

Intro to Deep Learning

IAP Computing, Optimization, and Statistics

Jan 13, 2022

Outline

1. Fundamental concepts on Neural Networks

1. Intro/Motivation
2. Perceptron
3. Neural Network Architecture
4. SGD and Back-Propagation

2. Tensorflow/Keras Tutorial

3. Convolutional Neural Networks

1. Convolutional layer
2. Pooling Layer

4. Image Classification using CNN

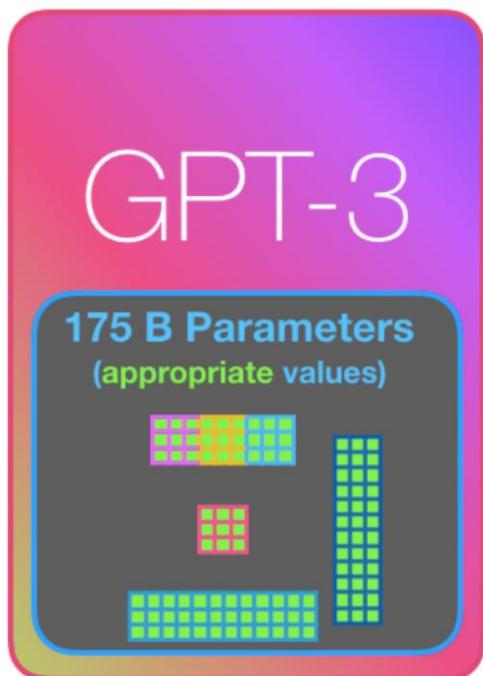
5. If time permits

1. Transfer Learning
2. Hyperparameter tuning and model selection
3. Data Augmentation



Let's start with some cool AI stuff

Open AI's GPT-3



Geoffrey Hinton
@geoffreyhinton

Extrapolating the spectacular performance of GPT3 into the future suggests that the answer to life, the universe and everything is just 4.398 trillion parameters.

2:26 PM · Jun 10, 2020 · Twitter Web App

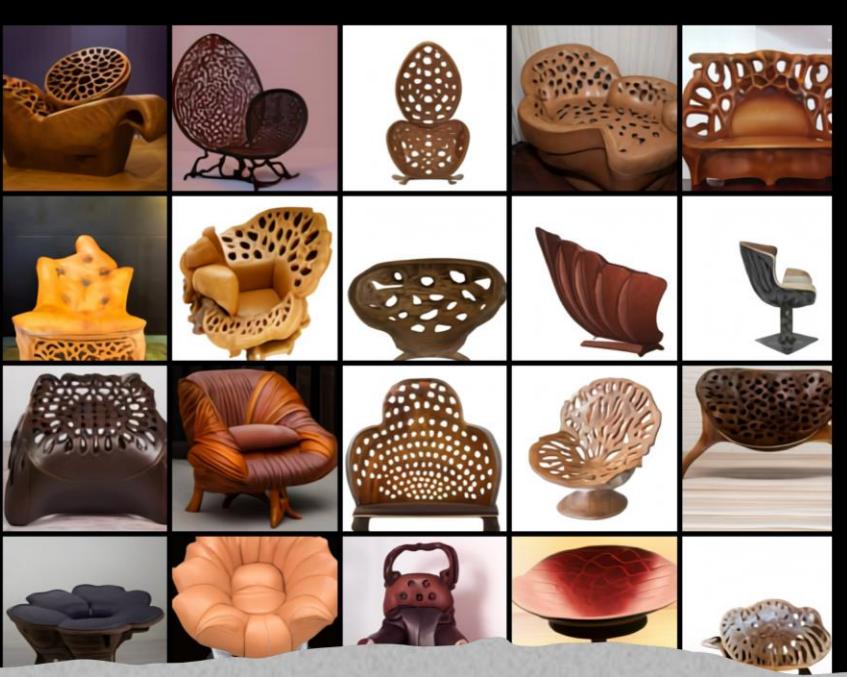
741 Retweets and comments 3.8K Likes



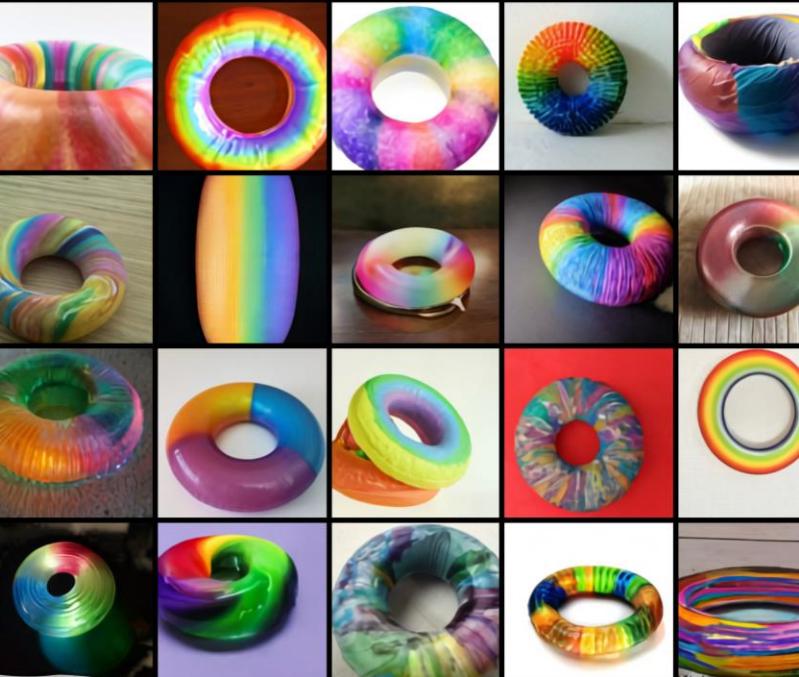
a peacock made of lettuce. a peacock with the texture of a lettuce.



a leather armchair in the style of a lotus root. a leather armchair imitating a lotus root.



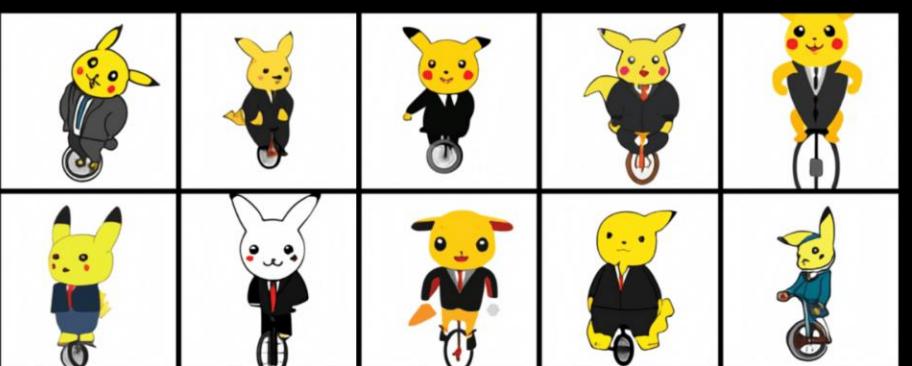
a torus made of rainbow. a torus with the texture of rainbow.



DALL-E

an illustration of a pikachu in a suit riding a unicycle

AI-GENERATED IMAGES



G google/pegasus-xsum

like 13

[Summarization](#)[PyTorch](#)[TensorFlow](#)[JAX](#)[Transformers](#)[en](#)[arxiv:1912.08777](#)[pegasus](#)[text2text-generation](#)[AutoNLP Compatible](#)

⚡ Hosted inference API ⓘ

[Summarization](#)

Examples

The "big data revolution" has placed added emphasis on computational techniques for decision-making with data. Large-scale optimization and machine learning are now commonplace among researchers and practitioners alike. More than ever, there is not only a need to develop techniques, but also to implement and use them in computational practice. 15.S60 is a multi-session workshop on software tools for informing decision-making using data, with a focus on contemporary methods in optimization and statistics. We concentrate on teaching elementary principles of computational practice using common software and practical methods. By the end of the course, students will possess a baseline technical knowledge for modern research practice. The course is divided into 8 self-contained modules. Each module consists of a 3-hour workshop where participants learn a specific software tool. The goal of each module is to enable participants to use the software and techniques covered in their own research.

[Compute](#)

Computation time on cpu: 4.8708 s

15.S60 is a multi-session workshop on software tools for informing decision-making using data, with a focus on contemporary methods in optimization and statistics.

[JSON Output](#)[Maximize](#)

multi_news

40.74/11.95/24.20

41.52/18.12/24.91

41.05/18.15/24.25

[Train](#)[Deploy](#)[Use in Transformers](#)

NEW Select AutoNLP in the "Train" menu to fine-tune this model automatically.

Downloads last month

254,220



⚡ Hosted inference API ⓘ

[Summarization](#)

Examples

The "big data revolution" has placed added emphasis on computational techniques for decision-making with data. Large-scale optimization and machine learning are now commonplace among researchers and practitioners alike. More than ever, there is not only a need to develop techniques, but also to implement and use them in computational practice. 15.S60 is a multi-session workshop on software tools for informing decision-making using data, with a focus on contemporary methods in optimization and statistics. We concentrate on teaching elementary principles of computational practice using common software and practical methods. By the end of the course, students will possess a baseline technical knowledge for modern research practice.

