Combine overview

Kyle Cormier (UZH), Aliya Nigamova (UHH) and Nicholas Wardle (IC)

Terascale Statistics School | 02/04/2024











Introduction

- Combine the most popular tool for statistical inference in lacksquareCMS
- Command-line interface to RooStats and RooFit methods, \bullet and even more:
 - Builds statistical model
 - Runs statistical test
 - Provides tool set for validation
- GitHub page: https://github.com/cms-analysis/ \bullet HiggsAnalysis-CombinedLimit









Analysis types: input

Template based model:

- When no simple analytic form can be assigned for the signal(background) description
- Templates (histograms) are provided for all signal and background processes and shape systematic uncertainties
- MC statistical uncertainties are modelled with Barlow-Beeston-lite approach (autoMCStats)







log scale



Analysis types: input





Kyle Cormier, Aliya Nigamova, Nick Wardle | Terascale Statistics School | 02/04/2024

Parametric model:

- Suitable for the analysis with analytically described bkg and signal: e.g. Gaussian signal on smoothly falling polynomial background
- In most cases used in the analyses with datadriven bkg description
- The systematic uncertainties assigned on the parameters of the model



Analysis types: input

Template based model:

- When no simple analytic form can be assigned for the signal(background) description
- Templates (histograms) are provided for all signal and background processes and shape systematic uncertainties
- MC statistical uncertainties are modelled with Barlow-Beeston-lite approach (autoMCStats)

Counting experiment:

signal, bkg and data yields



Kyle Cormier, Aliya Nigamova, Nick Wardle | Terascale Statistics School | 02/04/2024

Parametric model:

- Suitable for the analysis with analytically described bkg and signal: e.g. Gaussian signal on smoothly falling polynomial background
- Data-driven bkg description
- The systematic uncertainties assigned on the parameters of the model

• "One-bin analysis", i.e. the input parameters in the model are the



Likelihood function

Likelihood is a probability to observe data given parameters Φ , which we want to measure $\mathcal{L}_{\mathcal{M}}(\vec{\Phi}) = p_{\mathcal{M}}(\text{data}; \vec{\Phi})$

Often we factorize the likelihood into primary and auxiliary components

$$\mathcal{L} = \mathcal{L}_{ ext{primary}} \cdot \mathcal{L}_{ ext{auxiliary}}$$

If we have several independent measurements $\vec{x_d}$ (histogram bins) the total likelihood constructed with product of individual likelihoods

$$\mathscr{L} = \prod_{k} \mathscr{L}_{k}(\mathbf{x}, \mu, \vec{\nu})$$

Same applies for individual independent categories (each with multiple bins)



Kyle Cormier, Aliya Nigamova, Nick Wardle | Terascale Statistics School | 02/04/2024

 $\vec{\nu}$ auxiliary observables, from external measurement: E.g. systematic uncertainties from exp. corrections (b-tagging, reconstruction efficiency); theory uncertainties



Binned likelihood example

$$\mathcal{L} = \mathcal{L}_{ ext{primary}} \cdot \mathcal{L}_{ ext{auxiliary}} =$$

C runs over categories

$$= \prod_{c=1}^{N_c} \prod_{b=1}^{N_b^c} \text{Poiss}(n_{cb}; n_{cb}^{\exp}(\vec{\mu}, \vec{\nu})) \cdot \prod_{e=1}^{N_E} p_e(q)$$
Product over bins in a histogram
$$n_{cb}^{\exp} = \max(0, \sum_p M_{cp}(\vec{\mu}) N_{cp}(\nu_G, \vec{\nu}_L, \vec{\mu}))$$



Kyle Cormier, Aliya Nigamova, Nick Wardle | Terascale Statistics School | 02/04/2024







Constraint terms

When we want to model systematic uncertainties given by external measurements we introduce constraint terms:

