

# Classification

CISC 7026: Introduction to Deep Learning

University of Macau

We will have a make-up lecture later on for the missed lecture

# Admin

We will have a make-up lecture later on for the missed lecture

Assignment 1 grades were released on moodle

# Admin

We will have a make-up lecture later on for the missed lecture

Assignment 1 grades were released on moodle

The scores were very good, with a mean of approximately 90/100

# Admin

We will have a make-up lecture later on for the missed lecture

Assignment 1 grades were released on moodle

The scores were very good, with a mean of approximately 90/100

I am still grading quiz 2, but I had a look at the responses to question 4

Some requests from students:

# Admin

Some requests from students:

1. More coding, less theory

Some requests from students:

1. More coding, less theory
2. More math/theory



# Admin

Some requests from students:

1. More coding, less theory
2. More math/theory
3. Too easy, go faster

Some requests from students:

1. More coding, less theory
2. More math/theory
3. Too easy, go faster
4. Speak slower

Some requests from students:

1. More coding, less theory
2. More math/theory
3. Too easy, go faster
4. Speak slower
5. Course is too hard

Some requests from students:

1. More coding, less theory
2. More math/theory
3. Too easy, go faster
4. Speak slower
5. Course is too hard
6. Course is perfect

Some requests from students:

1. More coding, less theory
2. More math/theory
3. Too easy, go faster
4. Speak slower
5. Course is too hard
6. Course is perfect
7. Move captions to top of screen

Some requests from students:

1. More coding, less theory
2. More math/theory
3. Too easy, go faster
4. Speak slower
5. Course is too hard
6. Course is perfect
7. Move captions to top of screen
8. Upload powerpoint before lecture

Some requests from students:

1. More coding, less theory
2. More math/theory
3. Too easy, go faster
4. Speak slower
5. Course is too hard
6. Course is perfect
7. Move captions to top of screen
8. Upload powerpoint before lecture

There are conflicting student needs

[https://github.com/smorad/um\\_cisc\\_7026](https://github.com/smorad/um_cisc_7026)



1. Review
2. Torch optimization coding
3. Classification task
4. Probability review
5. Define model  $f$
6. Define loss function  $\mathcal{L}$
7. Find  $\theta$  that minimize  $\mathcal{L}$
8. Coding

1. **Review**
2. Torch optimization coding
3. Classification task
4. Probability review
5. Define model  $f$
6. Define loss function  $\mathcal{L}$
7. Find  $\theta$  that minimize  $\mathcal{L}$
8. Coding

# Review

Last time, we reviewed derivatives

# Review

Last time, we reviewed derivatives

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

# Review

Last time, we reviewed derivatives

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

and gradients

# Review

Last time, we reviewed derivatives

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

and gradients

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

# Review

Gradients are important in deep learning for two reasons:

# Review

Gradients are important in deep learning for two reasons:

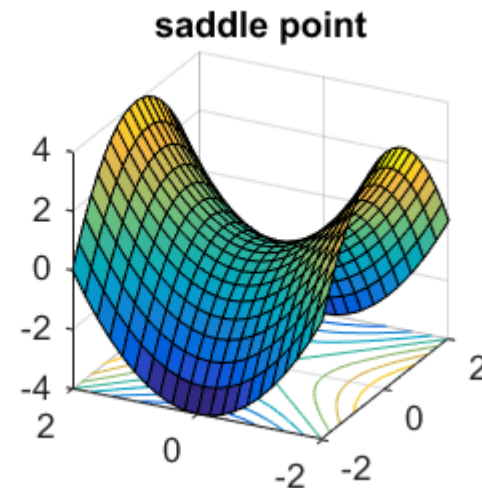
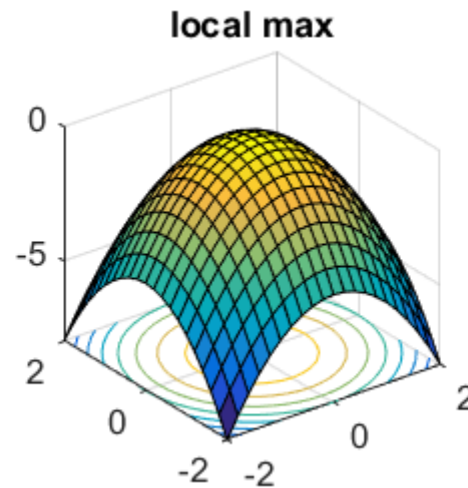
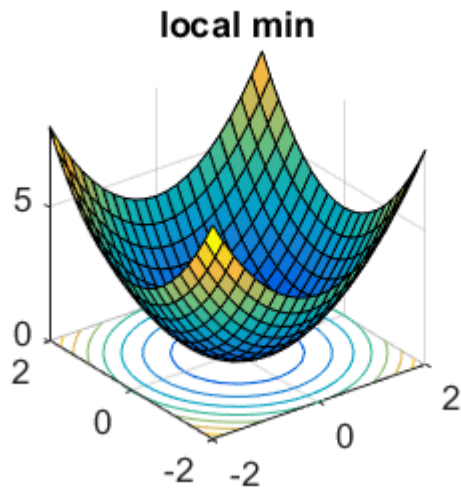
**Reason 1:**  $f(\boldsymbol{x})$  has critical points at  $\nabla_{\boldsymbol{x}} f(\boldsymbol{x}) = 0$



# Review

Gradients are important in deep learning for two reasons:

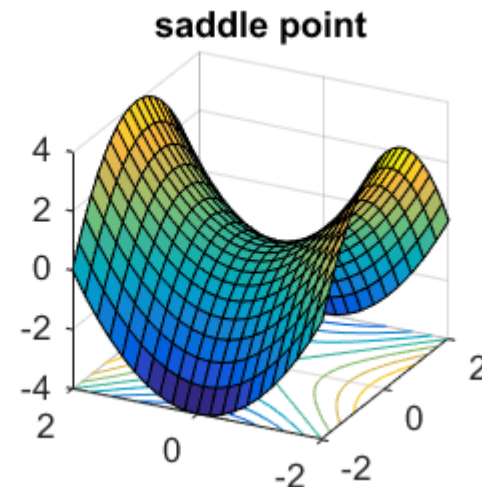
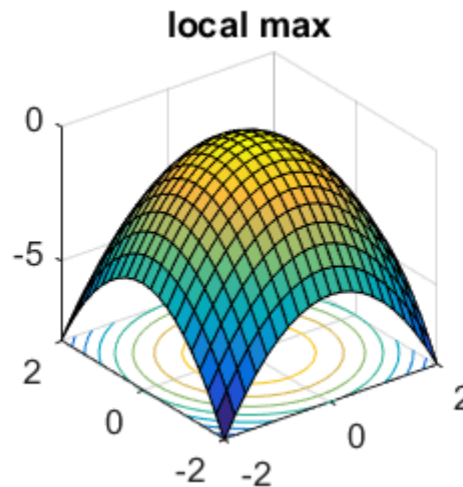
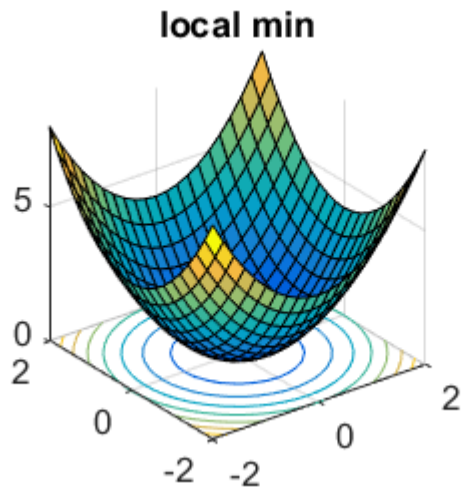
**Reason 1:**  $f(x)$  has critical points at  $\nabla_x f(x) = 0$



# Review

Gradients are important in deep learning for two reasons:

**Reason 1:**  $f(x)$  has critical points at  $\nabla_x f(x) = 0$



With optimization, we attempt to find minima of loss functions

# Review

Gradients are important in deep learning for two reasons:

