# Unsupervised Feature Selection Based on the Morisita Index

## Jean Golay and Mikhail Kanevski

Institute of Earth Surface Dynamics (IDYST)
Faculty of Geosciences and Environment (FGSE)
University of Lausanne, Switzerland
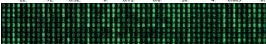
EGU General Assembly 2016

# Introduction

**High-Dimensional Data Sets**

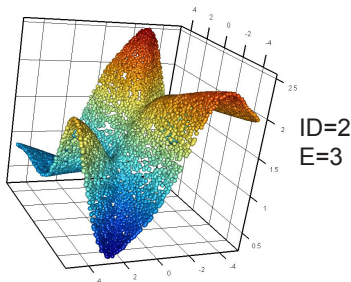←Variables / Features→



Instances

**Issues**

1. Curse of dimensionality
2. Computer performance
3. Data visualization
4. Interpretability of the results

**Solutions**

1. PCA
2. MDS
3. etc.

**A New Solution**

1. The concept of Intrinsic Dimension (ID)
2. The Morisita estimator of ID
3. An ID-based algorithm for selecting the smallest subset of features conveying all the information content of a data set
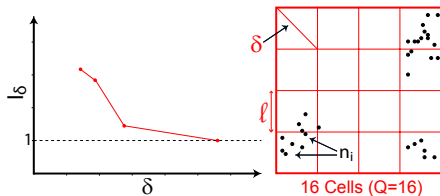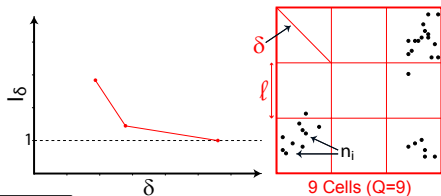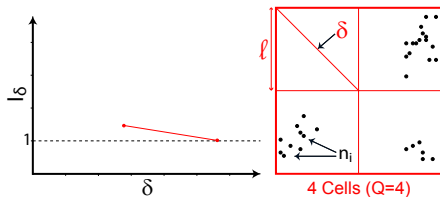


ID=2
E=3

# Outline

# Outline

# Calculation

$$I_\delta = Q \, \frac{\sum_{i=1}^{Q} n_i(n_i - 1)}{N(N-1)}$$

# Outline

## Calculation

$$M_2 = E - S_2$$

The concept of ID is extended to non-integer (fractal) dimensions.



$S_2 = 0.76$

$M_2 = 2 - 0.76 = 1.24$

9 cells

4 cells

16 cells

1 cell

$Log_e(l\ell)$

$Log_e(1/\ell)$

$F_2$

$F_1$

$n_i$

$\delta$

$\ell$

# Outline

## ID and Redundancy

$V_1$ and $V_2$ are two uniformly distributed variables and one has that:
$ID(V_1, V_2) \approx ID(V_1) + ID(V_2) \approx 1 + 1 = 2$ (see (A))
$ID(V_1, V_2) \approx ID(V_1) \approx ID(V_2) \approx 1$ (see (B) and (C))



No Redundancy     Redundancy (linear)     Redundancy (non-linear)

Redundant features (variables) do not contribute to the data ID
Idea (Traina's work): select the features which increase the data ID.

# The Proposed Algorithm

$A = \{F1, F2, F3, F4\}$ and $M_2(A) = 2.20$

Step 1
$|2.20 - M_2(F1)| = 1.20$
$|2.20 - M_2(F2)| = 1.34$
$|2.20 - M_2(F3)| = 1.30$
$|2.20 - M_2(F4)| = 1.19$

Step 2
$|2.20 - M_2(F4,F1)| = 1.14$
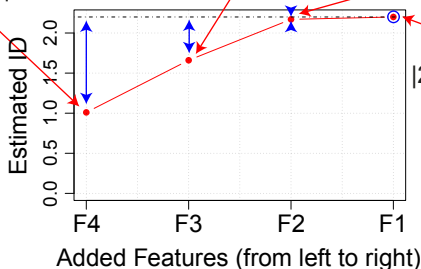$|2.20 - M_2(F4,F2)| = 0.59$
$|2.20 - M_2(F4,F3)| = 0.54$