Uncertainty type	Directive	Inputs	Multiplicative factor, $f(\nu)$	p(y; v)	Default values	
Log-normal	lnN	kappa	$\kappa^{ u}$	$\mathcal{N}(y;\nu,1)$	$\nu = y = 0$	-
Asymmetric log-normal	lnN	kappaDown, kappaUp	$(\kappa^{\text{Down}})^{-\nu}$ if $\nu < -0.5$, $(\kappa^{\text{Up}})^{\nu}$ if $\nu > 0.5$, $e^{\nu K(\kappa^{\text{Down}},\kappa^{\text{Up}},\nu)}$ otherwise.*	$\mathcal{N}(y;\nu,1)$	$\nu = y = 0$	$\kappa^{\mathrm{Up/Down}} = \frac{\sum_{b}}{\sum}$
Log-uniform	lnU	kappa	$\kappa^{ u}$	$\mathcal{U}\left(y,1/\kappa,\kappa ight)$	$\nu = y = \frac{1}{2} \left(\kappa + \frac{1}{\kappa} \right)$	
Gamma	gmN	N,alpha [†]	ν/N	$\mathcal{P}(y; \nu)$	$\nu = N + 1, y = N^{\ddagger}$	



Kyle Cormier, Aliya Nigamova, Nick Wardle | Terascale Statistics School | 02/04/2024

$$\mathcal{L} = \prod_{k} \mathcal{L}_{k}(\vec{\mathbf{x}}, \mu, \vec{\nu}) \prod_{j} \mathbf{p}_{j}(\nu)$$

With j running over NP





InN uncertainties

Uncertainty type	Directive	Inputs	Multiplicative factor, $f(\nu)$	<i>p</i> (
Log-normal	lnN	kappa	$\kappa^{ u}$	$\mathcal{N}(\mathfrak{z})$
Asymmetric log-normal	lnN	kappaDown, kappaUp	$(\kappa^{\text{Down}})^{-\nu}$ if $\nu < -0.5$, $(\kappa^{\text{Up}})^{\nu}$ if $\nu > 0.5$, $e^{\nu K (\kappa^{\text{Down}}, \kappa^{\text{Up}}, \nu)}$ otherwise.*	$\mathcal{N}(\mathfrak{z})$
Log-uniform	lnU	kappa	$\kappa^{ u}$	U (у,
Gamma	gmN	N,alpha [†]	ν/N	$\mathcal{P}($

$$p_e \propto \exp\left(-0.5 \left(rac{(
u_e-y_e)}{\sigma}
ight)^2
ight); \ p_e = ext{Poiss}(
u_e;y_e); \ ext{or} \ p_e \propto ext{constant}$$



Kyle Cormier, Aliya Nigamova, Nick Wardle | Terascale Statistics School | 02/04/2024





Likelihood definition

Complete picture

[More details here]





Interpolation between up and down variation

Kyle Cormier, Aliya Nigamova, Nick Wardle | Terascale Statistics School | 02/04/2024



