

# Modern Methods of Statistical Data Analysis

---

*From parameter estimation to deep learning – A guided tour of probability*

## Lecture 8

-

### Limit Setting & Unfolding

**P.-D. Dr. Roger Wolf**

[roger.wolf@kit.edu](mailto:roger.wolf@kit.edu)

**Dr. Pablo Goldenzweig**

[pablo.goldenzweig@kit.edu](mailto:pablo.goldenzweig@kit.edu)

**Dr. Jan Kieseler**

[jan.kieseler@kit.edu](mailto:jan.kieseler@kit.edu)

**Dr. Slavomira Stefkova**

[slavomira.stefkova@kit.edu](mailto:slavomira.stefkova@kit.edu)

# Today

- Limits near a physical boundary
  - Shifted and Bayesian approaches
  - Example: Upper limit on the mean of a Poisson variable with background
- Unfolding
  - Formulation of the problem
  - Matrix inversion
  - Method of correction factors
  - Regularized unfolding

## Higgs Challenge

Please mail [Sally Stefkova](#) if you plan to do it! We just want to gauge how many groups are working on this.

**Evaluations: [Lecture](#) & [Computerpraktikum](#).**

**Please take a few minutes to fill them out. Your feedback is greatly appreciated. We will take your comments into consideration in trying to improve the course.**

Evaluation period: through 22 June (lecture) & 15 July (Computerpraktikum)

# We've come a long way so far

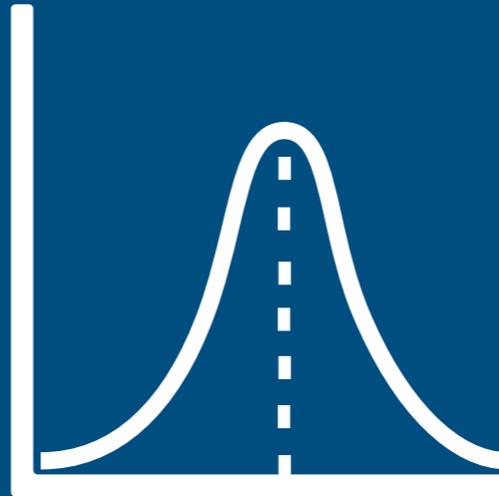
#	Lecture date	Lecture Topic
1	21.4	Fundamental concepts I
2	28.4	Fundamental concepts II
3	5.5	Monte Carlo method & production of random distributions
4	12.5	Parameter estimation & maximum likelihood
5	19.5	Chi-square method
6	26.5	Hypothesis tests & Neyman Pearson
7	9.6	Confidence intervals
8	16.6	Limit setting & unfolding
9	23.6	Event classification - Introduction and perceptron
10	30.6	Classification with the multilayer perceptron
11	7.7	Neural network training
12	14.7	Training algorithms & regularization methods
13	21.7	Training validation
14	28.7	Advanced neural networks



P.-D. Dr. Roger Wolf

Dr. Jan Kieseler

# Confidence Intervals



Statistical errors, confidence intervals and limits

Up to now: when discussing ‘error analysis’ we focused on estimating the (co)variances of estimators. This is not always adequate and other ways of communicating the statistical uncertainty of measurements have to be found.

$$\hat{\theta}_{\text{obs}} \pm \hat{\sigma}_{\hat{\theta}}$$

# Classical confidence intervals (CI)

Alternative (& often =) method of reporting the statistical uncertainty of a measurement

- Suppose you have  $n$  observations of a random variable  $X$ , which can be used to evaluate an estimator for an unknown true parameter  $\theta$ :

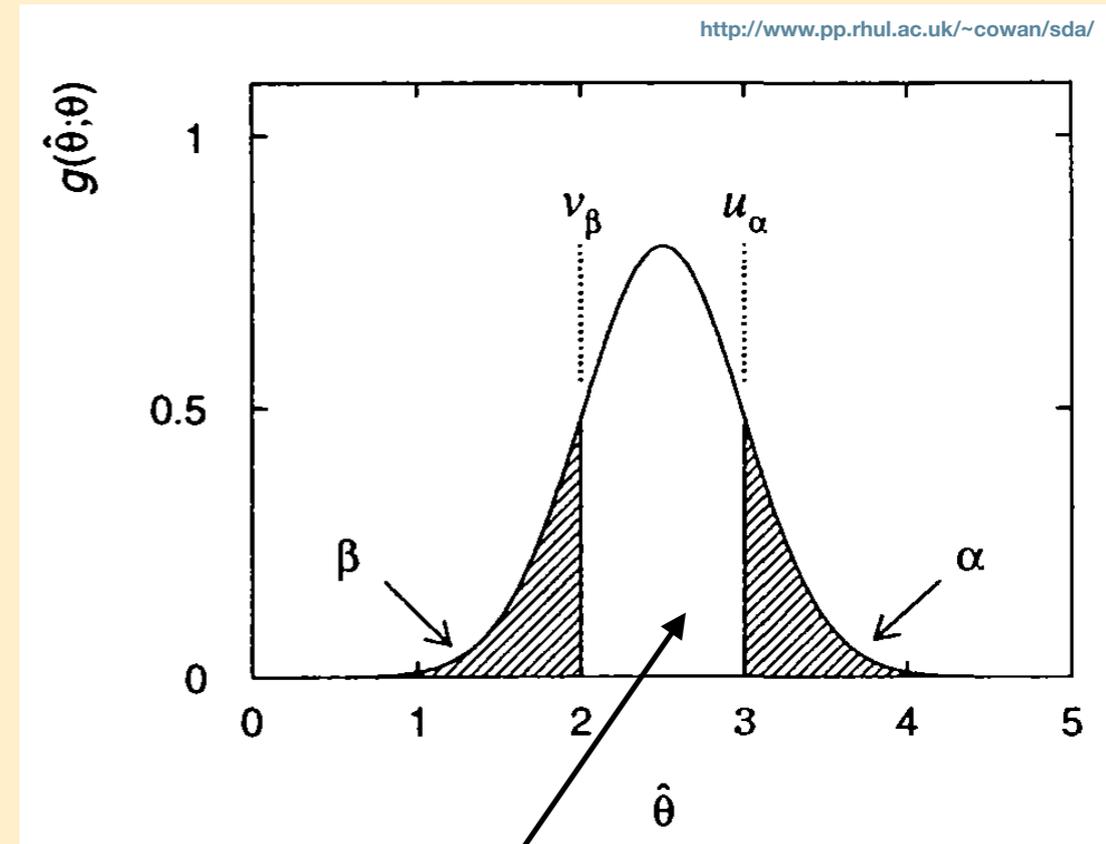
$$\hat{\theta}(x_1, \dots, x_n) = \hat{\theta}_{\text{obs}}$$

↑ value obtained

- Furthermore, suppose we know the PDF of  $\hat{\theta}$  denoted by  $g(\hat{\theta}; \theta)$

Real value of  $\theta$  unknown, BUT for a given  $\theta$  one knows what the PDF of  $\hat{\theta}$  would be

- From  $g(\hat{\theta}; \theta)$ , can determine  $\nu_\beta$  and  $u_\alpha$  such that there are fixed probabilities  $\beta$  and  $\alpha$  to observe  $\hat{\theta} < \nu_\beta$  or  $\hat{\theta} > u_\alpha$



$$P(\nu_\beta(\theta) \leq \hat{\theta} \leq u_\alpha(\theta)) = 1 - \alpha - \beta.$$

Shows the probability density for an estimator  $\hat{\theta}$  for a particular value of the true parameter  $\theta$

$u_\alpha$  and  $\nu_\beta$  depend on the true value  $\theta$  and are thus determined by

$$\beta = P(\hat{\theta} \leq \nu_\beta(\theta)) = \int_{-\infty}^{\nu_\beta(\theta)} g(\hat{\theta}; \theta) d\hat{\theta} = G(\nu_\beta(\theta); \theta),$$

$$\alpha = P(\hat{\theta} \geq u_\alpha(\theta)) = \int_{u_\alpha(\theta)}^{\infty} g(\hat{\theta}; \theta) d\hat{\theta} = 1 - G(u_\alpha(\theta); \theta),$$

CDF ... so  $\alpha$  and  $\beta$  are the probabilities!

Next: lets build the CI step by step...

# Classical confidence intervals (CI)

Alternative (& often =) method of reporting the statistical uncertainty of a measurement

- Suppose you have  $n$  observations of a random variable  $X$ , which can be used to evaluate an estimator for an unknown true parameter  $\theta$ :

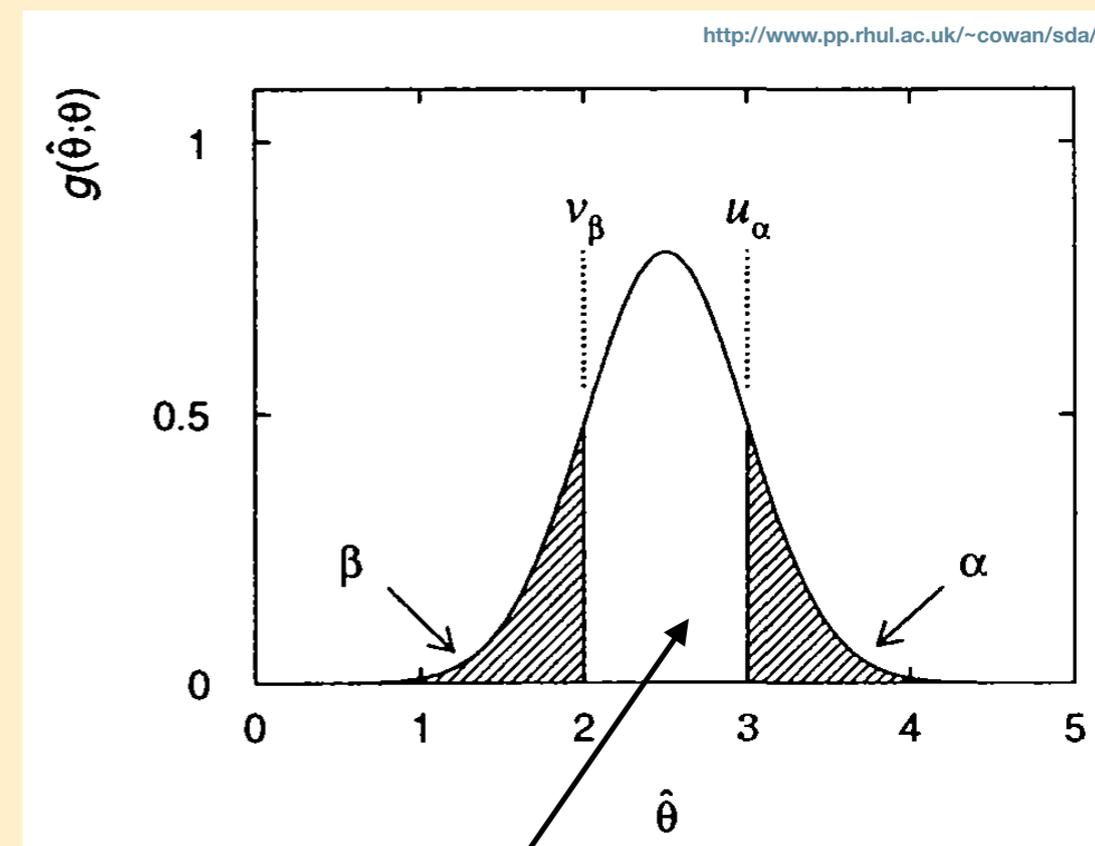
$$\hat{\theta}(x_1, \dots, x_n) = \hat{\theta}_{\text{obs}}$$

↑ value obtained

- Furthermore, suppose we know the PDF of  $\hat{\theta}$  denoted by  $g(\hat{\theta}; \theta)$

Real value of  $\theta$  unknown, BUT for a given  $\theta$  one knows what the PDF of  $\hat{\theta}$  would be

- From  $g(\hat{\theta}; \theta)$ , can determine  $\nu_\beta$  and  $u_\alpha$  such that there are fixed probabilities  $\beta$  and  $\alpha$  to observe  $\hat{\theta} < \nu_\beta$  or  $\hat{\theta} > u_\alpha$



$$P(\nu_\beta(\theta) \leq \hat{\theta} \leq u_\alpha(\theta)) = 1 - \alpha - \beta.$$

Shows the probability density for an estimator  $\hat{\theta}$  for a particular value of the true parameter  $\theta$

$u_\alpha$  and  $\nu_\beta$  depend on the true value  $\theta$  and are thus determined by

$$\beta = P(\hat{\theta} \leq \nu_\beta(\theta)) = \int_{-\infty}^{\nu_\beta(\theta)} g(\hat{\theta}; \theta) d\hat{\theta} = G(\nu_\beta(\theta); \theta),$$

$$\alpha = P(\hat{\theta} \geq u_\alpha(\theta)) = \int_{u_\alpha(\theta)}^{\infty} g(\hat{\theta}; \theta) d\hat{\theta} = 1 - G(u_\alpha(\theta); \theta),$$

CDF ... so  $\alpha$  and  $\beta$  are the probabilities!

Next: lets build the CI step by step...

# Classical confidence intervals (CI)

- Suppose you have  $n$  observations of a random variable  $X$ , which can be used to evaluate an estimator for an unknown true parameter  $\theta$ :

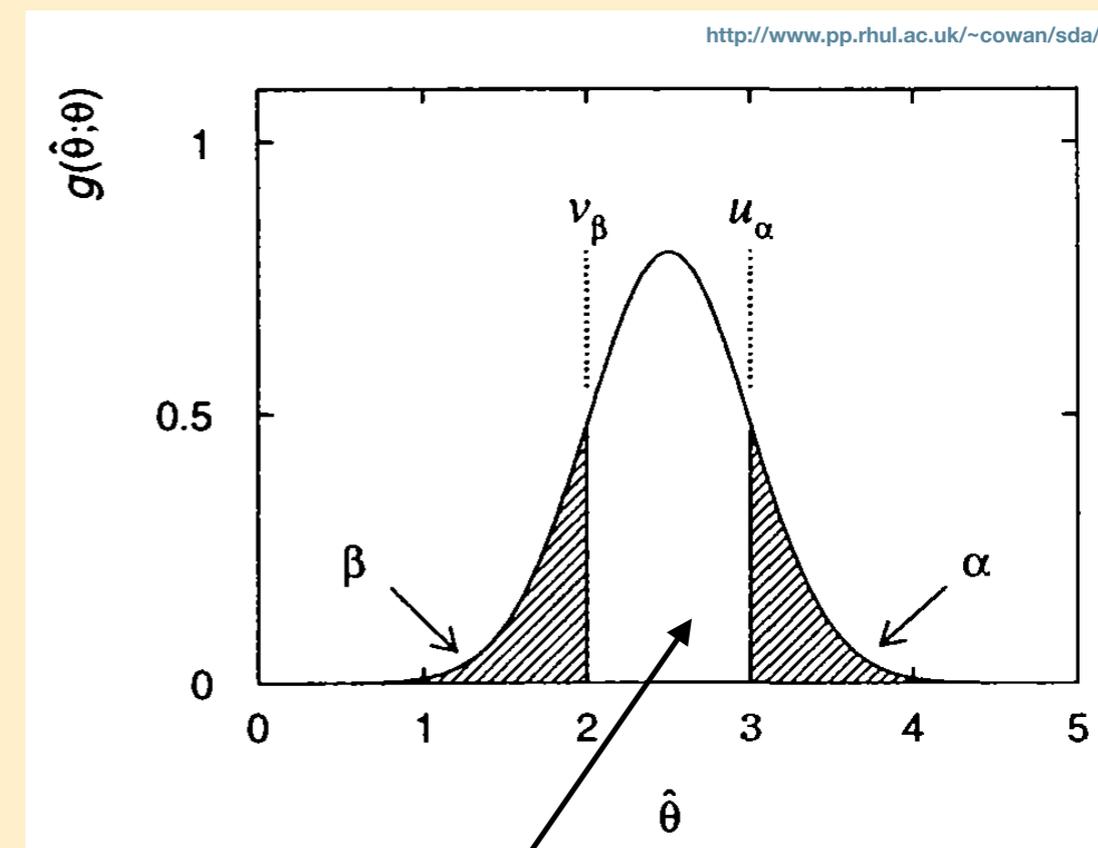
$$\hat{\theta}(x_1, \dots, x_n) = \hat{\theta}_{\text{obs}}$$

↑ value obtained

- Furthermore, suppose we know the PDF of  $\hat{\theta}$  denoted by  $g(\hat{\theta}; \theta)$

Real value of  $\theta$  unknown, BUT for a given  $\theta$  one knows what the PDF of  $\hat{\theta}$  would be

- From  $g(\hat{\theta}; \theta)$ , can determine  $\nu_\beta$  and  $u_\alpha$  such that there are fixed probabilities  $\beta$  and  $\alpha$  to observe  $\hat{\theta} < \nu_\beta$  or  $\hat{\theta} > u_\alpha$



$$P(\nu_\beta(\theta) \leq \hat{\theta} \leq u_\alpha(\theta)) = 1 - \alpha - \beta.$$

Shows the probability density for an estimator  $\hat{\theta}$  for a particular value of the true parameter  $\theta$

$u_\alpha$  and  $\nu_\beta$  depend on the true value  $\theta$  and are thus determined by

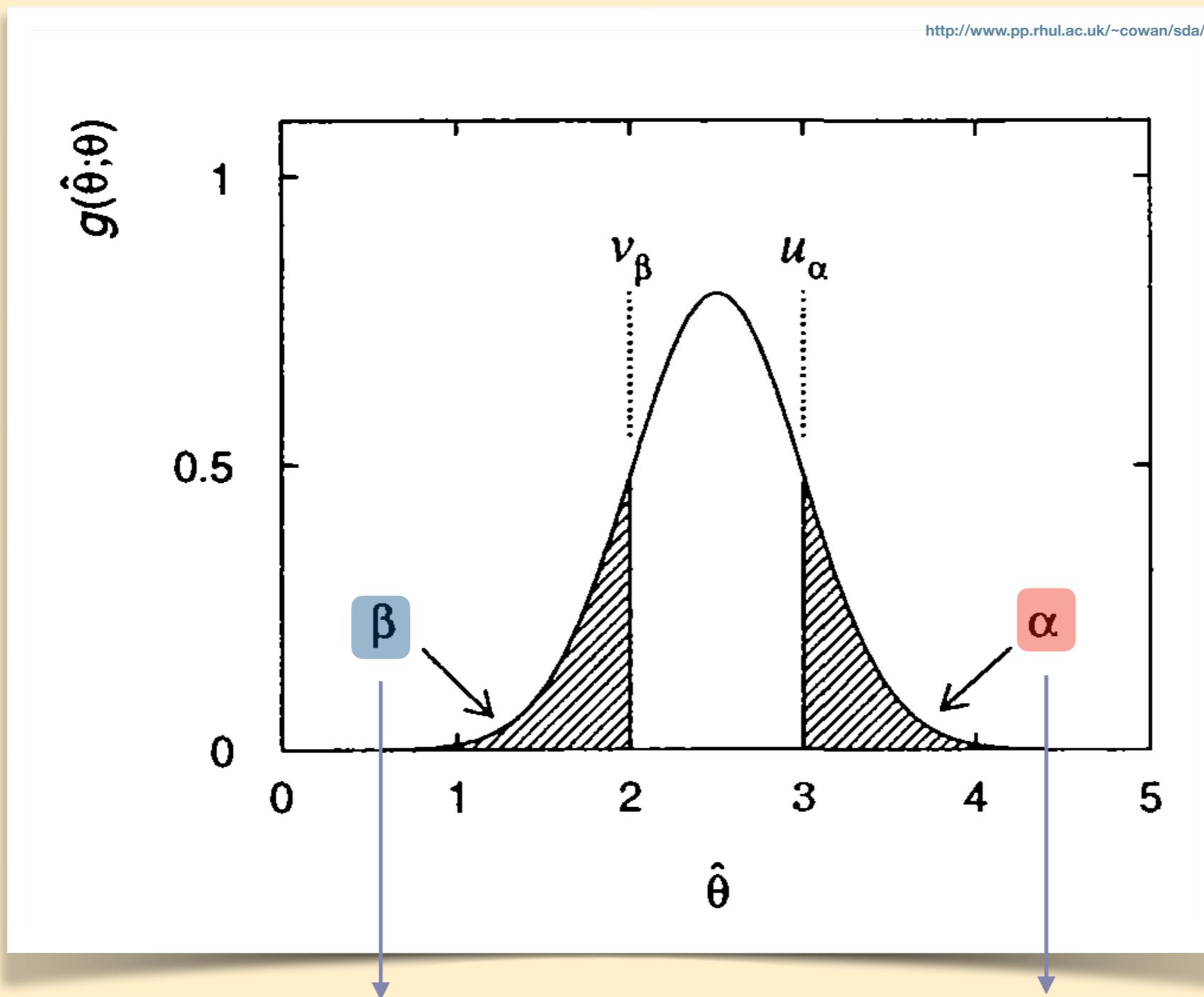
$$\beta = P(\hat{\theta} \leq \nu_\beta(\theta)) = \int_{-\infty}^{\nu_\beta(\theta)} g(\hat{\theta}; \theta) d\hat{\theta} = G(\nu_\beta(\theta); \theta),$$

$$\alpha = P(\hat{\theta} \geq u_\alpha(\theta)) = \int_{u_\alpha(\theta)}^{\infty} g(\hat{\theta}; \theta) d\hat{\theta} = 1 - G(u_\alpha(\theta); \theta),$$

CDF ... so  $\alpha$  and  $\beta$  are the probabilities!

Next: lets build the CI step by step...

