

HOLLOWAY **OF LONDON** 

03/07/24

#### Uncertainties on systematics in PDF fits

Enzo Canonero Glen Cowan

#### Motivation



- 1) Some systematic uncertainties can be well estimated:
  - Related to stat. error of control measurements
  - Related to size of MC event sample

- 2) But they can also be *quite uncertain*:
  - Theory systematics
  - Two points systematics ....



<u>see:</u> <u>G. Cowan, Eur. Phys. J. C (2019) 79:133; arXiv:1809.05778,</u> Canonero, E., Brazzale, A.R. & Cowan, G. *Eur. Phys. J. C* **83**, 1100 (2023)

### Motivation



- 1) Some systematic uncertainties can be well estimated:
  - Related to stat. error of control measurements
  - Related to size of MC event sample

- 2) But they can also be *quite uncertain*:
  - Theory systematics
  - Two points systematics ....
- Non-trivial consequences:
  - Fits are pulled less by incompatible data
  - Incompatible data are treated as an extra source of uncertainty resulting in inflated confidence intervals



see: G. Cowan, Eur. Phys. J. C (2019) 79:133; arXiv:1809.05778, Canonero, E., Brazzale, A.R. & Cowan, G. *Eur. Phys. J. C* 83, 1100 (2023)

## Formulation of the problem

ROYAL HOLLOWAY UNIVERSITY OF LONDON

- Suppose measurements y have a probability density  $P(y|\mu, \theta)$ 
  - $\mu$  = Parameters of interest
  - $\theta$  = Nuisance parameters
- Auxiliary Measurements *u* are used to provide info on nuisance parameters and are (often) assumed to be independently Gaussian distributed
- The resulting Likelihood is:

Can be a real measurement or just our best guess based on theoretical reasons

$$L(\boldsymbol{\mu}, \boldsymbol{\theta}) = P(\boldsymbol{y}, \boldsymbol{u} | \boldsymbol{\mu}, \boldsymbol{\theta}) = P(\boldsymbol{y} | \boldsymbol{\mu}, \boldsymbol{\theta}) \times \prod_{i} \frac{1}{\sqrt{2\pi}\sigma_{u_i}} e^{-(\boldsymbol{u}_i - \theta_i)^2/2\sigma_{u_i}^2}$$

## Formulation of the problem

ROYAL HOLLOWAY UNIVERSITY OF LONDON

- Suppose measurements y have a probability density  $P(y|\mu, \theta)$ 
  - $\mu$  = Parameters of interest
  - $\theta$  = Nuisance parameters
- Auxiliary Measurements *u* are used to provide info on nuisance parameters and are (often) assumed to be independently Gaussian distributed
- The resulting Likelihood is:

Can be a real measurement or just our best guess based on theoretical reasons

$$L(\boldsymbol{\mu},\boldsymbol{\theta}) = P(\boldsymbol{y},\boldsymbol{u}|\boldsymbol{\mu},\boldsymbol{\theta}) = P(\boldsymbol{y}|\boldsymbol{\mu},\boldsymbol{\theta}) \times \prod_{i} \frac{1}{\sqrt{2\pi}\sigma_{u_i}} e^{-(\boldsymbol{u}_i - \theta_i)^2/2\sigma_{u_i}^2}$$

• And the log Likelihood:

$$\log L(\boldsymbol{\mu}, \boldsymbol{\theta}) = \log P(\boldsymbol{y} | \boldsymbol{\mu}, \boldsymbol{\theta}) - \sum \frac{(\boldsymbol{u}_i - \boldsymbol{\theta}_i)^2}{2\sigma_{\boldsymbol{u}_i}^2}$$

## Formulation of the problem

ROYAL HOLLOWAY UNIVERSITY OF LONDON

- Suppose measurements y have a probability density  $P(y|\mu, \theta)$ 
  - $\mu$  = Parameters of interest
  - $\theta$  = Nuisance parameters
- Auxiliary Measurements *u* are used to provide info on nuisance parameters and are (often) assumed to be independently Gaussian distributed
- The resulting Likelihood is:

