

Deep Learning

Aymeric DIEULEVEUT

May 2024

1 Goals

2 The first Neural Network: the Perceptron

3 Convolutional Neural Networks

4 Applications

- Image classification
- Pose, action detection
- Object detection and segmentation
- Scene labeling - Semantic segmentation
- Object tracking - videos

Outline

1 Goals

2 The first Neural Network: the Perceptron

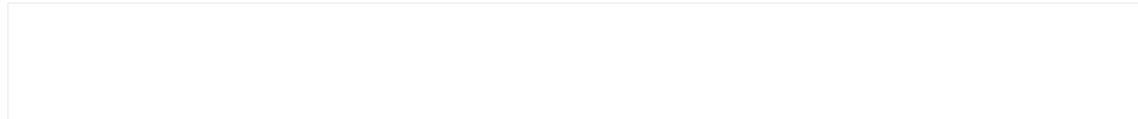
3 Convolutional Neural Networks

4 Applications

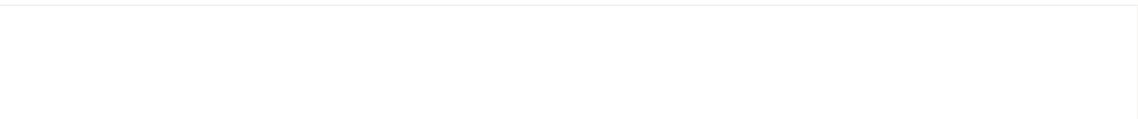
- Image classification
- Pose, action detection
- Object detection and segmentation
- Scene labeling - Semantic segmentation
- Object tracking - videos

Deep Learning in one slide + Goals of the afternoon

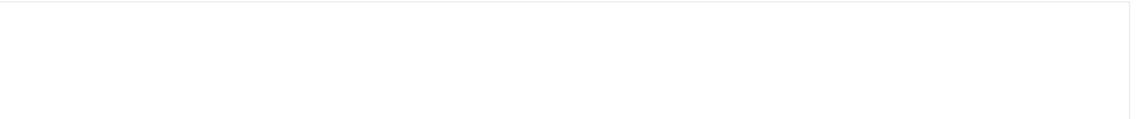
- How does it work:



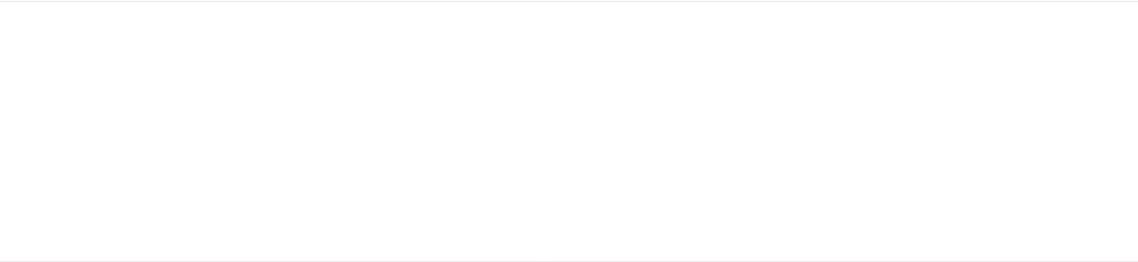
- What does it require:



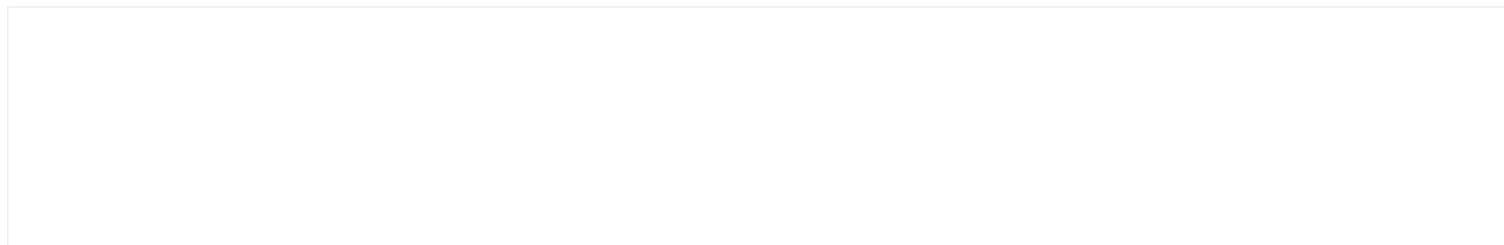
- Why now:



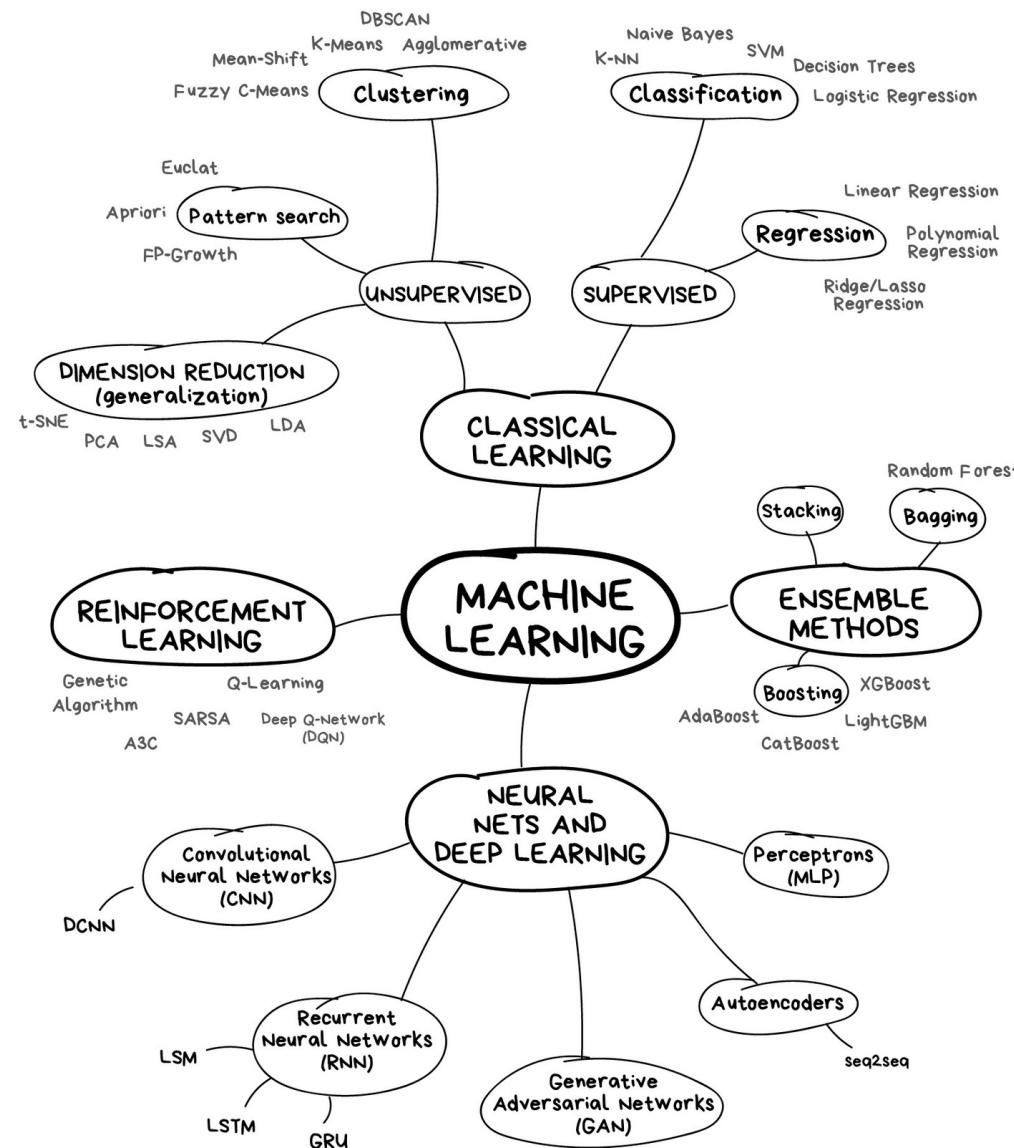
- Some Applications



Goals: Understanding core concepts of Deep Learning.



Machine Learning



Summary of previous lectures on Machine Learning

① General framework: loss, risk, generalization risk vs training risk

Training set
 $(X_{\text{train}}, Y_{\text{train}})$
User choices

Loss function, Generalization Risk and data-dependent predictors

Loss Function

- Loss function: $\ell(y, g(x))$ quantifies the quality of the prediction $g(x)$ of y , e.g.:
 - ▶ 0-1 Loss (classification): $\ell(y, g(x)) = 1_{y \neq g(x)}$, $y \in \mathcal{Y} = \{-1, 1\}$
 - ▶ Quadratic Loss (regression): $\ell(y, g(x)) = \|y - g(x)\|^2$, $y \in \mathbb{R}^d$

Risk of a Decision Rule = average loss

- Risk: $R(g) = \mathbb{E}[\ell(Y, g(x))] = \int \ell(y, g(x)) d\mathbb{P}(x, y)$, e.g.:
 - ▶ 0-1 Risk (classification): $R(g) = \mathbb{P}(Y \neq g(x))$
 - ▶ Quadratic Risk (regression): $R(g) = \mathbb{E}\|Y - g(x)\|^2$
- ↪ R is also referred to as *Generalization risk*, or *True risk*.

⚠ The distribution \mathbb{P} is unknown.

Data-dependent predictors

- Training set: $\mathcal{D}_n = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$ (*i.i.d. $\sim \mathbb{P}$*)
- Goal: build a **data dependent predictor / classifier \hat{g}_n** based on the training data
- That has a **minimal generalization risk $R(\hat{g}_n)$** .

$R(\hat{g}_n)$ = average loss on a "new point" $(X, Y) \sim \mathbb{P}$, independent from D_n

↪ Our goal is to predict well on inputs/outputs pairs that we have not seen in the dataset, i.e. *generalize well*.

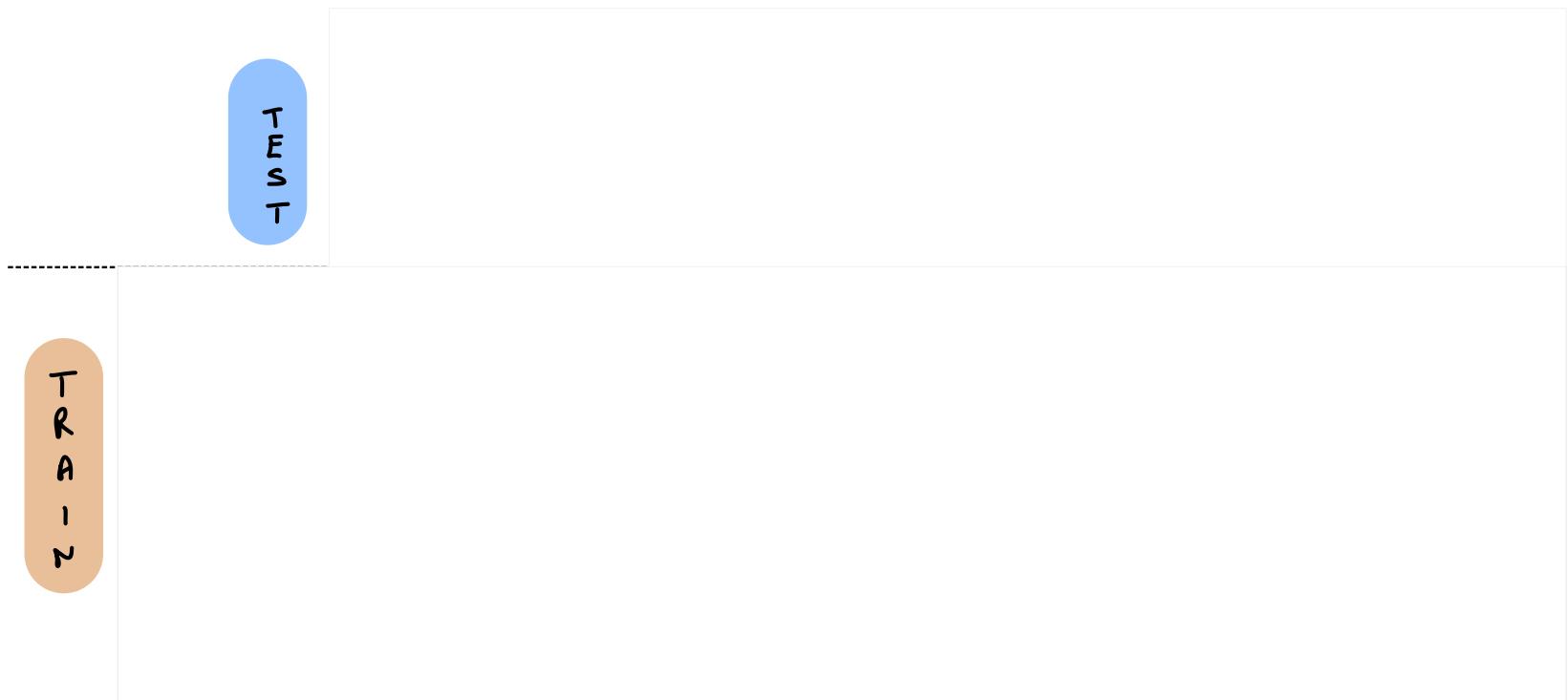
THM 1 .