# Confidence Belt (i)

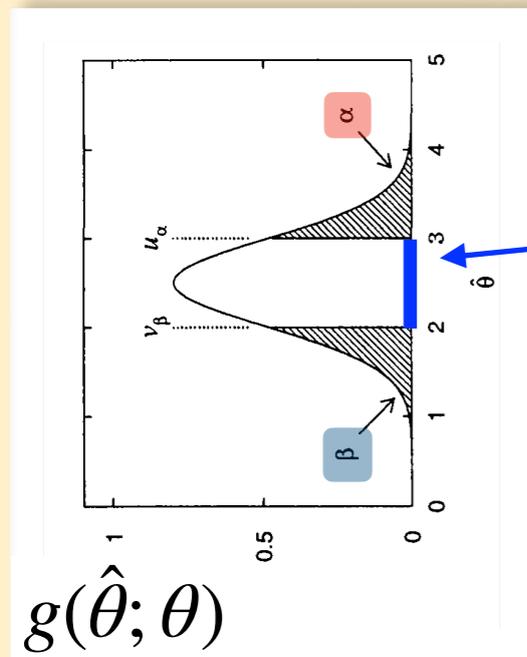


$$\beta = P(\hat{\theta} \leq v_{\beta}(\theta)) = \int_{-\infty}^{v_{\beta}(\theta)} g(\hat{\theta}; \theta) d\hat{\theta} = G(v_{\beta}(\theta); \theta),$$

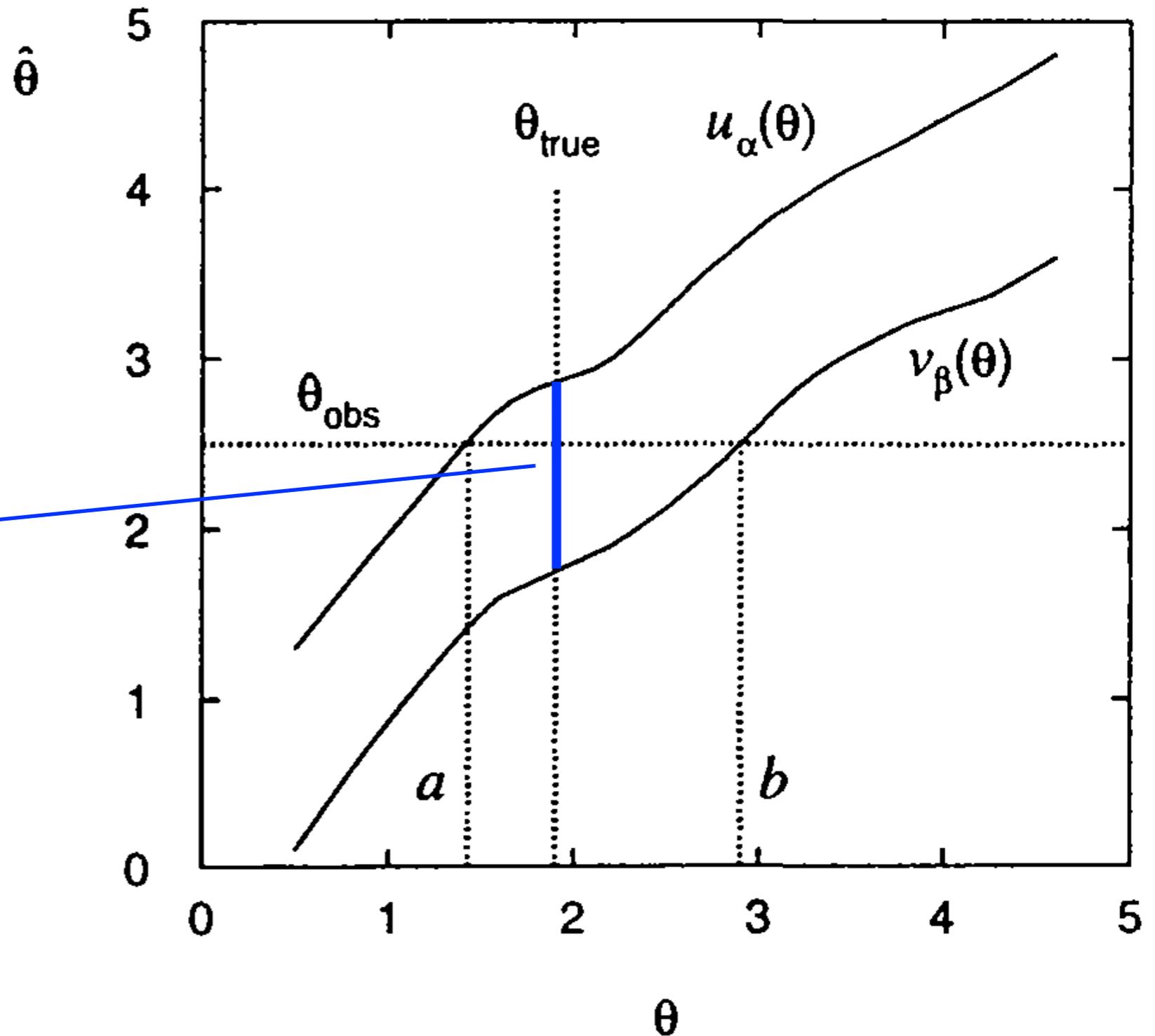
$$\alpha = P(\hat{\theta} \geq u_{\alpha}(\theta)) = \int_{u_{\alpha}(\theta)}^{\infty} g(\hat{\theta}; \theta) d\hat{\theta} = 1 - G(u_{\alpha}(\theta); \theta),$$

# Confidence Belt (ii)

<http://www.pp.rhul.ac.uk/~cowan/sda/>

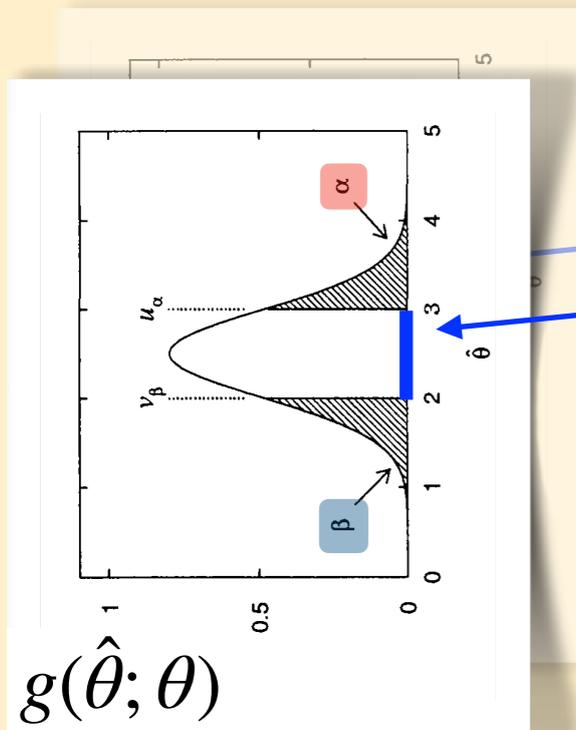


for a given value of (true)  $\theta$

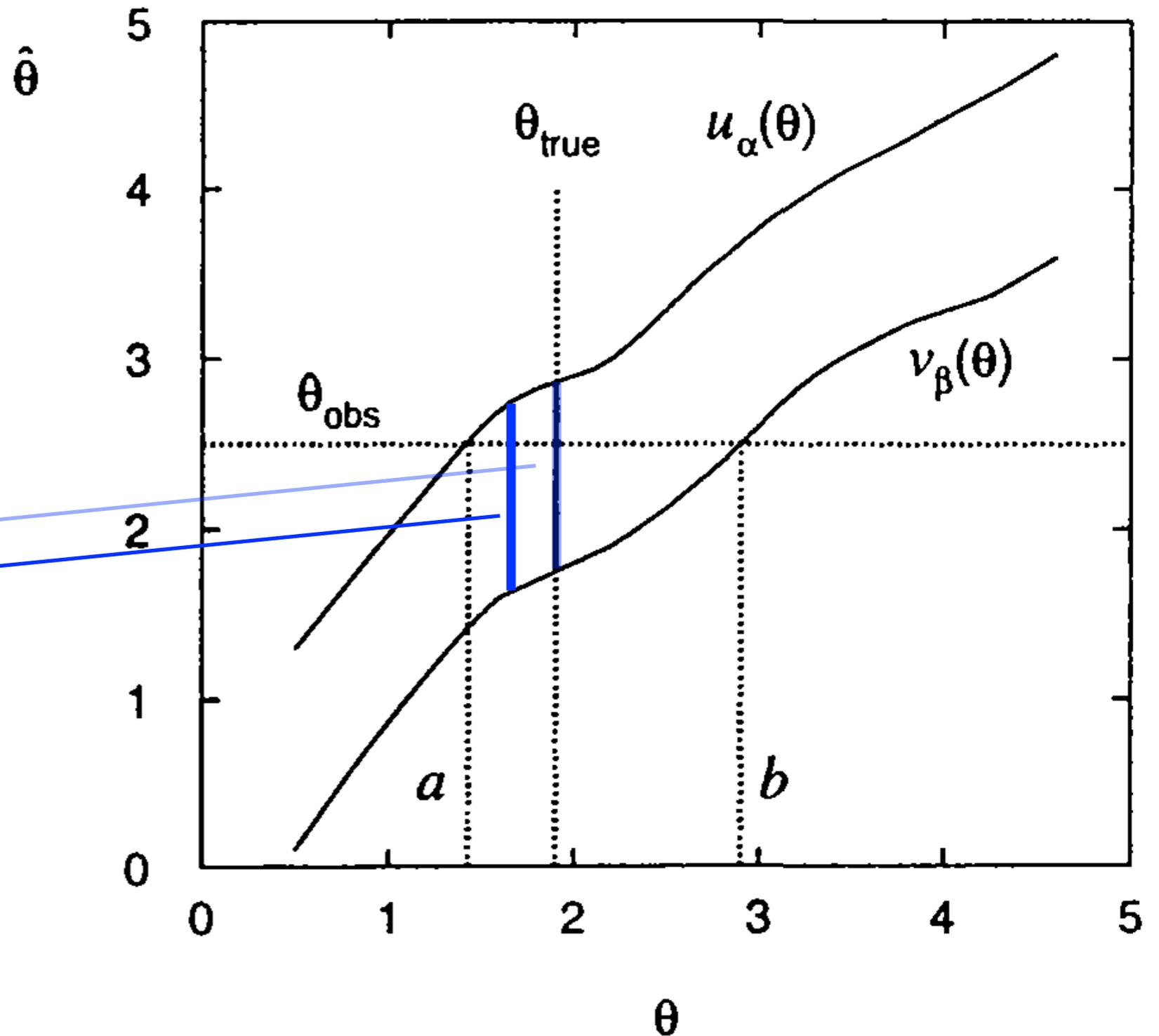


# Confidence Belt (iii)

<http://www.pp.rhul.ac.uk/~cowan/sda/>



for *another* given value of (true)  $\theta$



# Confidence Belt (iv)

<http://www.pp.rhul.ac.uk/~cowan/sda/>

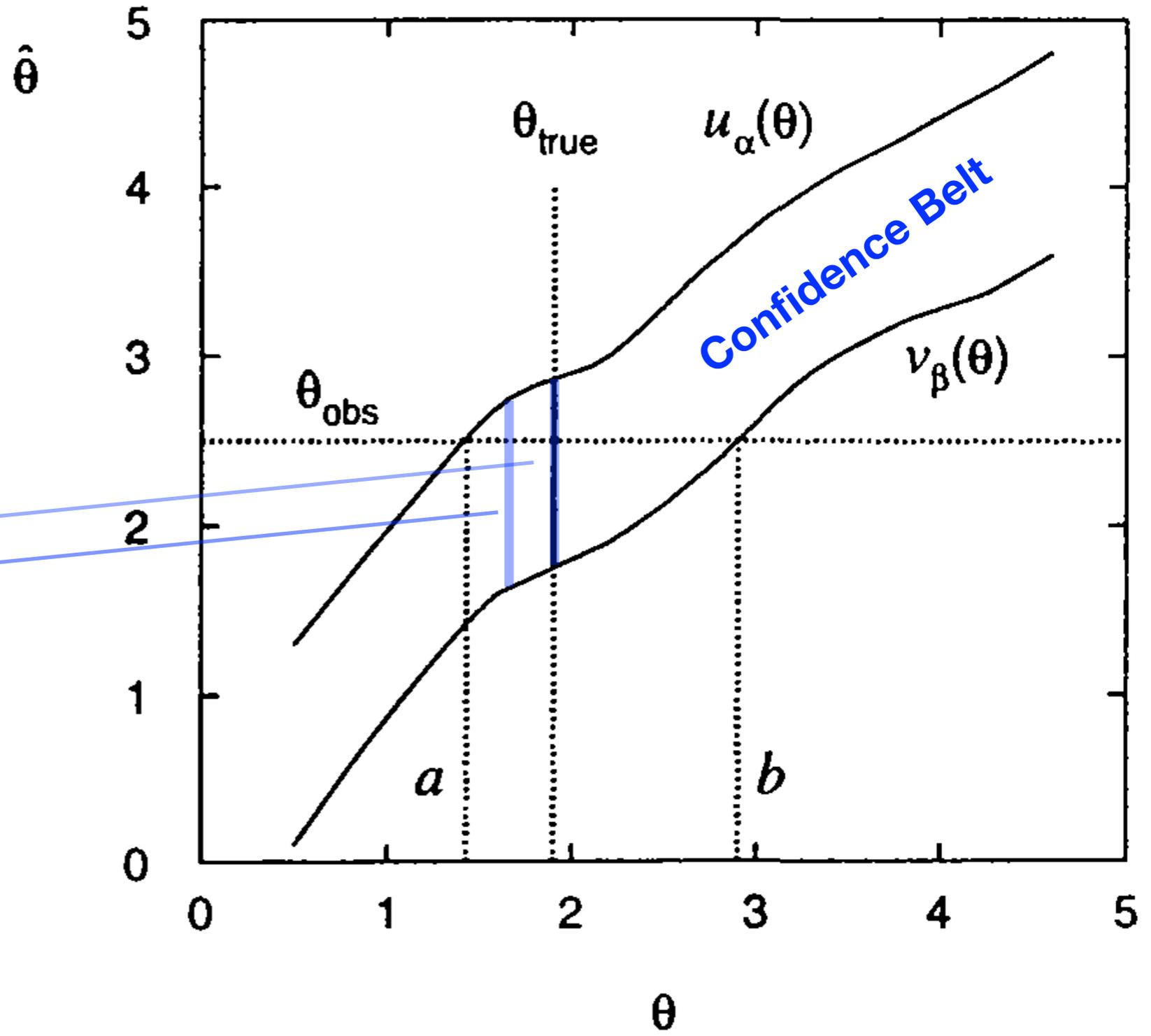
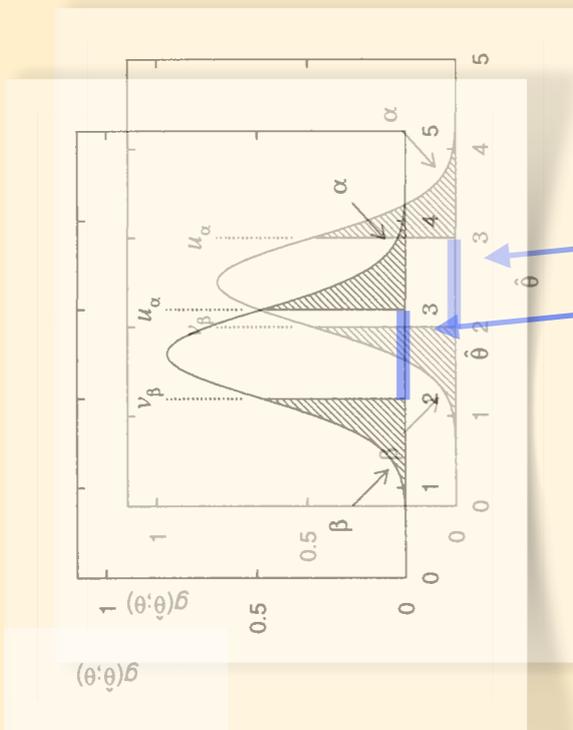
Region between the curves:

## Confidence Belt

(Neyman Belt)

$$P(v_\beta(\theta) \leq \hat{\theta} \leq u_\alpha(\theta)) = 1 - \alpha - \beta.$$

The probability for the estimator  $\hat{\theta}$  to be inside the belt, regardless of the value of  $\theta$



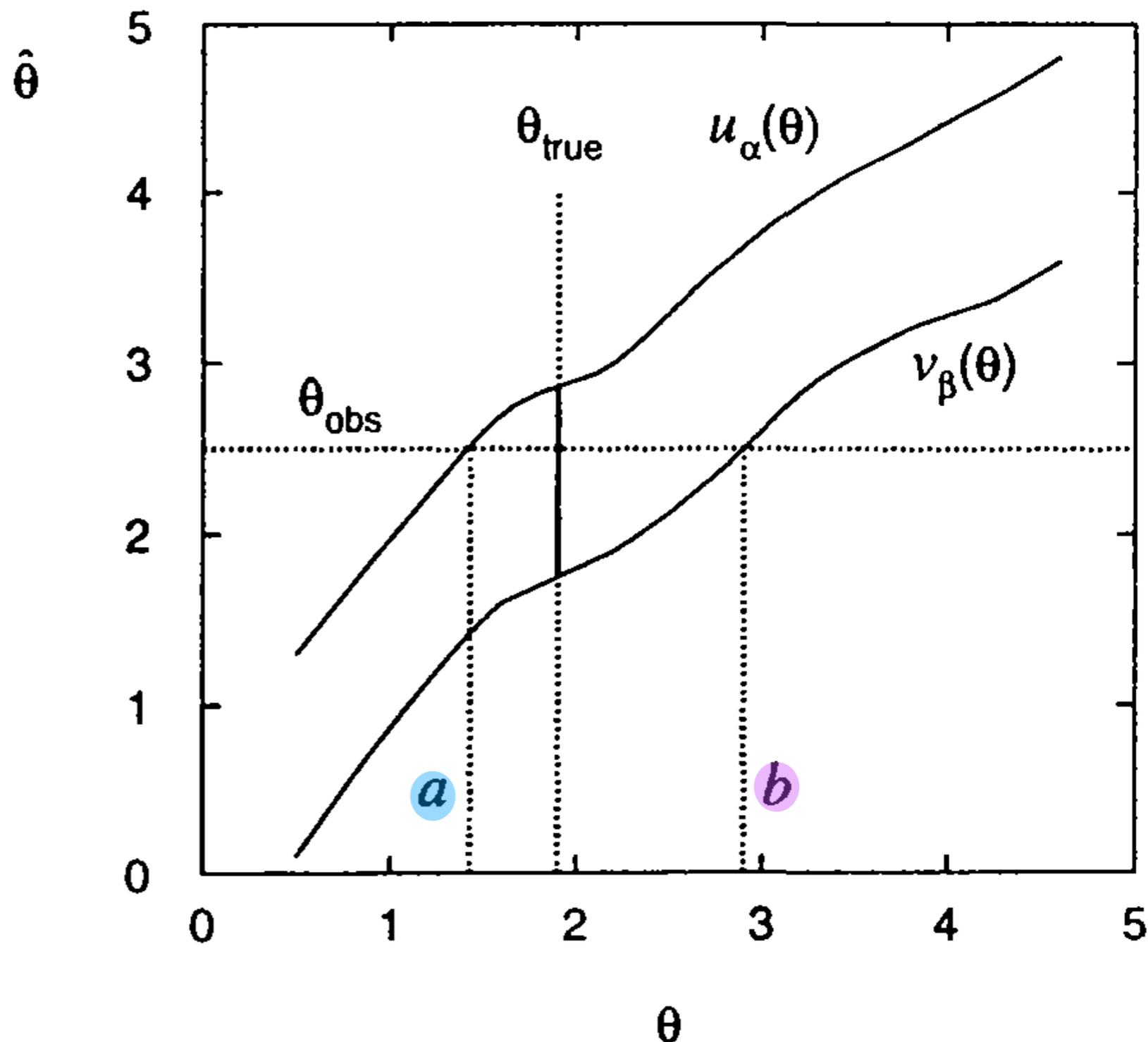
# Confidence Interval (i)

If  $u_\alpha(\theta)$  and  $v_\beta(\theta)$  are monotonically increasing functions of  $\theta$ , then one can determine the inverse functions

$$a(\hat{\theta}) \equiv u_\alpha^{-1}(\hat{\theta}),$$

$$b(\hat{\theta}) \equiv v_\beta^{-1}(\hat{\theta}).$$

(Should be the case if  $\hat{\theta}$  is a good estimator for  $\theta$ )



# Confidence Interval (ii)

$$a(\hat{\theta}) \equiv u_{\alpha}^{-1}(\hat{\theta}),$$

$$b(\hat{\theta}) \equiv v_{\beta}^{-1}(\hat{\theta}).$$

This then implies:

$$\begin{aligned} \hat{\theta} &\geq u_{\alpha}(\theta), \\ \hat{\theta} &\leq v_{\beta}(\theta), \end{aligned}$$

→  
invert

$$\begin{aligned} a(\hat{\theta}) &\geq \theta, \\ b(\hat{\theta}) &\leq \theta. \end{aligned}$$

$$\begin{aligned} P(a(\hat{\theta}) \geq \theta) &= \alpha, \\ P(b(\hat{\theta}) \leq \theta) &= \beta, \end{aligned}$$

or

$$P(v_{\beta}(\theta) \leq \hat{\theta} \leq u_{\alpha}(\theta)) = 1 - \alpha - \beta.$$

$$P(a(\hat{\theta}) \leq \theta \leq b(\hat{\theta})) = 1 - \alpha - \beta.$$

If the functions  $a(\hat{\theta})$  and  $b(\hat{\theta})$  are evaluated with the value of the estimator obtained in the experiment ( $\hat{\theta}_{\text{obs}}$ ), then this determines 2 values **[a, b]**

