



UNIVERSITÀ  
DEGLI STUDI  
DI PADOVA

# Machine Learning 2024/2025



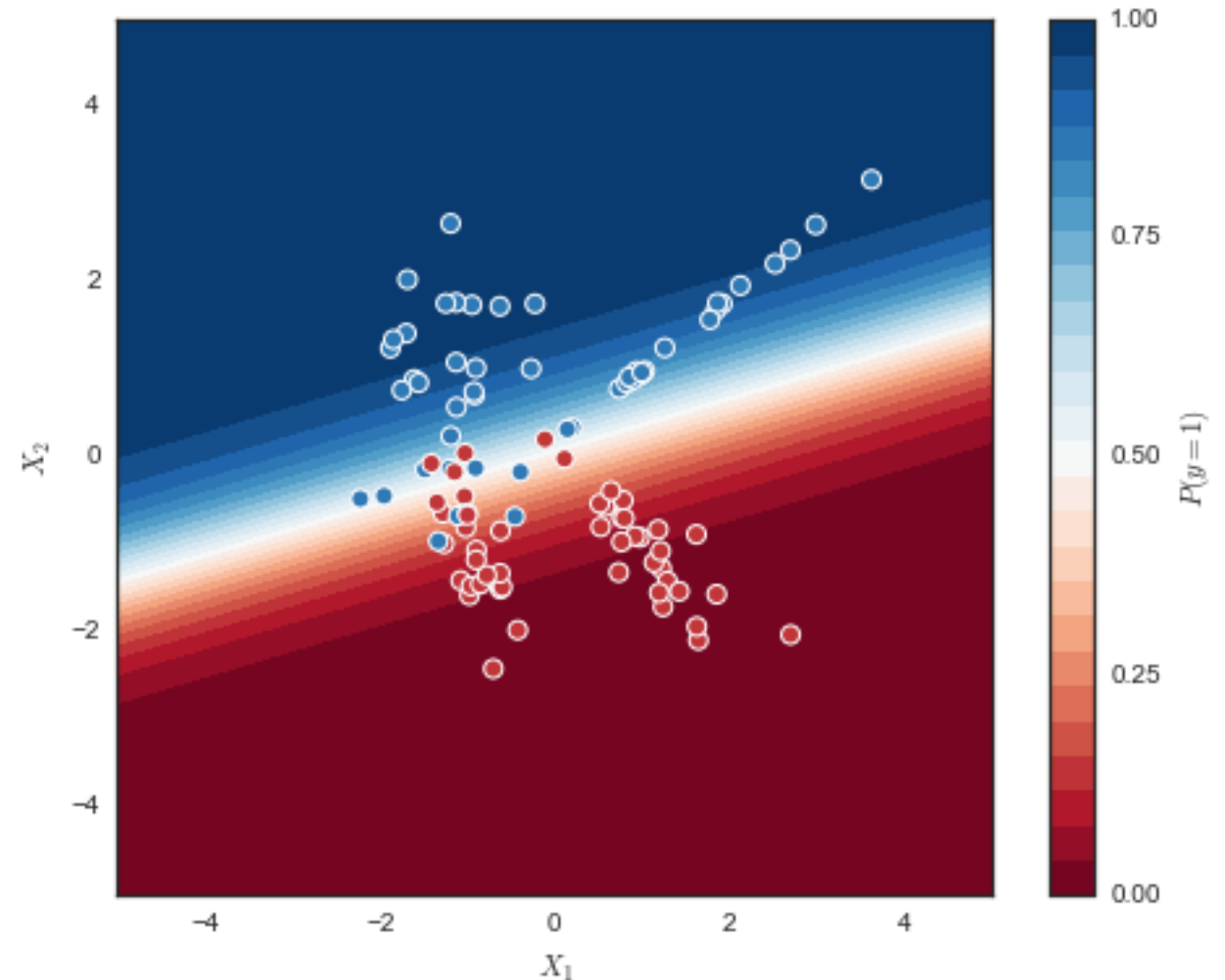
# Lecture #24 Neural Networks Training

Gian Antonio Susto



# Before Starting 1/2: clarifications on ROC curve

In many approaches (logistic regression, SVM, ...) we can derive a distance from the decision boundary, a probability of being classified to one class or another...



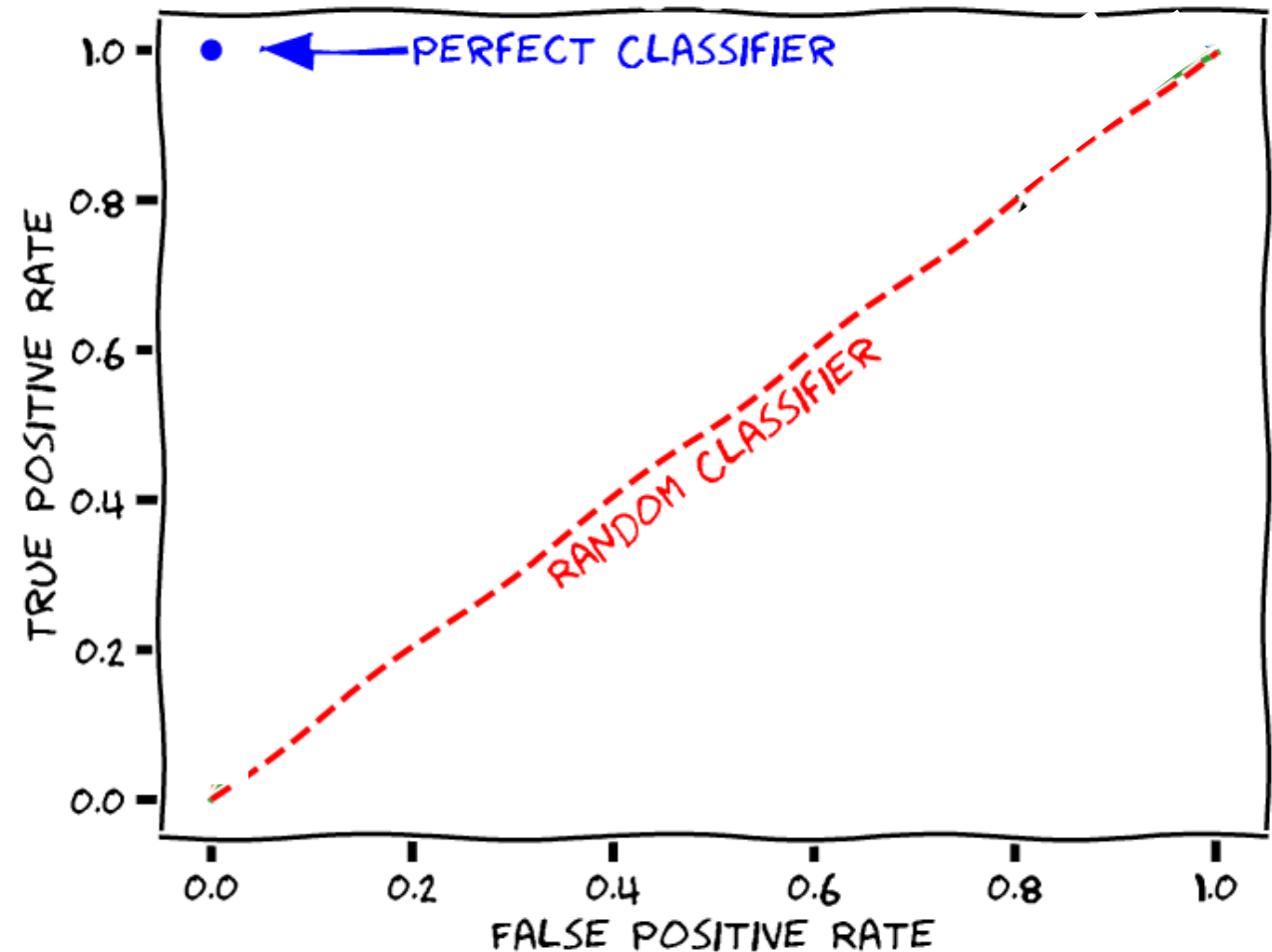
# Before Starting 1/2: clarifications on ROC curve

- True Positive Rate

$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

- False Positive Rate

$$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$



# Before Starting 1/2: clarifications on ROC curve

- True Positive Rate

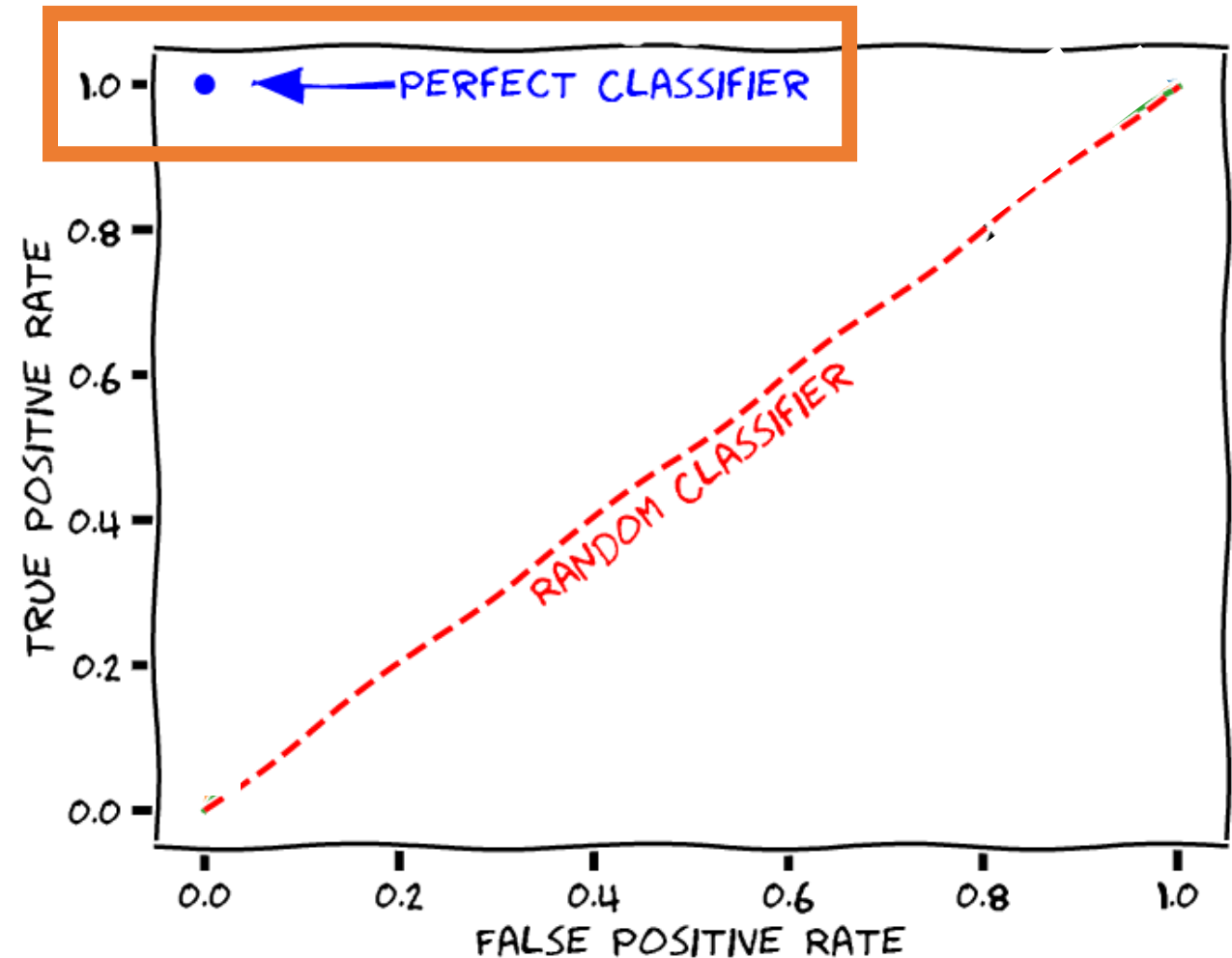
$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

- False Positive Rate

$$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

In this case:

- TPR = 1 (0 false negatives, ie all TP or TN)!
- FPR = 0 (0 false positives)



# Before Starting 1/2: clarifications on ROC curve

- True Positive Rate

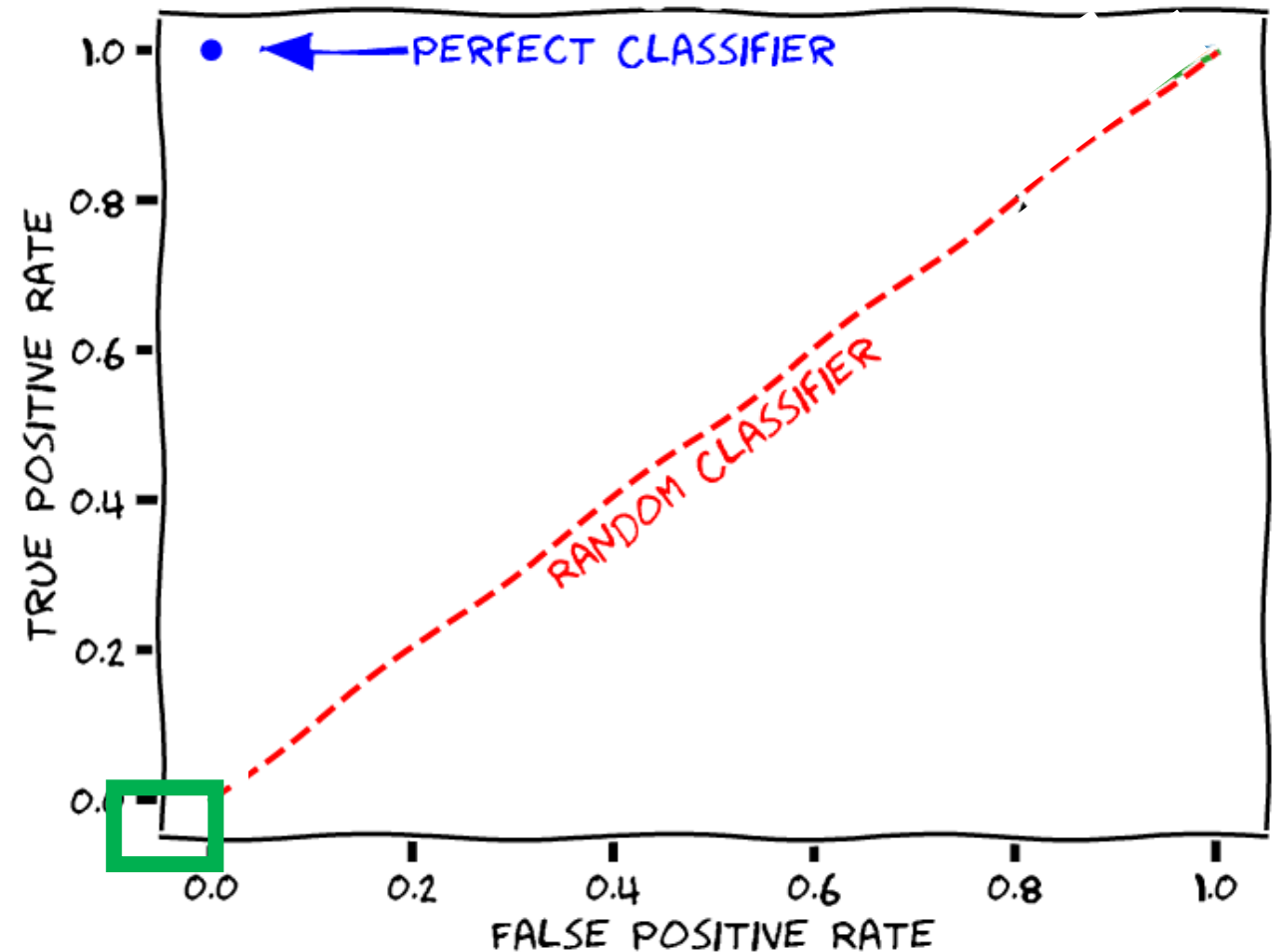
$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

- False Positive Rate

$$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

If a classifier always classifies data as negatives:

- TPR = 0 (as TP = 0)
- FPR = 0 (as FP = 0)



# Before Starting 1/2: clarifications on ROC curve

- True Positive Rate

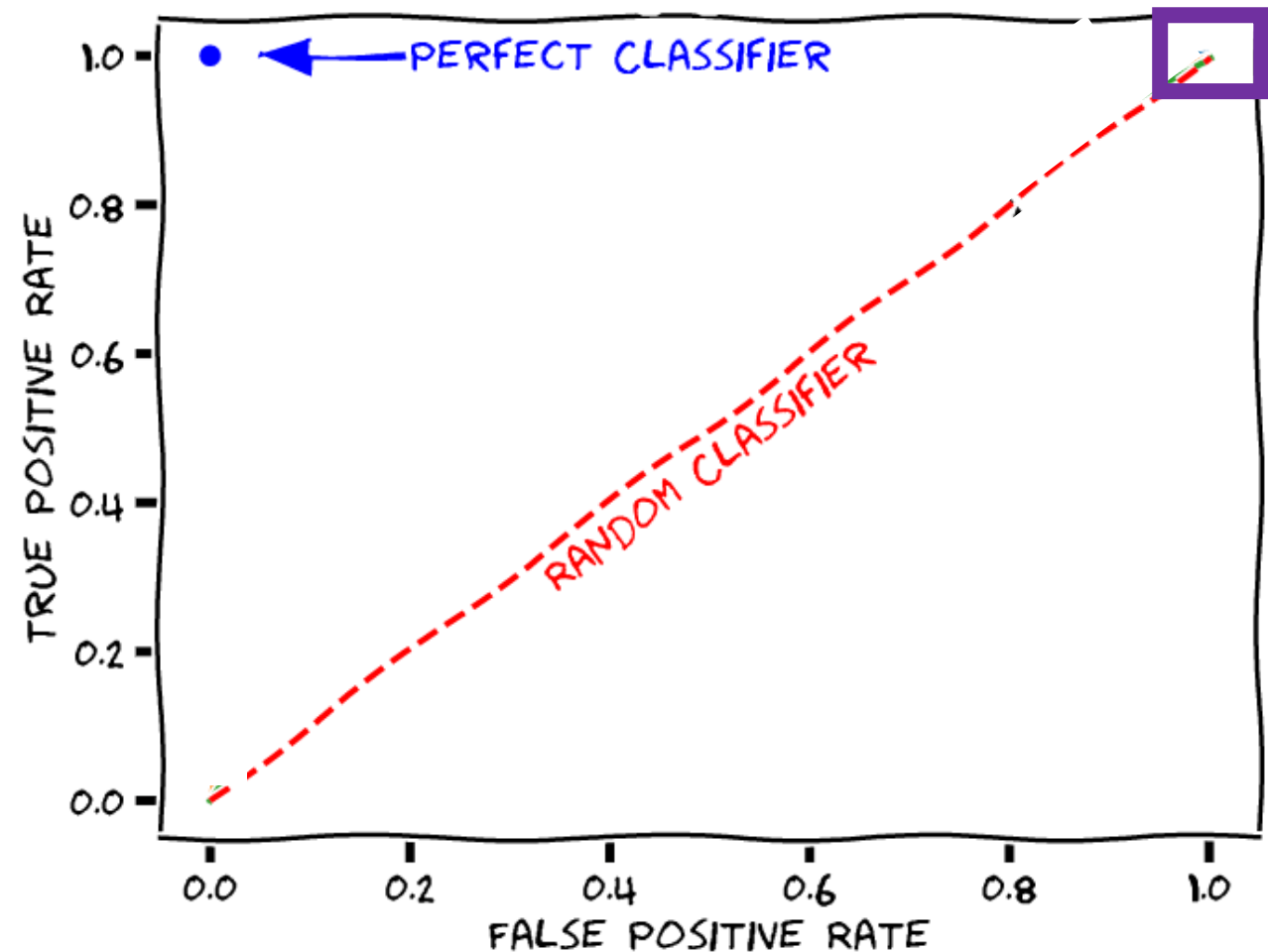
$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

- False Positive Rate

$$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

If a classifier always classifies data as positives:

- TPR = 1 (as FN = 0)
- FPR = 1 (as TN = 0)



# Before Starting 1/2: clarifications on ROC curve

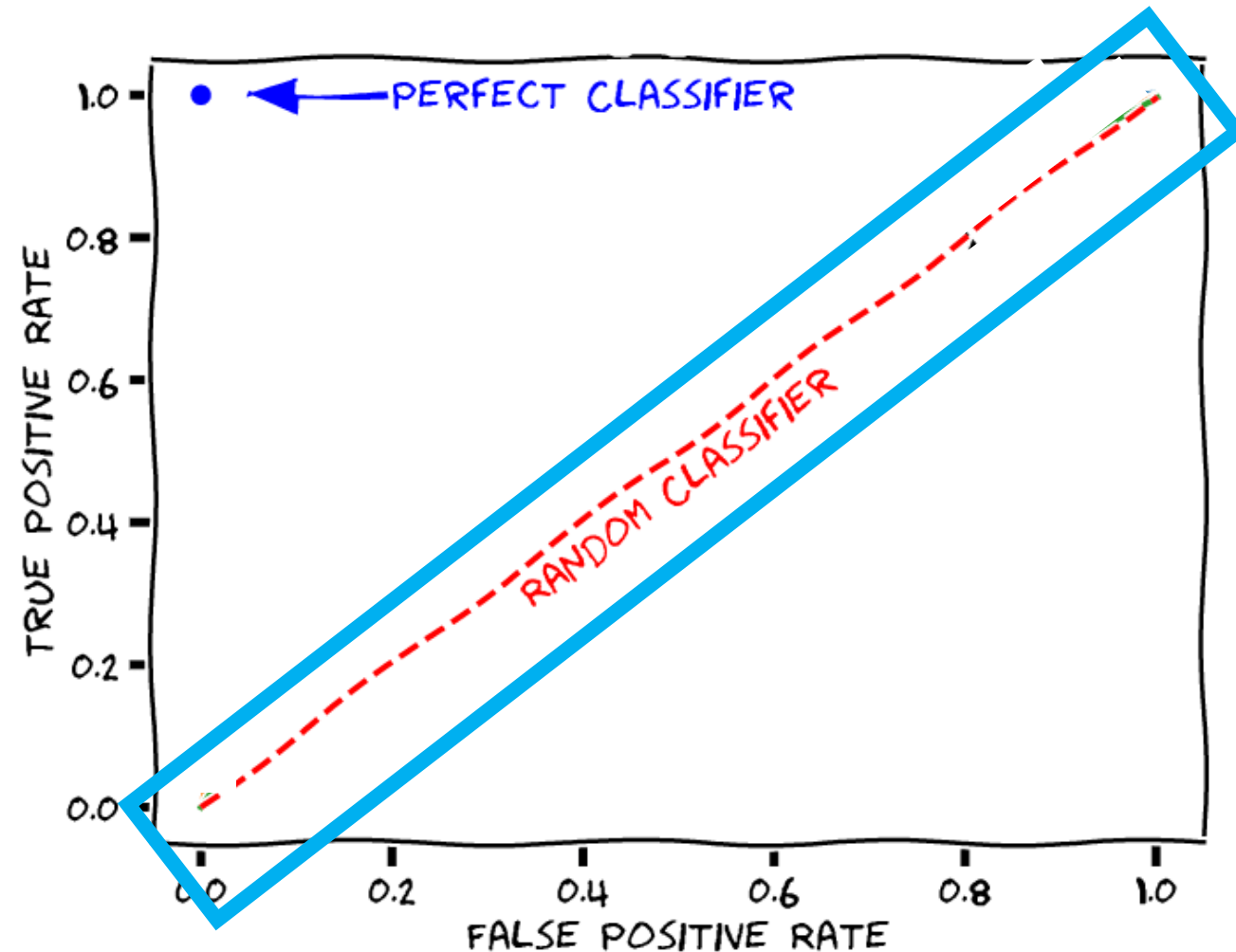
- True Positive Rate

$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

- False Positive Rate

$$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

A random classifier will have  $\text{TPR} = \text{FPR}$  since positives and negatives are assigned randomly





