

## 4. An Advanced Bayesian Toolbox



## Parameter estimation:

Posterior

Likelihood

Prior

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

## Parameter estimation: the Gaussian approximation

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

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

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

## Parameter estimation: the Gaussian approximation

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

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

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

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

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

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

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

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

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

## Parameter estimation: the Gaussian approximation

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

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

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

## Parameter estimation: the Gaussian approximation

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

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

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

## Parameter estimation: the Gaussian approximation

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

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

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

Can summarise inference from posterior by

$$\theta = \theta_0 \pm \sigma$$



**Question 9:** Neglecting the higher order terms in the log likelihood 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

**Question 9:** Neglecting the higher order terms in the log likelihood 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*

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

Extends to 2 variables - *covariance*

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

## Parameter estimation: 2-D case

Recall our definition of *variance*

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

Extends to 2 variables - *covariance*

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

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

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

## Parameter estimation: 2-D case

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

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

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

## Parameter estimation: 2-D case

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

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

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

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

## Parameter estimation: 2-D case

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

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

## Parameter estimation: 2-D case

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

$$\begin{aligned} \ell(\theta_1, \theta_2) = & \ell(\theta_{01}, \theta_{02}) + \frac{\partial \ell}{\partial \theta_1} \Big|_{\theta_j = \theta_{0j}} (\theta_1 - \theta_{01}) + \frac{\partial \ell}{\partial \theta_2} \Big|_{\theta_j = \theta_{0j}} (\theta_2 - \theta_{02}) + \\ & \frac{1}{2} \left[ \frac{\partial^2 \ell}{\partial \theta_1^2} \Big|_{\theta_j = \theta_{0j}} (\theta_1 - \theta_{01})^2 + \frac{\partial^2 \ell}{\partial \theta_2^2} \Big|_{\theta_j = \theta_{0j}} (\theta_2 - \theta_{02})^2 + 2 \frac{\partial^2 \ell}{\partial \theta_1 \partial \theta_2} \Big|_{\theta_j = \theta_{0j}} (\theta_1 - \theta_{01})(\theta_2 - \theta_{02}) \right] + \dots \end{aligned}$$

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

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

$$\chi^2 \nearrow$$



## Parameter estimation: 2-D case

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

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

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

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

$\chi^2$   $\nearrow$  Maximising likelihood  
 $\equiv$  Minimising  $\chi^2$

## Parameter estimation: 2-D case

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

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

← Quadratic form

where

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

## Parameter estimation: 2-D case

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

$$Q = (\theta_1 - \theta_{10} \quad \theta_2 - \theta_{20}) \begin{bmatrix} A & C \\ C & B \end{bmatrix} \begin{pmatrix} \theta_1 - \theta_{10} \\ \theta_2 - \theta_{20} \end{pmatrix} \quad \leftarrow \text{Quadratic form}$$

where  $A = \frac{\partial^2 \ell}{\partial \theta_1^2} \Big|_{\theta_j = \theta_{0j}}$        $B = \frac{\partial^2 \ell}{\partial \theta_2^2} \Big|_{\theta_j = \theta_{0j}}$        $C = \frac{\partial^2 \ell}{\partial \theta_1 \partial \theta_2} \Big|_{\theta_j = \theta_{0j}}$

This is a bivariate normal distribution with covariance matrix

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

Fisher information matrix  $\nearrow$

$$\mathbf{F} \equiv F_{ij} = \frac{\partial^2 \ell}{\partial \theta_i \partial \theta_j} \quad \text{is known as the Fisher information matrix}$$

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**.

$$\text{Then } \mathbf{F} = -\text{diag}(\sigma_1^{-2}, \dots, \sigma_n^{-2})$$

If the  $i^{\text{th}}$  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 6 )

We can write the log likelihood as

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

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

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

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

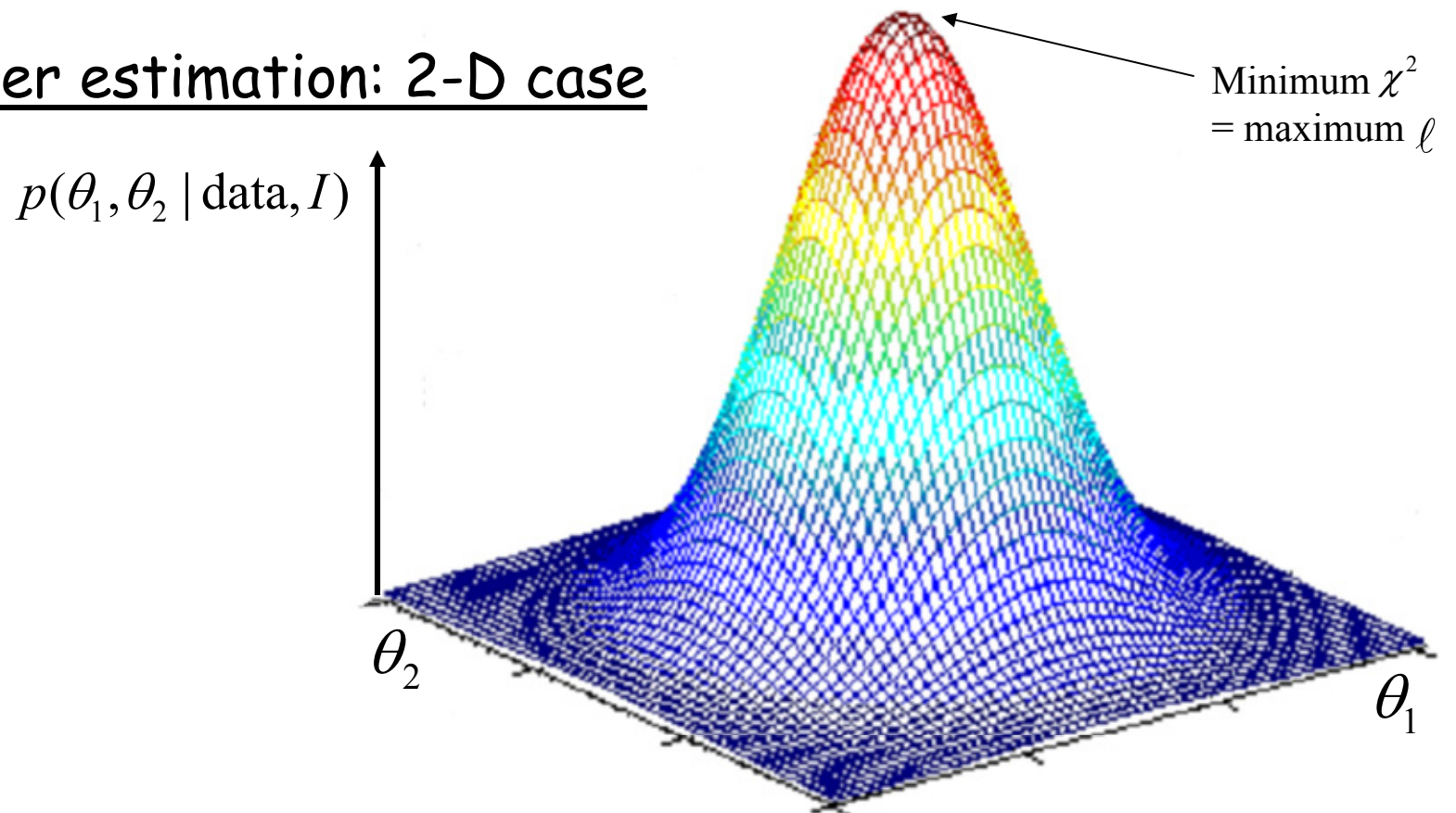
$$\ell(\theta_1, \theta_2) = \ell(\theta_{01}, \theta_{02}) - \frac{1}{2} \Delta\chi^2(\theta_1, \theta_2)$$

So that

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

Maximum of the posterior

## Parameter estimation: 2-D case

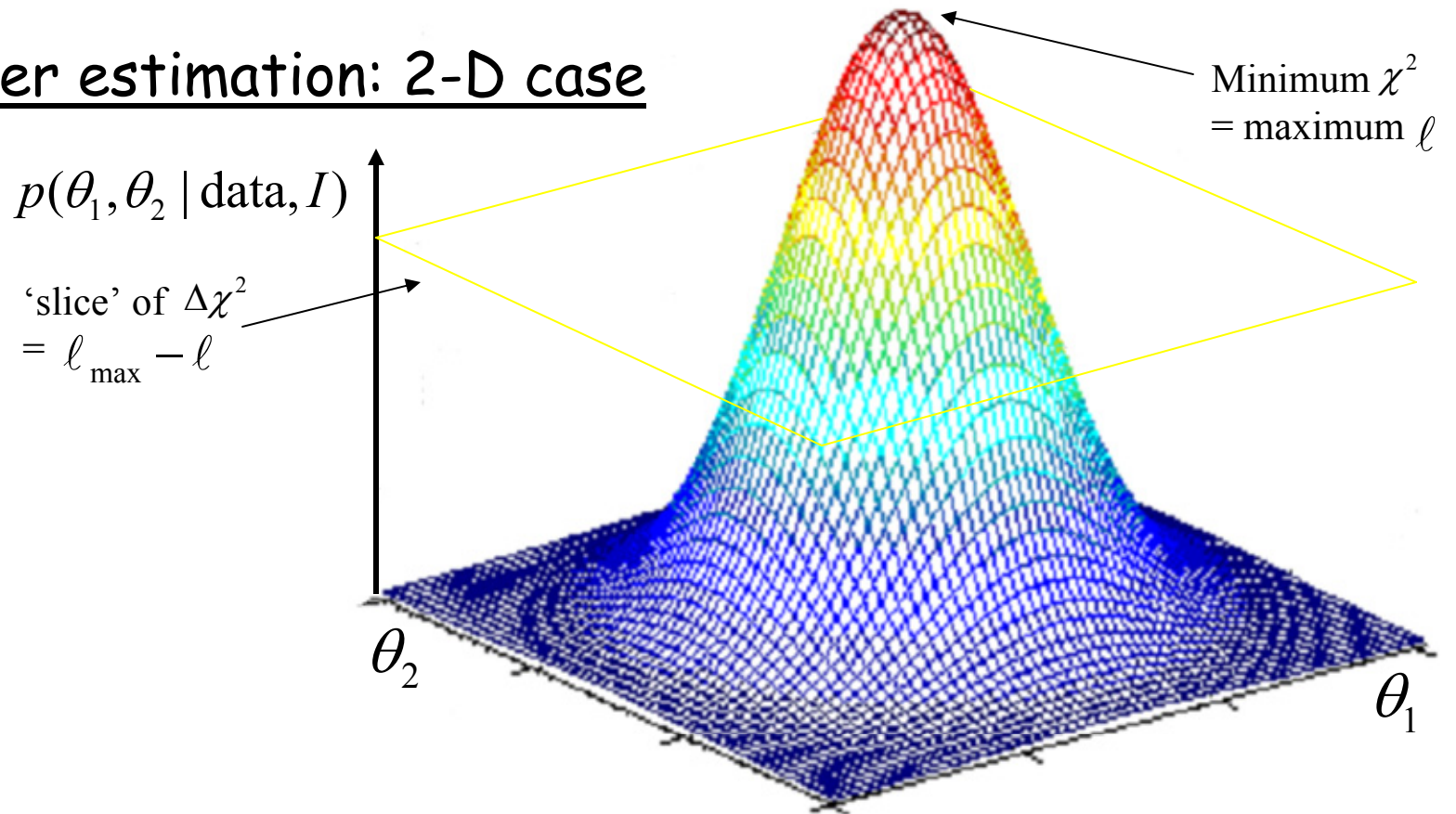


This is a **bivariate normal distribution with covariance matrix**

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

**Fisher information matrix**

# Parameter estimation: 2-D case



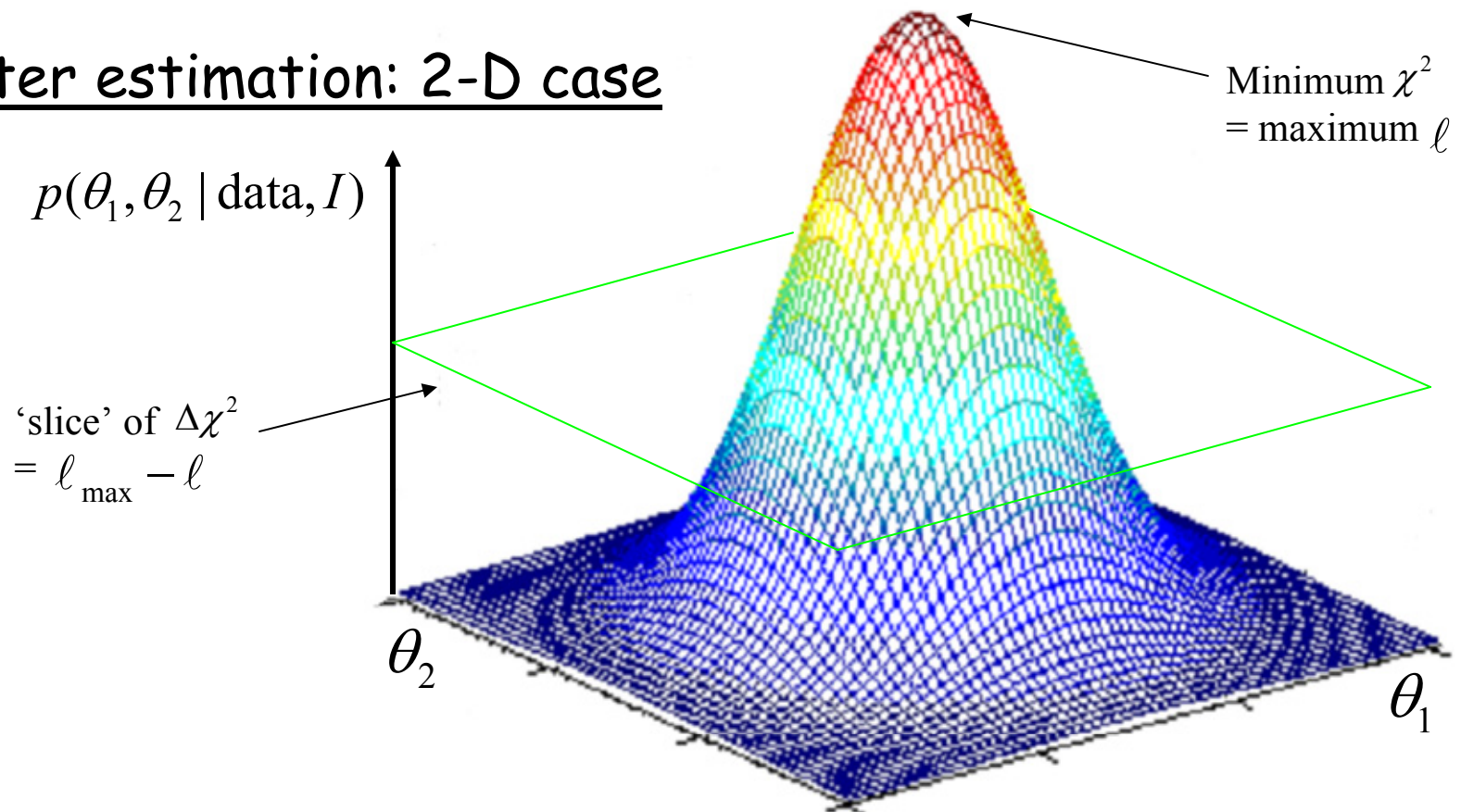
This is a **bivariate normal distribution with covariance matrix**

