

PREDIÇÃO NEURAL DE INTERNAÇÕES HOSPITALARES POR DOENÇA RESPIRATÓRIA

LEONARDO LIMA DA SILVA MAROTTA*, JOSÉ MANOEL DE SEIXAS†, LUIZ PEREIRA CALÔBA‡

**Laboratório de Processamento de Sinais
COPPE/Poli-UFRJ
Rio de Janeiro, RJ, Brasil*

Emails: marotta@lps.ufrj.br, seixas@lps.ufrj.br, caloba@lps.ufrj.br

Abstract— This article approaches the application of artificial neural networks in the prediction of the number of infantile hospital internments decurrent of respiratory illnesses in the city of Paris, France. Specifically the climatic conditions and the atmospheric pollution impact directly on the respiratory illness, in such a way had been used different series as entered of the similar neural network to predict the number of patients to be interned. The database used was removed of a study carried through in Paris, contends information throughout four winters, when the amount of illnesses increases. The gotten results show a good approach, what it would allow to alert the hospitals for eventual peaks of internments and to assist the management of information in collective health.

Keywords— Neural Networks, Temporal Series, Respiratory Illnesses, Hospital Internments.

Resumo— Este artigo aborda a aplicação de redes neurais artificiais na predição do número de internações hospitalares infantis decorrente de doenças respiratórias na cidade de Paris, França. Especificamente as condições climáticas e a poluição atmosférica impactam diretamente sobre a doença respiratória, desta forma foram utilizadas diferentes séries como entrada da rede neural afim de prever o número de pacientes a serem internados. A base de dados utilizada foi retirada de um estudo realizado em Paris, contendo informações ao longo de quatro invernos, quando a quantidade de doenças aumenta. Os resultados obtidos mostram uma boa aproximação, o que permitiria alertar os hospitais para eventuais picos de internações e auxiliar a gerência de informação em saúde coletiva.

Keywords— Redes Neurais, Séries Temporais, Doenças Respiratórias, Internações Hospitalares.

1 Introdução

A predição de séries temporais utilizando redes neurais vem se tornando um forte aliado em diversas áreas da ciência. Uma aplicação muito interessante e potencialmente eficaz está na área médica, tanto na previsão de picos de internação, como também no apoio ao diagnóstico clínico em geral.

No caso deste trabalho, foram desenvolvidos modelos preditivos em uma série de internações hospitalares causadas por broquiolite infantil em diversos hospitais de Paris. Um dos principais responsáveis pelo agravamento de doenças respiratórias, especialmente em crianças, é a poluição atmosférica junto com as condições climáticas desfavoráveis ZANOBETTI (2000).

Diversos episódios têm sido registrados em grandes centros urbanos em virtude das altas concentrações de poluentes na atmosfera, com sérias conseqüências para a população. A preocupação causada pelos efeitos decorrentes da degradação do ar, que é um problema de saúde pública, tem sido tema de discussão entre líderes de diversos países, levando-os a firmarem acordos para controle e redução da emissão de gases causadores de efeito estufa. O conhecimento científico é de grande utilidade no auxílio à gestão da informação relativa à saúde coletiva, subsidiando a elaboração de medidas voltadas à prevenção e à redução da poluição do ar.

Esse trabalho tem como objetivo propor

a utilização das redes neurais artificiais como metodologia alternativa na predição de internações hospitalares no período de inverno, que é quando os problemas respiratórios se agravam. Os modelos desenvolvidos em WILLEMS (2005), no qual foi utilizada a metodologia dos modelos aditivos generalizados (GAM), e o estudo de NASCIMENTO (2006), em que foram utilizadas redes neurais artificiais, visando identificar os fatores relevantes no número de internações hospitalares infantis por doença respiratória na cidade de Paris, serviram como base para o presente trabalho. Estas doenças respiratórias são causadas freqüentemente pelo vírus syncethial respiratory virus (RSV), e normalmente, o contato com este vírus causa um resfriado. Em crianças e, especialmente no início do inverno, o vírus pode ser responsável por uma grave doença respiratória, conduzindo a um elevado número de consultas e internações hospitalares.

A importância desse estudo se dá na medida em que a modelagem neural pode ser utilizada como um importante instrumento de auxílio à gestão da informação relativa à saúde coletiva, subsidiando análises e possibilitando a adoção de medidas, tendo como base o conhecimento científico.

