

GRADIENT-BASED OPTIMIZATION METHODS FOR NEURAL NETWORK TRAINING

PH.D. IN INDUSTRIAL ENGINEERING
UNIVERSITÀ DI FIRENZE, A.A. 2020/2021

Simone Rebegoldi
Dipartimento di Ingegneria Industriale
Università di Firenze



UNIVERSITÀ
DEGLI STUDI
FIRENZE

Neural network: ϕ_i -evaluation

$$\begin{cases} \mathbf{v}_1 = a \\ \mathbf{v}_j = \sigma_j(\mathbf{W}_j \mathbf{v}_{j-1} + \beta_j) \quad j = 2, \dots, m \quad m \text{ layers} \end{cases}$$

We have to minimize the empirical risk:

$$f(x) = \frac{1}{N} \sum_{i=1}^N \underbrace{\ell(\mathbf{v}_m(x; \mathbf{a}_i), b_i)}_{\phi_i} \quad x = (\text{vec}(\mathbf{W}_2), \beta_2^T, \dots, \text{vec}(\mathbf{W}_m), \beta_m^T)^T,$$

$$\mathbf{W}_j = \begin{bmatrix} w_{j,1,1} & w_{j,1,2} & \cdots & \cdots & w_{j,1,q_j} \\ w_{j,2,1} & w_{j,2,2} & \cdots & \cdots & w_{j,2,q_j} \\ \vdots & \vdots & & & \vdots \\ w_{j,q_{j-1},1} & w_{j,q_{j-1},2} & \cdots & \cdots & w_{j,q_{j-1},q_j} \end{bmatrix}$$

Forward propagation Algorithm

Input: $m \geq 2$, \mathbf{W}_j , β_j , σ_j for $j = 2, \dots, m$, $\ell(a_i, b_i)$ sample-label

Output: $L = \phi_i(x)$

Step 0: $v_1 \leftarrow a$

Inner loop:

For $j = 2, \dots, m$

$$\mathbf{r}_j \leftarrow \beta_j + \mathbf{W}_j \mathbf{v}_{j-1}$$

$$\mathbf{v}_j \leftarrow \sigma_j(\mathbf{r}_j)$$

m matrix-vector products

$L \leftarrow \ell(\mathbf{v}_m, b_i)$

Neural network: $\nabla\phi$ -evaluation

$$\begin{cases} \mathbf{v}_1 = a \\ \mathbf{v}_j = \boldsymbol{\sigma}_j (\mathbf{W}_j \mathbf{v}_{j-1} + \beta_j) \quad j = 2, \dots, m \end{cases} \quad \phi(x) = \ell(\mathbf{v}_m(x; \mathbf{a}_i), \mathbf{b}_i)$$

Backward propagation Algorithm

Input: $m \geq 2$, \mathbf{W}_j , β_j , $\boldsymbol{\sigma}_j$ for $j = 2, \dots, m$, ℓ , (a_i, b_i) sample-label,
 L , $\mathbf{v}_j \in \mathbf{r}_j$, $j = 2, \dots, m$ computed by the Forward propagation algorithm

Output: $\tilde{\nabla}_{\mathbf{W}_j} \phi_i(x)$ e $\nabla_{\beta_j} \phi_i(x)$ for $j = 2, \dots, m$

Step 0: $\mathbf{g} \leftarrow \nabla_{\mathbf{v}_m} L = \nabla_{\mathbf{v}_m} \ell(\mathbf{v}_m, b_i)$

Inner loop:

For $j = m, m-1, \dots, 2$

$$\mathbf{g} \leftarrow \left(\frac{d\boldsymbol{\sigma}_j(\mathbf{r}_j)}{d\mathbf{r}_j} \right)^T \mathbf{g}$$

$$\tilde{\nabla}_{\mathbf{W}_j} \phi_i(x) = \mathbf{g} \mathbf{v}_{j-1}^T$$

$$\nabla_{\beta_j} \phi_i(x) = \mathbf{g}$$

$$\mathbf{g} \leftarrow \mathbf{W}_j^T \mathbf{g}$$

DERIVATIVES CHAIN-RULE

$m-1$ matrix-vector products

where $\tilde{\nabla}_{\mathbf{W}_j} \phi(x) = \text{mat}(\nabla_{\text{vec}(\mathbf{W}_j)} \phi(x))$.

Binary classification problems - Logistic loss and derivative

The logistic loss with ℓ_2 regularization :

$$f(x) = \frac{1}{N} \sum_{i=1}^N \log(1 + e^{-b_i a_i^T x}) + \frac{1}{2N} \|x\|^2$$

