



Artificial Intelligence II: Deep learning methods

Dragoș Burileanu, Ana Neacșu & Vlad Vasilescu, Georgian Nicolae
Lecture 2: Linear Networks

National University of Science and Technology POLITEHNICA Bucharest, Romania
BIOSINF Master Program

March 2024

Overview

1 Why Deep Learning?

2 Convex Optimization Recap

3 What is a neural network?

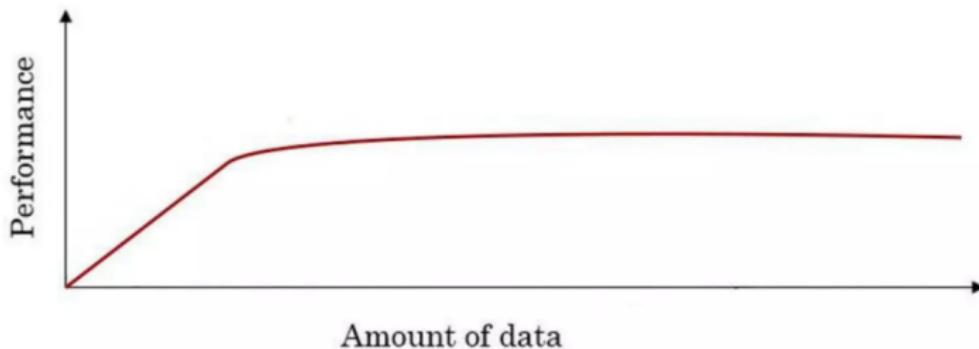
4 Linear Neural Networks

5 Gradient descent

Why Deep Learning?

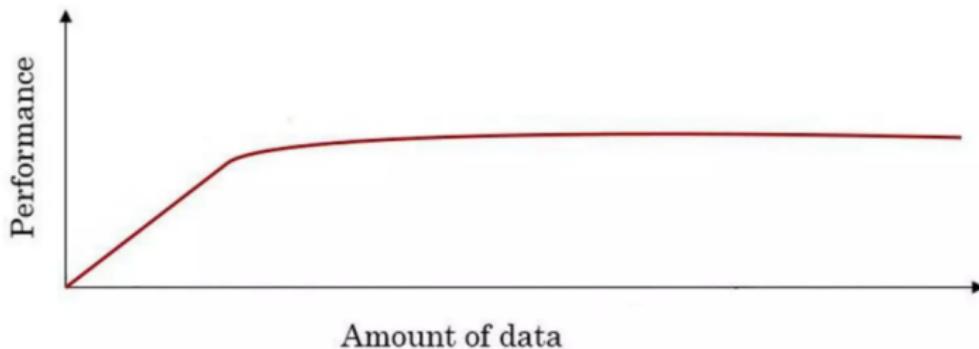
The rise of Deep Learning

- If the basic technical idea behind deep learning neural networks has been around for decades, why are they only taking off now?
- If we plot the performance of traditional algorithms such as SVM or Logistic Regression as function of the amount of data, we will get the following curve:



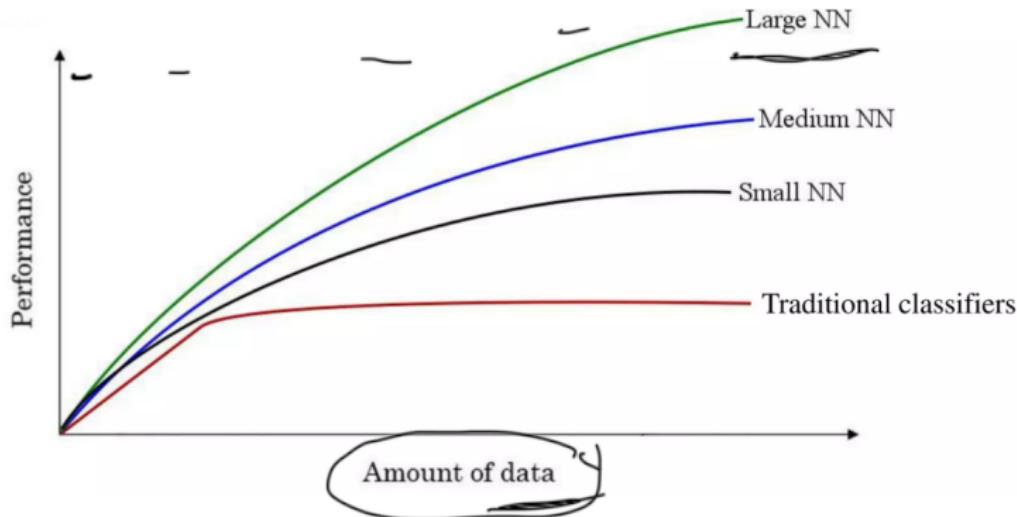
The rise of Deep Learning

- If the basic technical idea behind deep learning neural networks has been around for decades, why are they only taking off now?
- If we plot the performance of traditional algorithms such as SVM or Logistic Regression as function of the amount of data, we will get the following curve:



The rise of deep learning

- How to overcome performance plateau problem?



Why is Deep Learning working now

-
- Over the last 20 years we accumulated more data for applications than traditional learning algorithms were able to effectively take advantage of
 - GPU (speed of processing) + Data - ImageNet - 1M / 1000 class
 - Theoretical understandings of the difficulty of training deep networks (from 2006)

Libraries allow to easily implement/test/deploy neural networks :

- Torch (Lua) / PyTorch (Python/C++), Caffe(C++/Python), Caffe2 (RIP 2018)
- Microsoft CNTK
- Google Tensorflow / Keras
- Theano/Lasagne (Python, RIP 2017)
- CNTK, Chainer, Matlab, Mathematica,

Why is Deep Learning working now

- Over the last 20 years we accumulated more data for applications than traditional learning algorithms were able to effectively take advantage of
- GPU (speed of processing) + Data
- Theoretical understandings of the difficulty of training deep networks (from 2006)

Libraries allow to easily implement/test/deploy neural networks :

- Torch (Lua) / PyTorch (Python/C++), Caffe(C++/Python), Caffe2 (RIP 2018)
- Microsoft CNTK
- Google Tensorflow / Keras
- Theano/Lasagne (Python, RIP 2017)
- CNTK, Chainer, Matlab, Mathematica, ...

