

*Um tema importante no curso:*

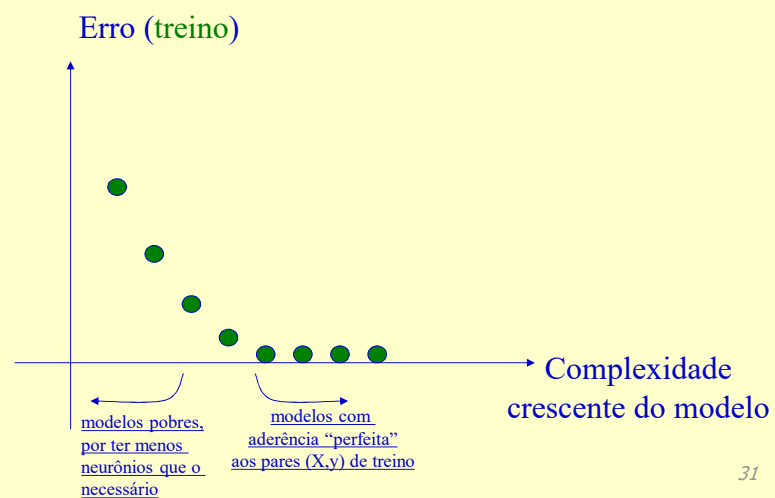
**Sobreapredizado / Sobreajuste**

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

© Prof. Emilio Del Moral – EPUSP

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

recordando

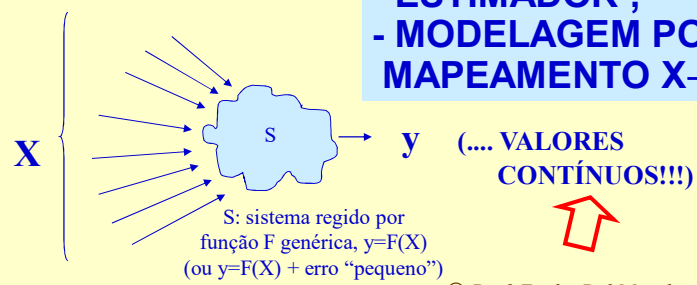
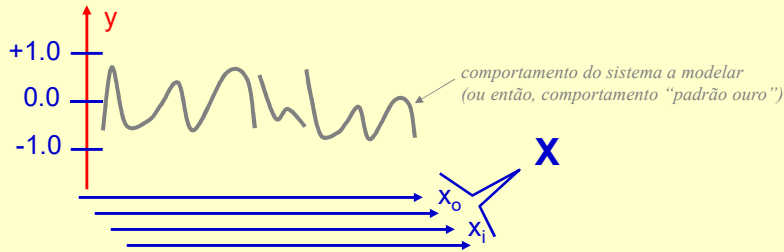


31

© Prof. Emilio Del Moral – EPUSP

## A função $y(X)$ "a descobrir", num caso geral de função analógica $y(X)$ ....

recordando

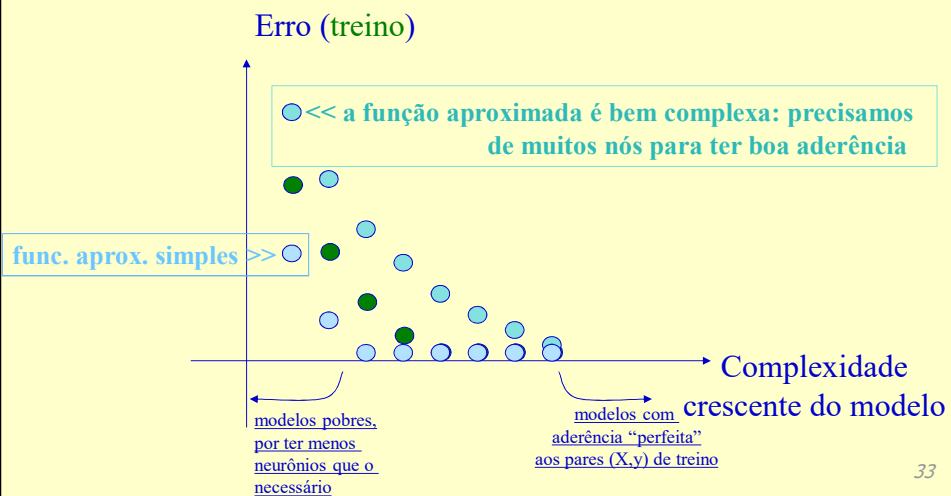


32

© Prof. Emilio Del Moral – EPUSP

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

recordando



33

© Prof. Emilio Del Moral – EPUSP

*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 ...*

35

© Prof. Emilio Del Moral – EPUSP

## **Sobreaprendizado:**

**Primeiro entendamos o conceito num contexto mais familiar (e simples), o de regressão polinomial univariada, para dados com tendência linear ou não linear de “grau” moderado, sujeitos a alguma flutuação não significativa em “y”**

...

**Depois, vocês pensem nos equivalentes dos nossos raciocínios, mas para as RNAs e outros tipos de modelos com número de parâmetros variável (complexidade variável) que você conheça ...**

37

© 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*

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

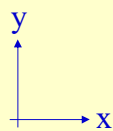
*falemos sobre regressão polinomial univariada, com o grau do polinômio aproximador podendo ser 1, 2, 3, ... 50, 51 etc.*

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

38

© Prof. Emilio Del Moral – EPUSP

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



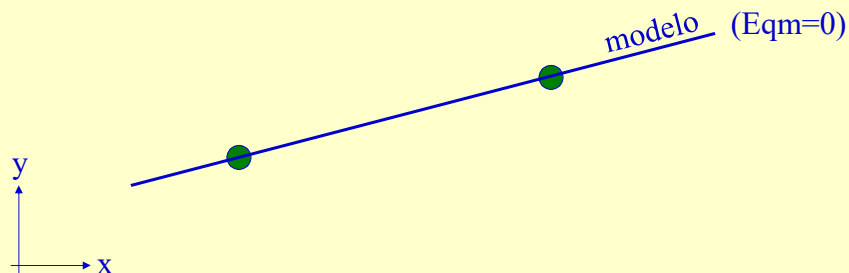
Os dados empíricos  $(x^i, y^i)$  estão em verde;

