고려대학교
빅데이터 연구회

# KU-BIG

## Outlier Detection

유현우 박정진 정희정 송예은 심정은 양수형

# 목 차

# PART. I Outlier / Novelty / Anomaly

**1**

**Outlier**

**Novelty**

**Anomaly**

**2**

**Outlier**

**(Anomaly)**

**Detection**

**3**

**Novelty**

**Detection**

# 1) Outlier / Novelty / Anomaly

**Outliers** are also referred to as abnormalities, discordants, deviants, or **anomalies** in the data mining and statistics literature.

(Source: "Outlier Analysis" (Springer), Charu Aggarwal, 2017, http://charuaggarwal.net/outlierbook.pdf)

## Outlier = Anomaly

# 2) Outlier(Anomaly) detection

The training data **contains outliers**

which are defined as observations that

are far from the others.



Copyright 2014, Laerd Statistics.

(Source: https://scikit-learn.org/stable/modules/outlier_detection.html)

## 2) Outlier(Anomaly) detection

i)  Unsupervised anomaly detection

ii) Supervised anomaly detection

iii) Semi-Supervised

# 3) Novelty detection

The training data **is not polluted by outliers** and we are interested in detecting whether a **new observation is an** outlier.

**Novelty**

# PART. Ⅱ Evaluate Measure

**1**

**Metric**

**For**

**Out-of-Distribution Detection**

**2**

**Better metric**

**For**

**Class-imbalanced data**

# 1) Metric for Out-of-Distribution Detection

| | PREDICTED CLASS | |
|---|---|---|
| | Class=Yes | Class=No |
| ACTUAL CLASS Class=Yes | a | b |
| Class=No | c | d |

**a: TP** (True Positive)

**b: FN** (False Negative)

**c: FP** (False Positive)

**d: TN** (True Negative)

$$\text{Accuracy} = \frac{a+d}{a+b+c+d} = \frac{TP+TN}{\text{TP}+\text{FN}+\text{FP}+TN}$$

(Source: JunGeol Baek, 2019 1st semester Data mining chapter 3. pp.79.)

# 1) Metric for Out-of-Distribution Detection

| | PREDICTED CLASS | |
|---|---|---|
| | Class=Yes | Class=No |
| **ACTUAL CLASS** Class=Yes | a | b |
| Class=No | c | d |

a: **TP** (True Positive)

b: **FN** (False Negative)

c: **FP** (False Positive)

d: **TN** (True Negative)

$$\text{Precision} = \frac{a}{a+c} = \frac{TP}{\text{TP}+\text{FP}}$$

$$\text{Recall} = \frac{a}{a+b} = \frac{TP}{\text{TP}+\text{F}N}$$

(Source: JunGeol Baek, 2019 1st semester Data mining chapter 3. pp.93.)

# 1) Metric for Out-of-Distribution Detection

| | PREDICTED CLASS | |
|---|---|---|
| | Class=Yes | Class=No |
| ACTUAL CLASS | | |
| Class=Yes | a | b |
| Class=No | c | d |

**a: TP** (True Positive)

**b: FN** (False Negative)

**c: FP** (False Positive)

**d: TN** (True Negative)

$$\text{TPR} = \frac{a}{a+c} = \frac{TP}{\text{TP}+\text{FP}}$$

$$\text{FPR} = \frac{c}{c+d} = \frac{FP}{\text{FP}+\text{T}N}$$

(Source: JunGeol Baek, 2019 1st semester Data mining chapter 3. pp.84.)

# 1) Metric for Out-of-Distribution Detection



X축 : FPR = FP / (FP + TN)

Y축 : TPR = TP /(TP + FN)

Diagonal line = Random Guessing

Area under ROC curve = AUC

# 1) Metric for Out-of-Distribution Detection



Area under ROC curve = AUC

AUC Range : [0, 1]

100% 맞는 예측 모델일 경우 AUC = 1.

# 2) Better metric for class-imbalanced data



PR_AUC Curve

X축 : Recall : TP / (TP + FN)

Y축 : Precision : TP /(TP + FP)

In Case of imbalanced-Data

⇒ Precision이 FPR에 비해 False Positive를

더 민감하게 잡아낼 수 있다.

⇒ Imbalanced data에서 효과적인 metric!

