

ЕМ-819(46)-85

ԵՐԵՎԱՆԻ ՖԻԶԻԿԱՅԻ ԻՆՍԻՏՈՒՏ

N.Z. AKOPOV, Sv.Kh.ARUTUNIAN, A.A.CHILINGARIAN,
S.Kh.GALFAYAN, V.Kh.MATEVOSIAN, M.Z.ZAZIAN

**THE DESIGN PRINCIPLE AND STRUCTURE
OF "ANI" DATA CENTRE**

ЦНИИатоминформ

ЕРЕВАН-1985

С. - структурный научно-исследовательский центр в г. Барнауле, подчиненный Ученому совету по экономико-экономическим исследованиям при Академии наук Сибири и Алтайского края. Центр занимается проблемами экономики и социальной политики Алтайского края.

Georgian, 17th century, 17th century

REFERENCES

Б.В.АКОПОВ, С.Е.ХАКУМЯН,¹ С.Г.ДАЛАЧЯН,²
Ж.З.ЗАЗЯК,³ В.Л.МАТЕВОСЯН,⁴ А.Л.ЛЕНТЬЕВ⁵

ОБЩИЕ ПРИНЦИПЫ И СТРУКТУРА ЦЕНТРА ДАННЫХ "АНИ"

Обосновываются основные принципы, лежащие в основе выбора методов и алгоритмов пакета прикладных программ статистического анализа "АНИ". Непараметрические процедуры позволяют создать эффективную методику совместного анализа результатов вычислительных экспериментов и данных экспериментов в космических лучах. Приводится описание реляционной базы данных, обеспечивающей унифицированное хранение информации, защиту, обновление и удаление данных, а также быстрый и удобный доступ к данным.

Ереванский физический институт

Ереван 1986

The computing and experiment is planned for combined statistical investigations in energy range up to 10^{10} ev. It is proposed for the first time it is proposed to join the information from such installations as huge ionization calorimeter, large emulsion chambers, high acceptance magnet spectrometer, and arrays of hodoscopic counters.

Therefore it is very actual to create a consistent multiple statistic analysis system to process both an information, concerning various characteristics of different particles /2/. It is necessary to mention that indirect nature of CR investigations requires a series of particle traversal estimations through the atmosphere and experimental setups for obtained data and theory predictions comparison.

Thus, only the combined analysis of simulation results and experimental data enables one to draw conclusions on the nature of particle interaction and the composition of the primary flux of cosmic radiation at superaccelerator energies.

The choice of statistical procedure for combined analysis is the problem of principle. It must be solved in the framework of information more general, than the theory of multinomial statistics, i.e. without the model and additional

In general, in applied statistics, the quality of the data is often measured directly in the determination of the shape of the distribution function of random variables, connected with measurement or simulation. The last column is the assumption of independently normally distributed random errors. However, as numerous investigations show, the reality is much complicated than any model, and deviations of processed data nature from proposed lead to considerable decrease in statistical procedures effectiveness.

As mentioned in /3/ in assumption making we are guided by considerations, represented a mixture of experience and physical intuition. The danger is that, if adopted assumptions turn to be wrong, the analysis remaining valid in the framework of chosen model may have highly indirect connection to reality.

This danger is overcome in the applied statistic - a technology in which the degree of adequateness of chosen model and reality is a crucial factor, that determines reliability of utilized data handling methods /4/. In applied statistic technology one tries to make more realistic assumptions. First of all it is expressed in the refuse of known parametric shape of distribution function assumption, as the complicated analytical expressions can with high accuracy describe a given random set of data, but - worse the desired intrinsic dependence of physical values /5/.

Much weaker assumption about data lies in the fact, that the features of the phenomenon under consideration are somehow

For the analysis of the problem of pattern recognition we have to take into account the following requirements:

1. The training sample must be large enough to provide a good approximation of the underlying probability distributions.
2. The training sample must be representative of the data to be analyzed.

3. The training sample must be nonparametric.

Nonparametric statistics (nonparametric methods) is a branch of statistics which does not assume that the underlying population distribution function is either well known or based on any simplifying assumptions. In practice it is often the case that the underlying distributions are not well known or are not even known at all.

3. As we see above the peculiarity of statistical decisions in cosmic ray physics consists in combined analysis of simulated and experimental data. The information on investigated physical problem contained in so called Training samples (TS) - a simulation results with controlled input parameters. TS presented all possible variations of physical values due to individual differences, admissible in the given state limits.

