

## Aprendizado do MLP por Error Back Propagation ...

$$\Delta \vec{W} = -\eta \cdot \vec{\nabla} E_{qm}$$

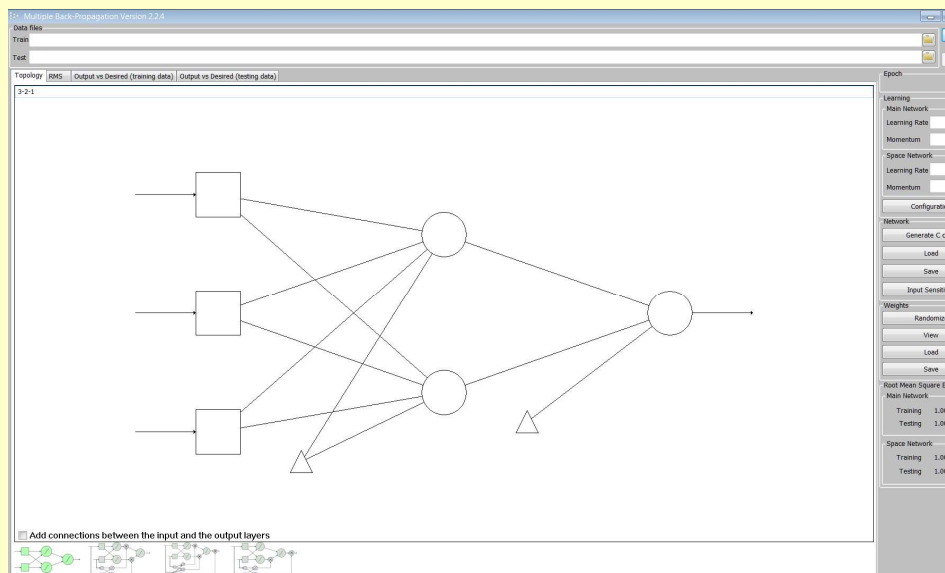
Gradiente de Eqm no espaço de pesos =  $(\partial E_{qm}(W)/\partial w_1, \partial E_{qm}(W)/\partial w_2, \partial E_{qm}(W)/\partial w_3, \dots)$

