# Lecture 3:

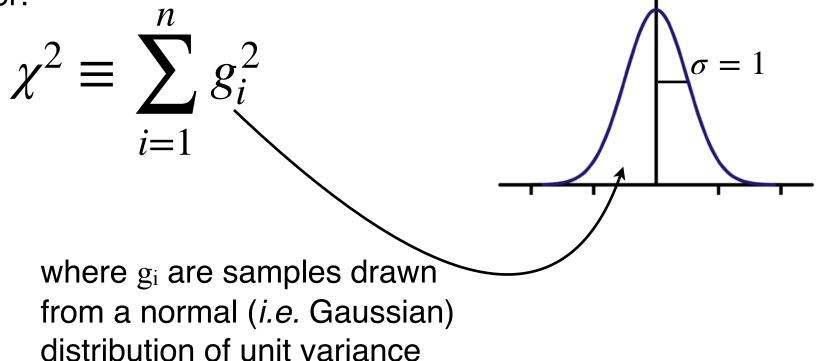
# **Testing Models**

- Chi-Squared
- "Scientific Method"
- Student's t
- Correlation test
- Non-parametric tests

# Testing Models



Consider:



Then the distribution of this quantity defines a  $\chi^2$  ("chi-squared") distribution with n <u>degrees of freedom</u>

effective number of independent samples contributing to the variance

The  $\chi^2$  probability density function for n degrees of freedom has the form:

$$P(\chi^2, n) = \frac{(\chi^2)^{\frac{n}{2} - 1} e^{-\frac{\chi^2}{2}}}{2^{\frac{n}{2}} \Gamma(\frac{n}{2})}$$

Where 
$$\Gamma(k) = (k-1)!$$

$$\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} dx$$
for any real value  $z$ 

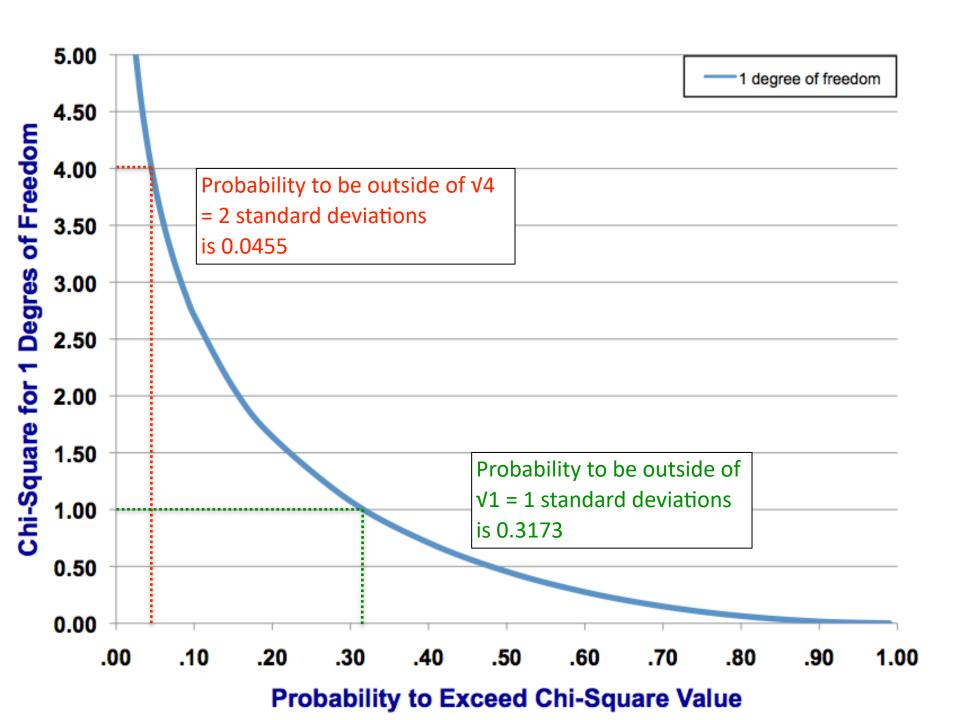
Note that:  $P(\chi^2,2) = \frac{1}{2}e^{-\frac{\chi^2}{2}}$  (will come back to this)

