

Introduction



Most generally, a machine learning algorithm can be thought of as a black box. It takes inputs and gives outputs.

The purpose of this course is to show you how to create this 'black box' and tailor it to your needs.

For example, we may create a model that predicts the weather tomorrow, based on meteorological data about the past few days.

The "black box" in fact is a mathematical model. The machine learning algorithm will follow a kind of trial-and-error method to determine the model that estimates the outputs, given inputs.

Once we have a model, we must **train** it. **Training** is the process through which, the model **learns** how to make sense of input data.

Types of machine learning

Supervised

It is called *supervised* as we provide the algorithm not only with the inputs, but also with the targets (desired outputs). This course focuses on supervised machine learning.

Based on that information the algorithm learns how to produce outputs as close to the **targets** as possible.

The objective function in supervised learning is called **loss function** (also cost or error). We are trying to minimize the loss as the lower the loss function, the higher the accuracy of the model.

Common methods:

- Regression
- Classification

Unsupervised

In *unsupervised* machine learning, the researcher feeds the model with inputs, but **not** with targets. Instead she asks it to find some sort of dependence or underlying logic in the data provided.

For example, you may have the financial data for 100 countries. The model manages to divide (cluster) them into 5 groups. You then examine the 5 clusters and reach the conclusion that the groups are: "Developed", "Developing but overachieving", "Developing but underachieving", "Stagnating", and "Worsening".

The algorithm divided them into 5 groups based on **similarities**, but you didn't know what similarities. It could have divided them by location instead.

Common methods:
• Clustering

Reinforcement

In reinforcement ML, the goal of the algorithm is to maximize its reward. It is inspired by human behavior and the way people change their actions according to incentives, such as getting a reward or avoiding punishment.

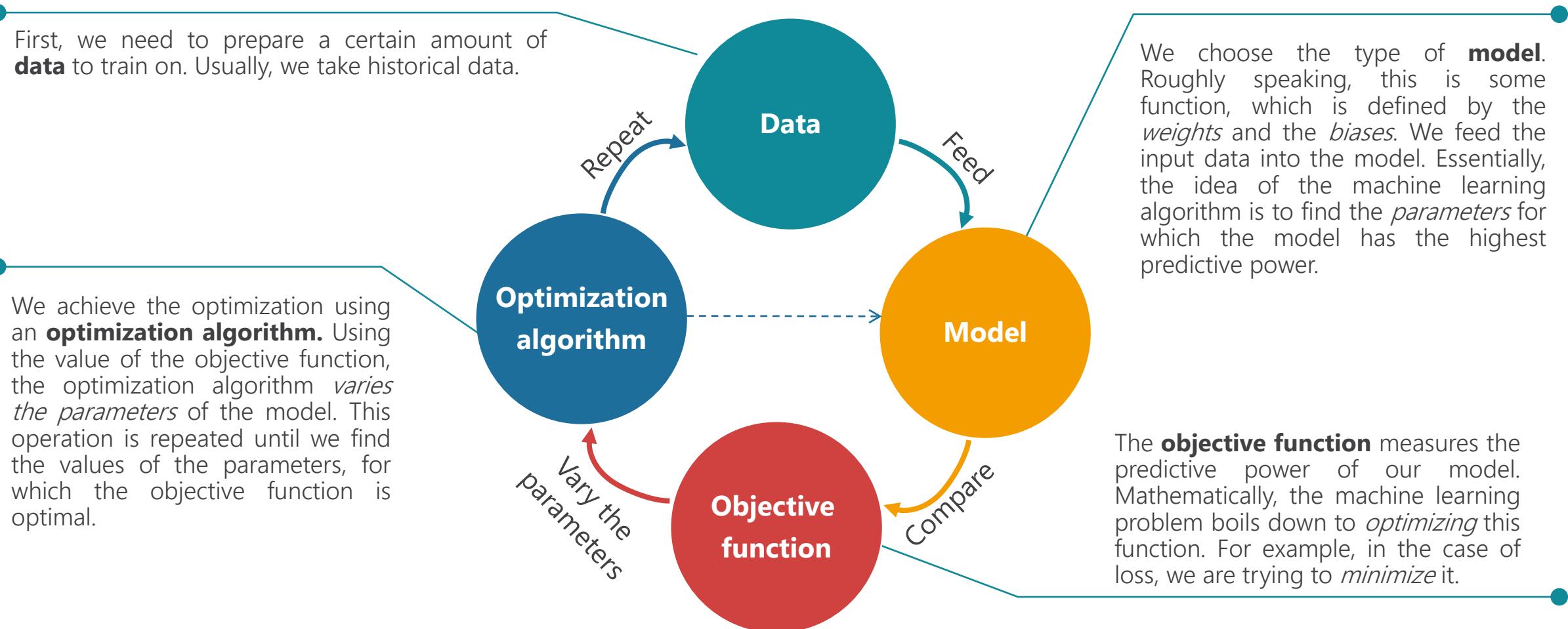
The objective function is called a **reward function**. We are trying to maximize the reward function.

An example is a computer playing Super Mario. The higher the score it achieves, the better it is performing. The score in this case is the objective function.

Common methods:
• Decision process
• Reward system

Building blocks of a machine learning algorithm

The basic logic behind training an algorithm involves four ingredients: data, model, objective function, and optimization algorithm. They are **ingredients**, instead of steps, as the process is iterative.



Types of supervised learning

Supervised learning could be split into two subtypes – **regression** and **classification**.

Regression

Regression outputs are continuous numbers.

Examples:

Predicting the EUR/USD exchange rate tomorrow. The output is a number, such as 1.02, 1.53, etc.

Predicting the price of a house, e.g. \$234,200 or \$512,414.

One of the main properties of the regression outputs is that they are ordered. As $1.02 < 1.53$, and $234200 < 512414$, we can surely say that one output is bigger than the other.

This distinction proves to be crucial in machine learning as the algorithm somewhat *gets additional information* from the outputs.

Classification

Classification outputs are labels from some sort of class.

Examples:

Classifying photos of animals. The classes are "cats", "dogs", "dolphins", etc.

Predicting conversions in a business context. You can have two classes of customers, e.g."will buy again" and "won't buy again".

In the case of classification, the labels are not ordered and cannot be compared at all. A dog photo is not "more" or "better" than a cat photo (objectively speaking) in the way a house worth \$512,414 is "more" (more expensive) than a house worth \$234,200.

This distinction proves to be crucial in machine learning, as the different classes are treated on **an equal footing**.

Model

The simplest possible model is **a linear model**. Despite appearing unrealistically simple, in the deep learning context, it is the basis of more complicated models.

$$y = \mathbf{X}\mathbf{w} + \mathbf{b}$$

The diagram shows the linear model equation $y = \mathbf{X}\mathbf{w} + \mathbf{b}$. Four arrows originate from labels below the equation and point to its components: 'output(s)' points to y , 'input(s)' points to \mathbf{X} , 'weight(s)' points to \mathbf{w} , and 'bias(es)' points to \mathbf{b} .

There are four elements.

- The input(s), x . That's basically the data that we feed to the model.
- The weight(s), w . Together with the biases, they are called **parameters**. The optimization algorithm will vary the weights and the biases, in order to produce the model that fits the data best.
- The bias(es), b . See weight(s).
- The output(s), y . y is a function of x , determined by w and b .

Each model is determined solely by its parameters (the weights and the biases). That is why we are interested in varying them using the (kind of) trial-and-error method, until we find a model that explains the data sufficiently well.

Model - Continued

$$y = Xw + b$$

Diagram illustrating the dimensions of the variables in the equation $y = Xw + b$:

- y : $n \times m$
- X : $n \times k$
- w : $k \times m$
- b : $1 \times m$

Where:

- n is the number of samples (observations)
- m is the number of output variables
- k is the number of input variables

A linear model can represent multidimensional relationships. The shapes of y , x , w , and b are given above (notation is arbitrary).

The simplest linear model, where $n = m = k = 1$.

$$\boxed{y} = \boxed{x} \boxed{w} + \boxed{b}$$

The simplest linear model, where $m = k = 1, n > 1$

$$\begin{array}{c|c|c|c|c} y_1 & = & x_1w + b & = & \boxed{x_1} \boxed{w} + \boxed{b} \\ y_2 & & x_2w + b & & \\ \dots & & \dots & & \\ y_n & & x_nw + b & & \end{array}$$

Examples:

$$\boxed{27} = \boxed{4} \boxed{6} + \boxed{3}$$

Since the weights and biases alone define a model, this example shows **the same model** as above but for many data points.

$$\begin{array}{c|c|c|c} 27 & = & \boxed{4} & \boxed{6} + \boxed{3} \\ 33 & & \boxed{5} & \\ \dots & & \dots & \\ -3 & & -1 & \end{array}$$

Note that we add the bias to each row, essentially simulating an $n \times 1$ matrix. That's how computers treat such matrices, that's why in machine learning we define the bias in this way.

Model multiple inputs

$$y = \mathbf{xw} + b$$

The diagram illustrates the dimensions of the variables in the equation $y = \mathbf{xw} + b$:

- x is labeled $n \times k$ with an arrow pointing to it.
- w is labeled $k \times m$ with an arrow pointing to it.
- b is labeled $1 \times m$ with an arrow pointing to it.

Where:

- n is the number of samples (observations)
 - m is the number of output variables
 - k is the number of input variables

We can extend the model to multiple inputs where $n, k > 1, m = 1$.

$$\begin{array}{l}
 \begin{array}{c}
 y_1 \\
 y_2 \\
 \dots \\
 \dots \\
 y_n
 \end{array}
 = \begin{array}{c}
 x_{11}w_1 + x_{12}w_2 + \dots + x_{1k}w_k + b \\
 x_{21}w_1 + x_{22}w_2 + \dots + x_{2k}w_k + b \\
 \dots \\
 \dots \\
 x_{n1}w_1 + x_{n2}w_2 + \dots + x_{nk}w_k + b
 \end{array}
 = \begin{array}{c}
 x_{11} & x_{12} & \dots & x_{1k} \\
 x_{21} & x_{22} & \dots & x_{2k} \\
 \dots & \dots & \dots & \dots \\
 \dots & \dots & \dots & \dots \\
 x_{n1} & x_{n2} & \dots & x_{nk}
 \end{array}
 \begin{array}{c}
 w_1 \\
 w_2 \\
 \dots \\
 \dots \\
 w_k
 \end{array}
 + \begin{array}{c}
 1 \\
 \\
 \\
 \\
 k \times 1
 \end{array}
 \end{array}$$

Note that we add the bias to each row, essentially simulating an $n \times 1$ matrix. That's how computers treat such matrices, that's why in machine learning we define the bias in this way.

Example on the next page.

Model multiple inputs

$$y = Xw + b$$

Diagram illustrating the dimensions of the variables in the equation $y = Xw + b$:

- y : $n \times m$
- X : $n \times k$
- w : $k \times m$
- b : $1 \times m$

Where:

- n is the number of samples (observations)
- m is the number of output variables
- k is the number of input variables

We can extend the model to multiple inputs where $n, k > 1$ $m = 1$. An example.

$$\begin{array}{c|c} y_1 & = 9x(-2) + 12x5 + \dots + 13x(-1) + 4 \\ y_2 & = 10x(-2) + 6x5 + \dots + 2x(-1) + 4 \\ \dots & \dots \\ \dots & \dots \\ y_n & = 7x(-2) + 7x5 + \dots + 1x(-1) + 4 \end{array} = \begin{array}{c|c|c|c} 9 & 12 & \dots & 13 \\ 10 & 6 & \dots & 2 \\ \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots \\ 7 & 7 & \dots & 1 \end{array} \begin{array}{c|c} -2 & + 4 \\ 5 & \\ \dots & \\ -1 & \end{array} \begin{array}{c} 1 \times 1 \\ k \times 1 \end{array}$$

$n \times 1$ $n \times k$

Note that we add the bias to each row, essentially simulating an $n \times 1$ matrix. That's how computers treat such matrices, that's why in machine learning we define the bias in this way.

Model multiple inputs and multiple outputs

$$y = Xw + b$$

Diagram illustrating the dimensions of the variables in the equation $y = Xw + b$:

- y : $n \times m$
- X : $n \times k$
- w : $k \times m$
- b : $1 \times m$

Where:

- n is the number of samples (observations)
- m is the number of output variables
- k is the number of input variables

We can extend the model to multiple inputs where $n, k, m > 1$.

$$\begin{array}{|c|c|c|} \hline y_{11} & \dots & y_{1m} \\ \hline y_{21} & \dots & y_{2m} \\ \hline \dots & \dots & \dots \\ \hline \dots & \dots & \dots \\ \hline y_{n1} & \dots & y_{nm} \\ \hline \end{array} = \begin{array}{|c|c|c|c|} \hline x_{11} & x_{12} & \dots & x_{1k} \\ \hline x_{21} & x_{22} & \dots & x_{2k} \\ \hline \dots & \dots & \dots & \dots \\ \hline \dots & \dots & \dots & \dots \\ \hline x_{n1} & x_{n2} & \dots & x_{nk} \\ \hline \end{array} + \begin{array}{|c|c|c|} \hline w_{11} & \dots & w_{1m} \\ \hline w_{21} & \dots & w_{2m} \\ \hline \dots & \dots & \dots \\ \hline \dots & \dots & \dots \\ \hline w_{k1} & \dots & w_{km} \\ \hline \end{array} \begin{array}{c} b_1 \dots b_m \\ 1 \times m \\ k \times m \end{array}$$

Note that we add the bias to each row, essentially simulating an $n \times m$ matrix. That's how computers treat such matrices, that's why in machine learning we define the bias in this way.

Example on the next page.

Model multiple inputs and multiple outputs

$$y = Xw + b$$

$n \times m$
 $n \times k$
 $k \times m$
 $1 \times m$

Where:

- n is the number of samples (observations)
- m is the number of output variables
- k is the number of input variables

We can extend the model to multiple inputs where $n, k, m > 1$.

$$\begin{array}{|c|c|c|} \hline
 y_{11} & \dots & y_{1m} \\ \hline
 y_{21} & \dots & y_{2m} \\ \hline
 \dots & \dots & \dots \\ \hline
 \dots & \dots & \dots \\ \hline
 y_{n1} & \dots & y_{nm} \\ \hline
 \end{array}
 = \begin{array}{|c|c|c|} \hline
 x_{11}w_{11} + x_{12}w_{21} + \dots + x_{1k}w_{k1} + b_1 & \dots & x_{11}w_{1m} + x_{12}w_{2m} + \dots + x_{1k}w_{km} + b_m \\ \hline
 x_{21}w_{11} + x_{22}w_{21} + \dots + x_{2k}w_{k1} + b_1 & \dots & x_{21}w_{1m} + x_{22}w_{2m} + \dots + x_{2k}w_{km} + b_m \\ \hline
 \dots & \dots & \dots \\ \hline
 \dots & \dots & \dots \\ \hline
 x_{n1}w_{11} + x_{n2}w_{21} + \dots + x_{nk}w_{k1} + b_1 & \dots & x_{n1}w_{1m} + x_{n2}w_{2m} + \dots + x_{nk}w_{km} + b_m \\ \hline
 \end{array}
 = \begin{array}{|c|c|c|c|} \hline
 x_{11} & x_{12} & \dots & x_{1k} \\ \hline
 x_{21} & x_{22} & \dots & x_{2k} \\ \hline
 \dots & \dots & \dots & \dots \\ \hline
 \dots & \dots & \dots & \dots \\ \hline
 x_{n1} & x_{n2} & \dots & x_{nk} \\ \hline
 \end{array}
 \begin{array}{|c|c|c|} \hline
 w_{11} & \dots & w_{1m} \\ \hline
 w_{21} & \dots & w_{2m} \\ \hline
 \dots & \dots & \dots \\ \hline
 \dots & \dots & \dots \\ \hline
 w_{k1} & \dots & w_{km} \\ \hline
 \end{array}
 +
 \begin{array}{|c|c|c|} \hline
 b_1 & \dots & b_m \\ \hline
 \end{array}$$

$n \times m$
 $n \times k$
 $1 \times m$
 $k \times m$

Note that we add the bias to each row, essentially simulating an $n \times m$ matrix. That's how computers treat such matrices, that's why in machine learning we define the bias in this way.

Objective function

The objective function is a measure of how well our model's outputs match the targets.

The **targets** are the "correct values" which we aim at. In the cats and dogs example, the targets were the "labels" we assigned to each photo (either "cat" or "dog").

Objective functions can be split into two types: **loss** (supervised learning) and **reward** (reinforcement learning). Our focus is supervised learning.



$$\sum_i (y_i - t_i)^2$$

The L2-norm of a vector, \mathbf{a} , (Euclidean length) is given by

$$\|\mathbf{a}\| = \sqrt{\mathbf{a}^T \cdot \mathbf{a}} = \sqrt{a_1^2 + \dots + a_n^2}$$

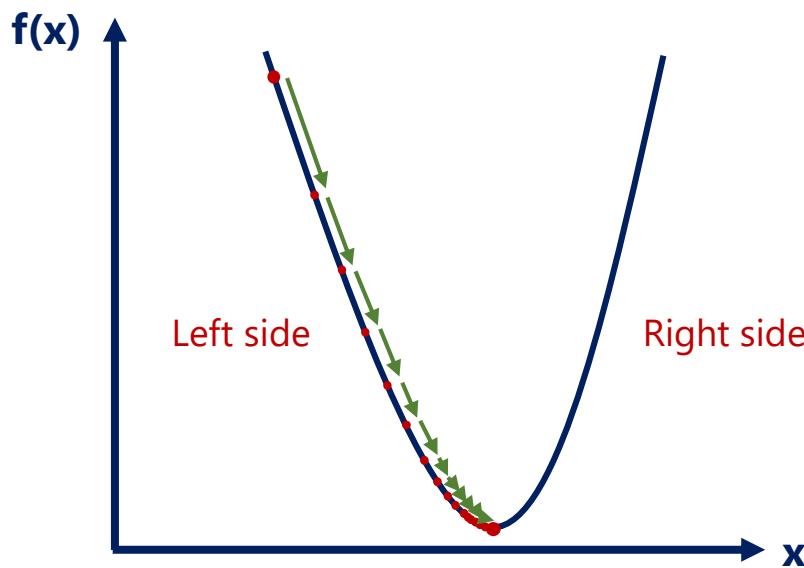
The main rationale is that the L2-norm loss is basically the distance from the origin (0). So, the closer to the origin is the difference of the outputs and the targets, the lower the loss, and the better the prediction.

$$-\sum_i t_i \ln y_i$$

The cross-entropy loss is mainly used for classification. Entropy comes from information theory, and measures how much information one is missing to answer a question. Cross-entropy (used in ML) works with probabilities – one is our opinion, the other – the true probability (the probability of a target to be correct is 1 by definition). If the cross-entropy is 0, then we are **not missing any information** and have a perfect model.

Gradient descent

The last ingredient is the optimization algorithm. The most commonly used one is the gradient descent. The main point is that we can find the minimum of a function by applying the rule: $x_{i+1} = x_i - \eta f'(x_i)$, where η is a small enough positive number. In machine learning, η , is called the learning rate. The rationale is that the first derivative at x_i , $f'(x_i)$ shows the slope of the function at x_i .



If the first derivative of $f(x)$ at x_i , $f'(x_i)$, is negative, then we are on the left side of the parabola (as shown in the figure). Subtracting a negative number from x_i (as η is positive), will result in x_{i+1} , that is bigger than x_i . This would cause our next *trial* to be on the right; thus, closer to the sought minimum.

Alternatively, if the first derivative is positive, then we are on the right side of the parabola. Subtracting a positive number from x_i will result in a lower number, so our next trial will be to the left (again closer to the minimum).

So, either way, using this rule, we are approaching the minimum. When the first derivative is 0, we have reached the minimum. Of course, our update rule won't update anymore ($x_{i+1} = x_i - 0$).

The learning rate η , must be low enough so we don't oscillate (bounce around without reaching the minimum) and big enough, so we reach it in rational time.

In machine learning, $f(x)$ is the **loss function**, which we are trying to minimize.



The variables that we are varying until we find the minimum are the **weights and the biases**. The proper update rules are:

$$\mathbf{w}_{i+1} = \mathbf{w}_i - \eta \nabla_{\mathbf{w}} L(\mathbf{w}_i) = \mathbf{w}_i - \eta \sum_i x_i \delta_i \quad \text{and} \quad b_{i+1} = b_i - \eta \nabla_b L(b_i) = b_i - \eta \sum_i \delta_i$$

Gradient descent. Multivariate derivation

The multivariate generalization of the gradient descent concept: $x_{i+1} = x_i - \eta f'(x_i)$ is given by: $\mathbf{w}_{i+1} = \mathbf{w}_i - \eta \nabla_{\mathbf{w}} L(\mathbf{w}_i)$ (using this example, as we will need it). In this new equation, w is a matrix and we are interested in the gradient of w , w.r.t. the loss function. As promised in the lecture, we will show the derivation of the gradient descent formulas for the L2-norm loss divided by 2.

Model: $y = \mathbf{x}\mathbf{w} + b$

Loss: $L = \frac{1}{2} \sum_i (y_i - t_i)^2$

Update rule: $\mathbf{w}_{i+1} = \mathbf{w}_i - \eta \nabla_{\mathbf{w}} L(\mathbf{w}_i)$

(opt. algorithm) $b_{i+1} = b_i - \eta \nabla_b L(b_i)$

Analogically, we find the update rule for the biases.

Please note that the division by 2 that we performed does not change the nature of the loss.

ANY function that holds the basic property of being higher for worse results and lower for better results can be a loss function.

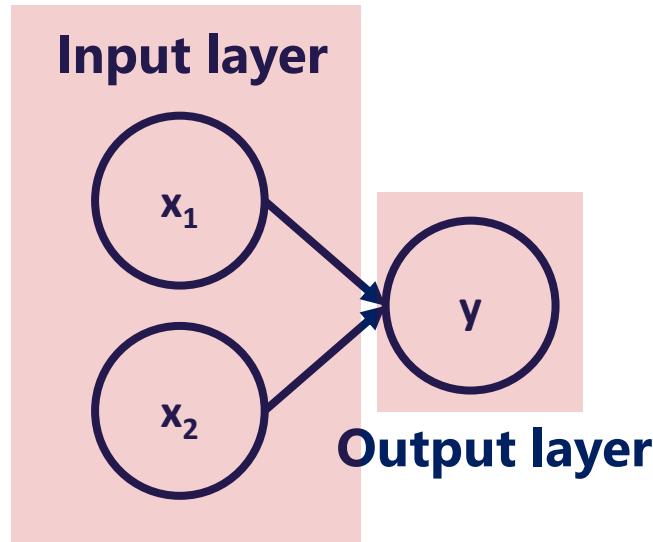
$$\begin{aligned}\nabla_{\mathbf{w}} L &= \nabla_{\mathbf{w}} \frac{1}{2} \sum_i (y_i - t_i)^2 = \\ &= \nabla_{\mathbf{w}} \frac{1}{2} \sum_i ((\mathbf{x}_i \mathbf{w} + b) - t_i)^2 = \\ &= \sum_i \nabla_{\mathbf{w}} \frac{1}{2} (\mathbf{x}_i \mathbf{w} + b - t_i)^2 = \\ &= \sum_i \mathbf{x}_i (\mathbf{x}_i \mathbf{w} + b - t_i) = \\ &= \sum_i \mathbf{x}_i (y_i - t_i) \equiv \\ &\equiv \sum_i \mathbf{x}_i \delta_i\end{aligned}$$

IT'S TIME TO DIG DEEPER

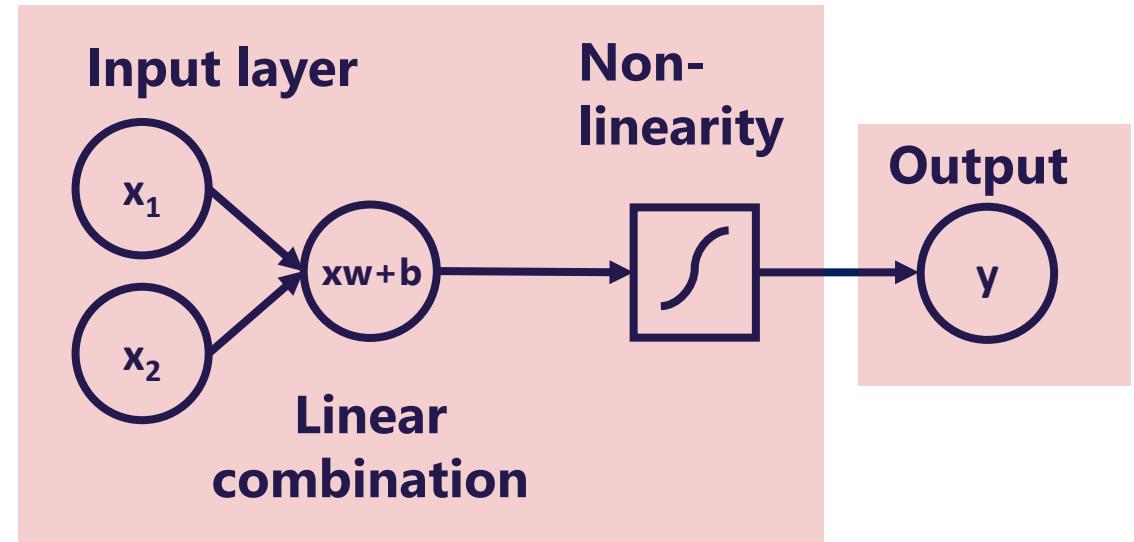
Layers

An initial linear combination and the added non-linearity form a **layer**. The layer is the building block of neural networks.

Minimal example (a simple neural network)



Neural networks



In the minimal example we trained a *neural network* which had no depth. There were solely an input layer and an output layer. Moreover, the output was simply **a linear combination** of the input.

Neural networks step on linear combinations, but add a non-linearity to each one of them. Mixing linear combinations and non-linearities allows us to model arbitrary functions.

A deep net

This is a deep neural network (deep net) with 5 layers.

How to read this diagram:



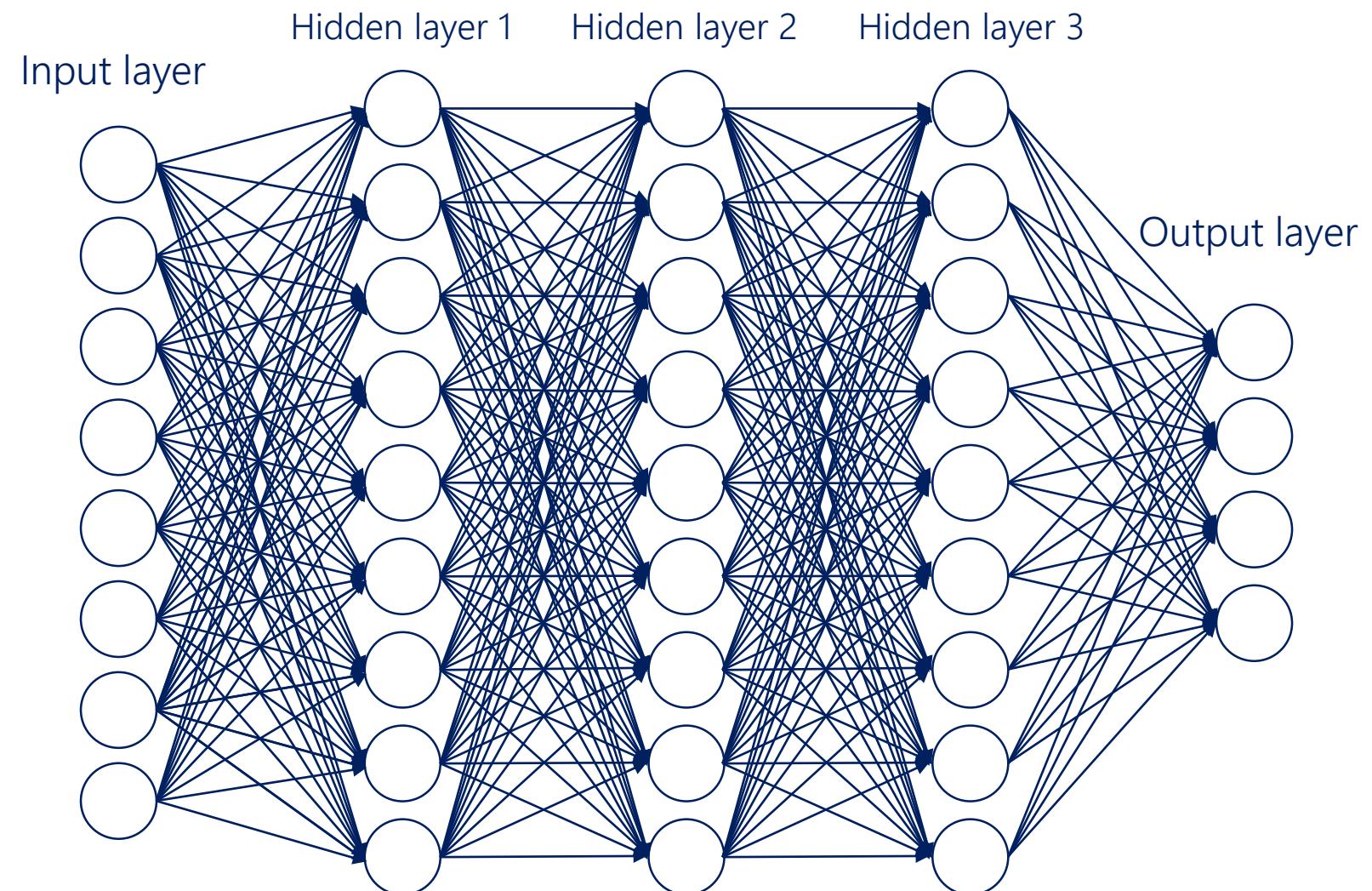
A layer



A unit (a neuron)



Arrows represent
mathematical transformations



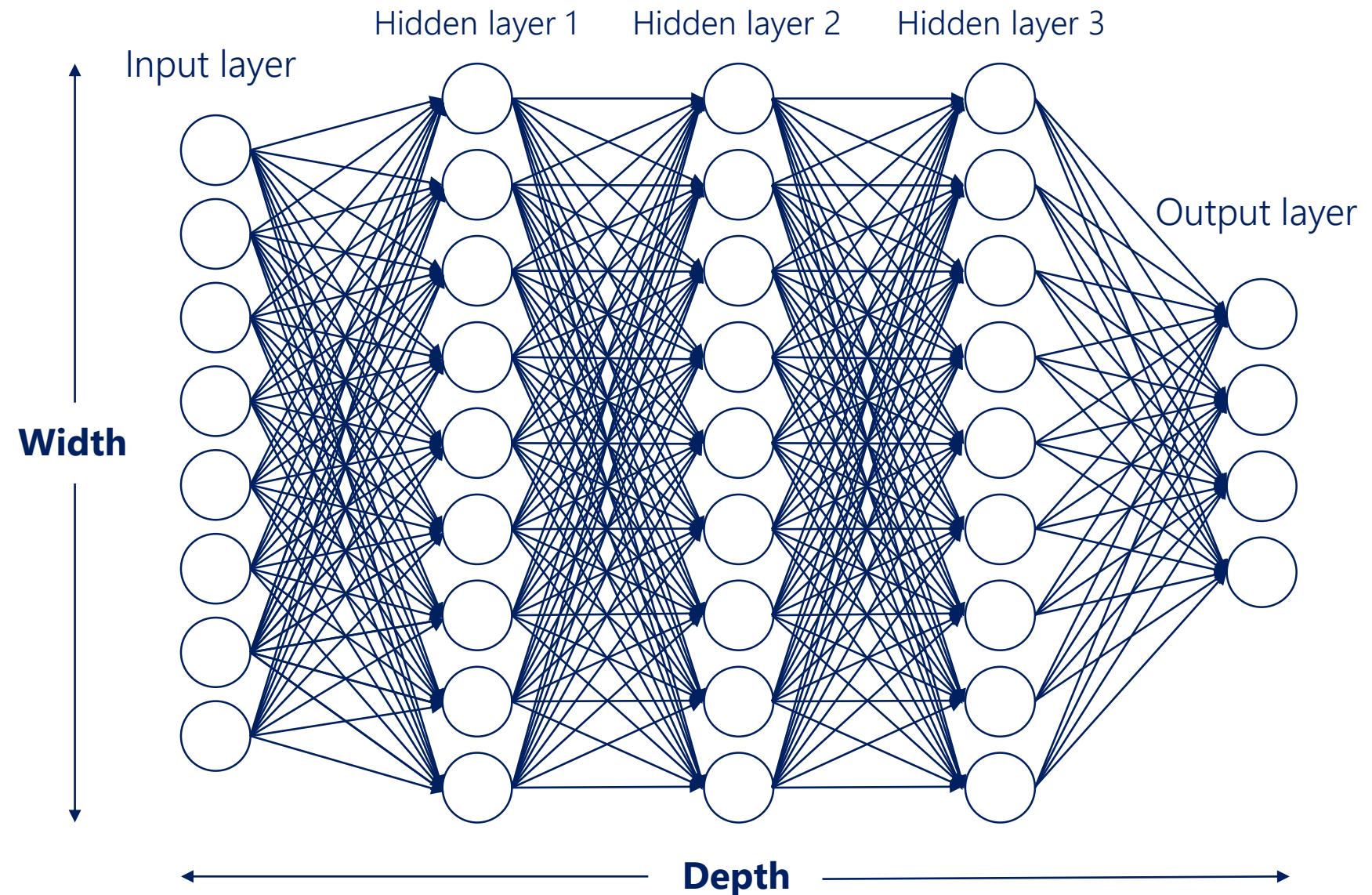
A deep net

The **width** of a layer is the number of units in that layer

The **width** of the net is the number of units of the biggest layer

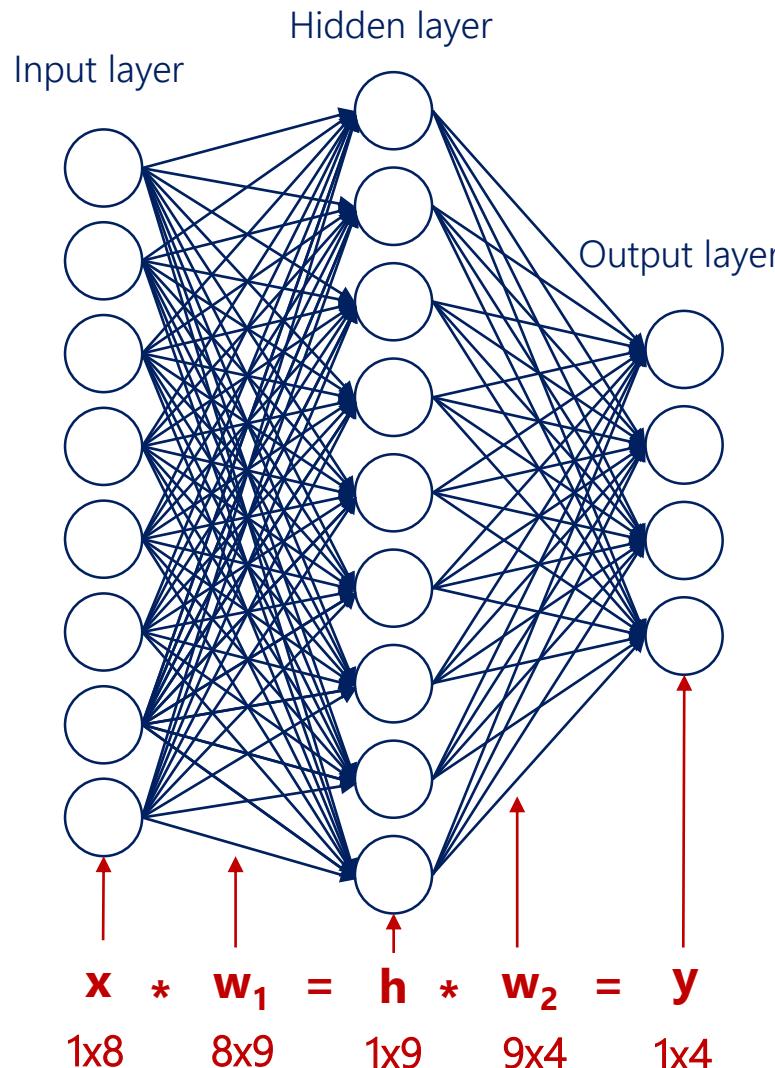
The **depth** of the net is equal to the number of layers or the number of hidden layers. The term has different definitions. More often than not, we are interested in the number of hidden layers (as there are always input and output layers).

The width and the depth of the net are called **hyperparameters**. They are values we manually chose when creating the net.



Why we need non-linearities to stack layers

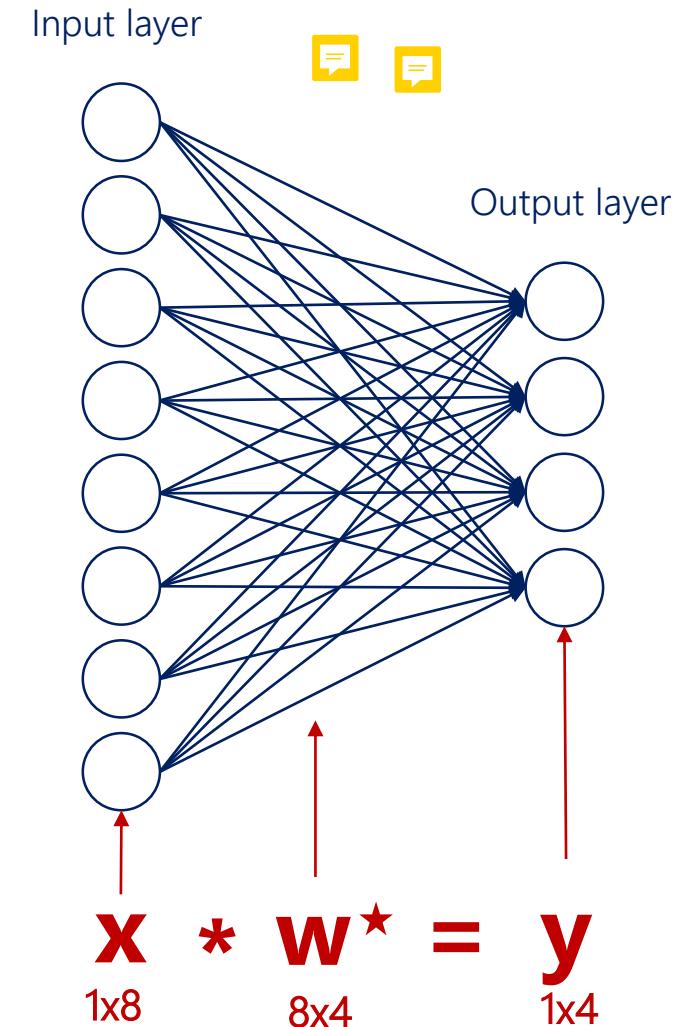
You can see a net with no non-linearities: just linear combinations.



Two consecutive linear transformations are equivalent to a single one.

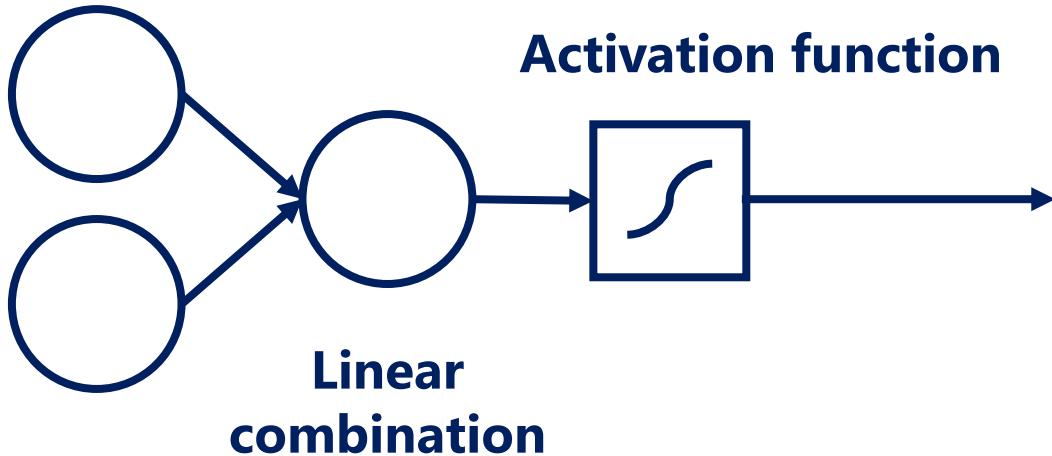
$$\begin{aligned} h &= x * w_1 \\ y &= h * w_2 \\ y &= x * w_1 * w_2 \\ &\quad 8 \times 9 \quad 9 \times 4 \\ y &= x * w^* \\ &\quad 8 \times 4 \end{aligned}$$

Two consecutive linear transformations are equivalent to a single one.



Activation functions

Input



In the respective lesson, we gave an example of temperature change. The temperature starts decreasing (which is a numerical change). Our brain is a kind of an 'activation function'. It tells us whether it is **cold enough** for us to put on a jacket.

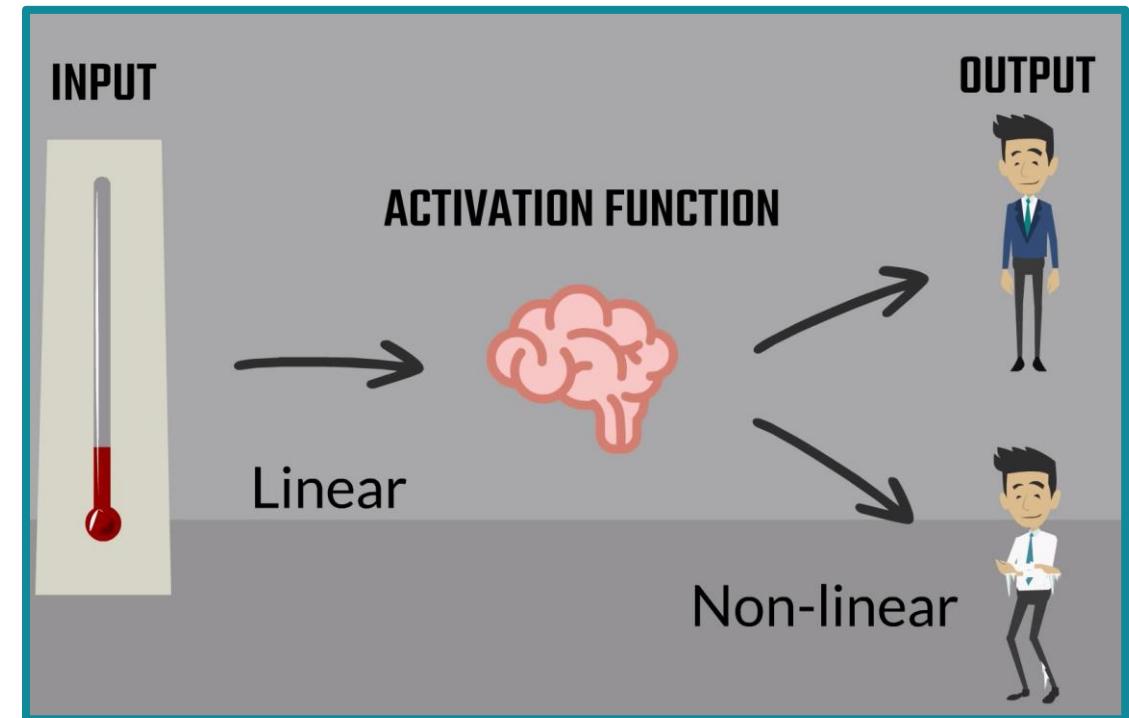
Putting on a jacket is a binary action: 0 (no jacket) or 1 (jacket).

This is a very intuitive and visual (yet not so practical) example of how activation functions work.

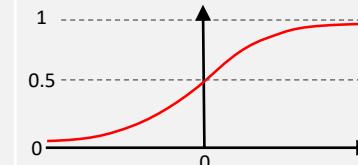
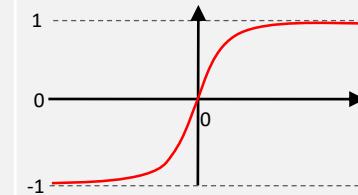
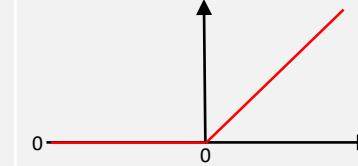
Activation functions (non-linearities) are needed so we can break the linearity and represent more complicated relationships.

Moreover, activation functions are required in order to **stack layers**.

Activation functions transform inputs into outputs of a different kind.



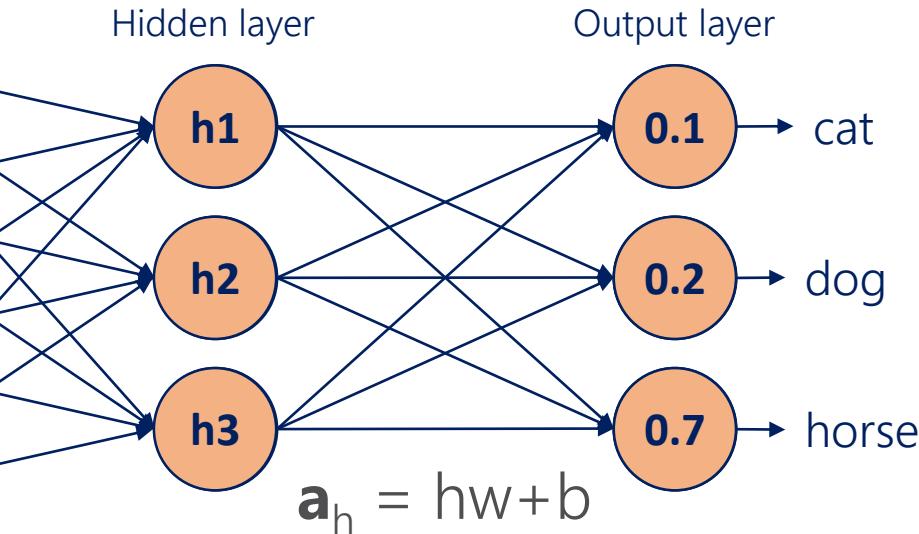
Common activation functions

Name	Formula	Derivative 	Graph	Range
sigmoid (logistic function)	$\sigma(a) = \frac{1}{1+e^{-a}}$	$\frac{\partial \sigma(a)}{\partial a} = \sigma(a)(1 - \sigma(a))$		(0,1)
TanH (hyperbolic tangent)	$\tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}}$	$\frac{\partial \tanh(a)}{\partial a} = \frac{4}{(e^a + e^{-a})^2}$		(-1,1)
ReLu (rectified linear unit)	$\text{relu}(a) = \max(0,a)$	$\frac{\partial \text{relu}(a)}{\partial a} = \begin{cases} 0, & \text{if } a \leq 0 \\ 1, & \text{if } a > 0 \end{cases}$		(0,∞)
softmax	$\sigma_i(\mathbf{a}) = \frac{e^{a_i}}{\sum_j e^{a_j}}$ 	$\frac{\partial \sigma_i(\mathbf{a})}{\partial a_j} = \sigma_i(\mathbf{a}) (\delta_{ij} - \sigma_j(\mathbf{a}))$ Where δ_{ij} is 1 if $i=j$, 0 otherwise		(0,1)

All common activation functions are: **monotonic**, **continuous**, and **differentiable**. These are important properties needed for the optimization.

Softmax activation

Input layer



The softmax activation transforms a bunch of arbitrarily large or small numbers into a valid probability distribution.

While other activation functions get an input value and transform it, regardless of the other elements, the softmax considers the information about the **whole set of numbers** we have.

The values that softmax outputs are in the range from 0 to 1 and their sum is exactly 1 (like probabilities).

