

A (Wild) Introduction to Deep Learning

Léonard Boussioux & Angelos Koulouras

15.S60, 2023

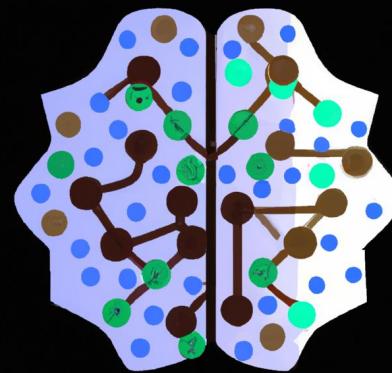


Image by DALLE

Plan

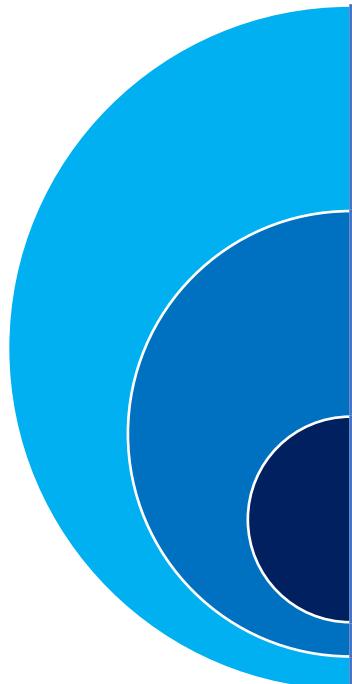
Concepts

- Perceptrons
- Gradient Descent & Training
- Convolutional Neural Networks

Code

- TensorFlow

What is Deep Learning?



Artificial Intelligence

- Techniques enabling computers to mimic human behaviour

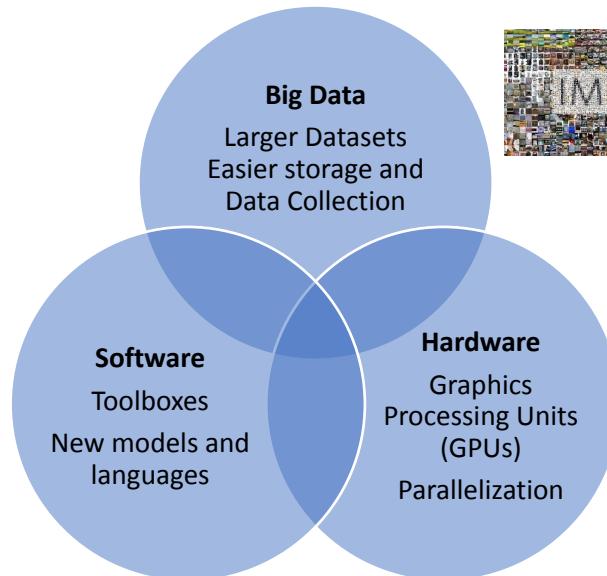
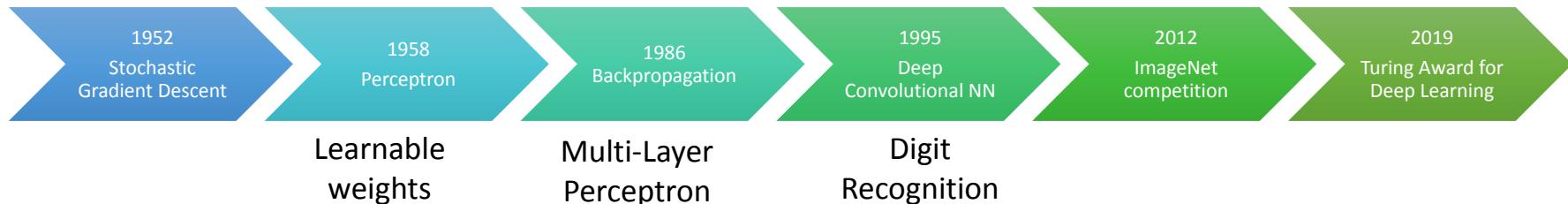
Machine Learning

- When computers can learn without being explicitly programmed

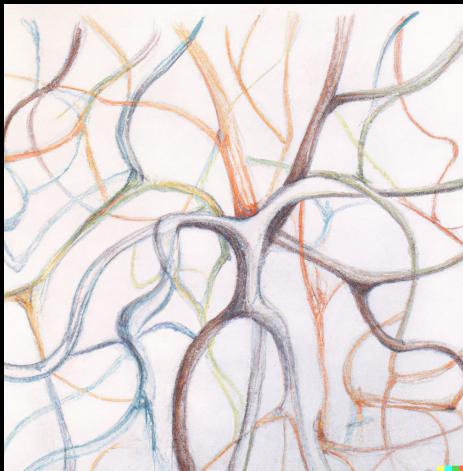
Deep Learning

- Extract patterns using (deep) neural networks

What is the history behind the current AI boom?



