

# NEURAL NETWORKS

MACHINE LEARNING 1 UE (INP.33761UF)

JOIN HERE: [fbr.io/ml1p7](https://fbr.io/ml1p7)

---

Thomas Wedenig

May 08, 2024

Institute of Theoretical Computer Science  
Graz University of Technology, Austria

Hello everyone and welcome to Machine Learning 1! Today, we're diving deep into the fascinating world of Neural Networks. In this session, we'll explore the fundamental concepts behind neural networks, their architecture, and how they enable machines to learn and make decisions akin to human brains.

Hello everyone and welcome to Machine Learning 1! Today, we're diving deep into the fascinating world of Neural Networks. In this session, we'll explore the fundamental concepts behind neural networks, their architecture, and how they enable machines to learn and make decisions akin to human brains.

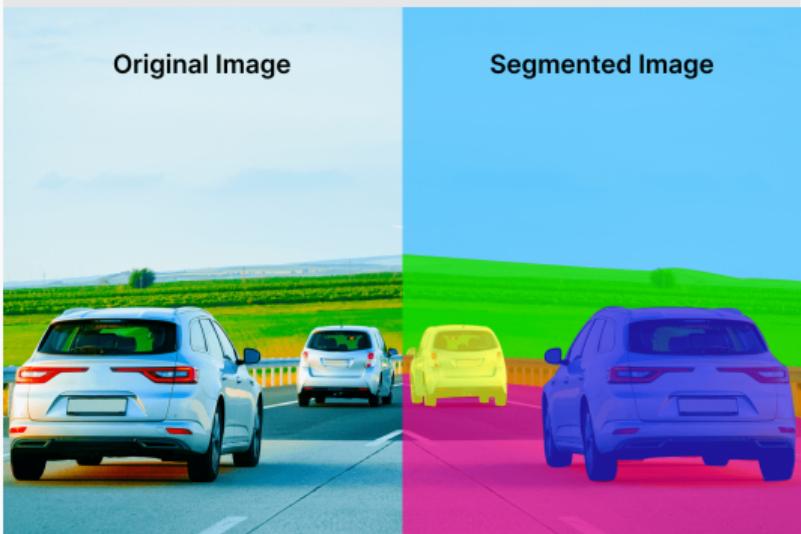
Fun fact: This text was written by a Large Language Model and was read by a voice cloning model that used a two-minute recording of my voice.

1. Motivation
2. Feed-Forward Neural Networks
3. Connection to previous methods
4. Loss Functions and Training

## MOTIVATION

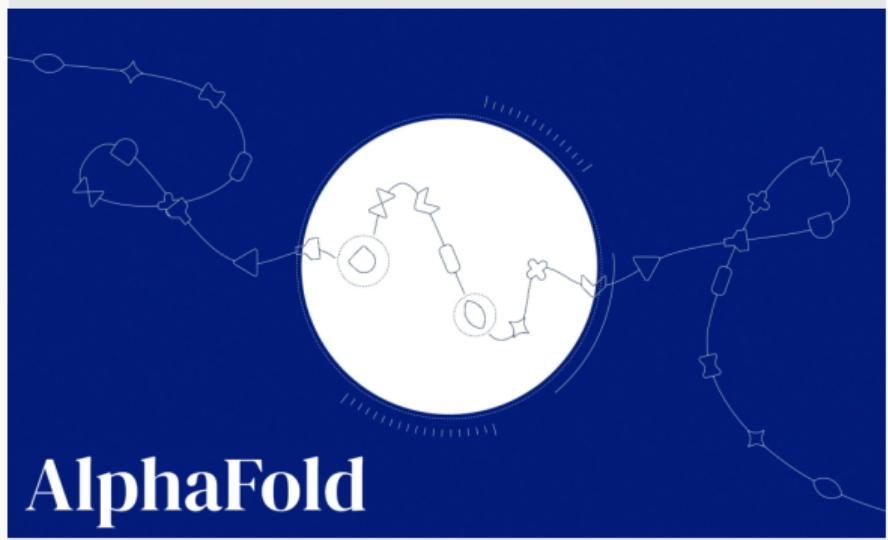
---

## Computer Vision



[https://deeplobe.ai/  
image-segmentation-the-most-interesting-applications/](https://deeplobe.ai/image-segmentation-the-most-interesting-applications/)

## AI 4 Science



[https://techcrunch.com/2021/07/22/  
deepmind-puts-the-entire-human-proteome-online-as-folded-by-alphafold](https://techcrunch.com/2021/07/22/deepmind-puts-the-entire-human-proteome-online-as-folded-by-alphafold)