# Review

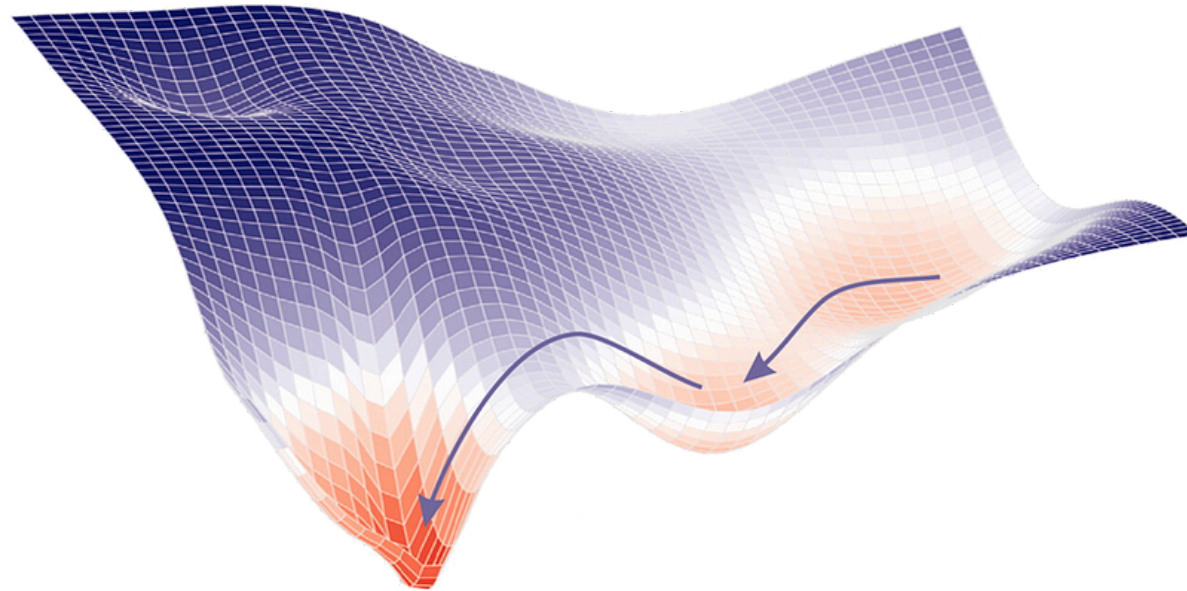
Gradients are important in deep learning for two reasons:

**Reason 2:** For problems without analytical solutions, the gradient (slope) is necessary for gradient descent

# Review

Gradients are important in deep learning for two reasons:

**Reason 2:** For problems without analytical solutions, the gradient (slope) is necessary for gradient descent



# Review

First, we derived the solution to linear regression

# Review

First, we derived the solution to linear regression

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

# Review

First, we derived the solution to linear regression

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

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



# Review

First, we derived the solution to linear regression

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

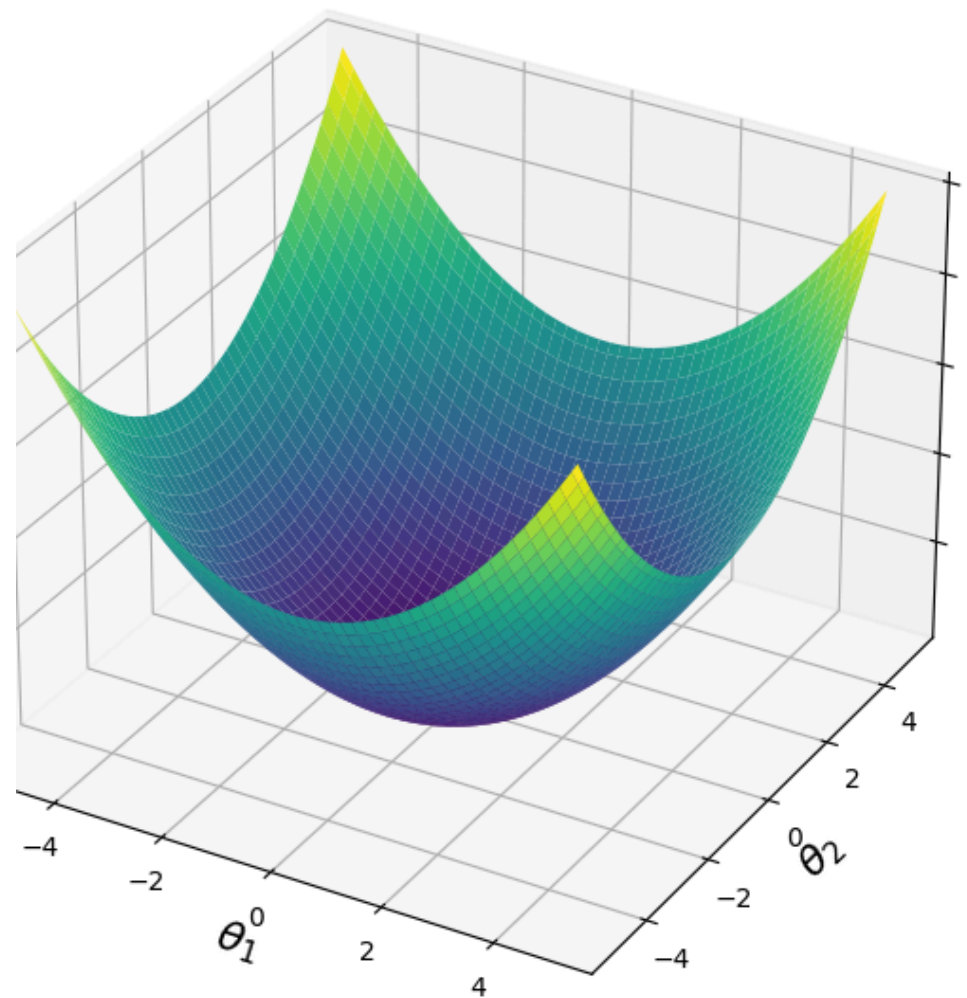
$$\mathcal{L}(\mathbf{X}, \mathbf{Y}, \boldsymbol{\theta}) = (\mathbf{Y} - \mathbf{X}_D \boldsymbol{\theta})^\top (\mathbf{Y} - \mathbf{X}_D \boldsymbol{\theta})$$

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

Quadratic function of  $\boldsymbol{\theta}$

# Review

A quadratic function has a single critical point, which must be a global minimum



# Review

We found the analytical solution for linear regression by finding where the gradient was zero and solving for  $\theta$

# Review

We found the analytical solution for linear regression by finding where the gradient was zero and solving for  $\theta$

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

# Review

We found the analytical solution for linear regression by finding where the gradient was zero and solving for  $\theta$

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

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

# Review

We found the analytical solution for linear regression by finding where the gradient was zero and solving for  $\theta$

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

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

Which solves

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

# Review

For neural networks, the square error loss is no longer quadratic

# Review

For neural networks, the square error loss is no longer quadratic

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

Loss function



# Review

For neural networks, the square error loss is no longer quadratic

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

Loss function

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

Neural network model

# Review

For neural networks, the square error loss is no longer quadratic

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

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

# Review

For neural networks, the square error loss is no longer quadratic

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

$$f(x, \boldsymbol{\theta}) = \sigma(\theta_0 + \theta_1 x) \quad \text{Neural network model}$$

Now, we plug the model  $f$  into the loss function

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

# Review

For neural networks, the square error loss is no longer quadratic

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

$$f(x, \boldsymbol{\theta}) = \sigma(\theta_0 + \theta_1 x) \quad \text{Neural network model}$$

Now, we plug the model  $f$  into the loss function

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

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

# Review

For neural networks, the square error loss is no longer quadratic

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

$$f(x, \boldsymbol{\theta}) = \sigma(\theta_0 + \theta_1 x) \quad \text{Neural network model}$$