$$\begin{array}{c} \text{Exp. events per bin} \\ \text{Poiss}(n_{cb}; n_{cb}^{\text{exp}}(\vec{\mu}, \vec{\nu})) \\ \text{Poiss}(n_{cb}; n_{cb}^{\text{exp}}(\vec{\mu}, \vec{\nu})) \\ \text{e=1} \\ \end{array} \\ \begin{array}{c} \text{Poiss}(n_{cb}; n_{cb}^{\text{exp}}(\vec{\mu}, \vec{\nu})) \\ \text{e=1} \\ \end{array} \\ \begin{array}{c} \text{P}_{e}(y_{e}; \nu_{e}) \\ \text{P}_{e}(y_{e}; \nu_{e}) \\ \text{P}_{e}(\vec{\mu}, \vec{\nu}, \nu_{e}) \\ \text{P}_{e}(\vec{\mu}, \vec{\nu}, \nu_{e}) \\ \text{P}_{e}(\vec{\mu}, \vec{\nu}, \nu_{e}) \\ \end{array} \\ \begin{array}{c} \text{P}_{e}(\vec{\mu}, \vec{\nu}, \nu_{e}) \\ \end{array} \\ \begin{array}{c} \text{P}_{e}(\vec{\mu}, \vec{\nu}, \nu_{e}) \\ \end{array} \\ \begin{array}{c} \text{P}_{e}(\vec{\mu}, \vec{\nu}, \nu_{e}) \\ \text{P}_{e}(\vec{\mu}, \vec{\nu}, \nu_{e}) \\ \text{P}_{e}(\vec{\mu}, \vec{\nu}, \nu_{e}) \\ \text{P}_{e}(\vec{\mu}, \vec{\nu}, \nu_{e}) \\ \end{array} \\ \begin{array}{c} \text{P}_{e}(\vec{\mu}, \vec{\nu}, \nu_{e}) \\ \text{P}_{e}(\vec{\mu}, \vec{\nu}, \nu_{e}) \\ \text{P}_{e}(\vec{\mu}, \vec{\nu}, \nu_{e}) \\ \end{array} \\ \begin{array}{c} \text{P}_{e}(\vec{\mu}, \vec{\nu}, \nu_{e}) \\ \text{P}_{e}(\vec{\mu}, \vec{\nu}, \nu_{e}) \\ \text{P}_{e}(\vec{\mu}, \vec{\nu}, \nu_{e}) \\ \end{array} \\ \begin{array}{c} \text{P}_{e}(\vec{\mu}, \vec{\nu}, \nu_{e}) \\ \text{P}_{e}(\vec{\mu}, \vec{\nu}, \nu_{e}) \\ \end{array} \\ \begin{array}{c} \text{P}_{e}(\vec{\mu}, \vec{\nu}, \nu_{e}) \\ \text{P}_{e}(\vec{\mu}, \vec{\nu}, \nu_{e}) \\ \end{array} \\ \begin{array}{c} \text{P}_{e}(\vec{\mu}, \vec{\nu}, \nu_{e}) \\ \text{P}_{e}(\vec{\mu}, \vec{\nu}, \nu_{e}) \\ \end{array} \\ \begin{array}{c} \text{P}_{e}(\vec{\mu}, \nu_{e}) \\ \end{array} \\ \end{array} \\ \begin{array}{c} \text{P}_{e}(\vec{\mu}, \nu_{e}) \\ \end{array} \\ \end{array} \\ \begin{array}{c} \text{P}_{e}(\vec{\mu}, \nu_{e}) \\ \end{array} \\ \begin{array}{c} \text{P}_{e}(\vec{\mu}, \nu_{e}) \\ \end{array} \\ \begin{array}{c} \text{P}_{e}(\vec{\mu}, \nu_{e}) \\ \end{array} \\ \end{array} \\ \begin{array}{c} \text{P}_{e}(\vec{\mu}, \nu_{e}) \\ \end{array} \\ \begin{array}{c} \text{P}_{e}(\vec{\mu}, \nu_{e}) \\ \end{array} \\ \begin{array}{c} \text{P}_{e}$$







Datacard structure



so divide the unit gaussian by 2 before doing the interpolation

11

Documentation

Combine v9.2.0 -

Introduction



These pages document the RooStats / RooFit - based software tool used for statistical analysis within the CMS experiment - COMBINE. Note that while this tool was originally developed in the Higgs Physics Analysis Group (PAG), its usage is now widespread within CMS.

COMBINE provides a command-line interface to many different statistical techniques, available inside RooFit/RooStats, that are used widely inside CMS.

