

# Heraverager

## Data combination package



Heraverager developers

August 15, 2016

### **Abstract**

The presented stand-alone tool is designed to combine the data of the measurements. The combination process is based on the  $\chi^2$  minimization with respect of the correlation model, assigned to uncertainties of the data points. The program performs a study of behaviour of the data-point uncertainties and tests the compatibility of the measurements. The package is available for downloads on the web-site <https://wiki-zeuthen.desy.de/HERAverager>

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Combination method</b>	<b>4</b>
2.1	Basic principal of the combination . . . . .	4
2.2	$\chi^2$ minimization . . . . .	6
2.3	Covariance matrix representation of the systematic uncertainties . . . . .	8
2.4	Bias correction for multiplicative uncertainties . . . . .	9
2.5	Bias correction for statistical uncertainties . . . . .	10
2.6	Treatment of off-set systematic uncertainties . . . . .	10
2.7	Treatment of asymmetric systematic uncertainties . . . . .	11
2.8	Different representation of combined uncertainties . . . . .	11
2.9	Combination with bias corrections . . . . .	11
<b>3</b>	<b>Program manual</b>	<b>12</b>
3.1	Program installation . . . . .	12
3.2	Code organization . . . . .	12
3.3	Input information . . . . .	13
3.3.1	Steering file . . . . .	13
3.3.2	Data file . . . . .	17
3.4	Output information . . . . .	18
<b>4</b>	<b>Examples</b>	<b>20</b>
<b>5</b>	<b>Python scripts</b>	<b>21</b>
5.0.1	Data reader . . . . .	21
5.0.2	Generation of random dataset . . . . .	22
5.1	Python-based averager . . . . .	23
5.2	Python-based plotting . . . . .	24

# 1 Introduction

If two or more measurements of the same physics quantity are statistically independent they can be combined to an average measurement. When these measurements have only one uncertainty and no information about correlation is known weighted-average algorithm can be used to calculate average value. However in many cases measurements are binned and performed with several sources of the correlated and uncorrelated systematic uncertainties. Presented package is designed in order to study the compatibility and perform bias-free combination of these measurements.

The combination procedure is based on the  $\chi^2$  minimization. This package was originally developed for combination of two data sets of the H1 data [3] and then used for combination of data from two HERA experiments: H1 and ZEUS [2]. Currently this package is widely used for combination of LHC data.

Following types of the data-uncertainties are considered in presented tool:

- Statistical uncertainties. Usually based on square root of number of measured events. Uncorrelated between data sets and between bins.
- Systematic uncertainties uncorrelated between bins. Can be correlated or uncorrelated between data sets.
- Bin-to-bin correlated systematic uncertainties. Can be correlated or uncorrelated between data sets.
  - Additive: not proportional to the measured values.
  - Multiplicative: proportional to the measured values.
  - Off-set: does not have any contribution to the  $\chi^2$  and therefore does not have any impact on the average value.

Tool support asymmetric uncertainties.

This manual describes a linear  $\chi^2$  definitions (see Sec. 2.1) and explains the basic linear-algebra manipulations in Sec. 2.2, which stays behind the combination procedure. Several different biases are discussed in Secs. 2.4–2.5. Bias-correction procedure introduces non-linear  $\chi^2$ , which minimized using linear approximation and iterative procedure is shown in Sec. 2.9.

User manual describes the installation procedure (Sec. 3.1) and code description (Sec. 3.2). The input and output of the combination are discussed in Secs. 3.3 and 3.4 respectively.

## 2 Combination method

### 2.1 Basic principal of the combination

For a measurement  $\mu$  with uncertainty  $\Delta$ , assuming a Gaussian shape of the uncertainty, the measurement can be considered as a probability distribution function for a “true” quantity  $m$ :

$$P(m) = \frac{1}{\sqrt{2\pi}\Delta} \exp\left(-\frac{(m-\mu)^2}{2\Delta^2}\right) \quad (1)$$

This can be written as a  $\chi^2$  function by taking  $-2 \log$  (constant term was skipped):

$$\chi^2(m) = \frac{(m-\mu)^2}{\Delta^2} \quad (2)$$

Minimum of  $\chi^2$  corresponds to  $m = \mu$ , while the change  $\Delta\chi^2 = 1$  corresponds to  $m = \mu \pm \Delta$ . In case of two statistically independent measurements of the case quantity  $m$ :  $\mu_1, \Delta_1$  and  $\mu_2, \Delta_2$ , the probability distribution function of  $m$  is given by the product of two:

$$P(m) \sim \exp\left(-\frac{(m-\mu_1)^2}{2\Delta_1^2}\right) \exp\left(-\frac{(m-\mu_2)^2}{2\Delta_2^2}\right), \quad (3)$$

which corresponds to  $\chi^2$  that is given by the sum of the two:  $\chi_{sum}^2 = \chi_1^2 + \chi_2^2$ .

Since  $\chi_{sum}^2$  is a positive definite quadratic form it can be re-written in the form of Eq. 2. In this case  $\mu$  is replaced by average  $\mu_{ave}$  and  $\Delta$  is replaced by the error on this average:

$$\chi^2(m) = \frac{(m-\mu_{ave})^2}{\Delta_{ave}^2} + \chi_0^2, \quad (4)$$

where the value of  $\chi_0^2$  measures consistency of the measurements,  $\chi^2/N_{DoF} \sim 1$  for consistent measurements.

The value of  $\mu_{ave}$  can be found by minimizing  $\chi_{sum}^2$  with respect to  $m$  (this leads to a usual averaging rule,  $1/\Delta^2$  weights).

Many experiments measures a number of independent quantities  $\mu_i$  which correspond to the underlying physics values  $m_i$  (e.g. cross-section measurement in bins of  $p_{T,Z}$ , where  $i$  refers to a bin number). In this case the  $\chi^2$  function is a simple sum over the measurements (bins):

$$\chi_{exp}^2(m_i) = \sum_i \frac{(m_i - \mu_i)^2}{\Delta_i^2} + \chi_0^2, \quad (5)$$

Where:

- $\mu_i$  - the measurement in the bin  $i$
- $m_i$  “truth” value in the bin  $i$ <sup>1</sup>
- $\Delta_i$  statistical uncertainty in bin  $i$

---

<sup>1</sup>Just to avoid confusion: “truth” here does not refers to MC-truth. This refers to a value, which we want to measure in experiment.

The systematic effects, which affects the measurement  $\mu_i$ , are often correlated across bins. Let's consider measurement binned in a certain variable, which is affected by up/down shift of certain parameter:

$$\mu_i \rightarrow \mu_i + \Gamma_i^+, \quad \mu_i \rightarrow \mu_i - \Gamma_i^-, \quad (6)$$

where  $\Gamma_i^\pm$  correspond to the variation up/down.

