

# Introduction to Deep Learning Topics: Multi-layer-Perceptron

Geoffroy Peeters

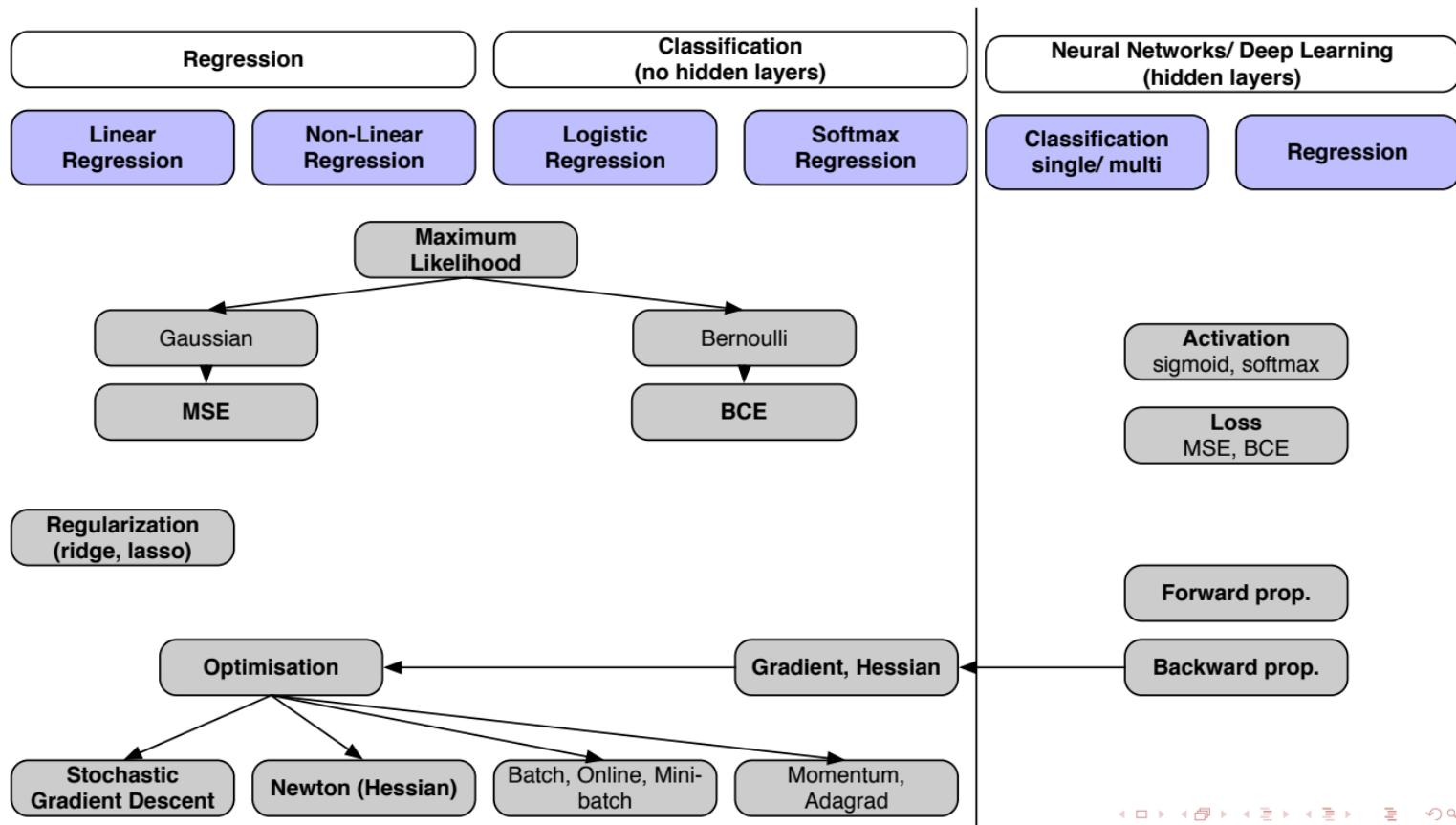
LTCI, Télécom Paris, IP Paris



- 1. Introduction**
- 2. Supervised classification**
- 3. Linear regression**
  - 3.1 ... 1 dimension
  - 3.2 ...  $d$  dimensions
  - 3.3 Parameter estimation
- 4. Reminder from probability**
  - 4.1 Univariate Gaussian distribution
  - 4.2 Expectation
  - 4.3 Covariance
  - 4.4 Multi-dimensional Gaussian distribution
  - 4.5 Bivariate Gaussian distribution
- 5. Maximum Likelihood Estimation (MLE)**
  - 5.1 Example
  - 5.2 Likelihood for linear Regression
  - 5.3 Making predictions (provides a confidence in prediction !)
  - 5.4 Bernoulli distribution
  - 5.5 Bernoulli distribution
- 6. Regression : Ridge, Lasso, Non-linear (basis, kernel)**
  - 6.1 Ridge regularization
  - 6.2 LASSO regularization
  - 6.3 Non-linear regression (regression with basis functions)
  - 6.4 Kernel regression with RBFs
  - 6.5 Choice of the kernel size  $\lambda$
- 7. Optimisation**
  - 7.1 Gradient
  - 7.2 Hessian
  - 7.3 Calculus background : Chain rule
  - 7.4 Gradient and Hessian for linear regression

- 7.5 Optimisation : 1) Steepest/Stochastic Gradient Descent (SGD) algorithm**
- 7.6 Optimisation : 2) Newton/Hessian algorithm**
- 7.7 Why "Stochastic" in SGD ?**
- 7.8 Variations**
- 8. Logistic Regression**
  - 8.1 Mc Culloch - Pitts model of a neuron
  - 8.2 Activation function
  - 8.3 Linear separating hyper-plane
  - 8.4 Cost function
  - 8.5 1) Optimization using Gradient
  - 8.6 2) Optimization using Hessian (Iteratively Reweighted Least Square - IRLS)
- 9. Softmax regression**
  - 9.1 Activation function
  - 9.2 Cost function
  - 9.3 Optimization : derivative of the cost w.r.t. the parameters (manual)
- 10. Backpropagation**
  - 10.1 Logistic Regression
  - 10.2 Layer specification
- 11. Neural Networks - Multi-Layer-Perceptron (MLP)**
  - 11.1 ... with 1 layer (0 hidden layer) / Binary classification
  - 11.2 ... with 2 layers (1 hidden layer) / Binary classification
  - 11.3 ... with 2 layers (1 hidden layer) / Regression
  - 11.4 ... with 2 layers (1 hidden layer) / Multi-Class classification
- 12. Deep Learning**
  - 12.1 Forward/ Backward propagation
  - 12.2 Two ways to implement gradient computation
  - 12.3 Linear layer
  - 12.4 ReLu layer

# Overview





# Supervised classification

## Training/Learning

- Given a set of  $m$  input/output examples :

- $\underline{x}^{(1:m)} = \{\underline{x}^{(1)}, \underline{x}^{(2)}, \dots, \underline{x}^{(m)}\}$ 
  - with  $\underline{x}^{(i)} \in \mathbb{R}^d$
- $\underline{y}^{(1:m)}$ 
  - $y^{(i)} \in \mathbb{R}$  (regression)
  - $y^{(i)} \in \{0, 1\}$  (binary classification)
  - $y^{(i)} \in \{0, 1\}^C$  (multi-classes classification)

- Estimate the parameters  $\underline{\theta}$  of the model :

$$\{x^{(1:m)}, y^{(1:m)}\} \rightarrow \boxed{\text{Learner}} \rightarrow \underline{\theta}$$

## Testing/ Prediction

- Given a new input  $x^{(m+1)}$  and the parameters  $\underline{\theta}$  of the model
- Predict the output  $y^{(m+1)}$

$$\{x^{(m+1)}, \underline{\theta}\} \rightarrow \boxed{\text{Prediction}} \rightarrow \hat{y}^{(m+1)}$$

# Linear regression

... 1 dimension

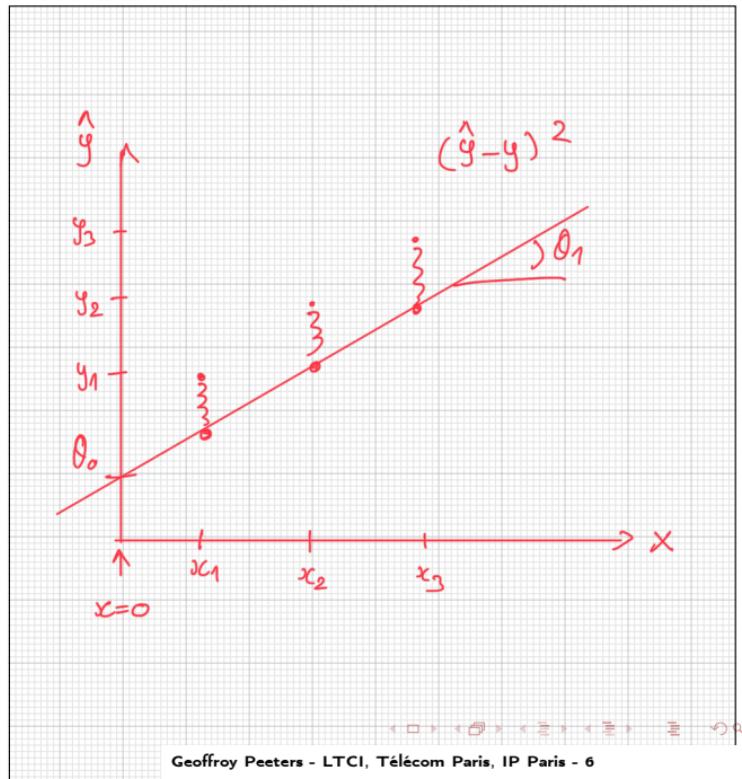
## Model

$$\hat{y}^{(i)} = \theta_0 + x^{(i)}\theta_1$$

Objective function (energy, loss) :

- Mean Square Error (MSE)

$$\begin{aligned} J(\theta) &= \sum_{i=1}^m (y^{(i)} - \hat{y}^{(i)})^2 \\ &= \sum_{i=1}^m (y^{(i)} - (\theta_0 + x^{(i)}\theta_1))^2 \end{aligned}$$



# Linear regression

...  $d$  dimensions

Model :

$$\begin{aligned}\hat{y}^{(i)} &= \theta_0 + \sum_{j=1}^d x_j^{(i)} \theta_j \\ &= \theta_0 + x_1^{(i)} \theta_1 + x_2^{(i)} \theta_2 + \dots + x_d^{(i)} \theta_d\end{aligned}$$

- If we note  $x_0^{(i)} = 1$ , then

$$\hat{y}^{(i)} = \sum_{j=0}^d x_j^{(i)} \theta_j$$

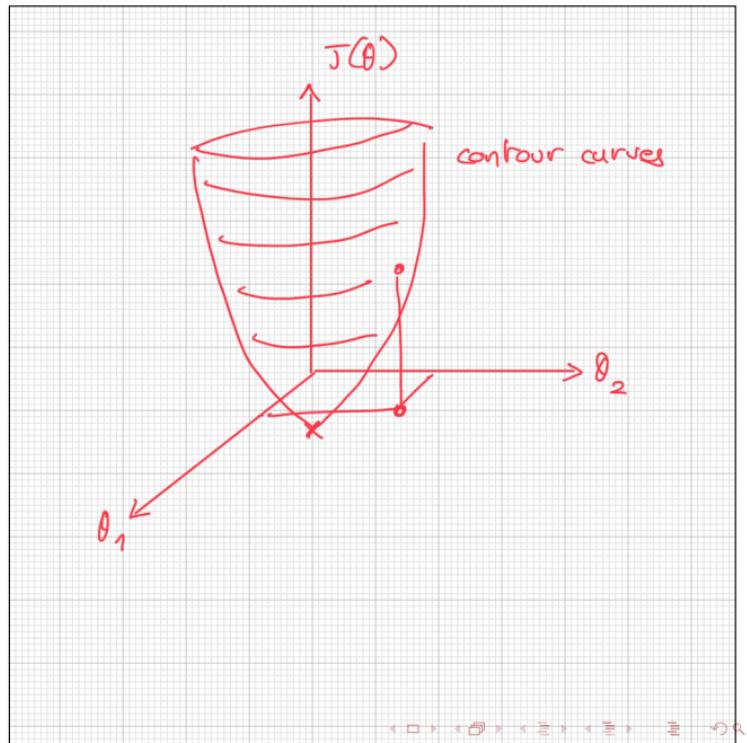
- In matrix form

$$\begin{bmatrix} \hat{y} \\ (\mathbb{R}^{(m,1)}) \end{bmatrix} = \begin{bmatrix} X \\ (\mathbb{R}^{(m,d+1)}) \end{bmatrix} \begin{bmatrix} \theta \\ (\mathbb{R}^{(d+1,1)}) \end{bmatrix}$$

$$\begin{bmatrix} \hat{y}^{(1)} \\ \vdots \\ \hat{y}^{(m)} \end{bmatrix} = \begin{bmatrix} x_0^{(1)} \dots x_d^{(1)} \\ \vdots \ddots \vdots \\ x_0^{(m)} \dots x_d^{(m)} \end{bmatrix} \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_d \end{bmatrix}$$

Objective function

$$\begin{aligned}J(\underline{\theta}) &= \sum_{i=1}^m (y^{(i)} - \underline{x}^T \underline{\theta})^2 \\ &= (\underline{y} - \underline{X} \underline{\theta})^T (\underline{y} - \underline{X} \underline{\theta})\end{aligned}$$



# Linear regression

## Parameter estimation

Find  $\underline{\theta}$  that minimizes  $J(\underline{\theta})$

$$\begin{aligned}\frac{\partial J(\underline{\theta})}{\partial \underline{\theta}} &= 0 = \frac{\partial}{\partial \underline{\theta}} (\underline{y} - \underline{\underline{X}} \underline{\theta})^T (\underline{y} - \underline{\underline{X}} \underline{\theta}) \\ &= \frac{\partial}{\partial \underline{\theta}} (\underline{y}^T \underline{y} - 2 \underline{y}^T \underline{\underline{X}} \underline{\theta} + \underline{\theta}^T \underline{\underline{X}}^T \underline{\underline{X}} \underline{\theta}) \\ &= 0 - 2 \underline{\underline{X}}^T \underline{y} + 2 \underline{\underline{X}}^T \underline{\underline{X}} \underline{\theta} \\ \underline{\theta} &= \boxed{(\underline{\underline{X}}^T \underline{\underline{X}})^{-1} \underline{\underline{X}}^T \underline{y}}\end{aligned}$$

Results from matrix differentiation

$$\begin{aligned}\frac{\partial A\theta}{\partial \theta} &= A^T \\ \frac{\partial \theta^T A\theta}{\partial \theta} &= 2A^T\theta\end{aligned}$$



## Univariate Gaussian distribution

### Probability Density Function (pdf)

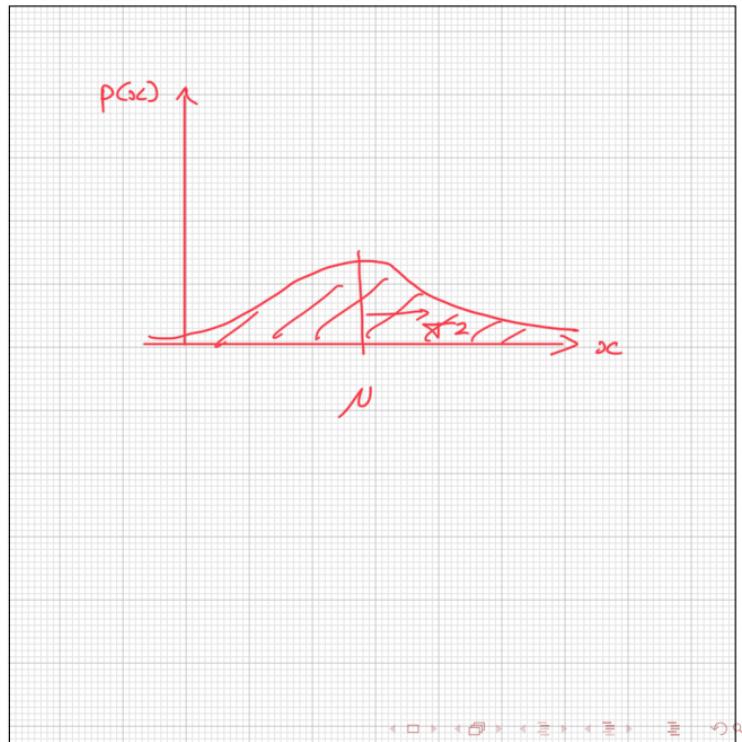
$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}$$

Parameters (sufficient statistics)

- $\mu$  : mean / center of mass
- $\sigma^2$  : variance

Properties

$$\int_{-\infty}^{+\infty} p(x) dx = 1$$



## Expectation

### Expectation

$$\mathbb{E}(X) = \int x p(x) dx = \mu \approx \frac{1}{m} \sum_{i=1}^m X^{(i)}$$

- When  $\underline{X}$  has  $d$  dimensions

$$\underline{\mu} = \mathbb{E}(\underline{X}) = \begin{bmatrix} \mu_{x_1} \\ \mu_{x_2} \\ \vdots \\ \mu_{x_d} \end{bmatrix}$$

$$\begin{aligned} p(x_c) &= \frac{1}{m} \sum_{i=1}^m \delta(x^{(i)}) \\ \int x \cdot \frac{1}{m} \sum_{i=1}^m \delta(x^{(i)}) dx \\ &= \frac{1}{m} \sum_{i=1}^m x^{(i)} \end{aligned}$$

