

Um tema importante no curso:

Sobreapredizado / Sobreajuste

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

© Prof. Emilio Del Moral – EPUSP

35

*... Passos preliminares ao tema: resgatemos
alguns aspectos d prova de Cybenko referentes
ao grau de ajuste entre modelo neural e função
sendo aproximada por esse modelo*

Prof. Emilio Del Moral Hernandez

36

36

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

© Prof. Emilio Del Moral Hernandez

37

37

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, Complexities.

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

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

where $y_j \in \mathbb{R}^n$ and $a_j, \theta_j \in \mathbb{R}$ are fixed. (y^T is the transpose of y so that $y^T x$ is the inner product of y 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:

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

Such functions arise naturally in neural network theory as the activation function of a neural node (or unit as is becoming the preferred term) [L1], [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-8619103, ONR Contract N00014-86-G-0202 and DOE Grant DE-FG02-85ER22001.

† 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 Gamba 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

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

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

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

and a set $D \subset I_n$ such that $m(D) \geq 1 - \epsilon$ and

$$|G(x) - f(x)| < \epsilon \quad \text{for } x \in D.$$

ed 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- volutionated and the number ound in [MSJ] and [BH], itness of the results of this : more attention.

n, Christopher Chase, Lee narov, Richard Lippmann, tences, and improvements

New York, 1972: uralization", *Neural Comput.* (to

tems and control, *IEEE Control*

3. Classifying learnable geometric

ndings of the 18th Annual *ACM*

s. 273–282.

1 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,

IEEE Trans. Acoust. Speech Signal

fferward networks are universal

a Neural Net and Conventional

87.

net Classifiers, Technical Report,

–475.

works by sigmoidal functions,

n, University of Lowell, 1988.

Métodos Numéricos e Re.

38

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

The screenshot shows the Wikipedia article for the 'Universal approximation theorem'. Key sections and annotations include:

- Formal statement**: The theorem states that for any continuous function f on a compact subset of \mathbb{R}^n , there exists a neural network with a single hidden layer of neurons that can approximate f to within any desired accuracy ϵ .
- Equation**:
$$F(x) = \sum_{i=1}^N \alpha_i \varphi(w_i^T x + b_i)$$
 - Annotations**:
 - $y_{rede}(X)$ points to $F(x)$.
 - X points to x .
 - número de nós escondidos** points to N .
 - sigmoidal** points to φ .
 - viés; viés do nó escondido i** points to b_i .
 - W_i : vetor de pesos do nó escondido i** points to w_i .
 - elementos do vetor de pesos do nó linear de saída W_s** points to α_i .

39

This close-up focuses on the 'Formal statement' section. The annotations provide a detailed breakdown of the mathematical components:

- Equation**:
$$F(x) = \sum_{i=1}^N \alpha_i \varphi(w_i^T x + b_i)$$
- Annotations**:
 - y_{rede}(X)**: Output of the network.
 - X**: Input vector.
 - número de nós escondidos**: Number of hidden nodes (N).
 - sigmoidal**: The activation function φ .
 - viés; viés do nó escondido i** : Bias for the hidden node i (b_i).
 - W_i : vetor de pesos do nó escondido i** : Weight vector for the hidden node i (w_i).
 - elementos do vetor de pesos do nó linear de saída W_s**: Elements of the weight vector for the output node (α_i).

40

recordando

[\[edit\]](#)

[2][3][4][5]

$$F(\mathbf{x}) = \sum_{i=1}^N \text{Fescondida_sistema}(\mathbf{x}) + b_i$$

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

Fescondida_sistema(X)

as an approximate realization of the function f where f is indepe

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

Limite de erro

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

41

- *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 da rede neural?*
- *O que ganhamos e o que perdemos se escolhermos usar MUITOS nós na construção da rede neural?*

47

42

... Um parênteses para discutirmos um pouco a aproximação universal usando tangentes hiperbólicas, sigmoides, etc ... / funções em formato de “S”

... ao menos no caso simples e bem particular de função escalar univariada ... aproximação de uma função y de uma única variável x_1 :

