

Interpretability of Big-Data Analytics

Presented by: Aindrila Ghosh

**** The research objective is based on our paper:***

A. Ghosh, M. Nashaat, J. Miller, S. Quader, and C. Marston, "A Comprehensive Review of Tools for Exploratory Analysis of Tabular Industrial Datasets.", In Visual Informatics, Elsevier, December 2018, Volume 2, Issue 4, Pages 235-253.



① Opportunity: Why Interpretability?

Interpret Black-Box Models



- Necessity to build trust for Analytics.
- Most machine learning models and data analysis algorithms are black-box.
- Users do not use what they don't understand or trust.

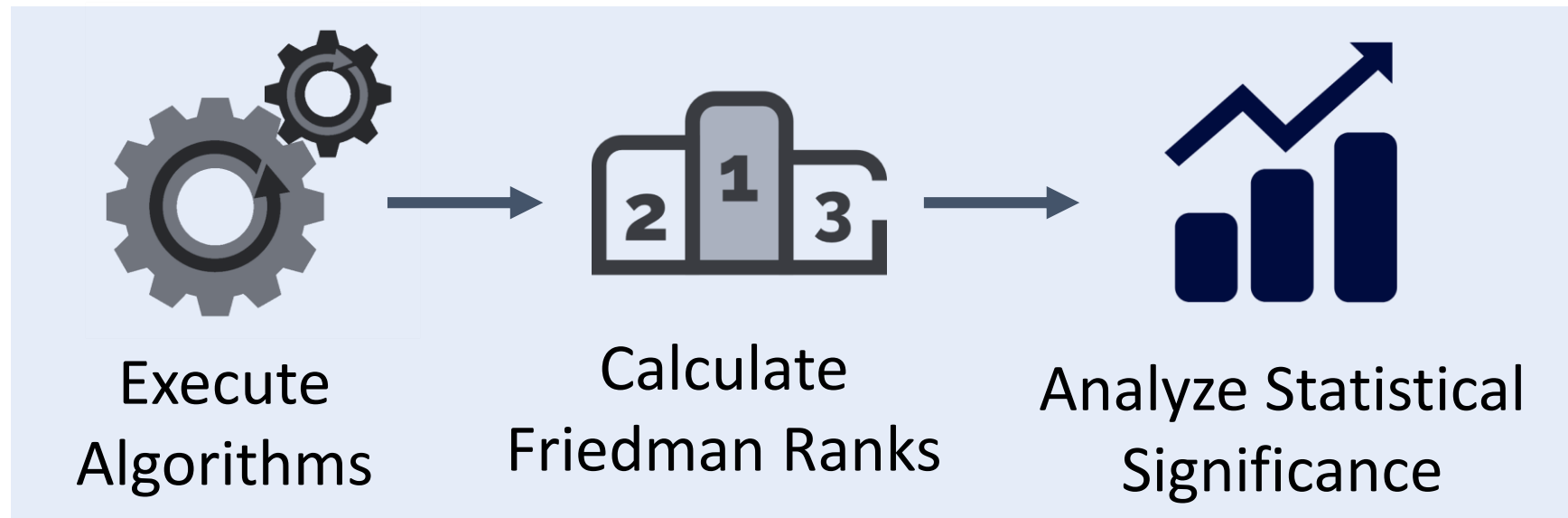
- Necessity for Human-in-the-loop.
- Need to incorporate user expertise in decisions.
- Essential to proactively guide users in the complex data analysis.



Engage Users in Decision Making

② Experimental Methodology

- Compare 15 Dimensionality Reduction Algorithms
- For 7 Contextual Evaluation Metrics
- Using over 30 Real-world Datasets
- Evaluate with 6 Statistical Significance Tests



