



# Modeling Heterogeneous Statistical Patterns In High-dimensional Data By Adversarial Distributions: An Unsupervised Generative Framework

Han Zhang<sup>1</sup>, Wenhao Zheng<sup>2</sup>, Charley Chen<sup>1</sup>, Kevin Gao<sup>1</sup>, Yao Hu<sup>2</sup>, Ling Huang<sup>4</sup>, Wei Xu<sup>1</sup>

<sup>1</sup>Institute for Interdisciplinary Information Sciences, Tsinghua University

<sup>2</sup>Alibaba Youku Cognitive and Intelligent Lab, <sup>4</sup>AHI Fintech



交叉信息研究院  
Institute for Interdisciplinary  
Information Sciences

## Motivation

- The intrinsic clusters in high-dimensional data may display heterogeneous statistical patterns.
- Specifically, different clusters display clustering patterns w.r.t. different features.
- For example, different fraud groups may share the IP address, phone number, ID card number etc.

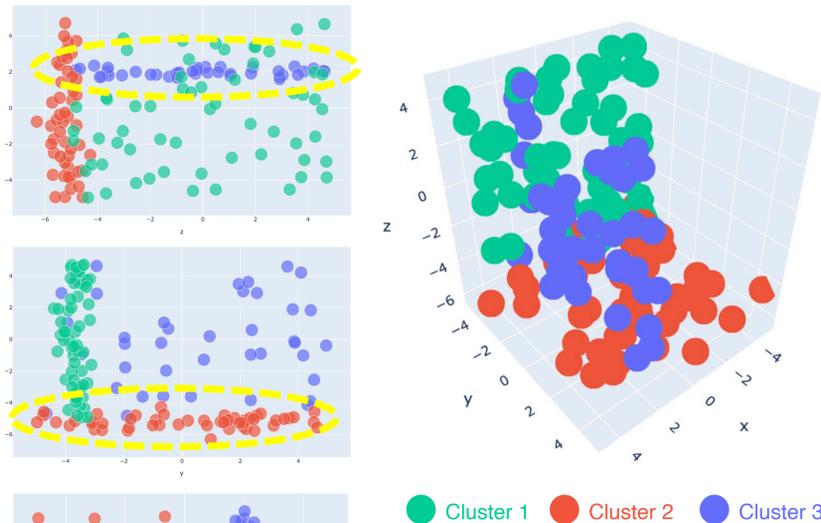


Fig 1. Heterogeneous patterns of high-dimensional data. Left: data projected to 2D planes. Right: data in the 3D space.

## FIRD: a generative framework

- **Feature Independent assumption:** within each data cluster, the features are independent with each other.
- **Adversarial Distributions:** for each feature within a cluster, a pair of distribution **compete** with each other for generating the observations.

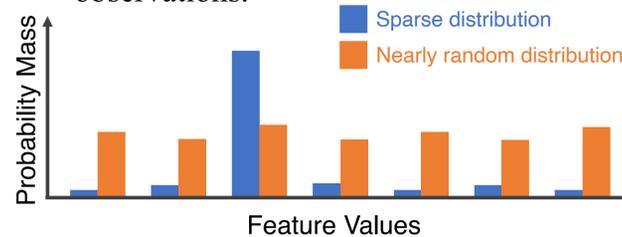


Fig 2. Example of discrete adversarial distributions. The **sparse distribution** models the fraud, and the **random part** models normal user.

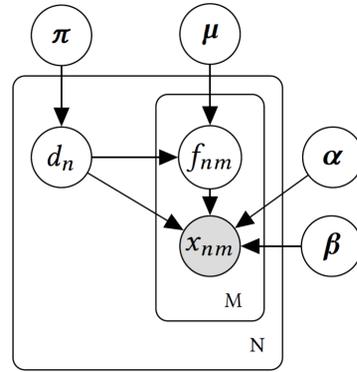


Fig 3. The plate representation of FIRD in discrete space. The parameter  $\pi$  is the mixture weight.  $\mu$  balances the adversarial distribution pairs. The synchronization and randomness are captured by the adversarial distribution pairs, whose parameters are  $\alpha$  and  $\beta$ , respectively.

- The generation process for an observation  $x_n$ :
  1. Choose cluster  $d_n \sim \text{Multinomial}(\pi)$
  2. For each feature  $F_m$ :
    - a) Choose the indicator variable  $f_{nm} \sim \text{Bernoulli}(\mu_{d_{nm}})$ ;
    - b) If  $f_{nm} = 1$ , choose  $x_{nm} \sim \text{Multinomial}(\alpha_{d_{nm}})$ ;
    - c) If  $f_{nm} = 0$ , choose  $x_{nm} \sim \text{Multinomial}(\beta_{d_{nm}})$ ;

- The objective function (log-likelihood):

$$\log \mathcal{L}(\pi, \mu, \alpha, \beta; \mathcal{D}, \lambda^{(1)}, \lambda^{(2)})$$

