

Statistics for Particle Physics

Lecture 2: Estimation and Errors


Roger Barlow
The University of Huddersfield

CERN European School on High Energy Physics
St Petersburg

15th September 2019

The 2019 European School of High-Energy Physics

St. Petersburg, Russia, 4-17 September 2019



Lecture 2: Estimation and Errors

1 Estimation

- Bias
- Efficiency
- Maximum Likelihood Estimation
- Least squares
- Straight line fits
- Fitting Histograms

2 Errors

- Errors from Likelihood
- Asymmetric Errors
- Systematic Errors

Estimation

Statistician-speak for 'Measurement'

The general problem

You know the probability (density) function $P(x; a)$

Take data $\{x_j\}$. What is the best value for a ?

x_j may be single values, or pairs, or higher-dimensional

a may be a single parameter or several. If more than one, sometimes split into 'parameters of interest' and 'nuisance parameters'

Occasionally estimate a property (e.g. the mean) rather than a parameter

Very Broad definition

An Estimator $\hat{a}(x_1 \dots x_N)$ is a function of the data that gives a value for the parameter a

A good estimator...

There is no 'correct' estimator - but some are better than others

A perfect estimator would be

Consistent $\hat{a}(x_1 \dots x_N) \rightarrow a$ as $N \rightarrow \infty$

Unbiased $\langle \hat{a} \rangle = a$

Efficient $\langle (\hat{a} - a)^2 \rangle$ is as small as possible

Invariant $\hat{f}(a) = f(\hat{a})$

No estimator is perfect - the goals are incompatible.

Examples of Bias

An unbiased estimator of the mean

Suppose we take $\hat{\mu} = \bar{x}$

$$\langle \hat{\mu} \rangle = \left\langle \frac{1}{N} \sum x_i \right\rangle = \frac{1}{N} \sum \langle x \rangle = \frac{1}{N} \sum \mu = \mu$$

A biased estimator of the Variance

Suppose we take $\hat{V} = \overline{x^2} - \bar{x}^2$

$$\text{So } \langle \hat{V} \rangle = \langle \overline{x^2} \rangle - \langle \bar{x}^2 \rangle$$

First term is just $\langle x^2 \rangle$. To make sense of second term, note $\langle x \rangle = \langle \bar{x} \rangle$ and add and subtract $\langle x \rangle^2$

$$\langle \hat{V} \rangle = \langle x^2 \rangle - \langle x \rangle^2 - (\langle \bar{x}^2 \rangle - \langle \bar{x} \rangle^2)$$

$$\langle \hat{V} \rangle = V(x) - V(\bar{x}) = V - \frac{V}{N} = \frac{N-1}{N} V$$

Estimator is biased! \hat{V} will, on average, give too small a value

Correct for the bias using $\hat{V} = \frac{N}{N-1}(\overline{x^2} - \bar{x}^2)$ and/or $\hat{\sigma} = \sqrt{\frac{\sum_i (x_i - \bar{x})^2}{N-1}}$

The Minimum Variance Bound

also known as the Cramer-Rao bound

Likelihood (again)

$$L(a; x_1, x_2, \dots, x_N) = P(x_1; a)P(x_2; a)\dots P(x_N; a)$$

Probability for the whole data sample, for a particular value of a

Will write $L(a; x_1, x_2, \dots, x_N)$ as $L(a; x)$ for simplicity

The Minimum Variance Bound (MVB)

For any unbiased estimator $\hat{a}(x)$, the variance is bounded

$$V(\hat{a}) \geq -\frac{1}{\left\langle \frac{d^2 \ln L}{da^2} \right\rangle} = \frac{1}{\left\langle \left(\frac{d \ln L}{da} \right)^2 \right\rangle} \quad (1)$$

Proof of the MVB

Proof.

