

Um tema importante no curso:

Sobreapredizado / Sobreajuste

Conceito, entendimento da sua origem e formas de limitá-lo

© Prof. Emilio Del Moral – EPUSP

recordando

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?*

© Prof. Emilio Del Moral Hernandez

34



A aproximação universal com RNAs do tipo MLP, segundo Cybenko (& Kolmogorov)

© Prof. Emilio Del Moral Hernandez

Cybenko – a prova matemática, disponível para download na internet, é bastante complexa



Math. Control Signals Systems (1989) 2: 303-314
 Mathematics of Control, Signals, and Systems
 © 1989 Springer-Verlag New York Inc.

310

G. Cybenko

313

Approximation by Superpositions of a Sigmoidal Function*

G. Cybenko†

Abstract. In this paper we demonstrate that finite linear combinations of compositions of a fixed, univariate function and a set of affine functionals can uniformly approximate any continuous function of n real variables with support in the unit hypercube, only mild conditions are imposed on the univariate function. Our results settle an open question about representability in the class of single hidden layer neural networks. In particular, we show that arbitrary decision regions can be arbitrarily well approximated by continuous feedforward neural networks with only a single internal, hidden layer and any continuous sigmoidal nonlinearity. The paper discusses approximation properties of other possible types of nonlinearities that might be implemented by artificial neural networks.

Key words. Neural networks, Approximation, Complexity.

1. Introduction

A number of diverse application areas are concerned with the representation of general functions of an n -dimensional real variable, $x \in \mathbb{R}^n$, by finite linear combinations of the form

$$\sum_{j=1}^m \alpha_j \sigma(y_j^T x + \theta_j), \quad (1)$$

where $y_j \in \mathbb{R}^n$ and $\alpha_j, \theta_j \in \mathbb{R}$ are fixed, (y_j^T is the transpose of y_j so that $y_j^T x$ is the inner product of y_j and x). Here the univariate function σ depends heavily on the context of the application. Our major concern is with so-called sigmoidal σ 's:

$$\sigma(t) \rightarrow \begin{cases} 1 & \text{as } t \rightarrow +\infty, \\ 0 & \text{as } t \rightarrow -\infty. \end{cases}$$

Such functions arise naturally in neural network theory as the activation function of a neural node (or unit as is becoming the preferred term) [LJ], [RHM]. The main result of this paper is a demonstration of the fact that sums of the form (1) are dense in the space of continuous functions on the unit cube if σ is any continuous sigmoidal

* Data received: October 21, 1988. Date revised: February 17, 1989. This research was supported in part by NSF Grant DCR-861903, ONR Contract N000186-G-0202 and DOE Grant DE-FG02-85ER25001.
 † Center for Supercomputing Research and Development and Department of Electrical and Computer Engineering, University of Illinois, Urbana, Illinois 61801, U.S.A.

303

4. Results for Other Activation Functions

In this section we discuss other classes of activation functions that have approximation properties similar to the ones enjoyed by continuous sigmoidals. Since these other examples are of somewhat less practical interest, we only sketch the corresponding proofs.

There is considerable interest in discontinuous sigmoidal functions such as hard limiters ($\sigma(x) = 1$ for $x \geq 0$ and $\sigma(x) = 0$ for $x < 0$). Discontinuous sigmoidal functions are not used as often as continuous ones (because of the lack of good training algorithms) but they are of theoretical interest because of their close relationship to classical perceptrons and Gamma networks [MP].

Assume that σ is a bounded, measurable sigmoidal function. We have an analog of Theorem 2 that goes as follows:

Theorem 4. Let σ be a bounded measurable sigmoidal function. Then finite sums of the form

$$G(x) = \sum_{j=1}^m \alpha_j \sigma(y_j^T x + \theta_j)$$

are dense in $L^1(I_n)$. In other words, given any $f \in L^1(I_n)$ and $\epsilon > 0$, there is a sum, $G(x)$, of the above form for which

$$\|G - f\|_{L^1} = \int_{I_n} |G(x) - f(x)| dx < \epsilon.$$

The proof follows the proof of Theorems 1 and 2 with obvious changes such as replacing continuous functions by integrable functions and using the fact that $L^1(I_n)$ is the dual of $L^\infty(I_n)$. The notion of being discriminatory accordingly changes to the following: for $h \in L^\infty(I_n)$ the condition that

$$\int_{I_n} \sigma(y^T x + \theta) h(x) dx = 0$$

for all y and θ implies that $h(x) = 0$ almost everywhere. General sigmoidal functions are discriminatory in this sense as already seen in Lemma 1 because measures of the form $h(x) dx$ belong to $M(I_n)$.

Since convergence in L^1 implies convergence in measure [A], we have an analog of Theorem 3 that goes as follows:

Theorem 5. Let σ be a general sigmoidal function. Let f be the decision function for any finite measurable partition of I_n . For any $\epsilon > 0$, there is a finite sum of the form

$$G(x) = \sum_{j=1}^m \alpha_j \sigma(y_j^T x + \theta_j)$$

and a set $D \subset I_n$, so that $m(D) \geq 1 - \epsilon$ and $|G(x) - f(x)| < \epsilon$ for $x \in D$.

led are quite powerful, we that remain to be answered imation (or equivalently, imation of a given quality? y) a role in determining the suspect quite strongly that i will require astronomical dimensionality that plagues Some recent progress con- proximated and the number ound in [MSJ] and [BH], iness of the results of this : more attention.

n, Christopher Chase, Lee marov, Richard Lippmann, 'tences, and improvements

New York, 1972. uralization?, *Neural Comput.* (to

tems and control, *IEEE Control* \, Classifying learnable geometric rdings of the 18th Annual ACM p. 273-282.

and the Pompeiu problem, *Ann.* i sets using the Radon transform,

wo Hidden Layers are Sufficient, University, 1988.

of linear combinations, *SIAM J.* us mappings by neural networks,

EE Trans. Acoust. Speech Signal stforward networks are universal a Neural Net and Conventional 87.

mal Classifiers, Technical Report, -475. tworks by sigmoidal functions, a, University of Lowell, 1988.

Cybenko – Enunciado da Prova ... (premissas + resultado)

The screenshot shows the Wikipedia page for the 'Universal approximation theorem'. The title is 'Universal approximation theorem'. The text explains that in the mathematical theory of artificial neural networks, the universal approximation theorem states that a feed-forward network with a single hidden layer containing a finite number of neurons (i.e., a multilayer perceptron) can approximate continuous functions on compact subsets of \mathbb{R}^n , under mild assumptions on the activation function. The theorem thus states that simple neural networks can represent a wide variety of interesting functions when given appropriate parameters; it does not touch upon the algorithmic learnability of those parameters.

One of the first versions of the theorem was proved by George Cybenko in 1989 for sigmoid activation functions.^[2] Kurt Hornik showed in 1991^[3] that it is not the specific choice of the activation function, but rather the multilayer feedforward architecture itself which gives neural networks the potential of being universal approximators. The output units are always assumed to be linear. For notational convenience, only the single output case will be shown. The general case can easily be deduced from the single output case.

Formal statement [edit]

The theorem^{[2][3][4][5]} in mathematical terms:

Let $\varphi(\cdot)$ be a nonconstant, bounded, and monotonically-increasing continuous function. Let I_m denote the m -dimensional unit hypercube $[0, 1]^m$. The space of continuous functions on I_m is denoted by $C(I_m)$. Then, given any function $f \in C(I_m)$ and $\epsilon > 0$, there exist an integer N and real constants $\alpha_i, b_i \in \mathbb{R}, w_i \in \mathbb{R}^m$ where $i = 1, \dots, N$ such that we may define:

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

as an approximate realization of the function f where f is independent of φ ; that is,

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

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

References [edit]

- ^[1] Balazs Csornai Csajj. Approximation with Artificial Neural Networks. Faculty of Sciences, Eötvös Loránd University, Hungary
- ^[2] ^[3] ^[4] ^[5] George Cybenko. (1989) "Approximation by superpositions of sigmoid functions". *Mathematics of Control, Signals, and Systems*, 2 (4), 303-314
- ^[6] Kurt Hornik (1991) "Approximation Capabilities of Multilayer Feedforward Networks". *Neural Networks*, 4(2), 251-257
- ^[7] Haykin, Simon (1998) *Neural Networks: A Comprehensive Foundation*, Volume 2. Prentice Hall. ISBN 0-13-27350-1
- ^[8] Hornik, M. (1993) *Fundamentals of Artificial Neural Networks* MIT Press, p. 43

⚠ This applied mathematics-related article is a stub. You can help Wikipedia by expanding it.

Categories: Theorems in discrete mathematics | Artificial neural networks | Neural networks | Network architecture | Networks | Information, knowledge, and uncertainty | Applied mathematics stubs

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

This block shows a close-up of the 'Formal statement' section from the previous image, with several annotations in blue boxes:

- A box labeled $y_{rede}(X)$ points to the function $F(x)$ in the equation.
- A box labeled X points to the input vector x in the equation.
- A box labeled "número de nós escondidos" points to the summation index $i=1$ to N .
- A box labeled "sigmoidal" points to the activation function φ .
- A box labeled "viés; viés do nó escondido i " points to the bias term b_i .
- A box labeled " W_i : vetor de pesos do nó escondido i " points to the weight vector w_i^T .
- A box labeled "elementos do vetor de pesos do nó linear de saída W_s " points to the output weights α_i .

The text below the equation states: "as an approximate realization of the function f where f is independent of φ ; that is, $|F(x) - f(x)| < \epsilon$ for all $x \in I_m$. In other words, functions of the form $F(x)$ are den

Kurt Hornik showed in 1991^[2] that it is not the specific choice of the ϕ assumed to be linear. For notational convenience, only the single out

Formal statement [edit]

The theorem^{[2][3][4][5]} in mathematical terms:

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

$$F(x) = \sum_{i=1}^N a_i \phi(a_i T_m + b_i)$$

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

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

$y_{\text{rede}}(X)$

Fescondida_sistema(X)

Limite de erro

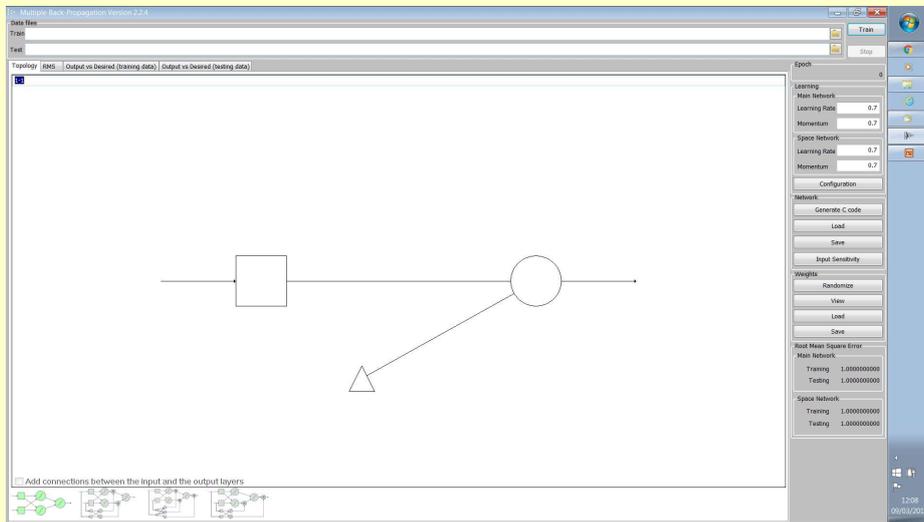
for all $x \in I_m$. In other words, functions of the form $F(x)$ are dense



... Um parênteses em lousa para discutirmos um pouco a aproximação universal com sigmoides / funções em formato de "S", no caso simples e bem particular univariado ... aproximação de uma função y de uma única variável x_1 :