# GENERATIVE AI - IMAGES

DALL·E 3



OpenAI DALL·E 3, "The inside of a neural network."

thispersondoesnotexist.com



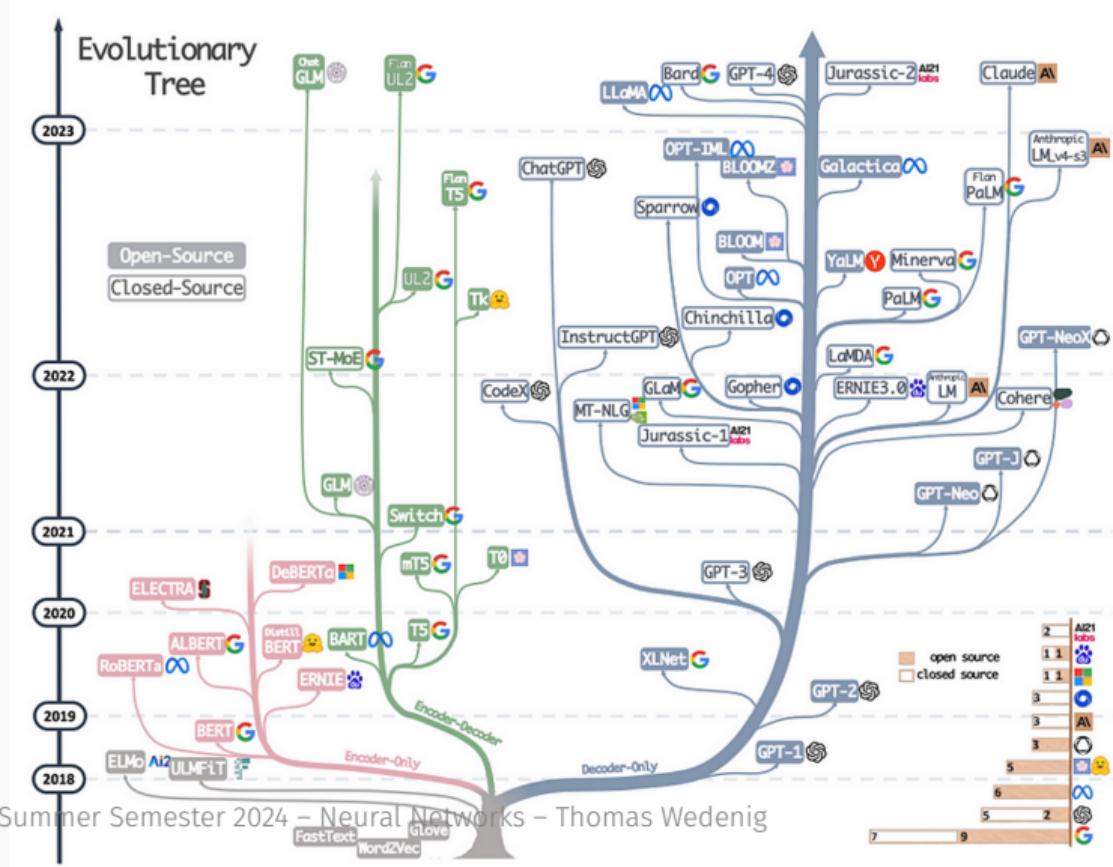
<https://thispersondoesnotexist.com>



e.g., <https://openai.com/index/sora>

# GENERATIVE AI - LARGE LANGUAGE MODELS

Source: [https://abiaryan.com/posts/intro\\_llms/](https://abiaryan.com/posts/intro_llms/)



- Neural Networks are **powerful, flexible, non-linear function approximators** !
- Most of these impressive applications utilize **large (and deep) neural networks**
- In this course, we will work with small and simple networks

- Neural Networks are **powerful, flexible, non-linear function approximators !**
- Most of these impressive applications utilize **large (and deep) neural networks**
- In this course, we will work with small and simple networks

## NN architectures

- Feed-forward NNs
- Recurrent NNs
- Convolutional NNs
- Transformers

Covered in **Deep Learning VO/KU**  
(Master's level)