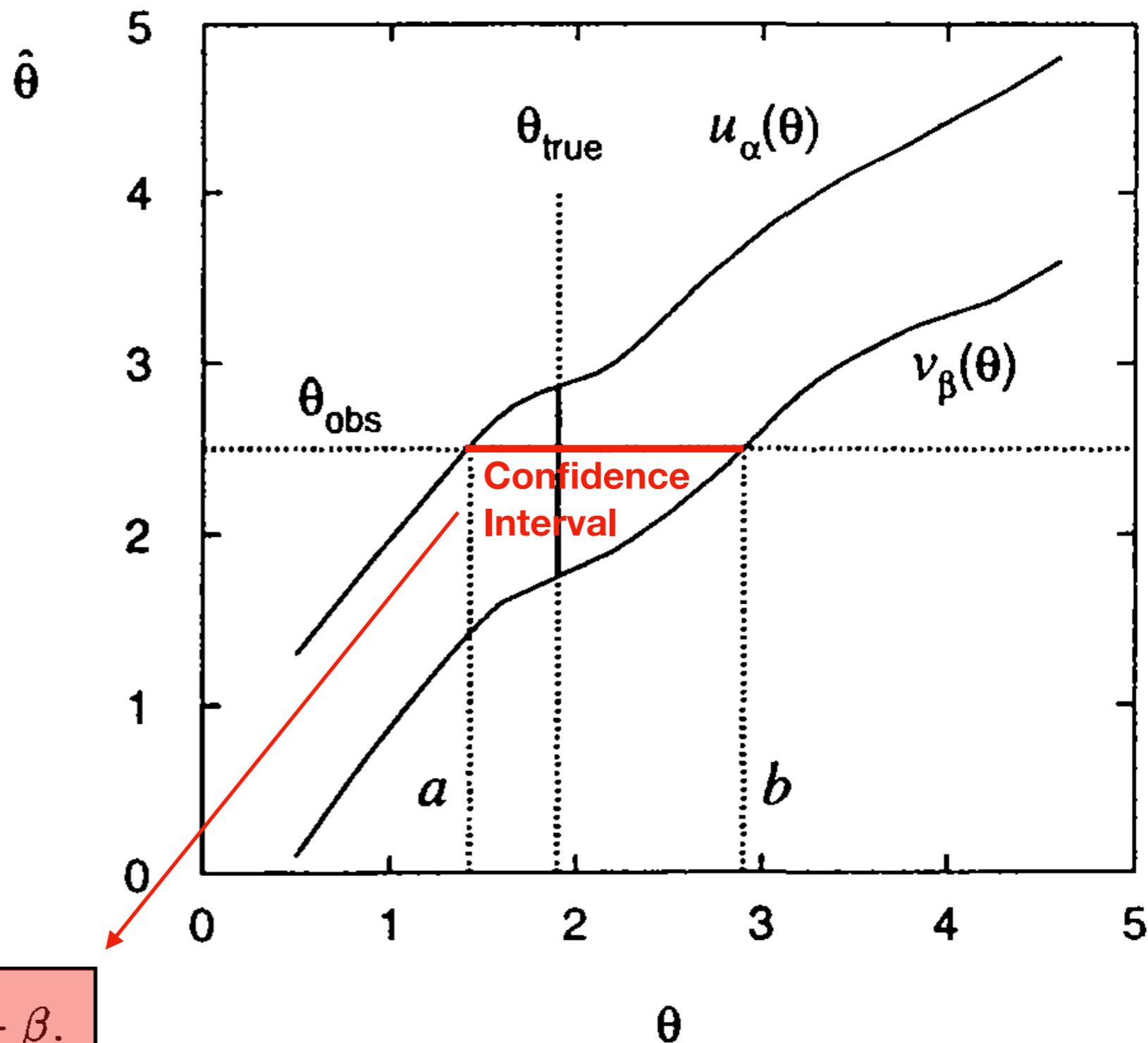
# Confidence Interval (iii)

Often chooses  $\alpha = \beta = \frac{\gamma}{2}$   
 giving a so-called **central CI**  
 with probability  $= 1 - \gamma$

**$[a, b]$ : Confidence Interval,**  
 at a **confidence level** (or  
**coverage probability**) of  
 $1 - \alpha - \beta$

$$P(a(\hat{\theta}) \leq \theta \leq b(\hat{\theta})) = 1 - \alpha - \beta.$$

<http://www.pp.rhul.ac.uk/~cowan/sda/>



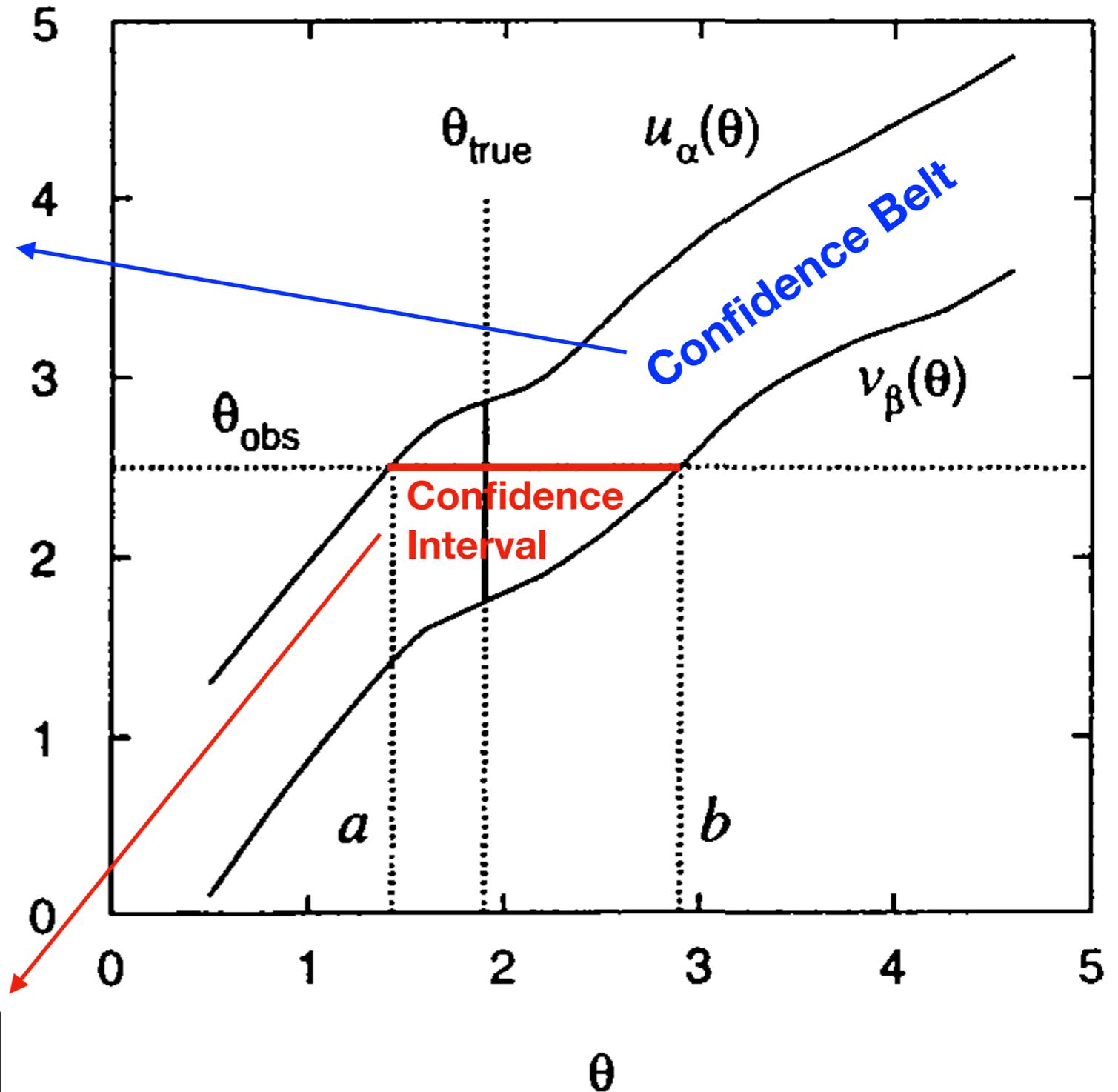
# All together now

$$P(v_{\beta}(\theta) \leq \hat{\theta} \leq u_{\alpha}(\theta)) = 1 - \alpha - \beta.$$

Note where the  $\hat{s}$  are in the 2 equations!

$[a, b]$ : Confidence Interval, at a confidence level (or coverage probability) of  $1 - \alpha - \beta$

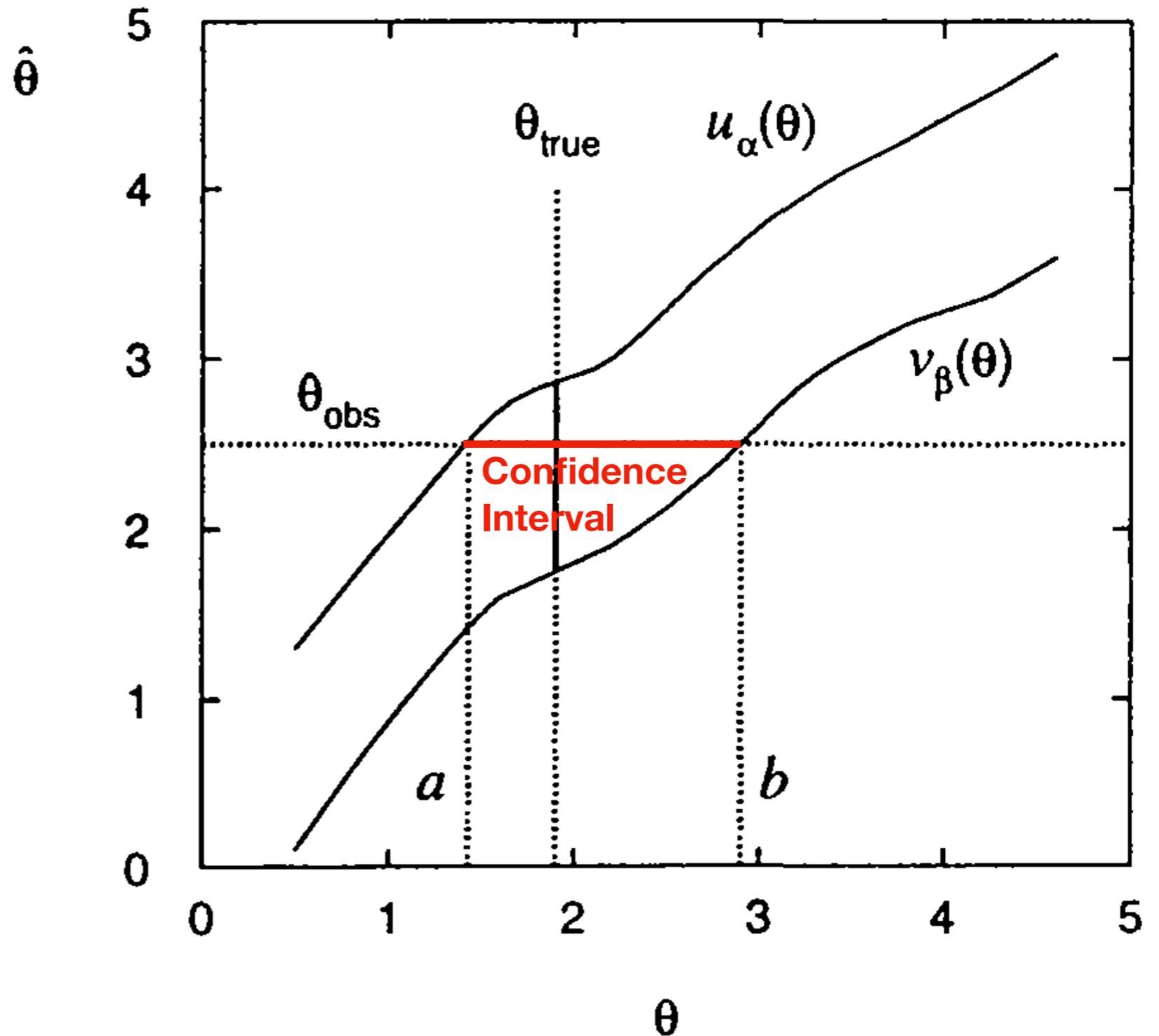
$\hat{\theta}$



$$P(a(\hat{\theta}) \leq \theta \leq b(\hat{\theta})) = 1 - \alpha - \beta.$$

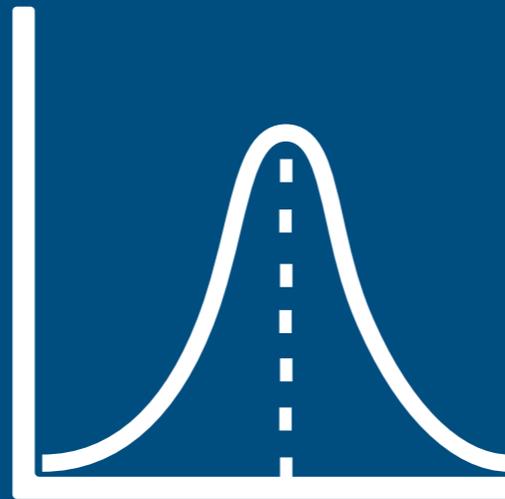
# Take home message

If the experiment were repeated many times, the interval  $[a, b]$  would include the true value of the parameter  $\theta$  in a fraction  $1 - \alpha - \beta$  of the experiments





# Gaussian Confidence Intervals



Let's apply what we've built up to the Gaussian limit

- **Simple and very important application:**

- $\hat{\theta}$  is Gaussian with mean  $\theta$  and standard deviation  $\sigma_{\hat{\theta}}$
- Cumulative distribution of  $\hat{\theta}$  is then

Commonly occurring situation, since according to the CLT, any estimator that is a linear function of a sum of RVs becomes Gaussian in the large sample limit

$$G(\hat{\theta}; \theta, \sigma_{\hat{\theta}}) = \int_{-\infty}^{\hat{\theta}} \frac{1}{\sqrt{2\pi\sigma_{\hat{\theta}}^2}} \exp\left(-\frac{(\hat{\theta}' - \theta)^2}{2\sigma_{\hat{\theta}}^2}\right) d\hat{\theta}'.$$

- Suppose that **the standard deviation is known** and that the experiment resulted in an estimate  $\hat{\theta}_{\text{obs}}$ . Then we can determine the **confidence interval**  $[a, b]$  by solving

$$\alpha = 1 - G(\hat{\theta}_{\text{obs}}; a, \sigma_{\hat{\theta}}) = 1 - \Phi\left(\frac{\hat{\theta}_{\text{obs}} - a}{\sigma_{\hat{\theta}}}\right),$$

$$\beta = G(\hat{\theta}_{\text{obs}}; b, \sigma_{\hat{\theta}}) = \Phi\left(\frac{\hat{\theta}_{\text{obs}} - b}{\sigma_{\hat{\theta}}}\right),$$

standard normal CDF  
 $\Phi = G(\hat{\mu}; \mu = 0, \sigma = 1)$

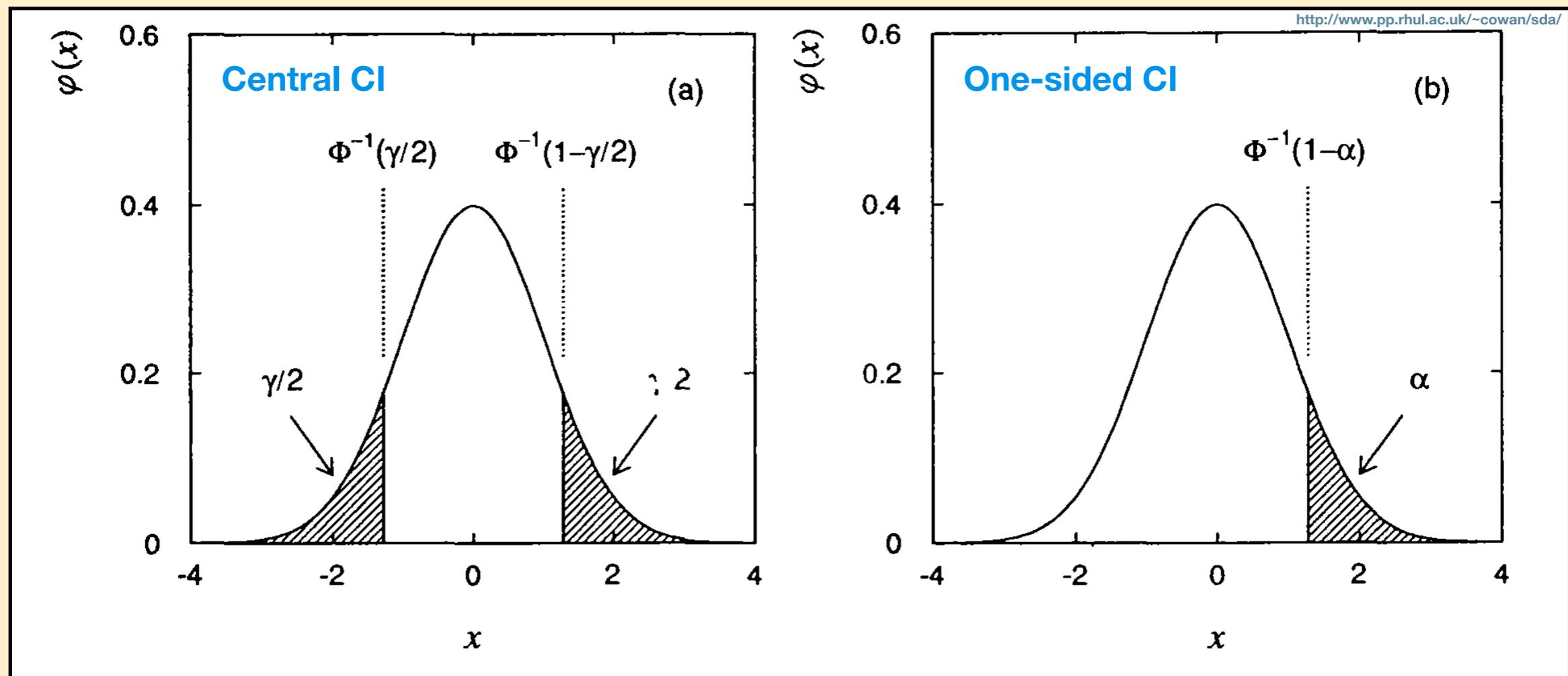
# CI for Gaussian distributed estimators

- This results in

$$a = \hat{\theta}_{\text{obs}} - \sigma_{\hat{\theta}} \Phi^{-1}(1 - \alpha),$$
$$b = \hat{\theta}_{\text{obs}} + \sigma_{\hat{\theta}} \Phi^{-1}(1 - \beta).$$

*i.e., the inverse function of  $\Phi$  equals the quantile of the std. Gaussian*

inverse of standard normal CDF



The relationship between the quantiles of the std. Gaussian distribution and the CI

- Consider a central confidence interval with  $\alpha = \beta = \gamma/2$ 
  - The confidence level  $(1 - \gamma)$  is often chosen, such that  $\Phi^{-1}(1 - \gamma/2)$  is a small integer (e.g., 1,2,3)
  - Similarly, one-sided intervals are often small integer values
    - Sometimes one also prefers to use a round value for  $1 - \alpha$  or  $1 - \gamma$

$\Phi^{-1}(1 - \gamma/2)$	$1 - \gamma$	$\Phi^{-1}(1 - \alpha)$	$1 - \alpha$
1	0.6827	1	0.8413
2	0.9544	2	0.9772
3	0.9973	3	0.9987

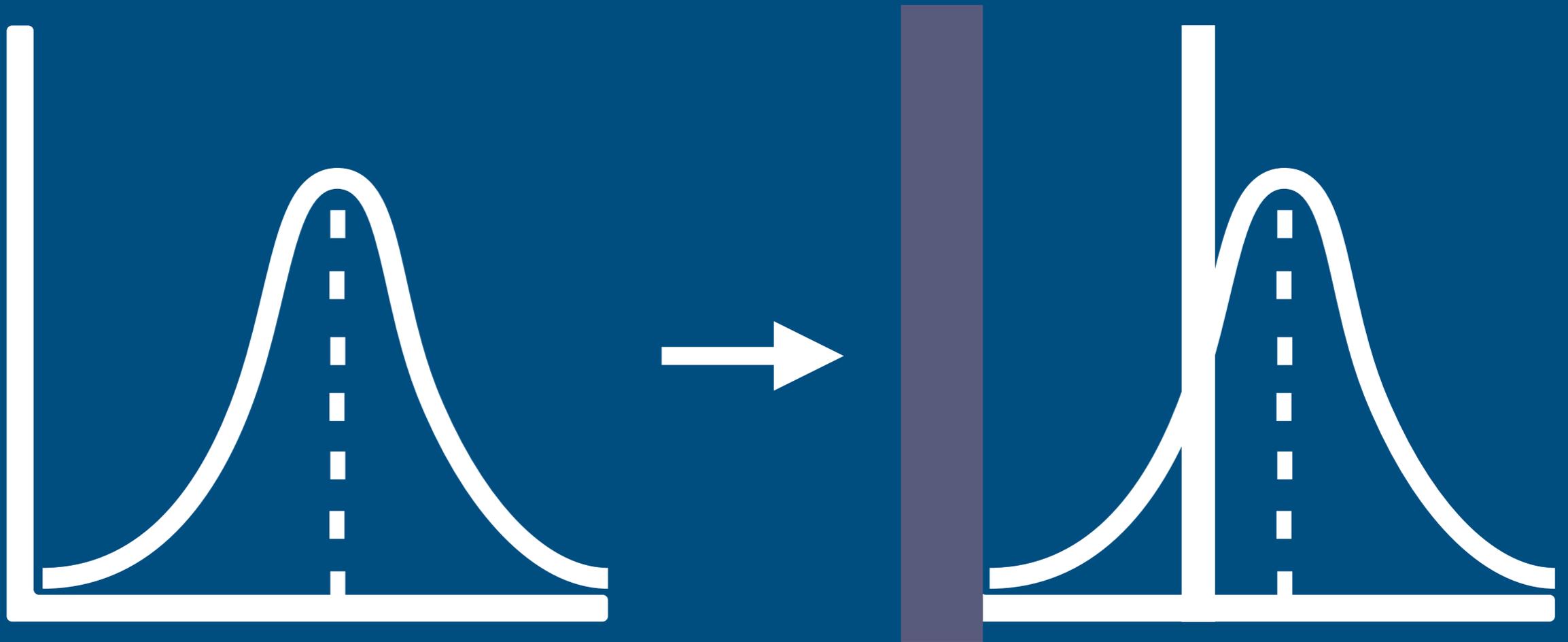
$1 - \gamma$	$\Phi^{-1}(1 - \gamma/2)$	$1 - \alpha$	$\Phi^{-1}(1 - \alpha)$
0.90	1.645	0.90	1.282
0.95	1.960	0.95	1.645
0.99	2.576	0.99	2.326

- For **conventional 68.3% CI** one has

$$[a, b] = [\hat{\theta}_{\text{obs}} - \sigma_{\hat{\theta}}, \hat{\theta}_{\text{obs}} + \sigma_{\hat{\theta}}].$$

- All of this is valid, if  $\sigma_{\hat{\theta}}$  is known
  - Often not the case, *but in large n limit* can use  $\sigma_{\hat{\theta}} \rightarrow \hat{\sigma}_{\hat{\theta}}$

# Limits near a physical boundary





# Limits near a physical boundary

- How to place a limit on  $m^2$  when the estimate is near an excluded or unphysical region?
- Let's make this more concrete with an example:
  - $\hat{\theta} = x - y$  with  $x, y$  Gaussian RVs with mean and variances  $\mu_x, \mu_y, \sigma_x^2, \sigma_y^2$
  - The difference is also a Gaussian variable with  $\theta = \mu_x - \mu_y$  and  $\sigma_{\hat{\theta}}^2 = \sigma_x^2 + \sigma_y^2$   
(see proof in characteristic functions chapter 10 Cowan)
  - Assume that  $\theta$  is known a priori to be non-negative (e.g. like the mass squared) and suppose the experiment resulted in a value  $\hat{\theta}_{\text{obs}}$  for the estimator  $\hat{\theta}$
  - According to what we discussed (S17 *Review*), the upper limit  $\theta_{\text{up}}$  at CL  $1 - \beta$  is

$$\theta_{\text{up}} = \hat{\theta}_{\text{obs}} + \sigma_{\hat{\theta}} \Phi^{-1}(1 - \beta)$$

↑ inverse of standard normal CDF

# Limits near a physical boundary

- For the commonly used 95% CL one obtains the quantile

$$\Phi^{-1}(0.95) = 1.645$$

Table 9.2 (L07, S33)

$1 - \gamma$	$\Phi^{-1}(1 - \gamma/2)$	$1 - \alpha$	$\Phi^{-1}(1 - \alpha)$
0.90	-1.645	0.90	1.282
0.95	1.960	0.95	1.645
0.99	2.576	0.99	2.326

- The interval  $(-\infty, \theta_{\text{up}}]$  is constructed to include the true value  $\theta$  with a probability of 95%, *independent of the true value*.
- Let's now suppose the standard deviation  $\sigma_{\hat{\theta}} = 1$  and the observed value from the experiment is  $\hat{\theta}_{\text{obs}} = -2.0$ 
  - Using  $\theta_{\text{up}} = \hat{\theta}_{\text{obs}} + \sigma_{\hat{\theta}} \Phi^{-1}(1 - \beta)$  we obtain  $\theta_{\text{up}} = -0.355$  at 95% CL
- Not only is the observed value in the unphysical region (half of the estimates actually should be if  $\theta$  is zero), *but the upper limit is below zero as well*
  - Not particularly unusual; we expect 5% of all experiments to report this if  $\theta$  is zero.