**Chegando às fórmulas das derivadas parciais, necessárias à Bússola do Gradiente**

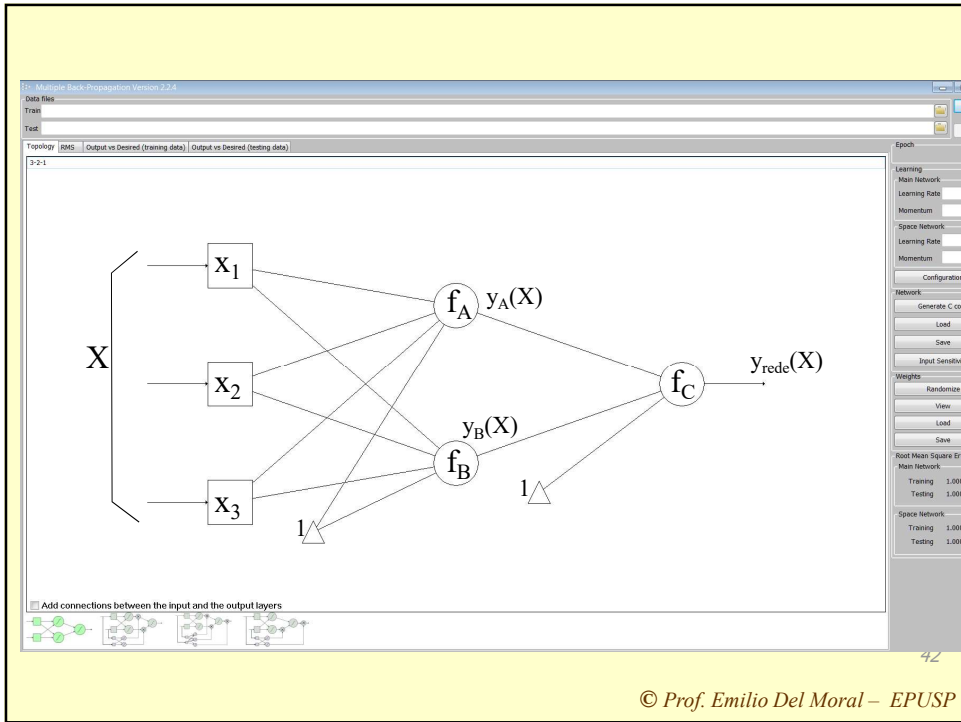
40

© Prof. Emilio Del Moral – EPUSP

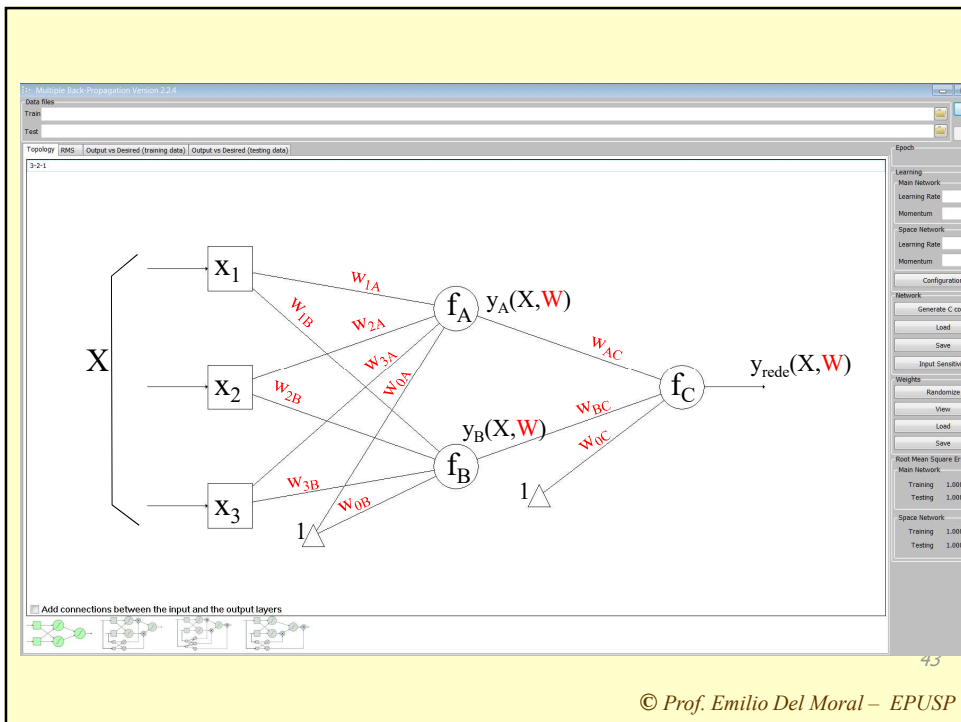
*Relembrando o que está por trás de um desenho como o que segue ...*



