In the frequentist approach, parameter estimation requires the definition of a lot of mathematical machinery

 \bullet $\ \mbox{Random}$ sample of size M , drawn from underlying pdf





In the frequentist approach, parameter estimation requires the definition of a lot of mathematical machinery

- \bullet $\ \mbox{Random}$ sample of size M , drawn from underlying pdf
- Sampling distribution, derived from underlying pdf (depends on underlying pdf, and on M)





In the frequentist approach, parameter estimation requires the definition of a lot of mathematical machinery

- \bullet Random sample of size M , drawn from underlying pdf
- Sampling distribution, derived from underlying pdf (depends on underlying pdf, and on M)
- Define an *estimator* function of sample used to estimate parameters of the pdf





In the frequentist approach, parameter estimation requires the definition of a lot of mathematical machinery

- \bullet Random sample of size M , drawn from underlying pdf
- Sampling distribution, derived from underlying pdf (depends on underlying pdf, and on M)
- Define an *estimator* function of sample used to estimate parameters of the pdf
- Hypothesis test to decide if estimator is 'acceptable', for the given sample size





In the frequentist approach, parameter estimation requires the definition of a lot of mathematical machinery

 \bullet Random sample of size M , drawn from underlying pdf

How do we decide what makes an 'acceptable' estimator?

estimate parameters of the pdf

 Hypothesis test – to decide if estimator is 'acceptable', for the given sample size





Compute sampling distribution for $\ \widehat{z}_1$ and $\ \widehat{z}_2$, modelling errors





Compute sampling distribution for $\ \widehat{z}_1$ and $\ \widehat{z}_2$, modelling errors





Compute sampling distribution for $\ \widehat{z}_1$ and $\ \widehat{z}_2$, modelling errors







Compute sampling distribution for $\ \widehat{z}_1$ and $\ \widehat{z}_2$, modelling errors



Better choice of estimator (if we can correct bias)





The Sample Mean

 ${x_1,...,x_M}$ = random sample from pdf p(x) with mean μ and variance σ^2

$$\widehat{\mu} = \frac{1}{M} \sum_{i=1}^{M} x_i$$
 = sample mean
Can show that $E(\widehat{\mu}) = \mu$ unbiased estimator

But bias is defined formally in terms of an infinite set of randomly chosen samples, each of size M.





The Sample Mean

 ${x_1,...,x_M}$ = random sample from pdf p(x) with mean μ and variance σ^2

$$\widehat{\mu} = \frac{1}{M} \sum_{i=1}^{M} x_i$$
 = sample mean
Can show that $E(\widehat{\mu}) = \mu$ unbiased estimator

But bias is defined formally in terms of an infinite set of randomly chosen samples, each of size M.

What can we say with a finite number of samples, each of finite size?





The Sample Mean

 $\{x_1, ..., x_M\}$ = random sample from pdf p(x) with mean μ and variance σ^2

$$\hat{\mu} = \frac{1}{M} \sum_{i=1}^{M} x_i = \text{sample mean}$$
Can show that
$$E(\hat{\mu}) = \mu \quad \text{unbiased estimator}$$
and
$$var[\hat{\mu}] = \frac{\sigma^2}{M} \quad \text{as sample size increases, sample mean increasingly concentrated near to true mean}$$





Linear correlation

Given sampled data $\{(x_i, y_i); i = 1, ..., n\}$ we can estimate the linear correlation between the variables as follows:







Linear correlation

Given sampled data $\{(x_i, y_i); i = 1, ..., n\}$ we can estimate the linear correlation between the variables as follows:



If p(x,y) is bivariate normal then r is an estimator of ho





The bivariate normal distribution

 $\mu_{y} + \frac{\sigma_{y}}{\sigma_{x}}\rho(x - \mu_{x})$ is often referred to as the **conditional expectation** (value) of y given x, and the equation

$$y = \mu_{\rm y} + \frac{\sigma_{\rm y}}{\sigma_{\rm x}}\rho(x - \mu_{\rm x})$$

is called the **regression line** of y on x.





Question 5: A correlation coefficient of r = 0.8 is calculated for a sample of paired data $\{(x, y)\}$ drawn from a bivariate normal distribution. Without being given any further information, which of the following statements can we say is correct?

