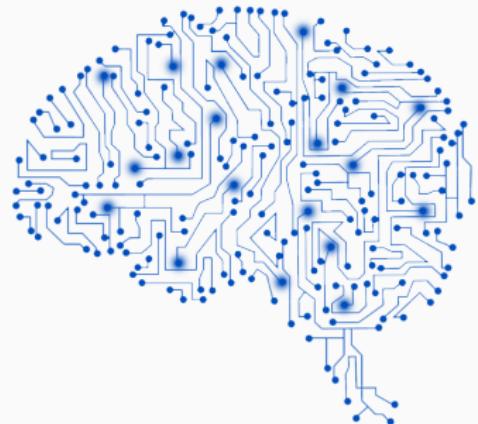


Apprentissage pour l'image Machine learning for image processing

Course I – Introduction to Artificial Neural Networks: Multiclass logistic classification

Emile Pierret

Jeudi 20 mars 2025



Credits

Most of the slides from **Charles Deledalle's** course “UCSD ECE285 Machine learning for image processing” (30 × 50 minutes course)



www.charles-deledalle.fr/

<https://www.charles-deledalle.fr/pages/teaching.php#learning>

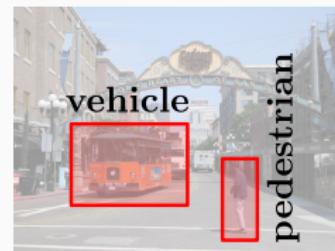
Programme du cours d'aujourd'hui

- ① Introduction aux capacités des réseaux de neurones
- ② Introductions aux "datasets"
- ③ Introduction aux réseaux de neurones
- ④ Introduction aux coûts
- ⑤ Etude de la regression logistique

Computer Vision and Machine Learning



Image



Symbols

Computer vision – Artificial Intelligence – Machine Learning

Definition (The British Machine Vision Association)

Computer vision (CV) consiste en l'extraction, analyse et compréhension d'information utile à partir d'une image ou d'un ensemble d'image.



CV is a subfield of Artificial Intelligence.

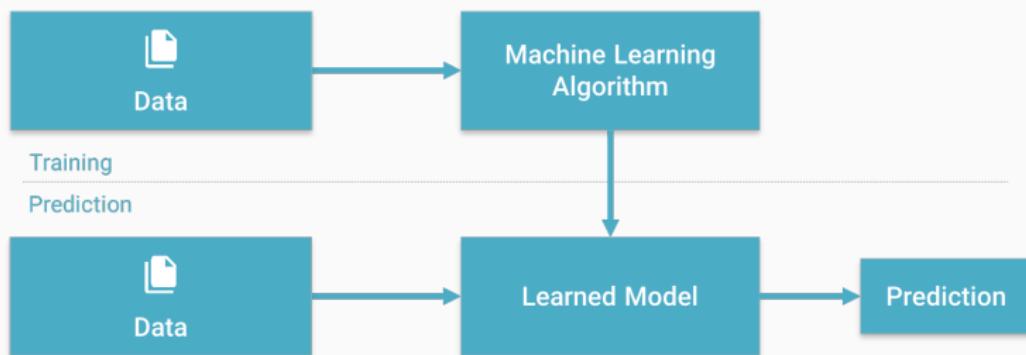
Definition (Oxford dictionary)

Artificial Intelligence, noun: La théorie et le développement de système informatique capable d'effectuer des tâches qui requièrent normalement l'intelligence humaine comme la perception visuelle, la reconnaissance de parole, la prise de décision ou la traduction.

Definition

Machine Learning, *noun*: Type d'intelligence artificielle qui permet à un système informatique d'apprendre une tâche sans être explicitement programmé pour.

L'objectif est de faire des prédictions à partir de données :



(Source: Lucas Masuch)

Deep learning: Academic actors

- Popularized by Hinton in 2006 with Restricted Boltzmann Machines



Geoffrey Hinton: University of Toronto & Google

- Developed by different actors:



Yann LeCun: New York University & Facebook



Andrew Ng: Stanford & Baidu



Yoshua Bengio: University of Montreal



Jürgen Schmidhuber: Swiss AI Lab & NNAISENSE

and many others...

Deep learning: Academic actors

- Popularized by Hinton in 2006 with Restricted Boltzmann Machines



Geoffrey Hinton: University of Toronto & Google

- Developed by different actors:



Yann LeCun: New York University & Facebook



Andrew Ng: Stanford & Baidu



Yoshua Bengio: University of Montreal



Jürgen Schmidhuber: Swiss AI Lab & NNAISENSE

and many others...

- Yoshua Bengio, Geoffrey Hinton, and Yann LeCun recipients of the 2018 ACM A.M. Turing Award for conceptual and engineering breakthroughs that have made deep neural networks a critical component of computing.

Actors and applications

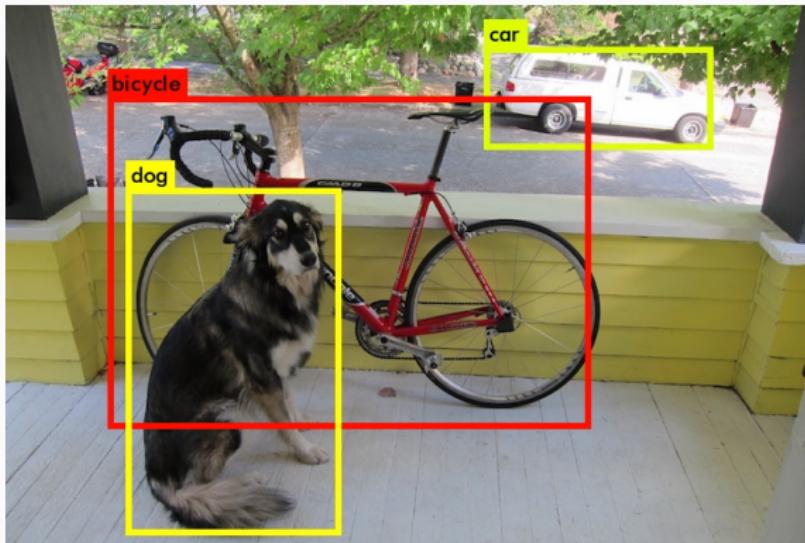
- Very active technology adopted by big actors



- Success story for many different academic problems
 - Image processing
 - Computer vision
 - Speech recognition
 - Natural language processing
 - Translation
 - etc
- Today all industries wonder if Machine learning can improve their process.

Computer vision – Object detection

Computer vision – Object detection



(Source: Joseph Redmon)

Goal: to detect instances of objects of a certain class (such as human).

Computer vision – Image segmentation

Computer vision – Image segmentation



(Source: Abhijit Kundu)

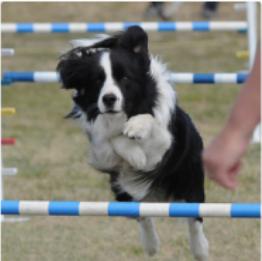
Goal: to partition an image into multiple segments such that pixels in a same segment share certain characteristics (color, texture or semantic).

Computer vision – Image captioning

Computer vision – Image captioning



"girl in pink dress is jumping in air."



"black and white dog jumps over bar."



"young girl in pink shirt is swinging on swing."



"man in blue wetsuit is surfing on wave."



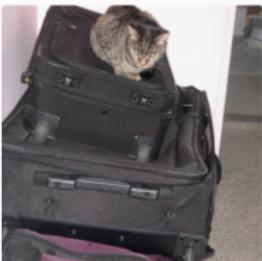
"little girl is eating piece of cake."



"baseball player is throwing ball in game."



"woman is holding bunch of bananas."



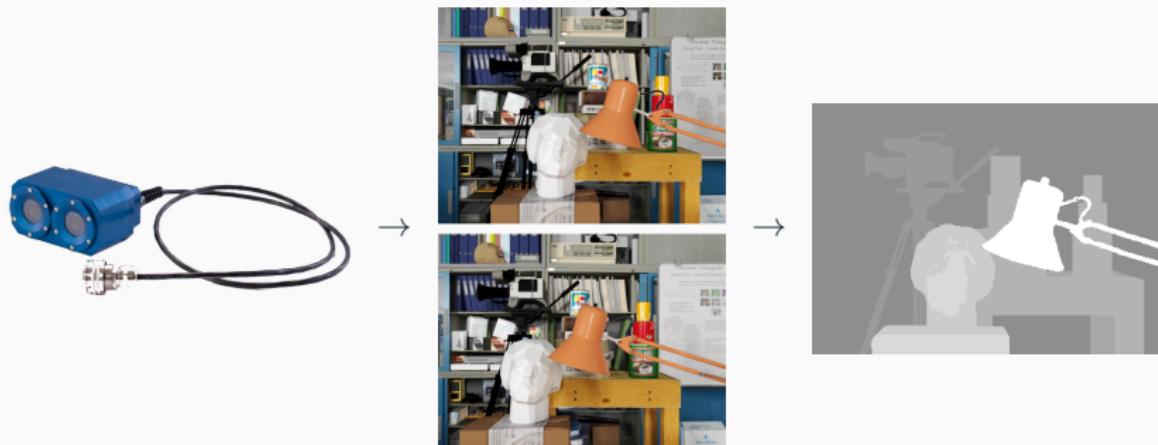
"black cat is sitting on top of suitcase."

(Karpathy, Fei-Fei, CVPR, 2015)

Goal: to write a sentence that describes what is happening.

Computer vision – Depth estimation

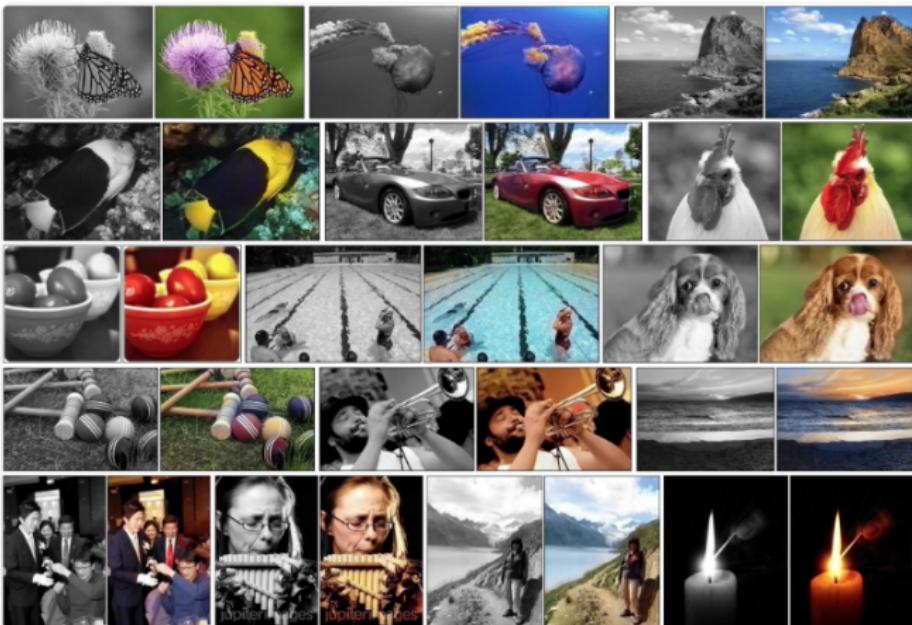
Computer vision – Depth estimation



(*Stereo-vision: from two images acquired with different views.*)

Goal: to estimate a depth map from one, two or several frames.

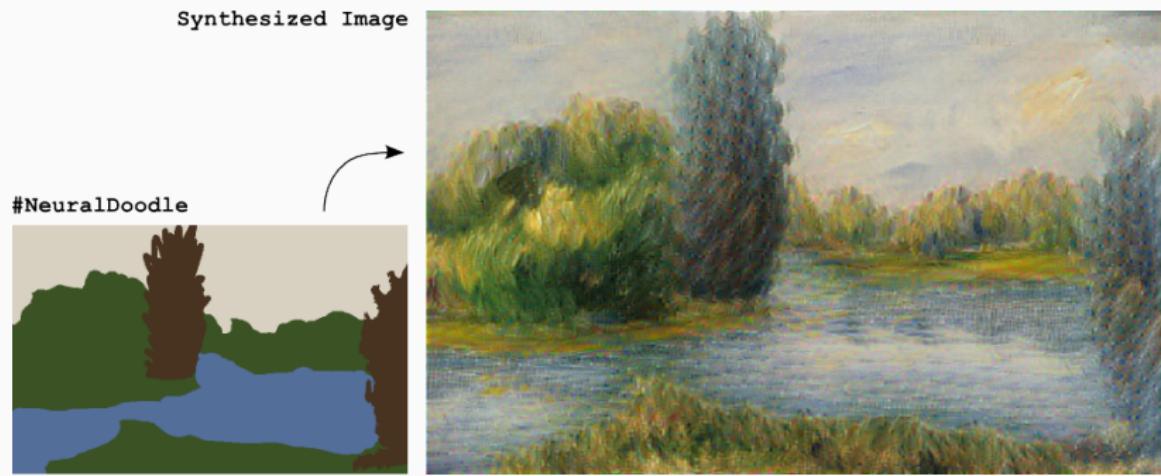
Image colorization



(Source: Richard Zhang, Phillip Isola and Alexei A. Efros, 2016)

Goal: to add color to grayscale photographs.

Image stylization



(Source: Neural Doodle, Champandard, 2016)

Goal: to create stylized images from rough sketches.

