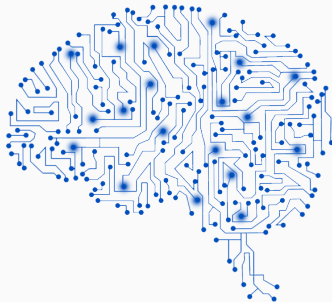


Introduction to Deep Learning

Course I – Introduction to Artificial Neural Networks

Bruno Galerne
2024-2025



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



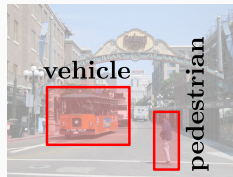
www.charles-deledalle.fr/

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

Computer Vision and Machine Learning



Image

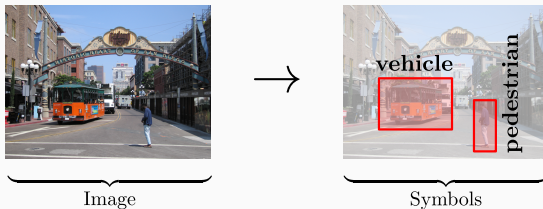


Symbols

Computer vision – Artificial Intelligence – Machine Learning

Definition (The British Machine Vision Association)

Computer vision (CV) is concerned with the automatic extraction, analysis and understanding of useful information from a single image or a sequence of images.



CV is a subfield of Artificial Intelligence.

Definition (Oxford dictionary)

Artificial Intelligence, *noun*: the theory and development of computer systems able to perform tasks normally requiring human intelligence, such as visual perception, speech recognition, decision-making, and translation.

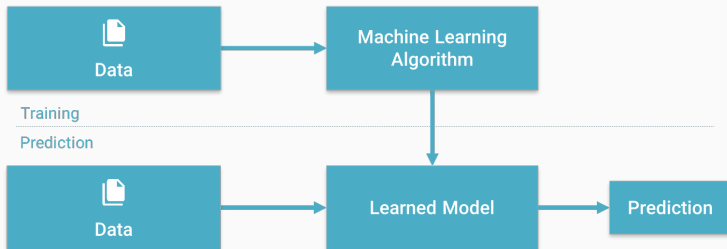
Computer vision – Artificial Intelligence – Machine Learning

CV is a subfield of AI, CV's new very best friend is **machine learning** (ML), ML is also a subfield of AI, but not all computer vision algorithms are ML.

Definition

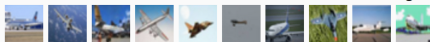
Machine Learning, *noun*: type of Artificial Intelligence that provides computers with the ability to **learn without being explicitly programmed**.

ML provides **various techniques** that can learn from and make predictions on data. Most of them follow the same general structure:

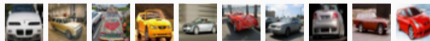


Computer vision – Image classification

airplane



automobile



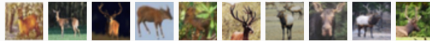
bird



cat



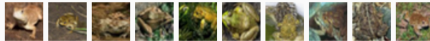
deer



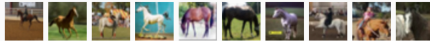
dog



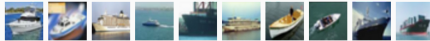
frog



horse



ship

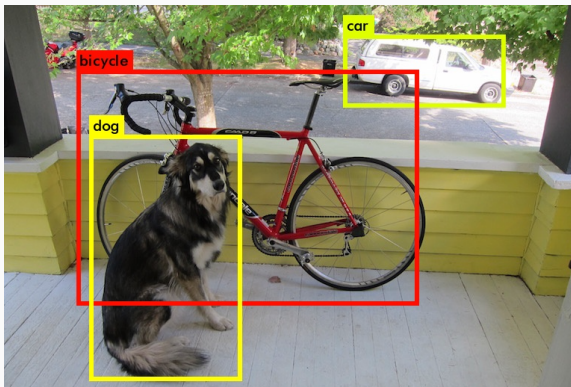


truck



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

Computer vision – Object detection



(Source: Joseph Redmon)

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

Computer vision – Image segmentation



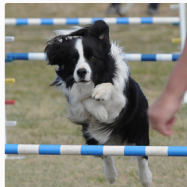
(Source: Abhijit Kundu)

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



"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."



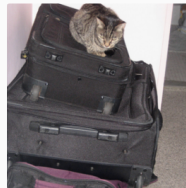
"little girl is eating piece of cake."



"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.

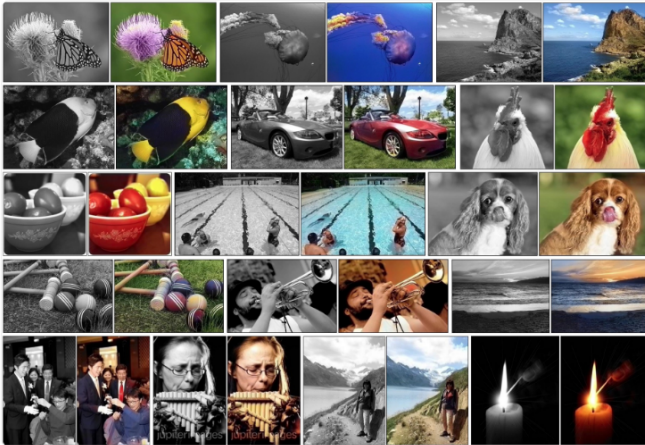
Computer vision – Depth estimation



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

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

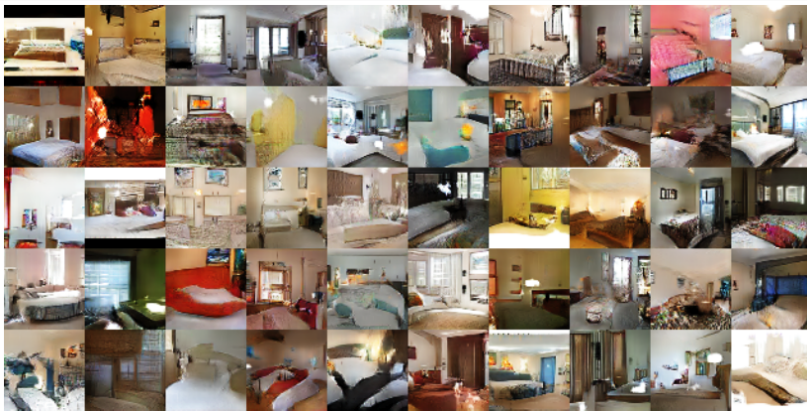
Image colorization



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

Goal: to add color to grayscale photographs.

Image generation



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

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

Image stylization

Synthesized Image

#NeuralDoodle



(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.

Learning from examples

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:

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

Learning from examples

Tools: $\left\{ \begin{array}{ll} \text{Data} & \leftrightarrow \text{Statistics} \\ \text{Loss} & \leftrightarrow \text{Optimization} \end{array} \right.$

Goal: to extract information from the training set

- relevant for the given task,
- relevant for other data of the same kind.

Terminology

Sample (Observation or Data): item to process (e.g., classify). *Example: an individual, a document, a picture, a sound, a video. . .*

Features (Input): set of distinct traits that can be used to describe each sample in a quantitative manner. Represented as a multi-dimensional vector usually denoted by x . *Example: size, weight, citizenship, . . .*

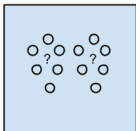
Training set: Set of data used to discover potentially predictive relationships.

Validation set: Set used to adjust the model hyperparameters.

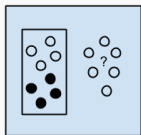
Testing set: Set used to assess the performance of a model.

Label (Output): The class or outcome assigned to a sample. The actual prediction is often denoted by y and the desired/targeted class by d or t . *Example: man/woman, wealth, education level, . . .*

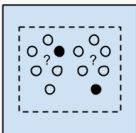
Learning approaches



Unsupervised Learning Algorithms



Supervised Learning Algorithms



Semi-supervised Learning Algorithms

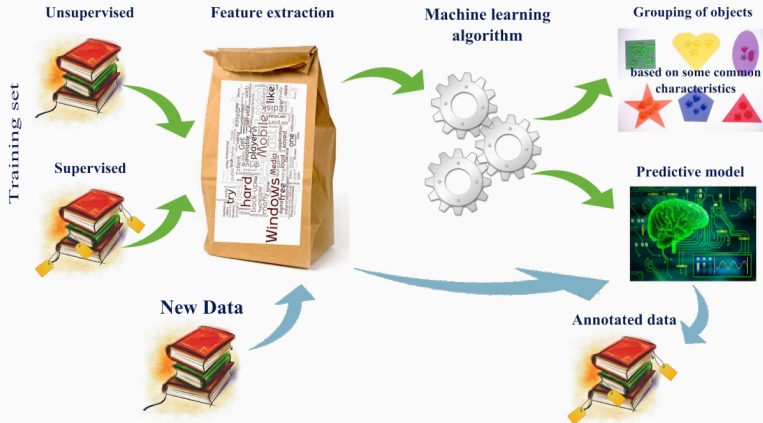
Unsupervised learning: Discovering patterns in unlabeled data. *Example: cluster similar documents based on the text content.*

Supervised learning: Learning with a labeled training set. *Example: email spam detector with training set of already labeled emails.*

Semisupervised learning: Learning with a small amount of labeled data and a large amount of unlabeled data. *Example: web content and protein sequence classifications.*

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

Machine learning workflow

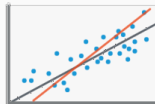


(Source: Michael Walker)

Problem types



Classification
(supervised – predictive)



Regression
(supervised – predictive)



Clustering
(unsupervised – descriptive)



Anomaly Detection
(unsupervised – descriptive)

(Source: Lucas Masuch)

What is deep learning?

- Part of the machine learning field of learning representations of data. Exceptionally effective at learning patterns.
- Utilizes learning algorithms that derive meaning out of data by using a hierarchy of multiple layers that mimic the neural networks of our brain.
- If you provide the system tons of information, it begins to understand it and respond in useful ways.
- Rebirth of artificial neural networks.

(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

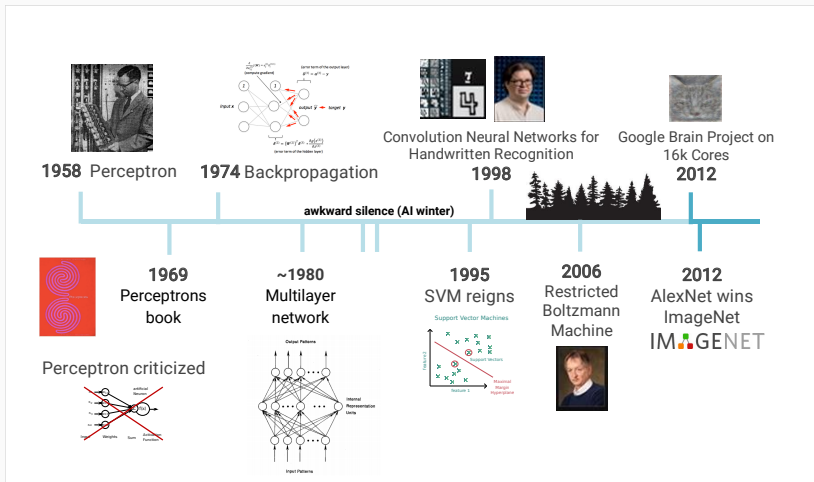
and many others...

- 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 developed by big actors: Facebook/Meta (PyTorch), Google (Tensorflow, Keras, JAX),...
- Success story for many different academic problems
 - Image processing
 - Computer vision
 - Speech recognition
 - Natural language processing
 - Translation
 - etc
- Today all industries wonder if AI/DL can improve their process.

Timeline of (deep) learning



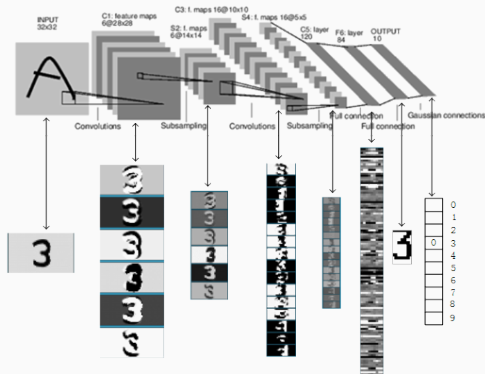
- ➊ Introduction to neural networks
- ➋ Convolutional neural networks for image classification
- ➌ Deep CNN for image classification, transfer learning
- ➍ Convolutional neural networks for image segmentation and image processing
- ➎ Deep generative models
- ➏ Transformers (if time allows)
- ➐ Diffusion models (if time allows)

Software: Python + PyTorch using Google Colab.

