## Second Step

# Define Cost Function

We get to the best fitting line by minimizing the cost function:

$$J(b_0, b_1) = \sum e_i^2 = \sum (y_i - \hat{y}_i)^2 = \sum (y_i - b_0 - b_1 x_i)^2$$

$$\min(J(b_0, b_1))$$

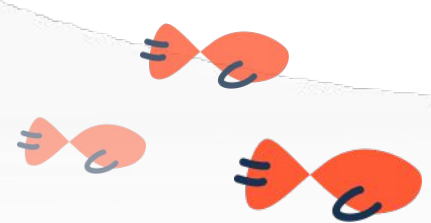
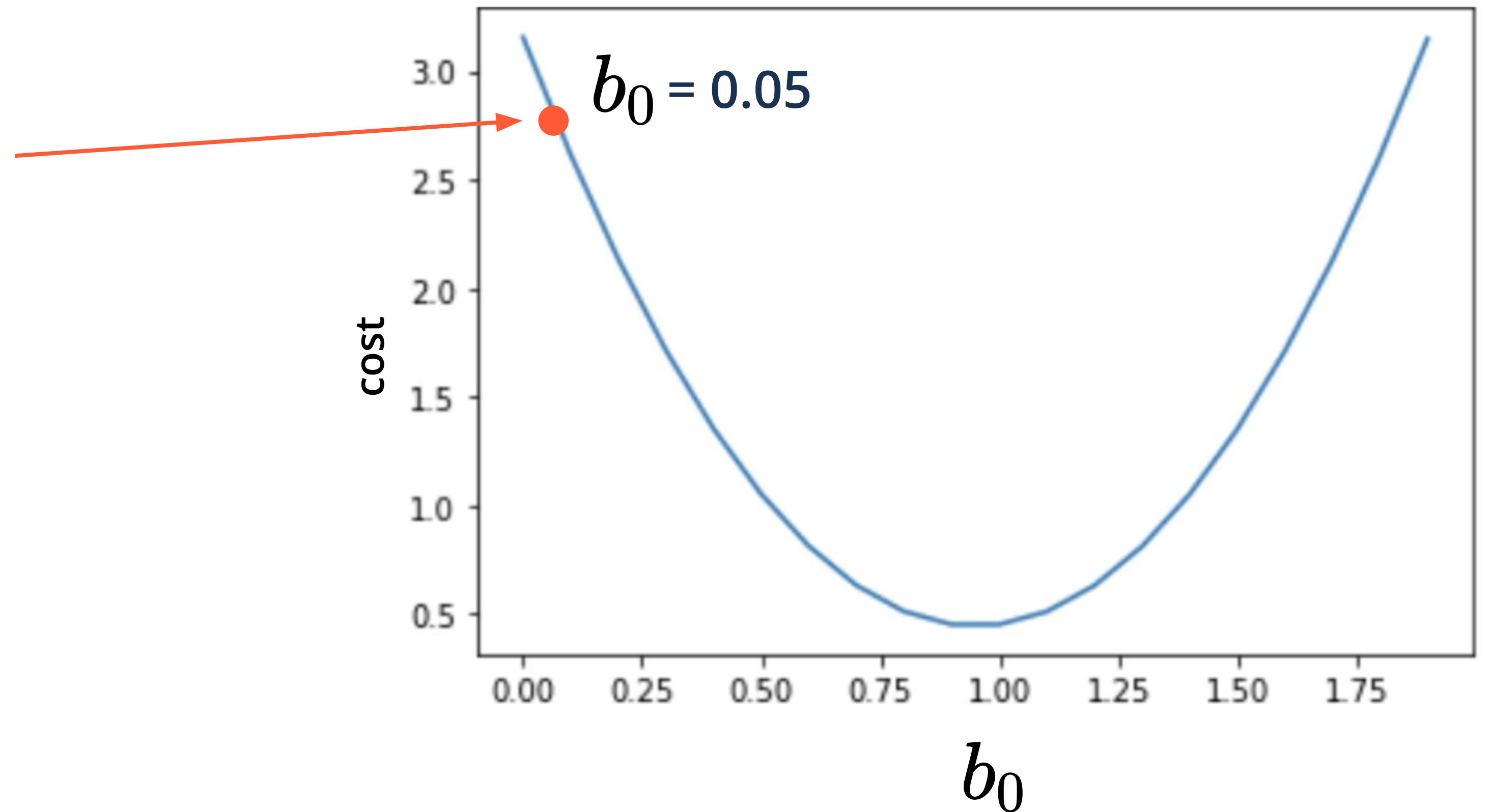