# Before Starting 1/2: clarifications on ROC curve

- True Positive Rate

$$\text{TPR} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

- False Positive Rate

$$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

Scenario: A dataset with 100 people (50 positives, 50 negatives)

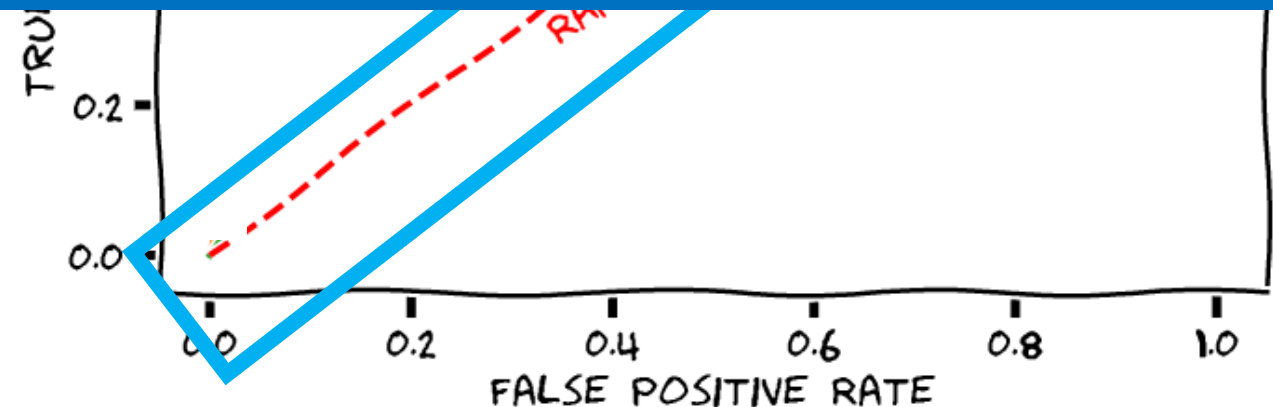
- CASE #01 (Threshold = 0.8, ie 20% are predicted as P, TP = 10, FP = 10, TN = 40, FN = 40)

$$\text{TPR} = 10/(10+40) = 0,2 \text{ \& } \text{FPR} = 10/(10+40) = 0,2$$

- CASE #02 (Threshold = 0.5, ie 50% are predicted as P, TP = 25, FP = 25, TN = 25, FN = 25)

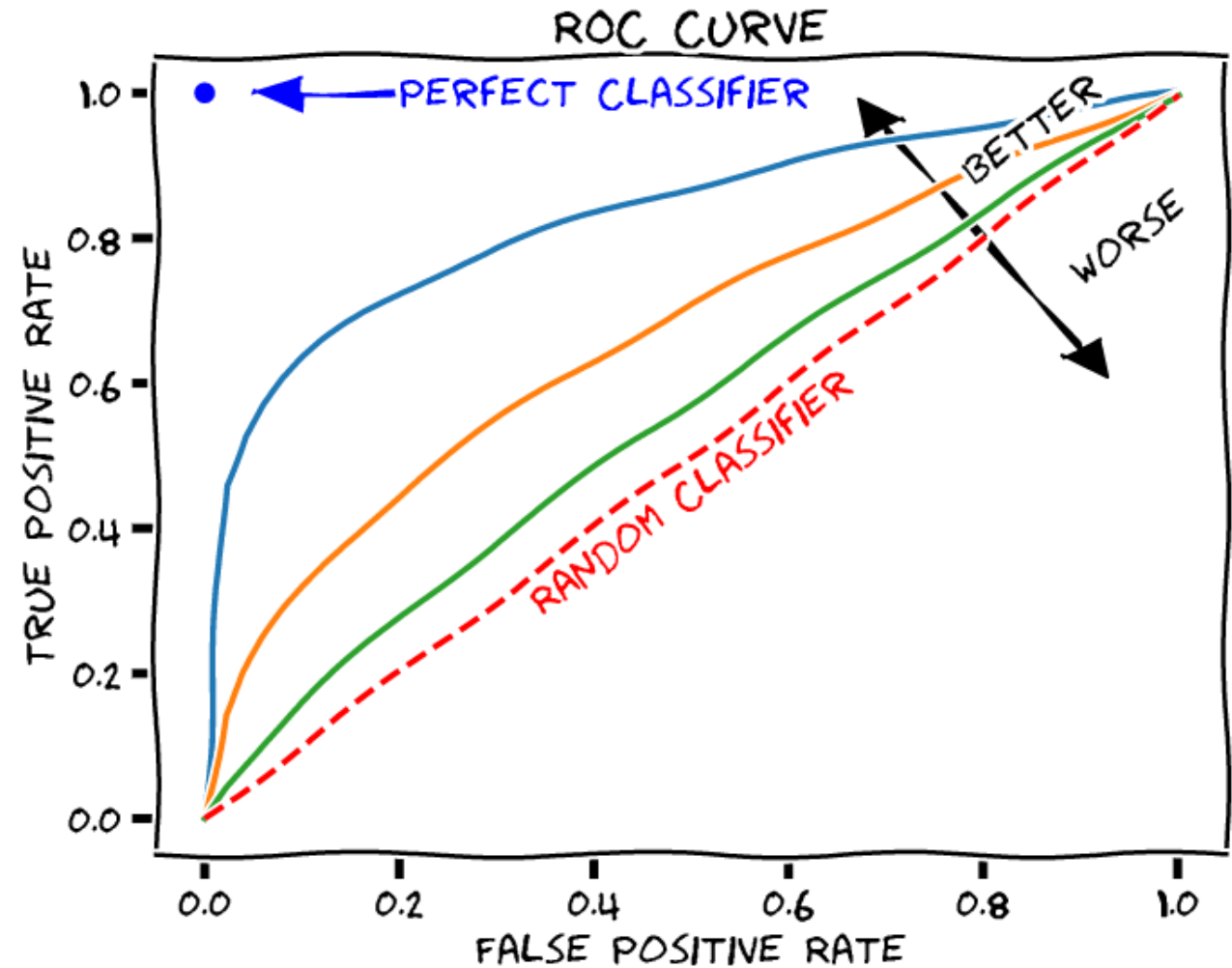
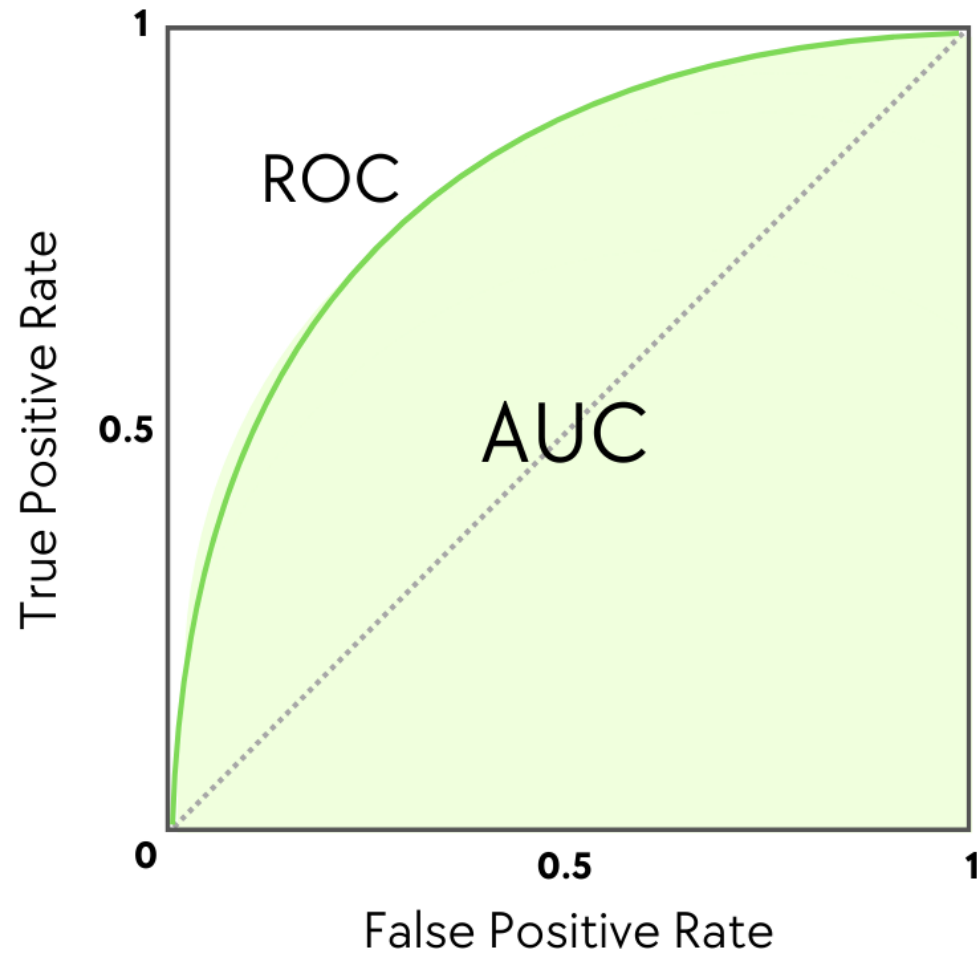
$$\text{TPR} = 25/(25+25) = 0,5 \text{ \& } \text{FPR} = 25/(25+25) = 0,5$$

A random classifier will have  $\text{TPR} = \text{FPR}$  since positives and negatives are assigned randomly

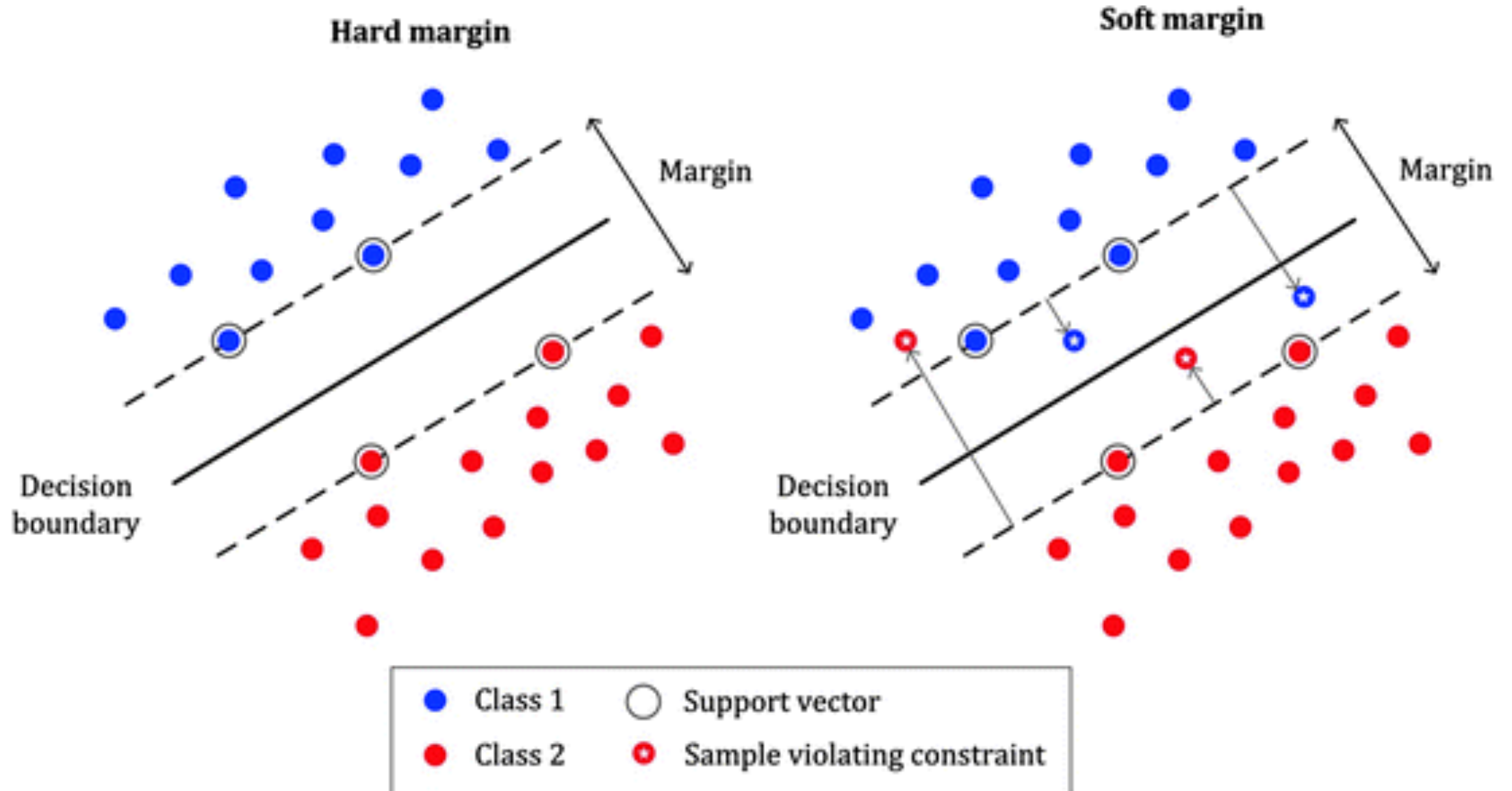




# Before Starting 1/2: clarifications on ROC curve



# Before Starting 2/2: on the role of C in SVM



# Before Starting 2/2: on the role of C in SVM

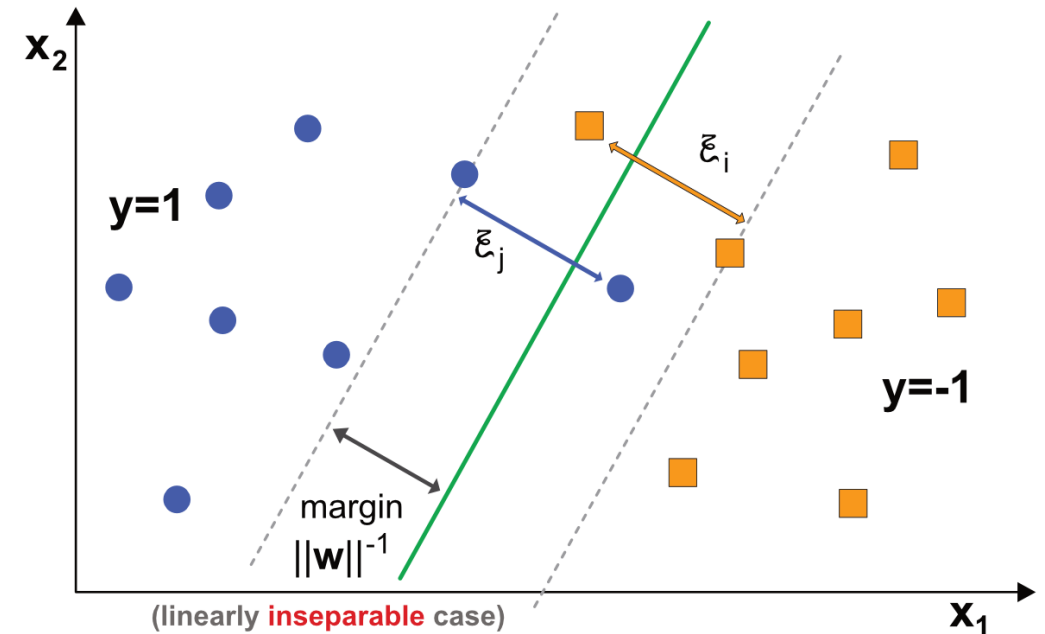
SVM solves a **convex optimization problem** (**linearly separable case**):

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{Subject to} \quad y_i(w \cdot x_i + b) \geq 1 \quad \forall i$$

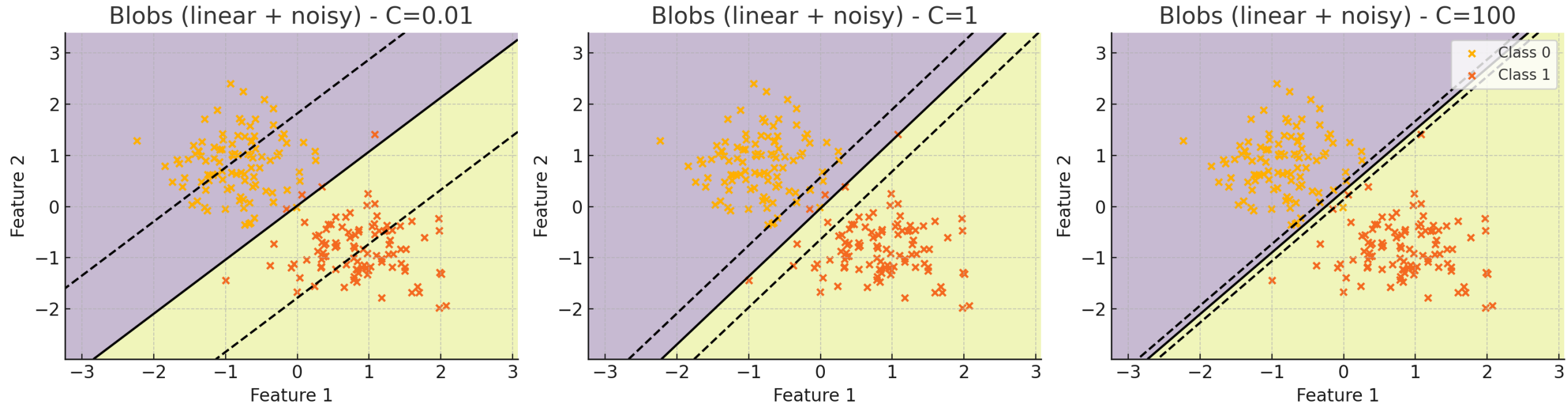
If **it is not separable**, we allow some "slack" (errors) by introducing variables  $\xi$ :

$$\min_{\mathbf{w}, b, \xi} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum \xi_i$$

$$\text{Subject to} \quad y_i(w \cdot x_i + b) \geq 1 - \xi_i, \quad \xi_i \geq 0$$



# Before Starting 2/2: on the role of C in SVM



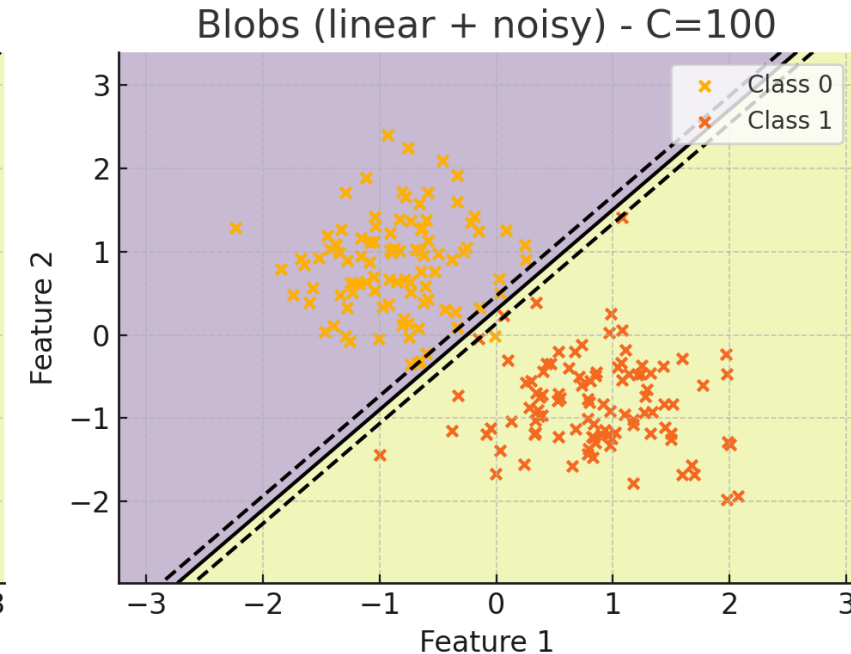
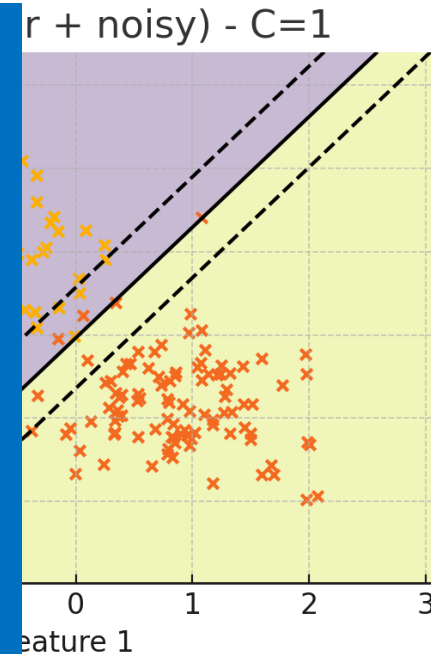
Here we allow the classifier to have some misclassification

