

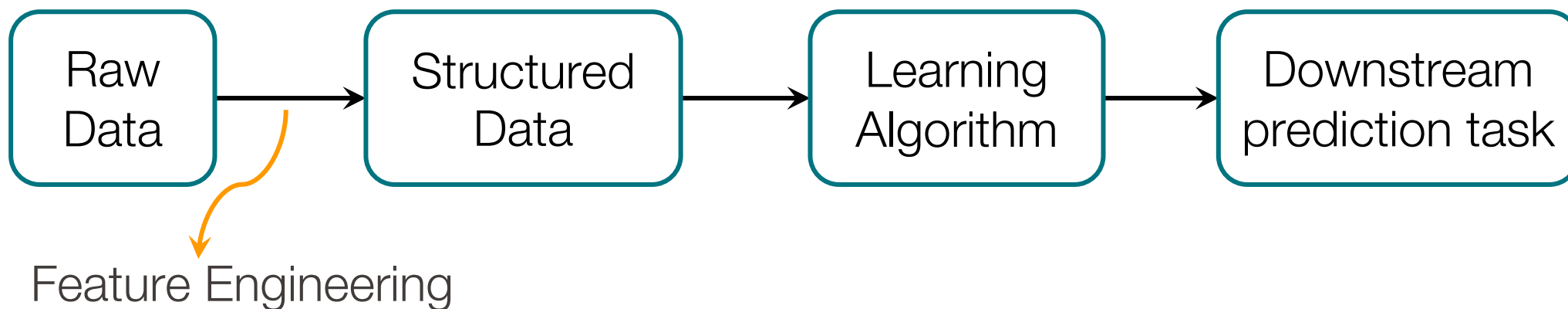
CS502: Deep Learning in Biomedicine

Deep Learning Basics

Maria Brbić
Fall 2023

Traditional Machine Learning

- Machine learning lifecycle



- Tedious manual effort on feature engineering
- Does not scale in practice and results in suboptimal features

Traditional Machine Learning: Example

DNA sequences

```
> F31.01
tattggacgggattgagtttaacgaacgccgttcctcaatgcgggtggttgcta
atgtag tgctaccgaacccaaagtccatcccgtagtcccttcaaagctttca
acgaggaactgg
```

```
> F31.02
aagttttcggcttcaaagggagataccaaatacggcggaataactttccgcag
ccgttttattgagtgtagcagaatgatctcccccatgttgataaggcttaac
tggttaccatctcttggggcacttgggtcctgctagacgggtacgtgcctcgcc
cgtagctaacgacttgcg agtatgcacacttgtccaataggcaacgcctggct
ttgtg
```

```
> F31.03
gccaaatcggtcatatcccgaggactggacaactaaatgtattaaagtctacag
tcagtgtccaccagccaccacagagctaccgcaccgtagattcattaggcgata
tatagacttgggctgagactgacacgggaccataactgattcgatctgtacctc
ctcccacgctattattaagattgtttcgtttgagatgctctacacggcgccatt
gcggtcataaactgaccaaccacaaacttgtcgcccatgcacgccagtaaggca
aactcctcggagatcggaagtcgaagct
```



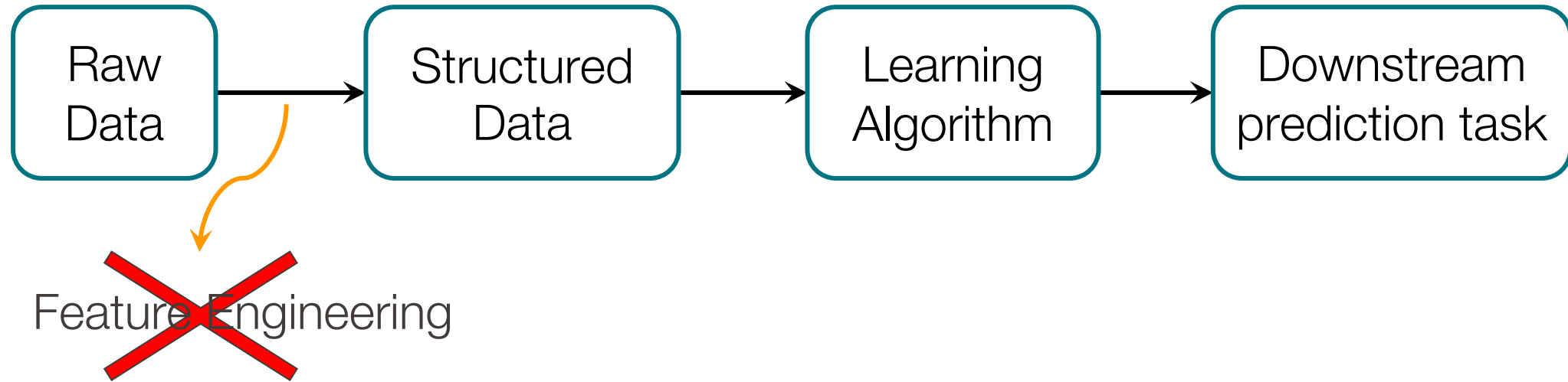
ID	AT	AC	AG	AA	TA	TC	TG	TT	...
F31.01									
F31.02									
F31.03									

What is the problem with such data representation?



