

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

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

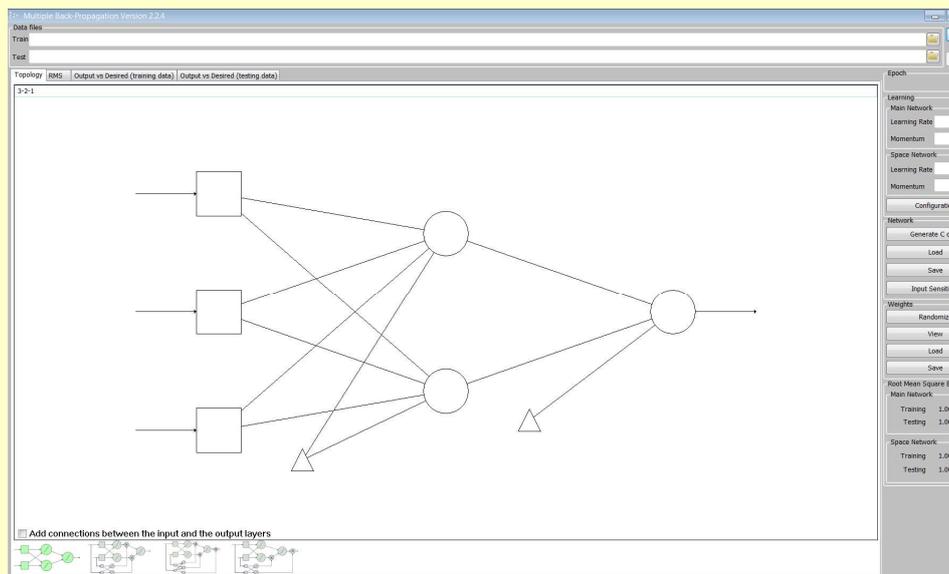
Gradiente de Eqm no espaço de pesos =  $(\partial E_{qm}(W)/\partial w_1, \partial E_{qm}(W)/\partial w_2, \partial E_{qm}(W)/\partial w_3, \dots)$

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

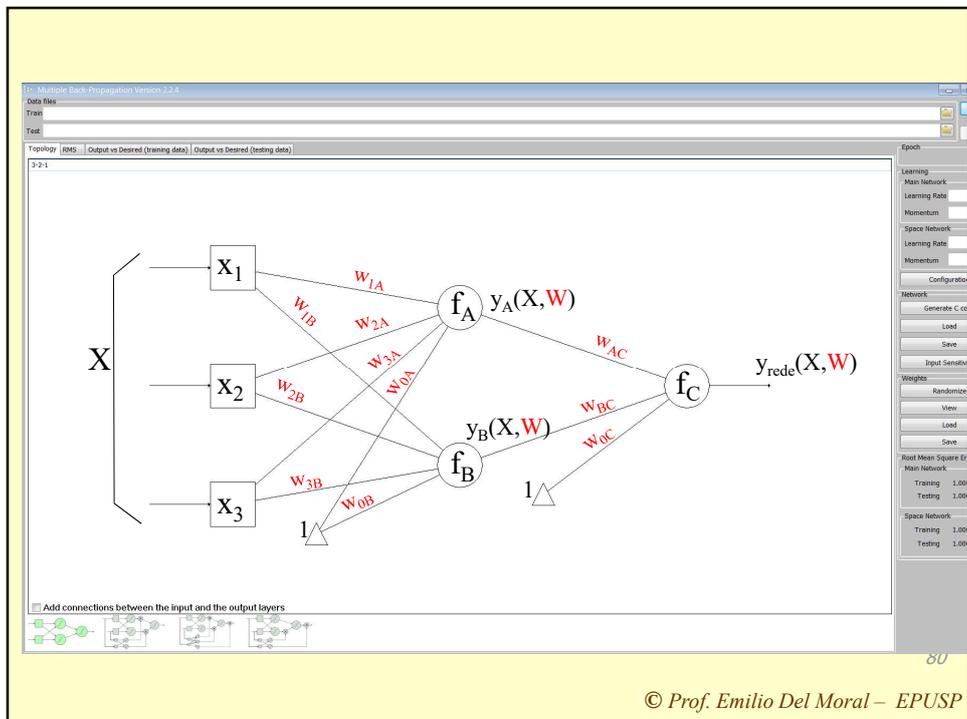
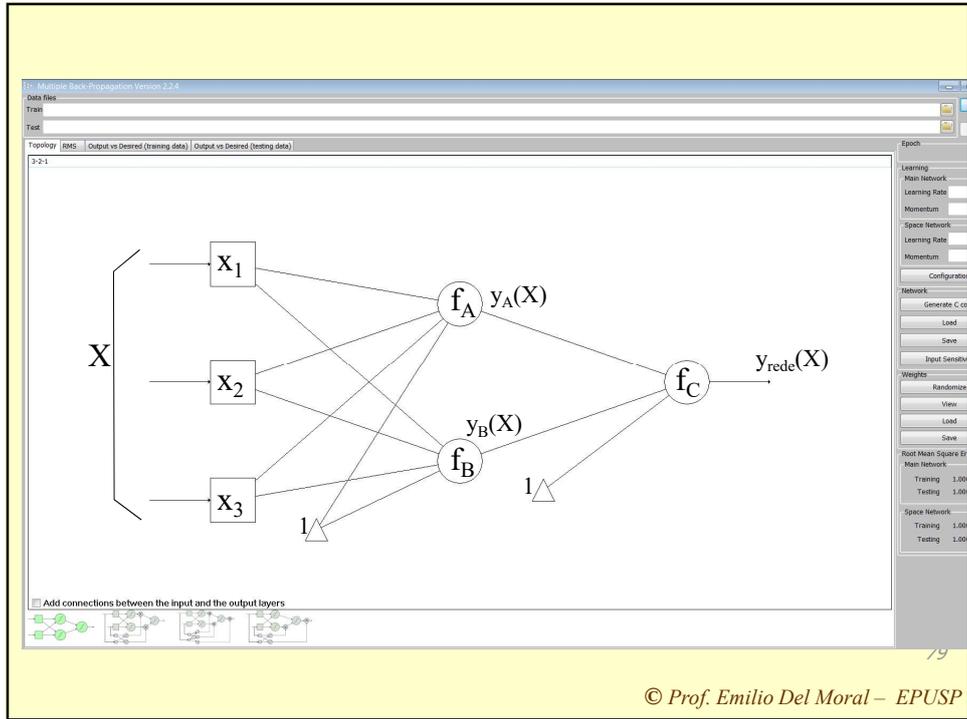
76

