

# HepData for BSM reinterpretation of LHC data

Andy Buckley, University of Glasgow  
HepData Advisory Meeting, 24 Nov 2017



# Intro

- ▶ **Quantity and quality of HepData content from LHC has been steadily improving – including ad hoc “auxiliary data”**
- ▶ BSM pheno community very happy about this – *vision of an automated limit re-setting toolchain with comprehensive data coverage*
- ▶ Include both “SM measurements” & dedicated BSM search data
- ▶ Primary data faithfully recorded, modulo format details. Issue is with *secondary data*:
  - (MC) background estimates
  - Correlation data
- ▶ **All data needs to “automatically” flow from experiments, through HepData, and into analysis tools**  
⇒ **standardise formats and conventions for data & aux data**

# Intro

- ▶ **Quantity and quality of HepData content from LHC has been steadily improving – including ad hoc “auxiliary data”**
- ▶ BSM pheno community very happy about this – *vision of an automated limit re-setting toolchain with comprehensive data coverage*
- ▶ Include both “SM measurements” & dedicated BSM search data
- ▶ Primary data faithfully recorded, modulo format details. Issue is with *secondary data*:
  - (MC) background estimates
  - Correlation data
- ▶ **All data needs to “automatically” flow from experiments, through HepData, and into analysis tools**  
⇒ standardise formats and conventions for data & aux data
- ▶ **Frames a potential work-plan to include in funding application**

# Correlations in fits/limit setting

## Many types of correlation:

- ▶ **Between bins/SRs**, introduced by experimental/theory systematics
- ▶ **Between bins/analyses**, introduced by sharing events (or normalisation)
- ▶ **Between systematic (nuisance) params**, induced by profile fitting

# Correlations in fits/limit setting

## Many types of correlation:

- ▶ **Between bins/SRs**, introduced by experimental/theory systematics
- ▶ **Between bins/analyses**, introduced by sharing events (or normalisation)
- ▶ **Between systematic (nuisance) params**, induced by profile fitting

## Possible approaches to providing this information:

- ▶ **full likelihood expression**, e.g. **HistFactory demo** ↗
- ▶ approximate: **express as independent error sources**, correlated across bins — *extensible*
- ▶ approximate: **simplified likelihoods**: drop connection to error sources, bkg systs only, express as (symm) bin covariance  
Actively used by CMS: <https://cds.cern.ch/record/2242860>

# Correlations in fits/limit setting

## Many types of correlation:

- ▶ **Between bins/SRs**, introduced by experimental/theory systematics
- ▶ **Between bins/analyses**, introduced by sharing events (or normalisation)
- ▶ **Between systematic (nuisance) params**, induced by profile fitting

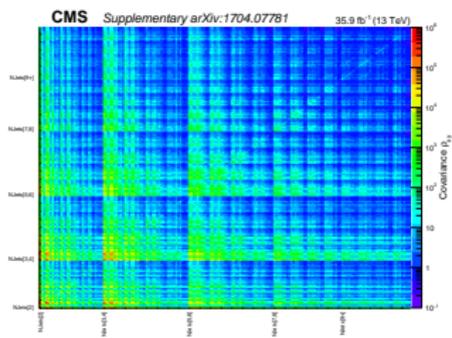
## Possible approaches to providing this information:

- ▶ **full likelihood expression**, e.g. **HistFactory demo** ↗
- ▶ approximate: **express as independent error sources**, correlated across bins — *extensible*
- ▶ approximate: **simplified likelihoods**: drop connection to error sources, bkg systs only, express as (symm) bin covariance  
Actively used by CMS: <https://cds.cern.ch/record/2242860>

**Not 100% clear that correlations are necessary, but without them there will always be questions of whether an analysis was too optimistic or conservative**

# Correlation formats: error sources vs. bin covariance

CMS  $0\ell$  cov matrix – note log-scale!



Error breakdown in a HepData record

NB. normal in *Standard Model* analyses

RE	P P → JETS
<b>COS PHI</b>	<b>TEEC</b>
-1 - -0.96	10.5165 ±0.00779481 stat +0.0117651 sys_jesHp1 +0.0034308 sys_jesHp2 + 71 more errors <a href="#">Show all</a>
-0.96 - -0.92	0.716955 ±0.00468718 stat +0.00357604 sys_jesHp1 +0.00145823 sys_jesHp2 + 71 more errors <a href="#">Show all</a>
-0.92 - -0.88	0.322052 ±0.00259636 stat +0.00184137 sys_jesHp1 +0.000814961 sys_jesHp2 + 71 more errors <a href="#">Show all</a>

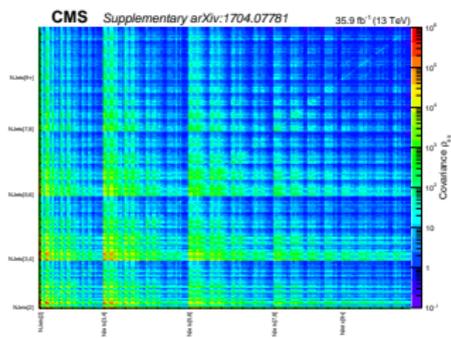
SL originally formalised as symm covariance

Simple to use:  $L(\mu, \vec{\theta}) = \prod_i \text{Pois}(n_i, \mu, \vec{\theta}) \cdot \text{Gaus}(\vec{\theta}, \mathbf{C})$

Dimensionality of cov fixed: uniform approach, scales well. But limited to symmetric errs and no correlations between analyses.