If the correlated systematic uncertainty is approximately symmetric, one can symmetrize them. For presented analysis following way was used:

$$\Gamma_i = \max(|\Gamma_i^+|, |\Gamma_i^-|) \frac{\Gamma_i^+}{|\Gamma_i^+|}, \quad (7)$$

i.e. the size of the uncertainty is taken as maximal of up and down variations and the sign from one of the variations.

The symmetrised correlated systematic uncertainties were included into the  $\chi^2$  function.

Systematic uncertainties, like energy scale, can be also viewed as a result of an experiment (e.g. measurement of the calibration): there is a "true" detector energy scale  $\alpha$ , measured detector calibration  $\alpha_0$  and its statistical uncertainty  $\Delta_\alpha$ . Therefore, it is natural to add term

$$\chi_{sys}^2(\alpha) = \frac{(\alpha - \alpha_0)^2}{\Delta_\alpha^2} \quad (8)$$

to the  $\chi^2$  function. The nuisance parameter  $b$ , defined as  $b = (\alpha - \alpha_0)/\Delta_\alpha$  corresponds to a coherent change of measurements  $\mu_i \rightarrow \mu_i + b\Gamma_i$ . This defines the combined  $\chi^2$  function:

$$\chi^2(\vec{m}, \vec{b})_{exp} = \sum_i \frac{(m_i - \mu_i - \sum_j \Gamma_i^j b_j)^2}{\Delta_i^2} + \sum_j b_j^2, \quad (9)$$

where

- $\vec{b}$  defines a vector of nuisance parameters  $b_j$  corresponding to each source of systematic uncertainty,
- summation over  $i$  runs over all data points, and summation over  $j$  runs over all correlated sources of systematic uncertainty,
- $\Gamma_i^j$  is the absolute correlated systematic uncertainty,
- $\Delta_i$  is the uncorrelated (statistical) uncertainty.

With this definition minimum  $\chi^2$  is obtained for all  $m_i = \mu_i$  and  $b_j = 0$ . If  $b_j = 0$  for all  $j$  except  $j = k$ ,  $b_k = 1$ , then  $\chi^2$  minimum is archived at  $m_i = \mu_i + \Gamma_i^k$  and it is equal to 1.

Total uncertainty for a parameter  $m_i$  defined by  $\Delta\chi_{exp}^2 = 1$  rule corresponds to the sum of correlated and uncorrelated uncertainties in quadrature:  $\Delta_{i,tot}^2 = \Delta_i^2 + \sum (\Gamma_i^j)^2$ .

## 2.2 $\chi^2$ minimization

Average of two data sets with systematic uncertainties follows the same ideas as for average of uncorrelated measurements: represent the sum of two  $\chi^2$  by a single  $\chi^2$ :

$$\chi^2(\vec{m}_1, \vec{b}_1)_{exp,1} + \chi^2(\vec{m}_2, \vec{b}_2)_{exp,2} = \chi_0^2 + \chi^2(\vec{m}_{ave}, \vec{b}_{ave}). \quad (10)$$

The dimension of  $\vec{m}_{ave}$  is equal to dimension of union set of  $\vec{m}_1$  and  $\vec{m}_2$ . e.g. if both experiments measure for the same binning,  $N_{M1} = N_{M2} = N_{M,ave} = N_M$ . Similarly, for the systematic uncertainties  $N_{S,ave} = N_{S1} + N_{S2} - N_{S,common} = N_S$ , where  $N_{S,common}$  is the number of common systematic error sources for the two measurements.

More explicitly, the sum of two  $\chi^2$ :

$$\chi^2(\vec{m}, \vec{b})_{sum} = \sum_e \sum_i \frac{(m_i - \mu_{i,e} - \sum_j^{N_S} \Gamma_{i,e}^j b_j)^2}{\Delta_{i,e}^2} W_{i,e} + \sum_j^{N_S} b_j^2, \quad (11)$$

where,

- $i$  runs over all measured points  $N_M$
- $j$  runs over all sources of systematic uncertainties  $N_S$
- symbol  $W_{i,e}$  is equal to 1 if data set  $e$  contributes to a measurement at the point  $i$ , otherwise it is 0.
- $\Gamma_{j,e}^i$  equals to 0 if the measurement  $i$  from the data set  $e$  is insensitive to the systematic source  $j$ .

This definition of  $\chi^2$  assumes that the data sets  $e$  are statistically uncorrelated. The systematic error sources  $b_j$ , however, may be either uncorrelated (separate sources) or correlated across data sets (different data sets sharing a common source).

Since  $\chi_{sum}^2$  is a quadratic form of  $\vec{m}$  and  $\vec{b}$ , it may be rearranged such that it takes a form similar to Eq. 4.

$$\begin{aligned} \chi^2(\vec{m}, \vec{b}) = \chi_{min}^2 + \sum_i^{N_{M,ave}} \frac{(m_i - \mu_{i,ave} - \sum_j^{N_{S,ave}} \Gamma_{i,ave}^j (\alpha_j - \alpha_{j,ave}))^2}{\Delta_{i,ave}^2} + \\ + \sum_j^{N_{S,ave}} \sum_k^{N_{S,ave}} (\alpha_j - \alpha_{j,ave})(\alpha_k - \alpha_{k,ave})(A'_S)_{ik}, \end{aligned} \quad (12)$$

where

- $\mu_{i,ave}$  are average values of measured quantities
- $\Delta_{i,ave}$  are their uncorrelated uncertainties

The values of  $\alpha_{j,ave}$ ,  $\Delta_{i,ave}$ ,  $\mu_{i,ave}$  and matrix  $A'_S$  are determined by minimization of  $\chi^2$  function in Eq. 11 with respect to  $m_i$  and  $b_j$ . The minimum of Eq. 11 is found by solving a system of linear equations obtained by requiring  $\partial\chi^2/\partial m_i = 0$  and  $\partial\chi^2/\partial b_j = 0$  which can be written in matrix form

$$\begin{pmatrix} A_M & A_{SM} \\ (A_{SM})^T & A_S \end{pmatrix} \begin{pmatrix} M_{ave} \\ B_{ave} \end{pmatrix} = \begin{pmatrix} C_M \\ C_S \end{pmatrix} \quad (13)$$

where

- 82 • vector  $M_{ave}$  corresponds to all measurements
- 83 • vector  $B_{ave}$  corresponds to all sources of the systematic uncertainties
- 84 • matrix  $A_M$  has a diagonal structure with  $N_{M,ave}$  diagonal elements  $A_M^{ii} = \sum_e \frac{W_{i,e}}{\Delta_{i,e}^2}$
- 85 •  $A_{SM}^{ij} = -\sum_e \frac{\Gamma_{i,e}^j}{\Delta_{i,e}^2} W_{i,e}$
- 86 •  $A_S^{ij} = \delta_{ij} + \sum_e \sum_k^{N_M} \frac{\Gamma_{i,e}^k \Gamma_{j,e}^k}{\Delta_{k,e}^2} W_{k,e}$
- 87 •  $C_M^i = \sum_e \frac{\mu_e^i}{\Delta_{i,e}^2} W_{i,e}$
- 88 •  $C_S^j = -\sum_e \sum_k^{N_M} \frac{\mu_e^k \Gamma_{j,e}^k}{\Delta_{k,e}^2} W_{k,e}$

89 Here  $\delta_{ij}$  is the Kronecker symbol. The matrix  $A_{SM}$  has dimension  $N_M \times N_S$  while the matrix  $A_S$  is  
 90 quadratic with  $N_S \times N_S$  elements.