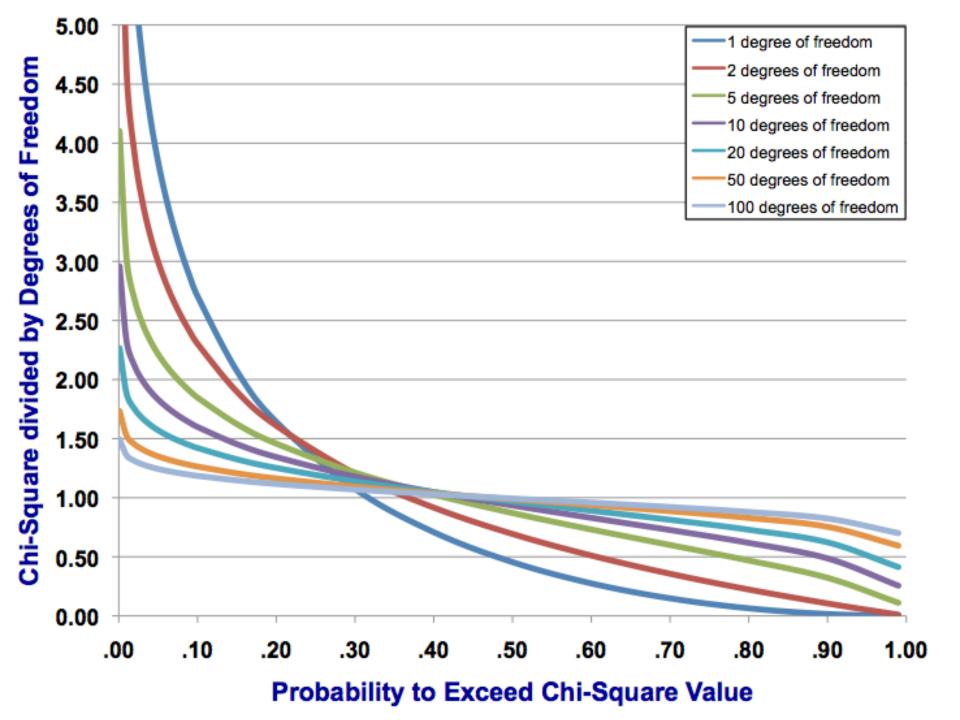
And the integral probability is given by:

if k = positive integer

$$P(>\chi^2, n) = 1 - \frac{\gamma\left(\frac{n}{2}, \frac{\chi^2}{2}\right)}{\Gamma(\frac{n}{2})}$$

where  $\gamma(z,\alpha) = \int_{0}^{\alpha} x^{z-1}e^{-x}dx$ 





# Pearson's $\chi^2$ Test

So, for example, if we have a model, m, involving k free parameters (determined by a fit to the data) that seeks to predict the values, x, of n data points, each with *normally distributed uncertainties*, we can construct the sum:

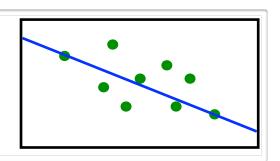
$$S \equiv \sum_{i=1}^{n} \frac{(x_i - m_i)^2}{\sigma_{m_i}^2}$$

normalises things to give a Gaussian distribution with unit variance for binned data (Poisson statistics)

$$\sigma_{m_i}^2 \cong m_i$$

S will then be distributed as a  $\chi^2$  distribution with (n-k) degrees of freedom, and can thus be used as a statistic to determine how well the model matches the data.

For example, imagine fitting a straight line (2 parameters: slope and intercept) to a set of data. You can always force the line to go through 2 of the data points exactly, so only n-2 of the data points will contribute to the variance around the model



S will then be distributed as a  $\chi^2$  distribution with n-k degrees of freedom, and can thus be used as a statistic to determine how well the model matches the data.

"If my model is correct, how often would a randomly drawn sample of data yield a value of  $\chi^2$  at least as large as this?"