39

© Prof. Emilio Del Moral – EPUSP

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

*façamos uso de modelagem linear ...*



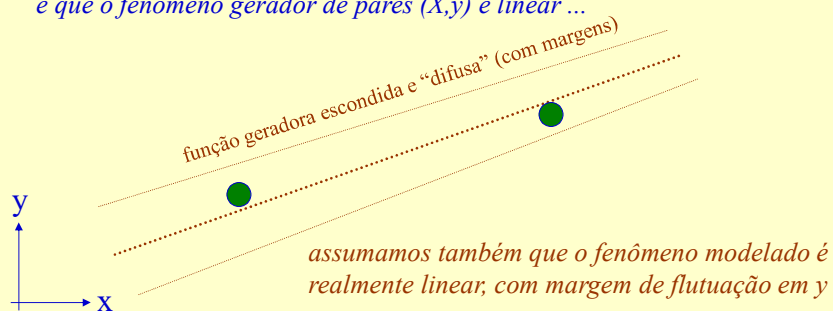
Os dados empíricos  $(x^i, y^i)$  estão em verde;  
O modelo linear gerado a partir dos dados, em azul.

40

© 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 ...*



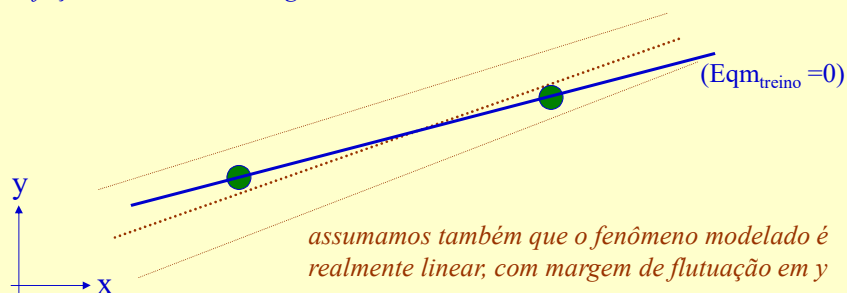
Os dados empíricos  $(x^i, y^i)$  estão em verde;  
O modelo linear gerado a partir dos dados, em azul.  
O fenômeno gerador de pares  $(x,y)$  é linear em essência, mas tem alguma flutuação randômica em  $y$ . A tendência e os limites da flutuação estão representados em marrom

41

© 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*

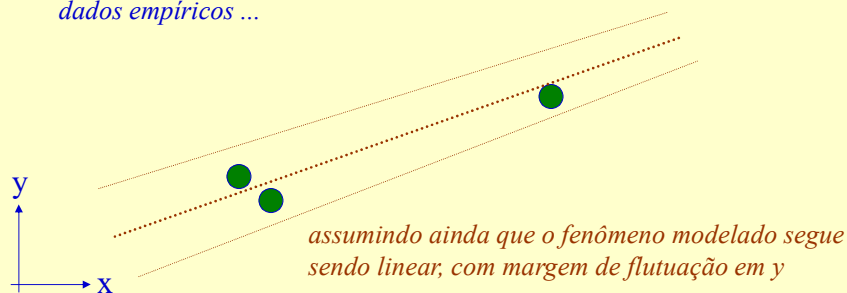
Os dados empíricos  $(x^h, y^h)$  estão em verde;  
O modelo linear gerado a partir dos dados, em azul.  
O fenômeno gerador de pares  $(x, y)$  é linear em essência, mas tem alguma flutuação randômica em y. A tendência e os limites da flutuação estão representados em marrom

42

© 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*

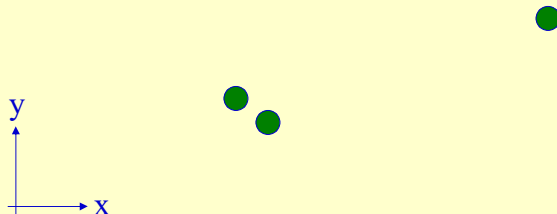
Os dados empíricos  $(x^h, y^h)$  estão em verde;

43

© 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 ...*



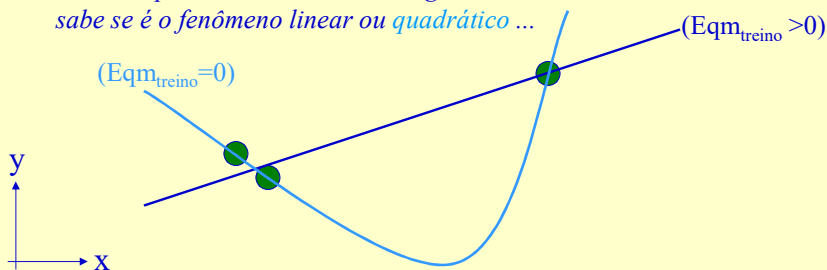
Os dados empíricos  $(x^i, y^i)$  estão em verde;

44

© 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 ...*



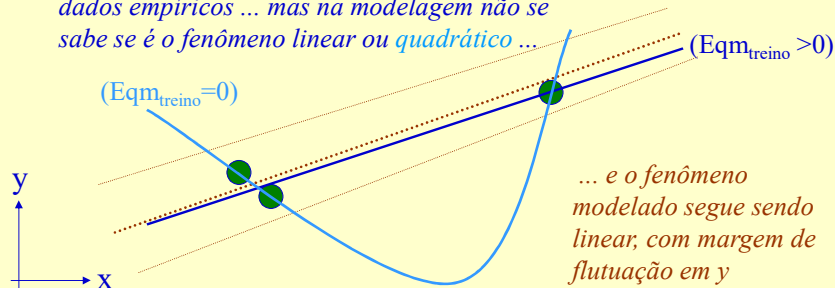
Os dados empíricos  $(x^i, y^i)$  estão em verde;  
Dois modelos polinomiais gerados a partir dos dados, em azuis.  
O fenômeno gerador de pares  $(x, y)$  é linear em essência, mas tem alguma flutuação randômica em  $y$ . A tendência e os limites da flutuação estão representados em marrom

45

© 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 ...

