

robROSE

An approach for dealing
with imbalanced data in
fraud detection

Sebastiaan Höppner

joint work with
Irène Ortner, Bart Baesens & Tim Verdonck

11 July 2019

The logo for KU Leuven, consisting of the text "KU LEUVEN" in white, bold, uppercase letters on a dark blue rectangular background.

KU LEUVEN

The logo for the Benford's Law Conference, featuring a stylized bar chart with blue bars and an orange line graph overlaid on it.

Benford's Law
Conference

10-12 July 2019 - Stresa, Italy

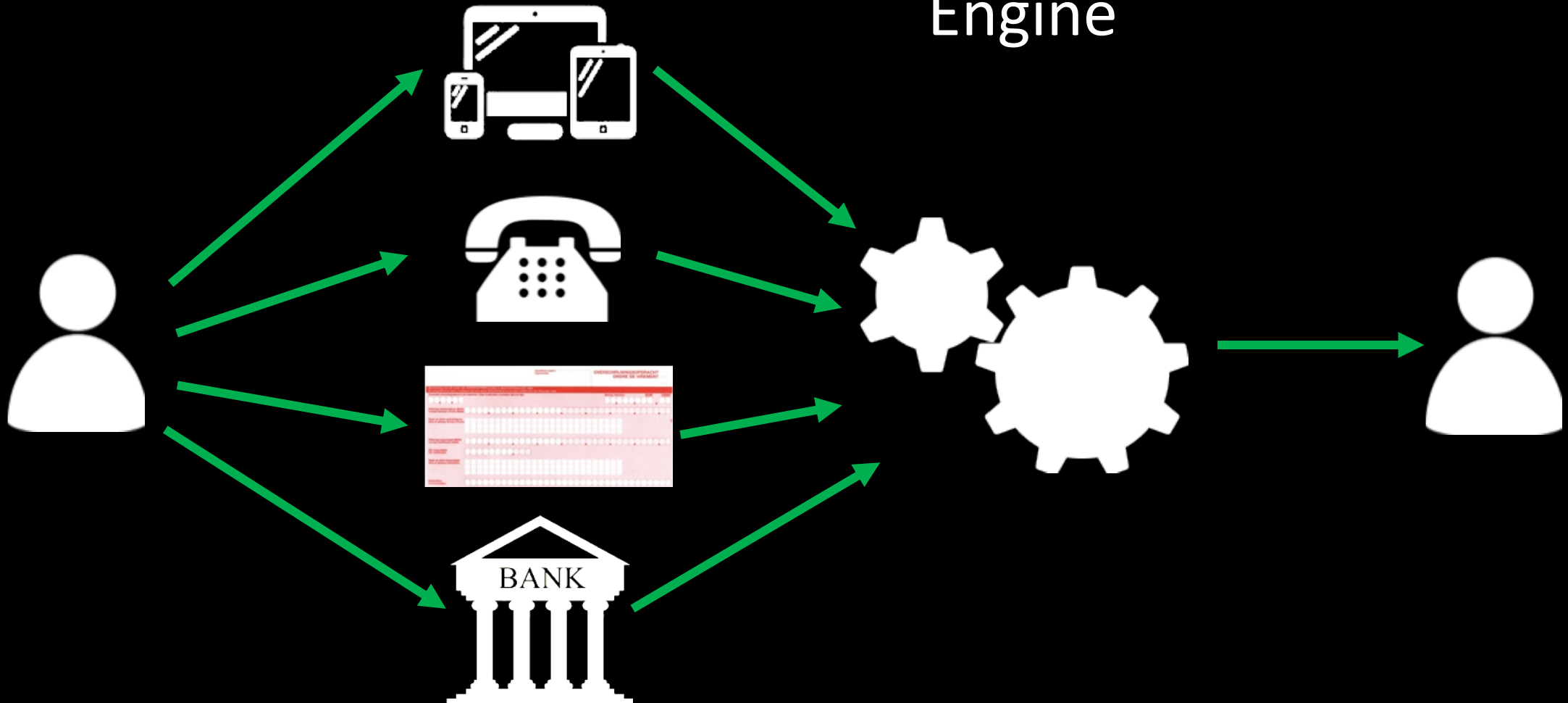
Credit transfers

Initiator

Channels

Payment
Engine

Beneficiary

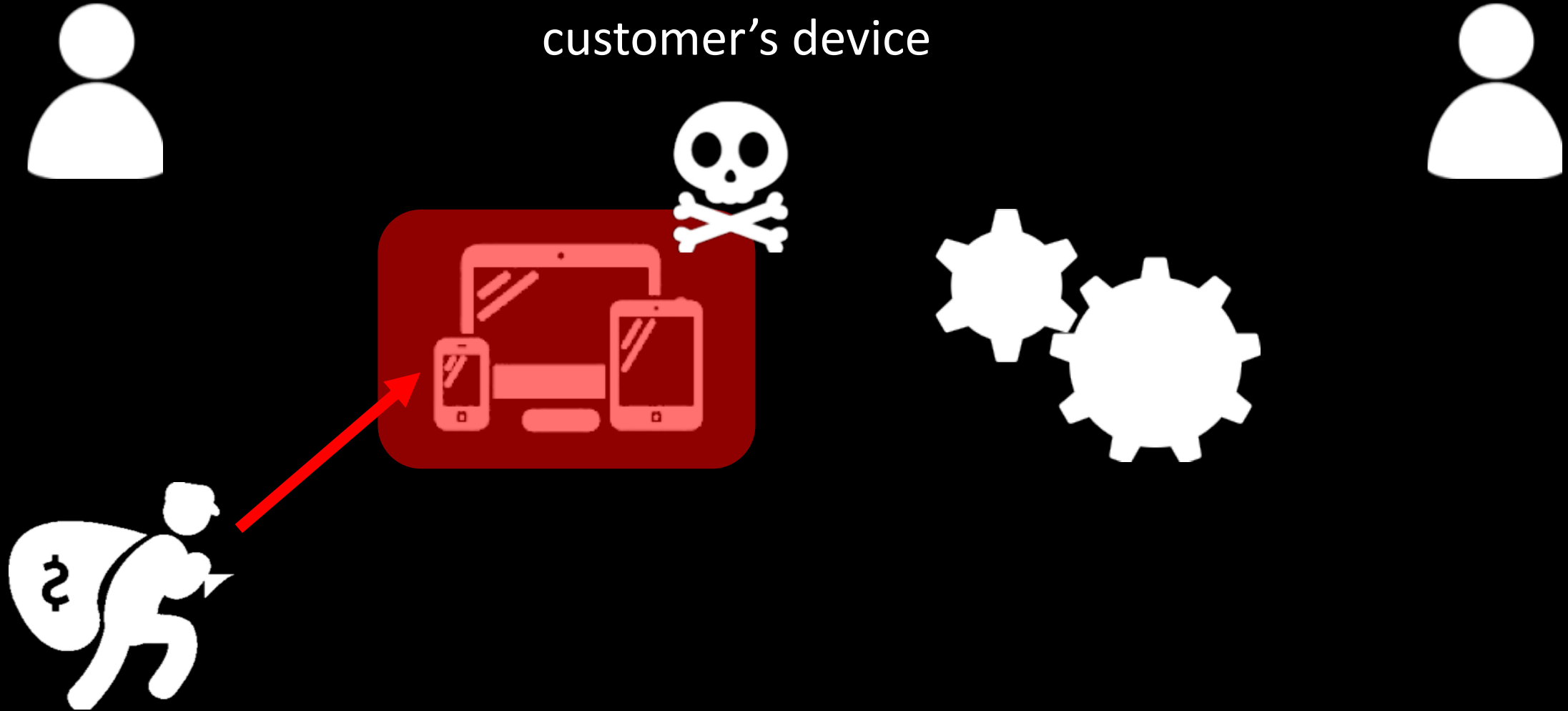


Hacking