### Third Step

## Initialize Gradient Descent

Deliberately set some random starting values  $b$



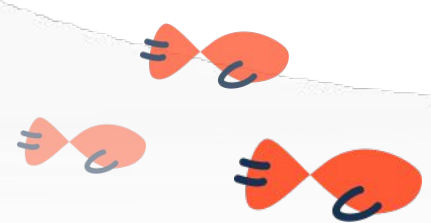
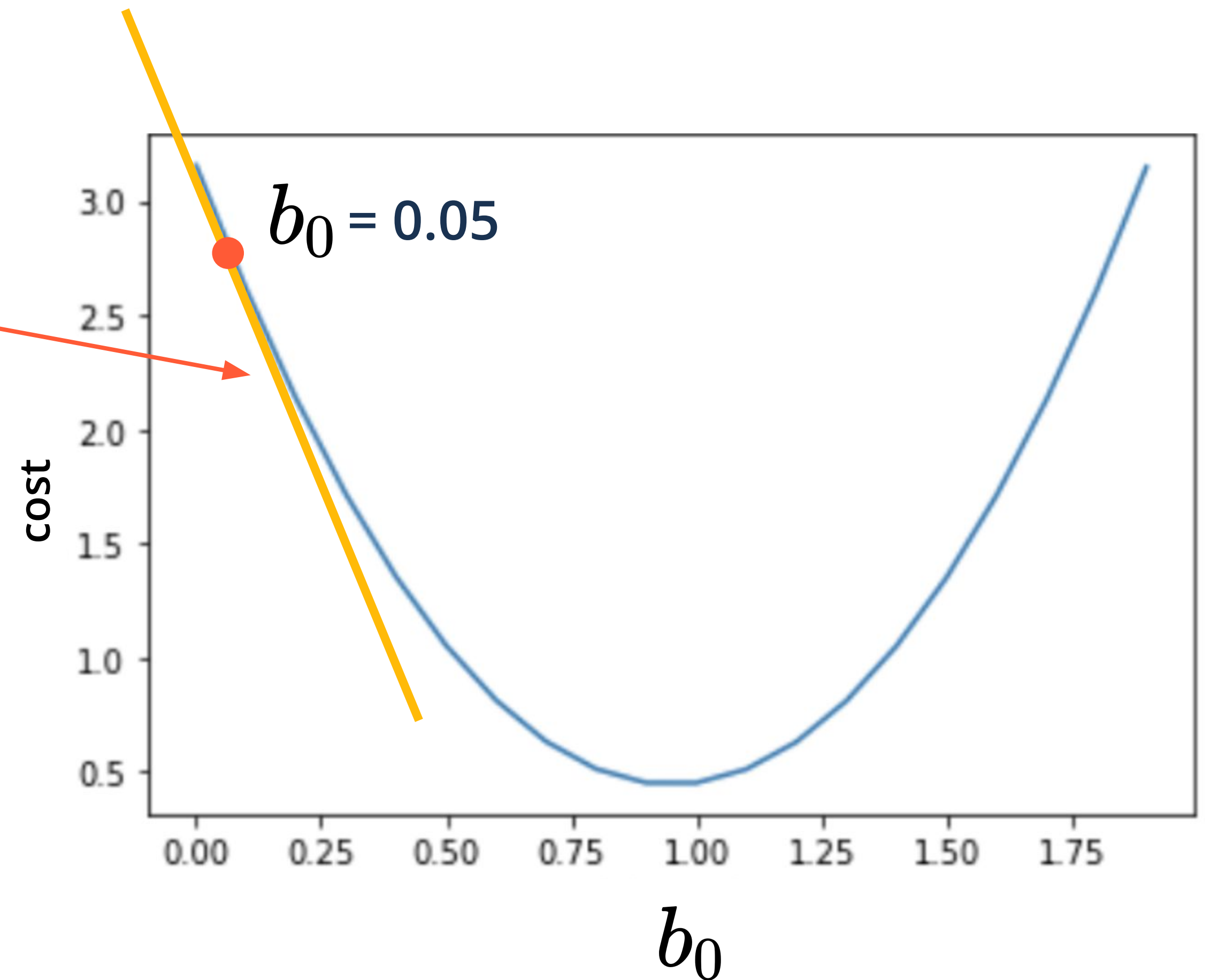
## Fourth Step

# Gradient Descent

Start descent:

- Take **derivatives** of the cost function with respect to your parameters  $b$

$$\text{slope} = -5.7$$



## Fourth Step

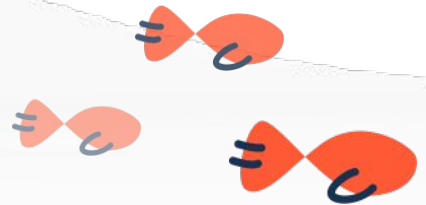
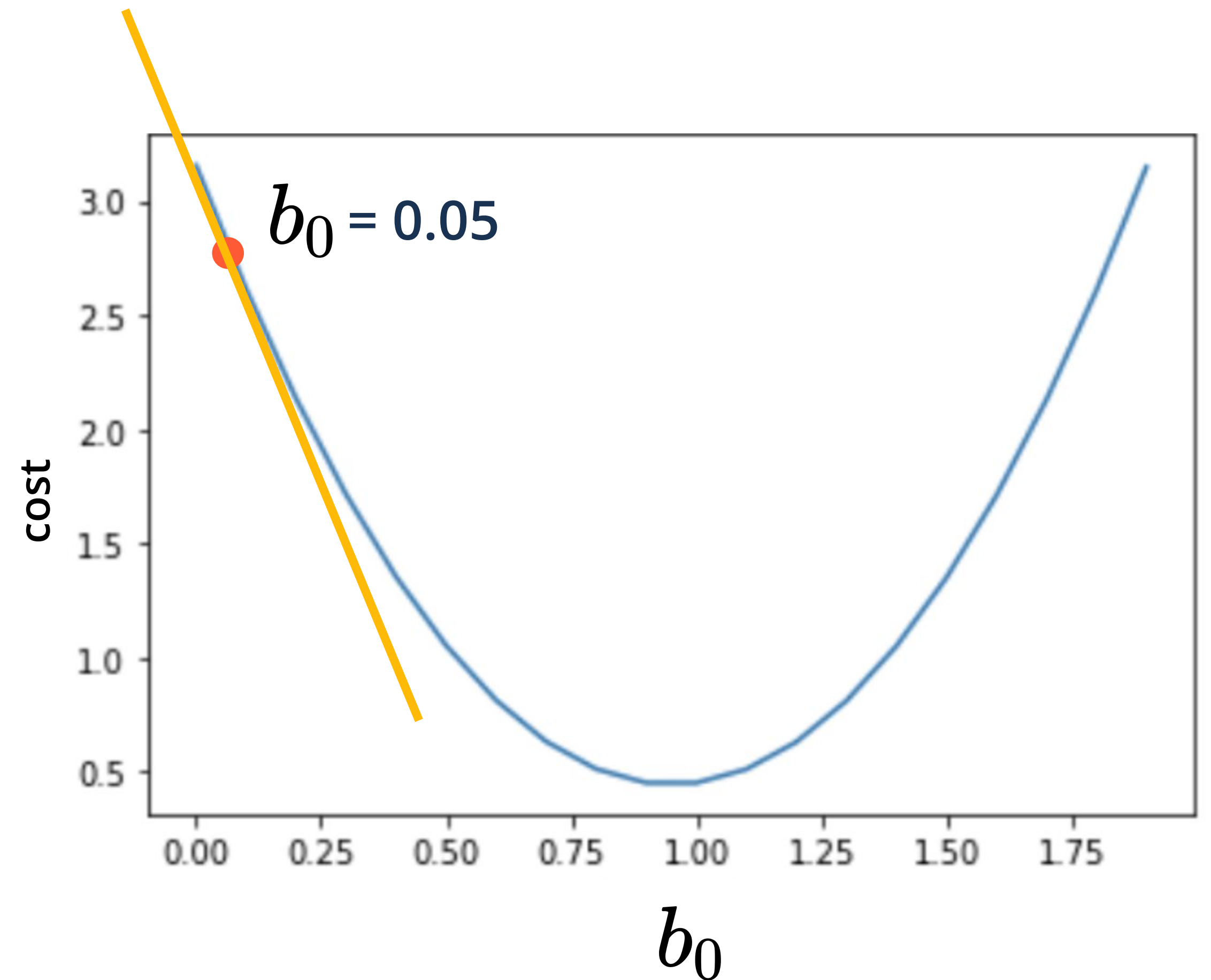
# Gradient Descent