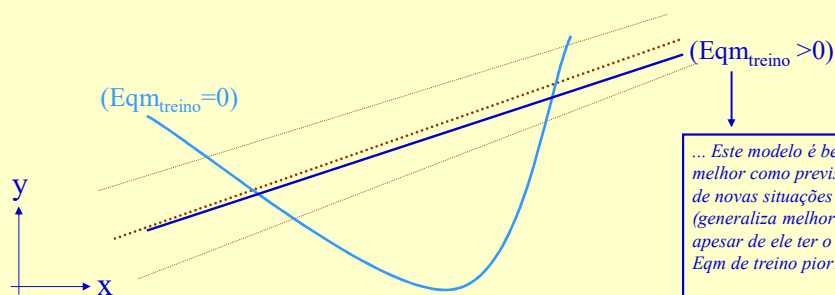


Os dados empíricos  $(x^i, y^i)$  estão em verde;  
Dois modelos polinomiais gerados a partir dos dados, em azuis.  
O fenômeno gerador de pares  $(x, y)$  é linear em essência, mas tem alguma flutuação randômica em y. A tendência e os limites da flutuação estão representados em marrom

46

© 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

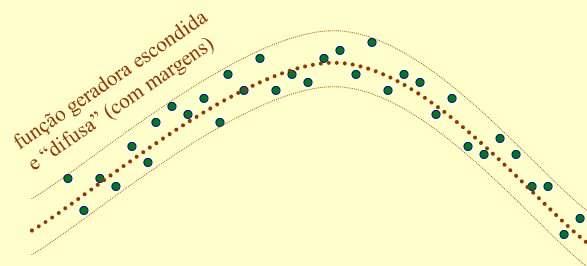
47

© 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;

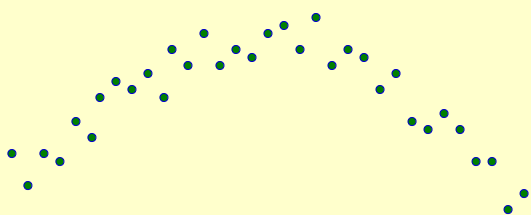
O fenômeno gerador de pares  $(x, y)$  é quadrático em essência, mas tem alguma flutuação randômica em  $y$ . A tendência e os limites da flutuação estão representados em marrom

48

© 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;

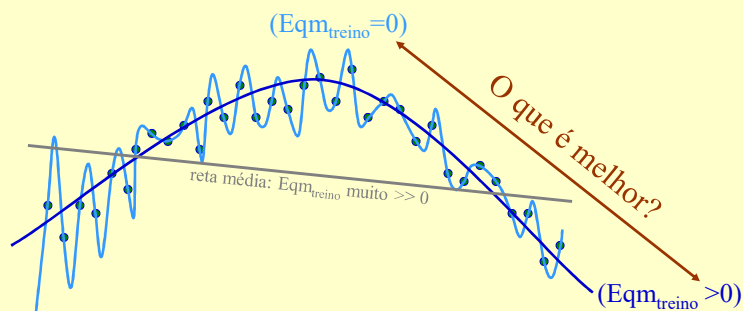
O fenômeno gerador de pares  $(x, y)$  é quadrático em essência, mas tem alguma flutuação randômica em  $y$ . A tendência e os limites da flutuação estão representados em marrom

