

Carlos E. Pedreira¹, Luiz Carlos Pedroza² e Mayte Fariñas¹

¹ DEE PUC-RIO CP. 38063 Rio de Janeiro, CEP 22452-970

² CEFET-RJ Av. Maracanã, 229 Rio de Janeiro, CEP 20271-110

E-mails: pedreira@ele.puc-rio.br, pedroza@cefet-rj.br, mayte@ele.puc-rio.br.

Abstract:

In this paper a new connectionist architecture is proposed. The proposed architecture is trained by a scheme based on partition of the function domain, approximating the generator function by a set of very simple supporting functions. This method has an interesting ability concerning interpolation. A synthetic experiment and a real data missing data application are presented.

1. Introdução

Muitos problemas de interesse prático estão relacionados a interpolação dos dados. Entre estes, de particular relevância está o problema de preenchimento de dados faltantes. Neste caso, o que se procura é uma forma de emular uma função em um intervalo do domínio onde apenas uma parte dos pontos é conhecida. A preocupação aqui não é a capacidade preditiva da rede, isto é, estimar pontos fora do domínio, mas sim de reproduzir a função da melhor forma possível dentro do domínio estabelecido. Neste artigo, propõe-se um algoritmo capaz de reconstruir a função a partir de estimativas locais ao longo do domínio de interesse.

A arquitetura conexonista proposta é não usual e utiliza uma metodologia de treinamento baseada no particionamento do domínio da função a ser emulada. Este função é aproximada, por um conjunto de funções de apoio muito simples, muitas vezes lineares. A idéia central é expressar o mapeamento entrada-saída através de uma função composta por partes. A estrutura básica é constituída pela combinação de vários pares compostos de funções de aproximação e funções de pertinência. As funções de pertinência definem em cada trecho do domínio a participação da função de ativação a essa associada. É possível a ocorrência de sobreposições parciais das funções de pertinência proporcionando uma maior riqueza do mapeamento pretendido. Desse modo o problema de aproximação de funções é focado especializando-se grupos de neurônios, formados pelos pares anteriormente descritos, que emulam a função geradora em cada setor do domínio. O grau de especialização em um determinado trecho é dado pelo nível da função de pertinência. Por exemplo, em um trecho aonde apenas

uma das funções de pertinência assume valor alto, haverá uma dominância da função de aproximação associada a esta.

2. A arquitetura Proposta

Consideremos uma rede com m nós ou neurônios. Seja $\{x_i\}_1^n$ a partição dos dados usada para treinamento. Por simplicidade algébrica e de notação iremos considerar o caso onde $x \in \mathbb{R}$ (o subscrito de x será omitido), a generalização para o caso onde $x \in \mathbb{R}^n$ é algebricamente direta. Define-se, para cada ponto x, m funções de pertinência do seguinte modo:

$$B_j(x) = -C_j \left[\frac{1}{1 + \exp(d_j(x - h_j^{(1)}))} - \frac{1}{1 + \exp(d_j(x - h_j^{(2)}))} \right],$$

$j=1, \dots, m$

onde C_j , d_j , $h_j^{(1)}$ e $h_j^{(2)}$ são parâmetros a serem ajustados. Note-se que o parâmetro C_j reflete o nível da função de pertinência, enquanto d_j está relacionado à declividade desta função. Os parâmetros $h_j^{(1)}$ e $h_j^{(2)}$ delimitam o setor do domínio no qual a função de aproximação associada a esta função de pertinência é mais ativa. (ver figura 1).