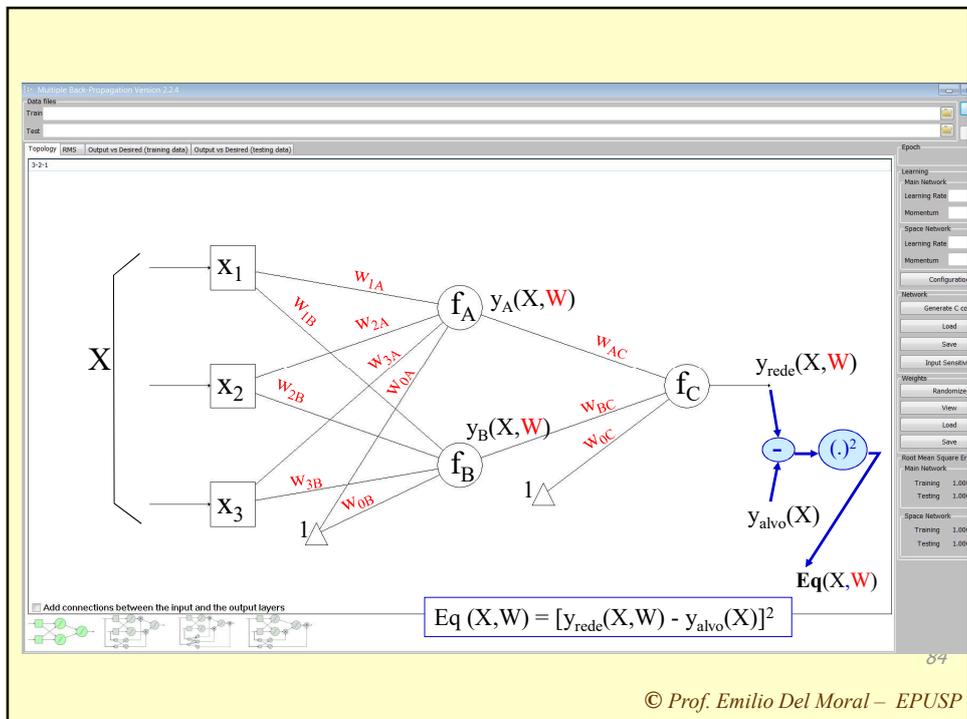
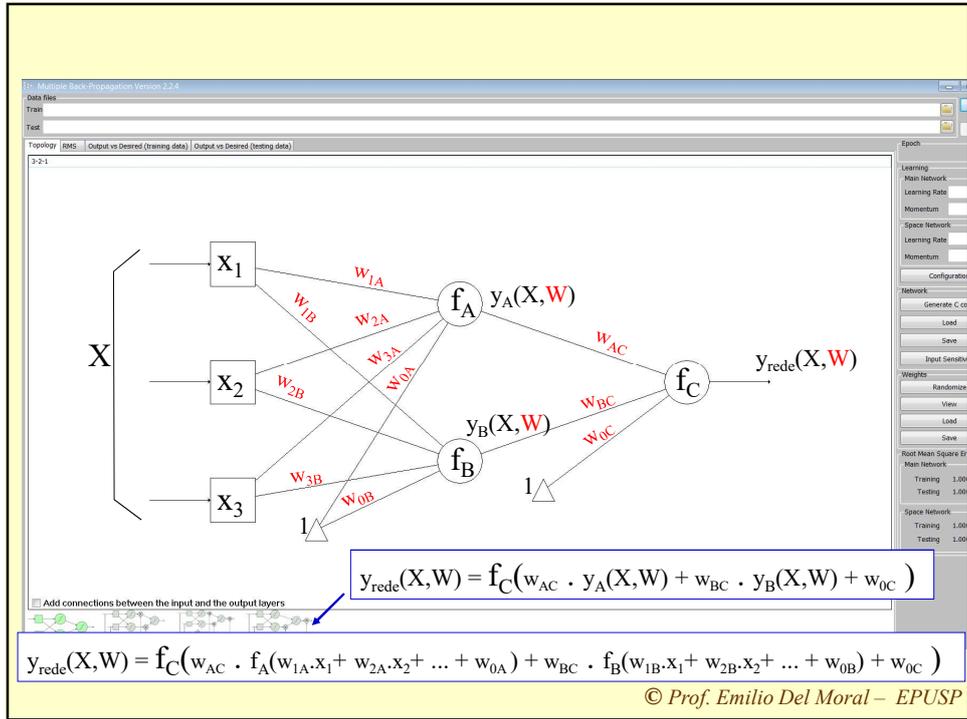
© Prof. Emilio Del Moral – EPUSP

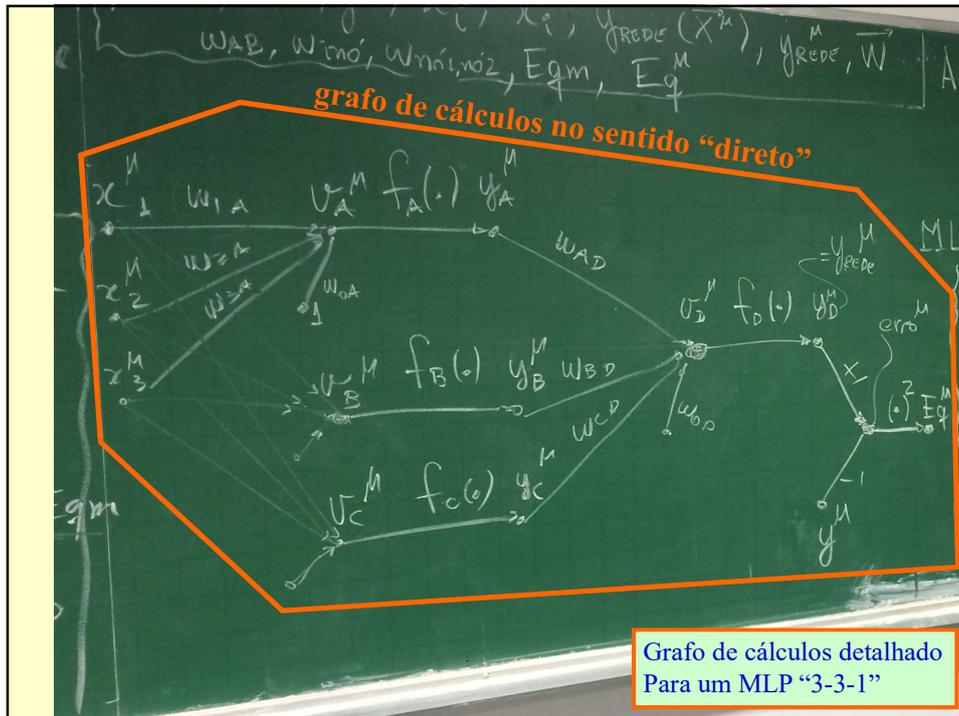
*Relembrando o que está por trás de um desenho como o que segue ...*



