

Temas da 3ª semana – PSI3471-2020 – Prof Emilio

59

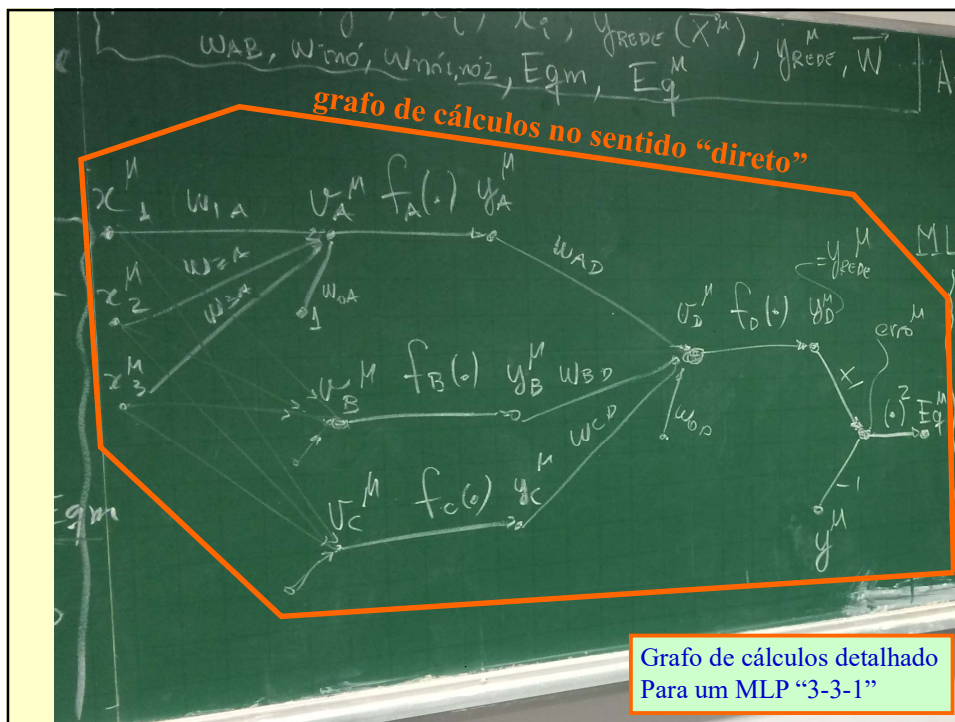
#5 (16/março – 2ªf) Foco da semana: aprendizado da Rede Neural MLP – O Gradiente descendente e a otimização de pesos sinápticos com base no conjunto de treino e EBP; dedução das fórmulas do EBP, em sala de aula em conjunto com os alunos: trabalho focado num peso sináptico específico da rede, escolhido pelo professor para máxima complexidade da dedução.

#6 (18/março – 4ªf) ... Discussão das extensões das deduções já feitas (para um peso no EBP) para os demais pesos sinápticos; redundâncias nos cálculos dos diversos pesos da rede neural e otimização do esforço computacional. Regra “Delta” de aprendizado de Widrow, para neurônio isolado; Aprendizado por EBP recursivo, camada a camada.

© Prof. Emilio Del Moral Hernandez

59

59



60

*Como escolhemos
os valores dos
diversos w 's ?*

61

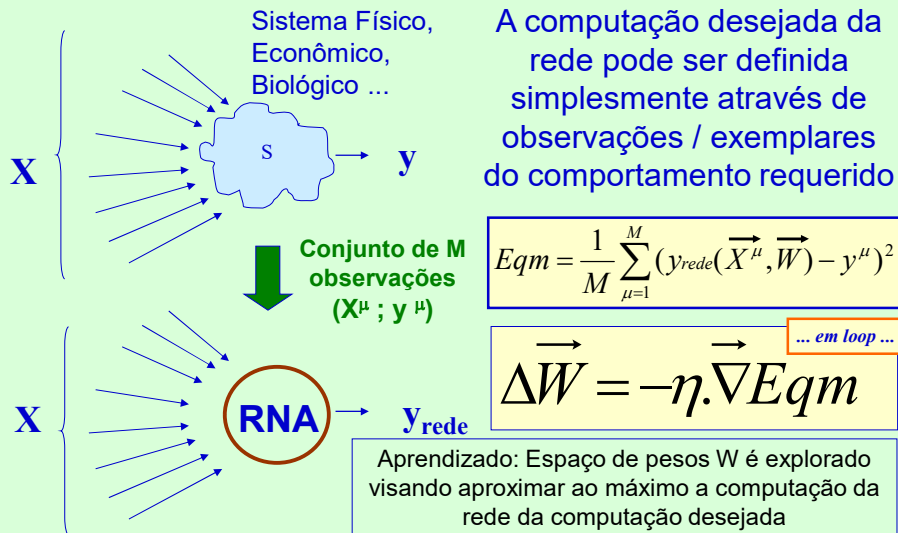
*Aprendizado em RNAs do
tipo MLP – Multi Layer
Perceptron – através do
algoritmo Error Back
Propagation*

(método do gradiente usado na otimização de w 's do MLP)