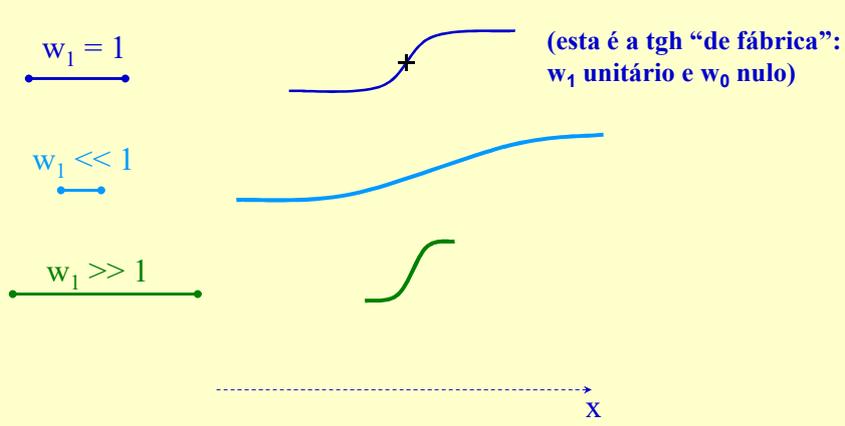
$$y(x_1)$$

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

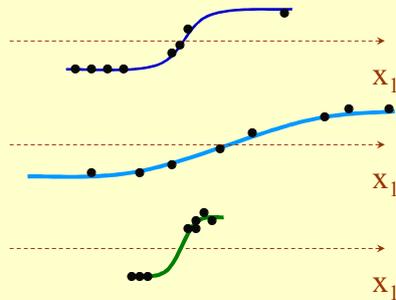


O que conseguimos fazer com **um único neurônio sigmoidal** $y(w_1 \cdot x_1)$ c/ escalamento de x_1 via w_1 e VIÉS NULO

recordando



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

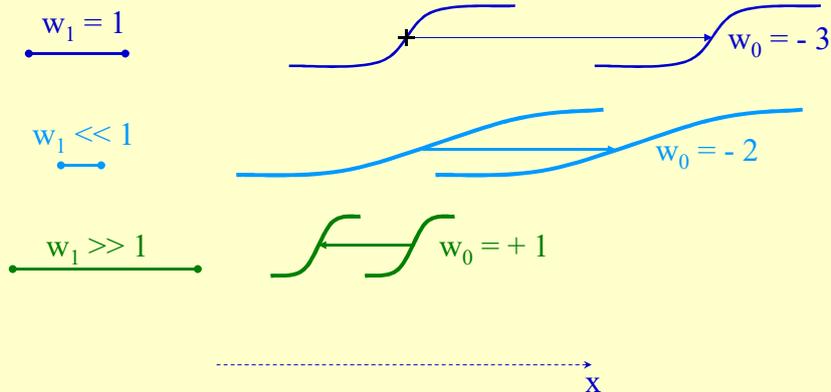


Os pontos pretos são pares empíricos (x^u, y^u) ; As curvas coloridas, são regressões sigmoidais aderentes a tais pares.

43

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 agora também com o viés, via viés w_0

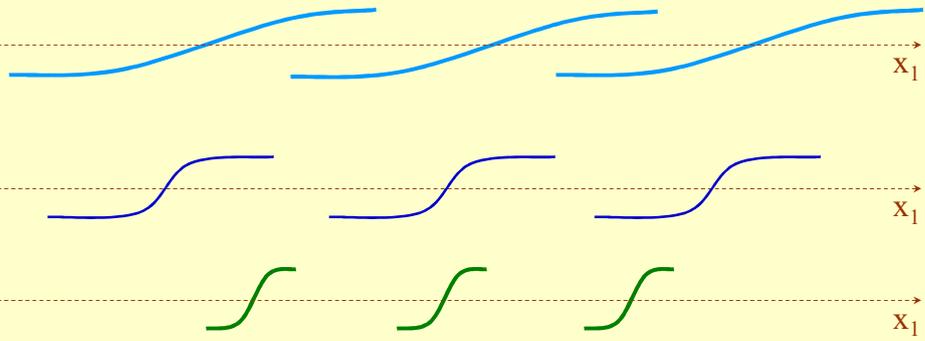


44

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

recordando

O que conseguimos fazer com **um único neurônio sigmoidal**, no caso de regressões (“ $y(x_1)$ contínuo”)?

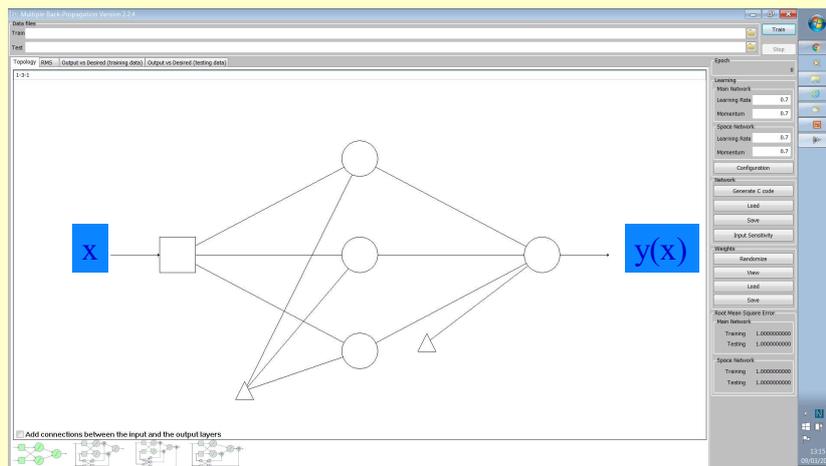


$$y = \text{tgh}(w_1 \cdot x_1 + \text{viés})$$

45

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

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



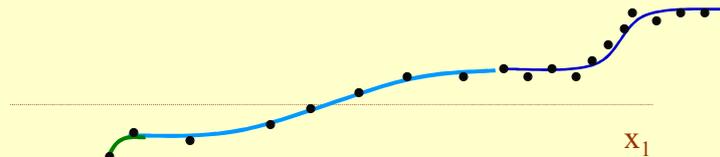
46

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?

47

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

Algumas discussões adicionais sobre o Cybenko “café com leite” da regressão univariada ...

- *Vimos acima como se comporta o regressor univariado de Cybenko “café com leite” quando o nó de saída tem função de ativação identidade, seus pesos ponderadores das saídas da primeira camada são todos unitários positivos e o peso de viés é nulo.*
- *O que ocorre se os esses pesos ponderadores não forem mais unitários? (podem ser agora positivos, negativos, encolhedores (módulo < 1) ou amplificadores (módulo > 1))*
- *E se o seu peso de viés do 4o nó não for mais nulo?*

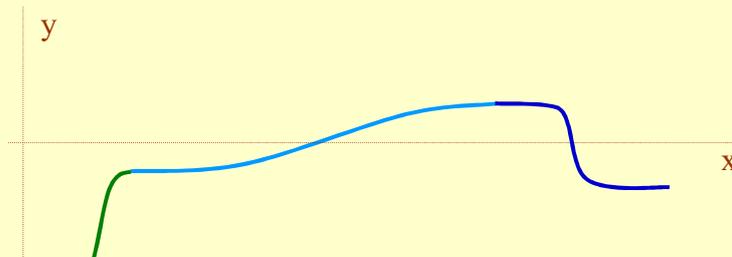
48

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

Cybenko “café com leite”, para aproximação universal de funções de 1 variável x apenas?



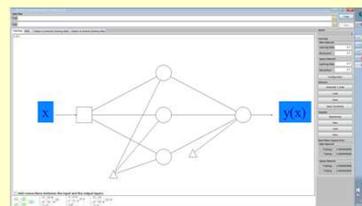
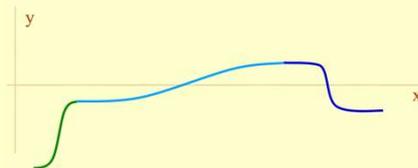
... superposição de várias sigmóides deslocadas e escaladas em x e em y ...



... Ponderadores das 3 tgh's da primeira camada, que são implementados nos pesos sinápticos do 4o nó, não são mais unitários nem necessariamente positivos

49

Isto indica claramente que ao menos no caso de funções univariadas no domínio e na imagem (uma única variável x no argumento e uma única variável y na “saída” da função) uma RNA de duas camadas (com vários neurônios na segunda, não apenas 3 como ilustrado) pode aproximar qualquer função contínua univariada com erro bem pequeno se necessário: se desejado, basta usarmos mais e mais nós na segunda camada do MLP, aumentando assim arbitrariamente a precisão da aproximação da função alvo da modelagem.



50

Cybenko foi além de mostrar a viabilidade de aproximação em casos unidimensionais, ele fez a prova de *Aproximação Universal* no âmbito de funções de múltiplas variáveis!

Qualquer Função(X) genérica pode ser aproximada por um MLP – O que é bom para Estimção / Regressão Contínua (um de alvos neste curso) !!!

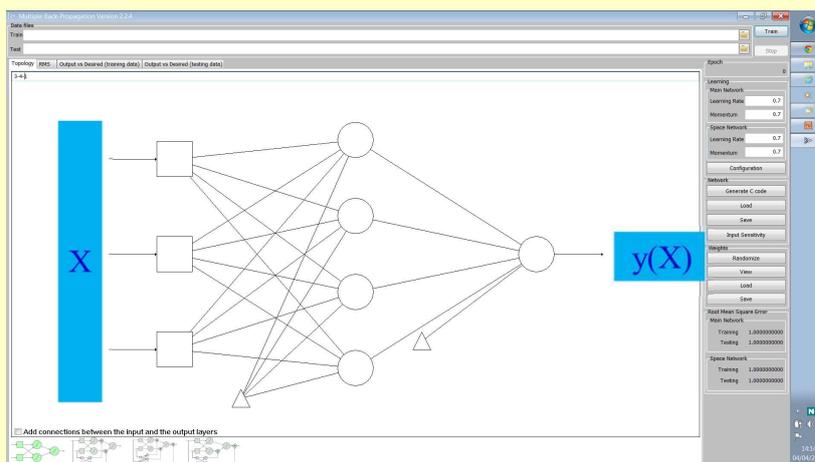
E ...

É também bom para o Reconhecimento de Padrões (nosso segundo alvo de modelagem) com o MLP ... Trata-se de um caso específico de função binária na saída !!!

51

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

Vamos neste caso para o terreno de vetor de entradas X em lugar de um x unidimensional

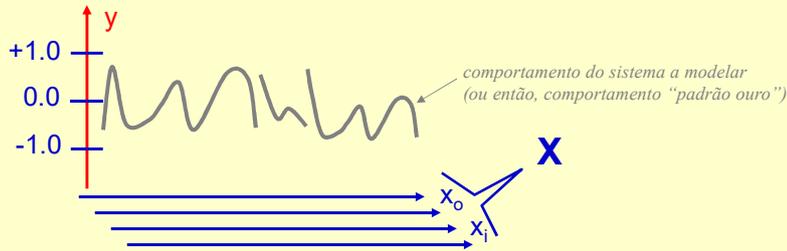


52

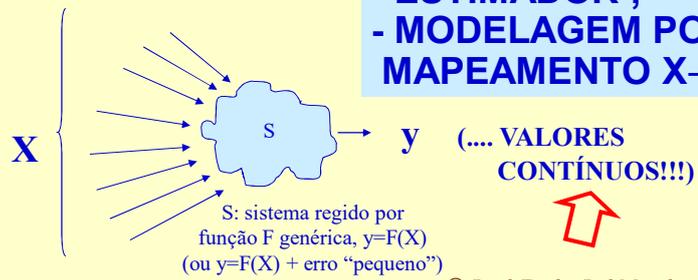
Métodos Numéricos e Reconhecimento de Padrões – © Prof. Emilio Del Moral – EPUSP

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

recordando



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

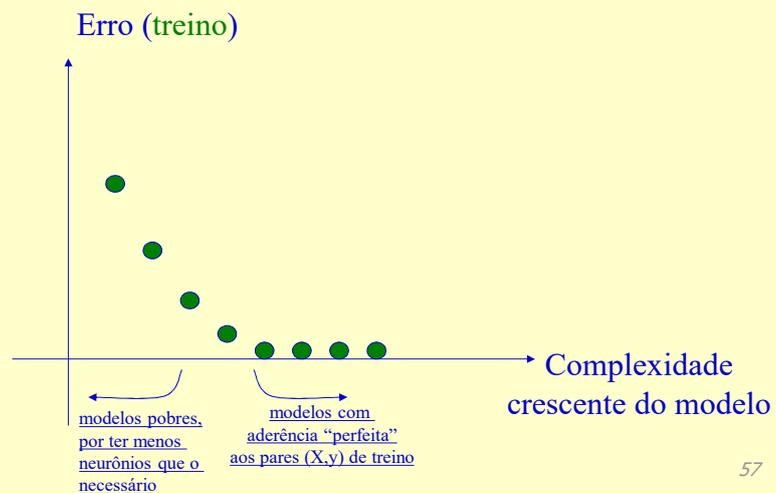


53

© Prof. Emilio Del Moral – EPUSP

Aumento de aderência aos dados de treino com o aumento de nós da RNA ...