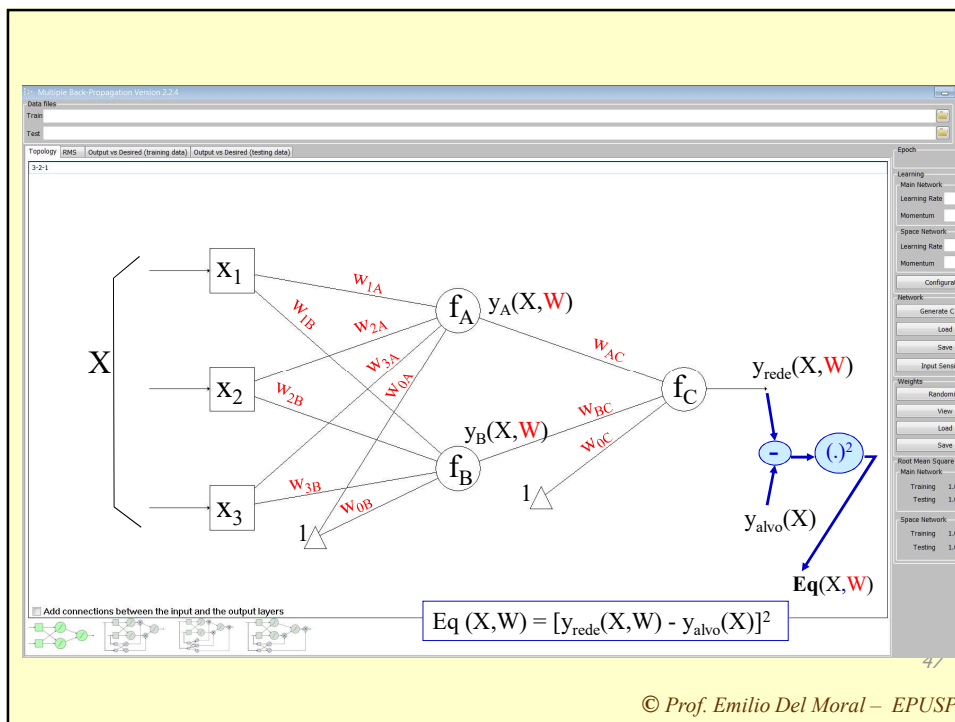
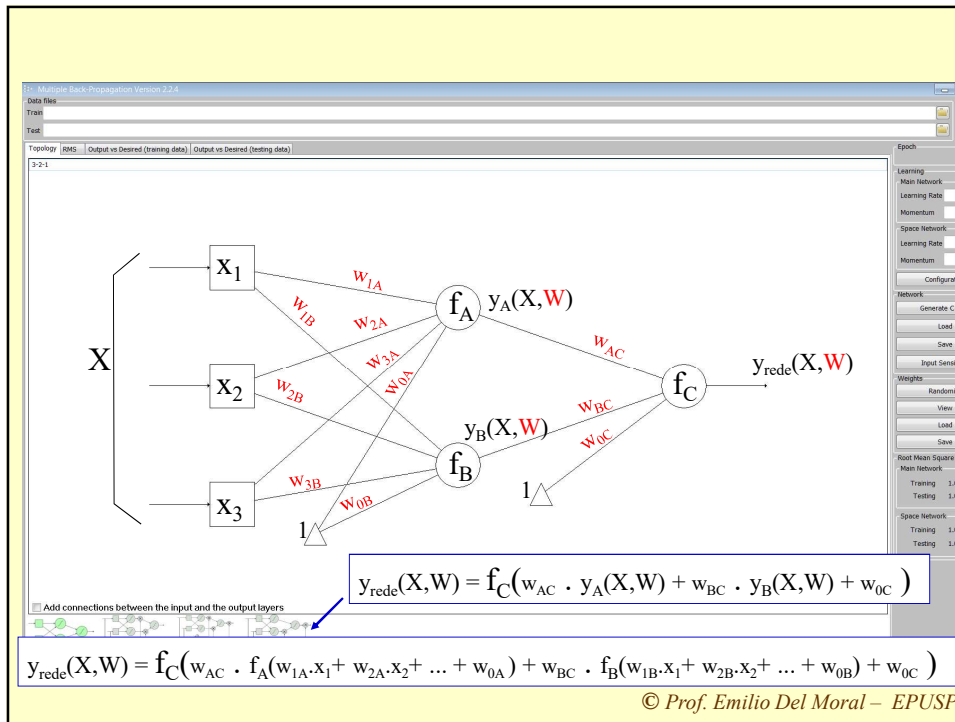
© Prof. Emilio Del Moral – EPUSP



© Prof. Emilio Del Moral – EPUSP



© Prof. Emilio Del Moral – EPUSP



Chamada oral sobre a lição de casa: estudar / reestudar os conceitos e a parte operacional de derivadas parciais, do vetor Gradiente ...

- Derivadas parciais (que são as componentes do gradiente):

$$\partial f(a,b,c)/\partial a \quad \partial f(a,b,c)/\partial b \quad \partial f(a,b,c)/\partial c$$