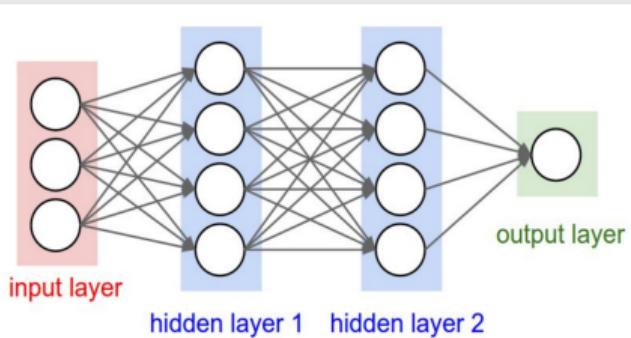
- Neural Networks are **powerful, flexible, non-linear function approximators !**
- Most of these impressive applications utilize **large (and deep) neural networks**
- In this course, we will work with small and simple networks

## NN architectures

- Feed-forward NNs
- Recurrent NNs
- Convolutional NNs
- Transformers

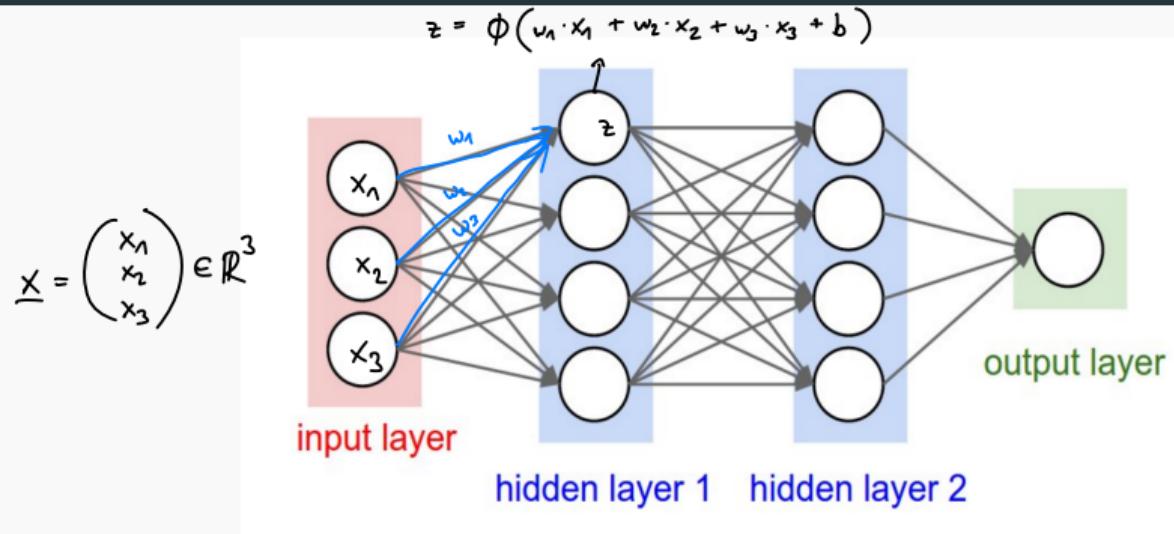
Covered in **Deep Learning VO/KU**  
(Master's level)

## Feed-forward Neural Network



## FEED-FORWARD NEURAL NETWORKS

---



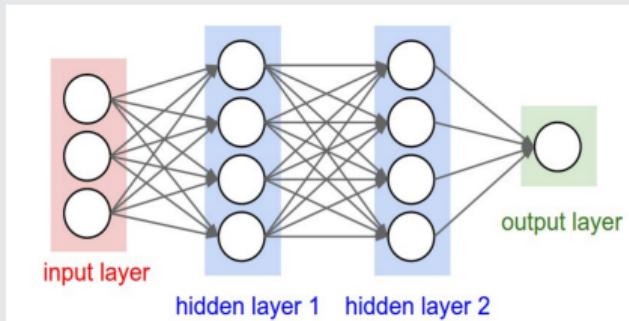
- Let  $f_{\theta}(x)$  denote the network output
- Popular activation function for hidden units  $\phi_h(z) = \text{ReLU}(z) = \max(0, z)$
- What about the activation at the output  $\phi_o$  ?

