

Machine Learning Tutorial

by Paulius Sasnauskas

sasnausk@ualberta.ca

2025-10-02

Part I: Fundamentals

Part II: Example Problems

Machine Learning

Machine Learning

- learn from data



Machine Learning

- learn from data
- generalize to unseen data



Machine Learning

- learn from data
- generalize to unseen data
- without explicit instructions



```
def has_ears(image):  
    # find oval in image center  
    circle = image.matchOval(image, x=100, y=100, w=30, h=20)  
  
    # find two triangles  
    triangle1 = image.matchTriangle(circle.x-20, circle.y)  
    triangle2 = image.matchTriangle(circle.x+20, circle.y)  
  
    if triangle1 and triangle2:  
        return True  
    else:  
        return False
```

What Problems Can ML Solve?

- Many labeled examples
- Patterns in the data
- Complex relationships

Model

Model – a simplification of a real process or system.

Model

Model – a simplification of a real process or system.

x	y
2	19.6
4.5	44.1
10	98

Model

Model – a simplification of a real process or system.

x	y
2	19.6
4.5	44.1
10	98

$$y = ?$$

Model

Model – a simplification of a real process or system.

Usually we want to know what is the process that generates the data.

x	y
2	19.6
4.5	44.1
10	98

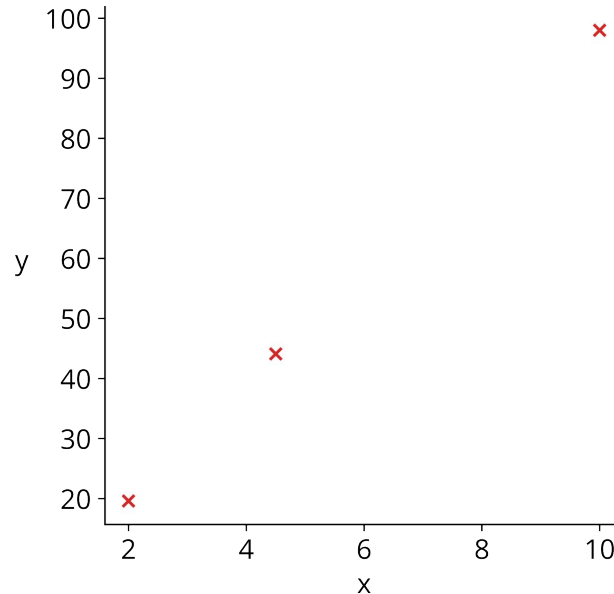
$$y = ?$$

Model

Model – a simplification of a real process or system.

Usually we want to know what is the process that generates the data.

x	y
2	19.6
4.5	44.1
10	98



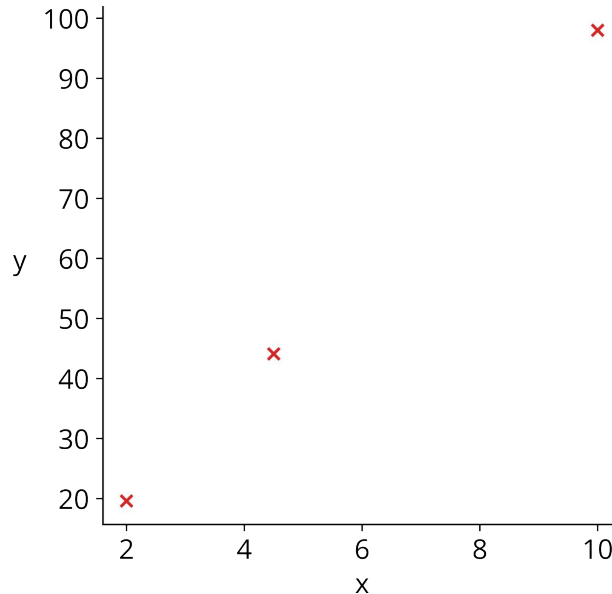
$y = ?$

Model

Model – a simplification of a real process or system.

Usually we want to know what is the process that generates the data.

x	y
2	19.6
4.5	44.1
10	98



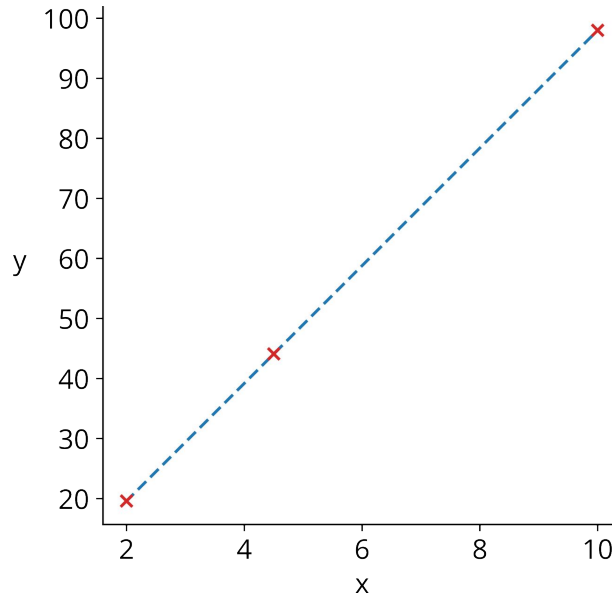
$$y = ax$$

Model

Model – a simplification of a real process or system.

Usually we want to know what is the process that generates the data.

x	y
2	19.6
4.5	44.1
10	98



$$y = ax$$

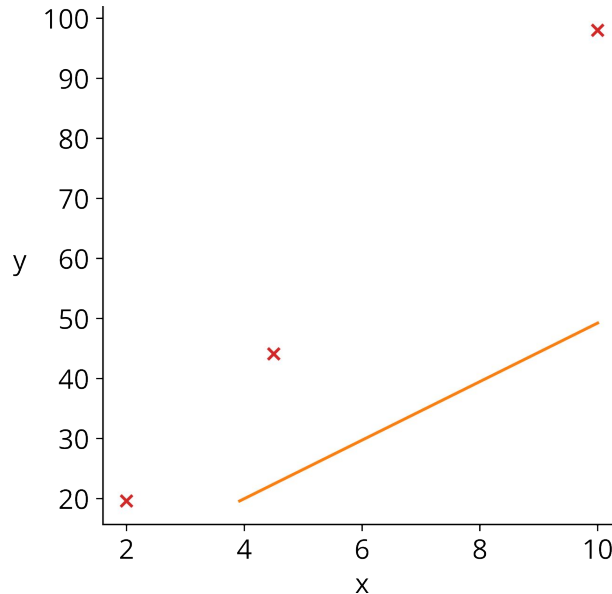
$$a = ?$$

Model

Model – a simplification of a real process or system.

Usually we want to know what is the process that generates the data.

x	y
2	19.6
4.5	44.1
10	98



$$y = ax$$

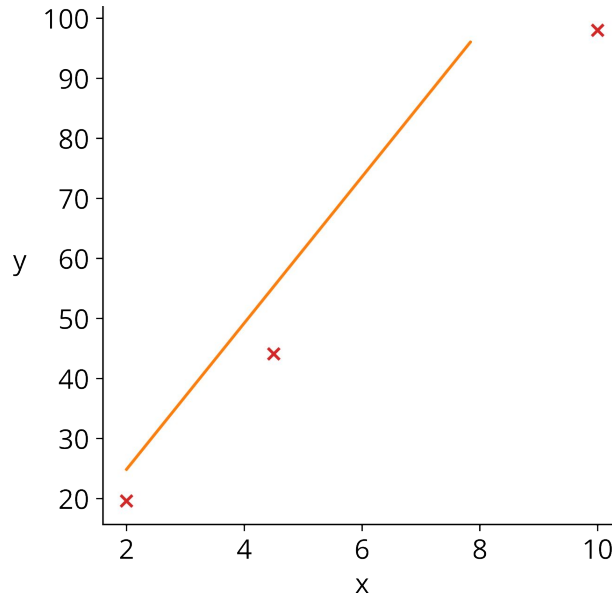
$$a = 5?$$

Model

Model – a simplification of a real process or system.

Usually we want to know what is the process that generates the data.

x	y
2	19.6
4.5	44.1
10	98



$$y = ax$$

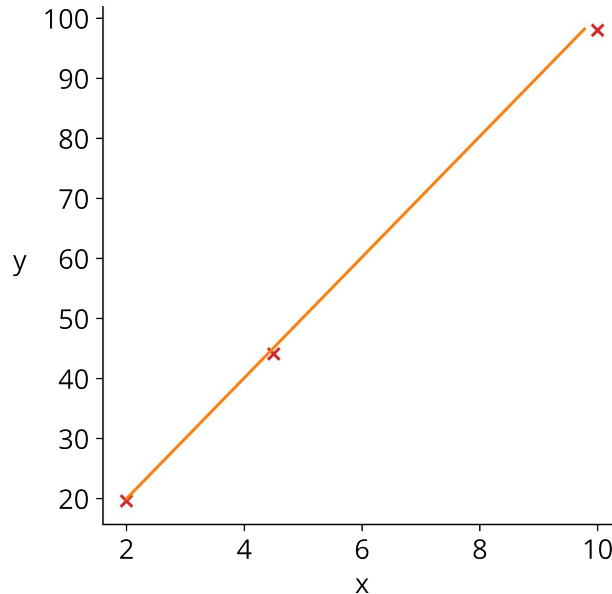
$$a = 12?$$

Model

Model – a simplification of a real process or system.

Usually we want to know what is the process that generates the data.

x	y
2	19.6
4.5	44.1
10	98



$$y = ax$$

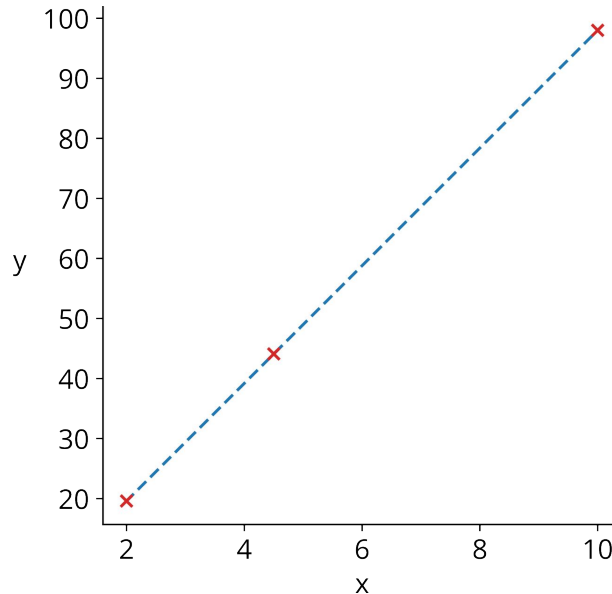
$$a = 10?$$

Model

Model – a simplification of a real process or system.

Usually we want to know what is the process that generates the data.

x	y
2	19.6
4.5	44.1
10	98



$$y = ax$$

$$a = 9.8$$

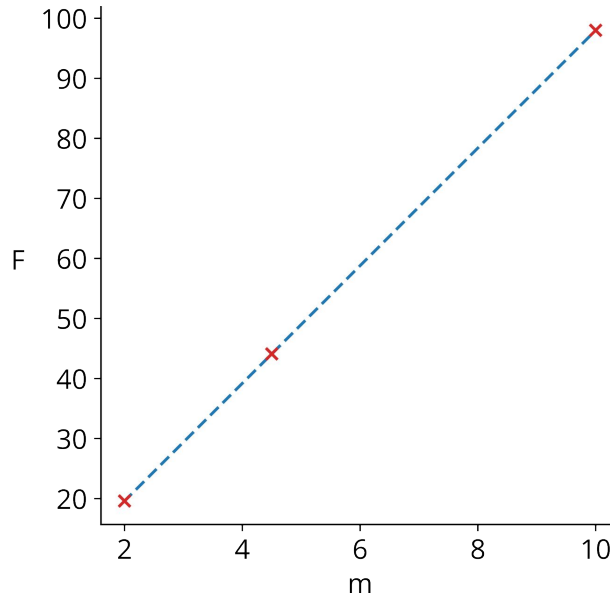
$$y = 9.8x$$

Model

Model – a simplification of a real process or system.

Usually we want to know what is the process that generates the data.

m	F
2	19.6
4.5	44.1
10	98



$$F = mg$$

$$g = 9.8$$

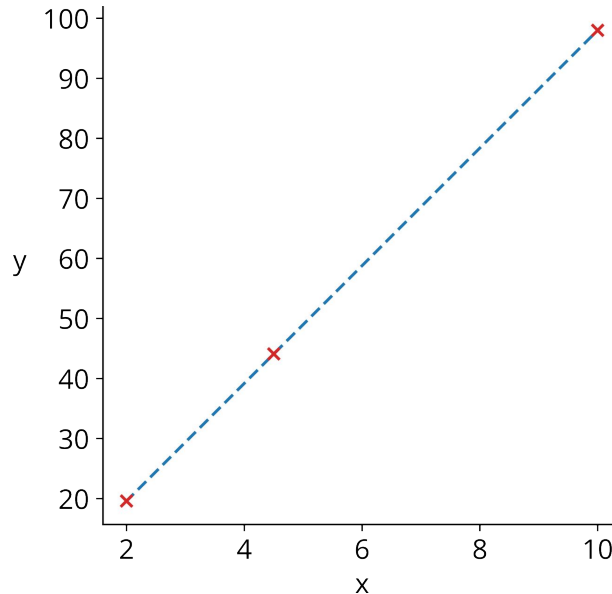
$$F = 9.8 m$$

Model

Model – a simplification of a real process or system.

Usually we want to know what is the process that generates the data.

x	y
2	19.6
4.5	44.1
10	98



$$y = ax$$

$$a = 9.8$$

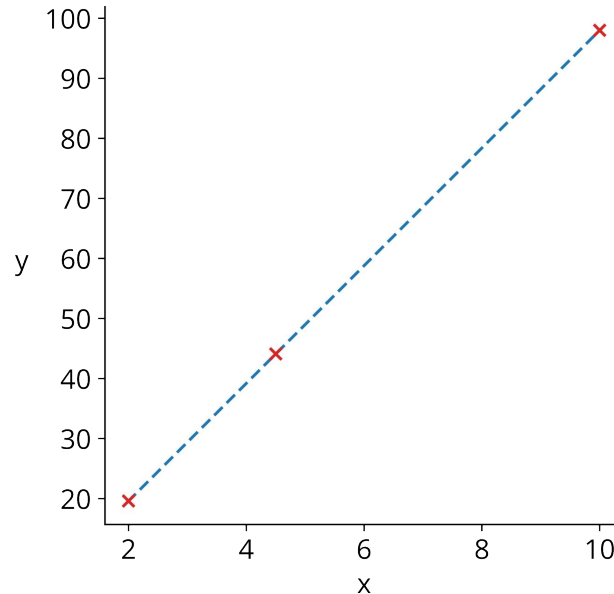
$$y = 9.8x$$

Model

Model – a simplification of a real process or system.

Usually we want to know what is the process that generates the data.

x	y
2	19.6
4.5	44.1
10	98



$$\hat{y} = ax$$

$$a = 9.8$$

$$\hat{y} = 9.8x$$

Model

(x, y) – data
 \hat{y} – prediction

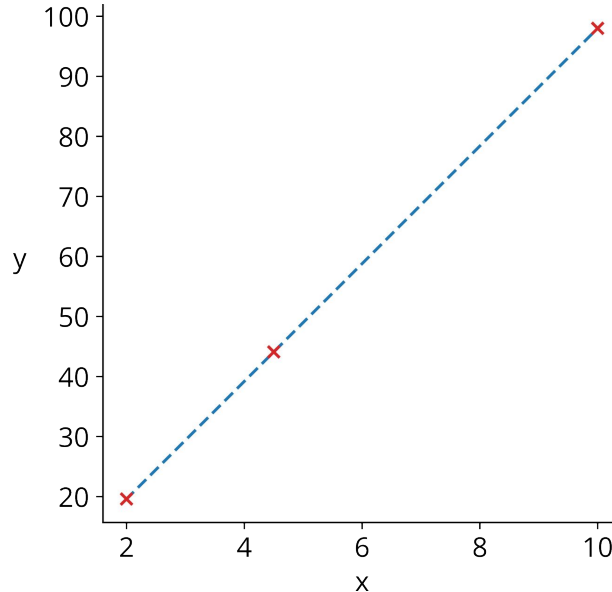


y_p ↗

Model – a simplification of a real process or system.

Usually we want to know what is the process that generates the data.

x	y
2	19.6
4.5	44.1
10	98



$$\hat{y} = ax$$

$$a = 9.8$$

$$\hat{y} = 9.8x$$

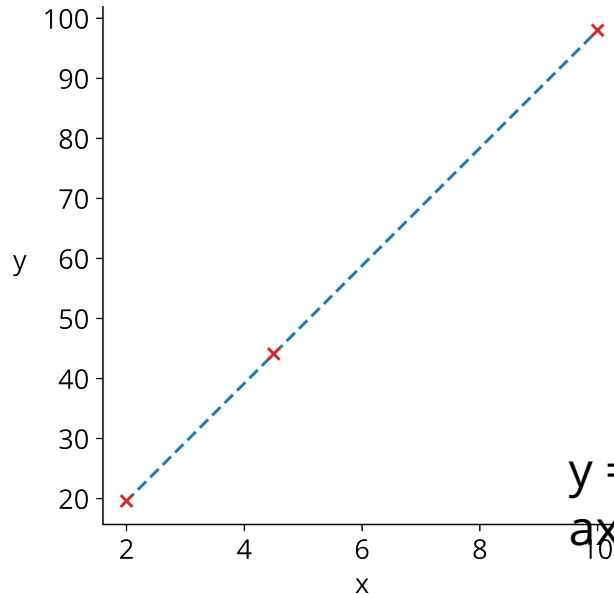
Model

(x, y) – data
 \hat{y} – prediction

Model – a simplification of a real process or system.

Usually we want to know what is the process that generates the data.

x	y	z
2	19.6	?
4.5	44.1	?
10	98	?



$$\hat{y} = ax$$

$$a = 9.8$$

$$\hat{y} = 9.8x$$

$$y = (1 + 1/(1 - z^2) z z^T / b^2) (1 - z^2)^{-1/2}$$

$$b \approx 2.99 \times 10^8$$

Model

(x, y) – data
 \hat{y} – prediction

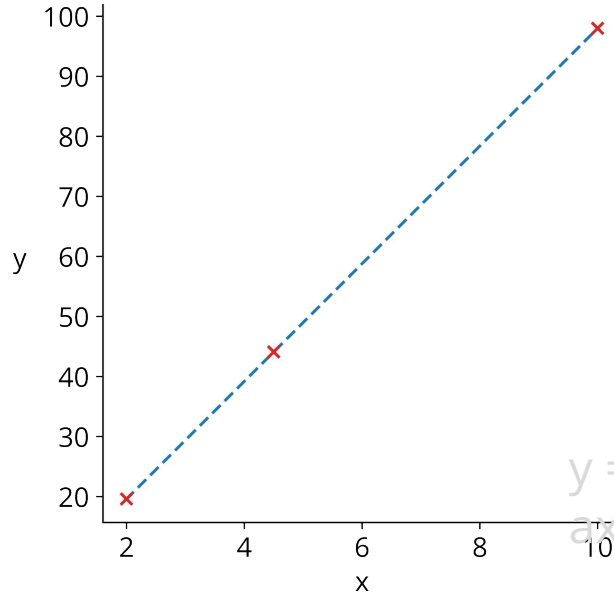
Model – a simplification of a real process or system.

Usually we want to know what is the process that generates the data.

x	y	z
2	19.6	?
4.5	44.1	?
10	98	?

$\hat{y} \approx y$

→



→

$$\hat{y} = ax$$

$$a = 9.8$$

$$\hat{y} = 9.8x$$

$$y = (1 + 1/(1 - z^2) z z^T / b^2) (1 - z^2)^{-1/2}$$

$$b \approx 2.99 \times 10^8$$

Model

(x, y) – data
 \hat{y} – prediction

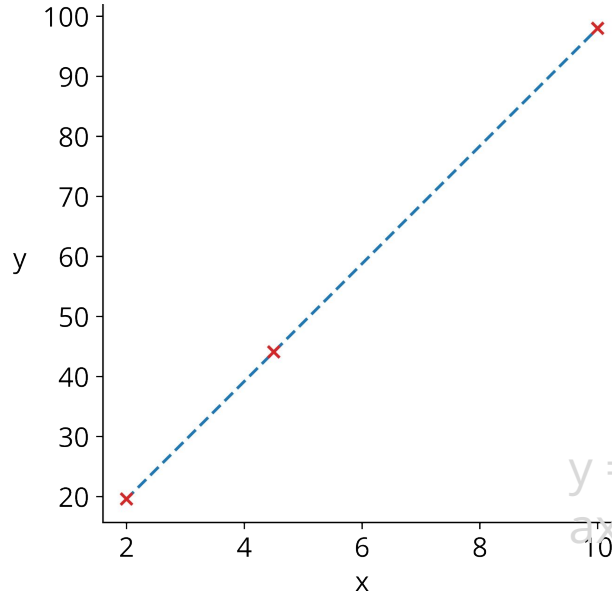
Model – a simplification of a real process or system.

Usually *it suffices to approximate* what is the process that generates the data.

x	y	z
2	19.6	?
4.5	44.1	?
10	98	?

$\hat{y} \approx y$

→



$$\hat{y} = ax$$

$$a = 9.8$$

$$\hat{y} = 9.8x$$

$$y = (1 + 1/(1 - z^2) z z^T / b^2) (1 - z^2)^{-1/2}$$

$$b \approx 2.99 \times 10^8$$

Loss Function

What is a good approximation?

Loss Function

What is a good approximation?

Is it possible to evaluate “goodness”?

Loss Function

What is a good approximation?

Is it possible to evaluate “goodness”?



x	y
2	19.6
4.5	44.1
10	98

Loss Function

What is a good approximation?

Is it possible to evaluate “goodness”?

Distance from the data to the model:

x	y
2	19.6
4.5	44.1
10	98

$$y_i - \hat{y}(x_i)$$

Loss Function

What is a good approximation?

Is it possible to evaluate “goodness”?

Distance from the data to the model:

x	y
2	19.6
4.5	44.1
10	98

$$\left| y_i - \hat{y}(x_i) \right|$$

Loss Function

What is a good approximation?

Is it possible to evaluate “goodness”?

Distance from the data to the model:

x	y
2	19.6
4.5	44.1
10	98

$$\left(y_i - \hat{y}(x_i)\right)^2$$

Loss Function

What is a good approximation?

Is it possible to evaluate “goodness”?

Distance from the data to the model:

x	y
2	19.6
4.5	44.1
10	98

$$\sum_i \left(y_i - \hat{y}(x_i) \right)^2$$

Loss Function

What is a good approximation?

Is it possible to evaluate “goodness”?

Distance from the data to the model:

x	y
2	19.6
4.5	44.1
10	98

$$\mathcal{L} = \sum_i \left(y_i - \hat{y}(x_i) \right)^2$$

Loss Function

What is a good approximation?

Is it possible to evaluate “goodness”?

Distance from the data to the model:

x	y
2	19.6
4.5	44.1
10	98

$$\mathcal{L} = \sum_i \left(y_i - \hat{y}(x_i) \right)^2 \quad \text{If } \mathcal{L} = 0 \text{ then } y_i - \hat{y}(x_i) = 0 \text{ for all } i$$

Loss Function

What is a good approximation?

Is it possible to evaluate “goodness”?

Distance from the data to the model:

x	y
2	19.6
4.5	44.1
10	98

$$\mathcal{L} = \sum_i \left(y_i - \hat{y}(x_i) \right)^2 \quad \text{If } \mathcal{L} = 0 \text{ then } y_i = \hat{y}(x_i) \text{ for all } i$$

Loss Function

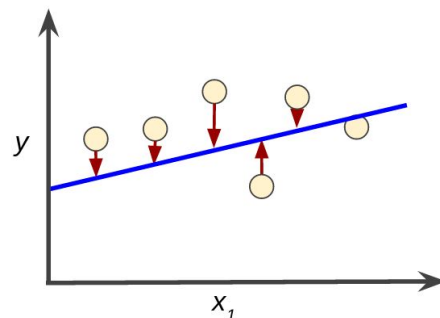
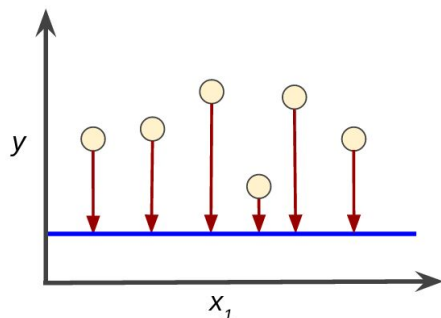
What is a good approximation?

Is it possible to evaluate “goodness”?

Distance from the data to the model:

x	y
2	19.6
4.5	44.1
10	98

$$\mathcal{L} = \sum_i \left(y_i - \hat{y}(x_i) \right)^2 \quad \text{If } \mathcal{L} = 0 \text{ then } y_i = \hat{y}(x_i) \text{ for all } i$$



Loss Function

What is a good approximation?

Is it possible to evaluate “goodness”?

Distance from the data to the model:

x	y
2	19.6
4.5	44.1
10	98

$$\mathcal{L} = \sum_i \left(y_i - \hat{y}(x_i) \right)^2 \quad \text{L2 Loss} \quad \img alt="pencil icon" data-bbox="565 548 595 620"/>$$

Loss Function

What is a good approximation?

Is it possible to evaluate “goodness”?

Distance from the data to the model:

x	y
2	19.6
4.5	44.1
10	98

$$\mathcal{L} = \sum_i \left(y_i - \hat{y}(x_i) \right)^2 \quad \text{L2 Loss} \quad \img alt="pencil icon" data-bbox="565 550 595 620"/>$$

$$\hat{y} = a x$$

Loss Function

What is a good approximation?

Is it possible to evaluate “goodness”?

Distance from the data to the model:

x	y
2	19.6
4.5	44.1
10	98

$$\mathcal{L} = \sum_i \left(y_i - \hat{y}(x_i) \right)^2 \quad \text{L2 Loss} \quad \img alt="pencil icon" data-bbox="565 548 595 622"/>$$

$$\hat{y} = a x \quad \mathcal{L} = (y_1 - a \cdot x_1)^2 + (y_2 - a \cdot x_2)^2 + (y_2 - a \cdot x_2)^2$$

Loss Function

What is a good approximation?

Is it possible to evaluate “goodness”?

Distance from the data to the model:

x	y
2	19.6
4.5	44.1
10	98

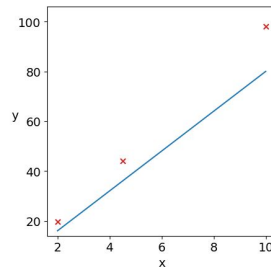
$$\mathcal{L} = \sum_i \left(y_i - \hat{y}(x_i) \right)^2 \quad \text{L2 Loss} \quad \img alt="pencil icon" data-bbox="565 548 596 621"/>$$

$$\hat{y} = a x \quad \mathcal{L} = (19.6 - a \cdot 2)^2 + (44.1 - a \cdot 4.5)^2 + (98 - a \cdot 10)^2$$

Loss Function

What is a good approximation?

Is it possible to evaluate “goodness”?



Distance from the data to the model:

x	y
2	19.6
4.5	44.1
10	98

$$\mathcal{L} = \sum_i \left(y_i - \hat{y}(x_i) \right)^2 \quad \text{L2 Loss} \quad \img alt="pencil icon" data-bbox="565 548 595 622"/>$$

$$\hat{y} = 8x \quad \mathcal{L} = (19.6 - 8 \cdot 2)^2 + (44.1 - 8 \cdot 4.5)^2 + (98 - 8 \cdot 10)^2$$

Loss Function

What is a good approximation?

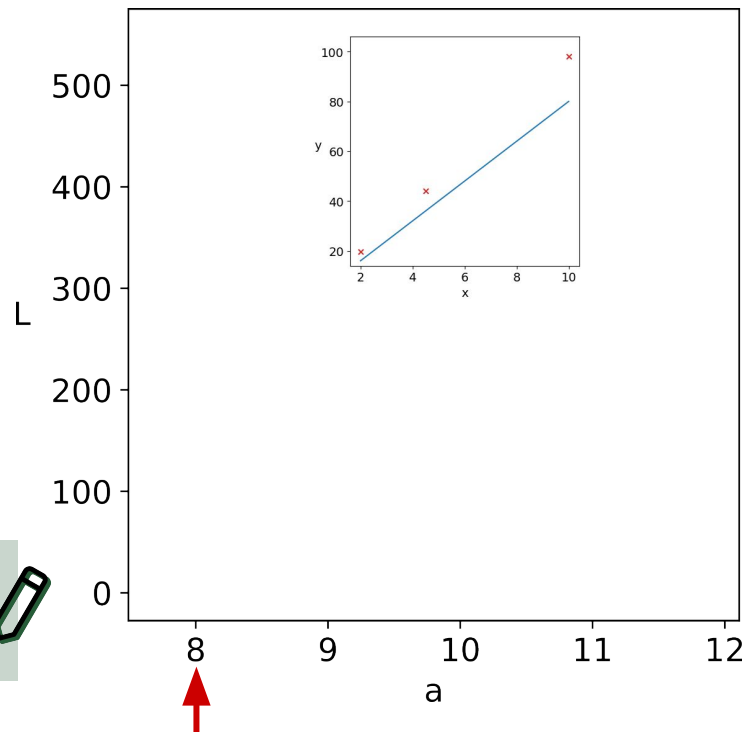
Is it possible to evaluate “goodness”?

Distance from the data to the model:

x	y
2	19.6
4.5	44.1
10	98

$$\mathcal{L} = \sum_i \left(y_i - \hat{y}(x_i) \right)^2 \quad \text{L2 Loss}$$

$$\hat{y} = 8x \quad \mathcal{L} = (19.6 - 8 \cdot 2)^2 + (44.1 - 8 \cdot 4.5)^2 + (98 - 8 \cdot 10)^2$$



Loss Function

What is a good approximation?

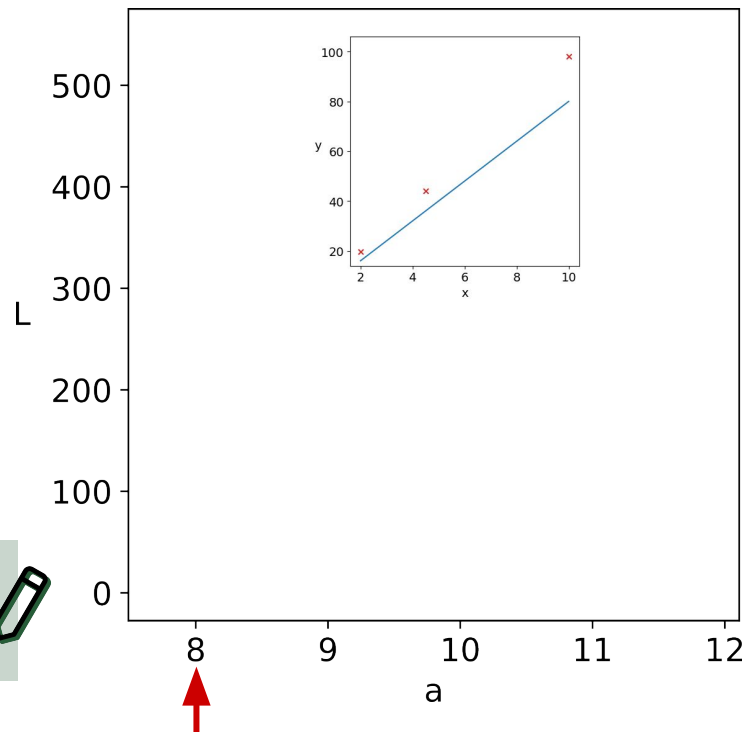
Is it possible to evaluate “goodness”?

