

Modern Techniques

CISC 7026: Introduction to Deep Learning

University of Macau

Agenda

1. Review
2. Dirty secret of deep learning
3. Optimization is hard
4. Deeper neural networks
5. Activation functions
6. Parameter initialization
7. Stochastic gradient descent
8. Modern optimization
9. Coding

1. Review
2. Dirty secret of deep learning
3. Optimization is hard
4. Deeper neural networks
5. Activation functions
6. Parameter initialization
7. Stochastic gradient descent
8. Modern optimization
9. Coding

Many problems in ML can be reduced to **regression** or **classification**

Many problems in ML can be reduced to **regression** or **classification**

Regression asks how many

Many problems in ML can be reduced to **regression** or **classification**

Regression asks how many

- How long will I live?

Many problems in ML can be reduced to **regression** or **classification**

Regression asks how many

- How long will I live?
- How much rain will there be tomorrow?

Many problems in ML can be reduced to **regression** or **classification**

Regression asks how many

- How long will I live?
- How much rain will there be tomorrow?
- How far away is this object?

Many problems in ML can be reduced to **regression** or **classification**

Regression asks how many

- How long will I live?
- How much rain will there be tomorrow?
- How far away is this object?

Classification asks which one

Many problems in ML can be reduced to **regression** or **classification**

Regression asks how many

- How long will I live?
- How much rain will there be tomorrow?
- How far away is this object?

Classification asks which one

- Is this a dog or muffin?

Many problems in ML can be reduced to **regression** or **classification**

Regression asks how many

- How long will I live?
- How much rain will there be tomorrow?
- How far away is this object?

Classification asks which one

- Is this a dog or muffin?
- Will it rain tomorrow? Yes or no?

Many problems in ML can be reduced to **regression** or **classification**

Regression asks how many

- How long will I live?
- How much rain will there be tomorrow?
- How far away is this object?

Classification asks which one

- Is this a dog or muffin?
- Will it rain tomorrow? Yes or no?
- What color is this object?

Many problems in ML can be reduced to **regression** or **classification**

Regression asks how many

- How long will I live?
- How much rain will there be tomorrow?
- How far away is this object?

Classification asks which one

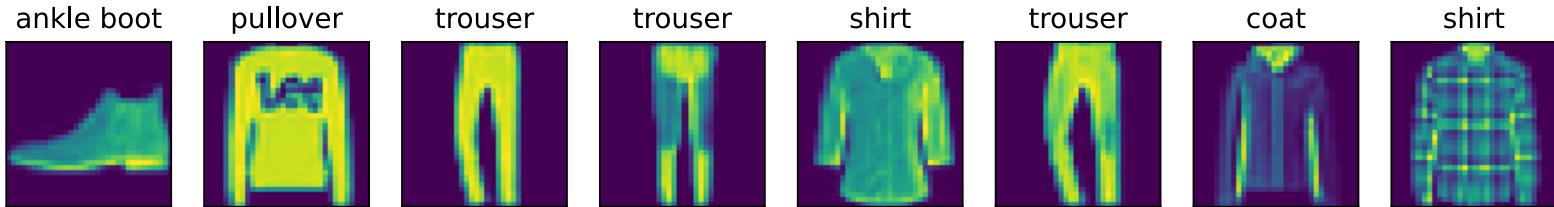
- Is this a dog or muffin?
- Will it rain tomorrow? Yes or no?
- What color is this object?

So far, we only looked at regression. Now, let us look at classification

Task: Given a picture of clothes, predict the text description

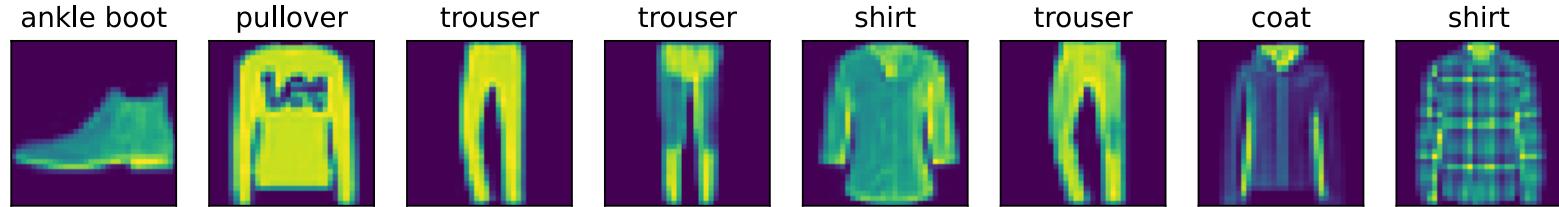
Task: Given a picture of clothes, predict the text description

$$X : \mathbb{Z}_{0,255}^{32 \times 32}$$



Task: Given a picture of clothes, predict the text description

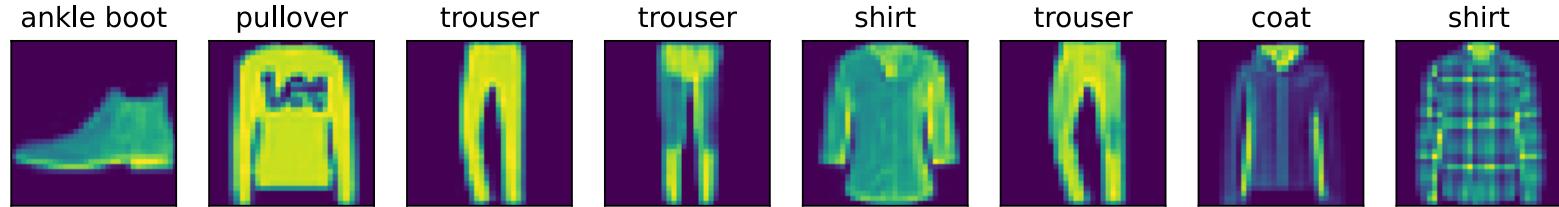
$$X : \mathbb{Z}_{0,255}^{32 \times 32}$$



$Y : \{\text{T-shirt, Trouser, Pullover, Dress, Coat, Sandal, Shirt, Sneaker, Bag, Ankle boot}\}$

Task: Given a picture of clothes, predict the text description

$$X : \mathbb{Z}_{0,255}^{32 \times 32}$$



$$Y : \{\text{T-shirt, Trouser, Pullover, Dress, Coat, Sandal, Shirt, Sneaker, Bag, Ankle boot}\}$$

Approach: Learn θ that produce **conditional probabilities**

Task: Given a picture of clothes, predict the text description


$$Y : \{\text{T-shirt, Trouser, Pullover, Dress, Coat, Sandal, Shirt, Sneaker, Bag, Ankle boot}\}$$

Approach: Learn θ that produce **conditional probabilities**

$$f(x, \theta) = P(y \mid x) = P\left(\begin{bmatrix} \text{T-Shirt} \\ \text{Trouser} \\ \vdots \end{bmatrix} \mid \begin{array}{c} \text{[Image of a pullover]} \\ \text{[Image of a pullover]} \end{array}\right) = \begin{bmatrix} 0.2 \\ 0.01 \\ \vdots \end{bmatrix}$$

If events A, B are not disjoint, they are **conditionally dependent**

If events A, B are not disjoint, they are **conditionally dependent**

$$P(\text{cloud}) = 0.2, P(\text{rain}) = 0.1$$

Review

If events A, B are not disjoint, they are **conditionally dependent**

$$P(\text{cloud}) = 0.2, P(\text{rain}) = 0.1$$

$$P(\text{rain} \mid \text{cloud}) = 0.5$$

Review

If events A, B are not disjoint, they are **conditionally dependent**

$$P(\text{cloud}) = 0.2, P(\text{rain}) = 0.1$$

$$P(\text{rain} \mid \text{cloud}) = 0.5$$

$$P(A \mid B) = \frac{P(A \cap B)}{P(B)}$$

Review

If events A, B are not disjoint, they are **conditionally dependent**

$$P(\text{cloud}) = 0.2, P(\text{rain}) = 0.1$$

$$P(\text{rain} \mid \text{cloud}) = 0.5$$

$$P(A \mid B) = \frac{P(A \cap B)}{P(B)}$$

$$P(\text{Rain} \cap \text{Cloud}) = 0.1$$

Walk outside

$$P(\text{Cloud}) = 0.2$$

$$P(\text{Rain} \mid \text{Cloud}) = \frac{0.1}{0.2} = 0.5$$

How can we represent a probability distribution for a neural network?

Review

How can we represent a probability distribution for a neural network?

$$\boldsymbol{v} = \left\{ \begin{bmatrix} v_1 \\ \vdots \\ v_{d_y} \end{bmatrix} \mid \sum_{i=1}^{d_y} v_i = 1; \quad v_i \in (0, 1) \right\}$$

How can we represent a probability distribution for a neural network?

$$\boldsymbol{v} = \left\{ \begin{bmatrix} v_1 \\ \vdots \\ v_{d_y} \end{bmatrix} \mid \sum_{i=1}^{d_y} v_i = 1; \quad v_i \in (0, 1) \right\}$$

There is special notation for this vector, called the **simplex**

How can we represent a probability distribution for a neural network?

$$\boldsymbol{v} = \left\{ \begin{bmatrix} v_1 \\ \vdots \\ v_{d_y} \end{bmatrix} \mid \sum_{i=1}^{d_y} v_i = 1; \quad v_i \in (0, 1) \right\}$$

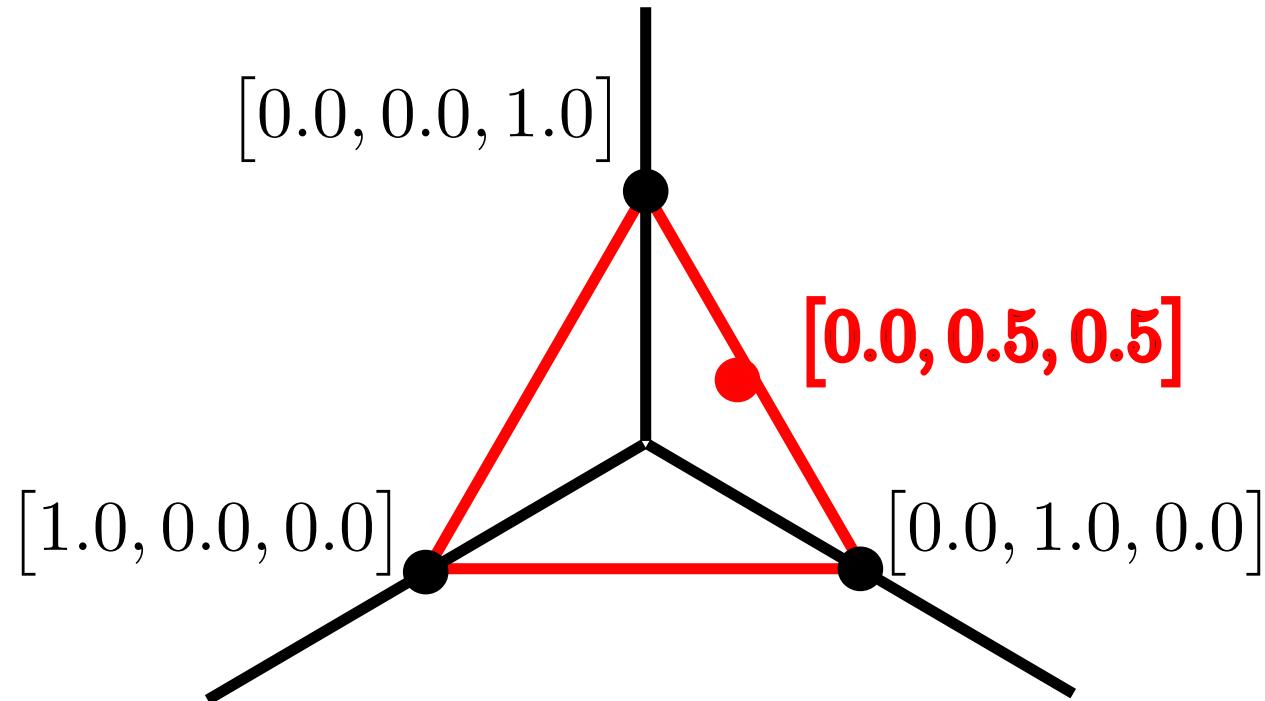
There is special notation for this vector, called the **simplex**

$$\Delta^{d_y-1}$$

The simplex Δ^k is an $k - 1$ -dimensional triangle in k -dimensional space

Review

The simplex Δ^k is an $k - 1$ -dimensional triangle in k -dimensional space



It has only $k - 1$ free variables, because $x_k = 1 - \sum_{i=1}^{k-1} x_i$

The softmax function maps real numbers to the simplex (probabilities)

$$\text{softmax} : \mathbb{R}^k \mapsto \Delta^{k-1}$$

The softmax function maps real numbers to the simplex (probabilities)

$$\text{softmax} : \mathbb{R}^k \mapsto \Delta^{k-1}$$

$$\text{softmax}\left(\begin{bmatrix} x_1 \\ \vdots \\ x_k \end{bmatrix}\right) = \frac{e^{\mathbf{x}}}{\sum_{i=1}^k e^{x_i}} = \begin{bmatrix} \frac{e^{x_1}}{e^{x_1} + e^{x_2} + \dots + e^{x_k}} \\ \frac{e^{x_2}}{e^{x_1} + e^{x_2} + \dots + e^{x_k}} \\ \vdots \\ \frac{e^{x_k}}{e^{x_1} + e^{x_2} + \dots + e^{x_k}} \end{bmatrix}$$

The softmax function maps real numbers to the simplex (probabilities)

$$\text{softmax} : \mathbb{R}^k \mapsto \Delta^{k-1}$$

$$\text{softmax}\left(\begin{bmatrix} x_1 \\ \vdots \\ x_k \end{bmatrix}\right) = \frac{e^{\mathbf{x}}}{\sum_{i=1}^k e^{x_i}} = \begin{bmatrix} \frac{e^{x_1}}{e^{x_1} + e^{x_2} + \dots + e^{x_k}} \\ \frac{e^{x_2}}{e^{x_1} + e^{x_2} + \dots + e^{x_k}} \\ \vdots \\ \frac{e^{x_k}}{e^{x_1} + e^{x_2} + \dots + e^{x_k}} \end{bmatrix}$$

If we attach it to our linear model, we can output probabilities!

$$f(\mathbf{x}, \boldsymbol{\theta}) = \text{softmax}(\boldsymbol{\theta}^\top \mathbf{x})$$

Question: Why do we output probabilities instead of binary values

$$f(x, \theta) = \begin{bmatrix} P(\text{Shirt} \mid \text{Image}) \\ P(\text{Bag} \mid \text{Image}) \\ \vdots \end{bmatrix} = \begin{bmatrix} 0.25 \\ 0.08 \\ \vdots \end{bmatrix}; \quad f(x, \theta) = \begin{bmatrix} 1 \\ 0 \\ \vdots \end{bmatrix}$$

Question: Why do we output probabilities instead of binary values

$$f(x, \theta) = \begin{bmatrix} P(\text{Shirt} | \text{Image}) \\ P(\text{Bag} | \text{Image}) \\ \vdots \end{bmatrix} = \begin{bmatrix} 0.25 \\ 0.08 \\ \vdots \end{bmatrix}; \quad f(x, \theta) = \begin{bmatrix} 1 \\ 0 \\ \vdots \end{bmatrix}$$

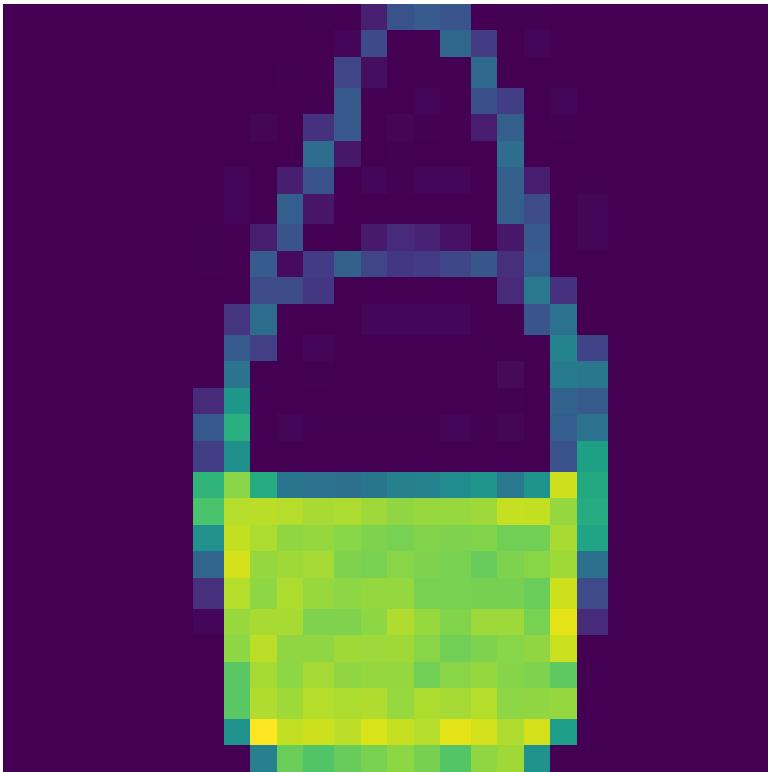
Answer 1: Outputting probabilities results in differentiable functions