A As the x values increase, the y values decrease

- **B** The data are scattered about a line of slope 0.8
- **C** The data are scattered about a line of unknown positive slope



Linear correlation

We can also rewrite the formula for γ in the slightly simpler forms:

$$r = \frac{\sum_{i=1}^{n} (x_i - \overline{x})(y_i - \overline{y})}{\sqrt{\sum_{i=1}^{n} (x_i - \overline{x})^2 \sum_{i=1}^{n} (y_i - \overline{y})^2}}$$

or

$$r = \frac{n \sum x_{i} y_{i} - \sum x_{i} \sum y_{i}}{\sqrt{n \sum x_{i}^{2} - (\sum x_{i})^{2}} \sqrt{n \sum y_{i}^{2} - (\sum y_{i})^{2}}}$$



SUPA)

Question 6: Estimate \mathcal{V} for the sample $\{(x, y)\}$ data shown in the graph below

$$y \uparrow \qquad \mathbf{A} \qquad r = 0$$

$$\mathbf{A} \qquad r = 0$$

$$\mathbf{B} \qquad r = 0.5$$

$$\mathbf{C} \qquad r = 1$$

D
$$r = -1$$

The Central Limit Theorem

For any pdf with finite variance σ^2 , as $M \to \infty$ $\hat{\mu}$ follows a normal pdf with mean μ and variance σ^2 / M



The Central Limit Theorem

For any pdf with finite variance σ^2 , as $M \to \infty$

 $\widehat{\mu}\,$ follows a normal pdf with mean $\,\mu\,$ and variance $\sigma^2\,/\,M\,$

Explains importance of normal pdf in statistics.

But still based on asymptotic behaviour of an infinite ensemble of samples that we didn't actually observe!





The Central Limit Theorem

For any pdf with finite variance σ^2 , as $M \to \infty$

 $\widehat{\mu}\,$ follows a normal pdf with mean $\,\mu\,$ and variance $\sigma^2\,/\,M\,$

Explains importance of normal pdf in statistics.

But still based on asymptotic behaviour of an infinite ensemble of samples that we didn't actually observe!

No 'hard and fast' rule for defining 'good' estimators. FPT invokes a number of principles – e.g. least squares, maximum likelihood



<u>Method of Least Squares</u>

- 'workhorse' method for fitting lines and curves to data in the physical sciences
- method often encountered (as a 'black box'?) in elementary courses
- useful demonstration of underlying statistical principles
- simple illustration of fitting straight line to (x,y) data





<u>Ordinary Linear Least Squares</u>

Suppose that the scatter in a plot of $\{x_i, y_i\}$ is assumed to arise from errors in only one of the two variables. This case is called **Ordinary Least Squares**. We then call x the **independent variable**, and y the **dependent variable**. Thus we suppose that we can write, for each data point:-

$$y_i = a + bx_i + \epsilon_i$$

where ϵ_i is known as the **residual** of the i^{th} data point – i.e. the difference between the observed value of y_i , and the value predicted by the best-fit straight line, characterised by parameters a and b.





Ordinary Linear Least Squares

We assume that the $\{\epsilon_i\}$ are an independently and identically distributed random sample from some underlying pdf with mean zero and variance σ^2 – i.e. the residuals are equally likely to be positive or negative and all have equal variance.

The least squares estimators of a and b minimise

$$S = \chi^2(a,b) = \sum_{i=1}^n [y_i - (a+bx_i)]^2$$

and \hat{a}_{LS} and \hat{b}_{LS} satisfy

$$\frac{\partial S}{\partial a} = 0 \quad \text{when} \quad a = \hat{a}_{\text{LS}} \qquad \frac{\partial S}{\partial b} = 0 \quad \text{when} \quad b = \hat{b}_{\text{LS}}$$





Ordinary Linear Least Squares

We assume that the $\{\epsilon_i\}$ are an independently and identically distributed random sample from some underlying pdf with mean zero and variance σ^2 – i.e. the residuals are equally likely to be positive or negative and all have equal variance. $S = \sum_{i=1}^{n} \varepsilon_i^2$

The least squares estimators of a and b minimise

$$S = \chi^2(a,b) = \sum_{i=1}^n [y_i - (a+bx_i)]^2$$

and \hat{a}_{LS} and \hat{b}_{LS} satisfy

$$\frac{\partial S}{\partial a} = 0$$
 when $a = \hat{a}_{\text{LS}}$ $\frac{\partial S}{\partial b} = 0$ when $b = \hat{b}_{\text{LS}}$





Solving these equations, \hat{a}_{LS} and \hat{b}_{LS} are given by

$$\hat{a}_{\text{LS}} = \frac{\sum y_i \sum x_i^2 - \sum y_i x_i \sum x_i}{n \sum x_i^2 - (\sum x_i)^2}$$

$$\hat{b}_{\text{LS}} = \frac{n \sum y_i x_i - \sum y_i \sum x_i}{n \sum x_i^2 - (\sum x_i)^2}$$

where n denotes the sample size and all sums are for i = 1, ..., n.