## Covariance

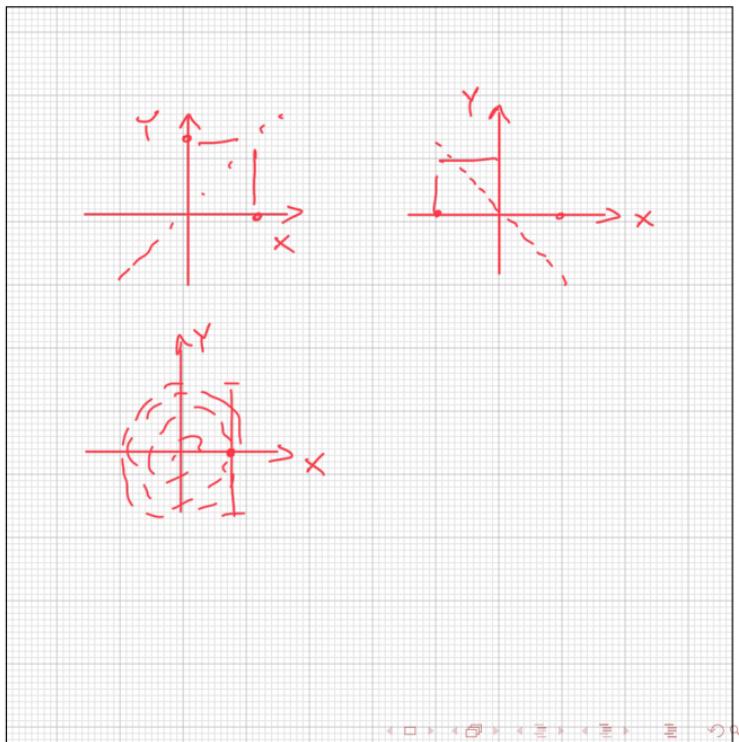
## Covariance

$$\text{cov}(X, Y) = \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$$

- When  $X$  has  $d$  dimensions

$$\begin{aligned}\underline{\Sigma} &= cov[\underline{X}] \\ &= \mathbb{E}[(\underline{X} - \mathbb{E}[\underline{X}])(\underline{X} - \mathbb{E}[\underline{X}])^T] \\ &= \begin{pmatrix} var(X_1) & cov(X_1, X_2) & \dots & cov(X_1, X_d) \\ cov(X_2, X_1) & var(X_2) & \dots & cov(X_2, X_d) \\ \vdots & \vdots & \ddots & \vdots \\ cov(X_d, X_1) & \dots & \dots & var(X_d) \end{pmatrix}\end{aligned}$$

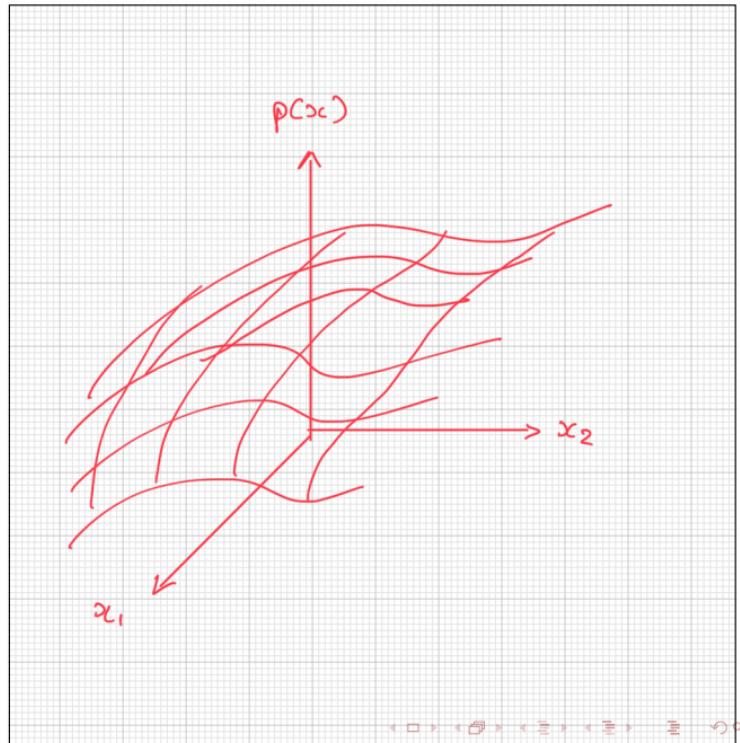
- with  $\sigma_{X_1}^2 = \text{var}(X_1)$



## Multi-dimensional Gaussian distribution

Probability Density Function (pdf)

$$\mathcal{N}(\underline{x}|\underline{\mu}, \underline{\Sigma}) = \frac{1}{(2\pi)^d/2 |\underline{\Sigma}|^{1/2}} e^{-\frac{1}{2}(\underline{x}-\underline{\mu})^T \underline{\Sigma}^{-1} (\underline{x}-\underline{\mu})}$$



## Bivariate Gaussian distribution

**Specific case** : we suppose that  $x_1$  and  $x_2$  are two **independent** Gaussian variables :

- $x_1 \sim \mathcal{N}(\mu_1, \sigma^2)$
- $x_2 \sim \mathcal{N}(\mu_2, \sigma^2)$

Their joint distribution is

$$p(x_1, x_2) = |2\pi\Sigma|^{-1/2} \exp\left(-\frac{1}{2}[\underline{x} - \underline{\mu}]^T \Sigma^{-1} [\underline{x} - \underline{\mu}]\right)$$

Proof

$$\begin{aligned} p(x_1, x_2) &= p(x_2|x_1)p(x_1) \\ &= p(x_2) \cdot p(x_1) \\ &= (2\pi\sigma^2)^{-1/2} e^{-\frac{1}{2\sigma^2}(x_1 - \mu_1)^2} \cdot (2\pi\sigma^2)^{-1/2} e^{-\frac{1}{2\sigma^2}(x_2 - \mu_2)^2} \\ &= |2\pi\Sigma|^{-1/2} e^{-\frac{1}{2}[\underline{x} - \underline{\mu}]^T \begin{pmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{pmatrix}^{-1} [\underline{x} - \underline{\mu}]} \end{aligned}$$

$$\text{with } \Sigma = \begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix}, |\Sigma| = (\sigma^2)^2, a\Sigma \rightarrow a^2|\Sigma|$$

# Maximum Likelihood Estimation (MLE)

# Maximum Likelihood Estimation (MLE)

## Example

Suppose we have  $n = 3$  data points :

- $y^{(1)} = 1,$
- $y^{(2)} = 0.5,$
- $y^{(3)} = 1.5$

which are independent and Gaussian with unknown mean  $\theta$  and variance 1.

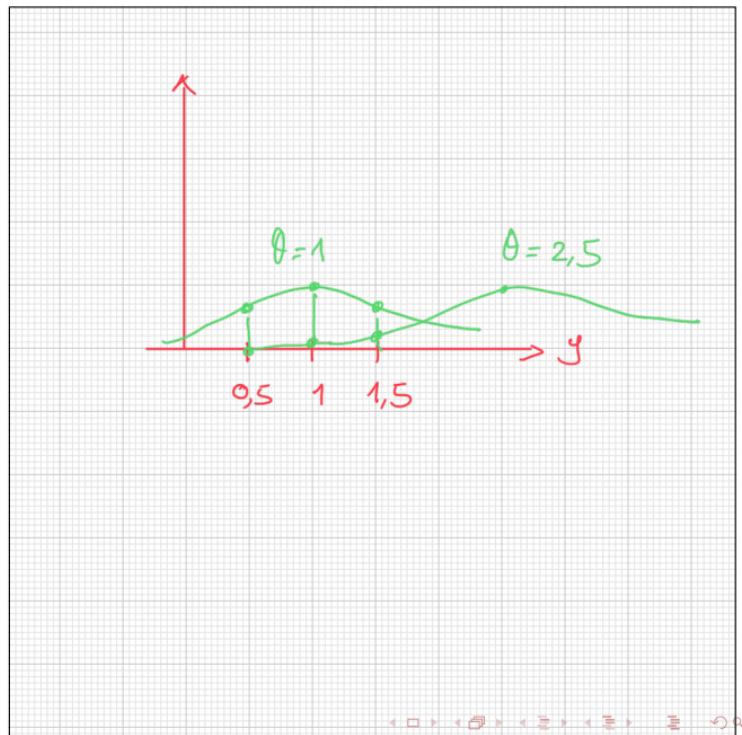
We can express the likelihood as :

- $y^{(i)} \sim \mathcal{N}(\theta, 1) = \theta + \mathcal{N}(0, 1)$
- $p(y^{(1)}, y^{(2)}, y^{(3)} | \theta) = p(y^{(1)} | \theta)p(y^{(2)} | \theta)p(y^{(3)} | \theta)$

Which guess of  $\theta$  is more likely ?

- $\theta = 1$
- $\theta = 2.5$

We take the  $\theta$  the maximizes the likelihood



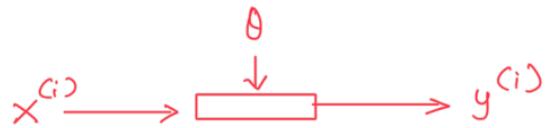
# Maximum Likelihood Estimation (MLE)

## Likelihood for linear Regression

Suppose  $y^{(i)}$  is Gaussian distributed with mean  $\underline{x}^{(i)} \underline{\theta}$  and variance  $\sigma^2$

$$y^{(i)} \sim \mathcal{N}(\underline{x}^{(i)} \underline{\theta}, \sigma^2)$$

$$= \underline{x}^{(i)} \underline{\theta} + \mathcal{N}(0, \sigma^2)$$



The likelihood is

$$\begin{aligned} p(\underline{y} | \underline{X}, \underline{\theta}, \sigma) &= \prod_{i=1}^m p(y^{(i)} | \underline{x}^{(i)}, \underline{\theta}, \sigma) \\ &= \prod_{i=1}^m (2\pi\sigma)^{-1/2} e^{-\frac{1}{2\sigma^2} (y^{(i)} - \underline{x}^{(i)} \underline{\theta})^2} \\ &= (2\pi\sigma)^{-m/2} e^{-\frac{1}{2\sigma^2} \sum_{i=1}^m (y^{(i)} - \underline{x}^{(i)} \underline{\theta})^2} \\ &= (2\pi\sigma)^{-m/2} e^{-\frac{1}{2\sigma^2} (\underline{y} - \underline{X} \underline{\theta})^T (\underline{y} - \underline{X} \underline{\theta})} \end{aligned}$$

This is equivalent to  $p(\text{data} | \text{parameters}) = \frac{1}{Z} e^{-\text{Loss}(\text{data}, \text{parameters})}$

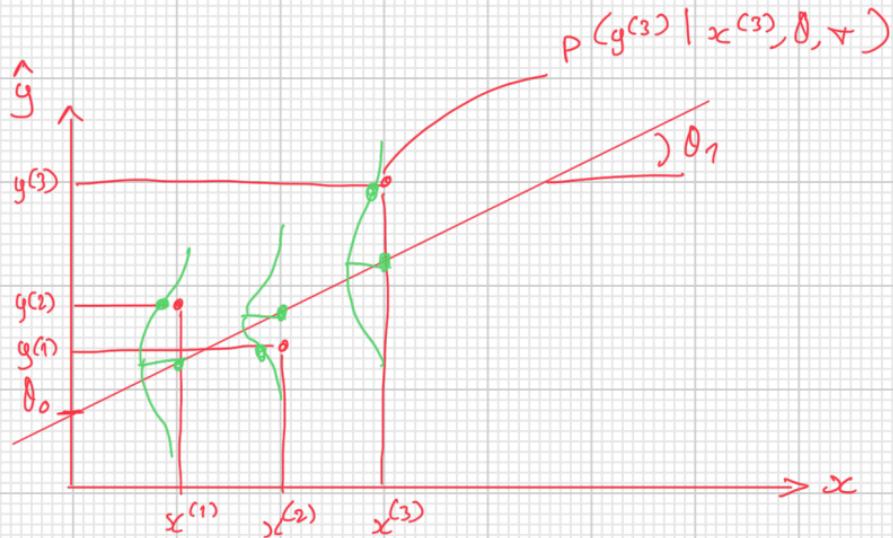
- with  $\text{Loss}(\text{data}, \text{parameters}) = J(\underline{\theta}) = \sum_{i=1}^m (y^{(i)} - \underline{x}^{(i)} \underline{\theta})^2$
- is the MSE (Mean Square Error) of Linear Regression

Maximizing the likelihood is equivalent to minimizing the Loss

# Maximum Likelihood Estimation (MLE)

Example (back to linear regression) :

$$p(y^{(i)} | \underline{x}^{(i)}, \underline{\theta}, \sigma) = (2\pi\sigma)^{-1/2} e^{-\frac{1}{2\sigma^2}(y^{(i)} - \theta_0 - x^{(i)}\theta_1)^2}$$



# Maximum Likelihood Estimation (MLE)

## Likelihood for linear Regression

Maximizing the **Likelihood**  $\Rightarrow$  minimizing the **Negative Log-Likelihood (NLL)**

$$p(\underline{y} | \underline{X}, \underline{\theta}, \sigma) = (2\pi\sigma^2)^{-m/2} e^{-\frac{1}{2\sigma^2}(\underline{y} - \underline{X}\underline{\theta})^T(\underline{y} - \underline{X}\underline{\theta})}$$

$$-\log p(\underline{y} | \underline{X}, \underline{\theta}, \sigma) = \frac{m}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2}(\underline{y} - \underline{X}\underline{\theta})^T(\underline{y} - \underline{X}\underline{\theta})$$

**MLE of  $\underline{\theta}$**  (same solution as MSE) :

$$\underline{\theta}_{ML} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{y}$$

Proof :

$$0 = \frac{\partial}{\partial \underline{\theta}} \left( \frac{1}{2\sigma^2}(\underline{y} - \underline{X}\underline{\theta})^T(\underline{y} - \underline{X}\underline{\theta}) \right)$$

**MLE of  $\sigma$**  :

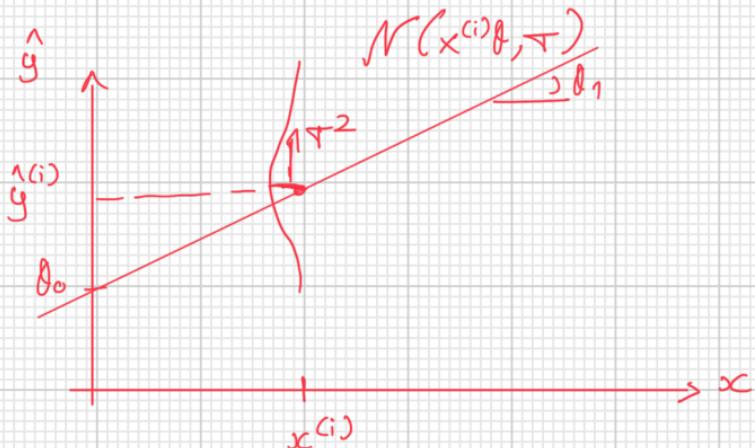
$$\begin{aligned} \sigma_{ML}^2 &= \frac{1}{m}(\underline{y} - \underline{X}\underline{\theta})^T(\underline{y} - \underline{X}\underline{\theta}) \\ &= \boxed{\frac{1}{m} \sum_{i=1}^n (y^{(i)} - \underline{x}^{(i)} \underline{\theta})^2} \end{aligned}$$

Proof

$$\begin{aligned} 0 &= \frac{\partial}{\partial \sigma} \left( \frac{m}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2}(\underline{y} - \underline{X}\underline{\theta})^T(\underline{y} - \underline{X}\underline{\theta}) \right) \\ &= m \frac{1}{\sigma} - \frac{1}{\sigma^3}(\underline{y} - \underline{X}\underline{\theta})^T(\underline{y} - \underline{X}\underline{\theta}) \end{aligned}$$

# Maximum Likelihood Estimation (MLE)

Making predictions (provides a confidence in prediction !)



