

Joint paper on discriminant analysis of the confinement database

A. Kus, A. Dinklage

*Max-Planck-Institut für Plasmaphysik, IPP-EURATOM Association,
Greifswald, Germany*

CWGM 8.0 @NIFS

16 – 17 Mar. 2011

Identification of variables causing clustering in the global energy confinement data by use of discriminant analysis

A. Kus¹, A. Dinklage¹, E. Ascasibar², C.D. Beidler¹, A.A. Beletskii³, B. D. Blackwell⁴,
T. Estrada², H. Funaba⁵, J. Geiger¹, J.H. Harris^{4,6}, C. Hidalgo², M. Hirsch¹, D. Lopez-Bruna²,
A. Lopez-Fraguas², H. Maaßberg¹, T. Minami⁵, T. Mizuuchi⁷, S. Murakami⁷, N. Nakajima⁵,
S. Okamura⁵, D. Pretty⁴, M. Ramisch⁸, S. Sakakibara⁵, F. Sano⁷, U. Stroth¹, Y. Suzuki⁵,
Y. Takeiri⁵, J. Talmadge⁹, K. Thomsen¹⁰, V. Tribaldos², Yu. A. Turkin¹, K.Y. Watanabe⁵,
A. Weller¹, R. Wolf¹, H. Yamada⁵, M. Yokoyama⁵

¹*Max-Planck-Institut für Plasmaphysik, Euratom Assoc., Garching & Greifswald, Germany,*