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

**Fisher information matrix**



## Parameter estimation: 2-D case

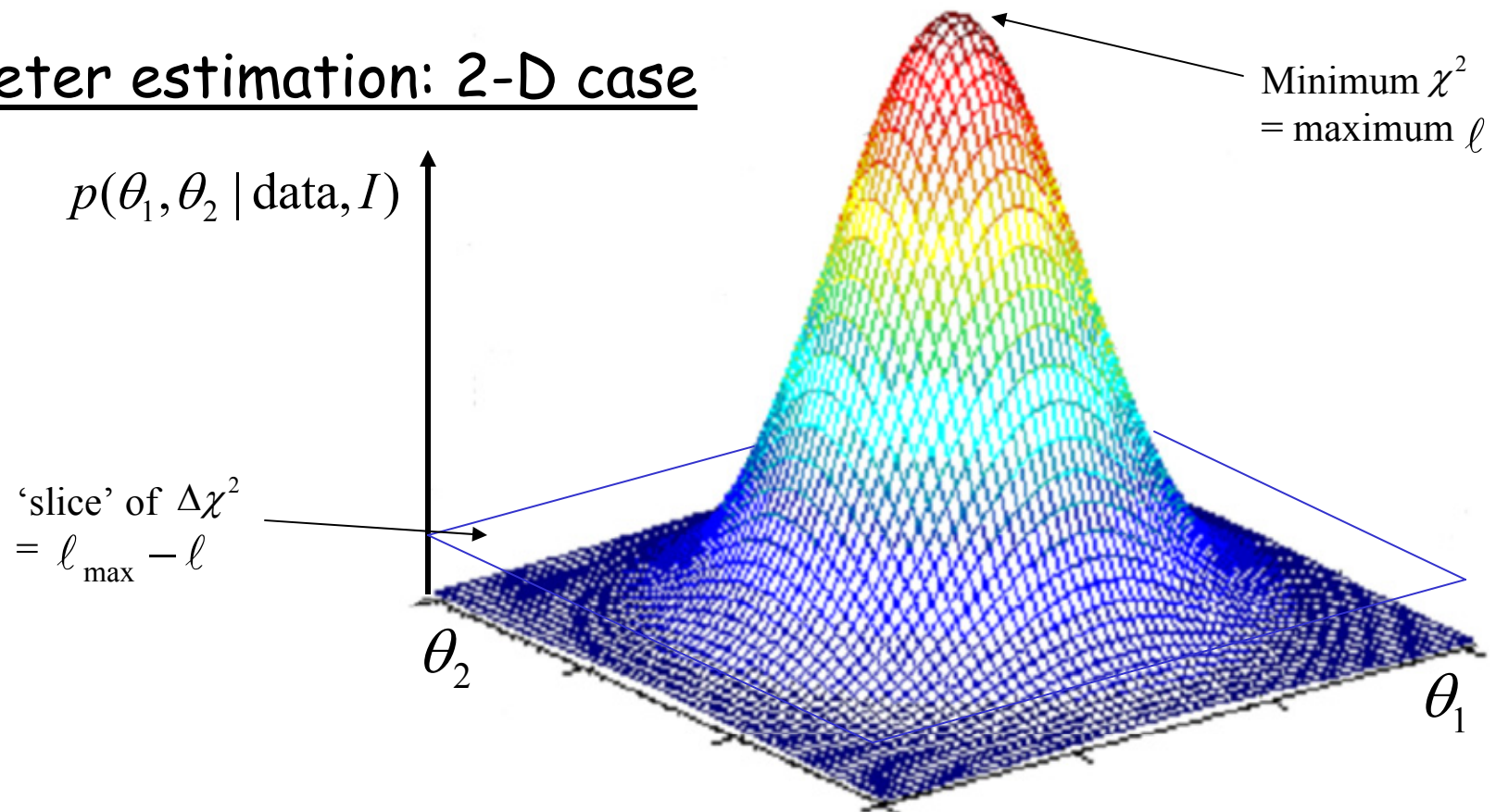


This is a **bivariate normal distribution with covariance matrix**

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

**Fisher information matrix**

## Parameter estimation: 2-D case



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

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

**Fisher information matrix**

## Parameter estimation: 2-D case

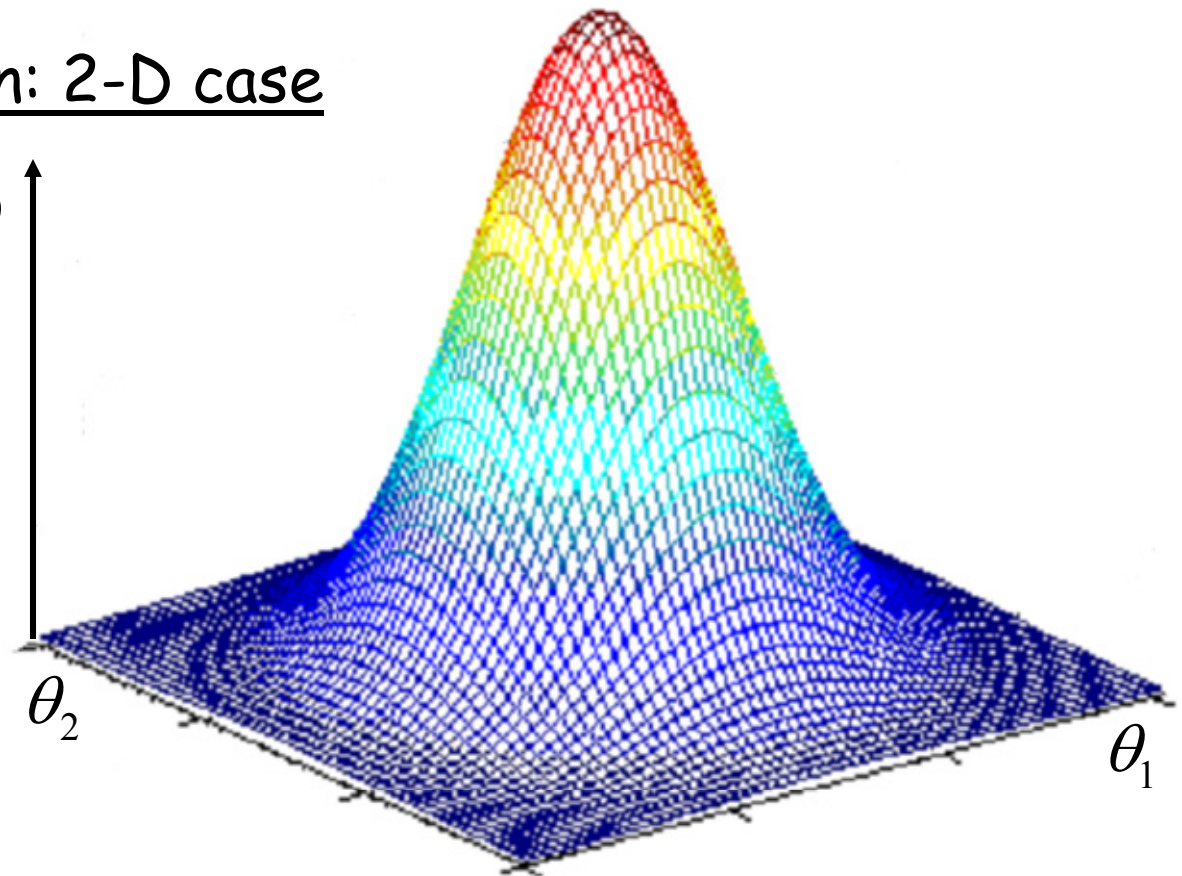
$$p(\theta_1, \theta_2 | \text{data}, I)$$

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

⇒ **Credible regions for the parameters**

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



# Parameter estimation: 2-D case

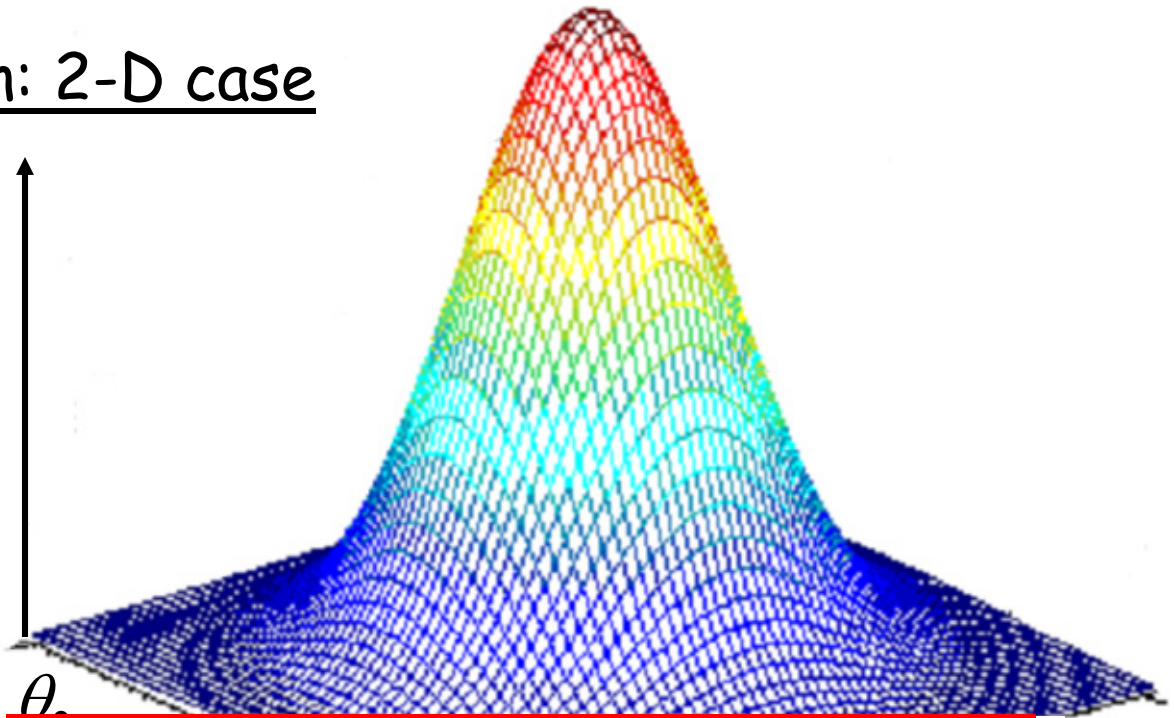
$$p(\theta_1, \theta_2 | \text{data}, I)$$

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

⇒ **Credible regions for the parameters**

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



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

From Numerical Recipes



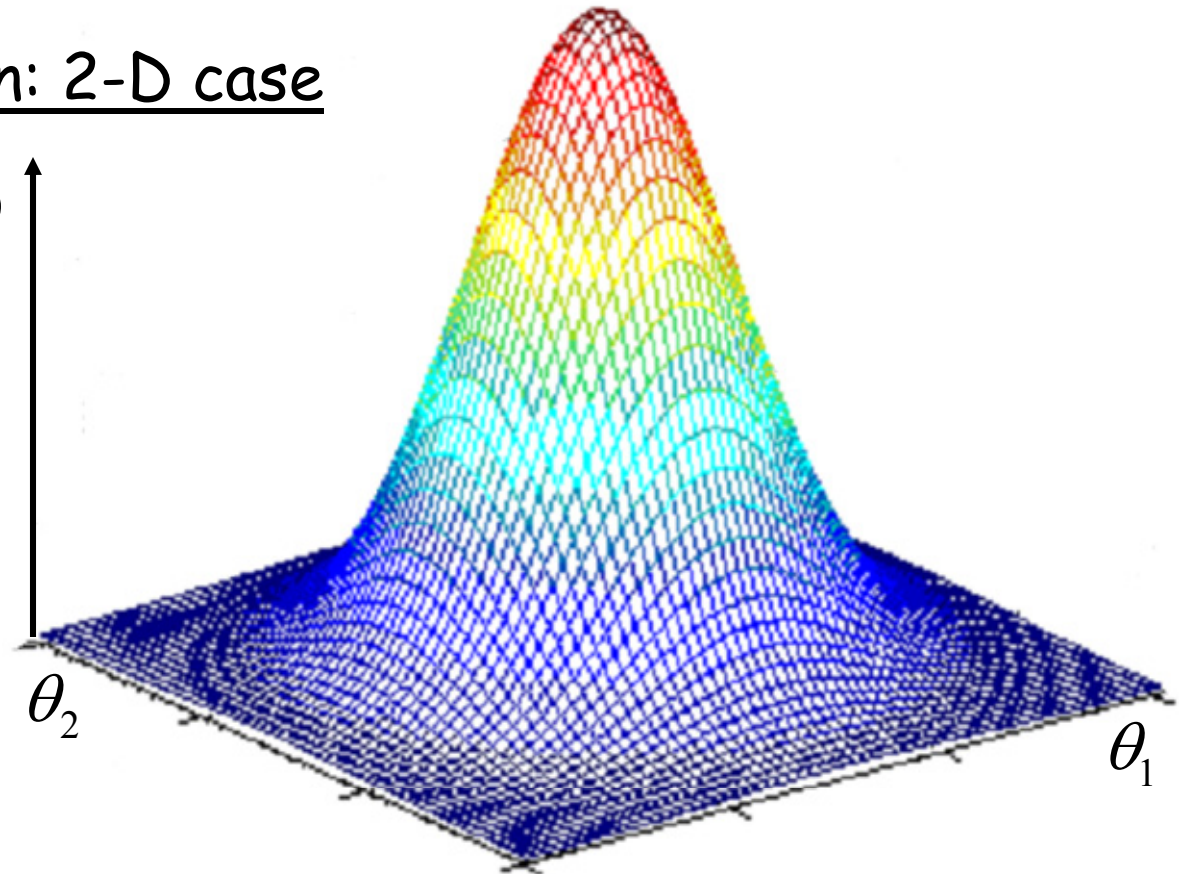
## Parameter estimation: 2-D case

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

Contours of constant probability are **ellipses**.

Covariance matrix is **not** in general diagonal

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



# Parameter estimation: 2-D case

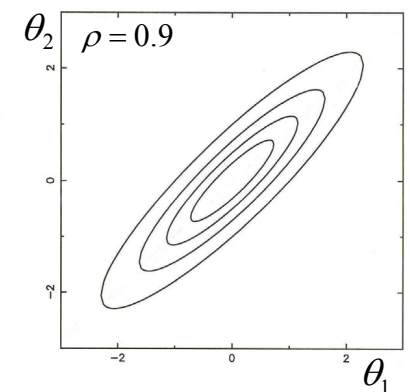
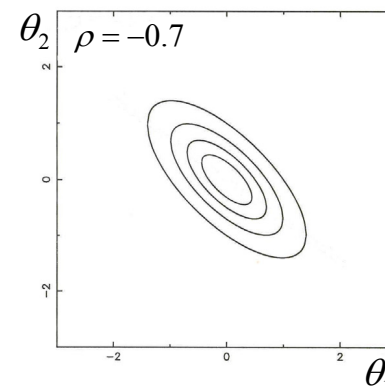
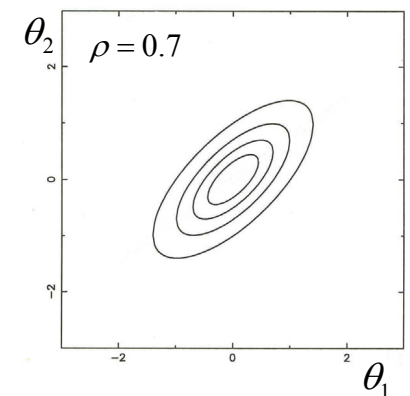
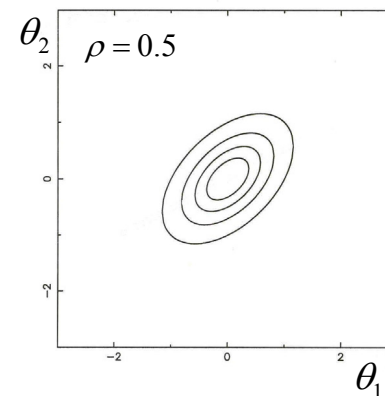
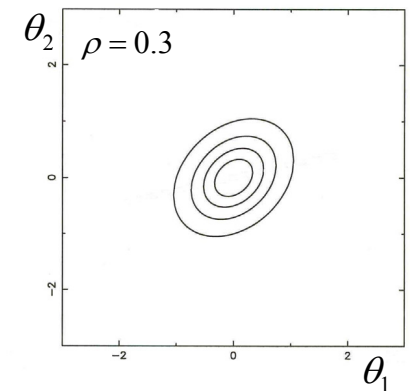
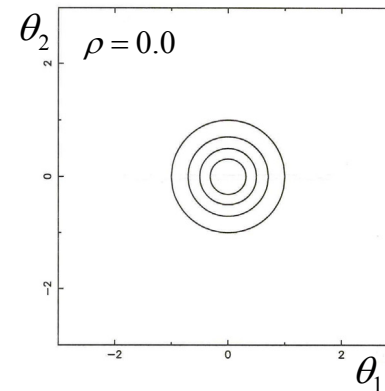
Can define *correlation coefficient*