Remark: Focus on image processing and computer vision, but deep learning works for many other applications:

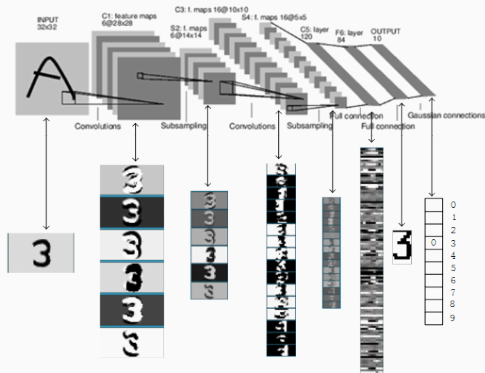
- Signal processing, speech recognition,...
- Text processing
- Graph processing (discrete geometry, social networks,...)
- Physics, chemistry,...

Neural networks for image classification



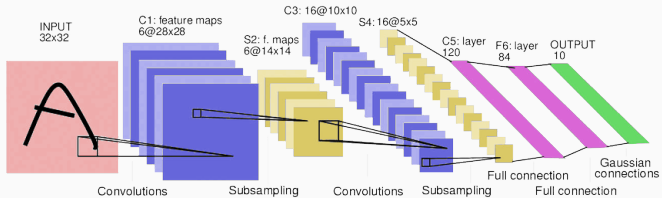
- **Goal:** Train a convolutional neural network for image classification

Neural networks for image classification

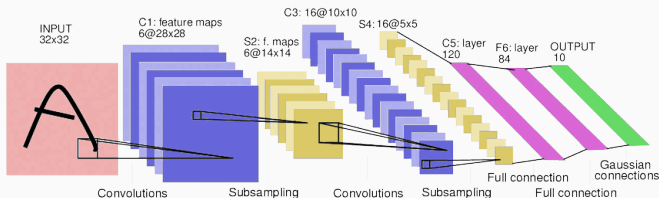


- **Goal:** Train a convolutional neural network for image classification
- **Goal:** Understand the training of a convolutional neural network for image classification

Understand the training of a convolutional neural network for image classification: A lot of notions: going backwards...

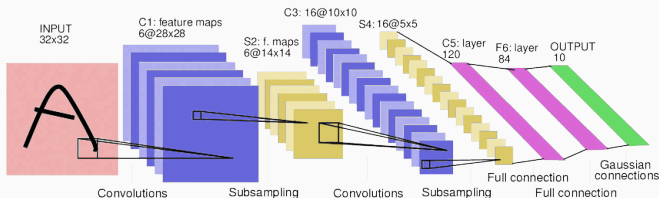


Understand the training of a convolutional neural network for image classification: A lot of notions: going backwards...



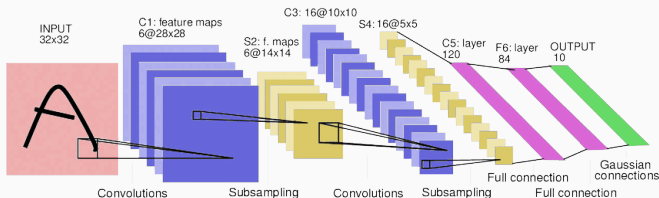
- **Convolutional neural networks:** Special neural networks for images that uses local convolutions (e.g. 3×3 filters) for the first layers.

Understand the training of a convolutional neural network for image classification: A lot of notions: going backwards...



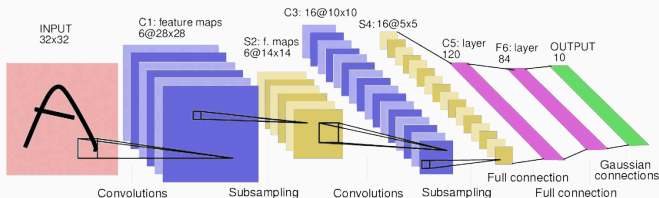
- **Convolutional neural networks:** Special neural networks for images that uses local convolutions (e.g. 3×3 filters) for the first layers.
- **Neural network:** A specific architecture to compute a classifier (or regression) having parameters=weights W to train at each layers.

Understand the training of a convolutional neural network for image classification: A lot of notions: going backwards...



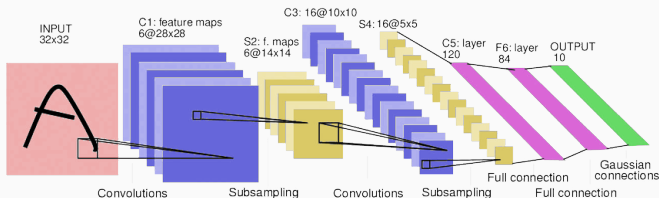
- **Convolutional neural networks:** Special neural networks for images that uses local convolutions (e.g. 3×3 filters) for the first layers.
- **Neural network:** A specific architecture to compute a classifier (or regression) having parameters=weights W to train at each layers.
- **Training** is done by optimizing a **classification loss** $L(W)$ on a training dataset: Typically this is a **linear classifier using cross-entropy**.

Understand the training of a convolutional neural network for image classification: A lot of notions: going backwards...



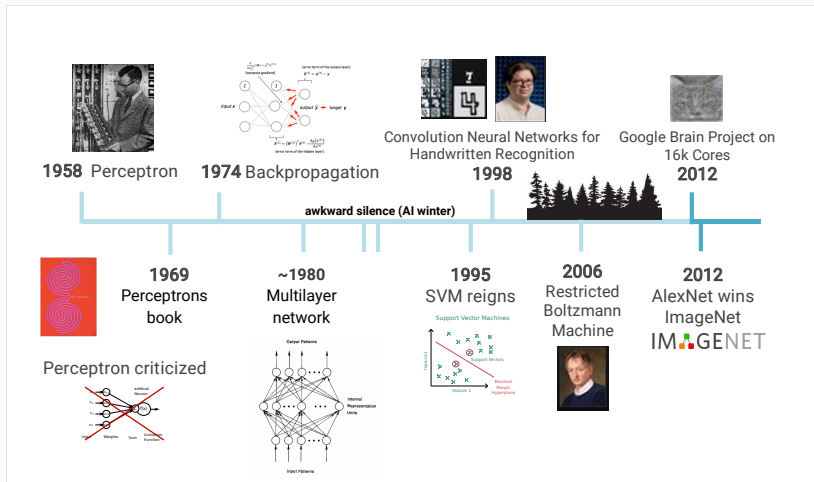
- **Convolutional neural networks:** Special neural networks for images that uses local convolutions (e.g. 3×3 filters) for the first layers.
- **Neural network:** A specific architecture to compute a classifier (or regression) having parameters=weights W to train at each layers.
- **Training** is done by optimizing a **classification loss** $L(W)$ on a training dataset: Typically this is a **linear classifier using cross-entropy**.
- The optimization of the classification loss is done using **stochastic gradient descent** on **batches of training data**.

Understand the training of a convolutional neural network for image classification: A lot of notions: going backwards...

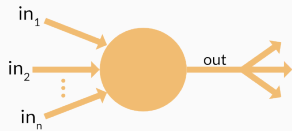
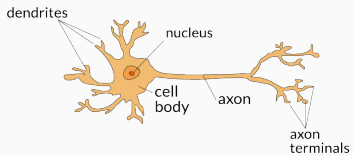


- **Convolutional neural networks:** Special neural networks for images that uses local convolutions (e.g. 3×3 filters) for the first layers.
- **Neural network:** A specific architecture to compute a classifier (or regression) having parameters=weights W to train at each layers.
- **Training** is done by optimizing a **classification loss** $L(W)$ on a training dataset: Typically this is a **linear classifier using cross-entropy**.
- The optimization of the classification loss is done using **stochastic gradient descent** on **batches of training data**.
- The gradient $\nabla L(W)$ is computed using **backpropagation**.

Timeline of (deep) learning



Perceptron



Perceptron



1958 Perceptron

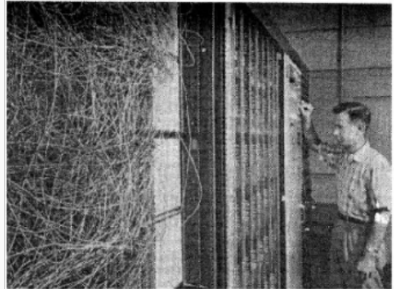
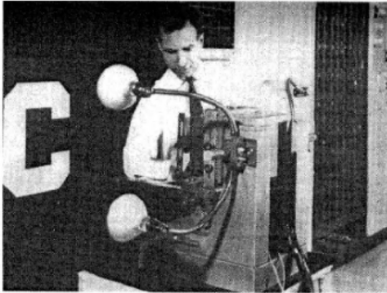


1969
Perceptrons
book

Perceptron criticized



Perceptron (Frank Rosenblatt, 1958)

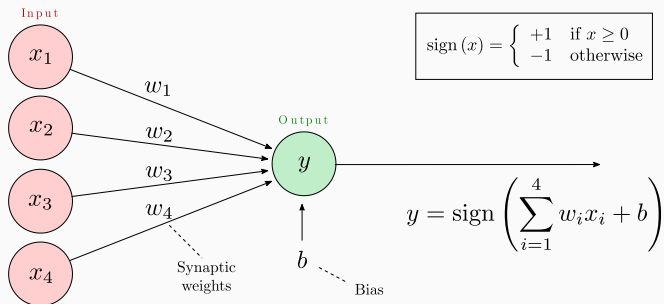


First binary classifier based on supervised learning (discrimination).

Foundation of modern artificial neural networks.

At that time: technological, scientific and philosophical challenges.

Representation of the Perceptron



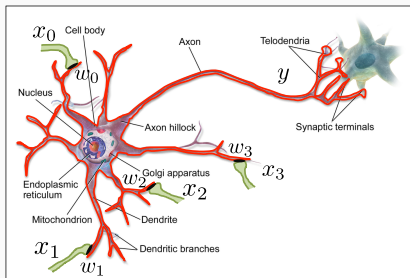
Parameters of the perceptron

- w_k : synaptic weights
 - b : bias
- } \leftarrow 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.



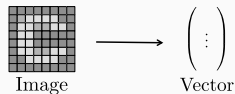
- When a neuron receives a stimulus with high enough voltage, it emits an **action potential** (aka, nerve impulse or spike). It is said to **fire**.
- The perceptron mimics this activation effect: it fires only when

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

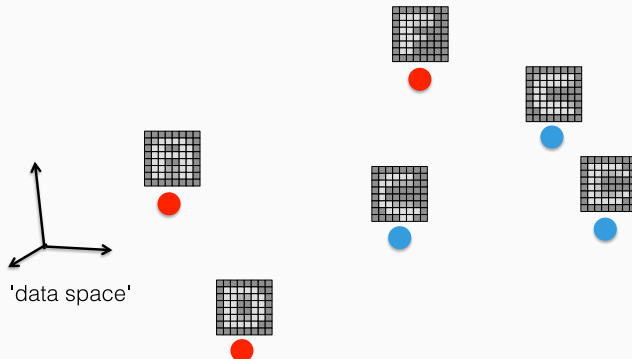
$$y = \underbrace{\text{sign}(w_0 x_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + b)}_{f(\mathbf{x}; \mathbf{w})} = \begin{cases} +1 & \text{for the first class} \\ -1 & \text{for the second class} \end{cases}$$

Machine learning – Perceptron – Principle

- 1 Data are represented as vectors:



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

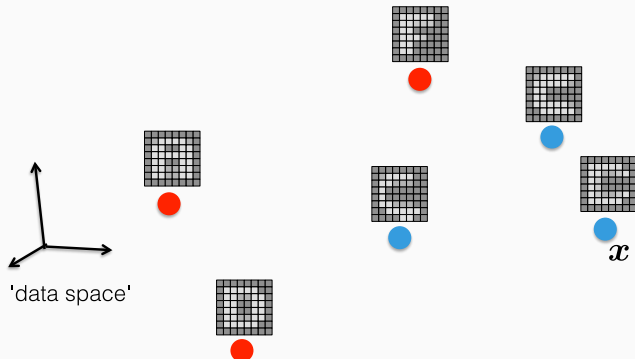


③ **Training:** find w and b so that:

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

Dot product:

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



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 .