²*Laboratorio Nacional de Fusión, Asociación Euratom/CIEMAT, Madrid, Spain,*

³*Institute of Plasma Physics, Kharkov, Ukraine,*

⁴*Australian National University, Canberra, Australia,*

⁵*National Institute for Fusion Science, Toki, Japan,*

⁶*Oak Ridge National Laboratory, Oak Ridge, USA,*

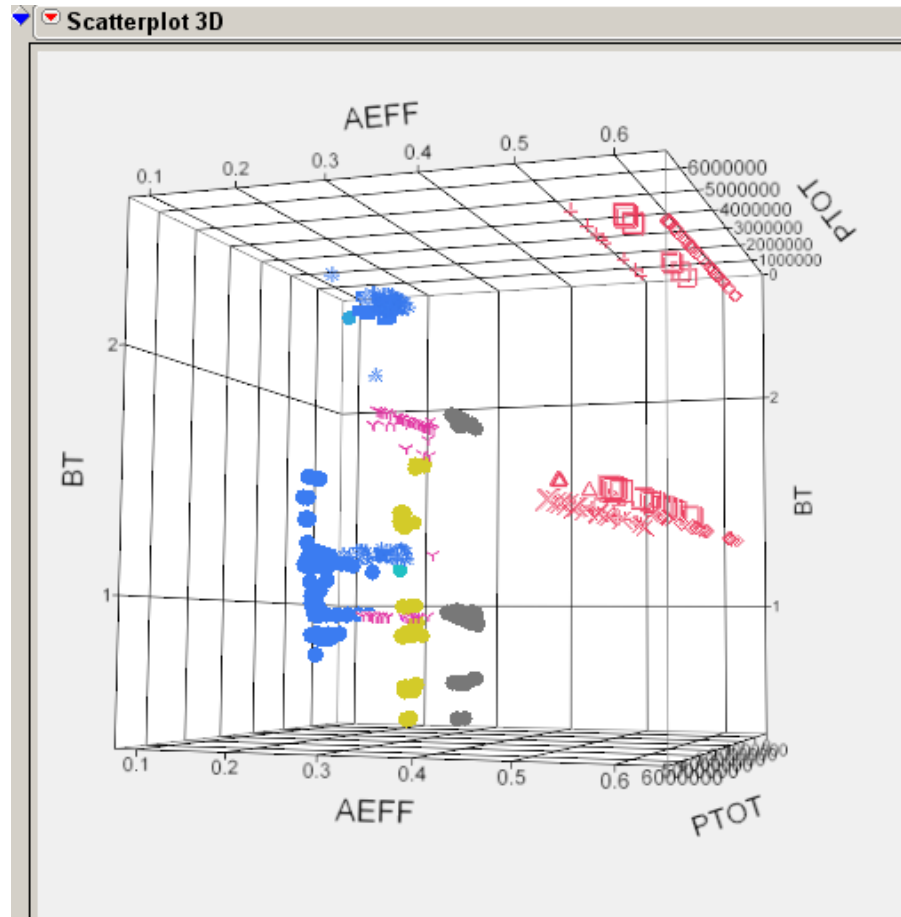
⁷*Kyoto University, Kyoto, Japan,*

⁸*Universität Stuttgart, Stuttgart, Germany,*

⁹*University of Wisconsin, Madison, USA,*

¹⁰*European Commission, DG Research J.4, Brussels, Belgium*

Example view of the 14 ISS04 subgroups



Questions:

- 1) Do exist any “natural” subgroups in the data?
- 2) Which quantities cause the clustering?

CLUSTER ANALYSIS

DISCRIMINANT ANALYSIS

Findings so far (a preliminary ISS04 datasets analysis)

- Cluster analysis (in the 7-dim. space spanned by LOG_TAU, LOG_A, ..., LOG_I) reveals the existence of cohesive subgroups (clusters) in the data
- Clusters differ from ISS04 subgroups but resemble ISS04 grouping roughly
- Saturation in regression coefficients when increasing the number of clusters
 - a. No clustering (all data in one cluster/one group): differences to ISS04
 - b. There is a kind of saturation in the number of clusters: a regression analysis on data divided into ca. 11-18 clusters provides results similar to ISS04 scaling formula

Reminder: Linear discriminant function analysis

Linear discriminant function analysis

aims at determination of linear combinations of independent variables (*predictors*) that discriminate among the *categories* of the grouping variable.

A linear combination of predictor variables X_1, \dots, X_p , called **discriminant function** is constructed such that it assigns its values into two subgroups that differ as much as possible.

In the most simple case, with two subgroups (e.g. L/H mode data), there exists only one discriminant function

$$D = b_1 X_1 + b_2 X_2 + \dots + b_p X_p \quad (\text{for data previously standardized})$$

A model must be correctly defined (all relevant variables involved, e.g. using „stepwise variable selection“ procedure).

Reminder: Linear discrim. function analysis (cont.)

In a general case the maximal number of discriminant functions is equal

- number of subgroups minus one, or
- number of used variables,

whatsoever is smaller.

Properties of discriminant functions:

- All discriminant functions are pairwise orthogonal (uncorrelated)
- Viewing the coefficients ***b***'s one can see how the predictor variables contribute to the discrimination: **the larger the *b* the larger the contribution**
(This is valid only for standardized data.)

Implementation of discriminant analysis

• Learning process

Determination of discriminant functions (model development):

Which quantities decide about the affiliation of an observation to the particular group (*what distinguish high-beta discharges from others?*)