62

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

63



© Prof. Emilio Del Moral Hernandez

63

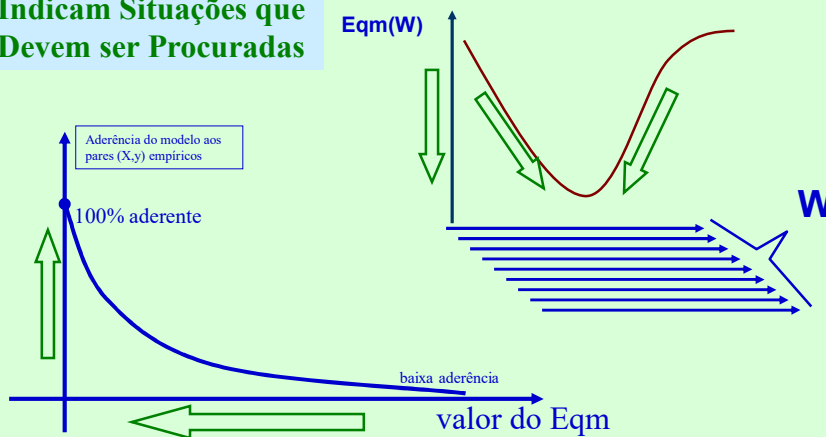
63

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

64

Devemos buscar Maximização da aderência = Mínimo Eqm possível

As Setas Verdes Indicam Situações que Devem ser Procuradas



64

64

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

- “Chute” um W inicial para o “ W corrente”, ou “ W melhor até agora”
- Em loop até obter E_{qm} zero, ou baixo o suficiente, ou estável:
 - Determine o vetor gradiente do E_{qm} , nesse espaço de W s
 - Em loop varrendo todos os M exemplos $(X^{\mu}; y^{\mu})$,
 - Calcule o gradiente de $E_{q^{\mu}}$ associado a um exemplo μ , e vá varrendo μ e somando os gradientes de cada $E_{q^{\mu}}$, para compor o vetor gradiente de E_{qm} , 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
 - Tire a média dos M gradientes individuais e dê um passo Delta ΔW nesse espaço, com direção e magnitude dados por $-\eta$ *vetor gradiente (E_{qm})

Prof. Emilio Del Moral Hernandez

65

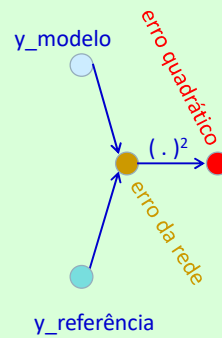
Processo de refinamentos graduais a cada iteração ...

$W\#0$	$E_{qm}\#0$	$GradE_{qm}(W\#0)$	$\Delta W\#0 = -n \cdot GradE_{qm}(W\#0)$
$W\#1$ (= $W\#0 + \Delta W\#0$)	$E_{qm}\#1$ (< $E_{qm}\#0$)	$GradE_{qm}(W\#1)$	$\Delta W\#1 = -n \cdot GradE_{qm}(W\#1)$
$W\#2$ (= $W\#1 + \Delta W\#1$)	$E_{qm}\#2$ (< $E_{qm}\#1$)	$GradE_{qm}(W\#2)$	$\Delta W\#2 = -n \cdot GradE_{qm}(W\#2)$
$W\#3$ (= $W\#2 + \Delta W\#2$)	$E_{qm}\#3$ (< $E_{qm}\#2$)	$GradE_{qm}(W\#3)$	$\Delta W\#3 = -n \cdot GradE_{qm}(W\#3)$
$W\#4$ (= $W\#3 + \Delta W\#3$)	$E_{qm}\#4$ (< $E_{qm}\#3$)	$GradE_{qm}(W\#4)$	$\Delta W\#4 = -n \cdot GradE_{qm}(W\#4)$
...
$W\#k$ (= $W\#k-1 + \Delta W\#k-1$)	$E_{qm}\#k$ (< $E_{qm}\#k-1$)	$GradE_{qm}(W\#k)$	$\Delta W\#k = -n \cdot GradE_{qm}(W\#k)$
...

Prof. Emilio Del Moral Hernandez

66

Modelo Neural no modo de Treinamento



Prof. Emilio Del Moral Hernandez

74

... erro da rede com relação ao conjunto de treinamento como um todo;
simbologia (X^μ ; y^μ); Erro quadrático de exemplar (Eq^μ); Erro quadrático
médio (Eqm)

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

μ identifica um de M exemplos de treinamento

Prof. Emilio Del Moral Hernandez

75

Deduzindo as Equações do Aprendizado em RNAs do tipo MLP – Multi Layer Perceptron – com o algoritmo Error Back Propagation (Gradiente Descendente)

© Prof. Emilio Del Moral – EPUSP

81

“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:

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

- 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

82

83

**Faça sua própria
revisão com base
no aprendizado
anterior !!**

(ou aprenda por si, se
saltou tópicos de
cursos anteriores como
cálculos, em que esses
temas são vistos e/ou
usados regularmente)

© Prof. Emilio Del Moral Hernandez

83

Derivada parcial - ilustração p/ função de 2 variáveis apenas

imagem exemplo extraída da internet

84

**Faça sua própria
revisão com base
no aprendizado
anterior !!**

(ou aprenda por si, se
saltou tópicos de
cursos anteriores como
cálculos, em que esses
temas são vistos e/ou
usados regularmente)

© Prof. Emilio Del Moral – EPUSP

84

Derivada parcial- ilustração p/ função de 2 variáveis apenas

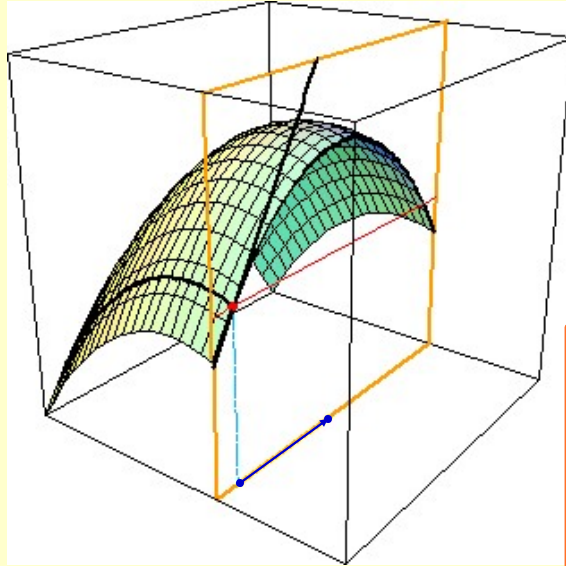


imagem base extraída da internet

Faça sua própria revisão com base no seu aprendizado anterior !!