49

© Prof. Emilio Del Moral – EPUSP

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

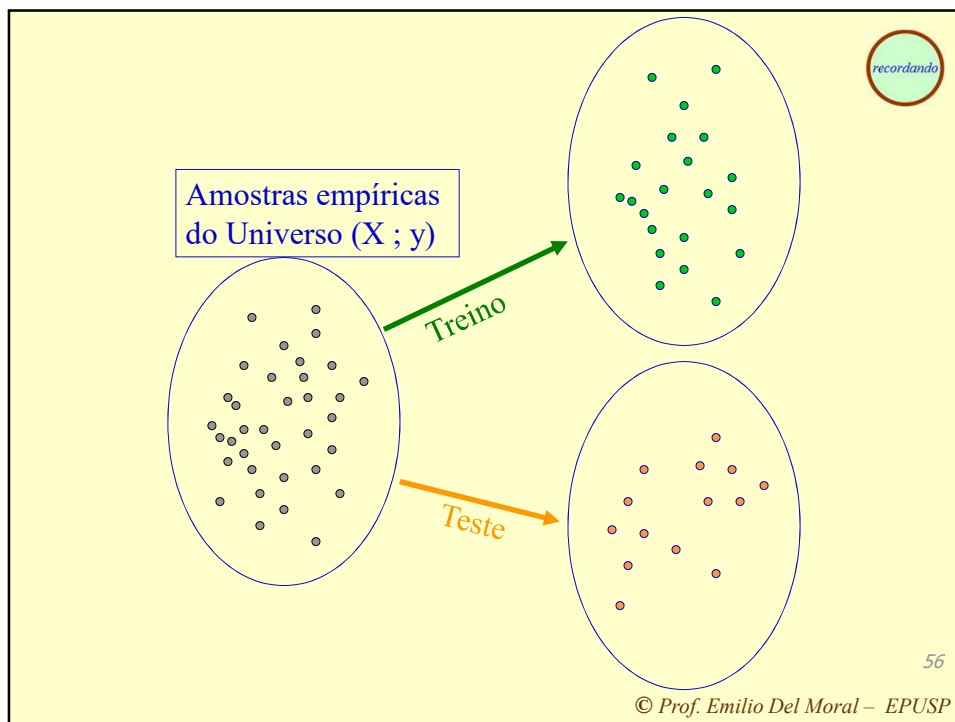
*um novo exemplo*



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

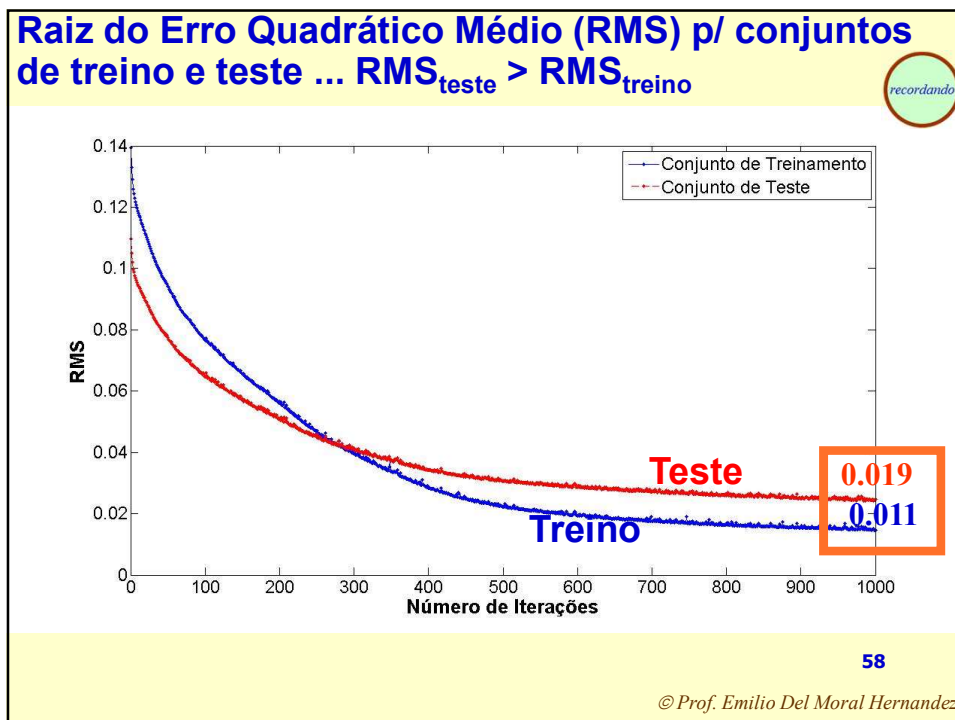
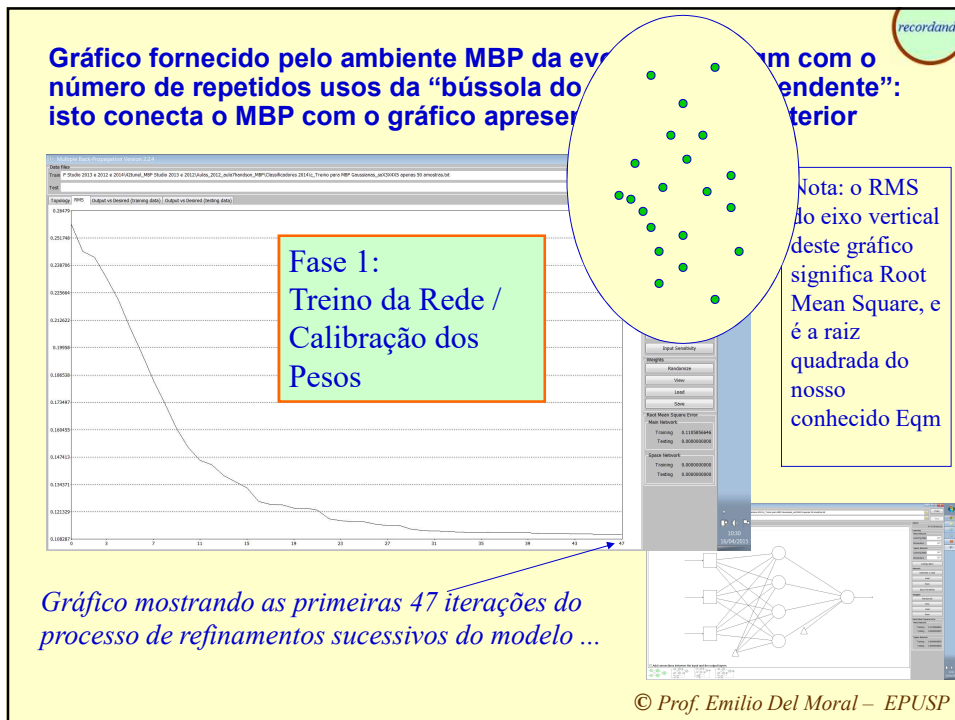
50

© Prof. Emilio Del Moral – EPUSP



56

© 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*

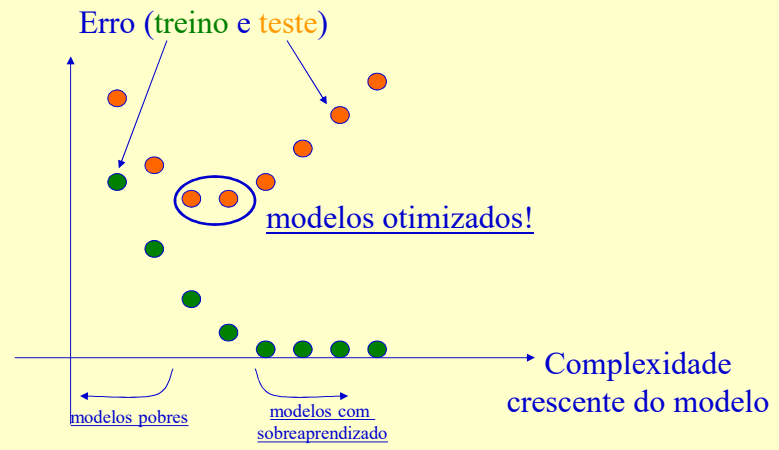
2) *Fase de TESTE / Caracterização da qualidade da RNA para generalizar: temos novos pares X e y, com y guardado "na gaveta", usado apenas para avaliação, não para re-calibração. É como um ensaio de uso final do modelo, com possibilidade de medir a sua qualidade com o y que foi guardado na gaveta.*