© 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):**

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

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

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

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

## 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:

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

87

© Prof. Emilio Del Moral – EPUSP

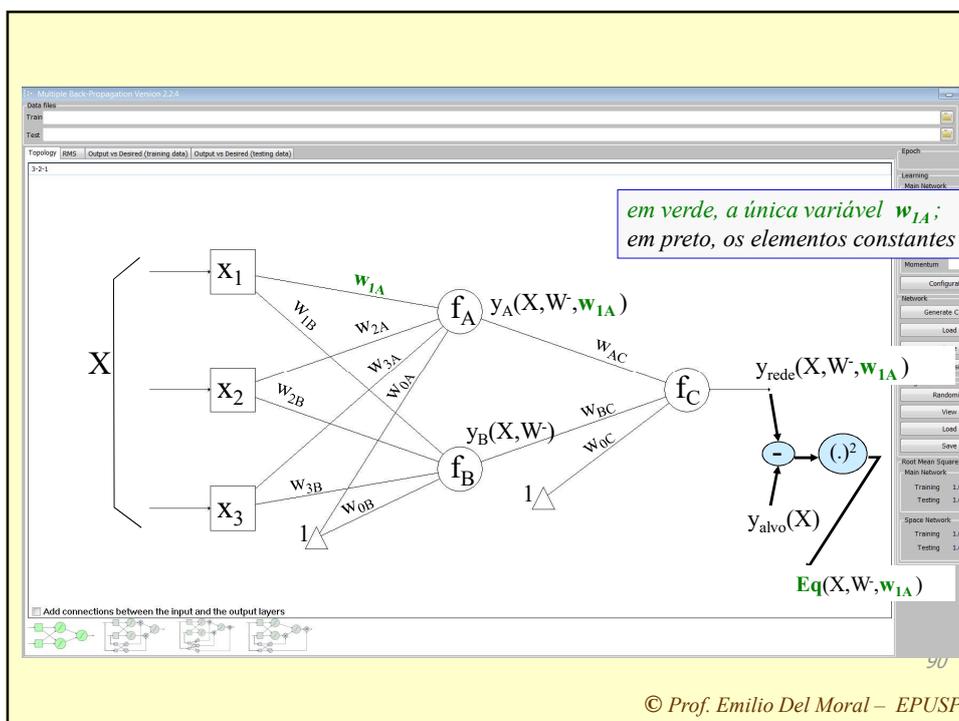
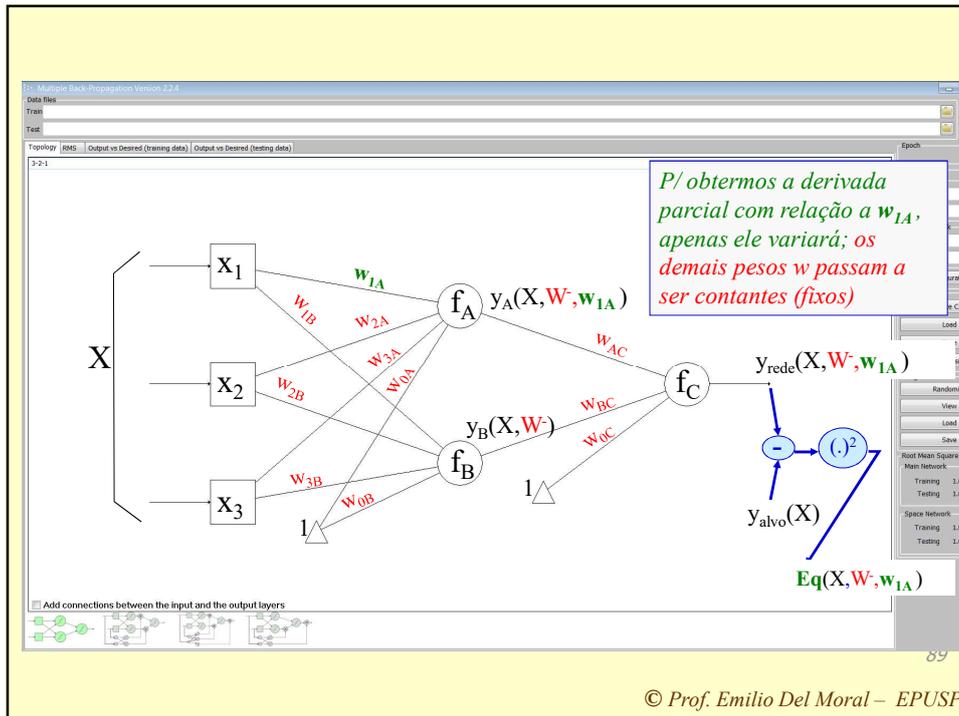
Note que a inversão do gradiente com a somatória nada mais é que usar de forma repetida – e em separado para cada dimensão

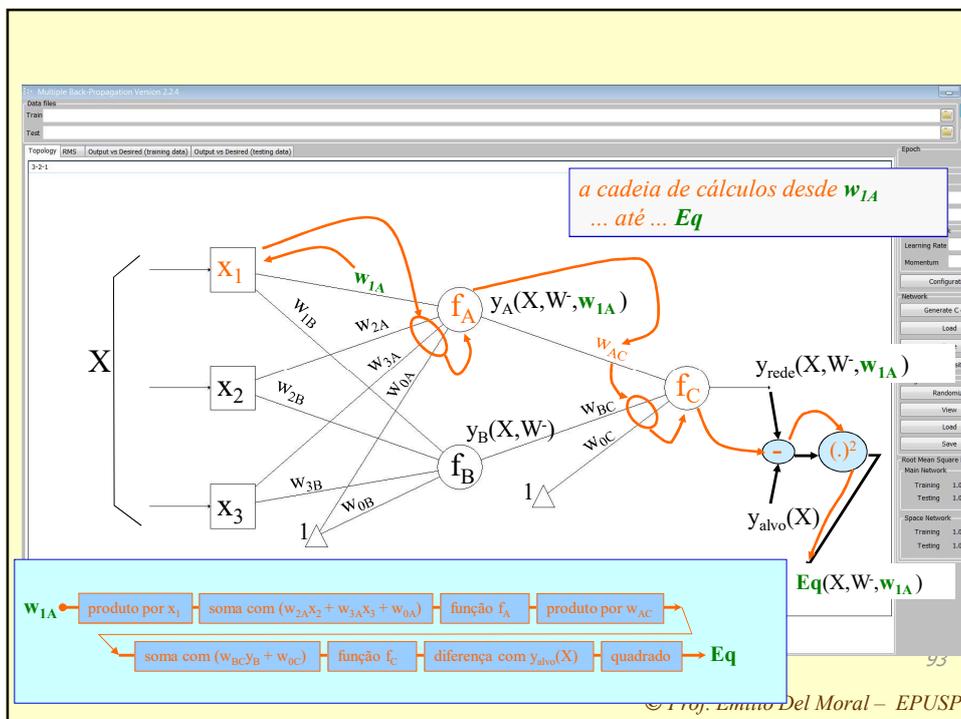
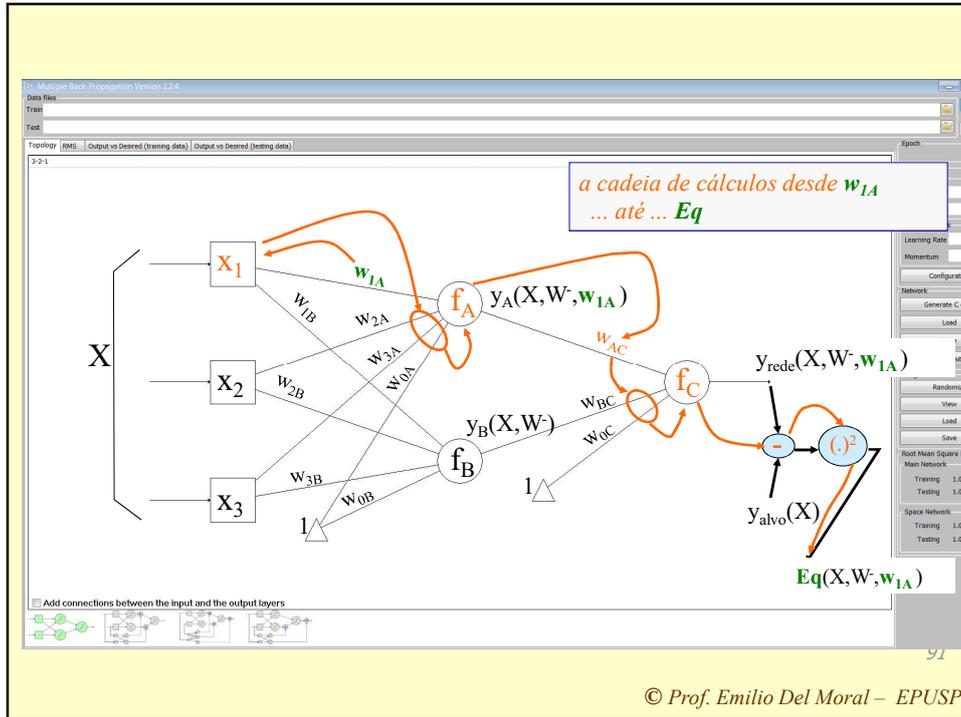
do vetor  $\mathbf{Grad}(\sum_{\mu} Eq^{\mu})$  – a seguinte propriedade simples e sua velha conhecida ...