Summary of previous lectures on Machine Learning

① General framework: loss, risk, generalization risk vs training risk

② SML Pipeline

- ▶ Train (Fit + validation) / Test / Deployment
- ▶ Evaluation of validation scores through cross validation to pick the best method/hyperparams



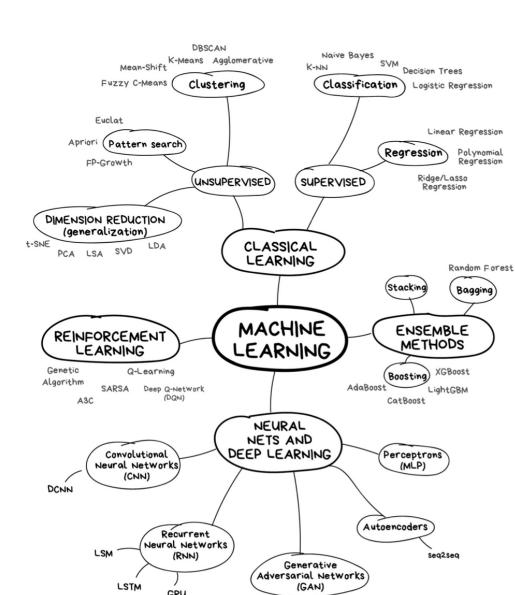
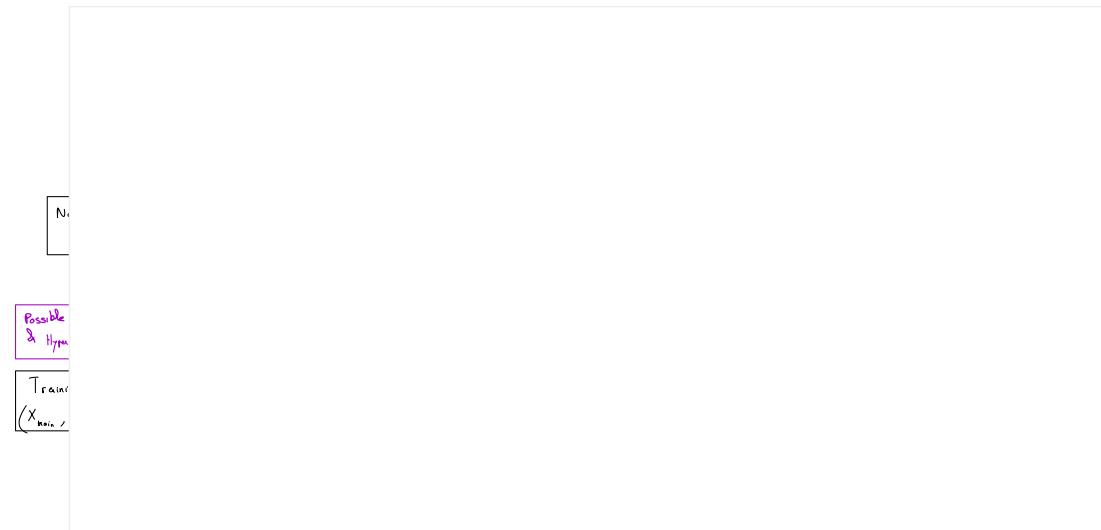
Summary of previous lectures on Machine Learning

① General framework: loss, risk, generalization risk vs training risk

② SML Pipeline

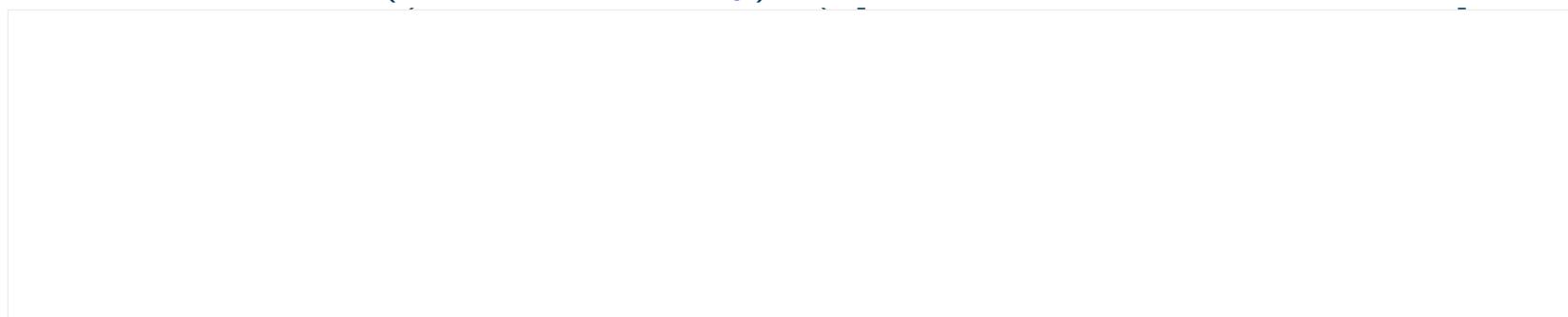
- ▶ Train (Fit + validation) / Test / Deployment
- ▶ Evaluation of validation scores through cross validation to pick the best method/hyperparams

③ Methods: Linear, Logistic regression, Decision Trees, Random Forests, Boosting



Summary of previous lectures on Machine Learning

- ① General framework: loss, risk, generalization risk vs training risk
- ② SML Pipeline
 - ▶ Train (Fit + validation) / Test / Deployment
 - ▶ Evaluation of validation scores through cross validation to pick the best method/hyperparams
- ③ Methods: Linear, Logistic regression, Decision Trees, Random Forests, Boosting
- ④ Opening the black box: for each method
 - ▶ How to predict (↪ interpretability)



$\hat{P}_n = \boxed{\quad}$: Lin. Reg	Polynomial Linear reg	Logistic Reg	Dec Trees	RF/ Boosting
How to predict					
How to fit					
Key hyperparameters					
\oplus					
\ominus					
Overfitting / Underfitting → regularization technique					
Approach					

Outline

1 Goals

2 The first Neural Network: the Perceptron

3 Convolutional Neural Networks

4 Applications

- Image classification
- Pose, action detection
- Object detection and segmentation
- Scene labeling - Semantic segmentation
- Object tracking - videos

From logistic regression to Multi Layer Perceptron 1

Supervised ML problem:

- ① Data $D_n = (X_i, Y_i)_{i=1,\dots,n}, X_i \in \mathbb{R}^d, Y_i \in \{0, \dots, K\}$
- ② Goal: Predict a the label for a new point.

2 approaches:

Generative approach

- Define a model on your data (e.g., logit model).
- Estimate parameter (likelihood maximization).

Discriminative approach

- No statistical model !
- Define a loss function ℓ
- Goal:

$$\min_{f \in \mathcal{F}} \mathbb{E}[\ell(f(X), Y)]$$

- Approach:

$$\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n [\ell(f(X_i), Y_i)]$$

From logistic regression to Multi Layer Perceptron 1

Supervised ML problem:

- ① Data $D_n = (X_i, Y_i)_{i=1,\dots,n}, X_i \in \mathbb{R}^d, Y_i \in \{0, \dots, K\}$
- ② Goal: Predict a the label for a new point.

2 approaches:

Generative approach

- Define a model on your data (e.g., logit model).
- Estimate parameter (likelihood maximization).

Discriminative approach

- No statistical model !
- Define a loss function ℓ
- Goal:

$$\min_{f \in \mathcal{F}} \mathbb{E}[\ell(f(X), Y)]$$

- Approach:

$$\min_{f \in \mathcal{F}} \frac{1}{n} \sum_{i=1}^n [\ell(f(X_i), Y_i)]$$

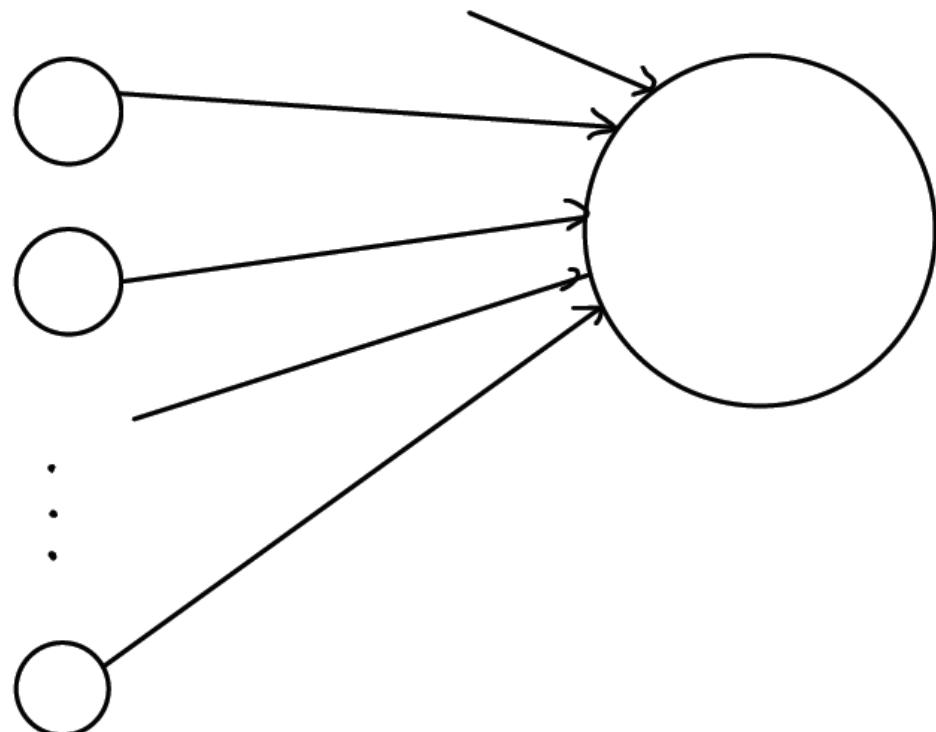
Logistic regression can be seen both ways !

$$\min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-Y_i \cdot w^T X_i))$$

Let's NN playground: [Link](#)

From logistic regression to Multi Layer Perceptron 3

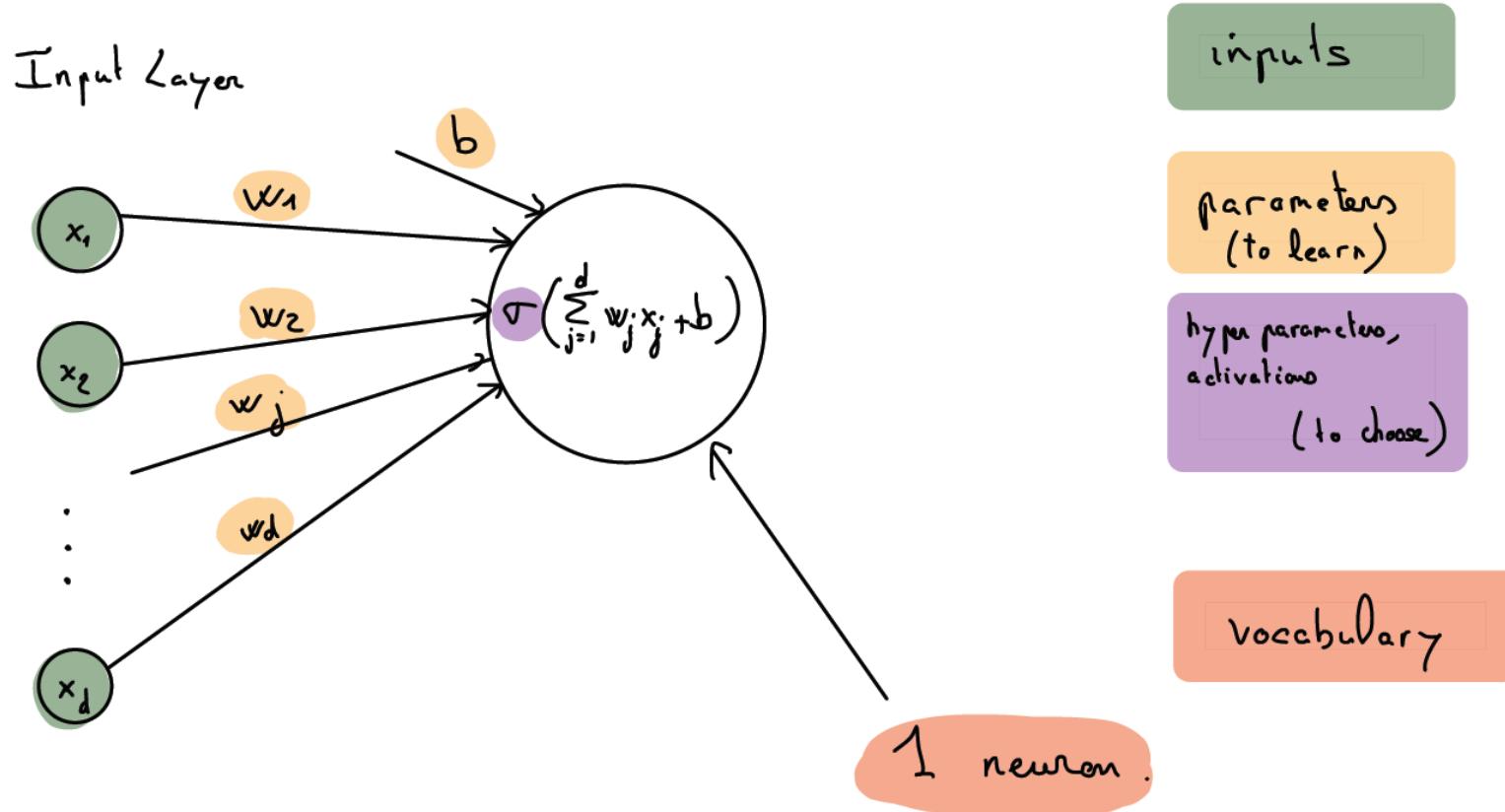
- Class of functions for logistic regression: $\mathcal{F} = \{X \mapsto w^T X, w \in \mathbb{R}^d\}$.
- Prediction for a new point is $1_{\sigma(X_i) > 1/2}$
- Can be seen as predicting a probability $\sigma(X_i)$ that the output is 1



This is a neural network, with one neuron !

From logistic regression to Multi Layer Perceptron 3

- Class of functions for logistic regression: $\mathcal{F} = \{X \mapsto w^T X, w \in \mathbb{R}^d\}$.
- Prediction for a new point is $1_{\sigma(X_i) > 1/2}$
- Can be seen as predicting a probability $\sigma(X_i)$ that the output is 1

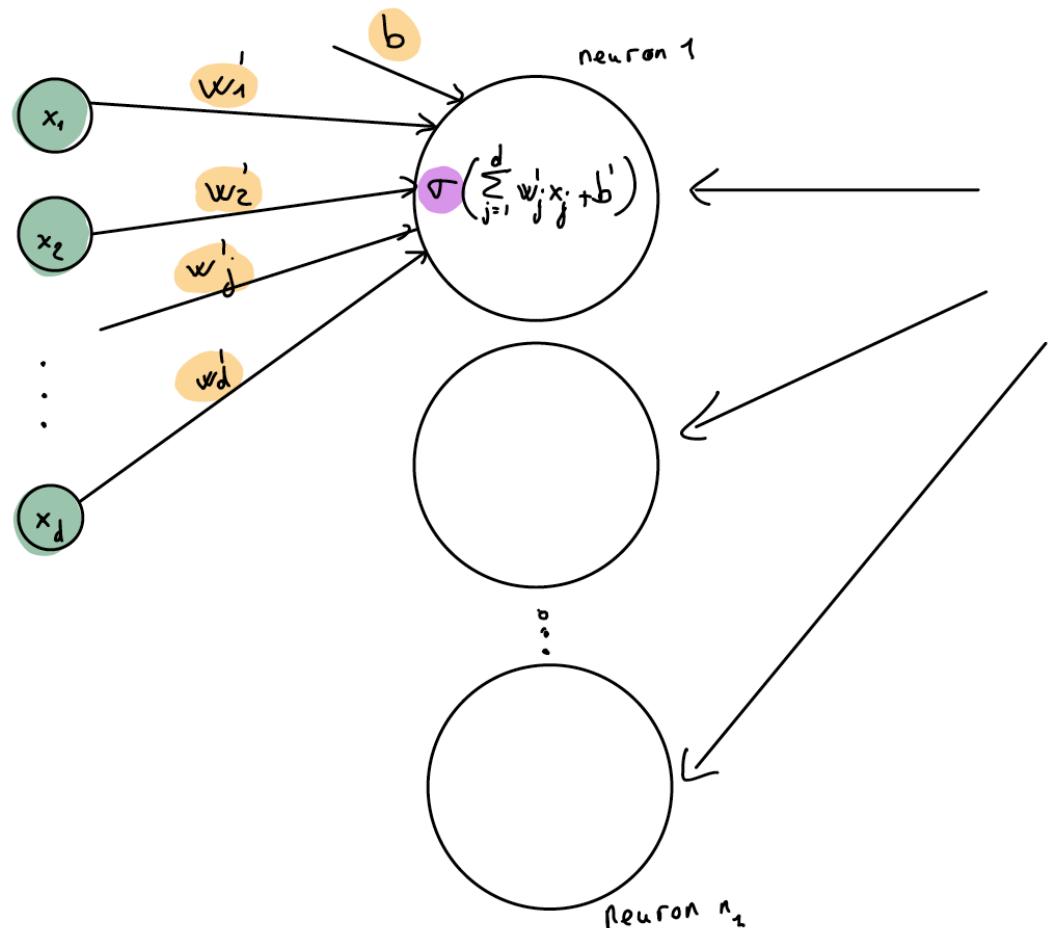


This is a neural network, with one neuron !