Can be a real measurement or just our best guess based on theoretical reasons

$$L(\boldsymbol{\mu}, \boldsymbol{\theta}) = P(\boldsymbol{y}, \boldsymbol{u} | \boldsymbol{\mu}, \boldsymbol{\theta}) = P(\boldsymbol{y} | \boldsymbol{\mu}, \boldsymbol{\theta}) \times \prod_{i} \frac{1}{\sqrt{2\pi}\sigma_{u_i}} e^{-(\boldsymbol{u}_i - \theta_i)^2/2\sigma_{u_i}^2}$$

• And the log Likelihood:

$$\log L(\mu, \theta) = \log P(\gamma | \mu, \theta) - \sum \frac{(u_i - \theta_i)^2}{2\sigma_{u_i}^2}$$
 Let systematic errors be potentially uncertain!

#### Gamma Variance Model (GVM)

- ROYAL HOLLOWAY UNIVERSITY OF LONDON
- The original quadratic terms in the log likelihood replaced by logarithmic terms:

$$\sum_{i} \frac{(u_i - \theta_i)^2}{2\sigma_{u_i}^2} \longrightarrow \sum_{i} \left(1 + \frac{1}{2\varepsilon_i^2}\right) \log\left(1 + 2\varepsilon_i^2 \frac{(u_i - \theta_i)^2}{\sigma_{u_i}^2}\right)$$

#### ε = error-on-error parameter

 $\epsilon$  = 0.3 means 30% uncertainty on  $\sigma$ 

#### Gamma Variance Model (GVM)

• The original quadratic terms in the log likelihood replaced by logarithmic terms:

$$\sum_{i} \frac{(u_i - \theta_i)^2}{2\sigma_{u_i}^2} \longrightarrow \sum_{i} \left(1 + \frac{1}{2\varepsilon_i^2}\right) \log\left(1 + 2\varepsilon_i^2 \frac{(u_i - \theta_i)^2}{\sigma_{u_i}^2}\right)$$

 $\epsilon$  = error-on-error parameter

 $\epsilon$  = 0.3 means 30% uncertainty on  $\sigma$ 

• Equivalent to switch from Gaussian constraints to Student's t constraints for systematics:





Suppose we want to average 4 measurements all with statistical and syst errors equal to 1.
 Also assume they all have equal errors-on-errors ε (auxiliary measurements set to zero):

$$\log L_P(\boldsymbol{\mu}, \boldsymbol{\theta}) = -\frac{1}{2} \sum_i \frac{(y_i - \mu - \theta_i)^2}{\sigma_{y_i}^2} - \frac{1}{2} \sum_i \left(1 + \frac{1}{2\boldsymbol{\varepsilon}_i^2}\right) \log \left(1 + 2\boldsymbol{\varepsilon}_i^2 \frac{\theta_i^2}{\sigma_{u_i}^2}\right)$$



Measurements internally compatible



Suppose we want to average 4 measurements all with statistical and syst errors equal to 1.
 Also assume they all have equal errors-on-errors ε (auxiliary measurements set to zero):

$$\log L_P(\boldsymbol{\mu}, \boldsymbol{\theta}) = -\frac{1}{2} \sum_i \frac{(y_i - \mu - \theta_i)^2}{\sigma_{y_i}^2} - \frac{1}{2} \sum_i \left(1 + \frac{1}{2\boldsymbol{\varepsilon}_i^2}\right) \log \left(1 + 2\boldsymbol{\varepsilon}_i^2 \frac{\theta_i^2}{\sigma_{u_i}^2}\right)$$







Suppose we want to average 4 measurements all with statistical and syst errors equal to 1.
 Also assume they all have equal errors-on-errors ε (auxiliary measurements set to zero):

$$\log L_P(\boldsymbol{\mu}, \boldsymbol{\theta}) = -\frac{1}{2} \sum_{i} \frac{(y_i - \mu - \theta_i)^2}{\sigma_{y_i}^2} - \frac{1}{2} \sum_{i} \left(1 + \frac{1}{2\boldsymbol{\varepsilon}_i^2}\right) \log\left(1 + 2\boldsymbol{\varepsilon}_i^2 \frac{\theta_i^2}{\sigma_{u_i}^2}\right)$$