## Bernoulli distribution

Bernoulli distribution :  $X$  takes value in  $\{0, 1\}$

- model of a coin (probability of head or tail)

$$P(X = x|\theta) = \begin{cases} p & \text{if } x = 1 \\ 1 - p & \text{if } x = 0 \end{cases}$$
$$= \boxed{p^x(1 - p)^{(1-x)}}$$

- with  $p \in [0, 1]$

# Maximum Likelihood Estimation (MLE)

## Bernoulli distribution

**Entropy** : is a measure of the uncertainty associated with a random variable :

$$H(X) = - \sum_x p(x|\theta) \log(p(x|\theta))$$

### Entropy for a Bernoulli distribution

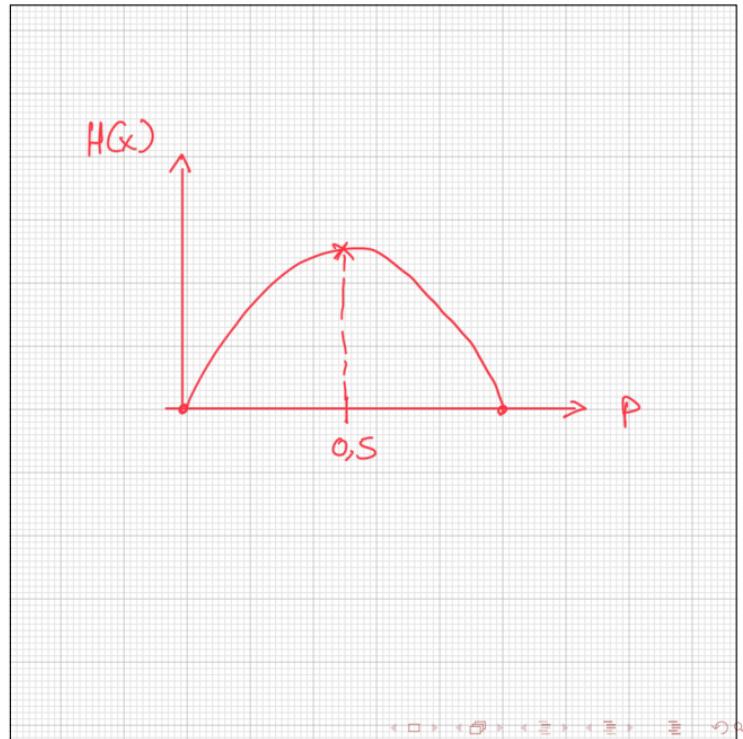
$$\begin{aligned} H(X) &= - \sum_{x=0}^1 p^x (1-p)^{(1-x)} \log(p^x (1-p)^{(1-x)}) \\ &= - \underbrace{[(1-p) \log(1-p) + p \log(p)]}_{\text{if } x=0} \end{aligned}$$

if  $x=1$

### Entropy of a Gaussian distribution at $D$ dimensions

$$H(\mathcal{N}(\mu, \Sigma)) = \frac{1}{2} \log((2\pi e)^D |\Sigma|)$$

- proportionnal to the determinant of the covariance matrix



## Regression : Ridge, Lasso, Non-linear (basis, kernel)

## Ridge regularization

### Least-square regression

$$J(\theta) = (\underline{y} - \underline{\underline{X}}\theta)^T(\underline{y} - \underline{\underline{X}}\theta)$$

$$\hat{\theta} = (\underline{\underline{X}}^T \underline{\underline{X}})^{-1} \underline{\underline{X}}^T \underline{y}$$

- inversion of  $(\underline{\underline{X}}^T \underline{\underline{X}})$  can lead to problems (if system poorly conditionned)
- solution : add a small constant  $\delta$  to the diagonal elements

### Ridge regression

$$\hat{\theta} = (\underline{\underline{X}}^T \underline{\underline{X}} + \delta^2 \mathbb{I})^{-1} \underline{\underline{X}}^T \underline{y} \quad \text{with } \mathbb{I} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

- Corresponds to the **regularized quadratic cost function**

Proof :

$$J(\theta) = (\underline{y} - \underline{\underline{X}}\underline{\theta})^T(\underline{y} - \underline{\underline{X}}\underline{\theta}) + \delta^2 \underline{\theta}^T \underline{\theta}$$

- equivalent to

$$\min_{\underline{\theta}: \underline{\theta}^T \underline{\theta} \leq t(\delta)} (\underline{y} - \underline{\underline{X}}\underline{\theta})^T(\underline{y} - \underline{\underline{X}}\underline{\theta})$$

- where  $t$  is an arbitrary function
- also named weight decay

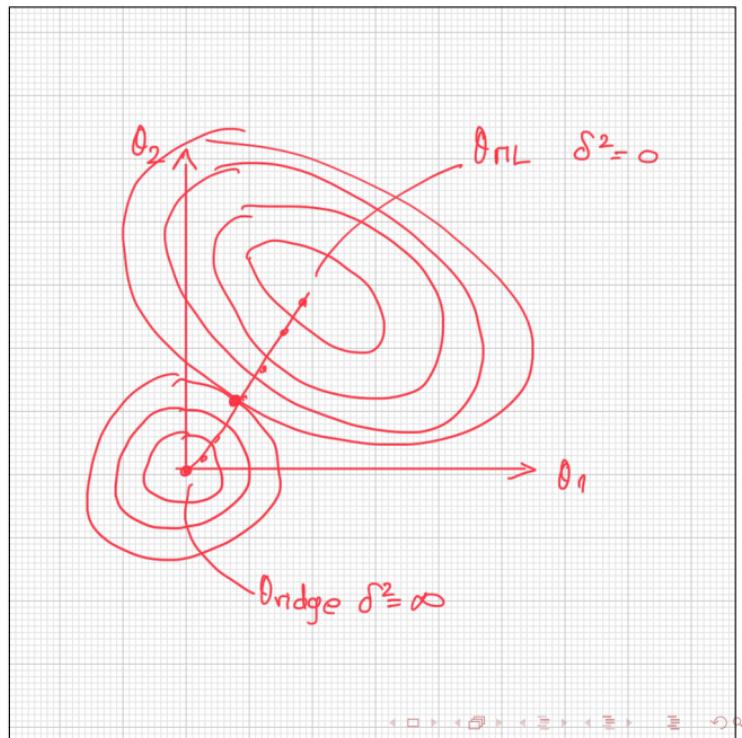
$$\begin{aligned} 0 &= \frac{\partial}{\partial \underline{\theta}} (\underline{\theta}^T \underline{\underline{X}}^T \underline{\underline{X}} \underline{\theta} - 2\underline{y}^T \underline{\underline{X}} \underline{\theta} + \underline{y}^T \underline{y} + \delta^2 \underline{\theta}^T \underline{\theta}) \\ &= 2\underline{\underline{X}}^T \underline{\underline{X}} \underline{\theta} - 2\underline{\underline{X}}^T \underline{y} + 2\delta^2 \mathbb{I} \underline{\theta} \\ &= 2(\underline{\underline{X}}^T \underline{\underline{X}} + \delta^2 \mathbb{I}) \underline{\theta} - 2\underline{\underline{X}}^T \underline{y} \\ \underline{\theta}_{ridge} &= (\underline{\underline{X}}^T \underline{\underline{X}} + \delta^2 \mathbb{I})^{-1} \underline{\underline{X}}^T \underline{y} \end{aligned}$$

## Ridge regularization

Example :

$$\begin{aligned}\underline{\theta}^T &= [\theta_1, \theta_2] \\ \rightarrow \underline{\theta}^T \underline{\theta} &= [\theta_1, \theta_2] \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \\ &= \theta_1^2 + \theta_2^2 = \text{const}\end{aligned}$$

- equation of a circle
- automatically remove irrelevant axes



## LASSO regularization

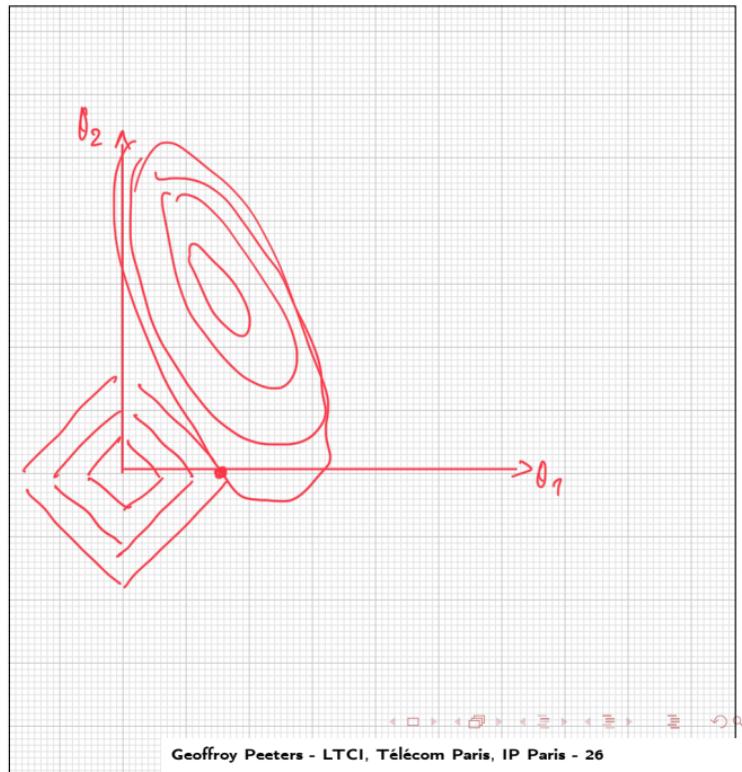
Ridge :

$$\|\underline{\theta}\|_2^2 = \underline{\theta}^T \underline{\theta}$$

Lasso :

$$\|\underline{\theta}\|_1 = \sum_d |\theta_{\text{d}}| = \theta_1 + \theta_2$$

- intersections happen at corners ( $\theta_2 = 0$ )
- sparse, compress-sensing



## Non-linear regression (regression with basis functions)

### Basis functions $\phi(\cdot)$

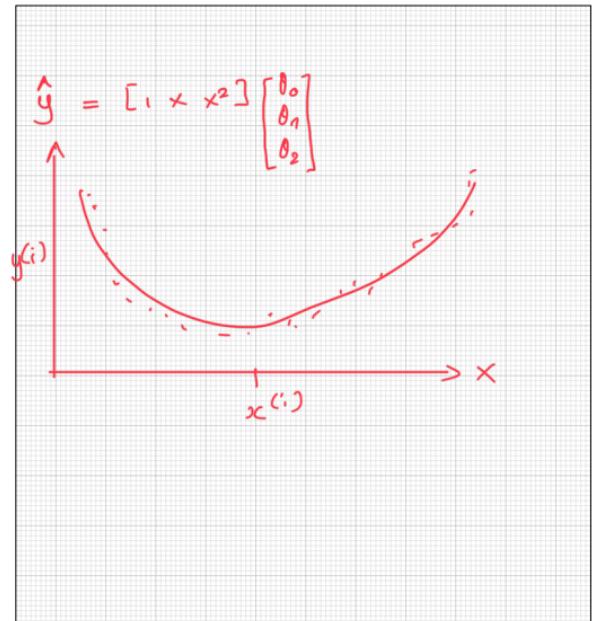
$$\hat{y}(\underline{x}) = \underline{\phi}(\underline{x}) \underline{\theta} + \epsilon$$

#### Examples

- in 1 dimensions :  $\underline{\phi}(x) = [1, x, x^2]$

$$\begin{aligned}\hat{y} &= \underline{\phi}(x) \underline{\theta} \\ &= \theta_0 + x\theta_1 + x^2\theta_2 \\ &= [1 \ x \ x^2] \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \end{bmatrix}\end{aligned}$$

- in 2 dimensions :  $\underline{\phi}(\underline{x}) = [1, x_1, x_2, x_1^2, x_2^2, x_1 x_2]$



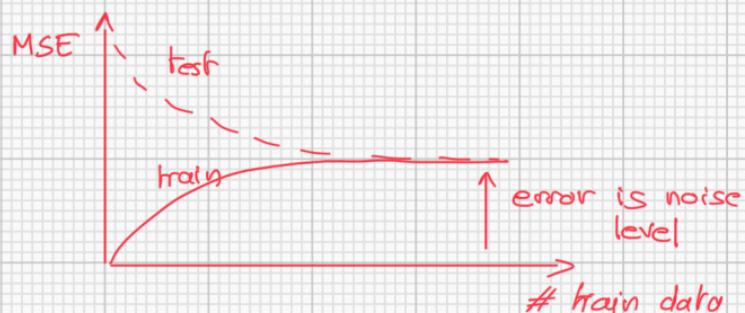
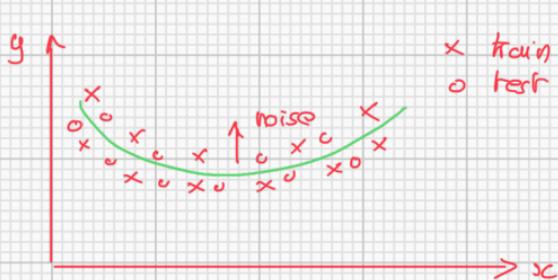
# Regression : Ridge, Lasso, Non-linear (basis, kernel)

Having the right model complexity / having large number of training data

Good model (bias equal noise level)

True model  $y^{(i)} = \theta_0 + \theta_1 x^{(i)} + \theta_2 x^{(i)2} + \underbrace{\mathcal{N}(0, \sigma^2)}_{\text{noise}}$   $\rightarrow$  model is OK

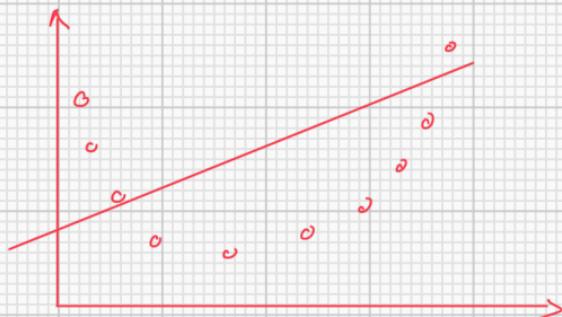
used model  $\hat{y} = \theta_0 + \theta_1 x + \theta_2 x^2$



# Regression : Ridge, Lasso, Non-linear (basis, kernel)

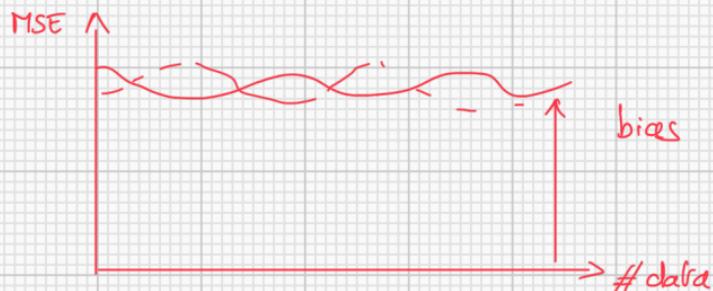
Having the right model complexity / having large number of training data

Too simplest model (bad results whatever number or data)



$$\hat{y}^{(i)} = \theta_0 + \theta_1 x^{(i)}$$

↳ model is too simple



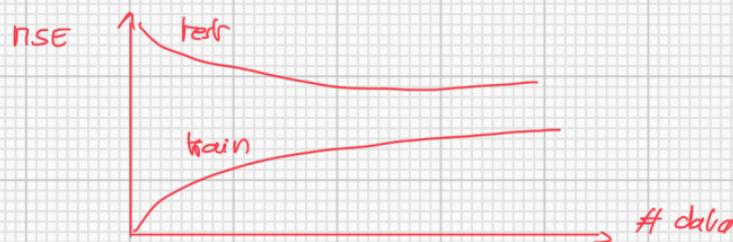
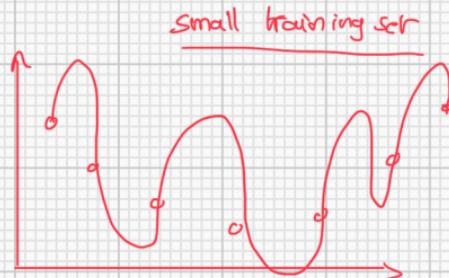
# Regression : Ridge, Lasso, Non-linear (basis, kernel)

## Having the right model complexity / having large number of training data

Very complex model (results on test-set depend on the number of data)

$$f(x) = \theta_0 + \theta_1 x^{(1)} + \dots + \theta_{25} x^{(25)}$$

↳ model is too complex



more data improves but  
only with the right  
complexity

# Regression : Ridge, Lasso, Non-linear (basis, kernel)

## Controlling model complexity with regularization $\delta$

$$J(\theta) = \text{MSE} + \delta^2 \frac{\theta^\top \theta}{\theta^\top \theta}$$



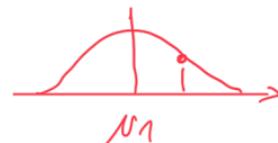
## Non-linear regression (regression with basis functions)

- **RBF kernel** :  $\phi(\underline{x}) = [k(\underline{x}, \underline{\mu}_1, \lambda), \dots, k(\underline{x}, \underline{\mu}_d, \lambda)]$

- with  $k(\underline{x}, \underline{\mu}_d, \lambda) = e^{-\frac{1}{\lambda} \|\underline{x} - \underline{\mu}_d\|^2}$

$$\begin{aligned}\hat{y}(\underline{x}^{(i)}) &= \phi(\underline{x}^{(i)}) \underline{\theta} \\ &= 1\theta_0 + k(\underline{x}^{(i)}, \underline{\mu}_1, \lambda)\theta_1 + \dots + k(\underline{x}^{(i)}, \underline{\mu}_d, \lambda)\theta_d\end{aligned}$$

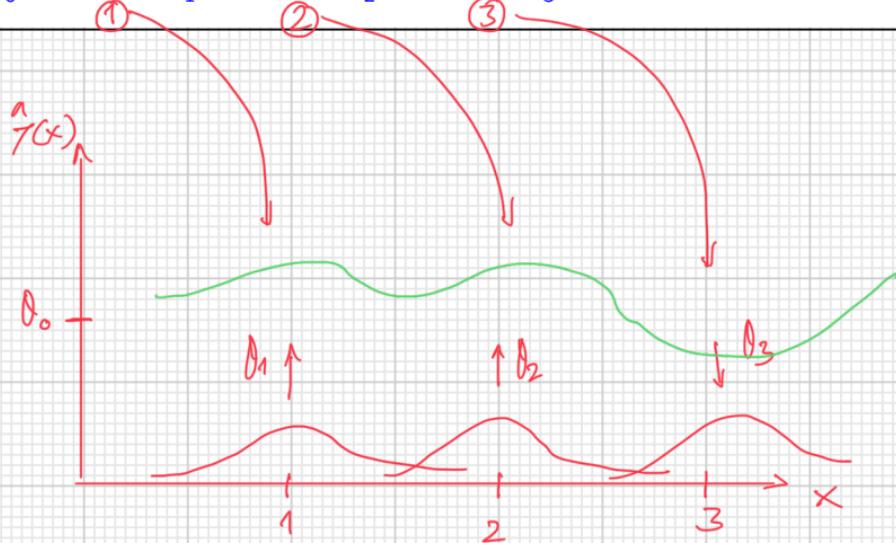
$$\begin{aligned}\Phi(\underline{x}) &= [1, \underline{x}, \underline{x}^2] \\ &= [\theta_0, \theta_1, \theta_2 \\ &\quad e^{-\frac{\|\underline{x}-\underline{\mu}_1\|^2}{\lambda}}, e^{-\frac{\|\underline{x}-\underline{\mu}_2\|^2}{\lambda}}, e^{-\frac{\|\underline{x}-\underline{\mu}_3\|^2}{\lambda}}]\end{aligned}$$