# Nothing went wrong!

- As far as the definition of CL is concerned, **nothing fundamental has gone wrong.**
  - The **interval** was **designed** to **cover** the **true value of  $\theta$**  in a **certain fraction of repeated experiments**, and we have obviously encountered one of those experiments where  $\theta$  is not in the interval
  - But many people **don't find this very satisfying**, since we already know from physical reasons that  $\theta$  is greater than zero (and certainly greater than  $\theta_{\text{up}} = -0.355$ ) without having to perform an experiment.
- **Regardless of the upper limit**, it is **important to report the actual value** of the estimate obtained and its standard deviation, i.e.

$$\hat{\theta}_{\text{obs}} \pm \sigma_{\hat{\theta}} \quad \text{or if Errors are non-Gaussian: the likelihood function } \mathcal{L}(\theta)$$

- In this way, the **average of many experiments will converge to the correct value** as long as the estimator is unbiased.

# Upper Limit Bonanza

Table 9.2 (L07, S33)

$1 - \gamma$	$\Phi^{-1}(1 - \gamma/2)$	$1 - \alpha$	$\Phi^{-1}(1 - \alpha)$
0.90	-1.645	0.90	1.282
0.95	1.960	0.95	1.645
0.99	2.576	0.99	2.326

- **Nevertheless**, most experimenters want to report some sort of upper limit, that takes into account the knowledge of the unphysical region.
  - Many different solutions have been proposed, but there is no established convention on how this should be done. **So it's imperative to state what procedure you used. Otherwise people will not be able to combine or use your result.**
- To come back to our example:  $\hat{\theta}_{\text{obs}} = -2.0$ ,  $\sigma_{\hat{\theta}} = 1$ 
  - One might **feel tempted** to just **quote a limit at a higher CL**, e.g. 99% would result in  $\theta_{\text{up}} = 0.326$  ( $\Phi^{-1}(0.99) = 2.326$ )
  - This would lead to an **upper limit better than the intrinsic resolution** of our **experiment** ( $\sigma_{\hat{\theta}} = 1$ ) at a **very high confidence level** of 99%
    - *This is a bit misleading...*
  - But even worse would be to adjust the CL to give an arbitrary small limit,  $\theta_{\text{up}} = 10^{-5}$  at 97.725 CL%

# Alternative approaches: Max method

- In order to avoid such difficulties, a commonly used technique is to simply shift a negative estimate to zero before determining the value, i.e.

$$\theta_{\text{up}} = \max(\hat{\theta}_{\text{obs}}, 0) + \sigma_{\hat{\theta}} \Phi^{-1}(1 - \beta)$$

- This way the upper limit is always at least the same order of magnitude as the resolution of the experiment
  - If  $\hat{\theta}_{\text{obs}}$  is positive, nothing changes and the upper limit coincides with the classical procedure. (See Fig. on slide 30.)
- This technique has a certain intuitive appeal and is often used, but the **interpretation as an interval that will cover the true parameter** with a probability  $1 - \beta$  **no longer applies**.
  - The coverage probability is clearly larger than  $1 - \beta$  (one speaks of **over-coverage**)

# Alternative approaches: Bayesian limit

- Another alternative is to report an interval based on the **Bayesian posterior PDF**  $p(\theta | \mathbf{x})$ , obtained via

$$p(\theta | \mathbf{x}) = \frac{\mathcal{L}(\mathbf{x} | \theta) \pi(\theta)}{\int \mathcal{L}(\mathbf{x} | \theta') \pi(\theta') d\theta'}$$

Likelihood function                      Prior PDF of  $\theta$

Observed data

*(Reflects the state of knowledge of  $\theta$  before consideration of the data)*

- We now can use  $p(\theta | \mathbf{x})$  to determine an interval  $[a, b]$  such that for given probabilities  $\alpha$  and  $\beta$  one has

$$\alpha = \int_{-\infty}^a p(\theta | \mathbf{x}) d\theta$$

$$\beta = \int_b^{\infty} p(\theta | \mathbf{x}) d\theta.$$

# Alternative approaches: Bayesian limit

- Choosing  $\alpha = \beta$  gives a central interval with e.g.  $1 - \alpha - \beta = 68.3\%$
- Another possibility is to choose  $\alpha$  and  $\beta$  s.t. all **values of  $p(\theta | \mathbf{x})$**  inside the interval  $[a, b]$  are **higher than any values outside**, which implies  $p(a | \mathbf{x}) = p(b | \mathbf{x})$ . One can show that this gives the **shortest possible interval**.
- One advantage of the Bayesian interval, is that the prior knowledge, e.g.  $\theta \geq 0$  can easily be incorporated **by setting the prior PDF to zero in the excluded region**.
  - Bayes' Theorem then gives a posterior probability  $p(\theta | \mathbf{x})$  with  $p(\theta | \mathbf{x}) = 0$  for  $\theta < 0$ . The upper limit thus is given by

$$1 - \beta = \int_{-\infty}^{\theta_{\text{up}}} p(\theta | \mathbf{x}) d\theta = \frac{\int_{-\infty}^{\theta_{\text{up}}} L(\mathbf{x} | \theta) \pi(\theta) d\theta}{\int_{-\infty}^{\infty} L(\mathbf{x} | \theta) \pi(\theta) d\theta}.$$

# Bayesian limit: constant prior

- The difficulties with this approach is that there is **no unique way** to **specify** the **prior density**  $\pi(\theta)$ . A **common choice** is:

$$\pi(\theta) = \begin{cases} 0 & \theta < 0 \\ 1 & \theta \geq 0. \end{cases}$$

- I.e.: Normalize the likelihood function to unit area in the physical region, and then integrate it out to  $\theta_{\text{up}}$  s.t. the fraction of the area covered is  $1 - \beta$ .
- **Although the method is simple, it has some conceptual drawbacks:**
  - For the case where one knows  $\theta \geq 0$  (e.g. Neutrino mass), one does not really believe that  $0 < \theta < 1$  has the same prior probability as  $10^{40} < \theta < 10^{40} + 1$
  - Furthermore the upper limit derived from  $\pi(\theta) = \text{const.}$  is **not invariant with respect to a nonlinear transformation of the parameter.**

# Bayesian limit: Jeffreys prior

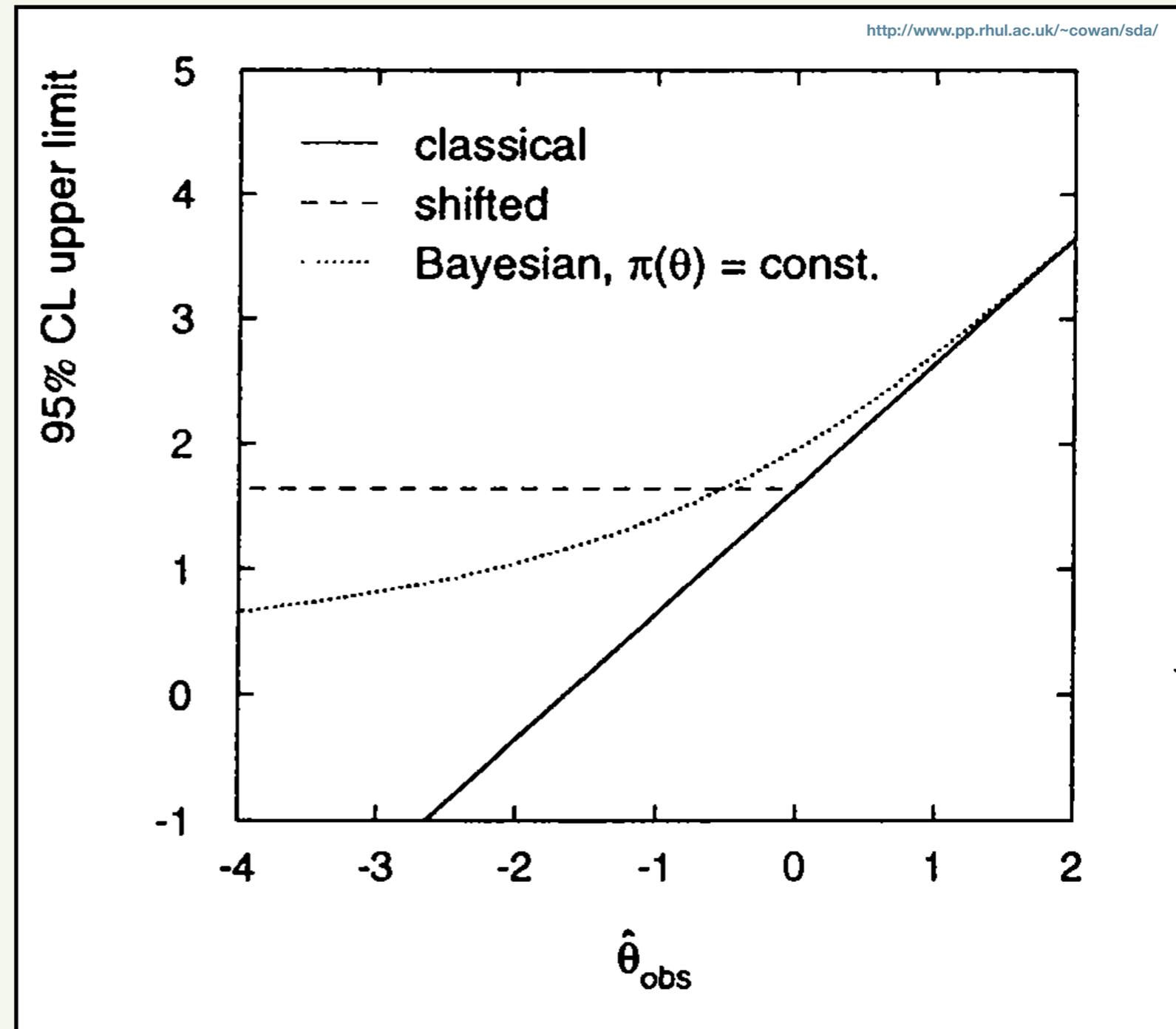
- It has been argued that in cases where  $\theta \geq 0$  *but no further information*, one should use

$$\pi(\theta) = \begin{cases} 0 & \theta \leq 0 \\ \frac{1}{\theta} & \theta > 0. \end{cases}$$

- This has the advantage that upper limits are invariant with respect to a transformation of the parameter by raising to an arbitrary power. This is equivalent to a uniform (improper) prior of previous form for  $\log \theta$ .
- For this to be usable, however, the **likelihood function** must go to **zero for  $\theta \rightarrow 0$  and  $\theta \rightarrow \infty$** , or else the **integrals diverge**. Thus this description is **often not applicable**.
- Therefore the uniform prior density (previous slide) is the most commonly used choice for setting limits on parameters.

# Different approaches compared

- Comparison of the three methods:
  - Classical and shifted are equal for  $\hat{\theta}_{\text{obs}} \geq 0$ ;
  - The Bayesian limit (here a constant prior is used) is always positive, and is always  $>$  the classical limit;
  - As the observed value grows, all limits approach each other.



# From one of our papers

PHYSICAL REVIEW D **101**, 032007 (2020)

## Search for $B^+ \rightarrow \mu^+ \nu_\mu$ and $B^+ \rightarrow \mu^+ N$ with inclusive tagging

M. T. Prim<sup>32</sup>, F. U. Bernlochner,<sup>3</sup> P. Goldenzweig,<sup>32</sup> M. Heck,<sup>32</sup> I. Adachi,<sup>18,15</sup> K. Adamczyk,<sup>62</sup> H. Aihara,<sup>85</sup> S. Al Said,<sup>79,34</sup> D. M. Asner,<sup>4</sup> H. Atmacan,<sup>76</sup> V. Aulchenko,<sup>5,65</sup> T. Aushev,<sup>54</sup> R. Ayad,<sup>79</sup> V. Babu,<sup>9</sup> A. M. Bakich,<sup>78</sup> V. Bansal,<sup>67</sup> P. Behera,<sup>25</sup> C. Beleño,<sup>14</sup> V. Bhardwaj,<sup>22</sup> B. Bhuyan,<sup>23</sup> T. Bilka,<sup>6</sup> J. Biswal,<sup>31</sup> A. Bobrov,<sup>5,65</sup> A. Bozek,<sup>62</sup> M. Bračko,<sup>48,31</sup> N. Braun,<sup>32</sup> T. E. Browder,<sup>17</sup> M. Campajola,<sup>29,57</sup> L. Cao,<sup>32</sup> D. Červenkov,<sup>6</sup> P. Chang,<sup>61</sup> V. Chekelian,<sup>49</sup> A. Chen,<sup>59</sup> B. G. Cheon,<sup>16</sup> K. Chilikin,<sup>42</sup> H. E. Cho,<sup>16</sup> K. Cho,<sup>36</sup> Y. Choi,<sup>77</sup> S. Choudhury,<sup>24</sup> D. Cinabro,<sup>89</sup> S. Cunliffe,<sup>9</sup> Z. Doležal,<sup>6</sup> S. Eidelman,<sup>5,65,42</sup> D. Epifanov,<sup>5,65</sup> J. E. Fast,<sup>67</sup> T. Ferber,<sup>9</sup> B. G. Fulsom,<sup>67</sup> R. Garg,<sup>68</sup> V. Gaur,<sup>88</sup> A. Garmash,<sup>5,65</sup> A. Giri,<sup>24</sup> O. Grzymkowska,<sup>62</sup> Y. Guan,<sup>8</sup> J. Haba,<sup>18,15</sup> T. Hara,<sup>18,15</sup> K. Hayasaka,<sup>64</sup> H. Hayashii,<sup>58</sup> W.-S. Hou,<sup>61</sup> T. Iijima,<sup>56,55</sup> K. Inami,<sup>55</sup> G. Inguglia,<sup>9</sup> A. Ishikawa,<sup>18</sup> M. Iwasaki,<sup>66</sup> Y. Iwasaki,<sup>18</sup> S. Jia,<sup>2</sup> Y. Jin,<sup>85</sup> D. Joffe,<sup>33</sup> K. K. Joo,<sup>7</sup> A. B. Kaliyar,<sup>25</sup> G. Karyan,<sup>9</sup> T. Kawasaki,<sup>35</sup> H. Kichimi,<sup>17</sup> C. Kiesling,<sup>49</sup> C. H. Kim,<sup>16</sup> D. Y. Kim,<sup>75</sup> K. T. Kim,<sup>37</sup> S. H. Kim,<sup>16</sup> K. Kinoshita,<sup>8</sup> P. Kodyš,<sup>6</sup> S. Korpar,<sup>48,31</sup> D. Kotchetkov,<sup>17</sup> P. Krizan,<sup>44,31</sup> R. Kroeger,<sup>51</sup> P. Krokovny,<sup>5,65</sup> T. Kuhr,<sup>45</sup> R. Kulasiri,<sup>33</sup> R. Kumar,<sup>70</sup> T. Kumita,<sup>87</sup> A. Kuzmin,<sup>5,65</sup> Y. J. Kwon,<sup>91</sup> K. Lalwani,<sup>47</sup> J. S. Lange,<sup>12</sup> I. S. Lee,<sup>16</sup> J. K. Lee,<sup>73</sup> J. Y. Lee,<sup>73</sup> S. C. Lee,<sup>39</sup> P. Lewis,<sup>17</sup> C. H. Li,<sup>43</sup> L. Li Gioi,<sup>49</sup> J. Libby,<sup>25</sup> K. Lieret,<sup>45</sup> D. Liventsev,<sup>88,18</sup> P.-C. Lu,<sup>61</sup> T. Luo,<sup>11</sup> J. MacNaughton,<sup>52</sup> M. Masuda,<sup>84</sup> T. Matsuda,<sup>52</sup> D. Matvienko,<sup>5,65,42</sup> M. Merola,<sup>29,57</sup> F. Metzner,<sup>32</sup> K. Miyabayashi,<sup>58</sup> R. Mizuk,<sup>42,54</sup> G. B. Mohanty,<sup>80</sup> R. Mussa,<sup>30</sup> M. Nakao,<sup>18,15</sup> G. De Nardo,<sup>29,57</sup> K. J. Nath,<sup>23</sup> Z. Natkaniec,<sup>62</sup> M. Nayak,<sup>89,18</sup> M. Niiyama,<sup>38</sup> N. K. Nisar,<sup>69</sup> S. Nishida,<sup>18,15</sup> K. Nishimura,<sup>17</sup> S. Ogawa,<sup>82</sup> H. Ono,<sup>63,64</sup> Y. Onuki,<sup>85</sup> P. Pakhlov,<sup>42,53</sup> G. Pakhlova,<sup>42,54</sup> B. Pal,<sup>4</sup> S. Pardi,<sup>29</sup> H. Park,<sup>39</sup> S.-H. Park,<sup>91</sup> S. Patra,<sup>22</sup> S. Paul,<sup>81</sup> T. K. Pedlar,<sup>46</sup> R. Pestotnik,<sup>31</sup> L. E. Piilonen,<sup>88</sup> V. Popov,<sup>42,54</sup> E. Prencipe,<sup>20</sup> M. Ritter,<sup>45</sup> A. Rostomyan,<sup>9</sup> M. Rozanska,<sup>62</sup> G. Russo,<sup>57</sup> D. Sahoo,<sup>80</sup> Y. Sakai,<sup>18,15</sup> L. Santelj,<sup>18</sup> V. Savinov,<sup>69</sup> O. Schneider,<sup>41</sup> G. Schnell,<sup>1,21</sup> J. Schueler,<sup>17</sup> C. Schwanda,<sup>27</sup> Y. Seino,<sup>64</sup> K. Senyo,<sup>90</sup> M. E. Sevir,<sup>50</sup> V. Shebalin,<sup>17</sup> J.-G. Shiu,<sup>61</sup> B. Shwartz,<sup>5,65</sup> F. Simon,<sup>49</sup> A. Sokolov,<sup>28</sup> E. Solovieva,<sup>42</sup> M. Starič,<sup>31</sup> J. F. Strube,<sup>67</sup> M. Sumihama,<sup>13</sup> T. Sumiyoshi,<sup>87</sup> W. Sutcliffe,<sup>32</sup> M. Takizawa,<sup>74,19,71</sup> U. Tamponi,<sup>30</sup> Y. Tao,<sup>10</sup> F. Tenchini,<sup>9</sup> K. Trabelsi,<sup>40</sup> M. Uchida,<sup>86</sup> T. Uglov,<sup>42,54</sup> Y. Unno,<sup>16</sup> S. Uno,<sup>18,15</sup> Y. Ushiroda,<sup>18,15</sup> Y. Usov,<sup>5,65</sup> S. E. Vahsen,<sup>17</sup> R. Van Tonder,<sup>32</sup> G. Varner,<sup>17</sup> K. E. Varvell,<sup>78</sup> A. Vinokurova,<sup>5,65</sup> B. Wang,<sup>49</sup> C. H. Wang,<sup>60</sup> M.-Z. Wang,<sup>61</sup> P. Wang,<sup>26</sup> S. Watanuki,<sup>83</sup> E. Won,<sup>37</sup> S. B. Yang,<sup>37</sup> H. Ye,<sup>9</sup> J. H. Yin,<sup>26</sup> Y. Yusa,<sup>64</sup> Z. P. Zhang,<sup>72</sup> V. Zhilich,<sup>5,65</sup> V. Zhukova,<sup>42</sup> and V. Zhulanov<sup>5,65</sup>

(Belle Collaboration)

<sup>1</sup>University of the Basque Country UPV/EHU, 48080 Bilbao

<sup>2</sup>Beihang University, Beijing 100191

<sup>3</sup>University of Bonn, 53115 Bonn

<sup>4</sup>Brookhaven National Laboratory, Upton, New York 11973

<sup>5</sup>Budker Institute of Nuclear Physics SB RAS, Novosibirsk 630090

<sup>6</sup>Faculty of Mathematics and Physics, Charles University, 121 16 Prague

<sup>7</sup>Chonnam National University, Kwangju 660-701

<sup>8</sup>University of Cincinnati, Cincinnati, Ohio 45221

<sup>9</sup>Deutsches Elektronen-Synchrotron, 22607 Hamburg

<sup>10</sup>University of Florida, Gainesville, Florida 32611

<sup>11</sup>Key Laboratory of Nuclear Physics and Ion-beam Application (MOE) and Institute of Modern Physics, Fudan University, Shanghai 200443