Example:

$$\mathbf{a} = [-0.21, 0.47, 1.72]$$

$$\text{softmax}(\mathbf{a}) = \frac{e^{a_i}}{\sum_j e^{a_j}}$$

$$\sum_j e^{a_j} = e^{-0.21} + e^{0.47} + e^{1.72} = 8$$

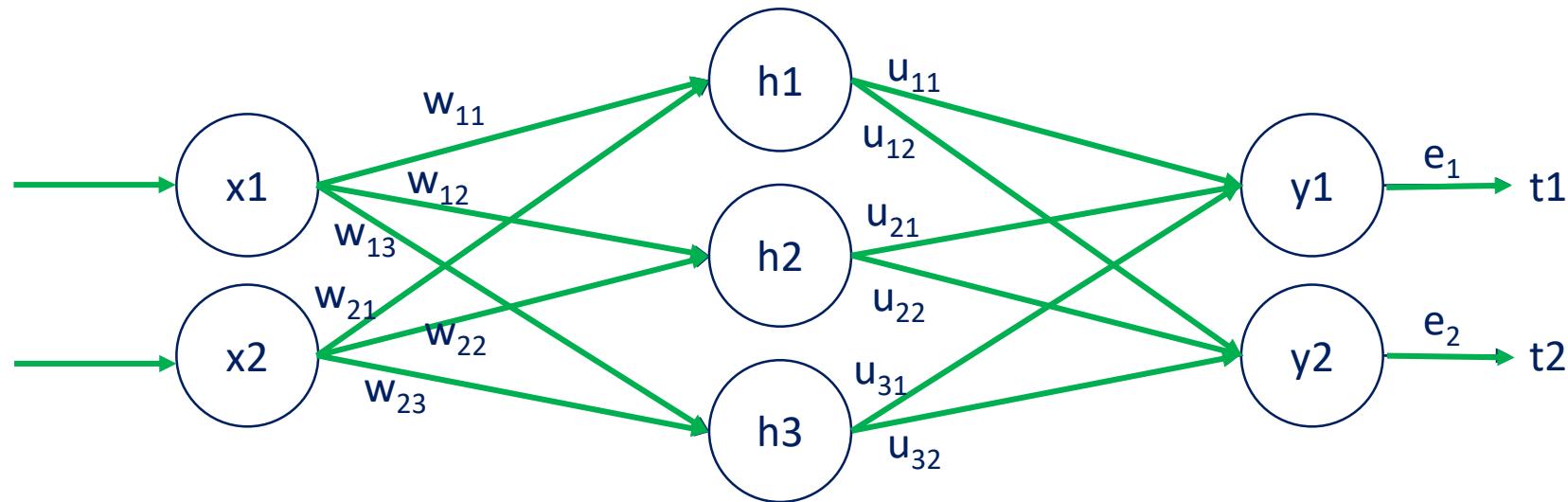
$$\text{softmax}(\mathbf{a}) = \left[\frac{e^{-0.21}}{8}, \frac{e^{0.47}}{8}, \frac{e^{1.72}}{8} \right]$$

$$y = [0.1, 0.2, 0.7] \rightarrow \text{probability distribution}$$

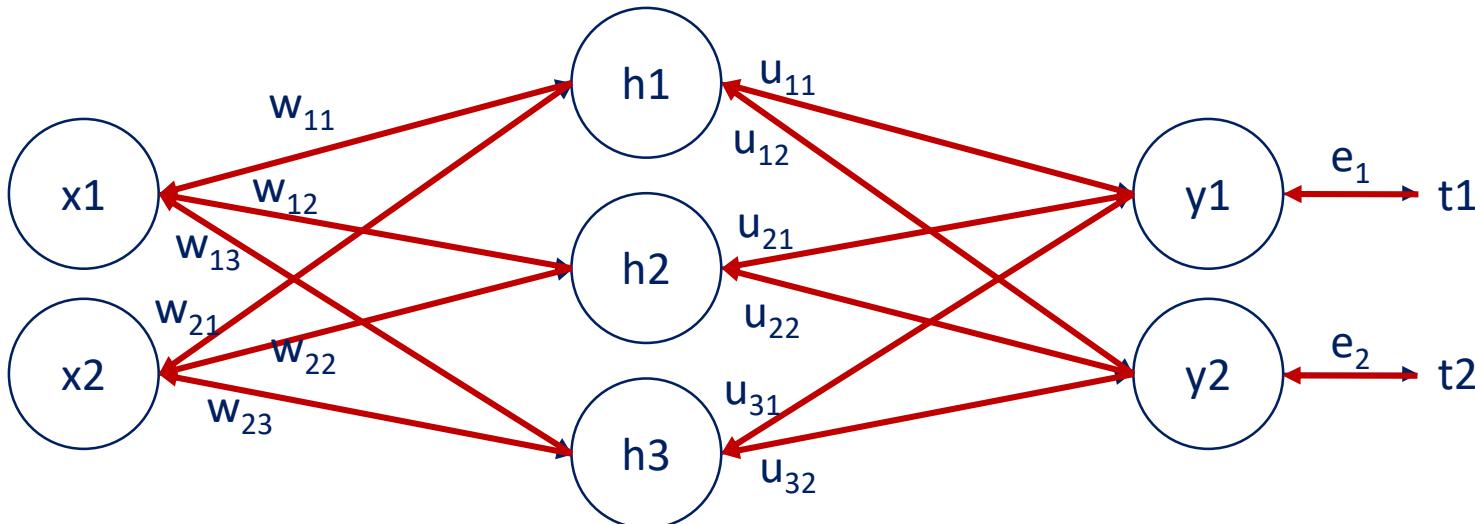
The property of the softmax to output probabilities is so useful and intuitive that it is often used as the activation function for the **final (output) layer**.

However, when the softmax is used prior to that (as the activation of a hidden layer), the results are not as satisfactory. That's because a lot of the information about the variability of the data is lost.

Backpropagation

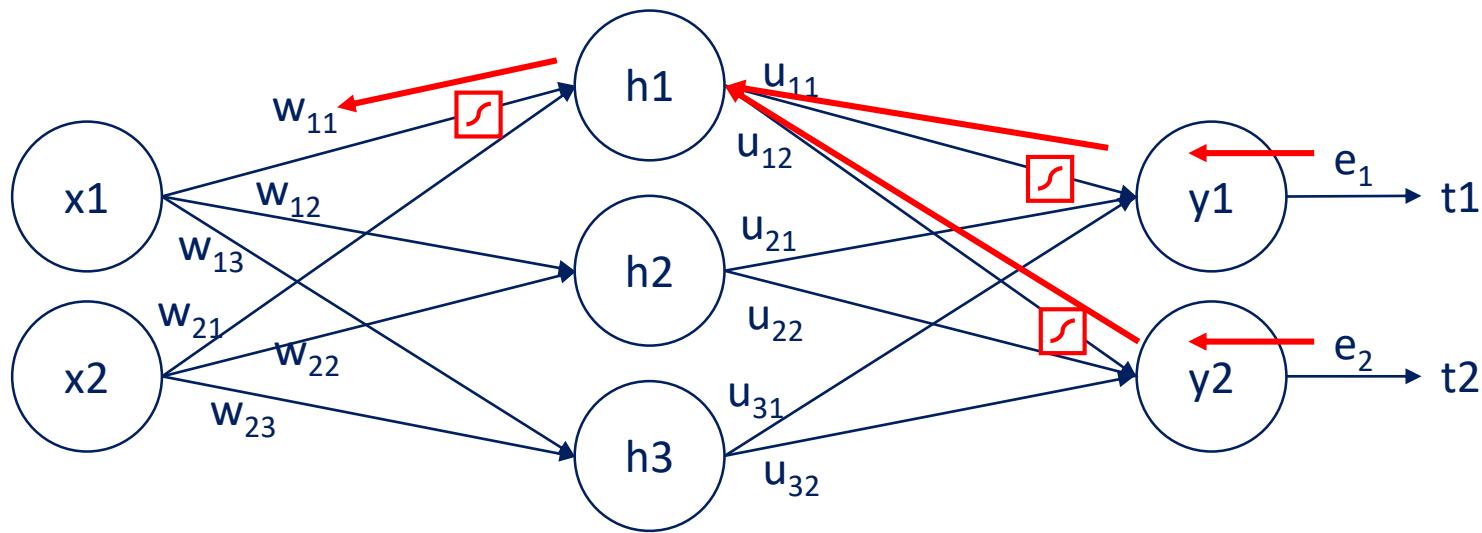


Forward propagation is the process of pushing inputs through the net. At the end of each epoch, the obtained outputs are compared to targets to form the errors.



Backpropagation of errors is an **algorithm** for neural networks using gradient descent. It consists of calculating the contribution of each **parameter** to the errors. We backpropagate the **errors** through the net and **update** the parameters (weights and biases) accordingly.

Backpropagation formula



$$\frac{\partial L}{\partial w_{ij}} = \delta_j x_i, \text{ where } \delta_j = \sum_k \delta_k w_{jk} y_j (1 - y_j)$$

If you want to examine the full derivation, please make use of the PDF we made available in the section: **Backpropagation. A peek into the Mathematics of Optimization.**

Backpropagation. A Peek into the Mathematics of Optimization

1 Motivation

In order to get a truly deep understanding of deep neural networks, one must look at the mathematics of it. As backpropagation is at the core of the optimization process, we wanted to introduce you to it. This is definitely not a necessary part of the course, as in TensorFlow, sk-learn, or any other machine learning package (as opposed to simply NumPy), will have backpropagation methods incorporated.

2 The specific net and notation we will examine

Here's our simple network:

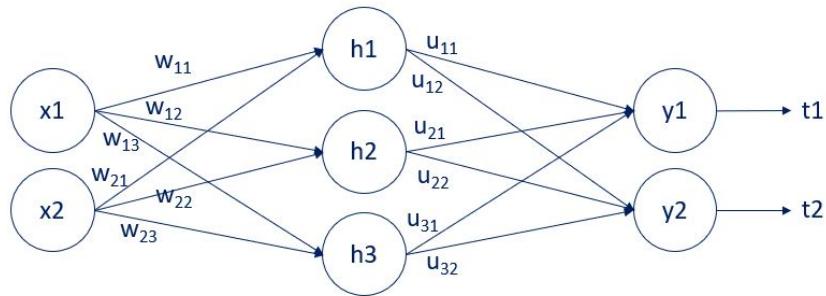


Figure 1: Backpropagation

We have two inputs: x_1 and x_2 . There is a single hidden layer with 3 units (nodes): h_1 , h_2 , and h_3 . Finally, there are two outputs: y_1 and y_2 . The arrows that connect them are the weights. There are two weights matrices: \mathbf{w} , and \mathbf{u} . The \mathbf{w} weights connect the input layer and the hidden layer. The \mathbf{u} weights connect the hidden layer and the output layer. We have employed the letters \mathbf{w} , and \mathbf{u} , so it is easier to follow the computation to follow.

You can also see that we compare the outputs y_1 and y_2 with the targets t_1 and t_2 .

There is one last letter we need to introduce before we can get to the computations. Let a be the linear combination prior to activation. Thus, we have: $\mathbf{a}^{(1)} = \mathbf{x}\mathbf{w} + \mathbf{b}^{(1)}$ and $\mathbf{a}^{(2)} = \mathbf{h}\mathbf{u} + \mathbf{b}^{(2)}$.

Since we cannot exhaust all activation functions and all loss functions, we will focus on two of the most common. A **sigmoid** activation and an **L2-norm loss**.

With this new information and the new notation, the output y is equal to the activated linear combination. Therefore, for the output layer, we have $\mathbf{y} = \sigma(\mathbf{a}^{(2)})$, while for the hidden layer: $\mathbf{h} = \sigma(\mathbf{a}^{(1)})$.

We will examine backpropagation for the output layer and the hidden layer separately, as the methodologies differ.

3 Useful formulas

I would like to remind you that:

$$\text{L2-norm loss: } L = \frac{1}{2} \sum_i (y_i - t_i)^2$$

The sigmoid function is:

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

and its derivative is:

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

4 Backpropagation for the output layer

In order to obtain the update rule:

$$\mathbf{u} \leftarrow \mathbf{u} - \eta \nabla_{\mathbf{u}} L(\mathbf{u})$$

we must calculate

$$\nabla_{\mathbf{u}} L(\mathbf{u})$$

Let's take a single weight u_{ij} . The partial derivative of the loss w.r.t. u_{ij} equals:

$$\frac{\partial L}{\partial u_{ij}} = \frac{\partial L}{\partial y_j} \frac{\partial y_j}{\partial a_j^{(2)}} \frac{\partial a_j^{(2)}}{\partial u_{ij}}$$

where i corresponds to the previous layer (input layer for this transformation) and j corresponds to the next layer (output layer of the transformation). The partial derivatives were computed simply following the chain rule.

$$\frac{\partial L}{\partial y_j} = (y_j - t_j)$$

following the L2-norm loss derivative.

$$\frac{\partial y_j}{\partial a_j^{(2)}} = \sigma(a_j^{(2)})(1 - \sigma(a_j^{(2)})) = y_j(1 - y_j)$$

following the sigmoid derivative.

Finally, the third partial derivative is simply the derivative of $\mathbf{a}^{(2)} = \mathbf{h}\mathbf{u} + \mathbf{b}^{(2)}$. So,

$$\frac{\partial a_j^{(2)}}{\partial u_{ij}} = h_i$$

Replacing the partial derivatives in the expression above, we get:

$$\frac{\partial L}{\partial u_{ij}} = \frac{\partial L}{\partial y_j} \frac{\partial y_j}{\partial a_j^{(2)}} \frac{\partial a_j^{(2)}}{\partial u_{ij}} = (y_j - t_j)y_j(1 - y_j)h_i = \delta_j h_i$$

Therefore, the update rule for a single weight for the output layer is given by:

$$u_{ij} \leftarrow u_{ij} - \eta \delta_j h_i$$

5 Backpropagation of a hidden layer

Similarly to the backpropagation of the output layer, the update rule for a single weight, w_{ij} would depend on:

$$\frac{\partial L}{\partial w_{ij}} = \frac{\partial L}{\partial h_j} \frac{\partial h_j}{\partial a_j^{(1)}} \frac{\partial a_j^{(1)}}{\partial w_{ij}}$$

following the chain rule.

Taking advantage of the results we have so far for transformation using the sigmoid activation and the linear model, we get:

$$\frac{\partial h_j}{\partial a_j^{(1)}} = \sigma(a_j^{(1)})(1 - \sigma(a_j^{(1)})) = h_j(1 - h_j)$$

and

$$\frac{\partial a_j^{(1)}}{\partial w_{ij}} = x_i$$

The actual problem for backpropagation comes from the term $\frac{\partial L}{\partial h_j}$. That's due to the fact that there is no "hidden" target. You can follow the solution for weight w_{11} below. It is advisable to also check Figure 1, while going through the computations.

$$\begin{aligned}\frac{\partial L}{\partial h_1} &= \frac{\partial L}{\partial y_1} \frac{\partial y_1}{\partial a_1^{(2)}} \frac{\partial a_1^{(2)}}{\partial h_1} + \frac{\partial L}{\partial y_2} \frac{\partial y_2}{\partial a_2^{(2)}} \frac{\partial a_2^{(2)}}{\partial h_1} = \\ &= (y_1 - t_1)y_1(1 - y_1)u_{11} + (y_2 - t_2)y_2(1 - y_2)u_{12}\end{aligned}$$

From here, we can calculate $\frac{\partial L}{\partial w_{11}}$, which was what we wanted. The final expression is:

$$\frac{\partial L}{\partial w_{11}} = [(y_1 - t_1)y_1(1 - y_1)u_{11} + (y_2 - t_2)y_2(1 - y_2)u_{12}] h_1(1 - h_1)x_1$$

The generalized form of this equation is:

$$\frac{\partial L}{\partial w_{ij}} = \sum_k (y_k - t_k)y_k(1 - y_k)u_{jk}h_j(1 - h_j)x_i$$

6 Backpropagation generalization

Using the results for backpropagation for the output layer and the hidden layer, we can put them together in one formula, summarizing backpropagation, in the presence of L2-norm loss and sigmoid activations.

$$\frac{\partial L}{\partial w_{ij}} = \delta_j x_i$$

where for a hidden layer

$$\delta_j = \sum_k \delta_k w_{jk} y_j (1 - y_j)$$

Kudos to those of you who got to the end.

Thanks for reading.

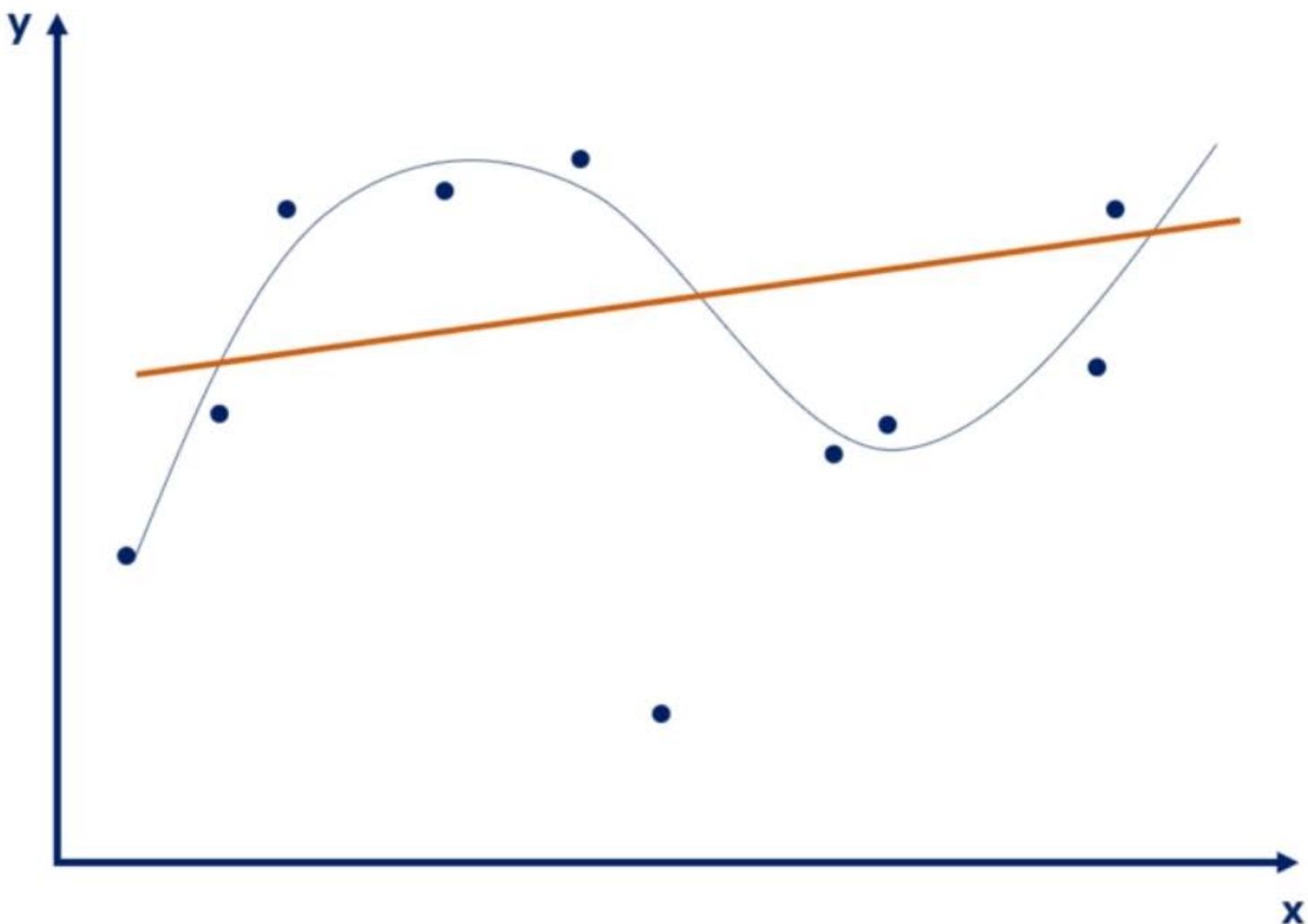
UNDERFITTING

The model has not captured the underlying logic of the data

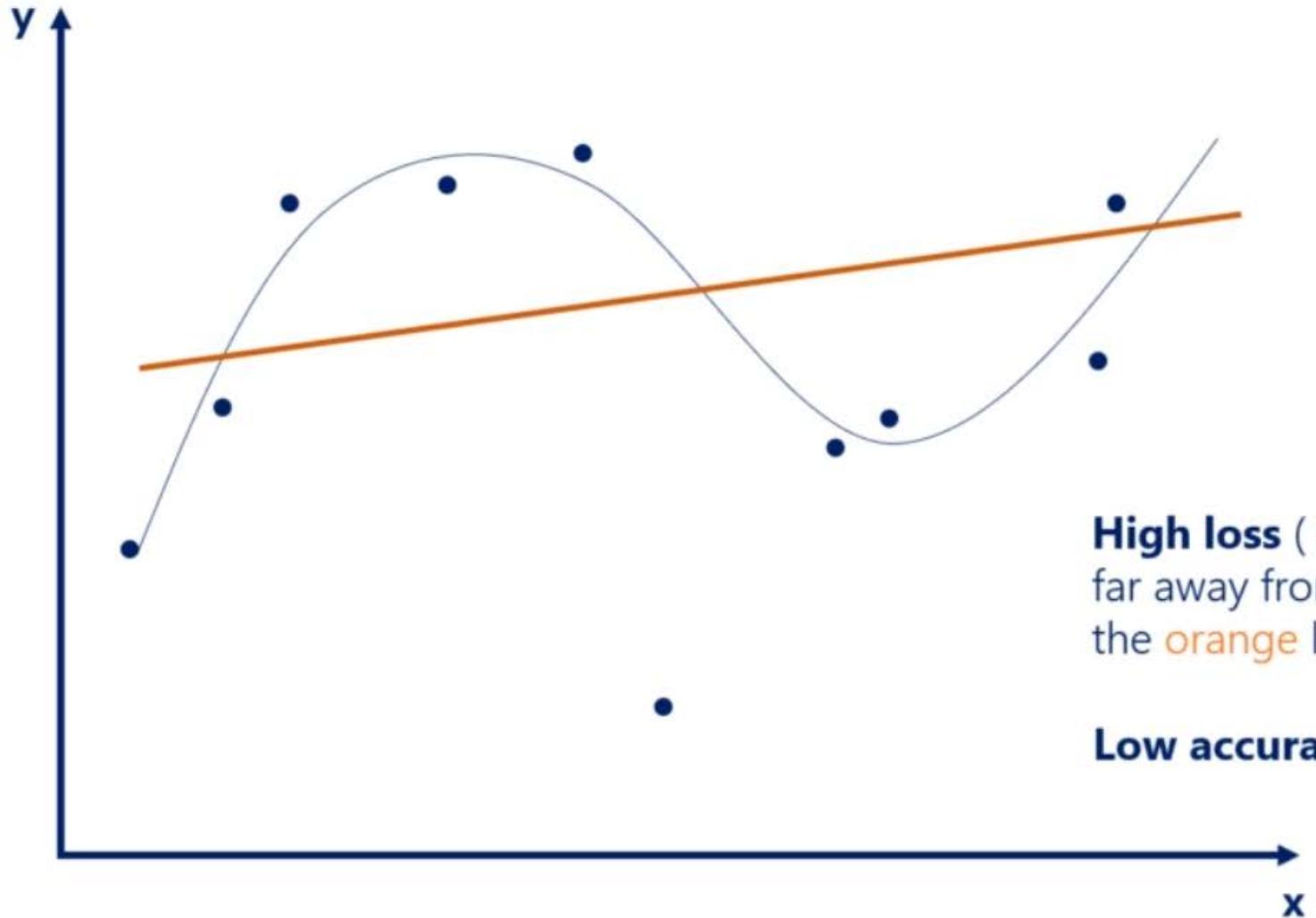
OVERTFITTING

Our training has focused on the particular training set so much, it has "missed the point"

Underfitting and overfitting



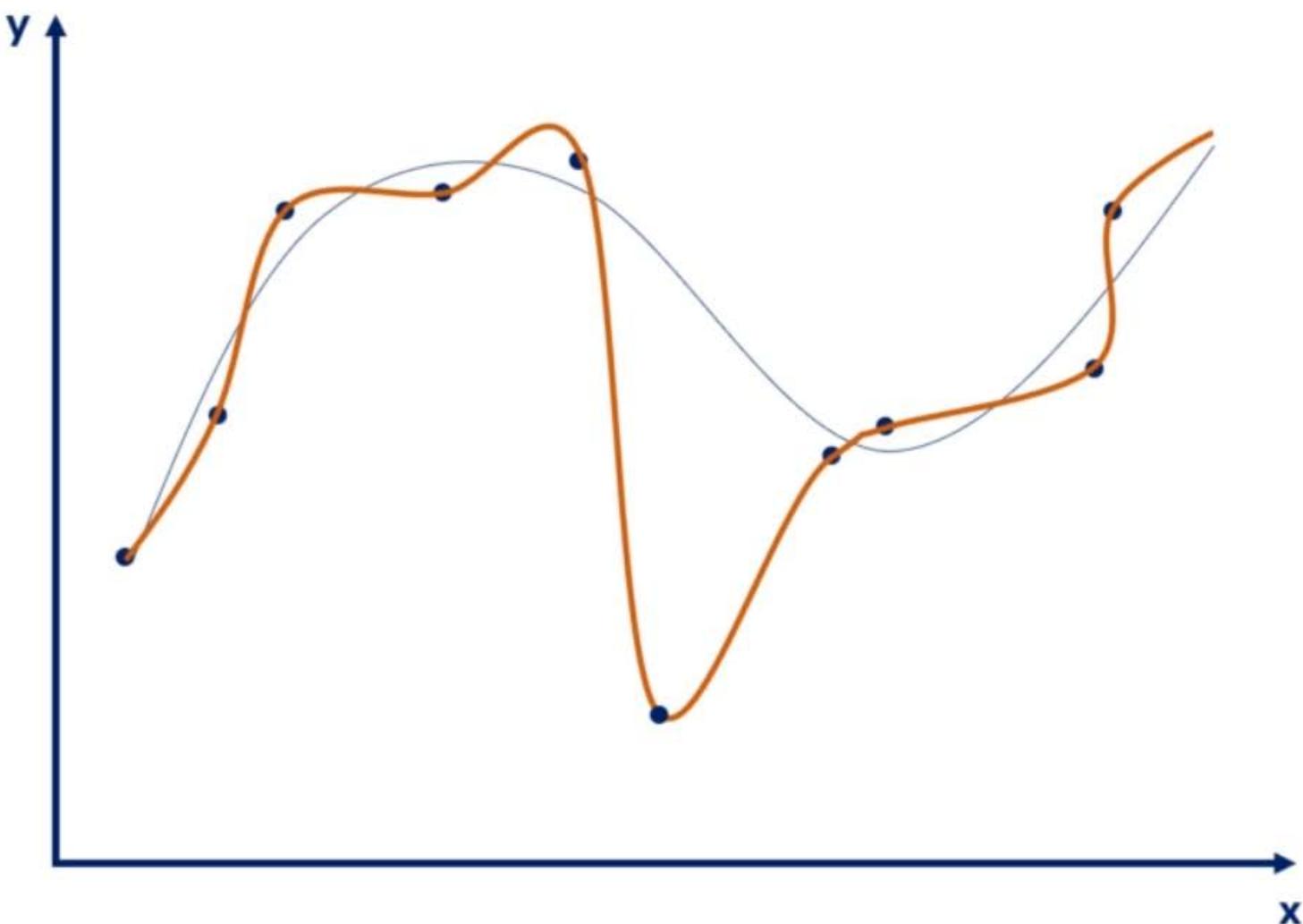
Underfitted model: provides an answer, but does not capture the underlying logic



High loss (targets, or points, are far away from the predictions, y , or the **orange** line)

Low accuracy

Underfitting and overfitting



Overfitted model: so super good at modeling the training data that it "misses the point"

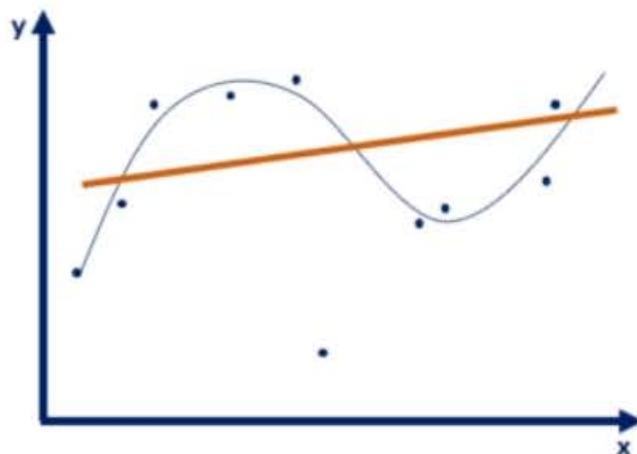
FIRST RULE OF PROGRAMMING



THE COMPUTER IS NEVER WRONG.
IT IS US, WHO MADE A MISTAKE

Underfitting and overfitting

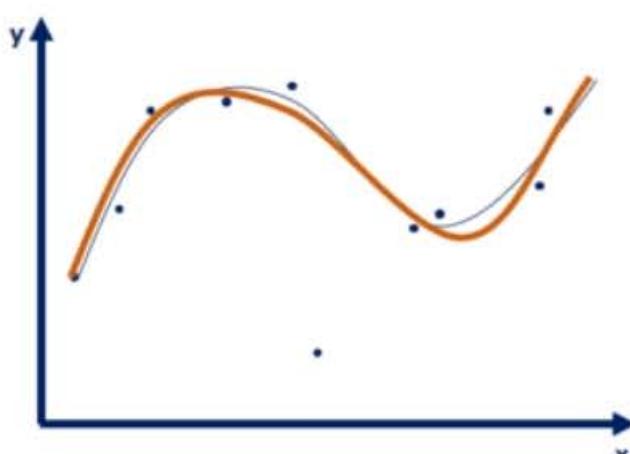
An **underfitted** model



Doesn't capture any logic

- High loss
- Low accuracy

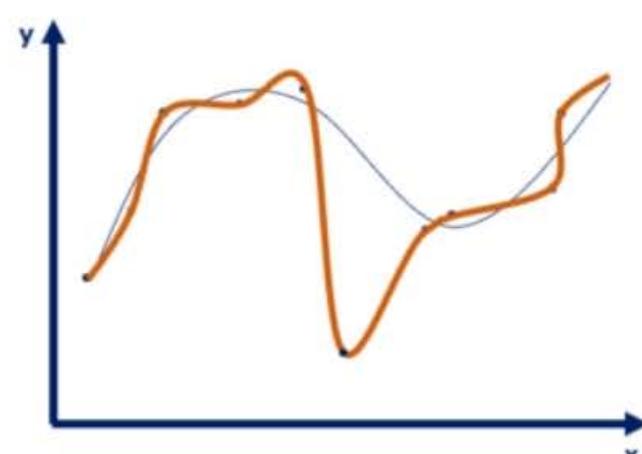
A **good** model



Captures the underlying logic of the dataset

- Low loss
- High accuracy

An **overfitted** model

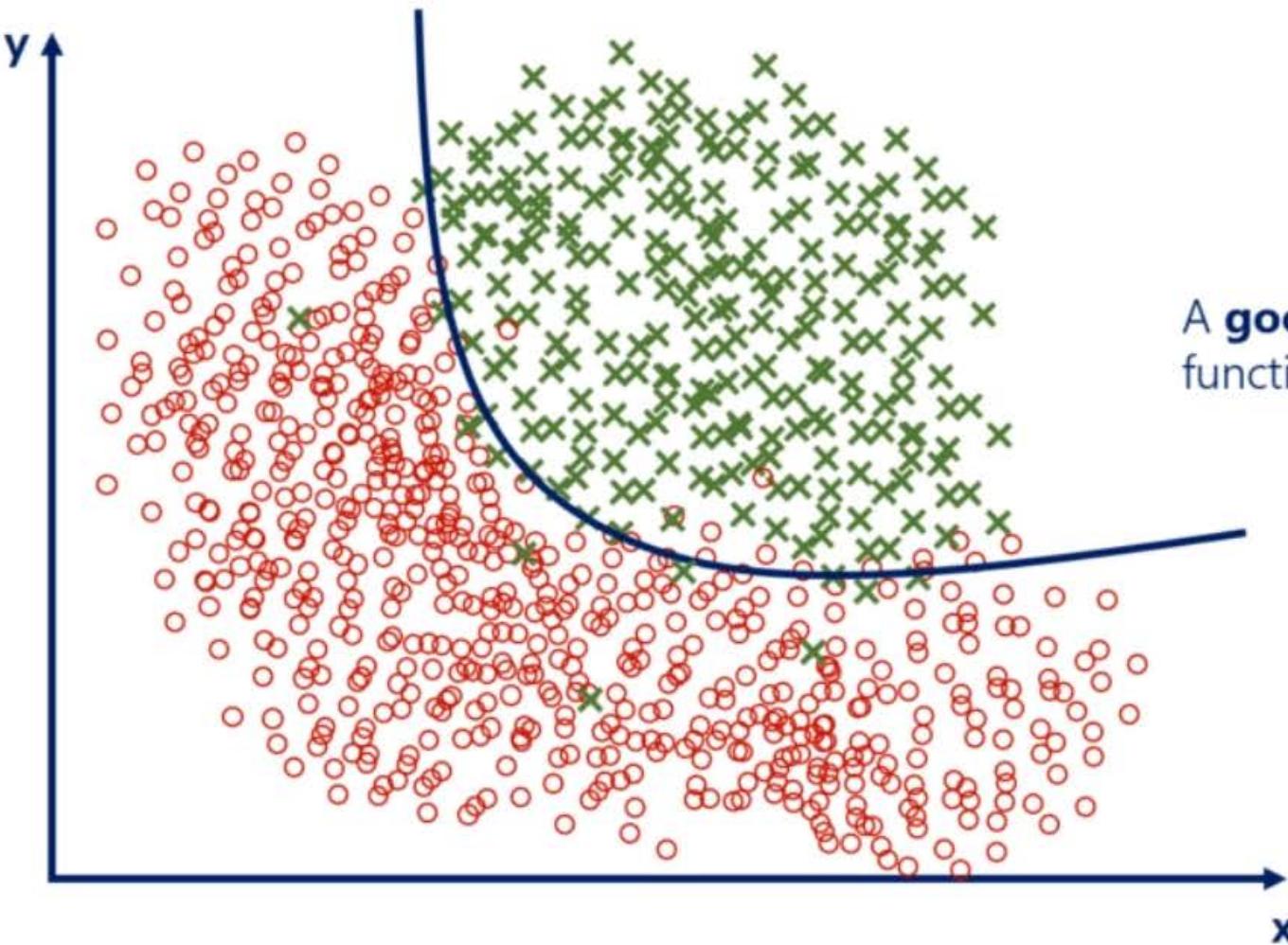


Captures all the noise, thus "missed the point"

- Low loss
- Low accuracy

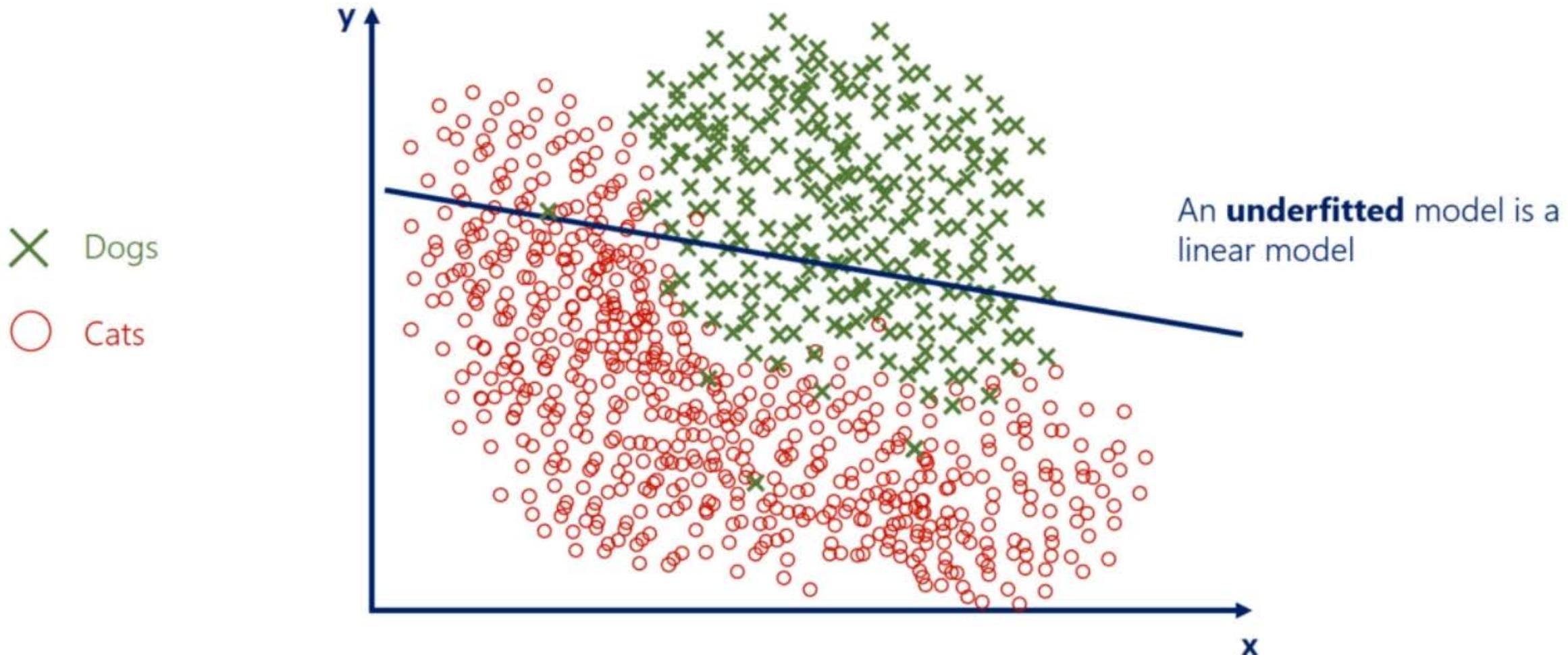
Underfitting and overfitting. A classification example

- ✖ Dogs
- Cats

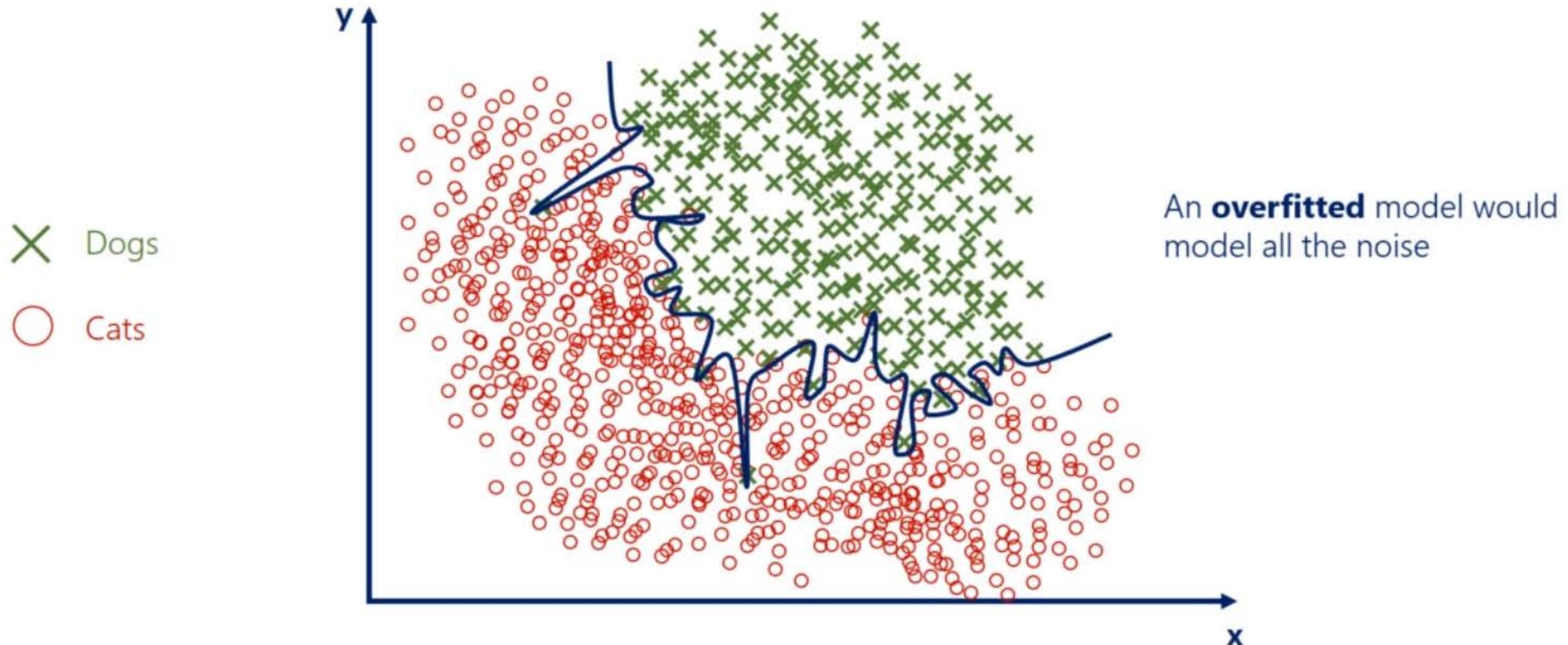


A **good** model is a quadratic function with a few errors

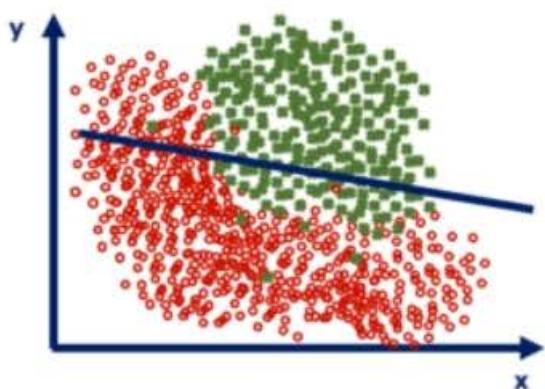
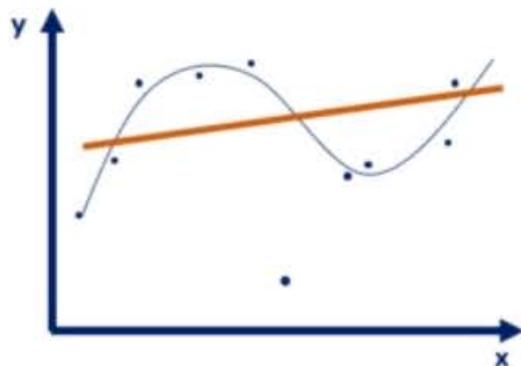
Underfitting and overfitting. A classification example



Underfitting and overfitting. A classification example



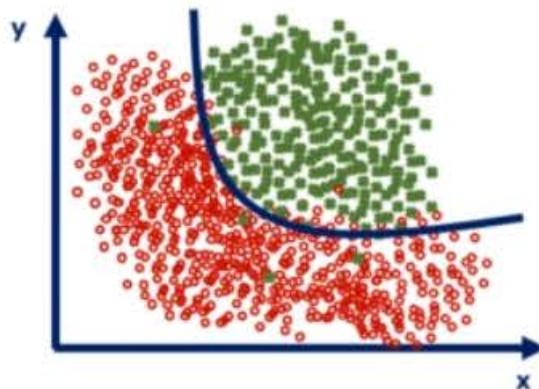
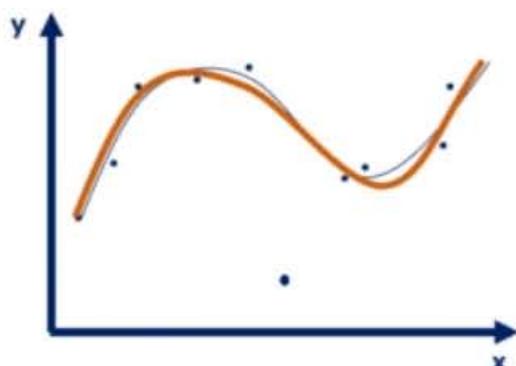
An **underfitted** model



Doesn't capture any logic

- High loss
- Low accuracy

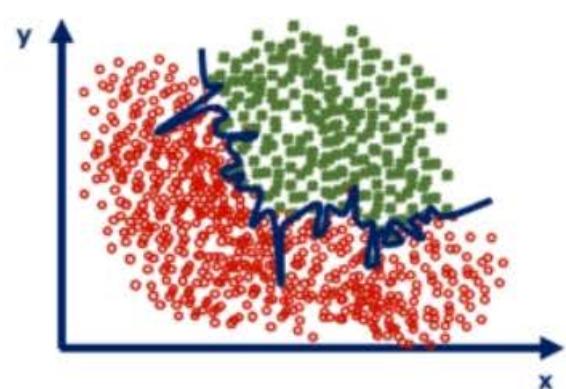
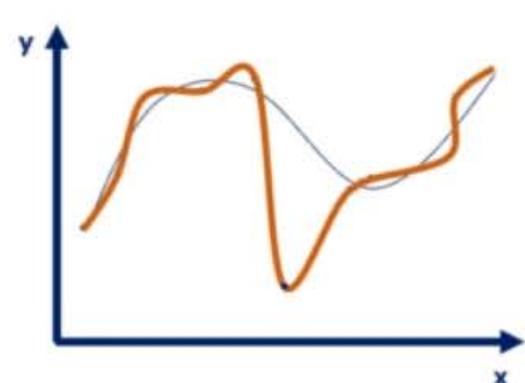
A **good** model



Captures the underlying logic
of the dataset

- Low loss
- High accuracy

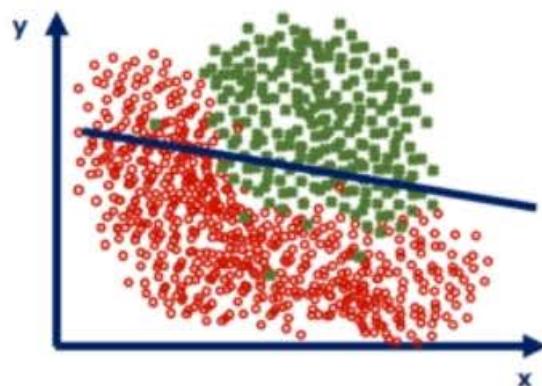
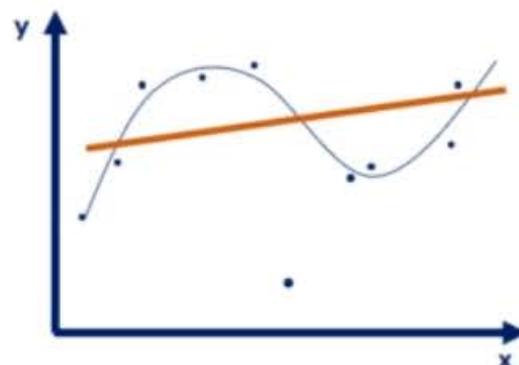
An **overfitted** model



Captures all the noise, thus
"missed the point"

- Low loss
- Low accuracy

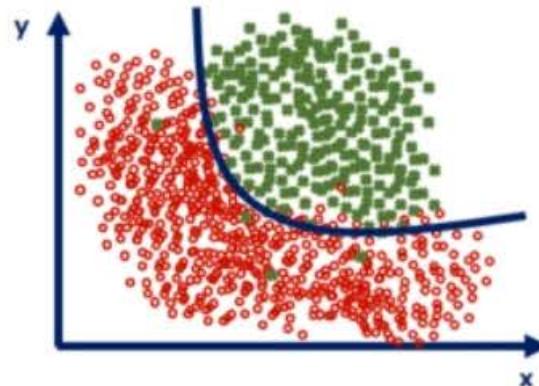
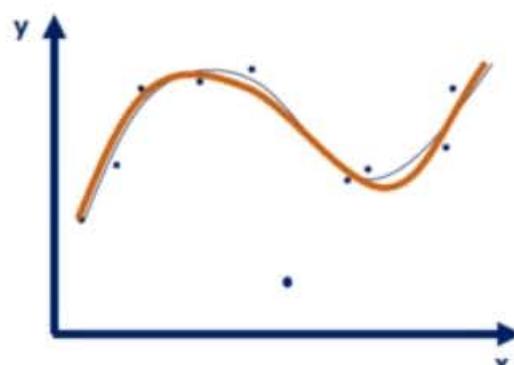
An **underfitted** model



Doesn't capture any logic

- High loss
- Low accuracy

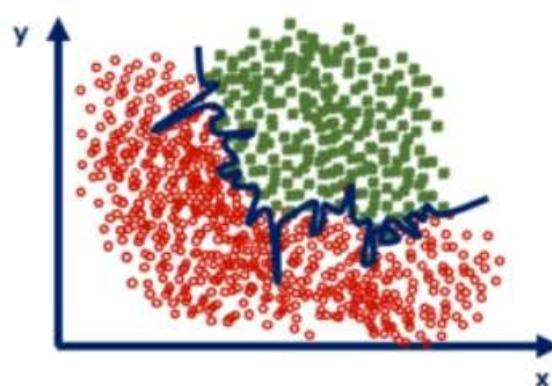
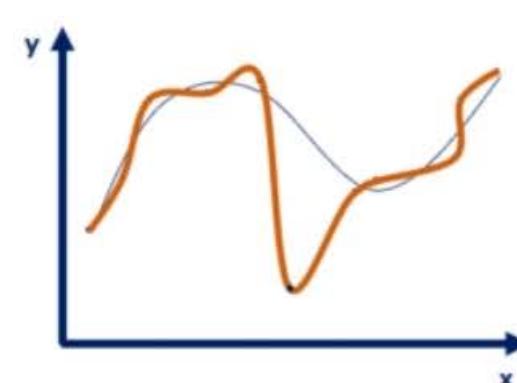
A **good** model