Here we try to have all historical data point correctly classified

# Before Starting 2/2: on the role of C in SVM

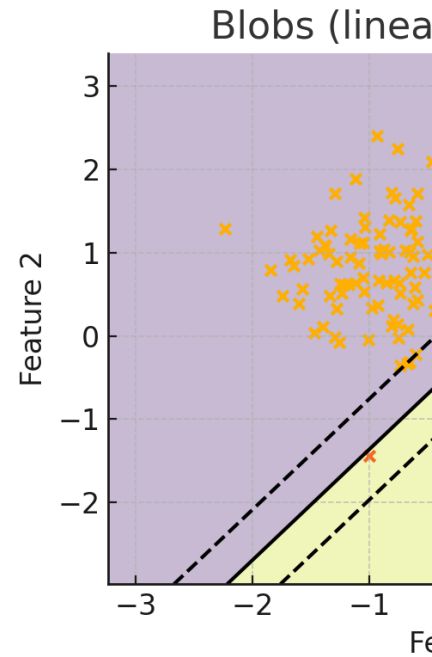
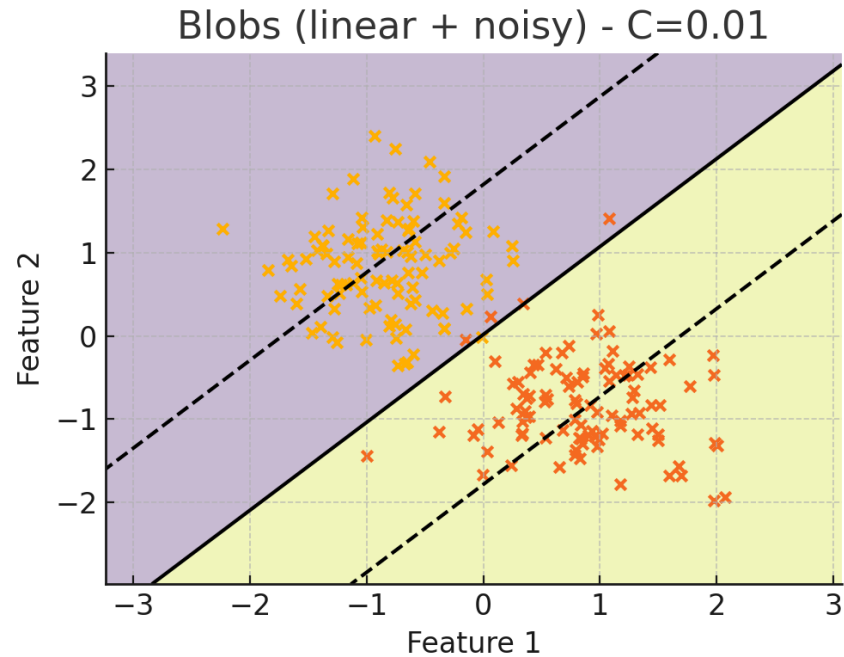
## C high

- Low regularization
- The model penalizes misclassifications heavily
- Tries to classify every point correctly (even noisy points)
- Margin becomes narrower to avoid errors
- Risk of overfitting



Here we try to have all historical data point correctly classified

# Before Starting 2/2: on the role of C in SVM



## C low

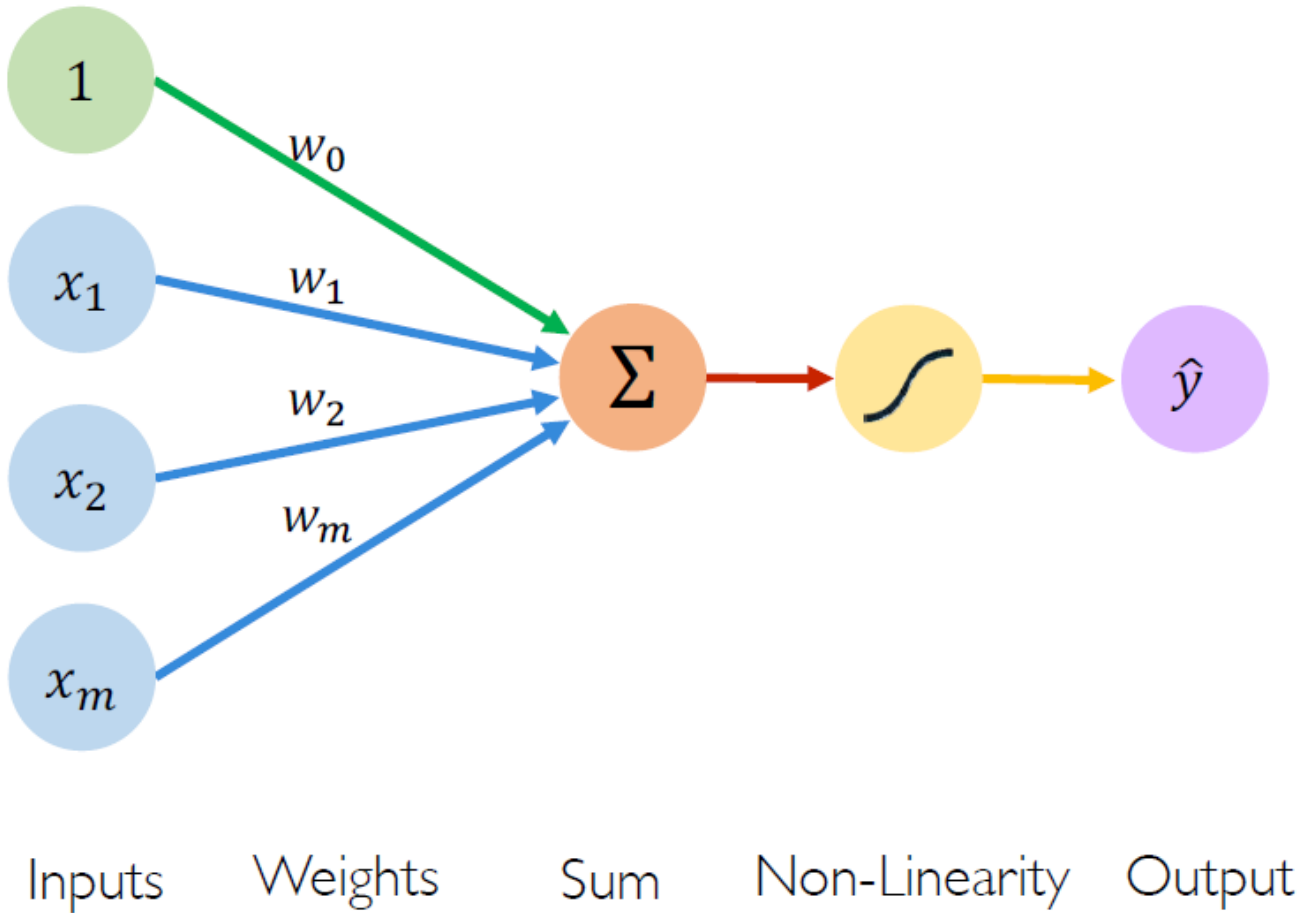
- High regularization
- Allows some classification errors to achieve a wider margin
- Focuses more on simplicity and generalization
- Risk of underfitting, but better performance on noisy data

Here we allow the classifier to have some misclassification

# Basics of Deep Learning (with some recap)



# The building block of Neural Networks (NN): the neuron

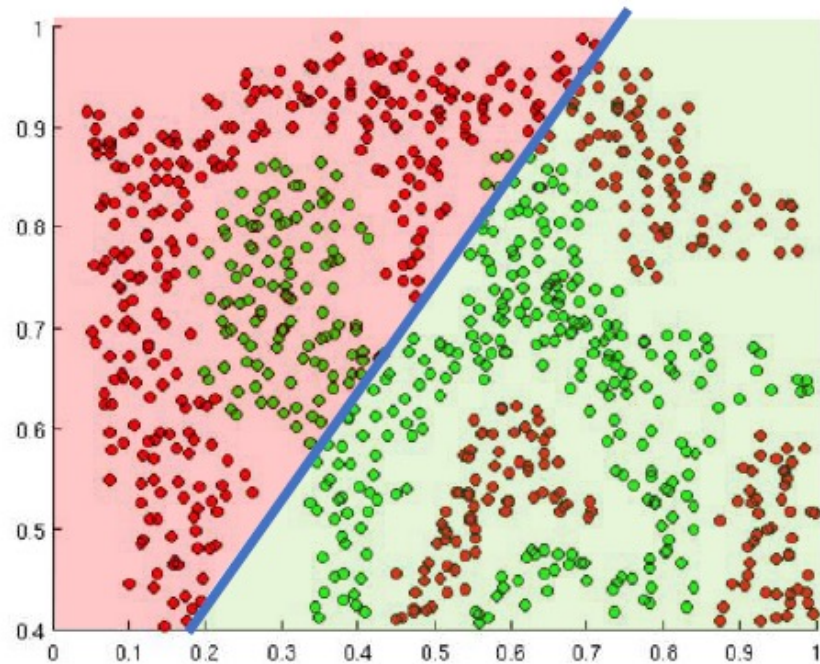


Activation function

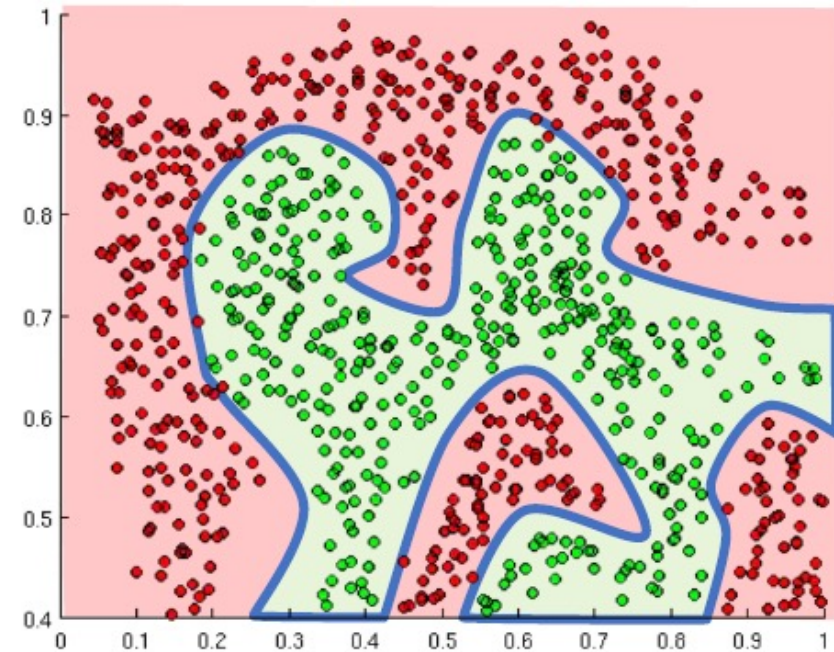
$$\hat{y} = g(w_0 + \mathbf{X}^T \mathbf{W})$$

$$\text{where: } \mathbf{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix} \text{ and } \mathbf{W} = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}$$

# Activation functions

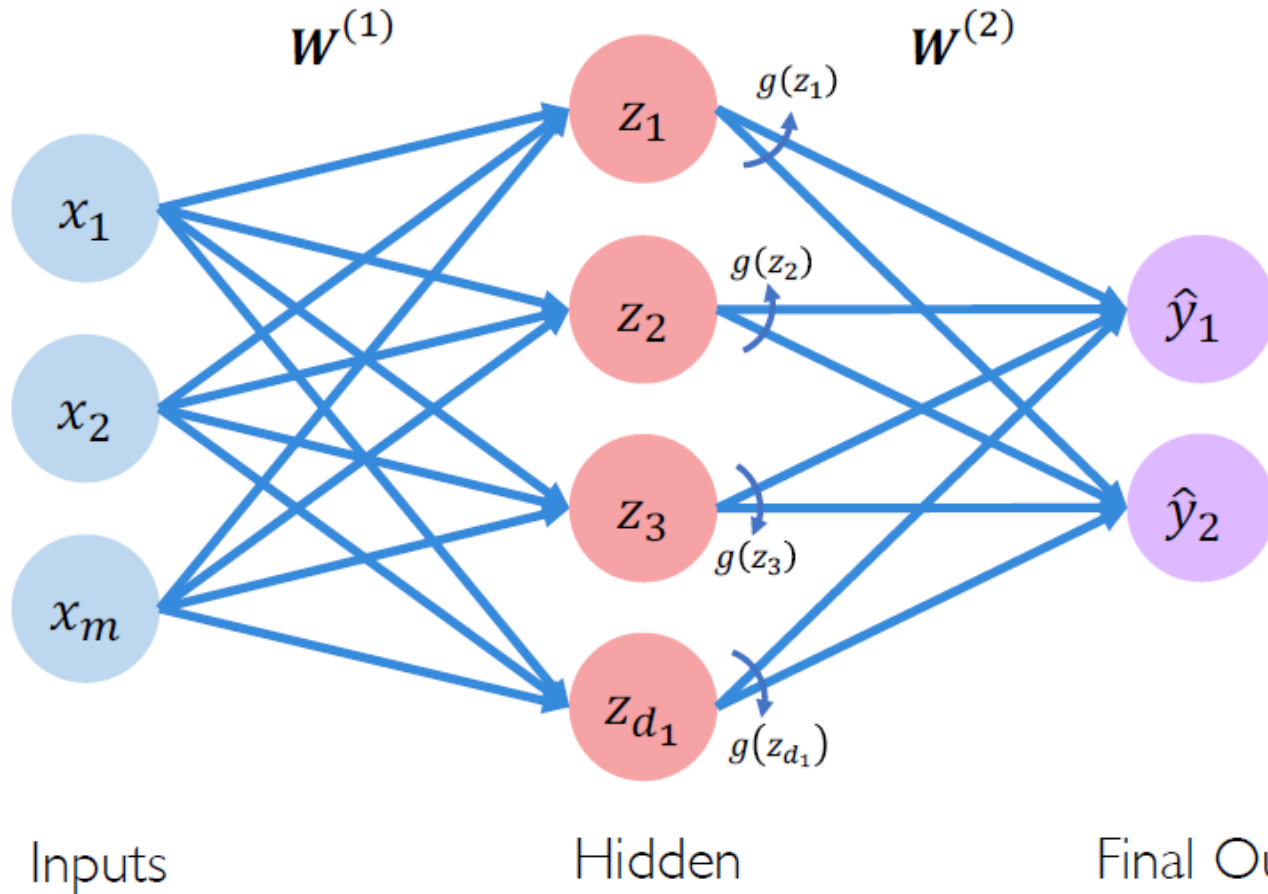


Linear Activation functions produce linear decisions no matter the network size



Non-linearities allow us to approximate arbitrarily complex functions

# Feedforward (Single Layer) Neural Network

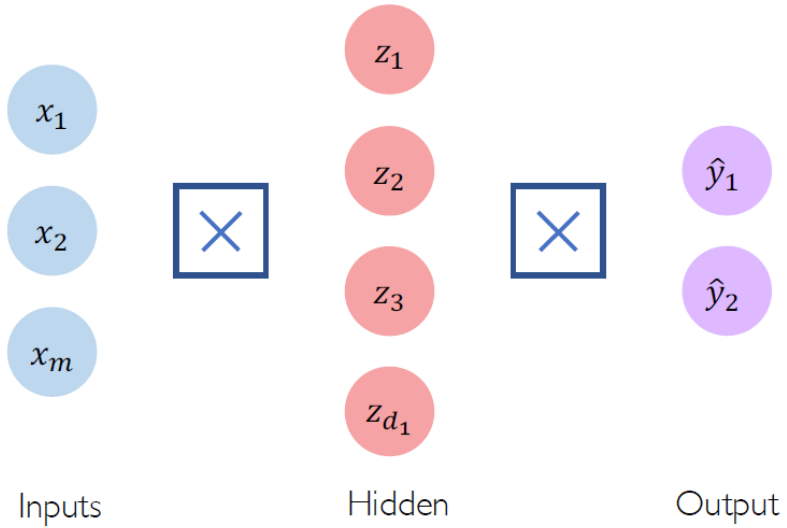


So-called *Vanilla Neural Network*

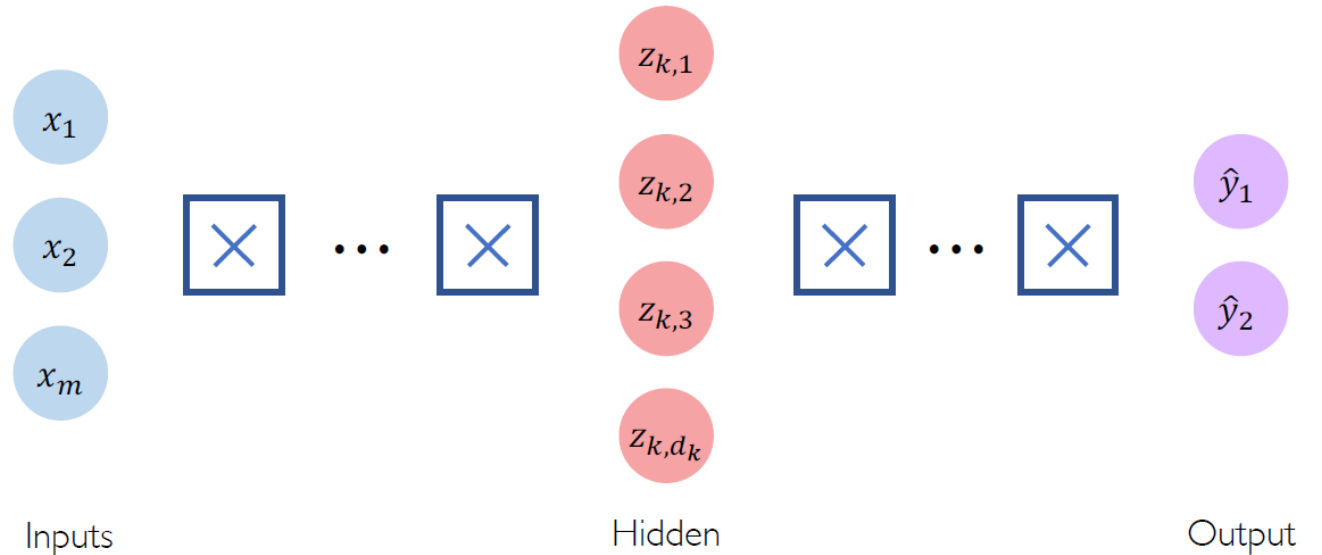
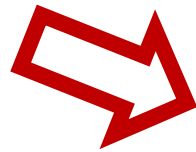
T. Hastie, R. Tibshirani, J. Friedman *The Elements of Statistical Learning*

$$z_i = w_{0,i}^{(1)} + \sum_{j=1}^m x_j w_{j,i}^{(1)} \quad \hat{y}_i = g \left( w_{0,i}^{(2)} + \sum_{j=1}^{d_1} z_j w_{j,i}^{(2)} \right)$$