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

88

© Prof. Emilio Del Moral – EPUSP





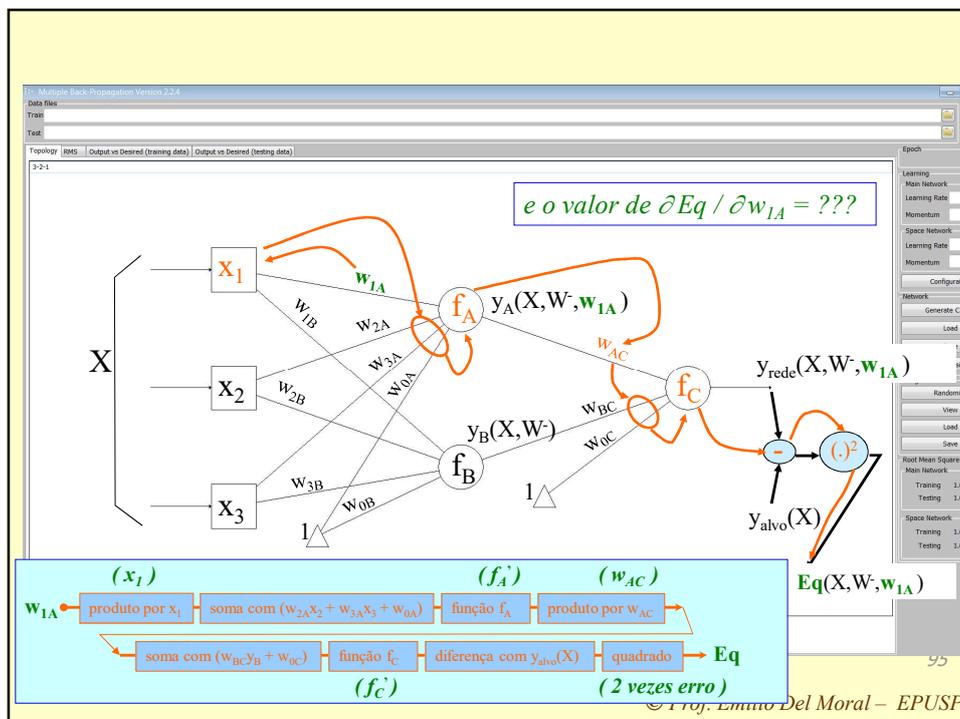
Note que aqui temos uma cadeia com muitos estágios que levam da variável  $w_{1A}$ , à variável  $Eq^u$ , e para a qual podemos calcular a derivada da saída ( $Eq^u$ ) com relação à entrada ( $w_{1A}$ ) aplicando de forma repetida a seguinte propriedade simples e sua velha conhecida ...

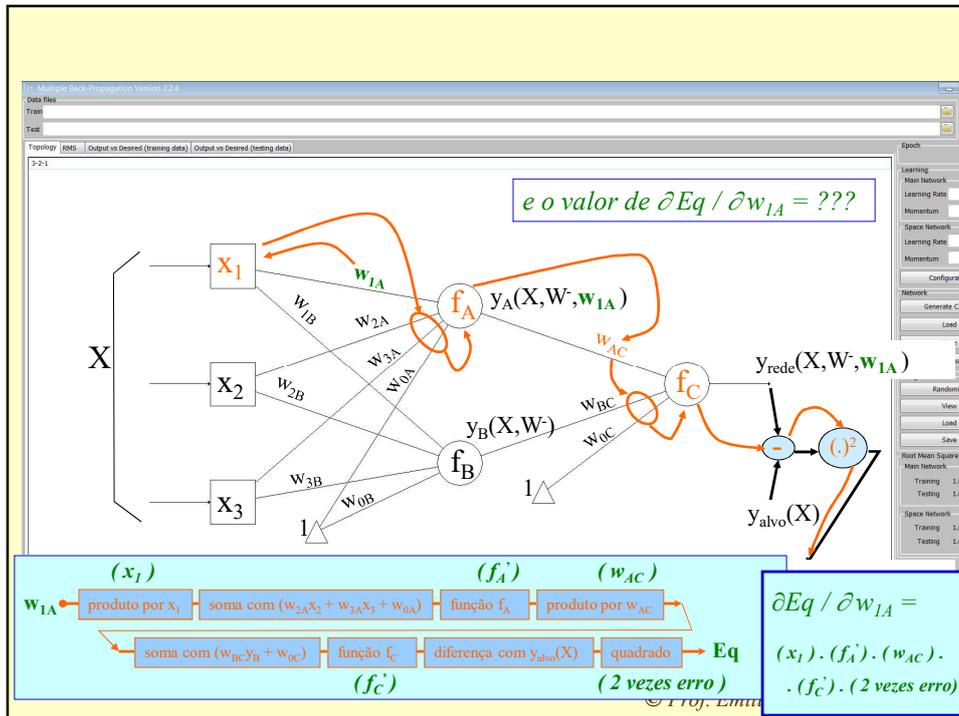
$$d(f_1(f_2(x)))/dx = df_1(x)/df_2 \cdot df_2(x)/dx$$

..., 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.

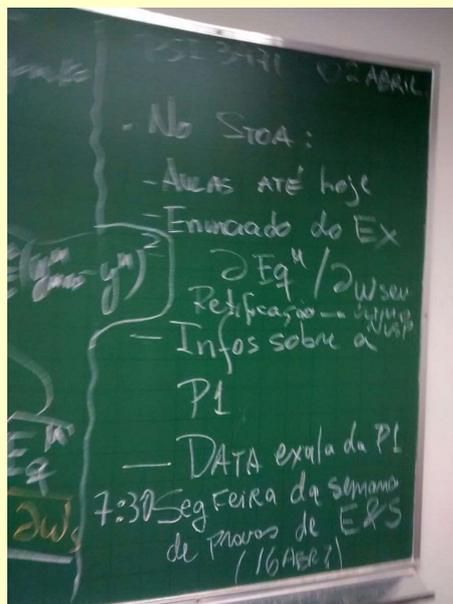
94

© Prof. Emilio Del Moral – EPUSP





## Avisos que detalhei na aula de 02 de abril



A nossa PI ocorrerá no horário de aula de PSI3471 e na semana de PI (de E&S), como no calendário divulgado no início do semestre; mais precisamente, ocorrerá no dia 16 de abril, 2ª feira às 7:30 (com duração prevista de 90 mins).

98

© Prof. Emilio Del Moral – EPUSP

## Avisos que comentei na aula de 02 de abril

Regras de consulta a folha A1 própria apenas; data e local.

- 1) Como conversado em sala de aula, a consulta para a P1 está restrita a uma folha tamanho A4 própria preparada por você com suas anotações de estudo que considere mais relevantes, identificadas claramente com seu nome completo seu NUSP e assinatura, e que será entregue junto com a prova; escaneie / fotografe essa folha de estudos antes do dia da prova para seus arquivos.
- 2) No dia da prova não será permitido nem necessário o uso de computadores, celulares ou quaisquer eletrônicos, apenas será necessária a consulta a sua folha A4 de estudos previamente preparada, identificada e entregue com a prova.
- 3) A prova ocorrerá dentro do horário de aulas na semana de provas do 4o ano de Eletrônica e Sistemas. Se tiver dúvidas sobre essas datas, consulte a tabela de aulas e provas PSI3471 disponível no STOA desde o início do semestre e projetada sala na primeira semana de aulas.
- 4) O local da prova será a sala de aula ou sala maior a ser definida mais proximamente à prova.

No STOA também está definido desde o início do semestre quem pode fazer PSubs por perda da P1 e os documentos que devem ser apresentados na secretaria para análise pelos professores da solicitação do aluno.

Última atualização: terça, 20 Mar 2018, 08:46

99

© Prof. Emilio Del Moral – EPUSP

## Lembretes ....

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

101

© Prof. Emilio Del Moral – EPUSP

## Lembretes ....

- O mesmo que foi feito para  $w_{1A}$  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  $w$ :  $Grad_w(Eq^\mu)$
- Essas 11 fórmulas devem ser aplicadas repetidamente aos  $M$  exemplares numéricos de  $X^\mu$  e  $y_{alvo}^\mu$ , calculando  $M$  gradientes!
- Com eles, se obtém o gradiente médio dos  $M$  pares empíricos:  $Grad_w(Eqm) = [\sum_\mu Grad_w(Eq^\mu)] / M$
- Esse gradiente médio é a Bussola do Gradiente!