From logistic regression to Multi Layer Perceptron 4

Extend to

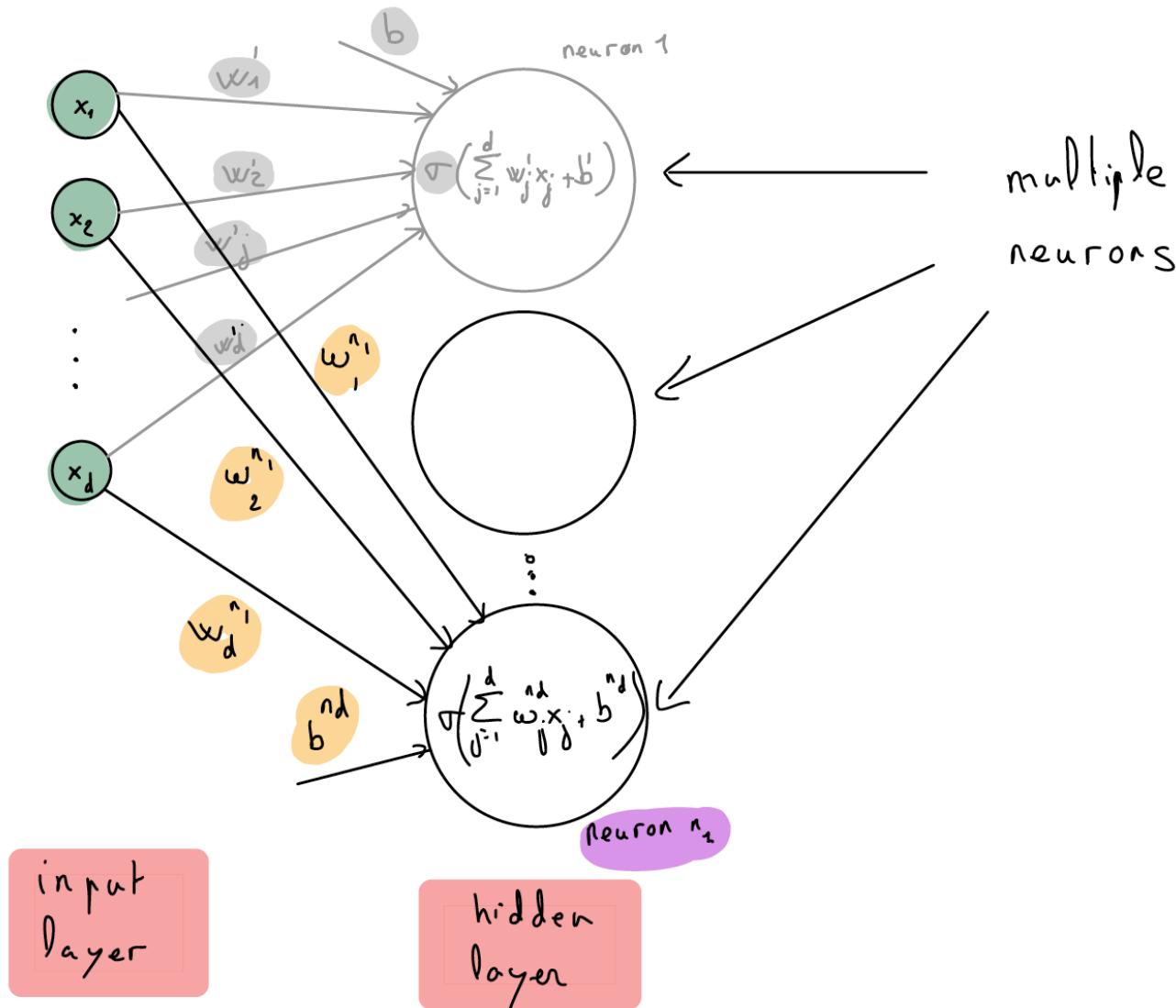
- Multiple Neurons
- Multiple Layers



From logistic regression to Multi Layer Perceptron 4

Extend to

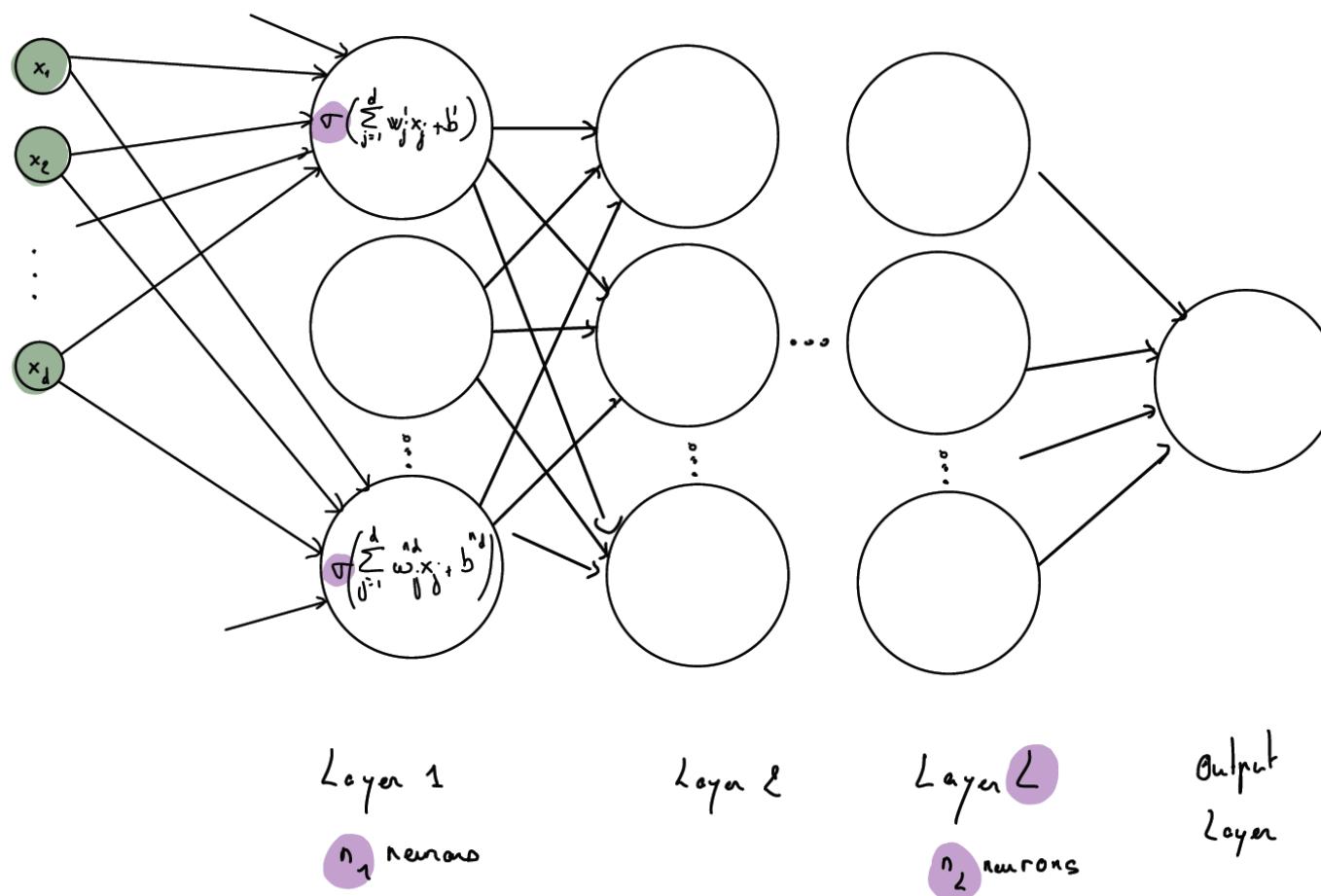
- Multiple Neurons
- Multiple Layers



From logistic regression to Multi Layer Perceptron 5

Extend to

- Multiple Neurons
- Multiple Layers



From logistic regression to Multi Layer Perceptron 5

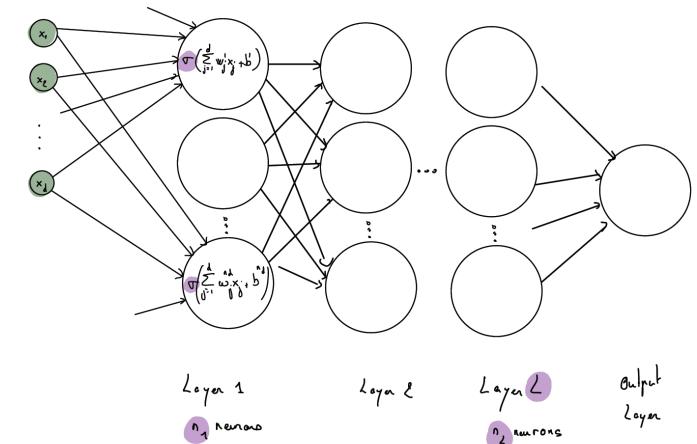
Neural network:

- Goal $\min_{f \in \mathcal{F}_{\text{MLP}}} \frac{1}{n} \sum_{i=1}^n [\ell(f(X_i), Y_i)]$

- Class of non linear functions.

$$\mathcal{F}_{\text{MLP}} = \{f(X; W, b) = \sigma(b_n + W_n \sigma(\dots(b_1 + W_1 x))),$$

$$W \in \mathbb{R}^{n \times n}, b \in \mathbb{R}^n\}$$



Summary: NN generalize regression by looking for candidates in a much larger class of functions than linear ones.

What needs to be learned

What needs to be chosen

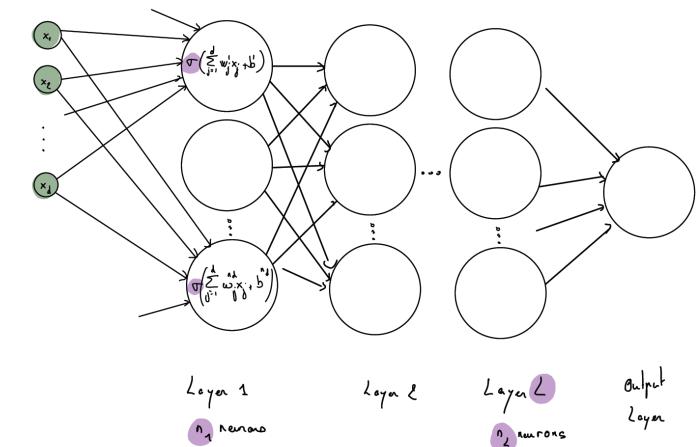
Vocabulary

From logistic regression to Multi Layer Perceptron 5

Neural network:

- Goal $\min_{f \in \mathcal{F}_{\text{MLP}}} \frac{1}{n} \sum_{i=1}^n [\ell(f(X_i), Y_i)]$
- Class of non linear functions.
$$\mathcal{F}_{\text{MLP}} = \{f(X; W, b) = \sigma(b_n + W_n \sigma(\dots(b_1 + W_1 X))),$$

$$W \in \mathbb{R}^{n \times n}, b \in \mathbb{R}^n\}$$



Summary: NN generalize regression by looking for candidates in a much larger class of functions than linear ones.

What needs to be learned

Parameters W, b

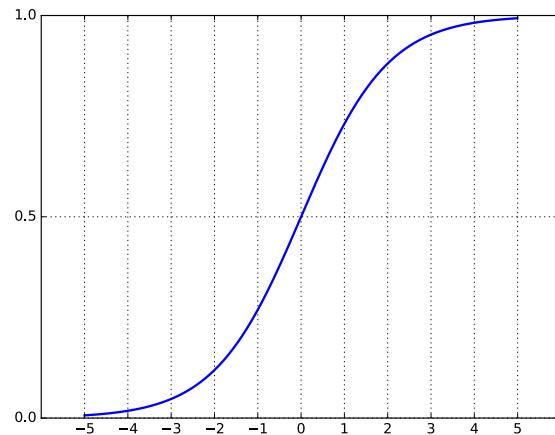
What needs to be chosen

- Activation σ
- Number of layers L
- Number of neurons per hidden layer
 $n_i, i = 1, \dots, L$.

Vocabulary

- Neuron
- Activation
- Hidden Layer
- Width, Depth
- Fully connected layer

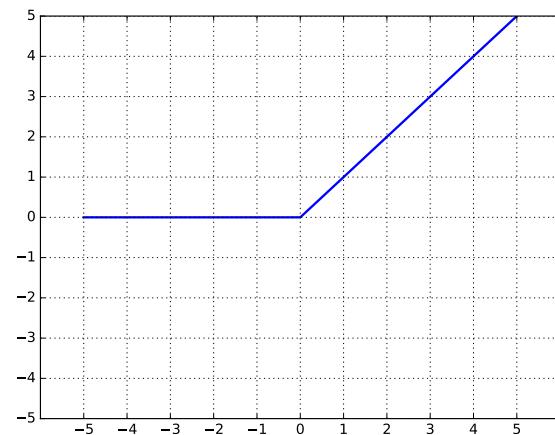
2 possible activations: Sigmoid and Rectified Linear Unit



Sigmoid function

- $x \mapsto \frac{\exp(x)}{1+\exp(x)}$
- Problems:

- ① Saturated function: gradient killer -> need for rescaling data



Rectified Linear Unit (ReLU)

- $x \mapsto \max(0, x)$
- Pros & cons:
 - ① Not a saturated function. Kills negative values.
 - ② Empirically, convergence is faster than sigmoid/tanh.
 - ③ Plus: biologically plausible

How to deal with multi-class output: softmax activation

When we do multi-class classification, i.e., $Y_i \in 1, \dots, K$, we output K different values:

- Each of them corresponds to the probability of belonging to the class j , $1 \leq j \leq K$
- To obtain probabilities (adding up to 1):
 - ① We have K neurons on the last layer
 - ② And use a **Softmax** activation to renormalize.