Question: Why do we output probabilities instead of binary values

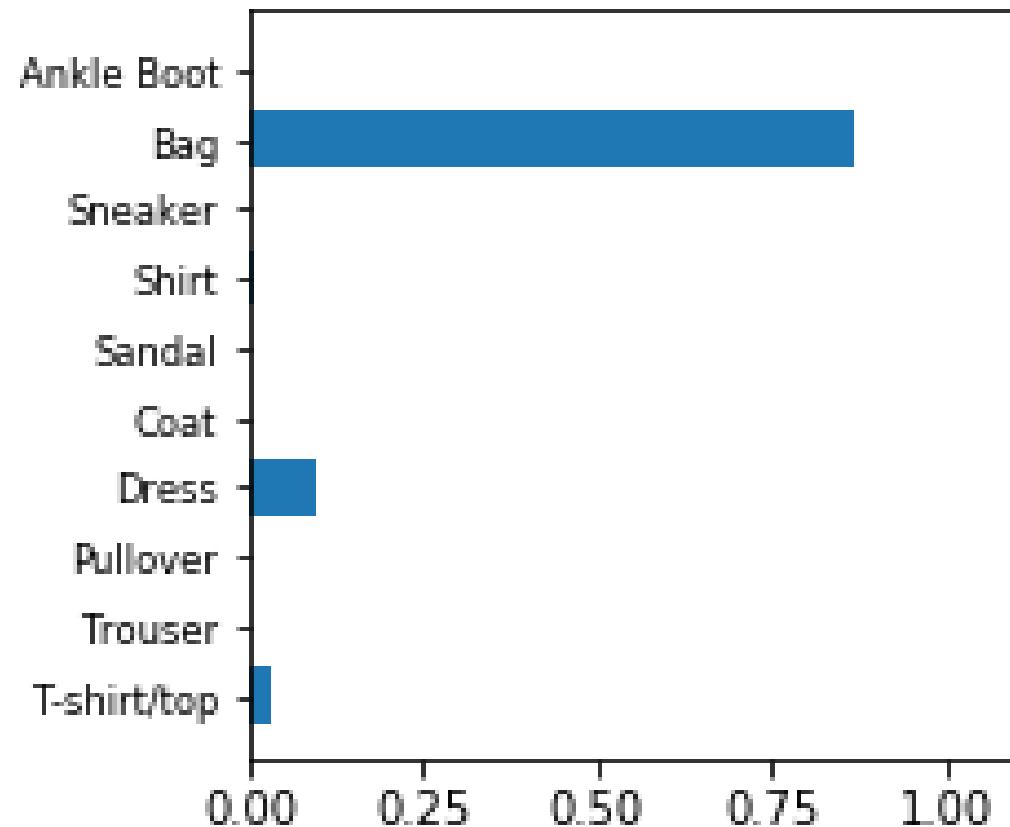
$$f(x, \theta) = \begin{bmatrix} P(\text{Shirt} \mid \text{Image}) \\ P(\text{Bag} \mid \text{Image}) \\ \vdots \end{bmatrix} = \begin{bmatrix} 0.25 \\ 0.08 \\ \vdots \end{bmatrix}; \quad f(x, \theta) = \begin{bmatrix} 1 \\ 0 \\ \vdots \end{bmatrix}$$

Answer 1: Outputting probabilities results in differentiable functions

Answer 2: We report uncertainty, which is useful in many applications



Class Probability



We consider the label $y_{[i]}$ as a conditional distribution

$$P(y_{[i]} \mid x_{[i]}) = \begin{bmatrix} P(\text{Shirt} \mid \text{Image}) \\ P(\text{Bag} \mid \text{Image}) \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

Review

We consider the label $y_{[i]}$ as a conditional distribution

$$P(y_{[i]} \mid x_{[i]}) = \begin{bmatrix} P(\text{Shirt} \mid \text{Image}) \\ P(\text{Bag} \mid \text{Image}) \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

What loss function should we use for classification?

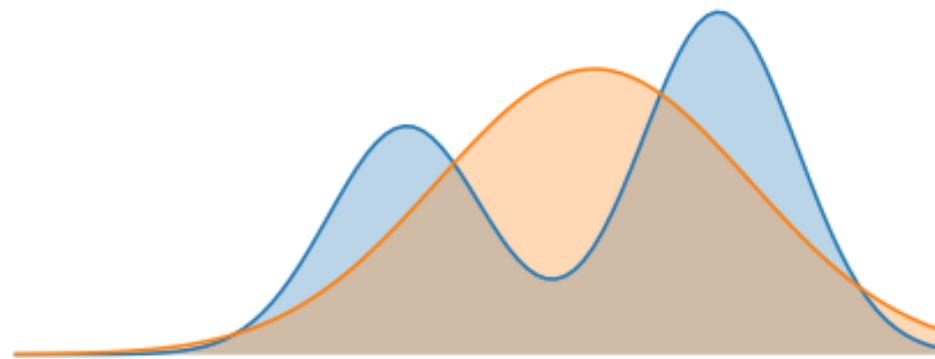
Since $f(x, \theta)$ and $P(y \mid x)$ are both distributions, we want to measure the difference between distributions

Since $f(x, \theta)$ and $P(y | x)$ are both distributions, we want to measure the difference between distributions

One measurement is the **Kullback-Leibler Divergence (KL)**

Since $f(x, \theta)$ and $P(y | x)$ are both distributions, we want to measure the difference between distributions

One measurement is the **Kullback-Leibler Divergence (KL)**



From the KL divergence, we derived the **cross-entropy loss** function, which we use for classification

Review

From the KL divergence, we derived the **cross-entropy loss** function, which we use for classification

$$= - \sum_{i=1}^{d_y} P(y_i \mid \boldsymbol{x}) \log f(\boldsymbol{x}, \boldsymbol{\theta})_i$$

Review

From the KL divergence, we derived the **cross-entropy loss** function, which we use for classification

$$= - \sum_{i=1}^{d_y} P(y_i \mid \mathbf{x}) \log f(\mathbf{x}, \boldsymbol{\theta})_i$$

$$\mathcal{L}(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) = \left[- \sum_{j=1}^n \sum_{i=1}^{d_y} P(y_{[j],i} \mid \mathbf{x}_{[j]}) \log f(\mathbf{x}_{[j]}, \boldsymbol{\theta})_i \right]$$

Finish coding exercise

https://colab.research.google.com/drive/1BGMIE2CjlLJOH-D2r9AariPDVgxjWlqG#scrollTo=AnHP-PHVhpW_

1. Review
2. Dirty secret of deep learning
3. Optimization is hard
4. Deeper neural networks
5. Activation functions
6. Parameter initialization
7. Stochastic gradient descent
8. Modern optimization
9. Coding

1. Review
2. **Dirty secret of deep learning**
3. Optimization is hard
4. Deeper neural networks
5. Activation functions
6. Parameter initialization
7. Stochastic gradient descent
8. Modern optimization
9. Coding

So far, I gave you the impression that deep learning is rigorous

Dirty secret of deep learning

So far, I gave you the impression that deep learning is rigorous

Biological inspiration, theoretical bounds and mathematical guarantees

Dirty secret of deep learning

So far, I gave you the impression that deep learning is rigorous

Biological inspiration, theoretical bounds and mathematical guarantees

For complex neural networks, deep learning is a **science** not **math**

Dirty secret of deep learning

So far, I gave you the impression that deep learning is rigorous

Biological inspiration, theoretical bounds and mathematical guarantees

For complex neural networks, deep learning is a **science** not **math**

There is no widely-accepted theory for why deep neural networks are so effective

Dirty secret of deep learning

So far, I gave you the impression that deep learning is rigorous

Biological inspiration, theoretical bounds and mathematical guarantees

For complex neural networks, deep learning is a **science** not **math**

There is no widely-accepted theory for why deep neural networks are so effective

In modern deep learning, we progress using trial and error

Dirty secret of deep learning

So far, I gave you the impression that deep learning is rigorous

Biological inspiration, theoretical bounds and mathematical guarantees

For complex neural networks, deep learning is a **science** not **math**

There is no widely-accepted theory for why deep neural networks are so effective

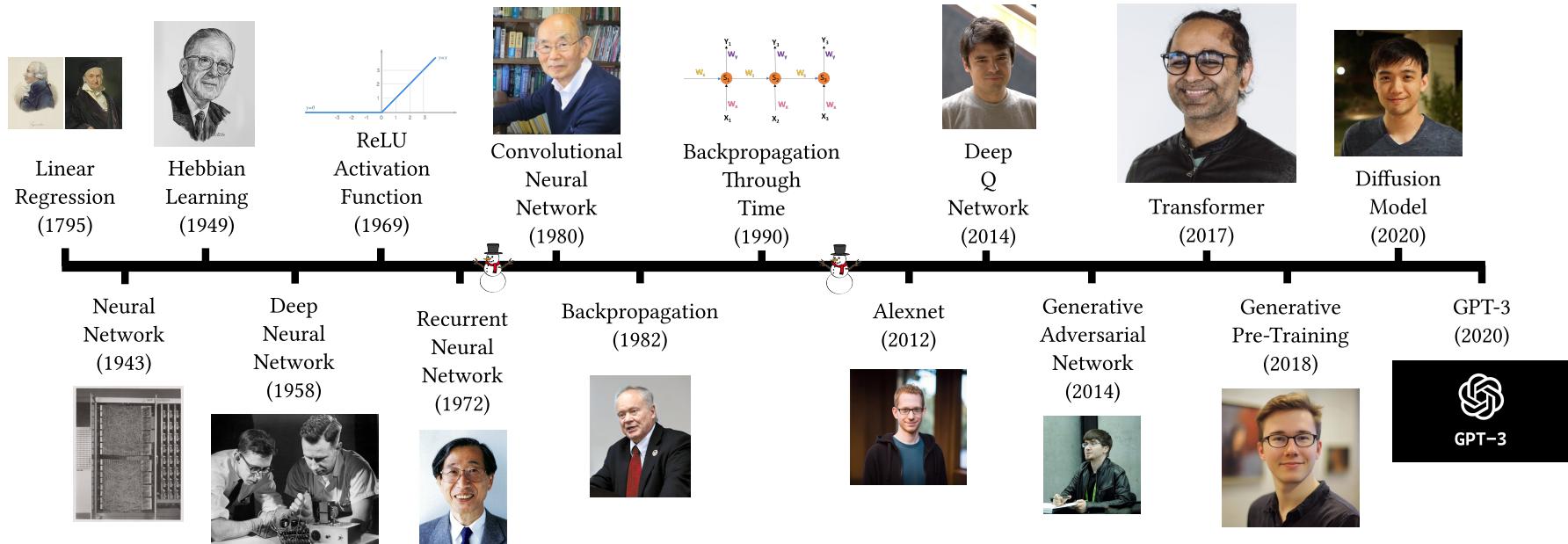
In modern deep learning, we progress using trial and error

Today we experiment, and maybe tomorrow we discover the theory

Similar to using neural networks for 40 years without knowing how to train them

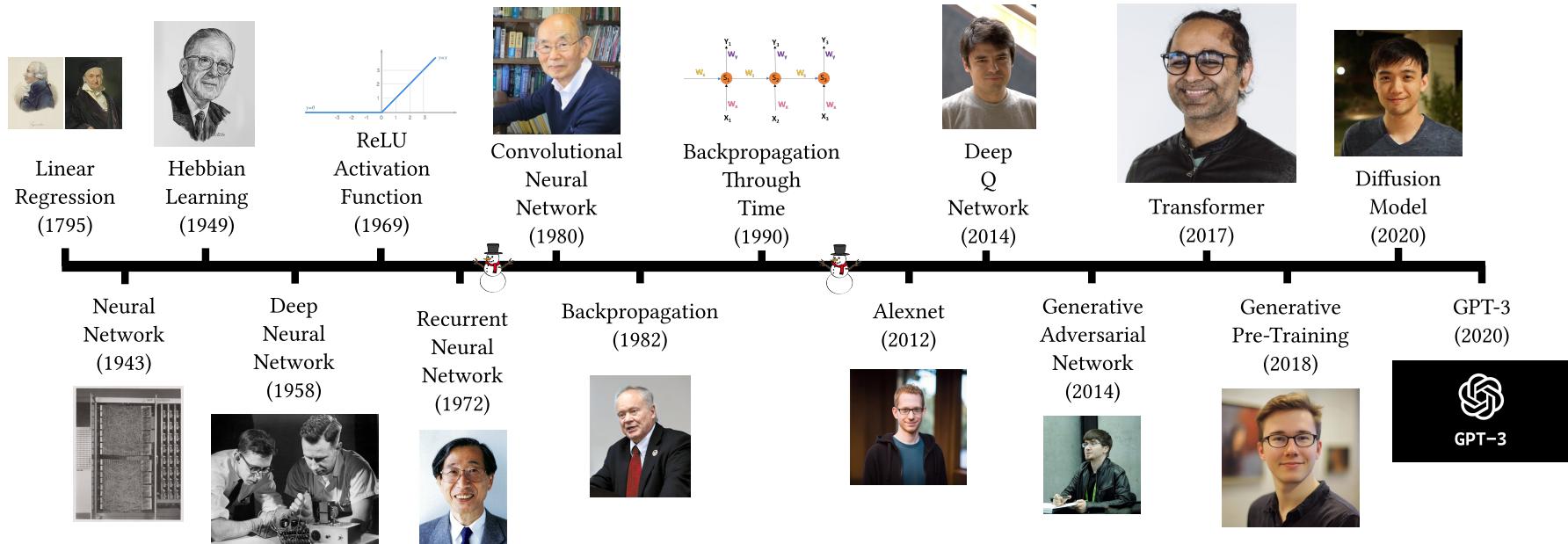
Dirty secret of deep learning

Similar to using neural networks for 40 years without knowing how to train them



Dirty secret of deep learning

Similar to using neural networks for 40 years without knowing how to train them



Are modern networks are too complex for humans to understand?

Scientific method:

Dirty secret of deep learning

Scientific method:

1. Collect observations

Dirty secret of deep learning

Scientific method:

1. Collect observations
2. Form hypothesis

Dirty secret of deep learning

Scientific method:

1. Collect observations
2. Form hypothesis
3. Run experiment

Dirty secret of deep learning

Scientific method:

1. Collect observations
2. Form hypothesis
3. Run experiment
4. Publish theory

Dirty secret of deep learning

Scientific method:

1. Collect observations
2. Form hypothesis
3. Run experiment
4. Publish theory

However, there is a second part:

Dirty secret of deep learning

Scientific method:

1. Collect observations
2. Form hypothesis
3. Run experiment
4. Publish theory

However, there is a second part:

1. Find theory

Dirty secret of deep learning

Scientific method:

1. Collect observations
2. Form hypothesis
3. Run experiment
4. Publish theory

However, there is a second part:

1. Find theory
2. Find counterexample

Dirty secret of deep learning

Scientific method:

1. Collect observations
2. Form hypothesis
3. Run experiment
4. Publish theory

However, there is a second part:

1. Find theory
2. Find counterexample
3. Publish counterexample

Dirty secret of deep learning

Scientific method:

1. Collect observations
2. Form hypothesis
3. Run experiment
4. Publish theory

However, there is a second part:

1. Find theory
2. Find counterexample
3. Publish counterexample
4. Falsify theory

Deep learning is new, so much of part 2 has not happened yet!

For many concepts, the **observations** are stronger than the **theory**

Dirty secret of deep learning

For many concepts, the **observations** are stronger than the **theory**

Observe that a concept improves many types of neural networks

Dirty secret of deep learning

For many concepts, the **observations** are stronger than the **theory**

Observe that a concept improves many types of neural networks

Then, try to create a theory

Dirty secret of deep learning

For many concepts, the **observations** are stronger than the **theory**

Observe that a concept improves many types of neural networks

Then, try to create a theory

Often, these theories are incomplete

Dirty secret of deep learning

For many concepts, the **observations** are stronger than the **theory**

Observe that a concept improves many types of neural networks

Then, try to create a theory

Often, these theories are incomplete

If you do not believe the theory, prove it wrong and be famous!

Dirty secret of deep learning

For many concepts, the **observations** are stronger than the **theory**

Observe that a concept improves many types of neural networks

Then, try to create a theory

Often, these theories are incomplete

If you do not believe the theory, prove it wrong and be famous!

Even if we do not agree on **why** a concept works, if we **observe** that it helps, we can still use it

Dirty secret of deep learning

For many concepts, the **observations** are stronger than the **theory**

Observe that a concept improves many types of neural networks

Then, try to create a theory

Often, these theories are incomplete

If you do not believe the theory, prove it wrong and be famous!

Even if we do not agree on **why** a concept works, if we **observe** that it helps, we can still use it

This is how medicine works (e.g., Anesthetics)!