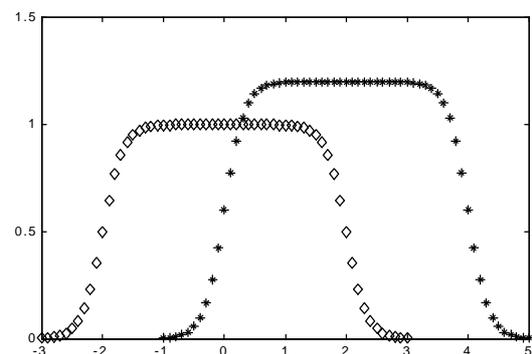


Figura 1 – Exemplos de funções de ativação

◇ $C=1$, $d=6$; $h_j^{(1)} = -2$; $h_j^{(2)} = 2$

* $C=1.2$, $d=6$; $h_j^{(1)} = 0$; $h_j^{(2)} = 4$

As funções de aproximação são tipicamente funções lineares ou quadráticas. Embora funções mais complexas possam ser usadas sem prejuízo da estrutura

teórica proposta, não parecem trazer contribuição significativa ao modelo. Consideremos funções de aproximação lineares:

$$\kappa_j(x) = a_j x + b_j \quad j=1, \dots, m$$

onde a_j e b_j são os parâmetros a serem estimados. Cada nó, ou neurônio, da rede é constituído de um par { função de pertinência ; função de aproximação } (ver figura 2). Então, para cada nó é necessário estimar 6 parâmetros (7 no caso de funções de aproximação quadráticas). Como usualmente, o número de nós indica a complexidade do modelo.

As entradas são conectadas ao nó onde é efetuado o produto da função de pertinência $B_j(x)$ e da função de aproximação $\kappa_j(x)$. A saída da rede é um somatório da saída de cada um destes nós. Note que não há pesos ligando a saída dos nós a saída da rede (veja figura 2). Deste modo a saída do j -ésimo nó é $B_j(x) \kappa_j(x)$, e a saída da rede é dada por:

$$g^m(x) = \sum_{j=1}^m B_j(x) \kappa_j(x) \quad (2.1)$$

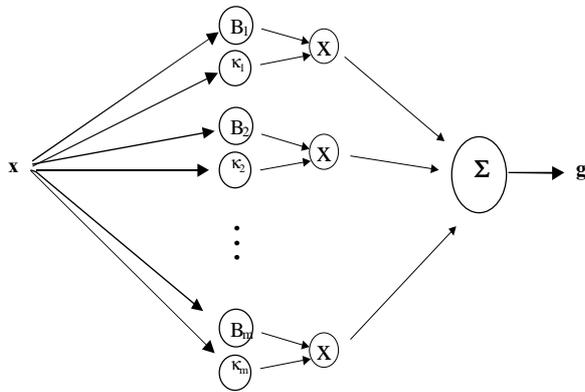


Figura 2 A arquitetura proposta

Definindo-se $B_{\text{vet}}(x) \equiv (B_1(x), B_2(x), \dots, B_m(x))$ e $\kappa_{\text{vet}}(x) \equiv (\kappa_1(x), \kappa_2(x), \dots, \kappa_m(x))$ pode-se escrever a saída da rede em forma de produto interno, i.e. $g^m(x) = \langle B_{\text{vet}}(x), \kappa_{\text{vet}}(x) \rangle$.

O objetivo central é projetar uma rede cuja saída aproxime uma função alvo, $f(x)$, da melhor forma possível. Define-se então, uma função de erro como uma combinação convexa do quadrado de duas medidas de erro E_1 e E_2 :

$$E \equiv a \sum_{i=1}^k E_1^2(x_i) + (1-a) E_2^2 \quad (2.2)$$

onde

$$E_1(x_i) \equiv g^m(x_i) - y(x_i) \quad \text{e} \quad E_2 \equiv 1 - \sum_{j=1}^m C_j \quad (2.3)$$

O termo E_1 está associado a qualidade da aproximação obtida enquanto E_2 é usado com a finalidade de manter as funções de pertinência limitadas. Na realidade se está penalizando soluções nas quais o somatório destas funções excede 1. A

escolha da unidade como valor limitante não é necessária embora confira mais interpretabilidade aos resultados.

Define-se, para cada neurônio, um vetor de parâmetros $\mathfrak{J}^j \equiv (C_j, d_j, h_j^{(1)}, h_j^{(2)}, a_j, b_j)$ e o objetivo central será encontrar \mathfrak{J}^j que minimize a função de erro E .