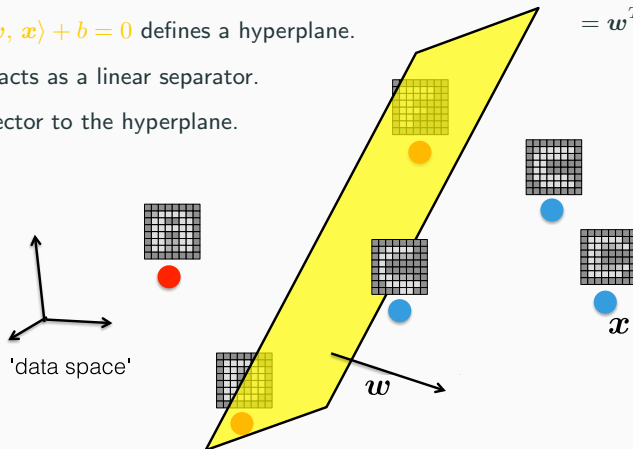
Dot product:

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

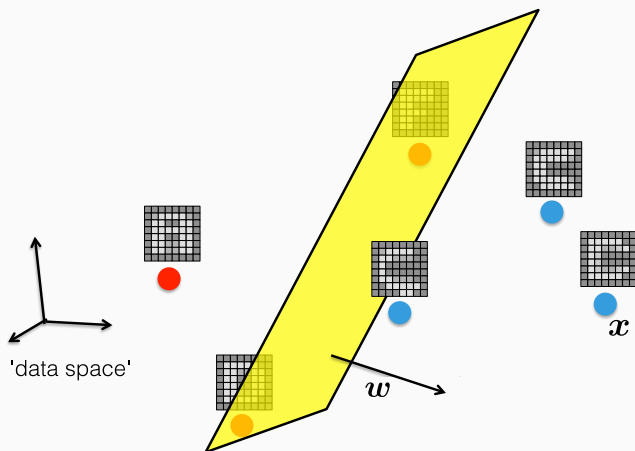
The equation $\langle w, x \rangle + b = 0$ defines a hyperplane.

The hyperplane acts as a linear separator.

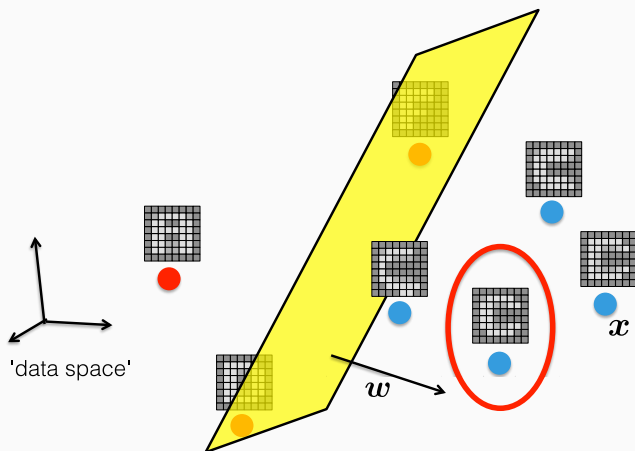
w is a normal vector to the hyperplane.



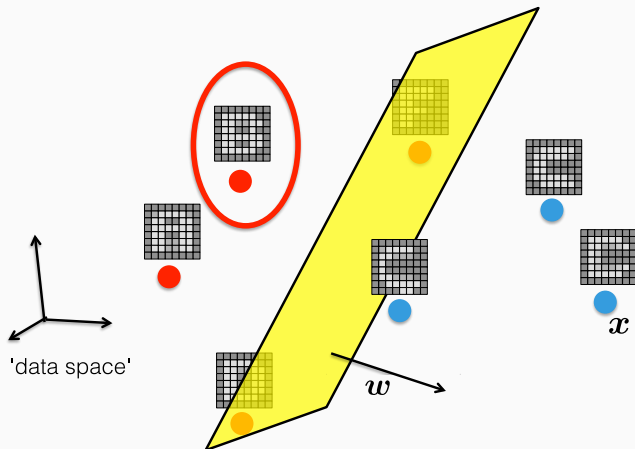
- ④ **Testing:** the perceptron can now classify new examples.



- ④ **Testing:** the perceptron can now classify new examples.
- A new example x is classified **positive** if $\langle w, x \rangle + b$ is **positive**,



- 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**.

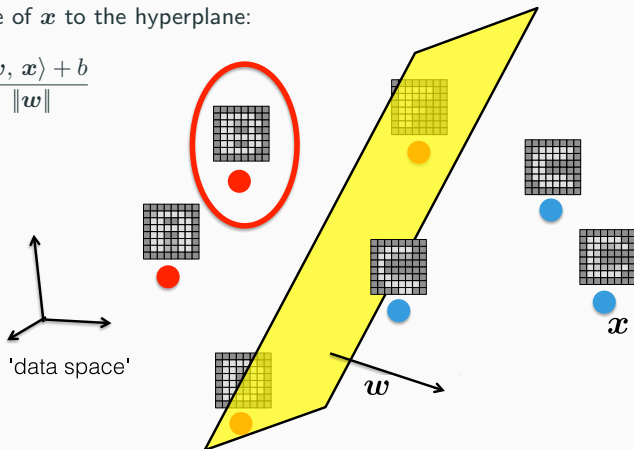


④ **Testing:** the perceptron can now classify new examples.

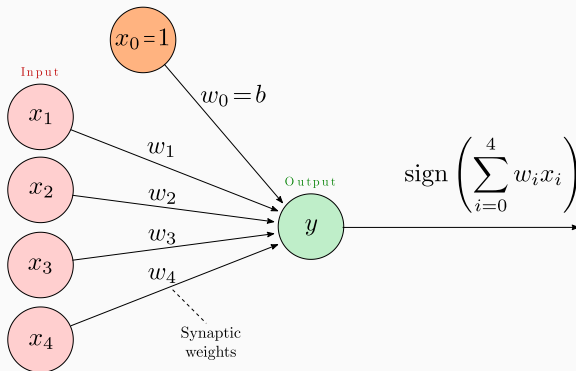
- A new example x is classified **positive** if $\langle w, x \rangle + b$ is **positive**,
- and **negative** if $\langle w, x \rangle + b$ is **negative**.

(signed) distance of x to the hyperplane:

$$r = \frac{\langle w, x \rangle + b}{\|w\|}$$



Alternative representation



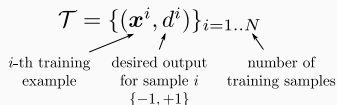
Use the zero-index to encode the bias as a synaptic weight.

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

Perceptron algorithm

Goal: find the vector of weights \mathbf{w} from a labeled training dataset \mathcal{T}

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



i -th training example desired output for sample i $\{-1, +1\}$ number of training samples

How: minimize classification errors

$$\min_{\mathbf{w}} E(\mathbf{w}) = - \sum_{\substack{(\mathbf{x}, d) \in \mathcal{T} \\ \text{st } y \neq d}} d \times \langle \mathbf{w}, \mathbf{x} \rangle = \sum_{(\mathbf{x}, d) \in \mathcal{T}} \max(-d \times \langle \mathbf{w}, \mathbf{x} \rangle, 0)$$

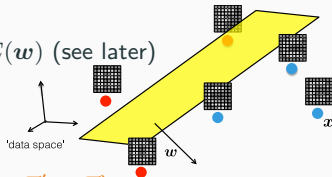
- penalize only misclassified samples ($y \neq d$) for which $d \times \langle \mathbf{w}, \mathbf{x} \rangle < 0$,
- zero if all samples are correctly classified.

Perceptron algorithm

- We assume that $\max(0, t)$ is derivable with derivative 1 if $t > 0$, 0 if $t \leq 0$.

Algorithm: (stochastic) gradient descent for $E(w)$ (see later)

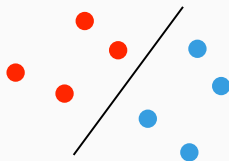
- Initialize w randomly
- Repeat until convergence
 - For all $(x, d) \in \mathcal{T}$ (or a random subset $\mathcal{T}' \subset \mathcal{T}$)
 - Compute: $y = \text{sign} \langle w, x \rangle$
 - If $y \neq d$:
Update: $w \leftarrow w + \gamma dx$



- Converges to some solution if the training data are linearly separable,
- But may pick any of many solutions of varying quality.
 \Rightarrow Poor generalization error, compared with SVM and logistic loss.

Perceptrons book (Minsky and Papert, 1969)

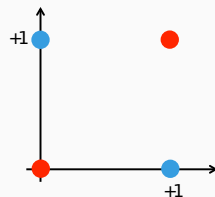
A perceptron can only classify data points that are linearly separable:



Linearly separable



Nonlinearly separable



The xor function

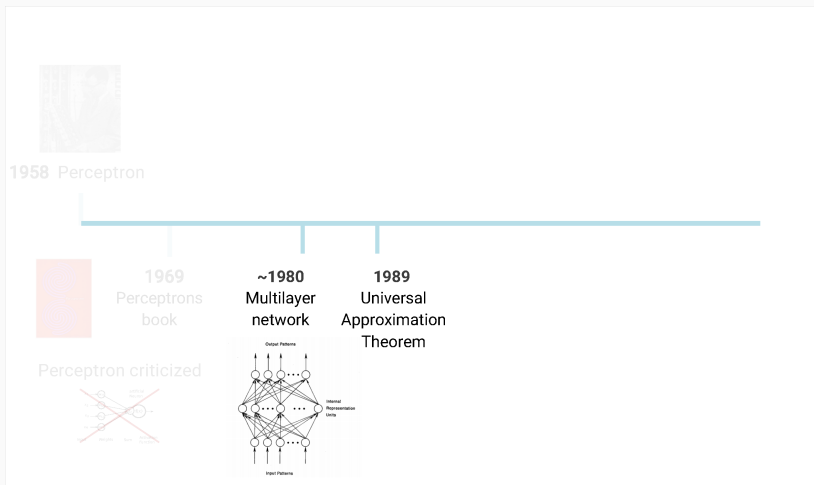
Seen by many as a justification to stop research on perceptrons.

(Source: Vincent Lepetit)

Artificial neural network



Artificial neural network

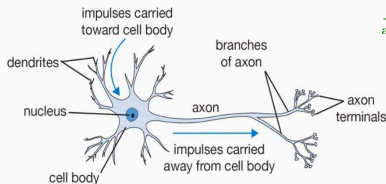


Artificial neural network

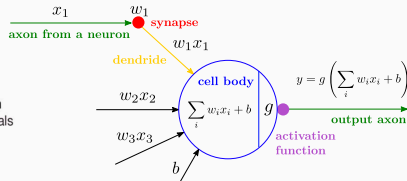


- Supervised learning method initially inspired by the behavior of the human brain.
- Consists of the inter-connection of several small units (just like in the human brain).
- Introduced in the late 50s, very popular in the 90s, reappeared in the 2010s with deep learning.
- Also referred to as **Multi-Layer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

Artificial neuron (McCulloch & Pitts, 1943)



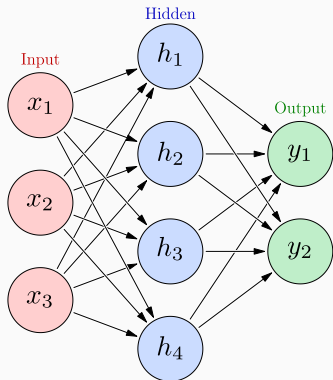
Biological neuron



Artificial neuron

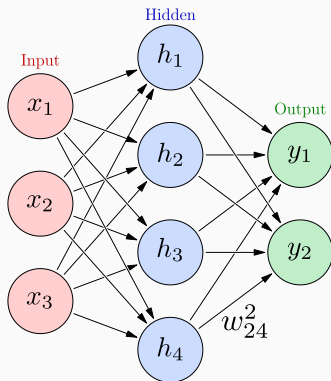
- An artificial neuron contains several incoming **weighted connections**, an outgoing connection and has a **nonlinear activation function** g .
- Neurons are **trained to filter and detect specific features** or patterns (e.g. edge, nose) by receiving weighted input, transforming it with the activation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



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

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

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

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

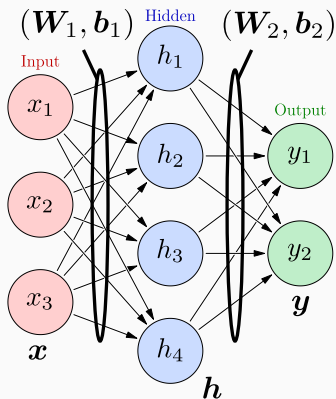
$$y_1 = g_2 (w_{11}^2 h_1 + w_{12}^2 h_2 + w_{13}^2 h_3 + w_{14}^2 h_4 + b_1^2)$$

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

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



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

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

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

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

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

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

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

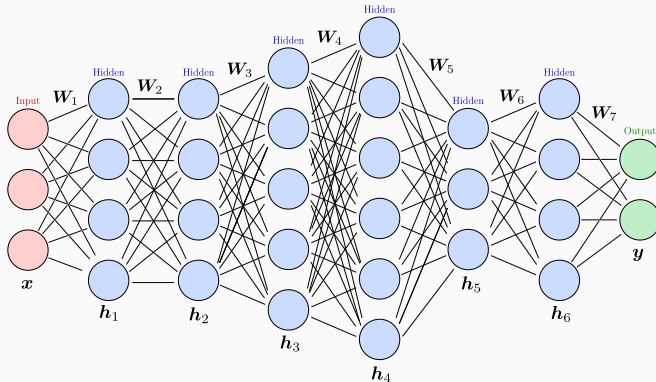
$$\mathbf{y} = g_2 (\mathbf{W}_2 \mathbf{h} + \mathbf{b}_2)$$

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: $g(a) = a$

$$\mathbf{y} = \mathbf{W}_L \mathbf{h}_{L-1} + \mathbf{b}_L$$

$$\mathbf{h}_{L-1} = \mathbf{W}_{L-1} \mathbf{h}_{L-2} + \mathbf{b}_{L-1}$$

$$\mathbf{y} = \mathbf{W}_L \mathbf{W}_{L-1} \mathbf{h}_{L-2} + \mathbf{W}_L \mathbf{b}_{L-1} + \mathbf{b}_L$$

$$\mathbf{y} = \mathbf{W}_L \dots \mathbf{W}_1 \mathbf{x} + \sum_{k=1}^{L-1} \mathbf{W}_L \dots \mathbf{W}_{k+1} \mathbf{b}_k + \mathbf{b}_L$$

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 Heaviside (aka, step) activation functions

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

Discontinuities in the hidden layers
make the optimization really difficult.