91 Using the method of the Schur complement, the solution is found as:

$$\begin{aligned} A'_S &= A_S - (A_{SM})^T A_M^{-1} A_{SM} \\ B_{ave} &= (A'_S)^{-1} (C_S - (A_{SM})^T A_M^{-1} C_M) \\ M_{ave} &= A_M^{-1} (C_M - A_{SM} B_{ave}) \end{aligned} \quad (14)$$

Given the components of the vector  $B_{ave}$ ,  $\beta_{j,ave} = \alpha_{j,ave} / \Delta_{\alpha_j}$ , the solution for  $\mu_{i,ave}$  can be written in explicit form:

$$\mu_{i,ave} = \frac{\sum_e \left( \mu_{i,e} + \sum_j \Gamma_{j,e}^i \beta_{j,ave} \right) \frac{W_{i,e}}{\Delta_{i,e}^2}}{\sum_e \frac{W_{i,e}}{\Delta_{i,e}^2}} \quad (15)$$

The uncorrelated uncertainty squared is determined by the inverse of the elements of the diagonal matrix  $A_M$ :

$$\Delta_{i,ave}^2 = \frac{1}{\sum_e \frac{W_{i,e}}{\Delta_{i,e}^2}} \quad (16)$$

92 Eq. 15 and 16 reproduce the standard formula for a statistically weighted average of several uncorre-  
 93 lated measurements when all shifts of the systematic error sources are set to zero. The values of  $\beta_{i,ave}$  in  
 94 Eq. 15 show, how the combined measurements  $\mu_{i,ave}$  are shifted, compared to initial measurements  $\mu_{i,e}$   
 95 in terms of systematic uncertainties  $\Gamma_{i,e}^j$ .

The non-diagonal nature of the matrix  $A'_S$  expresses the fact that the original sources of the systematic uncertainties are correlated with each other after averaging. The matrix  $A'_S$  can be decomposed to re-express Eq. 9 in terms of diagonalised sources of systematic uncertainties:

$$DD = U A'_S U^{-1} \quad \Gamma_{ave} = A_{SM} A_M^{-1} D^{-1} U^{-1} \quad (17)$$

96 Here  $U$  is an orthogonal matrix composed of the eigenvectors of  $A'_S$ ,  $D$  is a diagonal matrix with cor-  
 97 responding square roots of eigenvalues as diagonal elements and  $\Gamma_{ave}$  represents the sensitivity of the  
 98 average result to these new sources. Its elements are the  $\Gamma_{i,ave}^j$ .

99 After diagonalizability of matrix  $A'_S$ ,  $\chi^2$  function in Eq. 12 can be re-written in form, similar to Eq. 9:

$$\chi^2(\vec{m}, \vec{b}')_{tot} = \chi^2_{min} + \sum_i \frac{(m_i - \mu_{i,ave} - \sum_j^{N_s} \Gamma_{i,ave}^j b'_j)^2}{\Delta_{i,ave}^2} + \sum_j^{N_s} (b'_j)^2, \quad (18)$$

where  $b'_j = \sum_k U_{jk}(b_k - \beta_{k,ave})D_{jj}$ .

The orthogonal matrix  $U$  connecting the systematic sources before and after averaging with Eq. 17. Diagonal elements of matrix  $D$  shows, how the uncertainties of combined measurement  $\Gamma_{i,ave}^j$  are reduced, compared to initial systematic uncertainties.

The value of  $\chi^2_{min}$  corresponds to the minimum of Eq. 11 and calculated using values of  $\mu_{i,ave}$  and  $\beta_{j,ave}$  as a parameters  $\vec{m}$  and  $\vec{b}$ . The ratio  $\chi^2_{min}/N_{DoF}$  is a measure of the consistency of the data sets. The number of degrees of freedom,  $N_{DoF}$ , is calculated as the difference between the total number of measurements and the number of the measured points  $N_M$ . It is useful to note, the definition of  $\chi^2_{min}$  have two contributions, one is a usual shift if the measurement weighted with uncorrelated uncertainties (comes from  $\mu_{i,ave} - \mu_{i,e}$ ). Another contribution comes from correlated uncertainty term.

Another interesting parameter, which shown the compatibility of channels is pull of the central values:

$$p^{i,e} = \frac{\mu^{i,e} - \mu^{i,ave}(1 - \sum_j \gamma_{i,e}^j \beta_{j,ave})}{\sqrt{\Delta_{i,e}^2 - \Delta_{i,ave}^2}}, \quad (19)$$

where  $\gamma_{i,e}^j = \Gamma_{i,e}^j / \mu_{i,ave}$ . This definition is similar to the  $\chi^2$  definition, but not summed over bins. These pulls show how the average measurement are shifted compare to individual measurement and also have two contributions, similar to  $\chi^2$ .

The values, which reflect only correlated part are the shifts  $\beta$ . If the systematic uncertainties for measurement have rather similar size, than average fluctuation of shifts will reflect the correlated contribution to  $\chi^2$  compatibility. However we can always add such uncertainties to analysis, which will be not shifted after combination and therefore will just make smaller average fluctuation of shifts. Large shifts indicates, that corresponding central value of the systematic uncertainties were initially not correctly estimated (and therefore shifted during the combination).

The values coming out of matrix  $D$  show, how much the initial systematic uncertainties were reduced due to the combination. This parameter does not directly related to the channel compatibility, but show how much we gain out of the combination.

The pull for systematic uncertainties can be defined as:

$$p^i = \frac{\beta_{i,ave}}{\sqrt{1 - D_{ii}^2}}. \quad (20)$$

This value shown, how significant was the systematic shifted due to the combination. The large systematic pull suggests, that systematic uncertainty was not correctly estimated.