(ou aprenda por si, se saltou tópicos de cursos anteriores como cálculos, em que esses temas são vistos e/ou usados regularmente)

© Prof. Emilio Del Moral – EPUSP

85

Formação do vetor gradiente a partir de duas derivadas parciais

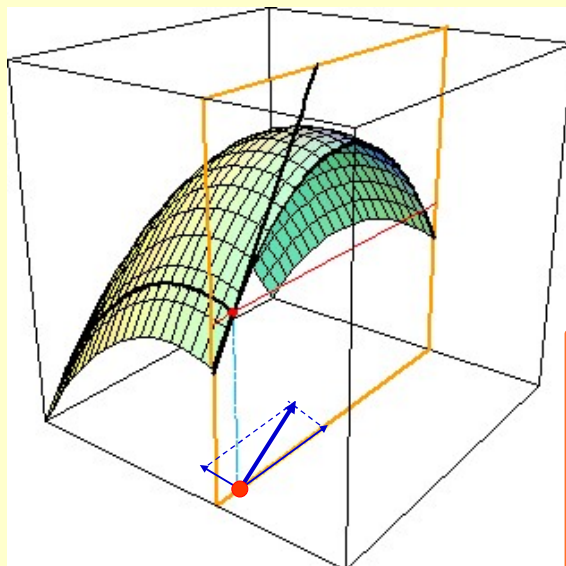


imagem base extraída da internet

Faça sua própria revisão com base no seu aprendizado anterior !!

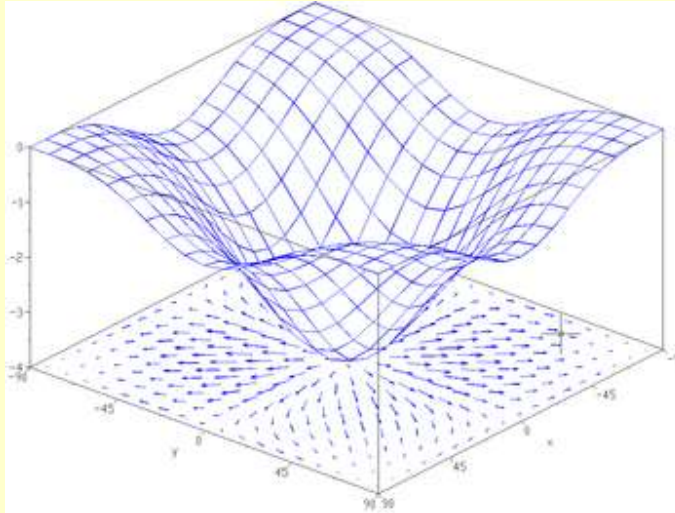
(ou aprenda por si, se saltou tópicos de cursos anteriores como cálculos, em que esses temas são vistos e/ou usados regularmente)

© Prof. Emilio Del Moral – EPUSP

86

<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



Faça sua própria revisão com base no seu aprendizado anterior !!

(ou aprenda por si, se saltou tópicos de cursos anteriores como cálculos, em que esses temas são vistos e/ou usados regularmente)

© Prof. Emilio Del Moral – EPUSP

87

Aprendizado do MLP por Error Back Propagation ...

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

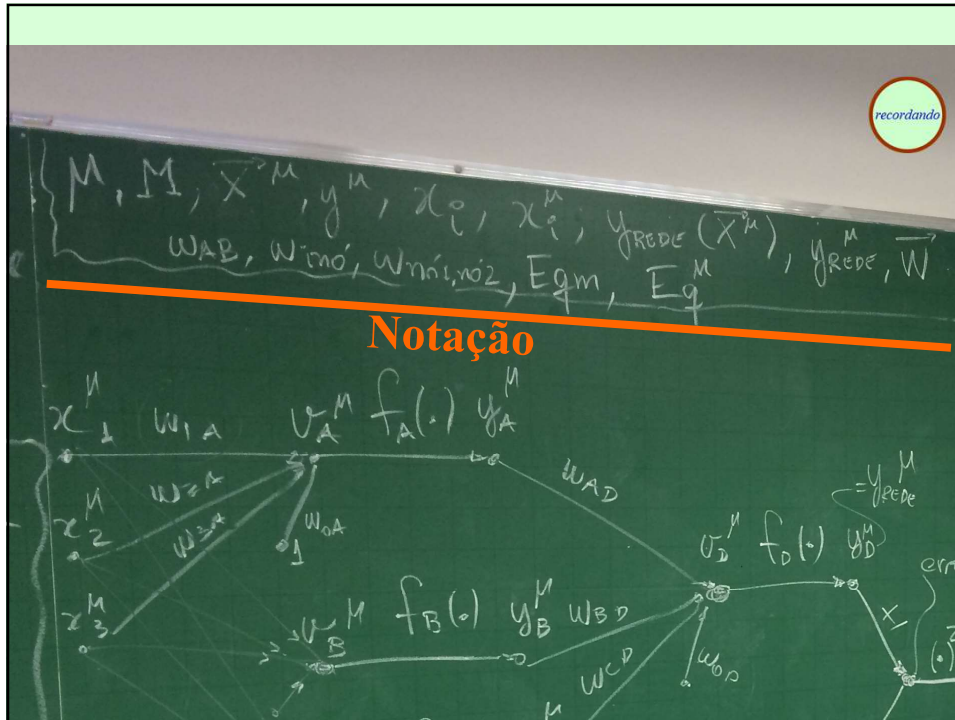
Gradiente de E_{qm} no espaço de pesos = $(\partial E_{qm}(W)/\partial w_1, \partial E_{qm}(W)/\partial w_2, \partial E_{qm}(W)/\partial w_3, \dots)$

