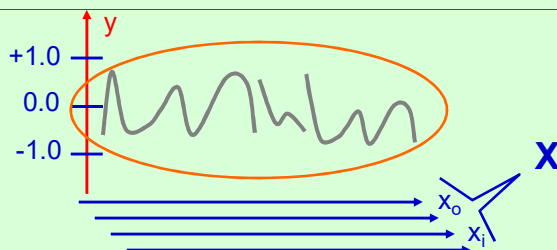


Aprendizado em RNAs do tipo MLP – Multi Layer Perceptron – através do algoritmo Error Back Propagation

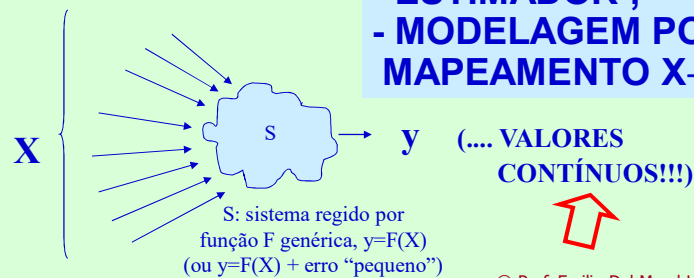
© Prof. Emilio Del Moral – EPUSP

1

A função $y(X)$ “a descobrir”, num caso geral de função contínua $y(X)$



- ESTIMADOR ;
- MODELAGEM POR
MAPEAMENTO $X \rightarrow y$



© Prof. Emilio Del Moral Hernandez

2

Conjunto de treino em arquiteturas supervisionadas (ex. clássico: MLP com Error Back Propagation)

Sistema Físico, Econômico, Biológico ...

X → s → y

Conjunto de M Amostras ($X^\mu; y^\mu$)

X → **RNA** → y_{rede}

A computação desejada da rede pode ser definida simplesmente através de amostras / exemplos do comportamento requerido

$$Eqm = \frac{1}{M} \sum_{\mu=1}^M (y_{rede}(\vec{X}^\mu, \vec{W}) - y^\mu)^2$$

... em loop ...

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

Aprendizado: Espaço de pesos W é explorado visando aproximar ao máximo a computação da rede da computação desejada

© Prof. Emilio Del Moral Hernandez

3

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

Devemos buscar Maximização da aderência = Mínimo Eqm possível

As Setas Verdes Indicam Situações que Devem ser Procuradas

Aderência do modelo aos pares (X,y) empíricos

100% aderente

baixa aderência

valor do Eqm

$Eqm(W)$

W

© Prof. Emilio Del Moral – EPUSP

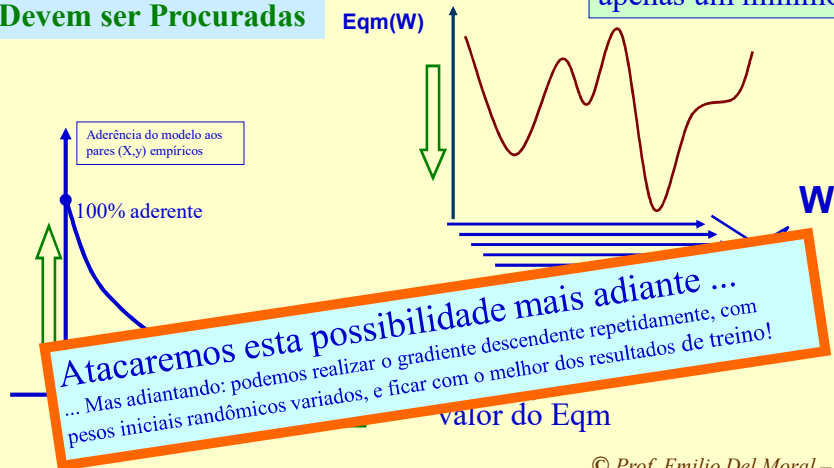
4

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??



5

© Prof. Emilio Del Moral – EPUSP

5

O treinamento mira minimizar o **Eqm** das amostras $(X ; y)$ de treino. (exclusivamente!)

sistema a modelar

Relação de amostragem ...
... E lembremos que as amostras sempre são uma representação parcial do comportamento mais geral do sistema que está sendo modelado.

Tabela de amostras $(X ; y)$

$$Eqm = \frac{\sum_{\mu} [y_{RNA}(X^{\mu}) - y_{sistema}^{\mu}]^2}{M}$$

RNA →
(modelo do sistema)

6

© Prof. Emilio Del Moral – EPUSP

6

Um Exemplo Ilustrativo para o Conceito de Conjunto de Treinamento e dos M pares (X,y) ...

7

© Prof. Emilio Del Moral – EPUSP

7

Exemplo de regressão multivariada para estimação contínua usando MLP

- O valor do y contínuo ... neste exemplo corresponde ao volume de consumo futuro num dado tipo de produto "A" a ser ofertado pela empresa a um cliente corrente já consumidor de outros produtos da empresa ("B" e "C"), volume esse previsto com base em várias medidas quantitativas que caracterizam tal indivíduo. ... Assim, $y = \text{Consumo do Produto A} = F(x_1, x_2, x_3, x_4, x_5)$.
- Consideremos 4 variáveis de entrada no modelo preditivo neural, ou seja, temos 5 medidas em X :
 - x_1 : Idade do indivíduo
 - x_2 : Renda mensal do indivíduo
 - x_3 : Volume de clicks do indivíduo no website de exibição de produtos oferecidos pela empresa
 - x_4 : Volume de consumo desse cliente observado para outro Produto B da mesma empresa
 - x_5 : Volume de consumo desse cliente Produto C da mesma empresa
- Problema: desenvolver uma MLP para regressão contínua multivariada que permita estimar esse volume de consumo futuro y com base no conhecimento dos X e numa base de dados de aprendizado com esses dados X e y para 350 já clientes de universo populacional similar ao do novo consumidor potencial.

8

© Prof. Emilio Del Moral – EPUSP

8

Em termos de Excel, teríamos ...

