The background features abstract geometric shapes in red, grey, and white. A large circle on the left contains the word 'Insper' in a serif font. Several curved lines of varying lengths and colors (red, grey, white) radiate from the top right corner across the slide.

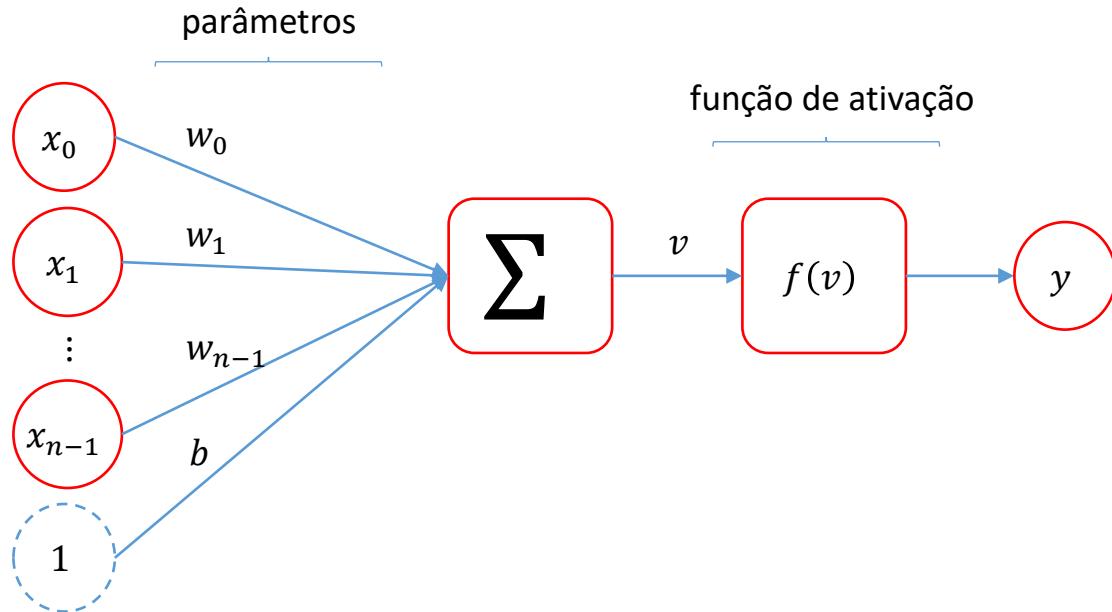
Insper

Machine Learning

Aula 21 – Treinando redes neurais

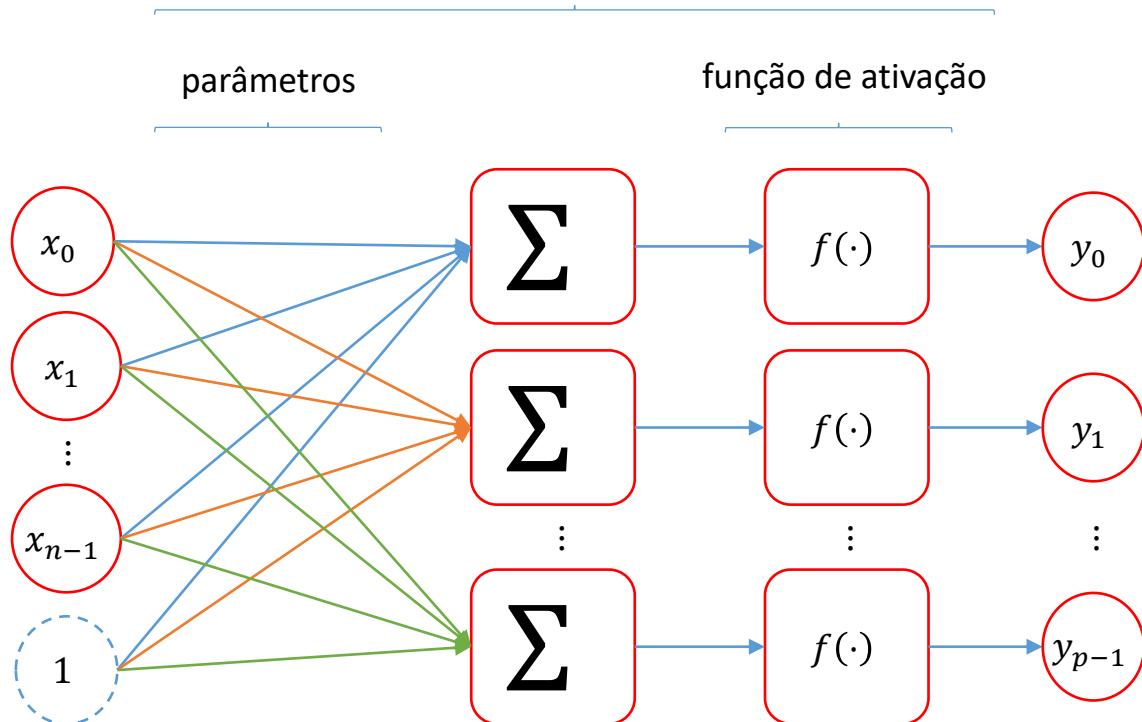
2021 – Engenharia
Fábio Ayres <fabioja@insper.edu.br>