102

© 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

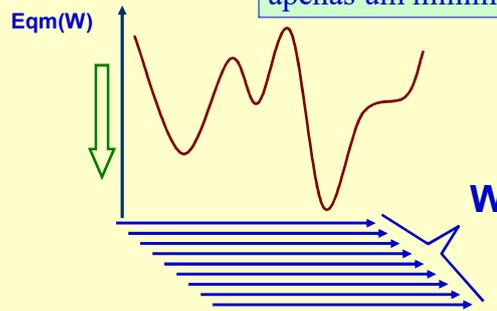
103

© 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  $E_{qm}$  possível*

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



Será que temos apenas um mínimo??

### Atacando a possibilidade de múltiplos mínimos locais:

... Executamos o gradiente descendente repetidamente, com a mesma arquitetura neural e os mesmos dados empíricos de treino, mas com os pesos iniciais randômicos sendo distintos em cada rodada; anotamos todos os pesos otimizados e valores de  $E_{qm}$  final de todas as rodadas; mantemos os  $w$ 's do melhor dos resultados de treino entre eles (aquele ensaio com menor valor do  $E_{qm}$ )!

## O Ciclo completo da modelagem:

0) *Formalização do problema, mapeamento quantitativo em um modelo neural inicial e ... 0b) coleta de pares empíricos  $(X,y)$*

1) *Fase de TREINO da RNA (MLP): com conhecimento dos  $X$  e dos  $y$ , que são ambos usados na calibração do modelo*

2) *Fase de TESTE / Caracterização da qualidade da RNA para generalizar: temos novos pares  $X$  e  $y$ , com  $y$  guardado "na gaveta", usado apenas para avaliação, não para re-calibração. É como um ensaio de uso final do modelo, com possibilidade de medir a sua qualidade com o  $y$  que foi guardado na gaveta.*

[Fase de refinamentos da RNA, dados e modelo, em ciclos, desde 0]

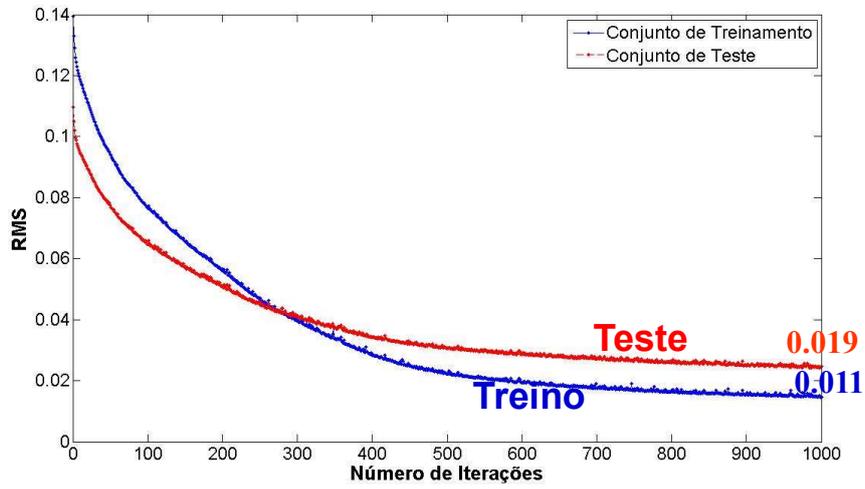
3) *Fase de USO FINAL da RNA, com  $y$  efetivamente não conhecido, e estimado com conhecimento dos  $X$  + uso do modelo calibrado.*

... Diferenças e semelhanças entre 1, 2 e 3

107

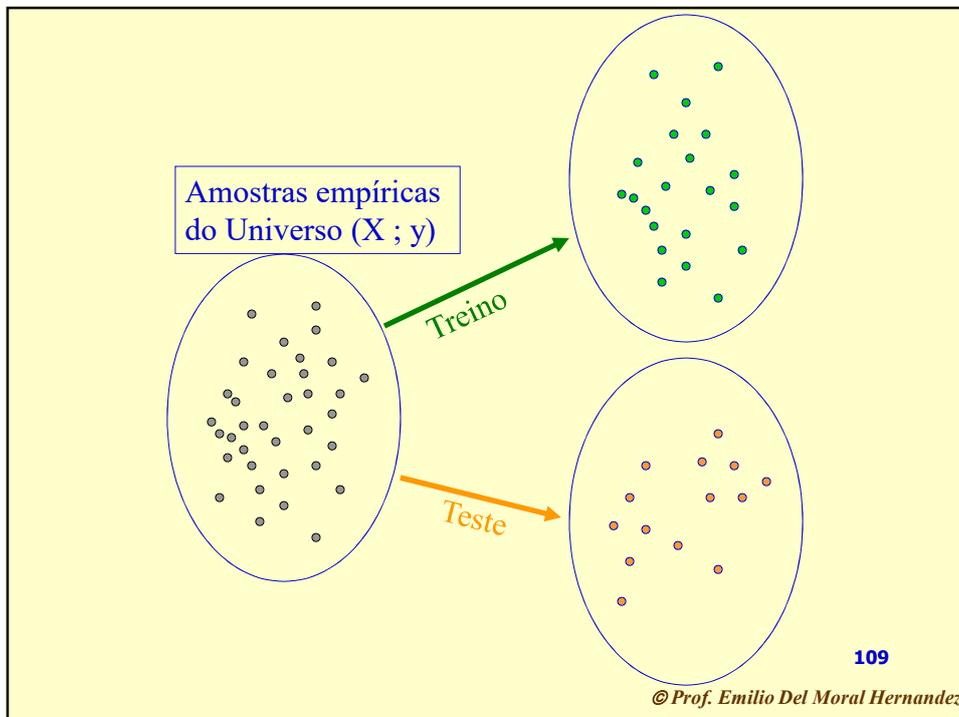
© Prof. Emilio Del Moral – EPUSP

**Raiz do Erro Quadrático Médio (RMS) p/ conjuntos de treino e teste ...  $RMS_{teste} > RMS_{treino}$**



108

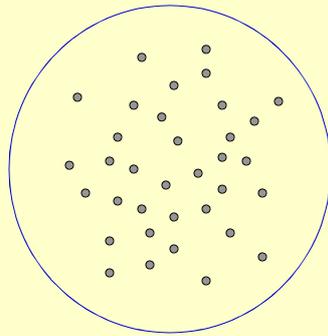
© Prof. Emilio Del Moral Hernandez



109

© Prof. Emilio Del Moral Hernandez

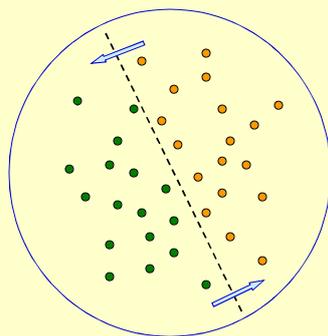
## Conjunto total de observações empíricas



110

© Prof. Emilio Del Moral Hernandez

*Cross validation / Validação cruzada ... e teremos muitos mais ensaios, se tivermos mudança na partição de observações*



2 fold cross  
validation:  
50% treino e  
50% teste

111

© Prof. Emilio Del Moral Hernandez

***k-fold Cross Validation:***  
***O conjunto total é “retalhado”***  
***em k partes, e uma delas apenas***  
***é usada para teste, com k***  
***ensaios distintos***

112

© Prof. Emilio Del Moral Hernandez

***Conceito geral que engloba***  
***as discussões anteriores***

...