24







- 1. If data are internally compatible results are only slightly modified
- 2. The estimate of the mean does not change when we increase  $\varepsilon$
- 3. The size of the confidence interval for the mean only slightly increases, reflecting the extra degree of uncertainty introduced by errors-on-errors



- Suppose one of the measurements is an outlier
- If data are internally incompatible important changes can be observed





- Suppose one of the measurements is an outlier
- If data are internally incompatible important changes can be observed





- Suppose one of the measurements is an outlier
- If data are internally incompatible important changes can be observed







- 1. With increasing  $\varepsilon$ , the estimate of mean is pulled less strongly by the outlier
- 2. The error bar grows more significantly: the GVM treats internal incompatibility as an additional source of uncertainty
- 3. The model is sensitive to internal compatibility of the data



• Gamma Variance Model:

$$\chi^{2} = \sum_{i} \frac{(y_{i} - f_{i}(\boldsymbol{a}) - \sum_{s} \Gamma_{i}^{s} \theta_{s})^{2}}{\sigma_{i}^{2}} + \sum_{s} \left(1 + \frac{1}{2\varepsilon_{i}^{2}}\right) \log(1 + 2\varepsilon_{i}^{2} \theta_{i}^{2})$$

• Appliable both to addictive and multiplicative systematics as only the systematic terms in the chi2 are being changed

- 1. Consider two measurements of the same distribution, analogous to results from two separate experiments.
- 2. Both distributions are subject to a normalization uncertainty, which is assumed to be itself uncertain.



ROYAL HOLLOWAY

- When considering errors-on-errors, the model gives greater weight to the more internally consistent distribution in the fit.
- The confidence interval is inflated to reflect the uncertainty coming from the conflicting scale factors.



#### Errors-on-errors: 0%



- When considering errors-on-errors, the model gives greater weight to the more internally consistent distribution in the fit.
- The confidence interval is inflated to reflect the uncertainty coming from the conflicting scale factors.



#### Errors-on-errors: 10%



- When considering errors-on-errors, the model gives greater weight to the more internally consistent distribution in the fit.
- The confidence interval is inflated to reflect the uncertainty coming from the conflicting scale factors.



#### Errors-on-errors: 20%



- When considering errors-on-errors, the model gives greater weight to the more internally consistent distribution in the fit.
- The confidence interval is inflated to reflect the uncertainty coming from the conflicting scale factors.



#### Errors-on-errors: 30%



- When considering errors-on-errors, the model gives greater weight to the more internally consistent distribution in the fit.
- The confidence interval is inflated to reflect the uncertainty coming from the conflicting scale factors.



#### Errors-on-errors: 40%



- When considering errors-on-errors, the model gives greater weight to the more internally consistent distribution in the fit.
- The confidence interval is inflated to reflect the uncertainty coming from the conflicting scale factors.



#### Errors-on-errors: 50%





• Treat uncertainties with 'errors-on-errors' as external parameters, since the minimization needs to be conducted numerically.

• Typically, only a few uncertainties will significantly impact the results if they are themselves uncertain, usually those with pulls greater than 1.

- Introduce an "eps parameter" input in the systematics namespace:
  - If eps = 0 or is not specified, use the standard quadratic term.
  - If eps > 0, treat it as an external parameter and modify the constraint in the chi2 accordingly.

#### Which file to modify?



For now I use this trivial approach (to be modified):

• Modify steering file

```
&Systematics
epsilon_value = 0.61
ListOfSources = 'sysHZComb1031:E', 'proc_tb21:E', 'sysHZComb1153:E', 'sysHZComb1098:E', 'sysHZComb1101:E'
&End
```

- Read epsilon value in read\_steer.f
- Modify GetChisquare.f

```
C Correlated chi2 part:
    fcorchi2_in = 0.d0
    do k=1, NSys
    if (SysForm(k) .eq. isExternal) then
        print *, epsilon_value, epsilon_bias
        temp_val = 2 * epsilon_value**2 * (rsys_in(k) - epsilon_bias)**2 * SysPriorScale(k)
        fcorchi2_in = fcorchi2_in + (1 + 1.0D0 / (2 * epsilon_value**2)) * log(1 + temp_val)
        else
        fcorchi2_in = fcorchi2_in + rsys_in(k)**2 * SysPriorScale(k)
        endif
```



- What else should be modified?
- Am I forgetting some other bits of the code (ex. Routines to compute uncertainties, ExtraParConstr.cc, ...)
- Other comments?