Softmax output unit, used to predict $\{1, \dots, K\}$:

$$\text{softmax}(z)_i = \frac{e^{z_i}}{\sum_{k=1}^K e^{z_k}}$$

Up to that point, we have seen :

- What a Neural Network was
- Why it is expected to learn better.

Next Question: how to implement it ?

Python - Keras

1. Hardware:

- CPU
- GPU
- TPU

2. Software (Python packages):

- Pytorch/Tensorflow
- → Keras : TF high level API. Ideal for applications.



Keras - A few lines !¹

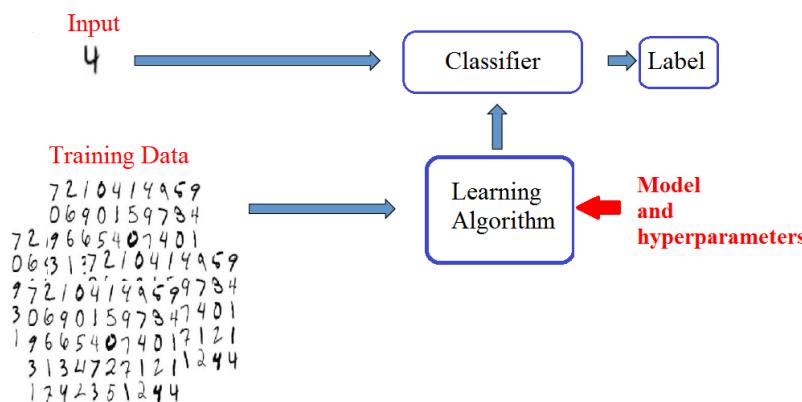
What do we need to specify?

What do we need to do next?

Keras - A few lines !¹

What do we need to specify?

What do we need to do next?



Intuition

Up to that point, we have seen :

- What a Neural Network was.
- How to implement it in Python.

Next Question: why and when does it work?

Intuition

Up to that point, we have seen :

- What a Neural Network was.
- How to implement it in Python.

Approximation Theorem

Goal

$$\min_{f \in \mathcal{F}_{\text{MLP}}} \frac{1}{n} \sum_{i=1}^n [\ell(f(X_i), Y_i)].$$

As said before, a MLP can output non-linear functions... But how powerful is it?

Approximation Theorem

Goal

$$\min_{f \in \mathcal{F}_{\text{MLP}}} \frac{1}{n} \sum_{i=1}^n [\ell(f(X_i), Y_i)].$$

As said before, a MLP can output non-linear functions... But how powerful is it?

Continuous Neural Networks with one single hidden layer and any bounded and non constant activation function can approximate any function in L^p , provided a sufficient number of hidden units.

- ↪ Very powerful in terms of approximations.
- Not very surprising: non linear function with millions of parameters!

Optimization - Fitting the network

Goal

$$\min_{f \in \mathcal{F}_{\text{MLP}}} \frac{1}{n} \sum_{i=1}^n [\ell(f(X_i), Y_i)].$$

How to minimize a function?

Gradient Methods. = cheapest way to update (learn) the weights iteratively to minimize the loss.

Optimization - Fitting the network

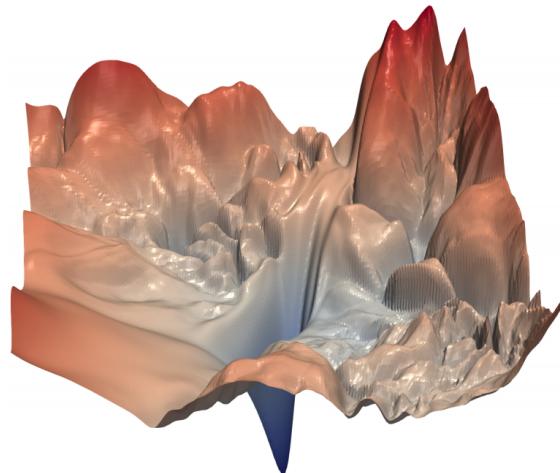
Goal

$$\min_{f \in \mathcal{F}_{\text{MLP}}} \frac{1}{n} \sum_{i=1}^n [\ell(f(X_i), Y_i)].$$

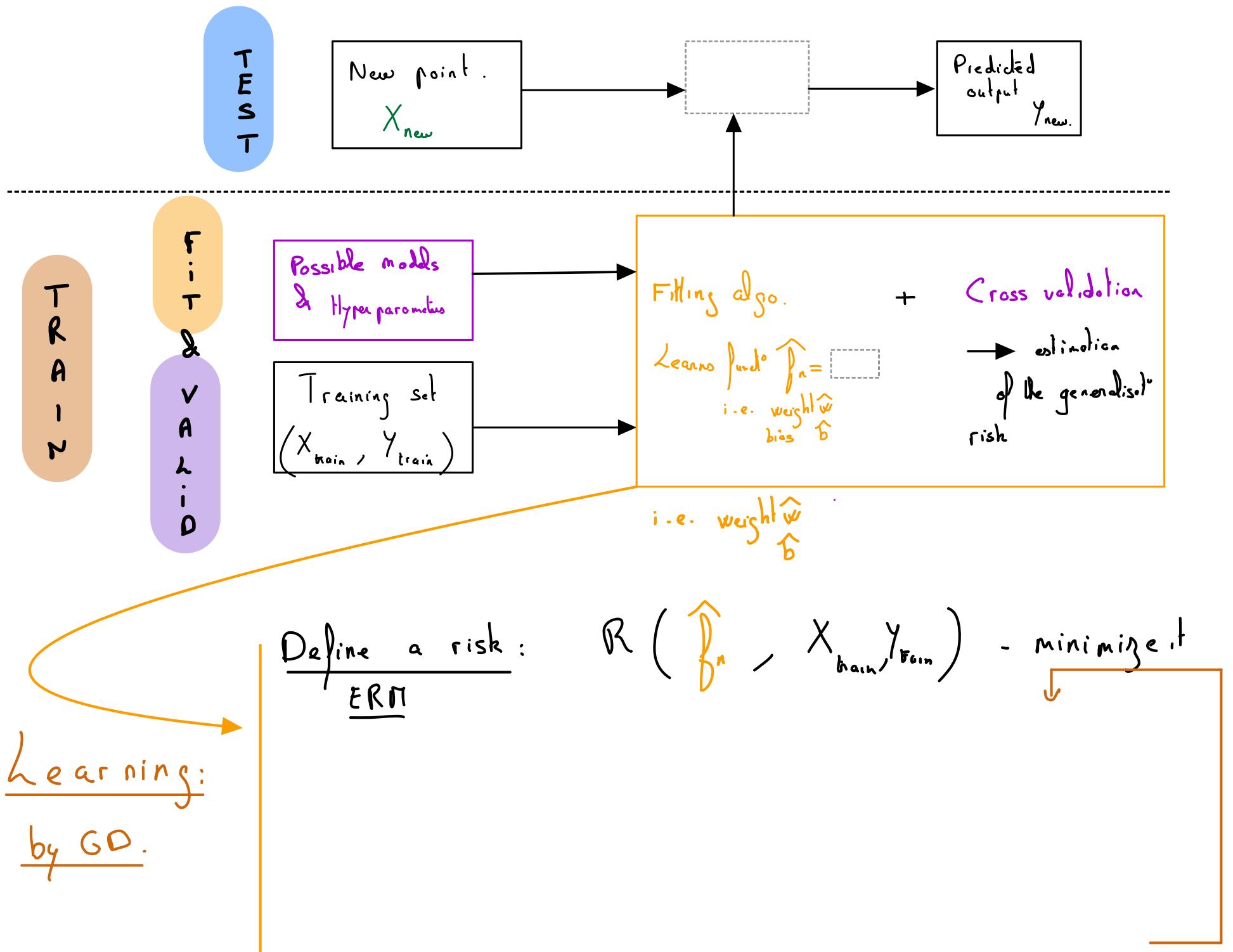
Gradient Methods. = cheapest way to update (learn) the weights iteratively to minimize the loss.

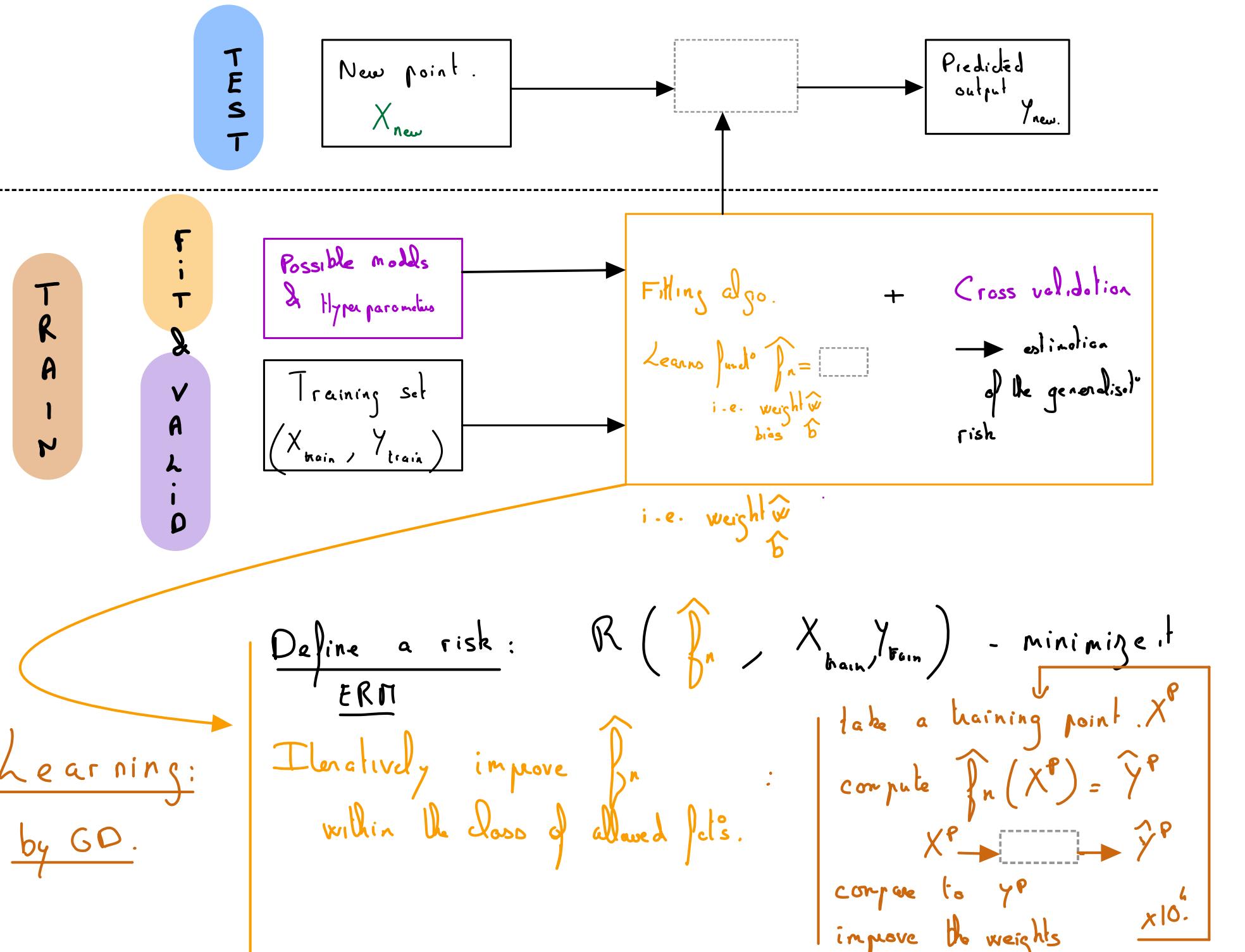
Two “miracles”:

- ① Autodifferentiation: even though the function is very complex and high dimensional, it is possible to compute gradients!
- ② Optimization seems to work ok, though the function is non convex - we do not end up in too bad local minima. (specialized optim algos: Adam, AdaDelta, etc.)



These 2 miracles make it possible to learn a good regressor/classifier with NN and to benefit from the powerful approximation properties





Conclusion

Neural networks :

- ① learn very complex non linear functions in high dimension
- ② can approximate nearly any function
- ③ can be optimized even though highly non convex.

Are thus expected to work:

- When depth or width increases.
- Thus with very large datasets
- Requiring a lot of computational power.

⇒ **Explains why they were so successful since 2010**

Conclusion

Neural networks :

- ① learn very complex non linear functions in high dimension
- ② can approximate nearly any function
- ③ can be optimized even though highly non convex.

Are thus expected to work:

- When depth or width increases.
- Thus with very large datasets
- Requiring a lot of computational power.

⇒ **Explains why they were so successful since 2010**

What my layer wants me to say: we haven't talked about many points !!

- Initialization
- Regularization

Conclusion

Neural networks :

- ① learn very complex non linear functions in high dimension
- ② can approximate nearly any function
- ③ can be optimized even though highly non convex.

Are thus expected to work:

- When depth or width increases.
- Thus with very large datasets
- Requiring a lot of computational power.

⇒ **Explains why they were so successful since 2010**

What my layer wants me to say: we haven't talked about many points !!

- Initialization
- Regularization

Next:

- Brief history
- How to use the structure?

A brief history of NN

- 1943: Neural networks
- 1957: Perceptron
- 1974-86: Backpropagation, RBM, RNN
- 1989-98: CNN, MNIST
- 2009: ImageNet
- 2012: AlexNet,
- 2014: GANss
- 2016: AlphaGo
- 2018: BERT

The Computational power made the major change (+ investment and creativity).

Directions

- When does it work?

Large or huge datasets. Structured tasks.

↗ **Images and text.**

⚠ do not try to apply DL everywhere

Each application requires a special architecture:

- Images : Convolutional Neural Network (CNN).
- Text : Recurrent Neural Networks (RNN).

Outline

1 Goals

2 The first Neural Network: the Perceptron

3 Convolutional Neural Networks

4 Applications

- Image classification
- Pose, action detection
- Object detection and segmentation
- Scene labeling - Semantic segmentation
- Object tracking - videos

We consider an **image classification task**: we want to recognize which object is on an image.