### 2.3 Covariance matrix representation of the systematic uncertainties

Another way representing bin-to-bin (point-to-point) correlations is by using a covariance matrix  $C$ :

$$\chi^2(\vec{m}) = \sum_{ik}^{N_M} (m_i - \mu_i)^T C_{ik}^{-1} (m_k - \mu_k), \quad C_{ik} = \sum_j^{N_s} \Gamma_i^j \Gamma_k^j, \quad (21)$$



where  $N_M$  is a number of bins and  $N_S$  is a number of sources of the systematic uncertainties. For Gaussian uncertainties covariance matrix and nuisance parameter representations are equivalent.

Matrix  $C$  can be written as:

$$C_{ik} = \sum_{lj}^{N_M} G_{il}^{-1} D_{lj} G_{jk}, \quad (22)$$

where columns of  $G^{-1}$  are made of eigenvectors of  $C$ , sorted by eigenvalues (largest to smallest) and  $D$  is a diagonal matrix made of the eigenvalues of  $C$ .

Since  $C$  is positively defined, the eigenvalues are real and greater zero and we can assume  $G' = \sqrt{D}G$ . Also  $G^{-1} = G'^T$ , since  $G$  is the orthogonal and then

$$C_{ik} = \sum_j^{N_M} G'_{ij}{}^T G'_{jk}. \quad (23)$$

The contribution of eigenvectors with small eigenvalues can be neglected by truncating the sum after  $N'_S < N_S$ .

$$C_{ik} \approx C'_{ik} = \delta_{ik} \Delta_{i,uncorr}^2 + \sum_{j=1}^{N'_S} G'_{ij}{}^T G'_{jk}, \quad \Delta_{i,uncorr}^2 = \sum_{j=N'_S+1}^{N_M} (G'_{ij})^2. \quad (24)$$

Here  $\delta_{ik}$  is the Kronecker symbol and  $\Delta_{i,uncorr}$  is uncorrelated systematic uncertainty. To preserve the total uncertainty,  $\Delta_{i,uncorr}$  are chosen such, that diagonal elements in  $C'$  are equal to the diagonal elements in  $C$ . The dimension  $N_M$  of covariance matrix  $C'$  is not reduced by this approximation.

The reduced summation allows for more compact representation using nuisance parameters. The  $\chi^2$  function takes form:

$$\chi^2(\vec{m}, \vec{b})_{exp} = \sum_i \frac{(m_i - \mu_i - \sum_j \Gamma_i^j b_j)^2}{\Delta_{i,stat}^2 + \Delta_{i,uncorr}^2} + \sum_j b_j^2, \quad (25)$$

## 2.4 Bias correction for multiplicative uncertainties

Many of the systematic uncertainties for the data measurements, correlated and uncorrelated are multiplicative, e.g. they are proportional to the measured values.

Consider two measurements  $\mu_1$  and  $\mu_2$  of  $m$ . Let's assume, that  $\mu_1 = m + mb$ ,  $\mu_2 = m - mb$ . Both measurements are performed with the same relative uncertainty  $\delta$ . An error weighted average of the two measurements returns

$$\mu_{ave} = m \frac{1 - b^2}{1 + b^2}, \quad (26)$$

which for  $b = 5\%$  corresponds to 0.5% bias.

The bias occurs because the measurement at smaller value  $\mu_2$  got smaller absolute uncertainty  $\delta(m - mb)$ .

Measurements with multiplicative uncertainties can be combined bias-free using expected values  $m_i$  instead of measured  $\mu_i$  to translate relative to absolute uncertainties. In this case Eq. 9 takes form:

$$\chi^2(\vec{m}, \vec{b})_{exp,mult} = \sum_i \left( \frac{m_i [1 - \sum_j \gamma_i^j b_j] - \mu_i}{\delta_i m_i} \right)^2 + \sum_j b_j^2, \quad (27)$$

## 2.5 Bias correction for statistical uncertainties

Let's consider the counting of number of arbitrary events. Two measurements  $\mu_1$  and  $\mu_2$  gives  $\mu_1 = N_1$ ,  $\mu_2 = N_2$ . Statistical uncertainties of the measurement are estimated as a square root of number of counts. Weighted average for these measurement returns:

$$\mu_{ave} = \frac{2N_1N_2}{N_1 + N_2}, \quad (28)$$

instead of

$$\mu_{ave} = \frac{N_1 + N_2}{2} \quad (29)$$

Bias for statistical average can removed by using expected instead of measured number of events. If statistical uncertainty for a measurement is quoted based on square root of number of measured events, then estimated unbiased relative statistical uncertainty  $\delta_{stat,cor} = \frac{\sqrt{m}}{\mu} = \delta_{stat} \sqrt{\frac{m}{\mu}}$ . Absolute unbiased statistical uncertainty can be expressed as:

$$\Delta_{stat,cor} = \delta_{stat} \sqrt{m\mu} \quad (30)$$

Finally, the number of observed events can be modified by the correlated systematic uncertainties. This modification can be taken into account by using

$$m(1 - \sum_j \gamma^j b_j) \quad (31)$$

instead of  $m$  in Eq. 30. This brings us to the  $\chi^2$  formula:

$$\chi^2(\vec{m}, \vec{b})_{exp,cor} = \sum_i \frac{(m_i[1 - \sum_j \gamma_i^j b_j] - \mu_i)^2}{\delta_{i,stat}^2 \mu_i m_i [1 - \sum_j \gamma_i^j b_j] + \delta_{i,uncorr}^2 m_i^2} + \sum_j b_j^2, \quad (32)$$

## 2.6 Treatment of off-set systematic uncertainties

Let's consider a source of systematic uncertainty on a certain measurement (so called off-set systematics), which does not have any contribution to the  $\chi^2$  and therefore does not have any impact on the average value. Systematic uncertainty on the combined value due to this source ( $\delta_{ave}^{off-set}$ ) can be calculated as:

$$\delta_{ave}^{off-set} = \frac{\mu_{ave}^{up} - \mu_{ave}^{down}}{\mu_{ave}}, \quad (33)$$

where  $\mu_{ave}^{up/down}$  the average value of the measurements  $e$ , shifted by considered off-set systematics  $\mu_e \pm \delta_{ave}^{off-set}$ . Therefore in case of  $N_o$  sources of the off-set systematics combination procedure performs  $2N_o + 1$  times: one nominal combination,  $N_o$  "up" and  $N_o$  "down" combinations for each sources of the off-set systematic respectively.