# From Neural Network to Deep Learning architecture

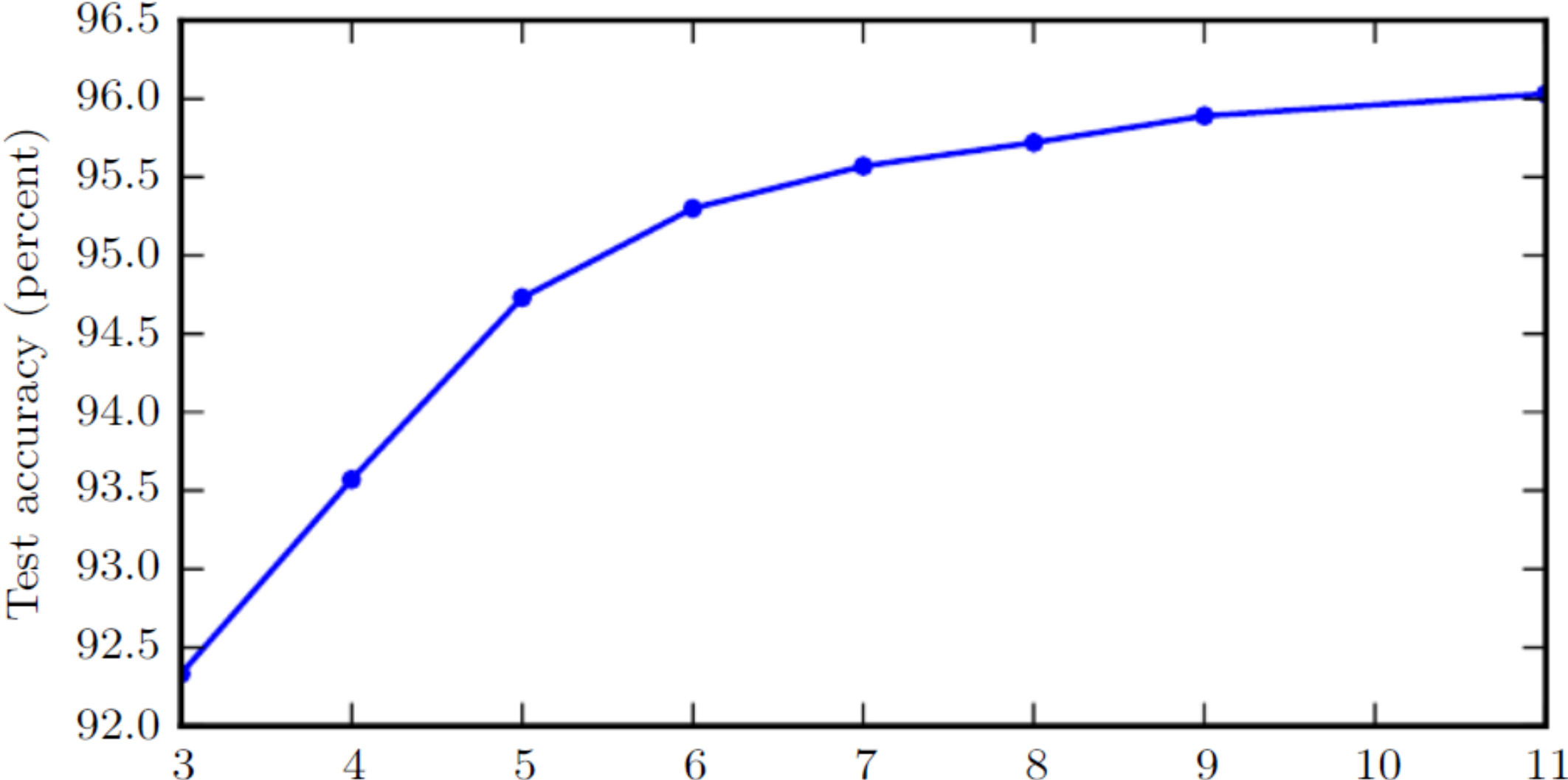


From 1 Layer to 'many' layers



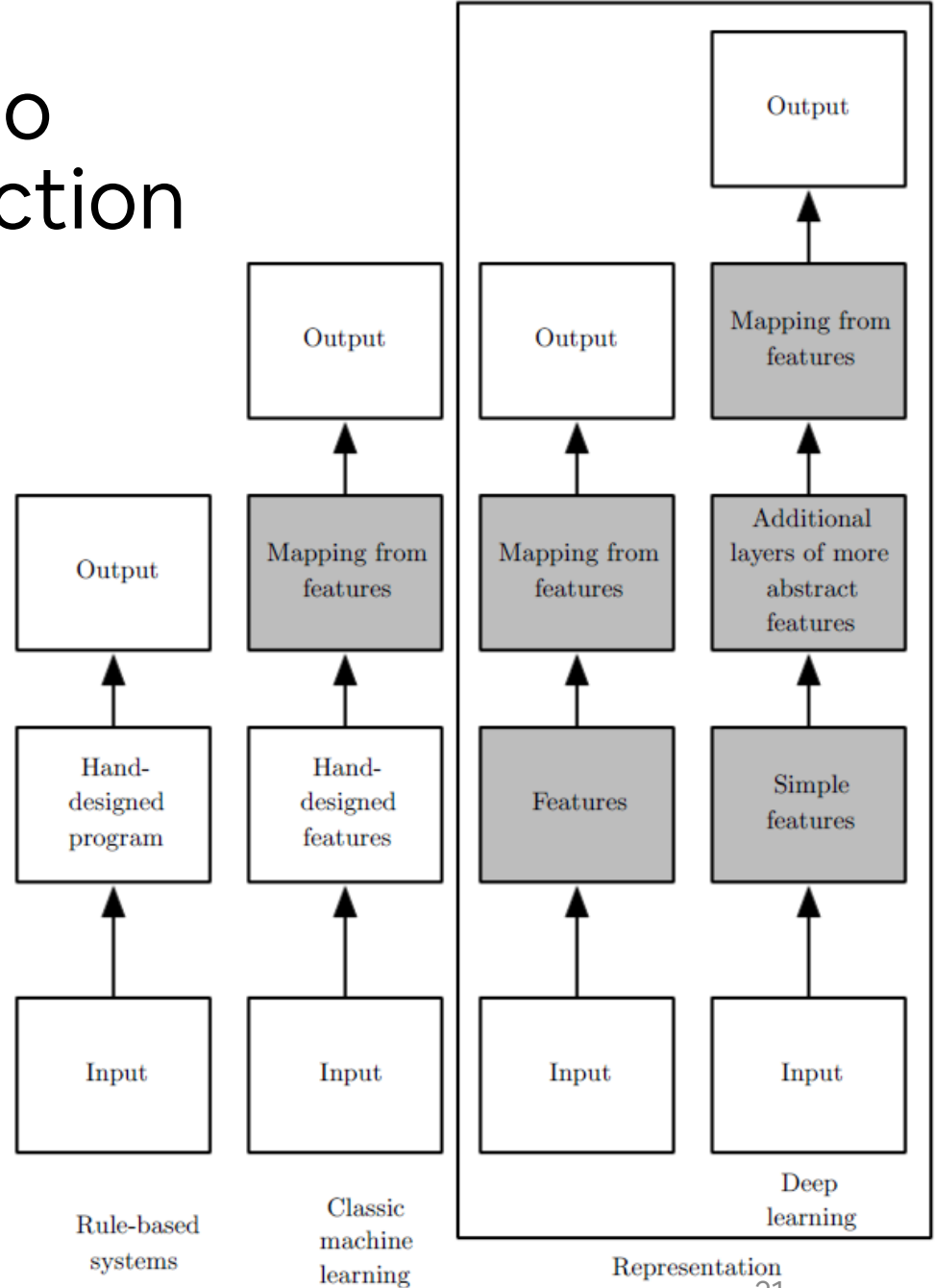
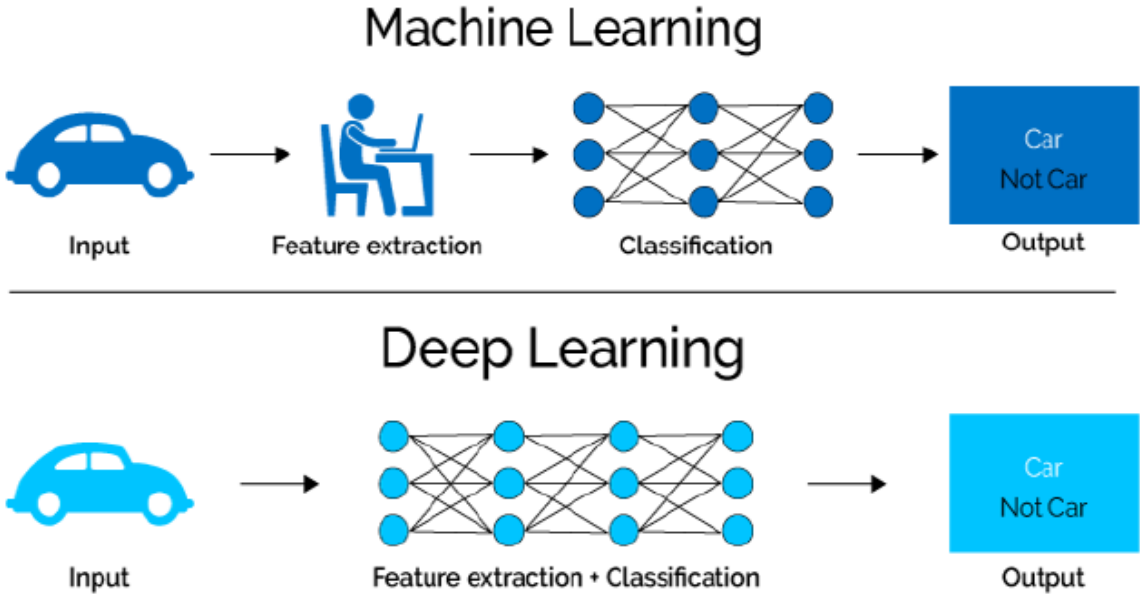
$$z_{k,i} = w_{0,i}^{(k)} + \sum_{j=1}^{d_{k-1}} g(z_{k-1,j}) w_{j,i}^{(k)}$$

# From Neural Network to Deep Learning architecture



Goodfellow, Ian, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.

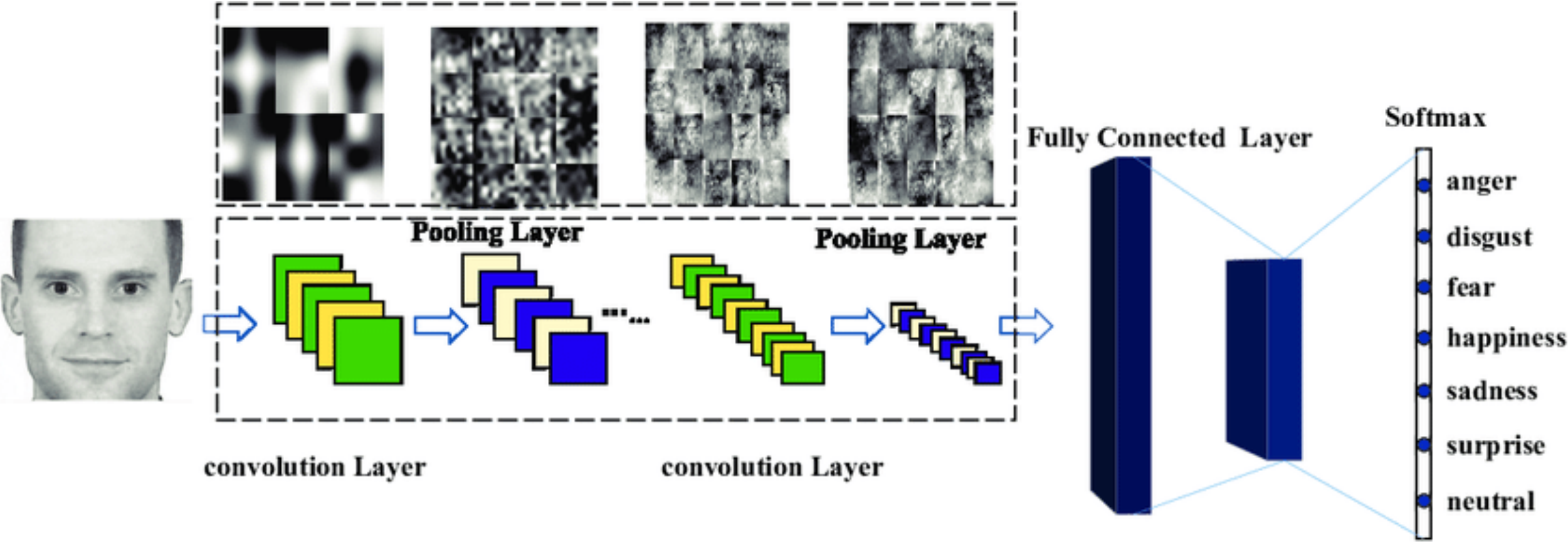
# One of the advantages of DL is to do 'embedded' both the feature extraction and the modeling part



L. Fridman MIT Deep Learning <https://deeplearning.mit.edu/>

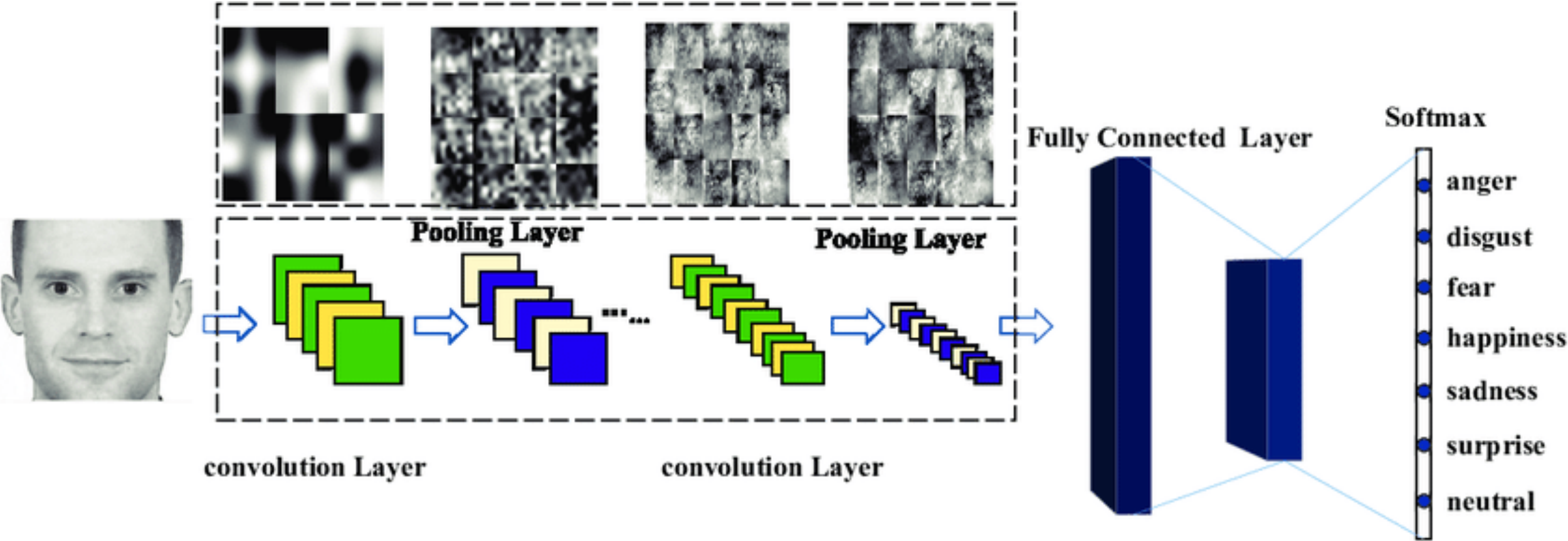


One of the advantages of DL is to do 'embedded' both the feature extraction and the modeling part



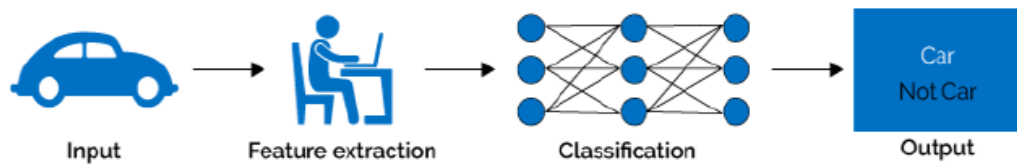


One of the advantages of DL is to do ‘embedded’ both the feature extraction and the modeling part



We are calling these, ‘learned features’

# Handcrafted Features vs Learned Features



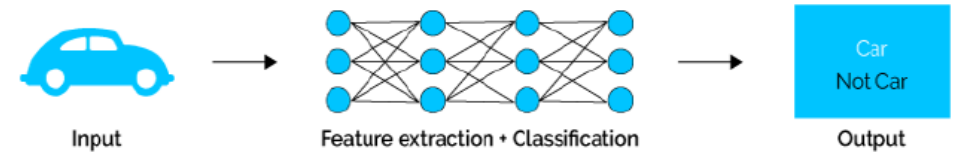
Handcrafted features are manually designed by humans based on domain knowledge and intuition.

Pros:

- Simple, interpretable
- Work well when you have strong domain expertise

Cons:

- Limited expressiveness — can't capture complex patterns
- Often task-specific — poor generalization
- Require manual effort, which doesn't scale



Learned features are automatically extracted by a neural network during training — the network learns how to represent the data in a way that is useful for the task.

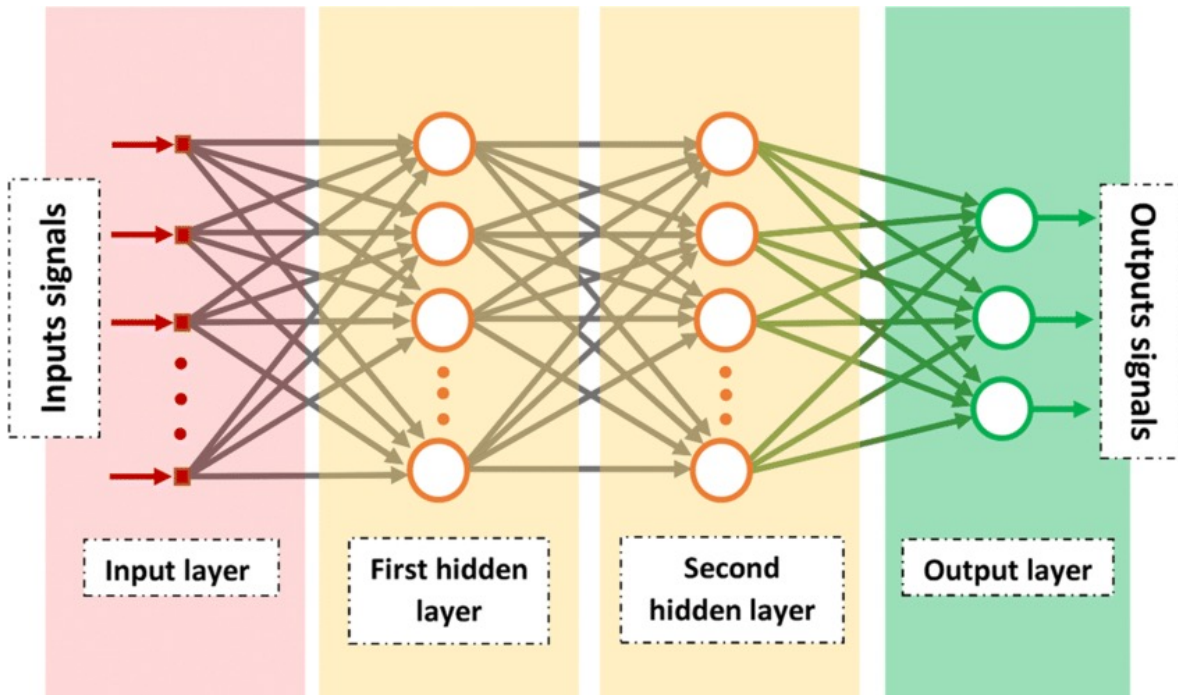
Pros:

- Can learn very complex patterns
- Adapt to the specific data and task
- Scale well to large datasets and varied problems

Cons:

- Less interpretable (but explainability tools help)
- Require more data and compute
- Can overfit if not properly regularized

# Multi-output DL are particularly relevant in some cases

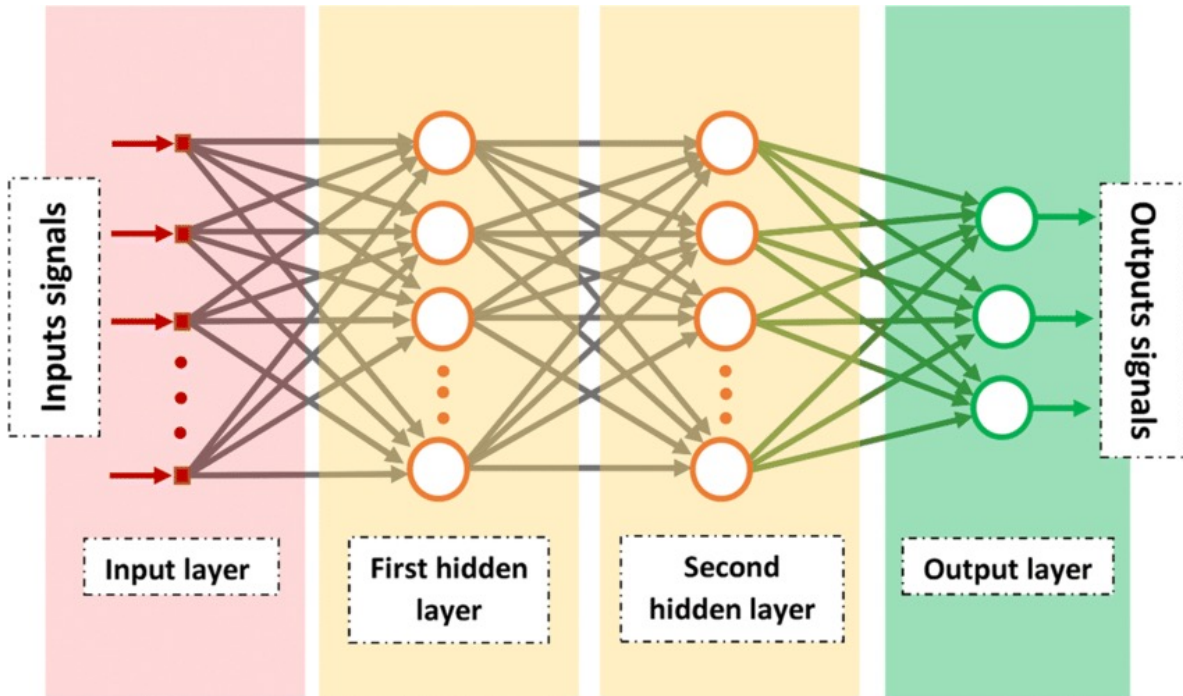


The learned features may be relevant for multiple objectives!

Some benefits:

- Efficiency (one model, one training process)
- Shared knowledge (task can benefit from each other's learned representations)
- Regularization (learning multiple tasks may prevent overfitting)

# Multi-output DL are particularly relevant in some cases



For example, with  $O$  outputs, a regression loss could simply be

$$\mathcal{L}_{\text{MSE}} = \frac{1}{N} \sum_{i=1}^N \sum_{o=1}^O \left( \hat{y}_o^{(i)} - y_o^{(i)} \right)^2$$

The learned features may be relevant for multiple objectives!

