

Rogério Custodio, João Carlos de Andrade e Fábio Augusto

Instituto de Química - Universidade Estadual de Campinas - UNICAMP - 13083 - 970 - Campinas - SP

Recebido em 16/1/96; aceito em 11/10/96

**CURVE FITTING OF MATHEMATICAL FUNCTIONS TO EXPERIMENTAL DATA. The least square method is analyzed. The basic aspects of the method are discussed. Emphasis is given in procedures that allow a simple memorization of the basic equations associated with the linear and non linear least square method, polynomial regression and multilinear method.**

**Keywords: least square method; linear regression; non linear regression.**

**INTRODUÇÃO**

Qualquer conjunto de pontos em um espaço multidimensional apresentando uma tendência regular pode ser representado por uma função matemática. Frequentemente esta função é escolhida através de um processo de ajuste conhecido como método de mínimos quadrados. Para o caso mais simples de um conjunto de pontos em um espaço bidimensional, cada ponto será representado por coordenadas  $x$  e  $y$ . A função matemática a ser construída pode representar a tendência de  $y$  como uma função de  $x$ ,  $y=f(x)$ , ou o inverso, com  $x$  sendo representado como uma função de  $y$ ,  $x=f(y)$ . Em outras palavras, pode-se ter  $y$  como uma variável dependente de  $x$  ou  $x$  como uma variável de  $y$ .

O emprego desta terminologia, matematicamente correta, pode ser pedagogicamente inadequada, pois pode fazer com que um aluno deixe de perceber a simplicidade intrínseca do método, recorrendo posteriormente a argumentos tais como: a) a minha calculadora não é programável, b) minha calculadora não faz regressão, c) não me lembro das fórmulas, etc. Estes argumentos não se justificam, uma vez que não é necessário uma bagagem matemática sofisticada para utilizar-se o método dos mínimos quadrados.

Neste artigo pretende-se abordar de modo simples alguns tópicos mais relevantes relacionados com a técnica de regressão linear, demonstrando que as equações necessárias para o emprego dos diferentes tipos de regressão podem ser obtidas a partir de conhecimentos básicos de matemática. Por outro lado, alguns exemplos mostram que o seu uso é uma pura questão de bom senso. Leituras especializadas<sup>1-4</sup> são recomendadas ao leitor interessado em um aprofundamento maior no assunto.

**MÉTODO DOS MÍNIMOS QUADRADOS**

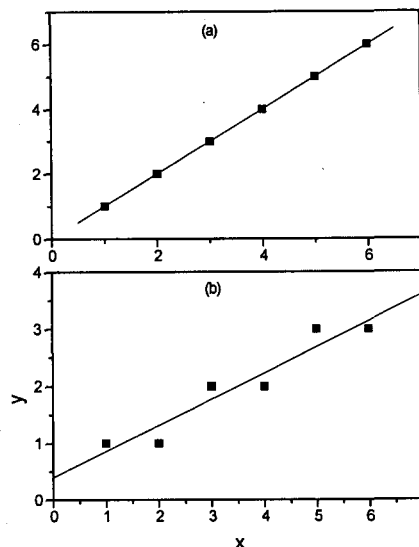
Considere o caso em que um determinado experimento produziu um conjunto de  $N$  pontos caracterizados pelas coordenadas  $\{(y_i, x_i)\}$ , onde  $i=1,2,\dots,N$ . Tomando-se a variável  $y$  como sendo a variável dependente e  $x$  como variável independente, caracterizamos o conjunto de pontos experimentais como:  $\{(y_i^{exp}, x_i)\}$ . Matematicamente pode-se determinar através de uma função matemática qualquer os valores de  $y_i$  para cada  $x_i$ . Aos valores de  $y_i$  estimados matematicamente denominaremos por  $y_i^{est}$ . O método de mínimos quadrados sugere que a função  $y^{est}$  deverá ser determinada de tal maneira que a diferença em relação ao conjunto de dados experimentais,  $y^{exp}$ , seja mínima. Este desvio pode ser representado pela equação:

$$Q = \sum_{i=1}^N (y_i^{exp} - y_i^{est})^2 \tag{1}$$

Esta equação deixa claro a natureza do nome do método. Pode-se questionar por que é necessário utilizar-se o quadrado dos desvios entre  $y_i^{exp}$  e  $y_i^{est}$ . A resposta é bastante simples. Suponha que o expoente quadrático da Eq.1 fosse eliminado e que o cálculo da diferença entre um conjunto de valores  $y_i^{exp}$  e  $y_i^{est}$  fosse realizado pela equação

$$Q' = \sum_{i=1}^N (y_i^{exp} - y_i^{est}) \tag{2}$$

O fato de se obter um valor de  $Q'$  igual a zero não implicaria em garantir que todas as diferenças calculadas,  $(y_i^{exp} - y_i^{est})$ , fossem iguais a zero. A figura 1 mostra dois gráficos contendo dados experimentais e uma função matemática (no caso, uma reta) estimada. Em ambos os casos o valor de  $Q'$  é igual a zero. Porém, verifica-se que, enquanto na figura 1.a existe uma concordância absoluta entre os valores experimentais e a equação estimada, na figura 1.b percebe-se um desvio significativo entre os valores experimentais e os valores calculados. Entretanto, neste segundo caso  $Q'$  também será igual a zero, uma vez que ocorre um cancelamento na grandeza dos desvios. Um conjunto de valores  $\Delta_i = (y_i^{exp} - y_i^{est})$  apresentam valores maiores do que zero e um outro conjunto apresenta valores menores do que zero. Somando-se todos os desvios obtém-se  $Q'=0$ .



**Figura 1.** Dois exemplos de regressão linear em conjuntos de dados experimentais com características diferentes a) completamente linear e b) distribuídos com padrão regular em torno de uma reta.

Para evitar situações como esta e ter-se uma medida da dispersão dos pontos ao redor da função matemática escolhida, utiliza-se o quadrado dos desvios. Com isto, não haverá o cancelamento dos desvios positivos e negativos e quando a função apresentar  $Q=0$ , isto indicará que a função estará passando exatamente sobre todos os pontos experimentais.

