

# Correlations between parameters of extended air showers and their proper use in analyses

Wolfgang Wittek, Harald Kornmayer  
*Max-Planck-Institut für Physik, München*  
 for the HEGRA Collaboration

## Abstract

In air shower experiments information about the initial cosmic ray particle or about the shower development is obtained by exploiting the correlations between the quantities of interest and the directly measurable quantities. It is shown how these correlations are properly treated in order to obtain unbiased results. As an example, the measurement of the average penetration depth as a function of the shower energy is presented.

## 1 Introduction

In air shower experiments cosmic ray particles are studied in an indirect way: the distributions of the interesting quantities ( $\vec{X}_{\text{orig}}$ ) of the initial cosmic ray particle (like its nature and energy) or of the air shower (like the penetration depth) have to be inferred from the distributions of measurable quantities ( $\vec{X}_{\text{meas}}$ ) (like particle and  $\check{C}$ -light densities at detector level). Using Monte Carlo (MC) simulations, in which the interaction of the cosmic ray particle with the atmosphere, the shower development and also the properties of the detector are simulated, one is able to establish the correlations between the measurable quantities  $\vec{X}_{\text{meas}}$  and the interesting quantities  $\vec{X}_{\text{orig}}$ . These correlations are then used to determine the distributions of  $\vec{X}_{\text{orig}}$  from the experimental distributions of  $\vec{X}_{\text{meas}}$ .

The aim of this paper is to demonstrate that it is essential to treat the correlations in a mathematically correct way in order to avoid biases in the results. As an example the determination of the average penetration depth  $X_{\text{max}}$  of air showers in the atmosphere as a function of the shower energy  $E$  is presented, using data from the HEGRA array of scintillator and wide-angle  $\check{C}$ -light detectors.

In this example, the number of shower particles ( $N_s$ ) and the  $\check{C}$ -light radius at detector level ( $R_L$ ) are taken as measurable quantities.  $N_s$  is determined from the particle densities as measured by the matrix of scintillation detectors and  $R_L$  is obtained as the inverse of the slope of the lateral  $\check{C}$ -light distribution as measured by the matrix of  $\check{C}$ -light detectors. The quantities of interest ("true" quantities) are the shower energy  $E$  and the penetration depth  $X_{\text{max}}$ .

In first approximation,  $N_s$  carries mainly information on  $E$ , and  $R_L$  mainly on  $X_{\text{max}}$ . While the correlation between  $R_L$  and  $X_{\text{max}}$  is quite independent of the nature (or atomic number  $A$ ) of the cosmic ray particle (Lindner, 1998), this is not the case for the correlation between  $N_s$  and  $E$ . It has been shown in (Cortina, 1997) that by using a modified  $N_s$  ( $N_s^c = N_s \cdot R_L^\alpha$ , with  $\alpha$  depending on the zenith angle) the correlation between  $N_s^c$  and  $E$  is quite independent of  $A$ .

## 2 The Method

The task is now to determine the 2-dimensional distribution of the variables ( $\log E, X_{\text{max}}$ ) from the 2-dimensional distribution of the variables ( $\log N_s^c, 1/R_L$ ). The simplest way of doing this is to use the average correlations

$$\langle \log N_s^c \rangle = f_1(\log E); \quad \langle 1/R_L \rangle = f_2(X_{\text{max}})$$

as determined from a sample of Monte Carlo events. In this procedure a one-to-one correlation is assumed between  $\log E$  and  $\log N_s^c$  and between  $X_{\text{max}}$  and  $1/R_L$  respectively. In addition, a possible  $X_{\text{max}}$  dependence of the  $\log E - \log N_s^c$  correlation, and a  $\log E$  dependence of the  $X_{\text{max}} - 1/R_L$

correlation is ignored. This procedure will later be referred to as the "one-to-one correlation procedure".

A mathematically correct approach is to use the full correlation between  $(\log E, X_{\max})$  and  $(\log N_s^c, 1/R_L)$ . It is convenient to define a grid in the  $(\log E, X_{\max})$  plane and a grid in the  $(\log N_s^c, 1/R_L)$  plane, and some numberings ( $i = 1$  to  $N_{xy}$ ) and ( $j = 1$  to  $N_{uv}$ ) of the resulting bins in the two planes respectively. The bin content of bin  $(j, i)$  of the distribution  $(\log N_s^c, 1/R_L)$  versus  $(\log E, X_{\max})$  for the sample of Monte Carlo events may be denoted by  $G_{ji}$ . The full correlation can then be written as

$$g_{ji} = \frac{G_{ji}}{(\sum_k G_{ki})} \quad (1)$$

$g_{ji}$  describes how a particular pair of values  $(\log E, X_{\max})$ , defined by a specific value  $i_o$  of  $i$ , is transformed into a distribution of  $(\log N_s^c, 1/R_L)$ , given by  $g_{ji_o}$  ( $j = 1$  to  $N_{uv}$ ). The division by  $(\sum_k G_{ki})$  was done to make  $g_{ji}$  independent of the  $(\log E, X_{\max})$  distribution in the Monte Carlo sample.