***Reamostragem /***  
***Data Resampling***

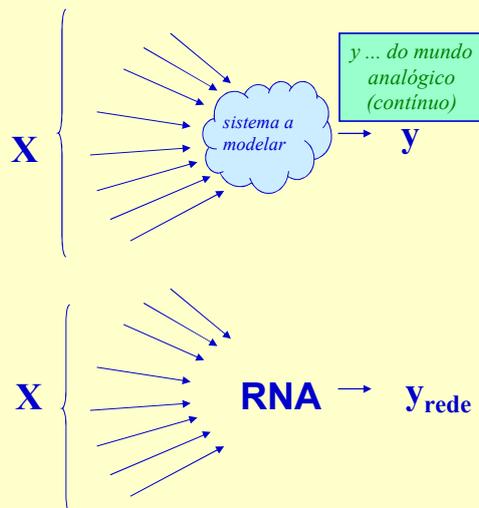
Pergunta ... Que impacto isto tem nas medidas de qualidade de regressores? Que impacto isto tem nas medidas de classificadores? Que informação ao cliente / usuário podemos fornecer com base neste conceito?

113

© Prof. Emilio Del Moral Hernandez

*Lembrete ... Para que  
queremos a rede  
neural com pesos  
otimizados?*

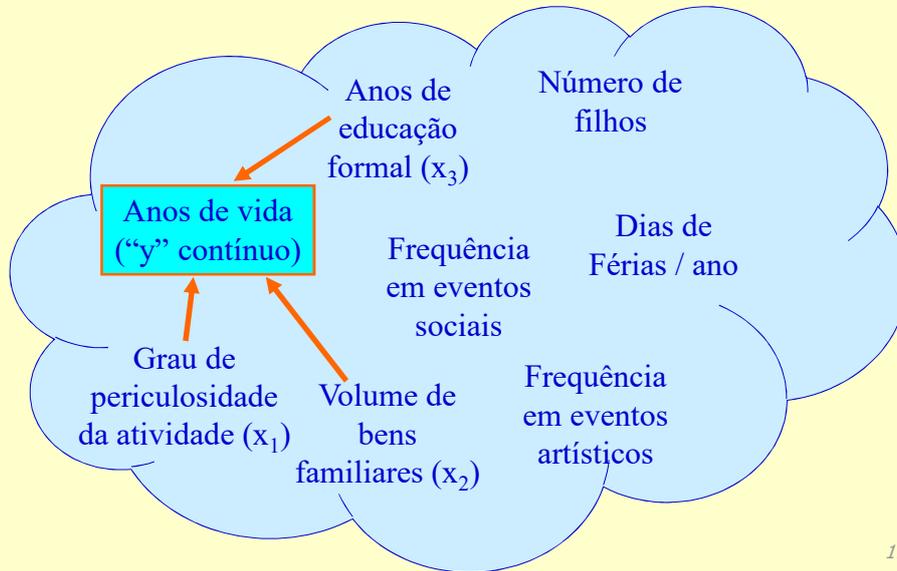
**Modelagem de um sistema por função de mapeamento  $X \rightarrow y$   
(a RNA como regressor contínuo não linear multivariável)**



**Assumimos que a variável  $y$  do sistema a modelar é uma função (normalmente desconhecida e possivelmente não linear) de diversas outras variáveis desse mesmo sistema**

**A RNA, para ser um bom modelo do sistema, deve reproduzir essa relação entre  $X$  e  $y$ , tão bem quanto possível**

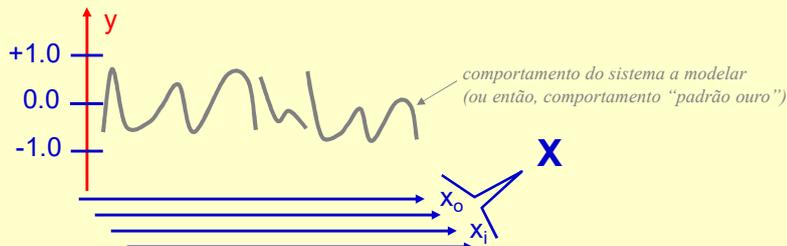
## Um hipotético universo de variáveis interdependentes, passível de modelagem/ens



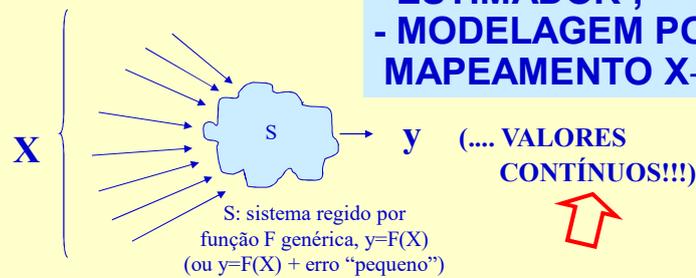
116

© Prof. Emilio Del Moral – EPUSP

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



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

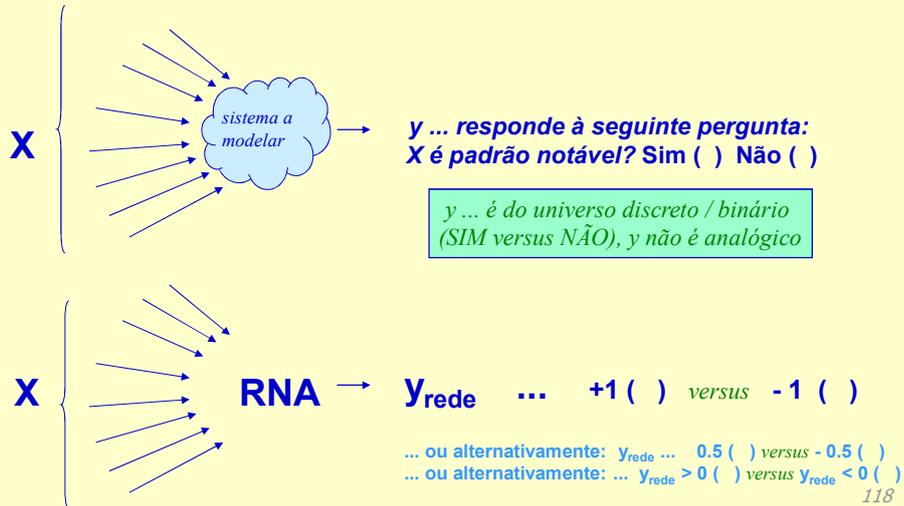


117

© Prof. Emilio Del Moral – EPUSP

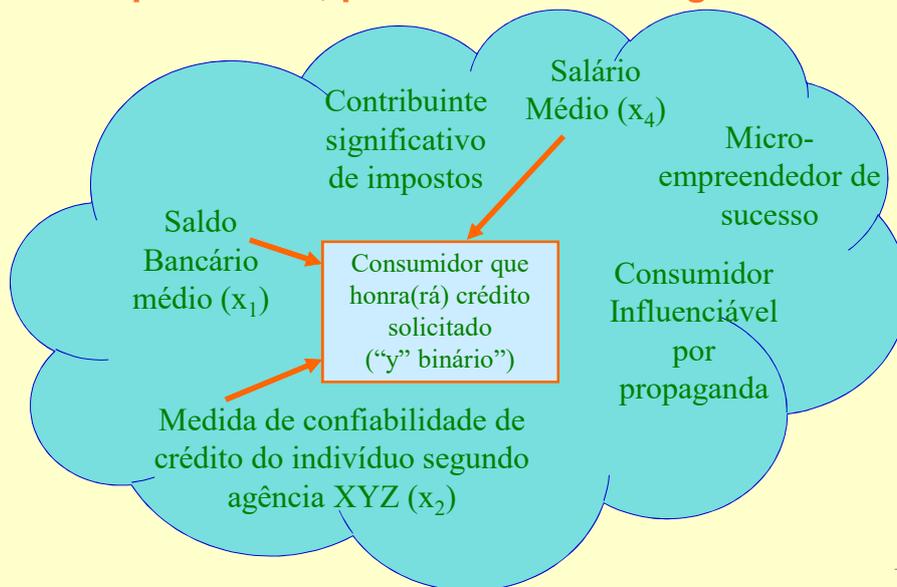
## RNAs como reconhecedor / detetor de padrões

...