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

88

© Prof. Emilio Del Moral – EPUSP

88



91

Entendendo símbolos que temos usado em nossos grafos da lousa e em alguns dos slides:

92

- X -- vetor X de entradas num MLP, com os valores genéricos nas suas componentes x_i
- x_i -- componente "i" do vetor de entradas de um MLP, com valor genérico, um número real entre " - infinito e + infinito"
- X^μ -- vetor X das entradas num MLP, mas com os valores de cada um dos x_i específicos da observação empírica μ
- μ -- identificador / indexador inteiro, entre 1 e M , usado para especificar um dos exemplares (uma das observações empíricas, aquela de "número" μ) que compõem o conjunto de pares empíricos (X, y) .
- M -- número total de exemplares empíricos que compõem o conjunto de treino; cardinalidade do conjunto de treino
- x_i^μ -- componente "i" do vetor de entradas de um MLP, com o valor específico referente à observação empírica μ
- y^μ -- saída alvo para a rede neural sendo treinada com aprendizado supervisionado, com o valor específico da observação empírica μ . Poderia ser chamada também de y_{alvo}^μ , explicitando mais claramente o significado

© Prof. Emilio Del Moral Hernandez

92

92

Entendendo símbolos que temos usado em nossos grafos da lousa e em alguns dos slides:

93

$y_{rede}(X^\mu)$ -- saída da rede neural quando a sua entrada corresponde aos valores empíricos X^μ

Ou ... y_{rede}^μ -- outra forma de representar $y_{rede}(X^\mu)$

W -- vetor de todos os pesos sinápticos de uma rede MLP, incluindo todos os vieses de todos os neurônios

w_{AB} -- peso sináptico específico que conecta a saída do neurônio A com uma das entradas do neurônio B

w_{iA} -- peso sináptico específico que conecta a entrada x_i da rede neural com o neurônio A da primeira camada do MLP

Eq_m -- Erro quadrático médio; média dos M valores de Eq^μ

Eq^μ -- Erro quadrático individual, referente especificamente à observação empírica μ ----- fórmula: $Eq^\mu = [y_{rede}(X^\mu) - y^\mu]^2$

© Prof. Emilio Del Moral Hernandez

93

93

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

95

© Prof. Emilio Del Moral – EPUSP

95

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. 96

© Prof. Emilio Del Moral – EPUSP

96

Exemplo de dados empíricos tabulados em Excel ...

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

97

© Prof. Emilio Del Moral – EPUSP

97

Exemplo de dados empíricos tabulados em Excel ...

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

M vetores X^μ (cada um deles é 5 dimensional) de entrada do MLP, referentes cada um deles a uma das M observações empíricas

M alvos y^μ

98

© Prof. Emilio Del Moral – EPUSP

98

Equivalente em .txt, em formato apropriado para o ambiente Multiple Back Propagation ...