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

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

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

## Sobreaprendizado em "sumário executivo"



**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  $E_{qm}$
- É legítimo ficarmos com a rede otimizada de menos  $E_{qm}$  final? Ou deveríamos ter um ataque de avaliar o  $E_{qm}$  médio de várias tentativas?
- *Outro assunto ... e quanto à variação associada ao k-fold Cross Validation? Média ou melhor  $E_{qm}$ ?*

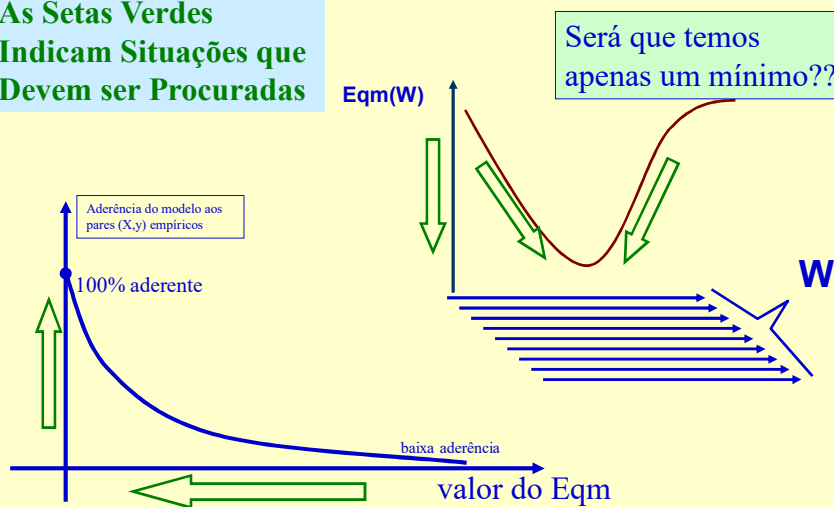
62

© 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**



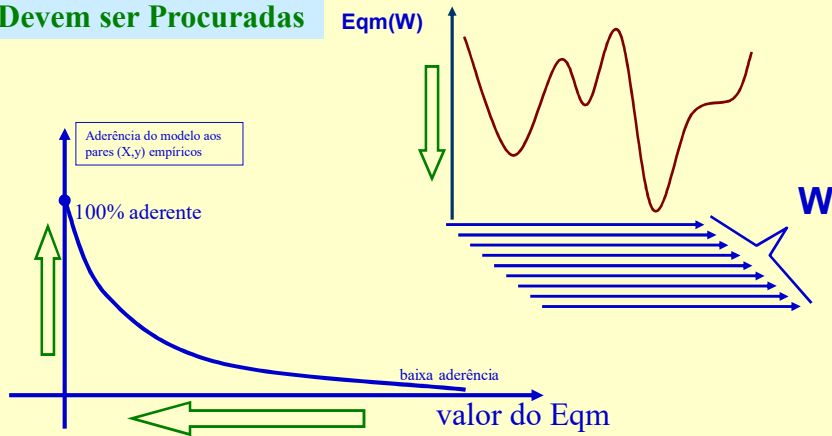
63

© 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*

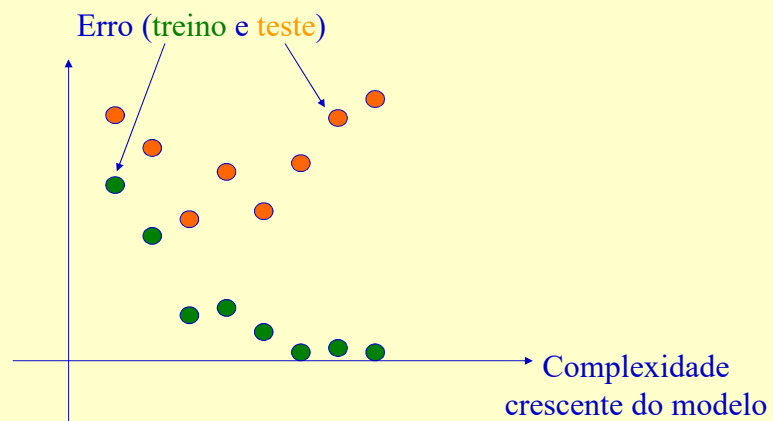
**As Setas Verdes Indicam Situações que Devem ser Procuradas**



64

© 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) ...



65

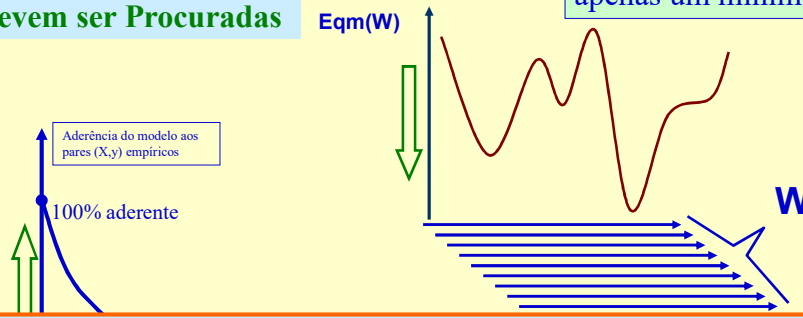
© Prof. Emilio Del Moral – EPUSP

## 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**

Será que temos apenas um mínimo??