Unitarity requires $\int P(x; a) dx = \int L(a; x) dx = 1$

Differentiate wrt a :

$$0 = \int \frac{dL}{da} dx = \int L \frac{d \ln L}{da} dx = \left\langle \frac{d \ln L}{da} \right\rangle \quad (2)$$

If \hat{a} is unbiased $\langle \hat{a} \rangle = \int \hat{a}(x) P(x; a) dx = \int \hat{a}(x) L(a; x) dx = a$

Differentiate wrt a : $1 = \int \hat{a}(x) \frac{dL}{da} dx = \int \hat{a} L \frac{d \ln L}{da} dx$

Subtract Eq 2 multiplied by a , and get $\int (\hat{a} - a) \frac{d \ln L}{da} L dx = 1$

Invoke the Schwarz Inequality $\int u^2 dx \int v^2 dx \geq (\int uv dx)^2$ with

$u \equiv (\hat{a} - a) \sqrt{L}$, $v \equiv \frac{d \ln L}{da} \sqrt{L}$

Hence $\int (\hat{a} - a)^2 L dx \int \left(\frac{d \ln L}{da} \right)^2 L dx \geq 1$

$$\langle (\hat{a} - a)^2 \rangle \geq 1 / \left\langle \left(\frac{d \ln L}{da} \right)^2 \right\rangle \quad (3)$$

Lemma

Differentiating Equation 2 again gives

$$\frac{d}{da} \int L \frac{d \ln L}{da} dx = \int \frac{dL}{da} \frac{d \ln L}{da} dx + \int L \frac{d^2 \ln L}{da^2} dx = \left\langle \left(\frac{d \ln L}{da} \right)^2 \right\rangle + \left\langle \frac{d^2 \ln L}{da^2} \right\rangle = 0$$

$$\text{Hence } \left\langle \left(\frac{d \ln L}{da} \right)^2 \right\rangle = - \left\langle \frac{d^2 \ln L}{da^2} \right\rangle$$

This is called the **Fisher Information**. Note how it is intrinsically positive.

Maximum Likelihood Estimation

Maximise the likelihood (actually the log likelihood)

$$\text{Maximise } \ln L = \sum_i \ln P(x_i; a) \quad (4)$$

$$\left. \frac{d \ln L}{da} \right|_{\hat{a}} = 0 \quad (5)$$

Is consistent.

Is biased, but bias falls like $1/N$

Is efficient for large N

Is invariant - doesn't matter if you reparametrise a

Particular problem may be solved in 3 ways depending on complexity

- 1 Solve Equation 5 algebraically
- 2 Solve Equation 5 numerically
- 3 Solve Equation 4 numerically

Least Squares Estimation

Gaussian measurements of y taken at various x values, with measurement error σ , and a prediction $y = f(x; a)$

$$P(y; x, a) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(y-f(x,a))^2/2\sigma^2}$$

$$\ln L = - \sum \frac{(y_i - f(x_i; a))^2}{2\sigma_i^2} + \text{constants}$$

To maximise $\ln L$, minimise $\chi^2 = \sum \frac{(y_i - f(x_i; a))^2}{\sigma_i^2}$

Differentiating gives the **Normal Equations**: $\sum \frac{(y_i - f(x_i; a))}{\sigma_i^2} f'(x_i; a) = 0$

If $f(x; a)$ is linear in a then these can be solved exactly.

Otherwise use an iterative method.

The Straight Line Fit

Function $y = mx + c$

Assume all σ_i the same (extension to general case straightforward)

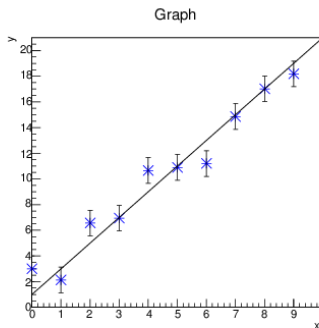
Normal Equations

$$\sum (y_i - mx_i - c)x_i = 0$$

$$\sum (y_i - mx_i - c) = 0$$

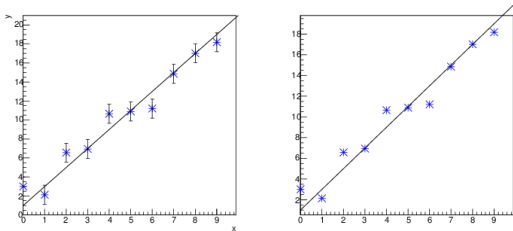
$$\text{Solution } m = \frac{\overline{xy} - \bar{x}\bar{y}}{\overline{x^2} - \bar{x}^2}$$

$$c = \bar{y} - m\bar{x}$$



Diversion: Regression

For most statisticians, 'Regression' = 'Straight Line fit'