A figura 1.b mostra ainda um outro ponto a ser discutido. Embora se tenha uma descrição correta da tendência dos pontos, observa-se ainda um desvio considerável entre valores experimentais e valores estimados. Pode-se ainda tentar ajustar outras funções que passem exatamente sobre o conjunto de dados experimentais. Por exemplo, a representação matemática dos dados experimentais por uma função que apresente  $Q=0$  (Fig.2). Pode-se verificar que existe uma diferença significativa entre o comportamento matemático da função apresentada na figura 1.a daquela mostrada pela figura 2. Resta a pergunta: qual a função que representa o conjunto de dados experimentais de maneira mais correta? A resposta não é óbvia. O fato da função apresentada na figura 2 possuir  $Q=0$  sugere que esta seja a melhor escolha em termos matemáticos. Porém, a escolha final deve levar também em conta a natureza física do conjunto de dados experimentais. Observando-se a natureza desse sistema pode-se concluir que, embora os desvios apresentados sejam maiores na figura 1.b do que na figura 2, fisicamente não existiria qualquer explicação razoável para admitir-se a escolha da função da figura 2 se o conjunto de dados estiver representando uma lei física linear, como por exemplo, a lei de Beer.

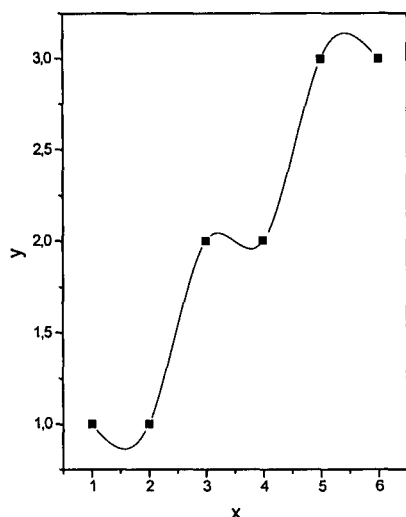


Figura 2. Ajuste não-linear a um conjunto arbitrário de dados experimentais.

A escolha de uma função, caso não se tenha nenhum modelo físico para representar  $y^{est}$ , também não é óbvia. Frequentemente observa-se a tendência dos dados experimentais e utiliza-se equações matemáticas que se ajustem a tal tendência. A escolha da melhor função, como mencionado acima, deve ser feita considerando-se uma função que apresente o menor  $Q$  e, quando possível, que não apresente nenhuma inconsistência física. Normalmente, esta escolha é efetuada de maneira criteriosa através da aplicação de testes estatísticos que avaliam não só a confiabilidade do ajuste matemático, mas também a confiabilidade dos dados estudados através de testes de rejeição de alguns pontos. Este é um tema complexo, que está associado ao erro no conjunto de dados disponível e envolve uma abordagem que não corresponde ao objetivo proposto por este artigo. Entretanto, o leitor interessado poderá recorrer a vasta literatura disponível, como por exemplo as refs.[1-4].

Mas, como uma função pode ser ajustada a dados experimentais? Pode-se analisar o problema de maneira genérica e

posteriormente particularizá-lo para casos específicos. Quando se escolhe uma função, estabelece-se uma relação entre  $y$  e  $x$  ou mais variáveis. Para ajustar-se uma função específica no espaço utilizam-se parâmetros que podem ser representados pelas letras  $a, b, c, \dots$ . Por exemplo, por um plano  $(x,y)$  passam infinitas retas representadas por  $y=ax+b$ . Porém, apenas uma reta específica será definida pela escolha apropriada de  $a$  e  $b$ .

No caso de um ajuste de mínimos quadrados, os parâmetros  $a, b, c, \dots$  devem ser ajustados de tal maneira que os valores de  $y^{est}$  aproximem-se ao máximo dos valores de  $y^{exp}$ , ou em outras palavras, que a diferença entre  $y^{exp}$  e  $y^{est}$  seja mínima. Na prática a busca de um mínimo envolve o uso de derivadas. Desta forma, deve-se derivar  $Q$  em relação a cada um dos parâmetros  $a, b, c, \dots$ , igualando cada uma dessas derivadas a zero. Isto é:

$$\frac{\partial Q}{\partial a} = \frac{\partial Q}{\partial b} = \frac{\partial Q}{\partial c} = \dots = 0 \quad (3)$$

Igualar as derivadas a zero é uma condição necessária para que se tenha o mínimo de uma função. A série de derivadas produzirá um conjunto de equações em função dos parâmetros  $a, b, c, \dots$  que, em alguns casos, possibilitará determinar os valores desses parâmetros em função do conjunto de pontos  $(x,y)$  experimentais.

Este procedimento pode ser ilustrado através da determinação dos parâmetros  $a$  e  $b$  usados para ajustar uma reta,  $y=ax+b$ , a um conjunto de  $N$  pontos experimentais  $(y_i, x_i)$ .

A Eq.1 para este caso específico deve ser escrita como:

$$Q = \sum_{i=1}^N [y_i^{exp} - y_i^{est}(x_i, a, b)]^2 \quad (4)$$

Nesta expressão são conhecidos os valores de  $y_i^{exp}$  e de  $x_i$  e pretende-se determinar os valores de  $a$  e  $b$ , de tal maneira que  $Q$  seja mínimo. Reescrevendo a Eq.4 em função dos dados experimentais e dos parâmetros  $a$  e  $b$  que se pretende determinar, tem-se que:

$$Q = \sum_{i=1}^N [y_i^{exp} - (a x_i + b)]^2 \quad (5)$$

Derivando  $Q$  em relação a  $a$  e  $b$  e considerando-se genericamente  $y_i^{exp} = y_i$ , obtêm-se as seguintes expressões:

$$\frac{\partial Q}{\partial a} = 0 = a \sum_{i=1}^N x_i^2 + b \sum_{i=1}^N x_i - \sum_{i=1}^N x_i y_i \quad (6)$$

$$\frac{\partial Q}{\partial b} = 0 = a \sum_{i=1}^N x_i + Nb - \sum_{i=1}^N y_i \quad (7)$$