Neurônio artificial



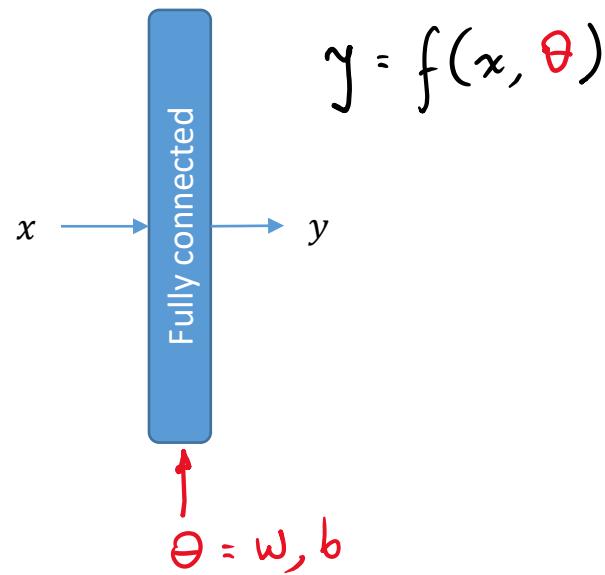
Camada (layer) de neurônios

“Fully-connected layer” \rightarrow Dense

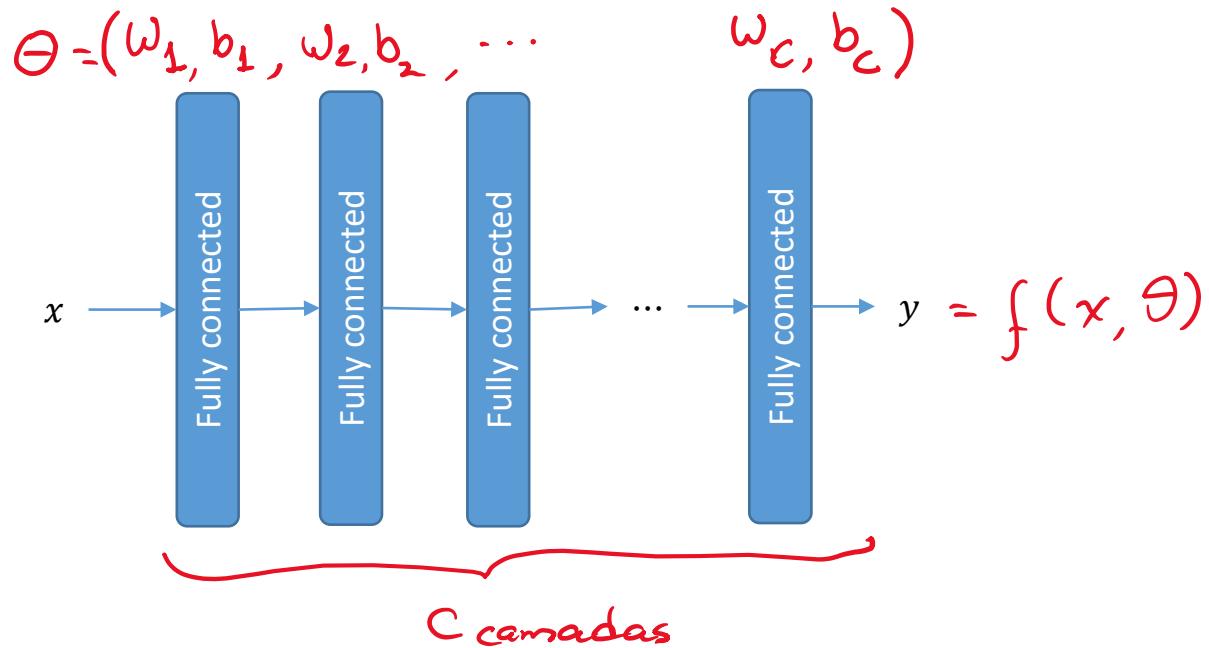


Camada de neuronios

- Geralmente representada de modo simples:



Rede neural multicamadas



Modelo: $y = f(x, \theta)$

↑ parametros treináveis

Dados: $X_{\text{train}}, y_{\text{train}}$

Função de perda: $L(\theta, X_{\text{train}}, y_{\text{train}}) \rightarrow L(\theta)$

fixos
↓
a verdadeira
variável desse
problema

Algoritmo de otimização

$$\theta_{\text{opt}} = \arg \min_{\theta} L(\theta)$$



Gradient descent

$$\theta^{(i+1)} = \theta^{(i)} - \eta \nabla_{\theta} L(\theta^{(i)})$$

↑
learning rate

Problemas com as redes neurais

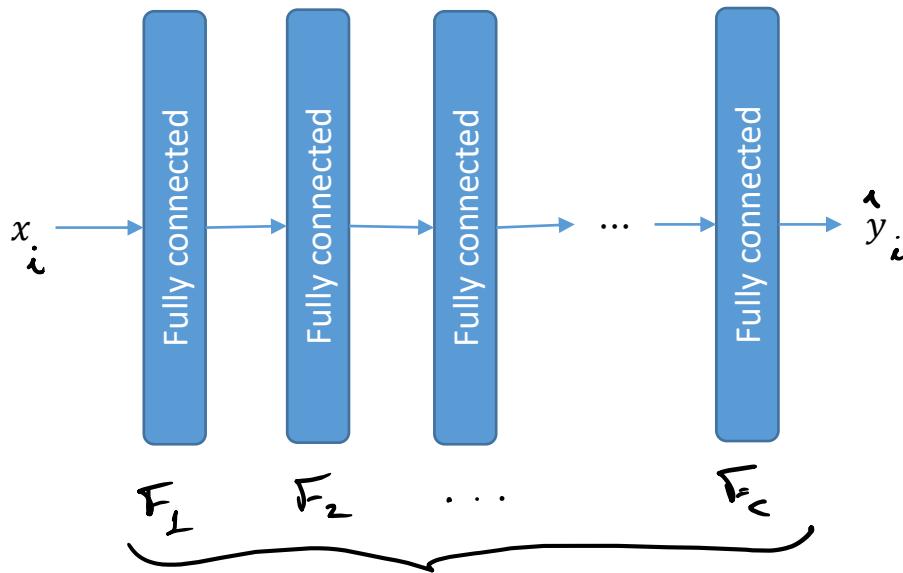
- Difícil de fazer convergir
- Overfitting
- Demora muito para treinar

Estes são os desafios que acabaram com os esforços de estudar redes neurais nos anos 90, inicio dos anos 2000.