<sup>12</sup>Justus-Liebig-Universität Gießen, 35392 Gießen

<sup>13</sup>Gifu University, Gifu 501-1193

<sup>14</sup>II. Physikalisches Institut, Georg-August-Universität Göttingen, 37073 Göttingen

<sup>15</sup>SOKENDAI (The Graduate University for Advanced Studies), Hayama 240-0193

<sup>16</sup>Hanyang University, Seoul 133-791

<sup>17</sup>University of Hawaii, Honolulu, Hawaii 96822

<sup>18</sup>High Energy Accelerator Research Organization (KEK), Tsukuba 305-0801

<sup>19</sup>J-PARC Branch, KEK Theory Center, High Energy Accelerator Research Organization (KEK), Tsukuba 305-0801

<sup>20</sup>Forschungszentrum Jülich, 52425 Jülich

<sup>21</sup>IKERBASQUE, Basque Foundation for Science, 48013 Bilbao

<sup>22</sup>Indian Institute of Science Education and Research Mohali, SAS Nagar, 140306

<sup>23</sup>Indian Institute of Technology Guwahati, Assam 781039

<sup>24</sup>Indian Institute of Technology Hyderabad, Telangana 502285

<sup>25</sup>Indian Institute of Technology Madras, Chennai 600036

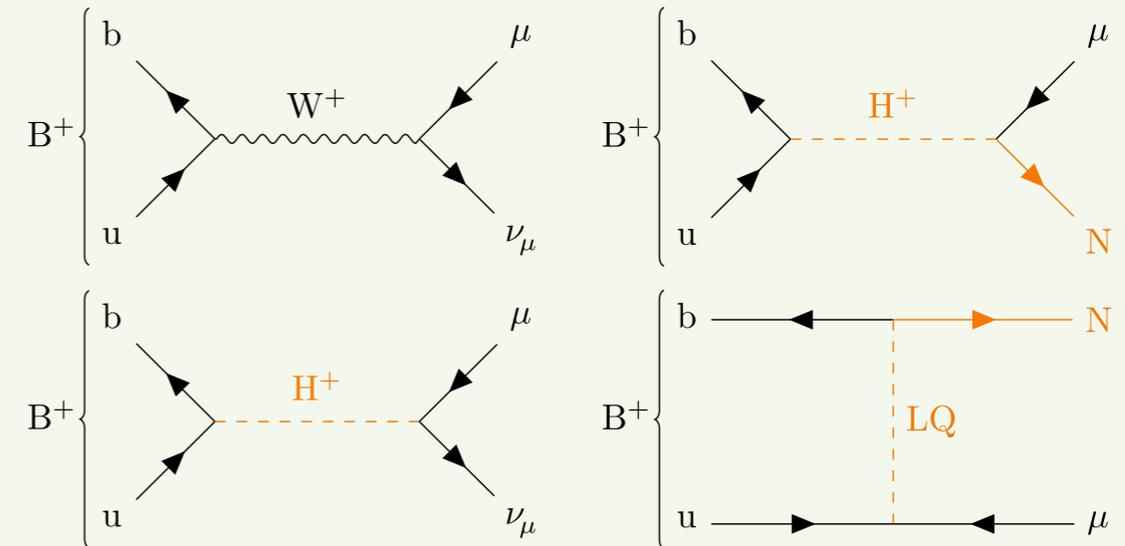
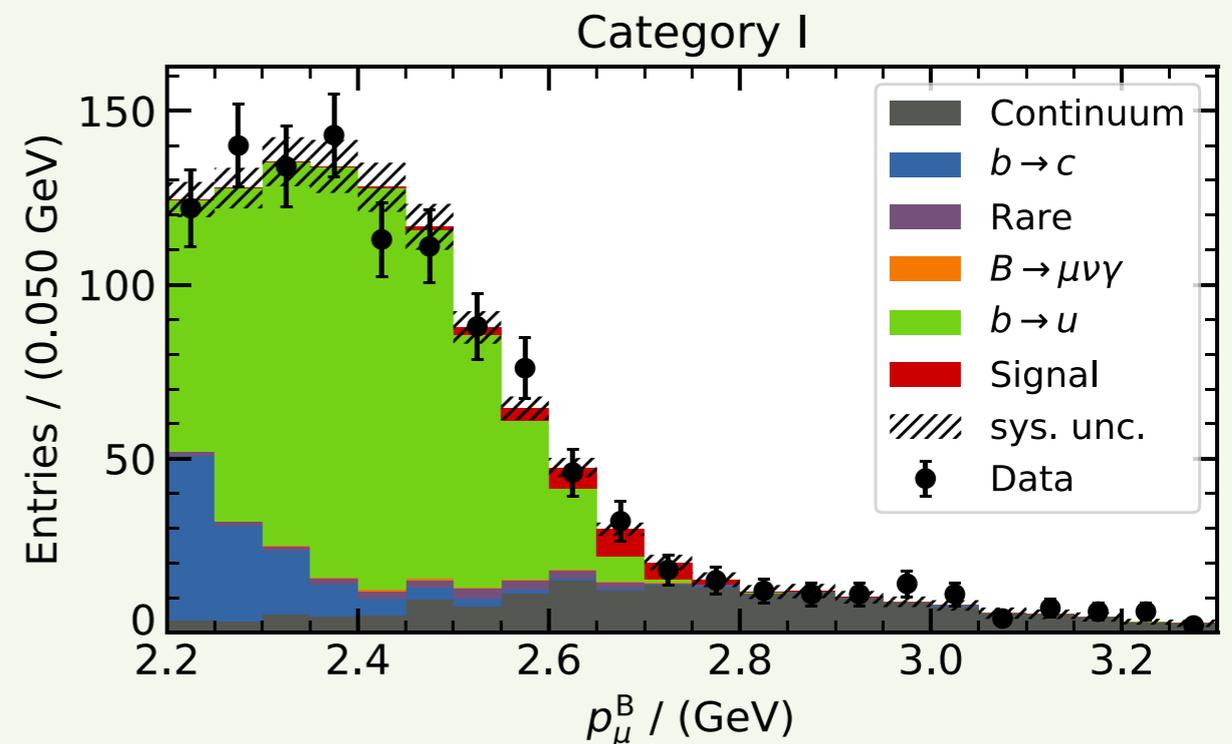
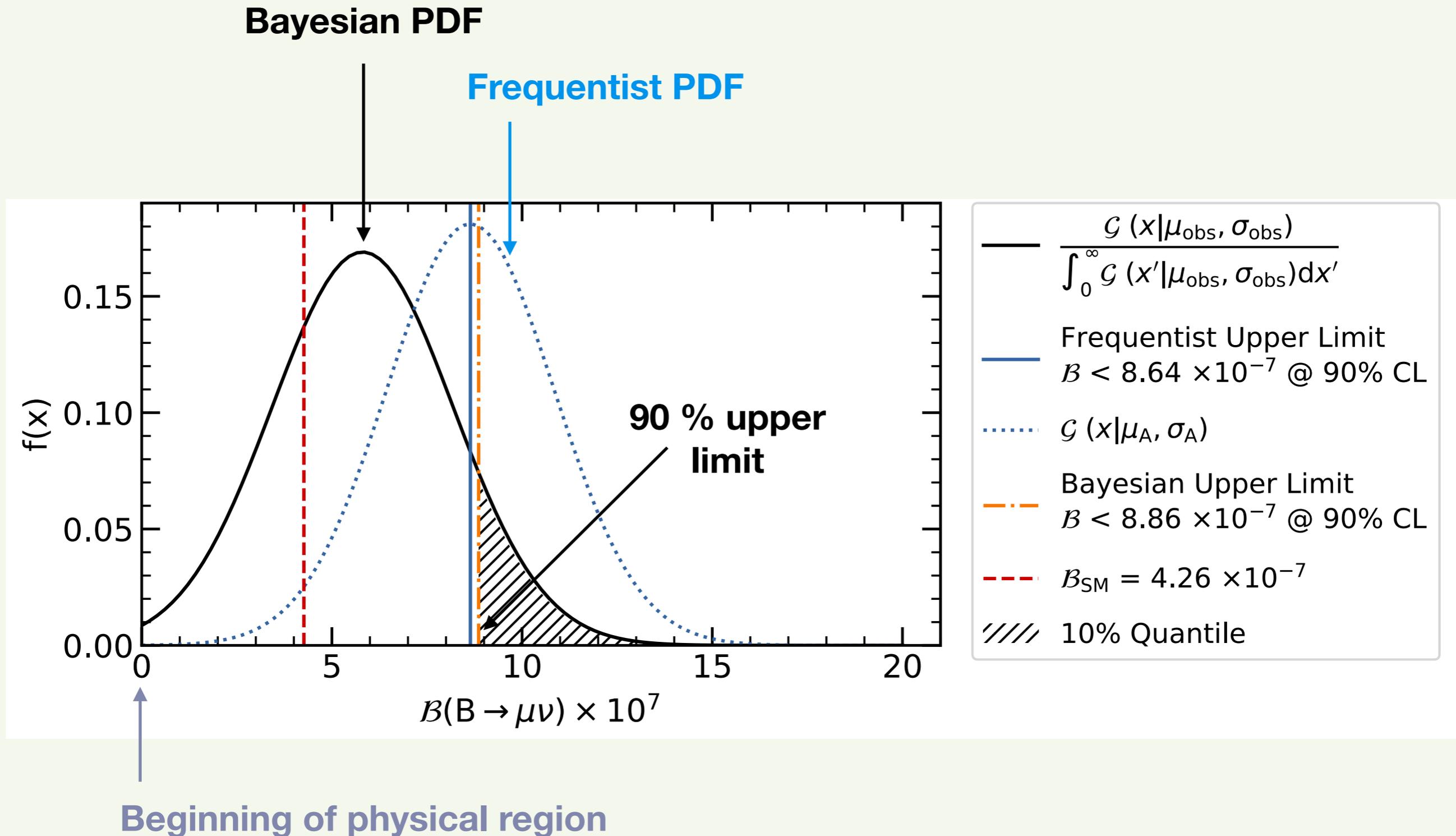


FIG. 1. The SM leptonic  $B^+ \rightarrow \mu^+ \nu_\mu$  decay process and possible BSM processes with and without a sterile neutrino  $N$  in the final state are shown.



# From one of our papers



# Upper limit on the mean of Poisson variable with background



Read at home

# UL of mean of Poisson variable with bkg

- Recall from last lecture the UL we placed on the mean  $\nu$  of a Poisson variable  $n$ . *(Last week we considered signal only though.)*
- Often one faces a somewhat more complicated situation, where the observed value of  $n$  is the sum of the desired **signal**  $n_s$ , as well as the **background events**  $n_b$ ,
  - $n = n_s + n_b$  where both  $n_s$  and  $n_b$  can be regarded as Poisson variables with **means**  $\nu_s$  and  $\nu_b$ , respectively.
  - Suppose for the moment, that the **mean of the background**  $\nu_b$  is **known without any uncertainty**.
  - For  $\nu_s$  one only **knows** a priori that  $\nu_s \geq 0$ .
- Our goal is to construct an UL for the signal parameter  $\nu_s$  given a measured value of  $n$ .

# Upper Limit

- Since  $n$  is the sum of two Poisson variables, one can show that it itself is a Poisson variable with the probability function

$$f(n; \nu_s, \nu_b) = \frac{(\nu_s + \nu_b)^n}{n!} e^{-(\nu_s + \nu_b)}.$$

- The ML estimator for  $\nu_s$  is

$$\hat{\nu}_s = n - \nu_b,$$

- It has zero bias since  $E[n] = \nu_s + \nu_b$ ,

- The equations determining the confidence interval become

$$\alpha = P(\hat{\nu}_s \geq \hat{\nu}_s^{\text{obs}}; \nu_s^{\text{lo}}) = \sum_{n \geq n_{\text{obs}}} \frac{(\nu_s^{\text{lo}} + \nu_b)^n e^{-(\nu_s^{\text{lo}} + \nu_b)}}{n!},$$
$$\beta = P(\hat{\nu}_s \leq \hat{\nu}_s^{\text{obs}}; \nu_s^{\text{up}}) = \sum_{n \leq n_{\text{obs}}} \frac{(\nu_s^{\text{up}} + \nu_b)^n e^{-(\nu_s^{\text{up}} + \nu_b)}}{n!}.$$

Can solve numerically for  $\nu_s^{\text{lo}}$  and  $\nu_s^{\text{up}}$

# Comparison with no background result

- **Comparing to our previous expressions, we see that these limits are related to the ones without background by**

$$\begin{aligned}\nu_s^{\text{lo}} &= \nu_s^{\text{lo}}(\text{no background}) - \nu_b, \\ \nu_s^{\text{up}} &= \nu_s^{\text{up}}(\text{no background}) - \nu_b.\end{aligned}$$

- The **difficulties** that can arise here are **similar to the example without background**, i.e. when the total number of events observed is not large compared to the expected number of background events.
  - **Because of these difficulties, the classical limit often causes problems**
    - *As previously mentioned, one should always report  $\hat{\nu}_s$  and an estimate for its variance to allow for meaningful combinations later*

# Bayesian Limit

- The **Bayesian method** can be **used here** as well, with for example a uniform prior. The **likelihood function** and **posterior probability** are given by

$$L(n_{\text{obs}}|\nu_s) = \frac{(\nu_s + \nu_b)^{n_{\text{obs}}}}{n_{\text{obs}}!} e^{-(\nu_s + \nu_b)}.$$

$$p(\nu_s|n_{\text{obs}}) = \frac{L(n_{\text{obs}}|\nu_s) \pi(\nu_s)}{\int_{-\infty}^{\infty} L(n_{\text{obs}}|\nu'_s) \pi(\nu'_s) d\nu'_s}.$$

- Taking  $\pi(\nu_s) = \text{const. for } \nu_s > 0$  and **zero otherwise**, the **upper limit**  $\nu_s^{\text{up}}$  at **CL**  $1 - \beta$  is

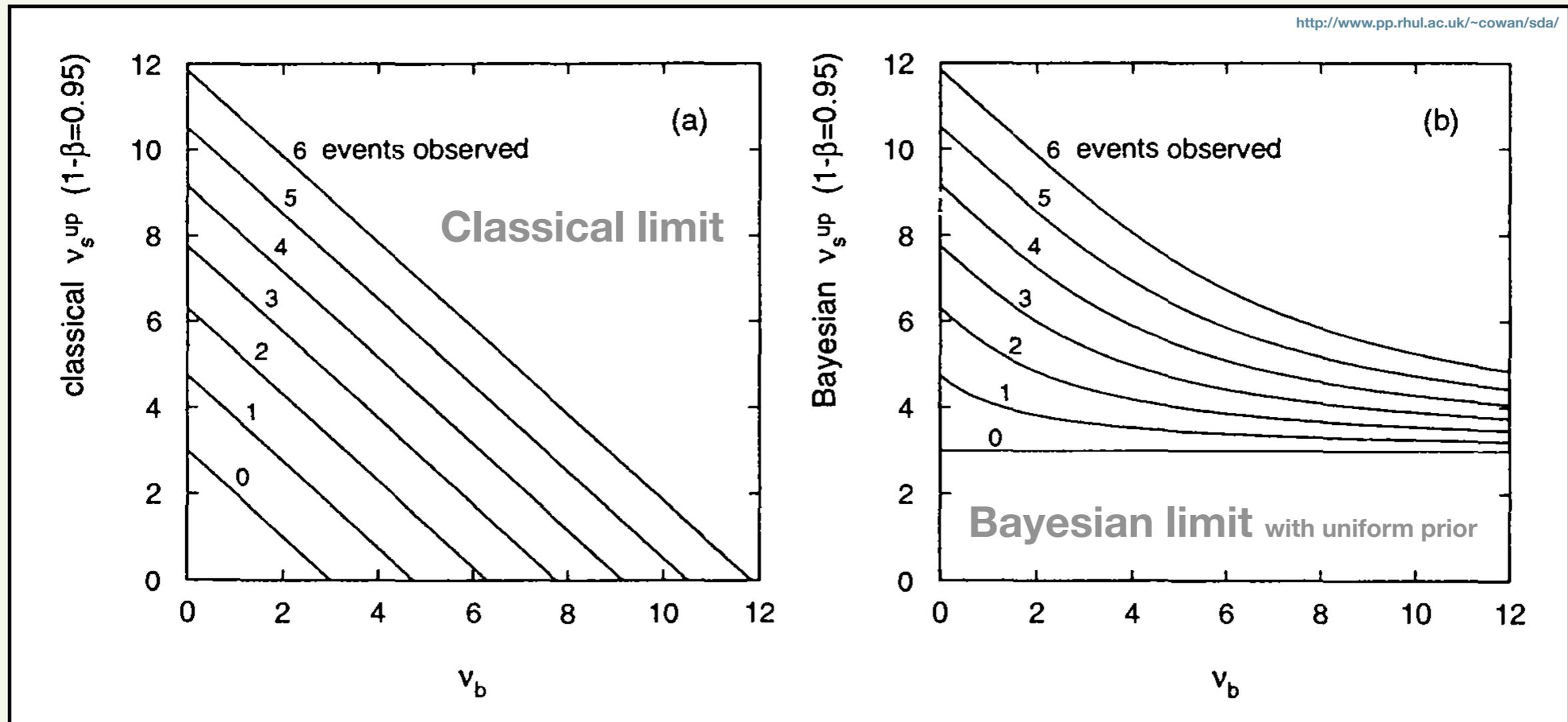
$$\begin{aligned} 1 - \beta &= \frac{\int_0^{\nu_s^{\text{up}}} L(n_{\text{obs}}|\nu_s) d\nu_s}{\int_0^{\infty} L(n_{\text{obs}}|\nu_s) d\nu_s} \\ &= \frac{\int_0^{\nu_s^{\text{up}}} (\nu_s + \nu_b)^{n_{\text{obs}}} e^{-(\nu_s + \nu_b)} d\nu_s}{\int_0^{\infty} (\nu_s + \nu_b)^{n_{\text{obs}}} e^{-(\nu_s + \nu_b)} d\nu_s}. \end{aligned}$$

Integrals can be related to incomplete gamma functions and one obtains:

$$\beta = \frac{e^{-(\nu_s^{\text{up}} + \nu_b)} \sum_{n=0}^{n_{\text{obs}}} \frac{(\nu_s^{\text{up}} + \nu_b)^n}{n!}}{e^{-\nu_b} \sum_{n=0}^{n_{\text{obs}}} \frac{\nu_b^n}{n!}}.$$

# Bayesian Limit

- Upper limits at  $CL\ 1 - \beta = 0.95$  for different number of observed events and as a function of the expected number of background events.



# More realistic scenario

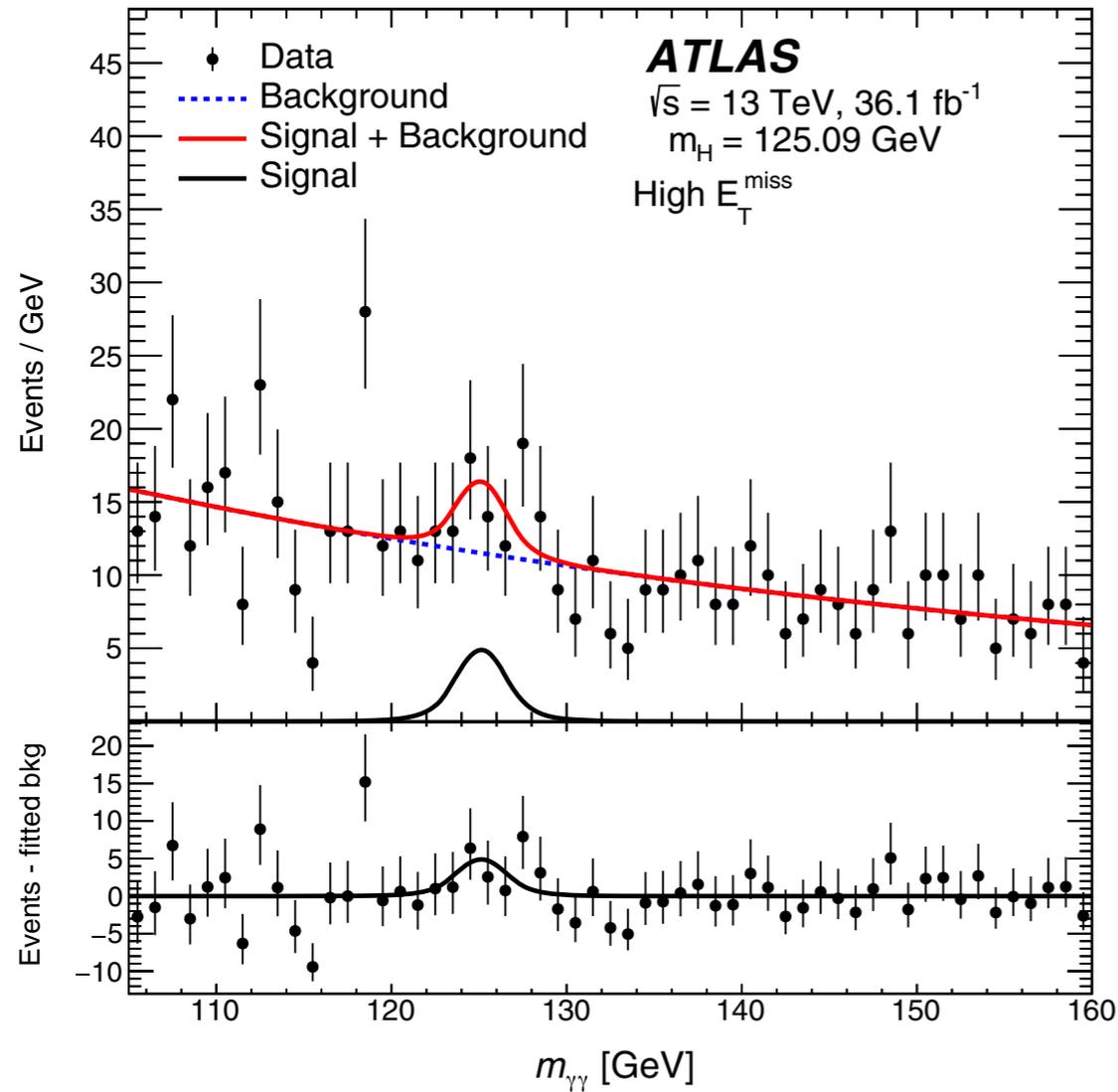
- Often the result of an experiment is not simply the number of  $n$  observed events, but includes in addition measured values  $x_1, x_2, \dots, x_n$  of some property of the events (e.g. mass).
- Suppose the probability density for  $x$  is

$$f(x; \nu_s, \nu_b) = \frac{\nu_s f_s(x) + \nu_b f_b(x)}{\nu_s + \nu_b},$$

- This information can be incorporated into the limit  $\nu_s$  by using the extended likelihood function