3. Um Resultado teórico

Um resultado interessante apresentado no teorema T-1 dá consistência teórica a metodologia proposta. Neste teorema mostra-se que qualquer função L^2 -integrável pode ser aproximada por funções da forma $g^m(x)$. Dois resultados auxiliares, apresentados a seguir facilitam a prova deste resultado.

Resultado auxiliar RA-1: As funções simples, da forma $S(x) = \sum_{j=1}^m a_j X_{A_j}(x)$ com α_j constante real e

$X_{A_j}(x)$ a função indicadora do conjunto A_j (isto é, toma-se o valor 1 se $x \in A_j$ e zero caso contrário) são densas em L^2 [3].

Resultado auxiliar RA-2: As funções $\{g^m(x)\}$ aproximam qualquer função simples. Isto é, $\forall S(x)$, existe uma seqüência de funções $g_n(x) \in \{g^m(x)\}$ tal que $g_n(x) \rightarrow S(x)$, onde esta convergência é referida à convergência em L^2

Prova do RA-2:

Seja $S(x) = X_{[0, 1]}$. Consideremos a seqüência onde $\kappa_n(x) = 1$, isto é: $g_n(x) = B_n(x)$ com $h_1=0$ $h_2=1$.

$$B_n(x) = -C_n \left[\frac{1}{(1 + e^{d_n x})} - \frac{1}{(1 + e^{d_n(x-1)})} \right]$$

onde $C_n = (e^{d_n/2} + 1) / (e^{d_n/2} - 1)$ e $d_n \rightarrow \infty$ (para garantir que o máximo de $B_n(x)$ em $[0, 1]$ seja igual a 1)

1) É fácil provar que, $B_n(x) \xrightarrow[n]{} S(x)$ pontualmente, isto é $B_n(x) \rightarrow \begin{cases} 0 & x \notin [0, 1] \\ 1 & x \in [0, 1] \end{cases}$

Para chegar a convergência em L^2 , partindo-se da convergência pontual, pode-se utilizar o teorema da convergência dominada de Lebesgue: (Se $f_n(x) \leq g(x) \in L^2$, e $f_n(x) \rightarrow f(x)$ pontualmente então $f_n(x) \rightarrow f(x)$ em L^2) [3].

Como já se tem a convergência pontual, basta provar que $B_n(x)$ é dominada por uma função L^2 -integrável.

É fácil verificar que a função L^2 -integrável $g(x)$ definida como:

Caso contrário, se algum m_i é complexo, dividir o intervalo em m intervalos uniformes disjuntos.

Passo 3: Para cada intervalo, ajustar regressão linear nos pares (x_i, y_i) , $x_i \in I_i$ para obter a_i e b_i .

5. Resultados numéricos

O primeiro experimento numérico (veja figuras 3a; 3b; 3c; 3d) consiste de 100 pontos gerados pela função $\text{sen}(x)+2$ com diferentes níveis de ruído. O ruído foi gerado adicionando-se a função geradora um sinal Gaussiano de média zero e desvio padrão igual a 0.1; 0.4 e 0.7 em cada uma dos 3 experimentos. A tabela 5.1 resume os resultados destes experimentos. Em cada uma destas simulações foram utilizadas 3 pares de funções de pertinência - aproximação, parâmetros de tolerância. Em todos os experimentos, exceto no experimento 4 foi usada a heurística exposta na seção 4 para inicialização.

No	Nível de ruído	No. de épocas	Mape do ruído	Mape treinamento	Mape generalização
1	0	111	0	0.14	0.157
2	0.1	137	4.79	4.64	1.04
3	0.4	78	16.13	16.80	4.91
4	0.7	37	30.88	52.36	7.87

Tabela 5.1 Resultados das simulações com a função $\text{sen}x+2$ com vários níveis de ruído.