Now, we plug the model  $f$  into the loss function

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

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

There is no analytical solution for  $\boldsymbol{\theta}$

# Review

Instead, we found the parameters of a neural network through gradient descent

# Review

Instead, we found the parameters of a neural network through gradient descent

Gradient descent is an optimization method for differentiable functions

# Review

Instead, we found the parameters of a neural network through gradient descent

Gradient descent is an optimization method for differentiable functions

We went over both the intuition and mathematical definitions



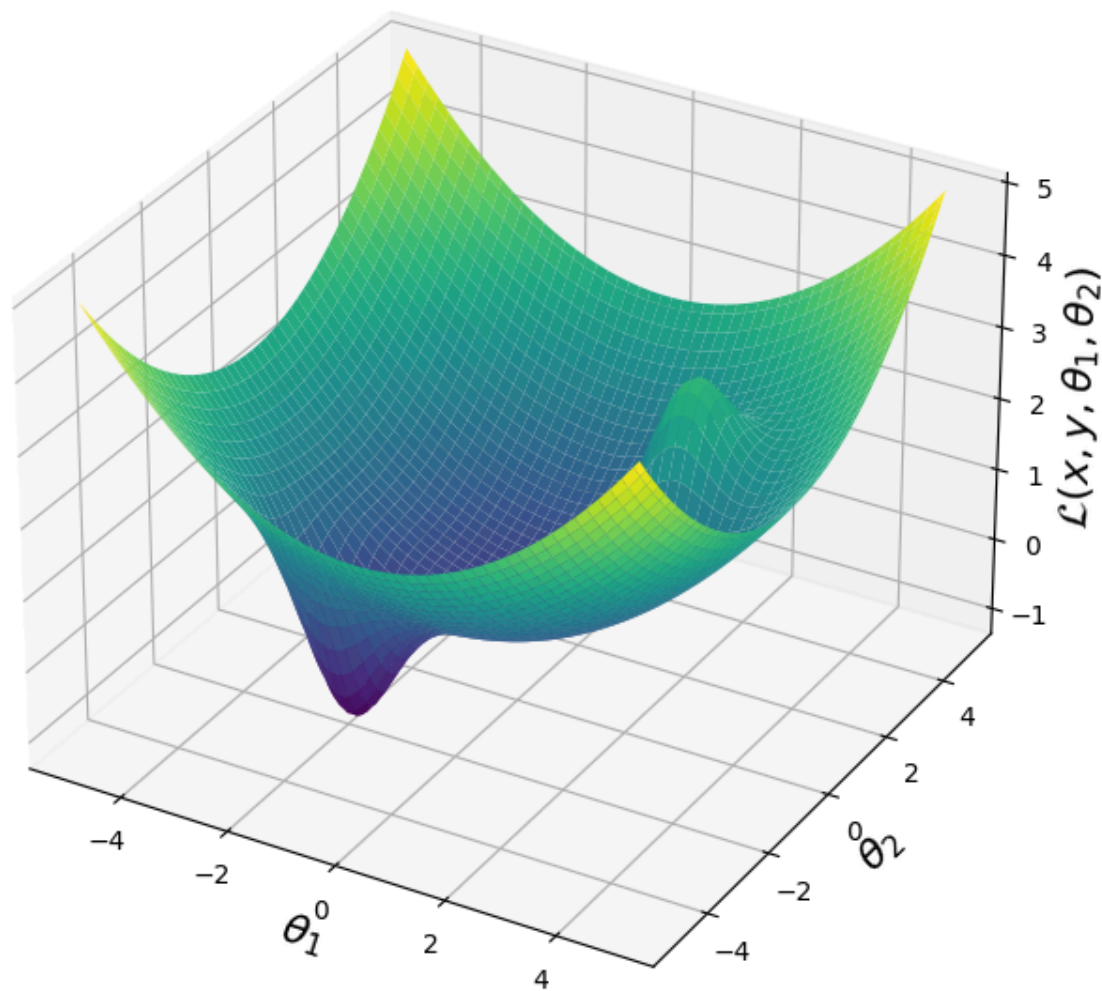
# Review



# Review



# Review



# Review

The gradient descent algorithm:

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

# Review

We derived the  $\nabla_{\theta} \mathcal{L}$  for deep neural networks using the chain rule

# Review

We derived the  $\nabla_{\theta} \mathcal{L}$  for deep neural networks using the chain rule

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

# Review

We derived the  $\nabla_{\theta} \mathcal{L}$  for deep neural networks using the chain rule

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

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

# Review

We derived the  $\nabla_{\theta} \mathcal{L}$  for deep neural networks using the chain rule

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

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

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



# Review

We ran into issues computing the gradient of a layer because of the Heaviside step function

# Review

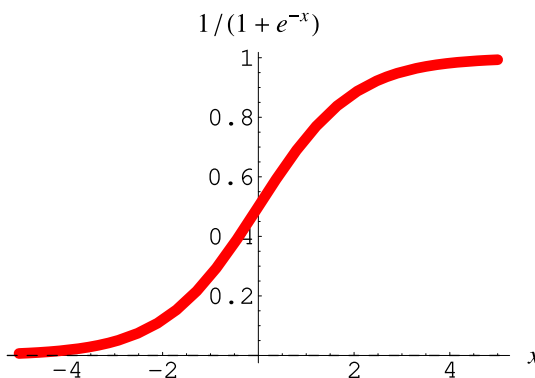
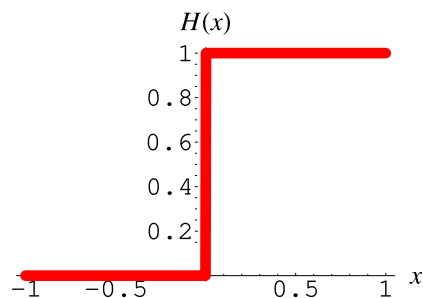
We ran into issues computing the gradient of a layer because of the Heaviside step function

We replaced it with a differentiable (soft) approximation called the sigmoid function

# Review

We ran into issues computing the gradient of a layer because of the Heaviside step function

We replaced it with a differentiable (soft) approximation called the sigmoid function



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

# Review

In `jax`, we compute the gradient using the `jax.grad` function

# Review

In jax, we compute the gradient using the `jax.grad` function

```
import jax
```

```
def L(theta, X, Y):
```

```
    ...
```

```
# Create a new function that is the gradient of L
```

```
# Then compute gradient of L for given inputs
```

```
J = jax.grad(L)(X, Y, theta)
```

```
# Update parameters
```

```
alpha = 0.0001
```

```
theta = theta - alpha * J
```

# Review

In torch, we backpropagate through a graph of operations

```
import torch
optimizer = torch.optim.SGD(lr=0.0001)

def L(model, X, Y):
    ...
    # Pytorch will record a graph of all operations
    # Everytime you do theta @ x, it stores inputs and outputs
    loss = L(X, Y, model) # compute loss
    # Traverse the graph backward and compute the gradient
    loss.backward() # Sets .grad attribute on each parameter
    optimizer.step() # Update the parameters using .grad
    optimizer.zero_grad() # Set .grad to zero, DO NOT FORGET!!
```

# Agenda

1. **Review**
2. Torch optimization coding
3. Classification task
4. Probability review
5. Define model  $f$
6. Define loss function  $\mathcal{L}$
7. Find  $\theta$  that minimize  $\mathcal{L}$
8. Coding

# Agenda

1. Review
2. **Torch optimization coding**
3. Classification task
4. Probability review
5. Define model  $f$
6. Define loss function  $\mathcal{L}$
7. Find  $\theta$  that minimize  $\mathcal{L}$
8. Coding



# Torch Optimization Coding

First, a video of one application of gradient descent

[https://youtu.be/kGDO2e\\_qiyI?si=ZopZKy-6WQ4B0csX](https://youtu.be/kGDO2e_qiyI?si=ZopZKy-6WQ4B0csX)

# Torch Optimization Coding

First, a video of one application of gradient descent

[https://youtu.be/kGDO2e\\_qiyI?si=ZopZKy-6WQ4B0csX](https://youtu.be/kGDO2e_qiyI?si=ZopZKy-6WQ4B0csX)

Time for some interactive coding

[https://colab.research.google.com/drive/1W8WVZ8n\\_9yJCcOqkPVURp\\_wJUx3EQc5w](https://colab.research.google.com/drive/1W8WVZ8n_9yJCcOqkPVURp_wJUx3EQc5w)

1. Review
2. **Torch optimization coding**
3. Classification task
4. Probability review
5. Define model  $f$
6. Define loss function  $\mathcal{L}$
7. Find  $\theta$  that minimize  $\mathcal{L}$
8. Coding

1. Review
2. Torch optimization coding
3. **Classification task**
4. Probability review
5. Define model  $f$
6. Define loss function  $\mathcal{L}$
7. Find  $\theta$  that minimize  $\mathcal{L}$
8. Coding

# Classification

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

# Classification

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

**Regression** asks how many

# Classification

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

**Regression** asks how many

- How long will I live?

# Classification

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?



# Classification

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

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

# Classification

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?

# Classification

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?

# Classification

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?

# Classification

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?

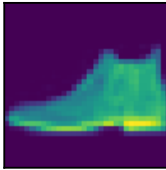
So far, we only looked at regression. Now, let us look at classification

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

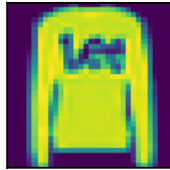
**Task:** Given a picture of clothes, predict the text description

