# 6. An Advanced Bayesian Toolbox – Part One

# Course Programme

**Lectures 6 to 10**

6.  **An Advanced Toolbox for Bayesian Inference**

7.  **An Advanced Toolbox for Bayesian Inference**

8.  **Bayesian Model Selection**

9.  **Monte Carlo Simulation Methods**

10. **Fourier Methods**

# Parameter estimation:

Posterior  Likelihood  Prior

$$p\,(\,\text{model}\;\,|\,\text{data},\;I\,)\quad\propto\quad p\,(\,\text{data}\;\,|\,\text{model}\;,\,I\,)\times p\,(\,\text{model}\;\,|\,I\,)$$

# Parameter estimation: the Gaussian approximation

$$p(\theta \mid \text{data}, I) \quad \propto \quad p(\text{data} \mid \theta, I) \times p(\theta \mid I)$$

'Best' estimator: $\left. \dfrac{\partial p(\theta \mid \text{data}, I)}{\partial \theta} \right|_{\theta=\theta_0} = 0$  $\longleftarrow$  Maximise posterior

Equivalently, we can define $\ell = \log p(\theta \mid \text{data}, I)$ and compute $\left. \dfrac{\partial \ell}{\partial \theta} \right|_{\theta=\theta_0} = 0$

University of Glasgow

Advanced Data Analysis Course, 2019-20

SUPA

# Parameter estimation: the Gaussian approximation

$$p(\theta \mid \text{data}, \ I) \quad \propto \quad p(\text{data} \mid \theta, I) \times p(\theta \mid I)$$

'Best' estimator: $\dfrac{\partial p(\theta \mid \text{data}, I)}{\partial \theta}\bigg|_{\theta = \theta_0} = 0$     ⟵   <span style="color:red">Maximise posterior</span>

Equivalently, we can define $\ell = \log p(\theta \mid \text{data}, I)$ and compute $\dfrac{\partial \ell}{\partial \theta}\bigg|_{\theta = \theta_0} = 0$

Taylor expand $\ell(\theta)$ around $\theta = \theta_0$ :

$$\ell(\theta) \ = \ \ell(\theta_0) + \frac{\partial \ell}{\partial \theta}\bigg|_{\theta = \theta_0} (\theta - \theta_0) + \frac{1}{2}\frac{\partial^2 \ell}{\partial \theta^2}\bigg|_{\theta = \theta_0} (\theta - \theta_0)^2 + \dots$$

# Parameter estimation: the Gaussian approximation

$$p(\theta \mid \text{data}, \ I) \quad \propto \quad p(\text{data} \mid \theta, I) \times p(\theta \mid I)$$

'Best' estimator: $\left.\dfrac{\partial p(\theta \mid \text{data}, I)}{\partial \theta}\right|_{\theta=\theta_0} = 0$ ← Maximise posterior

Equivalently, we can define $\ell = \log p(\theta \mid \text{data}, I)$ and compute $\left.\dfrac{\partial \ell}{\partial \theta}\right|_{\theta=\theta_0} = 0$

Taylor expand $\ell(\theta)$ around $\theta=\theta_0$ :

$$\ell(\theta) \ = \ \ell(\theta_0) + \left.\frac{\partial \ell}{\partial \theta}\right|_{\theta=\theta_0}(\theta-\theta_0) + \frac{1}{2}\left.\frac{\partial^2 \ell}{\partial \theta^2}\right|_{\theta=\theta_0}(\theta-\theta_0)^2 + \dots$$

# Parameter estimation: the Gaussian approximation

$$p(\theta \mid \text{data}, \ I) = \exp\left[\ell(\theta)\right]$$

Neglecting higher order terms in $\ell(\theta)$

$$p(\theta \mid \text{data}, \ I) \propto \exp\left(-\frac{A}{2}(\theta - \theta_0)^2\right)$$

where $A = -\left.\dfrac{\partial^2 \ell}{\partial \theta^2}\right|_{\theta=\theta_0}$

This is equivalent to a normal distribution, with $\sigma^{-2} = A = -\left.\dfrac{\partial^2 \ell}{\partial \theta^2}\right|_{\theta=\theta_0}$

## Parameter estimation: the Gaussian approximation

$$p(\theta \mid \text{data},\ I) = \exp\left[\ell(\theta)\right]$$

Neglecting higher order terms in $\ell(\theta)$ ← Gaussian approximation

$$p(\theta \mid \text{data},\ I) \propto \exp\left(-\frac{A}{2}(\theta - \theta_0)^2\right)$$

where $A = -\left.\dfrac{\partial^2 \ell}{\partial \theta^2}\right|_{\theta=\theta_0}$

This is equivalent to a normal distribution, with $\sigma^{-2} = A = -\left.\dfrac{\partial^2 \ell}{\partial \theta^2}\right|_{\theta=\theta_0}$

## Parameter estimation: the Gaussian approximation

$$p(\theta \mid \text{data}, I) = \exp\left[\ell(\theta)\right]$$

Neglecting higher order terms in $\ell(\theta)$ $\longleftarrow$ Gaussian approximation

$$p(\theta \mid \text{data}, I) \propto \exp\left(-\frac{A}{2}(\theta - \theta_0)^2\right)$$

where $A = -\dfrac{\partial^2 \ell}{\partial \theta^2}\bigg|_{\theta=\theta_0}$

This is equivalent to a **normal** distribution, with $\sigma^{-2} = A = -\dfrac{\partial^2 \ell}{\partial \theta^2}\bigg|_{\theta=\theta_0}$

Can summarise inference from posterior by

$$\boxed{\theta = \theta_0 \pm \sigma}$$

University of Glasgow

VIA VERITAS VITA

SUPA

Advanced Data Analysis Course, 2019-20

**Question 13:** Neglecting the higher order terms in the log posterior expansion produces a posterior which can be written as a normal pdf because

**A** The higher order moments of a Gaussian are all zero

**B** The Gaussian pdf is uniquely specified by its mean and variance

**C** The logarithm of a Gaussian pdf can be written in the form of a quadratic

**D** All of the above

# Parameter estimation: 2-D case

Recall our definition of *variance*

$$\mathrm{var}[x] \;=\; \int\limits_{-\infty}^{\infty} \left(x - \langle x \rangle\right)^2 p(x \mid I)\, dx$$

Extends to 2 variables – *covariance*

$$\mathrm{cov}[x, y] \;=\; \int\limits_{-\infty}^{\infty}\int\limits_{-\infty}^{\infty} \left(x - \langle x \rangle\right)\!\left(y - \langle y \rangle\right) p(x, y \mid I)\, dx\, dy$$

# Parameter estimation: 2-D case

Recall our definition of *variance*

$$\text{var}[x] = \int_{-\infty}^{\infty} \left(x - \langle x \rangle\right)^2 p(x \mid I) dx$$

Extends to 2 variables – *covariance*

$$\text{cov}[x, y] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \left(x - \langle x \rangle\right)\left(y - \langle y \rangle\right) p(x, y \mid I) dx dy$$

If $x$ and $y$ are independent, $\text{cov}[x, y] = 0$

This is because $p(x, y \mid I) = p(x \mid I) p(y \mid I)$

University of Glasgow
VIA VERITAS VITA

Advanced Data Analysis Course, 2019-20

SUPA

# Parameter estimation: 2-D case

$$p(\theta_1, \theta_2 \mid \text{data}, \ I) \quad \propto \quad p(\text{data} \mid \theta_1, \theta_2, I) \times p(\theta_1, \theta_2 \mid I)$$

'Best' estimator: $\left. \dfrac{\partial p(\theta_1, \theta_2 \mid \text{data}, I)}{\partial \theta_j} \right|_{\theta_j = \theta_{0j}} = 0$

Compute $\left. \dfrac{\partial \ell}{\partial \theta_j} \right|_{\theta_j = \theta_{0j}} = 0$ where $\ell = \log p(\theta_1, \theta_2 \mid \text{data}, I)$

# Parameter estimation: 2-D case

$$p(\theta_1, \theta_2 \mid \text{data}, \ I) \quad \propto \quad p(\text{data} \mid \theta_1, \theta_2, I) \times p(\theta_1, \theta_2 \mid I)$$

'Best' estimator: $\quad \dfrac{\partial p(\theta_1, \theta_2 \mid \text{data}, I)}{\partial \theta_j}\bigg|_{\theta_j = \theta_{0j}} = 0$

Compute $\quad \dfrac{\partial \ell}{\partial \theta_j}\bigg|_{\theta_j = \theta_{0j}} = 0 \quad$ where $\quad \ell = \log p(\theta_1, \theta_2 \mid \text{data}, I)$

Taylor expand $\ell(\theta_1, \theta_2)$ around $\theta_{0j}$ :