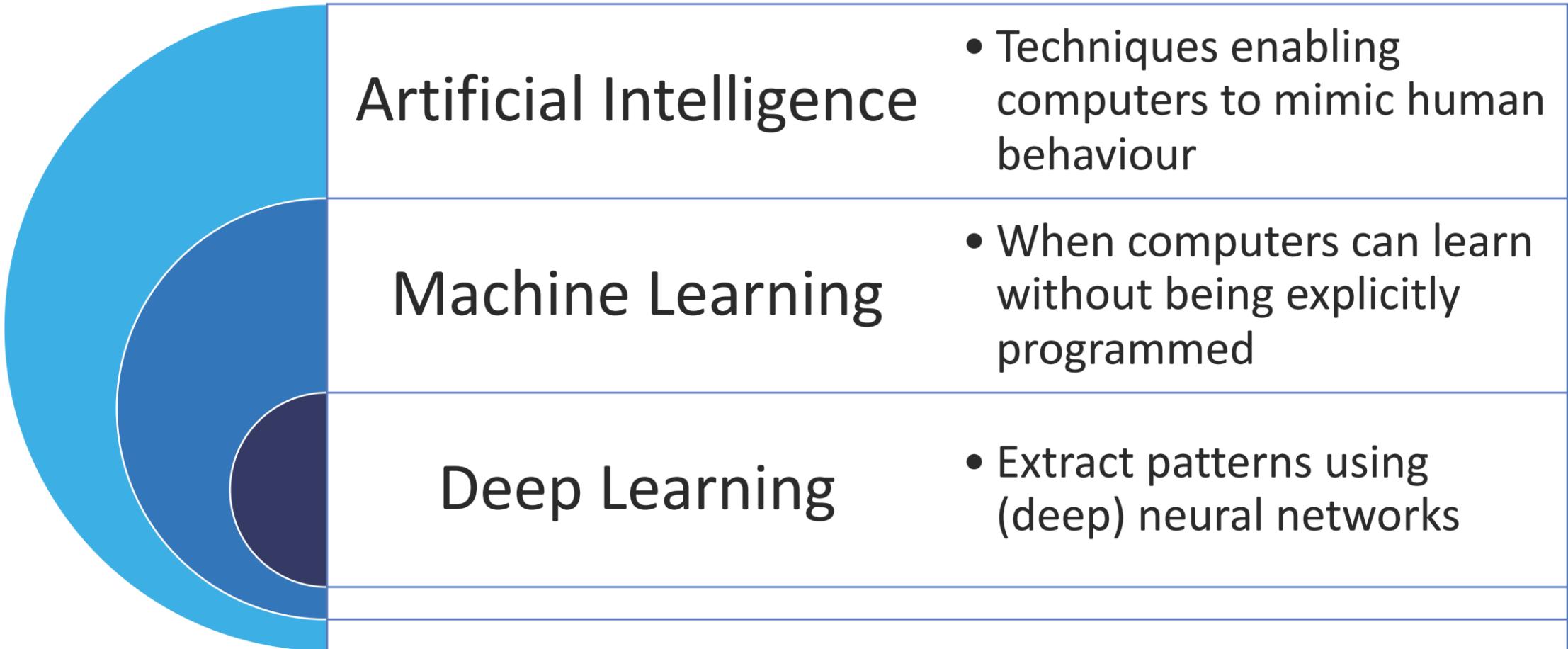
[Compute](#)

Computation time on cpu: 4.6536 s

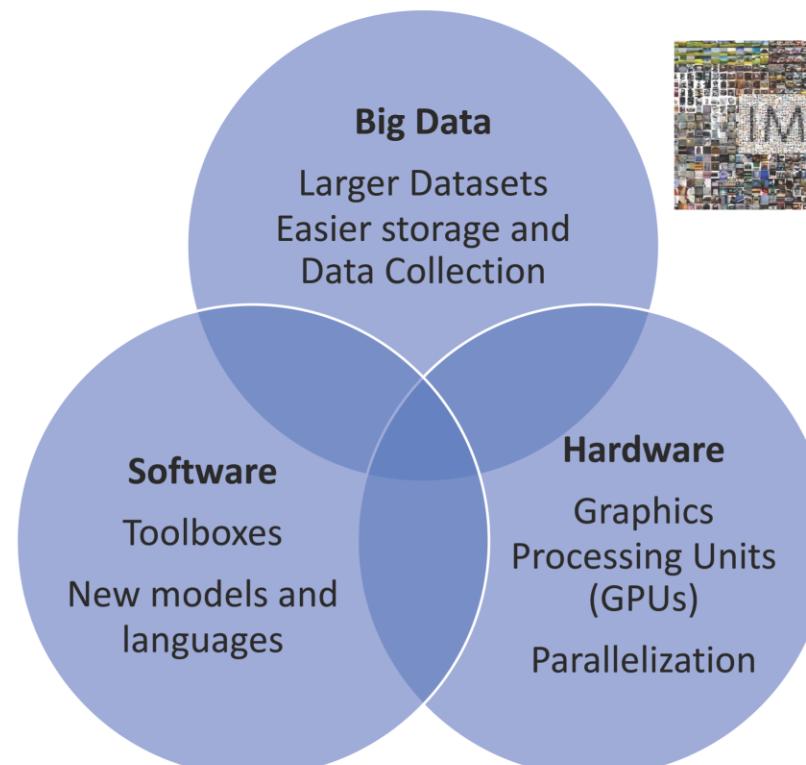
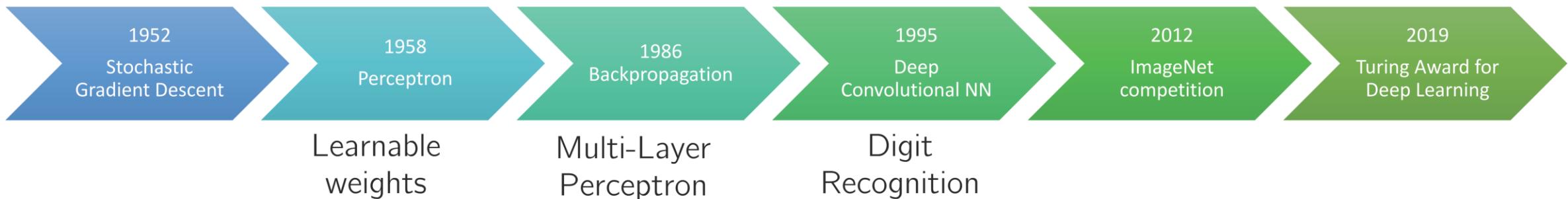
This course aims to provide students with a basic understanding of software tools for informing decision-making using data, with a focus on contemporary methods in optimization and statistics.

[JSON Output](#)[Maximize](#)

What is Deep Learning?

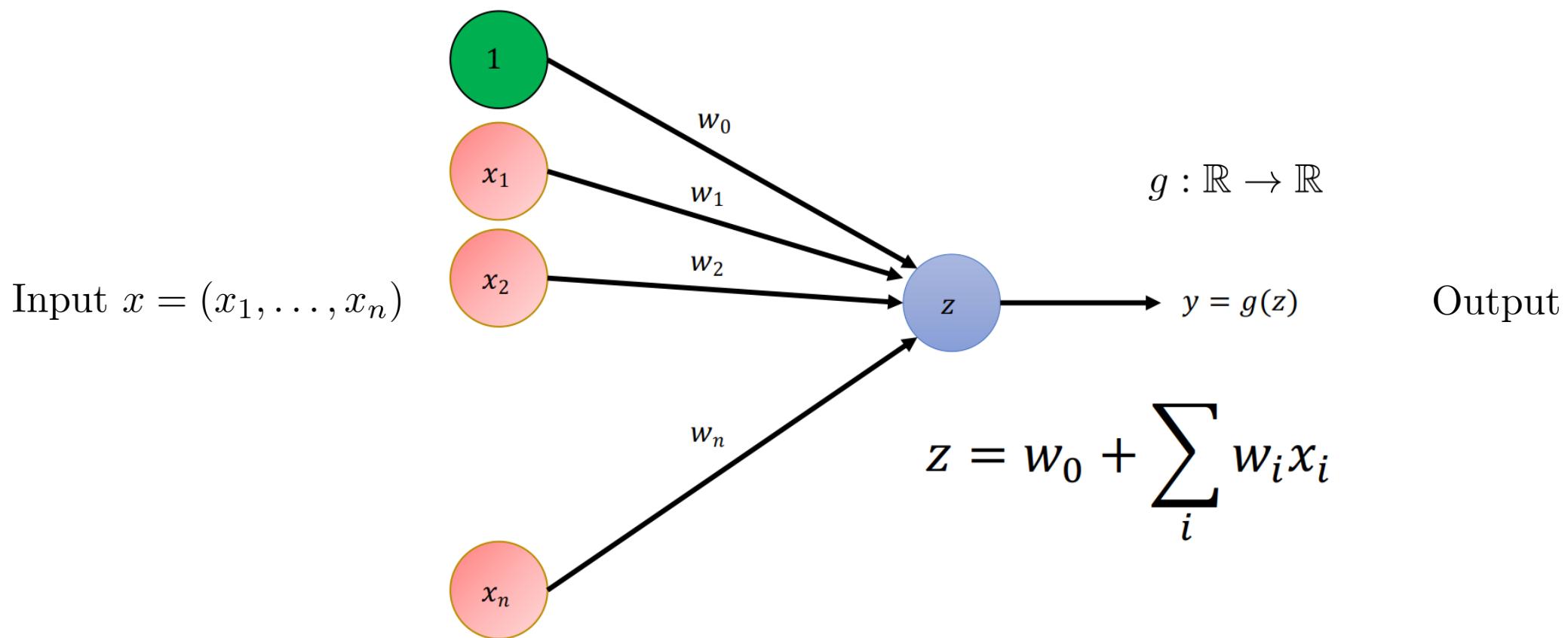


What is the history behind the current AI boom?

