

## Problems of classical statistics

- Statistics seems a disconnected series of receipts...

- Central limit theorem

Statistic is a theory from which if you have two cars and I have none, each has one car

- And therefore is “firmly” based on gaussian distribution

- And therefore the “solutions” are also “gaussians”

Let's stick to probability

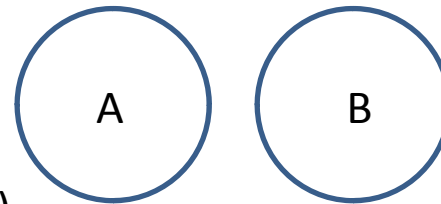
## Some basics of probability

A random variable A has a probability  $\text{prob}(X)$  of taking place between 0 and 1, being 1 when the event is sure.

If A and B are exclusive events then:

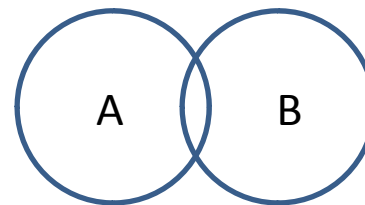
- $\text{prob}(X \text{ or } Y) = \text{prob}(X) + \text{prob}(Y)$
- $\text{prob}(X \text{ and } Y) \equiv \text{prob}(X, Y) = \text{prob}(A) \cdot \text{prob}(B)$
- $\text{prob}(X \text{ knowing } Y) \equiv \text{prob}(X | Y) = \text{prob}(X)$

(since the information of B does not affect to our knowledge of A, since they are independent)



If A and B are not exclusive events then:

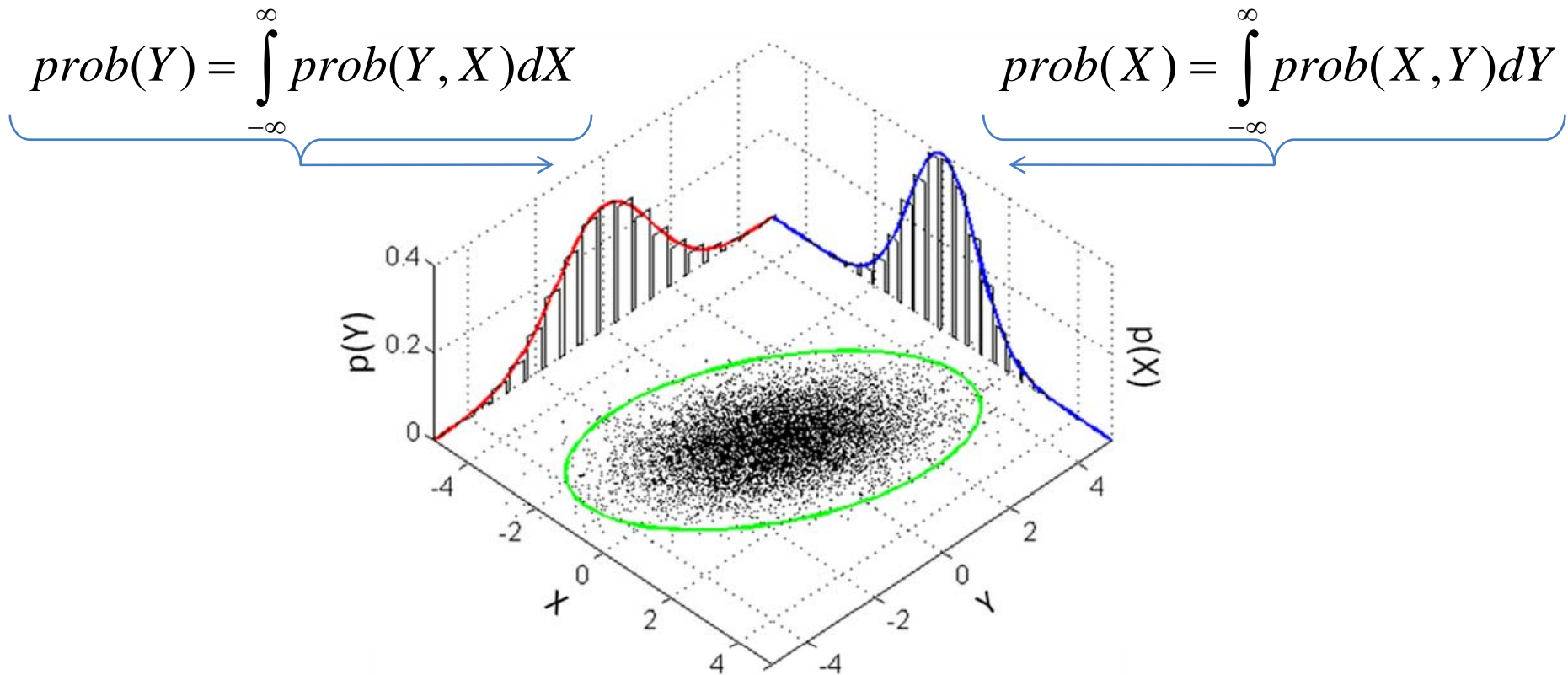
$$\text{prob}(X | Y) = \frac{\text{prob}(X, Y)}{\text{prob}(Y)}$$



## What can be do with Probability

### MARGINALIZATION

We want to know something about  $X$ , no matter what happens to  $Y$   
in other words  $Y$  represents a nuisance parameter... in the sense we are not interested in now!




It allows to “project” a N-Dimensional PDF in the variable we are interested in

## What can be do with Probability

### MARGINALIZATION

We want to know something about  $X$ , no matter what happens to  $Y$   
in other words  $Y$  represents a nuisance parameter... in the sense we are not interested in now!

$$prob(Y) = \int_{-\infty}^{\infty} prob(Y, X) dX = \int_{-\infty}^{\infty} prob(Y | X) prob(X) dX$$


$$prob(X | Y) = \frac{prob(X, Y)}{prob(Y)}$$

... for example t-student distribution can be obtained from the marginalization of  $\sigma$   
let  $D$  be the data

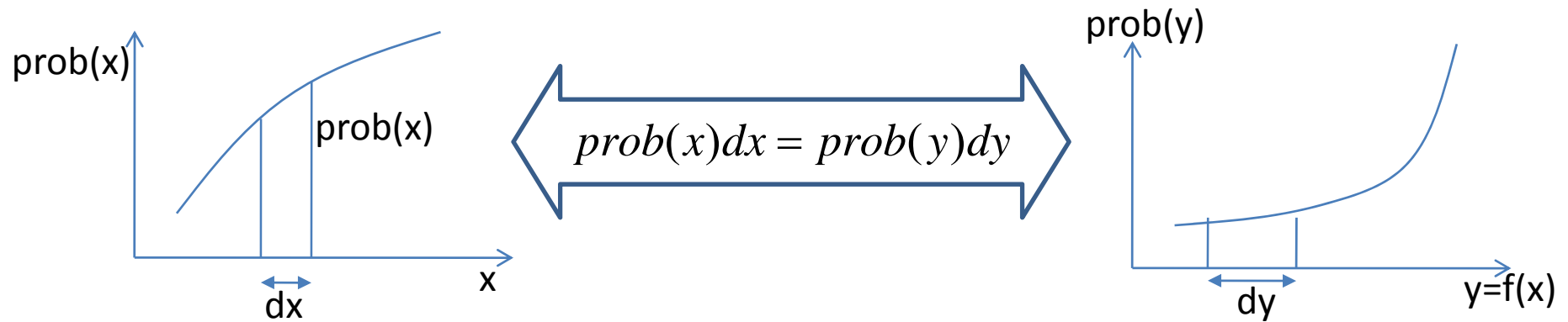
$$prob(\mu | D) = \int_{-\infty}^{\infty} prob(\mu, \sigma^2 | D) d\sigma^2 = \int_{-\infty}^{\infty} prob(\mu | D, \sigma^2) prob(\sigma^2 | D) d\sigma^2$$

## What can be do with Probability

### CHANGE OF VARIABLES

We want to know the PDF from a function of the original variable.

GOAL: after the change of variables the probability MUST NOT change



$$\text{prob}(x) = \text{prob}(y) \frac{dy}{dx}$$

## What can be do with Probability

### CHANGE OF VARIABLES

We want to know the PDF from a function of the original variable.

#### EXAMPLE

We have no idea about the “length scale”  $L$  of a problem .

It is reasonable that a maximum ignorance PDF is constant in log scale

$$prob(\ln(L)) = cte$$

What is the probability of  $L$

$$prob(L) = \frac{d(\ln(L))}{dL} prob(\ln(L)) \propto \frac{1}{L}$$

This is known as a Jeffrey's PDF (prior as we will see shortly)

## What can be do with Probability

### CHANGE OF VARIABLES

We want to know the PDF from a function of the original variable.

#### EXAMPLE

We have no idea about the “length scale”  $L$  of a problem .

It is reasonable that a maximum ignorance PDF is constant in log scale

$$prob(\ln(L)) = cte$$

What is the probability of  $L$

$$prob(L) = \frac{d(\ln(L))}{dL} prob(\ln(L)) \propto \frac{1}{L}$$

This is known as a Jeffrey's PDF (prior as we will see shortly)

## What can be do with Probability

### CHANGE OF VARIABLES

We want to know the PDF from a function of the original variable.

$$prob(x) = prob(y) \frac{dy}{dx} \quad \longrightarrow \quad prob(\{X_i\}) = prob(\{Y_i\}) \underbrace{\left| \frac{\partial(Y_1, Y_2, \dots, Y_n)}{\partial(X_1, X_2, \dots, X_n)} \right|}_{\text{Jacobian}}$$



## What can be do with Probability

### CHANGE OF VARIABLES

EXAMPLE: Changing from cartesian to polar coordinates

We want to express a bivariate gaussian  $prob(x, y) = \frac{1}{2\pi\sigma^2} \exp\left[-\frac{x^2 + y^2}{2\sigma^2}\right]$  in polar coordinates

We calculate the Jacobian  $\left| \frac{\partial(x, y)}{\partial(R, \theta)} \right| = \begin{vmatrix} \cos\theta & -R \sin\theta \\ \sin\theta & R \cos\theta \end{vmatrix} = R [\cos^2\theta + \sin^2\theta] = R$

$$prob(R, \theta) = \frac{R}{2\pi\sigma^2} \exp\left[-\frac{R^2}{2\sigma^2}\right]$$

Moreover, we can now extract  $prob(R)$  by marginalization!!!

$$prob(R) = \int_0^{2\pi} \frac{R}{2\pi\sigma^2} \exp\left[-\frac{R^2}{2\sigma^2}\right] d\theta = \frac{R}{\sigma^2} \exp\left[-\frac{R^2}{2\sigma^2}\right]$$

Let's think about the result...

$$prob(R) = 2\pi R \cdot \frac{A}{\sigma^2} \exp\left[-\frac{R^2}{2\sigma^2}\right] \propto R \cdot \exp\left[-\frac{R^2}{2\sigma^2}\right]$$

probability that R lies in a narrow  $\delta R$       probability at a distance R

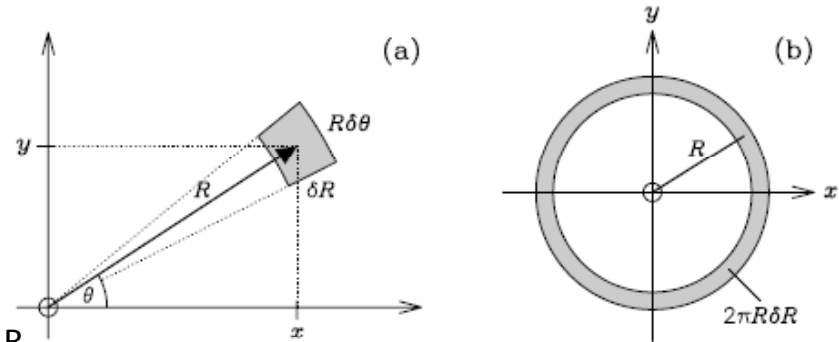


Fig. 3.14 Changing variables from Cartesian to polar coordinates.