- Vetor Gradiente, útil ao método do máximo declive:

$$(\partial \text{Eqm}(W)/\partial w_1, \partial \text{Eqm}(W)/\partial w_2, \partial \text{Eqm}(W)/\partial w_3, \dots)$$

$$\vec{\Delta W} = -\eta \cdot \vec{\nabla} \text{Eqm}_-$$

49

© Prof. Emilio Del Moral – EPUSP

## Invertamos o operador gradiente e a somatória

.. afinal, gradiente é uma derivada, e a derivada de um soma de várias funções é igual à soma das derivadas individuais de cada componente da soma:

$$\text{Grad}(\text{Eqm}) =$$

$$\text{Grad}(\sum_{\mu} \text{Eq}^{\mu}) / M$$

$$\sum_{\mu} \text{Grad}(\text{Eq}^{\mu}) / M$$

50

© Prof. Emilio Del Moral – EPUSP

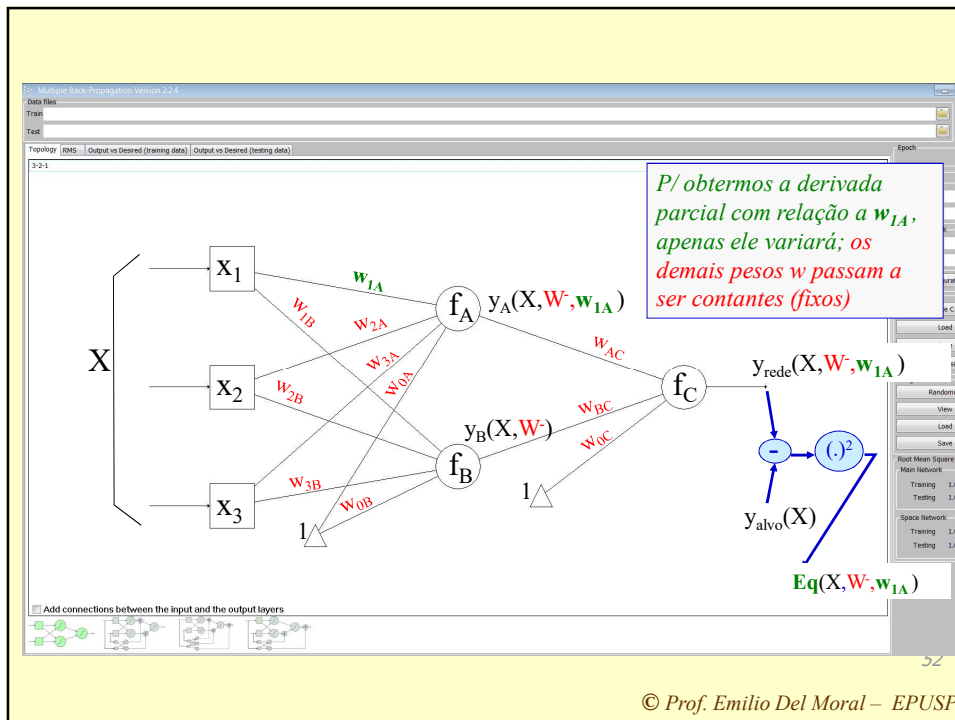
Note que a inversão do gradiente com a somatória nada mais é que usar de forma repetida – e em separado para cada dimensão

do vetor  $\mathbf{Grad}(\sum_{\mu} Eq^{\mu})$  – a seguinte propriedade simples e sua velha conhecida ...