Determining the best values for the model parameters by choosing them so as to minimise  $\chi^2$  is called the "Method of Least Squares."

Note that, if you vary one of the model parameters from its best fit value until  $\chi^2$  increases by 1, this therefore represents the change in the model parameter associated with 1 unit of variance in the fit quality (*i.e.* the "1 $\sigma$  uncertainty" in the model parameter).

#### Example 1:

Say we have n measurements of some quantity, with each measurement having a different Gaussian uncertainty. What is the best estimate for the mean value of this quantity?

$$\chi^2 = \sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma_i^2} \qquad \text{Let's find the value of } \mu \text{ that minimises this}$$

$$\frac{\partial \chi^2}{\partial \mu} = -2 \sum_{i=1}^n \frac{(x_i - \mu)}{\sigma_i^2} = -2 \sum_{i=1}^n \frac{x_i}{\sigma_i^2} + 2\mu \sum_{i=1}^n \frac{1}{\sigma_i^2} = 0$$

$$\mu_{best} = \frac{\sum_{i=1}^{n} \frac{x_i}{\sigma_i^2}}{\sum_{i=1}^{n} \frac{1}{\sigma_i^2}} = \frac{\sum_{i=1}^{n} w_i x_i}{\sum_{i=1}^{n} w_i} \quad \text{where } w_i \equiv \frac{1}{\sigma_i^2}$$

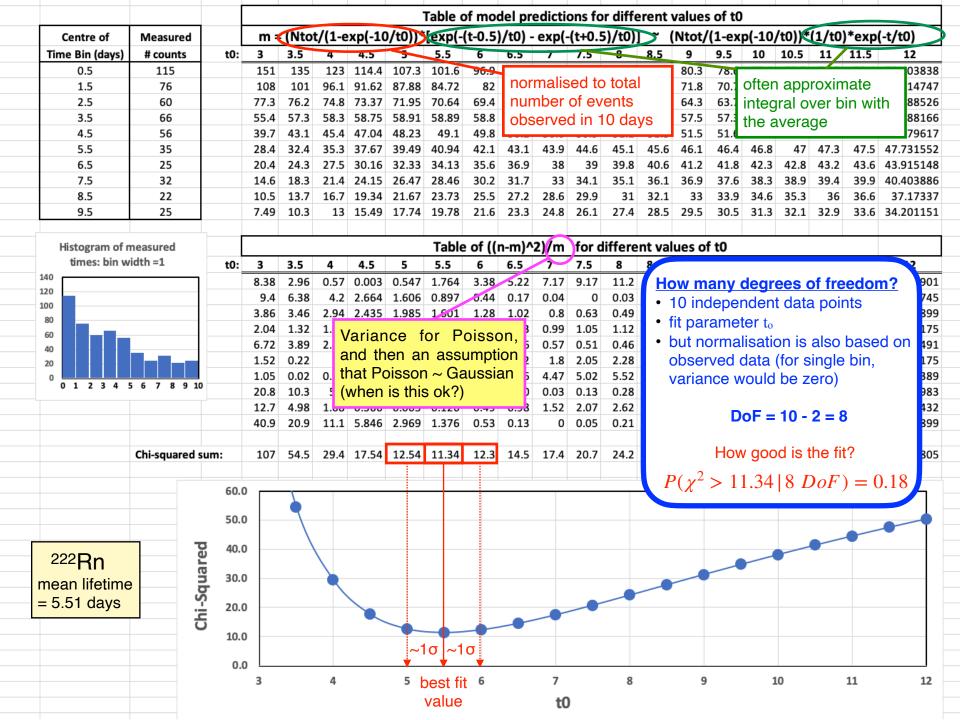
### Example 2:

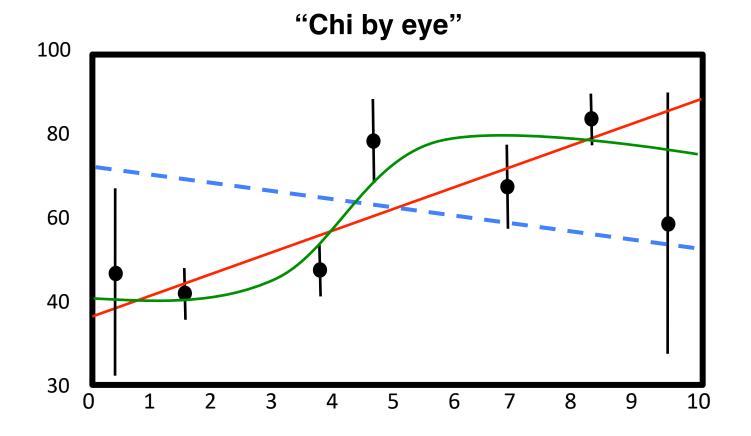
A newly commissioned underground neutrino detector sees a rate of internal radioactive contamination decreasing as a function of time. Measurements of the number of such events observed are taken on 10 consecutive days. Determine the best fit mean decay time in order to determine the source of the contamination.

decay probability:

$$P(t) = \frac{1}{t_o} e^{-\frac{t}{t_o}}$$

t<sub>o</sub> = mean decay lifetime

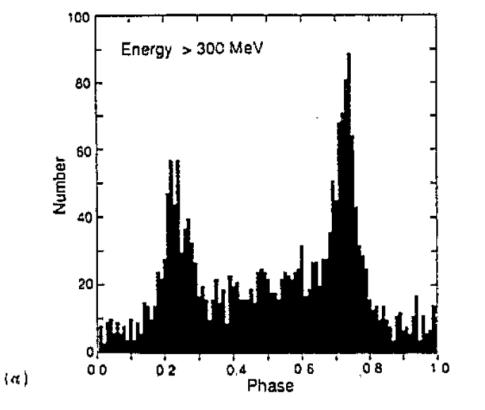




$$\chi^2 \sim (0.3)^2 + (0.1)^2 + (1.2)^2 + (1.7)^2 + (0.3)^2 + (0.8)^2 + (0.8)^2 = 5.8$$

Degrees of Freedom = 7 - 2 = 5

NOTE: This doesn't tell you which model is correct, but it can tell you which models don't fit well!



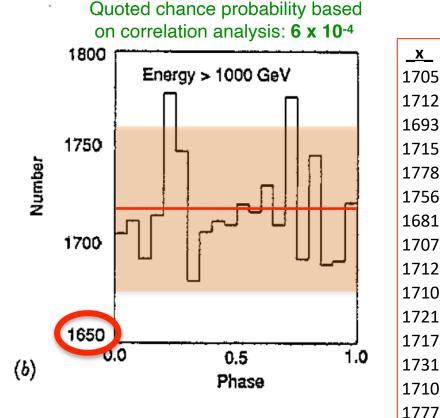


Figure 1. (a) Light curve for Geminga obtained with EGRET (b) The VHE  $\gamma$ -ray light curve of Geminga plotted at the GeV  $\gamma$ -ray phase, as derived from the COS-B ephemeris.

$$\chi^2 = \sum_{i=1}^{20} \frac{(x_i - 1718.45)^2}{1718.45} = 8.22$$

1693

17471690

1692

1722

**Wuant 'em Effect** 

DoF = 20 - 1 = 19 (98.4% chance of getting something larger)

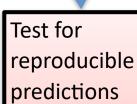
## Scientific Method:







Simplest and most predictive



to disprove

Rejected with high confidence

Not rejected with high confidence



Model

Next simplest & most predictive

A theory is judged not on what it can explain, but on what it can reproducibly predict!

We don't prove models correct; we reject those models that are wrong!

Test for reproducible

Rejected with high confidence



Don't state that data are "consistent" with a given model, but rather that they are "not inconsistent."