Let's extend the result to N-dimensions...

$$prob(D | H) \propto \exp\left[-\frac{r_1^2 + r_2^2 + \dots + r_N^2}{2}\right] \quad \text{being} \quad r_i = \frac{D_i - H_i}{\sigma_i}$$

This is almost a  $\chi^2$  distribution after some massage. Following the previous reasoning

$$prob(R | H) \propto R^{N-1} \exp\left[-\frac{R^2}{2}\right] \quad \text{being} \quad R = \sqrt{\sum_{i=1}^N r_i^2}$$

Finally, taking into account  $\chi^2 = R^2$  we get the probability for  $\chi^2$  (and not for R)

$$prob(\chi^2 | H) \propto R^{N-1} \exp\left[-\frac{R^2}{2}\right] \frac{dR}{d\chi^2} = \left(\sqrt{\chi^2}\right)^{N-1} \exp\left[-\frac{R^2}{2}\right] \frac{d\sqrt{\chi^2}}{d\chi^2} = (\chi^2)^{\frac{N-1}{2}} \exp\left(-\frac{\chi^2}{2}\right)$$

$$prob(\chi^2 | H) \propto (\chi^2)^{\frac{N-1}{2}} \exp\left(-\frac{\chi^2}{2}\right)$$

probability that  $\chi$  at a distance  $\delta\chi$

probability at a distance  $\chi$



And finally: Bayes theorem

From the “multiplication” of probability:

$$\text{prob}(A, B) = \text{prob}(A) \cdot \text{prob}(B | A)$$

$$\text{prob}(A, B) = \text{prob}(B) \cdot \text{prob}(A | B)$$

dividing the two equations we get:



## Bayes theorem

$$\text{prob}(A | B) = \frac{\text{prob}(B | A) \text{prob}(A)}{\text{prob}(B)}$$

That help us to reverse probabilities...

# Bayes theorem

$$\text{prob}(H | D) = \frac{\text{prob}(D | H) \text{prob}(H)}{\text{prob}(D)}$$

## Posterior prob(H|D):

What you want to know is the probability that your hypothesis is true given the data

## Likelihood (or L) prob(D|H):

What you know is your hypothesis. And therefore you can calculate the probability that your data “gathers” around your hypothesis

## Prior prob(H):

You might want to include any prior information about your hypothesis

## Evidence (E) prob(D):

Is “simply” a normalization factor... we will come back when doing model selection

# The ubiquitous $\chi^2$

## Bayes theorem



Likelihood  
Probability that our data  
describes the hypothesis

Prior  
Our forehand knowledge

$$P(H | D) = \frac{P(D | H) \cdot P(H)}{P(D)}$$

*T. Bayes.*

Posterior  
Probability that the hypothesis is true  
given the experimental data

Evidence  
Normalization factor

In our case is very simple...

# The ubiquitous $\chi^2$

## Bayes theorem



Likelihood  
Probability that our data  
describes the hypothesis

Maximum ignorance Prior

$$P(H | D) \propto \frac{P(D | H) \cdot P(H)}{P(D)}$$

*T. Bayes.*

Posterior  
Probability that the hypothesis is true  
given the experimental data

We only care about proportionality

Thus, assuming a maximum ignorance prior:

$$\textit{prob}(H | D) \propto \textit{prob}(D | H)$$

But what is exactly H, in practice?

Let's make clear some notation:

H is in our case a function of some parameters that describe your data

$$H_i(a, b) = a + bx_i$$

And we will find this written in many ways in the literature: argh!!

$$H \equiv H \{P_i\} \equiv \{P_i\} \equiv H, \vec{P} \equiv H, \vec{\alpha} \equiv H, \vec{\theta}$$



## Remember our goals

- 1.- Parameter estimation
- 2.- Hypothesis testing (model selection)

## Remember our goals

1.- Parameter estimation

2.- Hypothesis testing (model selection)

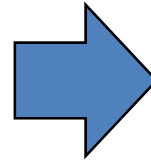
## Parameter estimation

We would like to estimate the best estimate of a quantity, given some data.

Let's define the posterior as:

$$prob(X | \{data\}) = P$$

The best estimate of X is given by a maximum in its probability and thus



$$\left. \frac{dP}{dX} \right|_{X_0} = 0$$

We want now to expand P as a Taylor series, since P varies too fast we take the logarithm of P

$$L = \ln(P) = \ln(prob(X | \{data\}))$$

The Taylor expansion of L is:

$$L = L(X_0) + \frac{1}{2} \left. \frac{d^2 L}{dX^2} \right|_{X_0} (X - X_0)^2 + \dots$$

Taking the first term we get for P:

$$prob(X | \{data\}) \approx A \exp \left[ \frac{1}{2} \left. \frac{d^2 L}{dX^2} \right|_{X_0} (X - X_0)^2 \right]$$

Let's compare with a gaussian

$$prob(X | \{data\}) \approx A \exp \left[ \frac{(X - X_0)^2}{2\sigma^2} \right]$$

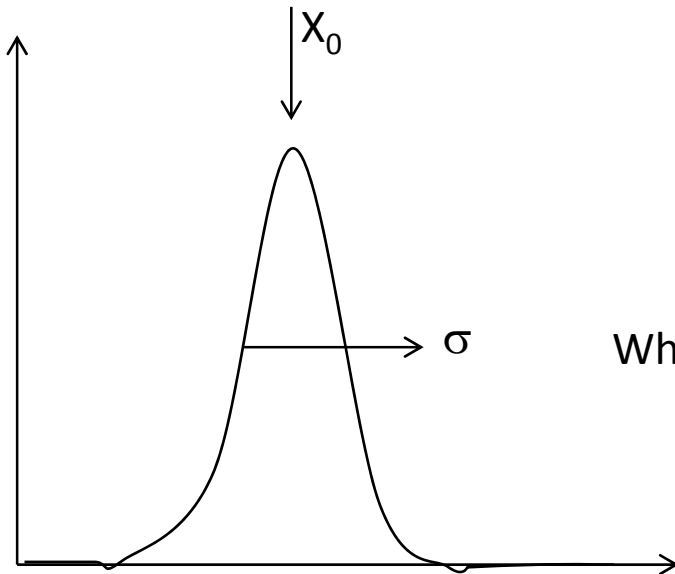
Taking the first term we get for P:

Comparing

$$\text{prob}(X | \{data\}) \approx A \exp \left[ \frac{1}{2} \frac{d^2 L}{dX^2} \Big|_{X_0} (X - X_0)^2 \right]$$

Let's compare with a gaussian

$$\text{prob}(X | \{data\}) \approx A \exp \left[ -\frac{(X - X_0)^2}{2\sigma^2} \right]$$



$$\sigma = \left( -\frac{d^2 L}{dX^2} \Big|_{X_0} \right)^{-1/2}$$

Where  $\frac{d^2 L}{dX^2} \Big|_{X_0}$  is always negative, so: no problem ;-)

$$X = X_0 \pm \sigma$$

Value plus error comes from a Taylor series  
(closing the eyes to the "real" shape... that sometimes is ok, and sometimes not)

## Remember our goals

1.- Parameter estimation

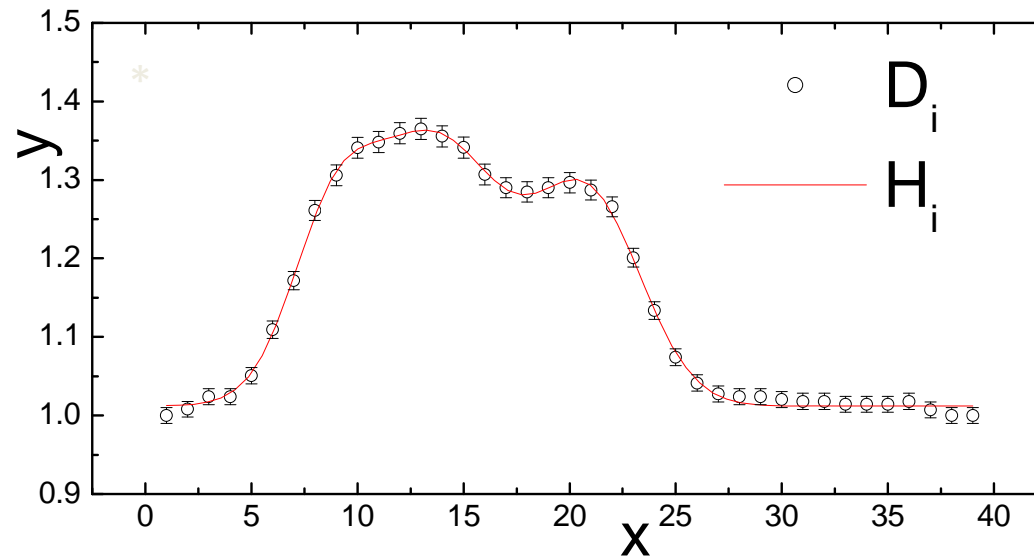
2.- Hypothesis testing (model selection)