# Regression : Ridge, Lasso, Non-linear (basis, kernel)

## Kernel regression with RBFs

- Example :  $\hat{y}(x) = \theta_0 + e^{-||x-1||^2} \theta_1 + e^{-||x-2||^2} \theta_2 + e^{-||x-3||^2} \theta_3$



## Kernel regression with RBFs

We note  $\underline{\phi}(x^{(i)})$  the 4-dimensional (3 bases) :

- For one data :

$$\underline{\phi}(x^{(i)}) = [1, k(x^{(i)}, \underline{\mu}_1, \lambda), k(x^{(i)}, \underline{\mu}_2, \lambda), k(x^{(i)}, \underline{\mu}_3, \lambda)]$$

- For  $m$  data

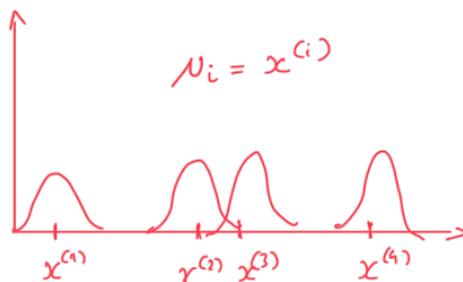
$$\begin{aligned}\hat{y} &= \underline{\phi} \underline{\theta} \\ \hat{\underline{\theta}}_{LS} &= (\underline{\phi}^T \underline{\phi})^{-1} \underline{\phi}^T \underline{y} \\ \hat{\underline{\theta}}_{Ridge} &= \left[ (\underline{\phi}^T \underline{\phi} + \delta^2 \mathbb{I})^{-1} \underline{\phi}^T \underline{y} \right]\end{aligned}$$

$$\underline{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(m)} \end{bmatrix}, \underline{\phi} = \begin{bmatrix} \underline{\phi}(x^{(1)}) \\ \vdots \\ \underline{\phi}(x^{(m)}) \end{bmatrix}, \underline{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

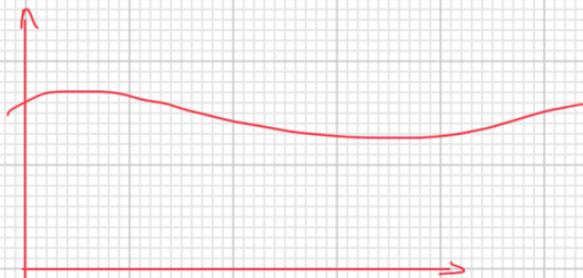
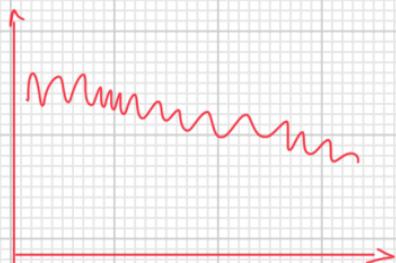
- $\Rightarrow$  It is still a regression !

### Specific case :

- choose the  $\underline{\mu}_i$  equal to the training data point  $x^{(i)}$
- $\Rightarrow$  one kernel on each training data point



## Choice of the kernel size $\lambda$





## Gradient

$$\frac{\partial f(\theta_1, \theta_2)}{\partial \theta_1} = \lim_{\Delta \theta_1 \rightarrow 0} \frac{f(\theta_1 + \Delta \theta_1, \theta_2) - f(\theta_1, \theta_2)}{\Delta \theta_1}$$

- Consider the function  $f(\theta_1, \theta_2) = \theta_1^2 + \theta_2^2$

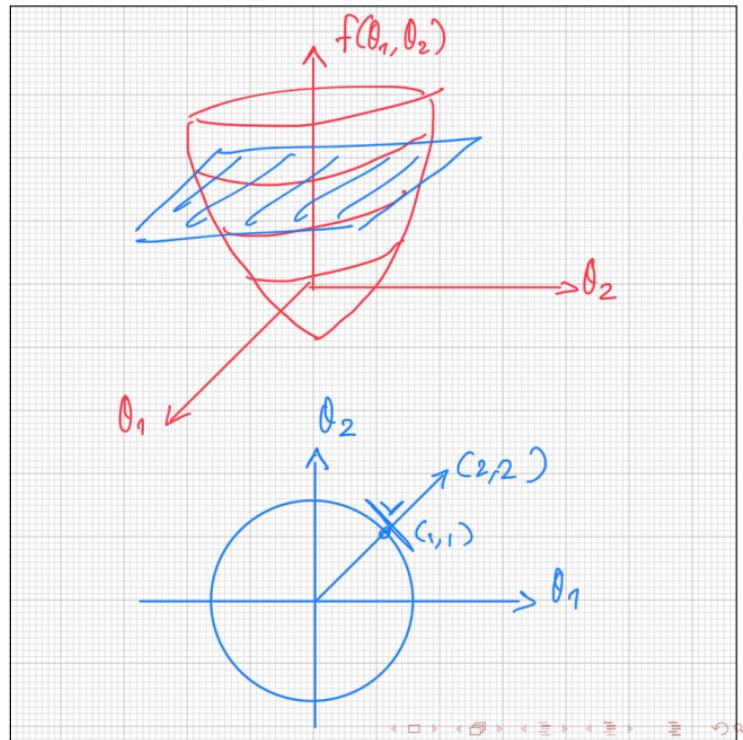
$$\frac{\partial f}{\partial \theta_1} = 2\theta_1$$

$$\frac{\partial f}{\partial \theta_2} = 2\theta_2$$

$$\nabla f(\theta) = \begin{bmatrix} 2\theta_1 \\ 2\theta_2 \end{bmatrix}$$

Gradient :

- always orthogonal to the tangent of the level curve
- indicates the direction of
  - the maximum change,
  - the steepest ascent



## Hessian

$$\frac{\partial}{\partial \theta_1} \left( \frac{\partial f(\theta)}{\partial \theta_1} \right) = \frac{\partial^2 f}{\partial \theta_1^2} = 2$$

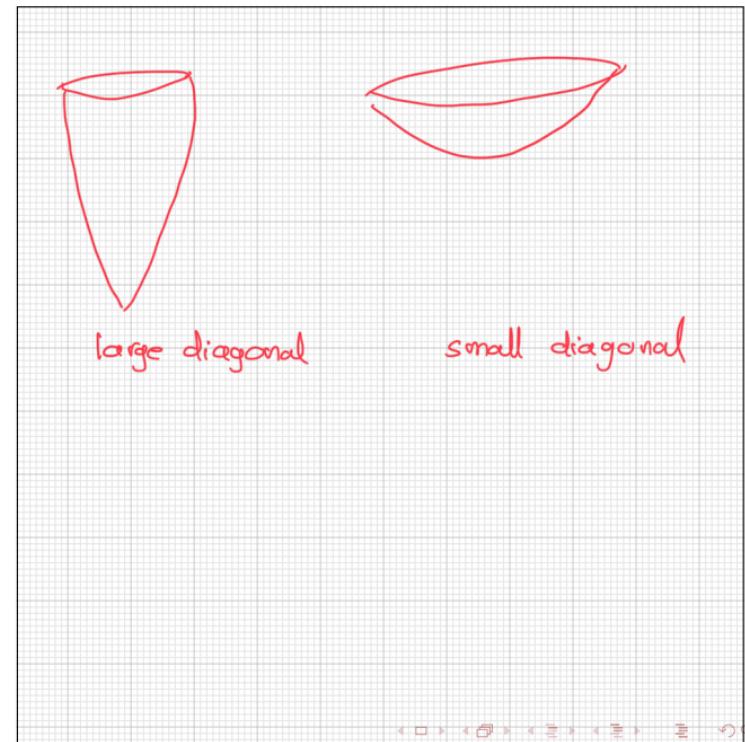
$$\frac{\partial^2 f}{\partial \theta_2^2} = 2$$

$$\frac{\partial^2 f}{\partial \theta_1 \partial \theta_2} = \frac{\partial^2 f}{\partial \theta_2 \partial \theta_1} = 0$$

$$H = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$

Hessian :

- represents the curvature,
- allows to check convexity (saddle point)
- if a minimum exists, all positive eigen values



# Optimisation

**Gradient** :  $\mathbb{R}^d \rightarrow \mathbb{R}$

$$\nabla_{\theta} f(\theta) = \begin{bmatrix} \frac{\partial f(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial f(\theta)}{\partial \theta_d} \end{bmatrix}$$

**Jacobian** :  $\mathbb{R}^d \rightarrow \mathbb{R}^{d'}$

$$J(\theta) = \begin{bmatrix} \frac{\partial^2 f_1(\theta)}{\partial \theta_1} & \dots & \frac{\partial f_1(\theta)}{\partial \theta_d} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f_{d'}(\theta)}{\partial \theta_1} & \dots & \frac{\partial f_{d'}(\theta)}{\partial \theta_d} \end{bmatrix}$$

**Hessian**  $\mathbb{R}^d \rightarrow \mathbb{R}$

$$\nabla_{\theta}^2 f(\theta) = \begin{bmatrix} \frac{\partial^2 f(\theta)}{\partial \theta_1^2} & \dots & \frac{\partial f(\theta)}{\partial \theta_1 \theta_d} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f(\theta)}{\partial \theta_d \theta_1} & \dots & \frac{\partial f(\theta)}{\partial \theta_d^2} \end{bmatrix}$$

**How to estimate them ?**

$$f(\theta) = f(\theta, x^{(1:m)}) = \frac{1}{m} \sum_{i=1}^m f(\theta, x^{(i)})$$

$$\underline{g}(\theta) = \nabla_{\theta} f(\theta) = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta} f(\theta, x^{(i)})$$

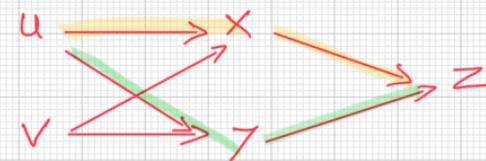
$$\underline{\underline{H}}(\theta) = \nabla_{\theta}^2 f(\theta) = \frac{1}{m} \sum_{i=1}^m \nabla_{\theta}^2 f(\theta, x^{(i)})$$

## Calculus background : Chain rule

$$z = f(x(u, v), y(u, v))$$

$$\frac{\partial z}{\partial u} = \frac{\partial z}{\partial x} \frac{\partial x}{\partial u} + \frac{\partial z}{\partial y} \frac{\partial y}{\partial u}$$

Yellow bar:  $\frac{\partial z}{\partial x}$   
Green bar:  $\frac{\partial z}{\partial y}$



## Gradient and Hessian for linear regression

- For linear regression :  $y^{(i)} = \theta_0 + x^{(i)}\theta_1$

### Scalar form

$$f(\theta) = \sum_i \left( y^{(i)} - \theta_0 - x^{(i)}\theta_1 \right)^2 + \delta^2\theta_1^2$$

Applying the chain rule :

$$\begin{aligned} \nabla_{\theta} f &= \begin{bmatrix} \frac{\partial}{\partial \theta_0} \\ \frac{\partial}{\partial \theta_1} \end{bmatrix} \\ &= \begin{bmatrix} \sum_i 2(y^{(i)} - \theta_0 - x^{(i)}\theta_1)(-1) \\ \sum_i 2(y^{(i)} - \theta_0 - x^{(i)}\theta_1)(-x^{(i)}) + 2\delta^2\theta_1 \end{bmatrix} \\ &= \begin{bmatrix} 2 \sum_i 1 \cdot ((\theta_0 + x^{(i)}\theta_1) - y^{(i)}) (-1) \\ 2 \sum_i x^{(i)} \cdot ((\theta_0 + x^{(i)}\theta_1) - y^{(i)}) + 2\delta^2\theta_1 \end{bmatrix} \end{aligned}$$

### Matrix form

$$\begin{aligned} f(\theta) &= (\underline{y} - \underline{X}\underline{\theta})^T (\underline{y} - \underline{X}\underline{\theta}) \\ \nabla_{\theta} f(\theta) &= \frac{\partial}{\partial \theta} (\underline{y}^T \underline{y} - 2\underline{y}^T \underline{X}\underline{\theta} + \underline{\theta}^T \underline{X}^T \underline{X}\underline{\theta}) \\ &= -2\underline{X}^T \underline{y} + 2\underline{X}^T \underline{X}\underline{\theta} \\ &= \boxed{2\underline{X}^T (\underline{X}\underline{\theta} - \underline{y})} \\ \nabla_{\theta}^2 f(\theta) &= 0 + 2\underline{X}^T \underline{X} \\ &= \boxed{2 \sum_{(d,m)(m,d)} \underline{X}} \end{aligned}$$

## Optimisation : 1) Steepest/Stochastic Gradient Descent (SGD) algorithm

Notation :  $\underline{\theta}_k$  value of  $\underline{\theta}$  at iteration  $k$

### Steepest/Stochastic Gradient Descent (SGD) algorithm

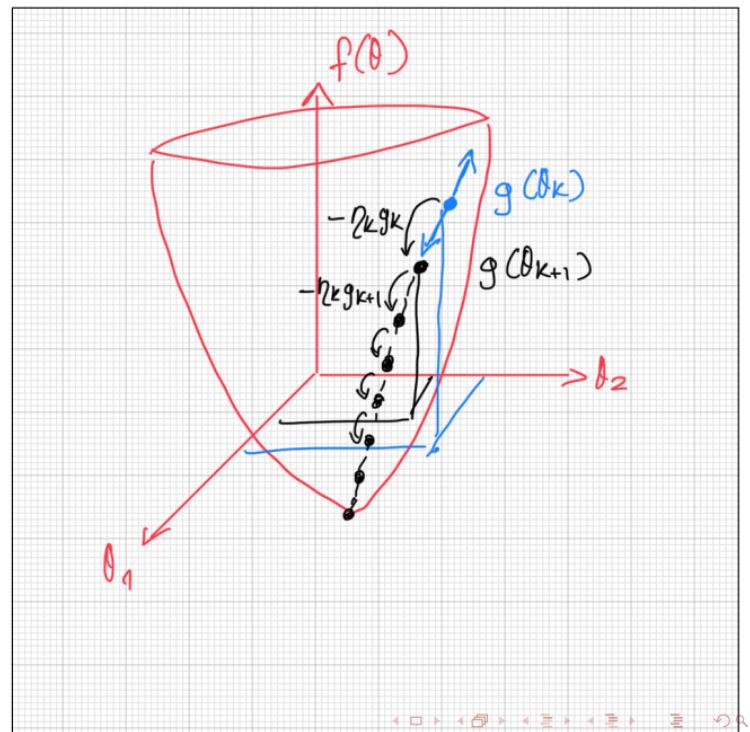
$$\underline{\theta}_{k+1} = \underline{\theta}_k - \eta_k \underline{g}_k$$

$$= \underline{\theta}_k - \eta_k \nabla_{\underline{\theta}} f(\underline{\theta}_k)$$

- where  $\eta_k$  is the learning rate, step size

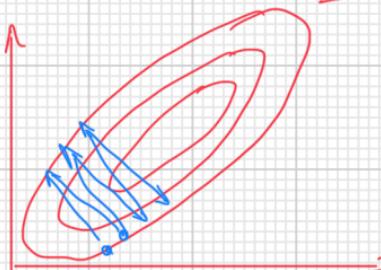
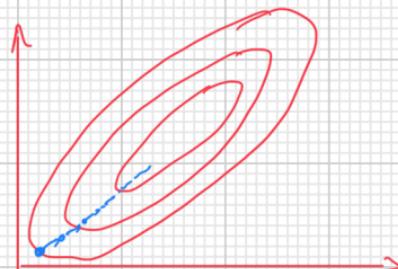
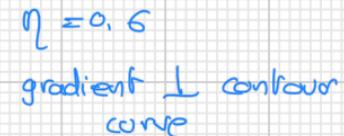
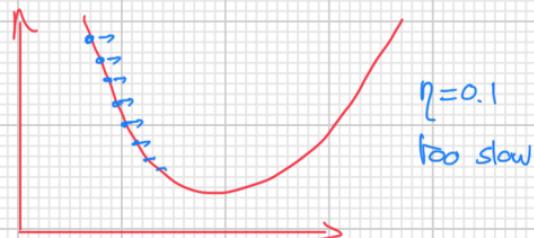
### SGD for Linear Regression

$$\underline{\theta}_{k+1} = \underline{\theta}_k - \eta_k \left[ 2\underline{X}^T (\underline{X} \underline{\theta}_k - \underline{y}) \right]$$



## Optimisation : 1) Steepest/Stochastic Gradient Descent (SGD) algorithm

- How to choose  $\eta_k$ ?
  - $\eta = 0.1$  too slow
  - $\eta = 0.6$  too large



## Optimisation : 2) Newton/Hessian algorithm

If we knew the curvature, we could adapt  $\eta_k$ !