Hacking

Step 1: the fraudster
installs malware on the
customer's device



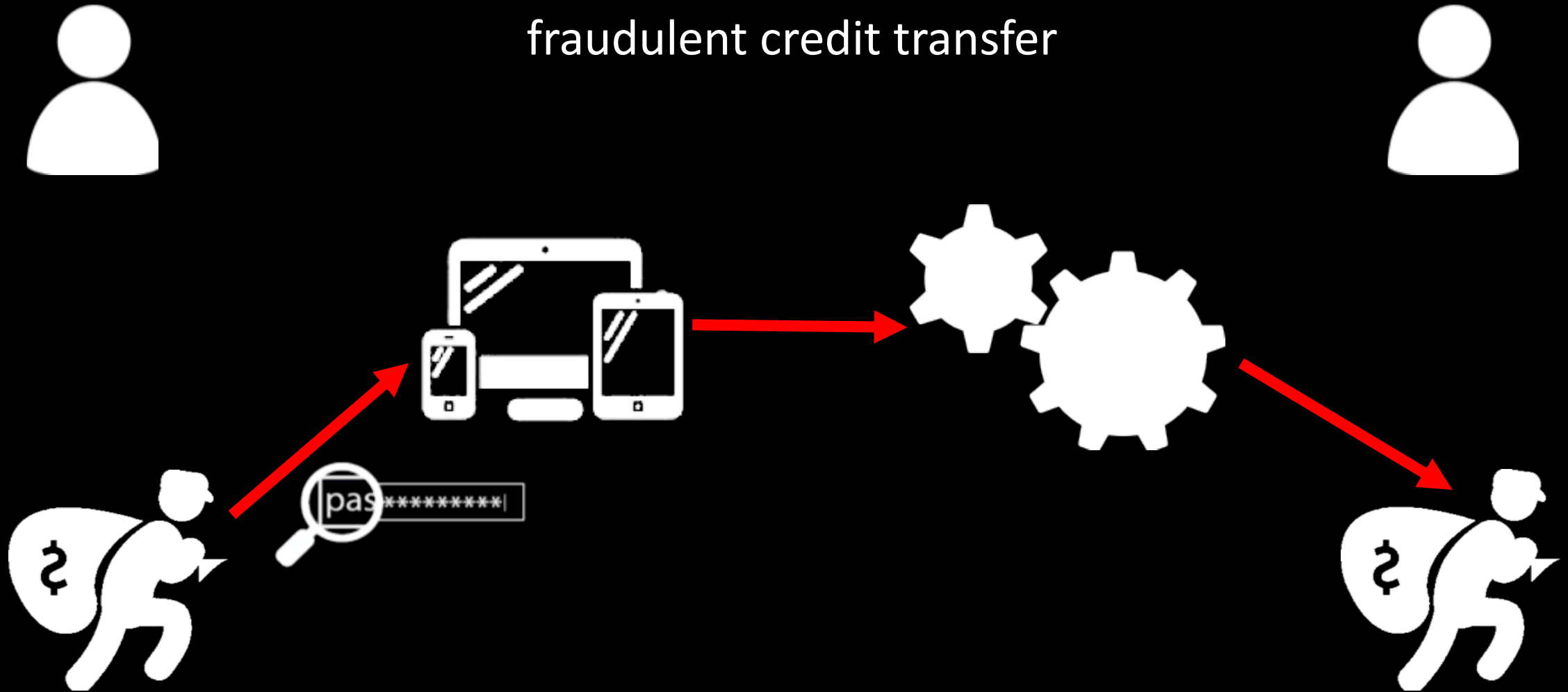
Hacking

Step 2: when the customer uses their device, the fraudster steals their credentials

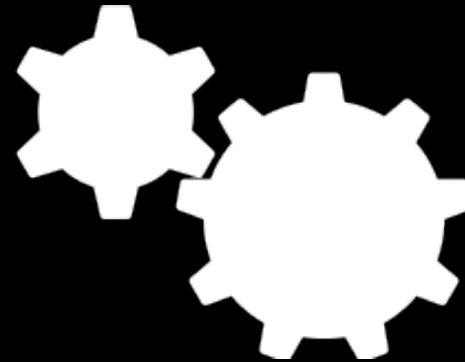
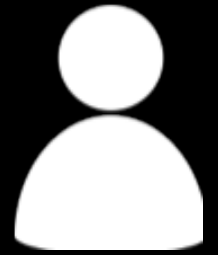


Hacking

Step 3: the fraudster uses the stolen credentials to book a fraudulent credit transfer



Phishing / vishing



Phishing / vishing

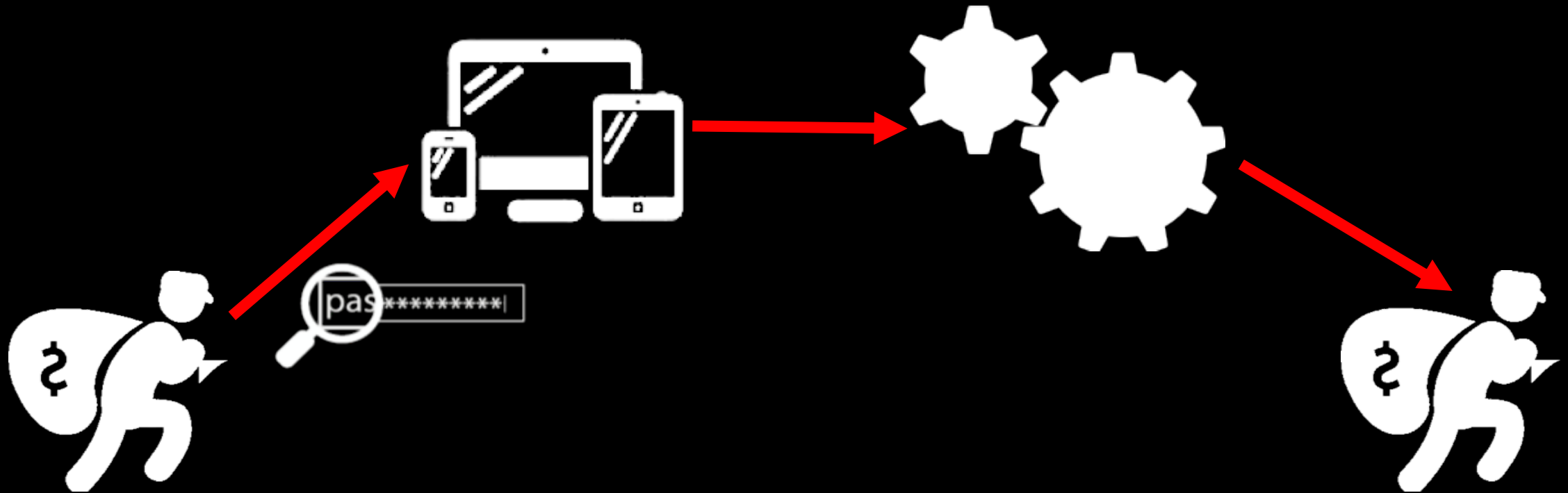
Step 1: a fraudster tricks a customer into sharing their credentials



Phishing / vishing



Step 2: the fraudster uses the stolen credentials to book a fraudulent credit transfer