Some benefits:

- Efficiency (one model, one training process)
- Shared knowledge (task can benefit from each other's learned representations)
- Regularization (learning multiple tasks may prevent overfitting)

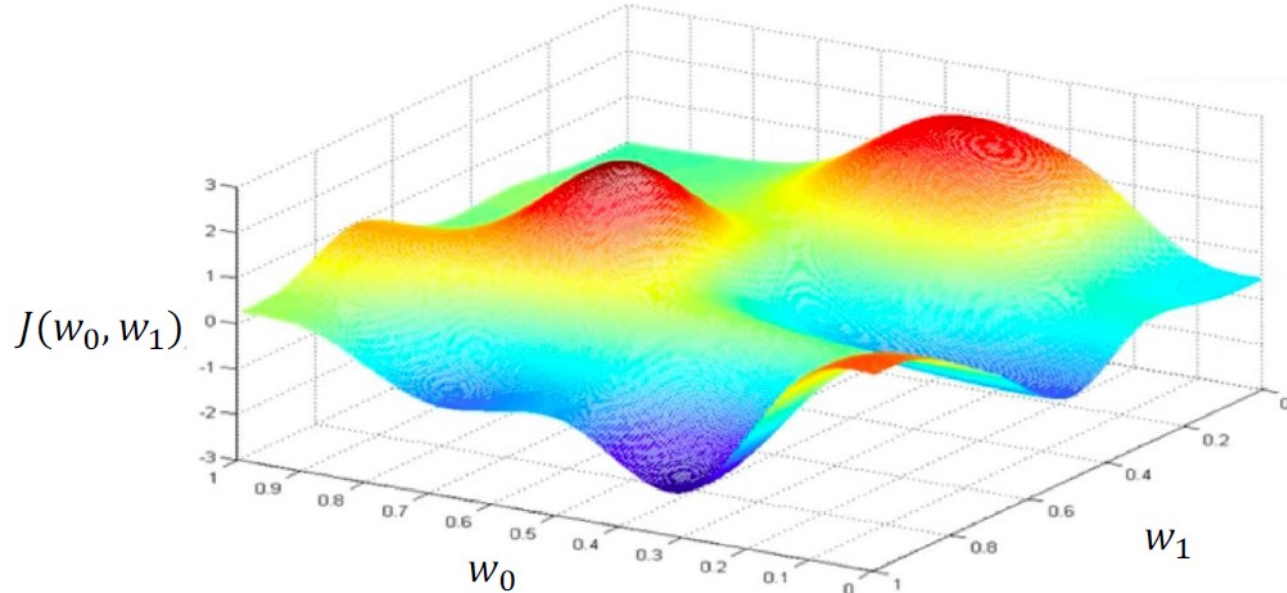
# Deep Learning Training

# Training a Neural Network = Minimizing a Loss

We seek for a set of weights that achieve minimal loss:

$$W^* = \operatorname{argmin}_W \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; W), y^{(i)})$$

$$W^* = \operatorname{argmin}_W J(W)$$



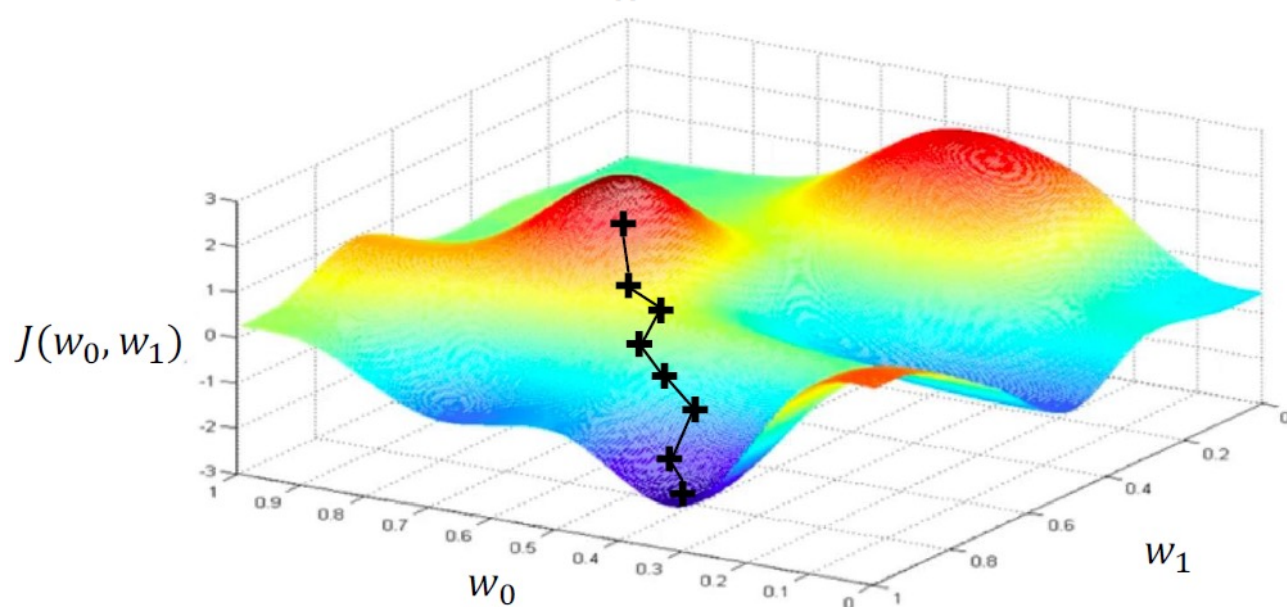


# Training a Neural Network = Minimizing a Loss

We seek for a set of weights that achieve minimal loss:

$$W^* = \operatorname{argmin}_W \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; W), y^{(i)})$$

$$W^* = \operatorname{argmin}_W J(W)$$

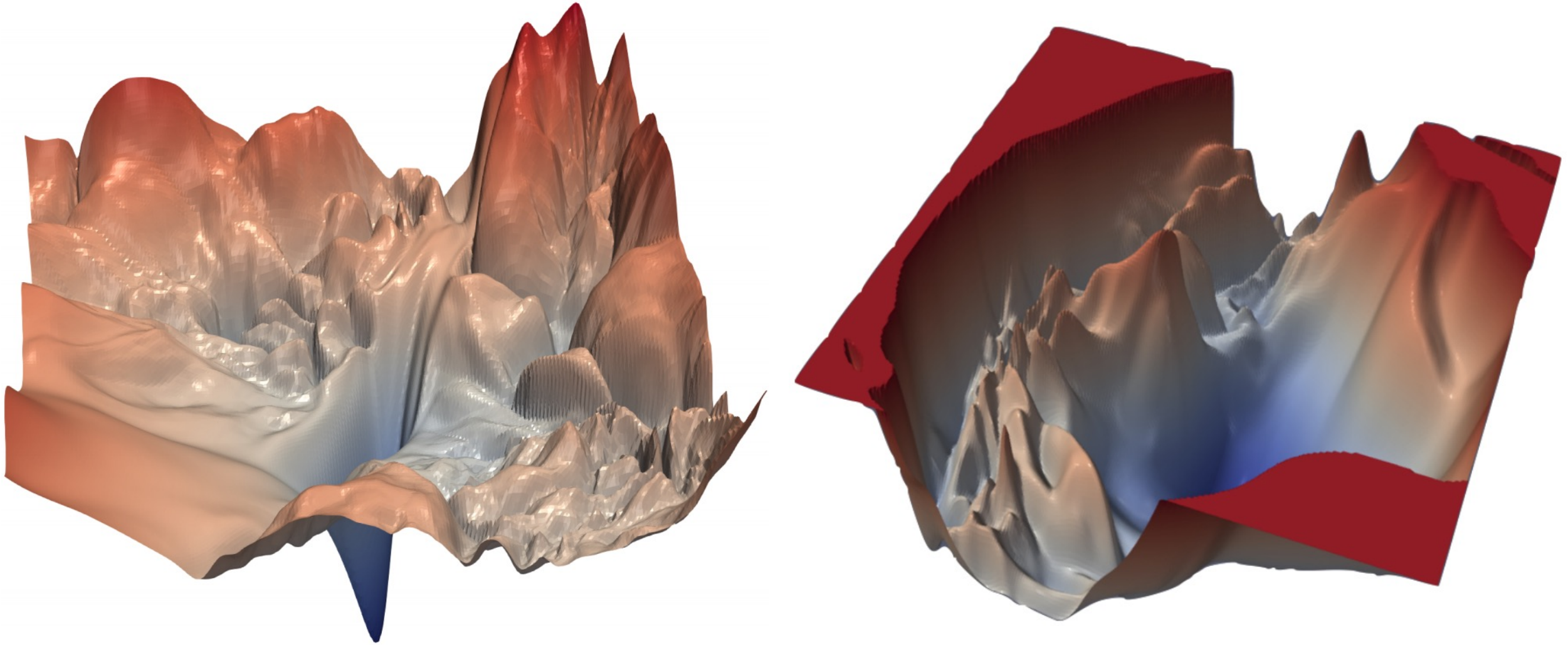


## Algorithm

1. Initialize weights randomly  $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient,  $\frac{\partial J(W)}{\partial W}$
4. Update weights,  $W \leftarrow W - \eta \frac{\partial J(W)}{\partial W}$
5. Return weights



Loss functions may be quite complex...

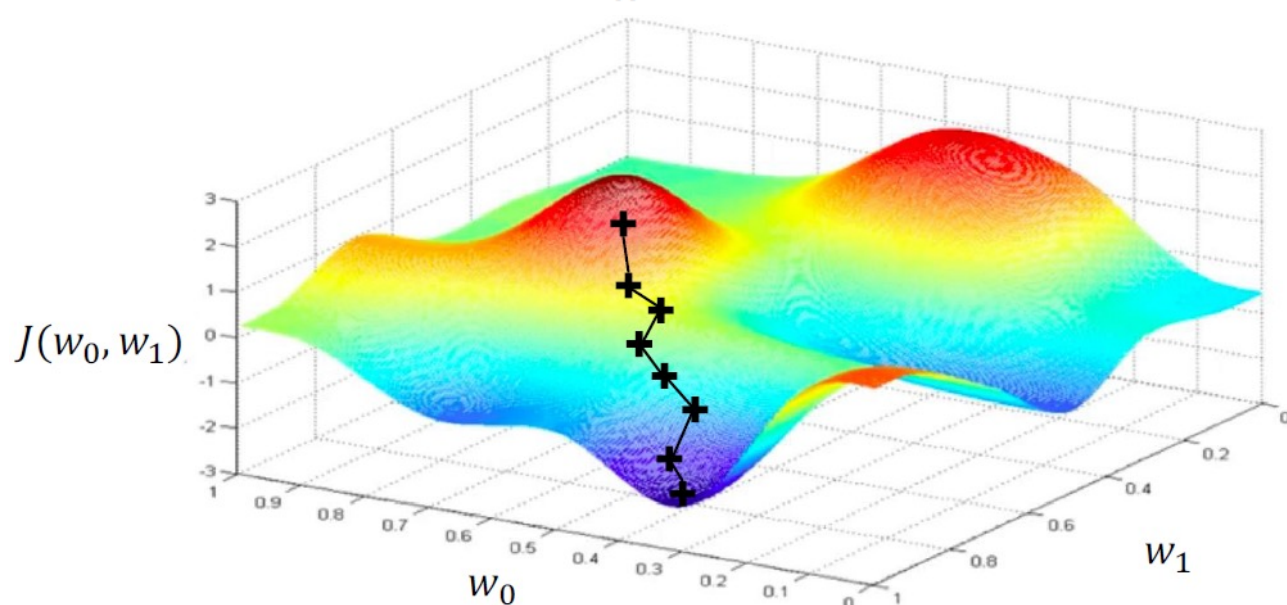


# Training a Neural Network = Minimizing a Loss

We seek for a set of weights that achieve minimal loss:

$$W^* = \operatorname{argmin}_W \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; W), y^{(i)})$$

$$W^* = \operatorname{argmin}_W J(W)$$



## Algorithm (Gradient Descent)

1. Initialize weights randomly  $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient,  $\frac{\partial J(W)}{\partial W}$
4. Update weights,  $W \leftarrow W - \eta \frac{\partial J(W)}{\partial W}$   
Learning Rate
5. Return weights

To cope with complex loss functions and to optimize training the following algorithm is generally sophisticated

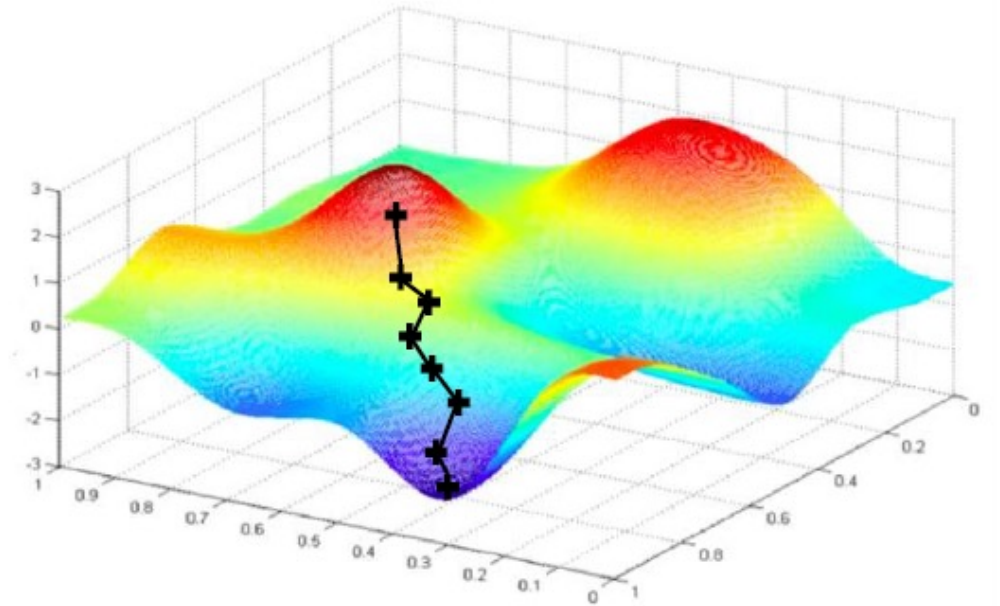
### Algorithm (Gradient Descent)

1. Initialize weights randomly  $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient,  $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4. Update weights,  $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$   
Learning Rate
5. Return weights

# Stochastic Gradient Descent

## Algorithm

1. Initialize weights randomly  $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient,  $\frac{\partial J(W)}{\partial W}$
4. Update weights,  $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(W)}{\partial W}$
5. Return weights



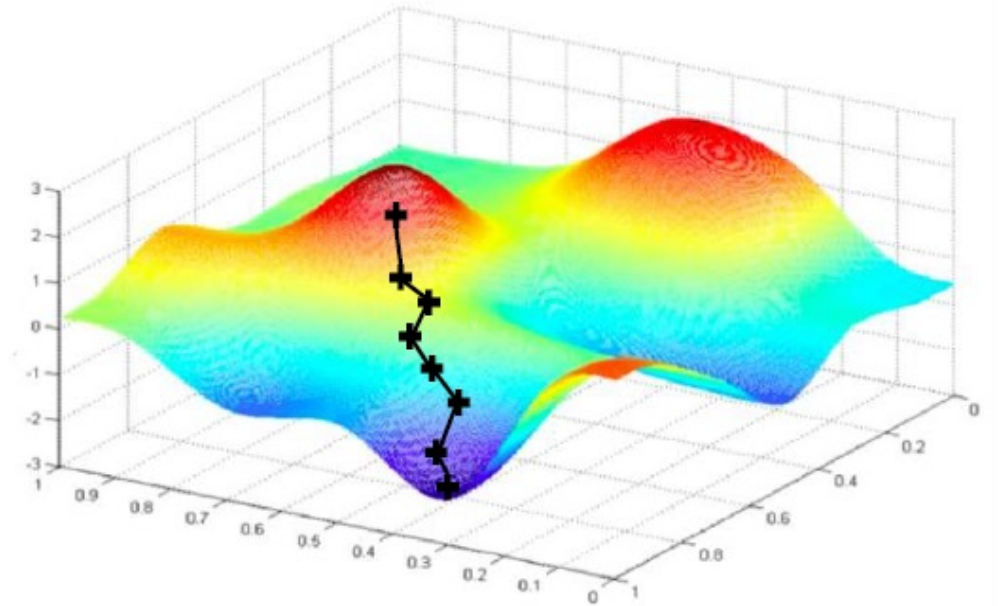
Can be very burdensome  
to compute...  
Average over all samples  
in the dataset!



# Stochastic Gradient Descent

## Algorithm

1. Initialize weights randomly  $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick single data point  $i$
4. Compute gradient,  $\frac{\partial J_i(W)}{\partial W}$
5. Update weights,  $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(W)}{\partial W}$
6. Return weights

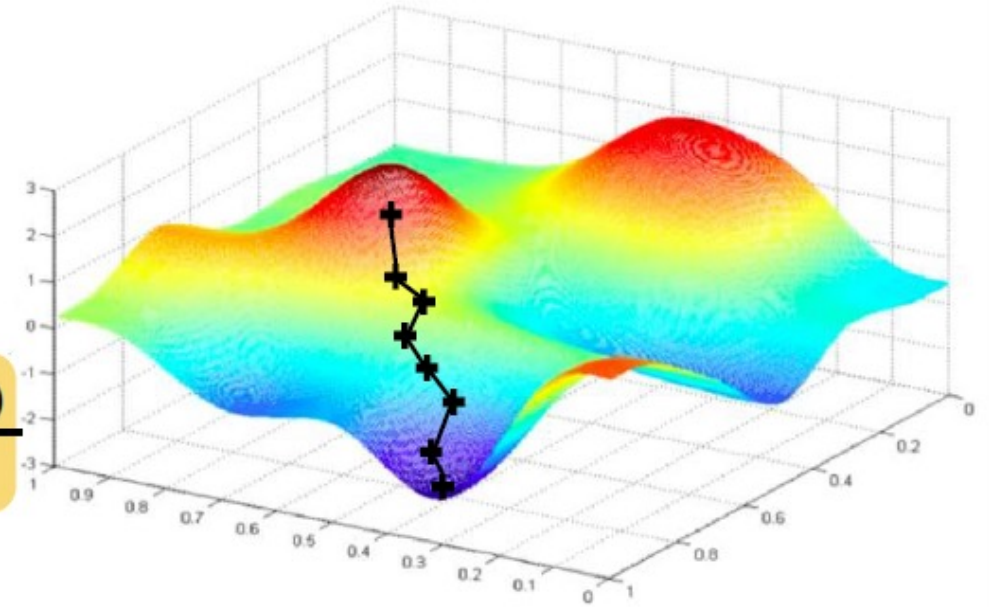


Easy to compute but  
**very noisy**  
(stochastic)!

# Stochastic Gradient Descent

## Algorithm

1. Initialize weights randomly  $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick batch of  $B$  data points
4. Compute gradient,  $\frac{\partial J(W)}{\partial W} = \frac{1}{B} \sum_{k=1}^B \frac{\partial J_k(W)}{\partial W}$
5. Update weights,  $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(W)}{\partial W}$
6. Return weights



Fast to compute and a much better estimate of the true gradient!

# Stochastic Gradient Descent: Mini-batches

The set of  $B$  data points is called *mini-batch*

Mini-batches allow to accurated estimation of gradient, smoother convergence, larger learning rates

Mini-batches lead to fast training: computation can be parallelized and significant speed increases can be obtained on GPUs

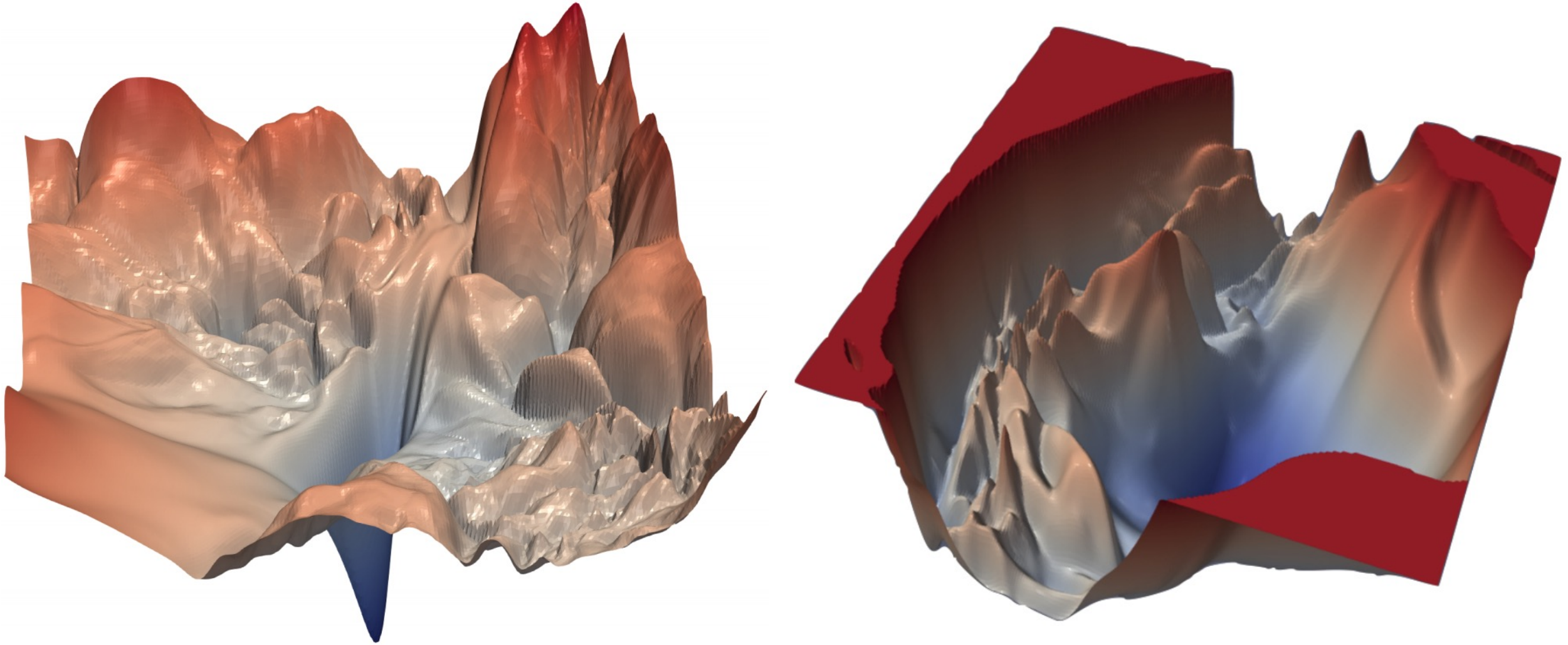




In large architectures, smart choices related to 'simple' things can have a huge impact



Loss functions may be quite complex...