The package exists on GitHub under https://github.com/cms-analysis/HiggsAnalysis-CombinedLimit

For more information about Git, GitHub and its usage in CMS, see http://cms-sw.github.io/cmssw/faq.html

The code can be checked out from GitHub and compiled on top of a CMSSW release that includes a recent RooFit/RooStats, or via standalone compilation without CMSSW dependencies. See the instructions for installation



Kyle Cormier, Aliya Nigamova, Nick Wardle | Terascale Statistics School | 02/04/2024

https://cms-analysis.github.io/HiggsAnalysis-CombinedLimit/latest/

Q Search

 istheta
 GitHub

 Sv9.2.0
 ☆ 68
 ¥ 369

Table of contents

Installation instructions

Within CMSSW (recommended for CMS users)

Combine v9 - recommended version

Combine v8: CMSSW_10_2_X release series

SLC6/CC7 release CMSSW_8_1_X

Oustide of CMSSW (recommended for non-CMS users)

Standalone compilation

Compilation on lxplus9

Standalone compilation with LCG

Standalone compilation with conda

Standalone compilation with CernVM

What has changed between tags? For developers

CombineHarvester/CombineTo...



Documentation



- Installation (w/ CMSSW, standalone, using LCG, conda, latest release) [link] \bullet \bullet
- Setting up the analysis (counting, template based, parametric) [link]
- Running Combine [link] \bullet
- Underlying statistics \bullet
 - Likelihood definition: [link] -
 - How the fits are performed (profiling, marginalization, confidence intervals): [link]
 - Statistical tests (test statistic, GOF): [link] -
- Tutorials: main features, parametric, unfolding, RooFit



https://cms-analysis.github.io/HiggsAnalysis-CombinedLimit/latest/



13

Main methods

- **Asymptotic** likelihood methods:
- AsymptoticLimits: limits calculated according to the asymptotic formulas in arxiv:1007.1727, valid for large event counts
- Significance: simple profile likelihood approximation for calculating significances.
- **Frequentist** or hybrid bayesian-frequentist methods:
 - HybridNew: compute modified frequentist limits with toys, significance/pvalues and confidence intervals with several options, --LHCmode LHC-limits is the recommended one
- **Bayesian** methods:
 - BayesianSimple: performing a classical numerical integration (for simple models only)
 - MarkovChainMC: performing Markov Chain integration (for arbitrarily complex models)



[Documentation]



 $egin{aligned} -2\log\left(rac{\mathcal{L}(\mu)}{\mathcal{L}(\mu=0)}
ight) & \hat{\mu} < 0 \ -2\log\left(rac{\mathcal{L}(\mu)}{\mathcal{L}(\hat{\mu})}
ight) & 0 < \hat{\mu} < \mu \ 0 & \mu < \hat{\mu} \end{aligned}$



Main methods

- Fitting
 - [MultiDimFit]: perform maximum likelihood fits with multiple POIs and likelihood scans
 - diagnostic tools

combine -M FitDiagnostics -d ws.root

- Other modules:
 - several GOF estimators (AD, KS, saturated)
 - are

 - [Impacts]: evaluate the shift in POI from $\pm \sigma_{\text{postfit}}$ variation for each NP



Kyle Cormier, Aliya Nigamova, Nick Wardle | Terascale Statistics School | 02/04/2024

combine -M MultiDimFit ---algo grid --points 50 -d ws.root - [FitDiagnostics]: performs maximum likelihood fits to extract the signal yield and provides

- [GoodnessOfFit]: perform a goodness of fit test for models including shape information using

- [ChannelCompatibiltyCheck]: check how consistent the individual channels of your analysis

- [GenerateOnly]: generate random or asimov toy datasets for use as input to other methods [Documentation]



Profiled likelihood



In a measurement where no need to measure the intervals for NP $\vec{\nu}$

→ Profiling: find a $\hat{\vec{\nu}}(\mu)$ which maximise the likelihood for each value POI



Kyle Cormier, Aliya Nigamova, Nick Wardle | Terascale Statistics School | 02/04/2024





FitDiagnostics

- **Provides more information than** –M MultiDimFit
- Runs background only fit first (r=0 fixed), followed by s+b (r is floating). combine -M FitDiagnostics -d ws.root; output: fitDiagnostics.root
 - Provides full list of NP constraints and pulls for both fits
 - Covariance matrix is saved, access to all correlations -
 - Using fit results from fitDiagnostic.root one can check the NP shifts and uncertainties wrt their input values:

python diffNuisances.py fitDiagnostics.root ---all [instructions]