## 2.7 Treatment of asymmetric systematic uncertainties

In case if assumption of the symmetric systematic uncertainty (expressed by Eq.7) is not valid for performed measurement, the  $\chi^2$  Eq. 9 can be written in more general form:

$$\chi^2(\vec{m}, \vec{b})_{exp} = \sum_i \frac{(m_i - \mu_i - \sum_j f_i(b_j))^2}{\Delta_i^2} + \sum_j b_j^2. \quad (34)$$

If  $f_i(b_j) = \Gamma_i^j b_j$ , Eq. 34 again back to Eq. 9. Asymmetric systematic uncertainties can be approximated as:

$$f_i(b_j) = \Gamma_i^j b_j + \omega_i^j b_j^2, \quad \omega_i = \frac{\Gamma_i^{j+} + \Gamma_i^{j-}}{2}. \quad (35)$$

$\chi^2$  definition in this case become non-linear. Instead of simple minimization procedure described in Sec. 2.2, iterative minimization procedure, shown in Sec. 2.9 is used.

## 2.8 Different representation of combined uncertainties

The systematic uncertainties of the averaged values are presented in orthogonal way (in terms of diagonalised sources of systematic uncertainties) by  $\Gamma_{ave}$  (see Eq.17). However in this representation each sources of the systematic uncertainty  $\Gamma_{ave}^j$  is a linear combination of the initial sources of the systematic uncertainties  $\Gamma^j$ .

In order to be able to compare systematic uncertainties of the combined measurement with  $\Gamma^j$  source by source, diagonal elements of matrix  $A'_S$  can be taken as it is (without diagonalisation). However in this case systematic uncertainties of the combined measurement will be not orthogonal and their quadratic sum will not give a total systematic uncertainty.

Alternative way is to use half-diagonal form of matrix  $A'_S$ . In this case systematic uncertainties of the combined measurement are orthogonal, however one of the course corresponds to one of the sources on the initial systematic uncertainties  $\Gamma^j$ .

## 2.9 Combination with bias corrections

For most practical situations the bias, described in Sec. 2.4 and 2.5 of the average is small. Therefore, the expectation  $m_i$  can be estimated in an iterative procedure starting from linear formula Eq. 9 and using  $m_i = \mu_{i,ave}$ . The key for the unbiased result is that the same expectation is used for all measurements. In most cases the convergence is observed after second iteration.

Similar iterative approach is applied to combine the measurements with asymmetric systematics uncertainties, introduces in Sec. 2.7. The first iteration is performed with linearised  $\chi^2$  Eq.9 using symmetrised uncertainties  $\Gamma$  (linear part of  $f(b_j)$ ). The next iterations are performed with corrected uncertainties  $\Gamma' = \Gamma + \omega * \beta_{ave}$ , e.g. correction depends on the systematic shift. The combination procedure require several iterations.

Described bias corrections and correction of the systematic uncertainties are not interfere with each other and therefore both are applied in simultaneously. E.g. at each iterations both bias corrections and correction for non-symmetric uncertainties are applied.

Iterative averaging with bias correction and only symmetrical uncertainties converging fast and requires 2–3 iterations. Presence of asymmetric uncertainties make convergence worse. In some cases iterative procedure will not converge at all. In order to monitor convergence of the iterative procedure systematic shifts  $\beta_{ave}$  for each iteration have to be considered.

## 3 Program manual

This is Fortran-base open-source software, which can be downloaded from: <https://wiki-zeuthen.desy.de/HERAverager>.  
Installation and basic pre-requirements are described in Sec. 3.1. Code organization is shown in Sec. 3.2.

The input and output formats are explained in Sec. 3.3 and Sec. 3.4 respectively.

### 3.1 Program installation

The program requires gfortran compiler to be compiled.

The program can be unpacked and installed with following lines:

```
tar -xvzf heraverager-1.0.0.tar.gz
cd heraverager-1.0.0
autoreconf --install
./configure
make
make install
```

After successful installation an executable file can be found in: `./bin/HERAverager`. All input parameters are defined in the steering file. Test run of the program, using test steering file (is in `./test/steering`) can be performed with lines:

```
cd ./bin
./HERAverager ../test/steering
```

The output information is located in `./output/`

### 3.2 Code organization

Package consist of following subdirectories:

- `include`: include-file with declaration of global variables
- `source`: source-code files
- `doc`: documentation to the program, including this manual.
- `test`: example of steering and data files
- `num_utils`: Utilities for the averaging.

Short description of source files is given below:

- `averaging.f`: Calling functions to prepare arrays and perform the combination
- `average_py.f90`: Interface to python (see Sec. 5)

- 226 • `error_logging.f`: print full error summary and close files
- 227 • `initave.f`: Reading the data and preparing for the combination
- 228 • `statrecalc.f`: Recalculate systematic uncertainties after combination and prepare values for the
- 229 next iterator.
- 230 • `common_tools.f`: Function, used in different steps of the combination
- 231 • `fillarrays.f`: Prepare arrays used for combination
- 232 • `output.f`: Print output of the combination in output files
- 233 • `covartouui.f`: Convert covariance matrix to nuisance param. representation
- 234 • `heraverager.f`: Main file, which is calling functions for initialization, combination and output.
- 235 • `readdata.f`: Read a table of measurement with their uncertainties from single experiment
- 236 • `toblockdiag.f`: Perform the combination by minimization of linearised  $\chi^2$

237 Summary of the code structure is shown in Fig. 1. 3 logic blocks: initialization, combination and  
 238 output are encapsulated in 3 functions called in `heraverager.f`.

239 At initialization step all files with data and options are read and values are stored in global internal  
 240 variables. First code read the steering file, which contain all options of the combination and list of data  
 241 files. During the loop over data files the name of systematic sources, values of the measurements and  
 242 their uncertainties are stored in global variables.

243 Combination block have a loop over all off-set systematics (2N+1 times, see Sec. 2.6) and over  
 244 iterations. Combination process starts with filling of auxiliary arrays (mainly elements of Eq. 13). Then  
 245 minimization is performed and results are stored in global internal variables. In case if it was not the last  
 246 iteration statistical and systematic uncertainties are recalculated (See Sec. 2.9).

247 The output written in different files. Systematic uncertainties can be presented in different ways (See  
 248 Sec. 2.8), depending on options in the steering file.

### 249 3.3 Input information

250 All input of the combination as well as the list of data files and supplementary information is given in  
 251 the steering file. The description of the steering file is given in Sec. 3.3.1, while data file is discussed in  
 252 Sec. 3.3.2.