Cliente (μ)	Idade (x_1)	Renda (x_2)	Clics (x_3)	Consumo do Produto B (x_4)	Consumo do Produto C (x_5)	Consumo do Produto A (y)
1	50	78	302	958	136	9800
2	65	128	186	985	196	8760
3	57	150	221	1093	35	520
....
M-2	16	19	51	707	131	11640
M-1	30	75	7	29	78	9640
M	19	47	116	285	124	5320

9

© Prof. Emilio Del Moral – EPUSP

9

The screenshot shows a presentation slide with the same table as above. A yellow text box with a red border is overlaid on the table, containing the text: *Equivalente em txt Para uso do MBP*. Below the table, a small window titled "treino em txt para exemplo de consumo A e B - Bloco de notas" is open, displaying the data from the table in a text format:

```

Idade  Renda  Clics  ConsumoA  ConsumoB  ConsumoA
50     78     302    958       136       9800
65     128    186    985       196       8760
57     150    221    1093      35        520
(....)
16     19     51     707       131       11640
30     75     7      29        78        9640
19     47     116    285       124       5320
    
```

10

*A estratégia de Aprendizado para o MLP
mais conhecida:*

Error Back Propagation (EBP)

= Propagação Reversa de Erro

*= Método do Gradiente personalizado
ao Eqm(W) do MLP*

11

© Prof. Emilio Del Moral – EPUSP

11

*Mas entendamos PRIMEIRO
o que é o método numérico do
gradiente ascendente /
gradiente descendente
genérico,*

*que pode ser aplicado tanto para se
chegar paulatinamente ao máximo de
uma função quanto para se chegar ao
mínimo de uma função
(ascendente / descendente)*

12

© Prof. Emilio Del Moral – EPUSP

12

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

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

$$\frac{\partial f(a,b,c)}{\partial a} \quad \frac{\partial f(a,b,c)}{\partial b} \quad \frac{\partial f(a,b,c)}{\partial c}$$

- Vetor Gradiente, útil ao método do máximo declive:

$$\left(\frac{\partial \text{Eqm}(W)}{\partial w_1}, \frac{\partial \text{Eqm}(W)}{\partial w_2}, \frac{\partial \text{Eqm}(W)}{\partial w_3}, \dots \right) \quad \Delta \vec{W} = -\eta \cdot \vec{\nabla} \text{Eqm}_-$$

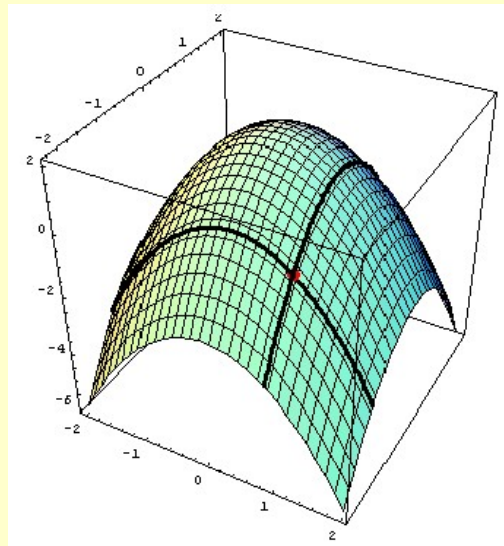
- Regra da cadeia, necessária ao cálculo de derivadas quando há encadeamento de funções:

$$\frac{\partial f(g(h(a)))}{\partial a} = \frac{\partial f}{\partial g} \cdot \frac{\partial g}{\partial h} \cdot \frac{\partial h}{\partial a}$$

© Prof. Emilio Del Moral – EPUSP

13

Derivada parcial- ilustração p/ função de 2 variáveis apenas

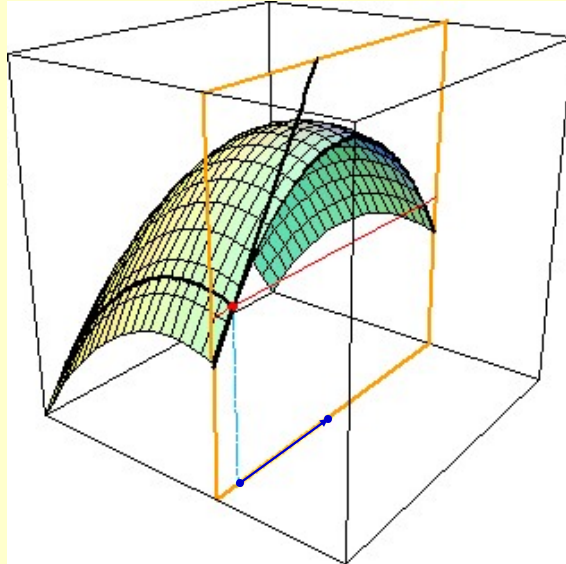


A visual model of the partial derivative

© Prof. Emilio Del Moral – EPUSP

14

Derivada parcial- ilustração p/ função de 2 variáveis apenas



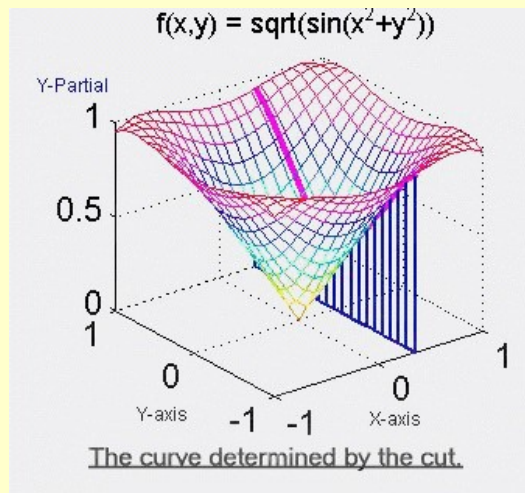
A visual model of the partial derivative