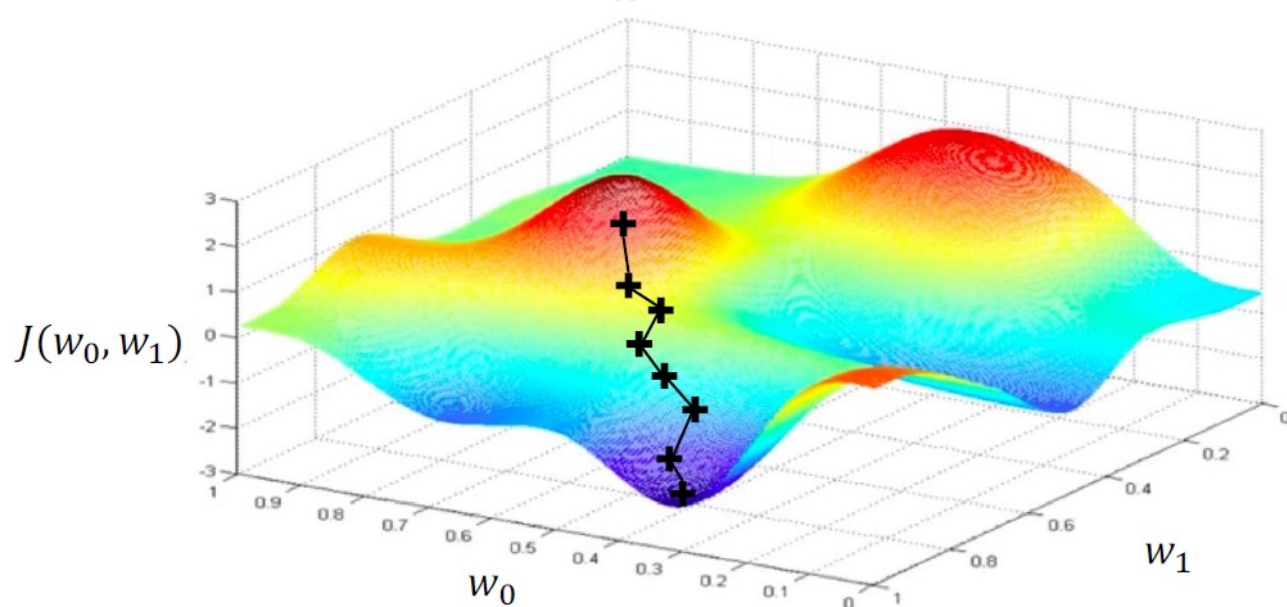


# Training a Neural Network = Minimizing a Loss

We seek for a set of weights that achieve minimal loss:

$$W^* = \operatorname{argmin}_W \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; W), y^{(i)})$$

$$W^* = \operatorname{argmin}_W J(W)$$

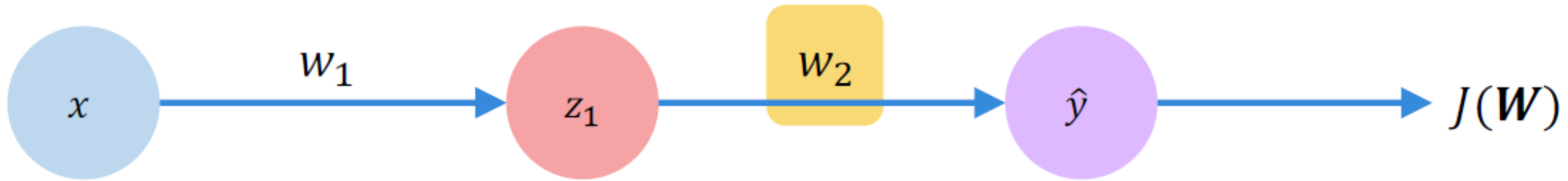


## Algorithm

1. Initialize weights randomly  $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient,  $\frac{\partial J(W)}{\partial W}$
4. Update weights,  $W \leftarrow W - \eta \frac{\partial J(W)}{\partial W}$
5. Return weights

# How to compute the gradient: **backpropagation**

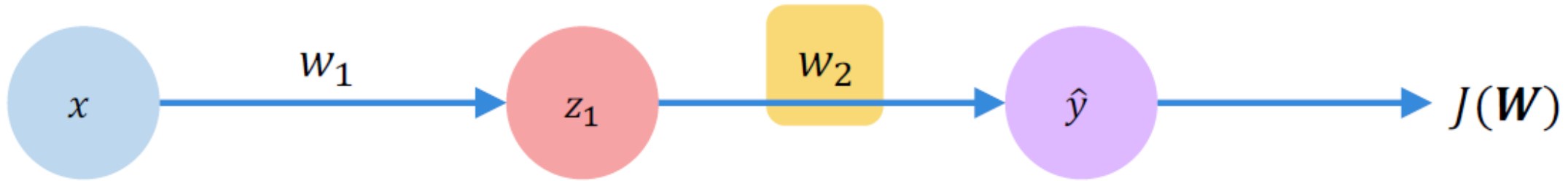
Backpropagation is about understanding how changing the weights and biases in a network changes the cost function.



Let's consider a simple NN with one node: how the final loss is affected by changes in  $w_2$ ?

# How to compute the gradient: backpropagation

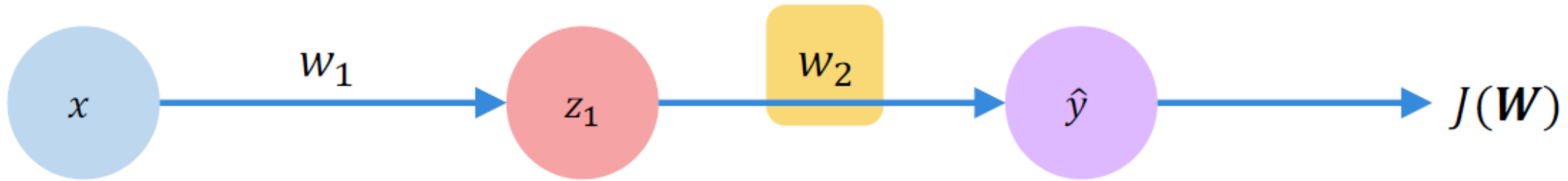
Let's consider a simple NN with one node: how the final loss is affected by changes in  $w_2$



$$\frac{\partial J(W)}{\partial w_2} = ?$$

# How to compute the gradient: **backpropagation**

Let's consider a simple NN with one node: how the final loss is affected by changes in  $w_2$



We can apply the chain rule:

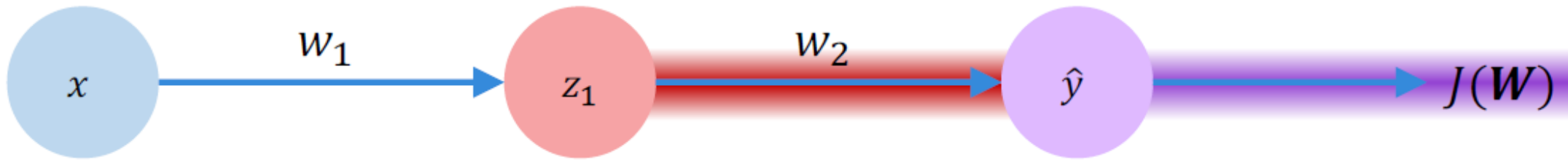
$$\begin{aligned} y &= f(u) \\ u &= g(x) \end{aligned} \quad \longrightarrow \quad \frac{dy}{dx} = \frac{df}{du} \cdot \frac{du}{dx}$$

$$\frac{\partial J(W)}{\partial w_2} = ?$$



# How to compute the gradient: **backpropagation**

Let's consider a simple NN with one node: how the final loss is affected by changes in  $w_2$



We can apply the chain rule:

$$\begin{aligned} y &= f(u) \\ u &= g(x) \end{aligned} \quad \longrightarrow \quad \frac{dy}{dx} = \frac{df}{du} \cdot \frac{du}{dx}$$

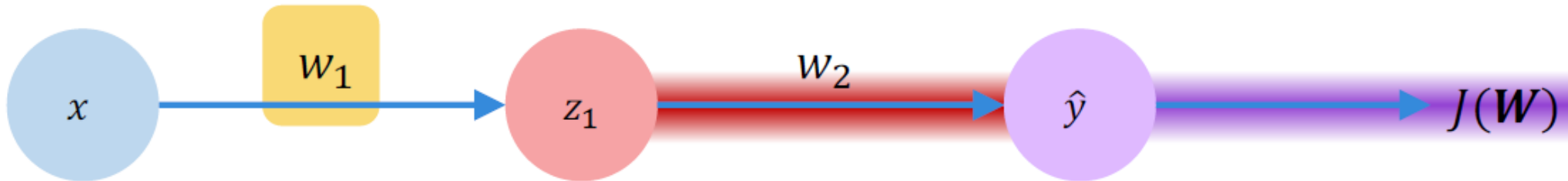
$$\frac{\partial J(W)}{\partial w_2} = \frac{\partial J(W)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w_2}$$



# How to compute the gradient: backpropagation

This simple network is characterized also by the weight

$w_1$



We can apply the chain rule:

$$\begin{aligned} y &= f(u) \\ u &= g(x) \end{aligned} \quad \longrightarrow \quad \frac{dy}{dx} = \frac{df}{du} \cdot \frac{du}{dx}$$

$$\frac{\partial J(W)}{\partial w_1} = \frac{\partial J(W)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w_1}$$

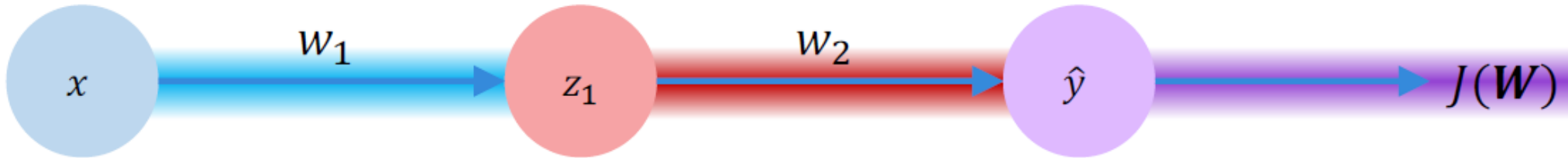
Apply chain rule!

Apply chain rule!

# How to compute the gradient: **backpropagation**

This simple network is characterized also by the weight

$W_1$



We can apply the chain rule:

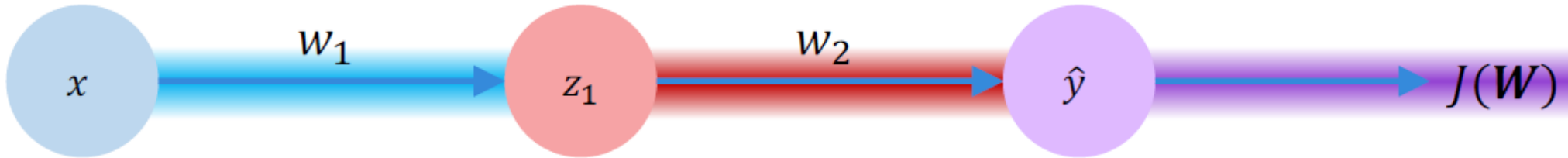
$$\begin{aligned} y &= f(u) \\ u &= g(x) \end{aligned} \quad \longrightarrow \quad \frac{dy}{dx} = \frac{df}{du} \cdot \frac{du}{dx}$$

$$\frac{\partial J(W)}{\partial w_1} = \frac{\partial J(W)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1} * \frac{\partial z_1}{\partial w_1}$$

# How to compute the gradient: **backpropagation**

This simple network is characterized also by the weight

$w_1$



We can apply the chain rule:

$$\begin{matrix} y = f(u) \\ u = g(x) \end{matrix} \quad \longrightarrow \quad \frac{dy}{dx} = \frac{df}{du} \cdot \frac{du}{dx}$$

$$\frac{\partial J(\mathbf{W})}{\partial w_1} = \frac{\partial J(\mathbf{W})}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1} * \frac{\partial z_1}{\partial w_1}$$

$$z_1 = \sigma(x \cdot w_1) \quad (\text{activation})$$

We also have to keep into account the non-linear transformations (activations)! In that case, we need to apply again the chain rule!

# How to compute the gradient: **backpropagation**

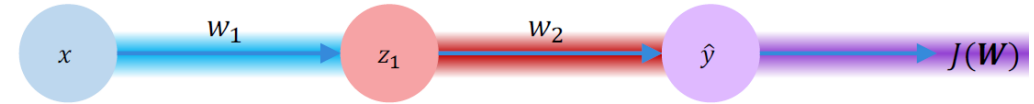
Let's consider an example with sigmoid activation

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

$$z_1 = \sigma(x \cdot w_1) \quad (\text{activation})$$

$$\hat{y} = z_1 \cdot w_2 \quad (\text{linear output})$$

$$L = \frac{1}{2}(\hat{y} - y)^2 \quad (\text{loss})$$



We can apply the chain rule:

$$\begin{aligned} y &= f(u) \\ u &= g(x) \end{aligned} \quad \longrightarrow \quad \frac{dy}{dx} = \frac{df}{du} \cdot \frac{du}{dx}$$

# How to compute the gradient: **backpropagation**

We can apply the chain rule:

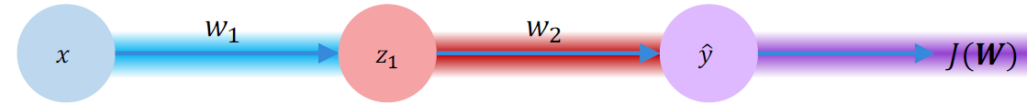
$$\begin{aligned} y &= f(u) \\ u &= g(x) \end{aligned} \quad \longrightarrow \quad \frac{dy}{dx} = \frac{df}{du} \cdot \frac{du}{dx}$$

Let's consider an example with sigmoid activation

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

$$z_1 = \sigma(x \cdot w_1) \quad (\text{activation})$$

$$\hat{y} = z_1 \cdot w_2 \quad (\text{linear output})$$



$$L = \frac{1}{2}(\hat{y} - y)^2 \quad (\text{loss})$$

With  $x=1$ ,  $w_1 = 0.5$ ,  $w_2=-1$ ,  $y=1$ :

$$z_1 = \sigma(x \cdot w_1) = \sigma(1.0 \cdot 0.5) = \sigma(0.5) \approx 0.622$$

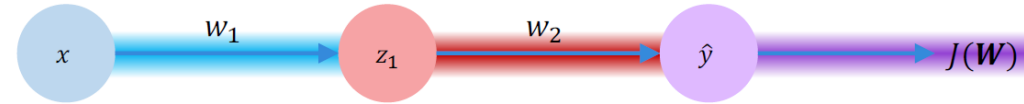
$$\hat{y} = z_1 \cdot w_2 = 0.622 \cdot (-1.0) = -0.622$$

$$L = \frac{1}{2}(\hat{y} - y)^2 = \frac{1}{2}(-0.622 - 1)^2 = \frac{1}{2}(-1.622)^2 \approx 1.315$$

# How to compute the gradient: **backpropagation**

We can apply the chain rule:

$$\begin{aligned} y &= f(u) \\ u &= g(x) \end{aligned} \quad \longrightarrow \quad \frac{dy}{dx} = \frac{df}{du} \cdot \frac{du}{dx}$$



$$\frac{dL}{dw_2} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dw_2}$$

$$\frac{dL}{d\hat{y}} = \hat{y} - y = -0.622 - 1 = -1.622$$

$$\frac{d\hat{y}}{dw_2} = z_1 = 0.622$$

$$\frac{dL}{dw_2} = -1.622 \cdot 0.622 \approx -1.009$$

# How to compute the gradient: **backpropagation**

$$\frac{dL}{dw_1} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dz_1} \cdot \frac{dz_1}{dw_1}$$

Breaking it down:

- $\frac{dL}{d\hat{y}} = -1.622$
- $\frac{d\hat{y}}{dz_1} = w_2 = -1.0$
- $\frac{dz_1}{dw_1} = \sigma'(x \cdot w_1) \cdot x = z_1(1 - z_1) \cdot x$

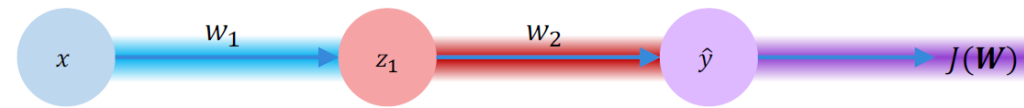
Using:

- $z_1 = 0.622$
- $\sigma'(z) = \sigma(z)(1 - \sigma(z)) \approx 0.622 \cdot (1 - 0.622) \approx 0.235$

Then:

$$\frac{dz_1}{dw_1} = 0.235 \cdot 1.0 = 0.235$$
$$\frac{dL}{dw_1} = (-1.622) \cdot (-1.0) \cdot 0.235 \approx 0.381$$

We can apply the chain rule:  
 $y = f(u)$   
 $u = g(x)$   $\rightarrow$   $\frac{dy}{dx} = \frac{df}{du} \cdot \frac{du}{dx}$



$$\frac{dL}{dw_2} = \frac{dL}{d\hat{y}} \cdot \frac{d\hat{y}}{dw_2}$$

$$\frac{dL}{d\hat{y}} = \hat{y} - y = -0.622 - 1 = -1.622$$

$$\frac{d\hat{y}}{dw_2} = z_1 = 0.622$$

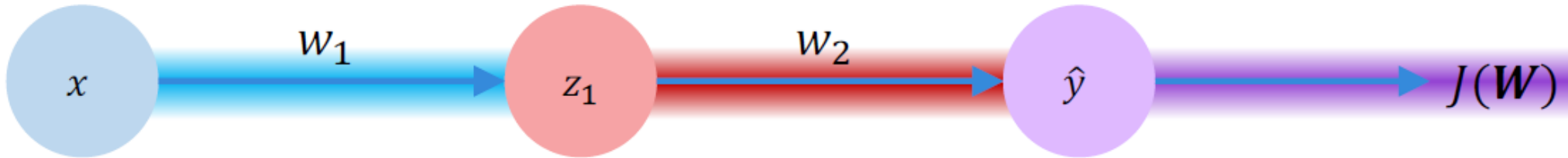
$$\frac{dL}{dw_2} = -1.622 \cdot 0.622 \approx -1.009$$



# How to compute the gradient: backpropagation

This simple network is characterized also by the weight

$w_1$



We can apply the chain rule:

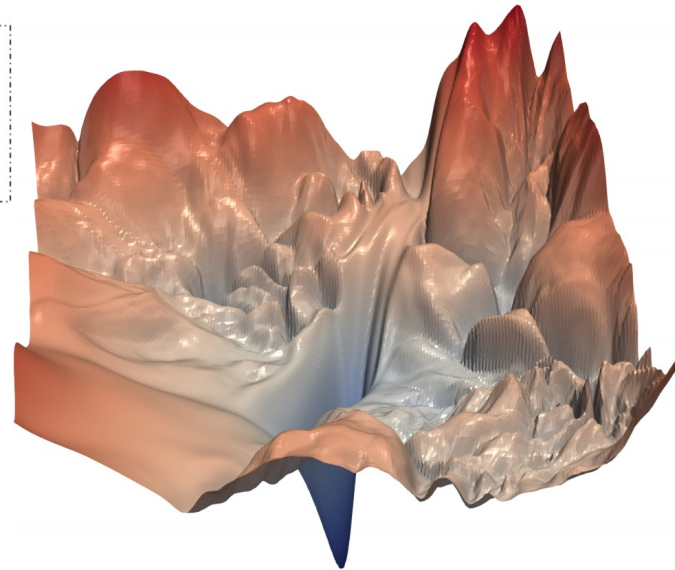
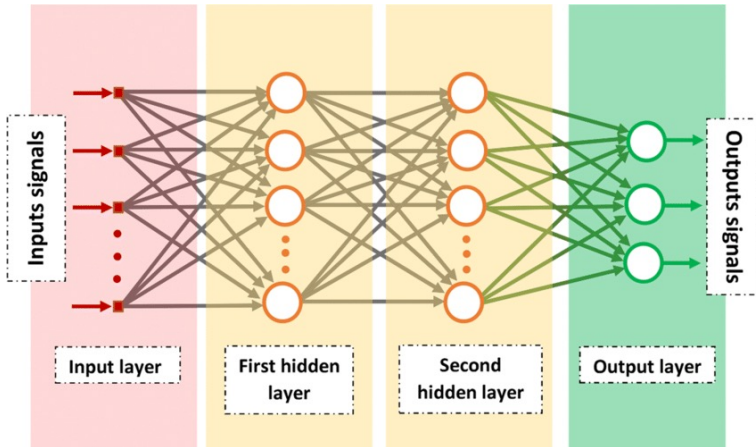
$$\begin{matrix} y = f(u) \\ u = g(x) \end{matrix} \longrightarrow \frac{dy}{dx} = \frac{df}{du} \cdot \frac{du}{dx}$$

$$\frac{\partial J(W)}{\partial w_1} = \frac{\partial J(W)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1} * \frac{\partial z_1}{\partial w_1}$$

You can have a look here for the math:  
<http://neuralnetworksanddeeplearning.com/chap2.html>

# When dealing with NN training, we are dealing with a 'beast'

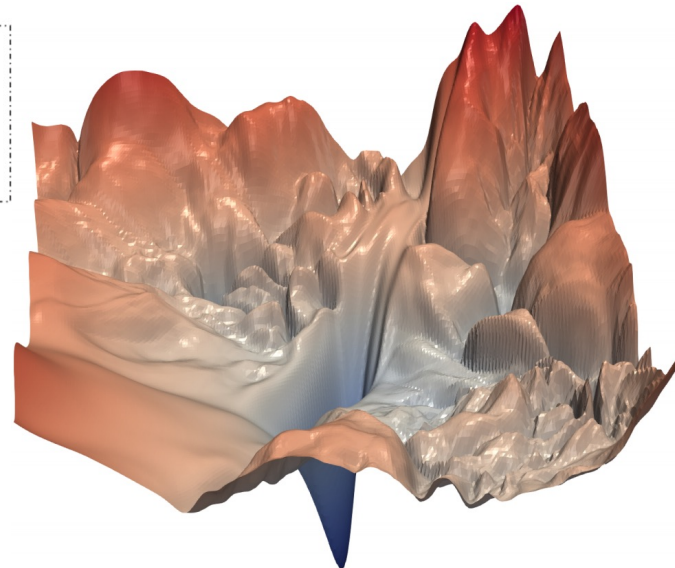
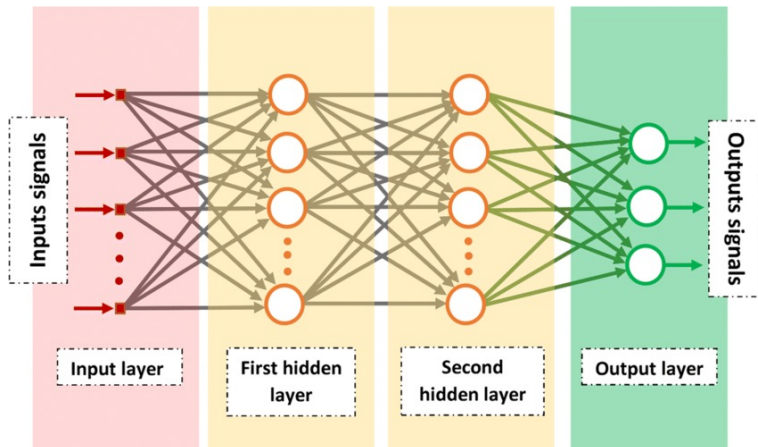
Stochastic Gradient Descent and Back-propagation have been there for decades, but we weren't able to properly train such architectures!