Nos anos 2010 esses desafios foram enfrentados de novo, com sucesso!

Convergência

- Vanishing/exploding gradients



$$\hat{y} = (F_c \circ F_{c-1} \circ \dots \circ F_1)(x, \theta_1, \theta_2, \dots, \theta_c)$$

$$\hat{y} = (F_c \circ F_{c-1} \circ \dots \circ F_1)(x, \theta_1, \theta_2, \dots, \theta_c)$$

$l(y, \hat{y})$: perda por amostra

Exemplo: $l(y, \hat{y}) = (y - \hat{y})^2$ no m.s.e

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

na entropia cruzada

$$L(\tilde{y}, \hat{\tilde{y}}) = \frac{1}{m} \sum_{i=1}^m l(y_i, \hat{y}_i)$$

$$\frac{\partial L}{\partial \theta_1} = \frac{1}{m} \sum_{i=1}^m \frac{\partial l}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial \theta}$$

$$\frac{\partial \hat{y}_i}{\partial \theta} = \frac{\partial}{\partial \theta} F_c(F_{c-1}(F_{c-2}(\dots F_1(x_i, \theta_1) \dots), \theta_{c-1}), \theta_c)$$

$$\frac{\partial \hat{y}_i}{\partial \theta_L} = \frac{\partial}{\partial \theta_L} F_c(F_{c-1}(F_{c-2}(\dots F_1(x_i, \theta_1) \dots), \theta_{c-1}), \theta_c)$$
10

$$\frac{\partial \hat{y}_i}{\partial \theta_L} = \frac{\partial F_c}{\partial F_{c-1}} \cdot \frac{\partial F_{c-1}}{\partial F_{c-2}} \cdot \dots \cdot \frac{\partial F_2}{\partial F_1} \cdot \frac{\partial F_1}{\partial \theta_L}$$

— “ —

$$E \approx \left| \frac{\partial F_k}{\partial F_{k-1}} \right| \approx 0 ?$$

I Num produto, se um termo é zero, o produto é zero!

$$\hat{\theta}^{(i+1)} = \theta^{(i)} - \eta \nabla_{\theta} L(\theta^{(i)})$$

$$\text{II Logo, } |\nabla_{\theta}| \approx 0 \Rightarrow \hat{\theta}^{(i+1)} \approx \theta^{(i)}$$

a rede parou de evoluir!

ser muito leve (massa ≈ 0)
descendo na gravidade
apenas

θ

!!!
otro

WHEEE!

$L(\theta)$

inclinacões: $\nabla_{\theta} L(\theta)$

tô devagar!
Praticamente
parei!

quase-plateau!

cadê ele?

!
otro

$$\frac{\partial \hat{y}_i}{\partial \theta_L} = \frac{\partial}{\partial \theta_L} F_c(F_{c-1}(F_{c-2}(\dots F_1(x_i, \theta_1) \dots), \theta_{c-1}), \theta_c)$$
(12)

$$\frac{\partial \hat{y}_i}{\partial \theta_L} = \frac{\partial F_c}{\partial F_{c-1}} \cdot \frac{\partial F_{c-1}}{\partial F_{c-2}} \cdot \dots \cdot \frac{\partial F_2}{\partial F_1} \cdot \frac{\partial F_1}{\partial \theta_L}$$

— “ —

$$E \approx \left| \frac{\partial F_k}{\partial F_{k-1}} \right| \rightarrow +\infty ?$$

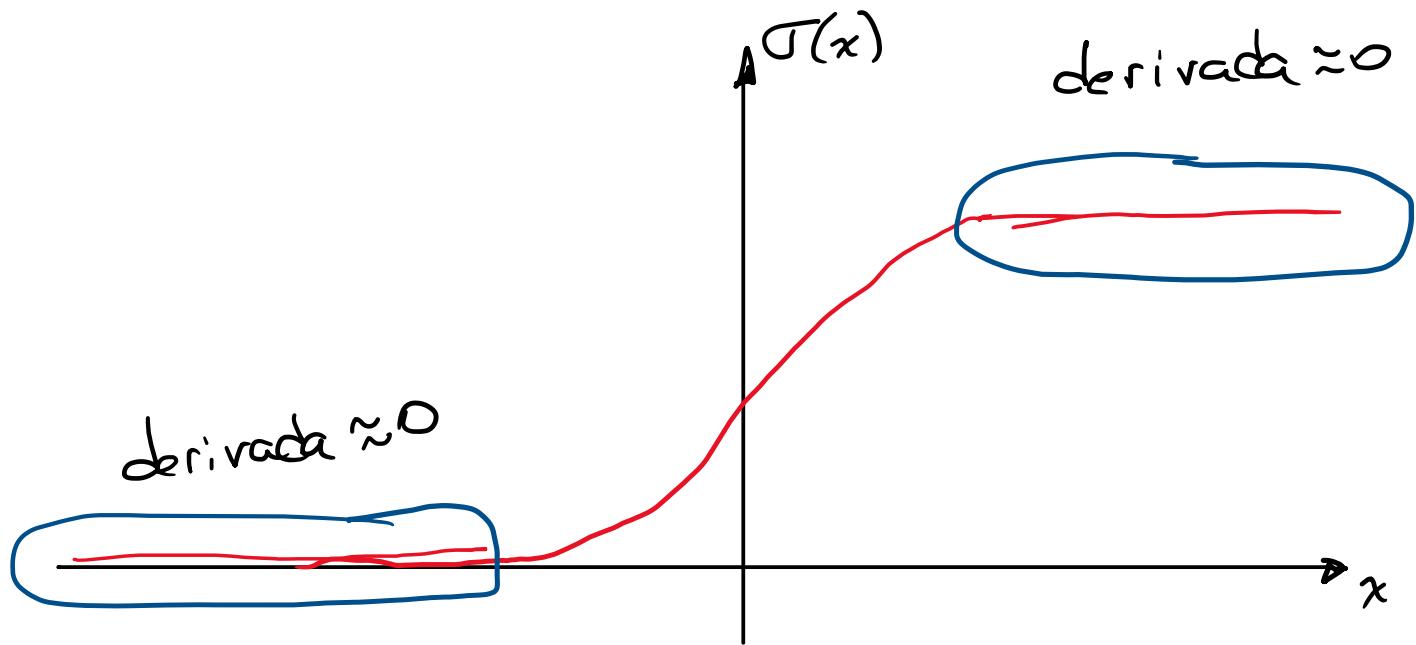
Num produto, se um termo c' too a magnitude do

$$\tilde{\theta}^{(i+1)} = \theta^{(i)} - \eta \frac{\nabla L(\theta^{(i)})}{\epsilon}$$

passo exagerado \Rightarrow rede diverge!

