

Inlierness, Outlierness, Hubness and Discriminability: an Extreme-Value-Theoretic Foundation

Michael E. HOULE

データの複雑さはどのように測定するのか？

現在、データマイニングは統一理論がまだ提案されていない。個々の問題のためには、分類また、クラスタリングなどの多くのアドホック技術が設計されている。我々は、さまざまな基本的な機械学習とデータマイニングタスクを結びつける理論的な枠組みを提案する。

THE PROBLEM

To date, no unifying theory of data mining has been proposed.

- ▶ Many ad-hoc techniques have been designed for individual problems, such as classification or clustering.
- ▶ Solutions involve much invention and reinvention, with few guidelines.

A theoretical framework that ties together different fundamental machine learning and data mining tasks (including indexing, clustering, classification, data discriminability, subspace methods, etc.) could help the discipline, and serve as a basis for future investigation.

THE CURSE OF DIMENSIONALITY

As the number of object features (data dimensionality) rises:

- ▶ Similarity values concentrate around their expected values.
- ▶ Items become less and less distinguishable.
- ▶ Data analysis based similarity (e.g. clustering and classification) becomes ineffective.

Some sets have higher *intrinsic dimensionality* (ID) than others.

- ▶ Intuitively, the minimum number of dimensions or features with which the data can be represented with minimal distortion.
- ▶ Many formalizations have been proposed (such as the Hausdorff dimension, in 1918!).

ID AND SCALABILITY

Implications for Big Data:

- ▶ Data mining is greatly concerned with what happens in neighborhoods of data (clustering, classification, outlier detection, ...).
- ▶ As the number of objects increases, the k -nearest neighbor (k -NN) distance tends to 0.
- ▶ Indiscriminability of neighborhood distances, and ID of k -NN query result, tend to $ID_{F_X}(0)$.

Limit effect characterizes the complexity of data.

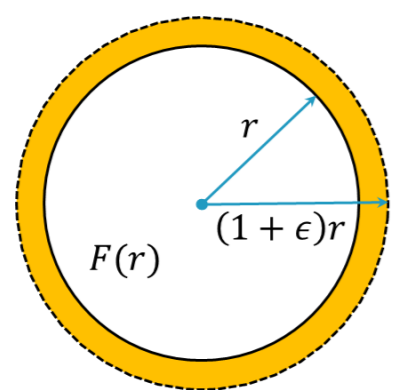
ID AND INLIERNESS / OUTLIERNESS

If $ID_{F_X}(r) < ID_{F_X}(0)$ within neighborhood $0 < r < \epsilon$ of some point \mathbf{p} , then:

- ▶ The growth rate at distance r is less than that which would be expected within a uniform distribution of dimension $ID_{F_X}(0)$.
- ▶ The drop in indiscriminability (rise in discriminability) indicates a decrease in local density as the distance from \mathbf{p} increases.
- ▶ The relationship between \mathbf{p} and its neighborhood is therefore that of an *inlier*.

If instead $ID_{F_X}(r) > ID_{F_X}(0)$, \mathbf{p} is an *outlier*.

DISCRIMINABILITY OF DISTANCES



- ▶ Let X be an absolutely continuous random distance variable with c.d.f. F_X and p.d.f. f_X .
- ▶ Discriminability of the distance measure can be regarded as a ratio between two quantities as the distance expands infinitesimally:
- ▶ (1) the relative increase in distance, and
- ▶ (2) in probability measure.

The indiscriminability of X at distance r is:

$$\text{InDiscr}_X(r) = \lim_{\epsilon \rightarrow 0^+} \left(\frac{F_X((1+\epsilon)r) - F_X(r)}{\epsilon \cdot F_X(r)} \right)$$

LOCAL ID

When $F_X(r) > 0$, the local intrinsic dimensionality of X at distance r is defined as:

$$\text{IntrDim}_X(r) = \lim_{\epsilon \rightarrow 0^+} \left(\frac{\ln F_X((1+\epsilon)r) - \ln F_X(r)}{\ln(1+\epsilon)} \right)$$

This definition is an extension for continuous distance distributions of the Expansion Dimension.

THEOREM (EQUIVALENCE OF ID AND INDISCRIMINABILITY)

Let X be an absolutely continuous random distance variable. If F_X is both positive and differentiable at r , then

$$\text{IntrDim}_X(r) = \text{InDiscr}_X(r) = \frac{r \cdot f_X(r)}{F_X(r)} =: ID_{F_X}(r).$$

THEOREM (ID REPRESENTATION FORMULA)

Let X be an absolutely continuous random distance variable such that $F_X(r) > 0$ whenever $r > 0$. Then for any $r, w \in (0, z)$,

$$F_X(r) = F_X(w) \cdot \left(\frac{r}{w} \right)^{ID_{F_X}(0)} \cdot G_{F_X,0,w}(r), \text{ where}$$

$$G_{F_X,0,w}(r) := \exp \left(\int_r^w \frac{ID_{F_X}(0) - ID_{F_X}(t)}{t} dt \right).$$