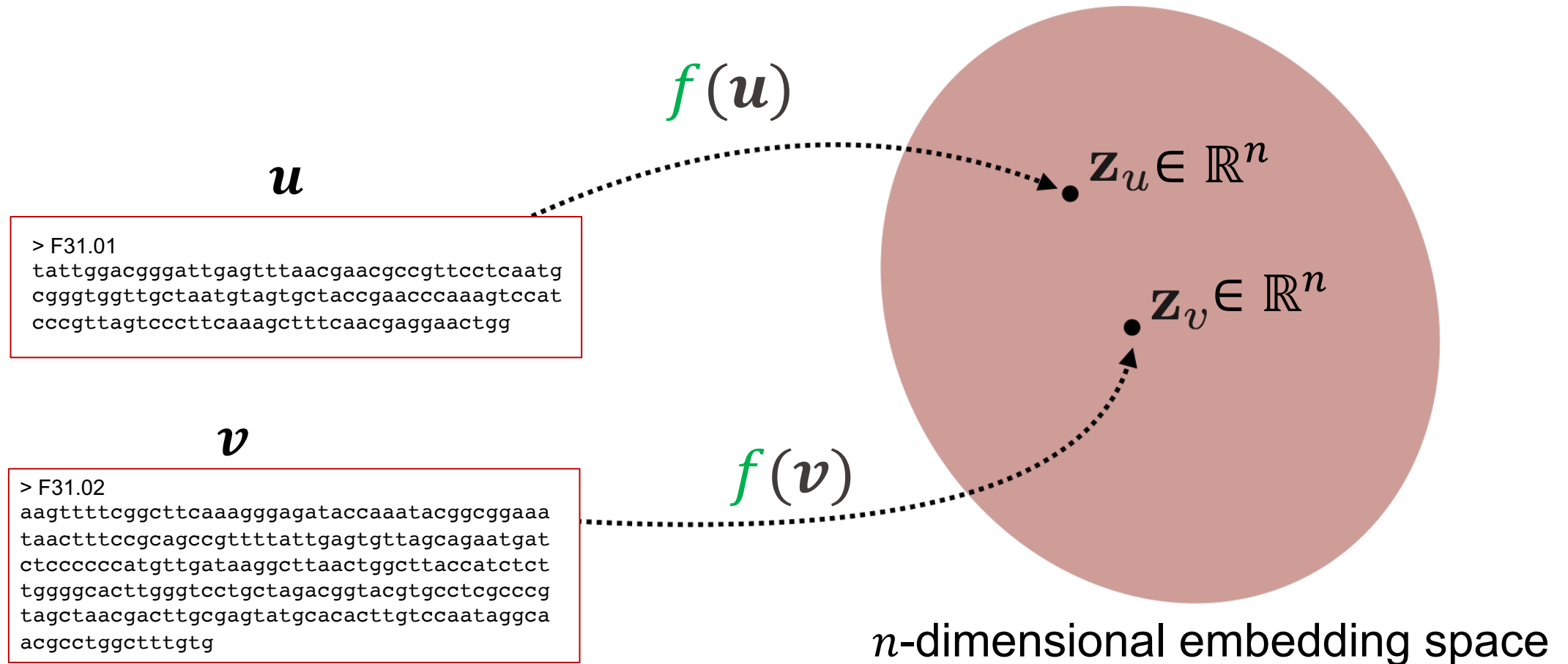
Deep Learning

- Machine learning lifecycle



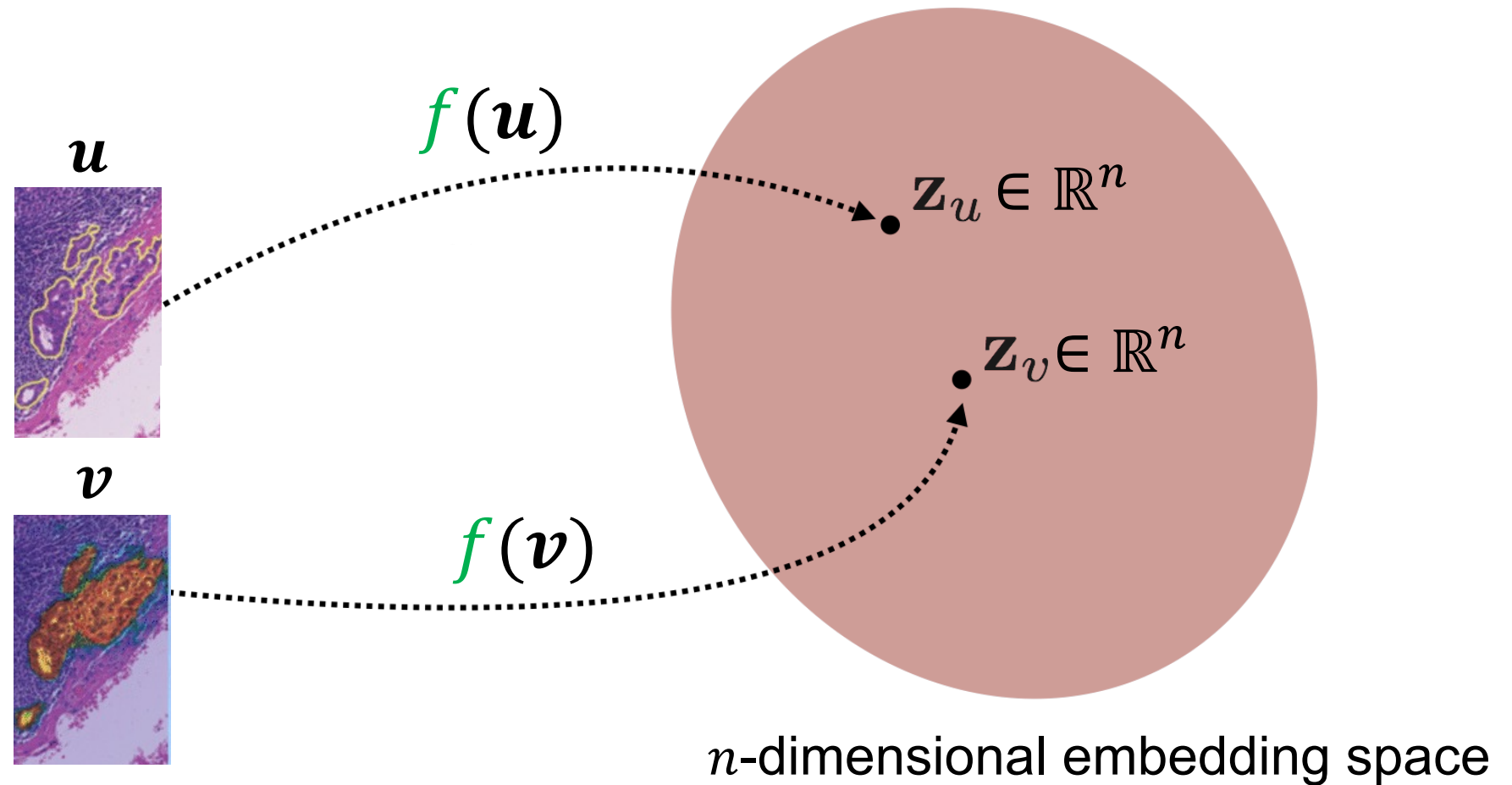
**Automatically
learn the features!**

Representation Learning



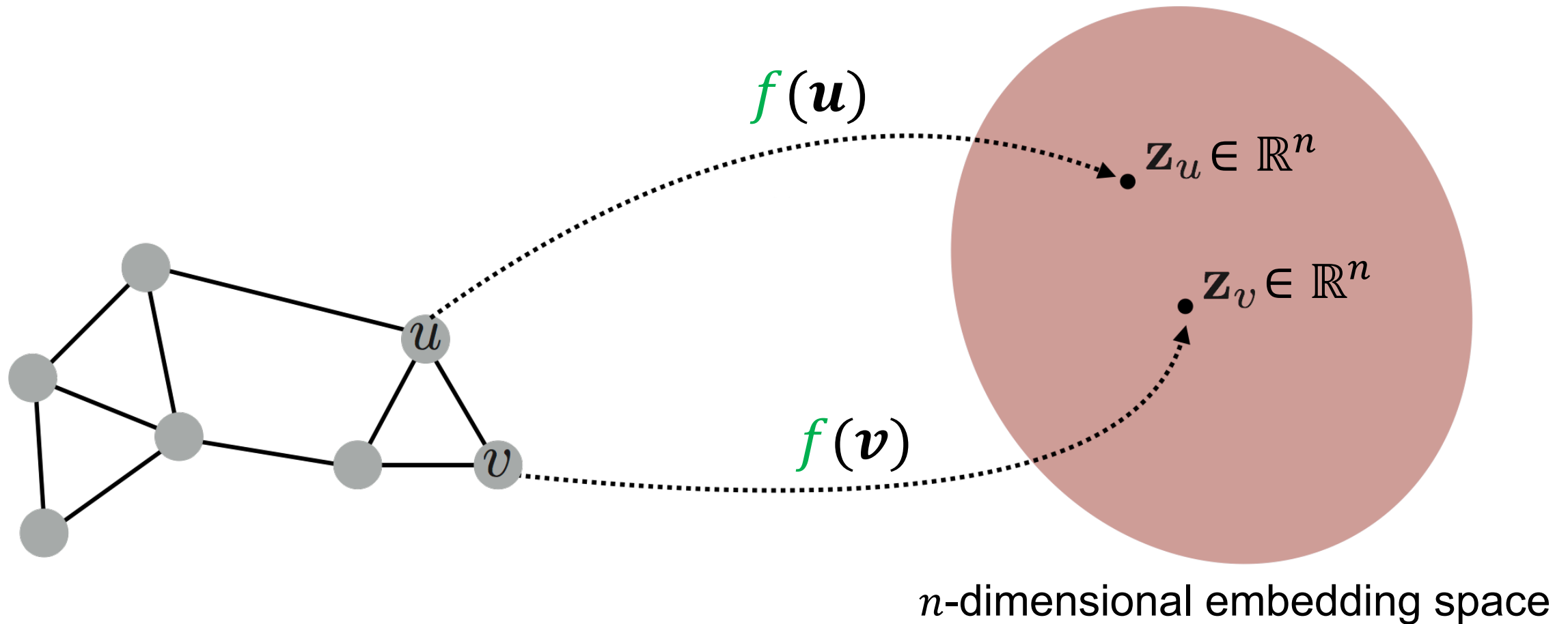
Key idea behind deep learning: Automatically learn underlying representations from the data!

Representation Learning



Key idea behind deep learning: Automatically learn underlying representations from the data!

Representation Learning

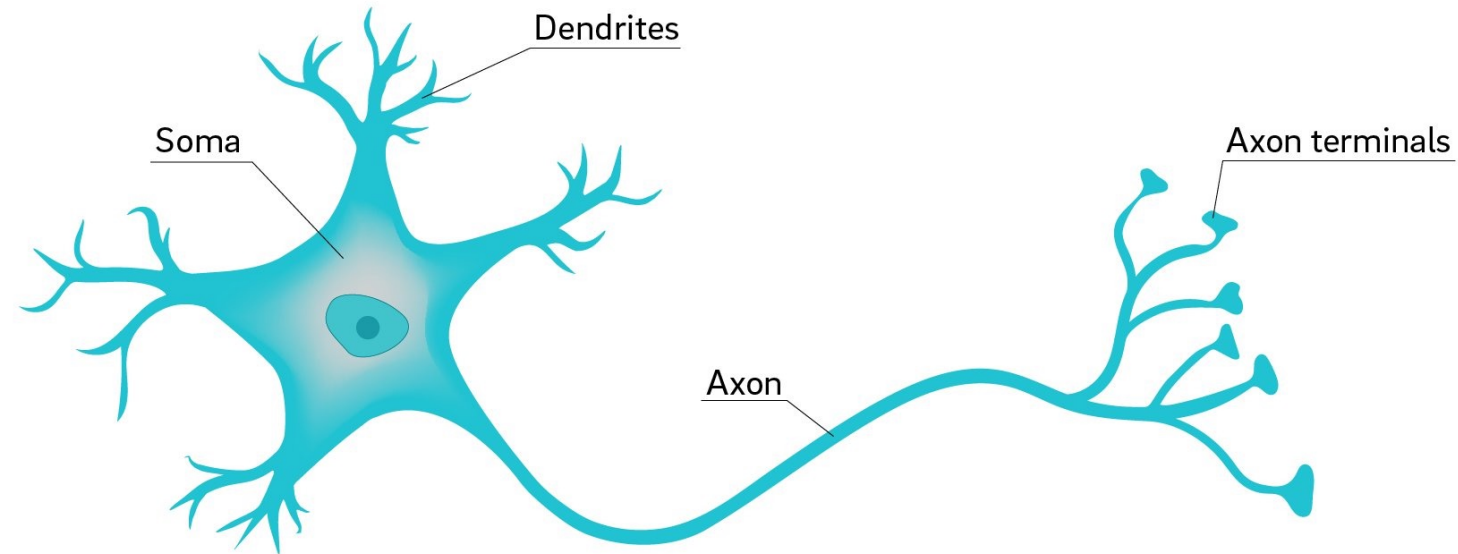


Key idea behind deep learning: Automatically learn underlying representations from the data!

Overview of neural networks

Inspiration

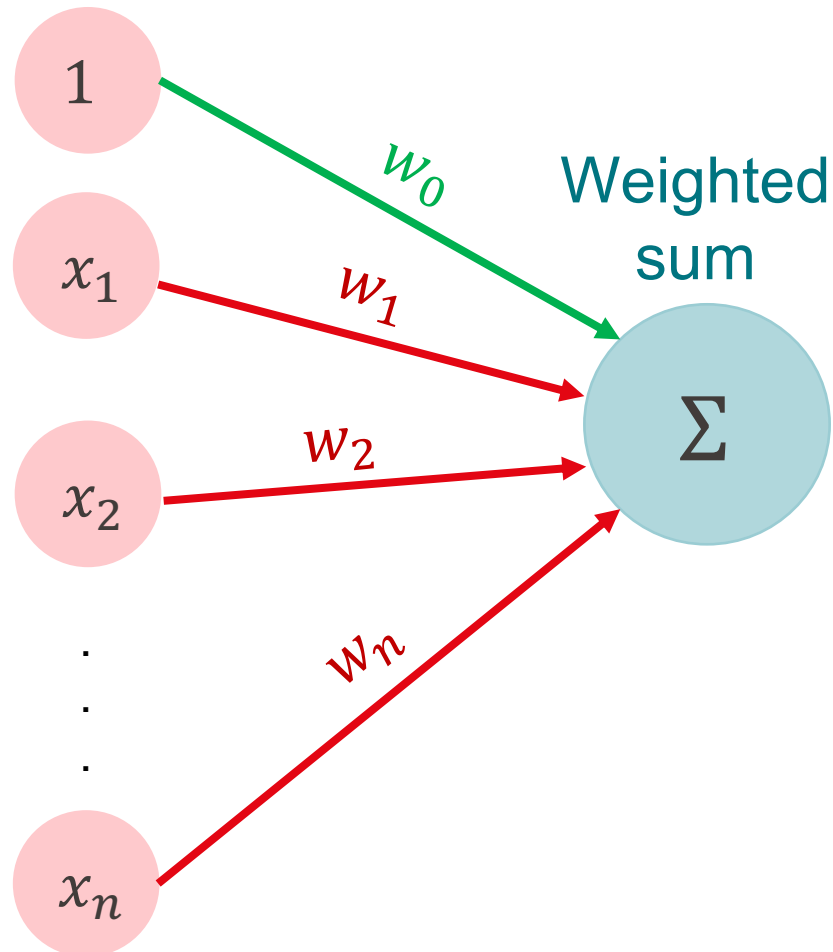
Neuron



How to transfer this idea to a computational model?

The Perceptron

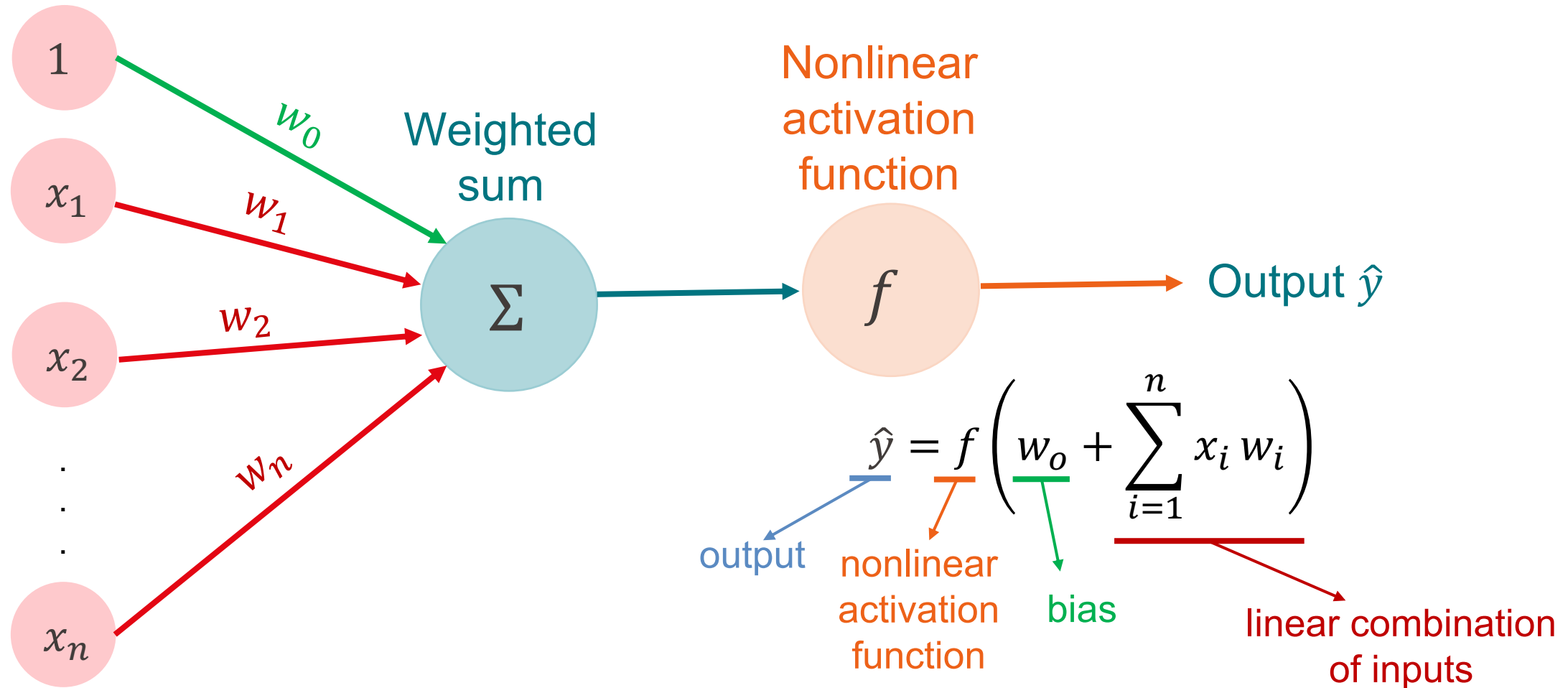
- The structural building block of deep learning



$$\underbrace{w_0}_{\text{bias}} + \underbrace{\sum_{i=1}^n x_i w_i}_{\text{linear combination of inputs}}$$

