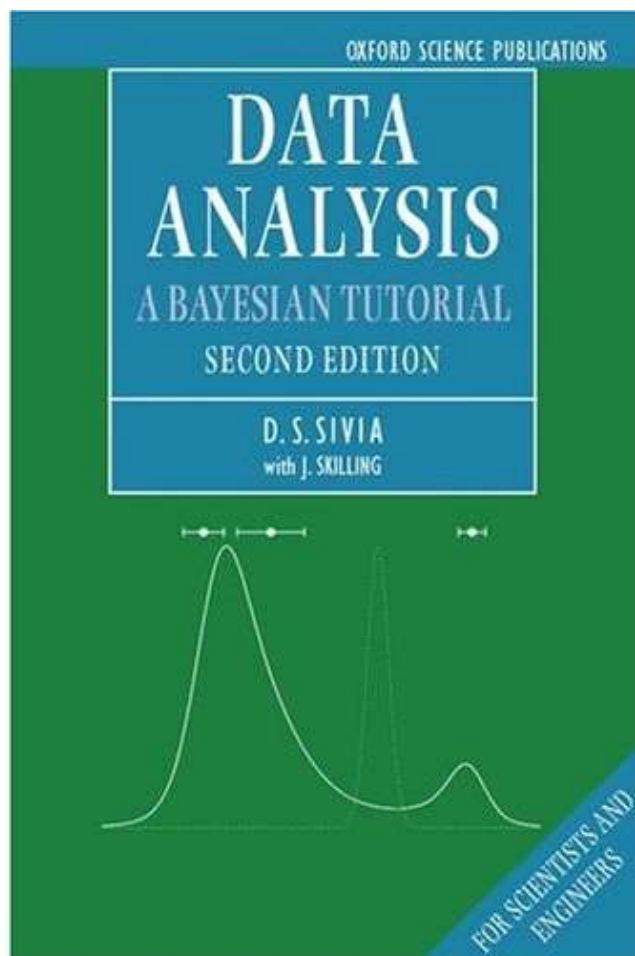


## 4. Parameter Estimation and Goodness of Fit - part two



Sivia Chapter 3 gives a very clear discussion of least squares fitting within a Bayesian framework.

In particular, contrasts, for Gaussian residuals:

- o known  $\sigma$
- o unknown  $\sigma \rightarrow$  Student's  $t$



# The principle of maximum likelihood

## Frequentist approach:

*A parameter is a fixed (but unknown) constant*

From actual data we can compute Likelihood,

$L$  = probability of obtaining the observed data, given the value of the parameter  $\theta$



# The principle of maximum likelihood

## Frequentist approach:

*A parameter is a fixed (but unknown) constant*

From actual data we can compute Likelihood,

$L$  = probability of obtaining the observed data, given the value of the parameter  $\theta$

Now define **likelihood function**: (infinite) family of curves generated by regarding  $L$  as a function of  $\theta$ , for data fixed.

## Principle of Maximum Likelihood

A good estimator of  $\theta$  maximises  $L$  -

$$\text{i.e. } \frac{\partial L}{\partial \theta} = 0 \quad \text{and} \quad \frac{\partial^2 L}{\partial \theta^2} < 0$$



# The principle of maximum likelihood

## Frequentist approach:

*A parameter is a fixed (but unknown) constant*

From actual data we can compute Likelihood,

$L$  = probability of obtaining the observed data, given the value of the parameter  $\theta$

Now define **likelihood function**: (infinite) family of curves generated by regarding  $L$  as a function of  $\theta$ , for data fixed.

## Principle of Maximum Likelihood

A good estimator of  $\theta$  maximises  $L$  -

$$\text{i.e. } \frac{\partial L}{\partial \theta} = 0 \quad \text{and} \quad \frac{\partial^2 L}{\partial \theta^2} < 0$$

We set the parameter equal to the value that makes the actual data sample we *did* observe - out of all the possible random samples we *could* have observed - the most likely.



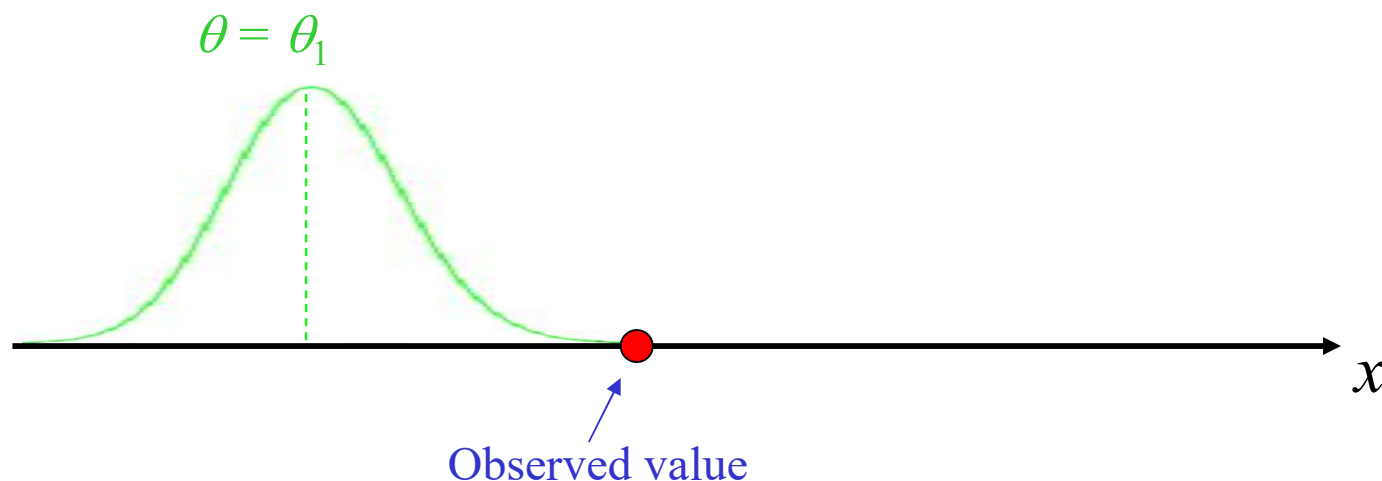
*Aside:* Likelihood function has same definition in Bayesian probability theory, but subtle difference in meaning and interpretation - no need to invoke idea of (infinite) ensemble of different samples.

---

## Principle of Maximum Likelihood

A good estimator of  $\theta$  maximises  $L$  -

$$\text{i.e. } \frac{\partial L}{\partial \theta} = 0 \quad \text{and} \quad \frac{\partial^2 L}{\partial \theta^2} < 0$$

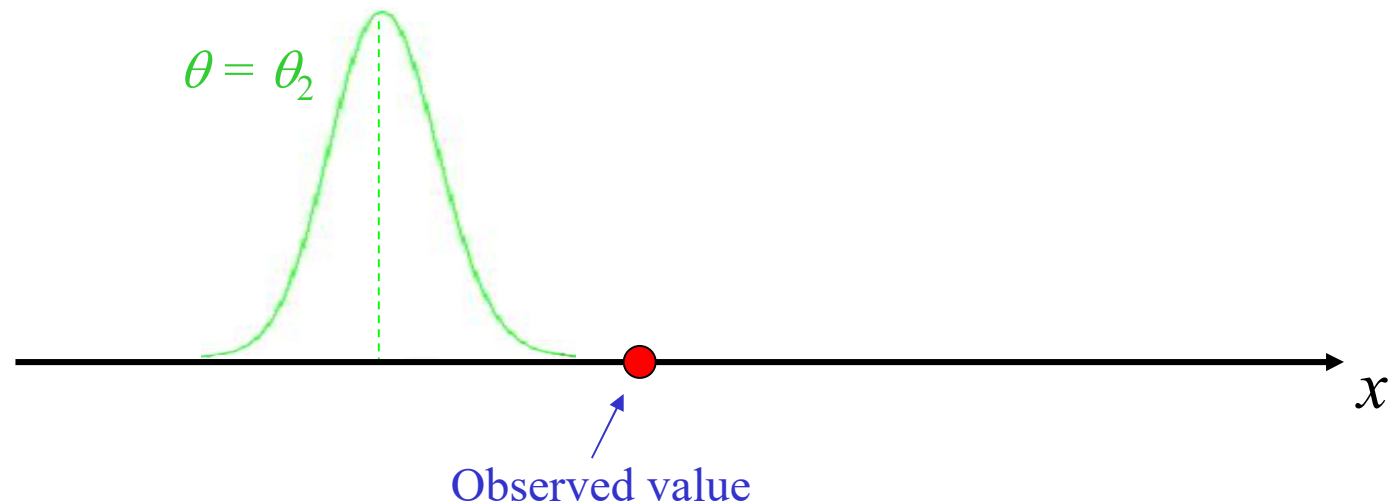




## Principle of Maximum Likelihood

A good estimator of  $\theta$  maximises  $L$  -

$$\text{i.e. } \frac{\partial L}{\partial \theta} = 0 \quad \text{and} \quad \frac{\partial^2 L}{\partial \theta^2} < 0$$

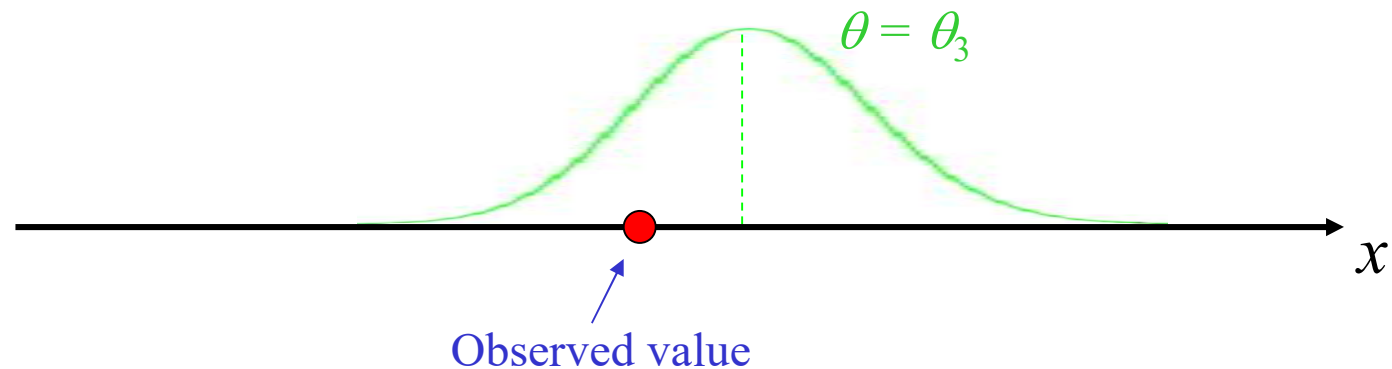




## Principle of Maximum Likelihood

A good estimator of  $\theta$  maximises  $L$  -

$$\text{i.e. } \frac{\partial L}{\partial \theta} = 0 \quad \text{and} \quad \frac{\partial^2 L}{\partial \theta^2} < 0$$

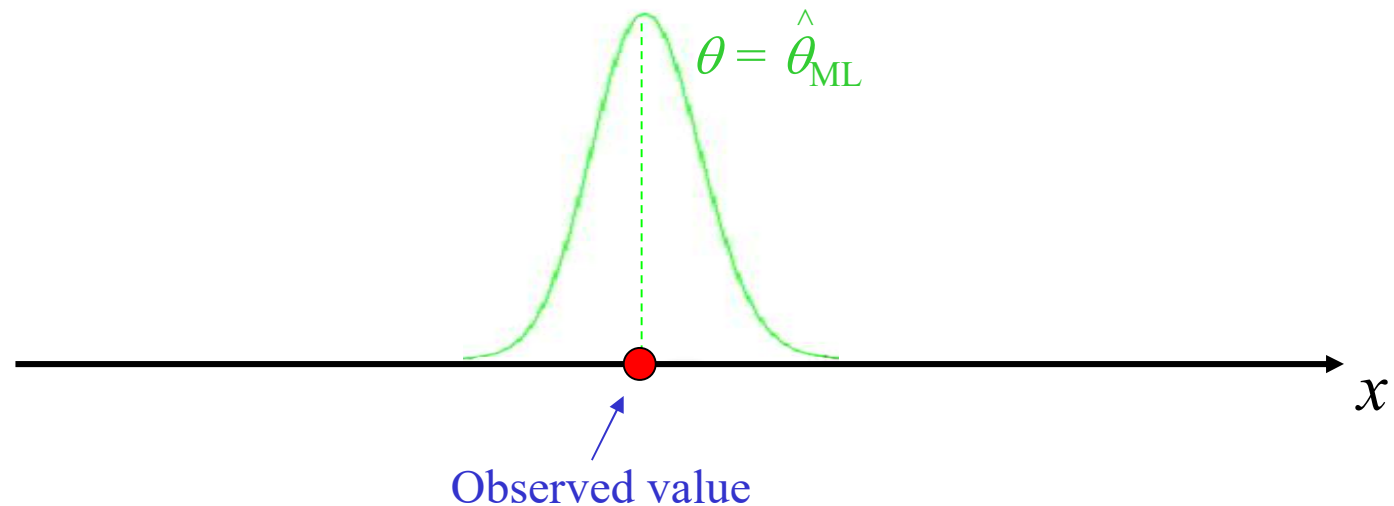




## Principle of Maximum Likelihood

A good estimator of  $\theta$  maximises  $L$  -

$$\text{i.e. } \frac{\partial L}{\partial \theta} = 0 \quad \text{and} \quad \frac{\partial^2 L}{\partial \theta^2} < 0$$





## Least squares as maximum likelihood estimators

To see the maximum likelihood method in action, let's consider again **weighted least squares** for the simple model  $y_i = a + bx_i + \epsilon_i$

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

Let's assume the pdf is a Gaussian



## Least squares as maximum likelihood estimators

To see the maximum likelihood method in action, let's consider again **weighted least squares** for the simple model  $y_i = a + bx_i + \epsilon_i$

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

Let's assume the pdf is a Gaussian

Likelihood

$$L = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left[-\frac{1}{2} \frac{\epsilon_i^2}{\sigma_i^2}\right]$$



**Question 7:** How can we justify writing the likelihood as a product?

$$L = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_i}} \exp\left[-\frac{1}{2} \frac{\varepsilon_i^2}{\sigma_i^2}\right]$$

- A** Because the residuals are all equal to each other
- B** Because the residuals are all Gaussian
- C** Because the residuals are all positive
- D** Because the residuals are all independent



## Least squares as maximum likelihood estimators

To see the maximum likelihood method in action, let's consider again **weighted least squares** for the simple model  $y_i = a + bx_i + \epsilon_i$

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

Let's assume the pdf is a Gaussian

Likelihood

$$L = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left[-\frac{1}{2} \frac{\epsilon_i^2}{\sigma_i^2}\right]$$