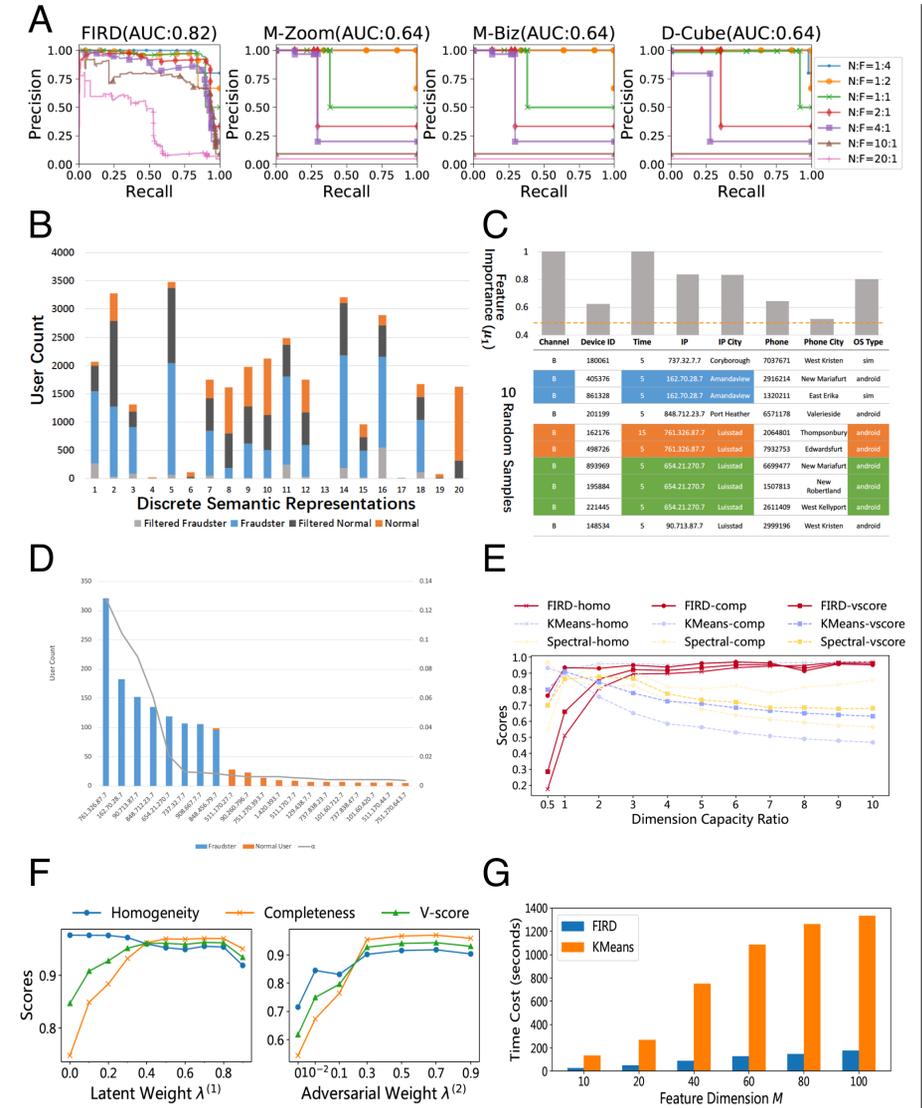
$$= \sum_{n=1}^N \log \left\{ \sum_{d_n, f_n} p(x_n | f_n, d_n, \alpha, \beta) p(f_n | d_n, \mu) p(d_n | \pi) \right\}$$

$$- \sum_{g=1}^G \lambda_g^{(1)} \log \pi_g - \sum_{g=1}^G \sum_{m=1}^M \sum_{i=1}^{D_m} \lambda_{gmi}^{(2)} (\log \alpha_{gmi} - \log \beta_{gmi})$$

Promote Sparsity
Promote Randomness

## Results

- We evaluate FIRD on 3 types of datasets:
  1. E-commerce platform registration dataset;
  2. anomaly detection benchmark dataset;
  3. synthetic datasets according to FIRD's generation process.
- We show that FIRD is able to:
  1. detect fraud groups much better than comparison methods;
  2. work as a general anomaly detection method;
  3. provide significant performance with low time cost.



## H

Dataset	FIRD	HBOS	IForest	OCSVM	LSCP
cardio	<b>0.949</b>	0.843.	0.924	0.938	0.901
musk	<b>1.000</b>	<b>1.000</b>	0.999	<b>1.000</b>	0.998
optdigits	<b>1.000</b>	0.865	0.714	0.500	-
satimage-2	<b>0.998</b>	0.977	0.993	0.997	0.9935
shuttle	0.990	0.986	<b>0.997</b>	0.992	0.5514
satellite	<b>0.900</b>	0.754	0.701	0.660	0.6015
ionosphere	<b>0.946</b>	0.5569	0.8529	0.8597	-
pendigits	<b>0.972</b>	0.9247	0.9435	0.931	0.8744
wbc	0.944	<b>0.954</b>	0.9325	0.9376	0.945