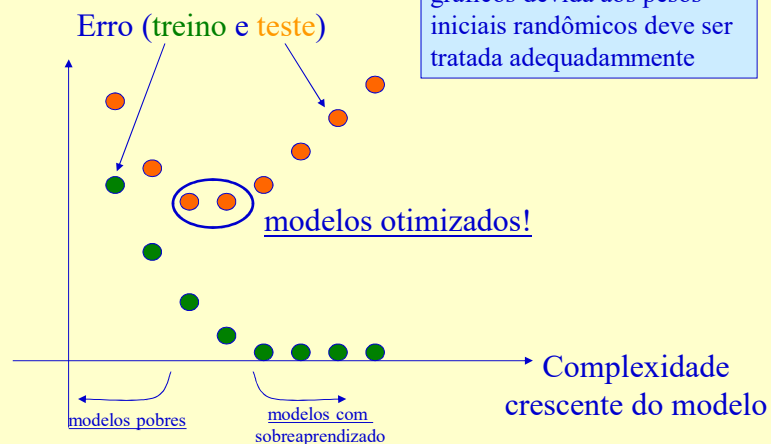


... Para não sermos reféns de mínimos locais com alto Eqm, podemos aplicar o gradiente descendente repetidamente na mesma RNA, com novos pesos iniciais randômicos em cada rodada, mantendo para o modelo final apenas os valores de pesos associados ao ensaio com o melhor dos resultados finais no Eqm!

© Prof. Emilio Del Moral – EPUSP

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

Ou seja, a flutuação nos gráficos devida aos pesos iniciais randômicos deve ser tratada adequadamente



67

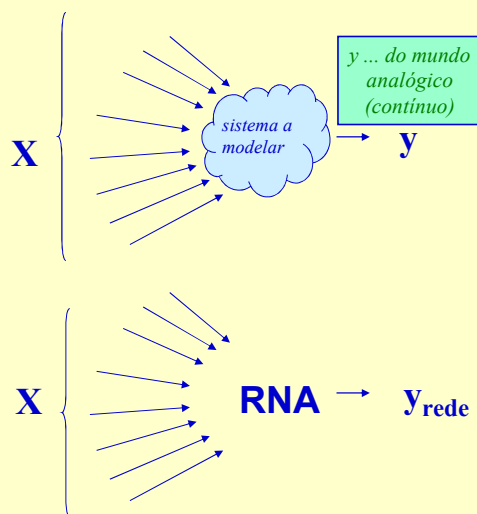
© Prof. Emilio Del Moral – EPUSP

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

72

© 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**

73

© Prof. Emilio Del Moral – EPUSP



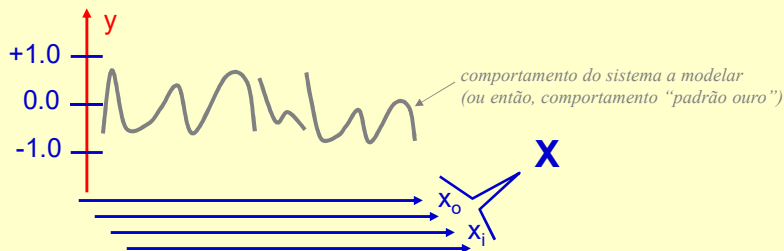
## Um hipotético universo de variáveis interdependentes, passível de modelagem/ens



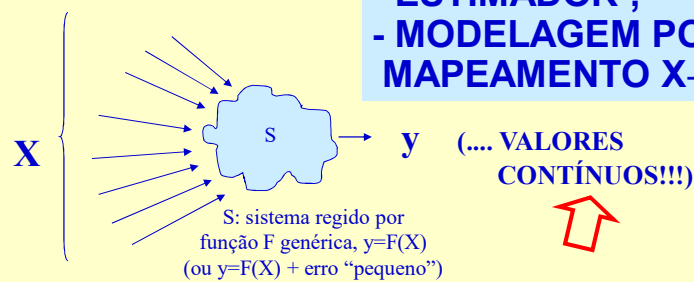
74

© Prof. Emilio Del Moral – EPUSP

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



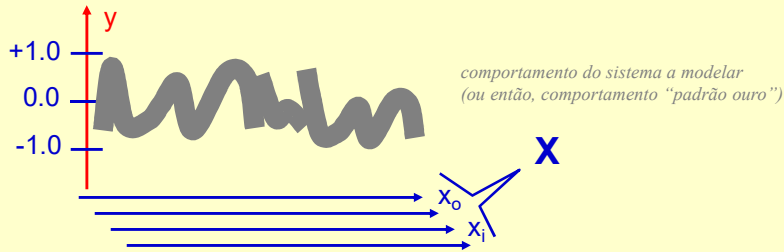
**- ESTIMADOR ;  
- MODELAGEM POR  
MAPEAMENTO  $X \rightarrow y$**



75

© Prof. Emilio Del Moral – EPUSP

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



.... 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  $y$ :  $y = F_{\text{médio}}(X) + \text{flutuação}$ .

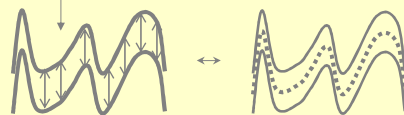
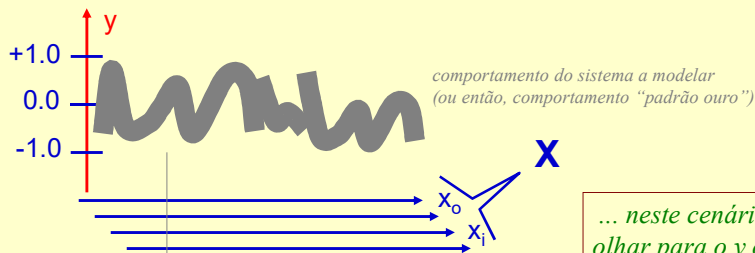
Neste cenário, buscamos que o modelo capture o comportamento médio das relações observadas entre  $X$  e  $y$ :