**Taylor serie expansion** (approximate  $f(\underline{\theta})$  around  $\underline{\theta}_k$  with a quadratic ball) :

$$f_{quad}(\underline{\theta}) = f(\underline{\theta}_k) + \underline{g}_k^T (\underline{\theta} - \underline{\theta}_k) + \frac{1}{2} (\underline{\theta} - \underline{\theta}_k)^T \underline{H}_k (\underline{\theta} - \underline{\theta}_k)$$

find the minimum of the quadratic ball :

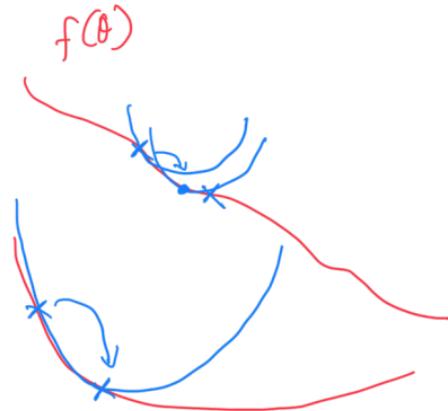
$$0 = \nabla_{\underline{\theta}} f_T(\underline{\theta})$$

$$= 0 + \underline{g}_k + \underline{H}_k (\underline{\theta} - \underline{\theta}_k)$$

$$-\underline{g}_k = \underline{H}_k (\underline{\theta} - \underline{\theta}_k)$$

$$\underline{\theta}_{k+1} = \boxed{\underline{\theta}_k - \underline{H}_k^{-1} \underline{g}_k}$$

- $\Rightarrow$  this is Newton's method



## Optimisation : 2) Newton/Hessian algorithm

### Newton's for Linear Regression :

$$\begin{aligned}\theta_{k+1} &= \theta_k - \underline{\underline{H}}_k^{-1} \underline{\underline{g}}_k \\&= \theta_k - (2\underline{\underline{X}}^T \underline{\underline{X}})^{-1} [2\underline{\underline{X}}^T (\underline{\underline{X}}\theta_k - \underline{\underline{y}})] \\&= \theta_k - (\underline{\underline{X}}^T \underline{\underline{X}})^{-1} \underline{\underline{X}}^T \underline{\underline{X}}\theta_k + (\underline{\underline{X}}^T \underline{\underline{X}})^{-1} \underline{\underline{X}}^T \underline{\underline{y}} \\&= \boxed{(\underline{\underline{X}}^T \underline{\underline{X}})^{-1} \underline{\underline{X}}^T \underline{\underline{y}}}\end{aligned}$$

using (as seen before)

$$\begin{aligned}f(\theta) &= (\underline{\underline{y}} - \underline{\underline{X}}\theta)^T (\underline{\underline{y}} - \underline{\underline{X}}\theta) \\g &= \nabla_\theta f(\theta) = 2\underline{\underline{X}}^T (\underline{\underline{X}}\theta - \underline{\underline{y}}) \\H &= \nabla_\theta^2 f(\theta) = 2\underline{\underline{X}}^T \underline{\underline{X}}\end{aligned}$$

- This is the theoretical solution (we get it in just one step) !
- In practice, it is very costly
  - $\Rightarrow$  need to store  $H$  (if we have 1.000.000 neurons this is a 1.000.000x1.000.000 matrix !)
  - BFGS (Broyden-Fletcher-Goldfarb-Shanno) algorithm, L-BFGS (Limited memory) algorithm

## Why "Stochastic" in SGD?

### Stochastic ?

- because we only have seen a fraction of the data, uncertainty

$$\nabla_{\underline{\theta}} f(\underline{\theta}) = \underbrace{\int \nabla_{\underline{\theta}} f(x, \underline{\theta}) p(x) dx}_{\mathbb{E}[\nabla_{\underline{\theta}} f(x, \underline{\theta})]} \stackrel{\text{we want}}{=} 0$$

SGD  
 ↑  
 Steepest

Can be re-written

$$\mathbb{E}[\nabla_{\underline{\theta}} f(x, \underline{\theta})] = 0$$

$$\eta \mathbb{E}[\nabla_{\underline{\theta}} f(x, \underline{\theta})] = 0 \cdot \eta = 0$$

$$\underline{\theta} + \eta \mathbb{E}[\nabla_{\underline{\theta}} f(x, \underline{\theta})] = \underline{\theta}$$

$$\underline{\theta}_{k+1} = \underline{\theta}_k - \eta \mathbb{E}[\nabla_{\underline{\theta}} f(x, \underline{\theta})]$$

$$\simeq \underline{\theta}_k - \eta \frac{1}{m} \sum_{i=1}^m \nabla_{\underline{\theta}} f(x^{(i)}, \underline{\theta}) \text{ with } x^{(i)} \sim p(x)$$

$$\simeq \underline{\theta}_k - \eta \nabla_{\underline{\theta}} f(x^{(k)}, \underline{\theta}_k)$$

↙ m=1

Therefore

$$\underline{\theta}_{k+1} = \underbrace{\underline{\theta}_k - \eta \mathbb{E}[\nabla_{\underline{\theta}} f(x, \underline{\theta})]}_{\text{what we want}} + \underbrace{\eta \left[ \mathbb{E}[\nabla_{\underline{\theta}} f] - \nabla_{\underline{\theta}} f(x^{(k)}, \underline{\theta}_k) \right]}_{\text{noise, change over time, stochastic, bounded}}$$

## Variations

Re-writting SGD as a sum :

$$\begin{aligned}\underline{\theta}_{k+1} &= \underline{\theta}_k - \eta_k \left[ 2\underline{X}^T (\underline{X} \underline{\theta}_k - \underline{y}) \right] \\ &= \underline{\theta}_k + \eta_k \underline{X}^T (\underline{y} - \underline{X} \underline{\theta}_k) \\ &= \underline{\theta}_k + \eta_k \sum_i \underline{x}^{(i)}^T (\underline{y}^{(i)} - \underline{x}^{(i)} \underline{\theta}_k)\end{aligned}$$

## Variations of SGD

- **Batch** (all  $m$  data points) :

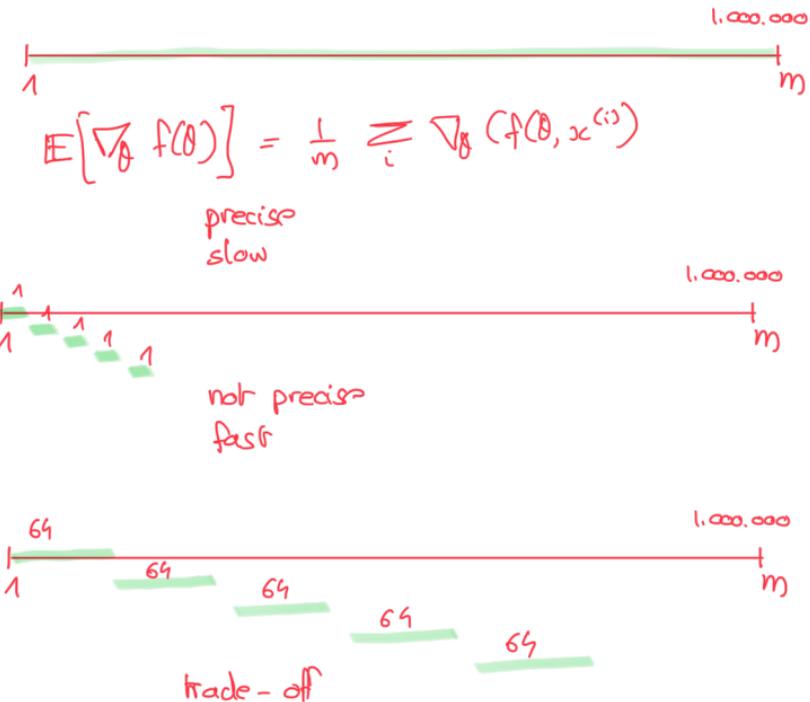
$$\theta_{k+1} = \theta_k + \eta_k \sum_{i=1}^m \underline{x}^{(i)} {}^T (y^{(i)} - \underline{x}^{(i)} \underline{\theta}_k)$$

- **Online** (a single data point) :

$$\theta_{k+1} = \theta_k + \eta_k (x^{(k)})^T (y^{(k)} - x^{(k)} \theta_k)$$

- **Mini-Batch** (a subset of data points) :

$$\theta_{k+1} = \theta_k + \eta_k \sum_{j=1}^{20} (\underline{x}^{(j)})^T (y^{(j)} - \underline{x}^{(j)} \underline{\theta}_k)$$



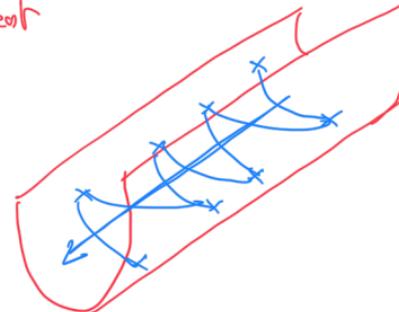
## Variations

### Momentum :

- Goal : accelerate descent in cumulated direction

$$\theta_{k+1} = \theta_k + \alpha(\theta_k - \theta_{k-1}) + (1 - \alpha)[- \eta \nabla J(\theta)]$$

$$\underbrace{\theta_{k+1} - \theta_k}_{\Delta \theta_{k+1}} = \alpha \underbrace{(\theta_k - \theta_{k-1})}_{\Delta \theta_k} + (1 - \alpha)[- \eta \nabla J(\theta)]$$



### Adagrad :

- Goal : put more weights on rare features ;

- For feature  $d$  at step  $k$

$$w_d^{k+1} = w_d^k - \frac{\eta}{\sum_{\tau=1}^k (g_d^\tau)^2} g_d^k \quad \text{d dimension}$$

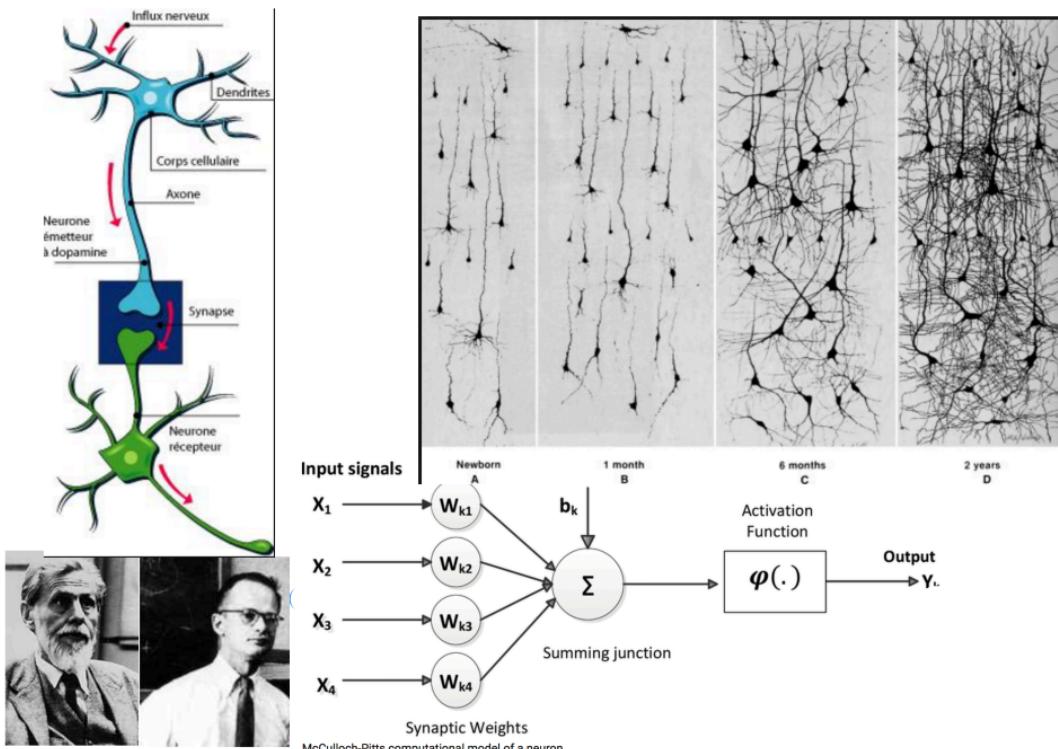
- $\sum (g_d^\tau)^2$  is the history of how much a feature has changed in the past
- if small (has been seen few times in the past) then put more weights

Momentum + Adagrad = ADAM



# Logistic Regression

## McCulloch - Pitts model of a neuron



# Logistic Regression

## Activation function

The logistic regression is a **binary classification** algorithm :

- output  $\pi^{(i)}$  = the probability that  $y^{(i)} \in \{0, 1\}$   
 $\pi^{(i)} = P(y^{(i)} = 1 | x^{(i)}, \theta) \in [0, 1]$

Activation function

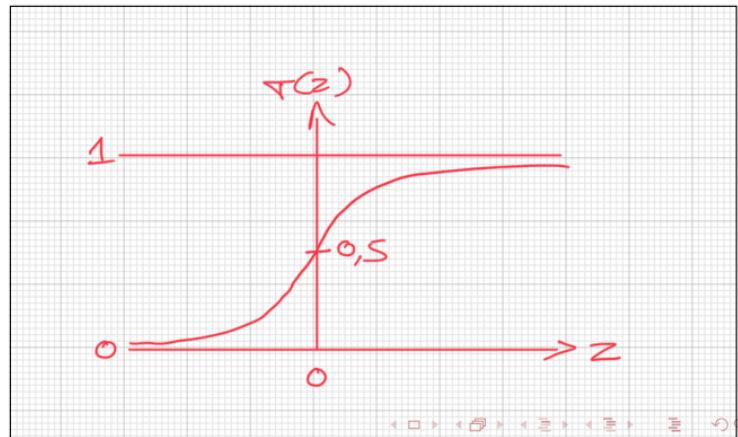
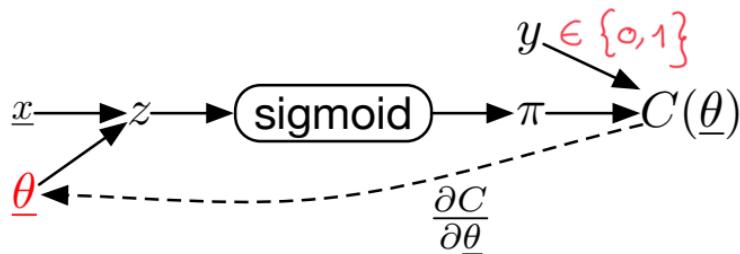
- **Sigmoid function** (Logistic or logit function)

$$\pi = \sigma(z) = \frac{1}{1 + e^{-z}} = \frac{e^z}{e^z + 1}$$

- with  $z = x \cdot \theta$

- Derivative of sigmoid function

$$\begin{aligned}\frac{\partial \sigma(z)}{\partial z} &= \sigma(z)(1 - \sigma(z)) \\ &= \pi^{(i)}(1 - \pi^{(i)})\end{aligned}$$



# Logistic Regression

## Linear separating hyper-plane

The logistic regression defines a linear separating hyper-plane

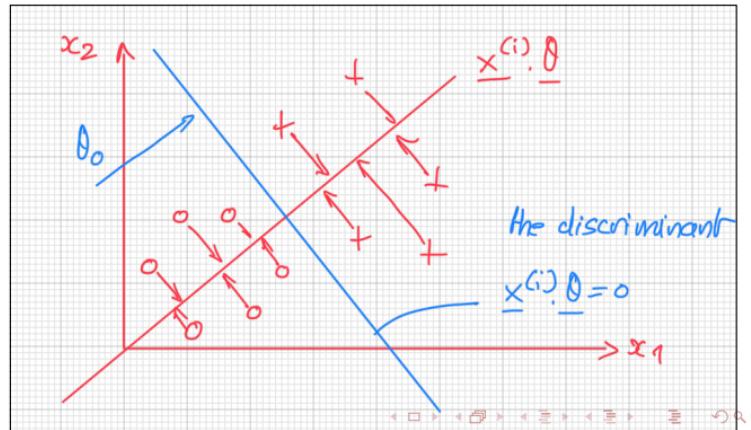
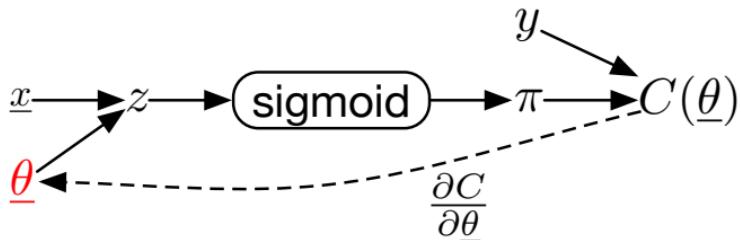
$$P(y^{(i)} = 1 | \underline{x}^{(i)}, \underline{\theta}) = \frac{1}{2}$$

$$\sigma(z) = \frac{1}{2}$$

$$\frac{1}{1 + e^{-z}} = \frac{1}{2}$$

$$z = \underline{x}^{(i)} \underline{\theta} = 0$$

$\underline{x}^{(i)} \underline{\theta} = 0$  is the equation of a plane



# Logistic Regression

## Cost function

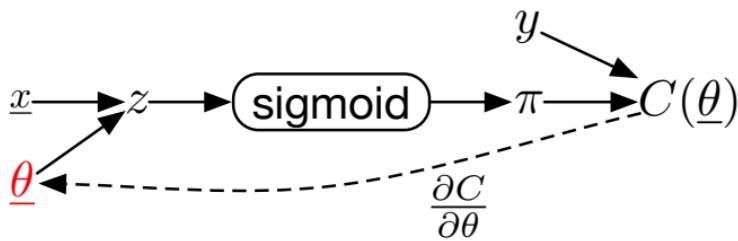
- For Linear Regression  $\Rightarrow$  the uncertainty (random variable) is model by a Gaussian distribution
- For Logistic Regression  $\Rightarrow$  the uncertainty (random variable) is model by a Bernoulli distribution
- Learning is about minimizing uncertainty