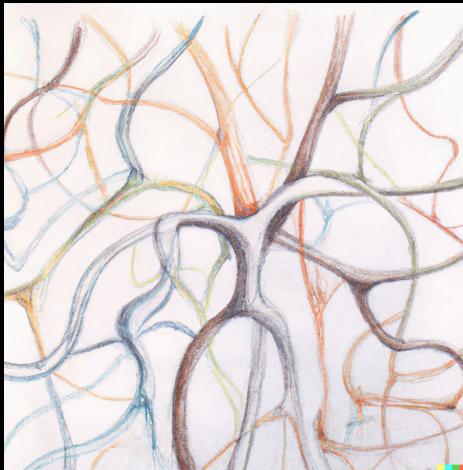
Neural Networks Basics



Images by DALLE



"A neural network"
impressionist style



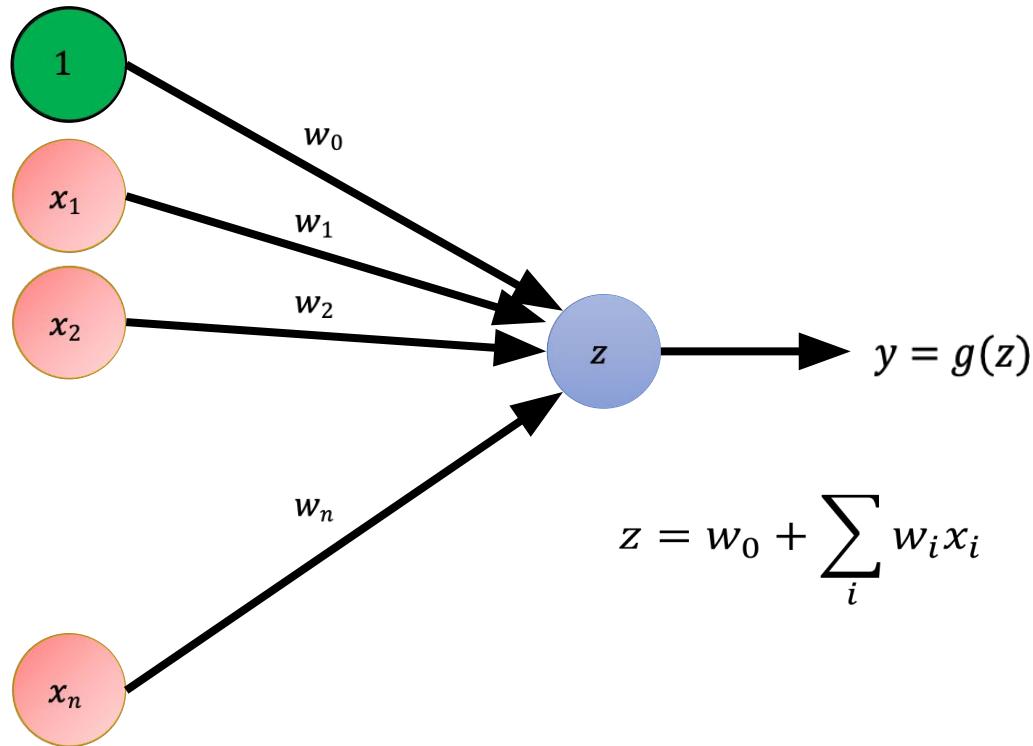
"A neural network"
abstract pencil with
watercolor