Convergência

- Vanishing/exploding gradients



Convergência

- Solução: novas funções de ativação

“Família -lu”

↑
também com SVM!

• relu : Rectified Linear Unit

• elu : Exponential Linear Unit

• selu

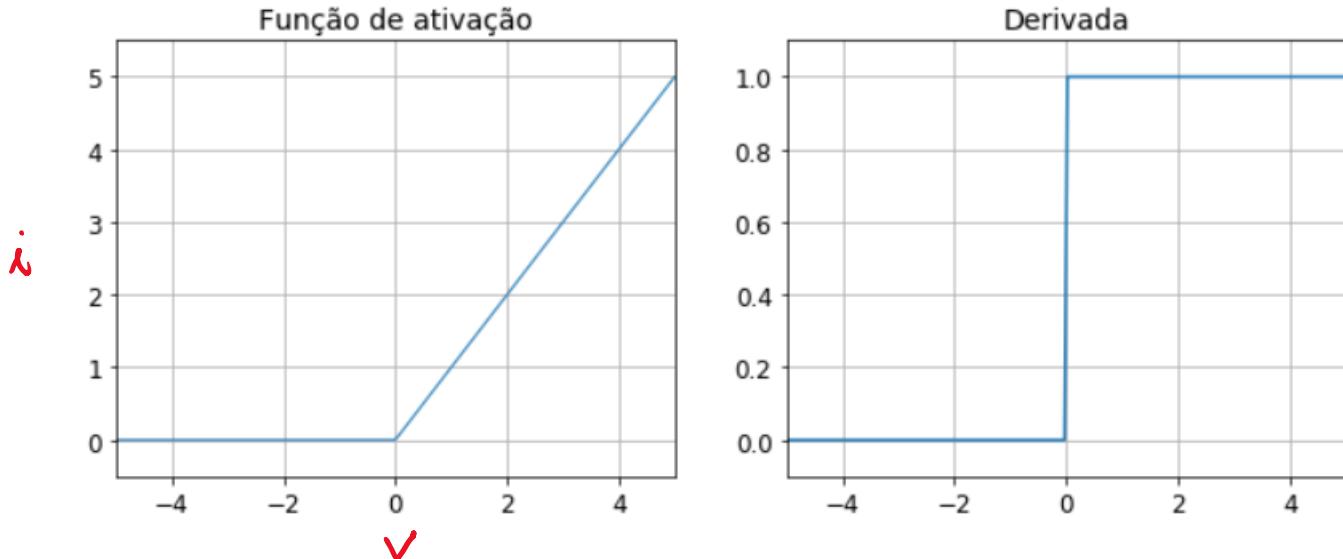
• gelu

• prelu

• etc

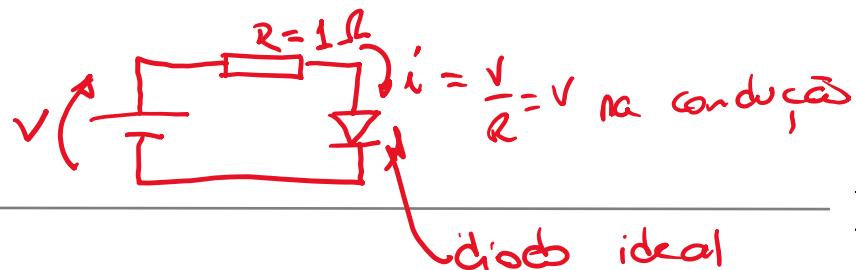
- simples de calcular
- Pelo menos em metade do domínio a deriv. não se anula.

Rectified linear unit

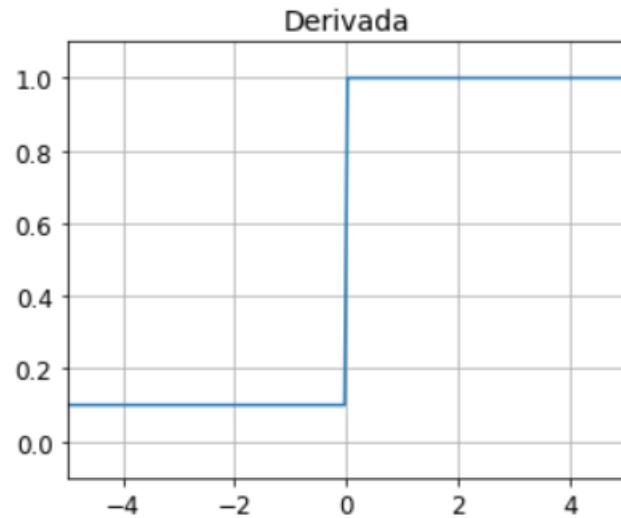
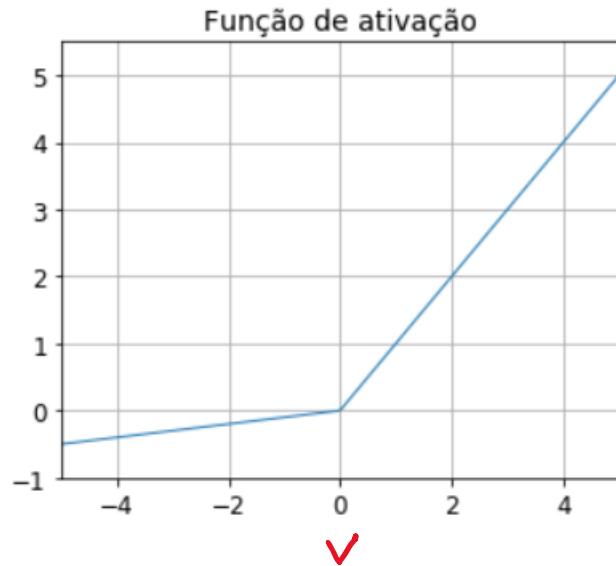