Supervised Learning: Dataset  $\mathcal{D} = \{(x^{(1)}, y^{(1)}), \dots, (x^{(M)}, y^{(M)})\}$  with  $x^{(i)} \in \mathbb{R}^N$

Supervised Learning: Dataset  $\mathcal{D} = \{(x^{(1)}, y^{(1)}), \dots, (x^{(M)}, y^{(M)})\}$  with  $x^{(i)} \in \mathbb{R}^N$

## Regression

$$y \in \mathbb{R} \implies f_{\theta}(x) \in \mathbb{R}$$



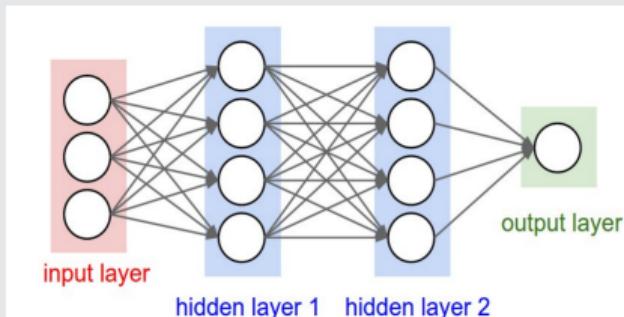
- No activation at output
  - “linear” activation:  $\phi_o(z) = z$

# REGRESSION VS. BINARY CLASSIFICATION

Supervised Learning: Dataset  $\mathcal{D} = \{(x^{(1)}, y^{(1)}), \dots, (x^{(M)}, y^{(M)})\}$  with  $x^{(i)} \in \mathbb{R}^N$

## Regression

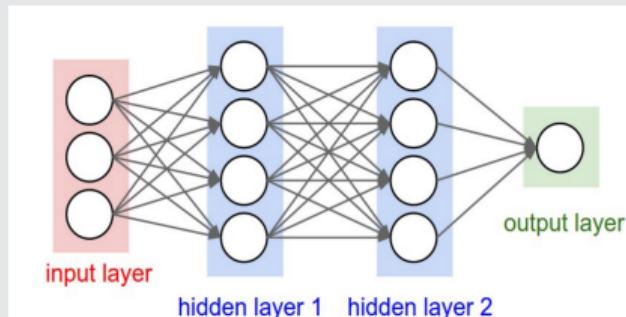
$$y \in \mathbb{R} \implies f_{\theta}(x) \in \mathbb{R}$$