1. Review
2. **Dirty secret of deep learning**
3. Optimization is hard
4. Deeper neural networks
5. Activation functions
6. Parameter initialization
7. Stochastic gradient descent
8. Modern optimization
9. Coding

1. Review
2. Dirty secret of deep learning
3. **Optimization is hard**
4. Deeper neural networks
5. Activation functions
6. Parameter initialization
7. Stochastic gradient descent
8. Modern optimization
9. Coding

A 2-layer neural network can represent **any** continuous function to arbitrary precision

Optimization is hard

A 2-layer neural network can represent **any** continuous function to arbitrary precision

$$| f(\mathbf{x}, \boldsymbol{\theta}) - g(\mathbf{x}) | < \varepsilon$$

Optimization is hard

A 2-layer neural network can represent **any** continuous function to arbitrary precision

$$| f(\mathbf{x}, \boldsymbol{\theta}) - g(\mathbf{x}) | < \varepsilon$$

$$\lim_{d_h \rightarrow \infty} \varepsilon = 0$$

Optimization is hard

A 2-layer neural network can represent **any** continuous function to arbitrary precision

$$| f(\mathbf{x}, \boldsymbol{\theta}) - g(\mathbf{x}) | < \varepsilon$$

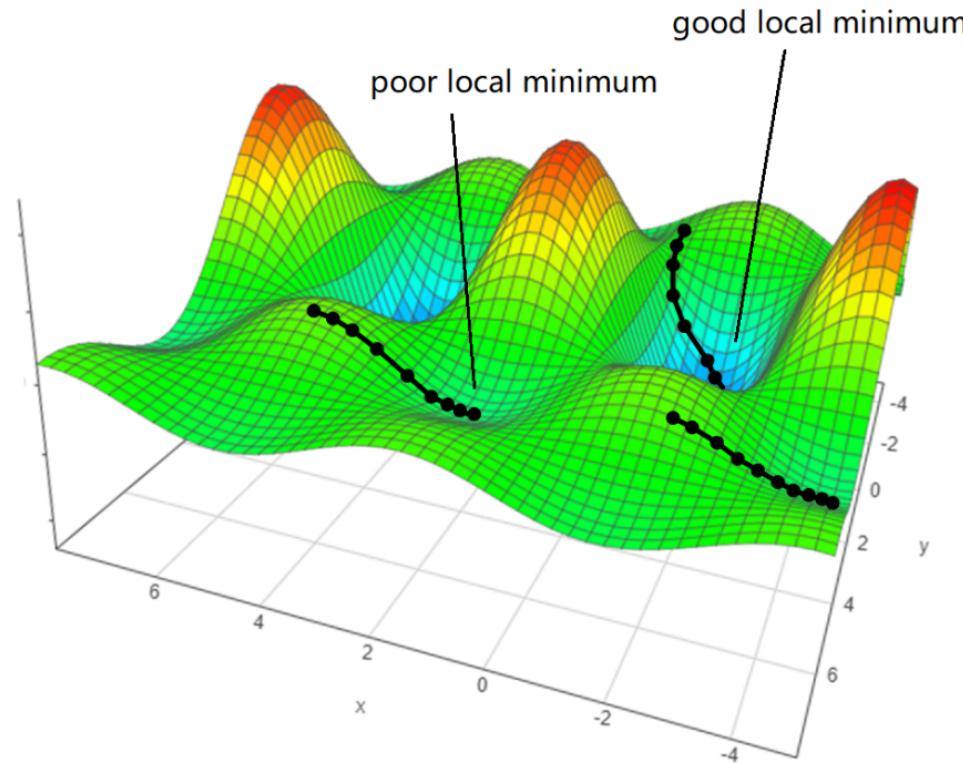
$$\lim_{d_h \rightarrow \infty} \varepsilon = 0$$

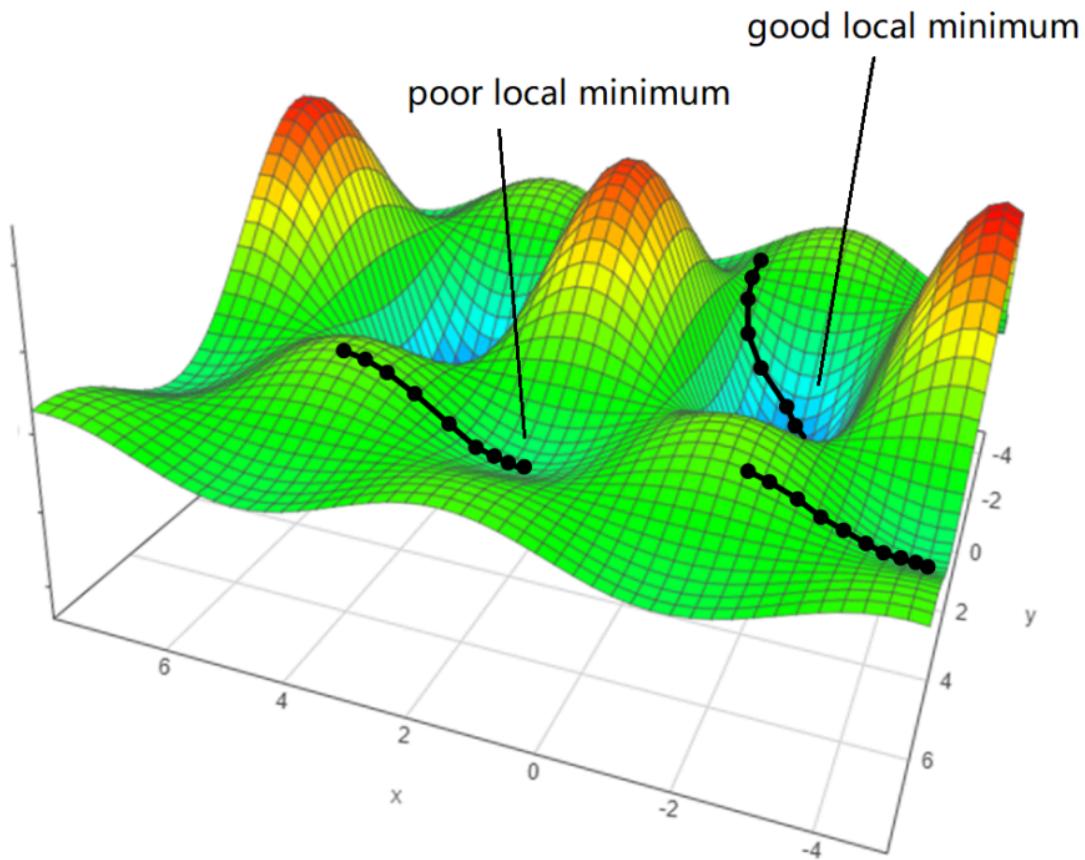
However, finding such $\boldsymbol{\theta}$ is a much harder problem

Gradient descent only guarantees convergence to a **local** optima

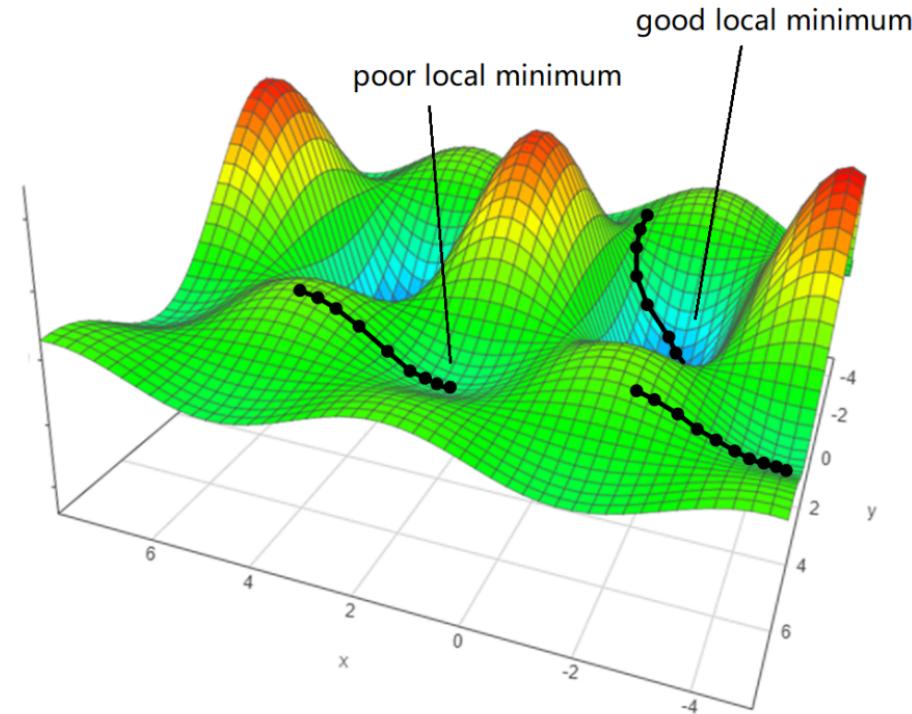
Optimization is hard

Gradient descent only guarantees convergence to a **local** optima





Optimization is hard

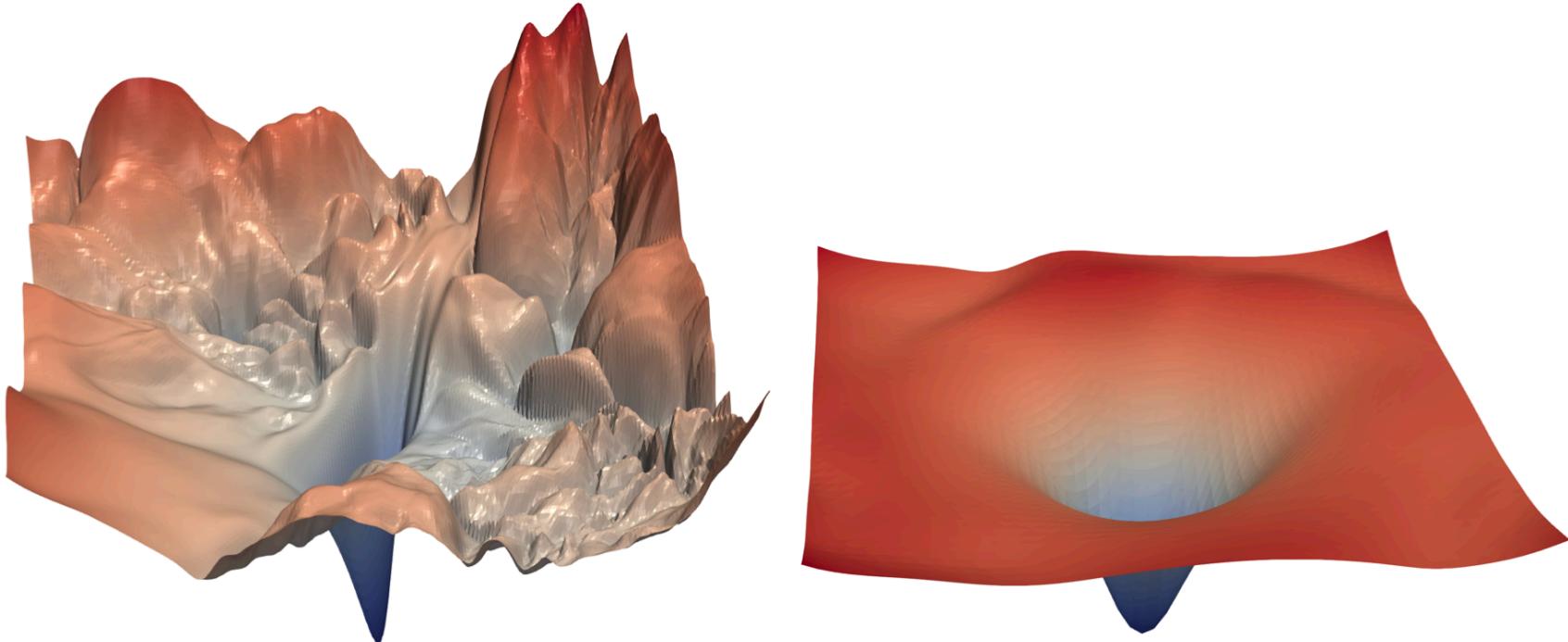


Harder tasks can have millions of local optima, and many of the local optima are not very good!

Many of the concepts today create a **flat loss landscape**

Optimization is hard

Many of the concepts today create a **flat loss landscape**



Gradient descent reaches a better optimum more quickly in these cases

1. Review
2. Dirty secret of deep learning
3. **Optimization is hard**
4. Deeper neural networks
5. Activation functions
6. Parameter initialization
7. Stochastic gradient descent
8. Modern optimization
9. Coding

1. Review
2. Dirty secret of deep learning
3. Optimization is hard
4. **Deeper neural networks**
5. Activation functions
6. Parameter initialization
7. Stochastic gradient descent
8. Modern optimization
9. Coding

A two-layer neural network is sufficient to approximate any continuous function to arbitrary precision

Deeper neural networks

A two-layer neural network is sufficient to approximate any continuous function to arbitrary precision

But only with infinite width $d_h \rightarrow \infty$

Deeper neural networks

A two-layer neural network is sufficient to approximate any continuous function to arbitrary precision

But only with infinite width $d_h \rightarrow \infty$

For certain problems, adding one more layer is equivalent to
exponentially increasing the width

- Eldan, Ronen, and Ohad Shamir. “The power of depth for feedforward neural networks.” Conference on learning theory. PMLR, 2016.

$2 \times 32 \times 32 \Rightarrow 2^{2 \times 32 \times 32} \approx 10^{616}$; universe has 10^{80} atoms

Deeper neural networks

A two-layer neural network is sufficient to approximate any continuous function to arbitrary precision

But only with infinite width $d_h \rightarrow \infty$

For certain problems, adding one more layer is equivalent to **exponentially increasing** the width

- Eldan, Ronen, and Ohad Shamir. “The power of depth for feedforward neural networks.” Conference on learning theory. PMLR, 2016.

$$2 \times 32 \times 32 \Rightarrow 2^{2 \times 32 \times 32} \approx 10^{616}; \text{ universe has } 10^{80} \text{ atoms}$$

We need more layers for harder problems

In fact, we do not just need **deeper** networks, but also **wider** networks

Deeper neural networks

In fact, we do not just need **deeper** networks, but also **wider** networks

The number of neurons in a deep neural network affects the quality of local optima

Deeper neural networks

In fact, we do not just need **deeper** networks, but also **wider** networks

The number of neurons in a deep neural network affects the quality of local optima

From Choromanska, Anna, et al. “The loss surfaces of multilayer networks.”:

Deeper neural networks

In fact, we do not just need **deeper** networks, but also **wider** networks

The number of neurons in a deep neural network affects the quality of local optima

From Choromanska, Anna, et al. “The loss surfaces of multilayer networks.”:

- “For large-size networks, most local minima are equivalent and yield similar performance on a test set.”