If the experimental distribution of  $(\log N_s^c, 1/R_L)$  is denoted by  $a_j$  ( $j = 1$  to  $N_{uv}$ ) and the distribution to be determined in  $(\log E, X_{\max})$  by  $b_i$  ( $i = 1$  to  $N_{xy}$ ), then  $b_i$  has to fulfill the condition

$$a_j = \sum_i (g_{ji} \cdot b_i) \quad (2)$$

The condition (2) ensures that the full correlations between the measured and the true quantities are taken into account.

Determining the distribution  $b_i$  from a measured distribution  $a_j$ , with known response matrix  $g_{ji}$ , is a typical unfolding problem. The main point of the unfolding methods is to impose, in addition to (2), certain smoothness conditions on the distribution  $b_i$  in order to avoid strong fluctuations of  $b_i$ , which arise from statistical fluctuations of  $a_j$ . In the example discussed here the method of reduced crossed entropy (MRX) is applied (Schmelling, 1994). In this method, a kind of smoothness condition is imposed by requiring the solution  $b_i$  to be close to a prior distribution  $b_i^{\text{prior}}$ .  $b_i^{\text{prior}}$  may be some guess of the true distribution. Usually the result  $b_i$  is quite independent of the choice of  $b_i^{\text{prior}}$ , so that  $b_i^{\text{prior}}$  may be set to a constant.

### 3 Results

The response matrix  $g_{ji}$  for the example discussed here is shown in Fig. 1. The ordinate corresponds to bins in the  $(\log E, X_{\max})$  plane, the abscissa to bins in the  $(\log N_s^c, 1/R_L)$  plane.  $g_{ji}$  was obtained by averaging the response matrices for different chemical elements ( $A = 1, 4, 16$  and  $56$ ) and by smoothing the average response matrix in the following way:  $g_{ji}$  was parametrized as a 2-dimensional Gaussian distribution in the variables  $(u, v) = (\log N_s^c, 1/R_L)$

$$g_{ji} = \frac{1}{2\pi\sigma_u\sigma_v\sqrt{1-\rho^2}} \cdot \exp\left\{-\frac{1}{2(1-\rho^2)}\left[\left(\frac{u-\bar{u}}{\sigma_u}\right)^2 - 2\rho\frac{(u-\bar{u})(v-\bar{v})}{\sigma_u\sigma_v} + \left(\frac{v-\bar{v}}{\sigma_v}\right)^2\right]\right\} \quad (3)$$

where the 5 parameters  $p = \bar{u}, \bar{v}, \sigma_u, \sigma_v$  and  $\rho$  were assumed to be linear functions of  $\log E$  and  $X_{\max}$  (3 parameters for each of the 5 parameters  $p$ ). The  $5 \times 3 = 15$  free parameters were determined by fitting expression (3) to the average response matrix. The fitted values of the parameters characterize in detail the behaviour of the correlations between  $(\log E, X_{\max})$  and  $(\log N_s^c, 1/R_L)$  and their dependence on  $\log E$  and  $X_{\max}$ . In particular one finds: In very good approximation,  $\langle 1/R_L \rangle (= \bar{v})$  is only a function of  $X_{\max}$ .  $\langle \log N_s^c \rangle (= \bar{u})$  is mainly a function of  $\log E$ , with some additional dependence on  $X_{\max}$ . The parameter  $\rho$ , which describes the correlation between  $\log N_s^c$  and  $1/R_L$  at fixed  $(\log E, X_{\max})$ , is a function of  $\log E$ : at low  $\log E$   $\log N_s^c$  and  $1/R_L$  are anti-correlated, whereas they are positively correlated at higher  $\log E$ . All these properties of the correlations are, of course, taken into account in the unfolding procedure.

The experimental distribution of  $(\log N_s^c, 1/R_L)$  is shown in Fig. 2 (Kornmayer, 1999). It should be noted that the measurements presented in this figure are based on preliminary data and are not the final official HEGRA results. By applying the MRX method a distribution of  $(\log E, X_{\max})$  is obtained ("unfolded" distribution) which is displayed in Fig. 3. Forming the average  $X_{\max}$  for each bin of  $\log E$  yields the result for the elongation plot  $\langle X_{\max} \rangle$  versus  $\log E$ , shown in Fig. 4a. In Fig. 4b the RMS of  $X_{\max}$  is plotted as a function of  $\log E$ .