i

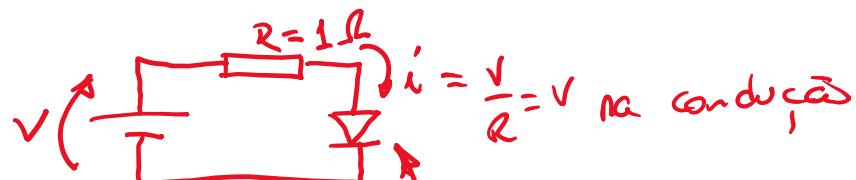
✓



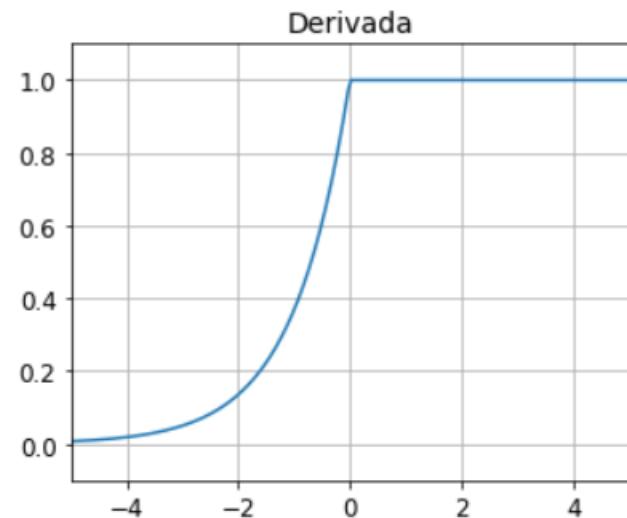
Leaky relu



✓



Exponential linear unit (elu)



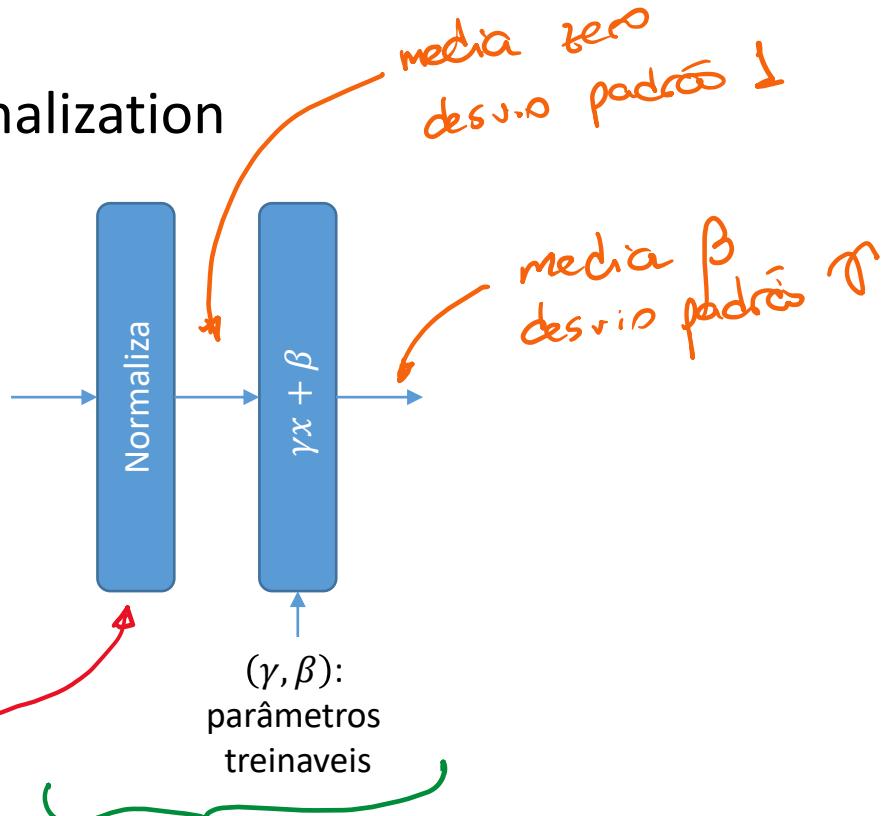
Convergência

- Solução: Batch normalization

treinamento: normalizar
teste ou uso: \bar{x}, s fixos

Para cada batch:

$$z_i = \frac{x_i - \bar{x}_B}{S_B}$$

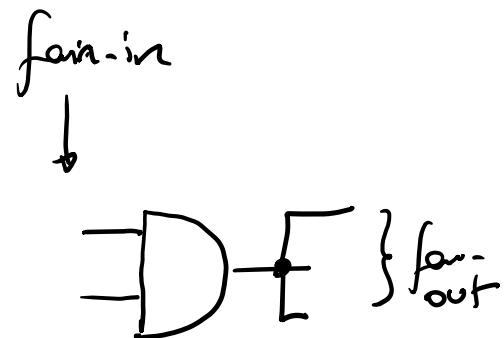
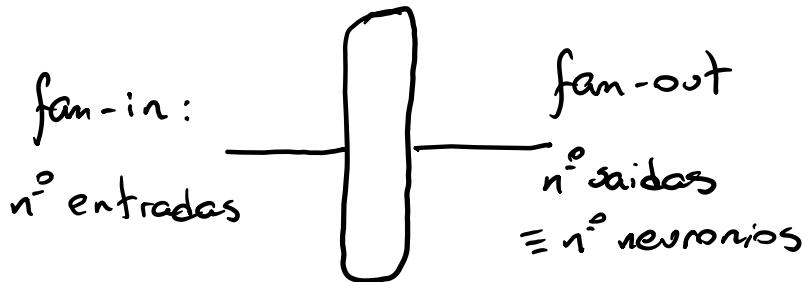


Convergência

- Solução: melhor inicialização dos pesos

→ Como definir $\Theta^{(0)}$?