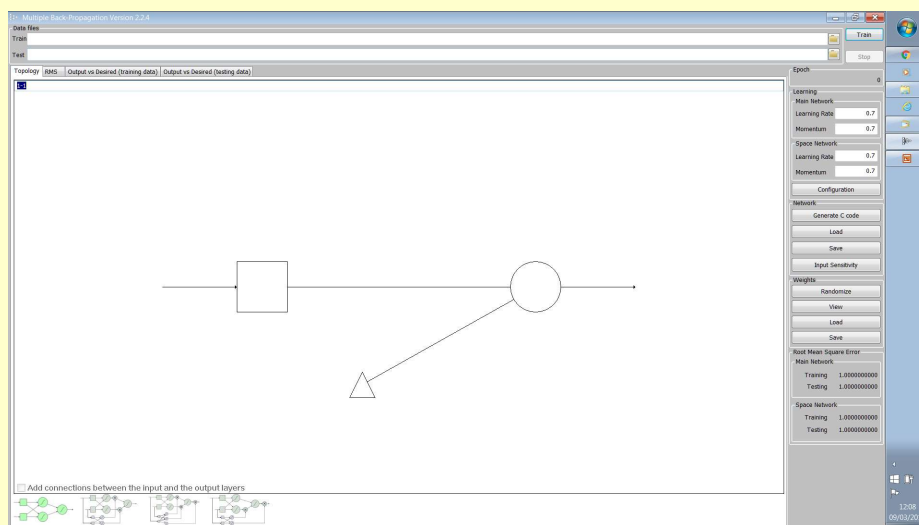
$$y(x_1)$$

43

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

43

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

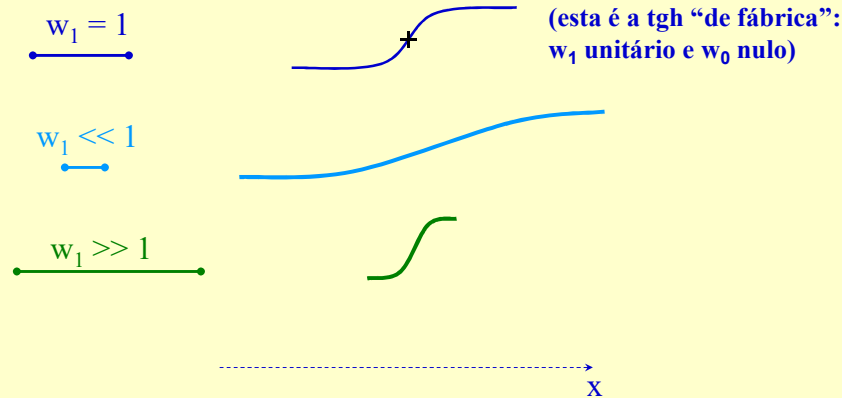


44

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

44

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**

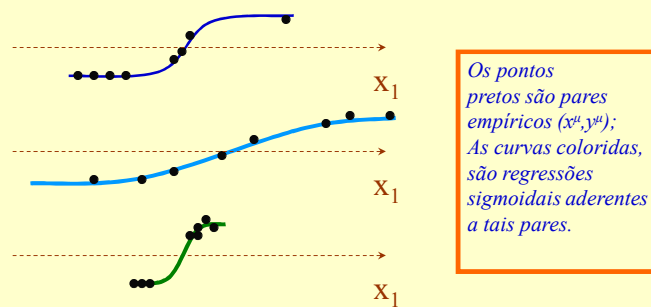


45

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

45

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

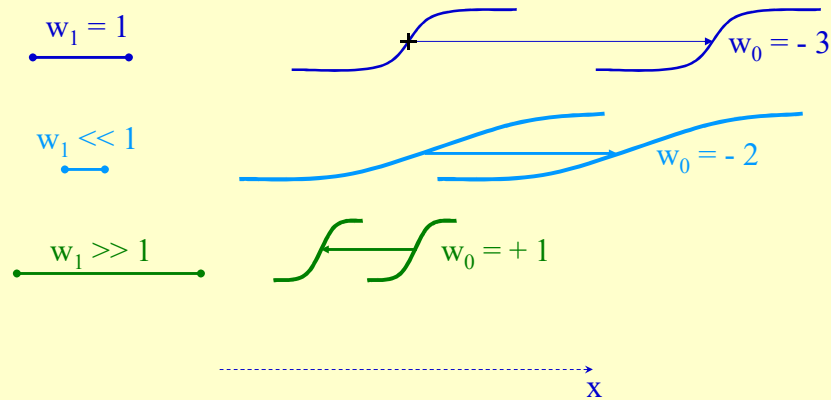


46

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

46

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

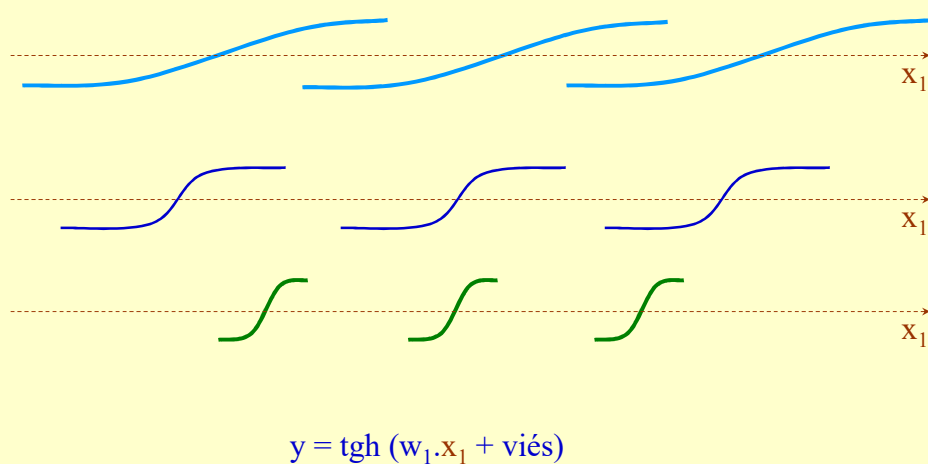


47

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

47

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

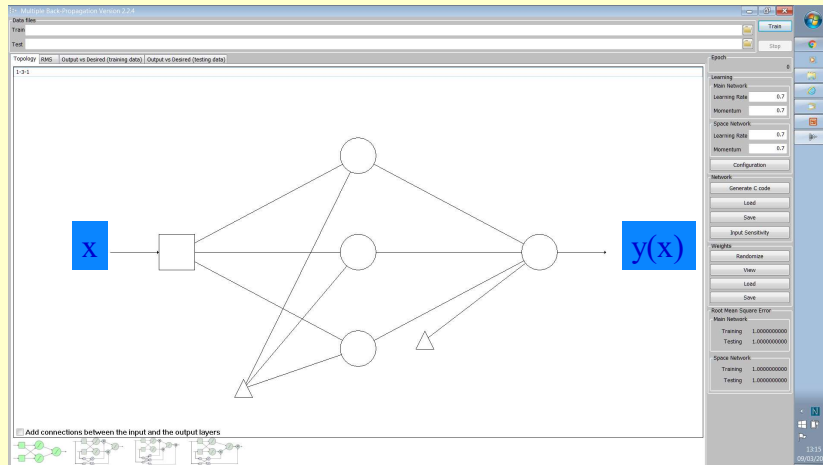


48

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

48

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



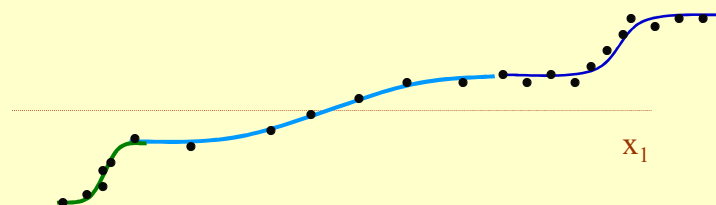
49

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

49

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 3 viéses distintos e 3 escaladores de x_1 distintos, e mais um 4o nó combinador (soma simples de 3 entradas) na camada de saída?

50

