# Chapter II – Preliminaries to deep learning

Charles Deledalle
July 10, 2019

## 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)*

## Brief history

- **First wave**
  - 1943. McCulloch and Pitts proposed the first neural model,
  - 1958. Rosenblatt introduced the Perceptron,
  - 1969. Minsky and Papert's book demonstrated the limitation of single layer perceptrons, and almost the whole field went into hibernation.
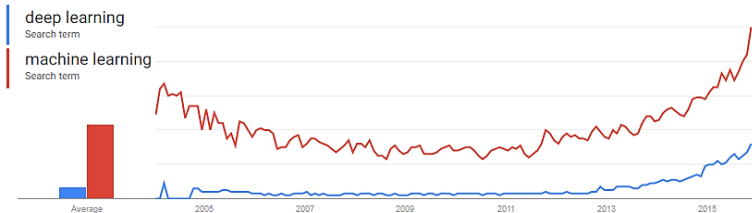
- **Second wave**
  - 1986. Backpropagation learning algorithm was rediscovered,
  - 1989. Yann LeCun (re)introduced Convolutional Neural Networks,
  - 1998. Breakthrough of CNNs in recognizing hand-written digits.

- **Third wave**
  - 2006. Deep (neural network) learning gains popularity,
  - 2012. Made significant breakthroughs in many applications,
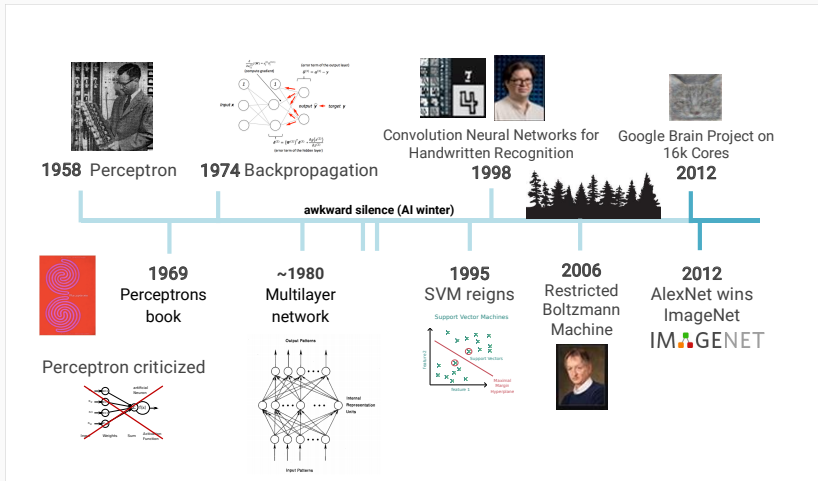  - 2015. AlphaGo first program to beat a professional Go player.

*(Source: Jun Wang)*
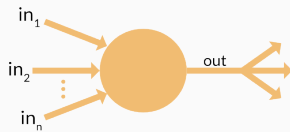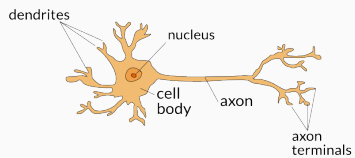
## Google NGRAM & Google Trends

# Timeline of (deep) learning



**1958** Perceptron   **1974** Backpropagation

Convolution Neural Networks for
Handwritten Recognition
**1998**

Google Brain Project on
16k Cores
**2012**

awkward silence (AI winter)

**1969**
Perceptrons
book

Perceptron criticized

**~1980**
Multilayer
network

**1995**
SVM reigns

**2006**
Restricted
Boltzmann
Machine

**2012**
AlexNet wins
ImageNet
IM✦GENET

*(Source: Lucas Masuch & Vincent Lepetit)*

# Perceptron

# Perceptron



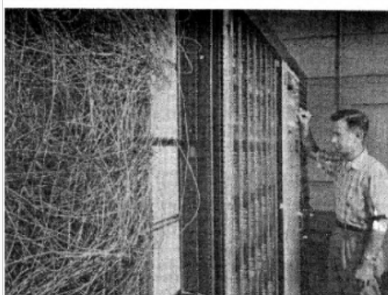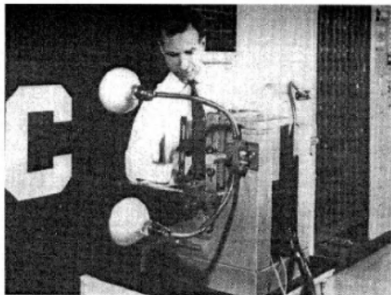**1958** Perceptron

**1969**
Perceptrons
book

Perceptron criticized

*(Source: Lucas Masuch & Vincent Lepetit)*
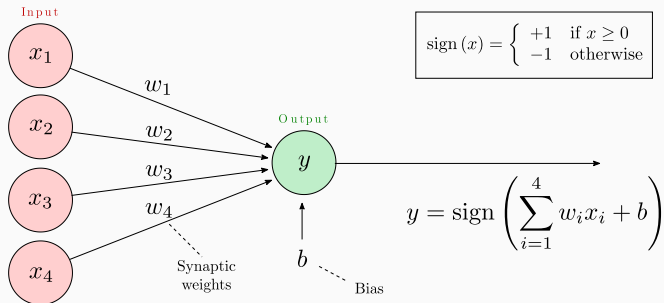
## Perceptron (Frank Rosenblatt, 1958)



First binary classifier based on supervised learning (discrimination).

Foundation of modern artificial neural networks.

At that time: technological, scientific and philosophical challenges.
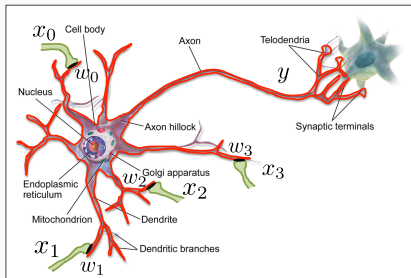
## Representation of the Perceptron



$$\text{sign}\,(x) = \left\{ \begin{array}{ll} +1 & \text{if } x \geq 0 \\ -1 & \text{otherwise} \end{array} \right.$$

$$y = \text{sign}\left( \sum_{i=1}^{4} w_i x_i + b \right)$$

**Parameters of the perceptron**

- $w_k$: synaptic weights
- $b$: bias

$\left. \right\}$ ⟵ real parameters to be estimated.

**Training = adjusting the weights and biases**

## The origin of the Perceptron

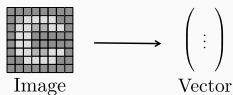**Takes inspiration from the visual system known for its ability to learn patterns.**



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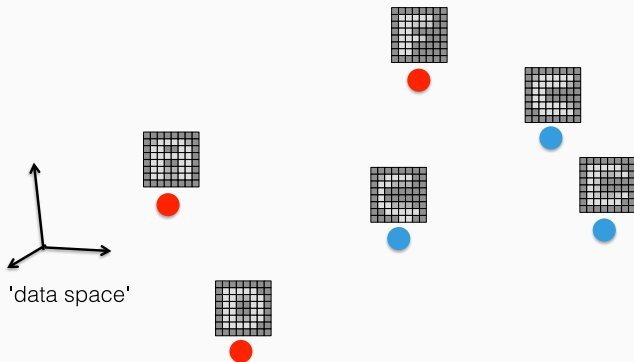
$$\sum_i w_i x_i + b > 0$$

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

❶ Data are represented as vectors:



Image → Vector

❷ Collect training data with positive and negative examples:



'data space'

❸ **Training:** find $\boldsymbol{w}$ and $b$ so that:

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

Dot product:

$$\langle \boldsymbol{w}, \boldsymbol{x} \rangle = \sum_{i=1}^{d} w_i x_i$$

$$= \boldsymbol{w}^T \boldsymbol{x}$$



'data space'

$\boldsymbol{x}$

10

❸ **Training:** find $w$ and $b$ so that:

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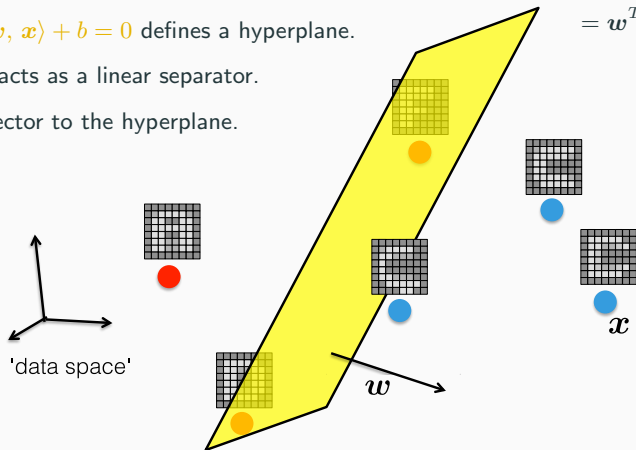- $\langle w, x \rangle + b$ is negative for negative samples $x$.

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.

Dot product:

$$\langle w, x \rangle = \sum_{i=1}^{d} w_i x_i$$

$$= w^T x$$



'data space'

$x$

$w$

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



'data space'

$w$

$x$

*(Source: Vincent Lepetit)*

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

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

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

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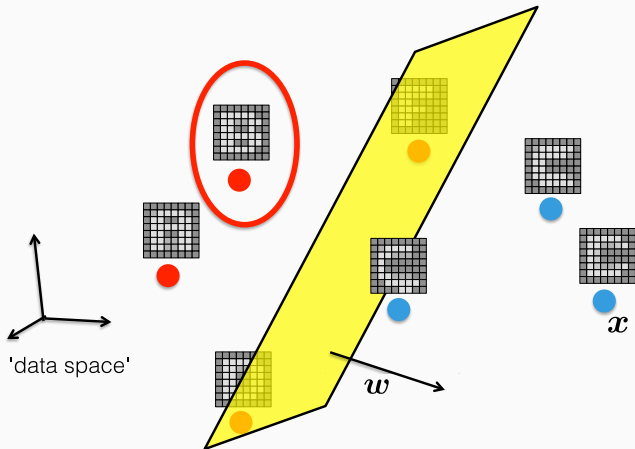- and negative if $\langle w, x \rangle + b$ is negative.



*(Source: Vincent Lepetit)*

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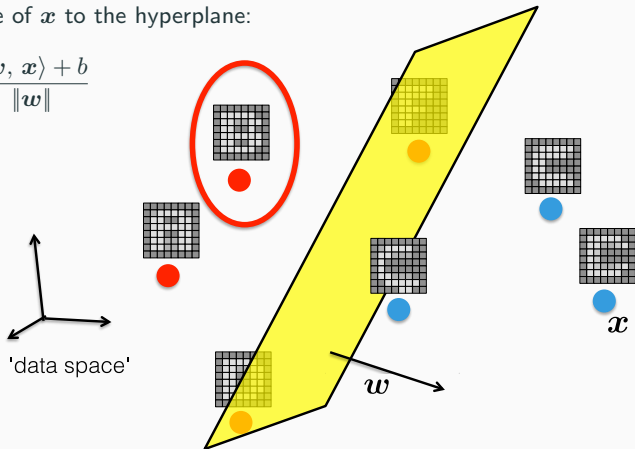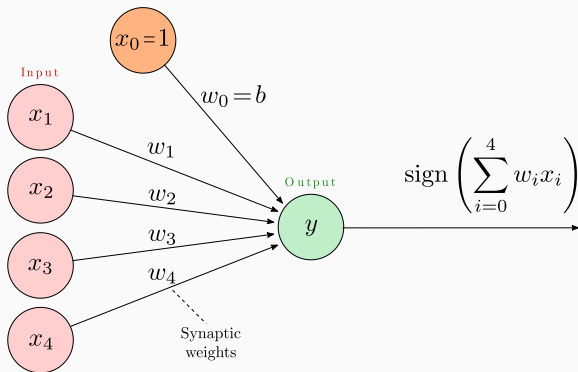
$$r = \frac{\langle w, x \rangle + b}{\|w\|}$$



'data space'

10

## 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 $w$ from a labeled training dataset $\mathcal{T}$

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

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

**How:** minimize classification errors

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

- penalize only misclassified samples ($y \neq d$),
- zero if all samples are correctly classified,
- proportional to the absolute distance to the hyperplane ($d \times \langle \boldsymbol{w}, \boldsymbol{x} \rangle < 0$)

## Perceptron algorithm

**Algorithm:**

- Initialize $\boldsymbol{w}$ randomly
- Repeat until convergence
  - For all $(\boldsymbol{x}, d) \in \mathcal{T}$
    - Compute: $y = \text{sign}\langle \boldsymbol{w}, \boldsymbol{x} \rangle$
    - If $y \neq d$:
      Update: $\boldsymbol{w} \leftarrow \boldsymbol{w} + d\boldsymbol{x}$



- Converges to some solution if the training data are linearly separable,
- Corresponds to stochastic gradient descent for $E(\boldsymbol{w})$ (see later),
- But may pick any of many solutions of varying quality.
  $\Rightarrow$ Poor generalization error.