15

© Prof. Emilio Del Moral – EPUSP

15

mais ilustrações p/ a derivada parcial em função de 2 variáveis

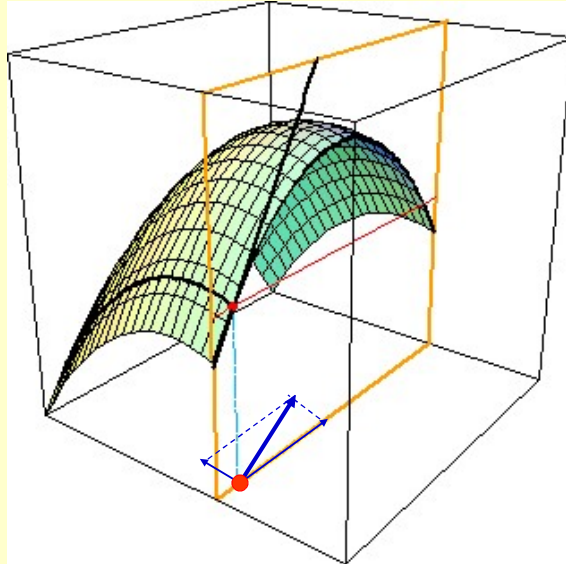


A visual model of the partial derivative with respect to y. 16

© Prof. Emilio Del Moral – EPUSP

16

Formação do vetor gradiente a partir de duas derivadas parciais



A visual model of the partial derivative

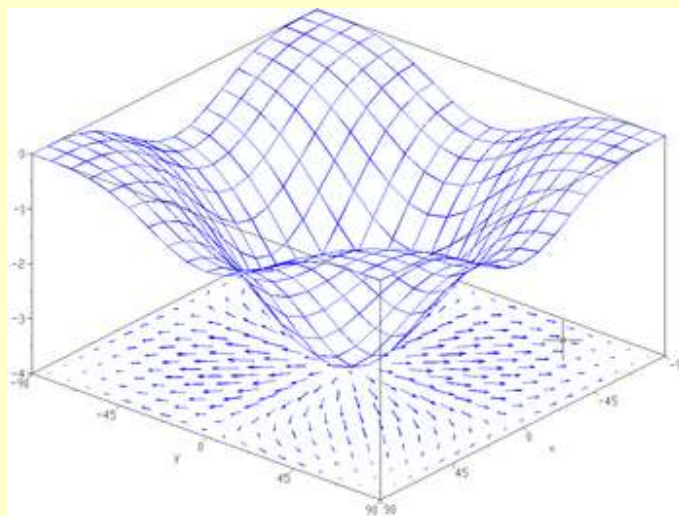
17

© Prof. Emilio Del Moral – EPUSP

17

<http://en.wikipedia.org/wiki/Gradient>

... O vetor gradiente indica a direção ascendente e seu módulo a magnitude de crescimento da função escalar – ilustração p/ função de 2 variáveis apenas



18

© Prof. Emilio Del Moral – EPUSP

18

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

- “Chute” um W inicial para o “Wcorrente”, 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^μ associado a um exemplo μ , e vá varrendo μ e somando os gradientes de cada Eq^μ , para compor o vetor gradiente de Eqm, assim que sair deste loop em μ ;
 - 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 ΔW nesse espaço, com direção e magnitude dados por $-\eta \cdot$ vetor gradiente médio para os M Exemplos $(X^\mu; y^\mu)$ de treino

19

© Prof. Emilio Del Moral – EPUSP

19

Processo de refinamentos graduais a cada iteração ...

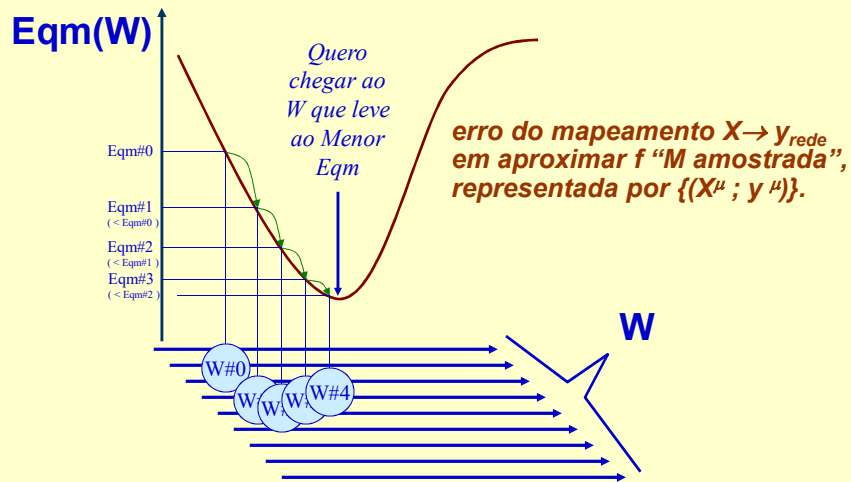
$W\#0$	$Eqm\#0$	$GradEqm(W\#0)$	$\Delta W\#0 = -n \cdot GradEqm(W\#0)$
$W\#1$ (= $W\#0 + \Delta W\#0$)	$Eqm\#1$ (< $Eqm\#0$)	$GradEqm(W\#1)$	$\Delta W\#1 = -n \cdot GradEqm(W\#1)$
$W\#2$ (= $W\#1 + \Delta W\#1$)	$Eqm\#2$ (< $Eqm\#1$)	$GradEqm(W\#2)$	$\Delta W\#2 = -n \cdot GradEqm(W\#2)$
$W\#3$ (= $W\#2 + \Delta W\#2$)	$Eqm\#3$ (< $Eqm\#2$)	$GradEqm(W\#3)$	$\Delta W\#3 = -n \cdot GradEqm(W\#3)$
$W\#4$ (= $W\#3 + \Delta W\#3$)	$Eqm\#4$ (< $Eqm\#3$)	$GradEqm(W\#4)$	$\Delta W\#4 = -n \cdot GradEqm(W\#4)$
...
$W\#k$ (= $W\#k-1 + \Delta W\#k-1$)	$Eqm\#k$ (< $Eqm\#k-1$)	$GradEqm(W\#k)$	$\Delta W\#k = -n \cdot GradEqm(W\#k)$
...