Distance from the data to the model:

x	y
2	19.6
4.5	44.1
10	98

$$\mathcal{L} = \sum_i \left(y_i - \hat{y}(x_i) \right)^2 \quad \text{L2 Loss}$$

$$\hat{y} = 8x \quad \mathcal{L} = (19.6 - 8 \cdot 2)^2 + (44.1 - 8 \cdot 4.5)^2 + (98 - 8 \cdot 10)^2 = 402.57$$



Loss Function

What is a good approximation?

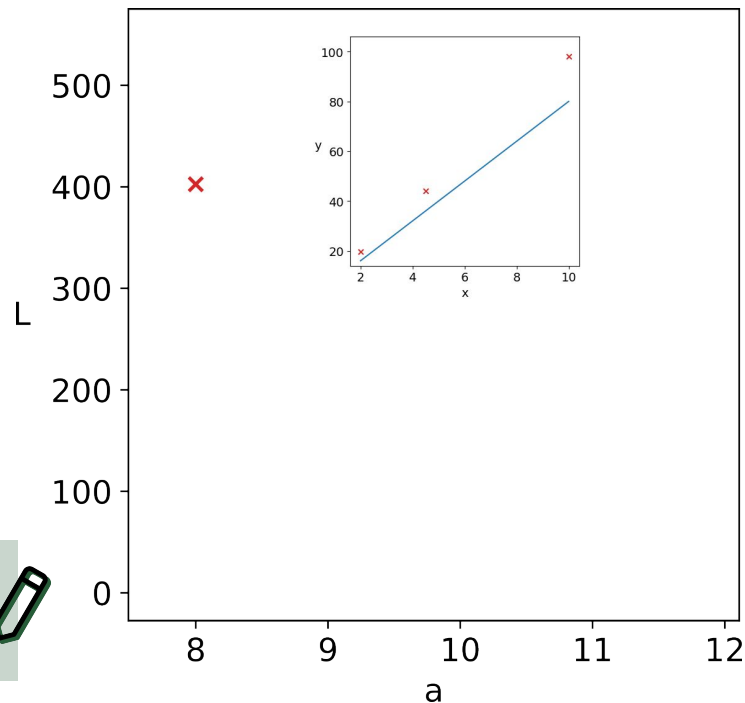
Is it possible to evaluate “goodness”?

Distance from the data to the model:

x	y
2	19.6
4.5	44.1
10	98

$$\mathcal{L} = \sum_i \left(y_i - \hat{y}(x_i) \right)^2 \quad \text{L2 Loss} \quad \text{📎}$$

$$\hat{y} = 8x \quad \mathcal{L} = (19.6 - 8 \cdot 2)^2 + (44.1 - 8 \cdot 4.5)^2 + (98 - 8 \cdot 10)^2 = 402.57$$



Loss Function

What is a good approximation?

Is it possible to evaluate “goodness”?

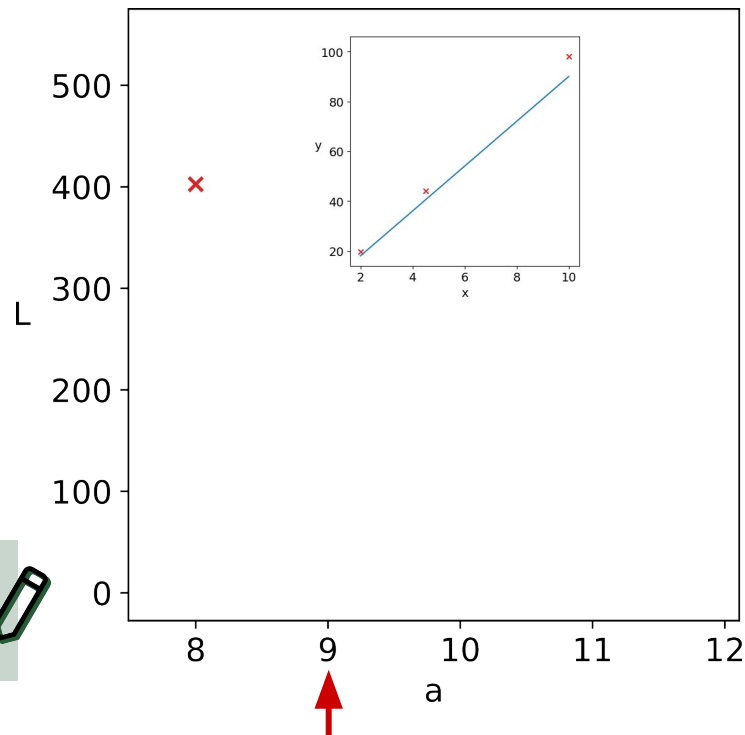
Distance from the data to the model:

x	y
2	19.6
4.5	44.1
10	98

$$\mathcal{L} = \sum_i \left(y_i - \hat{y}(x_i) \right)^2 \quad \text{L2 Loss} \quad \text{📎}$$

$$\hat{y} = 8x \quad \mathcal{L} = (19.6 - 8 \cdot 2)^2 + (44.1 - 8 \cdot 4.5)^2 + (98 - 8 \cdot 10)^2 = 402.57$$

$$\hat{y} = 9x$$



Loss Function

What is a good approximation?

Is it possible to evaluate “goodness”?

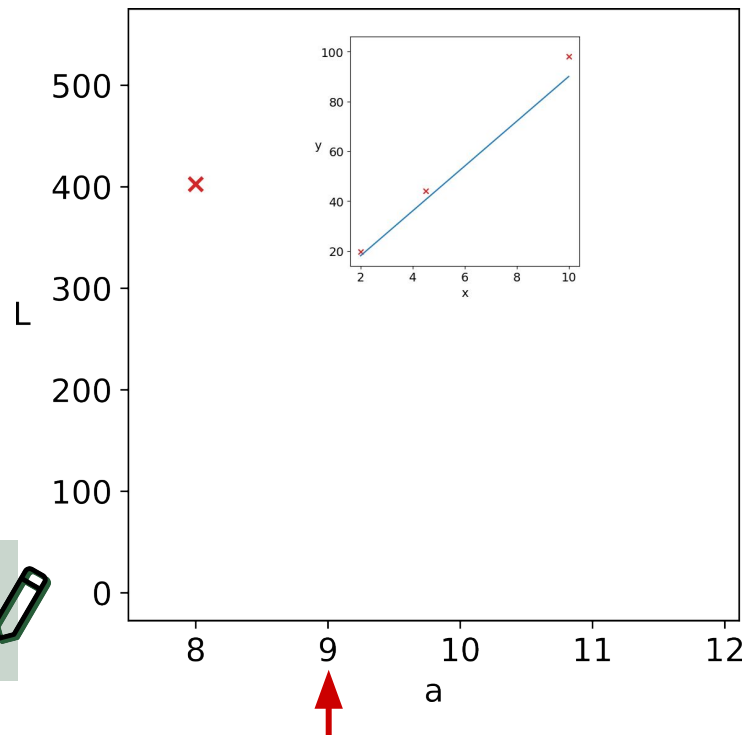
Distance from the data to the model:

x	y
2	19.6
4.5	44.1
10	98

$$\mathcal{L} = \sum_i (y_i - \hat{y}(x_i))^2 \quad \text{L2 Loss}$$

$$\hat{y} = 8x \quad \mathcal{L} = (19.6 - 8 \cdot 2)^2 + (44.1 - 8 \cdot 4.5)^2 + (98 - 8 \cdot 10)^2 = 402.57$$

$$\hat{y} = 9x \quad \mathcal{L} = (19.6 - 9 \cdot 2)^2 + (44.1 - 9 \cdot 4.5)^2 + (98 - 9 \cdot 10)^2$$



Loss Function

What is a good approximation?

Is it possible to evaluate “goodness”?

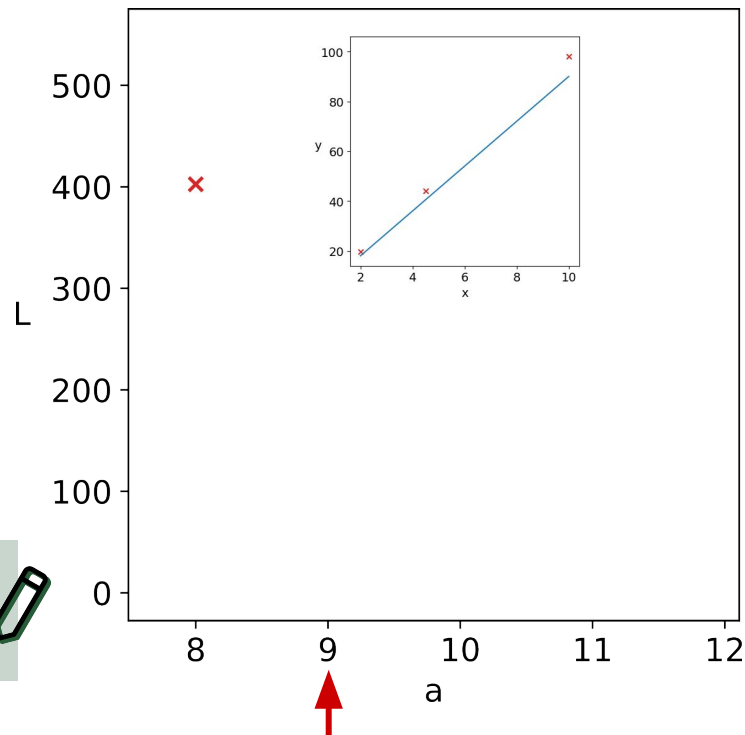
Distance from the data to the model:

x	y
2	19.6
4.5	44.1
10	98

$$\mathcal{L} = \sum_i \left(y_i - \hat{y}(x_i) \right)^2 \quad \text{L2 Loss} \quad \img alt="pencil icon" data-bbox="565 548 595 622"/>$$

$$\hat{y} = 8x \quad \mathcal{L} = (19.6 - 8 \cdot 2)^2 + (44.1 - 8 \cdot 4.5)^2 + (98 - 8 \cdot 10)^2 = 402.57$$

$$\hat{y} = 9x \quad \mathcal{L} = (19.6 - 9 \cdot 2)^2 + (44.1 - 9 \cdot 4.5)^2 + (98 - 9 \cdot 10)^2 = 79.52$$



Loss Function

What is a good approximation?

Is it possible to evaluate “goodness”?

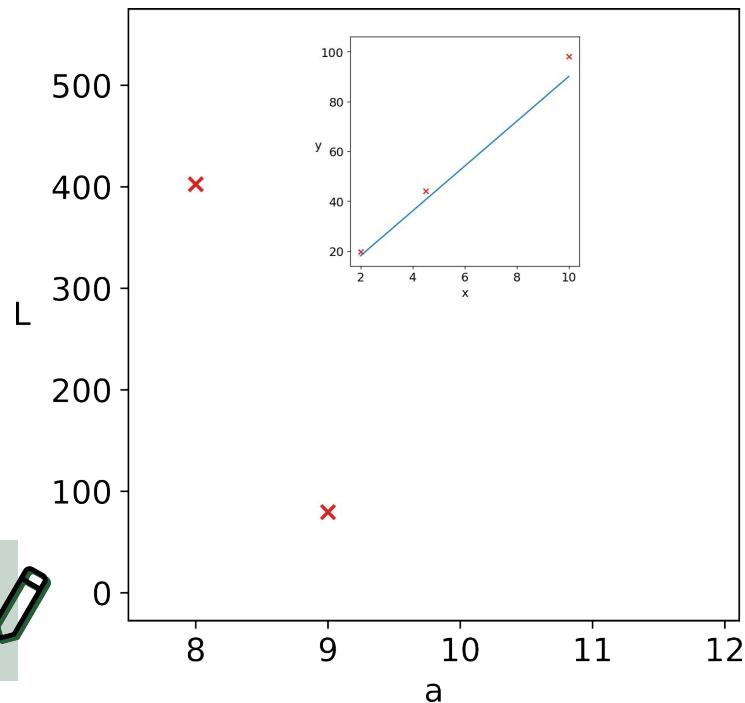
Distance from the data to the model:

x	y
2	19.6
4.5	44.1
10	98

$$\mathcal{L} = \sum_i \left(y_i - \hat{y}(x_i) \right)^2 \quad \text{L2 Loss}$$

$$\hat{y} = 8x \quad \mathcal{L} = (19.6 - 8 \cdot 2)^2 + (44.1 - 8 \cdot 4.5)^2 + (98 - 8 \cdot 10)^2 = 402.57$$

$$\hat{y} = 9x \quad \mathcal{L} = (19.6 - 9 \cdot 2)^2 + (44.1 - 9 \cdot 4.5)^2 + (98 - 9 \cdot 10)^2 = 79.52$$



Loss Function

What is a good approximation?

Is it possible to evaluate “goodness”?

Distance from the data to the model:

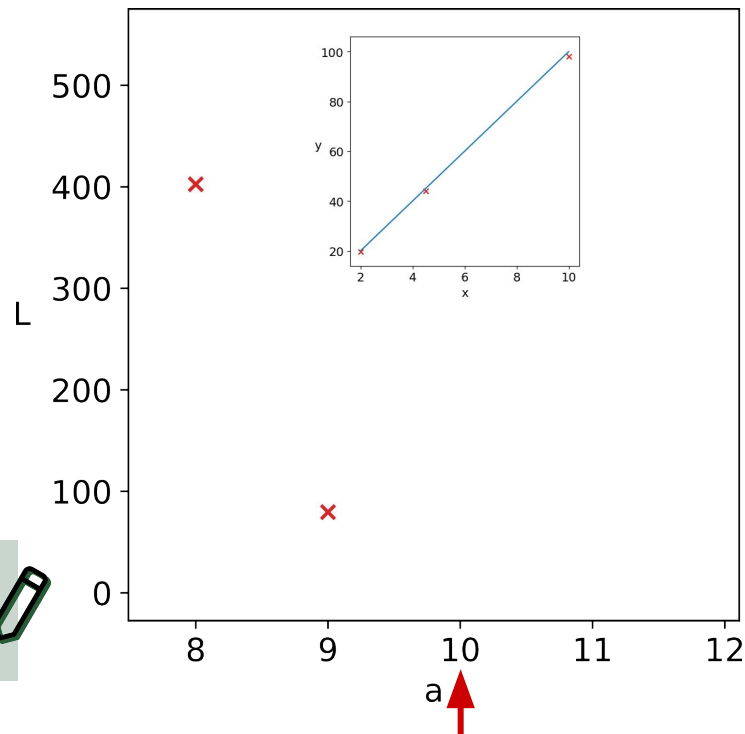
x	y
2	19.6
4.5	44.1
10	98

$$\mathcal{L} = \sum_i \left(y_i - \hat{y}(x_i) \right)^2 \quad \text{L2 Loss}$$

$$\hat{y} = 8x \quad \mathcal{L} = (19.6 - 8 \cdot 2)^2 + (44.1 - 8 \cdot 4.5)^2 + (98 - 8 \cdot 10)^2 = 402.57$$

$$\hat{y} = 9x \quad \mathcal{L} = (19.6 - 9 \cdot 2)^2 + (44.1 - 9 \cdot 4.5)^2 + (98 - 9 \cdot 10)^2 = 79.52$$

$$\hat{y} = 10x$$



Loss Function

What is a good approximation?

Is it possible to evaluate “goodness”?

Distance from the data to the model:

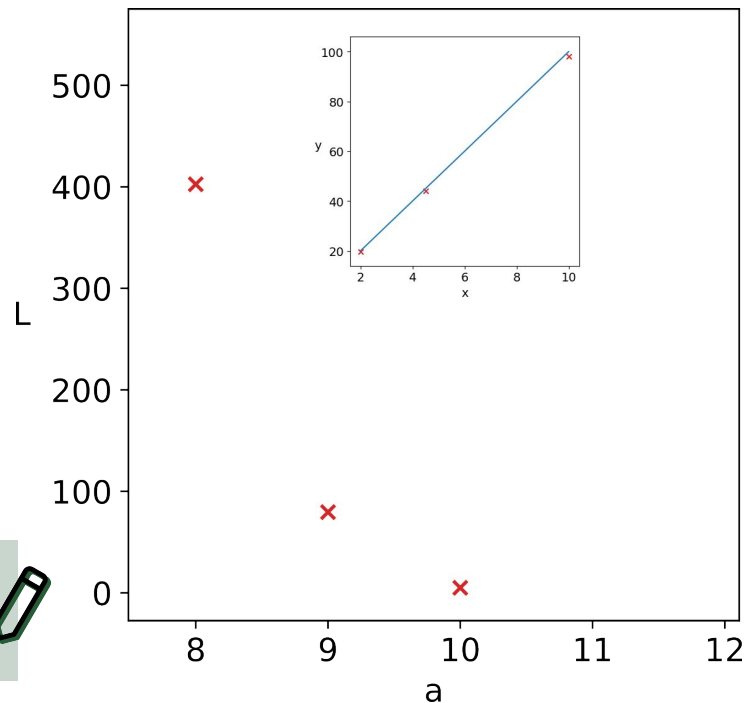
x	y
2	19.6
4.5	44.1
10	98

$$\mathcal{L} = \sum_i \left(y_i - \hat{y}(x_i) \right)^2 \quad \text{L2 Loss} \quad \text{📎}$$

$$\hat{y} = 8x \quad \mathcal{L} = (19.6 - 8 \cdot 2)^2 + (44.1 - 8 \cdot 4.5)^2 + (98 - 8 \cdot 10)^2 = 402.57$$

$$\hat{y} = 9x \quad \mathcal{L} = (19.6 - 9 \cdot 2)^2 + (44.1 - 9 \cdot 4.5)^2 + (98 - 9 \cdot 10)^2 = 79.52$$

$$\hat{y} = 10x \quad \mathcal{L} = (19.6 - 10 \cdot 2)^2 + (44.1 - 10 \cdot 4.5)^2 + (98 - 10 \cdot 10)^2 = 4.97$$



Loss Function

What is a good approximation?

Is it possible to evaluate “goodness”?

Distance from the data to the model:

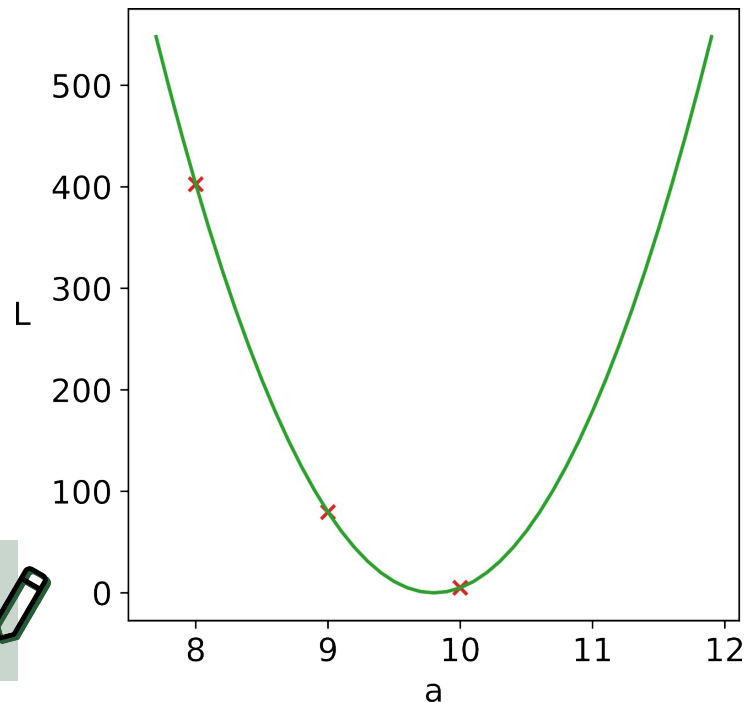
x	y
2	19.6
4.5	44.1
10	98

$$\mathcal{L} = \sum_i \left(y_i - \hat{y}(x_i) \right)^2 \quad \text{L2 Loss} \quad \img alt="pencil icon" data-bbox="565 548 595 622"/>$$

$$\hat{y} = 8x \quad \mathcal{L} = (19.6 - 8 \cdot 2)^2 + (44.1 - 8 \cdot 4.5)^2 + (98 - 8 \cdot 10)^2 = 402.57$$

$$\hat{y} = 9x \quad \mathcal{L} = (19.6 - 9 \cdot 2)^2 + (44.1 - 9 \cdot 4.5)^2 + (98 - 9 \cdot 10)^2 = 79.52$$

$$\hat{y} = 10x \quad \mathcal{L} = (19.6 - 10 \cdot 2)^2 + (44.1 - 10 \cdot 4.5)^2 + (98 - 10 \cdot 10)^2 = 4.97$$



Loss Function

What is a good approximation?

Is it possible to evaluate “goodness”?

Distance from the data to the model:

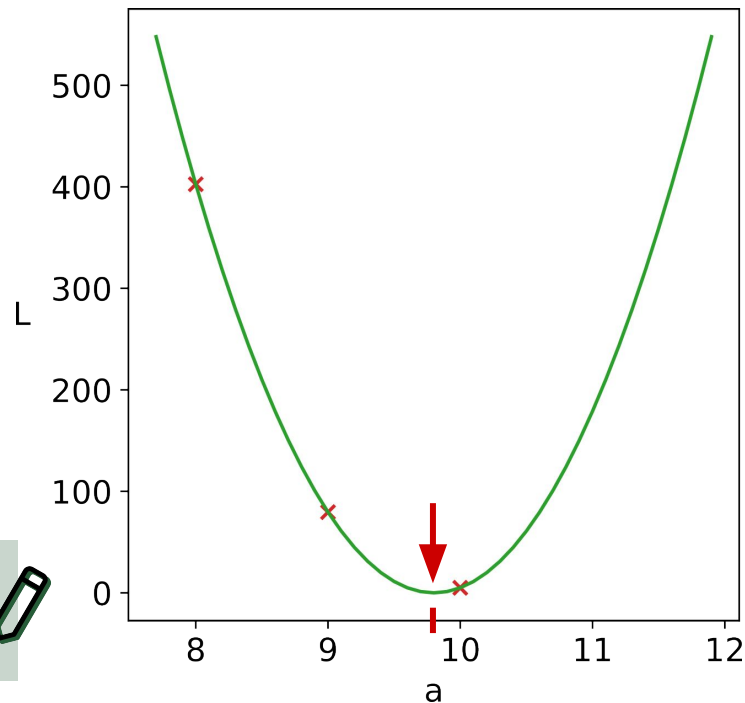
x	y
2	19.6
4.5	44.1
10	98

$$\mathcal{L} = \sum_i (y_i - \hat{y}(x_i))^2 \quad \text{L2 Loss}$$

$$\hat{y} = 8x \quad \mathcal{L} = (19.6 - 8 \cdot 2)^2 + (44.1 - 8 \cdot 4.5)^2 + (98 - 8 \cdot 10)^2 = 402.57$$

$$\hat{y} = 9x \quad \mathcal{L} = (19.6 - 9 \cdot 2)^2 + (44.1 - 9 \cdot 4.5)^2 + (98 - 9 \cdot 10)^2 = 79.52$$

$$\hat{y} = 10x \quad \mathcal{L} = (19.6 - 10 \cdot 2)^2 + (44.1 - 10 \cdot 4.5)^2 + (98 - 10 \cdot 10)^2 = 4.97$$



Gradient Descent

How do we find the value for **a**?

$$\hat{y} = ax$$



Gradient Descent

How to find the minimum?

\mathcal{L}

Gradient Descent

How to find the minimum? $\min \mathcal{L}$

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$ $a^* = \arg \min_a \mathcal{L}(a)$

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

?

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

\mathcal{L}'

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

$$\mathcal{L}' \quad \mathcal{L}'_x? \quad \mathcal{L}'_a?$$

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

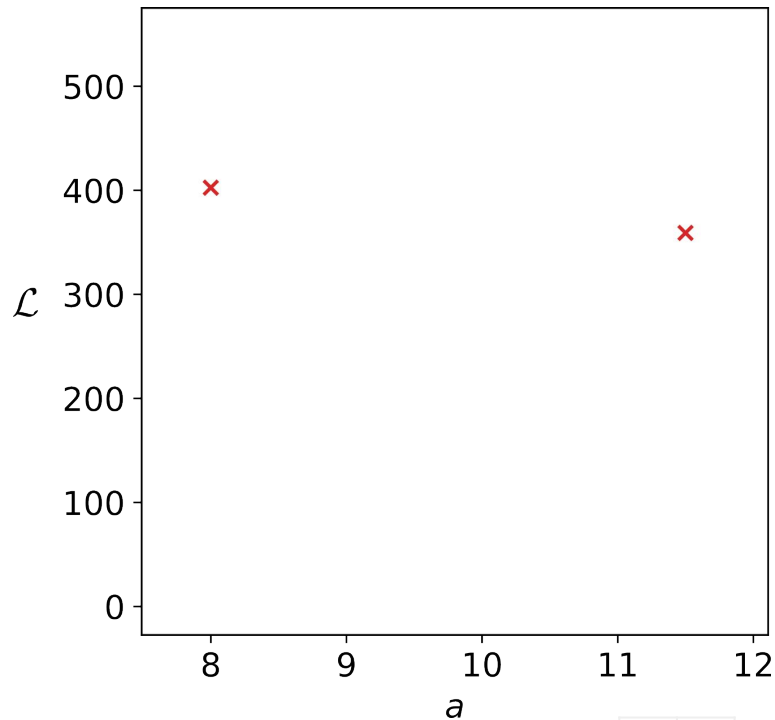
$$\frac{\partial \mathcal{L}}{\partial a}$$

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

$$\frac{\partial \mathcal{L}}{\partial a}$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

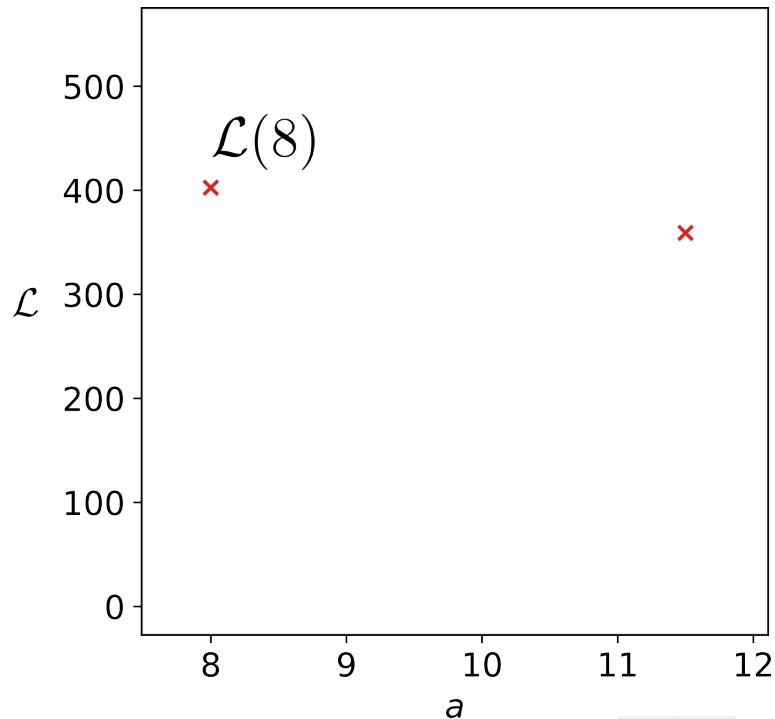
Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

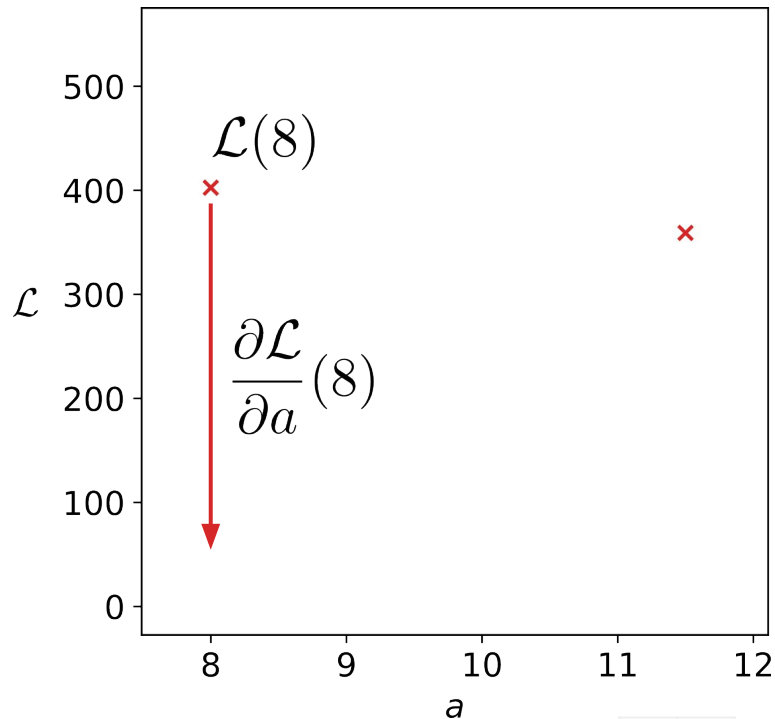
How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

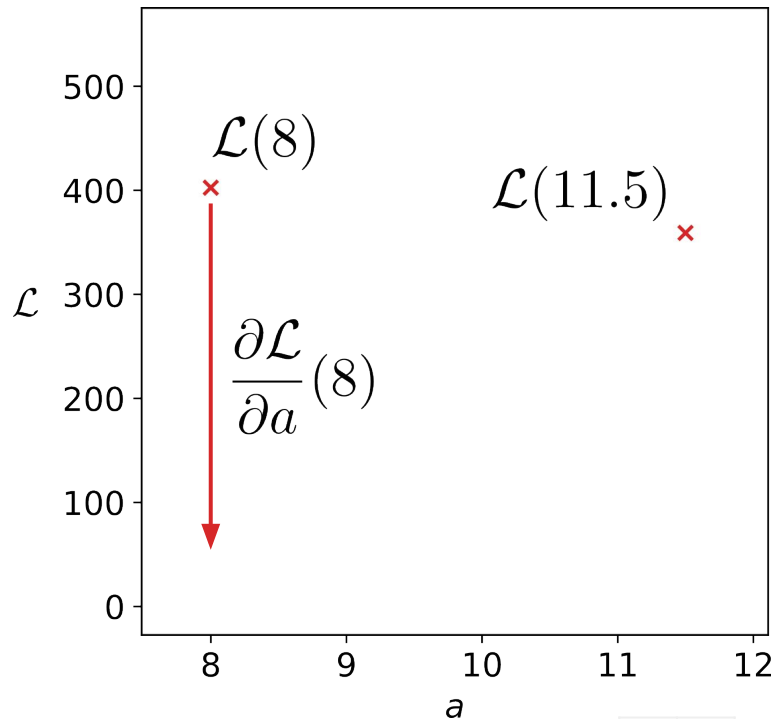
Derivatives show “influence” to functions