History: Galton and father/son heights

Tall fathers tend to have tall sons - but not that tall. 'Regression towards mediocrity'

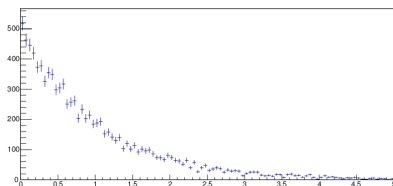
More accurate measurements would not decrease the spread

Ambiguity as to whether to plot x against y or y against x

Paradox: Tall sons tend to have tall fathers - but not that tall!

Fitting Histograms

Fitting a histogram - error given by Poisson statistics so $\sigma = \sqrt{N}$



4 methods - increasing accuracy, decreasing speed.

$$f_i(x_i; a) = P(x_i; a) \times \text{binwidth}$$

- 1 Minimise $\chi^2 = \sum_i \frac{(n_i - f_i)^2}{n_i}$. 'Neyman χ^2 '. . Breaks if $n_i = 0$
- 2 Minimise $\chi^2 = \sum_i \frac{(n_i - f_i)^2}{f_i}$. 'Pearson χ^2 '. **Only for histograms!**
- 3 Maximise $\ln L = \sum \ln(e^{-f_i} f_i^{n_i} / n_i!)$ $\sim \sum n_i \ln f_i - f_i$. "Binned ML"
- 4 Ignore bins and maximise likelihood. Sum runs over N_{events} not N_{bins} . Have to use for sparse data.

Errors

For large N , $\ln L$ curve is parabola

At the maximum,

$$\ln L(a) = \ln L(\hat{a}) + \frac{1}{2}(a - \hat{a})^2 \frac{d^2 \ln L}{da^2}$$

$$-1 / \left\langle \frac{d^2 \ln L}{da^2} \right\rangle \text{ gives } V(\hat{a})$$

(ML saturates MVB)

$$\text{Approximate } \left\langle \frac{d^2 \ln L}{da^2} \right\rangle \approx \left. \frac{d^2 \ln L}{da^2} \right|_{a=\hat{a}}$$

$$\sigma_{\hat{a}} = \sqrt{-\frac{1}{\frac{d^2 \ln L}{da^2}}}$$

When $a - \hat{a} = \pm \sigma_{\hat{a}}$,

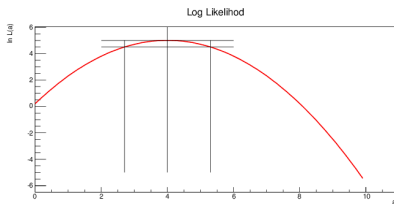
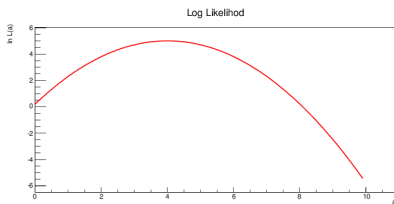
$$\ln L(a) = \ln L(\hat{a}) - \frac{1}{2}$$

Read off errors from $\Delta \ln L = -\frac{1}{2}$

See R.B. arXiv:physics/0403046 for small print

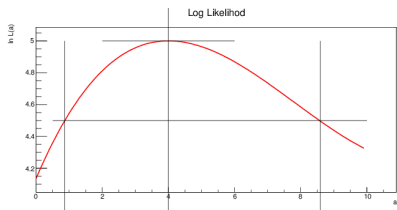
This gives σ , or 68% errors. Can also take $\Delta \ln L = -2$ to get $2\sigma=95\%$, etc.

If working with χ^2 , $L \propto e^{-\frac{1}{2}\chi^2}$ so take $\Delta\chi^2 = 1$



Asymmetric Errors

Typically arise in Poisson situations: say you see 1 event. $\lambda = 1.5$ is more likely to fluctuate down to 1 than $\lambda = 0.5$ fluctuate up to 1.



