Publication and Exchange of Data and Models or: making the most out of your analysis



CLUSTER OF EXCELLENCE QUANTUM UNIVERSE

PCD Data Science Basics – Oct 25, 2022





Sabine Kraml LPSC Grenoble

Outline

- Introduction and motivation
- Publication of statistical models
- Analysis preservation and reuse
- Specifics about ML models
- Conclusions



Motivation

Reproducibility and reuse of scientific results

- Reproducibility is a major principle underpinning the scientific method. It means that results obtained by an experiment or an observational study or in a statistical analysis of a data set should be achieved again with a high degree of reliability when the study is replicated.
- Goes back to Robert Boyle, a pioneer of the experimental method, who maintained that the foundations of knowledge should be constituted by experimentally produced facts, which can be made believable to a scientific community by their reproducibility.

Scientific results should be documented in such a way that their deduction is fully transparent.

- dataset and the code to calculate the results easily accessible.

Reproducibility can be distinguished from replication, as referring to reproducing the same results using the same dataset.



https://en.wikipedia.org/wiki/Scientific_method

Requires a detailed description of the methods used to obtain the data and making the full

In computational sciences: Any results should be documented by making all data and code available in such a way that the computations can be executed again with identical results.



Motivation

Reproducibility and reuse of scientific results

Reproducibility also enables reuse



e.g., updating constraints, testing new hypotheses, performing combinations and/or fits, etc.

 \rightarrow new research based on existing data and analyses

> → longer shelf life, more scientific impact



Snowmass 2021

US Community Study on the Future of Particle Physics

To achieve their full scientific impact, HEP experiments need to integrate extensive data and analysis preservation efforts into their publication processes, alongside the communication of results in reusable form and preservation of data products, and making event-level data publicly available.

Without this, the influence of the hundreds of published analyses from the LHC, HL-LHC, EIC, and other future experiments will be **limited mainly to the physics** ideas in vogue at the time the collaboration collected their data. The public investment in experimental programs underscores the importance of going beyond the original paper publication and ensuring that analyses continue providing scientific value in perpetuity.

> Executive summary from "Data and Analysis Preservation, Recasting and Reinterpretation" arXiv:2203.10057



On the individual scale

reusable form and preservation of data products.

Without this, the influence of your work will be limited to the physics idea en vogue \rightarrow the hypothesis pursued in the paper. The public investment \rightarrow your own *intellectual effort* in the study underscores the importance of going beyond the original paper publication and ensuring that analyses continue providing scientific value in perpetuity.

> Adaptation of Executive summary from "Data and Analysis Preservation, Recasting and Reinterpretation" arXiv:2203.10057

To achieve its full scientific impact, your analysis needs to integrate extensive data and analysis preservation efforts, alongside the communication of results in



Typical particle physics analysis



Typical particle physics ana



lysis	

Publication of statistical models

"Publishing statistical models: Getting the most out of particle physics experiments", K. Cranmer, SK, H. Prosper (eds) et al., arXiv:2109.04981

Statistical models

ascribing to it a probability, as specified by a statistical model p(data|theory)



In particle physics experiments, the statistical nature of the data is typically quantified by

Illustration courtesy Lukas Heinrich



Statistical models

ascribing to it a probability, as specified by a statistical model p(data|theory)

Primary measurements

Parameters of interest (POIs)

 $p(x, y | \mu, \theta) = p(x)$

auxiliary data

nuisance parameters

and the **nuisance parameters**.

In particle physics experiments, the statistical nature of the data is typically quantified by

Probability density of the auxiliary data

The values *y* are often estimates of corresponding nuisance parameters; their probability may be, e.g., a Gaussian with a specified standard deviation

Describes the probabilistic dependence of the observable data on the parameters of interest

When **observed data** are entered into the statistical model, this becomes the **likelihood function**.

w/o proper statistical model, a Gaussian approximation is forced onto the reuse of experimental results

$$\chi^{2}(\mu) = \frac{1}{N_{\text{dat}}} \sum_{ij=1}^{N_{\text{dat}}} \left(\mathcal{O}_{i}^{(\text{th})}(\mu) - \mathcal{O}_{i}^{(\mu)}(\mu) \right)$$

In principle OK if large enough statistics (\rightarrow CLT). However:

- Symmetric Gaussian uncertainties are often a simplification; does not hold when systematics dominate
- Often lack of information on correlations
- Non-positive-definite covariance matrices on HEPData
- Lack of breakdown of correlated systematic sources; different naming for systematic sources complicates combining processes
- Correlations between processes often not available
- In case of asymmetric uncertainties, assumptions have to be made about the shape of the likelihood function

Measurements reported as "observed value ± 1 sigma" (plus correlations)

 $\binom{(\exp)}{i} \left(\operatorname{cov}^{-1} \right)_{ij} \left(\mathcal{O}_j^{(\operatorname{th})}(\mu) - \mathcal{O}_j^{(\exp)} \right)$



11

Severely affects many areas:

- Parton distribution functions
- Effective field theory fits
- Higgs physics
- BSM searches
- Heavy flavour physics
- Global averages
- BSM Global fits

Examples discussed in arXiv:2109.04981

w/o proper statistical model, a Gaussian approximation is forced onto the reuse of experimental results

$$\chi^{2}(\mu) = \frac{1}{N_{\text{dat}}} \sum_{ij=1}^{N_{\text{dat}}} \left(\mathcal{O}_{i}^{(\text{th})}(\mu) - \mathcal{O}_{i}^{(\text{exp})} \right) \left(\text{cov}^{-1} \right)_{ij} \left(\mathcal{O}_{j}^{(\text{th})}(\mu) - \mathcal{O}_{j}^{(\text{exp})} \right)$$