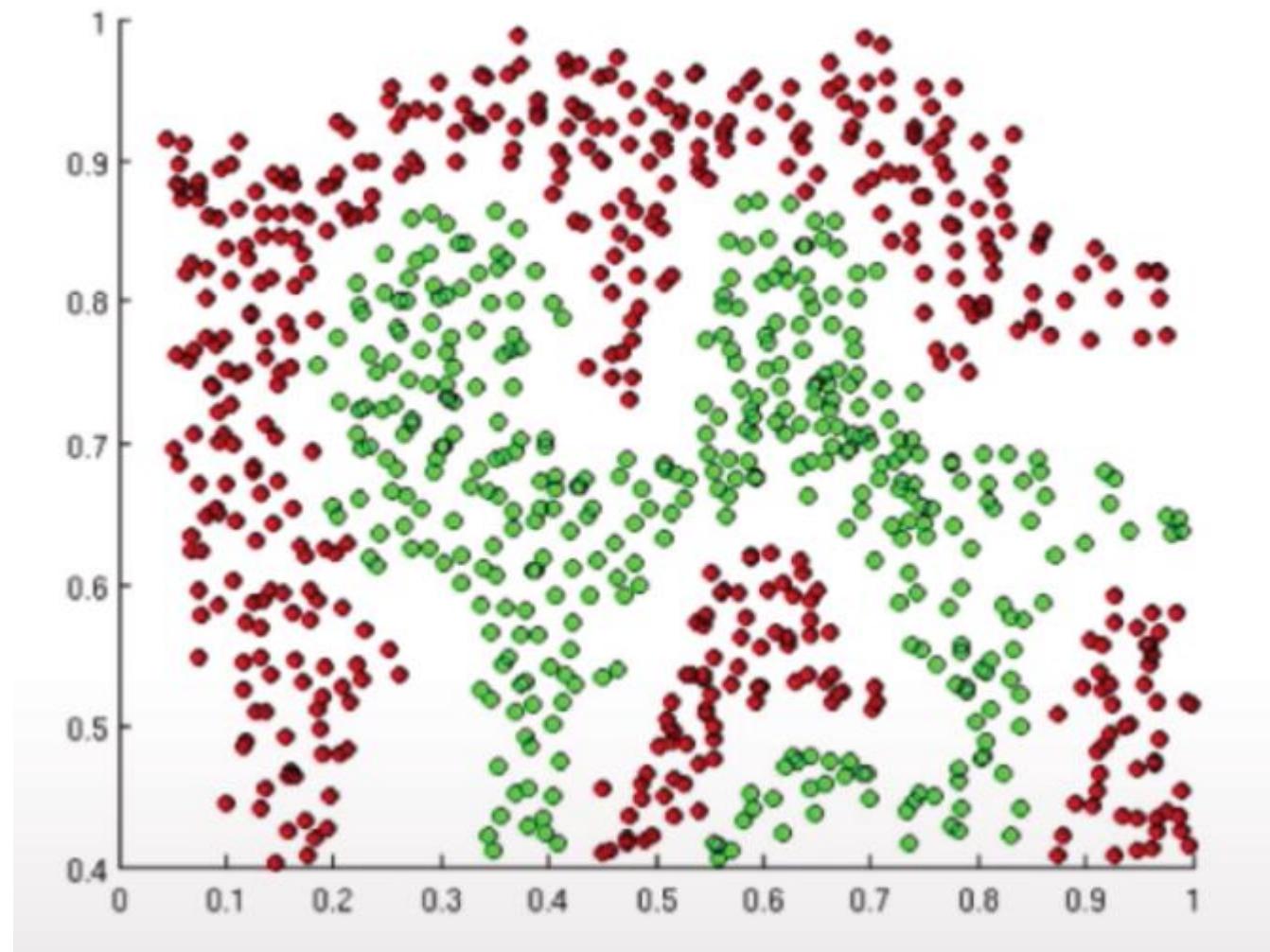


Neural Networks Basics

The Perceptron



Classification task

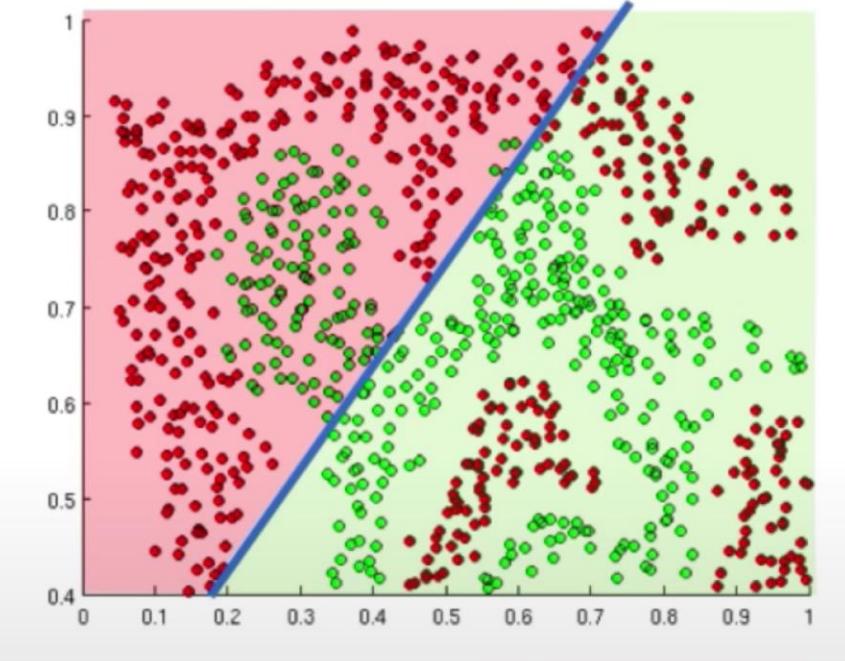


Activation functions

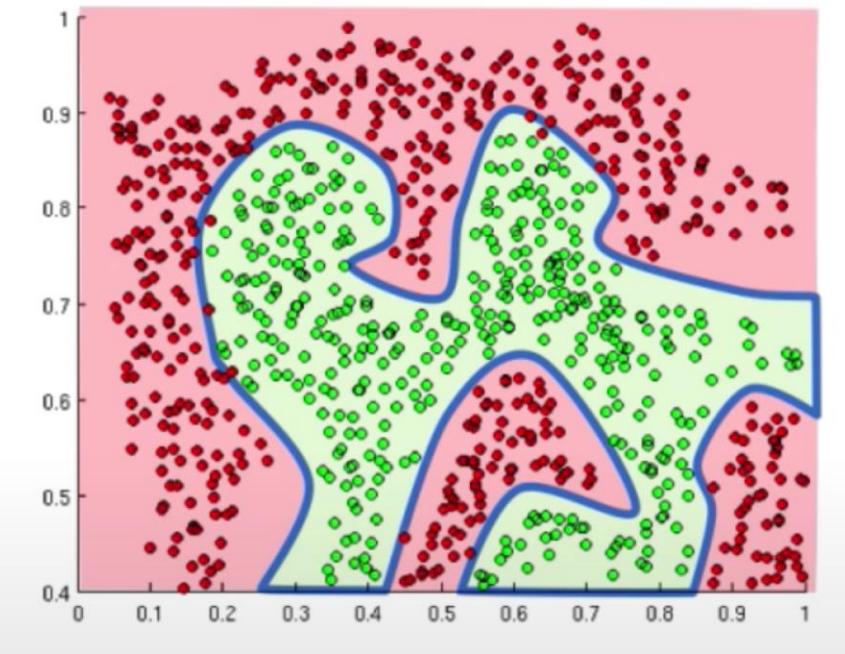
$$z = w_0 + \sum_i w_i x_i$$

$$y = g(z)$$

Linear activation functions produce linear decisions no matter the network size



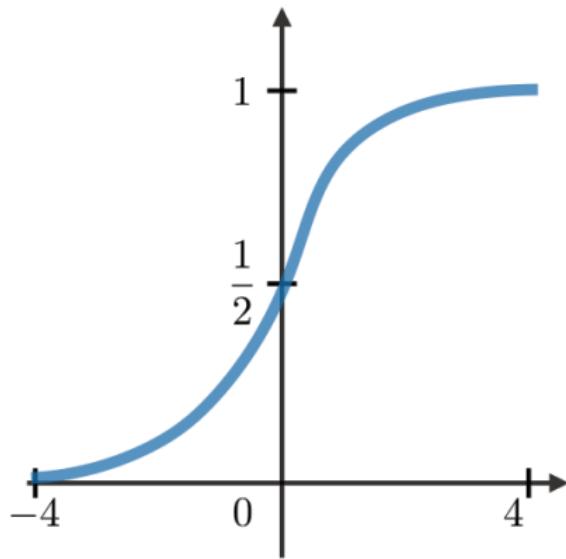
Non-linear activation can help approximate arbitrarily complex functions