Read off σ_+ and σ_- from the two $\Delta \ln L = -\frac{1}{2}$ crossings

Avoid if possible

Combination of Asymmetric Errors

Given $x \pm \sigma_x, y \pm \sigma_y$, (and $\rho_{xy} = 0$) the error on $f = x + y$ is $\sigma_f^2 = \sigma_x^2 + \sigma_y^2$ (Sum in quadrature)

Given $x_{-\sigma_x^-}^{+\sigma_x^+}, y_{-\sigma_y^-}^{+\sigma_y^+}$, (and $\rho_{xy} = 0$), what is the error on $f = x + y$?

Standard Recipe

Sum in quadrature separately: $\sigma_f^{+2} = \sigma_x^{+2} + \sigma_y^{+2}$, $\sigma_f^{-2} = \sigma_x^{-2} + \sigma_y^{-2}$

This is **manifestly wrong** as it breaks the central limit theorem

Counterexample

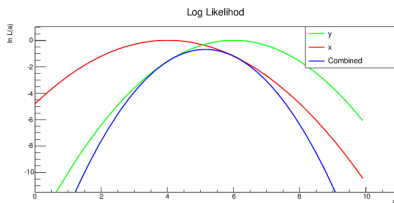
Add N i.i.d. variables with skew likelihood: $\sigma^+ = 2\sigma^-$.

Standard Recipe reduces both σ^+ and σ^- by factor $1/\sqrt{N}$ but still skew - and not Gaussian. Never will be.

Combining Asymmetric Measurements

Another approach

If you know the likelihood functions, you can do it
Here red and green curves are measurements of a
The log likelihood functions just add (blue)



But we don't know the full likelihood function: just 3 points (and that it had a maximum at the second)

Try various models (cubics, constrained quartic...) on likely instances
Two most plausible (for details see RB, arXiv:0406120)

$$\ln L = -\frac{1}{2} \left(\frac{x - \hat{x}}{\sigma_0 + \sigma'(x - \hat{x})} \right)^2 \quad \ln L = -\frac{1}{2} \frac{(x - \hat{x})^2}{V_0 + V'(x - \hat{x})}$$

Both pretty good. First does better with errors on log a , second does better with Poisson.

How to do it

For each measurement (x, σ^+, σ^-) find σ and σ' , or V and V'

$$\text{Given by } \sigma_0 = 2 \frac{\sigma^+ \sigma^-}{\sigma^+ + \sigma^-}, \quad \sigma' = \frac{\sigma^+ - \sigma^-}{\sigma^+ + \sigma^-}$$

$$\text{or } V_0 = \sigma^+ \sigma^-, \quad V' = \sigma^+ - \sigma^-$$

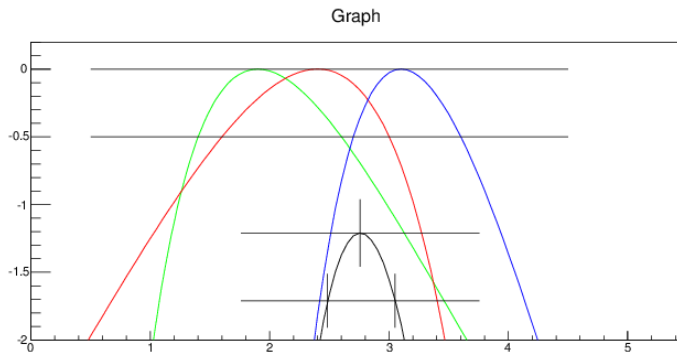
Find maximum of sum, numerically, and $\Delta \ln L = -\frac{1}{2}$ points

Programs in ROOT and R in

<https://doi.org/10.5281/zenodo.2576890>

Using the class

The result



Combining $1.9^{+0.7}_{-0.5}$, $2.4^{+0.6}_{-0.8}$ and $3.1^{+0.5}_{-0.4}$ gives $2.76^{+0.29}_{-0.27}$

Errors in 2 or more dimensions

For 2 (or more) dimensions,
define regions using contours
in $\Delta \ln L$ (or $\Delta \chi^2 \equiv -2\Delta \ln L$)

Levels change:

In 2D, cutting at 1σ square
would give $0.68^2 = 47\%$.

A 1σ contour gives 39%.