(note:  $L$  is a product of 1-D Gaussians because we are assuming the  $\epsilon_i$  are independent)



Substitute  $\varepsilon_i = y_i - a - bx_i$

$$\Rightarrow L = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left[-\frac{1}{2} \frac{(y_i - a - bx_i)^2}{\sigma_i^2}\right]$$

and the ML estimators of  $a$  and  $b$  satisfy  $\partial L / \partial a = 0$  and  $\partial L / \partial b = 0$



Substitute  $\varepsilon_i = y_i - a - bx_i$

$$\Rightarrow L = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left[-\frac{1}{2} \frac{(y_i - a - bx_i)^2}{\sigma_i^2}\right]$$

and the ML estimators of  $a$  and  $b$  satisfy  $\partial L / \partial a = 0$  and  $\partial L / \partial b = 0$

But maximising  $L$  is equivalent to maximising  $\ell = \ln L$

$$\begin{aligned} \text{Here } \ell &= -\frac{n}{2} \ln(2\pi) - \ln \sum_{i=1}^n \sigma_i - \frac{1}{2} \sum_{i=1}^n \left( \frac{y_i - a - bx_i}{\sigma_i} \right)^2 \\ &= \text{constant} - \frac{1}{2} S \end{aligned}$$

This is exactly the same  
sum of squares we  
defined earlier



Substitute  $\varepsilon_i = y_i - a - bx_i$


$$\Rightarrow L = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left[-\frac{1}{2} \frac{(y_i - a - bx_i)^2}{\sigma_i^2}\right]$$

and the ML estimators of  $a$  and  $b$  satisfy  $\partial L / \partial a = 0$  and  $\partial L / \partial b = 0$

But maximising  $L$  is equivalent to maximising  $\ell = \ln L$

Here

$$\begin{aligned} \ell &= -\frac{n}{2} \ln(2\pi) - \ln \sum_{i=1}^n \sigma_i - \frac{1}{2} \sum_{i=1}^n \left( \frac{y_i - a - bx_i}{\sigma_i} \right)^2 \\ &= \text{constant} - \frac{1}{2} S \end{aligned}$$

 This is exactly the same sum of squares we defined earlier

So in this case maximising  $L$  is **exactly equivalent** to minimising the sum of squares.

i.e. for Gaussian, independent errors, ML and weighted LS estimators are identical.



# Example application

## Measurement of the neutrino velocity with the OPERA detector in the CNGS beam\*

T. Adam<sup>a</sup>, N. Agafonova<sup>a</sup>, A. Aleksandrov<sup>c,1</sup>, O. Altinok<sup>d</sup>, P. Alvarez Sanchez<sup>e</sup>, A. Anokhina<sup>f</sup>, S. Aoki<sup>g</sup>, A. Ariga<sup>h</sup>, T. Ariga<sup>h</sup>, D. Autiero<sup>1</sup>, A. Badertscher<sup>1</sup>, A. Ben Dhabbi<sup>h</sup>, A. Bertolin<sup>k</sup>, C. Bozza<sup>1</sup>, T. Brugiè<sup>1</sup>, R. Brugnera<sup>m,k</sup>, F. Brunet<sup>g</sup>, G. Brunetti<sup>o,1,2</sup>, S. Buontempo<sup>o</sup>, B. Carls<sup>1</sup>, F. Cavanna<sup>h</sup>, A. Cazes<sup>1</sup>, L. Chaussard<sup>1</sup>, M. Chernyavsky<sup>1</sup>, V. Chiarella<sup>h</sup>, A. Chukanov<sup>1</sup>, G. Colosimo<sup>h</sup>, M. Crespi<sup>h</sup>, N. D'Ambrosio<sup>1</sup>, G. De Lellis<sup>w,e</sup>, M. De Serio<sup>h</sup>, Y. Déclais<sup>1</sup>, P. del Amo Sanchez<sup>h</sup>, F. Di Capua<sup>h</sup>, A. Di Crescenzo<sup>o,e</sup>, D. Di Ferdinando<sup>h</sup>, N. Di Marco<sup>o</sup>, S. Dmitrievsky<sup>1</sup>, M. Dracos<sup>1</sup>, D. Duchesneau<sup>h</sup>, S. Dusini<sup>h</sup>, J. Ebert<sup>1</sup>, I. Efthymiopoulos<sup>h</sup>, O. Egorov<sup>z</sup>, A. Ereditato<sup>h</sup>, L.S. Esposito<sup>1</sup>, J. Favier<sup>h</sup>, T. Ferber<sup>1</sup>, R.A. Fiani<sup>h</sup>, T. Fukuda<sup>ab</sup>, A. Garfagnini<sup>m,k</sup>, G. Giacomelli<sup>o,p</sup>, M. Giorgini<sup>o,p,3</sup>, M. Giovannozzi<sup>e</sup>, C. Girerd<sup>1</sup>, J. Goldberg<sup>ab</sup>, C. Göllnitz<sup>1</sup>, D. Golubkov<sup>z</sup>, L. Goncharova<sup>z</sup>, Y. Gornushkin<sup>1</sup>, G. Grella<sup>1</sup>, F. Grianti<sup>h,ac</sup>, E. Gschwendtner<sup>e</sup>, C. Guerin<sup>1</sup>, A.M. Güler<sup>d</sup>, C. Gustavino<sup>ad</sup>, C. Hagner<sup>1</sup>, K. Hamada<sup>ae</sup>, T. Hara<sup>g</sup>, M. Hierholzer<sup>1</sup>, A. Hollnagel<sup>1</sup>, M. Ieva<sup>h</sup>, H. Ishida<sup>h</sup>, K. Ishiguro<sup>ae</sup>, K. Jakovcic<sup>cd</sup>, C. Joller<sup>h</sup>, M. Jones<sup>1</sup>, F. Juget<sup>h</sup>, M. Kamiscioglu<sup>1</sup>, J. Kawada<sup>h</sup>, S.H. Kim<sup>ag,4</sup>, M. Kimura<sup>ah</sup>, E. Kiritsis<sup>ab</sup>, N. Kitagawa<sup>ae</sup>, B. Klicke<sup>cd</sup>, J. Knuesel<sup>h</sup>, K. Kodama<sup>ah</sup>, M. Komatsu<sup>ae</sup>, U. Kose<sup>1</sup>, I. Kreslo<sup>h</sup>, C. Lazzaro<sup>1</sup>, J. Lenkeit<sup>1</sup>, A. Ljubicic<sup>af</sup>, A. Longhin<sup>h</sup>, A. Malgin<sup>h</sup>, G. Mandrioli<sup>h</sup>, J. Marteau<sup>1</sup>, T. Matsuo<sup>ah</sup>, N. Mauri<sup>1</sup>, A. Mazzoni<sup>h</sup>, E. Medinaceli<sup>m,k</sup>, F. Meisel<sup>h</sup>, A. Mereaglia<sup>h</sup>, P. Migliozi<sup>h</sup>, S. Mikado<sup>ah</sup>, D. Missißen<sup>1</sup>, K. Morishima<sup>ae</sup>, U. Moser<sup>1</sup>, M.T. Muciaccia<sup>h,x</sup>, N. Naganawa<sup>ae</sup>, T. Naka<sup>ae</sup>, M. Nakamura<sup>ae</sup>, T. Nakano<sup>ae</sup>, Y. Nakatsuka<sup>ae</sup>, V. Nikitina<sup>h</sup>, F. Nitti<sup>ak</sup>, S. Ogawa<sup>aa</sup>, N. Okateva<sup>1</sup>, A. Olchevsky<sup>1</sup>, O. Palamara<sup>1</sup>, A. Paoloni<sup>1</sup>, B.D. Park<sup>ae,5</sup>, I.G. Park<sup>ae</sup>, A. Pastore<sup>h,x</sup>, L. Patrizi<sup>h</sup>, E. Pennacchio<sup>1</sup>, H. Pessard<sup>h</sup>, C. Pistillo<sup>h</sup>, N. Polukhina<sup>1</sup>, M. Pozzato<sup>o,p</sup>, K. Pretzl<sup>h</sup>, F. Pupilli<sup>1</sup>, R. Rescigno<sup>1</sup>, F. Riguzzi<sup>h</sup>, T. Roganova<sup>1</sup>, H. Rokujo<sup>h</sup>, G. Rosa<sup>am,ad</sup>, I. Rostovtseva<sup>1</sup>, A. Rubbia<sup>1</sup>, A. Russo<sup>o</sup>, O. Sato<sup>ae</sup>, Y. Sato<sup>ah</sup>, J. Schuler<sup>1</sup>, L. Scotto Lavina<sup>h,6</sup>, J. Serrano<sup>o</sup>, A. Sheshukov<sup>1</sup>, H. Shibuya<sup>ah</sup>, G. Shoziyoev<sup>1</sup>, S. Simone<sup>h,x</sup>, M. Sioli<sup>o,p</sup>, C. Sirignano<sup>h</sup>, G. Sirri<sup>h</sup>, J.S. Song<sup>ag</sup>, M. Spinetti<sup>h</sup>, L. Stanco<sup>h</sup>, N. Starkov<sup>1</sup>, S. Stellacci<sup>1</sup>, M. Stipcevic<sup>af</sup>, T. Strauss<sup>h</sup>, S. Takahashi<sup>h</sup>, M. Tenti<sup>o,p,1</sup>, F. Terranova<sup>h,ao</sup>, I. Tezuka<sup>ah</sup>, V. Tioukov<sup>1</sup>, P. Tolun<sup>1</sup>, N.T. Tran<sup>1</sup>, S. Tufanli<sup>h</sup>, P. Vilain<sup>op</sup>, M. Vladimirov<sup>1</sup>, L. Votano<sup>1</sup>, J.-L. Vuilleumier<sup>h</sup>, G. Wilquet<sup>op</sup>, B. Wonsak<sup>y</sup>, J. Wurtz<sup>1</sup>, C.S. Yoon<sup>ag</sup>, J. Yoshida<sup>ah</sup>, Y. Zaitsev<sup>z</sup>, S. Zemskova<sup>1</sup>, A. Zghiche<sup>h</sup>

<sup>a</sup> IPHC, Université de Strasbourg, CNRS/IN2P3, F-67037 Strasbourg, France

<sup>b</sup> INR-Institute for Nuclear Research of the Russian Academy of Sciences, RU-327312 Moscow, Russia

<sup>c</sup> INFN Sezione di Napoli, I-80125 Napoli, Italy

<sup>d</sup> METU-Middle East Technical University, TR-06532 Ankara, Turkey

<sup>e</sup> European Organization for Nuclear Research (CERN), Geneva, Switzerland

<sup>f</sup> (MSU SINP) Lomonosov Moscow State University Skobeltsyn Institute of Nuclear Physics, RU-119992 Moscow, Russia

<sup>g</sup> Kobe University, J-657-8501 Kobe, Japan

<sup>h</sup> Albert Einstein Center for Fundamental Physics, Laboratory for High Energy Physics (LHEP), University of Bern, CH-3012 Bern, Switzerland

<sup>i</sup> IPNL, Université Claude Bernard Lyon I, CNRS/IN2P3, F-69622 Villeurbanne, France

<sup>j</sup> ETH Zurich, Institute for Particle Physics, CH-8093 Zurich, Switzerland

\* Preprint submitted to the Journal of High Energy Physics (17 November 2011)

<sup>1</sup> Corresponding author [Dario.Autiero@cern.ch](mailto:Dario.Autiero@cern.ch)