13

## Variant: ADALINE (Adaptive Linear Neuron) algorithm
(Widrow & Hoff, 1960).

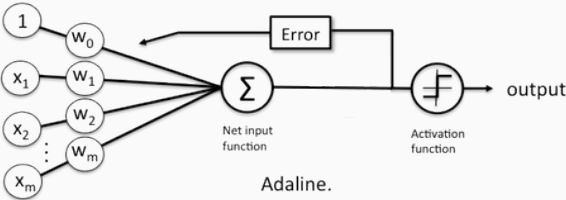**Loss:** minimize instead the least square error with the prediction without sign

$$\min_{\boldsymbol{w}} E(\boldsymbol{w}) = \sum_{(\boldsymbol{x}, d) \in \mathcal{T}} (\langle \boldsymbol{w},\, \boldsymbol{x} \rangle - d)^2$$

**Algorithm:**

- Initialize $\boldsymbol{w}$ randomly
- Repeat until convergence
    - For all $(\boldsymbol{x}, d) \in \mathcal{T}$
        - Compute: $y = \langle \boldsymbol{w},\, \boldsymbol{x} \rangle$
        - Update: $\boldsymbol{w} \leftarrow \boldsymbol{w} + \gamma(d - y)\boldsymbol{x},\ \gamma > 0$

Also corresponds to stochastic gradient descent for $E(\boldsymbol{w})$ (see later).

## Perceptron vs ADALINE algorithm



Perceptron rule.



Adaline.

Activation function = sign

## Perceptrons book (Minsky and Papert, 1969)
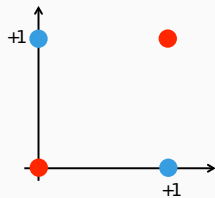
A perceptron can only classify data points that are linearly separable:



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



*(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).

- Essentially numerical but can handle classification and discrete inputs with appropriate coding.

- 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.

18

## 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



Input: $x_1$, $x_2$, $x_3$
Hidden: $h_1$, $h_2$, $h_3$, $h_4$
Output: $y_1$, $y_2$

**Quiz: how many layers does this network have?**

- Inter-connection of several artificial neurons.
- Each level in the graph is called a layer:
  - Input layer,
  - Hidden layer(s),
  - Output layer.
- Each neuron (also called node or unit) in the hidden layers acts as a classifier.
- 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.

20

# Artificial neural network / Multilayer perceptron / NeuralNet



$$h_1 = g_1 \left( w_{11}^1 x_1 + w_{12}^1 x_2 + w_{13}^1 x_3 + b_1^1 \right)$$
$$h_2 = g_1 \left( w_{21}^1 x_1 + w_{22}^1 x_2 + w_{23}^1 x_3 + b_2^1 \right)$$
$$h_3 = g_1 \left( w_{31}^1 x_1 + w_{32}^1 x_2 + w_{33}^1 x_3 + b_3^1 \right)$$
$$h_4 = g_1 \left( w_{41}^1 x_1 + w_{42}^1 x_2 + w_{43}^1 x_3 + b_4^1 \right)$$

$$y_1 = g_2 \left( w_{11}^2 h_1 + w_{12}^2 h_2 + w_{13}^2 h_3 + w_{14}^2 h_4 + b_1^2 \right)$$
$$y_2 = g_2 \left( w_{21}^2 h_1 + w_{22}^2 h_2 + w_{23}^2 h_3 + w_{24}^2 h_4 + b_2^2 \right)$$

$w_{ij}^k$ 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.

21

# Artificial neural network / Multilayer perceptron / NeuralNet



$(\boldsymbol{W}_1, \boldsymbol{b}_1)$    Hidden    $(\boldsymbol{W}_2, \boldsymbol{b}_2)$

$$h_1 = g_1 \left( w_{11}^1 x_1 + w_{12}^1 x_2 + w_{13}^1 x_3 + b_1^1 \right)$$
$$h_2 = g_1 \left( w_{21}^1 x_1 + w_{22}^1 x_2 + w_{23}^1 x_3 + b_2^1 \right)$$
$$h_3 = g_1 \left( w_{31}^1 x_1 + w_{32}^1 x_2 + w_{33}^1 x_3 + b_3^1 \right)$$
$$h_4 = g_1 \left( w_{41}^1 x_1 + w_{42}^1 x_2 + w_{43}^1 x_3 + b_4^1 \right)$$
$$\boldsymbol{h} = g_1 \left( \boldsymbol{W}_1 \boldsymbol{x} + \boldsymbol{b}_1 \right)$$
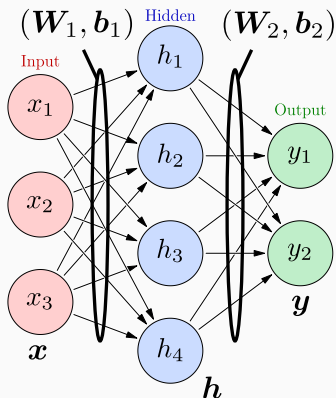
$$y_1 = g_2 \left( w_{11}^2 h_1 + w_{12}^2 h_2 + w_{13}^2 h_3 + w_{14}^2 h_4 + b_1^2 \right)$$
$$y_2 = g_2 \left( w_{21}^2 h_1 + w_{22}^2 h_2 + w_{23}^2 h_3 + w_{24}^2 h_4 + b_2^2 \right)$$
$$\boldsymbol{y} = g_2 \left( \boldsymbol{W}_2 \boldsymbol{h} + \boldsymbol{b}_2 \right)$$

$w_{ij}^k$ 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 $\boldsymbol{W_k}$ and biases $\boldsymbol{b}_k$ are learned from labeled training data.

21

## 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:

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

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

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

- It can be defined recursively as

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

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

## Activation functions

**Linear units:** $g(a) = a$

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

$$\frac{h_{L-1} = W_{L-1} h_{L-2} + b_{L-1}}{y = W_L W_{L-1} h_{L-2} + W_L b_{L-1} + b_L}$$

$$y = W_L \dots W_1 x + \sum_{k=1}^{L-1} W_L \dots W_{k+1} b_k + b_L$$

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

Sometimes used for linear dimensionality reduction (similarly to PCA).

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

$$g(a) = \left\{ \begin{array}{ll} -1 & \text{if} \quad a < 0 \\ +1 & \text{otherwise.} \end{array} \right.$$

or Heaviside (aka, step) activation functions

$$g(a) = \left\{ \begin{array}{ll} 0 & \text{if} \quad a < 0 \\ 1 & \text{otherwise.} \end{array} \right.$$

Discontinuities in the hidden layers
make the optimization really difficult.

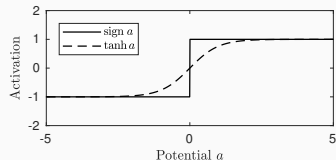We prefer functions that are continuous and differentiable.

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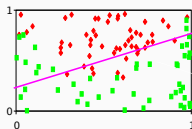
$$\tanh(a/2) = 2\text{logistic}(a) - 1$$

- Differentiable approximations of the sign and step functions, respectively.

- Act as threshold units for large values of $|a|$ and as linear for small values.

26

**Sigmoidal units**: logistic activation functions are used in binary classification (class $C_1$ vs $C_2$) as they can be interpreted as posterior probabilities:
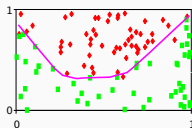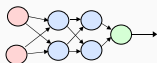
$$y = P(C_1|\boldsymbol{x}) \quad \text{and} \quad 1 - y = P(C_2|\boldsymbol{x})$$

The architecture of the network defines the shape of the separator
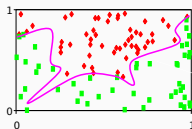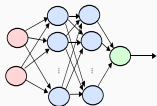
1 neuron

2+2+1 neurons

10+10+1 neurons



Separation
$$\{\boldsymbol{x} \setminus P(C_1|\boldsymbol{x}) = P(C_2|\boldsymbol{x})\}$$

Complexity/capacity of
the network
$$\Rightarrow$$
**Trade-off between
generalization and
overfitting**.

27

## 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 (see later).

*(Source: Lucas Masuch)*

28

## Neural networks solve non-linear separable problems

The x-or function



$$\boldsymbol{h} = g(\boldsymbol{W}_1 \boldsymbol{x} + \boldsymbol{b}_1)$$

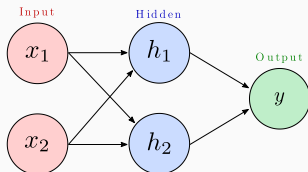$$y = \langle \boldsymbol{w}_2, \boldsymbol{h} \rangle + b_2$$

$$\boldsymbol{W}_1 = \begin{pmatrix} +1 & -1 \\ -1 & +1 \end{pmatrix}, \; \boldsymbol{b}_1 = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \; \boldsymbol{w}_2 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \; b_2 = 0$$

$$f(a) = \max(a, 0)$$

*(Source: Vincent Lepetit)*

29

## Universal Approximation Theorem
(Hornik et al, 1989; Cybenko, 1989)

Any continuous function can be approximated by a feedforward shallow network (*i.e.*, with 1-hidden layer only) with a sufficient number of neurons in the hidden layer.



$$\boldsymbol{h} = g(\boldsymbol{W}_1\boldsymbol{x} + \boldsymbol{b}_1)$$
$$\boldsymbol{y} = \boldsymbol{W}_2\boldsymbol{h} + \boldsymbol{b}_2$$

- 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.

# Tasks, architectures and loss functions

## Approximation – Least square regression

- **Goal:** Predict a real multivariate function.

- **How:** estimate the coefficients $\boldsymbol{W}$ of $\boldsymbol{y} = f(\boldsymbol{x}; \boldsymbol{W})$
  from labeled training examples where labels are real vectors:

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

$i$-th training example    desired output for sample $i$    number of training samples

- **Typical architecture:**



- Hidden layer:

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

- Linear output:

$$g(a) = a$$

31

## 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)

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

  and look for $\boldsymbol{W}^*$ such that $\nabla E(\boldsymbol{W}^*) = 0.$     *Recall: SSE ≡ SSD ≡ MSE*

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

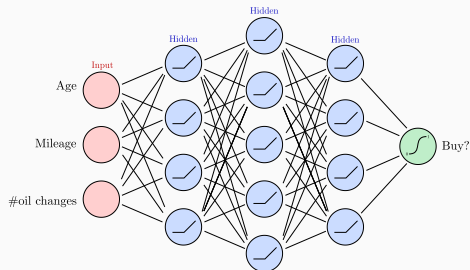$$\boldsymbol{y}^\star = f(\boldsymbol{x}; \boldsymbol{W}^\star) = \underbrace{\mathbb{E}[\boldsymbol{d}|\boldsymbol{x}] = \int \boldsymbol{d}p(\boldsymbol{d}|\boldsymbol{x})\, d\boldsymbol{d}}_{\text{posterior mean}}$$

## Binary classification – Logistic regression

- **Goal:** Classify object $\boldsymbol{x}$ into class $C_1$ or $C_2$.

