

Optimization

CISC 7026: Introduction to Deep Learning

University of Macau

Previous Quiz

Quiz 1 grades are on moodle (mean 2.75 / 4)

Previous Quiz

Quiz 1 grades are on moodle (mean 2.75 / 4)

When I have a function called g that maps some inputs $a \in A, b \in B, c \in C$ to outputs $d \in D, e \in E$ I would write

$$g : A \times B \times C \mapsto D \times E$$

or

$$g : A, B, C \mapsto D, E$$

Previous Quiz

Quiz 1 grades are on moodle (mean 2.75 / 4)

When I have a function called g that maps some inputs $a \in A, b \in B, c \in C$ to outputs $d \in D, e \in E$ I would write

$$g : A \times B \times C \mapsto D \times E$$

or

$$g : A, B, C \mapsto D, E$$

In my code, I would write

$$d, e = g(a, b, c)$$

Previous Quiz

\mathbb{R} is the set of real numbers

Previous Quiz

\mathbb{R} is the set of real numbers

$5, \pi, e, -3.2, \frac{3}{2}, 1.3333\dots$ are all real numbers

Previous Quiz

\mathbb{R} is the set of real numbers

$5, \pi, e, -3.2, \frac{3}{2}, 1.3333\dots$ are all real numbers

Unless the number is $\sqrt{-1}$, it is probably a real number

Previous Quiz

\mathbb{R} is the set of real numbers

$5, \pi, e, -3.2, \frac{3}{2}, 1.3333\dots$ are all real numbers

Unless the number is $\sqrt{-1}$, it is probably a real number

When I write \mathbb{R}^{d_x} , this corresponds to d_x real numbers

Previous Quiz

\mathbb{R} is the set of real numbers

$5, \pi, e, -3.2, \frac{3}{2}, 1.3333\dots$ are all real numbers

Unless the number is $\sqrt{-1}$, it is probably a real number

When I write \mathbb{R}^{d_x} , this corresponds to d_x real numbers

So $f : \mathbb{R}^{d_x} \mapsto \mathbb{R}^{d_y}$ is a function that maps d_x numbers to d_y numbers

Assignment one should be graded next week

Admin

Assignment one should be graded next week

Today's lecture will be very theoretical, but I must teach it

Admin

Assignment one should be graded next week

Today's lecture will be very theoretical, but I must teach it

I understand the homework might be difficult for some

Admin

Assignment one should be graded next week

Today's lecture will be very theoretical, but I must teach it

I understand the homework might be difficult for some

Assignment 2 is **optional** (it is **not** required)

Admin

Assignment one should be graded next week

Today's lecture will be very theoretical, but I must teach it

I understand the homework might be difficult for some

Assignment 2 is **optional** (it is **not** required)

I will give n assignments and you will receive the highest $n - 1$ grades

Admin

Assignment one should be graded next week

Today's lecture will be very theoretical, but I must teach it

I understand the homework might be difficult for some

Assignment 2 is **optional** (it is **not** required)

I will give n assignments and you will receive the highest $n - 1$ grades

E.g., $\text{asg1} = 60$, $\text{asg2} = 90$, $\text{asg3} = 70$, total score = $\frac{160}{200}$

Admin

Assignment one should be graded next week

Today's lecture will be very theoretical, but I must teach it

I understand the homework might be difficult for some

Assignment 2 is **optional** (it is **not** required)

I will give n assignments and you will receive the highest $n - 1$ grades

E.g., $\text{asg1} = 60$, $\text{asg2} = 90$, $\text{asg3} = 70$, total score = $\frac{160}{200}$

E.g., $\text{asg1} = 60$, $\text{asg2} = 0$, $\text{asg3} = 70$, total score = $\frac{130}{200}$

Agenda

1. Review
2. Quiz
3. Optimization
4. Calculus review
5. Deriving linear regression
6. Gradient descent
7. Backpropagation
8. Layer gradient
9. Full gradient
10. Autograd
11. Practical considerations

Agenda

1. **Review**
2. Quiz
3. Optimization
4. Calculus review
5. Deriving linear regression
6. Gradient descent
7. Backpropagation
8. Layer gradient
9. Full gradient
10. Autograd
11. Practical considerations

In lecture 1, we assumed a single-input system

In lecture 1, we assumed a single-input system

Years of education: $X \in \mathbb{R}$

In lecture 1, we assumed a single-input system

Years of education: $X \in \mathbb{R}$

But sometimes we want to consider multiple input dimensions

In lecture 1, we assumed a single-input system

Years of education: $X \in \mathbb{R}$

But sometimes we want to consider multiple input dimensions

Years of education, BMI, GDP: $X \in \mathbb{R}^3$

In lecture 1, we assumed a single-input system

Years of education: $X \in \mathbb{R}$

But sometimes we want to consider multiple input dimensions

Years of education, BMI, GDP: $X \in \mathbb{R}^3$

We can solve these problems using linear regression too

For multivariate problems, we will define the input dimension as d_x

For multivariate problems, we will define the input dimension as d_x

$$\mathbf{x} \in X; \quad X \in \mathbb{R}^{d_x}$$

For multivariate problems, we will define the input dimension as d_x

$$\mathbf{x} \in X; \quad X \in \mathbb{R}^{d_x}$$

We will write the vectors as

$$\mathbf{x}_{[i]} = \begin{bmatrix} x_{[i],1} \\ x_{[i],2} \\ \vdots \\ x_{[i],d_x} \end{bmatrix}$$

The design matrix for a **multivariate** linear system is

$$\mathbf{X}_D = \begin{bmatrix} x_{[1],d_x} & x_{[1],d_x-1} & \cdots & x_{[1],1} & 1 \\ x_{[2],d_x} & x_{[2],d_x-1} & \cdots & x_{[2],1} & 1 \\ \vdots & \vdots & \ddots & \vdots & \\ x_{[n],d_x} & x_{[n],d_x-1} & \cdots & x_{[n],1} & 1 \end{bmatrix}$$

The design matrix for a **multivariate** linear system is

$$\mathbf{X}_D = \begin{bmatrix} x_{[1],d_x} & x_{[1],d_x-1} & \cdots & x_{[1],1} & 1 \\ x_{[2],d_x} & x_{[2],d_x-1} & \cdots & x_{[2],1} & 1 \\ \vdots & \vdots & \ddots & \vdots & \\ x_{[n],d_x} & x_{[n],d_x-1} & \cdots & x_{[n],1} & 1 \end{bmatrix}$$

The solution is the same as before

$$\boldsymbol{\theta} = (\mathbf{X}_D^\top \mathbf{X}_D)^{-1} \mathbf{X}_D^\top \mathbf{y}$$

Limitations of Linear Regression

We combined **polynomial** and **multivariate** design matrices:

Limitations of Linear Regression

We combined **polynomial** and **multivariate** design matrices:

One-dimensional polynomial
functions

$$\mathbf{X}_D = \begin{bmatrix} x_{[1]}^m & x_{[1]}^{m-1} & \dots & x_{[1]} & 1 \\ x_{[2]}^m & x_{[2]}^{m-1} & \dots & x_{[2]} & 1 \\ \vdots & \vdots & \ddots & & \\ x_{[n]}^m & x_{[n]}^{m-1} & \dots & x_{[n]} & 1 \end{bmatrix}$$

Limitations of Linear Regression

We combined **polynomial** and **multivariate** design matrices:

One-dimensional polynomial
functions

$$\mathbf{X}_D = \begin{bmatrix} x_{[1]}^m & x_{[1]}^{m-1} & \dots & x_{[1]} & 1 \\ x_{[2]}^m & x_{[2]}^{m-1} & \dots & x_{[2]} & 1 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ x_{[n]}^m & x_{[n]}^{m-1} & \dots & x_{[n]} & 1 \end{bmatrix}$$

Multi-dimensional linear functions

$$\mathbf{X}_D = \begin{bmatrix} x_{[1],d_x} & x_{[1],d_x-1} & \dots & 1 \\ x_{[2],d_x} & x_{[2],d_x-1} & \dots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ x_{[n],d_x} & x_{[n],d_x-1} & \dots & 1 \end{bmatrix}$$

$$\mathbf{X}_D = [\mathbf{x}_{D,[1]} \quad \cdots \quad \mathbf{x}_{D,[n]}]^\top$$

$$\mathbf{X}_D = [\mathbf{x}_{D,[1]} \quad \cdots \quad \mathbf{x}_{D,[n]}]^\top$$

$$\mathbf{x}_{D,[i]} =$$

$$\left[\underbrace{x_{[i],d_x}^m x_{[i],d_x-1}^m \cdots x_{[i],1}^m}_{(d_x \Rightarrow 1, x^m)} \quad \underbrace{x_{[i],d_x}^m x_{[i],d_x-1}^m \cdots x_{[i],2}^m}_{(d_x \Rightarrow 2, x^m)} \quad \cdots \quad \underbrace{x_{[i],d_x}^{m-1} x_{[i],d_x-1}^{m-1} \cdots x_{[i],1}^m}_{(d_x \Rightarrow 1, x^{m-1})} \quad \cdots \right]$$

$$\mathbf{X}_D = [\mathbf{x}_{D,[1]} \quad \cdots \quad \mathbf{x}_{D,[n]}]^\top$$

$$\mathbf{x}_{D,[i]} =$$

$$\left[\underbrace{x_{[i],d_x}^m x_{[i],d_x-1}^m \cdots x_{[i],1}^m}_{(d_x \Rightarrow 1, x^m)} \quad \underbrace{x_{[i],d_x}^m x_{[i],d_x-1}^m \cdots x_{[i],2}^m}_{(d_x \Rightarrow 2, x^m)} \quad \cdots \quad \underbrace{x_{[i],d_x}^{m-1} x_{[i],d_x-1}^{m-1} \cdots x_{[i],1}^m}_{(d_x \Rightarrow 1, x^{m-1})} \quad \cdots \right]$$

The resulting design matrix is too large to solve

$$\mathbf{X}_D = [\mathbf{x}_{D,[1]} \quad \cdots \quad \mathbf{x}_{D,[n]}]^\top$$

$$\mathbf{x}_{D,[i]} =$$

$$\left[\underbrace{x_{[i],d_x}^m x_{[i],d_x-1}^m \cdots x_{[i],1}^m}_{(d_x \Rightarrow 1, x^m)} \quad \underbrace{x_{[i],d_x}^m x_{[i],d_x-1}^m \cdots x_{[i],2}^m}_{(d_x \Rightarrow 2, x^m)} \quad \cdots \quad \underbrace{x_{[i],d_x}^{m-1} x_{[i],d_x-1}^{m-1} \cdots x_{[i],1}^m}_{(d_x \Rightarrow 1, x^{m-1})} \quad \cdots \right]$$

The resulting design matrix is too large to solve

We introduced neural networks because they scale to larger problems

Brains and neural networks rely on **neurons**

Brains and neural networks rely on **neurons**

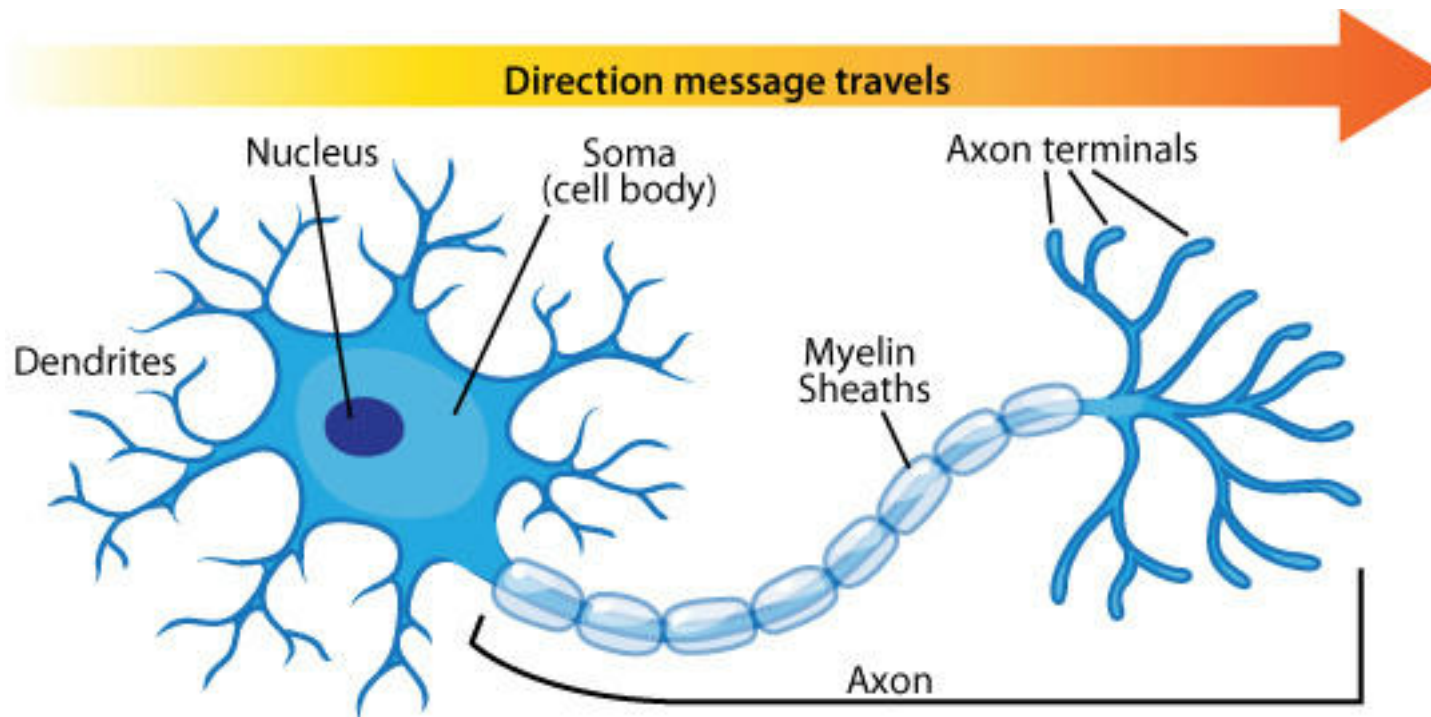
Brain: Biological neurons → Biological neural network

Brains and neural networks rely on **neurons**

Brain: Biological neurons \rightarrow Biological neural network

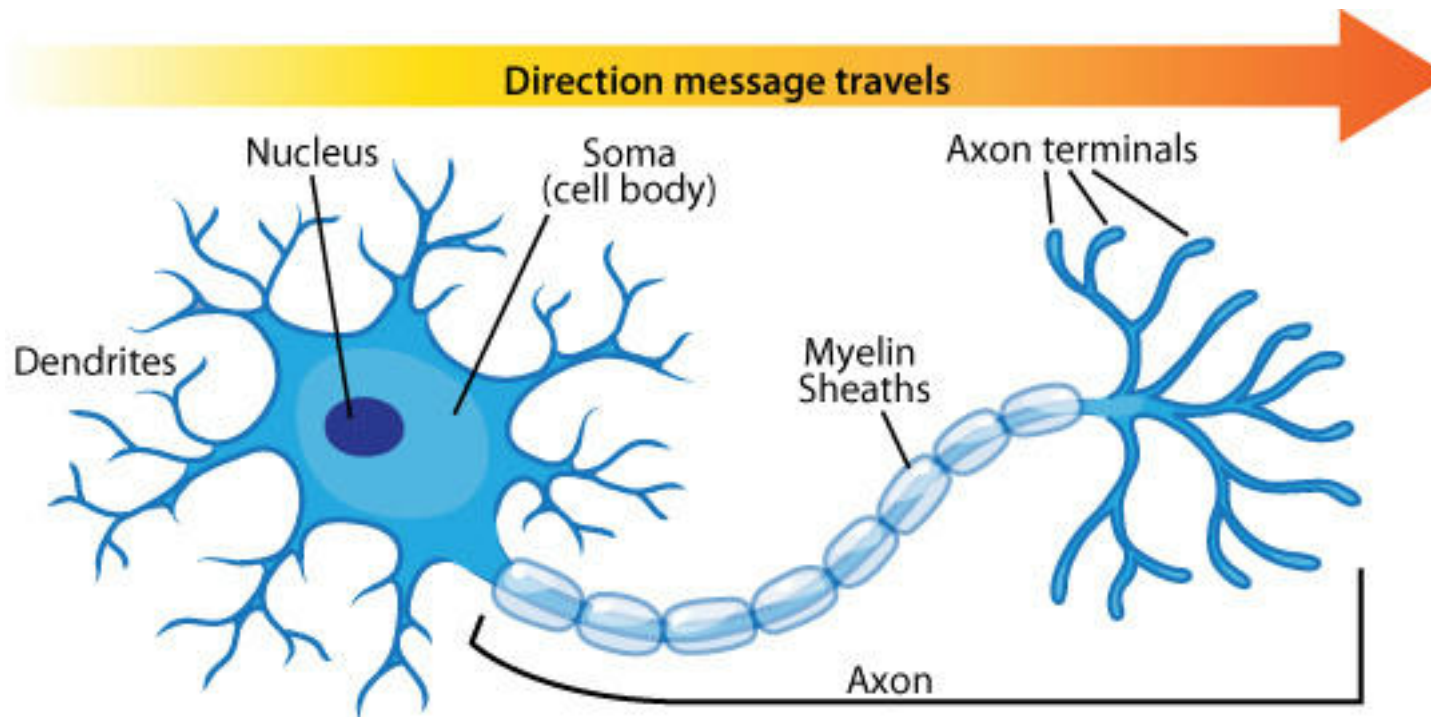
Computer: Artificial neurons \rightarrow Artificial neural network

Review



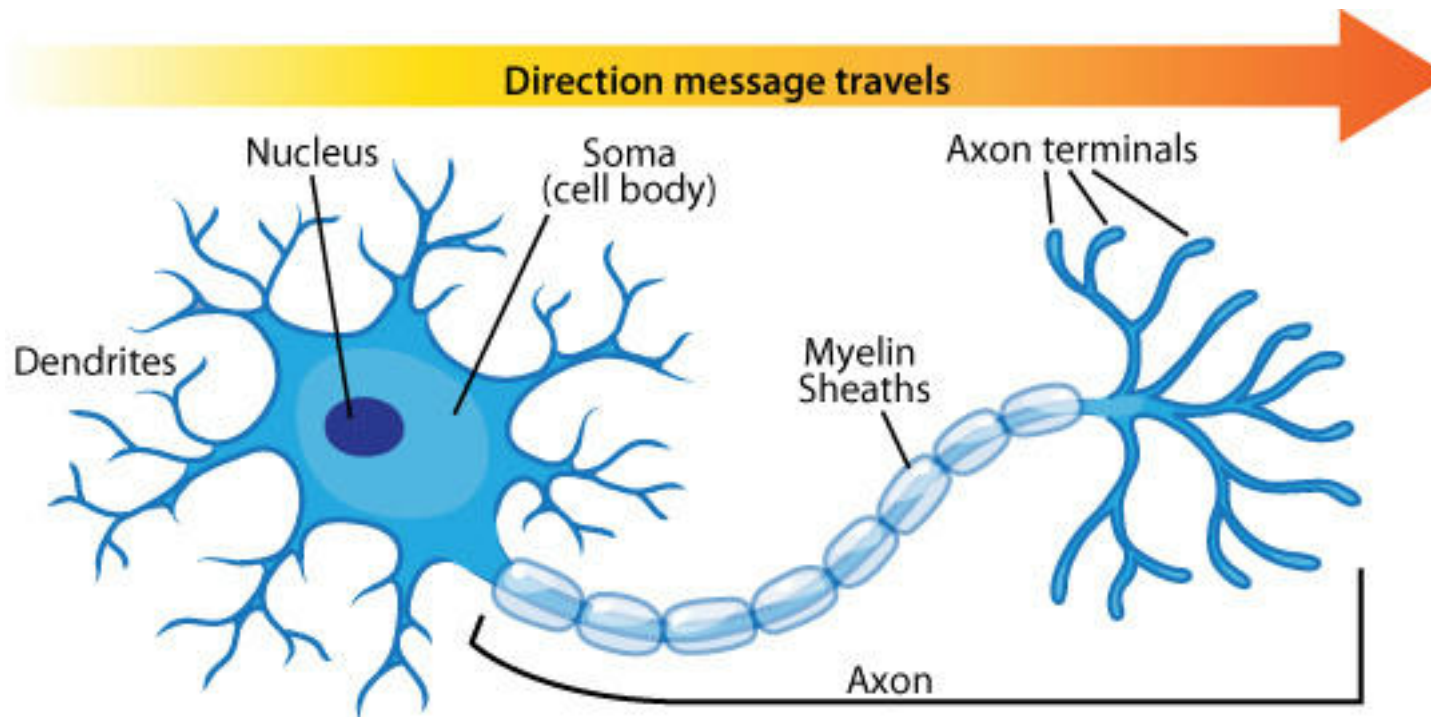
Neurons send messages based on messages received from other neurons

Review



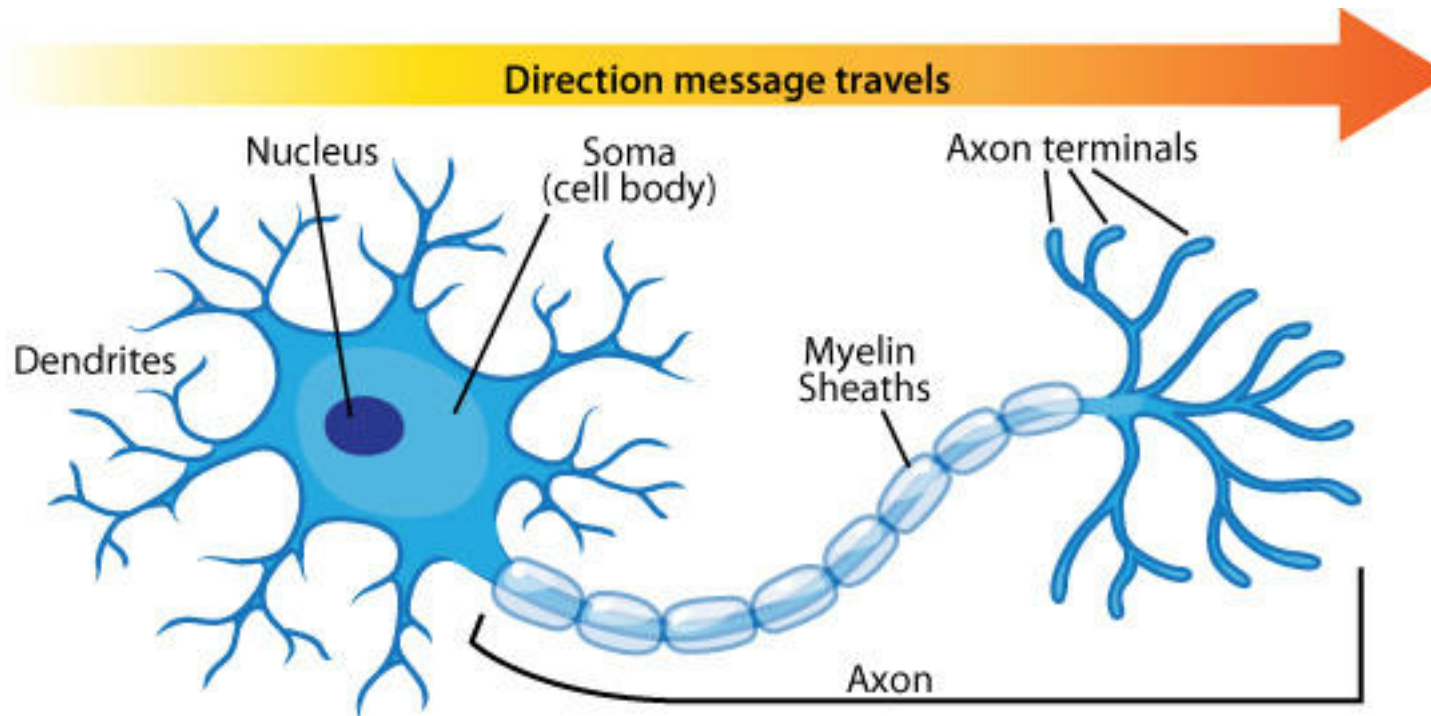
Incoming electrical signals travel along dendrites

Review



Electrical charges collect in the Soma (cell body)

Review



The axon outputs an electrical signal to other neurons

Biological Neurons

How does a neuron decide to send an impulse (“fire”)?