- Post(pre)-fit shapes can be saved with option --saveShapes, additional directories inside output file will be created for each category (can only be used when covariance matrix is properly estimated [see the debugging session in the afternoon])



Kyle Cormier, Aliya Nigamova, Nick Wardle | Terascale Statistics School | 02/04/2024

[Documentaion]

- Uncertainty on the measurement (r) should be estimated from full likelihood
- **scans** combine -M MultiDimFit ---algo grid --points 50 -d ws.root ...



17

Generating toys in combine

- lacksquarebe generated.
- \bullet
- lacksquareother methods (-t -1 or -t N)
- 2 options for systematic uncertainties: lacksquare

--toysNoSystematics nominal (prefit) NP values are used to generated toys. --toysFrequentist first the fit to data is performed, then the nuisance parameters in each toy are set to their post-fit values, with POIs fixed (eg. --setParameters r=1), before generating the data. The constraint terms are randomized within their Gaussian constraint pdfs around the post-fit NP values.

To read the toy dataset: -toysFile=higgsCombineTest.GenerateOnly.root \bullet



To compute expected significance, intervals, perform bias tests the toy experiments have to

Adding -t -1, for Asimov dataset, -t N for N random datasets, options to any method

The method –M GenerateOnly –-saveToys allows to generate toys that can be later used with



Batch submission



Works with many methods, e.g. with -M MultiDimFit to produce likelihood scans:

combineTool.py -d ws.root -M MultiDimFit --algo grid --points 50 --rMin 0 --rMax 1 --jobmode condor —split—points 10 ——sub—opts='+JobFlavour="workday"' ——task—name mytask —n mytask



Kyle Cormier, Aliya Nigamova, Nick Wardle | Terascale Statistics School | 02/04/2024

[Detailed instructions]

Useful when generating toys, running likelihood scans, other tasks that can be parallelised

Batch submission options

Particularly useful for GoodnessOfFit, Impacts, limits with toys



19

Plotting with combineTools

- Likelihood scans:
 - plot1DScan.py: uses the output of MultiDimFit
 - --other option can be used to load the scans with systematic uncertainties fixed, useful when producing scans with stat. and stat.+syst. uncertainties, full uncertainty breakdown to separate the contribution of various sources [documentaion]







Plotting with combineTools

- Observable distributions:
 - FitDiagnostics using --shapes option [link]
 - With PostFitShapesFromWorkspace
 - For post-fit fit shapes the output of FitDiagnostic should be provided with -f fitDiagnostic.root ; --sampling Option should be used to take into account NP correlations



Kyle Cormier, Aliya Nigamova, Nick Wardle | Terascale Statistics School | 02/04/2024







log scale



Questions?

Tutorial

I. Follow [Combine setup instructions]

Getting started

We need to set up a new CMSSW area and checkout the combine package:

```
source /cvmfs/cms.cern.ch/cmsset_default.sh
cmsrel CMSSW_11_3_4
cd CMSSW_11_3_4/src
cmsenv
git clone https://github.com/cms-analysis/HiggsAnalysis-CombinedLimit.git HiggsAnalysis/CombinedLimit
cd HiggsAnalysis/CombinedLimit
cd $CMSSW_BASE/src/HiggsAnalysis/CombinedLimit
git fetch origin
git checkout v9.2.0
cd $CMSSW_BASE/src/
```

We will also make use another package, CombineHarvester, which contains some high-level tools for working with combine. The following command will download the repository and checkout just the parts of it we need for this tutorial:

```
bash <(curl -s https://raw.githubusercontent.com/cms-analysis/CombineHarvester/main/CombineTools/scripts/sparse-checkout-https.sh)
```