- **How:** estimate the coefficients $\boldsymbol{W}$ of a real function $y = f(\boldsymbol{x}; \boldsymbol{W}) \in [0, 1]$ from training examples with labels $1$ (for class $C_1$) and $0$ (otherwise):

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

- **Typical architecture:**



- Hidden layer:

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

- Output layer:

$$\mathrm{logistic}(x) = \frac{1}{1 + e^{-a}}$$

## Binary classification – Logistic regression

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

$$E(\boldsymbol{W}) = -\sum_{i=1}^{N} d^i \log y^i + (1 - d^i) \log(1 - y^i) \quad \text{with} \quad y^i = f(\boldsymbol{x}^i; \boldsymbol{W})$$

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

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

$$y^\star = f(\boldsymbol{x}; \boldsymbol{W}^\star) = \underbrace{\mathbb{P}(C_1 | \boldsymbol{x})}_{\text{posterior probability}}$$

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

- **Goal:** Classify an object $x$ into one among $K$ classes $C_1, \ldots, C_K$.

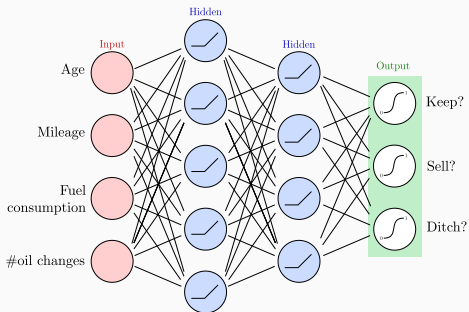- **How:** estimate the coefficients $W$ of a multivariate function

$$y = f(x; W) \in [0, 1]^K$$

from training examples $\mathcal{T} = \{(x^i, d^i)\}$ where $d^i$ is a 1-of-K (one-hot) code
   - Class 1: $\quad d^i = (1, 0, \ldots, 0)^T$ if $x^i \in C_1$
   - Class 2: $\quad d^i = (0, 1, \ldots, 0)^T$ if $x^i \in C_2$
   - $\ldots$
   - Class K: $\quad d^i = (0, 0, \ldots, 1)^T$ if $x^i \in C_K$

- **Remark:** do not use the class index $k$ directly as a scalar label.

35

## Multiclass classification – Multivariate logistic regression

- **Typical architecture:**



- Hidden layer:

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

- Output layer:

$$\mathsf{softmax}(\boldsymbol{a})_k = \frac{\exp(a_k)}{\sum_{l=1}^{K} \exp(a_l)}$$

Softmax guarantees the outputs $y_k$ to be positive and sum to $1$.

Generalization of the logistic sigmoid activation function.

Smooth version of winner-takes-all activation model (maxout).
(largest gets $+1$ others get $0$).

36

### Multiclass classification – Multivariate logistic regression

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

$$E(\boldsymbol{W}) = -\sum_{i=1}^{N}\sum_{k=1}^{K} d_k^i \log y_k^i \quad \text{with} \quad \boldsymbol{y}^i = f(\boldsymbol{x}^i; \boldsymbol{W})$$
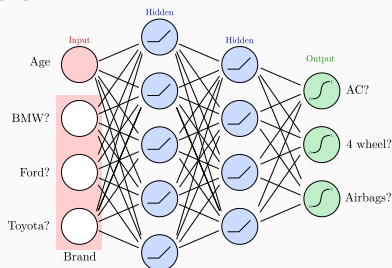
and look for $\boldsymbol{W}^*$ such that $\nabla E(\boldsymbol{W}^*) = 0$.

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

$$y_k^{\star} = f_k(\boldsymbol{x}; \boldsymbol{W}^{\star}) = \underbrace{\mathbb{P}(C_k|\boldsymbol{x})}_{\text{posterior probability}}$$
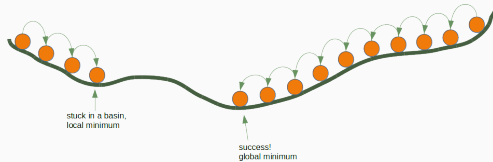
## Multi-label classification:
(aka, multi-output classification)

- **Goal:** Classify object $x$ into zero or several classes $C_1, \ldots, C_K$.
  The classes need to be non-mutually exclusive.

- **How:** Combination of binary-classification networks.

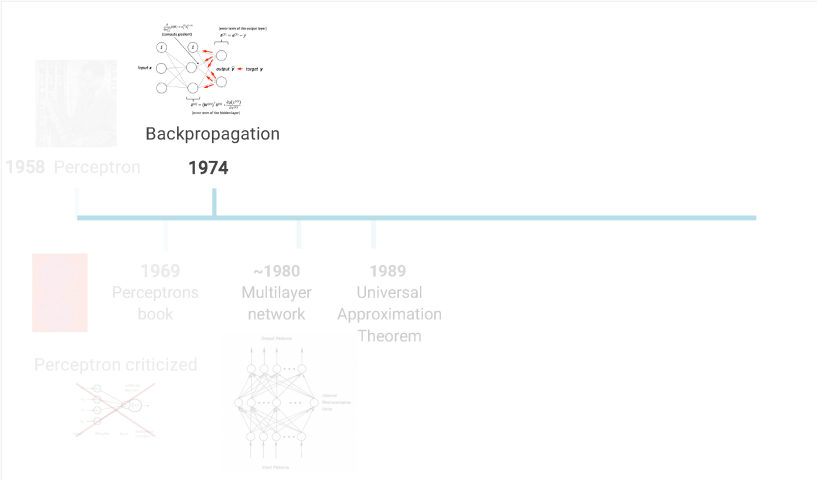- **Typical architecture:**



- **Remark:** For categorical inputs, use also 1-of-K codes.

# Backpropagation



stuck in a basin,
local minimum

success!
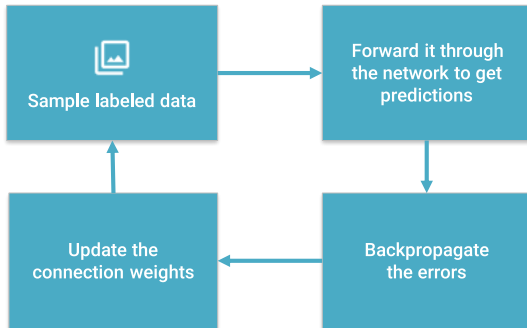global minimum

# Learning with backpropagation



*(Source: Lucas Masuch & Vincent Lepetit)*

## Training process

- To train a neural network over a large set of labeled data, you must continuously compute the difference between the network's predicted output and the actual output.

- This difference is measured by the loss $E$, and the process for training a net is known as backpropagation, or backprop.

- During backprop, weights and biases are tweaked slightly until the lowest possible loss is achieved.

$$\longrightarrow \textbf{Optimization: look for } \nabla E = 0$$

- The gradient is an important aspect of this process, as a measure of how much the loss changes with respect to a change in a weight or bias value.

*(Source: Caner Hazırbaş)*

## 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.

*(Source: Lucas Masuch)*

**Objective:** $\min_{\boldsymbol{W}} E(\boldsymbol{W}) \quad \Rightarrow \quad \nabla E(\boldsymbol{W}) = \left( \frac{\partial E(\boldsymbol{W})}{\partial \boldsymbol{W}_1} \quad \cdots \quad \frac{\partial E(\boldsymbol{W})}{\partial \boldsymbol{W}_L} \right)^T = 0$

**Loss functions:** recall that classical loss functions are

- Square error (for regression: $d_k \in \mathbb{R}$, $y_k \in \mathbb{R}$)

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

- Cross-entropy (for multi-class classification: $d_k \in \{0, 1\}$, $y_k \in [0, 1]$)

$$E(\boldsymbol{W}) = - \sum_{(\boldsymbol{x}, \boldsymbol{d}) \in \mathcal{T}} \sum_k d_k \log y_k$$

**Solution:** no closed-form solutions $\Rightarrow$ use gradient descent. **What is it?**

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

**Gradient descent**

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

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

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

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

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

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

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

## Gradient descent example

**Example (1d quadratic loss)**

- Consider: $F(x) = \frac{1}{2}(x - y)^2$

- We have: $F'(x) = x - y$

- And: $F''(x) = 1$

- Then: $L = \sup_x |F''(x)| = 1$

- Thus: The range for $\gamma$ is $0 < \gamma < 2/L = 2$

- And: $x^{t+1} = x^t - \gamma(x^t - y)$ converges towards $x^\star$, <u>whatever $x^0$</u>, such that $F'(x^\star) = 0 \Rightarrow x^\star = y$
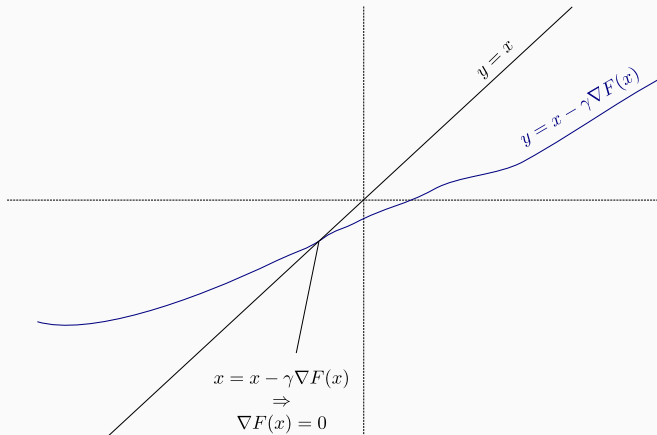
**Question: what is the best value for $\gamma$?**
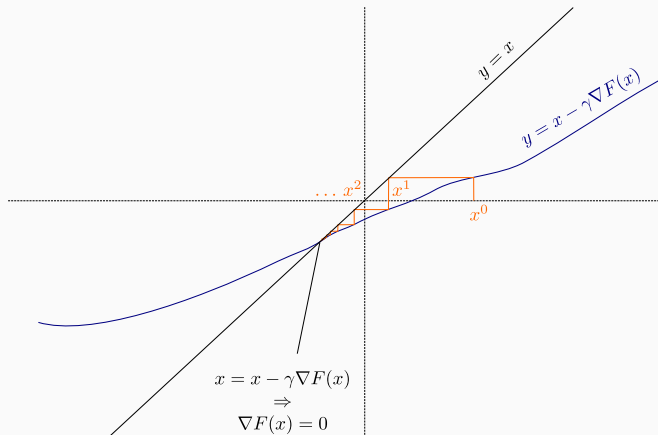
## Gradient descent example

**Example (1d quartic loss)**

- Consider: $F(x) = \frac{1}{4}(x - y)^4$
- We have: $F'(x) = (x - y)^3$
- And: $F''(x) = 3(x - y)^2$
- Then: $L = \sup_x |F''(x)| = \infty$
- Thus: There are no $\gamma$ satisfying $0 < \gamma < 2/L$
- And: $x^{t+1} = x^t - \gamma(x^t - y)^3$, for $\gamma > 0$, may diverge,
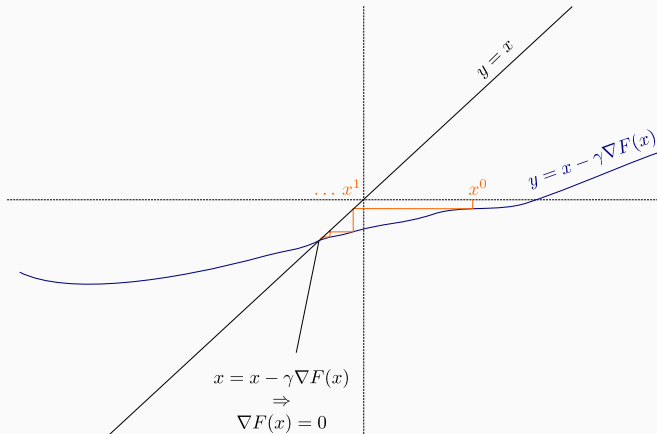
  oscillate forever or converge to the solution $y$.