## Cern test 'breaks speed of light'

**0.0024 seconds**

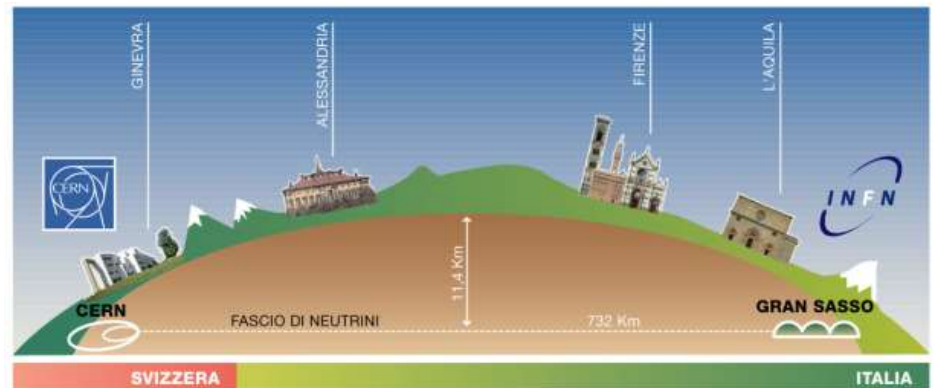
time taken by neutrinos

**0.00000006 seconds**

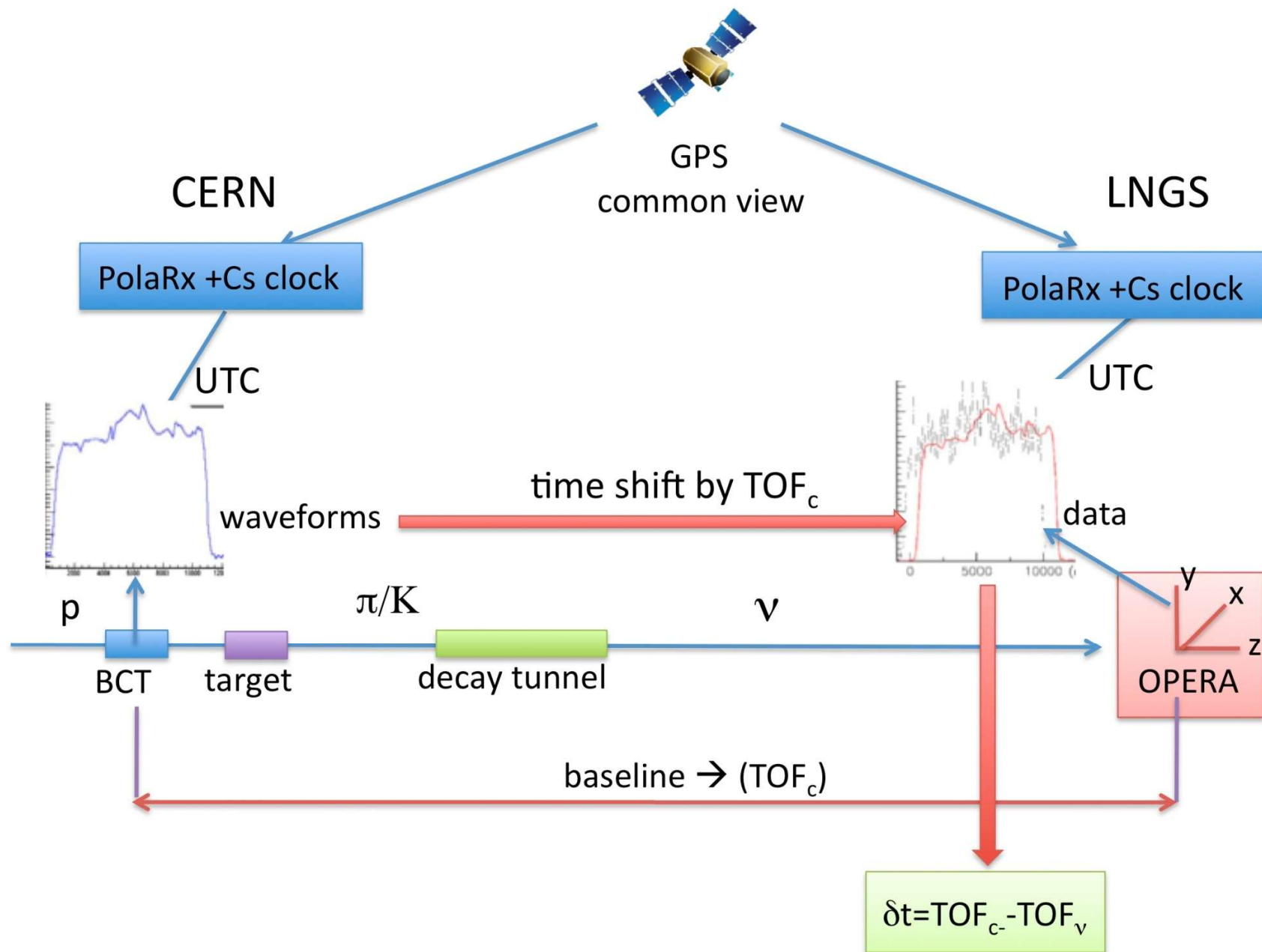
faster than the expected time

**732 km**

distance travelled through rock



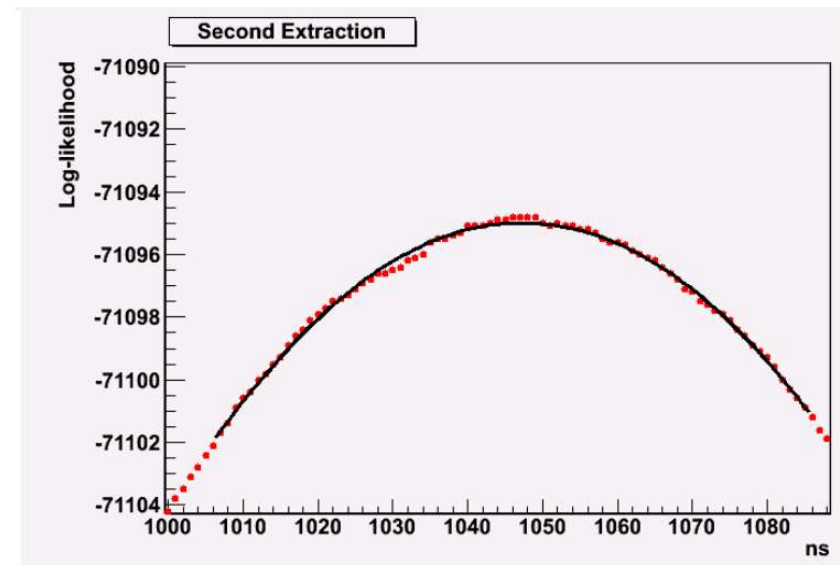
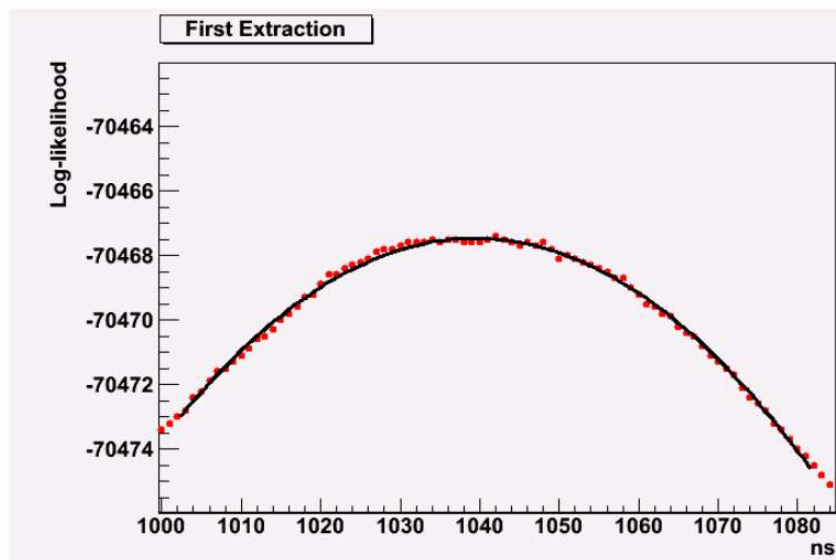






$$L_k(\delta t_k) = \prod_j w_k(t_j + \delta t_k) \quad k = 1, 2 \text{ extractions}$$

Near the maximum the likelihood function can be approximated by a Gaussian whose variance is a measure of the statistical uncertainty on  $\delta t$ . The data used for the maximum likelihood calculation are unbinned and the dependence on  $\delta t$  is computed by making a scan in steps of 1 ns. A parabolic fit is performed on the log-likelihood function for the evaluation of the maximum and of the statistical uncertainty (Fig. 10). As seen in Fig. 11, the PDF representing the time-structure of the proton extraction is not flat but exhibits a series of peaks and valleys, reflecting the features and the inefficiencies of the proton extraction from the PS to the SPS via the Continuous Transfer mechanism [41]. Such structures may well change with time. The way the PDF are built automatically accounts for the beam conditions corresponding to the neutrino interactions detected by OPERA.





# Hypothesis testing

In the previous section we have discussed how to estimate parameters of an underlying pdf model from sample data.

We now consider the closely related question:

*How good is our pdf model in the first place?*



# Hypothesis testing

In the previous section we have discussed how to estimate parameters of an underlying pdf model from sample data.

We now consider the closely related question:

*How good is our pdf model in the first place?*

## Simple example.

Null hypothesis:

sampled data are drawn from a normal pdf, with mean  $\mu_{\text{model}}$  and variance  $\sigma^2$ .

We want to **test** this null hypothesis: are our data consistent with it?



Assume (for the moment) that  $\sigma^2$  is known.

### Example

Measured data:  $\{x_i : i = 1, \dots, 10\}$   $\sum_{i=1}^{10} x_i = 47.8$

Null hypothesis:  $x \sim N(\mu, \sigma^2)$  with  $\mu_{\text{model}} = 4$

Assume:  $\sigma = 2$   $\sigma_{\mu}^2 = 0.4$

Under NH, sample mean

$$\bar{x}_{\text{model}} \sim N(4, 2^2/10)$$

Observed sample mean  $\bar{x}_{\text{obs}} = 4.78$



We transform to a standard normal variable

Under NH: 
$$Z = \left( \frac{\bar{x}_{\text{obs}} - \bar{x}_{\text{model}}}{\sigma_{\mu}} \right) \sim N(0,1)$$

From our measured data: 
$$Z_{\text{obs}} = \frac{4.78 - 4}{\sqrt{0.4}} = 1.233$$

**If** NH is true, how probable is it that we would obtain a value of  $Z_{\text{obs}}$  as large as this, or larger?

We call this probability the **p-value**



**Question 8:** Suppose that  $X$  is sampled from a normal distribution with mean  $\mu = 5$  and variance  $\sigma^2 = 9$ .

Which of the following is a standard normal variable?

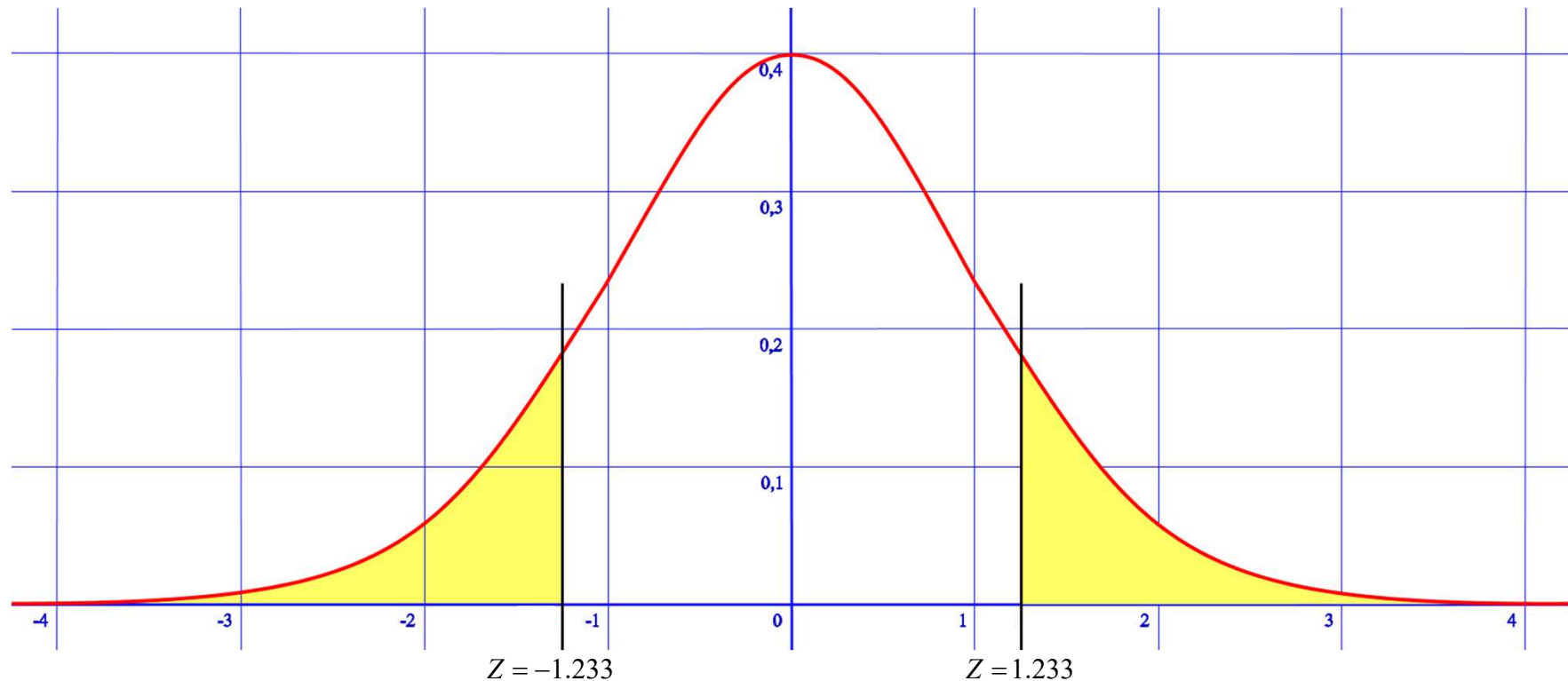
**A**  $Z = \frac{X - 5}{9}$

**B**  $Z = \frac{X - 5}{3}$

**C**  $Z = \frac{X - 3}{5}$

**D**  $Z = \frac{X - 3}{9}$



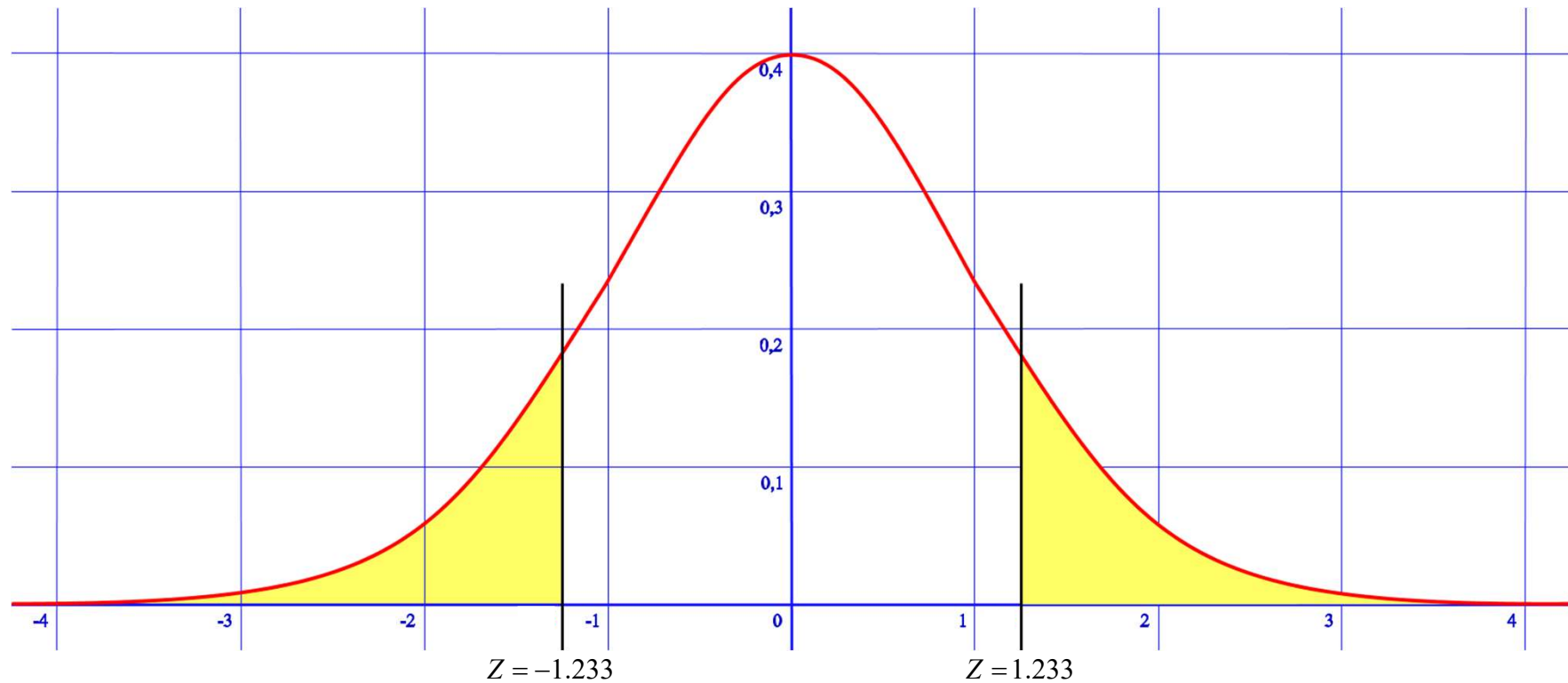


$$\text{p-value} = \text{Prob}(|Z| \geq |Z_{\text{obs}}|) = 1 - \int_{-Z_{\text{obs}}}^{Z_{\text{obs}}} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2} z^2\right) dz$$

Simple programs to perform this probability integral (and many others) can be found in numerical recipes, or built into e.g. MATLAB or MAPLE.