# Correlation formats: error sources vs. bin covariance

CMS  $0\ell$  cov matrix – note log-scale!



Error breakdown in a HepData record  
NB. normal in *Standard Model* analyses

RE	P P -> JETS
<b>COS PHI</b>	<b>TEEC</b>
-1 - -0.96	10.5165 ±0.00779481 stat +0.0117651 sys_jesthp1 +0.0034308 sys_jesthp2 + 71 more errors <a href="#">Show all</a>
-0.96 - -0.92	0.716955 ±0.00468718 stat +0.00357604 sys_jesthp1 +0.00165822 sys_jesthp2 + 71 more errors <a href="#">Show all</a>
-0.92 - -0.88	0.322052 ±0.00259636 stat +0.00184137 sys_jesthp1 +0.000814961 sys_jesthp2 + 71 more errors <a href="#">Show all</a>

SL originally formalised as symm covariance

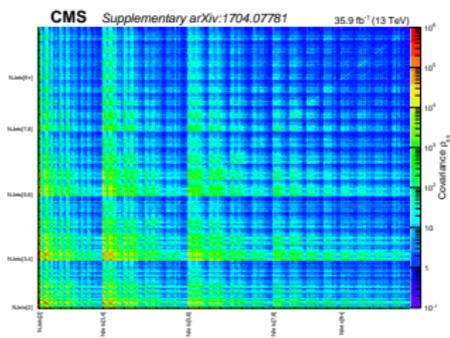
Simple to use:  $L(\mu, \vec{\theta}) = \prod_i \text{Pois}(n_i, \mu, \vec{\theta}) \cdot \text{Gaus}(\vec{\theta}, \mathbf{C})$

Dimensionality of cov fixed: uniform approach, scales well. But limited to symmetric errs and no correlations between analyses.

**HepData doesn't understand datasets semantics: would need "link" metadata to reliably connect correlation datasets to primary datasets**

# Correlation formats: error sources vs. bin covariance

CMS  $0\ell$  cov matrix – note log-scale!



Error breakdown in a HepData record  
NB. normal in *Standard Model* analyses

RE	P P -> JETS
<b>COS PHI</b>	<b>TEEC</b>
-1 -0.96	10.5165 ±0.00779481 stat +0.0117651 sys_jesHp1 +0.0034308 sys_jesHp2 + 71 more errors <a href="#">Show all</a>
-0.96 -0.92	0.716955 ±0.00468718 stat +0.00357604 sys_jesHp1 +0.00165823 sys_jesHp2 + 71 more errors <a href="#">Show all</a>
-0.92 -0.88	0.322052 ±0.00259636 stat +0.00184137 sys_jesHp1 +0.000814961 sys_jesHp2 + 71 more errors <a href="#">Show all</a>

SL originally formalised as symm covariance

Simple to use:  $L(\mu, \vec{\theta}) = \prod_i \text{Pois}(n_i, \mu, \vec{\theta}) \cdot \text{Gaus}(\vec{\theta}, \mathbf{C})$

Dimensionality of cov fixed: uniform approach, scales well. But limited to symmetric errs and no correlations between analyses.

**HepData doesn't understand datasets semantics: would need "link" metadata to reliably connect correlation datasets to primary datasets**

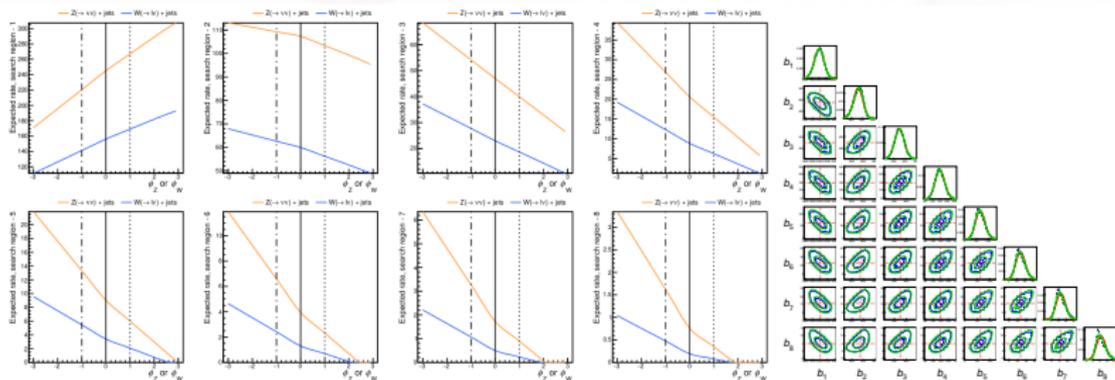
Error-source representation more flexible: can construct cov matrix