Activation Functions

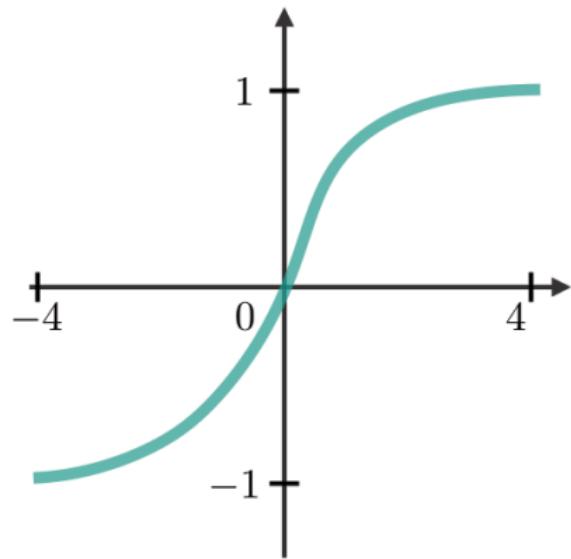
Sigmoid

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



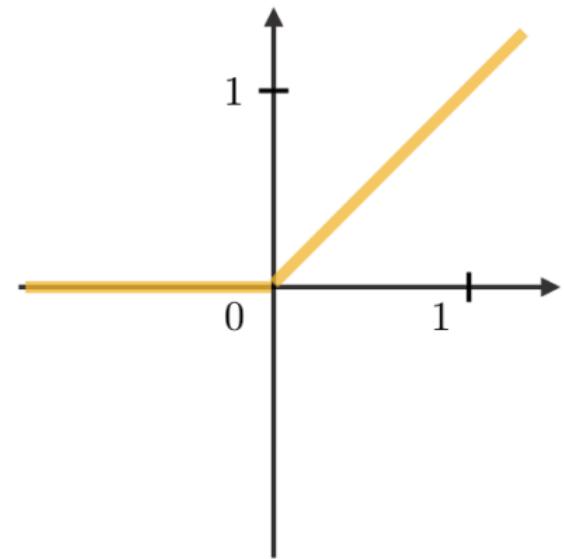
Tanh

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

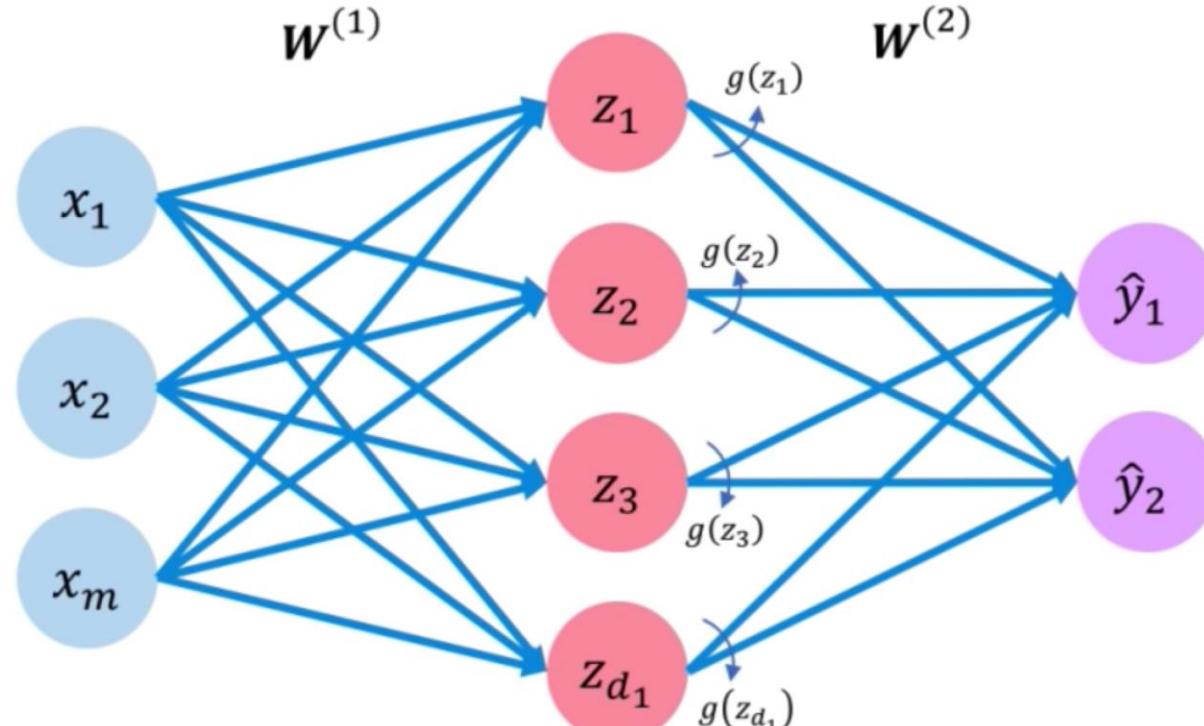


RELU

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



One-layer Neural Network



Inputs

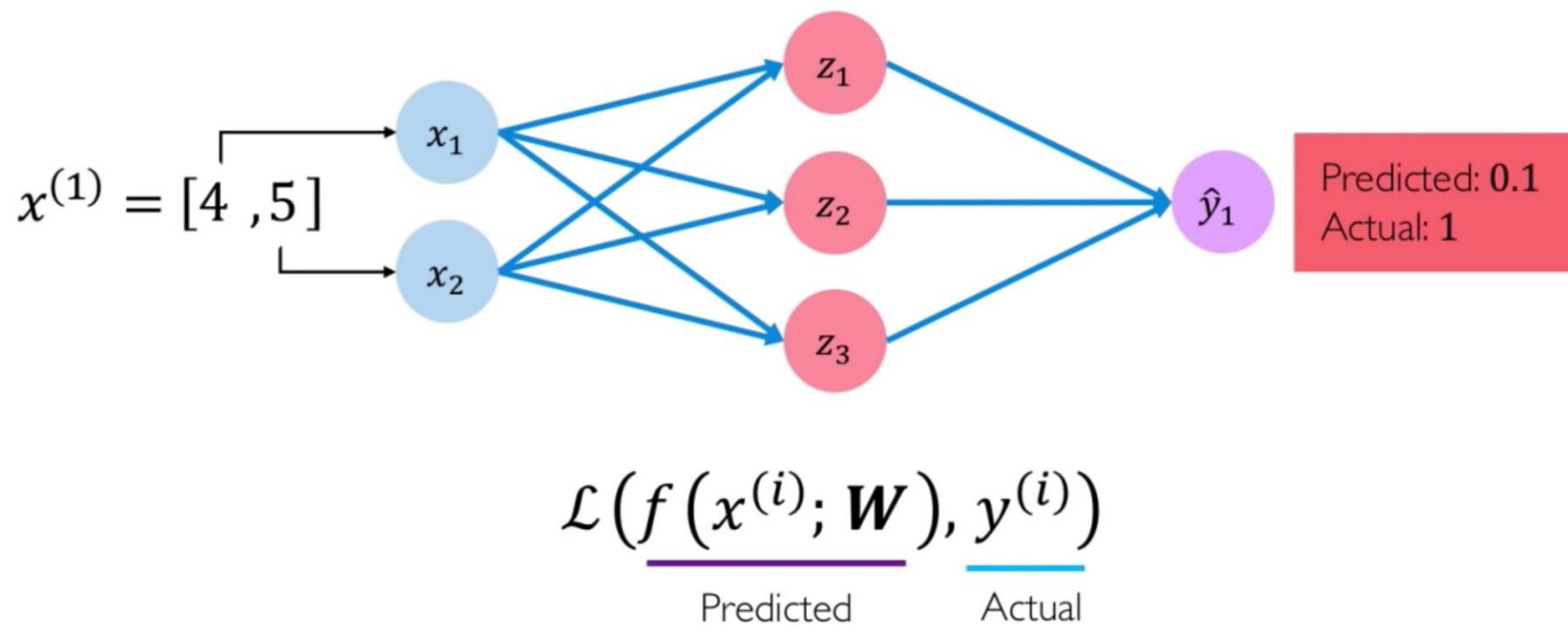
Hidden

Final Output

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

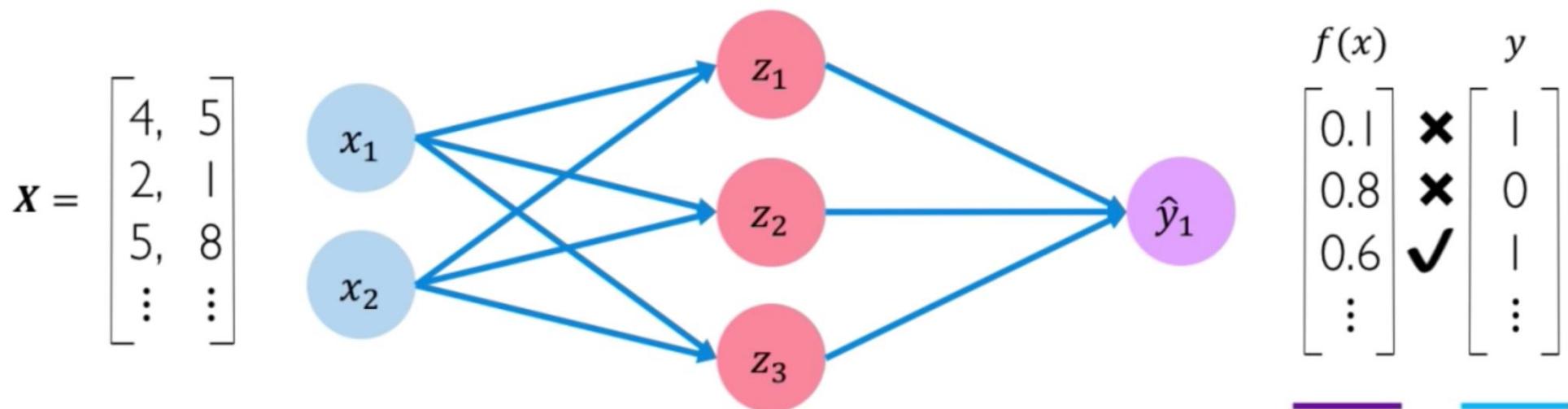
Loss function

The **loss** of our network measures the cost incurred from incorrect predictions.



Empirical Loss

The **empirical loss** measures the total loss over our entire dataset.



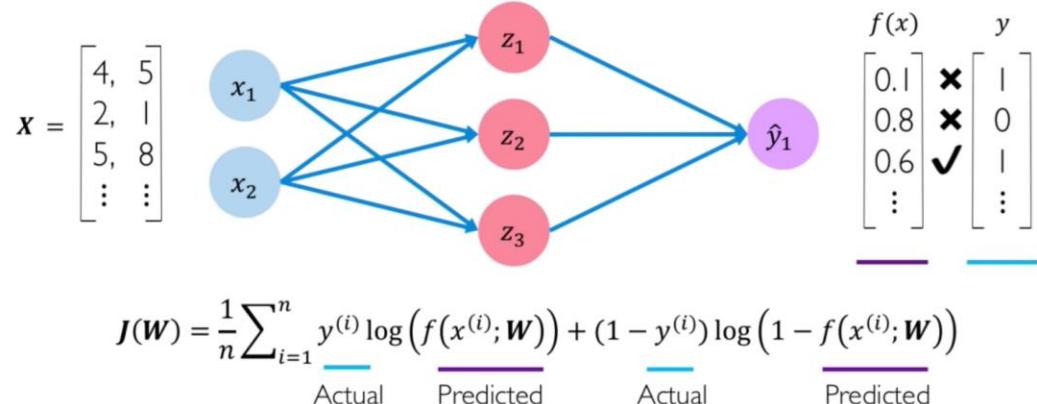
Also known as:

- Objective function
- Cost function
- Empirical Risk

Different loss functions

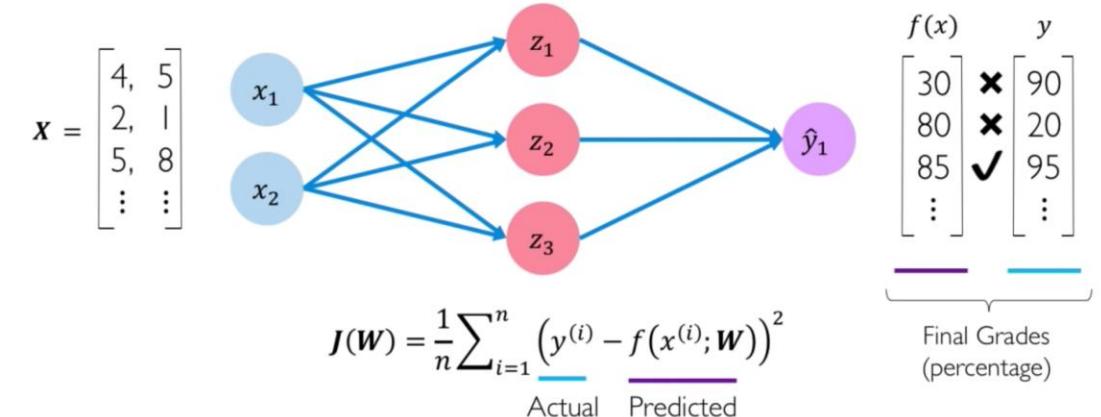
Binary Cross Entropy Loss

Cross entropy loss can be used with models that output a probability between 0 and 1



Mean Squared Error Loss

Mean squared error loss can be used with regression models that output continuous real numbers



```
TensorFlow loss = tf.reduce_mean(tf.nn.softmax_cross_entropy_with_logits(y, predicted))
```

```
TensorFlow loss = tf.reduce_mean(tf.square(tf.subtract(y, predicted)))
```

How to optimize the loss function?

We need to find the network weights leading to the lowest loss.

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

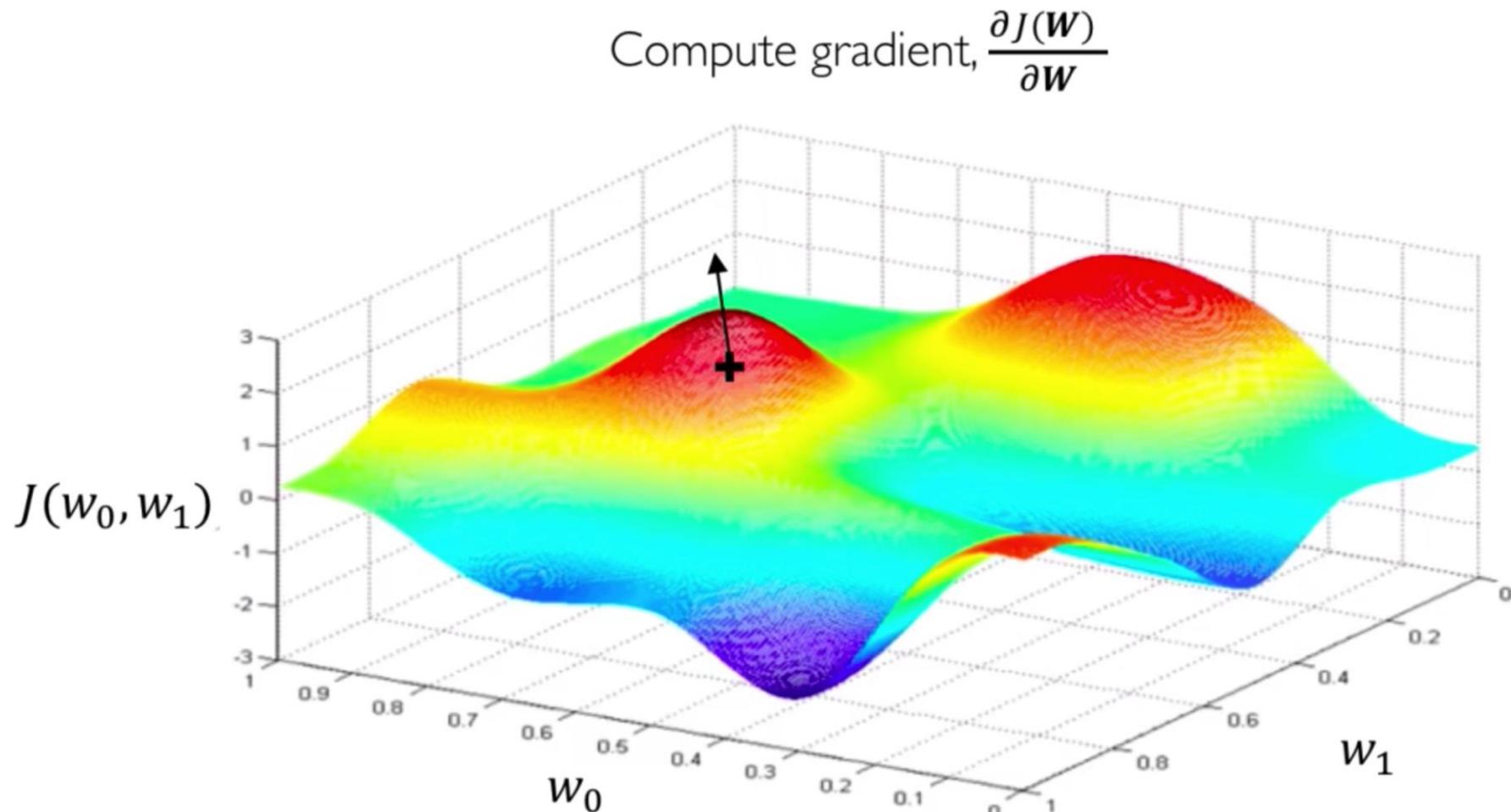
$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} J(\mathbf{W})$$



Remember:

$$\mathbf{W} = \{\mathbf{W}^{(0)}, \mathbf{W}^{(1)}, \dots\}$$

But it looks like this...



Fortunately there is Gradient Descent!

Algorithm

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, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights

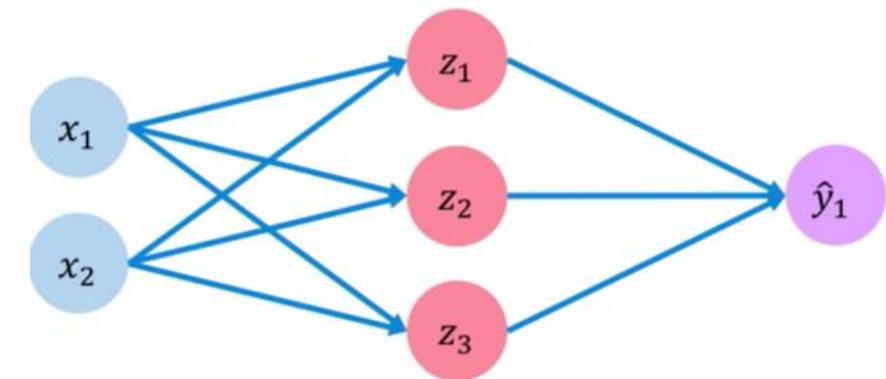
Next problem: compute gradients

We need to compute the gradient $\frac{\partial \mathcal{J}(W)}{\partial W}$ with respect to all parameters $W = (w_{ij}^{(l)})_{\{l; i,j\}}$.

Let's begin with a neural network with one layer:

$$\frac{\partial \mathcal{J}(w_{ji}^{(2)})}{\partial W} = \frac{1}{n} \frac{\partial \mathcal{L}(\hat{y}, y_i)}{\partial \hat{y}} \cdot \frac{\partial g}{\partial w_{ji}^{(2)}}$$

Chain rule!



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

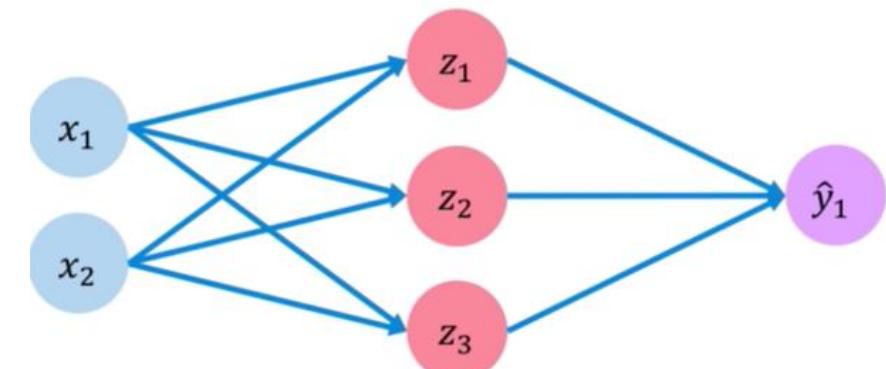
$$J(W) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \underline{W}), \underline{\hat{y}})$$

Next problem: compute gradients

We need to compute the gradient $\frac{\partial \mathcal{J}(W)}{\partial W}$ with respect to all parameters $W = (w_{ij}^{(l)})_{\{l; i,j\}}$.

Let's begin with a neural network with one layer:

$$\begin{aligned}\frac{\partial \mathcal{J}(w_{ji}^{(2)})}{\partial W} &= \frac{1}{n} \frac{\partial \mathcal{L}(\hat{y}, y_i)}{\partial \hat{y}} \cdot \frac{\partial g}{\partial w_{ji}^{(2)}} \\ &= \frac{1}{n} \frac{\partial \mathcal{L}(\hat{y})}{\partial \hat{y}} \cdot g'(w^{(2)\top} z) z_{ji}\end{aligned}$$



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

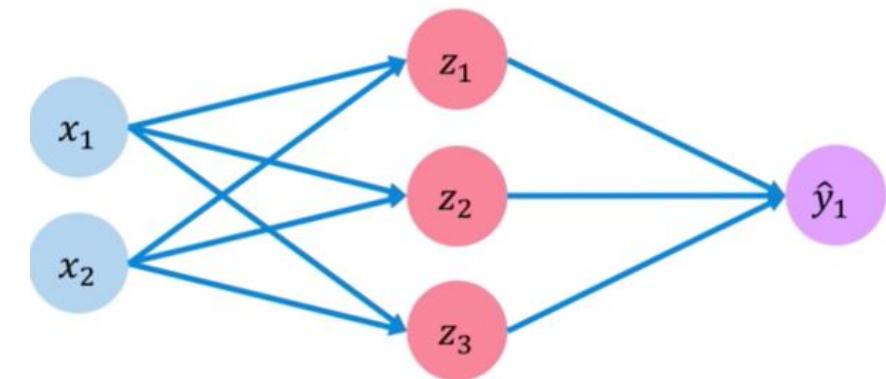
$$\begin{aligned}J(W) &= \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \underline{W}), \underline{\hat{y}}) \\ &= \underline{\hat{y}}\end{aligned}$$

Next problem: compute gradients

We need to compute the gradient $\frac{\partial \mathcal{J}(W)}{\partial W}$ with respect to all parameters $W = (w_{ij}^{(l)})_{\{l; i,j\}}$.