Java applets also available online at <http://statpages.org/pdfs.html> ([here](#)).





$$\text{p-value} = \text{Prob}(|Z| \geq |Z_{\text{obs}}|) = 0.2176$$

The **smaller** the p-value, the less credible is the null hypothesis.



We can also carry out a *one-tailed* hypothesis test, if appropriate, and for statistics with other sampling distributions.

**Question 9:** A one-tailed hypothesis test is carried out. Under the  $H_0$  the test statistic has a uniform distribution  $U[0,1]$ .

The observed value of the test statistic is 0.8.

**The p-value is:**

**A** 0.8

**B** 0.9

**C** 0.2

**D** 0.1



What if we don't assume that  $\sigma^2$  is known?

We can estimate it from our observed data (provided  $n \geq 2$  )

We form the statistic  $t_{\text{obs}} = \left( \frac{\bar{x}_{\text{obs}} - \bar{x}_{\text{model}}}{\hat{\sigma}_{\mu}} \right)$

where 
$$\hat{\sigma}_{\mu}^2 = \frac{1}{n(n-1)} \sum_{i=1}^n (x_i - \bar{x}_{\text{obs}})^2$$

Accounts for the fact that we don't know  $\mu$  , but must use  $\bar{x}_{\text{obs}}$  when we estimate  $\sigma_{\mu}$

However, now  $t_{\text{obs}}$  no longer has a normal distribution.

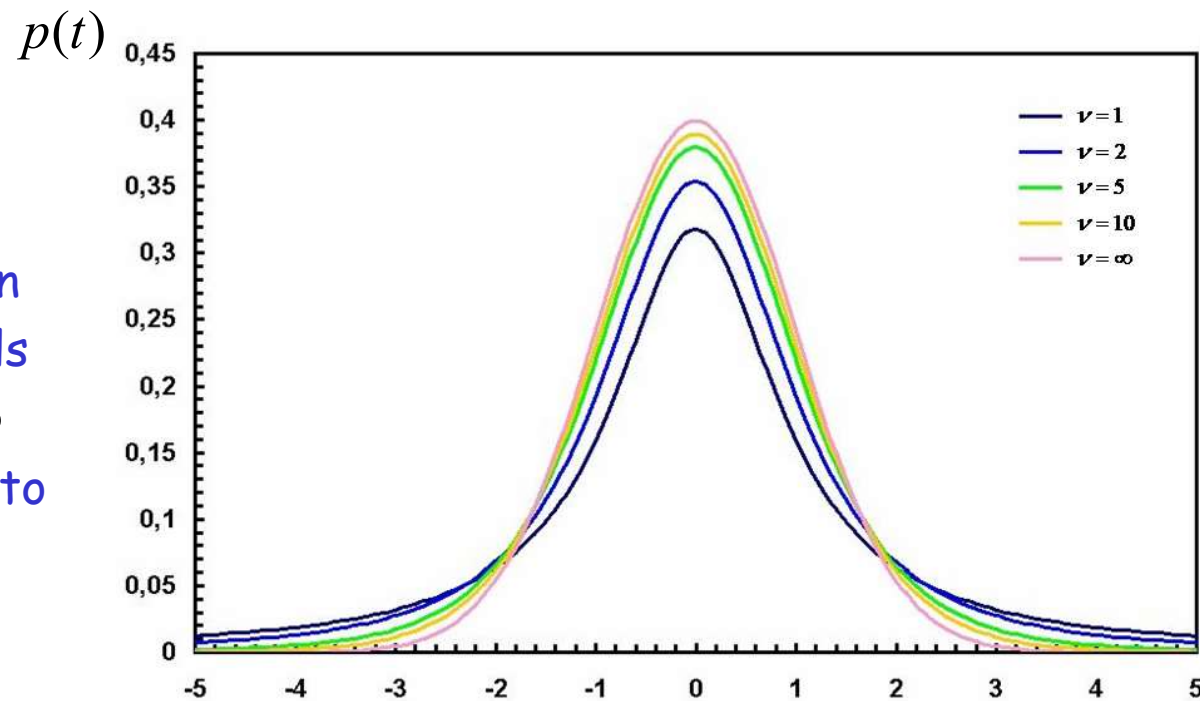


In fact  $t_{\text{obs}}$  has a pdf known as the **Student's t distribution**

$$p(t) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\nu\pi}\Gamma\left(\frac{\nu}{2}\right)} \left(1 + \frac{t^2}{\nu}\right)^{-\left(\frac{\nu+1}{2}\right)}$$

where  $\nu = n - 1$  is the **no. degrees of freedom** and  $\Gamma(\nu) = \int_0^{\infty} x^{\nu-1} e^{-x} dx$

For small  $n$  the Student's t distribution has more extended tails than  $Z$ , but as  $n \rightarrow \infty$  the distribution tends to  $N(0,1)$





**Question 10:** The more extended tails of the students' t distribution mean that, under the null hypothesis

- A** larger values of the test statistic are more likely
- B** larger values of the test statistic are less likely
- C** smaller values of the test statistic are more likely
- D** smaller values of the test statistic are less likely



# Hypothesis tests and decision theory

In a **simple hypothesis test**, we test our **null hypothesis** against a single **alternative hypothesis**.

We choose a **critical region**: set of values of the test statistic for which we choose to reject the **NH** and accept the **AH**.

This means we need to consider the distribution of our test statistic under the NH and the AH.



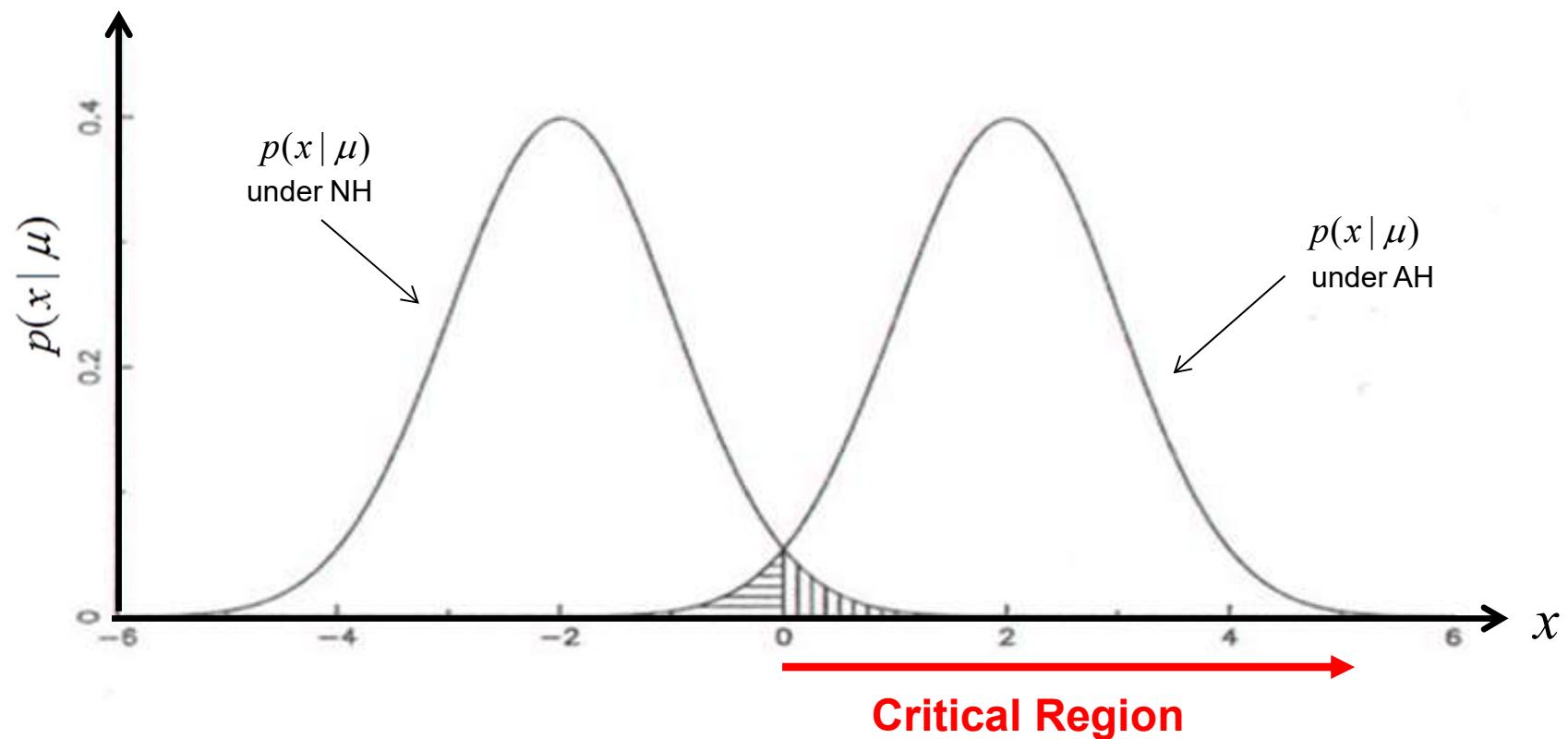
## Example

Measured data leads to test statistic  $x$ , estimator of  $\mu$  :

**NH:**  $\mu = -2$ ;  $x \sim N(-2, 1)$

**AH:**  $\mu = +2$ ;  $x \sim N(2, 1)$

Suppose we choose critical region to be  $x > 0$

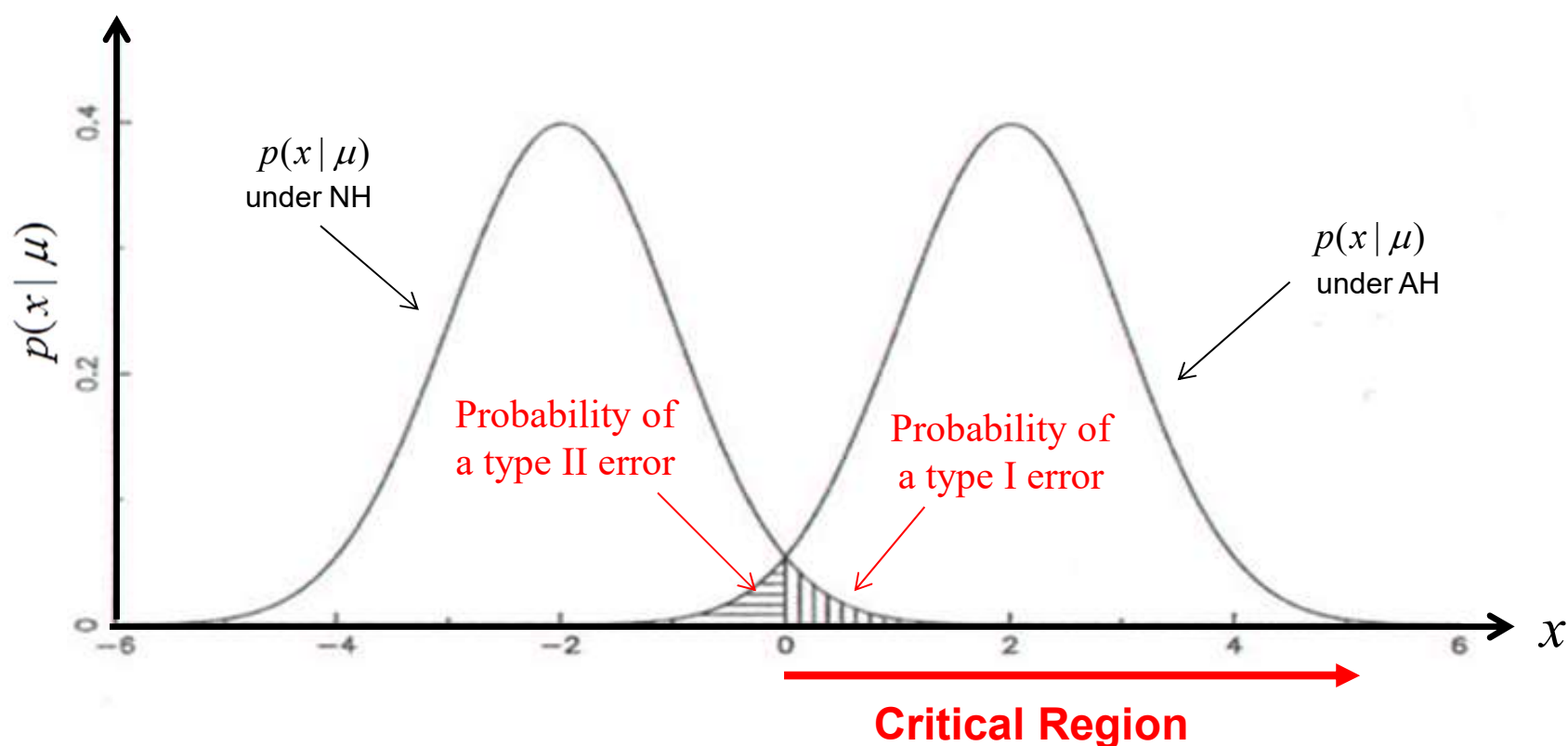




There are two ways in which we can make an incorrect decision:

**Type I error:** we reject the  $H_0$  when it is **TRUE**

**Type II error:** we accept the  $H_0$  when it is **FALSE**

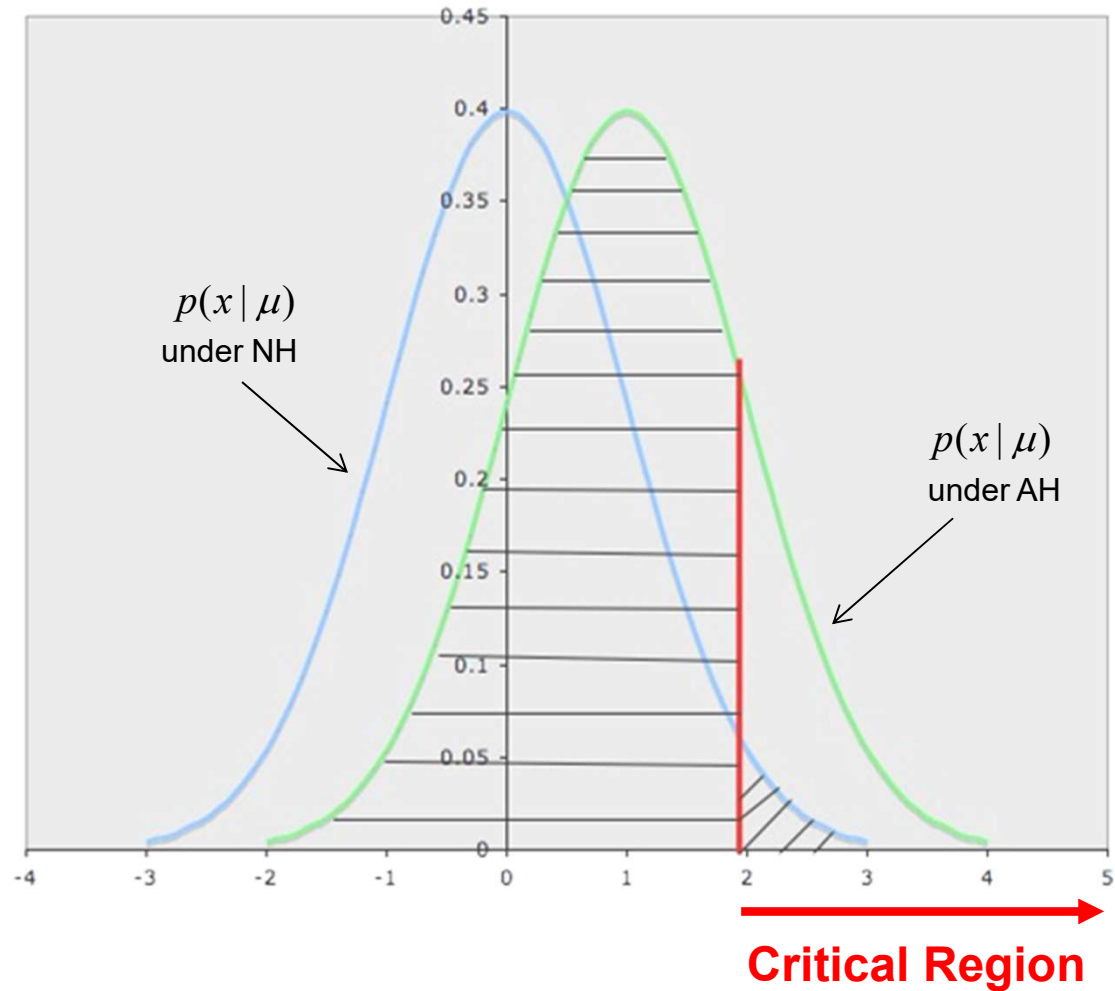




Consequences of incorrect decision will be worse when the sampling distributions under  $H_0$  and  $H_1$  have a greater overlap.