Given \hat{x} resulting from the classifier training, $\hat{a} \in \mathbb{R}^n$

$$\frac{1}{1 + e^{-\hat{a}^T \hat{x}}} \geq 0.5 \quad \hat{b} = 1$$

$$\frac{1}{1 + e^{-\hat{a}^T \hat{x}}} < 0.5, \quad \hat{b} = -1$$

- Labels in $\{-1, 1\}$

$$\phi_i(x) = \log(1 + c(a_i, b_i; x)) + \frac{1}{2} \|x\|^2$$

where $c(a_i, b_i; x) = e^{-b_i a_i^T x}$

-

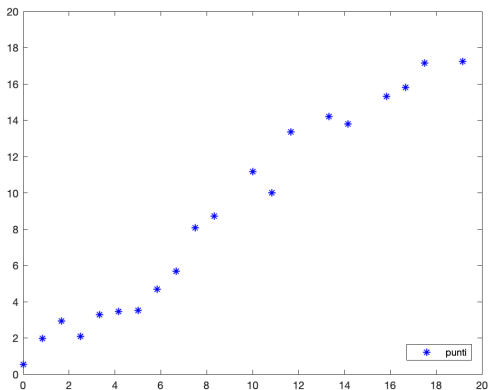
$$\nabla \phi_i(x) = -\frac{c(a_i, b_i; x)}{1 + c(a_i, b_i; x)} b_i a_i + x = \left(\frac{b_i}{1 + c(a_i, b_i; x)} - b_i \right) a_i + x$$

- **Epoch:** each set of N consecutive evaluations of a sample gradient $\nabla \phi_i(x)$.

Overfitting

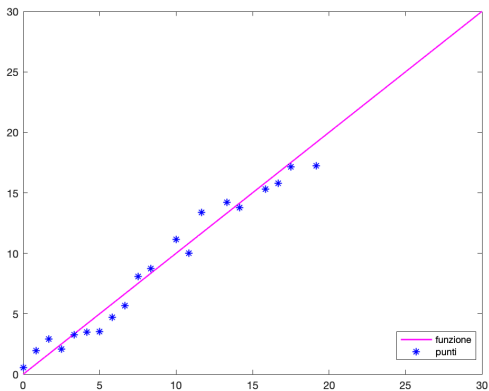
A regression problem

Dataset $(a_i, b_i) \in \mathbb{R} \times \mathbb{R}$, $N = 20$, b_i are values affected by noise



Overfitting

Linear model -square loss

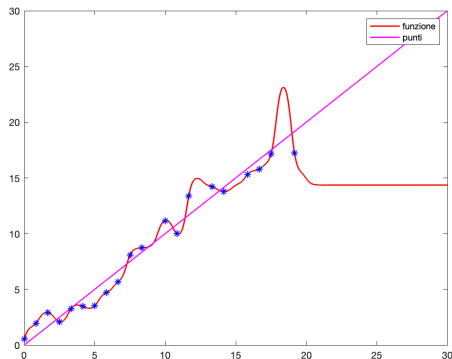


Overfitting

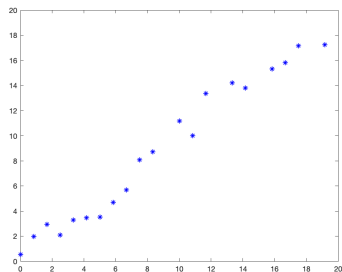
Neural Network-square loss

1 - 10 - 1
input hidden output

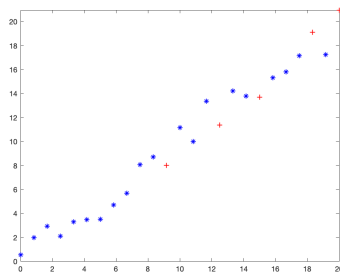
20 weights, 11 bias, 31 parameters.



Overfitting - Validation set



(a) Training set \mathcal{T}

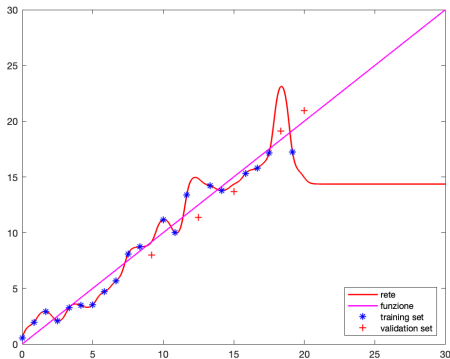


(b) Training set \mathcal{T} + Validation set \mathcal{V}

The data-set is split in

- training set (blue)
- validation set (red)
- testing set

Overfitting - training error and validation error

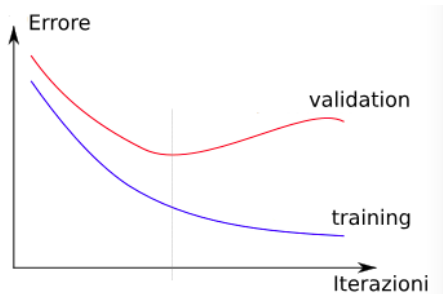


Errors

- $err_{\mathcal{T}}(x) = \frac{1}{|\mathcal{T}|} \sum_{(a_j, b_j) \in \mathcal{T}} (\mathbf{v}_3(x; a_j) - b_j)^2$
- $err_{\mathcal{V}}(x) = \frac{1}{|\mathcal{V}|} \sum_{(a_j, b_j) \in \mathcal{V}} (\mathbf{v}_3(x; a_j) - b_j)^2$

Overfitting - stopping criterion

- $err_{\mathcal{T}}(x) = \frac{1}{|\mathcal{T}|} \sum_{(a_j, b_j) \in \mathcal{T}} (\mathbf{v}_3(x; a_j) - b_j)^2$ (in blu)
- $err_{\mathcal{V}}(x) = \frac{1}{|\mathcal{V}|} \sum_{(a_j, b_j) \in \mathcal{V}} (\mathbf{v}_3(x; a_j) - b_j)^2$ (in rosso)



Stopping criterion: Training is stopped when the validation error increases.