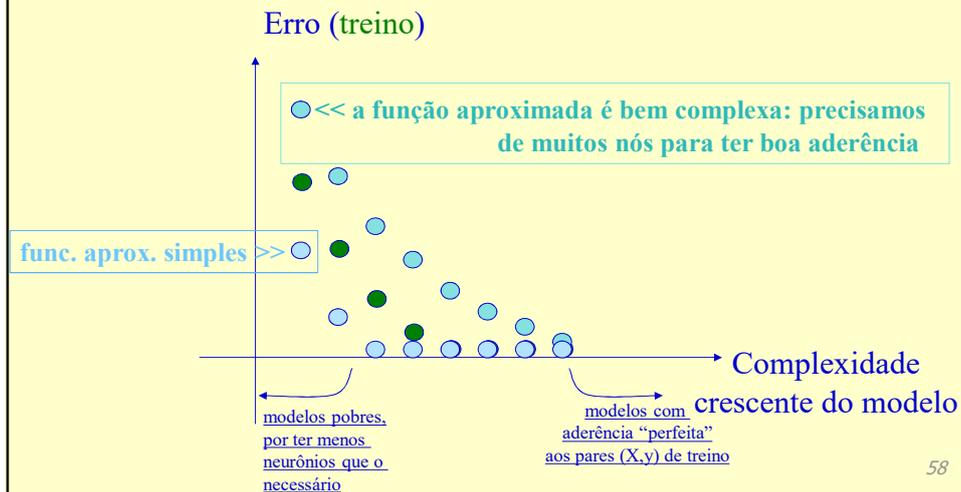
recordando



57

© Prof. Emilio Del Moral – EPUSP

Aumento de aderência aos dados de treino com o aumento de nós da RNA ...



Isto quer dizer que sempre é melhor termos um modelo com mais nós neurais que um modelo com menos nós neurais?

Afinal, da mesma maneira que a computação de um regressor polinomial de grau seis engloba a computação dos regressores polinomiais de graus menores, os modelos com mais nós neurais englobam os mais simples (em termos de capacidades de computações possíveis) correto?

Sim, correto! Mas há um limite no "lucro" em tal estratégia, dado pelo fenômeno de

Sobreaprendizado e perda de generalização ...

60

© Prof. Emilio Del Moral – EPUSP

O conceito de sobreaprendizado nos dará critérios adicionais para a definição do número de neurônios / grau de complexidade de uma rede neural

61

© Prof. Emilio Del Moral – EPUSP

Inicialmente façamos algumas reflexões sobre o conjunto de treino em regressores e sua característica de amostra da totalidade de comportamentos (X;y) possíveis do sistema sendo modelado

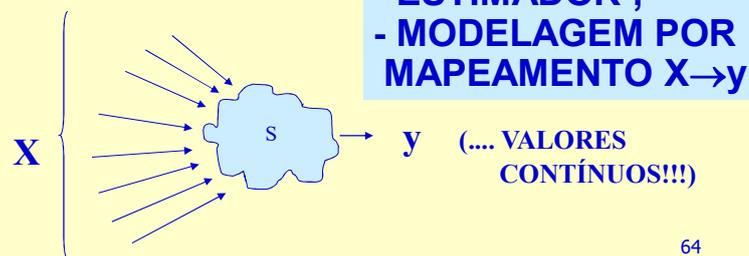
63

© Prof. Emilio Del Moral Hernandez

Hipótese:

$$y = f(x_1, x_2, x_3, x_4, x_5, x_6, x_7, x_8, x_9, x_{10})$$

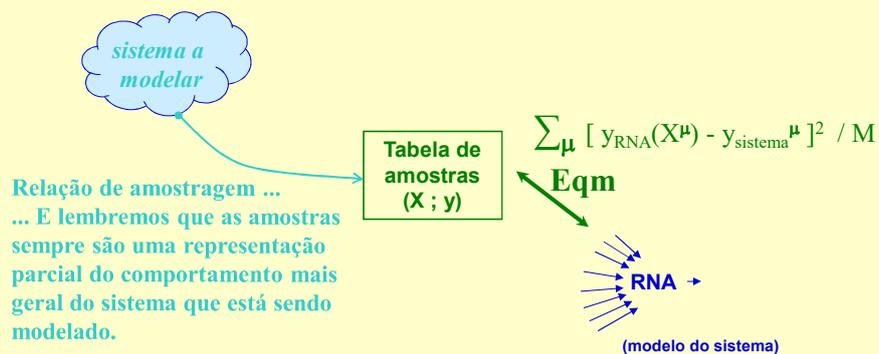
$$y = f(X)$$



64

© Prof. Emilio Del Moral Hernandez

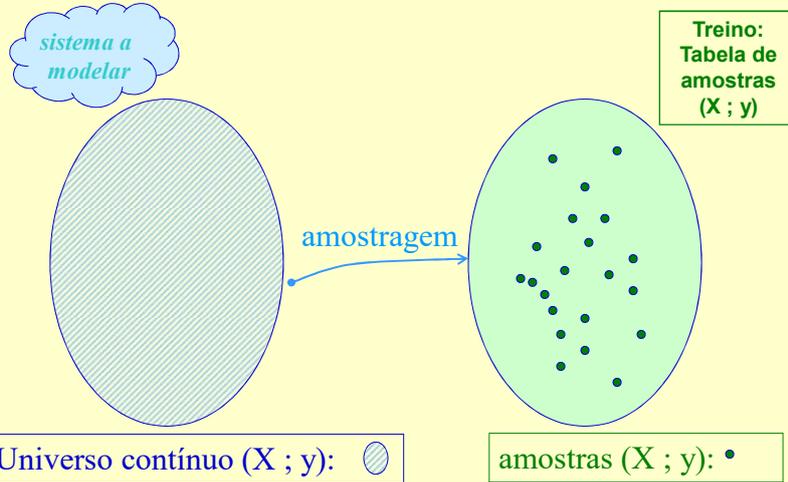
O treinamento mira minimizar o **Eqm** das amostras (X ; y) de treino. (exclusivamente!)



65

© Prof. Emilio Del Moral Hernandez

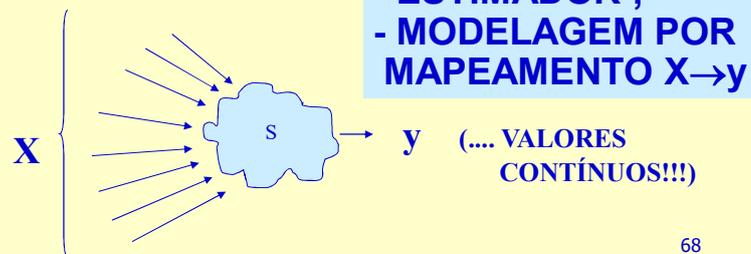
Conjunto de treino ... Composto de amostragens discretas do universo contínuo ($X ; y$)



66

© Prof. Emilio Del Moral Hernandez

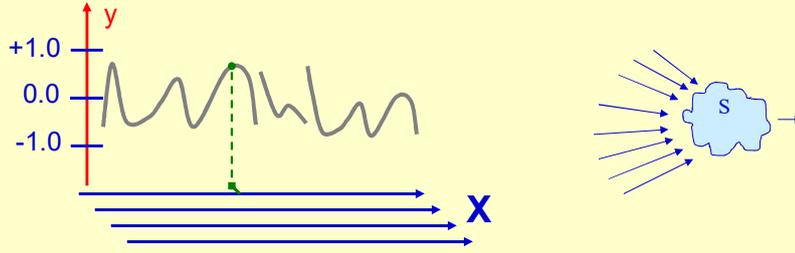
A função $f(X)$ (interna ao sistema a modelar) é **desconhecida** e assume valores contínuos



68

© Prof. Emilio Del Moral Hernandez

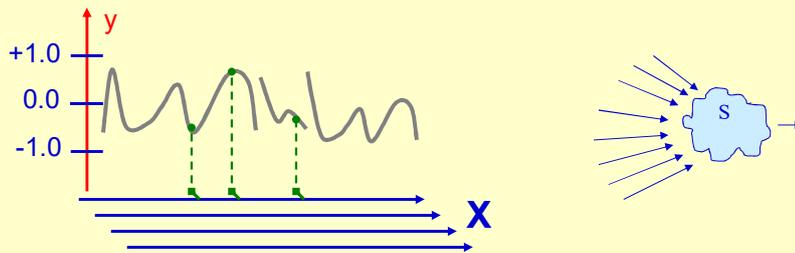
Mas podemos observar o sistema e registrar, por exemplo, um par (X;y) isolado ...



69

© Prof. Emilio Del Moral Hernandez

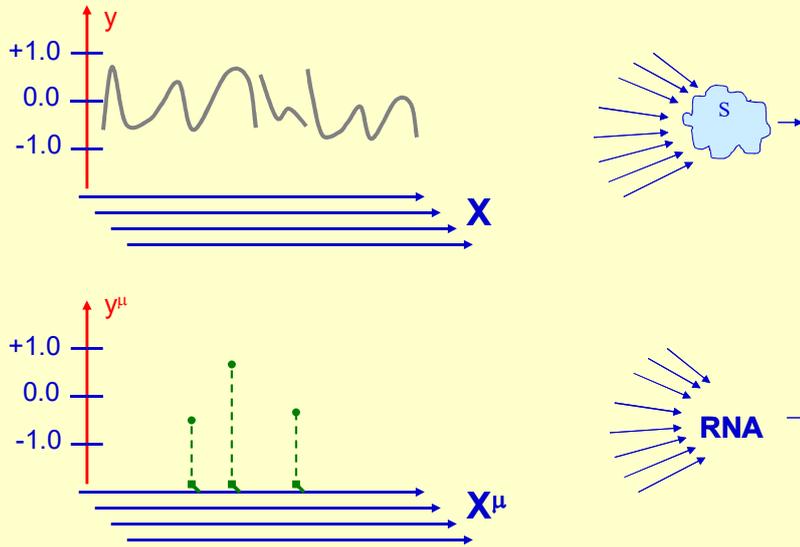
Ou mais pares (X;y) 3 pares, para começar



70

© Prof. Emilio Del Moral Hernandez

E podemos construir com esses 3 pares um conjunto de treino $\{(X^\mu ; y^\mu)\}$ para uma RNA ...



71

© Prof. Emilio Del Moral Hernandez

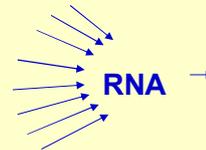
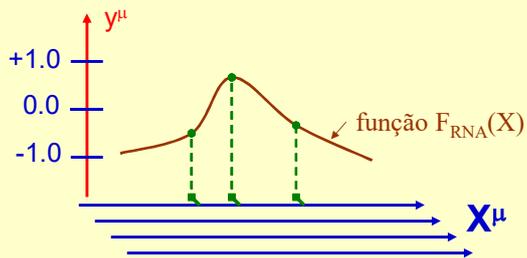
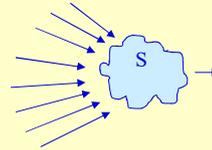
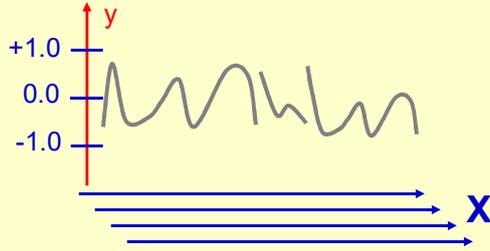
Gráfico fornecido pelo ambiente MBP da evolução do erro em função do número de repetidos usados da “bússola do erro” com o eixo vertical independente”: isto conecta o MBP com o gráfico apresentado anteriormente



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

© Prof. Emilio Del Moral Hernandez

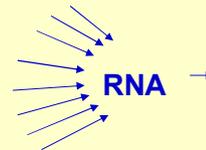
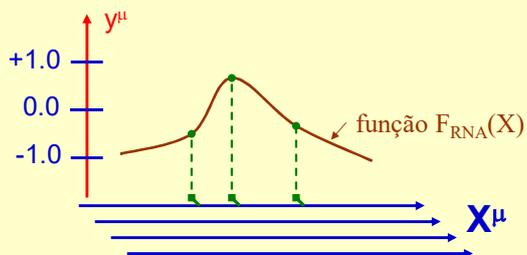
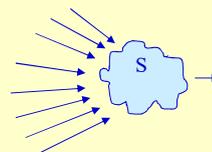
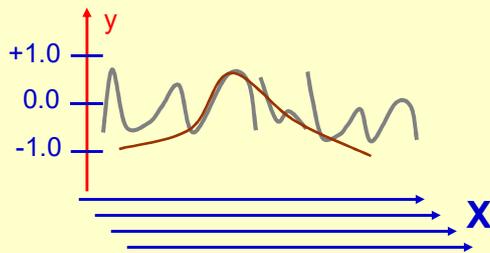
Após treino por EBP e minimização do Eqm, a RNA será um modelo fiel aos 3 pares



