University of
# BRISTOL

# Applied Statistics
## Lecture 14+15

## David Barton & Sabine Hauert

### Department of Engineering Mathematics

# Outline

≱ Linear regression

≱ Conditions for linear regression

≱ Residual plots

≱ Linear regression with more complicated functions

## OpenIntro Statistics

## Chapter 7

University of
# BRISTOL

# Linear regression

Data fitting is an important part of statistics.

- Given a particular model, what parameters best fit the data?
  - ▲ Modal analysis is a big part of structural engineering!

- More data than parameters: regular data fitting ← focus

- More parameters than data: inverse problems!
  - ▲ Most of medical imaging
  - ▲ Source reconstruction in acoustics
  - ▲ Optics, radar, communications, nondestructive testing,....
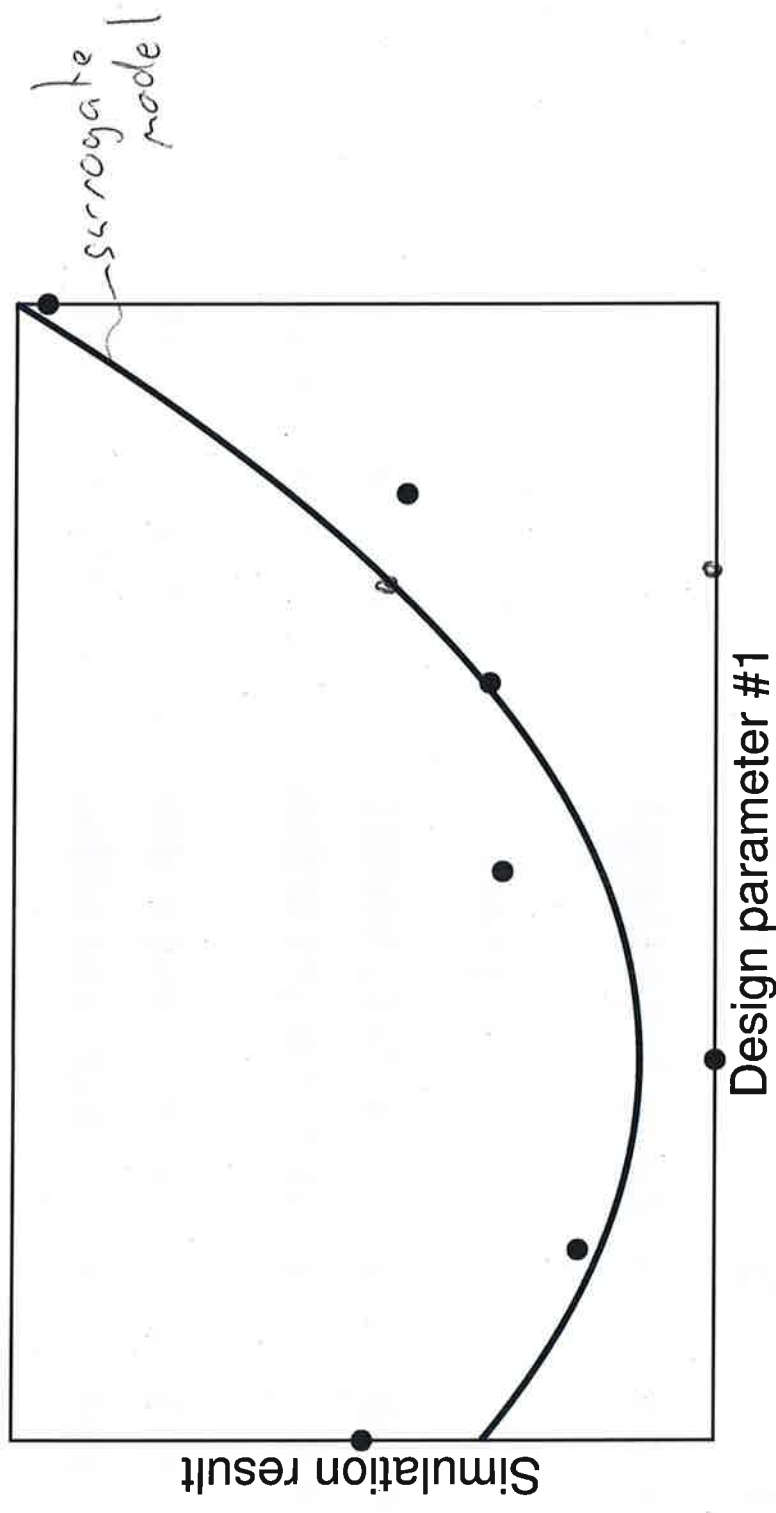
*quite challenging*

# How do you design large-scale structures?

only run costly
simulations a few
times and use the
data points to build
a model that allows
you to make predictions
about your system.

University of
**BRISTOL**

# Finite element simulations are slow...



*surrogate model*

Simulation result

Design parameter #1

Use a *surrogate model* for quick computation and the full model for final testing. (Also known as a metamodel, reduced-order model or emulator.)

# Linear regression

Linear regression is a means for estimating the parameters of a model of the form _approximation_

_general form_

_inputs_

$$y = \beta_1 u_1 + \beta_2 u_2 + \cdots + \beta_p u_p$$

$$y = \beta_1 + \beta_2 Y$$
$$u_1 = 1$$
$$u_2 = ?$$

where $y$ is the dependent (output) variable, $u_i$ are the independent (input) variables, and $\beta_i$ are the parameters to be estimated.

(There is lots of different terminology used in different textbooks — regressor variables, exogenous variables, explanatory variables, etc.)

This model is _linear in the parameters_; it can be that the model is nonlinear in terms of the independent variables! For example,

_linear in parameters_

$$y = \beta_1 x + \beta_2 x^2 + \beta_3 x^3$$

_find the best "$\beta$" values to match inputs to the outputs_

fits in the framework of linear regression.