$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

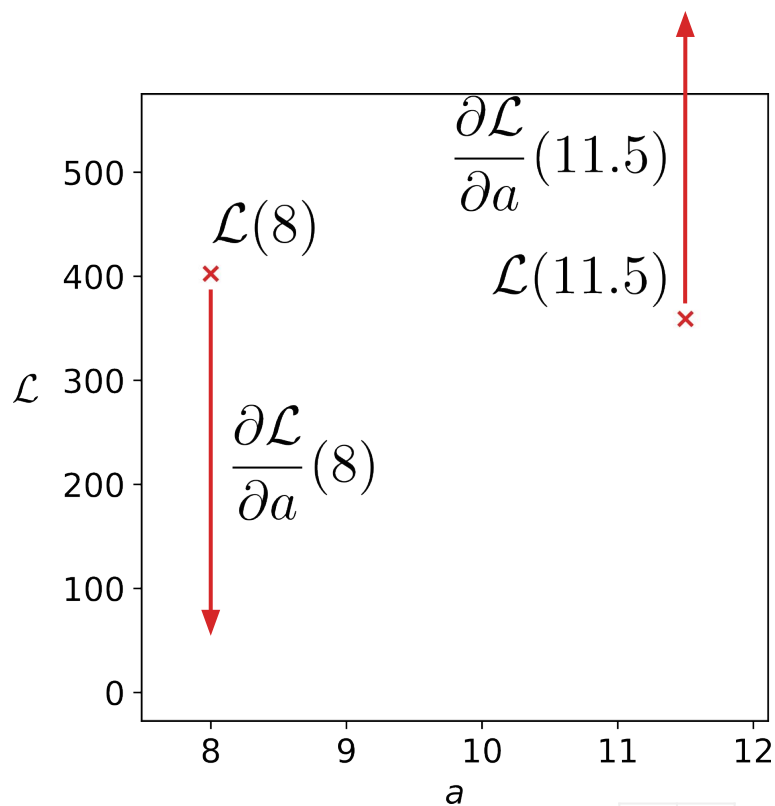
$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(11.5) \approx 422$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

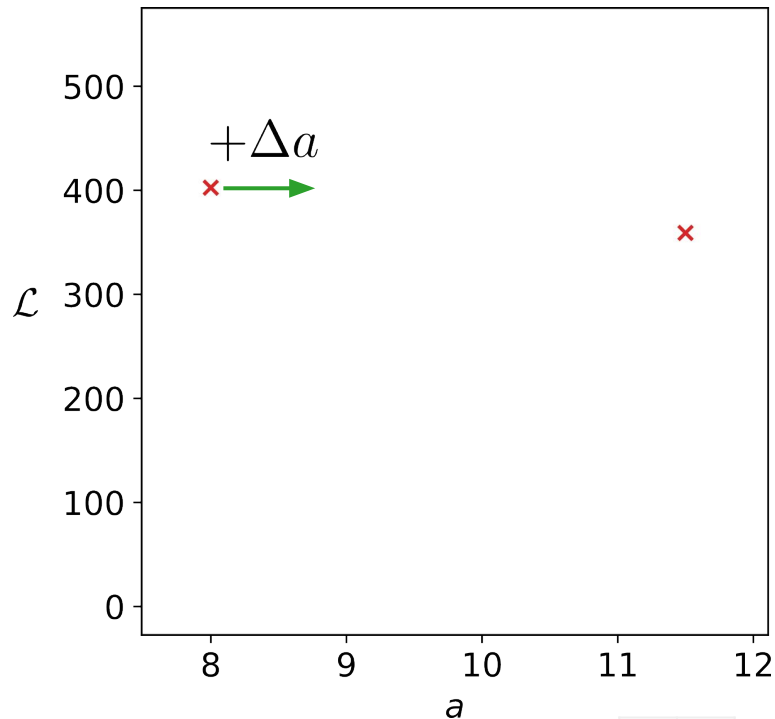
$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a) \quad \Delta a > 0$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

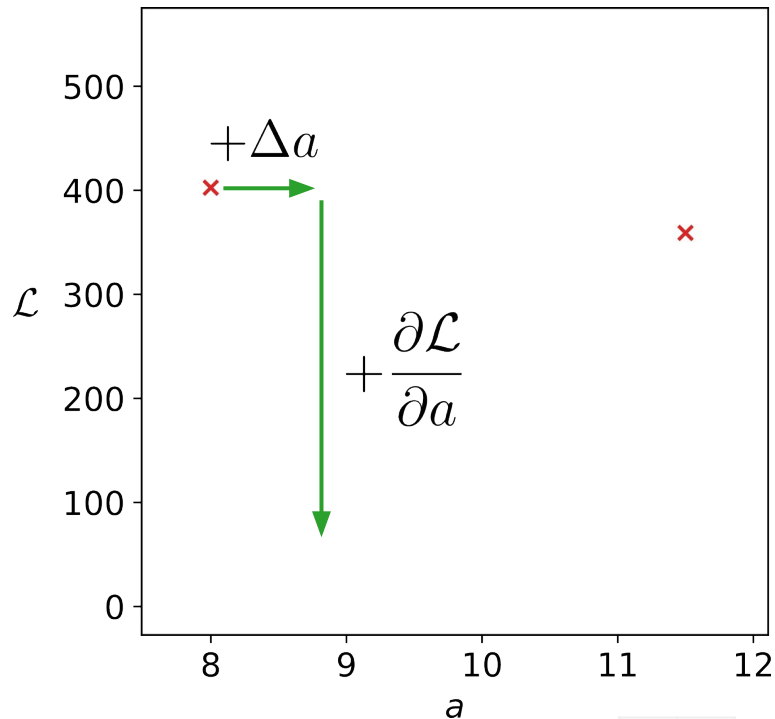
$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a) \quad \Delta a > 0$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

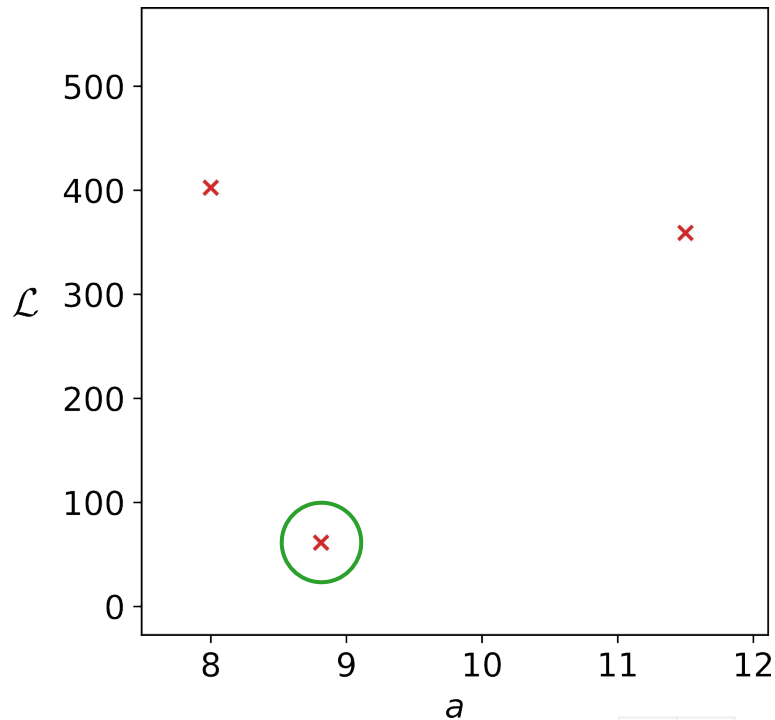
$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a) \quad \Delta a > 0$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

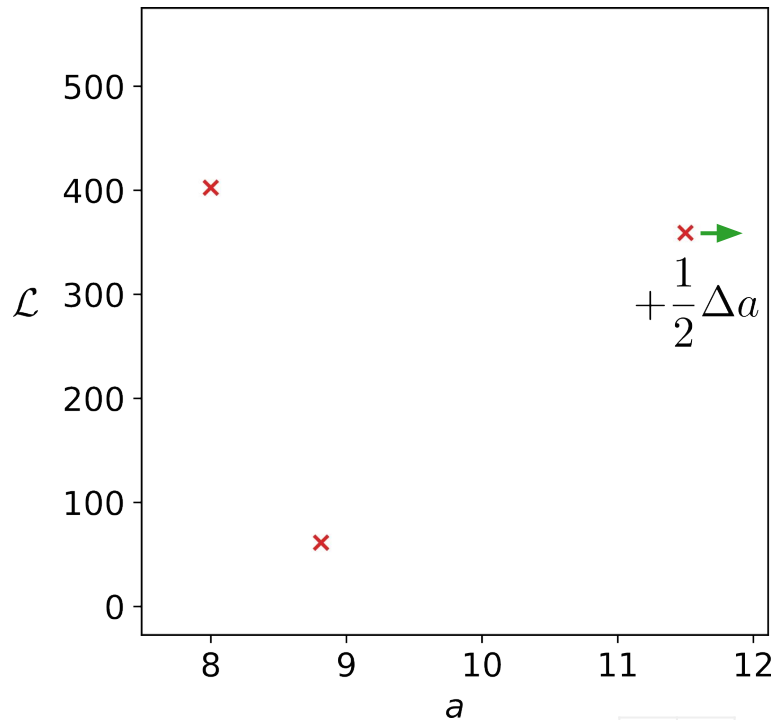
$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a) \quad \Delta a > 0$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

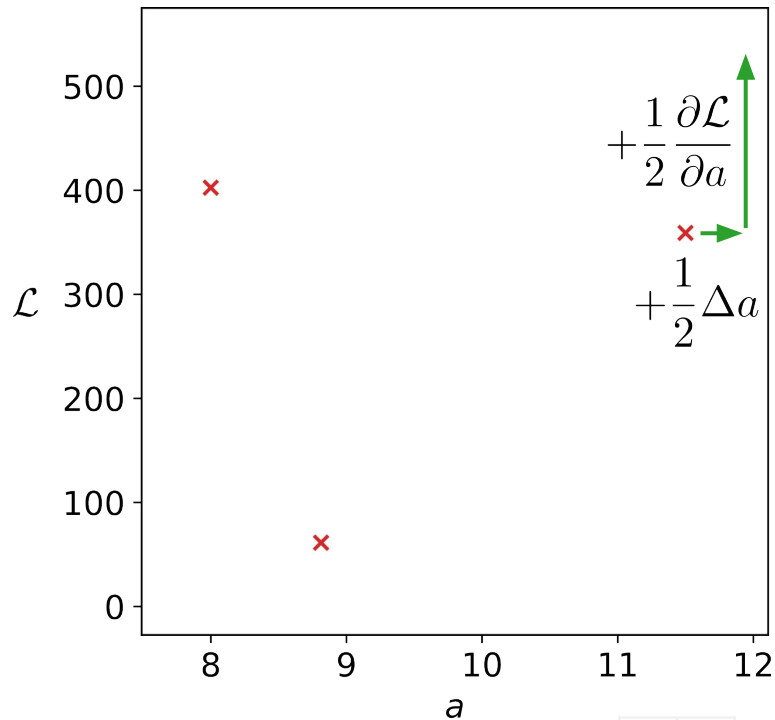
$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a) \quad \Delta a > 0$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

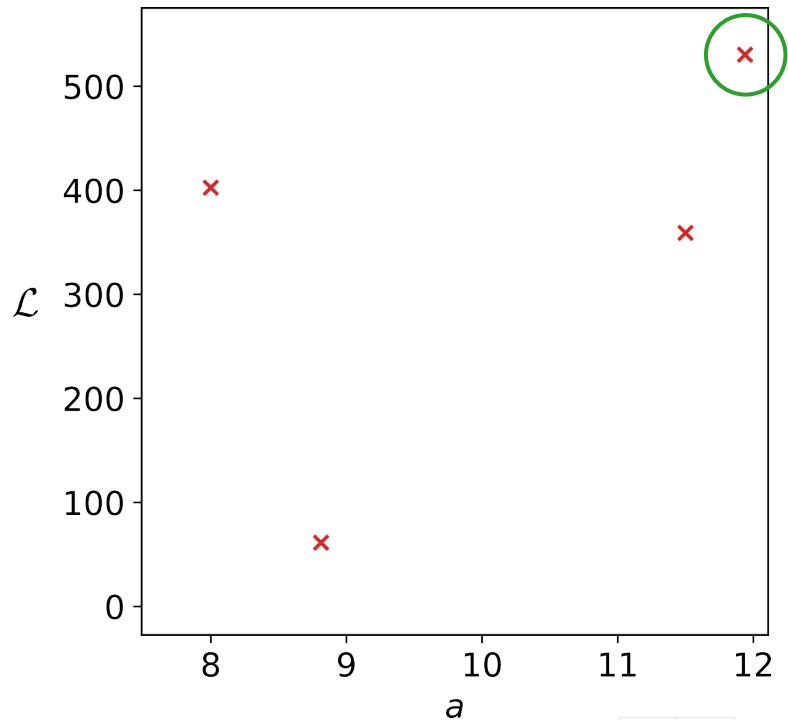
$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a) \quad \Delta a > 0$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

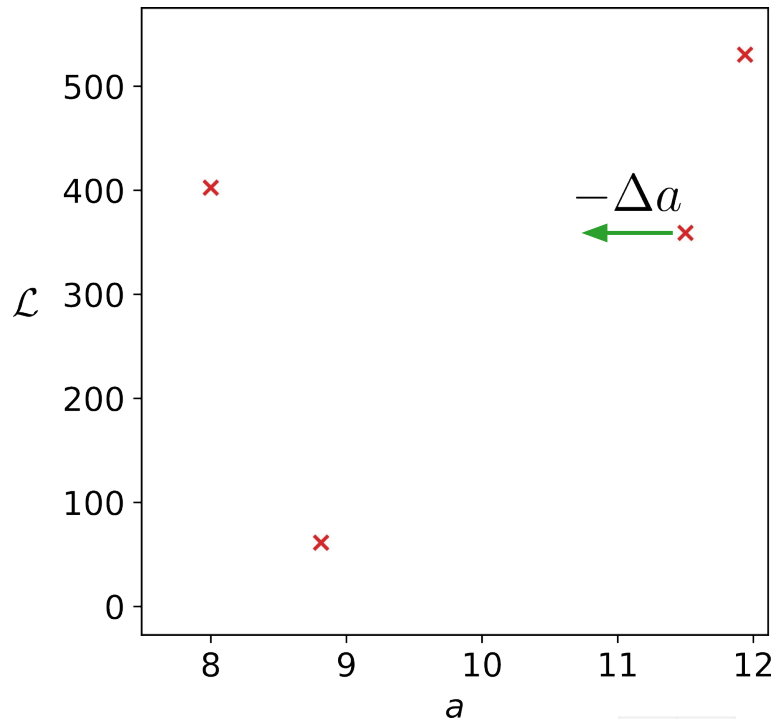
$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a) \quad \Delta a > 0$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

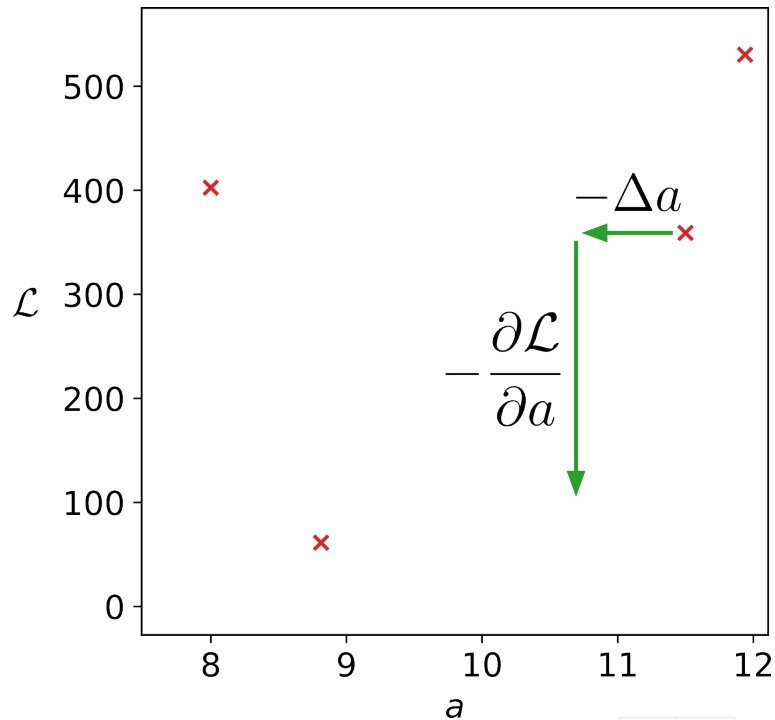
$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a) \quad \Delta a > 0$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

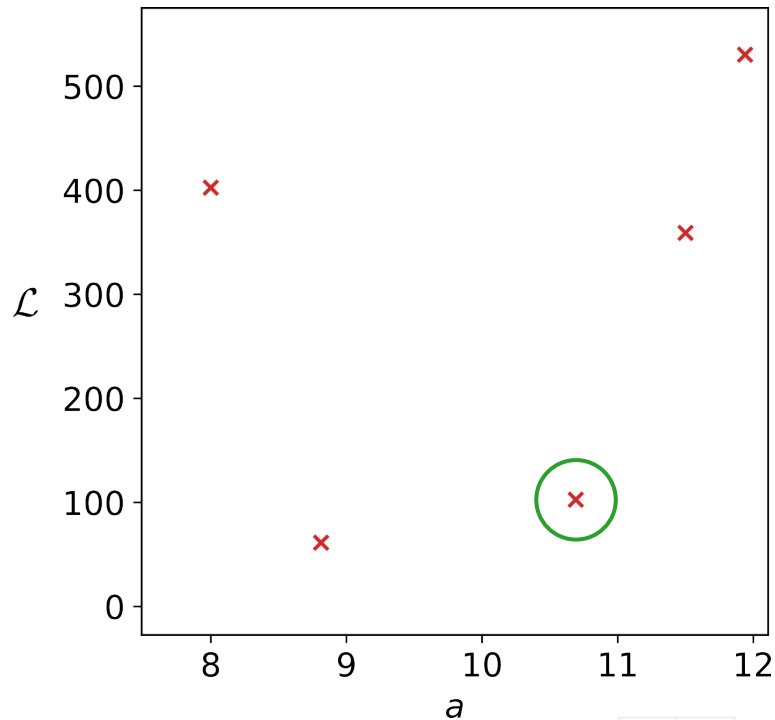
$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a) \quad \Delta a > 0$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

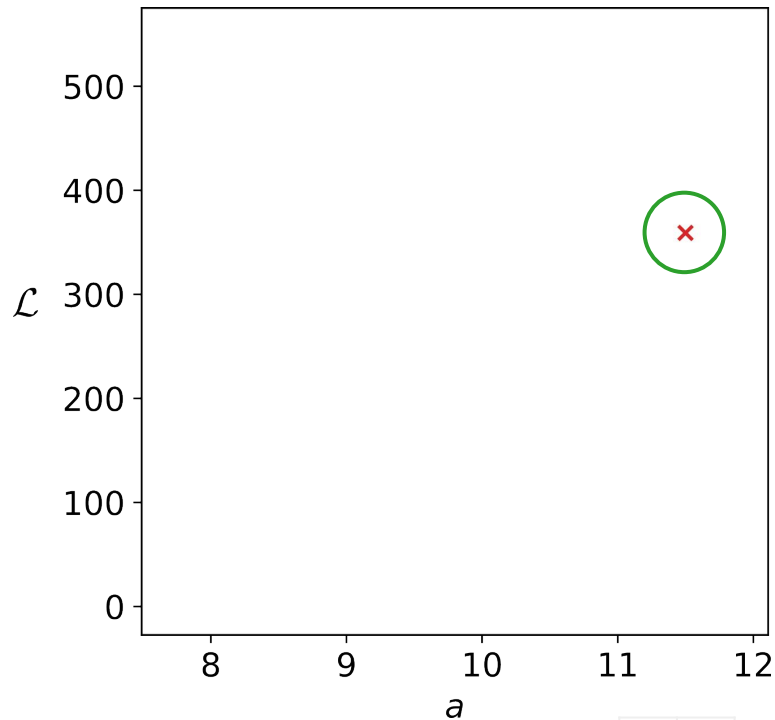
$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

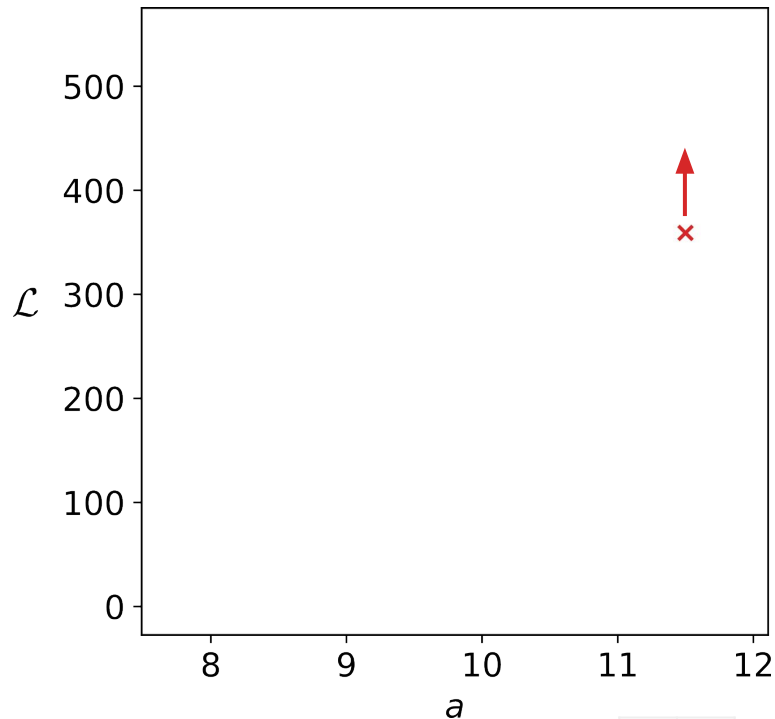
$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

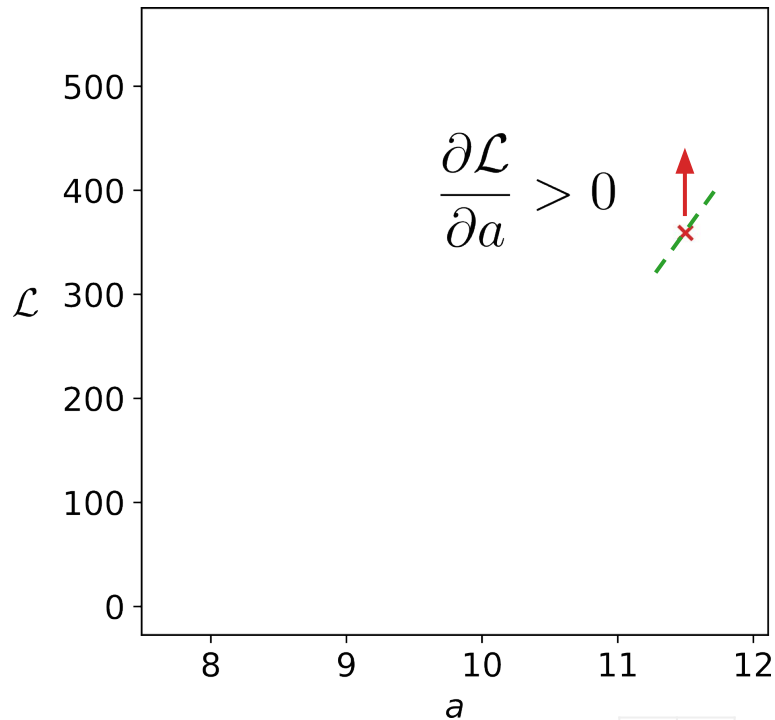
$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

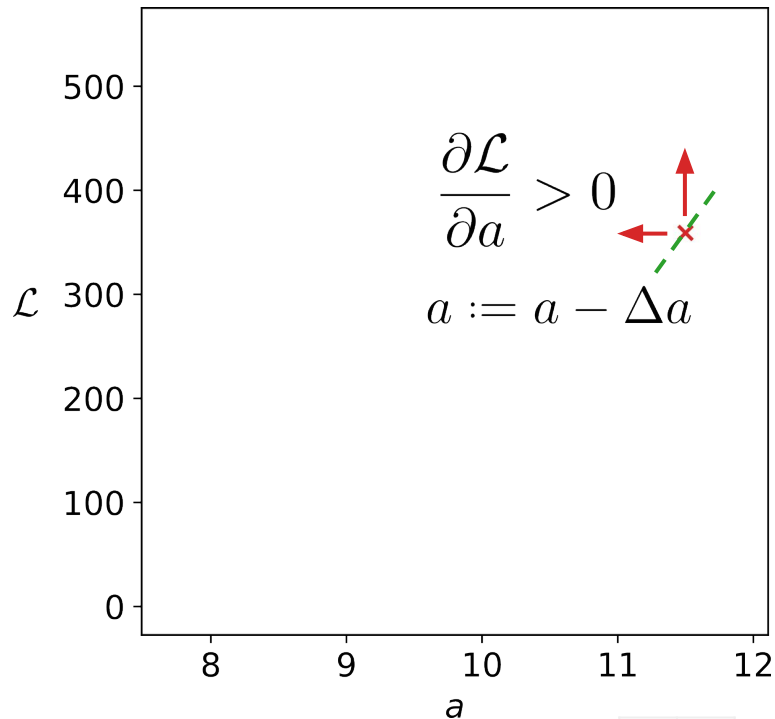
$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

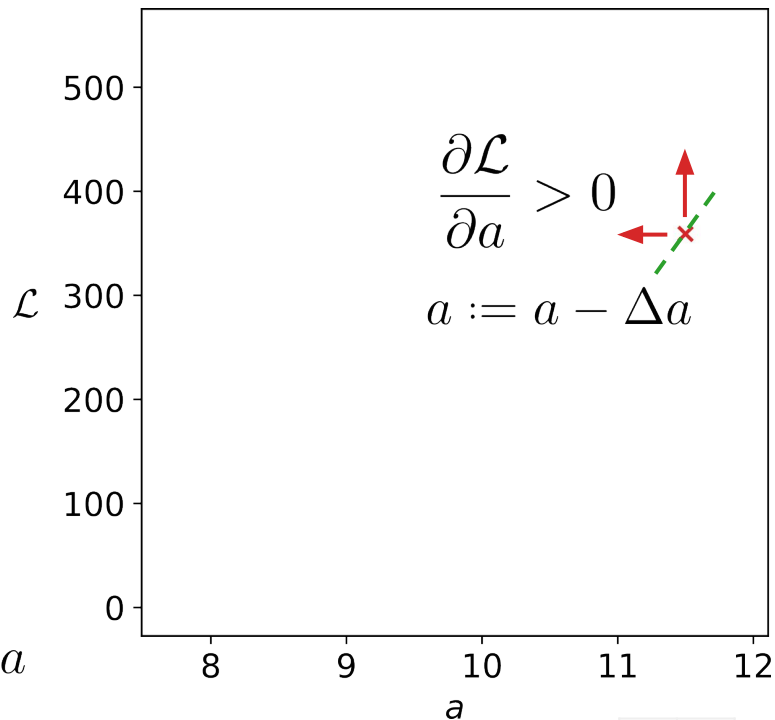
$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$

$$\frac{\partial \mathcal{L}}{\partial a} > 0$$

$$a := a - \Delta a$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

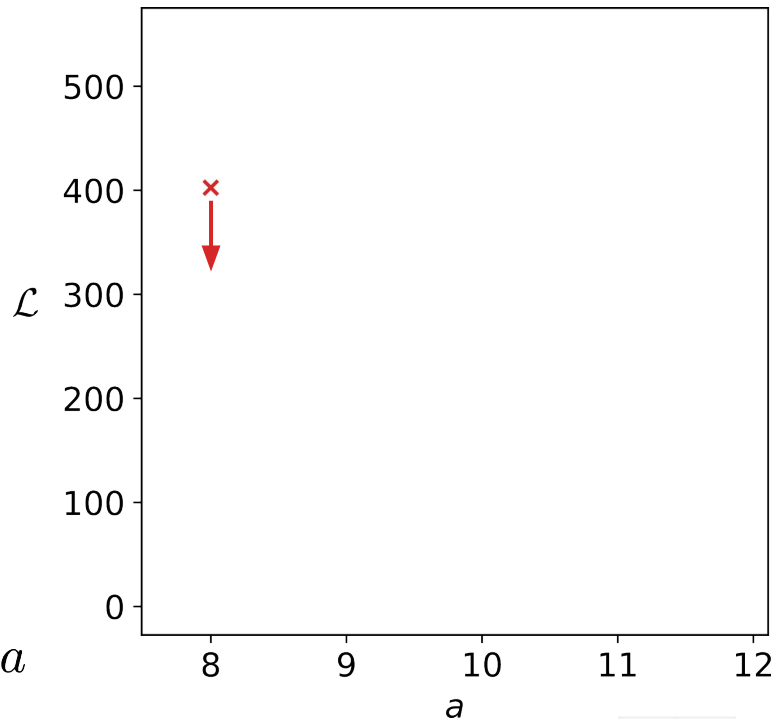
$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$

$$\frac{\partial \mathcal{L}}{\partial a} > 0$$

$$a := a - \Delta a$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

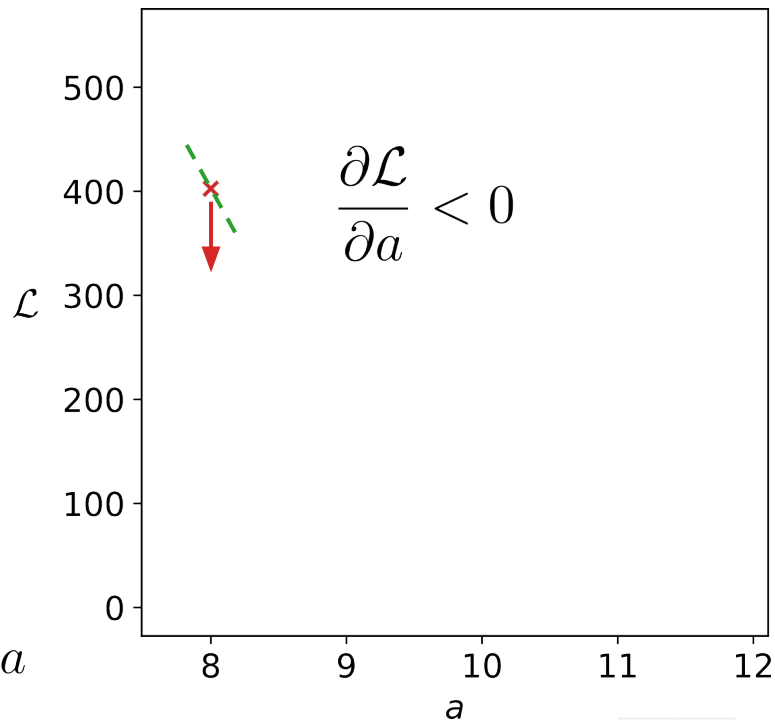
$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$