Como pode ser visto, temos duas equações e podemos determinar o valor das duas incógnitas  $a$  e  $b$ . As Eqs.6 e 7 podem ser rearranjadas fornecendo expressões para  $a$  e  $b$  em função apenas de  $y$  e  $x$ , ou seja:

$$a = \frac{N \sum_{i=1}^N x_i y_i - \sum_{i=1}^N x_i \sum_{i=1}^N y_i}{N \sum_{i=1}^N x_i^2 - (\sum_{i=1}^N x_i)^2} \quad (8)$$

$$b = \frac{\sum_{i=1}^N y_i \sum_{i=1}^N x_i^2 - \sum_{i=1}^N x_i \sum_{i=1}^N x_i y_i}{N \sum_{i=1}^N x_i^2 - (\sum_{i=1}^N x_i)^2} \quad (9)$$

A equação da reta descrita pelos parâmetros  $a$  e  $b$  determinadas pelas Eq.8 e 9 representa a reta que mais se aproxima de todos os pontos experimentais disponíveis.

O caso mais comum que se encontra para ilustrar o uso do método dos mínimos quadrados é a construção de curvas de calibração<sup>3</sup>, em especial aquelas envolvendo a lei de Beer. Como exemplo, pode-se considerar a curva de calibração para determinação de manganês como  $MnO_4^-$  em 545nm, usando uma cela de 10 mm de caminho óptico:

$x_i (\mu g \cdot ml^{-1})$	$y_i$	$x_i^2$	$x_i y_i$
2,00	0,103	4,00	0,206
4,00	0,185	16,00	0,740
6,00	0,257	36,00	1,542
8,00	0,331	64,00	2,648
10,00	0,422	100,00	4,220
15,00	0,601	225,00	9,015
20,00	0,803	400,00	16,060
$\Sigma$	65,00	2,702	845,00
			34,431

onde  $x$  é a concentração em termos de  $MnO_4^-$  e  $y$  é a absorbância lida. Desta forma:

$$a = \frac{[7(34,431) - (2,702)(65,00)]}{[7(845,00) - (65,00)^2]} = 0,03869... = 0,0387$$

$$b = \frac{[(845,00)(2,702) - (65,00)(34,431)]}{[7(845,00) - (65,00)^2]} = 0,02673... = 0,0267$$

Então:

$$y = 0,0387 + 0,0267 x$$

Embora tenha-se exemplificado o desenvolvimento formal para o ajuste da equação de uma reta, também conhecido como *método de regressão linear*, o mesmo procedimento pode ser aplicado para equações matemáticas com maior número de parâmetros. A dificuldade em ajustar-se equações mais complexas está obviamente na manipulação de um maior número de equações. Neste sentido, ao invés da manipulação analítica apresentada, pode-se utilizar técnicas numéricas. O objetivo final de todo procedimento acima é minimizar  $Q$  através dos melhores parâmetros  $a$  e  $b$ . Isto pode ser feito através de uma busca sistemática dos valores desses parâmetros. Esta busca frequentemente pode ser feita de duas maneiras, ou empregando-se as expressões de derivadas de  $Q$  ou então através de métodos que não empregam derivadas. Os métodos desenvolvidos para minimizar uma função baseados no uso das derivadas necessitam de expressões das derivadas da função  $Q$ , tais como as apresentadas nas Eqs.6 e 7. Conhecendo-se a expressão para  $Q$ , por exemplo a Eq.5, e as expressões das derivadas dos parâmetros envolvidos introduz-se parâmetros  $a$ ,  $b$ ,  $c$ , etc arbitrários e calculam-se as derivadas de  $Q$  em relação a cada parâmetro. Com esses valores o método em uso geralmente estabelece em que direção os valores dos parâmetros deve mudar para que todas as derivadas tendam a zero. Esses métodos, conhecidos como métodos de gradiente, são extremamente poderosos e normalmente devem ser utilizados em processos de otimização. Entretanto, pode-se empregar métodos que utilizam

apenas a expressão de  $Q$ . O procedimento é análogo ao mencionado acima, com a diferença de que não calculam-se as derivadas, mas apenas o valor de  $Q$  para um conjunto de valores arbitrários de  $a$ ,  $b$ ,  $c$ , etc. Em geral estes métodos produzem variações nos parâmetros seguindo determinada metodologia e procuram determinar os valores desses parâmetros de tal forma que  $Q$  seja mínimo. Um exemplo clássico em otimizações deste tipo é o método simplex<sup>5,6</sup>.

## MÉTODOS COMPUTACIONAIS PARA REGRESSÃO LINEAR E NÃO-LINEAR

Deve parecer evidente que os ajustes, seja empregando métodos baseados em gradiente ou não, tendem a facilitar o ajuste de funções. Um outro ponto que se deve ter em mente é que quanto maior o número de pontos experimentais e de parâmetros a serem ajustados, maiores as dificuldades de manipular as equações necessárias para a solução dos mínimos quadrados. A necessidade de computadores ou calculadoras torna-se então clara, de modo que os aspectos computacionais envolvidos neste tipo de ajuste deve ser real. Em geral, a necessidade do uso de computadores mascara a simplicidade da idéia de ajustes dos mínimos quadrados e leva em geral os alunos (de graduação ou mesmo de pós-graduação) a acreditarem que somente com o uso de equipamentos sofisticados é que o problema pode ser resolvido. Esta é uma visão completamente equivocada como mostrado a seguir.

## FORMA MNEMÔNICA PARA REGRESSÃO LINEAR

Independentemente de se saber que é preciso derivar  $Q$  com relação aos parâmetros  $a$ ,  $b$ ,  $c$ , etc e posteriormente remanejar-se todas as equações resultantes de tal forma que se tenha expressões de  $a$ , de  $b$ , de  $c$ , etc em função apenas de  $x$  e  $y$ , ou ainda de se encontrar um programa que procure (através de técnicas de gradiente ou não) os valores ótimos de  $a$ ,  $b$ ,  $c$ , etc, pode-se derivar um conjunto de equações de maneira extremamente mais simples, que levará à solução do ajuste dos mínimos quadrados e que são extremamente mais fáceis de serem memorizadas<sup>7</sup>. Vamos ao exemplo da regressão linear.