CEO fraud



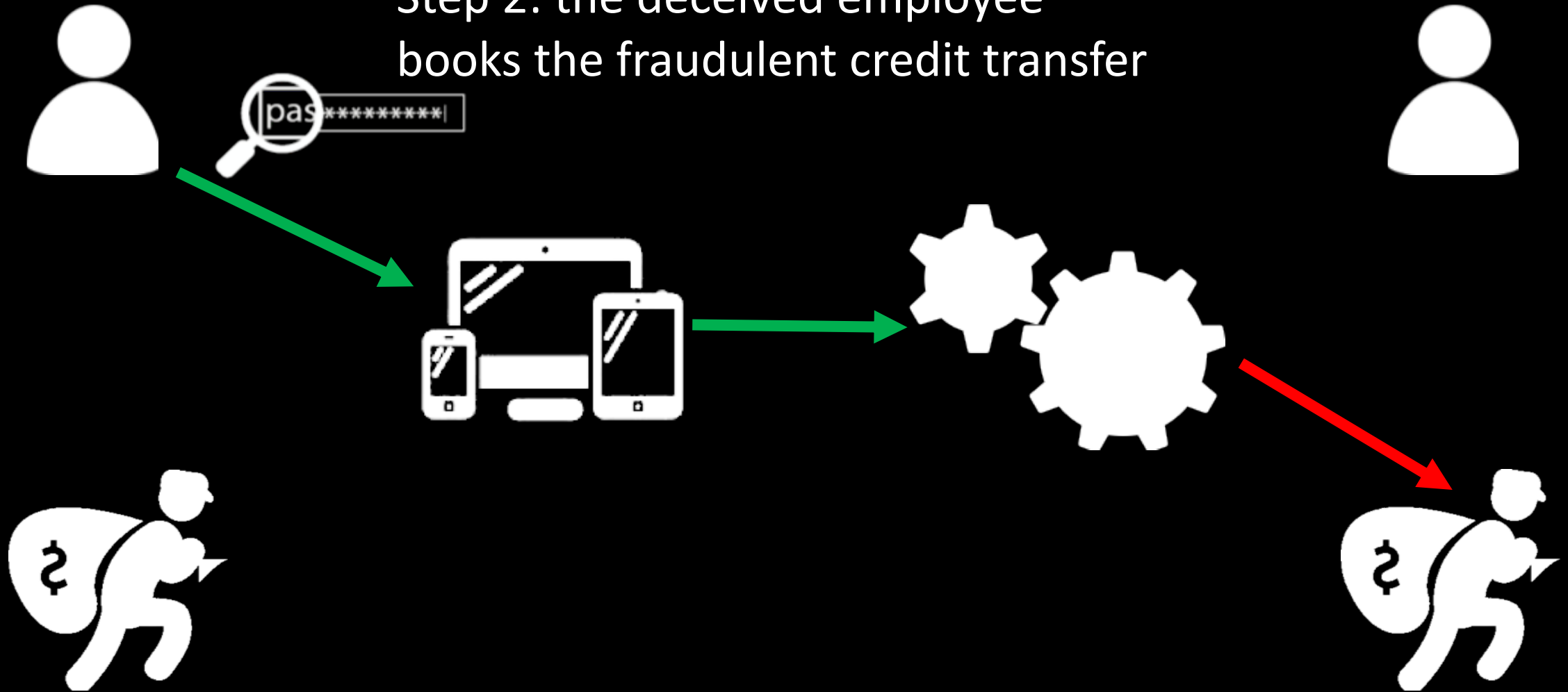
CEO fraud

Step 1: the fraudster impersonates the CEO and convinces an employee to book a credit transfer



CEO fraud

Step 2: the deceived employee books the fraudulent credit transfer



Problem: imbalanced data

- **Binary classification**
legitimate vs fraud



Problem: imbalanced data

- **Binary classification**
legitimate vs fraud
- **Imbalanced data**
(very) large difference in number of observations of both groups



Problem: imbalanced data

- **Binary classification**
legitimate vs fraud
- **Imbalanced data**
(very) large difference in number of observations of both groups
- **Credit card fraud**
less than 1 out 10m transactions
($< 0.00001\%$)