We prefer functions that are continuous and differentiable.

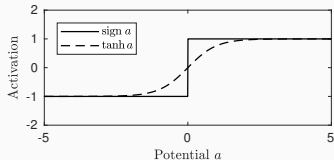
Activation functions

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 :

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

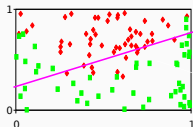
- Differentiable approximations of the sign and step functions, respectively.
- Act as threshold units for large values of $|a|$ and as linear for small values.

Sigmoidal units: logistic activation functions are used in binary classification (class C_1 vs C_2) as they can be **interpreted as posterior probabilities**:

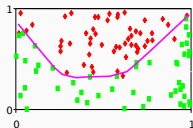
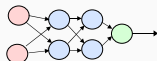
$$y = P(C_1|\mathbf{x}) \quad \text{and} \quad 1 - y = P(C_2|\mathbf{x})$$

The architecture of the network defines the shape of the separator

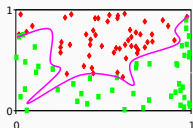
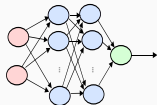
1 neuron



2+2+1 neurons



10+10+1 neurons



Separation
 $\{\mathbf{x} \text{ s.t. } P(C_1|\mathbf{x}) = P(C_2|\mathbf{x})\}$

Complexity/capacity of the
network

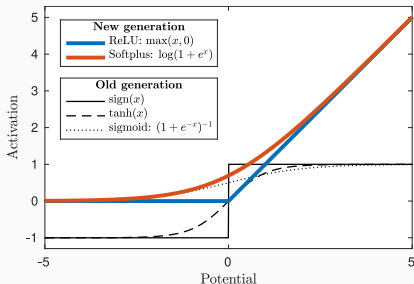
\Rightarrow

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

Activation functions

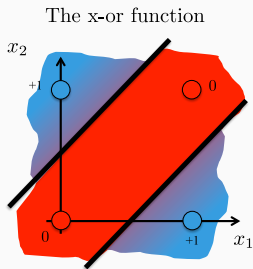
“Modern” units:

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



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.

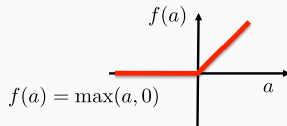
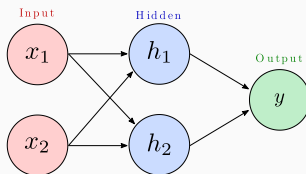
Neural networks solve non-linear separable problems



$$h = g(W_1 x + b_1)$$

$$y = \langle w_2, h \rangle + b_2$$

$$W_1 = \begin{pmatrix} +1 & -1 \\ -1 & +1 \end{pmatrix}, b_1 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, w_2 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, b_2 = 0$$



Tasks, architectures and loss functions



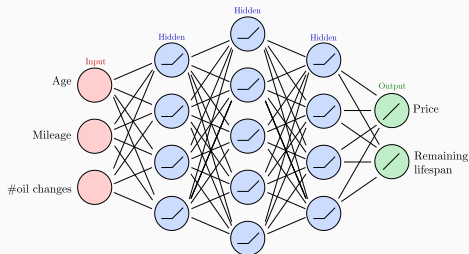
Approximation – Least square regression

- **Goal:** Predict a **real multivariate function**.
- **How:** estimate the coefficients \mathbf{W} of $\mathbf{y} = f(\mathbf{x}; \mathbf{W})$ from labeled training examples where labels are real vectors:

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

\swarrow \uparrow \searrow
 i -th training example desired output for sample i number of training samples

- **Typical architecture:**



- **Hidden layer:**

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

- **Linear output:**

$$g(a) = a$$

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

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

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

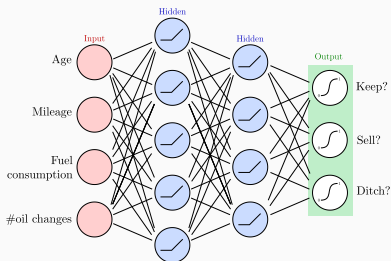
$$\mathbf{y} = f(\mathbf{x}; \mathbf{W}) \in [0, 1]^K \quad \text{s.t.} \quad \sum_{k=1}^K y_k = 1.$$

from training examples $\mathcal{T} = \{(\mathbf{x}^i, \mathbf{d}^i)\}$ where \mathbf{d}^i is a 1-of- K (one-hot) code

- Class 1: $\mathbf{d}^i = (1, 0, \dots, 0)^T$ if $\mathbf{x}^i \in C_1$
 - Class 2: $\mathbf{d}^i = (0, 1, \dots, 0)^T$ if $\mathbf{x}^i \in C_2$
 - ...
 - Class K : $\mathbf{d}^i = (0, 0, \dots, 1)^T$ if $\mathbf{x}^i \in C_K$
- $y_k = f(\mathbf{x}; \mathbf{W})$ is understood as the probability of $\mathbf{x} \in C_k$.
 - **Remark:** Do not use the class index k directly as a scalar label: The order of label is not informative.

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

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

Multiclass classification – Multivariate logistic regression

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

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

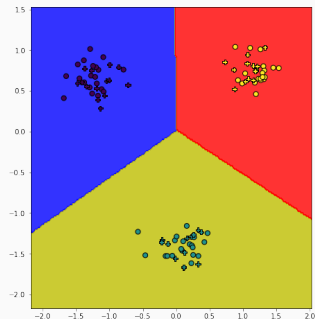
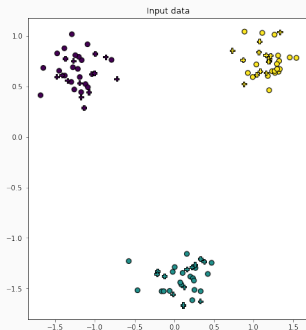
and look for \mathbf{W}^* such that $\nabla E(\mathbf{W}^*) = 0$.

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

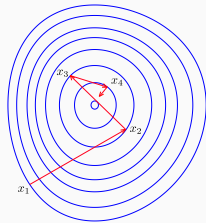
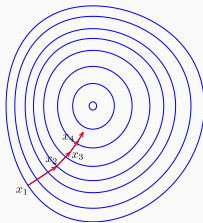
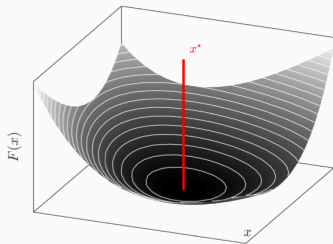
$$y_k^* = f_k(\mathbf{x}; \mathbf{W}^*) = \underbrace{\mathbb{P}(C_k | \mathbf{x})}_{\text{posterior probability}}$$

Multiclass classification – Multivariate logistic regression

- If there is just one layer (no hidden layer), we get linear separation for multiple classes: Each class region is the intersection of half-spaces regions.



Gradient descent



- The parameters of the neural networks are obtained by minimizing the training loss.
- This is done using (variants of) the standard optimization algorithm: **Gradient descent**.
- Recall that the **gradient** of function $F : \mathbb{R}^d \rightarrow \mathbb{R}$ is the vector of all its partial derivatives:

$$\nabla F(x) = \begin{pmatrix} \frac{\partial F}{\partial x_1}(x_1, \dots, x_d) \\ \frac{\partial F}{\partial x_2}(x_1, \dots, x_d) \\ \vdots \\ \frac{\partial F}{\partial x_d}(x_1, \dots, x_d) \end{pmatrix}$$

- It gives the steepest direction (local direction towards maximal increase of F values).
- Gradient descent consists in moving in the opposite direction $-\nabla F(x)$.

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

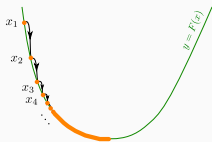
$$x^{(n+1)} = x^{(n)} - \underbrace{\gamma \nabla F(x^{(n)})}_{\text{direction of greatest descent}},$$

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

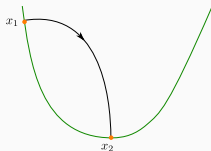
$$\nabla F(x^*) = 0.$$

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

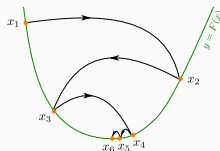
One dimension



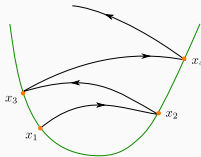
Small step size
Slow convergence



Good step size
Fast convergence

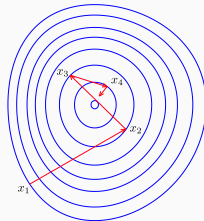
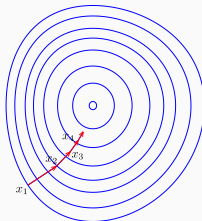
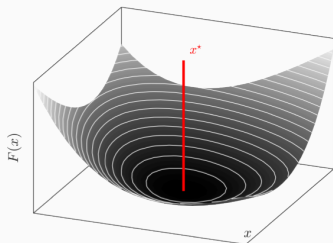


Large step size
Slow convergence



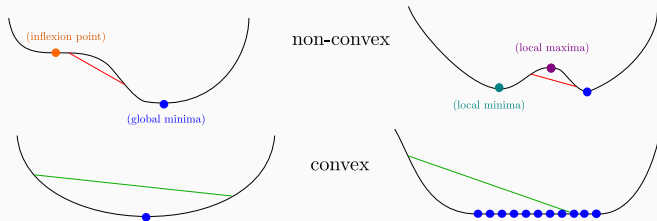
Too large step size
Divergence

Two dimensions

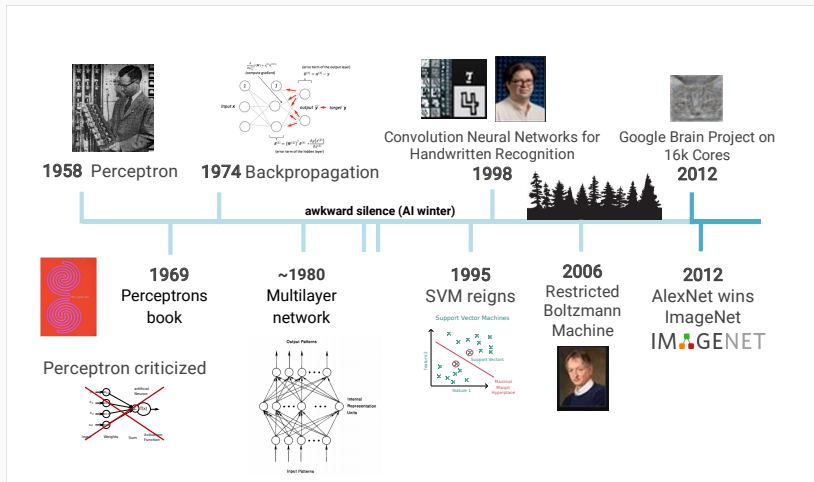


Non convexity in machine learning

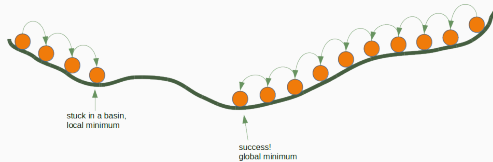
But for neural network the cost is **not convex**...