#### 253 3.3.1 Steering file

254 Steering file is organized in blocks, which corresponds to name-lists in the Fortran code. Lines, which  
 255 starts with ! are commented. Structure and meaning of these blocks are described below:

```

256 &InFiles
257   ! Specify data files to be averaged
258   NInputFiles = 2
259   InputFileNames(1) = '../test/h1460new.public.dat'
260   InputFileNames(2) = '../test/zeus460.public.dat'
261 &End

```

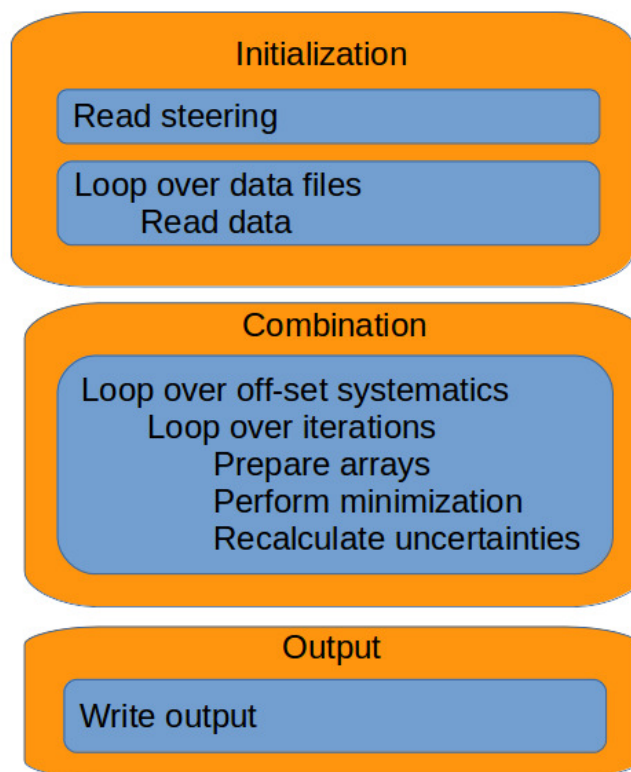


Figure 1: Code organization of the combination tool.

262 Block contains the list of data-files (see structure of the data files in Sec. 3.3.2). Path to these files have  
263 to be absolute and relative with respect to the place, where the program runs.

```
264 &CommonGrid
265   GridType = 'External' ! 'External' or 'Auto'
266   GridFiles = '../test/grid460575.dat'
267   AveSameExp = .true.
268 &End
```

269 Most of the measurements, for which tool is designed are binned measurements. The binning can be  
270 taken directly from the data file (option 'Auto') or can be written externally (option 'External') in  
271 grid of bins.

272 In case of external binning, files which describe binning have to be provided and have structure,  
273 similar to:

```
274 &Grid
275   Reaction = 'NC e+-p'
276   NDimension = 2
277   NPoints = 630
278   BinNames = 'Q2', 'x'
279 &End
280   1.5   0.000378678
281   1.5   0.000227206998
282   1.5   0.000134258007
283   1.5   9.52802002E-05
```

284 Where:

- 285 • Reaction is the name of measured process
- 286 • NDimension is the number of dimensions of the binning
- 287 • NPoints total number of bins
- 288 • BinNames name of the dimensions.

289 Numbers below the header define the center of bins. Each line represent one bin (number of lines  
290 should be equal to the total number of bins). Number of columns corresponds to the number of dimen-  
291 sions (two in this example) and define center of bins in each dimension ('Q2' and 'X' in example).

292 Each grid-file describe the binning for a certain measured process (reaction). In case of combination  
293 of the data for several different processes, several grid files have to be given.

294 The name of dimensions and the name of measured process should coincide with one given in data  
295 file (see Sec. 3.3.2). The measurements are considered as a measurements of the same physics quantity  
296 if bins, bin names and reactions are coincides.

297 In case if the bin-centers in the data file are different compared to bin-centers in the grid file nearest  
298 bin of the grid file will be used to define the bin. Parameter AveSameExp clarifies, how to treat the case  
299 when two bins from one experiment go to one grid bin:

- `.true.` - weighted average of these bins will be used for combination.
- `.false.` - measurements in these bins will be combined as different measurements of the same physics quantity.

```

303 &HERAverager
304   OutputMode   = 'ORTH'
305   OutputPrefix = 'Ave'
306   OutputFolder = '../output'
307   IDebug       = 0
308   WriteOriginal = .false.
309   WriteSysTexTable = .false.
310   PostRotateSyst = .true.
311 &End

```

312     `OutputMode` - is the output options for the systematics uncertainties (see Sec. 2.8).

- `'ORTH'` - orthogonal representation
- `'ORIG'` - original structure of the systematic uncertainties.

315     `OutputPrefix` Add a prefix for output file names

316     `OutputFolder` - folder, where the output information is stored

317     `PostRotateSyst = .true.` keep output systematic uncertainties align to the original sources as  
318 much as possible.

319     `IDebug` - Debug level. Higher value corresponds to more debug messages.

320     `WriteOriginal` - include original information to the output summary (file `ave Ave.dat`, see Sec. 3.4)

321     `WriteSysTexTable` - write output information about systematic uncertainties in tex format (file  
322 `sys.tex`, see Sec. 3.4)

```

323 &BiasCorrection
324   AverageType = 'MIXED'
325   Iteration   = 10
326   ! Rescale the stat and uncorr uncertainties separately:
327   RescaleStatSep = .false.
328   ! Correction of the syst bias for stat errors
329   CorrectStatBias = .false.
330   ! Keeping the stat errors fixed'
331   FixStat        = .false.
332 &End

```

333     Parameter `AverageType` define the type of the systematic uncertainties

- `'ADD'` - all systematic uncertainties are processed as additive
- `'MULT'` - all systematic uncertainties are processed as multiplicative



- 'MIXED' - type of systematic uncertainties is taken from the data file

Parameter `Iteration` set the number of iterations. In case, if all uncertainties are additive and symmetric  $\chi^2$  is linear, no additional iteration is required. In case of multiplicative systematic uncertainties 2-3 iterations are recommended to get stable combined value. If some of the uncertainties are asymmetric number of iteration have to increased (7-10 iterations).

Correction of the statistical uncertainties, discussed in Sec. 2.5 is normally performed for both statistical and uncorrelated systematic uncertainties. This correction can be done separately for each of them by setting flag `RescaleStatSep = .true..` Using flag `FixStat = .false.` statistical uncertainties are still uncorrected.

Setting flag `CorrectStatBias = .true.` statistical bias is corrected by implementing Eq. 31

### 3.3.2 Data file

The data-file is an ascii-file, which contains measurement values in a certain bins and their uncertainties.