University of Glasgow

SUPA

# Parameter estimation: 2-D case

Taylor expand $\ell(\theta_1, \theta_2)$ around $\theta_{0j}$ :

$$\ell(\theta_1, \theta_2) = \ell(\theta_{01}, \theta_{02}) + \left.\frac{\partial \ell}{\partial \theta_1}\right|_{\theta_j = \theta_{0j}} (\theta_1 - \theta_{01}) + \left.\frac{\partial \ell}{\partial \theta_2}\right|_{\theta_j = \theta_{0j}} (\theta_2 - \theta_{02}) +$$

$$\frac{1}{2}\left[ \left.\frac{\partial^2 \ell}{\partial \theta_1^2}\right|_{\theta_j = \theta_{0j}} (\theta_1 - \theta_{01})^2 + \left.\frac{\partial^2 \ell}{\partial \theta_2^2}\right|_{\theta_j = \theta_{0j}} (\theta_2 - \theta_{02})^2 + 2\left.\frac{\partial^2 \ell}{\partial \theta_1 \partial \theta_2}\right|_{\theta_j = \theta_{0j}} (\theta_1 - \theta_{01})(\theta_2 - \theta_{02}) \right] + \ldots$$

# Parameter estimation: 2-D case

Taylor expand $\ell(\theta_1, \theta_2)$ around $\theta_{0j}$ :

$$
\ell(\theta_1, \theta_2) = \ell(\theta_{01}, \theta_{02}) + \left.\frac{\partial \ell}{\partial \theta_1}\right|_{\theta_j = \theta_{0j}} (\theta_1 - \theta_{01}) + \left.\frac{\partial \ell}{\partial \theta_2}\right|_{\theta_j = \theta_{0j}} (\theta_2 - \theta_{02}) +
$$

$$
\frac{1}{2}\left[ \left.\frac{\partial^2 \ell}{\partial \theta_1^2}\right|_{\theta_j = \theta_{0j}} (\theta_1 - \theta_{01})^2 + \left.\frac{\partial^2 \ell}{\partial \theta_2^2}\right|_{\theta_j = \theta_{0j}} (\theta_2 - \theta_{02})^2 + 2\left.\frac{\partial^2 \ell}{\partial \theta_1 \partial \theta_2}\right|_{\theta_j = \theta_{0j}} (\theta_1 - \theta_{01})(\theta_2 - \theta_{02}) \right] + \ldots
$$

$$
p(\theta_1, \theta_2 \mid \text{data}, I) \quad \propto \quad \exp\left[\ell(\theta_1, \theta_2)\right]
$$

$$
\propto \quad \exp\left[-\tfrac{1}{2} Q\right] \quad \longleftarrow \quad \text{Gaussian approximation}
$$

$$
\chi^2
$$

University of Glasgow

SUPA

# Parameter estimation: 2-D case

Taylor expand $\ell(\theta_1, \theta_2)$ around $\theta_{0j}$ :

$$\ell(\theta_1, \theta_2) = \ell(\theta_{01}, \theta_{02}) + \left.\frac{\partial \ell}{\partial \theta_1}\right|_{\theta_j = \theta_{0j}} (\theta_1 - \theta_{01}) + \left.\frac{\partial \ell}{\partial \theta_2}\right|_{\theta_j = \theta_{0j}} (\theta_2 - \theta_{02}) +$$

$$\frac{1}{2}\left[ \left.\frac{\partial^2 \ell}{\partial \theta_1^2}\right|_{\theta_j = \theta_{0j}} (\theta_1 - \theta_{01})^2 + \left.\frac{\partial^2 \ell}{\partial \theta_2^2}\right|_{\theta_j = \theta_{0j}} (\theta_2 - \theta_{02})^2 + 2\left.\frac{\partial^2 \ell}{\partial \theta_1 \partial \theta_2}\right|_{\theta_j = \theta_{0j}} (\theta_1 - \theta_{01})(\theta_2 - \theta_{02}) \right] + \ldots$$

$$p(\theta_1, \theta_2 \mid \text{data}, \ I) \ \propto \ \exp\left[\ell(\theta_1, \theta_2)\right]$$

$$\propto \ \exp\left[-\tfrac{1}{2} Q\right] \ \longleftarrow \ \text{Gaussian approximation}$$

$\chi^2$

Maximising posterior
$\equiv$ Minimising $\chi^2$

Advanced Data Analysis Course, 2019-20

SUPA

# Parameter estimation: 2-D case

Taylor expand $\ell(\theta_1, \theta_2)$ around $\theta_{0j}$ :

$$Q = \begin{pmatrix} \theta_1 - \theta_{10} & \theta_2 - \theta_{20} \end{pmatrix} \begin{bmatrix} A & C \\ C & B \end{bmatrix} \begin{pmatrix} \theta_1 - \theta_{10} \\ \theta_2 - \theta_{20} \end{pmatrix}$$

<span style="color:red">Quadratic form</span>

where $\quad A = \dfrac{\partial^2 \ell}{\partial \theta_1^{\,2}}\bigg|_{\theta_j = \theta_{0j}} \qquad B = \dfrac{\partial^2 \ell}{\partial \theta_2^{\,2}}\bigg|_{\theta_j = \theta_{0j}} \qquad C = \dfrac{\partial^2 \ell}{\partial \theta_1 \partial \theta_2}\bigg|_{\theta_j = \theta_{0j}}$

# Parameter estimation: 2-D case

Taylor expand $\ell(\theta_1, \theta_2)$ around $\theta_{0j}$ :

$$Q = \begin{pmatrix} \theta_1 - \theta_{10} & \theta_2 - \theta_{20} \end{pmatrix} \begin{bmatrix} A & C \\ C & B \end{bmatrix} \begin{pmatrix} \theta_1 - \theta_{10} \\ \theta_2 - \theta_{20} \end{pmatrix}$$

<span style="color:red">Quadratic form</span>

where $A = \left. \dfrac{\partial^2 \ell}{\partial \theta_1^2} \right|_{\theta_j = \theta_{0j}}$     $B = \left. \dfrac{\partial^2 \ell}{\partial \theta_2^2} \right|_{\theta_j = \theta_{0j}}$     $C = \left. \dfrac{\partial^2 \ell}{\partial \theta_1 \partial \theta_2} \right|_{\theta_j = \theta_{0j}}$

This is a **bivariate normal distribution with covariance matrix**

$$\sigma_{ij}^2 = \text{cov}_{ij} = \left\langle (\theta_i - \theta_{i0})(\theta_j - \theta_{j0}) \right\rangle = \left[ -\frac{\partial^2 \ell}{\partial \theta_i \partial \theta_j} \right]^{-1}$$

<span style="color:red">Fisher information matrix</span>

$$\mathbf{F} \equiv F_{ij} = \frac{\partial^2 \ell}{\partial \theta_i \partial \theta_j}$$ is known as the **Fisher information matrix**
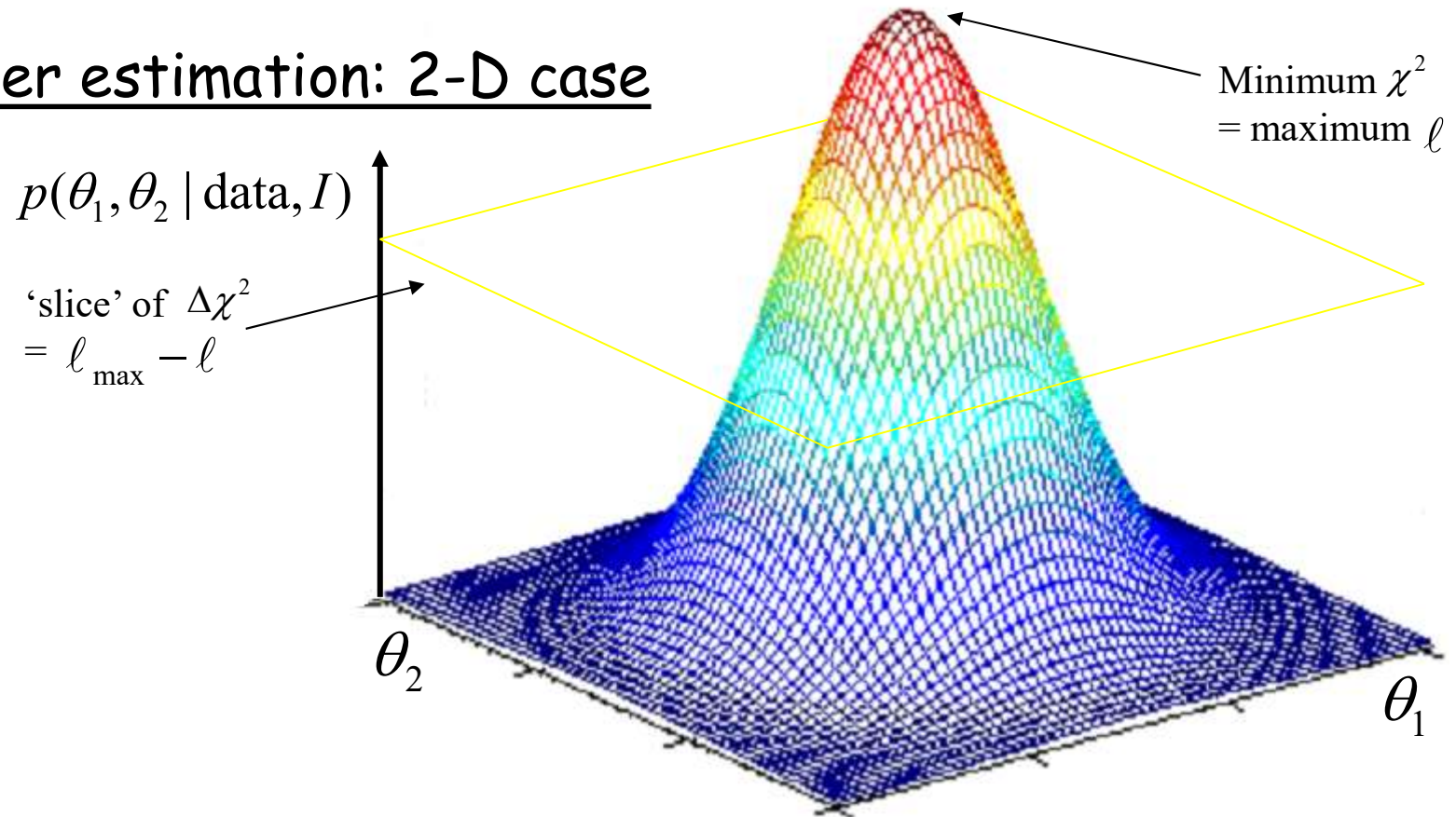
It provides a measure of how much information a given dataset can yield about the parameters of a model.