- Test the impact of "errors-on-errors" on the Hera fit using the initial simple implementation (for framework testing purposes).
- Enhance the implementation method.
- Improve the computation of confidence intervals, as "errors-onerrors" needs simulations for precise confidence interval analysis.



ROYAL HOLLOWAY UNIVERSITY OF LONDON

# Thank you for your attention



ROYAL HOLLOWAY UNIVERSITY OF LONDON

# Back-up slides

# Gamma Distributions



- Treat the systematic variances  $\sigma_{u_i}^2$  are *adjustable parameters* (*nuisance parameter*).
- Suppose their best estimates  $v_i$  are gamma distributed:



Gamma Distributions for Different  $\varepsilon$  Values 5  $\varepsilon = 0.5, \sigma = 1$  $\varepsilon = 0.2, \sigma = 1$ Probability Density  $\varepsilon = 0.1, \sigma = 1$  $\epsilon = 0.05, \sigma = 1$ 0.00 0.25 0.50 0.75 1.00 1.25 1.50 1.75 2.00

•  $\sigma_{u_i}$  Systematic Error

• 
$$\varepsilon_i = \frac{1}{2} \frac{\sigma_{v_i}}{\sigma_{u_i}^2} \cong \frac{\sqrt{v_i}}{\sigma_{u_i}}$$
 relative error on  $\sigma_{u_i}$ : "Error on error"

#### Gamma Variance Model (GVM)



• The likelihood is modified as follows:

$$L(\boldsymbol{\mu}, \boldsymbol{\theta}, \boldsymbol{\sigma_{u_i}^2}) = P(\boldsymbol{y}|\boldsymbol{\mu}, \boldsymbol{\theta}) \times \prod_{i} \frac{1}{\sqrt{2\pi}\sigma_{u_i}} e^{-(u_i - \theta_i)^2/2\sigma_{u_i}^2} \times \frac{\boldsymbol{\beta_i^{\alpha_i}}}{\boldsymbol{\Gamma(\alpha_i)}} \boldsymbol{v_i^{\alpha_i - 1}} e^{-\boldsymbol{\beta_i v_i}}$$

• One can profile over  $\sigma_{u_i}^2$  in closed form:

$$\log L_P(\boldsymbol{\mu}, \boldsymbol{\theta}) = \log P(\boldsymbol{y}|\boldsymbol{\mu}, \boldsymbol{\theta}) - \frac{1}{2} \sum_i \left( 1 + \frac{1}{2\varepsilon_i^2} \right) \log \left( 1 + 2\varepsilon_i^2 \frac{(\boldsymbol{u}_i - \boldsymbol{\theta}_i)^2}{\boldsymbol{v}_i} \right)$$

• We call this model the Gamma Variance Model (GVM)

(see: G. Cowan, Eur. Phys. J. C (2019) 79:133; arXiv:1809.05778)



- Gamma distributions allow to parametrize distributions of positive defined variables (like estimates of variances)
- Using Gamma distributions it is possible to profile in close form over  $\sigma_i^2$

#### Motivation for the GVM



• Gamma distributions include the case where the variance is estimate from a real dataset of control measurements:

$$v_i = \frac{1}{n_i - 1} \sum \left( u_{i,j} - \overline{u_i} \right)^2$$

•  $(n-1)v_i/\sigma_{u_i}^2$  follows a  $\chi_{n-1}^2$  distribution and  $v_i$  a Gamma distribution with:

$$\alpha_i = \frac{n_i - 1}{2}$$
$$\beta_i = \frac{n_i - 1}{2\sigma_{u_i}^2}$$

- ROYAL HOLLOWAY UNIVERSITY OF LONDON
- BLUE (Best Linear Unbiased Estimators) approach to combinations:

$$\chi^{2} = \sum_{i} (y_{i} - f(a)) V_{ij}^{-1} (y_{j} - f(a))$$

$$V_{ij} = V_{ij}^{(stat)} + V_{ij}^{(syst)}$$

- $V_{ij}^{(stat)}$ : Statistical covariance matrix.
- $V_{ij}^{(syst)}$ : Covariance matrix induced by systematic source.
- $V_{ij}^{(syst)} = \sum_{s} V_{ij}^{(s)}$