Why Deep Learning?
oooo

Convex Optimization Recap
●ooooo

What is a neural network?
oooo

Linear Neural Networks
oooooooooooo

Gradient descent
oooooooooooo

Convex Optimization Recap

The Cauchy-Schwarz inequality

$$\vec{a} \cdot \vec{b} = \|a\| \|b\| \cos \theta \Rightarrow \text{Perceptron}$$

 ω

Let $\mathbf{u}, \mathbf{v} \in \mathbb{R}^d$. The Cauchy-Schwartz inequality reads:

$$|\mathbf{u}^\top \mathbf{v}| \leq \|\mathbf{u}\| \|\mathbf{v}\|.$$

$$-1 \leq \frac{\mathbf{u}^\top \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|} \leq 1$$

Some notations

- $\mathbf{u} = (u_1, \dots, u_d)^\top, \mathbf{v} = (v_1, \dots, v_d)^\top$ d -dimensional column vectors with real entries.
- \mathbf{u}^\top , transpose of \mathbf{u} , a d -dimensional row vector
- $\mathbf{u}^\top \mathbf{v} = \sum_{i=1}^d u_i v_i$, scalar (or inner) product of \mathbf{u} and \mathbf{v} .
- $|\mathbf{u}^\top \mathbf{v}|$, absolute value of $\mathbf{u}^\top \mathbf{v}$
- $\|\mathbf{u}\|_2 = \sqrt{\mathbf{u}^\top \mathbf{u}} = \sqrt{\sum_{i=1}^d u_i^2}$, Euclidean (ℓ_2) norm of \mathbf{u} .

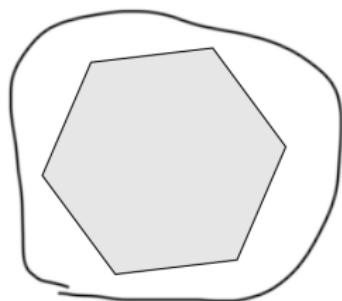
$$\mathbf{u} = \mathbf{v} \Rightarrow \mathbf{u}^\top \mathbf{u} = \|\mathbf{u}\|^2$$

$$\mathbf{u} = -\mathbf{v}$$

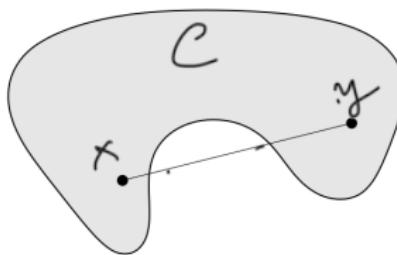
Convex Sets

A set C is **convex** if the line segment between any two points from C lies in C , i.e. for any $\mathbf{x}, \mathbf{y} \in C$ and any $\lambda \in [0, 1]$, we have:

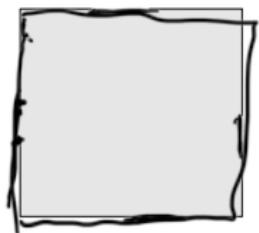
$$\lambda\mathbf{x} + (1 - \lambda)\mathbf{y} \in C$$



convex



non-convex



non-convex

*Figure 2.2 from S. Boyd, L. Vandenberghe

Properties of convex sets

- Let $(C_i)_{i \in I}$ be convex sets, where I is a (possibly infinite) index set. Then $\cap_{i \in I} C_i$ is a convex set.
- Projections onto nonempty closed convex sets are unique, and usually efficient to compute

$$P_C(\mathbf{x}) := \operatorname{argmin}_{\mathbf{y} \in C} \|\mathbf{y} - \mathbf{x}\|$$

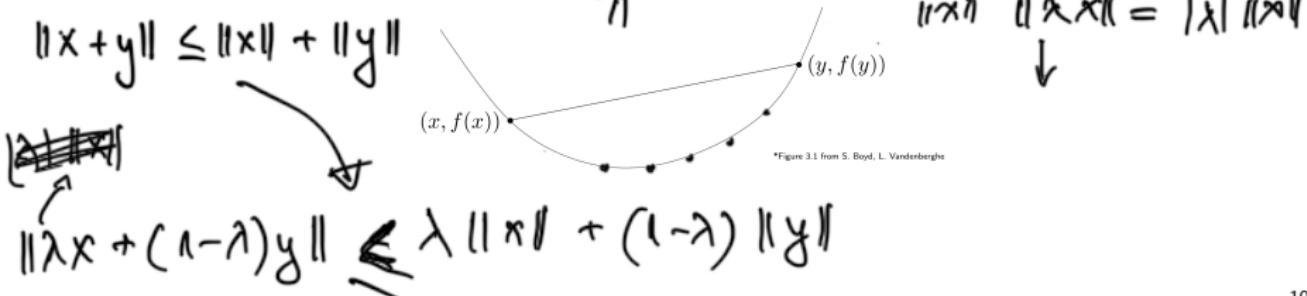
POSS, denoising

Convex functions

Definition: A function $f : \mathbb{R}^d \mapsto]-\infty, +\infty]$ is **convex** if:

i) $\operatorname{dom}(f) = \{x \in \mathbb{R}^d \mid f(x) < +\infty\}$ is a convex set;

ii) $\forall x, y \in \operatorname{dom}(f), \lambda \in [0, 1] \quad f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$



*Figure 3.1 from S. Boyd, L. Vandenberghe

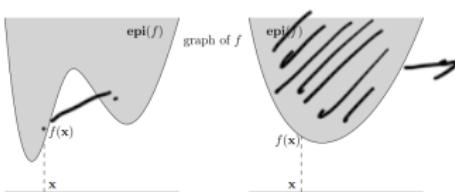
Convex functions and sets

- The **graph** of a function $f : \mathbb{R}^d \mapsto \mathbb{R}$ is defined as

$$\{(x, f(x)) \mid x \in \text{dom}(f)\}$$

- The **epigraph** of a function $f : \mathbb{R}^d \mapsto \mathbb{R}$ is defined as

$$\text{epi}(f) := \{(x, \alpha) \in \mathbb{R}^d \times \mathbb{R} \mid x \in \text{dom}(f), \alpha \geq f(x)\}$$



A function is convex iff its epigraph is a convex set.

Examples of convex functions: $\mathbf{u} \in \mathbb{R}^d, b \in \mathbb{R}$

- Linear functions : $f(\mathbf{x}) = \mathbf{u}^\top \mathbf{x}$ —
- Affine functions : $f(\mathbf{x}) = \underline{\mathbf{u}^\top \mathbf{x} + b}$
- Exponentials : $f(\mathbf{x}) = \underline{e^{\alpha \mathbf{x}}}$

Question: Is norm $\|\mathbf{x}\|$ convex?

Convex Optimization

$$\|x\|_2 = \sqrt{\sum_{i=1}^n |x_i|^2}$$

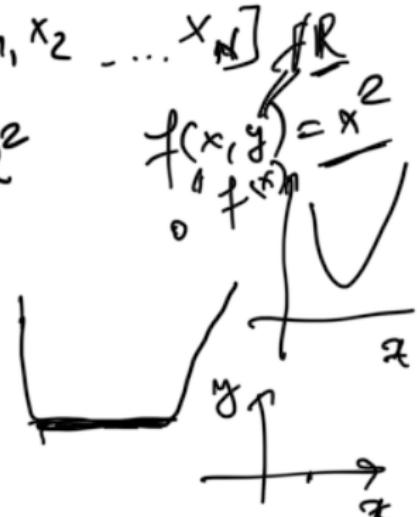
$$x = [x_1, x_2, \dots, x_n] \in \mathbb{R}^n$$

Convex Optimization Problems have the following form:

$$\min f(x) \quad \text{s.t. } x \in C,$$

where

- f is a convex function
- C is a convex set



Properties :

- Every local minimizer is a global minimizer.
- For convex optimization problems, assuming f differentiable, all properly designed algorithms (Gradient Descent, Stochastic Gradient Descent, Projected and Proximal Gradient Descent) do converge to a global minimizer (which is not necessarily unique)!