$$d(f_1(x)+f_2(x)) / dx = df_1(x)/dx + df_2(x)/dx$$

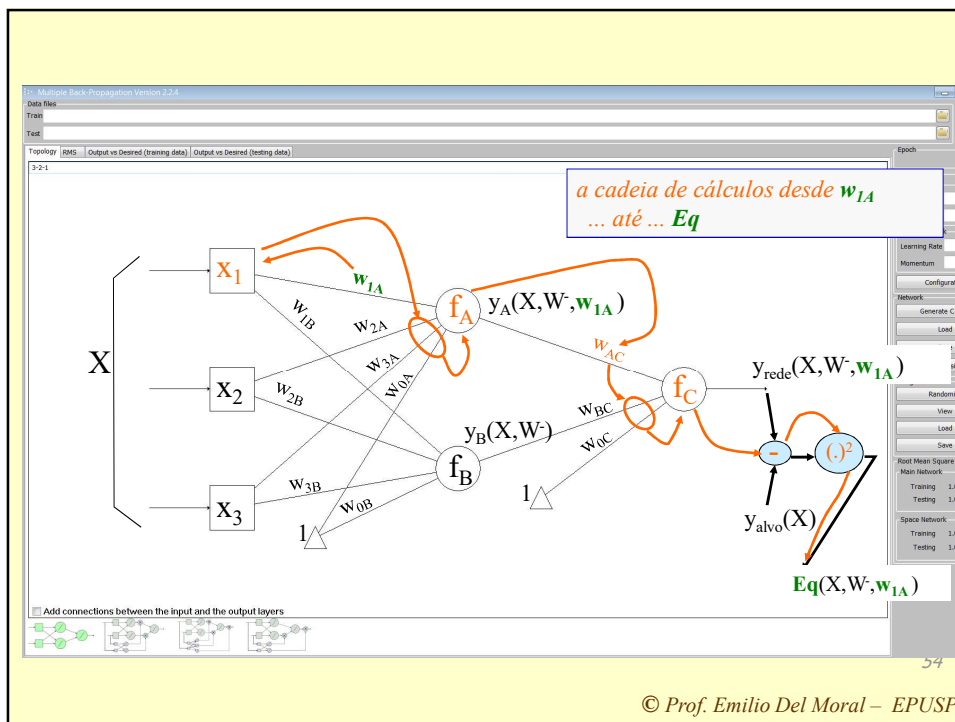
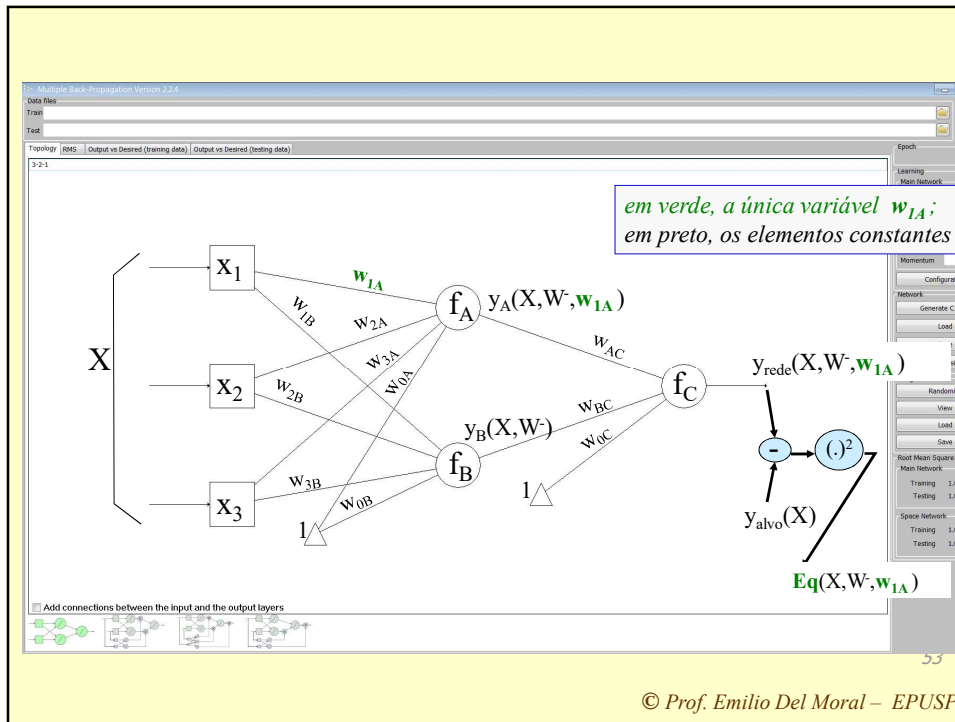
51

