

Linear Regression

50

In HW problem, wrote routine to fit a parabola to 3 points. In this case there was only one answer.

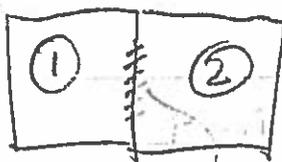
Often in data analysis you want to fit a curve to many data points

How do we do this? Use (usually) linear regression

This does not mean that we fit data with a line, rather that we use an algorithm which reduces the problem to a set of linear equations!

Let's look at a simple problem: measure mass transfer coefficients

Suppose we have a membrane between two reservoirs



↳ membrane area = A

51

We start with some concentration in reservoir ① of C_0 and we assume that reservoir ② is maintained at a constant C_{eq} . If the volume of reservoir ① is V , then:

$$V \frac{dC}{dt} = -Ah(C - C_{eq})$$

where $C \Big|_{t=0} = C_0$

We can solve this equation:

$$C = C_{eq} + (C_0 - C_{eq}) e^{-\left(\frac{hA}{V}\right)t}$$

We wish to determine h by measuring C as a function of time

In order for us to use linear regression, the model must be linear in the modelling parameter,

In this case, the modelling parameter is h . We can rewrite the eq'n:

52

$$(c - c_{eq}) = (c_0 - c_{eq}) e^{-\frac{hA}{V}t}$$

$$\ln(c - c_{eq}) = \ln(c_0 - c_{eq}) - \frac{hA}{V}t$$

So if we plot $\ln(c - c_{eq})$ vs $\frac{At}{V}$

we get a line w/ intercept of $\ln(c_0 - c_{eq})$ and a slope of h !

How do we get the best fit line?

We look at the deviation of the points from the line, and try to minimize this in some way.

Let's take:

$$y \equiv \ln(c - c_{eq}), \quad a \equiv \frac{Ah}{V}$$

$$b \equiv \ln(c_0 - c_{eq})$$

We want to fit a series of points (t_i, y_i) by the model $y = at + b$

(53)

We shall use the method of least squares. We form the sum:

$$\text{Sum} = \sum_{i=1}^N (y_i - (at_i + b))^2$$

which is the sum of the squared distance between data pt. and model in the y direction.

We pick a & b so that this is a minimum

We let:

$$\frac{\partial \text{sum}}{\partial a} = 0$$

$$\frac{\partial \text{sum}}{\partial b} = 0$$

} 2 eqns for a, b !

$$\frac{\partial \text{sum}}{\partial a} = \sum_{i=1}^N -2(y_i - (at_i + b))t_i$$

$$= \sum_{i=1}^N -2y_it_i + 2a \sum_{i=1}^N t_i^2 + 2b \sum_{i=1}^N t_i = 0$$

and:

$$\frac{\partial \text{sum}}{\partial b} = \sum_{i=1}^N -2(y_i - (at_i + b))$$

$$= \sum_{i=1}^N -2y_i + 2a \sum_{i=1}^N t_i + 2Nb = 0$$

This is a pair of linear equations for a & b !

We have the solution:

$$a\bar{t}^2 + b\bar{t} = \bar{y}\bar{t} \quad (\text{let } \bar{x} \equiv \frac{1}{N} \sum_{i=1}^N x_i)$$

$$a\bar{t} + b = \bar{y}$$

$$\therefore a = \frac{\bar{y}\bar{t} - \bar{y}\bar{t}}{\bar{t}^2 - \bar{t}^2}$$

$$\text{and } b = \bar{y} - a\bar{t}$$

which can be used to get the mass transfer coefficient!

55

This is not the only way to approach the problem. Let's look at it as a system of equations!

We are trying to satisfy:

$$y_1 = a t_1 + b$$

$$y_2 = a t_2 + b$$

$$\vdots$$

$$y_N = a t_N + b$$

Let's write this in matrix form:

$$\begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} \approx \begin{pmatrix} t_1 & 1 \\ t_2 & 1 \\ \vdots & \vdots \\ t_N & 1 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix}$$

we shall use X to be the vector containing the modelling parameters:

$$X_1 \equiv a \quad X_2 \equiv b$$

If we let the matrix of modelling functions be \underline{A} (e.g. $A_{11} = t_1$, etc.)

then we want to satisfy:

$$\underline{y} \approx \underline{A} \underline{x}$$

Actually, we want to minimize the residual:

$$\underline{r} = \underline{y} - \underline{A} \underline{x}$$

Or, for least squares, adjust \underline{x} s.t.:

$$\|\underline{r}\|_2 \equiv \|\underline{y} - \underline{A} \underline{x}\|_2 \text{ is as small as possible}$$

↑
Euclidean norm

We want to pick \underline{x} s.t.

$$\underline{r}^T \underline{r} \equiv (\underline{y} - \underline{A} \underline{x})^T (\underline{y} - \underline{A} \underline{x})$$

is a minimum.