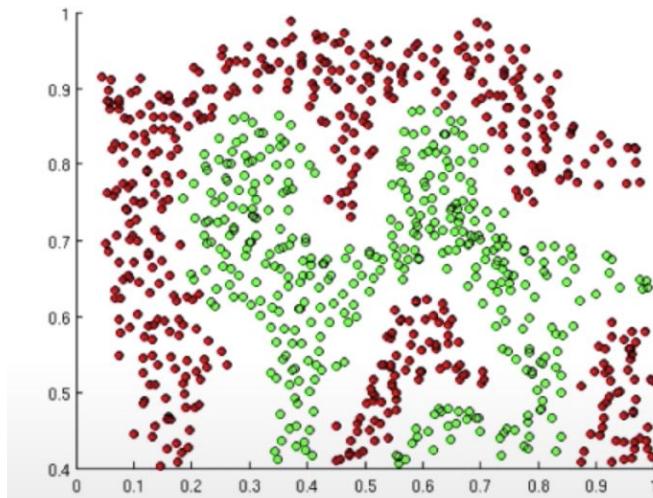
"A neural network"
by Matisse

Images by DALLE

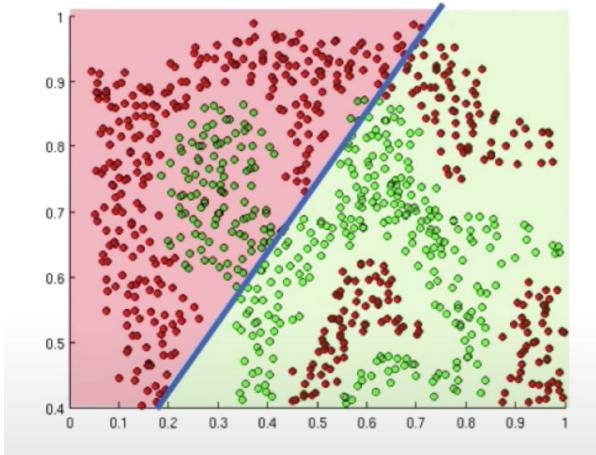
A simplified perceptron



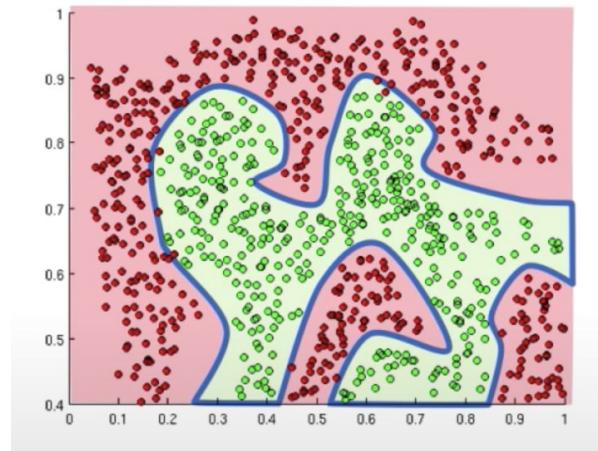
Classification Task



Activation functions

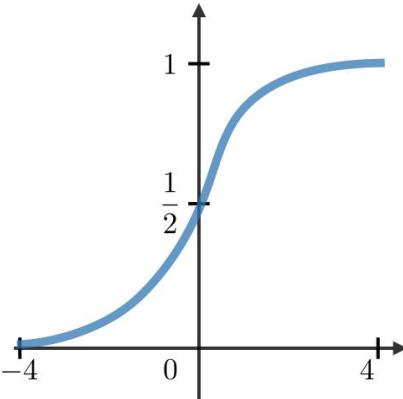
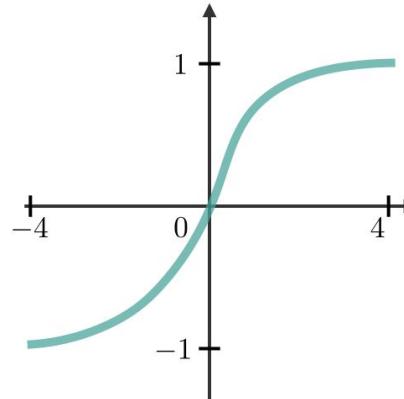
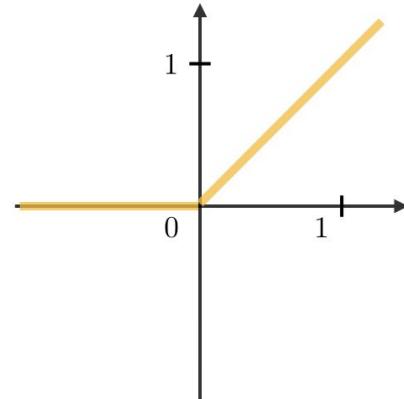


Linear activation functions -> linear decisions

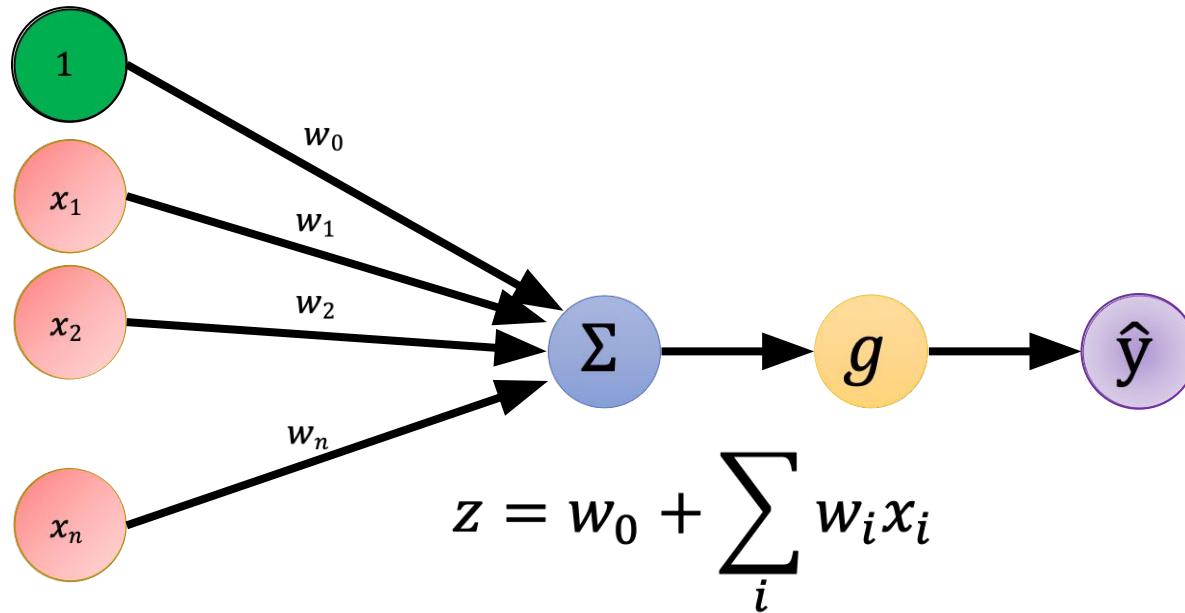


Non-linear activation -> approximate arbitrarily complex functions

Common Activation Functions

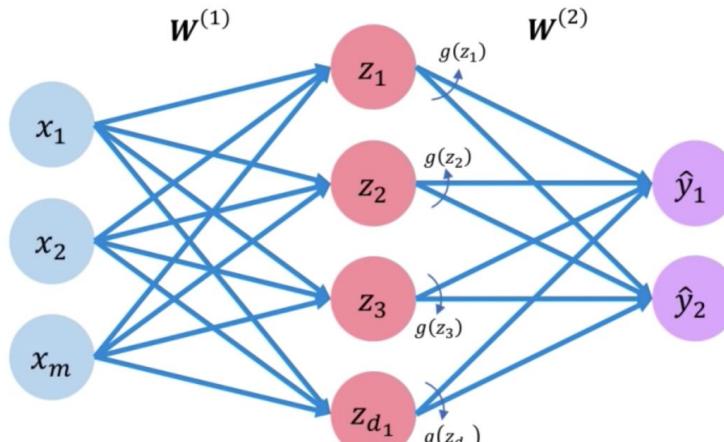
Sigmoid	Tanh	RELU
$g(z) = \frac{1}{1 + e^{-z}}$ 	$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$ 	$g(z) = \max(0, z)$ 