Captures the underlying logic
of the dataset

- Low loss
- High accuracy

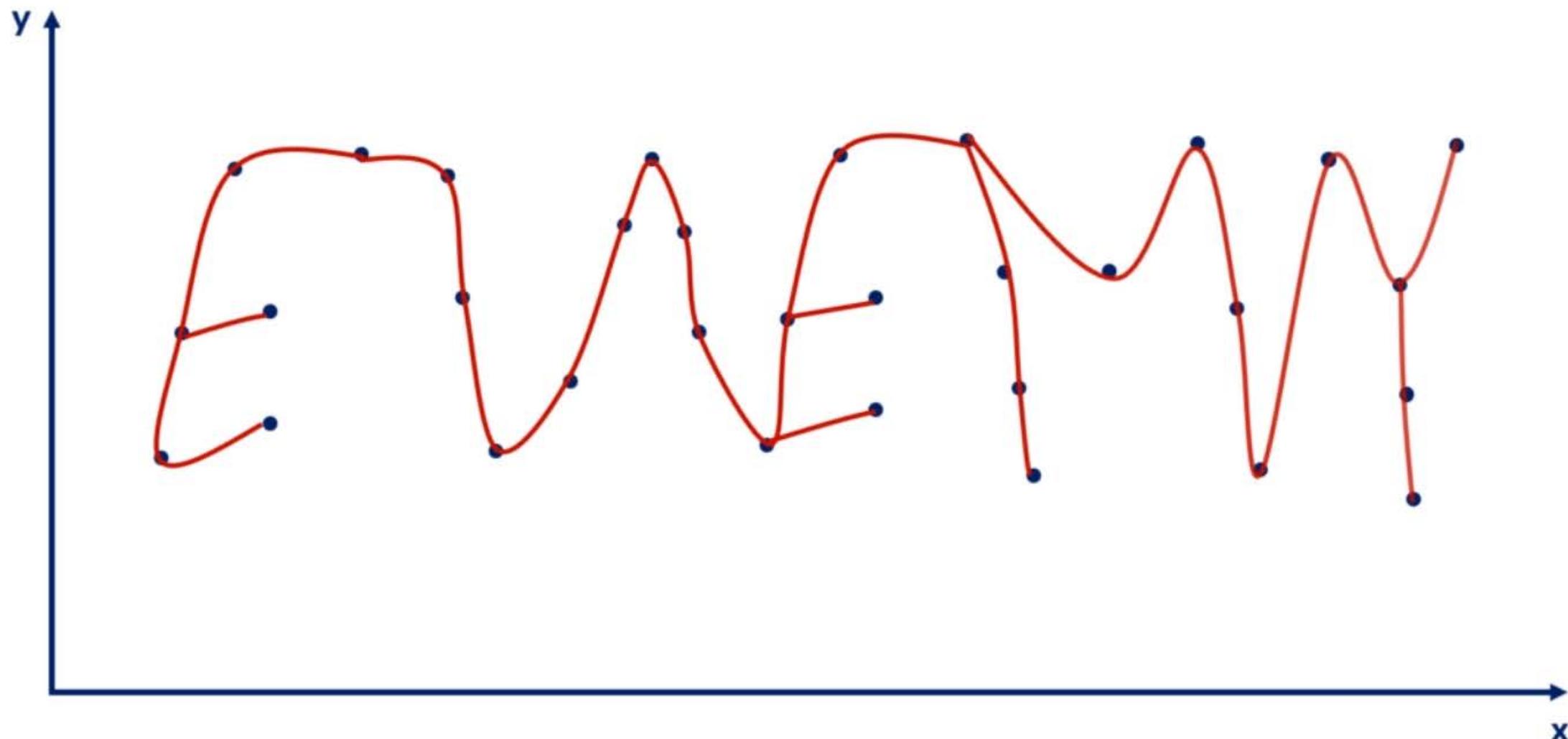
An **overfitted** model



Captures all the noise, thus
"missed the point"

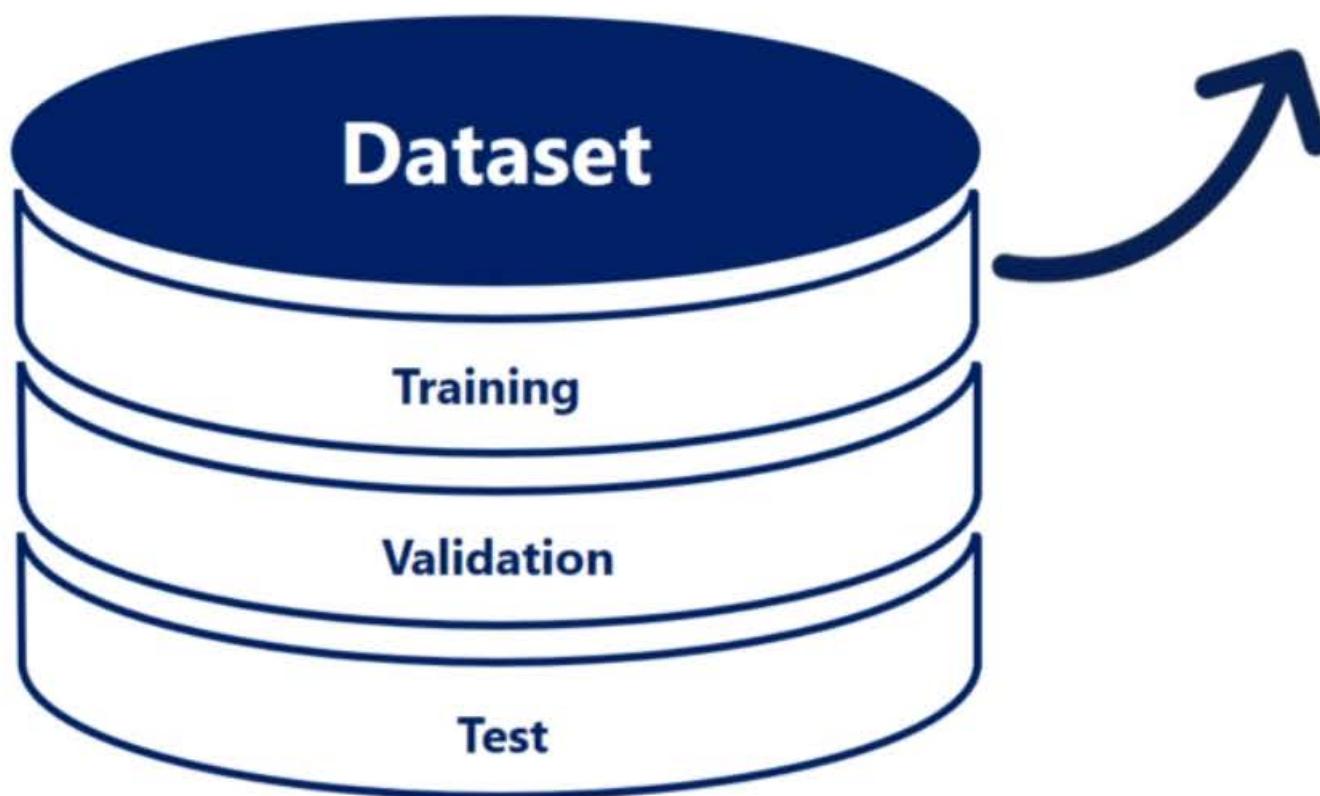
- Low loss
- Low accuracy

Overfitted model



Overfitting is the real enemy when it comes to machine learning

Training, validation, and test

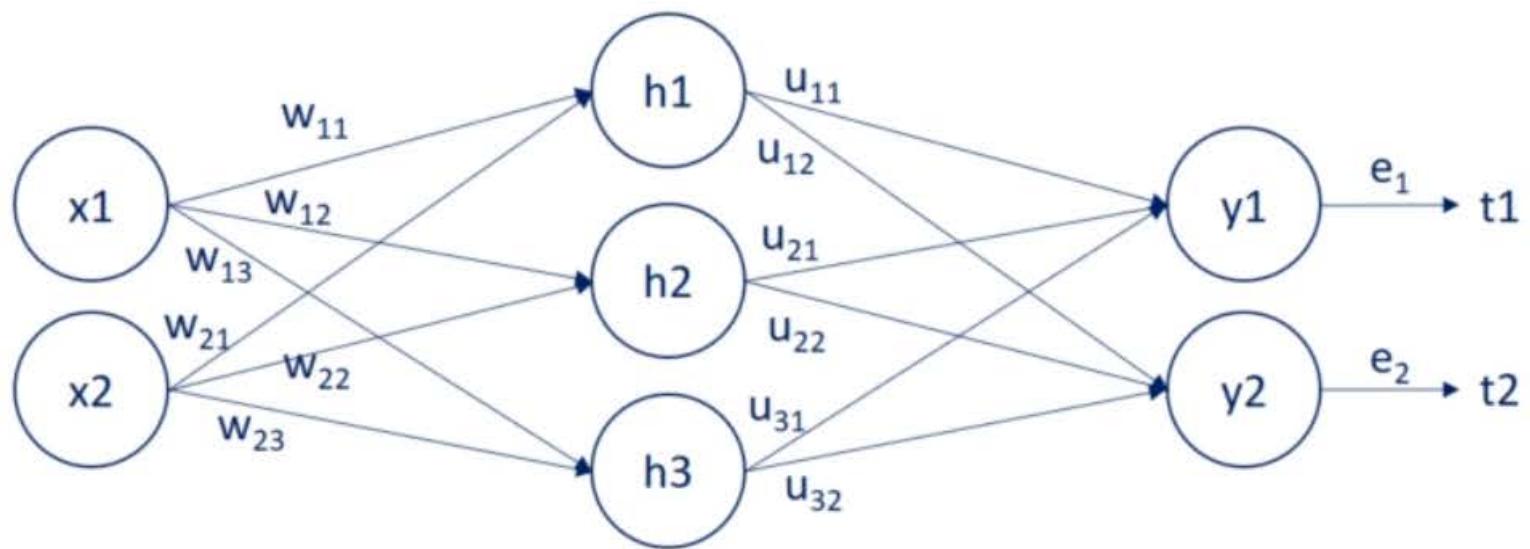


That's where the training happens

```
training_data = np.load('TF_intro.npz')
```



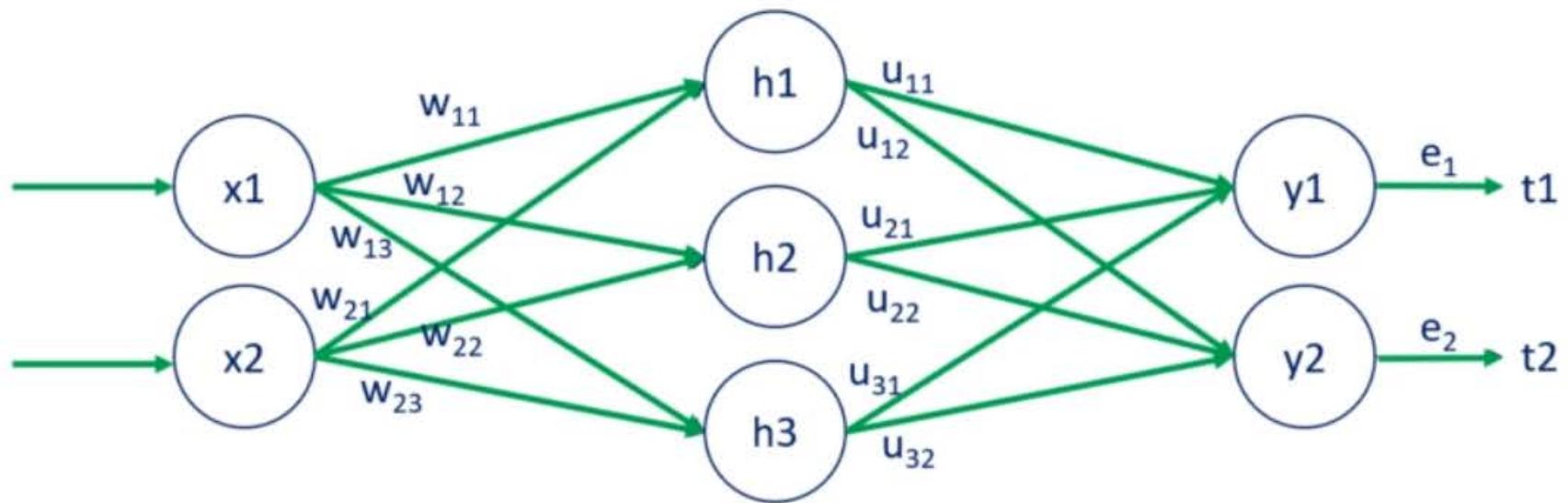
Training



We update the weights and biases for the training set **only**



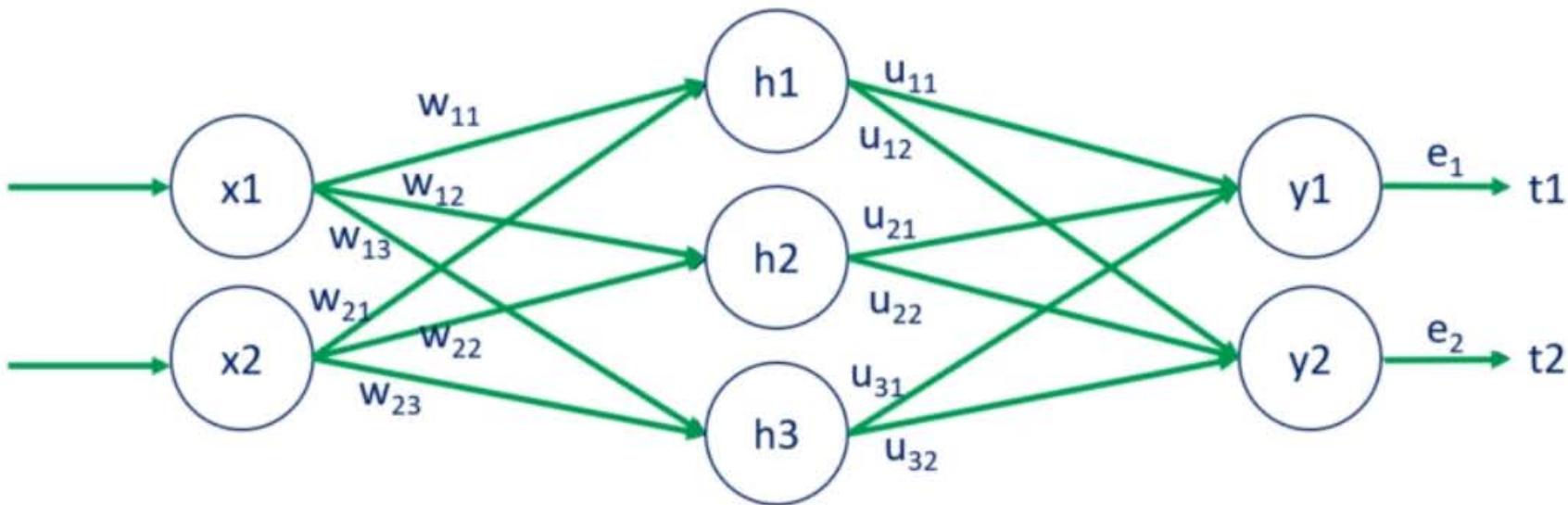
Validation



We run the model on the **validation dataset**, we only propagate **forward**



Validation



We just calculate the **validation loss**. On average it should equal the **training loss**



Training and validation



training_loss



validation_loss

```
Epoch 1. Training loss: 0.341. Validation loss: 0.192
Epoch 2. Training loss: 0.164. Validation loss: 0.143
Epoch 3. Training loss: 0.115. Validation loss: 0.106
Epoch 4. Training loss: 0.086. Validation loss: 0.093
Epoch 5. Training loss: 0.068. Validation loss: 0.087
Epoch 6. Training loss: 0.054. Validation loss: 0.084
Epoch 7. Training loss: 0.042. Validation loss: 0.074
Epoch 8. Training loss: 0.033. Validation loss: 0.077
Epoch 9. Training loss: 0.029. Validation loss: 0.090
Epoch 10. Training loss: 0.022. Validation loss: 0.072
Epoch 11. Training loss: 0.018. Validation loss: 0.074
Epoch 12. Training loss: 0.014. Validation loss: 0.075
Epoch 13. Training loss: 0.011. Validation loss: 0.081
Epoch 14. Training loss: 0.009. Validation loss: 0.085
Epoch 15. Training loss: 0.007. Validation loss: 0.085
Epoch 16. Training loss: 0.009. Validation loss: 0.091
Epoch 17. Training loss: 0.008. Validation loss: 0.099
Epoch 18. Training loss: 0.006. Validation loss: 0.100
Epoch 19. Training loss: 0.005. Validation loss: 0.092
Epoch 20. Training loss: 0.002. Validation loss: 0.084
Epoch 21. Training loss: 0.002. Validation loss: 0.092
Epoch 22. Training loss: 0.016. Validation loss: 0.085
Epoch 23. Training loss: 0.005. Validation loss: 0.094
Epoch 24. Training loss: 0.001. Validation loss: 0.085
Epoch 25. Training loss: 0.001. Validation loss: 0.088
Epoch 26. Training loss: 0.000. Validation loss: 0.087
Epoch 27. Training loss: 0.000. Validation loss: 0.089
Epoch 28. Training loss: 0.000. Validation loss: 0.091
Epoch 29. Training loss: 0.000. Validation loss: 0.092
```

Each subsequent training loss will be lower or equal to the previous one



Training and validation



training_loss

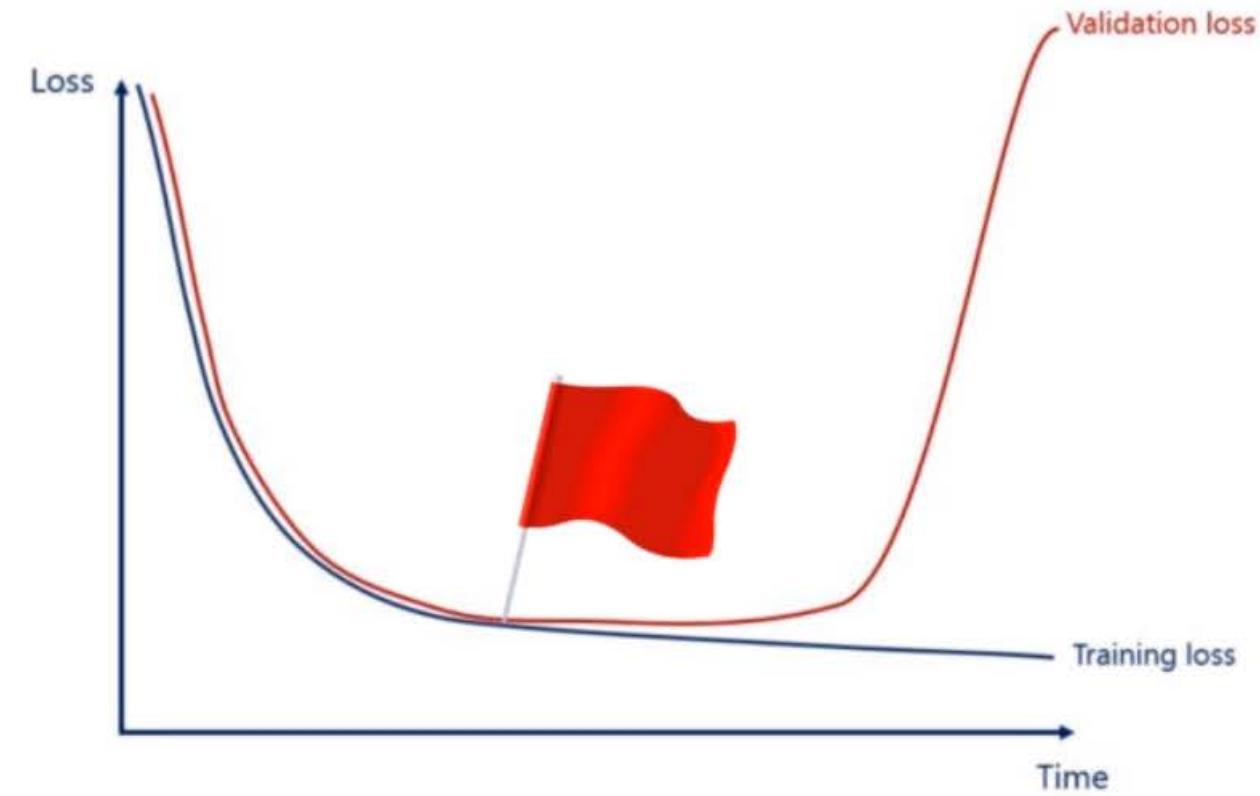
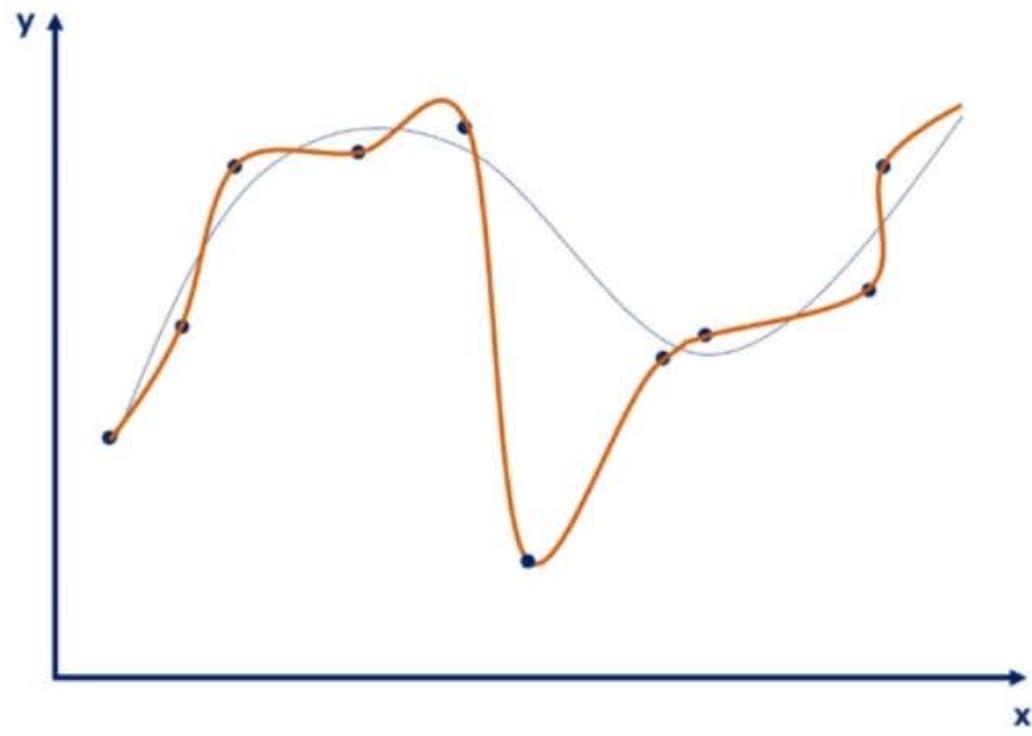


validation_loss



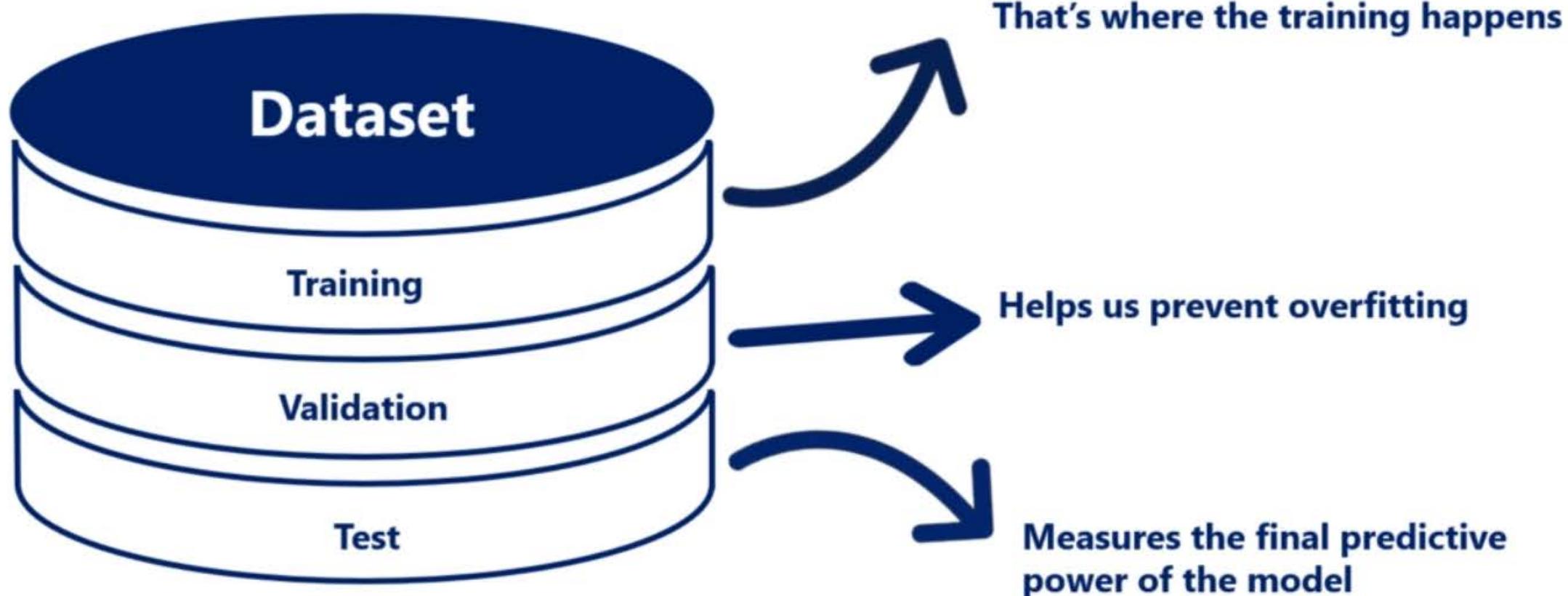
```
Epoch 1. Training loss: 0.341. Validation loss: 0.192
Epoch 2. Training loss: 0.164. Validation loss: 0.143
Epoch 3. Training loss: 0.115. Validation loss: 0.106
Epoch 4. Training loss: 0.086. Validation loss: 0.093
Epoch 5. Training loss: 0.068. Validation loss: 0.087
Epoch 6. Training loss: 0.054. Validation loss: 0.084
Epoch 7. Training loss: 0.042. Validation loss: 0.074
Epoch 8. Training loss: 0.033. Validation loss: 0.077
Epoch 9. Training loss: 0.029. Validation loss: 0.090
Epoch 10. Training loss: 0.022. Validation loss: 0.072
Epoch 11. Training loss: 0.018. Validation loss: 0.074
Epoch 12. Training loss: 0.014. Validation loss: 0.075
Epoch 13. Training loss: 0.011. Validation loss: 0.081
Epoch 14. Training loss: 0.009. Validation loss: 0.085
Epoch 15. Training loss: 0.007. Validation loss: 0.085
Epoch 16. Training loss: 0.009. Validation loss: 0.091
Epoch 17. Training loss: 0.008. Validation loss: 0.099
Epoch 18. Training loss: 0.006. Validation loss: 0.100
Epoch 19. Training loss: 0.005. Validation loss: 0.092
Epoch 20. Training loss: 0.002. Validation loss: 0.084
Epoch 21. Training loss: 0.002. Validation loss: 0.092
Epoch 22. Training loss: 0.016. Validation loss: 0.085
Epoch 23. Training loss: 0.005. Validation loss: 0.094
Epoch 24. Training loss: 0.001. Validation loss: 0.085
Epoch 25. Training loss: 0.001. Validation loss: 0.088
Epoch 26. Training loss: 0.000. Validation loss: 0.087
Epoch 27. Training loss: 0.000. Validation loss: 0.089
Epoch 28. Training loss: 0.000. Validation loss: 0.091
Epoch 29. Training loss: 0.000. Validation loss: 0.092
Epoch 30. Training loss: 0.000. Validation loss: 0.092
```

WE ARE OVERRFITTING!!!

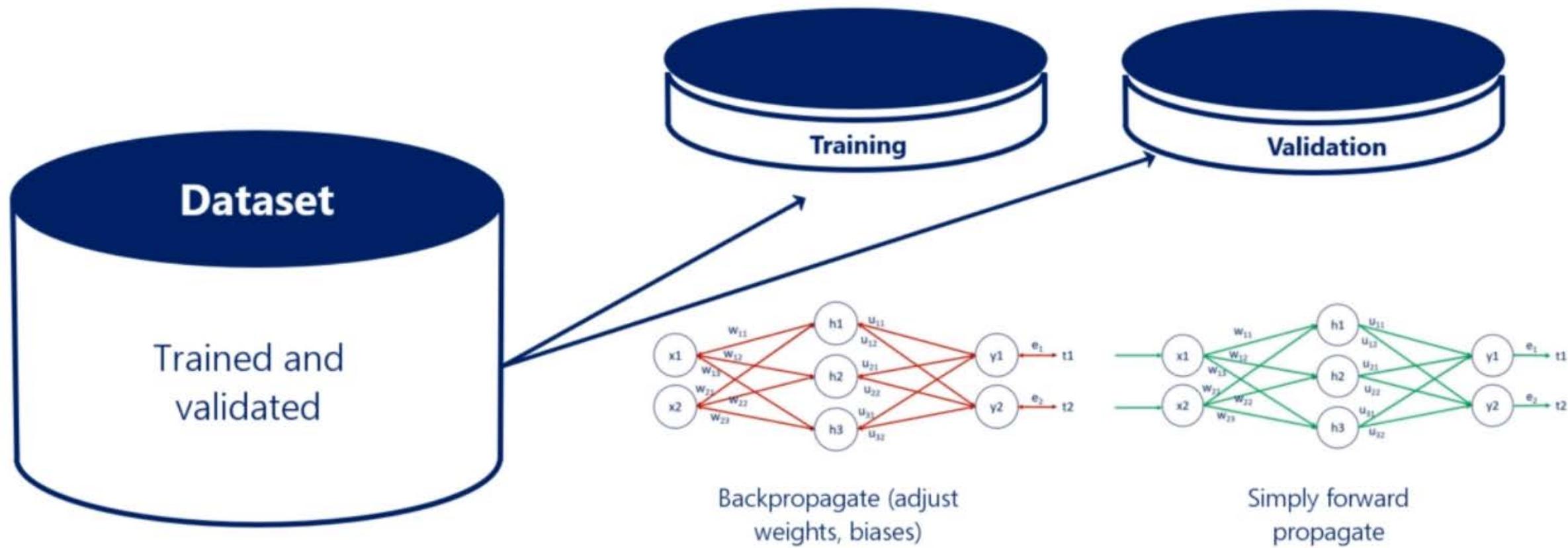


WE ARE OVERRFITTING!!!

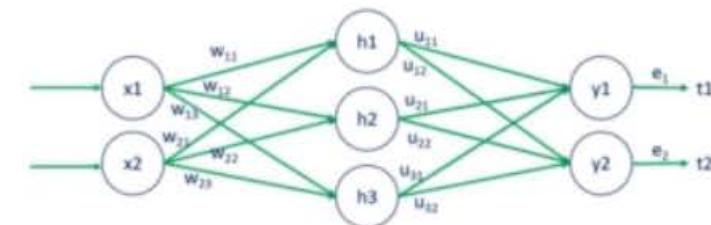
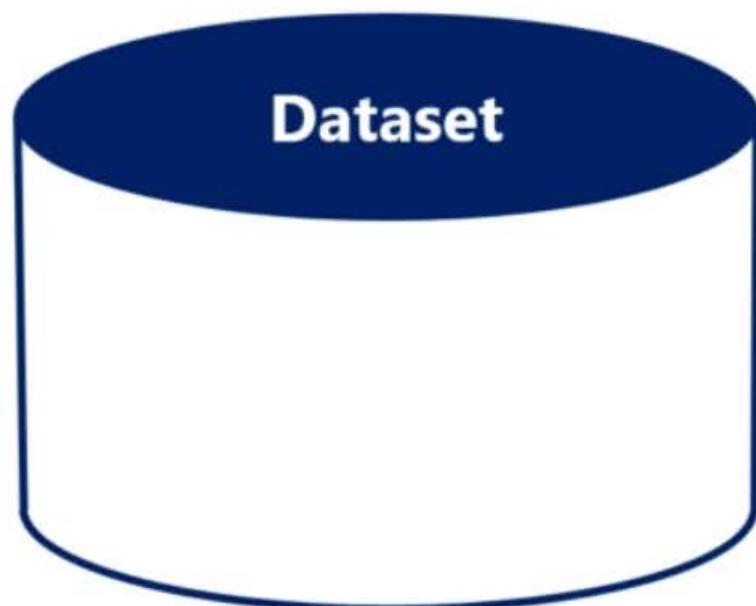
Training, validation, and test



Training, validation, and test

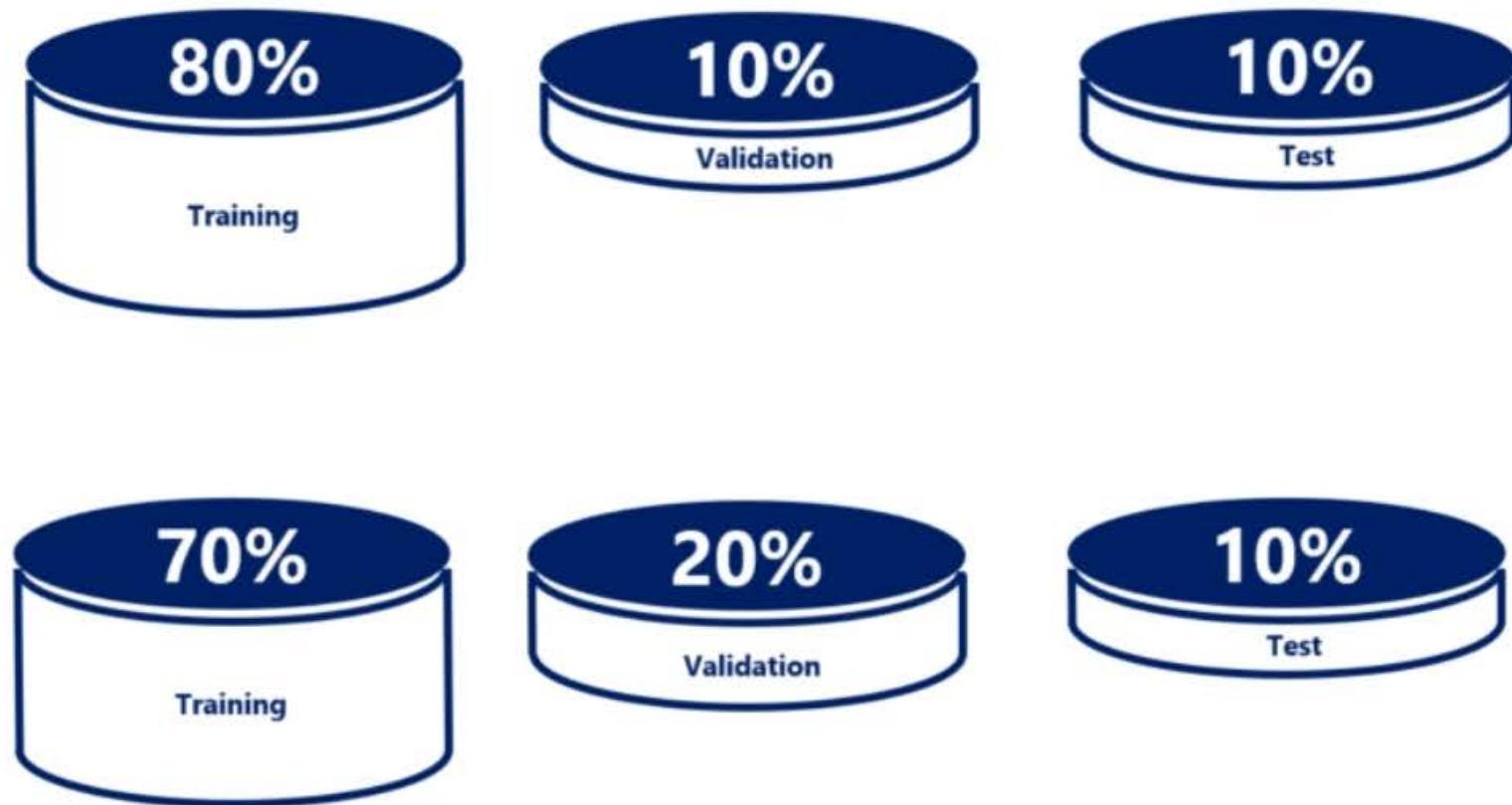
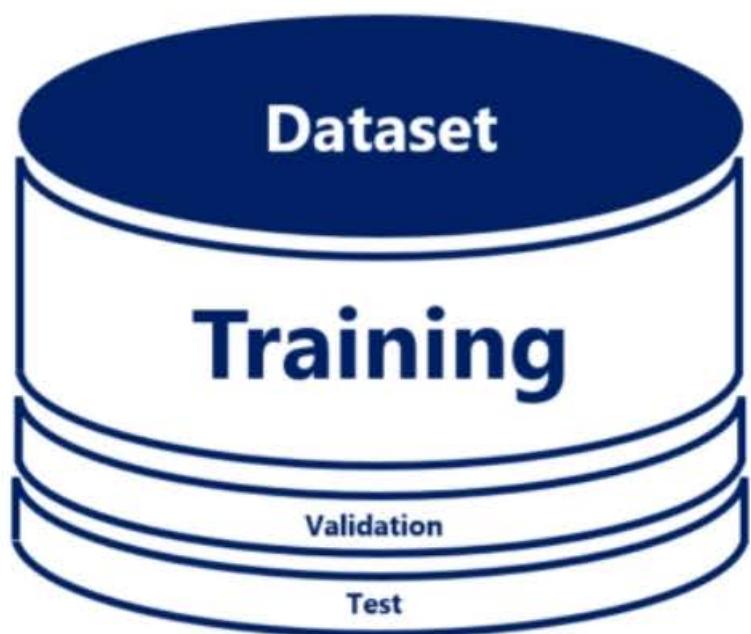


Training, validation, and test

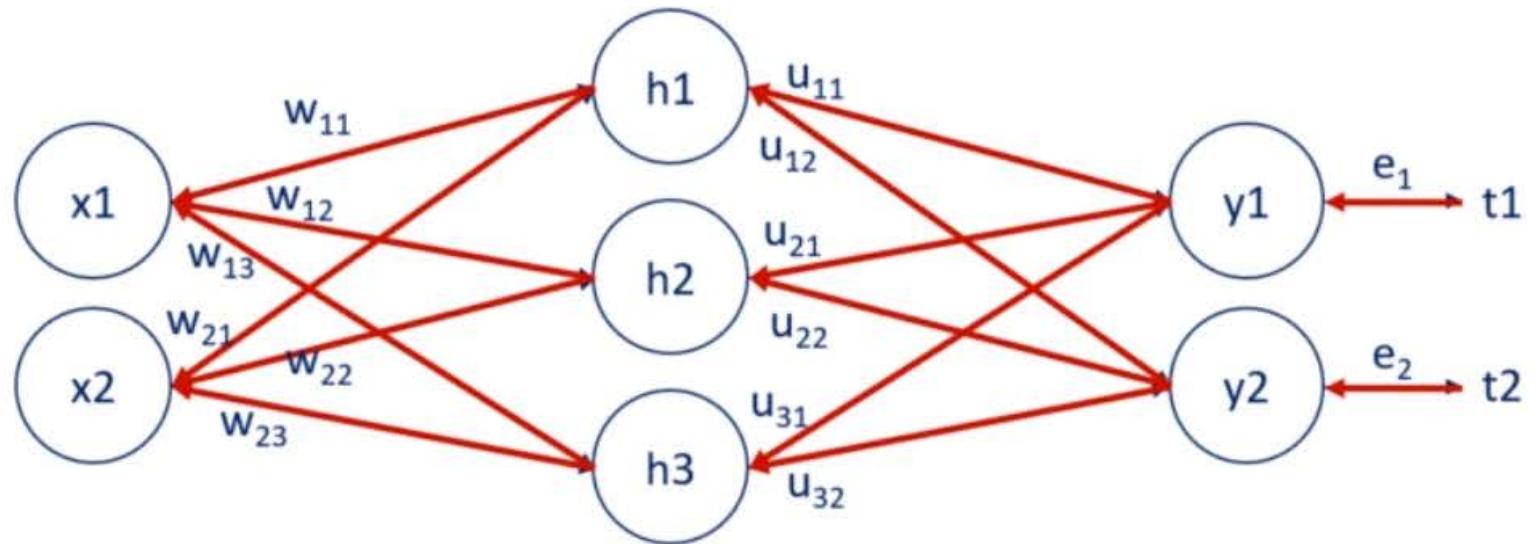


The accuracy that we get by forward propagating the test dataset, is the accuracy we expect the model to have if we deploy it in real life

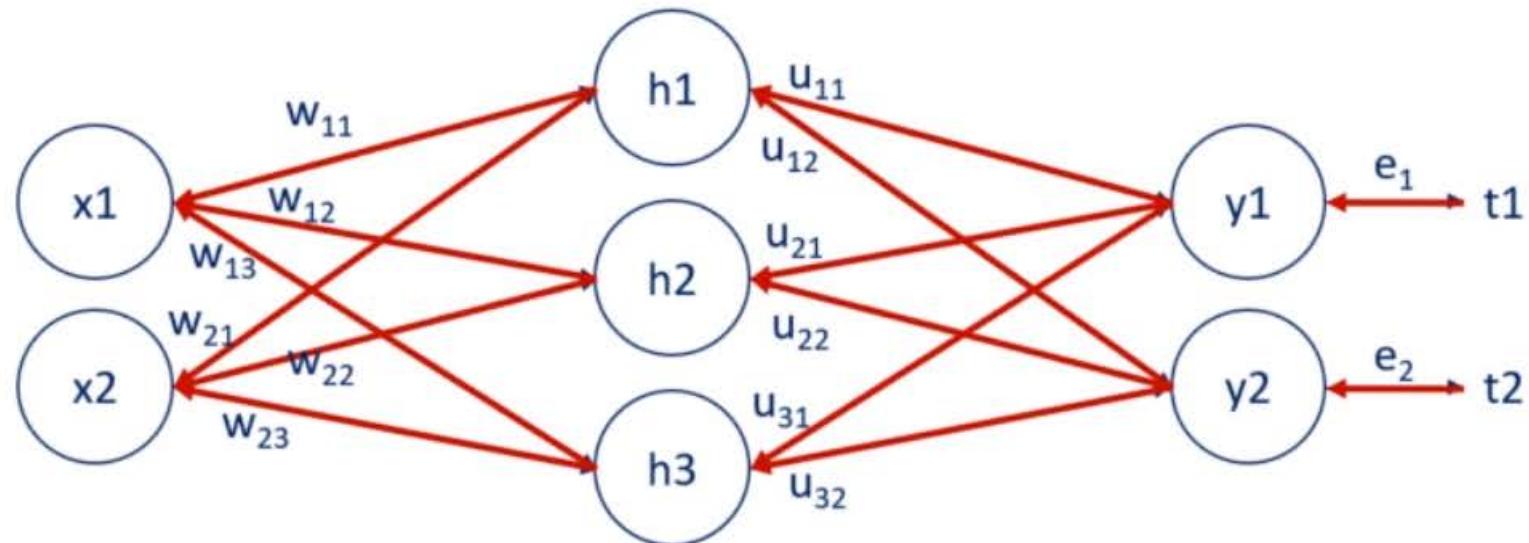
Summary



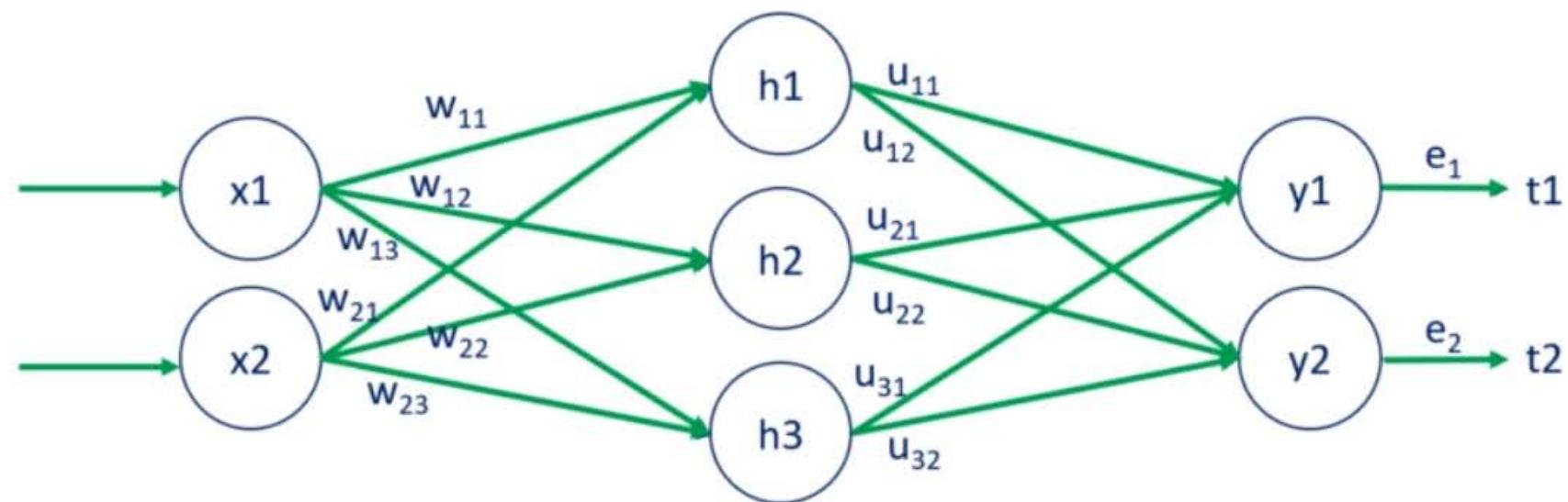
Summary



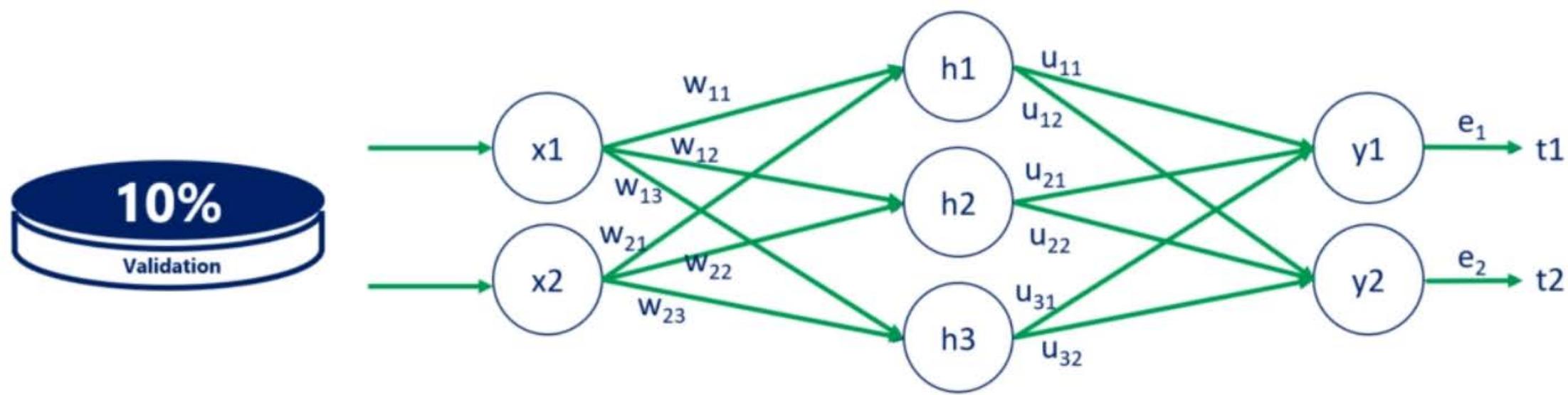
We train ONLY on the training set



4



What does "every now and then" mean?