© Prof. Emilio Del Moral – EPUSP

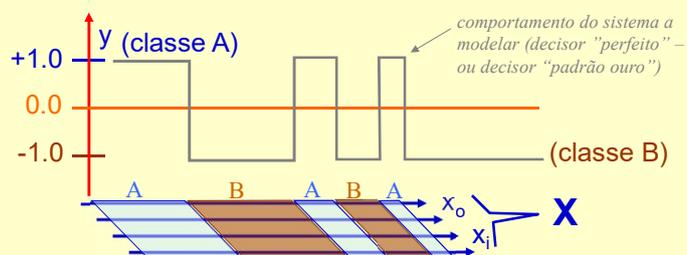
## Um hipotético universo de variáveis interdependentes, passível de modelagem/ens



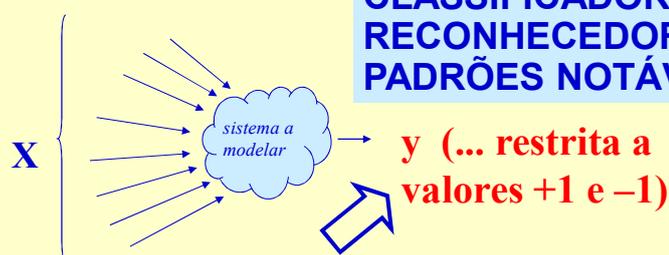
119

© Prof. Emilio Del Moral – EPUSP

## Caso de classificação binária / reconhecimento de padrões, será do tipo ...



**CLASSIFICADOR;  
RECONHECEDOR DE  
PADRÕES NOTÁVEIS**



120

© Prof. Emilio Del Moral – EPUSP

121

*E será que conseguiremos que a RNA calcule essas duas funções de X complicadas representadas nos slides anteriores (uma com y contínuo e a outra com y binário)?*

© Prof. Emilio Del Moral Hernandez

121

recordando

*De onde vem o grande poder do MLP?*

Kurt Hornik showed in 1991<sup>[2]</sup> that it is not the specific choice of the  $\varphi$  assumed to be linear. For notational convenience, only the single out

## Formal statement [\[edit\]](#)

recordando

The theorem<sup>[2][3][4][5]</sup> in mathematical terms:

Let  $\varphi(\cdot)$  be a nonconstant, bounded, and monotonically-increasing function in  $C(I_m)$  and  $\epsilon > 0$ , there exist an integer  $N$  and real constants  $\alpha_i$ ,  $w_i$

$$F(x) = \sum_{i=1}^N \alpha_i \varphi(w_i^T x + b_i)$$

as an approximate realization of the function  $f$  where  $f$  is independent of  $x$

$$|F(x) - f(x)| < \epsilon$$

for all  $x \in I_m$ . In other words, functions of the form  $F(x)$  are dense in  $C(I_m)$ .

Kurt Hornik showed in 1991<sup>[2]</sup> that it is not the specific choice of the  $\varphi$  assumed to be linear. For notational convenience, only the single out

## Formal statement [edit]



The theorem<sup>[2][3][4][5]</sup> in mathematical terms:

$Y_{rede}(X)$

$X$

Let  $\varphi(\cdot)$  be a nonconstant, bounded, and monotonically-increasing function in  $C(I_m)$  and  $\epsilon > 0$ , there exist an integer  $N$  and real constants  $\alpha_i$  and  $b_i$

$$F(x) = \sum_{i=1}^N \alpha_i \varphi(w_i^T x + b_i)$$

número de nós escondidos

sigmoidal

viés; : viés do nó escondido  $i$

$W_i$ : vetor de pesos do nó escondido  $i$

elementos do vetor de pesos do nó linear de saída  $W_s$

$$|F(x) - f(x)| < \epsilon$$

for all  $x \in I_m$ . In other words, functions of the form  $F(x)$  are dense in  $C(I_m)$ .

Kurt Hornik showed in 1991<sup>[2]</sup> that it is not the specific choice of the  $\varphi$  assumed to be linear. For notational convenience, only the single out

## Formal statement [edit]



The theorem<sup>[2][3][4][5]</sup> in mathematical terms:

Let  $\varphi(\cdot)$  be a nonconstant, bounded, and monotonically-increasing function in  $C(I_m)$  and  $\epsilon > 0$ , there exist an integer  $N$  and real constants  $\alpha_i$  and  $b_i$

$Y_{rede}(X)$

Fescondida\_sistema(X)

$$F(x) = \sum_{i=1}^N \alpha_i \varphi(w_i^T x + b_i)$$

as an approximate realization of the function  $f$  where  $f$  is independent of  $x$

Limite de erro

$$|F(x) - f(x)| < \epsilon$$

for all  $x \in I_m$ . In other words, functions of the form  $F(x)$  are dense in  $C(I_m)$ .

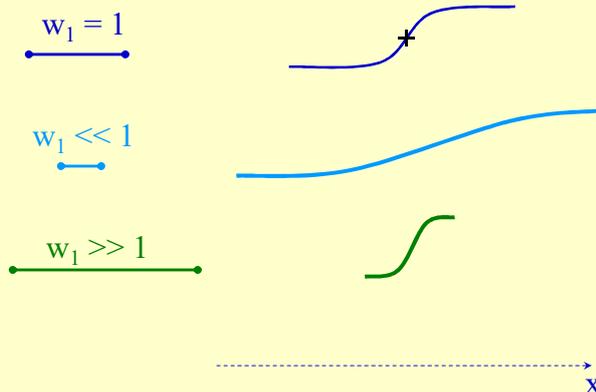
### **Questões intrigantes, p/ esta aula e p/ pensar em casa ...**

- *No que impacta escolhermos o “epsilon” de Cybenko de alto valor? O que muda na estrutura de Cybenko com isso?*
- *No que impacta escolhermos o “epsilon” de Cybenko de baixo valor?*
- *Como definimos o número de nós da primeira camada do MLP? Isto pode ser definido a priori, antes de testar o seu desempenho? (por exemplo com base no número de entradas da rede e/ou com base no número de exemplares de treino  $M$ ?)*
- *O que ganhamos e o que perdemos se escolhermos usar POUCOS nós na construção rede neural?*
- *O que ganhamos e o que perdemos se escolhermos usar MUITOS nós na construção da rede neural?*

**... Um parênteses em lousa para discutirmos um pouco a aproximação universal com sigmoides / funções em formato de “S”, num caso simples e bem particular ... aproximação de uma função  $y$  de uma única variável  $x_1$ :**

$$y(x_1)$$

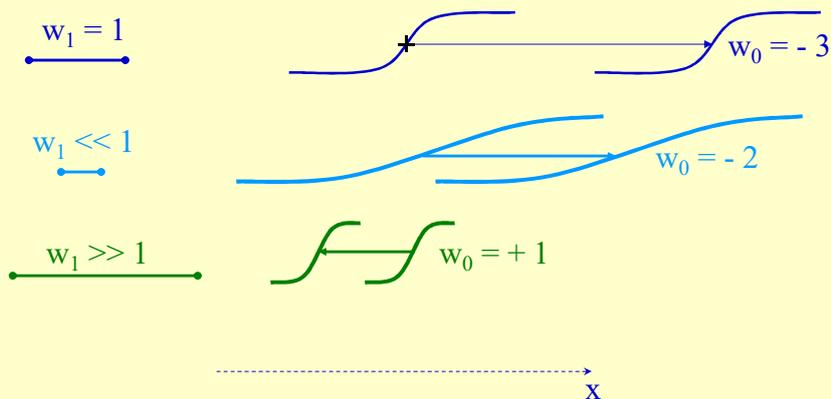
O que conseguimos fazer com **um único neurônio sigmoidal**  $y(w_1 \cdot x_1)$ , c/ escalamento de  $x_1$  via  $w_1$



