



#### Max-Planck-Institut für Physik (Werner-Heisenberg-Institut)

# How to tell good from bad

<u>Frederik Beaujean</u> http://arxiv.org/abs/1005.3233v2 http://arxiv.org/abs/1011.1674

Particle Physics Colloquium

Munich, 10.12.2010

# Outline



- Why goodness of fit?
- Statistical approach p-values
- Common statistics common pitfalls
- [A bit of number theory: runs and integer partitions]



+ Gaussian

#### Suppose:

- *N* measurements  $y_i(x_i)$  with uncertainty
- Standard Model (SM) background is quadratic
- New physics (NP) predicts signal peak (more than one NP model)



### Example problem





Fit function

$$f(x|\vec{\lambda}) = A + Bx + Cx^2 + \frac{D}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{\sigma^2}\right)$$

$$I : \vec{\lambda} = (A, B, C)$$
  

$$II : \vec{\lambda} = (A, D, \mu, \sigma)$$
  

$$III : \vec{\lambda} = (A, B, D, \mu, \sigma)$$
  

$$IV : \vec{\lambda} = (A, B, C, D, \mu, \sigma)$$

# Restriction



- Don't compare models directly here. Assume a model is true, then check consistency
- Often no alternative model known
- For model comparison (which one favored by data?) recommend Bayesian approach

#### Goodness of Fit: standard approach

#### **Requirement:**

• Assume a model *M* with parameters  $\lambda$ 

#### Test statistic:

- Any scalar function of data, T(D)
- Interpret: large T(D) = poor model

#### **Example:**

Familiar choice

• Prob. density of the data  $P(D|\vec{\lambda}) \propto \prod e^{i \lambda}$ 

$$P(D|\vec{\lambda}) \propto \prod \exp\left\{-\frac{\left(y_i - f(x_i|\vec{\lambda})\right)^2}{2\sigma_i^2}\right\} = \exp\left\{\frac{-\chi^2}{2}\right\}$$
$$T(D) \equiv \chi^2(D)$$

(

• Extension: discrepancy variable  $T(D|\lambda)$ . Fitting procedure important!



#### p-value





Warning: p-value not the P. that the model is true

#### Reasoning behind p-values





- Need prior knowledge about alternatives
- vague interpretation

#### Don't take p-value too seriously!

р

0.20





#### Frequent misunderstandings

There are several common misunderstandings about p-values.[3][4]

The p-value is not the probability that the null hypothesis is true.

In fact, frequentist statistics does not, and cannot, attach probabilities to hypotheses. Comparison of Bayesian and classical approaches shows that a p-value can be very close to zero while the posterior probability of the null is very close to unity (if there is no alternative hypothesis with a large enough *a priori* probability and which would explain the results more easily). This is the Jeffreys–Lindley paradox.

The p-value is not the probability that a finding is "merely a fluke."

As the calculation of a p-value is based on the *assumption* that a finding is the product of chance alone, it patently cannot also be used to gauge the probability of that assumption being true. This is different from the real meaning which is that the p-value is the chance of obtaining such results if the null hypothesis is true.

The p-value is *not* the probability of falsely rejecting the null hypothesis. This error is a version of the so-called prosecutor's fallacy. The p-value is *not* the probability that a replicating experiment would not yield the same conclusion.

1 – (p-value) is *not* the probability of the alternative hypothesis being true (see (1)).

The significance level of the test is not determined by the p-value.

The significance level of a test is a value that should be decided upon by the agent interpreting the data before the data are viewed, and is compared against the p-value or any other statistic calculated after the test has been performed. (However, reporting a p-value is more useful than simply saying that the results were or were not significant at a given level, and allows the reader to decide for himself whether to consider the results significant.)

The p-value does not indicate the size or importance of the observed effect (compare with effect size). The two do vary together however

- the larger the effect, the smaller the p-value will be, other things being equal.

#### From wikipedia

# Choice of statistic



Surprise: T arbitrary

Criteria:

- distribution known
- Easy to compute
- Power to reject alternatives
- Now consider statistics for Gauss, Poisson



- Goal: calculate p-value distribution for common statistics
- 10000 experiments
- Sample *N* data points from Model IV with fixed parameters
- Fit all models with (Markov chain+MINUIT) or MINUIT alone
- Plot the distribution of the p-value for the statistics chosen

# Test Statistics: Poisson

Neyman



Pearson

$$\chi_P^2 = \sum_i \frac{\left(n_i - \nu_i\right)^2}{\nu_i}$$

$$\chi_N^2 = \sum_i \frac{\left(n_i - \nu_i\right)^2}{n_i}$$

 $n_i$  observed events

 $u_i = 
u_i(ec{\lambda}, M) \,\, ext{expected events}$ 

• Uncertainty if  $n_i = 0$ ? Ignore bin or set uncertainty =1

- Asymptotically (i.e. infinite data, in **each** bin:  $n_i >>1$ ) know distribution of  $\chi^2_P$  .