20

© Prof. Emilio Del Moral – EPUSP

20

A estratégia do EBP / Gradiente Descendente no aprendizado do MLP



© Prof. Emilio Del Moral – EPUSP

21

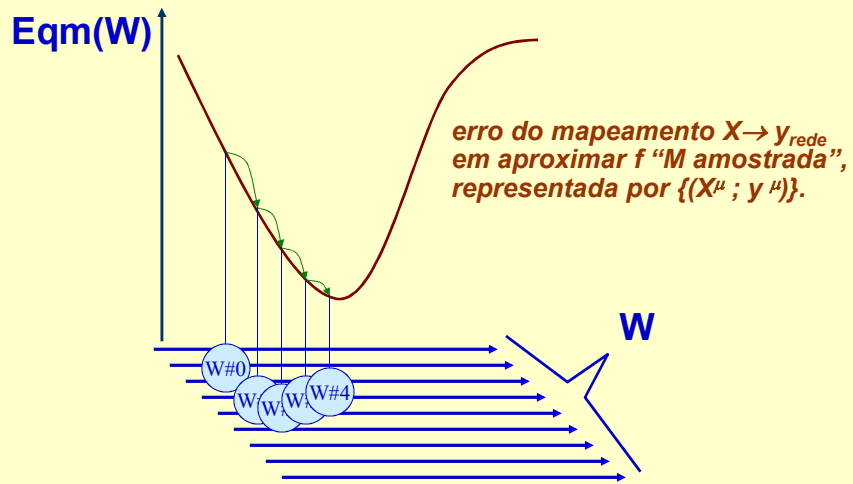
Este gráfico é denso e toca em muitos aspectos interrelacionados ... revisitemos alguns desses aspectos isoladamente com focos específicos nessas revisitas, assim teremos gráficos algo mais simples de interpretar ...

22

© Prof. Emilio Del Moral – EPUSP

22

A estratégia do EBP / Gradiente Descendente no aprendizado do MLP

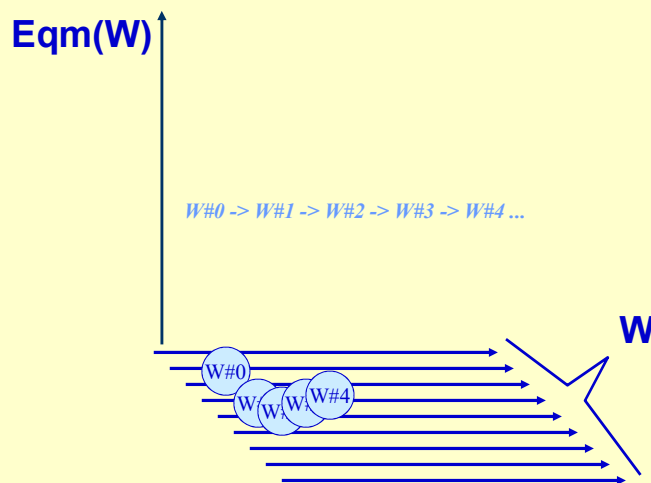


23

© Prof. Emilio Del Moral – EPUSP

23

Foco na evolução dos w 's com as iterações ...

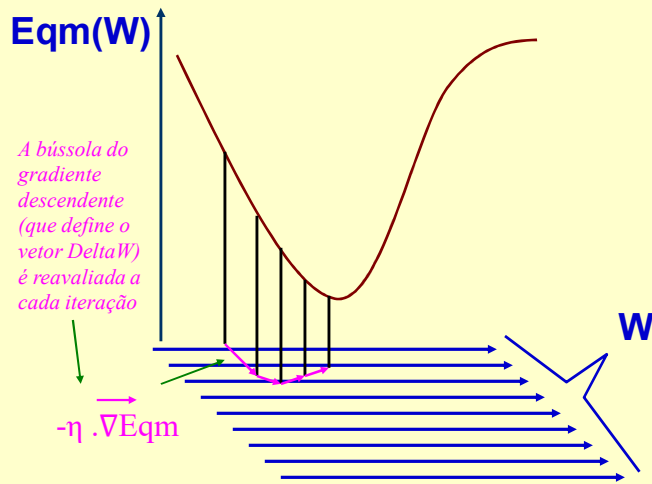


24

© Prof. Emilio Del Moral – EPUSP

24

Foco nos diferentes DeltaW de cada iteração ...

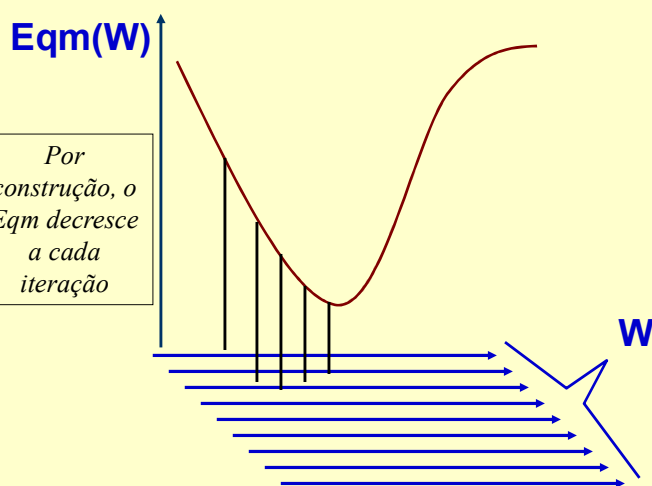


25

© Prof. Emilio Del Moral – EPUSP

25

Foco na evolução do Eqm com as iterações ...



26

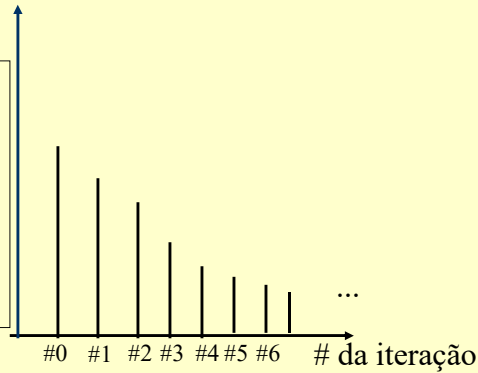
© Prof. Emilio Del Moral – EPUSP

26

Plotando a evolução do Eqm com as iterações ...