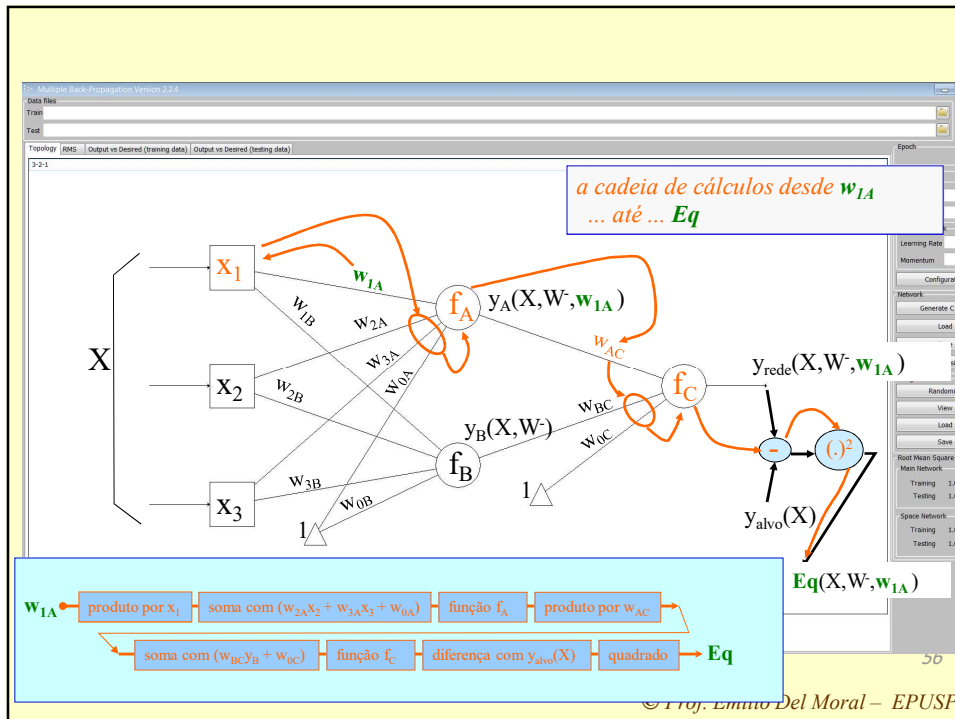
© Prof. Emilio Del Moral – EPUSP



52

© Prof. Emilio Del Moral – EPUSP

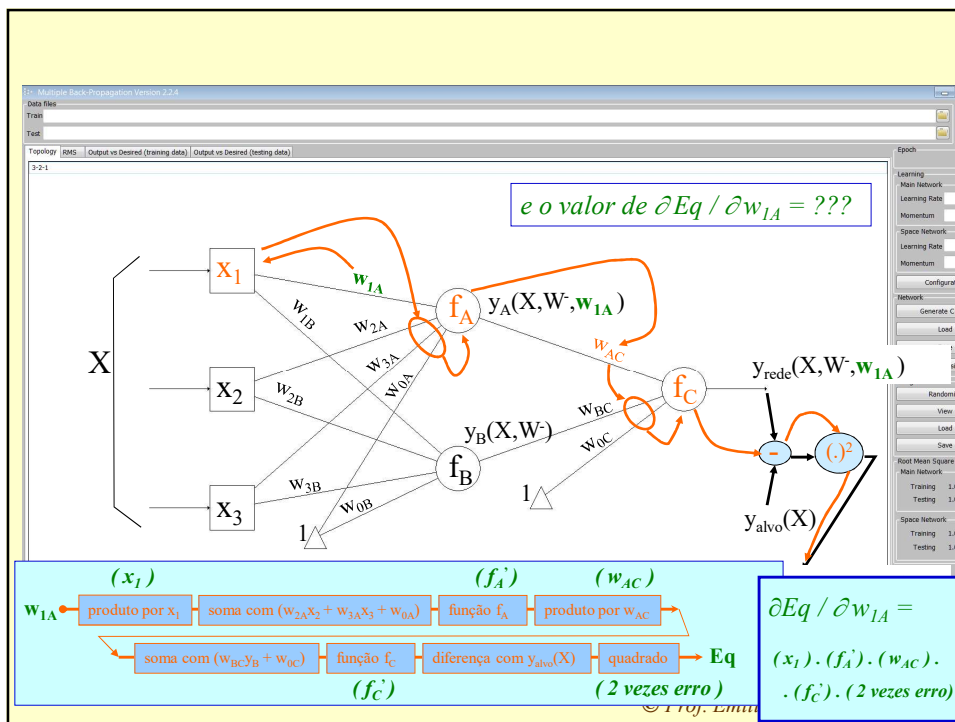
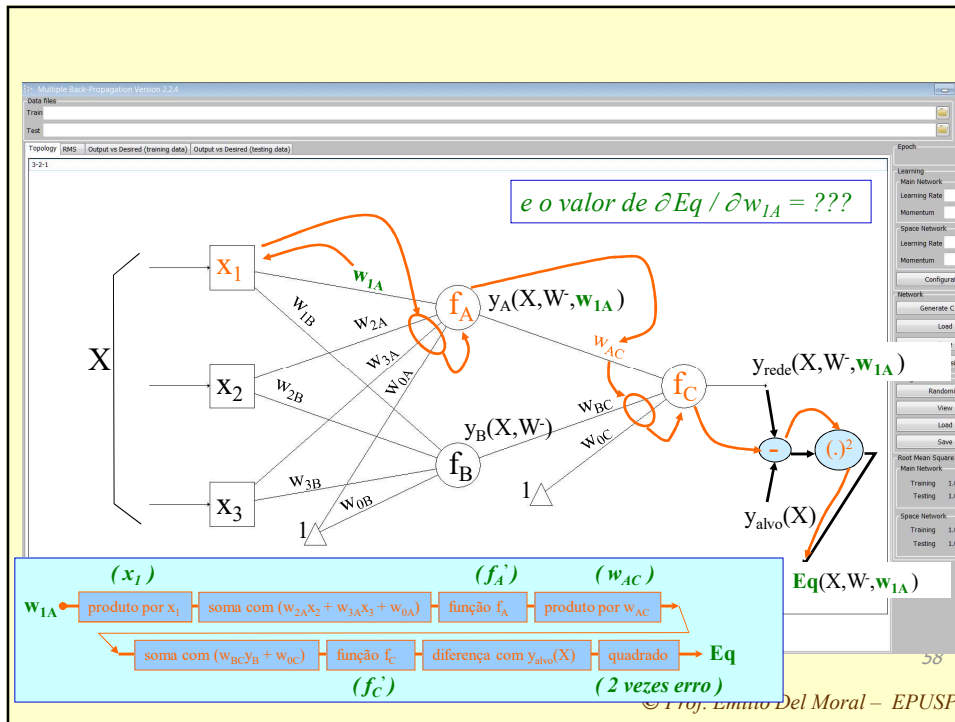