For comparison, in Fig. 4 the results from the "one-to-one correlation procedure" (see above) are also plotted. It can be seen that the latter procedure underestimates  $\langle X_{\max} \rangle$  by  $\sim 30 \text{ g/cm}^2$  at low  $\log E$  and by  $\sim 10 \text{ g/cm}^2$  at high  $\log E$ . The points in the bin of highest energy should be taken with care because they are based on low statistics both in the experimental data and in the MC sample. No systematic differences between the two methods are seen for the RMS of  $X_{\max}$  (Fig. 4b). Knowing that the one-to-one correlation procedure yields biased results one can try to correct the results by applying additional correction factors to  $\langle X_{\max} \rangle$ , which are determined from MC events. However, these correction factors will in general depend on the details of the MC simulation, in particular on the distribution of  $(\log E, X_{\max})$ .

By construction, the result of the unfolding procedure does not depend on the underlying MC distribution of  $(\log E, X_{\max})$ . By fulfilling the condition (2) (at least approximately, see Fig. 2), it takes into account the full correlations between the measured and the true quantities. Of course, these correlations and thus also the result for the elongation plot will depend on the model used in the MC simulation. How they depend on the MC model can be studied by doing the unfolding for different response matrices, corresponding to different MC models.

A study of the dependence of the results from the one-to-one correlation procedure on the MC model will be less conclusive because effects due to differences between the MC models and effects due to using a mathematical incorrect procedure are not well separated.

Since the average penetration depth of air showers depends on the nuclear mass number  $A$  of the cosmic ray particle inducing the air shower, a measurement of  $\langle X_{\max} \rangle$  as a function of  $E$  can be used to obtain information on the chemical composition of cosmic rays (see for example Roehring, 1999).

It should however be noted that from the same experimental data information about the chemical composition can also be obtained in a more direct way: one possibility is to start from the experimental 2-dimensional distribution of  $(\log N_s, 1/R_L)$  and apply the unfolding procedure to obtain the 2-dimensional distribution of  $(\log E, \log A)$ . In this case the response matrix would explicitly depend on  $\log E$  and  $\log A$  and one would not have to rely on an  $A$ -independence of the response matrix, as was the case for the example discussed in this paper. The  $A$ -independence was necessary in order to determine the distribution of the penetration depth. If one is only interested in the chemical composition a knowledge of this distribution is not required and the 2-dimensional distribution of  $(\log E, \log A)$ , which contains all the information about the chemical composition as a function of  $E$ , can be obtained directly.

## Acknowledgements

Fruitful discussions with M. Schmelling are gratefully acknowledged.

## References

- Lindner, A., 1998, *Astrop. Phys.* 8, 235
- Cortina, J., et al., 1997, *J. Phys. G: Nucl. Part. Phys.* 23, 1733
- Schmelling, M., 1994, *Nucl. Instr. Meth. A* 340, 400
- Kornmayer, H., 1999, PhD Thesis, Technische Universität München
- Roehring, A., et al., 1999, OG.1.2.09, Proc. 26th ICRC (Salt Lake City, 1999)

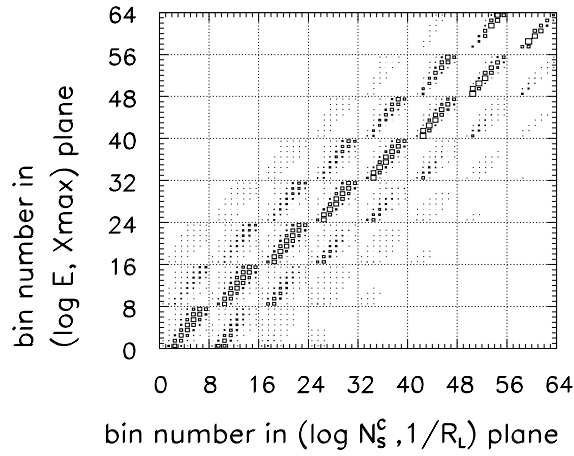


Figure 1: Response matrix  $g_{ji}$ . The abscissa corresponds to bins in the  $(\log N_s^c, 1/R_L)$  plane: In 8 consecutive bins  $1/R_L$  increases from 0.0 to  $0.035 \text{ m}^{-1}$ . Each block of 8 consecutive bins is for one bin in  $\log N_s^c$ . In 8 consecutive blocks  $\log N_s^c$  increases from 7.27 to 10.07. The ordinate corresponds to bins in the  $(\log E, X_{\max})$  plane: In 8 consecutive bins  $X_{\max}$  increases from 320 to  $800 \text{ g/cm}^2$ . Each block of 8 consecutive bins is for one bin in  $\log E$ . In 8 consecutive blocks  $\log E$  increases from 1.5 to 4.3.