In principle OK if large enough statistics (\rightarrow CLT). However:

- Symmetric Gaussian uncertainties are often a simplification; does not hold when systematics dominate
- Often lack of information on correlations
- Non-positive-definite covariance matrices on HEPData
- Lack of breakdown of correlated systematic sources; different naming for systematic sources complicates combining processes
- Correlations between processes often not available
- In case of asymmetric uncertainties, assumptions have to be made about the shape of the likelihood function

Measurements reported as "observed value ± 1 sigma" (plus correlations)







w/o proper statistical model, a Gaussian approximation is forced onto the reuse of experimental results

$$\chi^{2}(\mu) = \frac{1}{N_{\text{dat}}} \sum_{ij=1}^{N_{\text{dat}}} \left(\mathcal{O}_{i}^{(\text{th})}(\mu) - \mathcal{O}_{i}^{(\text{exp})} \right) \left(\text{cov}^{-1} \right)_{ij} \left(\mathcal{O}_{j}^{(\text{th})}(\mu) - \mathcal{O}_{j}^{(\text{exp})} \right) \right)$$

In principle OK if large enough statistics (\rightarrow CLT). However:

- Symmetric Gaussian uncertainties are often a simplification; does not hold when systematics dominate
- Often lack of information on correlations
- Non-positive-definite covariance matrices on HEPData
- Lack of breakdown of correlated systematic sources; different naming for systematic sources complicates combining processes
- Correlations between processes often not available
- In case of asymmetric uncertainties, assumptions have to be made about the shape of the likelihood function

 $\rho_{xy} = Correlation(x, y) = \frac{c}{\sqrt{var}}$



Examples discussed in arXiv:2109.04981





Higgs Measurements — Limitations of Reconstructed LLHs

Example: ATLAS $H \rightarrow ZZ \rightarrow 4\ell$ [ATLAS 2004.03447]

HiggsSignals implementation

measurements • (12-bin STXS)

- experimental correlations
- theory correlations [2017 Scheme]









Jonas Wittbrodt

Publication of Statistical Models Hands-on workshop 8-12 Nov 2021 https://indico.cern.ch/event/1088121/

ill-defined covariance matrices

One particularly worrisome consequence of the lack of open likelihoods is that results may become very sensitive to small variations of the official correlation model



Assess impact in fit by transforming the original covariance matrix into a matrix with the **same** eigenvectors but with clipped eigenvalues below some cut-off: stable PDFs with much lower x²

Dataset	$N_{ m dat}$	$Z_{ m orig}$	$\chi^2_{ m orig}$	$\chi^2_{ m reg}$	
ATLAS W, Z 7 TeV CC ($\mathcal{L} = 4.6 \text{ fb}^{-1}$)	46	9.01	1.89	0.93	minimal modification of
ATLAS W 8 TeV (*)	22	11.28	3.50	1.15	correlation model, large
CMS dijets 7 TeV	54	4.70	1.81	1.73	impact in fit quality,
ATLAS dijets 7 TeV	90	9.93	2.14	0.92	PDFs stable
CMS 3D dijets 8 TeV (*)	122	4.47	1.50	0.92	
		v			

Juan Rojo on PDF fits

Key component of predictions for particle, nuclear, and astro-particle experiments.

Address fundamental questions in QCD.

Publication of Statistical Models Hands-on workshop 8-12 Nov 2021 https://indico.cern.ch/event/1088121/





g at 100 GeV

having to drop datasets altogether due to **ill-defined covariance matrices** to implementing **ad***hoc* decorrelation models and the impossibility to account for constraints from search data

> PDF interpretations of the HL-LHC data may become seriously hampered, or even impossible altogether, unless experiments release their full statistical models

> > minimal modification of 0.93correlation model, large 1.15impact in fit quality, 1.73PDFs stable 0.920.92

Juan Rojo on PDF fits

Key component of predictions for particle,

> nd astro-particle nts.

fundamental in QCD.

Publication of Statistical Models Hands-on workshop 8-12 Nov 2021 https://indico.cern.ch/event/1088121/





gaussian)

Availability of open likelihoods would make possible PDF determinations where statistical and systematic uncertainties are always accounted for by means of the appropriate model

$$\sigma_{\text{syst,gauss,ij}} \pm \sum \sigma_{\text{theo,ij}} \pm \sum \sigma_{\text{model,ij}}$$

$$syst \, error \qquad theory \, error \qquad model \, error$$

$$(correlated, \qquad (gaussian? \qquad (correlated,))$$

uniform?)

(correlated, not gaussian)

$$\mathscr{U}_{\text{poiss}} \times \mathscr{U}_{\text{uniform}} \times \mathscr{U}_{\text{theory}} \times \dots$$

This is not only technically correct, but also allows including much more information on PDF

	atistical Models
0.92	Hands-on workshop 8-12 Nov 20 https://indico.cern.ch/event/1088121

Juan Rojo on PDF fits

ent of or particle, astro-particle

damental QCD.