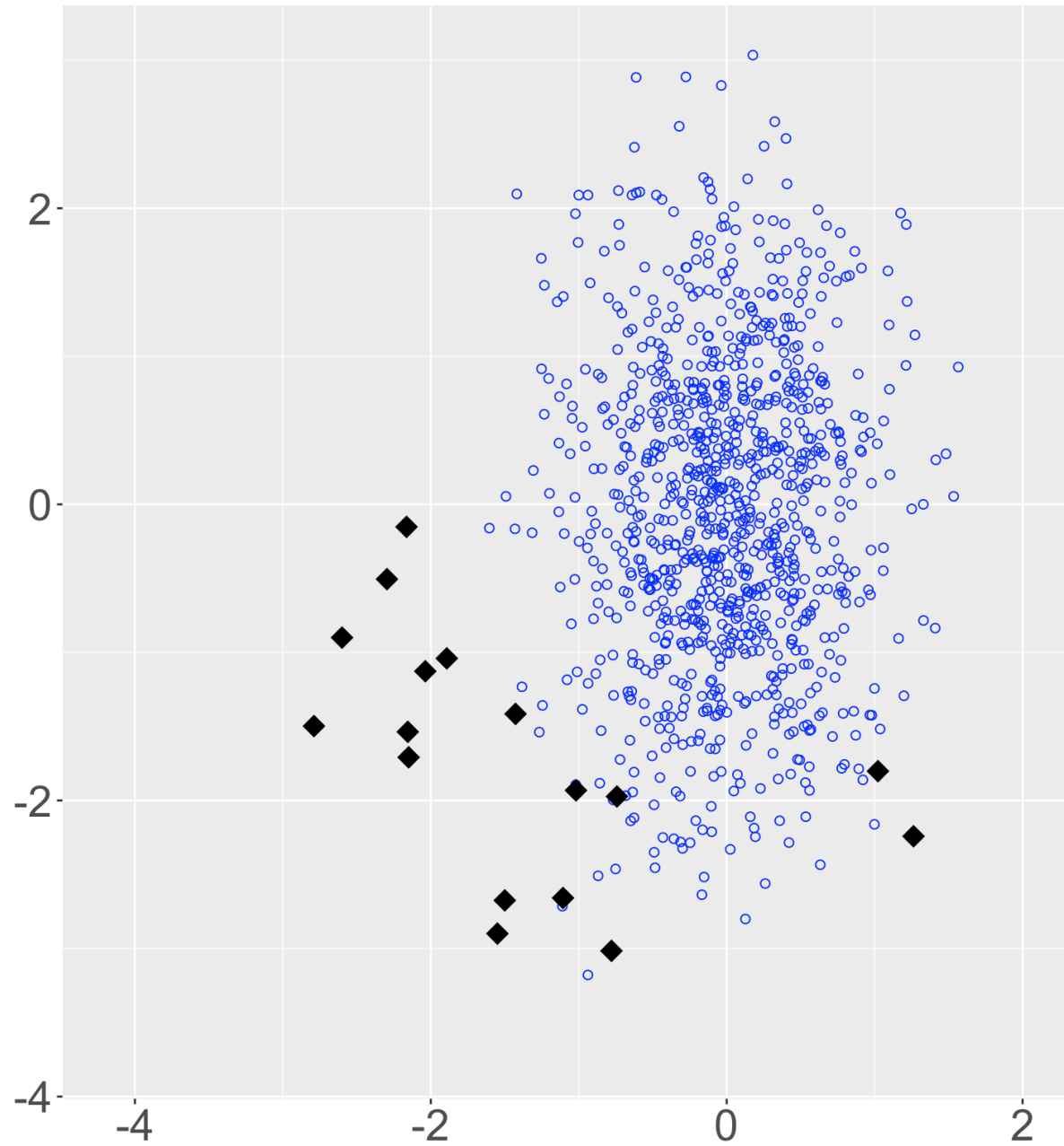


Problem: imbalanced data

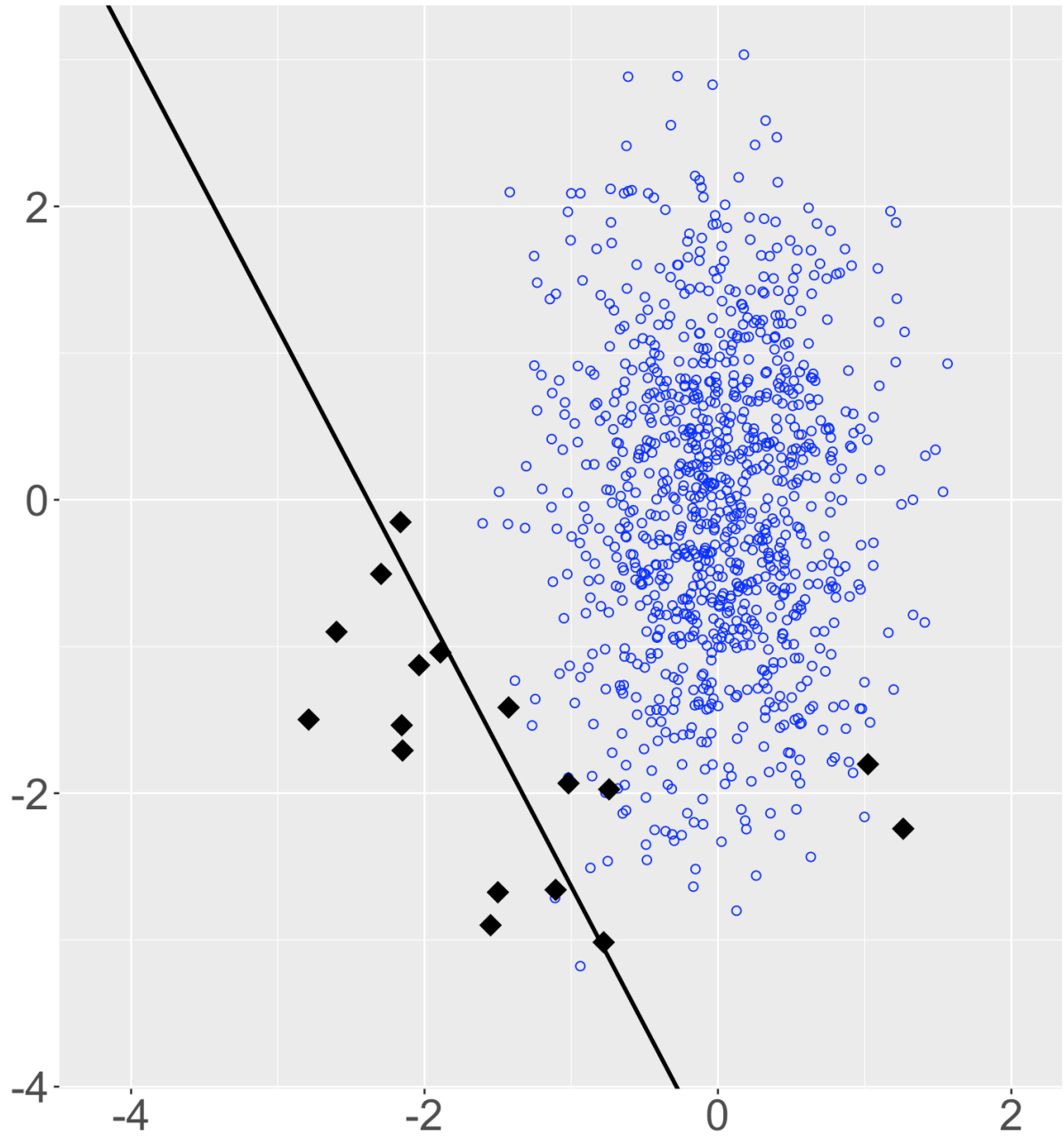
- **Binary classification**
legitimate vs fraud
- **Imbalanced data**
(very) large difference in number of observations of both groups
- **Credit card fraud**
less than 1 out 10m transactions
($< 0.00001\%$)
- Typically very few “cases of interest”
compared to legitimate observations
(20% - 0.01%)



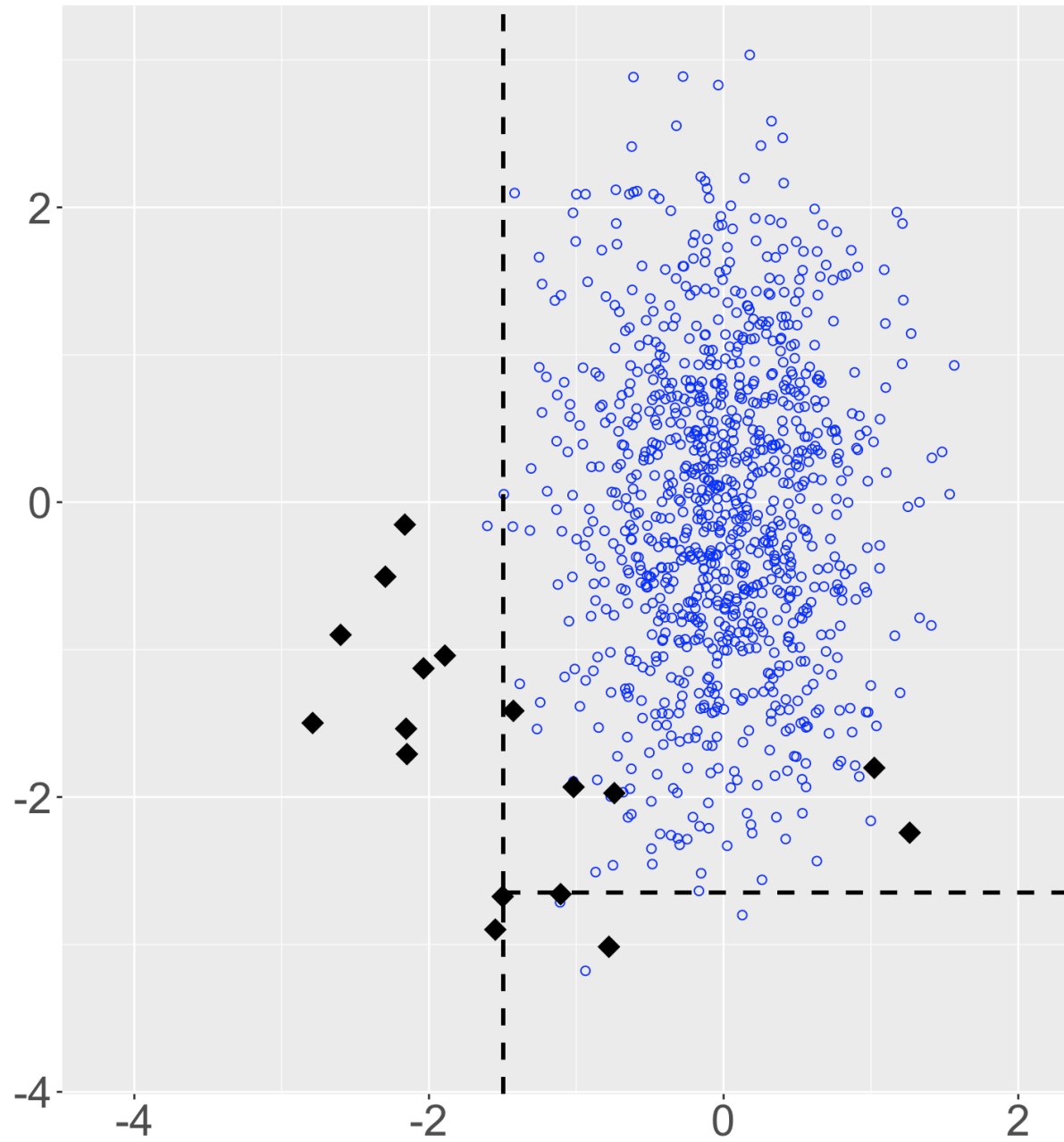
Original data



Logistic regression



Classification tree (CART)



“How to reduce imbalance in training data?”