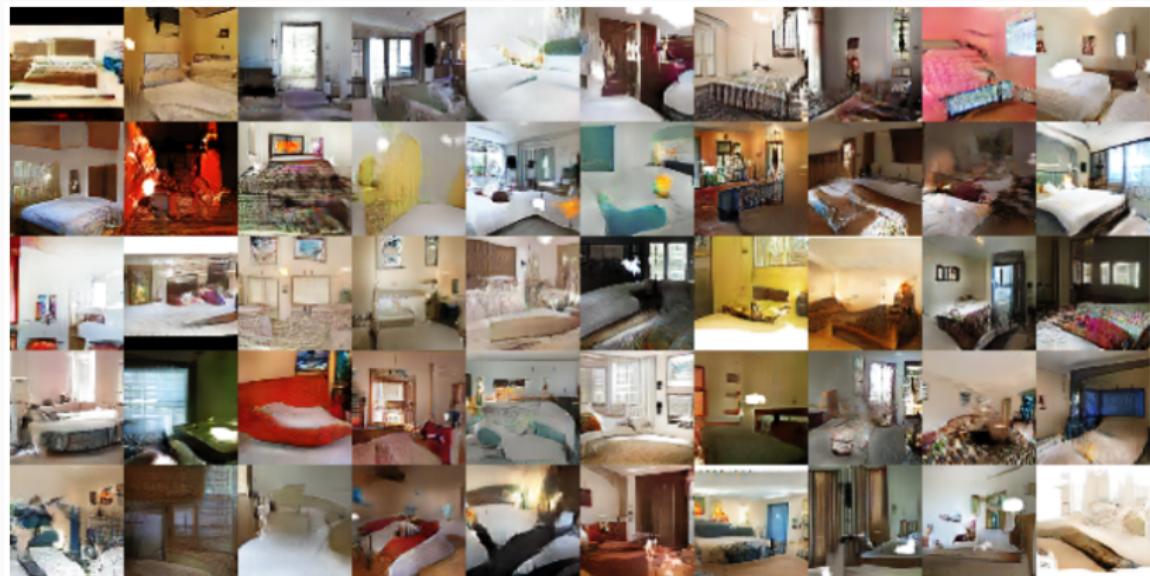
Style transfer



(Source: Gatys, Ecker and Bethge, 2015)

Goal: transfer the style of an image into another one.

Image generation



Generated images of bedrooms (Source: Alec Radford, Luke Metz, Soumith Chintala, 2015)

Goal: to automatically create realistic pictures of a given category.

Success stories

Face generation (Style GAN, (Karras et al., 2018) (NVIDIA)):

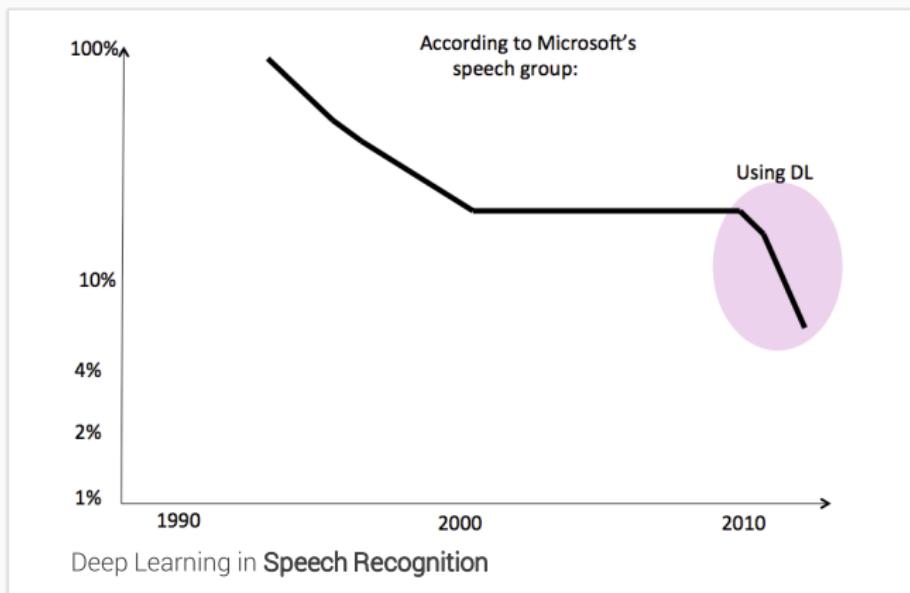
These people do not exist.



Deep learning – Success stories

Success stories

Speech recognition



Same improvements in other fields: automatic translation, Go, ...

Deep learning – Success stories

Success stories

AlexNet wins the ImageNet classification challenge (Krizhevsky, 2012)

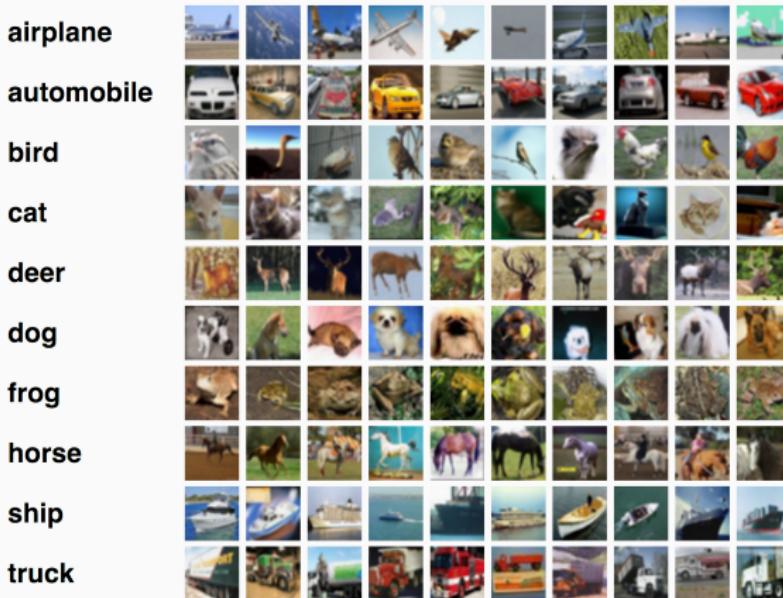


ImageNet: The "computer vision World Cup"

Since, Deep Learning has been the standard and achieved remarkable results.

Our goal: Image classification

Computer vision – Image classification



Goal: to assign a given image into one of the predefined classes.

Plan of the course: How to learn from examples

Plan of the course: How to learn from examples

3 main ingredients

- ① Training set / examples (Données d'entraînement):

$$\{\boldsymbol{x}_1, \boldsymbol{x}_2, \dots, \boldsymbol{x}_N\}$$

Plan of the course: How to learn from examples

3 main ingredients

- ① Training set / examples (Données d'entraînement):

$$\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$$

- ② Machine or model (Modèle):

$$\mathbf{x} \rightarrow \underbrace{f(\mathbf{x}; \theta)}_{\text{function / algorithm}} \rightarrow \underbrace{\mathbf{y}}_{\text{prediction}}$$

θ : parameters of the model

Plan of the course: How to learn from examples

3 main ingredients

- ① Training set / examples (Données d'entraînement):

$$\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$$

- ② Machine or model (Modèle):

$$\mathbf{x} \rightarrow \underbrace{f(\mathbf{x}; \theta)}_{\text{function / algorithm}} \rightarrow \underbrace{\mathbf{y}}_{\text{prediction}}$$

θ : parameters of the model

- ③ Loss, cost, objective function / energy (Fonction de coût):

$$\operatorname{argmin}_{\theta} E(\theta; \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$$

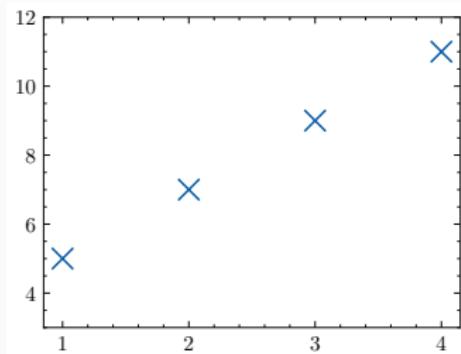
Exemple de la regression linéaire:

- Données d'entraînement: $(x_i, y_i)_{1 \leq i \leq N}$.

Machine learning – Learning from examples

Exemple de la regression linéaire:

- Données d'entraînement: $(x_i, y_i)_{1 \leq i \leq N}$.

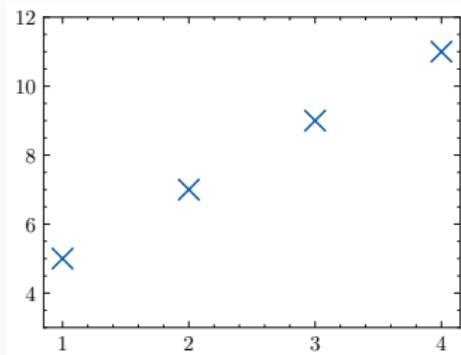


Machine learning – Learning from examples

Exemple de la regression linéaire:

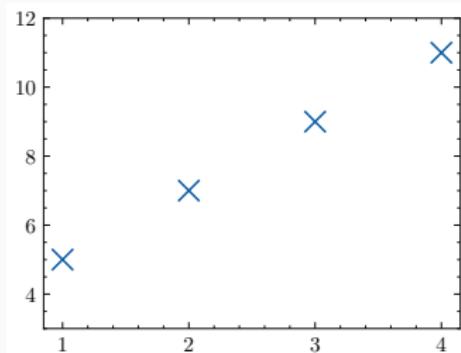
- Données d'entraînement: $(x_i, y_i)_{1 \leq i \leq N}$.
- Objectif : Trouver f_θ tel que:

$$y_i = f_\theta(x_i), \quad 1 \leq i \leq N$$



Machine learning – Learning from examples

Exemple de la regression linéaire:



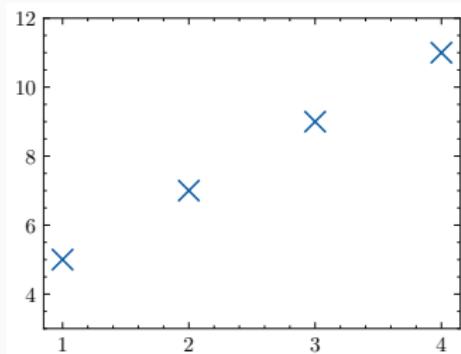
- Données d'entraînement: $(x_i, y_i)_{1 \leq i \leq N}$.
- Objectif : Trouver f_θ tel que:

$$y_i = f_\theta(x_i), \quad 1 \leq i \leq N$$

- Modèle retenu :
- $$\{f_\theta : x \mapsto ax + b\}_{\theta=(a,b) \in \mathbb{R}^2}$$

Machine learning – Learning from examples

Exemple de la regression linéaire:



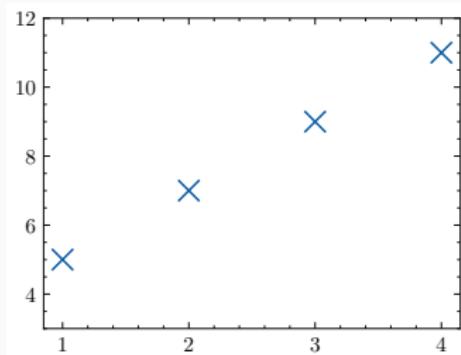
- Données d'entraînement: $(x_i, y_i)_{1 \leq i \leq N}$.
- Objectif : Trouver f_θ tel que:

$$y_i = f_\theta(x_i), \quad 1 \leq i \leq N$$

- Modèle retenu :
 - $\{f_\theta : x \mapsto ax + b\}_{\theta=(a,b) \in \mathbb{R}^2}$
 - Loss :
- $$E(\theta, (x_i, y_i)_{1 \leq i \leq N}) = \sum_{i=1}^N (f_\theta(x_i) - y_i)^2$$

Machine learning – Learning from examples

Exemple de la regression linéaire:

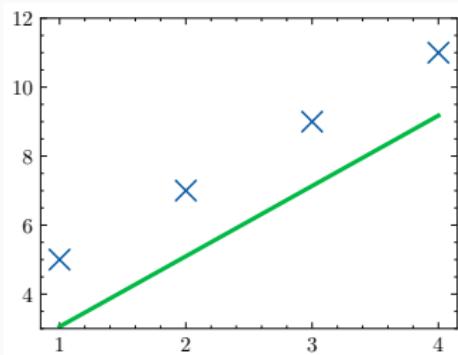


- Données d'entraînement: $(x_i, y_i)_{1 \leq i \leq N}$.
- Objectif : Trouver f_θ tel que:

$$y_i = f_\theta(x_i), \quad 1 \leq i \leq N$$

- Modèle retenu :
 $\{f_\theta : x \mapsto ax + b\}_{\theta=(a,b) \in \mathbb{R}^2}$
- Loss :
$$E(\theta, (x_i, y_i)_{1 \leq i \leq N}) = \sum_{i=1}^N (f_\theta(x_i) - y_i)^2$$
- Loss minimization:
$$\theta^* = \operatorname{argmin}_{\theta \in \mathbb{R}^2} E(\theta, (x_i, y_i)_{1 \leq i \leq N})$$

Exemple de la regression linéaire:



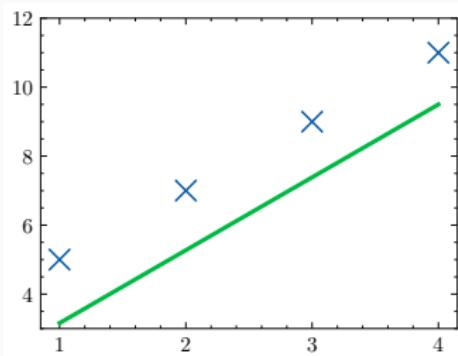
- Données d'entraînement: $(x_i, y_i)_{1 \leq i \leq N}$.
- Objectif : Trouver f_θ tel que:

$$y_i = f_\theta(x_i), \quad 1 \leq i \leq N$$

- Modèle retenu :
 $\{f_\theta : x \mapsto ax + b\}_{\theta=(a,b) \in \mathbb{R}^2}$
- Loss :
$$E(\theta, (x_i, y_i)_{1 \leq i \leq N}) = \sum_{i=1}^N (f_\theta(x_i) - y_i)^2$$
- Loss minimization:
$$\theta^* = \operatorname{argmin}_{\theta \in \mathbb{R}^2} E(\theta, (x_i, y_i)_{1 \leq i \leq N})$$