021

(Profile) likelihoods: very useful but not sufficient

- In the likelihood, the data is baked in
 - cannot evaluate likelihood on new data
 - cannot sample from the model (pseudo-data, toy MC)
- In profile likelihoods, nuisance parameter are fixed
 - cannot statistically combine profile likelihoods targeting parameters of interest if they share nuisances
 - cannot update constraint terms (auxiliary measurements)
- - risk of introducing dependencies which can result in a loss of information



Reparametrization in terms of different parameters of interest is not always possible

parametrization in terms of quantities such as masses, cross sections, widths, branching fractions, etc., is often more useful than a parametrization in terms of theory-model (Lagrangian) parameters

Full statistical model

The complete probability model for the analysis; includes dependence on the data x,y, the parameters of interest µ and nuisance parameters θ, access to the individual terms and the ability to generate pseudo-data ("toy Monte Carlo").

$$p(n, x, y | \mu, \theta) = \prod_{i=1}^{N_c} \left[\text{Pois}(n_i \mid \nu_i(\mu, \theta)) \prod_{j=1}^{n_i} p_i(x_{ij} | \mu, \theta) \right] p(y | \theta) \rightarrow L(\mu, \theta)$$
Likelihood: The value of the statistical model for a given fixed dataset as a function of the parameters
$$pdf \text{ of auxiliary data } y$$

$$pdf \text{ of auxiliary data } y$$

$$pdf \text{ of auxiliary data } y$$

$$p_i(x_{ij} | \mu, \theta) = \sum_k \frac{\nu_{ik}(\mu, \theta)}{\nu_i(\mu, \theta)} p_{ik}(x_{ij} | \mu, \theta), \quad \text{The probability to measure } x_{ij}$$

$$p_i(x_{ij} | \mu, \theta) = \sum_k \frac{\nu_{ik}(\mu, \theta)}{\nu_i(\mu, \theta)} p_{ik}(x_{ij} | \mu, \theta), \quad \text{The probability to measure } x_{ij}$$

$$p_i(\mu, \theta) = \sum_k \nu_{ik}(\mu, \theta)$$

$$p_i(\mu, \theta) = \sum_k \nu_{ik}(\mu, \theta)$$

Access to the types of rein

- changing process w considered
- updates o calculation

Data Science Basics • DESY & Uni Hamburg • 25 Oct 2022



ATLAS took the leap

... and started to publish plain-text serialisation of HistFactory workspaces in JSON format

Provides background estimates, changes under systematic variations, and observed data counts at the same fidelity as used in the experiment.

	Description	Modification	Constraint Term c_{χ}
ed	Uncorrelated Shape	$\kappa_{scb}(\gamma_b) = \gamma_b$	$\prod_{b} \operatorname{Pois} \left(r_{b} = \sigma_{b}^{-2} \middle \rho_{b} = \sigma_{b}^{-2} \gamma_{b} \right)$
constraine	Correlated Shape	$\Delta_{scb}(\alpha) = f_p\left(\alpha \middle \Delta_{scb,\alpha=-1}, \Delta_{scb,\alpha=1}\right)$	Gaus $(a = 0 \alpha, \sigma = 1)$
	Normalisation Unc.	$\kappa_{scb}(\alpha) = g_p\left(\alpha \middle \kappa_{scb,\alpha=-1}, \kappa_{scb,\alpha=1}\right)$	Gaus $(a = 0 \alpha, \sigma = 1)$
	MC Stat. Uncertainty	$\kappa_{scb}(\gamma_b) = \gamma_b$	$\prod_{b} \operatorname{Gaus} \left(a_{\gamma_{b}} = 1 \middle \gamma_{b}, \delta_{b} \right)$
	Luminosity	$\kappa_{scb}(\lambda) = \lambda$	Gaus $(l = \lambda_0 \lambda, \sigma_\lambda)$
free	Normalisation Data-driven Shape	$\kappa_{scb}(\mu_b) = \mu_b$ $\kappa_{scb}(\gamma_b) = \gamma_b$	
	Duite arriven bhape		

Rate modifications defined in HistFactory for bin b, sample s, channel c.

- Usage: RooFit, pyhf
- Target: long-term data/analysis preservation, reinterpretation purposes

Input
σ_b
$\Delta_{scb,\alpha=\pm 1}$
$\kappa_{scb,\alpha=\pm 1}$ $\delta_b^2 = \sum_s \delta_{sb}^2$ $\lambda_0, \sigma_\lambda$



ATLAS PUB Note ATL-PHYS-PUB-2019-029 21st October 2019



Reproducing searches for new physics with the ATLAS experiment through publication of full statistical likelihoods

The ATLAS Collaboration

The ATLAS Collaboration is starting to publicly provide likelihoods associated with statistical fits used in searches for new physics on HEPData. These likelihoods adhere to a specification first defined by the HistFactory p.d.f. template. This note introduces a JSON schema that fully describes the HistFactory statistical model and is sufficient to reproduce key results from published ATLAS analyses. This is per-se independent of its implementation in ROOT and it can be used to run statistical analysis outside of the ROOT and RooStats/RooFit framework. The first of these likelihoods published on HEPData is from a search for bottom-squark pair production. Using two independent implementations of the model, one in ROOT and one in pure Python, the limits on the bottom-squark mass are reproduced, underscoring the implementation independence and long-term viability of the archived data.



ATLAS took the leap

... and started to publish plain-text serialisation of HistFactory workspaces in JSON format

Provides background estimates, changes under systematic variations, and observed data counts at the same fidelity as used in the experiment.

	Description	Modification	Constraint Term c_{χ}
ed	Uncorrelated Shape	$\kappa_{scb}(\gamma_b) = \gamma_b$	$\prod_{b} \operatorname{Pois} \left(r_{b} = \sigma_{b}^{-2} \middle \rho_{b} = \sigma_{b}^{-2} \gamma_{b} \right)$
constraine	Correlated Shape	$\Delta_{scb}(\alpha) = f_p\left(\alpha \middle \Delta_{scb,\alpha=-1}, \Delta_{scb,\alpha=1}\right)$	Gaus $(a = 0 \alpha, \sigma = 1)$
	Normalisation Unc.	$\kappa_{scb}(\alpha) = g_p\left(\alpha \middle \kappa_{scb,\alpha=-1}, \kappa_{scb,\alpha=1}\right)$	Gaus $(a = 0 \alpha, \sigma = 1)$
	MC Stat. Uncertainty	$\kappa_{scb}(\gamma_b) = \gamma_b$	$\prod_{b} \operatorname{Gaus} \left(a_{\gamma_{b}} = 1 \middle \gamma_{b}, \delta_{b} \right)$
	Luminosity	$\kappa_{scb}(\lambda) = \lambda$	Gaus $(l = \lambda_0 \lambda, \sigma_\lambda)$
free	Normalisation Data-driven Shape	$\kappa_{scb}(\mu_b) = \mu_b$ $\kappa_{scb}(\gamma_b) = \gamma_b$	
	Duite arriven bhape		

Rate modifications defined in HistFactory for bin b, sample s, channel c.