Start descent:

- Take derivatives with respect to your parameters  $b$
- Set your **learning rate** (step-size)

$$\alpha = 0.1$$

$$\text{slope} = -5.7$$



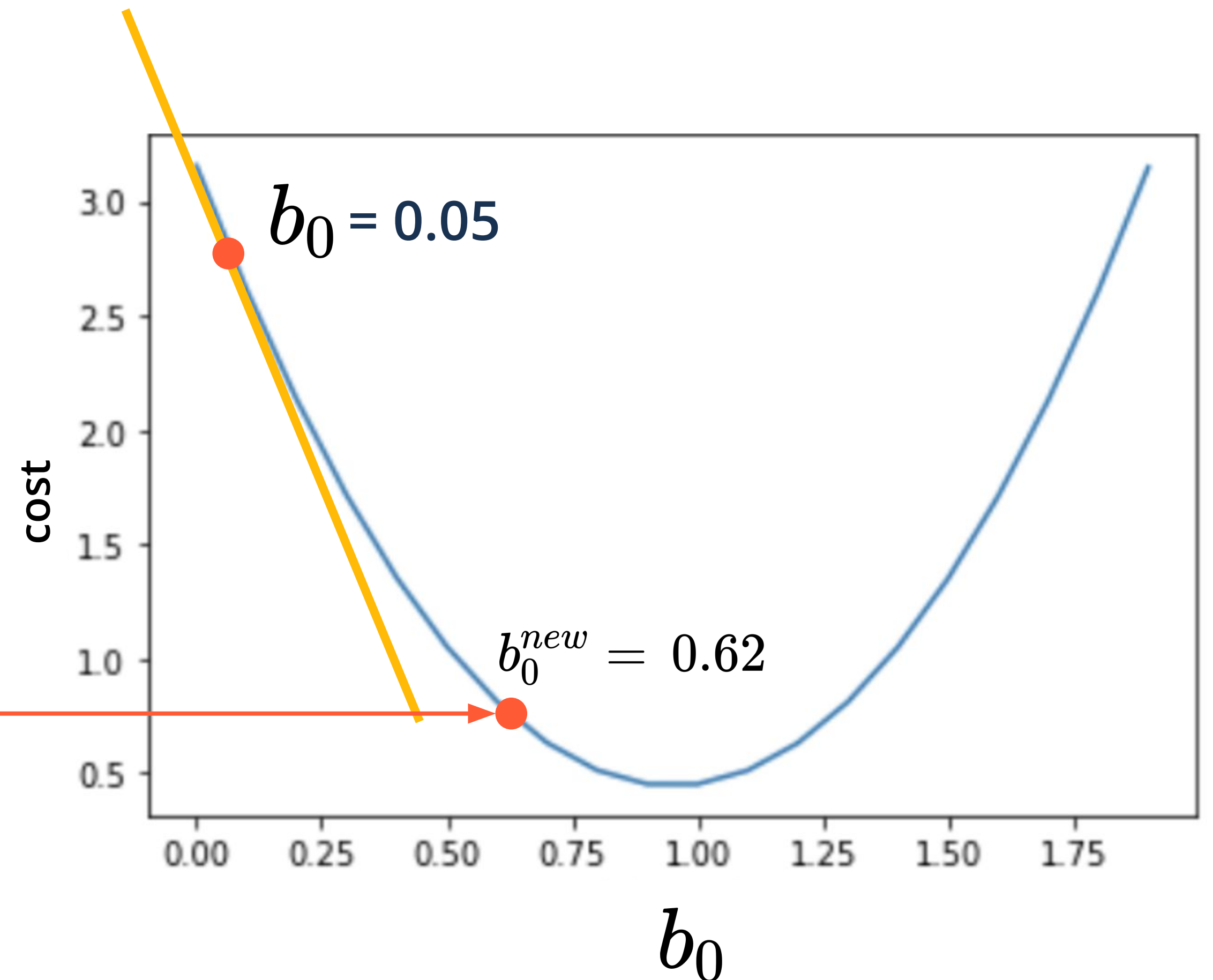
## Fourth Step

# Gradient Descent

Start descent:

- Take derivatives with respect to your parameters  $b$
- Set your learning rate (step-size)
- Adjust your parameters (step)

$$\begin{aligned} \alpha &= 0.1 & \text{slope} &= -5.7 \\ \text{step} &= \alpha \cdot \text{derivative of cost function} \\ b_0^{\text{new}} &= b_0^{\text{old}} - \text{step} \\ b_0^{\text{new}} &= 0.05 - (0.1 \cdot (-5.7)) \\ b_0^{\text{new}} &= 0.62 \end{aligned}$$