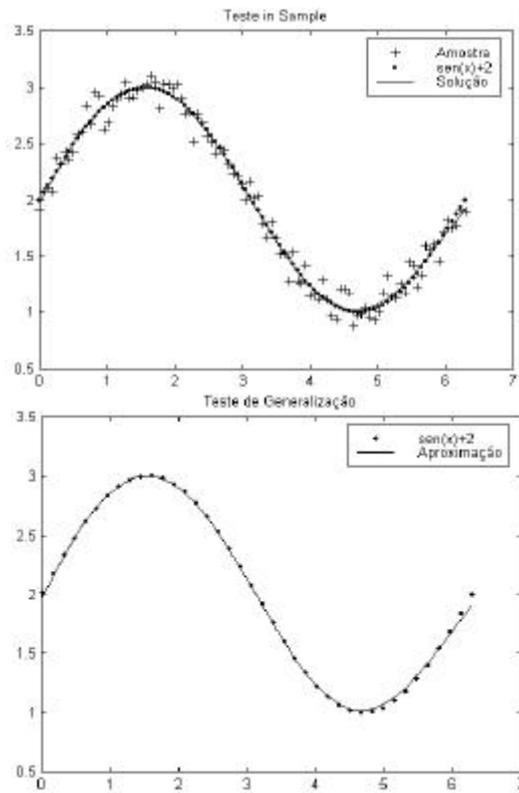


Figura 3b - Nível de ruído: 0.1

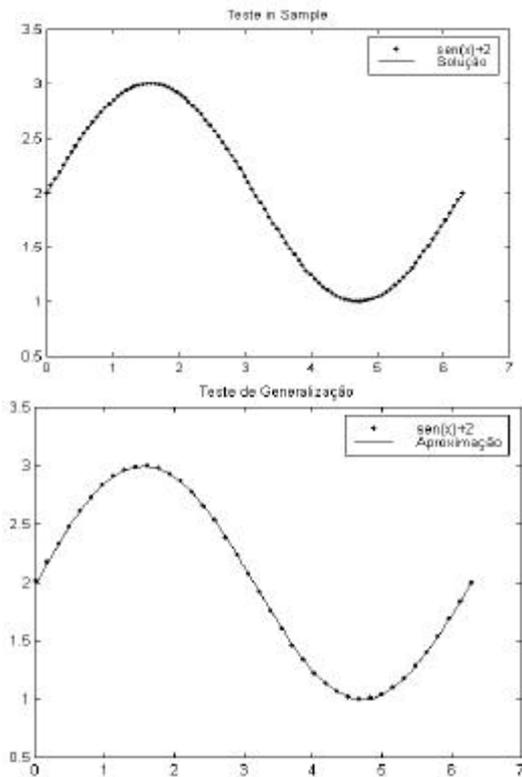


Figura 3a - Nível de ruído: 0

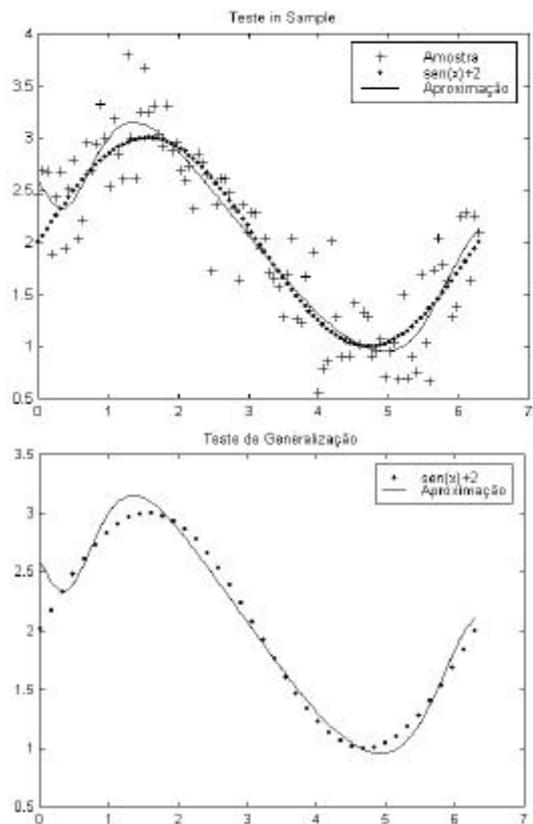


Figura 3c - Nível de ruído: 0.4

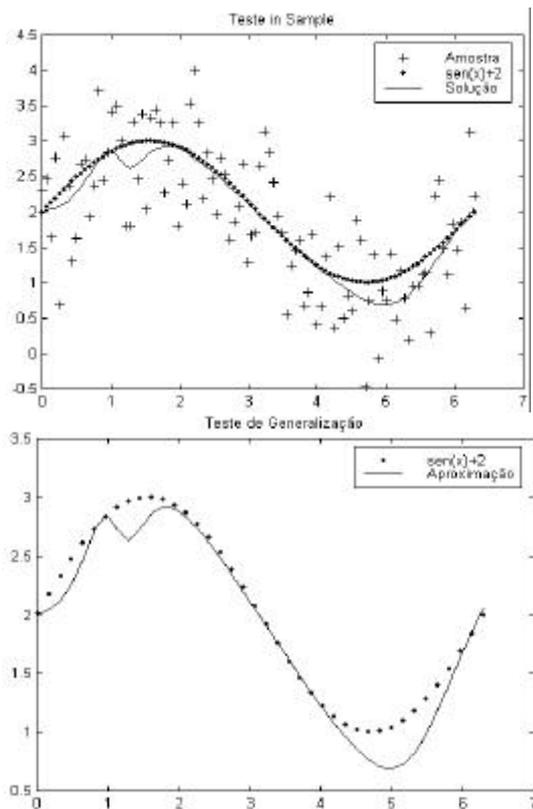


Figura 3d - Nível de ruído: 0.7

Utilizando os mesmos dados do experimento 4, tabela 5.1 inicializou-se o algoritmo, deliberadamente com uma condição inicial distante da solução esperada. Após 355 iterações obteve-se o Mape de 17.0na etapa de treinamento e 3.20 na etapa de generalização. A evolução da convergência neste caso é apresentada na figura 4.