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

76

© Prof. Emilio Del Moral – EPUSP

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



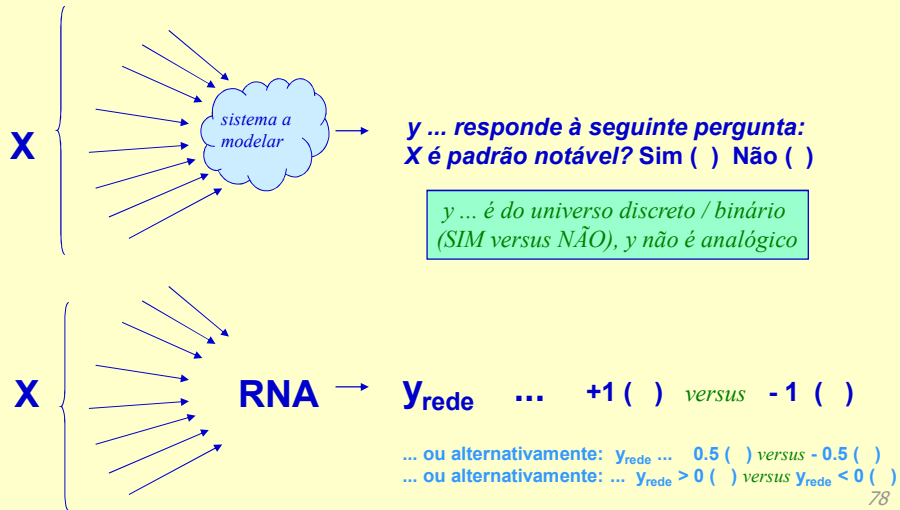
... neste cenário, podemos olhar para o  $y$  observado no sistema que se deseja modelar não mais como um valor específico bem definido, mas como um valor médio esperado (dado valor de  $X$ ) e uma faixa de valores em torno desse valor médio esperado.

77

© Prof. Emilio Del Moral – EPUSP

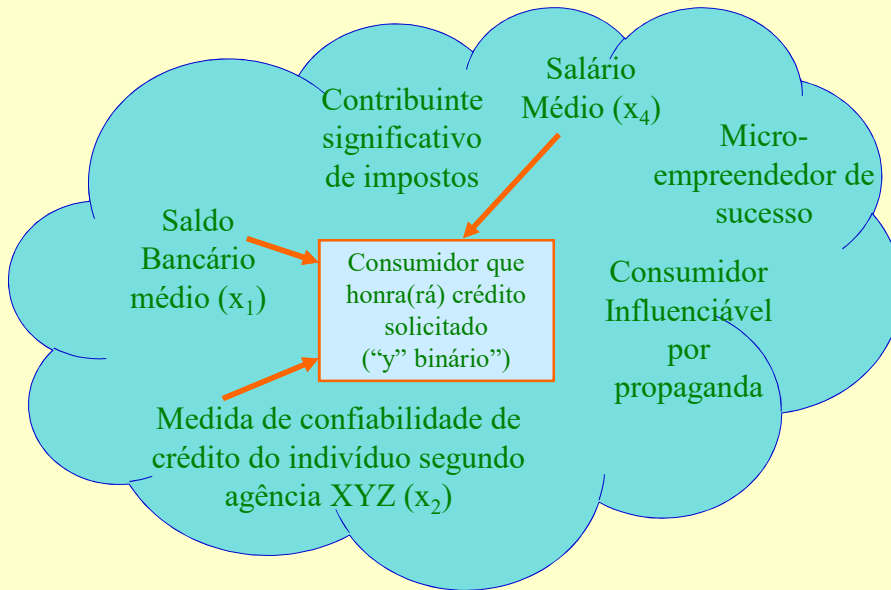
## RNAs como reconhecedor / detetor de padrões

...



© Prof. Emilio Del Moral – EPUSP

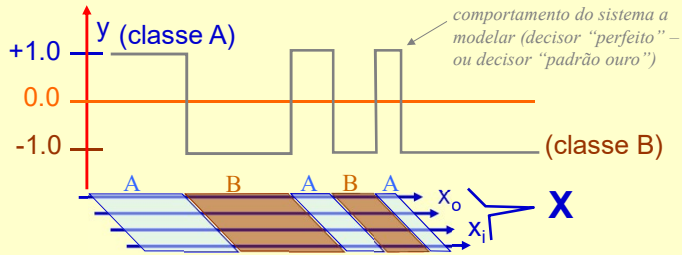
## Um hipotético universo de variáveis interdependentes, passível de modelagem/ens



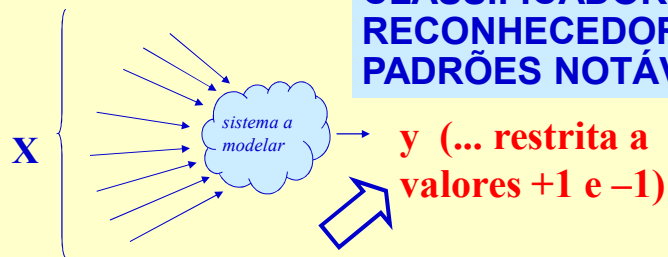
79

© Prof. Emilio Del Moral – EPUSP

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



**CLASSIFICADOR;  
RECONHECEDOR DE  
PADRÕES NOTÁVEIS**

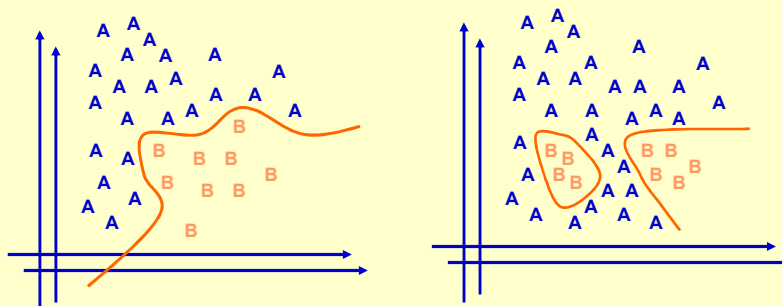


80

© Prof. Emilio Del Moral – EPUSP

## Capacidade de reconhecimento de padrões em casos complexos NÃO LINEARES

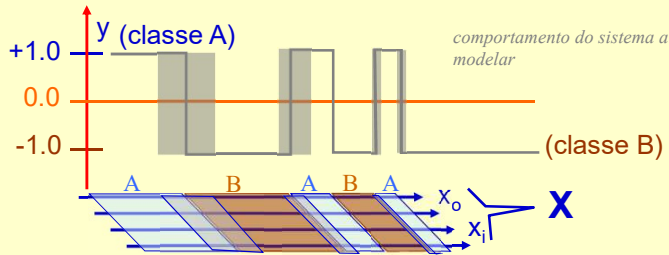
*Com as RNAs, a hipersuperfície de separação entre classes vai muito além dos hiperplanos*