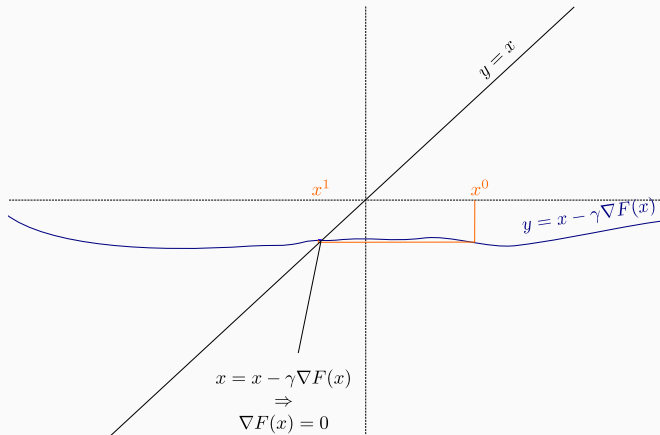**Convergence can be obtained for some specific choice
of $\gamma$ and initialization $x^0$.**

$y = x$

$y = x - \gamma \nabla F(x)$

$x = x - \gamma \nabla F(x)$
$\Rightarrow$
$\nabla F(x) = 0$

These two curves cross at $x^\star$ such that $\nabla F(x^\star) = 0$

Here $\gamma$ is small: slow convergence

$$y = x$$

$$y = x - \gamma \nabla F(x)$$

$$\dots x^1$$

$$x^0$$

$$x = x - \gamma \nabla F(x)$$
$$\Rightarrow$$
$$\nabla F(x) = 0$$

$\gamma$ a bit larger: faster convergence

$\gamma \approx 1/L$ even larger: around fastest convergence

$$y = x$$

$$x^1 \quad x^3 \; x^2 \quad x^0$$

$$y = x - \gamma \nabla F(x)$$

$$x = x - \gamma \nabla F(x)$$
$$\Rightarrow$$
$$\nabla F(x) = 0$$

$\gamma$ a bit too large: convergence slows down

$\gamma$ too large: convergence too slow again

$\gamma > 2/L$: divergence
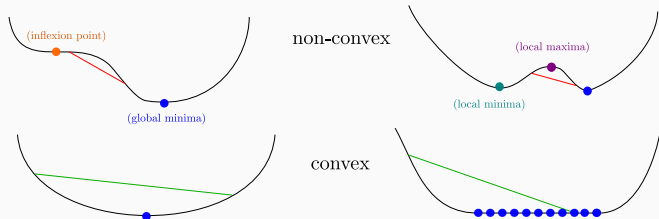
**Gradient descent for convex function**

- If moreover $F$ is convex

$$F(\lambda x_1 + (1 - \lambda)x_2) \leqslant \lambda F(x_1) + (1 - \lambda)F(x_2), \quad \forall x_1, x_2, \lambda \in [0, 1] \ ,$$
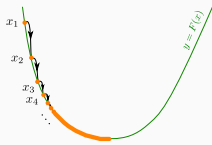
  then, the gradient descent converges towards a global minimum

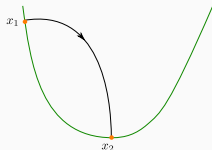$$x^\star \in \underset{x}{\operatorname{argmin}} \ F(x).$$

- For $0 < \gamma < 2/L$, the sequence $|F(x^k) - F(x^\star)|$ decays in $O(1/k)$.
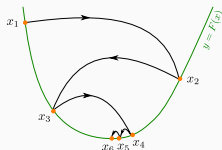- NB: All stationary points are global minima (not necessarily unique).
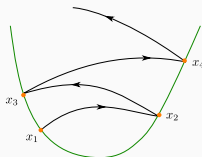
## One dimension



Small step size
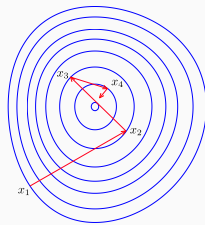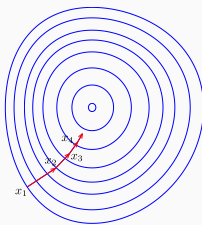Slow convergence

Good step size
Fast convergence

Large step size
Slow convergence

Too large step size
Divergence

## Two dimensions

## Let's start with a single artificial neuron

**Example ((Batch) Perceptron algorithm)**

- Model:
$$y = \text{sign}\langle \boldsymbol{w}, \boldsymbol{x} \rangle$$

- Loss:
$$E(\boldsymbol{w}) = - \sum_{\substack{(\boldsymbol{x},d)\in\mathcal{T} \\ \text{st } y\neq d}} d \times \langle \boldsymbol{w}, \boldsymbol{x} \rangle$$

- Gradient:
$$\nabla E(\boldsymbol{w}) = - \sum_{\substack{(\boldsymbol{x},d)\in\mathcal{T} \\ \text{st } y\neq d}} d \times \boldsymbol{x}$$

- Gradient descent:
$$\boldsymbol{w}^{t+1} \leftarrow \boldsymbol{w}^{t} + \gamma \sum_{\substack{(\boldsymbol{x},d)\in\mathcal{T} \\ \text{st } y^{t}\neq d}} d \times \boldsymbol{x}$$

**Convex or non-convex?**

### Let's start with a single artificial neuron

**Example ((Batch) ADALINE)**

- Model:
$$y = \langle \boldsymbol{w}, \boldsymbol{x} \rangle$$

- Loss:
$$E(\boldsymbol{w}) = \frac{1}{2} \sum_{(\boldsymbol{x},d)\in\mathcal{T}} \underbrace{(\langle \boldsymbol{w}, \boldsymbol{x} \rangle - d)^2}_{=d^2 + \boldsymbol{w}^T \boldsymbol{x}\boldsymbol{x}^T \boldsymbol{w} - 2d\boldsymbol{w}^T \boldsymbol{x}}$$

- Gradient:
$$\nabla E(\boldsymbol{w}) = \sum_{(\boldsymbol{x},d)\in\mathcal{T}} \boldsymbol{x}(\underbrace{\boldsymbol{x}^T \boldsymbol{w}}_{y} - d)$$

- Gradient descent:
$$\boldsymbol{w}^{t+1} \leftarrow \boldsymbol{w}^t + \gamma \sum_{(\boldsymbol{x},d)\in\mathcal{T}} (d - y^t)\boldsymbol{x}$$

**Convex or non-convex?**

51

### Let's start with a single artificial neuron

**Example ((Batch) ADALINE)**

- Model:
$$y = \langle \boldsymbol{w},\, \boldsymbol{x} \rangle$$

- Loss:
$$E(\boldsymbol{w}) = \frac{1}{2} \sum_{(\boldsymbol{x},d) \in \mathcal{T}} \underbrace{(\langle \boldsymbol{w},\, \boldsymbol{x} \rangle - d)^2}_{= d^2 + \boldsymbol{w}^T \boldsymbol{x}\boldsymbol{x}^T \boldsymbol{w} - 2d\boldsymbol{w}^T \boldsymbol{x}}$$

- Gradient:
$$\nabla E(\boldsymbol{w}) = \sum_{(\boldsymbol{x},d) \in \mathcal{T}} \boldsymbol{x}(\underbrace{\boldsymbol{x}^T \boldsymbol{w}}_{y} - d)$$

- Gradient descent:
$$\boldsymbol{w}^{t+1} \leftarrow \boldsymbol{w}^t + \gamma \sum_{(\boldsymbol{x},d) \in \mathcal{T}} (d - y^t)\boldsymbol{x}$$

**Convex or non-convex?**

If enough training samples: $\displaystyle \lim_{t \to \infty} \boldsymbol{w}^t = \Big( \underbrace{\sum_{(\boldsymbol{x},d) \in \mathcal{T}} \boldsymbol{x}\boldsymbol{x}^T}_{\text{Hessian}} \Big)^{-1} \Big( \sum_{(\boldsymbol{x},d) \in \mathcal{T}} d\boldsymbol{x} \Big)$
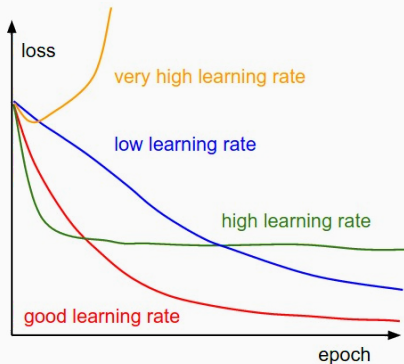
### Back to our optimization problem

In our case $W \mapsto E(W)$ is non-convex $\Rightarrow$ No guarantee of convergence.

Convergence will depend on $\left\{ \begin{array}{l} \bullet \text{ the initialization,} \\ \bullet \text{ the step size } \gamma. \end{array} \right.$

Because of this:

- Normalizing each data point $x$ in the range $[-1, +1]$ is important to control for the Hessian $\rightarrow$ the stability and the speed of the algorithm,

- The activation functions and the loss should be chosen to have a second derivative smaller than $1$ (when combined),

- The initialization should be random with well chosen variance
  (we will come back to this later),

- If all of these are satisfied, we can generally choose $\gamma \in [.001, 1]$.

**Influence of the step parameter in non-convex cases (learning rate)**

**Back to our optimization problem**

In our case $\boldsymbol{W} \mapsto E(\boldsymbol{W})$ is non-convex $\Rightarrow$ No guarantee of convergence.

Even if so, the limit solution depends on:
$$\left\{\begin{array}{l} \bullet \text{ the initialization,} \\ \bullet \text{ the step size } \gamma. \end{array}\right.$$

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

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

Gradient descent can be expressed coordinate by coordinate as:

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

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

$\Rightarrow$ The algorithm to compute $\dfrac{\partial E(\boldsymbol{W})}{\partial w_{i,j}}$ for ANNs is called backpropagation.

$$\text{Backpropagation: computation of } \frac{\partial E(\boldsymbol{W})}{\partial w_{i,j}}$$

**Feedforward least square regression context**

- **Model:** Feed-forward neural network.
  (for simplicity without bias)

- **Loss function:** $E(\boldsymbol{W}) = \dfrac{1}{2} \displaystyle\sum_{(\boldsymbol{x},\boldsymbol{d}) \in \mathcal{T}} \sum_{k} (y_k - d_k)^2$

We have:

$$E(\boldsymbol{W}) = \sum_{(\boldsymbol{x},\boldsymbol{d}) \in \mathcal{T}} \sum_{k} \underbrace{\frac{1}{2}(y_k - d_k)^2}_{e_k}$$

Apply linearity:

$$\frac{\partial E(\boldsymbol{W})}{\partial w_{i,j}} = \sum_{(\boldsymbol{x},\boldsymbol{d}) \in \mathcal{T}} \sum_{k} \frac{\partial e_k}{\partial w_{i,j}}$$

**1. Case where $w_{ij}$ is a synaptic weight for the output layer**



- $j$: neuron in the last hidden layer
- $h_j$: response of hidden neuron $j$
- $w_{i,j}$: synaptic weight between $j$ and $i$
- $y_i$: response of output neuron $i$

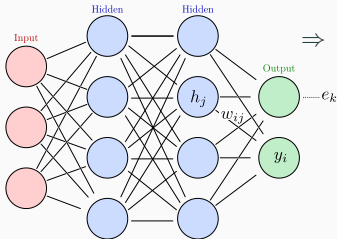$$y_i = g\left(a_i\right) \quad \text{with} \quad a_i = \sum_j w_{i,j} h_j$$

**Apply chain rule:** $\dfrac{\partial E(\boldsymbol{W})}{\partial w_{i,j}} = \displaystyle\sum_{(\boldsymbol{x},\boldsymbol{d})\in\mathcal{T}} \sum_k \dfrac{\partial e_k}{\partial w_{i,j}} = \sum_{(\boldsymbol{x},\boldsymbol{d})\in\mathcal{T}} \sum_k \dfrac{\partial e_k}{\partial y_i}\dfrac{\partial y_i}{\partial a_i}\dfrac{\partial a_i}{\partial w_{i,j}}$

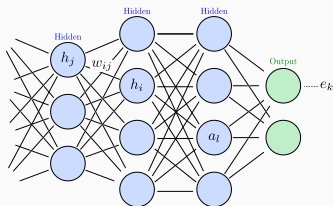**1. Case where $w_{ij}$ is a synaptic weight for the output layer**

$$e_k = \frac{1}{2}(y_k - d_k)^2 \quad \Rightarrow \quad \frac{\partial e_k}{\partial y_i} = \begin{cases} y_i - d_i & \text{if} \quad k = i \\ 0 & \text{otherwise} \end{cases}$$

$$y_i = g(a_i) \quad \Rightarrow \quad \frac{\partial y_i}{\partial a_i} = g'(a_i),$$

$$a_i = \sum_{j'} w_{i,j'} h_{j'} \quad \Rightarrow \quad \frac{\partial a_i}{\partial w_{i,j}} = h_j$$



$$\Rightarrow \quad \frac{\partial E(\boldsymbol{W})}{\partial w_{i,j}} = \sum_{(\boldsymbol{x},\boldsymbol{d}) \in \mathcal{T}} \sum_k \frac{\partial e_k}{\partial y_i} \frac{\partial y_i}{\partial a_i} \frac{\partial a_i}{\partial w_{i,j}}$$

$$= \sum_{(\boldsymbol{x},\boldsymbol{d}) \in \mathcal{T}} \underbrace{(y_i - d_i) g'(a_i)}_{\delta_i} h_j$$

$$= \sum_{(\boldsymbol{x},\boldsymbol{d}) \in \mathcal{T}} \delta_i h_j \quad \text{where} \quad \delta_i = \sum_k \frac{\partial e_k}{\partial a_i}$$