2 Os Dados

Para esse estudo, utilizou-se uma base de dados comum aos trabalhos já citados, obtida através do ERBUS (Epidémiologie et Recueil des Bronchiolites en Urgence pour Surveillance), referente a 43 hospitais localizados em Paris, no período compreendido entre 1997 e 2000, além de dados meteorológicos e de poluição atmosférica medidos pelas estações de AIRPARIF, nos anos de 1997 a 2001. Para evitar problemas com dados faltantes, foram considerados apenas 34 hospitais, cujos dados encontravam-se completos. Desta forma, a base de dados passou a contar com 419 observações.

A base de dados é composta por 17 variáveis explicativas, sendo 12 variáveis climáticas (temperatura mínima do dia, temperatura máxima do dia, temperatura média do dia, temperatura ao cair da tarde, precipitação diária, duração diária das chuvas, umidade mínima do dia, umidade máxima do dia, umidade média do dia, duração diária do período de sol, força do vento e pressão ao nível do mar) e 5 de poluição atmosférica (emissão de SO₂ pelas fábricas, emissão de CO pelos veículos, emissão de NO₂ pelos veículos, emissão média de O₃ por 8h e material particulado com diâmetro inferior a 10 μ m - PM₁₀). A série temporal do número de internações hospitalares representa o alvo em alguns dos modelos desenvolvidos neste trabalho e também entrada com atraso em alguns outros.

2.1 Tratamento dos Dados

A característica sazonal das variáveis deve ser levada em consideração, pois os dados devem ser dessazonalizados a fim de revelar o resíduo da série a ser modelada, permitindo assim a detecção das características escondidas NELSON (1999).

O projeto das redes neurais artificiais, no presente trabalho, consistiu em uma etapa inicial de normalização. As redes neurais foram utilizadas para prever, com um dia de antecedência, a série residual após o tratamento da série bruta de dados. Levando-se em consideração a característica sazonal das variáveis, avaliou-se a necessidade do pré-processamento dos dados. Esta avaliação foi feita após a implementação de diversas redes neurais, tendo como critério o maior poder de generalização do modelo.

A técnica utilizada para a previsão do número de internações nos hospitais concentra-se no estudo de séries temporais. O tratamento matemático de séries temporais tem o intuito de retirar todas as suas componentes determinísticas, tentando simplificar a análise estatística subsequente, e indicar a dependência ou não das mesmas, tanto no tempo quanto para diferentes séries.

As técnicas utilizadas foram:

- Normalização;
- Remoção da tendência;
- Remoção dos componentes senoidais;

2.2 Série original

A Figura 1 representa o número de internações hospitalares, onde são observados quatro períodos, cada um deles representando um inverno.

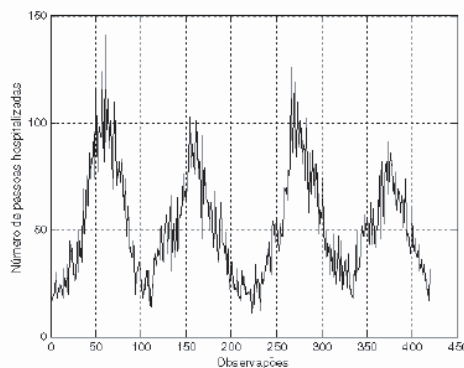


Figura 1: Série de Internações Hospitalares

2.3 Normalização

A normalização tem como finalidade adaptar os dados de entrada à faixa dinâmica das funções de ativação da rede neural, tornando o processo numericamente robusto e de fácil modelagem. O processo foi feito dividindo-se cada observação da série de uma dada variável pela sua média mais dois desvios-padrão.:

$$S_{norm}(n) = \frac{S(n)}{media[S(n)] + 2 \cdot desvioPadrao[S(n)]}$$

A Figura 2 mostra a série normalizada.