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

Covariance matrix becomes less diagonal

⇒  $|\rho|$  increases

⇒ isoprobability contours elongate



# Parameter estimation: 2-D case

Can define *correlation coefficient*

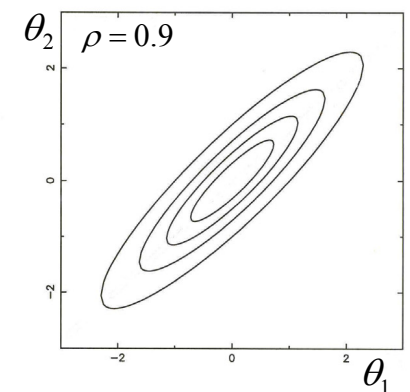
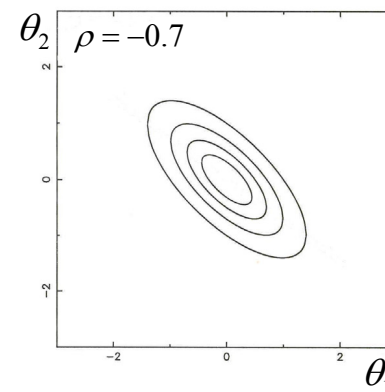
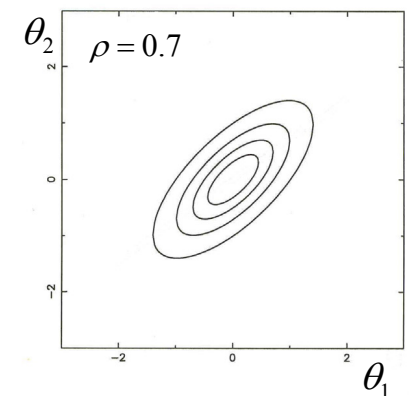
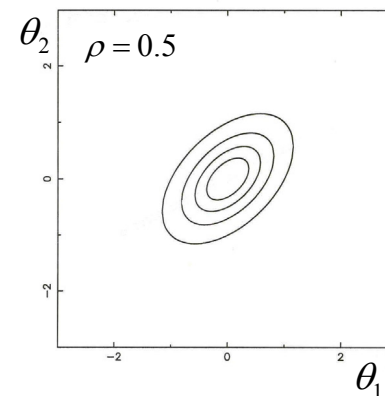
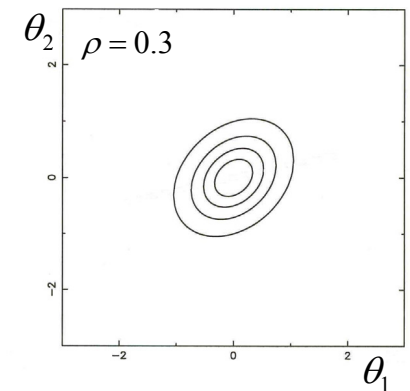
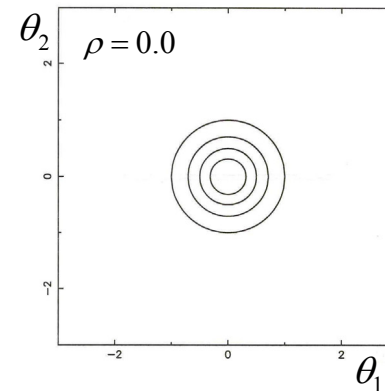
$$\rho = \frac{\text{cov}[\theta_1, \theta_2]}{\sqrt{\text{var}[\theta_1]} \sqrt{\text{var}[\theta_2]}} \quad -1 \leq \rho \leq 1$$

Covariance matrix becomes less diagonal

⇒  $|\rho|$  increases

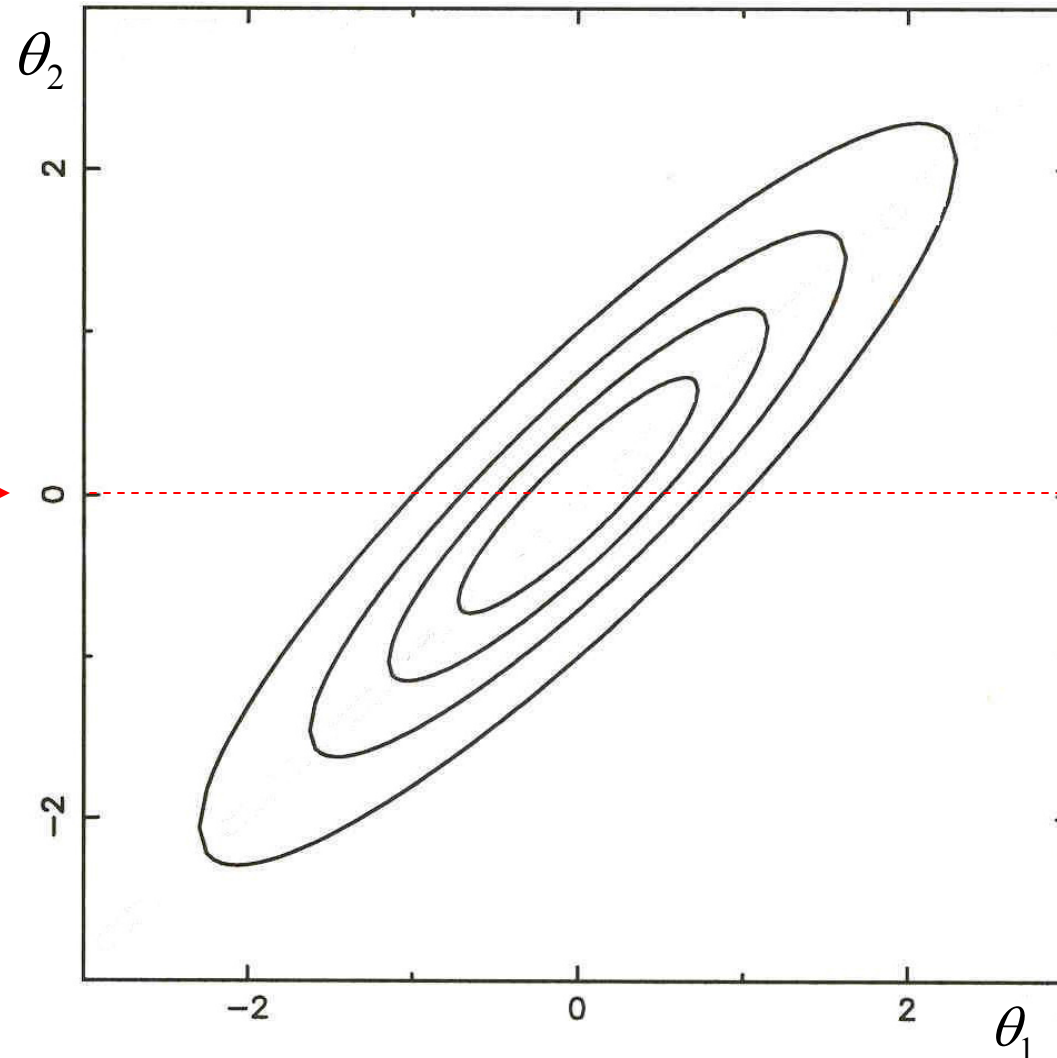
⇒ isoprobability contours elongate

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



## Parameter estimation: 2-D case

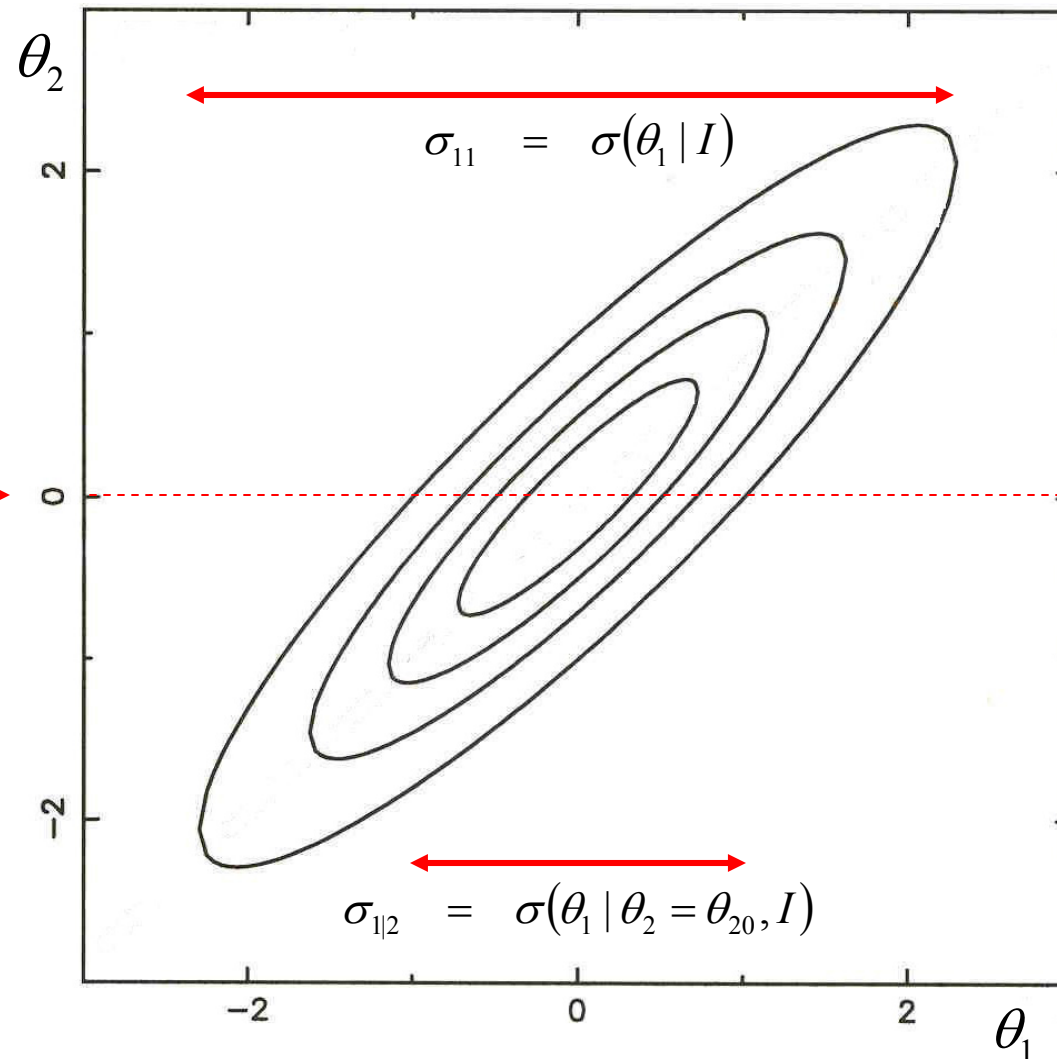
'Best-fit' value  
of  $\theta_2$ , found  
from  $\left. \frac{\partial \ell}{\partial \theta_j} \right|_{\theta_j = \theta_{0j}} = 0$





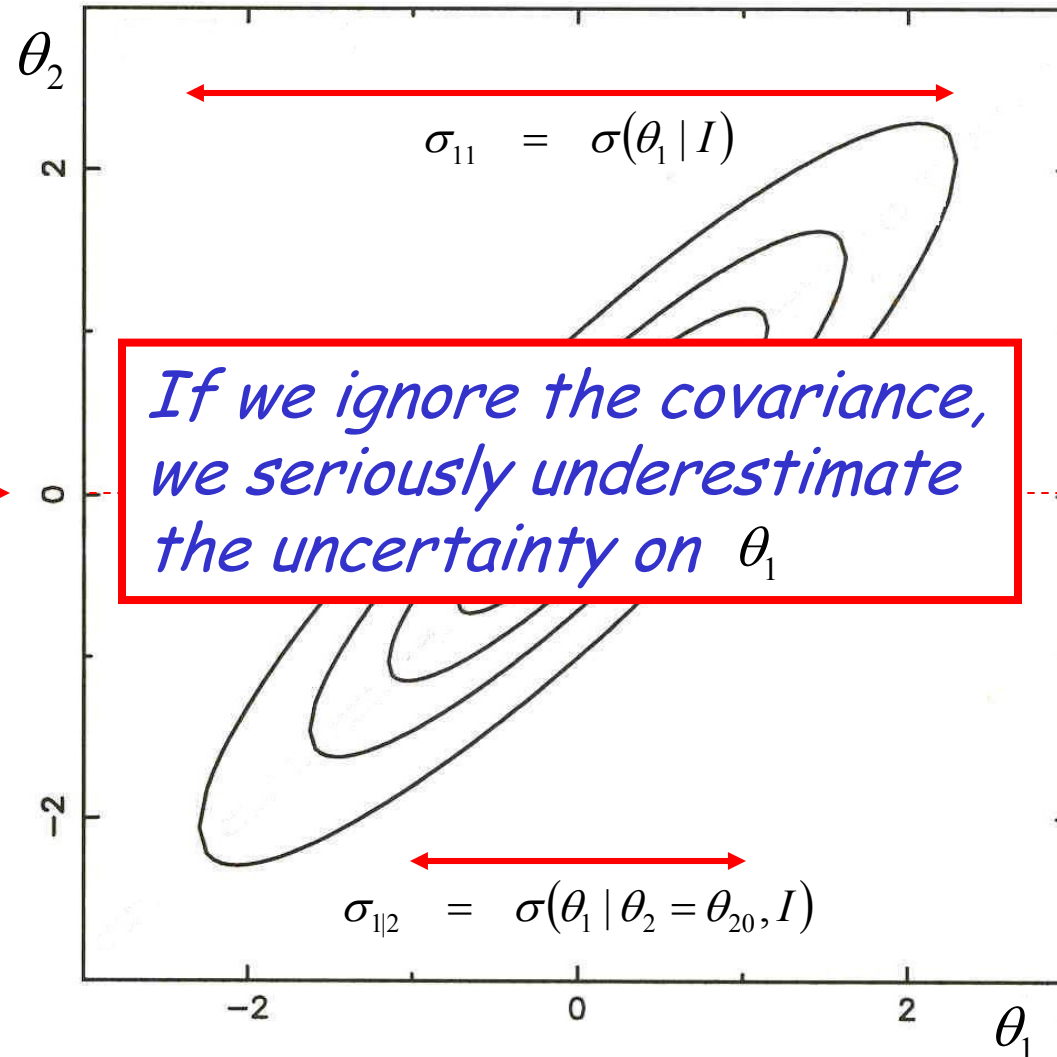
## Parameter estimation: 2-D case

'Best-fit' value  
of  $\theta_2$ , found  
from  $\left. \frac{\partial \ell}{\partial \theta_j} \right|_{\theta_j = \theta_{0j}} = 0$



# Parameter estimation: 2-D case

'Best-fit' value  
of  $\theta_2$ , found  
from  $\frac{\partial \ell}{\partial \theta_j} \Big|_{\theta_j = \theta_{0j}} = 0$



**Question 10:** The marginal and conditional error bars on  $\theta_1$  will be equal provided

**A**  $\text{cov}[\theta_1, \theta_2] = 0$

**B**  $\text{cov}[\theta_1, \theta_2] = 1$

**C**  $\text{cov}[\theta_1, \theta_2] = -1$

**D** None of the above



**Question 10:** The marginal and conditional error bars on  $\theta_1$  will be equal provided

**A**  $\text{cov}[\theta_1, \theta_2] = 0$

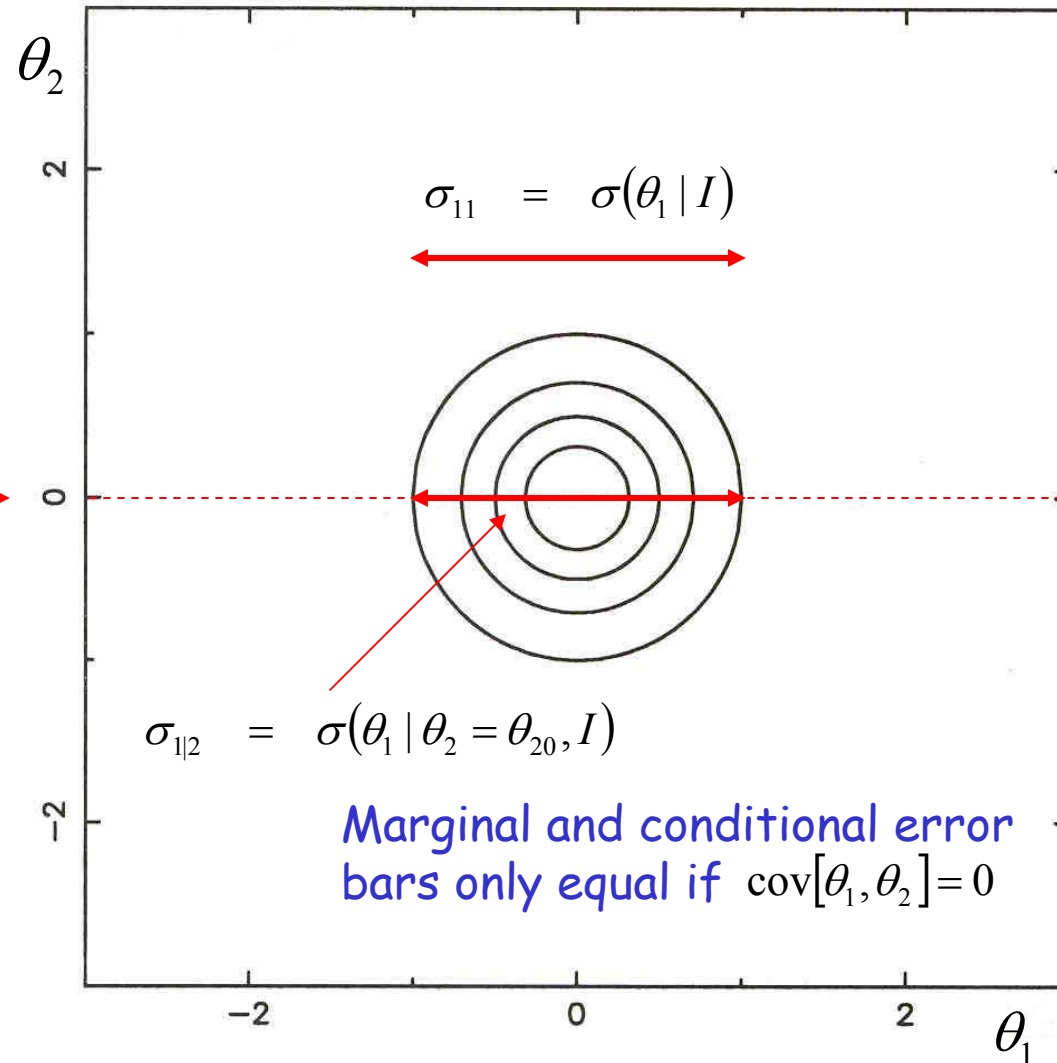
**B**  $\text{cov}[\theta_1, \theta_2] = 1$