• Connection:

$$V_{ij}^{(syst)} = \sum_{s} V_{ij}^{(s)}$$
$$V_{ij}^{(s)} = \Gamma_i^s \Gamma_j^s$$



• Gamma Variance Model:

$$\chi^{2} = \sum_{i} \frac{(y_{i} - f(\boldsymbol{a}) - \sum_{s} \Gamma_{i}^{s} \theta_{s})^{2}}{\sigma_{i}^{2}} + \sum_{s} \left(1 + \frac{1}{2\varepsilon_{i}^{2}}\right) \log\left(1 + 2\varepsilon_{i}^{2} \theta_{i}^{2}\right)$$

• What to do if we do not have access to the factors  $\Gamma_i^s$  (we only know  $V_{ij}^{(syst)}$ )?

$$V_{ij}^{(syst)} = \sum_{s} V_{ij}^{(s)}$$

• Switch to a nuisance parameters approach:

**Proof is non-trivial!** 

$$\chi^{2} = \sum_{i} \frac{(y_{i} - \mu - \theta_{i})^{2}}{\sigma_{i}^{2}} + \sum_{ij} \theta_{i} C_{ij}^{-1} \theta_{j}$$
$$C_{ij} = V_{ij}^{(s)}$$

• Substitute quadratic term with log-constraint:

$$\sum_{ij} \theta_i C_{ij}^{-1} \theta_j \longrightarrow \sum_i \left( N + \frac{1}{2\varepsilon_i^2} \right) \log \left( 1 + 2\varepsilon_i^2 \theta_i C_{ij}^{-1} \theta_j \right)$$





- The Hessian method is based on the assumption that the  $\chi^2$  follows a  $\chi^2$  distribution.
- Our "goodness-of-fit" statistics q is not a  $\chi^2$  so will will not follow exactly a  $\chi^2$  for large values of  $\epsilon^2$

*Large literature on the topic:* 

- Bartlett, M. S. (1937) Proceedings of the Royal Society A, 160, 268–282)
- Applied Asymptotics Case Studies in Small-Sample Statistics by A. R. Brazzale, A. C. Davison and N. Reid)
- Canonero, E., Brazzale, A.R. & Cowan, *Eur. Phys. J. C* 83, 1100 (2023).



• Modify the test statistic q so that its distribution is closer to a  $\chi^2$ :

$$q \quad \longrightarrow \quad q^* = q \, \frac{N_{dof}}{E[q]}$$

#### **Bartlett Correction**



• Modify the test statistic q so that its distribution is closer to a  $\chi^2$ :



#### **Bartlett Correction**



• Modify the test statistic q so that its distribution is closer to a  $\chi^2$ :



# Simplified Model (no real data)



#### GOAL:

- Construct a simplified toy model to test the implementations of errors-on-errors in a real PDF fit
- Choose a simple process that allows an easy and fast implementation.



# Simplified Model (no real data)



#### **GOAL**:

- Construct a simplified toy model to test the implementations of errors-on-errors in a real PDF fit
- Choose a simple process that allows an easy and fast implementation.



# Simplified Model (no real data)



- Construct a simplified toy model to test the implementations of errors-on-errors in a real PDF fit
- Choose a simple process that allows an easy and fast implementation.



# Simplified Model



- The aim of the exercise is to fit the gluon PDF, using fictious data points.
- The gluon PDF is parametrized as follow
  - $g(x) = C x^A (1-x)^B$

 $\begin{cases} A = -0.85 \\ B = 6 \end{cases}$ 

• 
$$C: \int_0^1 g(x) dx = 1/2$$



• We are assuming that this is the gluon PDF shape at  $Q^2$  close to  $t\bar{t}$  production scale.