$$\frac{\partial \mathcal{L}}{\partial a} > 0$$

$$a := a - \Delta a$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

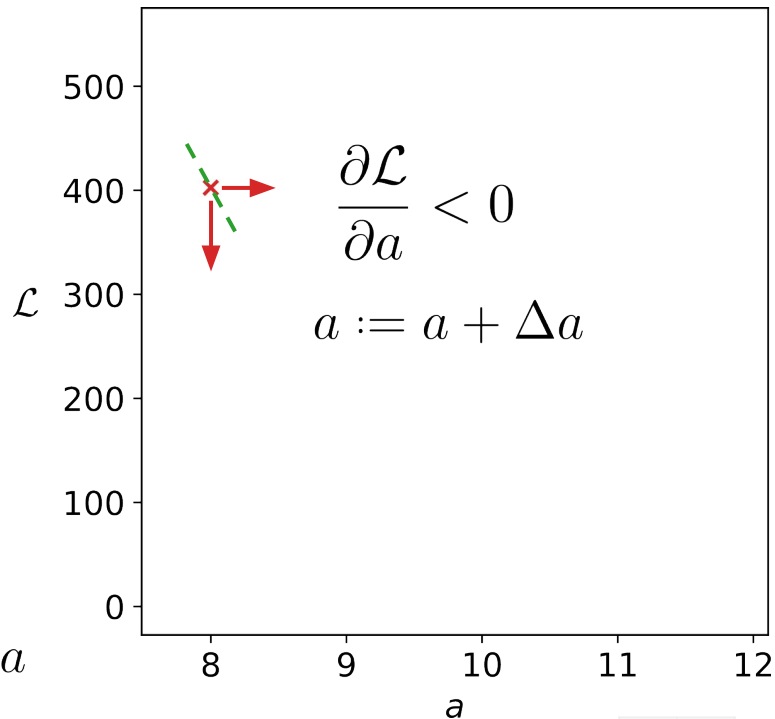
$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$

$$\frac{\partial \mathcal{L}}{\partial a} > 0$$

$$a := a - \Delta a$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

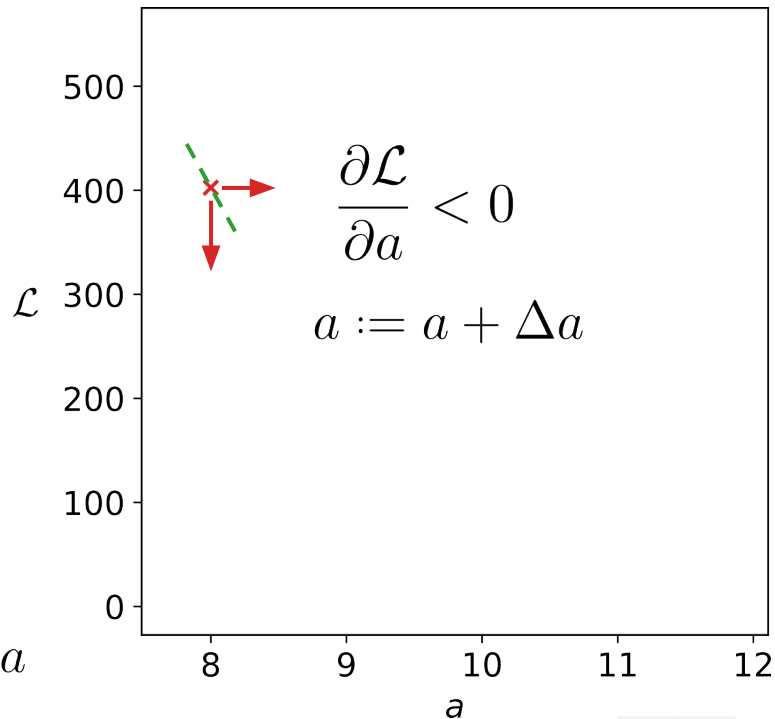
$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$

$$\frac{\partial \mathcal{L}}{\partial a} < 0$$

$$a := a + \Delta a$$

$$\frac{\partial \mathcal{L}}{\partial a} > 0$$

$$a := a - \Delta a$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

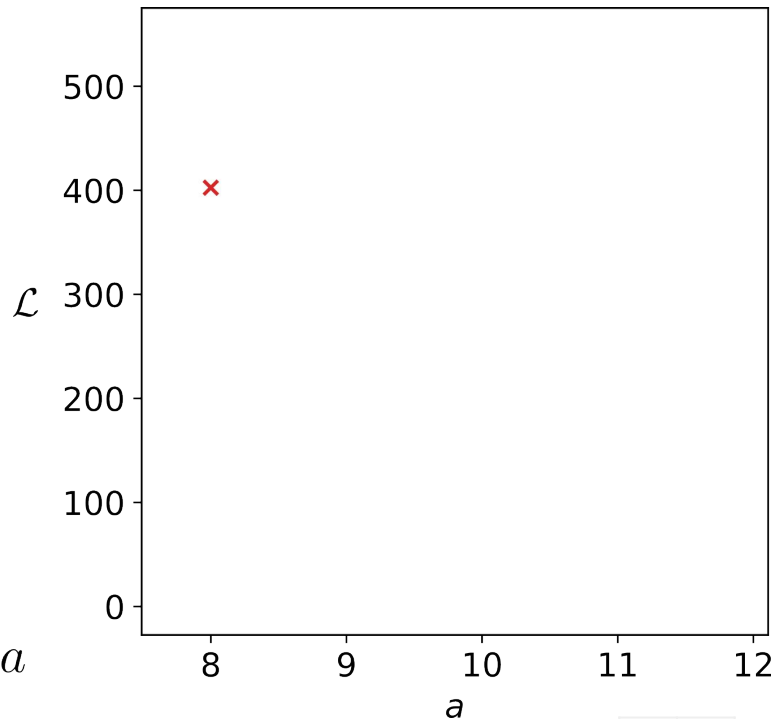
$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$

$$\frac{\partial \mathcal{L}}{\partial a} < 0$$

$$a := a + \Delta a$$

$$\frac{\partial \mathcal{L}}{\partial a} > 0$$

$$a := a - \Delta a$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

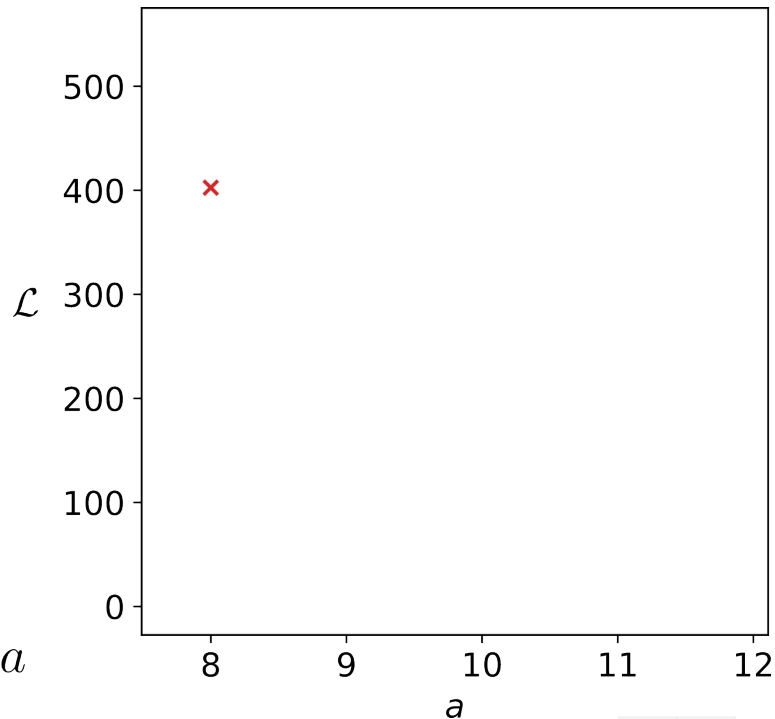
$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$

$$\frac{\partial \mathcal{L}}{\partial a} < 0 \quad \left| \quad \frac{\partial \mathcal{L}}{\partial a} > 0$$

$$a := a + \Delta a \quad \left| \quad a := a - \Delta a$$

$$a := a - \frac{\partial \mathcal{L}}{\partial a}$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$

$$\frac{\partial \mathcal{L}}{\partial a} < 0$$

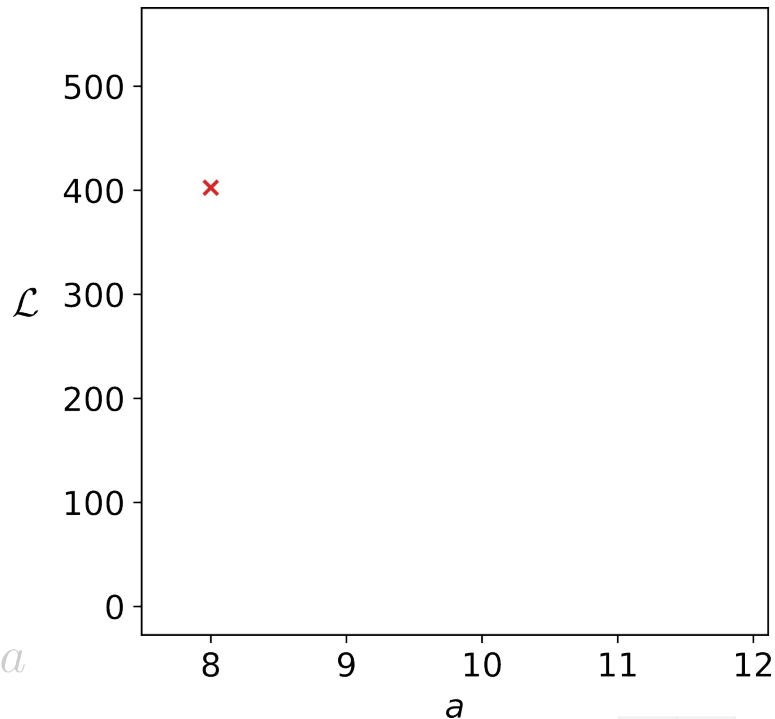
$$\frac{\partial \mathcal{L}}{\partial a} > 0$$

$$a := a + \Delta a$$

$$a := a - \Delta a$$

$$a := a - \gamma \frac{\partial \mathcal{L}}{\partial a}$$

γ chosen by us



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

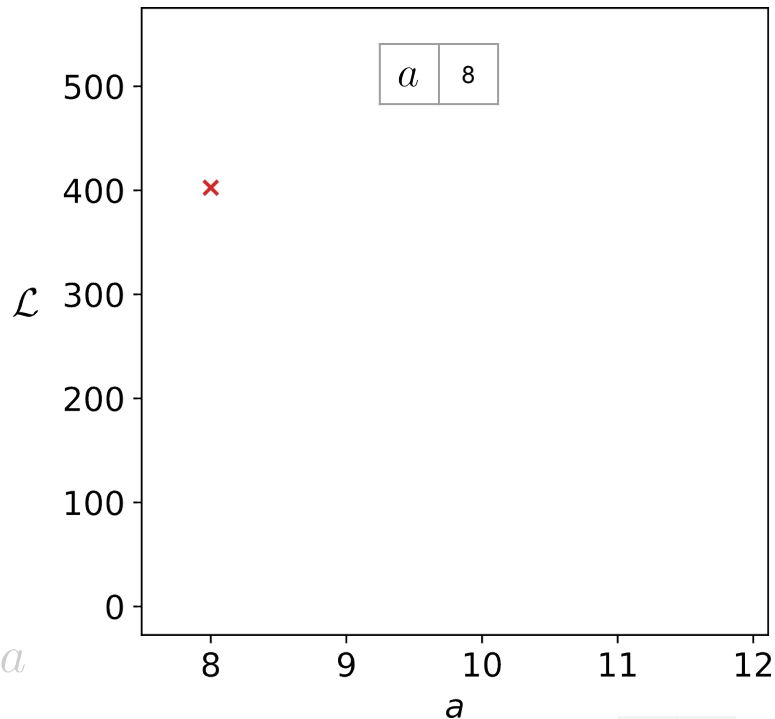
$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$

$$\frac{\partial \mathcal{L}}{\partial a} < 0 \quad \left| \quad \frac{\partial \mathcal{L}}{\partial a} > 0$$

$$a := a + \Delta a \quad \left| \quad a := a - \Delta a$$

$$a := a - \gamma \frac{\partial \mathcal{L}}{\partial a}$$

γ chosen by us



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show "influence" to functions

$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$

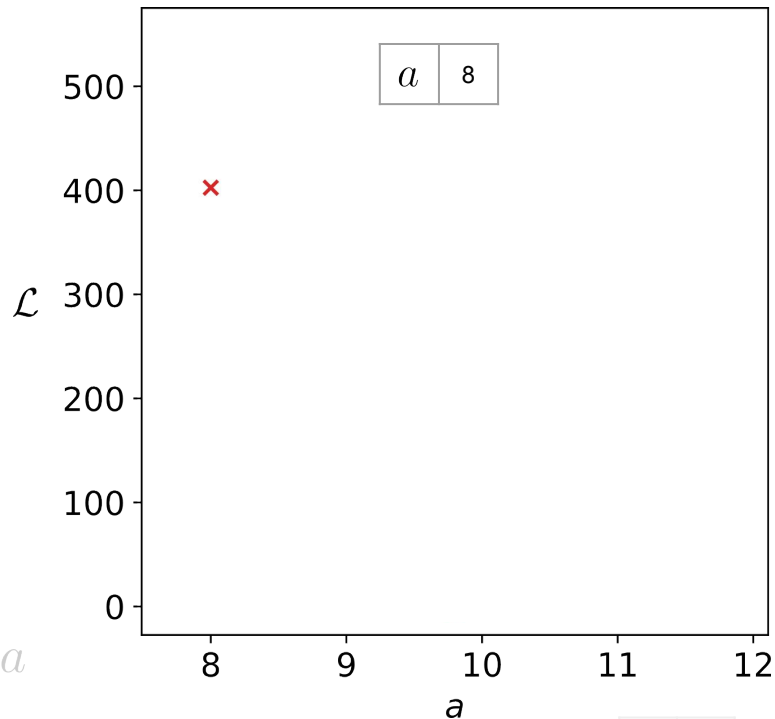
$$\frac{\partial \mathcal{L}}{\partial a} < 0 \quad \Bigg| \quad \frac{\partial \mathcal{L}}{\partial a} > 0$$

$$a := a + \Delta a \quad \Bigg| \quad a := a - \Delta a$$

$$a := a - \gamma \frac{\partial \mathcal{L}}{\partial a}$$

γ chosen by us

$$\gamma = 0.003$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$

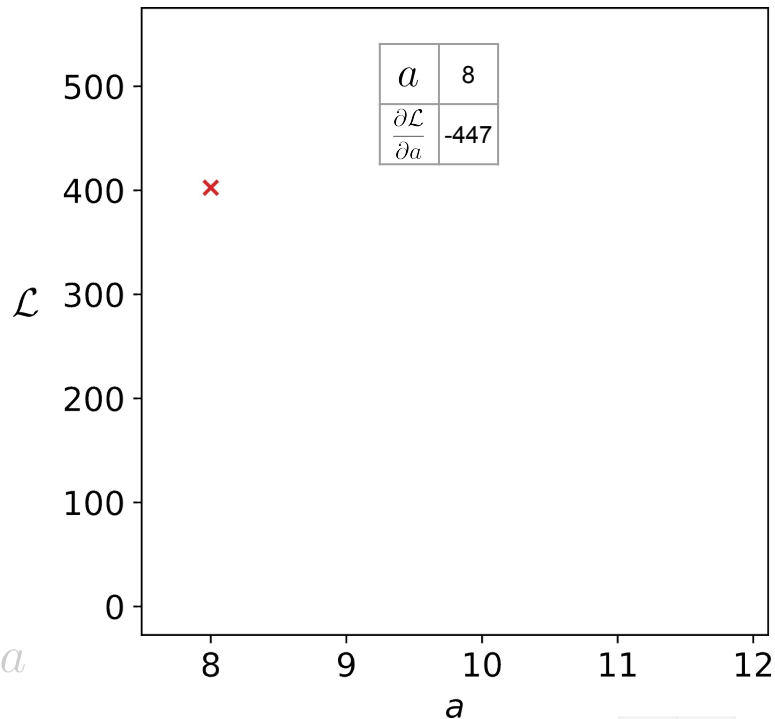
$$\frac{\partial \mathcal{L}}{\partial a} < 0 \quad \Bigg| \quad \frac{\partial \mathcal{L}}{\partial a} > 0$$

$$a := a + \Delta a \quad \Bigg| \quad a := a - \Delta a$$

$$a := a - \gamma \frac{\partial \mathcal{L}}{\partial a}$$

γ chosen by us

$$\gamma = 0.003$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show "influence" to functions

$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$

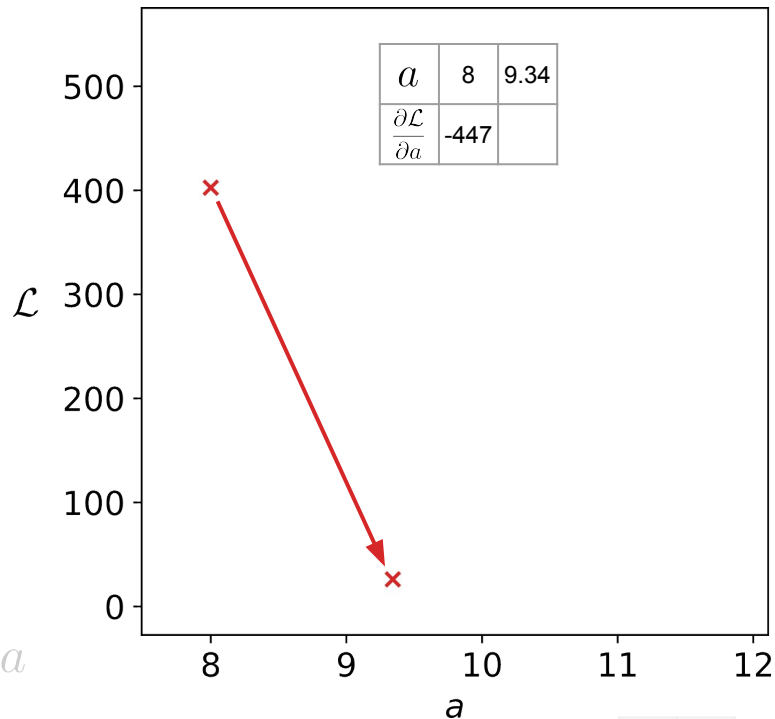
$$\frac{\partial \mathcal{L}}{\partial a} < 0 \quad \Bigg| \quad \frac{\partial \mathcal{L}}{\partial a} > 0$$

$$a := a + \Delta a \quad \Bigg| \quad a := a - \Delta a$$

$$a := a - \gamma \frac{\partial \mathcal{L}}{\partial a}$$

γ chosen by us

$$\gamma = 0.003$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show "influence" to functions

$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$

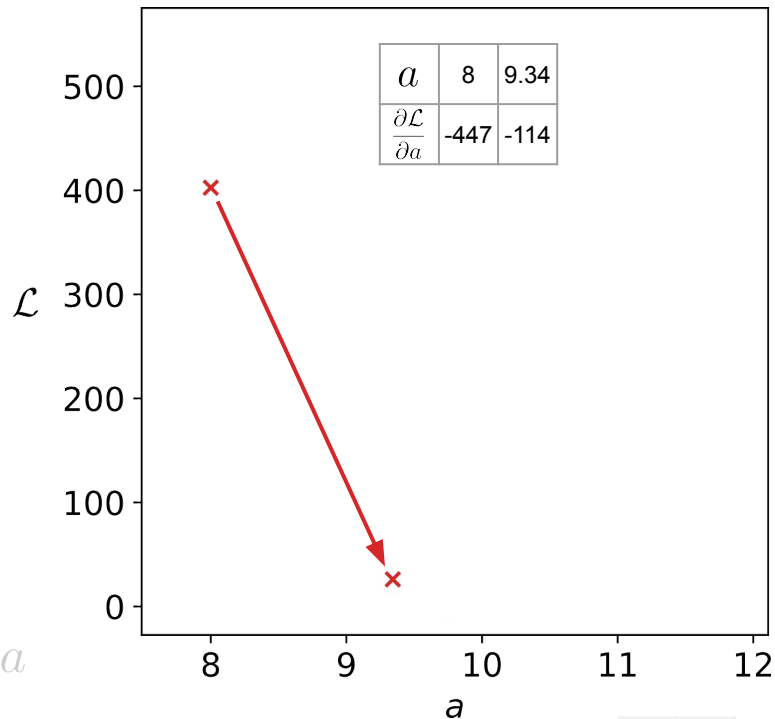
$$\frac{\partial \mathcal{L}}{\partial a} < 0 \quad \Bigg| \quad \frac{\partial \mathcal{L}}{\partial a} > 0$$

$$a := a + \Delta a \quad \Bigg| \quad a := a - \Delta a$$

$$a := a - \gamma \frac{\partial \mathcal{L}}{\partial a}$$

γ chosen by us

$$\gamma = 0.003$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show "influence" to functions

$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$

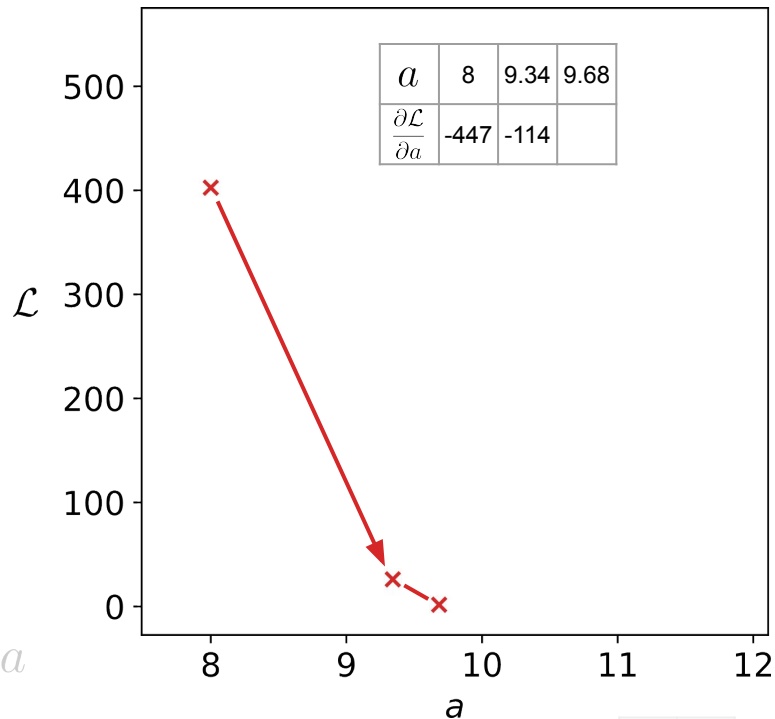
$$\frac{\partial \mathcal{L}}{\partial a} < 0 \quad \left| \quad \frac{\partial \mathcal{L}}{\partial a} > 0$$

$$a := a + \Delta a \quad \left| \quad a := a - \Delta a$$

$$a := a - \gamma \frac{\partial \mathcal{L}}{\partial a}$$

γ chosen by us

$$\gamma = 0.003$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show "influence" to functions

$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$

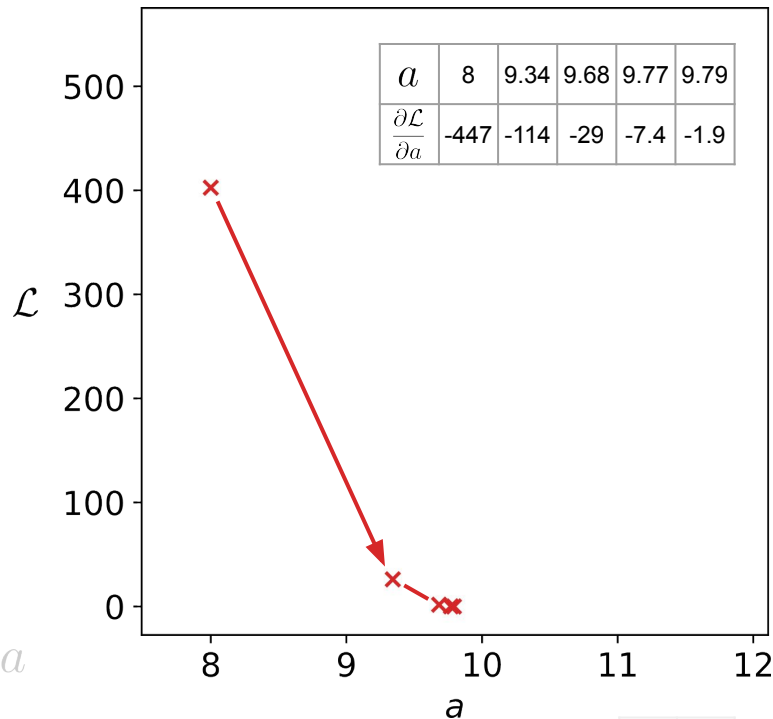
$$\frac{\partial \mathcal{L}}{\partial a} < 0 \quad \left| \quad \frac{\partial \mathcal{L}}{\partial a} > 0$$

$$a := a + \Delta a \quad \left| \quad a := a - \Delta a$$

$$a := a - \gamma \frac{\partial \mathcal{L}}{\partial a}$$

γ chosen by us

$$\gamma = 0.003$$



$$\hat{y} = ax$$

$$\mathcal{L} = \sum_i (\hat{y}_i - y_i)^2$$

x	y
2	19.6
4.5	44.1
10	98

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show "influence" to functions

$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$

$$\frac{\partial \mathcal{L}}{\partial a} < 0$$

$$\frac{\partial \mathcal{L}}{\partial a} > 0$$

$$a := a + \Delta a$$

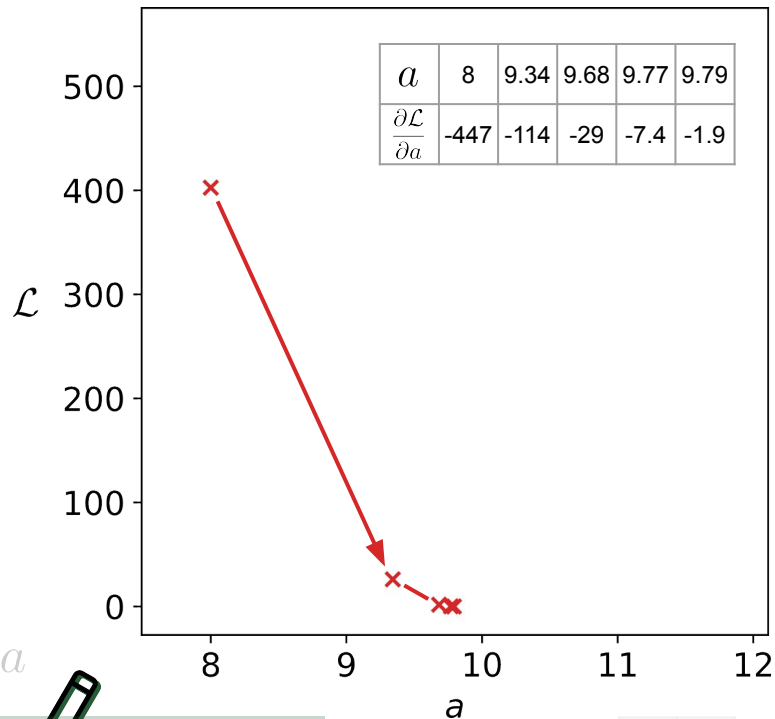
$$a := a - \Delta a$$

repeat:

$$a := a - \gamma \frac{\partial \mathcal{L}}{\partial a}$$

until minimum is reached

$$\gamma = 0.003$$



x	y
2	19.6
4.5	44.1
10	98

$$(\hat{y}_i - y_i)^2$$

Gradient Descent

How to find the minimum? $\min_a \mathcal{L}$

Derivatives show “influence” to functions

$$\frac{\partial \mathcal{L}}{\partial a} \quad \mathcal{L} = \mathcal{L}(a)$$

$$\mathcal{L}(8) \approx 402$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx -447$$

$$\mathcal{L}(11.5) \approx 359$$

$$\frac{\partial \mathcal{L}}{\partial a}(8) \approx 422$$

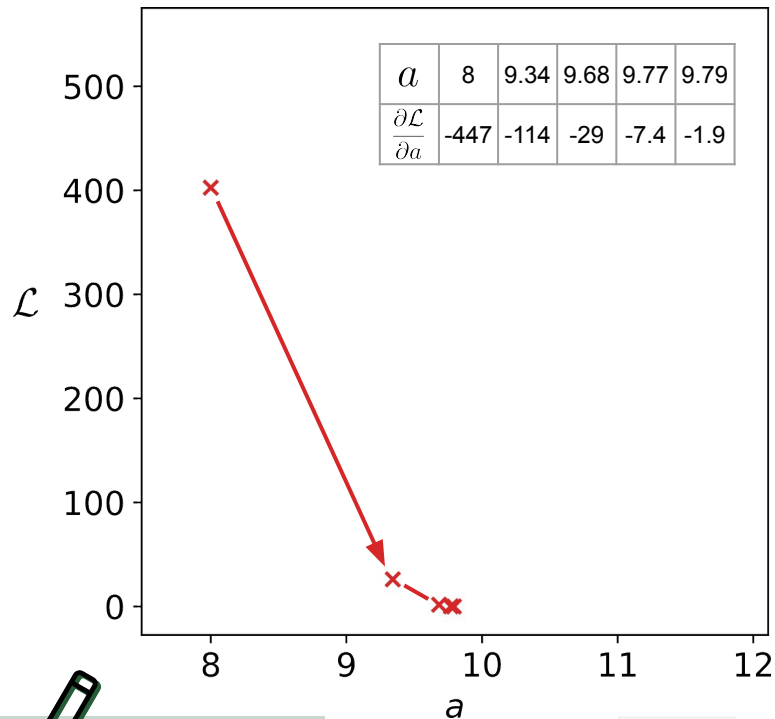
γ – learning rate

repeat:

$$a := a - \gamma \frac{\partial \mathcal{L}}{\partial a}$$

until minimum is reached

$\gamma = 0.003$

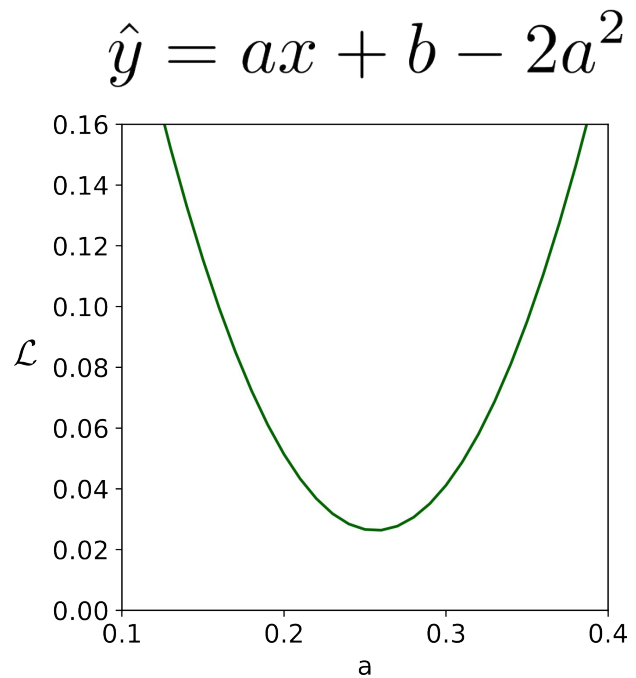


x	y
2	19.6
4.5	44.1
10	98

$$(\hat{y}_i - y_i)^2$$

Gradient Descent

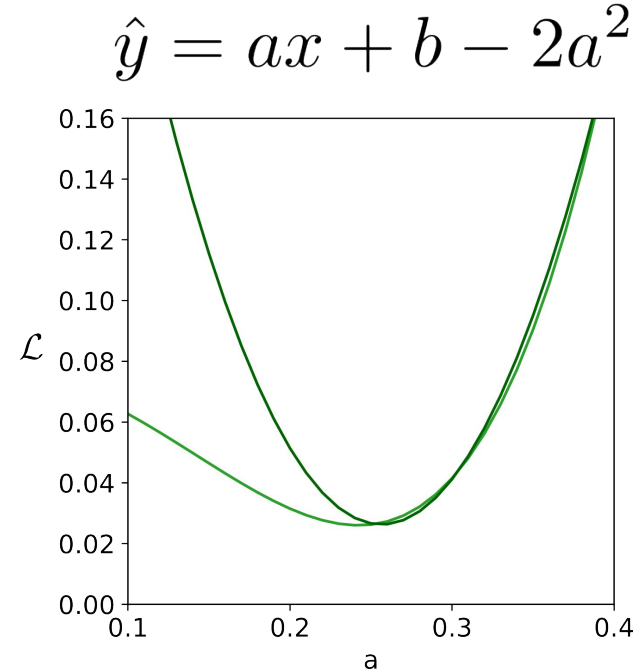
x	y
-0.3	1.9
-0.2	1.85
0.2	1.85
0.3	1.85