**C**  $\text{cov}[\theta_1, \theta_2] = -1$

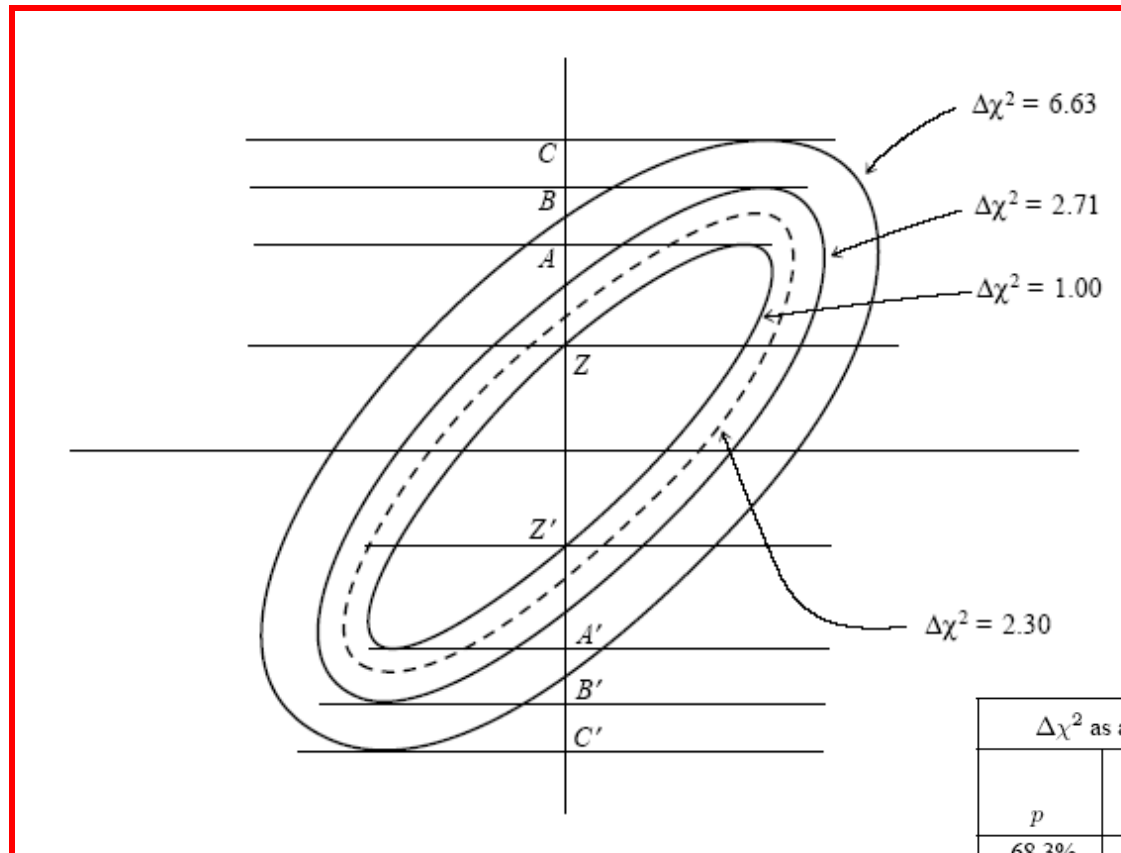
**D** None of the above

# Parameter estimation: 2-D case

'Best-fit' value of  $\theta_2$ , found from  $\frac{\partial \ell}{\partial \theta_j} \Big|_{\theta_j = \theta_{0j}} = 0$