The perceptron: forward propagation



$$\hat{y} = g\left(\sum_i w_i x_i\right)$$

Single Layer Neural Network



Inputs

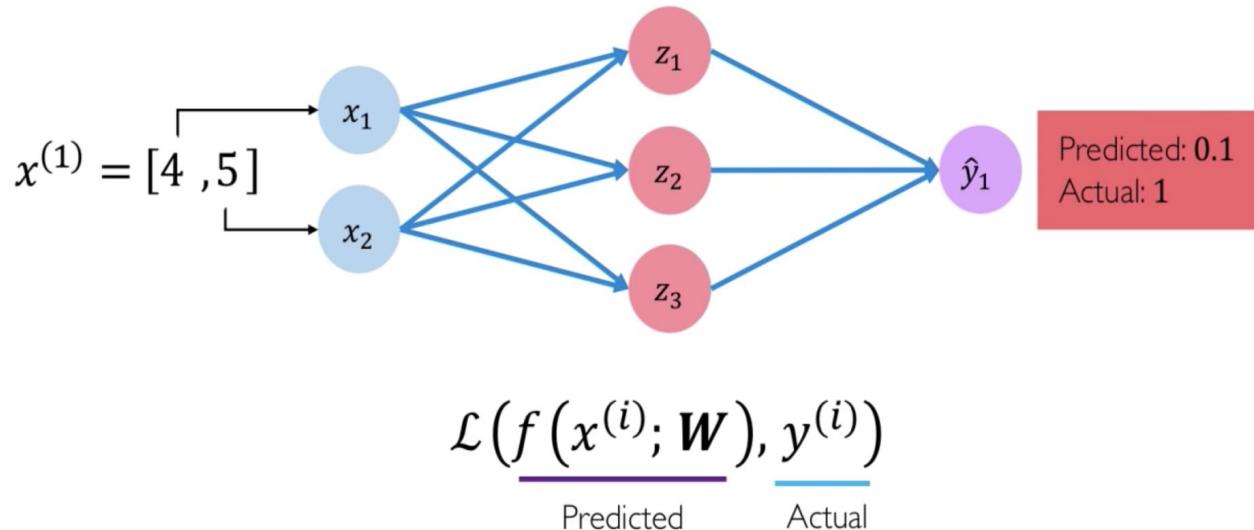
Hidden

Final Output

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

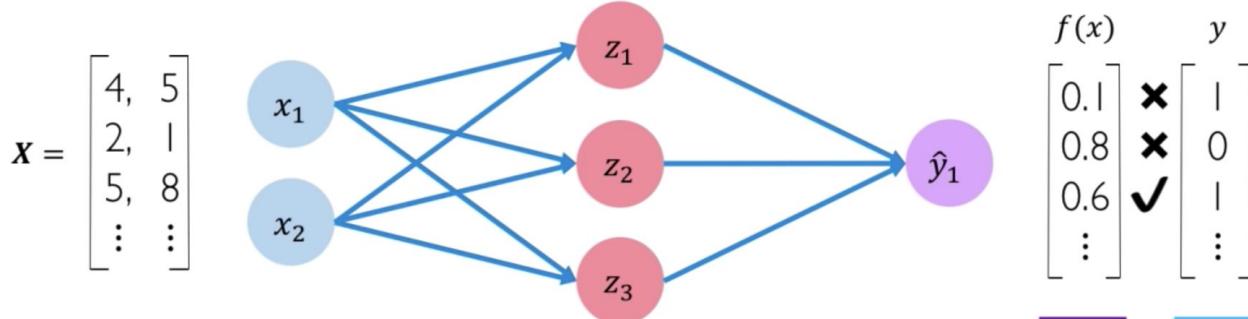
Loss function

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

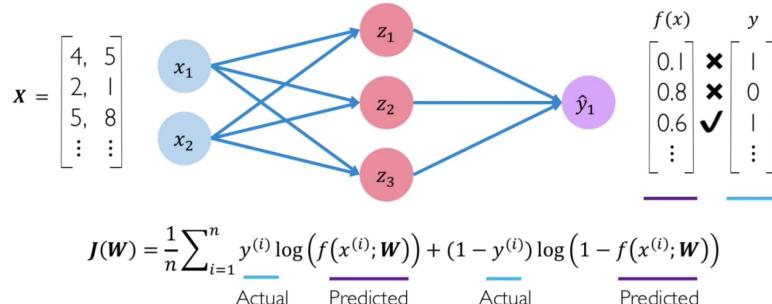
↗
$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \mathbf{W}), y^{(i)})$$