Furthermore, for any fixed $0 < c < 1$, we have

$$\lim_{\substack{w \rightarrow 0^+ \\ cw \leq r \leq w}} G_{F_X,0,w}(r) = 1.$$

EXTREME VALUE THEORY

- ▶ Profound importance in risk analysis, economics, civil engineering, operations research, material sciences, geophysics, ...
- ▶ Here, adapted for the lower tails of distance distributions.
- ▶ One of the three fundamental pillars of Extreme Value Theory, Karamata Characterization Theorem (1930): $F_X(x) = x^{\gamma_X} \ell_X(1/x)$ for some constant γ_X , where

$$\ell_X(1/x) = \exp \left(\eta_X(1/x) + \int_x^w \frac{\epsilon_X(1/t)}{t} dt \right).$$

ID AND EXTREME VALUE THEORY

- ▶ ID Representation is a more precise formulation of the Karamata Characterization, with:

$$\begin{aligned} \gamma_X &= ID_{F_X}(0); \\ \eta_X(1/x) &= \ln F_X(w) - ID_{F_X}(0) \ln w; \\ \epsilon_X(1/t) &= ID_{F_X}(0) - ID_{F_X}(t). \end{aligned}$$

- ▶ $ID_{F_X}(0)$ is the well-studied EVT index γ_X .
- ▶ Connections also exist between ID and Hausdorff dimension, and ID and the hubness phenomenon in data.

2ND-ORDER ID

- ▶ Inlierness / outlierness is determined by the sign of $ID'_{F_X}(r)$ as $r \rightarrow 0^+$.
- ▶ Strength is obtained by normalizing $ID'_{F_X}(r)$ for distance and intrinsic dimensionality:

$$ID_{ID_{F_X}}(r) = \frac{r \cdot ID'_{F_X}(r)}{ID_{F_X}(r)} = ID_{F_X}'(r) + 1 - ID_{F_X}(r),$$

- ▶ However, $ID_{ID_{F_X}}(0) = 0$ always ... need the growth rate $ID_{|ID_{F_X}|}(0)$ of $|ID_{F_X}(r)|$ instead.

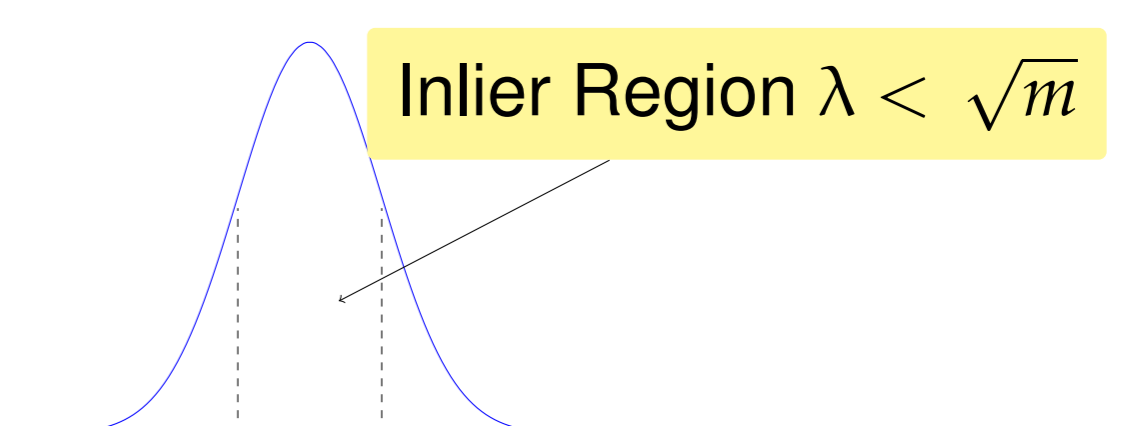
EXAMPLE: DISTANCES TO A GAUSSIAN

- ▶ Vector of normally distributed random variables with means μ_i and variances σ_i^2 .
- ▶ Distance from $\mathbf{0}$ to a point $\mathbf{X} = (X_1, X_2, \dots, X_m)$, defined as

$$Z = \sqrt{\sum_{i=1}^m \frac{X_i^2}{\sigma_i^2}}, \text{ and } \lambda = \sqrt{\sum_{i=1}^m \frac{\mu_i^2}{\sigma_i^2}}$$

is the normalized distance from $\mathbf{0}$ to the Gaussian mean.

- ▶ $ID_{F_Z}(0) = m$ and $ID_{|ID_{F_Z}|}(0) = 2$ whenever $\lambda \neq \sqrt{m}$. Also, as r tends to 0, $ID_{ID_{F_Z}}(r) > 0$ when $\lambda > \sqrt{m}$ (tail region \rightarrow outliers), and < 0 when $0 \leq \lambda < \sqrt{m}$ (central region \rightarrow inliers).
- ▶ Normalization works! Independent of λ & m .



REFERENCES

- [1] M. E. Houle. "Inlierness, Outlierness, Hubness and Discriminability: an Extreme-Value-Theoretic Foundation", NII Technical Report NII-2015-002E.
- [2] M. E. Houle. "Dimensionality, discriminability, density & distance distributions", ICDMW 2013.