# Linear regression

With linear regression, $y$ and $x_i$ are known from from $n$ different samples.
Hence we have

$$y_1 = \beta_1 u_{1,1} + \beta_2 u_{1,2} + \cdots + \beta_p u_{1,p}$$

$$y_2 = \beta_1 u_{2,1} + \beta_2 u_{2,2} + \cdots + \beta_p u_{2,p}$$

$$\vdots$$

$$y_n = \beta_1 u_{n,1} + \beta_2 u_{n,2} + \cdots + \beta_p u_{n,p}$$

$$y = \beta_1 + \beta_2 x$$

$$\begin{bmatrix} y \\ b \end{bmatrix} \begin{bmatrix} 1 \cdot 2 \\ 1 \end{bmatrix} \begin{bmatrix} 3/4 \\ 1,2 \end{bmatrix} \begin{bmatrix} 1 \\ 1,3 \end{bmatrix}$$

can't solve exact

$$\left. \begin{array}{l} 1 = \beta_1 + \beta_2 \cdot \frac{1}{2} \\ 1.2 = \beta_1 + \beta_2 \, 3/4 \\ 1.3 = \beta_1 + \beta_2 \, 1 \end{array} \right\} \begin{array}{l} 3 \text{ equation} \\ 2 \text{ unknown} \end{array}$$

In matrix-vector form

$$y = U\beta$$

vector of outputs
vector of parameters
matrix of inputs

$$y = \begin{bmatrix} 1 \cdot 2 \\ 1 \cdot 3 \end{bmatrix} \quad u = \begin{bmatrix} 1 & 1/2 \\ 1 & 3/4 \\ 1 & 1 \end{bmatrix} \quad \beta = \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix}$$

For a particular choice of $\beta$ construct the *residual or error vector*

$$e = y - U\beta$$

residual
error vector
actual output
output of the model

University of
**BRISTOL**

# Minimising error

What choice of $\beta$ minimises the error?

First, what is meant by small since the error is a vector? Obvious answer is the Euclidean norm

$$\|e\| = \left( \sum_i e_i^2 \right)^{\frac{1}{2}}$$

which gives us least-squares.
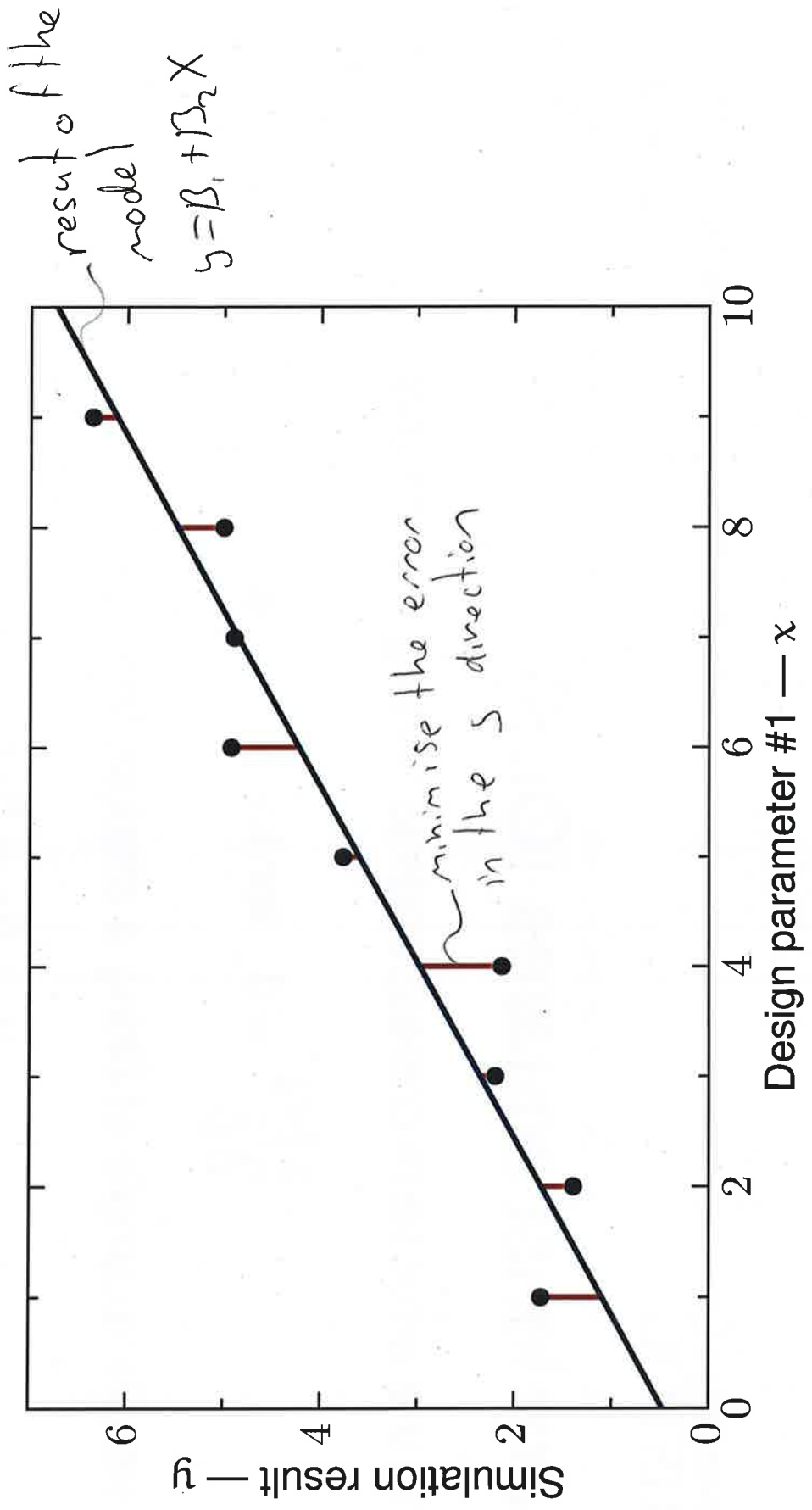
But it's not the only answer!

🔻 Least absolute deviation (more robust in some circumstances)

often better if you have
outliers

$$\|e\| = \sum_{i=1}^n |e_i|$$

🔻 For inverse problems: Tikhonov regularisation or the Lasso method

# Ordinary least squares (OLS)

Ordinary least squares corresponds to minimising the error in the y direction — errors in the x direction are ignored!



*result of the model*
$y = \beta_1 + \beta_2 X$

*minimise the error in the y direction*

Simulation result — y

Design parameter #1 — x

# Ordinary least squares (OLS)

Minimising the error term corresponds to finding $\beta$ such that

$$\frac{d\,\|e\|}{d\,\beta_i} = 0 \quad \text{for } i = 1, \ldots, p$$

Take the simple example of fitting a straight line

$$y = \beta_1 + \beta_2 x$$

Here we have

$$\|e\|^2 = \sum_{i=1}^{n} (\underbrace{y_i}_{\text{output}} - \underbrace{\beta_1 - \beta_2 x_i}_{\text{model output}})^2$$

Note that minimising $\|e\|^2$ gives the same results as minimising $\|e\|$ in this context.