- Usage: RooFit, **pyhf**
- Target: long-term data/analysis preservation, reinterpretation purposes

Input
σ_b
$\Delta_{scb,\alpha=\pm 1}$
$ \begin{aligned} \kappa_{scb,\alpha=\pm 1} \\ \delta_b^2 &= \sum_s \delta_{sb}^2 \\ \lambda_0, \sigma_\lambda \end{aligned} $

HEPData	Resources				
HistFactory File					
Archive of full likelihoods in the HistFactory JSON format described in ATL-PHYS-PUB-2019-029 Provided are 3 statiscal models labeled RegionA RegionB and RegionC respectively each in their own sub-directory. For each model the background-only model is found i the file named 'BkgOnly.json' For each model a set of patches for various signal points is provided					
10.17182/hepdata.89408.v3/r2					
Down	load				



ATLAS took the leap

... and started to publish plain-text serialisation of HistFactory workspaces in JSON format

Provides background estimates, changes under systematic variations, and observed data counts at the same fidelity as used in the experiment.

	Description	Modification	Constraint Term c_{χ}
p	Uncorrelated Shape	$\kappa_{scb}(\gamma_b) = \gamma_b$	$\prod_{b} \operatorname{Pois} \left(r_{b} = \sigma_{b}^{-2} \right \rho_{b} = \sigma_{b}^{-2} \gamma_{b} \right)$
constraine	Correlated Shape	$\Delta_{scb}(\alpha) = f_p\left(\alpha \middle \Delta_{scb,\alpha=-1}, \Delta_{scb,\alpha=1}\right)$	Gaus $(a = 0 \alpha, \sigma = 1)$
	Normalisation Unc.	$\kappa_{scb}(\alpha) = g_p\left(\alpha \middle \kappa_{scb,\alpha=-1}, \kappa_{scb,\alpha=1}\right)$	Gaus $(a = 0 \alpha, \sigma = 1)$
	MC Stat. Uncertainty	$\kappa_{scb}(\gamma_b) = \gamma_b$	$\prod_{b} \operatorname{Gaus} \left(a_{\gamma_{b}} = 1 \middle \gamma_{b}, \delta_{b} \right)$
	Luminosity	$\kappa_{scb}(\lambda) = \lambda$	Gaus $(l = \lambda_0 \lambda, \sigma_\lambda)$
e	Normalisation	$\kappa_{scb}(\mu_b) = \mu_b$	
fre	Data-driven Shape	$\kappa_{scb}(\gamma_b) = \gamma_b$	

Rate modifications defined in HistFactory for bin b, sample s, channel c.

- Usage: RooFit, **pyhf**
- Target: long-term data/analysis preservation, reinterpretation purposes

ATL-PHYS-PUB-2019-029

Input
σ_b
$\Delta_{scb,\alpha=\pm 1}$
$\kappa_{scb,\alpha=\pm 1} \\ \delta_b^2 = \sum_s \delta_{sb}^2 \\ \lambda_0, \sigma_\lambda$

	Search for charginos and neutralinos in all-hadronic final states	SUSY	Accepted by PRD	17-AUG-21	1
	4-top xsec measurement	TOPQ	Accepted by JHEP	22-JUN-21	1
	Search for gluinos, stops and electroweakinos in RPV models in final states with 1L and many jets	SUSY	Accepted by EPJC	17-JUN-21	1
	Search for charginos and neutralinos in final states with 3L and MET	SUSY	Accepted by EPJC	03-JUN-21	1
	Measurement of ttZ cross sections in Run 2	TOPQ	Eur. Phys. J. C 81 (2021) 737	23-MAR-21	1
	Search for third-generation scalar leptoquarks decaying to a top quark and a tau lepton	EXOT	JHEP 06 (2021) 179	27-JAN-21	1
	Search for squarks and gluinos in final states 1L, jets and MET	SUSY	Eur. Phys. J. C 81 (2021) 600	05-JAN-21	1
	Search for charginos and neutralinos in RPV models in final states with 3L (or more)	SUSY	Phys. Rev. D 103, (2021) 112003	20-NOV-20	1
	Search for displaced leptons	SUSY	Phys. Rev. Lett. 127 (2021) 051802	13-NOV-20	1
	Search for squarks and gluinos in final states with 0L, jets and MET	SUSY	JHEP 02 (2021) 143	27-OCT-20	1
	Measurement of the ttbar production cross-section in the lepton+jets channel at 13 TeV	TOPQ	Phys. Lett. B 810 (2020) 135797	24-JUN-20	1
	Stop pair, long-lived; displaced vertex and displaced muon	SUSY	Phys. Rev. D 102 (2020) 032006	26-MAR-20	1
	Chargino-neutralino pair; 3 leptons, weak-scale mass splittings	SUSY	Phys. Rev. D 101 (2020) 072001	18-DEC-19	1
	Chargino-neutralino pair, slepton pair; soft leptons	SUSY	Phys. Rev. D 101 (2020) 052005	28-NOV-19	1
а	Staus; taus	SUSY	Phys. Rev. D 101 (2020) 032009	15-NOV-19	1
	Chargino-neutralino pair; Higgs boson in final state, 2 b-jets and 1 lepton	SUSY	Eur. Phys. J. C 80 (2020) 691	19-SEP-19	1
	Stop pair, sbottom pair, gluino pair; two same-sign leptons or three leptons	SUSY	JHEP 06 (2020) 46	18-SEP-19	1
_	Sbottom; b-jets	SUSY	JHEP 12 (2019) 060	08-AUG-19	1