73

© Prof. Emilio Del Moral Hernandez

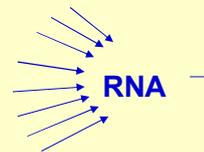
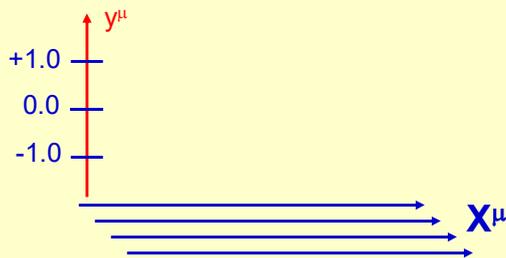
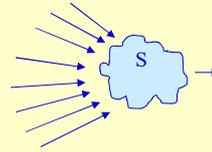
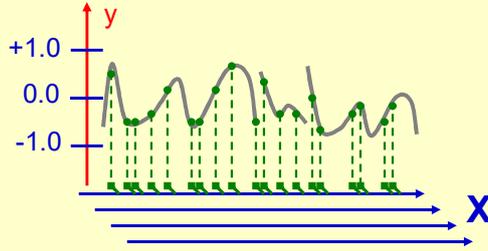
Mas não necessariamente fiel ao sistema sendo modelado



74

© Prof. Emilio Del Moral Hernandez

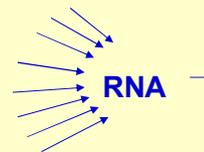
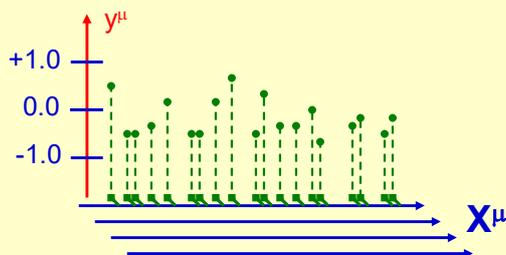
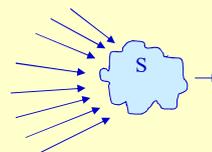
Melhor se tivermos mais pares, de forma a retratar melhor o sistema a modelar



75

© Prof. Emilio Del Moral Hernandez

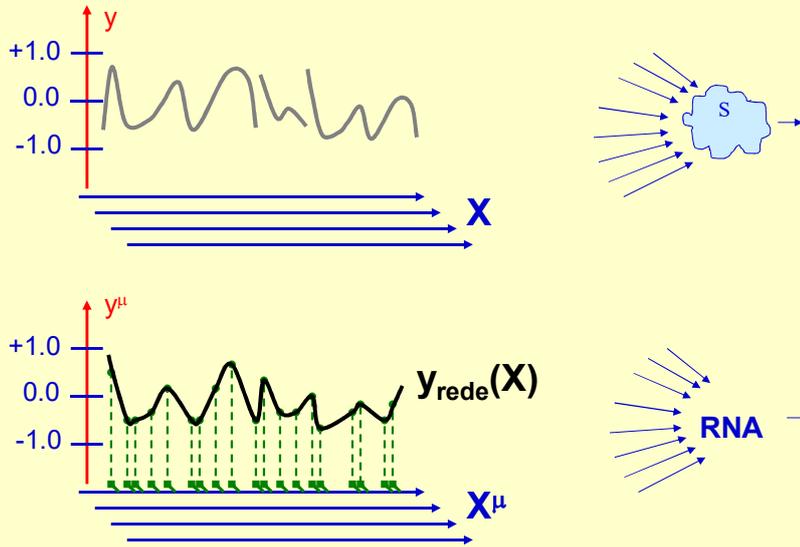
O tamanho M do conjunto de treino será maior



76

© Prof. Emilio Del Moral Hernandez

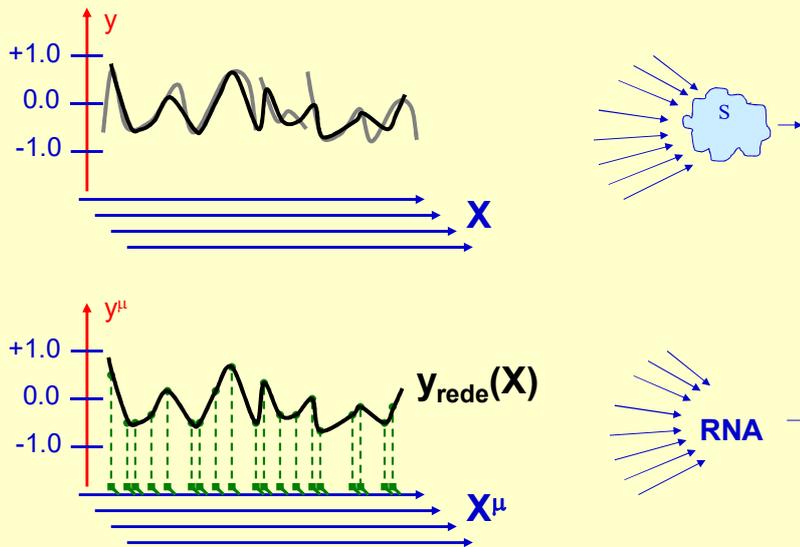
... e portanto o modelo será mais bem definido e mais próximo ao sistema modelado



77

© Prof. Emilio Del Moral Hernandez

... mas claro, nunca será totalmente idêntico ao sistema



78

© Prof. Emilio Del Moral Hernandez

E basta termos muitas amostras $(X;y)$ – ou seja, alto valor de M –, e ter um ajuste bom a elas – ou seja, alto número de nós neurais –, e a modelagem será de qualidade?

Não é bem assim!!!! ...

O fenômeno de Sobreaprendizado / Sobreajuste, que ocorre em modelos mais complexos leva à perda de generalização (falha do modelo em novos cenários), mesmo quando o número de observações usadas no treino desse modelo (M) é bastante grande ...

79

© Prof. Emilio Del Moral – EPUSP

Sobreaprendizado em polinômios:

Primeiro entendamos o conceito num universo mais familiar (e mais simples), o de regressão polinomial univariada, usada para representar dados com comportamento linear ou não linear e sujeitos a alguma flutuação em “y”

...

Depois, vocês mesmos podem pensar nos equivalentes dos nossos raciocínios feitos aqui para o universo de polinômios, mas para o universo de RNAs e mesmo de outros tipos de modelos com número de parâmetros variável (complexidade variável) que você conheça ...

82

© Prof. Emilio Del Moral – EPUSP

Falemos em lousa um pouco sobre a reta média para um conjunto de pares (x,y), a parábola média, a cúbica média ... etc

$$y \sim ax+b ; \quad y \sim ax^2 +bx +c ; \quad y \sim ax^3 +bx^2 +cx +d$$

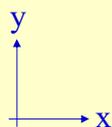
e mais além, falemos sobre regressão polinomial univariada, com o grau do polinômio aproximador podendo ser 1, 2, 3, ou mesmo graus bastante mais altos como 50, 51 etc.

$$y \sim ax^{51} +bx^{50} +cx^{49} + \dots$$

83

© Prof. Emilio Del Moral – EPUSP

Sobreaprendizado ilustrado em sua fenomenologia, na regressão polinomial (de vários graus) univariada



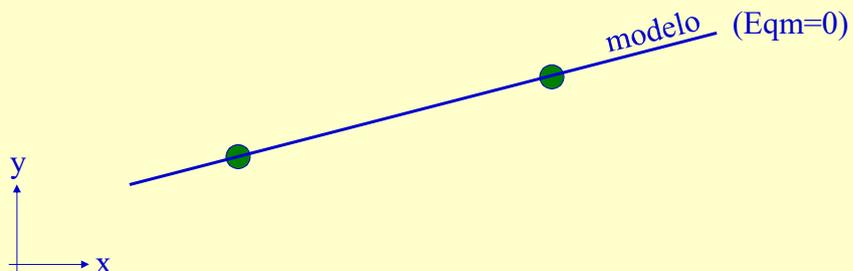
Os dados empíricos (x^i, y^i) estão em verde;

84

© Prof. Emilio Del Moral – EPUSP

Sobreaprendizado ilustrado em sua fenomenologia, na regressão polinomial (de vários graus) univariada

façamos uso de modelagem linear ...



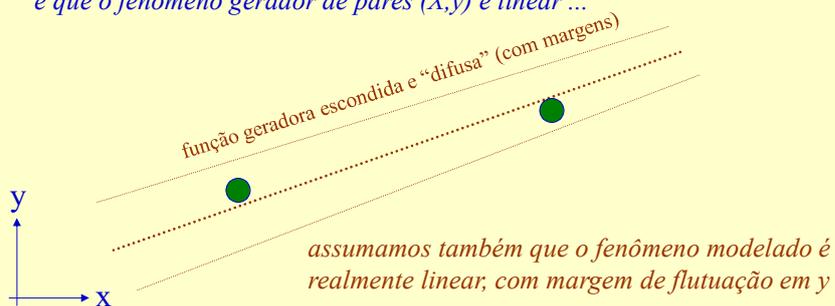
Os dados empíricos (x^i, y^i) estão em verde;
O modelo linear gerado a partir dos dados, em azul.

85

© Prof. Emilio Del Moral – EPUSP

Sobreaprendizado ilustrado em sua fenomenologia, na regressão polinomial (de vários graus) univariada

e que o fenômeno gerador de pares (X,y) é linear ...



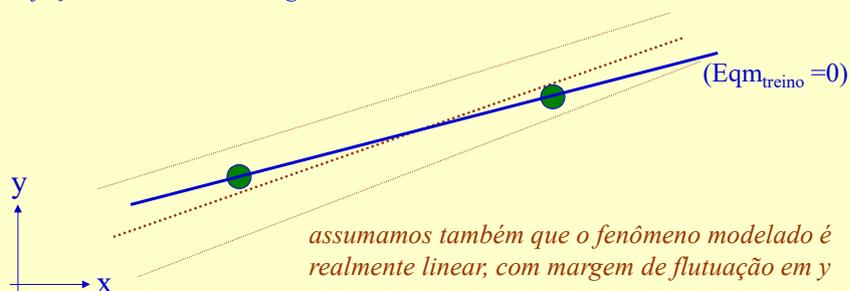
Os dados empíricos (x^i, y^i) estão em verde;
O modelo linear gerado a partir dos dados, em azul.
O fenômeno gerador de pares (x,y) é linear em essência, mas tem alguma flutuação randômica em y . A tendência e os limites da flutuação estão representados em marrom

86

© Prof. Emilio Del Moral – EPUSP

Sobreaprendizado ilustrado em sua fenomenologia, na regressão polinomial (de vários graus) univariada

façamos uso de modelagem linear ...



assumamos também que o fenômeno modelado é realmente linear, com margem de flutuação em y

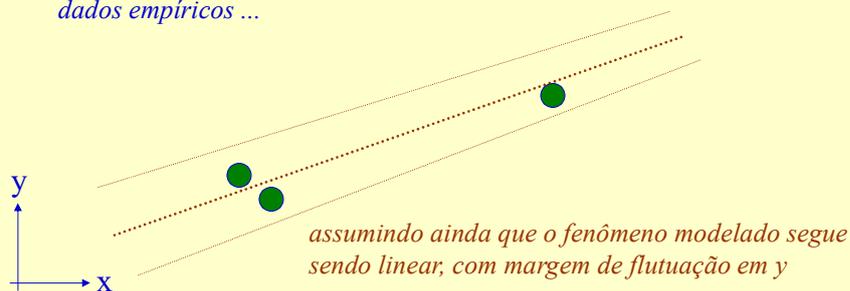
Os dados empíricos (x^h, y^h) estão em verde;
O modelo linear gerado a partir dos dados, em azul.
O fenômeno gerador de pares (x, y) é linear em essência, mas tem alguma flutuação randômica em y. A tendência e os limites da flutuação estão representados em marrom

87

© Prof. Emilio Del Moral – EPUSP

Sobreaprendizado ilustrado em sua fenomenologia, na regressão polinomial (de vários graus) univariada

Consideremos agora nova situação com mais dados empíricos ...



assumindo ainda que o fenômeno modelado segue sendo linear, com margem de flutuação em y

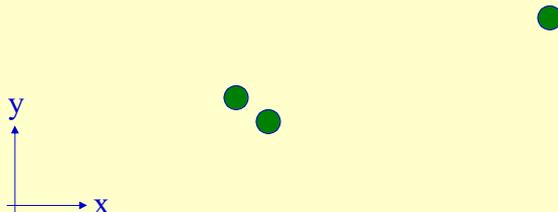
Os dados empíricos (x^h, y^h) estão em verde;

88

© Prof. Emilio Del Moral – EPUSP

Sobreaprendizado ilustrado em sua fenomenologia, na regressão polinomial (de vários graus) univariada

Consideremos agora nova situação com mais dados empíricos ... mas na modelagem não se sabe se é o fenômeno linear ou quadrático ...