Machine learning – Learning from examples

Exemple de la regression linéaire:



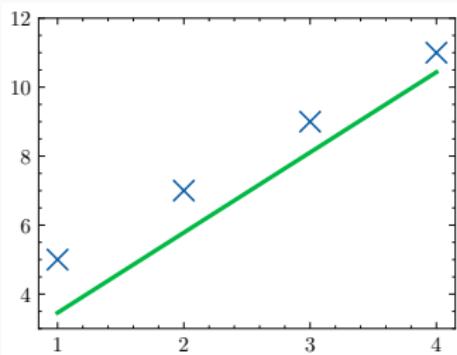
- Données d'entraînement: $(x_i, y_i)_{1 \leq i \leq N}$.
- Objectif : Trouver f_θ tel que:

$$y_i = f_\theta(x_i), \quad 1 \leq i \leq N$$

- Modèle retenu :
 $\{f_\theta : x \mapsto ax + b\}_{\theta=(a,b) \in \mathbb{R}^2}$
- Loss :
$$E(\theta, (x_i, y_i)_{1 \leq i \leq N}) = \sum_{i=1}^N (f_\theta(x_i) - y_i)^2$$
- Loss minimization:
$$\theta^* = \operatorname{argmin}_{\theta \in \mathbb{R}^2} E(\theta, (x_i, y_i)_{1 \leq i \leq N})$$

Machine learning – Learning from examples

Exemple de la regression linéaire:



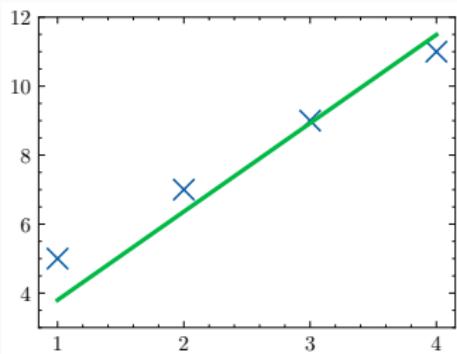
- Données d'entraînement: $(x_i, y_i)_{1 \leq i \leq N}$.
- Objectif : Trouver f_θ tel que:

$$y_i = f_\theta(x_i), \quad 1 \leq i \leq N$$

- Modèle retenu :
 $\{f_\theta : x \mapsto ax + b\}_{\theta=(a,b) \in \mathbb{R}^2}$
- Loss :
$$E(\theta, (x_i, y_i)_{1 \leq i \leq N}) = \sum_{i=1}^N (f_\theta(x_i) - y_i)^2$$
- Loss minimization:
$$\theta^* = \operatorname{argmin}_{\theta \in \mathbb{R}^2} E(\theta, (x_i, y_i)_{1 \leq i \leq N})$$

Machine learning – Learning from examples

Exemple de la regression linéaire:



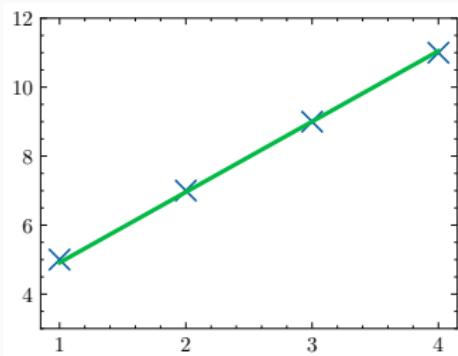
- Données d'entraînement: $(x_i, y_i)_{1 \leq i \leq N}$.
- Objectif : Trouver f_θ tel que:

$$y_i = f_\theta(x_i), \quad 1 \leq i \leq N$$

- Modèle retenu :
 $\{f_\theta : x \mapsto ax + b\}_{\theta=(a,b) \in \mathbb{R}^2}$
- Loss :
$$E(\theta, (x_i, y_i)_{1 \leq i \leq N}) = \sum_{i=1}^N (f_\theta(x_i) - y_i)^2$$
- Loss minimization:
$$\theta^* = \operatorname{argmin}_{\theta \in \mathbb{R}^2} E(\theta, (x_i, y_i)_{1 \leq i \leq N})$$

Machine learning – Learning from examples

Exemple de la regression linéaire:



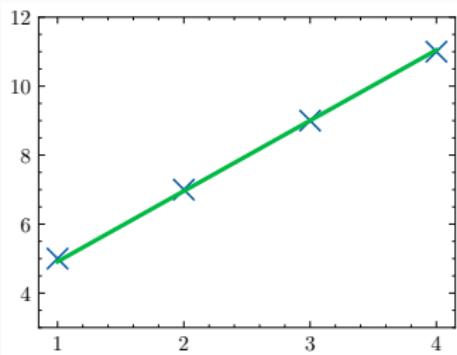
- Données d'entraînement: $(x_i, y_i)_{1 \leq i \leq N}$.
- Objectif : Trouver f_θ tel que:

$$y_i = f_\theta(x_i), \quad 1 \leq i \leq N$$

- Modèle retenu :
 $\{f_\theta : x \mapsto ax + b\}_{\theta=(a,b) \in \mathbb{R}^2}$
- Loss :
$$E(\theta, (x_i, y_i)_{1 \leq i \leq N}) = \sum_{i=1}^N (f_\theta(x_i) - y_i)^2$$
- Loss minimization:
$$\theta^* = \operatorname{argmin}_{\theta \in \mathbb{R}^2} E(\theta, (x_i, y_i)_{1 \leq i \leq N})$$

Machine learning – Learning from examples

Exemple de la regression linéaire:

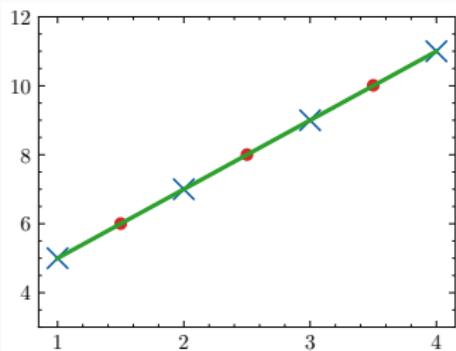


- Données d'entraînement: $(x_i, y_i)_{1 \leq i \leq N}$.
- Objectif : Trouver f_θ tel que:

$$y_i = f_\theta(x_i), \quad 1 \leq i \leq N$$

- Modèle retenu :
 $\{f_\theta : x \mapsto ax + b\}_{\theta=(a,b) \in \mathbb{R}^2}$
- Loss :
$$E(\theta, (x_i, y_i)_{1 \leq i \leq N}) = \sum_{i=1}^N (f_\theta(x_i) - y_i)^2$$
- Loss minimization:
$$\theta^* = \operatorname{argmin}_{\theta \in \mathbb{R}^2} E(\theta, (x_i, y_i)_{1 \leq i \leq N})$$
- Puis, phase de test.

Exemple de la regression linéaire:



- Données d'entraînement: $(x_i, y_i)_{1 \leq i \leq N}$.
- Objectif : Trouver f_θ tel que:

$$y_i = f_\theta(x_i), \quad 1 \leq i \leq N$$

- Modèle retenu :
 $\{f_\theta : x \mapsto ax + b\}_{\theta=(a,b) \in \mathbb{R}^2}$
- Loss :
$$E(\theta, (x_i, y_i)_{1 \leq i \leq N}) = \sum_{i=1}^N (f_\theta(x_i) - y_i)^2$$
- Loss minimization:
$$\theta^* = \operatorname{argmin}_{\theta \in \mathbb{R}^2} E(\theta, (x_i, y_i)_{1 \leq i \leq N})$$
- Puis, phase de test.

Terminology

Echantillon (Sample/Observation/Data): Elément à traiter (e.g., classify).

Example: un individu, un document, une image, un son, une vidéo... .

Terminology

Echantillon (Sample/Observation/Data): Elément à traiter (e.g., classify).

Example: un individu, un document, une image, un son, une vidéo... .

Features (Input): Ensemble de traits distincts qui peuvent être utilisés pour décrire chaque échantillon de manière quantitative. Habituellement représenté comme un vecteur multidimensionnel désigné par x . *Exemple: taille, poids, nationalité, rayon ...*

Terminology

Echantillon (Sample/Observation/Data): Elément à traiter (e.g., classify).

Example: un individu, un document, une image, un son, une vidéo... .

Features (Input): Ensemble de traits distincts qui peuvent être utilisés pour décrire chaque échantillon de manière quantitative. Habituellement représenté comme un vecteur multidimensionnel désigné par x . *Exemple: taille, poids, nationalité, rayon ...*

Training set: Ensemble de données utilisé pour découvrir des relations potentiellement prédictives.

Terminology

Echantillon (Sample/Observation/Data): Elément à traiter (e.g., classify).

Example: un individu, un document, une image, un son, une vidéo... .

Features (Input): Ensemble de traits distincts qui peuvent être utilisés pour décrire chaque échantillon de manière quantitative. Habituellement représenté comme un vecteur multidimensionnel désigné par x . *Exemple: taille, poids, nationalité, rayon ...*

Training set: Ensemble de données utilisé pour découvrir des relations potentiellement prédictives.

Validation set: Ensemble utilisé pour ajuster les hyperparamètres du modèle (à utiliser comme feedback).

Terminology

Echantillon (Sample/Observation/Data): Elément à traiter (e.g., classify).

Example: un individu, un document, une image, un son, une vidéo... .

Features (Input): Ensemble de traits distincts qui peuvent être utilisés pour décrire chaque échantillon de manière quantitative. Habituellement représenté comme un vecteur multidimensionnel désigné par x . *Exemple: taille, poids, nationalité, rayon ...*

Training set: Ensemble de données utilisé pour découvrir des relations potentiellement prédictives.

Validation set: Ensemble utilisé pour ajuster les hyperparamètres du modèle (à utiliser comme feedback).

Testing set: Ensemble utilisé pour tester les performances du modèles.

Terminology

Echantillon (Sample/Observation/Data): Elément à traiter (e.g., classify).

Example: un individu, un document, une image, un son, une vidéo... .

Features (Input): Ensemble de traits distincts qui peuvent être utilisés pour décrire chaque échantillon de manière quantitative. Habituellement représenté comme un vecteur multidimensionnel désigné par x . *Exemple: taille, poids, nationalité, rayon ...*

Training set: Ensemble de données utilisé pour découvrir des relations potentiellement prédictives.

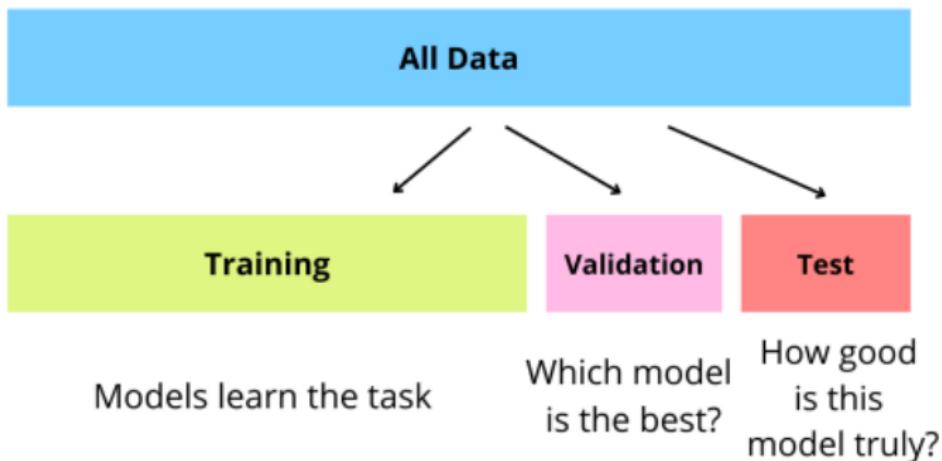
Validation set: Ensemble utilisé pour ajuster les hyperparamètres du modèle (à utiliser comme feedback).

Testing set: Ensemble utilisé pour tester les performances du modèles.

Label (Output):

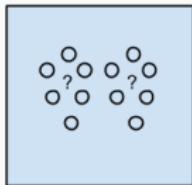
Classe ou résultat attribué à un échantillon. La prédiction réelle est souvent

Training set/Test set

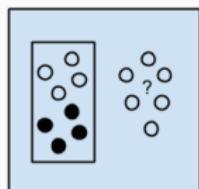


In our examples, there will be not validation dataset.

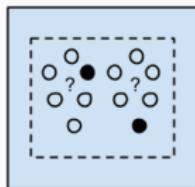
Plusieurs façons d'entraîner une fonction



Unsupervised Learning
Algorithms

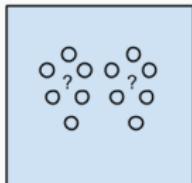


Supervised Learning
Algorithms

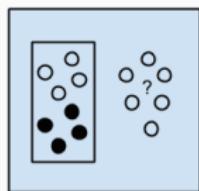


Semi-supervised
Learning Algorithms

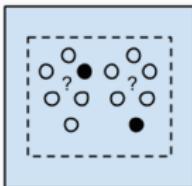
Plusieurs façons d'entraîner une fonction



Unsupervised Learning
Algorithms



Supervised Learning
Algorithms

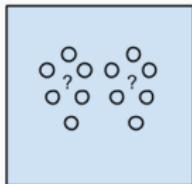


Semi-supervised
Learning Algorithms

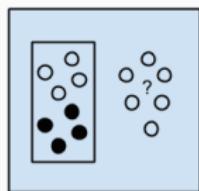
Apprentissage non supervisé: En cherchant à différencier les données sans qu'elles soient explicitement étiquetées.