O ajuste de uma equação linear por um conjunto de  $N$  pontos experimentais implica na determinação dos parâmetros  $a$  e  $b$  da equação:

$$y = a + bx \quad (10)$$

A Eq.10 mostra que é preciso determinar dois parâmetros e por enquanto tem-se apenas uma única equação. Seria necessário pelo menos mais uma segunda equação para que se pudesse determinar os parâmetros da reta. Como alternativa pode-se duplicar a Eq.10, resultando então duas equações e duas incógnitas. Porém, as duas equações seriam idênticas, o que impossibilitaria a determinação de  $a$  e  $b$ . Uma maneira de se diferenciar as duas equações seria multiplicar a segunda equação pela variável independente,  $x$ . Desta forma, as Eqs.11 e 12 resultantes poderiam ser escritas como

$$y = a + bx \quad (11)$$

$$yx = ax + bx^2 \quad (12)$$

Entretanto, deve-se lembrar que não se tem um único ponto experimental, mas em geral um número  $N$ , maior do que dois. Desta forma, é preciso identificar os pares de variáveis  $x$  e  $y$  por um subíndice  $i$ . Em outras palavras, cada ponto experimental deverá ser representado pelas equações:

$$y_i = a + bx_i \quad (13)$$

$$y_i x_i = ax_i + bx_i^2 \quad (14)$$

Uma vez que se tem uma única reta representando da melhor maneira possível todos os  $N$  dados experimentais, deve-se somar a contribuição de todos os pontos nas Eqs.13 e 14:

$$\sum_{i=1}^N y_i = Na + b \sum_{i=1}^N x_i \quad (15)$$

$$\sum_{i=1}^N y_i x_i = a \sum_{i=1}^N x_i + b \sum_{i=1}^N x_i^2 \quad (16)$$

A solução das duas equações com duas incógnitas apresentadas nas Eqs.15 e 16 levam a solução dos coeficientes de regressão linear. Rearranjando-se  $a$  e  $b$ , pode-se obter as equações apresentadas pelas Eqs.8 e 9. De outra forma, as Eqs.15 e 16 podem ser reescritas na forma matricial:

$$\begin{pmatrix} \sum_{i=1}^N y_i \\ \sum_{i=1}^N y_i x_i \end{pmatrix} = \begin{pmatrix} N & \sum_{i=1}^N x_i \\ \sum_{i=1}^N x_i & \sum_{i=1}^N x_i^2 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} \quad (17)$$

$$Y = XA \quad (18)$$

Um aspecto a ser observado é que a matriz  $X$  é uma matriz simétrica. Esta simetria possibilita o uso de eficientes algoritmos para a solução da Eq.18<sup>8</sup>. Uma possibilidade imediata ao alcance de qualquer aluno, seria o cálculo da matriz inversa de  $X$ ,  $X^{-1}$ . Assim, conhecendo-se  $X^{-1}$  pode-se multiplicar ambos os lados da Eq.17, obtendo-se os coeficientes  $A$  através da expressão:

$$X^{-1}Y = A \quad (19)$$

Considerando-se o exemplo anterior, tem-se que:

$$\begin{cases} Y = XA \\ \begin{bmatrix} 2,702 \\ 34,431 \end{bmatrix} = \begin{bmatrix} 7 & 65,00 \\ 65,00 & 845,00 \end{bmatrix} \begin{pmatrix} a \\ b \end{pmatrix} \end{cases}$$

ou:

$$\begin{cases} X^{-1}Y = A \\ \begin{bmatrix} 5,00 \cdot 10^{-1} & -3,85 \cdot 10^{-2} \\ -3,85 \cdot 10^{-2} & 4,14 \cdot 10^{-3} \end{bmatrix} \begin{bmatrix} 2,702 \\ 34,431 \end{bmatrix} = \begin{pmatrix} a \\ b \end{pmatrix} \end{cases}$$

## AJUSTES NÃO LINEARES EMPREGANDO REGRESSÃO LINEAR

O procedimento acima pode ser generalizado para ajustar não apenas funções lineares, mas também polinômios em geral. O mesmo recurso mnemônico pode ser empregado e encontra-se

convenientemente descrito em livros textos básicos<sup>7</sup>. Entretanto, vale a pena salientar que mesmo para polinômios de grau superior a 1 a Eq.19 apresenta exatamente a mesma estrutura básica a ser resolvida. A única diferença encontra-se nas matrizes  $A$  e  $Y$ , que para ajustar um polinômio de grau  $k$  apresentam-se sob a forma:

$$A = \begin{pmatrix} N & \sum_{i=1}^N x_i & \sum_{i=1}^N x_i^2 & \dots & \sum_{i=1}^N x_i^k \\ \sum_{i=1}^N x_i & \sum_{i=1}^N x_i^2 & \sum_{i=1}^N x_i^3 & \dots & \sum_{i=1}^N x_i^{k+1} \\ \sum_{i=1}^N x_i^2 & \sum_{i=1}^N x_i^3 & \sum_{i=1}^N x_i^4 & \dots & \sum_{i=1}^N x_i^{k+2} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^N x_i^k & \sum_{i=1}^N x_i^{k+1} & \sum_{i=1}^N x_i^{k+2} & \dots & \sum_{i=1}^N x_i^{2k} \end{pmatrix} \quad (20)$$

$$Y = \begin{pmatrix} \sum_{i=1}^N y_i \\ \sum_{i=1}^N y_i x_i \\ \sum_{i=1}^N y_i x_i^2 \\ \vdots \\ \sum_{i=1}^N y_i x_i^k \end{pmatrix} \quad (21)$$

Vale a pena chamar atenção para a forma extremamente simples das matrizes  $A$  e  $Y$ . Para ajustar-se um polinômio de grau  $k$ , as matrizes  $A$  e  $Y$  apresentam as dimensões do polinômio a ser ajustado. A forma simétrica da matriz  $A$  possibilita o emprego de métodos computacionais eficientes para a solução da Eq.19, como mencionado anteriormente.