Logistic regression specifies the probability of a binary output  $y^{(i)} \in \{0, 1\}$  given input  $\underline{x}^{(i)}$

$$\begin{aligned} p(y|\underline{x}, \underline{\theta}) &= \prod_{i=1}^n p(y^{(i)}|\underline{x}^{(i)}, \underline{\theta}) = \prod_{i=1}^n Ber(y^{(i)}|\underline{x}^{(i)}, \underline{\theta}) \\ &= \prod_{i=1}^n [\pi^{(i)}]^{y^{(i)}} [1 - \pi^{(i)}]^{1-y^{(i)}} \end{aligned}$$

$$\bullet \text{ with } \pi^{(i)} = \frac{1}{1+e^{-z^{(i)}}} \text{ and } z^{(i)} = \underline{x}^{(i)} \underline{\theta}$$

$$\left. \begin{array}{l} y=1 \rightarrow \pi \\ y=0 \rightarrow 1-\pi \end{array} \right\} P(Y=y) = \pi^y (1-\pi)^{1-y}$$



Cost to optimize (minimize) = Negative Log-Likelihood (NLL)

$$C(\underline{\theta}) = -\log p(y|\underline{x}, \underline{\theta}) = \boxed{-\sum_{i=1}^n y^{(i)} \log \pi^{(i)} + (1 - y^{(i)}) \log(1 - \pi^{(i)})}$$

- this is called the binary cross-entropy  $H(y^{(i)}, \pi^{(i)})$

Derivative of the cost w.r.t. to the output

$$\frac{\partial C(\underline{\theta})}{\partial \pi} = -\sum_{i=1}^n \left( \frac{y^{(i)}}{\pi^{(i)}} - \frac{(1 - y^{(i)})}{1 - \pi^{(i)}} \right) = \sum_{i=1}^n \frac{\pi^{(i)} - y^{(i)}}{\pi^{(i)}(1 - \pi^{(i)})}$$

# Logistic Regression

## 1) Optimization using Gradient

To minimize the cost  $C(\underline{\theta})$ , we compute the gradient of the  $C(\underline{\theta})$  w.r.t. to the parameters  $\underline{\theta}$

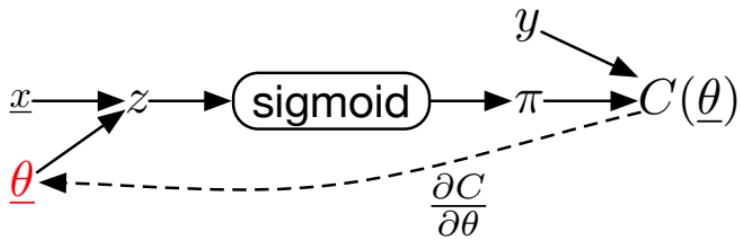
- we use the chain-rule :

$$\begin{aligned} \underline{g} &= \frac{\partial C(\underline{\theta})}{\partial \underline{\theta}} \\ &= \frac{\partial C(\underline{\theta})}{\partial \pi} \frac{\partial \pi}{\partial z} \frac{\partial z}{\partial \theta} \\ &= \sum_{i=1}^n \frac{\pi^{(i)} - y^{(i)}}{\pi^{(i)}(1 - \pi^{(i)})} \cdot \pi^{(i)}(1 - \pi^{(i)}) \cdot \underline{x}^{(i)} \\ &= \sum_{i=1}^n (\underline{x}^{(i)})^T (\pi^{(i)} - y^{(i)}) \\ &= \boxed{\underline{\underline{X}}^T (\underline{\pi} - \underline{y})} \end{aligned}$$

SGD update at iteration  $k + 1$  :

$$\begin{aligned} \underline{\theta}_{k+1} &= \underline{\theta}_k - \eta_k \underline{g}_k \\ &= \boxed{\underline{\theta}_k - \eta_k \underline{\underline{X}}^T (\underline{\pi}_k - \underline{y})} \end{aligned}$$

where  $\pi_k^{(i)} = \sigma(\underline{x}^{(i)} \underline{\theta}_k)$  at iteration  $k$



Note : it looks a bit like the least-square solution :

$$\underline{\theta}_{k+1} = \underline{\theta}_k - \eta_k \left[ 2 \underline{\underline{X}}^T (\underline{\underline{X}} \underline{\theta}_k - \underline{y}) \right]$$

## 2) Optimization using Hessian (Iteratively Reweighted Least Square - IRLS)

Newton update :

$$\theta_{k+1} = \theta_k - \underline{H}_k^{-1} \underline{g}_k$$

- with Gradient, Hessian

$$\underline{g}_k = \underline{X}^T (\underline{\pi}_k - \underline{y})$$

$$\underline{H}_k = \underline{X}^T \underline{S}_k \underline{X}$$

Proof

$$\begin{aligned} H &= \frac{\partial}{\partial \theta} g(\theta)^T \\ &= \sum_{i=1}^n \pi^{(i)} (1 - \pi^{(i)}) \underline{x}^{(i)} (\underline{x}^{(i)})^T \\ &= \underline{X}^T \text{diag}(\pi^{(i)} (1 - \pi^{(i)})) \underline{X} \end{aligned}$$

We note  $\underline{S}_k = \text{diag}(\pi_k^{(1)} (1 - \pi_k^{(1)}) \cdots \pi_k^{(n)} (1 - \pi_k^{(n)}))$

$H$  is positive definite, hence the NLL is convex and has a unique minimum

Newton update at iteration  $k + 1$  :

$$\begin{aligned} \theta_{k+1} &= \theta_k - H_k^{-1} g_k \\ &= \theta_k + \left( \underline{X}^T \underline{S}_k \underline{X} \right)^{-1} \underline{X}^T (\underline{y} - \underline{\pi}_k) \\ &= \left( \underline{X}^T \underline{S}_k \underline{X} \right)^{-1} \left[ (\underline{X}^T \underline{S}_k \underline{X}) \theta_k + \underline{X}^T (\underline{y} - \underline{\pi}_k) \right] \\ &= \boxed{\left( \underline{X}^T \underline{S}_k \underline{X} \right)^{-1} \underline{X}^T \left[ (\underline{S}_k \underline{X}) \theta_k + (\underline{y} - \underline{\pi}_k) \right]} \end{aligned}$$

Note : it looks a bit like the least-square solution :

$$\theta_{k+1} = \left( \underline{X}^T \underline{S}_k \underline{X} \right)^{-1} \underline{X}^T$$



# Softmax regression

## Activation function

The softmax regression is a **multi-class classification** algorithm :

- output  $\pi_c^{(i)} = \text{probability that } y^{(i)} = c$   
 $\pi_c^{(i)} = P(y^{(i)} = c | \underline{x}^{(i)}, \underline{\theta}) \in [0, 1]$

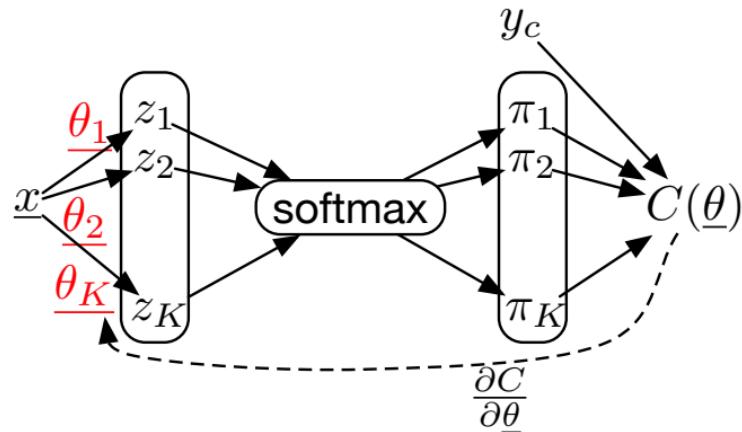
Activation function

- Softmax function

$$\begin{aligned}\pi_c^{(i)} &= \text{softmax}(z_c^{(i)}) = \frac{e^{z_c^{(i)}}}{\sum_{c'=1}^K z_{c'}^{(i)}} \\ &= \frac{e^{x^{(i)} \underline{\theta}_c}}{\sum_{c'=1}^K e^{x^{(i)} \underline{\theta}_{c'}}}\end{aligned}$$

- It is easy to show that

- $\pi_1^{(i)} + \cdots + \pi_K^{(i)} = 1$



# Softmax regression

## Cost function

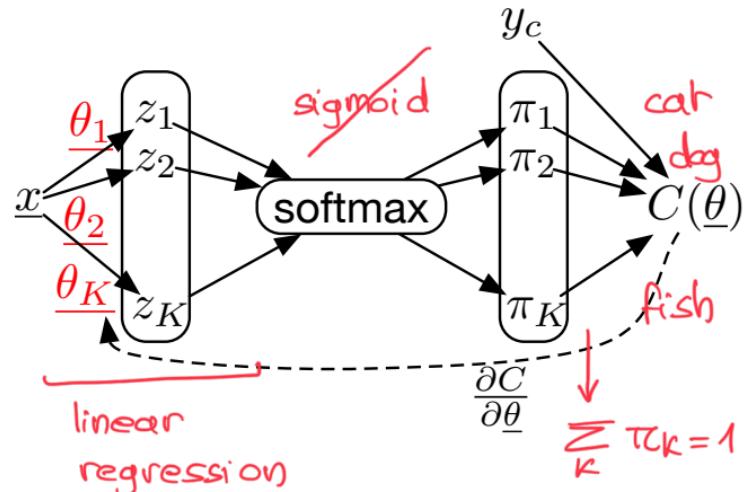
Indicator function  $\mathbb{1}_c(y^{(i)}) \triangleq \begin{cases} 1 & \text{if } y^{(i)} = c \\ 0 & \text{otherwise} \end{cases}$  Likelihood

function for Softmax

$$\begin{aligned} P(\underline{y}|\underline{x}, \underline{\theta}) &= \prod_{i=1}^m p(y^{(i)}|\underline{x}^{(i)}, \underline{\theta}) \\ &= \prod_{i=1}^m \prod_{c=1}^K (\pi_c^{(i)})^{\mathbb{1}_c(y^{(i)})} \end{aligned}$$

Cost to optimize : Negative Log-Likelihood (NLL) for Softmax

$$\begin{aligned} C(\underline{\theta}) &= -\log p(\underline{y}|\underline{x}, \underline{\theta}) \\ &= -\sum_{i=1}^m \sum_{c=1}^K \mathbb{1}_c(y^{(i)}) \log(\pi_c^{(i)}) \end{aligned}$$



# Softmax regression

## Cost function

For the specific case of 2 classes

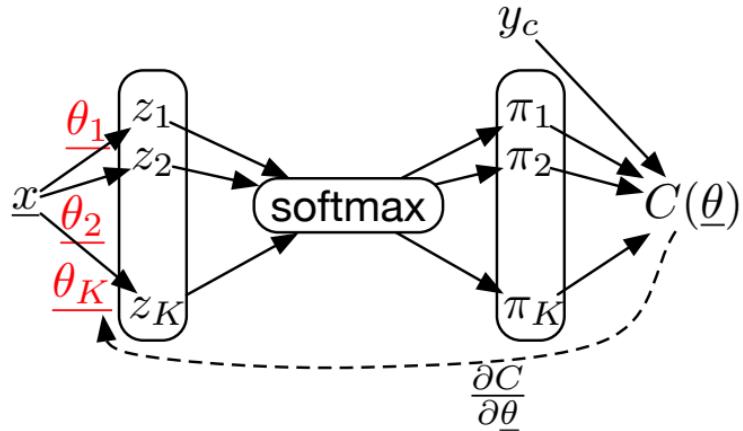
- Likelihood

$$P(y^{(i)} = 1 | \underline{x}^{(i)}, \underline{\theta}) = (\pi_1^{(i)})^1 (\pi_2^{(i)})^0 = \pi_1^{(i)}$$

$$= \frac{e^{\underline{x}^{(i)} \underline{\theta}_1}}{e^{\underline{x}^{(i)} \underline{\theta}_1} + e^{\underline{x}^{(i)} \underline{\theta}_2}}$$

$$P(y^{(i)} = 2 | \underline{x}^{(i)}, \underline{\theta}) = (\pi_1^{(i)})^0 (\pi_2^{(i)})^1 = \pi_2^{(i)}$$

$$= \frac{e^{\underline{x}^{(i)} \underline{\theta}_2}}{e^{\underline{x}^{(i)} \underline{\theta}_1} + e^{\underline{x}^{(i)} \underline{\theta}_2}}$$



- Negative Log-Likelihood (NLL)

$$C(\underline{\theta}) = -\log p(\underline{y} | \underline{x}, \underline{\theta})$$

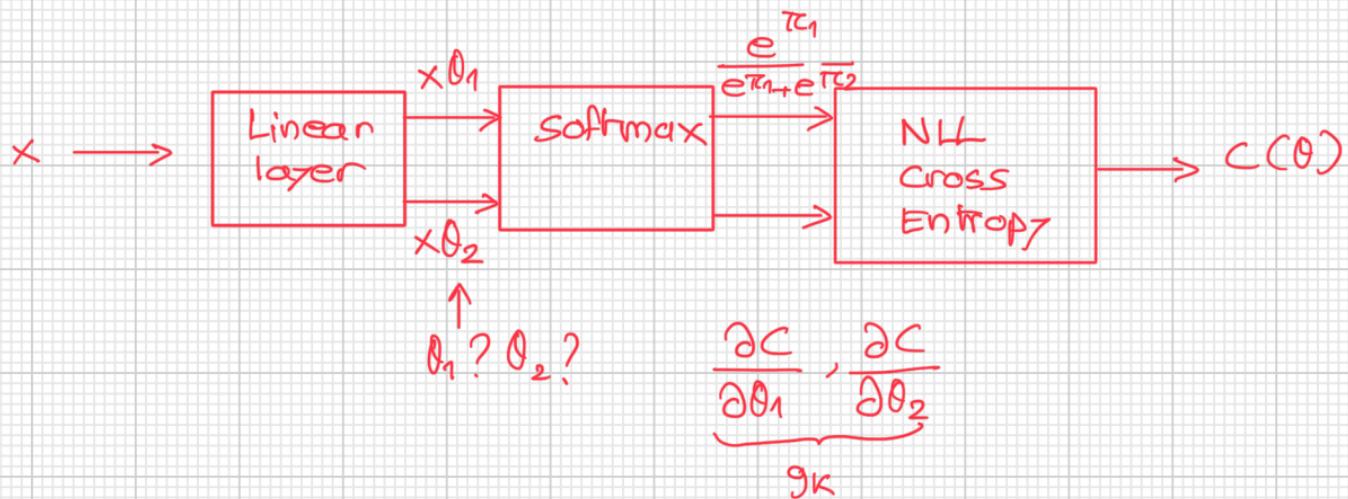
$$= - \sum_{i=1}^m \mathbb{1}_1(y^{(i)}) \log(\pi_1^{(i)}) + \mathbb{1}_2(y^{(i)}) \log(\pi_2^{(i)})$$

$$= - \sum_{i=1}^m [y_1 \log(\pi_1) + y_2 \log(\pi_2) + \dots]$$

# Softmax regression

## Cost function

We can consider the Cost as a new layer



# Softmax regression

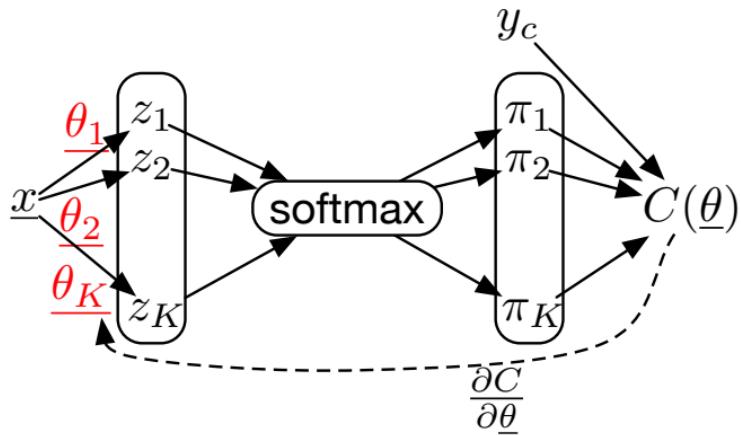
Optimization : derivative of the cost w.r.t. the parameters (manual)

For the specific case of 2 classes

$$\begin{aligned}\frac{\partial}{\partial \theta_2} C(\underline{\theta}) &= \frac{\partial}{\partial \theta_2} \left( -\sum_i \mathbb{1}_1(y^{(i)}) \log(\pi_1^{(i)}) + \mathbb{1}_2(y^{(i)}) \log(\pi_2^{(i)}) \right) \\ &= -\sum_i \mathbb{1}_1(y^{(i)}) \frac{\partial}{\partial \theta_2} \log(\pi_1^{(i)}) + \mathbb{1}_2(y^{(i)}) \frac{\partial}{\partial \theta_2} \log(\pi_2^{(i)})\end{aligned}$$

$$\begin{aligned}\frac{\partial}{\partial \theta_2} \log(\pi_1^{(i)}) &= \frac{\partial}{\partial \theta_2} \log \frac{e^{\underline{x}^{(i)} \underline{\theta}_1}}{e^{\underline{x}^{(i)} \underline{\theta}_1} + e^{\underline{x}^{(i)} \underline{\theta}_2}} \\ &= \frac{\partial}{\partial \theta_2} \left( \underline{x}^{(i)} \underline{\theta}_1 - \log(e^{\underline{x}^{(i)} \underline{\theta}_1} + e^{\underline{x}^{(i)} \underline{\theta}_2}) \right) \\ &= 0 - \frac{\underline{x}^{(i)} e^{\underline{x}^{(i)} \underline{\theta}_2}}{e^{\underline{x}^{(i)} \underline{\theta}_1} + e^{\underline{x}^{(i)} \underline{\theta}_2}} \\ &= -\underline{x}^{(i)} \pi_2^{(i)}\end{aligned}$$

$$\frac{\partial}{\partial \theta_2} \log(\pi_2^{(i)}) = \dots = \underline{x}^{(i)} (1 - \pi_2^{(i)})$$