Biological Neurons

How does a neuron decide to send an impulse (“fire”)?

Dendrites form a parallel circuit

Biological Neurons

How does a neuron decide to send an impulse (“fire”)?

Dendrites form a parallel circuit

In a parallel circuit, we can sum voltages together

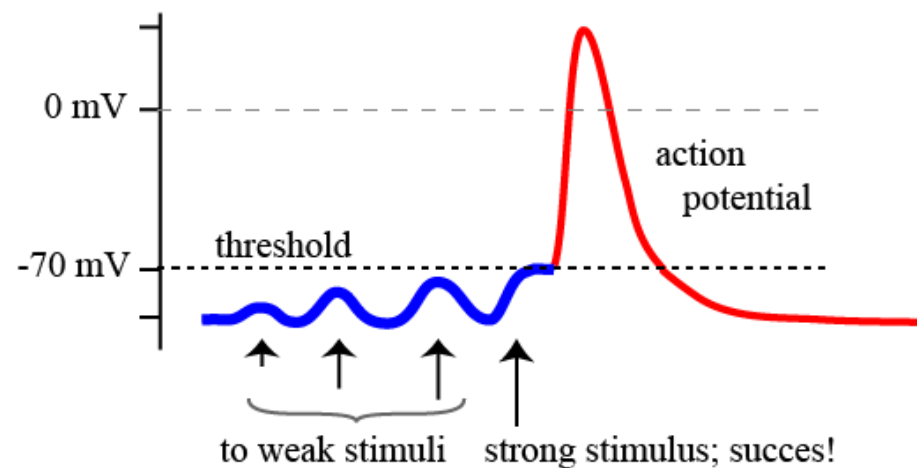
Biological Neurons

How does a neuron decide to send an impulse (“fire”)?

Dendrites form a parallel circuit

In a parallel circuit, we can sum voltages together

Incoming impulses (via dendrites)
change the electric potential of the
neuron



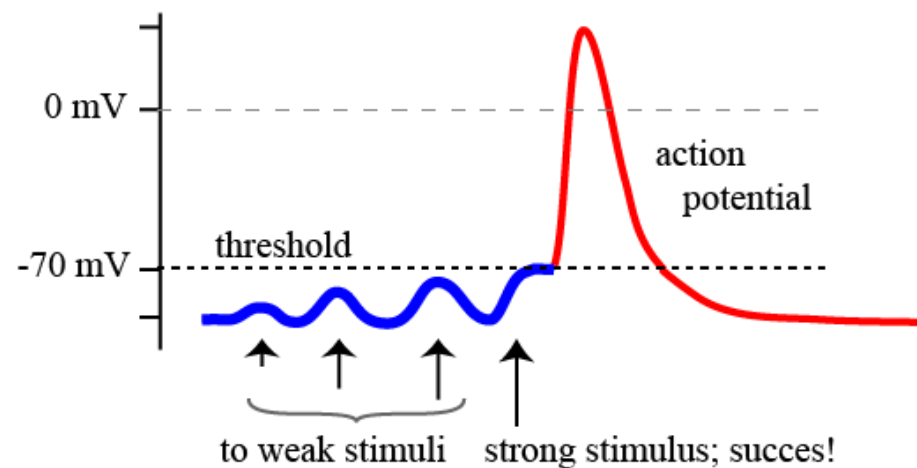
Biological Neurons

How does a neuron decide to send an impulse (“fire”)?

Dendrites form a parallel circuit

In a parallel circuit, we can sum voltages together

Incoming impulses (via dendrites)
change the electric potential of the
neuron



Many active dendrites will add together and trigger an impulse

Artificial Neurons

We model the neuron “firing” using an activation function σ

Artificial Neurons

We model the neuron “firing” using an activation function σ

Last time, we used the heaviside step function as the activation function

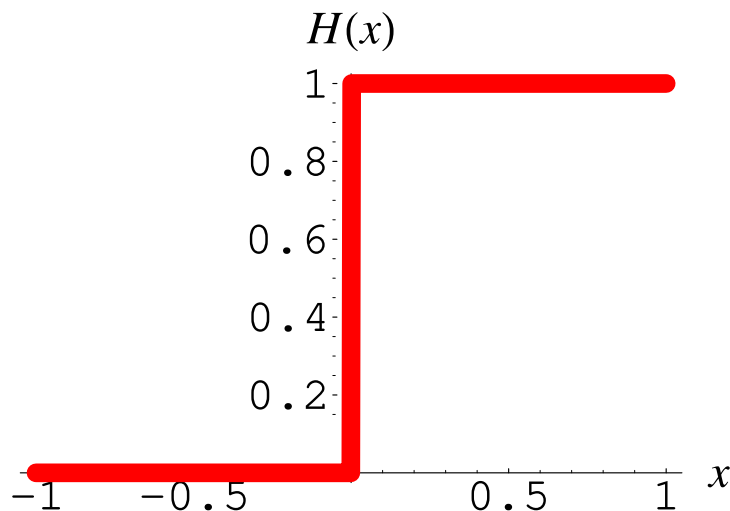
$$\sigma(x) = H(x) = \begin{cases} 0 & \text{if } x \leq 0 \\ 1 & \text{if } x > 0 \end{cases}$$

Artificial Neurons

We model the neuron “firing” using an activation function σ

Last time, we used the heaviside step function as the activation function

$$\sigma(x) = H(x) = \begin{cases} 0 & \text{if } x \leq 0 \\ 1 & \text{if } x > 0 \end{cases}$$



We modeled a neuron mathematically, creating an artificial neuron

We modeled a neuron mathematically, creating an artificial neuron

$$f(x, \theta) = \sigma(\theta^\top \bar{x}); \quad \bar{x} = \begin{bmatrix} 1 \\ x \end{bmatrix}$$

We modeled a neuron mathematically, creating an artificial neuron

$$f(\boldsymbol{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}^\top \bar{\boldsymbol{x}}); \quad \bar{\boldsymbol{x}} = \begin{bmatrix} 1 \\ \boldsymbol{x} \end{bmatrix}$$

$$f\left(\boldsymbol{x}, \begin{bmatrix} b \\ \boldsymbol{w} \end{bmatrix}\right) = \sigma(b + \boldsymbol{w}^\top \boldsymbol{x})$$

We modeled a neuron mathematically, creating an artificial neuron

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}^\top \bar{\mathbf{x}}); \quad \bar{\mathbf{x}} = \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix}$$

$$f\left(\mathbf{x}, \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix}\right) = \sigma(b + \mathbf{w}^\top \mathbf{x})$$

An artificial neuron is a linear model with an activation function σ

We modeled a neuron mathematically, creating an artificial neuron

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}^\top \bar{\mathbf{x}}); \quad \bar{\mathbf{x}} = \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix}$$

$$f\left(\mathbf{x}, \begin{bmatrix} b \\ \mathbf{w} \end{bmatrix}\right) = \sigma(b + \mathbf{w}^\top \mathbf{x})$$

An artificial neuron is a linear model with an activation function σ

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma \left(\underbrace{\theta_0 1 + \theta_1 x_1 + \dots + \theta_{d_x} x_{d_x}}_{\text{Linear model}} \right)$$

We can represent AND and OR boolean operators using a neuron

We can represent AND and OR boolean operators using a neuron

But a neuron cannot represent many other functions

We can represent AND and OR boolean operators using a neuron

But a neuron cannot represent many other functions

So we take many neurons and create a neural network

We can represent AND and OR boolean operators using a neuron

But a neuron cannot represent many other functions

So we take many neurons and create a neural network

We discussed **wide** neural networks and **deep** neural networks

Wide Neural Networks

How do we express a **wide** neural network mathematically?

Wide Neural Networks

How do we express a **wide** neural network mathematically?

A single neuron:

$$f : \mathbb{R}^{d_x} \times \Theta \mapsto \mathbb{R}$$

$$\Theta \in \mathbb{R}^{d_x+1}$$

Wide Neural Networks

How do we express a **wide** neural network mathematically?

A single neuron:

$$f : \mathbb{R}^{d_x} \times \Theta \mapsto \mathbb{R}$$

$$\Theta \in \mathbb{R}^{d_x+1}$$

d_y neurons (wide):

$$f : \mathbb{R}^{d_x} \times \Theta \mapsto \mathbb{R}^{d_y}$$

$$\Theta \in \mathbb{R}^{(d_x+1) \times d_y}$$

For a wide network (also called a layer):

$$f\left(\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{d_x} \end{bmatrix}, \begin{bmatrix} \theta_{0,1} & \theta_{0,2} & \cdots & \theta_{0,d_y} \\ \theta_{1,1} & \theta_{1,2} & \cdots & \theta_{1,d_y} \\ \vdots & \vdots & \ddots & \vdots \\ \theta_{d_x,1} & \theta_{d_x,2} & \cdots & \theta_{d_x,d_y} \end{bmatrix}\right) = \begin{bmatrix} \sigma\left(\sum_{i=0}^{d_x} \theta_{i,1} \bar{x}_i\right) \\ \sigma\left(\sum_{i=0}^{d_x} \theta_{i,2} \bar{x}_i\right) \\ \vdots \\ \sigma\left(\sum_{i=0}^{d_x} \theta_{i,d_y} \bar{x}_i\right) \end{bmatrix}$$

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}^\top \bar{\mathbf{x}}); \quad \boldsymbol{\theta} \in \mathbb{R}^{(d_x+1) \times d_y}$$

For a wide network (also called a layer):

$$f\left(\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{d_x} \end{bmatrix}, \begin{bmatrix} \theta_{0,1} & \theta_{0,2} & \cdots & \theta_{0,d_y} \\ \theta_{1,1} & \theta_{1,2} & \cdots & \theta_{1,d_y} \\ \vdots & \vdots & \ddots & \vdots \\ \theta_{d_x,1} & \theta_{d_x,2} & \cdots & \theta_{d_x,d_y} \end{bmatrix}\right) = \begin{bmatrix} \sigma\left(\sum_{i=0}^{d_x} \theta_{i,1} \bar{x}_i\right) \\ \sigma\left(\sum_{i=0}^{d_x} \theta_{i,2} \bar{x}_i\right) \\ \vdots \\ \sigma\left(\sum_{i=0}^{d_x} \theta_{i,d_y} \bar{x}_i\right) \end{bmatrix}$$

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}^\top \bar{\mathbf{x}}); \quad \boldsymbol{\theta} \in \mathbb{R}^{(d_x+1) \times d_y}$$

$$f\left(\mathbf{x}, \begin{bmatrix} \mathbf{b} \\ \mathbf{W} \end{bmatrix}\right) = \sigma(\mathbf{b} + \mathbf{W}^\top \mathbf{x}); \quad \mathbf{b} \in \mathbb{R}^{d_y}, \mathbf{W} \in \mathbb{R}^{d_x \times d_y}$$

Review

A **wide** neural network is also called a **layer**

Review

Review

A **wide** neural network is also called a **layer**

A layer is a linear operation and an activation function

$$f\left(x, \begin{bmatrix} \mathbf{b} \\ \mathbf{W} \end{bmatrix}\right) = \sigma(\mathbf{b} + \mathbf{W}^\top \mathbf{x})$$

Many layers makes a deep neural network

$$z_1 = f \left(x, \begin{bmatrix} b_1 \\ \mathbf{W}_1 \end{bmatrix} \right)$$

$$z_2 = f \left(z_1, \begin{bmatrix} b_2 \\ \mathbf{W}_2 \end{bmatrix} \right)$$

$$\mathbf{y} = f \left(z_2, \begin{bmatrix} b_2 \\ \mathbf{W}_2 \end{bmatrix} \right)$$

Agenda

1. Review
2. **Quiz**
3. Optimization
4. Calculus review
5. Deriving linear regression
6. Gradient descent
7. Backpropagation
8. Layer gradient
9. Full gradient
10. Autograd
11. Practical considerations

Quiz

Put away your phones and laptops

Quiz

Put away your phones and laptops

Take out a paper and pen, write your name and student ID

Quiz

Put away your phones and laptops

Take out a paper and pen, write your name and student ID

I will take away your quiz, give zero points, and refer you to the Dean if:

Quiz

Put away your phones and laptops

Take out a paper and pen, write your name and student ID

I will take away your quiz, give zero points, and refer you to the Dean if:

- Your phone or laptop is open during the quiz

Quiz

Put away your phones and laptops

Take out a paper and pen, write your name and student ID

I will take away your quiz, give zero points, and refer you to the Dean if:

- Your phone or laptop is open during the quiz
- You talk to your neighbor

Quiz

Put away your phones and laptops

Take out a paper and pen, write your name and student ID

I will take away your quiz, give zero points, and refer you to the Dean if:

- Your phone or laptop is open during the quiz
- You talk to your neighbor
- You look at your neighbor's quiz

Quiz

Put away your phones and laptops

Take out a paper and pen, write your name and student ID

I will take away your quiz, give zero points, and refer you to the Dean if:

- Your phone or laptop is open during the quiz
- You talk to your neighbor
- You look at your neighbor's quiz

After I explain the questions, you will have 15 minutes to finish the quiz

Agenda

1. Review
2. **Quiz**
3. Optimization
4. Calculus review
5. Deriving linear regression
6. Gradient descent
7. Backpropagation
8. Layer gradient
9. Full gradient
10. Autograd
11. Practical considerations

Agenda

1. Review
2. Quiz
3. **Optimization**
4. Calculus review
5. Deriving linear regression
6. Gradient descent
7. Backpropagation
8. Layer gradient
9. Full gradient
10. Autograd
11. Practical considerations

We previously introduced neural networks as universal function approximators

We previously introduced neural networks as universal function approximators

A deep neural network f can approximate **any** continuous function g to infinite precision

We previously introduced neural networks as universal function approximators

A deep neural network f can approximate **any** continuous function g to infinite precision

$$f(x, \theta) = g(x) + \varepsilon; \quad \varepsilon \rightarrow 0$$

We previously introduced neural networks as universal function approximators

A deep neural network f can approximate **any** continuous function g to infinite precision

$$f(x, \theta) = g(x) + \varepsilon; \quad \varepsilon \rightarrow 0$$

g can be a mapping from pictures to text, English to Chinese, etc

We previously introduced neural networks as universal function approximators

A deep neural network f can approximate **any** continuous function g to infinite precision

$$f(x, \theta) = g(x) + \varepsilon; \quad \varepsilon \rightarrow 0$$

g can be a mapping from pictures to text, English to Chinese, etc

But how do we find the θ that makes $f(x, \theta) = g(x) + \varepsilon$?

We previously introduced neural networks as universal function approximators

A deep neural network f can approximate **any** continuous function g to infinite precision

$$f(x, \theta) = g(x) + \varepsilon; \quad \varepsilon \rightarrow 0$$

g can be a mapping from pictures to text, English to Chinese, etc

But how do we find the θ that makes $f(x, \theta) = g(x) + \varepsilon$?

We said this θ exists, but never said how to find it

We previously introduced neural networks as universal function approximators

A deep neural network f can approximate **any** continuous function g to infinite precision

$$f(x, \theta) = g(x) + \varepsilon; \quad \varepsilon \rightarrow 0$$

g can be a mapping from pictures to text, English to Chinese, etc

But how do we find the θ that makes $f(x, \theta) = g(x) + \varepsilon$?

We said this θ exists, but never said how to find it

Goal: Find the parameters θ for a neural network

When we search for θ , we call it **optimization** or **training**

When we search for θ , we call it **optimization** or **training**

We optimize a loss function by computing

$$\arg \min_{\theta} \mathcal{L}(X, Y, \theta)$$

When we search for θ , we call it **optimization** or **training**

We optimize a loss function by computing

$$\arg \min_{\theta} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \theta)$$

This expression looks very simple, but it can be very hard to evaluate

When we search for θ , we call it **optimization** or **training**

We optimize a loss function by computing

$$\arg \min_{\theta} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \theta)$$

This expression looks very simple, but it can be very hard to evaluate

To start, let us consider how we find

$$\arg \min_{\theta} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \theta)$$

in linear regression

Recall how we solved the linear regression problem

Recall how we solved the linear regression problem

We define the square error loss function

Recall how we solved the linear regression problem

We define the square error loss function

$$\arg \min_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^n (f(x_i, \boldsymbol{\theta}) - y_i)^2$$

Recall how we solved the linear regression problem

We define the square error loss function

$$\arg \min_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^n (f(x_i, \boldsymbol{\theta}) - y_i)^2$$

Then, I gave you a magical solution to this optimization problem

Recall how we solved the linear regression problem

We define the square error loss function

$$\arg \min_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^n (f(x_i, \boldsymbol{\theta}) - y_i)^2$$

Then, I gave you a magical solution to this optimization problem

$$\boldsymbol{\theta} = (\mathbf{X}_D^\top \mathbf{X}_D)^{-1} \mathbf{X}_D^\top \mathbf{Y}$$

Recall how we solved the linear regression problem

We define the square error loss function

$$\arg \min_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^n (f(x_i, \boldsymbol{\theta}) - y_i)^2$$

Then, I gave you a magical solution to this optimization problem

$$\boldsymbol{\theta} = (\mathbf{X}_D^\top \mathbf{X}_D)^{-1} \mathbf{X}_D^\top \mathbf{Y}$$

Where does this solution come from? Can we do the same for neural networks?