The Perceptron

- The structural building block of deep learning



The Perceptron: Forward Propagation

$$\hat{y} = f \left(\underbrace{w_o}_{\text{bias}} + \underbrace{\sum_{i=1}^n x_i w_i}_{\text{linear combination of inputs}} \right)$$

output nonlinear activation function

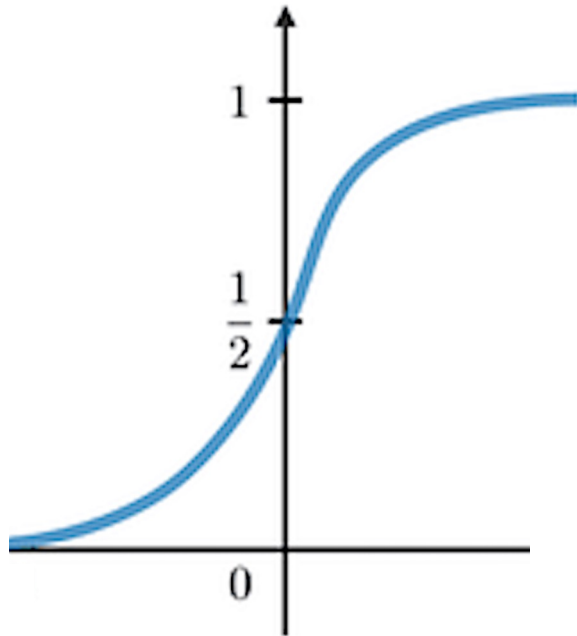
$$\hat{y} = f(w_o + \mathbf{x}^T \mathbf{w})$$

where $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{bmatrix}$ and $\mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \dots \\ w_n \end{bmatrix}$

What is activation function f ?

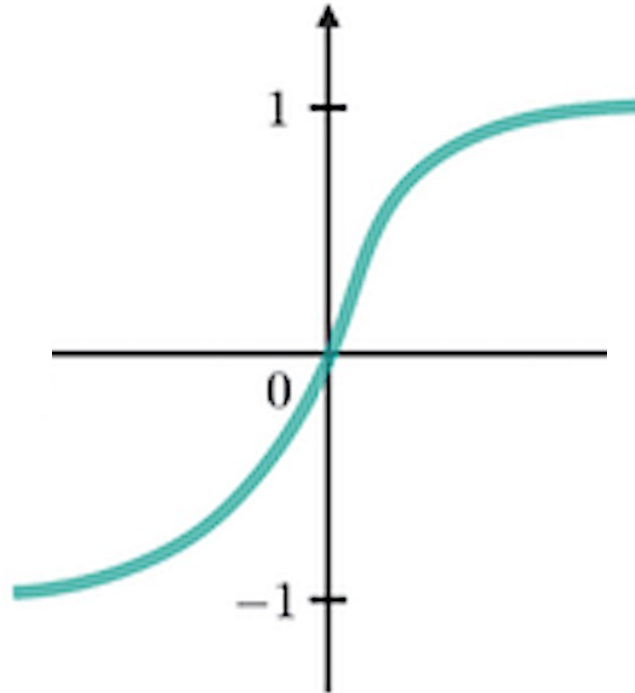
The Perceptron: Activation Functions Examples

Sigmoid



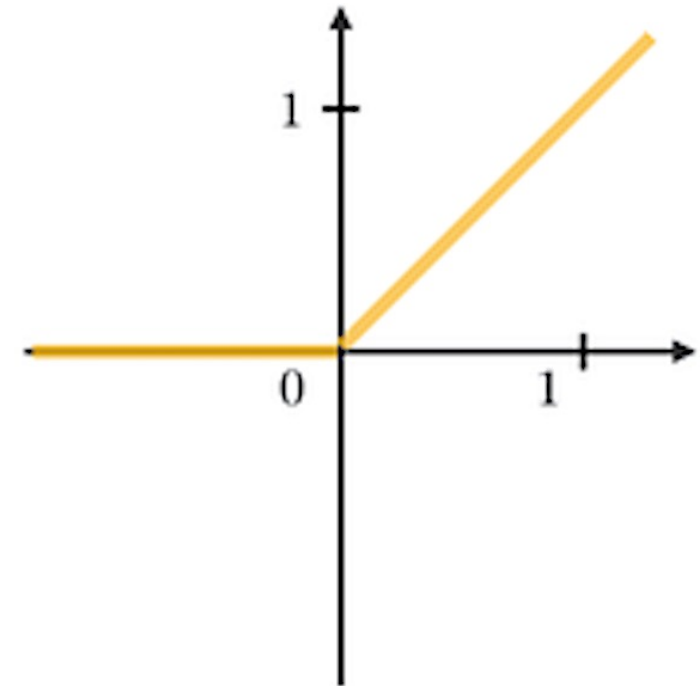
$$f(x) = \frac{1}{1 + e^{-x}}$$

Tanh



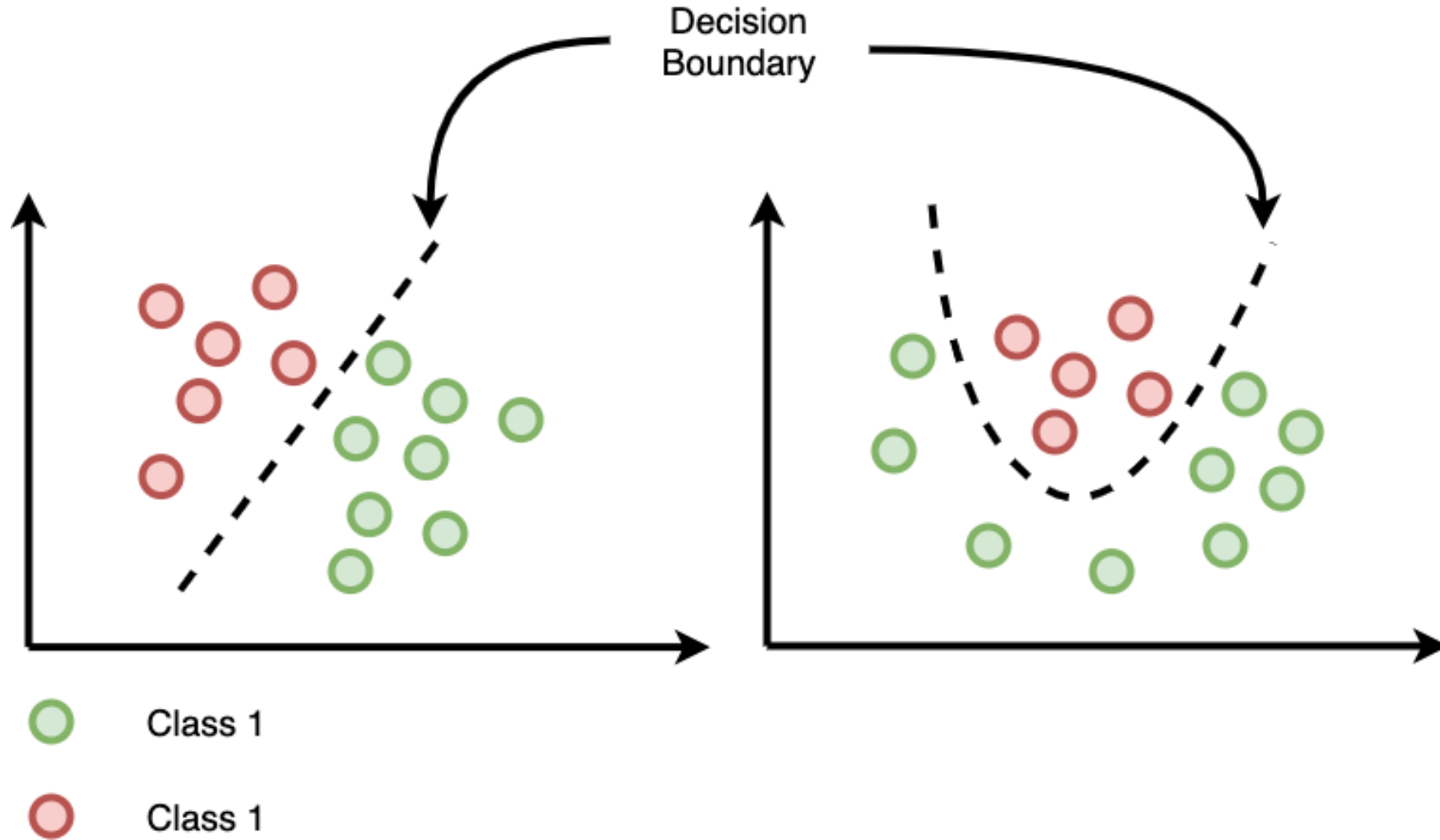
$$f(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

ReLU



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

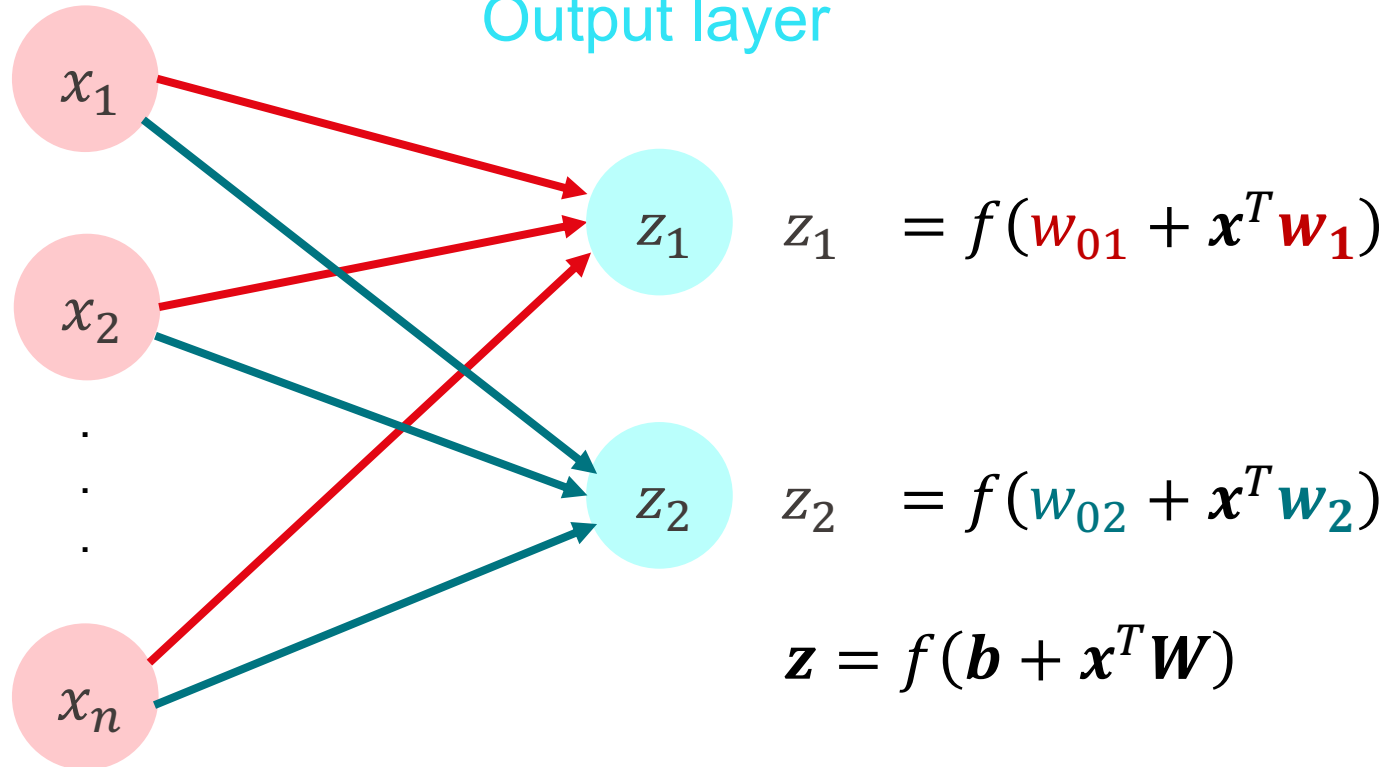
Why is Nonlinearity Important?