Example: cluster similar documents based on the text content.

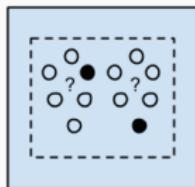
Plusieurs façons d'entraîner une fonction



Unsupervised Learning Algorithms



Supervised Learning Algorithms



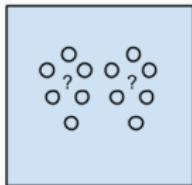
Semi-supervised Learning Algorithms

Apprentissage non supervisé: En cherchant à différencier les données sans qu'elles soient explicitement étiquetées.

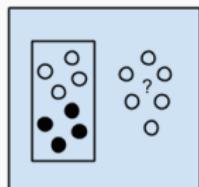
Example: cluster similar documents based on the text content.

Apprentissage supervisé: En apprenant à différencier les données à partir de données déjà étiquetées. **labeled training set.** *Example: email spam detector with training set of already labeled emails.*

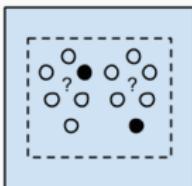
Plusieurs façons d'entraîner une fonction



Unsupervised Learning Algorithms



Supervised Learning Algorithms



Semi-supervised Learning Algorithms

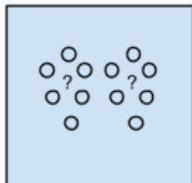
Apprentissage non supervisé: En cherchant à différencier les données sans qu'elles soient explicitement étiquetées.

Example: cluster similar documents based on the text content.

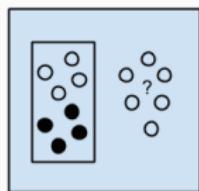
Apprentissage supervisé: En apprenant à différencier les données à partir de données déjà étiquetées. **labeled training set.** *Example: email spam detector with training set of already labeled emails.*

Apprentissage semi-supervisé: En apprenant avec une petite quantité de données étiquetées.. *Example: web content and protein sequence classifications.*

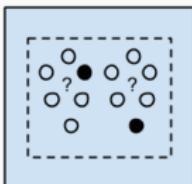
Plusieurs façons d'entraîner une fonction



Unsupervised Learning Algorithms



Supervised Learning Algorithms



Semi-supervised Learning Algorithms

Apprentissage non supervisé: En cherchant à différencier les données sans qu'elles soient explicitement étiquetées.

Example: cluster similar documents based on the text content.

Apprentissage supervisé: En apprenant à différencier les données à partir de données déjà étiquetées. **labeled training set.**

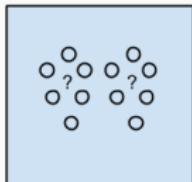
Example: email spam detector with training set of already labeled emails.

Apprentissage semi-supervisé: En apprenant avec une petite quantité de données étiquetées.. *Example: web content and protein sequence classifications.*

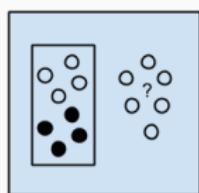
Reinforcement learning: Learning based on **feedback** or reward. *Example: learn to play chess by winning or losing.*

(Source: Jason Brownlee and Lucas Masuch)

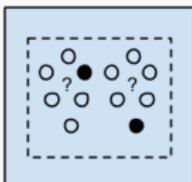
Plusieurs façons d'entraîner une fonction



Unsupervised Learning Algorithms



Supervised Learning Algorithms



Semi-supervised Learning Algorithms

Apprentissage non supervisé: En cherchant à différencier les données sans qu'elles soient explicitement étiquetées.
Example: cluster similar documents based on the text content.

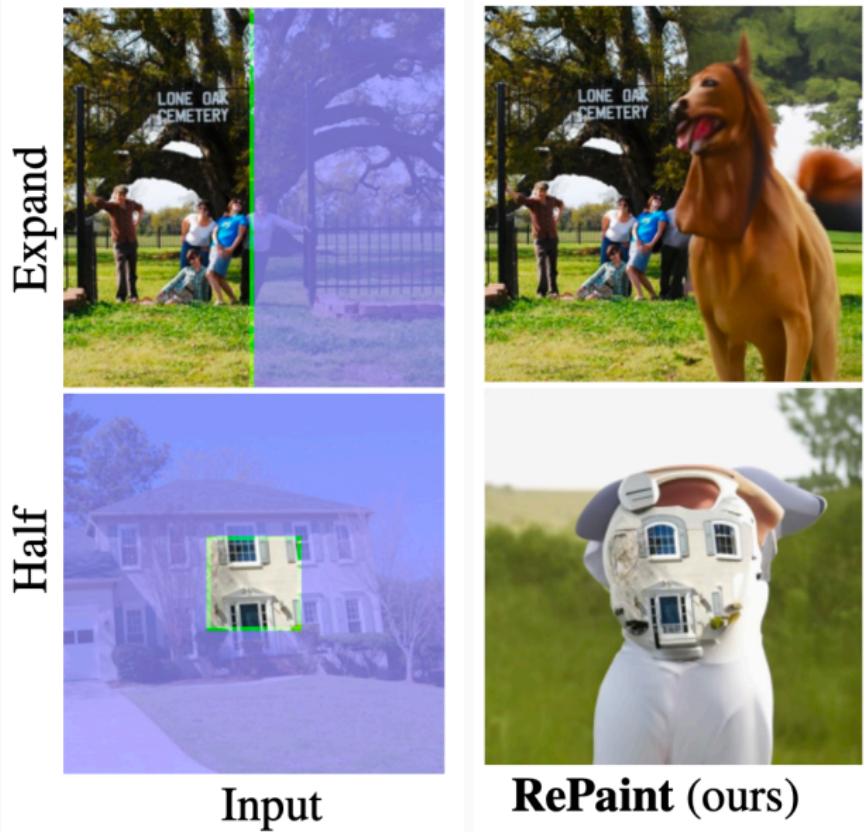
Apprentissage supervisé: En apprenant à différencier les données à partir de données déjà étiquetées. **labeled training set.** *Example: email spam detector with training set of already labeled emails.*

Apprentissage semi-supervisé: En apprenant avec une petite quantité de données étiquetées.. *Example: web content and protein sequence classifications.*

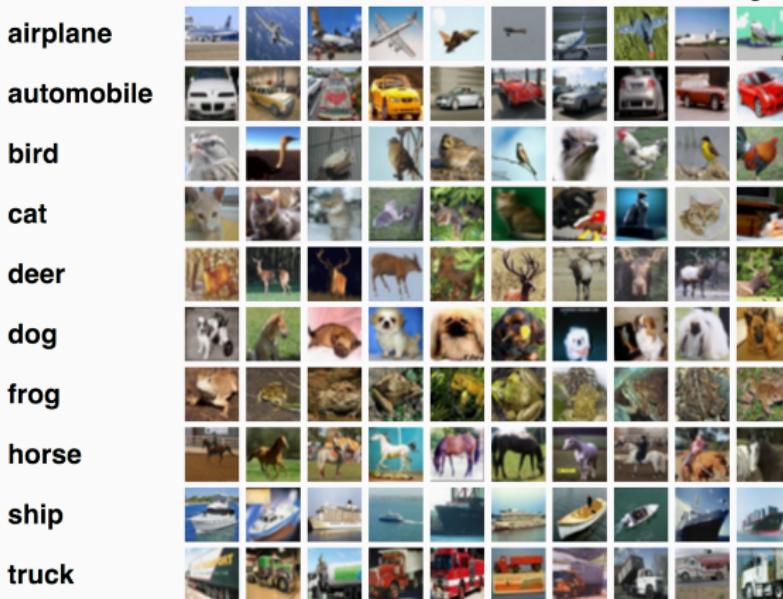
Reinforcement learning: Learning based on **feedback** or reward. *Example: learn to play chess by winning or losing.*

(Source: Jason Brownlee and Lucas Masuch)

Dataset problem



Classification example



Dataset: $(x_i, d_i)_{1 \leq i \leq N}$ where for $1 \leq i \leq N$, x_i is an image and x_i is in class d_i .

Learning from examples

3 main ingredients

- ① Training set / examples:

$$\{\boldsymbol{x}_1, \boldsymbol{x}_2, \dots, \boldsymbol{x}_N\}$$

Learning from examples

3 main ingredients

- ① Training set / examples:

$$\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$$

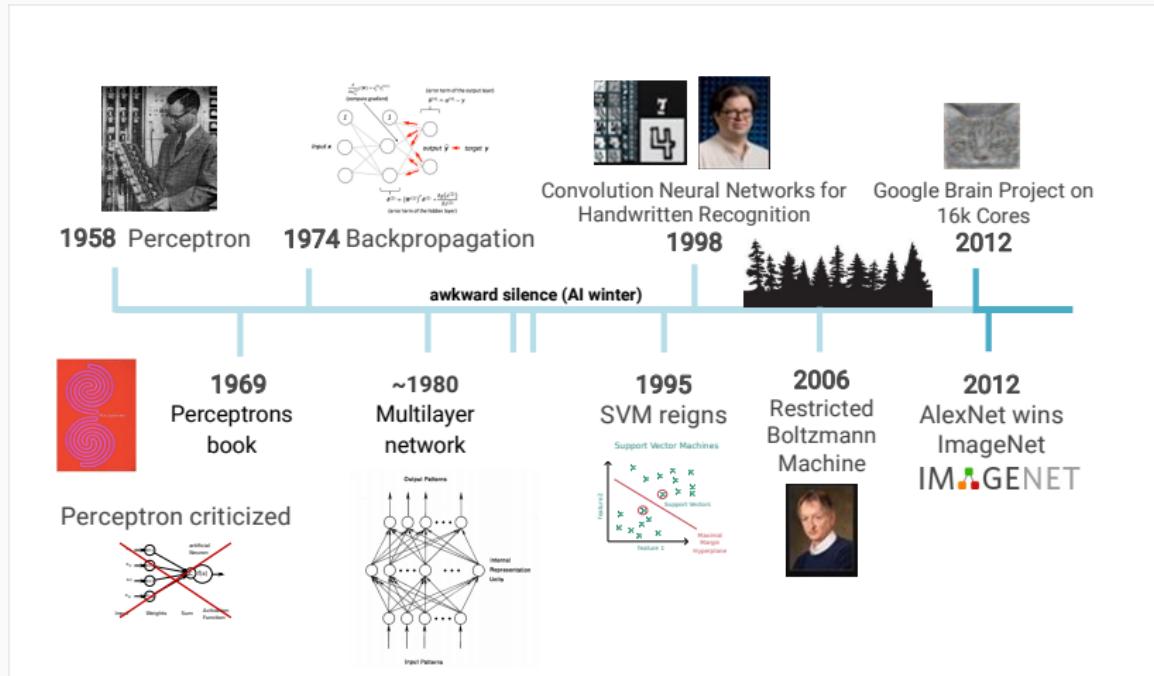
- ② Machine or model:

$$\mathbf{x} \rightarrow \underbrace{f(\mathbf{x}; \theta)}_{\text{function / algorithm}} \rightarrow \underbrace{\mathbf{y}}_{\text{prediction}}$$

θ : parameters of the model

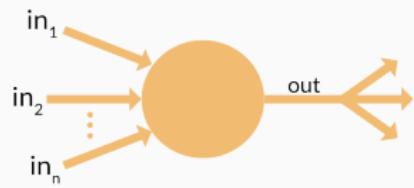
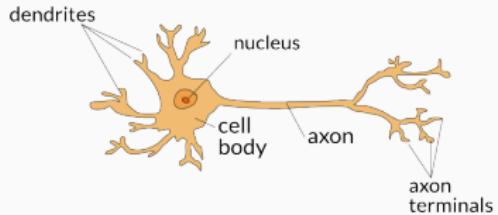
Machine learning – Timeline

Timeline of (deep) learning



(Source: Lucas Masuch & Vincent Lepetit)

Perceptron



Perceptron

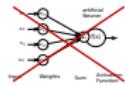


1958 Perceptron



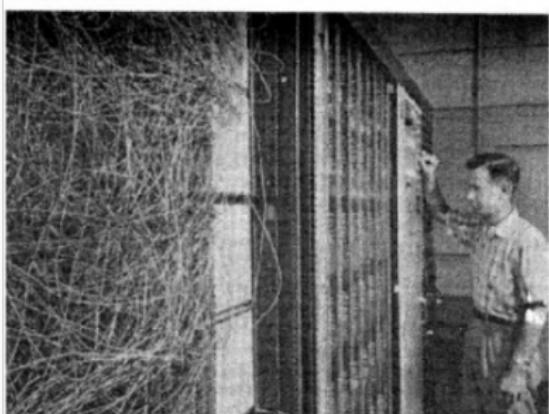
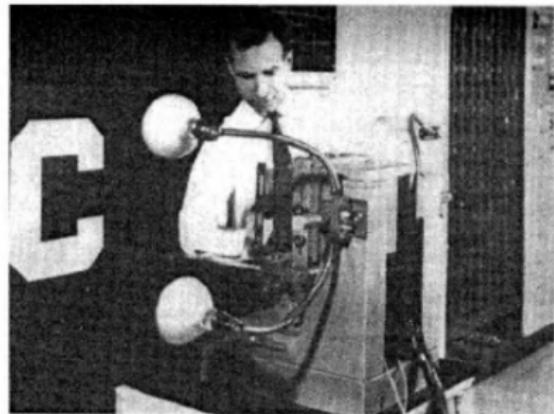
1969
Perceptrons
book

Perceptron criticized



(Source: Lucas Masuch & Vincent Lepetit)

Perceptron (Frank Rosenblatt, 1958)

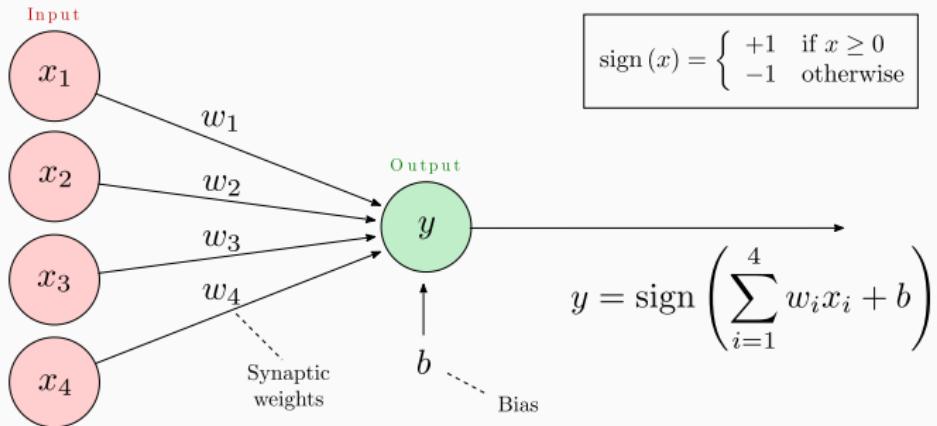


First binary classifier based on supervised learning (discrimination).