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

Equivalente em txt
Para uso do MBP

```

treino em txt para exemplo de consumo A e B - Bloco de notas
Arquivo Editar Formatar Exibir Ajuda
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

99

... erro da rede com relação ao conjunto de treinamento como um todo; simbologia (X^μ ; y^μ); Erro quadrático de exemplar (Eq^μ); Erro quadrático médio (Eqm)

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

μ identifica um de M exemplos de treinamento

Prof. Emilio Del Moral Hernandez

102

Inicialmente, invertamos o operador gradiente e a somatória

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

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

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

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

103

© Prof. Emilio Del Moral – EPUSP

103

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

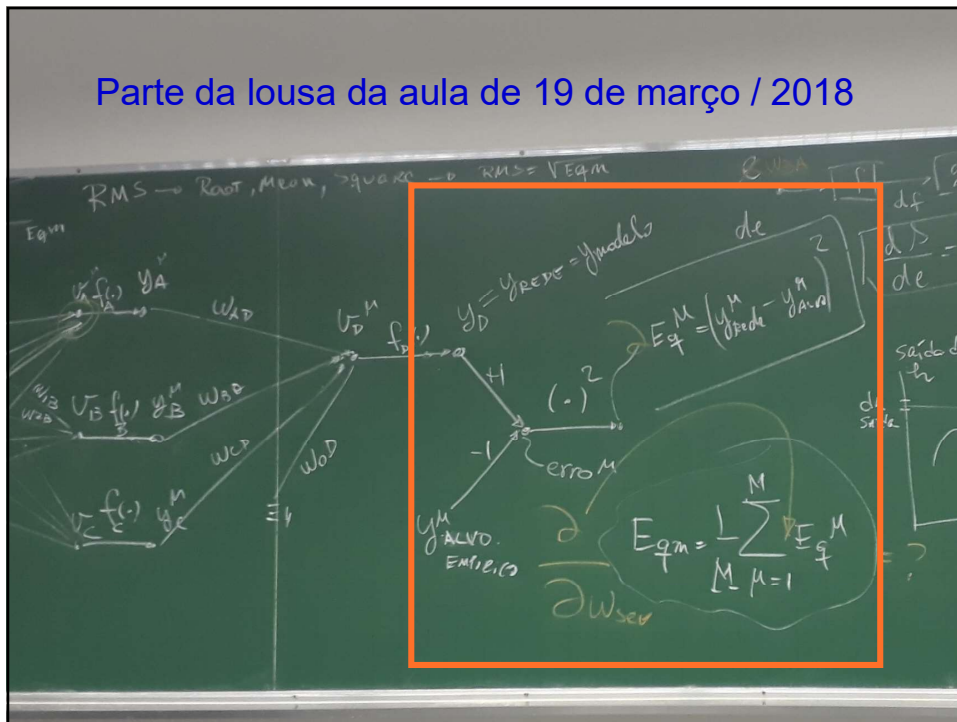
$$d(f_1(x)+f_2(x)) / dx = df_1(x)/dx + df_2(x)/dx$$

104

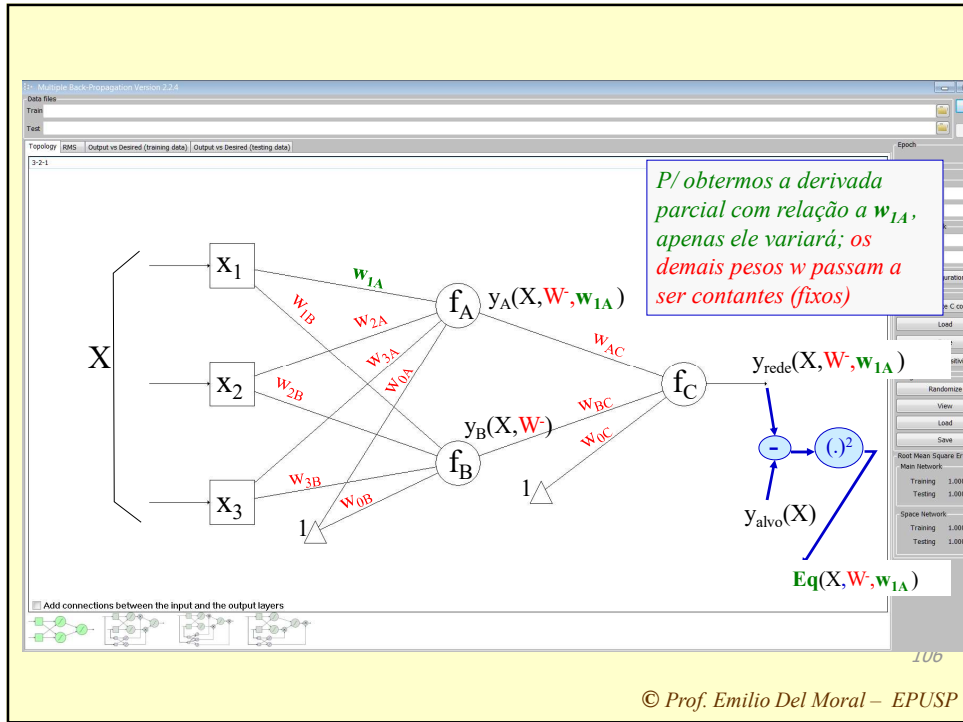
© Prof. Emilio Del Moral – EPUSP

104

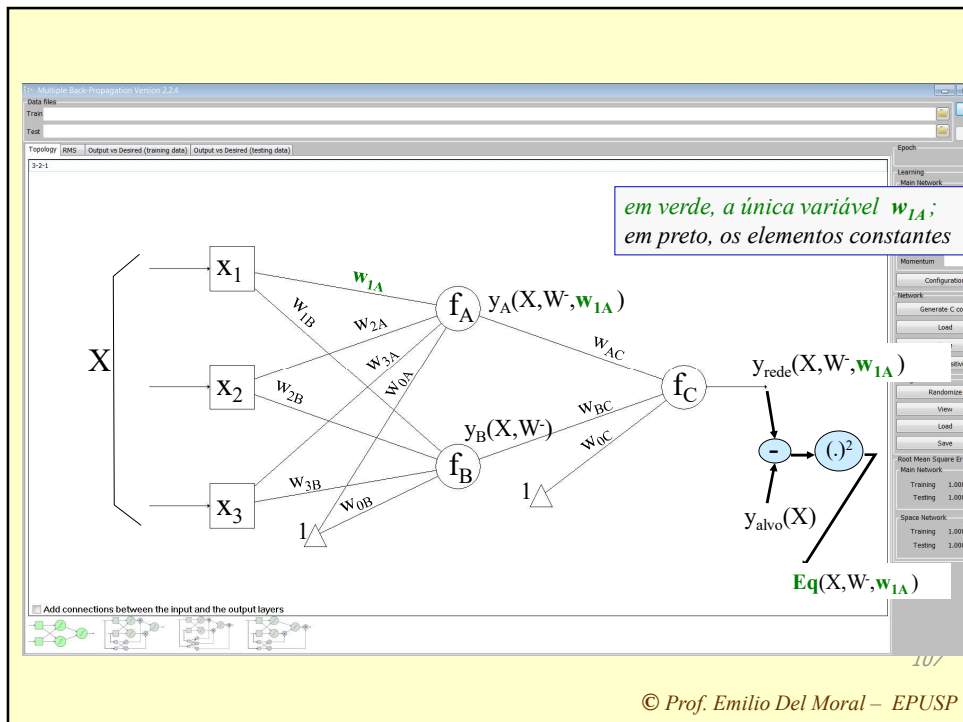
Parte da lousa da aula de 19 de março / 2018



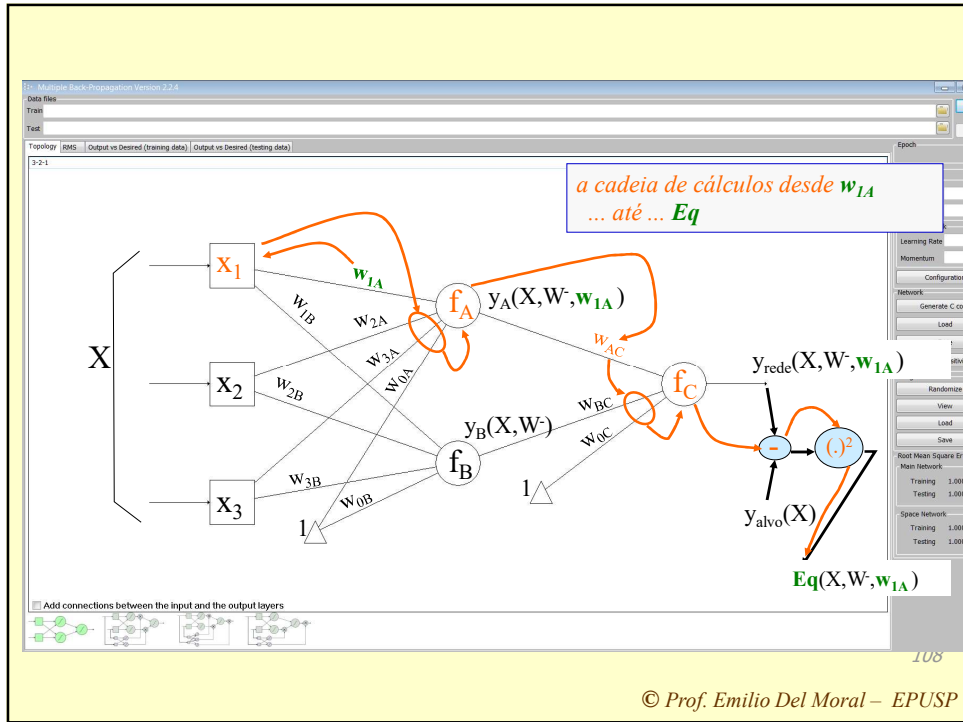
105



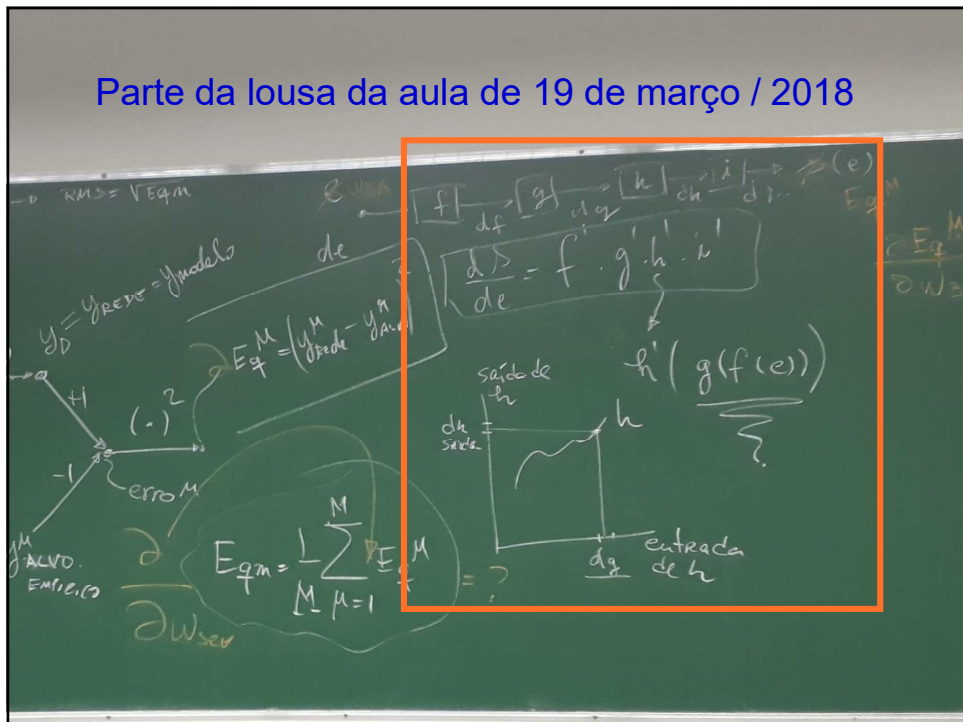
106



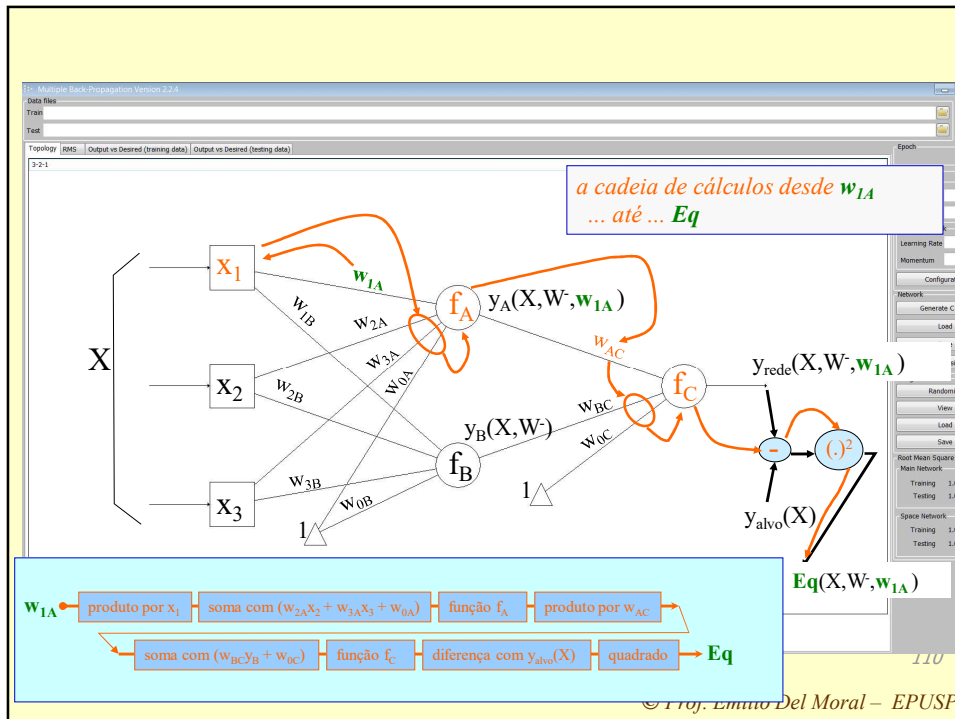
107



108



109



110

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 ...

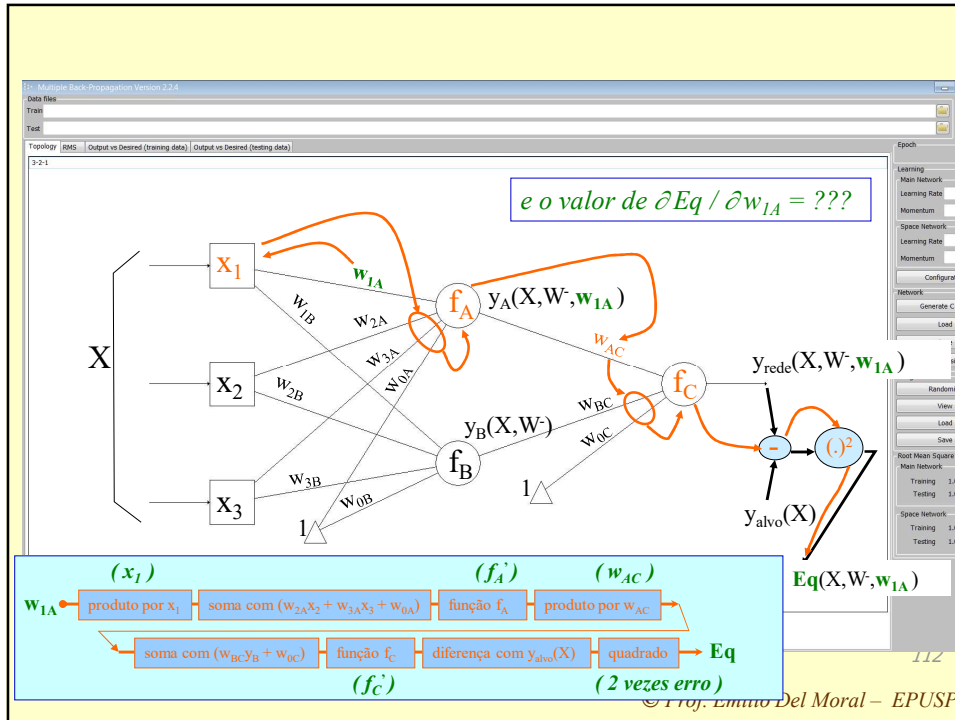
$$d(f_1(f_2(x))) / dx = df_1(x)/df_2 \cdot df_2(x)/dx$$

..., 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.

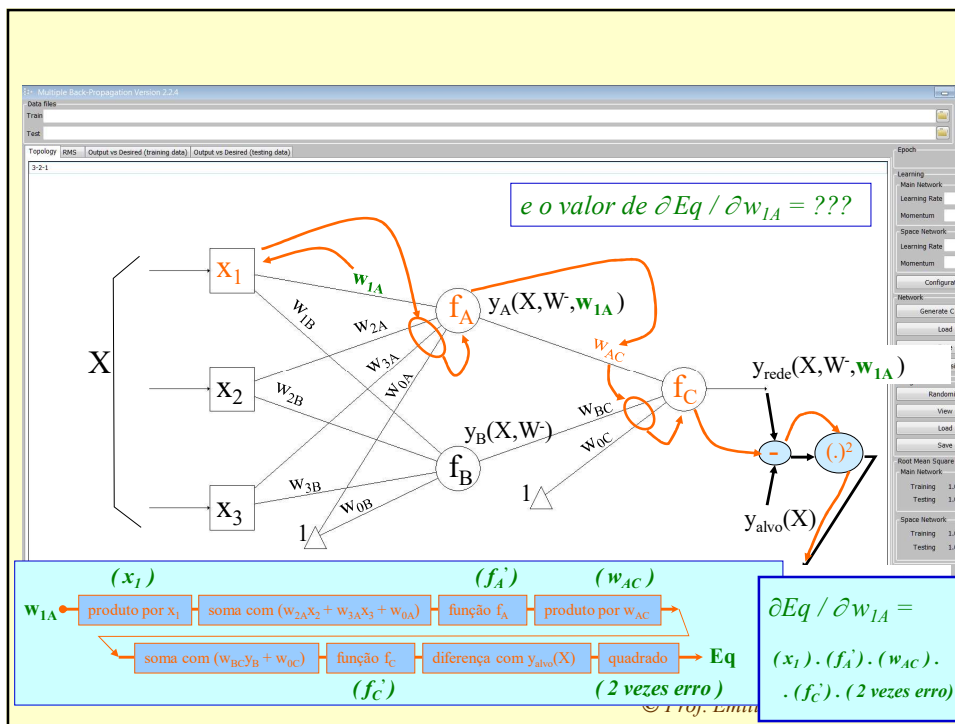
111

© Prof. Emilio Del Moral – EPUSP

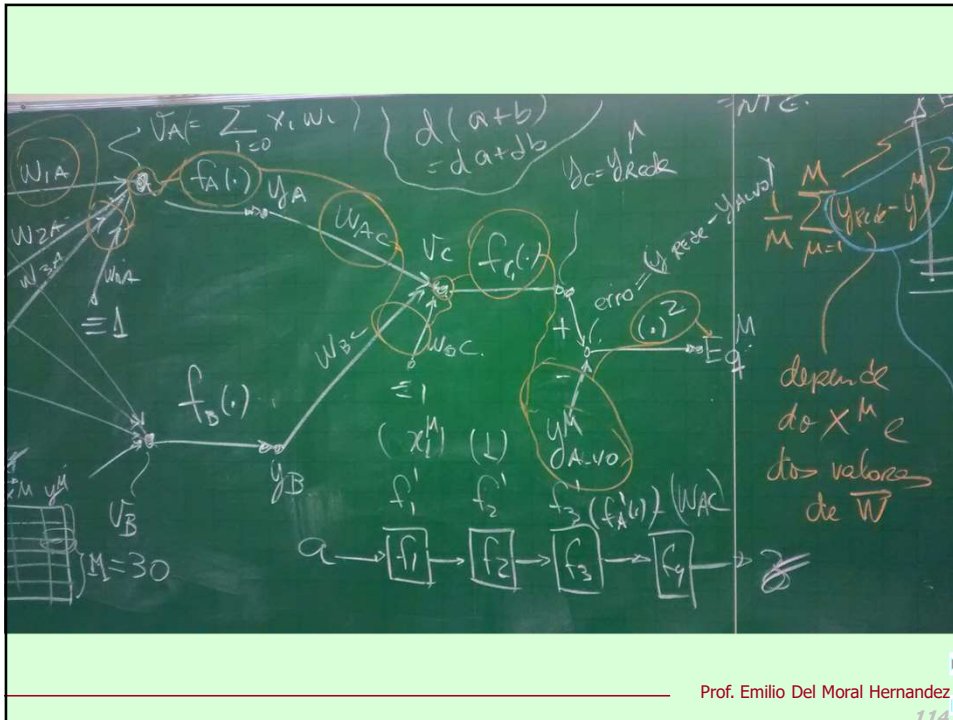
111



112



113



Prof. Emilio Del Moral Hernandez

114

Experimentem deduzir a $\partial Eq / \partial w_{AC}$ - w_{AC} na 2ª camada