Multi-dimensional Output

Input layer

Output layer

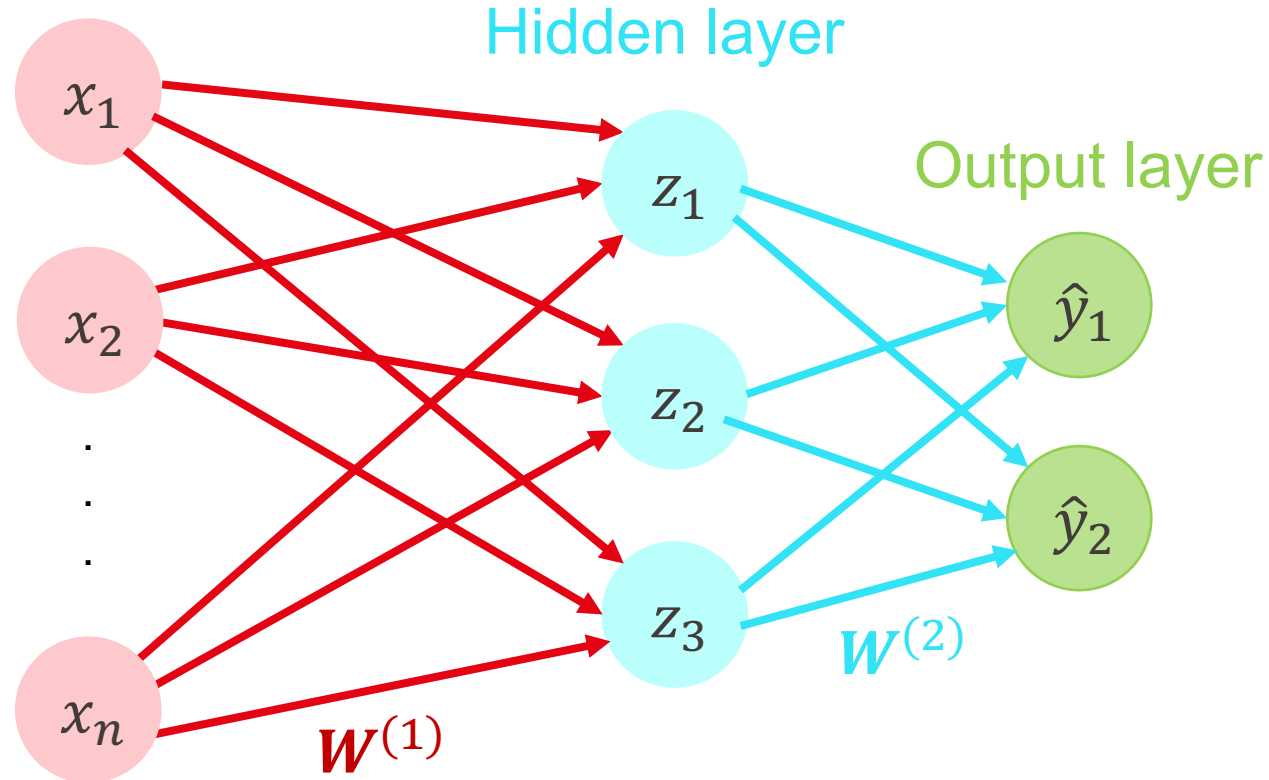


$$\mathbf{z} = f(\mathbf{b} + \mathbf{x}^T \mathbf{W})$$

$$\text{where } \mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}, \mathbf{b} = \begin{bmatrix} w_{01} \\ w_{02} \end{bmatrix} \text{ and } \mathbf{W} = [\mathbf{w}_1 \mathbf{w}_2]$$

Fully Connected Neural Network

Input layer



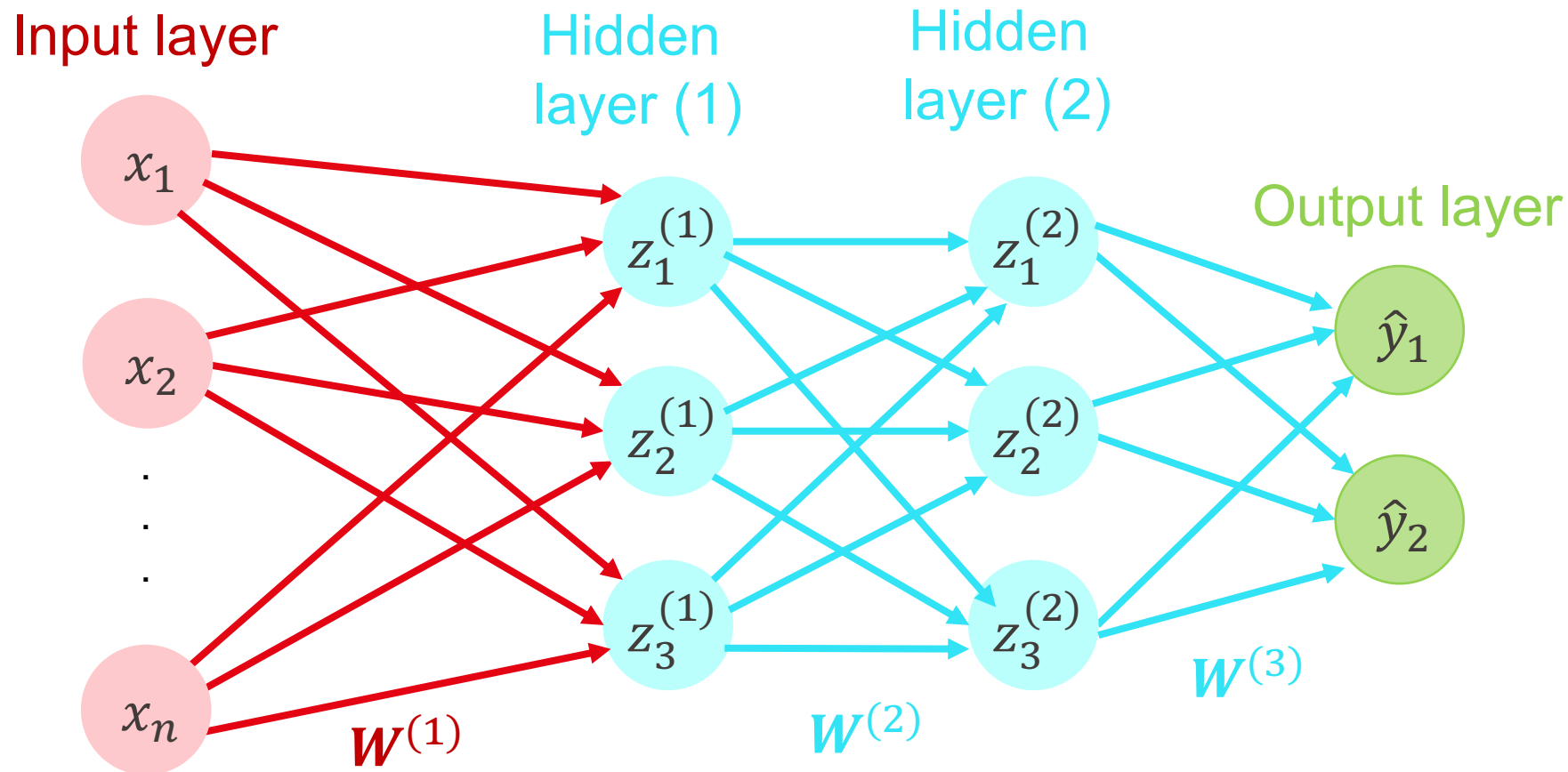
$$\mathbf{z} = f(\mathbf{b}^{(1)} + \mathbf{x}^T \mathbf{W}^{(1)})$$

$$\hat{\mathbf{y}} = f(\mathbf{b}^{(2)} + \mathbf{z}^T \mathbf{W}^{(2)})$$

Forward propagation

Deep Neural Network

■ Stacking multiple layers



Neural Network Learning as Optimization

- **Supervised learning:** we are given input x , and the goal is to predict label y
- **Input x can be:**
 - Vectors of real numbers
 - Sequences (natural language)
 - Matrices (images)
 - Graphs (potentially with node and edge features)
- **We formulate the task as an optimization problem**

Neural Network Learning as Optimization

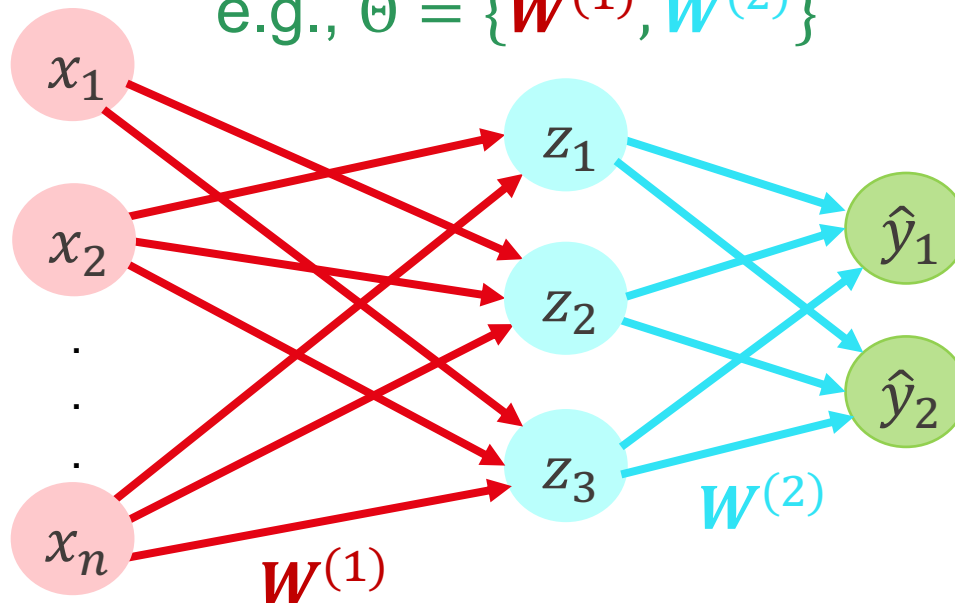
- Formulate the task as an optimization problem:

$$\min_{\Theta} \mathcal{L}(\mathbf{y}, f_{\Theta}(\mathbf{x}))$$

objective function

set of parameters we optimize

e.g., $\Theta = \{\mathbf{W}^{(1)}, \mathbf{W}^{(2)}\}$



$f_{\Theta}(\mathbf{x})$... our ML model

Neural Network Learning as Optimization: Loss Function

What is loss function \mathcal{L} ?