HEPData

Resources

Like	elihood	available	•
PRD	17-AUG-21	13	139 fb ⁻¹
	22-JUN-21	13	139 fb ⁻¹
	17-JUN-21	13	139 fb ⁻¹
	03-JUN-21	13	139 fb ⁻¹
<u>C 81</u>	23-MAR-21	13	139 fb ⁻¹
21)	27-JAN-21	13	139 fb ⁻¹
<u>C 81</u>	05-JAN-21	13	139 fb ⁻¹
<u>103,</u> 3	20-NOV-20	13	139 fb ⁻¹
ett.	13-NOV-20	13	139 fb ⁻¹
21)	27-OCT-20	13	139 fb ⁻¹
810 97	24-JUN-20	13	139 fb ⁻¹
102 06	26-MAR-20	13	136 fb ⁻¹
101)1	18-DEC-19	13	139 fb ⁻¹
101)5	28-NOV-19	13	139 fb ⁻¹
101)9	15-NOV-19	13	139 fb ⁻¹
C 80	19-SEP-19	13	139 fb ⁻¹
20)	18-SEP-19	13	139 fb ⁻¹
<u>19)</u>	08-AUG-19	13	139 fb ⁻¹

ATLAS full statistical models HistFactory JSON format



G. Alguero, J. Araz, B. Fuks, SK arXiv:2206.14870

> The <u>Simplify</u> python tool can be used to create simplified statistical models from full ones by merging all background contributions and combining all nuisance parameters into a single one; may yield equivalent results at much lower CPU cost — needs testing case-by-case!

SC-SC Section USING States States

\rightarrow statistical evaluation through JSON patching



Illustration by Lukas Heinrich Hands-on workshop 8 Nov 2021

Interfaced to pyhf since SModelS v1.2.4 (now v2.2) G. Alguero, SK, W. Waltenberger, arXiv:2009.01809



Improvements w.r.t. "best signal region" approach



G. Alguero, SK, W. Waltenberger, arXiv:2009.01809





Beyond the Standard Model (BSM) searches at the LHC

BSM searches are typically performed in **specific bins of** kinematic distributions, so-called signal regions (SRs)

- designed to maximise the number of events from the hypothesised signal with respect to the number of "background" events originating from Standard Model processes.
- control and validation regions are defined in the phase space where no or very little signal from new physics is expected.
- A statistical analysis is then performed to evaluate the confidence level of the hypothesised BSM scenario, and claim evidence for or set a limit on the new particles of this scenario.



Beyond the Standard Model (BSM) searches at the LHC

BSM searches are typically performed in **specific bins of** kinematic distributions, so-called signal regions (SRs)

- designed to maximise the number of events from the hypothesised signal with respect to the number of "background" events originating from Standard Model processes.
- control and validation regions are defined in the phase space where no or very little signal from new physics is expected.
- A statistical analysis is then performed to evaluate the confidence level of the hypothesised BSM scenario, and claim evidence for or set a limit on the new particles of this scenario.



Beyond the Standard Model (BSM) searches at the LHC

BSM searches are typically performed in **specific bins of** kinematic distributions, so-called signal regions (SRs)

- designed to maximise the number of events from the hypothesised signal with respect to the number of "background" events originating from Standard Model processes.
- control and validation regions are defined in the phase space where no or very little signal from new physics is expected.
- A statistical analysis is then performed to evaluate the confidence level of the hypothesised BSM scenario, and claim evidence for or set a limit on the new particles of this scenario.

The statistical combination of disjoint SRs in reinterpretation studies uses more of the data of an analysis and gives more robust results than the single (best) SR approach.





CMS: covariance matrices

 Best-SR approach: assuming a Poisson distribution for the data and a Gaussian with variance of δ^2 for the nuisances, p(θ)

$$\mathcal{L}(\mu,\theta) = \frac{(\mu s + b + \theta)^{n_{\text{obs}}} e^{-(\mu s + b + \theta)}}{n_{\text{obs}}!} exp\left(-\frac{\theta}{2}\right)$$

 CMS analyses sometimes provide a covariance matrix, which allow for the combination of disjoint SRs in a **simplified likelihood approach**

$$\mathcal{L}_S(\mu, \theta) = \prod_{i=1}^N \frac{(\mu \cdot s_i + b_i + \theta_i)^{n_i} e^{-(\mu \cdot s_i + b_i + \theta_i)}}{n_i!} \cdot \exp\left(-\frac{1}{n_i!}\right)$$

[CMS NOTE-2017/001]

Implemented in SModelS and GAMBIT since a while; recently also in MadAnalysis 5

Much(!) better than best-SR, but caveat are non-Gaussian effects e.g. when systematic uncertainties dominate





$$\delta^2 = \delta_s^2 + \delta_b^2$$
signal+background uncertainties



covariance matrix



CMS-SUS-16-039





CMS: covariance matrices

• Best-SR approach: assuming a Poisson distribution for the data and a Gaussian with variance of δ^2 for the nuisances, p(θ)

$$\mathcal{L}(\mu,\theta) = \frac{(\mu s + b + \theta)^{n_{\text{obs}}} e^{-(\mu s + b + \theta)}}{n_{\text{obs}}!} exp\left(-\frac{\theta}{2}\right)$$