Foundation of modern artificial neural networks.

At that time: technological, scientific and philosophical challenges.

Representation of the Perceptron



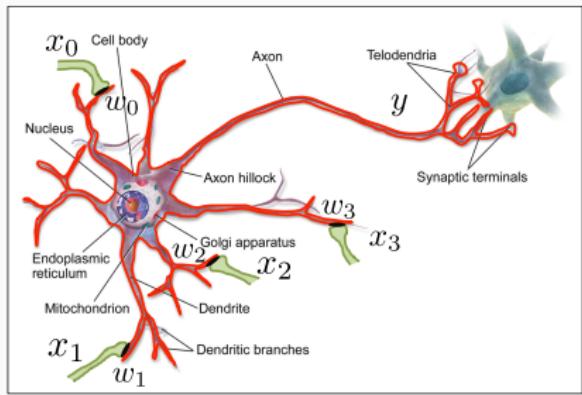
Parameters of the perceptron

- w_k : synaptic weights
 - b : bias
- $\left. \begin{array}{l} \\ \end{array} \right\}$ ← real parameters to be estimated.

Training = adjusting the weights and biases

The origin of the Perceptron

Takes inspiration from the visual system known for its ability to learn patterns.



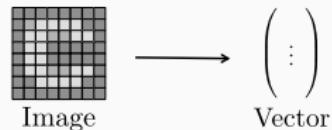
- When a neuron receives a stimulus with high enough voltage, it emits an **action potential** (aka, nerve impulse or spike). It is said to **fire**.
- The perceptron mimics this activation effect: it fires only when

$$\sum_i w_i x_i + b > 0$$

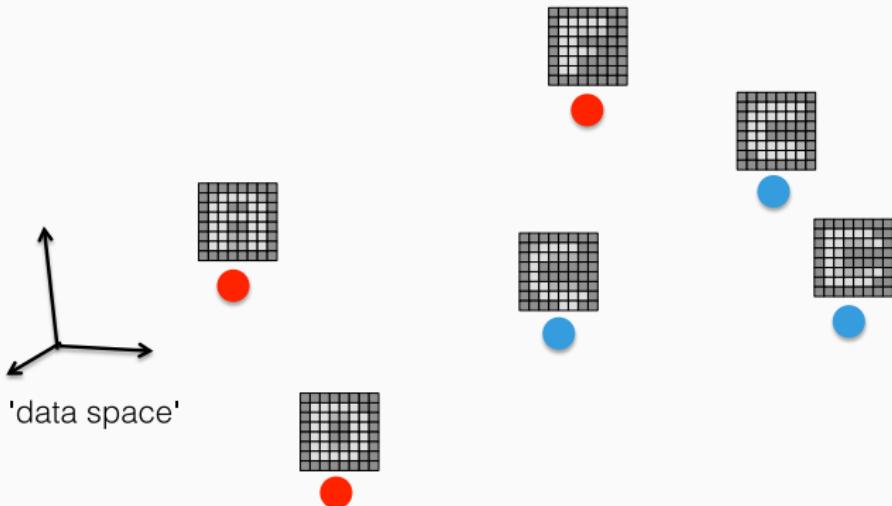
$$y = \underbrace{\text{sign}(w_0 x_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + b)}_{f(\mathbf{x}; \mathbf{w})} = \begin{cases} +1 & \text{for the first class} \\ -1 & \text{for the second class} \end{cases}$$

Machine learning – Perceptron – Principle

- ① Data are represented as vectors:



- ② Collect training data with **positive** and **negative** examples:



(Source: Vincent Lepetit)

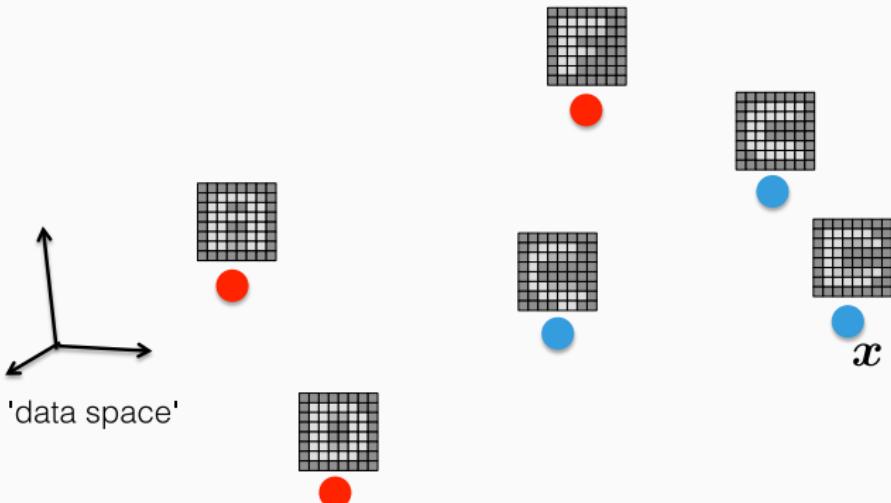
Machine learning – Perceptron – Principle

③ Training: find w and b so that:

- $\langle w, x \rangle + b$ is positive for **positive samples** x ,
- $\langle w, x \rangle + b$ is negative for **negative samples** x .

Dot product:

$$\begin{aligned}\langle w, x \rangle &= \sum_{i=1}^d w_i x_i \\ &= w^T x\end{aligned}$$



(Source: Vincent Lepetit)

Machine learning – Perceptron – Principle

③ Training: find w and b so that:

- $\langle w, x \rangle + b$ is positive for **positive samples** x ,
- $\langle w, x \rangle + b$ is negative for **negative samples** x .

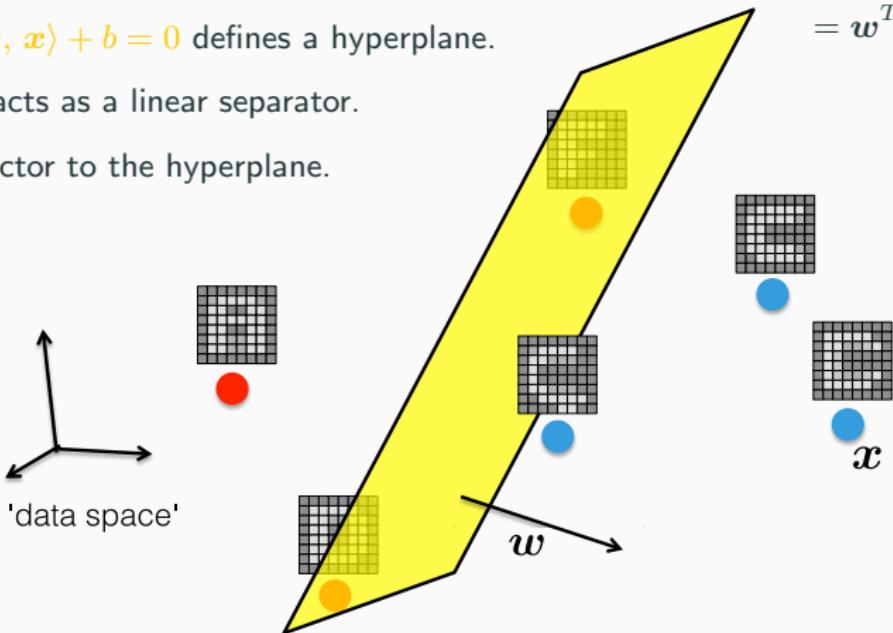
Dot product:

$$\begin{aligned}\langle w, x \rangle &= \sum_{i=1}^d w_i x_i \\ &= w^T x\end{aligned}$$

The equation $\langle w, x \rangle + b = 0$ defines a hyperplane.

The hyperplane acts as a linear separator.

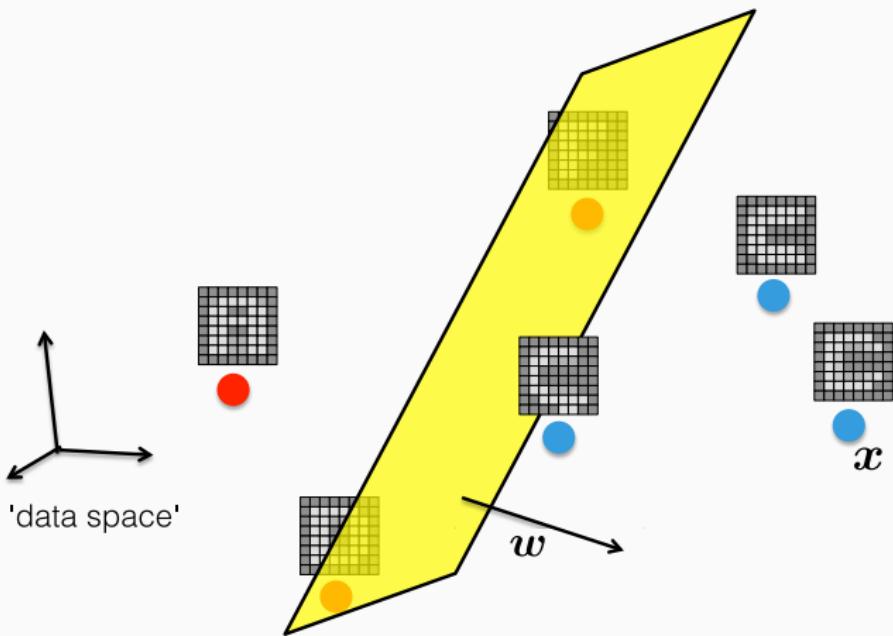
w is a normal vector to the hyperplane.



(Source: Vincent Lepetit)

Machine learning – Perceptron – Principle

- ④ **Testing:** the perceptron can now classify new examples.

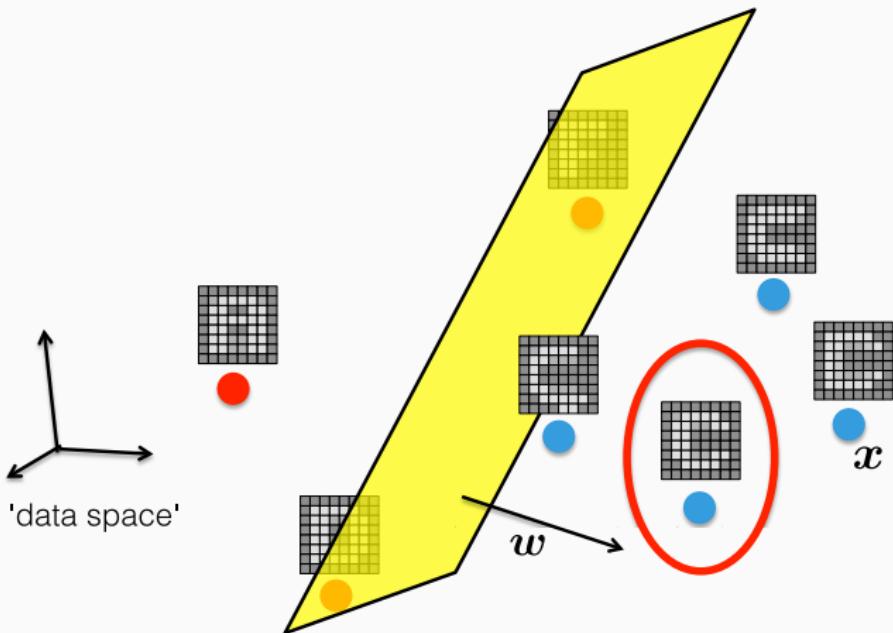


(Source: Vincent Lepetit)

Machine learning – Perceptron – Principle

④ **Testing:** the perceptron can now classify new examples.