SUPA)

We can show that

$$E(\hat{a}_{LS}) = a_{LS}$$
i.e. LS estimators are *unbiased*.

$$E(\hat{b}_{LS}) = b_{LS}$$
Also
$$var(\hat{a}_{LS}) = \frac{\sigma^2 \sum x_i^2}{n \sum x_i^2 - (\sum x_i)^2}$$

$$var(\hat{b}_{LS}) = \frac{\sigma^2 n}{n \sum x_i^2 - (\sum x_i)^2}$$
and
$$cov(\hat{a}_{LS}, \hat{b}_{LS}) = \frac{-\sigma^2 \sum x_i}{n \sum x_i^2 - (\sum x_i)^2}$$
Jniversity





We can show that
$$\begin{array}{l} E\left(\hat{a}_{LS}\right) = a_{LS} \\ E\left(\hat{b}_{LS}\right) = b_{LS} \end{array} \quad \text{i.e. LS estimators are unbiased.} \\ \end{array}$$
Also
$$\begin{array}{l} \operatorname{Var}(\hat{a}_{\mathrm{LS}}) &= & \frac{\sigma^2 \sum x_i^2}{n \sum x_i^2 - (\sum x_i)^2} \\ \operatorname{Var}(\hat{b}_{\mathrm{LS}}) &= & \frac{\sigma^2 n}{n \sum x_i^2 - (\sum x_i)^2} \end{array}$$
and
$$\begin{array}{l} \operatorname{cov}(\hat{a}_{\mathrm{LS}}, \hat{b}_{\mathrm{LS}}) &= & \frac{-\sigma^2 \sum x_i}{n \sum x_i^2 - (\sum x_i)^2} \end{array}$$
Choosing the $\{x_i\}$ so that $\sum x_i = 0$ we can make \hat{a}_{LS} and \hat{b}_{LS} independent.

<u>Weighted Linear Least Squares</u>

Suppose the i^{th} residual, $\{\epsilon_i\}$, is assumed to be drawn from some underlying pdf with mean zero and variance σ_i^2 , where the variance is allowed to be different for each residual.

Define

$$S = \chi^{2}(a, b) = \sum_{i=1}^{n} \left[\frac{y_{i} - (a + bx_{i})}{\sigma_{i}} \right]^{2}$$

Again we find Least Squares estimators of a and b satisfying

$$\frac{\partial S}{\partial a} = 0 \qquad \qquad \frac{\partial S}{\partial b} = 0$$



SUPA)

Solving, we find

$$\hat{a}_{\text{WLS}} = \frac{\sum \frac{y_i}{\sigma_i^2} \sum \frac{x_i^2}{\sigma_i^2} - \sum \frac{y_i x_i}{\sigma_i^2} \sum \frac{x_i}{\sigma_i^2}}{\sum \frac{1}{\sigma_i^2} \sum \frac{x_i^2}{\sigma_i^2} - \left(\sum \frac{x_i}{\sigma_i^2}\right)^2}$$

$$\hat{b}_{\text{WLS}} = \frac{\sum \frac{1}{\sigma_i^2} \sum \frac{y_i x_i}{\sigma_i^2} - \sum \frac{y_i}{\sigma_i^2} \sum \frac{x_i}{\sigma_i^2}}{\sum \frac{1}{\sigma_i^2} \sum \frac{x_i^2}{\sigma_i^2} - \left(\sum \frac{x_i}{\sigma_i^2}\right)^2}$$





Also

$$\operatorname{var}(\hat{a}_{\mathrm{WLS}}) = \frac{\sum \frac{x_i^2}{\sigma_i^2}}{\sum \frac{1}{\sigma_i^2} \sum \frac{x_i^2}{\sigma_i^2} - \left(\sum \frac{x_i}{\sigma_i^2}\right)^2}$$







Also

$$\operatorname{var}(\hat{a}_{\mathrm{WLS}}) = \frac{\sum \frac{x_i^2}{\sigma_i^2}}{\sum \frac{1}{\sigma_i^2} \sum \frac{x_i^2}{\sigma_i^2} - \left(\sum \frac{x_i}{\sigma_i^2}\right)^2}$$



In the case where σ_i^2 is constant, for all *i*, these formulae reduce to those for the unweighted case.