Usually, we validate on **every epoch**

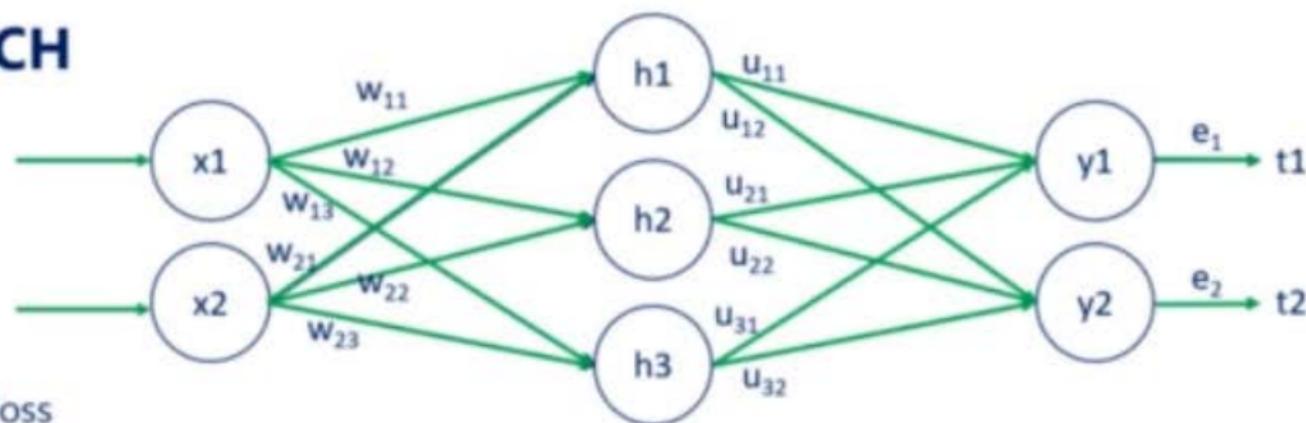
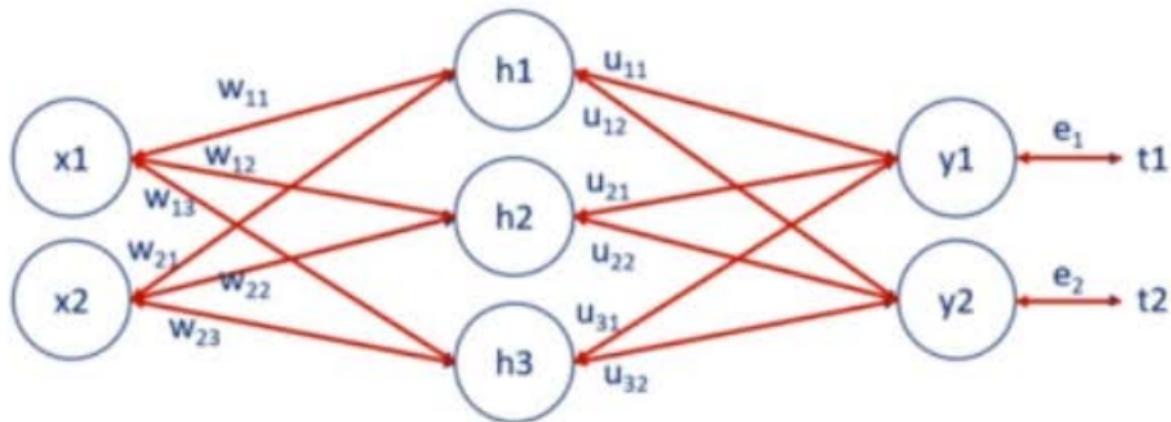


Calculate training_loss

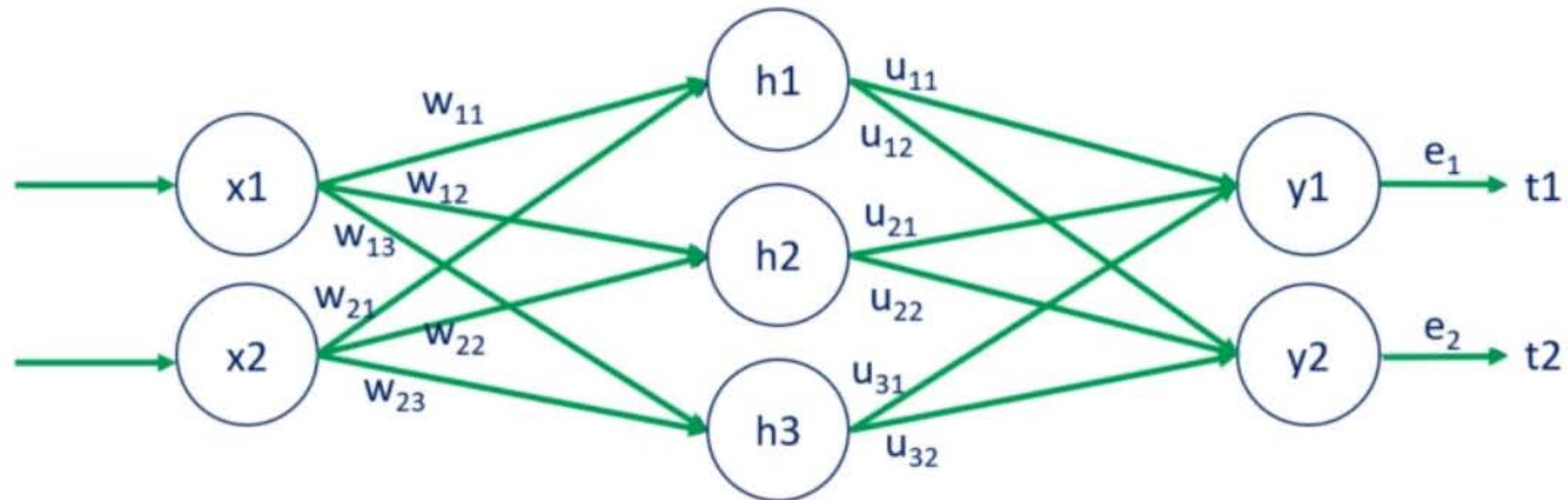
END OF EPOCH



Calculate validation_loss



10%
Test



The accuracy we obtain at this stage is the accuracy of the algorithm

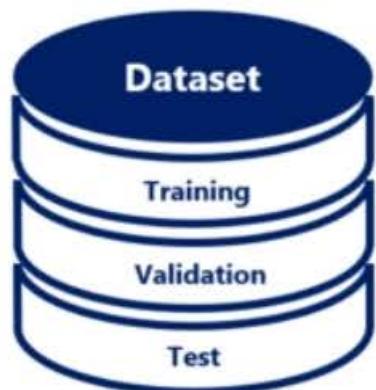
Training, validation, and test



We can't afford to split it in three...

Because the algorithm may not **learn anything**

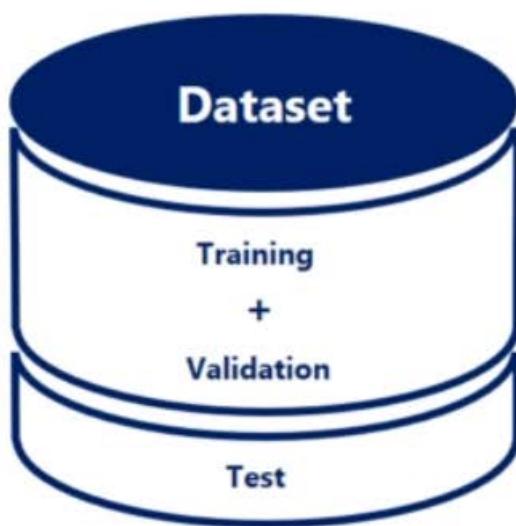
Training, validation, and test



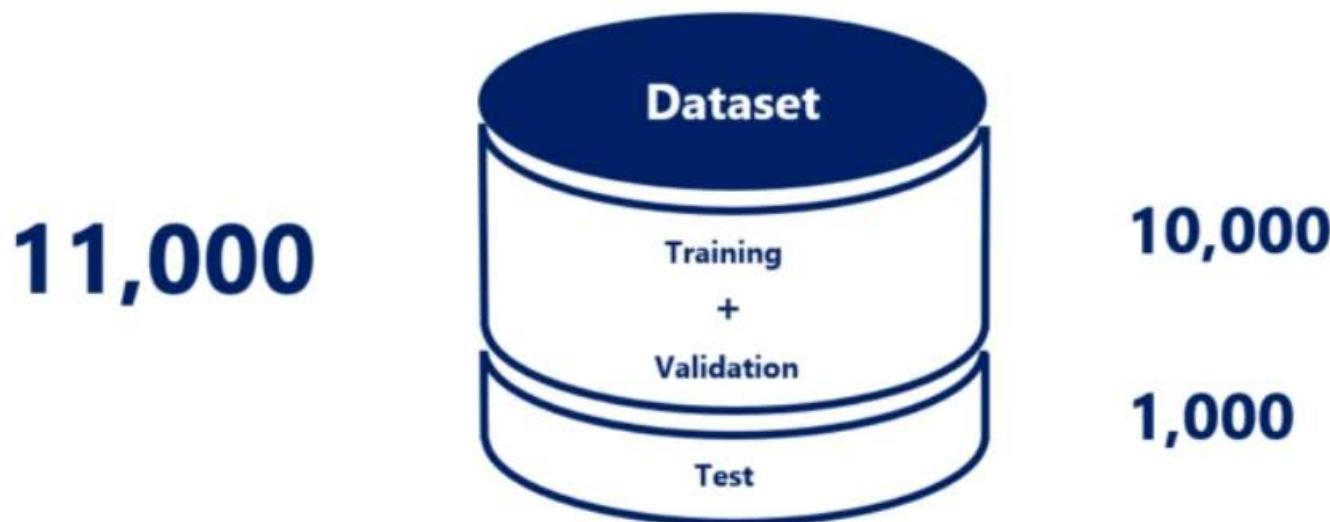
N-FOLD CROSS-VALIDATION

/ sometimes k-fold cross-validation /

N-FOLD CROSS-VALIDATION

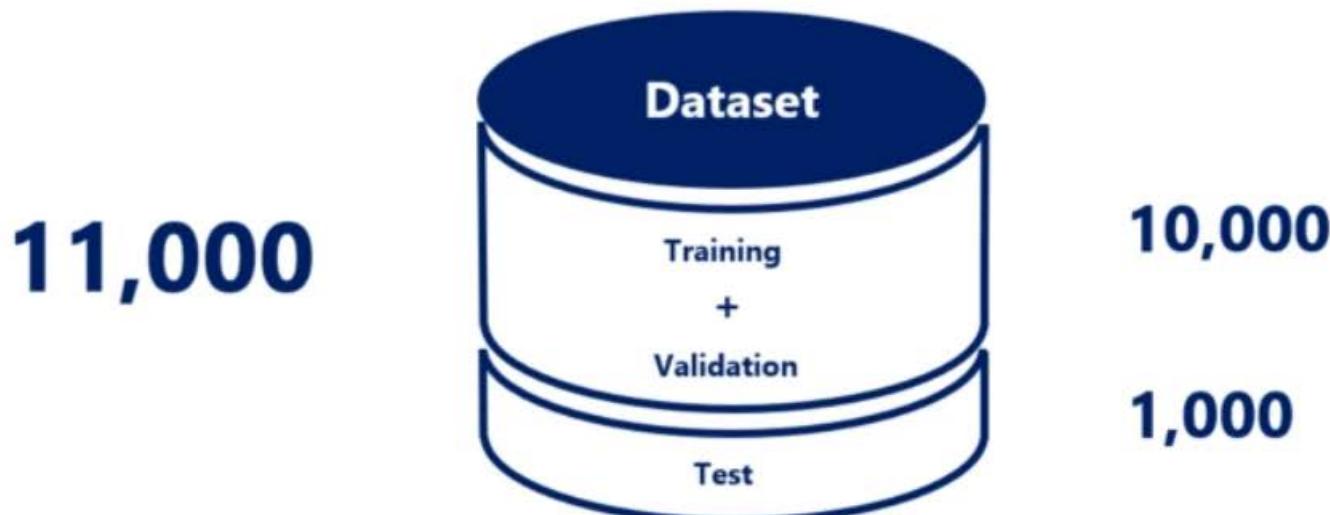


N-FOLD CROSS-VALIDATION



In data science we often deal with ginormous datasets

N-FOLD CROSS-VALIDATION



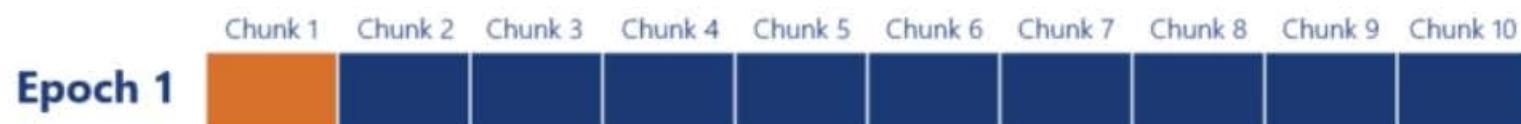
Side note:

"Ginormous" datasets have their own problems - being so large, they often have a lot of missing values. We refer to them as being: **sparse**

10-FOLD CROSS-VALIDATION



9,000
+
1,000



We fold it 10 times, so this is a **10-fold cross-validation**

10-FOLD CROSS-VALIDATION

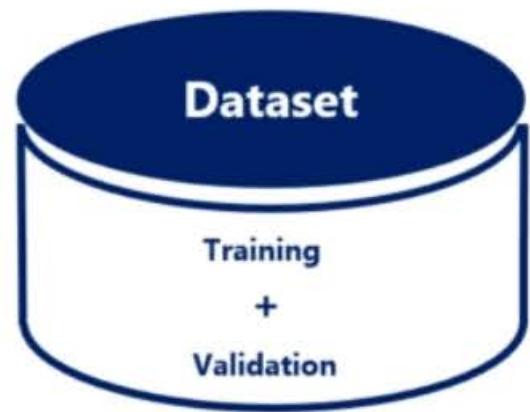


↑
validation

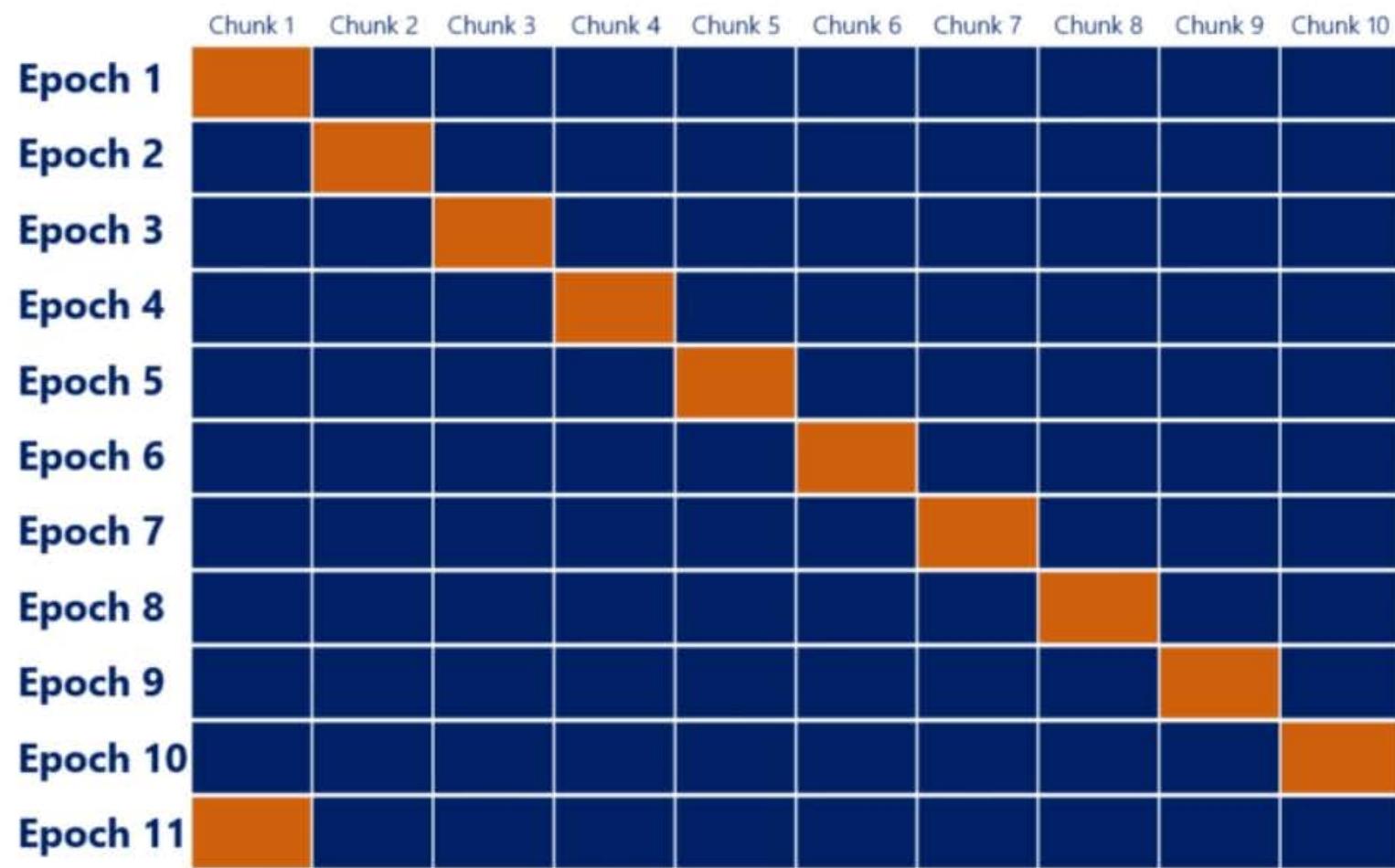
Training

Validation

10-FOLD CROSS-VALIDATION



 Training
 Validation



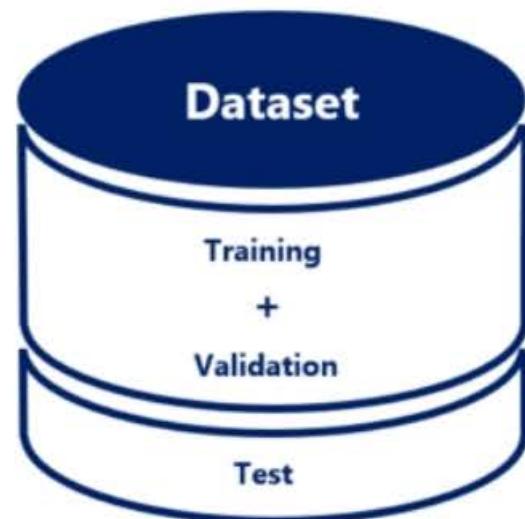
For each epoch, we don't overlap training and validation

N-FOLD CROSS-VALIDATION

Pros

Utilized more data

We have a model



Cons

Possibly overfitted a bit

Early stopping

Definition:

We want to...

... stop training early...

... i.e. before we overfit.

Early stopping

1

Preset number of epochs

Learning

```
In [9]: for e in range(100):
    _, curr_loss = sess.run([optimize, mean_loss],
                           feed_dict = {inputs: training_data['inputs'], targets: training_data['targets']})
    print (curr_loss)
```

Pros:

1. Eventually, solves the problem

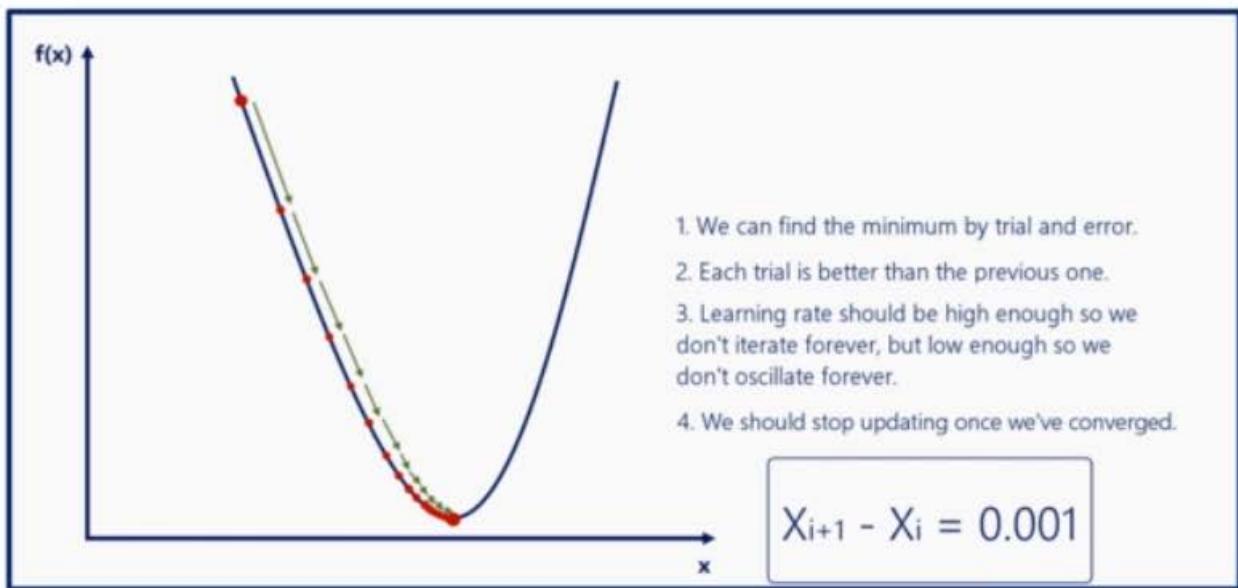
Cons:

1. No guarantee that the min is reached
2. Maybe doesn't minimize at all

Early stopping

2

Stop when updates become too small



Pros:

1. We are sure the loss is minimized

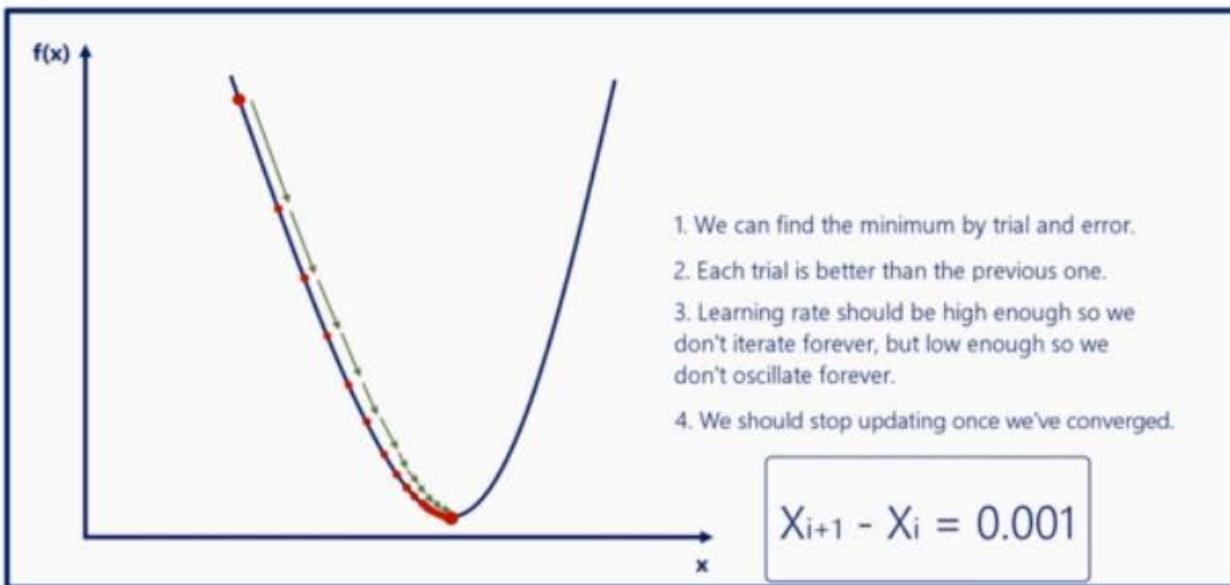
Cons:

The loss function stops changing, making the update rule yield same weights

Early stopping

2

Stop when updates become too small



Pros:

1. We are sure the loss is minimized
2. Saves computing power

Cons:

We save computing power by using as few iterations as possible; we don't iterate uselessly

Early stopping

	Preset epochs	Updates too small	
Solves the problem	✓		
Certain that loss is minimized	✗		
Doesn't iterate uselessly	✗		

We can't guess the number of required epochs

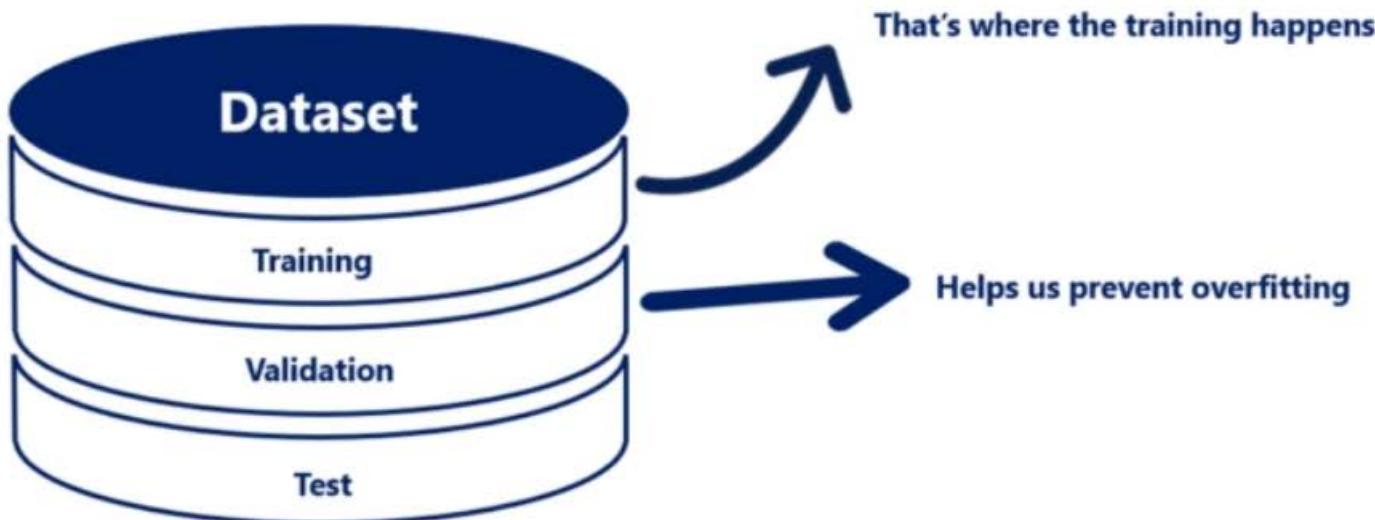
Early stopping

	Preset epochs	Updates too small	
Solves the problem	✓	✓	
Certain that loss is minimized	✗	✓	
Doesn't iterate uselessly	✗	✓	
Prevents overfitting	✗	✗	

Early stopping

3

Validation set strategy

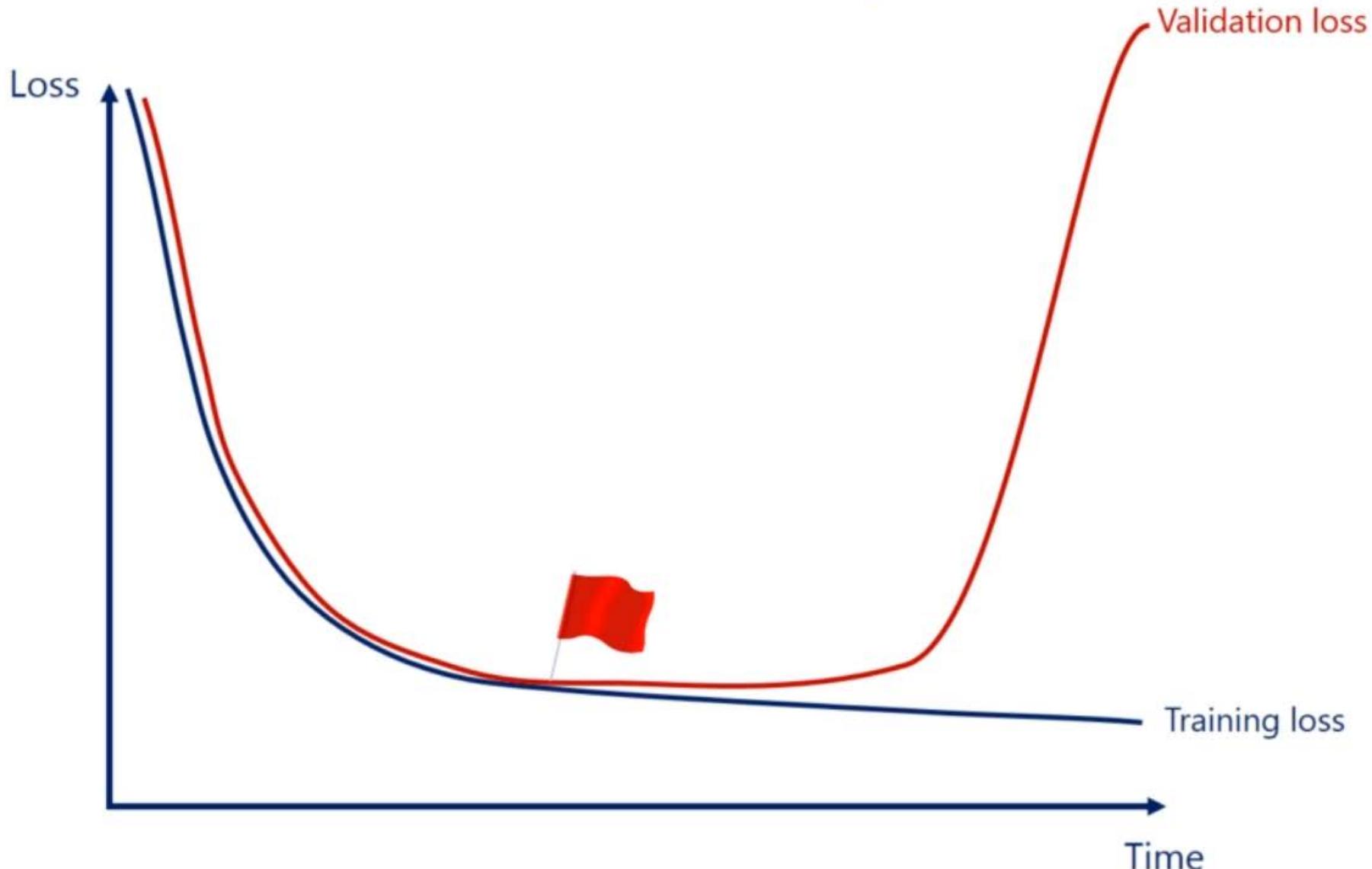


Pros:

1. We are sure the validation loss is minimized
2. Saves computing power
3. Prevents overfitting

Cons:

Typical training



Early stopping

	Preset epochs	Updates too small	Validation set
Solves the problem	✓	✓	✓
Certain that loss is minimized	✗	✓	✓
Doesn't iterate uselessly	✗	✓	✗
Prevents overfitting	✗	✗	✓

Early stopping

	Preset epochs	Updates too small	Validation set
Solves the problem	✓	✓	✓
Certain that loss is minimized	✗	✓	✓
Doesn't iterate uselessly	✗	✓	✗
Prevents overfitting	✗	✗	✓



Minimal example

Early stopping

	Preset epochs	Updates too small	Validation set
Solves the problem	✓	✓	✓
Certain that loss is minimized	✗	✓	✓
Doesn't iterate uselessly	✗	✓	✗
Prevents overfitting	✗	✗	✓

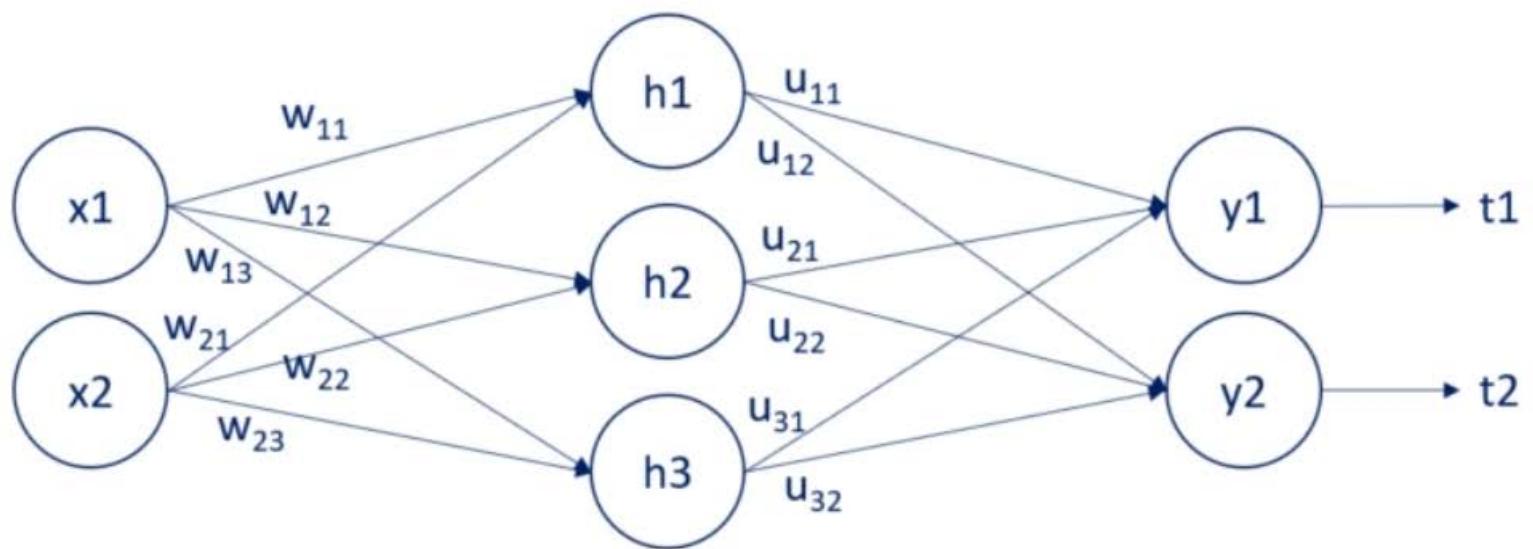
Minimal example

Best practice

INITIALIZATION

Initialization is the process in which we set the initial values of weights.

Initialization



Let's initialize our weights and biases in such a way that they **are all equal to a constant**

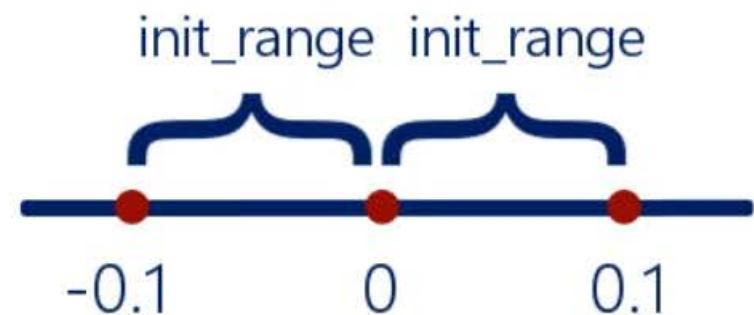
File Edit View Insert Cell Kernel Widgets Help

Python 3



CellToolbar

Our initial weights and biases will be picked randomly from the interval [-0.1,0.1]



Initialize variables

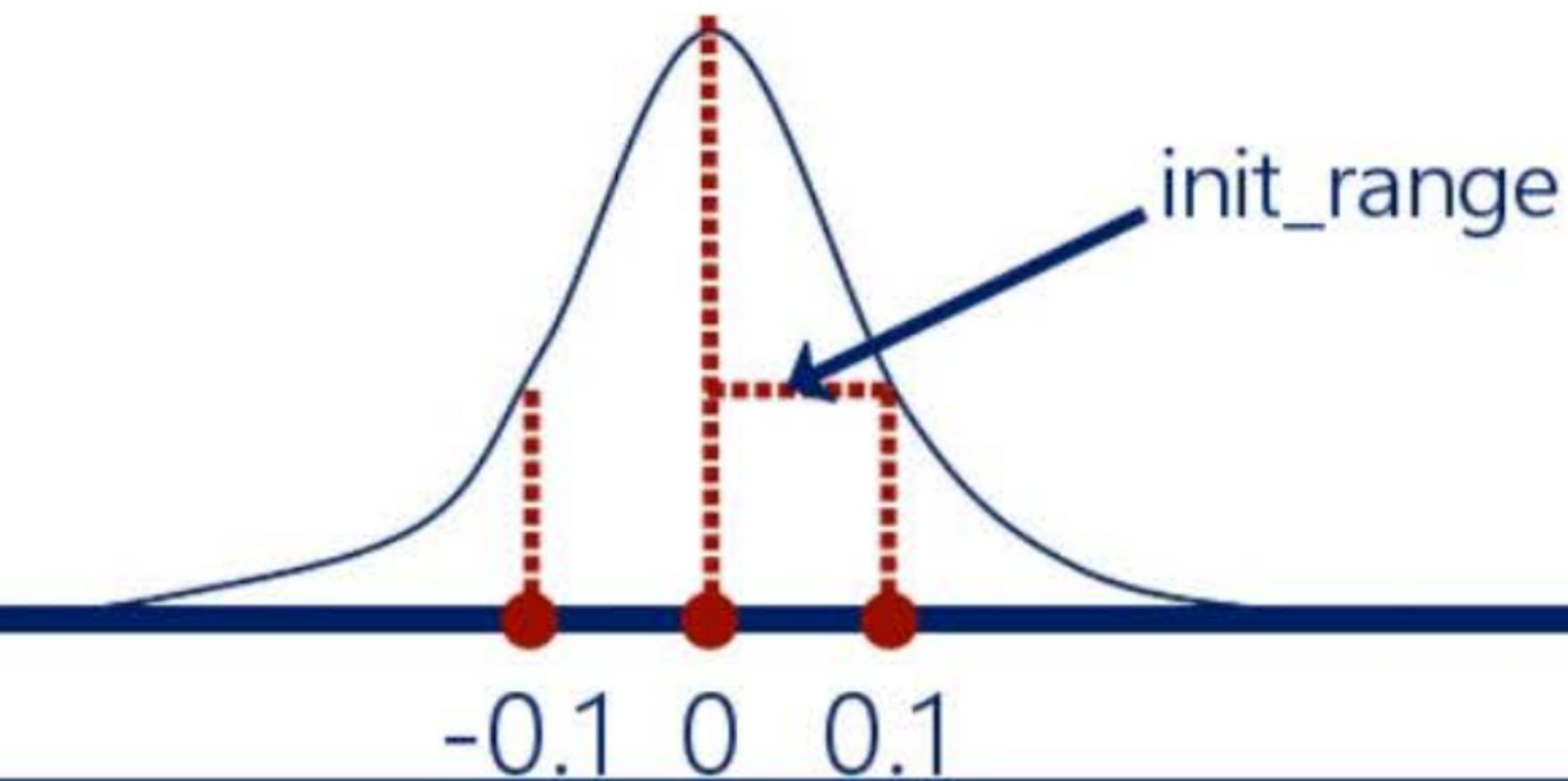
```
In [10]: init_range = 0.1  
weights = np.random.uniform(-init_range,init_range, size=(2,1))  
biases = np.random.uniform(-init_range,init_range,size=1)  
print(weights)  
print(biases)
```

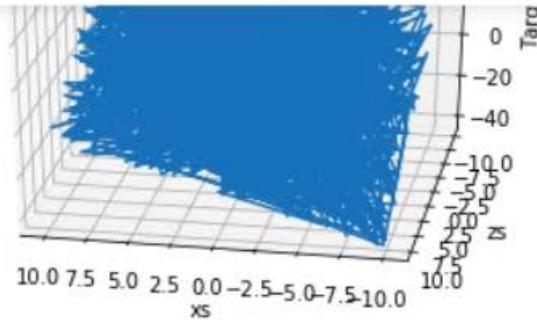
W
b

```
[[ -0.02758852]  
[ -0.03556078]  
[ -0.05488636]]
```

In []:

Our initial weights and biases will be picked randomly from the interval $[-0.1, 0.1]$ in a random normal manner, where the mean is 0 and the standard deviation is 0.1 (variance 0.01).





In []:

Initialize variables

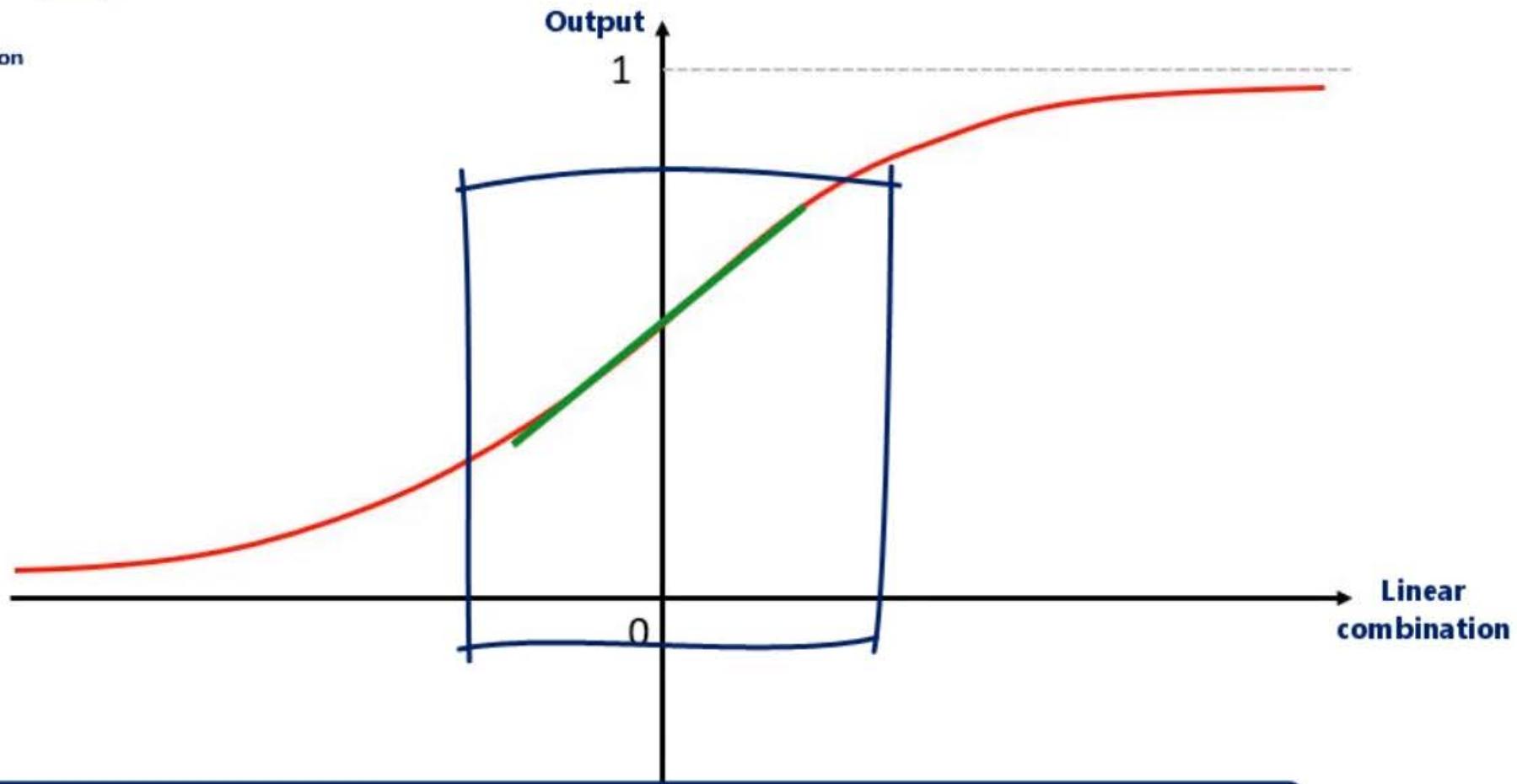
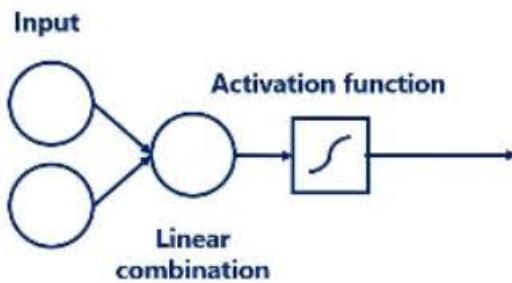
```
In [10]: init_range = 0.1  
  
weights = np.random.uniform(-init_range,init_range, size=(2,1))  
  
biases = np.random.uniform(-init_range,init_range,size=1)  
  
print(weights)  
print(biases)  
  
[[ -0.02758852]  
 [-0.03556078]]  
[-0.05488636]
```

In []:

In []:

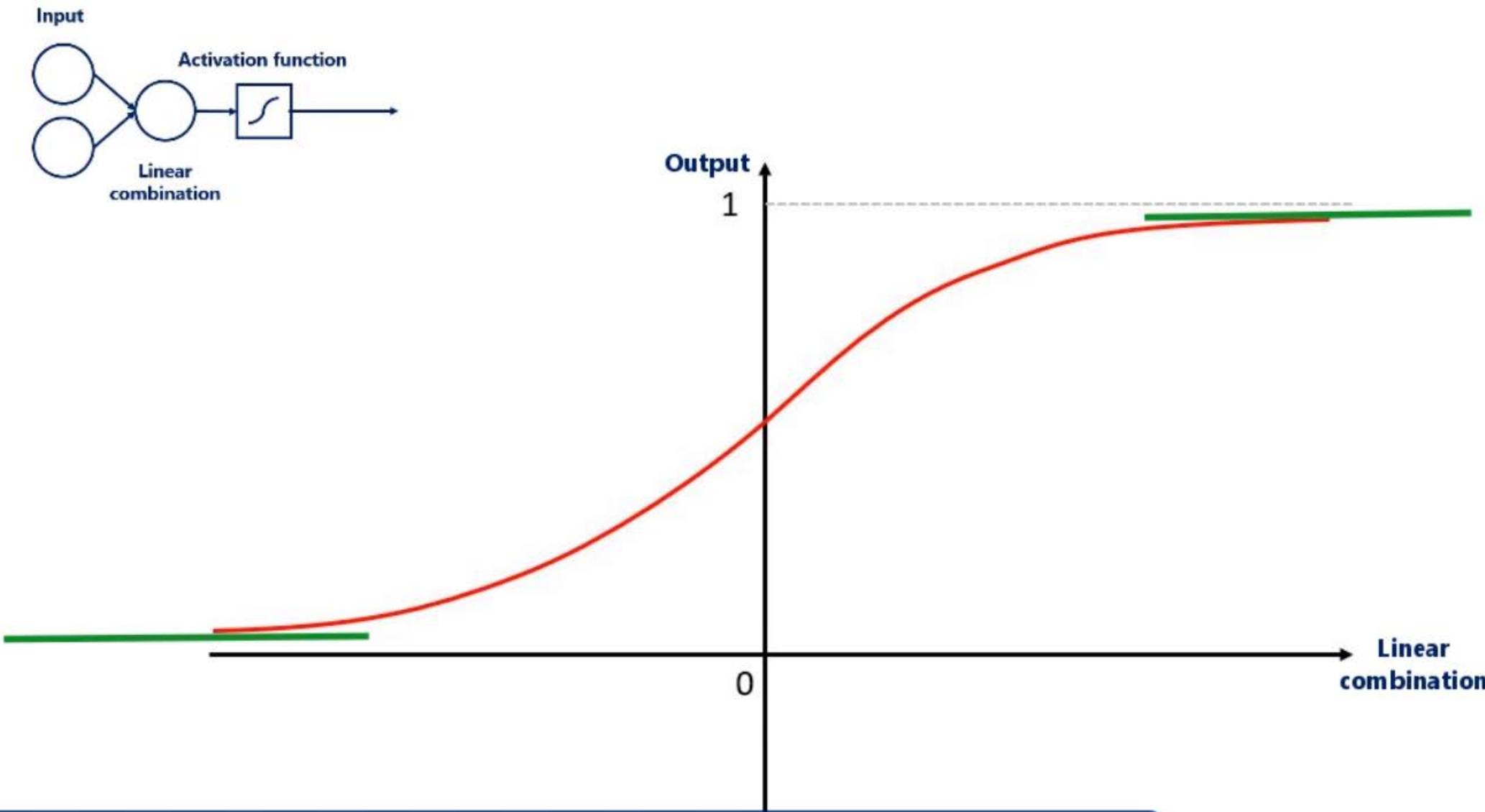
Both methods are somewhat **problematic**

Sigmoid activation function



If all the inputs for the sigmoid (the linear combinations) are in this range, the activation will be linear

Sigmoid activation function



A static output (only ~0s or ~1s) is also problematic, as the algorithm doesn't learn

Xavier initialization

Glorot initialization

Xavier Glorot

initialization

Xavier initialization

Original paper

Understanding the difficulty of training deep
feedforward neural networks

Xavier Glorot, Yoshua Bengio ;

*Proceedings of the Thirteenth International Conference on Artificial
Intelligence and Statistics, PMLR 9:249-256, 2010.*

Xavier initialization

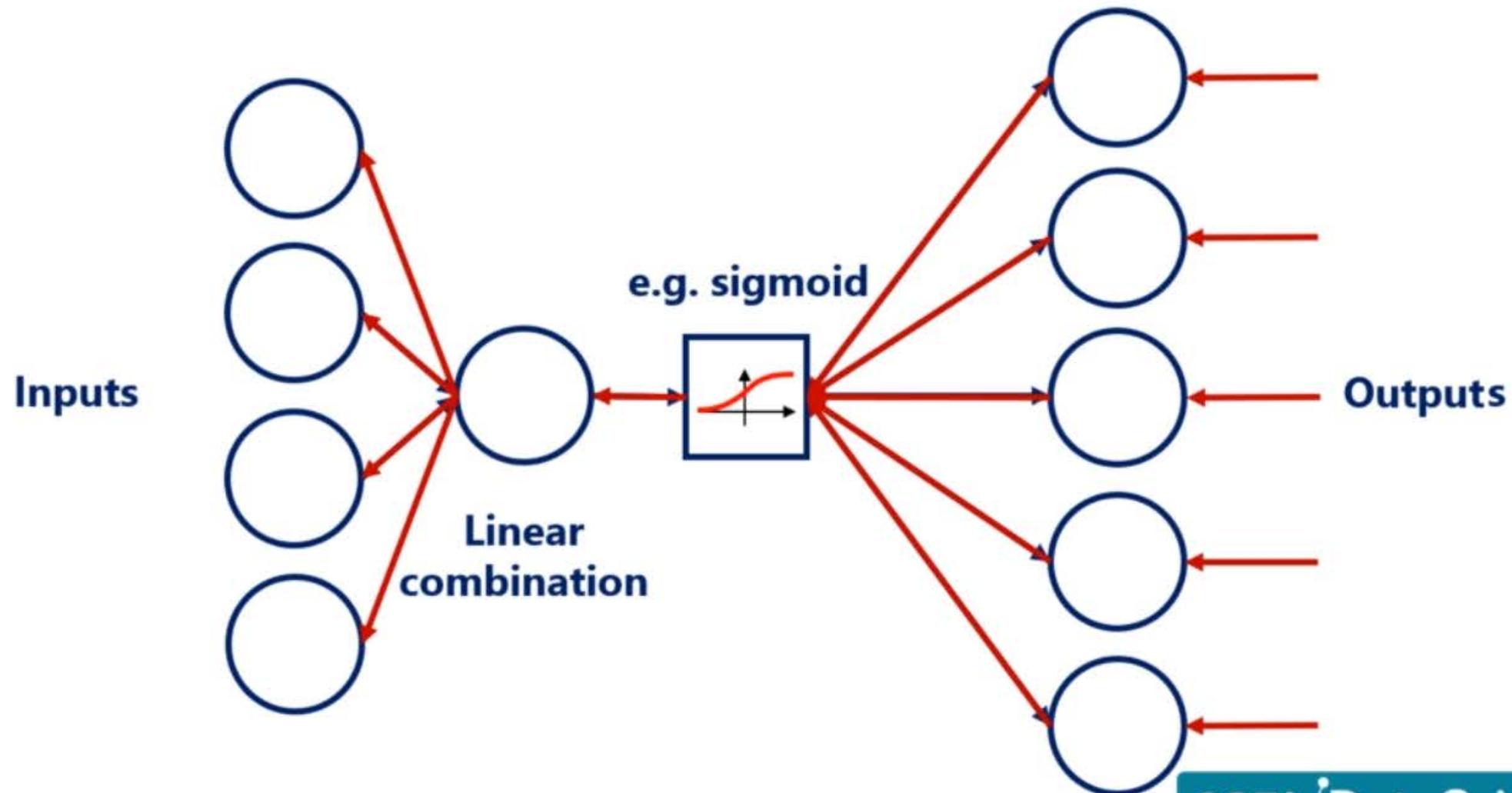
Uniform Xavier initialization: draw each weight, w , from a random uniform distribution

$$\text{in } [-x, x] \text{ for } x = \sqrt{\frac{6}{\text{inputs} + \text{outputs}}}$$

Normal Xavier initialization: draw each weight, w , from a normal distribution with a mean

$$\text{of 0, and a standard deviation } \sigma = \sqrt{\frac{2}{\text{inputs} + \text{outputs}}}$$

Why do inputs and outputs matter?