The data file consist of the header and lines of data:

```
&Data
  Name = 'H1'
  NData = 72
  NColumn = 9
  ColumnType = 2*'Bin','Sigma', 6*'Error'
  ColumnName = 'Q2', 'x', 'x-sec', 'stat', 'uncor', 'sys1', 'sys2:0', 'sys3+', 'sys3-'
  Reaction = 'NC e+-p'

  Percent = true,true,true,true,true,true

&END
  1.500  3.47999E-05  0.51952  8.10  4.96  0.69  8.00  0.45  0.58
  2.000  4.64000E-05  0.70449  4.57  4.31  1.10 -0.81  0.49  0.59
```

Header of the data file contain following information:

- Name: A name of the data file, which is used for user-friendly output
- NData: Number of data points (bins)
- NColumn: Total number of columns in the table, which includes bins, value of the measurement and uncertainties
- ColumnType: Type of the column. In case if several columns have the same type they can be grouped as e.g. 6\*'Error'. All columns should have a type. Following types are supported:
  - Bin: Bin-center. Number of these columns gives the number of dimensions of the measurement.
  - Sigma: Value of the measurement
  - Error: Uncertainty

- **ColumnName:** The name of the column. All columns should have a name. In case of 'Bin' type of the column name specify the name of the dimension. If case of using external grid file these names have to be mentioned in the grid file. For the 'Error' type of the column the name specify the name and type of systematic uncertainty. Following type of systematic uncertainties are supported:

- **stat:** Statistical uncertainty
- **uncor:** Bin-to-Bin uncorrelated systematic uncertainty
- **ignore:** Column ignored
- **somename:** Bin-to-Bin correlated multiplicative systematic uncertainty with the name “somename”. One data-file should not contain two columns with uncertainties with the same names. If systematic uncertainty with the same name is found in different data-files, they are considered as correlated between measurements.
- **somename:0:** Off-set systematic uncertainty
- **somename:A:** Bin-to-Bin correlated additive systematic uncertainty
- **somename+, somename-:** Bin-to-Bin correlated asymmetric systematic uncertainty (can be :0 or :A). Each asymmetric systematic uncertainty should contain both + and - part.

- **Reaction:** is the name of measured process. In case of using external grid file this name have to be mentioned in the grid file.

- **Percent:** Array of flags, shows is corresponding uncertainty is given in % or absolute. Number of entries in this field should corresponds to the number of uncertainties.

The data are given below the header as a table. The structure of the table corresponds to one described in the header. In a given example the table should contain 72 rows and 9 columns.

### 3.4 Output information

Output information are stored in ascii files inside of output folder. If this folder does not exist, the folder will be created automatically. If the files with the same names are exists, they will be overwritten.

output information consist of following files:

- **matrix.dat:** Covariance matrix  $\Gamma_{ave}$
- **eigvalues.dat:** Eigenvalues of matrix  $\Gamma_{ave}$
- **eigvectors.dat:** Eigenvectors of matrix  $\Gamma_{ave}$
- **tab.dat** – Summary table of the combination: combined value, statistical uncertainty, uncorrelated systematic uncertainty, all sources of the correlated systematic uncertainties.
- **ave\_Ave.dat:** – Short summary of the combination: Bin centers, combined value, statistical uncertainty, total systematic uncertainty, total uncertainty
- **ave\_proc\_Ave101.dat** – Summary of the combination, in the format of data file.
- **sys.txt** – Shift and reduction of the systematic uncertainties: count-number of systematic uncertainty, name of systematic uncertainty, shift of systematic  $\beta_{i,ave}$  (see Eq. 15), reduction of the systematic uncertainty  $D_{jj}$  (see Eq. 17), pull of systematic uncertainty (see Eq. 20).

- 409 • `sys.tex` same as `sys.txt`, but in latex format
- 410 • `chi2map.dat` – Information about pulls (see Eq. 19). The information is given in a table with  
411 columns: Count-number of the data point, number of degrees of freedom, bin1, .., binN, pull,  
412 data-file.
- 413 • `offsettab.dat` – Summary of off-set systematic uncertainties: combined value for all iterations  
414 over offset systematics (see Sec. 2.6)
- 415 • `ItrInfo.dat` – Information about systematics shifts and  $\chi^2$  compatibility for all performed it-  
416 erations. First line shows  $\chi^2$ , other lines show systematic shifts for different systematic sources.  
417 Different columns show values for different iterations.

## 4 Examples

An example of the input information can be found in folder `test` and consist of following files:

- `steering` - steering file
- `grid.dat` - grid file
- `elz.dat` - data file
- `elzFwd.dat` - data file
- `muz.dat` - data file

In this example data from 3 data-files are combined using grid from file. 2-dimensional grid (name of dimensions: 'Bla' and 'Y') contains 8 points in total for reaction: 'NC'. The dimension 'Bla' is a dummy dimension (only one bin with bin-center 0.0). Only one reaction is considered in this example. Bin-centers in some of the data files are not coincide with one described in the grid file. Following to the option in the steering weighted average of these bins will be used for combination.

All described uncertainties are multiplicative. Two iterations are used to correct for bias in statistical uncertainties.

The output of the combination is stored in folder `output` using orthogonal representation of the systematic uncertainties. The output prefix `HZComb` is used for files: `HZComb.dat` and `HZComb101.dat`

## 5 Python scripts

Basic functionality of the Fortran-based software described in Sec. 3 was compiled as a Python library (module). The installation is performed similar to 3.1, with additional key for configure script:

```
./configure --enable-python
```

After successful installation created library (which can be used as Python module) is available in `./bin/averager.so`

The code, which implements Fortran to Python interface is in file: `./source/average_py.f90` (written in Fortran90).

Several scripts written in python are provided to support Python-based averaging.

- `DatasetGen.py`: Script to generate random dataset, which can be used for testing `averager`
- `DataReader.py`: Module to parse data-files and read data in Python variables
- `test.py`: Script which shows an example of using Python-based `averager`

The `test.py` file demonstrates the work of data reader and Python-based `averager`. Also it contains a plotting part.

### 5.0.1 Data reader

Data reader module `DataReader` reads the data from all `.csv` files in current directory. The module consists of one function `paverager(bins,data,error)`, where:

- `bins`: names of bin columns ('bin1,bin2,...'). If the input string is empty all columns with substring 'bin' are considered as bins.
- `data`: name of measurement column (input as 'data1'). If the input string is empty column 'data' is considered as measurement.
- `error`: names of uncertainty columns. Columns with substring 'stat' are considered as bin-to-bin uncorrelated. At least one column should be bin-to-bin uncorrelated for statistical uncertainty ('error1:M,error2:stat,...'). Different uncertainties should have different names. If the input string is empty all columns with substring 'error' are considered as bin-to-bin correlated uncertainty, all columns with substring 'stat' are considered as bin-to-bin uncorrelated uncertainty.

The function returns 3 lists: measurements, their systematic and statistical uncertainties. All considered bin-to-bin uncorrelated uncertainties are summed quadratically and considered as statistical uncertainty.