Bayes theorem



$$P(H | D) \propto P(D | H) \equiv L$$

$D_i$  Data ( $i=1,n$ )       $H_i\{P_l\}$  Hypothesis ( $i=1,n$ ) using a parameter set  $\{P_l\}$  ( $l=1,m$ )



$$P(H_i\{P_l\} | D_i) \propto P(D_i | H_i\{P_l\})$$

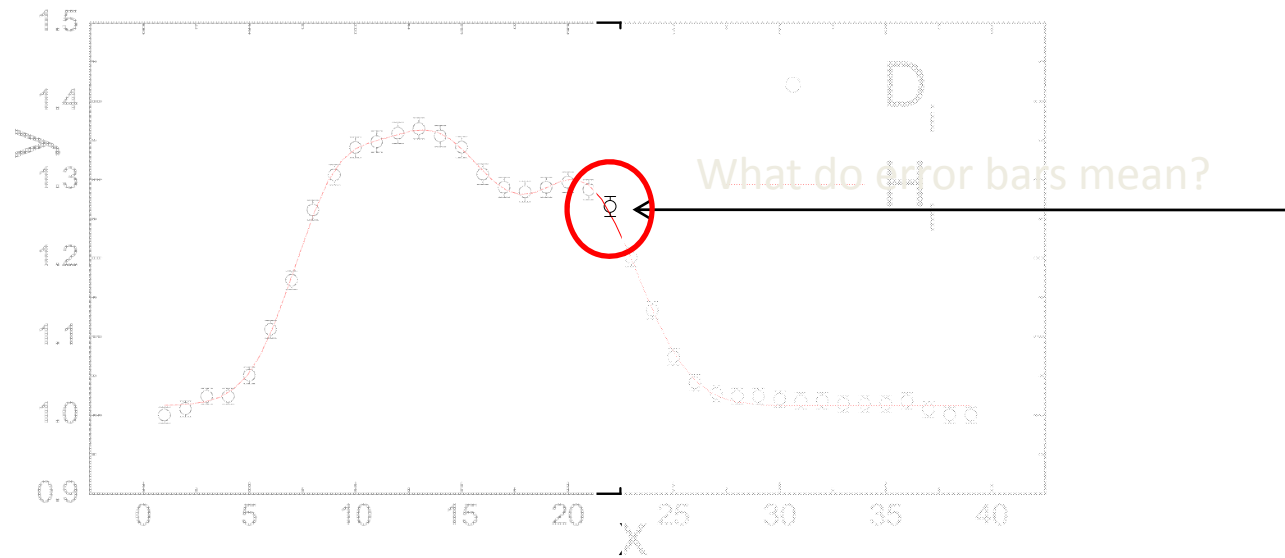
\* Figure adapted from „Le petit prince“ A. Saint Exupery (1943)

# Bayes theorem



$$P(H | D) \propto P(D | H)$$

$D_i$  Data ( $i=1,n$ )       $H_i \{P_l\}$  Hypothesis ( $i=1,n$ ) using a parameter set  $\{P_l\}$  ( $l=1,m$ )

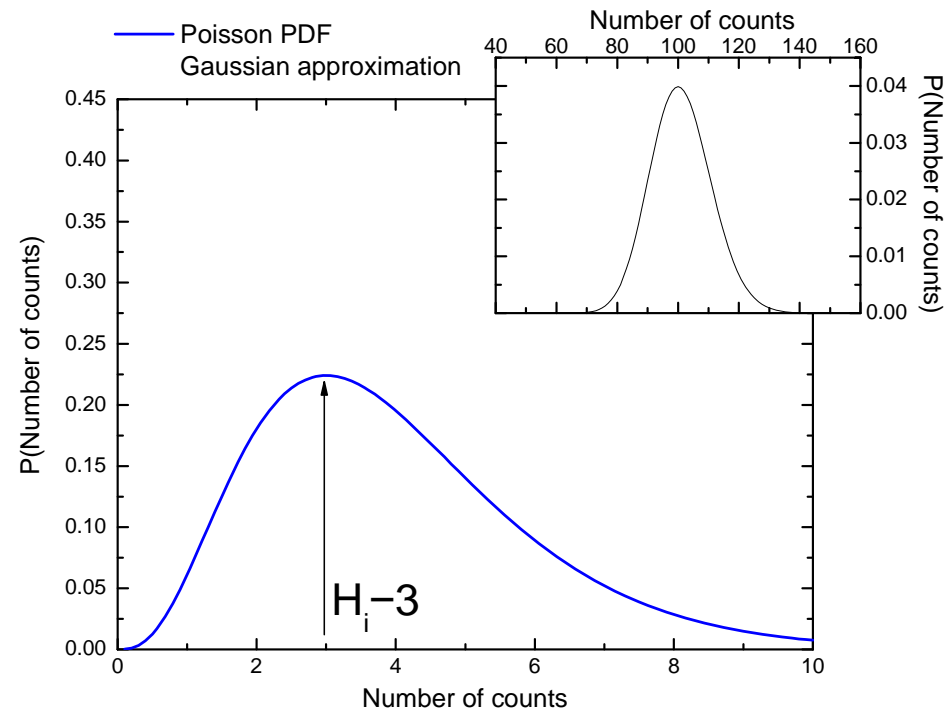


$$P(H_i \{P_l\} | D_i) \propto P(D_i | H_i \{P_l\})$$

## The ubiquitous $\chi^2$

In a counting experiment, if the expected value is  $H_i$   
measured values will follow a Poisson statistics

$$P(D_{i=k} | H_{i=k}) \propto H_i^{D_{i=k}} e^{-H_{i=k}}$$

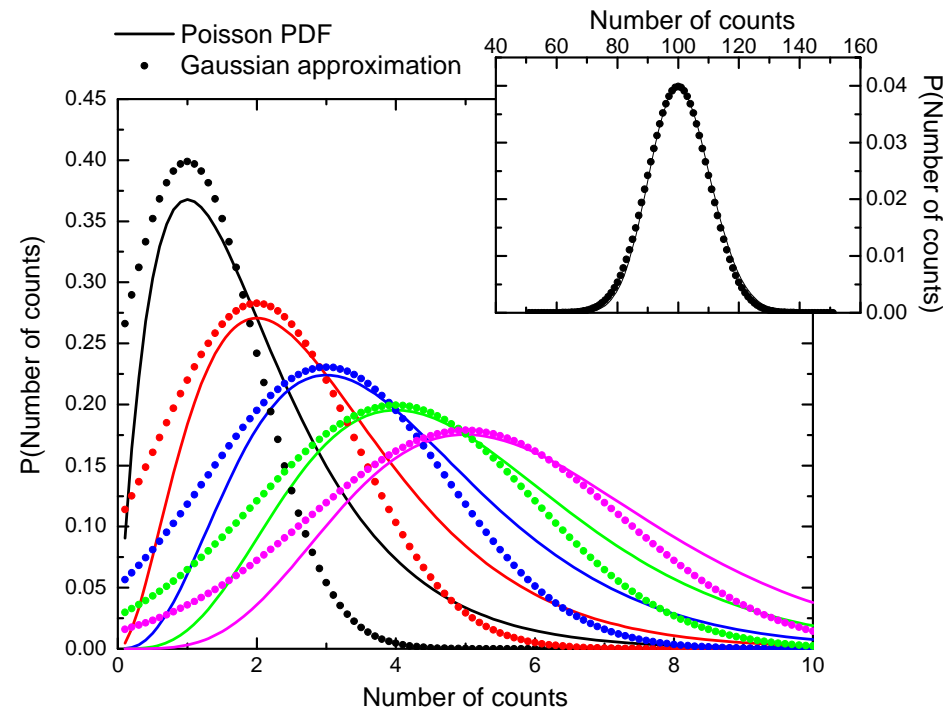




## The ubiquitous $\chi^2$

In a counting experiment, if the expected value is  $H_i$  measured values will follow a Poisson statistics

$$P(D_{i=k} | H_{i=k}) \propto H_i^{D_{i=k}} e^{-H_{i=k}}$$



If the number of counts is high enough poisson statistics = normal distribution with  $\sigma_i = \sqrt{D_i}$

$$P(D_{i=k} | H_{i=k}) \propto \exp\left[-\frac{(H_{i=k} - D_{i=k})^2}{2\sigma_{i=k}^2}\right]$$

The error is the square root of the variance  $\epsilon_i = \sigma_i$

## The ubiquitous $\chi^2$

Bayes theorem

$$P(H_i \{P_l\} | D_i) \propto P(D_i | H_i \{P_l\})$$



We now consider all the points  $i=1,2,\dots, n$

$$\begin{aligned} L = P(D_i | H_i \{P_l\}) &\propto \prod_{i=1}^n \exp - \frac{(H_i - D_i)^2}{2\sigma_i^2} \\ &= \exp \sum_{i=1}^n - \frac{(H_i - D_i)^2}{2\sigma_i^2} = \exp - \frac{\chi^2}{2} \end{aligned}$$

Therefore  $\chi^2$  is related to the likelihood:

$$\chi^2 \propto -2 \cdot \ln L$$

**So, we got the exact meaning on probability bases of  $\chi^2$   
Let's use it!**

## The ubiquitous $\chi^2$

Least squares method therefore:

$$L = \text{prob}(D | H, P_i) \propto \exp\left(-\chi^2/2\right)$$

Maximize the probability that the Data is described by your hypothesis

$$\text{prob}(H_i, P_i | D_i) \propto \exp\left(-\chi^2/2\right)$$

And therefore to maximize the probability of your hypothesis!!!

## Hypothesis testing

Does our model describe our data correctly?

$$\chi^2 \propto -2 \cdot \ln L$$

the logarithm of the likelihood answers the question  $\chi^2$  follows a  $\chi^2$  distribution...

... so we can go back to P values and company. But...

did we made all this way to go back to P values???