- Depends on a problem
- Quantifies discrepancy between the truth \mathbf{y} and prediction $f(\mathbf{x})$
- Example regression task: L2 loss

$$\mathcal{L}(\mathbf{y}, f(\mathbf{x})) = \|\mathbf{y} - f(\mathbf{x})\|_2$$

- Common loss for classification: **cross entropy** (CE)

$$CE(\mathbf{y}, f(\mathbf{x})) = - \sum_{i=1}^C y_i \log f(\mathbf{x})_i,$$

where C is the number of classes.

$\mathbf{y} =$ **labels vector**

0	0	1	0	0
---	---	---	---	---

$f(\mathbf{x}) =$ **predictions**

0.1	0.3	0.4	0.1	0.1
-----	-----	-----	-----	-----

Optimization: Gradient Descent

How to optimize the objective function?

- **Iterative algorithm:** Repeatedly update weights in the (opposite) direction of gradients until convergence

$$\Theta \leftarrow \Theta - \eta \nabla_{\Theta} \mathcal{L}$$

learning rate
(controls size of the gradient step)

- **Backpropagation:** compute gradient using a chain rule
 - We will derive it in the exercise session!
 - Nowadays done using automatic differentiation
- **Training:** Optimize Θ iteratively
 - **Iteration:** 1 step of gradient descent

Optimization: Gradient Descent

How to optimize the objective function?

Input: input dataset (X, y) , loss function \mathcal{L} , learning rate η , neural network f

1. Initialize parameters Θ
2. Loop until convergence
3. Compute gradient $\nabla_{\Theta} \mathcal{L}$
4. Update weights $\Theta \leftarrow \Theta - \eta \nabla_{\Theta} \mathcal{L}(y, f(X))$

Output: parameters Θ

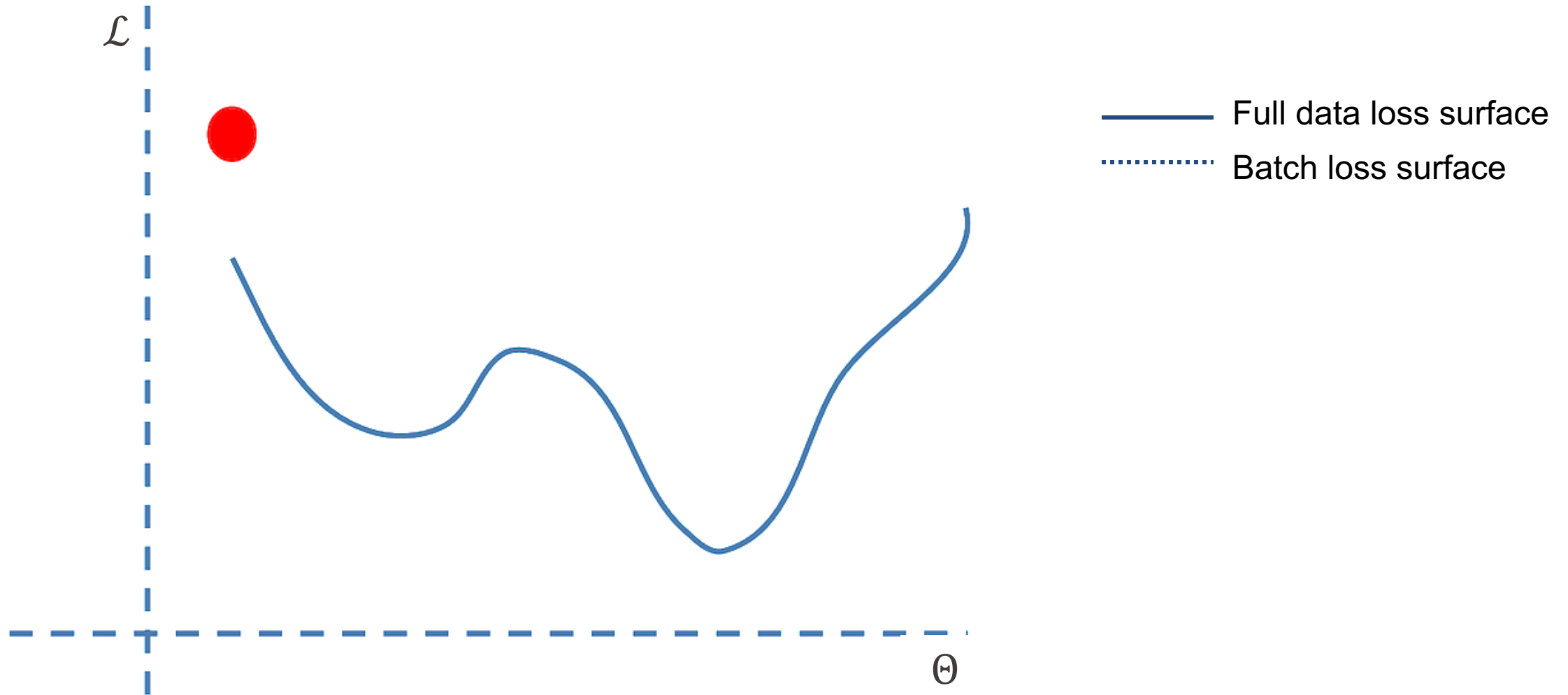
Problem with gradient descent:
Exact gradient requires computing
$$\nabla_{\Theta} \mathcal{L}(y, f(X)),$$

where X is the entire dataset.

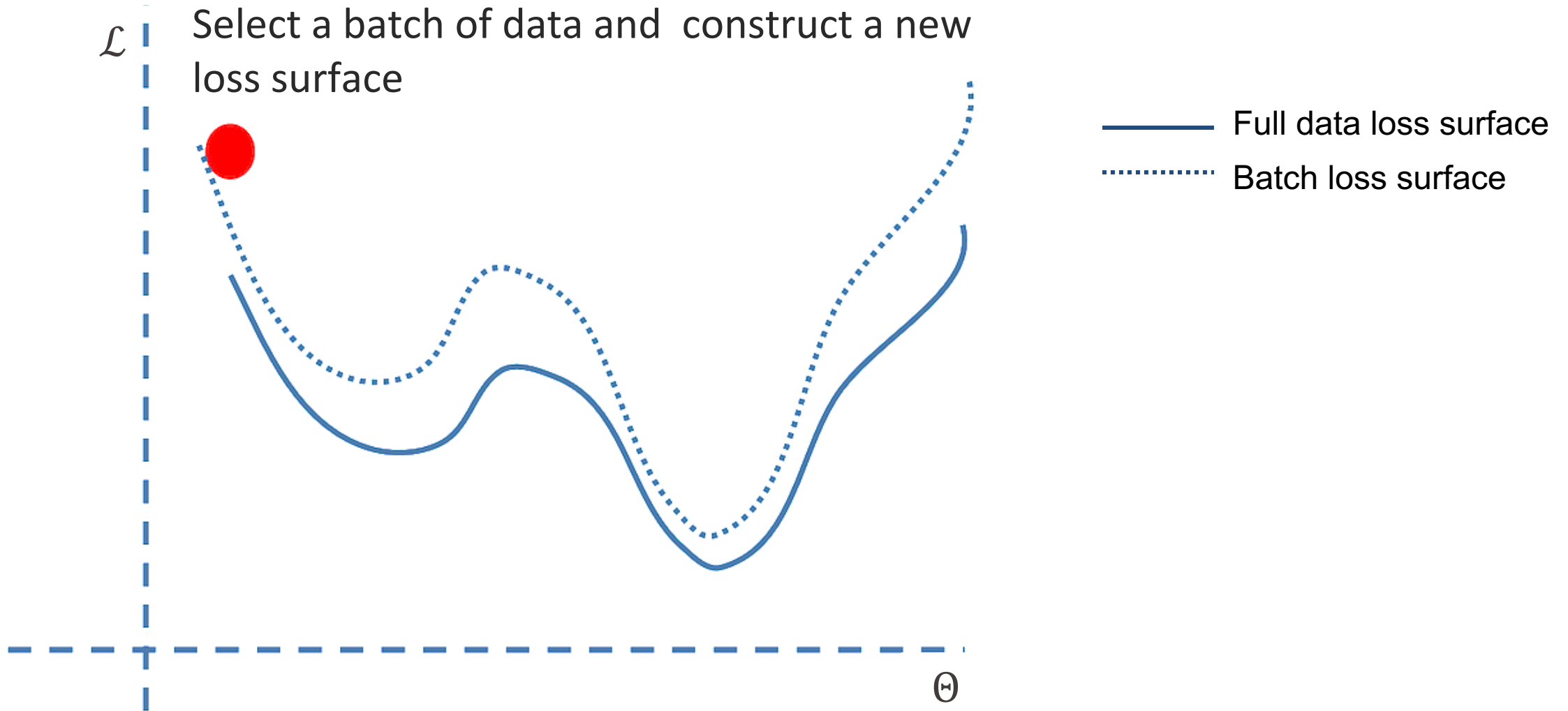
Stochastic Gradient Descent (SGD)

- **Solution:** Stochastic gradient descent
 - At every step, pick a different **minibatch** \mathcal{B} containing a subset of the dataset, use it as input x
- **Concepts:**
 - **Batch size:** the number of data points in a minibatch
 - **Iteration:** 1 step of SGD on a minibatch
 - **Epoch:** one full pass over the dataset (# iterations is equal to ratio of dataset size and batch size)