- get_collection
- get_collection_ref
- get_default_graph
- get_default_session
- get_local_variable
- get_seed

```
v1 = foo() # creates v.  
v2 = foo() # Gets the same, existing v.  
assert v1 == v2
```

If initializer is `None` (the default), the default initializer passed in the variable scope will be used. If that one is `None` too, a `glorot_uniform_initializer` will be used. The initializer can also be a Tensor, in which case the variable is initialized to this value and shape.

OPTIMIZATION

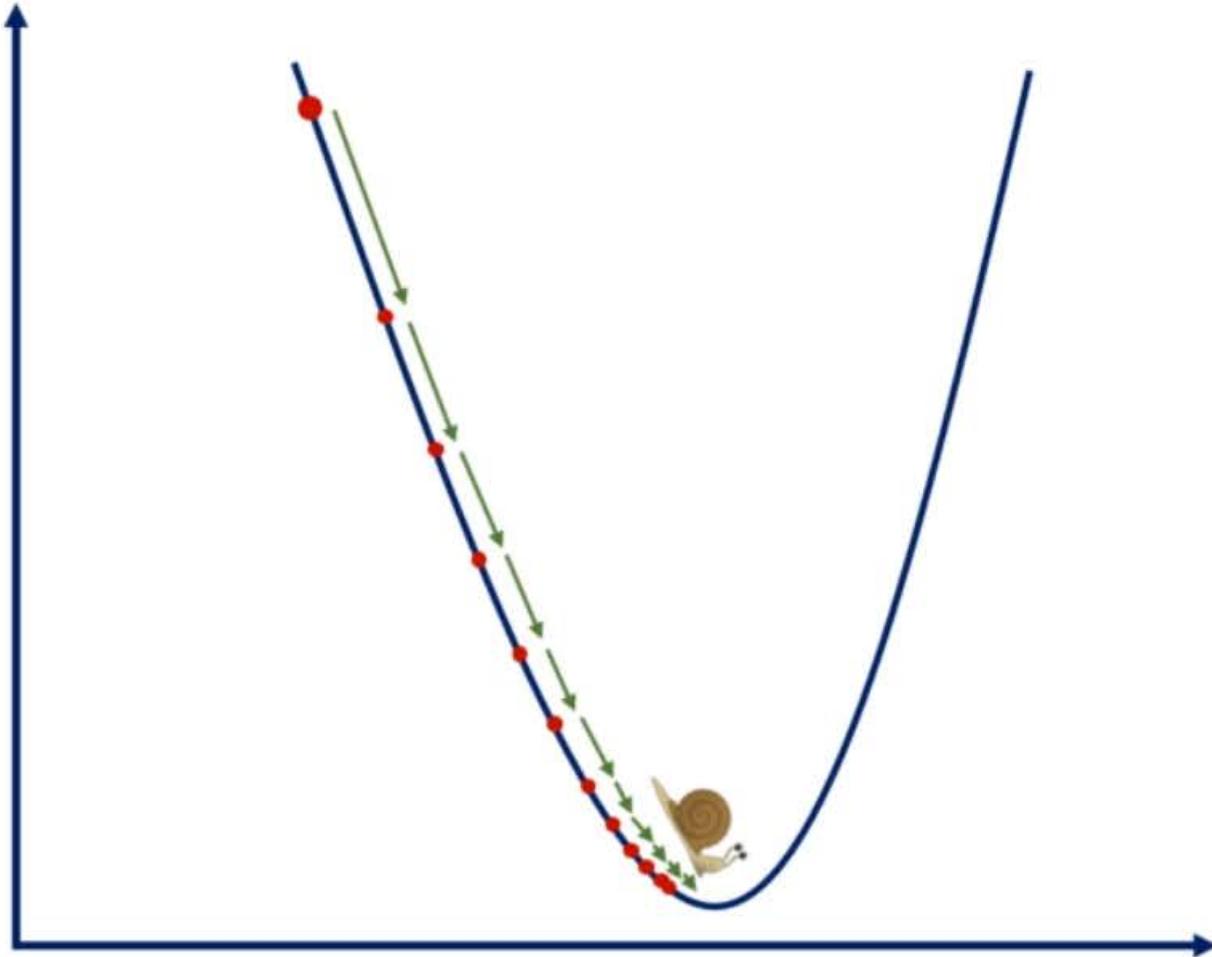


x_{11}	x_{12}	x_{13}	x_{14}	x_{15}	x_{16}	x_{17}	x_{18}
x_{21}	x_{22}	x_{23}	x_{24}	x_{25}	x_{26}	x_{27}	x_{28}
x_{31}	x_{32}	x_{33}	x_{34}	x_{35}	x_{36}	x_{37}	x_{38}
x_{41}	x_{42}	x_{43}	x_{44}	x_{45}	x_{46}	x_{47}	x_{48}
x_{51}	x_{52}	x_{53}	x_{54}	x_{55}	x_{56}	x_{57}	x_{58}
x_{61}	x_{62}	x_{63}	x_{64}	x_{65}	x_{66}	x_{67}	x_{68}
x_{71}	x_{72}	x_{73}	x_{74}	x_{75}	x_{76}	x_{77}	x_{78}
x_{81}	x_{82}	x_{83}	x_{84}	x_{85}	x_{86}	x_{87}	x_{88}

GRADIENT DESCENT

$n \times k$

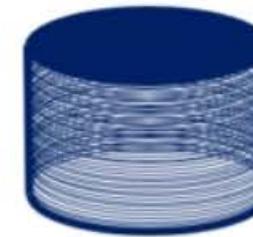
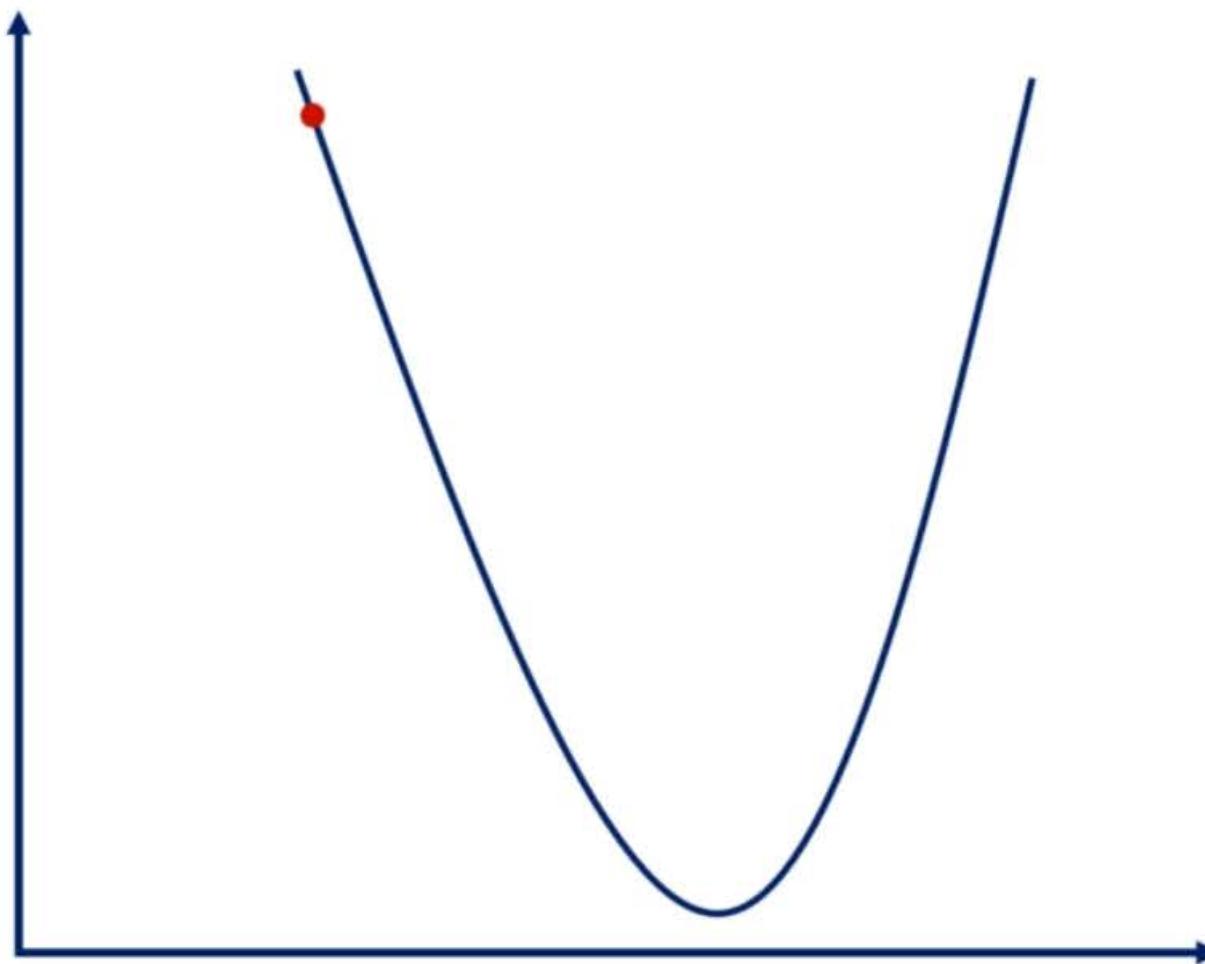
Gradient descent



X_{11}	X_{12}	X_{13}	X_{14}	X_{15}	X_{16}	X_{17}	X_{18}
X_{21}	X_{22}	X_{23}	X_{24}	X_{25}	X_{26}	X_{27}	X_{28}
X_{31}	X_{32}	X_{33}	X_{34}	X_{35}	X_{36}	X_{37}	X_{38}
-	-	-	-	-	-	-	-
-	-	-	-	-	-	-	-
X_{n1}	X_{n2}	X_{n3}	X_{n4}	X_{n5}	X_{n6}	X_{n7}	X_{n8}

$n \times k$

Stochastic gradient descent



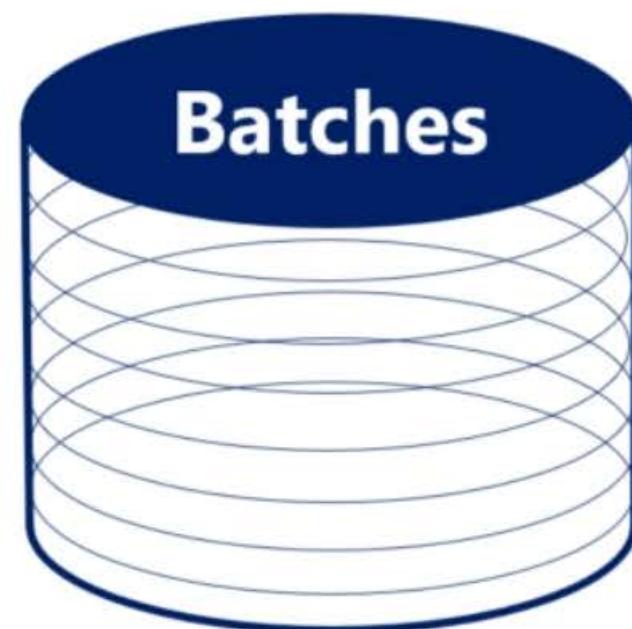
X_{11}	X_{12}	X_{13}	X_{14}	X_{15}	X_{16}	X_{17}	X_{18}
X_{21}	X_{22}	X_{23}	X_{24}	X_{25}	X_{26}	X_{27}	X_{28}
X_{31}	X_{32}	X_{33}	X_{34}	X_{35}	X_{36}	X_{37}	X_{38}
.
.
X_{n1}	X_{n2}	X_{n3}	X_{n4}	X_{n5}	X_{n6}	X_{n7}	X_{n8}

$n \times k$

It works in the same way, but updates the weights many times inside a single epoch

Batching

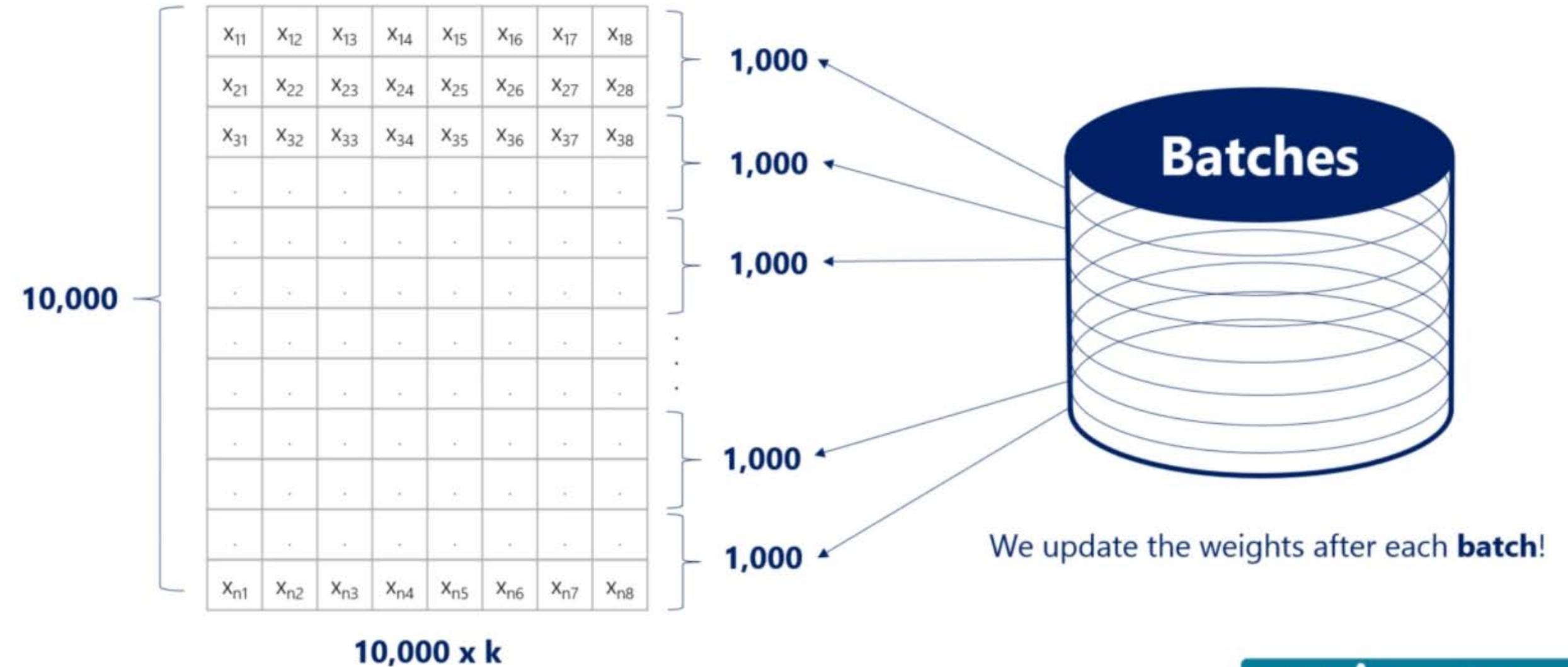
x_{11}	x_{12}	x_{13}	x_{14}	x_{15}	x_{16}	x_{17}	x_{18}
x_{21}	x_{22}	x_{23}	x_{24}	x_{25}	x_{26}	x_{27}	x_{28}
x_{31}	x_{32}	x_{33}	x_{34}	x_{35}	x_{36}	x_{37}	x_{38}
.
.
.
.
.
.
x_{n1}	x_{n2}	x_{n3}	x_{n4}	x_{n5}	x_{n6}	x_{n7}	x_{n8}



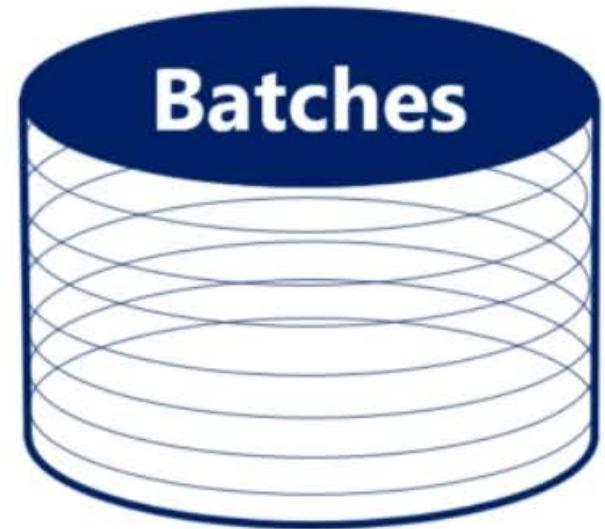
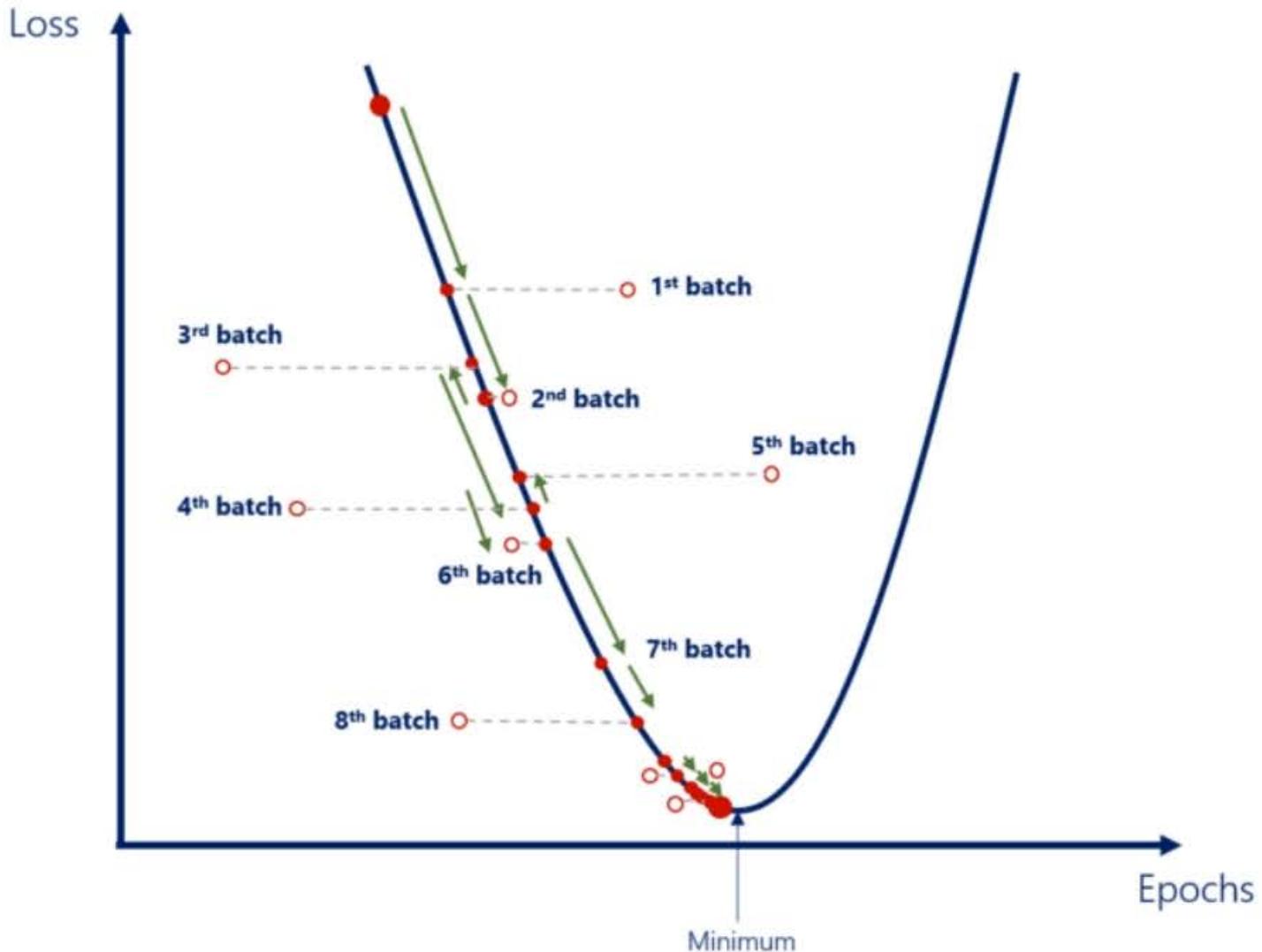
We update the weights after each **batch**!

Batching: the process of splitting the dataset in **n** batches (mini-batches)

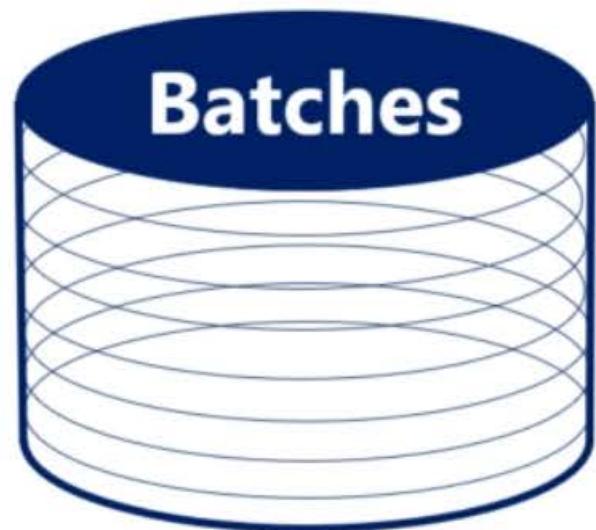
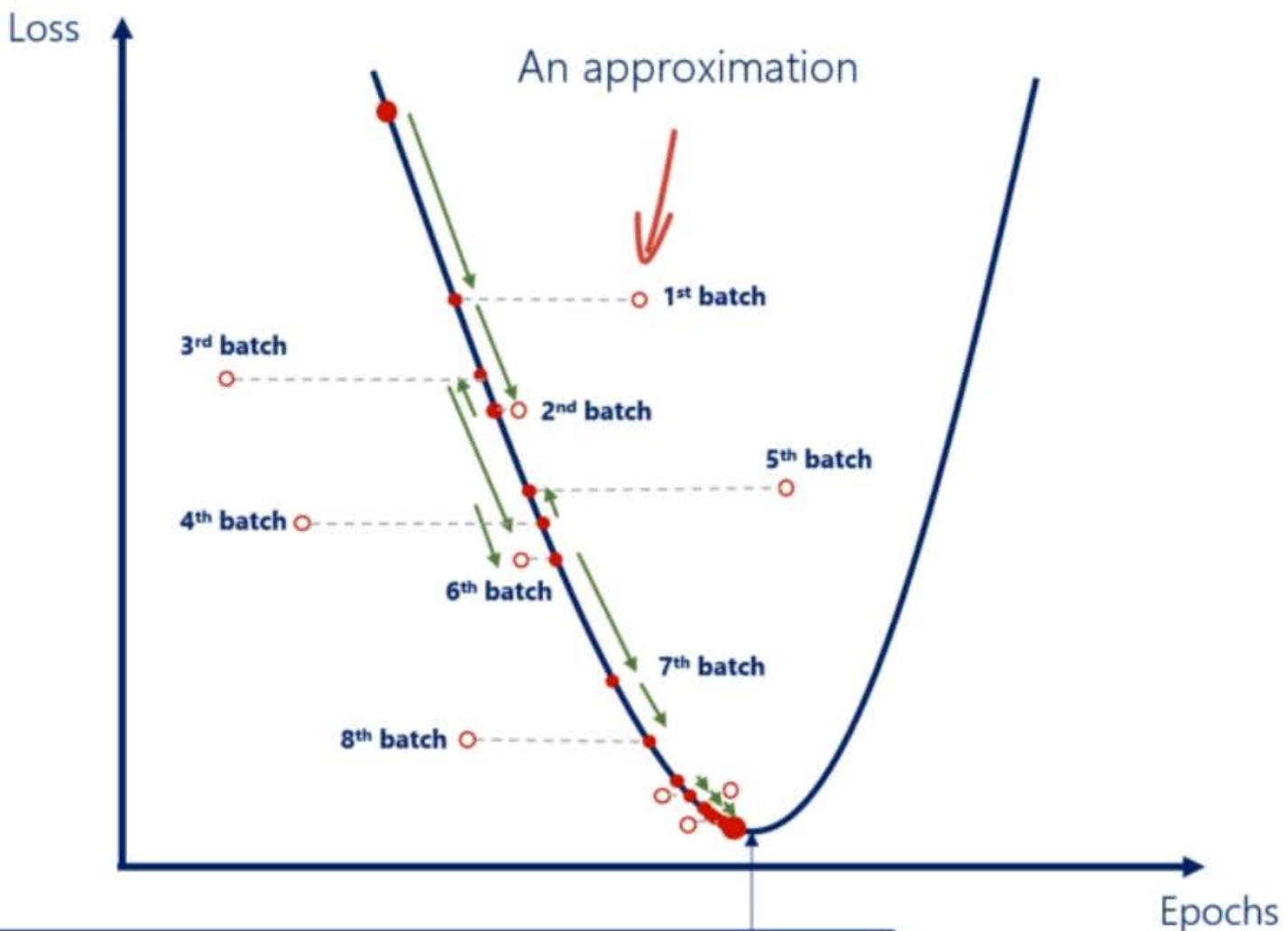
Batching



Batching

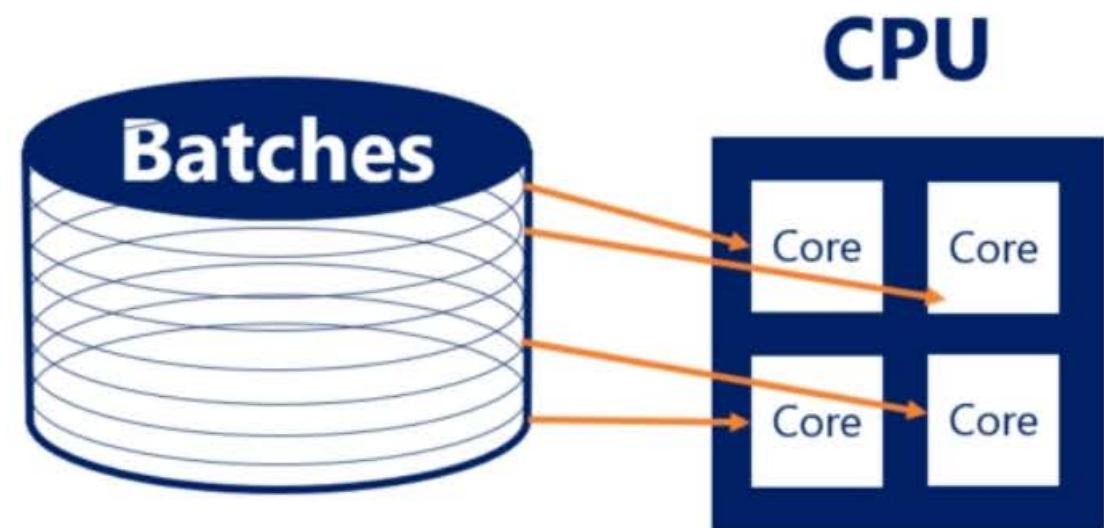
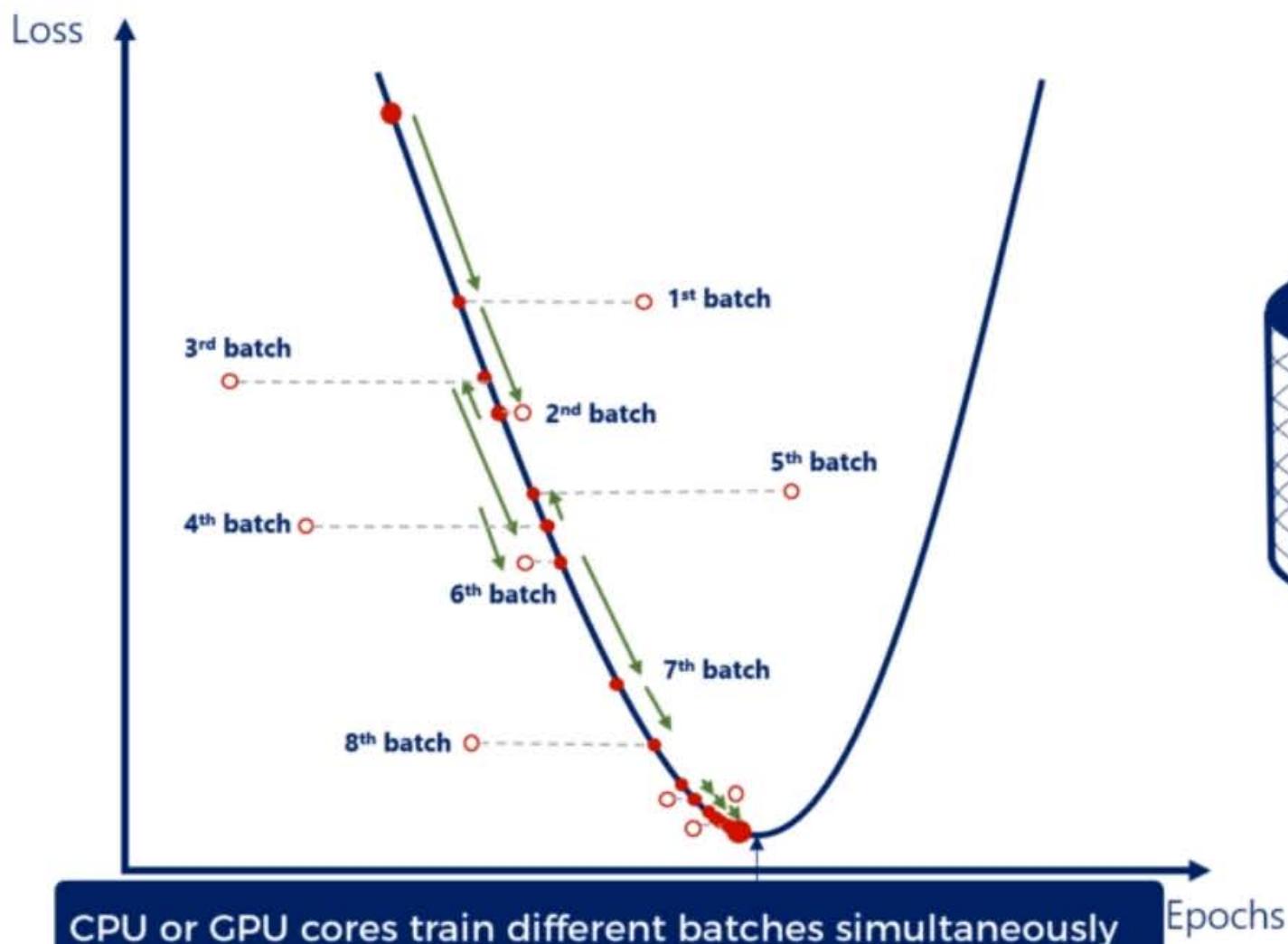


Batching



The SGD comes at a cost: it approximates things a bit

Batching

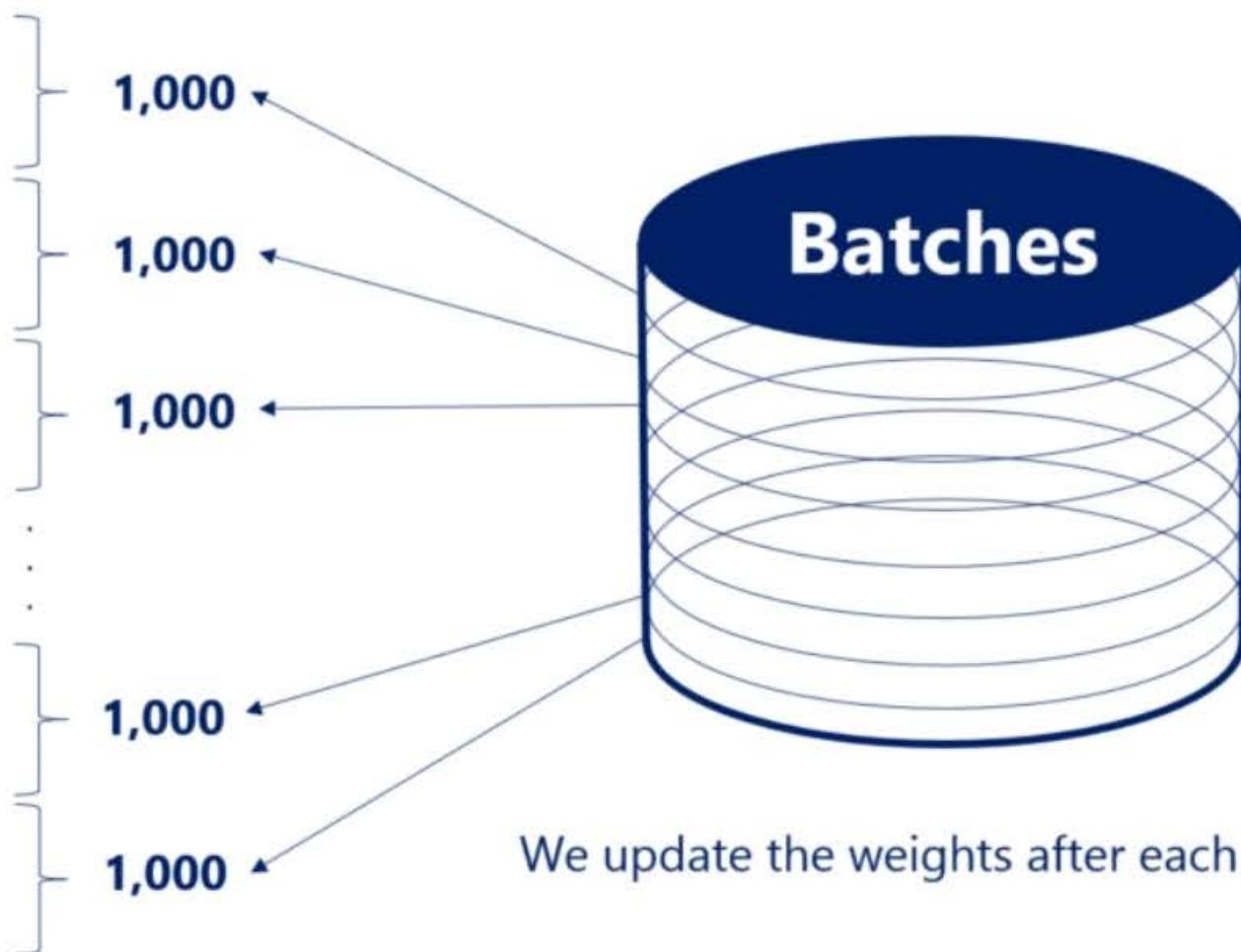




Mini-batch GD

10,000

	X ₁₂	X ₁₃	X ₁₄	X ₁₅	X ₁₆	X ₁₇	X ₁₈
X ₂₁	X ₂₂	X ₂₃	X ₂₄	X ₂₅	X ₂₆	X ₂₇	X ₂₈
X ₃₁	X ₃₂	X ₃₃	X ₃₄	X ₃₅	X ₃₆	X ₃₇	X ₃₈
.
.
.
.
.
.
X _{n1}	X _{n2}	X _{n3}	X _{n4}	X _{n5}	X _{n6}	X _{n7}	X _{n8}



Practitioners refer to the mini-batch GD as SGD

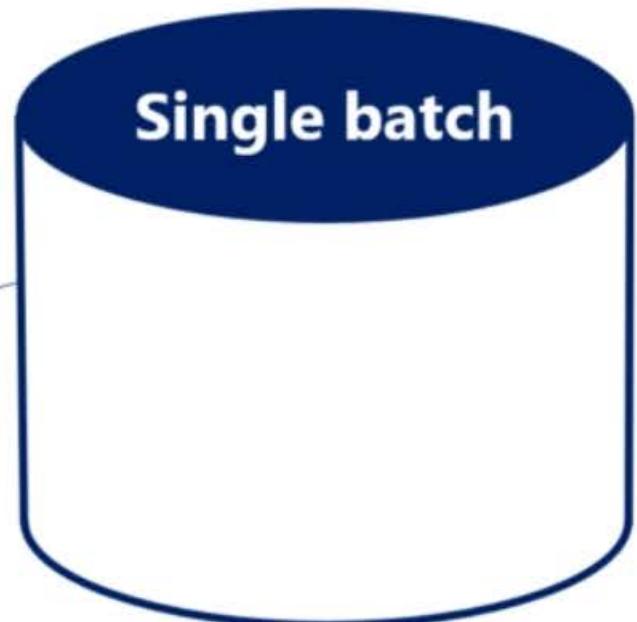
Batch GD

10,000

X_{11}	X_{12}	X_{13}	X_{14}	X_{15}	X_{16}	X_{17}	X_{18}
X_{21}	X_{22}	X_{23}	X_{24}	X_{25}	X_{26}	X_{27}	X_{28}
X_{31}	X_{32}	X_{33}	X_{34}	X_{35}	X_{36}	X_{37}	X_{38}
.
.
.
.
.
.
.
X_{n1}	X_{n2}	X_{n3}	X_{n4}	X_{n5}	X_{n6}	X_{n7}	X_{n8}

10,000 x k

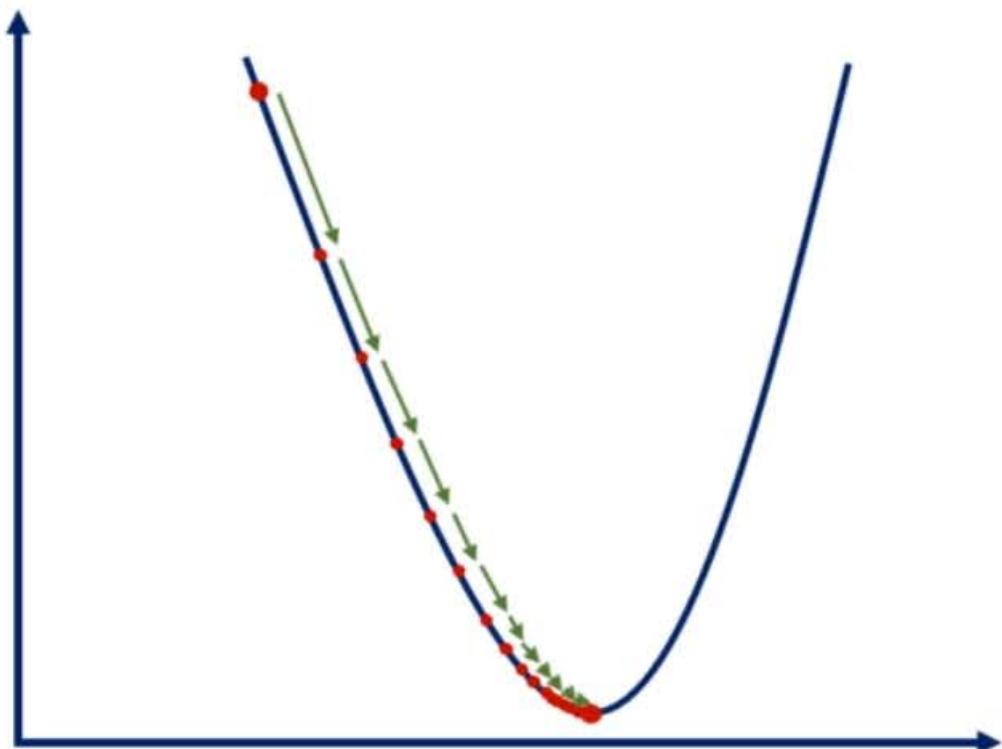
10,000



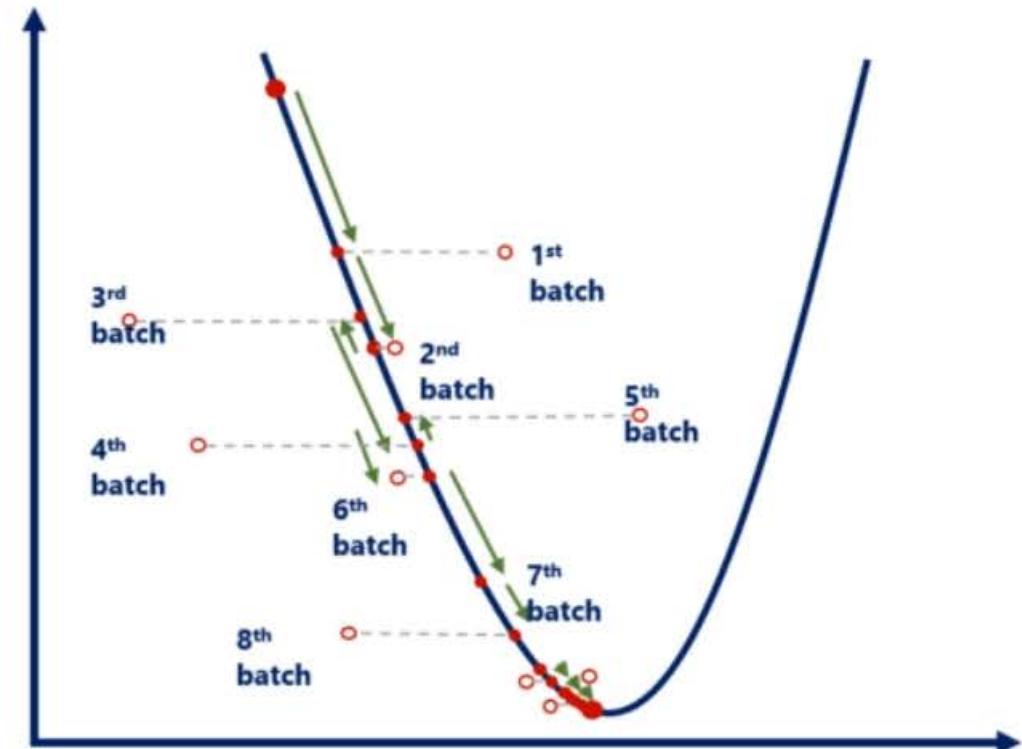
We update the weights once per epoch!