Deeper neural networks

In fact, we do not just need **deeper** networks, but also **wider** networks

The number of neurons in a deep neural network affects the quality of local optima

From Choromanska, Anna, et al. “The loss surfaces of multilayer networks.”:

- “For large-size networks, most local minima are equivalent and yield similar performance on a test set.”
- “The probability of finding a “bad” (high value) local minimum is non-zero for small-size networks and decreases quickly with network size”

To summarize, deeper and wider neural networks tend to produce better results

Deeper neural networks

To summarize, deeper and wider neural networks tend to produce better results

Add more layers to your network

Deeper neural networks

To summarize, deeper and wider neural networks tend to produce better results

Add more layers to your network

Increase the width of each layer

```
# Deep neural network
from torch import nn

d_x, d_y, d_h = 1, 1, 16
# Linear(input, output)
l1 = nn.Linear(d_x, d_h)
l2 = nn.Linear(d_h, d_y)
```

Deeper neural networks

```
# Deep neural network
from torch import nn

d_x, d_y, d_h = 1, 1, 16
# Linear(input, output)
l1 = nn.Linear(d_x, d_h)
l2 = nn.Linear(d_h, d_y)

# Deeper and wider neural
# network
from torch import nn

d_x, d_y, d_h = 1, 1, 256
# Linear(input, output)
l1 = nn.Linear(d_x, d_h)
l2 = nn.Linear(d_h, d_h)
l3 = nn.Linear(d_h, d_h)

...
l6 = nn.Linear(d_h, d_y)
```

```
import torch
d_x, d_y, d_h = 1, 1, 256
net = torch.nn.Sequential([
    torch.nn.Linear(d_x, d_h),
    torch.nn.Sigmoid(),
    torch.nn.Linear(d_h, d_h),
    torch.nn.Sigmoid(),
    ...
    torch.nn.Linear(d_h, d_y),
])
x = torch.ones((d_x,))
y = net(x)
```

```
import jax, equinox
d_x, d_y, d_h = 1, 1, 256
net = equinox.nn.Sequential([
    equinox.nn.Linear(d_x, d_h),
    equinox.nn.Lambda(jax.nn.sigmoid),
    equinox.nn.Linear(d_h, d_h),
    equinox.nn.Lambda(jax.nn.sigmoid),
    ...
    equinox.nn.Linear(d_h, d_y),
])
x = jax.numpy.ones((d_x,))
y = net(x)
```

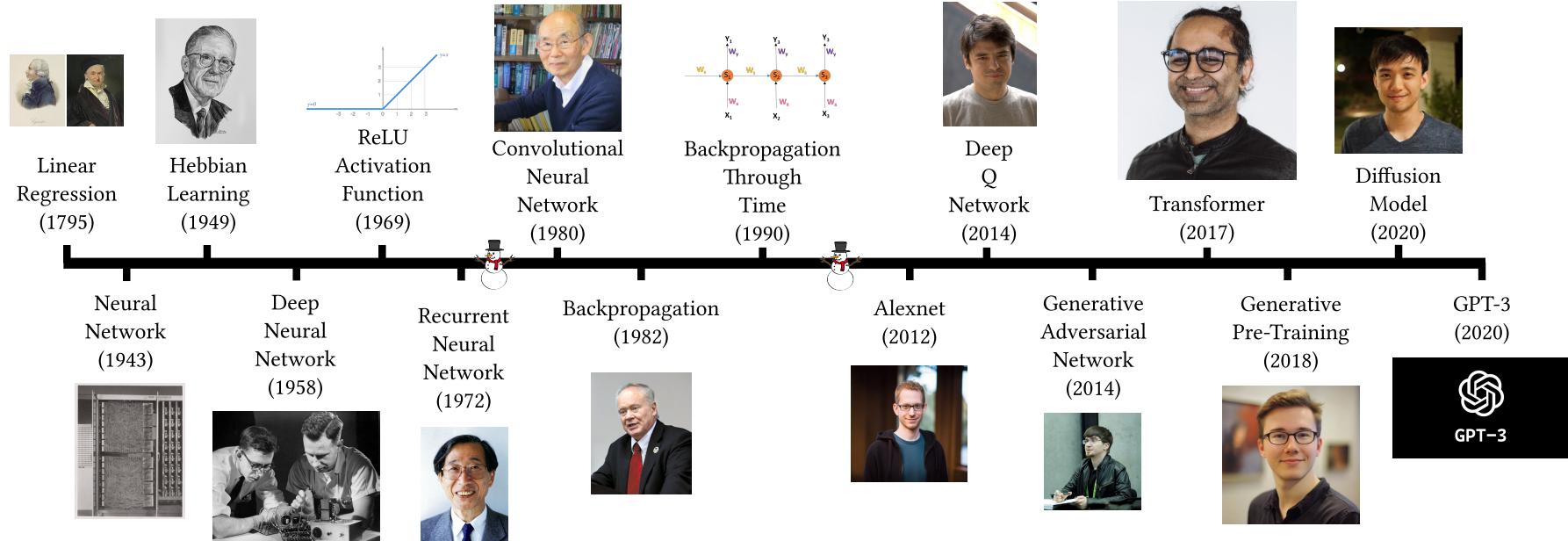
1. Review
2. Dirty secret of deep learning
3. Optimization is hard
4. **Deeper neural networks**
5. Activation functions
6. Parameter initialization
7. Stochastic gradient descent
8. Modern optimization
9. Coding

1. Review
2. Dirty secret of deep learning
3. Optimization is hard
4. Deeper neural networks
5. **Activation functions**
6. Parameter initialization
7. Stochastic gradient descent
8. Modern optimization
9. Coding

The sigmoid function was the standard activation function until ~ 2012

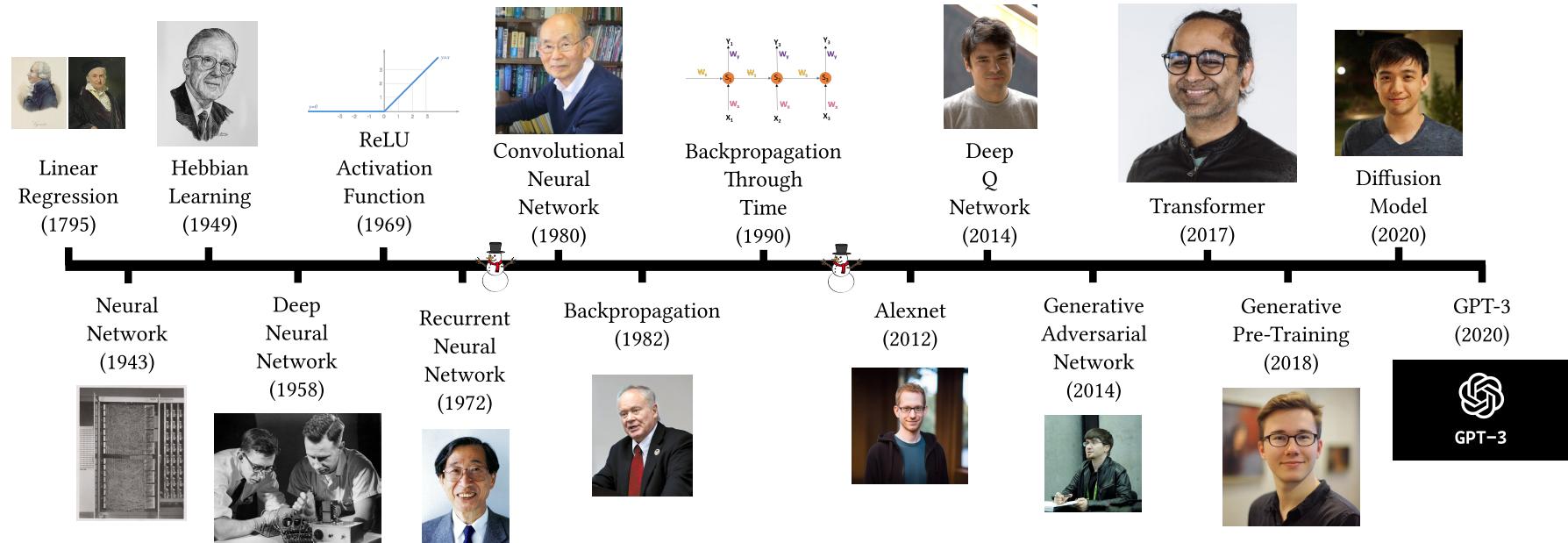
Activation functions

The sigmoid function was the standard activation function until ~ 2012

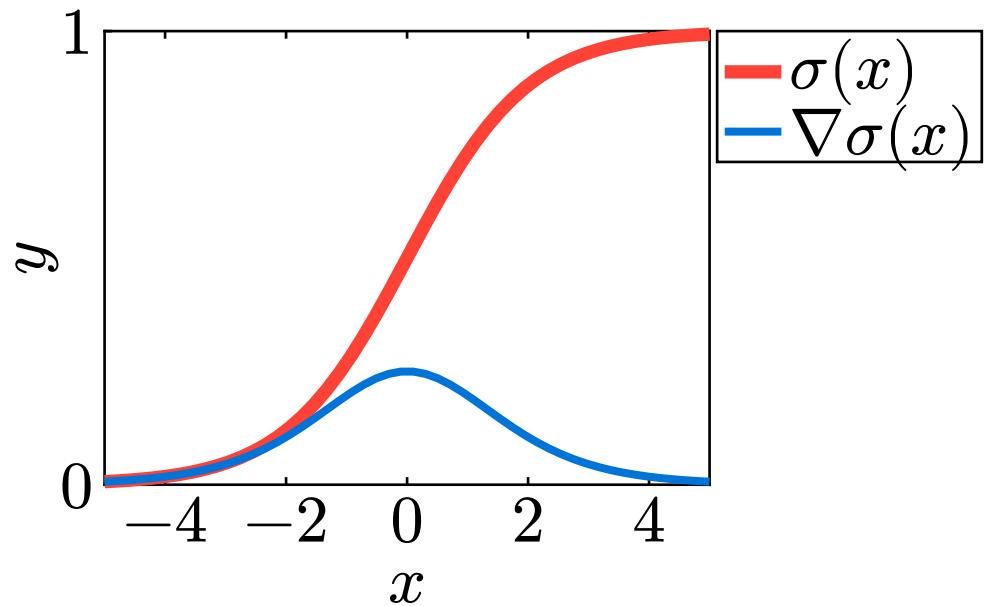


Activation functions

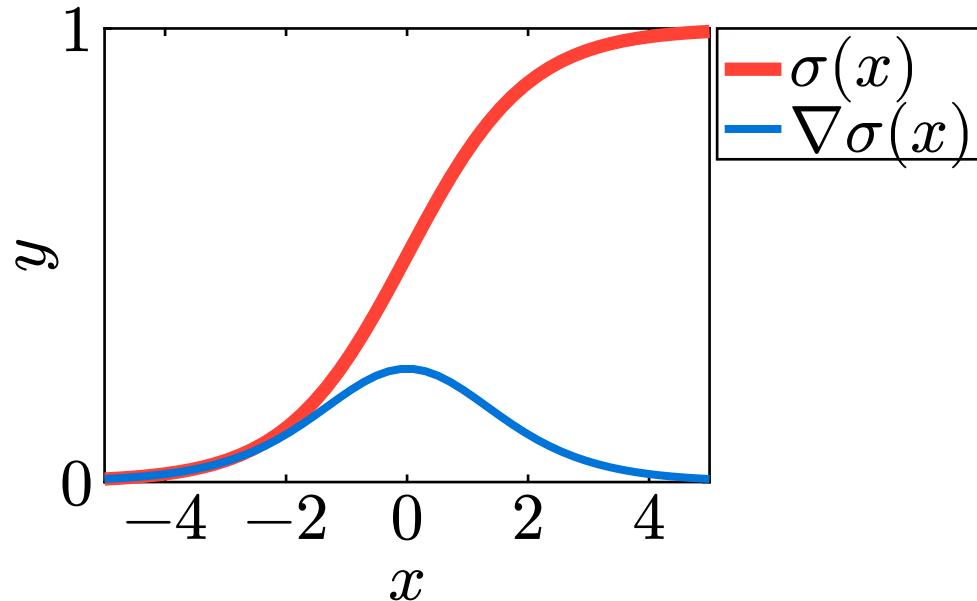
The sigmoid function was the standard activation function until ~ 2012



In 2012, people realized that ReLU activation performed much better

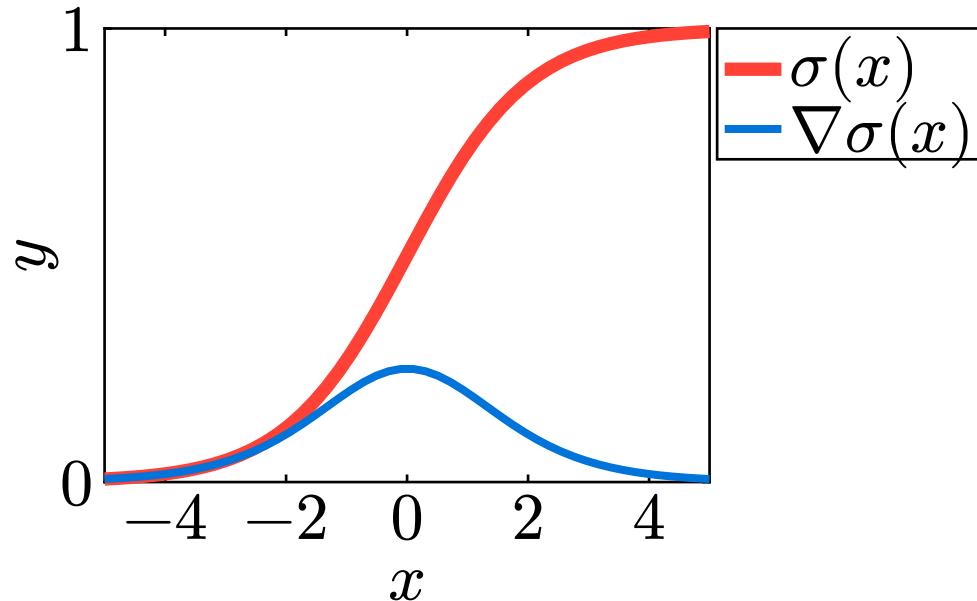


Activation functions



The sigmoid function can result in
a **vanishing gradient**

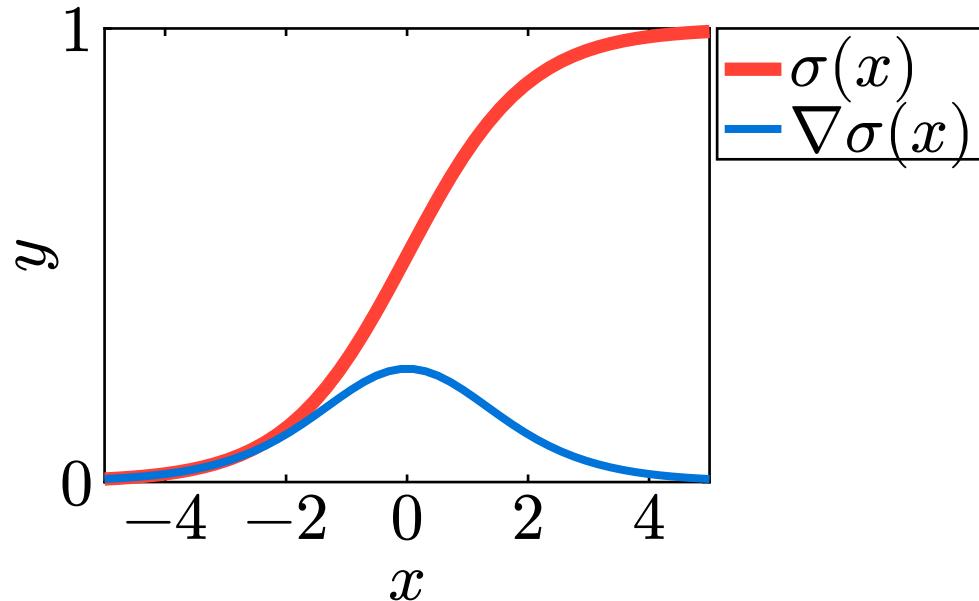
Activation functions



The sigmoid function can result in
a **vanishing gradient**

$$f(x, \theta) = \sigma(\theta_3^\top \sigma(\theta_2^\top \sigma(\theta_1^\top \bar{x})))$$

Activation functions

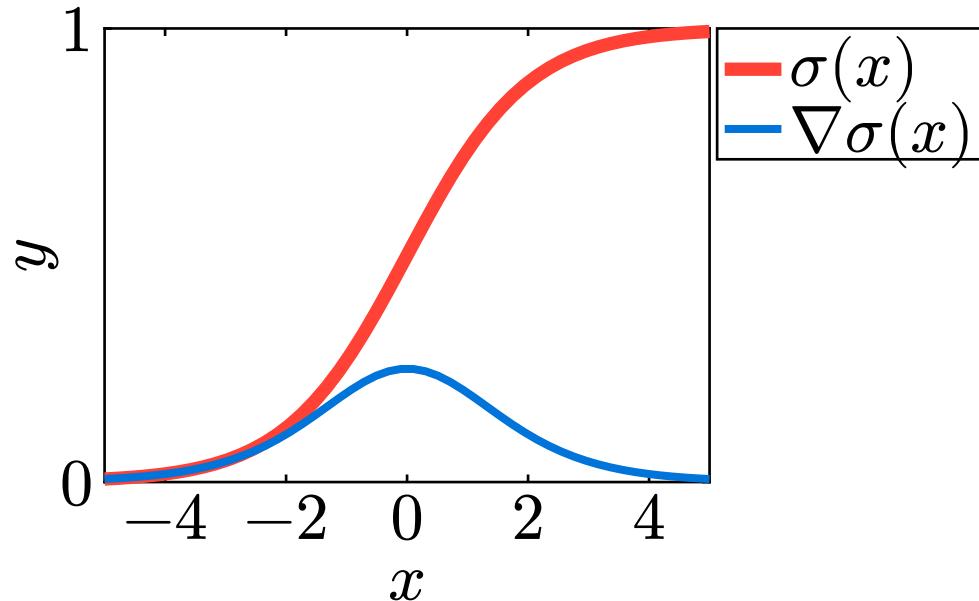


The sigmoid function can result in
a **vanishing gradient**

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}_3^\top \sigma(\boldsymbol{\theta}_2^\top \sigma(\boldsymbol{\theta}_1^\top \mathbf{x})))$$

$$\nabla_{\boldsymbol{\theta}} f(\mathbf{x}, \boldsymbol{\theta}) = \nabla[\sigma](\boldsymbol{\theta}_3^\top \sigma(\boldsymbol{\theta}_2^\top \sigma(\boldsymbol{\theta}_1^\top \mathbf{x}))) \cdot \nabla[\sigma](\boldsymbol{\theta}_2^\top \sigma(\boldsymbol{\theta}_1^\top \mathbf{x})) \cdot \nabla[\sigma](\boldsymbol{\theta}_1^\top \mathbf{x})$$