# Ordinary least squares (OLS)

Differentiating $\|e\|^2$ w.r.t. $\beta_1$ and $\beta_2$ and equating to zero gives

$$\frac{d\|e\|^2}{d\beta_1} = -2\sum_{i=1}^{n}(y_i - \beta_1 - \beta_2 x_i) = 0$$

$$\frac{d\|e\|^2}{d\beta_2} = -2\sum_{i=1}^{n}x_i(y_i - \beta_1 - \beta_2 x_i) = 0$$

These linear algebraic equations are called the *normal equations*

$$\beta_1 n + \beta_2 \sum_{i=1}^{n} x_i = \sum_{i=1}^{n} y_i$$

$$\beta_1 \sum_{i=1}^{n} x_i + \beta_2 \sum_{i=1}^{n} x_i^2 = \sum_{i=1}^{n} x_i y_i$$

known values

# Example calculation

$$y = B_1 + B_2 x + B_3 x^2 \quad \text{quadratic model}$$

$$\|e\|^2 = \sum_i (y_i - B_1 - B_2 x_i - B_3 x_i^2)^2$$

$$\frac{d\|e\|^2}{B_1} = -2 \sum_i (y_i - B_1 - B_2 y_i - B_3 x_i^2) = 0$$

$$\frac{d\|e\|^2}{B_2} = -2 \sum_i x_i (y_i - B_1 - B_2 x_i - B_3 x_i^2) = 0$$

$$\frac{d\|e\|^2}{B_3} = -2 \sum_i x_i^2 (y_i - B_1 - B_2 x_i - B_3 x_i^2) = 0$$

$$B_1 n + B_2 \sum x_i + B_3 \sum x_i^2 = \sum y_i$$
$$B_1 \sum x_i + B_2 \sum x_i^2 + B_3 \sum x_i^3 = \sum x_i y_i$$
$$B_1 \sum x_i^2 + B_2 \sum x_i^3 + B_3 \sum x_i^4 = \sum x_i^2 y_i$$

solve this
to find B

→ use matlab
polsfit

University of
## BRISTOL

✗

# Maximum likelihood

The errors in the $y$ direction are often assumed to be normally distributed i.i.d.; in this case least-squares is a maximum likelihood estimator (MLE)

Given a set of points $y_i$ (ignore $x$ for simplicity), what are the parameters of the normal distribution that are most likely to have generated that set of points?
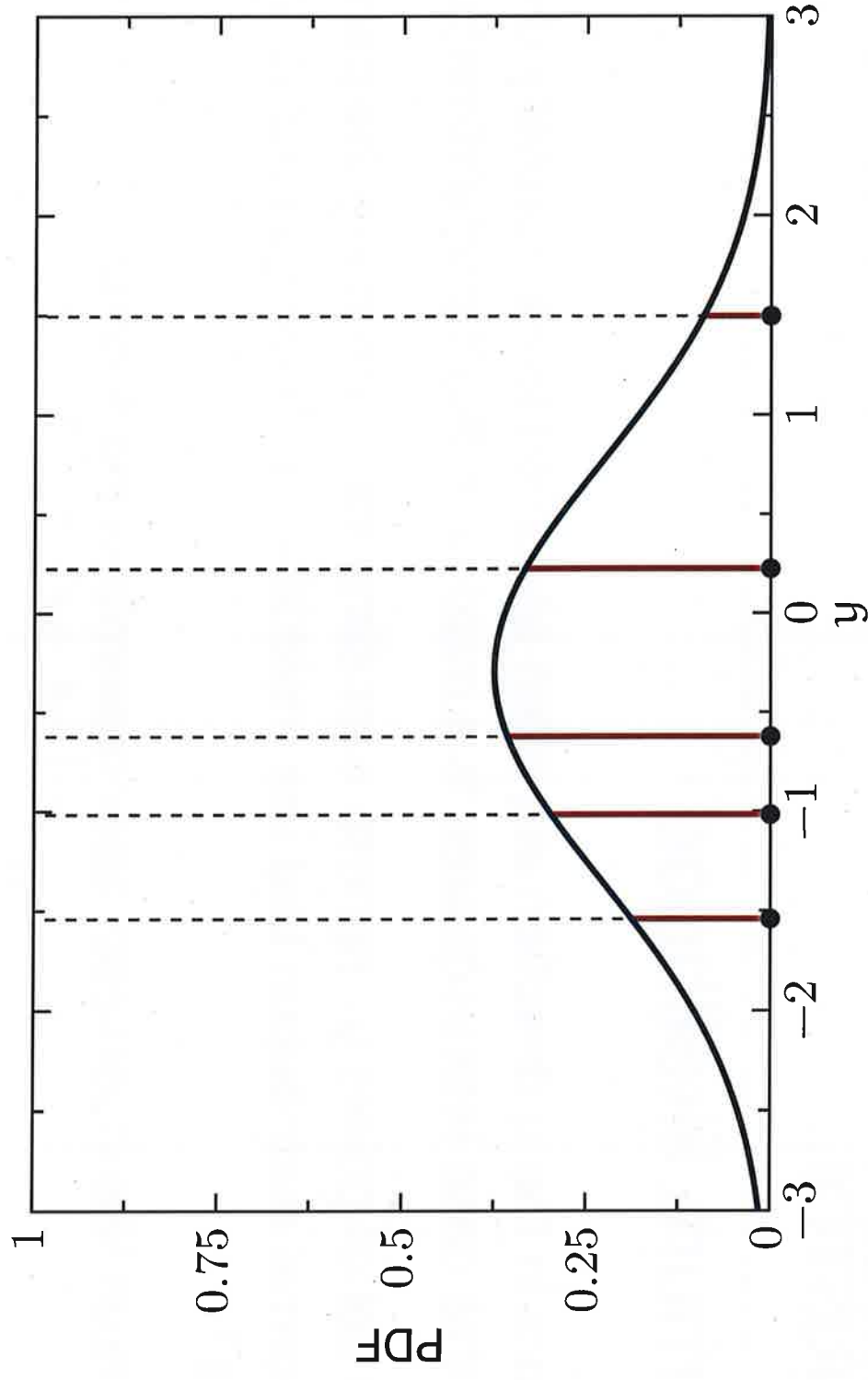
A *maximum likelihood estimator* maximises the function