## Fifth Step

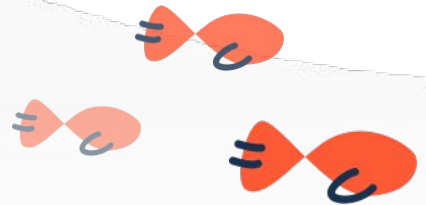
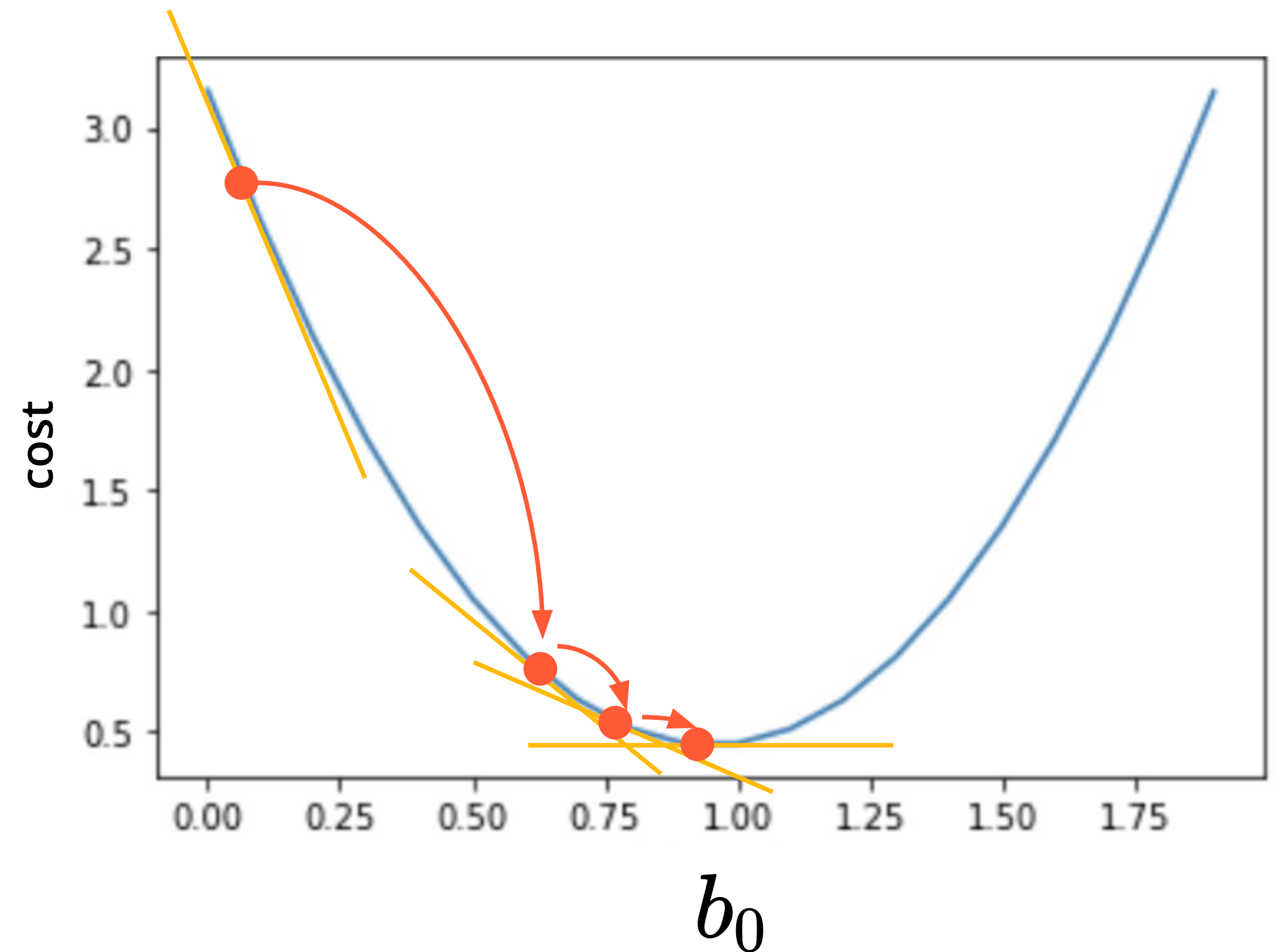
# Gradient Descent

Start descent:

- Take derivatives with respect to your parameters  $b$
- Set your learning rate (step-size)
- Adjust your parameters (step)

$$\text{step} = \alpha \cdot \text{derivative of cost function}$$
$$b_0^{\text{new}} = b_0^{\text{old}} - \text{step}$$

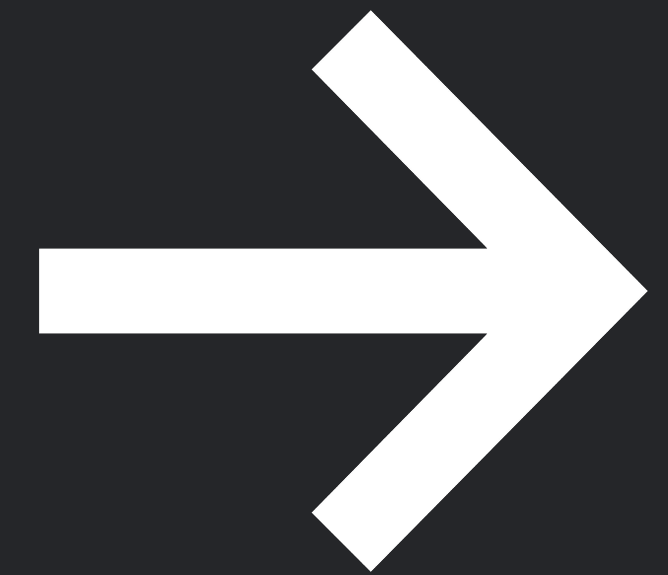
Repeat till there is no further improvement