Activation functions

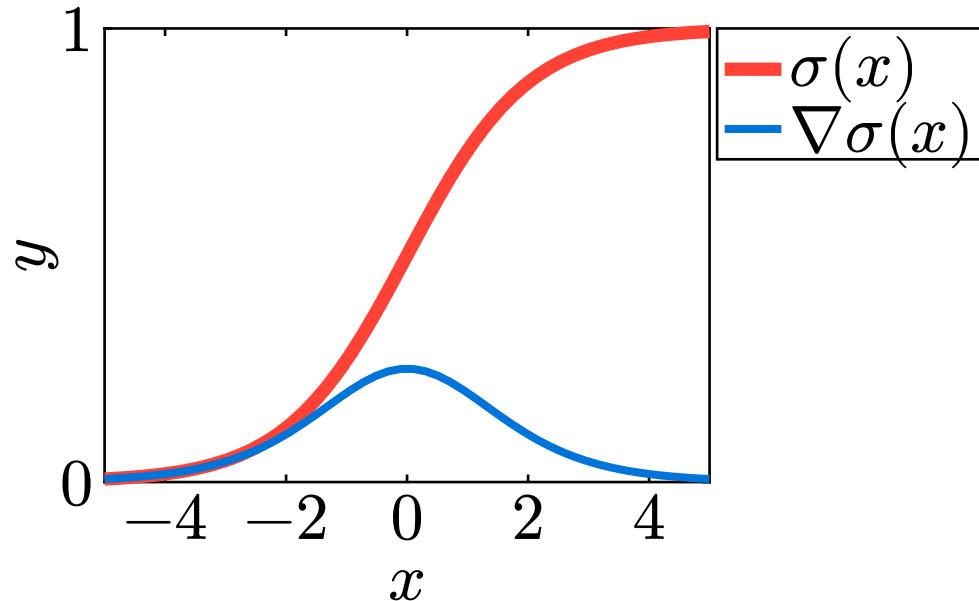


The sigmoid function can result in
a **vanishing gradient**

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}_3^\top \sigma(\boldsymbol{\theta}_2^\top \sigma(\boldsymbol{\theta}_1^\top \mathbf{x})))$$

$$\nabla_{\boldsymbol{\theta}} f(\mathbf{x}, \boldsymbol{\theta}) = \underbrace{\nabla[\sigma](\boldsymbol{\theta}_3^\top \sigma(\boldsymbol{\theta}_2^\top \sigma(\boldsymbol{\theta}_1^\top \mathbf{x})))}_{<1} \cdot \underbrace{\nabla[\sigma](\boldsymbol{\theta}_2^\top \sigma(\boldsymbol{\theta}_1^\top \mathbf{x}))}_{<1} \cdot \underbrace{\nabla[\sigma](\boldsymbol{\theta}_1^\top \mathbf{x})}_{<1}$$

Activation functions



The sigmoid function can result in
a **vanishing gradient**

$$f(\mathbf{x}, \boldsymbol{\theta}) = \sigma(\boldsymbol{\theta}_3^\top \sigma(\boldsymbol{\theta}_2^\top \sigma(\boldsymbol{\theta}_1^\top \mathbf{x})))$$

$$\nabla_{\boldsymbol{\theta}} f(\mathbf{x}, \boldsymbol{\theta}) = \underbrace{\nabla[\sigma](\boldsymbol{\theta}_3^\top \sigma(\boldsymbol{\theta}_2^\top \sigma(\boldsymbol{\theta}_1^\top \mathbf{x})))}_{<1} \cdot \underbrace{\nabla[\sigma](\boldsymbol{\theta}_2^\top \sigma(\boldsymbol{\theta}_1^\top \mathbf{x}))}_{<1} \cdot \underbrace{\nabla[\sigma](\boldsymbol{\theta}_1^\top \mathbf{x})}_{<1}$$

$$\nabla_{\boldsymbol{\theta}} f(\mathbf{x}, \boldsymbol{\theta}) \approx 0$$

To fix the vanishing gradient, researchers use the **rectified linear unit (ReLU)**

$$\sigma(x) = \max(0, x)$$

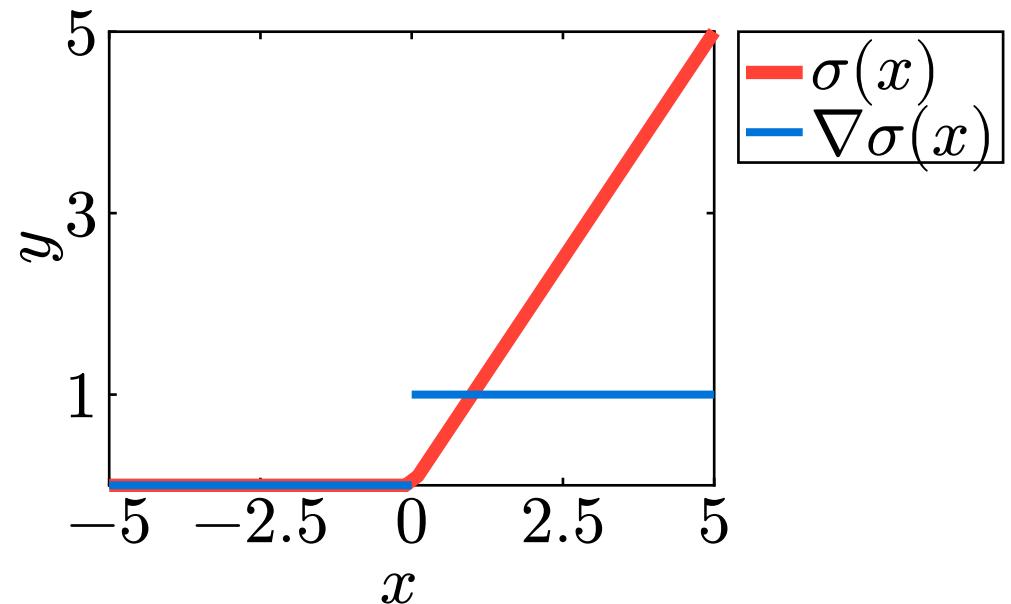
$$\nabla\sigma(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$

Activation functions

To fix the vanishing gradient, researchers use the **rectified linear unit (ReLU)**

$$\sigma(x) = \max(0, x)$$

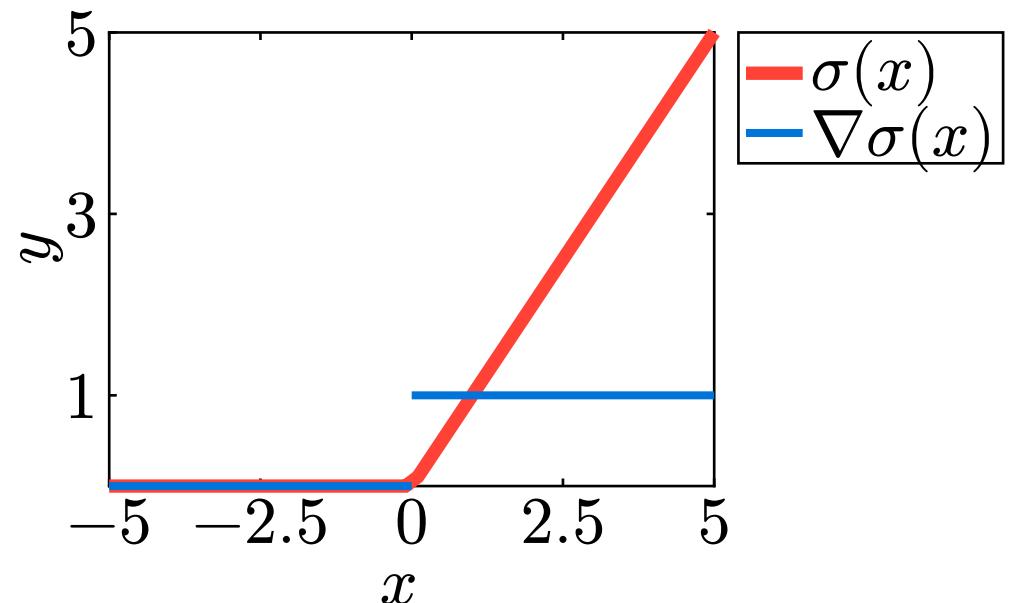
$$\nabla\sigma(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$



Activation functions

$$\sigma(x) = \max(0, x)$$

$$\nabla\sigma(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$

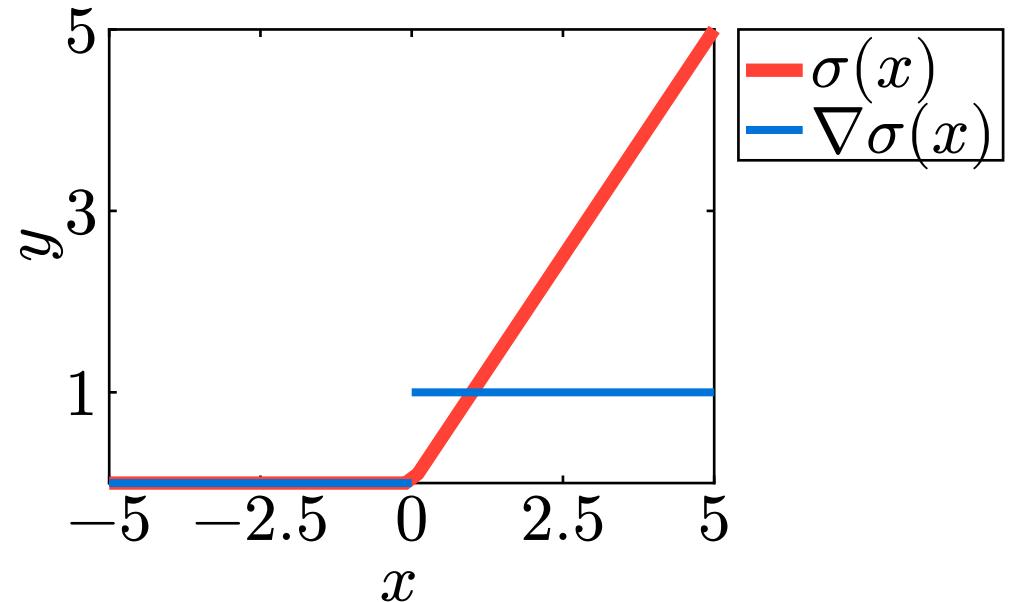


Looks nothing like a biological neuron

Activation functions

$$\sigma(x) = \max(0, x)$$

$$\nabla\sigma(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$



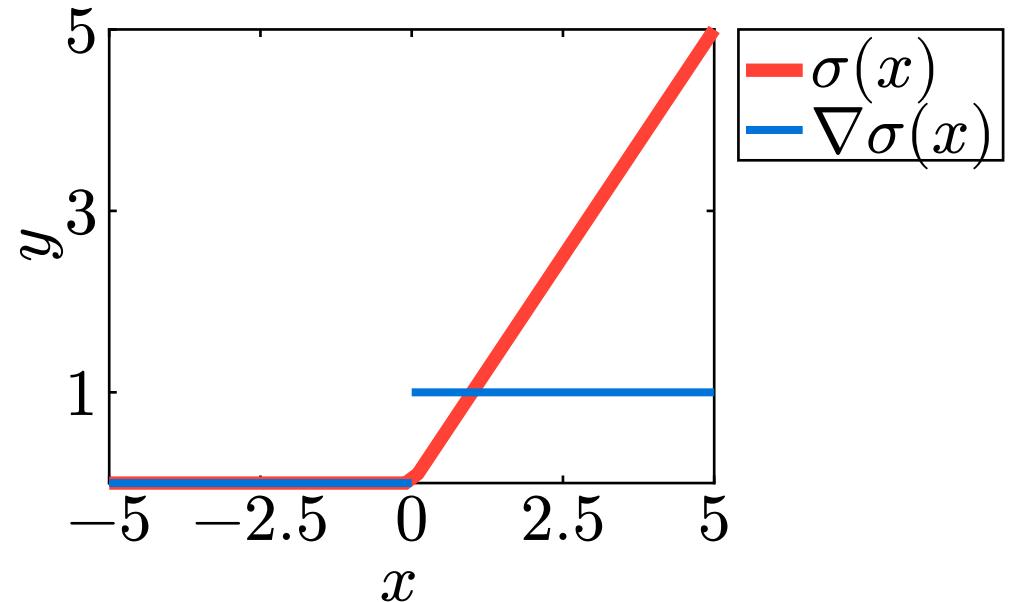
Looks nothing like a biological neuron

However, it works much better than sigmoid in practice

Activation functions

$$\sigma(x) = \max(0, x)$$

$$\nabla\sigma(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$



Looks nothing like a biological neuron

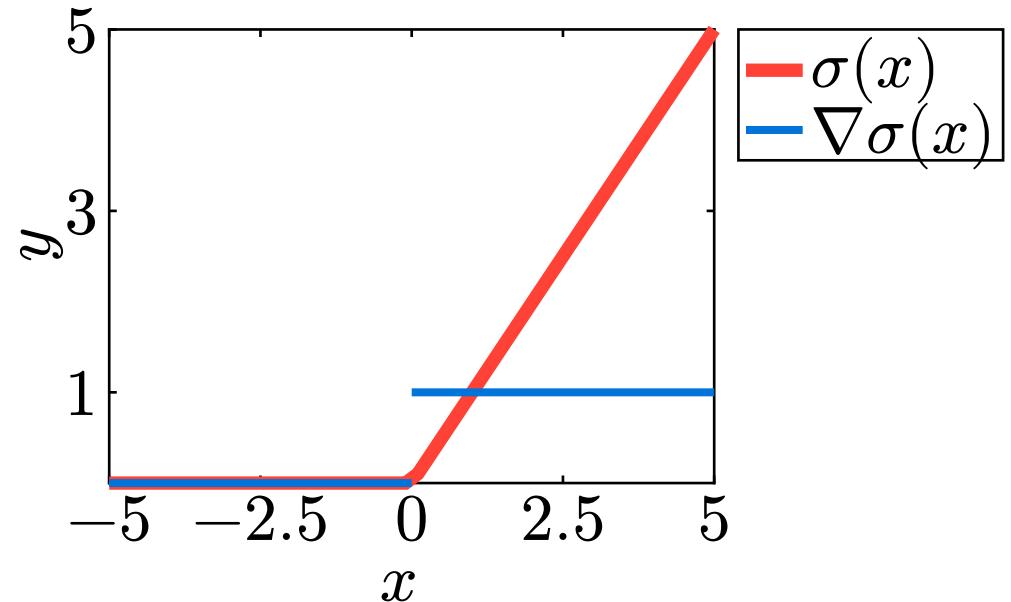
However, it works much better than sigmoid in practice

Via chain rule, gradient is $1 \cdot 1 \cdot 1 \dots$ which does not vanish

Activation functions

$$\sigma(x) = \max(0, x)$$

$$\nabla\sigma(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$



However, it works much better than sigmoid in practice

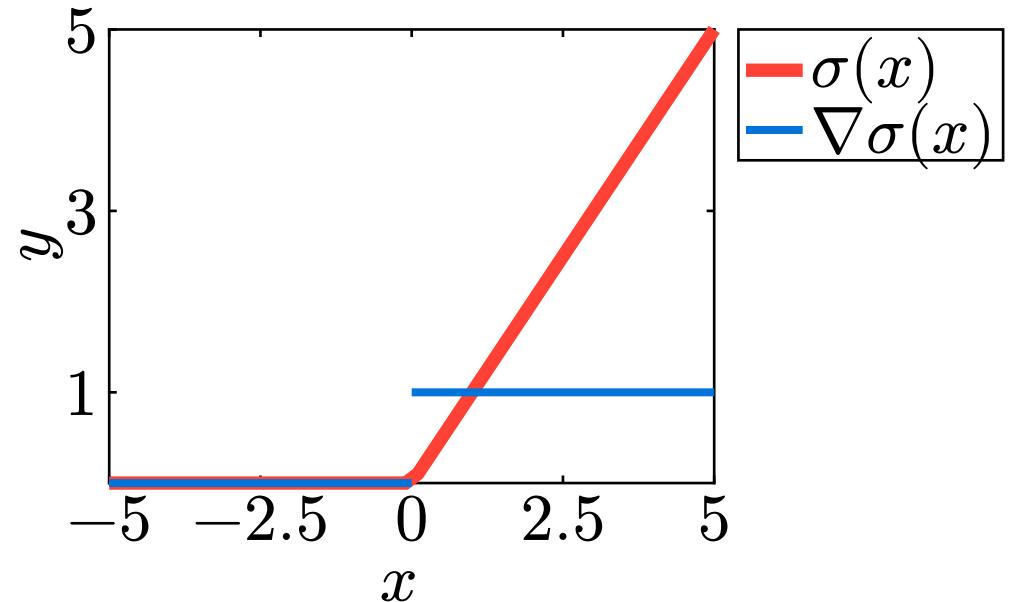
Via chain rule, gradient is $1 \cdot 1 \cdot 1 \dots$ which does not vanish

The gradient is constant, resulting in easier optimization

Activation functions

$$\sigma(x) = \max(0, x)$$

$$\nabla \sigma(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$



Via chain rule, gradient is $1 \cdot 1 \cdot 1 \dots$ which does not vanish

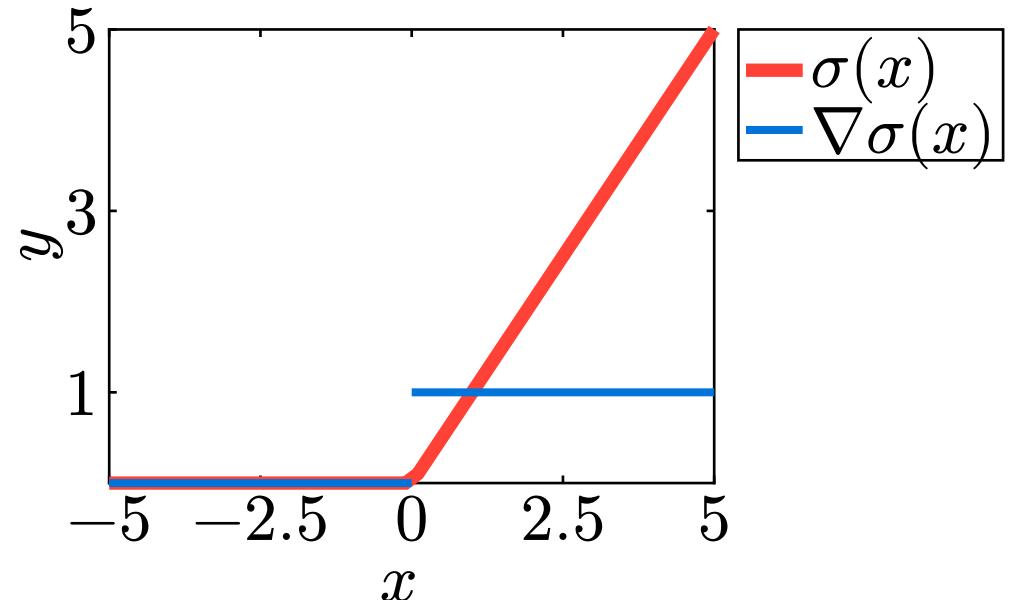
The gradient is constant, resulting in easier optimization

Question: Any problems?