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

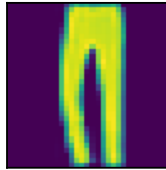
ankle boot



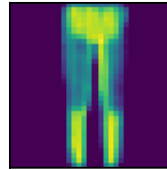
pullover



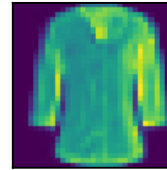
trouser



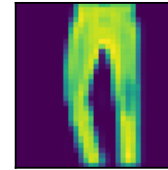
trouser



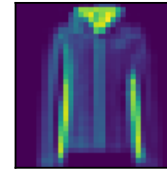
shirt



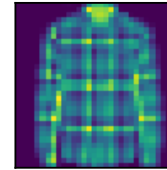
trouser



coat

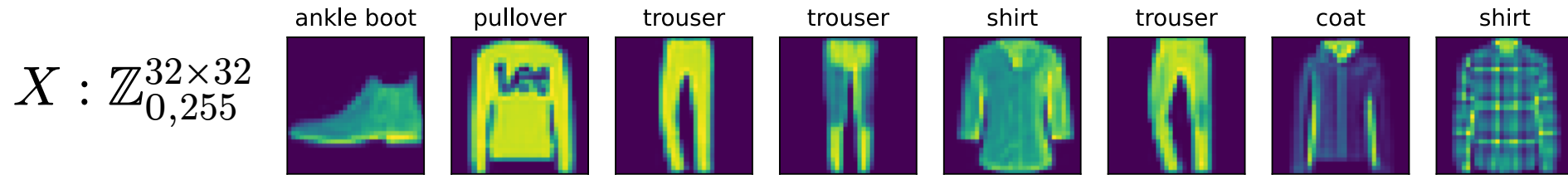


shirt



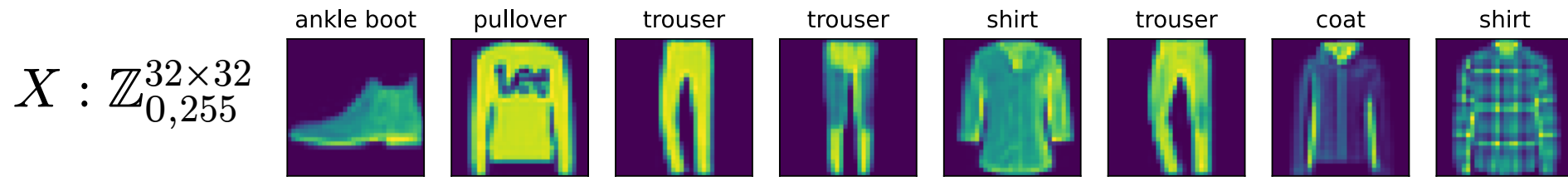


**Task:** Given a picture of clothes, predict the text description



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

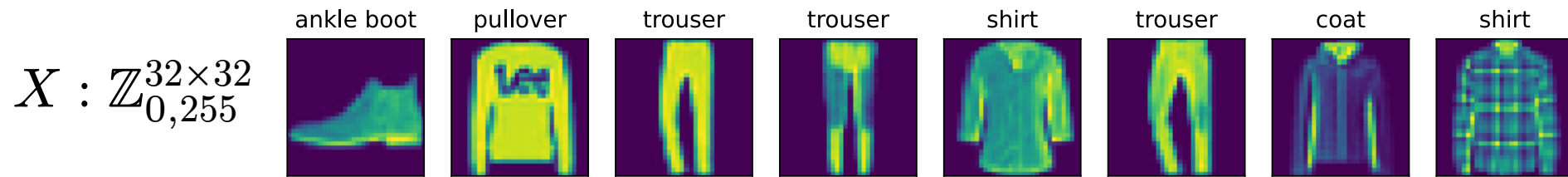
**Task:** Given a picture of clothes, predict the text description



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

**Approach:** Learn  $\theta$  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  $\theta$  that produce **conditional probabilities**

$$f(x, \theta) = P(y \mid x) = P\left(\begin{bmatrix} \text{T-Shirt} \\ \text{Trouser} \\ \vdots \end{bmatrix} \mid \begin{bmatrix} \text{img} \end{bmatrix}\right) = \begin{bmatrix} 0.2 \\ 0.01 \\ \vdots \end{bmatrix}$$

# Agenda

1. Review
2. Torch optimization coding
3. **Classification task**
4. Probability review
5. Define model  $f$
6. Define loss function  $\mathcal{L}$
7. Find  $\theta$  that minimize  $\mathcal{L}$
8. Coding

# Agenda

1. Review
2. Torch optimization coding
3. Classification task
4. **Probability review**
5. Define model  $f$
6. Define loss function  $\mathcal{L}$
7. Find  $\theta$  that minimize  $\mathcal{L}$
8. Coding

Classification tasks require **probability theory**, so let us review

Classification tasks require **probability theory**, so let us review

In probability, we have **experiments** and **outcomes**

Classification tasks require **probability theory**, so let us review

In probability, we have **experiments** and **outcomes**

An experiment yields one of many possible outcomes



Classification tasks require **probability theory**, so let us review

In probability, we have **experiments** and **outcomes**

An experiment yields one of many possible outcomes

Experiment

Outcome

Classification tasks require **probability theory**, so let us review

In probability, we have **experiments** and **outcomes**

An experiment yields one of many possible outcomes

Experiment

Outcome

Flip a coin

Heads

Classification tasks require **probability theory**, so let us review

In probability, we have **experiments** and **outcomes**

An experiment yields one of many possible outcomes

Experiment

Outcome

Flip a coin

Heads

Walk outside

Rain

Classification tasks require **probability theory**, so let us review

In probability, we have **experiments** and **outcomes**

An experiment yields one of many possible outcomes

Experiment	Outcome
Flip a coin	Heads
Walk outside	Rain
Grab clothing from closet	Coat

The **sample space**  $S$  defines all possible outcomes for an experiment

The **sample space**  $S$  defines all possible outcomes for an experiment

Experiment

Sample Space  $S$

The **sample space**  $S$  defines all possible outcomes for an experiment

Experiment

Sample Space  $S$

Flip a coin

$S = \{\text{heads}, \text{tails}\}$

The **sample space**  $S$  defines all possible outcomes for an experiment

Experiment

Sample Space  $S$

Flip a coin

$S = \{\text{heads, tails}\}$

Walk outside

$S = \{\text{rain, sun, wind, cloud}\}$



The **sample space**  $S$  defines all possible outcomes for an experiment

Experiment

Sample Space  $S$

Flip a coin

$S = \{\text{heads, tails}\}$

Walk outside

$S = \{\text{rain, sun, wind, cloud}\}$

Take clothing from closet

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

An **event**  $E$  is a specific subset of the sample space

An **event**  $E$  is a specific subset of the sample space

Experiment

Sample Space

Event

An **event**  $E$  is a specific subset of the sample space

Experiment

Sample Space

Event

Flip a coin

$S = \{\text{heads}, \text{tails}\}$

$E = \{\text{heads}\}$

An **event**  $E$  is a specific subset of the sample space

Experiment

Sample Space

Event

Flip a coin

$S = \{\text{heads, tails}\}$

$E = \{\text{heads}\}$

Walk outside

$S = \{\text{rain, sun, wind, cloud}\}$

$E = \{\text{rain, wind}\}$

An **event**  $E$  is a specific subset of the sample space

Experiment	Sample Space	Event
Flip a coin	$S = \{\text{heads, tails}\}$	$E = \{\text{heads}\}$
Walk outside	$S = \{\text{rain, sun, wind, cloud}\}$	$E = \{\text{rain, wind}\}$
Take from closet	$S = \{\text{T-shirt, Trouser, Pullover, Dress, Coat, Sandal, Shirt, Sneaker, Bag, Ankle boot}\}$	$E = \{\text{Shirt, T-Shirt, Coat}\}$

The **probability** measures how likely an event is to occur

The **probability** measures how likely an event is to occur

The probability must be between 0 (never occurs) and 1 (always occurs)



The **probability** measures how likely an event is to occur

The probability must be between 0 (never occurs) and 1 (always occurs)

$$0 \leq P(E) \leq 1$$

The **probability** measures how likely an event is to occur

The probability must be between 0 (never occurs) and 1 (always occurs)

$$0 \leq P(E) \leq 1$$

Experiment

Probabilities

The **probability** measures how likely an event is to occur

The probability must be between 0 (never occurs) and 1 (always occurs)

$$0 \leq P(E) \leq 1$$

Experiment

Probabilities

Flip a coin

$$P(\text{heads}) = 0.5$$

The **probability** measures how likely an event is to occur

The probability must be between 0 (never occurs) and 1 (always occurs)

$$0 \leq P(E) \leq 1$$

Experiment

Probabilities

Flip a coin

$$P(\text{heads}) = 0.5$$

Walk outside

$$P(\text{rain}) = 0.15$$

The **probability** measures how likely an event is to occur

The probability must be between 0 (never occurs) and 1 (always occurs)