$b = 2$

Gradient Descent

x	y
-0.3	1.9
-0.2	1.85
0.2	1.85
0.3	1.85

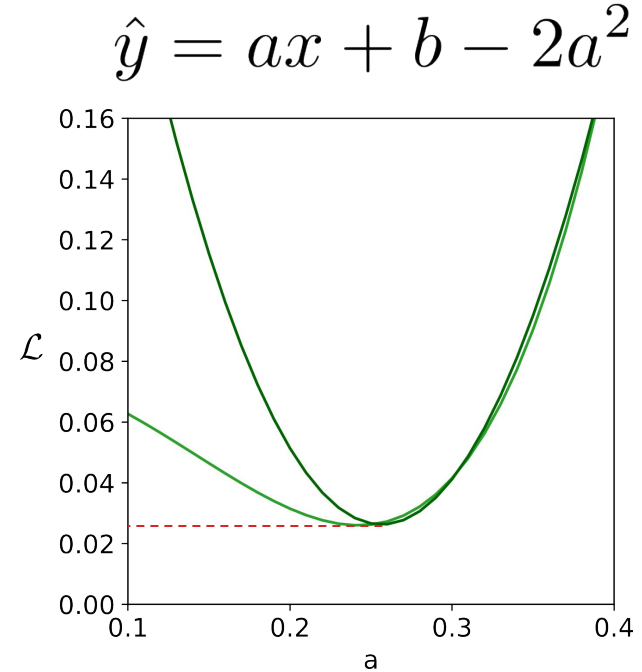


$$\hat{y} = -5.5(a - 0.246)x + b - 0.217$$

$$b = 2$$

Gradient Descent

x	y
-0.3	1.9
-0.2	1.85
0.2	1.85
0.3	1.85



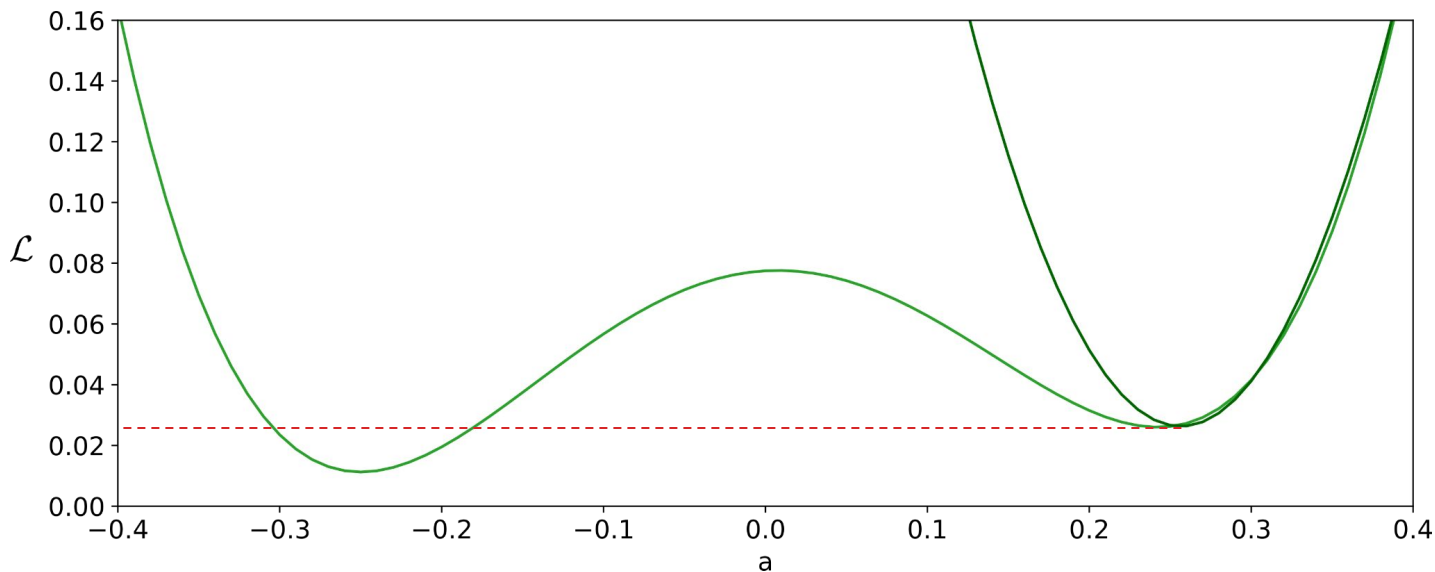
$$\hat{y} = -5.5(a - 0.246)x + b - 0.217$$

$$b = 2$$

Gradient Descent

x	y
-0.3	1.9
-0.2	1.85
0.2	1.85
0.3	1.85

$$\hat{y} = ax + b - 2a^2$$



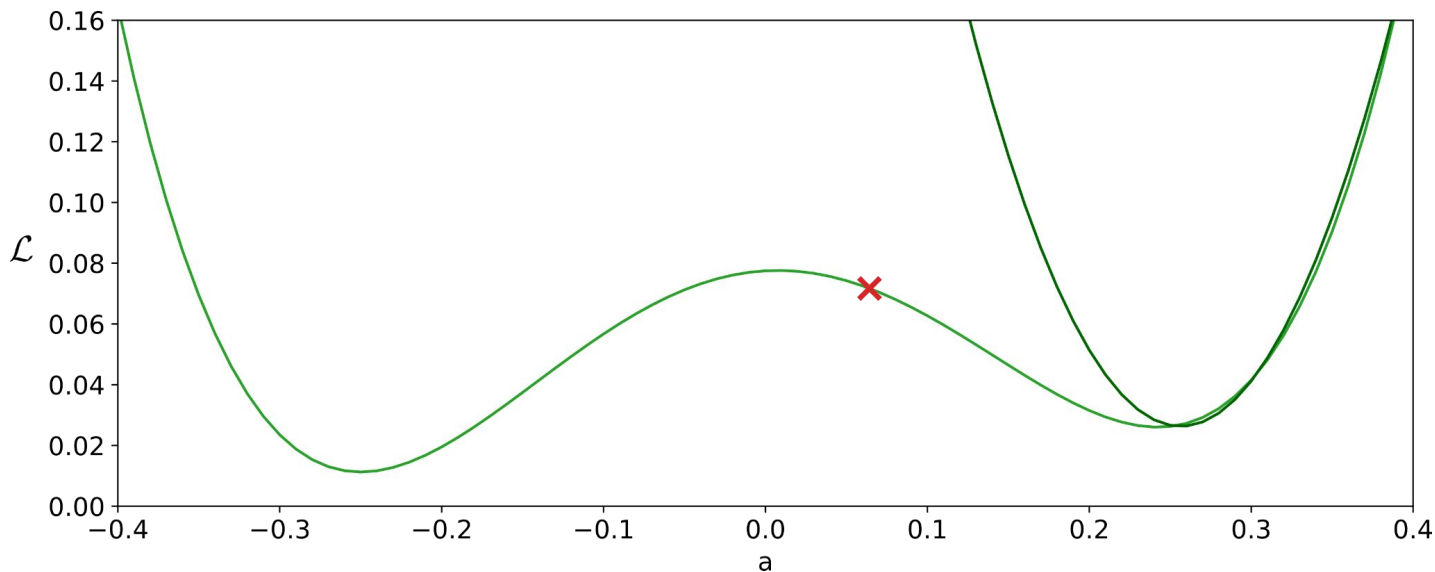
$$\hat{y} = -5.5(a - 0.246)x + b - 0.217$$

$$b = 2$$

Gradient Descent

x	y
-0.3	1.9
-0.2	1.85
0.2	1.85
0.3	1.85

$$\hat{y} = ax + b - 2a^2$$



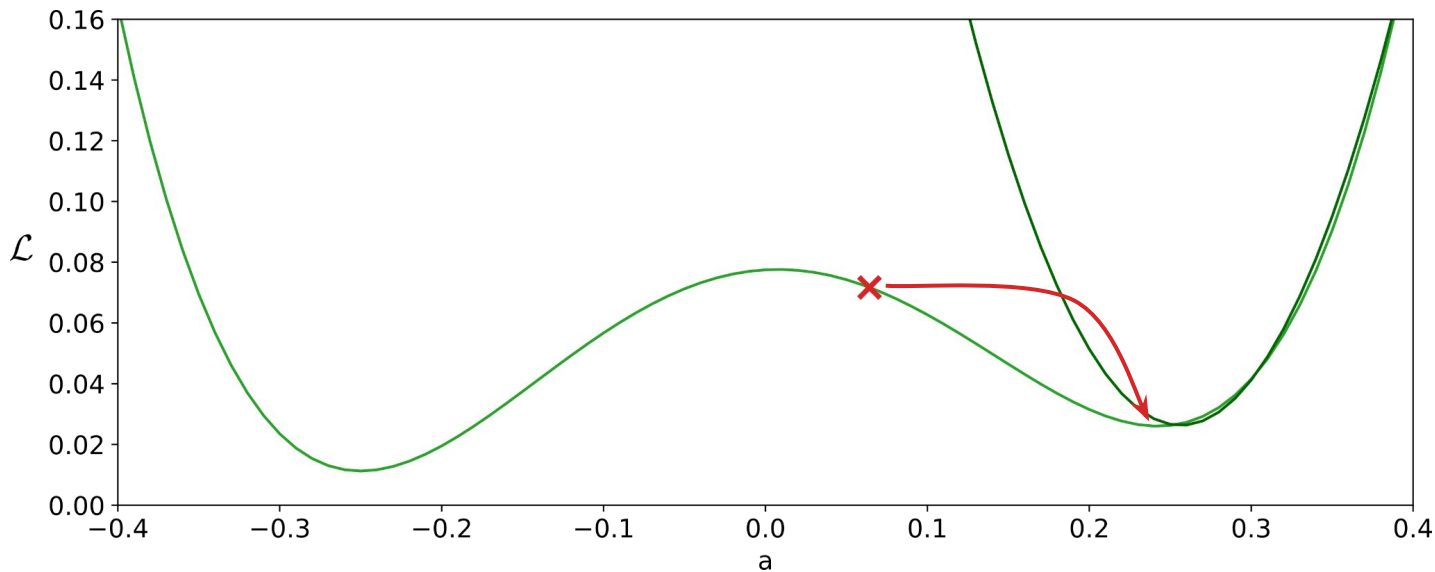
$$\hat{y} = -5.5(a - 0.246)x + b - 0.217$$

$b = 2$

Gradient Descent

x	y
-0.3	1.9
-0.2	1.85
0.2	1.85
0.3	1.85

$$\hat{y} = ax + b - 2a^2$$



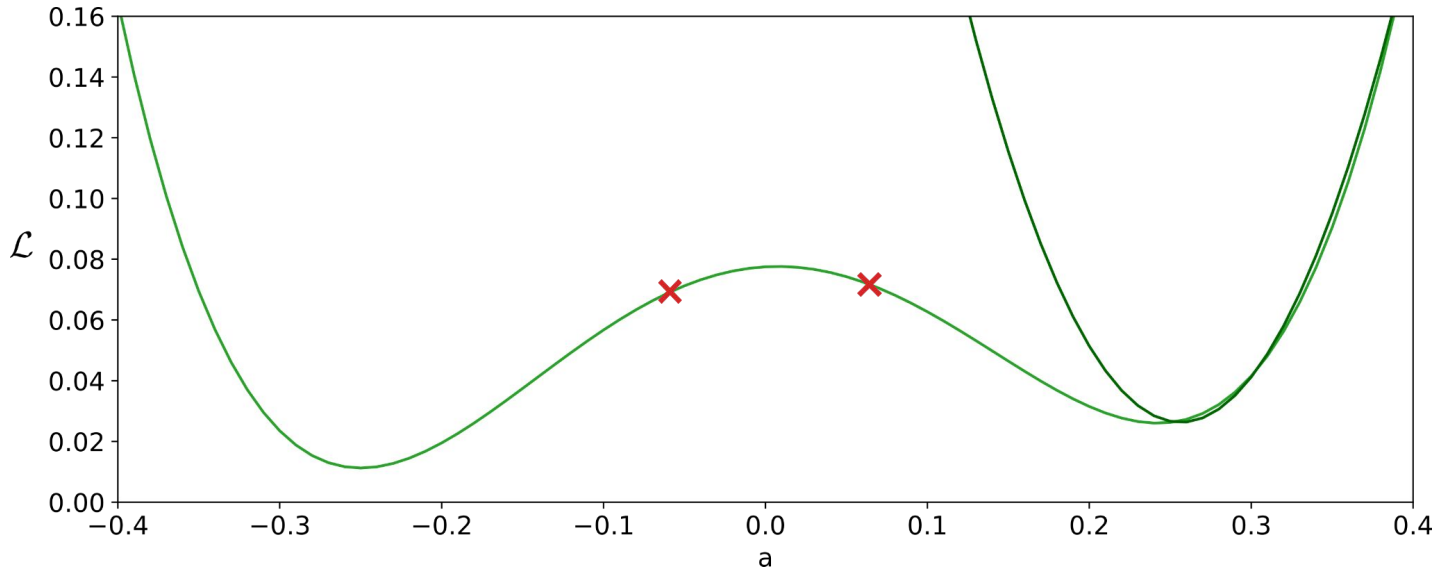
$$\hat{y} = -5.5(a - 0.246)x + b - 0.217$$

$b = 2$

Gradient Descent

x	y
-0.3	1.9
-0.2	1.85
0.2	1.85
0.3	1.85

$$\hat{y} = ax + b - 2a^2$$

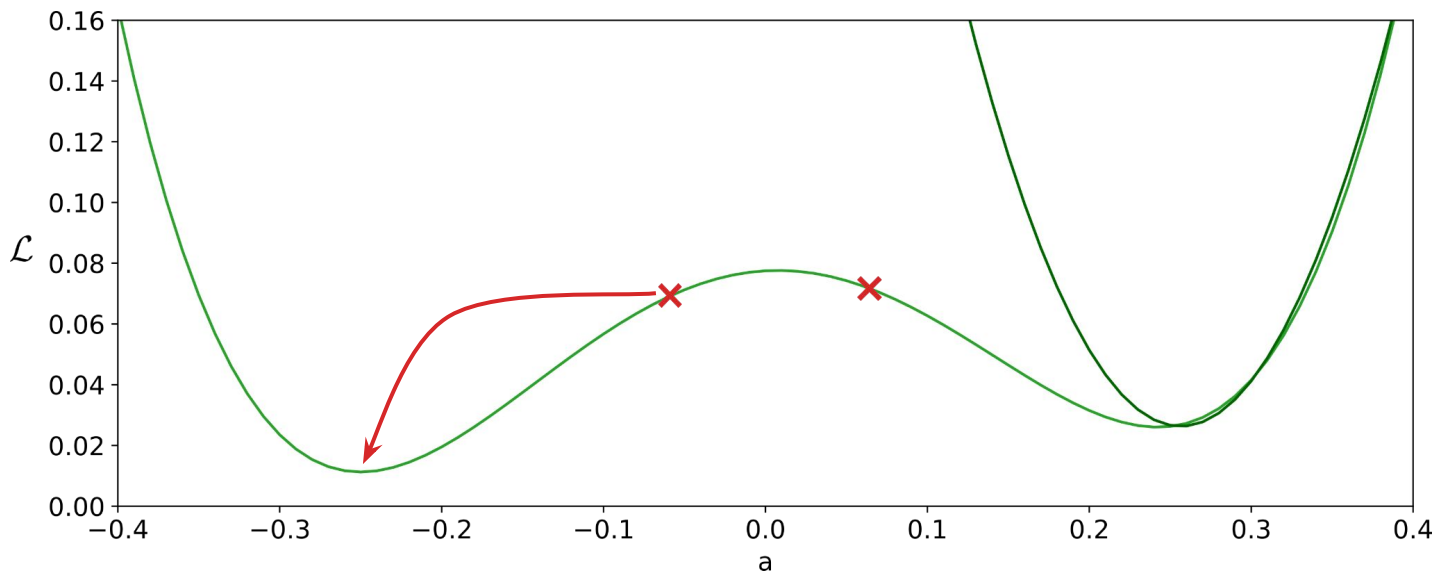


$$\hat{y} = -5.5(a - 0.246)x + b - 0.217$$

Gradient Descent

x	y
-0.3	1.9
-0.2	1.85
0.2	1.85
0.3	1.85

$$\hat{y} = ax + b - 2a^2$$



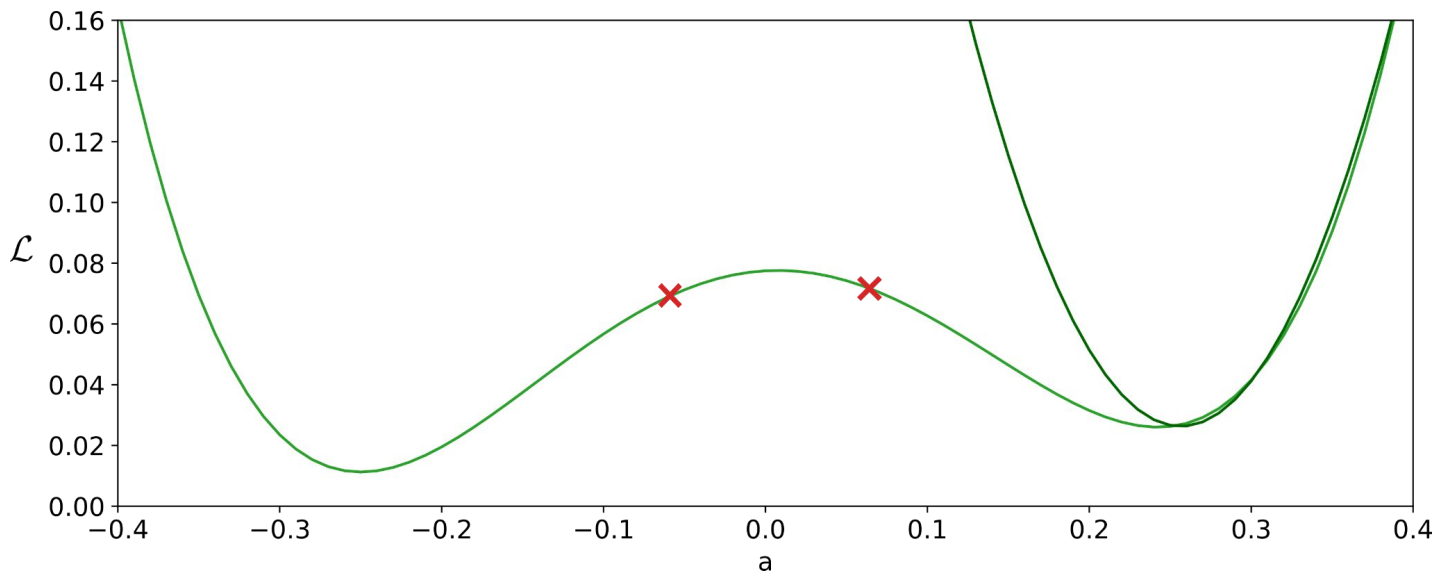
$$\hat{y} = -5.5(a - 0.246)x + b - 0.217$$