We can see this most easily in the case where the Fisher matrix is **diagonal**.

Then $$\mathbf{F} = -\mathrm{diag}\left(\sigma_1^{-2}, \ldots, \sigma_n^{-2}\right)$$

If the $i^{th}$ element of the Fisher matrix is large (negative), the **variance** of parameter $\theta_i$ is small (and positive).

In general the Fisher matrix (and covariance matrix) will **not** be diagonal; the Fisher matrix then tells us which **combinations** of the parameters are well constrained by the data. (see later).

So if, for our model:

o   the likelihood is Gaussian in shape (or if we can approximate it as Gaussian – i.e. if the higher order terms in the Taylor expansion of the log likelihood can be neglected);

o   the parameters have broad, uniform priors;

then the posterior will also be Gaussian.

If we can evaluate the first and second partial derivatives of the log likelihood, we can:

o   compute the **Fisher Information Matrix;**
o   compute the **Covariance Matrix** of the posterior.

We can also compute **credible regions** for the parameters (in fact for this we don't need the derivatives – see Section 9 )

We can write the log posterior as

$$\ell(\theta_1, \theta_2) = \text{const} - \tfrac{1}{2}\chi^2(\theta_1, \theta_2)$$

Now $\chi^2 = \chi^2_{\min}$ when $(\theta_1, \theta_2) = (\theta_{01}, \theta_{02})$

Maximising posterior $\equiv$ Minimising $\chi^2$

so we can write, for $\Delta\chi^2(\theta_1, \theta_2) = \chi^2(\theta_1, \theta_2) - \chi^2_{\min}$

$$\ell(\theta_1, \theta_2) = \ell(\theta_{01}, \theta_{02}) - \tfrac{1}{2}\Delta\chi^2(\theta_1, \theta_2)$$

So that

$$p(\theta_1, \theta_2 \mid \text{data}, I) = \underbrace{p(\theta_{01}, \theta_{02} \mid \text{data}, I)}_{\text{Maximum of the posterior}} \exp\left[-\tfrac{1}{2}\Delta\chi^2(\theta_1, \theta_2)\right]$$

Advanced Data Analysis Course, 2019-20

# Parameter estimation: 2-D case

$p(\theta_1, \theta_2 \mid \text{data}, I)$

Minimum $\chi^2$
= maximum $\ell$

$\theta_2$

$\theta_1$

This is a **bivariate normal distribution** with **covariance matrix**

$$\sigma_{ij}^2 \;=\; \text{cov}_{ij} \;=\; \left\langle (\theta_i - \theta_{i0})(\theta_j - \theta_{j0}) \right\rangle \;=\; \left[ -\frac{\partial^2 \ell}{\partial \theta_i \partial \theta_j} \right]^{-1}$$
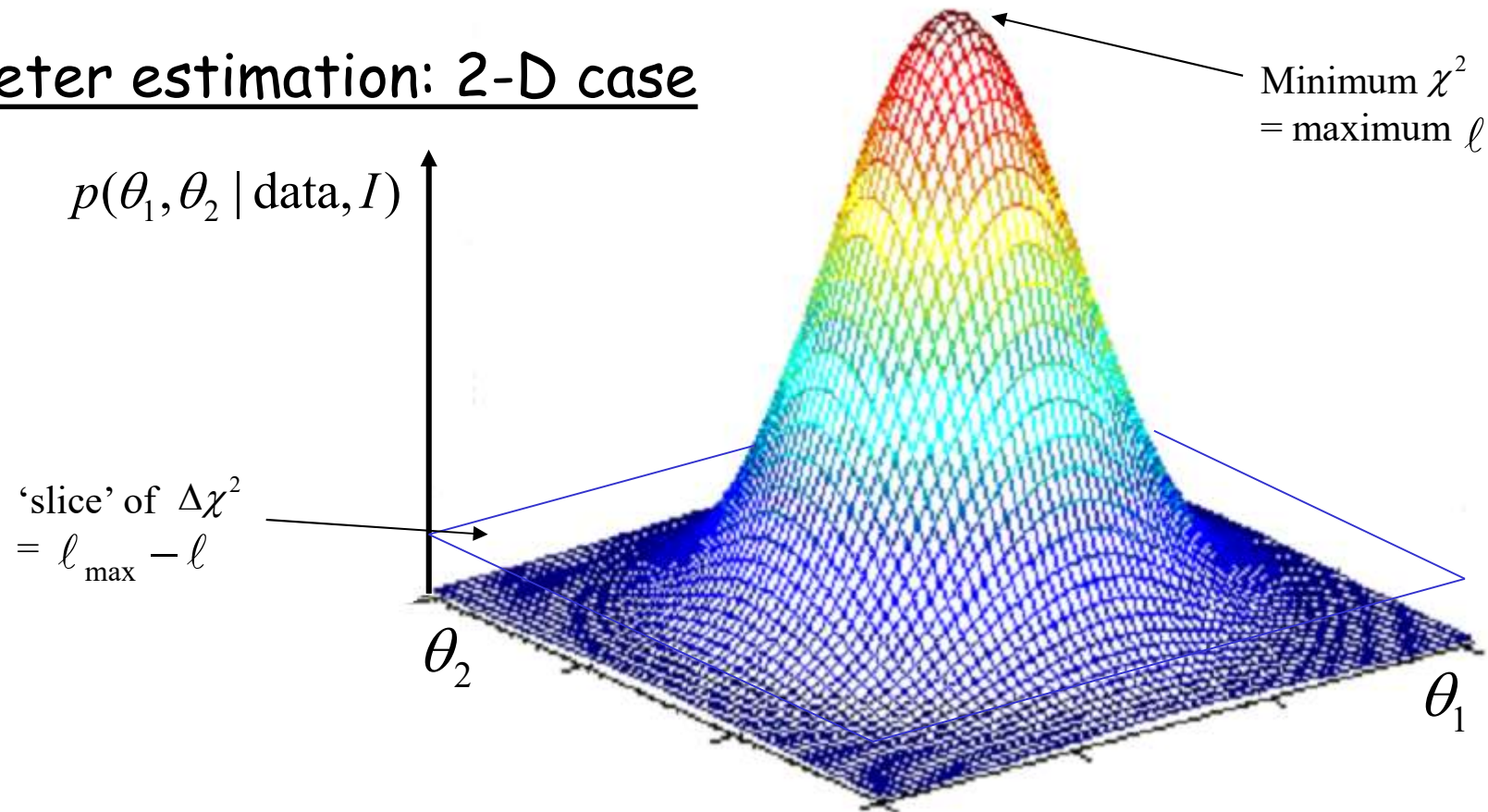
Fisher information matrix
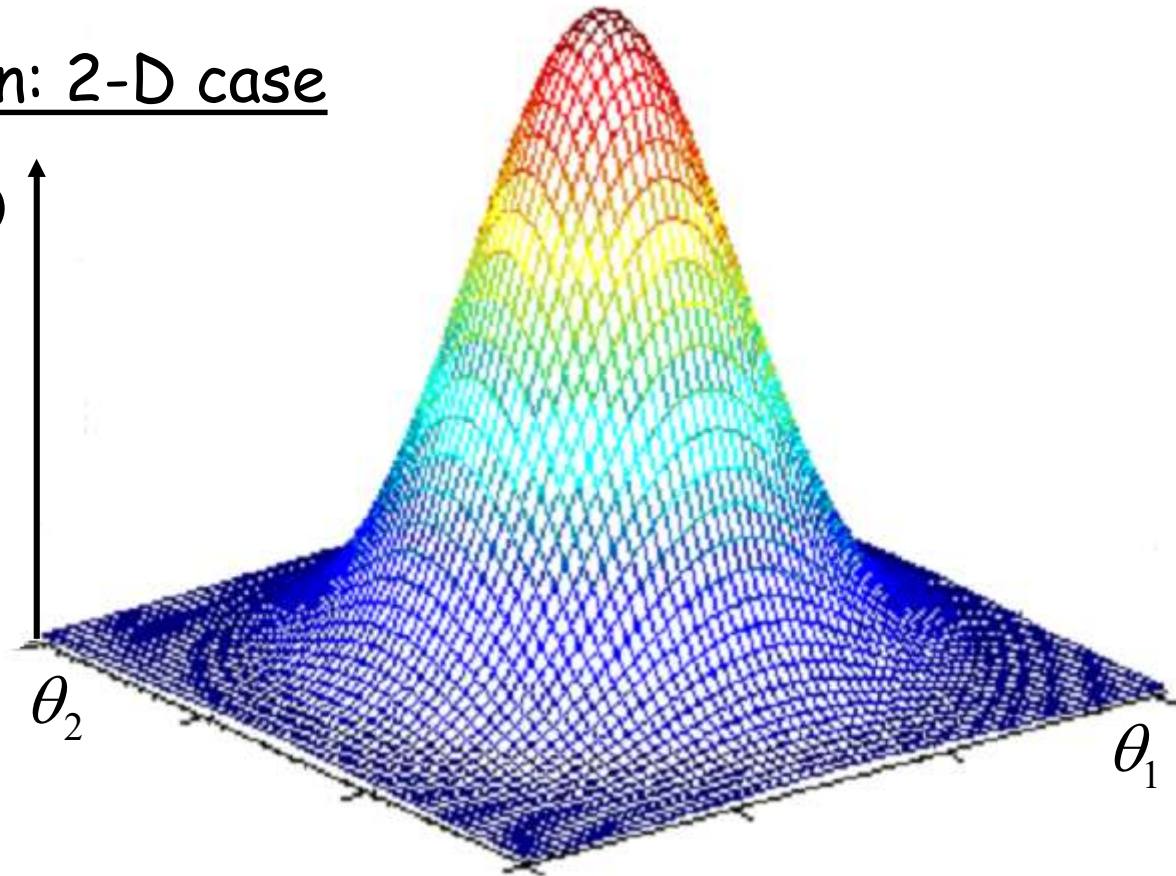
# Parameter estimation: 2-D case