## Model selection

Imagine two competing models or Hypothesis  $H_1$  and  $H_2$

$$P(H_1 | D) = \frac{P(D | H_1)P(H_1)}{P(D)} \qquad P(H_2 | D) = \frac{P(D | H_2)P(H_2)}{P(D)}$$

$$\underbrace{\frac{P(H_1 | D)}{P(H_2 | D)}}_{\text{posterior ratio}} = \underbrace{\frac{P(D | H_1)}{P(D | H_2)}}_{\text{Bayes factor: B}} \cdot \underbrace{\frac{P(H_1)}{P(H_2)}}_{\text{prior for theories (usually 1)}}$$

$P(D|H)$ : To assign probabilities we take the hypothesis **INDEPENDENT** of the particular parameter we obtain in the fitting: we marginalize the parameters:

$$P(D | H) = \int \text{prob}(D, \vec{\alpha} | H) d\vec{\alpha} = \int \text{prob}(D | \vec{\alpha}, H) \text{prob}(\vec{\alpha} | H) d\vec{\alpha}$$

Therefore

$$B = \frac{P(D | H_1)}{P(D | H_2)} = \frac{\int \text{prob}(D | \vec{\alpha}, H_1) \text{prob}(\vec{\alpha} | H_1) d\vec{\alpha}}{\int \text{prob}(D | \vec{\beta}, H_2) \text{prob}(\vec{\beta} | H_2) d\vec{\beta}}$$

## Model selection

Imagine two competing models or Hypothesis  $H_1$  and  $H_2$

$$B = \frac{P(D | H_1)}{P(D | H_2)} = \frac{\int \underbrace{\text{prob}(D | \vec{\alpha}, H_1)}_{\text{Probability of the data around best fit } \alpha_i^*} \underbrace{\text{prob}(\vec{\alpha} | H_1)}_{\text{prior from out parameter maximum ignorance:}} d\vec{\alpha}}{\int \underbrace{\text{prob}(D | \vec{\beta}, H_2)}_{\text{Probability of the data around best fit } \alpha_i^*} \underbrace{\text{prob}(\vec{\beta} | H_2)}_{\text{prior from out parameter maximum ignorance:}} d\vec{\beta}}$$

Probability of the data around best fit  $\alpha_i^*$

$$\text{prob}(D | \alpha_i, H) = \underbrace{\text{prob}(D | \alpha_i^*)}_{\text{Probability of the data around best fit } \alpha_i^*} \cdot \exp\left[-\frac{(\alpha_i - \alpha_i^*)^2}{2\delta\alpha^2}\right]$$

$$L = A \exp\left(-\frac{\hat{\chi}^2}{2}\right)$$

prior from out parameter maximum ignorance:

location

$$\text{prob}(\alpha_i) = \frac{1}{\alpha_{\max} - \alpha_{\min}}$$

scale factor

$$\text{prob}(\alpha_i) = \frac{1/\alpha_i}{\int_{\alpha_{\min}}^{\alpha_{\max}} 1/\alpha_i d\alpha_i} = \frac{1}{\alpha_i \ln\left(\frac{\alpha_{\max}}{\alpha_{\min}}\right)}$$

## Model selection

Imagine two competing models or Hypothesis  $H_1$  and  $H_2$

$$B = \frac{P(D | H_1)}{P(D | H_2)} = \frac{\int \text{prob}(D | \vec{\alpha}, H_1) \text{prob}(\vec{\alpha} | H_1) d\vec{\alpha}}{\int \text{prob}(D | \vec{\beta}, H_2) \text{prob}(\vec{\beta} | H_2) d\vec{\beta}}$$

Let's assume that each fit has only one parameter:

$$B = \frac{\exp\left(-\chi_{H_1}^2 / 2\right)}{\exp\left(-\chi_{H_2}^2 / 2\right)} \underbrace{\frac{\delta\alpha}{\alpha_{\max} - \alpha_{\min}} \frac{\beta_{\max} - \beta_{\min}}{\delta\beta}}_{\text{penalizes additional parameters}}$$

favours best fit

penalizes additional parameters

Ockham factor

$$\alpha_{\max} - \alpha_{\min} \gg \delta\alpha$$

# Model selection

Akaike and Bayesian information criteria

Akaike information criteria

gaussian distributed points

$$AIC = -2 \ln L_{\max} + 2k = \chi^2 + 2k$$

Bayesian Information Criteria

$$P(D | H) = \int \text{prob}(D, \vec{\alpha} | H) d\vec{\alpha} = \int \text{prob}(D | \vec{\alpha}, H) \text{prob}(\vec{\alpha} | H) d\vec{\alpha}$$

Can be approximated by:

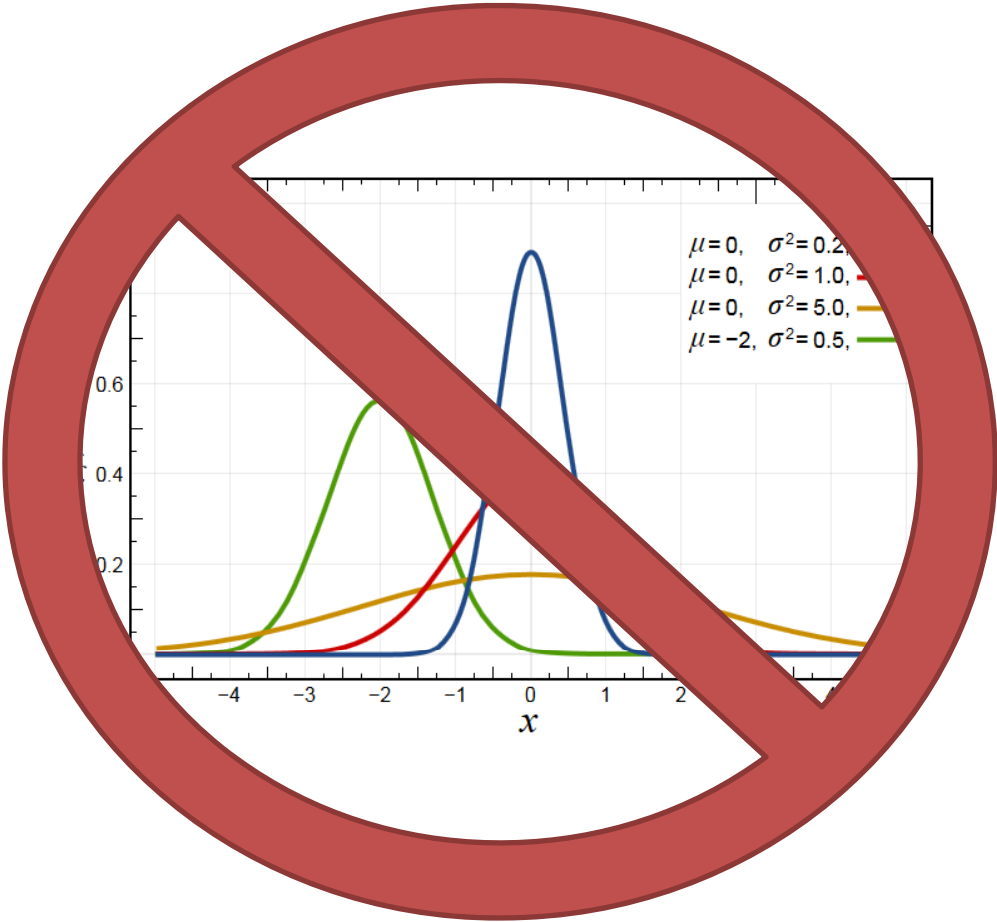
$$P(D | H) \approx BIC = -2 \ln L_{\max} + 2k = \chi^2 + k \ln N$$

where k is the number of parameters, and N the number of points.

We should then choose the one with minimum AIC or BIC



Beyond the Normal distribution



There is a lot still based on normal distribution: is there a method to avoid this????

$$L = P(D_i | H_i \{P_l\}) \propto \prod_{i=1}^n \underbrace{\exp\left[-\frac{(H_i - D_i)^2}{2\sigma_i^2}\right]}$$

You cannot avoid to assign a probability to each point, that might be gaussian as here...  
or can be a poisson distribution:

$$\log L = \sum_{i=1}^{i=Ndata} \log[P(H|Data_i)] = \sum_{i=1}^{i=Ndata} \log\left[\frac{e^{-H_i} H_i^{D_i}}{D_i!}\right] \propto \sum_{i=1}^{i=Ndata} [-H_i + D_i \cdot \log H_i]$$

or any other distribution...

There is a lot still based on normal distribution: is there a method to avoid this????

...also, gaussianity comes from exploring around the minimum:  
no matter if the minimum is flat, or multimodal...

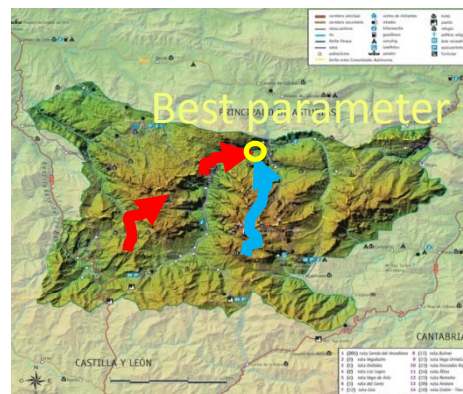
$$L = L(X_0) + \frac{1}{2} \frac{d^2 L}{dX^2} \Big|_{X_0} (X - X_0)^2 + \dots$$

Can we simply get L??

Markov Chain Montecarlo method

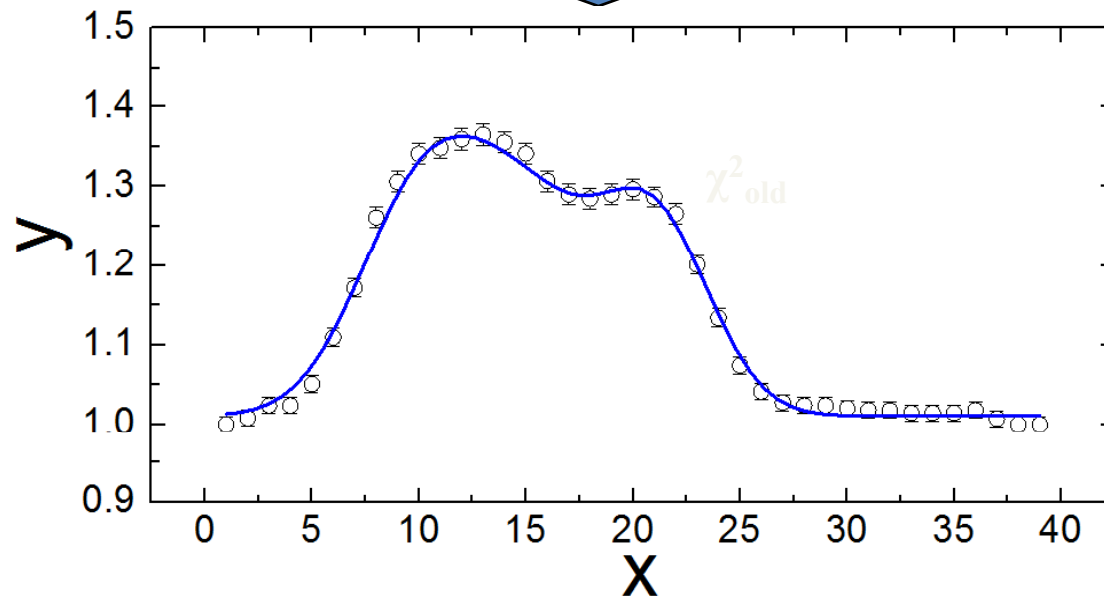
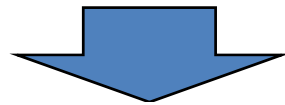
we can explore the likelihood, better its logarithm  $\chi^2$