- Therefore

$$\frac{\partial C(\underline{\theta})}{\partial \theta_2} = -\sum_i \mathbb{1}_1(y^{(i)}) (-\underline{x}^{(i)} \pi_2^{(i)}) + \mathbb{1}_2(y^{(i)}) \underline{x}^{(i)} (1 - \pi_2^{(i)})$$

# Softmax regression

## Optimization : derivative of the cost w.r.t. the parameters (manual)

### For the specific case of 2 classes

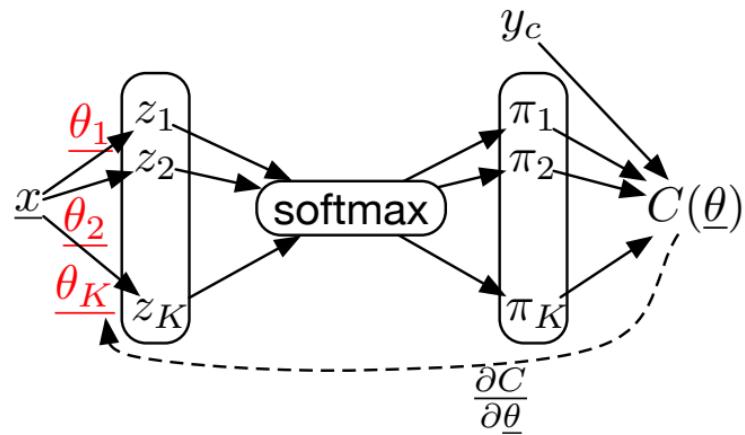
- $\Rightarrow$  equivalence with the logistic regression

Logistic regression  $y^{(i)} \in \{0, 1\}$

- using

$$\begin{aligned}\mathbb{1}_1(y^{(i)}) &= (1 - y^{(i)}) \\ \mathbb{1}_2(y^{(i)}) &= y^{(i)}\end{aligned}$$

$$\begin{aligned}\frac{\partial C(\underline{\theta})}{\partial \theta_2} &= - \sum_i \mathbb{1}_1(y^{(i)}) (-\underline{x}^{(i)} \pi_2^{(i)}) + \mathbb{1}_2(y^{(i)}) \underline{x}^{(i)} (1 - \pi_2^{(i)}) \\ &= - \sum_i (1 - y^{(i)}) (-\underline{x}^{(i)} \pi_2^{(i)}) + y^{(i)} \underline{x}^{(i)} (1 - \pi_2^{(i)}) \\ &= - \sum_i -\underline{x}^{(i)} \pi_2^{(i)} + y^{(i)} \underline{x}^{(i)} \pi_2^{(i)} + y^{(i)} \underline{x}^{(i)} - y^{(i)} \underline{x}^{(i)} \pi_2^{(i)} \\ &= \sum_i \underline{x}^{(i)} (\pi_2^{(i)} - y^{(i)})\end{aligned}$$





# Backpropagation

## ~~Logistic Regression~~

### ~~Softmax~~

$$C(\theta) = - \sum_{i=1}^n \mathbb{1}_1(y^{(i)}) \log \left( \frac{e^{\underline{\theta}_{1,:}\underline{x}^{(i)}}}{e^{\underline{\theta}_{1,:}\underline{x}^{(i)}} + e^{\underline{\theta}_{2,:}\underline{x}^{(i)}}} \right) + \mathbb{1}_2(y^{(i)}) \log \left( \frac{e^{\underline{\theta}_{2,:}\underline{x}^{(i)}}}{e^{\underline{\theta}_{1,:}\underline{x}^{(i)}} + e^{\underline{\theta}_{2,:}\underline{x}^{(i)}}} \right)$$

We consider the ~~logistic regression~~ as a succession of layers  $\underline{z}^{[l]}$

~~softmax~~

$$\underline{z}^{[0]} = \underline{x}^{(i)}$$

$$z^{[0](i)} \quad z_1^{[1]} = \underline{\theta}_{1,:}\underline{x} \quad z_1^{[2]} \quad z^{[3]}$$

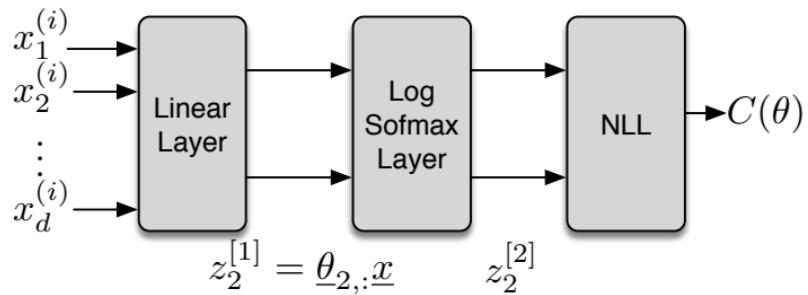
$$z_1^{[1]} = \underline{\theta}_{1,:}\underline{z}^{[0]}$$

$$z_2^{[1]} = \underline{\theta}_{2,:}\underline{z}^{[0]}$$

$$z_1^{[2]} = \log \left( \frac{e^{z_1^{[1]}}}{e^{z_1^{[1]}} + e^{z_2^{[1]}}} \right)$$

$$z_2^{[2]} = \log \left( \frac{e^{z_2^{[1]}}}{e^{z_1^{[1]}} + e^{z_2^{[1]}}} \right)$$

$$z^{[3]} = - \sum_i \mathbb{1}_1(y^{(i)}) z_1^{[2]} + \mathbb{1}_2(y^{(i)}) z_2^{[2]}$$



# Backpropagation

## Logistic Regression

We consider the logistic regression as a succession of layers  $z^{[l]}$

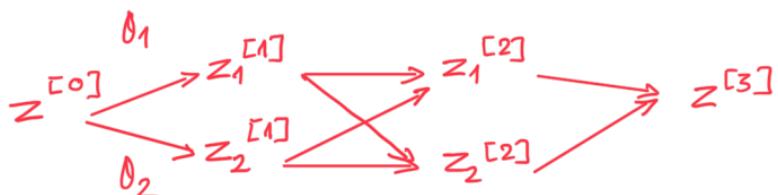
- **Composition of functions**

$$C(\theta) = z^{[3]} \left\{ z_1^{[2]} \left[ z_1^{[1]}(z^{[0]}, \underline{\theta}_1), z_2^{[1]}(z^{[0]}, \underline{\theta}_2) \right], z_2^{[2]} \left[ z_1^{[1]}(z^{[0]}, \underline{\theta}_1), z_2^{[1]}(z^{[0]}, \underline{\theta}_2) \right] \right\}$$

- **Computing derivatives using the chain rule**

- inefficient method : we compute several times the same things

$$\frac{\partial C(\theta)}{\partial \theta_1} = \frac{\partial z^{[3]}}{\partial z_1^{[2]}} \frac{\partial z_1^{[2]}}{\partial z_1^{[1]}} \frac{\partial z_1^{[1]}}{\partial \theta_1} + \frac{\partial z^{[3]}}{\partial z_1^{[2]}} \frac{\partial z_1^{[2]}}{\partial z_2^{[1]}} \frac{\partial z_2^{[1]}}{\partial \theta_1} + \frac{\partial z^{[3]}}{\partial z_2^{[2]}} \frac{\partial z_2^{[2]}}{\partial z_1^{[1]}} \frac{\partial z_1^{[1]}}{\partial \theta_1} + \frac{\partial z^{[3]}}{\partial z_2^{[2]}} \frac{\partial z_2^{[2]}}{\partial z_2^{[1]}} \frac{\partial z_2^{[1]}}{\partial \theta_1}$$



# Backpropagation

## Layer specification

We only have to design a **module** with a

- **Forward method :**

$$z^{[l]} = f_l(z^{[l-1]}, \theta^{[l]})$$

- **Backward method :**

(for  $i$  inputs and  $o$  outputs)

- Input :

$$\delta_o^{[l]} \triangleq \frac{\partial C}{\partial z_o^{[l]}}$$

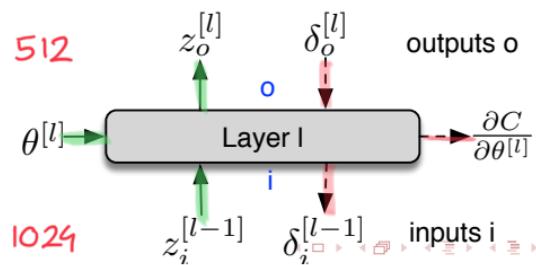
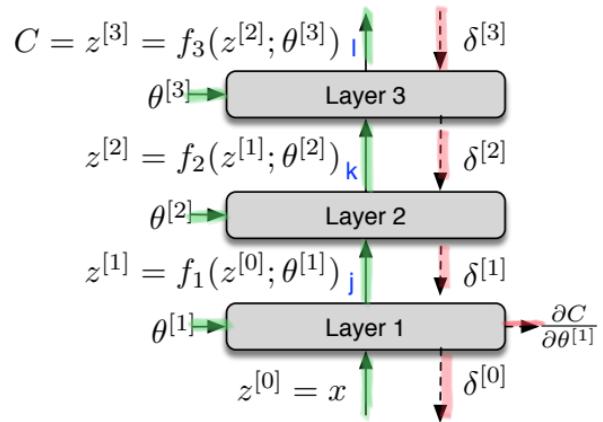
- Computation

$$\delta_i^{[l-1]} = \sum_o \frac{\partial C}{\partial z_o^{[l]}} \frac{\partial z_o^{[l]}}{\partial z_i^{[l-1]}}$$

$$= \sum_o \delta_o^{[l]} \frac{\partial z_o^{[l]}}{\partial z_i^{[l-1]}}$$

$$\frac{\partial C}{\partial \theta^{[l]}} = \sum_o \frac{\partial C}{\partial z_o^{[l]}} \frac{\partial z_o^{[l]}}{\partial \theta^{[l]}}$$

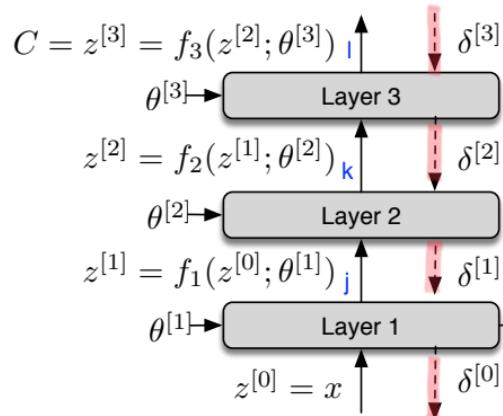
$$= \sum_o \delta_o^{[l]} \frac{\partial z_o^{[l]}}{\partial \theta^{[l]}}$$



# Backpropagation

## Derivative via layer-specification

$$\begin{aligned}
 \frac{\partial C}{\partial \theta^{[1]}} &= \sum_j \frac{\partial C}{\partial z_j^{[1]}} \frac{\partial z_j^{[1]}}{\partial \theta^{[1]}} \\
 &= \sum_j \left( \sum_k \frac{\partial C}{\partial z_k^{[2]}} \frac{\partial z_k^{[2]}}{\partial z_j^{[1]}} \right) \frac{\partial z_j^{[1]}}{\partial \theta^{[1]}} \\
 &= \sum_j \left( \sum_k \left( \sum_l \frac{\partial C}{\partial z_l^{[3]}} \frac{\partial z_l^{[3]}}{\partial z_k^{[2]}} \right) \frac{\partial z_k^{[2]}}{\partial z_j^{[1]}} \right) \frac{\partial z_j^{[1]}}{\partial \theta^{[1]}} \\
 &= \sum_j \left( \sum_k \left( \sum_{l=1}^L \frac{\partial z_l^{[3]}}{\partial z_k^{[2]}} \right) \frac{\partial z_k^{[2]}}{\partial z_j^{[1]}} \right) \frac{\partial z_j^{[1]}}{\partial \theta^{[1]}}
 \end{aligned}$$



- Note : this is equivalent to what we have seen above but without redundant computation

$$\frac{\partial C(\underline{\theta})}{\partial \theta_1} = \frac{\partial z^{[3]}}{\partial z_1^{[2]}} \frac{\partial z_1^{[2]}}{\partial z_1^{[1]}} \frac{\partial z_1^{[1]}}{\partial \theta_1} + \frac{\partial z^{[3]}}{\partial z_1^{[2]}} \frac{\partial z_1^{[2]}}{\partial z_2^{[1]}} \frac{\partial z_2^{[1]}}{\partial \theta_1} + \frac{\partial z^{[3]}}{\partial z_2^{[2]}} \frac{\partial z_2^{[2]}}{\partial z_1^{[1]}} \frac{\partial z_1^{[1]}}{\partial \theta_1} + \frac{\partial z^{[3]}}{\partial z_2^{[2]}} \frac{\partial z_2^{[2]}}{\partial z_2^{[1]}} \frac{\partial z_2^{[1]}}{\partial \theta_1}$$

Forward  $z^{[0]} = \underline{x}^{(i)} \rightarrow z^{[1]}(\underline{x}^{(i)}) \rightarrow z^{[2]}(\underline{x}^{(i)}) \rightarrow z^{[3]}(\underline{x}^{(i)}) = C$

Backward  $\delta^{[0]} \leftarrow \delta^{[1]} \leftarrow \delta^{[2]} \leftarrow \delta^{[3]} = 1$

# Neural Networks - Multi-Layer-Perceptron (MLP)

# Neural Networks - Multi-Layer-Perceptron (MLP)

... with 1 layer (0 hidden layer) / Binary classification

- Layer 1 :

$$z^{[1](i)} = \theta_0^{[1]} + \theta_1^{[1]} x_1^{(i)} + \theta_2^{[1]} x_2^{(i)}$$

$$z^{[2](i)} = \frac{1}{1 + e^{-z^{[1](i)}}}$$

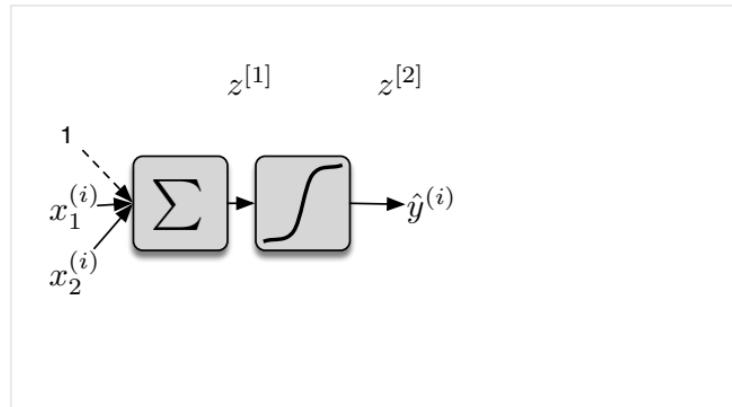
- Output

$$\hat{y}^{(i)} = z^{[2](i)}$$

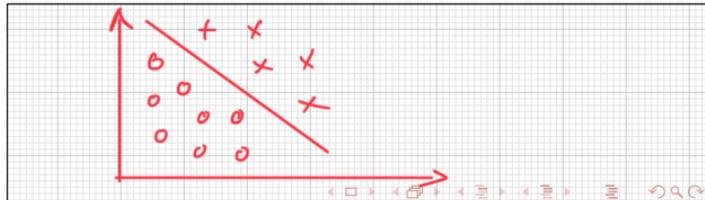
$$\hat{y}^{(i)} = p(y^{(i)} | x^{(i)}, \underline{\theta})$$

- Binary classification

$$\begin{aligned} p(y^{(i)} | x^{(i)}, \underline{\theta}) &= (\hat{y}^{(i)})^{y^{(i)}} (1 - \hat{y}^{(i)})^{1-y^{(i)}} \\ &= \begin{cases} \hat{y}^{(i)} & \text{when } y^{(i)} = 1 \\ 1 - \hat{y}^{(i)} & \text{when } y^{(i)} = 0 \end{cases} \end{aligned}$$



- Linear separating hyper-plane



# Neural Networks - Multi-Layer-Perceptron (MLP)

... with 2 layers (1 hidden layer) / Binary classification

- Layer 1 (hidden)

$$z_1^{[1]} = \theta_{0,1}^{[1]} + \theta_{1,1}^{[1]}x_1^{(i)} + \theta_{2,1}^{[1]}x_2^{(i)}$$

$$z_2^{[1]} = \theta_{0,2}^{[1]} + \theta_{1,2}^{[1]}x_1^{(i)} + \theta_{2,2}^{[1]}x_2^{(i)}$$

$$z_1^{[2]} = 1/(1 + e^{-z_1^{[1]}})$$

$$z_2^{[2]} = 1/(1 + e^{-z_2^{[1]}})$$

$$\text{e}^{-\frac{\|x-\mu\|}{\sigma}^2}$$

- Layer 2

$$z_1^{[3]} = \theta_{0,1}^{[2]} + \theta_{1,1}^{[2]}z_1^{[2]} + \theta_{2,1}^{[2]}z_2^{[2]}$$

$$z^{[4]} = \frac{1}{1 + e^{-z_1^{[3]}}}$$

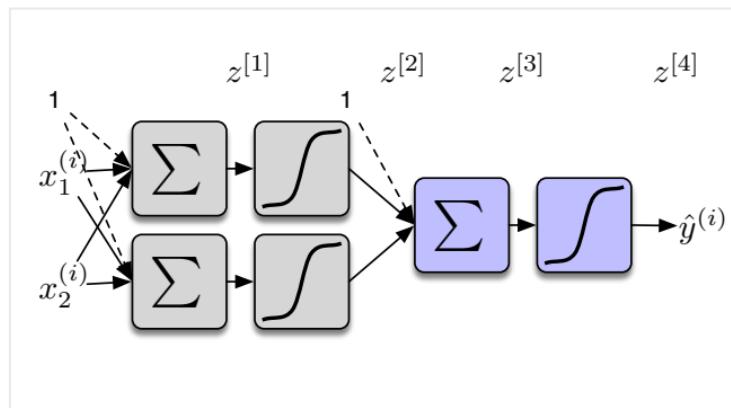
- Output

$$\hat{y}^{(i)} = z^{[4]}$$

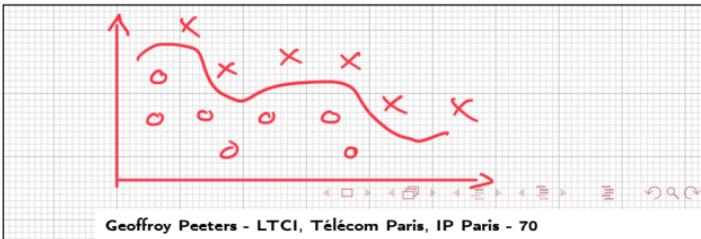
- Cost : minimize the binary cross-entropy

$$C(\theta) = -\log p(y|x, \theta)$$

$$= - \sum_{i=1}^n y^{(i)} \log \hat{y}^{(i)} + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)})$$