# PART. Ⅲ EDA

1

**Description of Data**

2

**Non-linear Relations**

# 1) Description of Data

Credit Card Dataset : (From Kaggle)

- Highly unbalanced data

- 492 frauds out of 284,807 transactions



Proportion of normal/abnormal transactions

```
0     284315
1        492
Name: Class, dtype: int64
```

# 1) Description of Data

Input variable :

- 28 input variable V, result of PCA transformation.

- Time : seconds ( about 48 hours )

- Amount : transaction amount.

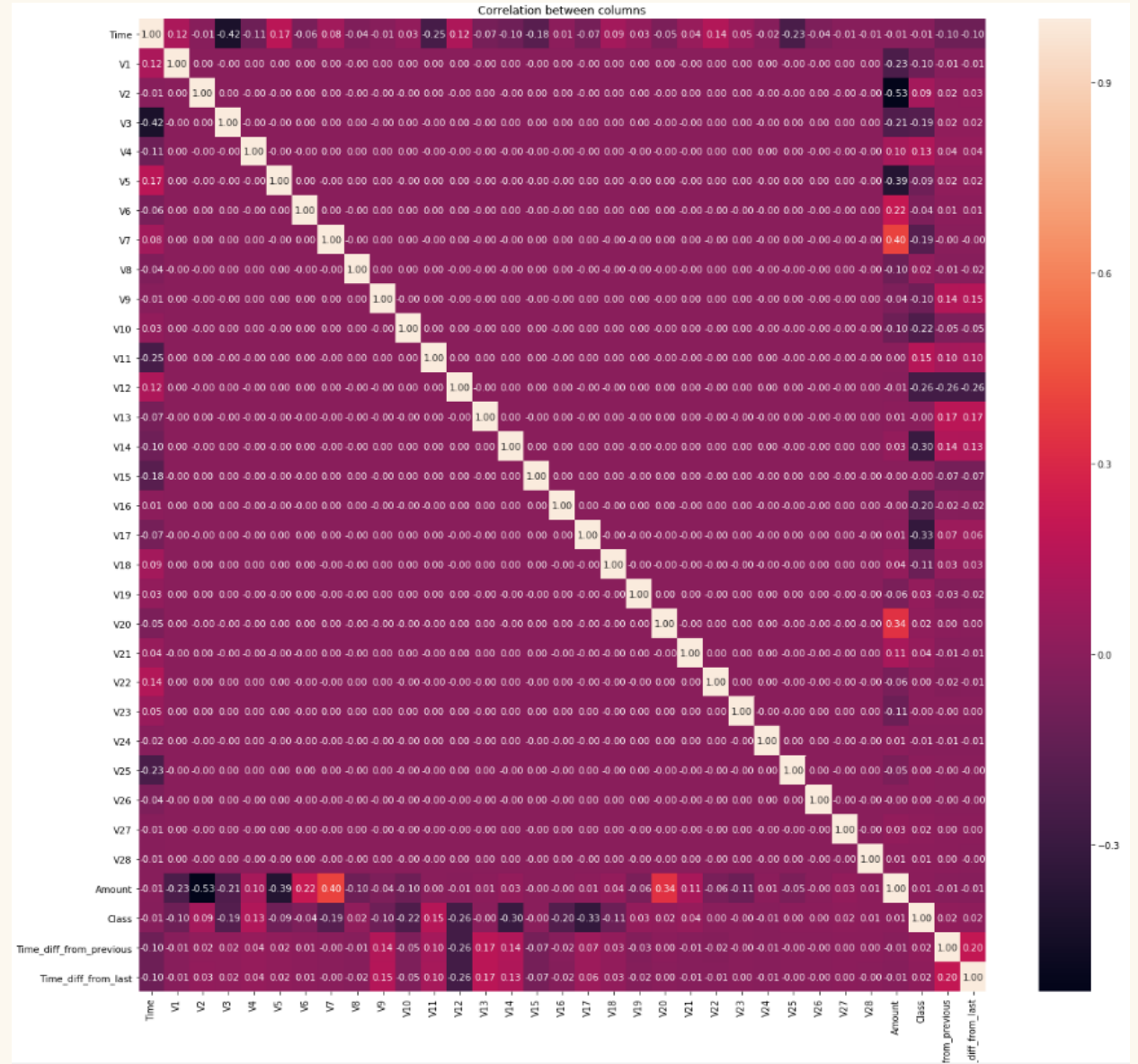# 1) Description of Data



Plot of Time and V1

# 1) Description of Data



Plot of V12 and V17

# 2) Linear Relations

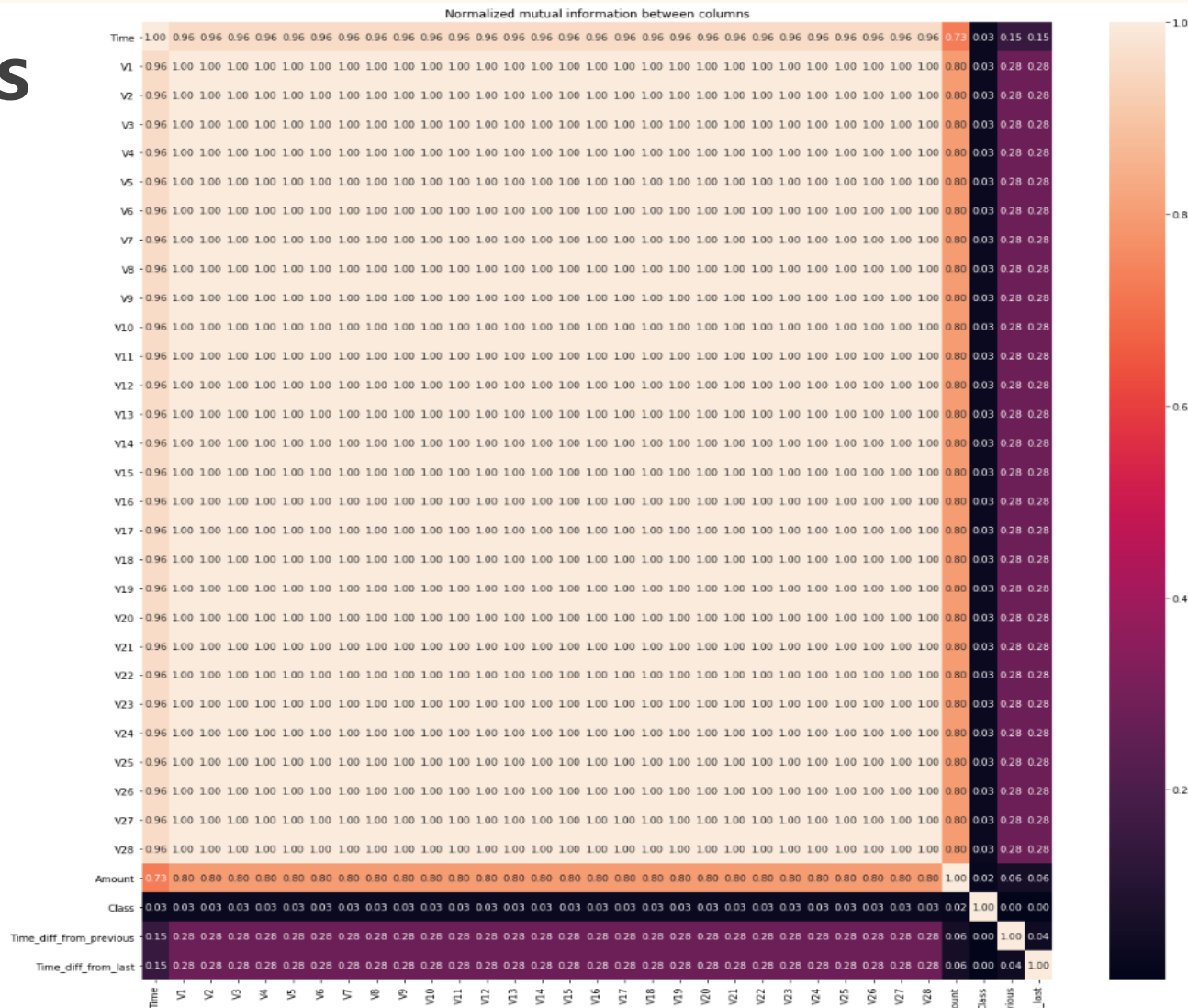## Correlation Matrix

- V's : Principal Components

⇒ Linearly independent

# 2) Nonlinear Relations

## Mutual Information Matrix

- All V's are nonlinearly dependent.
- All Vs and amount are nonlinearly dependent.



Normalized mutual information between columns