**2. Case where $w_{ij}$ is a synaptic weight for a hidden layer**



- $j$: neuron in the previous hidden layer
- $h_j$: response of hidden neuron $j$
- $w_{i,j}$: synaptic weight between $j$ and $i$
- $h_i$: response of hidden neuron $i$

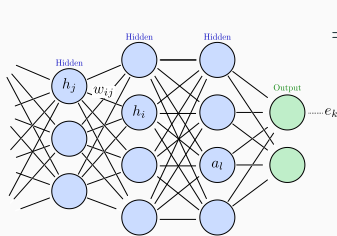$$h_i = g\left(a_i\right) \quad \text{with} \quad a_i = \sum_j w_{i,j} h_j$$

**Apply chain rule:**

$$\frac{\partial E(\boldsymbol{W})}{\partial w_{i,j}} = \sum_{(\boldsymbol{x},\boldsymbol{d})\in\mathcal{T}} \sum_k \frac{\partial e_k}{\partial w_{i,j}} = \sum_{(\boldsymbol{x},\boldsymbol{d})\in\mathcal{T}} \sum_k \frac{\partial e_k}{\partial h_i} \frac{\partial h_i}{\partial a_i} \frac{\partial a_i}{\partial w_{i,j}}$$

$$= \sum_{(\boldsymbol{x},\boldsymbol{d})\in\mathcal{T}} \sum_k \left(\sum_l \frac{\partial e_k}{\partial a_l} \frac{\partial a_l}{\partial h_i}\right) \frac{\partial h_i}{\partial a_i} \frac{\partial a_i}{\partial w_{i,j}}$$

$$= \sum_{(\boldsymbol{x},\boldsymbol{d})\in\mathcal{T}} \sum_l \underbrace{\left(\sum_k \frac{\partial e_k}{\partial a_l}\right)}_{\delta_l} \frac{\partial a_l}{\partial h_i} \frac{\partial h_i}{\partial a_i} \frac{\partial a_i}{\partial w_{i,j}}$$

**2. Case where $w_{ij}$ is a synaptic weight for a hidden layer**

$$a_l = \sum_{i'} w_{l,i'} h_{i'} \quad \Rightarrow \quad \frac{\partial a_l}{\partial h_i} = w_{l,i}$$

$$h_i = g(a_i) \quad \Rightarrow \quad \frac{\partial h_i}{\partial a_i} = g'(a_i),$$

$$a_i = \sum_{j'} w_{i,j'} h_{j'} \quad \Rightarrow \quad \frac{\partial a_j}{\partial w_{i,j}} = h_j$$



$$\Rightarrow \quad \frac{\partial E(\boldsymbol{W})}{\partial w_{i,j}} = \sum_{(\boldsymbol{x},\boldsymbol{d}) \in \mathcal{T}} \sum_l \delta_l \frac{\partial a_l}{\partial h_i} \frac{\partial h_i}{\partial a_i} \frac{\partial a_i}{\partial w_{i,j}}$$

$$= \sum_{(\boldsymbol{x},\boldsymbol{d}) \in \mathcal{T}} \underbrace{\left( \sum_l w_{l,i} \delta_l \right) g'(a_i)}_{\delta_i} h_j$$

$$= \sum_{(\boldsymbol{x},\boldsymbol{d}) \in \mathcal{T}} \delta_i h_j$$

### Backpropagation algorithm
(Werbos, 1974 & Rumelhart, Hinton and Williams, 1986)

$$\frac{\partial E(\boldsymbol{W})}{\partial w_{i,j}} = \sum_{(\boldsymbol{x},\boldsymbol{d})\in\mathcal{T}} \delta_i h_j \quad \text{where} \quad h_j = x_j \text{ if } j \text{ is an input node}$$

$$\text{where} \quad \delta_i = g'(a_i) \times \begin{cases} y_i - d_i & \text{if } i \text{ is the output node} \\ \sum_l w_{l,i}\delta_l & \text{otherwise} \end{cases}$$

For all input $\boldsymbol{x}$ and desired output $\boldsymbol{d}$

- Forward step:
  - $\rightarrow$ compute the response ($h_j$, $a_i$ and $y_i$) of all neurons,
  - $\rightarrow$ start from the first hidden layer and pursue towards the output one.

- Backward step:
  - $\rightarrow$ Retropropagate the error ($\delta_i$) from the output layer to the first layer.

Update $w_{i,j} \leftarrow w_{i,j} - \gamma \sum \delta_i h_j$, and repeat everything until convergence.

60

### Backpropagation algorithm with matrix-vector form

Easier to use **matrix-vector notations** for each layer:
($k$ denotes the layer)

$$\nabla_{\boldsymbol{W}_k} E(\boldsymbol{W}) = \boldsymbol{\delta}_k \boldsymbol{h}_{k-1}^T \quad \text{where} \quad \boldsymbol{h}_0 = \boldsymbol{x}$$

$$\text{where} \quad \boldsymbol{\delta}_k = \left[\frac{\partial g(\boldsymbol{a}_k)}{\partial \boldsymbol{a}_k}\right]^T \times \begin{cases} \boldsymbol{y} - \boldsymbol{d} & \text{if } k \text{ is an output layer} \\ \boldsymbol{W}_{k+1}^T \boldsymbol{\delta}_{k+1} & \text{otherwise} \end{cases}$$

- $\boldsymbol{x}$: matrix with all training input vectors in column,
- $\boldsymbol{d}$: matrix with corresponding desired target vectors in column,
- $\boldsymbol{y}$: matrix with all predictions in column,
- $\boldsymbol{a}_k = \boldsymbol{W}_k \boldsymbol{h}_{k-1}$: matrix with all weighted sums in column,
- $\boldsymbol{h}_k = g(\boldsymbol{a}_k)$: matrix with all hidden outputs in column,
- $\boldsymbol{W}_k$: matrix of weights at layer $k$,
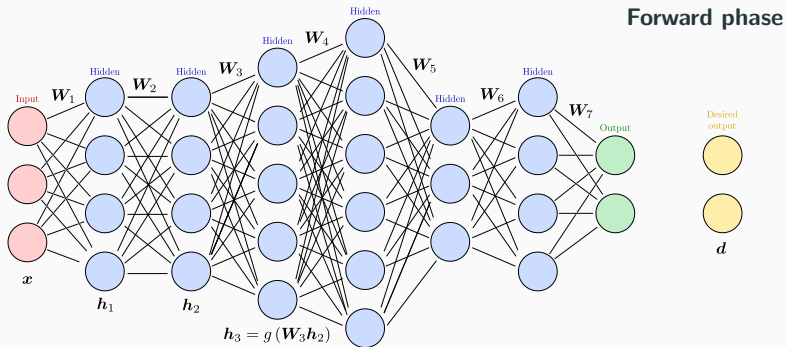
# Backpropagation algorithm



Forward phase

# Backpropagation algorithm



Forward phase

# Backpropagation algorithm



**Forward phase**

# Backpropagation algorithm



**Forward phase**

# Backpropagation algorithm



**Forward phase**

$$h_4 = g\left(\boldsymbol{W}_4 \boldsymbol{h}_3\right)$$

# Backpropagation algorithm



Forward phase

# Backpropagation algorithm



**Forward phase**

# Backpropagation algorithm



**Forward phase**

$$y = g(\underbrace{W_7 h_6}_{a_7})$$

# Backpropagation algorithm

# Backpropagation algorithm



**Backward phase**

# Backpropagation algorithm



**Backward phase**

$W_7 \leftarrow W_7 - \gamma \delta_7 h_6^T$

# Backpropagation algorithm



**Backward phase**

# Backpropagation algorithm



**Backward phase**

# Backpropagation algorithm



**Backward phase**

# Backpropagation algorithm



**Backward phase**

# Backpropagation algorithm



**Backward phase**

# Backpropagation algorithm



**Backward phase**

# Backpropagation algorithm



**Backward phase**

## Case where $g$ is the logistic sigmoid function

- Recall that the logistic sigmoid function is given by:

$$g(a) = \frac{1}{1 + e^{-a}} = \frac{1}{u(a)} \quad \text{where} \quad u(a) = 1 + e^{-a}$$

- We have:

$$u'(a) = -e^{-a}$$

- Then, we get

$$g'(a) = \frac{-u'(a)}{u(a)^2} = \frac{e^{-a}}{(1 + e^{-a})^2} = \frac{1 + e^{-a}}{(1 + e^{-a})^2} - \frac{1}{(1 + e^{-a})^2}$$

$$= \frac{1}{1 + e^{-a}} - \frac{1}{(1 + e^{-a})^2} = \frac{1}{1 + e^{-a}} \left(1 - \frac{1}{1 + e^{-a}}\right)$$

$$= g(a)\,(1 - g(a))$$

<div align="center">

**What if $g$ is non-differentiable?**

For instance: $g(a) = \text{ReLU}(a) = \max(a, 0)$

</div>

- ReLU is **continuous** and **differentiable almost everywhere**:

$$\frac{\partial \max(a, 0)}{\partial a} = \begin{cases} 1 & \text{if } a > 0 \\ 0 & \text{if } a < 0 \\ \text{undefined} & \text{otherwise} \end{cases}$$

- Gradient descent can handle this case with a simple modification:

$$\frac{\partial \max(a, 0)}{\partial a} = \begin{cases} 1 & \text{if } a > 0 \\ 0 & \text{if } a \leqslant 0 \end{cases}$$

- When $g$ is convex, this is called a sub-gradient, and gradient descent is called sub-gradient descent.

### Let's try to learn the x-or function

```python
import numpy as np

# Training set
x = np.array([[0,0], [0,1], [1,0], [1,1]]).T
d = np.array([  0,    +1,    +1,     0  ]).T

# Initialization for a 2 layer feedforward network
b1 = np.random.rand(2, 1)
W1 = np.random.rand(2, 2)
b2 = np.random.rand(1, 1)
W2 = np.random.rand(1, 2)

# Activation functions and their derivatives
def g1(a):  return a * (a > 0)          # ReLU
def g1p(a): return 1 * (a > 0)
def g2(a):  return a                    # Linear
def g2p(a): return 1
```

## Let's try to learn the x-or function

```python
import numpy as np

# Training set
x = np.array([[0,0], [0,1], [1,0], [1,1]]).T
d = np.array([  0,     +1,     +1,      0  ]).T

# Initialization for a 2 layer feedforward network
b1 = np.random.rand(2, 1)
W1 = np.random.rand(2, 2)
b2 = np.random.rand(1, 1)
W2 = np.random.rand(1, 2)

# Activation functions and their derivatives
def g1(a):  return a * (a > 0)            # ReLU
def g1p(a): return 1 * (a > 0)
def g2(a):  return 1 / (1 + np.exp(-a)) # Logistic
def g2p(a): return g2(a) * (1 - g2(a))
```