Eqm(#)

Por construção, Eqm decresce a cada iteração, até estabilização em ponto de mínimo

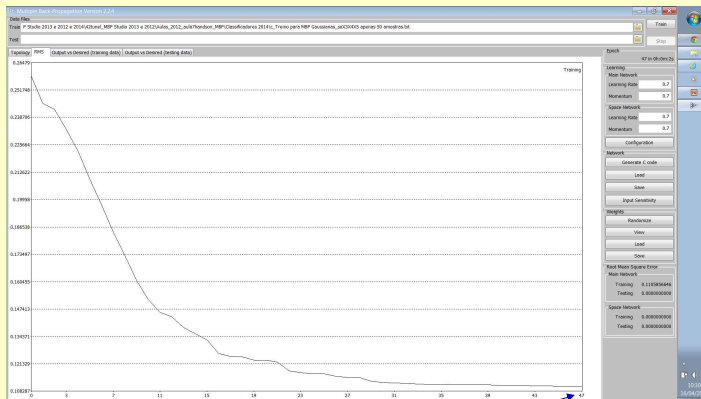


27

© Prof. Emilio Del Moral – EPUSP

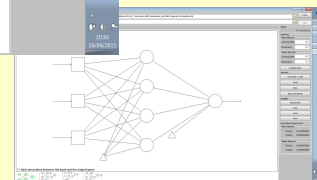
27

Gráfico fornecido pelo ambiente MBP da evolução do Eqm com o número de repetidos usos da “bússola do gradiente descendente”: isto conecta o MBP com o gráfico apresentado no slide anterior



Nota: o RMS do eixo vertical deste gráfico significa Root Mean Square, e é a raiz quadrada do nosso conhecido Eqm

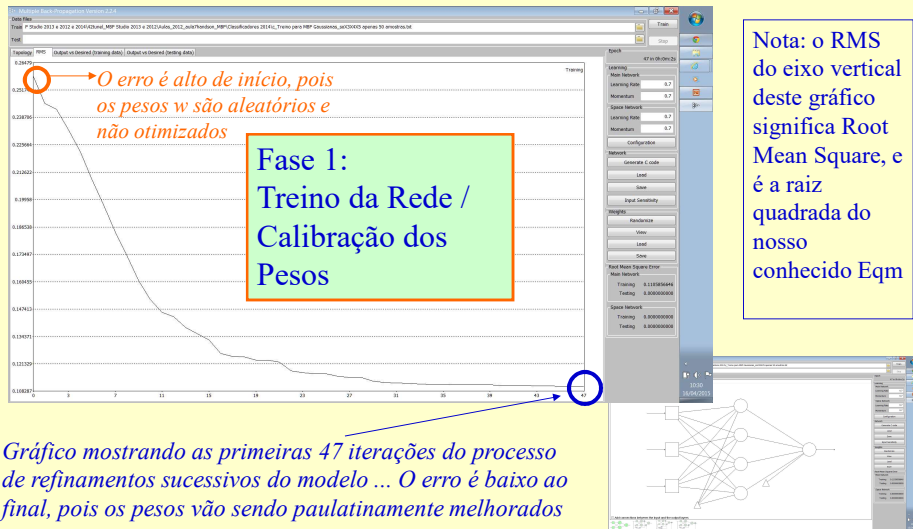
Gráfico mostrando as primeiras 47 iterações do processo de refinamentos sucessivos do modelo ...



© Prof. Emilio Del Moral – EPUSP

28

Gráfico fornecido pelo ambiente MBP da evolução do Eqm com o número de repetidos usos da “bússola do gradiente descendente”:
isto conecta o MBP com o gráfico apresentado no slide anterior



© Prof. Emilio Del Moral – EPUSP

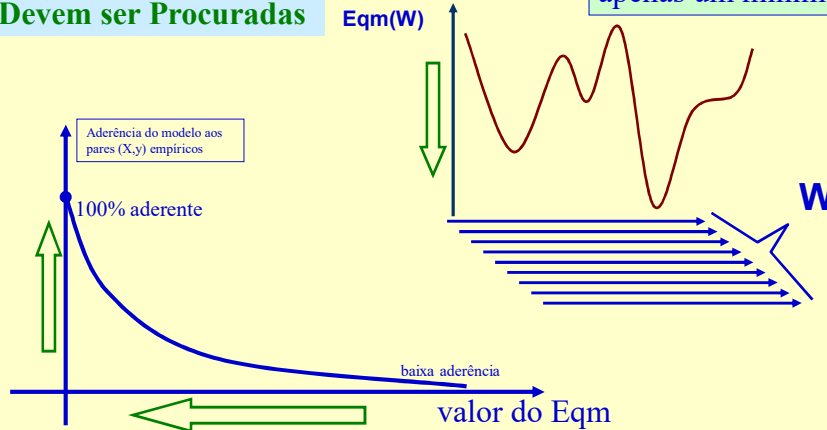
29

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??



30

© Prof. Emilio Del Moral – EPUSP

30

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

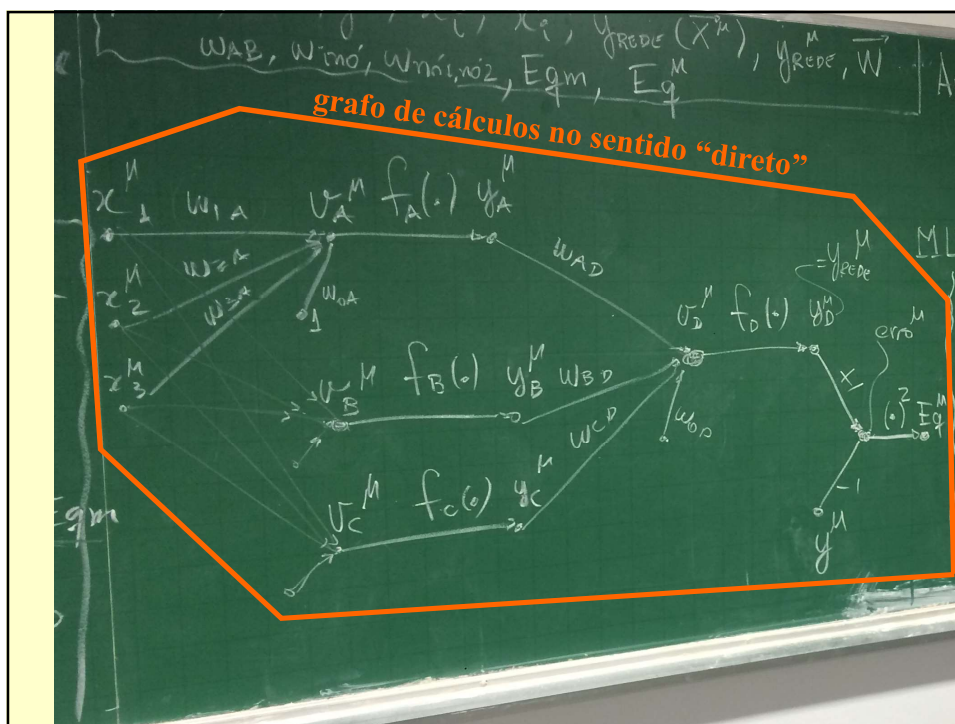
Gradiente de Eqm no espaço de pesos = $(\partial E_{qm}(W)/\partial w_1, \partial E_{qm}(W)/\partial w_2, \partial E_{qm}(W)/\partial w_3, \dots)$