Let's define a parameter space  $\chi^2\{\text{Pi}\}$



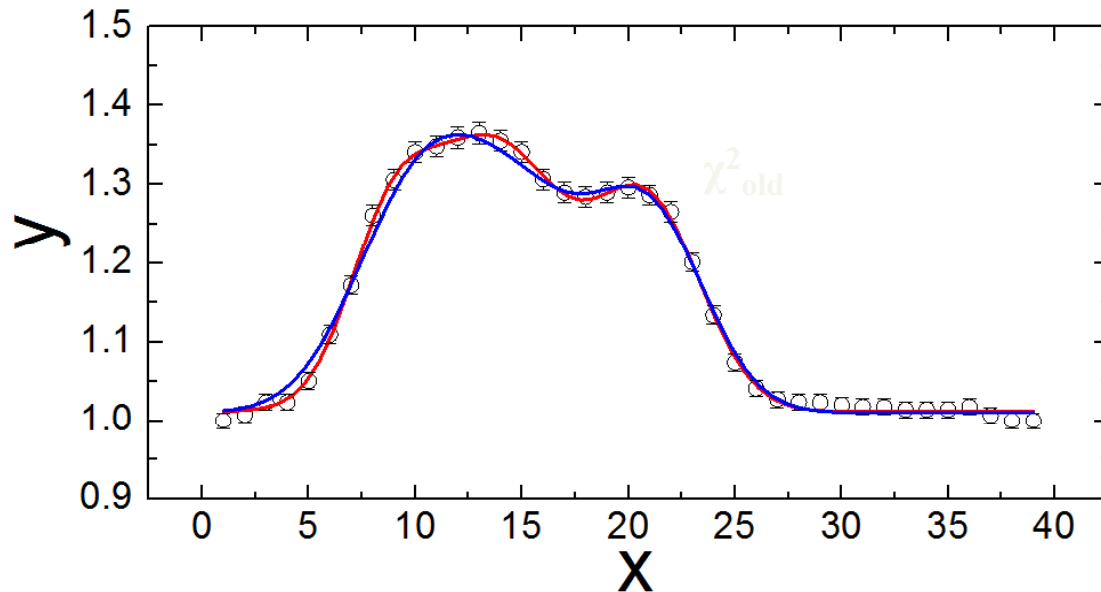
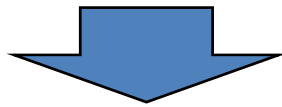
# Bayesian tool to explore $\chi^2$ landscape

Step	$\chi^2$	$P_1$	$P_2$	$P_3$
1	$1.368 \cdot 10^4$	1.341	0.00223	10.223
2				
3				
4				
...	...	...	...	...



## Bayesian tool to explore $\chi^2$ landscape

Step	$\chi^2$	$P_1$	$P_2$	$P_3$
1	$1.368 \cdot 10^4$	1.341	0.00223	10.223
2	$1.360 \cdot 10^4$	1.342	0.00223	10.223
3				
4				
...	...	...	...	...



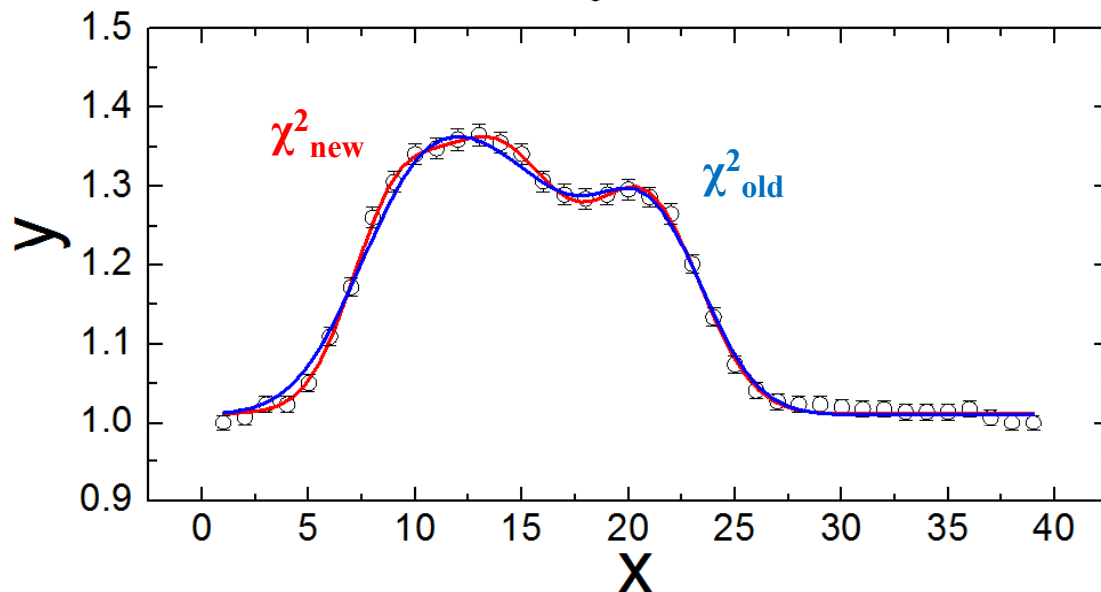
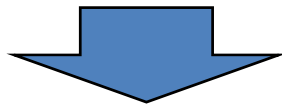
Acceptance rules:

- If better always
- If not, with probability

$$P_{old \rightarrow new} = \exp\left(-\frac{\chi^2_{new} - \chi^2_{old}}{2}\right)$$

## Bayesian tool to explore $\chi^2$ landscape

Step	$\chi^2$	$P_1$	$P_2$	$P_3$
1	$1.368 \cdot 10^4$	1.341	0.00223	10.223
2	$1.360 \cdot 10^4$	1.342	0.00223	10.223
3	$1.362 \cdot 10^4$	1.342	0.00212	10.223
4	$1.348 \cdot 10^4$	1.342	0.00212	10.170
...	...	...	...	...



Acceptance rules:

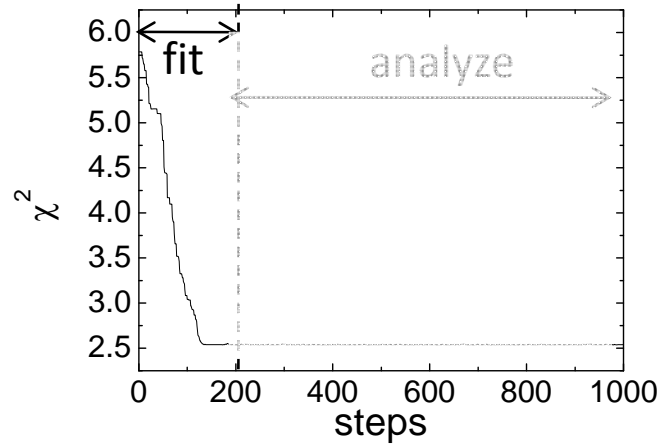
- If better always
- If not, with probability

$$P_{old \rightarrow new} = \exp\left(-\frac{\chi^2_{new} - \chi^2_{old}}{2}\right)$$

# Bayesian tool to explore $\chi^2$ landscape

☐ It does not get stuck

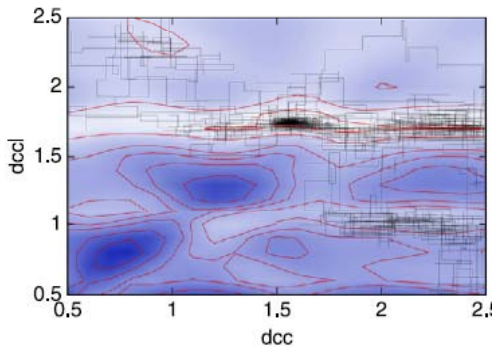
... and if it does  
we increase the errors



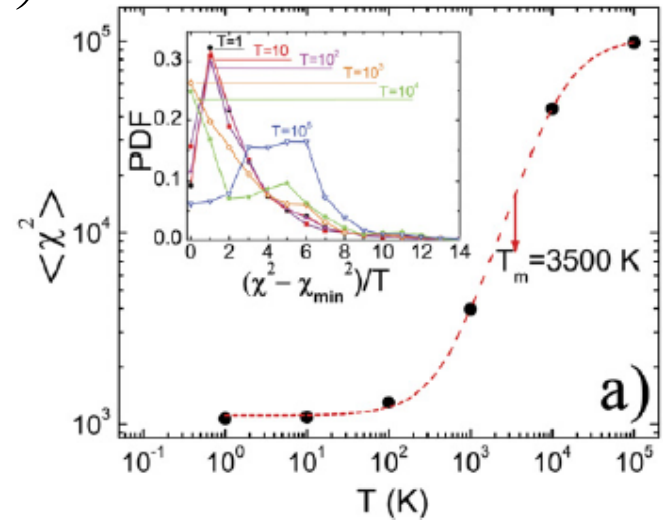
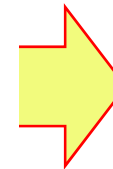
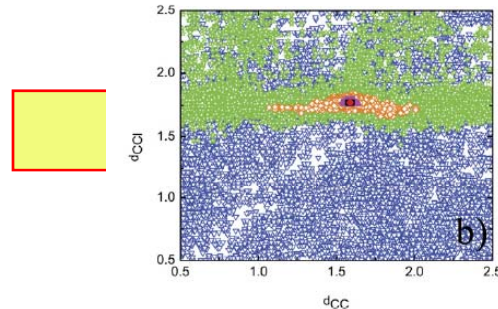
## Simulated annealing

$$P_{old \rightarrow new} = \exp\left(-\frac{\chi_{new}^2 - \chi_{old}^2}{2T}\right)$$