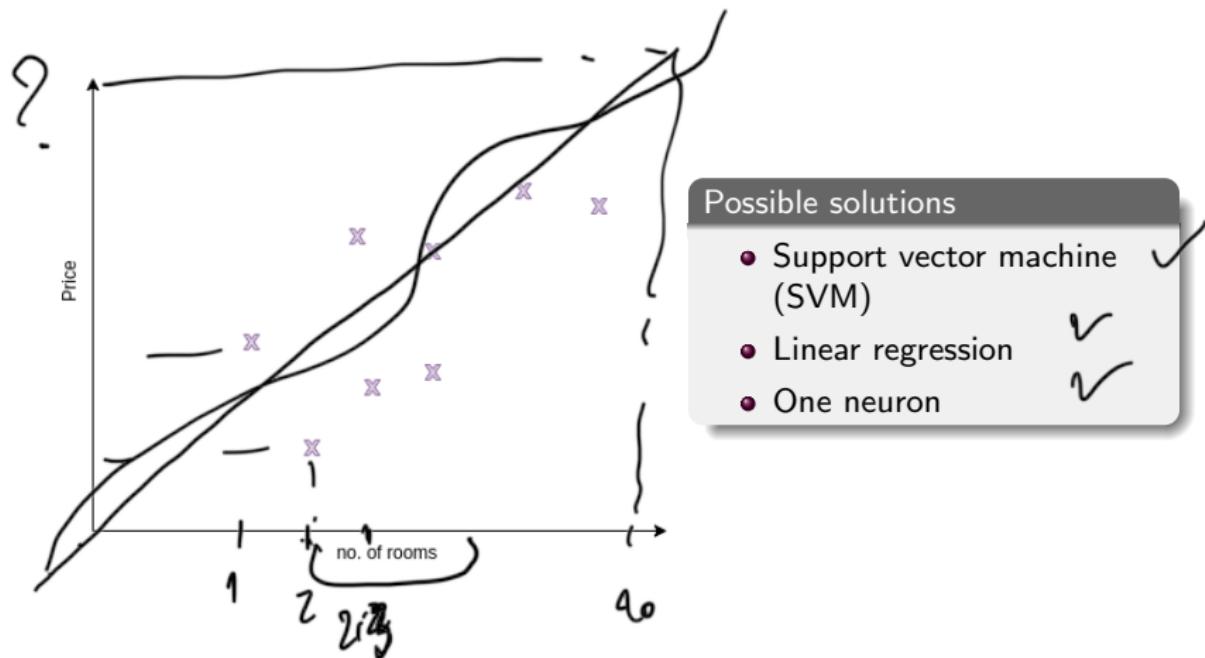
$$f(0, 30) = 0$$

$$f(0, 1524, 3) = 0$$

What is a neural network?

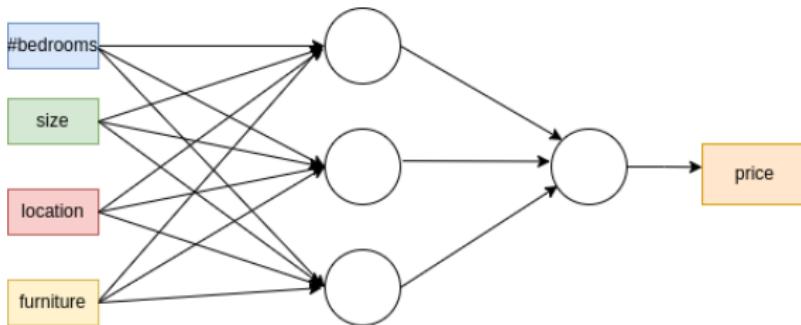
Housing price prediction – binary case

Problem : we want to predict the price of a house based on the number of rooms.



Housing price prediction – complex case

- The previous scenario was not realistic.
- Usually, there are many more factors we should take into consideration

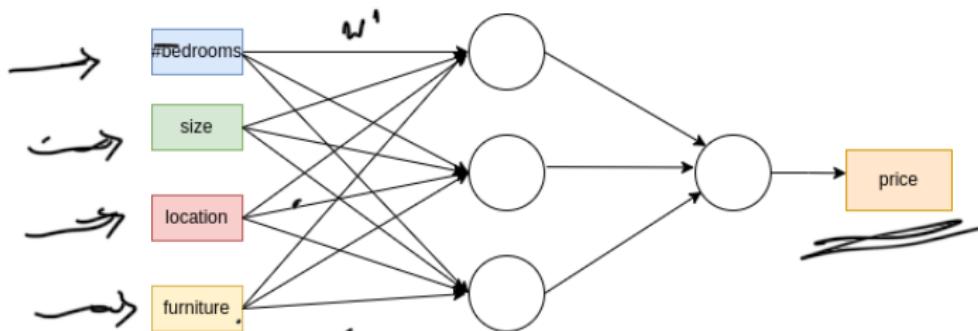


Observations:

- A neural network can combine all this information
- Each factor can influence differently the final decision

Housing price prediction – complex case

- The previous scenario was not realistic.
- Usually, there are many more factors we should take into consideration



Observations:

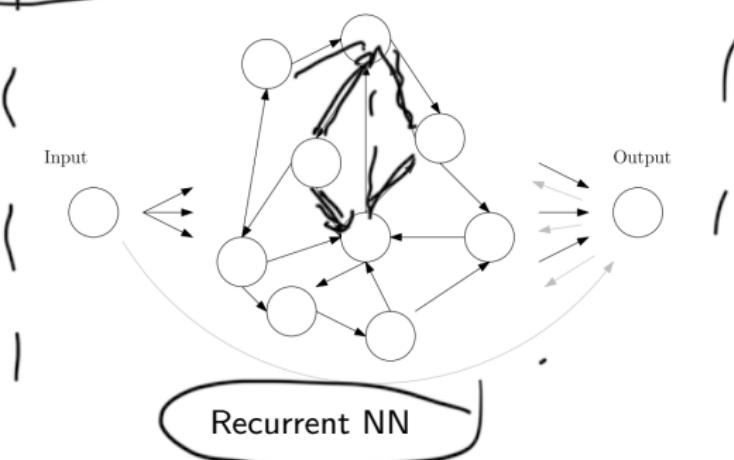
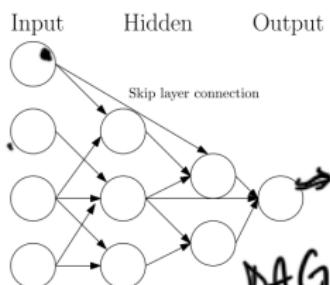
- A neural network can combine all this information
- Each factor can influence differently the final decision