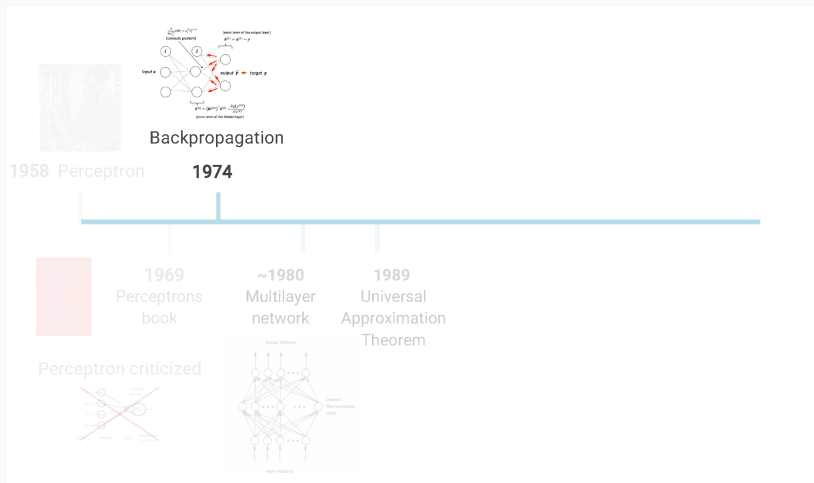
Timeline of (deep) learning



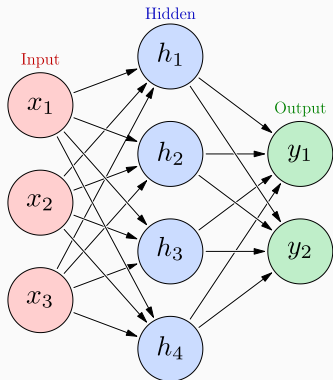
Backpropagation



Learning with backpropagation



Artificial neural network / Multilayer perceptron / NeuralNet



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

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

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

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

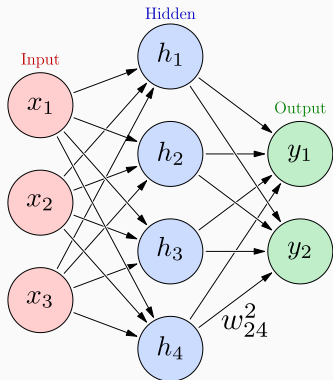
$$y_1 = g_2 (w_{11}^2 h_1 + w_{12}^2 h_2 + w_{13}^2 h_3 + w_{14}^2 h_4 + b_1^2)$$

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

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



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

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

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

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

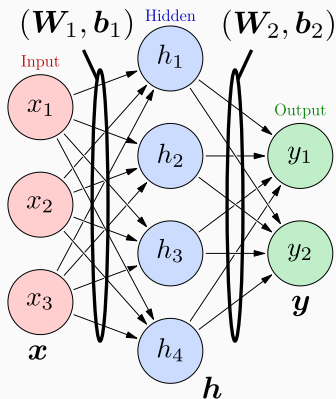
$$y_1 = g_2 (w_{11}^2 h_1 + w_{12}^2 h_2 + w_{13}^2 h_3 + w_{14}^2 h_4 + b_1^2)$$

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

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



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

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

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

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

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

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

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

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

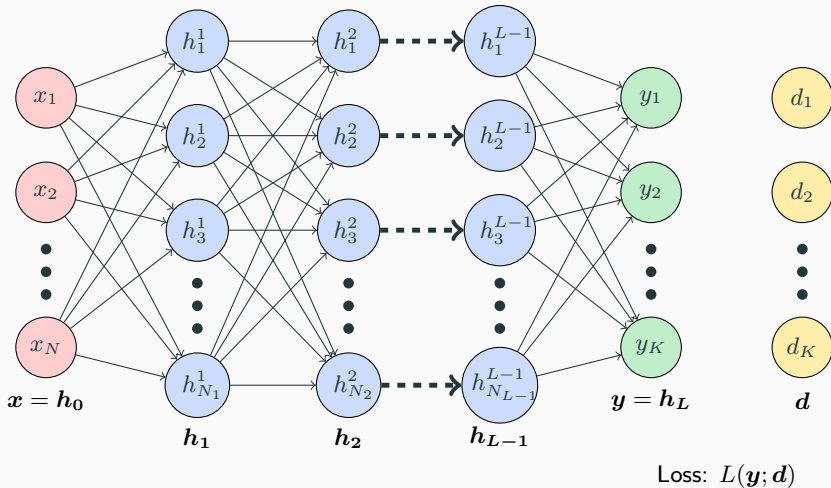
w_{ij}^k synaptic weight between previous node j and next node i at layer k .

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.

Feedforward Artificial Neural Network

Recall the feedforward structure



Input Layer

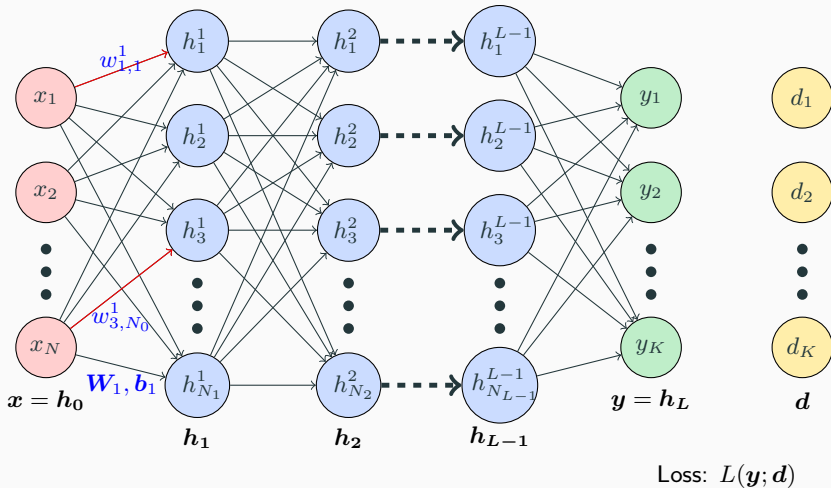
Hidden Layers

Output Layer

Label

Feedforward Artificial Neural Network

Recall the feedforward structure



Input Layer

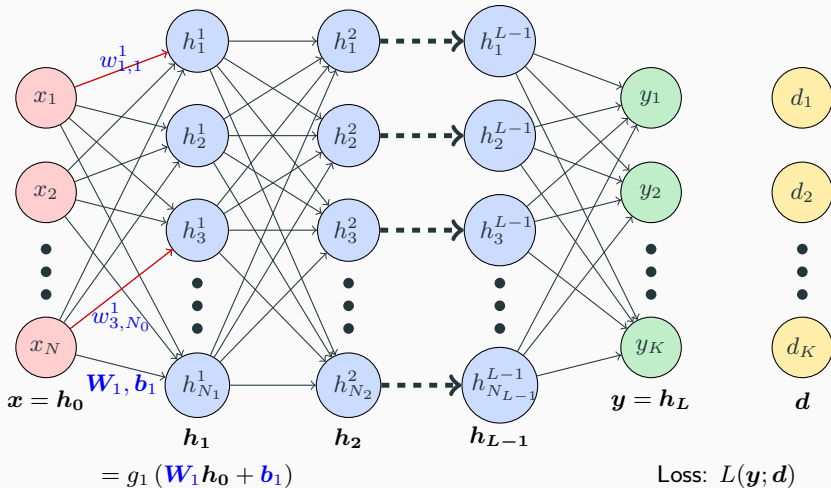
Hidden Layers

Output Layer

Label

Feedforward Artificial Neural Network

Recall the feedforward structure



Input Layer

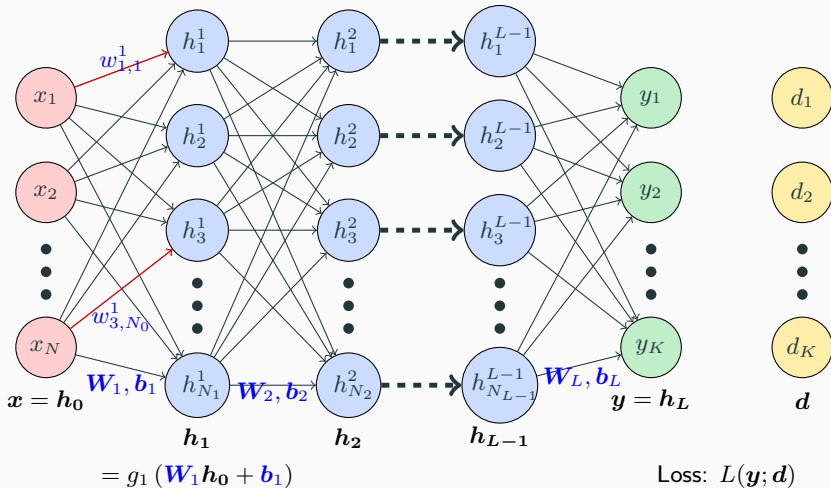
Hidden Layers

Output Layer

Label

Feedforward Artificial Neural Network

Recall the feedforward structure



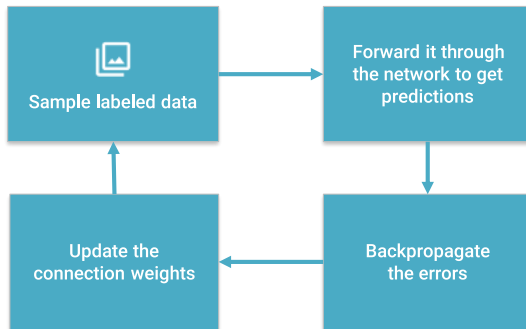
Input Layer

Hidden Layers

Output Layer

Label

Training process



Learns by generating an error signal that measures the difference between the predictions of the network and the desired values and then **using this error signal to change the weights** (or parameters) so that predictions get more accurate.

- The parameters of the neural network are

$$\mathbf{W} = (\mathbf{W}_1, \mathbf{b}_1, \mathbf{W}_2, \mathbf{b}_2, \dots, \mathbf{W}_L, \mathbf{b}_L)$$

- Training the network = minimizing the training loss $E(\mathbf{W})$

Objective: $\min_{\mathbf{W}} E(\mathbf{W})$ where $E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} L(\mathbf{y}^i; \mathbf{d}^i)$

$$\Rightarrow \nabla E(\mathbf{W}) = \left(\frac{\partial E(\mathbf{W})}{\partial \mathbf{W}_1} \quad \frac{\partial E(\mathbf{W})}{\partial \mathbf{b}_1} \quad \dots \quad \frac{\partial E(\mathbf{W})}{\partial \mathbf{W}_L} \quad \frac{\partial E(\mathbf{W})}{\partial \mathbf{b}_L} \right)^T = 0$$

- **Solution:** no closed-form solutions \Rightarrow use (stochastic) gradient descent.
- $\frac{\partial E(\mathbf{W})}{\partial \mathbf{W}_k}$ not really rigorous, we will use the notation

$$\nabla_{\mathbf{W}_k} E(\mathbf{W}) \quad \text{and} \quad \nabla_{\mathbf{b}_k} E(\mathbf{W}).$$

Minimizing training loss

For multilayer neural networks $\mathbf{W} \mapsto E(\mathbf{W})$ is non-convex

\Rightarrow No guarantee of convergence.

Even if convergence occurs, the solution depends on the initialization and the step size/learning rate γ .

Nevertheless, really good minima or saddle points are reached in practice by

$$\mathbf{W}^{t+1} \leftarrow \mathbf{W}^t - \gamma \nabla E(\mathbf{W}^t), \quad \gamma > 0$$

Gradient descent can be expressed coordinate by coordinate as:

$$w_{i,j}^{k,t+1} \leftarrow w_{i,j}^{k,t} - \gamma \frac{\partial E(\mathbf{W}^t)}{\partial w_{i,j}^k}$$

for all weights $w_{i,j}^k$ linking a node j to a node i in the next layer k .

\Rightarrow The algorithm to compute $\frac{\partial E(\mathbf{W})}{\partial w_{i,j}^k}$ for ANNs is called **backpropagation**.

- In practice we only use **stochastic gradient descent** with batch of training set.
- For some random small subset (e.g. batch) $\mathcal{S} \subset \mathcal{T}$, consider

$$E(\mathbf{W}; \mathcal{S}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{S}} L(\mathbf{y}^i; \mathbf{d}^i)$$

- Our **goal** is to compute the noisy gradient

$$\nabla_{\mathbf{W}_k} E(\mathbf{W}; \mathcal{S}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{S}} \nabla_{\mathbf{W}_k} L(\mathbf{y}^i; \mathbf{d}^i).$$

- In practice we only use **stochastic gradient descent** with batch of training set.
- For some random small subset (e.g. batch) $\mathcal{S} \subset \mathcal{T}$, consider

$$E(\mathbf{W}; \mathcal{S}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{S}} L(\mathbf{y}^i; \mathbf{d}^i)$$

- Our **goal** is to compute the noisy gradient

$$\nabla_{\mathbf{W}_k} E(\mathbf{W}; \mathcal{S}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{S}} \nabla_{\mathbf{W}_k} L(\mathbf{y}^i; \mathbf{d}^i).$$

- Why is this relevant to minimize $E(\mathbf{W}) = E(\mathbf{W}; \mathcal{T})$?