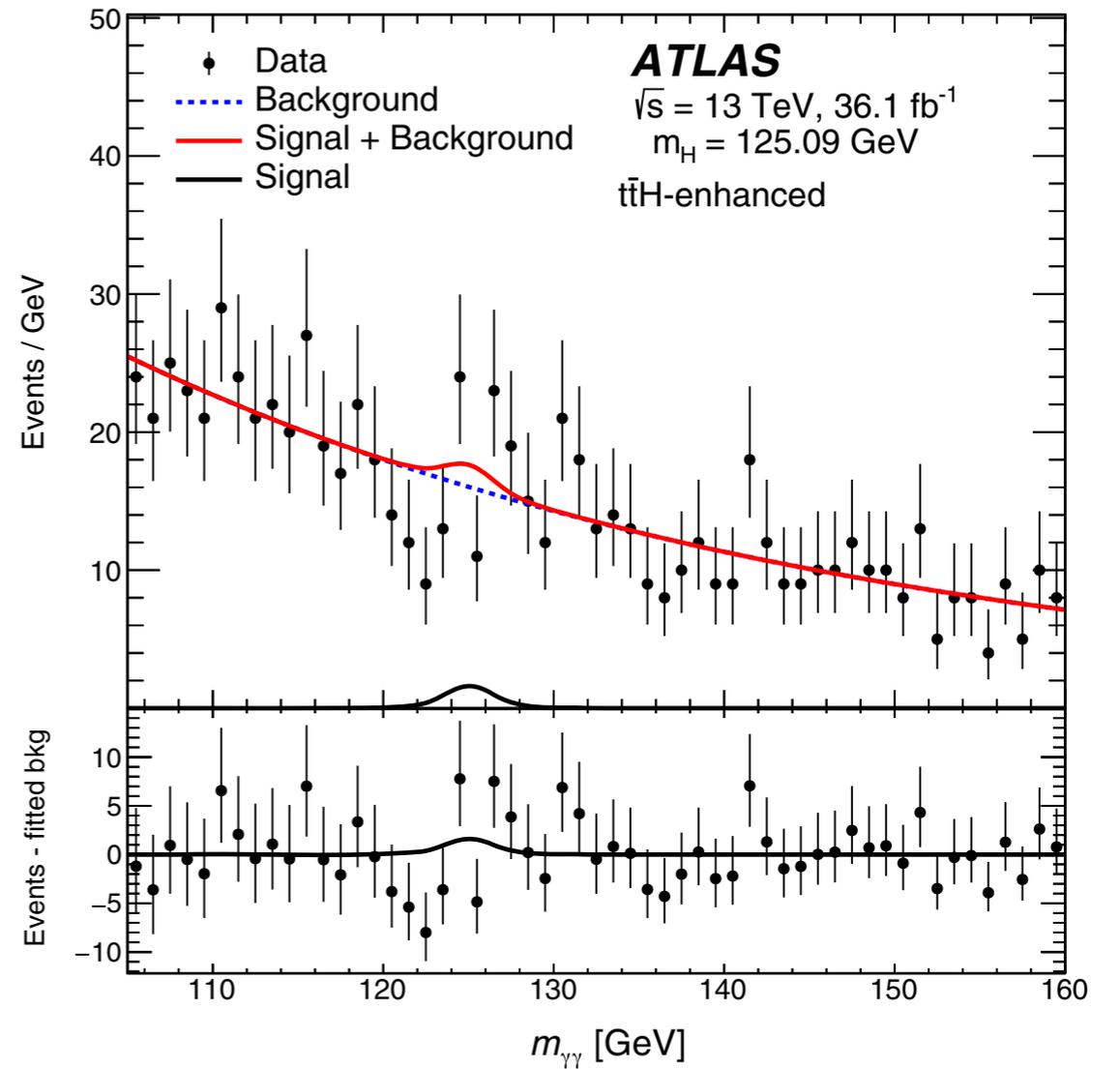
$$\begin{aligned} L(\nu_s) &= \frac{(\nu_s + \nu_b)^n}{n!} e^{-(\nu_s + \nu_b)} \prod_{i=1}^n \frac{\nu_s f_s(x_i) + \nu_b f_b(x_i)}{\nu_s + \nu_b} \\ &= \frac{e^{-(\nu_s + \nu_b)}}{n!} \prod_{i=1}^n [\nu_s f_s(x_i) + \nu_b f_b(x_i)], \end{aligned}$$

→ *limits in general must be determined numerically or via MC methods*



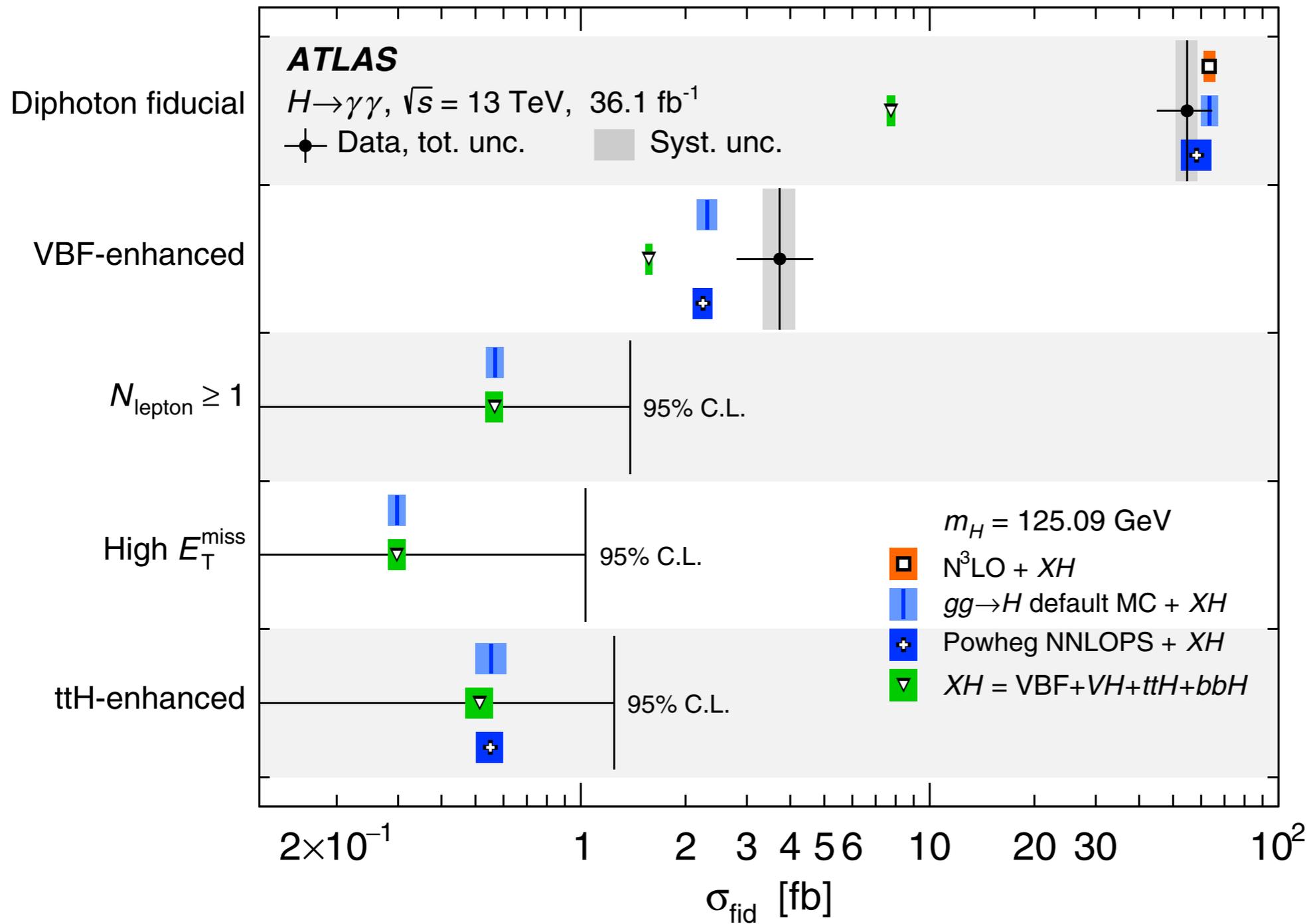
(c) High  $E_T^{\text{miss}}$

- (4) High  $E_T^{\text{miss}}$ : This region retains events with missing transverse momentum  $E_T^{\text{miss}} > 80$  GeV and  $p_T^{\gamma\gamma} > 80$  GeV is defined to study  $VH$  production and possible contributions of Higgs boson production with dark matter particles. The simultaneous requirement that the Higgs boson system balances the missing transverse momentum reduces the fraction of selected events at detector level without particle-level  $E_T^{\text{miss}} > 80$  GeV.



(d)  $t\bar{t}H$ -enhanced

- (5)  $t\bar{t}H$ -enhanced: This region retains events with either at least one lepton and three jets or no leptons and four jets to study Higgs boson production in association with top quarks. In addition, one of the jets needs to be identified as originating from a bottom quark.



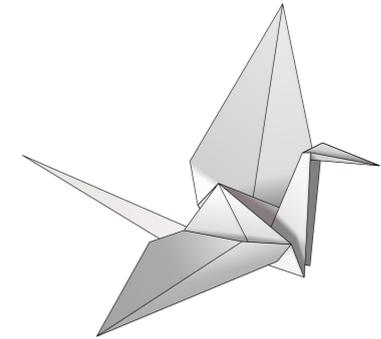
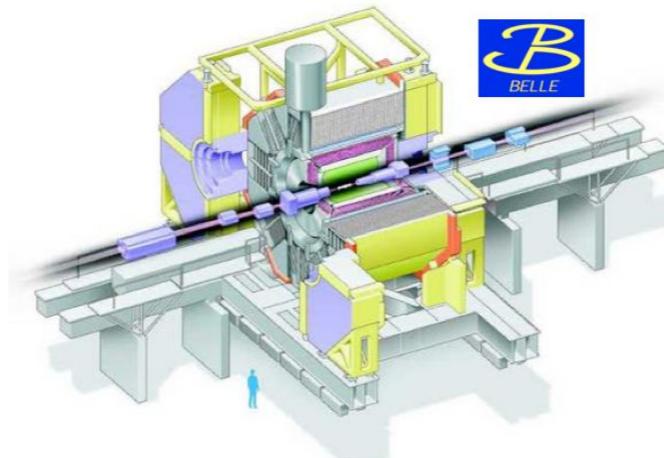
$$\theta_{\text{up}} = \hat{\theta}_{\text{obs}} + \sigma_{\hat{\theta}} \Phi^{-1}(1 - \beta) \quad \Phi^{-1}(0.95) = 1.645$$

# Take 5



# Unfolding

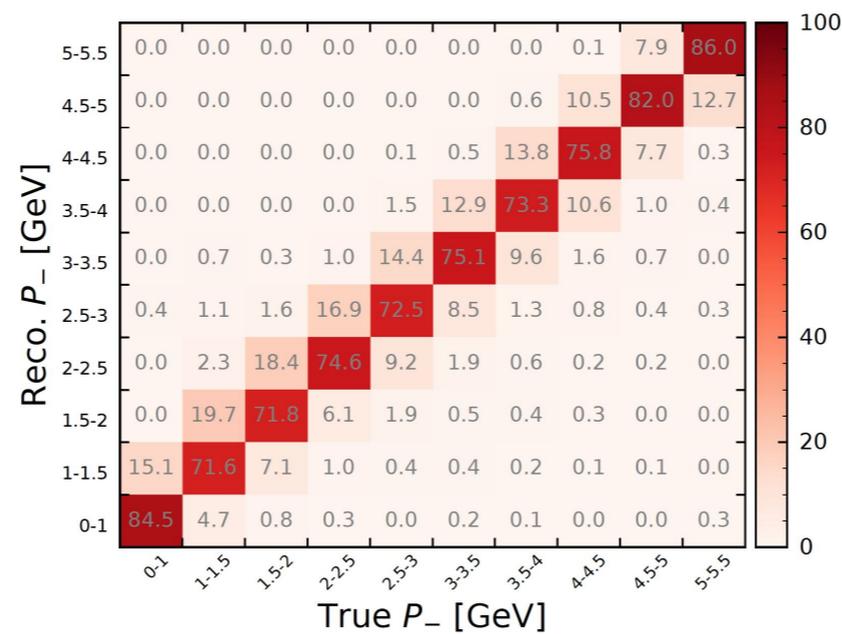
# In a nutshell



**$Y$** : true distribution

**$R$** : detector response

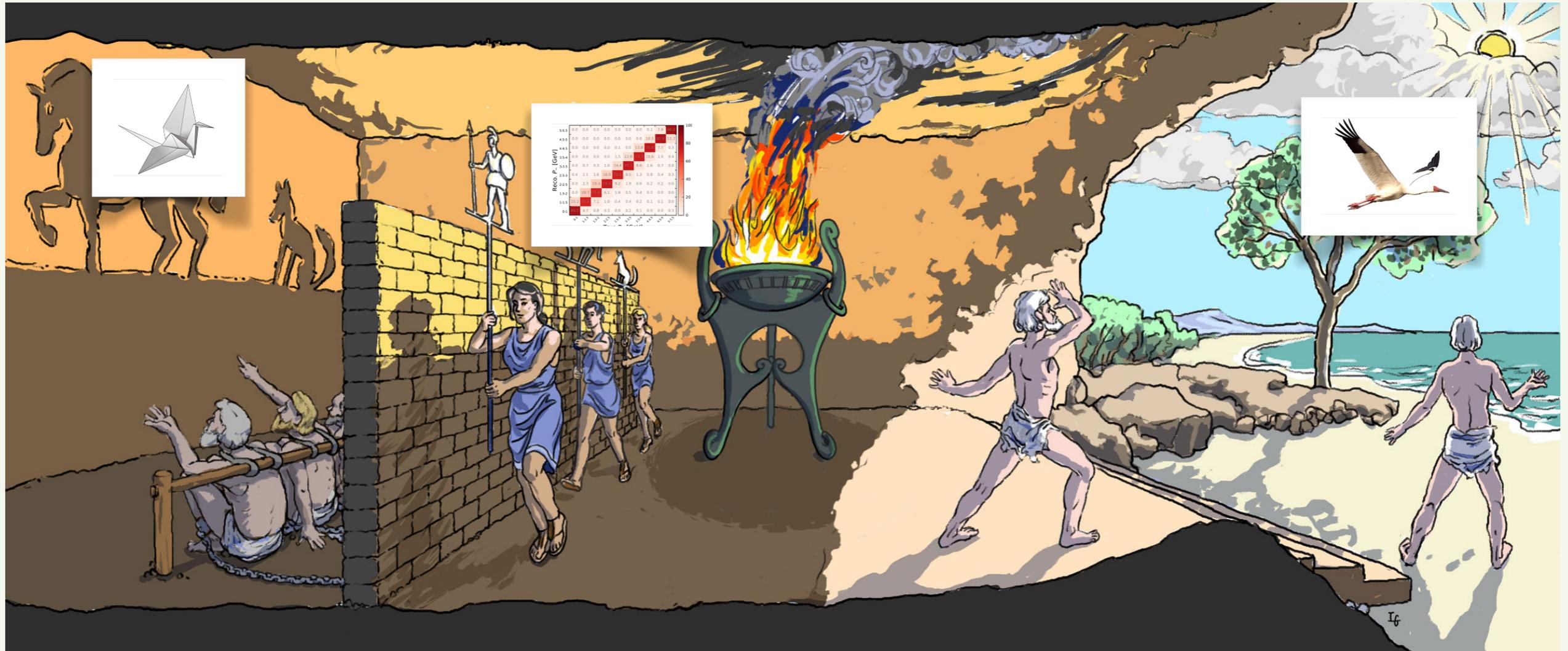
**$X$** : measured distribution



The detector response is represented by a migration (response) matrix  $R$ .

$R(i, j)$  indicates the probability to observe an event in bin  $i$  if it had generator-level value in bin  $j$ .

# Allegory of The Cave (Plato's Republic)



By 4edges - Own work, CC BY-SA 4.0, <https://commons.wikimedia.org/w/index.php?curid=73850232>

# The unfolding problem

- Up till now:
  - Have considered that RVs such as particle energies, decay times, etc., can be measured with **absolute precision**.
- In reality:
  - Every experimental apparatus has **finite resolution**.
    - This distorts measurements.
      - Correct for this = **Unfolding**

# Derivation (i)

- $f_{\text{true}}(y)$  = PDF of true value 'y'

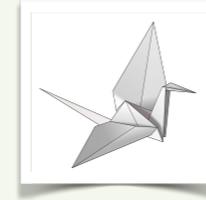


- To construct a usable estimator for  $f_{\text{true}}(y)$ , must represent it by means of some finite set of parameters.
- If no functional form for  $f_{\text{true}}(y)$  is known a priori, then it can still be represented as a normalized histogram with  $M$  bins.

$$p_j = \int_{\text{bin } j} f_{\text{true}}(y) dy \quad \text{is the probability to find } y \text{ in bin } j$$

- $\mu_{\text{tot}}$  = expectation value of total # of events.
  - $\mu_j = \mu_{\text{tot}} p_j$  is the expected # of events in bin  $j$
  - The vector  $\boldsymbol{\mu} = (\mu_1, \mu_2, \dots, \mu_M)$  is the 'true histogram'
    - **Careful**: not the actual number of events in each bin, but the *expectation values*

# Derivation (ii)



- Begin with a sample of measured values  $x$ 
  - Entered into a histogram of  $N$  bins:  $\mathbf{n} = (n_1, n_2, \dots, n_N)$ 
    - The # of bins  $N$  can be  $>$ ,  $<$ ,  $=$  to the # of bins in the true histogram  $M$
- Regard  $n_i$  as independent Poisson variables with expectation value  $\nu_i$

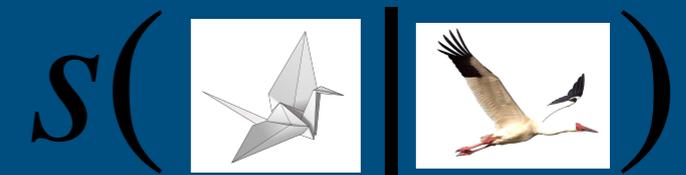
$$P(n_i; \nu_i) = \frac{\nu_i^{n_i} e^{-\nu_i}}{n_i!} \quad \nu_i = E[n_i]$$

- From the law of total probability:

$$\begin{aligned} \nu_i &= \mu_{\text{tot}} P(\text{event observed in bin } i) \\ &= \mu_{\text{tot}} \int dy P \left( \begin{array}{c} \text{observed} \\ \text{in bin } i \end{array} \middle| \begin{array}{c} \text{true value } y \text{ and} \\ \text{event detected} \end{array} \right) \varepsilon(y) f_{\text{true}}(y) \\ &= \mu_{\text{tot}} \int_{\text{bin } i} dx \int dy s(x|y) \varepsilon(y) f_{\text{true}}(y). \end{aligned}$$

Resolution function  
(point spread function in  
imaging applications)

Detection efficiency =  
the probability that an  
event leads to some  
measured value



- The **resolution function**  $s(x | y)$  is a conditional PDF:
  - For the measured value  $x$ , given the true value  $y$ 
    - Probability that an even leads to some measured value
- Sometimes also incorporates the **detection efficiency**  $\epsilon(y)$ 
  - $r(x | y) = s(x | y) \epsilon(y)$ 
    - $\uparrow$  **Response function:** includes the effect of limited efficiency
- One says that the true distribution is **folded** with the response function
  - i.e., expressing  $\nu_i$  as a function of  $s(x | y) = \text{folding}$
- **Unfolding** = the task of estimating  $f_{\text{true}}$

# Derivation (iv)

True:  $j, y, \mu, M$   
Measured:  $i, x, \nu, N$

- Take our integral for  $\nu_i$  (from slide 48)

$$= \mu_{\text{tot}} \int_{\text{bin } i} dx \int dy s(x|y) \varepsilon(y) f_{\text{true}}(y)$$



Break up the integral over  $y$  into a sum over bins  $j$   
Multiply numerator and denominator by  $\mu_j$

$$= \sum_{j=1}^M \frac{\int_{\text{bin } i} dx \int_{\text{bin } j} dy s(x|y) \varepsilon(y) f_{\text{true}}(y)}{(\mu_j / \mu_{\text{tot}})} \mu_j$$

$P(\text{observed in bin } i \text{ and true value in bin } j)$

$P(\text{true value in bin } j)$

$$= \sum_{j=1}^M R_{ij} \mu_j,$$

Response matrix =

The conditional probability that an event will be found with measured value  $x$  in bin  $i$ , given that the true value  $y$  was in bin  $j$



# Derivation (iv)

**True:**  $j, y, \mu, M$   
**Measured:**  $i, x, \nu, N$

- Take our integral for  $\nu_i$  (from slide 48)

$$= \mu_{\text{tot}} \int_{\text{bin } i} dx \int dy s(x|y) \varepsilon(y) f_{\text{true}}(y)$$