Definition

A **neural network** is a directed graph:

- nodes : computational units
- edges : weighted connections

distill, pub
cyclic



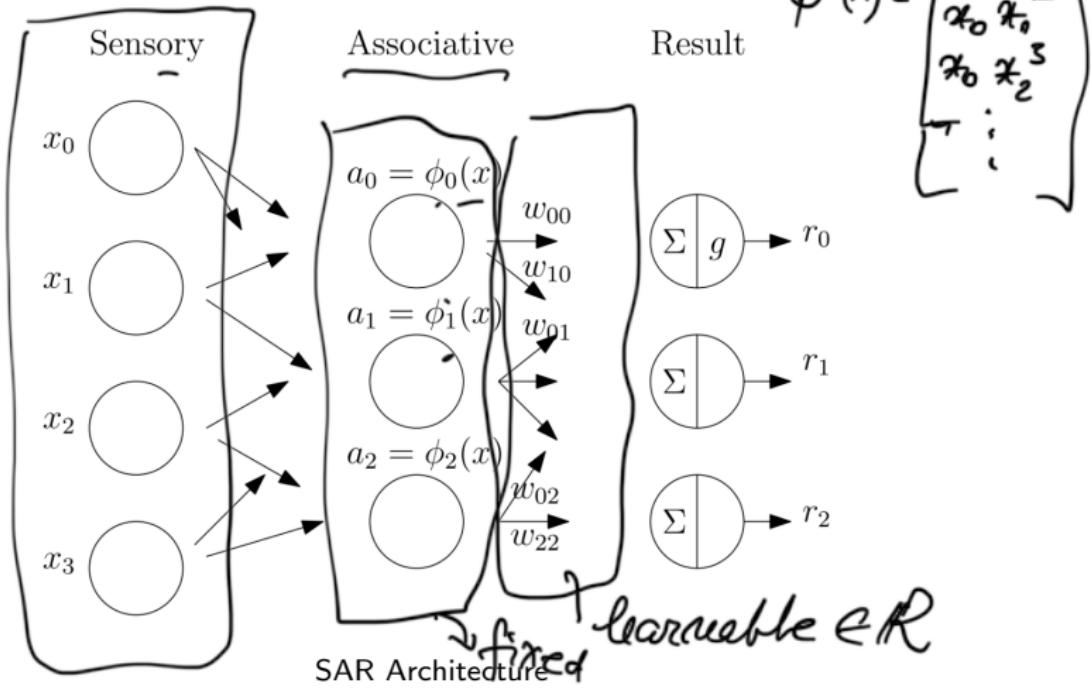
There are two possible types of graphs

- no cycle : feedforward neural network
- with at least one cycle : recurrent neural networks

Linear Neural Networks

The Perceptron (ROSENBLATT, 1958)

- Classification problem: given the pair $(x, y) \in \mathbb{R}^n \times \{-1, 1\}$
- Sensory - Associative - Response architecture, $\phi_j(x)$ with $\phi_0(x) = 1$
- The algorithm also has a geometrical interpretation



The classifier

Given **fixed, predefined** feature functions ϕ_j with $\phi_0(x) = 1, \forall x \in \mathbb{R}^n$, the perceptron classifies the input x as:

$$y = g(w_0 + w_1 \phi_1(x) + w_2 \phi_2(x) + \dots)$$

$y = g(\underline{w^\top \Phi(x)})$

\downarrow bias

$$g(x) = \begin{cases} -1 & \text{if } x < 0 \\ +1 & \text{if } x \geq 0, \end{cases}$$

$$\underline{w^\top v} \quad \underline{\phi(x)} = \begin{bmatrix} 1 \\ \underline{\phi_1} \\ \underline{\phi_2} \\ \vdots \\ \underline{\phi_N} \end{bmatrix} \quad (1)$$

(2)

$$\text{with } \Phi(x) \in \mathbb{R}^{n_a+1}, \phi(x) = \begin{bmatrix} 1 \\ \phi_1(x) \\ \phi_2(x) \\ \vdots \end{bmatrix}.$$

$$g(w^\top \phi(x_i)) = y_i$$

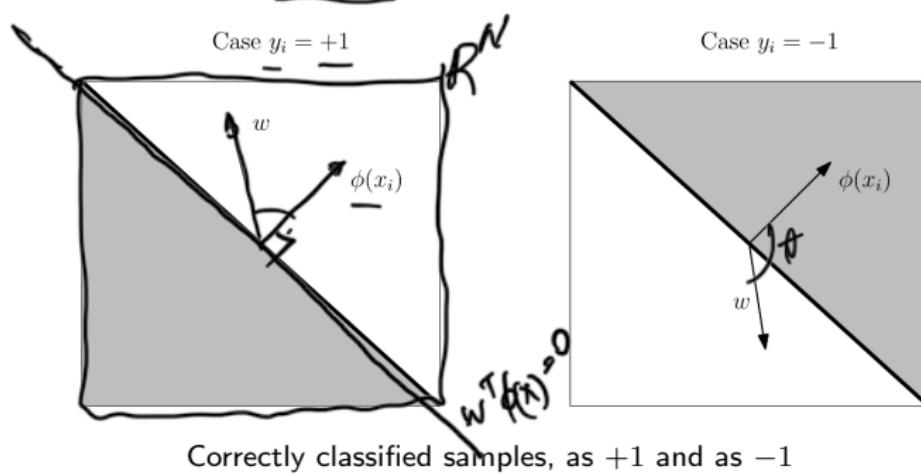
Training algorithm

Given $(x_i, y_i), y_i \in \{-1, 1\}$ the **perceptron learning rule** operates as follows:

$$w \leftarrow \begin{cases} w & \text{if the input is correctly classified} \\ w + \phi(x_i) & \text{if the input is } \underline{\text{incorrectly classified as } -1} \\ w - \phi(x_i) & \text{if the input is } \underline{\text{incorrectly classified as } +1} \end{cases} \quad (3)$$

Correct classification – geometrical interpretation

- Decision rule : $y = g(w^\top \Phi(x))$
- Algorithm : $w \leftarrow \begin{cases} w & \text{if the input is correctly classified} \\ w + \phi(x_i) & \text{if the input is incorrectly classified as } -1 \\ w - \phi(x_i) & \text{if the input is incorrectly classified as } +1 \end{cases}$



$$a \cdot b = \|a\| \|b\| \cos \theta$$

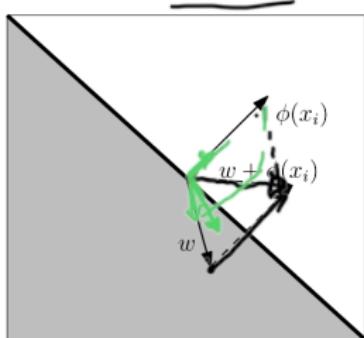
Incorrect classification – geometrical interpretation