$$0 \leq P(E) \leq 1$$

Experiment

Probabilities

Flip a coin

$$P(\text{heads}) = 0.5$$

Walk outside

$$P(\text{rain}) = 0.15$$

Take from closet

$$P(\text{Shirt}) = 0.1$$

When we define  $P$  as a function, we call it a **distribution**

When we define  $P$  as a function, we call it a **distribution**

$$P : E \mapsto (0, 1)$$

When we define  $P$  as a function, we call it a **distribution**

$$P : E \mapsto (0, 1)$$

The probabilities (distribution) over the sample space  $S$  must sum to one



When we define  $P$  as a function, we call it a **distribution**

$$P : E \mapsto (0, 1)$$

The probabilities (distribution) over the sample space  $S$  must sum to one

$$\sum_{x \in S} P(x) = 1$$

When we define  $P$  as a function, we call it a **distribution**

$$P : E \mapsto (0, 1)$$

The probabilities (distribution) over the sample space  $S$  must sum to one

$$\sum_{x \in S} P(x) = 1$$

Flip a coin  $\{P(\text{heads}) = 0.5, P(\text{tails}) = 0.5\}$

When we define  $P$  as a function, we call it a **distribution**

$$P : E \mapsto (0, 1)$$

The probabilities (distribution) over the sample space  $S$  must sum to one

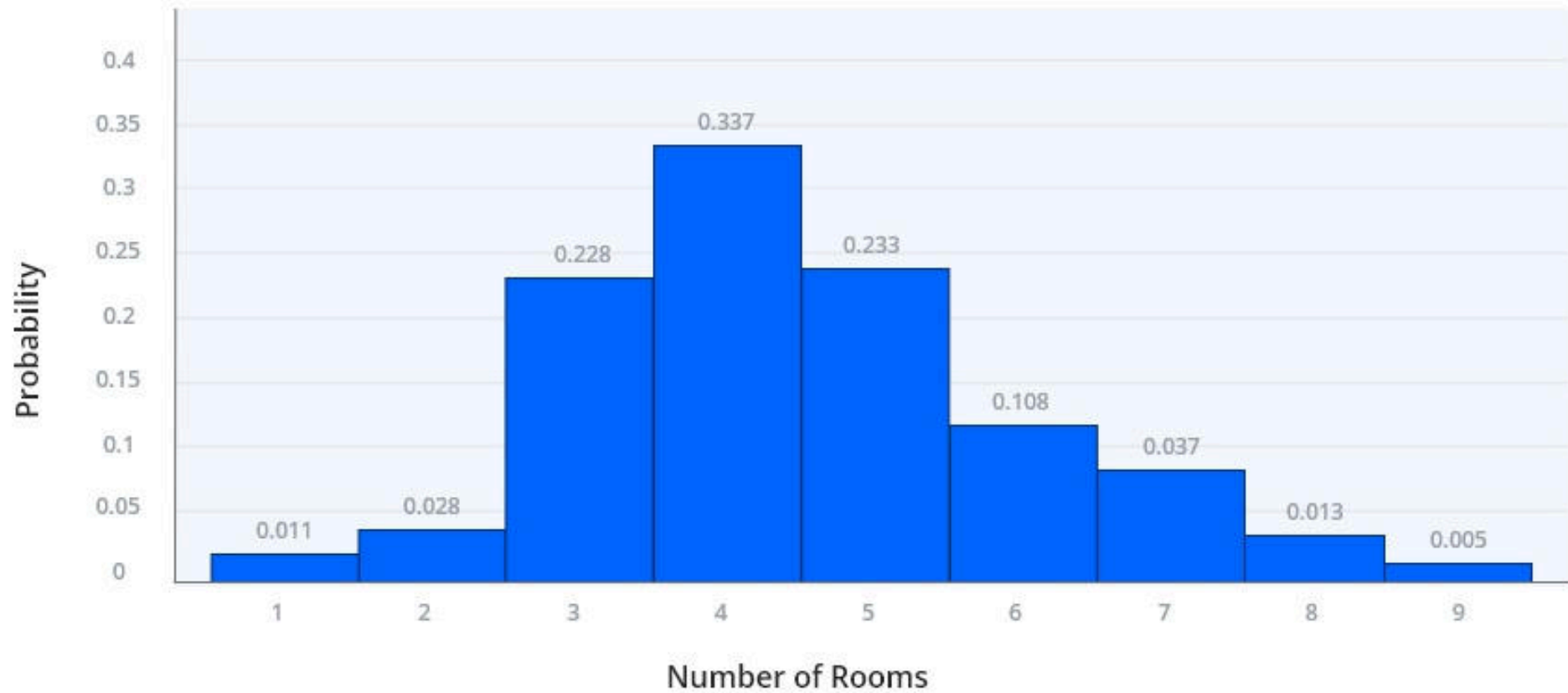
$$\sum_{x \in S} P(x) = 1$$

Flip a coin  $\{P(\text{heads}) = 0.5, P(\text{tails}) = 0.5\}$

Take clothing from closet  $\{P(\text{T-shirt}) = 0.1, P(\text{Trouser}) = 0.08,$   
 $P(\text{Pullover}) = 0.12, \dots\}$

The distribution is a function, so we can plot it

## Number of Rooms in Rental Unit



Events can overlap with each other

Events can overlap with each other

- Disjoint events

Events can overlap with each other

- Disjoint events
- Conditionally dependent events



Two events  $A, B$  are **disjoint** if either  $A$  occurs or  $B$  occurs, **but not both**

Two events  $A, B$  are **disjoint** if either  $A$  occurs or  $B$  occurs, **but not both**

With disjoint events,  $P(A \cap B) = 0$

Two events  $A, B$  are **disjoint** if either  $A$  occurs or  $B$  occurs, **but not both**

With disjoint events,  $P(A \cap B) = 0$

Flip a coin

$$P(\text{Heads}) = 0.5, P(\text{Tails}) = 0.5$$

$$P(\text{Heads} \cap \text{Tails}) = 0$$

Two events  $A, B$  are **disjoint** if either  $A$  occurs or  $B$  occurs, **but not both**

With disjoint events,  $P(A \cap B) = 0$

Flip a coin

$$P(\text{Heads}) = 0.5, P(\text{Tails}) = 0.5$$
$$P(\text{Heads} \cap \text{Tails}) = 0$$

Be careful!

Walk outside

$$P(\text{Rain}) = 0.1, P(\text{Cloud}) = 0.2$$
$$P(\text{Rain} \cap \text{Cloud}) \neq 0$$

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

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

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

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

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

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

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

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

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

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



If  $A, B$  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

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

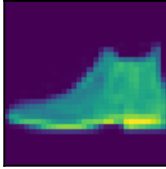
$$P(\text{Rain} \mid \text{Cloud}) = \frac{0.1}{0.2} = 0.5$$

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

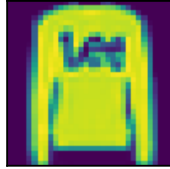
**Task:** Given a picture of clothes, predict the text description

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

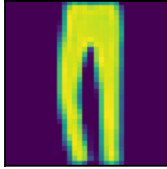
ankle boot



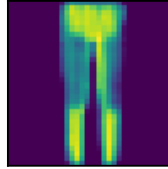
pullover



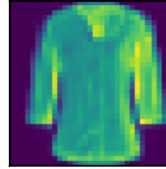
trouser



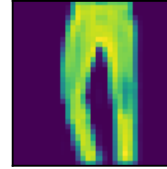
trouser



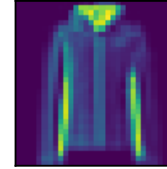
shirt



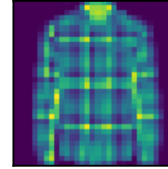
trouser



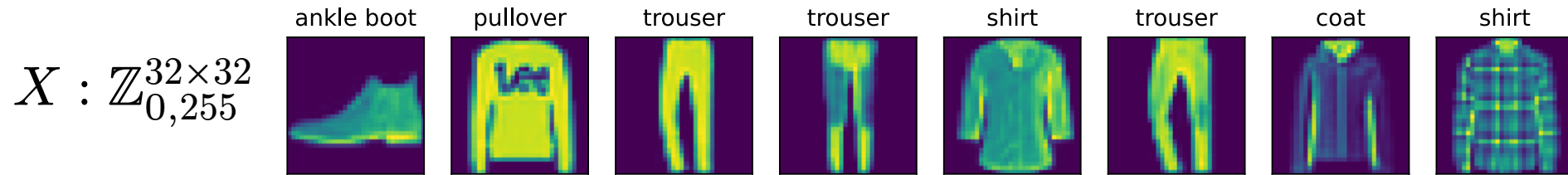
coat



shirt

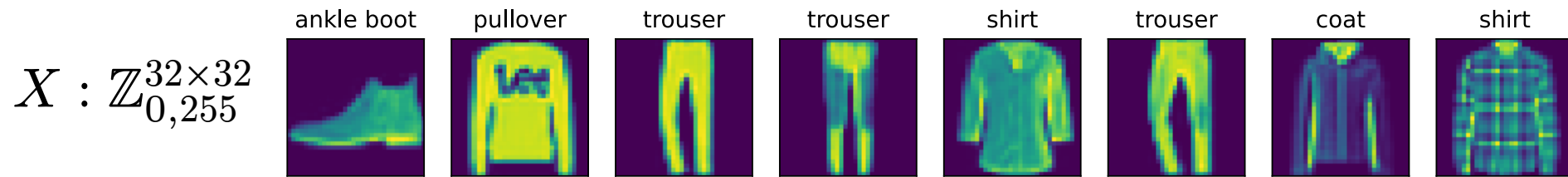


**Task:** Given a picture of clothes, predict the text description



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

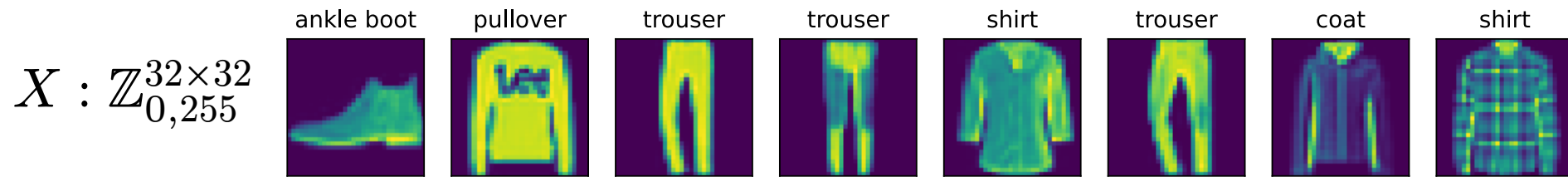
**Task:** Given a picture of clothes, predict the text description



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

**Approach:** Learn  $\theta$  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  $\theta$  that produce **conditional probabilities**

$$f(x, \theta) = P(y \mid x) = P\left(\begin{bmatrix} \text{T-Shirt} \\ \text{Trouser} \\ \vdots \end{bmatrix} \mid \begin{bmatrix} \text{img} \end{bmatrix}\right) = \begin{bmatrix} 0.2 \\ 0.01 \\ \vdots \end{bmatrix}$$