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

50

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 ($\text{módulo} < 1$) ou amplificadores ($\text{módulo} > 1$))*
- *E se o seu peso de viés do 4o nó não for mais nulo?*

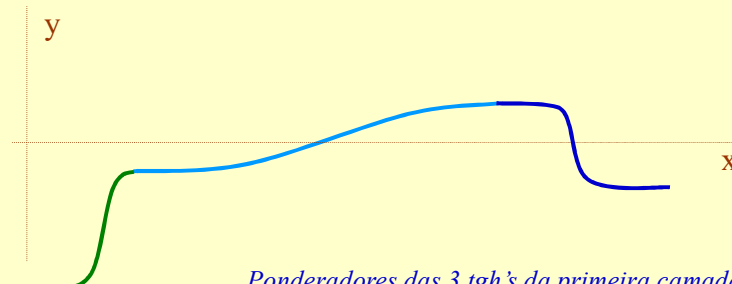
51

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

51

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 ...



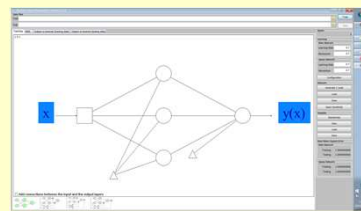
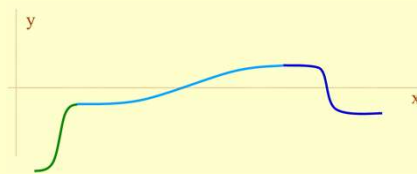
... Ponderadores das 3 tgh's da primeira camada, que são implementados nos pesos sinápticos do 4o nó, não são mais unitários nem necessariamente positivos

52

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

52

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.



53

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

53

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 Estimacão / Regressão Contínua (um de alvos neste curso) !!!

E ...

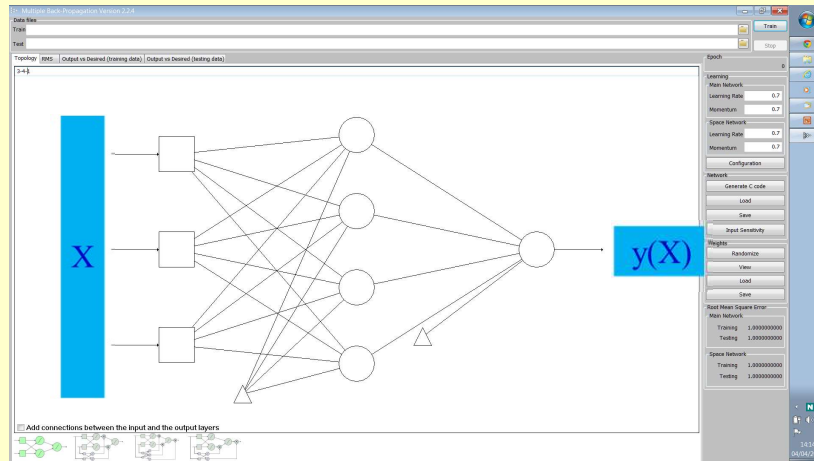
É também bom para o Reconhecimento de Padrões (nosso segundo alvo de modelagem) com o MLP ... Trata-se de um caso específico de função binária na saída !!!

54

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

54

Cybenko foi para um terreno mais complexo: temos um vetor de entradas X em lugar de um x unidimensional