$b = 2$

Gradient Descent

x	y
-0.3	1.9
-0.2	1.85
0.2	1.85
0.3	1.85

$$\hat{y} = ax + b - 2a^2$$

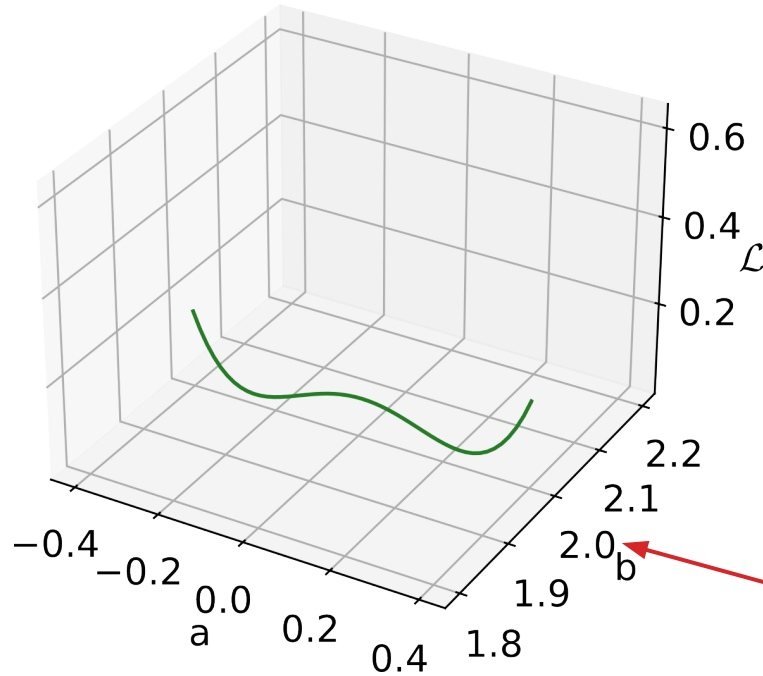


$$\hat{y} = -5.5(a - 0.246)x + b - 0.217$$

$$b = 2$$

Gradient Descent

x	y
-0.3	1.9
-0.2	1.85
0.2	1.85
0.3	1.85

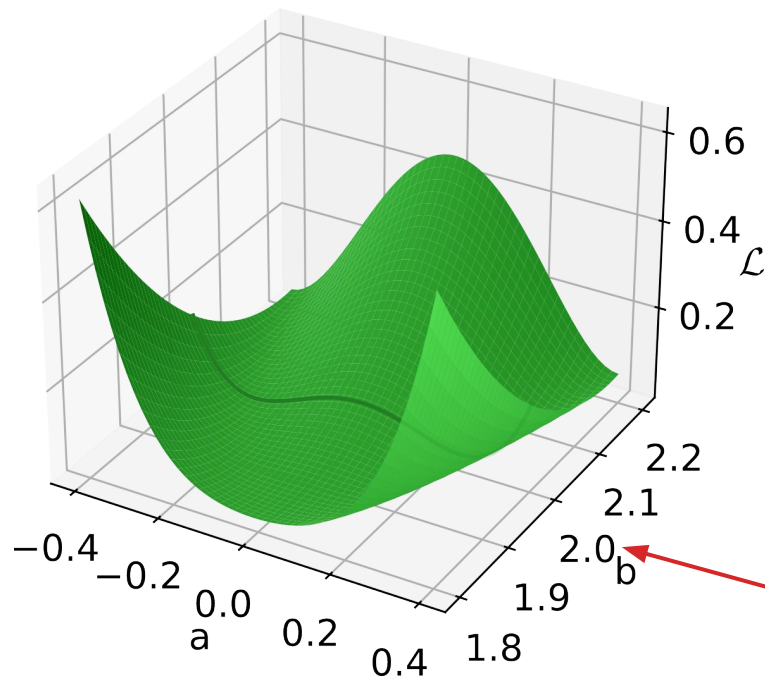


$$\hat{y} = -5.5(a - 0.246)x + b - 0.217$$

$b = 2$

Gradient Descent

x	y
-0.3	1.9
-0.2	1.85
0.2	1.85
0.3	1.85

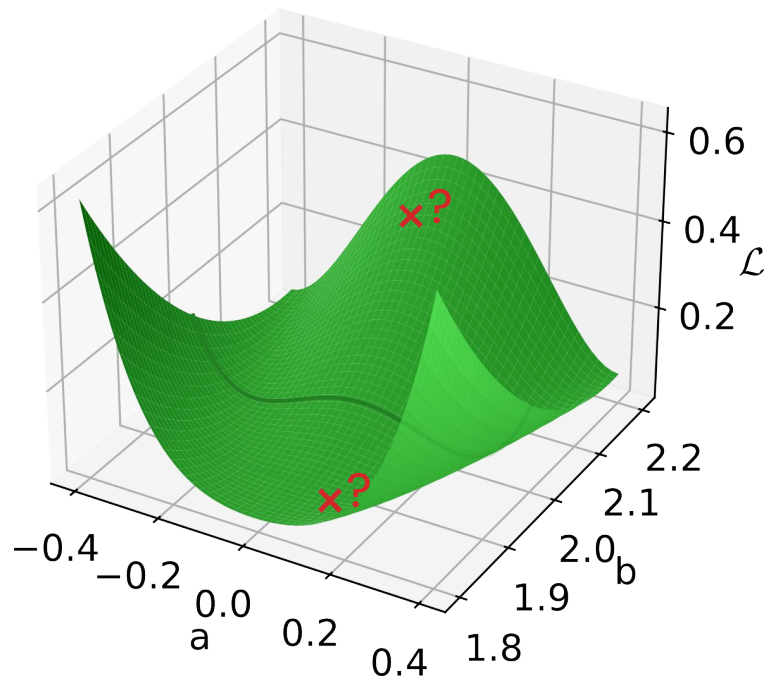


$$\hat{y} = -5.5(a - 0.246)x + b - 0.217$$

$b = 2$

Gradient Descent

x	y
-0.3	1.9
-0.2	1.85
0.2	1.85
0.3	1.85

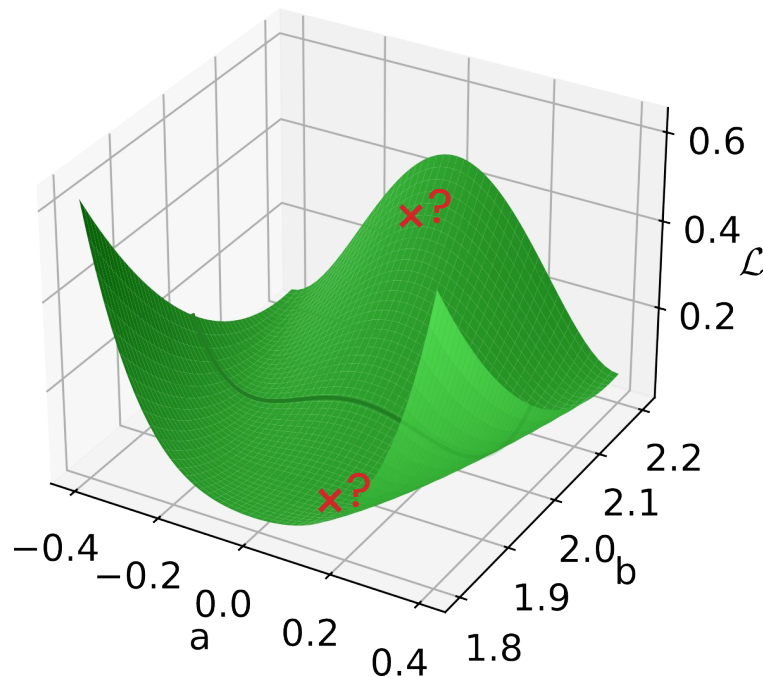


$$\hat{y} = -5.5(a - 0.246)x + b - 0.217$$

$$b = 2$$

Gradient Descent

x	y
-0.3	1.9
-0.2	1.85
0.2	1.85
0.3	1.85



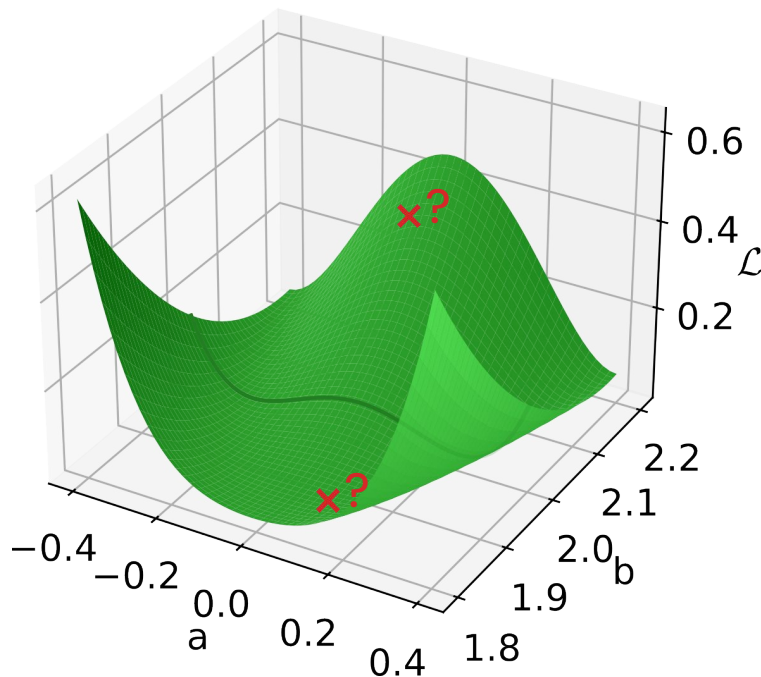
$$\theta = \{a, b\}$$

$$\hat{y} = -5.5(a - 0.246)x + b - 0.217$$

$$b = 2$$

Gradient Descent

x	y
-0.3	1.9
-0.2	1.85
0.2	1.85
0.3	1.85



$$\theta = \{a, b\}$$

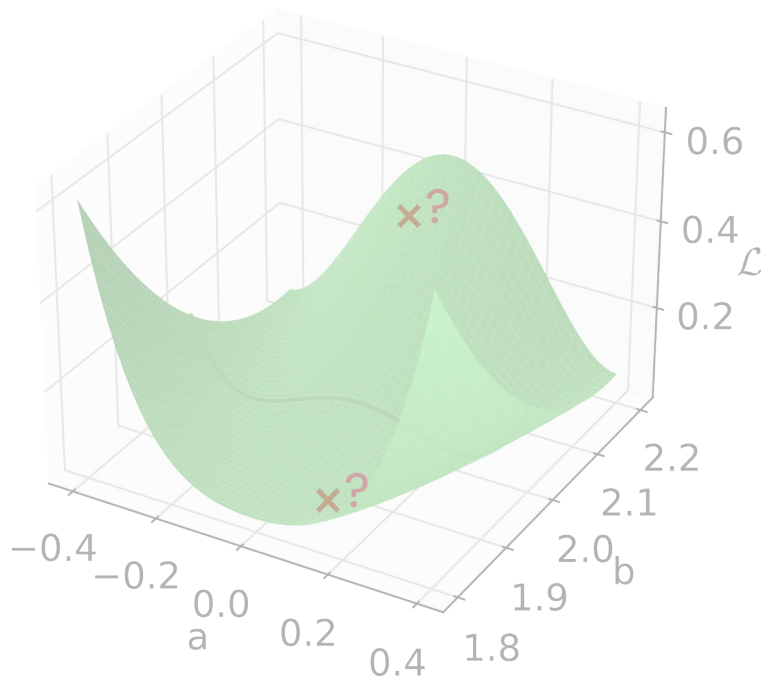
$$\theta = \{a, b, c, \dots\}$$

$$\hat{y} = -5.5(a - 0.246)x + b - 0.217$$

$$b = 2$$

Gradient Descent

x	y
-0.3	1.9
-0.2	1.85
0.2	1.85
0.3	1.85



$$\theta = \{a, b\}$$

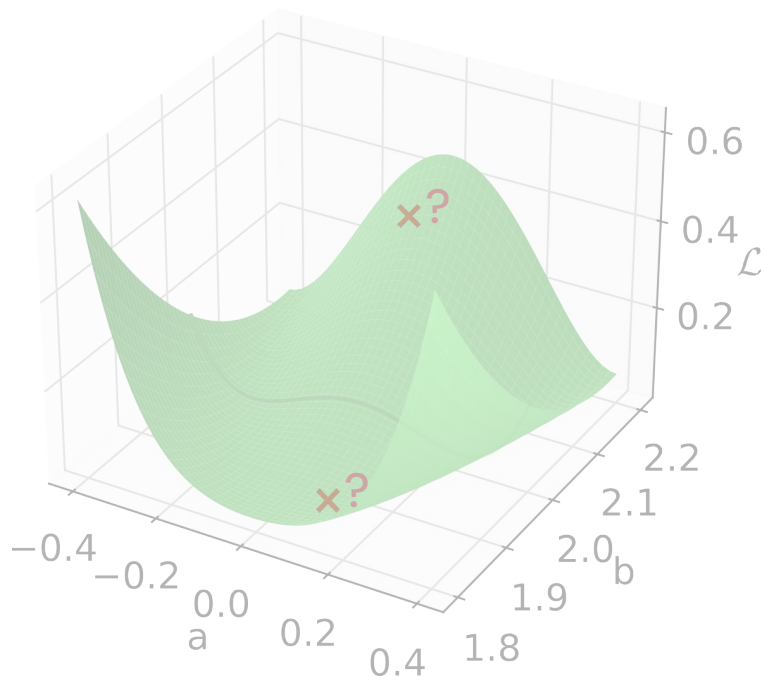
$$\theta = \{a, b, c, \dots\}$$

$$\hat{y} = -5.5(a - 0.246)x + b - 0.217$$

$$b = 2$$

Gradient Descent

x	y
-0.3	1.9
-0.2	1.85
0.2	1.85
0.3	1.85
...	...
...	...
...	...



$$\theta = \{a, b\}$$

$$\theta = \{a, b, c, \dots\}$$

$$\hat{y} = -5.5(a - 0.246)x + b - 0.217$$

$$b = 2$$

Loading the Data

size of dataset \gg *computer memory*

Loading the Data

size of dataset \gg *computer memory*

→ *batching*

Stochastic Gradient Descent

$$\mathcal{L}_\theta = \frac{1}{N} \sum_i^N \ell_\theta(x_i, y_i) \quad \text{e.g., } \ell_\theta(x_i, y_i) = (\hat{y}(x_i) - y_i)^2$$

Stochastic Gradient Descent

$$\mathcal{L}_\theta = \frac{1}{N} \sum_i^N \ell_\theta(x_i, y_i) \quad \text{e.g., } \ell_\theta(x_i, y_i) = (\hat{y}(x_i) - y_i)^2$$

Gradient Descent $\theta_{t+1} = \theta_t - \gamma \frac{\partial \mathcal{L}}{\partial \theta_t}$

Stochastic Gradient Descent

$$\mathcal{L}_\theta = \frac{1}{N} \sum_i^N \ell_\theta(x_i, y_i) \quad \text{e.g., } \ell_\theta(x_i, y_i) = (\hat{y}(x_i) - y_i)^2$$

(Full) Batch Gradient Descent $\theta_{t+1} = \theta_t - \gamma \frac{\partial \mathcal{L}}{\partial \theta_t}$

Stochastic Gradient Descent

$$\mathcal{L}_\theta = \frac{1}{N} \sum_i^N \ell_\theta(x_i, y_i) \quad \text{e.g., } \ell_\theta(x_i, y_i) = (\hat{y}(x_i) - y_i)^2$$

(Full) Batch Gradient Descent $\theta_{t+1} = \theta_t - \gamma \frac{\partial \mathcal{L}}{\partial \theta_t}$

Stochastic Gradient Descent $\theta_{t+1} = \theta_t - \gamma \frac{\partial \ell}{\partial \theta_t}$

Stochastic Gradient Descent

$$\mathcal{L}_\theta = \frac{1}{N} \sum_i^N \ell_\theta(x_i, y_i) \quad \text{e.g., } \ell_\theta(x_i, y_i) = (\hat{y}(x_i) - y_i)^2$$

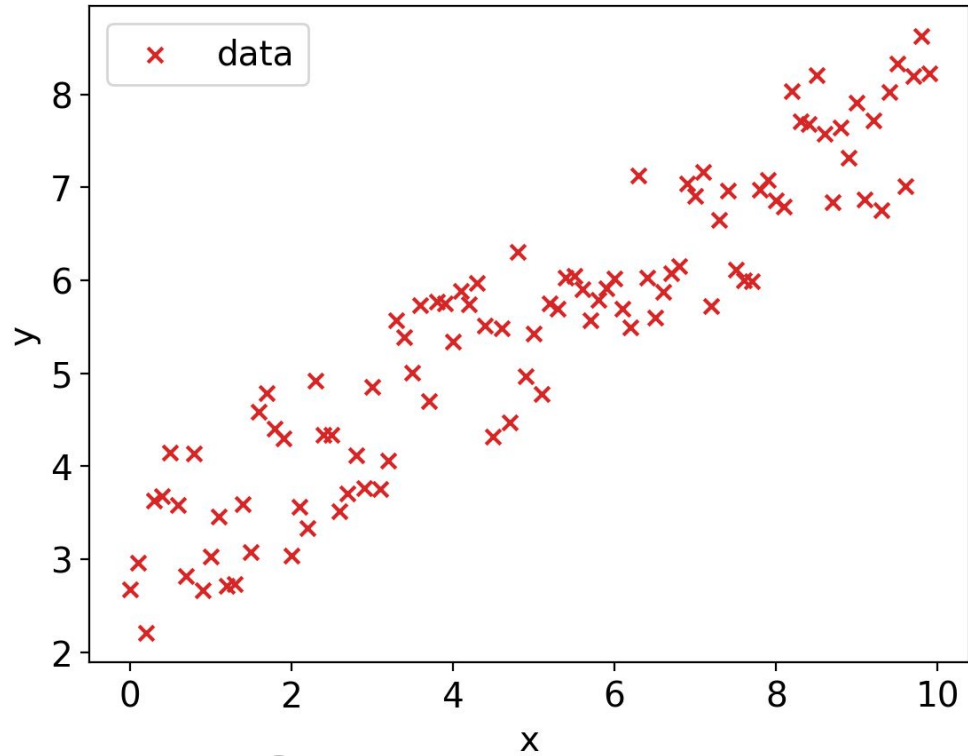
(Full) Batch Gradient Descent $\theta_{t+1} = \theta_t - \gamma \frac{\partial \mathcal{L}}{\partial \theta_t}$

Stochastic Gradient Descent $\theta_{t+1} = \theta_t - \gamma \frac{\partial \ell}{\partial \theta_t}$

Mini-Batch Gradient Descent $\theta_{t+1} = \theta_t - \gamma \frac{\partial \frac{1}{|B|} \sum_{b \in B} \ell(x_b, y_b)}{\partial \theta_t}$

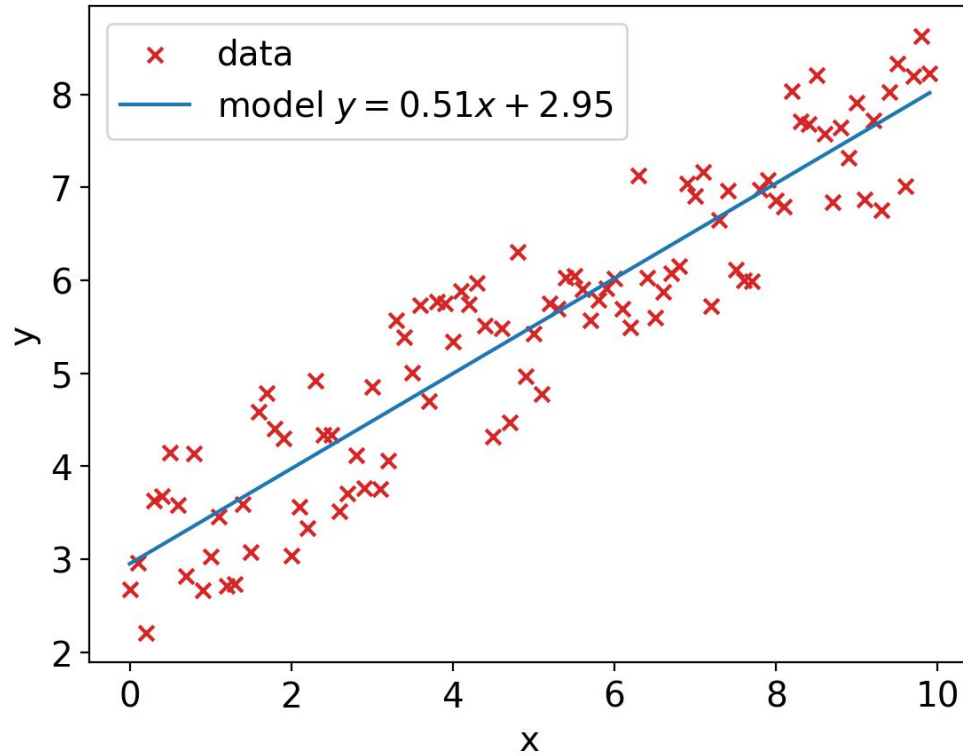
Models

Models



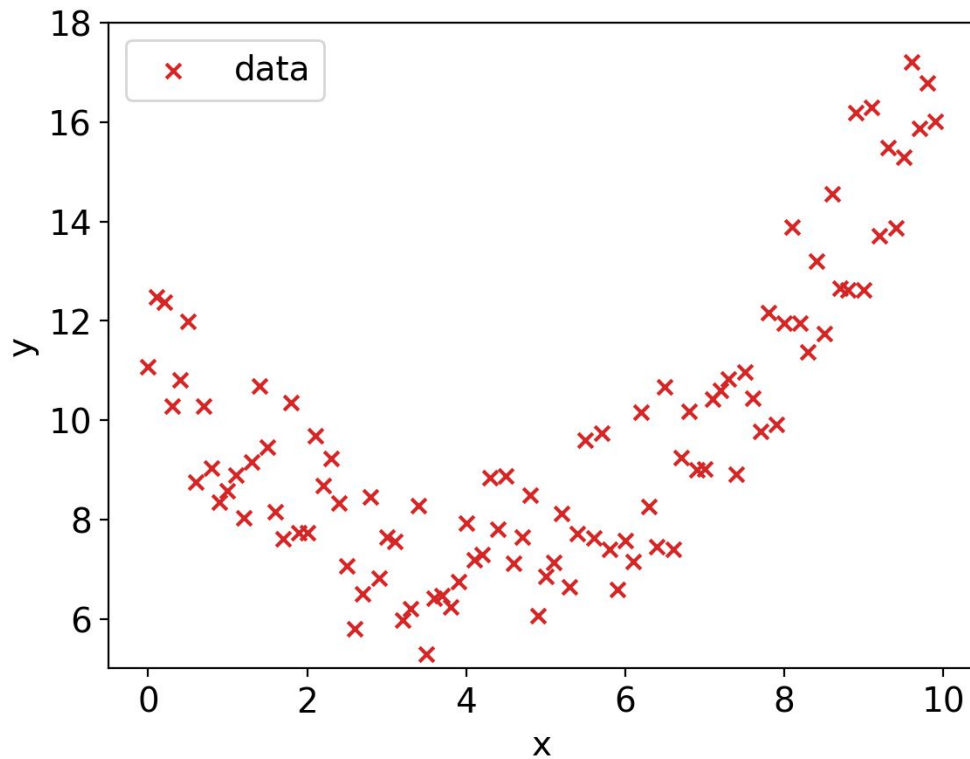
?

Models



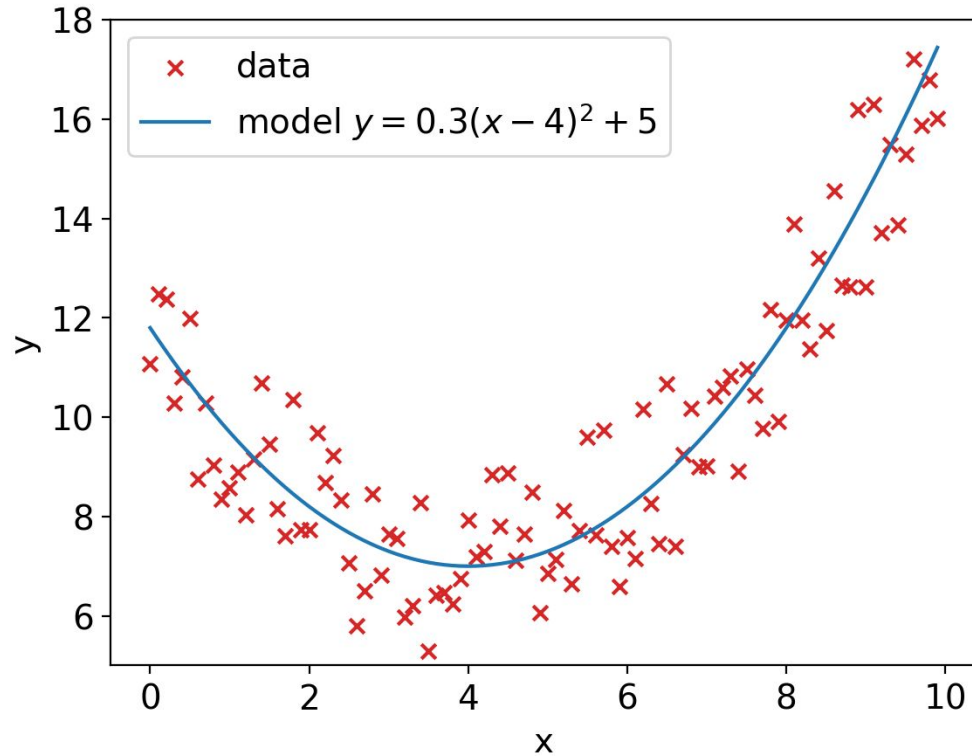
⇒ *Linear model*

Models



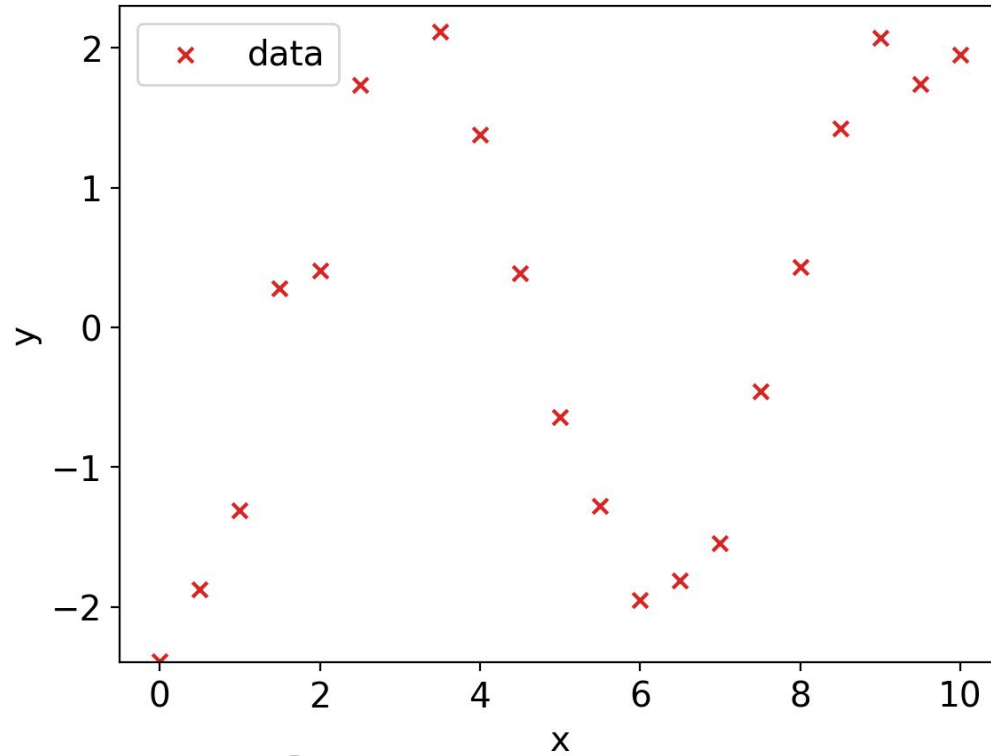
?

Models



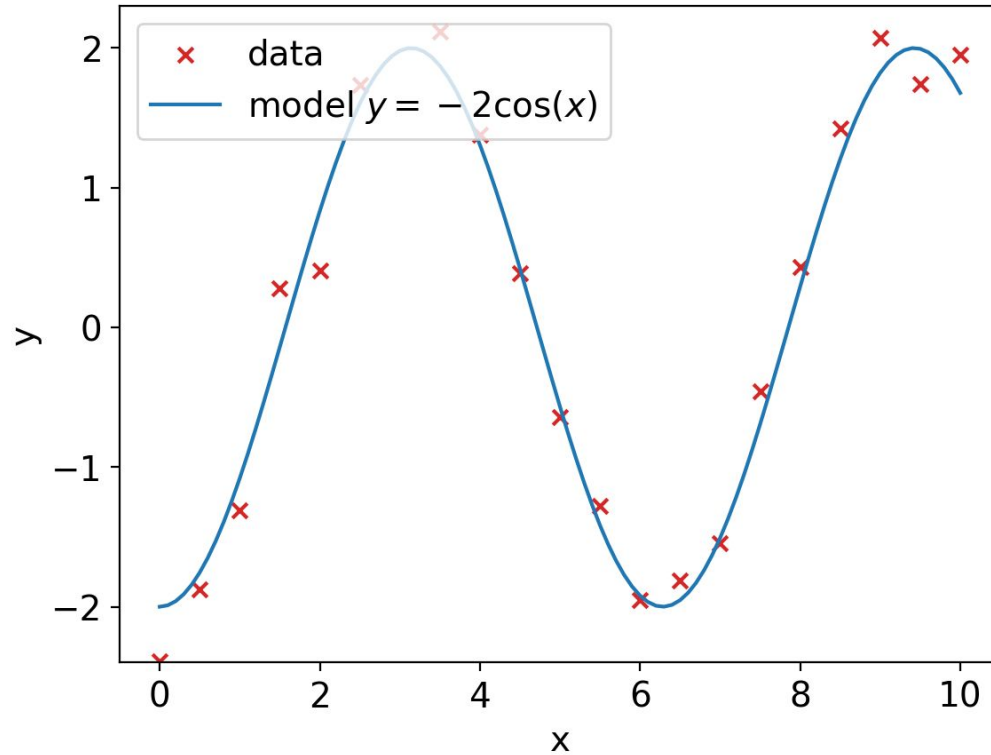
⇒ *Quadratic model*

Models



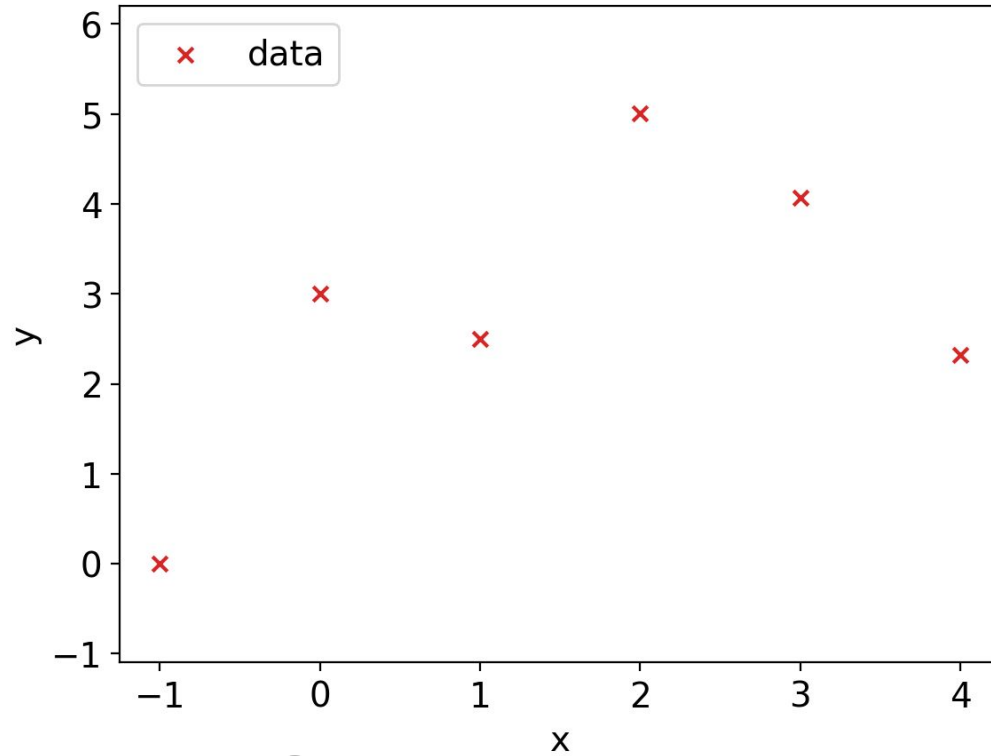
?

Models



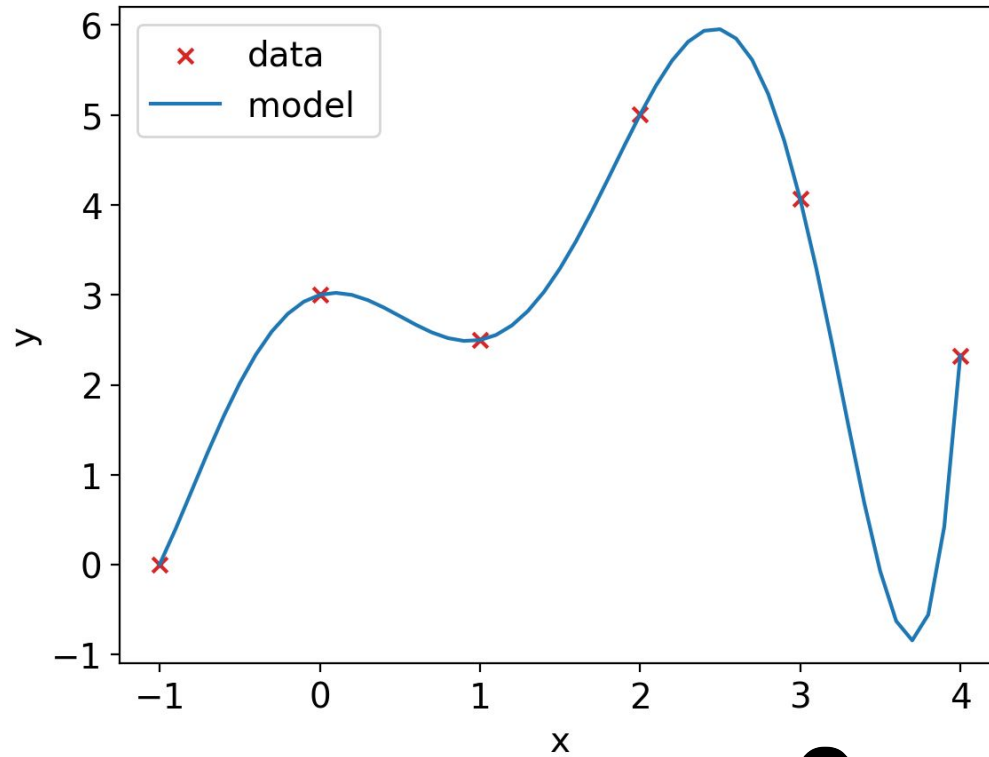
\Rightarrow *Cosine model*

Models



?

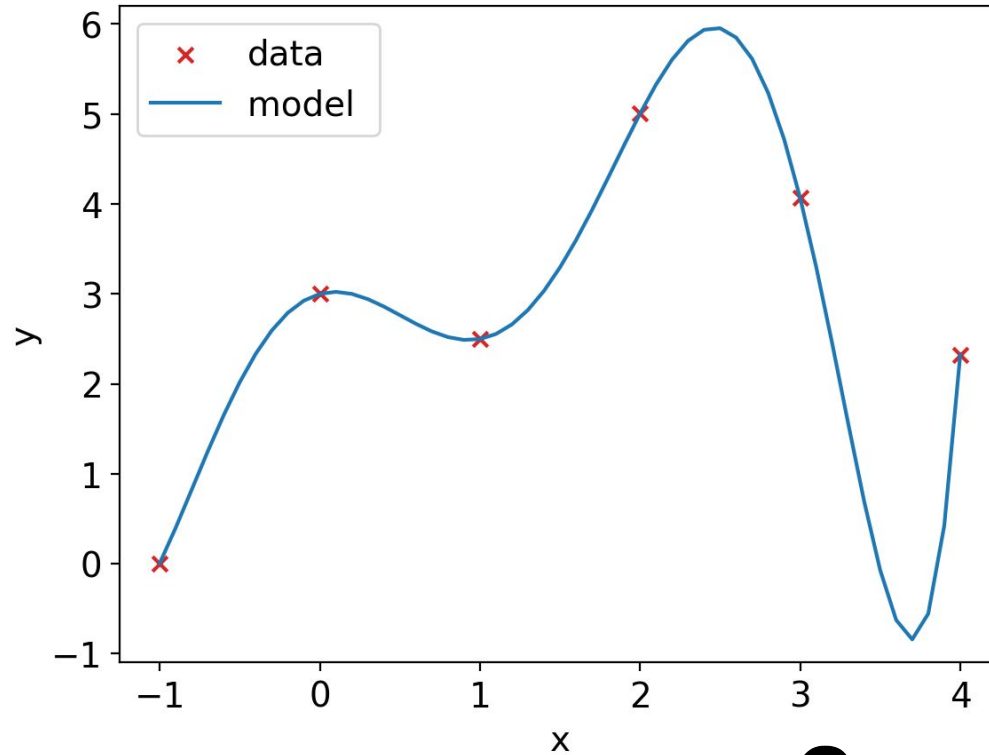
Models



⇒ *Polynomial model*



Models

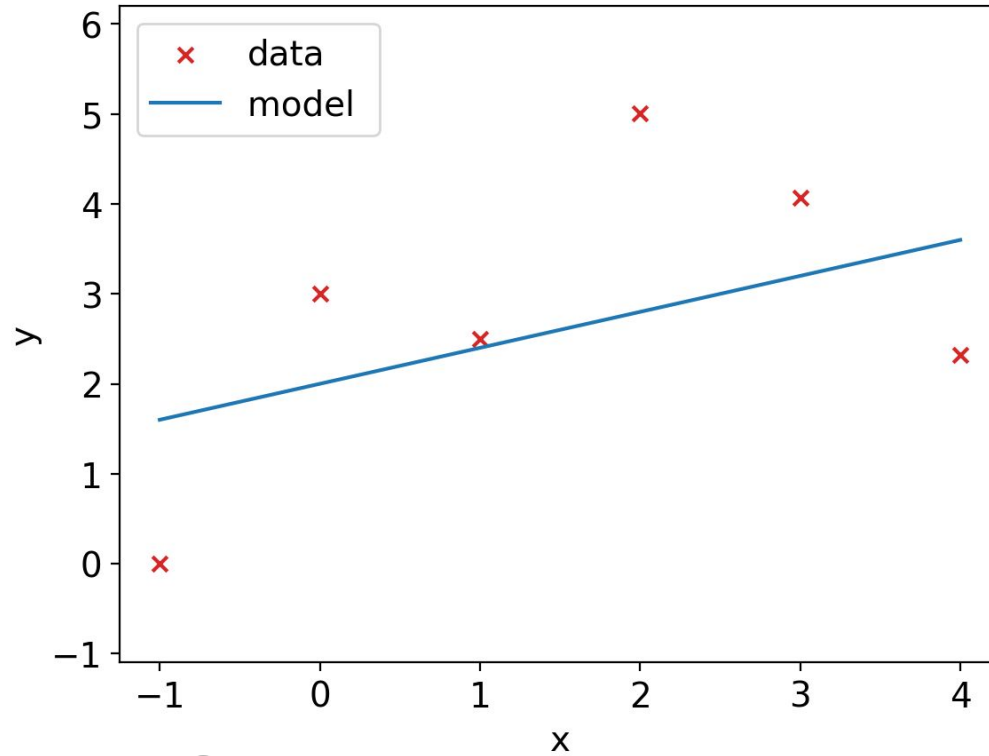


⇒ *Overfitting*

⇒ *Polynomial model*

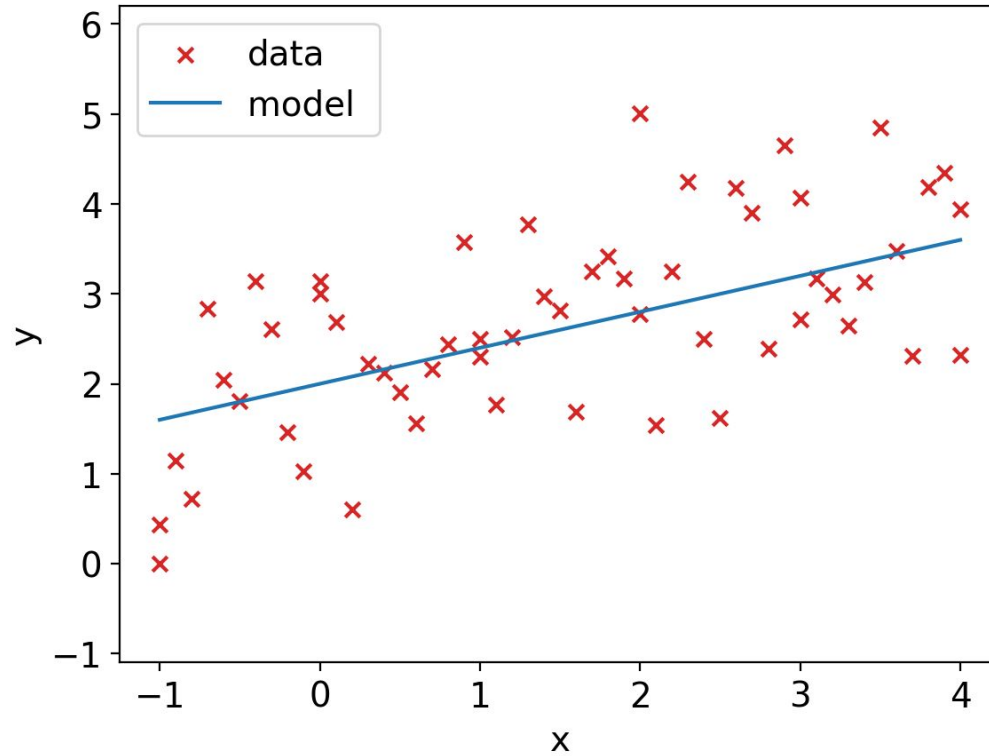


Models



⇒ *Linear model?*

Models



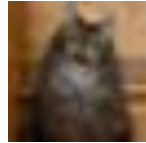
⇒ *Linear model*

Models

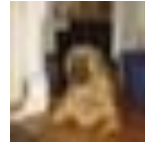
data:



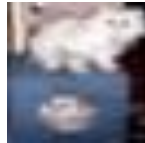
cat



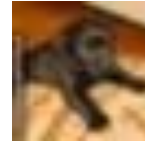
cat



dog



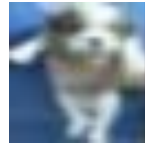
cat



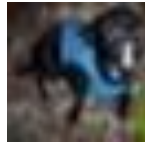
dog



cat



dog



dog



Models

data:



Linear model?
Quadratic model?