Gradient Descent

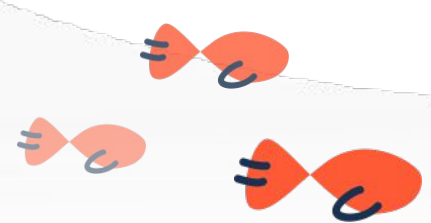
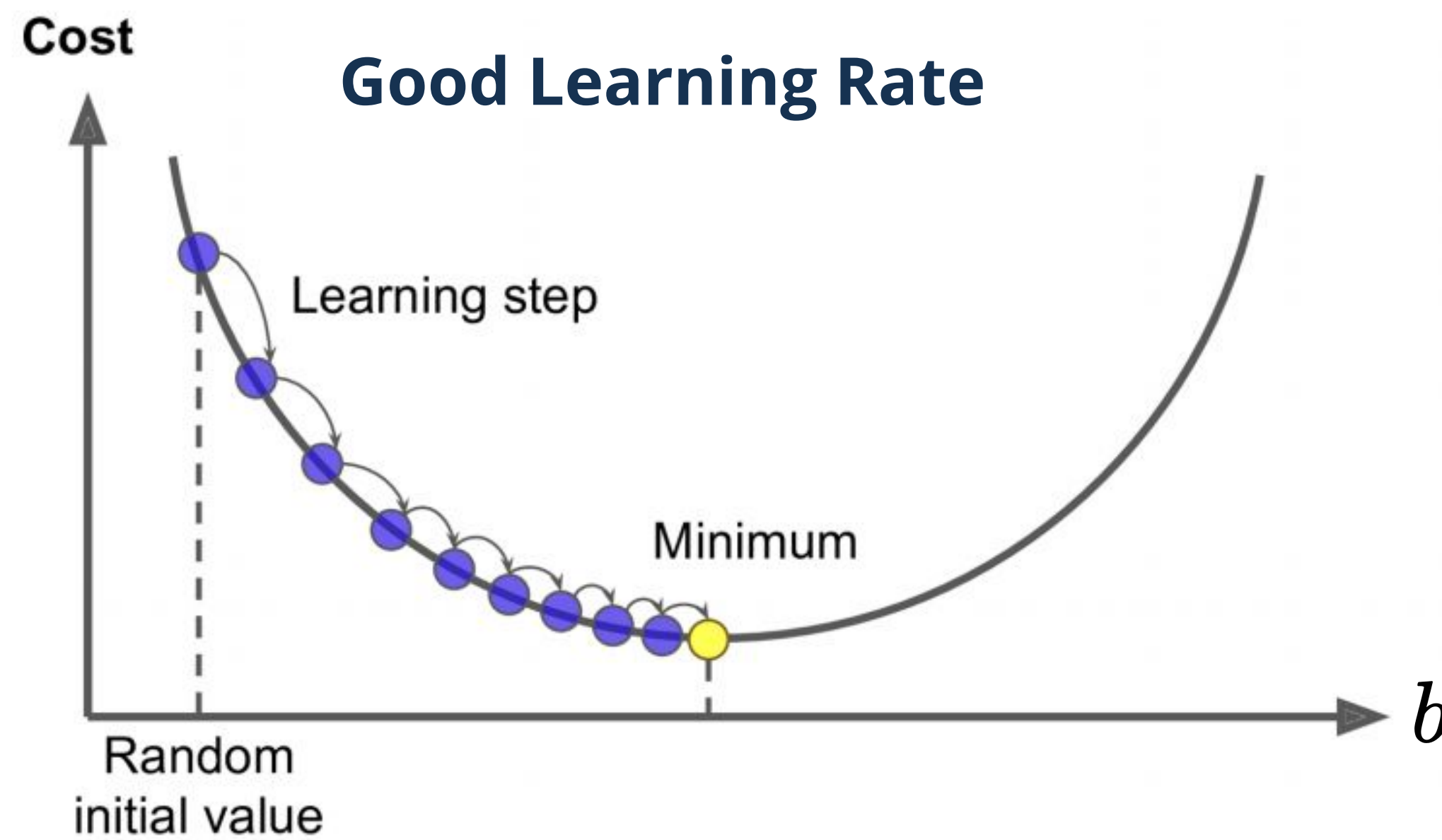
# Part 3