“How to reduce imbalance in training data?”

Reduce # legitimate cases

Increase # fraud cases

“How to reduce imbalance in training data?”

Reduce # legitimate cases

- Random under-sampling:
randomly sub-sample
legitimate cases

Increase # fraud cases

- Random over-sampling:
sampling with replacement
of fraud samples
- ***Generate synthetic
minority/fraud cases***

“How to reduce imbalance in training data?”

Reduce # legitimate cases

- Random under-sampling:
randomly sub-sample
legitimate cases

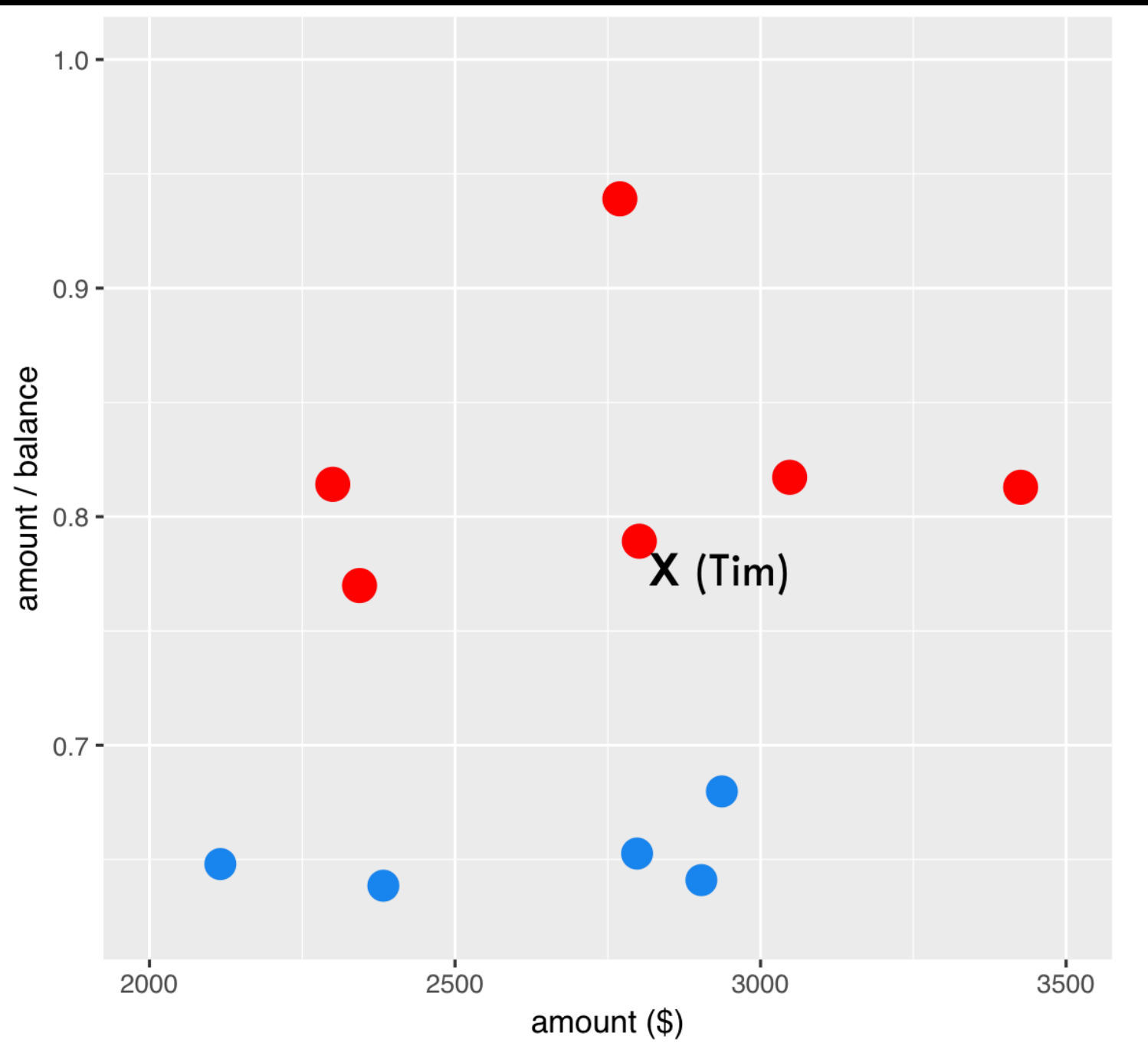
Increase # fraud cases

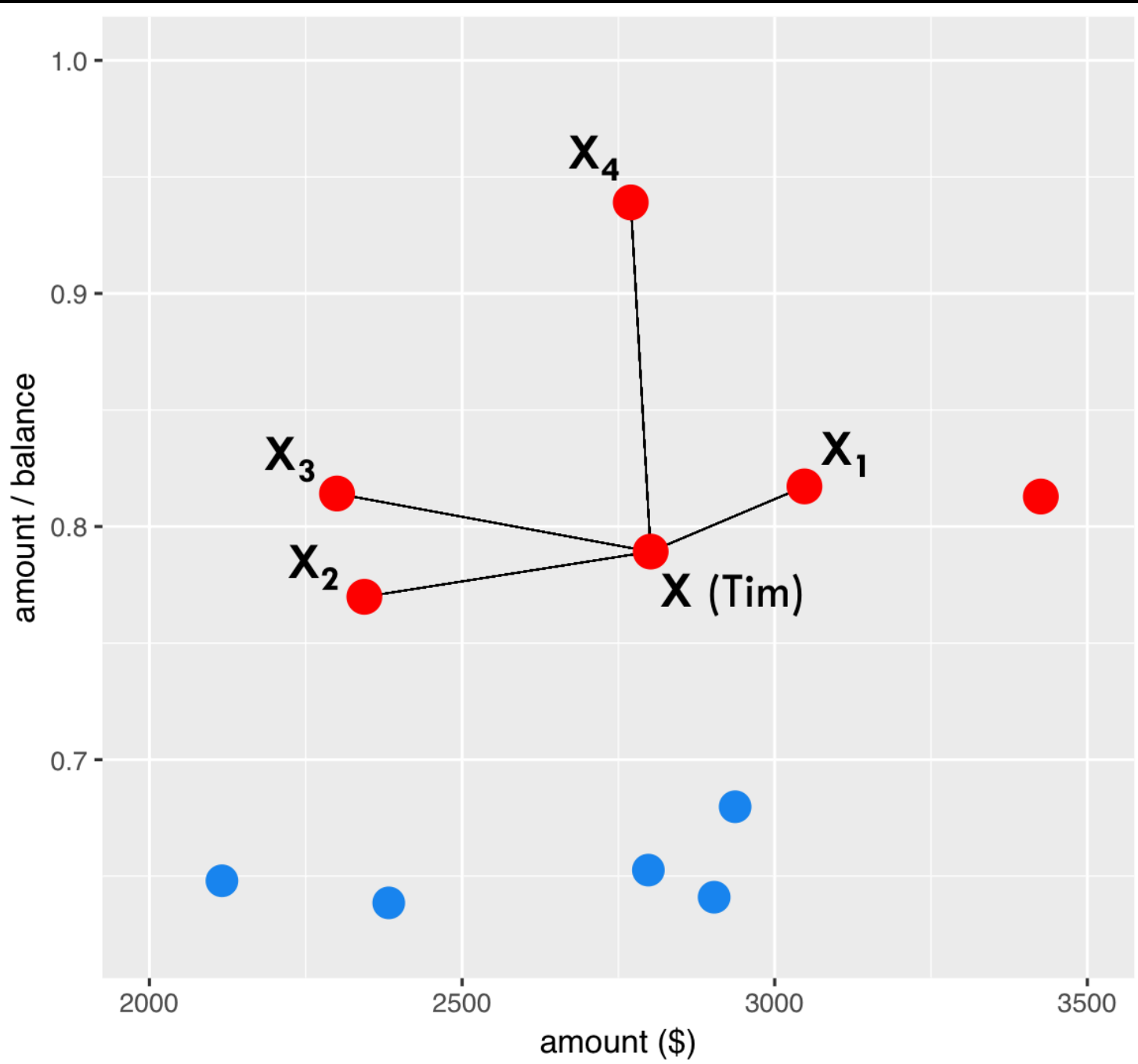
- Random over-sampling:
sampling with replacement
of fraud samples
- ***Generate synthetic
minority/fraud cases***