Python

`linear_model.py`

How to Choose a Model?



How to Choose a Model?

The Bitter Lesson (2019)



(Distilled version)

General methods that leverage computation and learning ultimately prove more effective than approaches relying on human knowledge and domain-specific expertise.

— Rich Sutton, ACM A. M. Turing Award 2024 recipient

<http://www.incompleteideas.net/InIdeas/BitterLesson.html>

How to Choose a Model?

The Bitter Lesson (2019)



(Distilled version)

General methods that leverage computation and learning ultimately prove more effective than approaches relying on human knowledge and domain-specific expertise.

— Rich Sutton, ACM A. M. Turing Award 2024 recipient

<http://www.incompleteideas.net/InIdeas/BitterLesson.html>

(\Rightarrow we should not assume models (e.g., linear or not), we should just have a general learning method leveraging computation, e.g., NNs + GD)

Neural Networks

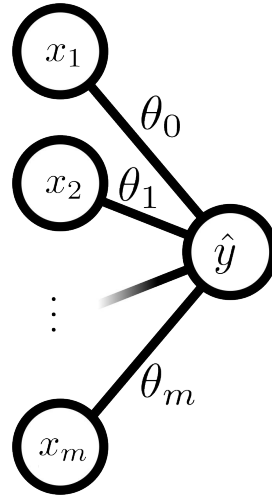


$$\hat{y} = \mathbf{x} \boldsymbol{\theta} + \mathbf{b} = \sum_{i=0}^m x_i \theta_i + b$$

$$\mathbf{x} = [x_1, x_2, \dots, x_m]$$

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_m \end{bmatrix}$$

$$\mathbf{b} \in \mathbb{R}$$



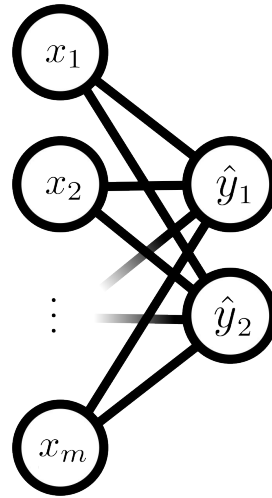
Neural Networks

$$\hat{y} = \mathbf{x} \boldsymbol{\theta} + \mathbf{b} = \left[\sum_{i=0}^m x_i \theta_{1i} + b_1 \quad \sum_{i=0}^m x_i \theta_{2i} + b_2 \right]$$

$$\mathbf{x} = [x_1, x_2, \dots, x_m]$$

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_{11} & \theta_{21} \\ \vdots & \vdots \\ \theta_{1m} & \theta_{2m} \end{bmatrix}$$

$$\mathbf{b} = [b_1 \quad b_2]$$



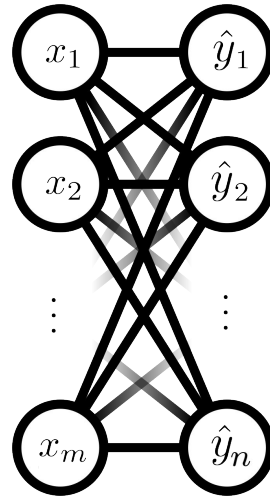
Neural Networks


$$\hat{y} = \mathbf{x} \boldsymbol{\theta} + \mathbf{b} = \left[\sum_{i=0}^m x_i \theta_{1i} + b_1 \quad \dots \quad \sum_{i=0}^m x_i \theta_{ni} + b_n \right]$$

$$\mathbf{x} = [x_1, x_2, \dots, x_m]$$

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_{11} & \theta_{21} & \dots & \theta_{n1} \\ \vdots & \vdots & \ddots & \vdots \\ \theta_{1m} & \theta_{2m} & \dots & \theta_{nm} \end{bmatrix}$$

$$\mathbf{b} = [b_1 \quad b_2 \quad \dots \quad b_n]$$



 (n → m)

?

Neural Networks

$$\hat{y} = \sigma(\mathbf{x} \boldsymbol{\theta}^{(1)} + \mathbf{b}^{(1)}) \boldsymbol{\theta}^{(2)} + \mathbf{b}^{(2)}$$



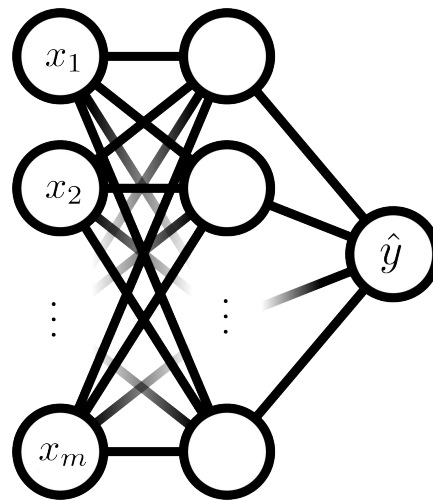
(nonlinearities)

$$\mathbf{x} = [x_1, x_2, \dots, x_m]$$

$$\boldsymbol{\theta} = (\boldsymbol{\theta}^{(1)}, \mathbf{b}^{(1)}, \boldsymbol{\theta}^{(2)}, \mathbf{b}^{(2)})$$

$$\boldsymbol{\theta}^{(1)} - M \times N \quad \mathbf{b}^{(1)} - N \times 1$$

$$\boldsymbol{\theta}^{(2)} - N \times 1 \quad \mathbf{b}^{(2)} \in \mathbb{R}$$



Neural Networks

$$\hat{y} = \sigma(\dots \sigma(\mathbf{x} \boldsymbol{\theta}^{(1)} + \mathbf{b}^{(1)}) \boldsymbol{\theta}^{(2)} + \mathbf{b}^{(2)} \dots) \boldsymbol{\theta}^{(L)} + \mathbf{b}^{(L)}$$

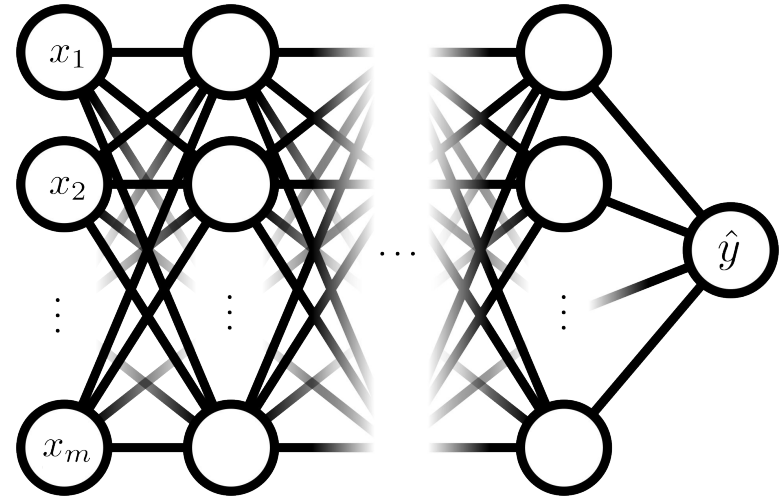
$$\mathbf{x} = [x_1, x_2, \dots, x_m]$$

$$\boldsymbol{\theta} = (\boldsymbol{\theta}^{(1)}, \mathbf{b}^{(1)}, \dots, \boldsymbol{\theta}^{(L)}, \mathbf{b}^{(L)})$$

$$\boldsymbol{\theta}^{(1)} - M \times N \quad \mathbf{b}^{(1)} - N \times 1$$

$$\boldsymbol{\theta}^{(2)} - N \times 1 \quad \mathbf{b}^{(2)} \in \mathbb{R}$$

...



Neural Networks

$$\hat{y} = \sigma(\dots \sigma(\mathbf{x} \boldsymbol{\theta}^{(1)} + \mathbf{b}^{(1)}) \boldsymbol{\theta}^{(2)} + \mathbf{b}^{(2)} \dots) \boldsymbol{\theta}^{(L)} + \mathbf{b}^{(L)}$$

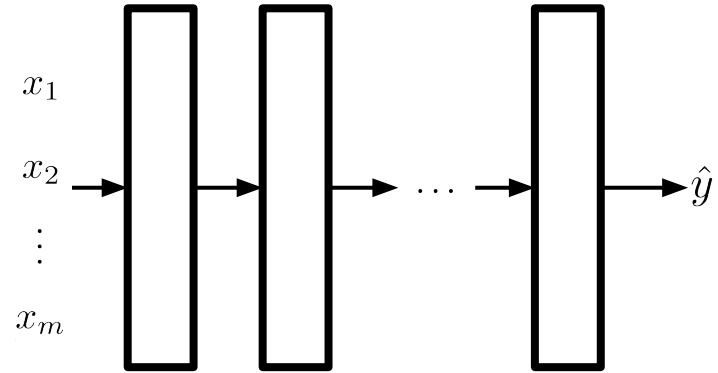
$$\mathbf{x} = [x_1, x_2, \dots, x_m]$$

$$\boldsymbol{\theta} = (\boldsymbol{\theta}^{(1)}, \mathbf{b}^{(1)}, \dots, \boldsymbol{\theta}^{(L)}, \mathbf{b}^{(L)})$$

$$\boldsymbol{\theta}^{(1)} - M \times N \quad \mathbf{b}^{(1)} - N \times 1$$

$$\boldsymbol{\theta}^{(2)} - N \times 1 \quad \mathbf{b}^{(2)} \in \mathbb{R}$$

...



Backpropagation Example

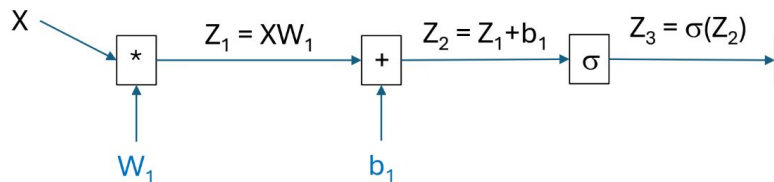
$$Y^p = \sigma(XW_1 + b_1)W_2 + b_2$$

Backpropagation Example

$$Y^p = \sigma(XW_1 + b_1)W_2 + b_2 \quad L(Y^p, Y) = \frac{1}{2}(Y^p - Y)^2$$

Backpropagation Example

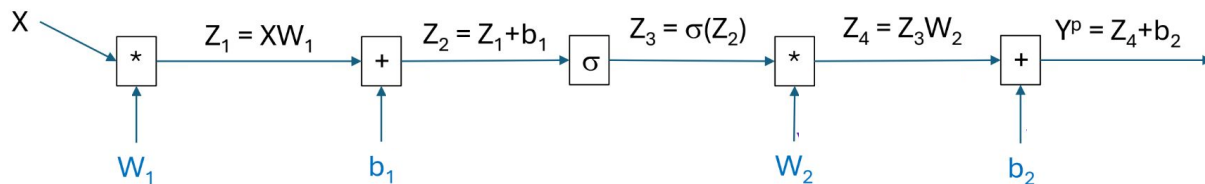
$$Y^p = \sigma(XW_1 + b_1)W_2 + b_2 \quad L(Y^p, Y) = \frac{1}{2}(Y^p - Y)^2$$



Backpropagation Example

$$Y^p = \sigma(XW_1 + b_1)W_2 + b_2$$

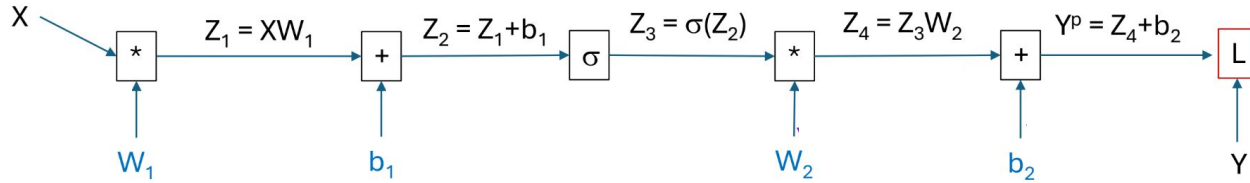
$$L(Y^p, Y) = \frac{1}{2} (Y^p - Y)^2$$



Backpropagation Example

$$Y^p = \sigma(XW_1 + b_1)W_2 + b_2$$

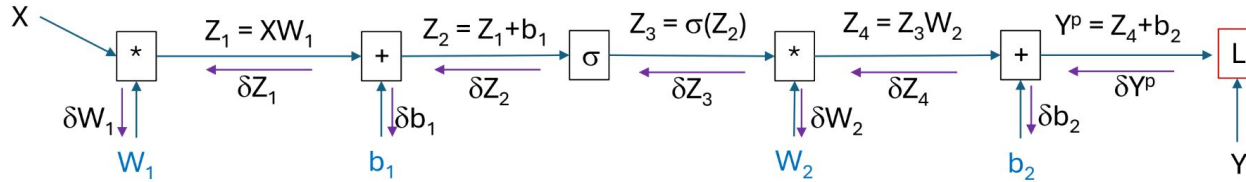
$$L(Y^p, Y) = \frac{1}{2} (Y^p - Y)^2$$



Backpropagation Example



$$Y^p = \sigma(XW_1 + b_1)W_2 + b_2 \quad L(Y^p, Y) = \frac{1}{2}(Y^p - Y)^2$$

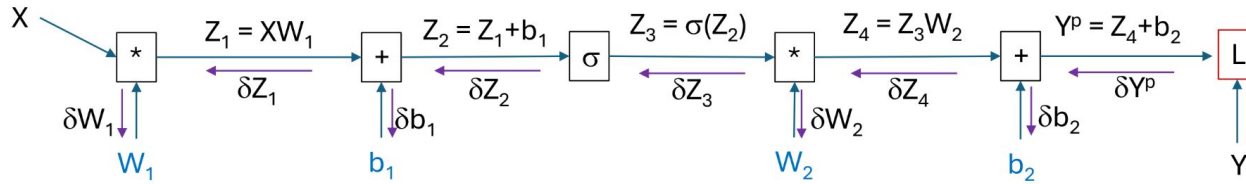


Backpropagation Implementation

See Andrej Karpathy's [micrograd](#) and [backpropagation lecture material](#)

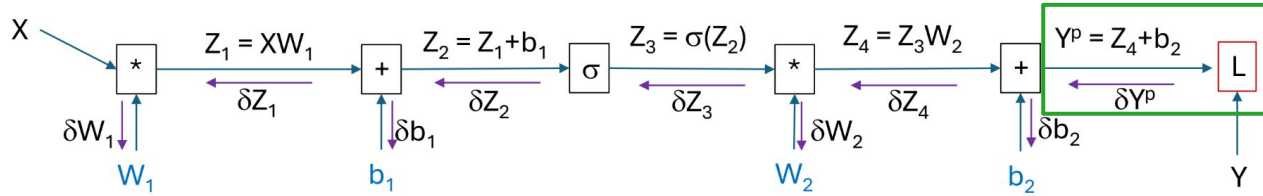
Backpropagation Example

$$Y^p = \sigma(XW_1 + b_1)W_2 + b_2 \quad L(Y^p, Y) = \frac{1}{2}(Y^p - Y)^2$$



Backpropagation Example

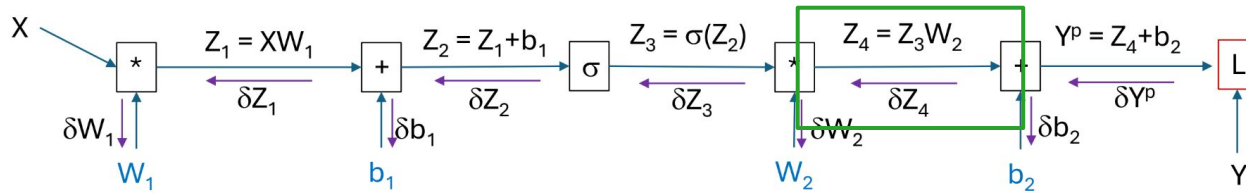
$$Y^p = \sigma(XW_1 + b_1)W_2 + b_2 \quad L(Y^p, Y) = \frac{1}{2}(Y^p - Y)^2$$



1
$$\delta Y^p \equiv \frac{\partial L}{\partial Y^p} = \frac{\partial}{\partial Y^p} \left[\frac{1}{2} \|Y^p - Y\|^2 \right] = y^p - y$$

Backpropagation Example

$$Y^p = \sigma(XW_1 + b_1)W_2 + b_2 \quad L(Y^p, Y) = \frac{1}{2} (Y^p - Y)^2$$



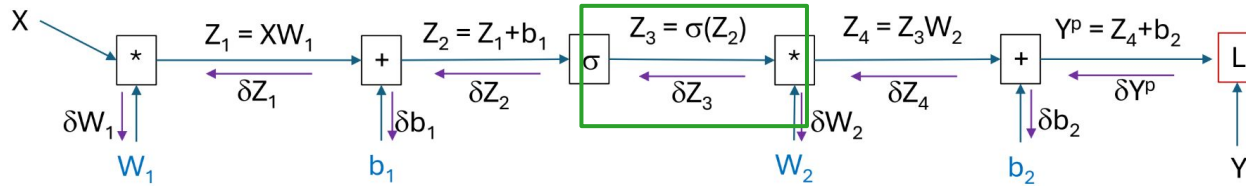
1
$$\delta Y^p \equiv \frac{\partial L}{\partial Y^p} = \frac{\partial}{\partial Y^p} \left[\frac{1}{2} \|Y^p - Y\|^2 \right] = y^p - y$$

2
$$\delta Z_4 \equiv \frac{\partial L}{\partial Z_4} = \frac{\partial Y^p}{\partial Z_4} \frac{\partial L}{\partial Y^p} = \delta Y^p$$

Because $Y^p = Z_4 + b_2$, $\frac{\partial Y^p}{\partial Z_4} = 1$

Backpropagation Example

$$Y^p = \sigma(XW_1 + b_1)W_2 + b_2 \quad L(Y^p, Y) = \frac{1}{2}(Y^p - Y)^2$$



$$1 \quad \delta Y^p \equiv \frac{\partial L}{\partial Y^p} = \frac{\partial}{\partial Y^p} \left[\frac{1}{2} \|Y^p - Y\|^2 \right] = y^p - y$$

$$2 \quad \delta Z_4 \equiv \frac{\partial L}{\partial Z_4} = \frac{\partial Y^p}{\partial Z_4} \frac{\partial L}{\partial Y^p} = \delta Y^p$$

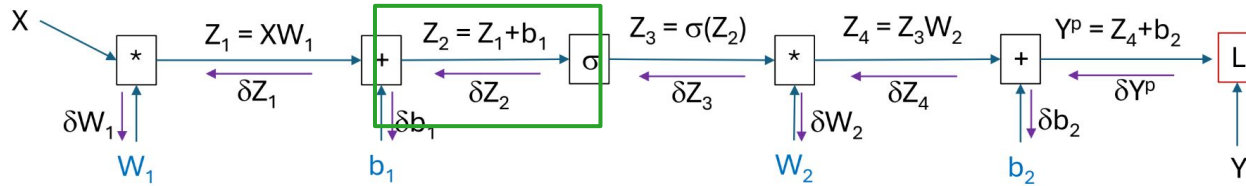
Because $Y^p = Z_4 + b_2$, $\frac{\partial Y^p}{\partial Z_4} = 1$

$$3 \quad \delta Z_3 \equiv \frac{\partial L}{\partial Z_3} = \frac{\partial Z_4}{\partial Z_3} \frac{\partial L}{\partial Z_4} = \delta Z_4 W_2^T$$

Because $Z_4 = Z_3 W_2$, $\frac{\partial Z_4}{\partial Z_3} = W_2^T$

Backpropagation Example

$$Y^p = \sigma(XW_1 + b_1)W_2 + b_2 \quad L(Y^p, Y) = \frac{1}{2}(Y^p - Y)^2$$



$$1 \quad \delta Y^p \equiv \frac{\partial L}{\partial Y^p} = \frac{\partial}{\partial Y^p} \left[\frac{1}{2} \|Y^p - Y\|^2 \right] = y^p - y$$

$$4 \quad \delta Z_2 \equiv \frac{\partial L}{\partial Z_2} = \frac{\partial Z_3}{\partial Z_2} \frac{\partial L}{\partial Z_3} = \sigma'(Z_2) \cdot \delta Z_3$$

$$2 \quad \delta Z_4 \equiv \frac{\partial L}{\partial Z_4} = \frac{\partial Y^p}{\partial Z_4} \frac{\partial L}{\partial Y^p} = \delta Y^p$$

Because $Z_3 = \sigma(Z_2)$, $\frac{\partial Z_3}{\partial Z_2} = \sigma'(Z_2)$

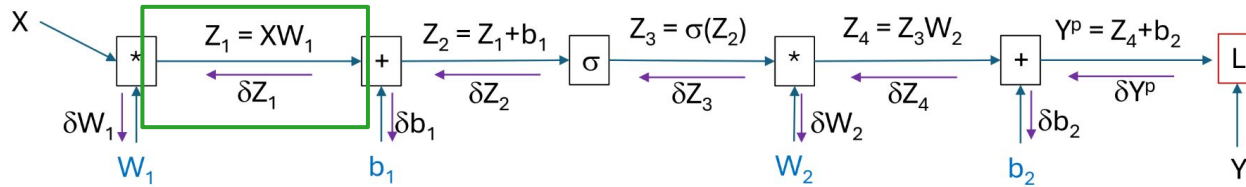
Because $Y^p = Z_4 + b_2$, $\frac{\partial Y^p}{\partial Z_4} = 1$

$$3 \quad \delta Z_3 \equiv \frac{\partial L}{\partial Z_3} = \frac{\partial Z_4}{\partial Z_3} \frac{\partial L}{\partial Z_4} = \delta Z_4 W_2^T$$

Because $Z_4 = Z_3 W_2$, $\frac{\partial Z_4}{\partial Z_3} = W_2^T$

Backpropagation Example

$$Y^p = \sigma(XW_1 + b_1)W_2 + b_2 \quad L(Y^p, Y) = \frac{1}{2} (Y^p - Y)^2$$



$$1 \quad \delta Y^p \equiv \frac{\partial L}{\partial Y^p} = \frac{\partial}{\partial Y^p} \left[\frac{1}{2} \|Y^p - Y\|^2 \right] = y^p - y$$

$$4 \quad \delta Z_2 \equiv \frac{\partial L}{\partial Z_2} = \frac{\partial Z_3}{\partial Z_2} \frac{\partial L}{\partial Z_3} = \sigma'(Z_2) \cdot \delta Z_3$$

$$2 \quad \delta Z_4 \equiv \frac{\partial L}{\partial Z_4} = \frac{\partial Y^p}{\partial Z_4} \frac{\partial L}{\partial Y^p} = \delta Y^p$$

Because $Z_3 = \sigma(Z_2)$, $\frac{\partial Z_3}{\partial Z_2} = \sigma'(Z_2)$

Because $Y^p = Z_4 + b_2$, $\frac{\partial Y^p}{\partial Z_4} = 1$

$$5 \quad \delta Z_1 \equiv \frac{\partial L}{\partial Z_1} = \frac{\partial Z_2}{\partial Z_1} \frac{\partial L}{\partial Z_2} = \delta Z_2$$

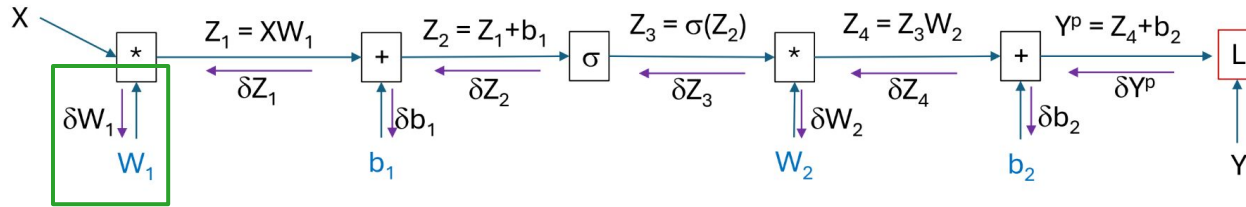
$$3 \quad \delta Z_3 \equiv \frac{\partial L}{\partial Z_3} = \frac{\partial Z_4}{\partial Z_3} \frac{\partial L}{\partial Z_4} = \delta Z_4 W_2^T$$

Because $Z_2 = Z_1 + b_1$, $\frac{\partial Z_2}{\partial Z_1} = 1$

Because $Z_4 = Z_3 W_2$, $\frac{\partial Z_4}{\partial Z_3} = W_2^T$

Backpropagation Example

$$Y^p = \sigma(XW_1 + b_1)W_2 + b_2 \quad L(Y^p, Y) = \frac{1}{2}(Y^p - Y)^2$$

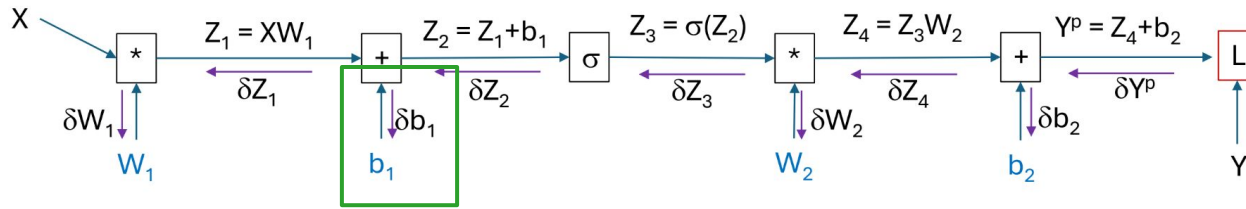


1 $\delta W_1 \equiv \frac{\partial L}{\partial W_1} = \frac{\partial Z_1}{\partial W_1} \frac{\partial L}{\partial Z_1} = \mathbf{X}^T \delta Z_1$

Because $Z_1 = XW_1$, $\frac{\partial Z_1}{\partial W_1} = \mathbf{X}^T$

Backpropagation Example

$$Y^p = \sigma(XW_1 + b_1)W_2 + b_2 \quad L(Y^p, Y) = \frac{1}{2}(Y^p - Y)^2$$



$$1 \quad \delta W_1 \equiv \frac{\partial L}{\partial W_1} = \frac{\partial Z_1}{\partial W_1} \frac{\partial L}{\partial Z_1} = \mathbf{X}^T \delta Z_1$$

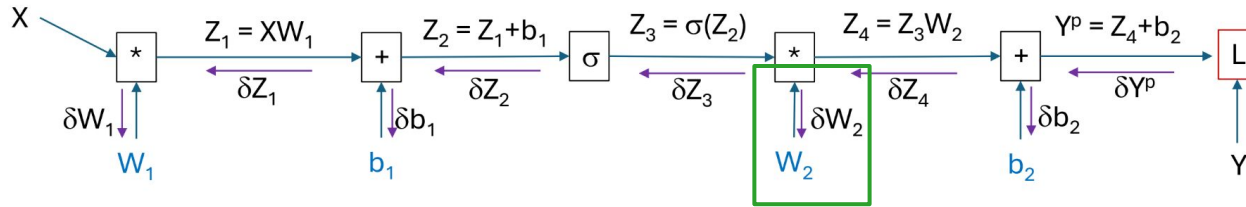
$$\text{Because } Z_1 = XW_1, \frac{\partial Z_1}{\partial W_1} = \mathbf{X}^T$$

$$2 \quad \delta b_1 \equiv \frac{\partial L}{\partial b_1} = \frac{\partial Z_2}{\partial b_1} \frac{\partial L}{\partial Z_2} = \sum_k (\delta Z_2)_k$$

$$\text{Because } Z_2 = Z_1 + b_1, \frac{\partial Z_2}{\partial b_1} = [\mathbf{1}, \dots, \mathbf{1}]$$