• Errors on both variables?

Need to modify merit function accordingly.

$$\chi^{2}(a,b) = \sum_{i=1}^{N} \frac{(y_{i} - a - bx_{i})^{2}}{\sigma_{y_{i}}^{2} + b^{2}\sigma_{x_{i}}^{2}}$$

Renders equations *non-linear*; no simple analytic solution!

See e.g. Numerical Recipes 15.3





• General linear models?

e.g.
$$y(x) = a_1 + a_2 x + a_3 x^2 + \dots + a_M x^{M-1}$$

We have

$$\chi^2 = \sum_{i=1}^{N} \left[\frac{y_i - \sum_{k=1}^{M} a_k X_k(x_i)}{\sigma_i} \right]^2$$

Can formulate as a matrix equation and solve for parameters

See e.g. Numerical Recipes 15.4





Define
$$a = \begin{bmatrix} a_1 \\ \vdots \\ a_M \end{bmatrix}$$
 Vector of model
parameters $y = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix}$ Vector of observations
 $X = \begin{bmatrix} X_1(x_1) \cdots X_1(x_M) \\ \vdots & \vdots \\ X_N(x_1) \cdots X_N(x_M) \end{bmatrix}$ Matrix of model
basis functions







$$\boldsymbol{\varepsilon} = \begin{bmatrix} \boldsymbol{\varepsilon}_1 \\ \vdots \\ \boldsymbol{\varepsilon}_N \end{bmatrix}$$

where we assume \mathcal{E}_i is drawn from some pdf with mean zero and variance σ^2





Weighting by errors

Define
$$\boldsymbol{a} = \begin{bmatrix} a_1 \\ \vdots \\ a_M \end{bmatrix}$$
 $\boldsymbol{b} = \begin{bmatrix} y_1/\sigma_1 \\ \vdots \\ y_N/\sigma_N \end{bmatrix}$ Vector of model parameters Vector of weighted observations $\boldsymbol{A} = \begin{bmatrix} \frac{X_1(x_1)}{\sigma_1} \cdots \frac{X_1(x_M)}{\sigma_1} \\ \vdots & \vdots \\ \frac{X_N(x_1)}{\sigma_N} \cdots \frac{X_N(x_M)}{\sigma_N} \end{bmatrix}$ Design matrix



SUPA)







We solve for the parameter vector \hat{a}_{LS} that minimises

$$\boldsymbol{S} = \boldsymbol{e}^T \cdot \boldsymbol{e} = \sum_{i=1}^n e_i^2$$

This has solution

solution
$$\hat{a}_{LS} = (A^T A)^{-1} A^T \cdot b$$

 $M \times M$ matrix

and
$$\operatorname{cov}(\hat{a}_{LS}) = (A^T A)^{-1}$$



SUPA)

Inverting $(A^T A)$ can be hazardous, particularly if A is a **sparse** matrix and/or close to singular.

Some inversion methods will break down, since they may give a formal solution, but are highly unstable to round-off error in the data.

Remedy: solution via Singular Value Decomposition.

From linear algebra theory:

Any $N \times M$ matrix can be decomposed as the product of an $N \times M$ column-orthogonal matrix U, an $M \times M$ diagonal matrix W with positive or zero elements (the *singular* values) and the transpose of an $M \times M$ orthogonal matrix V





From linear algebra theory:

Any $N \times M$ matrix can be decomposed as the product of an $N \times M$ column-orthogonal matrix U, an $M \times M$ diagonal matrix W with positive or zero elements (the *singular* values) and the transpose of an $M \times M$ orthogonal matrix V



Let the vectors $U_{(i)}$ $i = 1, \ldots, M$ denote the columns of U (each one is a vector of length N)

Let the vectors $\mathbf{V}_{(i)}; i = 1, \dots, M$ denote the columns of \mathbf{V} (each one is a vector of length M)

It can be shown that the solution to the general linear model satisfies

$$\hat{a}_{LS} = \sum_{i=1}^{M} \left(\frac{\mathbf{U}_{(i)} \cdot \mathbf{b}}{w_i} \right) \mathbf{V}_{(i)}$$





$$\hat{\pmb{a}}_{\pmb{LS}} = \sum_{i=1}^{M} \left(\frac{\mathbf{U}_{(i)} \cdot \mathbf{b}}{w_i} \right) \mathbf{V}_{(i)}$$

Very small values of w_i will amplify any round-off errors in \mathbf{b}

Solution: For these very small singular values, set $\frac{1}{w_i} = 0$.