Note que aqui temos uma cadeia com muitos estágios que levam da variável  $w_{1A}$ , à variável  $Eq^u$ , e para a qual podemos calcular a derivada da saída ( $Eq^u$ ) com relação à entrada ( $w_{1A}$ ) aplicando de forma repetida a seguinte propriedade simples e sua velha conhecida ...

$$d(f_1(f_2(x))) / dx = df_1(x)/df_2 \cdot df_2(x)/dx$$

..., ou seja, calculando isoladamente o valor da derivada para cada estágio da cadeia, e finalizando o cálculo de derivada de ponta a ponta nessa cadeia toda através do produto dos diversos valores de cada estágio.





## Lembretes ....

- Na maioria dos slides anteriores, onde aparece  $X$ , leia-se  $X^\mu$ , não incluído para não complicar demais os desenhos
- ... similarmente, onde aparece  $y_{\text{alvo}}$ , leia-se  $y_{\text{alvo}}^\mu$ . Idem para os Eq, leia-se Eq $^\mu$
- Nos itens de cadeia de derivadas ( $f_A'$ ) e ( $f_C'$ ), atenção para os valores dos argumentos, que devem ser os mesmos de  $f_A$  e  $f_C$  na cadeia original que leva  $w_{IA}$  a Eq.
- ... lembrando ... na cadeia original tínhamos ...
  - para  $f_C$ :  $f_C(w_{AC} \cdot f_A(w_{1A} \cdot x_1 + w_{2A} \cdot x_2 + w_{0A}) + w_{BC} \cdot f_B(w_{1B} \cdot x_1 + w_{2B} \cdot x_2 + w_{0B}) + w_{0C})$
  - para  $f_A$ :  $f_A(w_{1A} \cdot x_1 + w_{2A} \cdot x_2 + w_{0A})$
- Similarmente, para o bloco “quadrado”, cuja derivada é a função “2 vezes erro”, o argumento é  $[y_{\text{rede}}(X, W) - y_{\text{alvo}}(X)]$