Embora ajustes polinomiais sejam um assunto de interesse geral, vamos nos concentrar neste trabalho em um aspecto particular relacionado ao ajuste de funções não-lineares. O conjunto de equações utilizadas para efetuar regressões lineares pode ser utilizado para o ajuste um conjunto de funções não-lineares. Toda função não-linear que puder ser convertida na equação da reta poderá ser ajustada como uma função linear. A tabela 1 mostra alguns exemplo de funções que podem ser linearizadas. Um conjunto mais completo de equações pode ser encontrado na literatura<sup>9</sup>.

Entretanto, a utilização do método de regressão linear sobre funções como representadas pela tabela 1 devem ser precedidos de alguns cuidados. Um dos aspectos mais importantes a

**Tabela 1.** Exemplos de linearização de algumas funções não-lineares<sup>9</sup>.

Função não-linear ( $y=f(x)$ )	Função linearizada ( $Y=AX+B$ )
$y = a/x$	$Y=y, X=1/x, A=a, B=0$
$y = a \ln x$	$Y=y, X=\ln x, A=a, B=0$
$y = ae^{bx}$	$Y=\ln y, X=x, A=a, B=\ln b$
$y = ax^b$	$Y=\ln y, X=\ln x, A=\ln a, B=b$
$y = \frac{x}{a+bx}$	$Y=1/y, X=1/x, A=a, B=b$
$y = ab^{1/x}$	$Y=\ln y, X=1/x, A=\ln a, B=\ln b$

serem considerados nestes ajustes é a necessidade de utilização de um fator peso<sup>9,10</sup>. Por exemplo, suponhamos que um conjunto de dados experimentais ajusta-se adequadamente a uma função exponencial do tipo:

$$y = ae^{bx} \quad (22)$$

que ilustra o caso de uma reação cinética irreversível de primeira ordem.

Como mostrado na tabela 1, esta função pode ser transformada em uma representação linear se tomarmos o seu logaritmo, ou seja:

$$\ln y = \ln a + bx \quad (23)$$

Desta forma, assumindo-se  $Y = \ln y$ ,  $A = \ln a$ ,  $B = b$  e  $X = x$ , tem-se uma equação linear:

$$Y = A + BX \quad (24)$$

Pode-se assim utilizar o método de regressão linear para ajustar-se os dados experimentais convenientemente transformados para determinar  $A$  e  $B$ . Entretanto, este procedimento estará minimizando os desvios do conjunto de valores em  $Y$  e não em  $y$  como desejado, ou seja, o desvio entre os valores estimados ( $Y^{est}$ ) e experimentais ( $Y^{exp}$ ) estará sendo minimizado para:

$$Q'' = \sum_{i=1}^N (Y_i^{exp} - Y_i^{est})^2 = \sum_{i=1}^N (\Delta Y_i)^2 \quad (25)$$

e não para  $Q$  apresentado na Eq.1. A função que realmente se deseja ajustar aos dados experimentais está representada pela Eq.22 em termos de  $y$ , portanto deve-se procurar avaliar se esta troca de variáveis implicará em erros na função ajustada. Será que os desvios em  $Y$  são equivalentes aos desvios em  $y$ ? Para responder a esta pergunta deve-se verificar qual tipo de correlação existe entre  $\Delta Y_i$  e  $\Delta y_i$ . Sabe-se que  $Y_i = \ln y_i$ , portanto, pode-se dizer que:

$$\frac{dY_i}{dy_i} = \frac{d \ln y_i}{dy_i} = \frac{1}{y_i} \quad (26)$$

ou então:

$$dY_i = \frac{1}{y_i} dy_i \quad (27)$$

Desta forma, se os valores numéricos de  $\Delta Y_i$  e  $\Delta y_i$  forem muito pequenos, pode-se assumir que  $\Delta Y_i \approx dY_i$  e  $\Delta y_i \approx dy_i$ , ou seja, pode-se dizer que:

$$\Delta Y_i = \frac{1}{y_i} \Delta y_i \quad (28)$$

Para se ajustar a função exponencial utilizando-se regressão linear, deve-se então minimizar  $Q$  através da Eq.1, substituindo-se os valores de  $\Delta y_i$  por  $y_i \Delta Y_i$ . Reescrevendo-se a Eq.1 tem-se:

$$Q = \sum_{i=1}^N (\Delta y_i)^2 = \sum_{i=1}^N (y_i \Delta Y_i)^2 = \sum_{i=1}^N W_i (\Delta Y_i)^2 \quad (29)$$

onde  $W_i = y_i^2$  é o fator peso correspondente. A minimização de  $Q$  com o fator  $W_i$  irá introduzir alterações nas equações resultantes para a regressão linear. As Eq.8 e 9 se apresentarão sob a forma:

$$A = \frac{\sum_{i=1}^N W_i \sum_{i=1}^N W_i X_i Y_i - \sum_{i=1}^N W_i X_i \sum_{i=1}^N W_i Y_i}{\sum_{i=1}^N W_i \sum_{i=1}^N W_i X_i^2 - (\sum_{i=1}^N W_i X_i)^2} \quad (30)$$

$$B = \frac{\sum_{i=1}^N W_i X_i^2 \sum_{i=1}^N W_i Y_i - \sum_{i=1}^N W_i X_i \sum_{i=1}^N W_i X_i Y_i}{\sum_{i=1}^N W_i \sum_{i=1}^N W_i X_i^2 - (\sum_{i=1}^N W_i X_i)^2} \quad (31)$$

onde de acordo com a tabela 1,  $Y_i = \ln y_i$ ,  $X_i = x_i$ ,  $A = a$ ,  $B = \ln b$  e  $W_i = y_i^2$ . Comparando-se as Eqs.30 e 31 com as Eqs.8 e 9, pode-se verificar que a diferença básica entre as mesmas está no fato de que em todos os somatórios foram incorporados o fator peso  $W_i$  e que  $N$  foi substituído por  $\sum W_i$ . Estas modificações podem ser também apresentadas sob a forma matricial. Apenas para exemplificar, a Eq.17 seria reescrita como:

$$\begin{pmatrix} \sum_{i=1}^N W_i Y_i \\ \sum_{i=1}^N W_i Y_i X_i \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^N W_i & \sum_{i=1}^N W_i X_i \\ \sum_{i=1}^N W_i X_i & \sum_{i=1}^N W_i X_i^2 \end{pmatrix} \begin{pmatrix} A \\ B \end{pmatrix} \quad (32)$$