# Agenda

1. Review
2. Torch optimization coding
3. Classification task
4. **Probability review**
5. Define model  $f$
6. Define loss function  $\mathcal{L}$
7. Find  $\theta$  that minimize  $\mathcal{L}$
8. Coding

# Agenda

1. Review
2. Torch optimization coding
3. Classification task
4. Probability review
5. **Define model  $f$**
6. Define loss function  $\mathcal{L}$
7. Find  $\theta$  that minimize  $\mathcal{L}$
8. Coding



We will again start with a multivariate linear model

We will again start with a multivariate linear model

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

We will again start with a multivariate linear model

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

We want our model to predict the probability of each outcome

We will again start with a multivariate linear model

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

We want our model to predict the probability of each outcome

**Question:** What is the function signature of  $f$ ?

We will again start with a multivariate linear model

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

We want our model to predict the probability of each outcome

**Question:** What is the function signature of  $f$ ?

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

We will again start with a multivariate linear model

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

We want our model to predict the probability of each outcome

**Question:** What is the function signature of  $f$ ?

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

**Question:** Can we use this model to predict probabilities?

We will again start with a multivariate linear model

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

We want our model to predict the probability of each outcome

**Question:** What is the function signature of  $f$ ?

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

**Question:** Can we use this model to predict probabilities?

**Answer:** No! Because probabilities must be  $\in (0, 1)$  and sum to one

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



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

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

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

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

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

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

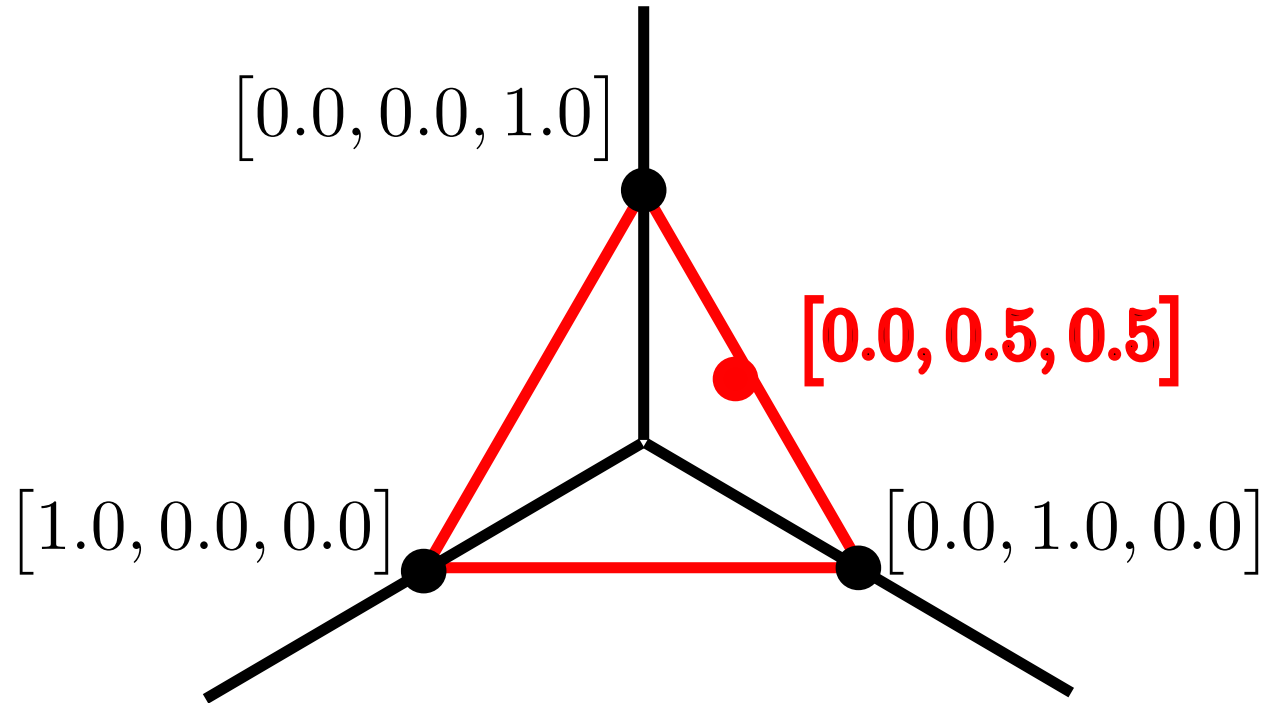
$$\mathbf{v} = \left\{ \begin{bmatrix} v_1 \\ \vdots \\ v_{d_y} \end{bmatrix} \mid \sum_{i=1}^{d_y} v_i = 1; \quad v_i \in (0, 1) \right\}$$

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

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

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

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



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

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

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

So we need a function that maps to the simplex

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

So we need a function that maps to the simplex

$$g : \mathbb{R}^{d_y} \mapsto \Delta^{d_y-1}$$



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

So we need a function that maps to the simplex

$$g : \mathbb{R}^{d_y} \mapsto \Delta^{d_y-1}$$

Then, we can combine  $f$  and  $g$

$$g(f) : \mathbb{R}^{d_x} \times \Theta \mapsto \Delta^{d_y-1}$$

We need a function  $g$

$$g : \mathbb{R}^{d_y} \mapsto \Delta^{d_y-1}$$

We need a function  $g$

$$g : \mathbb{R}^{d_y} \mapsto \Delta^{d_y-1}$$

There are many functions that can do this

We need a function  $g$

$$g : \mathbb{R}^{d_y} \mapsto \Delta^{d_y-1}$$

There are many functions that can do this

One example is dividing by the  $L_1$  norm:

$$g(\mathbf{x}) = \frac{\mathbf{x}}{\sum_{i=1}^{d_y} x_i}$$

We need a function  $g$

$$g : \mathbb{R}^{d_y} \mapsto \Delta^{d_y-1}$$

There are many functions that can do this

One example is dividing by the  $L_1$  norm:

$$g(\mathbf{x}) = \frac{\mathbf{x}}{\sum_{i=1}^{d_y} x_i}$$

In deep learning we often use the **softmax** function. When combined with the classification loss the gradient is linear, making learning faster

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)

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

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

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

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

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

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

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



And naturally, we can use the same method for a deep neural network

$$\begin{aligned} f_1(\mathbf{x}, \boldsymbol{\varphi}) &= \sigma(\boldsymbol{\varphi}^\top \overline{\mathbf{x}}) \\ &\vdots \\ f_\ell(\mathbf{x}, \boldsymbol{\xi}) &= \text{softmax}(\boldsymbol{\xi}^\top \overline{\mathbf{x}}) \end{aligned}$$

And naturally, we can use the same method for a deep neural network

$$\begin{aligned} f_1(\mathbf{x}, \boldsymbol{\varphi}) &= \sigma(\boldsymbol{\varphi}^\top \overline{\mathbf{x}}) \\ &\vdots \\ f_\ell(\mathbf{x}, \boldsymbol{\xi}) &= \text{softmax}(\boldsymbol{\xi}^\top \overline{\mathbf{x}}) \end{aligned}$$

Now, our neural network can output probabilities

$$f(\mathbf{x}, \boldsymbol{\theta}) = \begin{bmatrix} P(\text{Ankle boot} \mid \text{Ankle boot}) \\ P(\text{Bag} \mid \text{Ankle boot}) \\ \vdots \end{bmatrix} = \begin{bmatrix} 0.25 \\ 0.08 \\ \vdots \end{bmatrix}$$