**Recordemos como obtivemos
(em aula anterior) às fórmulas de
cada derivada parcial acima ...**

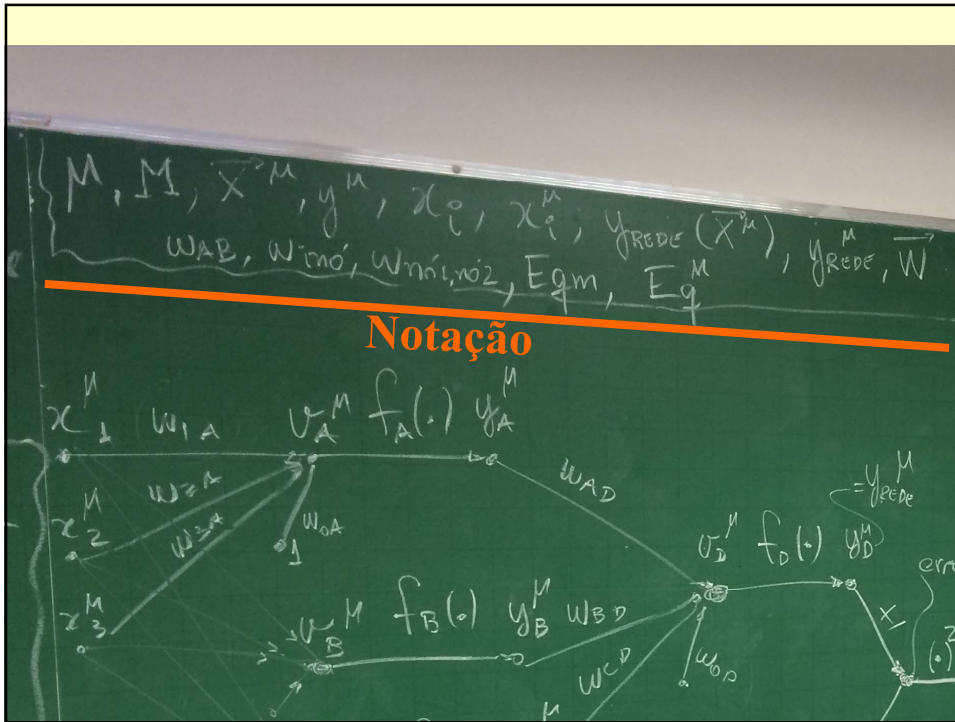
31

© Prof. Emilio Del Moral – EPUSP

31



32



33

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

34

© Prof. Emilio Del Moral – EPUSP

34

Importante ... O que está implícito nas imagens de redes apresentadas em lousa e nas telas do MBP ... (relembrando premissas colocadas nas primeiras aulas)

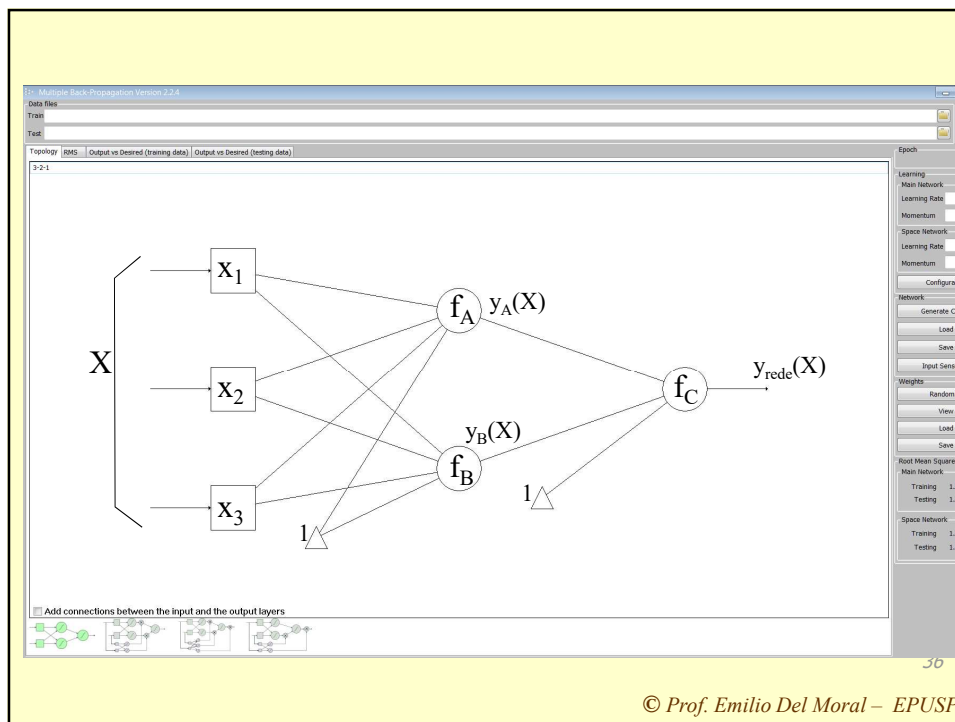
... O fato de que cada nó neural (cada "círculo" nos diagramas da RNA) que compõe a rede gera na saída um cálculo do tipo tgh (ou sigmoidal de sua escolha) que opera sobre a somatória ponderada (ponderação pelos seus w 's) das diversas entradas que lhe chegam (e incluindo também um viés nessa somatória)