↳ find a normal distribution that best fits, and
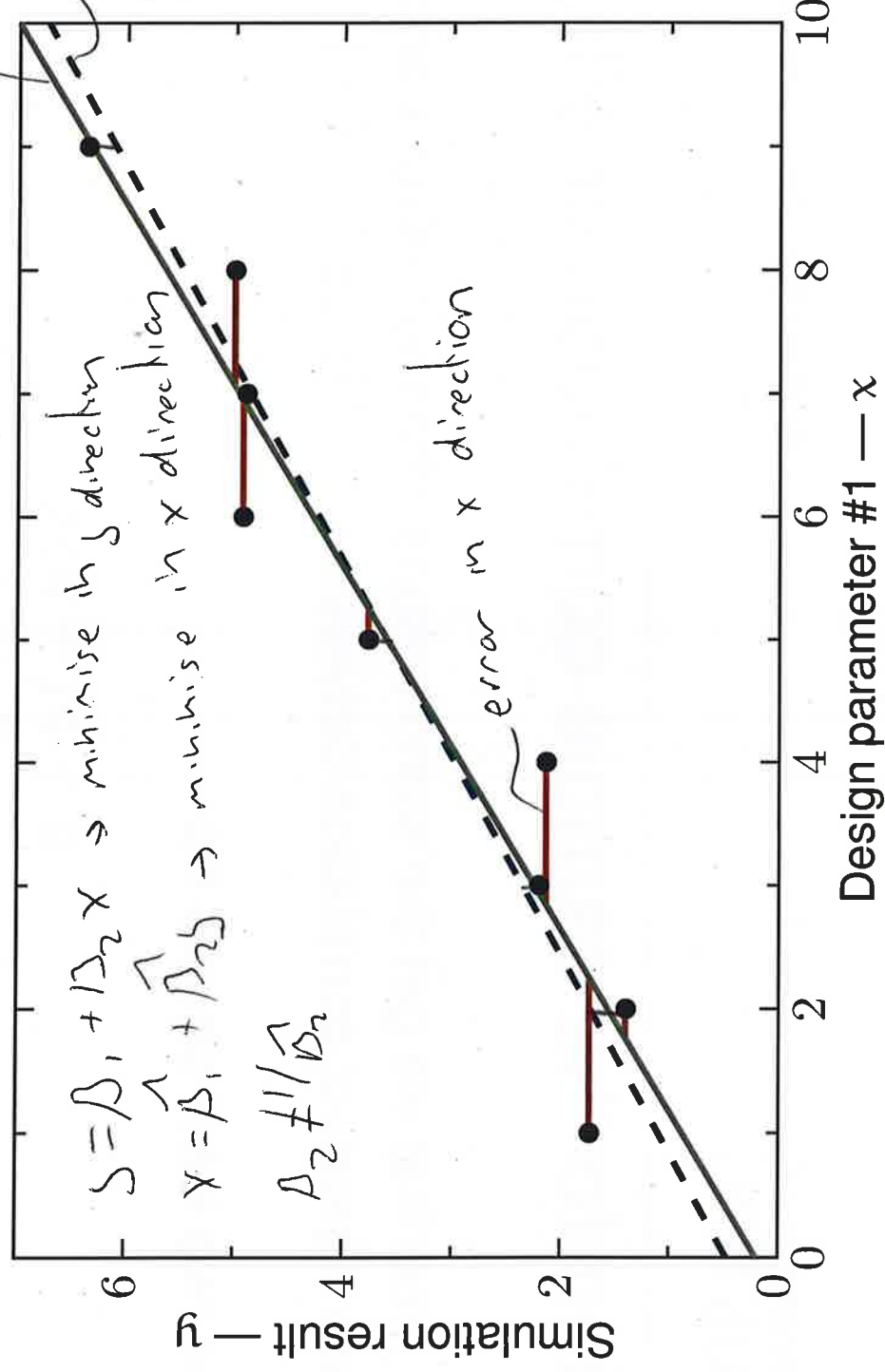
$$\prod_{i=1}^{n} p_Y(y_i \mid \mu_Y, \sigma_Y^2)$$

where $y_i$ are i.i.d. samples and $p_Y$ is the corresponding probability density function. (Often maximise the log of this function.)

University of
BRISTOL

# Maximum likelihood

University of
BRISTOL

# Errors in the x direction

We've assumed no errors in the x direction — swapping the x and y coordinates gives different answers!

University of
BRISTOL

# Relation to the sample correlation coefficient

The only time that we get the same results when we swap $x$ and $y$ is when there is perfect correlation between $x$ and $y$; i.e., $r = \pm 1$.

If we have two least-square fits

$$y = \beta_1 + \beta_2 x$$

$$x = \hat{\beta}_1 + \hat{\beta}_2 y$$

then it's possible to show that

$$\beta_2 \cdot \hat{\beta}_2 = r^2$$

Hence $\beta_2 = 1/\hat{\beta}_2$ only when $r = \pm 1$.