Training samples reflected a stochastic nature of CR interactions, existence of many strong and electromagnetic interactions channels and probabilistic character of detector operation.

It is necessary to check that both the mathematical definition of problem as well as the data handling method depend on available a priory information type. Training samples provide us only nonparametric mode of a priory information, and so our aim is to create a self-consistent nonparametric statistic technology to carry out multidimensional data processing .

Based on pattern recognition concept a package of analytical

programs for nonparametric multiple statistical analysis of CR data was developed. The use of package permits the computer realization of the following procedures:

1. Choice of theoretical model most adequately fitting experimental data.
2. Selection of events of definite type.
3. Identification of particles and interaction processes.
4. Search for cluster in multi-dimensional empirical distributions.

The principle features of package are:

- a) Utilization of Bayesian decision rules /8/.
- b) KNN (K-Nearest Neighbour) adaptive probability density estimation. /8/.
- c) Use of Bayesian risk as a closeness measure of multidimensional distributions. /8/.

The package consists of independent routines, characterized by common procedure of data definition, similar programming technique and matched resources.

The routines of package are realized in FORTRAN-4 and some modules of service data handling programs are in BEGM-6 Assembler.

The correctness and sound implementation of routines were provided by an extensive testing on problems, that are believed to be typical of those encountered in a general environment.

4. The relational type data base organization provides the unification of information storage, dependable protection, renewal of data. The information is organized in the form of named data arrays (files), and the names contain the

valuable information to make possible the associative access to data.

The main attributes of relation structures are the name of file, the number of simulation realizations (number of strings number of features, and the weight. The weight is used to compile the training samples according to definite energetic spectrum and chemical composition of primary flux.

The service software includes rather simple and convenient diagnostic message, that practically prevents incorrect utilization of data and besides allows one to efficiently make changes in the structure of stored data. There are also security facilities against unsanctional access to the information. The provision for interactive mode of operation and exploratory plotting analysis was made.

The applied program package and the relational data base are intended for use in "ANI" data analysis and technique development centre, where the information from most important experiments in the field of cosmic ray physics is planned to gather.

The accumulation and systematization of experimental and simulation data permit one to carry out the comparison of obtained data with up to date theories, to predict the behavior of physical parameters at impossible so far energies, to perform experiments optimal design.

REFERENCES.

1. Vasil'eva N.V., Sogolovskiy G.M., Yerushkin A.P. About a project of experimental investigation of particles with energy range 10^3 - 10^5 GeV (experiment "AN"), Izvest. Akad. of Scien. of Arm.SSR, vol.17, pp.12-21, 1981.
2. Chilingarian A.A. The development of statistical methods in cosmic ray physics 18 - ICRC, vol.5, p.51 - 17, Japan, 1983.
3. Bickel P.J., Doksum K.A. Mathematical statistics Holden-Day. INC. San-Francisco-Dusseldorf, Johannesburg-London-Panama-Singapure-Sydney
4. Aivazyan S.A., Yanaykov I.S., Meshalkin L.D. Applied Statistics, Finansy i statistika, Moscow, 1983.
5. Vapnik V.N. Algorithms and programmes of dependence reconstruction, Nauka, Moscow, 1984.
6. Granunder U. Regular Structures Springer-verlag New York-Heidelberg-Berlin, 1981.
7. Pearson E.S. Studies in the history of probability and statistics-Biometrika. 1965, vol.52, p.3 - 18.
8. Chilingarian A.A. On statistical methods of high energy particles identification. Proc. of Symp. on High Energy Particles Transitional Radiation, Yerevan 1984, pp. 47-434.

The manuscript was received 31 May 1985

Н.З.АКОНОВ, Св.Х.АРУТЮНЯН, С.Х.ГАЛФАЯН, М.З.ЗАЗЯН,
В.Х.МАТЕВОСЯН, А.А.ЧИЛИНГАРЯН

ОБЩИЕ ПРИНЦИПЫ И СТРУКТУРА ЦЕНТРА ДАННЫХ АНИ

Редактор Л.Н.Мукаян

Технический редактор А.С.Абрамян

Подписано в печать II/X-85г. ВФ-09038 Формат 60x84/16

Офсетная печать. Уч.изд.л. 0,5

Тираж 299 экз. Ц. 6 к.

Зак. тип. № 437

Индекс 3624

Отпечатано в Ереванском физическом институте
Ереван 36, Маркаряна 2

индекс 3624



ЕРЕВАНСКИЙ ФИЗИЧЕСКИЙ ИНСТИТУТ