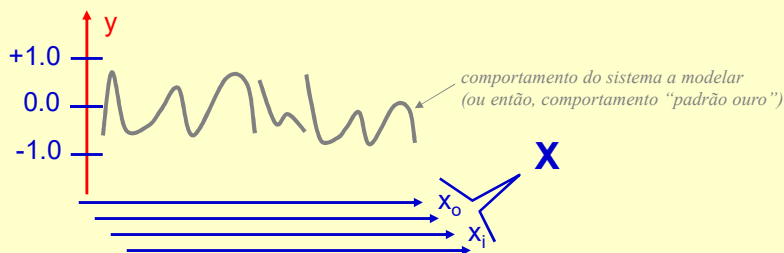
55

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

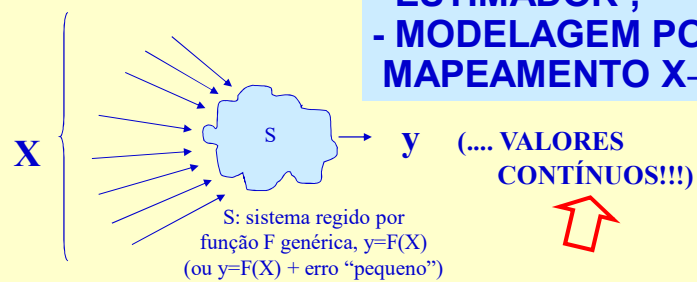
55

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

recordando



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

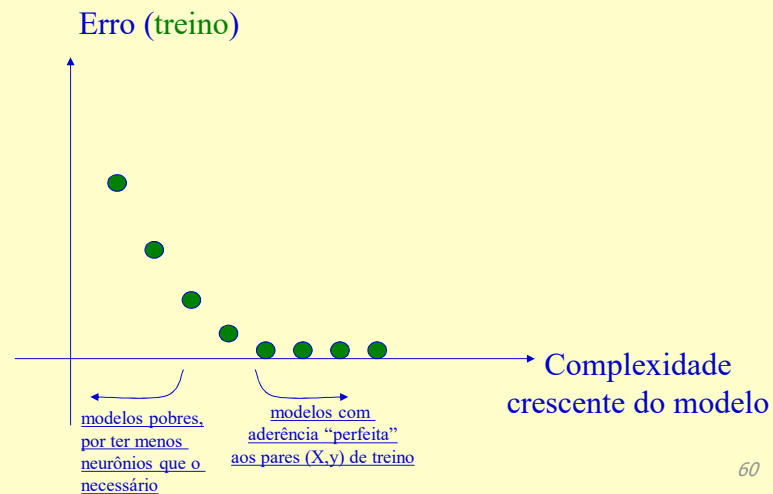


56

© Prof. Emilio Del Moral – EPUSP

56

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

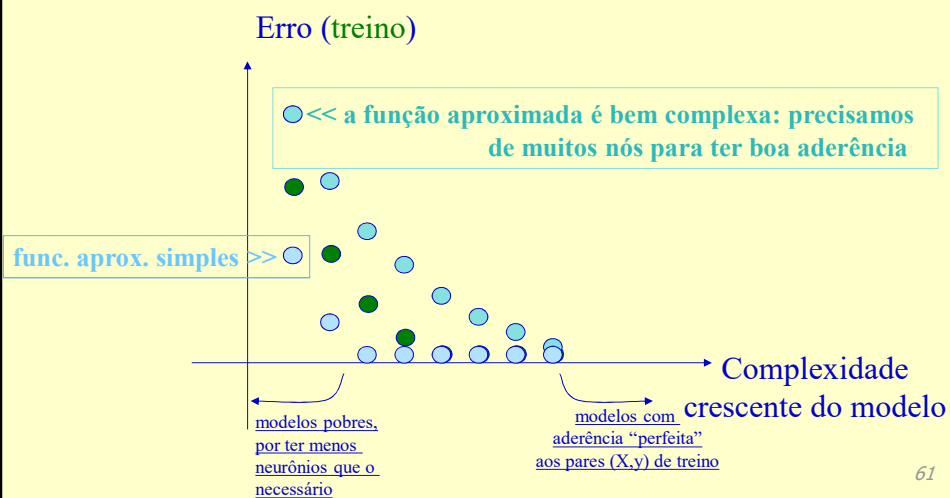


60

© Prof. Emilio Del Moral – EPUSP

60

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



61

© Prof. Emilio Del Moral – EPUSP

61

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 ...

63

© Prof. Emilio Del Moral – EPUSP

63

O conceito de sobreaprendizado nos dá critérios adicionais para a definição do número de neurônios / grau de complexidade de uma rede neural, critérios esses que vão bem além de simples econômica computacional – Esses critérios vão mais na direção de aumento de precisão na generalização

64

© Prof. Emilio Del Moral – EPUSP

64

Sobreaprendizado em polinômios:

Entendamos o conceito de sobreajuste / sobreaprendizado num universo mais familiar (e mais simples) que a modelagem com RNAs; o universo de modelagem por regressão polinomial univariada, usada para representar dados com comportamento linear ou não linear .

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 ... ⁸⁵

© Prof. Emilio Del Moral – EPUSP

85

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.

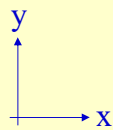
$$y \sim ax^{51} +bx^{50} +cx^{49} + \dots$$

⁸⁶

© Prof. Emilio Del Moral – EPUSP

86

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



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

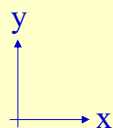
87

© Prof. Emilio Del Moral – EPUSP

87

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.

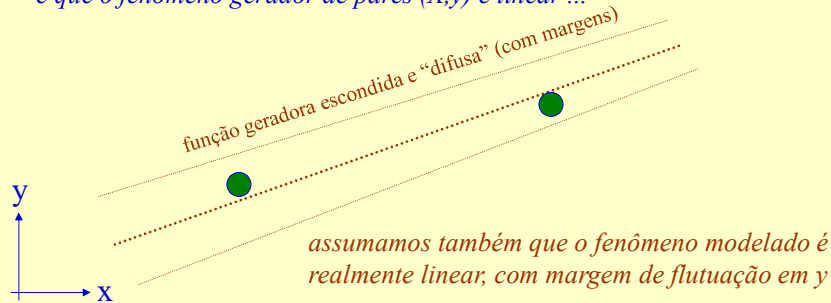
88

© Prof. Emilio Del Moral – EPUSP

88

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^u, y^u) 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