61

© Prof. Emilio Del Moral – EPUSP

## Lembretes ....

- O mesmo que foi feito para  $w_{IA}$  deve ser feito agora para os demais 10 pesos:  $w_{2A}$ ,  $w_{3A}$ ,  $w_{0A}$ ,  $w_{1B}$ ,  $w_{2B}$ ,  $w_{3B}$ ,  $w_{0B}$ ,  $w_{AC}$ ,  $w_{BC}$ , e  $w_{0C}$  !
- Assim compomos um gradiente de 11 dimensões, com as derivadas de Eq $^\mu$  com relação aos 11 diferentes pesos  $w$ :  $\text{Grad}_w(\text{Eq}^\mu)$
- Essas 11 fórmulas devem ser aplicadas repetidamente aos  $M$  exemplares numéricos de  $X^\mu$  e  $y_{\text{alvo}}^\mu$ , calculando  $M$  gradientes!
- Com eles, se obtém o gradiente médio dos  $M$  pares empíricos:  $\text{Grad}_w(\text{Eqm}) = [\sum_\mu \text{Grad}_w(\text{Eq}^\mu)] / M$
- Esse gradiente médio é a Bussola do Gradiente!

62

© Prof. Emilio Del Moral – EPUSP

Método do Gradiente Aplicado aos nossos MLPs: a partir de um  $W \neq 0$ , temos aproximações sucessivas ao Eqm mínimo, por repetidos pequenos passos  $\Delta W$ , sempre contrários ao gradiente ...

- “Chute” um  $W$  inicial para o “ $W$  corrente”, ou “ $W$  melhor até agora”
- Em loop até obter Eqm zero, ou baixo o suficiente, ou estável:
  - Determine o vetor gradiente do Eqm, nesse espaço de  $W$ s
  - Em loop varrendo todos os  $M$  exemplos  $(X^\mu; y^\mu)$ ,
    - Calcule o gradiente de  $Eq^\mu$  associado a um exemplo  $\mu$ , e vá varrendo  $\mu$  e somando os gradientes de cada  $Eq^\mu$ , para compor o vetor gradiente de Eqm, assim que sair deste loop em  $\mu$  ;
    - Cada cálculo como esse, envolve primeiro calcular os argumentos de cada tangente hiperbólica e depois usar esses argumentos na regra da cadeia das derivadas necessárias
  - Dê um passo Delta  $\Delta W$  nesse espaço, com direção e magnitude dados por  $-\eta \cdot$ vetor gradiente médio para os  $M$  Exemplos  $(X^\mu; y^\mu)$  de treino

63

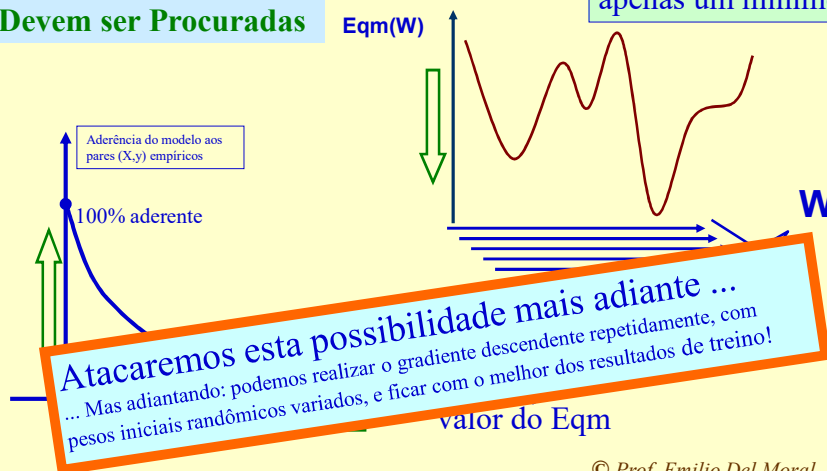
© Prof. Emilio Del Moral – EPUSP

O que devemos mirar quando exploramos o espaço de pesos  $W$  buscando que a RNA seja um bom modelo?

*Devemos mirar Maximização da aderência = Mínimo Eqm possível*

**As Setas Verdes Indicam Situações que Devem ser Procuradas**

Será que temos apenas um mínimo??



64

© Prof. Emilio Del Moral – EPUSP