The solution for linear regression and neural networks comes from
calculus

The solution for linear regression and neural networks comes from **calculus**

We will briefly review basic calculus concepts

Agenda

1. Review
2. Quiz
3. **Optimization**
4. Calculus review
5. Deriving linear regression
6. Gradient descent
7. Backpropagation
8. Layer gradient
9. Full gradient
10. Autograd
11. Practical considerations

Agenda

1. Review
2. Quiz
3. Optimization
4. **Calculus review**
5. Deriving linear regression
6. Gradient descent
7. Backpropagation
8. Layer gradient
9. Full gradient
10. Autograd
11. Practical considerations

We write the **derivative** of a function f with respect to an input x as

$$f'(x) = \frac{d}{dx} f = \frac{df}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

We write the **derivative** of a function f with respect to an input x as

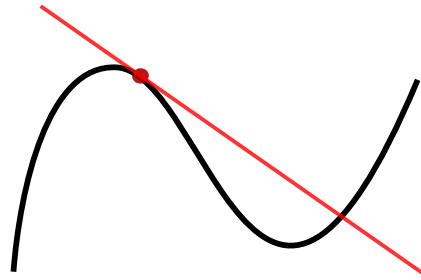
$$f'(x) = \frac{d}{dx} f = \frac{df}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

The derivative is the slope of a function

We write the **derivative** of a function f with respect to an input x as

$$f'(x) = \frac{d}{dx} f = \frac{df}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

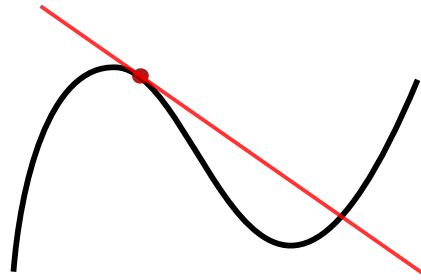
The derivative is the slope of a function



We write the **derivative** of a function f with respect to an input x as

$$f'(x) = \frac{d}{dx} f = \frac{df}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

The derivative is the slope of a function



$$f(x), f'(x = a)$$

It is easiest if you treat the derivative as a **function of functions**

It is easiest if you treat the derivative as a **function of functions**

$$\text{derivative}(f) = f'$$

It is easiest if you treat the derivative as a **function of functions**

$$\text{derivative}(f) = f'$$

$$\frac{d}{dx} : f \mapsto f'$$

It is easiest if you treat the derivative as a **function of functions**

$$\text{derivative}(f) = f'$$

$$\frac{d}{dx} : f \mapsto f'$$

The derivative takes a function f and outputs a new function f'

It is easiest if you treat the derivative as a **function of functions**

$$\text{derivative}(f) = f'$$

$$\frac{d}{dx} : f \mapsto f'$$

The derivative takes a function f and outputs a new function f'

$$\frac{d}{dx} : [f : X \mapsto Y] \mapsto [f' : X \mapsto Y]$$

There are formulas for computing the derivative of various operations

There are formulas for computing the derivative of various operations

Constant

$$\frac{d}{dx}c = 0$$

There are formulas for computing the derivative of various operations

Constant

$$\frac{d}{dx}c = 0$$

Power

$$\frac{d}{dx}x^n = nx^{n-1}$$

Sum/Difference

$$\frac{d}{dx}(f(x) + g(x)) = f'(x) + g'(x)$$

Sum/Difference

$$\frac{d}{dx}(f(x) + g(x)) = f'(x) + g'(x)$$

Product

$$\frac{d}{dx}[f(x)g(x)] = f'(x)g(x) + g'(x)f(x)$$

Sum/Difference

$$\frac{d}{dx}(f(x) + g(x)) = f'(x) + g'(x)$$

Product

$$\frac{d}{dx}[f(x)g(x)] = f'(x)g(x) + g'(x)f(x)$$

Chain

$$\frac{d}{dx}[f(g(x))] = f'(g(x))g'(x)$$

For example, consider the function

$$f(x) = x^2 - 3x$$

For example, consider the function

$$f(x) = x^2 - 3x$$

We can write the derivative as

$$\frac{d}{dx}(f(x)) = 2x - 3$$

For example, consider the function

$$f(x) = x^2 - 3x$$

We can write the derivative as

$$\frac{d}{dx}(f(x)) = 2x - 3$$

We can evaluate the derivative at specific points

For example, consider the function

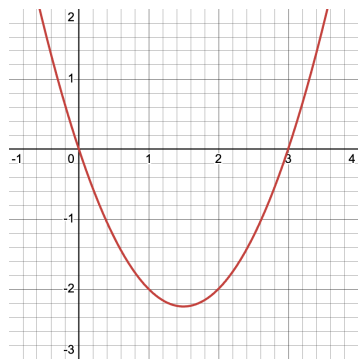
$$f(x) = x^2 - 3x$$

We can write the derivative as

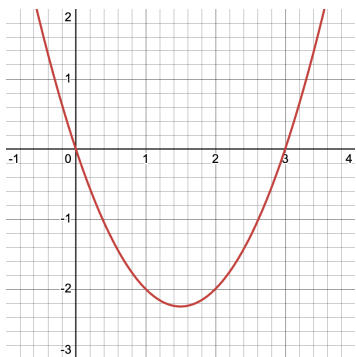
$$\frac{d}{dx}(f(x)) = 2x - 3$$

We can evaluate the derivative at specific points

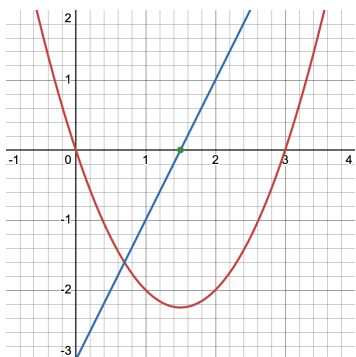
$$\left(\frac{d}{dx}f\right)(1) = 2 \cdot 1 - 3 = -1$$



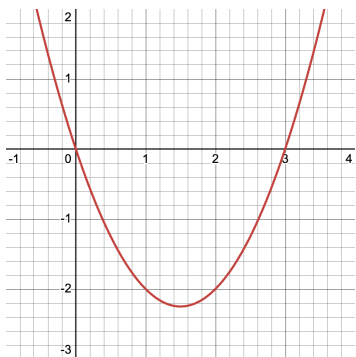
$$f(x) = x^2 - 3x$$



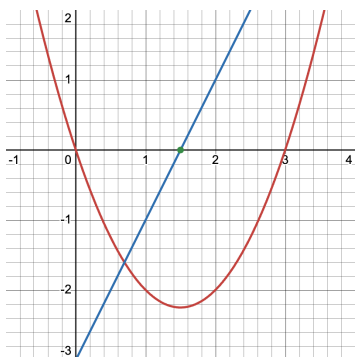
$$f(x) = x^2 - 3x$$



$$\frac{df}{dx} = 2x - 3$$



$$f(x) = x^2 - 3x$$



$$\frac{df}{dx} = 2x - 3$$

$$0 = 2x - 3 \quad \Rightarrow \quad x = \frac{3}{2}$$

We can expand the definition of derivative to multivariate functions. We call this the **gradient**

$$\nabla_{\mathbf{x}} f\left([x_1 \ x_2 \ \dots \ x_n]^\top\right) = \left[\frac{\partial f}{\partial x_1} \ \frac{\partial f}{\partial x_2} \ \dots \ \frac{\partial f}{\partial x_n}\right]^\top$$

We can expand the definition of derivative to multivariate functions. We call this the **gradient**

$$\nabla_{\mathbf{x}} f\left([x_1 \ x_2 \ \dots \ x_n]^\top\right) = \left[\frac{\partial f}{\partial x_1} \ \frac{\partial f}{\partial x_2} \ \dots \ \frac{\partial f}{\partial x_n}\right]^\top$$

For our purposes, we treat partial derivatives like standard derivatives

We can expand the definition of derivative to multivariate functions. We call this the **gradient**

$$\nabla_{\mathbf{x}} f\left([x_1 \ x_2 \ \dots \ x_n]^\top\right) = \left[\frac{\partial f}{\partial x_1} \ \frac{\partial f}{\partial x_2} \ \dots \ \frac{\partial f}{\partial x_n}\right]^\top$$

For our purposes, we treat partial derivatives like standard derivatives

$$\frac{\partial}{\partial x_1} f(x_1, \dots, x_n) = \frac{d}{dx_1} f(x_1, \dots, x_n)$$

We can expand the definition of derivative to multivariate functions. We call this the **gradient**

$$\nabla_{\mathbf{x}} f([x_1 \ x_2 \ \dots \ x_n]^\top) = \left[\frac{\partial f}{\partial x_1} \ \frac{\partial f}{\partial x_2} \ \dots \ \frac{\partial f}{\partial x_n} \right]^\top$$

For our purposes, we treat partial derivatives like standard derivatives

$$\frac{\partial}{\partial x_1} f(x_1, \dots, x_n) = \frac{d}{dx_1} f(x_1, \dots, x_n)$$

When computing $\frac{\partial}{\partial x_i} f(x_1, \dots, x_n)$, we treat $x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n$ as constant

For example, consider the function

$$f(x_1, x_2) = x_1^2 - 3x_1x_2$$

For example, consider the function

$$f(x_1, x_2) = x_1^2 - 3x_1x_2$$

We can write the gradient as

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \nabla_{x_1, x_2} f\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} \frac{\partial}{\partial x_1} f(x_1, x_2) \\ \frac{\partial}{\partial x_2} f(x_1, x_2) \end{bmatrix} = \begin{bmatrix} 2x_1 - 3x_2 \\ -3x_1 \end{bmatrix}$$

For example, consider the function

$$f(x_1, x_2) = x_1^2 - 3x_1x_2$$

We can write the gradient as

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \nabla_{x_1, x_2} f\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} \frac{\partial}{\partial x_1} f(x_1, x_2) \\ \frac{\partial}{\partial x_2} f(x_1, x_2) \end{bmatrix} = \begin{bmatrix} 2x_1 - 3x_2 \\ -3x_1 \end{bmatrix}$$

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \nabla_{x_1, x_2} f\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} \frac{\partial}{\partial x_1} f(x_1, x_2) \\ \frac{\partial}{\partial x_2} f(x_1, x_2) \end{bmatrix} = \begin{bmatrix} 2x_1 - 3x_2 \\ -3x_1 \end{bmatrix}$$

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \nabla_{x_1, x_2} f\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} \frac{\partial}{\partial x_1} f(x_1, x_2) \\ \frac{\partial}{\partial x_2} f(x_1, x_2) \end{bmatrix} = \begin{bmatrix} 2x_1 - 3x_2 \\ -3x_1 \end{bmatrix}$$

We can evaluate the gradient at specific points

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = \nabla_{x_1, x_2} f\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} \frac{\partial}{\partial x_1} f(x_1, x_2) \\ \frac{\partial}{\partial x_2} f(x_1, x_2) \end{bmatrix} = \begin{bmatrix} 2x_1 - 3x_2 \\ -3x_1 \end{bmatrix}$$

We can evaluate the gradient at specific points

$$\nabla_{\mathbf{x}} f\left(\begin{bmatrix} 1 \\ 0 \end{bmatrix}\right) = \nabla_{x_1, x_2} f\left(\begin{bmatrix} 1 \\ 0 \end{bmatrix}\right) = \begin{bmatrix} \frac{\partial}{\partial x_1} f(1, 0) \\ \frac{\partial}{\partial x_2} f(1, 0) \end{bmatrix} = \begin{bmatrix} 2 \cdot 1 - 3 \cdot 0 \\ -3 \cdot 1 \end{bmatrix} = \begin{bmatrix} -1 \\ -3 \end{bmatrix}$$

In calculus, we can find the local extrema of a function $f(x)$ by finding where the derivative is zero

$$f'(x) = \frac{d}{dx}f(x) = 0$$

In calculus, we can find the local extrema of a function $f(x)$ by finding where the derivative is zero

$$f'(x) = \frac{d}{dx} f(x) = 0$$

With a multivariate function, the extrema lies where the gradient is zero

$$\nabla_x f(x) = \left[\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \quad \cdots \quad \frac{\partial f}{\partial x_n} \right]^\top = [0 \quad 0 \quad \cdots \quad 0]^\top$$

Agenda

1. Review
2. Quiz
3. Optimization
4. **Calculus review**
5. Deriving linear regression
6. Gradient descent
7. Backpropagation
8. Layer gradient
9. Full gradient
10. Autograd
11. Practical considerations

Agenda

1. Review
2. Quiz
3. Optimization
4. Calculus review
5. **Deriving linear regression**
6. Gradient descent
7. Backpropagation
8. Layer gradient
9. Full gradient
10. Autograd
11. Practical considerations

Now that we remember calculus, let us revisit linear regression

Now that we remember calculus, let us revisit linear regression

If we can derive the solution for linear regression, maybe we can apply it to deep neural networks

In linear regression, our loss function is

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \sum_{i=1}^n \left(f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) - \mathbf{y}_{[i]} \right)^2$$

In linear regression, our loss function is

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \sum_{i=1}^n \left(f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) - \mathbf{y}_{[i]} \right)^2$$

We can write the square error loss in matrix form as

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = (\mathbf{Y} - \mathbf{X}_D \boldsymbol{\theta})^\top (\mathbf{Y} - \mathbf{X}_D \boldsymbol{\theta})$$

In linear regression, our loss function is

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \sum_{i=1}^n \left(f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) - \mathbf{y}_{[i]} \right)^2$$

We can write the square error loss in matrix form as

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = (\mathbf{Y} - \mathbf{X}_D \boldsymbol{\theta})^\top (\mathbf{Y} - \mathbf{X}_D \boldsymbol{\theta})$$

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \underbrace{(\mathbf{Y} - \mathbf{X}_D \boldsymbol{\theta})^\top}_{\text{Linear function of } \boldsymbol{\theta}} \underbrace{(\mathbf{Y} - \mathbf{X}_D \boldsymbol{\theta})}_{\text{Linear function of } \boldsymbol{\theta}} \\ \underbrace{\hspace{10em}}_{\text{Quadratic function of } \boldsymbol{\theta}}$$

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \underbrace{(\mathbf{Y} - \mathbf{X}_D \boldsymbol{\theta})^\top}_{\text{Linear function of } \boldsymbol{\theta}} \underbrace{(\mathbf{Y} - \mathbf{X}_D \boldsymbol{\theta})}_{\text{Linear function of } \boldsymbol{\theta}}$$

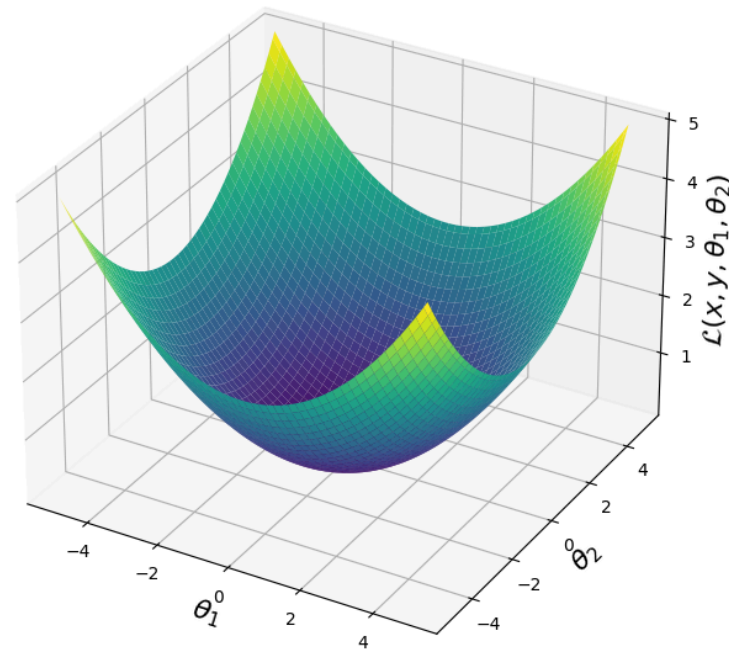
Quadratic function of $\boldsymbol{\theta}$

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \underbrace{(\mathbf{Y} - \mathbf{X}_D \boldsymbol{\theta})^\top}_{\text{Linear function of } \boldsymbol{\theta}} \underbrace{(\mathbf{Y} - \mathbf{X}_D \boldsymbol{\theta})}_{\text{Linear function of } \boldsymbol{\theta}}$$

Quadratic function of $\boldsymbol{\theta}$

A quadratic function has a single minima! The minima must be at

$$\nabla_{\boldsymbol{\theta}} \mathcal{L} = 0$$



Therefore, we know that the θ that solves

$$\nabla_{\theta} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \theta) = 0$$

Therefore, we know that the θ that solves

$$\nabla_{\theta} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \theta) = 0$$

Also solves

$$\arg \min_{\theta} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \theta)$$

Using calculus, let us derive the solution to linear regression

Using calculus, let us derive the solution to linear regression

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = (\mathbf{Y} - \mathbf{X}_D \boldsymbol{\theta})^\top (\mathbf{Y} - \mathbf{X}_D \boldsymbol{\theta})$$

Using calculus, let us derive the solution to linear regression

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = (\mathbf{Y} - \mathbf{X}_D \boldsymbol{\theta})^\top (\mathbf{Y} - \mathbf{X}_D \boldsymbol{\theta})$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \left[(\mathbf{Y} - \mathbf{X}_D \boldsymbol{\theta})^\top (\mathbf{Y} - \mathbf{X}_D \boldsymbol{\theta}) \right]$$

Using calculus, let us derive the solution to linear regression

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = (\mathbf{Y} - \mathbf{X}_D \boldsymbol{\theta})^\top (\mathbf{Y} - \mathbf{X}_D \boldsymbol{\theta})$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} [(\mathbf{Y} - \mathbf{X}_D \boldsymbol{\theta})^\top (\mathbf{Y} - \mathbf{X}_D \boldsymbol{\theta})]$$

$$= \nabla_{\boldsymbol{\theta}} [\mathbf{Y}^\top \mathbf{Y} - \mathbf{Y}^\top \mathbf{X}_D \boldsymbol{\theta} - (\mathbf{X}_D \boldsymbol{\theta})^\top \mathbf{Y} + (\mathbf{X}_D \boldsymbol{\theta})^\top \mathbf{X}_D \boldsymbol{\theta}]$$

$$= \mathbf{0} - \mathbf{Y}^\top \mathbf{X}_D \mathbf{I} - (\mathbf{X}_D \mathbf{I})^\top \mathbf{Y} + (\mathbf{X}_D \mathbf{I})^\top \mathbf{X}_D \boldsymbol{\theta} + (\mathbf{X}_D \boldsymbol{\theta})^\top \mathbf{X}_D \mathbf{I}$$

$$= \mathbf{0} - \mathbf{Y}^\top \mathbf{X}_D \mathbf{I} - (\mathbf{X}_D \mathbf{I})^\top \mathbf{Y} + (\mathbf{X}_D \mathbf{I})^\top \mathbf{X}_D \boldsymbol{\theta} + (\mathbf{X}_D \boldsymbol{\theta})^\top \mathbf{X}_D \mathbf{I}$$

$$\begin{aligned} &= \mathbf{0} - \mathbf{Y}^\top \mathbf{X}_D \mathbf{I} - (\mathbf{X}_D \mathbf{I})^\top \mathbf{Y} + (\mathbf{X}_D \mathbf{I})^\top \mathbf{X}_D \boldsymbol{\theta} + (\mathbf{X}_D \boldsymbol{\theta})^\top \mathbf{X}_D \mathbf{I} \\ &= -\mathbf{Y}^\top \mathbf{X}_D - \mathbf{X}_D^\top \mathbf{Y} + \mathbf{X}_D^\top \mathbf{X}_D \boldsymbol{\theta} + (\mathbf{X}_D \boldsymbol{\theta})^\top \mathbf{X}_D \end{aligned}$$

$$\begin{aligned}
&= \mathbf{0} - \mathbf{Y}^\top \mathbf{X}_D \mathbf{I} - (\mathbf{X}_D \mathbf{I})^\top \mathbf{Y} + (\mathbf{X}_D \mathbf{I})^\top \mathbf{X}_D \boldsymbol{\theta} + (\mathbf{X}_D \boldsymbol{\theta})^\top \mathbf{X}_D \mathbf{I} \\
&= -\mathbf{Y}^\top \mathbf{X}_D - \mathbf{X}_D^\top \mathbf{Y} + \mathbf{X}_D^\top \mathbf{X}_D \boldsymbol{\theta} + (\mathbf{X}_D \boldsymbol{\theta})^\top \mathbf{X}_D
\end{aligned}$$

Remember, $(\mathbf{AB})^\top = \mathbf{B}^\top \mathbf{A}^\top$, and so $\mathbf{Y}^\top \mathbf{X}_D = \mathbf{X}_D^\top \mathbf{Y}$

$$\begin{aligned}
&= \mathbf{0} - \mathbf{Y}^\top \mathbf{X}_D \mathbf{I} - (\mathbf{X}_D \mathbf{I})^\top \mathbf{Y} + (\mathbf{X}_D \mathbf{I})^\top \mathbf{X}_D \boldsymbol{\theta} + (\mathbf{X}_D \boldsymbol{\theta})^\top \mathbf{X}_D \mathbf{I} \\
&= -\mathbf{Y}^\top \mathbf{X}_D - \mathbf{X}_D^\top \mathbf{Y} + \mathbf{X}_D^\top \mathbf{X}_D \boldsymbol{\theta} + (\mathbf{X}_D \boldsymbol{\theta})^\top \mathbf{X}_D
\end{aligned}$$

Remember, $(\mathbf{AB})^\top = \mathbf{B}^\top \mathbf{A}^\top$, and so $\mathbf{Y}^\top \mathbf{X}_D = \mathbf{X}_D^\top \mathbf{Y}$

$$= -\mathbf{Y}^\top \mathbf{X}_D - \mathbf{Y}^\top \mathbf{X}_D + \mathbf{X}_D^\top \mathbf{X}_D \boldsymbol{\theta} + \mathbf{X}_D^\top \mathbf{X}_D \boldsymbol{\theta}$$

$$\begin{aligned}
&= \mathbf{0} - \mathbf{Y}^\top \mathbf{X}_D \mathbf{I} - (\mathbf{X}_D \mathbf{I})^\top \mathbf{Y} + (\mathbf{X}_D \mathbf{I})^\top \mathbf{X}_D \boldsymbol{\theta} + (\mathbf{X}_D \boldsymbol{\theta})^\top \mathbf{X}_D \mathbf{I} \\
&= -\mathbf{Y}^\top \mathbf{X}_D - \mathbf{X}_D^\top \mathbf{Y} + \mathbf{X}_D^\top \mathbf{X}_D \boldsymbol{\theta} + (\mathbf{X}_D \boldsymbol{\theta})^\top \mathbf{X}_D
\end{aligned}$$