Now make sure the CMSSW area is compiled:

scramv1 b clean; scramv1 b

Finally, move to the working directory for this tutorial, which contains all the inputs needed to run the exercises:

cd \$CMSSW_BASE/src/HiggsAnalysis/CombinedLimit/data/tutorials/longexercise/



Kyle Cormier, Aliya Nigamova, Nick Wardle | Terascale Statistics School | 02/04/2024

රී

II. Move to the first exercise here

Main Features of Combine (Long Exercises)

This exercise is designed to give a broad overview of the tools available for statistical analysis in CMS using the combine tool. COMBINE is a high-level tool for building RooFit / RooStats models and running common statistical methods. We will cover the typical aspects of setting up an analysis and producing the results, as well as look at ways in which we can diagnose issues and get a deeper understanding of the statistical model. This is a long exercise - expect to spend some time on it especially if you are new to COMBINE. If you get stuck while working through this exercise or have questions specifically about the exercise, you can ask them on this mattermost channel. Finally, we also provide some solutions to some of the questions that are asked as part of the exercise. These are available here.





Backup

Template morphing





Kyle Cormier, Aliya Nigamova, Nick Wardle | Terascale Statistics School | 02/04/2024

[details]



Model building: Physics Model

- Simple process scaling:

 - Multiple POIs can be assigned with [multiSignalModel]:

text2workspace.py -P HiggsAnalysis.CombinedLimit.PhysicsModel:multiSignalModel ---PO verbose --P0 'map=.*/sig1:r_sig1[1,0,10]' --P0 'map=.*/sig2:r_sig2[1,0,20]' datacard.txt -o ws.root —

- Mapping: ---P0 'map=bin/process:parameter'
- More complicated: Interference, κ , EFT (signal processes are parametrized) \bullet
 - Use existing model [location]:

text2workspace.py datacard -P HiggsAnalysis.CombinedLimit.PythonFile:modelName e.q for Higgs couplings κ-model: -P HiggsAnalysis.CombinedLimit.HiggsCouplings:c7 [link]

interference: [tutorial], [code]



Kyle Cormier, Aliya Nigamova, Nick Wardle | Terascale Statistics School | 02/04/2024

By default single POI assigned to all processes marked as signal in the datacards ("r")

Create your own model in CombinedLimit/python directory; example given here on a model with



Model building: datacards (CombineHarvester)

C++ package with [python interface] to create and modify datacards

- 1.Analysis categories ch::Categories
- 2.Signal and background processes ch::CombineHarvester::AddProcesses 3.Systematic uncertainties ch::CombineHarvester::AddSyst 4.Extract the related shape inputs from ROOT files ch::CombineHarvester::ExtractShapes
- 5. The input shapes/yields can be modified:
 - modifying signal processes to different cross sections (ch::Process::set rate), change types (signal/ background)
- (de)correlating systematic uncertainties (renaming) ch::CombineHarvester::RenameSystematic 6.Exporting to the text datacard format and creating the associated ROOT shape file(s): ch::CardWriter



[Details]

Examples: [Example 1] [Example 2] [Example 2 in python]



Model building: datacards (CombineHarvester)

- C++ package with [python interface] to create and modify datacards
- **1**.Analysis categories ch::Categories
- 2.Signal and background processes ch::CombineHarvester::AddProcesses
- 3.Systematic uncertainties ch::CombineHarvester::AddSyst
- 4.Extract the related shape inputs from ROOT files ch::CombineHarvester::ExtractShapes
- 5. The input shapes/yields can be modified:
 - background)

Parse already existing datacard and apply necessary changes: ch::CombineHarvester::ParseDatacard



[Details]

- modifying signal processes to different cross sections (ch::Process::set_rate), change types (signal/

- (de)correlating systematic uncertainties (renaming) ch::CombineHarvester::RenameSystematic 6.Exporting to the text datacard format and creating the associated ROOT shape file(s): ch::CardWriter

Examples: [Example 1] [Example 2]

[Example 2 in python]