- A new example x is classified **positive** if $\langle w, x \rangle + b$ is **positive**,

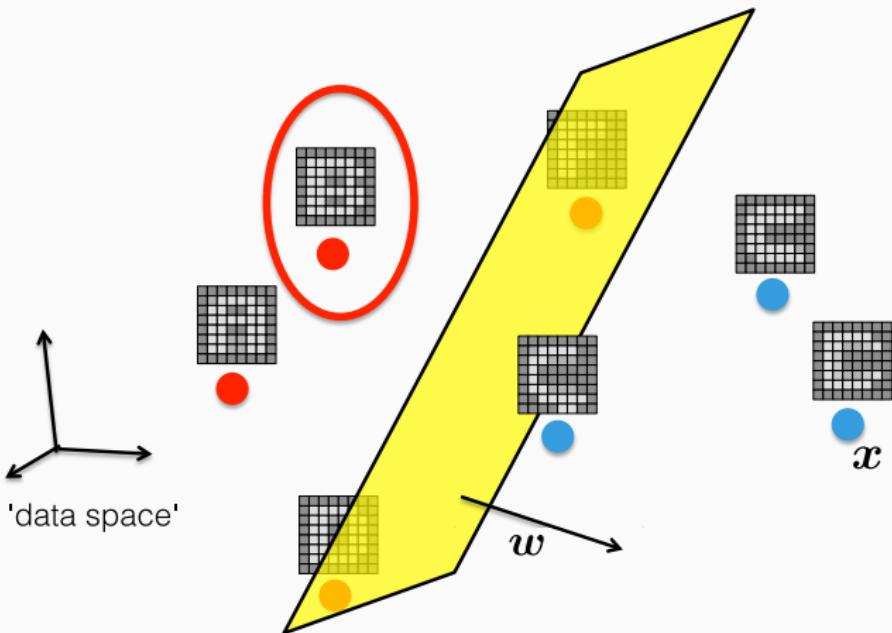


(Source: Vincent Lepetit)

Machine learning – Perceptron – Principle

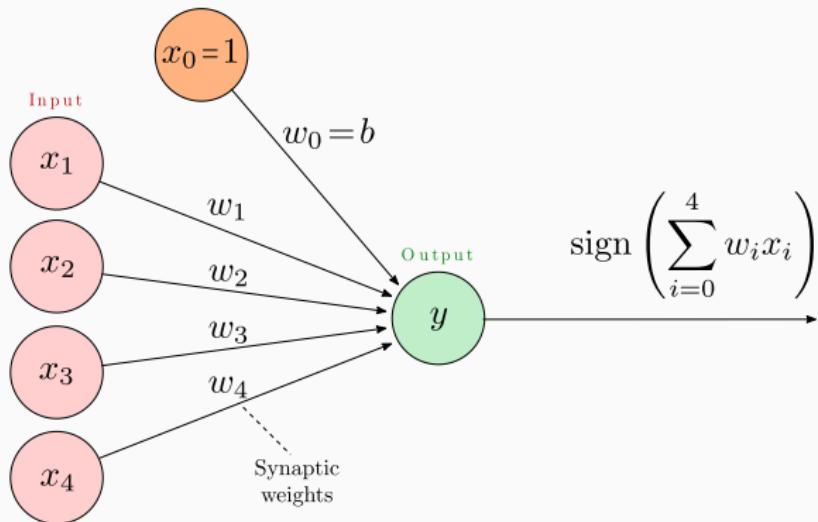
④ **Testing:** the perceptron can now classify new examples.

- A new example x is classified **positive** if $\langle w, x \rangle + b$ is **positive**,
- and **negative** if $\langle w, x \rangle + b$ is **negative**.



(Source: Vincent Lepetit)

Alternative representation

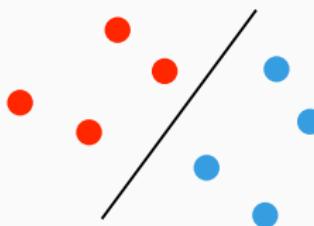


Use the zero-index to encode the bias as a synaptic weight.

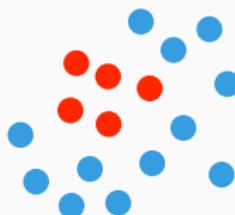
Simplifies algorithms as all parameters can now be processed in the same way.

Perceptrons book (Minsky and Papert, 1969)

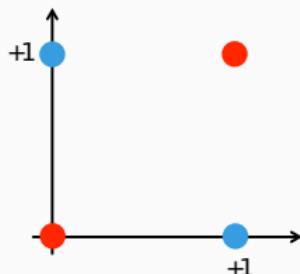
A perceptron can only classify data points that are linearly separable:



Linearly separable



Nonlinearly separable



The xor function

Seen by many as a justification to stop research on perceptrons.

(Source: Vincent Lepetit)

Artificial neural network



Artificial neural network



1958 Perceptron



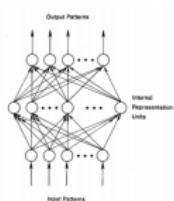
1969
Perceptrons
book

Perceptron criticized



~1980
Multilayer
network

1989
Universal
Approximation
Theorem



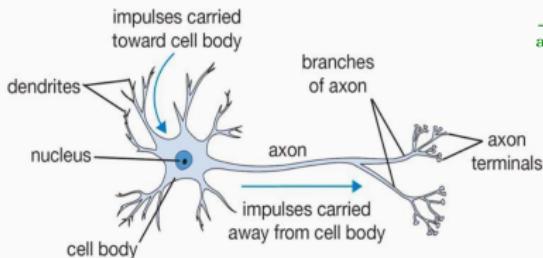
(Source: Lucas Masuch & Vincent Lepetit)

Artificial neural network

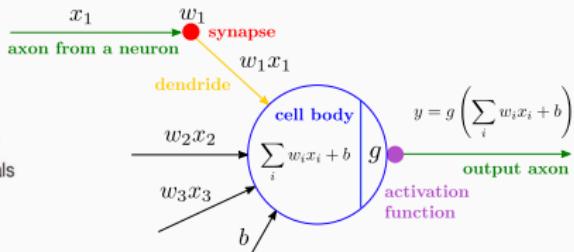


- Supervised learning method initially inspired by the behavior of the human brain.
- Consists of the inter-connection of several small units (just like in the human brain).
- Introduced in the late 50s, very popular in the 90s, reappeared in the 2010s with deep learning.
- Also referred to as **Multi-Layer Perceptron (MLP)**.
- Historically used after feature extraction.

Artificial neuron (McCulloch & Pitts, 1943)



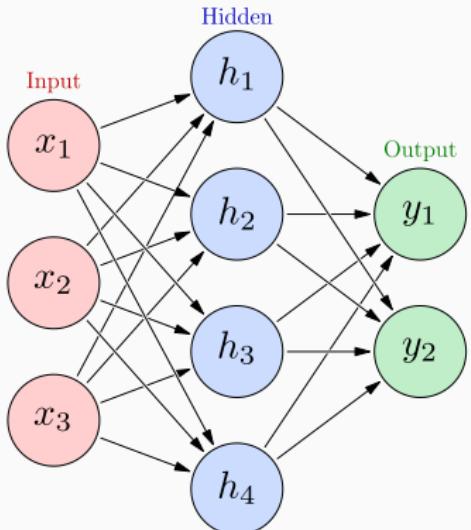
Biological neuron



Artificial neuron

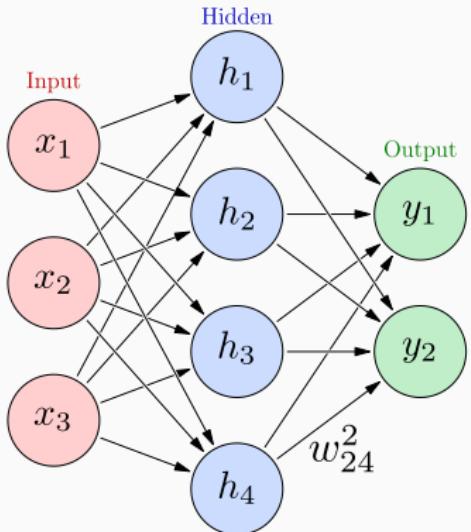
- An artificial neuron contains several incoming **weighted connections**, an outgoing connection and has a **nonlinear activation function** g .
- Neurons are **trained to filter and detect specific features** or patterns (e.g. edge, nose) by receiving weighted input, transforming it with the activation function and passing it to the outgoing connections.
- Unlike the perceptron, can be used for regression (with proper choice of g).

Artificial neural network / Multilayer perceptron / NeuralNet



- Inter-connection of several artificial neurons (also called nodes or units).
- Each level in the graph is called a layer:
 - Input layer,
 - Hidden layer(s),
 - Output layer.
- Each neuron in the hidden layers acts as a classifier / feature detector.
- Feedforward NN (no cycle)
 - first and simplest type of NN,
 - information moves in one direction.
- Recurrent NN (with cycle)
 - used for time sequences,
 - such as speech-recognition.

Artificial neural network / Multilayer perceptron / NeuralNet



$$h_1 = g_1 (w_{11}^1 x_1 + w_{12}^1 x_2 + w_{13}^1 x_3 + b_1^1)$$

$$h_2 = g_1 (w_{21}^1 x_1 + w_{22}^1 x_2 + w_{23}^1 x_3 + b_2^1)$$

$$h_3 = g_1 (w_{31}^1 x_1 + w_{32}^1 x_2 + w_{33}^1 x_3 + b_3^1)$$

$$h_4 = g_1 (w_{41}^1 x_1 + w_{42}^1 x_2 + w_{43}^1 x_3 + b_4^1)$$

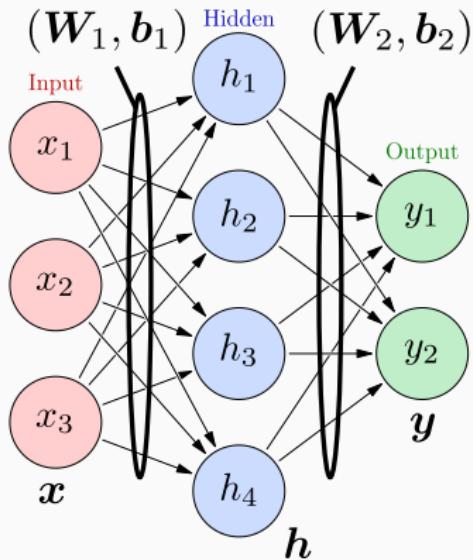
$$y_1 = g_2 (w_{11}^2 h_1 + w_{12}^2 h_2 + w_{13}^2 h_3 + w_{14}^2 h_4 + b_1^2)$$

$$y_2 = g_2 (w_{21}^2 h_1 + w_{22}^2 h_2 + w_{23}^2 h_3 + w_{24}^2 h_4 + b_2^2)$$

w_{ij}^k synaptic weight between previous node j and next node i at layer k .

g_k are any activation function applied to each coefficient of its input vector.

Artificial neural network / Multilayer perceptron / NeuralNet



$$h_1 = g_1 (w_{11}^1 x_1 + w_{12}^1 x_2 + w_{13}^1 x_3 + b_1^1)$$

$$h_2 = g_1 (w_{21}^1 x_1 + w_{22}^1 x_2 + w_{23}^1 x_3 + b_2^1)$$

$$h_3 = g_1 (w_{31}^1 x_1 + w_{32}^1 x_2 + w_{33}^1 x_3 + b_3^1)$$

$$h_4 = g_1 (w_{41}^1 x_1 + w_{42}^1 x_2 + w_{43}^1 x_3 + b_4^1)$$

$$\mathbf{h} = g_1 (\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)$$

$$y_1 = g_2 (w_{11}^2 h_1 + w_{12}^2 h_2 + w_{13}^2 h_3 + w_{14}^2 h_4 + b_1^2)$$

$$y_2 = g_2 (w_{21}^2 h_1 + w_{22}^2 h_2 + w_{23}^2 h_3 + w_{24}^2 h_4 + b_2^2)$$

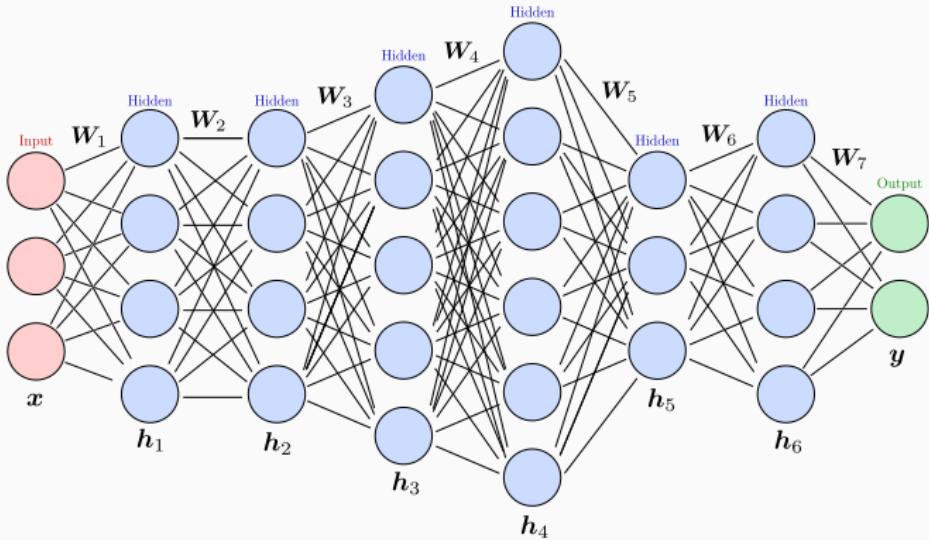
$$\mathbf{y} = g_2 (\mathbf{W}_2 \mathbf{h} + \mathbf{b}_2)$$

w_{ij}^k synaptic weight between previous node j and next node i at layer k .

g_k are any activation function applied to each coefficient of its input vector.