81

© Prof. Emilio Del Moral – EPUSP

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



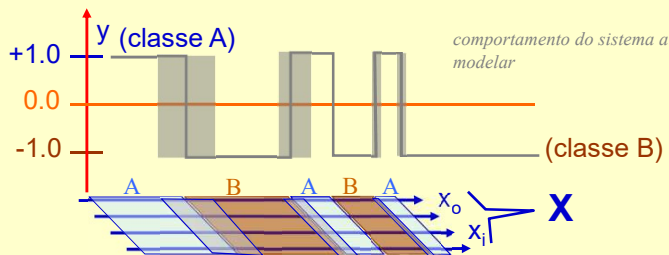
.... em problemas concretos / reais, há sempre alguma ambiguidade no mapeamento que leva valores de X aos valores discretos de y. Não temos uma função  $y=F(X)$  no sentido matemático exato, pois para uma dada ênupla de valores X fixado temos em alguns casos de fronteira a possibilidade de observar no 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

82

© Prof. Emilio Del Moral – EPUSP

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

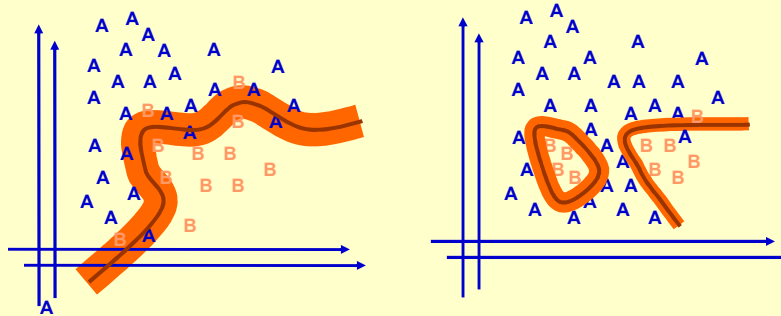


... podemos olhar para o y (classe A ou B) observado no sistema que se deseja modelar não mais como uma classe sempre bem definida e com fronteiras de separação entre A e B bem definidas no espaço de valores de X, mas como sendo delineadas na modelagem através de fronteiras com eventuais faixas de tolerância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

83

© Prof. Emilio Del Moral – EPUSP

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



84

© Prof. Emilio Del Moral – EPUSP

*Ciclos de refinamento num projeto de modelagem para reconhecimento e/ou regressão, com duração de médio / longo prazo e com interesse de otimização de desempenho máxima*

© 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*

90

© Prof. Emilio Del Moral – EPUSP

## ... Aspectos de refinamentos nos projetos ...

### Diferentes aspectos estudados no treinamento / otimização / caracterização da RNA ...

- Vários Delta W sequenciados (gradiente descendente)
- Re-sorteios de pesos iniciais (fugindo de mínimos locais)
- k-fold cross validation (avaliando sensibilidade aos dados empíricos)
- Diferentes graus de complexidade do modelo neural (evitando sobreaprendizado)
- Ensaio com vários Pré-Processamentos alternativos (aumentando desempenho)
- Descarte de algumas variáveis de menor relevância, p/ melhora do desempenho
- Aumento de M+M' c/ as mesmas variáveis (há custo extra com novas coletas X;y)
- Acréscimo de variáveis x incluídas no modelo (há custo extra com novas coletas)
- ... Outros ...

91

© Prof. Emilio Del Moral – EPUSP

## ... Multi-ciclos de refinamentos diversos que são aplicados em projetos mirem a otimização de desempenho ao máximo ...

### Alguns dos loops de treinamento e/ou otimização e/ou caracterização:

- Loop dos Delta W sequenciados (já é intrínseco ao gradiente descendente)
- Loop dos re-sorteios de pesos iniciais (fugindo de mínimos locais)
- Loop do k-fold cross validation (avaliando e medindo a sensibilidade aos dados empíricos)
- Loop de aumento gradual de complexidade do modelo neural (evitando sobreaprendizado)
- Ensaio com vários Pré-Processamentos alternativos (aumentando desempenho)
- Descarte de algumas variáveis de menor relevância (limitando ruído e sobreaprendizado)
- Aumento de  $M+M' c/$  as mesmas variáveis (há custo extra com novas coletas  $X; y$ )
- Acréscimo de variáveis  $x$  incluídas no modelo (há custo extra com novas coletas)

93

© Prof. Emilio Del Moral – EPUSP

## ... Multi-ciclos de refinamentos nos projetos ...

### Discutamos os loops de treinamento e/ou otimização e/ou caracterização ...

- 1- Loop dos Delta W sequenciados (intrínseco ao gradiente descendente)
- 2- Loop dos re-sorteios de pesos iniciais (fugindo de mínimos locais)
- 3- Loop do k-fold cross validation (avaliando sensibilidade aos dados empíricos)
- 4- Loop de aumento gradual de complexidade do modelo neural (evitando sobreaprendizado)
- 5- Ensaio com vários Pré-Processamentos alternativos (aumentando desempenho)
- 6- Descarte de algumas variáveis de menor relevância,  $p/$  melhora do desempenho
- 7- Aumento de  $M+M' c/$  as mesmas variáveis (há custo extra com novas coletas  $X; y$ )
- 8- Acréscimo de variáveis  $x$  incluídas no modelo (há custo extra com novas coletas)

*Essa ordem de aninhamentos se aplica bem ao seu projeto?  
Ou há outra/s ordem/ns que faz/em mais sentido para as SUAS circunstâncias específicas?*

94

© Prof. Emilio Del Moral – EPUSP