- Decision rule : $y = g(w^\top \Phi(x))$

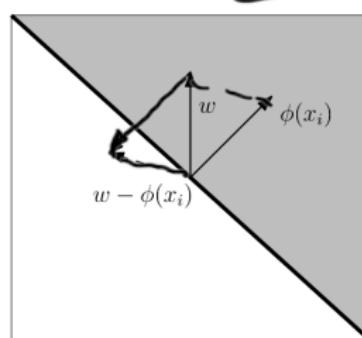
$1 \nexists i, y_i$

- Algorithm : $\underline{w} \leftarrow \begin{cases} w & \text{if the input is correctly classified} \\ w + \phi(x_i) & \text{if the input is incorrectly classified as } -1 \\ w - \phi(x_i) & \text{if the input is incorrectly classified as } +1 \end{cases}$

Case $y_i = +1$



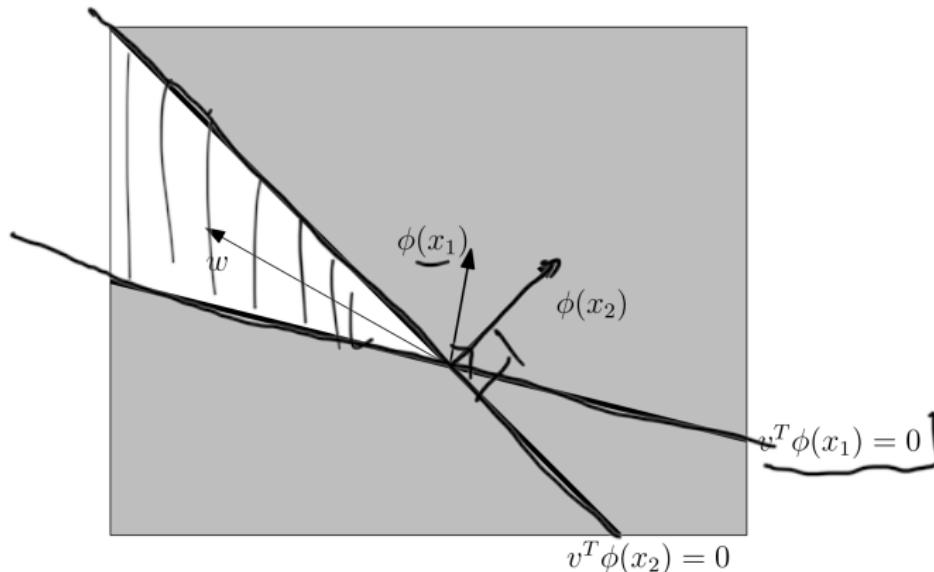
Case $y_i = -1$



Incorrectly classified samples, $+1 \rightarrow -1$ and $-1 \rightarrow +1$

Multiple samples – Geometrical interpretation

- Decision rule : $y = g(w^\top \Phi(x))$
- Cone of feasibility: The intersection of the valid halfspaces (it may be empty)
- We consider two samples x_1, x_2 and $y_1 = +1, y_2 = -1$



The cone of feasibility for $y_1 = +1$ and $y_2 = -1$

Towards a canonical learning rule – **delta rule**

Given $(x_i, y_i), y_i \in \{-1, 1\}$ the **perceptron learning rule** operates as follows:

$$w = \begin{cases} w & \text{if the input is correctly classified} \\ w + \phi(x_i) & \text{if the input is incorrectly classified as } -1 \\ w - \phi(x_i) & \text{if the input is incorrectly classified as } +1 \end{cases} \quad (4)$$


$$w \leftarrow \begin{cases} w & \text{if } g(w^\top \phi(x_i)) = y_i \\ w + \phi(x_i) & \text{if } g(w^\top \phi(x_i)) = -1 \text{ and } y_i = +1 \\ w - \phi(x_i) & \text{if } g(w^\top \phi(x_i)) = +1 \text{ and } y_i = -1 \end{cases} \quad (5)$$

$$w \leftarrow \begin{cases} w & \text{if } g(w^\top \phi(x_i)) = y_i \\ w + y_i \phi(x_i) & \text{if } g(w^\top \phi(x_i)) \neq y_i \end{cases} \quad (6)$$

$$w = w + \frac{1}{2}(y_i - \hat{y}_i)\phi(x_i), \quad (7)$$

with $\hat{y}_i = g(w^\top \phi(x_i))$. This is called the **Delta Rule**.

Towards a canonical learning rule – **delta rule**

Given $(x_i, y_i), y_i \in \{-1, 1\}$ the **perceptron learning rule** operates as follows:

$$w = \begin{cases} w & \text{if the input is correctly classified} \\ w + \phi(x_i) & \text{if the } \underline{\text{input is incorrectly classified}} \text{ as } -1 \\ w - \phi(x_i) & \text{if the input is incorrectly classified as } +1 \end{cases} \quad (4)$$

$$w \leftarrow \begin{cases} w & \text{if } g(w^\top \phi(x_i)) = y_i \\ w + \phi(x_i) & \text{if } g(w^\top \phi(x_i)) = -1 \text{ and } y_i = +1 \\ w - \phi(x_i) & \text{if } g(w^\top \phi(x_i)) = +1 \text{ and } y_i = -1 \end{cases} \quad (5)$$

$$w \leftarrow \begin{cases} w & \text{if } g(w^\top \phi(x_i)) = y_i \\ w + y_i \phi(x_i) & \text{if } g(w^\top \phi(x_i)) \neq y_i \end{cases} \quad (6)$$

$$w = w + \frac{1}{2}(y_i - \hat{y}_i)\phi(x_i), \quad (7)$$

with $\hat{y}_i = g(w^\top \phi(x_i))$. This is called the **Delta Rule**.