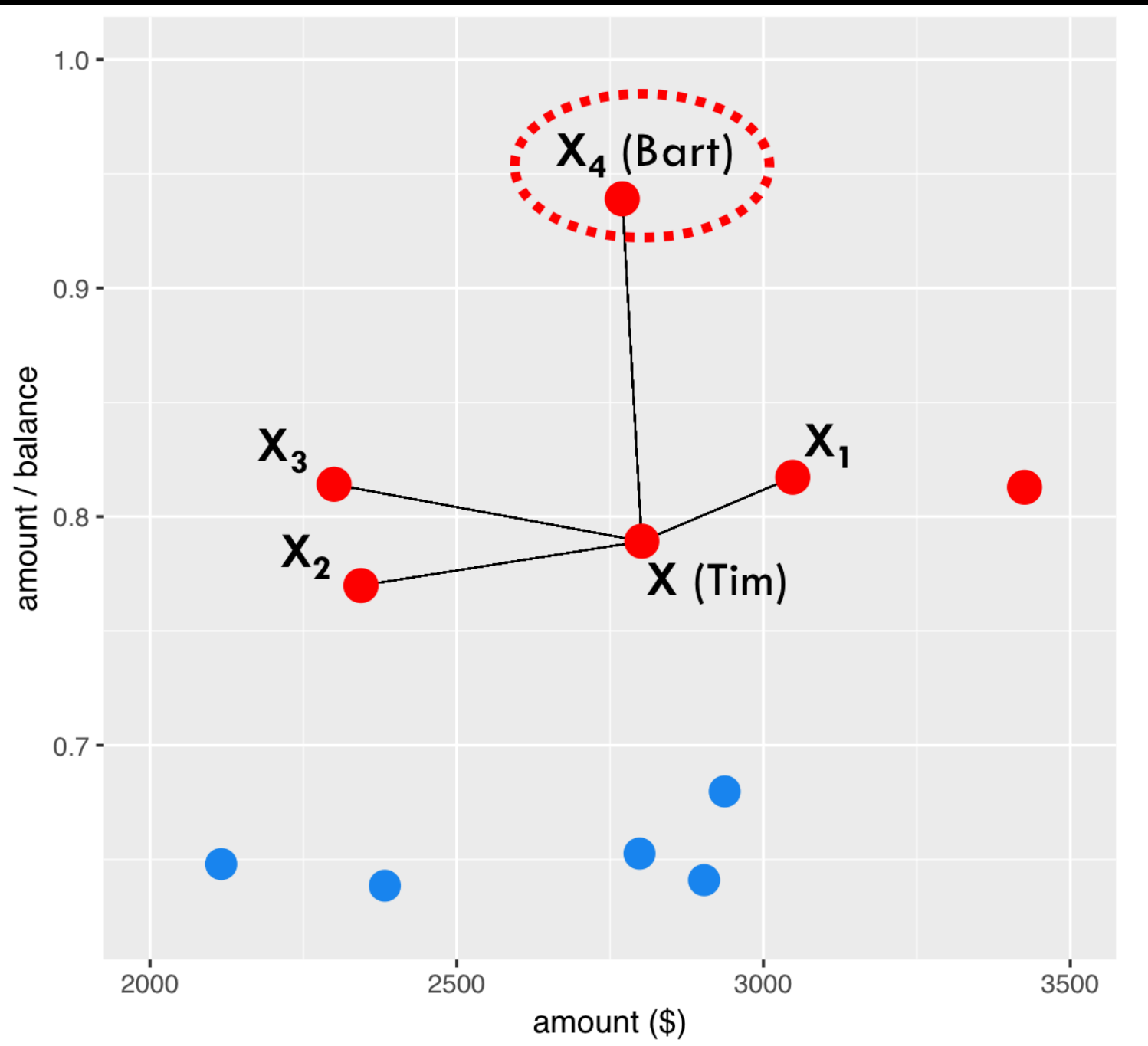
Evaluate model on imbalanced “original” test data !!!

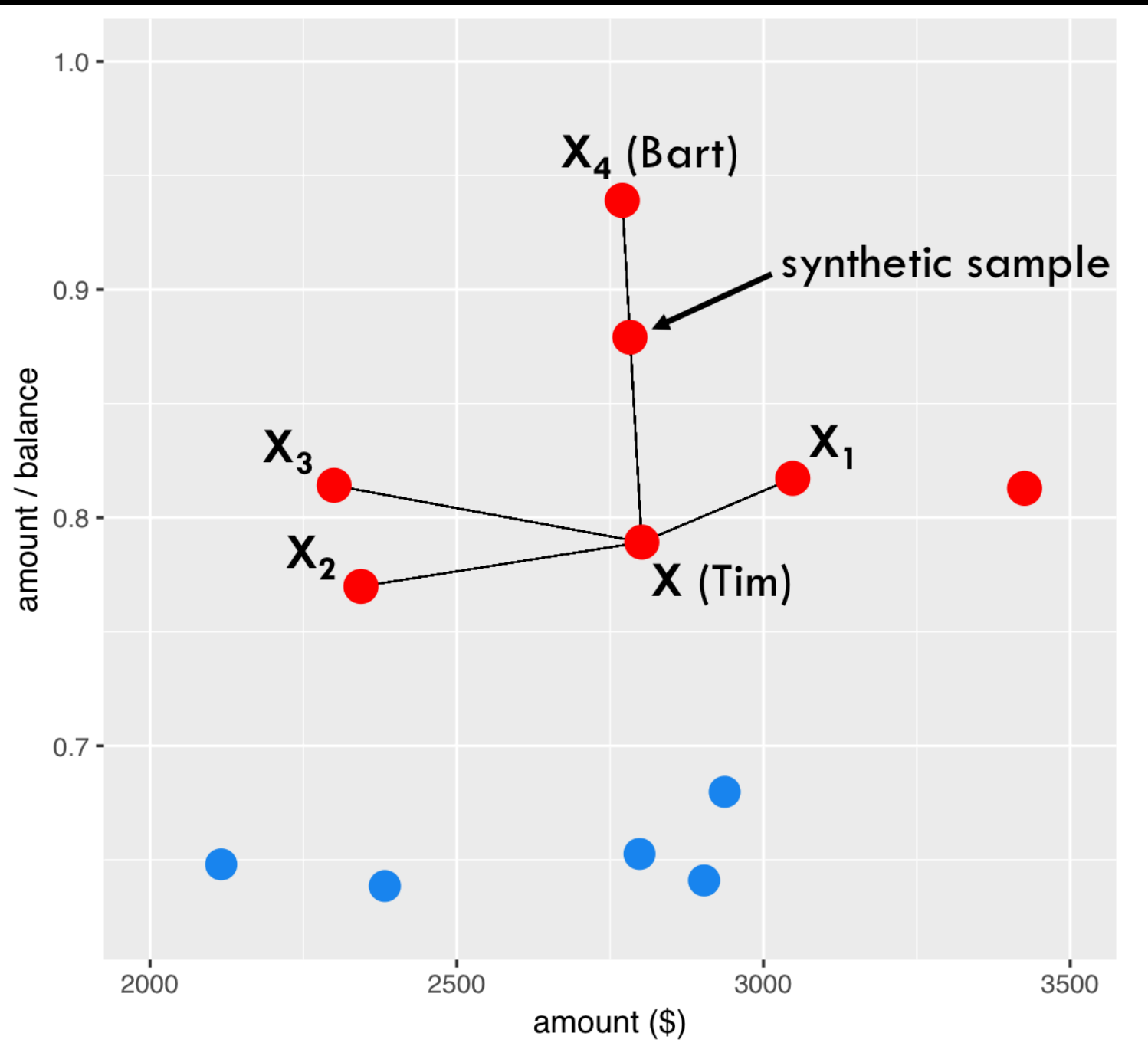
SMOTE - Synthetic Minority Over-sampling Technique