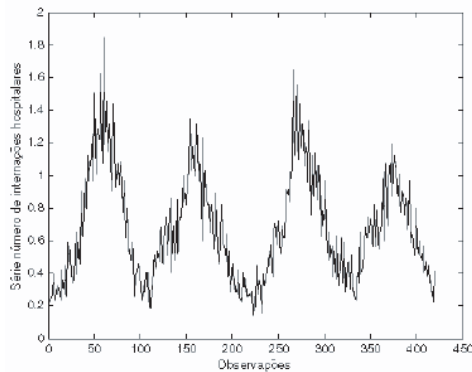


Figura 2: Série de Internações Hospitalares Normalizada

2.4 Remoção da tendência

A tendência de cada série pode ser estimada pelo método dos mínimos quadrados, através do seguinte cálculo:

$$S_2(n) = S_1(n) - (a + b \cdot n)$$

onde $S_1(n)$ é a série original normalizada, $S_2(n)$ é a série obtida após a eliminação da tendência, a é o intercepto e b é o coeficiente angular da reta de regressão. Neste caso, foi utilizada uma rede neural simples com apenas um neurônio e função de ativação linear, cujo par entrada-saída era $(t, S_1(n))$, minimizando-se o erro médio quadrático.

2.5 Remoção dos componentes senoidais

A remoção dos componentes senoidais de uma série é feita mediante a aplicação da Transformada de Fourier, e posterior observação das frequências que apresentam um pico significativo de amplitude relativamente às demais frequências da série KRISEIDA (1998). Desta forma, identificados os coeficientes dos senos e co-senos que compõem a informação da frequência, procede-se a retirada do pico correspondente, eliminando-se assim a oscilação naquela frequência específica.

Após a retirada dos componentes senoidais, tem-se a série S_3 :

$$S_3(n) = S_2(n) - (A \cdot \cos(2\pi w_o n) + B \cdot \sin(2\pi w_o n))$$

onde S_2 é a série original normalizada sem tendência linear, S_3 é a série obtida após a eliminação dos componentes senoidais, A e B são as amplitudes e w_o representa a frequência do componente senoidal identificado. Neste caso, foram retirados os ciclos, um a um, até que o valor RMS da série residual chegasse a um mínimo de 5% da série original. Também foi utilizada uma rede neural com um neurônio, tendo este a função de ativação linear com o fim de otimização no domínio do tempo para retirada dos ciclos, evitando assim problemas ocasionados por *aliasing*.

2.6 Séries Finais

Os mesmos procedimentos citados acima foram utilizados em todas as 17 séries de condições climáticas e de poluição atmosférica que foram usadas nos modelos, e também na série alvo. Apenas as informações dos três primeiros ciclos das séries foram utilizadas no desenvolvimento dos modelos e os resultados correspondentes foram retirados do quarto ciclo.

A Figura 3 mostra o resíduo da série de internações hospitalares, após o procedimento de pré-processamento.

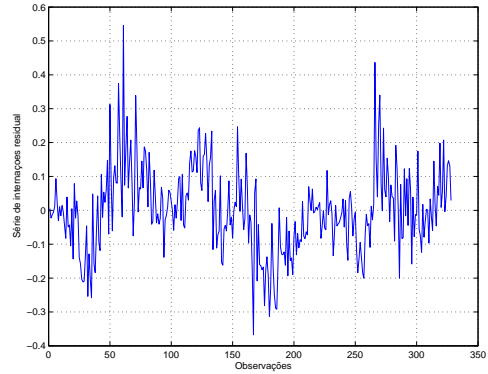


Figura 3: Série residual das Internações Hospitalares

3 Conjuntos de treino, teste e validação

Os modelos que visam a generalização da rede são restritos a um conjunto de treinamento, sendo a avaliação obtida através do teste de desempenho no conjunto restante de dados (conjunto de teste). O conjunto de treinamento deve conter um número suficiente de observações, para que possa constituir uma amostra representativa dos dados a serem analisados. Este conjunto é usado tanto para o desenvolvimento, quanto para o ajuste dos parâmetros da rede neural. Temos que contar também com um conjunto de validação que ajuda a evitar o *overtraining*.

O conjunto de teste não participa da fase treinamento da rede neural e, por ser um conjunto desconhecido, as informações nele contidas sofrem o mesmo tipo de tratamento dado àquelas que compõem o conjunto de treinamento. O conjunto de teste é usado apenas, para avaliar o poder de generalização da rede neural, ou seja, para avaliar se as informações pertencentes ao conjunto de treino são suficientes para uma análise mais geral. A generalização está associada à capacidade de predição de novos casos, indicando se o modelo neural conseguiu extrair as características principais da informação.

Do total de 419 observações, o conjunto de treinamento contém 328 amostras, correspondentes aos três primeiros invernos, e o conjunto de teste corresponde a 91 observações referentes ao último inverno observado. O conjunto de validação utilizado foi o de teste.