# Exercise

Calculate line of best fit

$x$ = death by entanglement

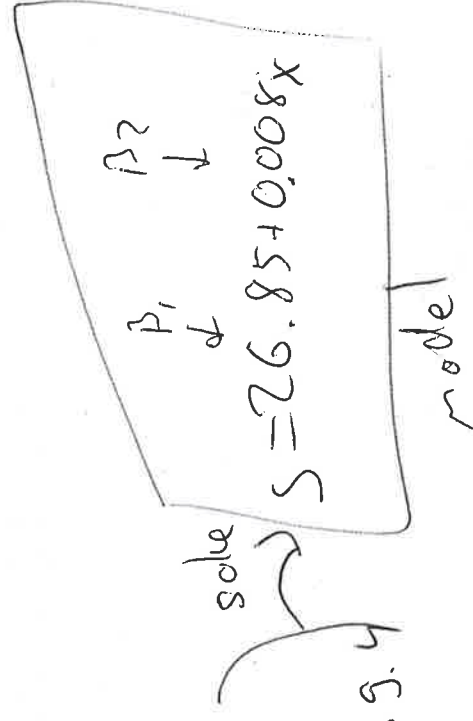$s$ = cheese consumption

$$s = \beta_1 + \beta_2 x$$

normal equations

$$n\beta_1 + \Sigma x_i \beta_2 = \Sigma s_i$$
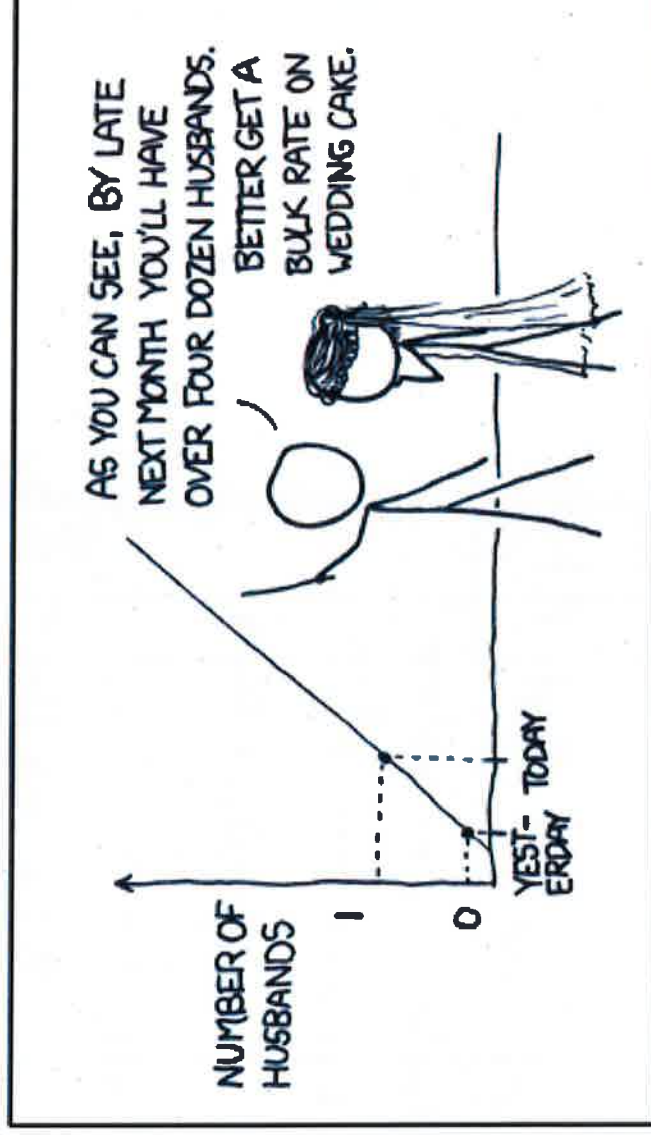
$$\Sigma x_i \beta_1 + \Sigma x_i^2 \beta_2 = \Sigma x_i y_i$$

$$10\beta_1 + 5886\beta_2 = 315.2$$

$$5886\beta_1 + 2659072\beta_2 = 187066.4$$

solve $\rightsquigarrow$ $\beta_1$, $\beta_2$

$\beta_1 \downarrow \quad \beta_2 \downarrow$

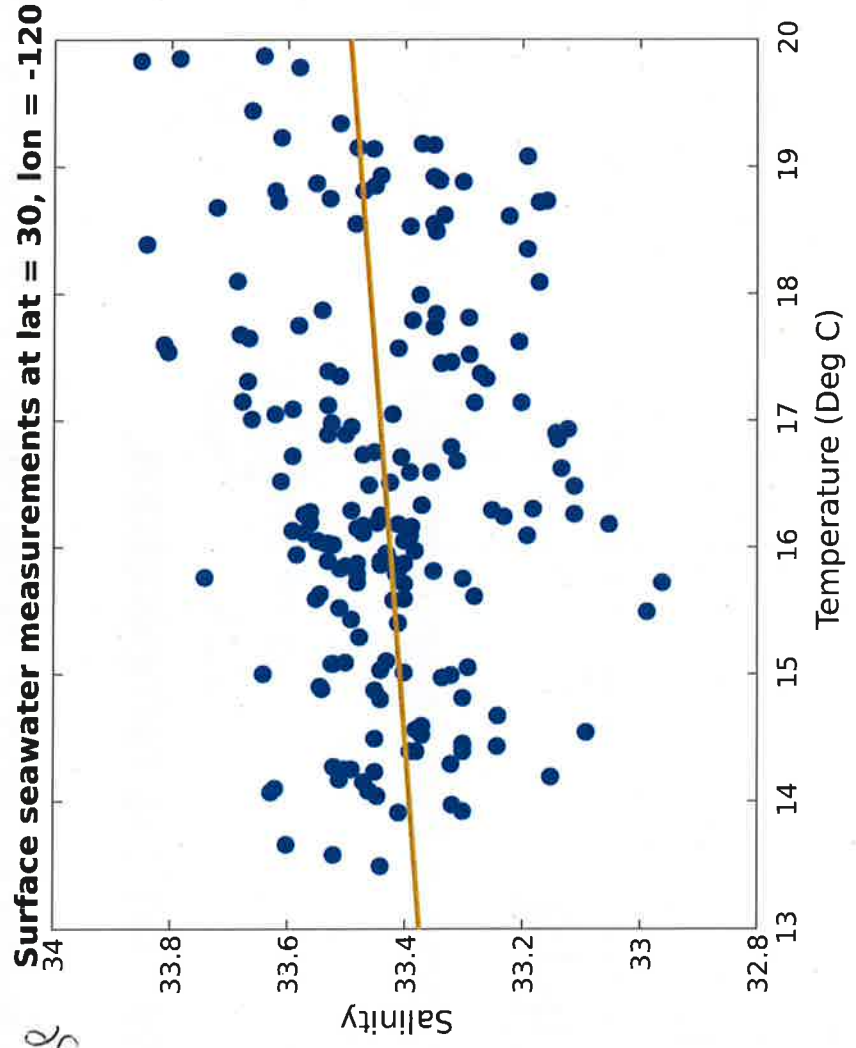$$s = 26.85 + 0.008x$$

$\leftarrow$ model

# Be careful with data fitting...



[XKCD #605]

# CalCOFI data

California Cooperative Oceanic Fisheries Investigations (CalCOFI) dataset is a database of oceanic measurements taken regularly since 1949. Amongst other aims, it collects data related to climate change.