# Parameter estimation: 2-D case

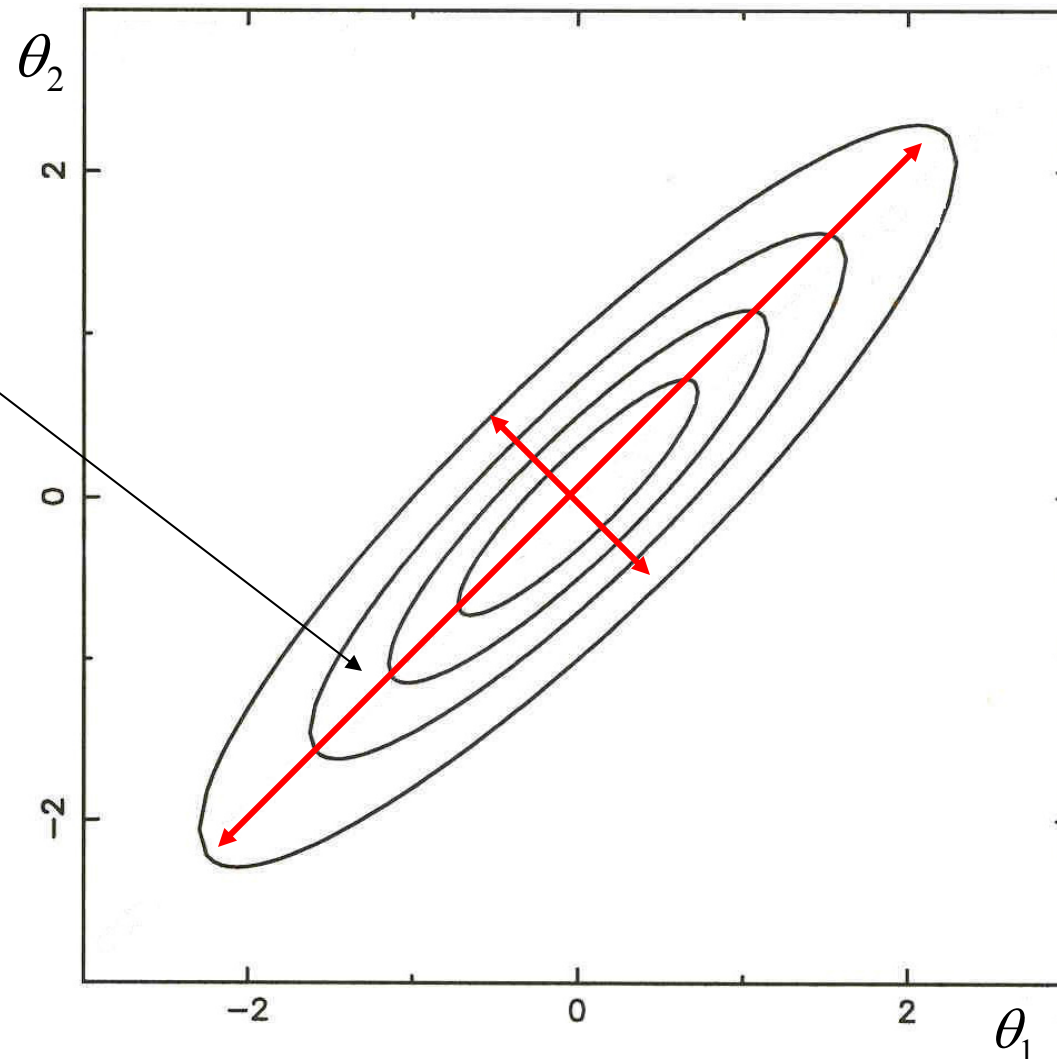


*From Numerical Recipes*

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



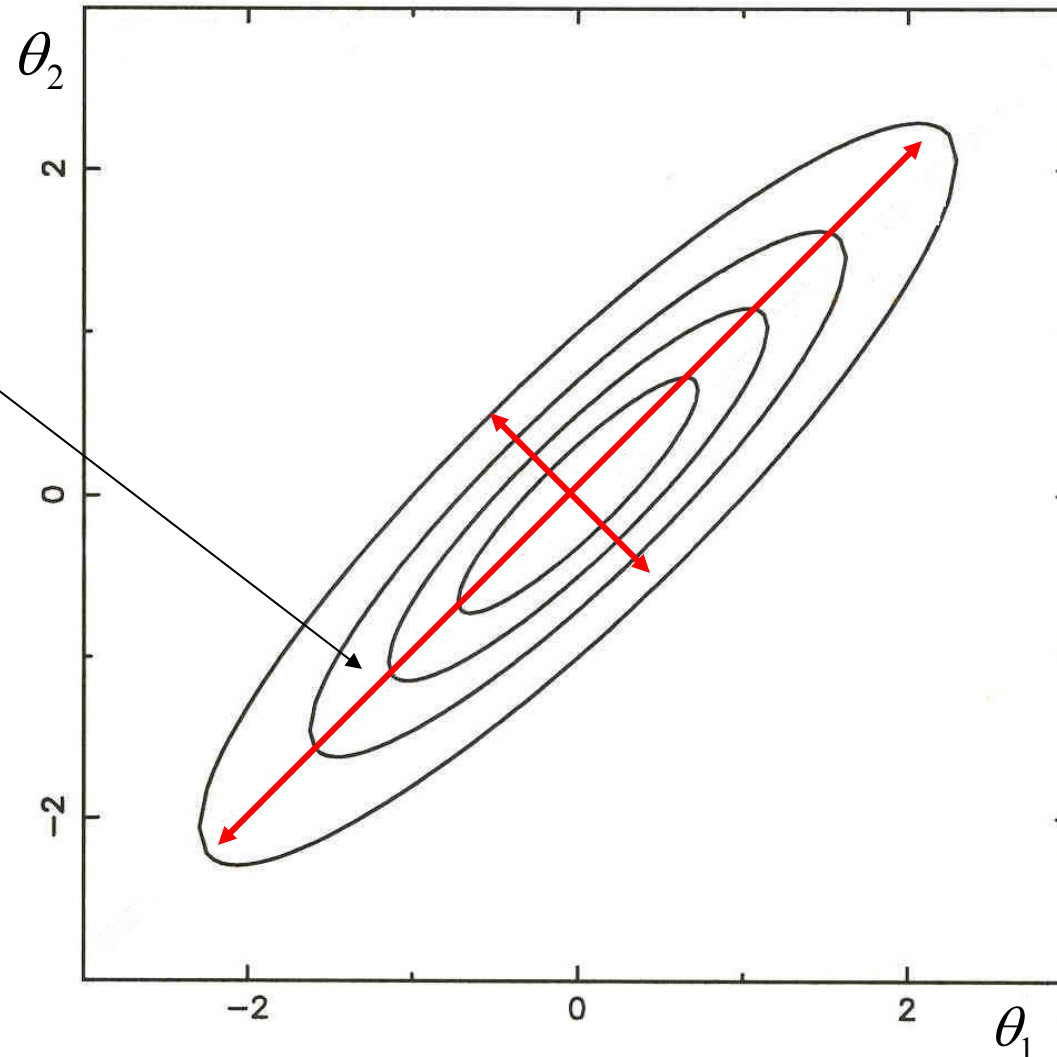
## Parameter estimation: 2-D case

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

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

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

$$\mathbf{F} \boldsymbol{\theta} = \lambda \boldsymbol{\theta}$$





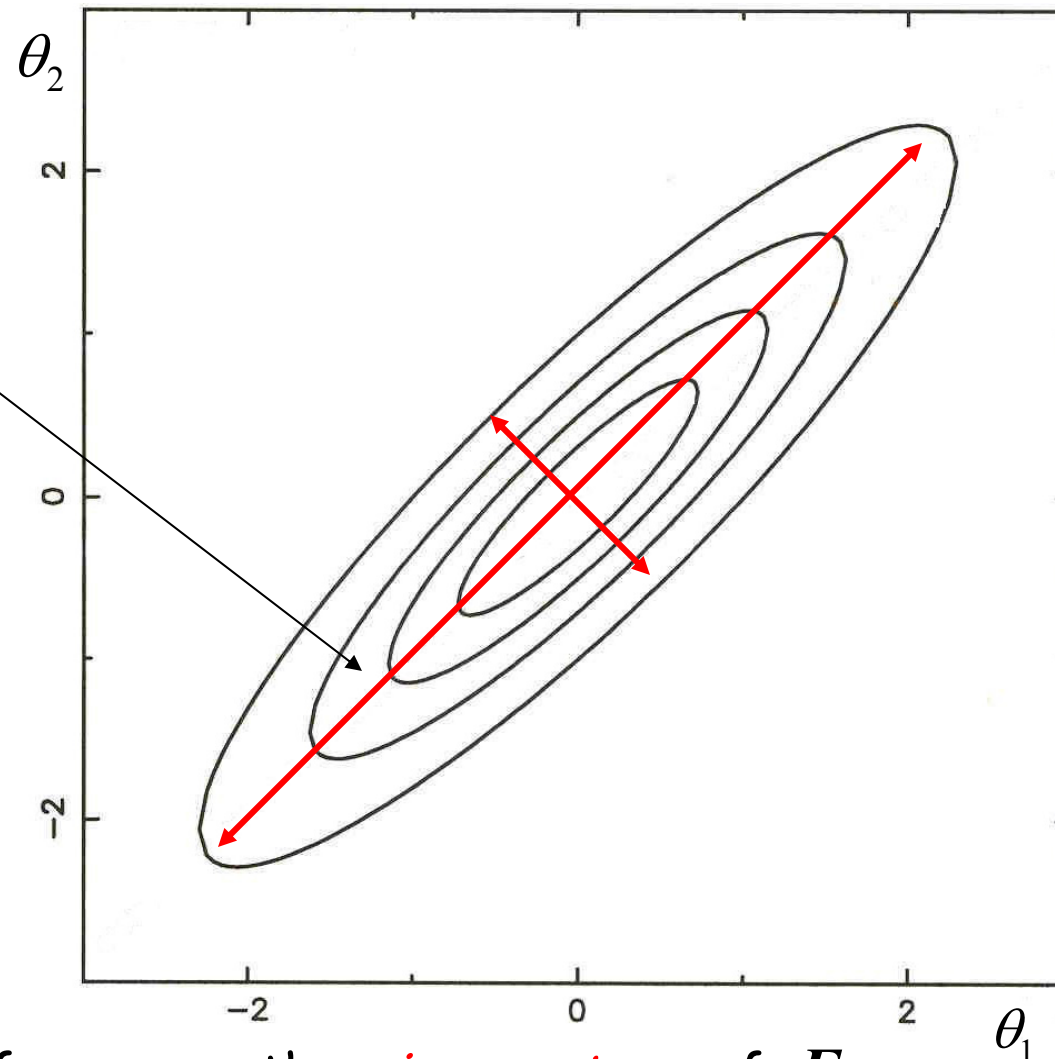
## Parameter estimation: 2-D case

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

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

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

$$F \theta = \lambda \theta$$



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

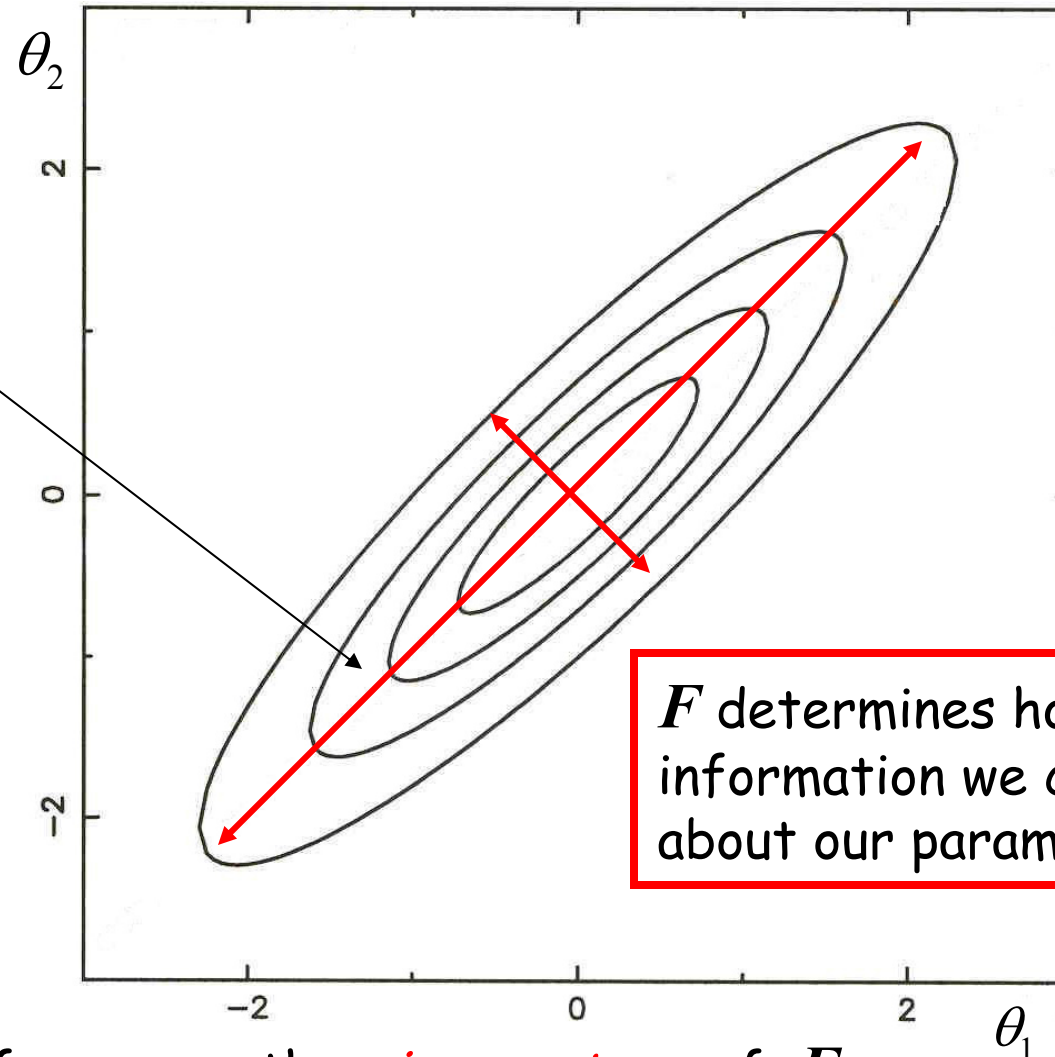
## Parameter estimation: 2-D case

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

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

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

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



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

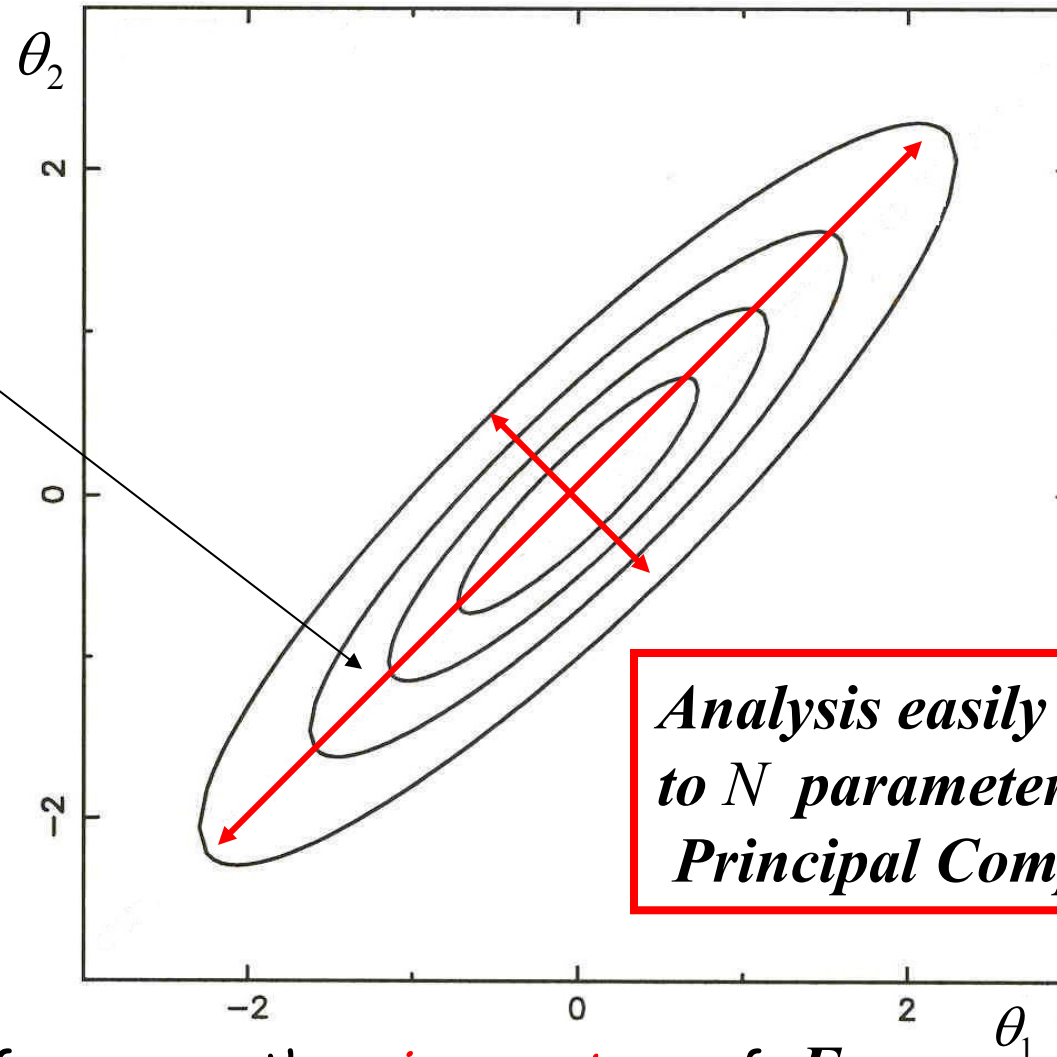
# Parameter estimation: N-dimensional case

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

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

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

$$F \theta = \lambda \theta$$



*Analysis easily extends to N parameters – Principal Components*

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

## What is PCA?

Principal Component Analysis:

*A method for transforming a multi-dimensional dataset, consisting of a number of statistically dependent (correlated) variables, into a set of uncorrelated variables: **principal components***

## What is PCA?

Principal Component Analysis:

**1<sup>st</sup> principal component** = *linear combination of the original variables that accounts for as much of the variation in the data as possible*

## What is PCA?

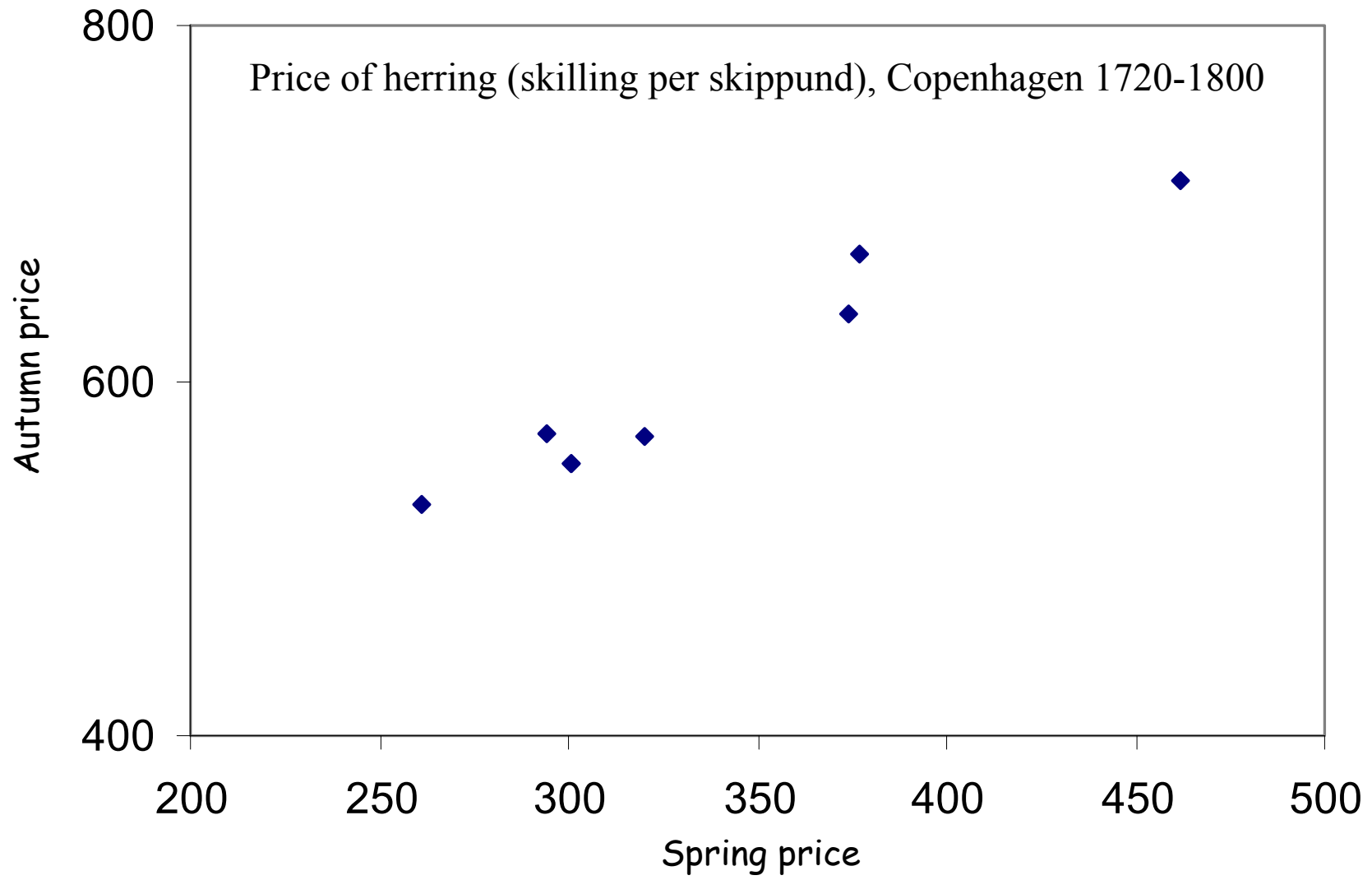
Principal Component Analysis:

**1<sup>st</sup> principal component** = *linear combination of the original variables that accounts for as much of the variation in the data as possible*

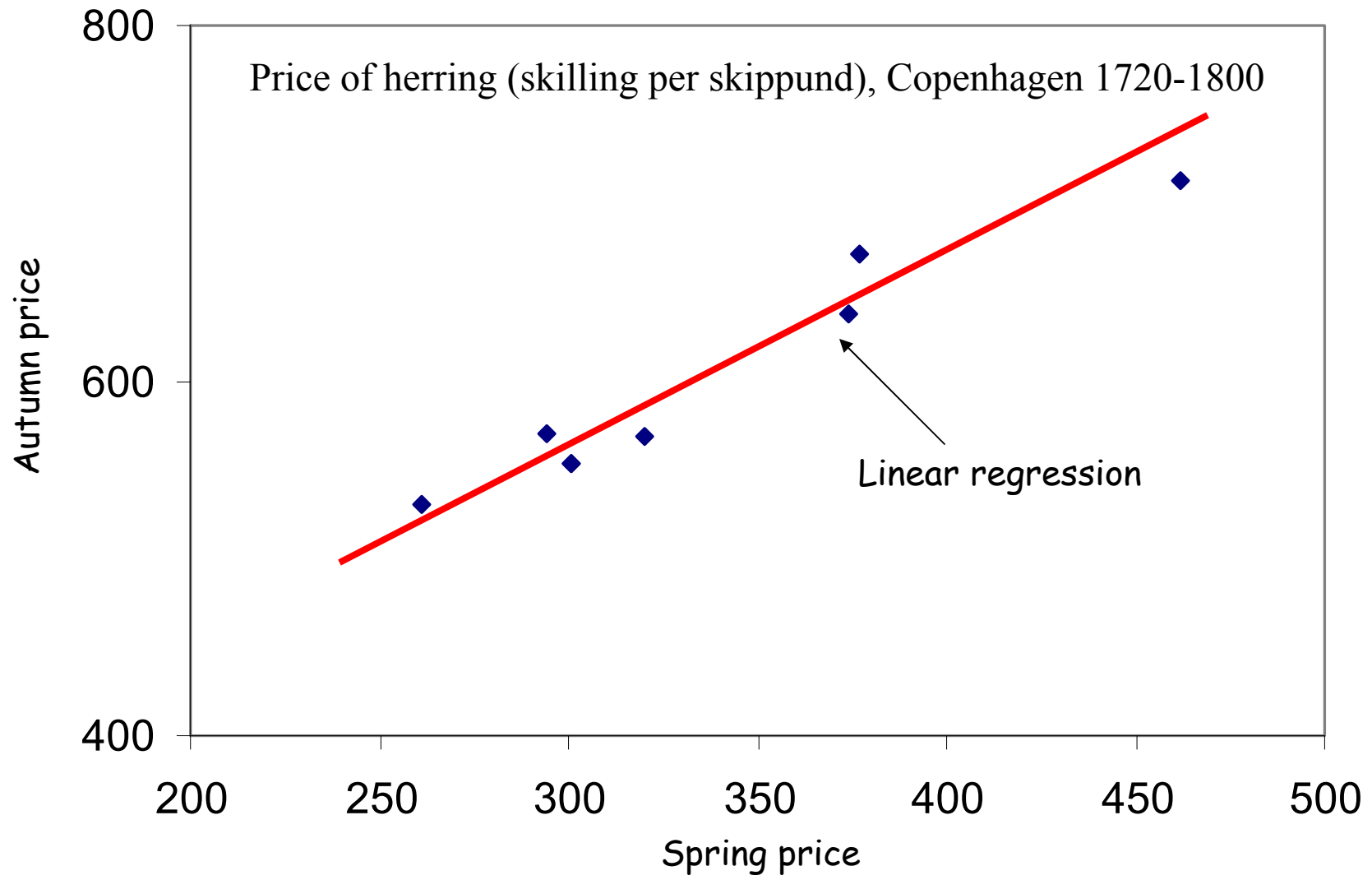
**2<sup>nd</sup> principal component** = *linear combination that accounts for as much of the remaining variation as possible, and is orthogonal to PC1*