Break up the integral over  $y$  into a sum over bins  $j$   
 Multiply numerator and denominator by  $\mu_j$

$$= \sum_{j=1}^M \frac{\int_{\text{bin } i} dx \int_{\text{bin } j} dy s(x|y) \varepsilon(y) f_{\text{true}}(y)}{(\mu_j / \mu_{\text{tot}})} \mu_j$$

$$= \sum_{j=1}^M R_{ij} \mu_j,$$

Response matrix =

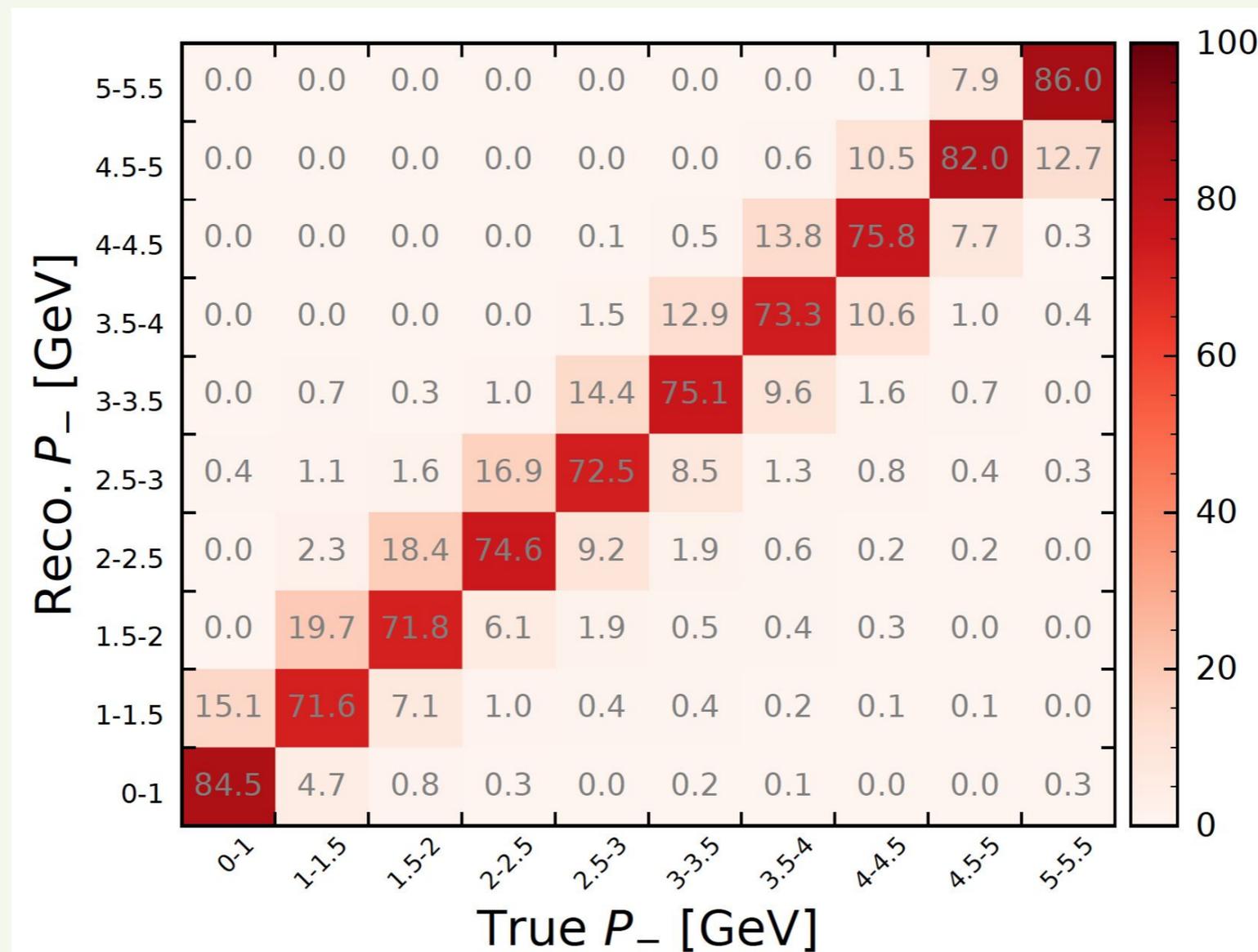
The conditional probability that an event will be found with measured value  $x$  in bin  $i$ , given that the true value  $y$  was in bin  $j$

$$\frac{P(\text{observed in bin } i \text{ and true value in bin } j)}{P(\text{true value in bin } j)} = P(\text{observed in bin } i \mid \text{true value in bin } j)$$



# Response matrix

- The effect of the off-diagonal elements in  $R$  is to smear out any fine structure
  - A peak in the true histogram concentrated mainly in 1 bin will be observed over several bins
  - 2 peaks separated by less than several bins will be merged into a single broad peak



*$R$  doesn't need to be symmetric*

# Efficiencies

**True:**  $j, y, \mu, M$   
**Measured:**  $i, x, \nu, N$

- Sum over the 'measured' index  $i$  and use  $\int s(x|y)dx = 1$

$$\begin{aligned}\sum_{i=1}^N R_{ij} &= \sum_{i=1}^N \frac{\int_{\text{bin } i} dx \int_{\text{bin } j} dy s(x|y) \varepsilon(y) f_{\text{true}}(y)}{(\mu_j / \mu_{\text{tot}})} \\ &= \frac{\int_{\text{bin } j} dy \varepsilon(y) f_{\text{true}}(y)}{\int_{\text{bin } j} f_{\text{true}}(y) dy} \\ &\equiv \varepsilon_j,\end{aligned}$$



The average value of the efficiency over bin  $j$

# Include background

- In addition to limited resolution and efficiency, must allow for the possibility of **background** processes
  - Measuring device produces a value when no true event of the type under study occurred
    - E.g., for  $\beta$ -decay, background = spurious signals in the detector, the presence of other radioactive nuclei in the sample, interactions due to cosmic rays, etc.

$$\nu_i = \sum_{j=1}^m R_{ij} \mu_j + \beta_i$$

↑  
The # of entries in bin  $i$  which originate from **background** processes

The uncertainty from the background is a source of **systematic error** in the unfolded result

# To summarize:

**True:**  $j, y, \mu, M$   
**Measured:**  $i, x, \nu, N$

- The vector  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_M)$  is the ‘true histogram’  $\rightarrow$  *Expectation values of true # of entries in each bin*
- The normalized true histogram  $\boldsymbol{p} = (p_1, \dots, p_M) = \boldsymbol{\mu} / \mu_{\text{tot}}$   $\rightarrow$  *Probabilities*
- The expectation values of the observed # of entries  $\boldsymbol{\nu} = (\nu_1, \dots, \nu_N)$
- The actual # of entries observed  $\boldsymbol{n} = (n_1, \dots, n_N)$   $\rightarrow$  *The data*
- Efficiencies  $\boldsymbol{\epsilon} = (\epsilon_1, \dots, \epsilon_M)$
- Expected background values  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_N)$
- Response matrix  $R_{ij}$ 
  - $i = 1, \dots, N$  represents the bin of the observed histogram
  - $j = 1, \dots, M$  gives the bin of the true histogram

## Related by

$$\boldsymbol{\nu} = R\boldsymbol{\mu} + \boldsymbol{\beta}$$

$\uparrow$  *Assume known*  $\uparrow$

## Goal

Construct estimators  $\hat{\boldsymbol{\mu}}$  for the true histogram, or estimators  $\hat{\boldsymbol{p}}$  for the probabilities

## Assume we either:

Know the form of the probability distribution for the data  $\boldsymbol{n}$   $\Rightarrow$  *Allow us to construct the  $\mathcal{L}$  function*

Have the covariance matrix  $V_{ij} = \text{cov}[n_i, n_j]$   $\Rightarrow$  *Used to construct a  $\chi^2$  function*

# Method 1: Invert the response matrix

True:  $j, y, \mu, M$   
Measured:  $i, x, \nu, N$

- Start with the matrix form (with  $M = N$ )

$$\boldsymbol{\nu} = R\boldsymbol{\mu} + \boldsymbol{\beta}$$

- Invert it to obtain

$$\boldsymbol{\mu} = R^{-1}(\boldsymbol{\nu} - \boldsymbol{\beta})$$

- Set the estimators for  $\boldsymbol{\nu}$  to be the data values  $\boldsymbol{n}$

$$\hat{\boldsymbol{\nu}} = \boldsymbol{n}$$

- The estimators for the  $\boldsymbol{\mu}$  are then

$$\hat{\boldsymbol{\mu}} = R^{-1}(\boldsymbol{n} - \boldsymbol{\beta})$$

# Properties

True:  $j, y, \mu, M$   
Measured:  $i, x, \nu, N$

- Expectation value of inversion:

$$\begin{aligned} E[\hat{\mu}_j] &= \sum_{i=1}^N (R^{-1})_{ji} E[n_i - \beta_i] = \sum_{i=1}^N (R^{-1})_{ji} (\nu_i - \beta_i) \\ &= \mu_j, \end{aligned}$$

Estimators  $\hat{\mu}_j$  are unbiased  
(Since by assumption  $\hat{\nu}_i = n_i$  is unbiased)

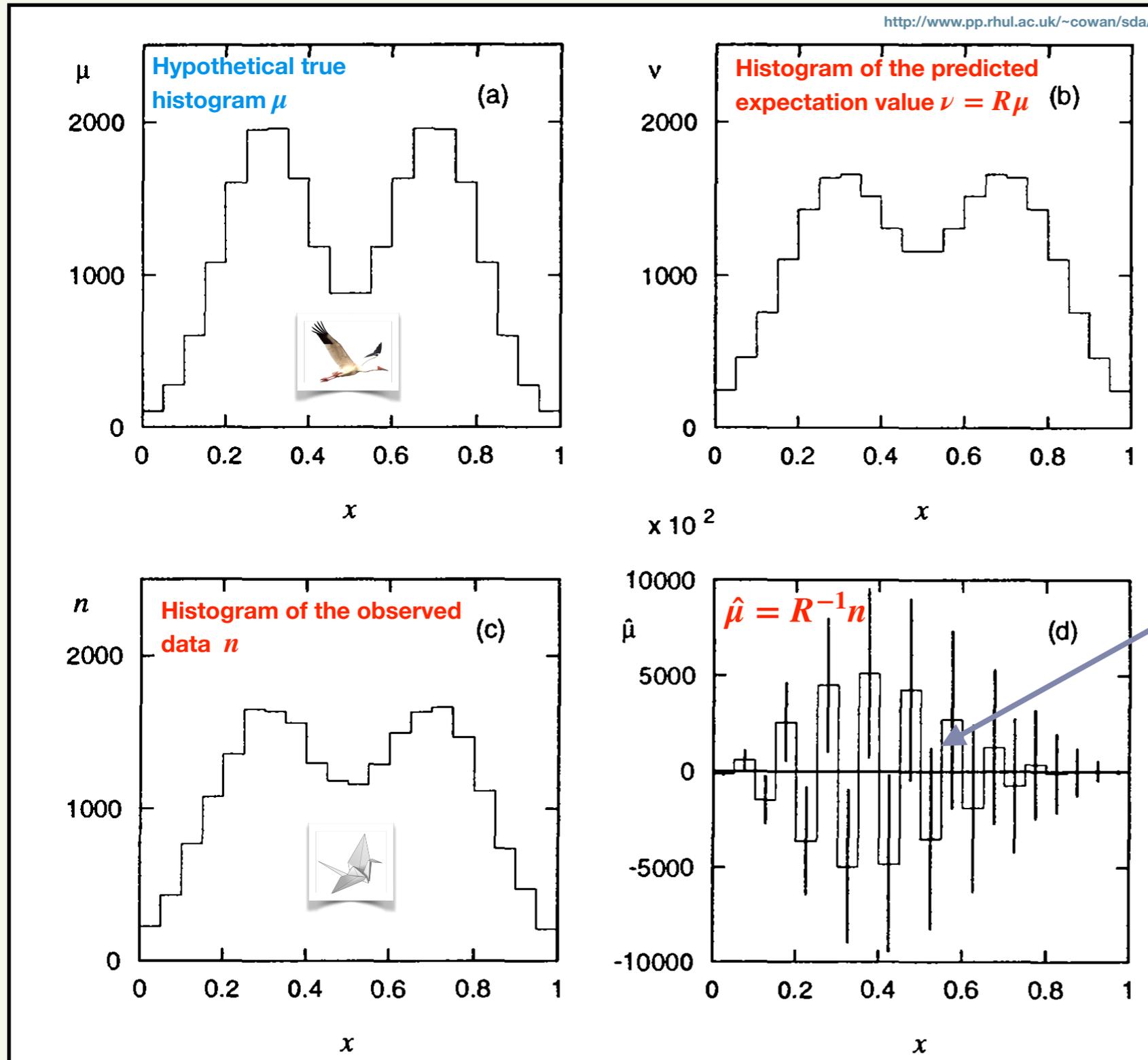
- Covariance of uncorrelated Poisson variables:

$$\begin{aligned} \text{cov}[\hat{\mu}_i, \hat{\mu}_j] &= \sum_{k,l=1}^N (R^{-1})_{ik} (R^{-1})_{jl} \text{cov}[n_k, n_l] \\ &= \sum_{k=1}^N (R^{-1})_{ik} (R^{-1})_{jk} \nu_k, \end{aligned}$$

- Covariance of correlated Gaussian variables:

$$U = R^{-1} V (R^{-1})^T.$$

# Ex. where matrix inversion goes “wrong”



Assume 0 background,  $\beta = 0$

Response matrix ( $R$ ) based on a Gaussian resolution function with  $\sigma = 1.5 \times \text{bin width}$

Assume  $\epsilon_i = 1$

Results in

$P(\text{event to remain in the bin created}) = 26\%$   
 $P(\text{event to migrate 1 bin}) = 21\%$   
 $P(\text{event to migrate 2 bins}) = 16\%$

**Very large anti-correlations!**

- Applying the response matrix  $R$  smears out fine structure
- Applying  $R^{-1}$  creates (often unwanted) structure
- We do not have the expectation values  $\nu$
- Only have the data  $n$ , which are RVs and subject to statistical fluctuations

# So what went “wrong”?

- **Nothing really**

- The resulting unfolded yields are **unbiased**, but **heavily correlated**
- They can be used to test hypothesis, given one takes into account the **full set of correlations**

$$\chi^2 = (\hat{\mu} - \mu_0)^T U^{-1} (\hat{\mu} - \mu_0),$$

Use to test the compatibility of the estimators  $\hat{\mu}$  with the hypothesis  $\mu_0$

- Can reduce such oscillations considerably by **making bin widths larger than the width of the resolution function**

- **Alternatives:**

- Either incorporate **prior knowledge** or do not rely on neighboring bins to determine resolution correction
- Both come at a price: **trading variance for bias**

# Method 2: Correction factors (i)

- Assume the bins of the true distribution ( $\mu$ ) are the same as the data ( $n$ )
- Determine the **correction factor** for each bin (e.g., from MC simulation)

$$\hat{\mu}_i = C_i(n_i - \beta_i) \qquad C_i = \frac{\mu_i^{MC}}{\nu_i^{MC}}$$

← Run MC program **w/out** detector simulation

← Run MC program **with** detector simulation

- Works well if bin-to-bin sharing (smearing) is negligible  $R_{ij} = \delta_{ij}\epsilon_j$

$$\nu_i^{\text{sig}} = \nu_i - \beta_i = \epsilon_i \mu_i$$

- Expectation value for corrected data

$$E[\hat{\mu}_i] = C_i E[n_i - \beta_i] = C_i(\nu_i - \beta_i) = \frac{\mu_i^{MC}}{\nu_i^{MC}} \nu_i^{\text{sig}}$$

# Method 2: Correction factors (ii)

- Rearrange to make the bias explicit (identical to previous expression for  $E[\hat{\mu}_i]$ )

$$E[\hat{\mu}_i] = \frac{\mu_i^{MC}}{\nu_i^{MC}} \nu_i^{\text{sig}} = \underbrace{\left( \frac{\mu_i^{MC}}{\nu_i^{MC}} - \frac{\mu_i}{\nu_i^{\text{sig}}} \right)}_{\text{Bias}} \nu_i^{\text{sig}} + \mu_i$$

Bias = 0 if MC = nature

- Covariance matrix for the estimators

$$\text{cov}[\hat{\mu}_i, \hat{\mu}_j] = C_i^2 \text{cov}[n_i, n_j] = C_i^2 \delta_{ij} \nu_i$$

Smearing fluctuations independent between bins

- Iterative bin-by-bin method:

- Begin with (plausible) guess of the true spectrum
- Apply correction to measurement
- Generate new  $C_i$  from corrected spectrum of previous iteration
- Repeat (for a few iterations)

**Drawback: Highly model dependent**

# Method 3: Regularized unfolding

- **Regularization** = impose a measure of **smoothness** on the estimators of the true histogram  $\mu$

- Matrix inversion **IS** the maximum likelihood solution (see page 162)

Independent  
Poisson  
fluctuations

$$\log \mathcal{L}(\mu) = \sum_{i=1}^N (n_i \log \nu_i - \nu_i)$$

ML estimator  
(same as s56)

$$\hat{\nu} = \mathbf{n}$$

$$\hat{\mu} = R^{-1}(\mathbf{n} - \beta)$$

- Accept solutions that are **close** to the ML estimate

$$\log \mathcal{L}(\mu) \geq \log \mathcal{L}(\mu_{\max}) - \Delta \log \mathcal{L}(\mu)$$

↑ determines trade-off between bias and variance in unfolded histogram

- Define a **regularization (aka smoothness) function**  $S$  that increases when the unfolded solution becomes smoother
  - Task: choose the solution with the highest degree of smoothness out of the acceptable solutions determined by above inequality

- Must maximize  $\Phi(\mu) = \alpha \log \mathcal{L}(\mu) + S(\mu)$

↑ Regularization parameter which depends on  $\Delta \log \mathcal{L}(\mu)$   
 $\alpha \rightarrow \infty$  gives ML solution

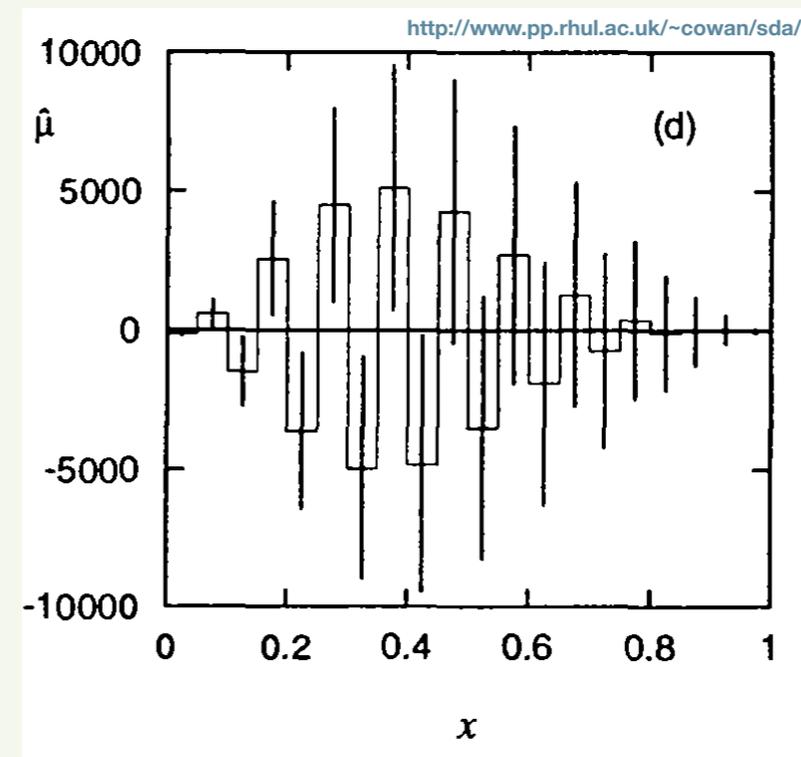
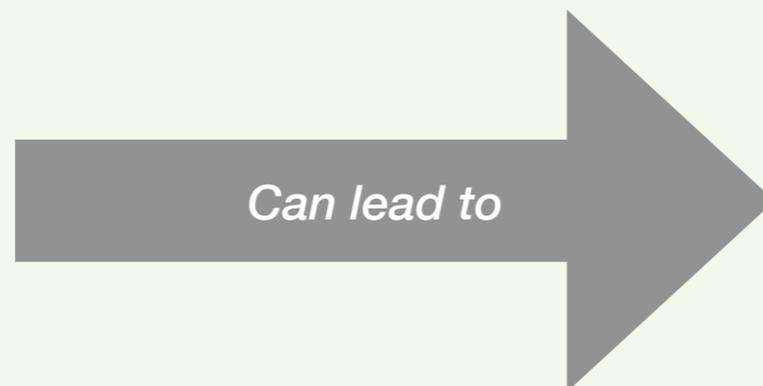
# Regularization functions

- There are several options (sections 11.5.1-11.5.4)
  - Tikhonov regularization:
    - Measure of smoothness is the mean value of the square of some derivative of the true distribution. Tikhonov regularization using the second derivative (so that  $S(\boldsymbol{\mu})$  is related to the avg. curvature) is widely used in particle physics.
  - Regularization functions based on entropy:
    - Interpret the entropy as a measure of the smoothness of a histogram. Estimators are constructed according to the principle of maximum entropy. Often developed in the framework of Bayesian statistics.
  - Regularization function based on cross-entropy:
    - Useful if we have prior knowledge that the true events approximately follow some distribution.