- **Stochastic gradient descent:** For some random small subset (e.g. batch) $\mathcal{S} \subset \mathcal{T}$, our **goal** is to compute the noisy gradient

$$\nabla_{\mathbf{W}_k} E(\mathbf{W}; \mathcal{S}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{S}} \nabla_{\mathbf{W}_k} L(\mathbf{y}^i; \mathbf{d}^i).$$

- **Unbiased approximation:** As soon as \mathcal{S} spans uniformly the whole training set \mathcal{T} ,

$$\begin{aligned} \mathbb{E}_{\mathcal{S}} (\nabla_{\mathbf{W}_k} E(\mathbf{W}; \mathcal{S})) &= \mathbb{E}_{\mathcal{S}} \left(\sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{S}} \nabla_{\mathbf{W}_k} L(\mathbf{y}^i; \mathbf{d}^i) \right) \\ &= \mathbb{E}_{\mathcal{S}} \left(\sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} \mathbf{1}_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{S}} \nabla_{\mathbf{W}_k} L(\mathbf{y}^i; \mathbf{d}^i) \right) \\ &= \frac{|\mathcal{S}|}{|\mathcal{T}|} \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} \nabla_{\mathbf{W}_k} L(\mathbf{y}^i; \mathbf{d}^i) = \frac{|\mathcal{S}|}{|\mathcal{T}|} \nabla_{\mathbf{W}_k} E(\mathbf{W}). \end{aligned}$$

- **Conclusion:** In expectation the noisy gradient is equal to the gradient using the whole training dataset (unbiased estimator).

Loss functions: Classical loss functions are:

For regression: $\mathbf{d}^i \in \mathbb{R}^K$

- Square error

$$E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} \frac{1}{2} \|\mathbf{y}^i - \mathbf{d}^i\|_2^2 = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} \frac{1}{2} \sum_k (y_k^i - d_k^i)^2$$

For multi-class classification: $d^i \in \{1, \dots, K\}$, coded by $\mathbf{d}^i \in \{0, 1\}^K$,

- Cross-entropy with softmax as the last layer

$$E(\mathbf{W}) = - \sum_{(\mathbf{x}^i, \mathbf{d}^i)} \sum_{k=1}^K d_k^i \log y_k^i \quad \text{with} \quad \mathbf{y}^i = f(\mathbf{x}^i; \mathbf{W}) = \text{softmax}(\mathbf{a}^i) \in (0, 1)^K.$$

- Cross-entropy with softmax included in loss (PyTorch convention):

$\mathbf{y}^i = \mathbf{a}^i$ is the output of the last linear layer:

$$E(\mathbf{W}) = - \sum_{(\mathbf{x}^i, \mathbf{d}^i)} \left[a_{d^i} - \log \left(\sum_{k=1}^K \exp(a_k) \right) \right] \quad \text{with } d^i \text{ the class of } \mathbf{x}^i.$$

- The loss functions are of the form

$$E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i)} L(\mathbf{y}^i; \mathbf{d}^i)$$

- By linearity,

$$\nabla E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i)} \nabla L(\mathbf{y}^i; \mathbf{d}^i)$$

- There the neural net output $\mathbf{y}^i = f(\mathbf{x}^i; \mathbf{W})$ is a function of the input data \mathbf{x}^i and the neural weights \mathbf{W} .
- We know the gradient of $L(\mathbf{y}^i; \mathbf{d}^i)$ with respect to the variable \mathbf{y}
 - Regression/Square error:

$$L(\mathbf{y}; \mathbf{d}) = \frac{1}{2} \|\mathbf{y} - \mathbf{d}\|_2^2 \quad \Rightarrow \quad \nabla_{\mathbf{y}} L(\mathbf{y}; \mathbf{d}) = \mathbf{y} - \mathbf{d}$$

- Multi-class classification/cross-entropy:

$$L(\mathbf{y}; \mathbf{d}) = -y_d + \log \left(\sum_{k=1}^K \exp(y_k) \right) \quad \Rightarrow \quad \nabla_{\mathbf{y}} L(\mathbf{y}; \mathbf{d}) = \text{softmax}(\mathbf{y}) - \mathbf{d}.$$

- The loss functions are of the form

$$E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i)} L(\mathbf{y}^i; \mathbf{d}^i)$$

- By linearity,

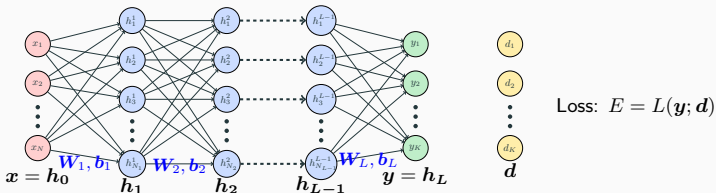
$$\nabla E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i)} \nabla L(\mathbf{y}^i; \mathbf{d}^i)$$

- There the neural net output $\mathbf{y}^i = f(\mathbf{x}^i; \mathbf{W})$ is a function of the input data \mathbf{x}^i and the neural weights \mathbf{W} .
- We know the gradient of $L(\mathbf{y}^i; \mathbf{d}^i)$ with respect to the variable \mathbf{y}
- We still need to compute

$$\nabla_{\mathbf{W}_k} L(\mathbf{y}; \mathbf{d}) \quad \text{and} \quad \nabla_{\mathbf{b}_k} L(\mathbf{y}; \mathbf{d}) \quad \text{for } k = 0, \dots, L.$$

- For simplicity above we will use the notation $E = L(\mathbf{y}; \mathbf{d})$, that is considering only one point.

ANN – Backpropagation



Forward pass

Initialization:

$$h_0 = x$$

for layer $k = 1$ **to** L **do**

Linear unit:

$$a_k = W_k h_{k-1} + b_k$$

Componentwise non-linear activation:

$$h_k = g_k(a_k)$$

end

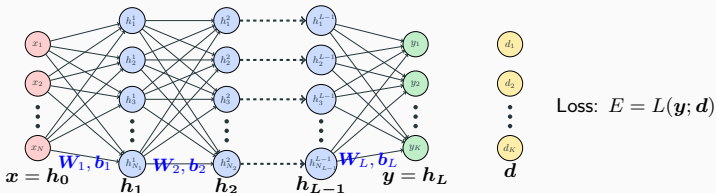
Output layer:

$$y = h_L$$

Compute loss:

$$E = L(y; d)$$

ANN – Backpropagation



Forward pass

Initialization:

$$\mathbf{h}_0 = \mathbf{x}$$

for layer $k = 1$ **to** L **do**

Linear unit:

$$\mathbf{a}_k = \mathbf{W}_k \mathbf{h}_{k-1} + \mathbf{b}_k$$

Componentwise non-linear activation:

$$\mathbf{h}_k = g_k(\mathbf{a}_k)$$

end

Output layer:

$$\mathbf{y} = \mathbf{h}_L$$

Compute loss:

$$E = L(\mathbf{y}; \mathbf{d})$$

Backward pass

Goal: Compute the gradient with respect to all parameters

$$\frac{\partial E}{\partial w_{i,j}^k} = ? \quad \frac{\partial E}{\partial b_i^k} = ?$$

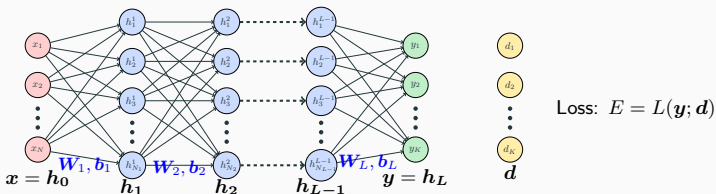
for all

$$k \in \{1, \dots, L\},$$

$$i \in \{1, \dots, N_k\},$$

$$j \in \{1, \dots, N_{k-1}\}.$$

ANN – Backpropagation

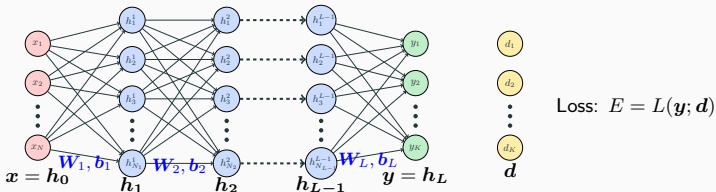


Going backward

- We know how to compute the loss function and its gradient:

$$\nabla_{h_L} E = \nabla L(y; d)$$

ANN – Backpropagation



Gradient with respect to last linear unit output a_L

$$h_L = g_L(a_L)$$

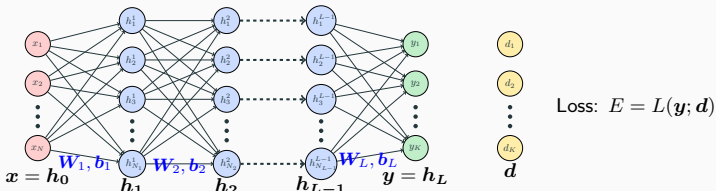
That is for all $i \in \{1, \dots, N_L\}$, $h_i^L = g_L(a_i^L)$. By the chain rule,

$$\frac{\partial E}{\partial a_i^L} = \frac{\partial E}{\partial h_i^L} \frac{\partial h_i^L}{\partial a_i^L} = [\nabla_{h_L} E]_i g'_L(a_i^L)$$

Vector formula: $\nabla_{a_L} E = \nabla_{h_L} E \odot g'_L(a_L)$

where \odot is the componentwise product between vectors, ie Hadamard product.

ANN – Backpropagation



Gradient with respect to bias of last linear unit b_L

$$a_L = W_L h_{L-1} + b_L$$

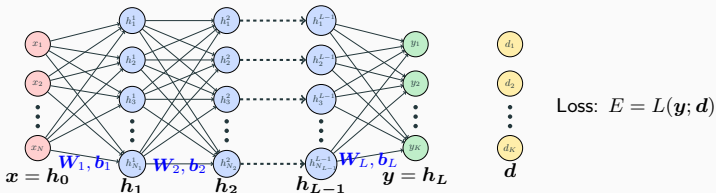
That is for all $i \in \{1, \dots, N_L\}$, $a_i^L = \sum_{j=1}^{N_{L-1}} w_{i,j}^L h_j^{L-1} + b_i^L$.

By the chain rule, for all $i \in \{1, \dots, N_L\}$,

$$\frac{\partial E}{\partial b_i^L} = \frac{\partial E}{\partial a_i^L} \underbrace{\frac{\partial a_i^L}{\partial b_i^L}}_{=1} = \frac{\partial E}{\partial a_i^L} = [\nabla_{a_L} E]_i$$

Vector formula: $\nabla_{b_L} E = \nabla_{a_L} E$

ANN – Backpropagation



Gradient with respect to weights of last linear unit W_L

$$a_L = W_L h_{L-1} + b_L$$

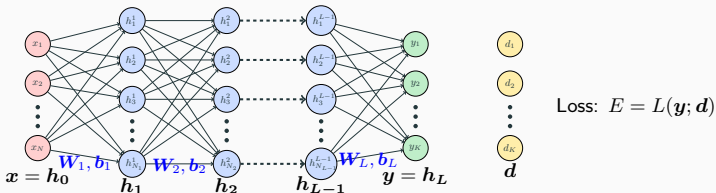
That is for all $i \in \{1, \dots, N_L\}$, $a_i^L = \sum_{j=1}^{N_{L-1}} w_{i,j}^L h_j^{L-1} + b_i^L$.