**Surface seawater measurements at lat = 30, lon = -120**

# CalCOFI data

Can fit a linear function between temperature and salinity such that

$$s = mt + c$$

salinity, temp

$$e_i = s_i - mt_i - c$$

prediction, actual value

where s is salinity and t is temperature.

minimise $\sum_i \|e\|^2$ wrt $m, c$

Finding the constants $m$ and $c$ amounts to solving the normal equations

$$\frac{d\|e\|^2}{dc} = \sum -2(s_i - mt_i - c) = 0$$

$$\frac{d\|e\|^2}{dm} = \sum -2t(s_i - mt_i - c) = 0$$

Normal equations

$$nc + m\sum_i t_i = \sum s_i$$

$$c\sum_i t_i + m\sum_i t_i^2 = \sum_i s_i t_i$$

#samples

$$\begin{bmatrix} n & \sum_i t_i \\ \sum_i t_i & \sum_i t_i^2 \end{bmatrix} \begin{bmatrix} c \\ m \end{bmatrix} = \begin{bmatrix} \sum_i s_i \\ \sum_i s_i t_i \end{bmatrix}$$

In this case

$$\begin{bmatrix} 185 & 3048.1 \\ 3048.1 & 50676.9 \end{bmatrix} \begin{bmatrix} c \\ m \end{bmatrix} = \begin{bmatrix} 6184.4 \\ 101899.6 \end{bmatrix}$$

$$s = 0.0136t + 37.205$$

giving $c = 33.205$ and $m = 0.0136$

University of
BRISTOL

# Making inferences from the data

Might want to make inferences, e.g., what is the probability that the salinity will be over 33.5 given a temperature of 16°C?

Under what conditions can we use linear regression to answer questions like this?

*errors are normally distributed*

*we don't need to do this to do this*

Three conditions to check

* Normal residuals — *errors or normally distributed*
* Constant variability (homoscedasticity) *variability in the errors doesn't depend on the independent variable*
* Independence of observations

  *data points are independent*

  *e.g error in variables doesn't depend on the temperature*

University of
**BRISTOL**

# Normal residuals

The residuals are

$$e_i = y_i - \sum_j \beta_j x_{i,j}, \quad \text{e.g.,} \quad e_i = y_i - \beta_1 - \beta_2 x_i$$

*actual* *prediction*

$S = \beta_1 + \beta_2 X$

$e_i = S_i - \eta_i - c$

$y$ residual for every observation

i.e., the difference between the line of best fit and the actual observation.

Plot a histogram of the residuals.

distribution of residuals



residual plot

predts good to normal distribution
close to normal distribution

normal distribution
with the same
mean and std as
the observed residuals

flier

an impact on least squares
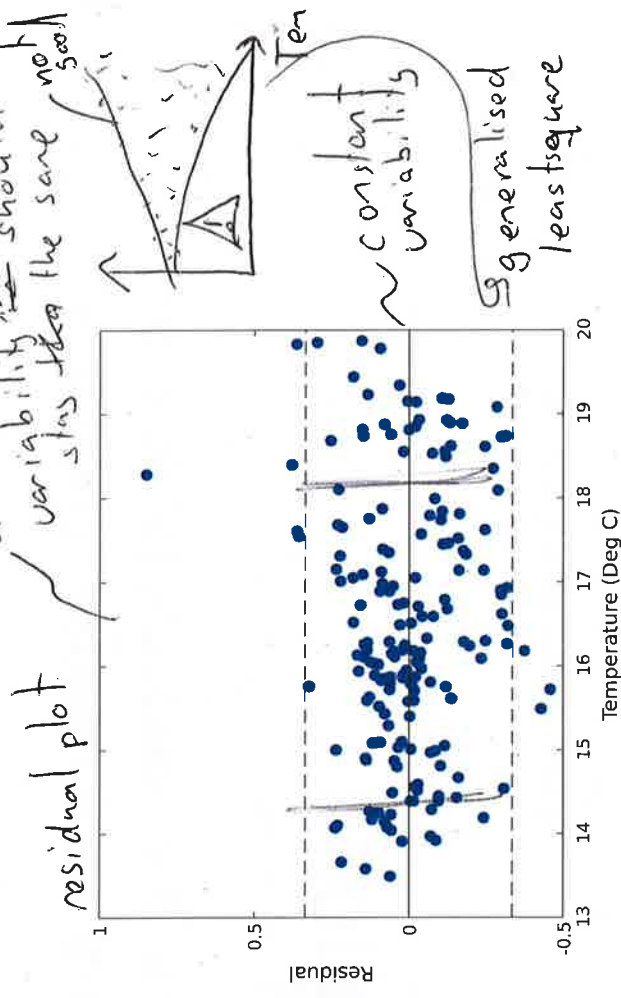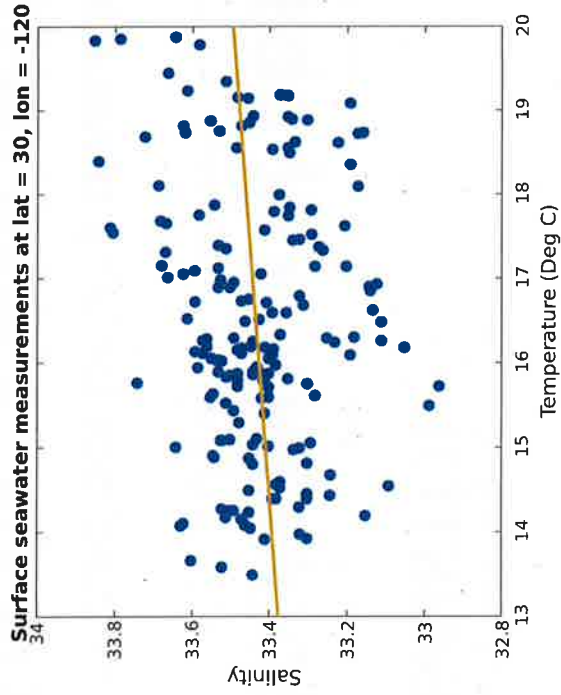
→ filter your data

→ generalised linear model
GLM

University of
BRISTOL

# Constant variability – homoscedasticity

There should not be any obvious trends in the residuals (should be normal random variables) and variability (e.g., standard deviation) should not change as the independent variable does.

Plot the residuals against the independent variable — this is the most common sort of residual plot.
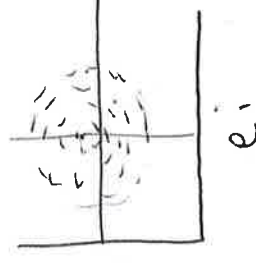


Surface seawater measurements at lat = 30, lon = -120

residual plot

as temperature increases variability is residual should stays the same not so

constant variability

generalised least square

Ten

# University of BRISTOL

# Independence of observations
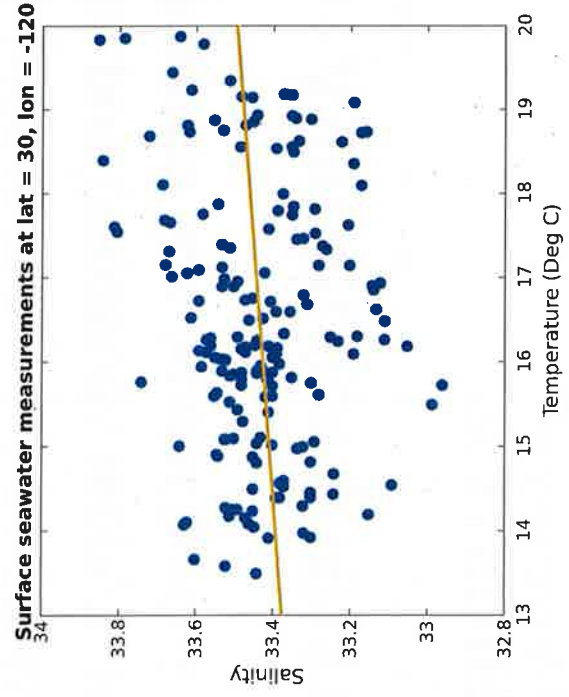
└ challenging when you have time series

There should not be any obvious trends in the residuals (should be normal random variables) and the value of each residual should not depend on another.

