

Talent course July 22

Classification & Logistic regression

$$\frac{\partial C(\beta)}{\partial \beta} = 0 = -X^T(y - p)$$

$$p(y_i | x_i, \beta) = \frac{\exp(\beta_0 + \beta_1 x_i)}{1 + \exp(\beta_0 + \beta_1 x_i)}$$

$$\underline{y_i = 1}$$

$$p(y_i = 0 | x_i, \beta) = 1 - p(y_i = 1 | x_i, \beta)$$

$$\boxed{\frac{\partial^2 C}{\partial \beta \partial \beta^T}} > 0, \text{ convex function}$$

Iterative solver like
Newton-Raphson:

$$\underline{x_{n+1}} = x_n - \frac{f(x_n)}{f'(x_n)}$$

$$f(s) = 0 = f(x) + (s-x)f'(x) + \frac{(s-x)^2}{2}f''(x) + \dots$$

$$f(x) + (S-x) f'(x) \stackrel{!}{=} 0$$

$$S = x - \frac{f(x)}{f'(x)} \sim \frac{\partial C}{\partial \beta} = 0$$

$$|x_{n+1} - x_n| < \epsilon \sim 10^{-10}$$

$$\beta^{\text{new}} = \beta^{\text{old}} - \frac{1}{\left(\frac{\partial^2 C}{\partial \beta \partial \beta^T} \right)_{\beta^{\text{old}}}} \times \left(\frac{\partial C}{\partial \beta} \right)_{\beta^{\text{old}}}$$

$$= \beta^{\text{old}} - \boxed{\frac{(X^T W X)^{-1}}{\text{Hessian matrix}}} \left(-X^T (y - p) \right)_{\beta^{\text{old}}}$$

Conjugate gradient method
(gradient descent)

$$(X^T W X)^{-1} \equiv \gamma = \text{constant}$$

$$\boxed{\text{new} = \text{old} - \gamma \nabla C(\beta)}$$

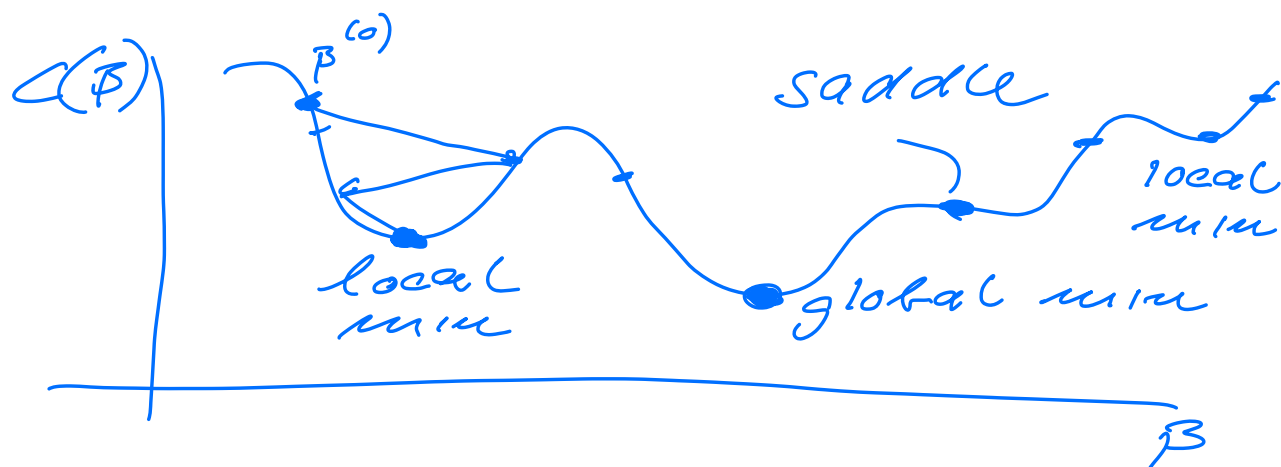
$$\beta = \beta - \gamma \text{VCLP}|_{\beta^{\text{old}}}$$

parameter in ML

$$\gamma \in [10^{-5}, 10^{-9}, \dots, 10^{-1}]$$

γ = LEARNING RATE

To avoid local minima one uses the family of stochastic gradient descent methods.



$\nabla C(\beta)$ using automatic differentiation (autograd)