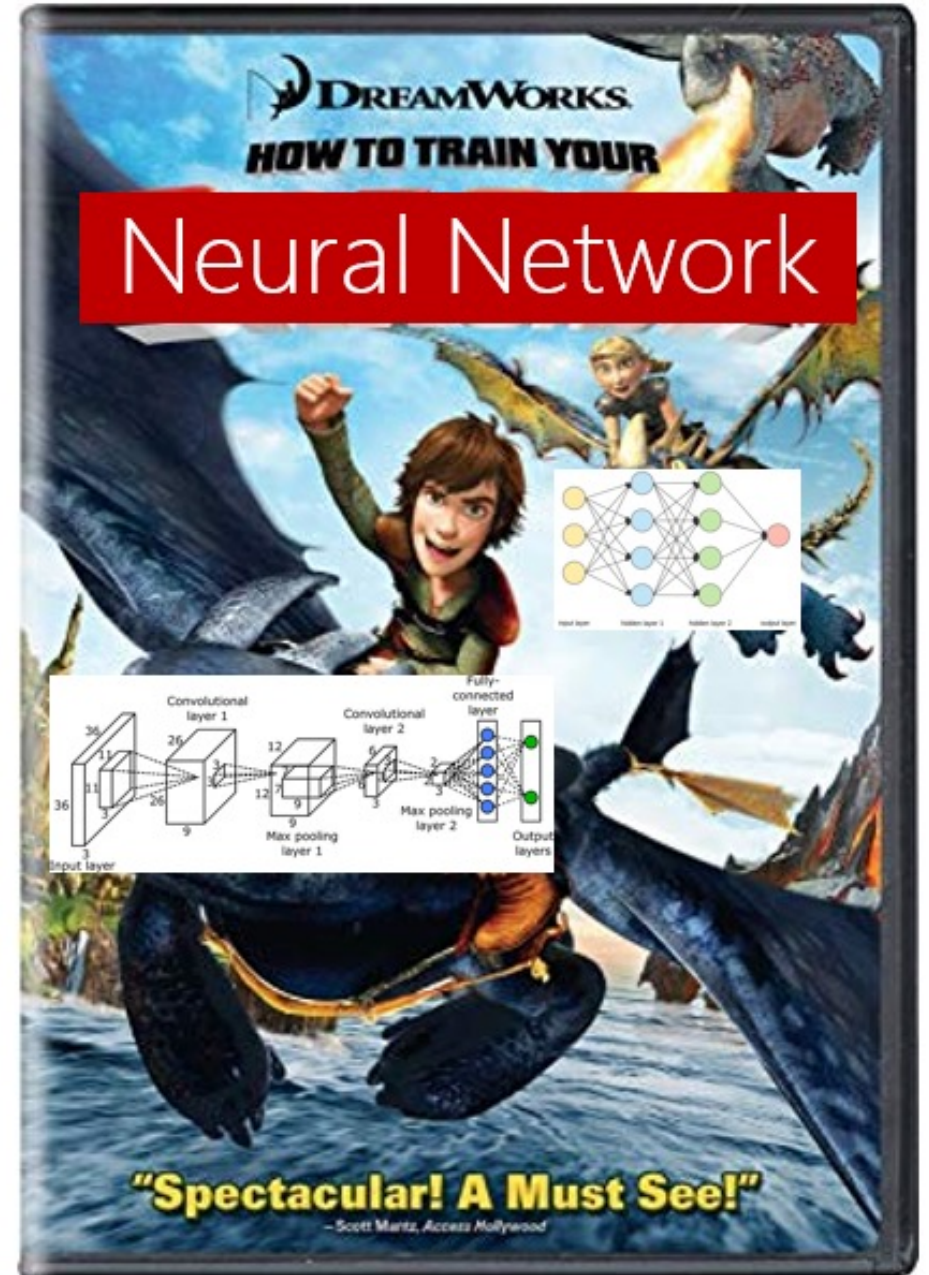
## When dealing with NN training, we are dealing with a 'beast'

Stochastic Gradient Descent and Back-propagation have been there for decades, but we weren't able to properly train such architectures!

In past recent years we have developed many tricks to 'tame the beast'

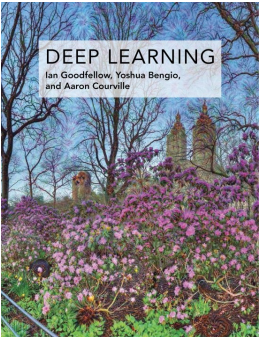


1. Activation functions
2. Weight Initialization
3. Batch Normalization
4. Optimization
5. Learning Rate
6. Regularization

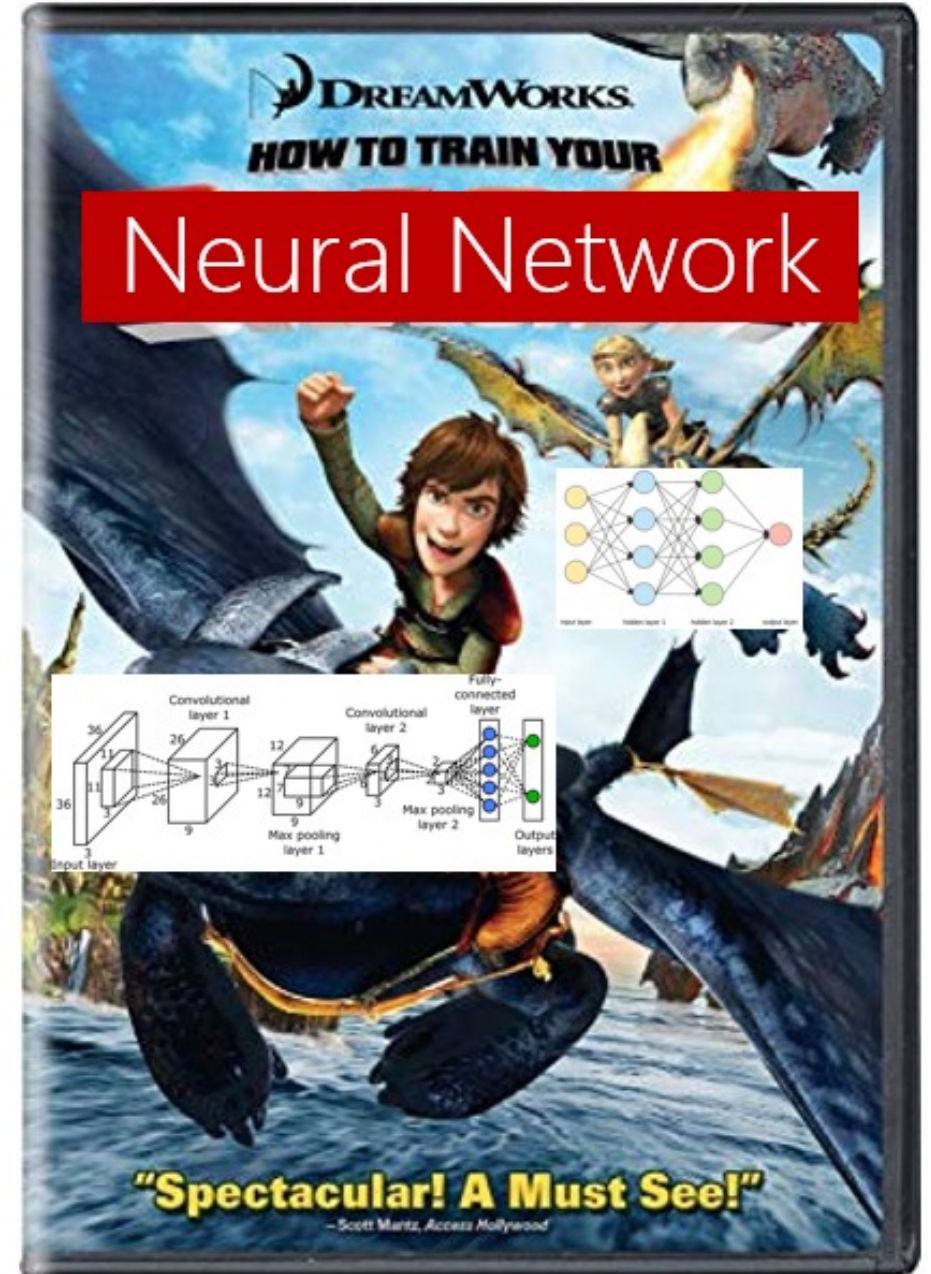




1. Activation functions
2. Weight Initialization
3. Batch Normalization
4. Optimization
5. Learning Rate
6. Regularization



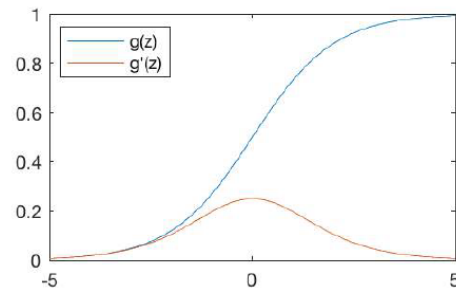
Chapter 6.3 Hidden Units



# Activation functions

Beside introducing non-linearities, it is important for Activation functions to have an easy way to compute the gradients (we need this in backpropagation)

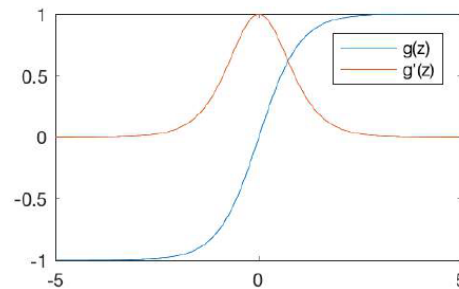
Sigmoid Function



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

$$g'(z) = g(z)(1 - g(z))$$

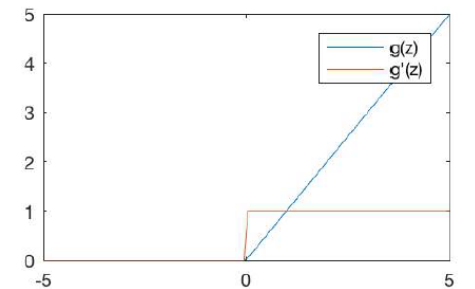
Hyperbolic Tangent



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

$$g'(z) = 1 - g(z)^2$$

Rectified Linear Unit (ReLU)

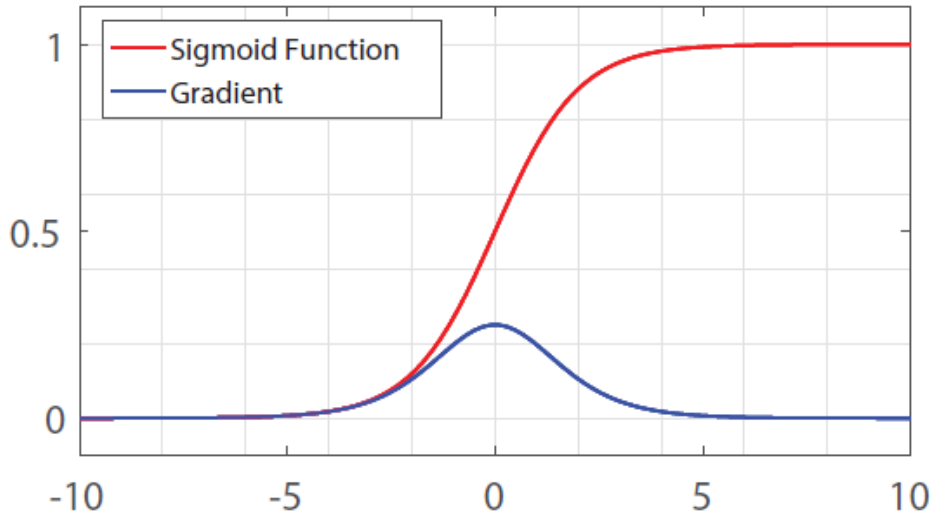


$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

# Sigmoid

Simple interpretation (probability)



*G. Roffo Ranking to Learn and Learning to Rank: On the Role of Ranking in Pattern Recognition Applications*

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

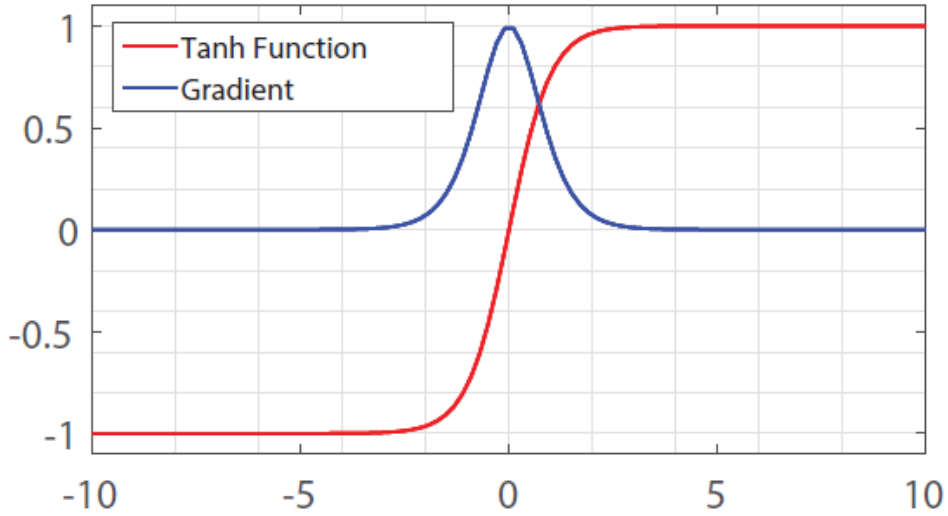
$$g'(z) = g(z)(1 - g(z))$$

Issues:

- **Vanishing gradient problem:** gradient becomes increasingly small as the absolute value of the input increases (it is a problem in backpropagation)
- No zero-centered output (zig-zagging dynamics in the gradient updates)
- Exp function can be expensive to compute



# Tanh



*G. Roffo Ranking to Learn and Learning to Rank: On the Role of Ranking in Pattern Recognition Applications*

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

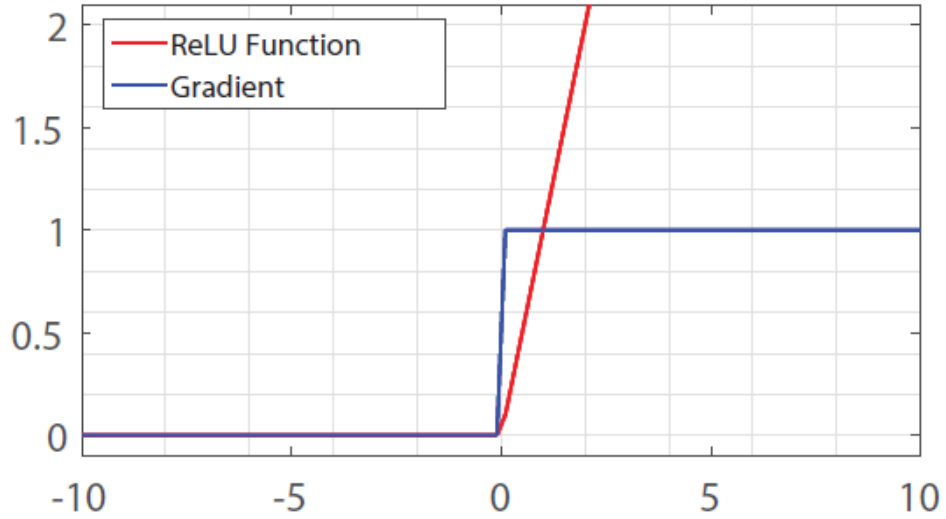
$$g'(z) = 1 - g(z)^2$$

Zero-centered

Issues:

- Vanishing gradient problem
- Exp function can be expensive to compute

# ReLU



G. Roffo *Ranking to Learn and Learning to Rank: On the Role of Ranking in Pattern Recognition Applications*

$$g(z) = \max(0, z)$$

$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

Simple implementation

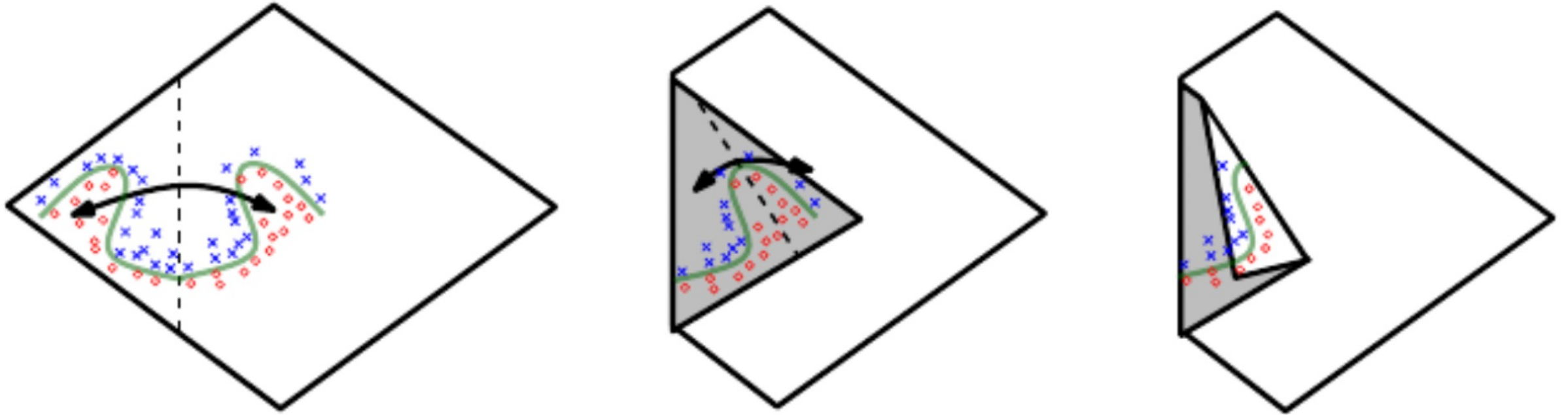
Does not saturate (in + region we have no vanishing gradient)

Faster convergence of SGD than sigmoid/tanh

Issues:

- Not zero-centered
- Gradient is zero for negative values: dead ReLU (as much as 40% of the network never activate if the learning rate is too high)

# ReLU: Intuition

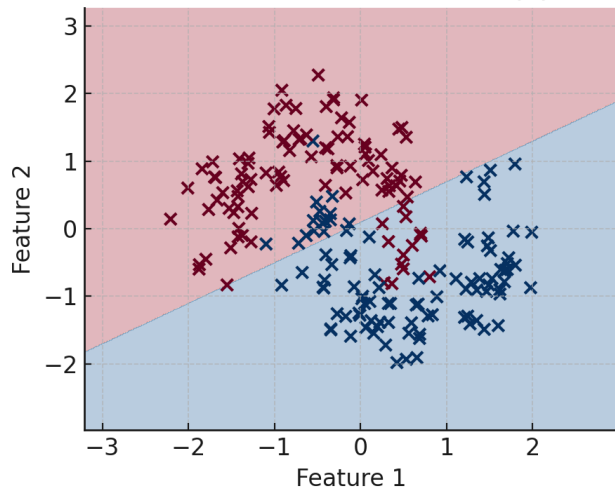


ReLU gives neural networks the power to ‘fold’ the input space.

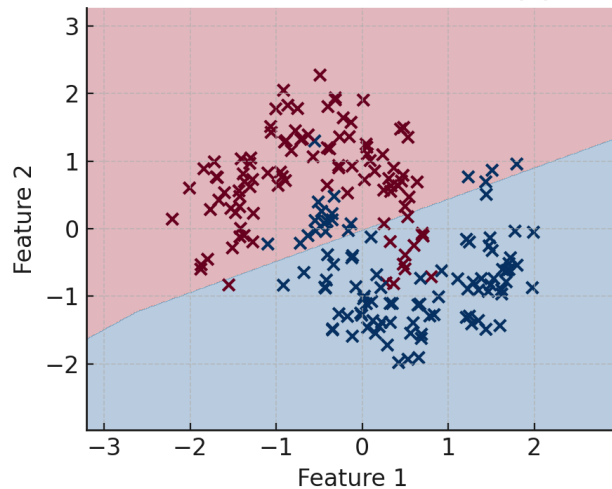
This lets them transform complex, tangled data into something that even a straight line can separate.