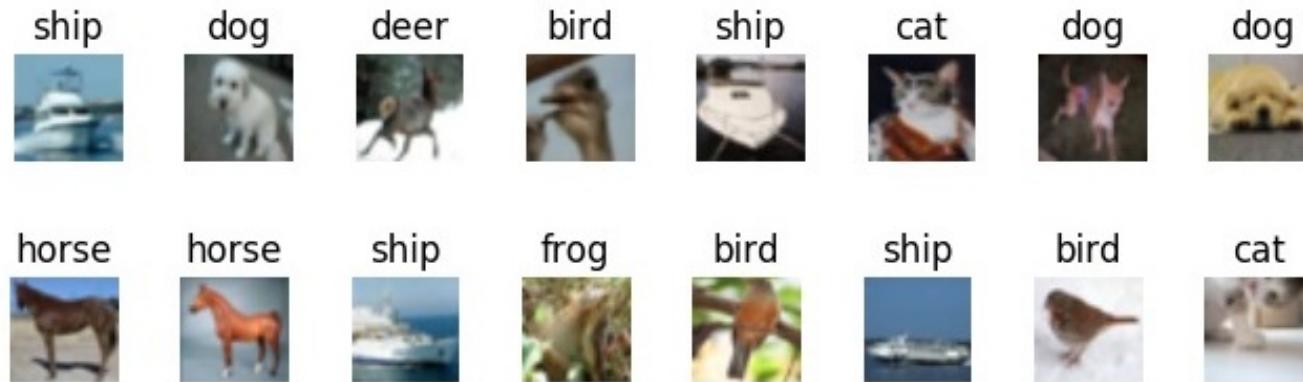


Figure: CIFAR dataset

The input is an image: about 10^3 to 10^7 pixels: these are our inputs.

Problems: The multi-layer perceptron:

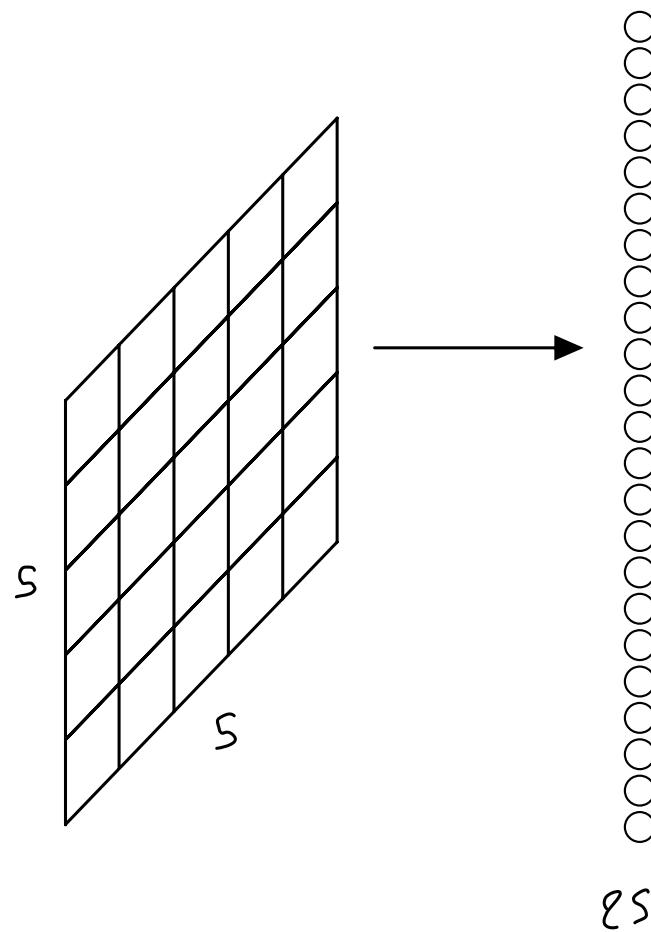
- Does not take into account local information. It would work similarly on any permutation of the pixels!
- Has, just on the first layer, as many parameters as there are inputs!

Convolutional Layer

Instead of making a linear combination of all the pixels, we consider a weighted average of a moving window of pixels, of size 3×3 , or 5×5 ...

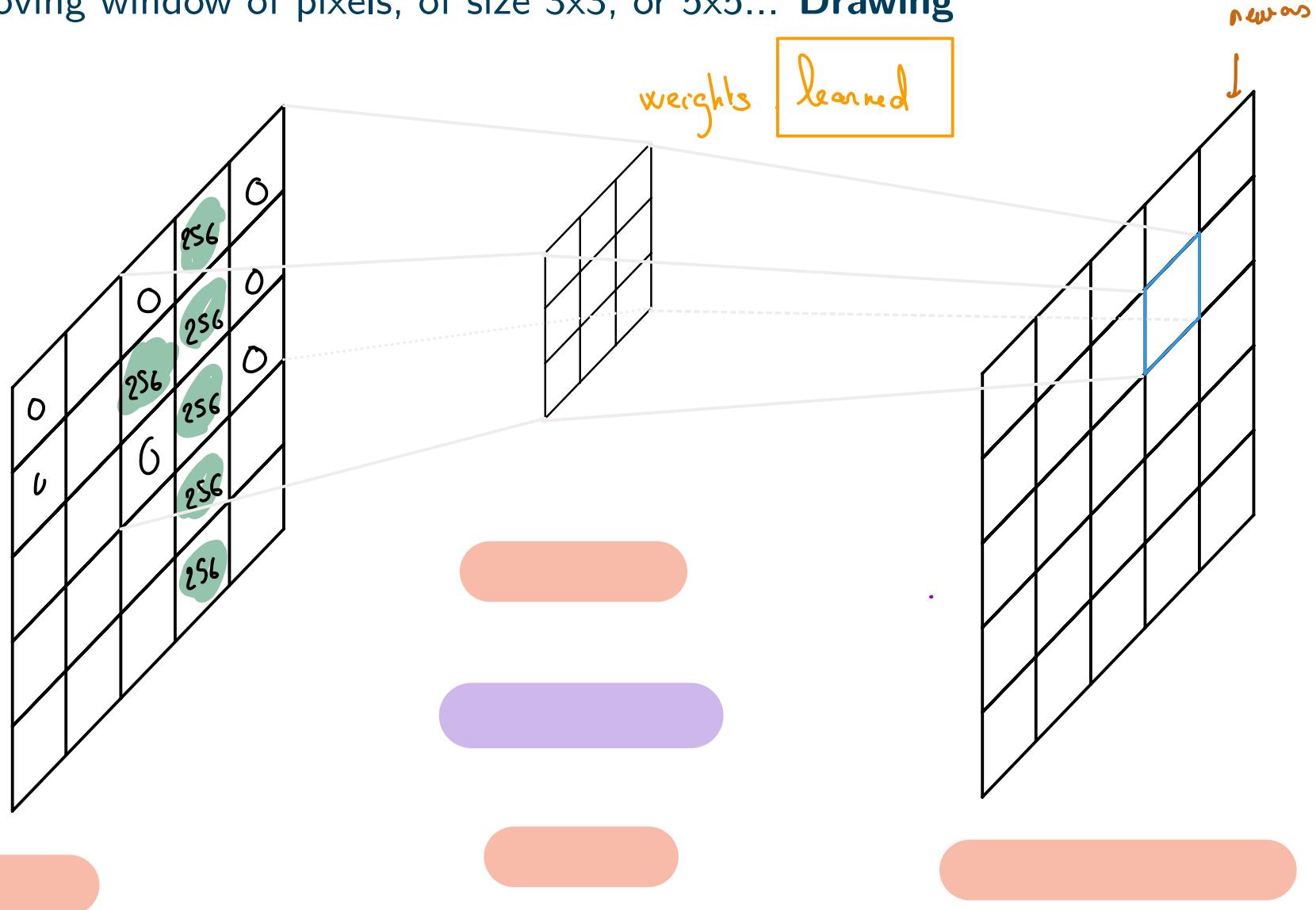
Convolutional Layer

Instead of making a linear combination of all the pixels, we consider a weighted average of a moving window of pixels, of size 3×3 , or 5×5 ...



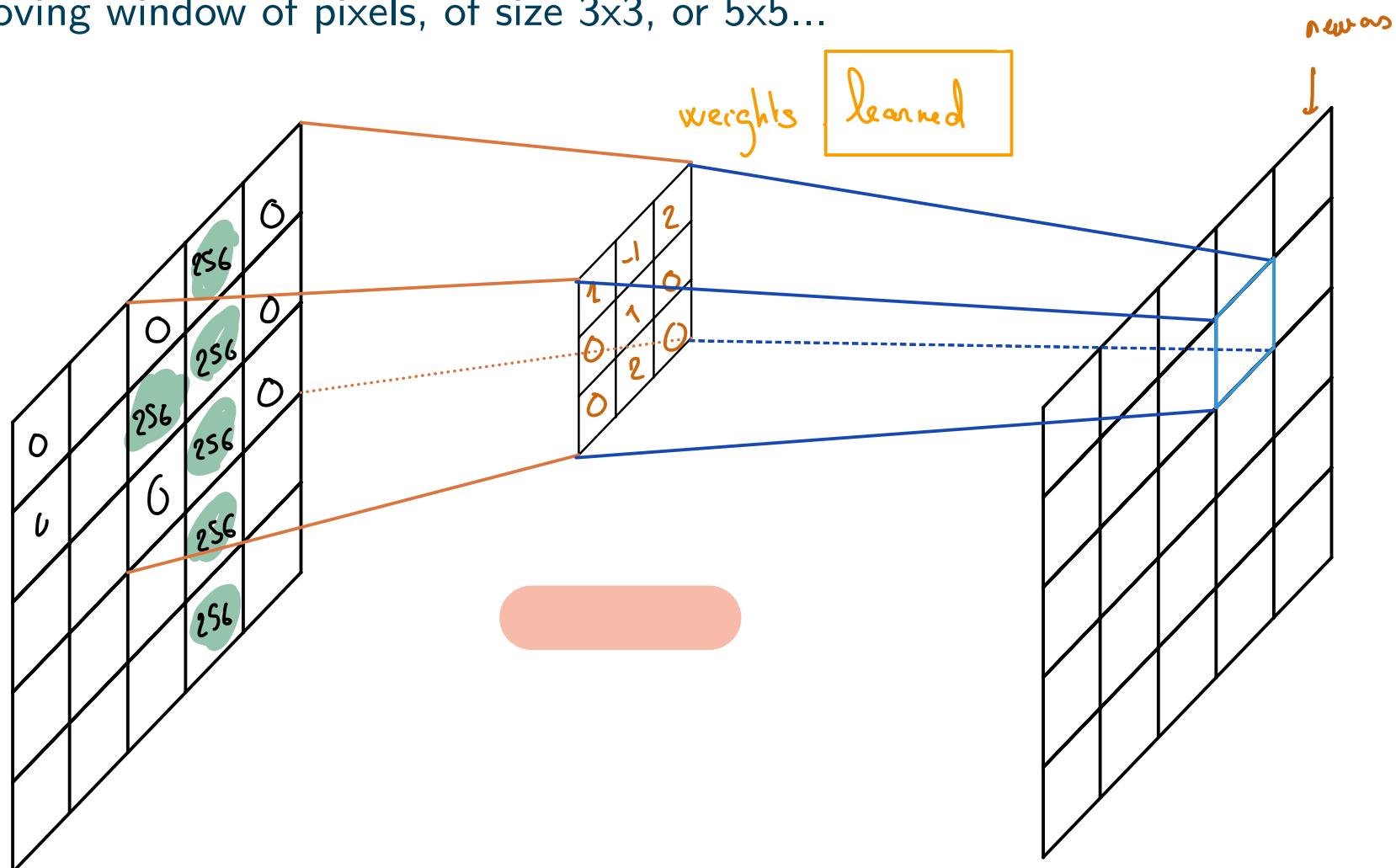
Convolutional Layer

Instead of making a linear combination of all the pixels, we consider a weighted average of a moving window of pixels, of size 3x3, or 5x5... **Drawing**



Convolutional Layer

Instead of making a linear combination of all the pixels, we consider a weighted average of a moving window of pixels, of size 3x3, or 5x5...



Convolutional Layer

0	0	0	0	0	0	0	...
0	156	155	156	158	158	...	
0	153	154	157	159	159	...	
0	149	151	155	158	159	...	
0	146	146	149	153	158	...	
0	145	143	143	148	158	...	
...

Input Channel #1 (Red)

0	0	0	0	0	0	0	...
0	167	166	167	169	169	...	
0	164	165	168	170	170	...	
0	160	162	166	169	170	...	
0	156	156	159	163	168	...	
0	155	153	153	158	168	...	
...

Input Channel #2 (Green)

0	0	0	0	0	0	0	...
0	163	162	163	165	165	...	
0	160	161	164	166	166	...	
0	156	158	162	165	166	...	
0	155	155	158	162	167	...	
0	154	152	152	157	167	...	
...

Input Channel #3 (Blue)

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

Kernel Channel #1

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

Kernel Channel #2

0	1	1
0	1	0
1	-1	1

Kernel Channel #3

308

+

-498

164

+ 1 = -25

Bias = 1

Output

-25					...
					...
					...
					...
...



Convolutional Layer

0	0	0	0	0	0	...
0	156	155	156	158	158	...
0	153	154	157	159	159	...
0	149	151	155	158	159	...
0	146	146	149	153	158	...
0	145	143	143	148	158	...
...

Input Channel #1 (Red)

0	0	0	0	0	0	...
0	167	166	167	169	169	...
0	164	165	168	170	170	...
0	160	162	166	169	170	...
0	156	156	159	163	168	...
0	155	153	153	158	168	...
...

Input Channel #2 (Green)

0	0	0	0	0	0	...
0	163	162	163	165	165	...
0	160	161	164	166	166	...
0	156	158	162	165	166	...
0	155	155	158	162	167	...
0	154	152	152	157	167	...
...

Input Channel #3 (Blue)

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

Kernel Channel #1



308

+

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

Kernel Channel #2



-498

0	1	1
0	1	0
1	-1	1

Kernel Channel #3



+

164

$$+ 1 = -25$$

Bias = 1

Output

-25					...
					...
					...
					...
...

Convolutional Layer

0	0	0	0	0	0	...
0	156	155	156	158	158	...
0	153	154	157	159	159	...
0	149	151	155	158	159	...
0	146	146	149	153	158	...
0	145	143	143	148	158	...
...

Input Channel #1 (Red)

0	0	0	0	0	0	...
0	167	166	167	169	169	...
0	164	165	168	170	170	...
0	160	162	166	169	170	...
0	156	156	159	163	168	...
0	155	153	153	158	168	...
...