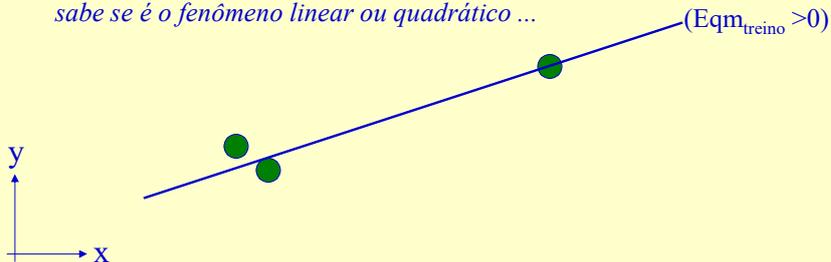
Os dados empíricos (x^i, y^i) estão em verde;

89

© Prof. Emilio Del Moral – EPUSP

Sobreaprendizado ilustrado em sua fenomenologia, na regressão polinomial (de vários graus) univariada

Consideremos agora nova situação com mais dados empíricos ... mas na modelagem não se sabe se é o fenômeno linear ou quadrático ...



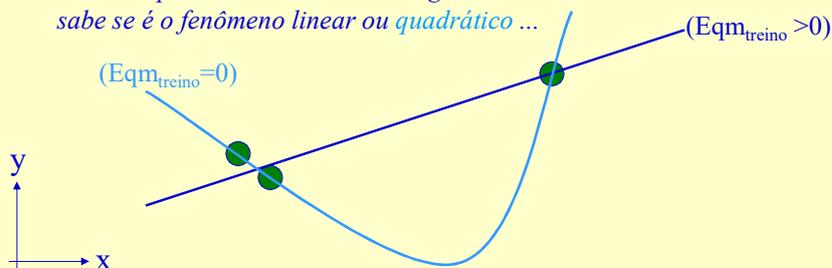
Os dados empíricos (x^i, y^i) estão em verde;

90

© Prof. Emilio Del Moral – EPUSP

Sobreaprendizado ilustrado em sua fenomenologia, na regressão polinomial (de vários graus) univariada

Consideremos agora nova situação com mais dados empíricos ... mas na modelagem não se sabe se é o fenômeno linear ou quadrático ...



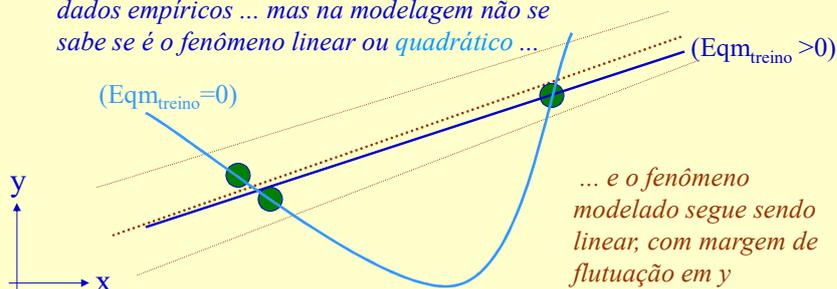
Os dados empíricos (x^i, y^i) estão em verde;
Dois modelos polinomiais gerados a partir dos dados, em azuis.
O fenômeno gerador de pares (x, y) é linear em essência, mas tem alguma flutuação randômica em y . A tendência e os limites da flutuação estão representados em marrom

91

© Prof. Emilio Del Moral – EPUSP

Sobreaprendizado ilustrado em sua fenomenologia, na regressão polinomial (de vários graus) univariada

Consideremos agora nova situação com mais dados empíricos ... mas na modelagem não se sabe se é o fenômeno linear ou quadrático ...

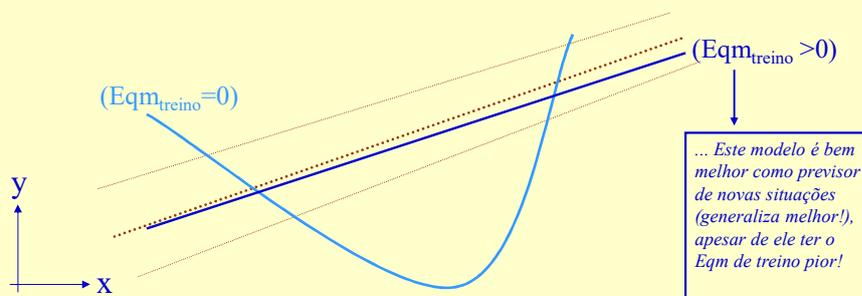


Os dados empíricos (x^i, y^i) estão em verde;
Dois modelos polinomiais gerados a partir dos dados, em azuis.
O fenômeno gerador de pares (x, y) é linear em essência, mas tem alguma flutuação randômica em y . A tendência e os limites da flutuação estão representados em marrom

92

© Prof. Emilio Del Moral – EPUSP

Sobreaprendizado ilustrado em sua fenomenologia, na regressão polinomial (de vários graus) univariada



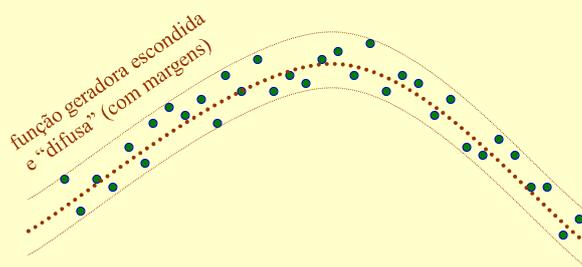
Os dados empíricos (x^i, y^i) estão em verde;
 Dois modelos polinomiais gerados a partir dos dados, em azuis.
 O fenômeno gerador de pares (x, y) é linear em essência, mas tem alguma flutuação randômica em y . A tendência e os limites da flutuação estão representados em marrom

93

© Prof. Emilio Del Moral – EPUSP

Sobreaprendizado ilustrado em sua fenomenologia, na regressão polinomial (de vários graus) univariada

um novo exemplo



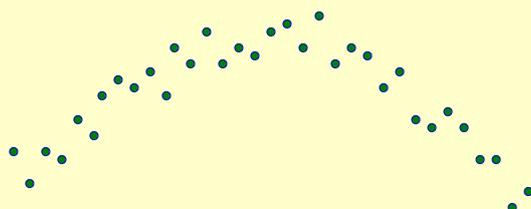
Os dados empíricos (x^i, y^i) estão em verde;
 O fenômeno gerador de pares (x, y) é quadrático em essência, mas tem alguma flutuação randômica em y . A tendência e os limites da flutuação estão representados em marrom

94

© Prof. Emilio Del Moral – EPUSP

Sobreaprendizado ilustrado em sua fenomenologia, na regressão polinomial (de vários graus) univariada

um novo exemplo



Os dados empíricos (x^{μ}, y^{μ}) estão em verde;

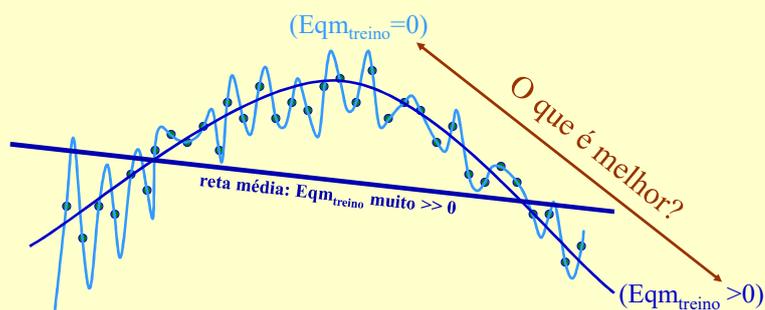
O fenômeno gerador de pares (x, y) é quadrático em essência, mas tem alguma flutuação randômica em y . A tendência e os limites da flutuação estão representados em marrom

95

© Prof. Emilio Del Moral – EPUSP

Sobreaprendizado ilustrado em sua fenomenologia, na regressão polinomial (de vários graus) univariada

um novo exemplo



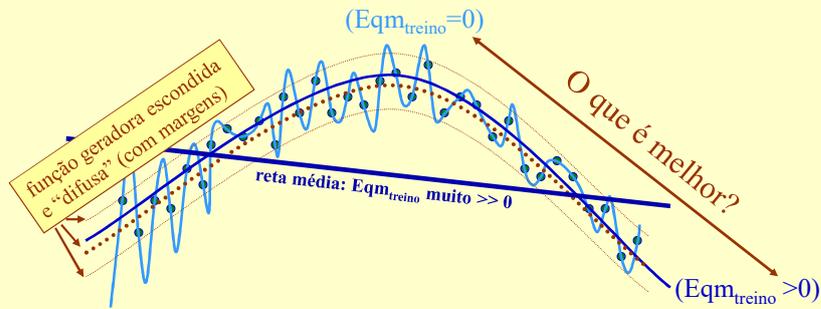
Os dados empíricos (x^{μ}, y^{μ}) estão em verde;
Três modelos polinomiais gerados a partir dos dados, em azuis.

96

© Prof. Emilio Del Moral – EPUSP

Sobreaprendizado ilustrado em sua fenomenologia, na regressão polinomial (de vários graus) univariada

um novo exemplo

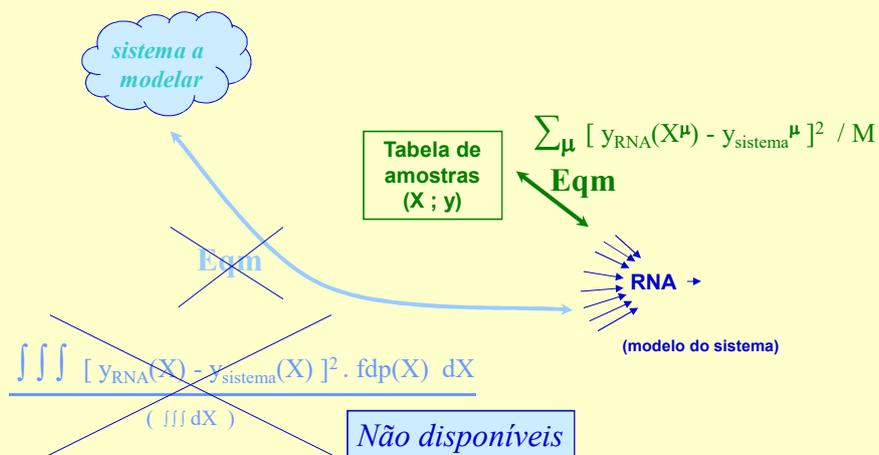


Os dados empíricos (x^μ, y^μ) estão em verde;
Três modelos polinomiais gerados a partir dos dados, em azuis.