**Let's try to learn the x-or function**

```
gamma = .01 # step parameter (learning rate)
for t in range(0, 10000):
    # Forward phase
    a1      = W1.dot(x)
    h1      = g1(a1)
    a2      = W2.dot(h1)
    y       = g2(a2)
    # Error gradient evaluation
    e       = y - d
    # Backward phase
    delta2 = g2p(a2) * e
    delta1 = g1p(a1) * W2.T.dot(delta2)
    # gradient update
    W2      = W2 - gamma * delta2.dot(h1.T)
    W1      = W1 - gamma * delta1.dot(x.T)
    -
    -
```
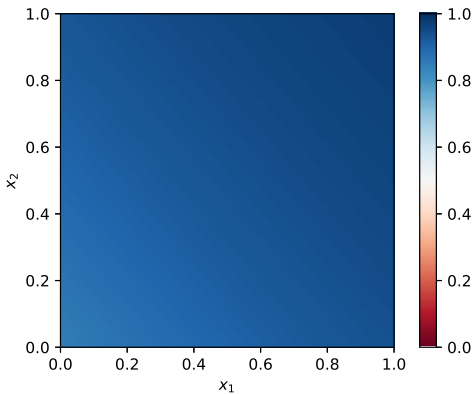
**Let's try to learn the x-or function**

```
gamma = .01 # step parameter
for t in range(0, 10000):
    # Forward phase
    a1      = W1.dot(x)  + b1
    h1      = g1(a1)
    a2      = W2.dot(h1) + b2
    y       = g2(a2)
    # Error gradient evaluation
    e       = y - d
    # Backward phase
    delta2 = g2p(a2) * e
    delta1 = g1p(a1) * W2.T.dot(delta2)
    # gradient update
    W2      = W2 - gamma * delta2.dot(h1.T)
    W1      = W1 - gamma * delta1.dot(x.T)
    b2      = b2 - gamma * delta2.sum(axis=1, keepdims=True)
    b1      = b1 - gamma * delta1.sum(axis=1, keepdims=True)
```

**Remark 1:** dealing with the bias is similar

**Let's try to learn the x-or function**

```
gamma = .01 # step parameter
for t in range(0, 10000):
    # Forward phase
    a1     = W1.dot(x)  + b1
    h1     = g1(a1)
    a2     = W2.dot(h1) + b2
    y      = g2(a2)
    # Error gradient evaluation
    e      = -d / y + (1 - d) / (1 - y)
    # Backward phase
    delta2 = g2p(a2) * e
    delta1 = g1p(a1) * W2.T.dot(delta2)
    # gradient update
    W2     = W2 - gamma * delta2.dot(h1.T)
    W1     = W1 - gamma * delta1.dot(x.T)
    b2     = b2 - gamma * delta2.sum(axis=1, keepdims=True)
    b1     = b1 - gamma * delta1.sum(axis=1, keepdims=True)
```

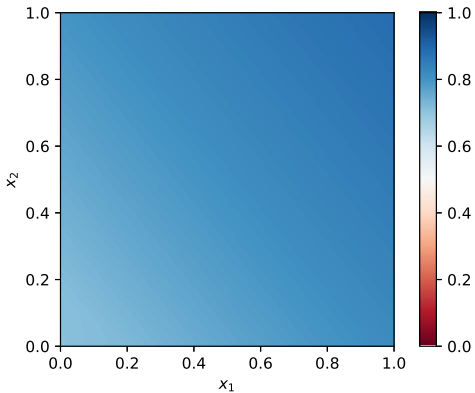**Remark 1:** dealing with the bias is similar

**Remark 2:** using cross-entropy is simple too
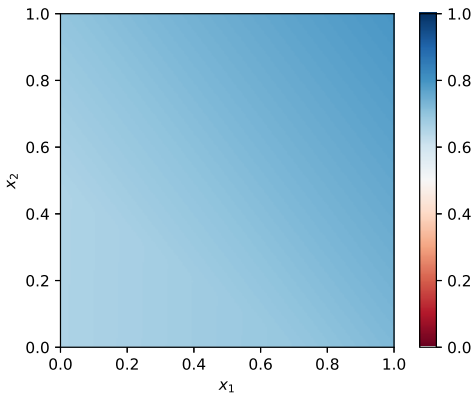
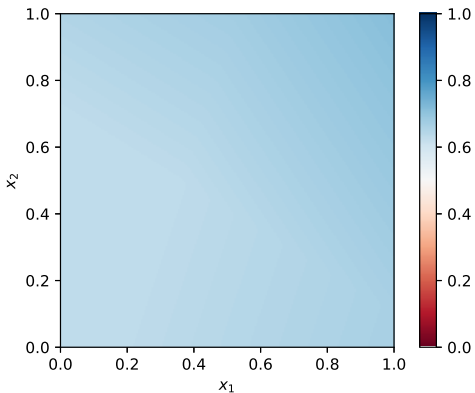**Let's see how it works**



Iteration 1

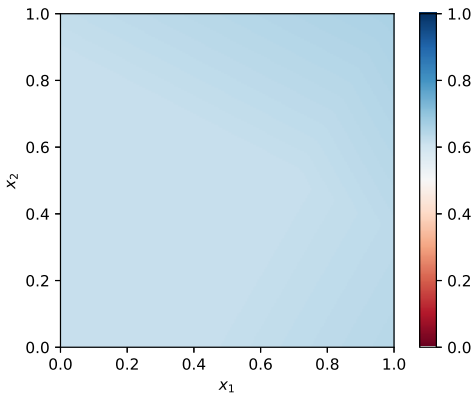**Let's see how it works**



Iteration 11

**Let's see how it works**



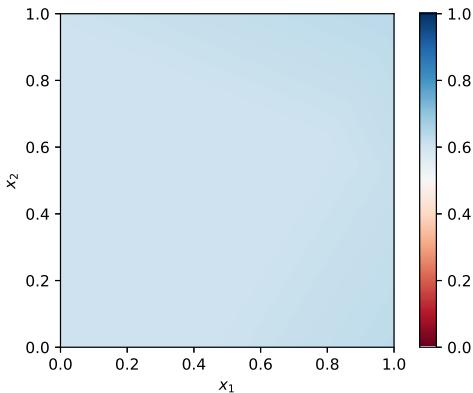Iteration 21

**Let's see how it works**



Iteration 31

**Let's see how it works**



Iteration 41

**Let's see how it works**



Iteration 51

**Let's see how it works**



Iteration 61

**Let's see how it works**



Iteration 71

**Let's see how it works**



Iteration 81

**Let's see how it works**



Iteration 91

**Let's see how it works**



Iteration 101

**Let's see how it works**



Iteration 201

**Let's see how it works**



Iteration 301
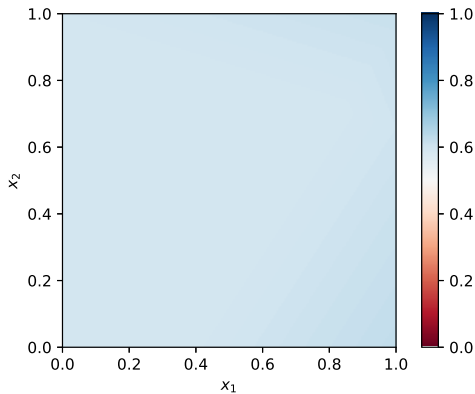
## Let's see how it works



Iteration 401
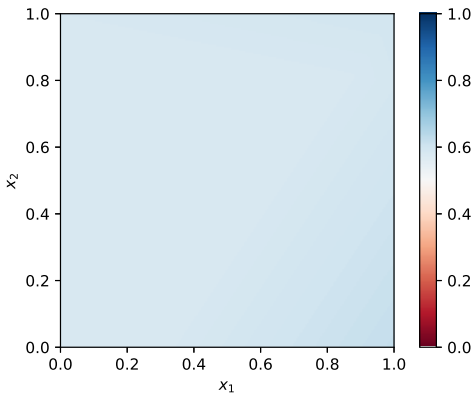
## Let's see how it works



Iteration 501
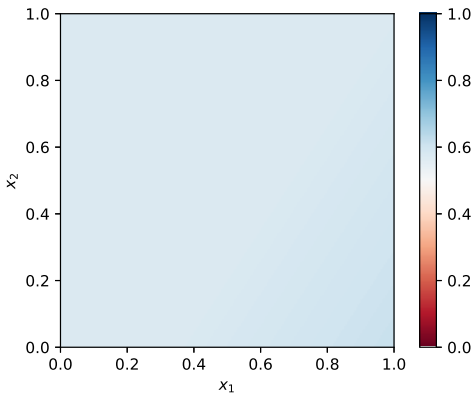
**Let's see how it works**



Iteration 601

## Let's see how it works



Iteration 701

**Let's see how it works**



Iteration 801

## Let's see how it works



Iteration 901

**Let's see how it works**



Iteration 1001

**Let's see how it works**



Iteration 2001

**Let's see how it works**



Iteration 3001

**Let's see how it works**



Iteration 4001

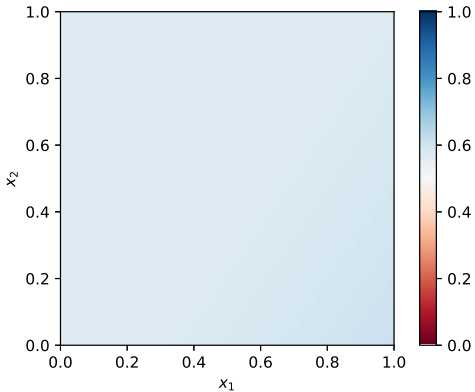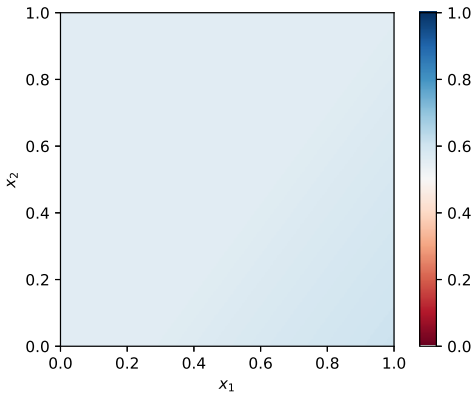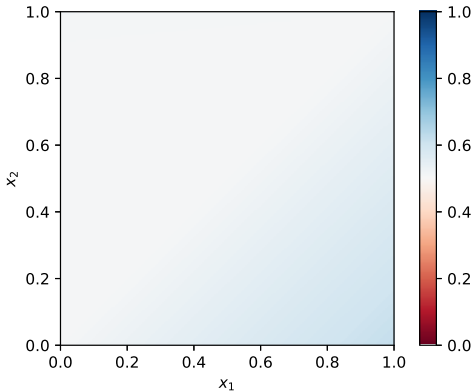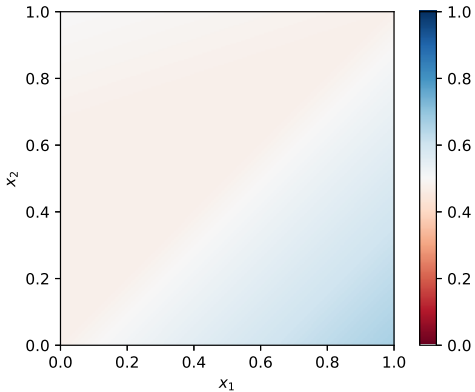## Let's see how it works
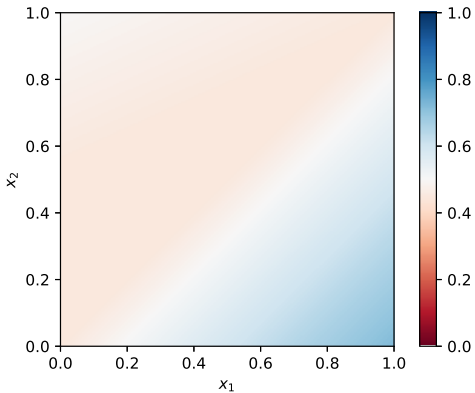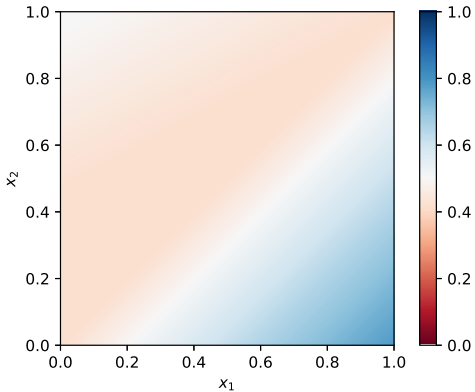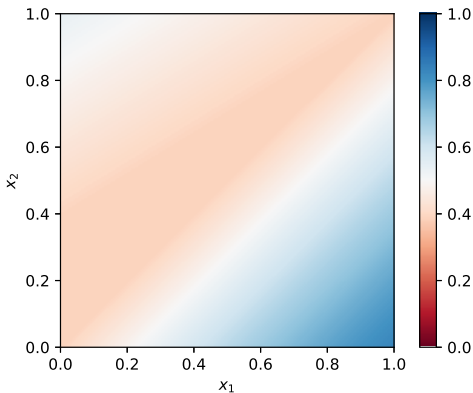


Iteration 5001

**Let's see how it works**



Iteration 6001

**Let's see how it works**



Iteration 7001

**Let's see how it works**



Iteration 8001

**Let's see how it works**



Iteration 9001

**Comparisons – different initializations & activations**

**Success** (Global minima)

**Comparisons – different initializations & activations**

**Failures** (Local minima or saddle points)

## The 1990s view of ANN and back-propagation

**Advantages**

- Universal approximators,
- May be very accurate.