Activation functions

$$\sigma(x) = \max(0, x)$$

$$\nabla\sigma(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$



The gradient is constant, resulting in easier optimization

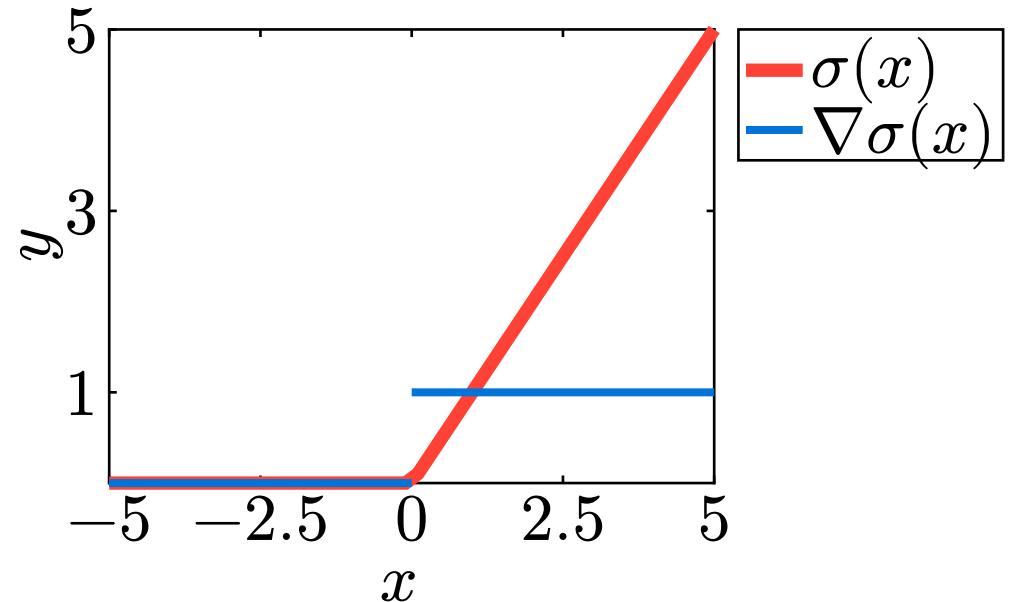
Question: Any problems?

Answer: Zero gradient region!

Activation functions

$$\sigma(x) = \max(0, x)$$

$$\nabla\sigma(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$



Question: Any problems?

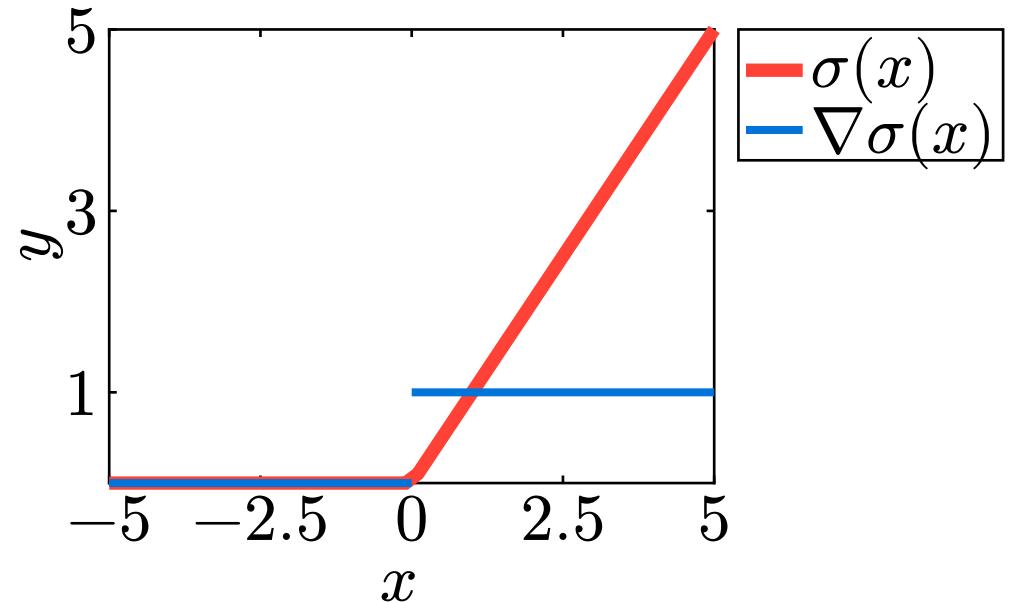
Answer: Zero gradient region!

Neurons can get “stuck”, always output 0

Activation functions

$$\sigma(x) = \max(0, x)$$

$$\nabla\sigma(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$



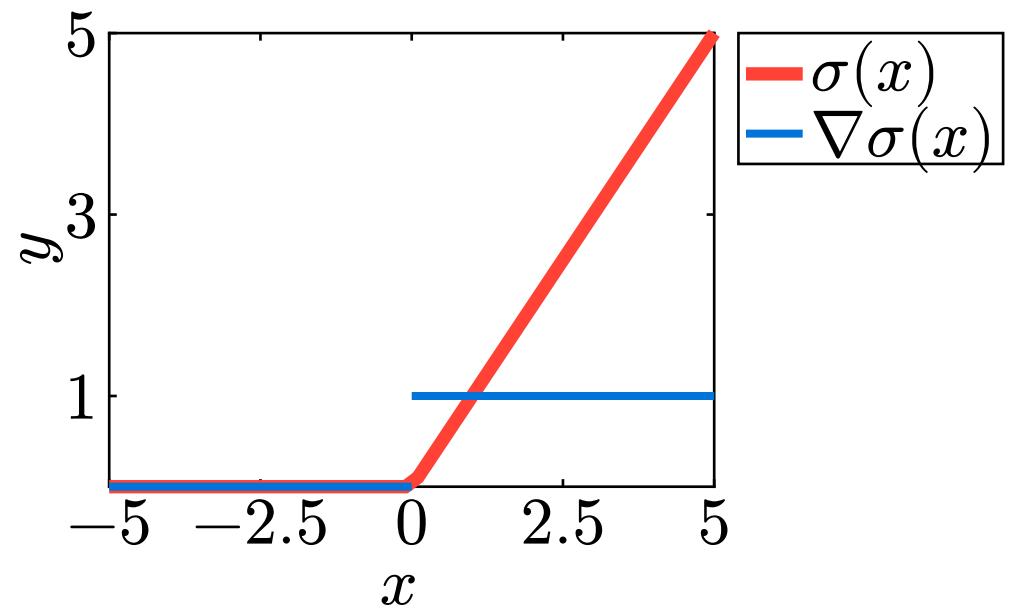
Answer: Zero gradient region!

Neurons can get “stuck”, always output 0

These neurons cannot recover, they are **dead neurons**

$$\sigma(x) = \max(0, x)$$

$$\nabla\sigma(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$

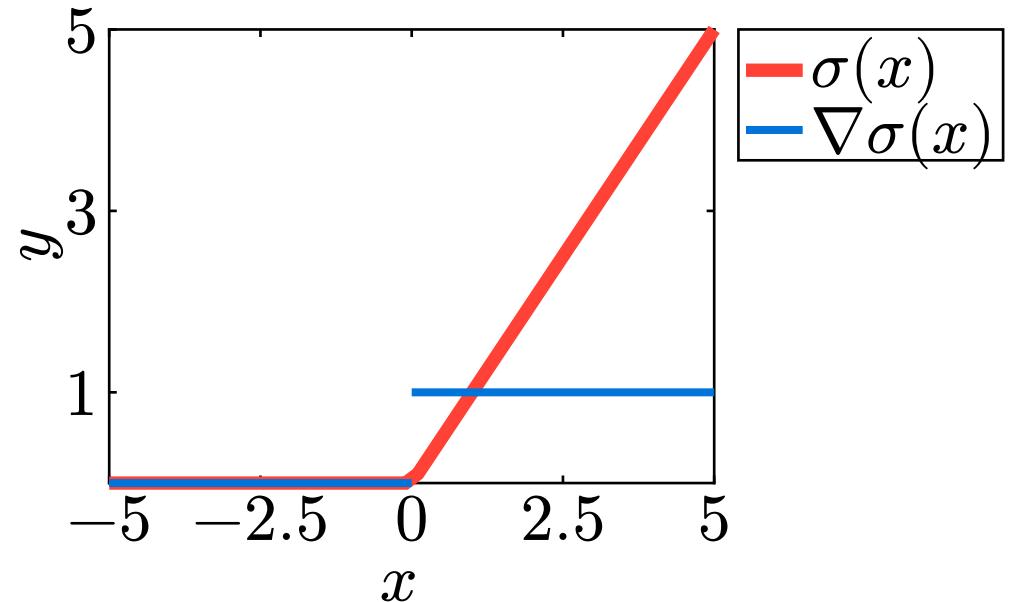


These neurons cannot recover, they are **dead neurons**

Activation functions

$$\sigma(x) = \max(0, x)$$

$$\nabla\sigma(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$



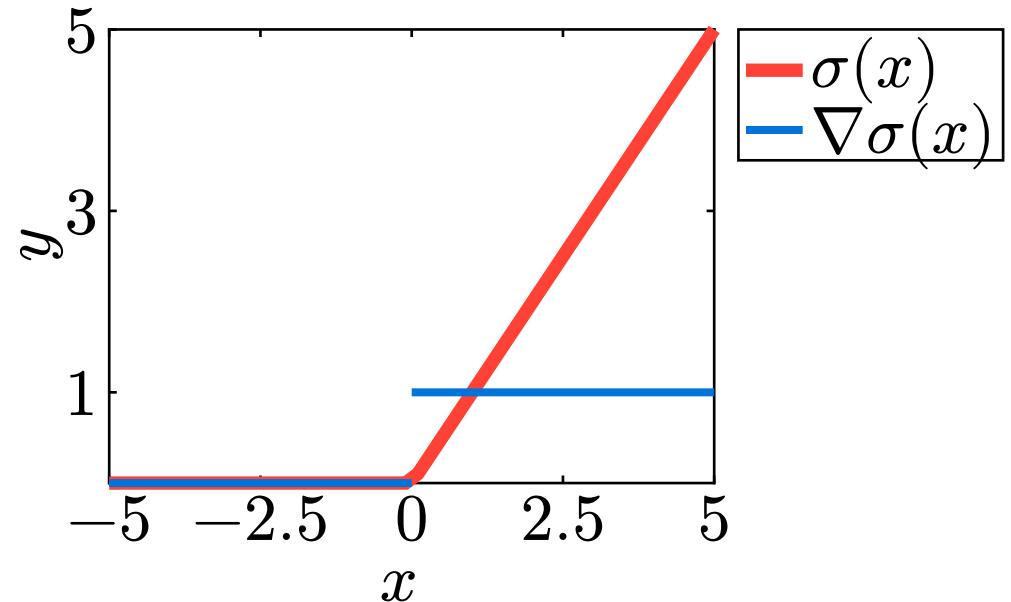
These neurons cannot recover, they are **dead neurons**

Training for longer results in more dead neurons

Activation functions

$$\sigma(x) = \max(0, x)$$

$$\nabla\sigma(x) = \begin{cases} 0 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$



These neurons cannot recover, they are **dead neurons**

Training for longer results in more dead neurons

Dead neurons hurt your network!

To fix dying neurons, use **leaky ReLU**

Activation functions

To fix dying neurons, use **leaky ReLU**

$$\sigma(x) = \max(0.1x, x)$$

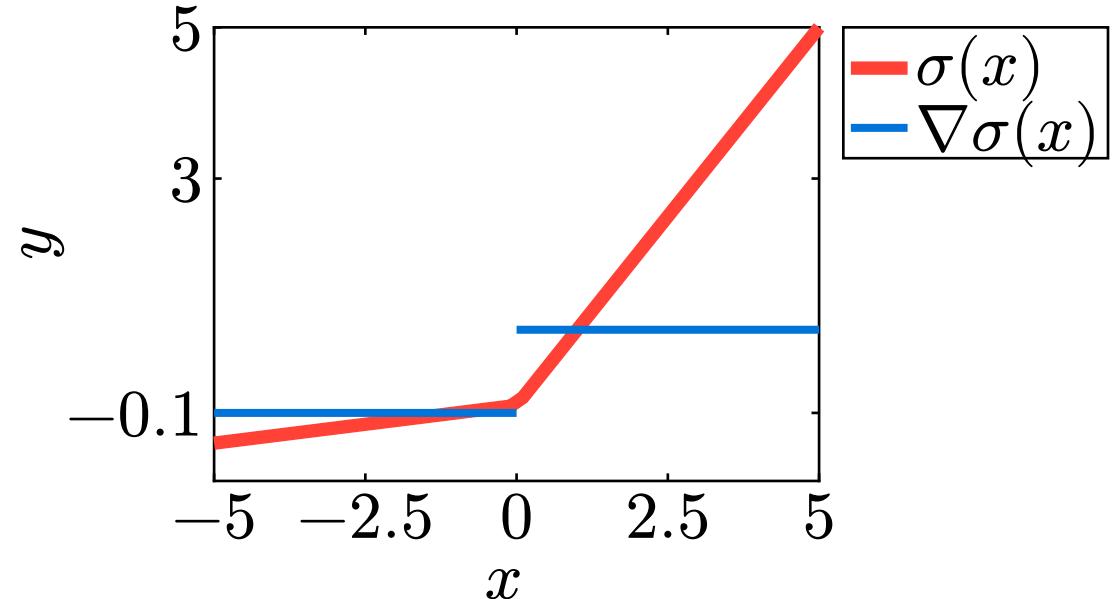
$$\nabla\sigma(x) = \begin{cases} 0.1 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$

Activation functions

To fix dying neurons, use **leaky ReLU**

$$\sigma(x) = \max(0.1x, x)$$

$$\nabla\sigma(x) = \begin{cases} 0.1 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$

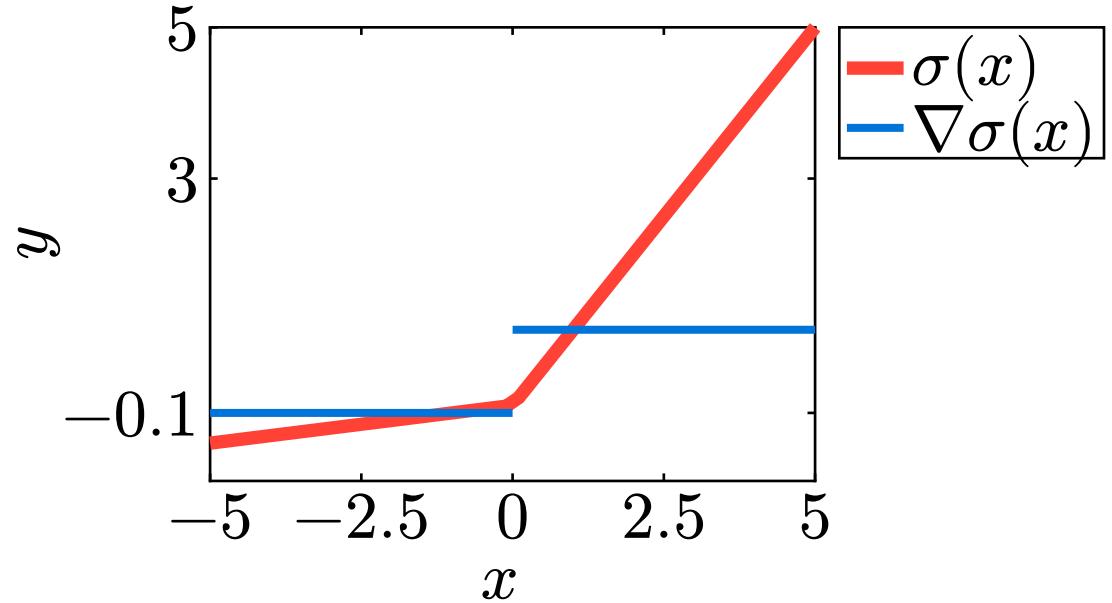


Activation functions

To fix dying neurons, use **leaky ReLU**

$$\sigma(x) = \max(0.1x, x)$$

$$\nabla\sigma(x) = \begin{cases} 0.1 & \text{if } x < 0 \\ 1 & \text{if } x \geq 0 \end{cases}$$