4 Treinamento da Rede Neural

A implementação foi feita através de redes neurais completamente conectadas, modelo MLP (Multi-layer Perceptron), com uma camada intermediária e sem realimentação. Os neurônios na camada escondida são do tipo tangente hiperbólica. Na camada de saída, foi utilizado um neurônio com função de ativação linear. O treinamento foi feito no modo batelada através do algoritmo Resilient

Backpropagation, que tem a vantagem de convergir mais rapidamente.

Como medida de desempenho no processo de treinamento da rede neural, adotou-se o MSE (erro médio quadrático), sendo calculado como:

$$mse = \frac{\sum (y - y')^2}{n}$$

onde y é o valor observado, y' é a saída da rede e n é o tamanho do vetor de entrada da rede neural.

A função objetivo teve como finalidade a minimização do MSE no conjunto de treinamento.

O erro percentual absoluto médio (MAPE), foi calculado em paralelo ao MSE e foi adotado como critério de parada de treinamento e também para avaliar o poder de generalização do modelo. O MAPE é definido como:

$$mape = \frac{\sum \left| \frac{y - y'}{y} \right|}{n}$$

5 Resultados

Antes da predição com redes neurais a série foi reconstruída utilizando-se apenas a média, a tendência e os ciclos senoidais extraídos das séries de treino. Isto foi feito de forma a obter um parâmetro de comparação com o preditor neural. Após a reconstrução o erro MAPE encontrado foi de 34,93%.

Todos os modelos de redes neurais utilizados tinham como objetivo prever, com um dia de antecedência, o número de internações hospitalares. O primeiro modelo utilizou como entrada apenas as séries auxiliares de condições climáticas e poluição atmosférica. O segundo teve como entrada as séries auxiliares e a série alvo no dia D para prever o dia D+1. O terceiro recebeu as séries auxiliares e a série alvo no dia D e D-1. O quarto e último modelo recebeu apenas a série alvo no dia D e D-1.

Em todos os modelos foram treinadas redes neurais com o número dos neurônios da camada intermediária variando de 1 a 15. No primeiro modelo o menor erro MAPE encontrado foi de 28,11%. Este utilizou oito neurônios na camada escondida. A Figura 4 mostra a série de internações hospitalares original em linha contínua, e a reconstruída em tracejado obtida por este modelo.

O modelo número 2, com apenas um neurônio, apresentou um erro MAPE de 11,12%. A Figura 5 mostra a série de internações hospitalares original em linha contínua, e a reconstruída por este modelo em tracejado.

O terceiro modelo obteve um erro MAPE mínimo de 13,45% com 13 neurônios na camada escondida. A Figura 6 mostra a série de internações hospitalares original em linha contínua, e a reconstruída por este modelo em tracejado.