Plot the i-th residual against the i + 1-residual.

independent

$e_i$

not independent $e_{i+1}$

residual plot

△ linear regression

→ AR autoregressive model



Surface seawater measurements at lat = 30, lon = -120

University of
BRISTOL

# Least squares with arbitrary functions

Often data contains non-polynomial trends that we want to investigate

# Least squares with arbitrary functions

Take for example the function

$linear\ trend$

$seasonal\ trend$

$$y = \beta_1 + \beta_2 t + \beta_3 \sin(2\pi t)$$

Notice that the frequency is specified! Otherwise it would be a nonlinear regression problem — much harder and requires a numerical solution.

Define the error as

$residual$

$actual$

$prediction$

$$\|e\|^2 = \sum_{i=1}^{n} (y_i - \beta_1 - \beta_2 t_i - \beta_3 \sin(2\pi t_i))^2$$

and derive the normal equations in the usual way. (Find derivatives, set them equal to zero. . . )

University of
BRISTOL

# Least squares with arbitrary functions

The normal equations in this case are

$$\beta_1 n + \beta_2 \sum_{i=1}^{n} t_i + \beta_3 \sum_{i=1}^{n} S_i = \sum_{i=1}^{n} y_i \qquad \text{derivative of } \|e\|^2 \text{ w.r.t } \beta_1 = 0$$

$$\beta_1 \sum_{i=1}^{n} t_i + \beta_2 \sum_{i=1}^{n} t_i^2 + \beta_3 \sum_{i=1}^{n} t_i S_i = \sum_{i=1}^{n} t_i y_i \qquad \text{derivative of } \|e\|^2 \text{ w.r.t } \beta_2$$

$$\beta_1 \sum_{i=1}^{n} S_i + \beta_2 \sum_{i=1}^{n} t_i S_i + \beta_3 \sum_{i=1}^{n} S_i^2 = \sum_{i=1}^{n} S_i y_i \qquad \text{'' } \beta_3$$

where $S_i = \sin(2\pi t_i)$. (Notice the symmetry again!)

Again in the form $A\beta = b$.

$$A = \begin{bmatrix} n & \sum t_i & \sum S_i \\ \sum t_i & \sum t_i^2 & \sum t_i S_i \\ \sum S_i & \sum t_i S_i & \sum S_i^2 \end{bmatrix} \qquad \beta = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix} \qquad b = \begin{bmatrix} \sum y_i \\ \sum t_i y_i \\ \sum S_i y_i \end{bmatrix}$$

University of
BRISTOL

# General approach

Deriving the normal equations each time is straightforward but tedious. Is there a general equation? Yes!

General linear regression uses

$$y_i = \beta_1 x_{i,1} + \beta_2 x_{i,2} + \cdots + \beta_p x_{i,p}$$

the $x_{i,j}$ terms can be whatever we want! In the last example $p = 3$ and

$$x_{i,1} = 1$$

$$x_{i,2} = t_i$$

$$x_{i,3} = \sin(2\pi t_i)$$

to give

$$y_i = \beta_1 + \beta_2 t_i + \beta_3 \sin(2\pi t_i)$$

# General approach

Hence writing this in matrix-vector form gives

$$y = X\beta$$

where

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}, \quad X = \begin{bmatrix} 1 & t_1 & \sin(2\pi t_1) \\ \vdots & \vdots & \vdots \\ 1 & t_n & \sin(2\pi t_n) \end{bmatrix}, \quad \beta = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \end{bmatrix}$$

The corresponding error vector is given by

$$\|e\|^2 = \underbrace{(y - X\beta)^{\mathsf{T}}}_{\substack{\text{prediction} \\ \text{vector} \\ \text{actual values}}} \underbrace{(y - X\beta)}_{\substack{\text{vector} \\ \text{= scalar}}}$$

# General approach

Expanding out the error vector and differentiating gives the *general normal equations*

$$\underbrace{(X^{\mathsf{T}}X)}_{A}\underline{\beta} = X^{\mathsf{T}}y$$

i.e., the same $A\beta = b$ form as before, and so

$$\beta = \underbrace{(X^{\mathsf{T}}X)^{-1}}_{A^{-1}}X^{\mathsf{T}}y$$

In Matlab this is achieved very simply with the command

```
beta = X \ y;
```

For example with $y = \beta_1 + \beta_2 t + \beta_3 \sin(2\pi t)$ with $y$ and $t$ as column vectors we have

```
beta = [ones(size(t)), t, sin(2*pi*t)] \ y;
```

# Exercise

A particular process follows a daily cycle which suggests a least-squares fit using