This is exactly equivalent to what we did before!

Let's multiply this out:

$$r^2 \equiv \underbrace{y^T y}_{\sim} - \underbrace{y^T A}_{\sim} \underbrace{x}_{\sim} - \underbrace{x^T A^T}_{\sim} \underbrace{y}_{\sim} + \underbrace{x^T A^T A}_{\sim} \underbrace{x}_{\sim}$$

We take the gradient of this w.r.t. \underbrace{x}_{\sim} and set it equal to zero!

$$\nabla_{\underbrace{x}_{\sim}} r^2 = 0 = -2 \underbrace{A^T}_{\sim} \underbrace{y}_{\sim} + 2 \underbrace{A^T A}_{\sim} \underbrace{x}_{\sim} = 0$$

Thus we have the problem:

$$\underbrace{A^T A}_{\sim} \underbrace{x}_{\sim} = \underbrace{A^T}_{\sim} \underbrace{y}_{\sim} \quad !$$

This gives us a set of n eq'ns for the n unknowns in \underbrace{x}_{\sim} !

These are known as the normal equations.

Linear Regression Error

98

Let's apply all this back to linear regression. Suppose we have some data and we are trying to fit it, and calc. the uncertainty in the modelling parameters.

We measure a lead weight's pos'n a

We get

t_i (s)	
0	
1	1.8
2	20.7
3	15.9
4	2.0

skip this

(go straight to vector approach)

We fit this to the equation:

$$h = h_0 + v_0 t - \frac{1}{2} g t^2$$

We want to both calculate g , and determine the uncertainty in g .

First let's get g . We have the problem:

$$h_i = (1) h_0 + (t) V_0 + (-\frac{1}{2}t^2) g$$

\uparrow \uparrow \uparrow
 x_1 x_2 x_3

So :

$$\underset{\approx}{A} \begin{pmatrix} 1 & 0 & 0 \\ 1 & 1 & -0.5 \\ 1 & 2 & -2 \\ 1 & 3 & -\frac{9}{2} \\ 1 & 4 & -8 \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} 0.3 \\ 14.8 \\ 20.7 \\ 15.9 \\ 2.0 \end{pmatrix} \underset{\approx}{h}$$

Using the normal eq's we get:

$$\underset{\approx}{A}^T \underset{\approx}{A} \underset{\approx}{x} = \underset{\approx}{A}^T \underset{\approx}{h}$$

$$\underset{\approx}{x} = \left[\underset{\approx}{A}^T \underset{\approx}{A} \right]^{-1} \underset{\approx}{A}^T \underset{\approx}{h}$$

100

Or numerically:

$$X \approx \begin{pmatrix} .886 & .257 & -.086 & -.143 & .086 \\ -.771 & .186 & .571 & .386 & -.371 \\ -.286 & .143 & .286 & .143 & -.286 \end{pmatrix} \begin{pmatrix} h_1 \\ h_2 \\ h_3 \\ h_4 \\ h_5 \end{pmatrix}$$

This yields:

$$x_1 = 0.197, \quad x_2 = 19.7, \quad x_3 = 9.64$$

Now we want to examine the error in the position measurements h_i .

To do this we look at the deviation:

$$S_h^2 \approx \frac{1}{N-3} \sum_{i=1}^N (h_i - \underset{\substack{\uparrow \\ \text{using fitted} \\ \text{parameters}}}{h(t_i)})^2$$

because 3 adjustable parameters det. from data

The quantity $h_i - h(t_i)$ is just the residual (10)

$$\tilde{w} = \tilde{h} - \tilde{A} \tilde{x} !$$

$$\therefore S_h^2 = \frac{1}{N-3} \|\tilde{w}\|_2^2 = 0.110$$

↳ you should really have more pts for accuracy!

Ok, how do we use this to get the error in $g(x_3)$?

We had:

$$\tilde{x} = \left\{ \left[\tilde{A}^T \tilde{A} \right]^{-1} \tilde{A}^T \right\} \tilde{h}$$

$$\equiv \tilde{K}$$

$$\text{Thus } x_i \equiv \sum_{j=1}^N k_{ij} h_j$$

i 'th row of \tilde{K} !