Predicted Actual

Different loss functions

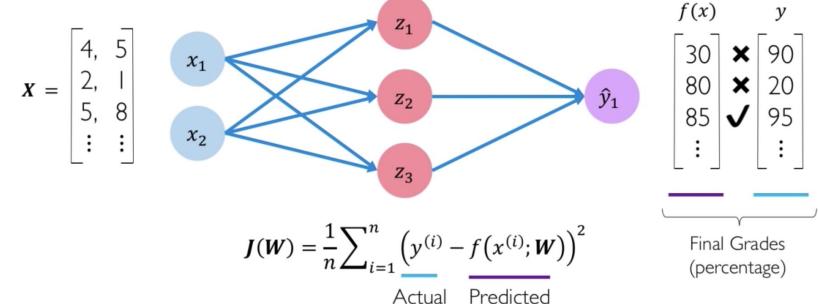
Binary Cross Entropy Loss

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

```
TensorFlow icon: 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.

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

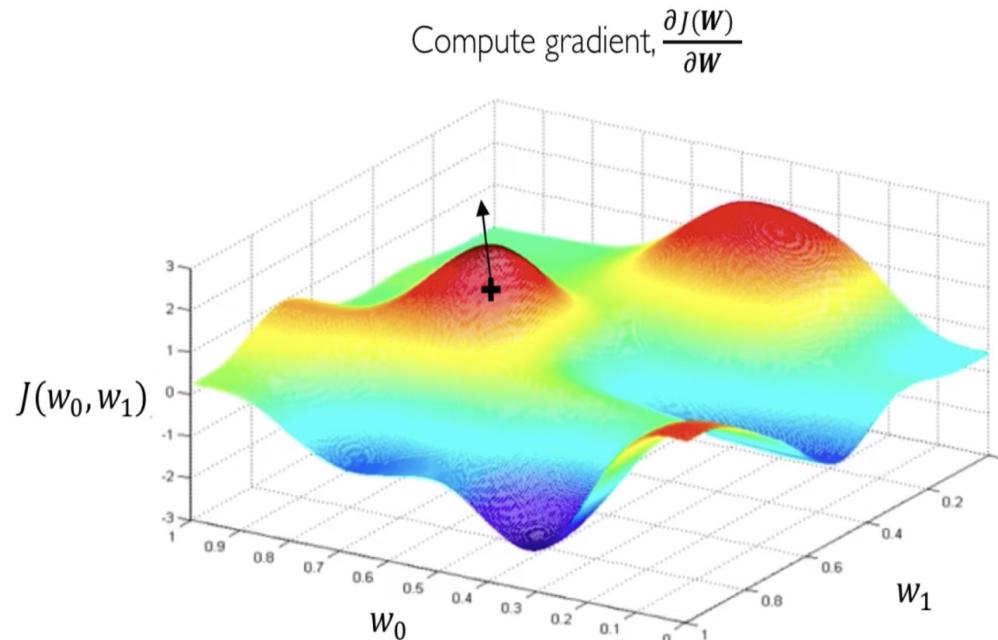
$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} J(\mathbf{W})$$



Remember:

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

But it looks like this...

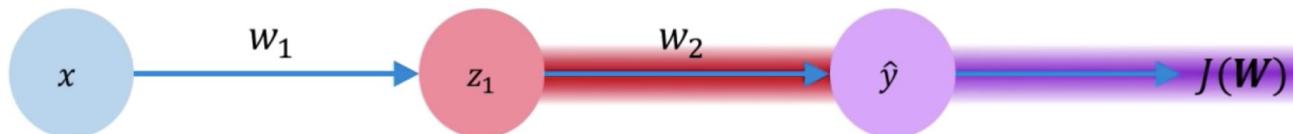


Fortunately, there is Gradient Descent!

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

And backpropagation with chain rule



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