Gradient descent pitfalls

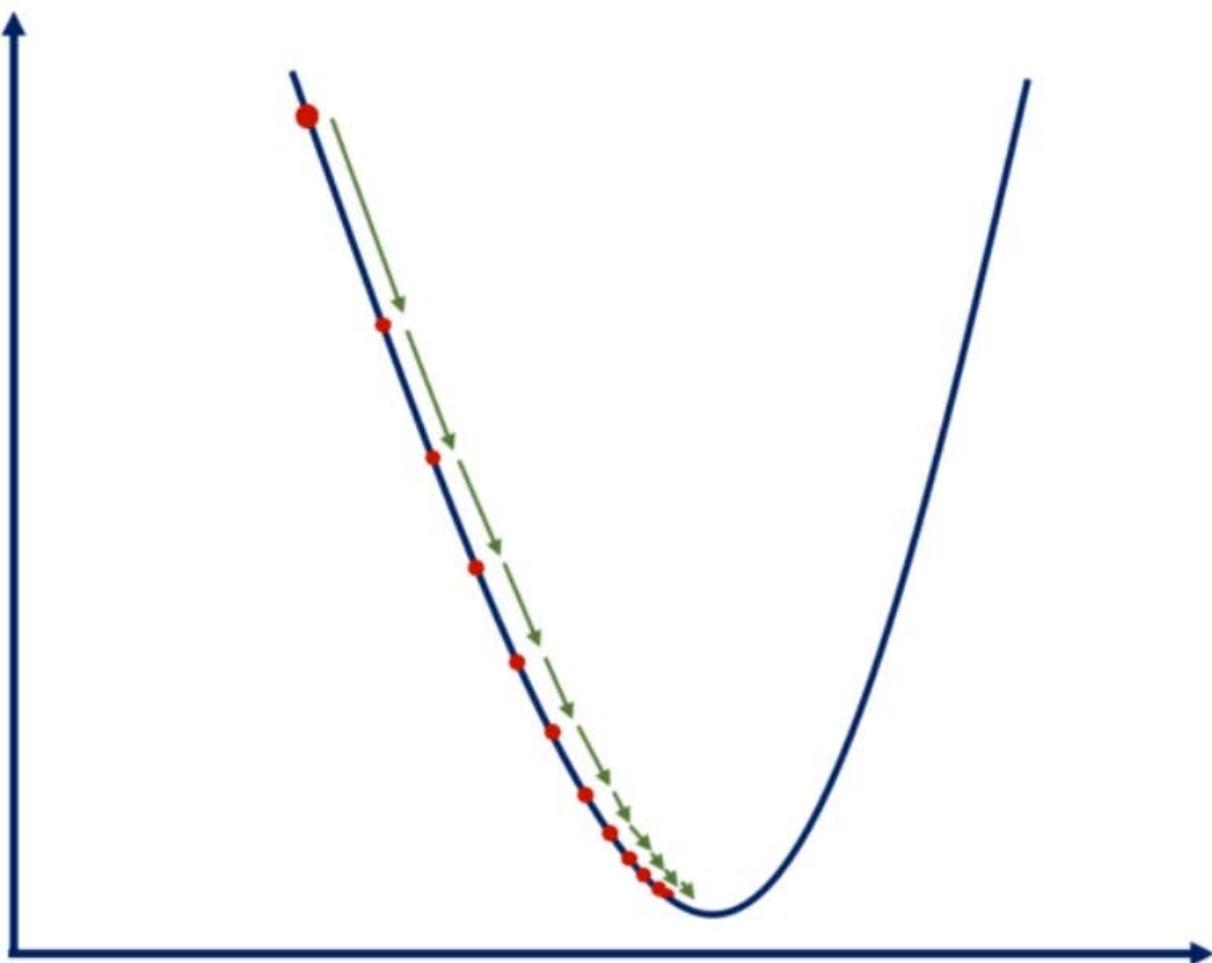
GD



SGD

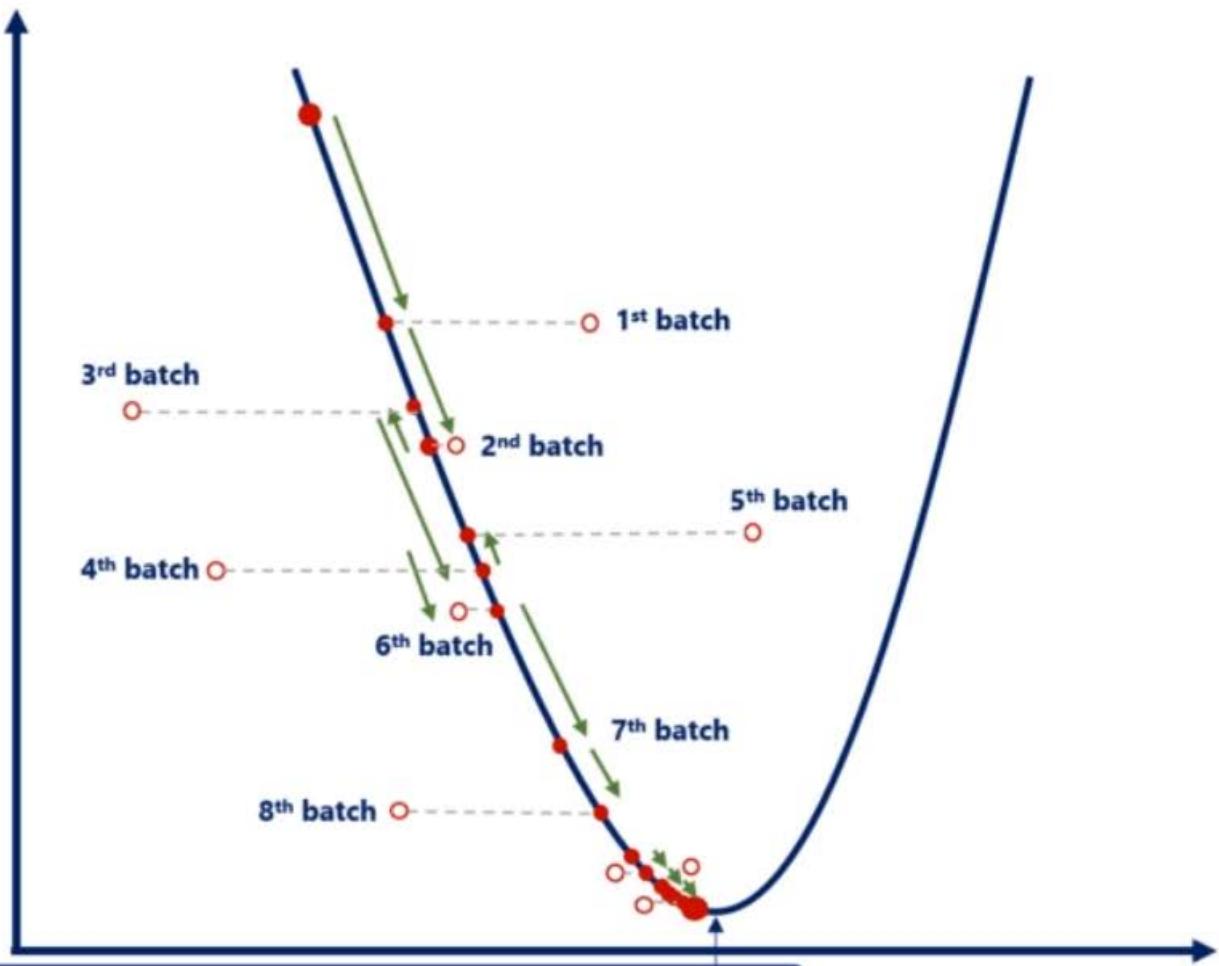


Gradient descent pitfalls



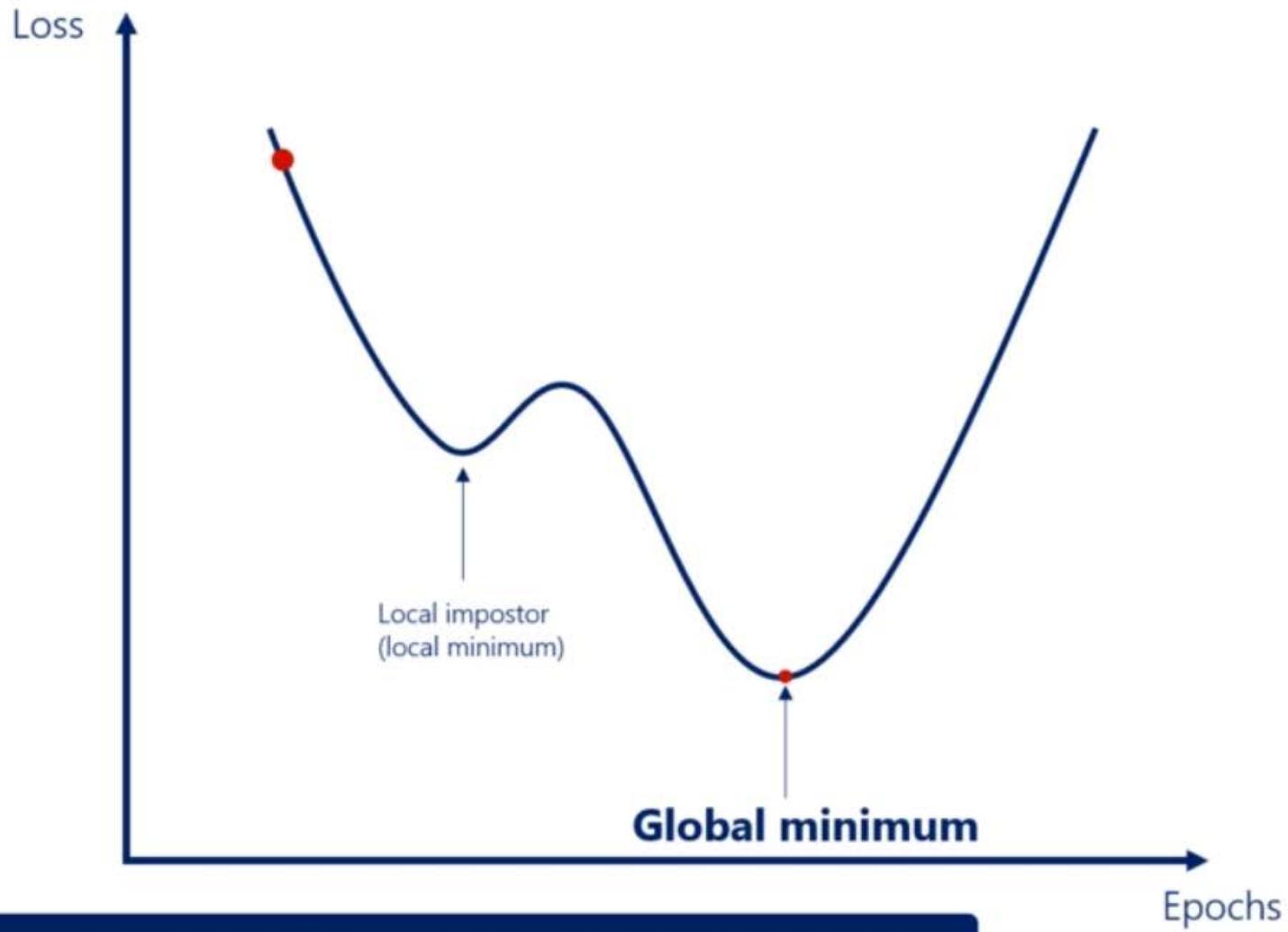
single batch GD would be slow, but will eventually reach the minimum

Gradient descent pitfalls



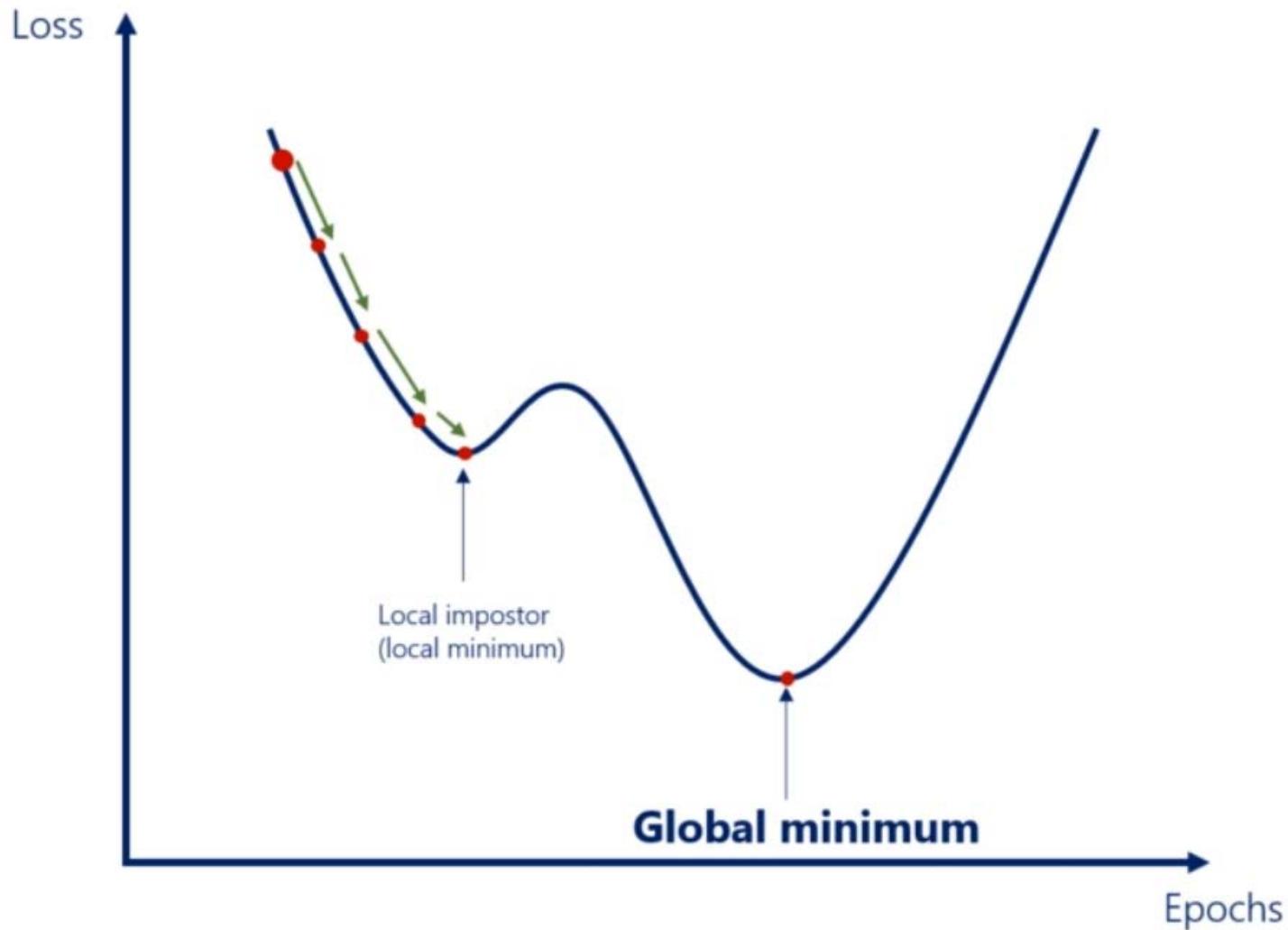
SGD would be much faster, but will give us an approximate answer

Gradient descent pitfalls

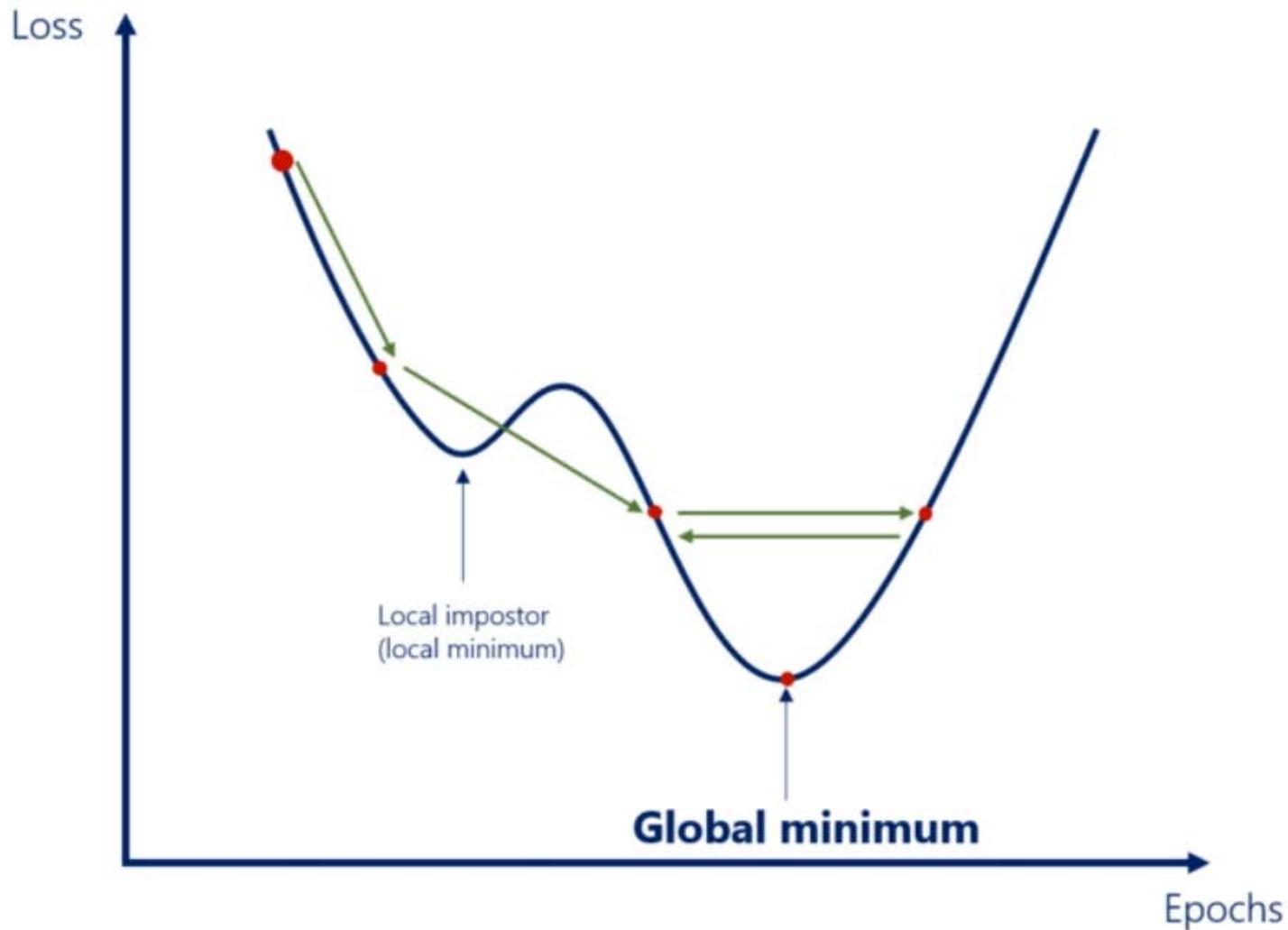


Each local minimum is a suboptimal solution to the optimization problem

Gradient descent pitfalls



Gradient descent pitfalls



MOMENTUM



MOMENTUM

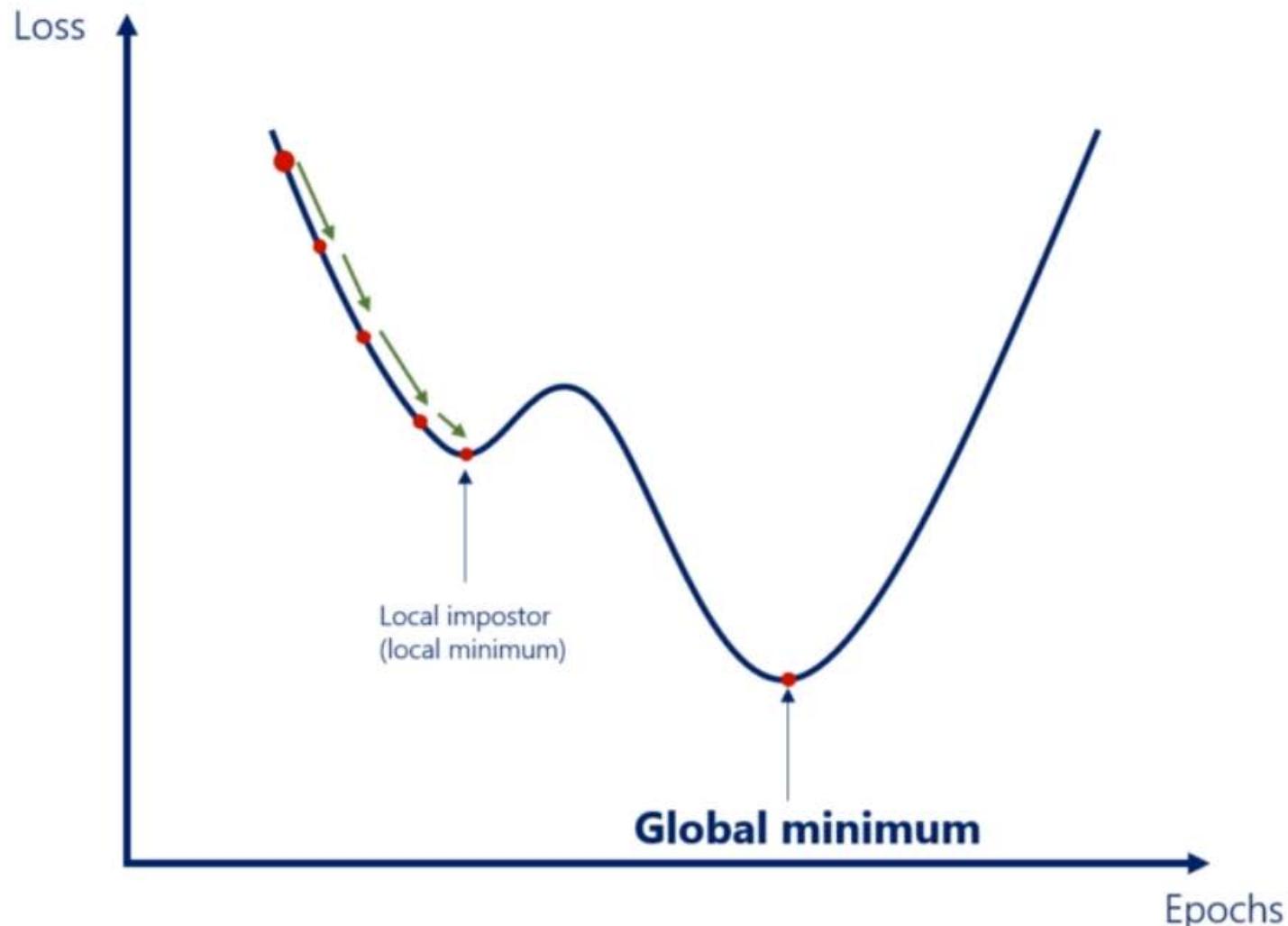


MOMENTUM

The momentum accounts for the fact that the ball actually going downhill

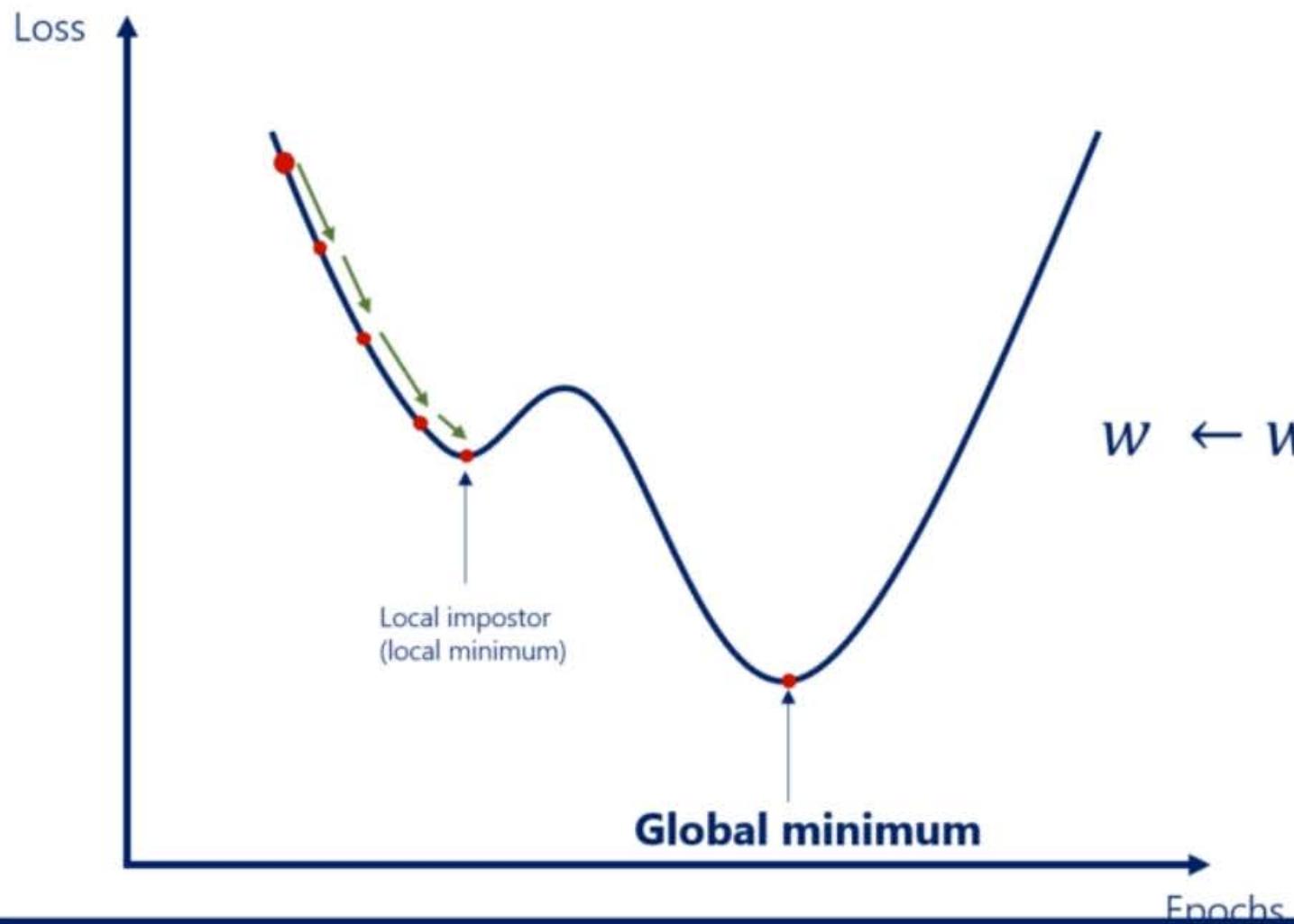


Momentum



$$w \leftarrow w - \eta \frac{\partial L}{\partial w}$$

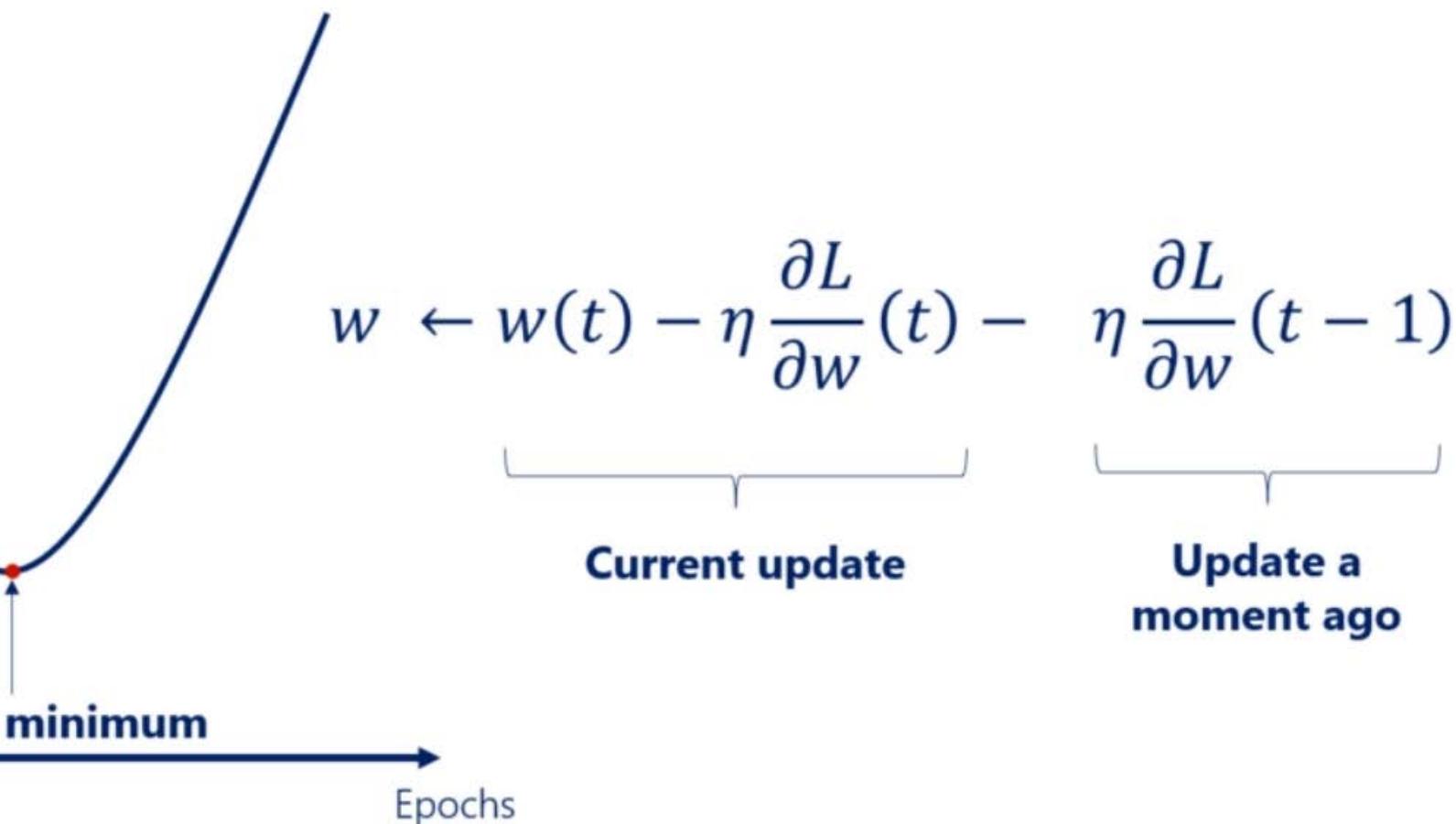
Momentum



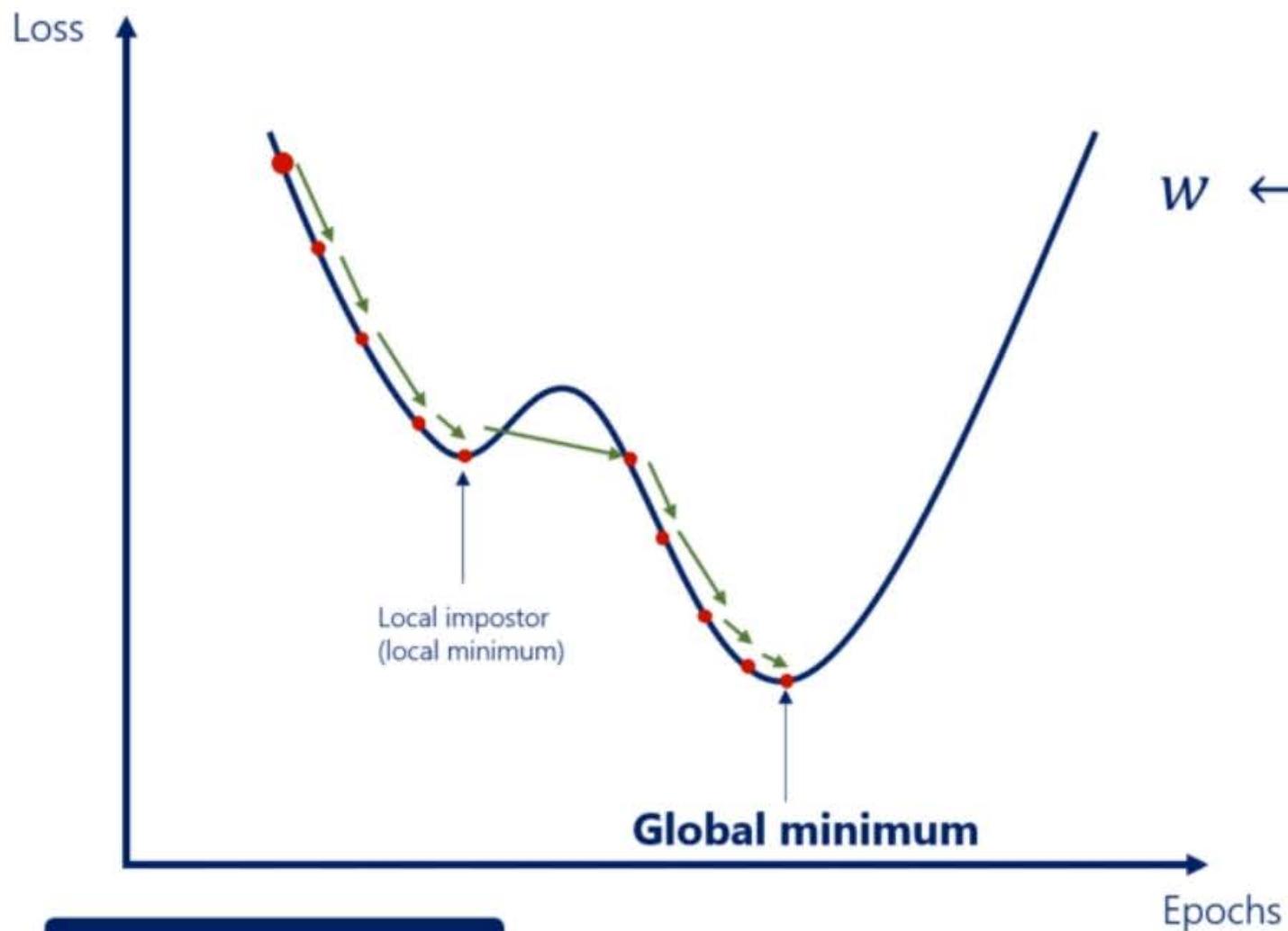
$$w \leftarrow w(t) - \eta \frac{\partial L}{\partial w}(t) - \eta \frac{\partial L}{\partial w}(t-1)$$

The best way to check how fast the ball rolls, is to check how fast it rolled a **moment** ago

Momentum



Momentum



α is a hyperparameter

$$w \leftarrow w(t) - \eta \frac{\partial L}{\partial w}(t) - \alpha \eta \frac{\partial L}{\partial w}(t-1)$$

Current update

Update a moment ago

$\alpha = 0.9$ is conventional

Hyperparameters

vs

Parameters

pre-set by us

found by optimizing

Width

Weights (w)

Depth

Biases (b)

Learning rate (η)

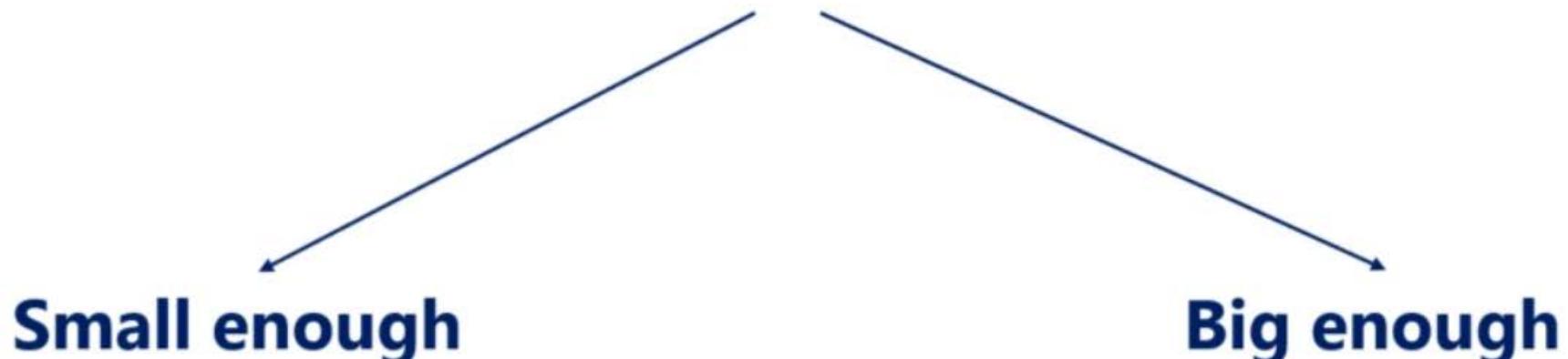
Learning rate (η)

Small enough so we gently descend, instead of oscillating or diverging

Big enough so we reach it in a rational amount of time

#datascience

Learning rate schedules



1. We start from a high initial learning rate
2. At some point we lower the rate to avoid oscillation
3. Around the end we pick a very small rate to get a precise answer

Learning rate schedules

1. We start from a high initial learning rate

First 5 epochs
 $\eta = 0.1$

2. At some point we lower the rate to avoid oscillation

Next 5 epochs
 $\eta = 0.01$

3. Around the end we pick a very small rate to get a precise answer

Until the end
 $\eta = 0.001$

Learning rate schedules

1. We start from a high initial learning rate

2. At some point we lower the rate to avoid oscillation

3. Around the end we pick a very small rate to get a precise answer

Learning rate schedules

1. We start from a high initial learning rate
2. At some point we lower the rate to avoid oscillation
3. Around the end we pick a very small rate to get a precise answer

$$\eta_0 = 0.1$$

current epoch

$$\eta = \eta_0 e^{-n/c}$$

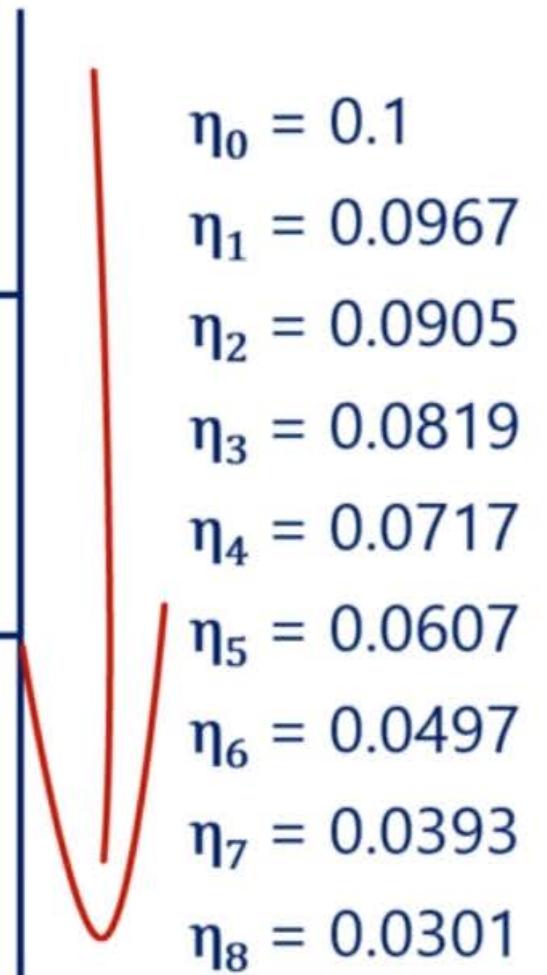
some constant

Learning rate schedules

1. We start from a high initial learning rate

2. At some point we lower the rate to avoid oscillation

3. Around the end we pick a very small rate to get a precise answer



current epoch

$$\eta = \eta_0 e^{-n/20}$$

some constant

Learning rate schedules

$$\eta_0 = 0.1$$

$$\eta_1 = 0.0967$$

$$\eta_2 = 0.0905$$

$$\eta_3 = 0.0819$$

$$\eta_4 = 0.0717$$

$$\eta_5 = 0.0607$$

$$\eta_6 = 0.0497$$

$$\eta_7 = 0.0393$$

$$\eta_8 = 0.0301$$

$$\eta = \eta_0 e^{-n/20}$$

current epoch
some constant

No set rule, but same order of magnitude

e.g. if we need:

100 epochs, $50 < c < 500$

1000 epochs, $500 < c < 5000$

$c \sim 20$ is ok

The exact value is not important. **The presence of the learning schedule does**

Hyperparameters

pre-set by us

vs

Parameters

found by optimizing

Width

Depth

Learning rate (η)

Batch size

Momentum coefficient (α)

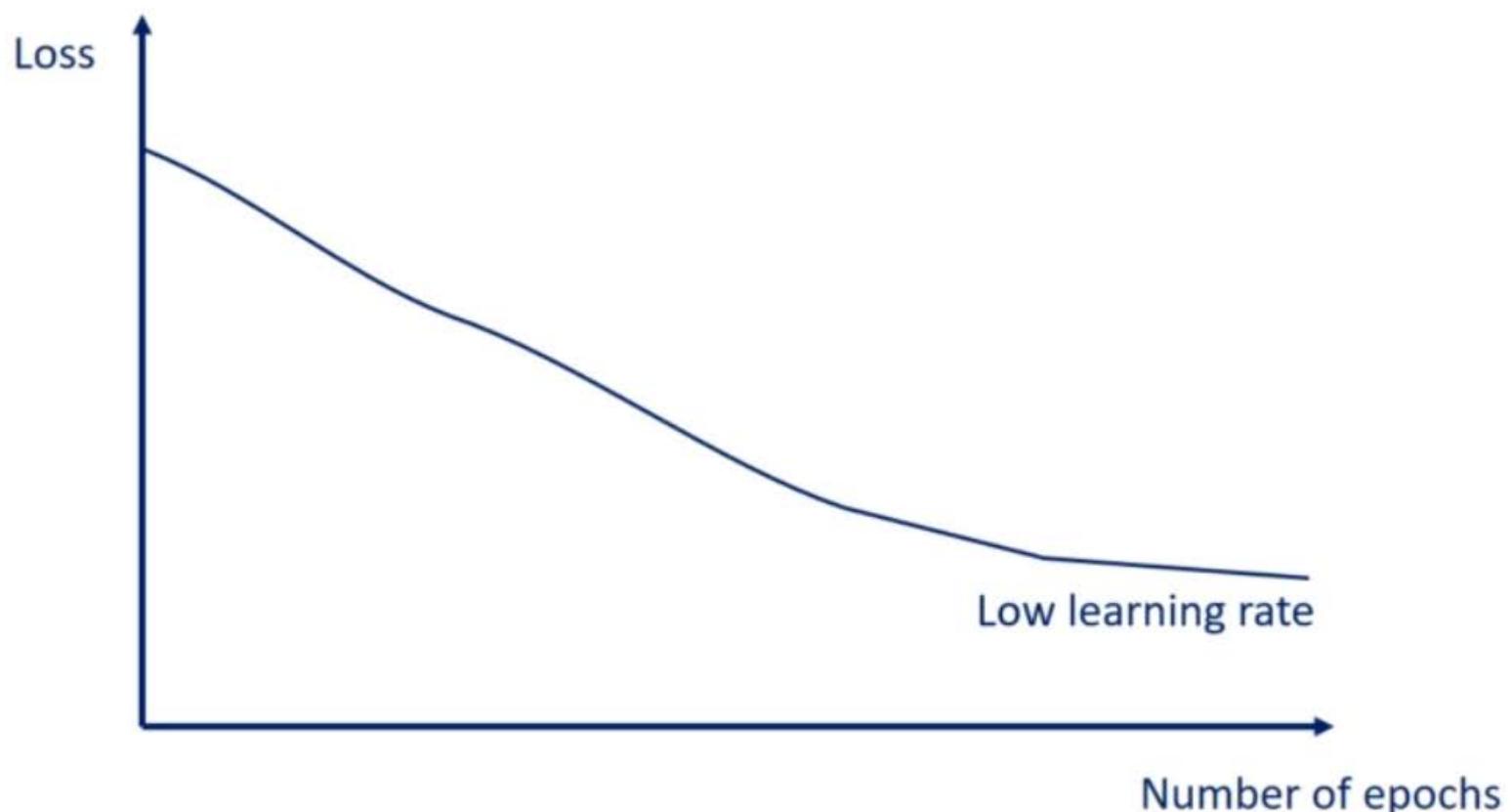
Decay coefficient (c)

Weights (w)

Biases (b)

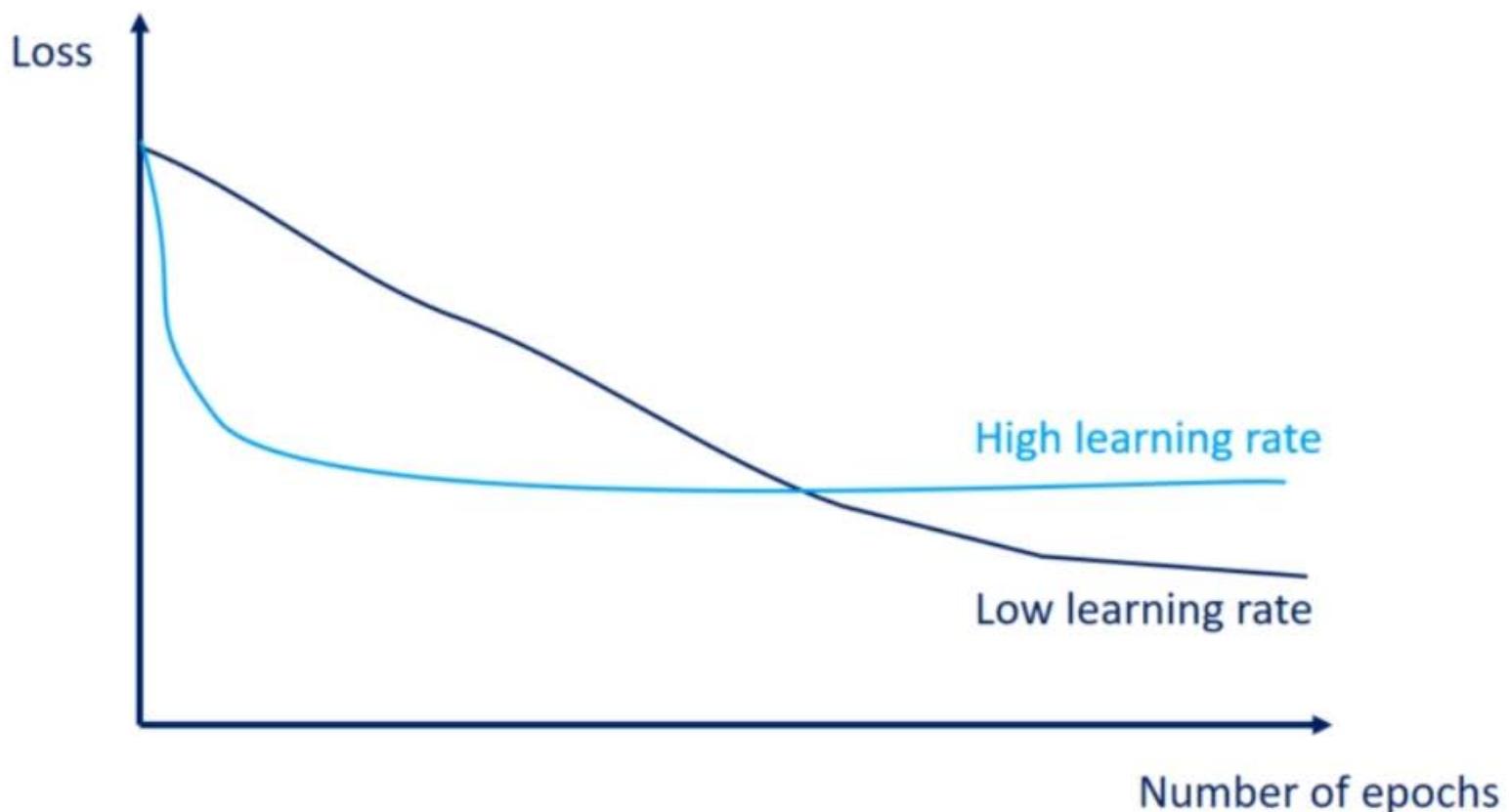
We increase the number of hyperparameters we must pick values for

Learning rate. A picture



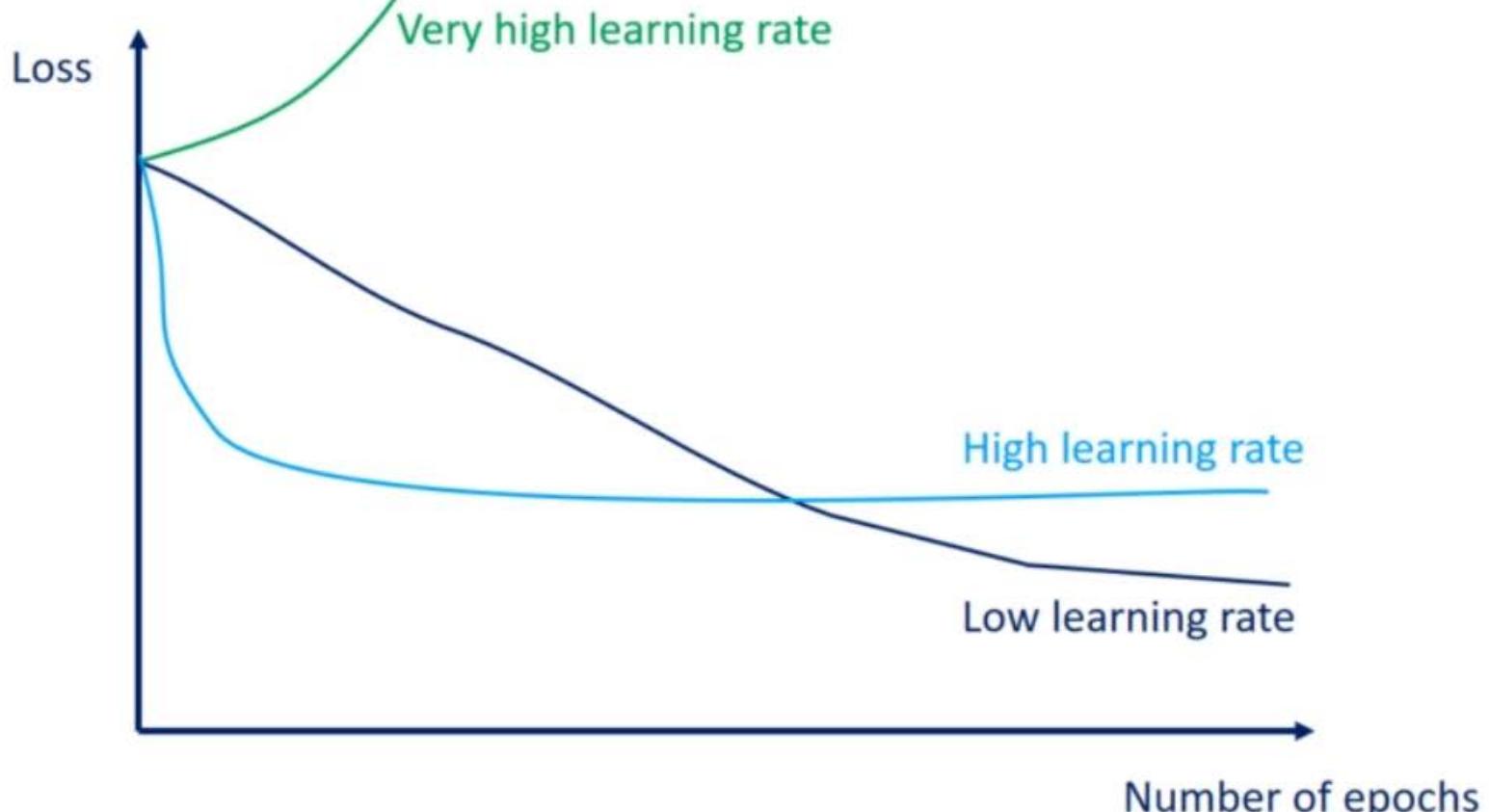
A low learning rate would minimize the loss but quite slowly

Learning rate. A picture



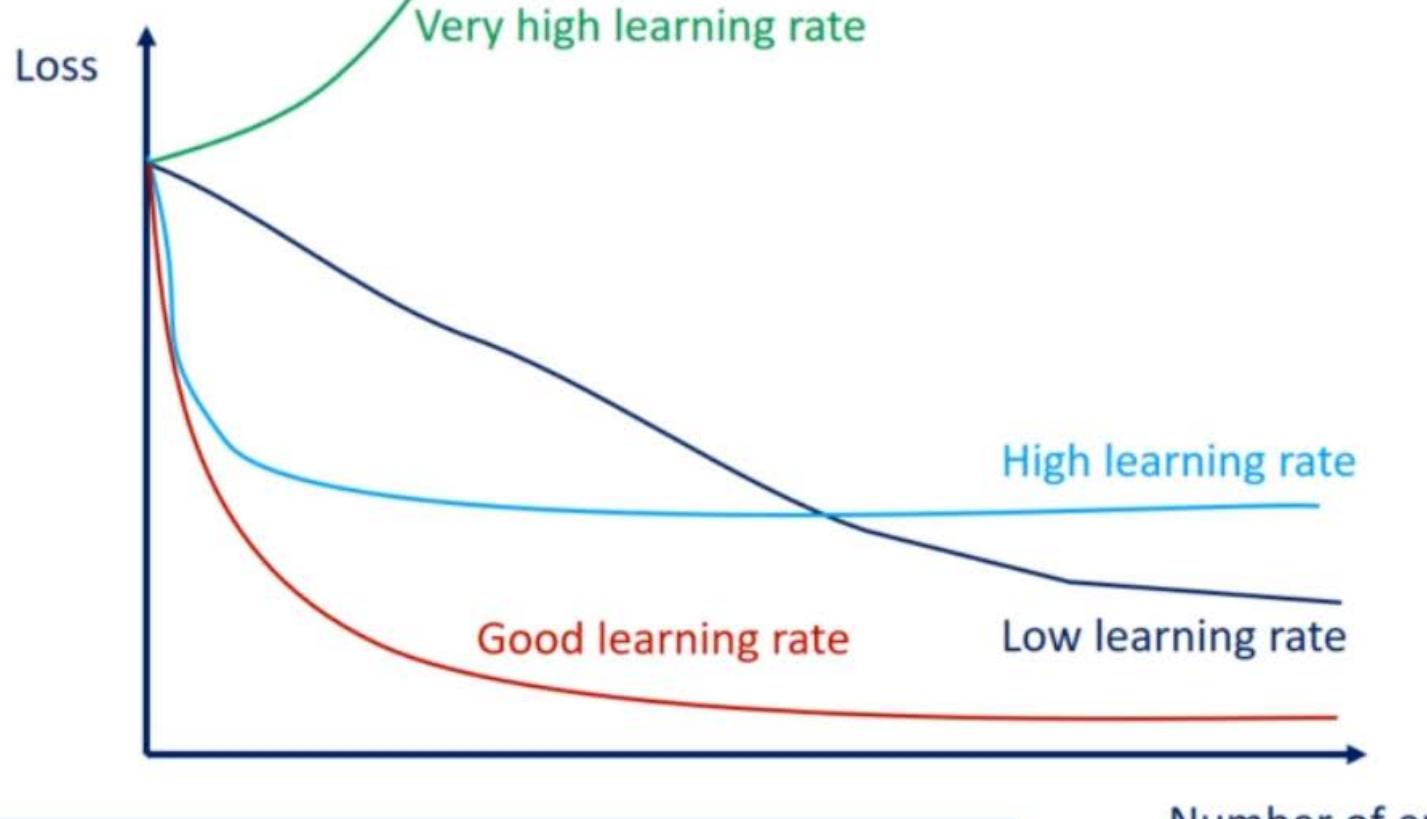
With a high learning rate, the loss is minimized faster but only to a certain extent

Learning rate. A picture



A **very** high learning rate would not even minimize the loss!