- Non-Linear hyper-plane



# Neural Networks - Multi-Layer-Perceptron (MLP)

... with 2 layers (1 hidden layer) / Regression

Neural networks are function approximation tools

- Model

$$\begin{aligned}\hat{y}^{(i)} &= \theta_{0,1}^{[2]} \\ &+ \theta_{1,1}^{[2]} \frac{1}{1 - e^{-(\theta_{0,1}^{[1]} + \theta_{1,1}^{[1]} x^{(i)})}} \\ &+ \theta_{2,1}^{[2]} \frac{1}{1 - e^{-(\theta_{0,2}^{[1]} + \theta_{1,2}^{[1]} x^{(i)})}}\end{aligned}$$

- Output

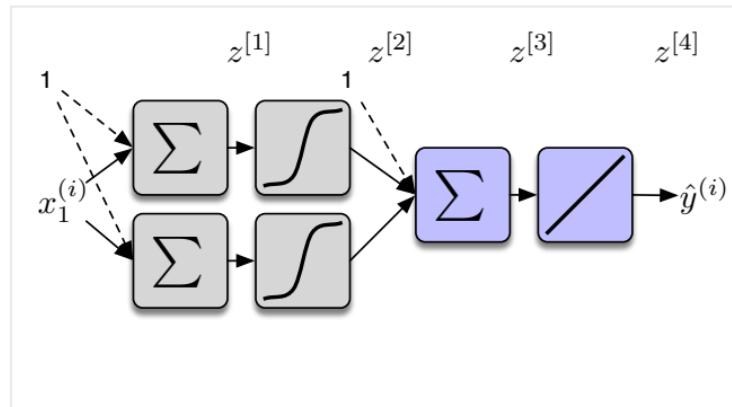
$$\hat{y}^{(i)} = \mathbb{E}(y^{(i)} | x^{(i)}, \theta)$$

- Likelihood/ Cost

$$p(y^{(i)} | x^{(i)}, \theta) = (2\pi\sigma^2)^{-1/2} e^{-\frac{1}{2\sigma^2}(\hat{y}^{(i)} - y^{(i)})^2}$$

$$C(\theta) = \sum_{i=1}^n (y^{(i)} - \hat{y}^{(i)})^2$$

- with  $\hat{y}^{(i)} = f(x^{(i)}, \theta)$

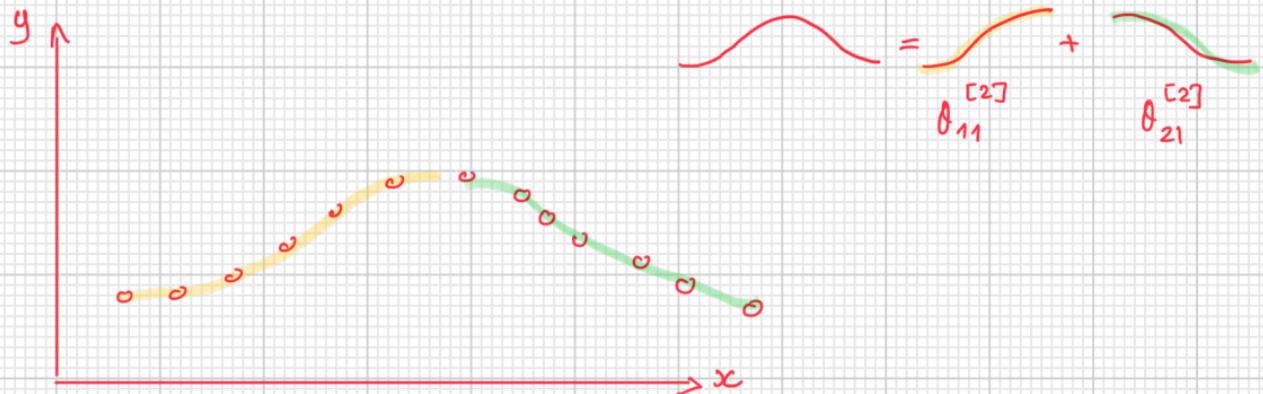


# Neural Networks - Multi-Layer-Perceptron (MLP)

... with 2 layers (1 hidden layer) / Regression

$$\hat{y}^{(i)} = \theta_{0,1}^{[2]} + \theta_{1,1}^{[2]} \frac{1}{1 - e^{-(\theta_{0,1}^{[1]} + \theta_{1,1}^{[1]} x^{(i)})}} + \theta_{2,1}^{[2]} \frac{1}{1 - e^{-(\theta_{0,2}^{[1]} + \theta_{1,2}^{[1]} x^{(i)})}} \quad \text{RBF regression}$$

$$\frac{\|x - N_d\|^2}{n}$$



# Neural Networks - Multi-Layer-Perceptron (MLP)

... with 2 layers (1 hidden layer) / Multi-Class classification

- Output :

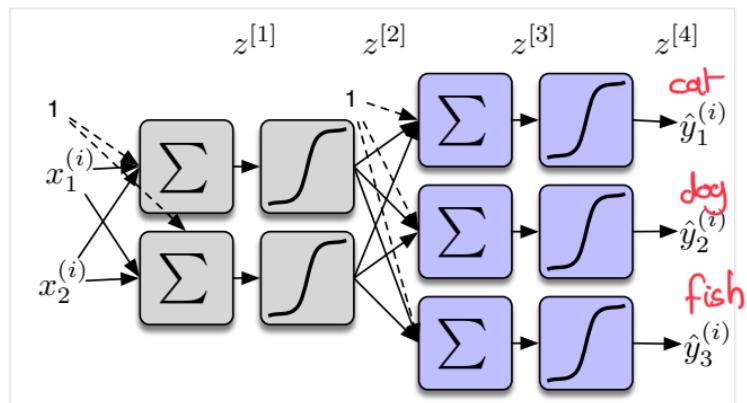
$$\begin{aligned} p(\underline{y}^{(i)} = [010] | \underline{x}^{(i)}, \theta) &= p(y^{(i)} = 2 | \underline{x}^{(i)}, \theta) \\ &= \text{softmax}(z_2^{(i)}) \\ &= \frac{e^{z_2^{(i)}}}{e^{z_1^{(i)}} + e^{z_2^{(i)}} + e^{z_3^{(i)}}} \end{aligned}$$

- Likelihood

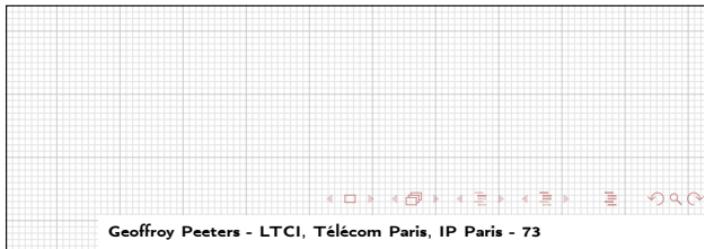
$$p(\underline{y}^{(i)} | \underline{x}^{(i)}, \theta) = \left[ \frac{e^{z_1^{(i)}}}{\sum} \right]^{\mathbb{1}_1(y^{(i)})} + \left[ \frac{e^{z_2^{(i)}}}{\sum} \right]^{\mathbb{1}_2(y^{(i)})} + \left[ \frac{e^{z_3^{(i)}}}{\sum} \right]^{\mathbb{1}_3(y^{(i)})}$$

- Cost : minimize the cross-entropy

$$\begin{aligned} C(\theta) &= -\log p(\underline{y} | \underline{x}, \theta) \\ &= -\sum_{i=1}^n \sum_{c=1}^3 \mathbb{1}_c(y^{(i)}) \log \frac{e^{z_c^{(i)}}}{\sum} \end{aligned}$$



- Separating hyper-planes

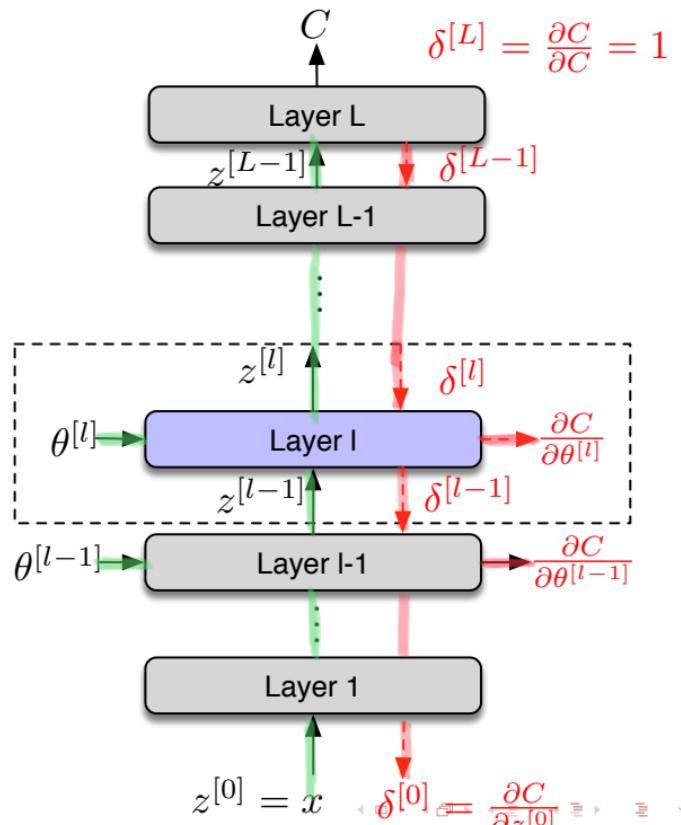




# Deep Learning

## Deep Learning ?

- neural network with many many layers
- vanishing gradient problem (sigmoid, ReLu)
- sequential  $\neq$  recursive



## Forward/ Backward propagation

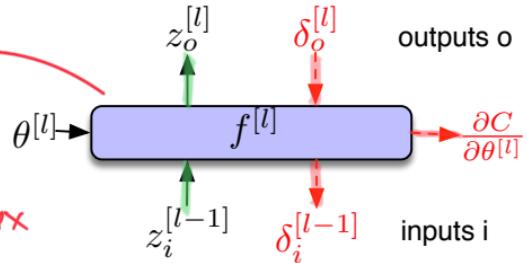
### Layer $l$ :

- Input :  $\underline{z}^{[l-1]}$  / Output :  $\underline{z}^{[l]}$

### Forward :

$$\underline{z}^{[l]} = f^{[l]}(\underline{z}^{[l-1]}, \theta^{[l]})$$

- Linear  $\underline{\theta} \times$
- ReLU
- Sigmoid, Softmax



### Backward :

### Gradients to propagate back

- matrix form

$$\underline{\delta}^{[l-1]} = \frac{\partial C}{\partial \underline{z}^{[l-1]}} = \frac{\partial C}{\partial \underline{z}^{[l]}} \frac{\partial \underline{z}^{[l]}}{\partial \underline{z}^{[l-1]}} = \underline{\delta}^{[l]} \frac{\partial \underline{z}^{[l]}}{\partial \underline{z}^{[l-1]}}$$

• where  $\frac{\partial \underline{z}^{[l]}}{\partial \underline{z}^{[l-1]}}$  is a Jacobian matrix

- scalar form

$$\begin{aligned} \delta_i^{[l-1]} &= \sum_o \frac{\partial C}{\partial z_o^{[l]}} \frac{\partial z_o^{[l]}}{\partial z_i^{[l-1]}} \\ &= \sum_o \delta_o^{[l]} \frac{\partial [f_o^{[l]}(\underline{z}^{[l-1]}, \underline{\theta}^{[l]})]}{\partial z_i^{[l-1]}} \end{aligned}$$

### Local gradients for parameters updating

- matrix form

$$\frac{\partial C}{\partial \underline{\theta}^{[l]}} = \frac{\partial C}{\partial \underline{z}^{[l]}} \frac{\partial \underline{z}^{[l]}}{\partial \underline{\theta}^{[l]}} = \underline{\delta}^{[l]} \frac{\partial f^{[l]}(\underline{z}^{[l-1]}, \underline{\theta}^{[l]})}{\partial \underline{\theta}^{[l]}}$$

- scalar form

$$\begin{aligned} \frac{\partial C}{\partial \theta_i^{[l]}} &= \sum_o \frac{\partial C}{\partial z_o^{[l]}} \frac{\partial z_o^{[l]}}{\partial \theta_i^{[l]}} \\ &= \sum_o \delta_o^{[l]} \frac{\partial [f_o^{[l]}(\underline{z}^{[l-1]}, \underline{\theta}^{[l]})]}{\partial \theta_i^{[l]}} \end{aligned}$$

## Two ways to implement gradient computation

Suppose we have the following composition of functions

$$f^{[L]}(\underline{f}^{[L-1]}(\cdots \underline{f}^{[2]}(\underline{f}^{[1]}(\underline{x}))))$$

- **Solution 1**

$$\frac{\partial f^{[L]}(\underline{f}^{[L-1]}(\cdots \underline{f}^{[2]}(\underline{f}^{[1]}(\underline{x}))))}{\partial \underline{x}} = \frac{\partial f^{[L]}}{\partial \underline{f}^{[L-1]}} \cdot \frac{\partial \underline{f}^{[L-1]}}{\partial \underline{f}^{[L-2]}} \cdots \frac{\partial \underline{f}^{[2]}}{\partial \underline{f}^{[1]}} \cdot \frac{\partial \underline{f}^{[1]}}{\partial \underline{x}}$$

- chain rule
- need a lot of memory (need to store all matrices)

- **Solution 2**

$$\frac{\partial f^{[L]}(\underline{f}^{[L-1]}(\cdots \underline{f}^{[2]}(\underline{f}^{[1]}(\underline{x}))))}{\partial \underline{x}} = \left( \left( \left( \left( \frac{\partial \underline{f}^{[L]}}{\partial \underline{f}^{[L-1]}} \cdot \frac{\partial \underline{f}^{[L-1]}}{\partial \underline{f}^{[L-2]}} \right) \cdots \right) \frac{\partial \underline{f}^{[2]}}{\partial \underline{f}^{[1]}} \right) \cdot \frac{\partial \underline{f}^{[1]}}{\partial \underline{x}} \right)$$

Jacobian

- compute first (-) then second (-) then third (-)
- less memory (the intermediate results are always vectors : vector x matrix = vector)
- $\Rightarrow$  this is **the back-propagation algorithm**

## Linear layer

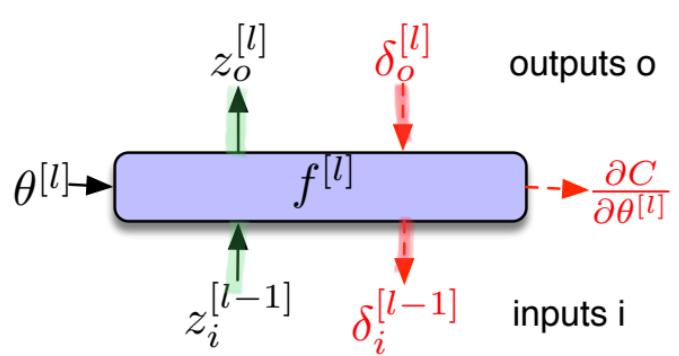
- Forward

$$z_o^{[l]} = f_o(\underline{z}^{[l-1]}, \underline{\theta}) = \sum_i z_i^{[l-1]} \theta_{io}$$

- Backward

$$\delta_i^{[l-1]} = \sum_o \delta_o^{[l]} \frac{\partial [f_o(z^{[l-1]}, \theta)]}{\partial z^{[l-1]}} = \sum_o \delta_o^{[l]} \theta_{io}$$

$$\frac{\partial C}{\partial \theta_{io}} = \sum_o \delta_o^{[I]} \frac{\partial [f_o(\underline{z}^{[I-1]}, \underline{\theta})]}{\partial \theta_{io}} = \delta_o^{[I]} z_i^{[I-1]}$$



## ReLU layer

- Forward**

$$z_o = f_o(z_o^{[l-1]}) = \max(0, z_o^{[l-1]})$$

- Backward**

$$\begin{aligned} \delta_i^{[l-1]} &= \sum_o \delta_o^{[l]} \frac{\partial [f_o(z_o^{[l-1]})]}{\partial z_i^{[l-1]}} \\ &= \delta_i^{[l]} \mathbb{1}_{[z_i^{[l-1]} > 0]} \end{aligned}$$