Input Channel #2 (Green)

0	0	0	0	0	0	...
0	163	162	163	165	165	...
0	160	161	164	166	166	...
0	156	158	162	165	166	...
0	155	155	158	162	167	...
0	154	152	152	157	167	...
...

Input Channel #3 (Blue)

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

Kernel Channel #1



310

+

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

Kernel Channel #2



-170

0	1	1
0	1	0
1	-1	1

Kernel Channel #3



+

325

+ 1 = 466

Bias = 1

Output

-25	466				...
					...
					...
					...
...

Convolutional Layer

0	0	0	0	0	0	0	...
0	156	155	156	158	158	158	...
0	153	154	157	159	159	159	...
0	149	151	155	158	159	159	...
0	146	146	149	153	158	158	...
0	145	143	143	148	158	158	...
...

Input Channel #1 (Red)

0	0	0	0	0	0	0	...
0	167	166	167	169	169	169	...
0	164	165	168	170	170	170	...
0	160	162	166	169	170	170	...
0	156	156	159	163	168	168	...
0	155	153	153	158	168	168	...
...

Input Channel #2 (Green)

0	0	0	0	0	0	0	...
0	163	162	163	165	165	165	...
0	160	161	164	166	166	166	...
0	156	158	162	165	166	166	...
0	155	155	158	162	167	167	...
0	154	152	152	157	167	167	...
...

Input Channel #3 (Blue)

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

Kernel Channel #1



314

+

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

Kernel Channel #2



-175

0	1	1
0	1	0
1	-1	1

Kernel Channel #3



+ 326 + 1 = 466

Bias = 1

Output

-25	466	466	...
...
...
...
...

Convolutional Layer

0	0	0	0	0	0	0	...
0	156	155	156	158	158	158	...
0	153	154	157	159	159	159	...
0	149	151	155	158	159	159	...
0	146	146	149	153	158	158	...
0	145	143	143	148	158	158	...
...

Input Channel #1 (Red)

0	0	0	0	0	0	0	...
0	167	166	167	169	169	169	...
0	164	165	168	170	170	170	...
0	160	162	166	169	170	170	...
0	156	156	159	163	168	168	...
0	155	153	153	158	168	168	...
...

Input Channel #2 (Green)

0	0	0	0	0	0	0	...
0	163	162	163	165	165	165	...
0	160	161	164	166	166	166	...
0	156	158	162	165	166	166	...
0	155	155	158	162	167	167	...
0	154	152	152	157	167	167	...
...

Input Channel #3 (Blue)

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

Kernel Channel #1

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

Kernel Channel #2

0	1	1
0	1	0
1	-1	1

Kernel Channel #3

↓
318

+

↓
-173

+

329

+ 1 = 475

↑
Bias = 1

Output

-25	466	466	475	...
...
...
...
...

Convolutional Layer

0	0	0	0	0	0	...
0	156	155	156	158	158	...
0	153	154	157	159	159	...
0	149	151	155	158	159	...
0	146	146	149	153	158	...
0	145	143	143	148	158	...
...

Input Channel #1 (Red)

0	0	0	0	0	0	...
0	167	166	167	169	169	...
0	164	165	168	170	170	...
0	160	162	166	169	170	...
0	156	156	159	163	168	...
0	155	153	153	158	168	...
...

Input Channel #2 (Green)

0	0	0	0	0	0	...
0	163	162	163	165	165	...
0	160	161	164	166	166	...
0	156	158	162	165	166	...
0	155	155	158	162	167	...
0	154	152	152	157	167	...
...

Input Channel #3 (Blue)

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

Kernel Channel #1



298

+

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

Kernel Channel #2



-491

0	1	1
0	1	0
1	-1	1

Kernel Channel #3



487

+

+ 1 = 295

Bias = 1

Output

-25	466	466	475	...
295				...
				...
				...
...

Convolutional Layer

0	0	0	0	0	0	...
0	156	155	156	158	158	...
0	153	154	157	159	159	...
0	149	151	155	158	159	...
0	146	146	149	153	158	...
0	145	143	143	148	158	...
...

Input Channel #1 (Red)

0	0	0	0	0	0	...
0	167	166	167	169	169	...
0	164	165	168	170	170	...
0	160	162	166	169	170	...
0	156	156	159	163	168	...
0	155	153	153	158	168	...
...

Input Channel #2 (Green)

0	0	0	0	0	0	...
0	163	162	163	165	165	...
0	160	161	164	166	166	...
0	156	158	162	165	166	...
0	155	155	158	162	167	...
0	154	152	152	157	167	...
...

Input Channel #3 (Blue)

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

Kernel Channel #1



148

+

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

Kernel Channel #2



-8

0	1	1
0	1	0
1	-1	1

Kernel Channel #3



646

+

+ 1 = 787

Bias = 1
↑

-25	466	466	475	...
295	787			...
				...
				...
...

Output

Summary

CNN:

- allow to take spacial information into account
- reduce the number of parameters per layer. We can have deeper networks!

Summary

CNN:

- allow to take spacial information into account
- reduce the number of parameters per layer. We can have deeper networks!

LeNet 1998:

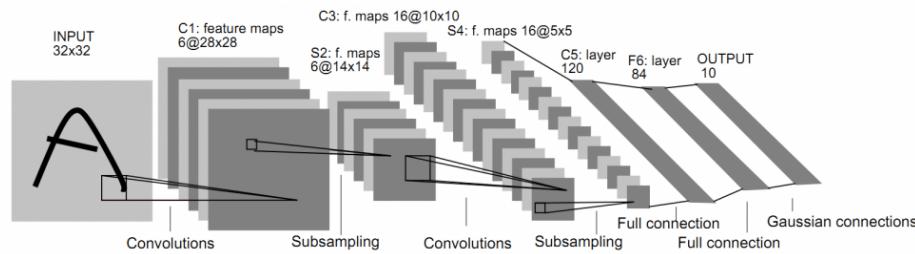
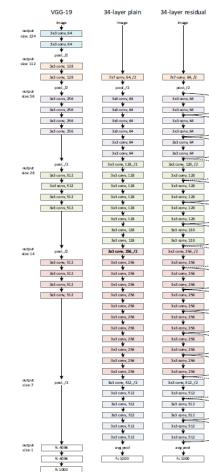


Fig. 2. Architecture of LeNet-5, a Convolutional Neural Network, here for digits recognition. Each plane is a feature map, i.e. a set of units whose weights are constrained to be identical.

Resnet 2016



8 layers, 60k parameters.

~ 40 Layers, 140M parameters

Outline

1 Goals

2 The first Neural Network: the Perceptron

3 Convolutional Neural Networks

4 Applications

- Image classification
- Pose, action detection
- Object detection and segmentation
- Scene labeling - Semantic segmentation
- Object tracking - videos

Outline

1 Goals

2 The first Neural Network: the Perceptron

3 Convolutional Neural Networks

4 Applications

- **Image classification**
- Pose, action detection
- Object detection and segmentation
- Scene labeling - Semantic segmentation
- Object tracking - videos

Image classification - CNN Tree

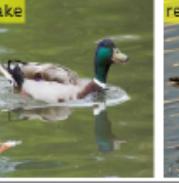
Category	Confusion Set						
tench							
indigo bunting							
red-breasted merganser							
echidna							
shopping basket							

Figure: Confusion set outputs by AlexNet softmax prediction on validation set of ILSVRC 2015.

Image classification - CNN Tree

Category	Example Validation Images					
barracouta						
church						
spaghetti squash						
espresso						
trolleybus						

Figure: Top label is given by basic AlexNet CNN while bottom one is given by CNNTree (green color corresponds to a correct prediction)

Outline

1 Goals

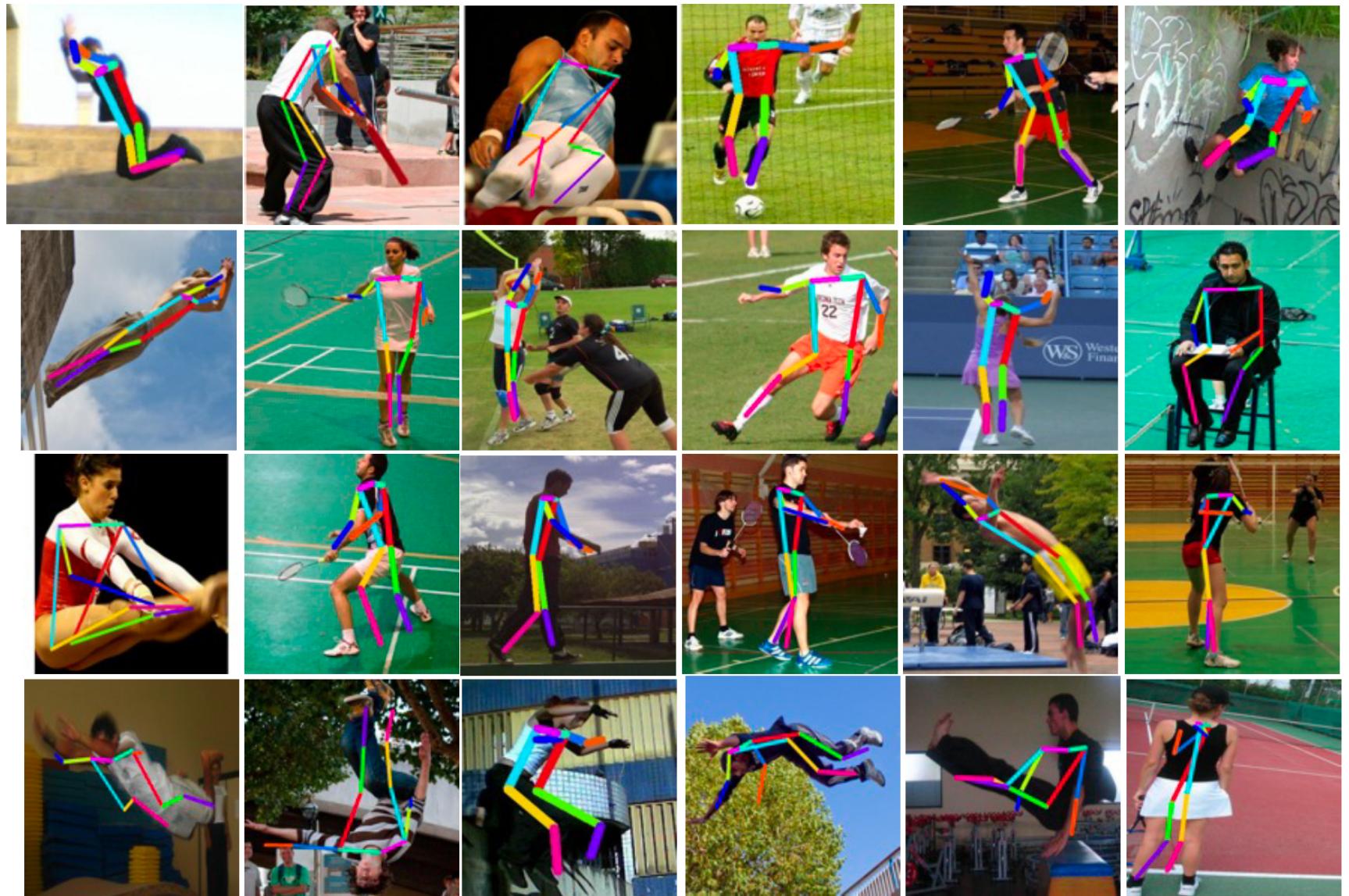
2 The first Neural Network: the Perceptron

3 Convolutional Neural Networks

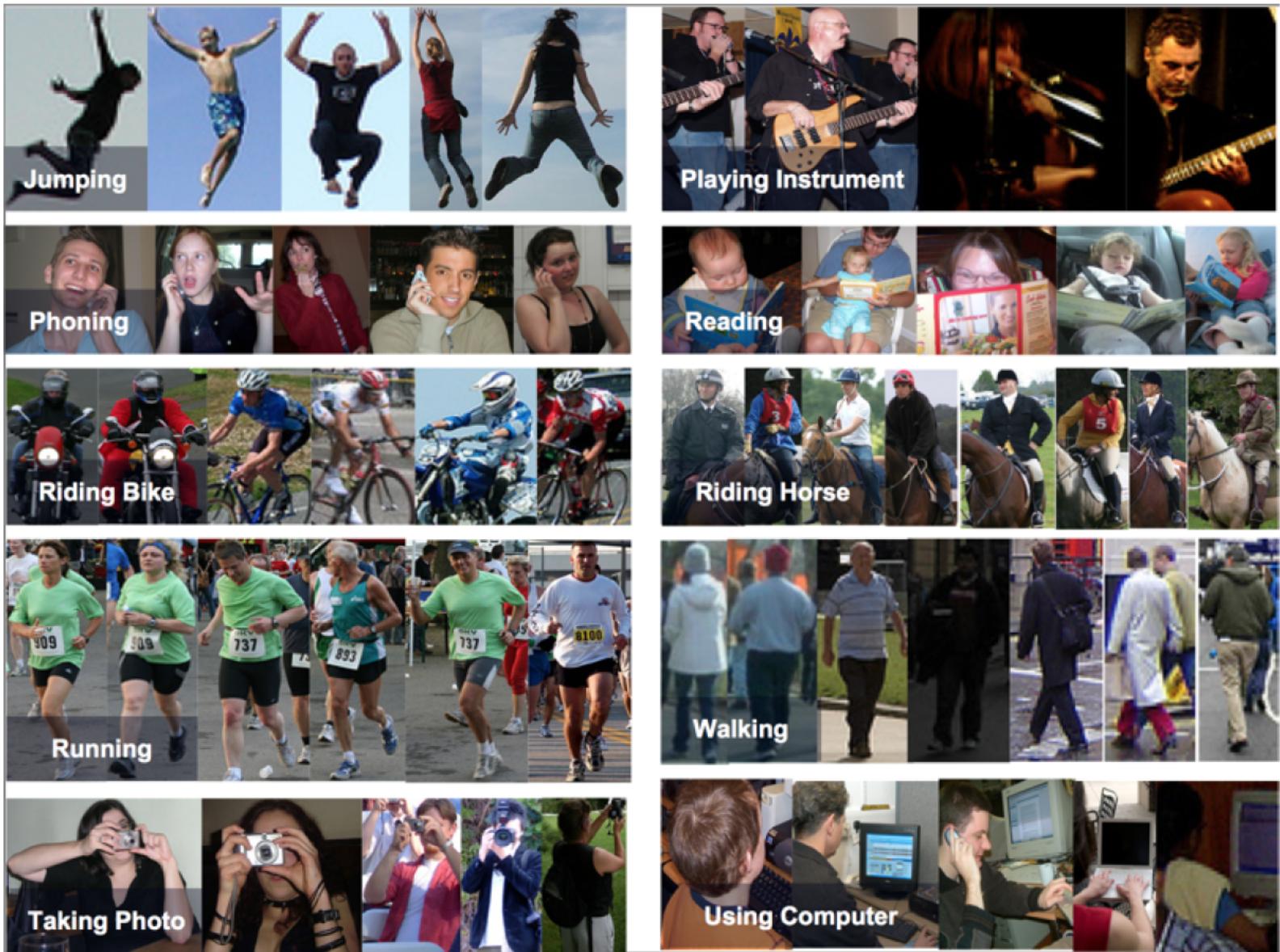
4 Applications

- Image classification
- Pose, action detection
- Object detection and segmentation
- Scene labeling - Semantic segmentation
- Object tracking - videos