A 2D  $\chi^2$  landscape

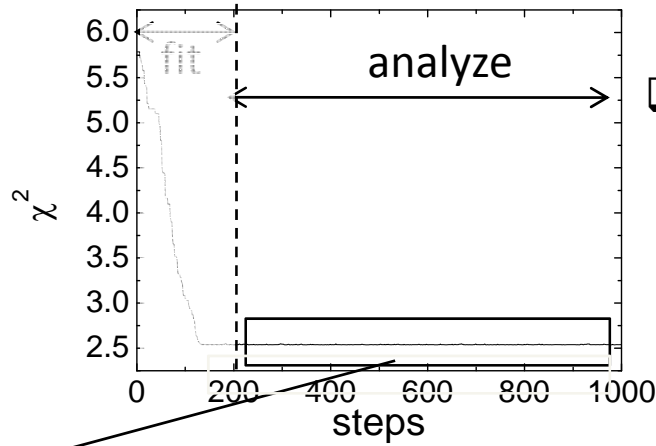


Warming up the fitting

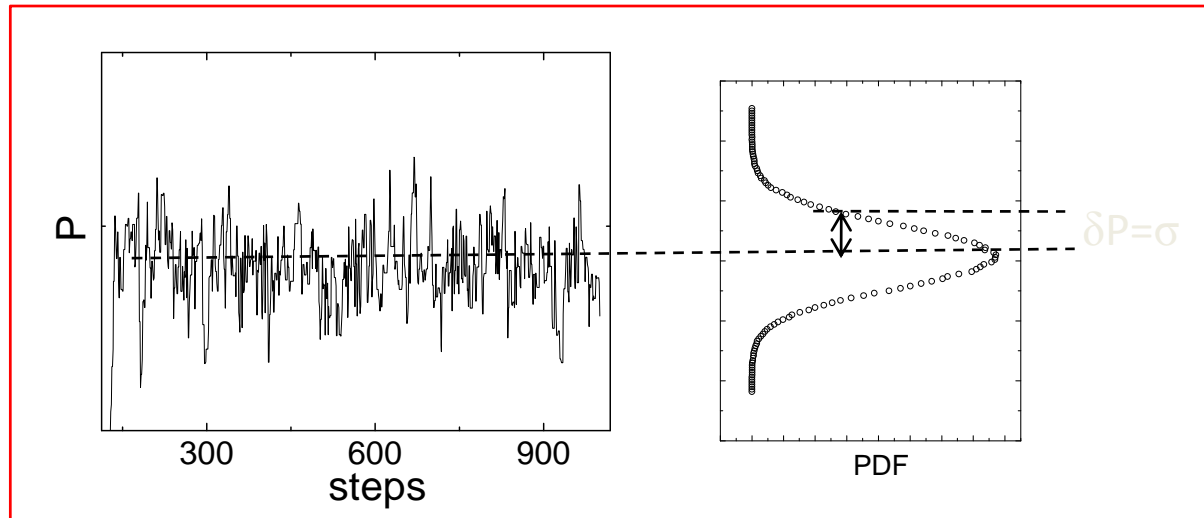
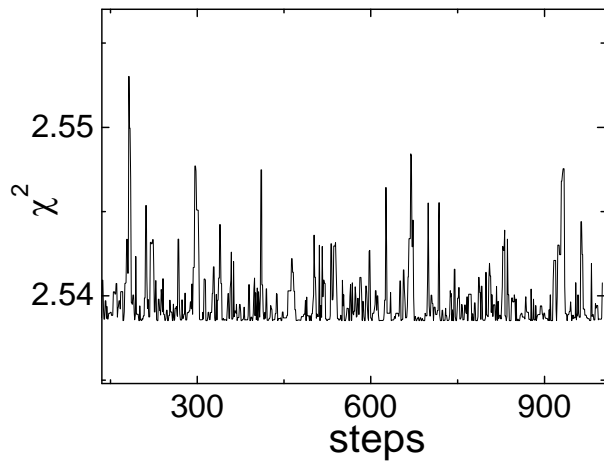


# Bayesian tool to explore $\chi^2$ landscape

☐ It does not get stuck



☐ Robust parameter estimation



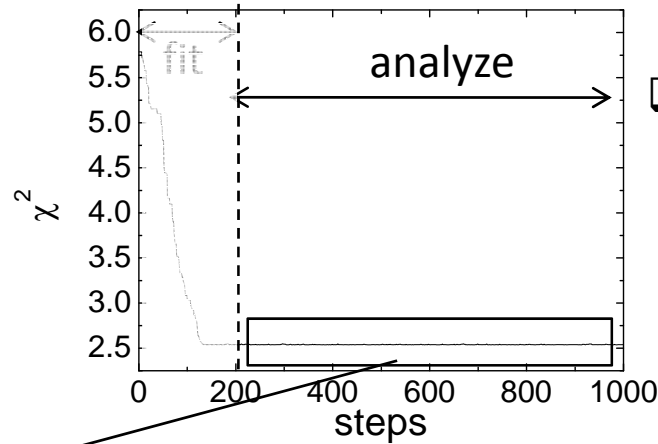
Parameter estimation:

Parameters are obtained as PDF's not as  $P \pm \delta P$

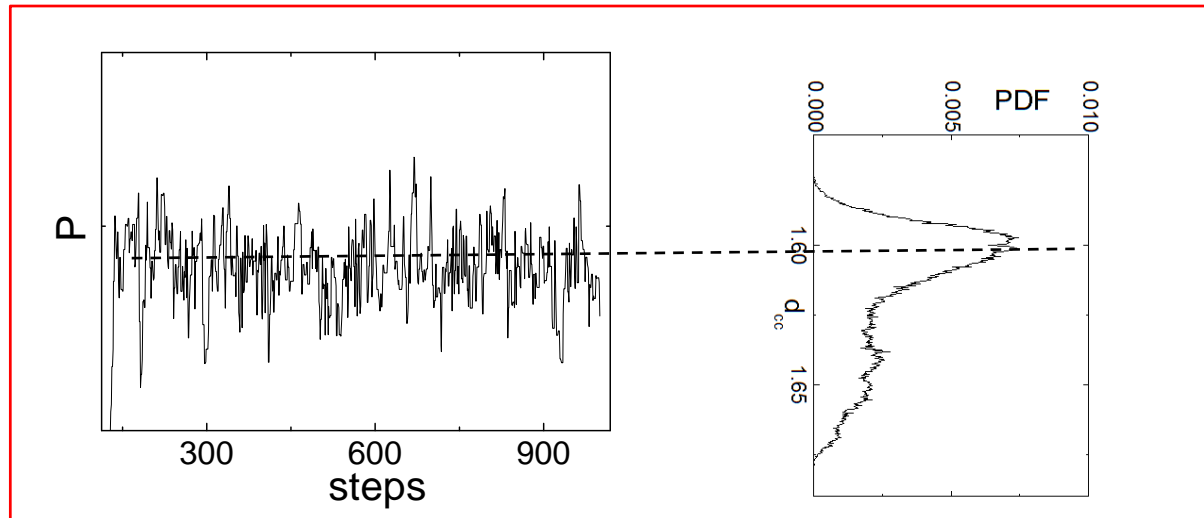
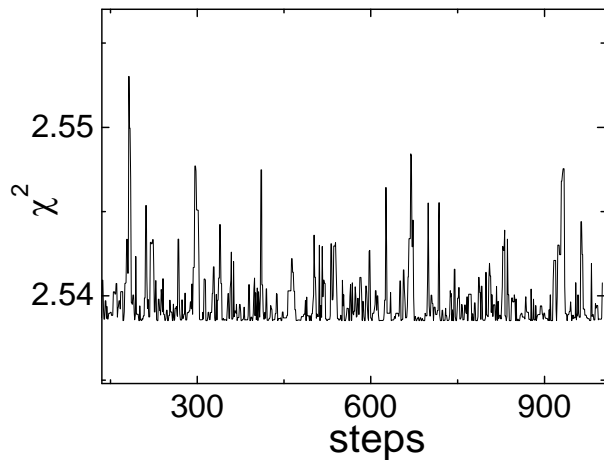


# Bayesian tool to explore $\chi^2$ landscape

☐ It does not get stuck



☐ Robust parameter estimation

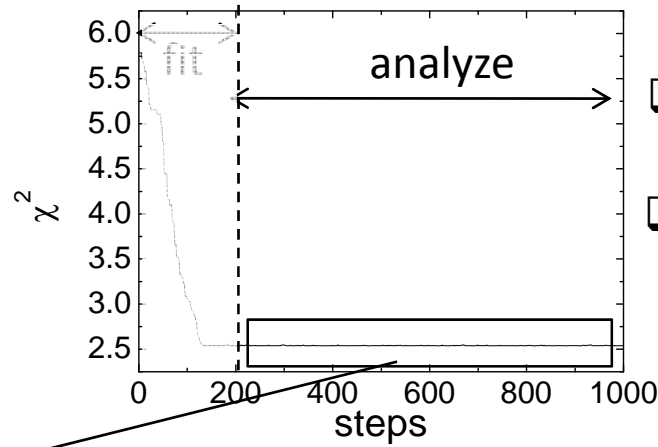


Parameter estimation: Parameters are obtained as PDF's not as  $P \pm \delta P$

You are marginalizing all parameters except the one you are interested in

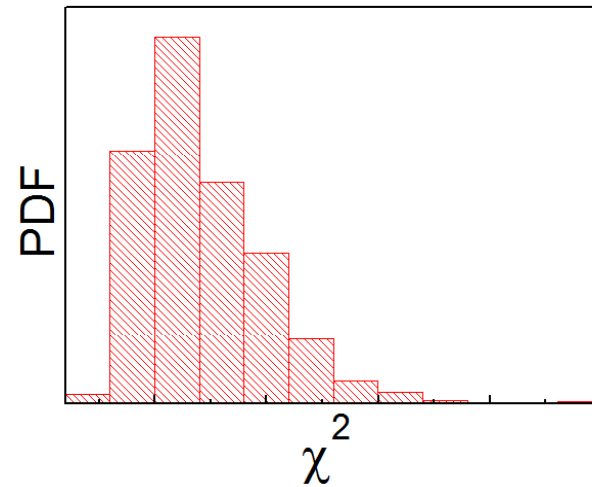
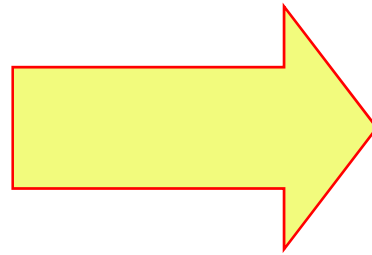
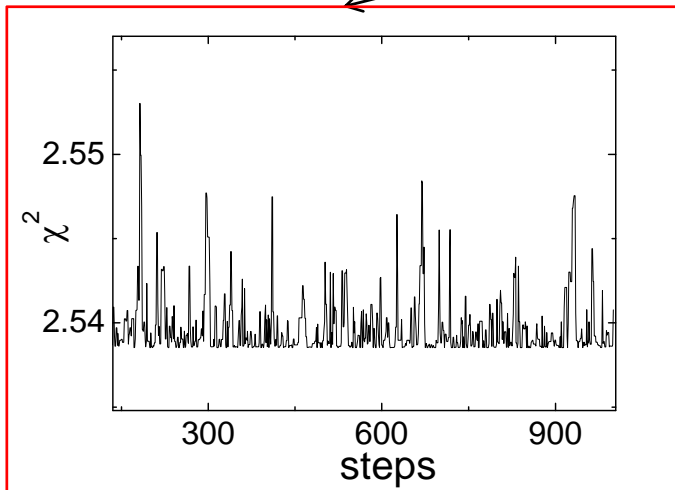
# Bayesian tool to explore $\chi^2$ landscape