$$p(\theta_1, \theta_2 \mid \text{data}, I)$$

Minimum $\chi^2$
= maximum $\ell$

'slice' of $\Delta\chi^2$
$= \ell_{max} - \ell$

$\theta_2$

$\theta_1$

This is a bivariate normal distribution with covariance matrix

$$\sigma_{ij}^2 \;=\; \mathrm{cov}_{ij} \;=\; \left\langle (\theta_i - \theta_{i0})(\theta_j - \theta_{j0}) \right\rangle \;=\; \left[ -\frac{\partial^2 \ell}{\partial \theta_i \partial \theta_j} \right]^{-1}$$

Fisher information matrix

Advanced Data Analysis Course, 2019-20

SUPA

# Parameter estimation: 2-D case

$$p(\theta_1, \theta_2 \mid \text{data}, I)$$

Minimum $\chi^2$
= maximum $\ell$

'slice' of $\Delta\chi^2$
= $\ell_{max} - \ell$

$\theta_2$

$\theta_1$

This is a bivariate normal distribution with covariance matrix

$$\sigma_{ij}^2 \;=\; \text{cov}_{ij} \;=\; \left\langle (\theta_i - \theta_{i0})(\theta_j - \theta_{j0}) \right\rangle \;=\; \left[ -\frac{\partial^2 \ell}{\partial \theta_i \partial \theta_j} \right]^{-1}$$

Fisher information matrix

University of Glasgow

SUPA

# Parameter estimation: 2-D case

$p(\theta_1, \theta_2 \mid \text{data}, I)$

Minimum $\chi^2$ = maximum $\ell$

'slice' of $\Delta\chi^2$ = $\ell_{max} - \ell$

$\theta_2$

$\theta_1$

This is a **bivariate normal distribution** with **covariance matrix**

$$\sigma_{ij}^2 \quad = \quad \text{cov}_{ij} \quad = \quad \left\langle (\theta_i - \theta_{i0})(\theta_j - \theta_{j0}) \right\rangle \quad = \quad \left[ -\frac{\partial^2 \ell}{\partial \theta_i \partial \theta_j} \right]^{-1}$$

Fisher information matrix

Advanced Data Analysis Course, 2019-20

SUPA

# Parameter estimation: 2-D case

$$p(\theta_1, \theta_2 \mid \text{data}, I)$$

We can compute the $\Delta\chi^2$ that corresponds to e.g. 68%, 95%, 99% of the posterior pdf.

We can draw contours of equal probability

$\Rightarrow$ **Credible regions for the parameters**

Extends easily to $N$ parameters – or *degrees of freedom*

$\theta_2$

$\theta_1$

University of Glasgow

SUPA

Advanced Data Analysis Course, 2019-20

# Parameter estimation: 2-D case

$p(\theta_1, \theta_2 \mid \text{data}, I)$

We can compute the $\Delta\chi^2$ that corresponds to e.g. 68%, 95%, 99% of the posterior pdf.

We can draw contours of equal probability

$\Rightarrow$ **Credible regions for the parameters**

Extends easily to $N$ parameters – or *degrees of freedom*

$\theta_2$

$\theta_1$

| | | | | $\nu$ | | | |
|---|---|---|---|---|---|---|
| $p$ | 1 | 2 | 3 | 4 | 5 | 6 |
| 68.3% | 1.00 | 2.30 | 3.53 | 4.72 | 5.89 | 7.04 |
| 90% | 2.71 | 4.61 | 6.25 | 7.78 | 9.24 | 10.6 |
| 95.4% | 4.00 | 6.17 | 8.02 | 9.70 | 11.3 | 12.8 |
| 99% | 6.63 | 9.21 | 11.3 | 13.3 | 15.1 | 16.8 |
| 99.73% | 9.00 | 11.8 | 14.2 | 16.3 | 18.2 | 20.1 |
| 99.99% | 15.1 | 18.4 | 21.1 | 23.5 | 25.7 | 27.8 |

$\Delta\chi^2$ as a Function of Confidence Level and Degrees of Freedom

*From Numerical Recipes*

Advanced Data Analysis Course, 2019-20

# Parameter estimation: 2-D case

$$p(\theta_1, \theta_2 \mid \text{data}, I)$$

Contours of constant probability are **ellipses**.

Covariance matrix is **not** in general diagonal

$\Rightarrow$ What we infer about $\theta_1$ and $\theta_1$ is **not** independent
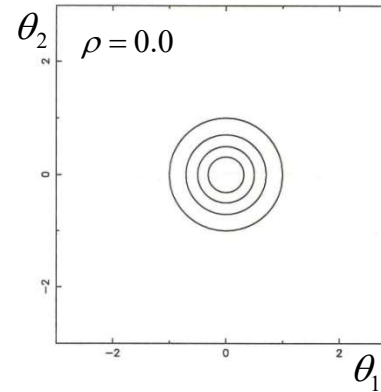
$\theta_2$

$\theta_1$

# Parameter estimation: 2-D case
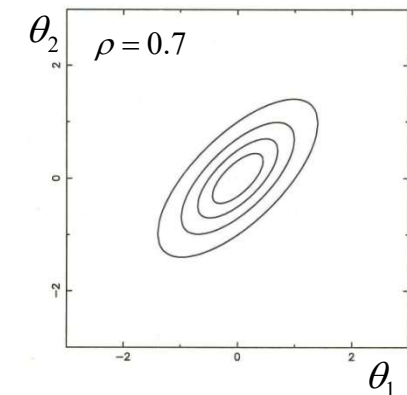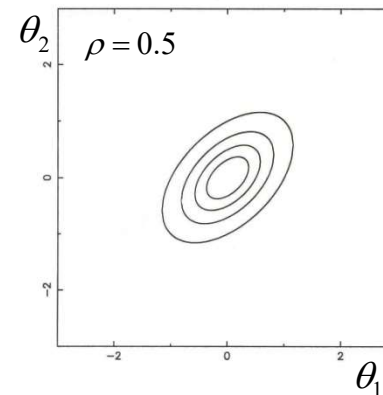
Can define *correlation coefficient*

$$\rho = \frac{\mathrm{cov}[\theta_1, \theta_2]}{\sqrt{\mathrm{var}[\theta_1]} \sqrt{\mathrm{var}[\theta_2]}} \qquad -1 \le \rho \le 1$$