## REGRESSÃO MULTILINEAR

Todos os casos acima levam em consideração o fato de que está sendo feita uma correlação entre uma única variável  $x$  com uma função  $y$  através de uma equação linear ou de uma equação que possa ser linearizada. Entretanto, pode-se imaginar um modelo mais complexo envolvendo regressão linear, mas não envolvendo uma única variável  $x$ , mas um conjunto com várias variáveis independentes,  $\{x^{(1)}, x^{(2)}, \dots, x^{(k)}\}$ , representadas por funções lineares. Este poderia ser o caso da dependência do rendimento de um processo químico com a temperatura, pressão e concentração do catalizador<sup>2</sup>. Em outras palavras, pode-se imaginar que em determinada situação pode-se ter uma função  $y$  que dependa de  $x^{(1)}$ ,  $x^{(2)}$ ,  $x^{(3)}$ , ...,  $x^{(k)}$  e que a relação que exista entre  $y$  e os diferentes conjuntos de  $x$  seja do tipo:

$$y = a + b_1 x^{(1)} + b_2 x^{(2)} + b_3 x^{(3)} + \dots + b_k x^{(k)} \quad (33)$$

Da mesma maneira como apresentado neste trabalho, o método dos mínimos quadrados pode ser aplicado aqui para ajustar esta função a um conjunto de dados experimentais através da determinação de  $a$ ,  $b_1$ ,  $b_2$ ,  $b_3$ , ...,  $b_k$ . A situação neste caso é mais complicada devido ao fato de que para cada tipo de variável  $x$ , tem-se um conjunto de pontos com um número arbitrário qualquer, ou seja, fixando-se os valores de  $x^{(2)}$ ,  $x^{(3)}$ , ...,  $x^{(k)}$  pode-se variar  $x^{(1)}$ , obtendo-se diferentes valores de  $y$ . O mesmo pode ser feito com cada um dos outros tipos de variáveis. Desta forma, a equação acima com  $k$  elementos poderá ter  $N$  pontos. A solução para a equação acima pode ser obtida de maneira análoga a apresentada anteriormente para regressão linear ou polinomial. Uma vez que tem-se um conjunto de  $k+1$  parâmetros a ser determinado, deve-se necessariamente utilizar um conjunto de  $k+1$  equações. Escreve-se a equação acima  $k+1$  vezes e multiplica-se a primeira equação por 1, a segunda por  $x^{(1)}$ , a terceira por  $x^{(2)}$ , a quarta por  $x^{(3)}$  e assim sucessivamente obtendo-se:

$$\begin{aligned} y &= a + b_1 x^{(1)} + b_2 x^{(2)} + \dots + b_k x^{(k)} \\ yx^{(1)} &= ax^{(1)} + b_1 x^{(1)} x^{(1)} + b_2 x^{(2)} x^{(1)} + \dots + b_k x^{(k)} x^{(1)} \\ yx^{(2)} &= ax^{(2)} + b_1 x^{(1)} x^{(2)} + b_2 x^{(2)} x^{(2)} + \dots + b_k x^{(k)} x^{(2)} \\ &\vdots \\ yx^{(k)} &= ax^{(k)} + b_1 x^{(1)} x^{(k)} + b_2 x^{(2)} x^{(k)} + \dots + b_k x^{(k)} x^{(k)} \end{aligned} \quad (34)$$

Tendo-se  $N$  pontos experimentais deve-se efetuar o somatório sobre esses  $N$  pontos para as  $k+1$  equações, obtendo-se assim:

$$\begin{aligned} \sum_{i=1}^N y_i &= aN + \sum_{i=1}^N b_1 x_i^{(1)} + \sum_{i=1}^N b_2 x_i^{(2)} + \dots + \sum_{i=1}^N b_k x_i^{(k)} \\ \sum_{i=1}^N y_i x_i^{(1)} &= a \sum_{i=1}^N x_i^{(1)} + b_1 \sum_{i=1}^N x_i^{(1)} x_i^{(1)} + b_2 \sum_{i=1}^N x_i^{(2)} x_i^{(1)} + \dots + b_k \sum_{i=1}^N x_i^{(k)} x_i^{(1)} \\ \sum_{i=1}^N y_i x_i^{(2)} &= a \sum_{i=1}^N x_i^{(2)} + b_1 \sum_{i=1}^N x_i^{(1)} x_i^{(2)} + b_2 \sum_{i=1}^N x_i^{(2)} x_i^{(2)} + \dots + b_k \sum_{i=1}^N x_i^{(k)} x_i^{(2)} \\ &\vdots \\ \sum_{i=1}^N y_i x_i^{(k)} &= a \sum_{i=1}^N x_i^{(k)} + b_1 \sum_{i=1}^N x_i^{(1)} x_i^{(k)} + b_2 \sum_{i=1}^N x_i^{(2)} x_i^{(k)} + \dots + b_k \sum_{i=1}^N x_i^{(k)} x_i^{(k)} \end{aligned} \quad (34)$$

ou matricialmente:

$$\begin{pmatrix} \sum_{i=1}^N y_i \\ \sum_{i=1}^N y_i x_i^{(1)} \\ \sum_{i=1}^N y_i x_i^{(2)} \\ \vdots \\ \sum_{i=1}^N y_i x_i^{(k)} \end{pmatrix} = \begin{pmatrix} N & \sum_{i=1}^N x_i^{(1)} & \sum_{i=1}^N x_i^{(2)} & \dots & \sum_{i=1}^N x_i^{(k)} \\ \sum_{i=1}^N x_i^{(1)} & \sum_{i=1}^N x_i^{(1)} x_i^{(1)} & \sum_{i=1}^N x_i^{(2)} x_i^{(1)} & \dots & \sum_{i=1}^N x_i^{(k)} x_i^{(1)} \\ \sum_{i=1}^N x_i^{(2)} & \sum_{i=1}^N x_i^{(1)} x_i^{(2)} & \sum_{i=1}^N x_i^{(2)} x_i^{(2)} & \dots & \sum_{i=1}^N x_i^{(k)} x_i^{(2)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^N x_i^{(k)} & \sum_{i=1}^N x_i^{(1)} x_i^{(k)} & \sum_{i=1}^N x_i^{(2)} x_i^{(k)} & \dots & \sum_{i=1}^N x_i^{(k)} x_i^{(k)} \end{pmatrix} \begin{pmatrix} a \\ b_1 \\ b_2 \\ \vdots \\ b_k \end{pmatrix} \quad (35)$$