☐ It does not get stuck



☐ Robust parameter estimation

☐ Robust model selection



Model selection:

all possible combinations of parameters are investigated  
and their  $\chi^2$  calculated

### **Molecular spectroscopy and Bayesian spectral analysis—how many lines are there?**

D. S. Sivia and C. J. Carlile

*ISIS Pulsed Neutron Facility, Rutherford Appleton Laboratory, Chilton, Didcot, Oxfordshire, OX11 0QX,  
United Kingdom*

(Received 24 July 1991; accepted 24 September 1991)

#### Models tested:

- Stretched exponential (KWW)
- Stretched exponential + gaussian distribution of lorentz lines
- Two gaussian distribution of lorentz lines
- Two lorentzians
- Three lorentzians

## Molecular spectroscopy and Bayesian spectral analysis—how many lines are there?

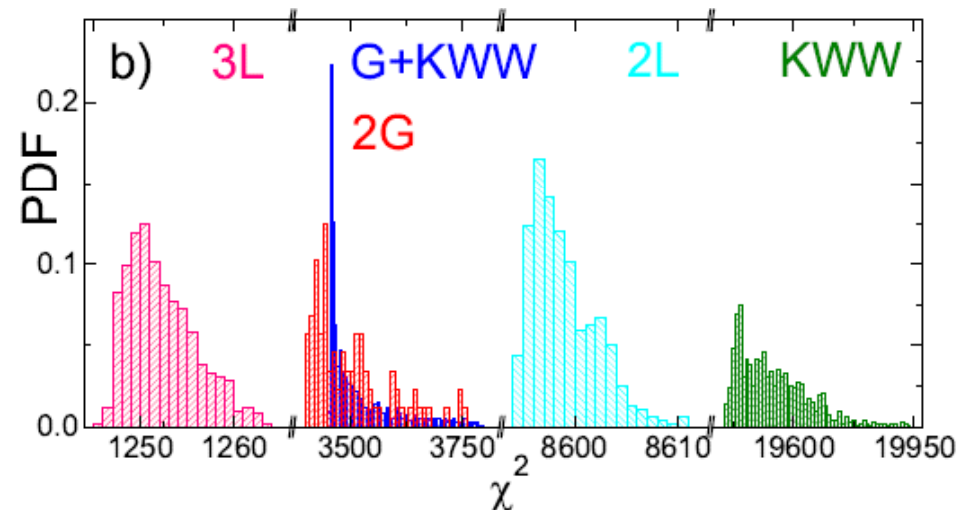
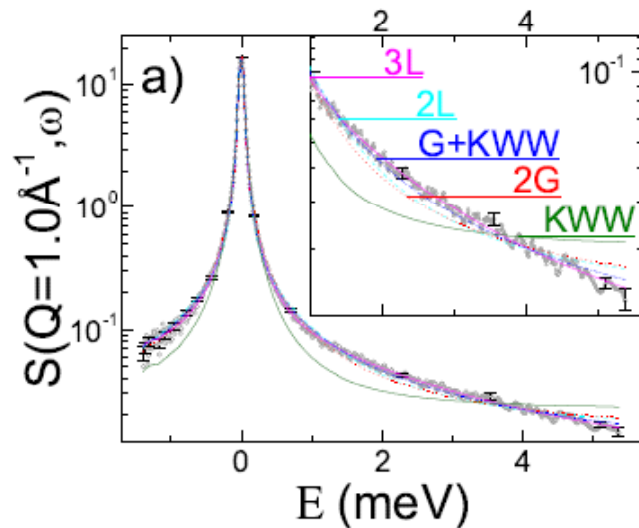
D. S. Sivia and C. J. Carlile

*ISIS Pulsed Neutron Facility, Rutherford Appleton Laboratory, Chilton, Didcot, Oxfordshire, OX11 0QX, United Kingdom*

(Received 24 July 1991; accepted 24 September 1991)

### Models tested:

- Stretched exponential (KWW)
- Stretched exponential + gaussian distribution of lorentz lines
- Two gaussian distribution of lorentz lines
- Two lorentzians
- Three lorentzians



## Molecular spectroscopy and Bayesian spectral analysis—how many lines are there?

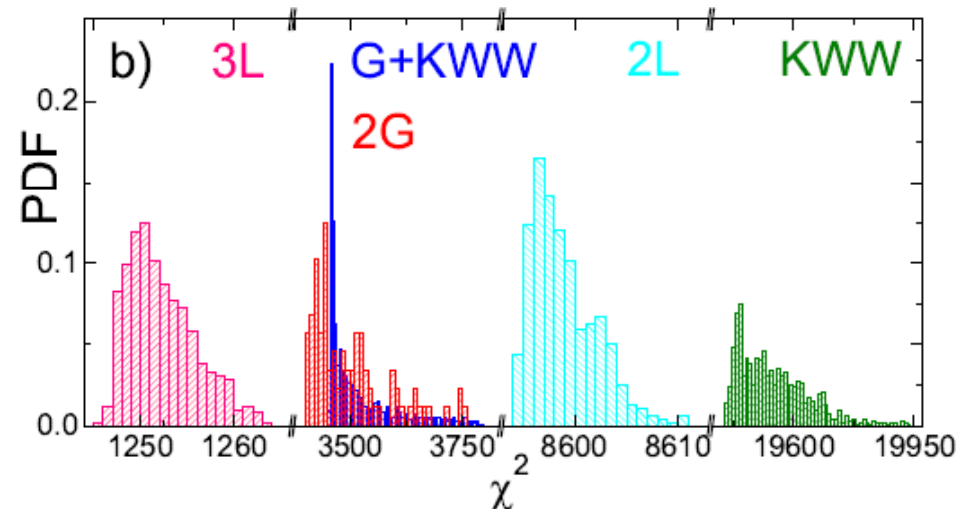
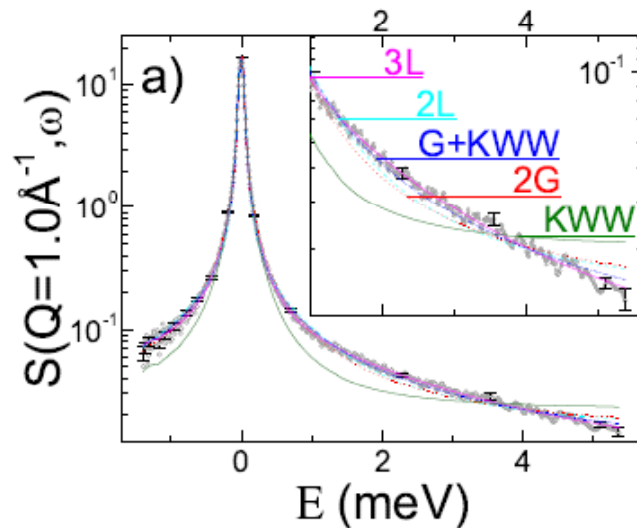
D. S. Sivia and C. J. Carlile

*ISIS Pulsed Neutron Facility, Rutherford Appleton Laboratory, Chilton, Didcot, Oxfordshire, OX11 0QX, United Kingdom*

(Received 24 July 1991; accepted 24 September 1991)

### Models tested:

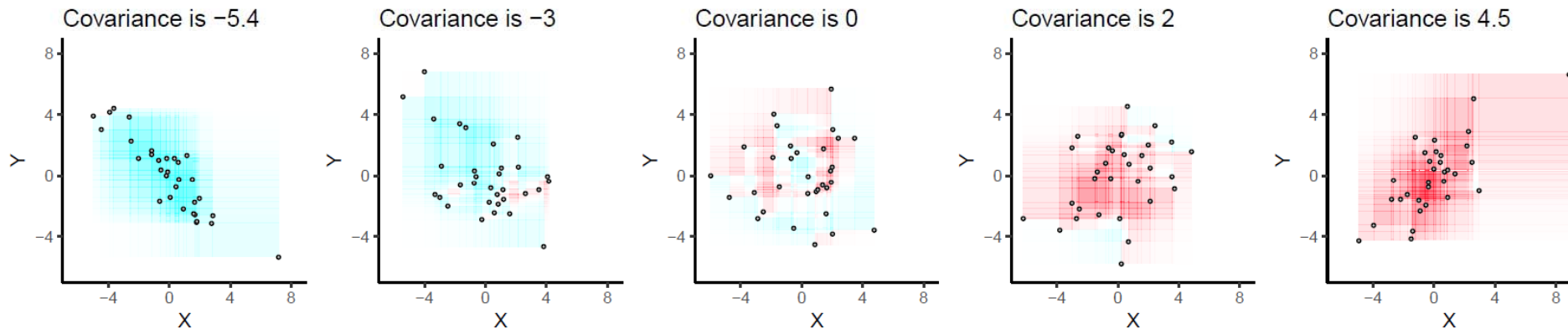
- Stretched exponential (KWW)
- Stretched exponential + gaussian distribution of lorentz lines
- Two gaussian distribution of lorentz lines
- Two lorentzians
- **Three lorentzians THE WINNER**



## PARAMETER CORRELATION

Usually only **linear** correlation, and **twobody** correlations: **covariance**

$$\sigma_{XY}^2 = \sum_{i=1}^N (X_i - \bar{X})(Y_i - \bar{Y})$$



Indeed, we can define a covariance matrix

$$\begin{pmatrix} \sigma_X^2 & \sigma_{XY}^2 \\ \sigma_{XY}^2 & \sigma_Y^2 \end{pmatrix}$$

## PARAMETER CORRELATION

Usually only correlation around the minimum, and **twobody** correlations: **covariance**

The bayesian counterpart is (not assuming gaussianity, but only the shape of the minimum)