Let's begin with a neural network with one layer:

$$\frac{\partial \mathcal{J}(w_{ji}^{(1)})}{\partial W} = \frac{1}{n} \frac{\partial \mathcal{L}(\hat{y}, y_i)}{\partial \hat{y}} \cdot \frac{\partial g}{\partial w_{ji}^{(1)}}$$



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

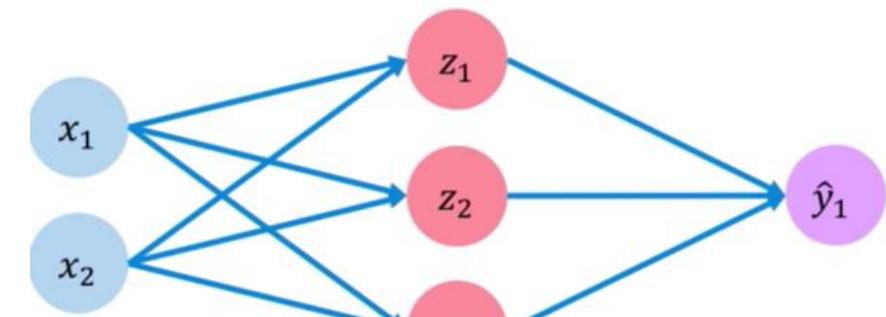
$$J(W) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \underline{W}), \underline{\hat{y}})$$

Next problem: compute gradients

We need to compute the gradient $\frac{\partial \mathcal{J}(W)}{\partial W}$ with respect to all parameters $W = (w_{ij}^{(l)})_{\{l; i,j\}}$.

Let's begin with a neural network with one layer:

$$\begin{aligned}\frac{\partial \mathcal{J}(w_{ji}^{(1)})}{\partial W} &= \frac{1}{n} \frac{\partial \mathcal{L}(\hat{y}, y_i)}{\partial \hat{y}} \cdot \frac{\partial g}{\partial w_{ji}^{(1)}} \\ &= \frac{1}{n} \frac{\partial \mathcal{L}(\hat{y})}{\partial \hat{y}} \cdot g'(w^{(2)\top} z) \frac{\partial (w^{(2)\top} z)}{\partial w_{ji}^{(1)}}\end{aligned}$$



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

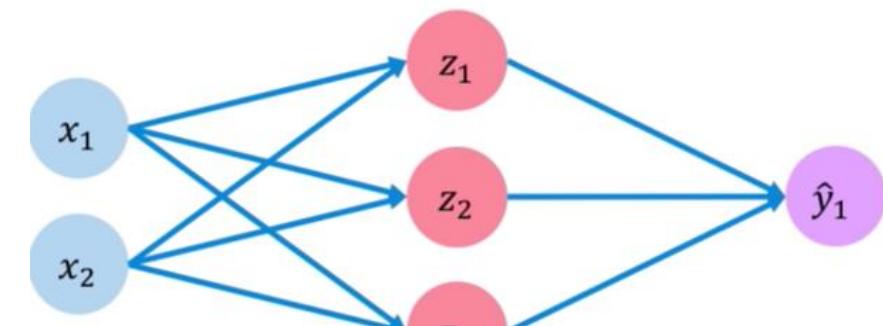
$$\begin{aligned}J(W) &= \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \underline{W}), \underline{\hat{y}}) \\ &= \underline{\hat{y}}\end{aligned}$$

Next problem: compute gradients

We need to compute the gradient $\frac{\partial \mathcal{J}(W)}{\partial W}$ with respect to all parameters $W = (w_{ij}^{(l)})_{\{l; i,j\}}$.

Let's begin with a neural network with one layer:

$$\begin{aligned}\frac{\partial \mathcal{J}(w_{ji}^{(1)})}{\partial W} &= \frac{1}{n} \frac{\partial \mathcal{L}(\hat{y}, y_i)}{\partial \hat{y}} \cdot \frac{\partial g}{\partial w_{ji}^{(1)}} \\ &= \frac{1}{n} \frac{\partial \mathcal{L}(\hat{y})}{\partial \hat{y}} \cdot g'(w^{(2)\top} z) \frac{\partial (w^{(2)\top} z)}{\partial w_{ji}^{(1)}} \\ &= \frac{1}{n} \frac{\partial \mathcal{L}(\hat{y})}{\partial \hat{y}} \cdot g'(w^{(2)\top} z) w_{ji}^{(2)} \frac{\partial z_j}{\partial w_{ji}^{(1)}}\end{aligned}$$



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

$$\begin{aligned}J(W) &= \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \underline{W}), \underline{y}^{(i)}) \\ &= \hat{y}\end{aligned}$$

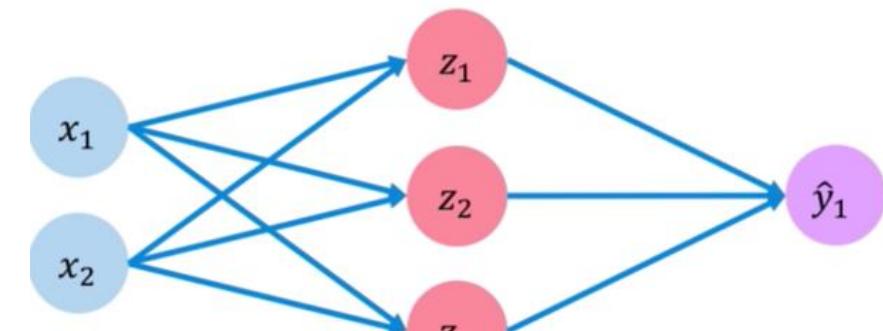
Next problem: compute gradients

We need to compute the gradient $\frac{\partial \mathcal{J}(W)}{\partial W}$ with respect to all parameters $W = (w_{ij}^{(l)})_{\{l; i,j\}}$.

Let's begin with a neural network with one layer:

$$\begin{aligned}\frac{\partial \mathcal{J}(w_{ji}^{(1)})}{\partial W} &= \frac{1}{n} \frac{\partial \mathcal{L}(\hat{y}, y_i)}{\partial \hat{y}} \cdot \frac{\partial g}{\partial w_{ji}^{(1)}} \\ &= \frac{1}{n} \frac{\partial \mathcal{L}(\hat{y})}{\partial \hat{y}} \cdot g'(w^{(2)\top} z) \frac{\partial (w^{(2)\top} z)}{\partial w_{ji}^{(1)}} \\ &= \frac{1}{n} \frac{\partial \mathcal{L}(\hat{y})}{\partial \hat{y}} \cdot g'(w^{(2)\top} z) w_{ji}^{(2)} \frac{\partial z_j}{\partial w_{ji}^{(1)}}\end{aligned}$$

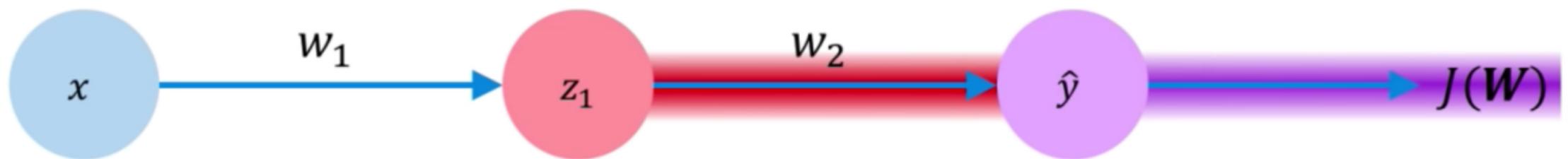
We already computed this!



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

$$J(W) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \underline{W}), \underline{\hat{y}})$$

Backpropagation (aka chain rule)



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

But how do you choose the learning rate?

Remember:

Optimization through gradient descent

$$\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}}$$

↑
How can we set the
learning rate?

Ideas?

Adaptative Learning Rate

Learning rate is not fixed anymore, instead:

You can adapt it depending on:

- How large gradients are
- How fast the learning is happening (momentum)
- ...

In Deep Learning, Adam is an extremely popular method and does all this automatically. But you still need to choose a base value (typically between $10^{-3}/5.10^{-5}$ but depends on the task and architecture.)

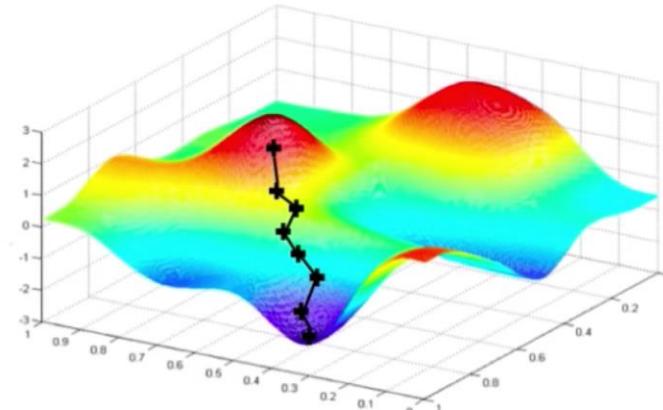
Another problem: computing gradients is expensive!

Gradient Descent

Algorithm

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, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights

Can be very
computationally
intensive to compute!