This may influence our choice of **critical region**.

We want to reduce the probability of type I and type II errors - but we can't do *both* at the same time...





Lots of terminology:

**Type I error:** also known as **false alarm** or **false positive**

**Type II error:** also known as **false negative** or **"miss"**

**Sensitivity** = probability (rate) of obtaining **true positive** ("hit")

**Specificity** = probability (rate) of obtaining **true negative**

Sensitivity (or power) =  $1 - \text{prob}(\text{type II error})$

Specificity =  $1 - \text{prob}(\text{type I error})$



In many fields (particularly medical applications) a **Receiver Operating Characteristic (ROC)**

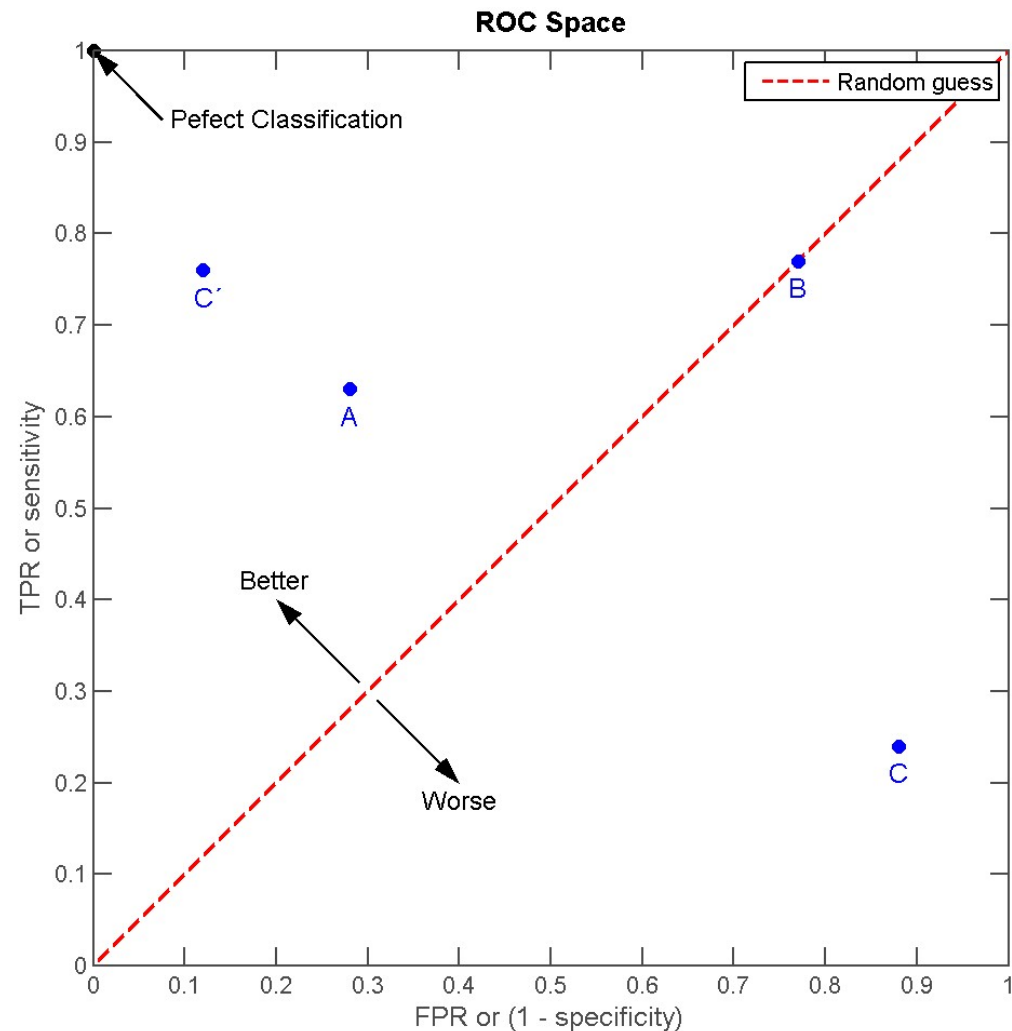
plot can be used to assess the performance of a simple hypothesis test.

x-axis = **(1 - specificity)**

y-axis = **sensitivity**

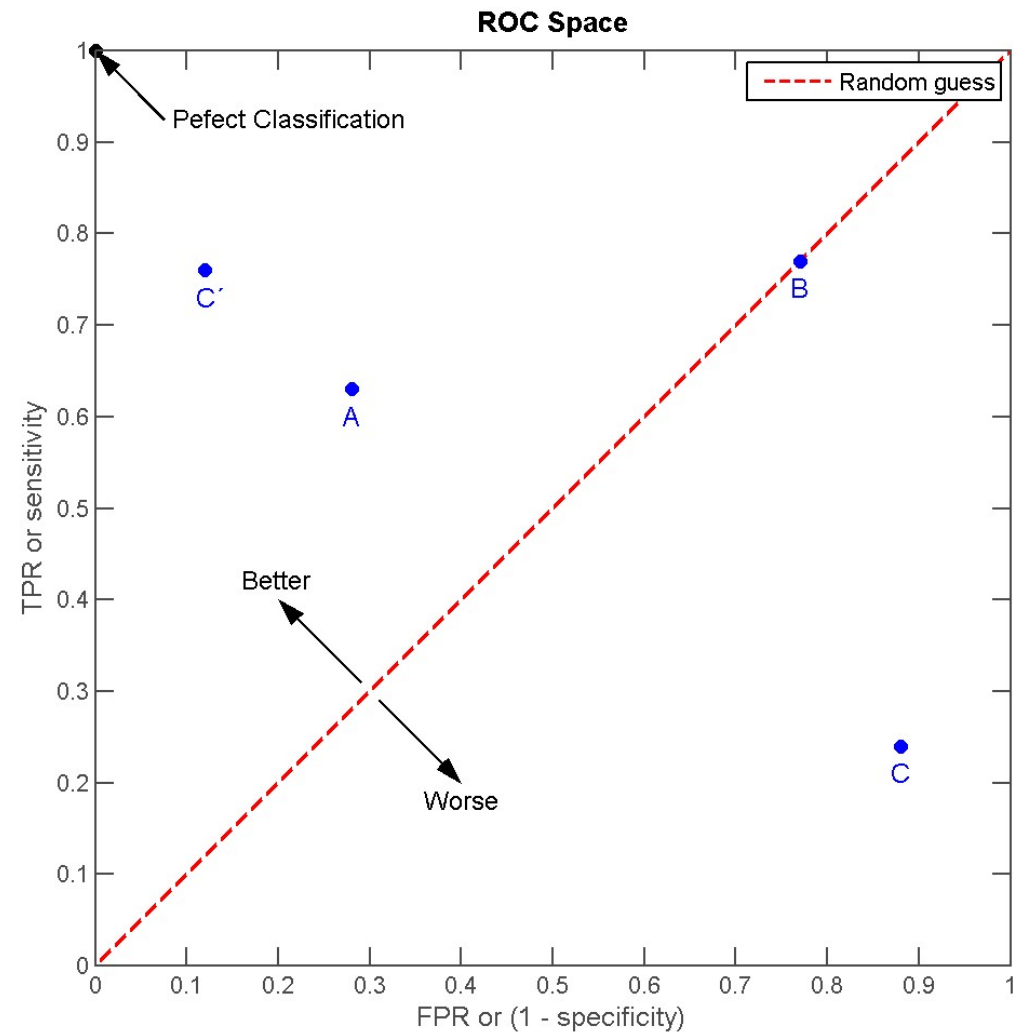
**Blue dots** = results for different hypothesis tests.

**Red diagonal** = what we'd expect from a random guess alone.





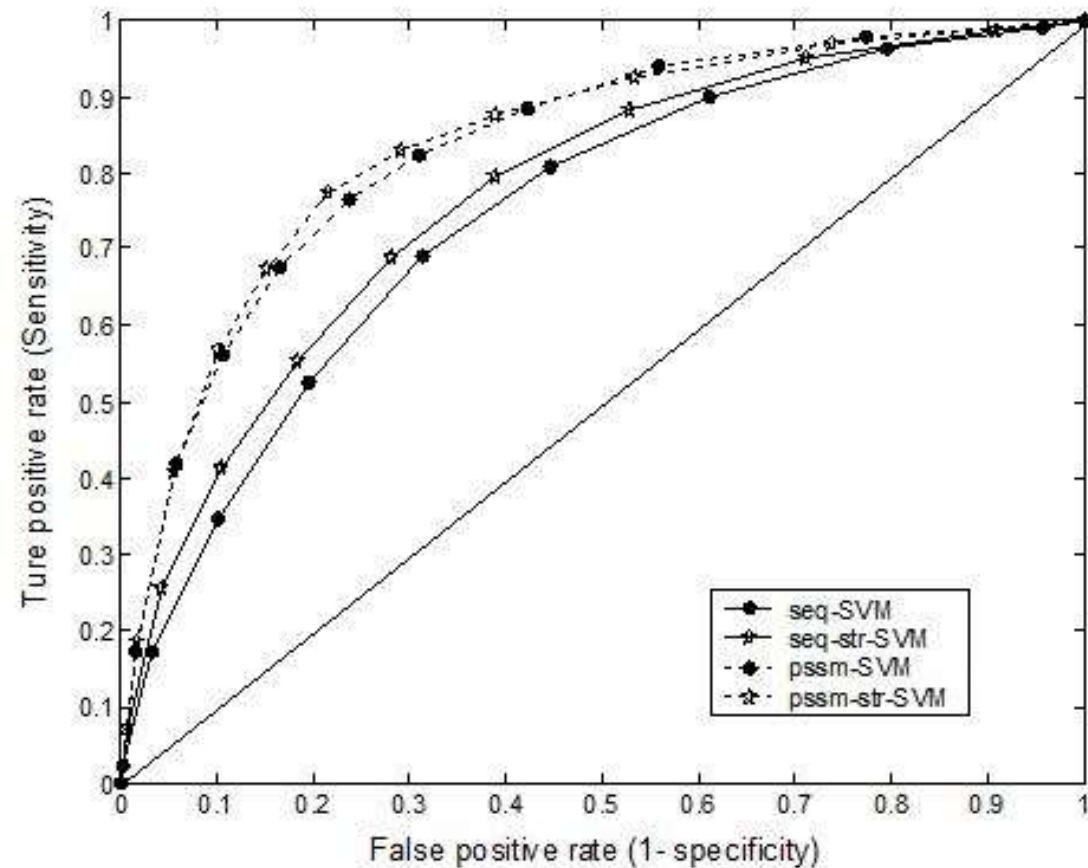
We can generate a ROC curve for a given test by varying the **critical region**. This can help to optimise our choice of CR - trading off type I and type II error, and getting us as close as possible to a perfect decision / classification.





We can generate a ROC curve for a given test by varying the **critical region**. This can help to optimise our choice of CR - trading off type I and type II error, and getting us as close as possible to a perfect decision / classification.

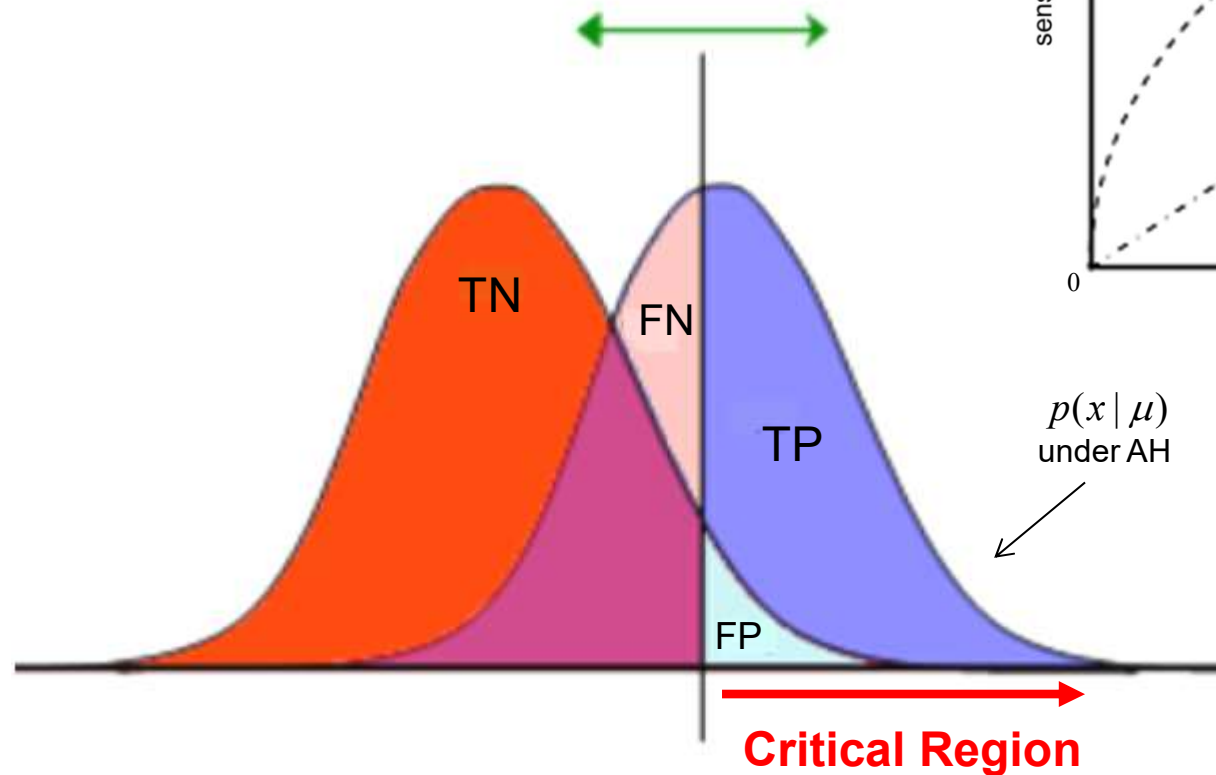
The area under the ROC curve can be used as a measure to compare different tests and choose the best.