e o valor de $\partial Eq / \partial w_{1A} = ???$

(x_1) (f_A) (w_{AC})
 w_{1A} → produto por x_1 → soma com $(w_{2A}x_2 + w_{3A}x_3 + w_{0A})$ → função f_A → produto por w_{AC} →
 soma com $(w_{BC}y_B + w_{0C})$ → função f_C → diferença com $y_{alvo}(X)$ → quadrado → Eq
 (f_C') (2 vezes erro)

$\partial Eq / \partial w_{1A} =$
 $(x_1) \cdot (f_A') \cdot (w_{AC}) \cdot$
 $(f_C') \cdot (2 \text{ vezes erro})$

118

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 $^\mu$
- 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 : $f'_C(w_{AC} \cdot f'_A(w_{1A} \cdot x_1 + w_{2A} \cdot x_2 + \dots + w_{0A}) + w_{BC} \cdot f'_B(w_{1B} \cdot x_1 + w_{2B} \cdot x_2 + \dots + w_{0B}) + w_{0C})$
 - 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 $[y_{\text{rede}}(X, W) - y_{\text{alvo}}(X)]$

123

© Prof. Emilio Del Moral – EPUSP

123

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^μ e y_{alvo}^μ , 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!

124

© Prof. Emilio Del Moral – EPUSP

124

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

126

- “Chute” um W inicial para o “ $W_{corrente}$ ”, ou “ W melhor até agora”
- Em loop até obter E_{qm} zero, ou baixo o suficiente, ou estável:
 - Determine o vetor gradiente do E_{qm} , nesse espaço de W s
 - Em loop varrendo todos os M exemplos $(X^{\mu}; y^{\mu})$,
 - Calcule o gradiente de $E_{q^{\mu}}$ associado a um exemplo μ , e vá varrendo μ e somando os gradientes de cada $E_{q^{\mu}}$, para compor o vetor gradiente de E_{qm} , 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
 - Tire a média dos M gradientes individuais e dê um passo Delta ΔW nesse espaço, com direção e magnitude dados por $-\eta \cdot \text{vetor gradiente} (E_{qm})$

Prof. Emilio Del Moral Hernandez

126