Towards a canonical learning rule – delta rule

Given $(x_i, y_i), y_i \in \{-1, 1\}$ the **perceptron learning rule** operates as follows:

$$w = \begin{cases} w & \text{if the input is correctly classified} \\ w + \phi(x_i) & \text{if the input is incorrectly classified as } -1 \\ w - \phi(x_i) & \text{if the input is incorrectly classified as } +1 \end{cases} \quad (4)$$

$$w \leftarrow \begin{cases} w & \text{if } g(w^\top \phi(x_i)) = y_i \\ w + \phi(x_i) & \text{if } g(w^\top \phi(x_i)) = -1 \text{ and } y_i = +1 \\ w - \phi(x_i) & \text{if } g(w^\top \phi(x_i)) = +1 \text{ and } y_i = -1 \end{cases} \quad (5)$$

$$w \leftarrow \begin{cases} w & \text{if } g(w^\top \phi(x_i)) = y_i \\ w + y_i \phi(x_i) & \text{if } g(w^\top \phi(x_i)) \neq y_i \end{cases} \quad (6)$$

$$\boxed{w = w + \frac{1}{2}(y_i - \hat{y}_i)\phi(x_i)}, \quad (7)$$

with $\hat{y}_i = g(w^\top \phi(x_i))$. This is called the **Delta Rule**.

$$\hat{y}_i - y_i$$

$\Rightarrow \pm 2$

Convergence theorem *perceptron*

Definition

A binary classification problem $(x_i, y_i) \in \mathbb{R}^d \times \{-1, 1\}$ and $i \in [1 \dots N]$ is linearly separable if $\exists w \in \mathbb{R}^d$ such that

$$\forall i \quad \text{sign}(w^\top x_i) = y_i, \quad (8)$$

①

Theorem (Perception convergence theorem)

A binary classification problem $(x_i, y_i) \in \mathbb{R}^d \times \{-1, 1\}$ and $i \in [1 \dots N]$ is linearly separable iff perceptron learning rule converges to an optimal solution in a finite number of steps.

Proof : \Leftarrow : easy; \Rightarrow : we upper/lower bound $\|w_t\|_2^2$, where t is the index of the current iteration

Convergence theorem

Definition

A binary classification problem $(x_i, y_i) \in \mathbb{R}^d \times \{-1, 1\}$ and $i \in [1 \dots N]$ is linearly separable if $\exists w \in \mathbb{R}^d$ such that

$$\forall i \text{ sign}(w^\top x_i) = y_i, \quad (8)$$

Theorem (Perception convergence theorem)

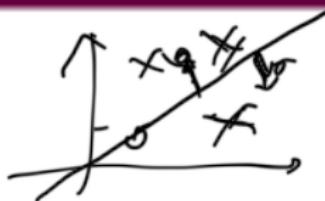
A binary classification problem $(x_i, y_i) \in \mathbb{R}^d \times \{-1, 1\}$ and $i \in [1 \dots N]$ is linearly separable iff perceptron learning rule converges to an optimal solution in a finite number of steps.

Proof : \Leftarrow : easy; \Rightarrow : we upper/lower bound $\|w_t\|_2^2$, where t is the index of the current iteration



Observations

$$w \leftarrow w + \eta x_i, y_i$$



- $w_t = w_0 + \sum_{i \in S_t} y_i \phi(x_i)$, with S_t the set of misclassified samples.
- The cost function to be minimised is: $J(w) = \frac{1}{M} \sum_i \max(0, -y_i w^\top \phi(x_i))$
- The solution:

$$w_t = w_0 + \sum_i \frac{1}{2} (y_i - \hat{y}_i) \phi(x_i),$$

where $(y_i - \hat{y}_i)$ is called the **prediction error**.

$$\hat{y}$$

Hinge loss
w/
margin

ADALINE - Regression

$w^T x + b$

- Introduced by Widrow & Hoff in 1962.
- Let us consider the linear regression problem analytically.

Problem : Given $(x_i, y_i), x_i \in \mathbb{R}^{n+1}, y_i \in \mathbb{R}$, minimize

$$J(w) = \frac{1}{N} \sum_i \|y_i - w^T x_i\|^2$$

$$X \in \mathbb{R}^{N \times (n+1)}$$

$$\begin{bmatrix} 1 & x_0 & x_1 & \dots & x_N \end{bmatrix}^T \in \mathbb{R}^{(n+1) \times N}$$

Here-above we assume that $\forall i \quad x_i[0] = 1$ and $w[0]$ accounts for the bias term, $n \in \mathbb{N}^*$ is the input dimension, while $N \in \mathbb{N}^*$ is the total number of samples.

Analytically, we can vectorize the expression: Assume $X = [x_0 | x_1 | \dots]$, then

$$J(w) = \|y - X^T w\|^2$$

$$\nabla_w J(w) = 0 \Leftrightarrow -2(y - X^T w)^T X^T = 0 \Rightarrow X X^T w = X y$$

eigen(XX^T) \Rightarrow XX^T non-singular : $w = (XX^T)^{-1} X y \Leftrightarrow \det(XX^T) \neq 0$

$$XX^T \in \mathbb{R}^{n \times n}$$

$$XX^T \geq 0$$

$$XX^T \in \mathbb{R}^{n \times n}$$

$$XX^T \geq 0$$

$$XX^T \in \mathbb{R}^{n \times n}$$

- XX^T singular (e.g., points along a line in 2D) \rightarrow infinite no. solutions
 - Use regularized least square method to solve the problem:

$$\nabla_w G = \nabla_w J + \alpha w$$

$$\nabla_w G(w) = 0 \Rightarrow (XX^T + \alpha I_d)w = X y$$

$$\text{if } \alpha \in (0, +\infty[\Rightarrow (XX^T + \alpha I_d) \text{ is non-singular}$$

$$X = \begin{bmatrix} x_0 & x_1 & \dots \end{bmatrix}^T$$

$$XX^T = \sum$$

Observation : We need to compute XX^T over the whole training set!

Linear regression with stochastic gradient descent

~~ADALINE~~

Problem : Given $(x_i, y_i), x_i \in \mathbb{R}^{n+1}, y_i \in \mathbb{R}$, minimize

$$J(w) = \frac{1}{N} \sum_i \|y_i - w^\top x_i\|^2.$$

$$x_{k+1} = x_k - \eta f'(x_k)$$

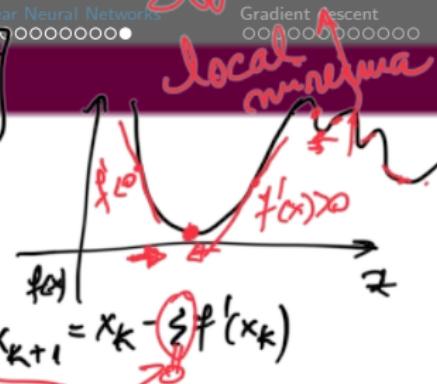
Here-above we assume that $\forall i \quad x_i[0] = 1$ and $w[0]$ accounts for the bias term.

~~Algorithm 1: Stochastic Gradient Descent Algorithm~~

Output: Trained weights w

- 1 Initialize w_0 randomly;
- 2 **for** $t = 1$ **to** T **do** — # iterations
- 3 **for** $i = 1$ **to** N **do** — # example
- 4 $\hat{y}_i \leftarrow w_t^\top x_i$;
- 5 $w_t \leftarrow w_{t-1} - \epsilon \nabla_w J(w_{t-1}) = w_{t-1} + \epsilon(y_i - \hat{y}_i)x_i$;
- 6 **endfor**
- 7 **endfor**

where $\epsilon \in]0, \infty[$ is the learning rate.