*ROC curves for predictors of DNA-binding sites.  
SUNY Albany Center for Excellence in Cancer Genomics*



We can generate a ROC curve for a given test by varying the **critical region**. This can help to optimise our choice of CR - trading off type I and type II error, and getting us as close as possible to a perfect decision / classification.



As we move the CR from left to right, we move along the ROC curve (dashed line)

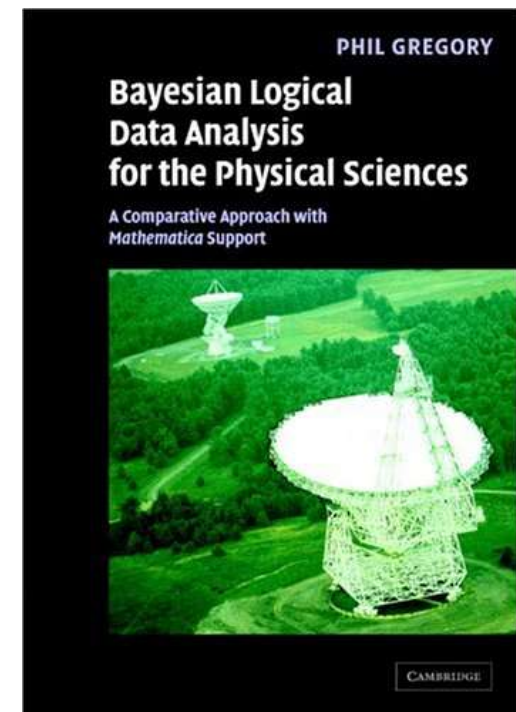


# Hypothesis testing

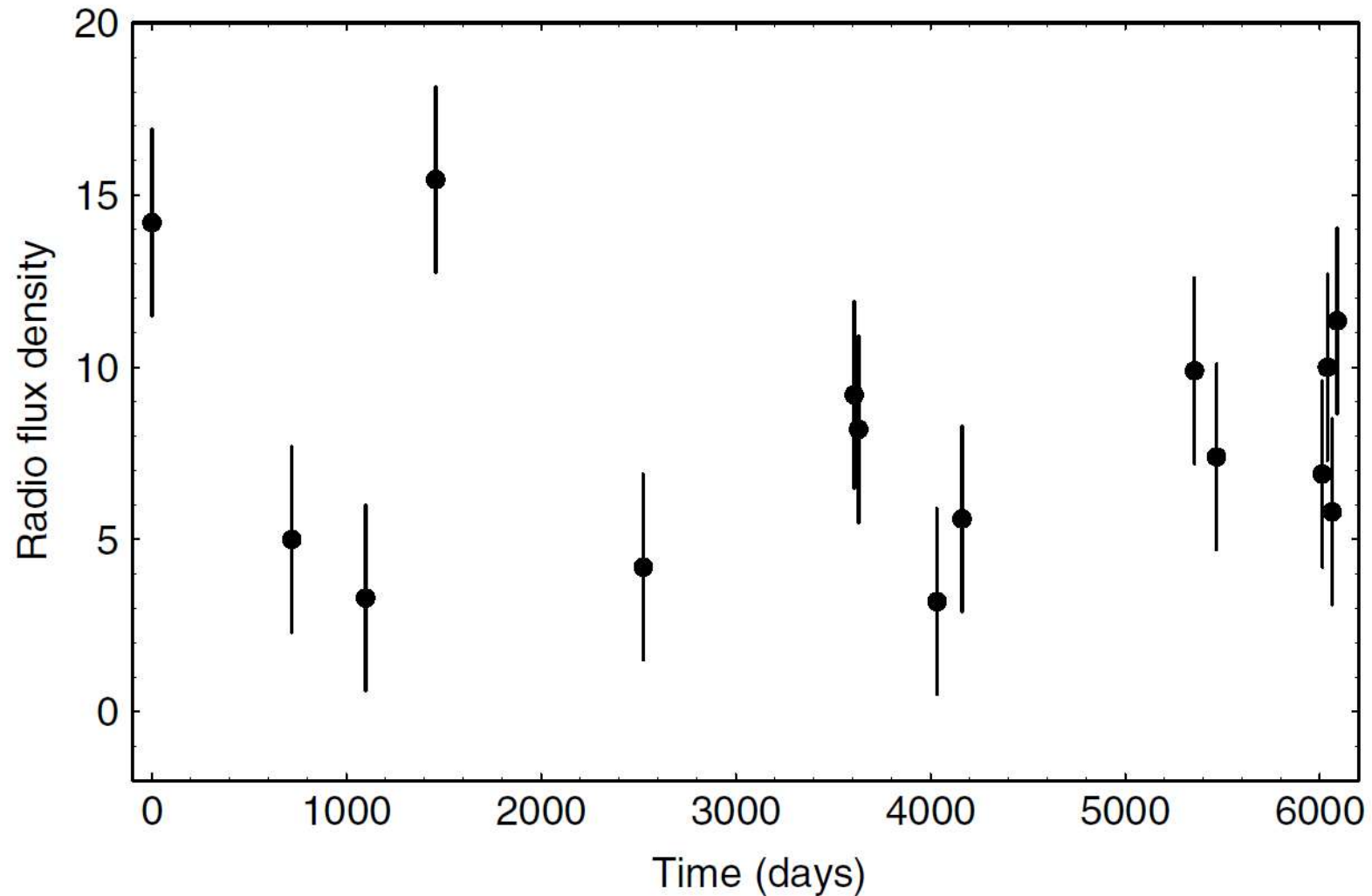
More generally, we now illustrate the frequentist approach to the question of how good is the fit to our model, using the **Chi-squared goodness of fit test**.

We take an example from Gregory (Chapter 7)

*(book focusses mainly on Bayesian probability, but is very good on frequentist approach too)*







Model: radio emission from a galaxy is constant in time.

Assume residuals are iid, drawn from  $N(0, \sigma)$



## Goodness-of-fit Test: the basic ideas

1. Choose as our null hypothesis that the galaxy has an unknown but constant flux density. If we can demonstrate that this hypothesis is absurd at say the 95% confidence level, then this provides indirect evidence that the radio emission is variable. Previous experience with the measurement apparatus indicates that the measurement errors are independently normal with a  $\sigma = 2.7$ .
2. Select a suitable statistic that (a) can be computed from the measurements, and (b) has a predictable distribution. More precisely, (b) means that we can predict the distribution of values of the statistic that we would expect to obtain from an infinite number of repeats of the above set of radio measurements under identical conditions. We will refer to these as our hypothetical reference set. More specifically, we are predicting a probability distribution for this reference set.

To refute the null hypothesis, we will need to show that scatter of the individual measurements about the mean is larger than would be expected from measurement errors alone.

3. Evaluate the  $\chi^2$  statistic from the measured data. Let's start with the expression for the  $\chi^2$  statistic for our data set:

$$\chi^2 = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{\sigma^2}$$

*From Gregory, pg. 164*

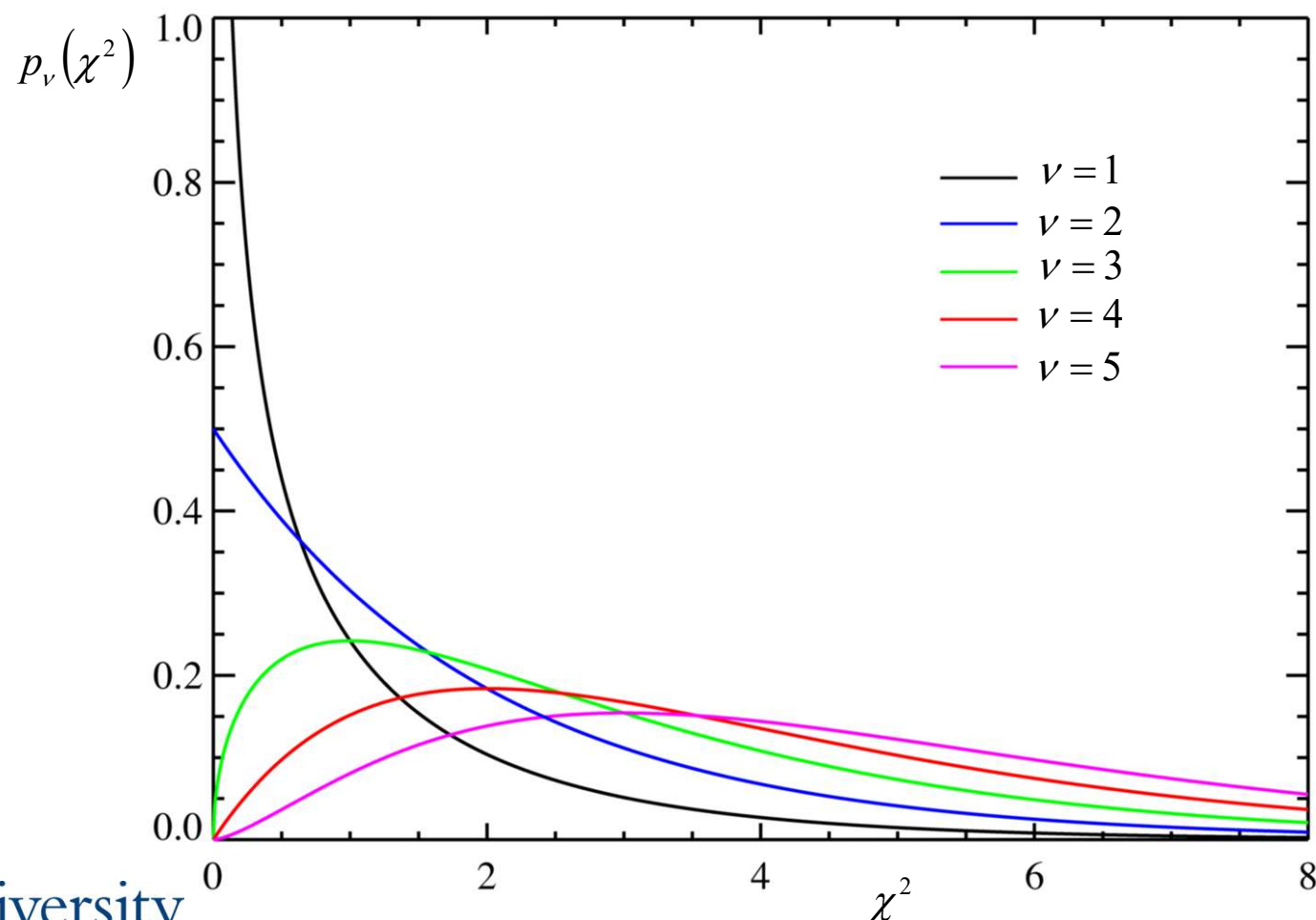


## The $\chi^2$ pdf

$$p_{\nu}(\chi^2) = p_0 \times (\chi^2)^{\frac{\nu}{2}-1} e^{-\chi^2/2}$$

Here  $\nu$  is known as the number of degrees of freedom of the pdf.

The mean value of the pdf is  $\nu$  and the variance is  $2\nu$ .



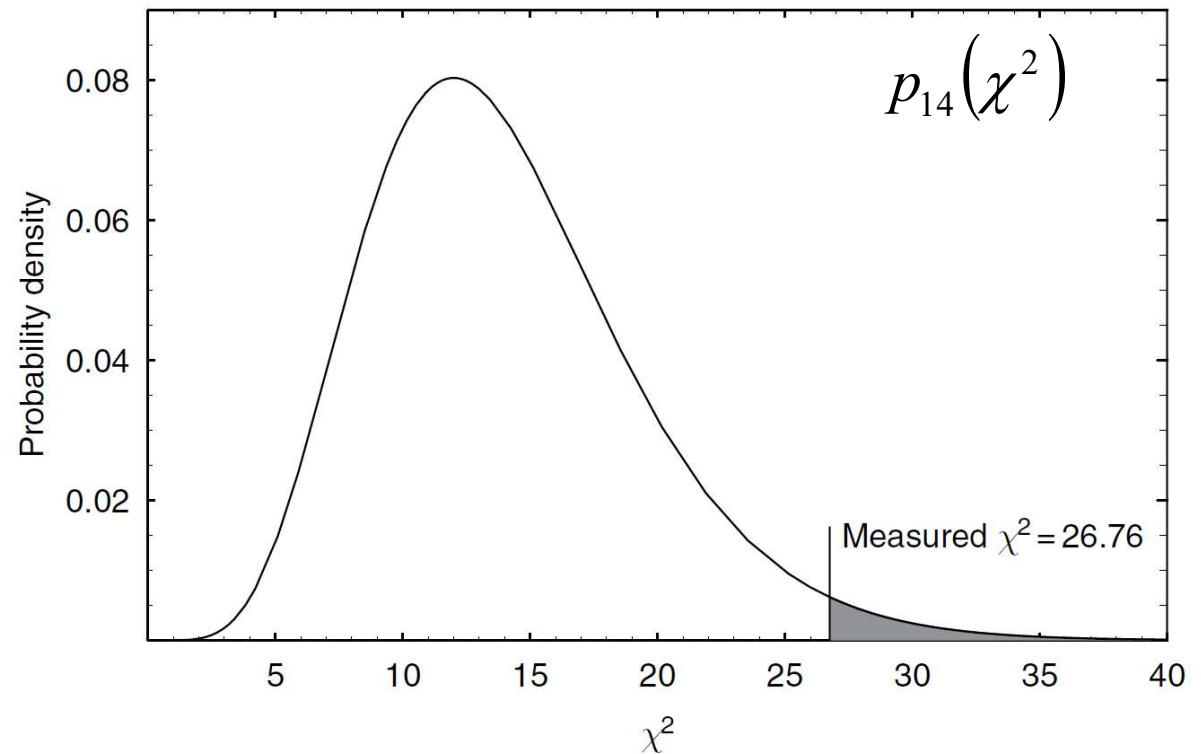


| Day Number | Flux Density (mJy) |
|------------|--------------------|
| 0.0        | 14.2               |
| 718.0      | 5.0                |
| 1097.0     | 3.3                |
| 1457.1     | 15.5               |
| 2524.1     | 4.2                |
| 3607.7     | 9.2                |
| 3630.1     | 8.2                |
| 4033.1     | 3.2                |
| 4161.3     | 5.6                |
| 5355.9     | 9.9                |
| 5469.1     | 7.4                |
| 6012.4     | 6.9                |
| 6038.3     | 10.0               |
| 6063.2     | 5.8                |
| 6089.3     | 11.4               |

$$\chi^2 = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{\sigma^2} = \sum_{i=1}^n \frac{(x_i - 7.98)^2}{2.7^2} = 26.76$$

$n = 15$  data points, but  $\nu = 14$  degrees of freedom, because  $\chi^2$  statistic involves the *sample mean* and not the true mean.

