



Data science in particle physics

Lydia Brenner

Trieste, 16-18 November 2020

Introduction

Challenge in particle physics:

Handling of large amounts of data





Big data and the standard model



Some examples

Particle	Expected number of events in ATLAS in completed Run 2
Higgs Boson	7.7 million
Top quark	275 million
Z Boson	2.8 billion
W Boson	12 billion

Big data processing

The data flow from all four experiments for Run 2 is about 25 GB/s

- ALICE: 4 GB/s (Pb-Pb running)
- ATLAS: 800 MB/s
- CMS: 600 MB/s
- LHCb: 750 MB/s

Data selection in two stages

- 1) Hardware filter reduces 600 million events per second to 100 000 per second
- 2) Software filter uses algorithms to reduce to 100-200 events per second

End up with around 1050 MB/s of raw data

Grid computing

Using the Worldwide LHC Computing Grid (WLCG) to

- store,
- distribute
- process

More than 170 collaborating centres in 36 countries around the world, linked to CERN



Measurement roadmap



Particle identification

Event reconstruction based on

- Detector response
- Object classification
- Rejection of background processes

On real and simulated data



Detector response

Event reconstruction considering insensitive areas in the detector due to

- Readout electronics
- Detector dead-time
- Detector efficiencies

On real and simulated data



Background rejection



Remove predicted contribution in each bin

Fit analytical function to subtract

160

myy [GeV]

Background rejection

Use control regions to correct for simulation mismatch

Needs region that is close to the signal region in terms of physics model, without signal



Fit analytical function to subtract

Particle discovery



Particle discovery



Limit setting





Precision measurements

Uncertainty bands

ATLAS+CMS Preliminary LHC <i>top</i> WG	m _{top} summary, √ s = 7-13 TeV	September 2018
World comb. (Mar 2014) [2] stat	total stat	
total uncertainty	m _{top} ± total (stat ± syst)	s Ref.
LHC comb. (Sep 2013) LHCtopWG	173.29 \pm 0.95 (0.35 \pm 0.88)	7 TeV [1]
World comb. (Mar 2014)	173.34 \pm 0.76 (0.36 \pm 0.67)	1.96-7 TeV [2]
ATLAS, I+jets	172.33 ± 1.27 (0.75 ± 1.02)	7 TeV [3]
ATLAS, dilepton	173.79 ± 1.41 (0.54 ± 1.30)	7 TeV [3]
ATLAS, all jets	175.1±1.8 (1.4±1.2)	7 TeV [4]
ATLAS, single top	172.2 ± 2.1 (0.7 ± 2.0)	8 TeV [5]
ATLAS, dilepton	$172.99 \pm 0.85 \; (0.41 \pm 0.74)$	8 TeV [6]
ATLAS, all jets	$173.72 \pm 1.15 \ (0.55 \pm 1.01)$	8 TeV [7]
ATLAS, I+jets	172.08 \pm 0.91 (0.38 \pm 0.82)	8 TeV [8]
ATLAS comb. (Sep 2017) H=H	172.51 \pm 0.50 (0.27 \pm 0.42)	7+8 TeV [8]
CMS, I+jets	173.49 ± 1.06 (0.43 ± 0.97)	7 TeV [9]
CMS, dilepton	172.50 ± 1.52 (0.43 ± 1.46)	7 TeV [10]
CMS, all jets	173.49 ± 1.41 (0.69 ± 1.23)	7 TeV [11]
CMS, I+jets	172.35 \pm 0.51 (0.16 \pm 0.48)	8 TeV [12]
CMS, dilepton	172.82 ± 1.23 (0.19 ± 1.22)	8 TeV [12]
CMS, all jets	$172.32 \pm 0.64 \; (0.25 \pm 0.59)$	8 TeV [12]
CMS, single top	172.95 ± 1.22 (0.77 ± 0.95)	8 TeV [13]
CMS comb. (Sep 2015)	172.44 \pm 0.48 (0.13 \pm 0.47)	7+8 TeV [12]
CMS, I+jets	$172.25 \pm 0.63 \; (0.08 \pm 0.62)$	13 TeV [14]
CMS, all jets	172.34 ± 0.79 (0.20 ± 0.76)	13 TeV [15]
ען דען נון אלא פון גרשה 19 גרשה 19 גרדה (בן	SCONF-2013-102 [6] Phys.Lett.10251 (2016) 350 Phys.J.C 75 (2015) 330 [9] ATU-Ke CONF-2017-071 Phys.J.C 75 (2015) 158 [9] HFP 12 (2012) 106 AS-CONF-2014-065 [10] Eur. Phys.J.C72 (2012) 2202	[11] Eur.Phys.J. 074 (2014) 2758 [12] Phys.Rev.D33 (2016) 072004 [13] EPJC 77 (2017) 354 [14] arXiv:1805.01428 [15] CMS PAS TOP-17-008
165 170 175	5 180	185
m _{top}	[GeV]	1 Martin Dat



Construction of the likelihood

Start with the formulation of a likelihood: $L(\vec{x} \mid \vec{\mu}, \vec{\theta})$

- (B)SM physics model * Soft physics model * Detector description * Analysis reconstruction



Interpolation of uncertainties

Additional terms in the likelihood for modelling uncertainties only known for certain points in parameter phase-space



Fitting sum of analytical functions

Create a template based on physics model

- Create summed function for fitting
- Allow variations within modelling uncertainties



Histogramming

Monte carlo simulation created from physics model

- Often impossible to create analytical function
- Simplify into counting-experiment by binning the distribution
- Single events saved to allow binning in the analysis step



Classification

Use simulation to identify regions of interest based on

- Physics variables
- Phase-space regions

Multivariate analysis using simple cuts or machine learning



Machine learning

Used to differentiate between signal and background, different models, model detector responses or object reconstruction and particle identification





Alternative measurement roadmap: Unfolding

Detector reconstruction

22



The future of data analysis at the LHC



Amount of simulated data needed for the next run is too big

- Clever solutions needed for data manipulation

Moving from statistically limited measurements, to measurements limited by modelling of the physics model and the detector response







The standard model of particle physics

Hadron building blocks



Computing languages used

Special data-processing framework: ROOT

C++ for performance

More recently: python

Previously: Fortran

Others used: Java, Perl, html

Not used: R -> memory issues due to too large datasets

Morphing based on physics model