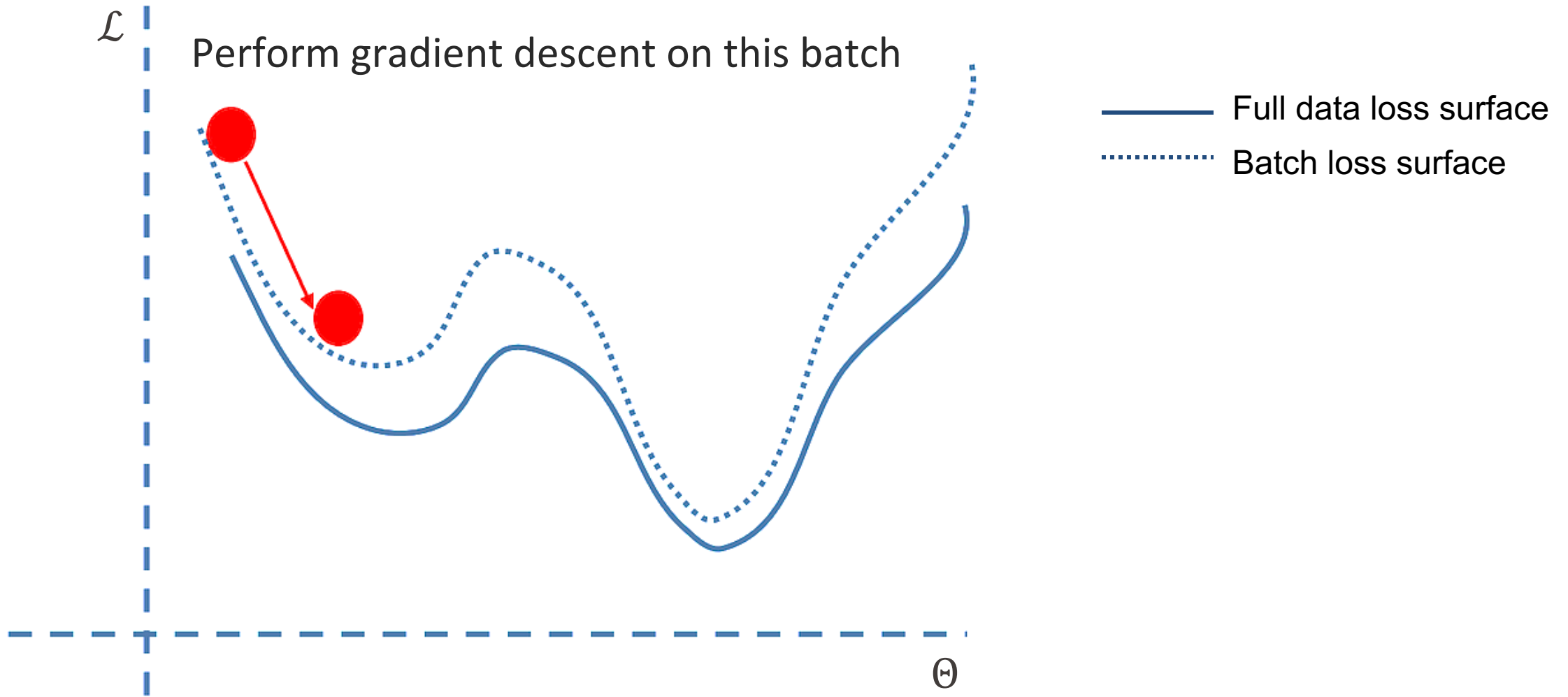
SGD example



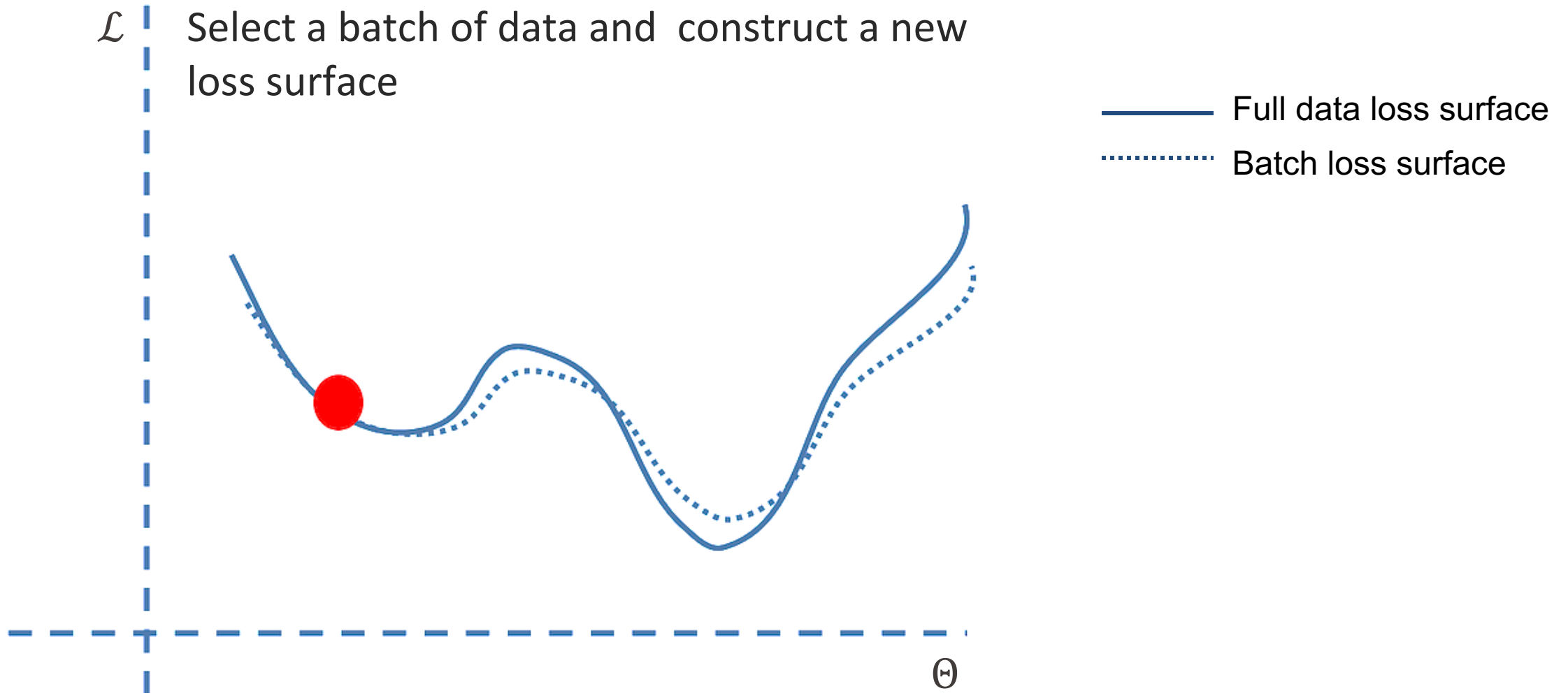
SGD example



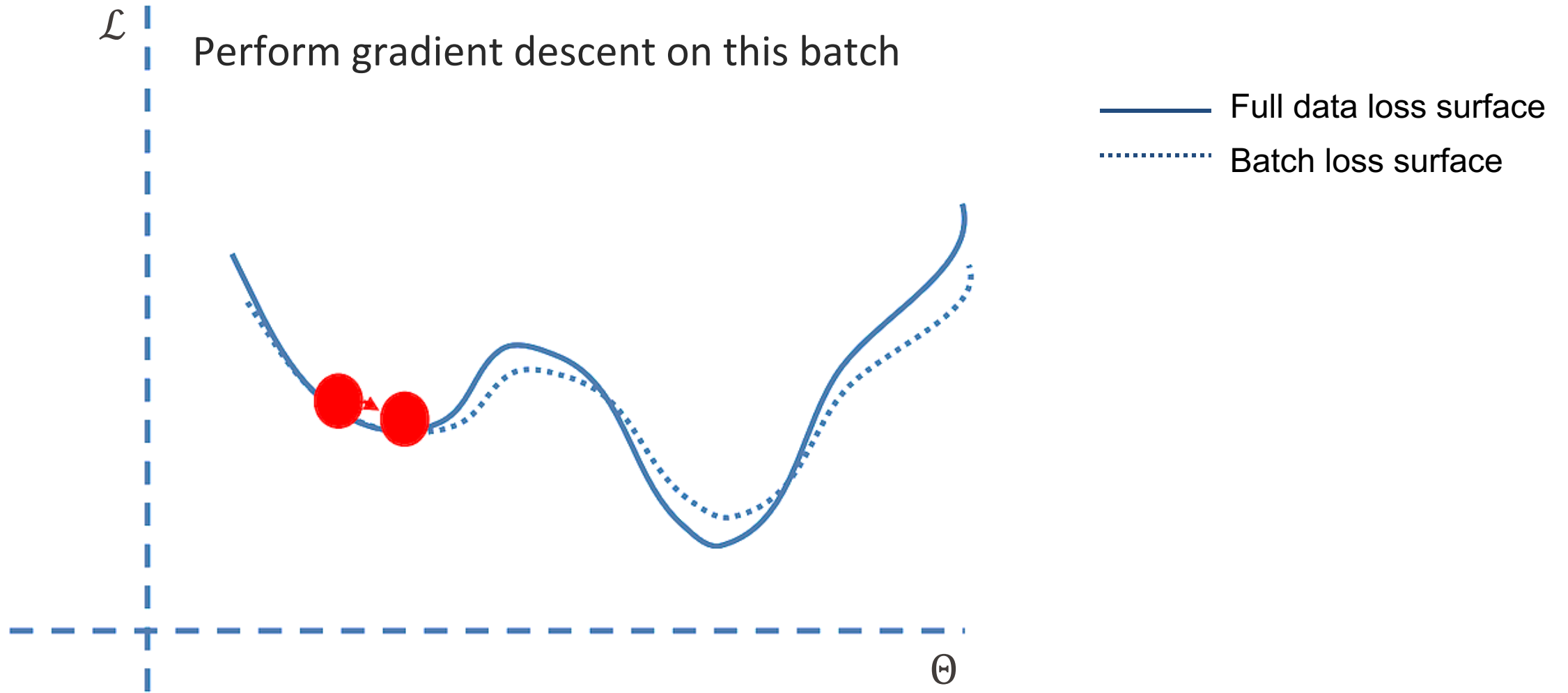
SGD example



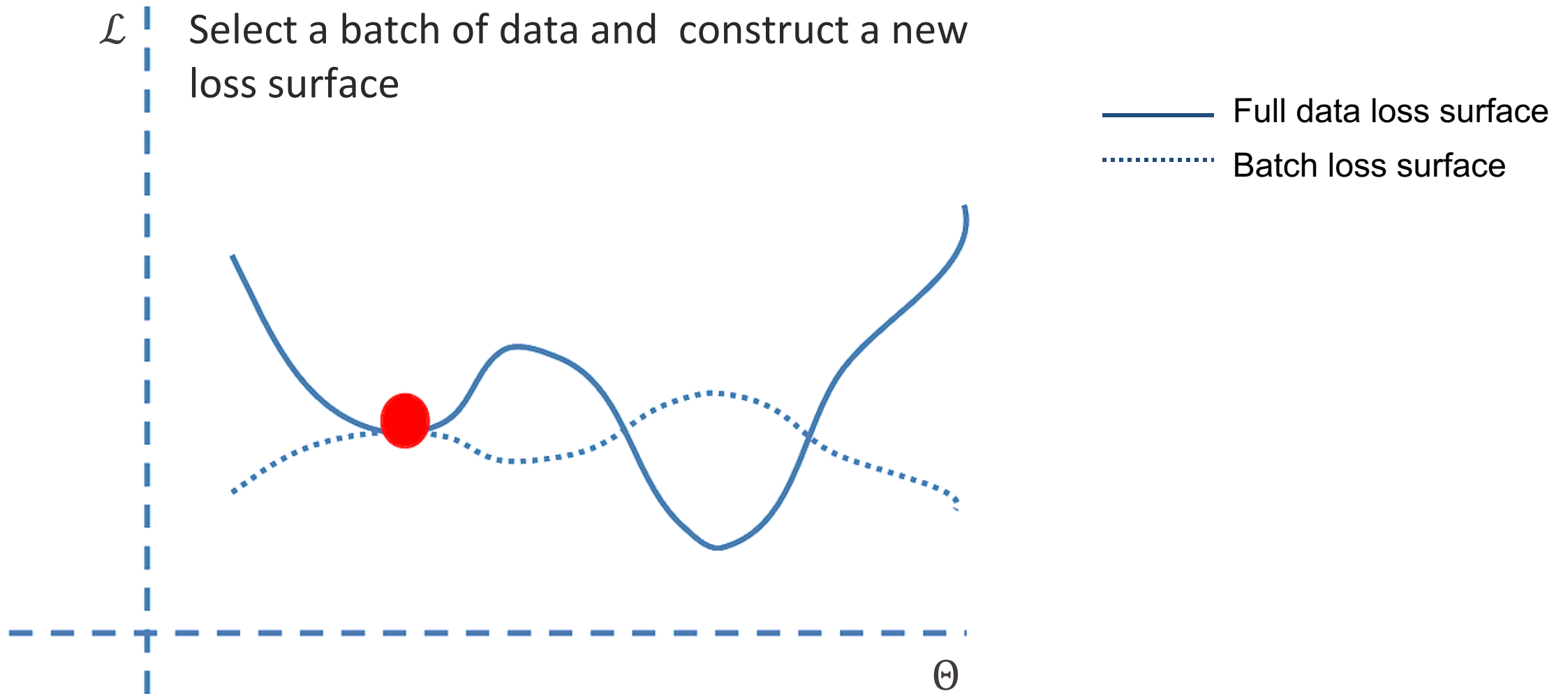
SGD example



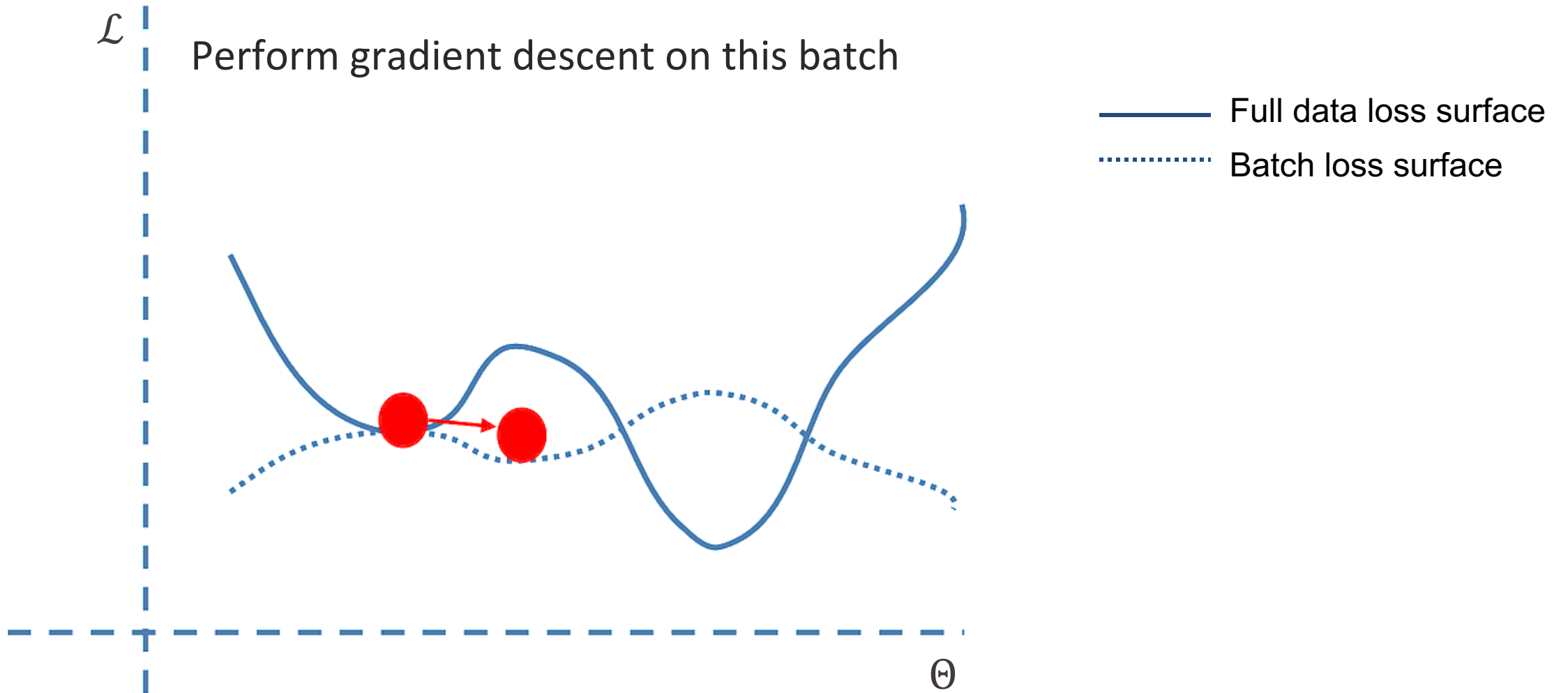
SGD example



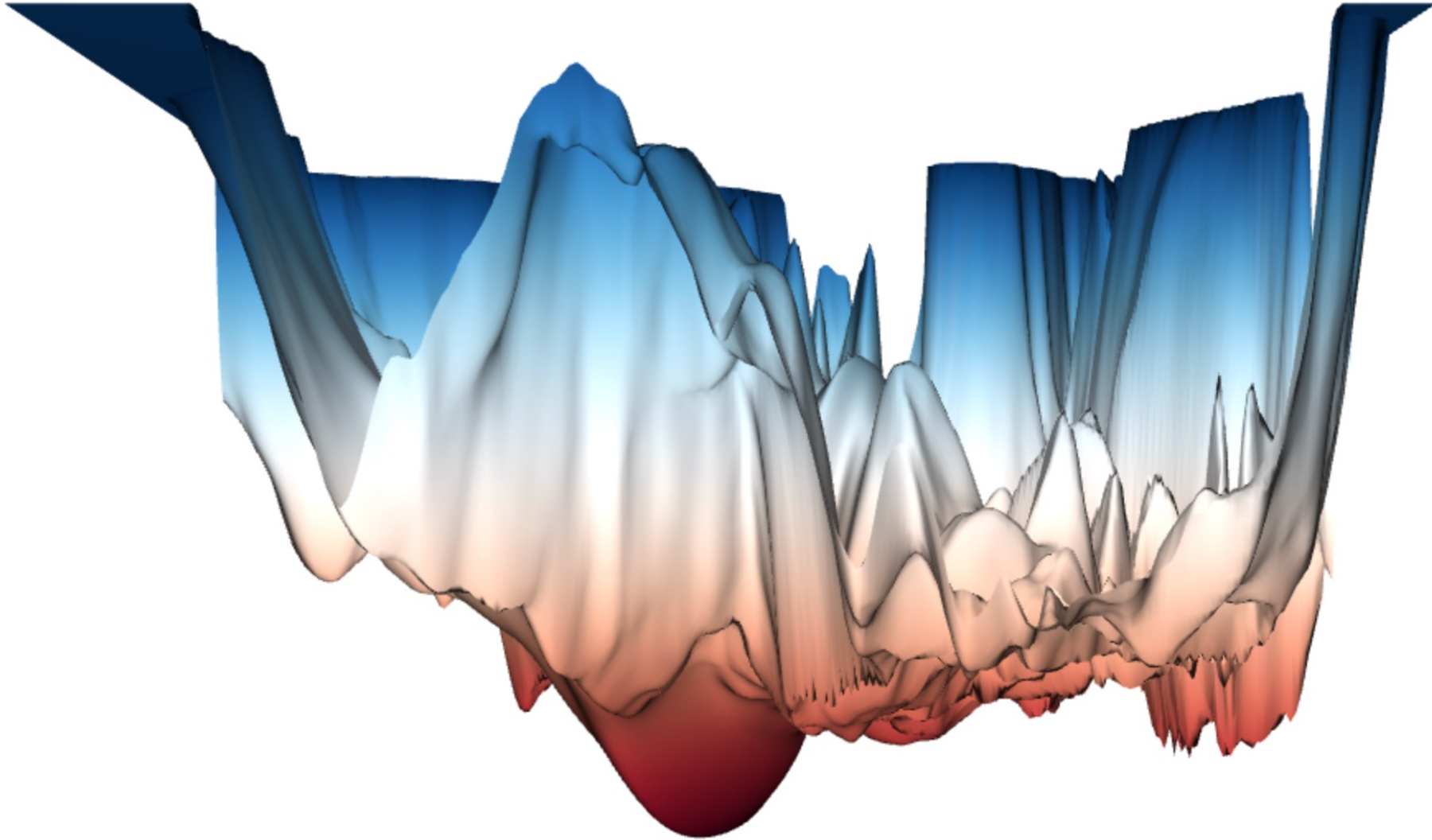
SGD example



SGD example



Loss Surface: Reality



Visualizations from: <http://www.telesens.co/loss-landscape-viz/viewer.html>



33