We subtract one d.o.f. to account for this.





**Question 11:** Given that the mean and variance of a chi-squared distribution with  $n$  degrees of freedom are  $n$  and  $2n$  respectively, in the Gregory example (with 14 degrees of freedom) estimate the number of sigma by which the value  $\chi^2_{\text{obs}} = 26.76$  exceeds the expected value.

- A** between 0 and 1 sigma
- B** between 1 and 2 sigma
- C** between 2 and 3 sigma
- D** between 3 and 4 sigma

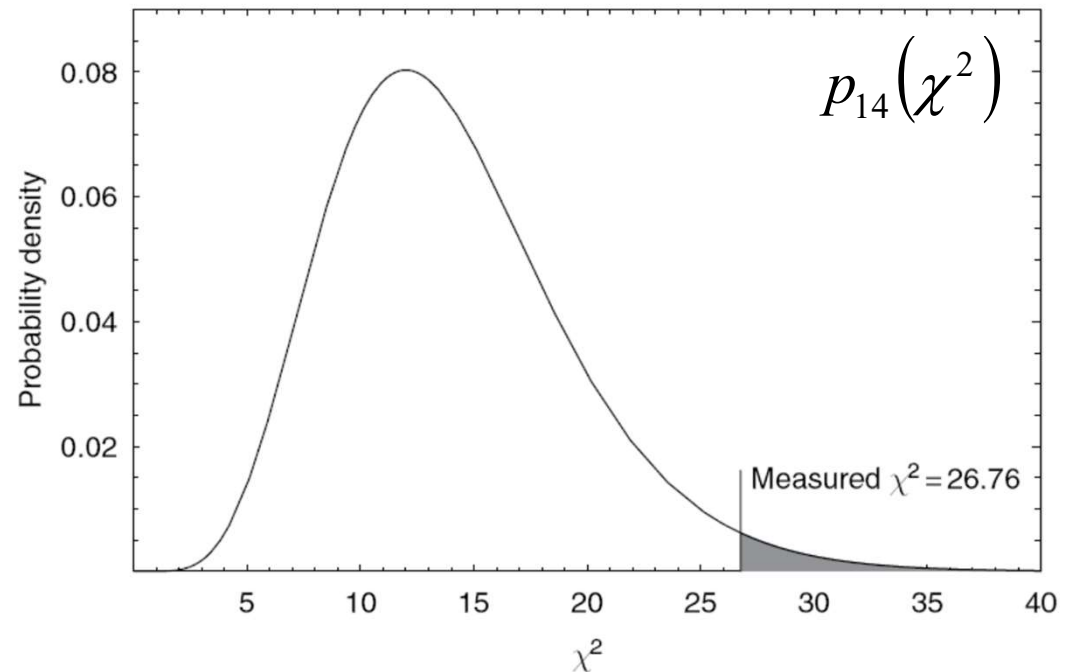


$$\chi^2 = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{\sigma^2} = \sum_{i=1}^n \frac{(x_i - 7.98)^2}{2.7^2} = 26.76.$$

$n = 15$  data points, but  $\nu = 14$  degrees of freedom, because  $\chi^2$  statistic involves the *sample mean* and not the true mean.

We subtract one d.o.f. to account for this.

| Day Number | Flux Density (mJy) |
|------------|--------------------|
| 0.0        | 14.2               |
| 718.0      | 5.0                |
| 1097.0     | 3.3                |
| 1457.1     | 15.5               |
| 2524.1     | 4.2                |
| 3607.7     | 9.2                |
| 3630.1     | 8.2                |
| 4033.1     | 3.2                |
| 4161.3     | 5.6                |
| 5355.9     | 9.9                |
| 5469.1     | 7.4                |
| 6012.4     | 6.9                |
| 6038.3     | 10.0               |
| 6063.2     | 5.8                |
| 6089.3     | 11.4               |



*If the null hypothesis is true, how probable is it that we would measure as large, or larger, a value of  $\chi^2$ ?*

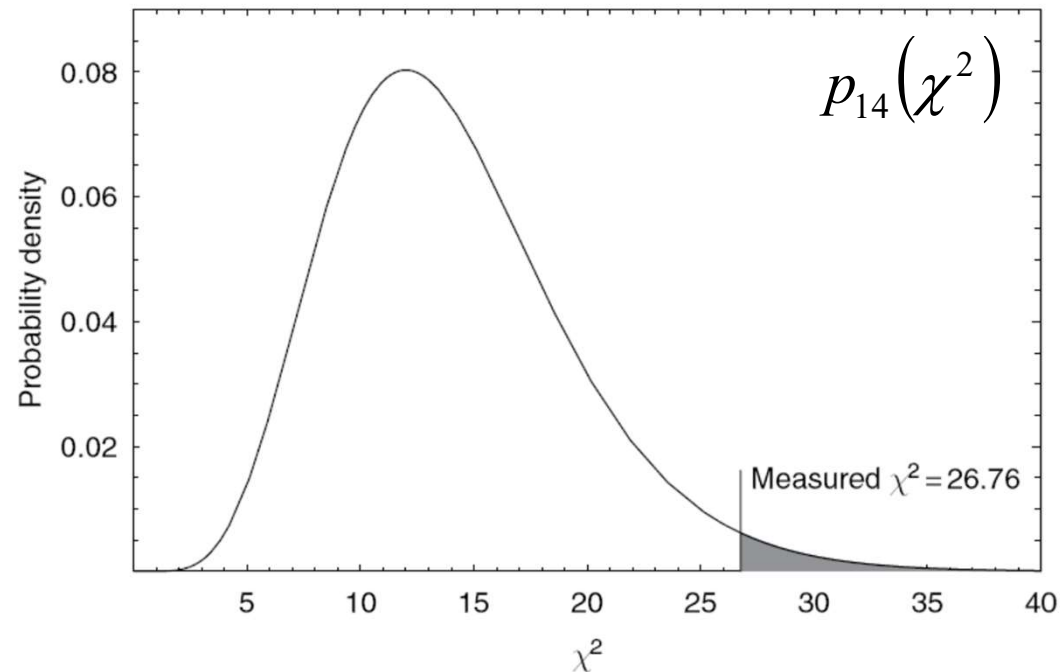


$$\chi^2 = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{\sigma^2} = \sum_{i=1}^n \frac{(x_i - 7.98)^2}{2.7^2} = 26.76.$$

$n = 15$  data points, but  $\nu = 14$  degrees of freedom, because  $\chi^2$  statistic involves the *sample mean* and not the true mean.

We subtract one d.o.f. to account for this.

| Day Number | Flux Density (mJy) |
|------------|--------------------|
| 0.0        | 14.2               |
| 718.0      | 5.0                |
| 1097.0     | 3.3                |
| 1457.1     | 15.5               |
| 2524.1     | 4.2                |
| 3607.7     | 9.2                |
| 3630.1     | 8.2                |
| 4033.1     | 3.2                |
| 4161.3     | 5.6                |
| 5355.9     | 9.9                |
| 5469.1     | 7.4                |
| 6012.4     | 6.9                |
| 6038.3     | 10.0               |
| 6063.2     | 5.8                |
| 6089.3     | 11.4               |



*If the null hypothesis is true, how probable is it that we would measure as large, or larger, a value of  $\chi^2$ ?*



If the null hypothesis were true, how probable is it that we would measure as large, or larger, a value of  $\chi^2$  ?

Recall that we refer to this important quantity as the **p-value**

$$\text{p-value} = 1 - P(\chi_{\text{obs}}^2) = 1 - \int_0^{\chi_{\text{obs}}^2} p_0 x^{\frac{\nu}{2}-1} \exp\left(-\frac{x}{2}\right) dx = 0.02$$



If the null hypothesis were true, how probable is it that we would measure as large, or larger, a value of  $\chi^2$  ?

Recall that we refer to this important quantity as the **p-value**

$$\text{p-value} = 1 - P(\chi_{\text{obs}}^2) = 1 - \int_0^{\chi_{\text{obs}}^2} p_0 x^{\frac{\nu}{2}-1} \exp\left(-\frac{x}{2}\right) dx = 0.02$$

What precisely does the p-value mean?

“If the galaxy flux density really *is* constant, and we repeatedly obtained sets of 15 measurements under the same conditions, then only 2% of the  $\chi^2$  values derived from these sets would be expected to be greater than our one actual measured value of 26.76”

*From Gregory, pg. 165*

If we obtain a very small p-value (e.g. a few percent?) we can interpret this as providing little support for the null hypothesis, which we may then **choose to reject**.

*(Ultimately this choice is subjective, but  $\chi^2$  may provide objective ammunition for doing so)*



If the null hypothesis were true, how probable is it that we would measure as large, or larger, a value of  $\chi^2$  ?

Recall that we refer to this important quantity as the **p-value**

$$\text{p-value} = 1 - P(\chi_{\text{obs}}^2) = 1 - \int_0^{\chi_{\text{obs}}^2} p_0 x^{\frac{\nu}{2}-1} \exp\left(-\frac{x}{2}\right) dx = 0.02$$

What precisely does the p-value mean?

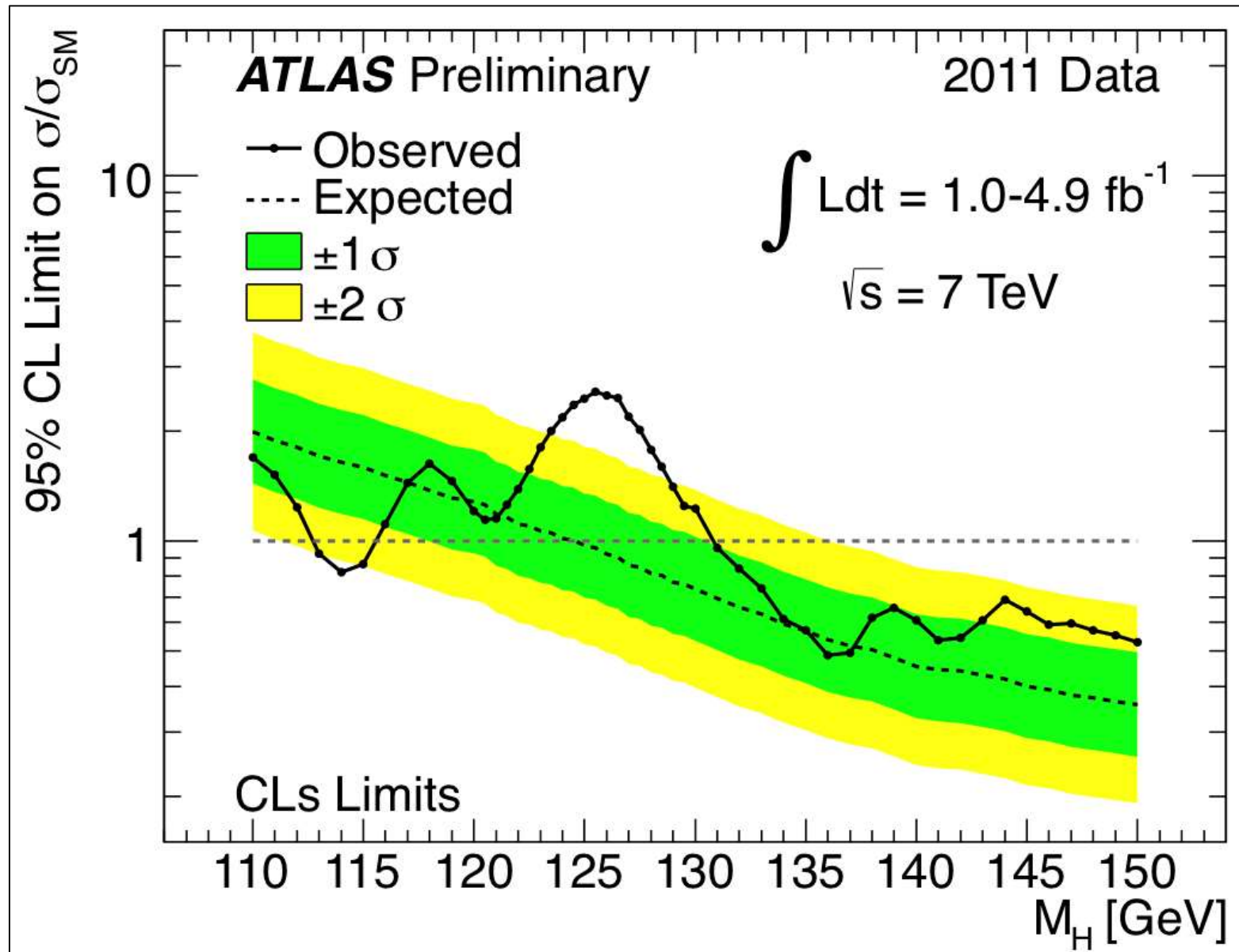
“If the galaxy flux density really *is* constant, and we repeatedly obtained sets of 15 measurements under the same conditions, then only 2% of the  $\chi^2$  values derived from these sets would be expected to be greater than our one actual measured value of 26.76”