• Prediction

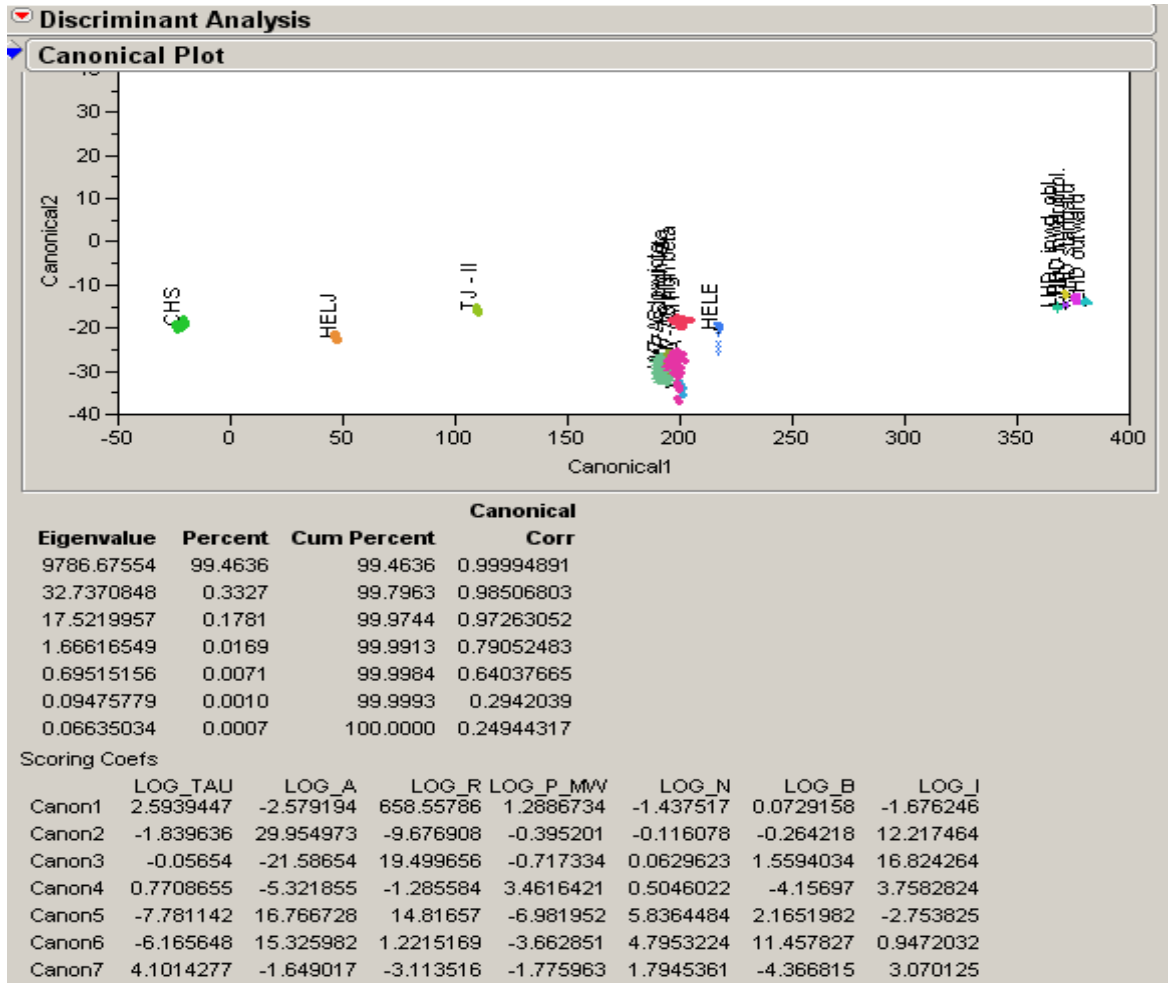
Application (after a discrim. model was developed):