Covariance matrix becomes less diagonal

$\Rightarrow$ $|\rho|$ increases

$\Rightarrow$ isoprobability contours elongate

University of Glasgow

SUPA

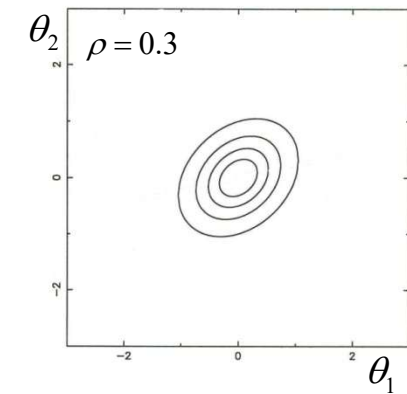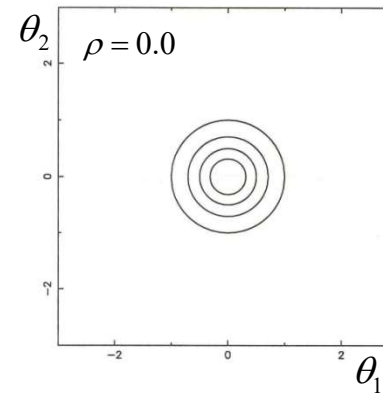# Parameter estimation: 2-D case

Can define *correlation coefficient*

$$\rho = \frac{\text{cov}[\theta_1, \theta_2]}{\sqrt{\text{var}[\theta_1]}\sqrt{\text{var}[\theta_2]}} \qquad -1 \le \rho \le 1$$

Covariance matrix becomes less diagonal

$\Rightarrow$ $|\rho|$ increases

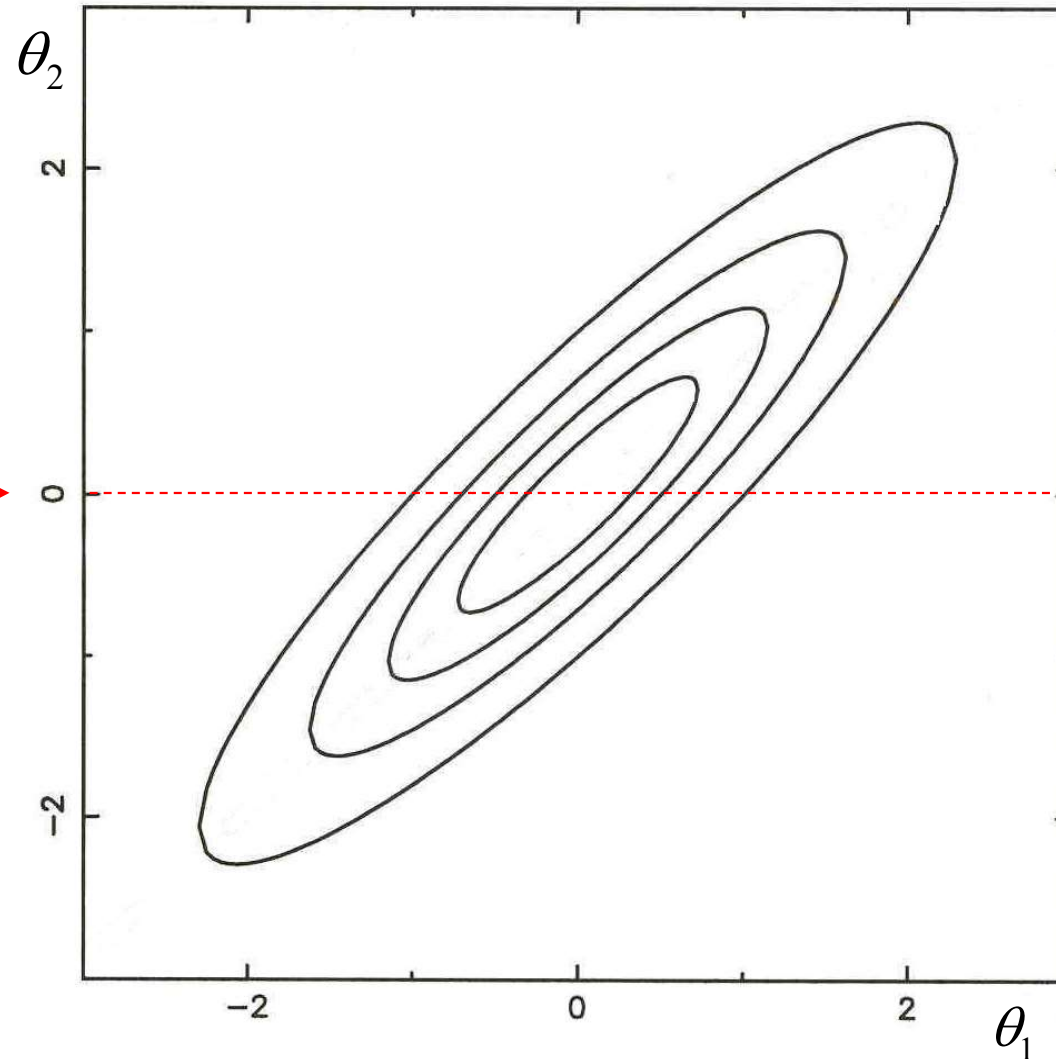$\Rightarrow$ isoprobability contours elongate