• CMS analyses sometimes provide a covariance matrix, which allows for the combination of disjoint SRs in a **simplified likelihood approach**

$$\mathcal{L}_S(\mu, \theta) = \prod_{i=1}^N \frac{(\mu \cdot s_i + b_i + \theta_i)^{n_i} e^{-(\mu \cdot s_i + b_i + \theta_i)}}{n_i!} \cdot \exp\left(-\frac{1}{n_i!}\right)$$

[CMS NOTE-2017/001]

covariance matrix

Implemented in SModelS and GAMBIT since a while; recently also in MadAnalysis 5

Much(!) better than best-SR, but caveat are non-Gaussian effects e.g. when systematic unc. dominate

G. Alguero, J. Araz, B. Fuks, SK, arXiv:2206.14870







Data Science Basics • DESY & Uni Hamburg • 25 Oct 2022







back to full statistical models

Simplify: from full to simplified likelihoods

Computing e.g. a CLs value from the full statistical model with tens to hundreds of nuisance parameters is accurate but too CPU intensive for some use cases \rightarrow can we simplify this?

The <u>Simplify</u> python tool creates simplified statistical models from full ones by merging all background contributions and combining all nuisance parameters into a single one; same HistFactory JSON format



much lower CPU cost but not always a good approximation





Machine-learn likelihoods

- Machine-learn likelihoods from full statistical models as functions of signal counts in each SRs
- Neural networks: sequential multilayer perceptrons using TensorFlow2.
 - loss function: mean squared error
 - activation function: Leaky Relu with alpha=0.2.
- Example for ATLAS-SUSY-2019-08 (9 SRs):
 - 3 layers, 512 neurons, MAPE = 0.2504, Max PE: 5.429
- The NN models are
 - saved in ONNX (Open Neural Network Exchange) format using the tf2onnx converter,
 - then interfaced from SModelS using ONNX Runtime









ATLAS-SUSY-2019-08 -- mape: 0.2504





Machine-learn likelihoods

- Machine-learn likelihoods from full statistical models as functions of signal counts in each SRs
- Neural networks: sequential multilayer perceptrons using TensorFlow2.
 - loss function: mean squared error
 - activation function: Leaky Relu with alpha=0.2.
- Example for ATLAS-SUSY-2019-08 (9 SRs):
 - 3 layers, 512 neurons, MAPE = 0.2504, Max PE: 5.429
- The NN models are
 - saved in ONNX (Open Neural Network Exchange) format using the tf2onnx converter,
 - then interfaced from SModelS using ONNX Runtime











ATLAS-SUSY-2019-08 (efficiencyMap)







Analysis preservation and reuse

"Data and Analysis Preservation, Recasting and Reinterpretation" S. Bailey et al., arXiv:2203.10057

Analysis preservation



Full preservation

Exact software chain used to perform the analysis in the experiment.

- Full post-generation software stack from detector simulation and reconstruction to the physics analysis code.
- Complicated by the diversity of analysis software frameworks, even within a given experiment \rightarrow container images
- High computing power requirements on any large-scale reuse



Preservation of analysis logic and workflows enabling the reuse of the original analysis process and associated data products.

Lightweight preservation





Snowmass white paper on data and analysis preservation and reinterpretation

Analysis Preservation Recommendations

- **3.1:** Ensure use of interoperable systems to maximise the preservability and reusablility of experiment simulation and analysis software chains. This includes the use of version control, archival systems, containerisation, common software interfaces and data formats, and commitments from experimental collaborations and their host laboratories to maintain documentation and provide long-term support.
- **3.2:** Ensure that all operational and in-preparation experiments have a planned and resourced programme for capture and long-term reproduction of their complete computational processing chain, including validation regression-tests.
- **3.3:** Ensure commu process to maximise analysis impact. umenta

for community consumption.

3.4: Support continuing development and uptake of new technologies for increasingly framework-independent analysis specifications, such as via declarative domainspecific analysis description languages.

S. Bailey et al., arXiv:2203.10057

Ensure that release of analysis preservation logic via public frameworks for the community to use is integrated with experiment publication and data-release processes,



Lightweight, public

Simulation-based reinterpretation ("recasting")

- Aims at reproducing experimental analyses in Monte Carlo simulation
- Nowadays well established for traditional cut-based analyses. Information needed: cf. arXiv:2003.07868

object definitions; identification, tagging, reconstruction efficiencies

detailed preselection and signal (+control) region cuts

However, more and more analyses exploit ML techniques to gain in sensitivity

e.g. ML-based taggers, signal/bkg discrimination with ML classifiers

Pb: how can we reuse those?

workflow



* except for detector-unfolded results (Rivet/Contur)



ML as a bottleneck for reinterpretation?



Conventional cut-based analysis. All needed information is provided, recasting works very well

Data Science Basics • DESY & Uni Hamburg • 25 Oct 2022

Illustrative example from G. Alguero, J. Araz, B. Fuks, SK, arXiv:2206.14870



Analysis employs a BDT for tau tagging; resulting efficiencies given only approximately \rightarrow serious differences in the recasting (final weights ~30% too high)

MAD Analysis 5



ML as a bottleneck for reinterpretation?