Small negative slope allows dead neurons to recover

There are other activation
functions that are better than
leaky ReLU

Activation functions

There are other activation functions that are better than leaky ReLU

- Mish

Activation functions

There are other activation functions that are better than leaky ReLU

- Mish
- Swish

Activation functions

There are other activation functions that are better than leaky ReLU

- Mish
- Swish
- ELU

Activation functions

There are other activation functions that are better than leaky ReLU

- Mish
- Swish
- ELU
- GeLU

Activation functions

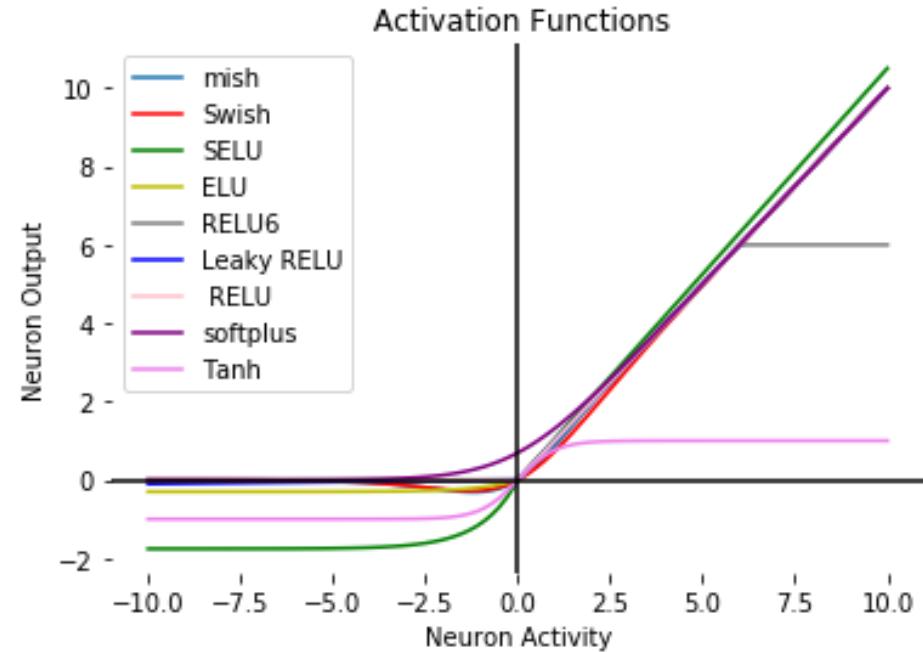
There are other activation functions that are better than leaky ReLU

- Mish
- Swish
- ELU
- GeLU
- SeLU

Activation functions

There are other activation functions that are better than leaky ReLU

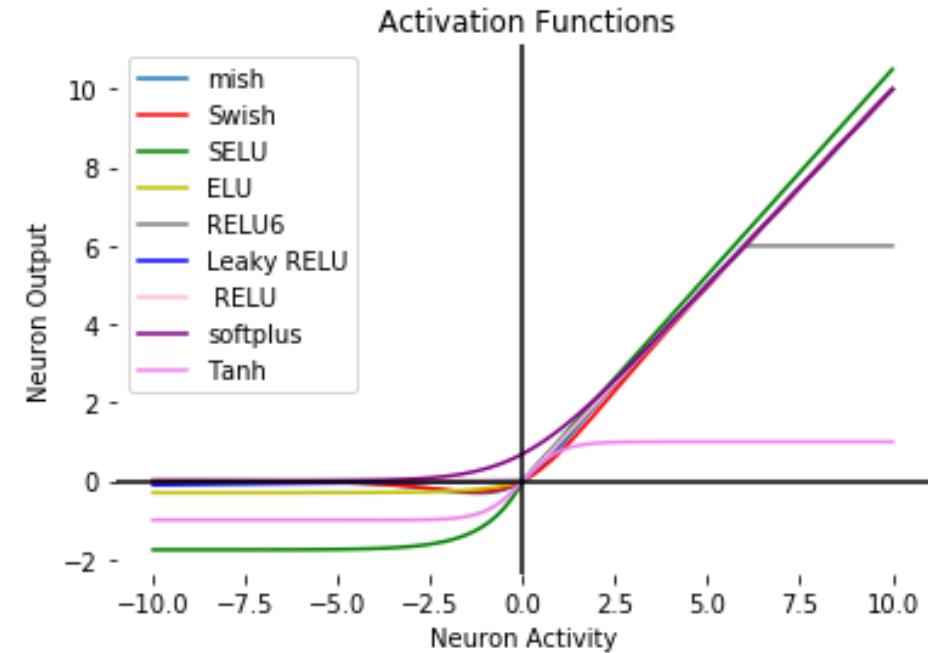
- Mish
- Swish
- ELU
- GeLU
- SeLU



Activation functions

There are other activation functions that are better than leaky ReLU

- Mish
- Swish
- ELU
- GeLU
- SeLU

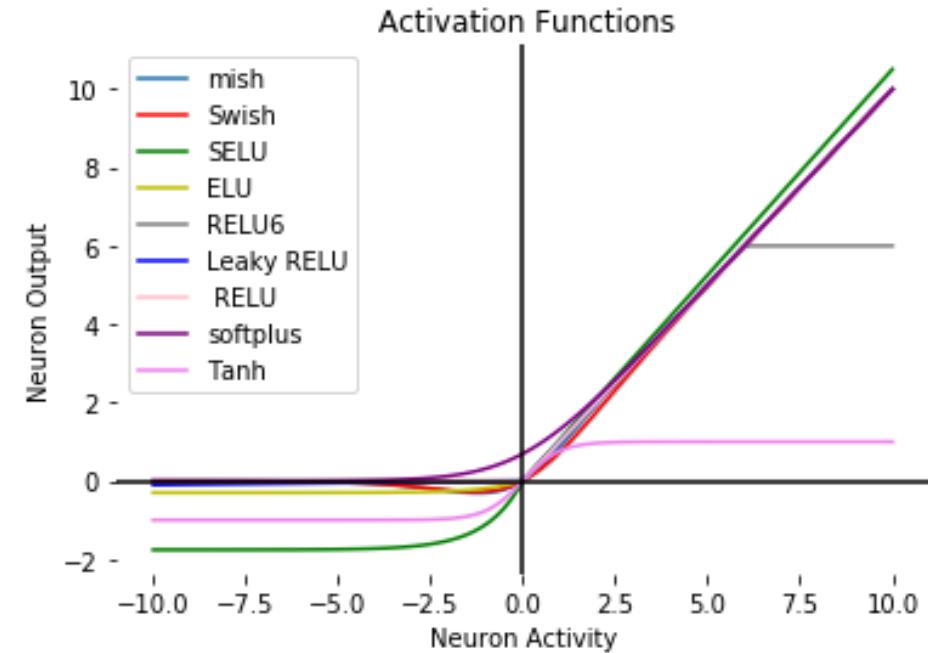


They are all very similar

Activation functions

There are other activation functions that are better than leaky ReLU

- Mish
- Swish
- ELU
- GeLU
- SeLU



They are all very similar

I usually use leaky ReLU because it works well enough

<https://pytorch.org/docs/stable/nn.html#non-linear-activations-weighted-sum-nonlinearity>

<https://jax.readthedocs.io/en/latest/jax.nn.html#activation-functions>

1. Review
2. Dirty secret of deep learning
3. Optimization is hard
4. Deeper neural networks
5. **Activation functions**
6. Parameter initialization
7. Stochastic gradient descent
8. Modern optimization
9. Coding

1. Review
2. Dirty secret of deep learning
3. Optimization is hard
4. Deeper neural networks
5. Activation functions
6. **Parameter initialization**
7. Stochastic gradient descent
8. Modern optimization
9. Coding

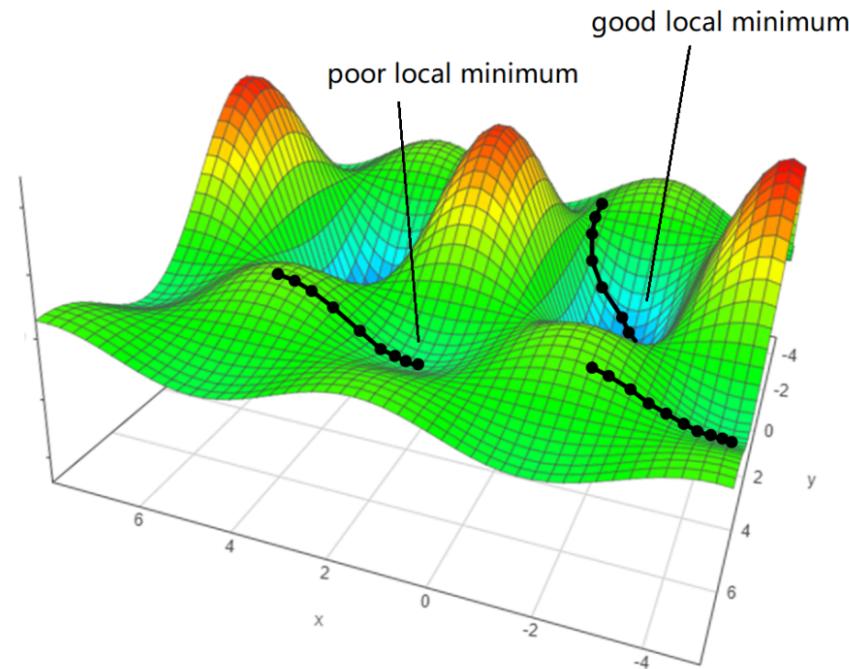
Recall the gradient descent algorithm

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

Initial θ is starting position for gradient descent

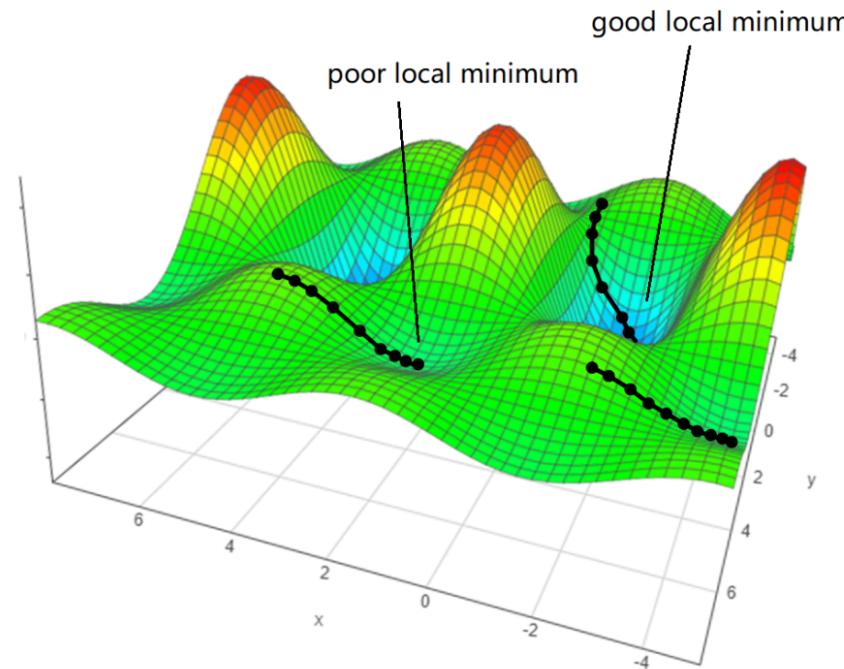
Parameter initialization

Initial θ is starting position for gradient descent



Parameter initialization

Initial θ is starting position for gradient descent



Pick θ that results in good local minima

Start simple, initialize all parameters to 0

$$\theta = \begin{bmatrix} 0 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{bmatrix}, \begin{bmatrix} 0 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{bmatrix}, \dots$$

Parameter initialization

Start simple, initialize all parameters to 0

$$\theta = \begin{bmatrix} 0 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{bmatrix}, \begin{bmatrix} 0 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{bmatrix}, \dots$$

Question: Any issues?

Parameter initialization

Start simple, initialize all parameters to 0

$$\theta = \begin{bmatrix} 0 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{bmatrix}, \begin{bmatrix} 0 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{bmatrix}, \dots$$

Question: Any issues?

Answer: The gradient will always be zero

$$\nabla_{\theta_1} f = \sigma(\theta_2^\top \sigma(\theta_1^\top \bar{x})) \ \sigma(\theta_1^\top \bar{x}) \ \bar{x}$$

Parameter initialization

Start simple, initialize all parameters to 0

$$\theta = \begin{bmatrix} 0 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{bmatrix}, \begin{bmatrix} 0 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \end{bmatrix}, \dots$$

Question: Any issues?

Answer: The gradient will always be zero

$$\nabla_{\theta_1} f = \sigma(\theta_2^\top \sigma(\theta_1^\top \bar{x})) \sigma(\theta_1^\top \bar{x}) \bar{x}$$

$$\nabla_{\theta_1} f = \sigma(\mathbf{0}^\top \sigma(\theta_1^\top \bar{x})) \sigma(\theta_1^\top \bar{x}) \bar{x} = 0$$

Ok, so initialize $\theta = 1$

Parameter initialization

Ok, so initialize $\theta = 1$

$$\theta = \begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{bmatrix}, \begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{bmatrix}, \dots$$

Parameter initialization

Ok, so initialize $\theta = 1$

$$\theta = \begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{bmatrix}, \begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{bmatrix}, \dots$$

Question: Any issues?

Parameter initialization

Ok, so initialize $\theta = 1$

$$\theta = \begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{bmatrix}, \begin{bmatrix} 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ 1 & \dots & 1 \end{bmatrix}, \dots$$

Question: Any issues?

All neurons in a layer will have the same gradient, and so they will always be the same (useless)

$$z_i = \sigma \left(\sum_{j=1}^{d_x} \theta_j \cdot \bar{x}_j \right) = \sigma \left(\sum_{j=1}^{d_x} \bar{x}_j \right)$$

θ must be randomly initialized for neurons

$$\theta = \begin{bmatrix} -0.5 & \dots & 2 \\ \vdots & \ddots & \vdots \\ 0.1 & \dots & 0.6 \end{bmatrix}, \begin{bmatrix} 1.3 & \dots & 1.2 \\ \vdots & \ddots & \vdots \\ -0.8 & \dots & -1.1 \end{bmatrix}, \dots$$

Parameter initialization

θ must be randomly initialized for neurons

$$\theta = \begin{bmatrix} -0.5 & \dots & 2 \\ \vdots & \ddots & \vdots \\ 0.1 & \dots & 0.6 \end{bmatrix}, \begin{bmatrix} 1.3 & \dots & 1.2 \\ \vdots & \ddots & \vdots \\ -0.8 & \dots & -1.1 \end{bmatrix}, \dots$$

But what scale? If $\theta \ll 0$ the gradients will vanish to zero, if $\theta \gg 0$ the gradients explode to infinity

Parameter initialization

θ must be randomly initialized for neurons

$$\theta = \begin{bmatrix} -0.5 & \dots & 2 \\ \vdots & \ddots & \vdots \\ 0.1 & \dots & 0.6 \end{bmatrix}, \begin{bmatrix} 1.3 & \dots & 1.2 \\ \vdots & \ddots & \vdots \\ -0.8 & \dots & -1.1 \end{bmatrix}, \dots$$

But what scale? If $\theta \ll 0$ the gradients will vanish to zero, if $\theta \gg 0$ the gradients explode to infinity

Almost everyone initializes following a single paper from 2010:

Parameter initialization

θ must be randomly initialized for neurons

$$\theta = \begin{bmatrix} -0.5 & \dots & 2 \\ \vdots & \ddots & \vdots \\ 0.1 & \dots & 0.6 \end{bmatrix}, \begin{bmatrix} 1.3 & \dots & 1.2 \\ \vdots & \ddots & \vdots \\ -0.8 & \dots & -1.1 \end{bmatrix}, \dots$$

But what scale? If $\theta \ll 0$ the gradients will vanish to zero, if $\theta \gg 0$ the gradients explode to infinity

Almost everyone initializes following a single paper from 2010:

- Glorot, Xavier, and Yoshua Bengio. “Understanding the difficulty of training deep feedforward neural networks.”