Gradient descent

Batch Gradient Descent

The cost function to be minimized is

$$J(w, x, y) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(w, \underline{x}_i, \underline{y}_i),$$

where \mathcal{L} is the *loss function*, e.g. $\mathcal{L}(w, x_i, y_i) = \|y_i - w^\top x_i\|^2$

Batch gradient descent

- compute the gradient of the cost $J(w)$ over the whole training set
 - performs one step in direction of $-\nabla_w J(w, x, y)$, so the update rule becomes:

$$w_{t+1} = w_t - \epsilon_t \nabla_w J(w, x, y)$$

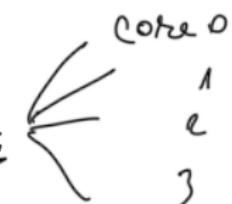
Algorithm 2: Batch Gradient Descent Algorithm

Output: Trained weights w

- ```

1 Initialize w_0 randomly;
2 for $t = 1$ to T do # iterations
3 | Compute gradient $\nabla_w J(w_t; x, y)$ using the entire dataset
4 | Update $w_{t+1} \leftarrow w_t - \epsilon \nabla_w J(w_t; x, y)$;
5 endfor

```



# Stochastic Gradient Descent

SGD

→ converges faster -

→ Not GPU friendly

The cost function to be minimized is



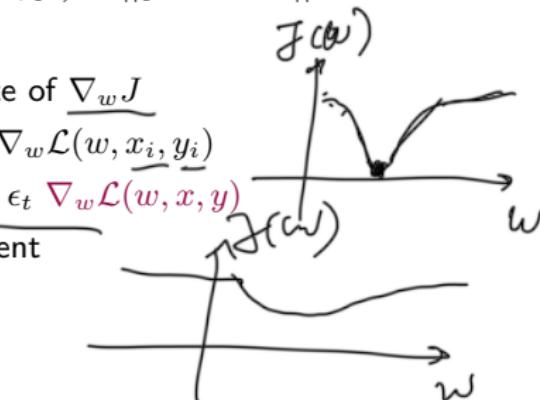
$$J(w, x, y) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(w, x_i, y_i),$$

where  $\mathcal{L}$  is the *loss function*, e.g.  $\mathcal{L}(w, x_i, y_i) = \|y_i - w^\top x_i\|^2$

## Stochastic gradient descent (SGD)

- one sample at a time, noisy estimate of  $\nabla_w J$
- performs one step in direction of  $-\nabla_w \mathcal{L}(w, x_i, y_i)$
- converges faster than gradient descent

$$\underbrace{w_{t+1} = w_t - \epsilon_t \nabla_w \mathcal{L}(w, x, y)}$$

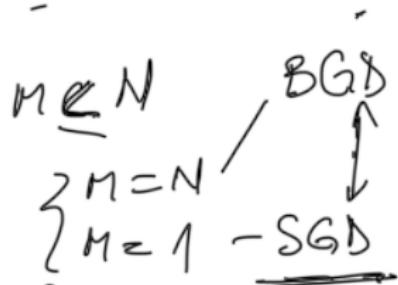


# Minibatch gradient descent



The cost function to be minimized is

$$J(w, x, y) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(w, x_i, y_i),$$



## Minibatch gradient descent

- noisy estimate of the true gradient with  $M$  samples, with  $M$  being the *minibatch size*
- randomize  $\mathcal{I}$  with  $|\mathcal{I}| = M$ , one set at a time

$$w_{t+1} = w_t - \epsilon_t \frac{1}{M} \sum_{j \in \mathcal{I}} \nabla_w \mathcal{L}(w, x_j, y_j)$$

- creates a smoother estimate than SGD
- great for parallel architectures (GPU)

**Observation :** if the batch size is too large, there is a generalization gap.

# Why use gradient descent?

$$\nabla f(x) = \begin{pmatrix} 8 \\ -1 \end{pmatrix}$$

$$\nabla^2 f(x) = \begin{pmatrix} 2 & 1 \\ 1 & 3 \end{pmatrix}$$

Convexity



$x \in \mathbb{R}$

A function  $f : \mathbb{R}^n \leftarrow \mathbb{R}$  is convex :

- ①  $\iff \forall x_1, x_2 \in \mathbb{R}^n, \forall \alpha \in [0; 1] f(\alpha x_1 + (1 - \alpha)x_2) \leq \alpha f(x_1) + (1 - \alpha)f(x_2)$
- ② with  $f$  twice diff.,  $\iff \forall x \in \mathbb{R}^n, \mathbf{H} = \nabla^2 f(x)$  is positive semidefinite, i.e.  $\forall x \in \mathbb{R}^n, x^\top \mathbf{H} x \geq 0$ .

## Observations

- For a convex function  $f$ , all local minima are global minima
- Under mild conditions gradient descent and stochastic gradient descent converge
- Typically,  $\sum \epsilon_t = \infty$  and  $\sum \epsilon_t^2 \leq \infty$

# Linear Regression – Recap

**Problem :** Given  $(x_i, y_i), x_i \in \mathbb{R}^{n+1}, y_i \in \mathbb{R}$

- We assume that  $x[0] = 1$  to encompass the bias term.
- We have a **linear model** :  $\hat{y} = w^\top x$
- We consider loss :  $\mathcal{L}(\hat{y}, y) = \|\hat{y} - y\|^2$ .
- Solve using gradient descent,

$$\phi(x)$$

$$\nabla_w \mathcal{L}(w, x_i, y_i) = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w} = -\frac{1}{2} (y_i - \hat{y}_i) x_i$$

## Observations

- Other choices for the loss function may also be considered, e.g. Huber loss, MAE, MSE ...
- We can also include a regularization term (we'll discuss this later).
- Linear regression with  $\ell_2$  loss is convex:

$$\mathcal{L}(w) = \frac{1}{2} (w^\top x_i - y_i)^2$$

$$\nabla_w \mathcal{L} = (w^\top x_i - y_i) x_i$$

$$\nabla_w^2 \mathcal{L} = x_i x_i^\top$$

$U^\top U = \|U\|^2 \geq 0$

$\forall x \in \mathbb{R}^n x^\top x_i x_i^\top x = (x_i^\top x)^2 \geq 0$

# Linear classification – Recap

$$\overbrace{? \rightarrow \cdot \cdot \cdot \cdot \cdot}$$



Consider the Maximum Likelihood in binary classification task:

**Problem :** Given  $(x_i, y_i), x_i \in \mathbb{R}^{n+1}, y_i \in \{0, 1\}$

We consider our samples to be independent and we consider the conditional probability to be parametrized by  $w$ ,  $P(y = 1|x) = p(x, w)$ .