# 1 hidden layer

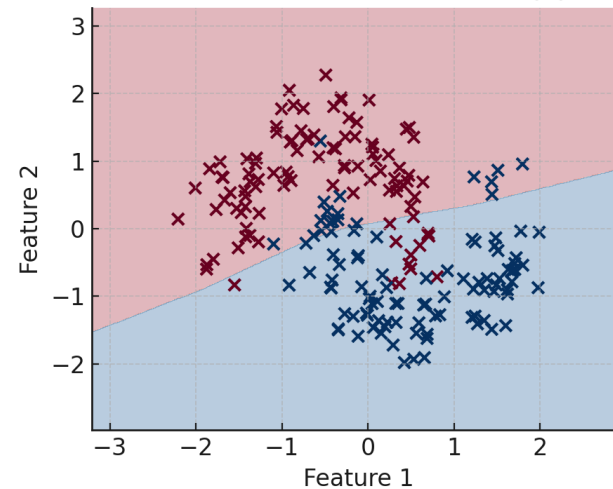
ReLU NN with 1 neuron(s)



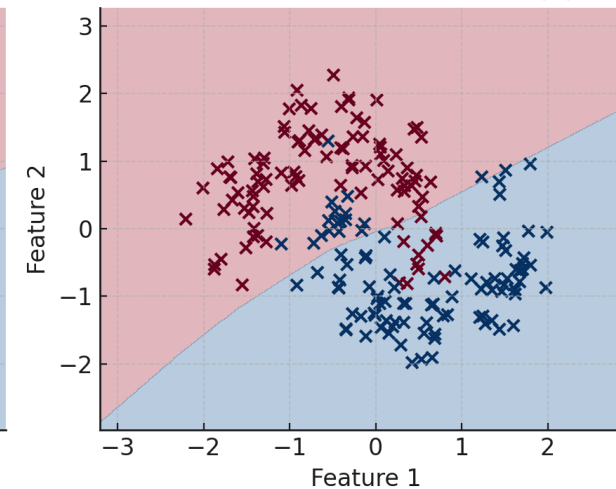
ReLU NN with 3 neuron(s)



ReLU NN with 10 neuron(s)

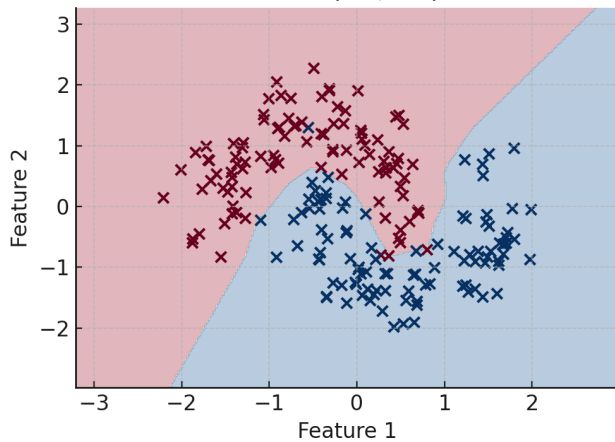


ReLU NN with 20 neuron(s)

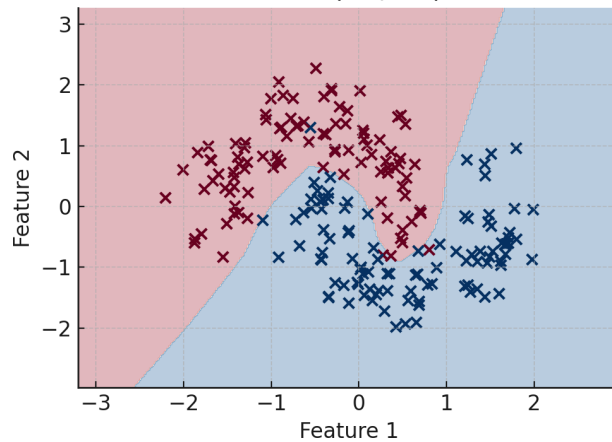


# 2 hidden layers

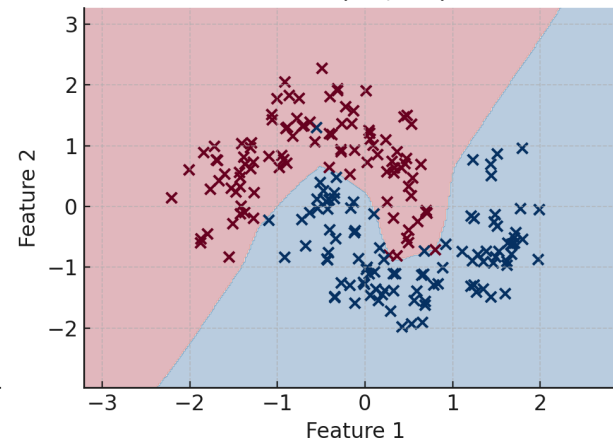
ReLU NN with (10, 10) neurons



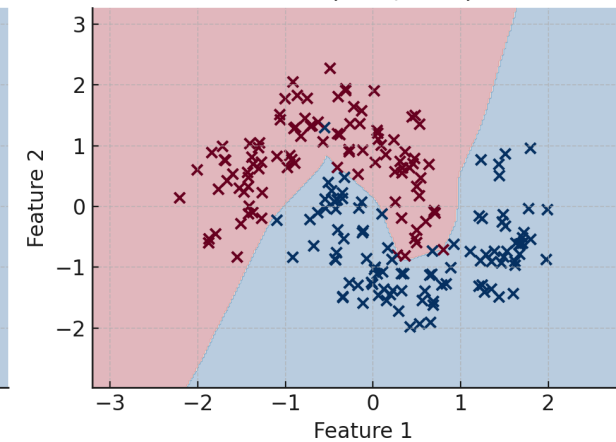
ReLU NN with (20, 20) neurons



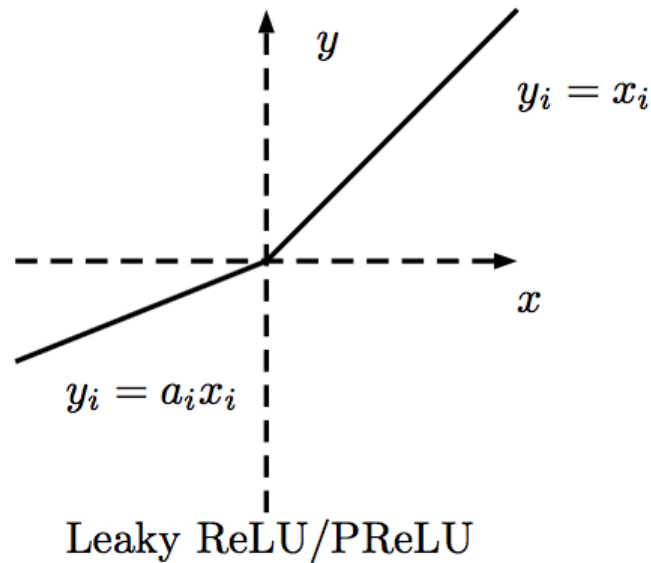
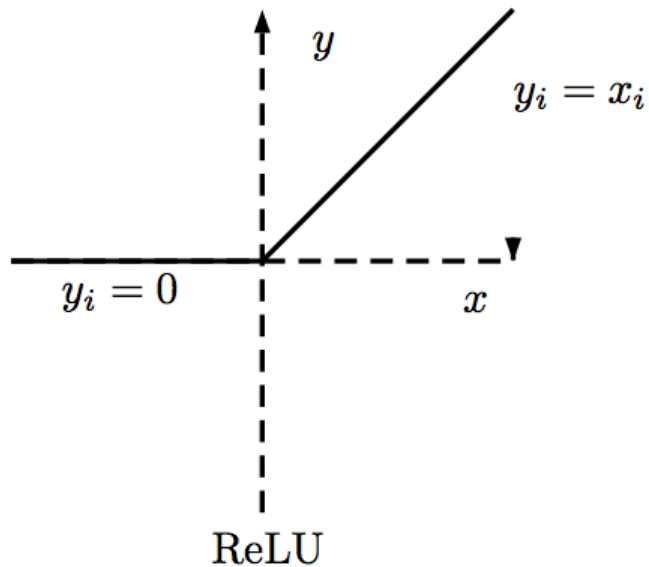
ReLU NN with (50, 50) neurons



ReLU NN with (100, 100) neurons



# Leaky ReLU, Parametric ReLU, Randomized Leaky ReLU



$$g(z) = 0.01z, (z < 0)$$

$$g(z) = z, (z \geq 0)$$

$$g(z) = \alpha z, (z < 0)$$

$$g(z) = z, (z \geq 0)$$

Very easy to compute

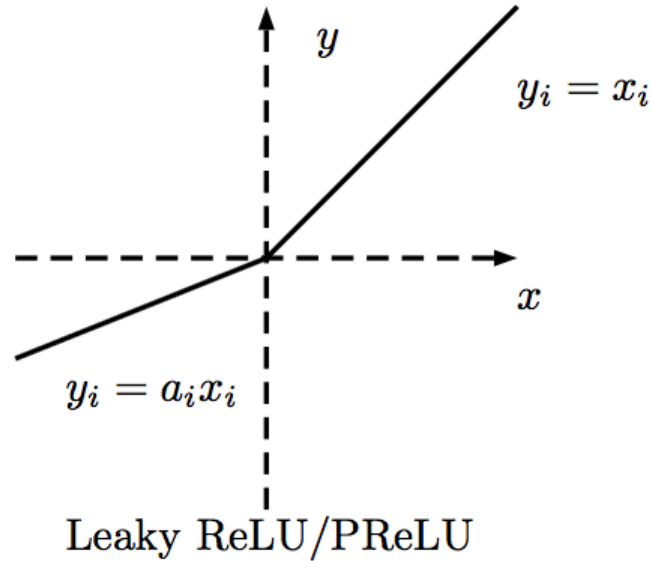
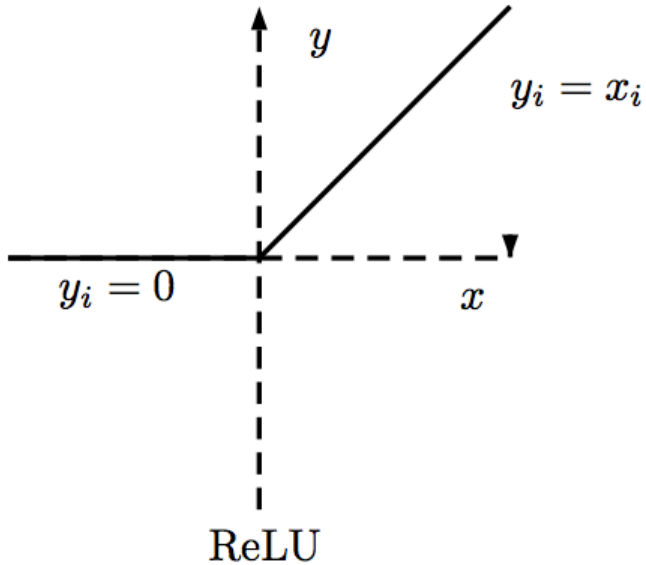
Does not saturate (in + region)

Faster convergence than sigmoid/tanh

**Does not die!**

In the Parametric ReLU (PReLU) the parameter alpha is learned along with the other network parameters

# Leaky ReLU, Parametric ReLU, Randomized Leaky ReLU

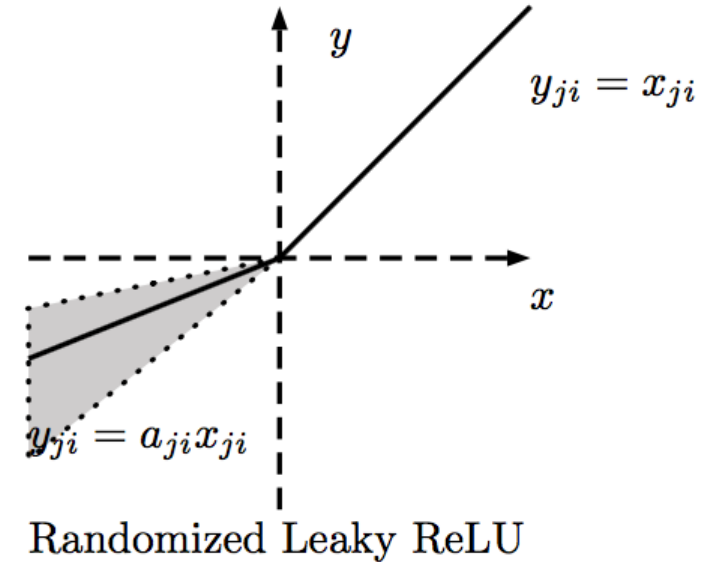


$$g(z) = 0.01z, (z < 0)$$

$$g(z) = z, (z \geq 0)$$

$$g(z) = \alpha z, (z < 0)$$

$$g(z) = z, (z \geq 0)$$



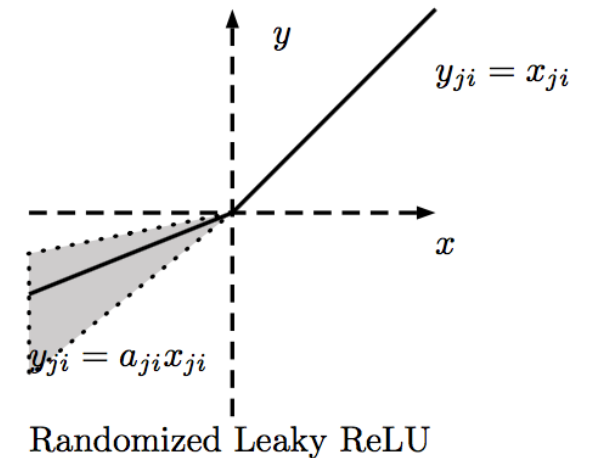
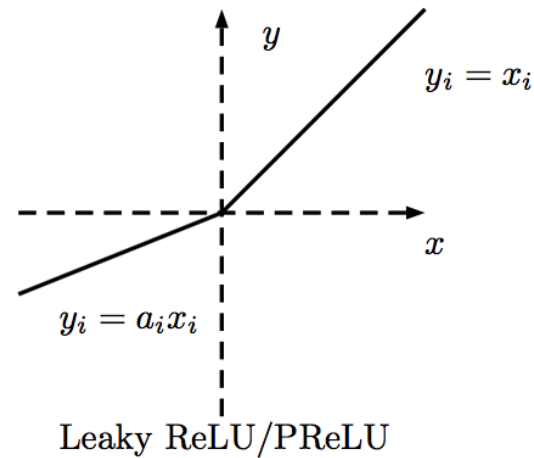
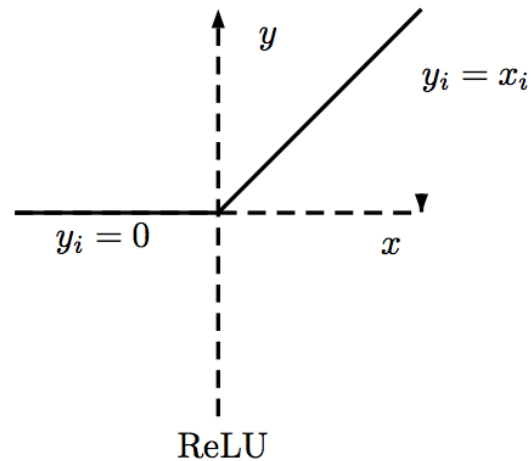
Alpha is chosen at random in a given range in training and it is then fixed in test



# TLDR: Activations

Main issues with activation functions:

- Vanishing gradients
- Non-centered on zero outputs
- Costly computations



[https://isaacchanghau.github.io/post/activation\\_functions/](https://isaacchanghau.github.io/post/activation_functions/)

Name	Plot	Equation	Derivative (with respect to x)	Range	Order of continuity	Monotonic	Monotonic derivative	Approximates identity near the origin
Identity		$f(x) = x$	$f'(x) = 1$	$(-\infty, \infty)$	$C^\infty$	Yes	Yes	Yes
Binary step		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x \neq 0 \\ ? & \text{for } x = 0 \end{cases}$	$\{0, 1\}$	$C^{-1}$	Yes	No	No
Logistic (a.k.a. Sigmoid or Soft step)		$f(x) = \sigma(x) = \frac{1}{1 + e^{-x}}$ [1]	$f'(x) = f(x)(1 - f(x))$	$(0, 1)$	$C^\infty$	Yes	No	No
Tanh		$f(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$	$f'(x) = 1 - f(x)^2$	$(-1, 1)$	$C^\infty$	Yes	No	Yes
Arctan		$f(x) = \tan^{-1}(x)$	$f'(x) = \frac{1}{x^2 + 1}$	$(-\frac{\pi}{2}, \frac{\pi}{2})$	$C^\infty$	Yes	No	Yes
ArSinh		$f(x) = \sinh^{-1}(x) = \ln(x + \sqrt{x^2 + 1})$	$f'(x) = \frac{1}{\sqrt{x^2 + 1}}$	$(-\infty, \infty)$	$C^\infty$	Yes	No	Yes
EluSeg [8][9][10] Softsign [11][12]		$f(x) = \frac{x}{1 +  x }$	$f'(x) = \frac{1}{(1 +  x )^2}$	$(-1, 1)$	$C^1$	Yes	No	Yes
Inverse square root unit (ISRU) [13]		$f(x) = \frac{x}{\sqrt{1 + \alpha x^2}}$	$f'(x) = \left( \frac{1}{\sqrt{1 + \alpha x^2}} \right)'$	$(-\frac{1}{\sqrt{\alpha}}, \frac{1}{\sqrt{\alpha}})$	$C^\infty$	Yes	No	Yes
Inverse square root linear unit (ISRLU) [13]		$f(x) = \begin{cases} \frac{x}{\sqrt{1 - \alpha x^2}} & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} \left( \frac{1}{\sqrt{1 - \alpha x^2}} \right)'$ & for $x < 0$ \\ 1 & for $x \geq 0$ \end{cases}	$(-\frac{1}{\sqrt{\alpha}}, \infty)$	$C^0$	Yes	Yes	Yes
Square Nonlinearity (SQNL) [10]		$f(x) = \begin{cases} 1 & : x > 2.0 \\ x - \frac{x^2}{4} & : 0 \leq x \leq 2.0 \\ x + \frac{x^2}{4} & : -2.0 \leq x < 0 \\ -1 & : x < -2.0 \end{cases}$	$f'(x) = 1 \mp \frac{x}{2}$	$(-1, 1)$	$C^0$	Yes	No	Yes
Rectified linear unit (ReLU) [14]		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$[0, \infty)$	$C^0$	Yes	Yes	No
Bipolar rectified linear unit (BReLU) [15]		$f(x_i) = \begin{cases} \text{ReLU}(x_i) & \text{if } i \bmod 2 = 0 \\ -\text{ReLU}(-x_i) & \text{if } i \bmod 2 \neq 0 \end{cases}$	$f'(x_i) = \begin{cases} \text{ReLU}'(x_i) & \text{if } i \bmod 2 = 0 \\ -\text{ReLU}'(-x_i) & \text{if } i \bmod 2 \neq 0 \end{cases}$	$(-\infty, \infty)$	$C^0$	Yes	Yes	No
Leaky rectified linear unit (Leaky ReLU) [16]		$f(x) = \begin{cases} 0.01x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(x) = \begin{cases} 0.01 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$(-\infty, \infty)$	$C^0$	Yes	Yes	No
Parametric rectified linear unit (PReLU) [17]		$f(\alpha, x) = \begin{cases} \alpha x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(\alpha, x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$(-\infty, \infty)$ [20]	$C^0$	Yes if $\alpha \geq 0$	Yes	Yes if $\alpha = 1$
Randomized leaky rectified linear unit (RReLU) [18]		$f(\alpha, x) = \begin{cases} \alpha x & \text{for } x < 0 \text{ [3]} \\ x & \text{for } x \geq 0 \end{cases}$	$f'(\alpha, x) = \begin{cases} \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$(-\infty, \infty)$	$C^0$	Yes	Yes	No
Exponential linear unit (ELU) [19]		$f(\alpha, x) = \begin{cases} \alpha(e^x - 1) & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$	$f'(\alpha, x) = \begin{cases} f(\alpha, x) + \alpha & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$(-\alpha, \infty)$	$\begin{cases} C_1 & \text{when } \alpha = 1 \\ C_0 & \text{otherwise} \end{cases}$	Yes if $\alpha \geq 0$	Yes if $0 \leq \alpha \leq 1$	Yes if $\alpha = 1$
Scaled exponential linear unit (SELU) [20]		$f(\alpha, x) = \lambda \begin{cases} \alpha(e^x - 1) & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$ with $\lambda = 1.0507$ and $\alpha = 1.67326$	$f'(\alpha, x) = \lambda \begin{cases} \alpha(e^x) & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$	$(-\lambda\alpha, \infty)$	$C^0$	Yes	No	No
S-shaped rectified linear unit (SReLU) [21]		$f_{t_1, a_1, t_2, a_2}(x) = \begin{cases} t_1 + a_1(x - t_1) & \text{for } x \leq t_1 \\ x & \text{for } t_1 < x < t_2 \\ t_2 + a_2(x - t_2) & \text{for } x \geq t_2 \end{cases}$ $t_1, a_1, t_2, a_2$ are parameters.	$f'_{t_1, a_1, t_2, a_2}(x) = \begin{cases} a_1 & \text{for } x \leq t_1 \\ 1 & \text{for } t_1 < x < t_2 \\ a_2 & \text{for } x \geq t_2 \end{cases}$	$(-\infty, \infty)$	$C^0$	No	No	No
Adaptive piecewise linear (APL) [22]		$f(x) = \max(0, x) + \sum_{i=1}^k \alpha_i \max(0, -x + b_i)$	$f'(x) = H(x) - \sum_{i=1}^k \alpha_i H(-x + b_i)$ [24]	$(-\infty, \infty)$	$C^0$	No	No	No
SoftPlus [23]		$f(x) = \ln(1 + e^x)$	$f'(x) = \frac{1}{1 + e^{-x}}$	$(0, \infty)$	$C^\infty$	Yes	Yes	No
Bent identity		$f(x) = \frac{\sqrt{x^2 + 1} - 1}{2} + x$	$f'(x) = \frac{x}{2\sqrt{x^2 + 1}} + 1$	$(-\infty, \infty)$	$C^\infty$	Yes	Yes	Yes
Sigmoid Linear Unit (SLU) [24] (AKA SLU) [25] and Ssbsh-1 [26]		$f(x) = x \cdot \sigma(x)$ [25]	$f'(x) = f(x) + \sigma(x)(1 - f(x))$ [25]	$[-\infty, \infty)$	$C^\infty$	No	No	Approximates identity/2
SoftExponential [27]		$f(\alpha, x) = \begin{cases} \frac{-\ln(-\alpha x + 1)}{\alpha} & \text{for } \alpha < 0 \\ \frac{x}{\alpha} & \text{for } \alpha = 0 \\ \frac{e^{\alpha x} - 1}{\alpha} + \alpha & \text{for } \alpha > 0 \end{cases}$	$f'(\alpha, x) = \begin{cases} \frac{1}{1 - \alpha x} & \text{for } \alpha < 0 \\ \frac{1}{e^{\alpha x}} & \text{for } \alpha = 0 \\ e^{\alpha x} & \text{for } \alpha > 0 \end{cases}$	$(-\infty, \infty)$	$C^\infty$	Yes	Yes	Yes if $\alpha = 0$
Soft Clipping [28]		$f(\alpha, x) = \frac{1}{\alpha} \log \frac{1 + e^{\alpha x}}{1 + e^{\alpha(x-1)}}$	$f'(\alpha, x) = \frac{1}{2} \sinh\left(\frac{\alpha x}{2}\right) \operatorname{sech}\left(\frac{\alpha x}{2}\right) \operatorname{sech}\left(\frac{\alpha}{2}(1 - x)\right)$	$(0, 1)$	$C^\infty$	Yes	No	No
Sinuid [29]		$f(x) = \sin(x)$	$f'(x) = \cos(x)$	$[-1, 1]$	$C^\infty$	No	No	Yes
Sinc		$f(x) = \begin{cases} 1 & \text{for } x = 0 \\ \frac{\sin(x)}{x} & \text{for } x \neq 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x = 0 \\ \frac{\cos(x) - \sin(x)}{x^2} & \text{for } x \neq 0 \end{cases}$	$[-\infty, 1]$	$C^\infty$	No	No	No
Gaussian		$f(x) = e^{-x^2}$	$f'(x) = -2xe^{-x^2}$	$(0, 1]$	$C^\infty$	No	No	No

# Which activation?

In practice:

- prefer ReLU. Use slightly positive initial bias to avoid dead Relu issue.
- Try out Leaky ReLU/PRelu
- No sigmoid!



UNIVERSITÀ  
DEGLI STUDI  
DI PADOVA

# Machine Learning 2024/2025

**AMCO**  
ARTIFICIAL INTELLIGENCE, MACHINE  
LEARNING AND CONTROL RESEARCH GROUP

# Thank you!

# Gian Antonio Susto