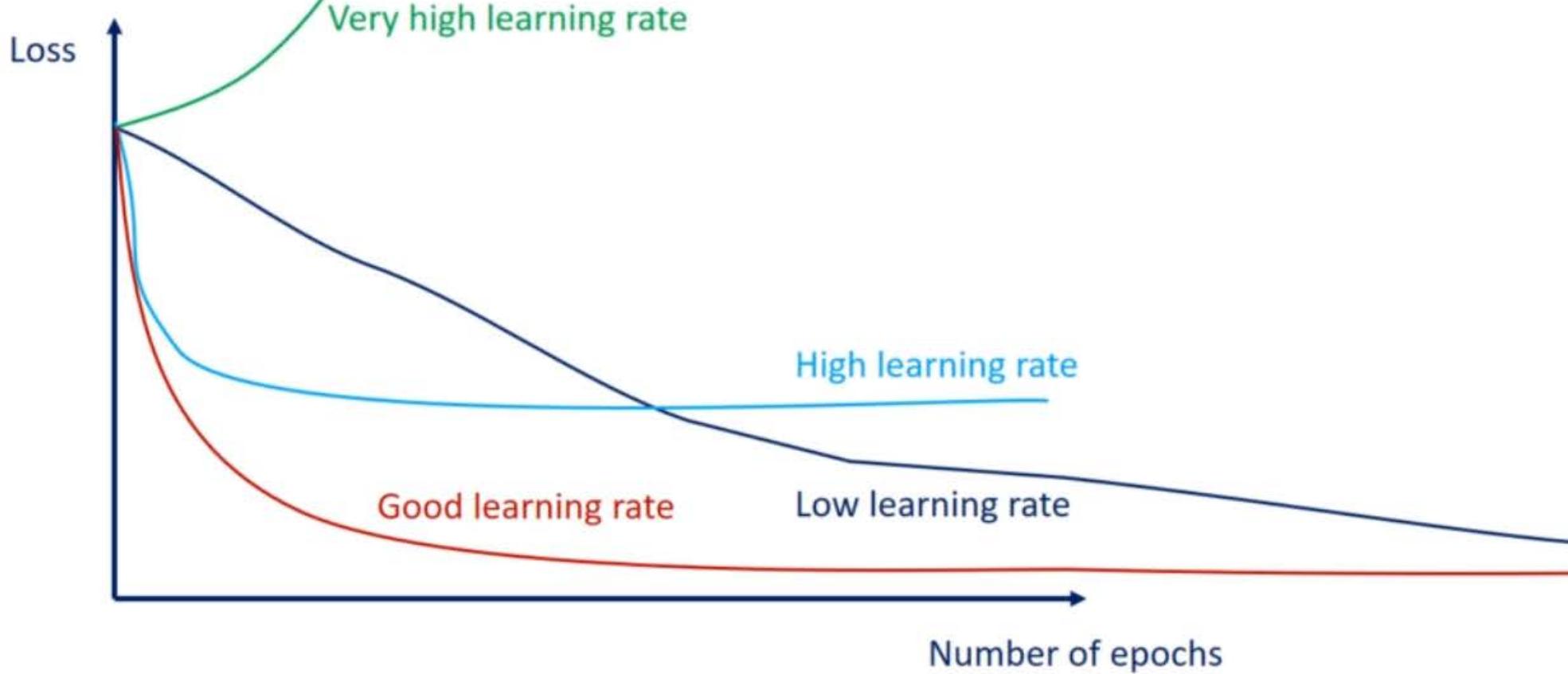
Learning rate. A picture



A learning rate following a schedule would minimize the loss faster than a low learning rate, and more accurately than a high learning rate

Number of epochs

Learning rate. A picture



A low learning rate eventually converges with the good learning rate

Learning rate schedules

AdaGrad

RMSProp



TensorFlow

Just use
AdaGrad plz

AdaGrad

/adaptive gradient algorithm/

2011

It dynamically varies the learning rate at each update and for each weight **individually**

AdaGrad

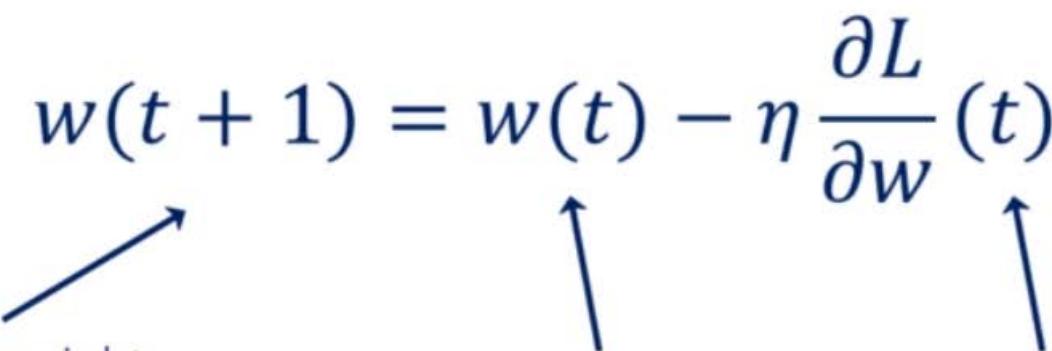
/adaptive gradient algorithm/

$$w(t + 1) = w(t) - \eta \frac{\partial L}{\partial w}(t)$$

The next weight
(at time t+1)

The previous
weight (at time t)

Follow the update
rule



AdaGrad

/adaptive gradient algorithm/

$$w(t+1) = w(t) - \eta \frac{\partial L}{\partial w}(t)$$

$$w(t+1) - w(t) = -\eta \frac{\partial L}{\partial w}(t)$$

$$\Delta w = -\eta \frac{\partial L}{\partial w}(t)$$

AdaGrad

/adaptive gradient algorithm/

$$\Delta w = -\eta \frac{\partial L}{\partial w}(t)$$

$$\Delta w_i(t) = -\frac{\eta}{\sqrt{G_i(t)} + \epsilon} \frac{\partial L}{\partial w_i}(t)$$

AdaGrad

/adaptive gradient algorithm/

$$\Delta w_i(t) = -\frac{\eta}{\sqrt{G_i(t)} + \epsilon} \frac{\partial L}{\partial w_i}(t)$$

Weight index. The update rule is individual for each weight.

Iteration (time) at which we are updating

dynamically varies the learning rate at each update and for each weight **individually**

AdaGrad

/adaptive gradient algorithm/

$$\Delta w_i(t) = -\frac{\eta}{\sqrt{G_i(t)} + \epsilon} \frac{\partial L}{\partial w_i}(t)$$

$$G_i(t) = G_i(t-1) + \left(\frac{\partial L}{\partial w_i}(t) \right)^2$$

with beginning point $G_i(0) = 0$

$$\begin{aligned} G_i(0) &= 0 \\ G_i(1) &= 0 + \text{non-neg} \\ G_i(2) &= [0 + \text{non-neg}] + \text{non-neg} \\ &\vdots \end{aligned}$$

$G(t)$ is monotonously increasing function (each consequent G is bigger or equal to the previous one)

AdaGrad

effective learning rate

$$\Delta w_i(t) = -\frac{\eta}{\sqrt{G_i(t)} + \epsilon} \frac{\partial L}{\partial w_i}(t)$$

$$G_i(t) = G_i(t-1) + \left(\frac{\partial L}{\partial w_i}(t) \right)^2$$

with beginning point $G_i(0) = 0$

/adaptive gradient algorithm/

$$G_i(0) = 0$$

$$G_i(1) = 0 + \text{non-neg}$$

$$G_i(2) = [0 + \text{non-neg}] + \text{non-neg}$$

.

.

.

AdaGrad

/adaptive gradient algorithm/

$$\Delta w_i(t) = -\frac{\eta}{\sqrt{G_i(t)} + \epsilon} \frac{\partial L}{\partial w_i}(t)$$

$$G_i(t) = G_i(t-1) + \left(\frac{\partial L}{\partial w_i}(t) \right)^2$$

with beginning point $G_i(0) = 0$

- Smart
- Adaptive learning rate schedule
- Based on the training itself
- **Per weight**

AdaGrad

/adaptive gradient algorithm/

$$\Delta w_i(t) = -\frac{\eta}{\sqrt{G_i(t)} + \epsilon} \frac{\partial L}{\partial w_i}(t)$$

$$G_i(t) = G_i(t-1) + \left(\frac{\partial L}{\partial w_i}(t) \right)^2$$

with beginning point $G_i(0) = 0$

- Smart
- Adaptive learning rate schedule
- Based on the training itself
- **Per weight**

AdaGrad

/adaptive gradient algorithm/

$$\Delta w_i(t) = -\frac{\eta}{\sqrt{G_i(t)} + \epsilon} \frac{\partial L}{\partial w_i}(t)$$

$$G_i(t) = G_i(t-1) + \left(\frac{\partial L}{\partial w_i}(t) \right)^2$$

with beginning point $G_i(0) = 0$

RMSprop

/root mean square propagation/

$$\Delta w_i(t) = -\frac{\eta}{\sqrt{G_i(t)} + \epsilon} \frac{\partial L}{\partial w_i}(t)$$

$$G_i(t) = \beta G_i(t-1) + (1-\beta) \left(\frac{\partial L}{\partial w_i}(t) \right)^2$$

with beginning point $G_i(0) = 0$

β – yet another hyperparameter
usually, around ~0.9

AdaGrad

/adaptive gradient algorithm/

$$\Delta w_i(t) = -\frac{\eta}{\sqrt{G_i(t)} + \epsilon} \frac{\partial L}{\partial w_i}(t)$$

$$G_i(t) = G_i(t-1) + \left(\frac{\partial L}{\partial w_i}(t) \right)^2$$

with beginning point $G_i(0) = 0$

RMSprop

/root mean square propagation/

$$\Delta w_i(t) = -\frac{\eta}{\sqrt{G_i(t)} + \epsilon} \frac{\partial L}{\partial w_i}(t)$$

$$G_i(t) = \beta G_i(t-1) + (1-\beta) \left(\frac{\partial L}{\partial w_i}(t) \right)^2$$

with beginning point $G_i(0) = 0$

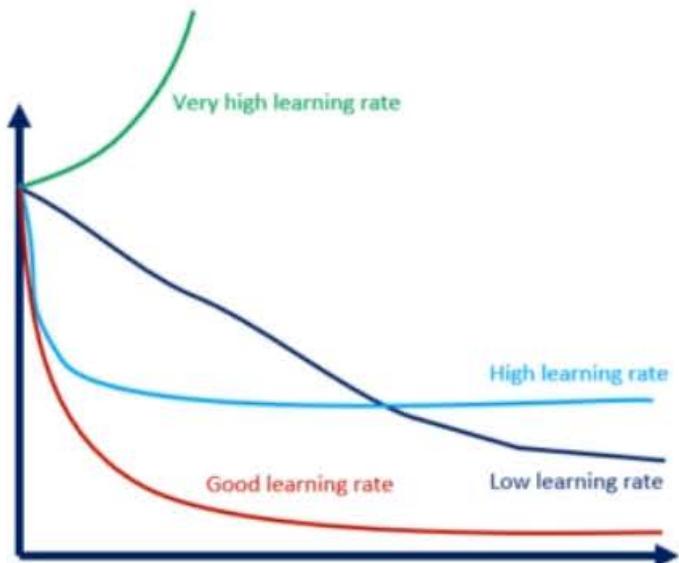
β – yet another hyperparameter
usually, around ~0.9

AdaGrad, RMSprop

$$\Delta w_i(t) = -\frac{\eta}{\sqrt{G_i(t)} + \epsilon} \frac{\partial L}{\partial w_i}(t)$$

- Smart
- Adaptive learning rate schedule
- Based on the training itself
- **Per weight**

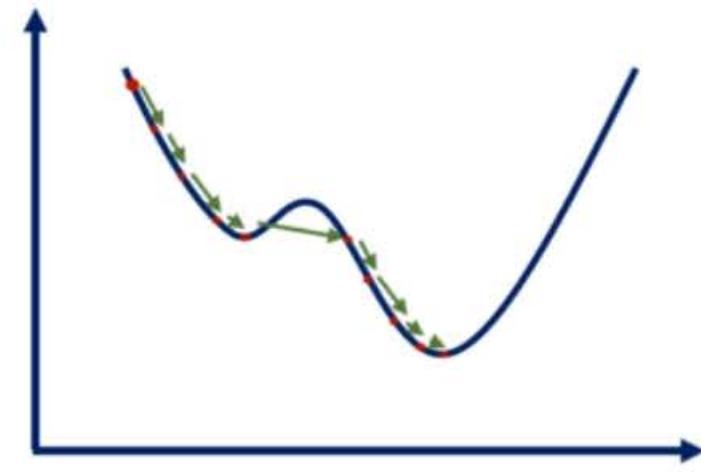
Learning rate schedules



$$\Delta w_i(t) = -\frac{\eta}{\sqrt{G_i(t)} + \epsilon} \frac{\partial L}{\partial w_i}(t)$$

AdaGrad, RMSprop

Momentum



$$w \leftarrow w(t) - \eta \frac{\partial L}{\partial w}(t) - \alpha \eta \frac{\partial L}{\partial w}(t-1)$$

Current update

Update a moment ago

Adam

/adaptive moment estimation/

~2015



The most advanced optimizer (very fast and efficient)

Adam

/adaptive moment estimation/

RMSprop

$$\Delta w_i(t) = -\frac{\eta}{\sqrt{G_i(t)} + \epsilon} \frac{\partial L}{\partial w_i}(t)$$

Momentum

$$w \leftarrow w(t) - \eta \frac{\partial L}{\partial w}(t) - \alpha \eta \frac{\partial L}{\partial w}(t-1)$$

$$\Delta w_i(t) = -\frac{\eta}{\sqrt{G_i(t)} + \epsilon} M_i(t)$$

$$M_i(t) = \alpha M_i(t-1) + (1-\alpha) \frac{\partial L}{\partial w_i}(t)$$

$$M_i(0) = 0$$

PREPROCESSING

Any manipulation of the dataset before
running it through the model



PREPROCESSING

motivation

- Compatibility



TensorFlow



PREPROCESSING

motivation

- Compatibility

Exchange rate

~ 1

Volume

$\sim 100,000$

- Orders of magnitude

Mathematically 1 is negligible w.r.t. 100,000

$$(\sim 100,000) + (\sim 1) = \sim 100,000$$

PREPROCESSING

motivation

- Compatibility

Exchange rate

~ 1

Volume

$\sim 100,000$

- Orders of magnitude

The diagram shows two vectors, x and w , being multiplied. Vector x is represented by a horizontal rectangle divided into two equal halves: the left half contains the number 1 and the right half contains the number 100000. Vector w is represented by a vertical rectangle divided into two equal halves: the top half contains the number w_1 and the bottom half contains the number w_2 . Below the multiplication sign (\times) is a large black letter w , which is identical in size to the vector w . To the right of the multiplication sign is an equals sign (=) followed by a blank rectangular box.

$$\begin{matrix} 1 & 100000 \end{matrix} \times \begin{matrix} w_1 \\ w_2 \end{matrix} = \boxed{\quad}$$
$$w$$

100,000 affects the result much more than the 1

PREPROCESSING

motivation

- Compatibility

- Orders of magnitude

- Generalization



Same model,
different issue

All

News

Images

Videos

Maps

More

Settings

Tools

About 14,000,000 results (0.69 seconds)

Apple Inc.

NASDAQ: AAPL

169.98 USD **+0.17 (0.10%)**After-hours: 170.59 **+0.61 (0.36%)**

1 day

5 day

1 month

3 months

1 year

5 years

max

After hours



Open 170.29

High 170.56

Low 169.56

Mkt cap 872.73B

P/E ratio 18.5

Div yield 1.48%

Google Finance - Yahoo Finance - MSN Money

Disclaimer

Top stories

Relative metrics are especially useful when we have time-series data

Apple

Technology company

Apple Inc. is an American multinational technology company headquartered in Cupertino, California that designs, develops, and sells consumer electronics, computer software, and online services.

Customer service: 80060027753

Headquarters: Cupertino, California, United States

Founded: April 1, 1976, Cupertino, California, United States

Founders: Steve Jobs, Steve Wozniak, Ronald Wayne

Subsidiaries: Beats Electronics, Apple Store

Did you know: Apple is the world's ninth-largest company by revenue. [wikipedia.org](#)

Profiles



YouTube



LinkedIn



Twitter

People also search for

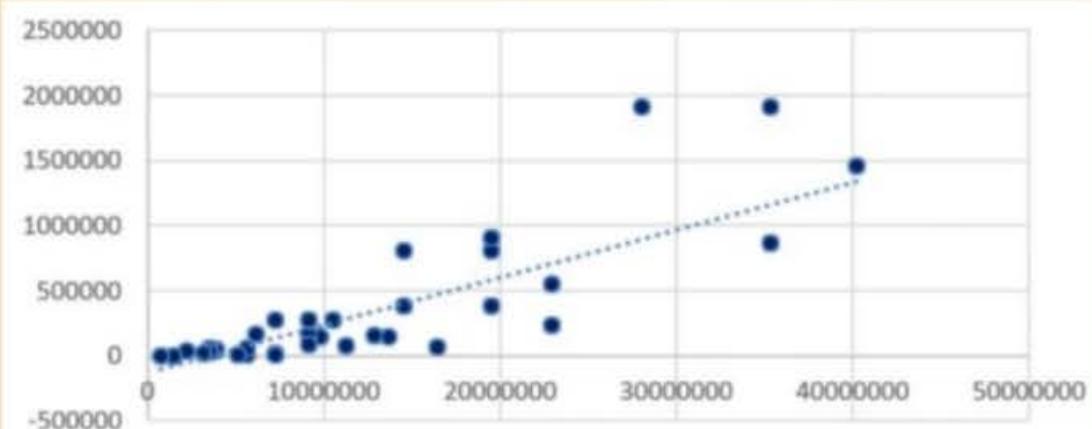


NOKIA SAMSUNG

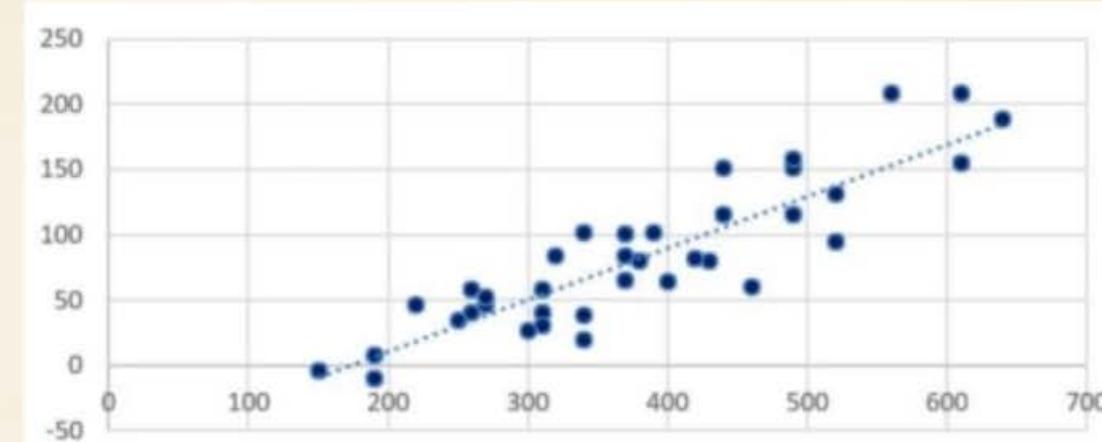
TESLA

LOGARITHMS

Financial data



Log transformed financial data



Advantages:

Faster computation

Lower order of magnitude

Clearer relationships

Homogeneous variance

STANDARDIZATION

FEATURE SCALING

NORMALIZATION

STANDARDIZATION

FEATURE SCALING

The process of transforming data
into a standard scale

STANDARDIZATION

FEATURE SCALING

original variable

standardized variable = $\frac{x - \mu}{\sigma}$

mean of original variable

standard deviation of original variable

The diagram illustrates the standardization formula $\frac{x - \mu}{\sigma}$. It features four labels: "original variable" above x , "mean of original variable" above μ , "standard deviation of original variable" below σ , and "standardized variable" to the left of the equals sign. Three arrows point from these labels to their corresponding terms in the formula: one arrow points from "original variable" to x , another from "mean of original variable" to μ , and a third from "standard deviation of original variable" to σ .

AACI	60.754	0.00%	0.00%	0.394	0.88%	16,319	0.00
ABCO	4.344	0.30%	0.57%	0.204	0.81%	NA	0.00
AMCI	55.543	0.13%	0.00%	0.394	0.88%	NA	0.00
ANOT	83.486	0.00%	0.00%	0.394	0.88%	NA	0.00
ANOT	33.364	0.34%	0.02%	0.394	0.88%	NA	0.00
ANOT	12.334	0.00%	0.00%	0.394	0.88%	NA	0.00
ANOT	32.795	0.00%	0.00%	0.394	0.88%	NA	0.00
ANOT	21.734	0.34%	0.02%	0.394	0.88%	NA	0.00

STANDARDIZATION

FEATURE SCALING

day

Exchange rate

Daily trading volume

1

1.3

110,000

2

1.34

98,700

3

1.25

135,000

mean: 1.3

std: 0.045

	60.754	-0.002	-0.001	0.308	3.03%	16,319	3.00
ABC	4.384	0.204	0.071	0.304	3.03%	NA	0.00
ABCD	33.543	0.134	0.007	0.304	3.03%	NA	0.00
ADT	83.446	0.008	0.001	0.304	3.03%	NA	0.00
ADTCH	32.384	0.384	2.03%	0.304	3.03%	NA	0.00
ADTCH	12.334	0.008	0.001	0.304	3.03%	NA	0.00
ADTCH	32.795	0.008	0.001	0.304	3.03%	NA	0.00
ADTCH	21.734	0.384	2.03%	0.304	3.03%	NA	0.00

STANDARDIZATION

FEATURE SCALING

$$\frac{x - \mu}{\sigma}$$

day

Exchange rate

Daily trading volume

1

0.07

(1.3)

110,000

2

0.96

(1.34)

98,700

3

-1.03

(1.25)

135,000

JPM	40.754	+0.08	+0.20%	0.394	+0.03%	16.310	0.00
MSFT	4.344	+0.04	+0.93%	0.394	+0.07%	NA	0.00
AMZN	95.543	+0.16	+0.17%	0.394	+0.03%	NA	0.00
GOOG	13.446	-0.09	-0.67%	0.394	+0.03%	NA	0.00
FB	12.394	+0.04	+0.33%	0.394	+0.03%	NA	0.00
DIS	13.324	-0.08	-0.59%	0.394	+0.03%	NA	0.00
AMC	32.769	-0.08	-0.24%	0.394	+0.03%	NA	0.00
AMZN	91.734	+0.04	+0.05%	0.394	+0.03%	NA	0.00

STANDARDIZATION

FEATURE SCALING

$$\frac{x - \mu}{\sigma}$$

day

Exchange rate

Daily trading volume

1	0.07	(1.3)	-0.25	(110,000)
2	0.96	(1.34)	-0.85	(98,700)
3	-1.03	(1.25)	1.1	(135,000)



STANDARDIZATION

FEATURE SCALING

day

Exchange rate

Daily trading volume

1

0.07

(1.3)

-0.25

(110,000)

2

0.96

(1.34)

-0.85

(98,700)

3

-1.03

(1.25)

1.1

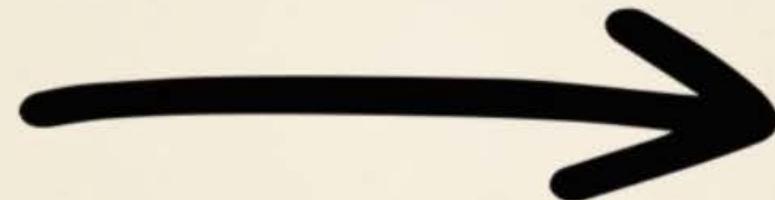
(135,000)

We have forced the features to appear similar

NORMALIZATION

5	100
120	-10
-1	0

normalize



using
L2-norm

0.042	0.995
0.999	-0.1
-0.008	0

PCA

/principal components analysis/

Dimension reduction technique used to combine several variables into a bigger (latent) variable

PCA

/principal components analysis/

Religion



voting history



participation in
associations



upbringing



attitude
towards
immigration

mean = 0

std = 1

STANDARDIZATION

NORMALIZATION

PCA

WHITENING

...

CATEGORICAL DATA



Categories or groups (non-numerical data)



numbers



numbers

CATEGORICAL DATA

How to encode categories in a way useful for ML



One-hot encoding



Binary encoding

ONE-HOT ENCODING

	bread	yogurt	muffins
	1	0	0
	0	1	0
	0	0	1

BINARY ENCODING

	ordinal	binary
	1	01
	2	10
	3	11

BINARY ENCODING

	ordinal	var1	var2
	1	0	1
	2	1	0
	3	1	1

BINARY ENCODING

	var 1	var 2
bread	0	1
yogurt	1	0
cupcake	1	1

By splitting the numbers in two variables we have removed the order

...but...

there are some implied correlation between them

BINARY ENCODING

	var 1	var 2	
	0	1	Bread is the opposite of yogurt
	1	0	Whatever is bread is not yogurt and vice versa
	1	1	

12,000
products

	ONE-HOT	BINARY
NEW COLUMNS	12,000	16

12,000 IN BINARY: 10111011100000

ONE-HOT



FEW
CATEGORIES

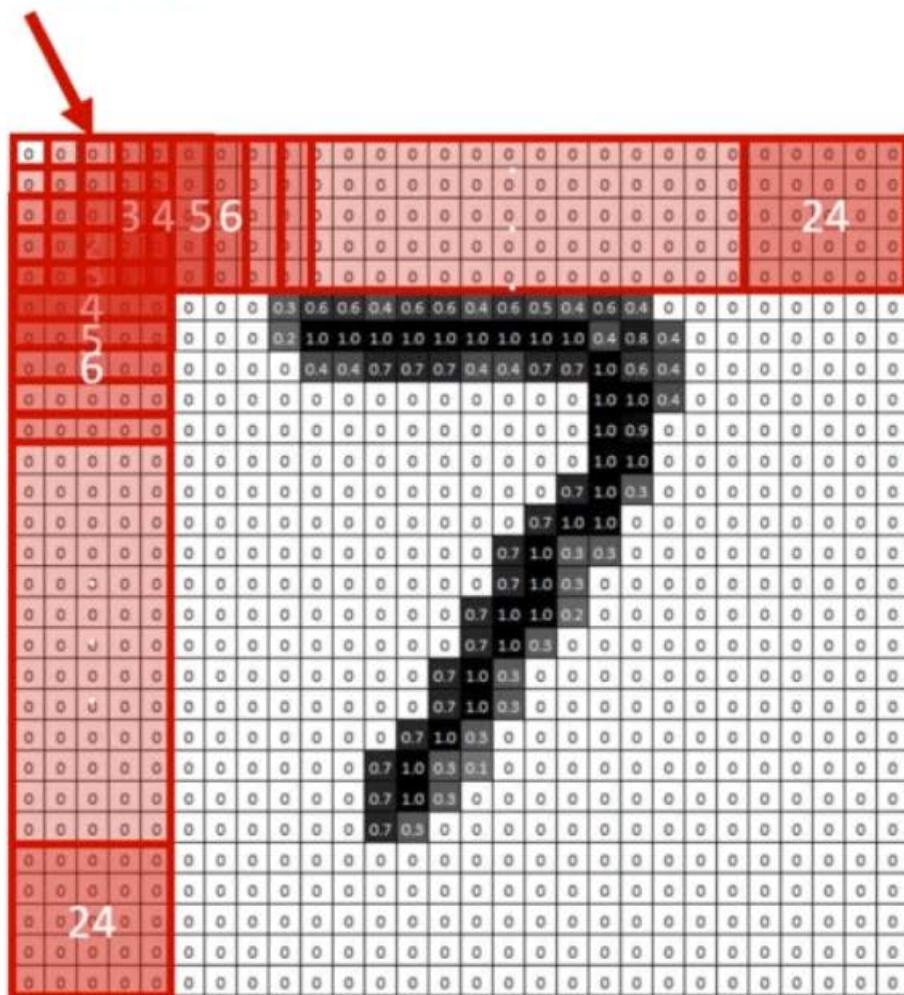
BINARY



MANY
CATEGORIES

Convolutional Neural Networks (CNNs)

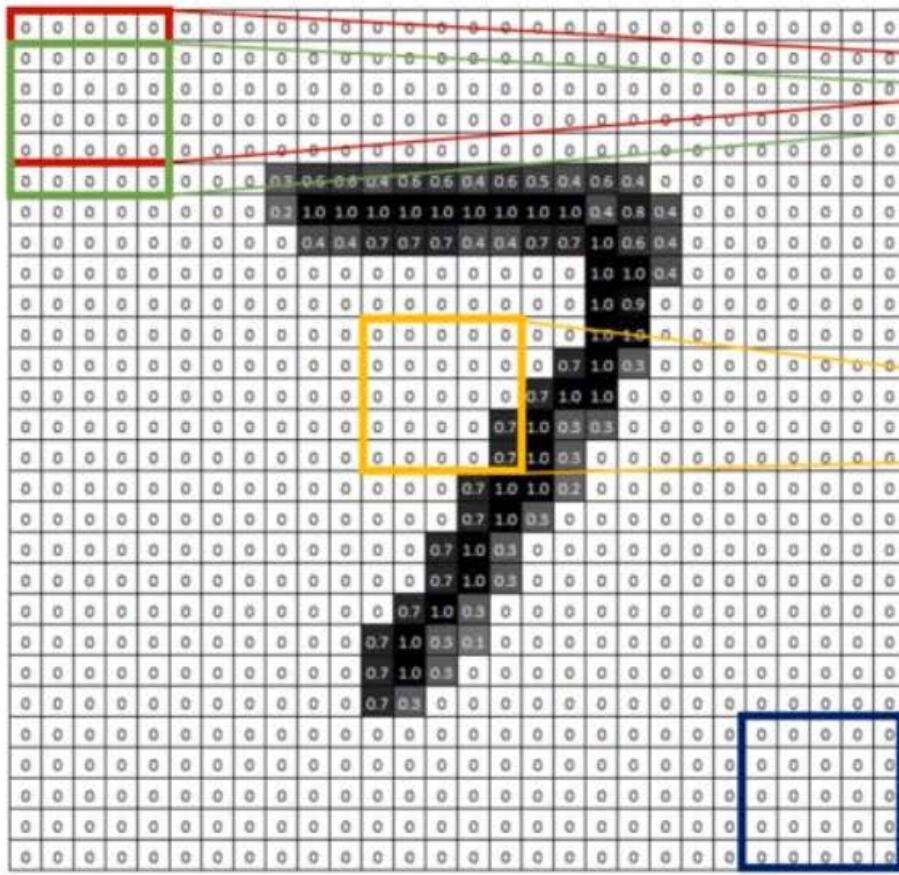
A 5×5 kernel



28x28

Convolutional Neural Networks (CNNs)

Input layer



Convolutional layer

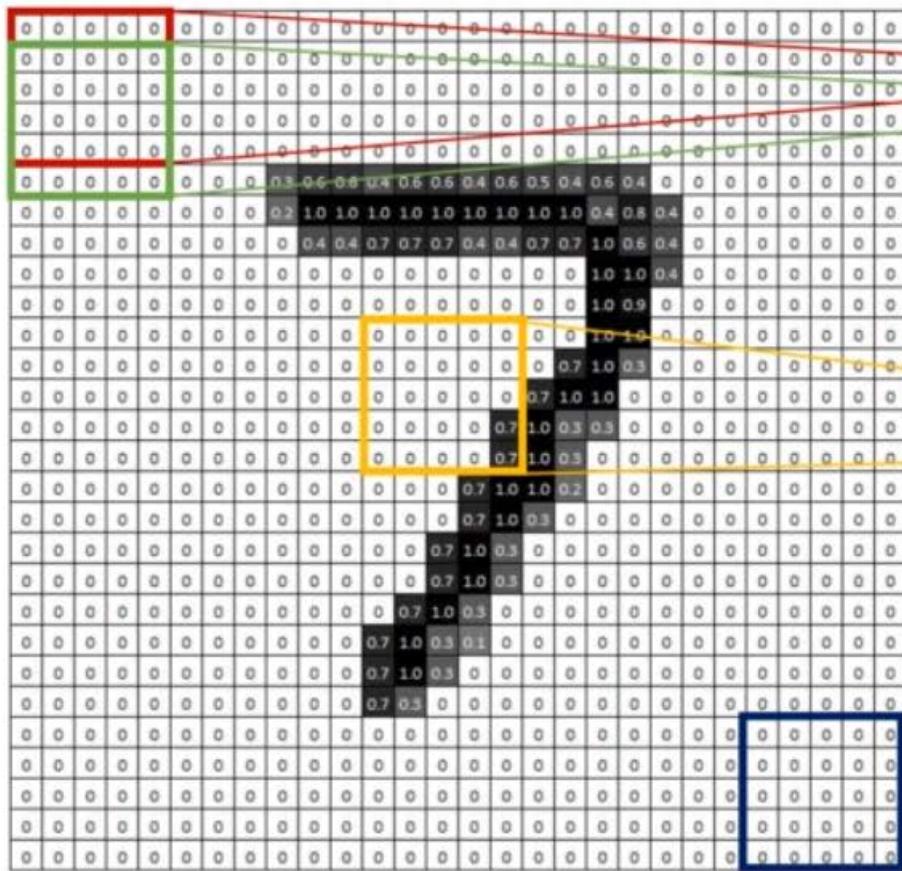
24x24

28x28

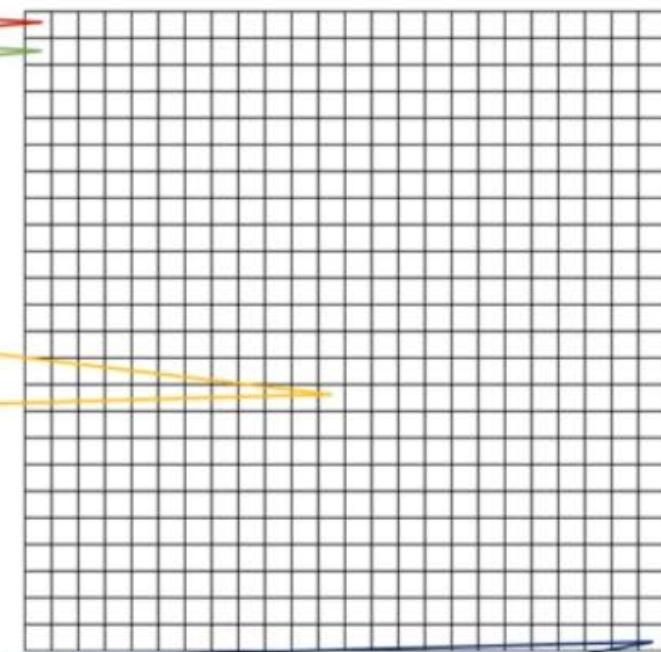
This is called **convolution**

Convolutional Neural Networks (CNNs)

Input layer



Convolutional layer



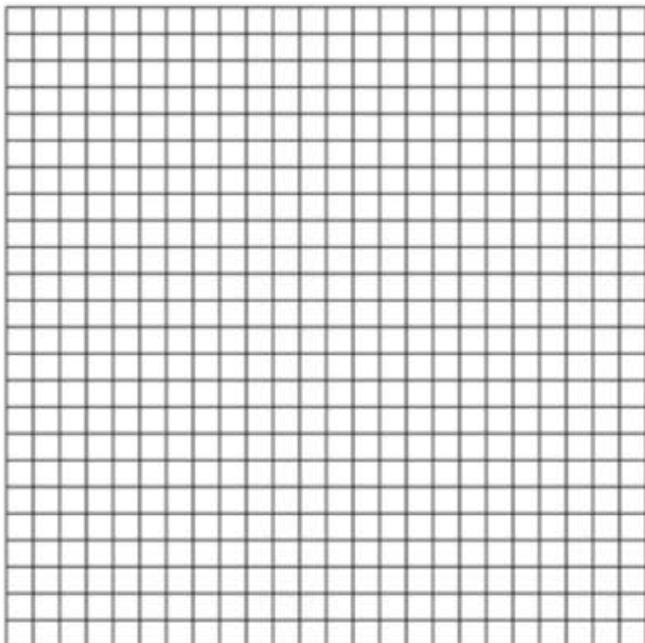
24x24

28x28

The number of kernels you choose is a hyperparameter

Convolutional Neural Networks (CNNs)

Convolutional layer

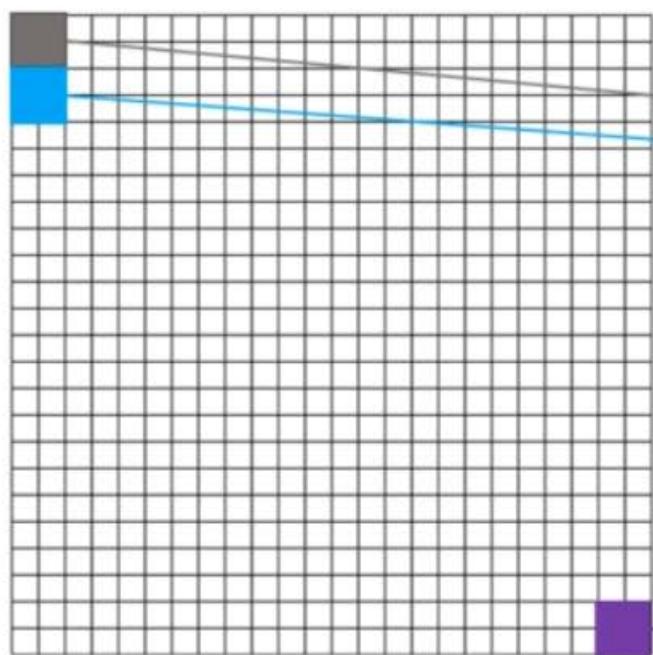


24x24

Apart from convolution there is another main step - **pooling**

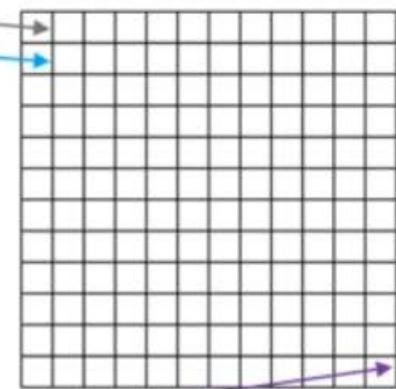
Convolutional Neural Networks (CNNs)

Convolutional layer



24x24

Pooling layer

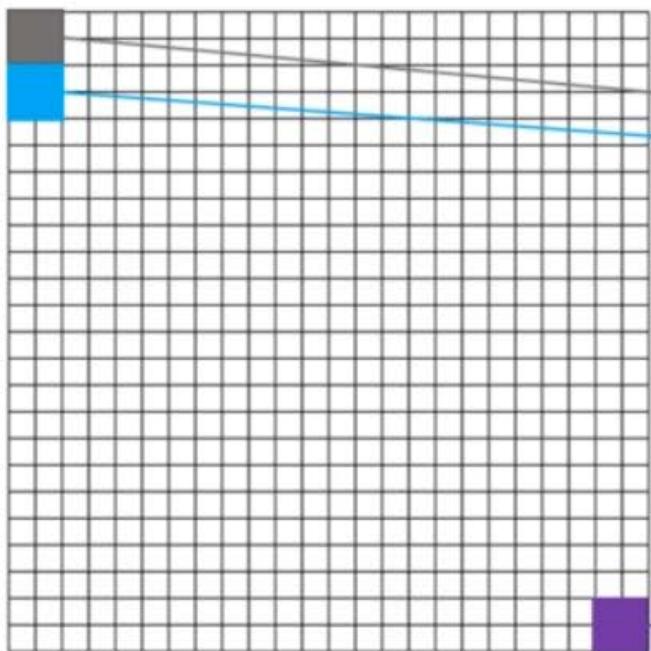


12x12

We would divide the convolutional layer into small squares **without overlapping**

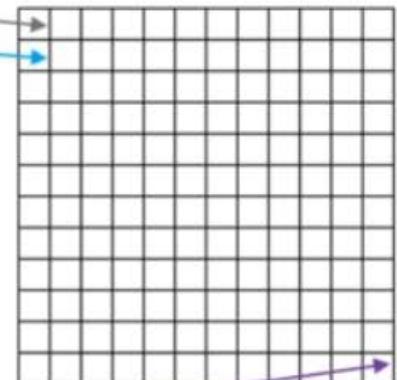
Convolutional Neural Networks (CNNs)

Convolutional layer



24x24

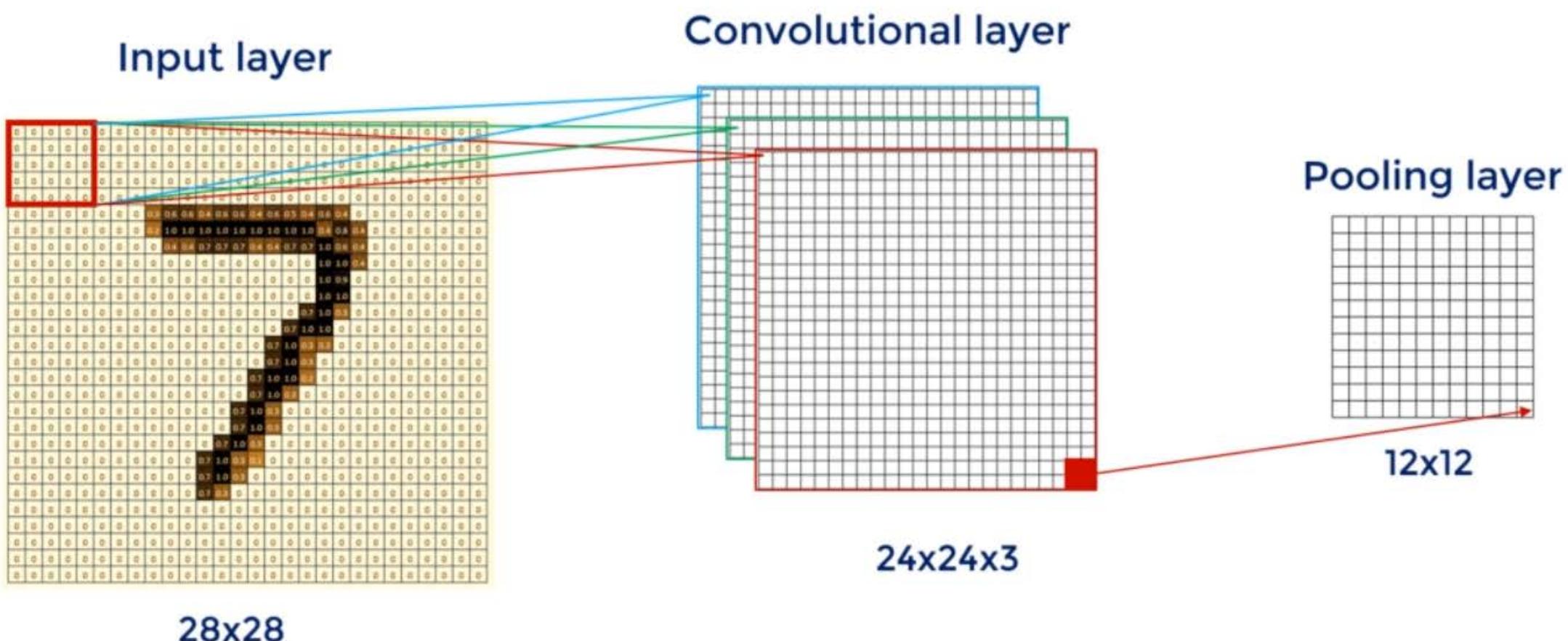
Pooling layer



12x12

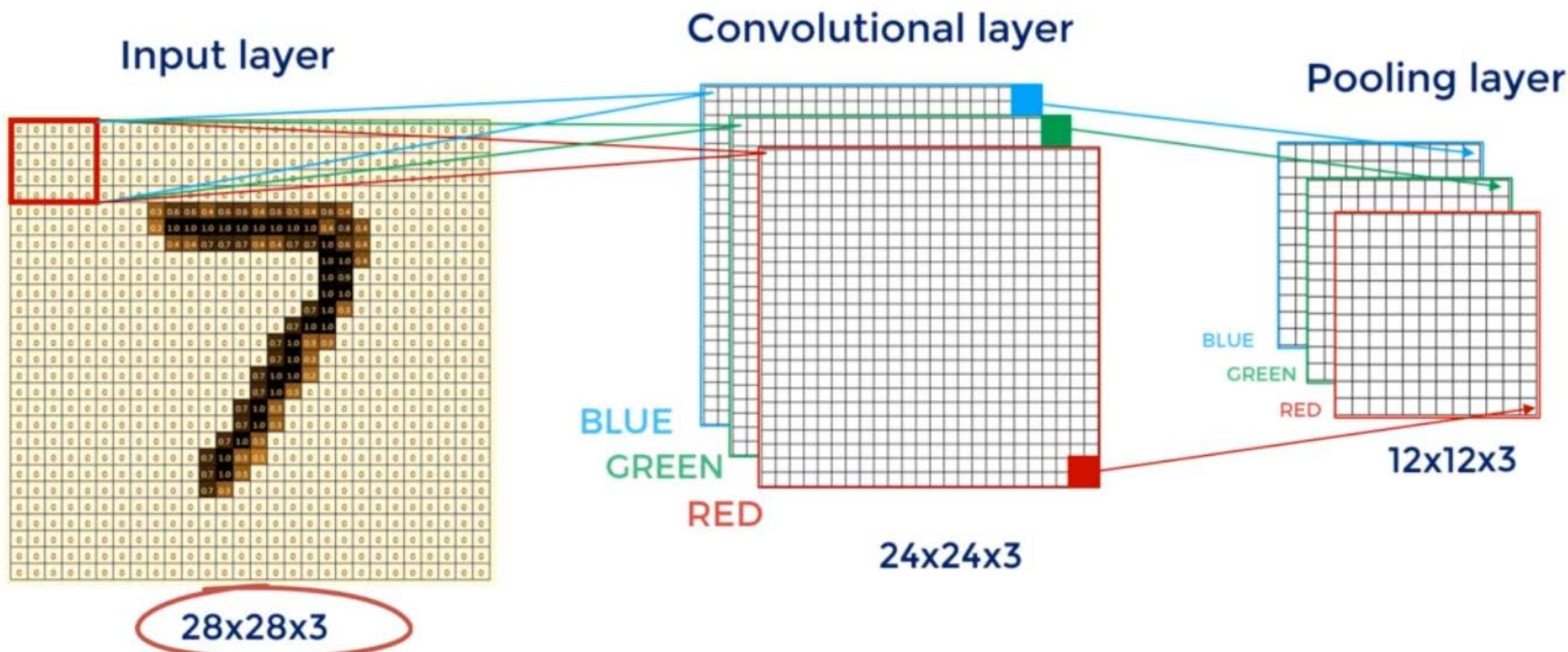
Normally we will take (and keep) the largest number from the square as it is the **strongest detail**

Convolutional Neural Networks (CNNs)