Uma outra situação a ser considerada durante um processo de ajuste é obtenção de diferentes retas empregando-se exatamente o mesmo conjunto de variáveis. Em outras palavras, suponha que experimentalmente um indivíduo controle seu experimento através de um conjunto de variáveis  $x$ , pertencentes a uma única substância, e de alguma outra variável  $z$  de tal maneira que independente de  $z$ , a relação entre  $y$  e  $x$  será uma função linear. Neste caso, pode-se desejar estabelecer explicitamente a relação entre  $y$  e  $x$  especificando-se posteriormente o valor de  $z$ . Matematicamente pode-se sugerir a seguinte possibilidade, para 3 valores de  $z$  e uma única substância:

**experimento 1:** para um conjunto de dados  $\{y, x, z\}$ , fixa-se o valor de  $z^{(1)}$ , varia-se o valor de  $x$  e obtêm-se diferentes valores de  $y$ . Para este caso, pode-se ter como função ajustável uma equação linear (Eq.10).

Pode-se repetir o mesmo experimento para o mesmo conjunto de variáveis  $x$  modificando-se apenas o valor de  $z^{(1)}$  para  $z^{(2)}$  e  $z^{(3)}$ . Certamente serão obtidas retas ajustadas como:

**experimento 2:** dados  $\{y', x, z^{(2)}\}$ , função linear:

$$y' = a' + b'x \quad (36)$$

**experimento 3:** dados  $\{y'', x, z^{(3)}\}$ , função linear:

$$y'' = a'' + b''x \quad (37)$$

Para  $m$  valores de  $z$ , outros experimentos devem ser repetidos com o mesmo conjunto de variáveis  $x$ , mas para ilustrar um tratamento genérico serão empregados apenas os três experimentos acima. Os valores de  $a, b, a', b', a''$  e  $b''$  podem ser determinados por ajustes das funções lineares para cada caso, ou podem ser determinados simultaneamente se as equações específicas para cada regressão for escrita de maneira adequada. Como visto anteriormente, a solução das equações de regressão linear podem ser escritas segundo as Eq.17 e 18. Estas diferentes equações, como as representadas pelos três experimentos acima, podem ser reescritas na forma matricial como:

$$\begin{pmatrix} \sum_{i=1}^N y_i & \sum_{i=1}^N y'_i & \sum_{i=1}^N y''_i \\ \sum_{i=1}^N y_i x_i & \sum_{i=1}^N y'_i x_i & \sum_{i=1}^N y''_i x_i \end{pmatrix} = \begin{pmatrix} N & \sum_{i=1}^N x_i \\ \sum_{i=1}^N x_i & \sum_{i=1}^N x_i^2 \end{pmatrix} \begin{pmatrix} a & a' & a'' \\ b & b' & b'' \end{pmatrix} \quad (38)$$

A observação desta equação mostra que o método de solução é idêntico ao apresentado anteriormente, invertendo-se a matriz quadrada contendo apenas as variáveis independentes e multiplicando-se esta matriz inversa do lado esquerdo da equação acima. O resultado será uma matriz contendo os valores dos coeficientes de regressão linear.

Para a utilização de mais de um elemento o raciocínio utilizado é o mesmo. Deve-se entretanto utilizar a Eq.35 para representar os experimentos necessários. Embora o processo tenha sido exemplificado apenas para o ajuste linear simples, este método pode também ser empregado para ajustes polinomiais ou regressões multilíneas. Este método foi testado utilizando-se os dados de um trabalho no qual determinou-se simultaneamente as concentrações de cobalto, cobre e níquel por regressão multilinear<sup>12</sup> e os resultados foram bastante satisfatórios. A diferença entre os diferentes ajustes está na forma das matrizes, que podem ser facilmente construídas observando-se a Eq.38 acima e as formas adequadas para cada caso.

Para ilustrar este exemplo específico utilizou-se o processo acima descrito empregando-se os dados de concentração e absorvância da ref.[12]. Na tabela 2 encontram-se os dados referentes a sete misturas com concentrações de Co(II), Cu(II) e Ni(II) e na tabela 3 as respectivas absorvâncias das sete misturas. Observando-se as tabelas 2 e 3 verifica-se que são empregados 3 elementos químicos e que procurar-se-á estabelecer uma relação linear entre a absorvância medida e a concentração dos três componentes simultaneamente. O experimento foi realizado em 5 comprimentos de onda diferentes. Desta forma, deve-se procurar um conjunto de 5 equações do tipo apresentado na Eq.33 para cada comprimento de onda. Uma vez que deseja-se resolver as cinco equações simultaneamente, um sistema matricial semelhante ao da Eq.35 pode ser construído levando-se em conta os cinco experimentos. Para esta situação, o sistema matricial resultante adquire a seguinte expressão:

$$\begin{pmatrix} \sum_{i=1}^N y_i^I & \sum_{i=1}^N y_i^{II} & \sum_{i=1}^N y_i^{III} & \sum_{i=1}^N y_i^{IV} & \sum_{i=1}^N y_i^V \\ \sum_{i=1}^N y_i^I x_i^{(1)} & \sum_{i=1}^N y_i^{II} x_i^{(1)} & \sum_{i=1}^N y_i^{III} x_i^{(1)} & \sum_{i=1}^N y_i^{IV} x_i^{(1)} & \sum_{i=1}^N y_i^V x_i^{(1)} \\ \sum_{i=1}^N y_i^I x_i^{(2)} & \sum_{i=1}^N y_i^{II} x_i^{(2)} & \sum_{i=1}^N y_i^{III} x_i^{(2)} & \sum_{i=1}^N y_i^{IV} x_i^{(2)} & \sum_{i=1}^N y_i^V x_i^{(2)} \\ \sum_{i=1}^N y_i^I x_i^{(3)} & \sum_{i=1}^N y_i^{II} x_i^{(3)} & \sum_{i=1}^N y_i^{III} x_i^{(3)} & \sum_{i=1}^N y_i^{IV} x_i^{(3)} & \sum_{i=1}^N y_i^V x_i^{(3)} \end{pmatrix} = \begin{pmatrix} N & \sum_{i=1}^N x_i^{(1)} & \sum_{i=1}^N x_i^{(2)} & \sum_{i=1}^N x_i^{(3)} \\ \sum_{i=1}^N x_i^{(1)} & \sum_{i=1}^N x_i^{(1)} x_i^{(1)} & \sum_{i=1}^N x_i^{(2)} x_i^{(1)} & \sum_{i=1}^N x_i^{(3)} x_i^{(1)} \\ \sum_{i=1}^N x_i^{(2)} & \sum_{i=1}^N x_i^{(1)} x_i^{(2)} & \sum_{i=1}^N x_i^{(2)} x_i^{(2)} & \sum_{i=1}^N x_i^{(3)} x_i^{(2)} \\ \sum_{i=1}^N x_i^{(3)} & \sum_{i=1}^N x_i^{(1)} x_i^{(3)} & \sum_{i=1}^N x_i^{(2)} x_i^{(3)} & \sum_{i=1}^N x_i^{(3)} x_i^{(3)} \end{pmatrix} \begin{pmatrix} a^I & a^{II} & a^{III} & a^{IV} & a^V \\ b_1^I & b_1^{II} & b_1^{III} & b_1^{IV} & b_1^V \\ b_2^I & b_2^{II} & b_2^{III} & b_2^{IV} & b_2^V \\ b_3^I & b_3^{II} & b_3^{III} & b_3^{IV} & b_3^V \end{pmatrix} \quad (39)$$

**Tabela 2.** Composição de soluções constituídas por mistura de Co(II), Cu(II) e Ni(II)<sup>12</sup>. Concentrações expressas em mmol/l.

Soluções	Co(II)	Cu(II)	Ni(II)
1	0.000	9.052	0.000
2	9.958	0.000	0.000
3	0.000	0.000	16.17
4	5.975	9.052	3.233
5	9.958	5.431	3.233
6	5.975	1.810	16.17
7	1.992	1.810	16.17

**Tabela 3.** Absorbâncias medidas em diferentes comprimentos de onda para um conjunto de soluções constituídas por mistura de Co(II), Cu(II) e Ni(II) com concentrações especificadas na Tabela 2<sup>12</sup>.

Soluções	868,9nm	732,0nm	585,6nm	462,9nm	378,7nm
1	0,414	0,862	0,169	0,021	0,070
2	0,031	0,021	0,042	0,157	0,031
3	0,221	0,049	0,145	0,029	0,219
4	0,475	0,897	0,214	0,108	0,111
5	0,319	0,547	0,163	0,164	0,099
6	0,326	0,237	0,199	0,117	0,240
7	0,312	0,225	0,183	0,059	0,230

Na Eq.38 os algarismos romanos indicam os cinco comprimentos de onda (I=868,9nm, II=732,0nm, III=585,6nm, IV=462,9nm e V=378,7nm) e certamente os dados com estes expoentes referem-se a dados coletados no respectivo comprimento de onda. Os expoentes (1), (2) e (3) referem-se a índices que identificam o elemento químico a que estão associados. Por exemplo, (1) pode referir-se aos dados do Co, (2) ao Cu e (3) ao Ni. Os valores de  $y$  são os valores de absorbância e os valores de  $x$  correspondem aos valores de concentração. Uma vez que foram empregadas 7 soluções com composições diferentes, o valor de  $N$  deve ser 7. Substituindo-se os valores apresentados nas tabelas 2 e 3 na Eq.38 e resolvendo-se o sistema matricial de acordo com a Eq.19 obtém-se os valores dos coeficientes  $a$  e  $b$  em perfeita concordância com os valores da ref.[12]. Ou seja:

	868,9nm	732,0nm	585,6nm	462,9nm	378,7nm
$a$	0,0006	-0,0065	0,0095	0,0144	0,0197
$b_1$	2,9	2,6	3,2	14,3	1,0
$b_2$	45,5	96,2	17,5	0,6	5,3
$b_3$	13,8	3,3	8,4	0,9	12,4

## REFERÊNCIAS

1. Draper, N e Smith, H; *Applied Regression Analysis*; 2a. edição; John Wiley & Sons, New York, 1981, pp.1-136.
2. Montgomery, D. C.; *Design and Analysis of Experiments*; 3a.edição; John Wiley & Sons, New York, 1991, pp.479-520.
3. Miller, J. N.; *Analyst* **1991**, *116*, 3.
4. Pimentel, M. F. e Barros Neto, B.; *Quím. Nova* **1996**, *19*, 268.
5. Barros Neto, B.; Scarmínio, I. S. e Bruns, R. E.; *Planejamento e Otimização de Experimentos*; Campinas, Editora da Unicamp, 1995.
6. Eiras, S. P. e Andrade, J. C.; *Quím. Nova* **1996**, *19*, 24.
7. Spiegel, M. R., *Estatística*; 11ª. edição; McGraw-Hill do Brasil, São Paulo, 1979.
8. Press, W. H.; Teukolsky, S. A.; Vetterling, W. T. e Flannery, B. P., *Numerical Recipes in Fortran*; 2nd. edition, Cambridge University Press, New York, 1992.
9. de Levie, R.; *J. Chem. Educ.* **1986**, *63*, 10.
10. Sands, D. E.; *J. Chem. Educ.* **1974**, *51*, 473.
11. O'Neil, R. T. e Flaspohler, D. C.; *J. Chem. Educ.* **1990**, *61*, 40.
12. Dado, G. e Rosenthal, J.; *J. Chem. Educ.* **1990**, *67*, 797.