The conditional likelihood of the labels is:

$$L(w) = \prod_i P(y = y_i | x_i) = \prod_i \underbrace{p(x_i; w)^{y_i}}_{\text{underbrace}} \underbrace{(1 - p(x_i; w))^{1-y_i}}_{\text{underbrace}}$$

Usually, we prefer to minimize the averaged negative log-likelihood:

$$J(w) = -\frac{1}{N} \log(L(w)) = \underbrace{\frac{1}{N} \sum_i}_{\text{underbrace}} -y_i \log(p(x_i; w)) - (1 - y_i) \log(1 - p(x_i; w))$$

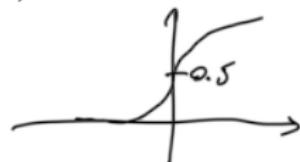
# Binary classification

**Problem :** Given  $(x_i, y_i), x_i \in \mathbb{R}^{n+1}, y_i \in \{0, 1\}$

- **Linear logit model :**  $g(x) = w^\top x \in \mathbb{R}$
- Use Sigmoid transfer function :  $\hat{y}(x) = \sigma(g(x)) = \sigma(w^\top x)$ , where

$$\sigma : \mathbb{R} \mapsto [0, 1], \quad \sigma(x) = \frac{1}{1 + e^{-x}}$$

$$\frac{\partial}{\partial x} \sigma(x) = \sigma(x)(1 - \sigma(x))$$



- Cross-entropy loss (also called negative log-likelihood) :

$$\mathcal{L}(\hat{y}, y) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$$

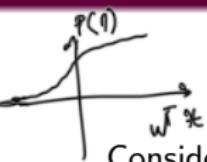
$$\|\hat{y} - y\|^2$$

- The gradient of CE is easy to compute :

$$\nabla_w \mathcal{L}(w, x_i, y_i) = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w} = -(y_i - \hat{y}_i) x_i$$

**Question:** Is Logistic regression convex for  $\ell_1$  and  $\ell_2$  norms?

# Should we use $\ell_2$ as a loss?



Consider  $\ell_2$  loss  $\mathcal{L} = \frac{1}{2} \|\hat{y} - y\|^2$  and the "linear" model

$$\hat{y} = \sigma(w^\top x_i)$$

$$y = 1 \\ w^\top x \leq 0$$

Let us compute the gradient w.r.t.  $w$ :

$$\nabla_w \mathcal{L}(w, x_i, y_i) = \frac{\partial \mathcal{L}}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial w} = -(y_i - \hat{y}_i) \underbrace{\sigma(w^\top x_i)(1 - \sigma(w^\top x_i))}_{\text{check}} x_i$$

## Observations

- If  $x_i \rightarrow$  strongly misclassified (e.g.  $y_i = 1, w^\top x_i \rightarrow +\infty$ ), then  $\sigma(w^\top x_i)(1 - \sigma(w^\top x_i)) \approx 0$  and  $\nabla_w \mathcal{L}(w, x_i, y_i) \approx 0$   
step-size is very small.
- If we use CE loss, then  $\nabla_w \mathcal{L}(w, x_i, y_i)$  is proportional with the error.

flat gradient problem

# Multiclass Classification (Softmax regression)

**Problem :** Given  $(x_i, y_i), x_i \in \mathbb{R}^{n+1}, y_i \in \{0, \dots, K-1\}$ , where  $K > 2$  is the total number of classes.

Assume that the samples are independent and the conditional probability for a class  $c$  is  $P(y = c|x) = \frac{e^{(w_c^\top x)}}{\sum_k e^{(w_k^\top x)}}$ , parameterized by  $(w_k)_{1 \leq k \leq K}$ .

The conditional likelihood of the labels is:

$$L(w) = \prod_i P(y = y_i|x_i)$$

$$\underbrace{k \left[ \begin{array}{c} \\ \\ \\ \end{array} \right] \bar{z} = 1}_{\in [0, 1]}$$

We can also minimize the averaged negative log-likelihood:

$$\overbrace{J(w) = -\frac{1}{N} \log(L(w))} = -\frac{1}{N} \sum_i \log(P(y = y_i|x_i))$$

Usually we use one-hot encoding of the target class (i.e.

$y_i = [0, \dots, 0, \underbrace{1}_\text{at index } i, 0, \dots, 0]$ ), the cost function can be expressed as:

$$J(w) = -\frac{1}{N} \log(L(w)) = -\frac{1}{N} \sum_i \sum_c y_{ci} \log(P(y = c|x_i))$$

# Softmax vs. CE loss – numerical stability issues

## Large exponentials

- if we compute naively the softmax  $\rightarrow \exp(\cdot)$  may have large values
- We can use the following property:

$$\text{softmax}(g_1, g_2, \dots) = \text{softmax}(g_1 - g', g_2 - g', \dots) = \left( \frac{e^{g_i - g'}}{\sum_j e^{g_j - g'}} \right) \underbrace{\frac{e^{g_i}}{e^{g_j}}}^{\sim}$$

, where  $g'$  can be a constant, but usually  $g' = \max(x)$ , in order to have  $g_j - g' \leq 0$

## Avoiding some exponentials with the log-sum-exp trick

$$\log\left(\sum_j e^{g_j}\right) = g' + \log\left(\sum_j e^{g_j - g'}\right)$$

- No need to compute  $\log(\hat{y}_i) = \log(\text{softmax}_j(x))$ :

$$\log(\hat{y}_i) = \log\left(\frac{e^{g_i - g'}}{\sum_j e^{g_j - g'}}\right) = g_i - g' - \log\left(\sum_j e^{g_j - g'}\right)$$

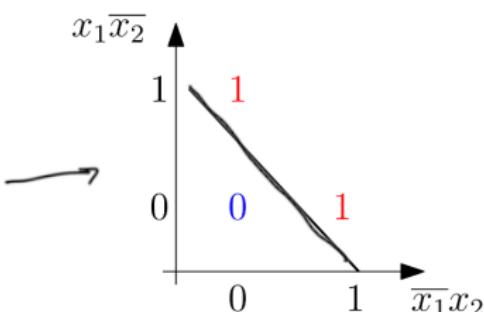
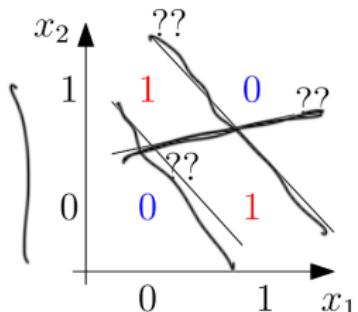
- This is why in practice we use CE loss with logits rather than Softmax + negative log-likelihood.

# The limitations of Linear classification

Perceptrons and logistic regression perform linear separation in a **predefined, fixed** feature space.

$$\text{xor}(x_1, x_2)$$

$$\text{xor}(x_1, x_2) = x_1 \overline{x_2} + \overline{x_1} x_2$$



XOR and its transformation

Can we learn these features ?

$$\varphi(\vec{x}) = \begin{bmatrix} \bar{x}_1 & x_2 \\ x_1 & \bar{x}_2 \end{bmatrix}$$