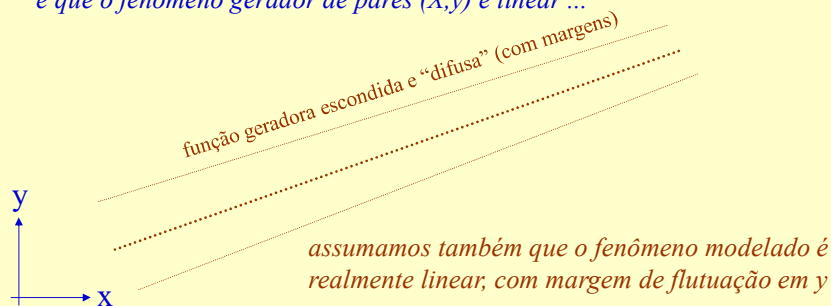
89

© Prof. Emilio Del Moral – EPUSP

89

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

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



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

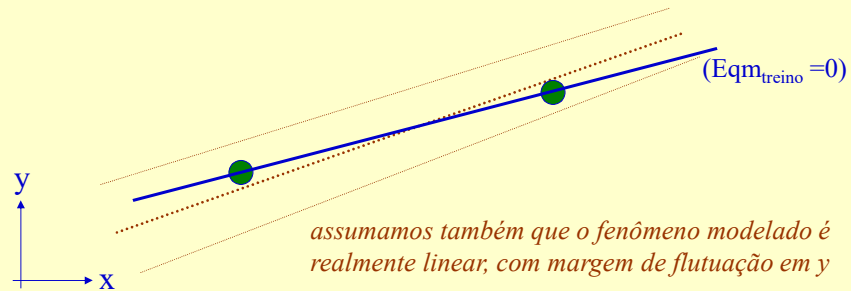
90

© Prof. Emilio Del Moral – EPUSP

90

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

façamos uso de modelagem linear ...



Os dados empíricos (x^u, y^u) 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

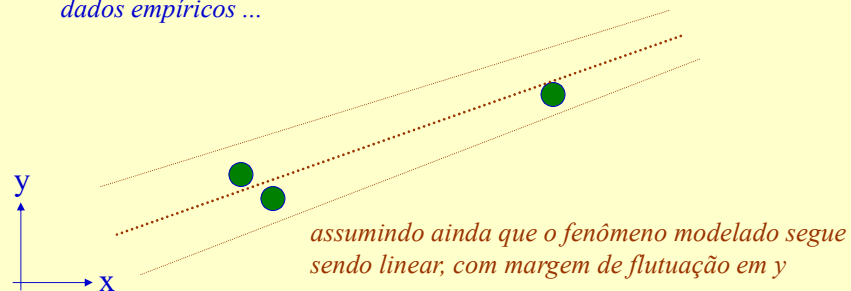
91

© Prof. Emilio Del Moral – EPUSP

91

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

Consideremos agora nova situação com mais dados empíricos ...



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

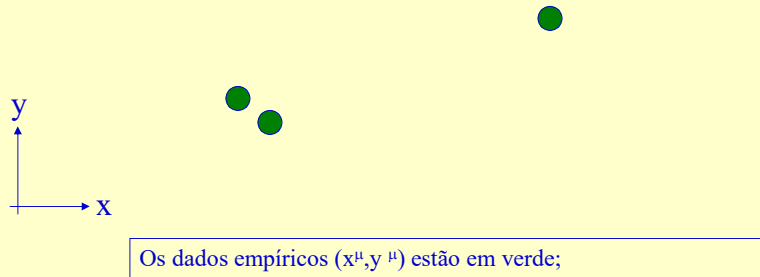
92

© Prof. Emilio Del Moral – EPUSP

92

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 ...



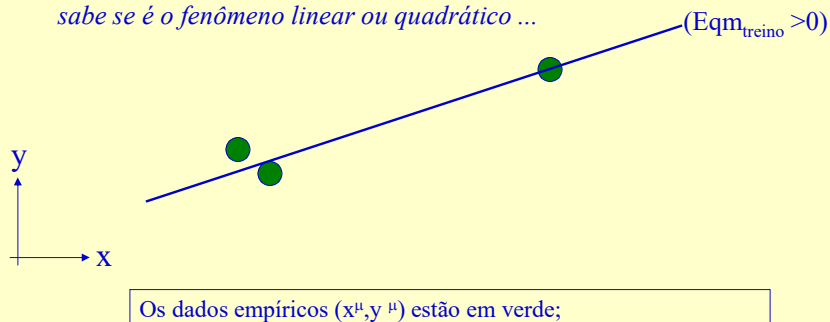
93

© Prof. Emilio Del Moral – EPUSP

93

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 ...



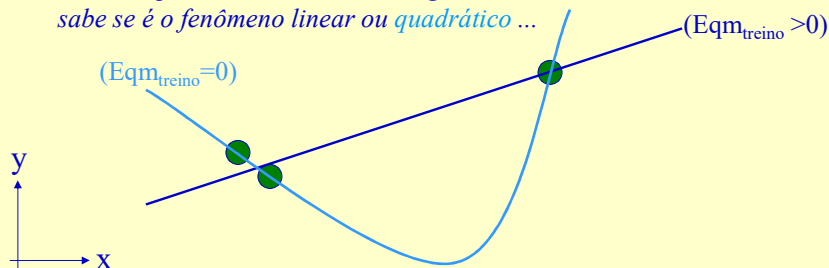
94

© Prof. Emilio Del Moral – EPUSP

94

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^u, y^u) 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

