

Trabalho 1: Regressão Linear com Múltiplas Variáveis

Reconhecimento de Padrões – 2021/1

Nesse trabalho sobre regressão linear múltipla, utilizaremos a base de dados hospital do matlab (disponível na pasta Bases de Dados como hospital.xls), com 4 variáveis de predição (sexo, idade, peso, fumante/não fumante), para prever o valor da pressão máxima de um paciente. Essa base de dados contém 100 observações. Para tanto é desejável que:

- 1) Seja empregada validação cruzada com 5 pastas;
- 2) Sejam determinados os valores dos coeficientes de Pearson e o valores dos erros médio quadráticos para cada pasta (5 ao total);
- 3) Utilizando os valores obtidos em cada pasta determine o valor médio do coeficiente de Pearson e do erro médio quadrático.

Sugestão para elaboração de relatório:

1. Introdução
Expor o problema a ser resolvido
2. Fundamentação Teórica
Mostrar o regressor de múltiplas variáveis utilizando a pseudo-inversa
3. Metodologia
Descrever como o experimento será feito utilizando as 5 pastas: quantos registros cada pasta terá, como os resultados serão avaliados, etc
4. Resultados
Apresentar os 5 valores do coeficiente de Pearson e do erro médio quadrático, bem como o valor médio dos mesmos.
5. Conclusões
Avaliar os resultados obtidos.

Fazer o upload no sistema de EAD até 05/04/2021

Apêndice: Método da PseudoInversa

Suponha que em um problema estejam disponíveis p variáveis para serem utilizadas na predição do valor de uma variável y de saída. Por exemplo, na determinação da pressão sanguínea de um paciente, podem ser utilizadas as variáveis; sexo, idade, peso e a condição de ser fumante ou não. Na regressão múltipla, a variável a ser predita é expressa através de uma equação linear:

$$y_k = x_{k1}w_1 + x_{k2}w_2 + \dots + x_{kp}w_p \quad (1)$$

$$y_k = \mathbf{w}^T \mathbf{x}_k$$

Em que:

\mathbf{x}_k – *padrão* $_k$

\mathbf{w} – vetor de pesos

y_k – *valor a ser predito*

Para N observações, temos o seguinte conjunto com N equações:

$$\begin{aligned} y_1 &= x_{11}w_1 + x_{12}w_2 + \dots + x_{1p}w_p + k \\ y_2 &= x_{21}w_1 + x_{22}w_2 + \dots + x_{2p}w_p + k \\ &\dots\dots\dots \\ &\dots\dots\dots \\ y_N &= x_{N1}w_1 + x_{N2}w_2 + \dots + x_{Np}w_p + k \end{aligned} \quad (2)$$

O Sistema linear de (2) é resumido a seguir:

$$\mathbf{y} = \mathbf{X}\mathbf{w} \quad (3)$$

E que:

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} & 1 \\ x_{21} & x_{22} & \dots & x_{2p} & 1 \\ \dots & \dots & \dots & \dots & \dots \\ x_{N1} & x_{N2} & \dots & x_{Np} & 1 \end{pmatrix} \quad (4)$$

$$\mathbf{w} = \begin{pmatrix} w_1 \\ w_2 \\ \dots \\ w_p \\ k \end{pmatrix} \quad (5)$$

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_N \end{pmatrix} \quad (6)$$

Se o número de variáveis p for igual ao número de observações N , o sistema linear mostrado em (3) pode ser solucionado através da matriz inversa de \mathbf{X} :

$$\mathbf{w} = \mathbf{X}^{-1} \cdot \mathbf{y} \quad (7)$$

Normalmente, a dimensão de p é muito menor do que a dimensão de N , de tal forma que o sistema linear mostrado em (3) não pode ser resolvido por (7). Conforme demonstrado em sala de aula, para minimizar o erro médio quadrático mostrado em (8), a solução é dada pela equação (9):

$$E = \sum_{i=1}^N \|y_i - \mathbf{X}_i \cdot \mathbf{w}\|^2 \quad (8)$$

$$\mathbf{w} = \mathbf{X}^+ \cdot \mathbf{y} \quad (9)$$

Em que:

$$X^+ - \text{matriz pseudo inversa (PI)} = (X^T \cdot X)^{-1} X^T \quad (10)$$

$$X_i = (x_{i1} \ x_{i2} \ \dots \ x_{ip} \ 1) \quad (11)$$

Com base na matriz PI o seguinte método de previsão da variável y pode ser proposto:

1. Defina um conjunto de treinamento com N observações;
2. Para esse conjunto determine a matriz PI através de (10);
3. Determinar os pesos através de (9)
4. Utilizar a matriz PI para a predição através da equação $y = \mathbf{w}^T \mathbf{x}$