130

Métodos Numéricos e Reconhecimento de Padrões – © Prof. Emilio Del Moral – EPUSP

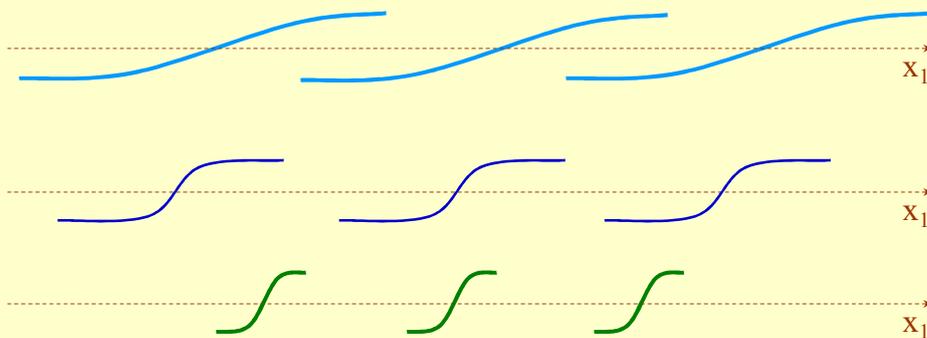
O que conseguimos fazer com **um único neurônio sigmoidal**  $y(w_1 \cdot x_1 + w_0 \cdot 1)$ , c/ escalamento de  $x_1$  via  $w_1$  ... e também com o viés, via viés  $w_0$



131

Métodos Numéricos e Reconhecimento de Padrões – © Prof. Emilio Del Moral – EPUSP

O que conseguimos fazer com **um único neurônio sigmoidal**, no caso de regressões (“y(x<sub>1</sub>) contínuo”)?



$$y = \text{tgh}(w_1 \cdot x_1 + \text{viés})$$

132

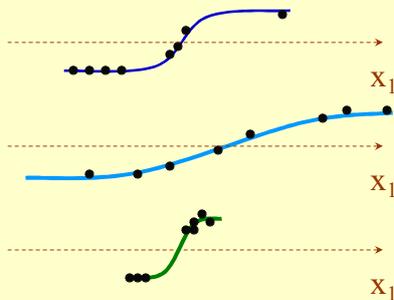
Métodos Numéricos e Reconhecimento de Padrões – © Prof. Emilio Del Moral – EPUSP

O que conseguimos fazer com **um único neurônio sigmoidal**, no caso de regressões (“y contínuo”)?

133

Métodos Numéricos e Reconhecimento de Padrões – © Prof. Emilio Del Moral – EPUSP

Que tipo de dados empíricos modelamos com **um único neurônio sigmoidal** em regressões (“ $y(x_1)$  contínuo”)?



Os pontos pretos são pares empíricos  $(x^i, y^i)$ ; As curvas coloridas, são regressões sigmoidais aderentes a tais pares.

134

Métodos Numéricos e Reconhecimento de Padrões – © Prof. Emilio Del Moral – EPUSP

Em termos de Excel, teríamos ...

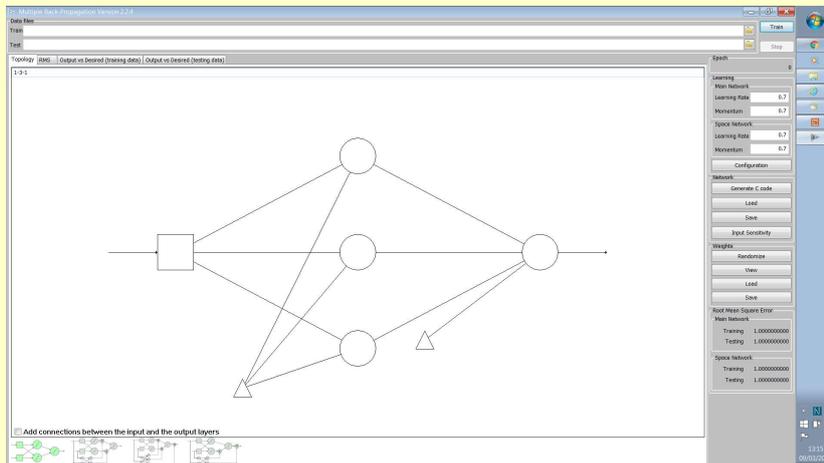
Cliente ( $\mu$ )	Idade ( $x_1$ )	Renda ( $x_2$ )	Clics ( $x_3$ )	Consumo do Produto B ( $x_4$ )	Consumo do Produto C ( $x_5$ )	Consumo do Produto A ( $y$ )
			302	958	136	9800
			186	985	196	8760
						520
M-2	16					11640
M-1	30					9640
M	19					5320

Equivalente em txt  
Para uso do MBP

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

75  
Moral – EPUSP

## Regressão univariada com Cybenko “café com leite” de 3 nós na primeira camada ...

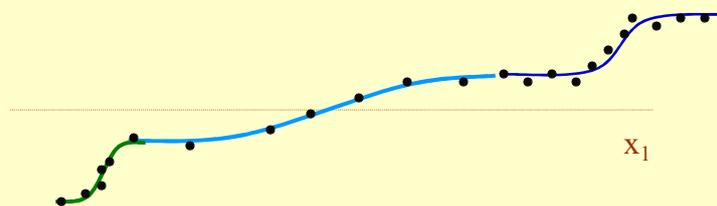


136

Métodos Numéricos e Reconhecimento de Padrões – © Prof. Emilio Del Moral – EPUSP

Cybenko “café com leite” (regressão genérica univariada), para aproximação universal de funções de 1 variável  $x_1$  apenas?

... superposição de várias sigmóides deslocadas e escaladas



Vocês enxergam acima 3 nós “tgh” na primeira camada, com com 3 viéses distintos e 3 escaladores de  $x_1$  distintos, e mais um 4o nó combinador (somatória simples de 3 entradas) na camada de saída?

137

Métodos Numéricos e Reconhecimento de Padrões – © Prof. Emilio Del Moral – EPUSP

Cybenko “café com leite” (regressão genérica univariada), para aproximação universal de funções de 1 variável  $x_1$  apenas?

... várias sigmóides deslocadas e escaladas



Acima ... Comportamento individual dos 3 nós “tgh” na primeira camada, com 3 viéses distintos e 3 escaladores de  $x_1$  distintos.

138

Métodos Numéricos e Reconhecimento de Padrões – © Prof. Emilio Del Moral – EPUSP

139

### Questões intrigantes, p/ esta aula e p/ pensar em casa ...

- No que impacta escolhermos o “epsilon” de Cybenko de alto valor? O que muda na estrutura de Cybenko com isso?
- No que impacta escolhermos o “epsilon” de Cybenko de baixo valor?
- Como definimos o número de nós da primeira camada do MLP? Isto pode ser definido a priori, antes de testar o seu desempenho? (por exemplo com base no número de entradas da rede e/ou com base no número de exemplares de treino  $M$ )?
- O que ganhamos e o que perdemos se escolhermos usar POUCOS nós na construção rede neural?
- O que ganhamos e o que perdemos se escolhermos usar MUITOS nós na construção da rede neural?

recordando

© Prof. Emilio Del Moral Hernandez

139