$$y = \beta_1 \sin(2\pi t) + \beta_2 \cos(2\pi t)$$

where $t$ is measured in days, with the data

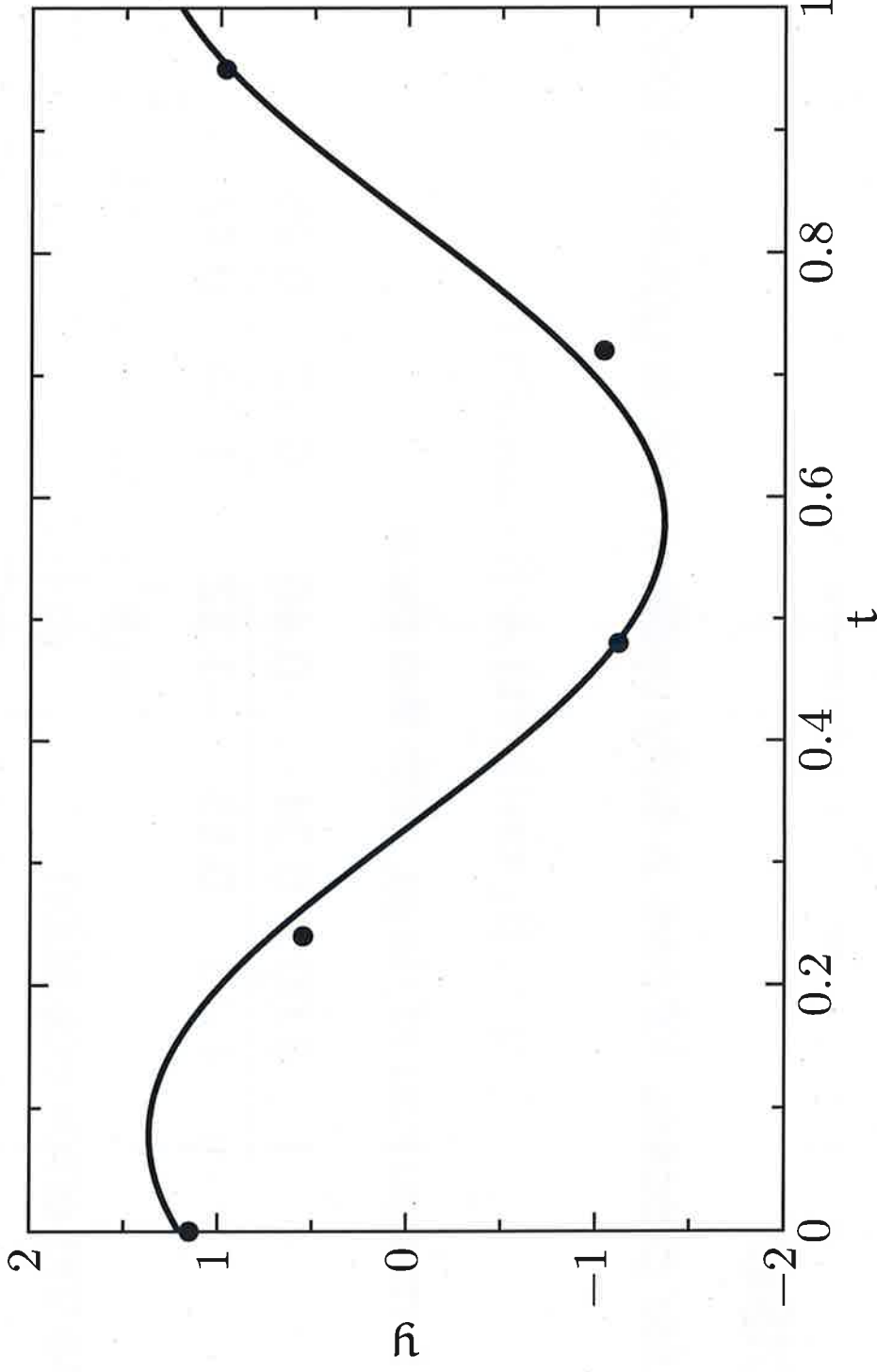| $t$ | 0.00 | 0.24 | 0.48 | 0.72 | 0.95 |
|---|---|---|---|---|---|
| $y$ | 1.15 | 0.55 | $-1.12$ | $-1.04$ | 0.97 |

$$\|e\|^2 = \sum \left( b_i - \beta_1 \sin(2\pi t_i) - \beta_2 \cos(2\pi t_i) \right)^2$$

$$X = \begin{bmatrix} \sin(2\pi t_*) & \cos(2\pi t_*) \\ \sin(2\pi t_*) & \cos(2\pi t_*) \end{bmatrix}$$

$$\|e\|^2 = (y - X\beta)^T (y - X\beta)$$

1. State the error equation

2. Derive the normal equations by differentiating w.r.t. $\beta_1$ and $\beta_2$

3. Calculate the required quantities

4. Solve the normal equations for $\beta_1$ and $\beta_2$

(Or use the general equation but that's probably harder by hand... )

University of
BRISTOL

# Answers



In this case, least-squares is a good alternative to an FFT (small number of data points and sampling freq is incommensurate with the period).

# † Uncertainty in least-squares estimates

When calculating least squares estimates for $\beta$ of the form

$$y = X\beta$$

it is important to remember that $\beta$ is a random variable with its own distribution.

The mean of the distribution is the value of $\beta$ calculated, but what is the variance? Multiple linked variables means we have to compute the *covariance matrix*.

With this information, we can say how confident we are about any estimates.

University of
BRISTOL

# † Uncertainty in least-squares estimates

The general solution for the least-squares estimator gives the estimated value for $\beta$

$$\hat{\beta} = (X^{\mathrm{T}}X)^{-1}X^{\mathrm{T}}y$$

From this expression it's possible to show that

$$\mathrm{covar}(\hat{\beta}) = (X^{\mathrm{T}}X)^{-1}\sigma^2$$

where $\sigma^2$ is the variance of the noise.

Can estimate the variance of the noise from the data

$$\hat{\sigma}^2 = \frac{1}{n-p}\|y - X\beta\|^2 = \frac{1}{n-p}\|e\|^2$$

where $p$ is the number of parameters to be estimated.

University of
**BRISTOL**

# † Uncertainty in least-squares estimates

For the previous exercise with

$$y = \beta_1 \sin(2\pi t) + \beta_2 \cos(2\pi t)$$

we have

$$X = \begin{bmatrix} 0 & 1.0000 \\ 0.9980 & 0.0628 \\ 0.1253 & -0.9921 \\ -0.9823 & -0.1874 \\ -0.3090 & 0.9511 \end{bmatrix}$$

$$\hat{\sigma}^2 = \frac{1}{n-p} \|y - X\beta\|^2 = 0.0216$$

Hence

$$\mathrm{covar}(\beta) = \begin{bmatrix} 0.0105 & 0.0006 \\ 0.0006 & 0.0074 \end{bmatrix}$$

University of
BRISTOL

# † Exercise

Given the covariance matrix for $\beta$ (from the previous exercise)

$$\text{covar}(\beta) = \begin{bmatrix} 0.0105 & 0.0006 \\ 0.0006 & 0.0074 \end{bmatrix}$$

and that $\beta_1 = 0.6451$, what is the range of maximum value that $\beta_1$ could take and still pass a hypothesis test to 5% significance?

Ignore the covariance in this case — just use the first element of the matrix as the variance of $\beta_1$.

A two-tailed t test with 3 degrees of freedom at 5% significance gives a critical value of 3.182.