Antigamente: $\Theta^{(0)} \sim N(0, 1)$



Convergência

- Solução: gradient clipping

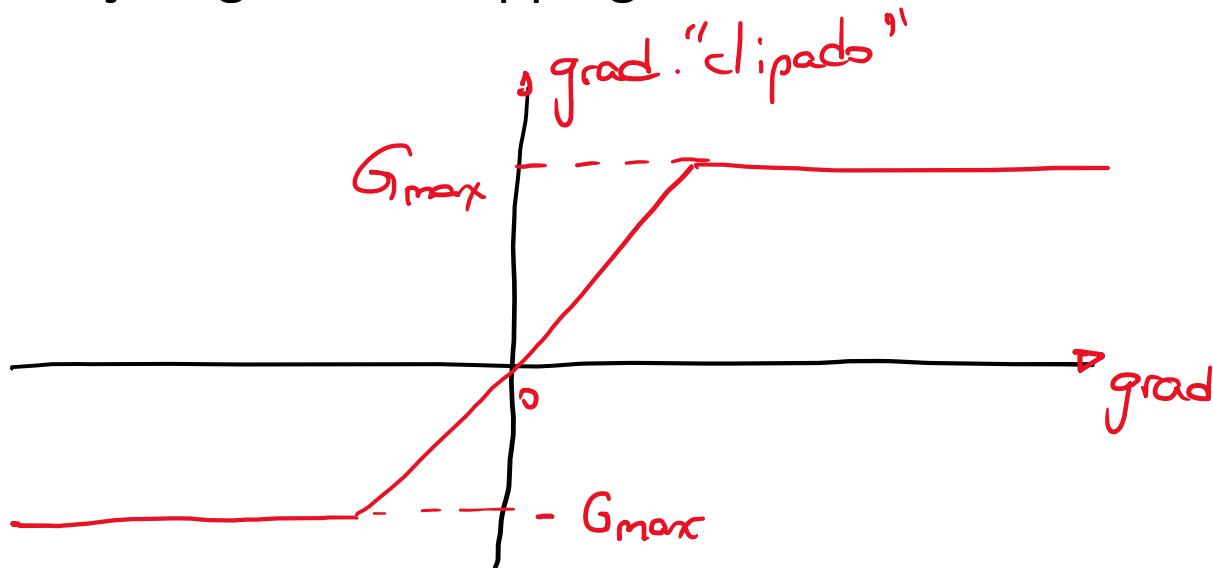
$$\text{fan-avg} = \frac{\text{fan-int} + \text{fan-out}}{2}$$

Glorot e Bengio

$$\text{pesos} \sim N(0, \frac{1}{\text{fan-avg}})$$

Convergência

- Solução: gradient clipping



Overfitting

- Solução: regularização
 - Cada camada de rede neural não é uma regressão logística metida a besta? Então tasca Lasso e Ridge nelas!