By the chain rule, for all $i \in \{1, \dots, N_L\}$ and $j \in \{1, \dots, N_{L-1}\}$

$$\frac{\partial E}{\partial w_{i,j}^L} = \frac{\partial E}{\partial a_i^L} \underbrace{\frac{\partial a_i^L}{\partial w_{i,j}^L}}_{=h_j^{L-1}} = \frac{\partial E}{\partial a_i^L} h_j^{L-1} = [\nabla_{a_L} E]_i [h_{L-1}]_j$$

Matrix formula: $\nabla_{W_L} E = \nabla_{a_L} E h_{L-1}^T$

ANN – Backpropagation

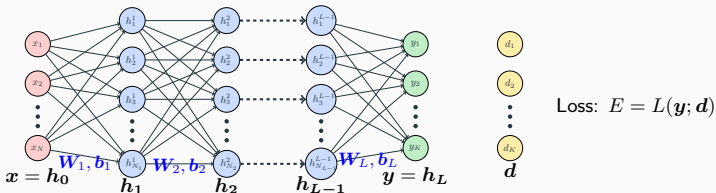


Gradients for last layer parameters

Given the gradient with respect to the output layer $\nabla_{h_L} E$, so far we can compute:

- $\nabla_{a_L} E = \nabla_{h_L} E \odot g'_L(a_L)$
- $\nabla_{b_L} E = \nabla_{a_L} E$
- $\nabla_{W_L} E = \nabla_{a_L} E h_{L-1}^T$

ANN – Backpropagation



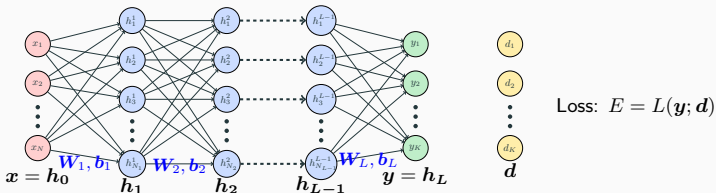
Gradients for last layer parameters

Given the gradient with respect to the output layer $\nabla_{\mathbf{h}_L} E$, so far we can compute:

- $\nabla_{\mathbf{a}_L} E = \nabla_{\mathbf{h}_L} E \odot g'_L(\mathbf{a}_L)$
- $\nabla_{\mathbf{b}_L} E = \nabla_{\mathbf{a}_L} E$
- $\nabla_{\mathbf{W}_L} E = \nabla_{\mathbf{a}_L} E \mathbf{h}_{L-1}^T$

How can we compute the gradients for the parameters of layer $L - 1$?

ANN – Backpropagation



Gradients for last layer parameters

Given the gradient with respect to the output layer $\nabla_{\mathbf{h}_L} E$, so far we can compute:

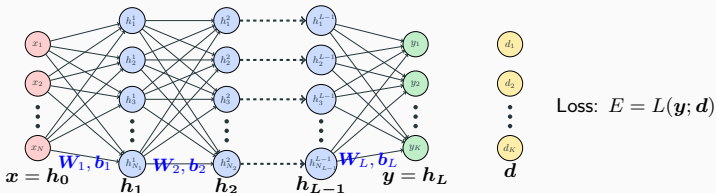
- $\nabla_{\mathbf{a}_L} E = \nabla_{\mathbf{h}_L} E \odot g'_L(\mathbf{a}_L)$
- $\nabla_{\mathbf{b}_L} E = \nabla_{\mathbf{a}_L} E$
- $\nabla_{\mathbf{W}_L} E = \nabla_{\mathbf{a}_L} E \mathbf{h}_{L-1}^T$

How can we compute the gradients for the parameters of layer $L - 1$?

We need the expression of the gradient with respect to the last but one hidden layer \mathbf{h}_{L-1} ... and then the same formulas apply!

$$\nabla_{\mathbf{h}_{L-1}} E = ?$$

ANN – Backpropagation

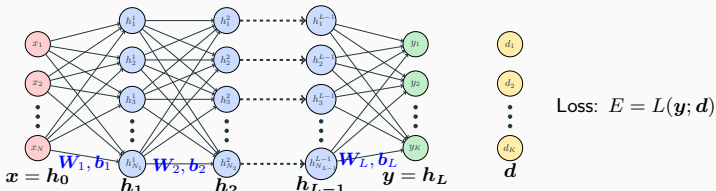


Gradient with respect to the last but one hidden layer h_{L-1}

Here, even to compute the scalar partial derivative $\frac{\partial E}{\partial h_j^{L-1}}$, we need to use differential calculus for multivariate functions since h_j^{L-1} appears in each component of \mathbf{a}_L :

$$\text{For all } i \in \{1, \dots, N_L\}, a_i^L = \sum_{j=1}^{N_{L-1}} w_{i,j}^L h_j^{L-1} + b_i^L.$$

ANN – Backpropagation



Gradient with respect to the last but one hidden layer h_{L-1}

Let us recall the derivative rule for composition with affine maps:

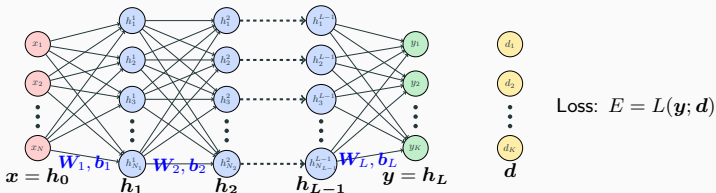
$$\text{For } \varphi(\mathbf{x}) = f(\mathbf{A}\mathbf{x} + \mathbf{b}) \quad \text{one has} \quad \nabla \varphi(\mathbf{x}) = \mathbf{A}^T \nabla f(\mathbf{A}\mathbf{x} + \mathbf{b}).$$

Using the decomposition

$$\begin{aligned} \mathbb{R}^{N_{L-1}} &\rightarrow \mathbb{R}^{N_L} \rightarrow \mathbb{R} \\ \mathbf{h}_{L-1} &\mapsto \mathbf{a}_L = \mathbf{W}_L \mathbf{h}_{L-1} + \mathbf{b}_L \mapsto E \end{aligned}$$

$$\text{Vector formula: } \nabla_{\mathbf{h}_{L-1}} E = \mathbf{W}_L^T \nabla_{\mathbf{a}_L} E$$

ANN – Backpropagation



Forward pass

Initialization:

$$\mathbf{h}_0 = \mathbf{x}$$

for layer $k = 1$ **to** L **do**

Linear unit:

$$\mathbf{a}_k = \mathbf{W}_k \mathbf{h}_{k-1} + \mathbf{b}_k$$

Componentwise non-linear activation:

$$\mathbf{h}_k = g_k(\mathbf{a}_k)$$

end

Output layer:

$$\mathbf{y} = \mathbf{h}_L$$

Compute loss:

$$E = L(\mathbf{y}; \mathbf{d})$$

Backward pass

Initialization: Gradient of output layer:

$$\nabla_{\mathbf{h}_L} E = \nabla L(\mathbf{y}; \mathbf{d})$$

for layer $k = L$ **to** 1 **do**

Componentwise gain of error:

$$\delta_k = \nabla_{\mathbf{a}_k} E = \nabla_{\mathbf{h}_k} E \odot g'_k(\mathbf{a}_k)$$

Gradient of layer bias:

$$\nabla_{\mathbf{b}_k} E = \delta_k$$

Gradient of weights:

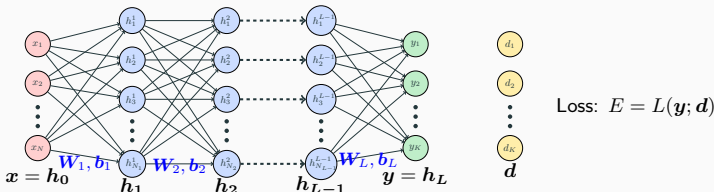
$$\nabla_{\mathbf{W}_k} E = \delta_k \mathbf{h}_{k-1}^T$$

Gradient of previous hidden layer:

$$\nabla_{\mathbf{h}_{k-1}} E = \mathbf{W}_k^T \delta_k$$

end

ANN – Backpropagation



Forward pass

Initialization:

$$\mathbf{h}_0 = \mathbf{x}$$

for layer $k = 1$ **to** L **do**

Linear unit:

$$\mathbf{a}_k = \mathbf{W}_k \mathbf{h}_{k-1} + \mathbf{b}_k \text{ (stored)}$$

Componentwise non-linear activation:

$$\mathbf{h}_k = g_k(\mathbf{a}_k) \text{ (stored)}$$

end

Output layer:

$$\mathbf{y} = \mathbf{h}_L$$

Compute loss:

$$E = L(\mathbf{y}; \mathbf{d})$$

Backward pass

Initialization: Gradient of output layer:

$$\nabla_{\mathbf{h}_L} E = \nabla L(\mathbf{y}; \mathbf{d})$$

for layer $k = L$ **to** 1 **do**

Componentwise gain of error:

$$\delta_k = \nabla_{\mathbf{a}_k} E = \nabla_{\mathbf{h}_k} E \odot g'_k(\mathbf{a}_k)$$

Gradient of layer bias:

$$\nabla_{\mathbf{b}_k} E = \delta_k$$

Gradient of weights:

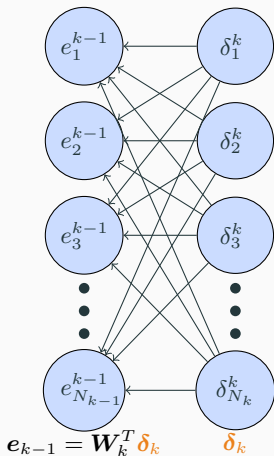
$$\nabla_{\mathbf{W}_k} E = \delta_k \mathbf{h}_{k-1}^T$$

Gradient of previous hidden layer:

$$\nabla_{\mathbf{h}_{k-1}} E = \mathbf{W}_k^T \delta_k$$

end

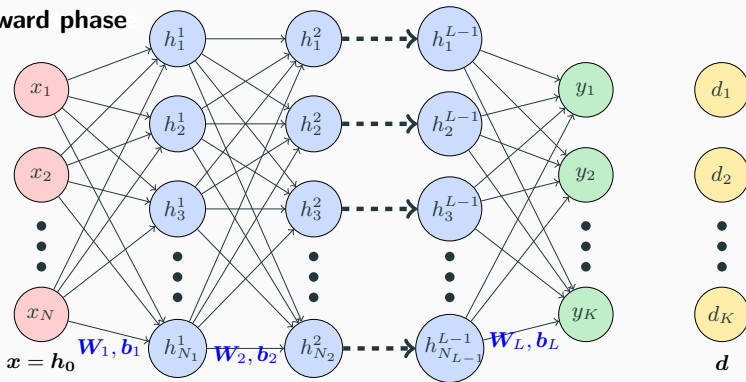
Error backpropagation



- Gradient of previous hidden layer:
$$e_{k-1} = \nabla_{h_{k-1}} E = \mathbf{W}_k^T \delta_k$$
- Multiplying by \mathbf{W}_k^T corresponds to passing to the linear layer in reverse order.
- The error is backpropagated layer by layer to compute the gradient with respect to each layer parameters.