Remember, $(\mathbf{AB})^\top = \mathbf{B}^\top \mathbf{A}^\top$, and so $\mathbf{Y}^\top \mathbf{X}_D = \mathbf{X}_D^\top \mathbf{Y}$

$$\begin{aligned}
&= -\mathbf{Y}^\top \mathbf{X}_D - \mathbf{Y}^\top \mathbf{X}_D + \mathbf{X}_D^\top \mathbf{X}_D \boldsymbol{\theta} + \mathbf{X}_D^\top \mathbf{X}_D \boldsymbol{\theta} \\
&= -2\mathbf{X}_D^\top \mathbf{Y} + 2\mathbf{X}_D^\top \mathbf{X}_D \boldsymbol{\theta}
\end{aligned}$$

$$\begin{aligned}
&= \mathbf{0} - \mathbf{Y}^\top \mathbf{X}_D \mathbf{I} - (\mathbf{X}_D \mathbf{I})^\top \mathbf{Y} + (\mathbf{X}_D \mathbf{I})^\top \mathbf{X}_D \boldsymbol{\theta} + (\mathbf{X}_D \boldsymbol{\theta})^\top \mathbf{X}_D \mathbf{I} \\
&= -\mathbf{Y}^\top \mathbf{X}_D - \mathbf{X}_D^\top \mathbf{Y} + \mathbf{X}_D^\top \mathbf{X}_D \boldsymbol{\theta} + (\mathbf{X}_D \boldsymbol{\theta})^\top \mathbf{X}_D
\end{aligned}$$

Remember, $(\mathbf{AB})^\top = \mathbf{B}^\top \mathbf{A}^\top$, and so $\mathbf{Y}^\top \mathbf{X}_D = \mathbf{X}_D^\top \mathbf{Y}$

$$\begin{aligned}
&= -\mathbf{Y}^\top \mathbf{X}_D - \mathbf{Y}^\top \mathbf{X}_D + \mathbf{X}_D^\top \mathbf{X}_D \boldsymbol{\theta} + \mathbf{X}_D^\top \mathbf{X}_D \boldsymbol{\theta} \\
&= -2\mathbf{X}_D^\top \mathbf{Y} + 2\mathbf{X}_D^\top \mathbf{X}_D \boldsymbol{\theta}
\end{aligned}$$

And so, the gradient of the loss is

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = -2\mathbf{X}_D^\top \mathbf{Y} + 2\mathbf{X}_D^\top \mathbf{X}_D \boldsymbol{\theta}$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{X}, \boldsymbol{Y}, \boldsymbol{\theta}) = -2\boldsymbol{X}_D^{\top} \boldsymbol{Y} + 2\boldsymbol{X}_D^{\top} \boldsymbol{X}_D \boldsymbol{\theta}$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{X}, \boldsymbol{Y}, \boldsymbol{\theta}) = -2\boldsymbol{X}_D^\top \boldsymbol{Y} + 2\boldsymbol{X}_D^\top \boldsymbol{X}_D \boldsymbol{\theta}$$

We want to find the $\boldsymbol{\theta}$ that makes the gradient of the loss zero

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{X}, \boldsymbol{Y}, \boldsymbol{\theta}) = -2\boldsymbol{X}_D^{\top} \boldsymbol{Y} + 2\boldsymbol{X}_D^{\top} \boldsymbol{X}_D \boldsymbol{\theta}$$

We want to find the $\boldsymbol{\theta}$ that makes the gradient of the loss zero

$$\mathbf{0} = -2\boldsymbol{X}_D^{\top} \boldsymbol{Y} + 2\boldsymbol{X}_D^{\top} \boldsymbol{X}_D \boldsymbol{\theta}$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{X}, \boldsymbol{Y}, \boldsymbol{\theta}) = -2\boldsymbol{X}_D^\top \boldsymbol{Y} + 2\boldsymbol{X}_D^\top \boldsymbol{X}_D \boldsymbol{\theta}$$

We want to find the $\boldsymbol{\theta}$ that makes the gradient of the loss zero

$$\mathbf{0} = -2\boldsymbol{X}_D^\top \boldsymbol{Y} + 2\boldsymbol{X}_D^\top \boldsymbol{X}_D \boldsymbol{\theta}$$

$$2\boldsymbol{X}_D^\top \boldsymbol{Y} = 2\boldsymbol{X}_D^\top \boldsymbol{X}_D \boldsymbol{\theta}$$

$$\nabla_{\theta} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \theta) = -2\mathbf{X}_D^{\top} \mathbf{Y} + 2\mathbf{X}_D^{\top} \mathbf{X}_D \theta$$

We want to find the θ that makes the gradient of the loss zero

$$\mathbf{0} = -2\mathbf{X}_D^{\top} \mathbf{Y} + 2\mathbf{X}_D^{\top} \mathbf{X}_D \theta$$

$$2\mathbf{X}_D^{\top} \mathbf{Y} = 2\mathbf{X}_D^{\top} \mathbf{X}_D \theta$$

$$\mathbf{X}_D^{\top} \mathbf{Y} = \mathbf{X}_D^{\top} \mathbf{X}_D \theta$$

$$\nabla_{\theta} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \theta) = -2\mathbf{X}_D^{\top} \mathbf{Y} + 2\mathbf{X}_D^{\top} \mathbf{X}_D \theta$$

We want to find the θ that makes the gradient of the loss zero

$$\mathbf{0} = -2\mathbf{X}_D^{\top} \mathbf{Y} + 2\mathbf{X}_D^{\top} \mathbf{X}_D \theta$$

$$2\mathbf{X}_D^{\top} \mathbf{Y} = 2\mathbf{X}_D^{\top} \mathbf{X}_D \theta$$

$$\mathbf{X}_D^{\top} \mathbf{Y} = \mathbf{X}_D^{\top} \mathbf{X}_D \theta$$

$$(\mathbf{X}_D^{\top} \mathbf{X}_D)^{-1} \mathbf{X}_D^{\top} \mathbf{Y} = \theta$$

$$\left(\mathbf{X}_D^\top \mathbf{X}_D\right)^{-1} \mathbf{X}_D^\top \mathbf{Y} = \boldsymbol{\theta}$$

$$(\mathbf{X}_D^\top \mathbf{X}_D)^{-1} \mathbf{X}_D^\top \mathbf{Y} = \boldsymbol{\theta}$$

This was the “magic” solution I gave you for linear regression

$$(\mathbf{X}_D^\top \mathbf{X}_D)^{-1} \mathbf{X}_D^\top \mathbf{Y} = \boldsymbol{\theta}$$

This was the “magic” solution I gave you for linear regression

$$\boldsymbol{\theta} = (\mathbf{X}_D^\top \mathbf{X}_D)^{-1} \mathbf{X}_D^\top \mathbf{Y}$$

Great! We derived the solution to linear regression

Great! We derived the solution to linear regression

Now, we will do the same approach for neural networks

Great! We derived the solution to linear regression

Now, we will do the same approach for neural networks

To make it simple, we assume $d_x = 1, d_y = 1, n = 1$

Great! We derived the solution to linear regression

Now, we will do the same approach for neural networks

To make it simple, we assume $d_x = 1, d_y = 1, n = 1$

One input dimension, one output dimension, one datapoint

Great! We derived the solution to linear regression

Now, we will do the same approach for neural networks

To make it simple, we assume $d_x = 1, d_y = 1, n = 1$

One input dimension, one output dimension, one datapoint

Step 1: Write the loss function for a neural network

Like linear regression, we can use square error for a neural network

Like linear regression, we can use square error for a neural network

$$\mathcal{L}(x, y, \boldsymbol{\theta}) = (f(x, \boldsymbol{\theta}) - y)^2$$

Like linear regression, we can use square error for a neural network

$$\mathcal{L}(x, y, \boldsymbol{\theta}) = (f(x, \boldsymbol{\theta}) - y)^2$$

All that changes is the model f

Like linear regression, we can use square error for a neural network

$$\mathcal{L}(x, y, \boldsymbol{\theta}) = (f(x, \boldsymbol{\theta}) - y)^2$$

All that changes is the model f

Linear regression:

$$f(x, y, \boldsymbol{\theta}) = \theta_0 + \theta_1 x$$

Like linear regression, we can use square error for a neural network

$$\mathcal{L}(x, y, \boldsymbol{\theta}) = (f(x, \boldsymbol{\theta}) - y)^2$$

All that changes is the model f

Linear regression:

$$f(x, y, \boldsymbol{\theta}) = \theta_0 + \theta_1 x$$

Perceptron:

$$f(x, y, \boldsymbol{\theta}) = \sigma(\theta_0 + \theta_1 x)$$

$$\mathcal{L}(x, y, \boldsymbol{\theta}) = (f(x, \boldsymbol{\theta}) - y)^2$$

Loss function

$$\mathcal{L}(x, y, \boldsymbol{\theta}) = (f(x, \boldsymbol{\theta}) - y)^2$$

Loss function

$$f(x, \boldsymbol{\theta}) = \sigma(\theta_0 + \theta_1 x)$$

Neural network model

$$\mathcal{L}(x, y, \boldsymbol{\theta}) = (f(x, \boldsymbol{\theta}) - y)^2$$

Loss function

$$f(x, \boldsymbol{\theta}) = \sigma(\theta_0 + \theta_1 x)$$

Neural network model

Now, we plug the model f into the loss function

$$\mathcal{L}(x, y, \boldsymbol{\theta}) = (f(x, \boldsymbol{\theta}) - y)^2$$

Loss function

$$f(x, \boldsymbol{\theta}) = \sigma(\theta_0 + \theta_1 x)$$

Neural network model

Now, we plug the model f into the loss function

$$\mathcal{L}(x, y, \boldsymbol{\theta}) = (\sigma(\theta_0 + \theta_1 x) - y)^2$$

$$\mathcal{L}(x, y, \boldsymbol{\theta}) = (f(x, \boldsymbol{\theta}) - y)^2$$

Loss function

$$f(x, \boldsymbol{\theta}) = \sigma(\theta_0 + \theta_1 x)$$

Neural network model

Now, we plug the model f into the loss function

$$\mathcal{L}(x, y, \boldsymbol{\theta}) = (\sigma(\theta_0 + \theta_1 x) - y)^2$$

Rewrite

$$\mathcal{L}(x, y, \boldsymbol{\theta}) = (f(x, \boldsymbol{\theta}) - y)^2$$

Loss function

$$f(x, \boldsymbol{\theta}) = \sigma(\theta_0 + \theta_1 x)$$

Neural network model

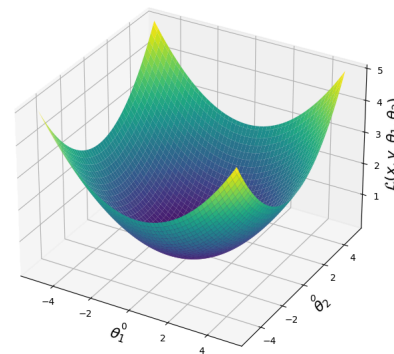
Now, we plug the model f into the loss function

$$\mathcal{L}(x, y, \boldsymbol{\theta}) = (\sigma(\theta_0 + \theta_1 x) - y)^2$$

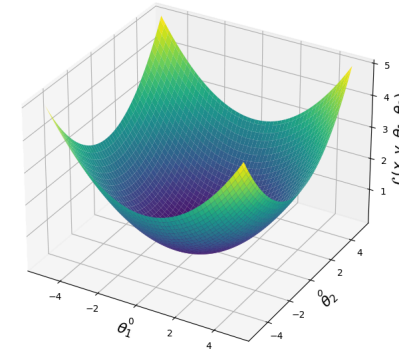
Rewrite

$$\mathcal{L}(x, y, \boldsymbol{\theta}) = \underbrace{(\sigma(\theta_0 + \theta_1 x) - y)}_{\text{Nonlinear function of } \theta} \underbrace{(\sigma(\theta_0 + \theta_1 x) - y)}_{\text{Nonlinear function of } \theta}$$

Linear regression loss function
was quadratic with one minima



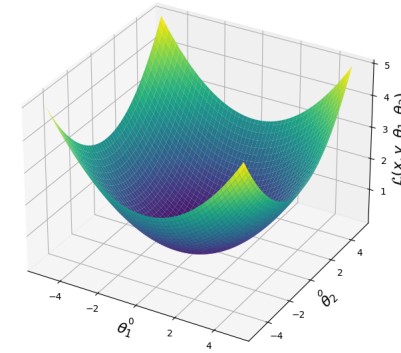
Linear regression loss function
was quadratic with one minima



With a neural network, this is our loss function

$$\mathcal{L}(x, y, \theta) = \underbrace{(\sigma(\theta_0 + \theta_1 x) - y)}_{\text{Nonlinear function of } \theta} \underbrace{(\sigma(\theta_0 + \theta_1 x) - y)}_{\text{Nonlinear function of } \theta}$$

Linear regression loss function
was quadratic with one minima

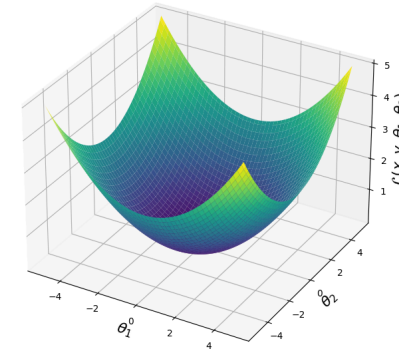


With a neural network, this is our loss function

$$\mathcal{L}(x, y, \theta) = \underbrace{(\sigma(\theta_0 + \theta_1 x) - y)}_{\text{Nonlinear function of } \theta} \underbrace{(\sigma(\theta_0 + \theta_1 x) - y)}_{\text{Nonlinear function of } \theta}$$

Question: How many minima does this function have?

Linear regression loss function
was quadratic with one minima



With a neural network, this is our loss function

$$\mathcal{L}(x, y, \theta) = \underbrace{(\sigma(\theta_0 + \theta_1 x) - y)}_{\text{Nonlinear function of } \theta} \underbrace{(\sigma(\theta_0 + \theta_1 x) - y)}_{\text{Nonlinear function of } \theta}$$

Question: How many minima does this function have?

Answer: We do not know

The nonlinearity/activation function σ in the neural network means we cannot find an analytical solution

The nonlinearity/activation function σ in the neural network means we cannot find an analytical solution

$$f(x, \boldsymbol{\theta}) = \sigma(\theta_0 + \theta_1 x)$$

The nonlinearity/activation function σ in the neural network means we cannot find an analytical solution

$$f(x, \theta) = \sigma(\theta_0 + \theta_1 x)$$

Question: Can we remove the activation function σ ?

The nonlinearity/activation function σ in the neural network means we cannot find an analytical solution

$$f(x, \theta) = \sigma(\theta_0 + \theta_1 x)$$

Question: Can we remove the activation function σ ?

Answer: Yes, but the result is linear regression

The nonlinearity/activation function σ in the neural network means we cannot find an analytical solution

$$f(x, \boldsymbol{\theta}) = \sigma(\theta_0 + \theta_1 x)$$

Question: Can we remove the activation function σ ?

Answer: Yes, but the result is linear regression

$$f(x, \boldsymbol{\theta}) = \theta_0 + \theta_1 x$$

The nonlinearity/activation function σ in the neural network means we cannot find an analytical solution

$$f(x, \boldsymbol{\theta}) = \sigma(\theta_0 + \theta_1 x)$$

Question: Can we remove the activation function σ ?

Answer: Yes, but the result is linear regression

$$f(x, \boldsymbol{\theta}) = \theta_0 + \theta_1 x$$

Activation functions make the neural network powerful

Linear regression: analytical solution for θ

Linear regression: analytical solution for θ

Neural network: no analytical solution for θ

Linear regression: analytical solution for θ

Neural network: no analytical solution for θ

So how do to find θ for a neural network?

To find θ for a neural network, we use **gradient descent**

To find θ for a neural network, we use **gradient descent**

Gradient descent optimizes **differentiable** functions

To find θ for a neural network, we use **gradient descent**

Gradient descent optimizes **differentiable** functions

We must be able to take the derivative or gradient of the loss function to use gradient descent

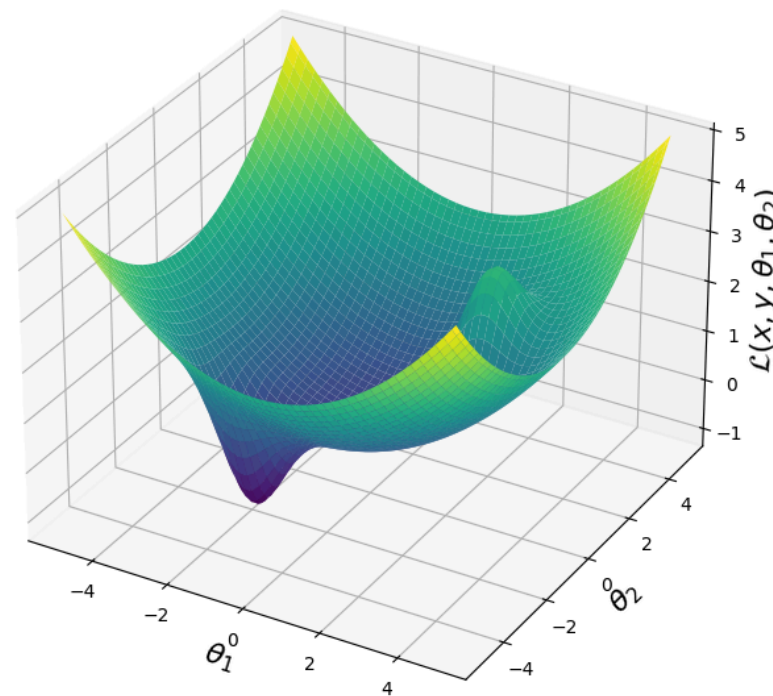
To find θ for a neural network, we use **gradient descent**

Gradient descent optimizes **differentiable** functions

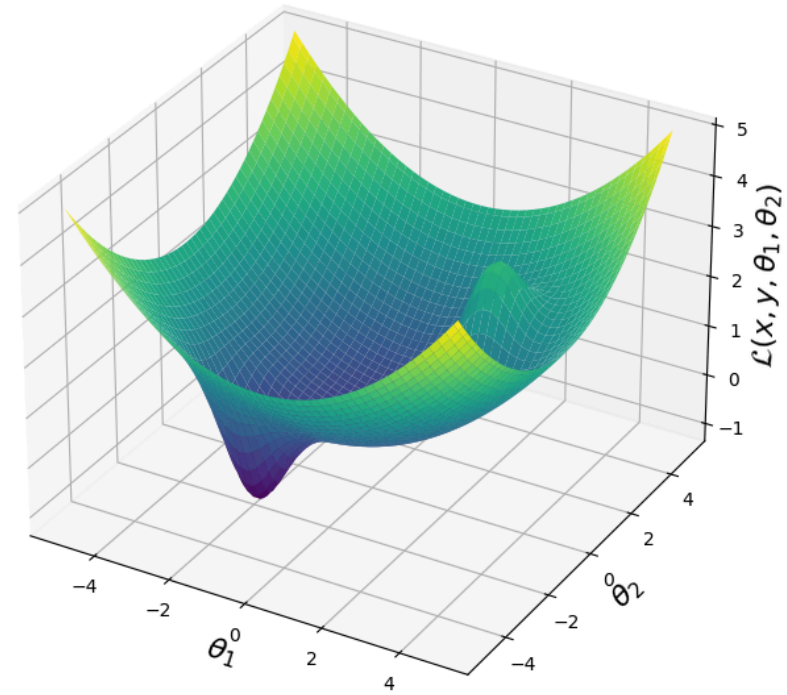
We must be able to take the derivative or gradient of the loss function to use gradient descent

How does gradient descent work?