**Drawbacks**

- Black box models (very difficult to interpret).
- The learning time does not scale well
  - it was very slow in networks with multiple hidden layers.
- It got stuck at local optima or saddle points
  - these can be nevertheless often surprisingly good.
- All global minima do not have the same quality (generalization error)
  - strong variations in prediction errors on different testing datasets.

# Support Vector Machine



$\varnothing$

## Support Vector Machine



*(Source: Lucas Masuch & Vincent Lepetit)*

### Support Vector Machine
(Vapnik, 1995)

- Very successful methods in the mid-90's.

- Developed for binary classification.

- Clever type of perceptron.

- Based on two major ideas
  1. large margin,
  2. kernel trick.

**Many NNs researchers switched to SVMs in the 1990s
because they (used to) work better.**

## Shattering points with oriented hyperplanes

- **Goal:** Build hyperplanes that separate points in two classes,

- **Question:** Which is the best separating line?



Remember, hyperplane: $\langle \boldsymbol{w},\, \boldsymbol{x} \rangle + b = 0$

## Classification margin

- Signed distance from an example to the separator: $r = \dfrac{\langle \boldsymbol{x},\, \boldsymbol{w} \rangle + b}{\|\boldsymbol{w}\|}$
- Examples closest to the hyperplane are support vectors.
- Margin $\rho$ is the distance between the separator and the support vector(s).



- Is this one good?

## Classification margin

- Signed distance from an example to the separator: $r = \dfrac{\langle \boldsymbol{x}, \, \boldsymbol{w} \rangle + b}{\|\boldsymbol{w}\|}$
- Examples closest to the hyperplane are support vectors.
- Margin $\rho$ is the distance between the separator and the support vector(s).



Margin $\rho$

Support
vector

Distance $r$

- Is this one good?

- Consider two testing samples, where are they assigned?

## Classification margin

- Signed distance from an example to the separator: $r = \dfrac{\langle \boldsymbol{x},\, \boldsymbol{w} \rangle + b}{\|\boldsymbol{w}\|}$

- Examples closest to the hyperplane are support vectors.

- Margin $\rho$ is the distance between the separator and the support vector(s).



Margin $\rho$

Support vector

Distance $r$

- Is this one good?

- Consider two testing samples, where are they assigned?

## Classification margin

- Signed distance from an example to the separator: $r = \dfrac{\langle \boldsymbol{x}, \boldsymbol{w} \rangle + b}{\|\boldsymbol{w}\|}$
- Examples closest to the hyperplane are support vectors.
- Margin $\rho$ is the distance between the separator and the support vector(s).



Really?

- Is this one good?
- Consider two testing samples, where are they assigned?
- This separator may have large generalization error.

## Largest margin = Support Vector Machine

- Maximizing the margin reduces generalization error (see, PAC learning).
- Then, there are necessary support vectors on both sides of the hyperplane.
- Only support vectors are important $\Rightarrow$ other examples can be discarded.



Poor separation                                  Optimal separation

## Training

- As for neural networks, the parameters $\boldsymbol{w}$ and $b$ are obtained by optimizing a loss function (the margin).

- What is the formula for the margin?

---

- Training dataset: feature vectors $\boldsymbol{x}^i$, targeted class $d^i \in \{-1, +1\}$
- Assume the training samples to be **linearly separable**,
  *i.e.*, there exist $\boldsymbol{w}$ and $b$ such that

$$\begin{cases} \langle \boldsymbol{x}^i, \boldsymbol{w} \rangle + b \geqslant +1 & \text{if} \quad d^i = +1 \\ \langle \boldsymbol{x}^i, \boldsymbol{w} \rangle + b \leqslant -1 & \text{if} \quad d^i = -1 \end{cases}$$

---

$$d^i(\langle \boldsymbol{x}^i, \boldsymbol{w} \rangle + b) \geqslant 1$$

**Then:** $\quad |\langle \boldsymbol{x}^i, \boldsymbol{w} \rangle + b| \geqslant 1$

## Training

- For support vectors, the inequality can be forced to be an equality

$$|\langle \boldsymbol{x}^i, \ \boldsymbol{w} \rangle + b| = 1$$

- Then, the distance between any support vector and the hyperplane is

$$|r| = \frac{|\langle \boldsymbol{x}^i, \ \boldsymbol{w} \rangle + b|}{\|\boldsymbol{w}\|} = \frac{1}{\|\boldsymbol{w}\|}$$

- So, the margin is (by definition)

$$\rho = |r| = \frac{1}{\|\boldsymbol{w}\|}$$

- Maximizing the margin is then equivalent to minimizing $\|\boldsymbol{w}\|^2$, provided the hyperplane $(\boldsymbol{w}, b)$ separates the data.

## Training

- Therefore, the problem can be recast as

$$\min_{\boldsymbol{w},b} \frac{1}{2} \|\boldsymbol{w}\|^2 \quad \text{subject to} \quad d^i(\langle \boldsymbol{x}^i, \boldsymbol{w} \rangle + b) \geqslant 1, \quad \text{for all } i.$$

⇒ Quadratic (convex) optimization problem subject to linear constraints,

⇒ No local minima! Only a single global one.

**Equivalent formulation:**

$$\min_{\boldsymbol{w},b} E(\boldsymbol{w}, b)$$

$$\text{where} \quad E(\boldsymbol{w}, b) = \begin{cases} \frac{1}{2}\|\boldsymbol{w}\|^2 & \text{if} \quad d^i(\langle \boldsymbol{x}^i, \boldsymbol{w} \rangle + b) \geqslant 1, \quad \text{for all } i \\ +\infty & \text{otherwise} \end{cases}$$

## **Training**

$$\min_{\boldsymbol{w},b} E(\boldsymbol{w}, b)$$

$$\text{where} \quad E(\boldsymbol{w}, b) = \begin{cases} \frac{1}{2}\|\boldsymbol{w}\|^2 & \text{if} \quad d^i(\langle \boldsymbol{x}^i, \, \boldsymbol{w}\rangle + b) \geqslant 1, \quad \text{for all } i \\ +\infty & \text{otherwise} \end{cases}$$

- We can use the techniques of Lagrange multipliers:
  1. Introduce a Lagrange multiplier $\alpha_i \geqslant 0$ for each constraint.
  2. Let $\boldsymbol{\alpha} = (\alpha_1, \alpha_2, \dots)^T$ be the vector of Lagrange multipliers.
  3. Define the (primal) function as

$$\mathcal{L}_P(\boldsymbol{w}, b, \boldsymbol{\alpha}) = \frac{1}{2}\|\boldsymbol{w}\|^2 - \sum_i \alpha^i \underbrace{(d^i(\langle \boldsymbol{x}^i, \, \boldsymbol{w}\rangle + b) - 1)}_{\text{inequality constraint}}$$

  4. Solve the saddle point problem

$$E(\boldsymbol{w}, b) = \max_{\boldsymbol{\alpha}} \mathcal{L}_P(\boldsymbol{w}, b, \boldsymbol{\alpha}) \quad \Rightarrow \quad \min_{\boldsymbol{w},b} E(\boldsymbol{w}, b) = \min_{\boldsymbol{w},b} \max_{\boldsymbol{\alpha}} \mathcal{L}_P(\boldsymbol{w}, b, \boldsymbol{\alpha})$$

## Training

- When the original objective-function is convex (and only then), we can interchange the minimization and maximization

$$\min_{\boldsymbol{w},b} \max_{\boldsymbol{\alpha}} \mathcal{L}_P(\boldsymbol{w},b,\boldsymbol{\alpha}) = \max_{\boldsymbol{\alpha}} \min_{\boldsymbol{w},b} \mathcal{L}_P(\boldsymbol{w},b,\boldsymbol{\alpha})$$

- The following is called dual problem

$$\mathcal{L}_D(\boldsymbol{\alpha}) = \min_{\boldsymbol{w},b} \left\{ \mathcal{L}_P(\boldsymbol{w},b,\boldsymbol{\alpha}) = \frac{1}{2} \underbrace{\|\boldsymbol{w}\|^2}_{\langle \boldsymbol{w},\,\boldsymbol{w}\rangle} - \sum_i \alpha^i (d^i (\langle \boldsymbol{x}^i,\,\boldsymbol{w}\rangle + b) - 1) \right\}$$

- For a fixed $\boldsymbol{\alpha}$, the minimum is achieved by simultaneously canceling the gradients with respect to $\boldsymbol{w}$ and $b$

$$\nabla_{\boldsymbol{w}} \mathcal{L}_P(\boldsymbol{w},b,\boldsymbol{\alpha}) = 0 \Rightarrow \boldsymbol{w} - \sum_i \alpha^i d^i \boldsymbol{x}^i = 0 \Rightarrow \boldsymbol{w} = \sum_i \alpha^i d^i \boldsymbol{x}^i$$

$$\text{and} \quad \nabla_b \mathcal{L}_P(\boldsymbol{w},b,\boldsymbol{\alpha}) = 0 \Rightarrow \sum_i \alpha^i d^i = 0$$

## Training

The dual is obtained by plugging these equations in $\mathcal{L}_P(\boldsymbol{w}, b, \boldsymbol{\alpha})$:

$$\mathcal{L}_D(\boldsymbol{\alpha}) = \frac{1}{2} \underbrace{\|\boldsymbol{w}\|^2}_{\langle \boldsymbol{w}, \, \boldsymbol{w} \rangle} - \sum_i \alpha^i (d^i (\langle \boldsymbol{x}^i, \, \boldsymbol{w} \rangle + b) - 1)$$

$$= \frac{1}{2} \langle \sum_i \alpha^i d^i \boldsymbol{x}^i, \, \sum_j \alpha^j d^j \boldsymbol{x}^j \rangle - \sum_i \alpha^i d^i (\langle \boldsymbol{x}^i, \, \sum_j \alpha^j d^j \boldsymbol{x}^j \rangle + b) + \sum_i \alpha^i$$

$$= \frac{1}{2} \sum_i \sum_j \alpha^i \alpha^j d^i d^j \langle \boldsymbol{x}^i, \, \boldsymbol{x}^j \rangle - \sum_i \sum_j \alpha^i \alpha^j d^i d^j \langle \boldsymbol{x}^i, \, \boldsymbol{x}^j \rangle + b \underbrace{\sum_i \alpha^i d^i}_{=0} + \sum_i \alpha^i$$

$$= \sum_i \alpha^i - \frac{1}{2} \sum_i \sum_j \alpha^i \alpha^j d^i d^j \langle \boldsymbol{x}^i, \, \boldsymbol{x}^j \rangle$$

## SVM training algorithm

① Maximize the dual (*e.g.*, using coordinate projected gradient descent):

$$\max_{\boldsymbol{\alpha}} \left\{ \mathcal{L}_D(\boldsymbol{\alpha}) = \sum_i \alpha^i - \frac{1}{2} \sum_i \sum_j \alpha^i \alpha^j d^i d^j \langle \boldsymbol{x}^i, \boldsymbol{x}^j \rangle \right\}$$

subject to $\quad \alpha^i \geqslant 0 \quad$ and $\quad \sum_i \alpha^i d^i = 0, \quad$ for all $i$.

② Each non-zero $\alpha^i$ indicates that corresponding $\boldsymbol{x}^i$ is a support vector.

③ Deduce the parameters from these support vectors (s.v)

$$\boldsymbol{w} = \sum_{i \text{ (s.v.)}} \alpha^i d^i \boldsymbol{x}^i \quad \text{and} \quad b = d^i - \sum_{j \text{ (s.v.)}} \alpha^j d^j \langle \boldsymbol{x}^i, \boldsymbol{x}^j \rangle \quad \text{for any s.v. } i$$

- The dual is in practice more efficient to solve than the original problem, and it allows identifying the support vectors.