3 Experimental Results

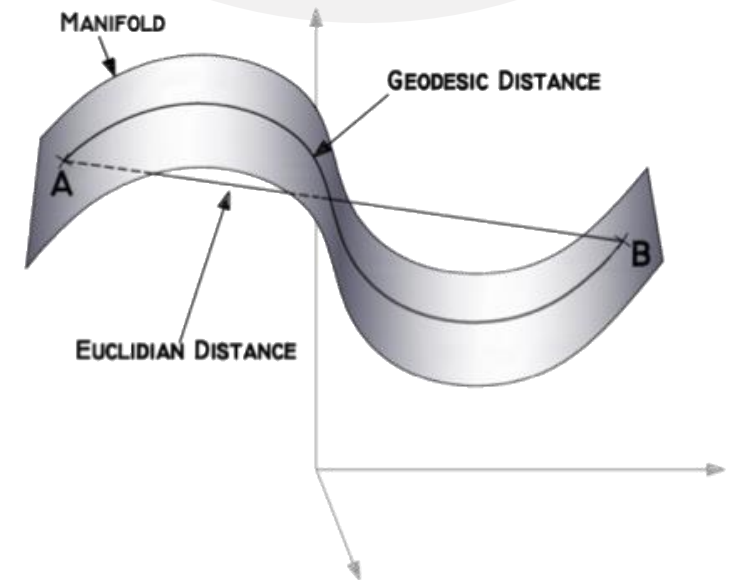
Evaluation Metric	Best	Mediocre	Worst
ML Accuracy	KernelPCA, PCA	Fit-SNE, LEM	LTSA, HLLC
Execution Time	PCA, Isomap	openTSNE, LTSA	MDS, LEM
Local Structure	MDS, openTSNE	Fit-SNE, UMAP	LLE, Isomap
Global Structure	MDS, KernelPCA	LEM, HLLC	Trimap, t-SNE
Outlier Effects	LTSA, Isomap	t-SNE, openTSNE	LLE, MLLC
Duplicate Effects	t-SNE, Trimap	HLLC, LEM	MDS, KernelPCA
Partial Records	PCA, KernelPCA	UMAP, Trimap	Fit-SNE, t-SNE

④ New Algorithm - IDLE

- Even with generic guidelines, dimensionality reduction lacks in interpretability.
- Most real-world data sets are distributed over non-linear manifolds.
- Hence linear distance among data-points does not project their actual distances.

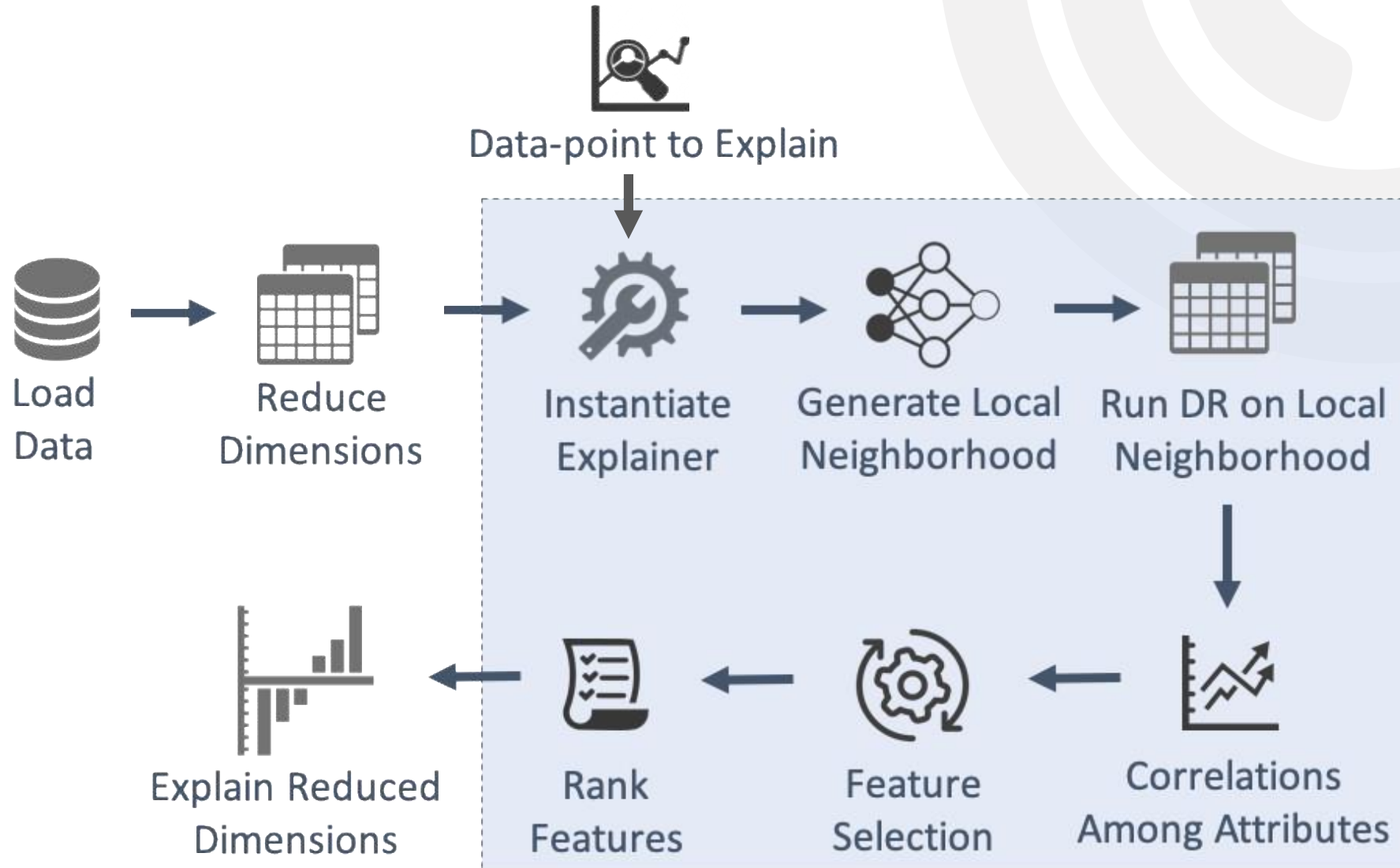
Proposed Solution:

IDLE: Interactive Descriptions for Low-dimensional Embedding

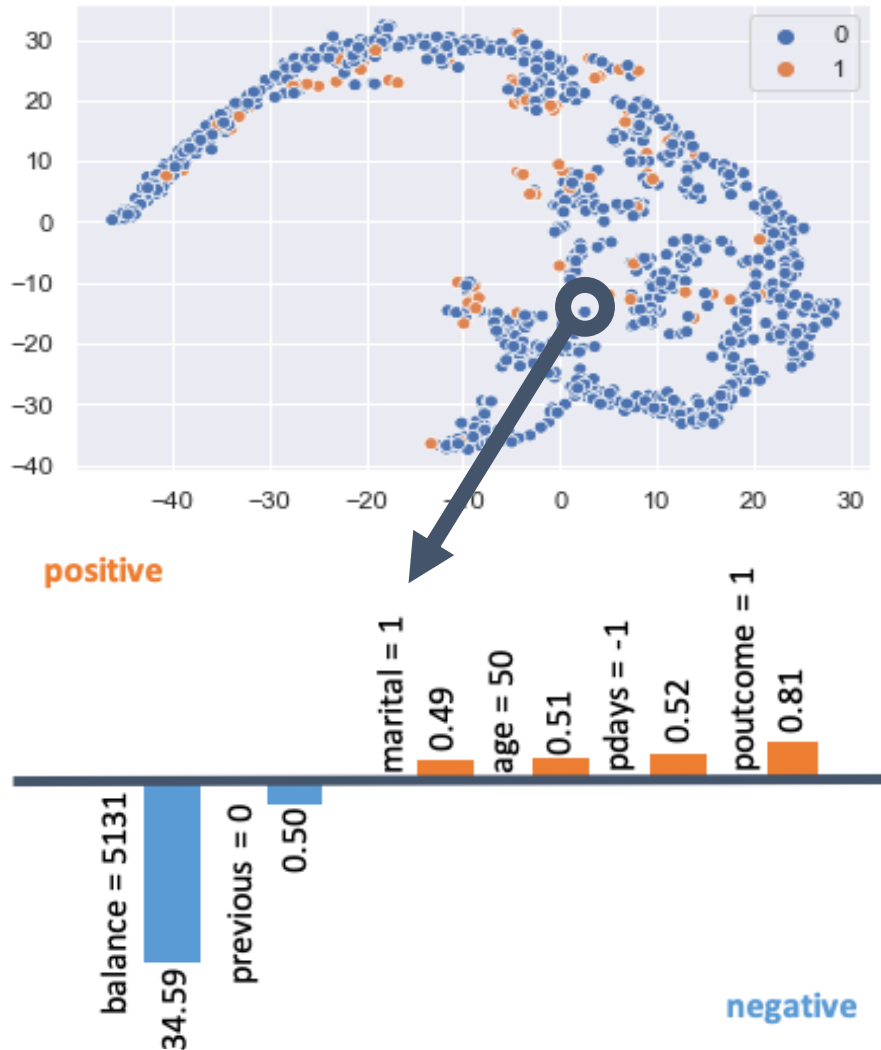


5

IDLE: Interactive Descriptions for Low-dimensional Embedding



⑥ Example – Explaining t-SNE*



Dataset: Bank (Source: UCI)

- 45,211 data-points
- 17 attributes
- Interactive Selection of one data-point
- Highly Influencing Attributes: 6
- Positive Influences: 4 attributes
- Negative Influences: 2 attributes

* *t-Distributed Stochastic Neighbor Embedding*



THANK YOU

 aindrila@ualberta.ca