If you remember from parameter estimation:

$$L = \ln(\text{prob}(X | \{\text{data}\})) = \ln(\text{likelihood}) \equiv \ln(\ell) \quad \sigma_x^2 = \left( - \frac{d^2 L}{dX^2} \Big|_{x_0} \right)^{-1}$$

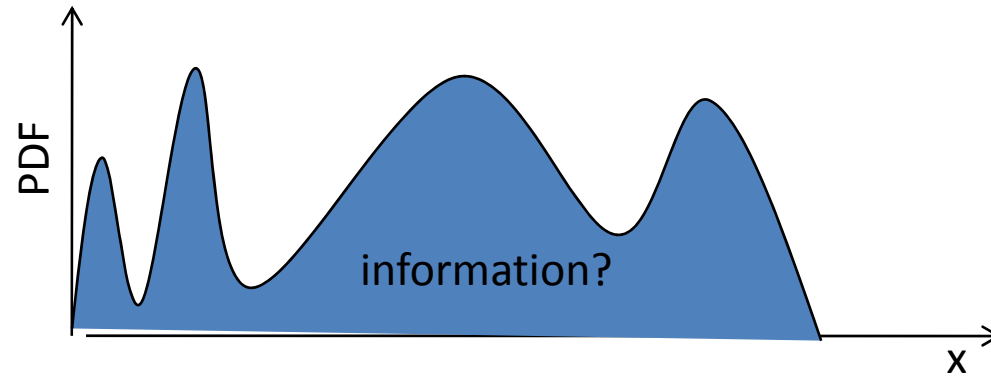
Bayesianilly speaking, we can thus define the matrix

$$\begin{pmatrix} \frac{\partial^2 \ln \ell}{\partial X^2} & \frac{\partial^2 \ln \ell}{\partial X \partial Y} \\ \frac{\partial^2 \ln \ell}{\partial Y \partial X} & \frac{\partial^2 \ln \ell}{\partial Y^2} \end{pmatrix} \text{ Fischer information matrix}$$

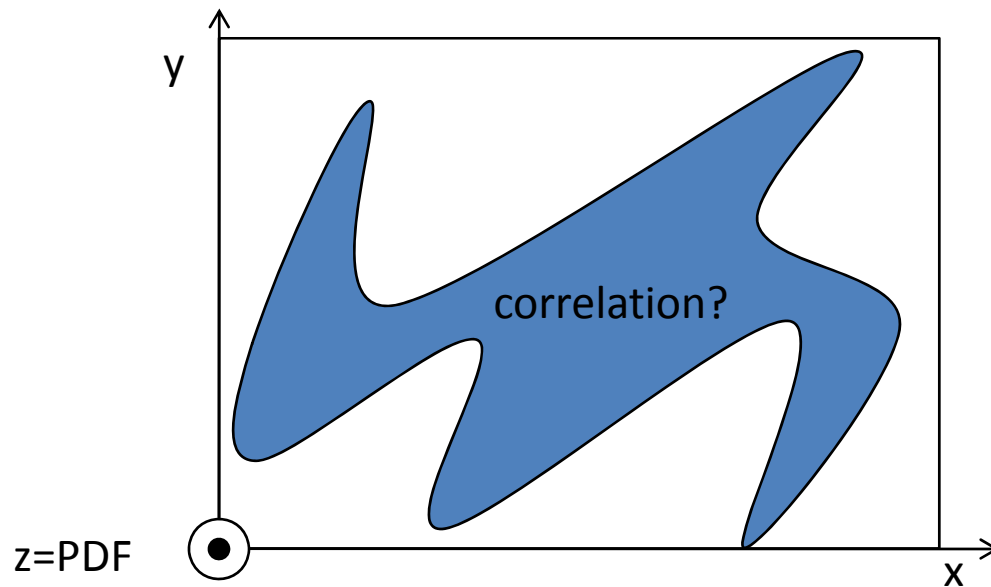
But we can go to information theory to avoid these problems, rejoice!!!

Information theory tells us:

- How many information is carrying a PDF



- How correlated is the information in a multidimensional PDF



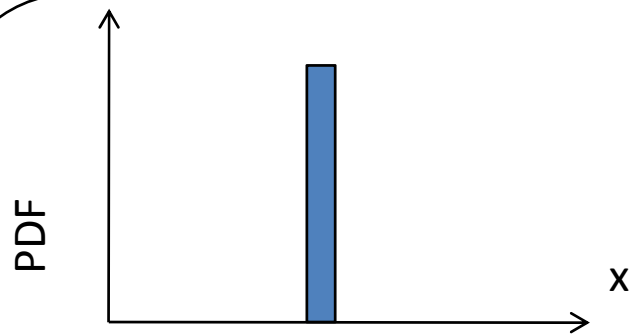


Information theory tells us:

□ How many information is carrying a PDF:

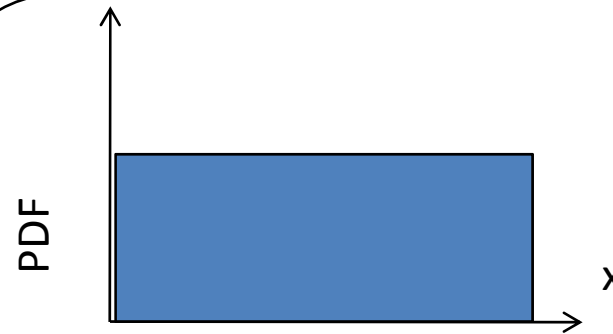
Shannon  
ENTROPY

$$H(x) = -\sum_{i=1}^n P_i(x) \ln P_i(x) = 0$$



$$H_{\min} = -\sum_{i=1}^n 1 \ln 1 = 0$$

"SHARP" PDF



$$H_{\max} = -\sum_{i=1}^n \frac{1}{n} \ln \frac{1}{n} = \ln n$$

Non informative PDF

Information theory tells us:

- How correlated are  $x, y$  in a PDF

MUTUAL INFORMATION

$$I(x, y) = H(x) + H(y) - H(x, y)$$

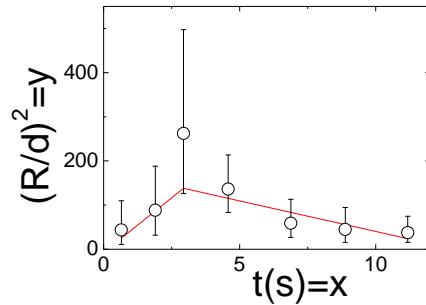
$I(x, y) = H(x) + H(y) - H(x, y)$   
 $I_{\min}(x, y) = 0 + 0 - 0 = 0$   
 **$x$  and  $y$  NOT correlated**

$I(x, y) = H(x) + H(y) - H(x, y)$   
 $I_{\min}(x, y) = n + n - 2n = 0$

$I(x, y) = H(x) + H(y) - H(x, y)$   
 $I_{\max}(x, y) = n + n - n = n$   
 **$x$  and  $y$  strongly correlated**

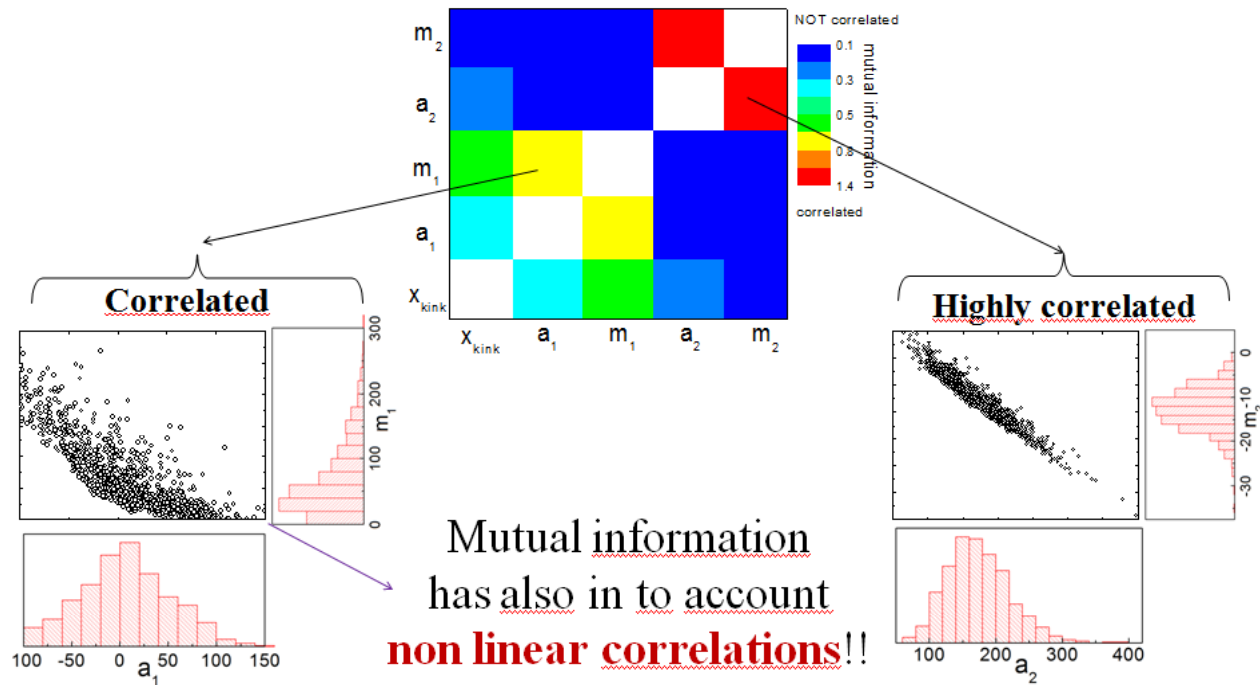
**MUTUAL INFORMATION**  
 $I(x, y) = H(x) + H(y) - H(x, y)$

model  
(astrophysics)



For  $x < x_{\text{kink}}$       $y_1 = a_1 + m_1 x$   
 For  $x > x_{\text{kink}}$       $y_2 = a_2 + m_2 x$   
 Continuity:      $a_2 = a_1 + (m_1 - m_2) x_{\text{kink}}$

We calculate MI for all combinations of two parameters



Information theory tells us:

- How many information is carrying a PDF

Shannon ENTROPY

$$H(x) = -\sum_{i=1}^n P_i(x) \ln P_i(x) = 0$$

- How correlated is the information in a multidimensional PDF

MUTUAL INFORMATION

$$I_n(A_1, \dots, A_n) = \sum_{k=1}^n (-1)^{k+1} \sum_{i_1 < \dots < i_k} S(A_{i_1}, \dots, A_{i_k})$$