# More ways to test models...



## Student's t

Often misinterpreted as referring to being from or for "a student," rather than the fact that the name of the author happens to be "Student"

Except this was actually a pseudonym used by William Sealy Gosset in his 1908 paper, who was couching himself as "a student"!

where

Recall that the rms deviation in the estimated mean from a set of n samples is given by :

$$\sigma_m = rac{\sigma}{\sqrt{n}}^{rms ext{ of the full distribution.}}$$

But what if we don't know  $\sigma$  *a priori* and all we have are the sampled estimators?

$$t \equiv \frac{\bar{x} - \mu}{\left(s/\sqrt{n}\right)}$$

Want to find the distribution of t

$$\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

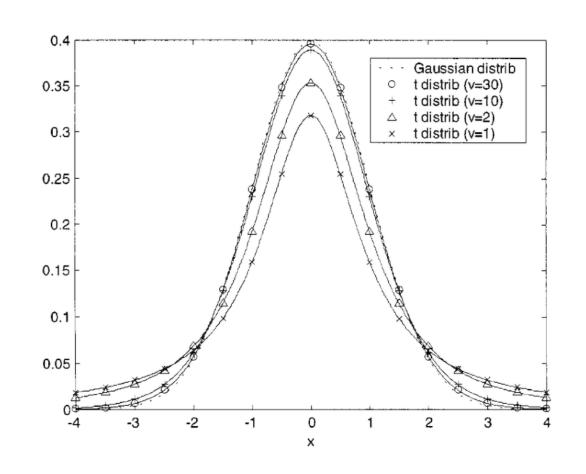
$$s^{2} = \frac{1}{(n-1)} \sum_{i=1}^{n} (x_{i} - \bar{x})^{2}$$

$$t \equiv \frac{\bar{x} - \mu}{\left(s/\sqrt{n}\right)}$$

$$f(t) = rac{\Gamma(rac{
u+1}{2})}{\sqrt{
u\pi}\,\Gamma(rac{
u}{2})}igg(1+rac{t^2}{
u}igg)^{-(
u+1)/2}$$

v = # degrees of freedom

As you would expect, this approaches the shape of a Gaussian distribution as the sample size grows:



## Pearson Correlation Coefficient

A test of linear correlation between two sets of data

$$r_{xy} = \frac{1}{n-1} \sum_{i=1}^{n} \left( \frac{x_i - \bar{x}}{s_x} \right) \left( \frac{y_i - \bar{y}}{s_y} \right)$$

$$=rac{\sum_{i=1}^{n}(x_{i}-ar{x})(y_{i}-ar{y})}{\sqrt{\sum_{i=1}^{n}(x_{i}-ar{x})^{2}}\sqrt{\sum_{i=1}^{n}(y_{i}-ar{y})^{2}}} \;\; =rac{\sum_{i}x_{i}y_{i}-nar{x}ar{y}}{\sqrt{\sum_{i}x_{i}^{2}-nar{x}^{2}}\,\sqrt{\sum_{i}y_{i}^{2}-nar{y}^{2}}}$$

This is just the covariance normalised by the sample rms deviations.

The value of this quantity runs from 1 (completely correlated) to -1 (completely anti-correlated), with zero indicating no correlation.

The statistics provides a *relative* measure of linear correlation but, in general, the probability distribution for r will depend on the distributions of x and y.

**IF** x and y are uncorrelated and each drawn from a normal distribution (such that, jointly, they can be described by a 2-D Gaussian), then:

$$\sigma_r = \sqrt{\frac{1 - r^2}{n - 2}}$$
DoF for 2 free parameters in linear fit

From which it is possible to define a t statistic for r:

$$t_r = r\sqrt{\frac{n-2}{1-r^2}}$$

# Spearman Rank-Order Correlation Coefficient

A non-parametric test of correlation between two sets of data (i.e. linearity is not assumed)

Define  $R_i$  as the 'rank' of  $x_i$  (i.e. the numerical position in an ordered list of the n data points from lowest to highest x value).

Define  $S_i$  as the 'rank' of  $y_i$  (i.e. the numerical position in an ordered list of the n data points from lowest to highest y value).