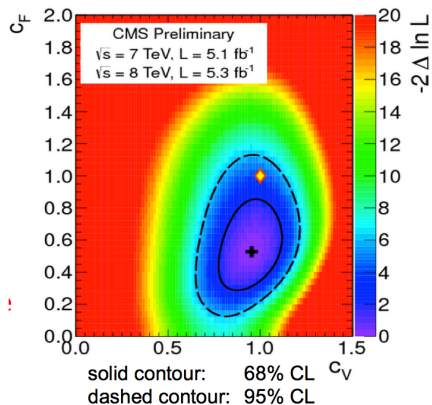
For 68% use $\Delta \chi^2 = 2.27$

$\Delta \ln L = -1.14$

For 95% use $\Delta \chi^2 = 5.99$

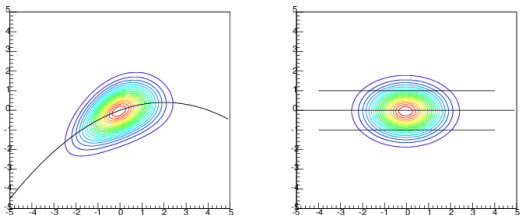
$\Delta \ln L = -3.00$

(Values from χ^2 distribution -
coming later)



Nuisance Parameters I

Profile Likelihood - motivation (not very rigorous)



You have a 2D likelihood plot with axes a_1 and a_2 . You are interested in a_1 but not in a_2 ('Nuisance parameter')

Different values of a_2 give different results (central and errors) for a_1

Suppose it is possible to transform to $a'_2(a_1, a_2)$ so L factorises, like the one on the right. $L(a_1, a'_2) = L_1(a_1)L_2(a'_2)$

Whatever the value of a'_2 , get same result for a_1

So can present this result for a_1 , independent of anything about a'_2 .

Path of central a'_2 value as fn of a_1 , is peak - path is same in both plots

So no need to factorise explicitly: plot $L(a_1, \hat{a}_2)$ as fn of a_1 and read off 1D values.

$\hat{a}_2(a_1)$ is the value of a_2 which maximises $\ln L$ for this a_1

Nuisance Parameters 2

Marginalised likelihoods

Instead of profiling, just integrate over a_2 .

Can be very helpful alternative, specially with many nuisance parameters

But be aware - this is strictly Bayesian

Frequentists are not allowed to integrate likelihoods wrt the parameter

$\int P(x; a) dx$ is fine, but $\int P(x; a) da$ is off limits

Reparametrising a_2 (or choosing a different prior) will give different values for a_1

Systematic Errors

Caution! This contains material some people may find offensive.

There is a lot of bad practice out there. Muddled thinking and following traditional procedures without understanding.

When statistical errors dominated, this didn't matter much. In the days of particle factories and big data samples, it does.

People are ignorant - ignorance leads to fear. They follow familiar rituals they hope will keep them safe.



- What is a Systematic Error?
- How to deal with them
- How to evaluate them
- Checking your analysis
- Conclusions and recommendations

What is a Systematic Error?

Systematic error: reproducible inaccuracy introduced by faulty equipment, calibration, or technique.

Systematic effects is a general category which includes effects such as background, scanning efficiency, energy resolution, variation of counter efficiency with beam position, and energy, dead time, etc. The uncertainty in the estimation of such a systematic effect is called a systematic error.

Bevington

Orear

What is a Systematic Error?

Systematic error: reproducible inaccuracy introduced by faulty equipment, calibration, or technique.

Systematic effects is a general category which includes effects such as background, scanning efficiency, energy resolution, variation of counter efficiency with beam position, and energy, dead time, etc. The uncertainty in the estimation of such a systematic effect is called a systematic error.

Bevington

Orear

These are contradictory

What is a Systematic Error?

Systematic error: reproducible inaccuracy introduced by faulty equipment, calibration, or technique.

Systematic effects is a general category which includes effects such as background, scanning efficiency, energy resolution, variation of counter efficiency with beam position, and energy, dead time, etc. The uncertainty in the estimation of such a systematic effect is called a systematic error.

Bevington

Orear

These are contradictory

Orear is **RIGHT**

What is a Systematic Error?