⋮

# The price of fish...

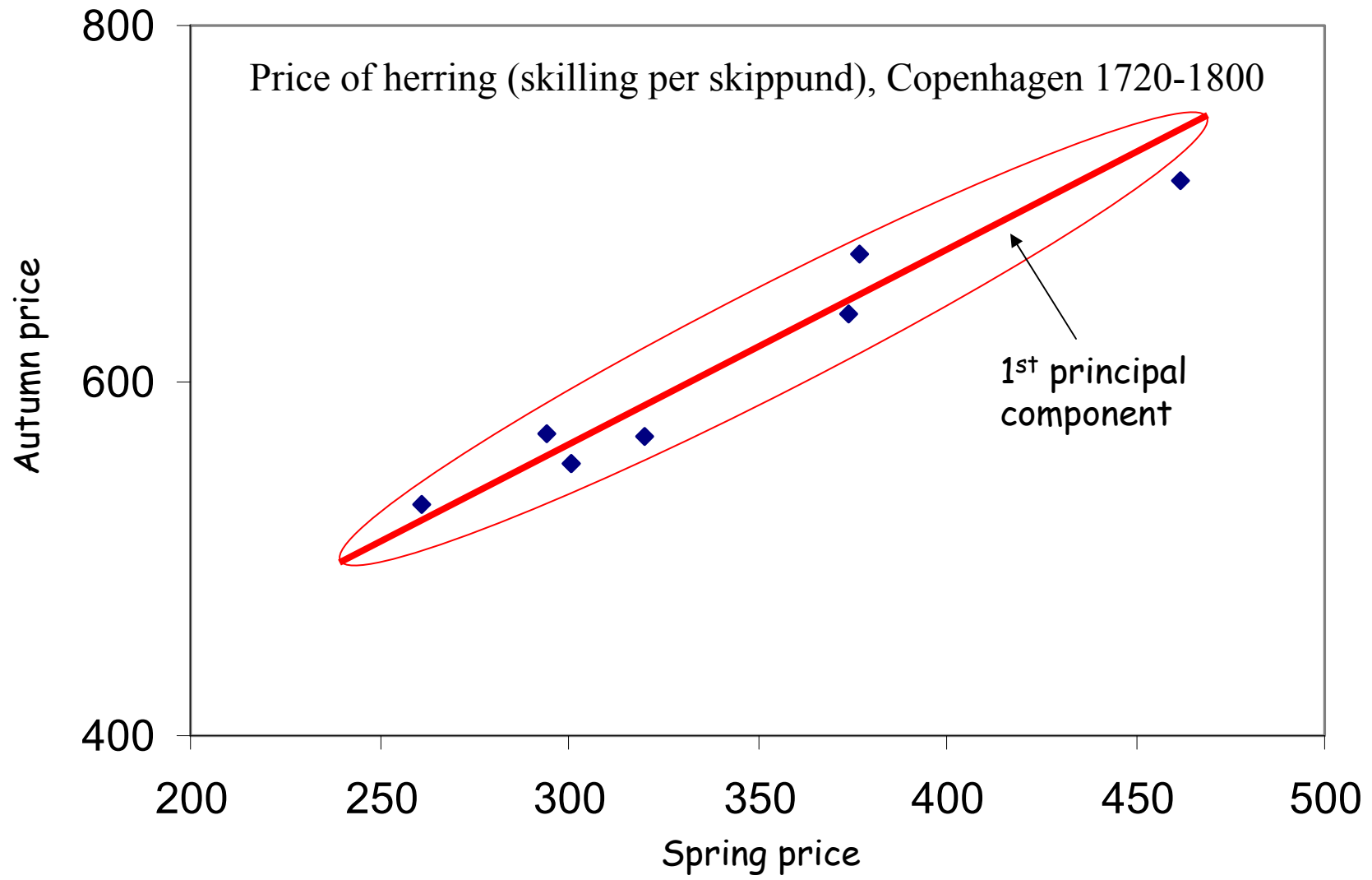


# The price of fish...

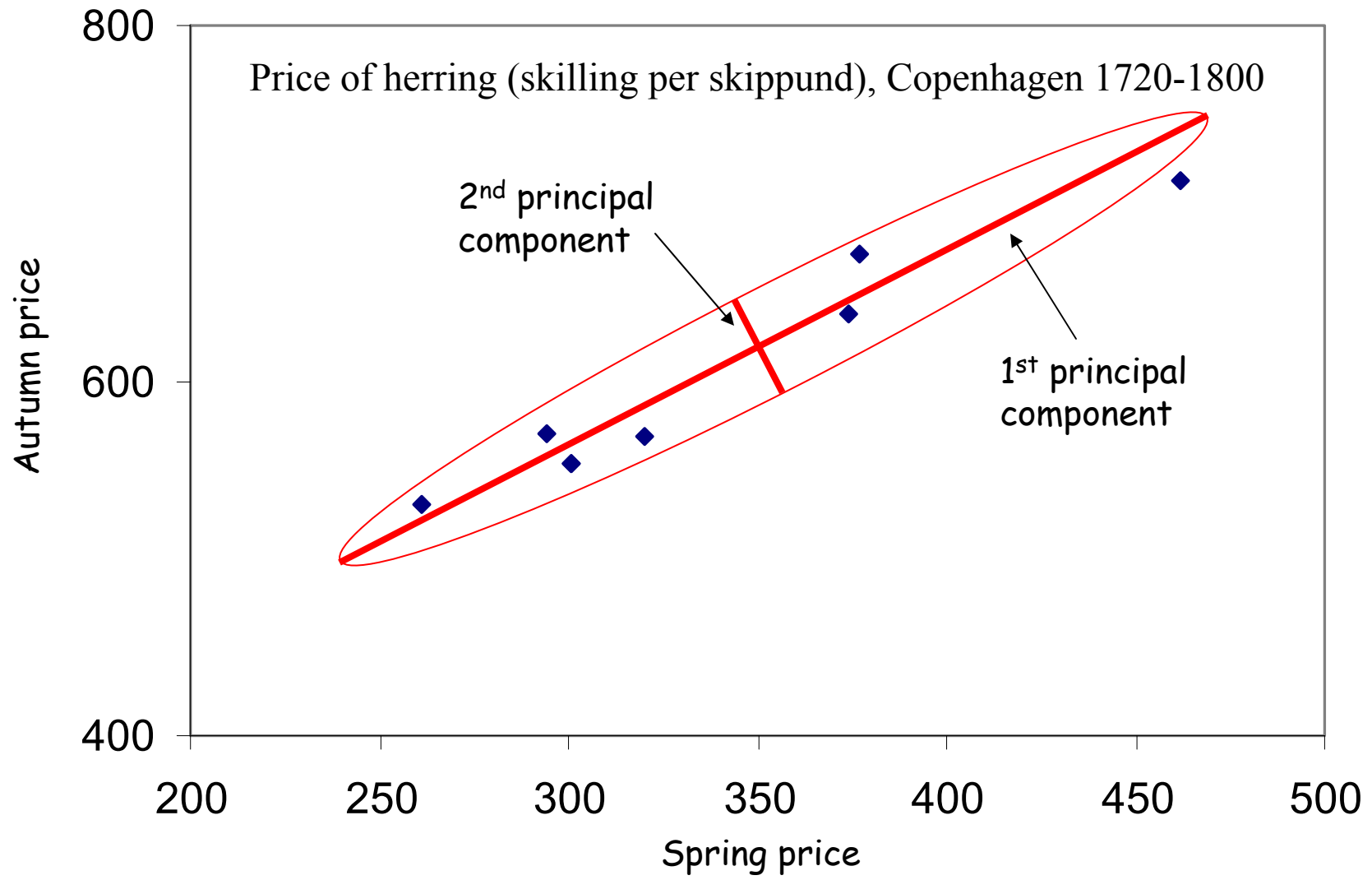




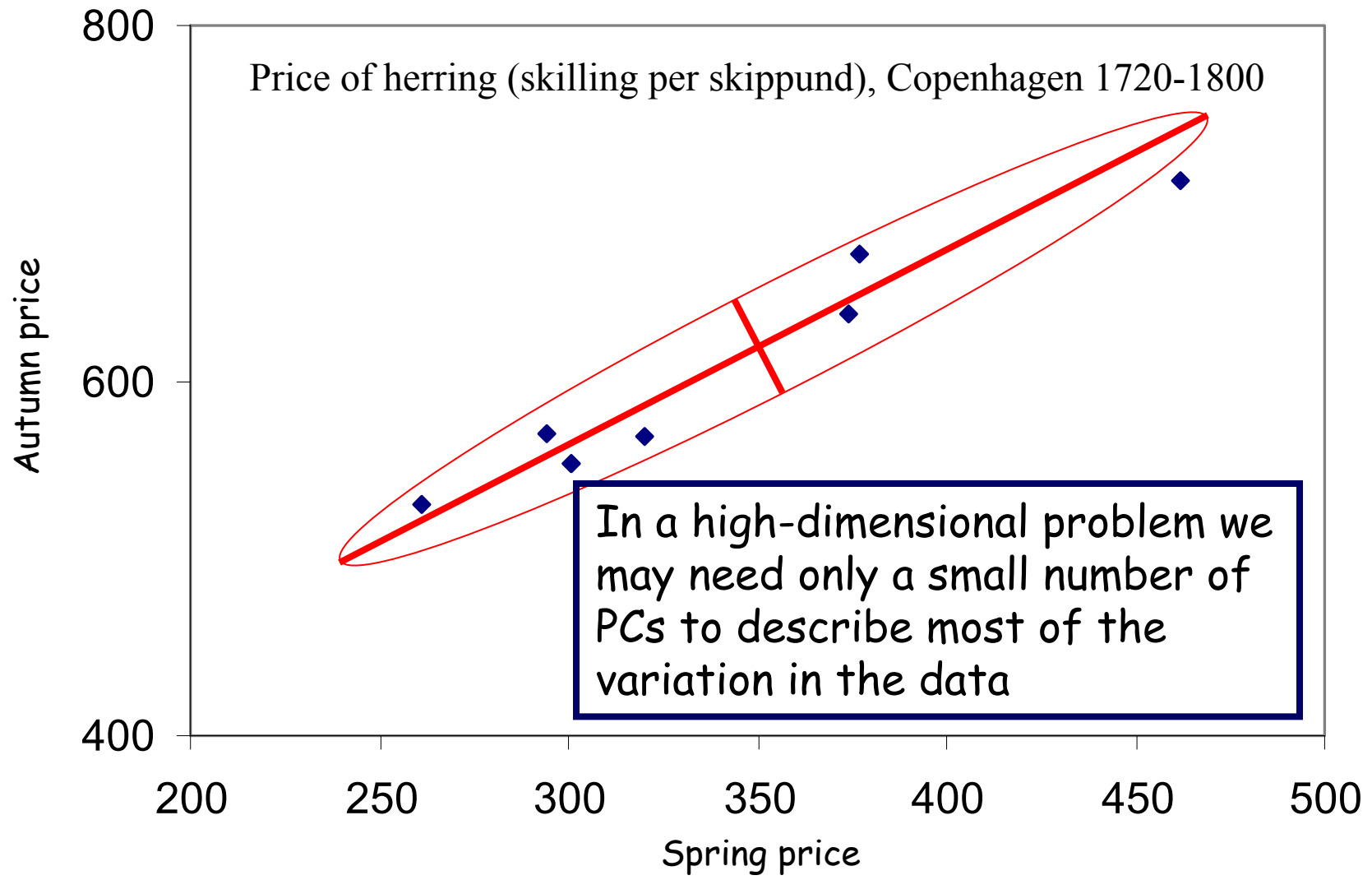
# The price of fish...



# The price of fish...



# The price of fish...



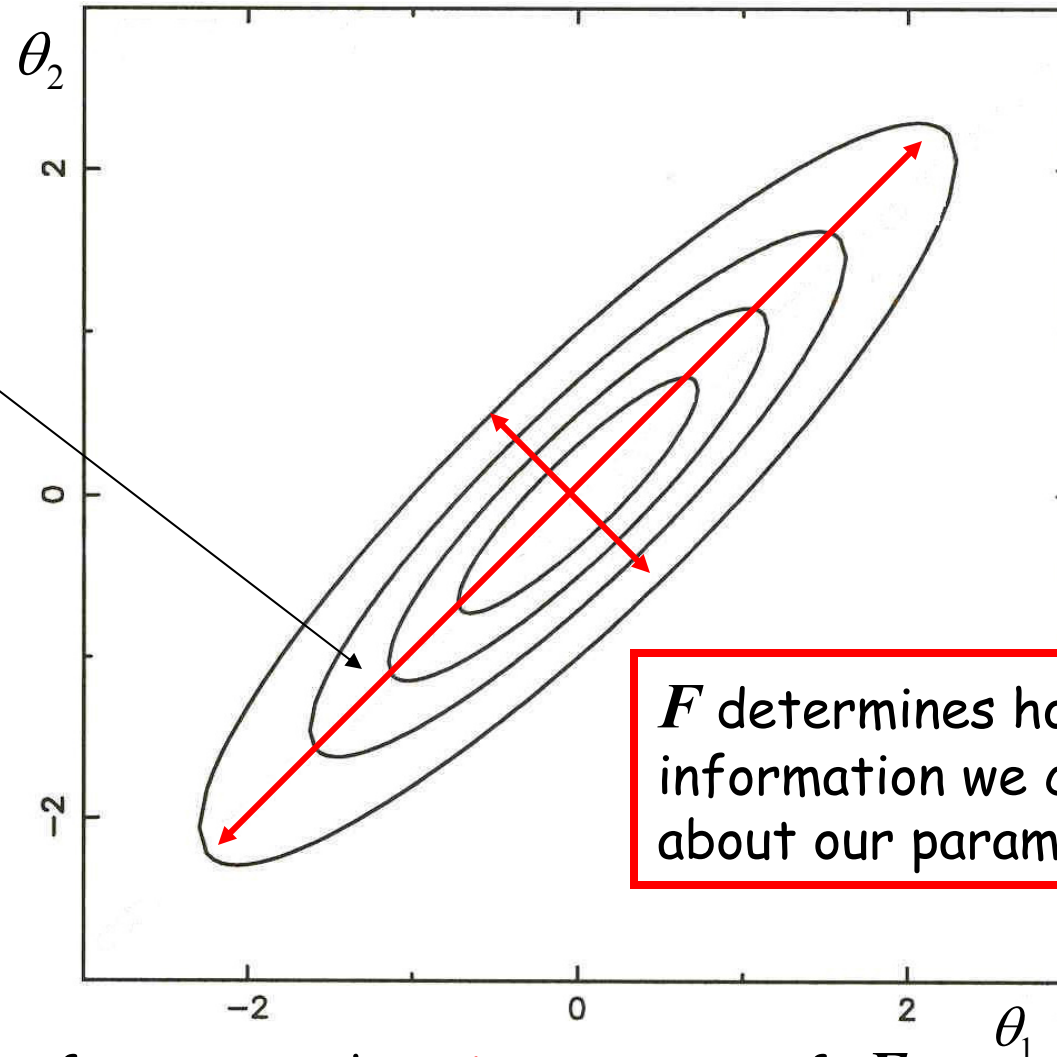
# Parameter estimation: 2-D case

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

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

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

$$F \theta = \lambda \theta$$



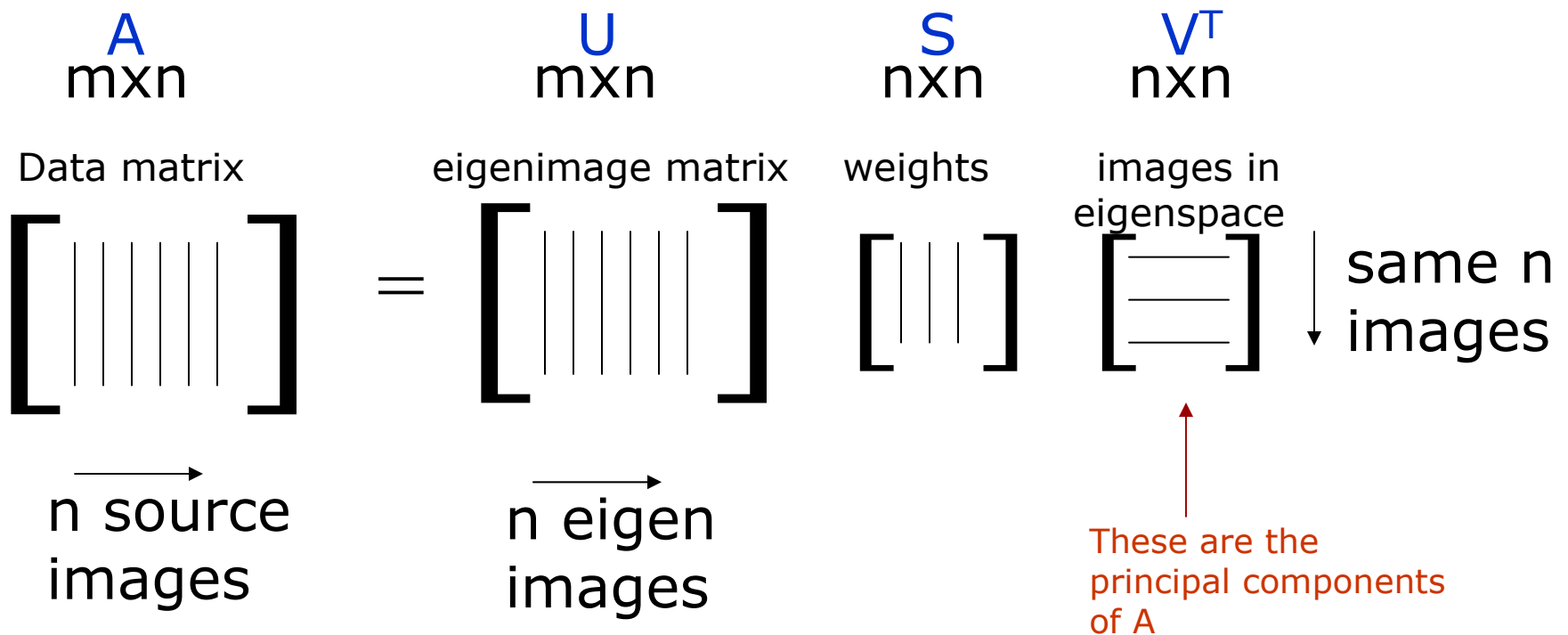
Direction of axes are the *eigenvectors* of  $F$

# More advanced PCA

- PCA can be applied to any number of datasets of any dimension
- Widely used for data compression and data characterisation by only considering the highest principal components, and throwing away those that contain little correlation information.
- E.g., “eigenfaces”...