#### Results







Worrisome peak for Neyman
 in model III

# Approximations





- No parameters fit, just plug in values used to generate the data sets
- Expect flat p-value
- Pearson good approximation
- Neyman bad, gets worse for smaller event numbers

Use Pearson. No uncertainties without model

# Gaussian linear regression





Predictions depend on parameters:

 $f(x_i|\vec{\lambda})$  linear in  $\vec{\lambda} \Rightarrow \nabla \chi^2 = 0$  linear in  $z_i \Rightarrow (N-n)$  DoF

Example: 
$$f(x|\vec{\lambda}) = A + Bx + Cx^2 + \frac{D}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu)^2}{\sigma^2}\right)$$

In real life, usually don't know exact form of  $P(\chi^2 | \vec{\lambda}^*, N, n)$ 







Likelihood of model IV for particular data set and *small* range

• Physics motivates *small* parameter range: e.g. C>0,  $\sigma$ >0.2 ...

- Global minimum may be elsewhere. Compare with *large* range
- Gradient based optimization (MINUIT/MIGRAD): need good starting point

• Clever user guess (difficult) or Markov chain (preferred). [mpp.mpg.de/bat/]

Fitting procedure not trivial!

# Results: p-value distribution for $\chi^2$





Fitting procedure and parameter ranges affect distribution of p-value

# Local vs Global Minimum





Use  $\chi^2$ -distribution with (N-n) DoF

True model, global minimum, but distribution not flat. → Nonlinear problem

#### Constraining parameter range = prior belief

# Conclusions



- P-values useful, but need understanding
- Fitting can make big difference
- Choice of statistic crucial
- Beware: distributions usually approximate





# FINIS

# Backup





- Most statistics disrespect order of data, information wasted
- Human brain good for simple problems

#### **Example:**

- N=25 datapoints
- Each Gaussian with mean = 0 and variance = 1



# Can we combine information about order and magnitude of deviation?

# Runs statistic





# **Runs distribution**



#### **Gaussian case:**

- Distribution of *T* exactly calculated for any *N* (nonparametric)
- Requires sum over integer partitions

N = 25

15

20



0.15

0.10

0.05

0.00

0

5

10

 $T_{\rm obs}$ 

 $PDF(T_{obs})$ 



• Number of ways to write integer *n* as sum of [*k*] integers

• Example 
$$5 = 5$$
  
 $= 4+1$   
 $= 3+2$   
 $= 3+1+1$   
 $= 2+2+1 \Rightarrow Part (5,3) = 2$   
 $= 2+1+1+1$   
 $= 1+1+1+1+1$ 

 $\Rightarrow \operatorname{Part}(5) = 7$ 

- Investigated by famous people (Leibniz, Euler ...)
- Related to strings, solid state, group theory ...

#### Partition triangle





Proposition:  $Part(n) - 1 = \sum_{r=1}^{N} \sum_{M=1}^{\min(r, N-r+1)} Part(n, k)$ 

# **Runs distribution**





- Good model:
- a) fitting bias towards p=1
- b) Success and failure similar
- Bad model:
- a) Success and failure different
- b) Bias towards p=0
- c) Missed a peak: failures OK



#### Model selection:

- Need explicit alternatives M<sub>1</sub>, M<sub>2</sub>
- Posterior odds

$$\frac{P(\boldsymbol{M}_1|\boldsymbol{D})}{P(\boldsymbol{M}_2|\boldsymbol{D})} = \frac{P(\boldsymbol{M}_1)}{P(\boldsymbol{M}_2)} \times \frac{P(\boldsymbol{D}|\boldsymbol{M}_1)}{P(\boldsymbol{D}|\boldsymbol{M}_2)}$$

#### **Bayes factor:**

• (very) sensitive to parameter range

$$P(D|M_1) = \int p(D|\vec{\lambda}) p_0(\vec{\lambda}) d\vec{\lambda}$$

• Occam's razor built in

#### **Example:**

• Six (NP) vs three (SM) parameters

 $\frac{P(\mathsf{SM}|D)}{P(\mathsf{NP}|D)} = \frac{P(\mathsf{SM})}{P(\mathsf{NP})} \times 61.7$