A differentiable loss function
produces a manifold

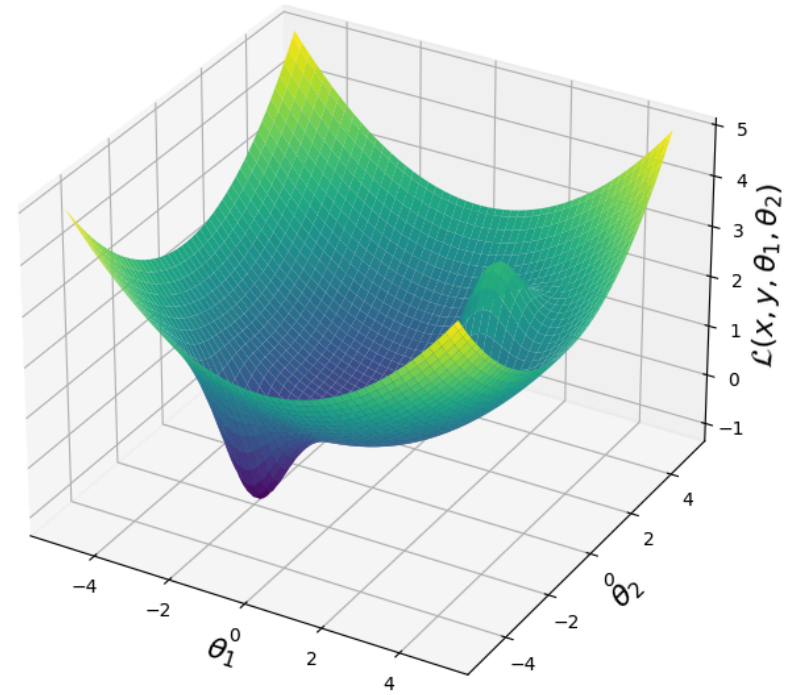


A differentiable loss function
produces a manifold



Our goal is to find the lowest point on this manifold

A differentiable loss function
produces a manifold

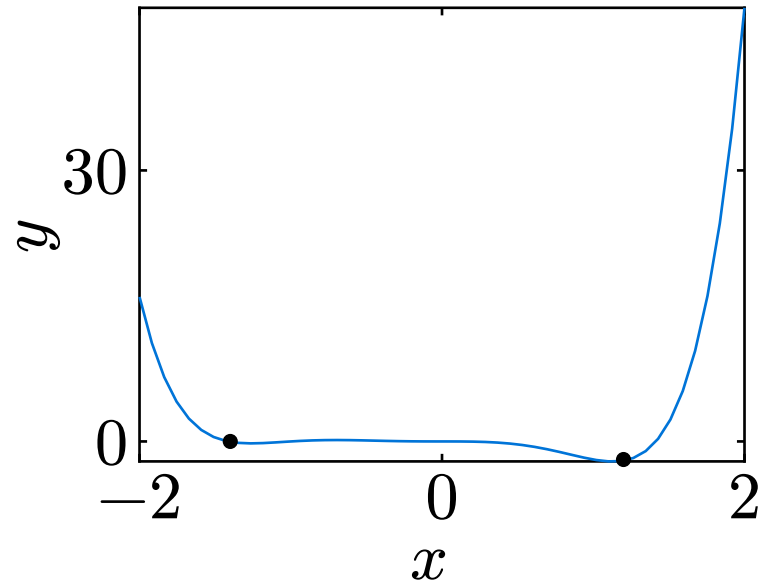


Our goal is to find the lowest point on this manifold

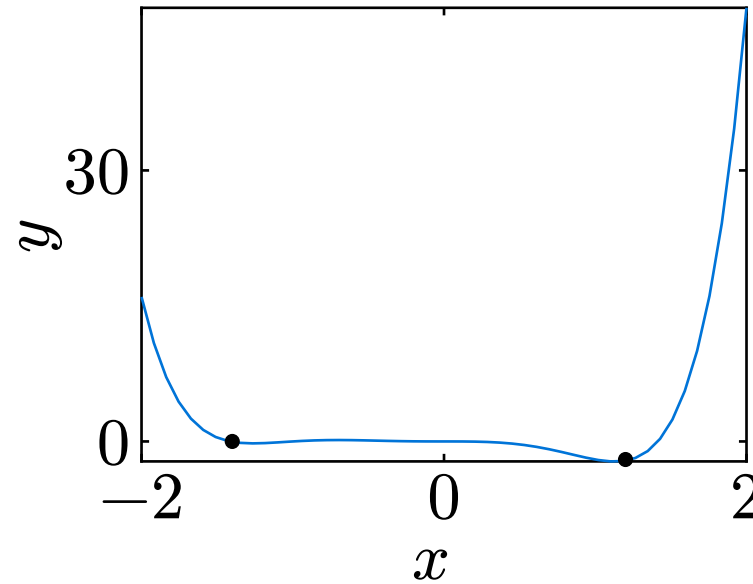
The lowest point solves $\arg \min_{\theta} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \theta)$

Gradient descent provides a **local** optima, not necessarily a **global** optima

Gradient descent provides a **local** optima, not necessarily a **global** optima



Gradient descent provides a **local** optima, not necessarily a **global** optima



In practice, a local optima provides a good enough model

Let us define gradient descent without math

Let us define gradient descent without math

You are on the top of a mountain and there is lightning storm

Let us define gradient descent without math

You are on the top of a mountain and there is lightning storm



Let us define gradient descent without math

You are on the top of a mountain and there is lightning storm



For safety, you should walk down the mountain to escape the lightning

But you do not know the path down!

But you do not know the path down!



You see this, which way do you walk next?



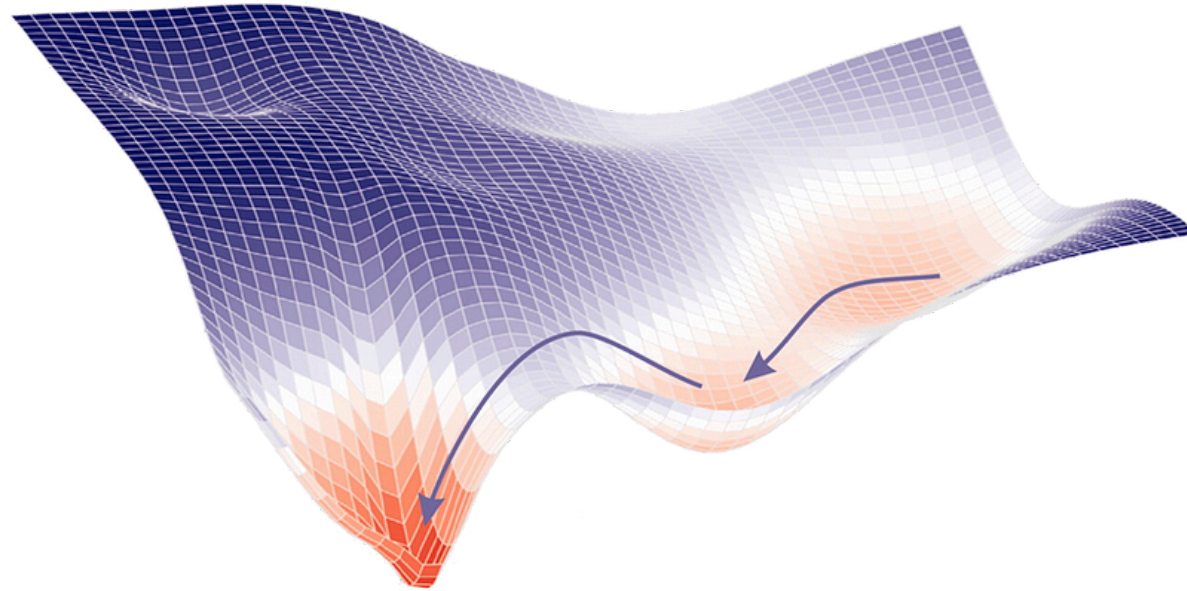
This is gradient descent

In gradient descent, we look at the **slope** of the loss function

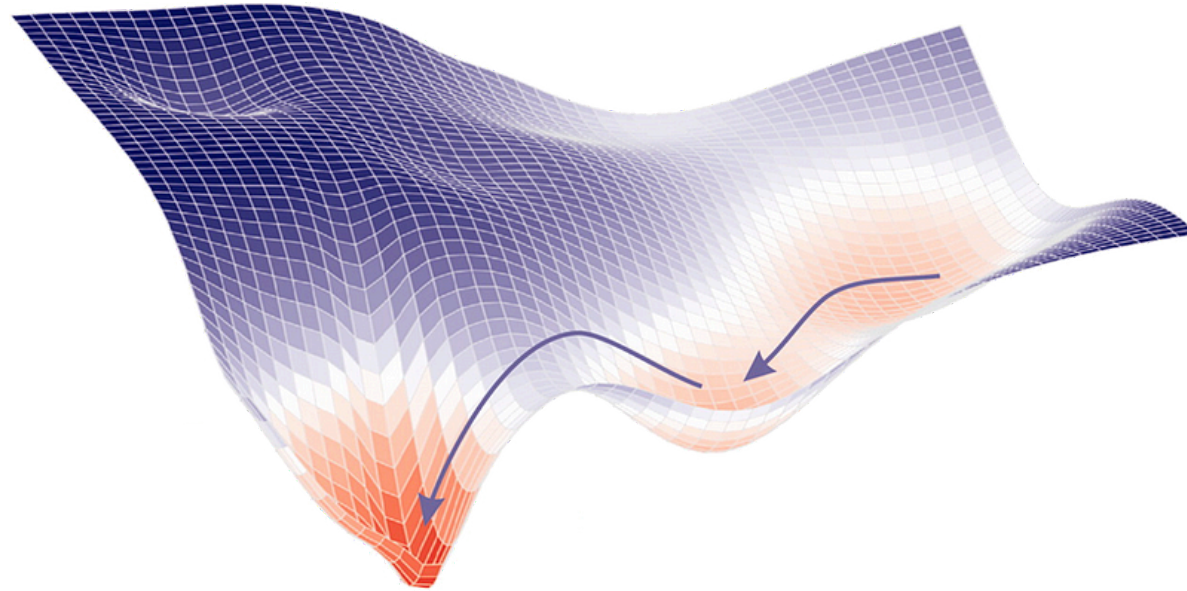
In gradient descent, we look at the **slope** of the loss function

And we walk in the steepest direction

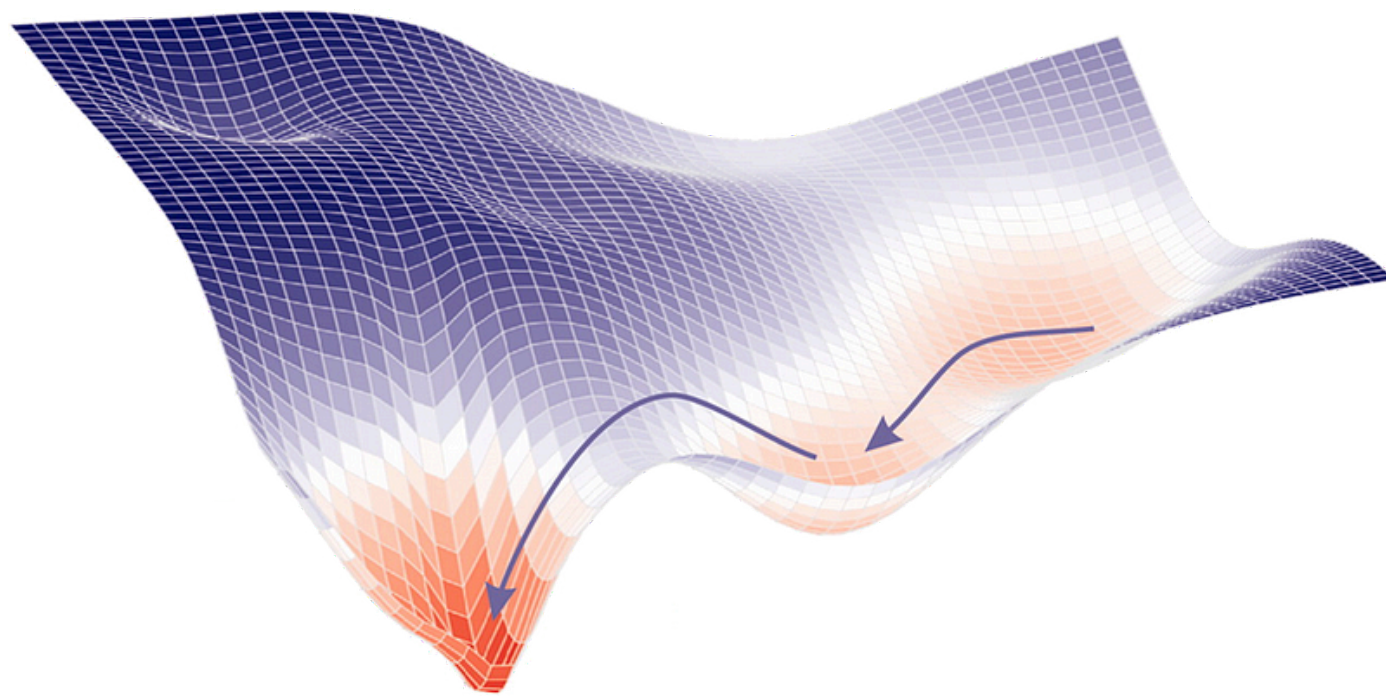
In gradient descent, we look at the **slope** of the loss function
And we walk in the steepest direction

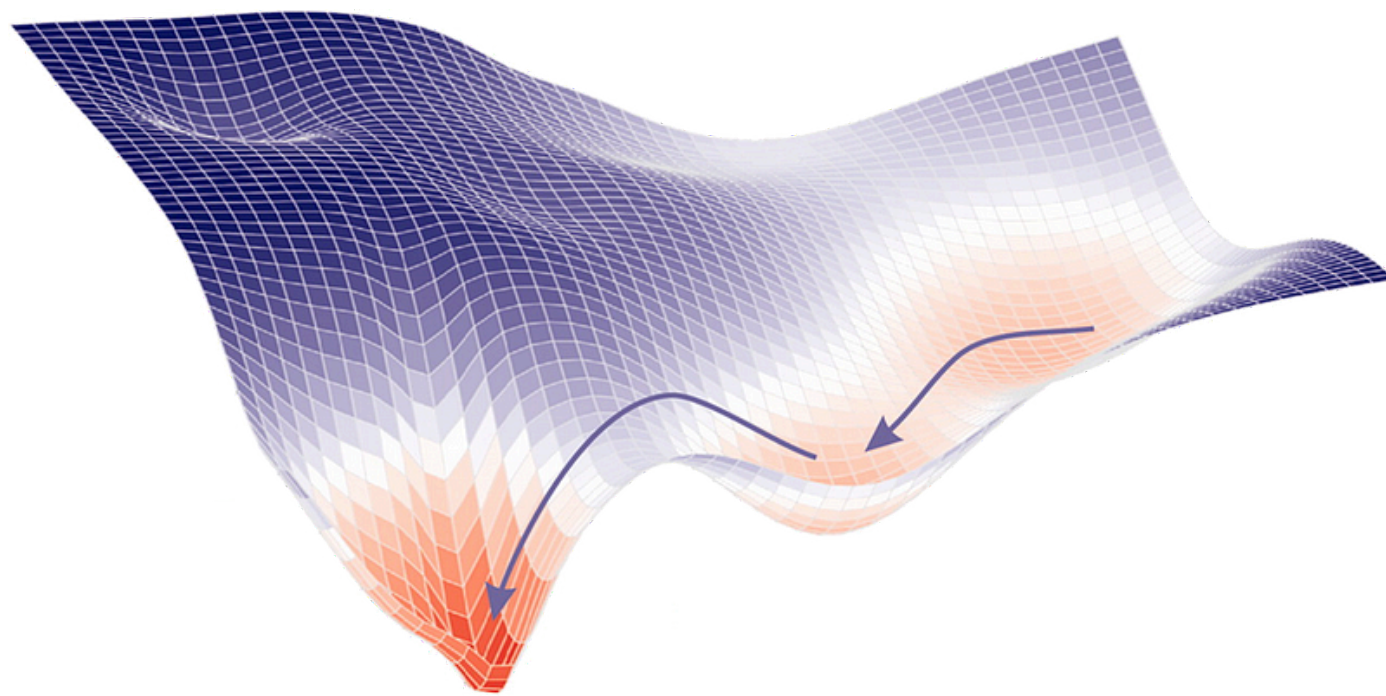


In gradient descent, we look at the **slope** of the loss function
And we walk in the steepest direction

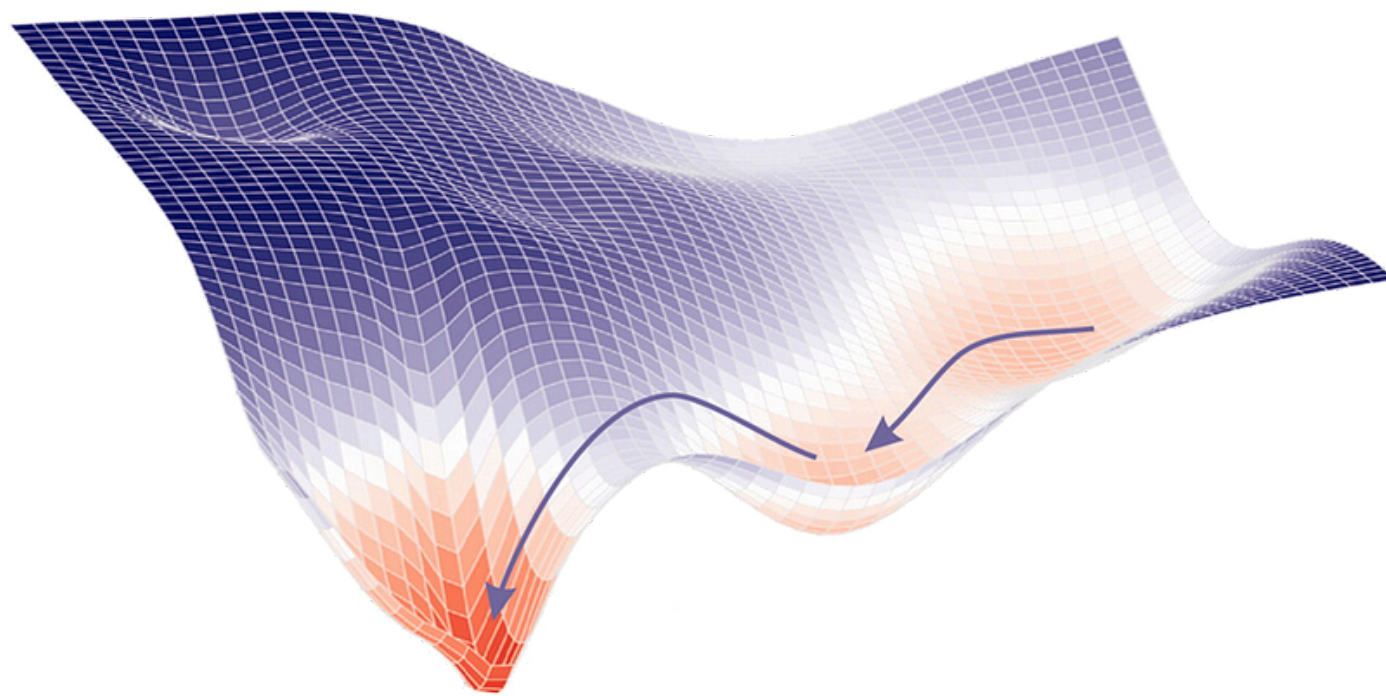


And then we repeat



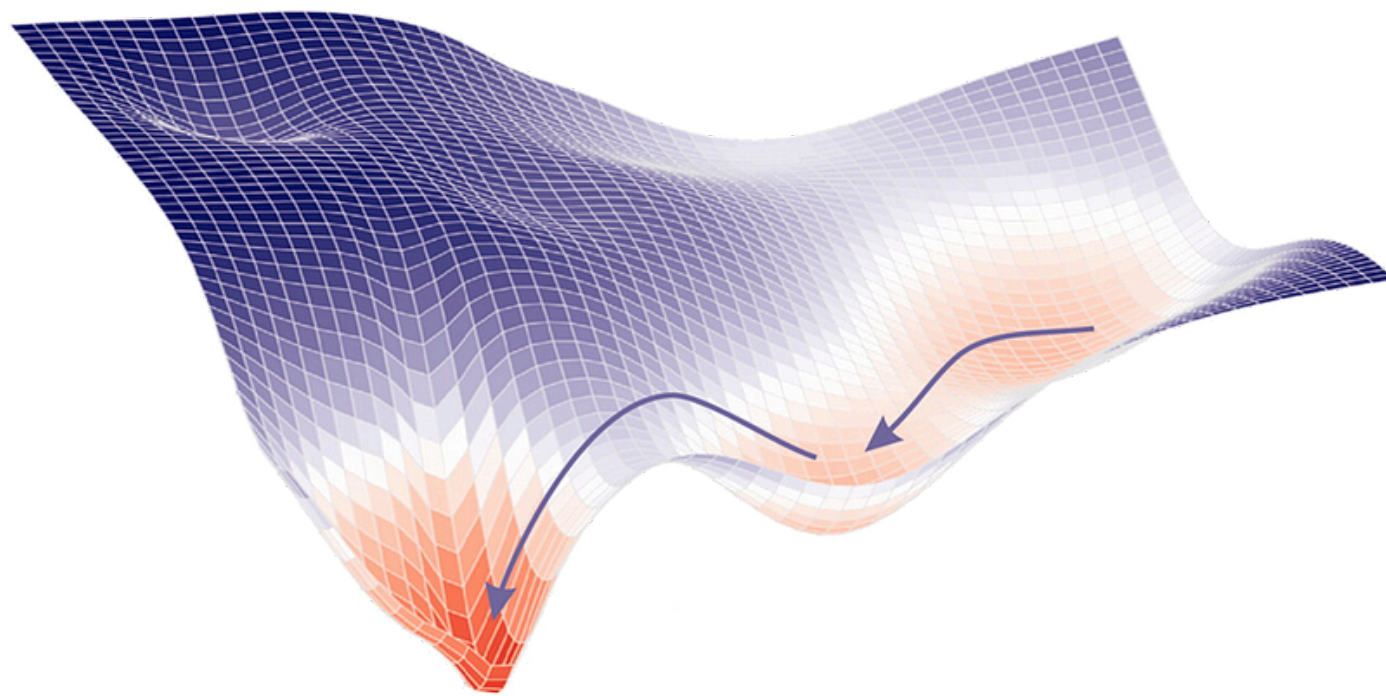


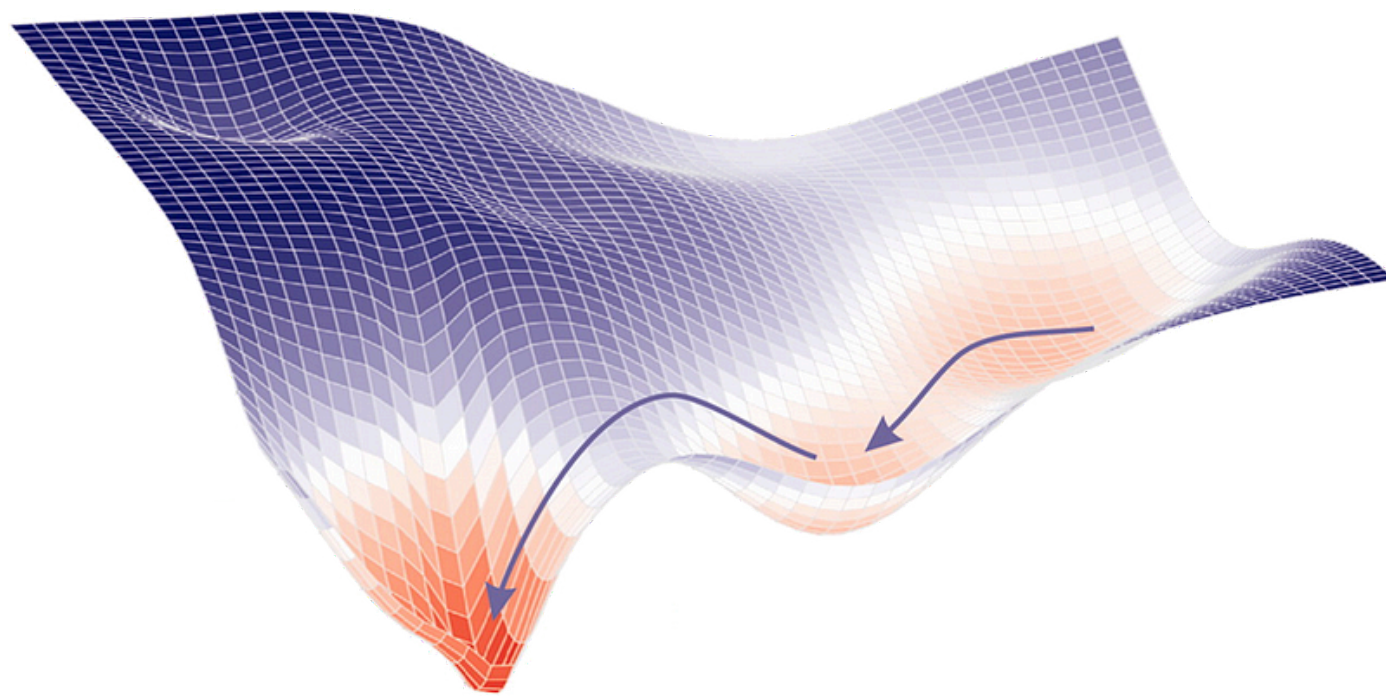
We find the gradient $\nabla_{\theta} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \theta)$