$C_{ij} = \sum_e \sigma_i \sigma_j$ , or asymm by toy-sampling

**Extensible! Supported already. HD preference. But...**

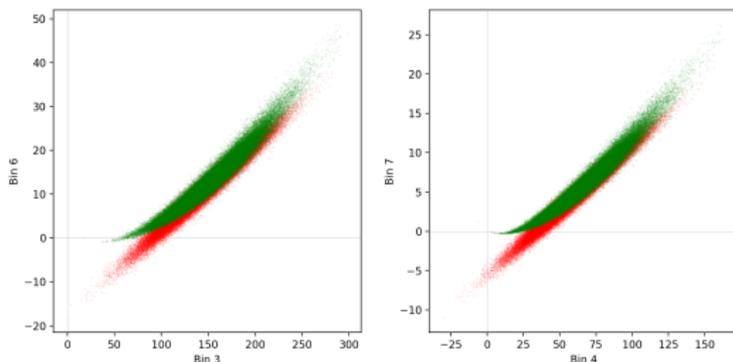
# Logistical issues & extensions

- ▶ Need standard names, esp. to distinguish uncorr stat errors
- ▶ Also need groupings, e.g. to separate theory/MC errors from experimental/detector resolutions  $\Rightarrow$  future reinterpretations with theory improvements. Easier with explicit cov matrices?
- ▶ Error-sources are naturally usable in an asymmetric way. But **current activity**  $\Rightarrow$  on use of skew moments to implement asymm parametrisation: **how to store this in HD?!**



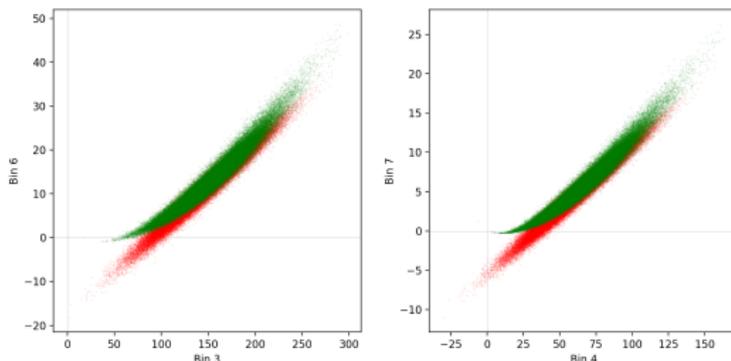
# Logistical issues & extensions

- ▶ Need standard names, esp. to distinguish uncorr stat errors
- ▶ Also need groupings, e.g. to separate theory/MC errors from experimental/detector resolutions  $\Rightarrow$  future reinterpretations with theory improvements. Easier with explicit cov matrices?
- ▶ Error-sources are naturally usable in an asymmetric way. But **current activity**  $\checkmark$  on use of skew moments to implement asymm parametrisation: **how to store this in HD?!**



# Logistical issues & extensions

- ▶ Need standard names, esp. to distinguish uncorr stat errors
- ▶ Also need groupings, e.g. to separate theory/MC errors from experimental/detector resolutions  $\Rightarrow$  future reinterpretations with theory improvements. Easier with explicit cov matrices?
- ▶ Error-sources are naturally usable in an asymmetric way. But **current activity**  $\square$  on use of skew moments to implement asymm parametrisation: **how to store this in HD?!**



**Possibility for HD to have *semantic* understanding of correlations?**

## MC and background data

- ▶ Correlations are the most technically complex demand, since the data objects are semantically different from “normal” datasets

# MC and background data

- ▶ Correlations are the most technically complex demand, since the data objects are semantically different from “normal” datasets
- ▶ Not the only requirement for scalable recasting, though: background estimates are also crucial

# MC and background data

- ▶ Correlations are the most technically complex demand, since the data objects are semantically different from “normal” datasets
- ▶ Not the only requirement for scalable recasting, though: background estimates are also crucial
- ▶ Typical BSM reinterpretations only have the capacity to generate (maybe LO) signal events

# MC and background data

- ▶ Correlations are the most technically complex demand, since the data objects are semantically different from “normal” datasets
- ▶ Not the only requirement for scalable recasting, though: background estimates are also crucial
- ▶ Typical BSM reinterpretations only have the capacity to generate (maybe LO) signal events
- ▶ Backgrounds computed by experiments using vast MC datasets with very complex and CPU-intensive high-sophistication modelling: not reproducible, so needs to be published

# MC and background data

- ▶ Correlations are the most technically complex demand, since the data objects are semantically different from “normal” datasets
- ▶ Not the only requirement for scalable recasting, though: background estimates are also crucial
- ▶ Typical BSM reinterpretations only have the capacity to generate (maybe LO) signal events
- ▶ Backgrounds computed by experiments using vast MC datasets with very complex and CPU-intensive high-sophistication modelling: not reproducible, so needs to be published
- ▶ This has started, but – again – **how to make HD (and its API) *semantically aware of what is data and what’s the corresponding MC?***

And background process breakdown? And pre-/post-fit? ...