- It only depends on the dot products $\langle \boldsymbol{x}^i, \boldsymbol{x}^j \rangle$ between all training points.

## SVM classification algorithm

- Given the parameters of the hyperplane $\boldsymbol{w}$ and $b$,
  the SVM classifies a new sample $\boldsymbol{x}$ as

$$y = \langle \boldsymbol{w}, \boldsymbol{x} \rangle + b \quad \lessgtr \quad 0$$

(same as for the perceptron).

- But (very important), this can be reformulated as

$$y = \sum_{i \text{ (s.v.)}} \alpha^i d^i \langle \boldsymbol{x}, \boldsymbol{x}^i \rangle + b$$

which only depends on the dot products $\langle \boldsymbol{x}, \boldsymbol{x}^i \rangle$ between $\boldsymbol{x}$ and the support vectors – we will return to this later.

## Non-separable case

What if the training set is not linearly separable?



**The optimization problem does not admit solutions.**

## Non-separable case

What if the training set is not linearly separable?



**The optimization problem does not admit solutions.**
**Solution: allows the system to make errors $\varepsilon_i$.**

## Non-separable case – Soft-margin SVM

- Just relax the constraints by permitting errors to some extent

$$\min_{\boldsymbol{w}, b, \boldsymbol{\varepsilon}} \frac{1}{2} \|\boldsymbol{w}\|^2 + C \sum_i \varepsilon_i$$

subject to $\quad d^i(\langle \boldsymbol{x}^i, \boldsymbol{w} \rangle + b) \geqslant 1 - \varepsilon_i \quad$ and $\quad \varepsilon_i \geqslant 0, \quad$ for all $i$.

- Quantities $\varepsilon_i$ are called slack variables.

- The parameter $C > 0$ is chosen by the user and controls overfitting.

- The dual problem becomes

$$\max_{\boldsymbol{\alpha}} \left\{ \mathcal{L}_D(\boldsymbol{\alpha}) = \sum_i \alpha^i - \frac{1}{2} \sum_i \sum_j \alpha^i \alpha^j d^i d^j \langle \boldsymbol{x}^i, \boldsymbol{x}^j \rangle \right\}$$

subject to $\quad 0 \leqslant \alpha_i \leqslant C \quad$ and $\quad \sum_i \alpha^i d^i = 0, \quad$ for all $i$.

and does not depend on the slack variables $\varepsilon_i$.

## Linear SVM – Overview

- SVM finds the separating hyperplane maximizing the margin.

- Soft-margin SVM: trade-off between large margin and errors.

- These optimal hyperplanes are only defined in terms of support vectors.

- Lagrangian formulation allows identifying the support vectors with non-zero Lagrange multipliers $\alpha^i$.

- The training and classification algorithms both depend only on dot products between data points.

**So far, the SVM is a linear separator (the best in some sense).
But, it cannot solve non-linear problems (such as xor),
as done by multi-layer neural networks.**

## Non-linear SVM – Overview



What happens if the separator is non-linear?

## Non-linear SVM – Map to higher dimensions



Input space          $x \to \varphi(x)$          Feature space

Here: $(x_1, x_2) \mapsto (x_1^2, x_2^2, \sqrt{2}x_1 x_2)$

The input space (original feature space) can always be mapped to some higher-dimensional feature space where the training data set is linearly separable, via some (non-linear) transformation $x \to \varphi(x)$.

## Non-linear SVM – Map to higher dimensions

Replace all occurrences of $\boldsymbol{x}$ by $\varphi(\boldsymbol{x})$, for the training step

$$\mathcal{L}_D(\boldsymbol{\alpha}) = \sum_i \alpha^i - \frac{1}{2} \sum_i \sum_j \alpha^i \alpha^j d^i d^j \langle \varphi(\boldsymbol{x}^i), \, \varphi(\boldsymbol{x}^j) \rangle$$

and the classification step

$$y = \sum_{i \text{ (s.v.)}} \alpha^i d^i \langle \varphi(\boldsymbol{x}), \, \varphi(\boldsymbol{x}^i) \rangle + b$$

- May require a feature space of really high (even infinite) dimension.
- In this case, manipulating $\varphi(\boldsymbol{x})$ might be tough, impossible, or lead to intense computation: complexity depends on the feature space dimension.

## Non-linear SVM – Kernel trick

- SVMs do not care about feature vectors $\varphi(\boldsymbol{x})$.
- They rely only on dot products between them. Define

$$K(\boldsymbol{x}, \boldsymbol{x}') = \langle \varphi(\boldsymbol{x}),\ \varphi(\boldsymbol{x}') \rangle$$

- $K$ is called kernel function: dot product in a higher dimensional space.
- Even if $\varphi(\boldsymbol{x})$ is tough to manipulate, $K(\boldsymbol{x}, \boldsymbol{x}')$ can be rather simple.

- **Mercer's theorem:** $K$ continuous and symmetric positive semi-definite (*i.e.*, the matrix $\boldsymbol{K} = (K(\boldsymbol{x}^i, \boldsymbol{x}^j))_{i,j}$ is symmetric positive semi-definite for all finite sequences $\boldsymbol{x}^1, \ldots, \boldsymbol{x}^n$) then there exists a mapping $\varphi$ such that

$$K(\boldsymbol{x}, \boldsymbol{x}') = \langle \varphi(\boldsymbol{x}),\ \varphi(\boldsymbol{x}') \rangle$$

- We do not even need to know $\varphi$, pick a continuous symmetric positive semi-definite kernel $K$ and use it instead of dot products.

## Non-linear SVM – Kernel trick

- Replace all occurrences of $\langle \boldsymbol{x}, \boldsymbol{x}' \rangle$ by $K(\boldsymbol{x}, \boldsymbol{x}')$, for the training step

$$\mathcal{L}_D(\boldsymbol{\alpha}) = \sum_i \alpha^i - \frac{1}{2} \sum_i \sum_j \alpha^i \alpha^j d^i d^j K(\boldsymbol{x}^i, \boldsymbol{x}^j)$$

and the classification step

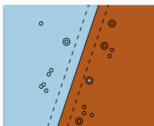$$y = \sum_{i \text{ (s.v.)}} \alpha^i d^i K(\boldsymbol{x}, \boldsymbol{x}^i) + b$$

- Python implementation available in `scikit-learn`.
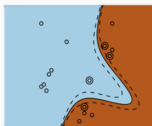
**Complexity depends on the input space dimension.**

## Non-linear SVM – Standard kernels

- Linear: $$K(\boldsymbol{x}, \boldsymbol{x}') = \langle \boldsymbol{x},\ \boldsymbol{x}' \rangle$$

- Polynomial: $$K(\boldsymbol{x}, \boldsymbol{x}') = (\gamma \langle \boldsymbol{x},\ \boldsymbol{x}' \rangle + \beta)^p$$

- Gaussian: $$K(\boldsymbol{x}, \boldsymbol{x}') = \exp\left(-\gamma \|\boldsymbol{x} - \boldsymbol{x}'\|^2\right)$$

  $$\longrightarrow \text{Radial basis function (RBF) network.}$$

- Sigmoid: $$K(\boldsymbol{x}, \boldsymbol{x}') = \tanh(\gamma \langle \boldsymbol{x},\ \boldsymbol{x}' \rangle + \beta)$$
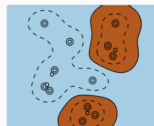
**Linear Kernel**     **Polynomial Kernel**     **RBF Kernel**



**In practice:** $K$ **is usually chosen by trial and error.**

**A linear SVM is an optimal perceptron but what is a non-linear SVM?**

## Non-linear SVMs are shallow ANNs

Consider: $K(\boldsymbol{x}, \boldsymbol{x}') = \tanh(\gamma \langle \boldsymbol{x}, \boldsymbol{x}' \rangle + \beta)$

We get:
$$y = \sum_{i \text{ (s.v.)}} \alpha^i d^i K(\boldsymbol{x}, \boldsymbol{x}^i) + b \qquad \leftarrow \text{SVM}$$

$$y = \sum_{i \text{ (s.v.)}} \alpha^i d^i \tanh(\gamma \langle \boldsymbol{x}, \boldsymbol{x}^i \rangle + \beta) + b$$

$$y = \langle \boldsymbol{w}_2, \, g(\boldsymbol{W}_1 \boldsymbol{x} + \boldsymbol{b}_1) \rangle + b_2 \qquad \leftarrow \text{Shallow ANN}$$

where $\boldsymbol{W}_1 = \gamma \underbrace{\left( \boldsymbol{x}^1 \quad \boldsymbol{x}^2 \quad \dots \right)}_{\text{support vectors}}^T, \quad \boldsymbol{b}_1 = \beta \boldsymbol{1},$

$\boldsymbol{w}_2 = \underbrace{\left( \alpha^1 d^1 \quad \alpha^2 d^2 \quad \dots \right)}_{\substack{\text{Lagrange multipliers and desired} \\ \text{class of support vectors}}}, \quad b_2 = b, \quad \text{and} \quad g = \tanh.$

**The number of support vectors is the number of hidden units.**

## Difference between SVMs and ANNs

- SVM:
  - based on stronger mathematical theory,
  - avoid over-fitting,
  - not trapped in local-minima,
  - works well with fewer training samples.

- But limited to binary classification
  $\rightarrow$ extensions to regression in (Vapnik et al., 1997)
    and PCA in (Schölkopf et al, 1999).

- Tuning SVMs is tough:
  $\rightarrow$ selecting a specific kernel (feature space) and
    regularization parameter $C$ done by trial and errors.

## Similarity between SVMs and ANNs
(in binary classification)

- Similar formalisms:

$$y = \langle \boldsymbol{w},\, \varphi(\boldsymbol{x}) \rangle + b \quad \lessgtr \quad 0 \quad ?$$

**SVM**

- $\varphi(\boldsymbol{x})$: feature vector,

- $\varphi$: non-linear function implicitly defined by the kernel $K$,

- UAT: $\varphi$ could be approximated by a shallow NN.
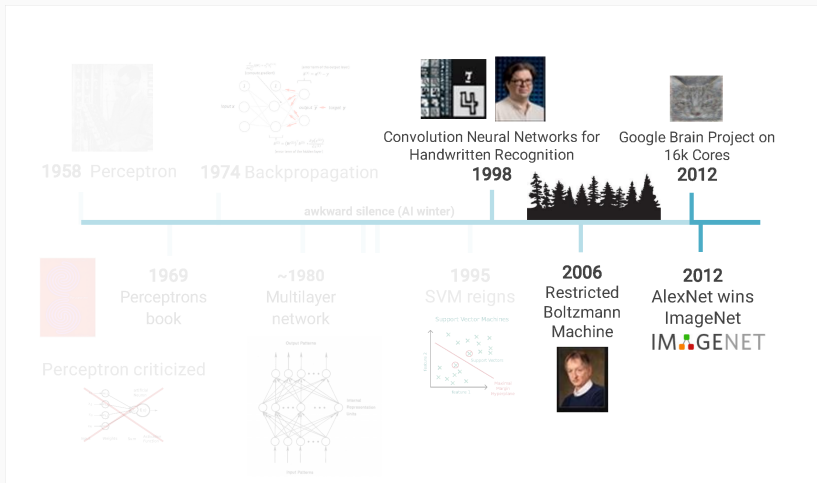
**ANN**

- $\varphi(\boldsymbol{x})$: output of the last hidden layer,

- $\varphi$: non-linear recursive function learned by backpropagation,

- the output of hidden layers can be interpreted as feature vectors.

## Similarity between SVMs and ANNs
(in binary classification)

- Both are linear separators in a suitable feature space,

- Mapping to the feature space is obtained by a non-linear transform,

- For SVM, the non-linear transform is chosen by the user.
  Instead, ANNs learn it with backpropagation.

- ANNs can also do it recursively layer by layer:
    - $\rightarrow$ this is called a feature hierarchy,
    - $\rightarrow$ this is the foundation of deep learning.

# What's next?



*(Source: Lucas Masuch & Vincent Lepetit)*

# Questions?

## Next class: Introduction to deep learning