- Isto ocorre para todos os nós neurais da RNA: se tivermos por exemplo 2+1 (3 no total) nós neurais, como na RNA exemplo do slide anterior, cada um deles realizará um cálculo desse tipo, ou seja, tgh(soma ponderada), empregando os valores específicos de seus pesos ponderadores e seus vieses exclusivos, valores esses que podem ser distintos para cada nó
- Em particular, os nós das camadas mais adiante na RNA têm como suas entradas as variáveis de saída dos nós da camada anterior (ou seja, operam sobre as saídas de tangentes hiperbólicas das camadas anteriores)
- Isto que digo estar implícito no MBP já faz parte de qualquer rede neural (*revise os slides das primeiras aulas*), por isso não precisa ser detalhado em cada figura
- De qualquer forma, para ajudar a visualizar o que está implícito e que é premissa em RNAs, nos slides que seguem eu adiciono os significados assumidos na rede exemplo 3-2-1 do slide anterior.

35

© Prof. Emilio Del Moral – EPUSP

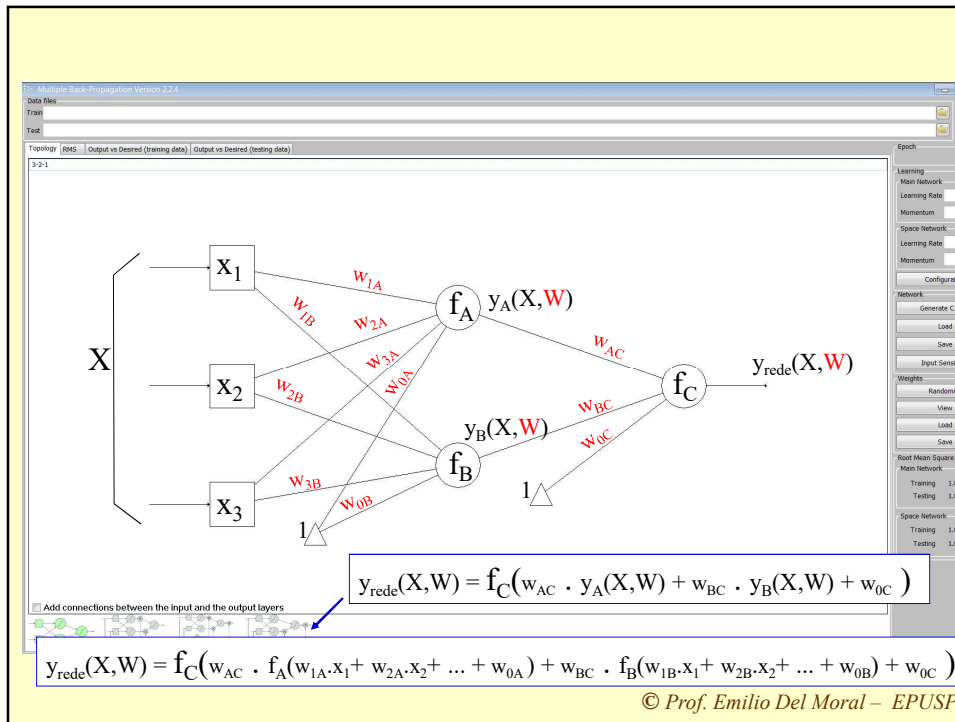
35



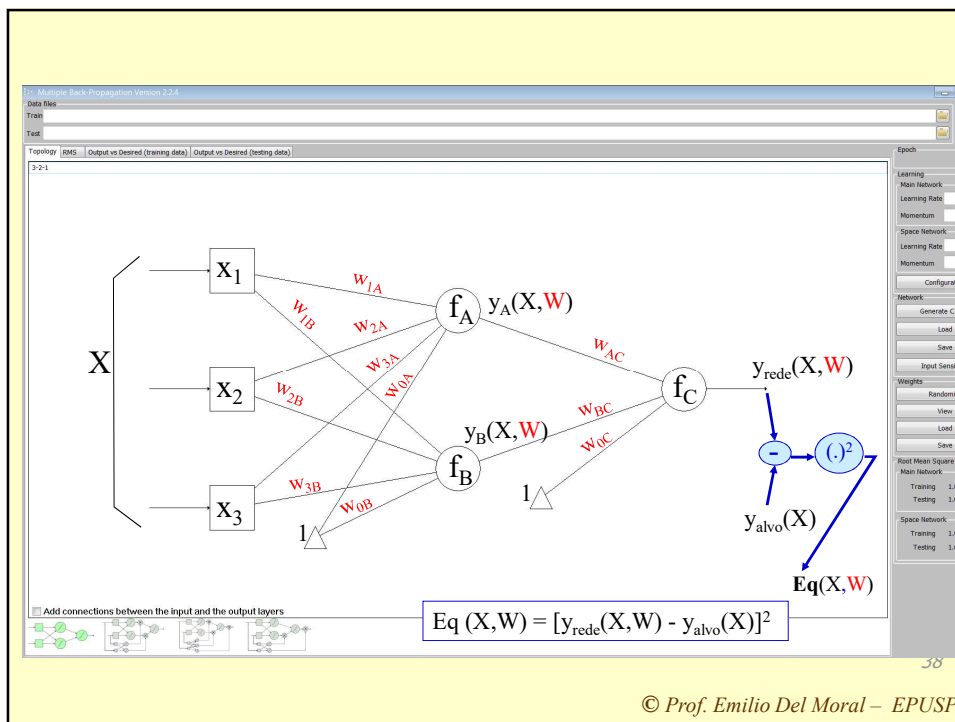
36

© Prof. Emilio Del Moral – EPUSP

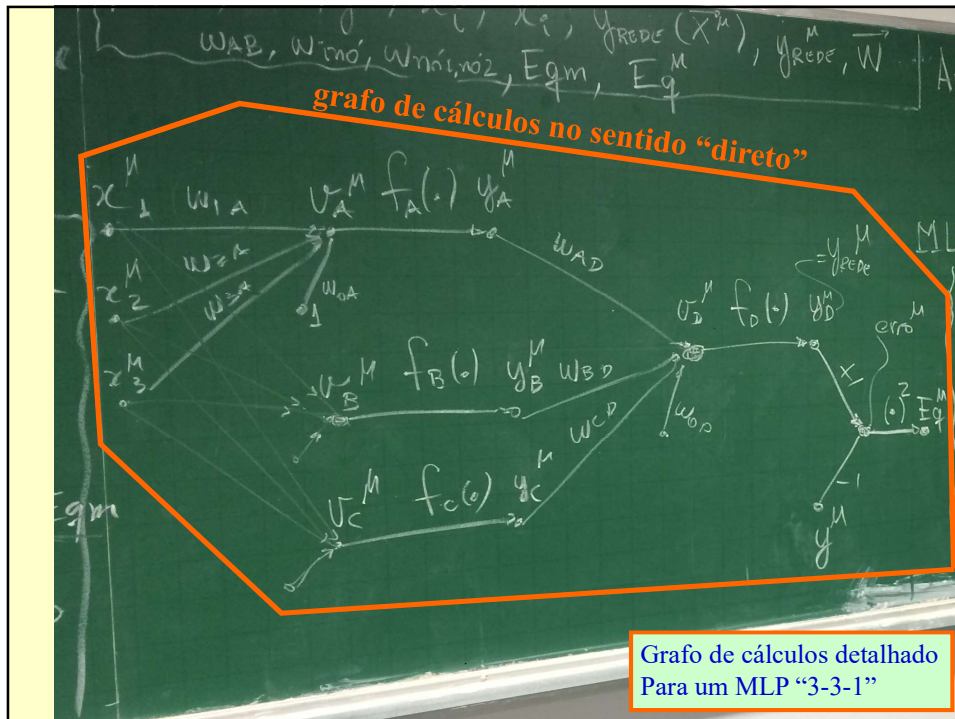
36



37



38



39

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):**

$$\frac{\partial f(a,b,c)}{\partial a} \quad \frac{\partial f(a,b,c)}{\partial b} \quad \frac{\partial f(a,b,c)}{\partial c}$$

- **Vetor Gradiente, útil ao método do máximo declive:**

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

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

40