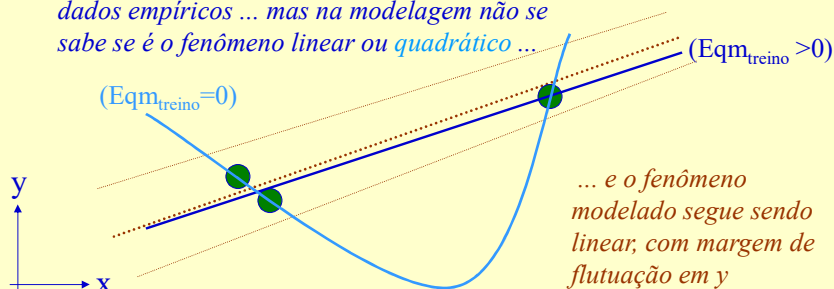
95

© Prof. Emilio Del Moral – EPUSP

95

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* ...



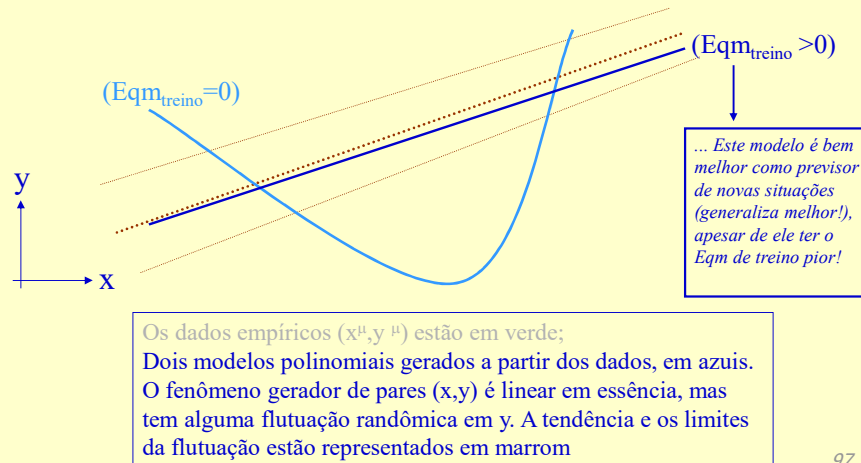
Os dados empíricos (x^u, y^u) estão em verde;
Dois modelos polinomiais gerados a partir dos dados, em azuis.
O fenômeno gerador de pares (x, y) é linear em essência, mas tem alguma flutuação randômica em y . A tendência e os limites da flutuação estão representados em marrom

96

© Prof. Emilio Del Moral – EPUSP

96

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



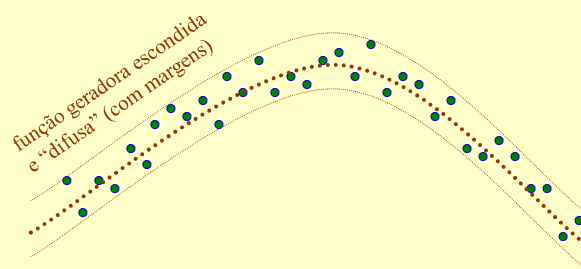
97

© Prof. Emilio Del Moral – EPUSP

97

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

um novo exemplo



Os dados empíricos (x^u, y^u) 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

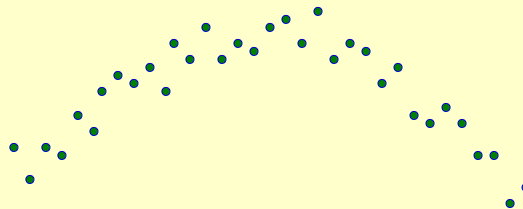
98

© Prof. Emilio Del Moral – EPUSP

98

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

um novo exemplo



Os dados empíricos (x^u, y^u) 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