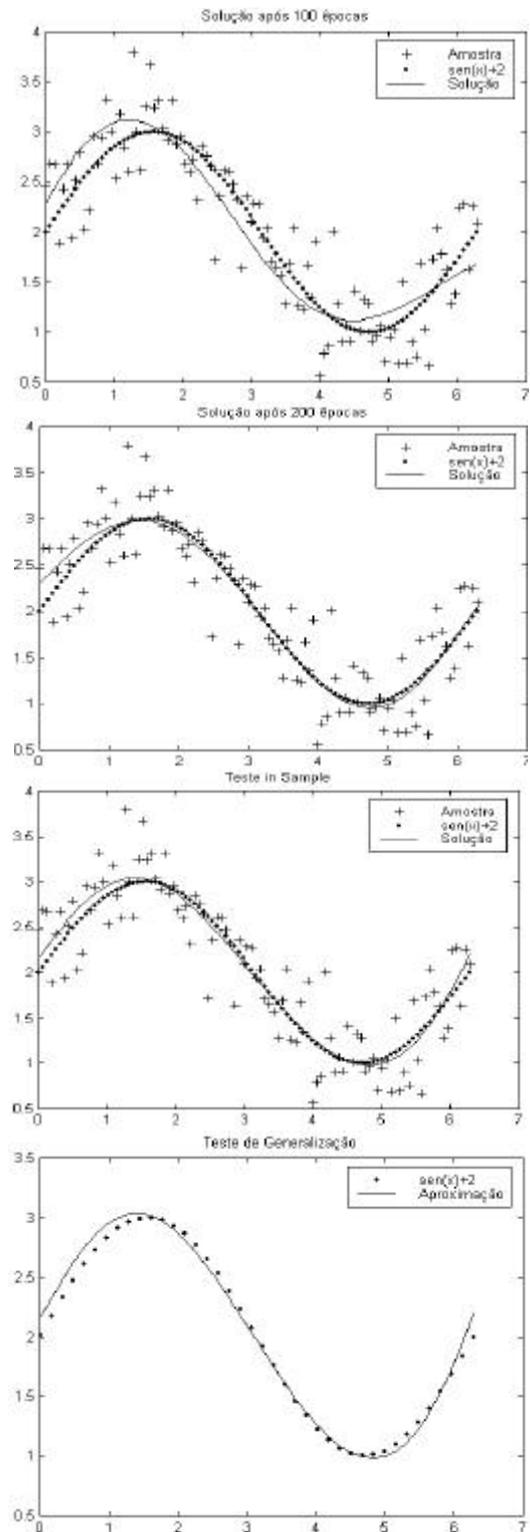
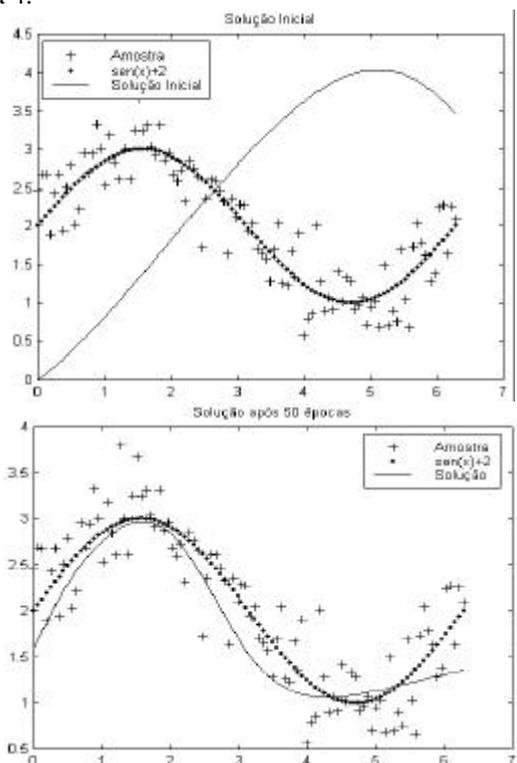


Figura 4 Evolução da convergência com inicialização sem uso da heurística

6. Uma aplicação em um problema de dados faltantes

Nesta seção apresenta-se uma aplicação da metodologia proposta em um problema de dados faltantes. Foram usados aqui dados reais de uma série de carga elétrica. A série considerada refere-se aos

dados de carga minuto a minuto para o dia 1 de Julho de 1999 da concessionária Light. Neste tipo de série o problema de dados faltantes é bastante comum.

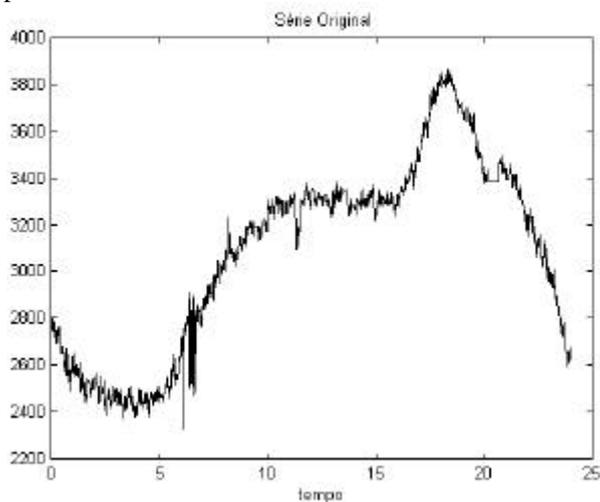


Figura 5 - Dados de carga minuto a minuto para o dia 1 de Julho de 1999

O experimento consiste em retirar da série original uma quantidade de pontos. A série é então recomposta pelo método proposto. Desta forma simula-se o problema de dados faltantes que é muito comum em várias aplicações e em particular no problema de previsão de carga. Foram retirados da série, aleatoriamente de modo uniforme, 5; 10; 20; 30 e 40 % de pontos em cada um dos 5 experimentos realizados. A tabela 6.1 resume o resultado encontrado.