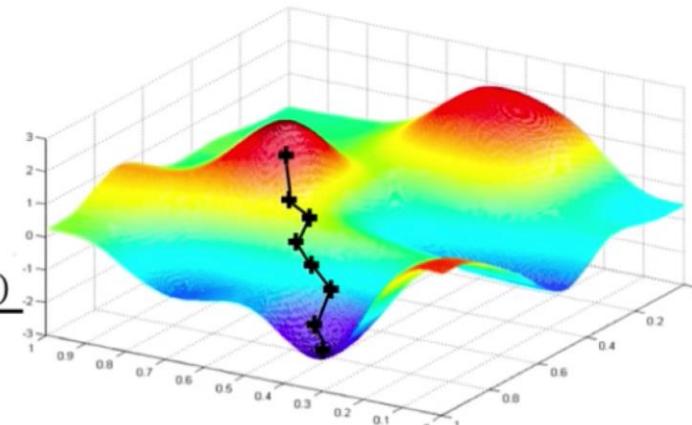


Another problem: computing gradients is expensive!

Stochastic Gradient Descent

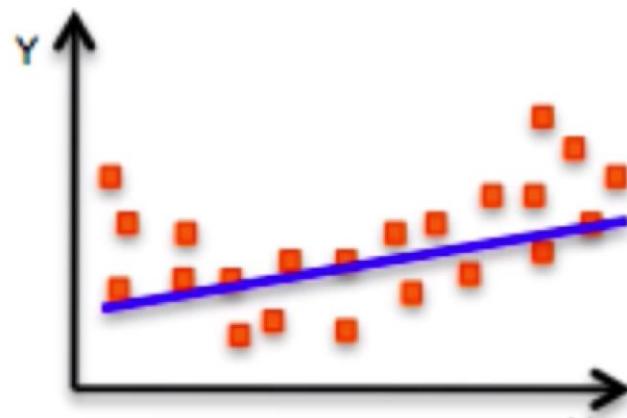
Algorithm

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

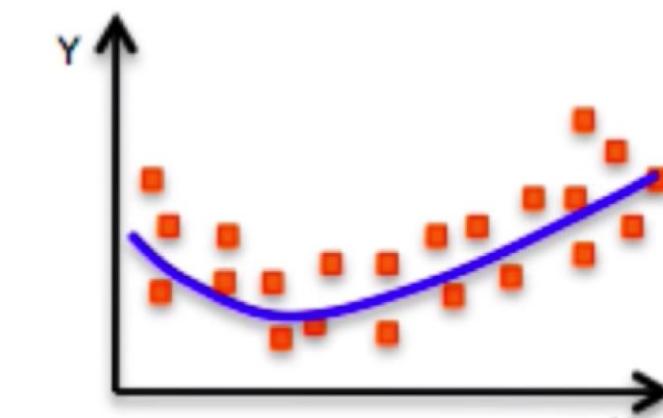


We use **mini-batches** while training allowing for a smoother convergence and larger learning rates.

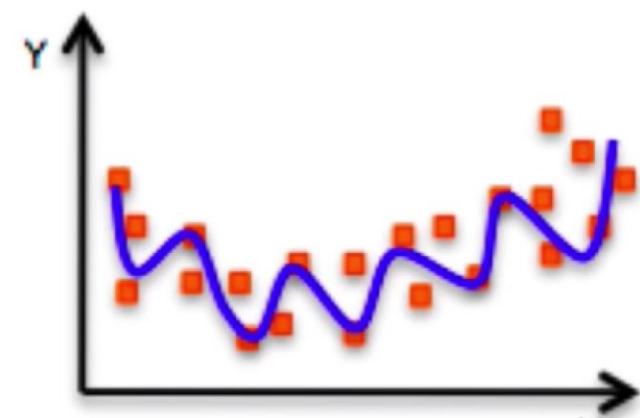
Another classic ML problem: overfitting



Underfitting
Model does not have capacity
to fully learn the data



← **Ideal fit** →



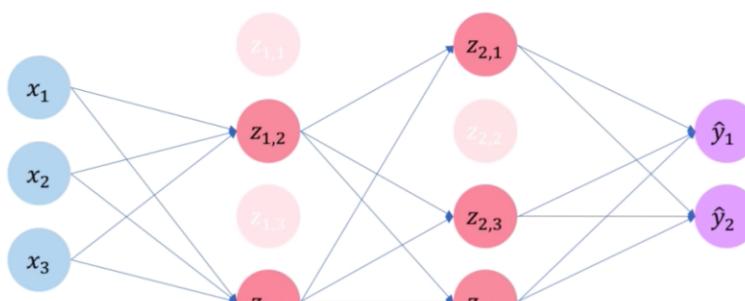
Overfitting
Too complex, extra parameters,
does not generalize well

Let's overcome this! Ideas?

Regularization I: Dropout

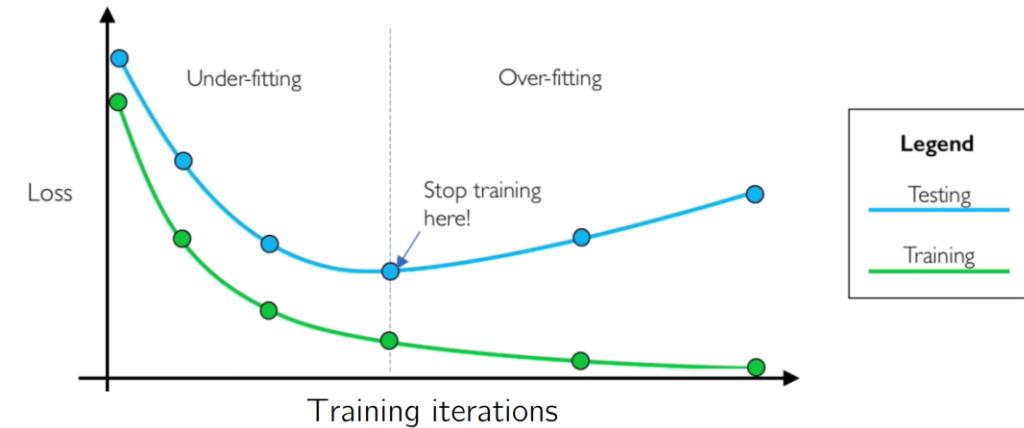
- During training, randomly set some activations to 0
 - Typically 'drop' 50% of activations in layer
 - Forces network to not rely on any 1 node

 `tf.keras.layers.Dropout(p=0.5)`



Regularization 2: Early Stopping

- Stop training before we have a chance to overfit

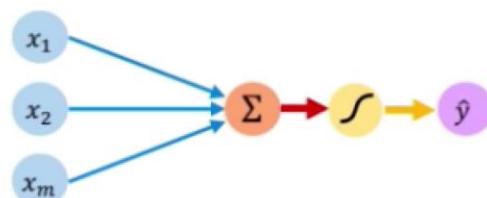


Regularization 3: L1 or L2 penalty on the weights

Review

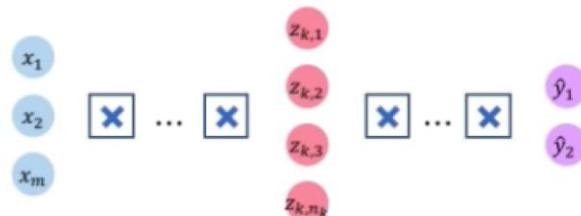
The Perceptron

- Structural building blocks
- Nonlinear activation functions



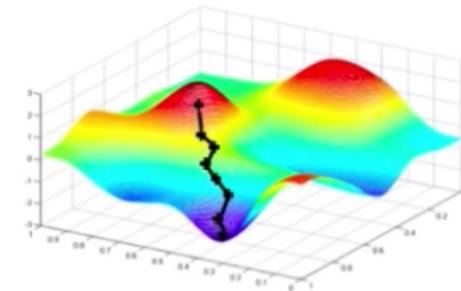
Neural Networks

- Stacking Perceptrons to form neural networks
- Optimization through backpropagation



