

(11)

t cdf - t cum dist  $F^n$

$\Rightarrow$  can do hyp. testing w/  
z test, t test, ttest2  
 $\uparrow$   
2 samples

Ok, what if we want to  
test many samples? Suppose

we look to see if there is a  
difference among many groups

We could do a t-test between  
each sample in a combinatoric  
fashion, but you are guaranteed  
a false positive! In fact,  
talks about looking at relation between

8 groups  $\Rightarrow 7 \times 8 / 2 = 28$  combinations

Even random chance at 95% conf.

level would lead to at least 1 false pos!

(12)

A better way is ANOVA, analysis of Variance

Basically, you look at the Variance of "supergroup" and Variance of all the subgroups.

If the variance of supergroups is larger than expected from subgroups, then at least one of the subgroups doesn't belong.

This ratio yields the F-statistic which has a dist. depending on the number of degrees of freedom.

Matlab has a canned routine `anova1(A)` which does this to a matrix of values.

(13)

Most of these formulas  
require independent normally  
distributed data

If the data is not indep,  
 $S_x$  will be underestimated!  
(Effectively,  $n$  is reduced)

If the data isn't normally  
dist, sometimes you can fix it!

Example : Dist of particle size

- often have ~~the~~ density of small particles  $\gg$  larger ones
- could look at radius vs mass fraction, not ~~#~~ of particles
- Also - log normal dist -
  - work w/ log of variable rather than direct meas
  - plot dist & see what works!

## Linear Regression

(50)

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

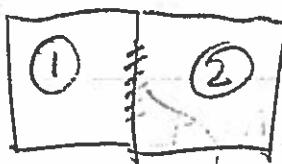
Often in data analysis you want to fit a curve to many data points

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:

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

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

We can solve this equation:

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

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:

$$y \equiv \ln(C - C_{eq}), a \equiv \frac{Ah}{V}$$

$$b \equiv \ln(C_0 - C_{eq})$$

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 :

$$\left. \begin{aligned} \frac{\partial \text{sum}}{\partial a} &= 0 \\ \frac{\partial \text{sum}}{\partial b} &= 0 \end{aligned} \right\} \text{2 eqns for } a, b!$$

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

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

54

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:

$$at^2 + bt = \bar{yt} \quad (\text{let } \bar{x} = \frac{1}{N} \sum_{i=1}^N x_i)$$

$$at + b = \bar{y}$$

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

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

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:

$$y_1 = a t_1 + b$$

$$y_2 = a t_2 + b$$

$$\vdots \quad \vdots \quad \vdots$$

$$y_N = a t_N + b$$

Let's write this in matrix form:

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

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

$$x_1 \equiv a \quad x_2 \equiv b$$

56

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

then we want to satisfy:

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

Actually, we want to minimize the residual:

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

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

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

↑  
Euclidean norm

We want to pick  $\tilde{x}$  s.t.

$$\tilde{r}^T \tilde{r} = (\tilde{y} - \tilde{A} \tilde{x})^T (\tilde{y} - \tilde{A} \tilde{x})$$

is a minimum.

(57)

This is exactly equivalent to what we did before!

Let's multiply this out:

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

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

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

Thus we have the problem:

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

This gives us a set of  $n$  eq'n's for the  $n$  unknowns in  $\tilde{x}$ !

These are known as the normal equations.