*From Gregory, pg. 165*

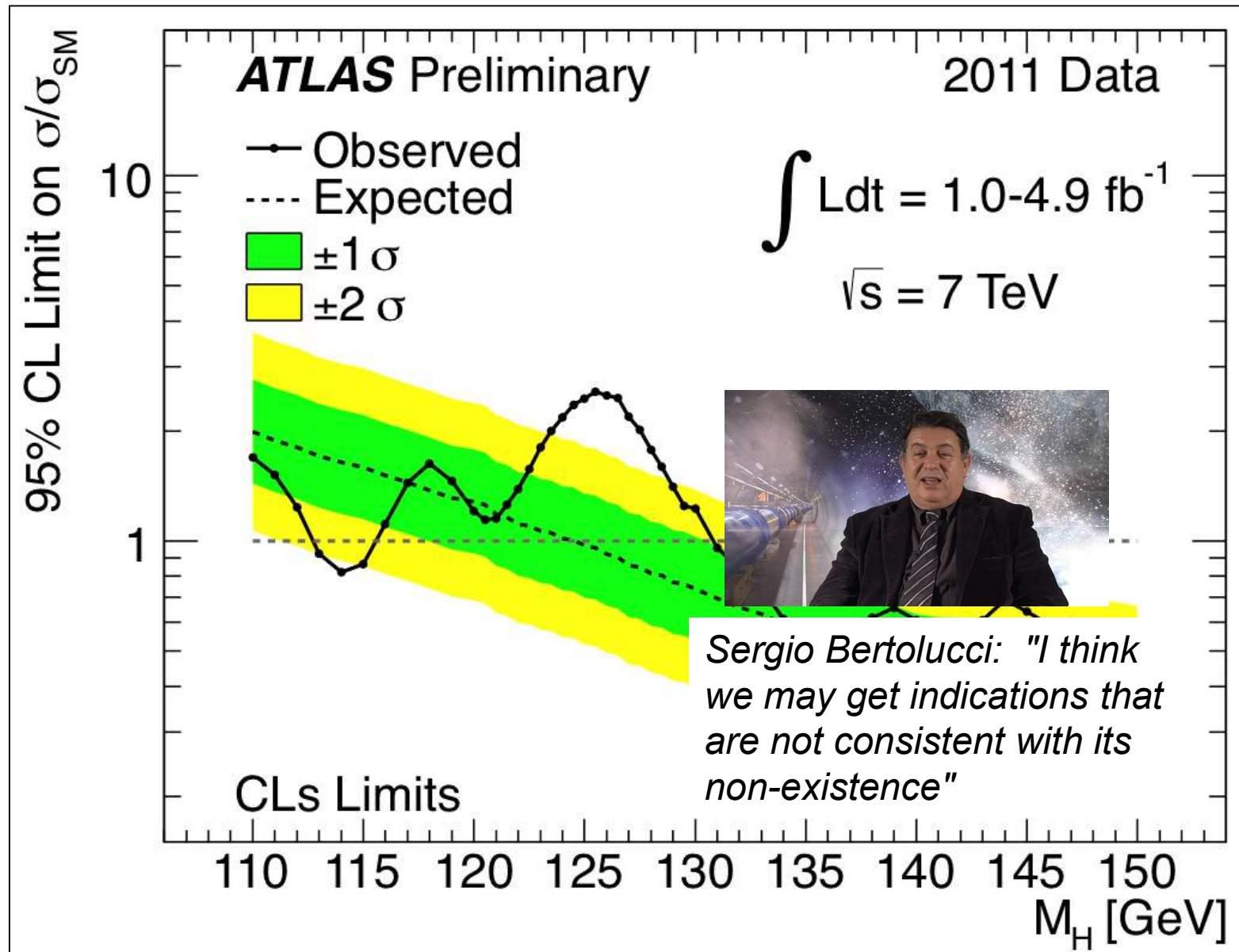
“At this point you may be asking yourself why we should care about a probability involving results never actually obtained”

*From Gregory, pg. 166*











Nevertheless, p-value based frequentist hypothesis testing remains very common in the literature:

| Type of problem                | test                                 | References     |
|--------------------------------|--------------------------------------|----------------|
| Line and curve goodness-of-fit | $\chi^2$ test                        | NR: 15.1-15.6  |
| Difference of means            | Student's $t$                        | NR: 14.2       |
| Ratio of variances             | $F$ test                             | NR: 14.2       |
| Sample CDF                     | K-S test<br>Rank sum tests           | NR: 14.3, 14.6 |
| Correlated variables?          | Sample correlation coefficient       | NR: 14.5, 14.6 |
| Discrete RVs                   | $\chi^2$ test /<br>contingency table | NR: 14.4       |



Nevertheless, p-value based frequentist hypothesis testing remains very common in the literature:

| Type of problem                | test                                 | References     |
|--------------------------------|--------------------------------------|----------------|
| Line and curve goodness-of-fit | $\chi^2$ test                        | NR: 15.1-15.6  |
| Difference of means            | Student's $t$                        | NR: 14.2       |
| Ratio of variances             | $F$ test                             | NR: 14.2       |
| Sample CDF                     | K-S test<br>Rank sum tests           | NR: 14.3, 14.6 |
| Correlated variables?          | Sample correlation coefficient       | NR: 14.5, 14.6 |
| Discrete RVs                   | $\chi^2$ test /<br>contingency table | NR: 14.4       |

*See also supplementary notes on my.SUPA and Moodle*



In the Bayesian approach, we can test our model, in the light of our data (e.g. rolling a die) and see how our knowledge of its parameters evolves, for any sample size, considering only the data that we *did* actually observe

$$\begin{array}{ccc} \text{Posterior} & & \text{Likelihood} & & \text{Prior} \\ \downarrow & & \downarrow & & \downarrow \\ p(\text{model} \mid \text{data}, I) & \propto & p(\text{data} \mid \text{model}, I) \times & p(\text{model} \mid I) \\ \text{What we know now} & & \text{Influence of our} & & \text{What we knew} \\ & & \text{observations} & & \text{before} \end{array}$$



What do we choose as our prior?

Good question!

Source of much argument between  
Bayesians and frequentists



Blood on the walls



What do we choose as our prior?

Good question!

Source of much argument between Bayesians and frequentists



Blood on the walls

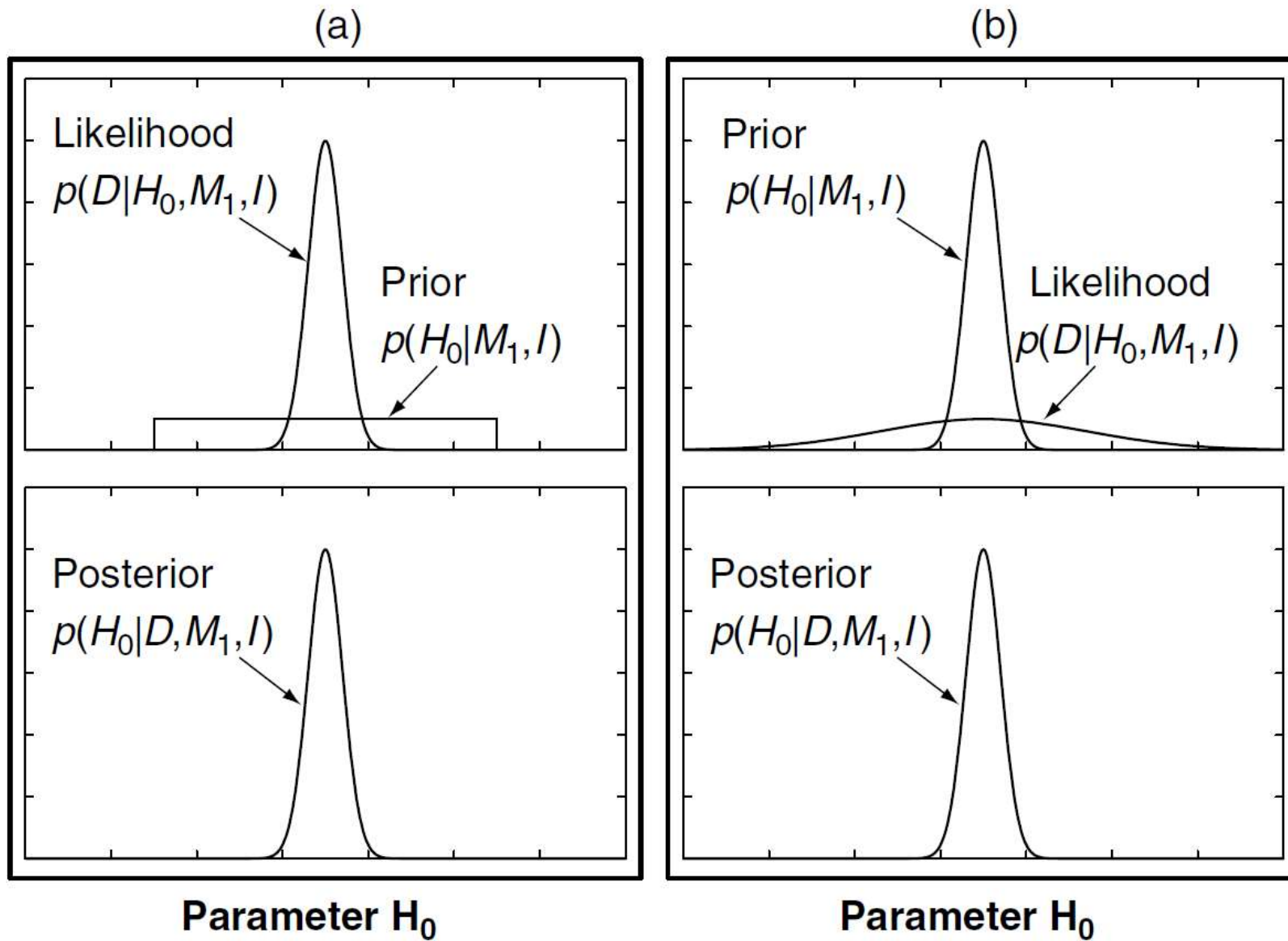
If our data are good enough, it shouldn't matter

$$\begin{array}{ccc} \text{Posterior} & & \text{Likelihood} \quad \text{Prior} \\ \downarrow & & \downarrow \quad \downarrow \\ p(\text{model} \mid \text{data}, I) & \propto & p(\text{data} \mid \text{model}, I) \times p(\text{model} \mid I) \end{array}$$

└──────────────────┘

*Dominates*

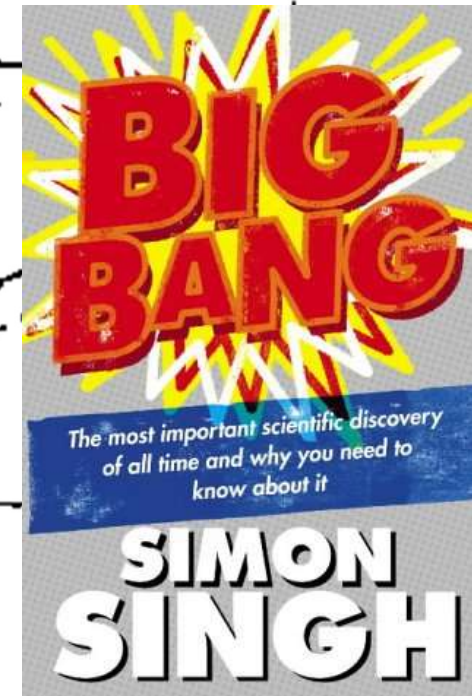
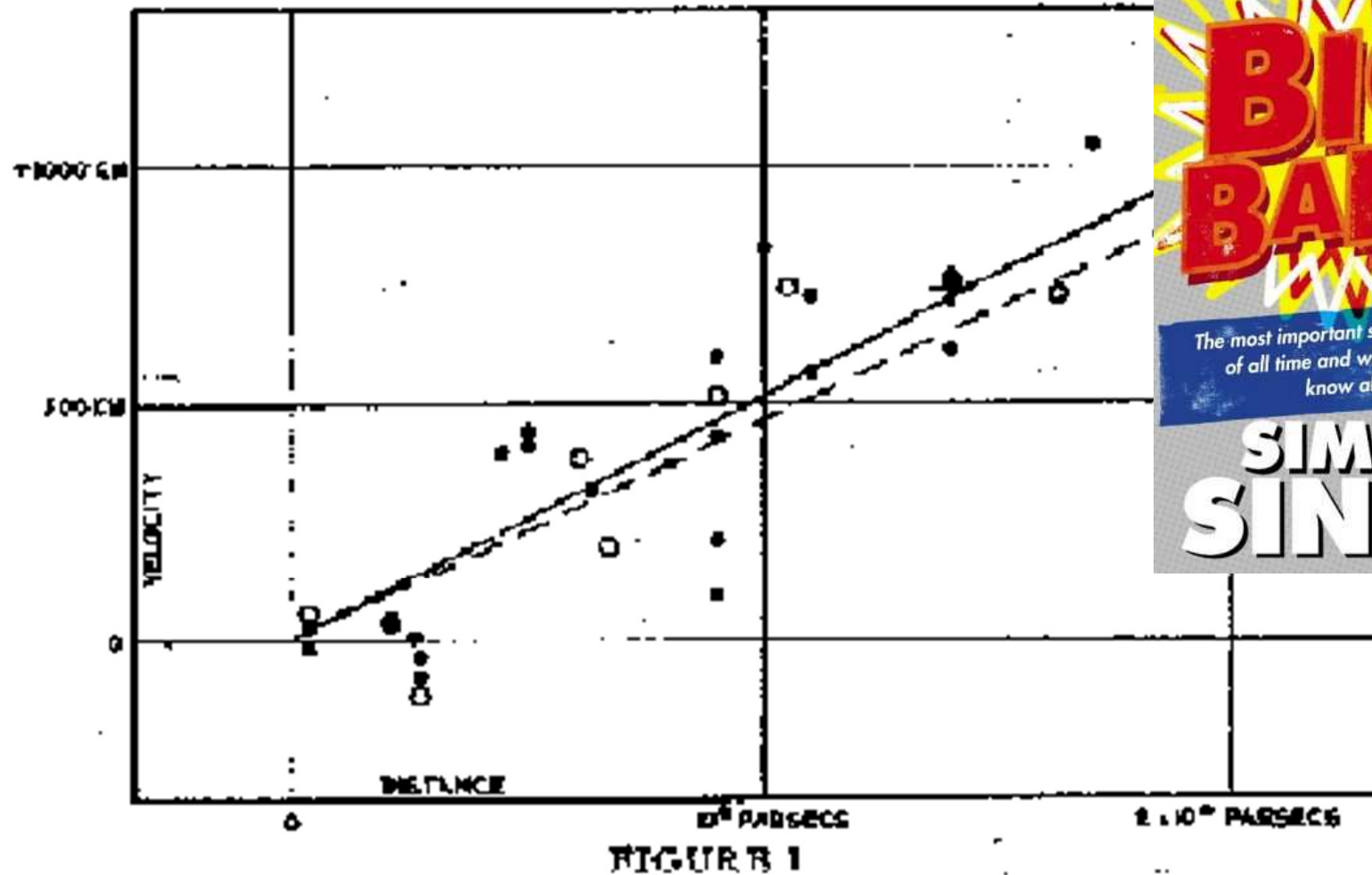




*From Gregory, pg 8.*



# Hubble's Law: 1929



Hubble parameter = expansion rate of the Universe  
= slope of Hubble's law