- More and more analyses exploit ML techniques to gain in sensitivity.
- Serious difficulty for analysis preservation and reuse unless
 - resulting id/reco efficiencies can be (and are!) parametrised in terms of quantities accessible in a simulation, e.g., p_t , η , ...
 - the actual ML model is published in appropriate form. Caveat: input variables need to be physics quantities that can be matched in a simulation

Two analyses where the latter has been attempted:

ATLAS-SUSY-2018-22 (0-lepton gluino/squark search) published BDT weights as XML file

→ RAMP seminar by Kenta Uno

ATLAS-SUSY-2019-04 (1-2 leptons + jets RPV search) published neural network as **ONNX** file

→ RAMP seminar by Javier Montejo Berlingen

RAMP: Reinterpretation Auxiliary Material Presentation



13 different algorithms: image-based (2), 4-vector-based (5), theory-inspired (6) taggers

"The Machine Learning Landscape of Top Taggers" G. Kasieczka, T. Plehn et al., arXiv:1902.09914



Kenta Uno on ATLAS-SUSY-2018-22





Javier Montejo Berlingen on ATLAS-SUSY-2019-04



HEPData Q Sear	Additional Publication Re	esources		Ab K
Publication Information ch for R-parity violating ersymmetry in a final state aining leptons and many je ATLAS experiment using \sqrt{s} proton-proton collision data	 ♥ filter Common Resources Distribution: 1 Distribution: 2 Distribution: 3 Distribution: 4 	Charlen External Link web page with auxiliary material View Resource	C++ File Code snippet with the implementation of analysis selection at truth-level Download	of the eren L whi
eorges , Abbott, Braden Keim , Abbott, D. dam , Abeling, Kira , Abhayasinghe, Des Haider , Abramowicz, Halina , Abreu, He i, Yiming EP-2021-066, 2021. /dear good ata.104860 (Resources) ct (data abstract) .HC.	Distribution: 52Distribution: 62Distribution: 72Distribution: 82Distribution: 92Distribution: 102Distribution: 112Distribution: 122Distribution: 132	tgz File SLHA files for benchmark signals Download	tgz File ONNX files for the neural networks for th analysis Download	Electroweak produ 6 or 8 jets, 4 b-tag p
pton and either zero or at least three <i>b</i> -tag ted. The search uses 139 fb ^{-1} of $\sqrt{s} = 13$ collision data collected by the ATLAS expe Run 2 of the Large Hadron Collider. The re eted in the context of R-parity-violating	Distribution: 14 2	pos	-2741.3 22698	$ \begin{array}{c} \tilde{\chi}_{1}^{0} \\ \tilde{\chi}_{323}^{+} \\ p \\ \tilde{\chi}_{1}^{\pm} \\ p \\ p \\ p \\ p \\ \tilde{\chi}_{2}^{0} \\ \tilde{\chi}_{2}^{0} \\ p \\ \tilde{\chi}_{2}^{0} \\ \tilde$





Snowmass white paper on data and analysis preservation and reinterpretation

Reinterpretation and Recasting Recommendations 5.1: Encourage that reinterpretability and reuse be kept in mind early on in the ar

- \mathbf{te}
- nc
- th
- 5.2: De

ucts, such as statistical models, with reinterpretation use-cases in mind.

- 5.3: Improve the coordination among the different public reinterpretation frameworks with the goal of a centralised database of recast codes, common input/output formats, and a unified statistical treatment.
- **5.4:** Encourage the FAIR-ification of codes and data products from (theory) reinterpretation studies outside the experimental collaborations at the same level of sophistication as asked for experimental analyses and results. Suitable repositories are, e.g., GitHub and Zenodo; appropriate versioning is essential.

S. Bailey et al., arXiv:2203.10057

Encourage that reinterpretability and reuse be kept in mind early on in the analysis design. This concerns, for instance, the **choice of input parameters in ML models**, the full specification of the fiducial phase space of a measurement in terms of the final state, including any vetos applied, and generally the **choice of non-overlapping regions** and standard naming of shared nuisances to facilitate the combination of analyses.





Some more comments on reproducibility and reuse of ML models

Open Neural Network Exchange

"defines all the necessary operations a machine learning model needs to implement its inference function"

- ONNX is an open format built to represent ML models.
 - aims at providing a common language any ML framework can use to describe its models.
 - makes it possible to deploy a model independent from the learning framework used to build it.

The deployment of a ML model usually requires replicating the entire ecosystem used to train the model, most of the time with a *docker*. Once a model is converted into ONNX, the production environment only needs a runtime (C, java, python, javascript,) to execute the graph defined with ONNX operators.

- Converters exist for scikit-learn, tensorflow, pytorch, and others NB must be updated every time ONNX or the library they support have a new released version.
- Beware of custom layers, experimental features, etc.! may be troublesome for converter and/or runtime (interpreter)

For pertinent reuse, input variables must be clearly documented

Runtime (interpreter) must match ONXX version \rightarrow possible issue for preservation?





Ltwnn

Lightweight Trained Neural Network

build passing coverity passed DOI 10.5281/zenodo.597221

What is this?

The code comes in two parts:

- A set of scripts to convert saved neural networks to a standard JSON format
- 2. A set of classes which reconstruct the neural network for application in a C++ production environment

The main design principles are:

- Minimal dependencies: The C++ code depends on C++11, Eigen, and boost PropertyTree. The converters have additional requirements (Python3 and h5py) but these can be run outside the C++ production environment.
- Easy to extend: Should cover 95% of deep network architectures we would realistically consider.
- Hard to break: The NN constructor checks the input NN for consistency and fails loudly if anything goes wrong.

We also include converters from several popular formats to the lwtnn JSON format. Currently the following formats are supported:

- Scikit Learn
- Keras (most popular, see below)

David Hohn > lwtnn



"Our underlying assumption is that training and inference happen in very different environments: we assume that the training environment is flexible enough to support modern and frequently-changing libraries, and that the inference environment is much less flexible."