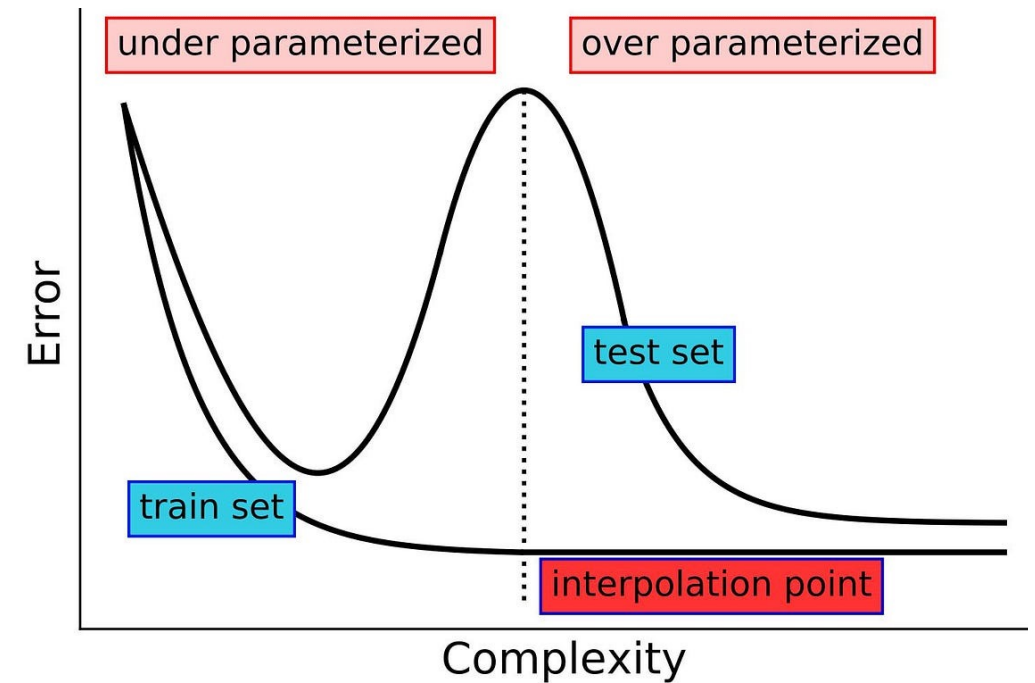
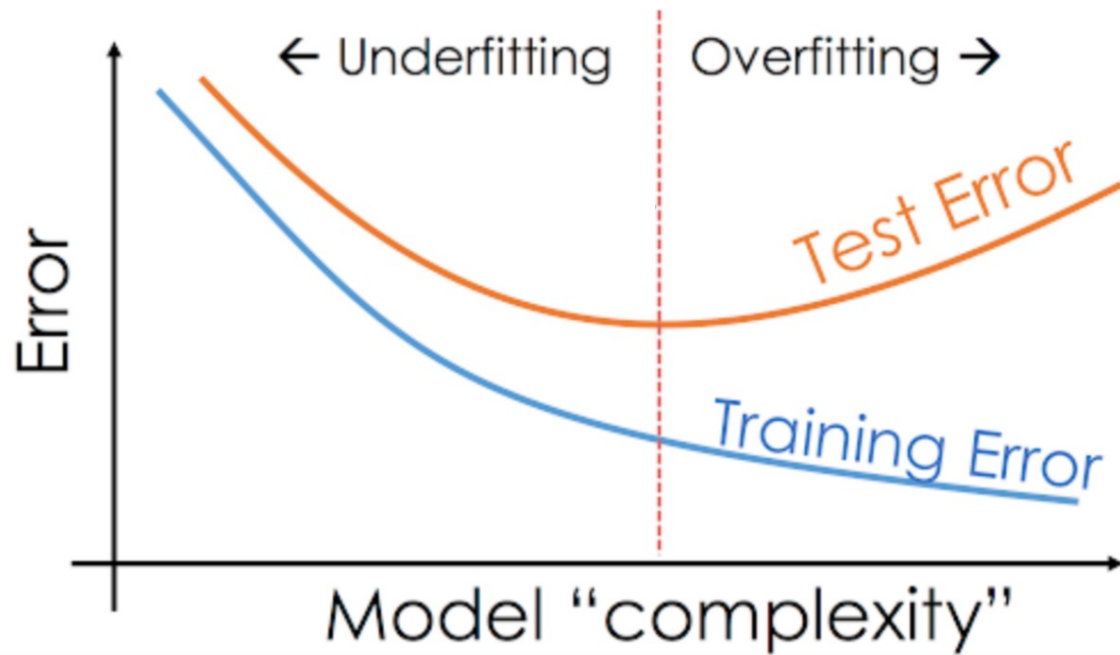
Optimization in Practice

- **How to choose learning rate**
 - Too small → slow convergence
 - Too large → can cause the loss function to fluctuate around the minimum or even diverge
- Common strategies for improvement:
 - **Momentum**: helps accelerate gradients in the right direction
 - **Adaptive learning rate**
- **Common optimizers**: Adam, Adagrad, Adadelata, RMSprop...