© Prof. Emilio Del Moral – EPUSP

40

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:

$$\begin{aligned} \mathbf{Grad}(Eqm) &= \\ \mathbf{Grad}(\sum_{\mu} Eq^{\mu}) / M & \\ \sum_{\mu} \mathbf{Grad}(Eq^{\mu}) / M & \end{aligned}$$

41

© Prof. Emilio Del Moral – EPUSP

41

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

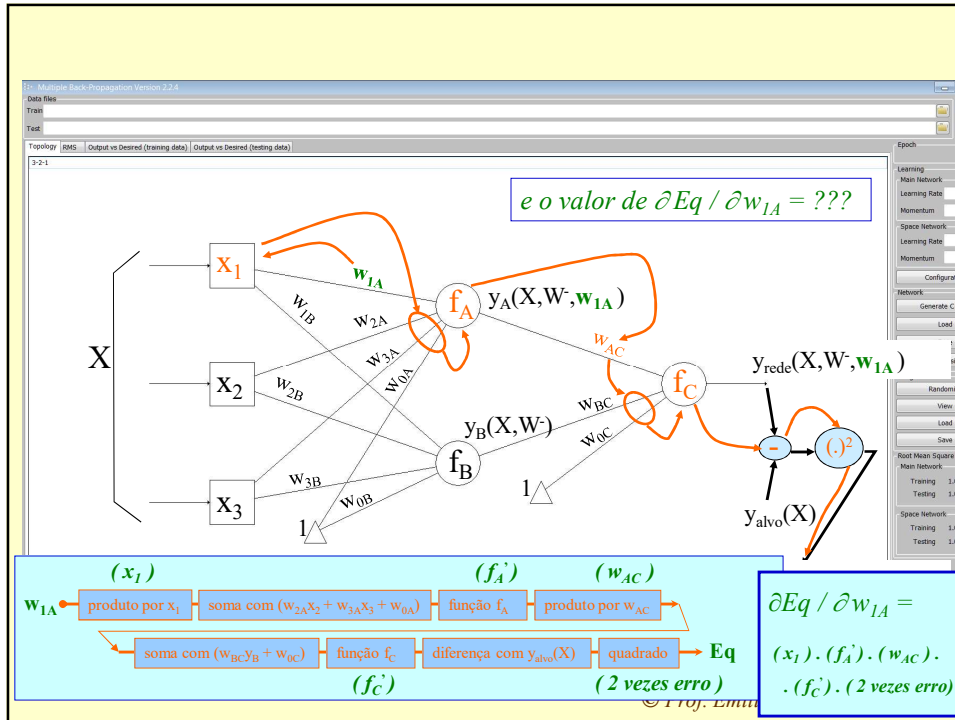
do vetor $\mathbf{Grad}(\sum_{\mu} Eq^{\mu})$ – a seguinte propriedade simples e sua velha conhecida ...

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

42

© Prof. Emilio Del Moral – EPUSP

42



43

Lembretes

- Na maioria dos slides anteriores, onde aparece X , leia-se X^μ , não incluído para não complicar demais os desenhos
- ... similarmente, onde aparece y_{alvo} , leia-se y_{alvo}^μ . Idem para os Eq, leia-se Eq^μ
- 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_{1A} 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 $[y_{rede}(X, W) - y_{alvo}(X)]$

44

© Prof. Emilio Del Moral – EPUSP

44

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_{IB} , 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^μ com relação aos 11 diferentes pesos w : $\text{Grad}_w(Eq^\mu)$
- Essas 11 fórmulas devem ser aplicadas repetidamente aos M exemplares numéricos de X^μ e y_{alvo}^μ , calculando M gradientes!
- Com eles, se obtém o gradiente médio dos M pares empíricos: $\text{Grad}_w(Eq_m) = [\sum_\mu \text{Grad}_w(Eq^\mu)] / M$
- Esse gradiente médio é a Bussola do Gradiente!

45

© Prof. Emilio Del Moral – EPUSP