```
graph LR; J_W["J(W)"] --- P1["purple bar"]; P1 --- P2["purple bar"]; P2 --- J_W;
```

But how do you choose the learning rate?

Remember:

Optimization through gradient descent

Ideas?

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

How can we set the
learning rate?

Adaptive 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.)

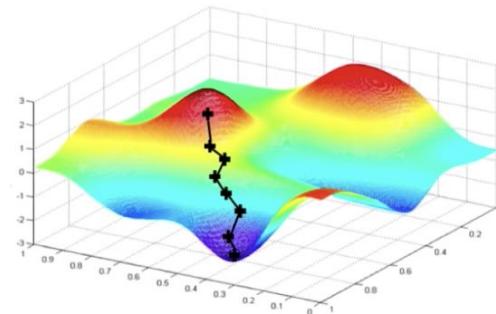
Computing gradients is expensive!

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!

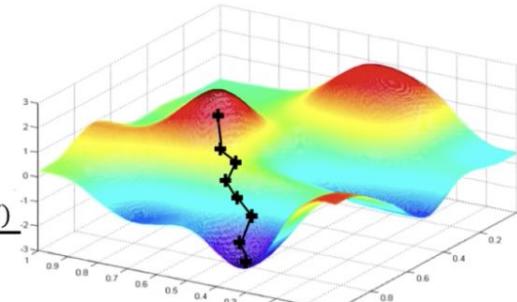


Computing gradients is expensive!

Stochastic Gradient Descent

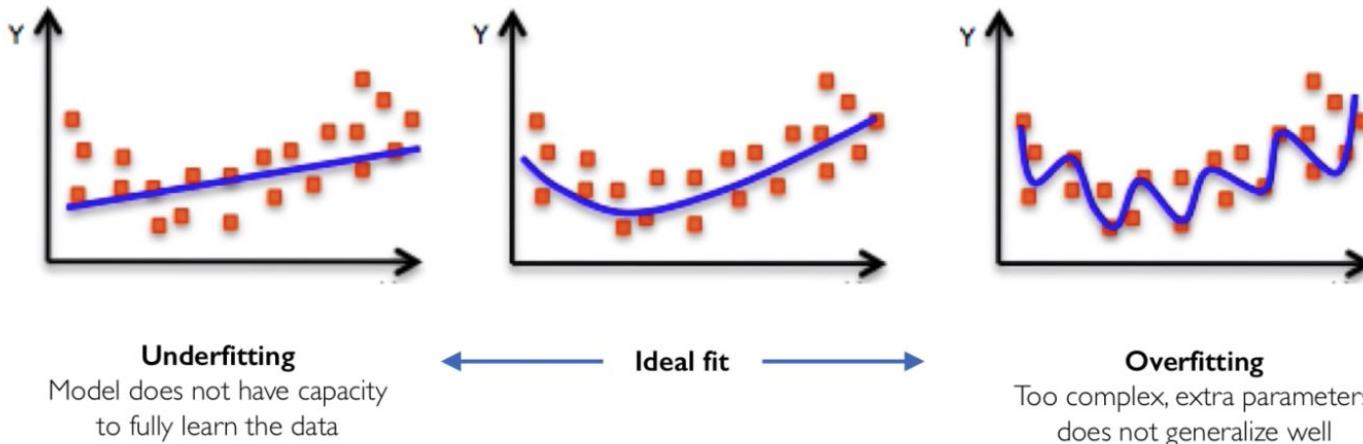
Algorithm

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

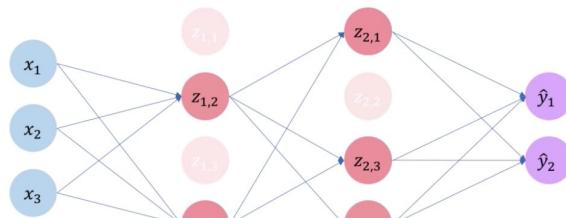


Let's overcome this! Ideas?

Regularization I: Dropout

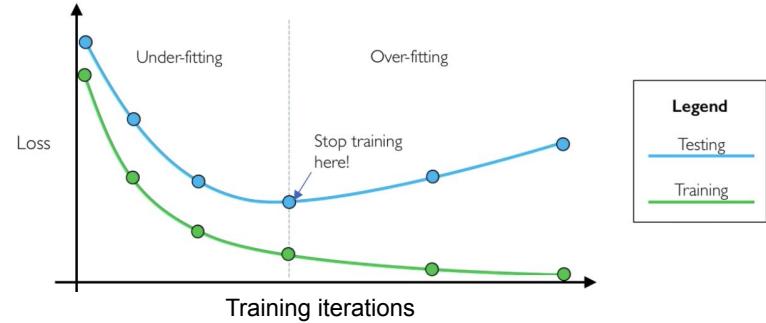
- During training, randomly set some activations to 0