Nice overview here: <https://www.ruder.io/optimizing-gradient-descent/>

How to Avoid Overfitting?

Use regularization!



Double descent phenomena

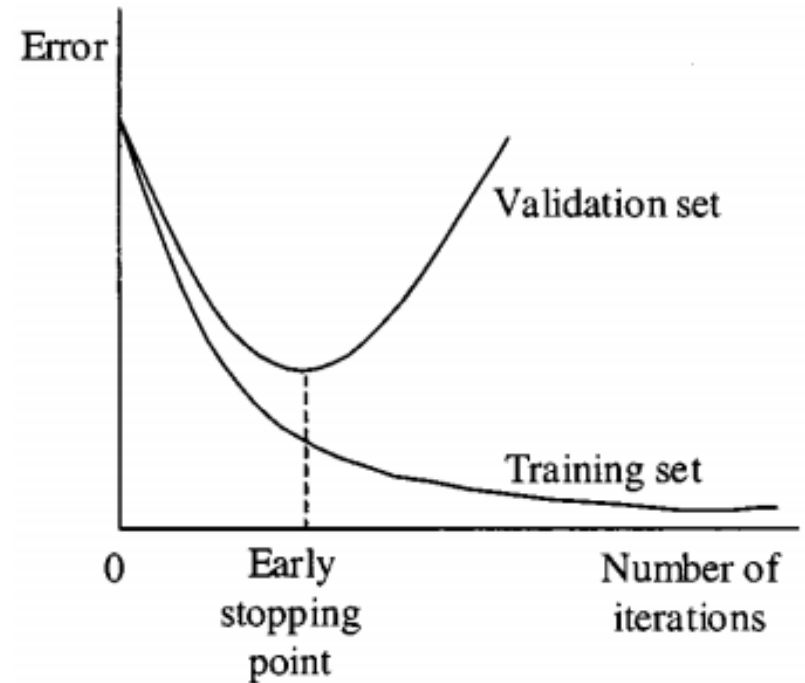
Regularization Techniques

■ Early stopping

- Stop training when performance on a validation dataset starts to degrade

■ Weight regularization

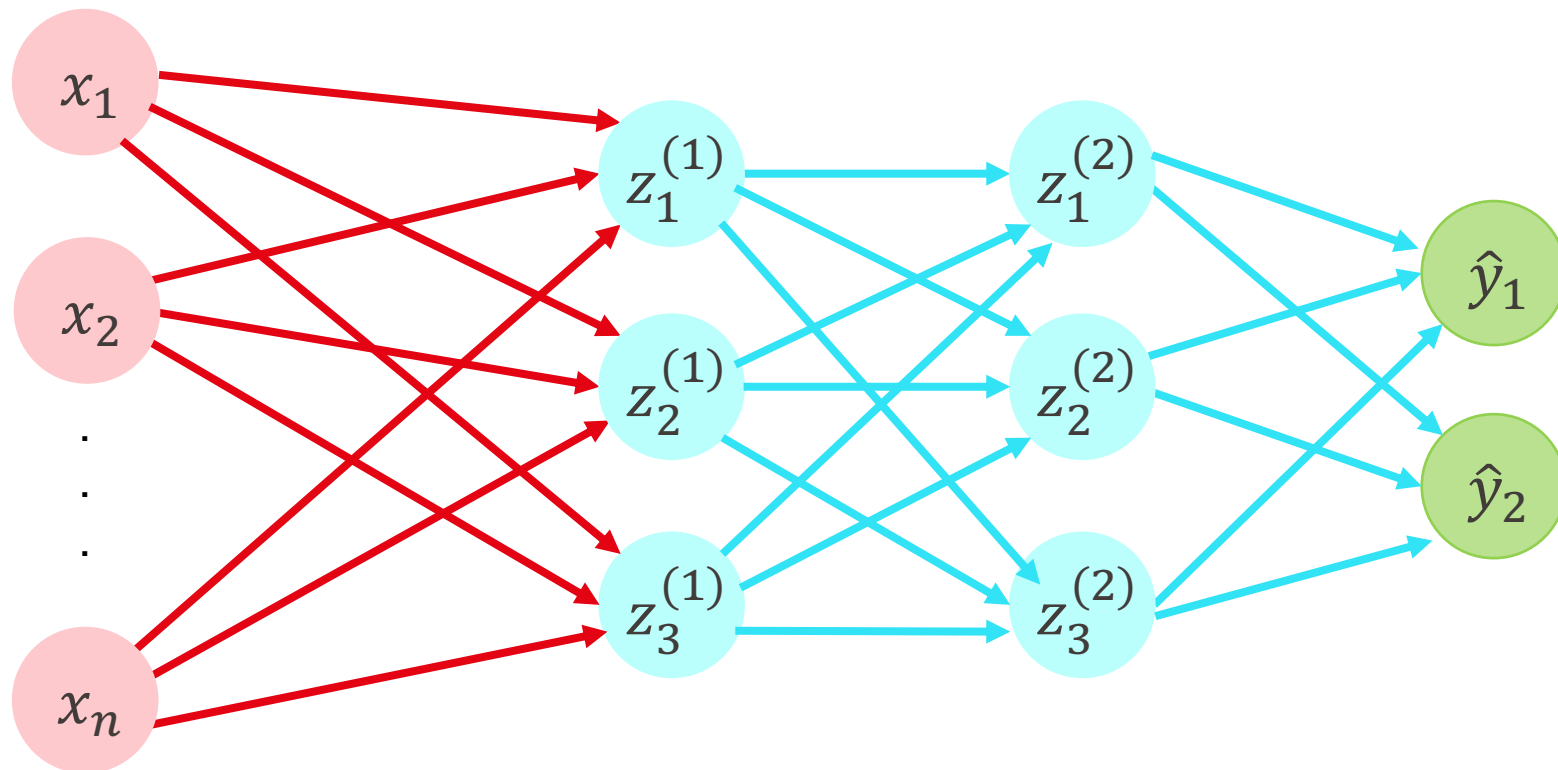
- Penalizing a network based on the size of the network weights during training
- E.g., L1 or L2 vector norm penalty



Regularization Techniques

■ Dropout

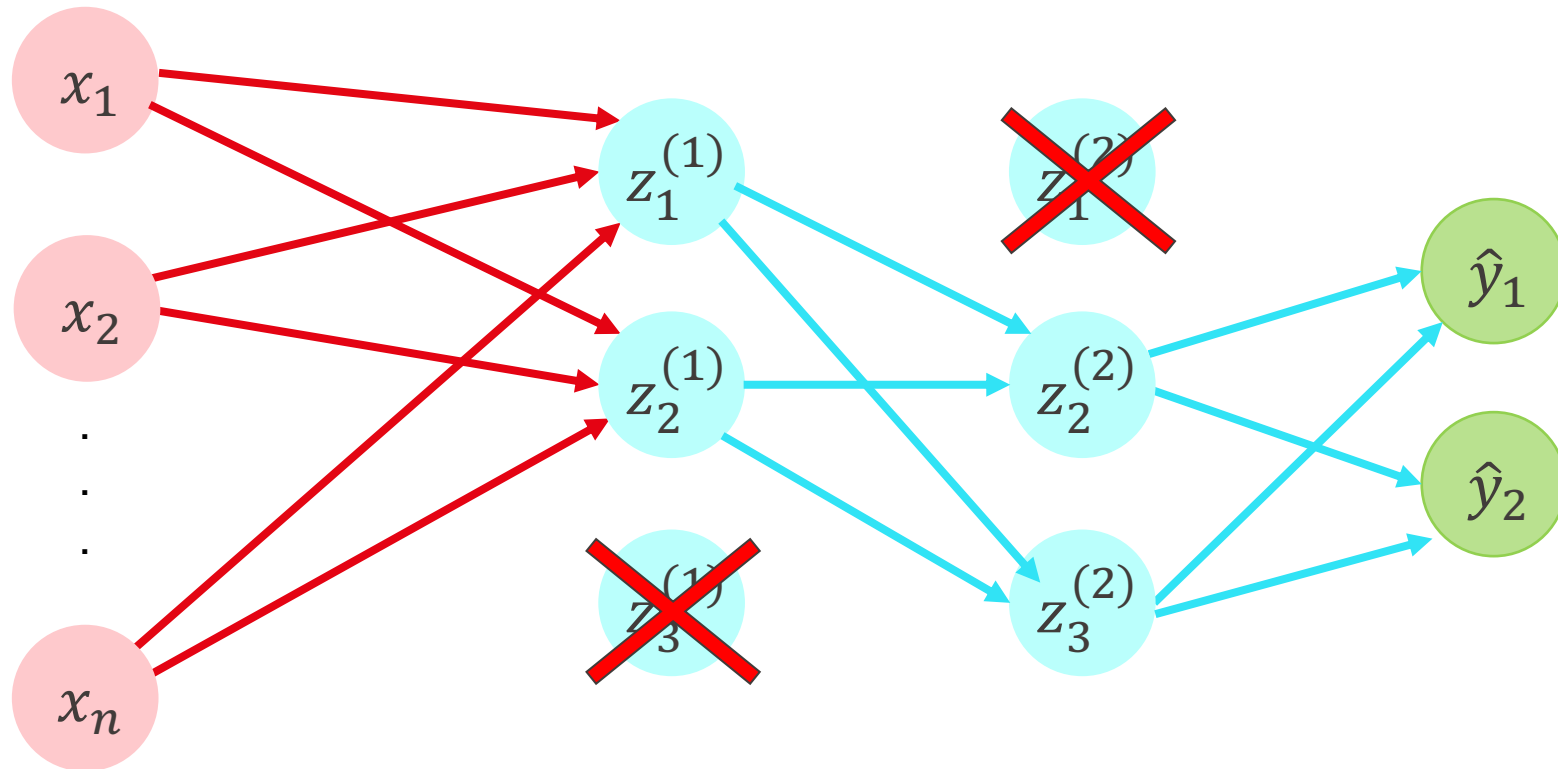
- Randomly drop out neurons during training with some probability



Regularization Techniques

■ Dropout

- Randomly drop out neurons during training with some probability



Regularization Techniques

More details in
exercise session!

■ Batch normalization

- More than “just regularization”
- Allows higher learning rates
- Reduces the strong dependence on initialization
- Usually applied after a fully connected/convolutional layer and before a non-linearity

$$x_i \leftarrow \gamma \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} + \beta$$

hyperparameters

mean and standard-deviation
calculated over the mini-batches

Hyperparameter Optimization

- **Hyperparameters:** Parameters used to control the learning the learning process; they are **not learnt**
 - E.g., learning rate, number of neural network layers, mini-batch size
- How to find hyperparameters:
 - Grid search, random search, Bayesian optimization
- **Cross-validation** is often used to estimate generalization performance of the model
 - Partition the data in non-overlapping subsets of train, validation and test data and evaluate model's performance on the held-out data

Recap

- **Objective function:**

$$\min_{\Theta} \mathcal{L}(\mathbf{y}, f(\mathbf{x}))$$

- Sample a minibatch of input \mathbf{x}
- **Forward propagation:** compute \mathcal{L} given \mathbf{x}
- **Back-propagation:** obtain gradient $\nabla_{\Theta} \mathcal{L}$ using a chain rule
- Use **stochastic gradient descent (SGD)** to optimize for Θ over many iterations
- f can be any neural network

Any Feedback?

Give us feedback on the lecture:

- <https://docs.google.com/forms/d/e/1FAIpQLSdgn11OkIAvx9iIKLpviRwJ9gruCdLGUREFWPIpYR0QycY6jA/viewform>

Next lectures: Different neural network architectures

- **CNNs** for images
- **GNNs** for graphs
- **Transformers** for sequence data