Chawla, Bowyer, Hall & Kegelmeyer (2002)

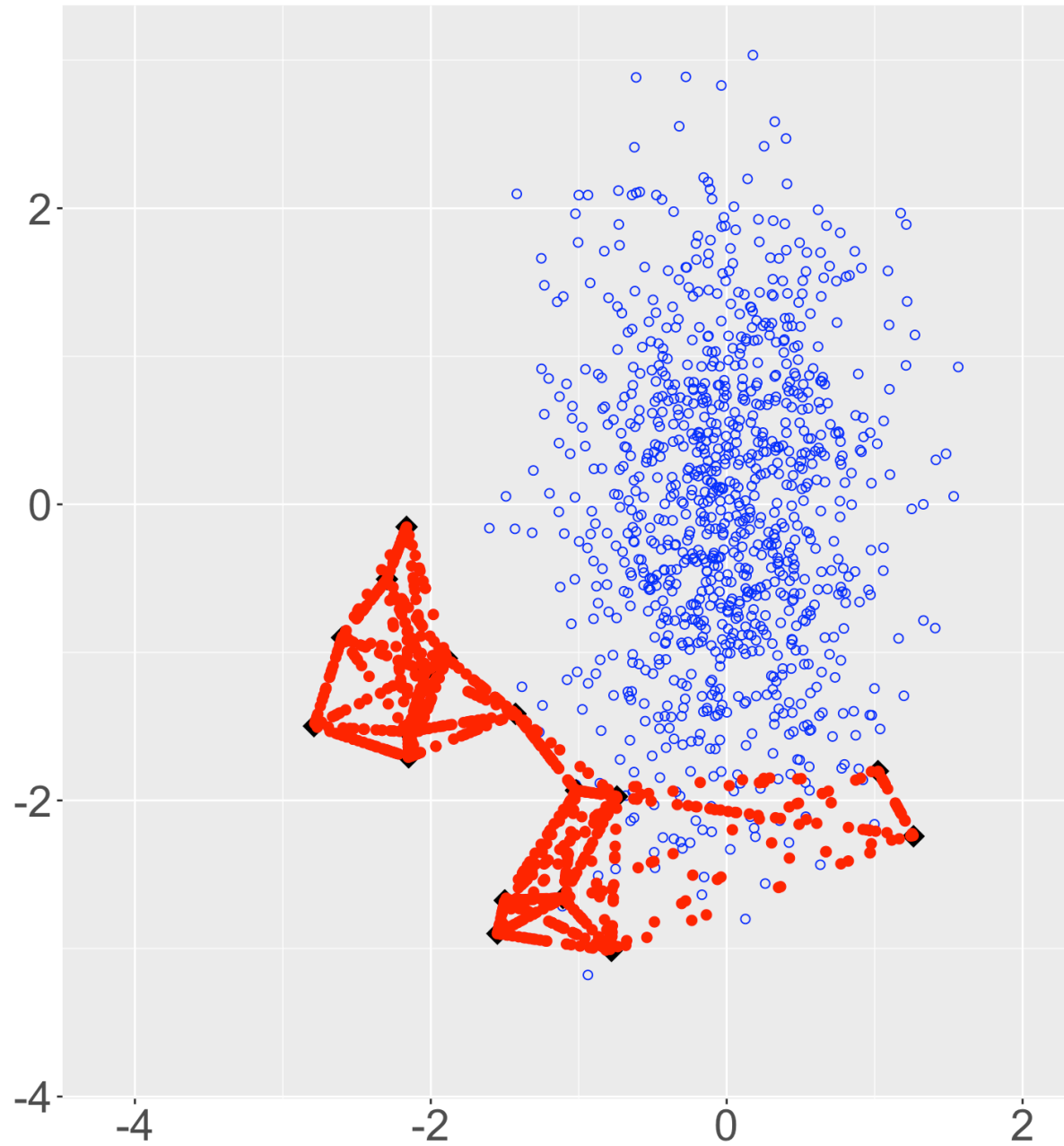








SMOTE



ROSE - Random Over-Sampling Examples

ROSE - Random Over-Sampling Examples

1. Randomly select an observation x_i from the minority group

ROSE - Random Over-Sampling Examples

1. Randomly select an observation \mathbf{x}_i from the minority group
2. Estimate normal density distribution $\mathcal{N}(\mathbf{x}_i, H)$

ROSE - Random Over-Sampling Examples

1. Randomly select an observation \mathbf{x}_i from the minority group
2. Estimate normal density distribution $\mathcal{N}(\mathbf{x}_i, H)$
 - selected observation \mathbf{x}_i as center

ROSE - Random Over-Sampling Examples

1. Randomly select an observation \mathbf{x}_i from the minority group
2. Estimate normal density distribution $\mathcal{N}(\mathbf{x}_i, H)$
 - selected observation \mathbf{x}_i as center
 - *smoothing* matrix $H = \text{diag}(h_1, \dots, h_d)$

$$h_q = \left(\frac{4}{(d+2)n} \right)^{1/(d+4)} \hat{\sigma}_q \quad (q = 1, \dots, d)$$

$\hat{\sigma}_q$ = sample standard deviation of q-th variable of minority group

ROSE - Random Over-Sampling Examples

1. Randomly select an observation \mathbf{x}_i from the minority group
2. Estimate normal density distribution $\mathcal{N}(\mathbf{x}_i, H)$
 - selected observation \mathbf{x}_i as center
 - *smoothing* matrix $H = \text{diag}(h_1, \dots, h_d)$

$$h_q = \left(\frac{4}{(d+2)n} \right)^{1/(d+4)} \hat{\sigma}_q \quad (q = 1, \dots, d)$$