**Question:** Why do we output probabilities instead of a binary values

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

**Question:** Why do we output probabilities instead of a binary values

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

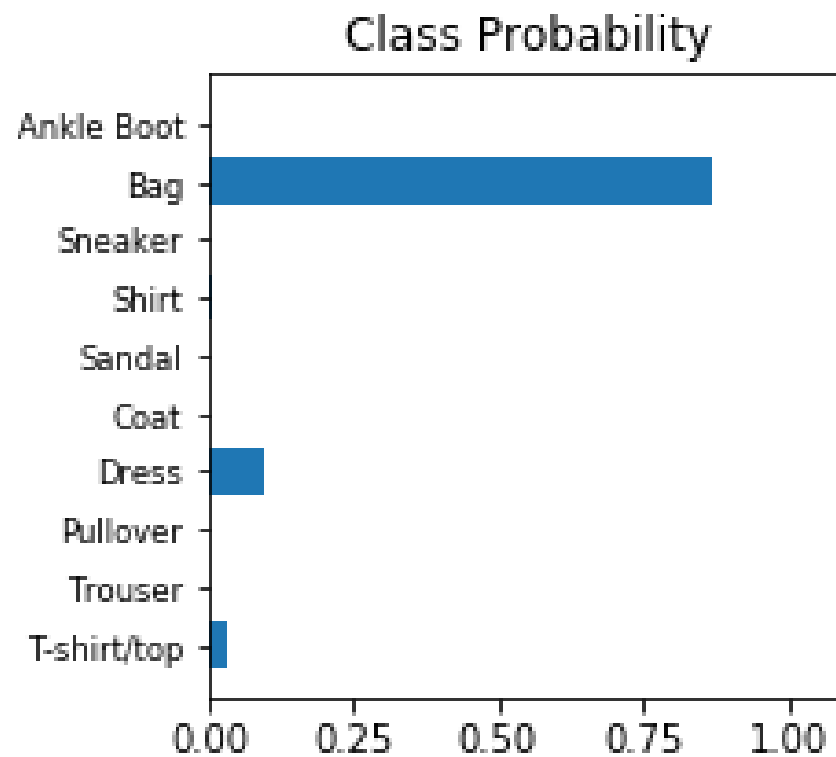
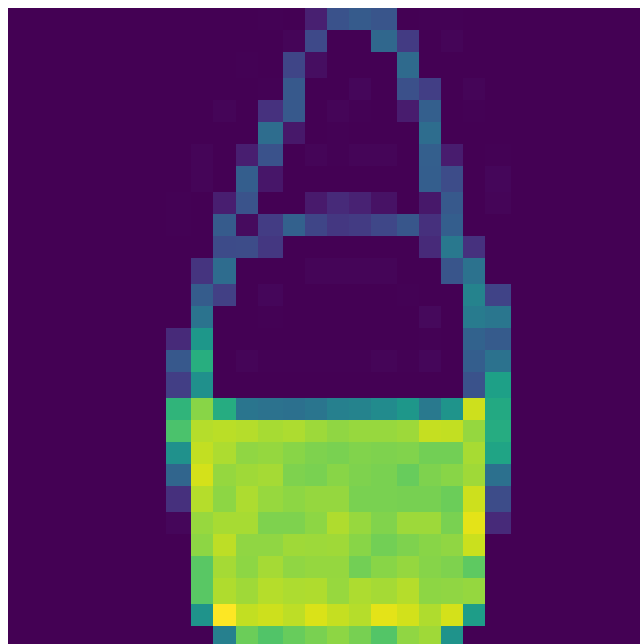
**Answer 1:** Outputting probabilities results in differentiable functions

**Question:** Why do we output probabilities instead of a binary values

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

**Answer 1:** Outputting probabilities results in differentiable functions

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



# Agenda

1. Review
2. Torch optimization coding
3. Classification task
4. Probability review
5. **Define model  $f$**
6. Define loss function  $\mathcal{L}$
7. Find  $\theta$  that minimize  $\mathcal{L}$
8. Coding

# Agenda

1. Review
2. Torch optimization coding
3. Classification task
4. Probability review
5. Define model  $f$
6. **Define loss function  $\mathcal{L}$**
7. Find  $\theta$  that minimize  $\mathcal{L}$
8. Coding



We consider the label  $\mathbf{y}_{[i]}$  as a conditional distribution

$$P(\mathbf{y}_{[i]} \mid \mathbf{x}_{[i]}) = \begin{bmatrix} P(\text{Shirt} \mid \text{img}) \\ P(\text{Bag} \mid \text{img}) \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

We consider the label  $\mathbf{y}_{[i]}$  as a conditional distribution

$$P(\mathbf{y}_{[i]} \mid \mathbf{x}_{[i]}) = \begin{bmatrix} P(\text{Shirt} \mid \text{img}) \\ P(\text{Bag} \mid \text{img}) \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

Our neural network also outputs some distribution

$$f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) = \begin{bmatrix} P(\text{Shirt} \mid \text{img}) \\ P(\text{Bag} \mid \text{img}) \end{bmatrix} = \begin{bmatrix} 0.6 \\ 0.4 \end{bmatrix}$$

We consider the label  $\mathbf{y}_{[i]}$  as a conditional distribution

$$P(\mathbf{y}_{[i]} \mid \mathbf{x}_{[i]}) = \begin{bmatrix} P(\text{Shirt} \mid \text{Image}) \\ P(\text{Bag} \mid \text{Image}) \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

Our neural network also outputs some distribution

$$f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) = \begin{bmatrix} P(\text{Shirt} \mid \text{Image}) \\ P(\text{Bag} \mid \text{Image}) \end{bmatrix} = \begin{bmatrix} 0.6 \\ 0.4 \end{bmatrix}$$

What loss function should we use for classification?

$$f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) = \begin{bmatrix} 0.6 \\ 0.4 \end{bmatrix}; \quad \mathbf{y}_{[i]} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) = \begin{bmatrix} 0.6 \\ 0.4 \end{bmatrix}; \quad \mathbf{y}_{[i]} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

We could use the square error like linear regression

$$f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) = \begin{bmatrix} 0.6 \\ 0.4 \end{bmatrix}; \quad \mathbf{y}_{[i]} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

We could use the square error like linear regression

$$(0.6 - 1)^2 + (0.4 - 0)^2$$

$$f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) = \begin{bmatrix} 0.6 \\ 0.4 \end{bmatrix}; \quad \mathbf{y}_{[i]} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

We could use the square error like linear regression

$$(0.6 - 1)^2 + (0.4 - 0)^2$$

This can work, but in reality it does not work well

$$f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) = \begin{bmatrix} 0.6 \\ 0.4 \end{bmatrix}; \quad \mathbf{y}_{[i]} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

We could use the square error like linear regression

$$(0.6 - 1)^2 + (0.4 - 0)^2$$

This can work, but in reality it does not work well

Instead, we use the **cross-entropy loss**



$$f(\mathbf{x}_{[i]}, \boldsymbol{\theta}) = \begin{bmatrix} 0.6 \\ 0.4 \end{bmatrix}; \quad \mathbf{y}_{[i]} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

We could use the square error like linear regression

$$(0.6 - 1)^2 + (0.4 - 0)^2$$

This can work, but in reality it does not work well

Instead, we use the **cross-entropy loss**

Let us derive it

We can model  $f(\boldsymbol{x}, \boldsymbol{\theta})$  and  $\boldsymbol{y}$  as probability distributions

We can model  $f(\mathbf{x}, \boldsymbol{\theta})$  and  $\mathbf{y}$  as probability distributions

How do we measure the difference between probability distributions?

We can model  $f(x, \theta)$  and  $y$  as probability distributions

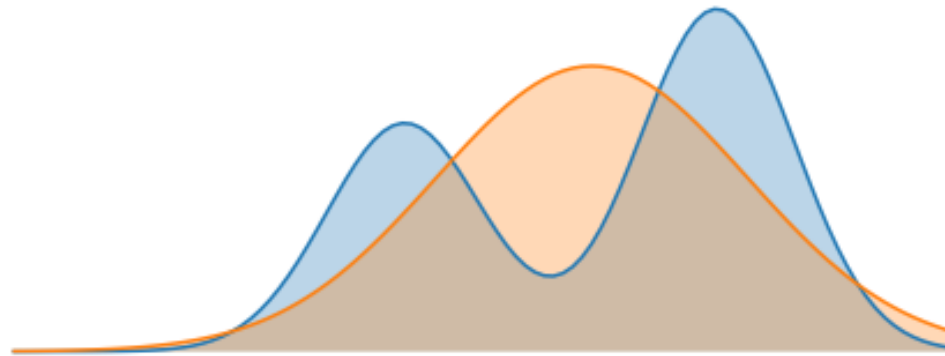
How do we measure the difference between probability distributions?

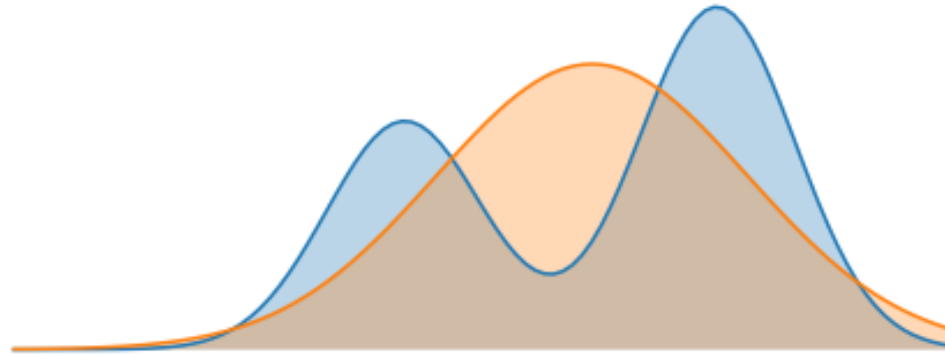
We use the **Kullback-Leibler Divergence (KL)**

We can model  $f(x, \theta)$  and  $y$  as probability distributions

How do we measure the difference between probability distributions?

We use the **Kullback-Leibler Divergence (KL)**





$$\text{KL}(P, Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)}$$

First, write down KL-divergence

$$\text{KL}(P, Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)}$$

First, write down KL-divergence

$$\text{KL}(P, Q) = \sum_i P(i) \log \frac{P(i)}{Q(i)}$$

Plug in our two distributions  $P = f$  and  $Q = y$

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



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

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

Rewrite the logarithm using the sum rule of logarithms

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

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

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

Split the sum into two parts

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

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

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

The first term is constant, and we will minimize the loss. So  $\arg \min_{\boldsymbol{\theta}} \mathcal{L} + k = \arg \min_{\boldsymbol{\theta}} \mathcal{L}$ . Therefore, we can ignore the first term.