Exp.	pontos retirados	Num. de épocas	Mape treinamento	Mape generalização
1	5%	75	0.99	1.21
2	10%	333	1.01	1.15
3	20%	106	1.06	1.116
4	30%	379	1.00	1.119
5	40%	105	1.07	1.117

Tabela 6.1 Simulação de dados faltantes com dados reais de carga elétrica

Em seguida fez-se uma comparação deste método com uma interpolação através de spline cúbico suavizado.

Nota-se (ver figura 7) que o spline não apresenta bons resultados quando se tem um número considerável de valores faltantes consecutivos. Este resultado era esperado uma vez que o spline fornece soluções com curvatura acentuada distanciando-se do padrão esperado de curva de carga diária.

Utilizou-se neste experimento comparativo o dia 1 de Julho omitindo-se 150 valores consecutivos. Após cerca de 400 épocas chegou-se ao resultado apresentado na figura 7. O MAPE obtido pelo método proposto foi de 1.02 na fase de treino e 2.19 no teste out of sample. Já a interpolação spline resultou em um MAPE de 9.31.

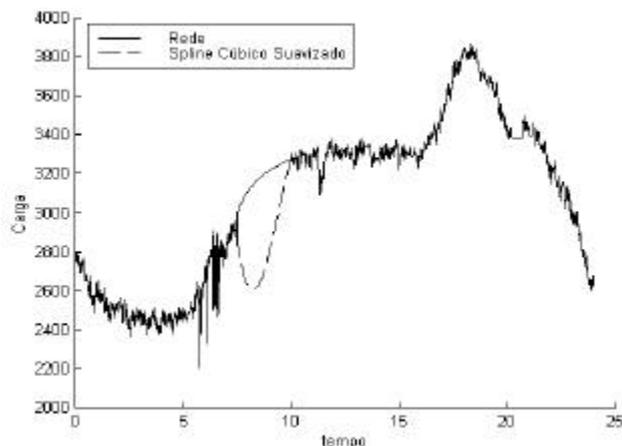


Figura 7 Comparação do método proposto com uma interpolação através de Spline Cúbico Suavizado

7. Conclusões

Neste artigo foi proposta uma arquitetura conexista que possui uma interessante capacidade interpolativa. Sua principal característica é a capacidade intrínseca de gerar soluções regulares. Este tipo de arquitetura abre também a possibilidade de interpretabilidade dos resultados uma vez que a localização das funções de pertinência pode indicar uma mudança no modelo. Mudanças na estrutura da função geradora dos dados devem se refletir em mudanças de posicionamento e de nível das funções de pertinência.

Os resultados simulados em ambientes ruidosos foram particularmente bons. Embora em alguns experimentos uma quantidade razoável de ruído tenha sido adicionada, o método mostrou-se robusto produzindo ótimos resultados em especial na fase de generalização.

Do ponto de vista de aplicações reais o experimento para dados faltantes apresentou resultados muito animadores, mais uma vez demonstrando a capacidade do método de interpolar dados produzindo soluções regulares e consistentes.

Referências:

- [1] Haykin S. "Neural Networks – A Comprehensive Foundation", Prentice Hall, second edition, 1999.
- [2] Pedroza L.C e Pedreira C.E. "Multilayer Neural Networks and Function Reconstruction by Using a priori Knowledge" International Journal of Neural Systems, Volume 9, number 3, pp 251-256, 1999.
- [3] Bartle, R.G. Elements of integration. Wiley. New York, 1966.
- [4] De Barra, G. (1974) Introduction to measure theory, New York, Van Nostrand Reinhold .