97

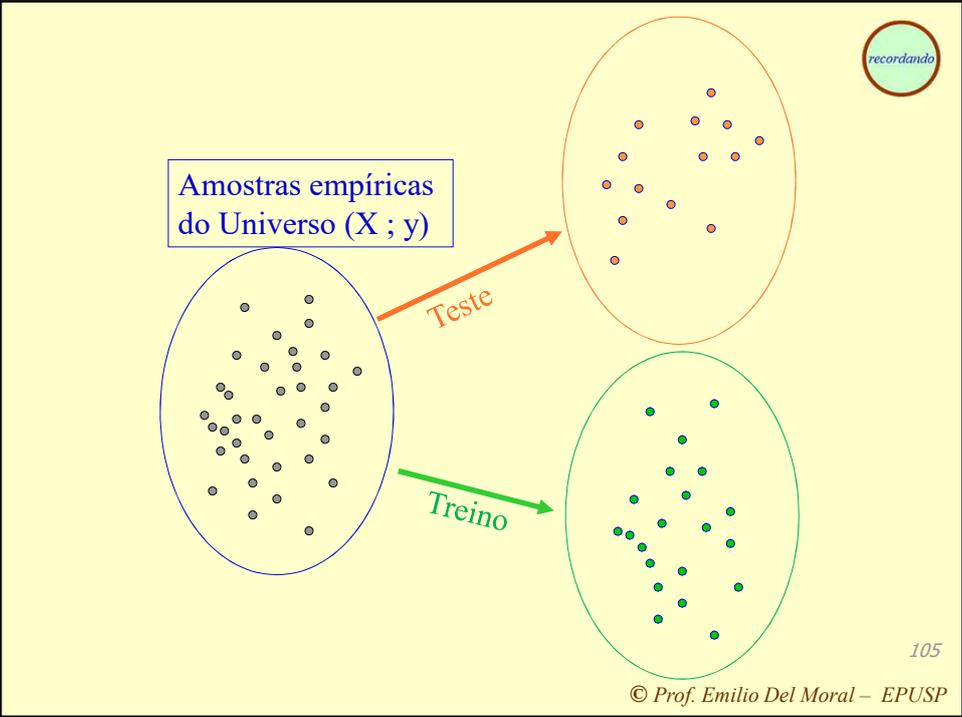
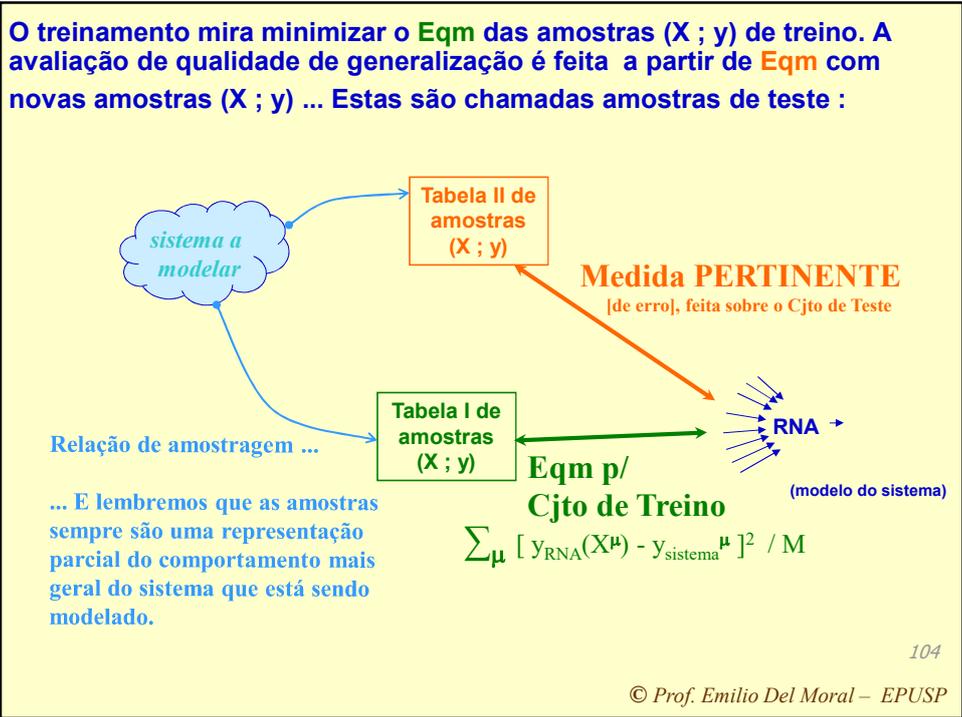
© Prof. Emilio Del Moral – EPUSP

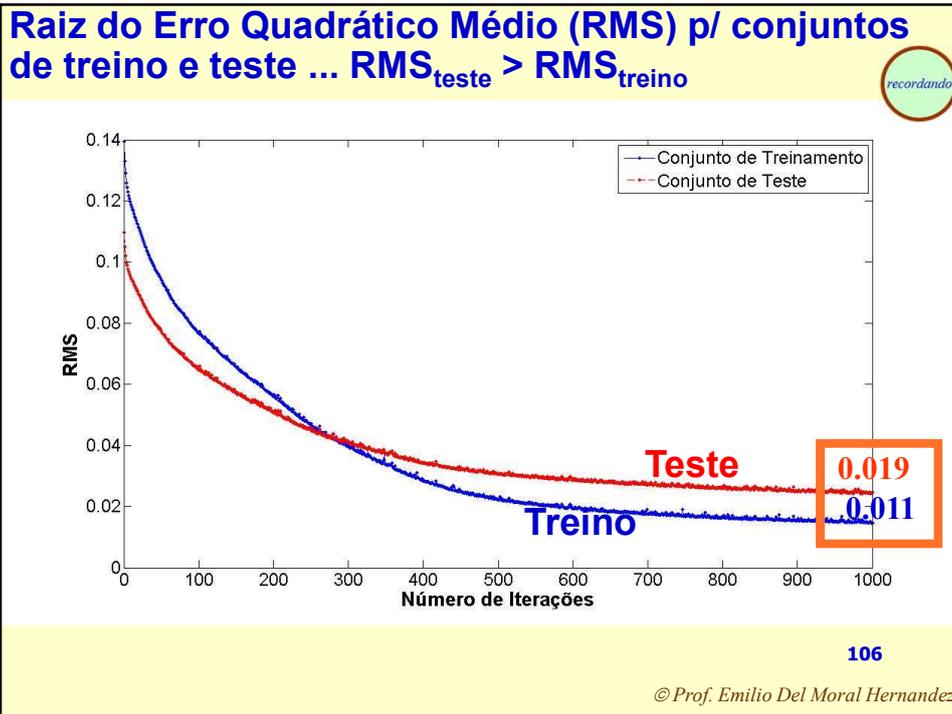
Sistema ... Amostras de treino ... RNA ...



103

© Prof. Emilio Del Moral Hernandez





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 sucessivos da RNA e/ou dos dados e/ou do modelo, em ciclos diversos, começando desde o passo 0 ou do passo 1]

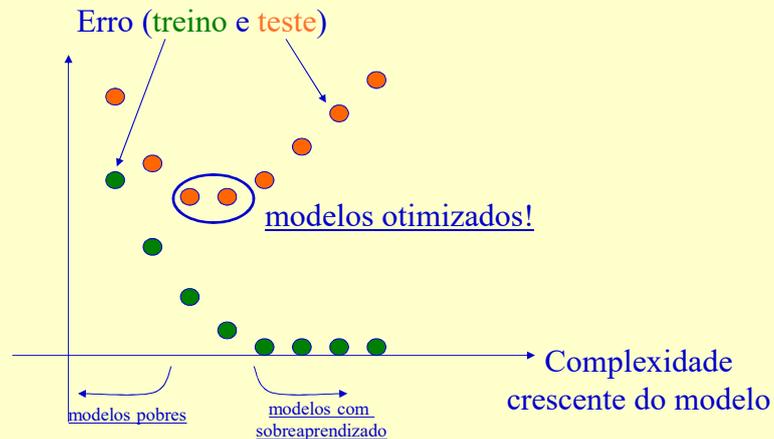
3) *Fase de USO FINAL da RNA, com y efetivamente não conhecido, e estimado com conhecimento dos X + uso do modelo calibrado.*

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

107

© Prof. Emilio Del Moral – EPUSP

Sobreaprendizado em “sumário executivo”



109

© Prof. Emilio Del Moral – EPUSP

Mas cuidado ... você experimentou rodar mais de uma vez o MBP sobre os mesmos dados de treino e ver se o resultado final é o mesmo?

- Não há garantia de que o W randômico, seguido de Gradiente Descendente leve ao mínimo global quando há vários mínimos, só há garantia de que levará a um dos mínimos locais, que pode não ser o mínimo global
- Isto faz necessário rodar várias vezes a otimização de pesos e selecionar a configuração com melhor Eqm
- É legítimo ficarmos com a rede otimizada de menos Eqm final? Ou deveríamos ter um ataque de avaliar o Eqm médio de várias tentativas?
- *Outro assunto ... e quanto à variação associada ao k-fold Cross Validation? Média ou melhor Eqm?*

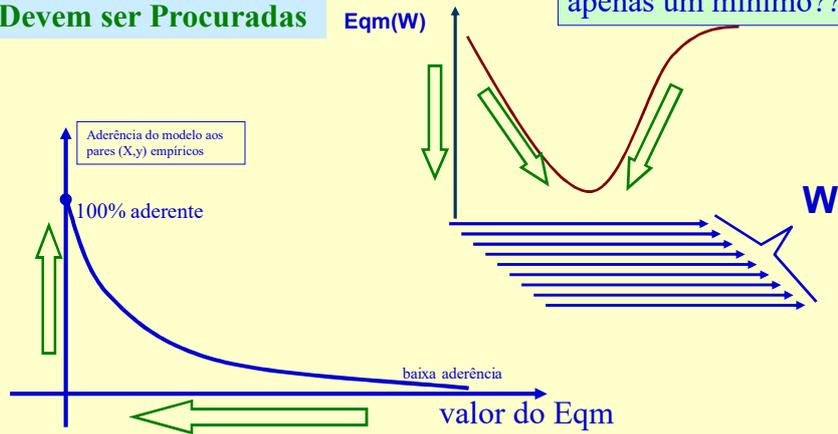
111

© 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 buscar Maximização da aderência = Mínimo Eqm possível

As Setas Verdes Indicam Situações que Devem ser Procuradas



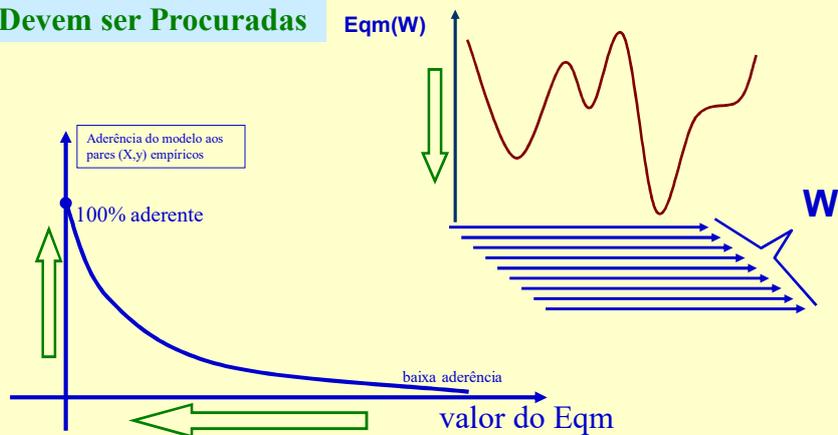
112

© 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 buscar Maximização da aderência = Mínimo Eqm possível

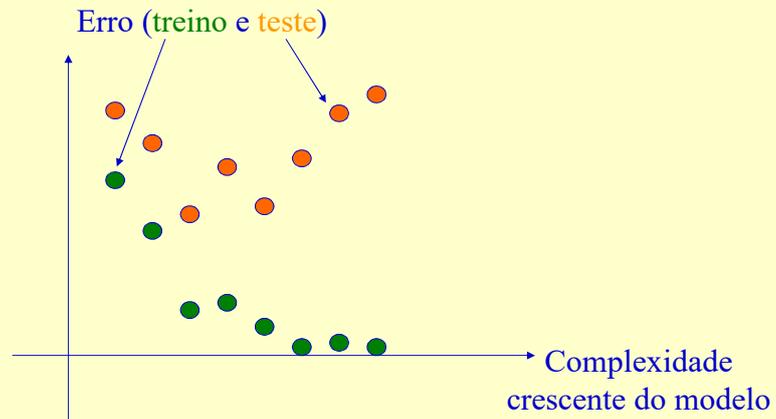
As Setas Verdes Indicam Situações que Devem ser Procuradas



113

© Prof. Emilio Del Moral – EPUSP

Atenção para componentes randômicas que impactam muito quando se faz um único ensaio de medida de erro, para cada tamanho de rede específico (um ensaio apenas, para cada grau de complexidade) ...