The matrices \mathbf{W}_k and biases \mathbf{b}_k are learned from labeled training data.

Artificial neural network / Multilayer perceptron



It can have 1 hidden layer only (**shallow network**),
It can have more than 1 hidden layer (**deep network**),
each layer may have a different size, and
hidden and output layers often have different activation functions.

Artificial neural network / Multilayer perceptron

- As for the perceptron, the biases can be integrated into the weights:

$$\mathbf{W}_k \mathbf{h}_{k-1} + \mathbf{b}_k = \underbrace{\begin{pmatrix} \mathbf{b}_k & \mathbf{W}_k \end{pmatrix}}_{\tilde{\mathbf{W}}_k} \underbrace{\begin{pmatrix} 1 \\ \mathbf{h}_{k-1} \end{pmatrix}}_{\tilde{\mathbf{h}}_{k-1}} = \tilde{\mathbf{W}}_k \tilde{\mathbf{h}}_{k-1}$$

- A neural network with L layers is a function of \mathbf{x} parameterized by $\tilde{\mathbf{W}}$:

$$\mathbf{y} = f(\mathbf{x}; \tilde{\mathbf{W}}) \quad \text{where} \quad \tilde{\mathbf{W}} = (\tilde{\mathbf{W}}_1, \tilde{\mathbf{W}}_2, \dots, \tilde{\mathbf{W}}_L)^T$$

- It can be defined recursively as

$$\mathbf{y} = f(\mathbf{x}; \tilde{\mathbf{W}}) = \mathbf{h}_L, \quad \mathbf{h}_k = g_k \left(\tilde{\mathbf{W}}_k \tilde{\mathbf{h}}_{k-1} \right) \quad \text{and} \quad \mathbf{h}_0 = \mathbf{x}$$

- For simplicity, $\tilde{\mathbf{W}}$ will be denoted \mathbf{W} (when no possible confusions).

Activation functions

Linear units: $g(a) = a$

$$y = \mathbf{W}_L h_{L-1} + b_L$$

$$\begin{array}{c} \mathbf{h}_{L-1} = \mathbf{W}_{L-1} \mathbf{h}_{L-2} + \mathbf{b}_{L-1} \\ \hline \mathbf{y} = \mathbf{W}_L \mathbf{W}_{L-1} \mathbf{h}_{L-2} + \mathbf{W}_L \mathbf{b}_{L-1} + \mathbf{b}_L \\ \hline \mathbf{y} = \mathbf{W}_L \dots \mathbf{W}_1 \mathbf{x} + \sum_{k=1}^{L-1} \mathbf{W}_L \dots \mathbf{W}_{k+1} \mathbf{b}_k + \mathbf{b}_L \end{array}$$

We can always find an equivalent network without hidden units,
because compositions of affine functions are affine.

In general, **non-linearity** is needed to learn complex (non-linear)
representations of data, otherwise the NN would be just a linear function.
Otherwise, back to the problem of nonlinearly separable datasets.

Activation functions

Threshold units: for instance the sign function

$$g(a) = \begin{cases} -1 & \text{if } a < 0 \\ +1 & \text{otherwise.} \end{cases}$$

or ReLu activation function

$$g(a) = \max(0, a) = \begin{cases} 0 & \text{if } a < 0 \\ 1 & \text{otherwise.} \end{cases}$$

Discontinuities in the hidden layers
make the optimization really difficult.

We prefer functions that are continuous and differentiable.

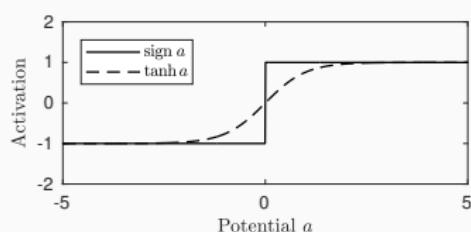
Activation functions

Sigmoidal units: for instance the hyperbolic tangent function

$$g(a) = \tanh a = \frac{e^a - e^{-a}}{e^a + e^{-a}} \in [-1, 1]$$

or the logistic sigmoid function

$$g(a) = \frac{1}{1 + e^{-a}} \in [0, 1]$$



- In fact equivalent by linear transformations :

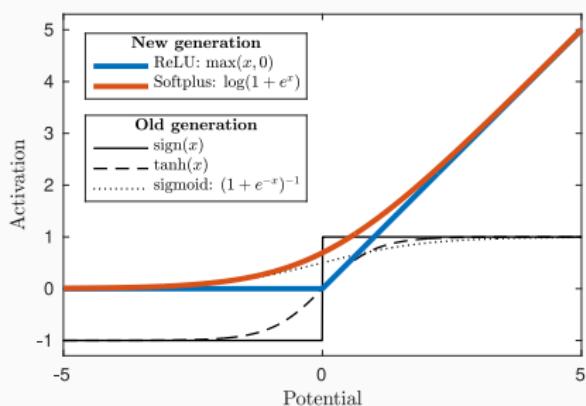
$$\tanh(a/2) = 2\text{logistic}(a) - 1$$

- Differentiable approximations of the sign and step functions, respectively.
- Act as threshold units for large values of $|a|$ and as linear for small values.

Activation functions

“Modern” units:

$$\underbrace{g(a) = \max(a, 0)}_{\text{ReLU}} \quad \text{or} \quad \underbrace{g(a) = \log(1 + e^a)}_{\text{Softplus}}$$



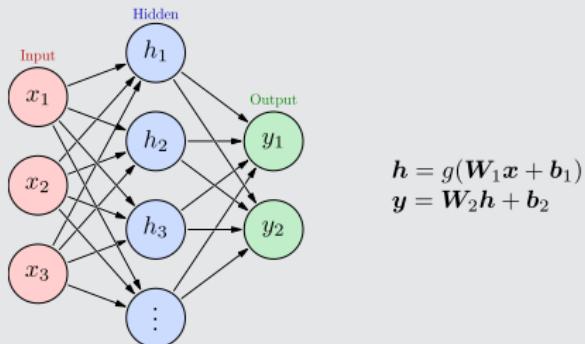
Most neural networks use **ReLU** (Rectifier linear unit) – $\max(a, 0)$ – nowadays for hidden layers, since it trains much faster, is more expressive than logistic function and prevents the gradient vanishing problem.

(Source: Lucas Masuch)

Universal Approximation Theorem

(Hornik et al, 1989; Cybenko, 1989)

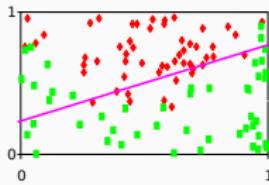
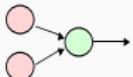
Any continuous function $f : K \subset \mathbb{R}^N \rightarrow \mathbb{R}^K$ can be uniformly approximated by a feedforward **shallow network** (i.e., with 1-hidden layer only) with a sufficient number of neurons in the hidden layer.



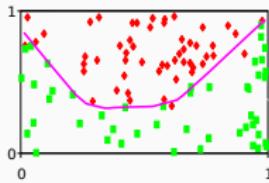
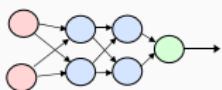
- Works if and only if g is not polynomial (and thus non linear).
- The theorem does not say how large the network needs to be.
- No guarantee that the training algorithm will be able to train the network.

Machine learning – ANN

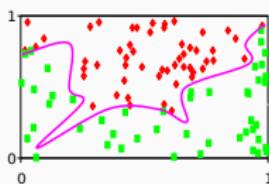
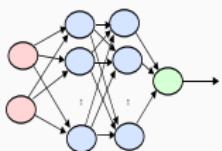
1 neuron



2+2+1 neurons



10+10+1 neurons



Complexity/capacity of the

network

⇒

**Trade-off between
generalization and overfitting.**

Learning from examples

3 main ingredients

- ① Training set / examples:

$$\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$$

- ② Machine or model:

$$\mathbf{x} \rightarrow \underbrace{f(\mathbf{x}; \theta)}_{\text{function / algorithm}} \rightarrow \underbrace{\mathbf{y}}_{\text{prediction}}$$

θ : parameters of the model

- ③ Loss, cost, objective function / energy:

$$\operatorname{argmin}_{\theta} E(\theta; \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$$

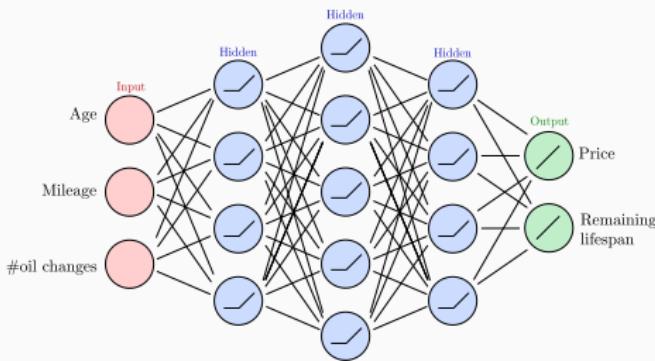
Approximation – Least square regression

- **Goal:** Predict a **real multivariate function**.
- **How:** estimate the coefficients $\theta = (\mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2, \dots)$ of $d_i = f(\mathbf{x}_i; \mathbf{W})$, $1 \leq i \leq d$?
from labeled training examples where labels are real vectors:

$$\mathcal{T} = \{(\mathbf{x}^i, \mathbf{d}^i)\}_{i=1..N}$$

\nearrow \downarrow \searrow
i-th training example desired output for sample i number of training samples

- **Typical architecture:**



Approximation – Least square regression

- **Loss:** As for the polynomial curve fitting, it is standard to consider the sum of square errors (assumption of Gaussian distributed errors)

$$E(\mathbf{W}) = \sum_{i=1}^N \|\mathbf{y}^i - \mathbf{d}^i\|_2^2 = \sum_{i=1}^N \|f(\mathbf{x}^i; \mathbf{W}) - \mathbf{d}^i\|_2^2$$

and look for \mathbf{W}^* such that $\nabla E(\mathbf{W}^*) = 0$.

Approximation – Least square regression

- **Loss:** As for the polynomial curve fitting, it is standard to consider the sum of square errors (assumption of Gaussian distributed errors)

$$E(\mathbf{W}) = \sum_{i=1}^N \|\mathbf{y}^i - \mathbf{d}^i\|_2^2 = \sum_{i=1}^N \|f(\mathbf{x}^i; \mathbf{W}) - \mathbf{d}^i\|_2^2$$

and look for \mathbf{W}^* such that $\nabla E(\mathbf{W}^*) = 0$.

- **Solution:** Provided the network has enough flexibility and the size of the training set grows to infinity

$$\mathbf{y}^* = f(\mathbf{x}; \mathbf{W}^*) = \underbrace{\mathbb{E}[\mathbf{d}|\mathbf{x}]}_{\text{posterior mean}} = \int d p(\mathbf{d}|\mathbf{x}) d \mathbf{d}$$

Approximation – Least square regression

- **Loss:** As for the polynomial curve fitting, it is standard to consider the sum of square errors (assumption of Gaussian distributed errors)

$$E(\mathbf{W}) = \sum_{i=1}^N \|\mathbf{y}^i - \mathbf{d}^i\|_2^2 = \sum_{i=1}^N \|f(\mathbf{x}^i; \mathbf{W}) - \mathbf{d}^i\|_2^2$$

and look for \mathbf{W}^* such that $\nabla E(\mathbf{W}^*) = 0$.

- **Solution:** Provided the network has enough flexibility and the size of the training set grows to infinity

$$\mathbf{y}^* = f(\mathbf{x}; \mathbf{W}^*) = \underbrace{\mathbb{E}[\mathbf{d}|\mathbf{x}]}_{\text{posterior mean}} = \int d p(\mathbf{d}|\mathbf{x}) d \mathbf{d}$$

- **Problem:** Inefficient for image classification. (Cf TD for a new cost function).

Training

We want to have $f(x, \theta)$ such that for all $(x, d) \in \mathcal{T}$,

$$f(x, \theta) \approx d$$

by minimizing a loss function $E(\theta)$.

Algorithm: (Stochastic) Gradient Descent (SGD)

- Initialize $\theta = (\mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2, \dots)$ randomly
- Repeat until convergence
 - For all $(x, d) \in \mathcal{T}$
 - Compute: $y = f(x, \theta)$
 - Update: $\theta \leftarrow \theta + \gamma \nabla_{\theta} E(\theta)$

An iterative algorithm trying to find a minimum of a real function.

Gradient descent

- Let F be a real function, coercive, and twice-differentiable such that:

$$\underbrace{\|\nabla^2 F(x)\|_2}_{\text{Hessian matrix of } F} \leq L, \quad \text{for some } L > 0.$$

- Then, whatever the initialization x^0 , if $0 < \gamma < 2/L$, the sequence

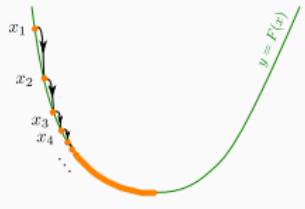
$$x^{(n+1)} = x^{(n)} - \underbrace{\gamma \nabla F(x^{(n)})}_{\text{direction of greatest descent}},$$

converges to a **stationary point** x^* (i.e., it cancels the gradient)

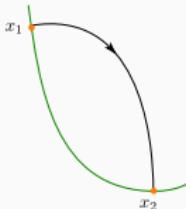
$$\nabla F(x^*) = 0 .$$

- The parameter γ is called the step size (or **learning rate** in ML field).
- A too small step size γ leads to slow convergence.