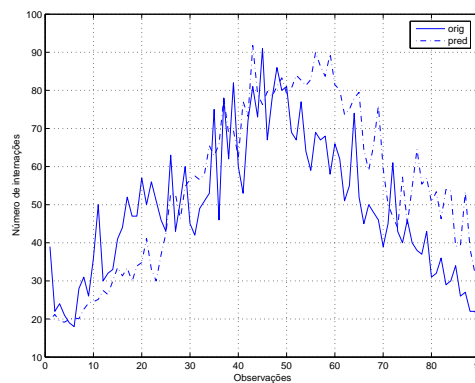


Figura 4: Série comparativa das Internações Hospitalares do Modelo 1

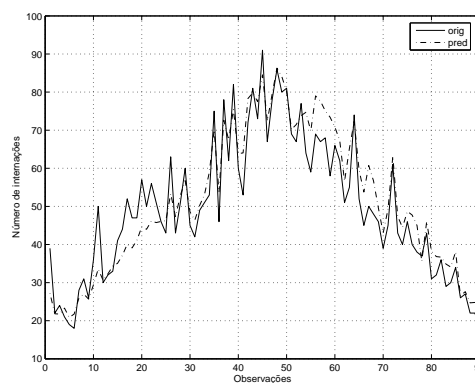


Figura 5: Série comparativa das Internações Hospitalares do Modelo 2

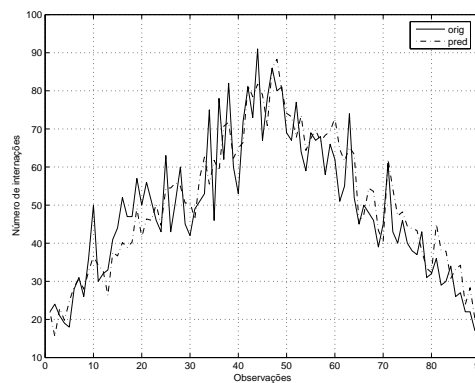


Figura 6: Série comparativa das Internações Hospitalares do Modelo 3

O quarto modelo com quatro neurônios na camada escondida apresentou um erro MAPE de 10,29%. A Figura 7 mostra a série de internações hospitalares original em linha contínua, e a reconstruída por este modelo em tracejado.

6 Conclusão

Sem dúvida alguma o estudo de predição de séries temporais não é uma tarefa fácil, pois não basta colocar como entrada em uma rede neural os dados brutos, é necessário retirar todas as car-

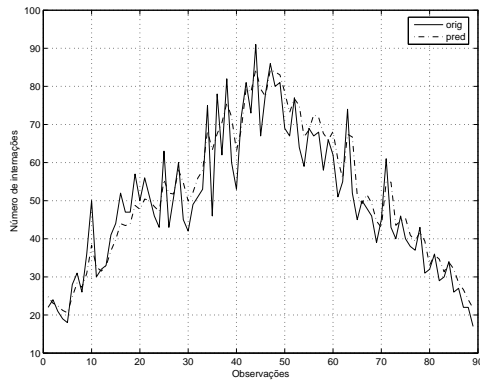


Figura 7: Série comparativa das Internações Hospitalares do Modelo 4

acterísticas possíveis normalmente representadas por funções determinísticas das séries. O objetivo deste trabalho é formar uma série residual que contenha apenas aquilo que não é determinado facilmente por métodos clássicos, e que normalmente só podem ser representados por modelagens não-lineares. As redes neurais têm um papel fundamental após esse pré-processamento.

Os resultados obtidos são promissores mas ainda podem melhorar com a implementação de outros modelos mais complexos. A inclusão de médias móveis de janelas variadas seria uma opção a mais no pré-processamento. Redes neurais realimentadas, tomados os devidos cuidados com a instabilidade, podem surtir efeitos positivos. Um estudo de relevância pode mostrar quais variáveis impactam de fato na predição da série alvo.

Os modelos neuronais se mostraram bastante eficazes quando comparados a reconstrução utilizando-se apenas a média, a tendência e os ciclos senoidais das séries de treino.

Este trabalho desenvolveu quatro modelos: o primeiro deles apresentou um erro relativamente grande se comparado aos outros modelos, nota-se que este não possui uma eficiência boa na representação da magnitude das internações mas, acompanha suas variações; o segundo modelo, já com um erro bem menor, é falho na previsão antes do pico da série, no entanto, após este evento, a previsão acompanha muito bem tanto na magnitude quanto nas variações; o terceiro modelo apresentou um erro relativamente próximo ao do segundo modelo, o que indica que a entrada da série alvo atrasada é relevante no processo de predição, e da mesma forma acompanha bem a série após o pico; o quarto modelo apresentou o melhor dos resultados, mostrando que a série alvo por si só já contém informações importantes para a predição.

Ainda assim, falta uma melhor exploração nas séries auxiliares. O trabalho de NASCIMENTO (2006) mostra que o uso de médias móveis de 6 dias apresenta uma melhora nos modelos representados pelas séries auxiliares, isto seria uma alter-

nativa para aproveitar melhor estas informações.

Referências

- KRISEIDA (1998). *Modelos Neurais em Engenharia de Transporte*, Tese de Doutorado - COPPE/UFRJ.
- NASCIMENTO (2006). *Redes Neurais Artificiais: Uma Aplicação no Estudo da Poluição Atmosférica e seus Efeitos Adversos à Saúde*, Dissertação de Mestrado - COPPE/UFRJ.
- NELSON (1999). *Times series forecasting using neural networks: should the data be deseasonalized first?*, J. Forecast.
- WILLEMS (2005). *Longitudinal analysis of short term bronchiolitis air pollution association using semi parametric models*, Paris: Conseil Scientifique de l'ANTADIR.
- ZANOBETTI (2000). *Generalized additive distributed lag models: quantifying mortality displacement*, Biostatistics.