114

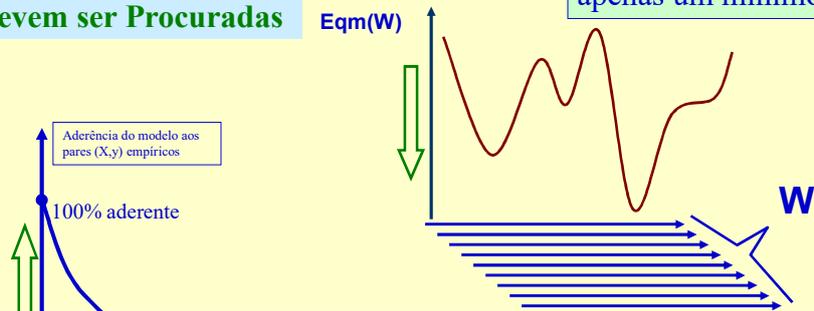
© Prof. Emilio Del Moral – EPUSP

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

Devemos buscar Maximização da aderência = Mínimo E_{qm} possível

As Setas Verdes Indicam Situações que Devem ser Procuradas

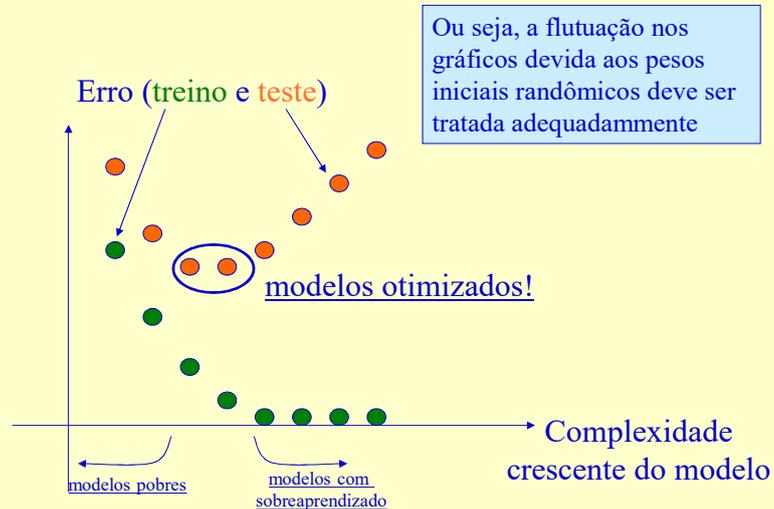
Será que temos apenas um mínimo??



... Para não sermos reféns de mínimos locais com alto E_{qm} , podemos aplicar o gradiente descendente repetidamente na mesma RNA, com novos pesos iniciais randômicos em cada rodada, mantendo para o modelo final apenas os valores de pesos associados ao ensaio com o melhor dos resultados finais no E_{qm} !

© Prof. Emilio Del Moral – EPUSP

Com repetidos ensaios em cada grau de complexidade os mínimos locais são evitados e detectamos adequadamente o sobreaprendizado



116

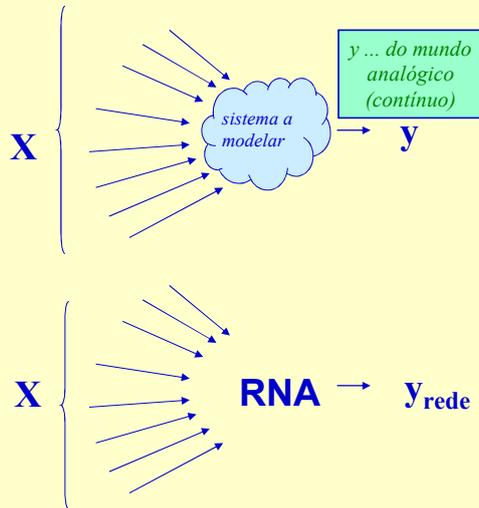
© Prof. Emilio Del Moral – EPUSP

Identificando os ingredientes para o risco de sobreaprendizado nos contextos de regressão multivariada e de reconhecimento de padrões multivariado

121

© Prof. Emilio Del Moral – EPUSP

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



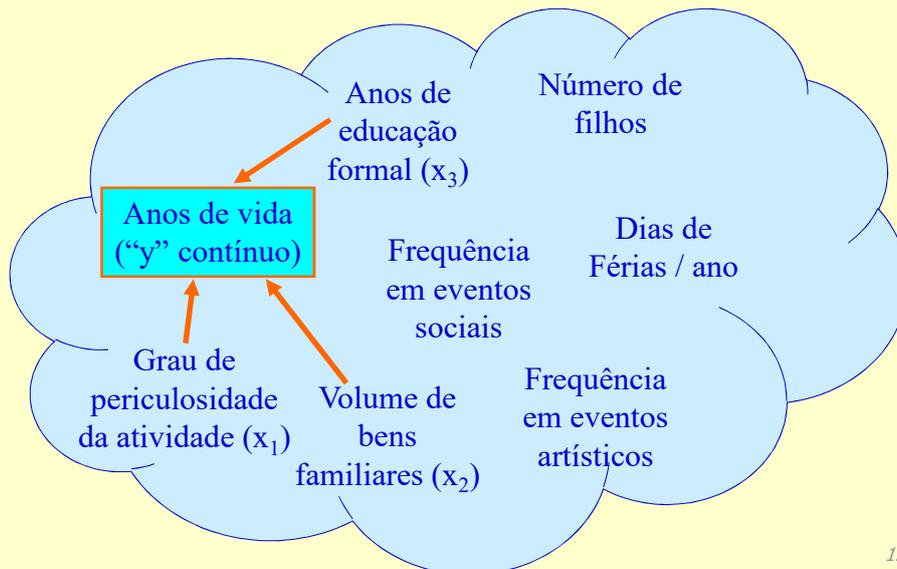
Assumimos que a variável y do sistema a modelar é uma função (normalmente desconhecida e possivelmente não linear) de diversas outras variáveis desse mesmo sistema

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

122

© Prof. Emilio Del Moral – EPUSP

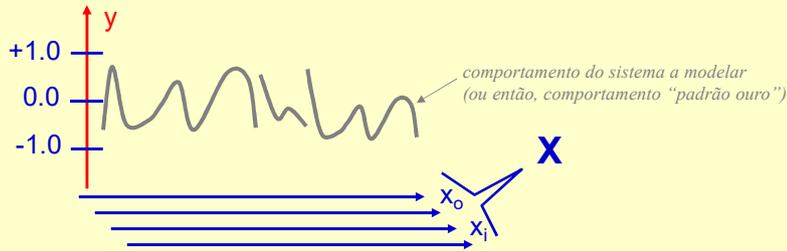
Um hipotético universo de variáveis interdependentes, passível de modelagem/ens



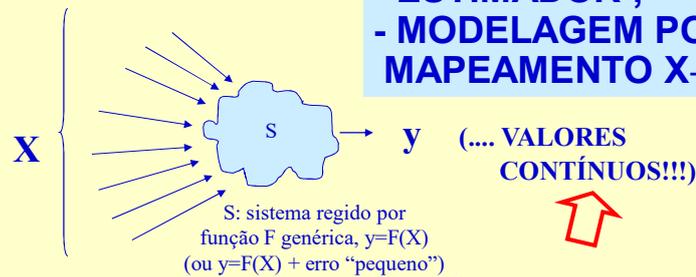
123

© Prof. Emilio Del Moral – EPUSP

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



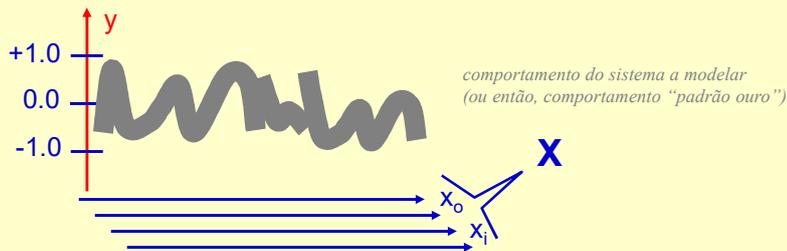
- ESTIMADOR ;
- MODELAGEM POR
MAPEAMENTO $X \rightarrow y$



124

© Prof. Emilio Del Moral – EPUSP

Cenário mais real: a “função” $y(X)$ do sistema modelado é “difusa”: $y=F_{\text{médio}}(X) + \text{flutuação}$



.... em problemas concretos / reais, há sempre alguma ambiguidade no mapeamento que leva valores de X a valores de y . Para decepção de Cybenko, não temos uma função $y=F(X)$ no sentido matemático exato, pois para uma dada ênupla de valores X fixados, temos tipicamente uma faixa de valores que podem ser observados para a variável y : $y=F_{\text{médio}}(X) + \text{flutuação}$.

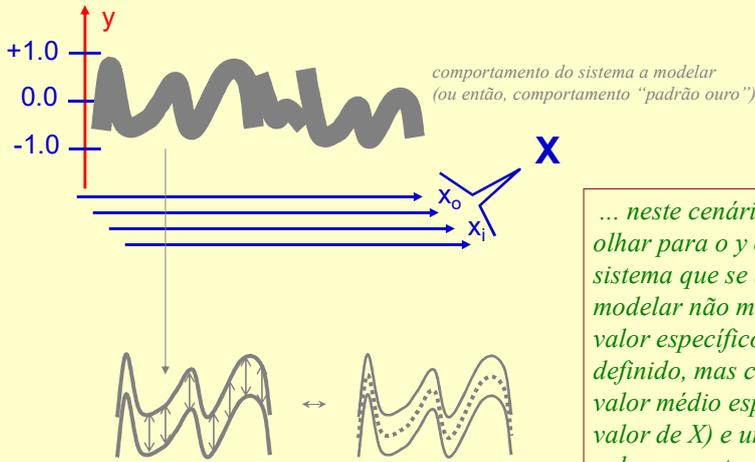
Neste cenário, buscamos que o modelo capture o comportamento médio das relações observadas entre X e y :

$$\dots y_{\text{rede}} \sim y_{\text{médio}} \text{ esperado para um dado } X$$

125

© Prof. Emilio Del Moral – EPUSP

Cenário mais real: a “função” $y(X)$ do sistema modelado é “difusa”: $y = F_{\text{médio}}(X) + \text{flutuação} \dots$



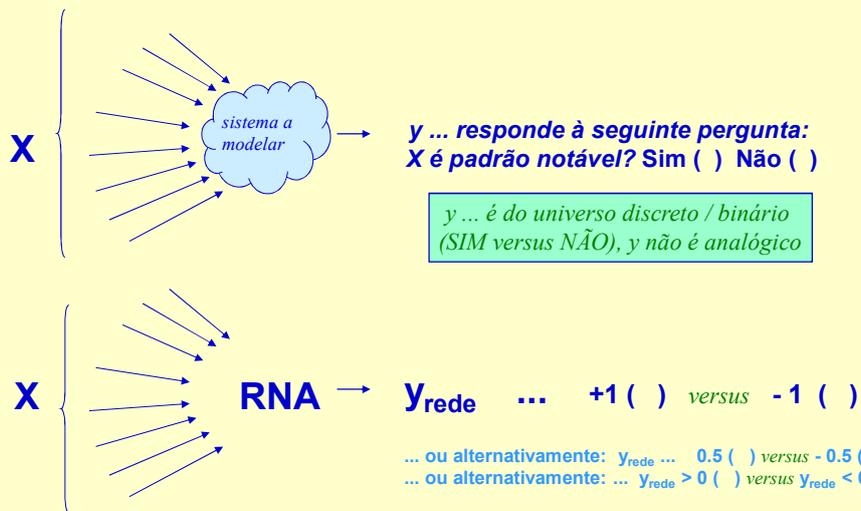
... neste cenário, podemos olhar para o y observado no sistema que se deseja modelar não mais como um valor específico bem definido, mas como um valor médio esperado (dado valor de X) e uma faixa de valores em torno desse valor médio esperado.

126

© Prof. Emilio Del Moral – EPUSP

RNAs como reconhecedor / detetor de padrões

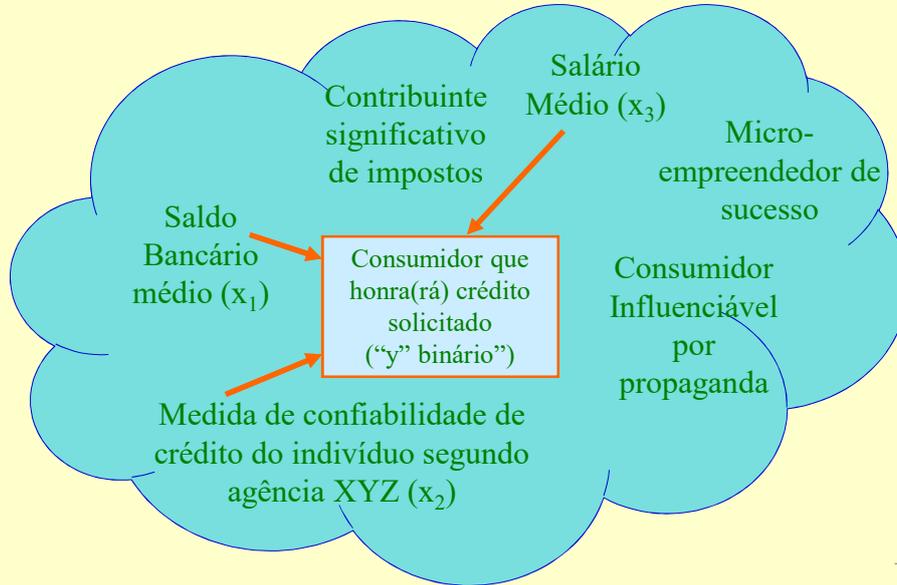
...