We find the gradient $\nabla_{\theta} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \theta)$

And update θ in the steepest direction





Eventually, we arrive at the bottom

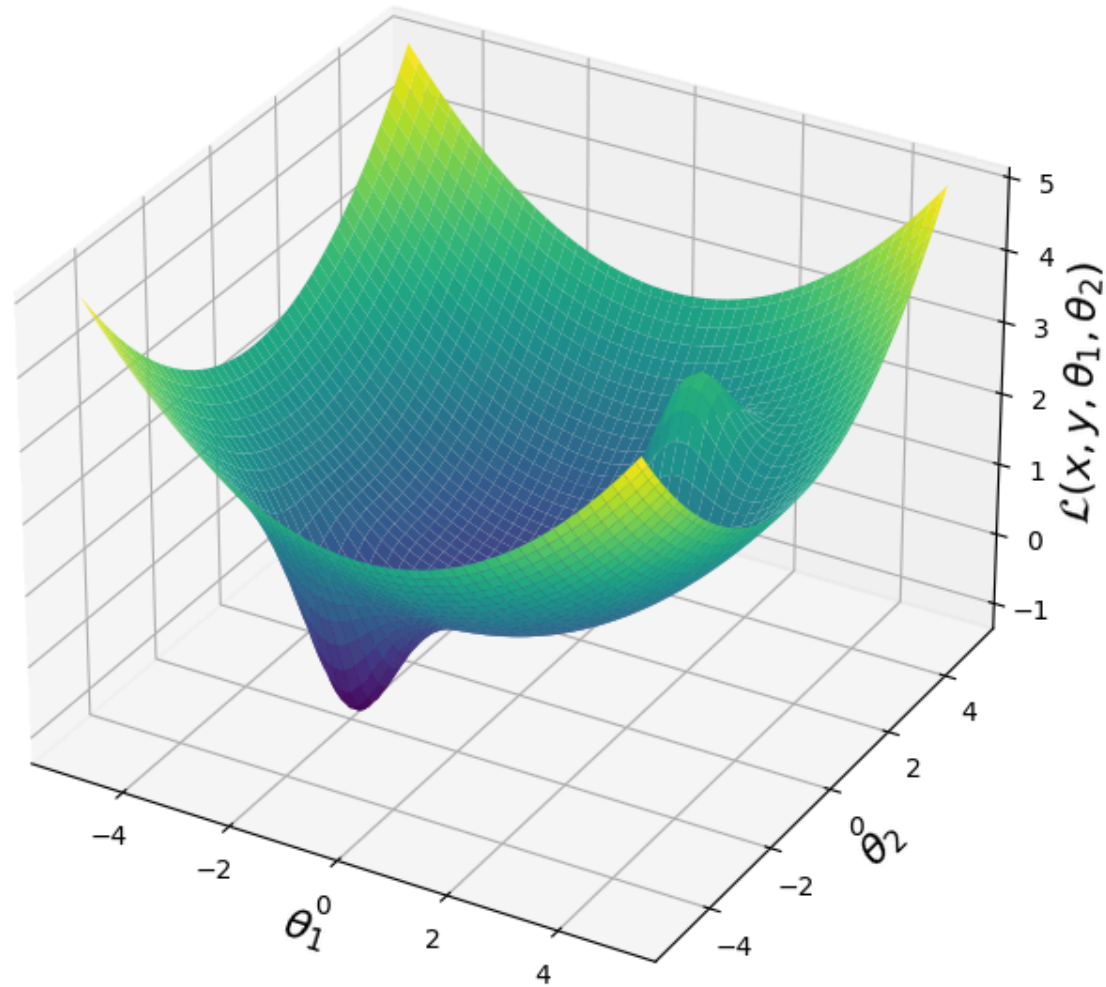
With gradient descent, the loss function must be differentiable

With gradient descent, the loss function must be differentiable

If we cannot compute the derivative/gradient, then we do not know which way to walk!

The gradient descent algorithm is as follows:

```
1: function GRADIENT DESCENT( $x, y, \mathcal{L}, t, \alpha$ )
2:     ▷ Randomly initialize parameters
3:      $\theta \leftarrow \mathcal{N}(0, 1)$ 
4:     for  $i \in 1 \dots t$  do
5:         ▷ Compute the gradient of the loss
6:          $J \leftarrow \nabla_{\theta} \mathcal{L}(x, y, \theta)$ 
7:         ▷ Update the parameters using the negative gradient
8:          $\theta \leftarrow \theta - \alpha J$ 
9:     return  $\theta$ 
```



Two main steps in gradient descent:

Two main steps in gradient descent:

Step 1: Compute the gradient of the loss

Two main steps in gradient descent:

Step 1: Compute the gradient of the loss

Step 2: Update the parameters using the gradient

Two main steps in gradient descent:

Step 1: Compute the gradient of the loss

Step 2: Update the parameters using the gradient

Let us start with step 1

Agenda

1. Review
2. Quiz
3. Optimization
4. Calculus review
5. **Deriving linear regression**
6. Gradient descent
7. Backpropagation
8. Layer gradient
9. Full gradient
10. Autograd
11. Practical considerations

Agenda

1. Review
2. Quiz
3. Optimization
4. Calculus review
5. Deriving linear regression
6. **Gradient descent**
7. Backpropagation
8. Layer gradient
9. Full gradient
10. Autograd
11. Practical considerations

Goal: Compute the gradient of the loss $\nabla_{\theta} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \theta)$

Goal: Compute the gradient of the loss $\nabla_{\theta} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \theta)$

We call this process **backpropagation**

Goal: Compute the gradient of the loss $\nabla_{\theta} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \theta)$

We call this process **backpropagation**

We propagate errors from the loss function **backward** through each layer of the neural network

Goal: Compute the gradient of the loss $\nabla_{\theta} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \theta)$

We call this process **backpropagation**

We propagate errors from the loss function **backward** through each layer of the neural network

Let us propagate errors through a deep neural network

Finding the gradient is necessary to use gradient descent!

Finding the gradient is necessary to use gradient descent!

We will use the gradient of layers to find the gradient of a deep neural network

Finding the gradient is necessary to use gradient descent!

We will use the gradient of layers to find the gradient of a deep neural network

First, let us compute the gradient of a neural network layer

Agenda

1. Review
2. Quiz
3. Optimization
4. Calculus review
5. Deriving linear regression
6. Gradient descent
7. **Backpropagation**
8. Layer gradient
9. Full gradient
10. Autograd
11. Practical considerations

Agenda

1. Review
2. Quiz
3. Optimization
4. Calculus review
5. Deriving linear regression
6. Gradient descent
7. Backpropagation
8. **Layer gradient**
9. Full gradient
10. Autograd
11. Practical considerations

Start with the equation of a neural network layer

$$f(\boldsymbol{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}^\top \overline{\boldsymbol{x}})$$

Start with the equation of a neural network layer

$$f(\boldsymbol{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}^\top \overline{\boldsymbol{x}})$$

Take the gradient of both sides

$$\nabla_{\boldsymbol{\theta}} f(\boldsymbol{x}, \boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \sigma(\boldsymbol{\theta}^\top \overline{\boldsymbol{x}})$$

Start with the equation of a neural network layer

$$f(x, \theta) = \sigma(\theta^\top \bar{x})$$

Take the gradient of both sides

$$\nabla_{\theta} f(x, \theta) = \nabla_{\theta} \sigma(\theta^\top \bar{x})$$

$$\text{Chain: } \frac{d}{dx} f(g(x)) = f'(g(x))g'(x)$$

Start with the equation of a neural network layer

$$f(x, \theta) = \sigma(\theta^\top \bar{x})$$

Take the gradient of both sides

$$\nabla_{\theta} f(x, \theta) = \nabla_{\theta} \sigma(\theta^\top \bar{x})$$

$$\text{Chain: } \frac{d}{dx} f(g(x)) = f'(g(x)) g'(x)$$

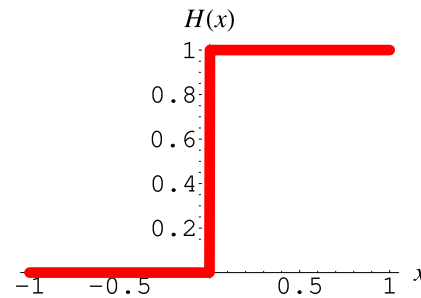
$$\nabla_{\theta} f(x, \theta) = \frac{\partial \sigma}{\partial \theta}(\theta^\top \bar{x}) \cdot \nabla_{\theta}(\theta^\top \bar{x})$$

$$\nabla_{\boldsymbol{\theta}} f(\boldsymbol{x}, \boldsymbol{\theta}) = \frac{\partial \sigma}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}^\top \bar{\boldsymbol{x}}) \nabla_{\boldsymbol{\theta}}(\boldsymbol{\theta}^\top \bar{\boldsymbol{x}})$$

What is $\frac{\partial \sigma}{\partial x}(z)$?

$$\nabla_{\theta} f(x, \theta) = \frac{\partial \sigma}{\partial \theta}(\theta^{\top} \bar{x}) \nabla_{\theta}(\theta^{\top} \bar{x})$$

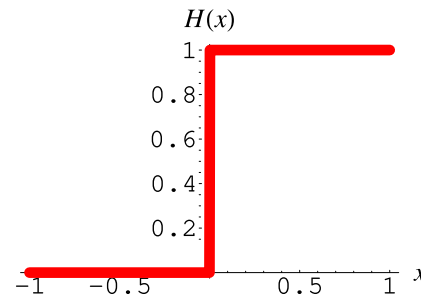
What is $\frac{\partial \sigma}{\partial x}(z)$?



Derivative is zero everywhere and infinity at $x = 0$, so the derivative for a layer is either infinity or zero

$$\nabla_{\theta} f(x, \theta) = \frac{\partial \sigma}{\partial \theta}(\theta^{\top} \bar{x}) \nabla_{\theta}(\theta^{\top} \bar{x})$$

What is $\frac{\partial \sigma}{\partial x}(z)$?



Derivative is zero everywhere and infinity at $x = 0$, so the derivative for a layer is either infinity or zero

We use a differentiable approximation of the heaviside step function

We use a differentiable approximation of the heaviside step function

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

We use a differentiable approximation of the heaviside step function

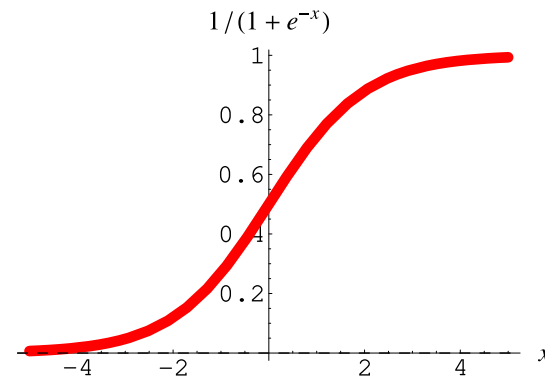
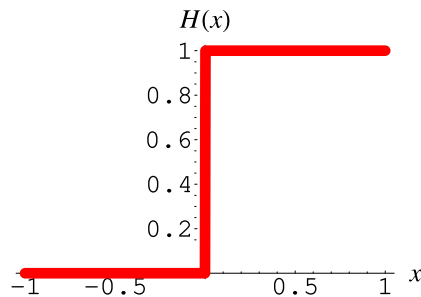
$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

We call this approximation the **sigmoid function**

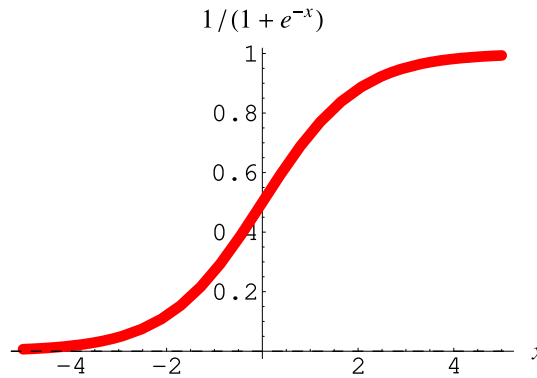
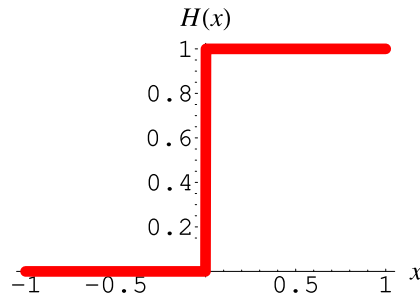
We use a differentiable approximation of the heaviside step function

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

We call this approximation the **sigmoid function**



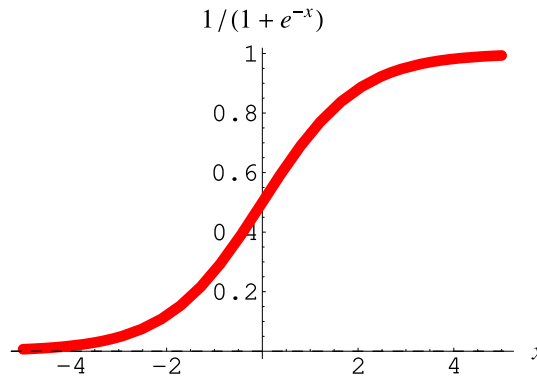
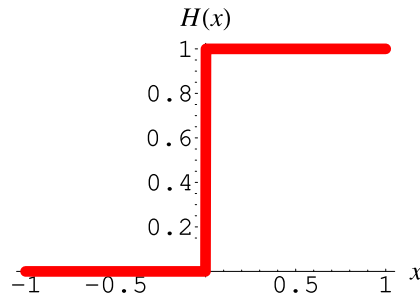
The sigmoid function has finite and nonzero derivative everywhere



$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

The derivative of the sigmoid function is

$$\frac{d}{dz}\sigma(z) = \sigma(z) \cdot (1 - \sigma(z))$$

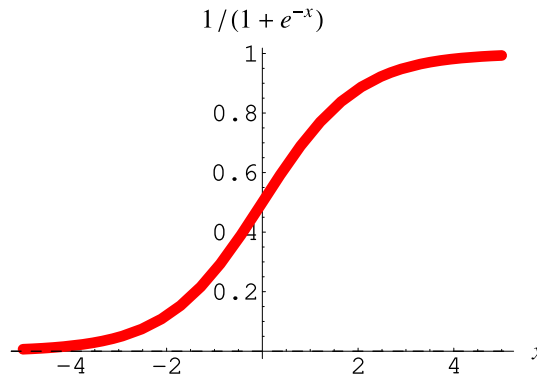
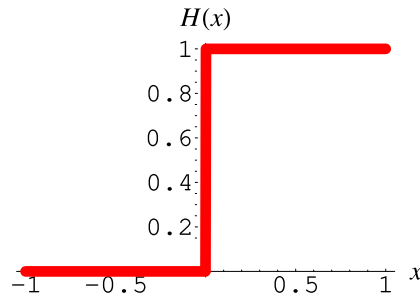


$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

The derivative of the sigmoid function is

$$\frac{d}{dz}\sigma(z) = \sigma(z) \cdot (1 - \sigma(z))$$

$$\frac{\partial}{\partial z}\sigma(z) = \sigma(z) \odot (1 - \sigma(z))$$



$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

The derivative of the sigmoid function is

$$\frac{d}{dz}\sigma(z) = \sigma(z) \cdot (1 - \sigma(z))$$

$$\frac{\partial}{\partial \mathbf{z}}\sigma(\mathbf{z}) = \sigma(\mathbf{z}) \odot (1 - \sigma(\mathbf{z}))$$

Back to our layer

Back to our layer

$$\nabla_{\theta} f(x, \theta) = \frac{\partial \sigma}{\partial \theta}(\theta^{\top} \bar{x}) \nabla_{\theta}(\theta^{\top} \bar{x})$$

Back to our layer

$$\nabla_{\theta} f(x, \theta) = \frac{\partial \sigma}{\partial \theta}(\theta^{\top} \bar{x}) \nabla_{\theta}(\theta^{\top} \bar{x})$$

Plug in the derivative of our new activation function

Back to our layer

$$\nabla_{\theta} f(x, \theta) = \frac{\partial \sigma}{\partial \theta}(\theta^{\top} \bar{x}) \nabla_{\theta}(\theta^{\top} \bar{x})$$

Plug in the derivative of our new activation function

$$\nabla_{\theta} f(x, \theta) = (\sigma(\theta^{\top} \bar{x}) \odot (1 - \sigma(\theta^{\top} \bar{x}))) (\theta^{\top} \bar{x}) \nabla_{\theta}(\theta^{\top} \bar{x})$$

Evaluate the final term

Back to our layer

$$\nabla_{\theta} f(x, \theta) = \frac{\partial \sigma}{\partial \theta}(\theta^{\top} \bar{x}) \nabla_{\theta}(\theta^{\top} \bar{x})$$

Plug in the derivative of our new activation function

$$\nabla_{\theta} f(x, \theta) = (\sigma(\theta^{\top} \bar{x}) \odot (1 - \sigma(\theta^{\top} \bar{x}))) (\theta^{\top} \bar{x}) \nabla_{\theta}(\theta^{\top} \bar{x})$$

Evaluate the final term

$$\nabla_{\theta} f(x, \theta) = (\sigma(\theta^{\top} \bar{x}) \odot (1 - \sigma(\theta^{\top} \bar{x}))) (\theta^{\top} \bar{x}) \bar{x}^{\top}$$

$$\nabla_{\theta} f(x, \theta) = (\sigma(\theta^{\top} \bar{x}) \odot (1 - \sigma(\theta^{\top} \bar{x}))) (\theta^{\top} \bar{x}) \bar{x}^{\top}$$

$$\nabla_{\theta} f(x, \theta) = (\sigma(\theta^{\top} \bar{x}) \odot (1 - \sigma(\theta^{\top} \bar{x}))) (\theta^{\top} \bar{x}) \bar{x}^{\top}$$

This is the gradient for the layer of a neural network!

$$\nabla_{\theta} f(x, \theta) = (\sigma(\theta^{\top} \bar{x}) \odot (1 - \sigma(\theta^{\top} \bar{x}))) (\theta^{\top} \bar{x}) \bar{x}^{\top}$$

This is the gradient for the layer of a neural network!