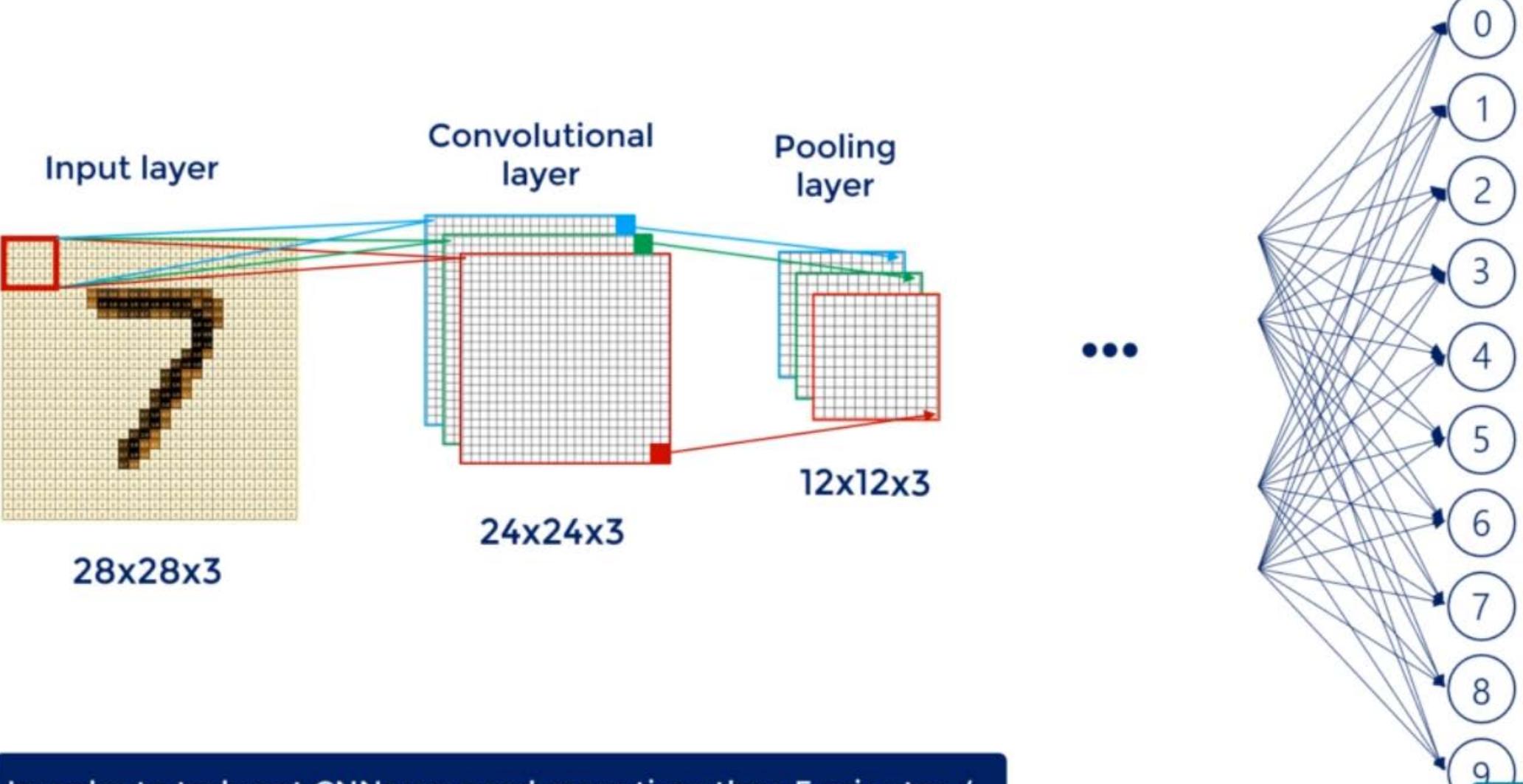
Images have height and width, but also... color (**depth**)

Convolutional Neural Networks (CNNs)



That's where the whole **tensor approach** fits perfectly

Convolutional Neural Networks (CNNs)



CNNs

- spatial proximities
are preserved
- a detail is looked for everywhere in a photo



CNNs



Image recognition



RNNs



sequential data

trading

ABCD	▼	60.754	0.304	0.03%	0.304	2.03%	16.310	0.00
ABCD	▲	4.344	0.304	0.87%	0.304	0.87	N/A	0.00
ABCD	▲	55.543	0.130	0.50%	0.304	2.03%	N/A	0.00
ABCD	▼	63.446	0.304	0.03%	0.304	2.03%	N/A	0.00
ABCD	▲	32.304	0.304	2.03%	0.304	2.03%	N/A	0.00
ABCD	▼	12.324	0.304	0.03%	0.304	2.03%	N/A	0.00
ABCD	▼	32.765	0.304	0.03%	0.304	2.03%	N/A	0.00
ABCD	▲	21.734	0.304	2.03%	0.304	2.03%	N/A	0.00

music

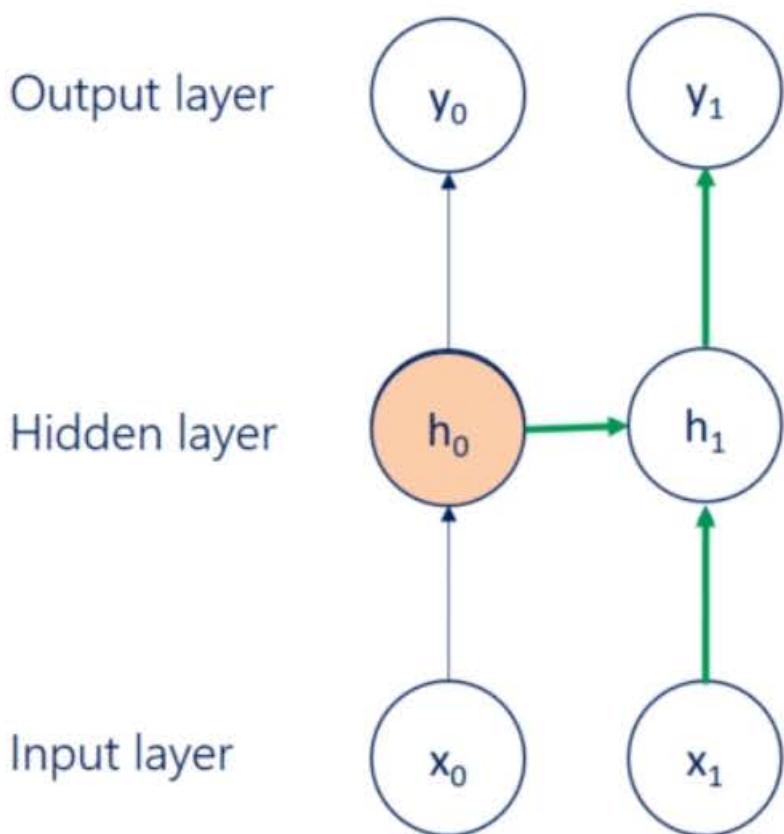


NLP



depends on what you've
said before

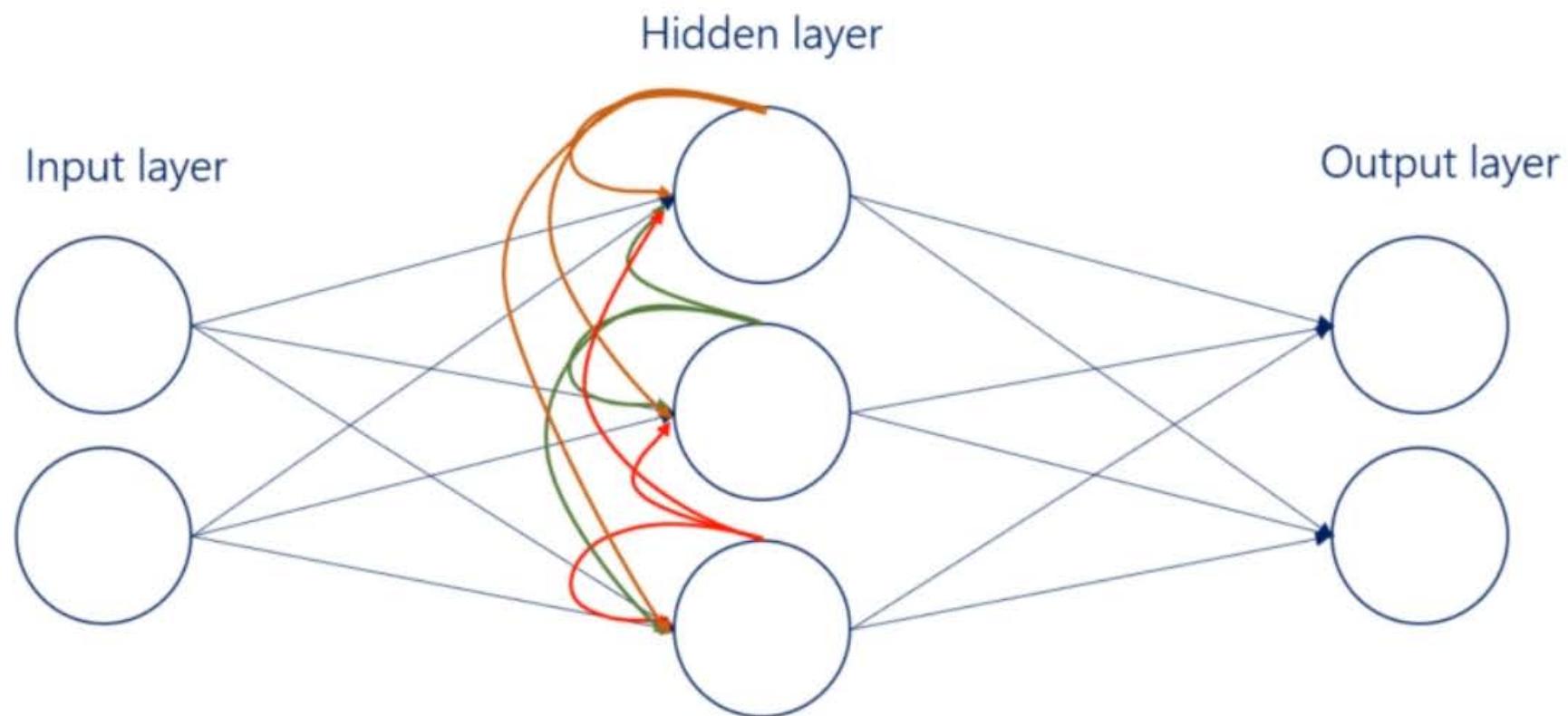
Recurrent neural networks (RNNs)



The hidden layer for the second sample is a function of the new inputs, the updated weights, **and the hidden layer from the previous iteration.**

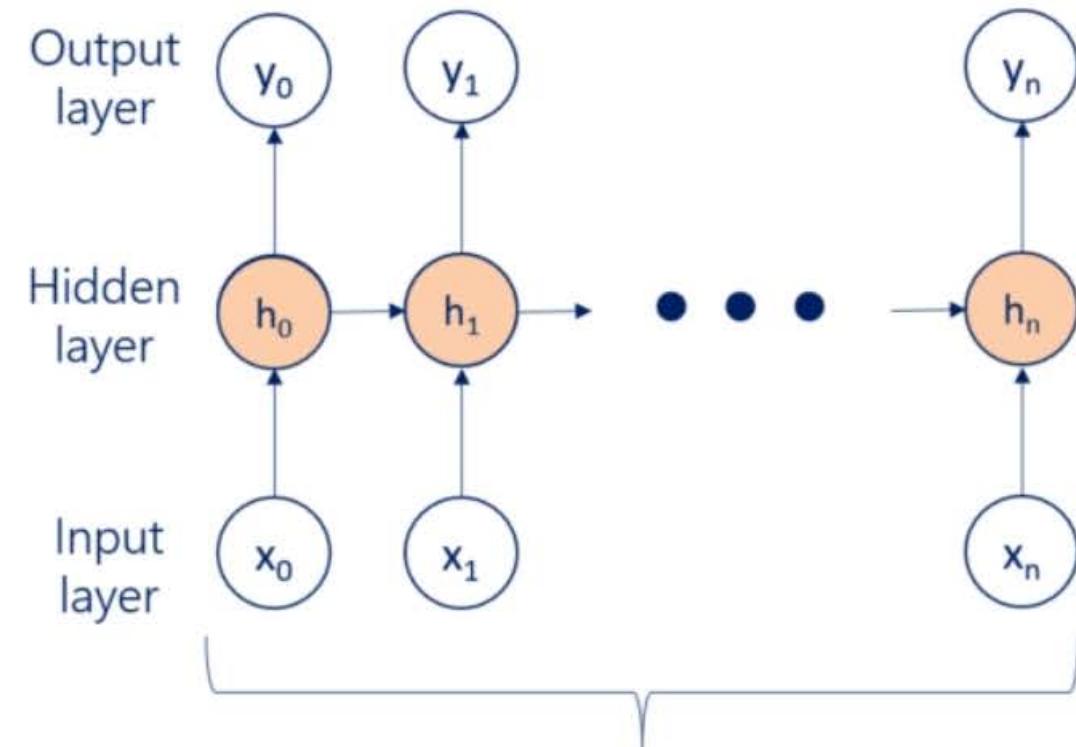
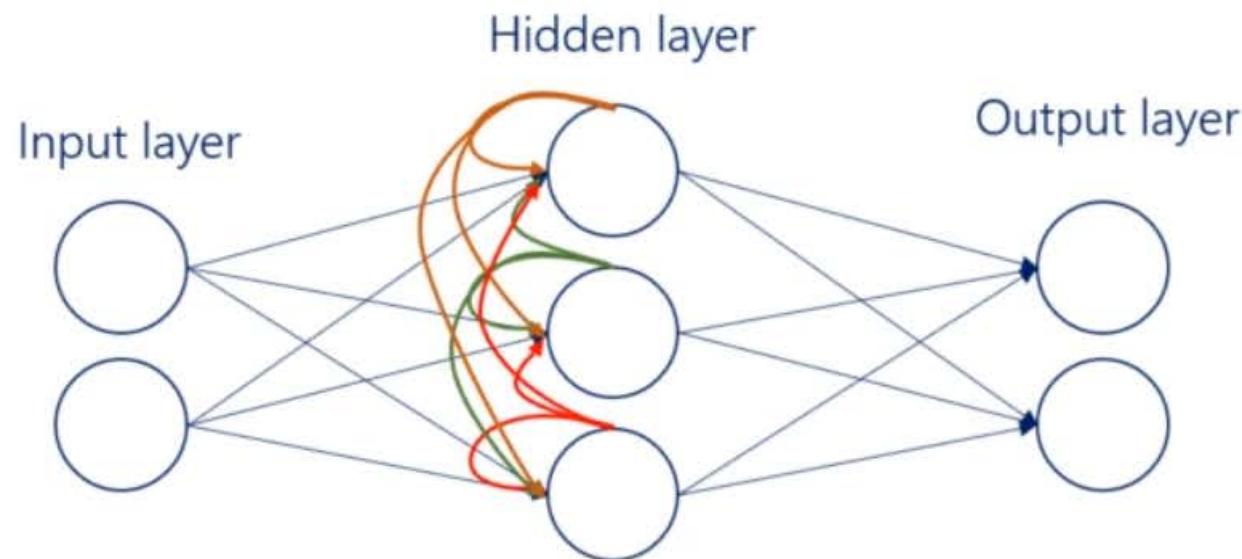
*Indices here indicate time, rather than a specific unit. An index of 0, implies that layer at time 0.

Recurrent neural networks (RNNs)



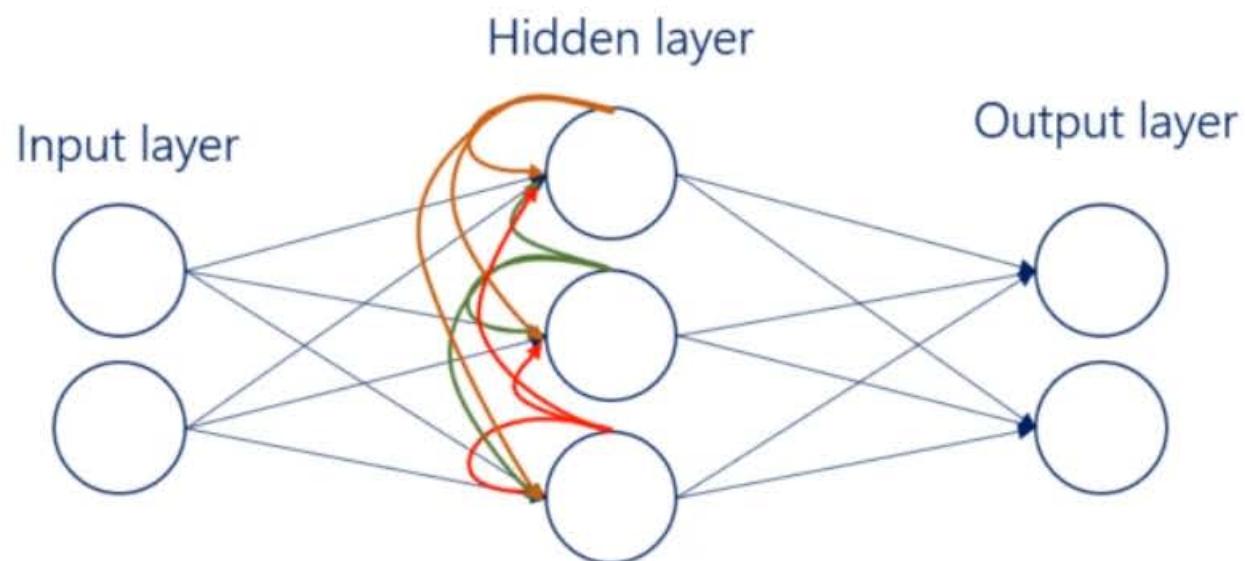
There are additional weights that lead from the hidden units to **themselves**

Recurrent neural networks (RNNs)

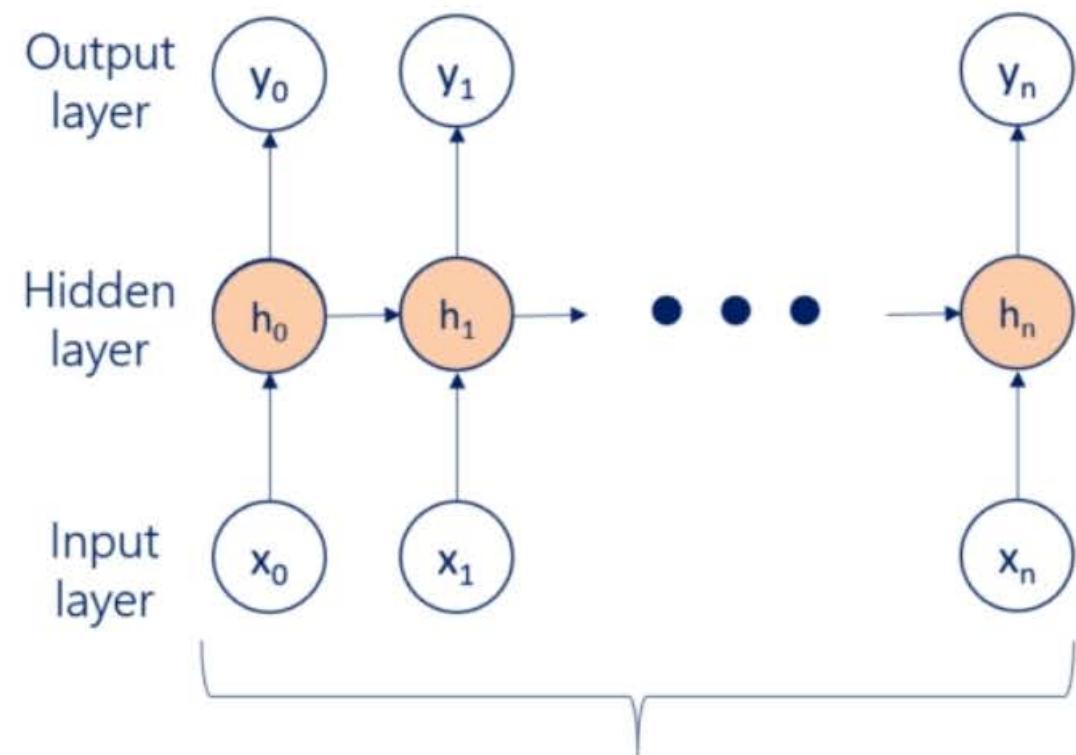


An extremely deep feedforward NN

Recurrent neural networks (RNNs)



“unrolling” the RNN



An extremely deep feedforward NN

RNNs are very computationally expensive

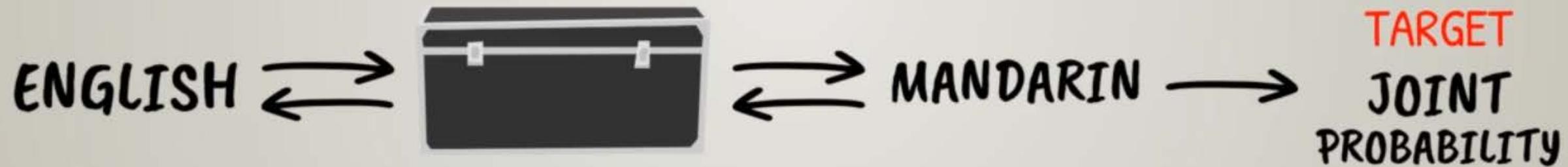


MACHINE LEARNING

NON-NEURAL NETWORK APPROACHES

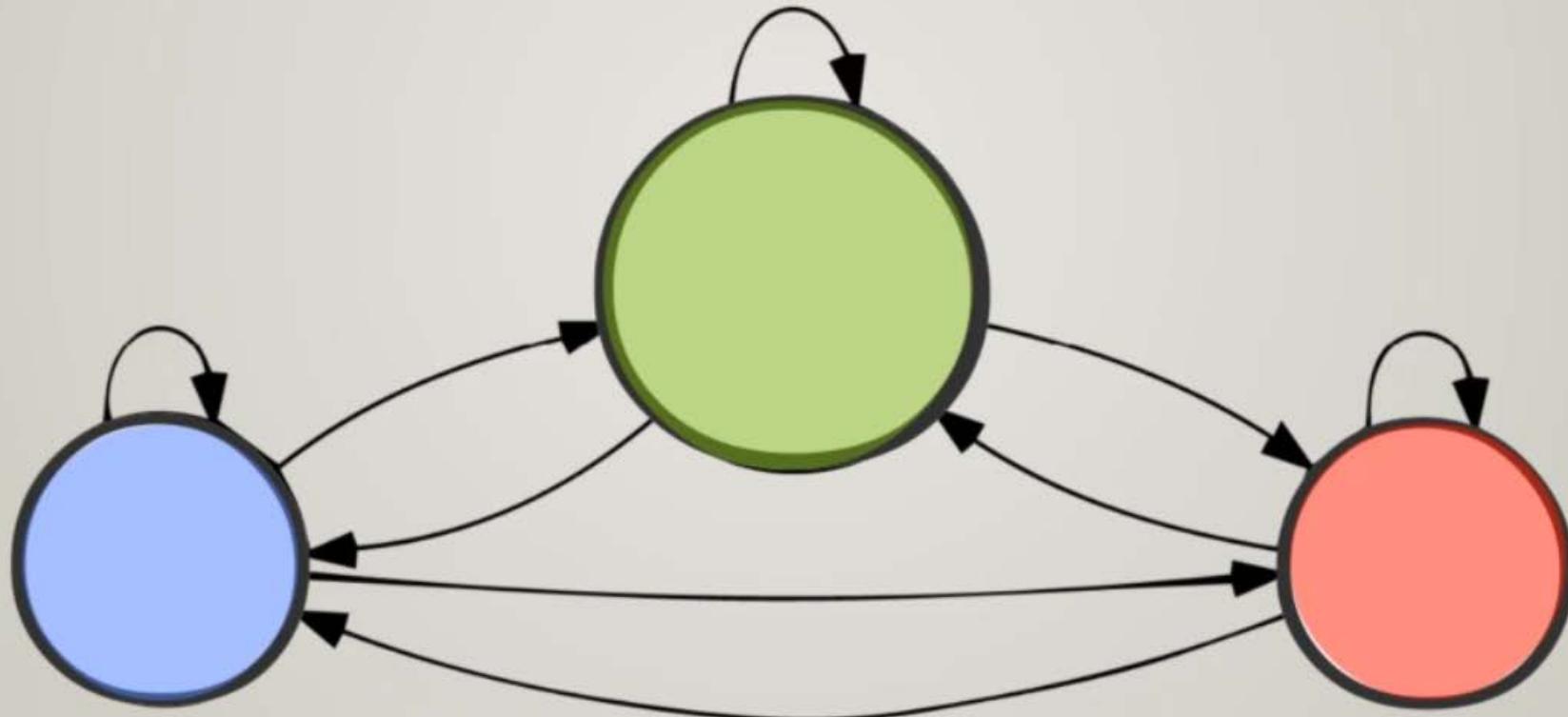


Generative models



Generative models

Hidden Markov models



Generative models

Bayesian networks

Likelihood of the evidence 'E' if
the Hypothesis 'H' is true

Prior Probability

$$p(H|E) = \frac{p(E|H) * p(H)}{p(E)}$$

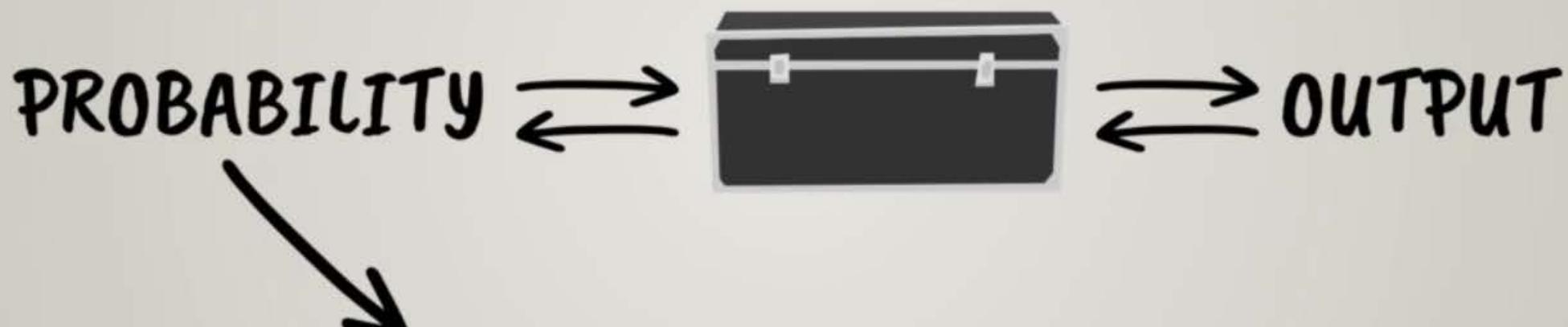
Posterior Probability of 'H'
given the evidence

Prior probability that the
evidence itself is true

Generative models

Bayesian networks

inputs are..



More information than in NN

Each problem, no matter if an image, a business problem, or a Shakespeare sonnet, tells a story. With the right tools you can reveal it and take advantage of it.



A matrix is a collection of numbers
ordered in rows and columns.
Here is one.

element

$$\begin{bmatrix} 5 & 12 & 6 \\ -3 & 0 & 14 \end{bmatrix}$$

A is a 2-by-3 matrix

$$\underline{\mathbf{A = \begin{bmatrix} 5 & 12 & 6 \\ -3 & 0 & 14 \end{bmatrix}}}$$

2 rows

3 columns

Matrices are main characters
in mathematical operations

◆ addition

◆ subtraction

◆ multiplication

$$\begin{bmatrix} 5 & 12 & 6 \\ -3 & 0 & 14 \end{bmatrix}$$

A matrix can only contain numbers,
symbols, or expressions

$$A = \begin{bmatrix} 5 & 12 & 6 \\ -3 & 0 & 14 \end{bmatrix}$$

2x3

$$B = \begin{bmatrix} a & b & c \\ d & e & f \end{bmatrix}$$

2x3

$$C = \begin{bmatrix} a-x & b+5 \\ d+e & e \end{bmatrix}$$

2x2

$A =$
 $m \times n$

[

a_{ij}

j-th
column

]

i-th row

MATHEMATICS

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} & \cdots & a_{1n} \\ a_{21} & a_{22} & a_{23} & \cdots & a_{2n} \\ a_{31} & a_{32} & a_{33} & \cdots & a_{3n} \\ \vdots & \vdots & \vdots & a_{ij} & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ a_{m1} & a_{m2} & a_{m3} & \cdots & a_{mn} \end{bmatrix}$$

PROGRAMMING

$$\begin{bmatrix} a_{00} & a_{01} & a_{02} & \cdots & a_{0(n-1)} \\ a_{10} & a_{11} & a_{12} & \cdots & a_{1(n-1)} \\ a_{20} & a_{21} & a_{22} & \cdots & a_{2(n-1)} \\ \vdots & \vdots & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ a_{(m-1)0} & a_{(m-1)1} & a_{(m-1)2} & \cdots & a_{(m-1)(n-1)} \end{bmatrix}$$

MATRICES

$$A = m \times n \left[\begin{array}{cccc|c} a_{11} & a_{12} & a_{13} & \dots & a_{1n} \\ a_{21} & a_{22} & a_{23} & \dots & a_{2n} \\ a_{31} & a_{32} & a_{33} & \dots & a_{3n} \\ \vdots & \vdots & \vdots & a_{ij} & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ a_{m1} & a_{m2} & a_{m3} & \dots & a_{mn} \end{array} \right] \quad \left. \right\} \text{rows}$$

columns

2D

SCALARS

$$[a_{11}] \}$$

1 row

1 column

SCALARS

All numbers we know from algebra are referred to as scalars in linear algebra

[15]; [1]; [2]; [-5]; []

SCALARS HAVE 0 DIMENSIONS

VECTORS

$$U = \begin{bmatrix} 5 \\ -2 \\ 4 \end{bmatrix}$$

3 rows

1 column

3×1

VECTORS

A vector is practically the
simplest linear algebraic object

$\begin{bmatrix} 5 \\ -2 \\ 4 \end{bmatrix}$

$\begin{bmatrix} 5 & 12 & 6 \\ -3 & 0 & 14 \end{bmatrix}$

vect. vect. vect.

1 2 3

TYPES OF VECTORS

column

$$\begin{bmatrix} 5 \\ -2 \\ 4 \end{bmatrix}$$

row

$$\begin{pmatrix} 3 & 4 & 5 & 8 \end{pmatrix}$$

LENGTH

Df: The number of elements in a vector

$$\begin{bmatrix} 5 \\ -2 \\ 4 \end{bmatrix}$$

Length = 3

Length = 4

$$[3 4 5 8]$$

SUMMARY

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} & \dots & a_{1n} \\ a_{21} & a_{22} & a_{23} & \dots & a_{2n} \\ a_{31} & a_{32} & a_{33} & \dots & a_{3n} \\ \vdots & \vdots & \vdots & a_{ij} & \vdots \\ a_{m1} & a_{m2} & a_{m3} & \dots & a_{mn} \end{bmatrix}_{m \times n}$$

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}_{m \times 1}$$

$$\begin{bmatrix} x \end{bmatrix}_{1 \times 1}$$

2D Matrix
 $m \times n$

1D Vector
 $m \times 1$

0D Scalar
 1×1

LINEAR ALGEBRA AND GEOMETRY

[x]

0D

Scalar
 1×1



0D

Point

LINEAR ALGEBRA AND GEOMETRY

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix}$$

1D

Vector
 $m \times 1$

1D

Line



LINEAR ALGEBRA AND GEOMETRY

$$\begin{bmatrix} x_1 & y_1 \\ x_2 & y_2 \\ \vdots & \vdots \\ x_m & y_m \end{bmatrix}$$

2D

Matrix
 $m \times 2$

2D

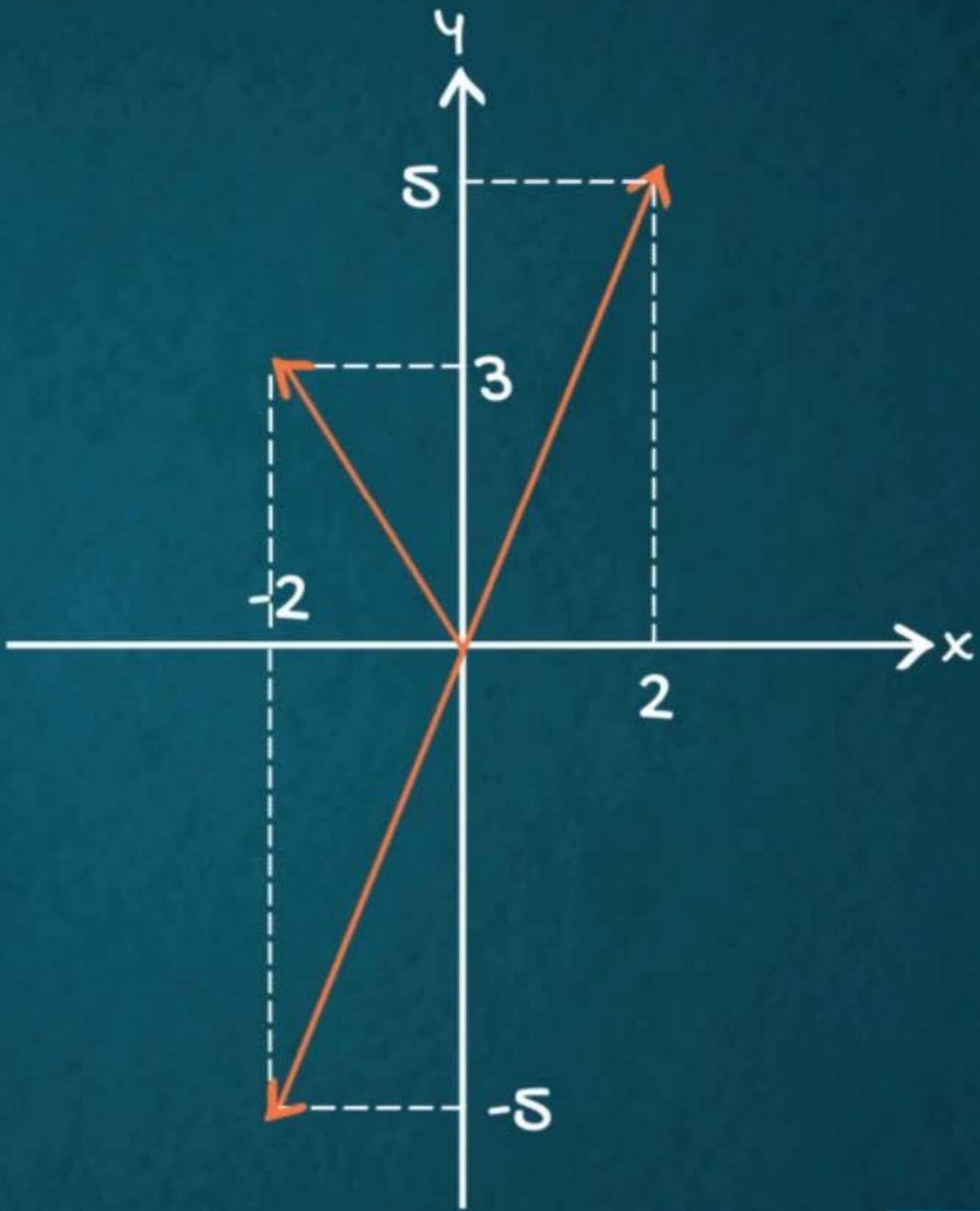
Line



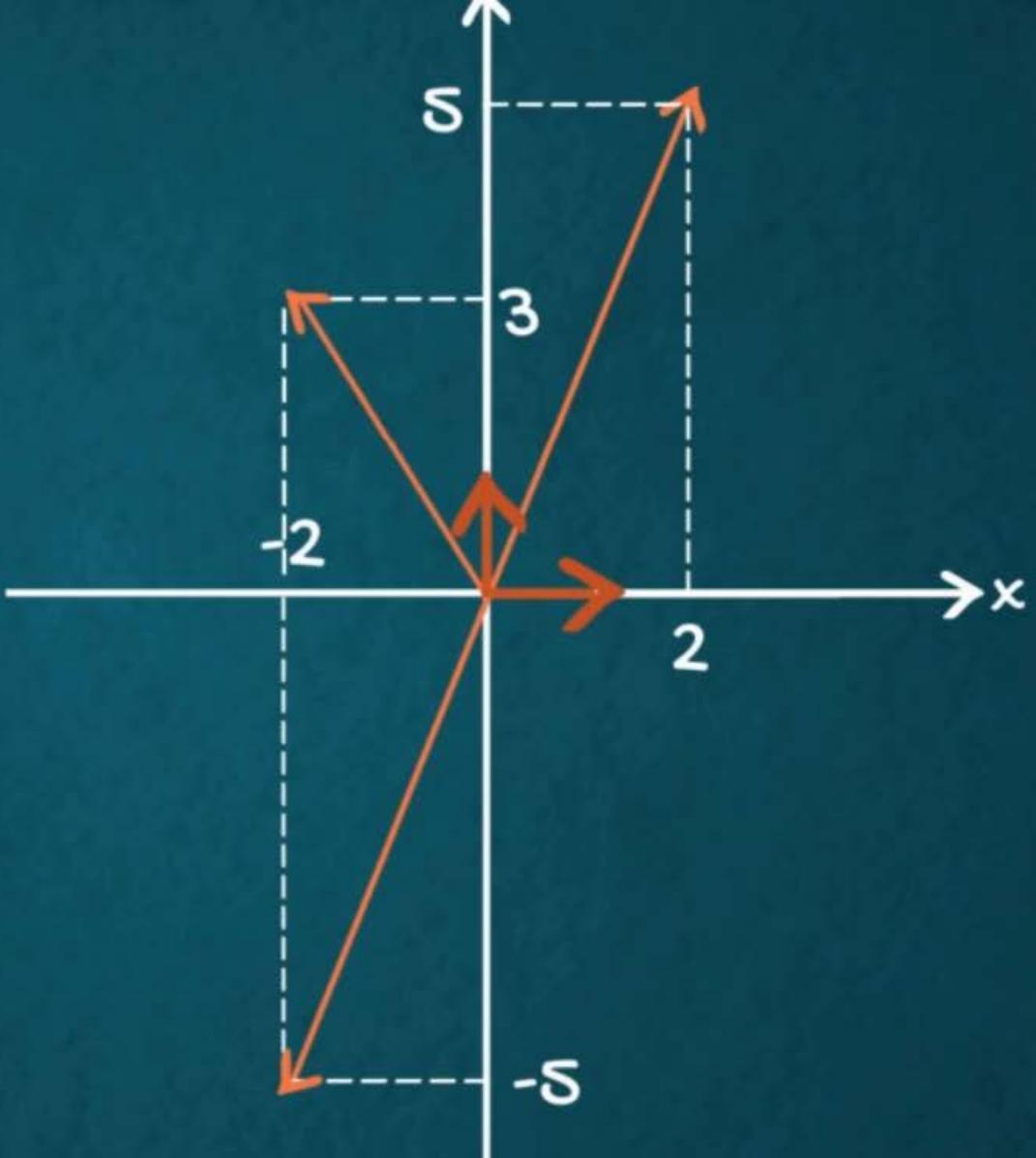
$$\begin{bmatrix} 2 \\ 5 \end{bmatrix}$$

$$\begin{bmatrix} -2 \\ -5 \end{bmatrix}$$

$$\begin{bmatrix} -2 \\ 3 \end{bmatrix}$$



$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



TRANSPOSING VECTORS

$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

x

\top

$$\begin{bmatrix} 1 & 2 & 3 \end{bmatrix}$$

x^\top

TRANSPOSING VECTORS

- ▲ The values are not changing or transforming;
only their position is
- ▲ Transposing the same vector (object) twice
yields the initial vector (object)
- ▲ A 3×1 matrix transposed is a 1×3 matrix

TRANSPOSING MATRICES

$$\underline{\underline{A = \begin{bmatrix} 5 & 12 & 6 \\ -3 & 0 & 14 \end{bmatrix}}}$$

$$\underline{\underline{A^T = \begin{bmatrix} 5 & -3 \\ 12 & 0 \\ 6 & 14 \end{bmatrix}}}$$

SCALAR MULTIPLICATION

$$[6] \cdot [5] = [30]$$

$$[10] \cdot [-2] = [-20]$$

DOT PRODUCT

$$\begin{bmatrix} 2 \\ 8 \\ -4 \end{bmatrix} \times \begin{bmatrix} 1 \\ -7 \\ 3 \end{bmatrix}$$

(inner product)

DOT PRODUCT

$$\begin{bmatrix} 2 \\ 8 \\ -4 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ -7 \\ 3 \end{bmatrix} = [2 \times 1 + 8 \times (-7) + (-4) \times 3]$$

DOT PRODUCT

$$\begin{bmatrix} 2 \\ 8 \\ -4 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ -7 \\ 3 \end{bmatrix} = [-66]$$

Vector

Vector

Scalar

DOT PRODUCT

$$\begin{bmatrix} 2 \\ 8 \\ -4 \end{bmatrix} \cdot \begin{bmatrix} 1 \\ -7 \\ 3 \end{bmatrix} = [-66]$$

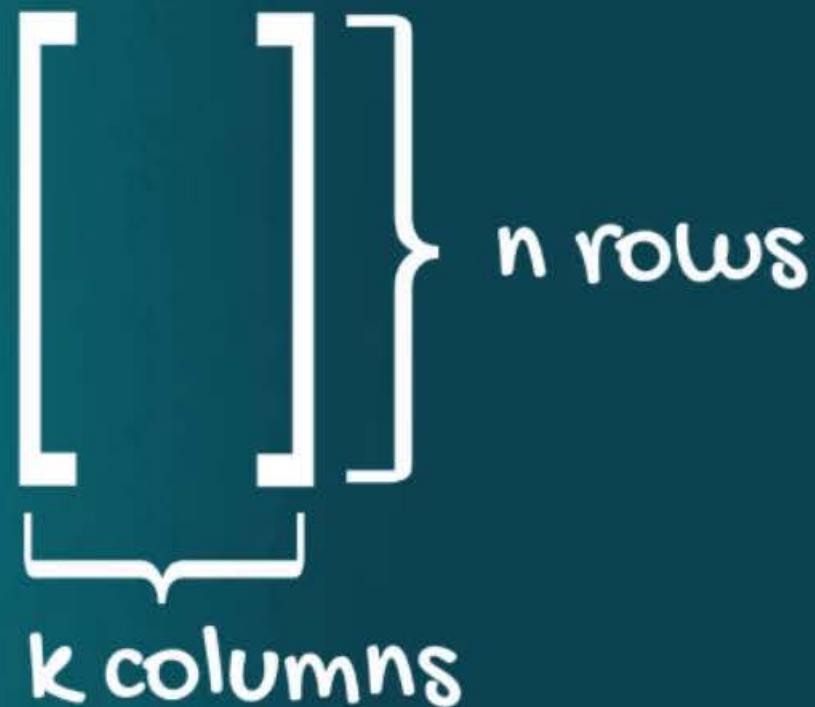
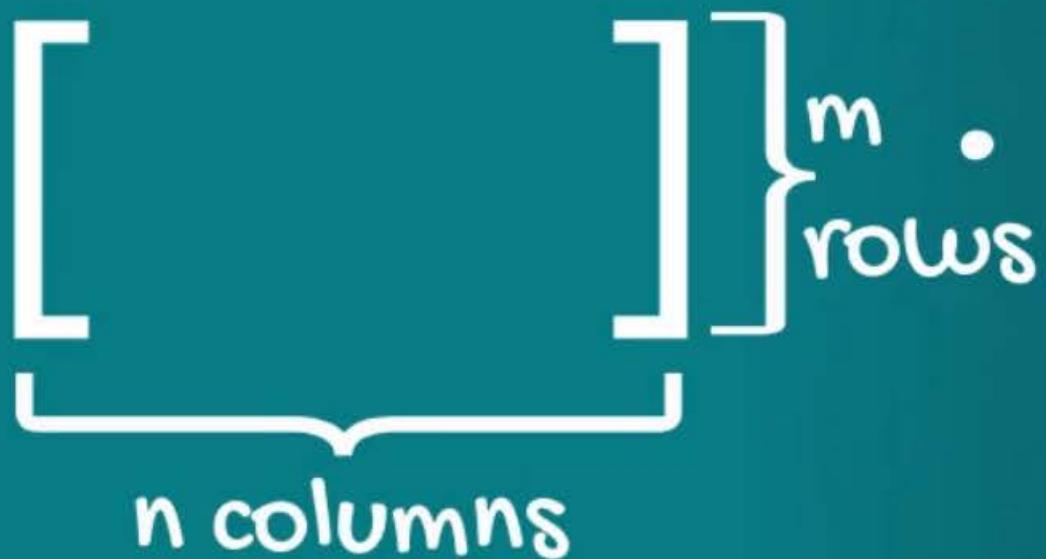
Vector

Vector

Scalar

the 'sum of the products of the corresponding elements'

DOT PRODUCT



Condition: We can only multiply an $m \times n$ with an $n \times k$ matrix

DOT PRODUCT

$$\begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}$$

2×3

$$\cdot \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}$$

$3 \times k$

Condition: We can only multiply an $m \times n$ with an $n \times k$ matrix

DOT PRODUCT

$$[\underset{m \times n}{\square}] \cdot [\underset{n \times k}{\square}] = [\underset{m \times k}{\square}]$$

$$\begin{matrix} \mathbf{A} & \cdot & \mathbf{B} & = & \mathbf{C} \\ 2 \times 3 & & 3 \times 6 & & 2 \times 6 \end{matrix}$$

WHY IS LINEAR ALGEBRA ACTUALLY USEFUL?

Applications in data science

Vectorized code

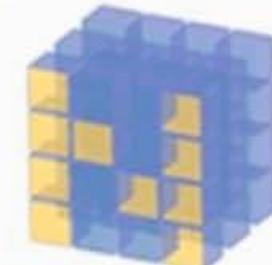
Image recognition

Dimensionality reduction

VECTORIZED CODE

Whenever we are using linear algebra
to compute many values simultaneously

much
faster



NUMPY

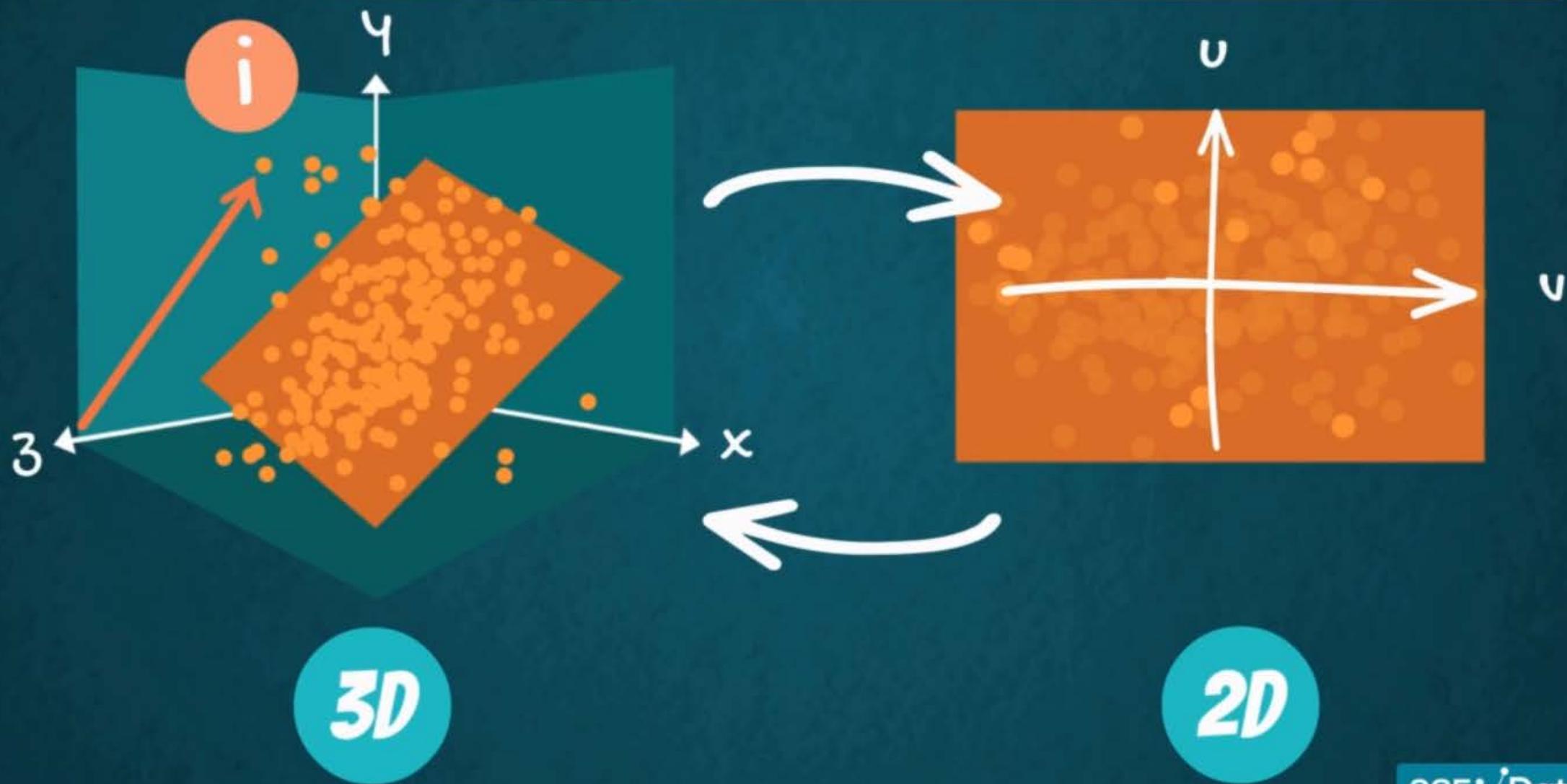
DIMENSIONALITY REDUCTION

◆ Eigenvalues

◆ Eigenvector



DIMENSIONALITY REDUCTION



DIMENSIONALITY REDUCTION

INITIAL

i

$$\begin{bmatrix} x & y & z \\ \vdots & \vdots & \vdots \\ x_i & y_i & z_i \\ \vdots & \vdots & \vdots \end{bmatrix}$$

$m \times 3$

3D

TRANSFORMED

$$\begin{bmatrix} u & v \\ \vdots & \vdots \\ u_i & v_i \\ \vdots & \vdots \end{bmatrix}$$

$m \times 2$

2D