Very important if we are interested only in *one* parameter

# Parameter estimation: 2-D case



'Best-fit' value of $\theta_2$, found from

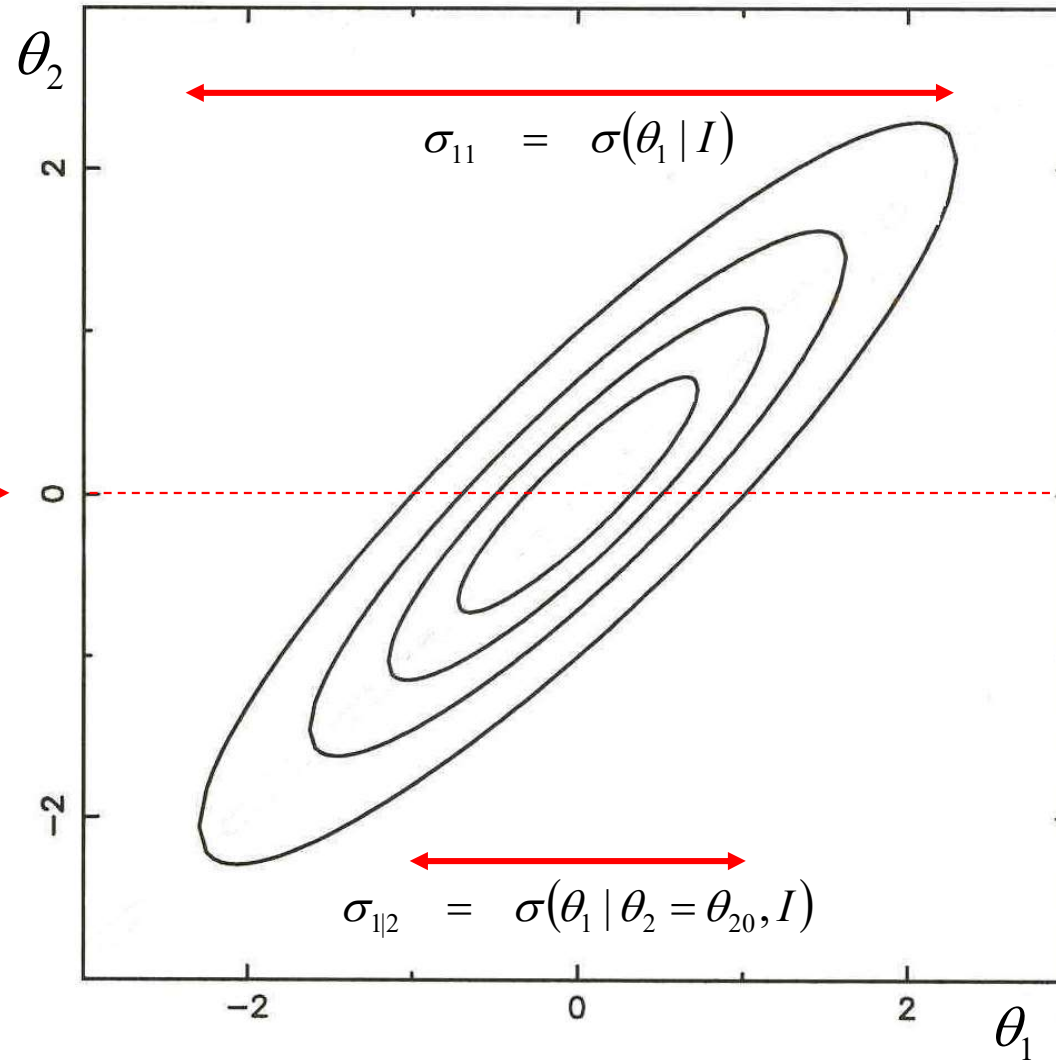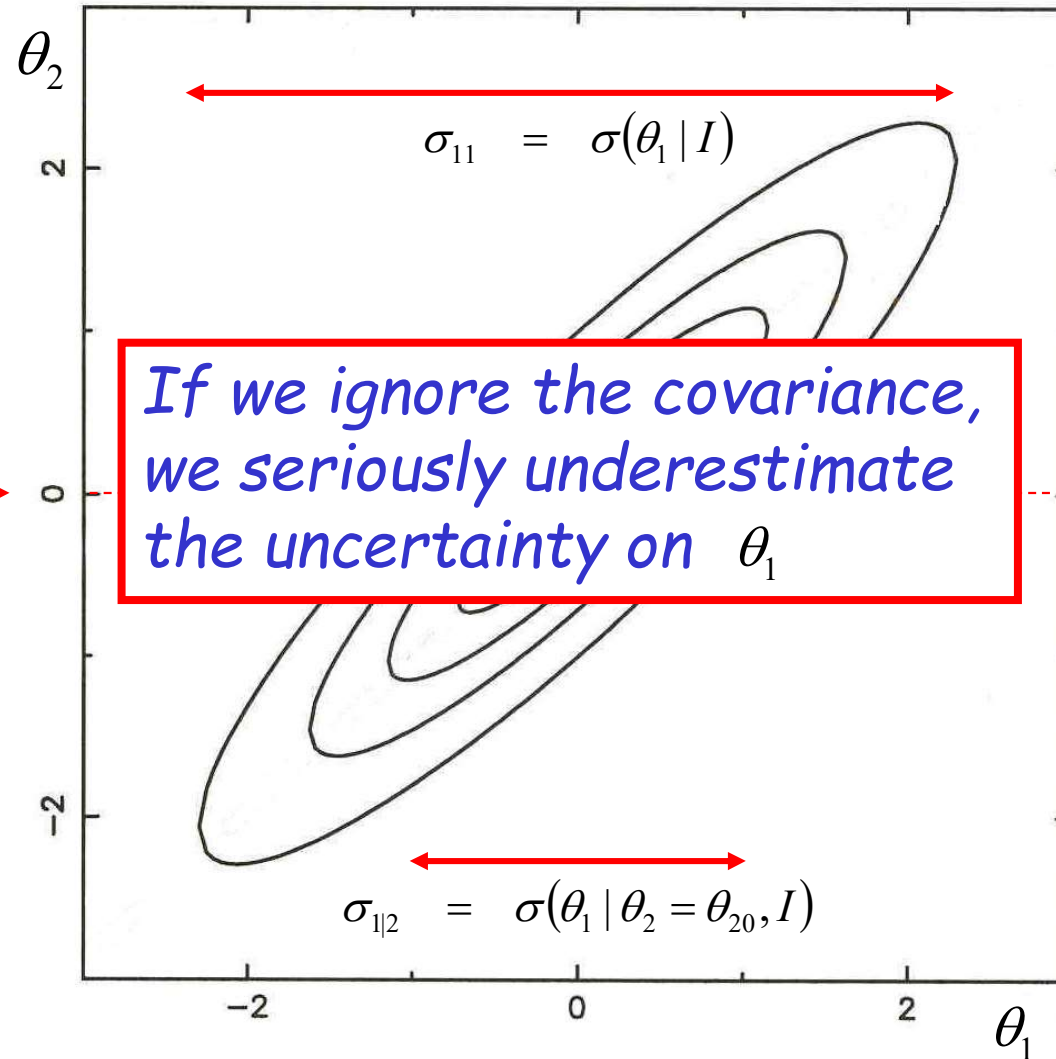$$\left.\frac{\partial \ell}{\partial \theta_j}\right|_{\theta_j = \theta_{0j}} = 0$$

# Parameter estimation: 2-D case



'Best-fit' value of $\theta_2$, found from

$$\left.\frac{\partial \ell}{\partial \theta_j}\right|_{\theta_j = \theta_{0j}} = 0$$

$\sigma_{11} = \sigma(\theta_1 \mid I)$

$\sigma_{1|2} = \sigma(\theta_1 \mid \theta_2 = \theta_{20}, I)$

# Parameter estimation: 2-D case

$$\theta_2$$

$$\sigma_{11} \;\; = \;\; \sigma(\theta_1 \,|\, I)$$

'Best-fit' value of $\theta_2$, found from

$$\left. \frac{\partial \ell}{\partial \theta_j} \right|_{\theta_j = \theta_{0j}} = 0$$

**If we ignore the covariance, we seriously underestimate the uncertainty on $\theta_1$**

$$\sigma_{1|2} \;\; = \;\; \sigma(\theta_1 \,|\, \theta_2 = \theta_{20}, I)$$

$$\theta_1$$

University of Glasgow

Advanced Data Analysis Course, 2019-20

SUPA

**Question 14:** The marginal and conditional error bars on $\theta_1$ will be equal provided

**A** $\quad \mathrm{cov}\left[\theta_1, \theta_2\right] = 0$

**B** $\quad \mathrm{cov}\left[\theta_1, \theta_2\right] = 1$

**C** $\quad \mathrm{cov}\left[\theta_1, \theta_2\right] = -1$
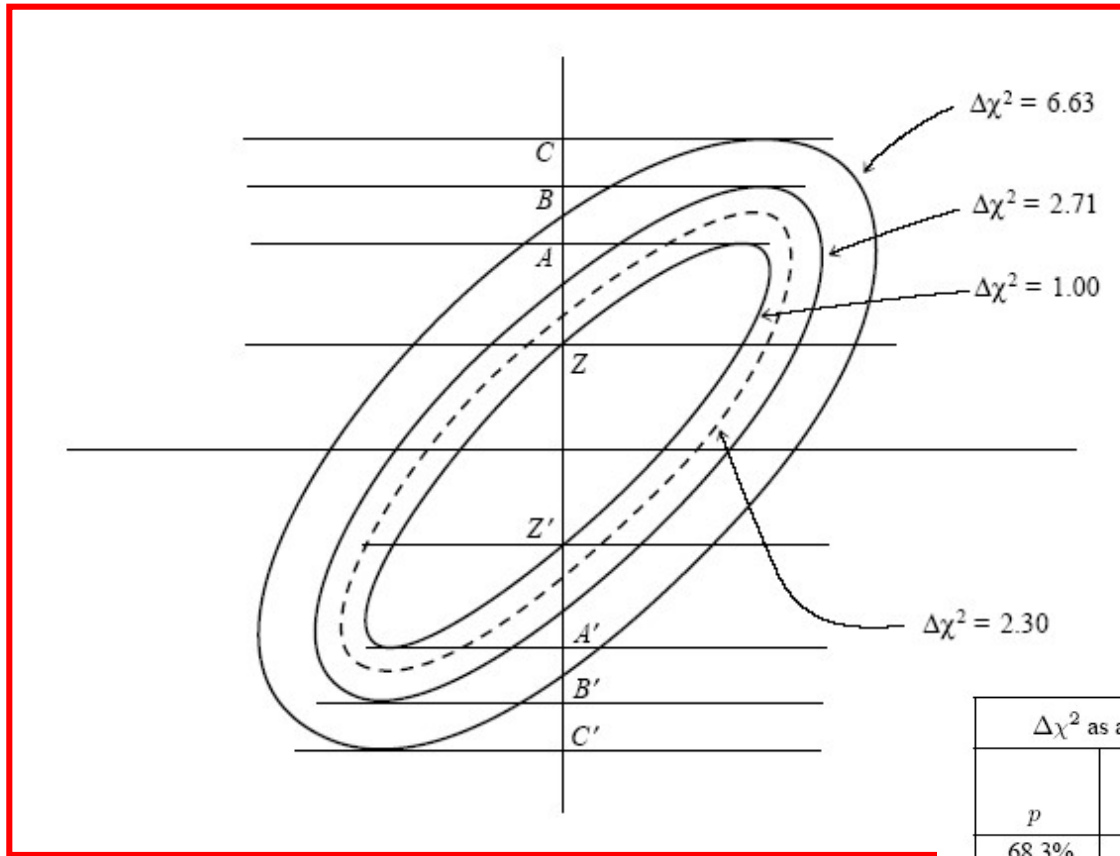
**D** $\quad$ None of the above