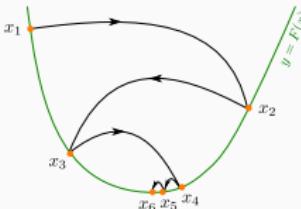
One dimension



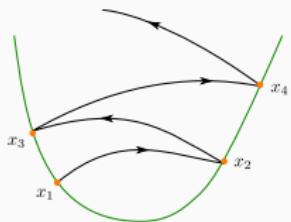
Small step size
Slow convergence



Good step size
Fast convergence

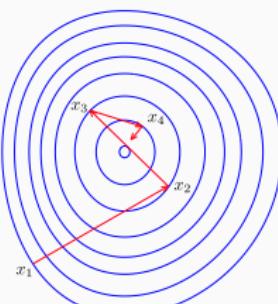
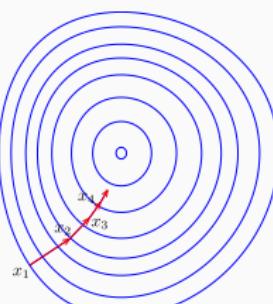
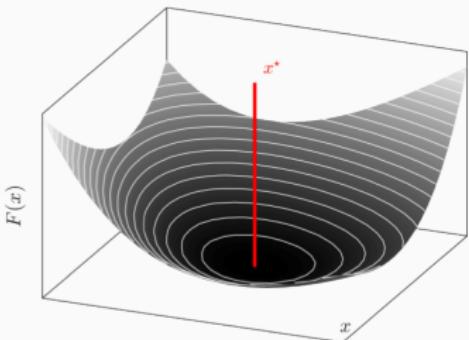


Large step size
Slow convergence



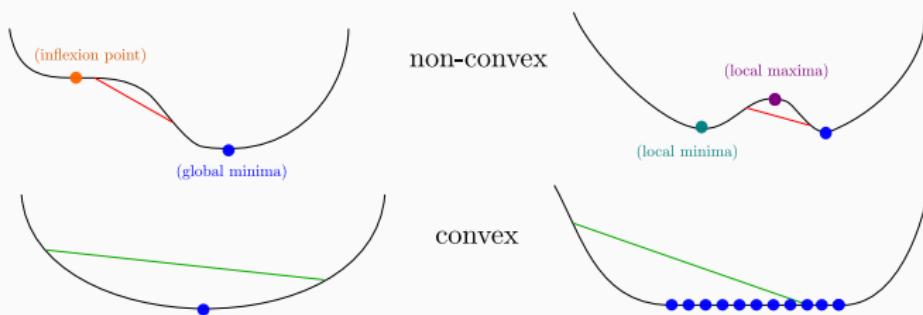
Too large step size
Divergence

Two dimensions



Non convexity in machine learning

But for neural network the cost is **not convex**...



Learning from examples

- ① Training set / examples:

$$\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$$

- ② Machine or model:

$$\mathbf{x} \rightarrow \underbrace{f(\mathbf{x}; \theta)}_{\text{function / algorithm}} \rightarrow \underbrace{\mathbf{y}}_{\text{prediction}}$$

θ : parameters of the model

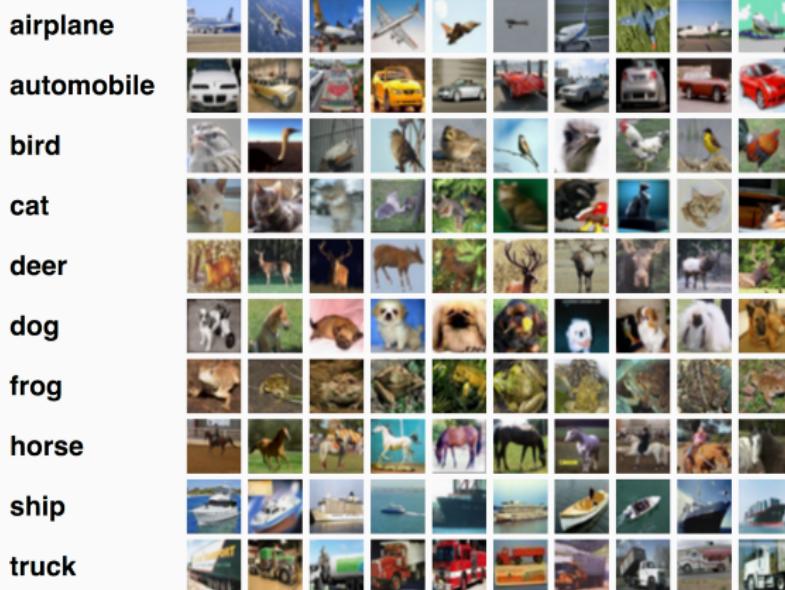
- ③ Loss, cost, objective function / energy:

$$\operatorname{argmin}_{\theta} E(\theta; \mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N)$$

- ④ Minimization of E via SGD.

- ⑤ Control of $f(\cdot, \theta)$ with a test dataset.

Classification example



Questions?

Next class: Backpropagation

Slides from Charles Deledalle

Sources, images courtesy and acknowledgment

K. Chatfield

P. Gallinari,

C. Hazırbaş

A. Horodniceanu

Y. LeCun

V. Lepetit

L. Masuch

A. Ng

M. Ranzato

Focus on classification: the multivariate logistic regression



Notations

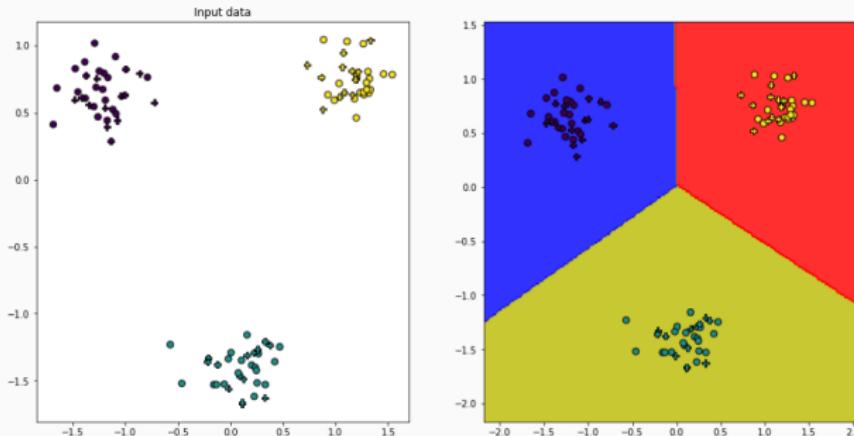
- **Goal:** Classify an object \boldsymbol{x} into one among K classes C_1, \dots, C_K .
- Training set: $\mathcal{T} = \{(\boldsymbol{x}_n, t_n), n = 1, \dots, N\}$, $t_n \in \{1, \dots, K\}$ encodes the class of \boldsymbol{x}_n .
- - Class 1: $t_n = 1$ if $\boldsymbol{x}_n \in C_1$,
 - Class 2: $t_n = 2$ if $\boldsymbol{x}_n \in C_2$,
 - ...
 - Class K: $t_n = K$ if $\boldsymbol{x}_n \in C_K$,

Feature transform

- We apply a feature transform $\varphi : \mathbb{R}^p \rightarrow \mathbb{R}^D$ to each x_n :

$$\varphi_n = \varphi(x_n), \quad n = 1, \dots, N.$$

- Depending on context it allows to increase ($D > p$) or decrease ($D < p$) the dimension in a way to favor class discrimination (e.g. PCA...).
- This is a non linear map that should make the classes linearly separable.



Multiclass classification – Multivariate logistic regression (aka, multinomial classification)

- **Goal:** Classify an object \mathbf{x} into one among K classes C_1, \dots, C_K .
- **How:** Estimate the coefficients \mathbf{W} of a multivariate function

$$\mathbf{y} = f(\mathbf{x}; \mathbf{W}) \in [0, 1]^K \quad \text{s.t.} \quad \sum_{k=1}^K y_k = 1.$$

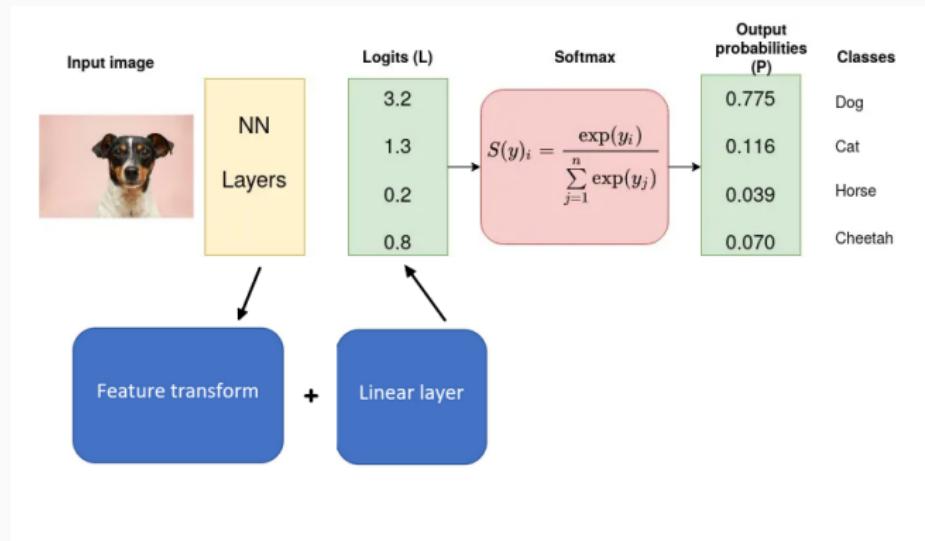
from training examples $\mathcal{T} = \{(\mathbf{x}^i, \mathbf{d}^i)\}$ where \mathbf{d}^i is the class of \mathbf{x}^i

- Class 1: $\mathbf{d}^i = 1$ if $\mathbf{x}^i \in C_1$
- Class 2: $\mathbf{d}^i = 2$ if $\mathbf{x}^i \in C_2$
- ...
- Class K: $\mathbf{d}^i = K$ if $\mathbf{x}^i \in C_K$

- $y_k = f(\mathbf{x}; \mathbf{W})$ is understood as the probability of $\mathbf{x} \in C_k$.

Machine learning with a NN

In real-life

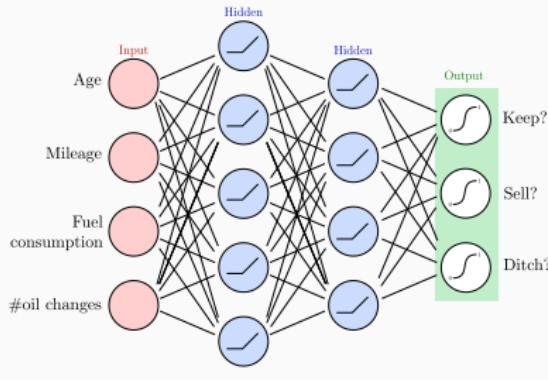


In real-life, all the network have to be trained. The feature transform is trained with its linear separation.

For this session, we will just train the last layer. We suppose that we have a good feature transform to have linearly separable data.

Multiclass classification – Multivariate logistic regression

- Typical architecture:



- Hidden layer:

$$\text{ReLU}(a) = \max(a, 0)$$

- Output layer:

$$\text{softmax}(\mathbf{a})_k = \frac{\exp(a_k)}{\sum_{\ell=1}^K \exp(a_\ell)}$$

- Softmax maps \mathbb{R}^K to the set of probability vectors $\{\mathbf{y} \in (0, 1)^K, \sum_{k=1}^K \mathbf{y}_k = 1\}$.
- Smooth version of winner-takes-all activation model (maxout).
- The final decision function is winner-takes-all

$$\operatorname{argmax}_k \text{softmax}(\mathbf{a}) = \operatorname{argmax}_k \mathbf{a}$$

Multiclass classification – Multivariate logistic regression

- **Loss:** it is standard to consider the **cross-entropy** for K classes (assumption of multinomial distributed data)

$$\begin{aligned} E(\mathbf{W}) &= - \sum_{i=1}^N \sum_{k=1}^K d_k^i \log y_k^i \quad \text{with} \quad \mathbf{y}^i = f(\mathbf{x}^i; \mathbf{W}) = \text{softmax}(\mathbf{a}) \in (0, 1)^K. \\ &= - \sum_{i=1}^N \left[a_{d^i} - \log \left(\sum_{k=1}^K \exp(a_k) \right) \right] \quad \text{with } d^i \text{ the class of } \mathbf{x}^i. \end{aligned}$$

and look for \mathbf{W}^* such that $\nabla E(\mathbf{W}^*) = 0$.

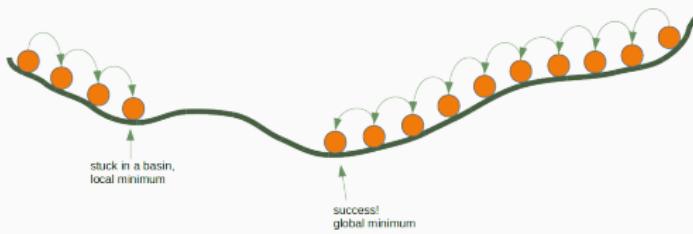
- **Solution:** Provided the network has enough flexibility and the size of the training set grows to infinity

$$y_k^* = f_k(\mathbf{x}; \mathbf{W}^*) = \underbrace{\mathbb{P}(C_k | \mathbf{x})}_{\text{posterior probability}}$$

The end

Questions ?

Next course: Backpropagation



Questions?

Next class: Backpropagation

Slides from Charles Deledalle

Sources, images courtesy and acknowledgment

K. Chatfield

P. Gallinari,

C. Hazırbaş

A. Horodniceanu

Y. LeCun

V. Lepetit

L. Masuch

A. Ng

M. Ranzato