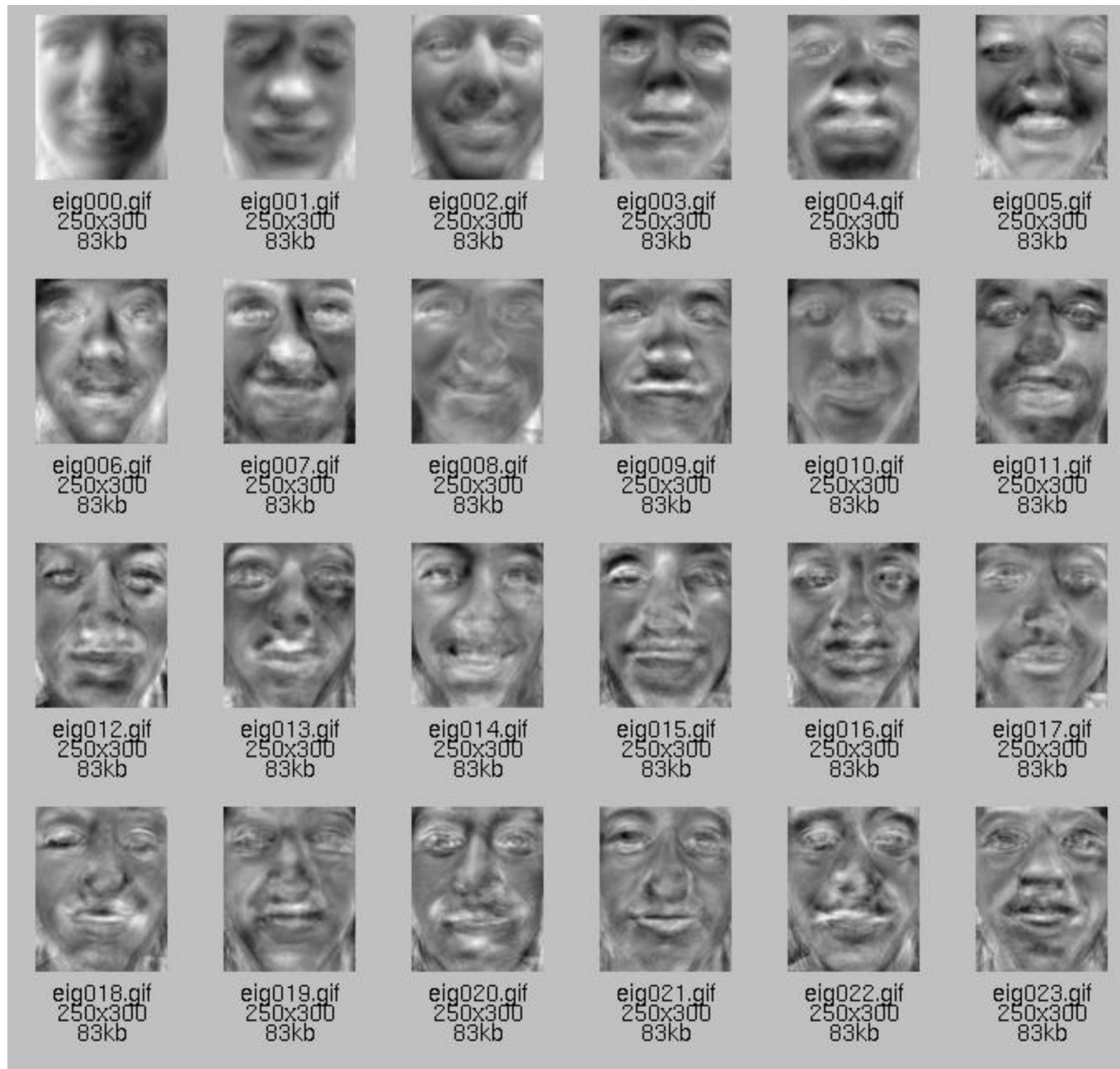
# Generating eigenvectors by SVD

- A useful way of generating the principal components is by carrying out a singular value decomposition (SVD) on the data matrix:



$$A = USV^T$$







## Application of Singular Value Decomposition to the Analysis of Time-Resolved Macromolecular X-Ray Data

Marius Schmidt<sup>\*,†</sup>, Sudarshan Rajagopal<sup>†</sup>, Zhong Ren<sup>†,‡,§</sup> and Keith Moffat<sup>†,‡</sup>

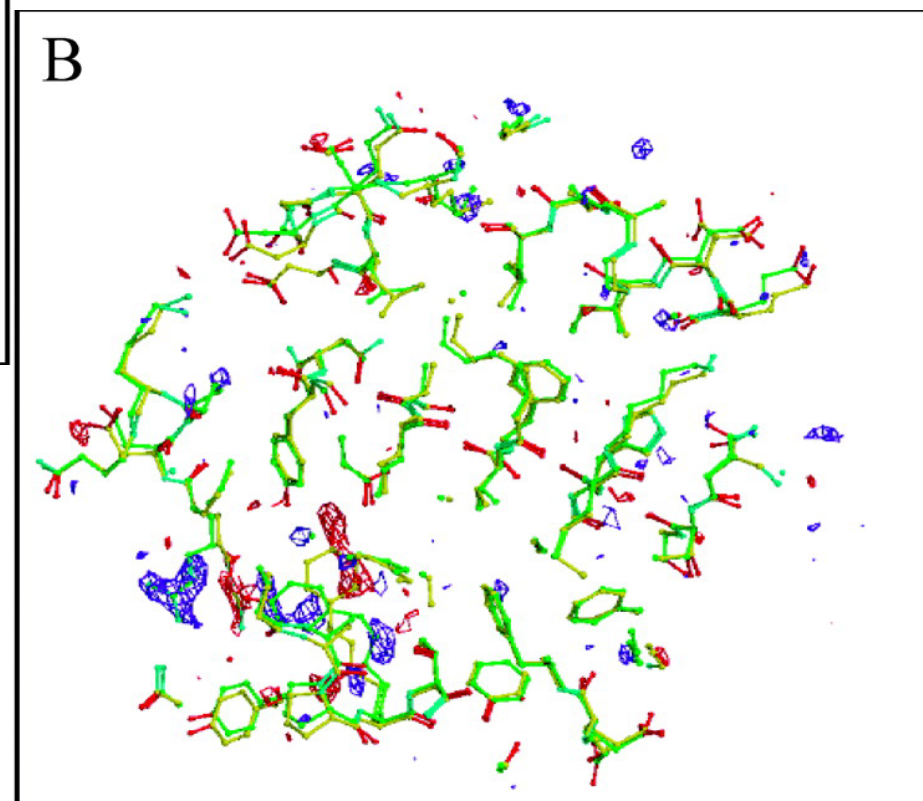
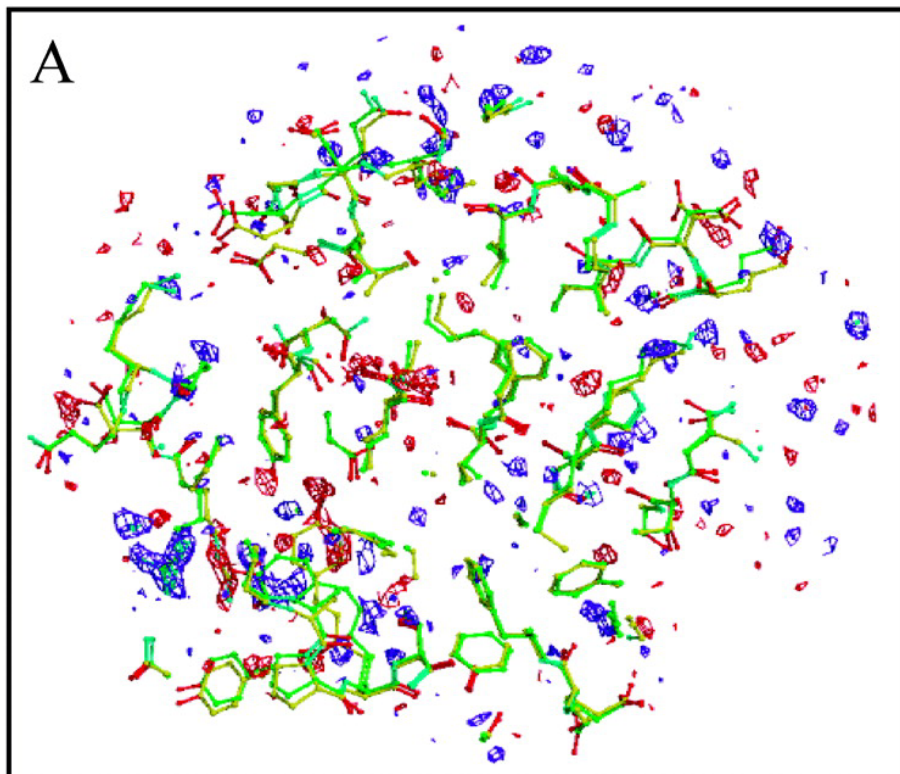
\* Physik-Department E17, Technische Universitaet Muenchen, 85747 Garching, Germany; <sup>†</sup> Department of Biochemistry and Molecular Biology, University of Chicago, Chicago, Illinois 60637 USA; <sup>‡</sup> BioCARS, Argonne National Laboratory, Argonne, Illinois 60439 USA; and <sup>§</sup> Renz Research, Des Plaines, Illinois 60018 USA

Correspondence: Address reprint requests to Marius Schmidt, E-mail: [marius@hexa.e17.physik.tu-muenchen.de](mailto:marius@hexa.e17.physik.tu-muenchen.de).

### ▶ ABSTRACT

Singular value decomposition (SVD) is a technique commonly used in the analysis of spectroscopic data that both acts as a noise filter and reduces the dimensionality of subsequent least-squares fits. To establish the applicability of SVD to crystallographic data, we applied SVD to calculated difference Fourier maps simulating those to be obtained in a time-resolved crystallographic study of photoactive yellow protein. The atomic structures of one dark state and three intermediates were used in

- ▲ [TOP](#)
- [ABSTRACT](#)
- ▼ [INTRODUCTION](#)
- ▼ [A GENERAL GUIDE TO...](#)
- ▼ [APPLICATION OF SVD TO...](#)
- ▼ [GENERATION OF THE MOCK...](#)
- ▼ [APPLICATION OF SVD TO...](#)
- ▼ [SVD AS A NOISE...](#)
- ▼ [EXTRACTION OF MECHANISM FROM...](#)
- ▼ [DISCUSSION](#)
- ▼ [ACKNOWLEDGEMENTS](#)



# Parameter estimation: Gaussian approximation

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

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

$$p(\theta_1, \theta_2 \mid \text{data}, I) = \exp[\ell(\theta_1, \theta_2)]$$
$$= \exp\left[-\frac{1}{2} Q\right] \quad \leftarrow \text{Gaussian approximation}$$

Is the Gaussian approximation a good idea?

# Parameter estimation: Gaussian approximation

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

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

Is the Gaussian approximation a good idea?

- o Greatly simplifies calculations - only need to compute the elements of the Fisher matrix (covariance matrix)

## Parameter estimation: Gaussian approximation

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

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

Is the Gaussian approximation a good idea?

- o Greatly simplifies calculations - only need to compute the elements of the Fisher matrix (covariance matrix)
- o Nowadays, however, we can compute **full posterior** pdf. Not too hard with present-day computers, even for large  $N$

## Parameter estimation: Gaussian approximation

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

$$\begin{aligned} \ell(\theta_1, \theta_2) = & \ell(\theta_{01}, \theta_{02}) + \frac{\partial \ell}{\partial \theta_1} \Big|_{\theta_j = \theta_{0j}} (\theta_1 - \theta_{01}) + \frac{\partial \ell}{\partial \theta_2} \Big|_{\theta_j = \theta_{0j}} (\theta_2 - \theta_{02}) + \\ & \frac{1}{2} \left[ \frac{\partial^2 \ell}{\partial \theta_1^2} \Big|_{\theta_j = \theta_{0j}} (\theta_1 - \theta_{01})^2 + \frac{\partial^2 \ell}{\partial \theta_2^2} \Big|_{\theta_j = \theta_{0j}} (\theta_2 - \theta_{02})^2 + 2 \frac{\partial^2 \ell}{\partial \theta_1 \partial \theta_2} \Big|_{\theta_j = \theta_{0j}} (\theta_1 - \theta_{01})(\theta_2 - \theta_{02}) \right] + \dots \end{aligned}$$

Is the Gaussian approximation a good idea?

- o Greatly simplifies calculations - only need to compute the elements of the Fisher matrix (covariance matrix)
- o Nowadays, however, we can compute **full posterior** pdf. Not too hard with present-day computers, even for large  $N$

*Markov Chain Monte Carlo Methods - see later*

# Beyond PCA: ICA

- PCA works by diagonalising the covariance matrix of a dataset, producing new linear combinations of the data which are uncorrelated.
- If the data are **Gaussian** then PCA will produce combinations which are statistically independent (all higher moments are zero for a Gaussian).
- For **non-Gaussian** data, uncorrelated does not in general imply independent. New method called **independent component analysis**: linear transformation of data vector to minimise statistical dependence of components.

Interesting applications to time series data



## COCKTAIL PARTY PROBLEM

Imagine you're at a cocktail party. For you it is no problem to follow the discussion of your neighbours, even if there are lots of other sound sources in the room: other discussions in English and in other languages, different kinds of music, etc.. You might even hear a siren from the passing-by police car.

It is not known exactly how humans are able to separate the different sound sources. **Independent component analysis** is able to do it, if there are at least as many microphones or 'ears' in the room as there are different simultaneous sound sources. In this demo, you can select which sounds are present in your cocktail party. ICA will separate them without knowing anything about the different sound sources or the positions of the microphones.

### ORIGINAL SOUND SOURCES

By clicking the icons you can listen to the original [sound sources](#).



### SAMPLES AT THE COCKTAIL PARTY

Listen to the mixtures by clicking the microphones.



### FOUND SOUND SOURCES

Below are the the sound sources separated by ICA. Note that they might be in different order than the original ones.



[http://www.cis.hut.fi/projects/ica/cocktail/cocktail\\_en.cgi](http://www.cis.hut.fi/projects/ica/cocktail/cocktail_en.cgi)

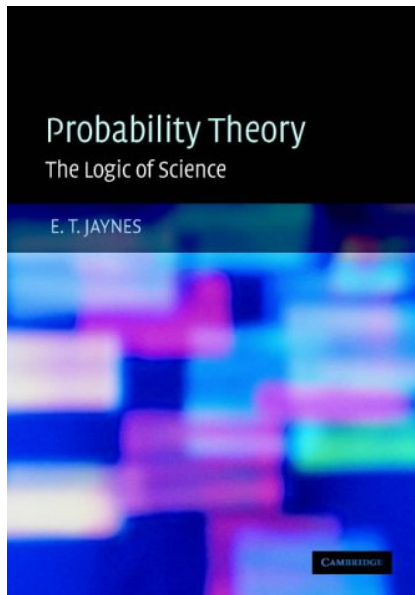


## 4b. Defining Probabilities



# Bayesian versus Frequentist statistics: Who is right?

Frequentists are correct to worry about subjectiveness of assigning probabilities - Bayesians worry about this too!!!