We will use this to compute the gradient of a deep neural network

Agenda

1. Review
2. Quiz
3. Optimization
4. Calculus review
5. Deriving linear regression
6. Gradient descent
7. Backpropagation
8. **Layer gradient**
9. Full gradient
10. Autograd
11. Practical considerations

Agenda

1. Review
2. Quiz
3. Optimization
4. Calculus review
5. Deriving linear regression
6. Gradient descent
7. Backpropagation
8. Layer gradient
9. **Full gradient**
10. Autograd
11. Practical considerations

Recall the deep neural network has many layers

$$f_1(\mathbf{x}, \boldsymbol{\varphi}) = \sigma(\boldsymbol{\varphi}^\top \overline{\mathbf{x}}) \quad f_2(\mathbf{x}, \boldsymbol{\psi}) = \sigma(\boldsymbol{\psi}^\top \overline{\mathbf{x}}) \quad \dots \quad f_\ell(\mathbf{x}, \boldsymbol{\xi}) = \sigma(\boldsymbol{\xi}^\top \overline{\mathbf{x}})$$

Recall the deep neural network has many layers

$$f_1(\mathbf{x}, \boldsymbol{\varphi}) = \sigma(\boldsymbol{\varphi}^\top \overline{\mathbf{x}}) \quad f_2(\mathbf{x}, \boldsymbol{\psi}) = \sigma(\boldsymbol{\psi}^\top \overline{\mathbf{x}}) \quad \dots \quad f_\ell(\mathbf{x}, \boldsymbol{\xi}) = \sigma(\boldsymbol{\xi}^\top \overline{\mathbf{x}})$$

And that we call them in series

$$z_1 = f_1(\mathbf{x}, \boldsymbol{\varphi})$$

$$z_2 = f_2(z_1, \boldsymbol{\psi})$$

$$\vdots$$

$$z_\ell = f_\ell(z_{\ell-1}, \boldsymbol{\xi})$$

Recall the deep neural network has many layers

$$f_1(\mathbf{x}, \boldsymbol{\varphi}) = \sigma(\boldsymbol{\varphi}^\top \overline{\mathbf{x}}) \quad f_2(\mathbf{x}, \boldsymbol{\psi}) = \sigma(\boldsymbol{\psi}^\top \overline{\mathbf{x}}) \quad \dots \quad f_\ell(\mathbf{x}, \boldsymbol{\xi}) = \sigma(\boldsymbol{\xi}^\top \overline{\mathbf{x}})$$

And that we call them in series

$$z_1 = f_1(\mathbf{x}, \boldsymbol{\varphi})$$

$$z_2 = f_2(z_1, \boldsymbol{\psi})$$

$$\vdots$$

$$z_\ell = f_\ell(z_{\ell-1}, \boldsymbol{\xi})$$

Take the gradient of both sides

$$\nabla_{\varphi, \psi, \dots, \xi} z_1 = \nabla_{\varphi, \psi, \xi} f_1(x, \varphi)$$

$$\nabla_{\varphi, \psi, \dots, \xi} z_2 = \nabla_{\varphi, \psi, \xi} f_2(z_1, \psi)$$

$$\vdots$$

$$\nabla_{\varphi, \psi, \dots, \xi} z_\ell = \nabla_{\varphi, \psi, \xi} f_\ell(z_{\ell-1}, \xi)$$

Take the gradient of both sides

$$\begin{aligned}\nabla_{\varphi, \psi, \dots, \xi} z_1 &= \nabla_{\varphi, \psi, \xi} f_1(x, \varphi) \\ \nabla_{\varphi, \psi, \dots, \xi} z_2 &= \nabla_{\varphi, \psi, \xi} f_2(z_1, \psi) \\ &\vdots \\ \nabla_{\varphi, \psi, \dots, \xi} z_\ell &= \nabla_{\varphi, \psi, \xi} f_\ell(z_{\ell-1}, \xi)\end{aligned}$$

Each layer only uses one set of parameters

$$\begin{aligned}\nabla_{\varphi} z_1 &= \nabla_{\varphi} f_1(x, \varphi) \\ \nabla_{\psi} z_2 &= \nabla_{\psi} f_2(z_1, \psi) \\ &\vdots \\ \nabla_{\xi} z_\ell &= \nabla_{\xi} f_\ell(z_{\ell-1}, \xi)\end{aligned}$$

Take the gradient of both sides

$$\begin{aligned}\nabla_{\varphi, \psi, \dots, \xi} z_1 &= \nabla_{\varphi, \psi, \xi} f_1(x, \varphi) \\ \nabla_{\varphi, \psi, \dots, \xi} z_2 &= \nabla_{\varphi, \psi, \xi} f_2(z_1, \psi) \\ &\vdots \\ \nabla_{\varphi, \psi, \dots, \xi} z_\ell &= \nabla_{\varphi, \psi, \xi} f_\ell(z_{\ell-1}, \xi)\end{aligned}$$

Each layer only uses one set of parameters

$$\begin{aligned}\nabla_{\varphi} z_1 &= \nabla_{\varphi} f_1(x, \varphi) \\ \nabla_{\psi} z_2 &= \nabla_{\psi} f_2(z_1, \psi) \\ &\vdots \\ \nabla_{\xi} z_\ell &= \nabla_{\xi} f_\ell(z_{\ell-1}, \xi)\end{aligned}$$

The gradient of a deep neural network is

$$\nabla_{\theta} f(x, \theta) = \nabla_{\varphi, \psi, \dots, \xi} f(x, [\varphi \ \psi \ \dots \ \xi]^{\top}) = \begin{bmatrix} \nabla_{\varphi} f_1(x, \varphi) \\ \nabla_{\psi} f_2(z_1, \psi) \\ \vdots \\ \nabla_{\xi} f_{\ell}(z_{\ell-1}, \xi) \end{bmatrix}$$

The gradient of a deep neural network is

$$\nabla_{\theta} f(x, \theta) = \nabla_{\varphi, \psi, \dots, \xi} f(x, [\varphi \ \psi \ \dots \ \xi]^{\top}) = \begin{bmatrix} \nabla_{\varphi} f_1(x, \varphi) \\ \nabla_{\psi} f_2(z_1, \psi) \\ \vdots \\ \nabla_{\xi} f_{\ell}(z_{\ell-1}, \xi) \end{bmatrix}$$

Where each layer gradient is

$$\nabla_{\xi} f_{\ell}(z_{\ell-1}, \xi) = (\sigma(\xi^{\top} \bar{z}_{\ell-1}) \odot (1 - \sigma(\xi^{\top} \bar{z}_{\ell-1}))) (\xi^{\top} \bar{z}_{\ell-1}) \bar{z}_{\ell-1}^{\top}$$

We computed the gradient of a neural network layer

We computed the gradient of a neural network layer

We computed the gradient of the neural network

We computed the gradient of a neural network layer

We computed the gradient of the neural network

Now, we must compute gradient of the loss function

We computed the gradient of a neural network layer

We computed the gradient of the neural network

Now, we must compute gradient of the loss function

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \sum_{i=1}^n \left(f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) - \mathbf{y}_{[i]} \right)^2$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \sum_{i=1}^n \left(f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) - \mathbf{y}_{[i]} \right)^2$$

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \sum_{i=1}^n \left(f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) - \mathbf{y}_{[i]} \right)^2$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \sum_{i=1}^n \left(f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) - \mathbf{y}_{[i]} \right)^2$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \sum_{i=1}^n \nabla_{\boldsymbol{\theta}} \left(f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) - \mathbf{y}_{[i]}^2 \right)$$

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \sum_{i=1}^n \left(f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) - \mathbf{y}_{[i]} \right)^2$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \sum_{i=1}^n \left(f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) - \mathbf{y}_{[i]} \right)^2$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \sum_{i=1}^n \nabla_{\boldsymbol{\theta}} \left(f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) - \mathbf{y}_{[i]}^2 \right)$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \sum_{i=1}^n \left(f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) - \mathbf{y}_{[i]} \right) \nabla_{\boldsymbol{\theta}} f(\mathbf{x}_{[i]}, \boldsymbol{\theta})$$

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \sum_{i=1}^n \left(f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) - \mathbf{y}_{[i]} \right)^2$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \sum_{i=1}^n \left(f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) - \mathbf{y}_{[i]} \right)^2$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \sum_{i=1}^n \nabla_{\boldsymbol{\theta}} \left(f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) - \mathbf{y}_{[i]}^2 \right)$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \sum_{i=1}^n \left(f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) - \mathbf{y}_{[i]} \right) \nabla_{\boldsymbol{\theta}} f(\mathbf{x}_{[i]}, \boldsymbol{\theta})$$

To summarize:

To summarize:

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{X}, \boldsymbol{Y}, \boldsymbol{\theta}) = \sum_{i=1}^n \left(f(\boldsymbol{x}_{[i]}, \boldsymbol{\theta}) - \boldsymbol{y}_{[i]} \right) \nabla_{\boldsymbol{\theta}} f(\boldsymbol{x}_{[i]}, \boldsymbol{\theta})$$

To summarize:

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{X}, \boldsymbol{Y}, \boldsymbol{\theta}) = \sum_{i=1}^n \left(f(\boldsymbol{x}_{[i]}, \boldsymbol{\theta}) - \boldsymbol{y}_{[i]} \right) \nabla_{\boldsymbol{\theta}} f(\boldsymbol{x}_{[i]}, \boldsymbol{\theta})$$

$$\nabla_{\boldsymbol{\theta}} f(\boldsymbol{x}, \boldsymbol{\theta}) = \nabla_{\boldsymbol{\varphi}, \boldsymbol{\psi}, \dots, \boldsymbol{\xi}} f\left(\boldsymbol{x}, [\boldsymbol{\varphi} \ \boldsymbol{\psi} \ \dots \ \boldsymbol{\xi}]^{\top}\right) = \begin{bmatrix} \nabla_{\boldsymbol{\varphi}} f_1(\boldsymbol{x}, \boldsymbol{\varphi}) \\ \nabla_{\boldsymbol{\psi}} f_2(\boldsymbol{z}_1, \boldsymbol{\psi}) \\ \vdots \\ \nabla_{\boldsymbol{\xi}} f_{\ell}(\boldsymbol{z}_{\ell-1}, \boldsymbol{\xi}) \end{bmatrix}$$

To summarize:

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \sum_{i=1}^n \left(f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) - \mathbf{y}_{[i]} \right) \nabla_{\boldsymbol{\theta}} f(\mathbf{x}_{[i]}, \boldsymbol{\theta})$$

$$\nabla_{\boldsymbol{\theta}} f(\mathbf{x}, \boldsymbol{\theta}) = \nabla_{\boldsymbol{\varphi}, \boldsymbol{\psi}, \dots, \boldsymbol{\xi}} f\left(\mathbf{x}, [\boldsymbol{\varphi} \ \boldsymbol{\psi} \ \dots \ \boldsymbol{\xi}]^{\top}\right) = \begin{bmatrix} \nabla_{\boldsymbol{\varphi}} f_1(\mathbf{x}, \boldsymbol{\varphi}) \\ \nabla_{\boldsymbol{\psi}} f_2(\mathbf{z}_1, \boldsymbol{\psi}) \\ \vdots \\ \nabla_{\boldsymbol{\xi}} f_{\ell}(\mathbf{z}_{\ell-1}, \boldsymbol{\xi}) \end{bmatrix}$$

$$\nabla_{\boldsymbol{\xi}} f_{\ell}(\mathbf{z}_{\ell-1}, \boldsymbol{\xi}) = (\sigma(\boldsymbol{\xi}^{\top} \bar{\mathbf{z}}_{\ell-1}) \odot (1 - \sigma(\boldsymbol{\xi}^{\top} \bar{\mathbf{z}}_{\ell-1}))) (\boldsymbol{\xi}^{\top} \bar{\mathbf{z}}_{\ell-1}) \bar{\mathbf{z}}_{\ell-1}^{\top}$$

Question: Why did we spend all this time computing gradients?

Question: Why did we spend all this time computing gradients?

Answer: The gradient is necessary for gradient descent

Question: Why did we spend all this time computing gradients?

Answer: The gradient is necessary for gradient descent

1: **for** $i \in 1 \dots t$ **do**

2: \triangleright Compute the gradient of the loss

3: $\mathbf{J} \leftarrow \partial \mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) / \partial \boldsymbol{\theta}$

4: \triangleright Update the parameters using the negative gradient

5: $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \mathbf{J}$

Now that we have the gradient, we can do gradient descent

Now that we have the gradient, we can do gradient descent

How do we actually compute this in torch or jax?

Agenda

1. Review
2. Quiz
3. Optimization
4. Calculus review
5. Deriving linear regression
6. Gradient descent
7. Backpropagation
8. **Layer gradient**
9. Full gradient
10. Autograd
11. Practical considerations

Agenda

1. Review
2. Quiz
3. Optimization
4. Calculus review
5. Deriving linear regression
6. Gradient descent
7. Backpropagation
8. Layer gradient
9. **Full gradient**
10. Autograd
11. Practical considerations

Both jax and torch compute gradients using

$$f(\boldsymbol{x}, \boldsymbol{\theta}) = f_{\ell}(\dots f_2(f_1(\boldsymbol{x}, \boldsymbol{\varphi}), \boldsymbol{\psi}) \dots \boldsymbol{\xi})$$

$$f(\mathbf{x}, \boldsymbol{\theta}) = f_{\ell}(\dots f_2(f_1(\mathbf{x}, \boldsymbol{\varphi}), \boldsymbol{\psi}) \dots \boldsymbol{\xi})$$

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \sum_{i=1}^n \left(f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) - \mathbf{y}_{[i]} \right)^2$$

Like before, plug f into the loss function

$$f(\mathbf{x}, \boldsymbol{\theta}) = f_\ell(\dots f_2(f_1(\mathbf{x}, \boldsymbol{\varphi}), \boldsymbol{\psi}) \dots \boldsymbol{\xi})$$

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \sum_{i=1}^n \left(f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) - \mathbf{y}_{[i]} \right)^2$$

Like before, plug f into the loss function

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \sum_{i=1}^n \left(f_\ell(\dots f_2(f_1(\mathbf{x}_{[i]}, \boldsymbol{\varphi}), \boldsymbol{\psi}) \dots \boldsymbol{\xi}) - \mathbf{y}_{[i]} \right)^2$$

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \sum_{i=1}^n \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right) \right) - \mathbf{y}_{[i]} \right)^2$$

We have fully written the loss function, now we take the gradient

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \sum_{i=1}^n \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right) - \mathbf{y}_{[i]} \right)^2$$

We have fully written the loss function, now we take the gradient

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \nabla_{\boldsymbol{\varphi}, \boldsymbol{\psi}, \boldsymbol{\xi}} \sum_{i=1}^n \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right) - \mathbf{y}_{[i]} \right)^2$$

$$\mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \sum_{i=1}^n \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right) - \mathbf{y}_{[i]} \right)^2$$

We have fully written the loss function, now we take the gradient

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \nabla_{\boldsymbol{\varphi}, \boldsymbol{\psi}, \boldsymbol{\xi}} \sum_{i=1}^n \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right) - \mathbf{y}_{[i]} \right)^2$$

Use the sum rule $\frac{d}{dx}(f(x) + g(x)) = f'(x) + g'(x)$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = \sum_{i=1}^n \nabla_{\boldsymbol{\varphi}, \boldsymbol{\psi}, \boldsymbol{\xi}} \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right) - \mathbf{y}_{[i]} \right)^2$$

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{X}, \boldsymbol{Y}, \boldsymbol{\theta}) = \sum_{i=1}^n \nabla_{\boldsymbol{\varphi}, \boldsymbol{\psi}, \boldsymbol{\xi}} \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(\boldsymbol{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right) - \boldsymbol{y}_{[i]} \right)^2$$

$$\nabla_{\theta} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \theta) = \sum_{i=1}^n \nabla_{\varphi, \psi, \xi} \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(x_{[i]}, \varphi \right), \psi \right) \dots \xi \right) - y_{[i]} \right)^2$$

$$\nabla_{\varphi, \psi, \xi} f = \begin{bmatrix} \frac{\partial}{\partial \varphi} f \\ \frac{\partial}{\partial \psi} f \\ \frac{\partial}{\partial \xi} f \end{bmatrix}$$

$$\nabla_{\theta} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \theta) = \sum_{i=1}^n \nabla_{\varphi, \psi, \xi} \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(x_{[i]}, \varphi \right), \psi \right) \dots \xi \right) - y_{[i]} \right)^2$$

$$\nabla_{\varphi, \psi, \xi} f = \begin{bmatrix} \frac{\partial}{\partial \varphi} f \\ \frac{\partial}{\partial \psi} f \\ \frac{\partial}{\partial \xi} f \end{bmatrix}$$

We will determine the gradient for the ξ term first

$$\nabla_{\theta} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \theta) = \sum_{i=1}^n \nabla_{\varphi, \psi, \xi} \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \varphi \right), \psi \right) \dots \xi \right) - \mathbf{y}_{[i]} \right)^2$$

$$\nabla_{\varphi, \psi, \xi} f = \begin{bmatrix} \frac{\partial}{\partial \varphi} f \\ \frac{\partial}{\partial \psi} f \\ \frac{\partial}{\partial \xi} f \end{bmatrix}$$

We will determine the gradient for the ξ term first

$$\nabla_{\theta} \mathcal{L}(\mathbf{X}, \mathbf{Y}, \theta) = \sum_{i=1}^n \nabla_{\xi} \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \varphi \right), \psi \right) \dots \xi \right) - \mathbf{y}_{[i]} \right)^2$$

Power: $\frac{d}{dx}x^n = nx^{n-1}$ Chain: $\frac{d}{dx}f(g(x)) = f'(g(x))g'(x)$

Power: $\frac{d}{dx}x^n = nx^{n-1}$ Chain: $\frac{d}{dx}f(g(x)) = f'(g(x))g'(x)$

$$\begin{aligned}
 &= \sum_{i=1}^n 2 \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right) - \mathbf{y}_{[i]} \right) \\
 &\quad \cdot \nabla_{\boldsymbol{\xi}} \left[f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right) - \mathbf{y}_{[i]} \right] \\
 &= \sum_{i=1}^n 2 \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right) - \mathbf{y}_{[i]} \right) \\
 &\quad \cdot \nabla_{\boldsymbol{\xi}} f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right)
 \end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^n 2 \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right) - \mathbf{y}_{[i]} \right) \\
&\quad \cdot \nabla_{\boldsymbol{\xi}} f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right)
\end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^n 2 \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right) - \mathbf{y}_{[i]} \right) \\
&\quad \cdot \nabla_{\boldsymbol{\xi}} f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right)
\end{aligned}$$

We already know the gradient for the layer f_{ℓ}

$$= \sum_{i=1}^n 2 \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right) - \mathbf{y}_{[i]} \right) \\ \cdot \nabla_{\boldsymbol{\xi}} f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right)$$

We already know the gradient for the layer f_{ℓ}

$$= \sum_{i=1}^n 2 \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right) - \mathbf{y}_{[i]} \right) \\ \cdot \nabla_{\boldsymbol{\xi}} f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right)$$

$$= \sum_{i=1}^n \nabla_{\varphi, \psi, \xi} \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(x_{[i]}, \varphi \right), \psi \right) \dots \xi \right) - y_{[i]} \right)^2$$

$$= \sum_{i=1}^n \nabla_{\varphi, \psi, \xi} \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(x_{[i]}, \varphi \right), \psi \right) \dots \xi \right) - y_{[i]} \right)^2$$

Power: $\frac{d}{dx} x^n = nx^{n-1}$ Chain: $\frac{d}{dx} f(g(x)) = f'(g(x))g'(x)$

$$= \sum_{i=1}^n \nabla_{\varphi, \psi, \xi} \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \varphi \right), \psi \right) \dots \xi \right) - \mathbf{y}_{[i]} \right)^2$$

Power: $\frac{d}{dx} x^n = nx^{n-1}$ Chain: $\frac{d}{dx} f(g(x)) = f'(g(x))g'(x)$

$$= \sum_{i=1}^n 2 \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \varphi \right), \psi \right) \dots \xi \right) - \mathbf{y}_{[i]} \right) \cdot \nabla_{\varphi, \psi, \xi} f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \varphi \right), \psi \right) \dots \xi \right) - \mathbf{y}_{[i]}$$

$$\begin{aligned}
&= \sum_{i=1}^n 2 \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right) - \mathbf{y}_{[i]} \right) \\
&\quad \cdot \nabla_{\boldsymbol{\varphi}, \boldsymbol{\psi}, \boldsymbol{\xi}} f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right) - \mathbf{y}_{[i]}
\end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^n 2 \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right) - \mathbf{y}_{[i]} \right) \\
&\quad \cdot \nabla_{\boldsymbol{\varphi}, \boldsymbol{\psi}, \boldsymbol{\xi}} f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right) - \mathbf{y}_{[i]}
\end{aligned}$$

Use the chain rule again

$$\begin{aligned}
&= \sum_{i=1}^n 2 \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right) - \mathbf{y}_{[i]} \right) \\
&\quad \cdot \nabla_{\boldsymbol{\varphi}, \boldsymbol{\psi}, \boldsymbol{\xi}} f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right) - \mathbf{y}_{[i]}
\end{aligned}$$