# Learning Rate



# Learning Rate

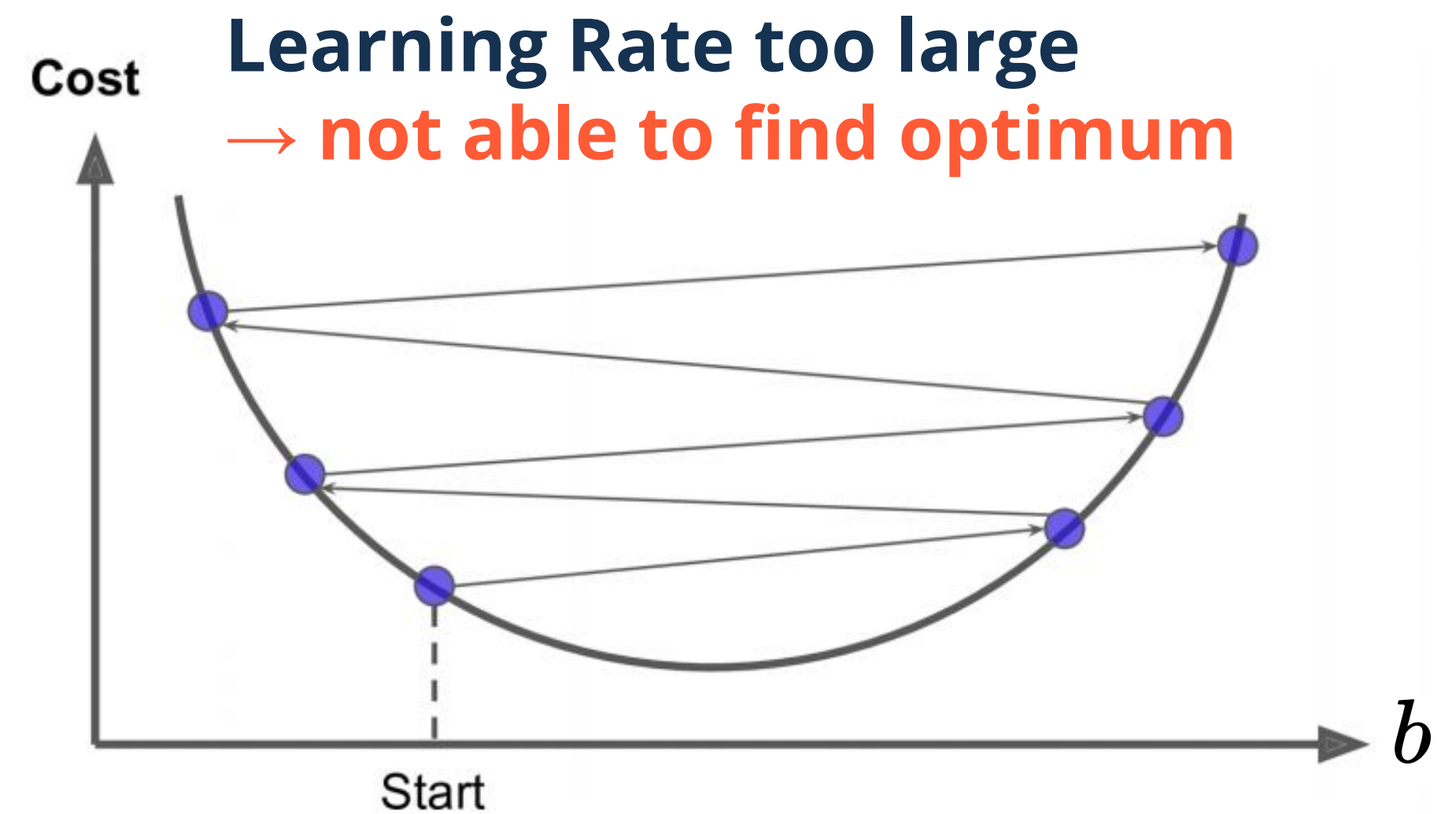
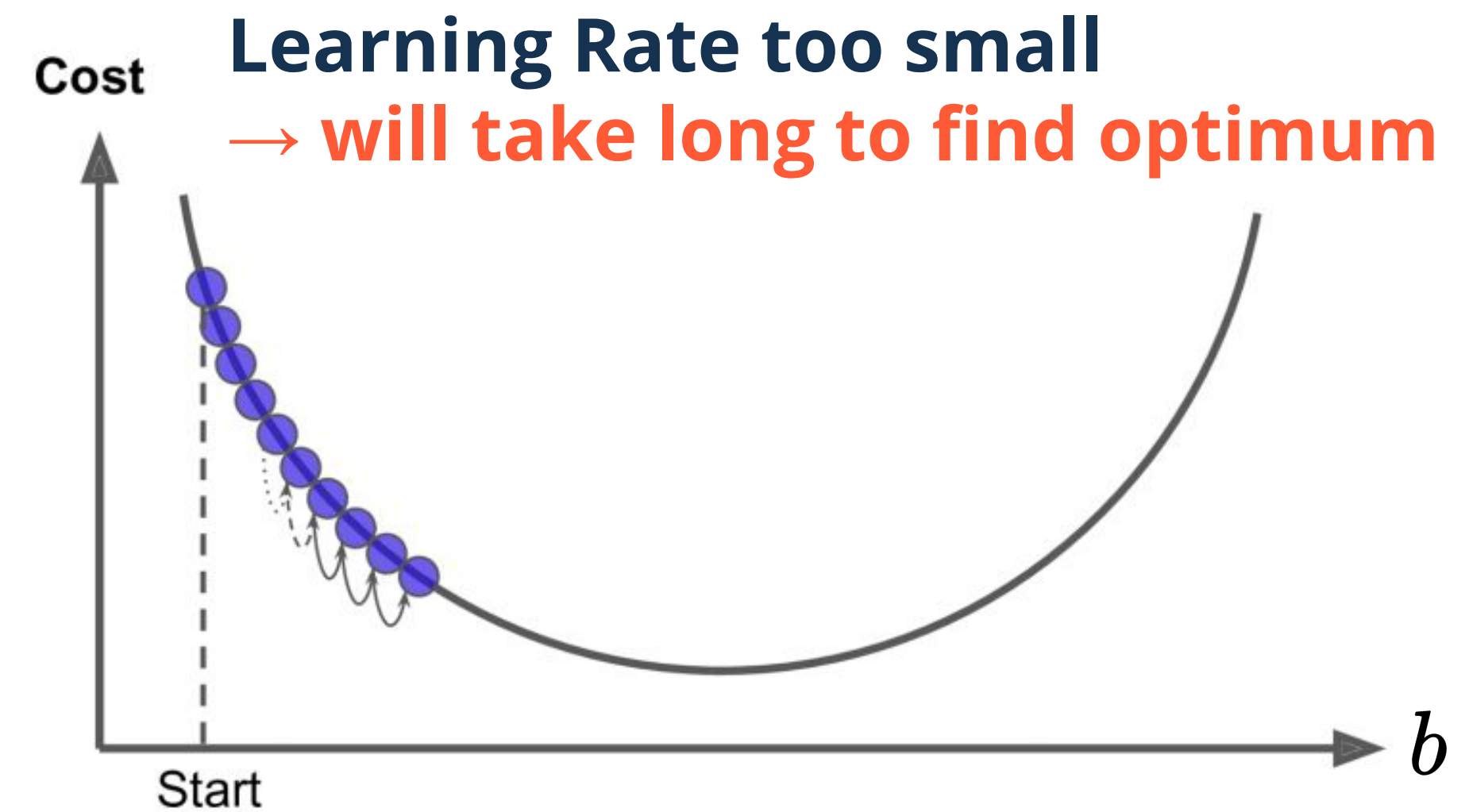
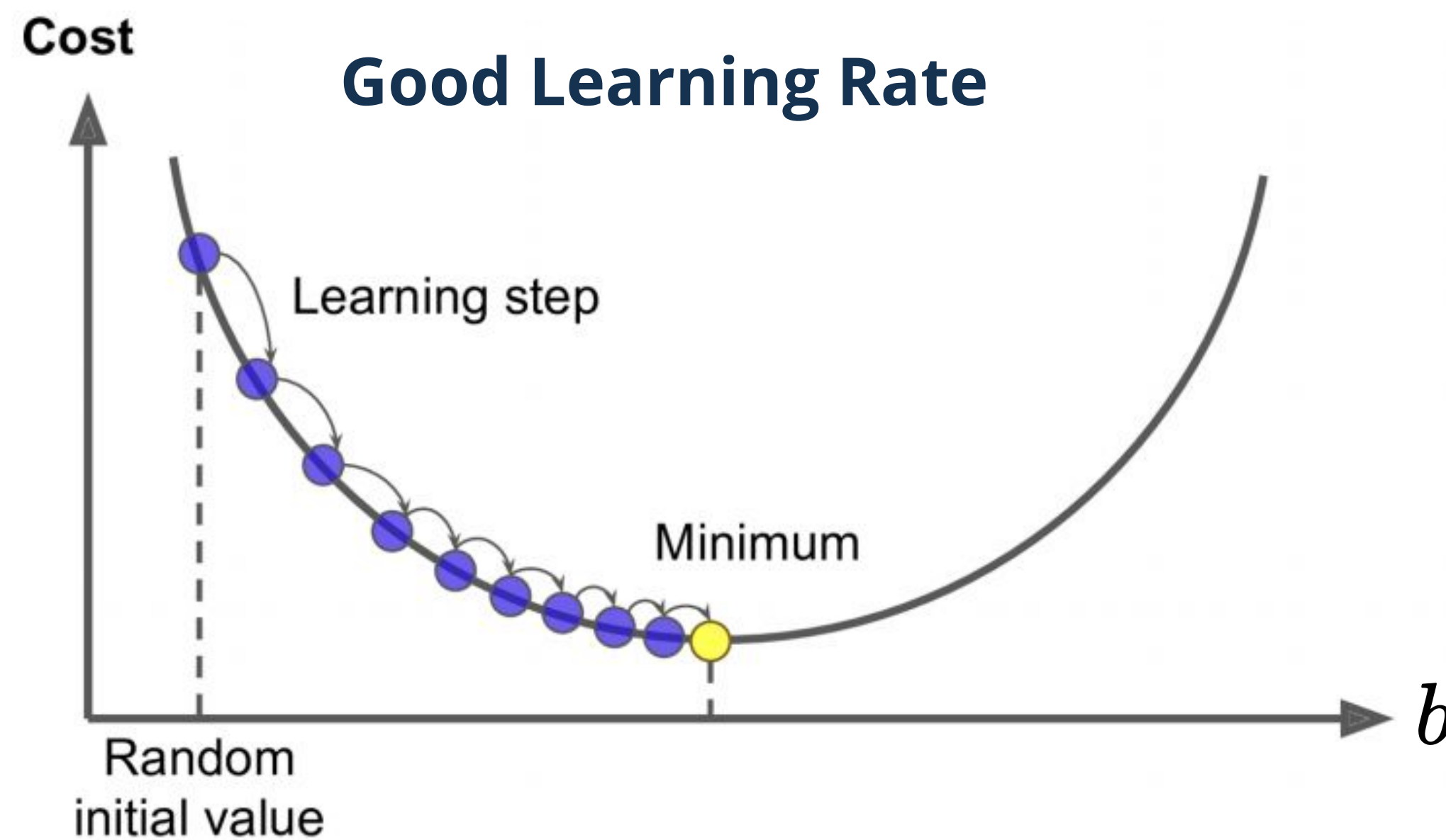
It is important to find a good value for the learning rate!