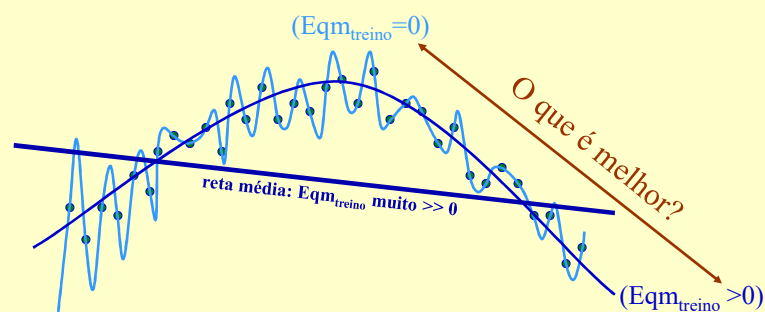
99

© Prof. Emilio Del Moral – EPUSP

99

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

um novo exemplo



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

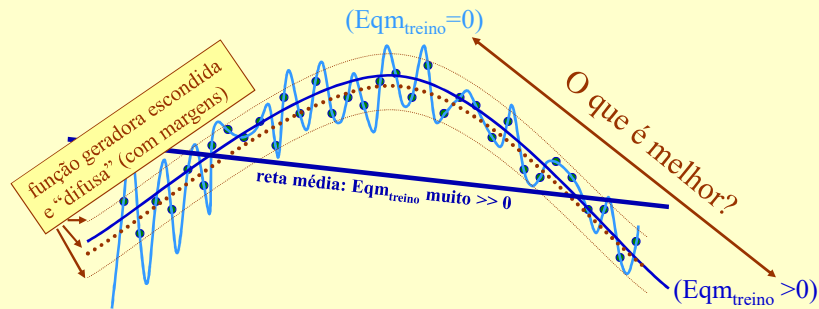
100

© Prof. Emilio Del Moral – EPUSP

100

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

um novo exemplo



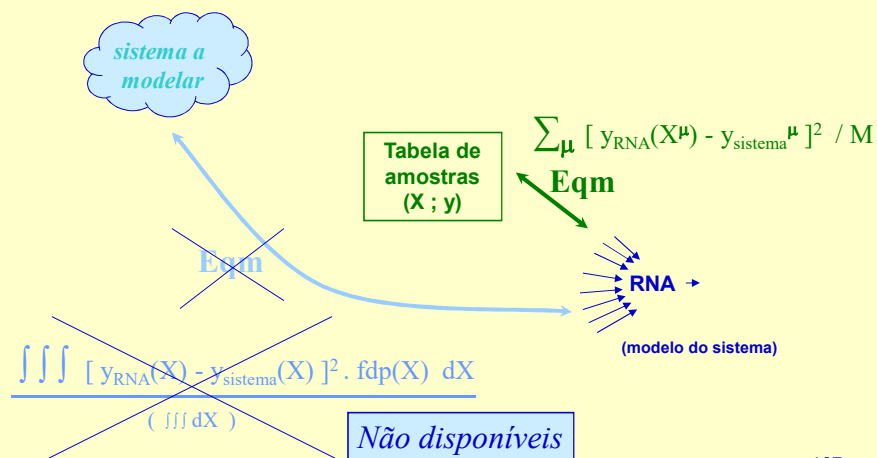
Os dados empíricos (x^μ, y^μ) estão em verde;
Três modelos polinomiais gerados a partir dos dados, em azuis.

101

© Prof. Emilio Del Moral – EPUSP

101

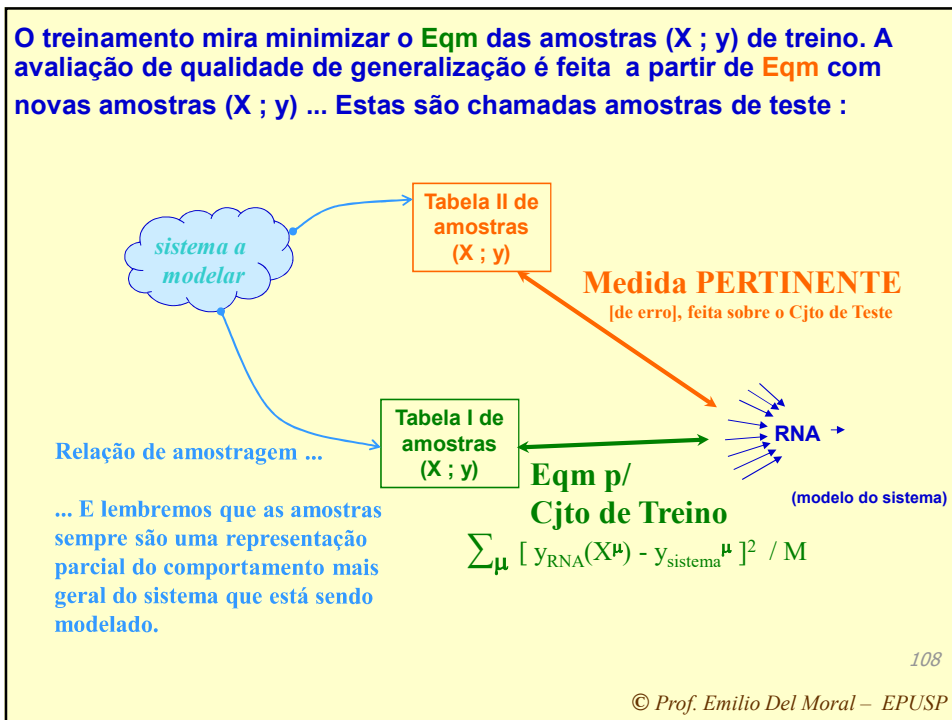
Sistema ... Amostras de treino ... RNA ...