Reproducibility of ML models

To reproduce a ML model, one needs

The precise ecosystem used to train the model

- Tools and their versions -----
- Exact architecture (layers, neurons, ...) -
- Activation and loss functions
- Anything else specifying the algorithm
- The training and validation data
- The **initial conditions** (random seeds) for the weights

"A learning algorithm can be viewed as searching a space H of hypotheses to identify the best hypothesis in the space."



Thomas G. Dietterich Ensemble Methods in Machine Learning

Data Science Basics • DESY & Uni Hamburg • 25 Oct 2022



Effect of different initialisation \rightarrow ensemble methods





Conclusions

- * The data and analyses from particle physics experiments are unique and of immense scientific value.
- Impact can go much beyond original paper publication.
- Proper preservation, enabling long-term future reuse, maximises the scientific return.
- It is worthwhile to keep reinterpretability and reuse in mind early on in analysis design.



Next Reinterpretation-Forum workshop

(Re)interpretation of the LHC results for new physics

12–15 Dec 2022 CERN

Europe/Zurich timezone

Overview

Timetable

Registration

Call for Abstracts

Participant List

topics of this workshop will be: i) the publication and reuse of statistical models, iii) global analyses and global fits.

Continuing the conversation from the last workshop session, we would like to include general best practices for reinterpretation/reuse of experimental results beyond the LHC, and particularly welcome contributions regarding results from precision or astrophysical experiments.

Enter your search term

м.

This is the 7th general workshop of the "Forum on the interpretation of the LHC results for BSM studies", or LHC Reinterpretation Forum (RIF) for short. Its aim is to review new developments on the tools, pheno, and the experimental sides, and to prepare for the Run 3 results of the LHC. In this context, major

- ii) the reinterpretation of analyses that employ machine learning, and

https://indico.cern.ch/event/1197680/





Backup

CMS simplified likelihood

$$\mathcal{L}(\mu, \theta) = \mathcal{P}(\text{data}|\mu \cdot s(\theta) + b(\theta)) \cdot p(\tilde{\theta}|\theta)$$

Some simplifying assumptions must be made to reduce the complexity of the full probability density function $p(\mathbf{q} | \mathbf{q})$.

- background events is symmetric about the expectation, b, and its variance is independent of q. Often, the each region is sufficient to approximate $p(\vec{q} | q)$ at least for values of q which are close to \vec{q} .
- region.
- contributions.

$$\mathcal{L}_{S}(\mu, \theta) = \prod_{i=1}^{N} \frac{(\mu \cdot s_{i} + b_{i} + \theta_{i})^{n_{i}} e^{-(\mu \cdot s_{i} + b_{i} + \theta_{i})}}{n_{i}!} \cdot \exp\left(-\frac{1}{2}\theta^{T} \mathbf{V}^{-1}\theta\right)$$

The constraints on the background contributions are Gaussian such that the distribution of the number of background contributions are estimated from control regions in data with large sample sizes, which makes this assumption valid. The covariance, and therefore only the linear correlation, between the background contribution in

The numbers of events, *n*, are statistically independent from one another. This is true when there are no events which are included in more than one search region and the estimates of the background contributions, b_i, and covariance matrix V have not been obtained from data which are statistically dependent on the data from any search

The systematic uncertainties in the signal model can be neglected. The validity of this assumption will strongly depend on the specific BSM physics model being considered. Systematic uncertainties on the signal can be accounted for by adding appropriate nuisance parameters with Gaussian constraints as for the background



Reinterpretation: two approaches

Reproduce experimental analysis in a Monte Carlo simulation

"recasting"



Sabine Kraml

Reuse simplified model results $(\sigma_{95}, \text{ signal } A \times \varepsilon)$

Assumes that A×ɛ doesn't change too much w.r.t. original model

Test of **BSM hypothesis**



Data Science Basics • DESY & Uni Hamburg • 25 Oct 2022



Workflow

simulation-based recasting



* except for detector-unfolded results (Rivet/Contur)



[‡] in case exp. result is σ_{95} : only allowed/excluded



Pro's and con's

simulation-based recasting



- More generic and often more precise than simplified model results; in principle applicable to any new signal caveat: control regions typically not included in react codes!
- Need to take care to **simulate all relevant processes** (not always obvious e.g. in scans of complex parameters spaces where dominant processes can change)

• Very CPU expensive

- So far only cut-and-count analyses are recasted
- ATLAS / CMS as well as Run1 (8 Tev) / Run 2 (13 TeV) analyses need to be **run separately**
- So far, prompt and long-lived signatures need to be treated separately
 - \rightarrow careful separation needed in models featuring both
 - \rightarrow response of prompt analyses to LLPs unclear / wrong
- Implementation and validation of new analyses is timeconsuming and sometimes quite difficult

→ Detailed information needed from experiment analysis logic, object definitions, cuts, efficiencies, cut-flows, etc.



Pro's and con's

 Assumes that signal acceptances are to good approximation the same as in original experimental result.

Valid for **simple rescaling** of production and decay rates ($\sigma \times BR$); other cases need to be **verified**, e.g. spin or production mode dependence.

- Applicable beyond cut & count analyses (ML techniques)
- Advantages are simplicity and speed!
 - → very fast b/c no MC simulation needed
 - \rightarrow well suited for large scans and model surveys
- Large database of experimental results
- ATLAS and CMS, Run1 and Run2, prompt and longlived results all treated simultaneously
- Easy classification of unconstrained cross section, missing topologies
- Often conservative: coverage depends on variety of available simplified-model results

simplified model approach (SModelS)



[‡] in case exp. result is σ_{95} : only allowed/excluded