$$\nabla C(\mu) = -\sqrt{T}(\mu - \mu^*)$$

$$V(\beta) = \frac{1}{2} (Y - X\beta)^T (Y - X\beta)$$

$$X \in \mathbb{R}^{n \times p}$$

$$y \in \mathbb{R}^n$$

$$\beta \in \mathbb{R}^p$$

$$\frac{\partial^2 V(\beta)}{\partial \beta \partial \beta^T} = X^T W X$$

$$W \in \mathbb{R}^{n \times n}$$

$$\rightarrow \in \mathbb{R}^{p \times p}$$

$$VC(\beta) \in \mathbb{R}^p$$

$$n \gg p, \quad n \geq p$$

$$p \gg n$$

$$n \sim 10^3 - 10^5$$

$$p \sim 10^1 - 10^2$$

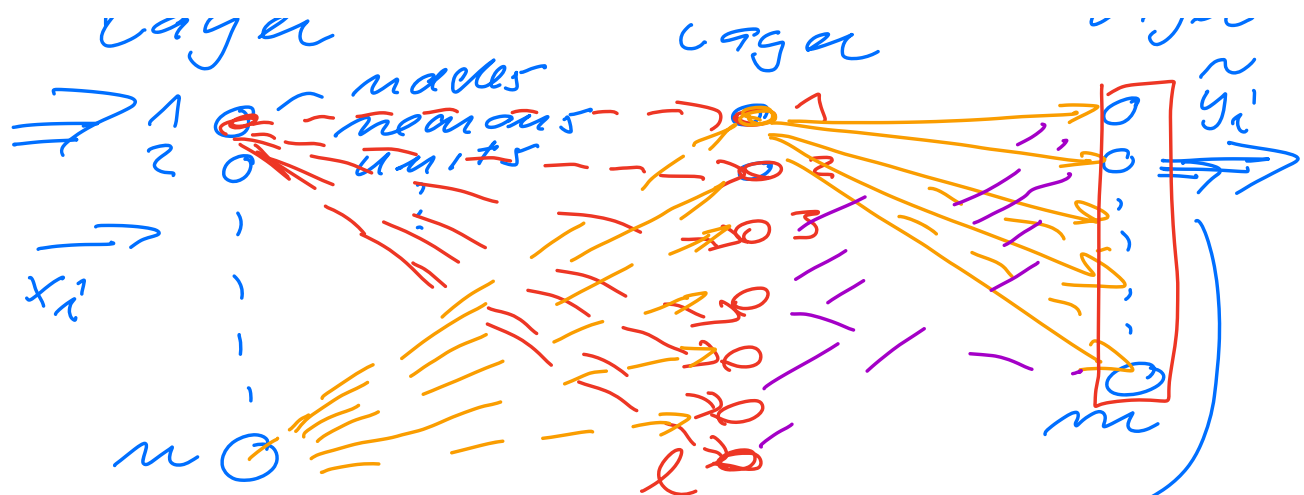
Neural Networks (NN)

Feed Forward NN

input
layer

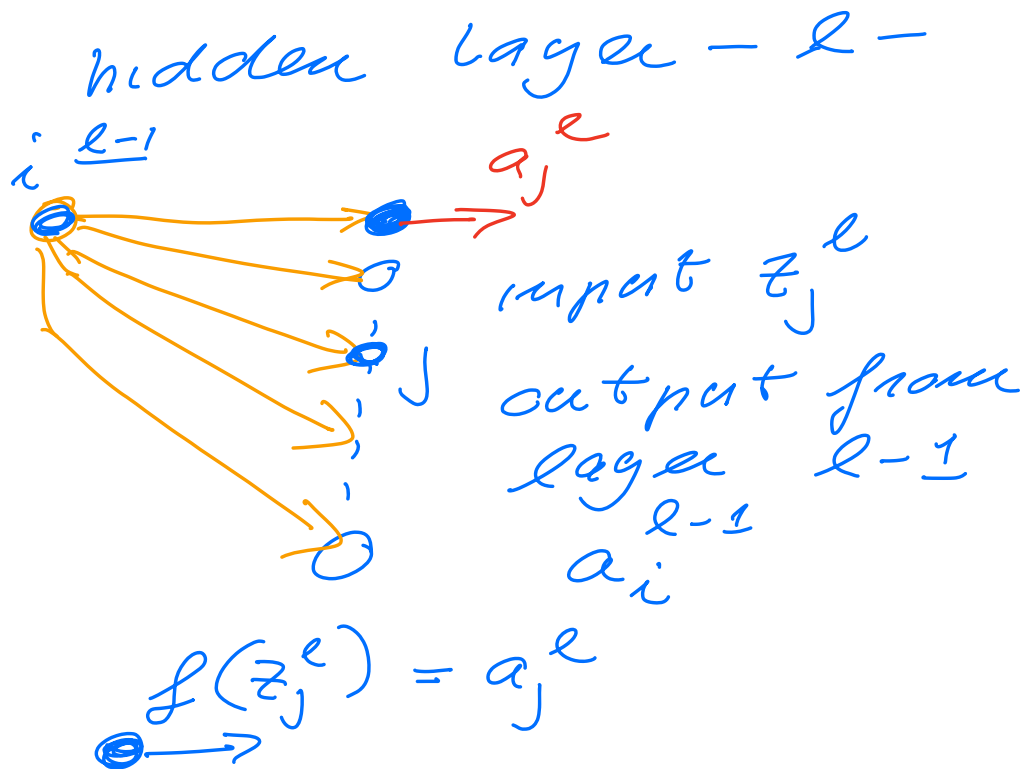
hidden

output
layer



Compared with the target $y_i(t_i)$

Basic parameters



$$z_j^l = \sum_{i=1}^{M_{l-1}} w_{ij}^l a_i^{l-1} + b_j^l$$

\nwarrow \nearrow
weights biases-
parameters-
 to fit in opti
 mization.

z_j^l produces the output from a node, a_j^l

$$a_j^l = \underline{f(z_j^l)} = \frac{1}{1 + e^{-z_j^l}}$$

\downarrow
 activation function

Gradients of the cost function to optimize weights w_{ij}^l and b_j^l

Algorithm to compute the gradients =

Back propagation algo (chain rule),

Parameters:

- w_{ij}^l for each layer,
- b_j^l for each layer
- activation function
 - Logit
 - tanh
 - RELU
 - ELU
 - ...
- Number of hidden layers
- Number of nodes (except input and output)
- Learning rate & (grid search)

- hyperparameters λ
 norm-2 , norm-1 - ...
- choice of gradient descent method.