## Gradient Descent - Hyperparameter

# Learning Rate

It is important to find a good value for the learning rate!





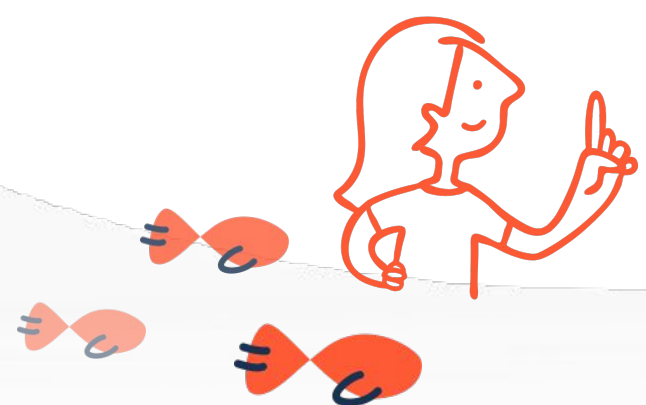
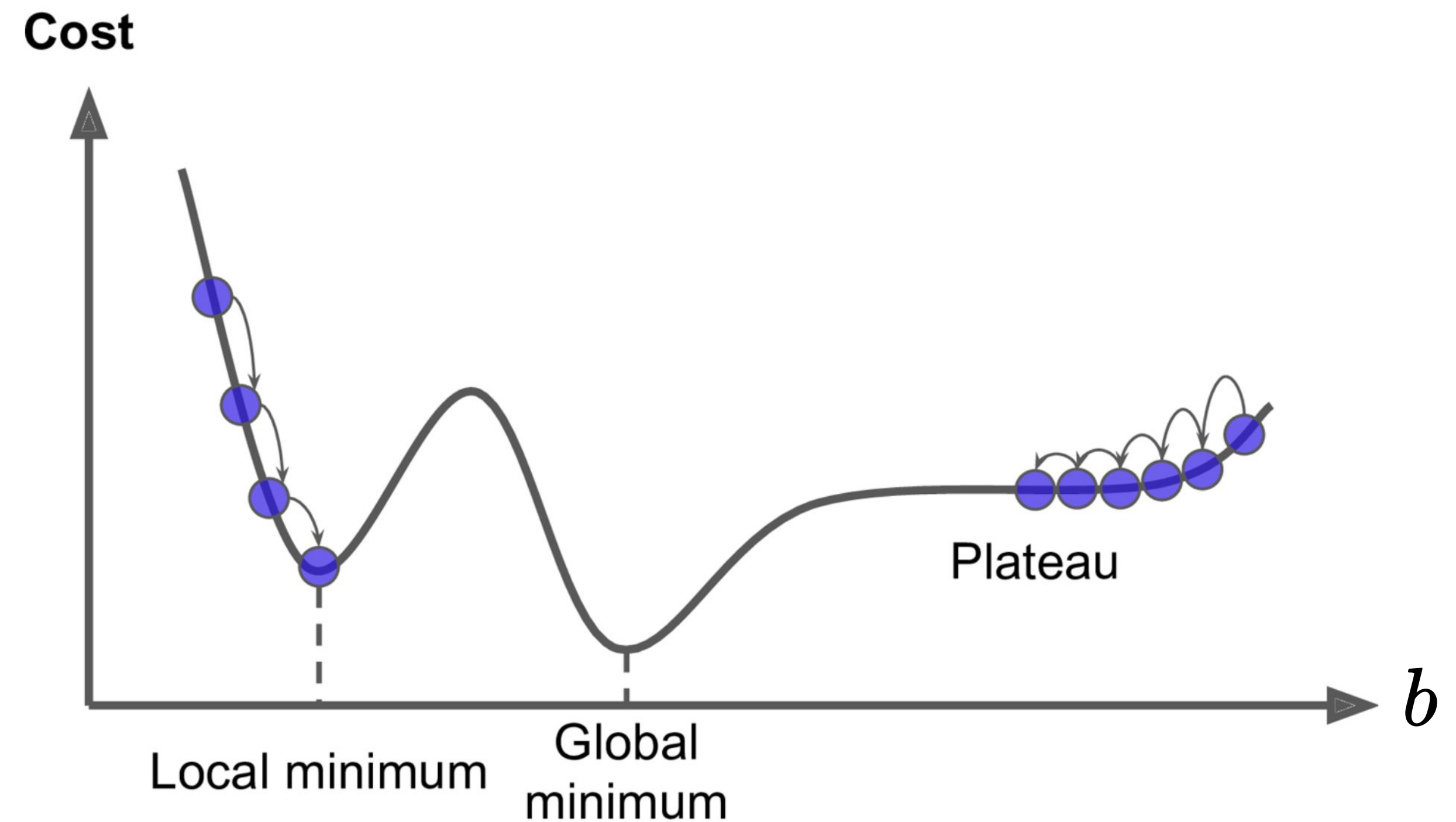
## Problems

# Two main challenges of Gradient Descent

The MSE cost function for linear regression is a convex function, it just has one global minimum.