Pose estimation - Deeppose



Action recognition



Outline

1 Goals

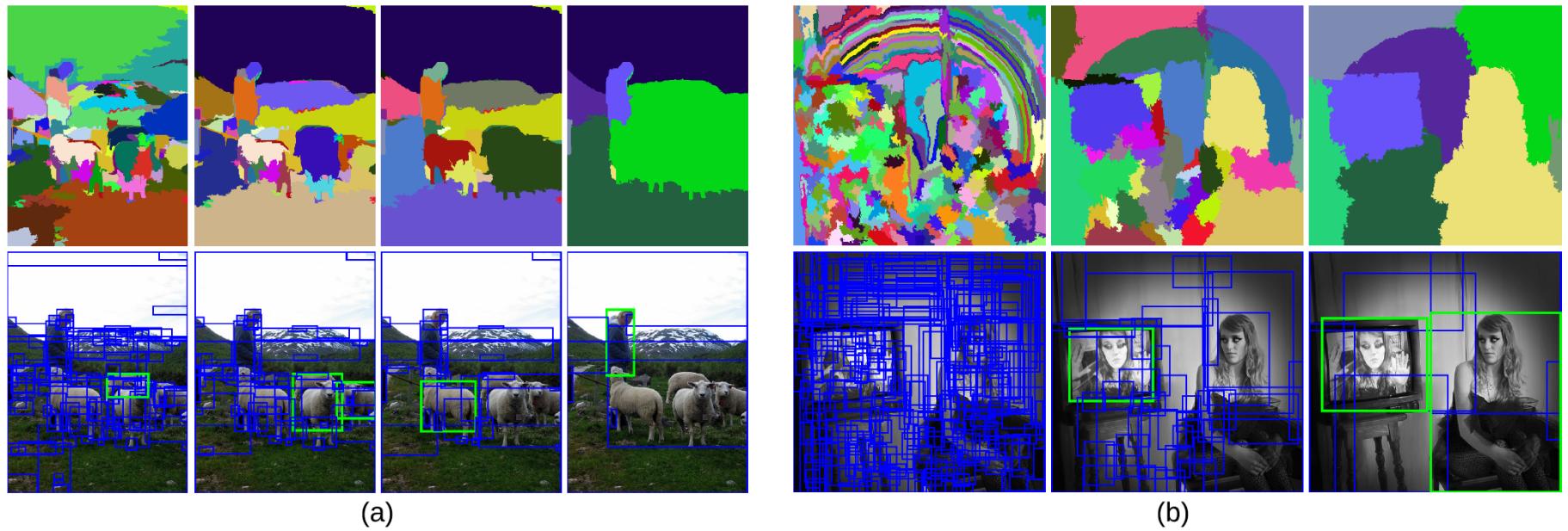
2 The first Neural Network: the Perceptron

3 Convolutional Neural Networks

4 Applications

- Image classification
- Pose, action detection
- **Object detection and segmentation**
- Scene labeling - Semantic segmentation
- Object tracking - videos

Object detection and segmentation



Object detection - YOLO, SDD

More recently, YOLO (You Only Look Once) and SSD (Single Shot Detector) allow single pipeline detection that directly predicts class labels.

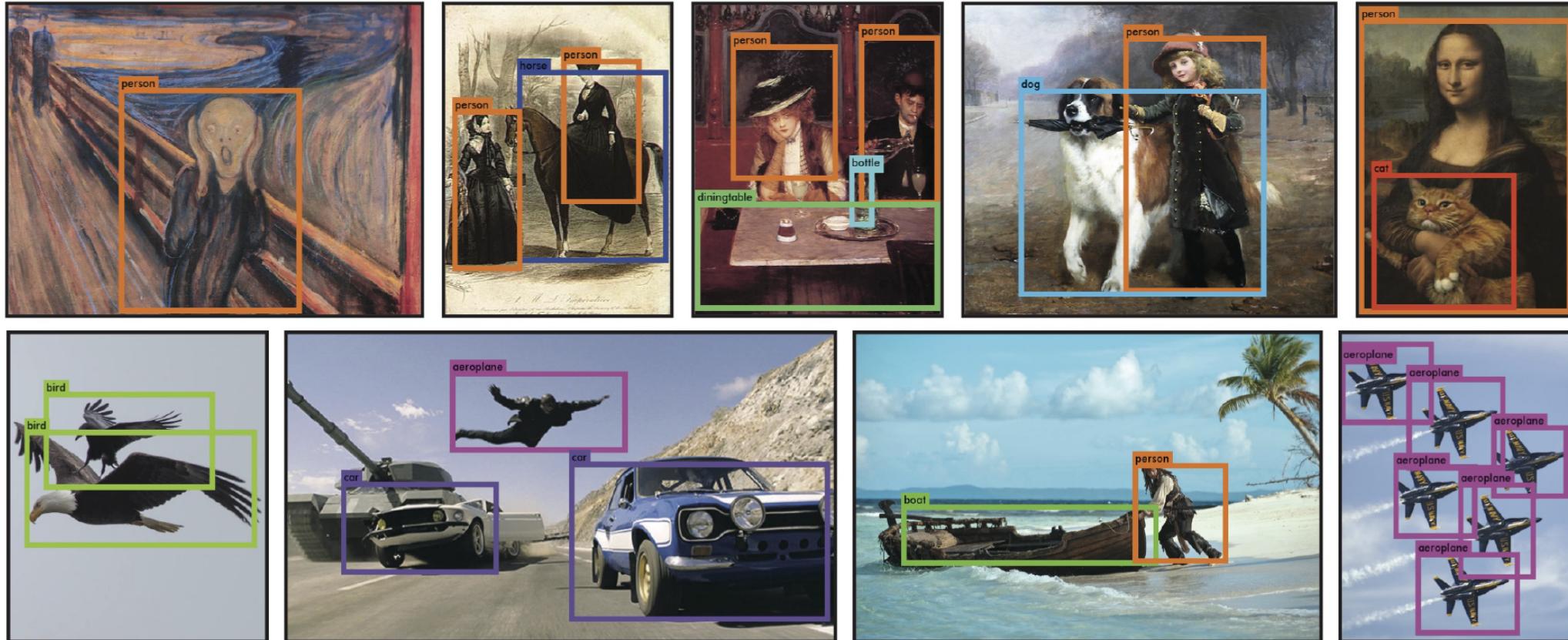
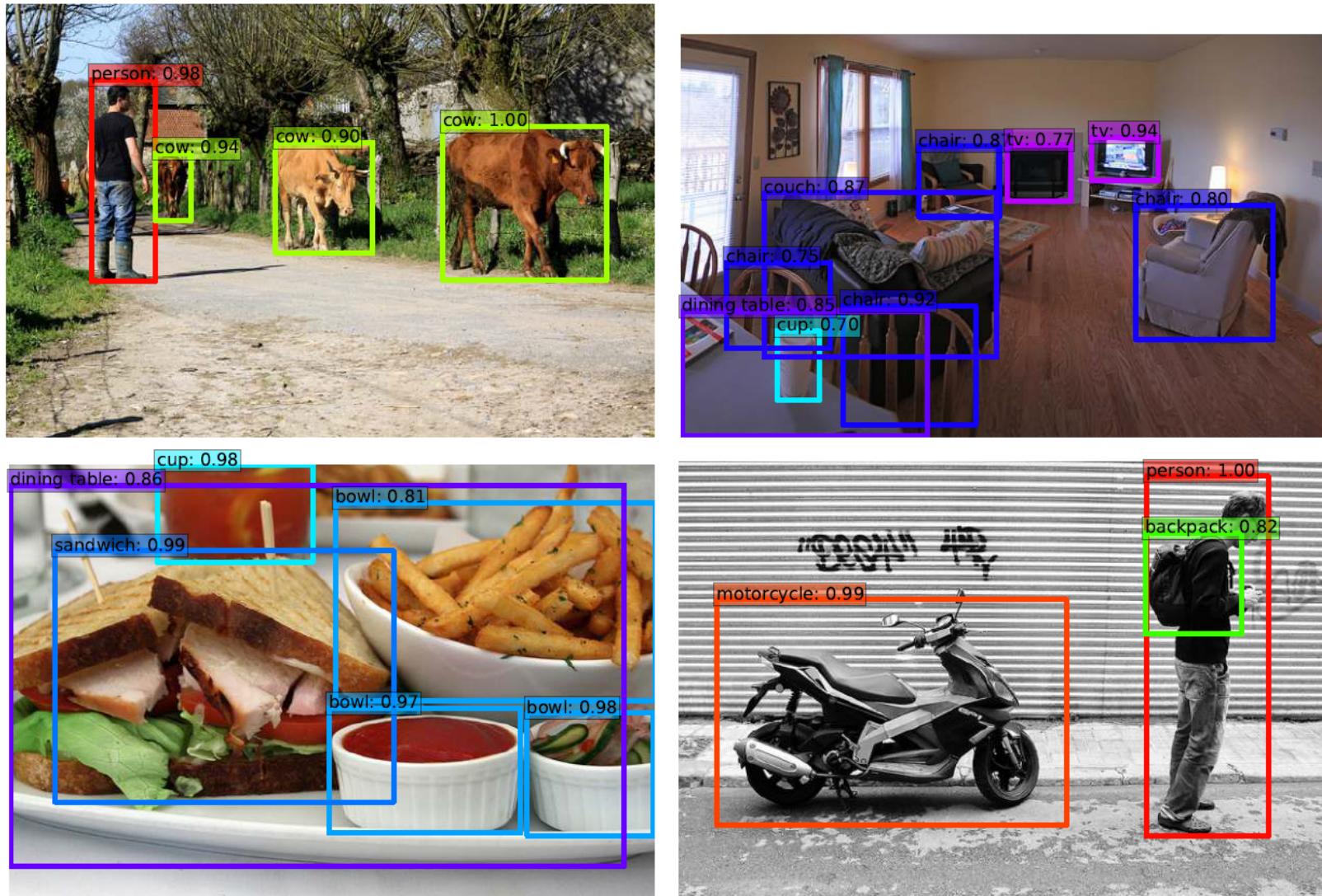


Figure: YOLO results

SSD



Outline

1 Goals

2 The first Neural Network: the Perceptron

3 Convolutional Neural Networks

4 Applications

- Image classification
- Pose, action detection
- Object detection and segmentation
- **Scene labeling - Semantic segmentation**
- Object tracking - videos

Scene labeling - DAG-RNN

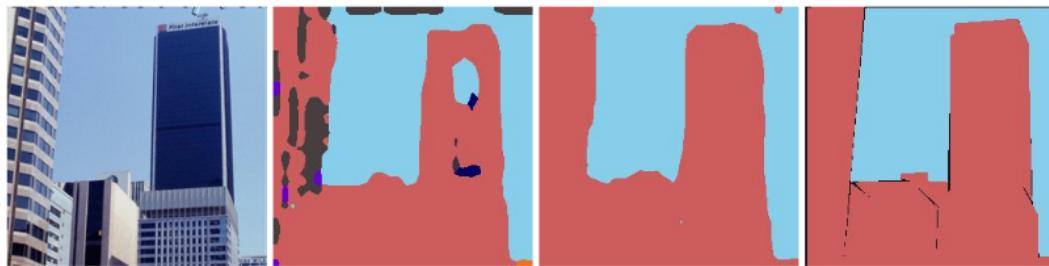


Input Image

CNN

DAG-RNN

Ground Truth



Input Image

CNN

DAG-RNN

Ground Truth

Outline

1 Goals

2 The first Neural Network: the Perceptron

3 Convolutional Neural Networks

4 Applications

- Image classification
- Pose, action detection
- Object detection and segmentation
- Scene labeling - Semantic segmentation
- Object tracking - videos

Object tracking

Object Detection in a video:

- YOLO <https://pjreddie.com/darknet/yolo/>
- YOLO - James Bond <https://www.youtube.com/watch?v=V0C3huqHrss>

Deep Learning - Many many other applications!

- ① Image and video (super-resolution, 3D, Image captioning), medical applications...
- ② Self Driving cars (scene segmentation, etc.)
- ③ In NLP (language, text), automatic translation, voice recognition, text generation, next word completion ... (Virtual assistants)
- ④ Recommendation systems...
- ⑤ For time series, complex datasets too...

But also beyond the supervised settings: using adversarial networks (GANS) for generation of images, faces, voice, etc.

Click on the person who is real.



<http://www.whichfaceisreal.com/>

Limits and problems

Limits of Deep Learning: → Wooclap!

Limits and problems

Limits of Deep Learning: → Wooclap!

- Interpretability
- Vulnerability, failures, biases

Limits and problems

Limits of Deep Learning: → Wooclap!

- Interpretability
- Vulnerability, failures, biases
- Size of the models: Applicability, Ethical, Energy & Ecological concerns (Ai has both a positive and a huge negative carbon footprint)!

Limits and problems

Limits of Deep Learning: → Wooclap!

- Interpretability
- Vulnerability, failures, biases
- Size of the models: Applicability, Ethical, Energy & Ecological concerns (Ai has both a positive and a huge negative carbon footprint)!
- Theory!

Limits and problems

Limits of Deep Learning: → Wooclap!

- Interpretability
- Vulnerability, failures, biases
- Size of the models: Applicability, Ethical, Energy & Ecological concerns (Ai has both a positive and a huge negative carbon footprint)!
- Theory!
- Ethical and practical concerns

Limits and problems

Limits of Deep Learning: → Wooclap!

- Interpretability
- Vulnerability, failures, biases
- Size of the models: Applicability, Ethical, Energy & Ecological concerns (Ai has both a positive and a huge negative carbon footprint)!
- Theory!
- Ethical and practical concerns
- Privacy

Limits and problems

Limits of Deep Learning: → Wooclap!

- Interpretability
- Vulnerability, failures, biases
- Size of the models: Applicability, Ethical, Energy & Ecological concerns (Ai has both a positive and a huge negative carbon footprint)!
- Theory!
- Ethical and practical concerns
- Privacy
- Source tracking, Ownership, IP

Limits and problems

Limits of Deep Learning: → Wooclap!