Systematic error: reproducible inaccuracy introduced by faulty equipment, calibration, or technique.

Systematic effects is a general category which includes effects such as background, scanning efficiency, energy resolution, variation of counter efficiency with beam position, and energy, dead time, etc. The uncertainty in the estimation of such a systematic effect is called a systematic error.

Bevington

Orear

These are contradictory

Orear is **RIGHT**

Bevington is **WRONG**

What is a Systematic Error?

Systematic error: reproducible inaccuracy introduced by faulty equipment, calibration, or technique.

Systematic effects is a general category which includes effects such as background, scanning efficiency, energy resolution, variation of counter efficiency with beam position, and energy, dead time, etc. The uncertainty in the estimation of such a systematic effect is called a systematic error.

Bevington

Orear

These are contradictory

Orear is **RIGHT**

Bevington is **WRONG**

So are a lot of other books and websites

An error is not a mistake

We teach undergraduates the difference between *measurement errors*, which are part of doing science, and *mistakes*.

If you measure a potential of 12.3 V as 12.4 V, with a voltmeter accurate to 0.1V, that is fine. Even if you measure 12.5 V

If you measure it as 124 V, that is a mistake.

Bevington is describing *Systematic mistakes*

Orear is describing *Systematic uncertainties* - which are 'errors' in the way we use the term.

Avoid using 'systematic error' and always use 'uncertainty' or 'mistake'?
Probably impossible. But should **always** know which you mean

Examples

Track momenta from $p_i = 0.3B\rho_i$ have statistical errors from ρ and systematic errors from B

Calorimeter energies from $E_i = \alpha D_i + \beta$ have statistical errors from light signal D_i and systematic errors from calibration α, β

Branching ratios from $Br = \frac{N_D - B}{\eta N_T}$ have statistical error from N_D and systematics from efficiency η , background B , total N_T

Bayesian or Frequentist?

Can be either

Frequentist: Errors determined by an *ancillary experiment* (real or simulated)

E.g. magnetic field measurements, calorimeter calibration in a testbeam, efficiency from Monte Carlo simulation

Sometimes the ancillary experiment is also the main experiment - e.g. background from sidebands.

Bayesian: theorist thinks the calculation is good to 5% (or whatever).
Experimentalist affirms calibration will not have shifted during the run by more than 2% (or whatever)

Some analysis techniques use hybrid of frequentist and Bayesian.

How to handle them: Correlation

Actually straightforward. Systematic uncertainties obey the same rules as statistical uncertainties

We write $x = 12.2 \pm 0.3 \pm 0.4$ but we could write $x = 12.2 \pm 0.5$.
For single measurement extra information is small.

For multiple measurements e.g. $x_a = 12.2 \pm 0.3$, $x_b = 17.1 \pm 0.4$, *all* ± 0.5 extra information important, as results correlated.

Example: cross sections with common luminosity error, branching ratios with common efficiency ...

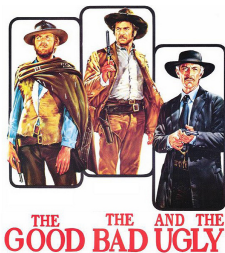
Taking more measurements and averaging does not reduce the error.

Consequence

No way to estimate σ_{sys} from the data - hence no check from χ^2 test etc
Not because systematic errors are unusually hostile - but because statistical errors are unusually friendly

Handling Systematic Errors in your analysis

3 types



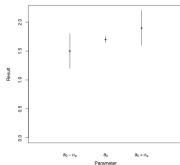
1) Uncertainty in an explicit continuous parameter:

E.g. uncertainty in efficiency, background and luminosity in branching ratio or cross section

Standard combination of errors formula and algebra, just like undergraduate labs. Have to include correlations but this is all handled by matrices.

Handling Systematic Errors (2)

Uncertainty in an implicit continuous parameter such as: MC tuning numbers (σ_{p_T} , polarisation.....)



Not amenable to algebra

Method: vary parameter by $\pm\sigma$ and look at what happens to your analysis result (directly, or through efficiency, background etc.)