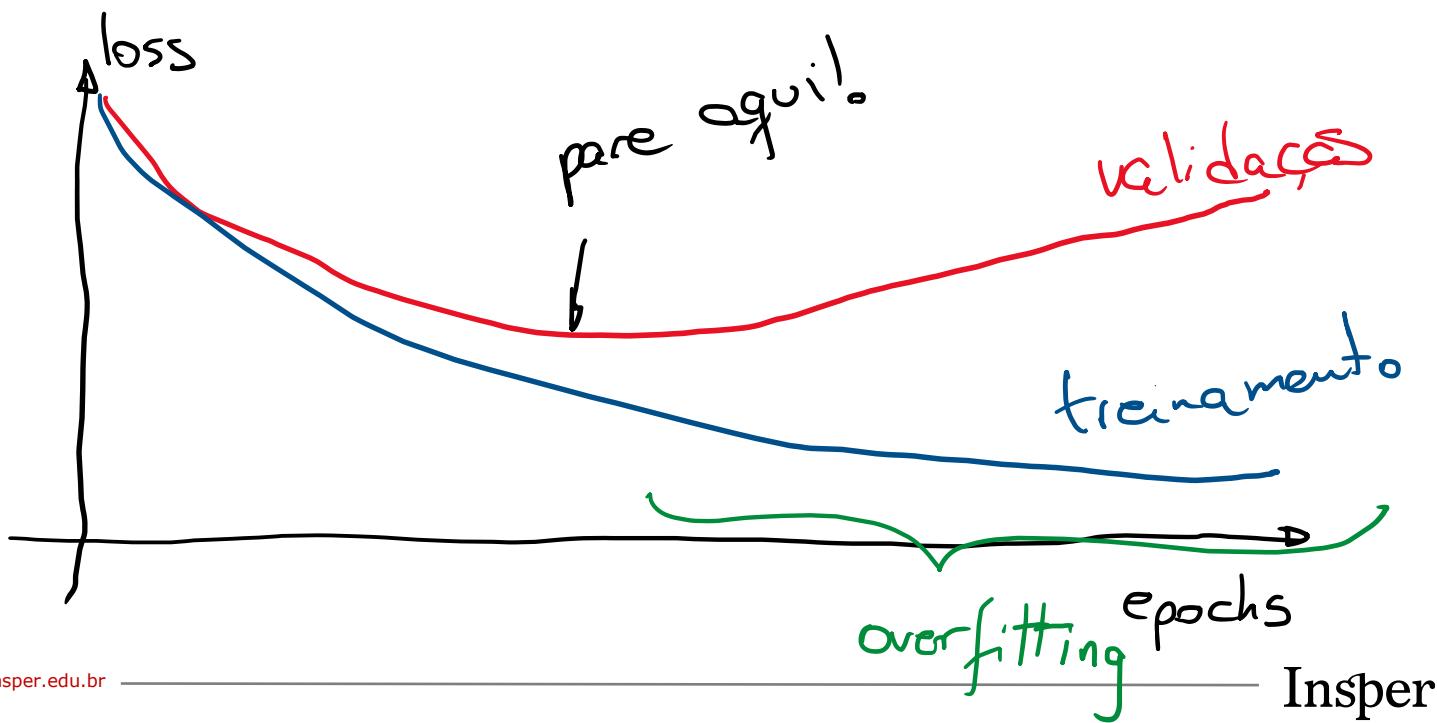
$$\ell_1 : \sum |\omega_i|$$

$$\ell_2 = \sum \omega_i^2$$

Joga como penalidade extra

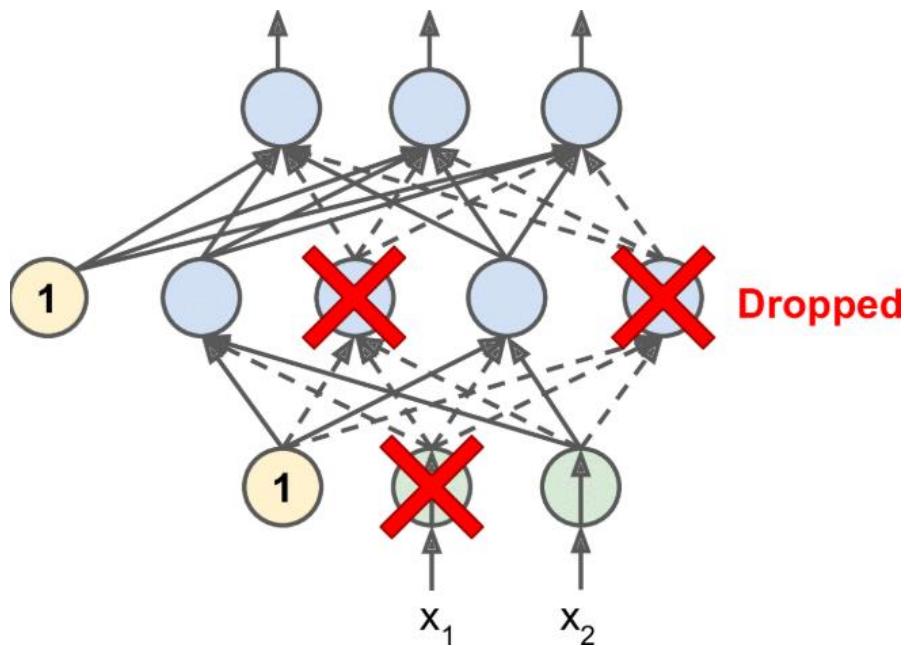
Overfitting

- Solução: Early stopping



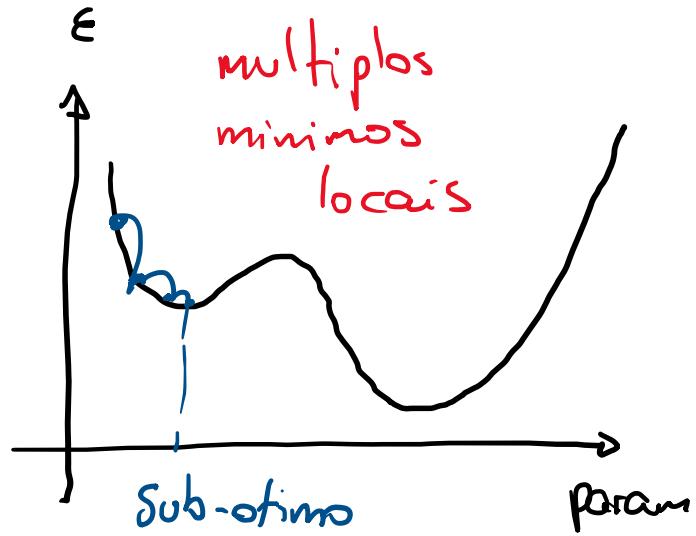
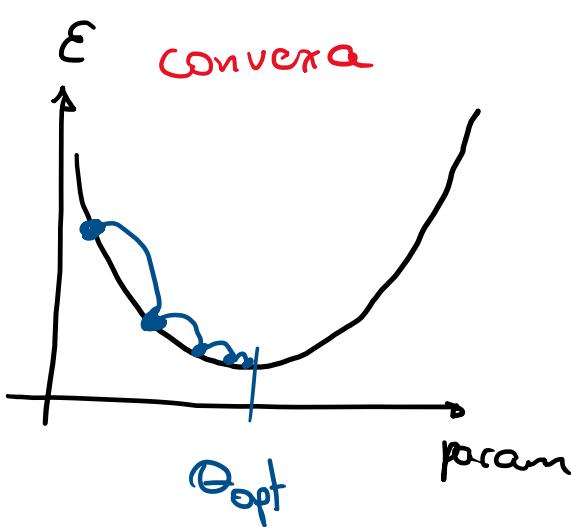
Overfitting

- Solução: Dropout!