# Parameter estimation: 2-D case



'Best-fit' value of $\theta_2$, found from

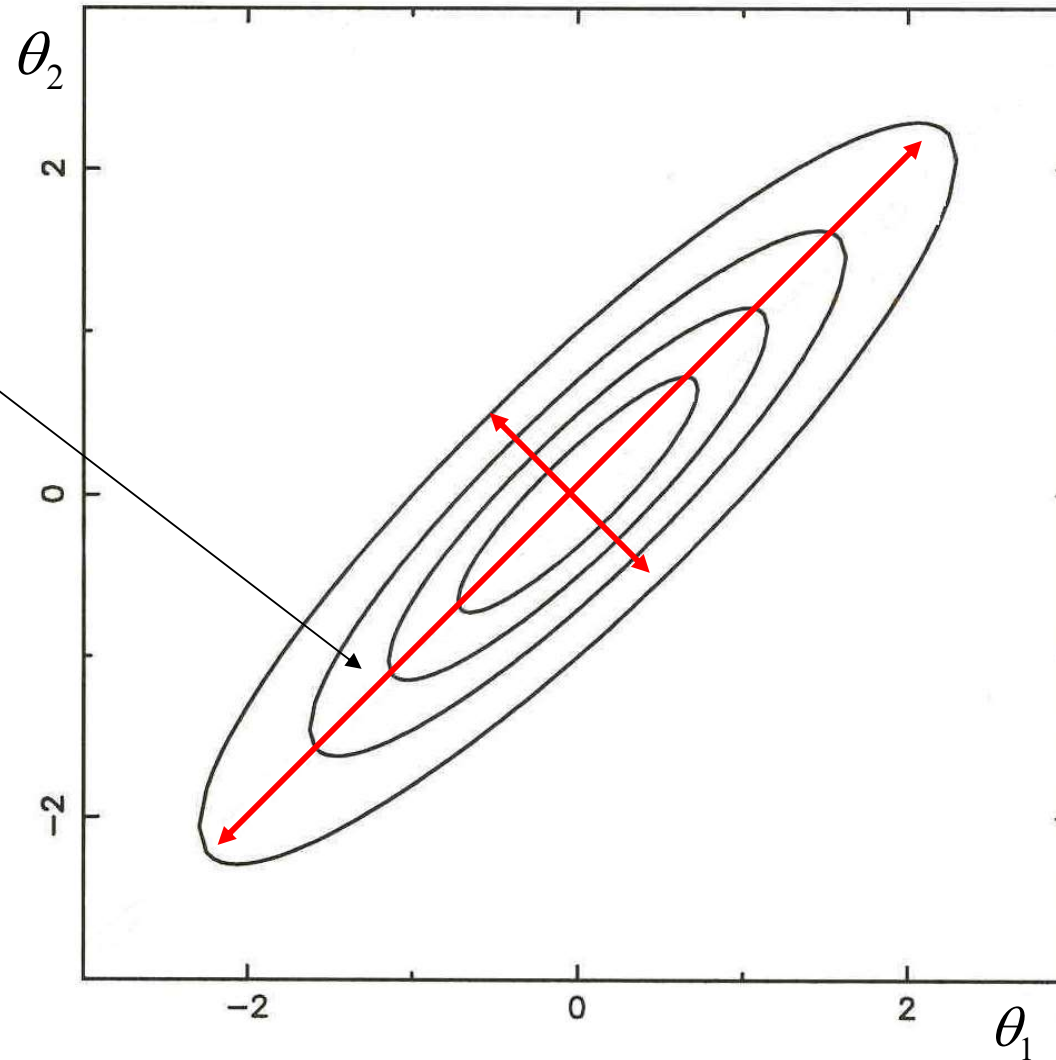$$\left.\frac{\partial \ell}{\partial \theta_j}\right|_{\theta_j = \theta_{0j}} = 0$$

$\sigma_{11} = \sigma(\theta_1 \mid I)$

$\sigma_{1|2} = \sigma(\theta_1 \mid \theta_2 = \theta_{20}, I)$

Marginal and conditional error bars only equal if $\mathrm{cov}[\theta_1, \theta_2] = 0$

# Parameter estimation: 2-D case



*From Numerical Recipes*

| $\Delta\chi^2$ as a Function of Confidence Level and Degrees of Freedom | | | | | | |
|---|---|---|---|---|---|---|
| | | | $\nu$ | | | |
| $p$ | 1 | 2 | 3 | 4 | 5 | 6 |
| 68.3% | 1.00 | 2.30 | 3.53 | 4.72 | 5.89 | 7.04 |
| 90% | 2.71 | 4.61 | 6.25 | 7.78 | 9.24 | 10.6 |
| 95.4% | 4.00 | 6.17 | 8.02 | 9.70 | 11.3 | 12.8 |
| 99% | 6.63 | 9.21 | 11.3 | 13.3 | 15.1 | 16.8 |
| 99.73% | 9.00 | 11.8 | 14.2 | 16.3 | 18.2 | 20.1 |
| 99.99% | 15.1 | 18.4 | 21.1 | 23.5 | 25.7 | 27.8 |

# Parameter estimation: 2-D case

Linear combination of $\theta_1$ and $\theta_2$ well constrained by data
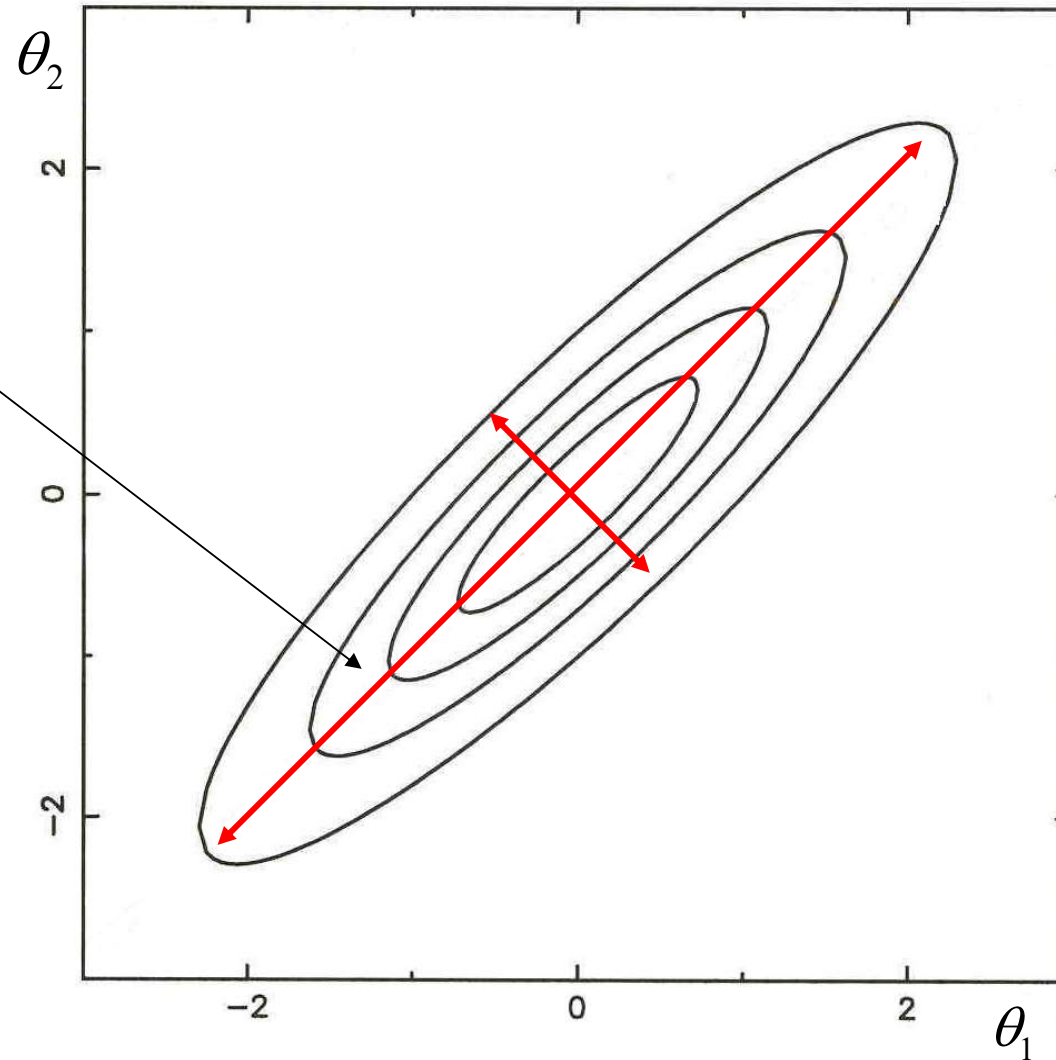
# Parameter estimation: 2-D case