Note: it's possible to have identical ranks if a data set contains multiple identical values! In which case you should ascribe to each of these an 'average' rank value

Then define the rank coefficient as:

$$r = \frac{1}{n-1} \sum_{i=1}^n \left( \frac{R_i - \overline{R}}{s_R} \right) \left( \frac{S_i - \overline{S}}{s_S} \right) = \frac{\sum_i (R_i - \overline{R})(S_i - \overline{S})}{\sqrt{\sum_i (R_i - \overline{R})^2} \sqrt{\sum_i (S_i - \overline{S})^2}}$$

Similarly, the probability distribution can be approximated by the t statistic:

$$t_r = r\sqrt{\frac{n-2}{1-r^2}}$$

Generally pretty good and no longer depends on the actual distributions of x & y

# Kolmogorov-Smirnov (and the like)

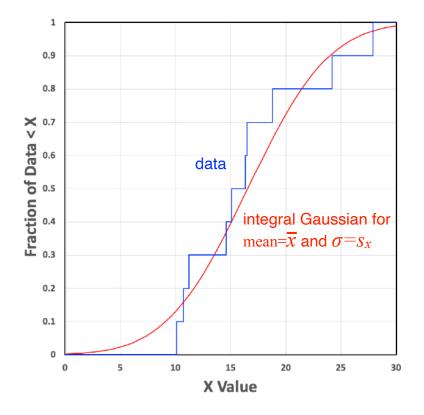
#### A non-parametric test of distributions

Plot the cumulative fraction of events less than or equal to a particular value of x as a function of x, along with the cumulative distribution for some model:

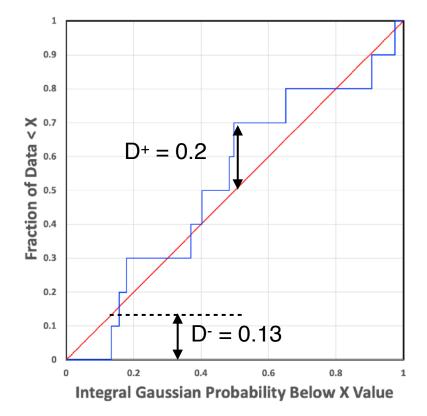
data point	x value	fraction ≤ x
1	10.1	0.1
2	10.7	0.2
3	11.2	0.3
4	14.6	0.4
5	15.1	0.5
6	16.3	0.6
7	16.5	0.7
8	18.8	0.8
9	24.2	0.9
10	27.9	1.0

$$\bar{x} = 16.54$$
$$s_x = 5.8$$

For example: Is this data Normally distributed?



Equivalently:



A more clearly defined span on the x-axis with a visually simple model expectation that is independent of the test distribution

There are several statistics that can be used to assess the level of agreement:

#### K-S statistics (Clustering)

 $D^+$  = maximum positive deviation from the model line  $D^-$  = maximum negative deviation from the model line  $D = max(D^+, D^-)$ 

$$V = D^+ + D^-$$
 (Kuiper test)

#### **Cramer-von Mises** (Variance)

$$W^{2} = \sum_{i=1}^{n} \left( y_{i} - \frac{2i - 1}{2n} \right)^{2} + \frac{1}{12n}$$

$$U^{2} = W^{2} - n \left( \bar{y} - \frac{1}{2} \right)^{2}$$

In general, the probability distributions for these statistics need to be determined by Monte Carlo calculations. However, for continuous variables tested against a well-defined model distribution under the null hypothesis, tables and approximate parameterisations exist to obtain p-values:

Test Statistic (T)	Modified Test Statistic (T*)	"High Tail" Approximate Parameterisation for P(T* > z)
$D^+ \ D^- \ D \ V \ W^2 \ U^2$	$D^{+}(\sqrt{n}+0.12+0.11/\sqrt{n})$ $D^{-}(\sqrt{n}+0.12+0.11/\sqrt{n})$ $D(\sqrt{n}+0.12+0.11/\sqrt{n})$ $V(\sqrt{n}+0.155+0.24/\sqrt{n})$ $(W^{2}-0.4/n+0.6/n^{2}) (1.0+0.1)$ $(U^{2}-0.1/n+0.1/n^{2}) (1.0+0.1)$	