Note 1: Hopefully effect is equal but opposite - if not then can introduce asymmetric error, but avoid if you can. Rewrite $^{+0.5}_{-0.3}$ as ± 0.4

Note 2. Your analysis results will have errors due to e.g. MC statistics. Some people add these (in quadrature). This is **wrong**. Technically correct thing to do is subtract them in quadrature, but this is not advised.

Handling Systematic Errors (3)

Discrete uncertainties, typically in model choice

Situation depends on status of model. Sometimes one preferred, sometimes all equal (more or less)

With 1 preferred model and one other, quote $R_1 \pm |R_1 - R_2|$

With 2 models of equal status, quote $\frac{R_1+R_2}{2} \pm \left| \frac{R_1-R_2}{\sqrt{2}} \right|$

N models: take $\bar{R} \pm \sqrt{\frac{N}{N-1}(\bar{R}^2 - \overline{R^2})}$ or similar mean value

2 extreme models: take $\frac{R_1+R_2}{2} \pm \frac{|R_1-R_2|}{\sqrt{12}}$

These are just ballpark estimates. Do not push them too hard. If the difference is not small, you have a problem - which can be an opportunity to study model differences.

Checking the analysis



“As we know, there are known knowns. There are things we know that we know. There are known unknowns. That is to say, there are things that we know we don’t know. But there are also unknown unknowns. There are things we don’t know we don’t know.”

Donald H Rumsfeld

Checking the analysis: Errors are not mistakes - but mistakes still happen.

Statistical tools can help find them - though not always give the solution. Check by repeating analysis with changes which *should* make no difference:

- Data subsets
- Magnet up/down
- Different selection cuts
- Changing histogram bin size and fit ranges
- Changing parametrisation (including order of polynomial)
- Changing fit technique
- Looking for impossibilities
- ...

Example: the BaBar CP violation measurement “.. consistency checks, including separation of the decay by decay mode, tagging category and B_{tag} flavour... We also fit the samples of non-CP decay modes for $\sin 2\beta$ with no statistically significant difference found.”

If it passes the test

Tick the box and move on

Do **not** add the discrepancy to the systematic error



- It's illogical
- It penalises diligence
- Errors get inflated

The more tests the better. You cannot prove the analysis is correct. But the more tests it survives the more likely your colleagues¹ will be to believe the result.

¹and eventually even you

If it fails the test



Worry!

- Check the test. Very often this turns out to be faulty.
- Check the analysis. Find mistake, enjoy improvement.
- Worry. Consider whether the effect might be real. (E.g. June's results are different from July's. Temperature effect? If so can (i) compensate and (ii) introduce implicit systematic uncertainty)
- Worry harder. Ask colleagues, look at other experiments

Only as a last resort, add the term to the systematic error. Remember that this could be a hint of something much bigger and nastier

Clearing up a possible confusion

What's the difference between?

Evaluating implicit systematic errors: vary lots of parameters, see what happens to the result, and include in systematic error

Checks: vary lots of parameters, see what happens to the result, and don't include in systematic error

(1) Are you expecting to see an effect? If so, it's an evaluation, if not, it's a check

(2) Do you clearly know how much to vary them by? If so, it's an evaluation. If not, it's a check.

Cover cases such as trigger energy cut where the energy calibration is uncertain - may be simpler to simulate the effect by varying the cut.

So finally:

- 1 Thou shalt never say 'systematic error' when thou meanest 'systematic effect' or 'systematic mistake'.
- 2 Thou shalt know at all times whether what thou performest is a check for a mistake or an evaluation of an uncertainty.
- 3 Thou shalt not incorporate successful check results into thy total systematic error and make thereby a shield to hide thy dodgy result.
- 4 Thou shalt not incorporate failed check results unless thou art truly at thy wits' end.
- 5 Thou shalt not add uncertainties on uncertainties in quadrature. If they are larger than chickenfeed thou shalt generate more Monte Carlo until they shrink to become so.
- 6 Thou shalt say what thou doest, and thou shalt be able to justify it out of thine own mouth; not the mouth of thy supervisor, nor thy colleague who did the analysis last time, nor thy local statistics guru, nor thy mate down the pub.

Do these, and thou shalt flourish, and thine analysis likewise.