Training in Practice

- Adaptive learning
- Batching
- Regularization



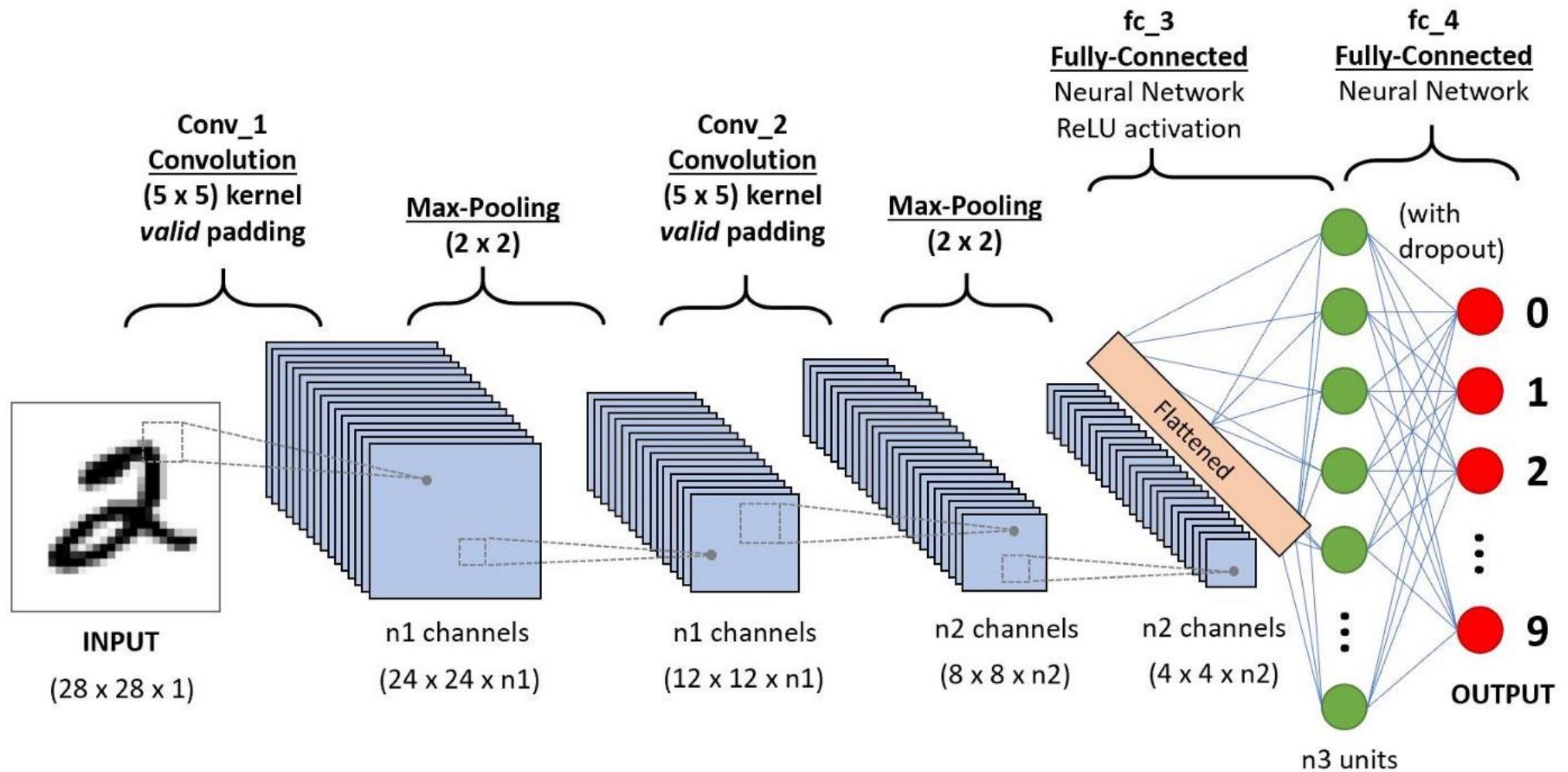
Let's get hands on

<https://colab.research.google.com/drive/1abnkCeDxtwQPqauK7Xjxrl6ewZesPjzv?usp=sharing>

Break

Convolutional Neural Networks

CNN are well suited for images as inputs



Convolutional Layer

7	2	3	3	8
4	5	3	8	4
3	3	2	8	4
2	8	7	2	7
5	4	4	5	4

*

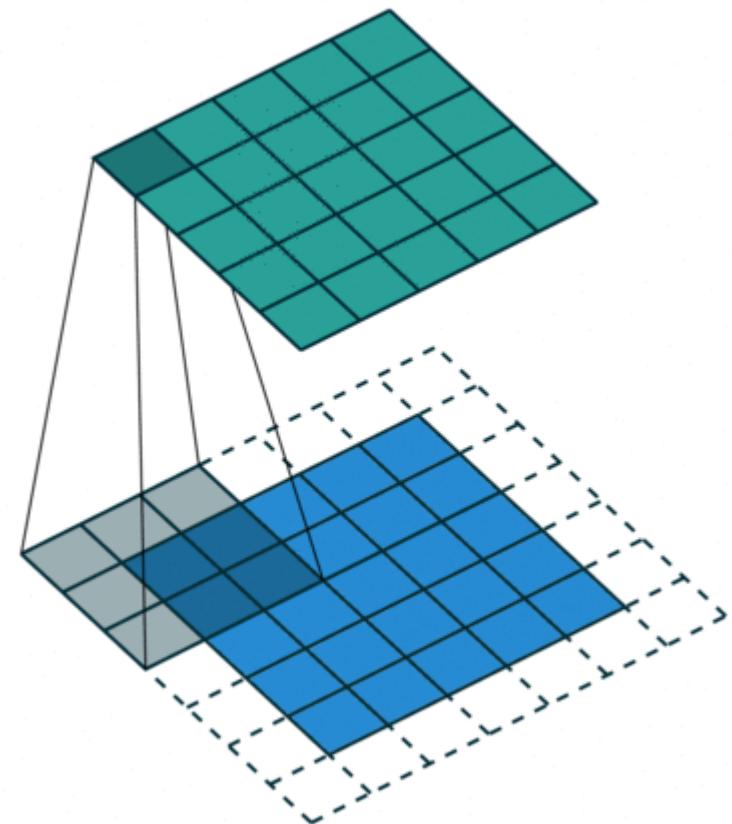
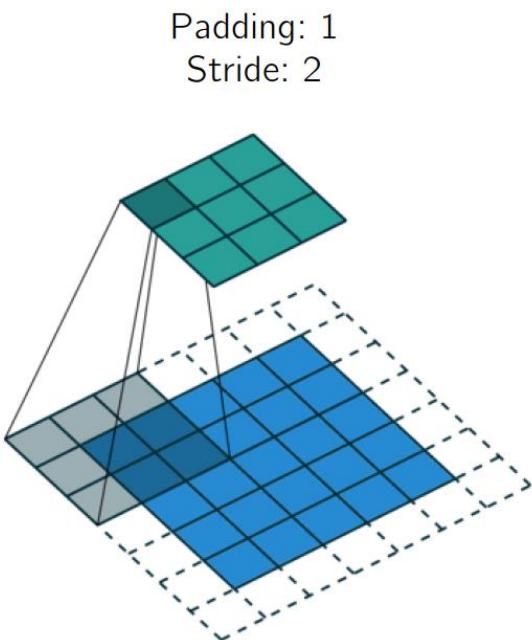
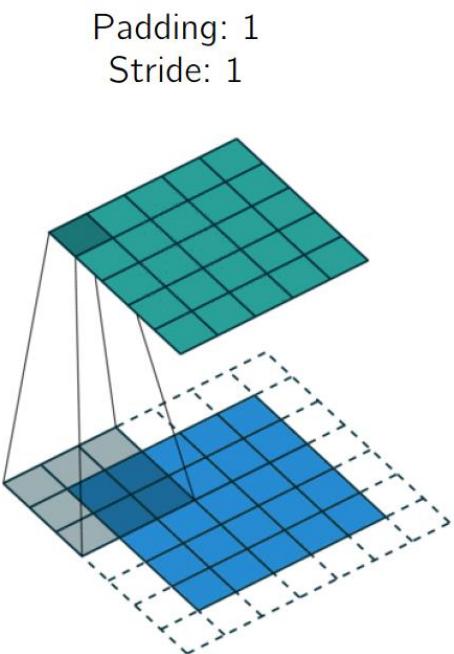
1	0	-1
1	0	-1
1	0	-1

=

6		

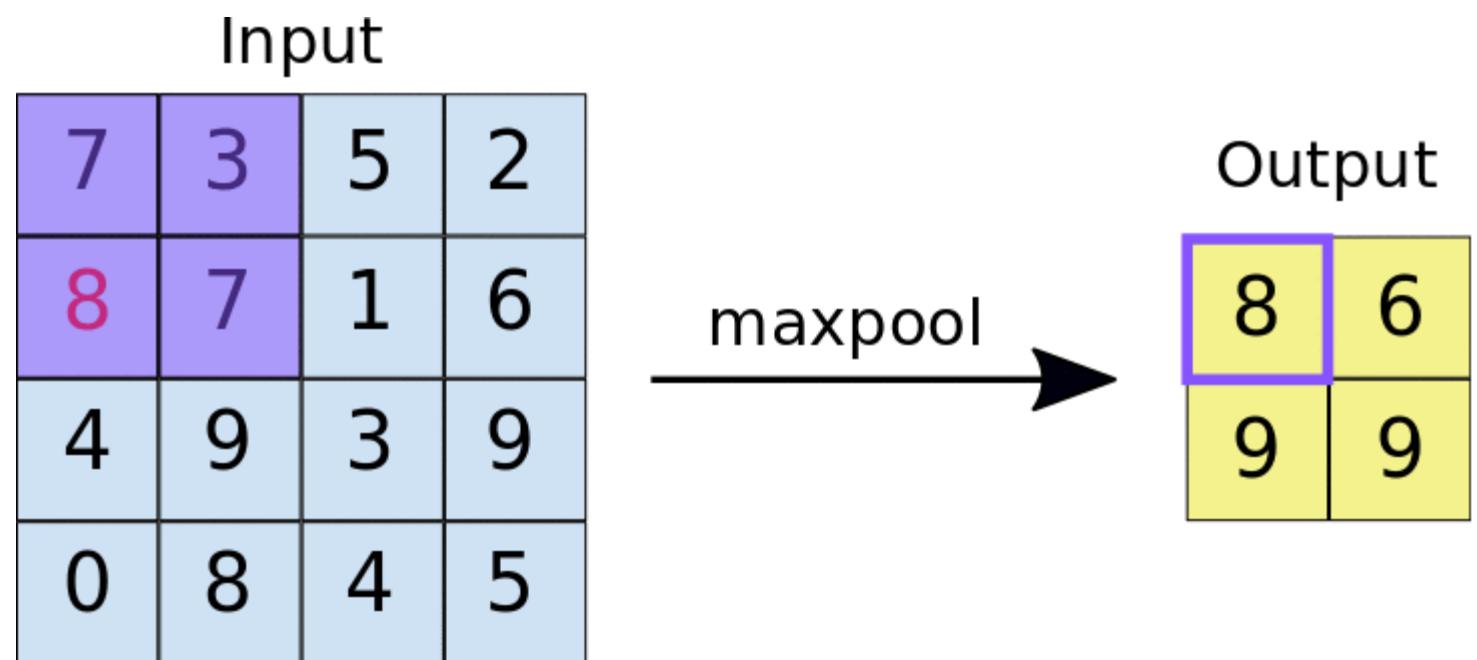
$$\begin{aligned} & 7 \times 1 + 4 \times 1 + 3 \times 1 + \\ & 2 \times 0 + 5 \times 0 + 3 \times 0 + \\ & 3 \times -1 + 3 \times -1 + 2 \times -1 \\ & = 6 \end{aligned}$$

Convolutional Layer



Pooling Layer

- Decreases dimensionality.
- Highlights relevant features.
- Some sort of ‘regularization’



Let's try it in Fashion MNIST!