- Interpretability
- Vulnerability, failures, biases
- Size of the models: Applicability, Ethical, Energy & Ecological concerns (Ai has both a positive and a huge negative carbon footprint)!
- Theory!
- Ethical and practical concerns
- Privacy
- Source tracking, Ownership, IP
- Alignment

Deep Learning: attacks

CNN are not always “robust” to adversarial attacks!

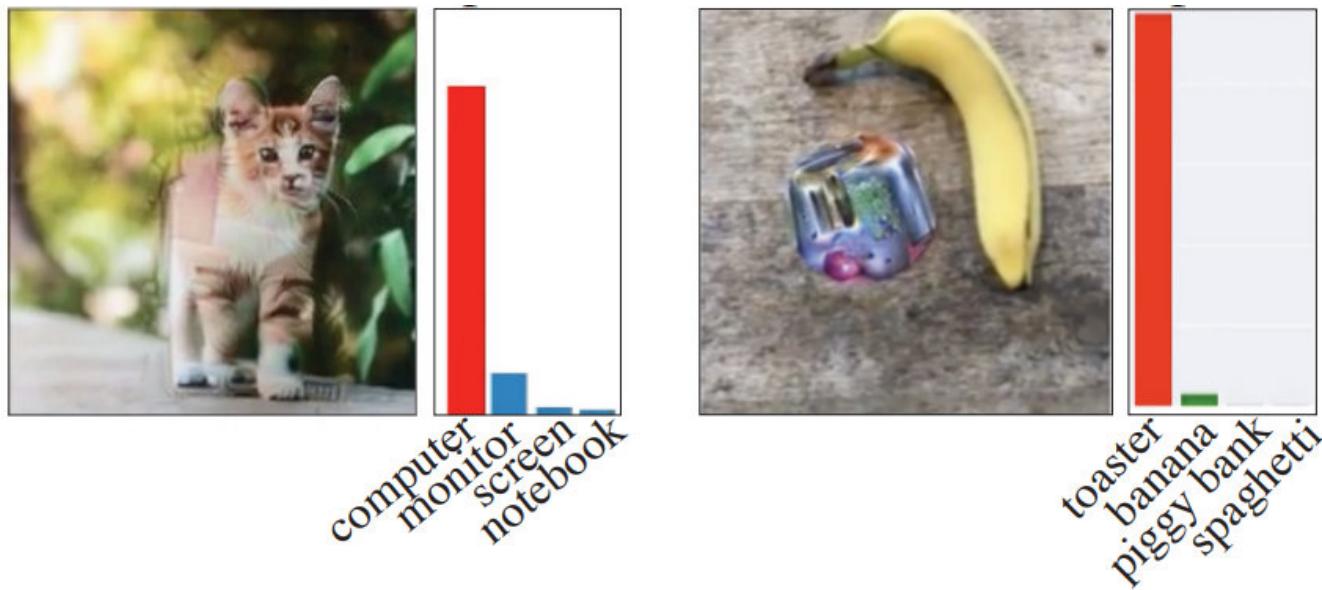


Figure: What happened here??

Deep Learning: attacks

CNN are not always “robust” to adversarial attacks!



Adding a small well chosen noise can completely fool a CNN!

↗ Can we trust deep networks on cars, medical applications, planes...?

Deep Learning: failures

① Biases in Learning sets that are transferred

- ① Chatbot becomes racist
- ② Google apologises for Photos app's racist blunder

Deep Learning: failures

- ① Biases in Learning sets that are transferred
 - ① Chatbot becomes racist
 - ② Google apologises for Photos app's racist blunder
- ② Privacy leaks:
 - ① The size of the models often enable memorization of the training sets.
 - ② Need for guarantees and regulation

Deep Learning: failures

- ① Biases in Learning sets that are transferred
 - ① Chatbot becomes racist
 - ② Google apologises for Photos app's racist blunder
- ② Privacy leaks:
 - ① The size of the models often enable memorization of the training sets.
 - ② Need for guarantees and regulation
- ③ False information/prediction, hallucinations...

Deep Learning: failures

- ① Biases in Learning sets that are transferred
 - ① Chatbot becomes racist
 - ② Google apologises for Photos app's racist blunder
- ② Privacy leaks:
 - ① The size of the models often enable memorization of the training sets.
 - ② Need for guarantees and regulation
- ③ False information/prediction, hallucinations...
- ④ Overall, lack of control on the model behavior

PIXELS • VIE PRIVÉE

Nouvelle plainte contre OpenAI pour infraction au RGPD

ChatGPT enfreint le Règlement général sur la protection des données, estime l'association autrichienne None of Your Business (NOYB), car il peut donner de fausses informations concernant un individu, sans possibilité de les faire corriger ou supprimer.

Le Monde avec AFP
Publié hier à 10h44 - ⏲ Lecture 1 min.

Figure: Link

Alignment



FAST COMPANY

CO.DESIGN TECH WORK LIFE NEWS IMPACT PODCASTS VIDEO RECOMMENDER INNOVATION FESTI

11-06-17 | PLATFORM WARS

Netflix CEO Reed Hastings: Sleep Is Our Competition

For Netflix, the battle for domination goes far beyond which TV remote to pick up.



CEO Reed Hastings [Photo: Mondileinchen/Wikimedia Commons]

Deep Learning: societal risk and questionable applications

Are some applications just bad (they all already exist :()?

- Deep Fakes: insert someone in a video.
- Underground exploration: finding new fossil resources.
- Voice imitation: what if I cannot check who is calling me.
- Military applications + automatic weapons.
- Mass surveillance.

The likelihood of each of those applications being largely deployed increases every year.

Deep Learning: societal risk and questionable applications

Are some applications just bad (they all already exist :()?

- Deep Fakes: insert someone in a video.
- Underground exploration: finding new fossil resources.
- Voice imitation: what if I cannot check who is calling me.
- Military applications + automatic weapons.
- Mass surveillance.

The likelihood of each of those applications being largely deployed increases every year.

Deep Learning: questionable applications

A challenging problem: [chain of responsibility](#). The same tools are used for positive and negative applications, e.g.:

- mass surveillance
- tumor recognition

are both based on image recognition, and the same ML techniques.

Deep Learning: questionable applications

A challenging problem: [chain of responsibility](#). The same tools are used for positive and negative applications, e.g.:

- mass surveillance
- tumor recognition

are both based on image recognition, and the same ML techniques.

[Good news!](#) the research community is very aware of the situation.

Some references:

- [Asilomar AI Principles](#) [Ai principles](#)
- [AI Act](#)
- [European guidelines](#)
- [AI for Good](#)



We need you to be aware and active on those challenges!!

Conclusion: Deep Learning in one slide

- **How does it work:**

- ▶ Automatically learn representations of observations
- ▶ Learn highly non-linear models.

- **What does it require:**

- ▶ Large datasets with structure
- ▶ Computational power

- **Why now:**

- ▶ Combination of the 2 points above
- ▶ investment!

- **Some Applications**

- ▶ Image classification; object / face recognition
- ▶ Self driving cars
- ▶ Automatic Translation, Information extraction
- ▶ Caption Generation
- ▶ Ads, recommendation systems, etc.

Goals: Understanding core concepts of Deep Learning.

- When and how it works on paper (data, models, architecture)
- How to implement a simple neural network with Python.
- Overview of some of the main applications and challenges.

You cannot vote anymore

Which of the following things need to be chosen (hyperparameters of the NN)?

1 neuron 13% 5 13%
 2 activation functions 95% 37 ✓ 2 activation functions 95% 37 ✓
 3 number of neurons per layer 92% 36 ✓ 3 number of neurons per layer 92% 36 ✓
 4 weights 18% 7 18%
 5 input size 46% 18 46%
 6 number of layers 100% 39 ✓ 6 number of layers 100% 39 ✓
 7 optimization algorithm 38% 15 ✓ 7 optimization algorithm 38% 15 ✓
 8 biases 8% 3 8%

wooclap Questions 1 / 10 Messages 🔒 100 % 🔍 Exit
 california.ipynb Tout afficher

You cannot vote anymore

Which of the following things need to be learned ?

1 neuron 9% 3 9%
 2 activation functions 3% 1 3%
 3 number of neurons per layer 3% 1 3%
 4 weights 97% 32 ✓ 4 weights 97% 32 ✓
 5 input size 9% 3 9%
 6 number of layers 0% 0 0%
 7 optimization algorithm 9% 3 9%
 8 biases 76% 25 ✓ 8 biases 76% 25 ✓

wooclap Questions 2 / 10 Messages 🔒 100 % 🔍 Exit
 california.ipynb Tout afficher

You cannot vote anymore

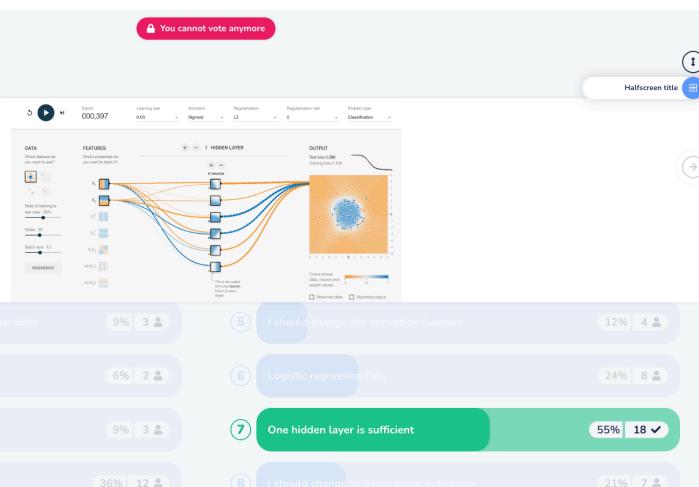
What happens here?

1 Logistic Regression works (data is linearly separable) 58% 15 ✓ 1 Logistic Regression works (data is linearly separable) 58% 15 ✓
 2 I should use more neurons 0% 0 0%
 3 The network overfits 4% 1 4%
 4 I should add more layers 0% 0 0%
 5 I should change the activation function 4% 1 4%
 6 Logistic regression fails 0% 0 0%
 7 One hidden layer is sufficient 31% 8 31%
 8 I should change to a non linear activation 4% 1 4%

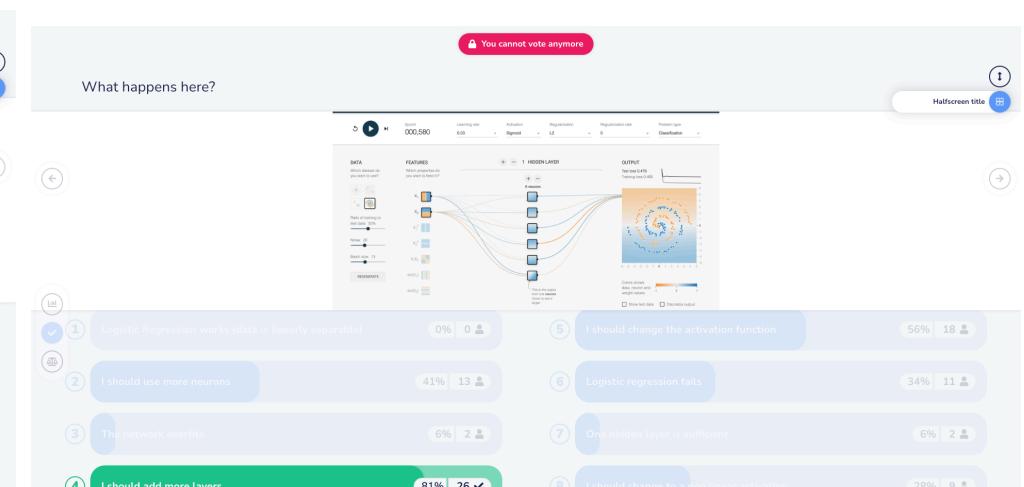
1 Logistic Regression works (data is linearly separable) 13% 3 13%
 2 I should use more neurons 46% 11 46%
 3 The network overfits 8% 2 8%
 4 I should add more layers 54% 13 54%
 5 I should change the activation function 21% 5 21%
 6 Logistic regression fails 75% 18 75%
 7 One hidden layer is sufficient 4% 1 4%
 8 I should change to a non linear activation 25% 6 25%

wooclap Questions 3 / 10 Messages 🔒 100 % 🔍 Exit
 california.ipynb Tout afficher

What happens here?



What happens here?



What happens here?

You cannot vote anymore

This screenshot shows a Wooclap poll titled "What happens here?". At the top, there's a message "You cannot vote anymore". Below it is a neural network diagram with 4 hidden layers. The poll has 8 options, each with a progress bar:

- ① Logistic Regression works (data is linearly separable) - 0% (0 people)
- ② I should use more neurons - 12% (3 people)
- ③ The network overfits - 8% (2 people)
- ④ I should add more layers - 15% (4 people)
- ⑤ I should change the activation function - 85% (22 people) [highlighted in green]
- ⑥ Logistic regression fails - 12% (3 people)
- ⑦ One hidden layer is sufficient - 0% (0 people)
- ⑧ I should change to a non-linear activation - 0% (0 people)

At the bottom, the Wooclap interface shows 7 questions, 100% messages, and 26 / 45 participants.

What happens here?

You cannot vote anymore

This screenshot shows a Wooclap poll titled "What happens here?". At the top, there's a message "You cannot vote anymore". Below it is a neural network diagram with 4 hidden layers. The poll has 8 options, each with a progress bar:

- ① Logistic Regression works (data is linearly separable) - 0% (0 people)
- ② I should use more neurons - 4% (1 person)
- ③ The network overfits - 83% (19 people) [highlighted in green]
- ④ I should add more layers - 0% (0 people)
- ⑤ I should change the activation function - 4% (1 person)
- ⑥ Logistic regression fails - 0% (0 people)
- ⑦ One hidden layer is sufficient - 0% (0 people)
- ⑧ I should change to a non-linear activation - 9% (2 people)

At the bottom, the Wooclap interface shows 7 questions, 100% messages, and 23 / 45 participants.

What happens here?

You cannot vote anymore

This screenshot shows a Wooclap poll titled "What happens here?". At the top, there's a message "You cannot vote anymore". Below it is a neural network diagram with 6 hidden layers. The poll has 8 options, each with a progress bar:

- ① Logistic Regression works (data is linearly separable) - 0% (0 people)
- ② I should use more neurons - 0% (0 people)
- ③ The network overfits - 0% (0 people)
- ④ I should add more layers - 0% (0 people)
- ⑤ I should change the activation function - 0% (0 people)
- ⑥ Logistic regression fails - 0% (0 people)
- ⑦ It works! - 100% (21 people) [highlighted in green]
- ⑧ I should change to a non-linear activation - 0% (0 people)

At the bottom, the Wooclap interface shows 10 questions, 100% messages, and 21 / 45 participants.