Ed Jaynes  
(1922 - 1998)

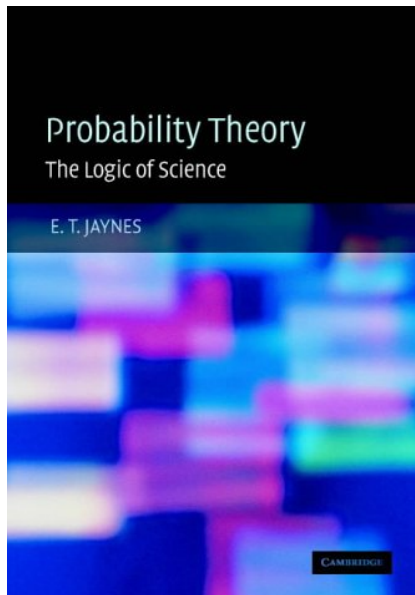
Probability *is* subjective;  
it depends on the available  
information

**Subjective  $\neq$  arbitrary**

Given the *same* background  
information, two observers should  
assign the *same* probabilities

# Bayesian versus Frequentist statistics: Who is right?

Frequentists are correct to worry about subjectiveness of assigning probabilities - Bayesians worry about this too!!!



Ed Jaynes  
(1922 - 1998)

Probability *is* subjective;  
it depends on the available  
information

**Subjective  $\neq$  arbitrary**

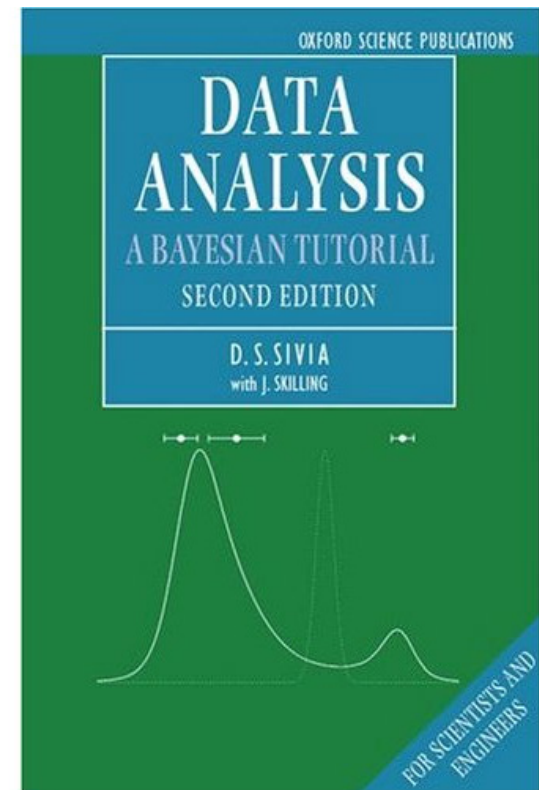
Given the *same* background  
information, two observers should  
assign the *same* probabilities

*But what should they be?...*

Bernoulli (1713) 'Principle of insufficient reason'

Keynes (1921) 'Principle of indifference'

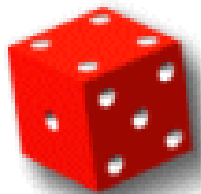
If we can enumerate a set of basic mutually exclusive possibilities, and we have no reason to believe that any one of these is more likely to be true than another, then we should assign the same probability to all.



Bernoulli (1713) 'Principle of insufficient reason'

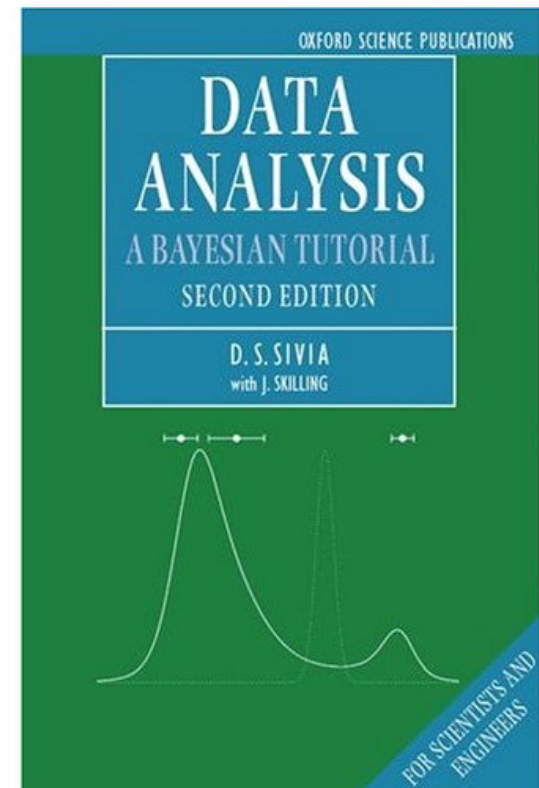
Keynes (1921) 'Principle of indifference'

If we can enumerate a set of basic mutually exclusive possibilities, and we have no reason to believe that any one of these is more likely to be true than another, then we should assign the same probability to all.



$X_i \equiv$  face on top has  $i$  dots

$$p(X_i | I) = \frac{1}{6} \quad \text{for all } i$$

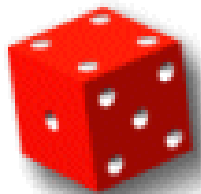




Bernoulli (1713) 'Principle of insufficient reason'

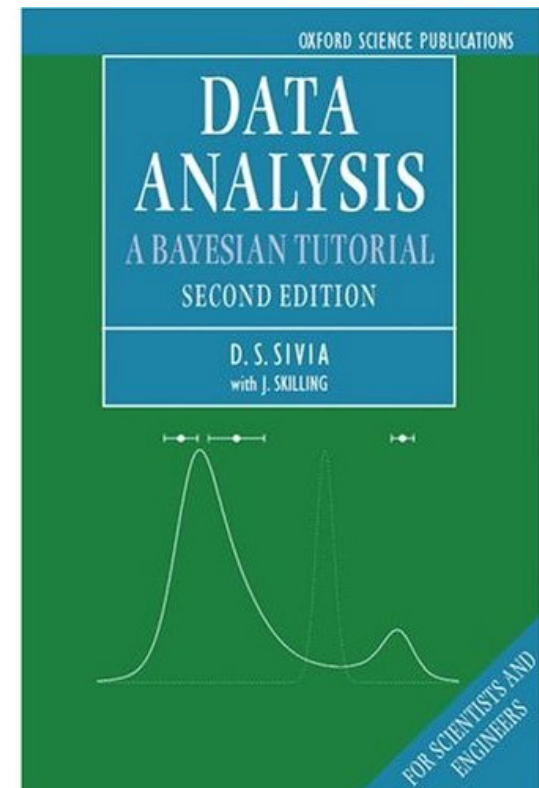
Keynes (1921) 'Principle of indifference'

If we can enumerate a set of basic mutually exclusive possibilities, and we have no reason to believe that any one of these is more likely to be true than another, then we should assign the same probability to all.



$X_i \equiv$  face on top has  $i$  dots

$$p(X_i | I) = \frac{1}{6} \quad \text{for all } i$$

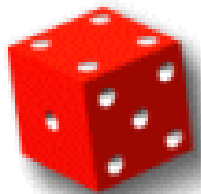


Agrees with common sense, but can we justify more fundamentally?

Bernoulli (1713) 'Principle of insufficient reason'

Keynes (1921) 'Principle of indifference'

If we can enumerate a set of basic mutually exclusive possibilities, and we have no reason to believe that any one of these is more likely to be true than another, then we should assign the same probability to all.



$X_i$  are just labels, e.g. suppose we define

$X_i \equiv$  face on top has  $7 - i$  dots

Should still have  $p(X_i | I) = \frac{1}{6}$  for all  $i$

Extending to continuum case,

Let  $x$  be a **location parameter**.

Principle of indifference means we should have

$$p(x | I)dx = p(x + \Delta | I)d(x + \Delta)$$

where  $\Delta$  is a constant

Since  $dx = d(x + \Delta)$  we must have

$$p(x | I) = \text{constant}$$



Similarly,

Let  $L$  be a **scale parameter**.

Principle of indifference means we should have

$$p(L | I)dL = p(\beta L | I)d(\beta L)$$

where  $\beta$  is a positive constant

Since  $d(\beta L) = \beta dL$  we must have

$$p(L | I) \propto 1/L$$

Jeffreys' prior

A Jeffreys' prior represents complete ignorance about the value of a scale parameter.

It is equivalent to a uniform pdf for the logarithm of  $L$

i.e.

$$p(\log L | I)dL = \text{constant}$$

A Jeffreys' prior represents complete ignorance about the value of a scale parameter.

It is equivalent to a uniform pdf for the logarithm of  $L$

i.e.  $p(\log L | I)dL = \text{constant}$

In fact what is referred to as a Jeffreys prior  $p(L | I) \propto 1/L$  is just the special case of a more general result.

Suppose our inference problem is described by a likelihood with parameter(s)  $\vec{\theta}$ .

The **Jeffreys prior** is a non-informative (objective) prior defined as:

$$p(\vec{\theta}) \propto [\det I(\vec{\theta})]^{1/2}$$

Here  $I(\vec{\theta})$  is the **Fisher Information** defined as

$$I(\vec{\theta})_{i,j} = E \left[ \frac{\partial}{\partial \theta_i} \ln L(\vec{\theta}) \frac{\partial}{\partial \theta_j} \ln L(\vec{\theta}) \right]$$

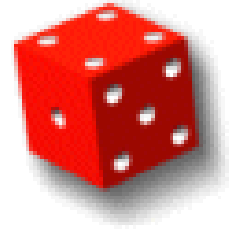
[ Note this expression reduces to that for the Fisher matrix given in Section 4a for the special case of a Gaussian likelihood. ]

Key feature: the Jeffreys prior is **invariant** under *any* re-parameterisation of  $\vec{\theta}$

## Testable information

How do we deal with more complicated situations?

e.g. suppose we know that, when our die was rolled many times, the average result was 4.5 (and not 3.5)



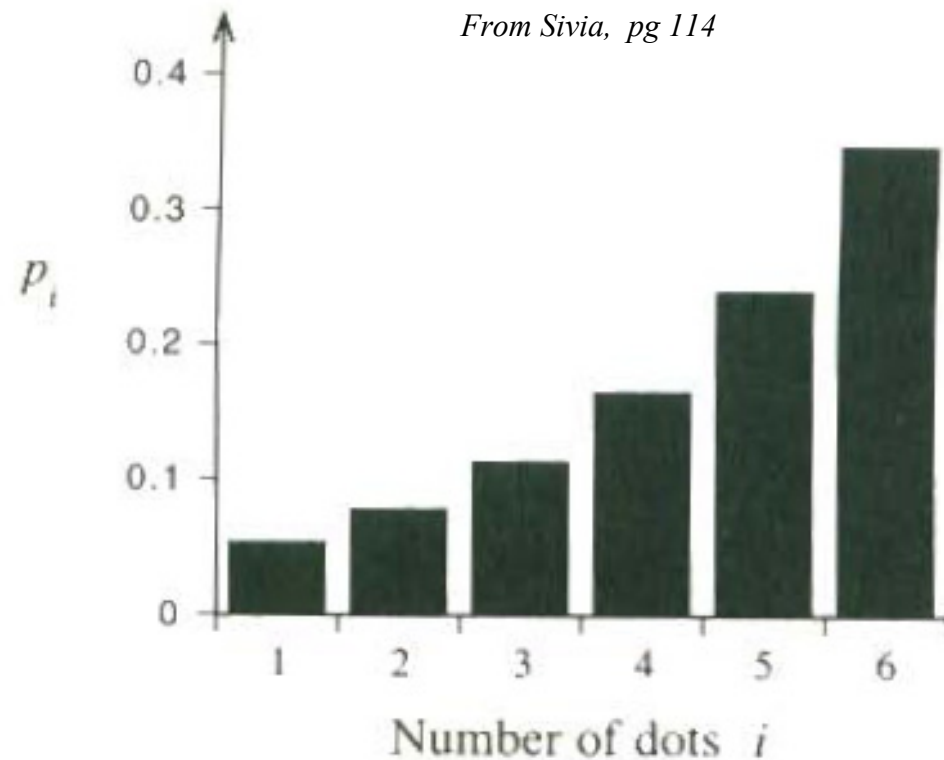
How do we use this information to constrain  $p(X_i | I)$  ?

Jaynes (1957) suggests maximising the **Entropy**

$$S = -\sum_{i=1}^6 p_i \log[p_i] \quad \text{subject to} \quad \sum_{i=1}^6 p_i = 1 \quad \text{and} \quad \sum_{i=1}^6 i p_i = 4.5$$

We can solve for the  $p_i$  using Lagrange Multipliers.

But why MAXENT?



We can justify the importance of MAXENT via two approaches:

- 1) Independence argument (the kangaroo problem)
- 2) Shannon's Theorem and multiplicity (see notes on website)

## Consider the Kangaroo problem!

*Information:*      $1/3$  of all kangaroos have blue eyes;  
                          $1/3$  of all kangaroos are left-handed

*Question:*        *On the basis of the above information alone,  
                         what proportion of kangaroos are both blue  
                         eyed and left-handed?*

**Question 11:** Assuming that eye-colour and handedness are uncorrelated for kangaroos (humans?), we expect the proportion of kangaroos that are both blue-eyed and left-handed to be:

**A** zero

**B** 100%

**C**  $1/9$

**D**  $1/3$





**Question 11:** Assuming that eye-colour and handedness are uncorrelated for kangaroos (humans?), we expect the proportion of kangaroos that are both blue-eyed and left-handed to be:

**A** zero

**B** 100%

**C**  $1/9$

**D**  $1/3$

## Consider the Kangaroo problem!

*Information: 1/3 of all kangaroos have blue eyes;  
1/3 of all kangaroos are left-handed*

*Question: On the basis of the above information alone,  
what proportion of kangaroos are both blue eyed and left-handed?*

		Left-Handed	
		True	False
Blue-Eyed	True	$p_1$	$p_2$
	False	$p_3$	$p_4$

		Left-Handed	
		True	False
Blue-Eyed	True	$0 \leq x \leq \frac{1}{3}$	$\frac{1}{3} - x$
	False	$\frac{1}{3} - x$	$\frac{1}{3} + x$

We know that:  $p_1 + p_2 + p_3 + p_4 = 1$

$$p_1 + p_2 = 1/3 \quad p_1 + p_3 = 1/3$$

What is  $x$  ?

Independence arguments favour  $x = 1/9$

		<i>Left-Handed</i>	
		True	False
<i>Blue-Eyed</i>	True	$p_1$	$p_2$
	False	$p_3$	$p_4$

		<i>Left-Handed</i>	
		True	False
<i>Blue-Eyed</i>	True	$0 \leq x \leq \frac{1}{3}$	$\frac{1}{3} - x$
	False	$\frac{1}{3} - x$	$\frac{1}{3} + x$

We know that:  $p_1 + p_2 + p_3 + p_4 = 1$

$$p_1 + p_2 = 1/3 \quad p_1 + p_3 = 1/3$$

What is  $x$  ?

Independence arguments favour  $x = 1/9$

Variational function	Optimal $x$	Implied correlation
<b>MAXENT</b> $\rightarrow -\sum p_i \log_c(p_i)$	0.1111	None
$-\sum p_i^2$	0.0833	Negative
$\sum \log_c(p_i)$	0.1301	Positive
$\sum \sqrt{p_i}$	0.1218	Positive

## MAXENT and common pdfs

Suppose we only know the expected value,  $\mu$ , of a continuous physical quantity,  $x$

What should we assign as  $p(x | I)$  ?

Using MAXENT it can be shown that

$$p(x | \mu) = \frac{1}{\mu} \exp\left(-\frac{x}{\mu}\right)$$

Exponential  
distribution

## MAXENT and common pdfs

Suppose we only know the expected value,  $\mu$ , of a **discrete** physical quantity,  $N$

What should we assign as  $p(x | I)$  ?

Using MAXENT it can be shown that

$$p(N | \mu) = \frac{\mu^N e^{-\mu}}{N!}$$

Poisson  
distribution

## MAXENT and common pdfs

Suppose we only know the expected value,  $\mu$ , and  $\langle x - \mu \rangle^2 = \sigma^2$  of a continuous physical quantity,  $x$

What should we assign as  $p(x | I)$  ?

Using MAXENT it can be shown that

$$p(x | \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

Normal  
distribution

## MAXENT and common pdfs

Suppose we only know the expected value,  $\mu$ , and  $\langle x - \mu \rangle^2 = \sigma^2$  of a continuous physical quantity,  $x$

What should we assign as  $p(x | I)$  ?

Using MAXENT it can be shown that

$$p(x | \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

Normal  
distribution

**MAXENT justifies the relevance of common pdfs**