- No activation at output
  - “linear” activation:  $\phi_o(z) = z$

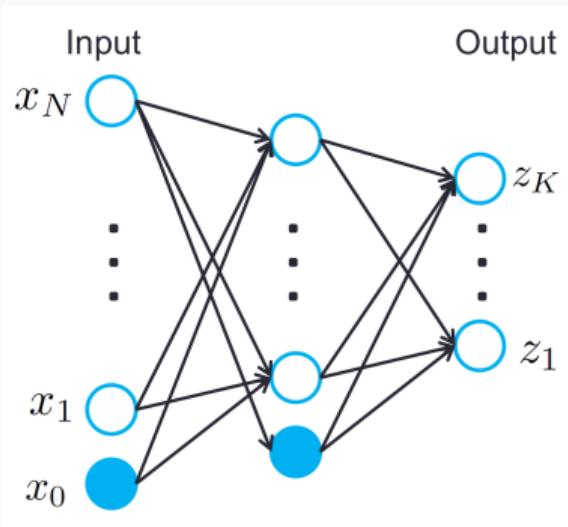
## Binary Classification

$$y \in \{-1, 1\} \implies f_{\theta}(x) \in [0, 1]$$



- Sigmoid activation at output
  - $\phi_o(z) = \sigma(z) = 1/(1 + \exp(-z))$

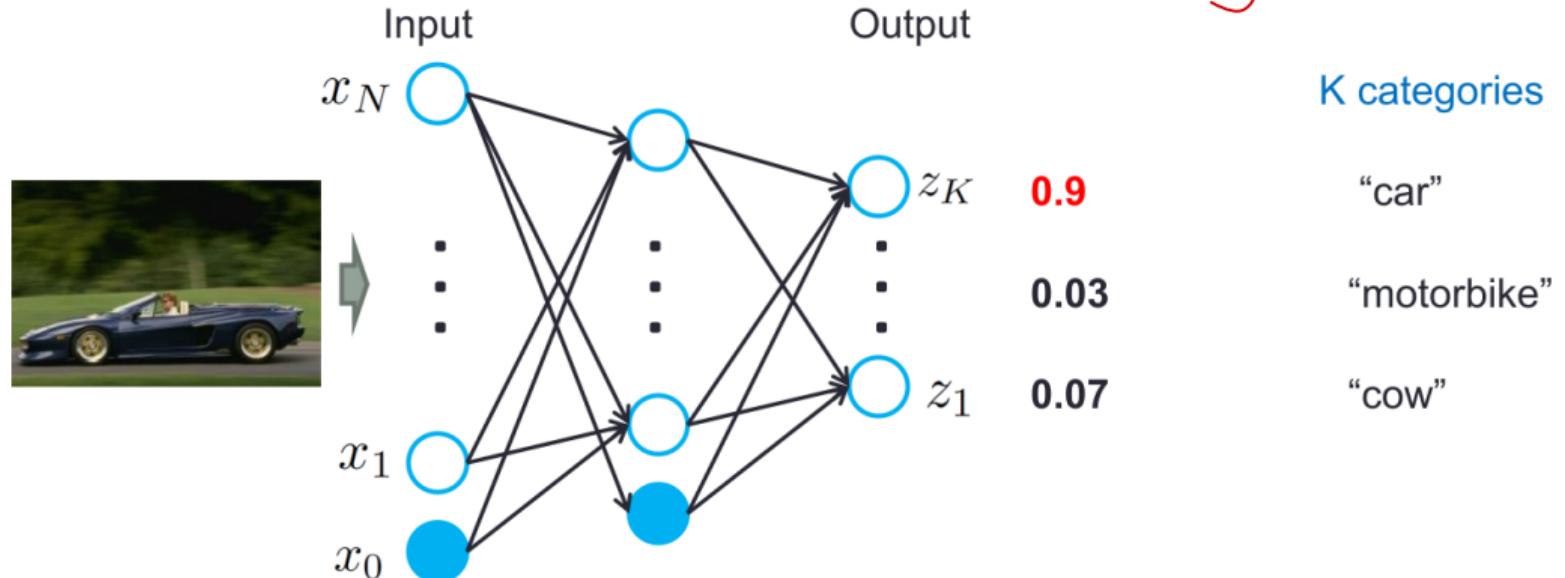
- What if  $y \in \{1, \dots, K\}$  ?
  - K-class classification problem 😊



- NN should output a **probability distribution** over  $K$  classes
- Let  $\mathbf{z} = (z_1, \dots, z_K)^T$  be the output **before applying the output activation**
  - Called “pre-activation” or “logits”
- We use  $\phi_o(\mathbf{z}) = \text{softmax}(\mathbf{z})$  with

$$\text{softmax}(\mathbf{z})_k = \frac{\exp(z_k)}{\sum_{j=1}^K \exp(z_j)}$$

$$-\log(0.9)$$

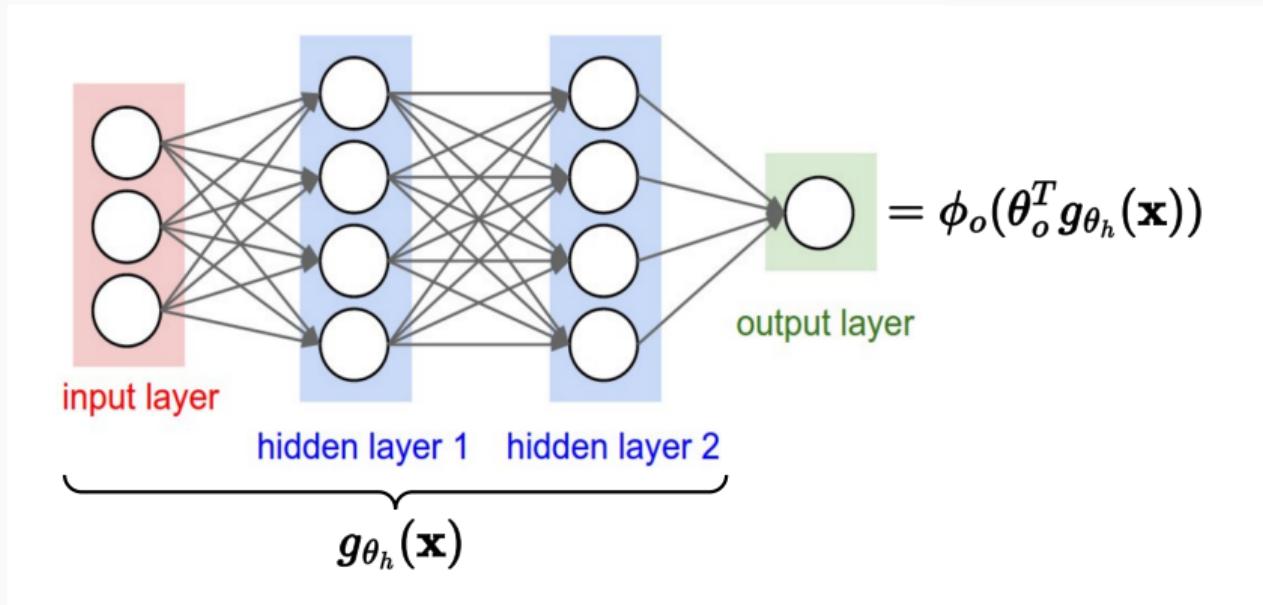


## CONNECTION TO PREVIOUS METHODS

---

Previously, our prediction function  $f_{\theta}(x)$  looked like this:

- **Linear Regression:**  $f_{\theta}(x) = \theta^T x$ 
  - A linear function in  $x$  and  $\theta$
- **Linear Regression with non-linear feature transform:**  $f_{\theta}(x) = \theta^T g(x)$ 
  - A linear function in  $\theta$ , but not in  $x$
- **Logistic Regression:**  $f_{\theta}(x) = \sigma(\theta^T g(x))$ 
  - As above, pushed through a sigmoid non-linearity
- **Note:** Feature transform  $g$  has no parameters that are learned



- Non-linear feature transform  $g$  now has learnable parameters  $\theta_h$  !

## LOSS FUNCTIONS AND TRAINING

---

# LOSS FUNCTIONS

- We use the same loss functions as always:

$$\mathcal{L}(\theta) = \frac{1}{M} \sum_{i=1}^M \ell(y^{(i)}, \underbrace{f_\theta(x^{(i)})}_{\hat{y}^{(i)}})$$

where  $\ell(y, \hat{y})$  is a per-example loss

- Regression uses Mean Squared Error:

$$\ell(y, \hat{y}) = (y - \hat{y})^2$$

- Binary Classification uses Binary Cross-Entropy:

$$\ell(y, \hat{y}) = \begin{cases} 0 & \text{if } y = 1 \\ 1 & \text{if } y = -1 \end{cases} (-\log(\hat{y})) + \begin{cases} 1 & \text{if } y = 1 \\ 0 & \text{if } y = -1 \end{cases} (-\log(1 - \hat{y}))$$

- Multi-class Classification uses Multi-class Cross-Entropy:

$$\ell(y, \hat{y}) = -\log(\hat{y}_y), \quad \hat{y} \in \mathbb{R}^K, y \in \{1, \dots, K\}$$

- We want

$$\theta^* = \underset{\theta}{\operatorname{argmin}} \mathcal{L}(\theta)$$

- $\mathcal{L}(\theta)$  is highly non-convex in  $\theta$  😞
- But it's **differentiable** almost everywhere 😎
- Use **Gradient Descent** (or some variant) to find a local optimum of  $\mathcal{L}(\theta)$
- Gradients  $\nabla_{\theta} \mathcal{L}(\theta)$  can be computed via **backpropagation**
- However...

- Recall that

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{M} \sum_{i=1}^M \ell\left(y^{(i)}, f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})\right) \quad \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) = \frac{1}{M} \sum_{i=1}^M \nabla_{\boldsymbol{\theta}} \ell\left(y^{(i)}, f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})\right)$$

- What could be the **problem from a computational perspective** ?

- Recall that

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{M} \sum_{i=1}^M \ell\left(y^{(i)}, f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})\right) \quad \nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) = \frac{1}{M} \sum_{i=1}^M \nabla_{\boldsymbol{\theta}} \ell\left(y^{(i)}, f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})\right)$$

- What could be the **problem from a computational perspective** ?
- In many cases,  $M = |\mathcal{D}|$  is **very large**
- We have to loop over our **entire training set** to get a single gradient  $\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta})$ 
  - Which amounts to a **single step** in gradient descent 😞

### Solution: Mini-Batch Gradient Descent

For each step, use a **small subset** of our training data  $\mathcal{B} \subset \mathcal{D}$  to estimate the gradient

$$\nabla_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}) = \frac{1}{M} \sum_{i=1}^M \nabla_{\boldsymbol{\theta}} \ell \left( y^{(i)}, f_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \right) \approx \frac{1}{|\mathcal{B}|} \sum_{j=1}^{|\mathcal{B}|} \nabla_{\boldsymbol{\theta}} \ell \left( y^{(j)}, f_{\boldsymbol{\theta}}(\mathbf{x}^{(j)}) \right)$$

where  $(\mathbf{x}^{(j)}, y^{(j)}) \in \mathcal{B}$ .

- $\mathcal{B}$  is called a **mini-batch**
- This version of Gradient Descent is often called **Mini-Batch Gradient Descent**

# NEURAL NETWORK REGRESSION DEMO