# Choice of $\alpha$

$$\Phi(\boldsymbol{\mu}) = \alpha \log \mathcal{L}(\boldsymbol{\mu}) + S(\boldsymbol{\mu})$$

- The choice of  $\alpha$  determines the trade off between the **bias** and **variance** of the estimators  $\hat{\boldsymbol{\mu}}$ 
  - If  $\alpha$  is very large, solution is dominated by the likelihood function and one has  $\log \mathcal{L}(\boldsymbol{\mu}) = \log \mathcal{L}_{\max}$  and very **large variances**



- If  $\alpha$  is small, leads to a perfectly smooth distribution (since all of the weight is put on the regularization function  $S$ )

# Choice of $\alpha$

- Recall the mean square error from L04, S11:

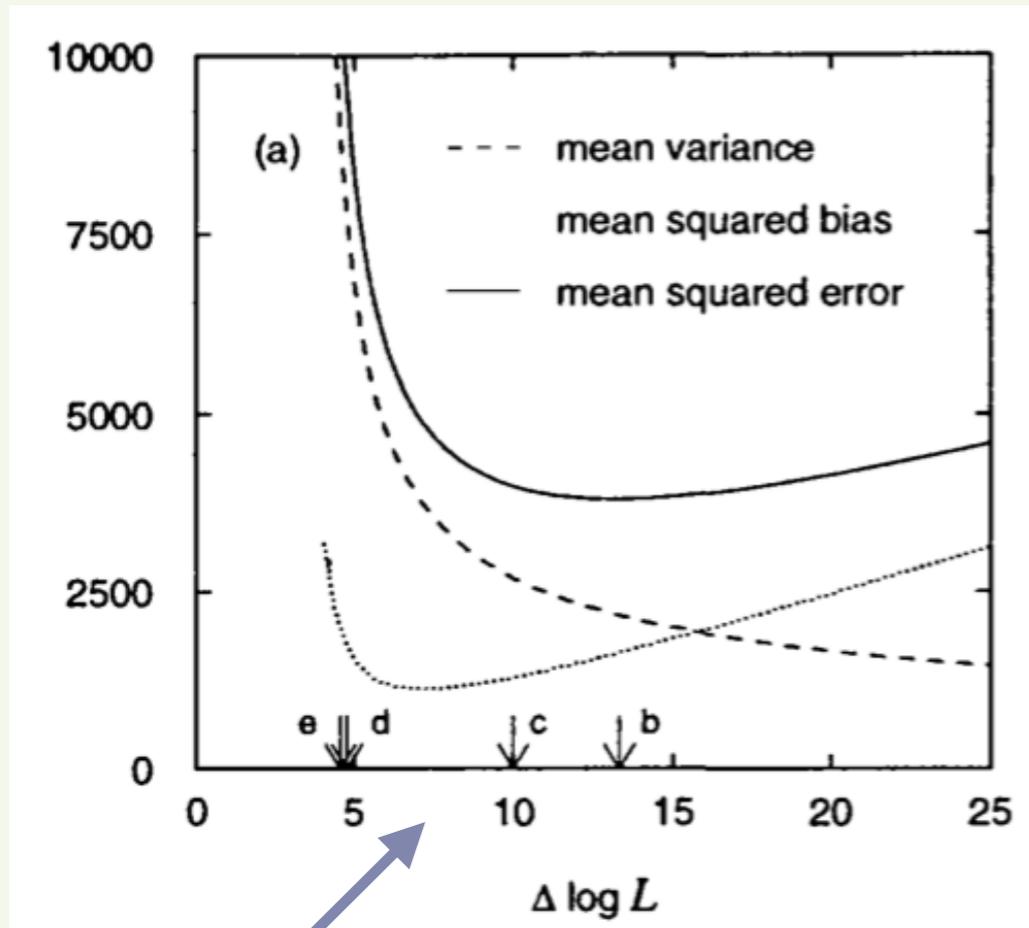
$$\begin{aligned} \text{MSE} &= E \left[ \left( \hat{\theta} - \theta \right)^2 \right] = E \left[ \left( \hat{\theta} - E[\hat{\theta}] \right)^2 \right] + \left( E[\hat{\theta}] - \theta \right)^2 \\ &= \underbrace{V[\hat{\theta}]}_{\text{variance}} + \underbrace{b^2}_{\text{bias}^2} \quad \text{i.e., sum of variance and bias}^2 \end{aligned}$$

*Interpret: sum of squares of statistical and systematic uncertainties*

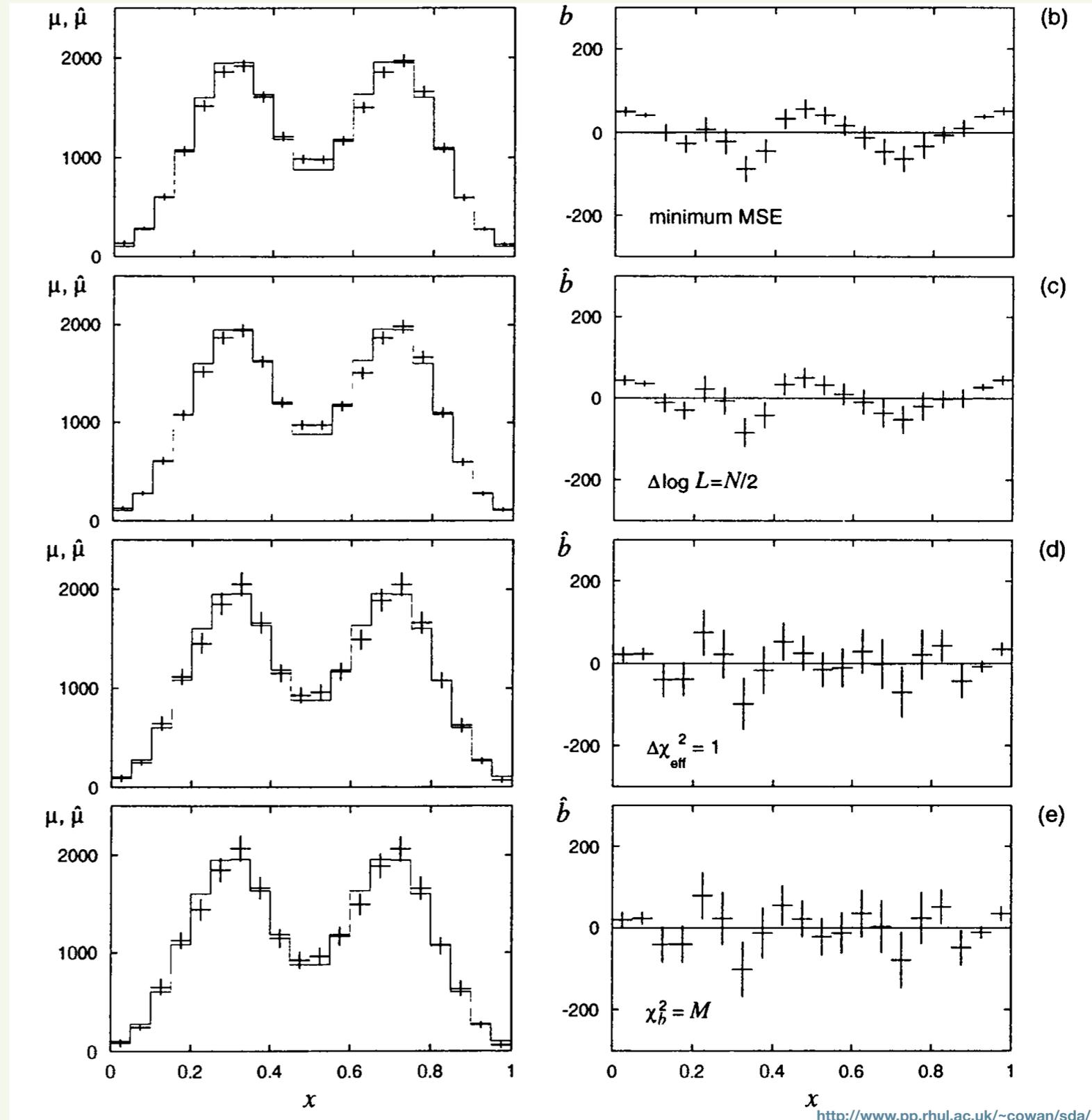
- Take the **MSE averaged over all bins** as the measure of the goodness of the final result. *One can determine  $\alpha$  so as to obtain a particular value of the MSE.*
- Can also use:
  - $\Delta \log \mathcal{L} = \log \mathcal{L}_{\max} - \log \mathcal{L} = N/2$
  - $\Delta \chi_{\text{eff}}^2 = 1$
  - $\chi_b^2 = M$

# Example with Maximum Entropy

- Return to our original example (which was unfolded using matrix inversion in s58)
- Now try with **Maximum Entropy regularization**



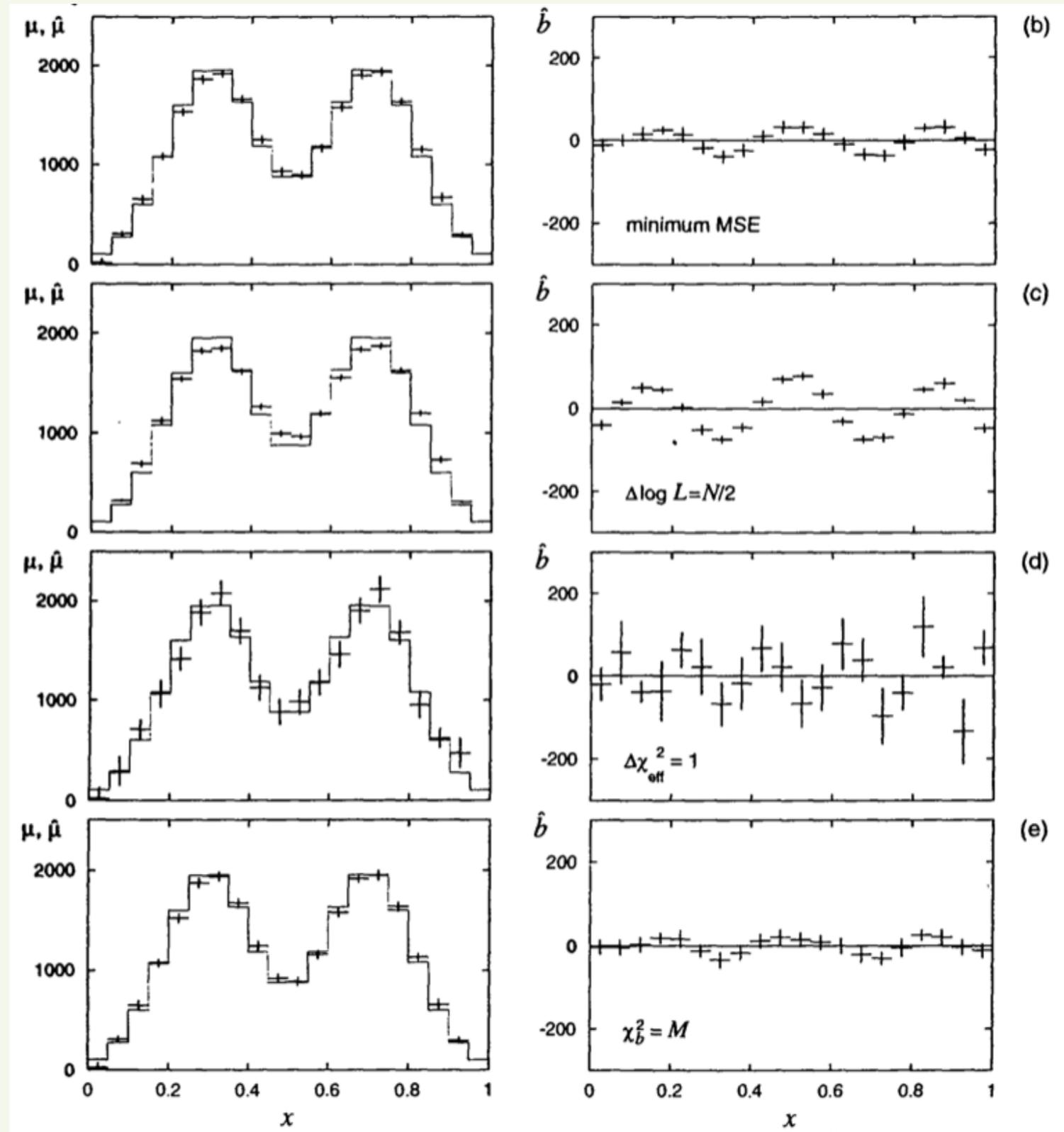
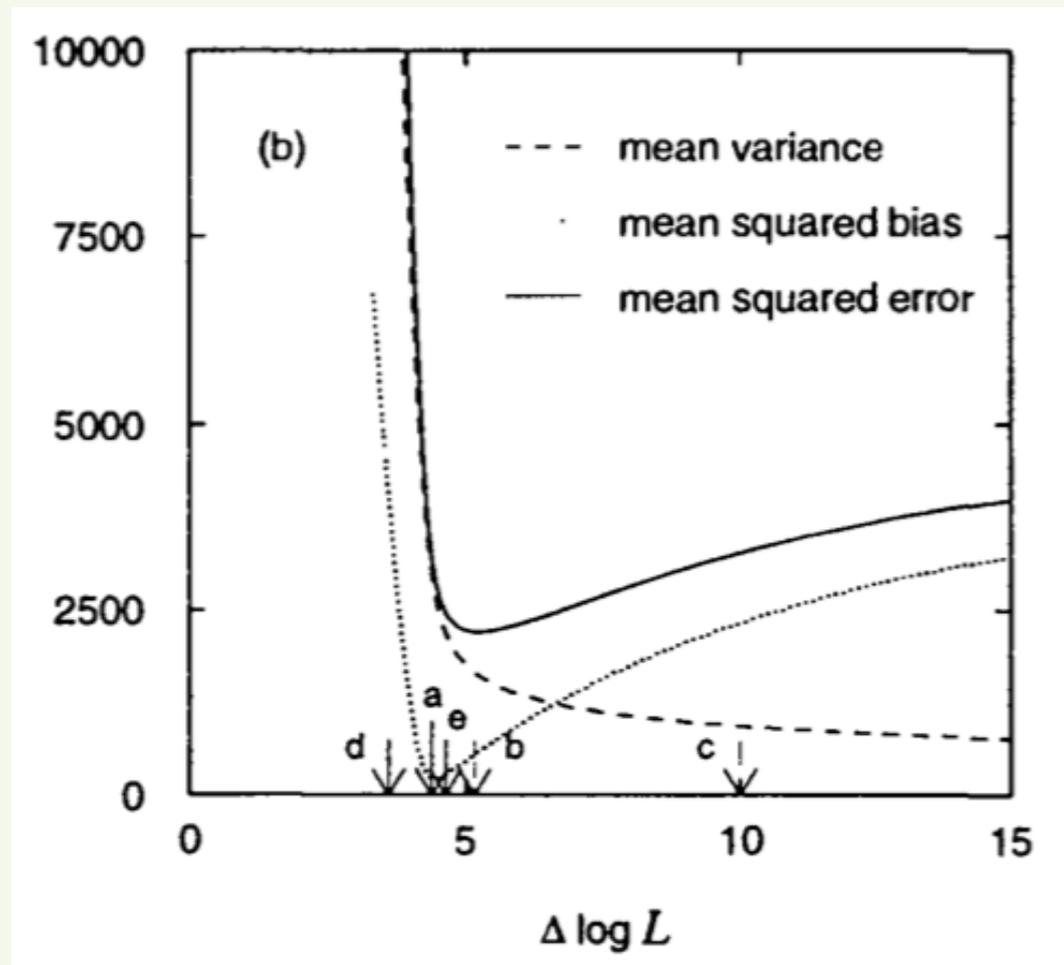
Arrows indicate solutions based on the criteria in the last slide



<http://www.pp.rhul.ac.uk/~cowan/sda/>

# Example with Tikhonov

- Return to our original example (which was unfolded using matrix inversion in s58)
- Now try with **Tikhonov regularization**



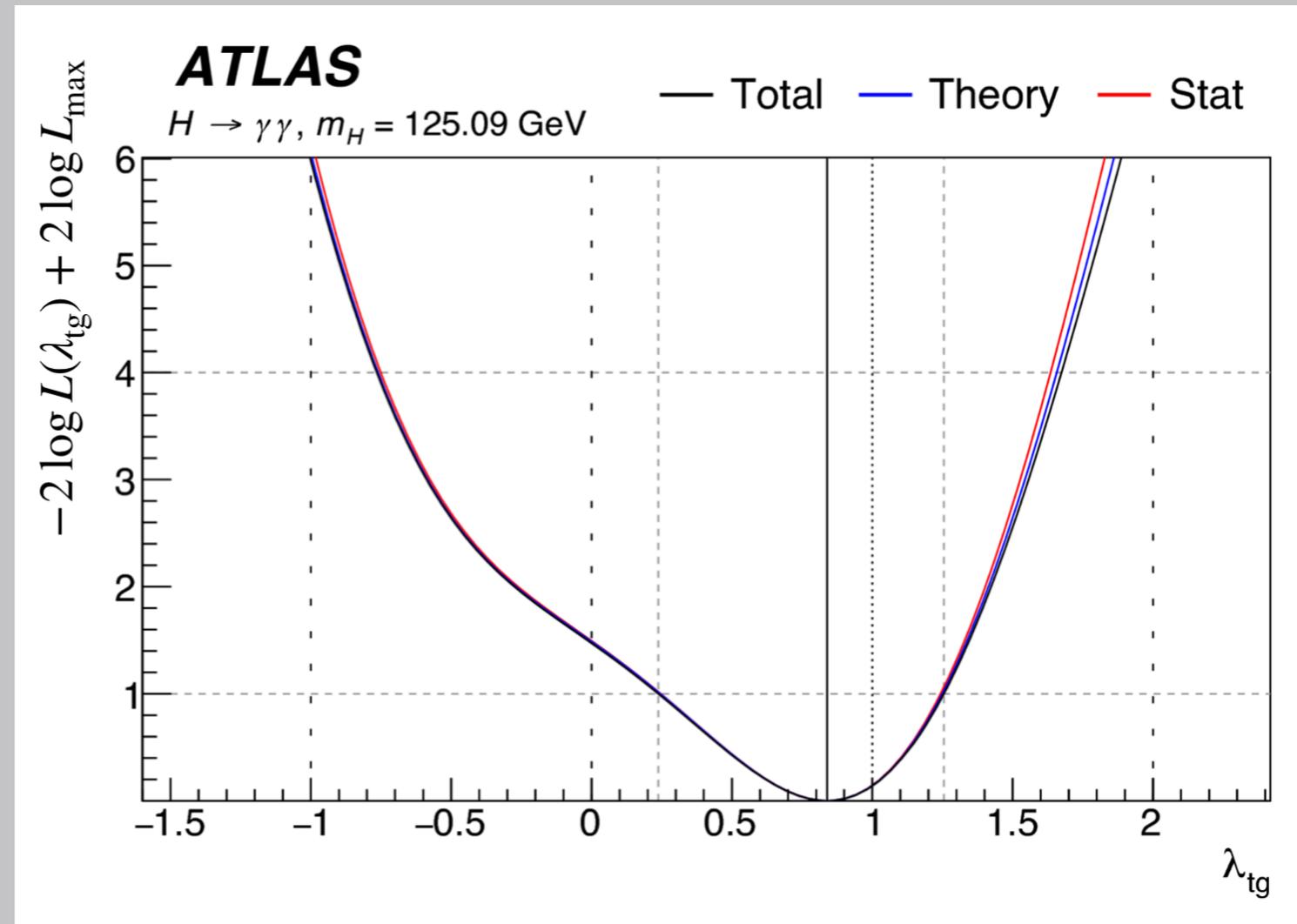
# For next time

- Required reading
  - Cowan textbook: chapters 9 (9.8-9.9), 10, and 11
- Extra reading for fun: /Reading material / L08 /
  - Search for  $B^+ \rightarrow \mu^+ \nu_\mu$  and  $B^+ \rightarrow \mu^+ N$  at Belle

**Quiz Time:** 8<sup>th</sup> Round

# Lower and upper limit

- Determine the lower and upper limit at 90% CL of the parameter  $\lambda_{tg}$  using the scan of the likelihood function and the provided values for  $Q_\gamma$ . Note that the likelihood is already normalised with respect to its maximal value.



$1 - \gamma$	$Q_\gamma$				
	$n = 1$	$n = 2$	$n = 3$	$n = 4$	$n = 5$
0.683	1.00	2.30	3.53	4.72	5.89
0.90	2.71	4.61	6.25	7.78	9.24
0.95	3.84	5.99	7.82	9.49	11.1
0.99	6.63	9.21	11.3	13.3	15.1

# The problem with priors

2. Write down the posterior probability density function of a parameter  $\theta$  as a function of the Likelihood of some data  $x$  and prior probability density function  $\pi(x)$ . What is the problem of using Bayesian priors when you quote a limit and which functional form for the prior somehow remedies them?

# Comprehension about unfolding

1. What is the unfolding problem? Write down the relevant equations for measurement  $\boldsymbol{\nu} = (\nu_1, \dots, \nu_N)$  with background  $\boldsymbol{\beta} = (\beta_1 \dots \beta_N)$ , which should be unfolded in yields  $\boldsymbol{\mu} = (\mu_1 \dots \mu_N)$  using a response matrix  $R$ .
2. Describe what the response matrix element  $R_{ij}$  means in terms of a conditional probability. Does the response matrix need to be a square or maybe even a symmetric matrix?
3. Describe two methods to solve the unfolding problem that do not involve regularization. Sketch out in detail what steps need to be taken.

# Correction factors and regularized unfolding

4. What is the method of correction factors? What are the drawbacks of using this method?

5. What is the idea behind regularized unfolding?



## KCETA Colloquium

# The muon $g-2$ window discrepancy and GeV-scale new physics

Thursday, June 22, 2023

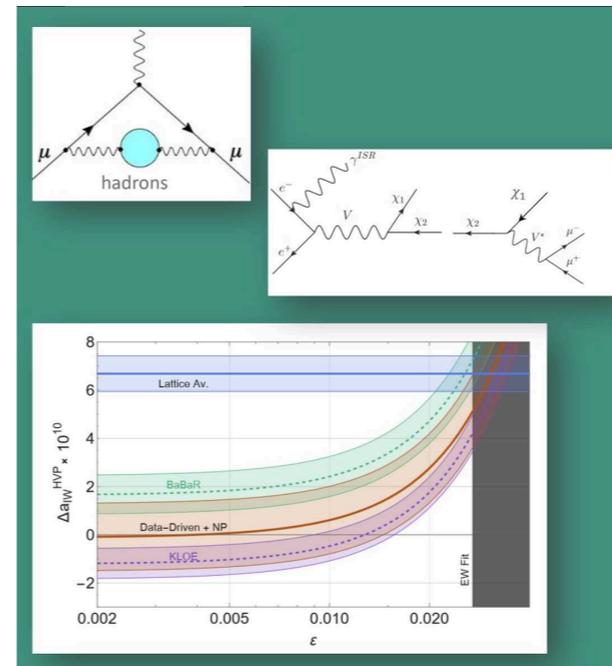
Kleiner Hörsaal A (CS) 15:45 - 17:00

Dr. Luc Darmé

(Institut de Physique des 2 Infinis de Lyon (IP2I), CNRS/IN2P3)

The decade-old discrepancy between the Standard Model prediction of the muon anomalous magnetic moment and the experimental results has seen striking developments in the past two years. In particular, recent lattice determinations of the hadronic vacuum polarization contribution deviate from the established data-driven ones at almost  $5\sigma$ . This new anomaly can be also seen as a tension between ab-initio lattice calculations and experimental measurements of  $e+e\rightarrow$  hadrons processes at and below the GeV scale.

We will review this puzzling situation and show how new processes beyond the standard model can affect indirectly the hadronic data around this scale, reconciling the lattice and data-driven results while complying with current phenomenological constraints. We will finally present a simple dark matter-motivated model as an explicit example.



Please note:

The colloquium will also be live-streamed to B402 SR 224 (CN).

# Bibliography

- Part of the material presented in this lecture is taken from the following sources. See the active links (when available) for a complete reference
  - **Statistical Data Analysis** textbook by G. Cowan (U. London): [all figures & equations with white background](#)