Example of usage:

```
mes,stat,syst=DataReader.paverager('',' ','')
```

The format of the input data-file is following:

```
bin1,bin2,data,stat,error00000,error00001,error00002,error00004
0,1,14.01,0.261,1.499,1.43,-1.293,1.334
5,1,37.7,0.552,3.829,3.461,-2.908,4.197
6,1,1.84,0.048,0.167,0.167,-0.15,0.157
7,1,69.75,1.918,7.331,7.194,-5.481,6.44
8,1,6.44,0.169,0.502,0.659,-0.665,0.514
9,1,58.61,1.456,6.055,7.045,-5.182,6.405
10,1,3.03,0.041,0.293,0.275,-0.217,0.281
12,1,6.46,0.089,0.63,0.649,-0.534,0.611
13,1,23.79,0.345,2.44,2.279,-2.026,2.233
14,1,11.67,0.303,1.293,1.293,-0.839,0.582
17,1,9.64,0.253,1.12,0.75,-1.042,0.974
```

Here header contains the name of the columns (separated by “,”).

Lists of systematic-names `snames` and file-names `fnames` can be extracted using corresponding global variables: `oerror` and `fnames`.

```
snames=DataReader.oerror
fnames=DataReader.fnames
```

## 5.0.2 Generation of random dataset

Parameters of the data-set generator are given in the beginning of the script in following variables:

- `nData` (integer): Number of data points
- `nMes` (list of 3 integers): Total number of measurements, minimal and maximal number of measurements per data point (have to be smaller as total number of measurements). In case if minimal, maximal and total are the same, each data point will have the same number of measurements.
- `nSyst` (list of 3 integers): Total number of systematic uncertainties, minimal and maximal number of correlated systematic uncertainties per measurement (have to be smaller as total number of systematic uncertainties). In case if minimal, maximal and total are the same, each measurement will have the same systematic uncertainties.
- `vStat` (list of 2 floats): Minimal and maximal relative statistical uncertainty
- `vSyst` (list of 2 floats): Minimal and maximal relative systematic uncertainty

The output of the generator is stored in text files (Each measurement is in separate file):

- `test0.dat`, `test1.dat`, ... – suitable for the Fortran-based averager.
- `test0.csv`, `test1.csv`, ... – suitable for the Python-based averager.

## 5.1 Python-based averager

Python-based averager can be imported as module `averager`. The function, which performs the averaging have 3 input and 3 output parameters:

```
ave, statOut, systOut = averager.averager(mes, stat, syst)
```

where:

- `mes` – measurements
- `stat, syst` – statistical and systematic uncertainties of the measurements.
- `ave` – average values
- `statOut, systOut` – statistical and systematic uncertainties of the average values.

Input parameters are prepared by `DataReader`.

Default parameters of the averager can be modified by changing global variables of the `averager`. In this case `averager` have to be initialized before these changes. Otherwise all parameters will be overwritten by their default values.

```
#initialization
averager.avin.initvariables()

averager.avin.indebug = 0
averager.avin.inwriteoriginal = .false.
averager.avin.inwritesystextable = .false.
averager.avin.inpostrotatesyst = .false.

averager.avin.setoutputprefix('Ave')
averager.avin.setoutputfolder('../output')
averager.avin.setsnames(snames)

averager.avin.initeration = 10
averager.avin.inrescalestatsep = .false.
averager.avin.inborroctbtatbias = .false.
averager.avin.infixstat = .false.
```

Parameters given here are default parameters. All variable have the same meaning as for the Fortran-base program (see Sec. 3.3.1). Output mode is always orthogonal.

Names of the systematic uncertainties are given by the list of names `snames`. If names are not given all systematics are considered as symmetric, multiplicative.

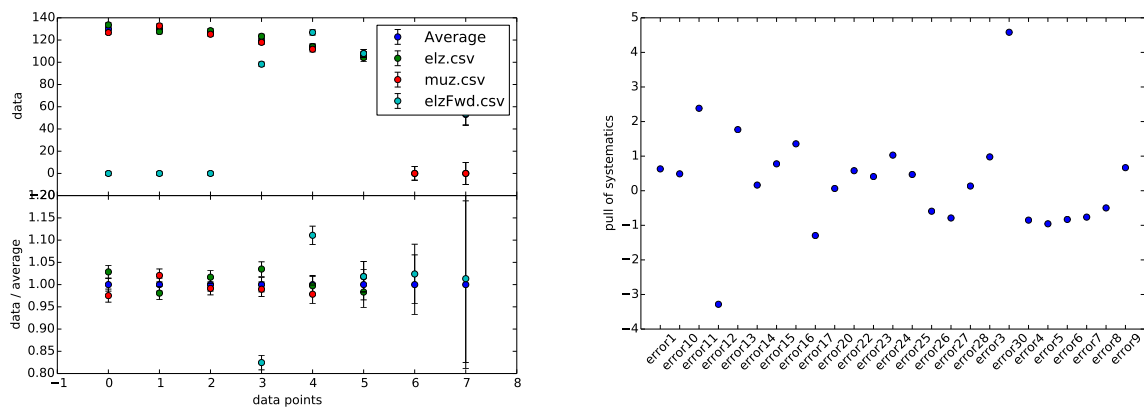
Additional output information of the averaging is available with following variables

- `averager.avout.pulldata` – pulls of data

- `averager.avout.pullsys` – pulls of systematic uncertainties
- `averager.avout.shiftsys` – shift of systematic uncertainties
- `averager.avout.squeezesyst` – reduction of the systematic uncertainties
- `averager.avout.chi2` –  $\chi^2$
- `averager.avout.ndof` – number degrees of freedom

## 5.2 Python-based plotting

Some examples of plotting functionality is shown in the `test.py`. Script plots the averaged data and pull of systematic uncertainties. Results are stored in pdf files. An example of output plot is shown below.



Additionally a script `plot.py` plot histograms of data pulls and systematic pulls using text files as input. Therefore this script can be used also to visualise output of the fortran-base averager. Using of the script is following:

`./plot.py output`

where `output` is the path to the output folder of the averager. This script reads files: `sys.txt`, `tab.dat`, `chi2map.dat` and save plots as pdf files.



## References

- [1] A. Glazov, “Averaging of DIS cross section data,” AIP Conf. Proc. 792 (2005) 237–240.
- [2] H1 and ZEUS collaborations, “Combined measurement and QCD analysis of the inclusive  $e^\pm p$  scattering cross sections at HERA”, JHEP01 (2010) 109
- [3] H1 Collaboration, “Measurement of the Inclusive ep Scattering Cross Section at Low  $Q^2$  and  $x$  at HERA”, Eur. Phys. J. C **63** (2009) 625
- [4] S. Glazov, “Data Combination”, Statistics School, DESY, October 2011. Slides.