Step 3
$|2.20 - M_2(F4,F3,F1)| = 0.52$
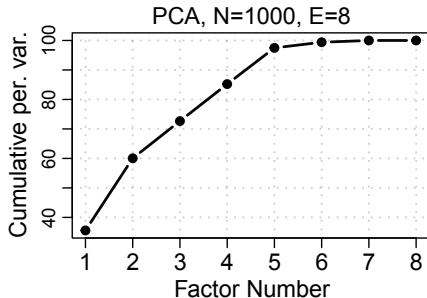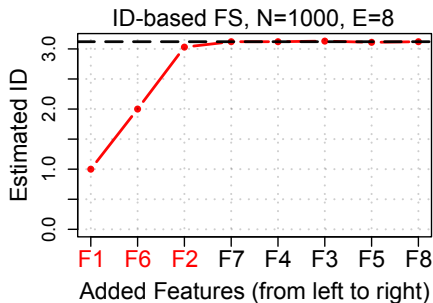$|2.20 - M_2(F4,F3,F2)| = 0.03$

Step 4
$|2.20 - M_2(F4,F3,F2,F1)| = 0.00$



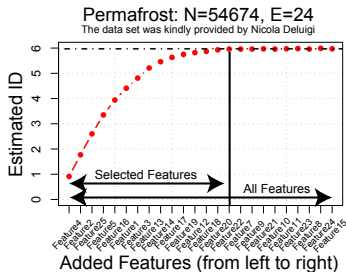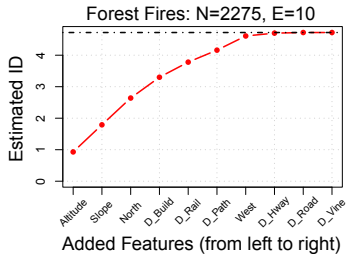F1 is redundant, since it hardly contributes to increasing the data ID

# A Simulated Case Study: the Input Space of the Butterfly Data Set

$$F_1, F_2, F_6 \in ]-5, 5[ \qquad F_3 = \log_{10}(F_1 + 5) \qquad F_4 = F_1^2 - F_2^2$$
$$F_5 = F_1^4 - F_2^4 \qquad F_7 = \log_{10}(F_6 + 5) \qquad F_8 = F_6 + F_7$$

# Real Case Studies I



Windspeed: N=127, E=11

Forest Fires: N=2275, E=10

Permafrost: N=54674, E=24
The data set was kindly provided by Nicola Deluigi

The information content was assessed by means of random forest

|  | All Features | Selected Features |
|---|---|---|
| Error rate (%) | 14.93 (3.13) | 14.40 (2.63) |

# Real Case Studies II



Hyperspectral Image of Pavia: N=50000 (sample), E=102

# Outline

# Conclusion

### Take-Home Message

The concept of Intrinsic Dimension (ID) can help find solutions to the issues raised by large data sets.

📄 C. Traina Jr., A. J. M. Traina, L. Wu, C. Faloutsos, Fast feature selection using fractal dimension, *Proceedings of the XV Brazilian Symposium on Databases (SBBD)*, pp. 158-171, 2000.

📄 J. Golay, M. Leuenberger, M. Kanevski, Feature Selection for Regression Problems Based on the Morisita Estimator of Intrinsic Dimension, *arXiv:1602.00216*, 2016.

📄 J. Golay and M. Kanevski, A new estimator of intrinsic dimension based on the multipoint Morisita index, *Pattern Recognition*, 48(12):4070-4081, 2015.

📄 J. Golay, M. Kanevski, C. D. Vega Orozco and M. Leuenberger, The multipoint Morisita index for the analysis of spatial patterns, *Physica A*, 406:191-202, 2014.