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

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

The first term is constant, and we will minimize the loss. So  $\arg \min_{\boldsymbol{\theta}} \mathcal{L} + k = \arg \min_{\boldsymbol{\theta}} \mathcal{L}$ . Therefore, we can ignore the first term.

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

This is the loss for a classification task! We call this the **cross-entropy** loss function

# Example:



## Example:

$$\mathcal{L}(\boldsymbol{x}, \boldsymbol{y}, \boldsymbol{\theta}) = - \sum_{i=1}^{d_y} P(y_i \mid \boldsymbol{x}) \log f(\boldsymbol{x}, \boldsymbol{\theta})_i; \quad f(\boldsymbol{x}, \boldsymbol{\theta}) = \begin{bmatrix} 0.6 \\ 0.4 \end{bmatrix}; \quad \boldsymbol{y} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

## Example:

$$\begin{aligned}\mathcal{L}(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) &= - \sum_{i=1}^{d_y} P(y_i \mid \mathbf{x}) \log f(\mathbf{x}, \boldsymbol{\theta})_i; \quad f(\mathbf{x}, \boldsymbol{\theta}) = \begin{bmatrix} 0.6 \\ 0.4 \end{bmatrix}; \quad \mathbf{y} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ &= - \left[ P(y_1 \mid \mathbf{x}) \log f(\mathbf{x}, \boldsymbol{\theta})_1 + P(y_2 \mid \mathbf{x}) \log f(\mathbf{x}, \boldsymbol{\theta})_2 \right]\end{aligned}$$

## Example:

$$\begin{aligned}\mathcal{L}(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) &= - \sum_{i=1}^{d_y} P(y_i \mid \mathbf{x}) \log f(\mathbf{x}, \boldsymbol{\theta})_i; \quad f(\mathbf{x}, \boldsymbol{\theta}) = \begin{bmatrix} 0.6 \\ 0.4 \end{bmatrix}; \quad \mathbf{y} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ &= - \left[ P(y_1 \mid \mathbf{x}) \log f(\mathbf{x}, \boldsymbol{\theta})_1 + P(y_2 \mid \mathbf{x}) \log f(\mathbf{x}, \boldsymbol{\theta})_2 \right] \\ &= - [1 \cdot \log 0.6 + 0 \cdot \log 0.4]\end{aligned}$$

## Example:

$$\begin{aligned}\mathcal{L}(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) &= - \sum_{i=1}^{d_y} P(y_i \mid \mathbf{x}) \log f(\mathbf{x}, \boldsymbol{\theta})_i; \quad f(\mathbf{x}, \boldsymbol{\theta}) = \begin{bmatrix} 0.6 \\ 0.4 \end{bmatrix}; \quad \mathbf{y} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ &= - \left[ P(y_1 \mid \mathbf{x}) \log f(\mathbf{x}, \boldsymbol{\theta})_1 + P(y_2 \mid \mathbf{x}) \log f(\mathbf{x}, \boldsymbol{\theta})_2 \right] \\ &= - [1 \cdot \log 0.6 + 0 \cdot \log 0.4] \\ &= - [-0.22]\end{aligned}$$

## Example:

$$\begin{aligned}\mathcal{L}(\mathbf{x}, \mathbf{y}, \boldsymbol{\theta}) &= - \sum_{i=1}^{d_y} P(y_i \mid \mathbf{x}) \log f(\mathbf{x}, \boldsymbol{\theta})_i; \quad f(\mathbf{x}, \boldsymbol{\theta}) = \begin{bmatrix} 0.6 \\ 0.4 \end{bmatrix}; \quad \mathbf{y} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \\ &= - \left[ P(y_1 \mid \mathbf{x}) \log f(\mathbf{x}, \boldsymbol{\theta})_1 + P(y_2 \mid \mathbf{x}) \log f(\mathbf{x}, \boldsymbol{\theta})_2 \right] \\ &= - [1 \cdot \log 0.6 + 0 \cdot \log 0.4] \\ &= - [-0.22] \\ &= 0.22\end{aligned}$$

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

By minimizing the loss, we make  $f(\mathbf{x}, \boldsymbol{\theta}) = P(\mathbf{y} \mid \mathbf{x})$

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

By minimizing the loss, we make  $f(\mathbf{x}, \boldsymbol{\theta}) = P(\mathbf{y} \mid \mathbf{x})$

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

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

By minimizing the loss, we make  $f(\mathbf{x}, \boldsymbol{\theta}) = P(\mathbf{y} \mid \mathbf{x})$

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

$$f(\mathbf{x}, \boldsymbol{\theta}) = P(\mathbf{y} \mid \mathbf{x}) = P \left( \begin{bmatrix} \text{boot} \\ \text{dress} \\ \vdots \end{bmatrix} \mid \begin{bmatrix} \text{img} \end{bmatrix} \right)$$



Our loss was just for a single image

Our loss was just for a single image

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

Our loss was just for a single image

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

Find  $\boldsymbol{\theta}$  that minimize the loss over the whole dataset

Our loss was just for a single image

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

Find  $\boldsymbol{\theta}$  that minimize the loss over the whole dataset

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

# Agenda

1. Review
2. Torch optimization coding
3. Classification task
4. Probability review
5. Define model  $f$
6. **Define loss function  $\mathcal{L}$**
7. Find  $\theta$  that minimize  $\mathcal{L}$
8. Coding

# Agenda

1. Review
2. Torch optimization coding
3. Classification task
4. Probability review
5. Define model  $f$
6. Define loss function  $\mathcal{L}$
7. **Find  $\theta$  that minimize  $\mathcal{L}$**
8. Coding

# Classification

Find  $\theta$  just like before, using gradient descent

# Classification

Find  $\theta$  just like before, using gradient descent

The gradients are the same as before except the last layer



# Classification

Find  $\theta$  just like before, using gradient descent

The gradients are the same as before except the last layer

I will not derive any more gradients, but the softmax gradient nearly identical to the sigmoid function gradient

# Classification

Find  $\theta$  just like before, using gradient descent

The gradients are the same as before except the last layer

I will not derive any more gradients, but the softmax gradient nearly identical to the sigmoid function gradient

$$\nabla_{\theta} \text{softmax}(\mathbf{z}) = \text{softmax}(\mathbf{z}) \odot (1 - \text{softmax}(\mathbf{z}))$$

# Classification

Find  $\theta$  just like before, using gradient descent

The gradients are the same as before except the last layer

I will not derive any more gradients, but the softmax gradient nearly identical to the sigmoid function gradient

$$\nabla_{\theta} \text{softmax}(\mathbf{z}) = \text{softmax}(\mathbf{z}) \odot (1 - \text{softmax}(\mathbf{z}))$$

This is because softmax is a multi-class generalization of the sigmoid function

# Agenda

1. Review
2. Torch optimization coding
3. Classification task
4. Probability review
5. Define model  $f$
6. Define loss function  $\mathcal{L}$
7. **Find  $\theta$  that minimize  $\mathcal{L}$**
8. Coding

Relax

# Agenda

1. Review
2. Torch optimization coding
3. Classification task
4. Probability review
5. Define model  $f$
6. Define loss function  $\mathcal{L}$
7. Find  $\theta$  that minimize  $\mathcal{L}$
8. **Coding**

# Coding

You have everything you need to solve deep learning tasks!

## 1. Regression

# Coding

You have everything you need to solve deep learning tasks!

1. Regression
2. Classification



# Coding

You have everything you need to solve deep learning tasks!

1. Regression
2. Classification

Every interesting task (chatbot, self driving car, etc):

# Coding

You have everything you need to solve deep learning tasks!

1. Regression
2. Classification

Every interesting task (chatbot, self driving car, etc):

1. Construct a deep neural network

# Coding

You have everything you need to solve deep learning tasks!

1. Regression
2. Classification

Every interesting task (chatbot, self driving car, etc):

1. Construct a deep neural network
2. Compute regression or classification loss

# Coding

You have everything you need to solve deep learning tasks!

1. Regression
2. Classification

Every interesting task (chatbot, self driving car, etc):

1. Construct a deep neural network
2. Compute regression or classification loss
3. Optimize with gradient descent

# Coding

You have everything you need to solve deep learning tasks!

1. Regression
2. Classification

Every interesting task (chatbot, self driving car, etc):

1. Construct a deep neural network
2. Compute regression or classification loss
3. Optimize with gradient descent

The rest of this course will examine neural network architectures

# Coding

<https://colab.research.google.com/drive/1BGMIE2CjlLJOH-D2r9AariPDVgxjWlqG?usp=sharing>

# Admin

Homework 3 is released, you have two weeks to complete it

# Admin

Homework 3 is released, you have two weeks to complete it

<https://colab.research.google.com/drive/1LainS20p6c3YVRFM4XgHRBODh6dvAaT2?usp=sharing#scrollTo=q8pJST5xFt-p>