(102)

We can thus use the formula for a linear combination of random variables!

$$S_{X_{ij}}^2 = \sum_{i=1}^N K_{ij}^2 S_{h_i}^2 = 0.0315$$

↑
these are all the same
in this problem!

Thus $g = 9.64 \pm 0.18$

to the 69% confidence level (1σ)

Often we want to develop correlations for their predictive value, for example what is the expected ht. of the lead wt. at some time t ?

$$h_m = x_1 + x_2 t - \frac{1}{2} t^2 x_3 !$$

But what is the error? We have a linear combination of x_1, x_2, x_3 each of which have some error.

You might expect that:

$$S_{hm}^2 = S_{x_1}^2 + (t)^2 S_{x_2}^2 + (-\frac{1}{2}t^2)^2 S_{x_3}^2$$

But this is incorrect !!

The problem is that x_1, x_2, x_3 are not independent! Their covariance is non-zero!

Thus we should have:

$$\begin{aligned} S_{hm}^2 &= S_{x_1}^2 + (t)^2 S_{x_2}^2 + (-\frac{1}{2}t^2)^2 S_{x_3}^2 \\ &+ 2(1)(t) S_{x_1x_2}^2 + 2(1)(-\frac{1}{2}t^2) S_{x_1x_3}^2 \\ &+ 2(t)(-\frac{1}{2}t^2) S_{x_2x_3}^2 \end{aligned}$$

How do we calculate the covariance?

Recall:

$$x_i \equiv \sum_{j=1}^N K_{ij} h_j$$

Where $K \equiv [A^T A]^{-1} A^T$

So:

$$\sigma_{x_1, x_2}^2 \equiv E \left(\left(\sum_{j=1}^N \alpha_j (h_j - \mu_{h_j}) \right) \left(\sum_{j=1}^N \beta_j (h_j - \mu_{h_j}) \right) \right)$$

Now if all of the h_j are independent then

$$E \left((h_j - \mu_{h_j}) (h_i - \mu_{h_i}) \right) = \underline{\underline{0}}$$

for $i \neq j$

$$\therefore \sigma_{x_1, x_2}^2 = E \left(\sum_{j=1}^N \alpha_j \beta_j (h_j - \mu_{h_j})^2 \right)$$

$$= \sum_{j=1}^N \alpha_j \beta_j \sigma_h^2$$

↑
provided all the h_j have the same uncertainty

$$\therefore \sigma_{x_1, x_2}^2 = \alpha \beta^T \sigma_h^2$$

So for this problem if we let

$$\gamma \equiv \text{3rd row of } K$$

106

$$x_1 = \underset{\sim}{\alpha} \underset{\sim}{h}, \quad x_2 = \underset{\sim}{\beta} \underset{\sim}{h}, \quad x_3 = \underset{\sim}{\gamma} \underset{\sim}{h}$$

$$S_h^2 = \frac{1}{N-3} \|\underset{\sim}{r}\|_2^2, \quad \underset{\sim}{r} = \underset{\sim}{h} - \underset{\sim}{A} \underset{\sim}{X}$$

$$S_{x_1}^2 = \left(\underset{\sim}{\alpha} \underset{\sim}{\alpha}^T \right) S_h^2$$

$$S_{x_2}^2 = \left(\underset{\sim}{\beta} \underset{\sim}{\beta}^T \right) S_h^2 \quad \text{etc.}$$

and:

$$h_m = (1)x_1 + (t)x_2 + \left(-\frac{1}{2}t^2\right)x_3$$

$$\begin{aligned} \sigma_{h_m}^2 &= (1)^2 \sigma_{x_1}^2 + (t)^2 \sigma_{x_2}^2 + \left(-\frac{1}{2}t^2\right)^2 \sigma_{x_3}^2 \\ &\quad + 2(1)(t) \underset{\sim}{\alpha} \underset{\sim}{\beta}^T \sigma_h^2 + 2(1)\left(-\frac{1}{2}t^2\right) \underset{\sim}{\alpha} \underset{\sim}{\gamma}^T \sigma_h^2 \\ &\quad + 2(t)\left(-\frac{1}{2}t^2\right) \underset{\sim}{\beta} \underset{\sim}{\gamma}^T \sigma_h^2 \\ &= \left\| \left[(1) \underset{\sim}{\alpha} + (t) \underset{\sim}{\beta} + \left(-\frac{1}{2}t^2\right) \underset{\sim}{\gamma} \right] \right\|_2^2 \sigma_h^2 \end{aligned}$$