Parameter initialization

θ must be randomly initialized for neurons

$$\theta = \begin{bmatrix} -0.5 & \dots & 2 \\ \vdots & \ddots & \vdots \\ 0.1 & \dots & 0.6 \end{bmatrix}, \begin{bmatrix} 1.3 & \dots & 1.2 \\ \vdots & \ddots & \vdots \\ -0.8 & \dots & -1.1 \end{bmatrix}, \dots$$

But what scale? If $\theta \ll 0$ the gradients will vanish to zero, if $\theta \gg 0$ the gradients explode to infinity

Almost everyone initializes following a single paper from 2010:

- Glorot, Xavier, and Yoshua Bengio. “Understanding the difficulty of training deep feedforward neural networks.”
- Maybe there are better options?

Here is the magic equation, given the input and output size of the layer is d_h

$$\theta \sim \mathcal{U} \left[-\frac{\sqrt{6}}{\sqrt{2d_h}}, \frac{\sqrt{6}}{\sqrt{2d_h}} \right]$$

Parameter initialization

Here is the magic equation, given the input and output size of the layer is d_h

$$\theta \sim \mathcal{U} \left[-\frac{\sqrt{6}}{\sqrt{2d_h}}, \frac{\sqrt{6}}{\sqrt{2d_h}} \right]$$

If you have different input or output sizes, such as d_x, d_y , then the equation is

$$\theta \sim \mathcal{U} \left[-\frac{\sqrt{6}}{\sqrt{d_x + d_y}}, \frac{\sqrt{6}}{\sqrt{d_x + d_y}} \right]$$

These equations are designed for ReLU and similar activation functions

Parameter initialization

These equations are designed for ReLU and similar activation functions

They prevent vanishing or exploding gradients

Usually torch and jax/equinox will automatically use this initialization when you create nn.Linear

```
layer = nn.Linear(d_x, d_h) # Uses Glorot init
```

You can find many initialization functions at <https://pytorch.org/docs/stable/nn.init.html>

For JAX it is <https://jax.readthedocs.io/en/latest/jax.nn.initializers.html>

```
import torch
d_h = 10
# Manually
theta = torch.zeros((d_h + 1, d_h))
torch.nn.init.xavier_uniform_(theta)
theta = torch.nn.Parameter(theta)

# Using nn.Linear
layer = torch.nn.Linear(d_h, d_h)
# Use .data, to bypass autograd security
torch.nn.init.xavier_uniform_(layer.weight.data)
torch.nn.init.xavier_uniform_(layer.bias.data)
```

```
import jax
d_h = 10

init = jax.nn.initializers.glorot_uniform()
theta = init(jax.random.key(0), (d_h + 1, d_h))
```

```
import jax, equinox
d_h = 10

layer = equinox.nn.Linear(d_h, d_h, key=jax.random.key(0))
# Create new bias and weight
new_weight = init(jax.random.key(1), (d_h, d_h))
new_bias = init(jax.random.key(2), (d_h,))

# Use a lambda function to save space
# tree_at creates a new layer with the new weight
layer = equinox.tree_at(lambda l: l.weight, layer,
new_weight)
layer = equinox.tree_at(lambda l: l.bias, layer, new_weight)
```

Remember, in equinox and torch, nn.Linear will already be initialized correctly!

1. Review
2. Dirty secret of deep learning
3. Optimization is hard
4. Deeper neural networks
5. Activation functions
6. **Parameter initialization**
7. Stochastic gradient descent
8. Modern optimization
9. Coding

1. Review
2. Dirty secret of deep learning
3. Optimization is hard
4. Deeper neural networks
5. Activation functions
6. Parameter initialization
7. **Stochastic gradient descent**
8. Modern optimization
9. Coding

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

Stochastic gradient descent

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

Gradient descent computes $\nabla \mathcal{L}$ over all X

This works for our small datasets, where $n = 1000$

Stochastic gradient descent

This works for our small datasets, where $n = 1000$

Question: How many GB are the LLM datasets?

Stochastic gradient descent

This works for our small datasets, where $n = 1000$

Question: How many GB are the LLM datasets?

Answer: About 774,000 GB according to *Datasets for Large Language Models: A Comprehensive Survey*

Stochastic gradient descent

This works for our small datasets, where $n = 1000$

Question: How many GB are the LLM datasets?

Answer: About 774,000 GB according to *Datasets for Large Language Models: A Comprehensive Survey*

This is just the dataset size, the gradient is orders of magnitude larger

$$\nabla_{\theta} \mathcal{L}(x_{[i]}, y_{[i]}, \theta) = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_\ell}{\partial x_1} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_n}{\partial x_1} & \cdots & \frac{\partial f_\ell}{\partial x_1} \end{bmatrix}_{[i]}$$

Question: We do not have enough memory to compute the gradient.
What can we do?

Stochastic gradient descent

Question: We do not have enough memory to compute the gradient.
What can we do?

Answer: We approximate the gradient using a subset of the data

First, we sample random datapoint indices

$$i, j, k, \dots \sim \mathcal{U}[1, n]$$

Stochastic gradient descent

First, we sample random datapoint indices

$$i, j, k, \dots \sim \mathcal{U}[1, n]$$

Then construct a **batch** of training data

$$\begin{bmatrix} \mathbf{x}_{[i]} \\ \mathbf{x}_{[j]} \\ \mathbf{x}_{[k]} \\ \vdots \end{bmatrix}; \quad \begin{bmatrix} \mathbf{y}_{[i]} \\ \mathbf{y}_{[j]} \\ \mathbf{y}_{[k]} \\ \vdots \end{bmatrix}$$

Stochastic gradient descent

First, we sample random datapoint indices

$$i, j, k, \dots \sim \mathcal{U}[1, n]$$

Then construct a **batch** of training data

$$\begin{bmatrix} \mathbf{x}_{[i]} \\ \mathbf{x}_{[j]} \\ \mathbf{x}_{[k]} \\ \vdots \end{bmatrix}; \quad \begin{bmatrix} \mathbf{y}_{[i]} \\ \mathbf{y}_{[j]} \\ \mathbf{y}_{[k]} \\ \vdots \end{bmatrix}$$

We call this **stochastic gradient descent**

```

1:function STOCHASTIC GRADIENT DESCENT( $\mathbf{X}$ ,  $\mathbf{Y}$ ,  $\mathcal{L}$ ,  $t$ ,  $\alpha$ )
2:     $\theta \leftarrow$  Glorot()
3:    for  $i \in 1 \dots t$  do
4:         $\mathbf{X}, \mathbf{Y} \leftarrow$  Shuffle( $\mathbf{X}$ ), Shuffle( $\mathbf{Y}$ )
5:        for  $j \in 0 \dots \frac{n}{B} - 1$  do
6:             $\mathbf{X}_j \leftarrow [\mathbf{x}_{[jB]} \ \mathbf{x}_{[jB+1]} \ \dots \ \mathbf{x}_{[(j+1)B]}]$ 
7:             $\mathbf{Y}_j \leftarrow [\mathbf{y}_{[jB]} \ \mathbf{y}_{[jB+1]} \ \dots \ \mathbf{y}_{[(j+1)B]}]$ 
8:             $\mathbf{J} \leftarrow \nabla_{\theta} \mathcal{L}(\mathbf{X}_j, \mathbf{Y}_j, \theta)$ 
9:             $\theta \leftarrow \theta - \alpha \mathbf{J}$ 
10:   return  $\theta$ 

```

Stochastic gradient descent (SGD) is useful for saving memory

Stochastic gradient descent

Stochastic gradient descent (SGD) is useful for saving memory

But it can also improve performance

Stochastic gradient descent

Stochastic gradient descent (SGD) is useful for saving memory

But it can also improve performance

Since the “dataset” changes every update, so does the loss manifold

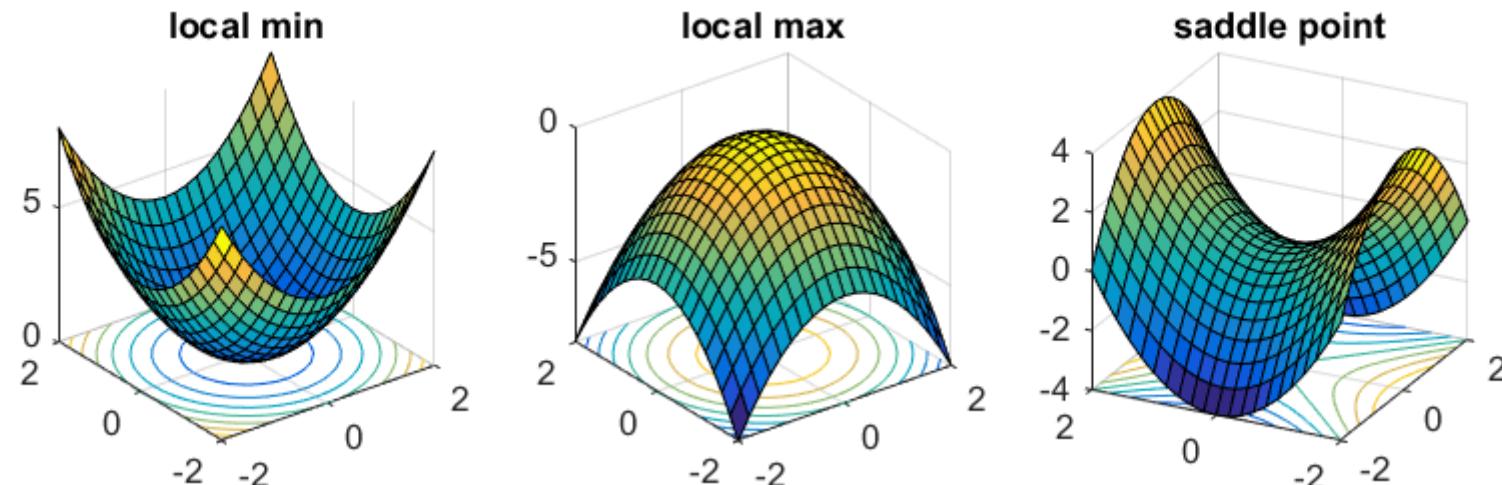
Stochastic gradient descent

Stochastic gradient descent (SGD) is useful for saving memory

But it can also improve performance

Since the “dataset” changes every update, so does the loss manifold

This makes it less likely we get stuck in bad optima



There is `torch.utils.data.DataLoader` to help

Stochastic gradient descent

There is `torch.utils.data.DataLoader` to help

```
import torch
dataloader = torch.utils.data.DataLoader(
    training_data,
    batch_size=32, # How many datapoints to sample
    shuffle=True, # Randomly shuffle each epoch
)
for epoch in number_of_epochs:
    for batch in dataloader:
        X_j, Y_j = batch
        loss = L(X_j, Y_j, theta)
        ...
    
```

1. Review
2. Dirty secret of deep learning
3. Optimization is hard
4. Deeper neural networks
5. Activation functions
6. Parameter initialization
7. **Stochastic gradient descent**
8. Modern optimization
9. Coding

1. Review
2. Dirty secret of deep learning
3. Optimization is hard
4. Deeper neural networks
5. Activation functions
6. Parameter initialization
7. Stochastic gradient descent
8. **Modern optimization**
9. Coding

Gradient descent is a powerful tool, but it has issues

Modern optimization

Gradient descent is a powerful tool, but it has issues

1. It can be slow to converge

Modern optimization

Gradient descent is a powerful tool, but it has issues

1. It can be slow to converge
2. It can get stuck in poor local optima

Modern optimization

Gradient descent is a powerful tool, but it has issues

1. It can be slow to converge
2. It can get stuck in poor local optima

Many researchers work on improving gradient descent to converge more quickly, while also preventing premature convergence

Modern optimization

Gradient descent is a powerful tool, but it has issues

1. It can be slow to converge
2. It can get stuck in poor local optima

Many researchers work on improving gradient descent to converge more quickly, while also preventing premature convergence

It is hard to teach adaptive optimization through math

Modern optimization

Gradient descent is a powerful tool, but it has issues

1. It can be slow to converge
2. It can get stuck in poor local optima

Many researchers work on improving gradient descent to converge more quickly, while also preventing premature convergence

It is hard to teach adaptive optimization through math

So first, I want to show you a video to prepare you

<https://www.youtube.com/watch?v=MD2fYip6QsQ&t=77s>

The video simulations provide an intuitive understanding of adaptive optimizers

Modern optimization

The video simulations provide an intuitive understanding of adaptive optimizers

The key behind modern optimizers is two concepts:

- Momentum

Modern optimization

The video simulations provide an intuitive understanding of adaptive optimizers

The key behind modern optimizers is two concepts:

- Momentum
- Adaptive learning rate

Modern optimization

The video simulations provide an intuitive understanding of adaptive optimizers

The key behind modern optimizers is two concepts:

- Momentum
- Adaptive learning rate

Let us discuss the algorithms more slowly

Review gradient descent again, because we will be making changes to it

Modern optimization

Review gradient descent again, because we will be making changes to it

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

Introduce **momentum** first

Modern optimization

Introduce **momentum** first

```
1:function Momentum GRADIENT DESCENT( $X, Y, \mathcal{L}, t, \alpha, \beta$ )
2:     $\theta \leftarrow$  Glorot()
3:     $M \leftarrow \mathbf{0}$  # Init momentum
4:    for  $i \in 1 \dots t$  do
5:         $J \leftarrow \nabla_{\theta} \mathcal{L}(X, Y, \theta)$  # Represents acceleration
6:         $M \leftarrow \beta \cdot M + (1 - \beta) \cdot J$  # Momentum and acceleration
7:         $\theta \leftarrow \theta - \alpha M$ 
8:    return  $\theta$ 
```

Now adaptive learning rate

Modern optimization

Now adaptive learning rate

```
1:function RMSProp( $X, Y, \mathcal{L}, t, \alpha, \beta, \varepsilon$ )
2:     $\theta \leftarrow$  Glorot()
3:     $V \leftarrow 0$  # Init variance
4:    for  $i \in 1 \dots t$  do
5:         $J \leftarrow \nabla_{\theta} \mathcal{L}(X, Y, \theta)$  # Represents acceleration
6:         $V \leftarrow \beta \cdot V + (1 - \beta) \cdot J \odot J$  # Magnitude
7:         $\theta \leftarrow \theta - \alpha J \oslash \sqrt[V]{V + \varepsilon}$  # Rescale grad by prev updates
8:    return  $\theta$ 
```

Combine **momentum** and **adaptive learning rate** to create **Adam**

Modern optimization

Combine **momentum** and **adaptive learning rate** to create **Adam**

```
1:function ADAPTIVE MOMENT ESTIMATION( $X, Y, \mathcal{L}, t, \alpha, \beta_1, \beta_2, \varepsilon$ )
2:     $\theta \leftarrow$  Glorot()
3:     $M, V \leftarrow 0$ 
4:    for  $i \in 1 \dots t$  do
5:         $J \leftarrow \nabla_{\theta} \mathcal{L}(X, Y, \theta)$ 
6:         $M \leftarrow \beta_1 M + (1 - \beta_1) J$  # Compute momentum
7:         $V \leftarrow \beta_2 \cdot V + (1 - \beta_2) \cdot J \odot J$  # Magnitude
8:         $\theta \leftarrow \theta - \alpha M \oslash \sqrt[V]{V + \varepsilon}$  # Adaptive param update
9:    return  $\theta$  # Note, we use biased  $M, V$  for clarity
```

```
import torch
betas = (0.9, 0.999)
net = ...
theta = net.parameters()

sgd = torch.optim.SGD(theta, lr=alpha)
momentum = torch.optim.SGD(
    theta, lr=alpha, momentum=betas[0])
rmsprop = torch.optim.RMSprop(
    theta, lr=alpha, momentum=betas[1])
adam = torch.optim.Adam(theta, lr=alpha, betas=betas)

...
sgd.step(), momentum.step(), rmsprop.step(), adam.step()
```

```
import optax
betas = (0.9, 0.999)
theta = ...

sgd = optax.sgd(lr=alpha)
momentum = optax.sgd(lr=alpha, momentum=betas[0])
rmsprop = optax.rmsprop(lr=alpha, decay=betas[1])
adam = optax.adam(lr=alpha, b1=betas[0], b2=betas[1])

v = rmsprop.init(theta)
theta, v = rmsprop.update(J, v, theta)
mv = adam.init(theta) # contains M and V
theta, mv = mv.update(J, mv, theta)
```

1. Review
2. Dirty secret of deep learning
3. Optimization is hard
4. Deeper neural networks
5. Activation functions
6. Parameter initialization
7. Stochastic gradient descent
8. **Modern optimization**
9. Coding

1. Review
2. Dirty secret of deep learning
3. Optimization is hard
4. Deeper neural networks
5. Activation functions
6. Parameter initialization
7. Stochastic gradient descent
8. Modern optimization
9. **Coding**

https://colab.research.google.com/drive/1qTNSvB_JEMnMJfcAwsLJTullfxa_kyTD#scrollTo=YVkCyz78x4Rp