127

© Prof. Emilio Del Moral – EPUSP

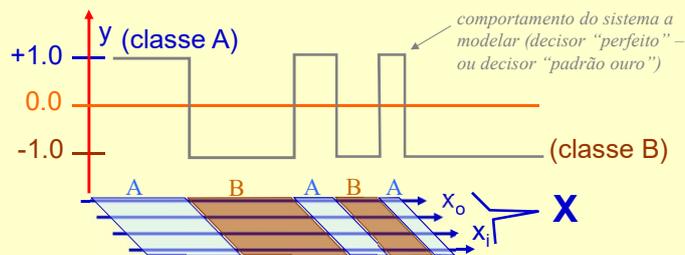
Um hipotético universo de variáveis interdependentes, passível de modelagem/ens



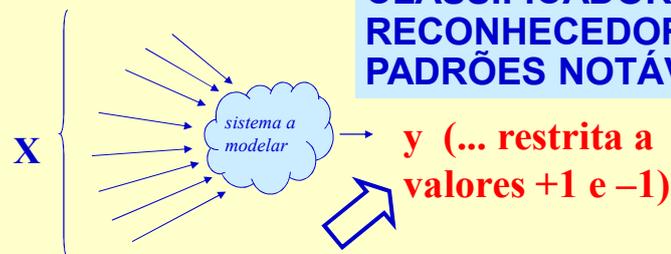
128

© Prof. Emilio Del Moral – EPUSP

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



**CLASSIFICADOR;
RECONHECEDOR DE
PADRÕES NOTÁVEIS**

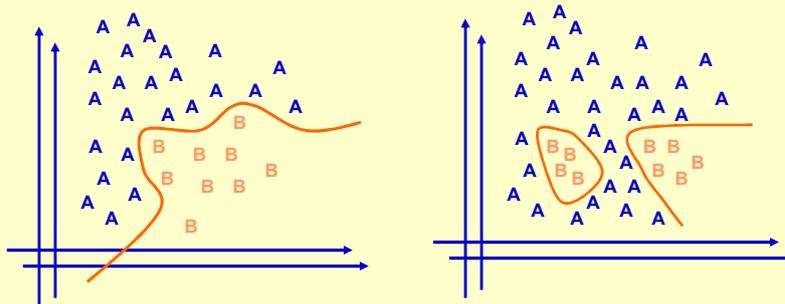


129

© Prof. Emilio Del Moral – EPUSP

Capacidade de reconhecimento de padrões em casos complexos NÃO LINEARES

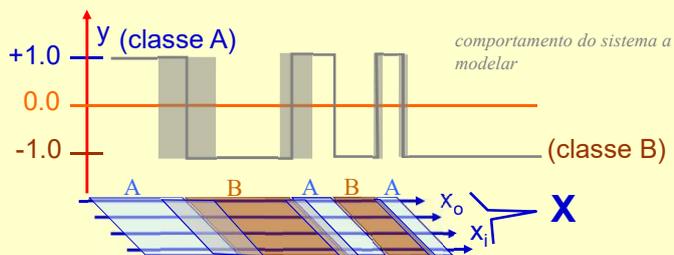
Com as RNAs, a hipersuperfície de separação entre classes vai muito além dos hiperplanos



130

© Prof. Emilio Del Moral – EPUSP

Cenário mais real: a separação entre regiões do espaço de X não é perfeitamente definida



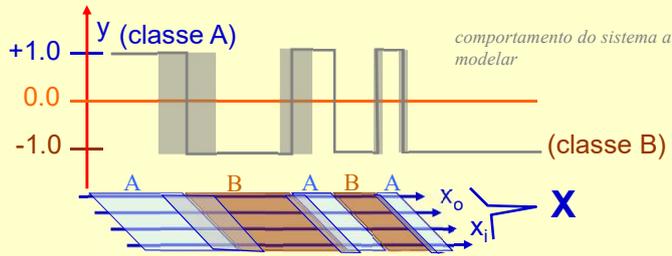
.... em problemas concretos / reais, há sempre alguma ambiguidade no mapeamento que leva valores de X aos valores discretos de y . Não temos uma função $y=F(X)$ no sentido matemático exato, pois para uma dada ênupla de valores X fixado temos em alguns casos de fronteira a possibilidade de observar no y empírico tanto a classe A quanto a classe B : $y=A$ ou B , com maior ou menor probabilidade para cada classe de acordo com o X . Neste desejamos que o modelo capture o comportamento médio das relações observadas entre X e y :

... $y_{rede} \sim$ classe 'mais esperada' para um dado X

131

© Prof. Emilio Del Moral – EPUSP

Cenário mais real: a separação entre regiões do espaço de X não é perfeitamente definida

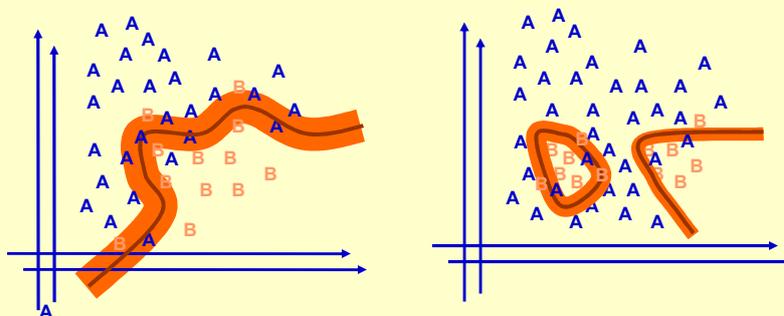


... podemos olhar para o y (classe A ou B) observado no sistema que se deseja modelar não mais como uma classe sempre bem definida e com fronteiras de separação entre A e B bem definidas no espaço de valores de X, mas como sendo delineadas na modelagem através de fronteiras com eventuais faixas de tolerância e com sobreposição parcial das classes no espaço de X

132

© Prof. Emilio Del Moral – EPUSP

Situações de classes com sobreposição parcial no espaço de atributos X ; situações de fronteiras de separação difusas ...



134

© Prof. Emilio Del Moral – EPUSP

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 sucessivos da RNA e/ou dos dados e/ou do modelo, em ciclos diversos, recomeçando desde o passo 0 ou do passo 1]

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

140

© Prof. Emilio Del Moral – EPUSP

141

Revisitando os Conjuntos de Dados Empíricos ...

Treino + Teste ...

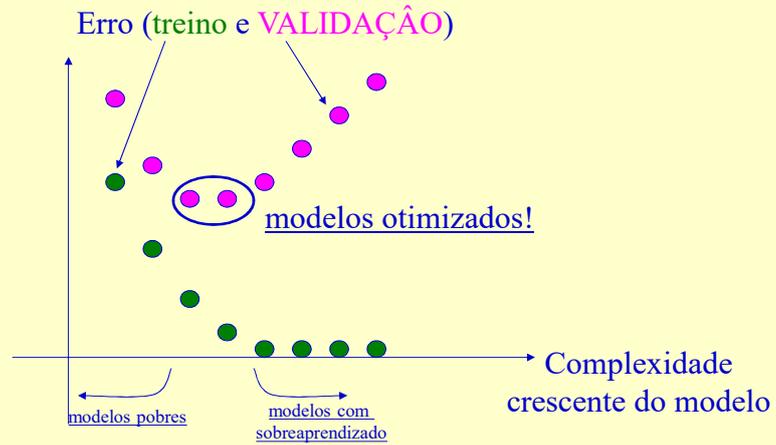
versus

Treino + **Validação** + Teste

© Prof. Emilio Del Moral Hernandez

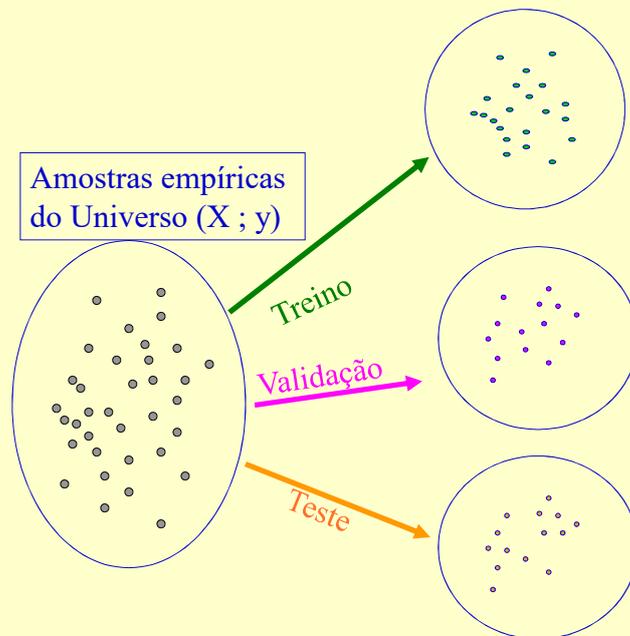
141

Sobreaprendizado em “sumário executivo”



142

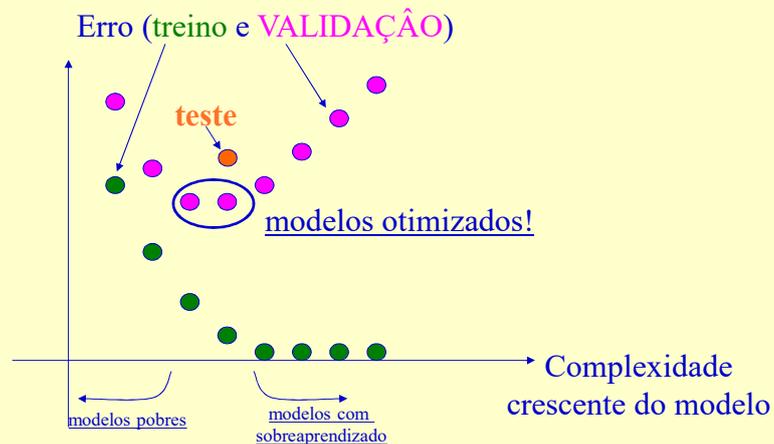
© Prof. Emilio Del Moral – EPUSP



143

© Prof. Emilio Del Moral – EPUSP

Sobreaprendizado em “sumário executivo”



144

© Prof. Emilio Del Moral – EPUSP

Alguns dos diversos usos do conceito de conjunto de validação, adicional ao conjunto de treino e de teste:

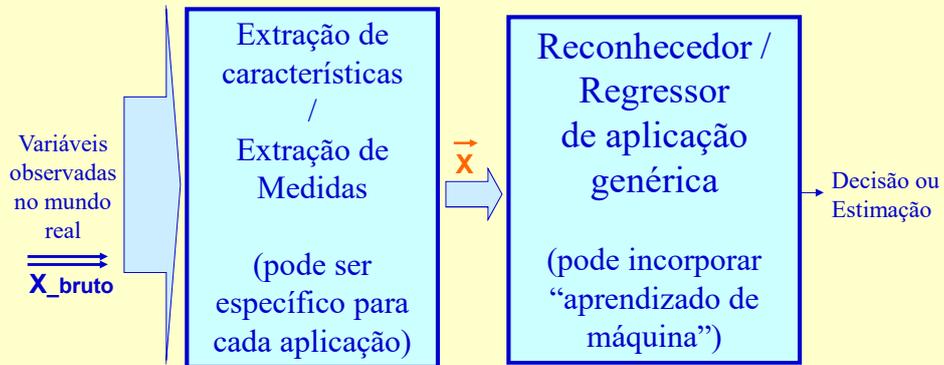
- *Seleção de complexidade do modelo neural para limitação de sobreaprendizado*
- *Ativação do early stop no aprendizado (Matlab): critério de parada adicional no processo de refinamento de pesos sinápticos*
- *Balizador no processo de seleção de estratégias de pré-processamento / extração de medidas X alternativos*
- *... etc*

145

© Prof. Emilio Del Moral – EPUSP

recordando

Elaborando uma Solução em dois estágios



146

© Prof. Emilio Del Moral – EPUSP