This suppresses their noisy contribution to the least-squares solution for the parameters \hat{a}_{LS} .

SVD acts as a noise filter – see also Section 7

o Non-linear models?
$$y_i^{\text{model}} \equiv y^{\text{model}}(x_i; \theta_1, ..., \theta_k)$$
Model parameters

Suppose $y_i^{\text{obs}} = y_i^{\text{model}} + \epsilon_i$

 \mathcal{E}_i drawn from pdf with mean zero, variance σ_i^2

Then

$$S = \chi^2 = \sum_{i=1}^n \left[\frac{y_i^{\text{obs}} - y_i^{\text{model}}}{\sigma_i} \right]^2$$



SUPA)

• Non-linear models?
$$y_i^{\text{model}} \equiv y^{\text{model}}(x_i; \theta_1, ..., \theta_k)$$
Model parameters

Suppose
$$y_i^{\text{obs}} = y_i^{\text{model}} + \epsilon_i$$

 \mathcal{E}_i drawn from pdf with mean zero, variance σ_i^2

Then

$$S = \chi^2 = \sum_{i=1}^n \left[\frac{y_i^{\text{obs}} - y_i^{\text{model}}}{\sigma_i} \right]^2$$

But no simple analytic method to minimise sum of squares (e.g. no analytic solutions to $\partial S/\partial \theta_i = 0$)

• Non-linear models?

Methods of solution often involve assuming Taylor expansion of χ^2 around minimum, and solving by gradient descent



See e.g. Numerical Recipes 15.5 and Section 6





• Correlated errors?

We need to define a covariance matrix $C_{ii} = cov(x_i, x_i)$

$$\chi^{2} = \sum_{i} \sum_{j} \left(y_{i} - y_{i}^{\text{model}} \right) \left[C_{ij} \right]^{-1} \left(y_{j} - y_{j}^{\text{model}} \right)$$

See e.g. Gregory, Chapter 10





[Note: not-examinable, and not included in video and audio files]

• More general non-linear approaches?

Suppose we want to derive a function y = f(x) with errors from observed data \mathcal{D} .

We might want to use this function to interpolate or extrapolate from our data.



A Gaussian process (GP) defines a distribution over functions p(f):

$$p(f|\mathcal{D}) = \frac{p(f)p(\mathcal{D}|f)}{p(\mathcal{D})}$$





[Note: not-examinable, and not included in video and audio files]

GPs are parametrized by a mean function $\mu(x)$ and a kernel function K(x, x')

$$p(f(x), f(x')) = \mathsf{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$
$$\boldsymbol{\mu} = \begin{bmatrix} \mu(x) \\ \mu(x') \end{bmatrix} \boldsymbol{\Sigma} = \begin{bmatrix} K(x, x) & K(x, x') \\ K(x', x) & K(x', x') \end{bmatrix}$$

We can learn about these functions from the data themselves → strong connections to machine learning and neural networks





[Note: not-examinable, and not included in video and audio files]

Vast, and rapidly growing, literature on GPs across the physical sciences.

Beyond the scope of ADA Course but very interesting!

See e.g. Sivia Chapter 6 for some introductory ideas, or some nice recent articles / lecture notes:

https://arxiv.org/pdf/1505.02965v2.pdf

http://www.cs.toronto.edu/~hinton/csc2515/notes/gp_slides_fall08.pdf





[Note: not-examinable, and not included in video and audio files]

Some really excellent books available (all free, downloadable!)



http://www.inference.phy.cam.ac.uk/itprnn/book.pdf

Gaussian Processes for Machine Learning



Carl Edward Rasmussen and Christopher K. I. Williams



http://www.gaussianprocess.org/gpml/ http://web4.cs.ucl.ac.uk

http://web4.cs.ucl.ac.uk/staff/D.Barber/textbook/090310.pdf