Backpropagation Example

$$Y^p = \sigma(XW_1 + b_1)W_2 + b_2 \quad L(Y^p, Y) = \frac{1}{2} (Y^p - Y)^2$$



$$1 \quad \delta W_1 \equiv \frac{\partial L}{\partial W_1} = \frac{\partial Z_1}{\partial W_1} \frac{\partial L}{\partial Z_1} = \mathbf{X}^T \delta Z_1$$

$$\text{Because } Z_1 = XW_1, \frac{\partial Z_1}{\partial W_1} = \mathbf{X}^T$$

$$3 \quad \delta W_2 \equiv \frac{\partial L}{\partial W_2} = \frac{\partial Z_4}{\partial W_2} \frac{\partial L}{\partial Z_4} = \mathbf{Z}_3^T \delta Z_4$$

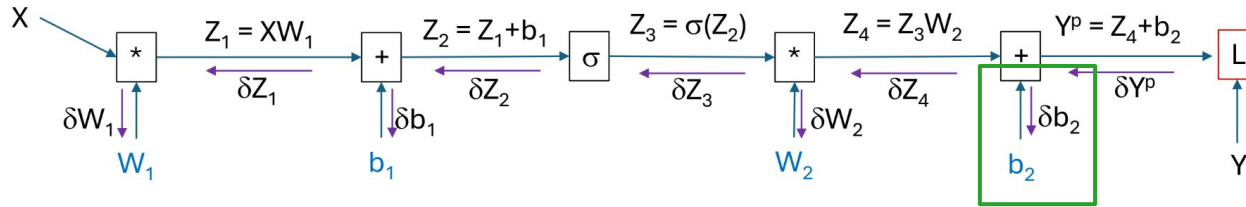
$$\text{Because } Z_4 = Z_3W_2, \frac{\partial Z_4}{\partial W_2} = \mathbf{Z}_3^T$$

$$2 \quad \delta b_1 \equiv \frac{\partial L}{\partial b_1} = \frac{\partial Z_2}{\partial b_1} \frac{\partial L}{\partial Z_2} = \sum_k (\delta Z_2)_k$$

$$\text{Because } Z_2 = Z_1 + b_1, \frac{\partial Z_2}{\partial b_1} = [\mathbf{1}, \dots, \mathbf{1}]$$

Backpropagation Example

$$Y^p = \sigma(XW_1 + b_1)W_2 + b_2 \quad L(Y^p, Y) = \frac{1}{2} (Y^p - Y)^2$$



$$1 \quad \delta W_1 \equiv \frac{\partial L}{\partial W_1} = \frac{\partial Z_1}{\partial W_1} \frac{\partial L}{\partial Z_1} = \mathbf{X}^T \delta Z_1$$

$$\text{Because } Z_1 = XW_1, \frac{\partial Z_1}{\partial W_1} = \mathbf{X}^T$$

$$2 \quad \delta b_1 \equiv \frac{\partial L}{\partial b_1} = \frac{\partial Z_2}{\partial b_1} \frac{\partial L}{\partial Z_2} = \sum_k (\delta Z_2)_k$$

$$\text{Because } Z_2 = Z_1 + b_1, \frac{\partial Z_2}{\partial b_1} = [\mathbf{1}, \dots, \mathbf{1}]$$

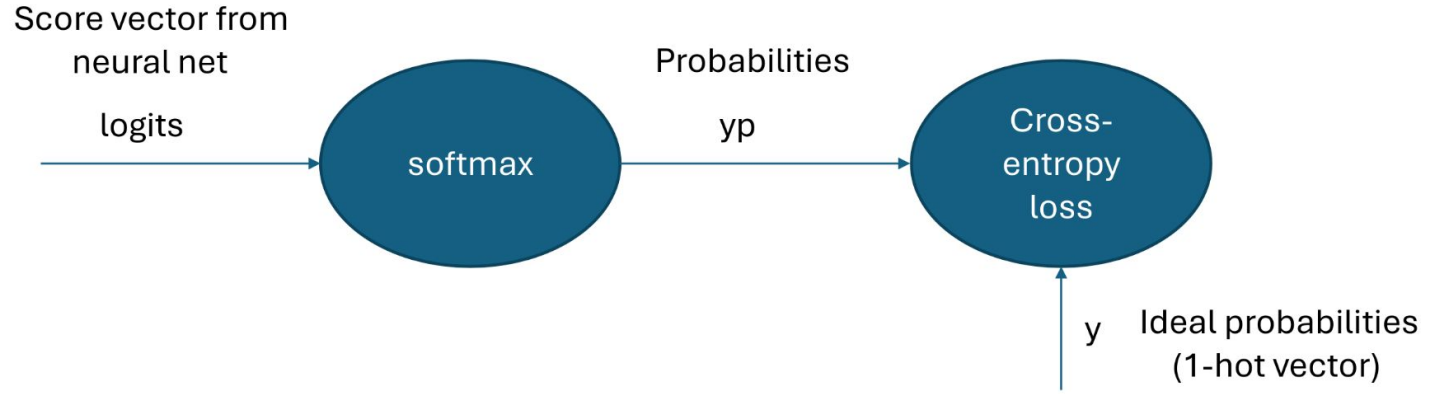
$$3 \quad \delta W_2 \equiv \frac{\partial L}{\partial W_2} = \frac{\partial Z_4}{\partial W_2} \frac{\partial L}{\partial Z_4} = \mathbf{Z}_3^T \delta Z_4$$

$$\text{Because } Z_4 = Z_3W_2, \frac{\partial Z_4}{\partial W_2} = \mathbf{Z}_3^T$$

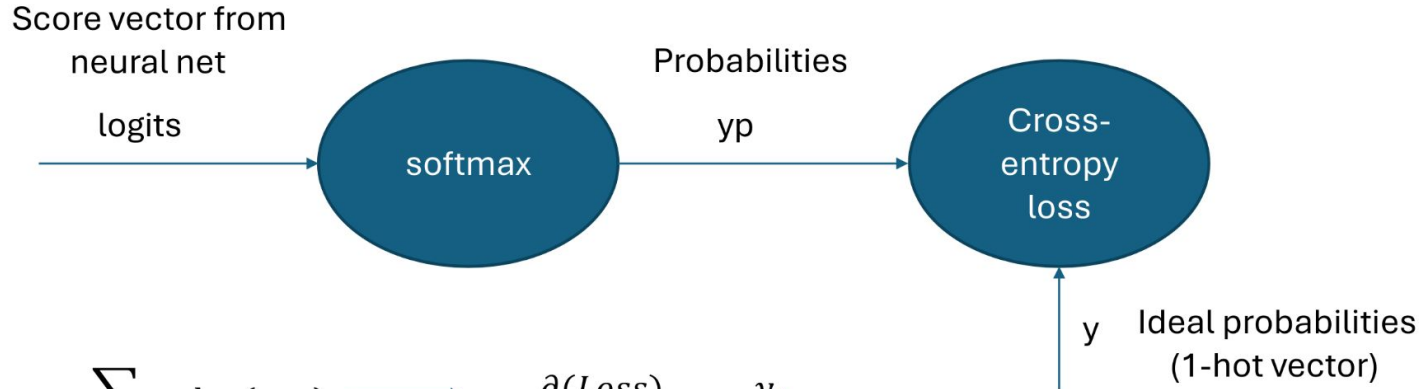
$$4 \quad \delta b_2 \equiv \frac{\partial L}{\partial b_2} = \frac{\partial Y^p}{\partial b_2} \frac{\partial L}{\partial Y^p} = \sum_k (\delta Y^p)_k$$

$$\text{Because } Y^p = Z_4 + b_2, \frac{\partial Y^p}{\partial b_2} = [\mathbf{1}, \dots, \mathbf{1}]$$

Backpropagation (Softmax & CE Loss)

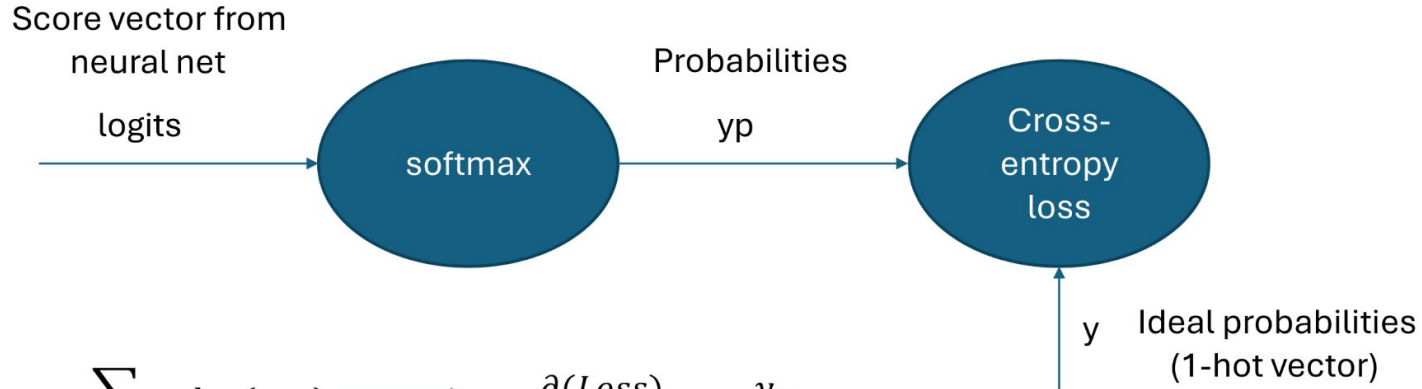


Backpropagation (Softmax & CE Loss)



$$Loss = - \sum_k y_k \log(yp_k) \longrightarrow \frac{\partial(Loss)}{\partial(yp)_k} = - \frac{y_k}{yp_k}$$

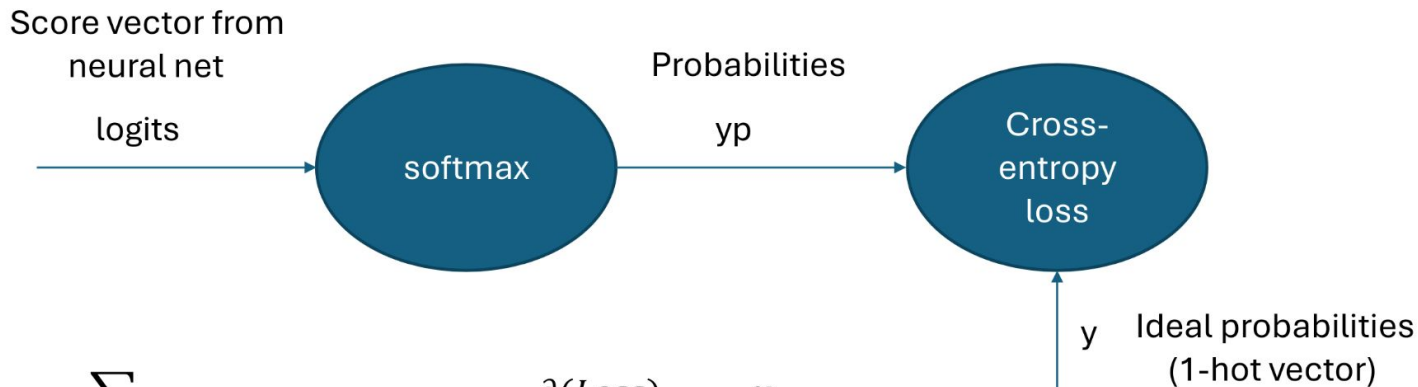
Backpropagation (Softmax & CE Loss)



$$Loss = - \sum_k y_k \log(yp_k) \implies \frac{\partial(Loss)}{\partial(yp)_k} = - \frac{y_k}{yp_k}$$

$$yp_i = \frac{\exp(logits_i)}{\sum_k \exp(logits_k)} \implies \frac{\partial(yp)_k}{\partial(logits)_i} = \begin{cases} yp_i(1 - yp_i), & \text{if } i = k, \\ -yp_i yp_k, & \text{otherwise.} \end{cases}$$

Backpropagation (Softmax & CE Loss)



$$Loss = - \sum_k y_k \log(yp_k) \implies \frac{\partial(Loss)}{\partial(yp)_k} = - \frac{y_k}{yp_k}$$

$$yp_i = \frac{\exp(logits_i)}{\sum_k \exp(logits_k)} \implies \frac{\partial(yp)_k}{\partial(logits)_i} = \begin{cases} yp_i(1 - yp_i), & \text{if } i = k, \\ -yp_i yp_k, & \text{otherwise.} \end{cases}$$

Using the above two results in the chain rule, $\delta(logits)_i \equiv \frac{\partial(Loss)}{\partial(logits)_i} = \sum_k \frac{\partial(yp)_k}{\partial(logits)_i} \frac{\partial(Loss)}{\partial(yp)_k} = yp_i - y_i$

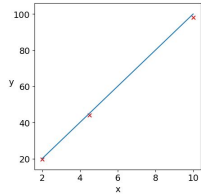
Backpropagation (More Explanations)

See the [Back Propagation Practice](#) folder

Recap

Recap

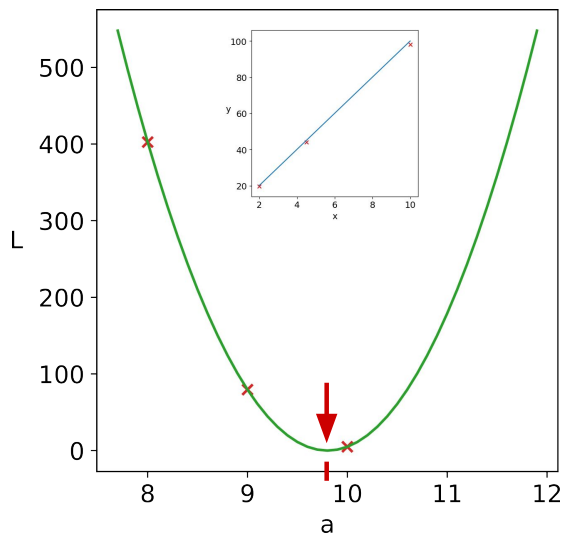
$$\hat{y} = ax$$



Recap

$$\mathcal{L} = \sum_i (y_i - \hat{y}(x_i))^2 \quad \text{L2 Loss} \quad \text{📎}$$

$$\hat{y} = ax$$



Recap

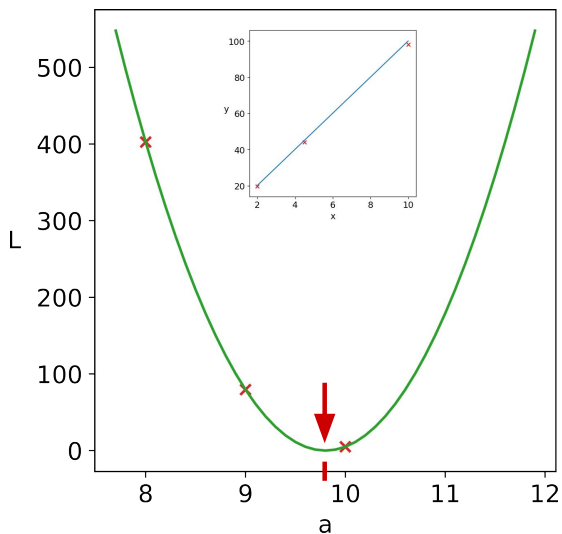
$$\mathcal{L} = \sum_i (y_i - \hat{y}(x_i))^2 \quad \text{L2 Loss}$$

repeat:

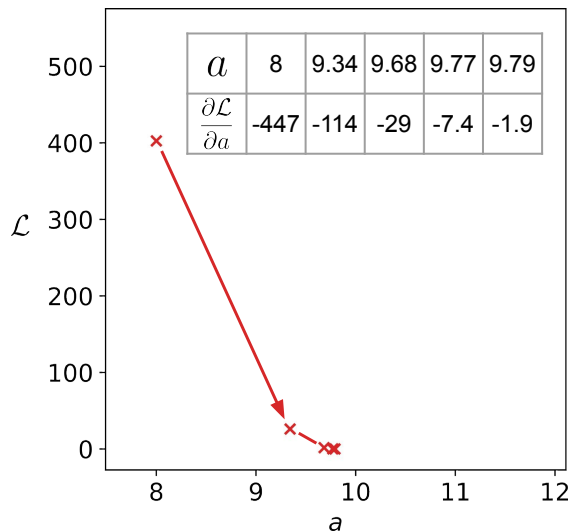
$$a := a - \gamma \frac{\partial \mathcal{L}}{\partial a}$$

until minimum is reached

$$\hat{y} = ax$$



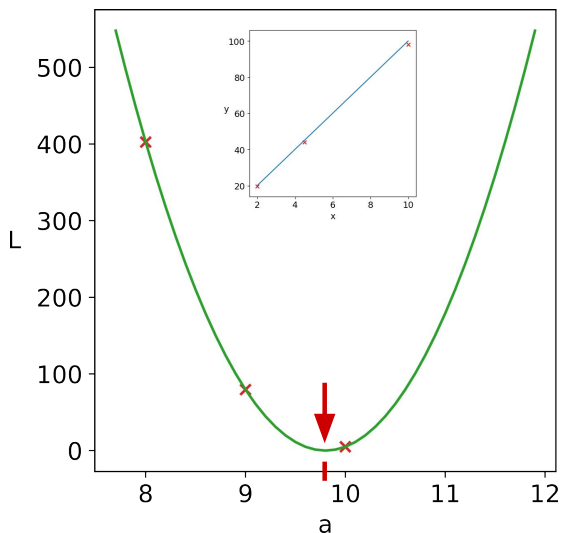
$$\gamma = 0.003$$



Recap

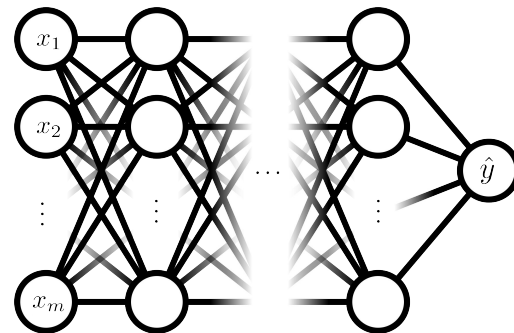
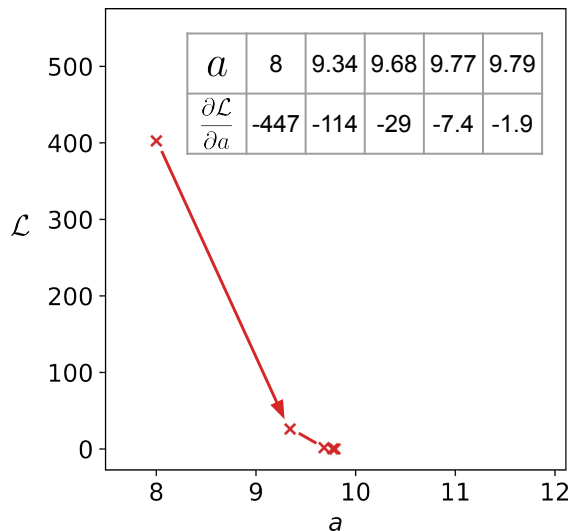
$$\mathcal{L} = \sum_i (y_i - \hat{y}(x_i))^2 \quad \text{L2 Loss}$$

$$\hat{y} = ax$$



repeat:
 $a := a - \gamma \frac{\partial \mathcal{L}}{\partial a}$
 until minimum is reached

$$\gamma = 0.003$$

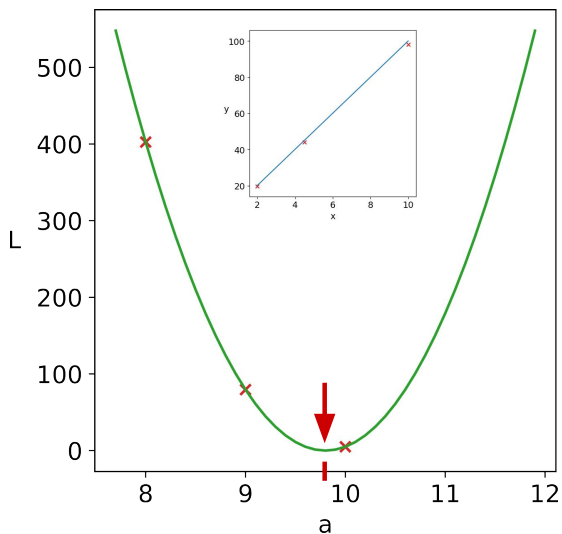


$$\hat{y} = \sigma(\dots \sigma(\mathbf{x}\boldsymbol{\theta}^{(1)} + \mathbf{b}^{(1)})\boldsymbol{\theta}^{(2)} + \mathbf{b}^{(2)} \dots)\boldsymbol{\theta}^{(L)} + \mathbf{b}^{(L)}$$

Recap

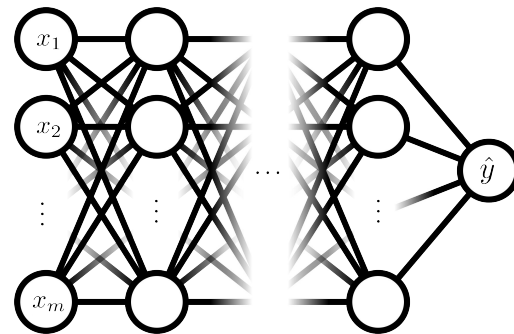
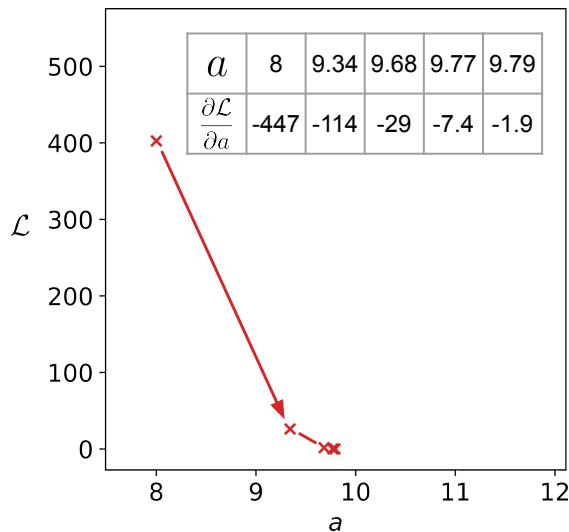
$$\mathcal{L} = \sum_i (y_i - \hat{y}(x_i))^2 \quad \text{L2 Loss}$$

$$\hat{y} = ax$$

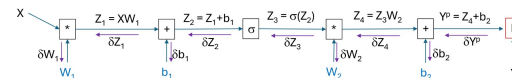


repeat:
 $a := a - \gamma \frac{\partial \mathcal{L}}{\partial a}$
 until minimum is reached

$$\gamma = 0.003$$



$$\hat{y} = \sigma(\dots \sigma(\mathbf{x}\boldsymbol{\theta}^{(1)} + \mathbf{b}^{(1)})\boldsymbol{\theta}^{(2)} + \mathbf{b}^{(2)} \dots)\boldsymbol{\theta}^{(L)} + \mathbf{b}^{(L)}$$



Part I: Fundamentals

Part II: Example Problems

Q1 (30). Consider an autoencoder and a classifier acting together for semi-supervised learning. Let us denote the symbols first. X is a batch of input images. $W_1, b_1, W_2, b_2, W_3, b_3$ are learnable parameters. The latent code is $Z = \text{ReLU}(X * W_1 + b_1)$, where $*$ denotes matrix-matrix multiplication and $+$ denotes broadcast addition. The decoder output is $Y = Z * W_2 + b_2$. The classifier is attached to the latent code as $C = \text{softmax}(Z * W_3 + b_3)$. The reconstruction loss for the autoencoder is the MSE loss, $L_2(X, Y)$, whereas the classification loss is cross entropy, $\text{CE}(C, C_{\text{gt}})$, where C_{gt} is the ground truth class probabilities with 1-hot encoding. Write the expressions for $dW_1, db_1, dW_2, db_2, dW_3$ and db_3 , where these symbols mean gradient of the combined loss, $L = L_2(X, Y) + \text{CE}(C, C_{\text{gt}})$, with respect to W_1, b_1, W_2, b_2, W_3 , and b_3 , respectively. Draw a computational graph for the entire model with loss nodes. Unlike a realistic image batch in a semi-supervised learning, for the sake of simplicity, the entire batch of images X here has ground truth C_{gt} available.

Q1 distilled:

Problem:

$$Z = \text{ReLU}(XW_1 + b_1)$$

$$Y = ZW_2 + b_2$$

$$C = \text{softmax}(ZW_3 + b_3)$$

$$L_2(X, Y) = \frac{1}{2}(X - Y)^2$$

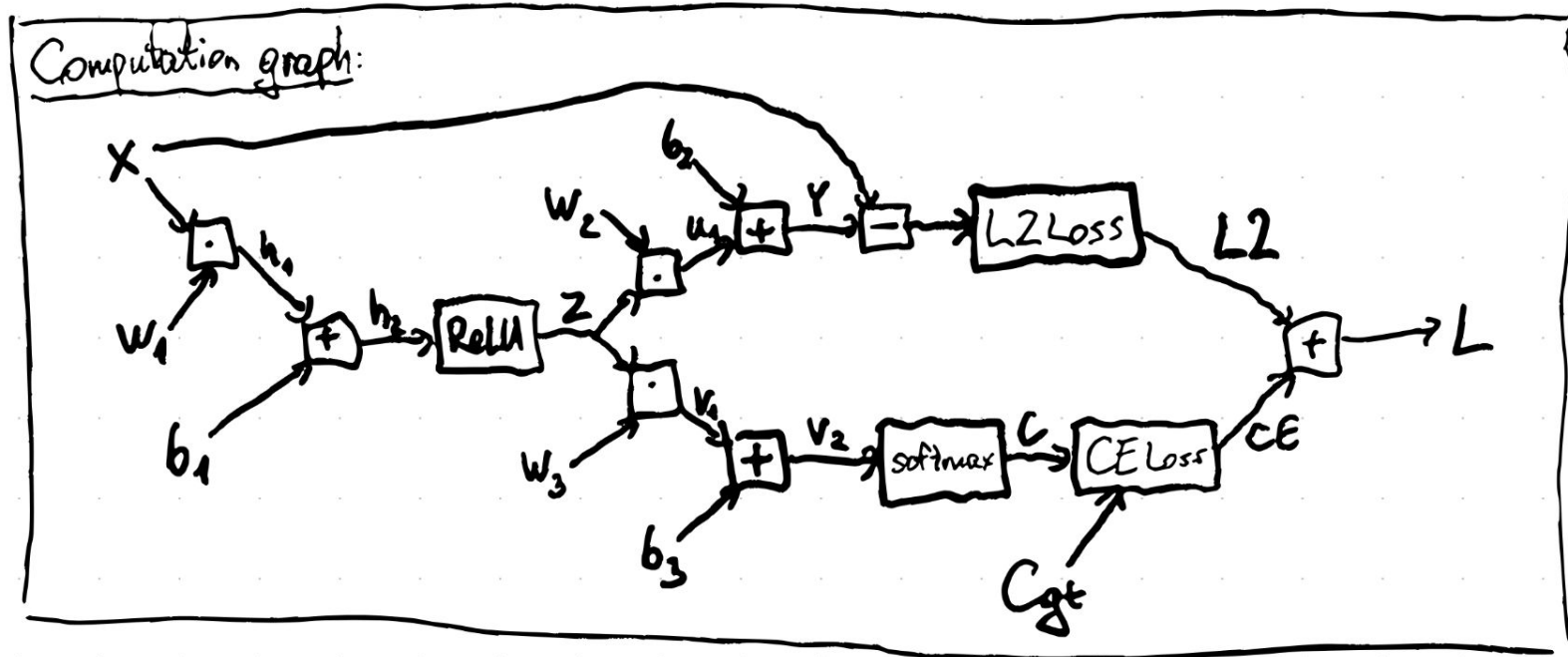
$$CE(C, C_{gt})$$

$$L = L_2(X, Y) + CE(C, C_{gt})$$

Tasks: - draw a computation graph for the model with loss nodes

- find $\delta W_1, \delta b_1, \delta W_2, \delta b_2, \delta W_3, \delta b_3$, where $\delta A = \frac{\partial L}{\partial A}$ for some A

Q1 solution (1)



^ one of multiple possible

Q1 solution (2)

Answer:

(also expanding δZ and h_2)

$$\delta W_1 = \begin{cases} X(W_2(Y-X) + W_3(C-C_{gt})), & \text{if } XW_1 + b_1 > 0, \\ 0, & \text{otherwise} \end{cases}$$

$$\delta b_1 = \begin{cases} W_2(Y-X) + W_3(C-C_{gt}), & \text{if } XW_1 + b_1 > 0, \\ 0, & \text{otherwise} \end{cases}$$

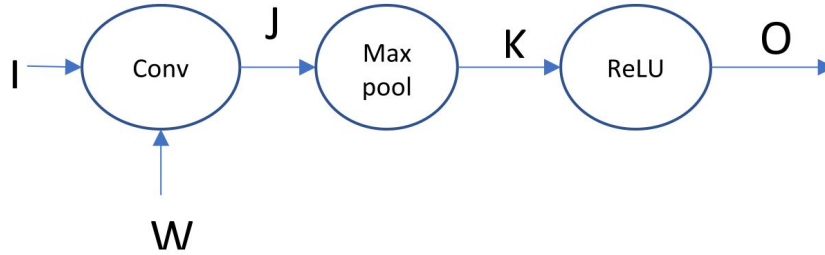
$$\delta W_2 = Z(Y-X)$$

$$\delta b_2 = Y-X$$

$$\delta W_3 = Z(C-C_{gt})$$

$$\delta b_3 = C-C_{gt}$$

Q3 (30). Compute the forward pass of the following model:



If I is the following image patch

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

and W is the following 3-by-3 filter matrix

1	2	1
0	0	0
-1	-2	-1

Compute J , K and O . Do not zero pad while doing the convolution. Assume stride size 1 for the convolution. Assume a 2-by-2 max pooling with stride 2. Write J , K and O below.

Q3 solution

Problem:

$$J = I * W$$

(no padding, stride=1)

$$K = \text{max pool}(J)$$

(2x2, stride=2)

$$O = \text{ReLU}(K)$$

$$I = \begin{bmatrix} 1 & 3 & 2 & 4 & 6 & 4 \\ 4 & 8 & 3 & 1 & 0 & 2 \\ 2 & 1 & 4 & 3 & 9 & 1 \\ 4 & 7 & 2 & 3 & 9 & 2 \end{bmatrix}$$

$$W = \begin{bmatrix} 1 & 2 & 1 \\ 0 & 0 & 0 \\ -1 & -2 & -1 \end{bmatrix}$$

Tasks:

- Write J, K, O

Answer:

$$J = \begin{bmatrix} 1 & -1 & -3 & -2 \\ 3 & 1 & -12 & -20 \end{bmatrix}$$

$$K = [3 \quad -2]$$

$$O = [3 \quad 0]$$



sasnausk@ualberta.ca