Use the chain rule again

$$\begin{aligned}
&= \sum_{i=1}^n 2 \left(f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right) - \mathbf{y}_{[i]} \right) \\
&\quad \cdot \frac{\partial}{\partial \boldsymbol{\xi}} f_{\ell} \left(\dots f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right) \dots \boldsymbol{\xi} \right) \\
&\quad \dots \\
&\quad \cdot \nabla_{\boldsymbol{\varphi}, \boldsymbol{\psi}, \boldsymbol{\xi}} f_2 \left(f_1 \left(\mathbf{x}_{[i]}, \boldsymbol{\varphi} \right), \boldsymbol{\psi} \right)
\end{aligned}$$

And again

$$\begin{aligned} &= \sum_{i=1}^n 2 \left(f_{\ell}(\dots f_2(f_1(x, \varphi), \psi) \dots \xi) - y_{[i]} \right) \\ &\quad \cdot \frac{\partial}{\partial \xi} f_{\ell}(\dots f_2(f_1(x, \varphi), \psi) \dots \xi) \\ &\quad \dots \\ &\quad \cdot \frac{\partial}{\partial \psi} f_2(f_1(x, \varphi), \psi) \\ &\quad \cdot \nabla_{\varphi, \psi, \xi} f_1(x, \varphi) \end{aligned}$$

And again

$$\begin{aligned} &= \sum_{i=1}^n 2 \left(f_{\ell}(\dots f_2(f_1(\mathbf{x}, \boldsymbol{\varphi}), \boldsymbol{\psi}) \dots \boldsymbol{\xi}) - \mathbf{y}_{[i]} \right) \\ &\quad \cdot \frac{\partial}{\partial \boldsymbol{\xi}} f_{\ell}(\dots f_2(f_1(\mathbf{x}, \boldsymbol{\varphi}), \boldsymbol{\psi}) \dots \boldsymbol{\xi}) \\ &\quad \dots \\ &\quad \cdot \frac{\partial}{\partial \boldsymbol{\psi}} f_2(f_1(\mathbf{x}, \boldsymbol{\varphi}), \boldsymbol{\psi}) \\ &\quad \cdot \frac{\partial}{\partial \boldsymbol{\varphi}} f_1(\mathbf{x}, \boldsymbol{\varphi}) \end{aligned}$$

This is the gradient for a deep neural network

$$\begin{aligned} &= \sum_{i=1}^n 2 \left(f_{\ell}(\dots f_2(f_1(x, \varphi), \psi) \dots \xi) - y_{[i]} \right) \\ &\quad \cdot \frac{\partial}{\partial \xi} f_{\ell}(\dots f_2(f_1(x, \varphi), \psi) \dots \xi) \\ &\quad \dots \\ &\quad \cdot \frac{\partial}{\partial \psi} f_2(f_1(x, \varphi), \psi) \\ &\quad \cdot \frac{\partial}{\partial \varphi} f_1(x, \varphi) \end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^n 2 \left(f_{\ell}(\dots f_2(f_1(x, \varphi), \psi) \dots \xi) - y_{[i]} \right) \\
&\quad \cdot \frac{\partial}{\partial \xi} f_{\ell}(\dots f_2(f_1(x, \varphi), \psi) \dots \xi) \\
&\quad \dots \\
&\quad \cdot \frac{\partial}{\partial \psi} f_2(f_1(x, \varphi), \psi) \\
&\quad \cdot \frac{\partial}{\partial \varphi} f_1(x, \varphi)
\end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^n 2 \left(f_{\ell}(\dots f_2(f_1(x, \varphi), \psi) \dots \xi) - y_{[i]} \right) \\
&\quad \cdot \frac{\partial}{\partial \xi} f_{\ell}(\dots f_2(f_1(x, \varphi), \psi) \dots \xi) \\
&\quad \dots \\
&\quad \cdot \frac{\partial}{\partial \psi} f_2(f_1(x, \varphi), \psi) \\
&\quad \cdot \frac{\partial}{\partial \varphi} f_1(x, \varphi)
\end{aligned}$$

When we call `torch.backward()` or `jax.grad`, they compute this

With a **forward** pass, we use inputs and compute the loss

With a **forward** pass, we use inputs and compute the loss

With a **backward** pass, we compute the gradients from the loss

The representational power of neural networks comes from the nonlinear activation function σ

The representational power of neural networks comes from the nonlinear activation function σ

Unfortunately, this also means we do not have a closed-form solution to find the parameters θ

The representational power of neural networks comes from the nonlinear activation function σ

Unfortunately, this also means we do not have a closed-form solution to find the parameters θ

We need to find another way to find θ

The representational power of neural networks comes from the nonlinear activation function σ

Unfortunately, this also means we do not have a closed-form solution to find the parameters θ

We need to find another way to find θ

We will use **gradient descent** to find θ

Gradient descent is an optimization method for **differentiable** functions

Gradient descent is an optimization method for **differentiable** functions

More formally, gradient descent approximates the θ that solves

$$\arg \min_{\theta} \mathcal{L}(x, y, \theta)$$

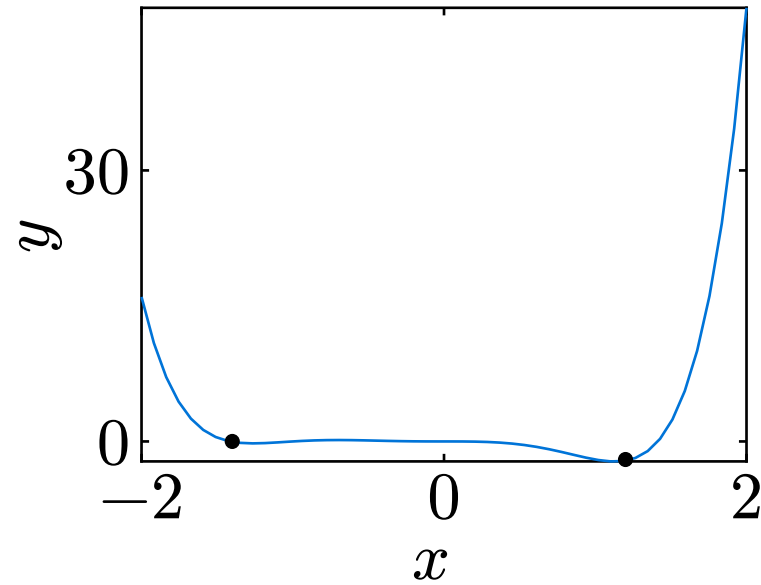
Gradient descent is an optimization method for **differentiable** functions

More formally, gradient descent approximates the θ that solves

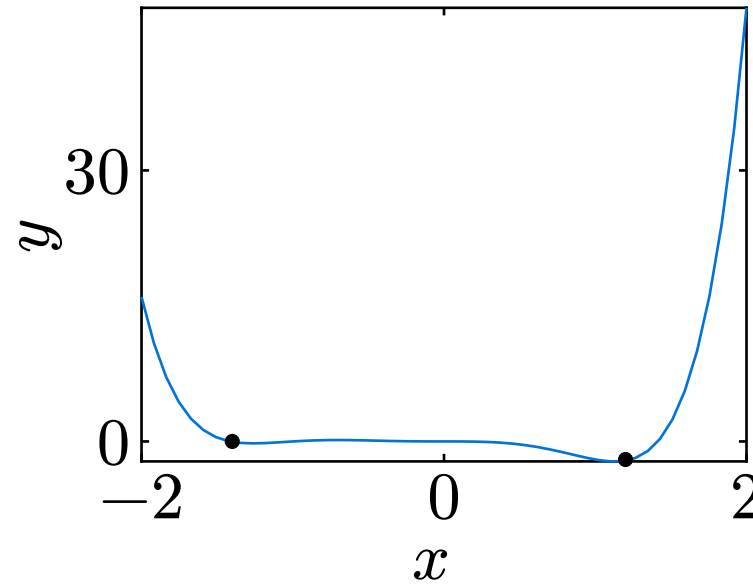
$$\arg \min_{\theta} \mathcal{L}(x, y, \theta)$$

Gradient descent provides a **local** optima, not necessarily a **global** optima

Gradient descent provides a **local** optima, not necessarily a **global** optima



Gradient descent provides a **local** optima, not necessarily a **global** optima



In practice, a local optima provides a good enough model

Gradient descent relies on the **gradient** of \mathcal{L} , therefore \mathcal{L} must be differentiable

Gradient descent relies on the **gradient** of \mathcal{L} , therefore \mathcal{L} must be differentiable

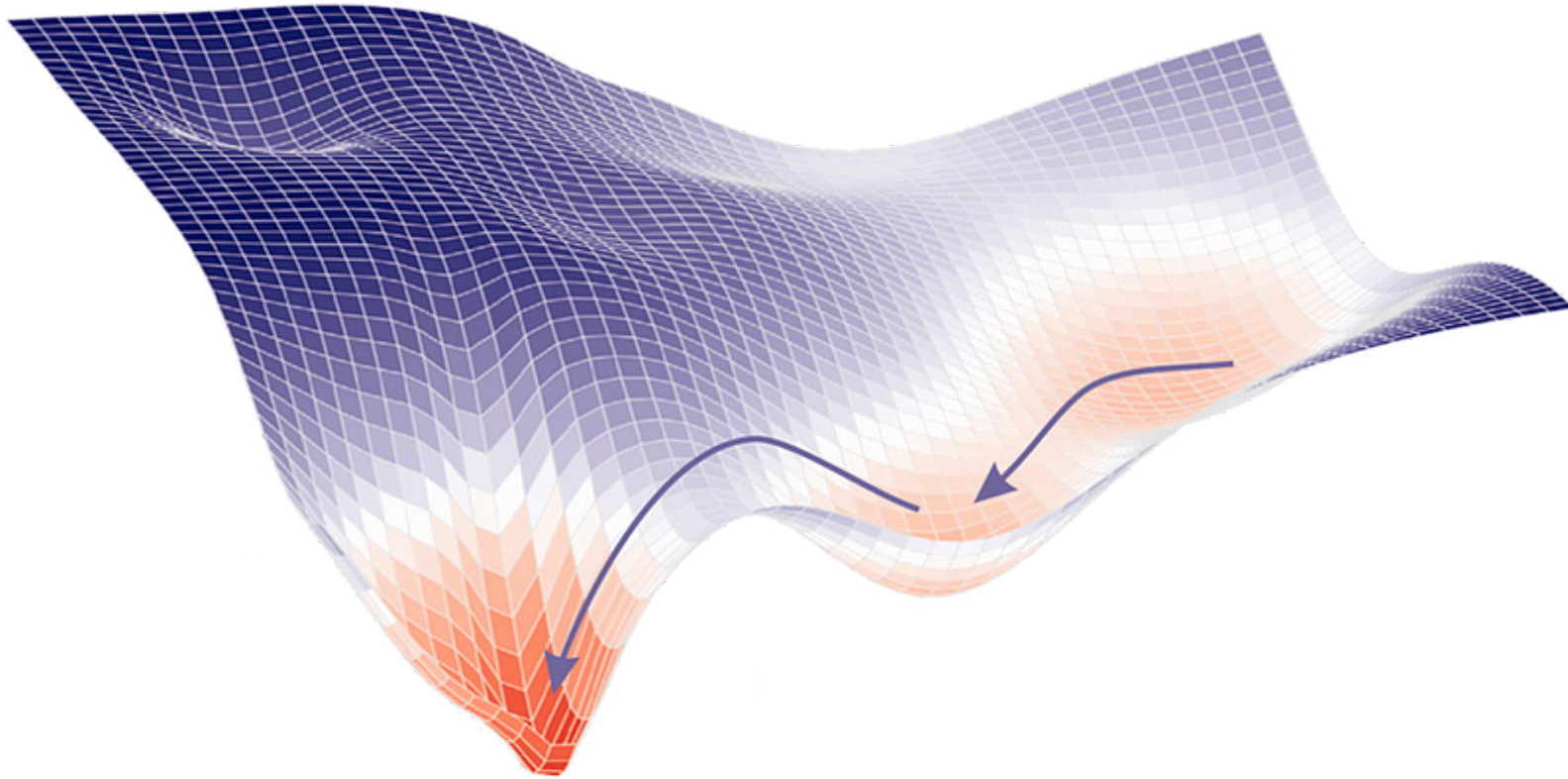
In gradient descent, we update the parameters in the direction of the negative gradient of the loss

The gradient descent algorithm is as follows:

```
1: function GRADIENT DESCENT( $\mathbf{X}, \mathbf{Y}, \mathcal{L}, t, \alpha$ )
2:     ▷ Randomly initialize parameters
3:      $\boldsymbol{\theta} \leftarrow \mathcal{N}(0, 1)$ 
4:     for  $i \in 1 \dots t$  do
5:         ▷ Compute the gradient of the loss
6:          $\mathbf{J} \leftarrow \partial \mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) / \partial \boldsymbol{\theta}$ 
7:         ▷ Update the parameters using the negative gradient
8:          $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \mathbf{J}$ 
9:     return  $\boldsymbol{\theta}$ 
```

We can visualize gradient descent of a bivariate function

We can visualize gradient descent of a bivariate function



Let us do a very simple example

Let us do a very simple example

Goal: Solve $\arg \min_{\theta} \mathcal{L}(x, y, \theta) = \theta_1^2 x + \theta_0$ using gradient descent

Dataset: One item in the dataset, $x = 1, y = 2$

1. Randomly initialize θ

Let us do a very simple example

Goal: Solve $\arg \min_{\theta} \mathcal{L}(x, y, \theta) = \theta_1^2 x + \theta_0$ using gradient descent

Dataset: One item in the dataset, $x = 1, y = 2$

1. Randomly initialize θ

$$\theta = [2.5 \ 1.1]^\top$$

Let us do a very simple example

Goal: Solve $\arg \min_{\theta} \mathcal{L}(x, y, \theta) = \theta_1^2 x + \theta_0$ using gradient descent

Dataset: One item in the dataset, $x = 1, y = 2$

1. Randomly initialize θ

$$\theta = [2.5 \ 1.1]^\top$$

2. Compute the gradient for our dataset

$$J = \frac{\partial \mathcal{L}(x, y, \theta)}{\partial \theta}$$

$$\mathbf{J} = \frac{\partial \mathcal{L}(x, y, \theta)}{\partial \theta}$$

From the last slide

$$\mathbf{J} = \frac{\partial \mathcal{L}(x, y, \theta)}{\partial \theta}$$

From the last slide

$$\mathbf{J} = \left[\frac{\partial}{\partial \theta_1} \mathcal{L}(x, y, \theta) \quad \frac{\partial}{\partial \theta_0} \mathcal{L}(x, y, \theta) \right]^\top$$

Write out gradient w.r.t. all
parameters

$$\mathbf{J} = \frac{\partial \mathcal{L}(x, y, \theta)}{\partial \theta}$$

From the last slide

$$\mathbf{J} = \left[\frac{\partial}{\partial \theta_1} \mathcal{L}(x, y, \theta) \quad \frac{\partial}{\partial \theta_0} \mathcal{L}(x, y, \theta) \right]^\top$$

Write out gradient w.r.t. all parameters

$$\mathbf{J} = \left[\frac{\partial}{\partial \theta_1} (\theta_1^2 x + \theta_0) \quad \frac{\partial}{\partial \theta_0} (\theta_1^2 x + \theta_0) \right]^\top$$

Plug in \mathcal{L}

$$\mathbf{J} = [2\theta_1 x \quad 1]^\top$$

Differentiate

$$\mathbf{J} = [2 \cdot 2.5 \cdot 1 \quad 1]^\top = [5 \quad 1]^\top$$

Plug in x, y, θ

3. Update θ using J

$$\theta \leftarrow \theta - \alpha J$$

From algorithm

$$\begin{bmatrix} \theta_1 \\ \theta_0 \end{bmatrix} \leftarrow \begin{bmatrix} 2.5 \\ 1.1 \end{bmatrix} - \alpha \begin{bmatrix} 5 \\ 1 \end{bmatrix}$$

Plug in θ and J

$$\begin{bmatrix} \theta_1 \\ \theta_0 \end{bmatrix} \leftarrow \begin{bmatrix} 2.5 \\ 1.1 \end{bmatrix} - 0.1 \begin{bmatrix} 5 \\ 1 \end{bmatrix}$$

Let $\alpha = 0.1$

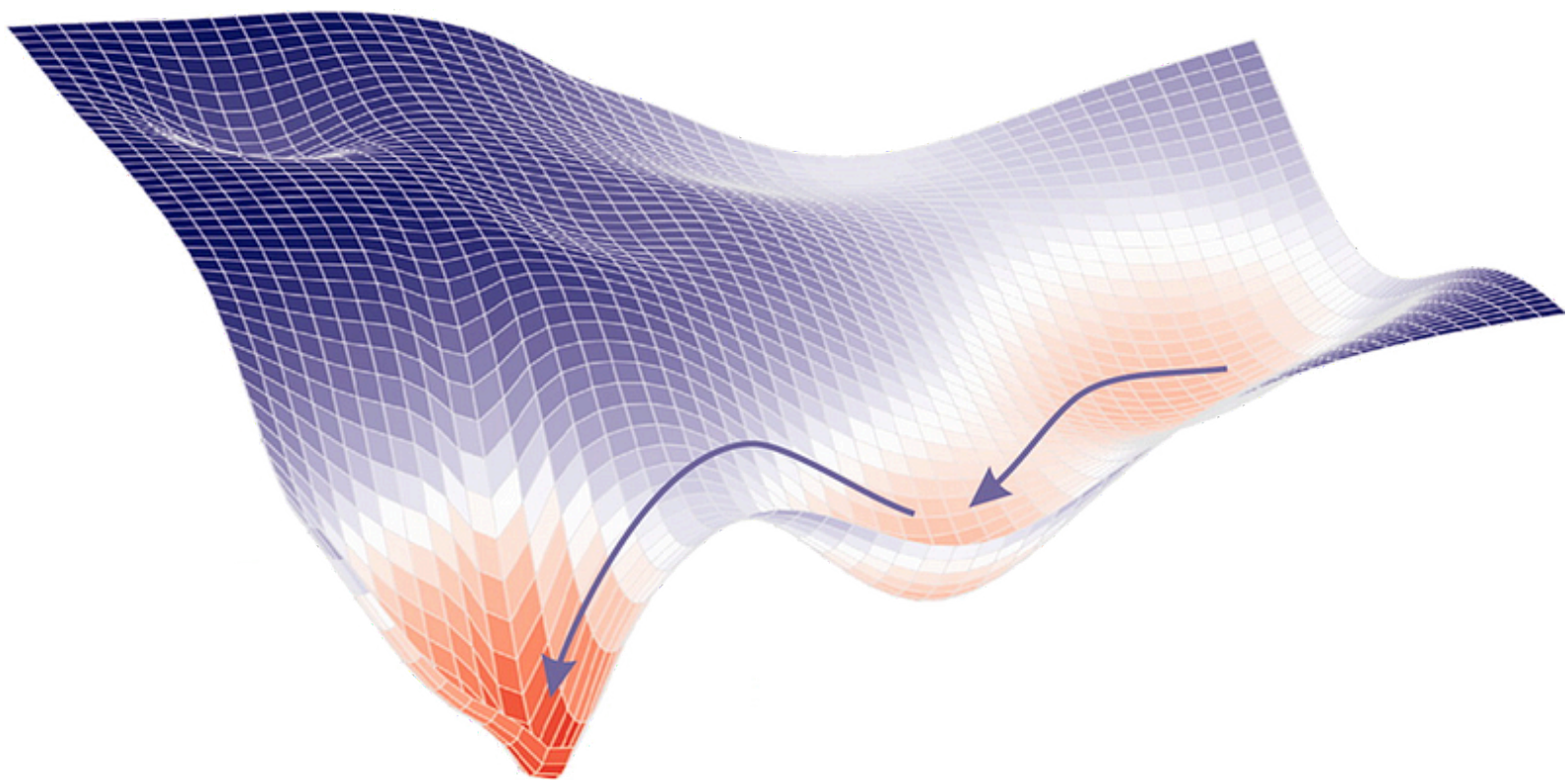
$$\begin{bmatrix} \theta_1 \\ \theta_0 \end{bmatrix} \leftarrow \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

Evaluate

4. Repeat this process until convergence (loss no longer decreases)

$$\mathbf{J} = \frac{\partial \mathcal{L}(x, y, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \mathbf{J}$$



TODO: Summary

There are different types of gradient descent, depending on how much training data we use:

Batch gradient descent: uses the entire dataset $\mathcal{L} \left(\begin{bmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_n \end{bmatrix}, \begin{bmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_n \end{bmatrix}, \boldsymbol{\theta} \right)$

Stochastic gradient descent: Gradient descent over one datapoint at a time $\mathcal{L}(\mathbf{x}_i, \mathbf{y}_i, \boldsymbol{\theta})$

Minibatch gradient descent: Gradient descent over many (but not all) datapoints $\mathcal{L} \left(\begin{bmatrix} \mathbf{x}_i \\ \vdots \\ \mathbf{x}_j \end{bmatrix}, \begin{bmatrix} \mathbf{y}_i \\ \vdots \\ \mathbf{y}_j \end{bmatrix}, \boldsymbol{\theta} \right)$

In practice, we do not have enough computer memory for batch gradient descent

In practice, we do not have enough computer memory for batch gradient descent

Instead, we use stochastic gradient descent or minibatch gradient descent

In practice, we do not have enough computer memory for batch gradient descent

Instead, we use stochastic gradient descent or minibatch gradient descent