Assignment of a new observation to the particular group

Examples

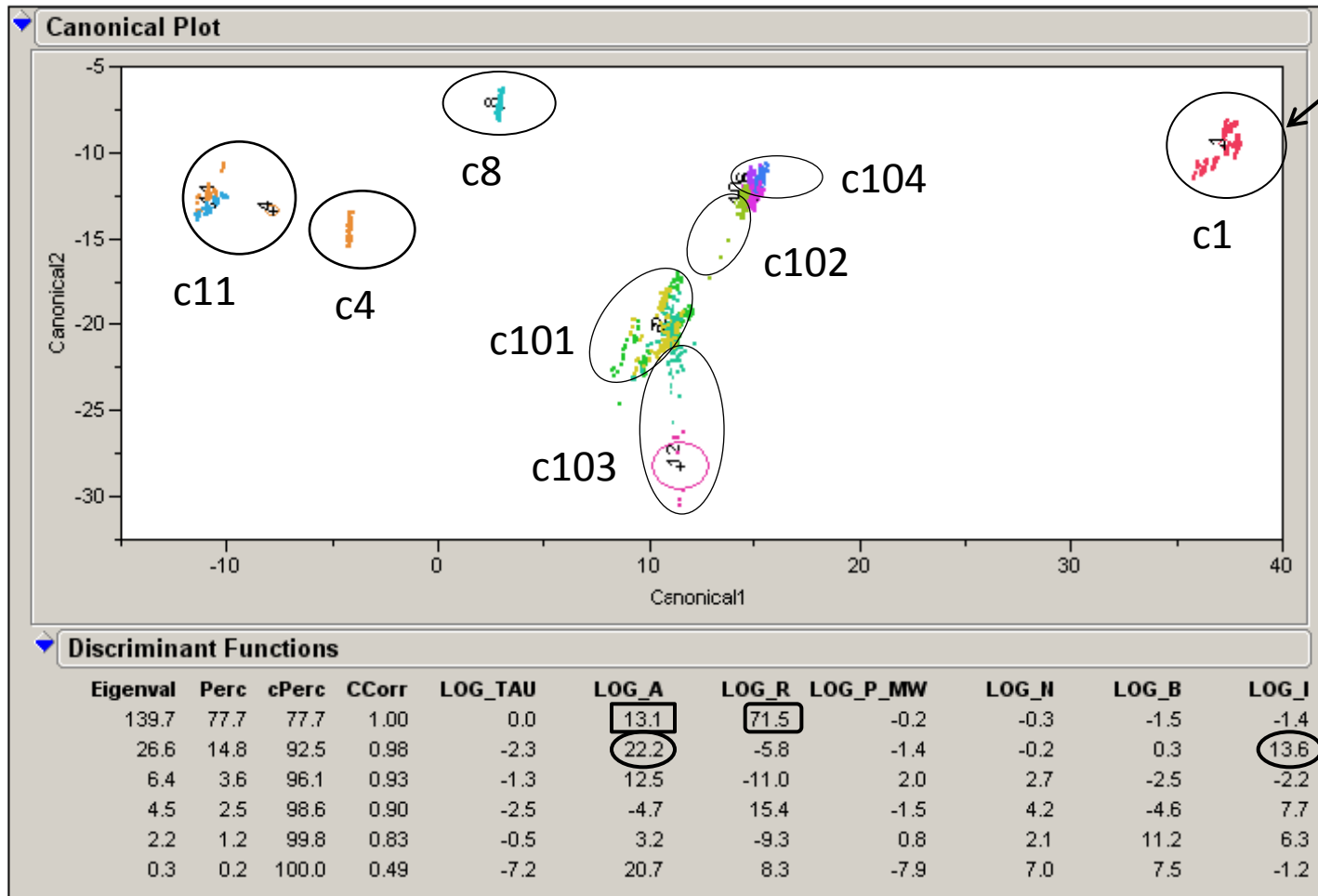
- Plasma physics: an observation with certain parameters is assigned to L or H-Mode
- Medicine: particular risk factors cause a specific disease
- (Psychology, war against terrorism, ...)

ISS04 subgroups in the first two discrim. dimensions



5-6 larger groups,
determined by LOG_R

ISS04 data clusters along the first two discrim. functions



- Cluster c1
LHD
- c4
CHS (68), HELJ
- c8
TJ-II
- c11
CHS (128)
- c101
W7-A,
W7-AS h.+l. iota
- c102
HELE
- c103
W7-AS h. beta
- c104
ATF

- Geometry is crucial for the discrimination (thesis valid only for the used model)
- 92.5 percent of the total variation happens in plane spanned by the first two discriminant functions
- Cluster c103 (W7-AS high-beta) overlaps other W7-AS data (An additional/latent discrim. factor exist?)

To do till EPS:

- Define a subset of „suspected“ variables that may cause clustering
An automatic procedure is exponentially time consuming:
10 variables → 1023 (= $2^{10}-1$) combinations (models),
20 variables → 1 048 575 combinations,
104 numerical variables in ISHCDB_25 → $2 \cdot 10^{31}$ possible models!
- Select a dataset for analysis
Which group of data should indentified at first?
(High beta?)
- Perform analysis