 - Typically 'drop' 50% of activations in layer
 - Forces network to not rely on any 1 node

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



Regularization 2: Early Stopping

- Stop training before we have a chance to overfit

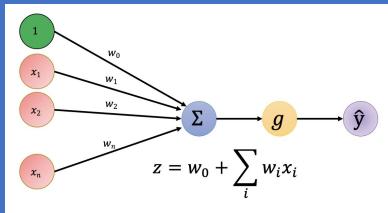


Regularization 3: L1 or L2 penalty on the weights

Review

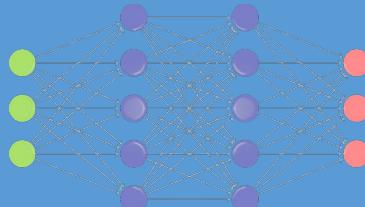
The Perceptron

- Structure
- Nonlinear activation functions



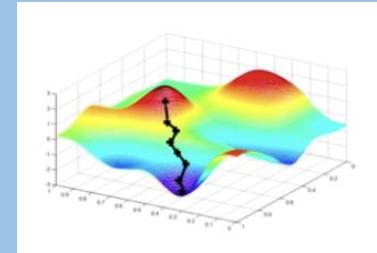
Neural Networks

- Stacking layers of perceptrons
- Backpropagation and gradient descent



Training in Practice

- Adaptive learning
- Batching
- Regularization

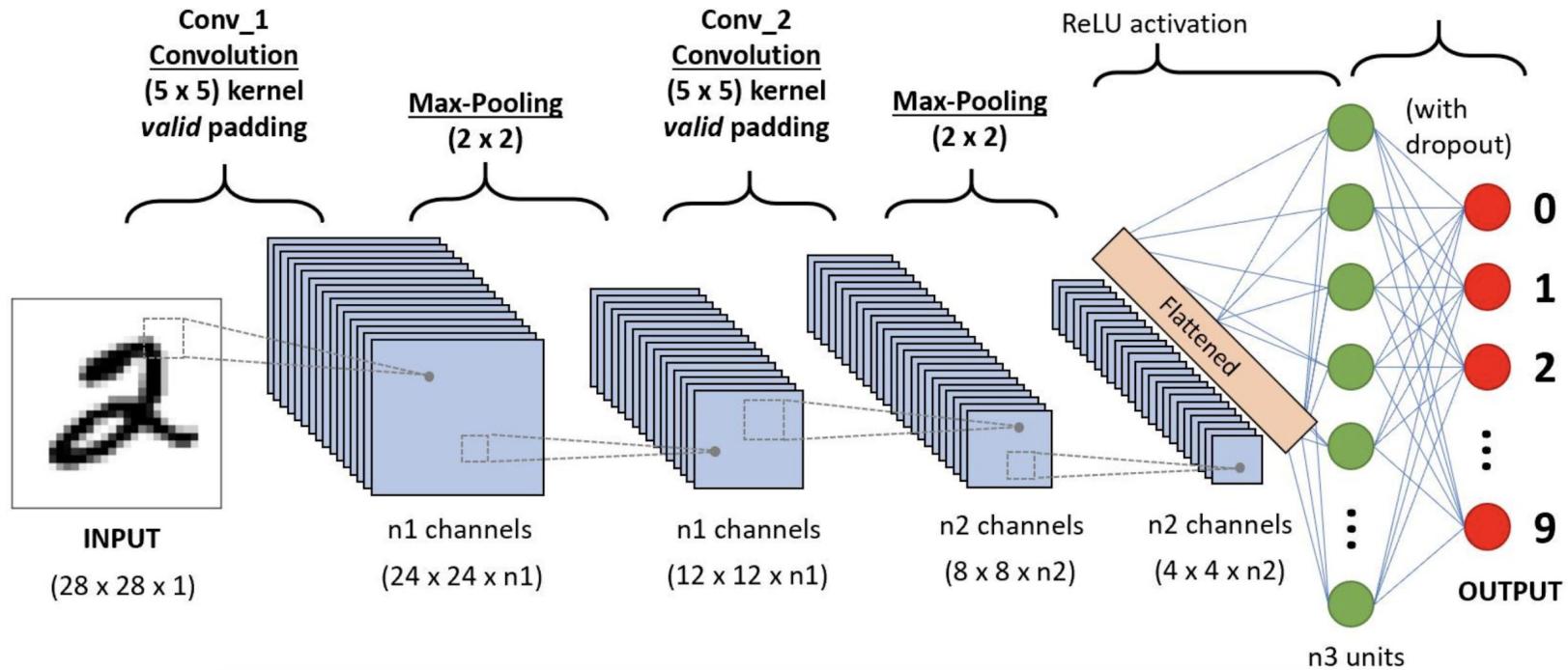


Please ask any questions!

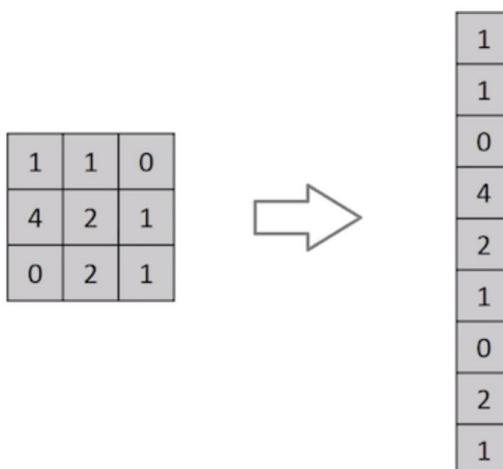
Break

Convolutional Neural Networks

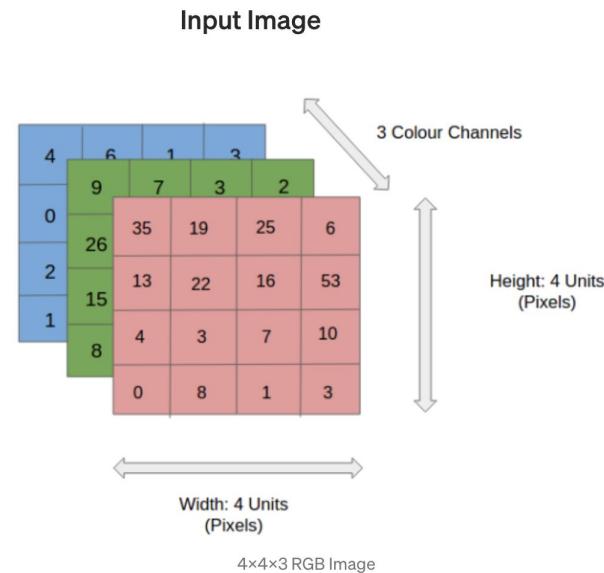
Let's now deal with images!



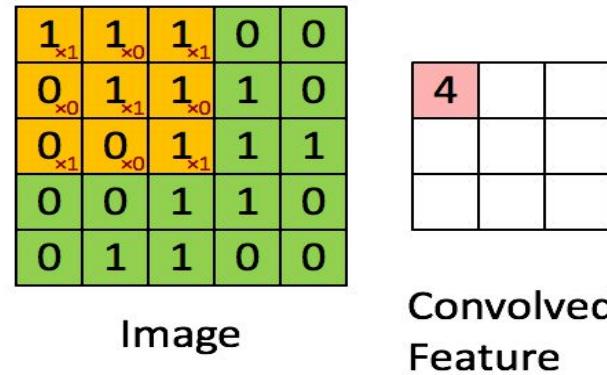