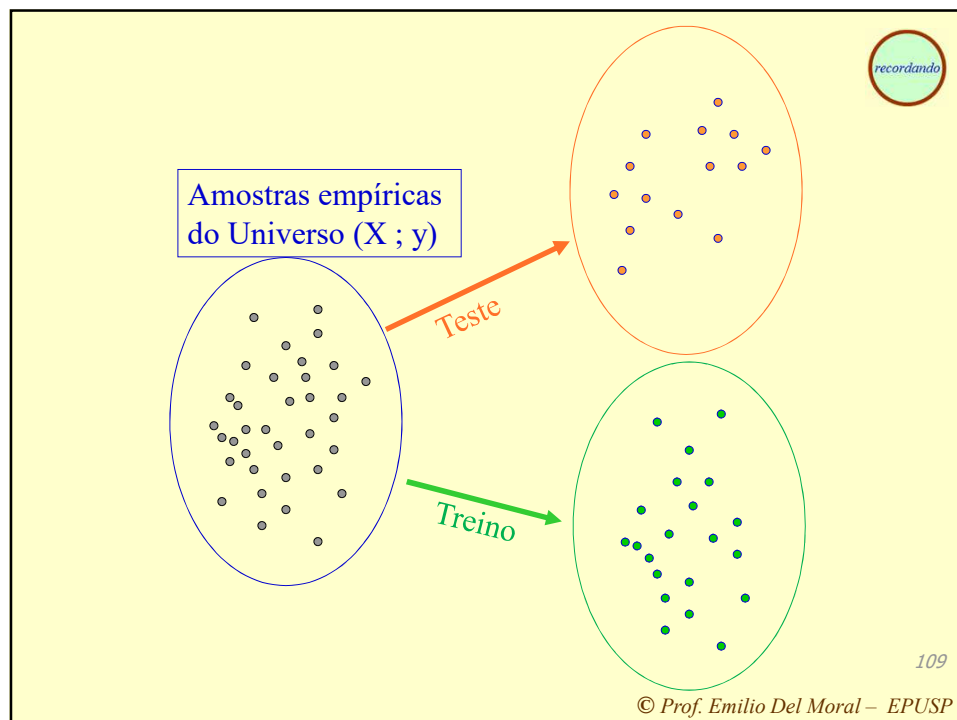
107

© Prof. Emilio Del Moral Hernandez

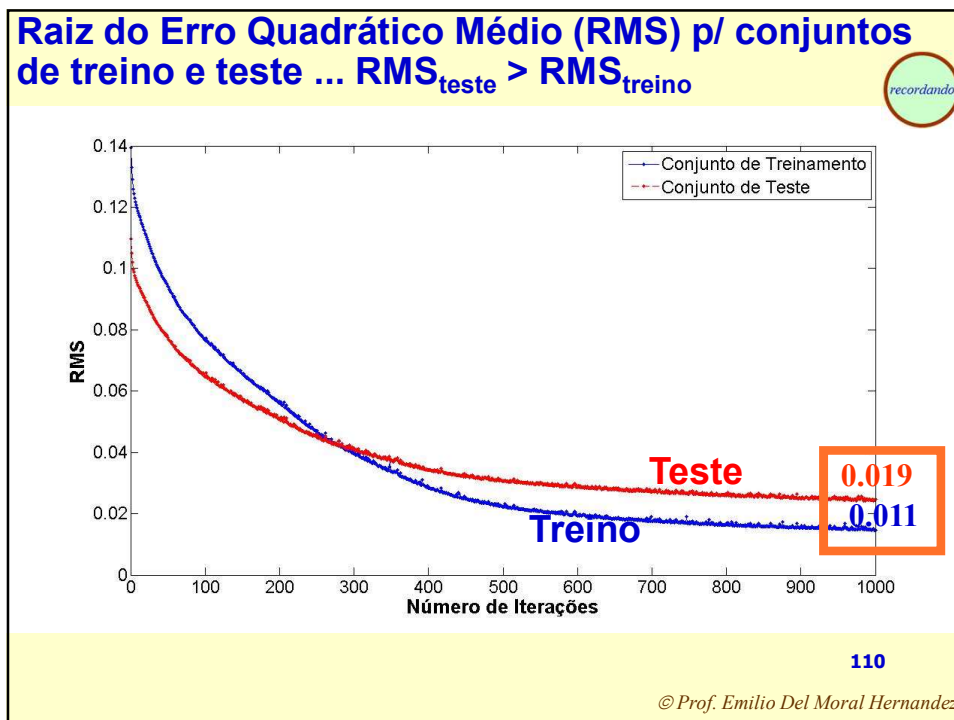
107



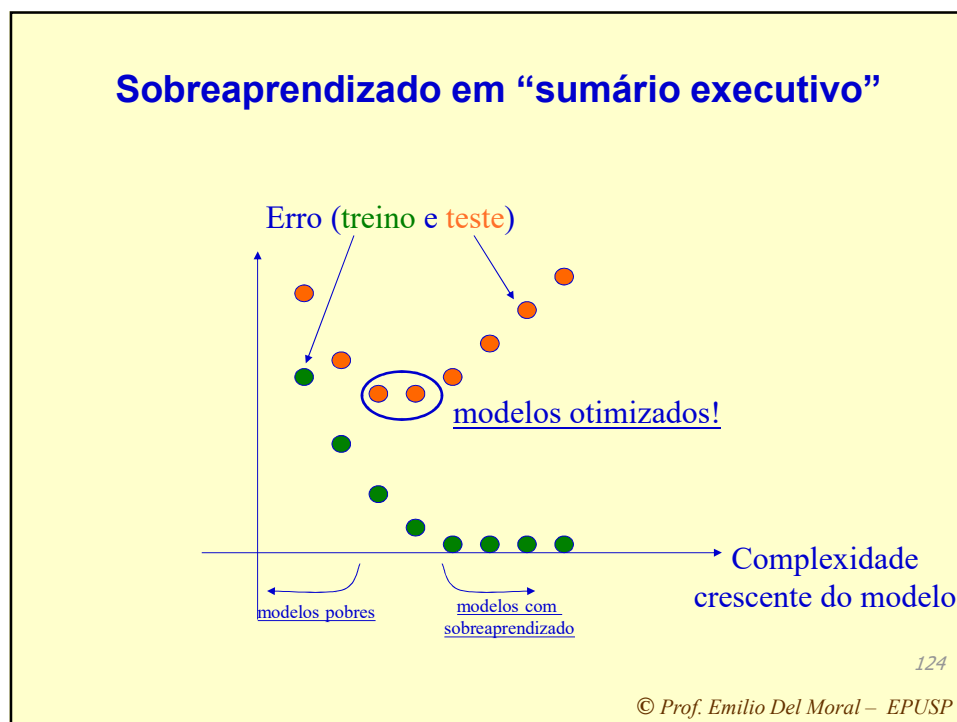
108



109



110



124