This is not always the case. Problems are:

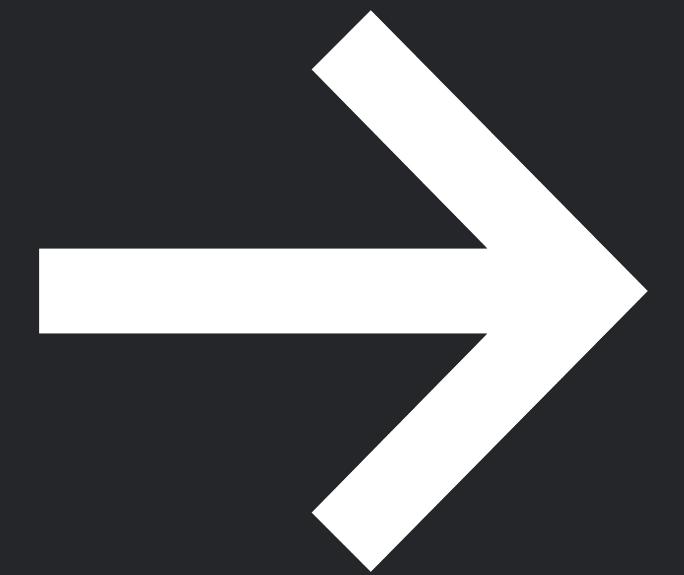
- local minima
- plateaus



*Convex: If you pick two random points the line that is connecting them will never be below the curve.*

Gradient Descent

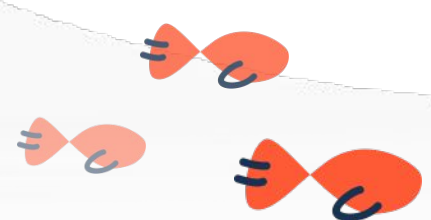
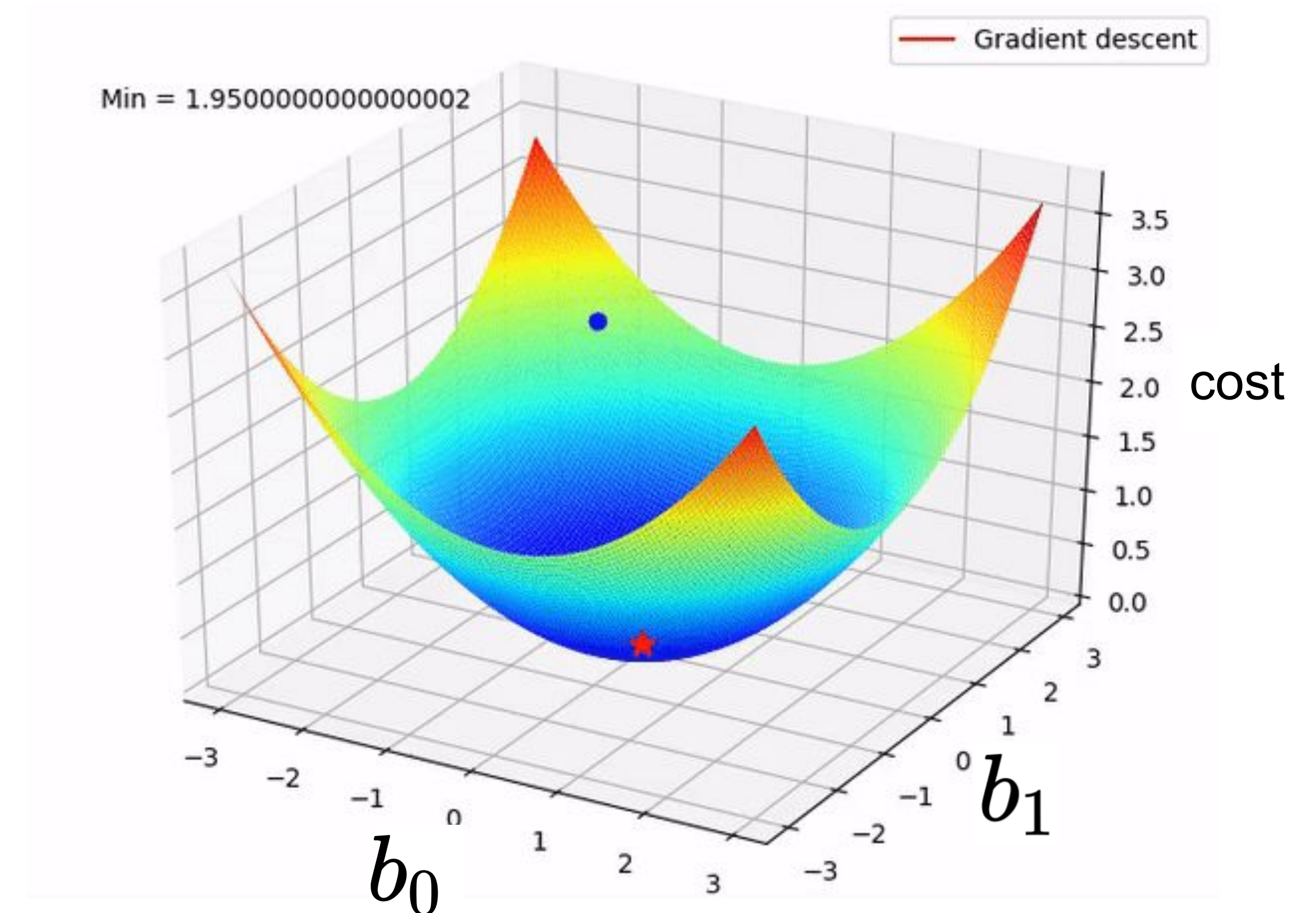
# Part 4 Summary



## Algorithm - Summary

# Gradient Descent summed up

1. Define model
2. Define the cost function
3. Deliberately set some starting values
4. Start descent:
  - Take derivatives with respect to parameters
  - Set your learning rate (step-size)
  - Adjust your parameters (step)
5. Repeat 4. till there is no further improvement





# References

[Gradient Descent Step by Step](#)