3-dimensional data



Flattening of a 3×3 image matrix into a 9×1 vector



Convolution Operation

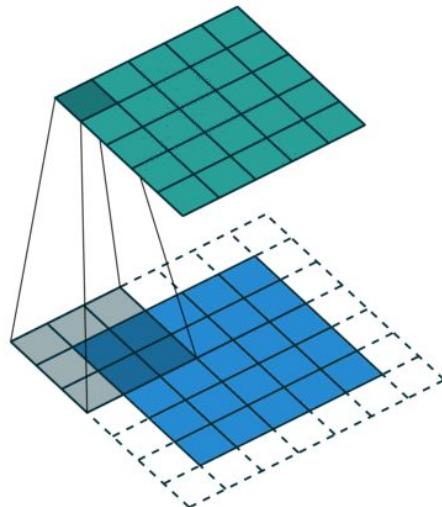


Images source:

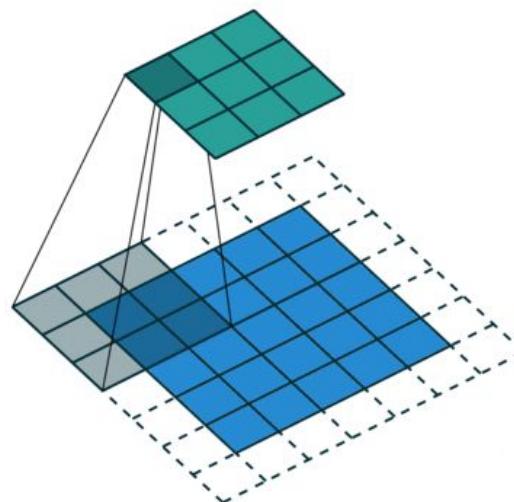
<https://towardsdatascience.com/a-comprehensive-guide-to-convolutional-neural-networks-the-eli5-way-3bd2b1164a53>

Convolution Operation

Padding: 1
Stride: 1



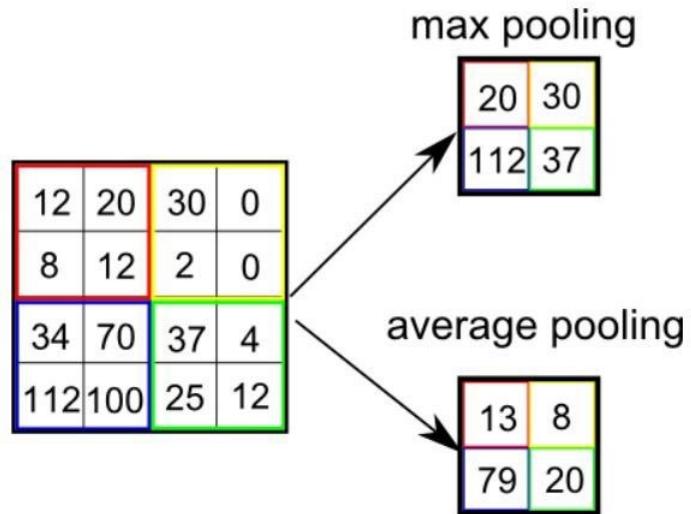
Padding: 1
Stride: 2



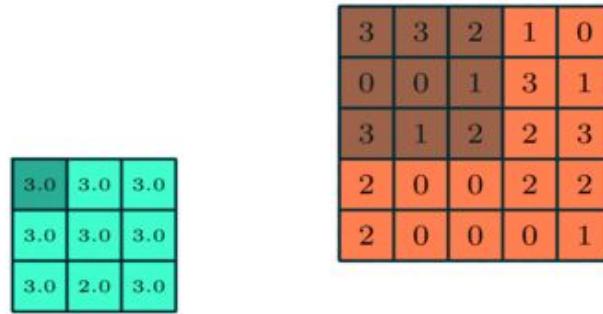
Images source:

<https://towardsdatascience.com/a-comprehensive-guide-to-convolutional-neural-networks-the-eli5-way-3bd2b1164a53>

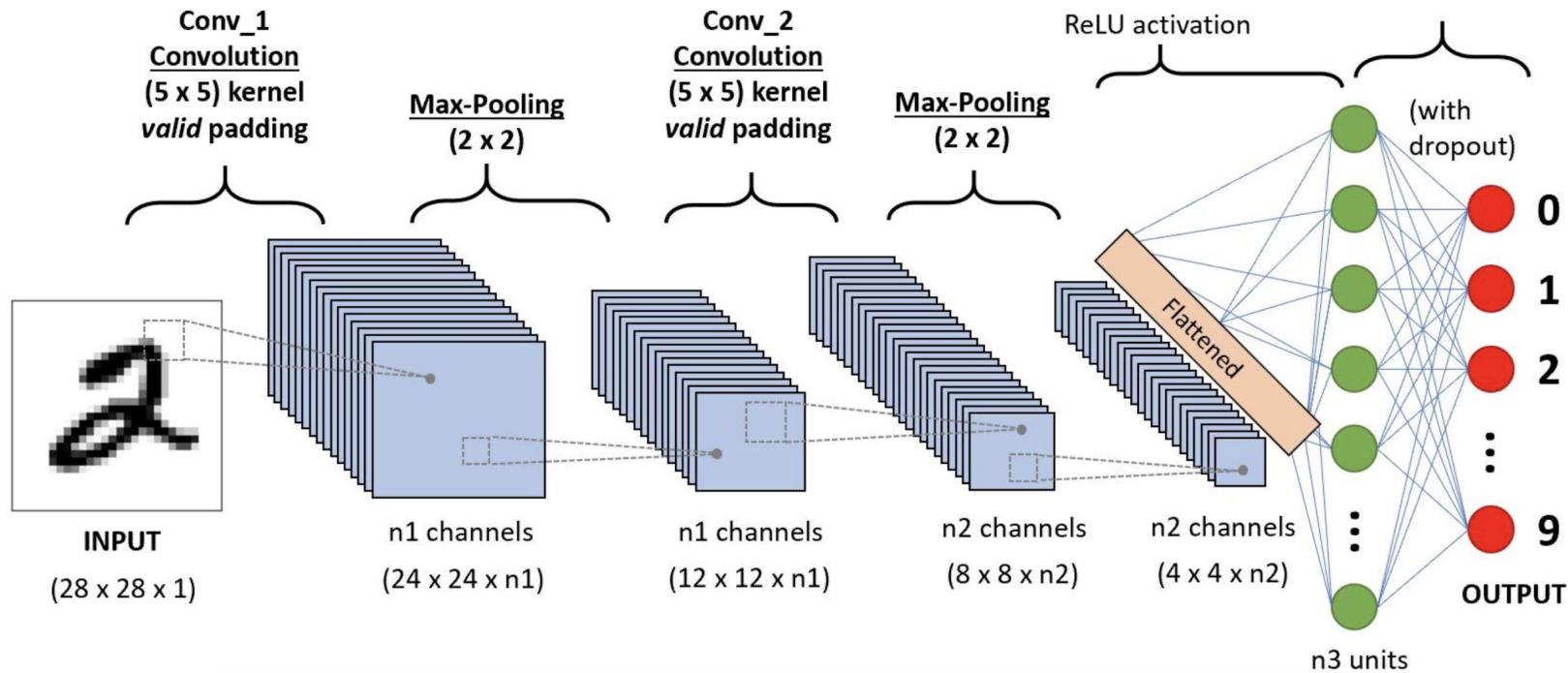
Pooling



Pooling



Now we can deal with images!



What kind of tasks can we perform?

- Classification, Regression
- Segmentation
- Prediction of next image
- Image generation
- ...

Please ask any questions!

Thank you for your attention!

+

.

o