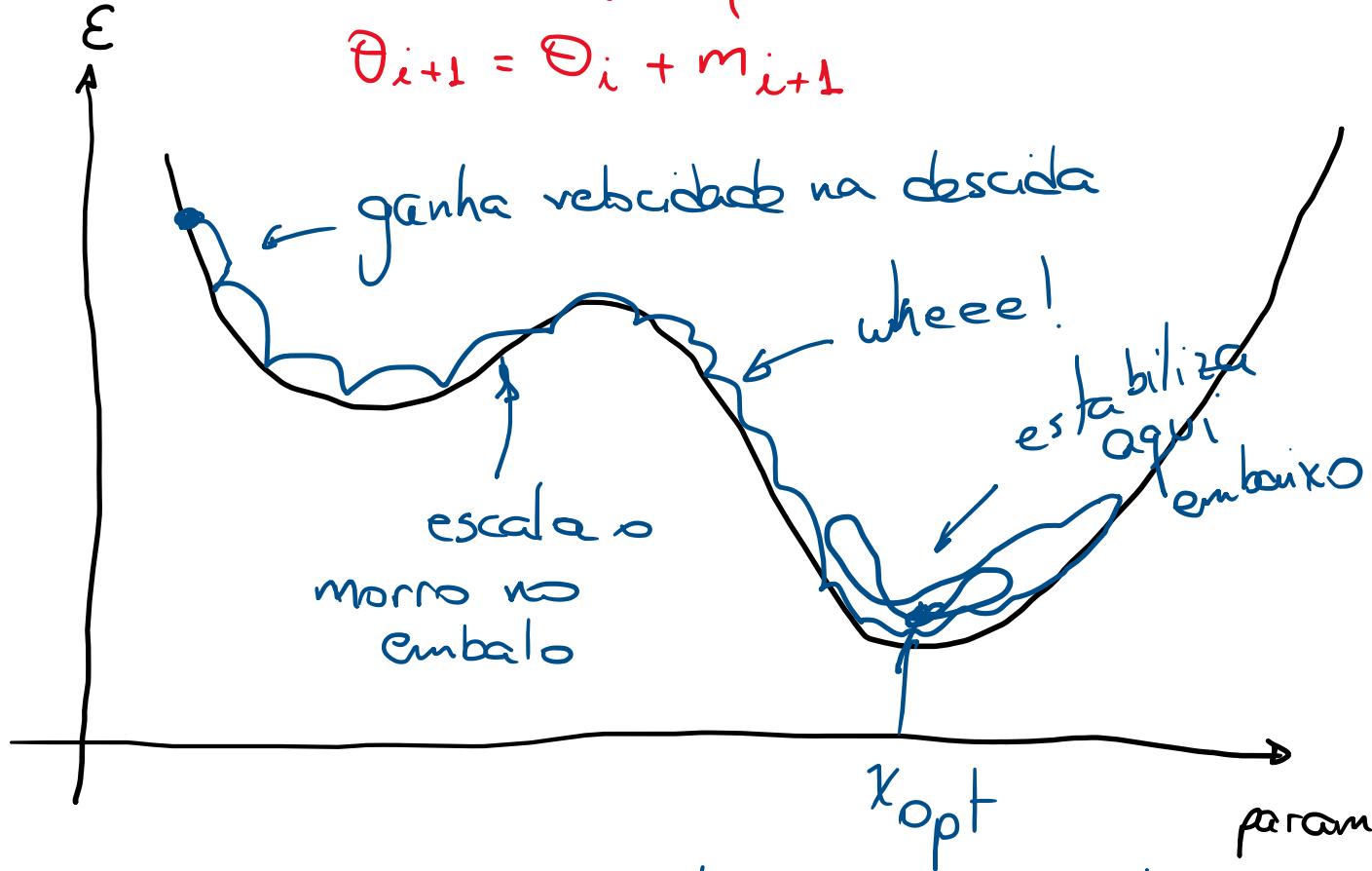
Treinamento

- Solução: melhores otimizadores
 - Melhora também o problema dos ótimos locais da função de erro médio



$$m_{i+1} = m_i - \eta \nabla \mathcal{E}(\theta_i)$$

$$\theta_{i+1} = \theta_i + m_{i+1}$$

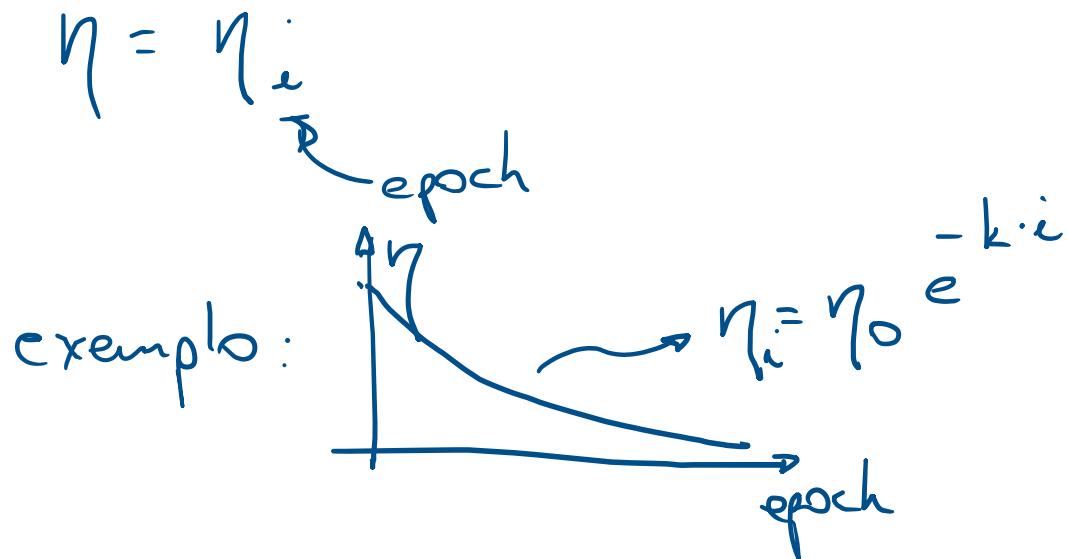


Gradient descent com momento

Insp^{er}
massa x velocidade

Treinamento

- Solução: mexer na taxa de aprendizado







Insper