Error backpropagation

Forward phase



Input Layer

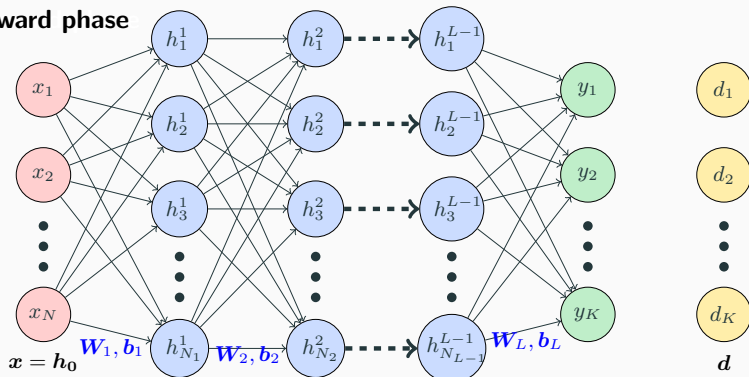
Hidden Layers

Output Layer

Label

Error backpropagation

Forward phase



Input Layer

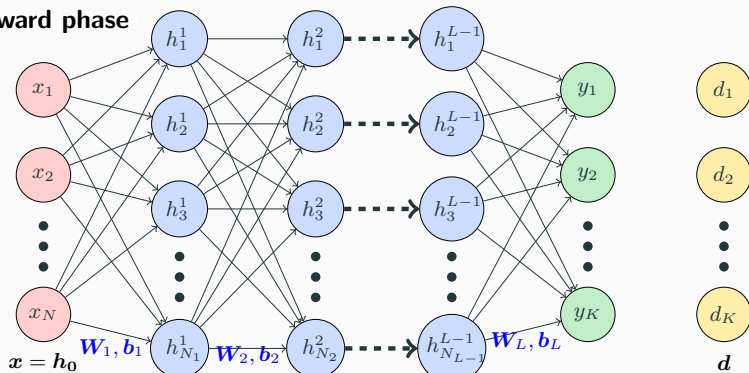
Hidden Layers

Output Layer

Label

Error backpropagation

Forward phase



Input Layer

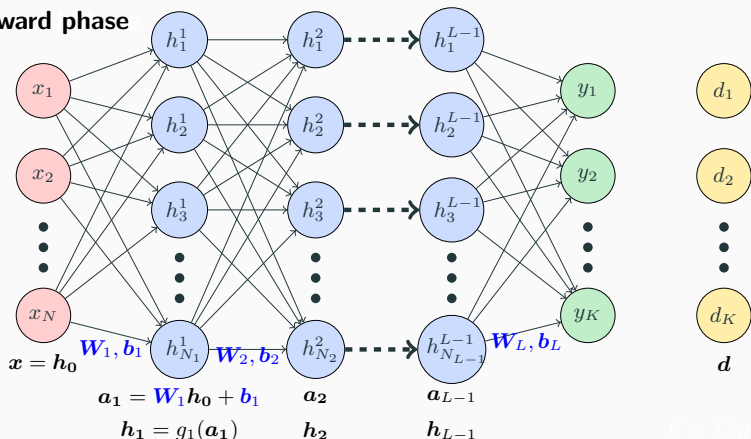
Hidden Layers

Output Layer

Label

Error backpropagation

Forward phase



Input Layer

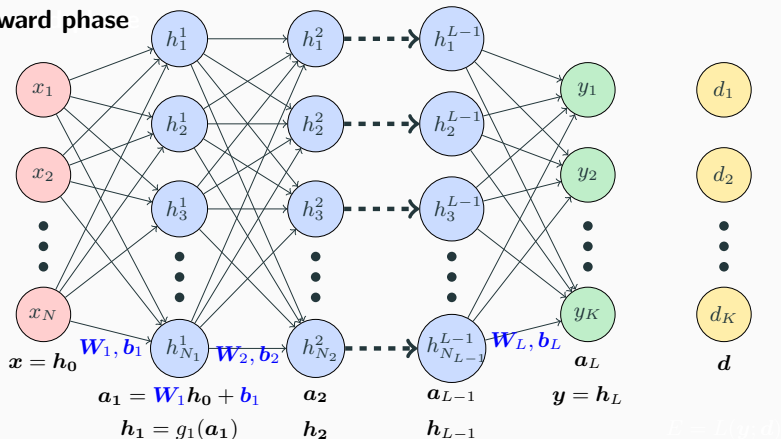
Hidden Layers

Output Layer

Label

Error backpropagation

Forward phase



Input Layer

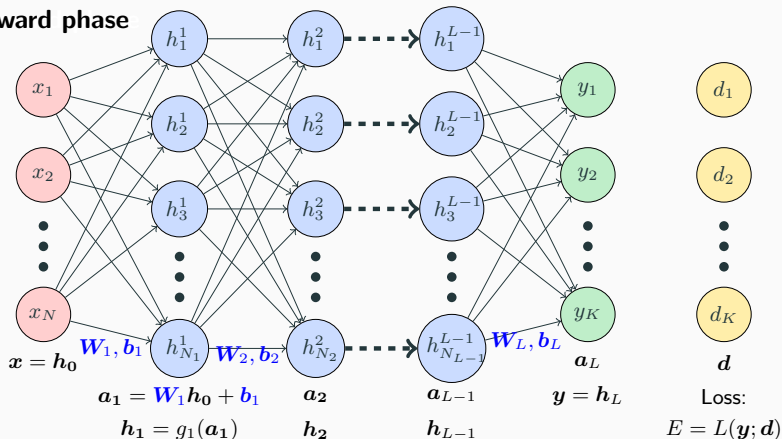
Hidden Layers

Output Layer

Label

Error backpropagation

Forward phase



Input Layer

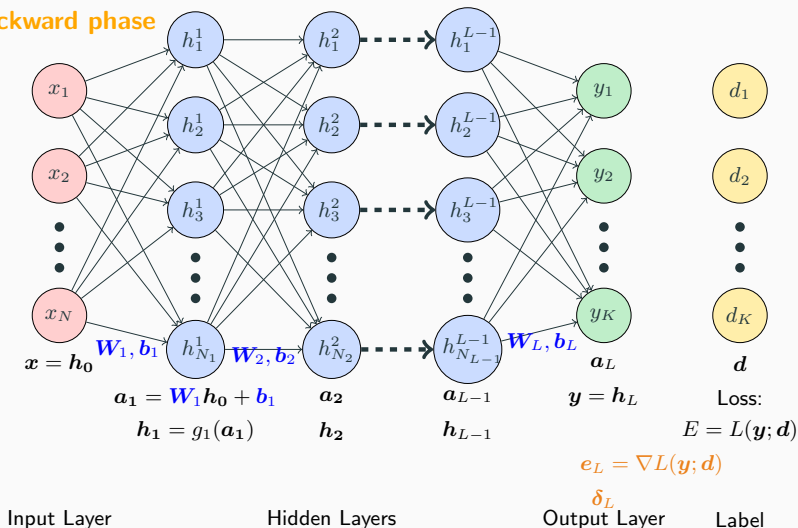
Hidden Layers

Output Layer

Label

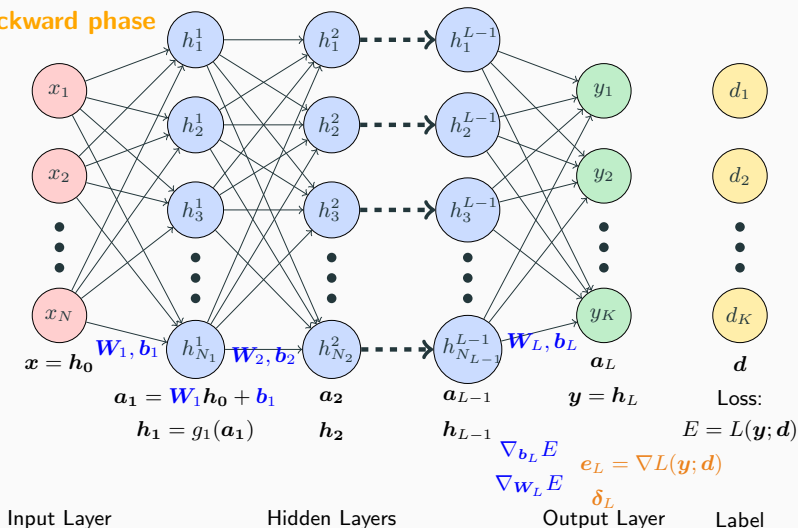
Error backpropagation

Backward phase



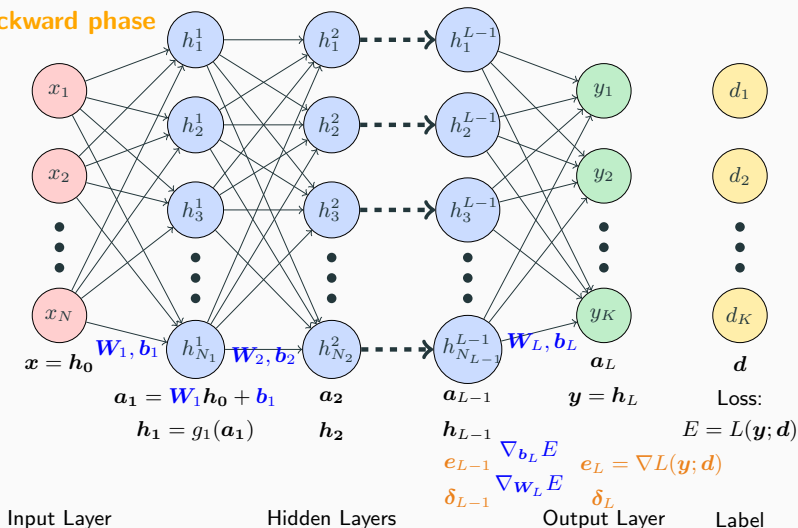
Error backpropagation

Backward phase



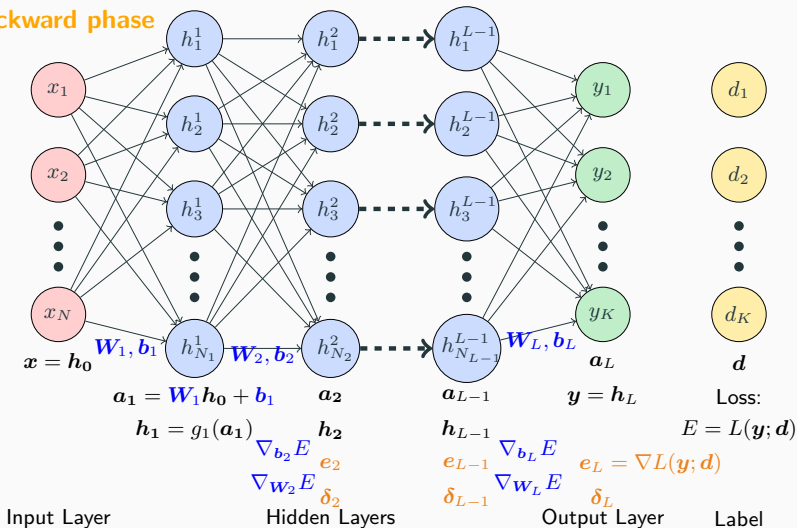
Error backpropagation

Backward phase



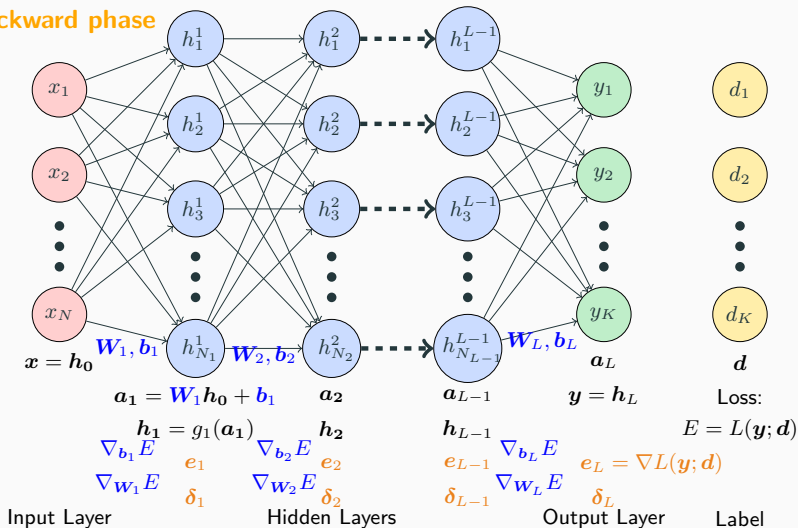
Error backpropagation

Backward phase



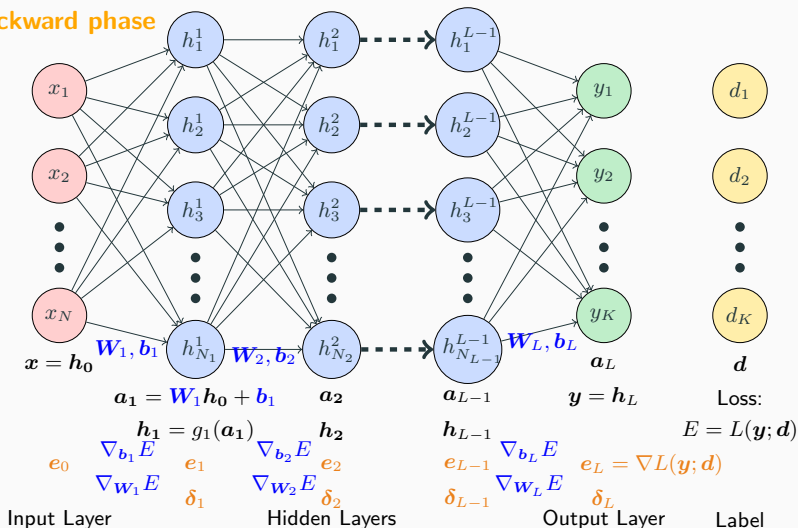
Error backpropagation

Backward phase



Error backpropagation

Backward phase



Error backpropagation in practice

Training loss:

$$E(\mathbf{W}) = \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{T}} L(\mathbf{y}^i; \mathbf{d}^i)$$

- The backpropagation procedure computes $\nabla_{\mathbf{W}} L(\mathbf{y}^i; \mathbf{d}^i) = \nabla_{\mathbf{W}} L(f(\mathbf{x}^i; \mathbf{W}); \mathbf{d}^i)$.
- This has to be done for each data point $\mathbf{x}^i \in \mathcal{T}$.
- By linearity, the final gradient $\nabla E(\mathbf{W})$ is the sum of all individual gradients $\nabla_{\mathbf{W}} L(\mathbf{y}^i; \mathbf{d}^i)$.
- These gradients are summed sequentially (no need to store each individual gradients).
- In general we do not compute the exact gradient...

Error backpropagation in practice

Batch loss:

$$E(\mathbf{W}) \approx \sum_{(\mathbf{x}^i, \mathbf{d}^i) \in \mathcal{S}} L(\mathbf{y}^i; \mathbf{d}^i), \quad \text{with } \mathcal{S} \subset \mathcal{T}$$

- The backpropagation has to be done for each visited data point $\mathbf{x}^i \in \mathcal{S}$ of the batch.
- The gradient for each point \mathbf{x}^i is added to the running gradient = current gradient estimation.
- Once the noisy estimated gradient is used as a gradient step, one needs to set the gradients to zero: See PyTorch `torch.zero_grad()` procedure.

Questions?

Sources, images courtesy and acknowledgment

Charles Deledalle

V. Lepetit

L. Masuch