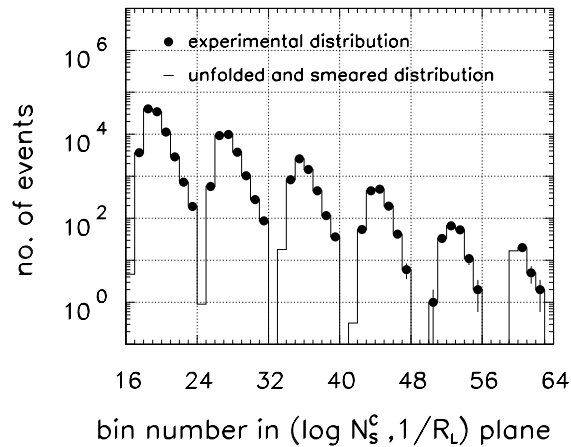


Figure 2: Experimental distribution of  $1/R_L$  in different bins of  $\log N_s^c$  (full circles). The histogram represents the result from applying the response matrix to the unfolded distribution of  $(\log E, X_{\max})$ . The good agreement between the two distributions shows that relation (2) between the measured and unfolded distribution is well fulfilled. The abscissa in this figure corresponds to the abscissa in Fig. 1.

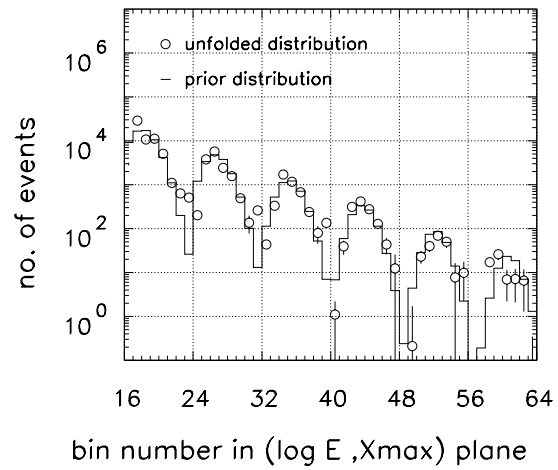


Figure 3: Unfolded distribution of  $X_{\max}$  in different bins of  $\log E$  (open circles). The histogram represents the prior distribution used in the unfolding procedure. The abscissa in this figure corresponds to the ordinate in Fig. 1.

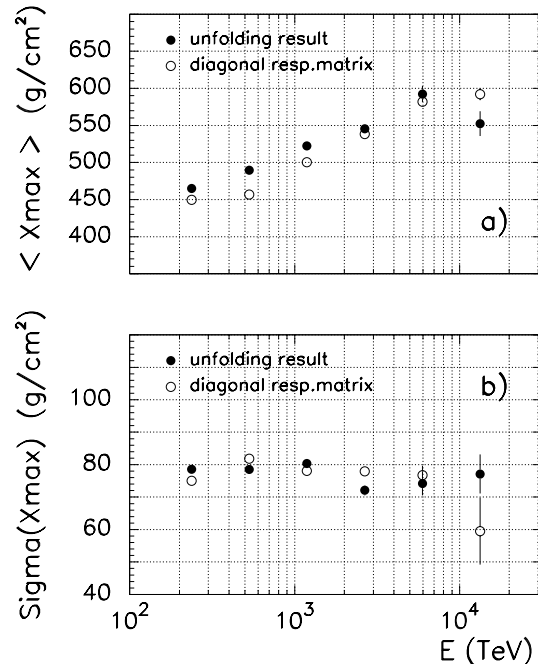


Figure 4: a) Average penetration depth  $\langle X_{\max} \rangle$  and b) RMS of  $X_{\max}$  as functions of  $E$ . The full circles represent the results of the unfolding procedure using the full response matrix. If one-to-one correlations are assumed between  $1/R_L$  and  $X_{\max}$ , and  $\log N_s^c$  and  $\log E$  respectively, the points represented by open circles are obtained ("diagonal response matrix"). While both methods yield consistent results for RMS of  $X_{\max}$ , the one-to-one correlation procedure in general underestimates  $\langle X_{\max} \rangle$  by 10 to  $30 \text{ g/cm}^2$ .