Linear combination
of $\theta_1$ and $\theta_2$ well
constrained by data

Length of axes
determined by the
**eigenvalues** of the
Fisher information
matrix

$$F_{ij} = \frac{\partial^2 \ell}{\partial \theta_i \partial \theta_j} = \left[ -\sigma_{ij}^2 \right]^{-1}$$

$$\boxed{\boldsymbol{F}\,\boldsymbol{\theta} = \lambda\,\boldsymbol{\theta}}$$
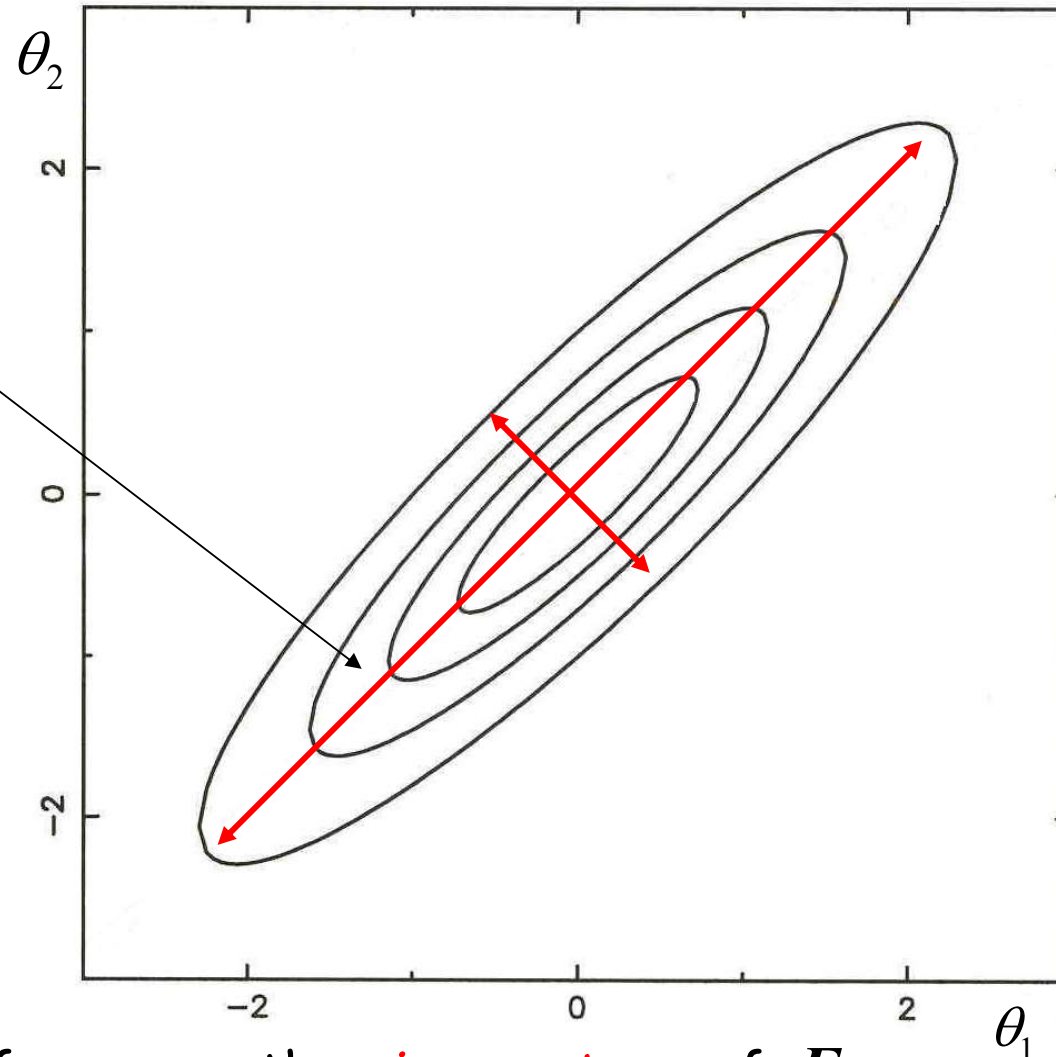


University
of Glasgow

# Parameter estimation: 2-D case

Linear combination
of $\theta_1$ and $\theta_2$ well
constrained by data

Length of axes
determined by the
**eigenvalues** of the
Fisher information
matrix

$$F_{ij} = \frac{\partial^2 \ell}{\partial \theta_i \partial \theta_j} = \left[-\sigma_{ij}^2\right]^{-1}$$

$$\boxed{F\,\boldsymbol{\theta} = \lambda\,\boldsymbol{\theta}}$$

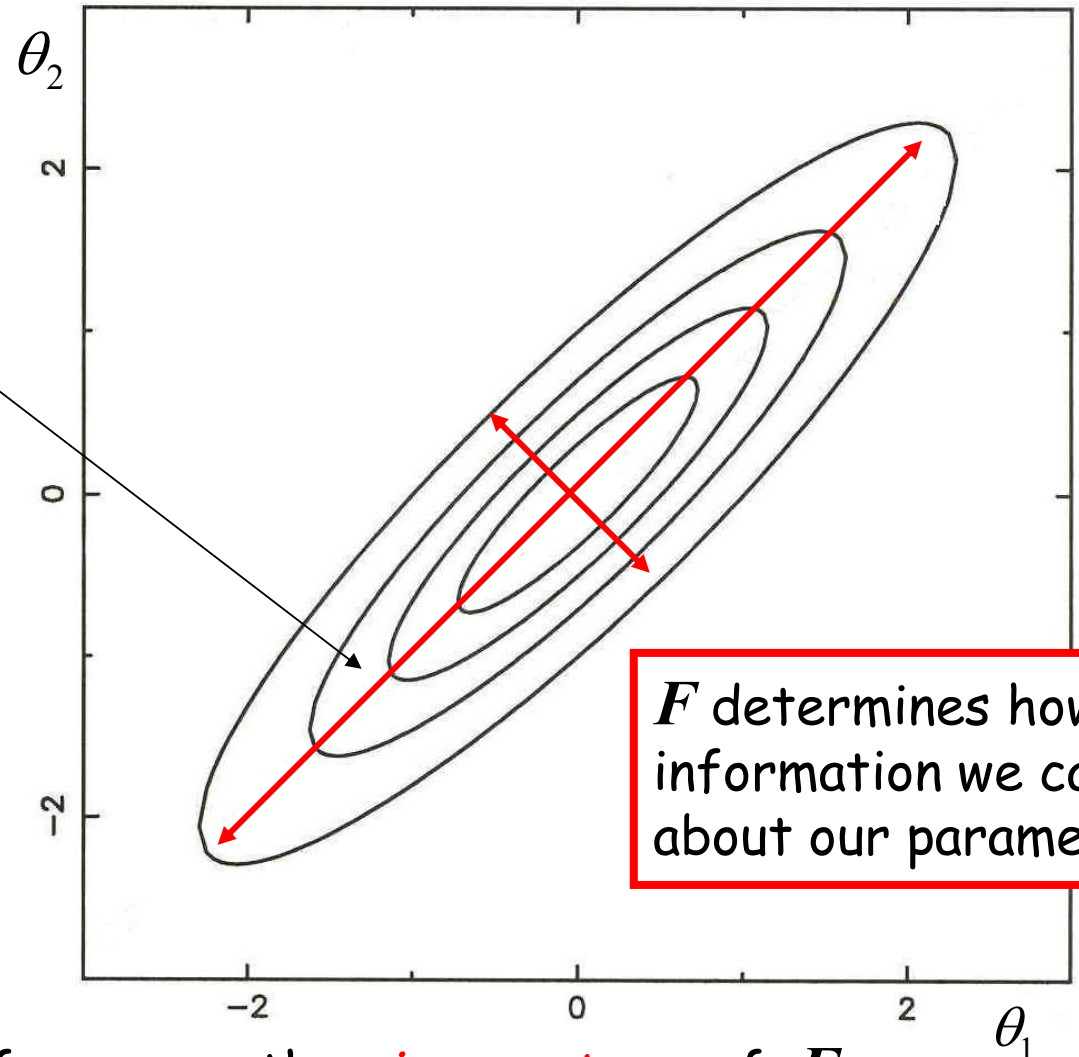Direction of axes are the *eigenvectors* of $F$

# Parameter estimation: 2-D case

Linear combination of $\theta_1$ and $\theta_2$ well constrained by data

Length of axes determined by the **eigenvalues** of the Fisher information matrix

$$F_{ij} = \frac{\partial^2 \ell}{\partial \theta_i \partial \theta_j} = \left[-\sigma_{ij}^2\right]^{-1}$$

$$\boxed{F\,\boldsymbol{\theta} = \lambda\,\boldsymbol{\theta}}$$

$\boldsymbol{F}$ determines how much information we can learn about our parameters

Direction of axes are the *eigenvectors* of $\boldsymbol{F}$

University of Glasgow
VIA VERITAS VITA

SUPA

Advanced Data Analysis Course, 2019-20