$\hat{\sigma}_q$ = sample standard deviation of q-th variable of minority group

ROSE - Random Over-Sampling Examples

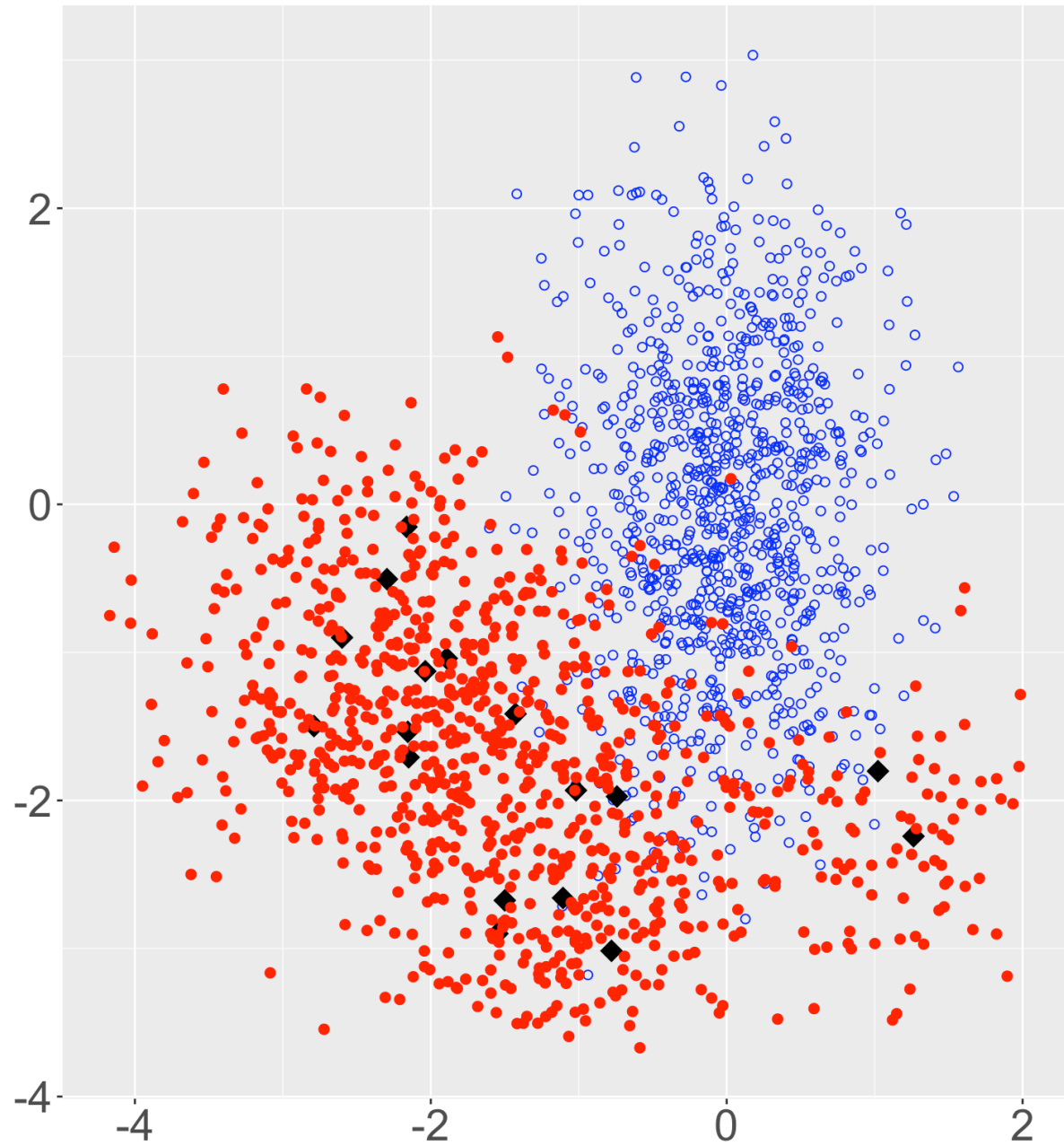
1. Randomly select an observation \mathbf{x}_i from the minority group
2. Estimate normal density distribution $\mathcal{N}(\mathbf{x}_i, H)$
 - selected observation \mathbf{x}_i as center
 - *smoothing* matrix $H = \text{diag}(h_1, \dots, h_d)$

$$h_q = \left(\frac{4}{(d+2)n} \right)^{1/(d+4)} \hat{\sigma}_q \quad (q = 1, \dots, d)$$

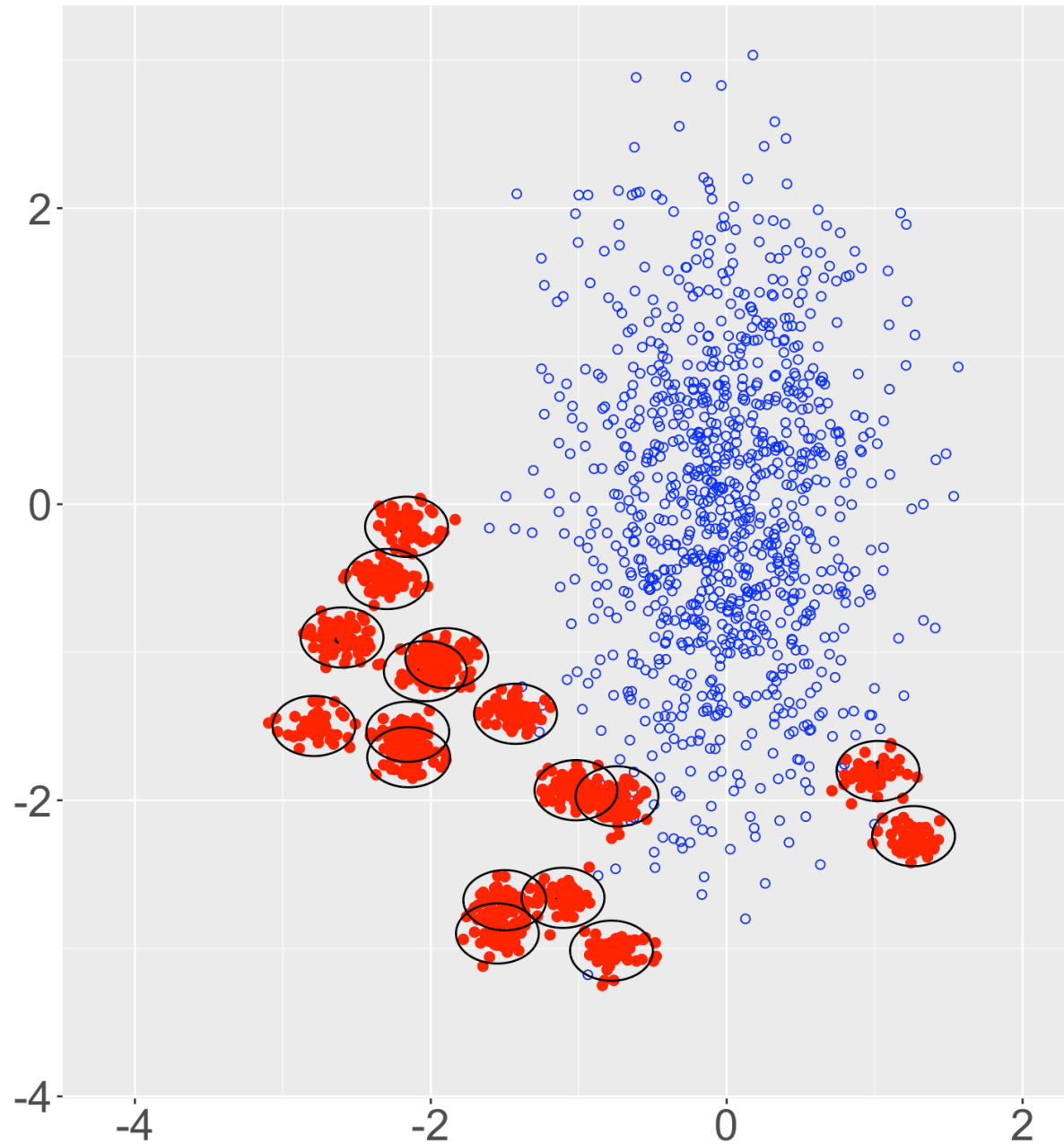
$\hat{\sigma}_q$ = sample standard deviation of q-th variable of minority group

3. Generate a new observation from this normal density estimate

ROSE (h.minor = 1) (default)



ROSE (h.minor = 0.15)



Robust ROSE

- Identify “outlying” minority cases based on Mahalanobis distance (MD) using the robust MCD¹ estimator

¹ Rousseeuw and Van Driessen (1999)

Robust ROSE

- Identify “outlying” minority cases based on Mahalanobis distance (MD) using the robust MCD¹ estimator
- Synthetic cases generated only for minority cases x_i with $MD^2(x_i) < \chi_d^2(1 - \alpha)$, e.g. $\alpha = 1\% \rightarrow$ “non-outlying” minority cases

¹ Rousseeuw and Van Driessen (1999)

Robust ROSE

1. Select x_i from the “non-outlying” minority group with probability $\propto 1 /$ (shortest distance from x_i to majority case)

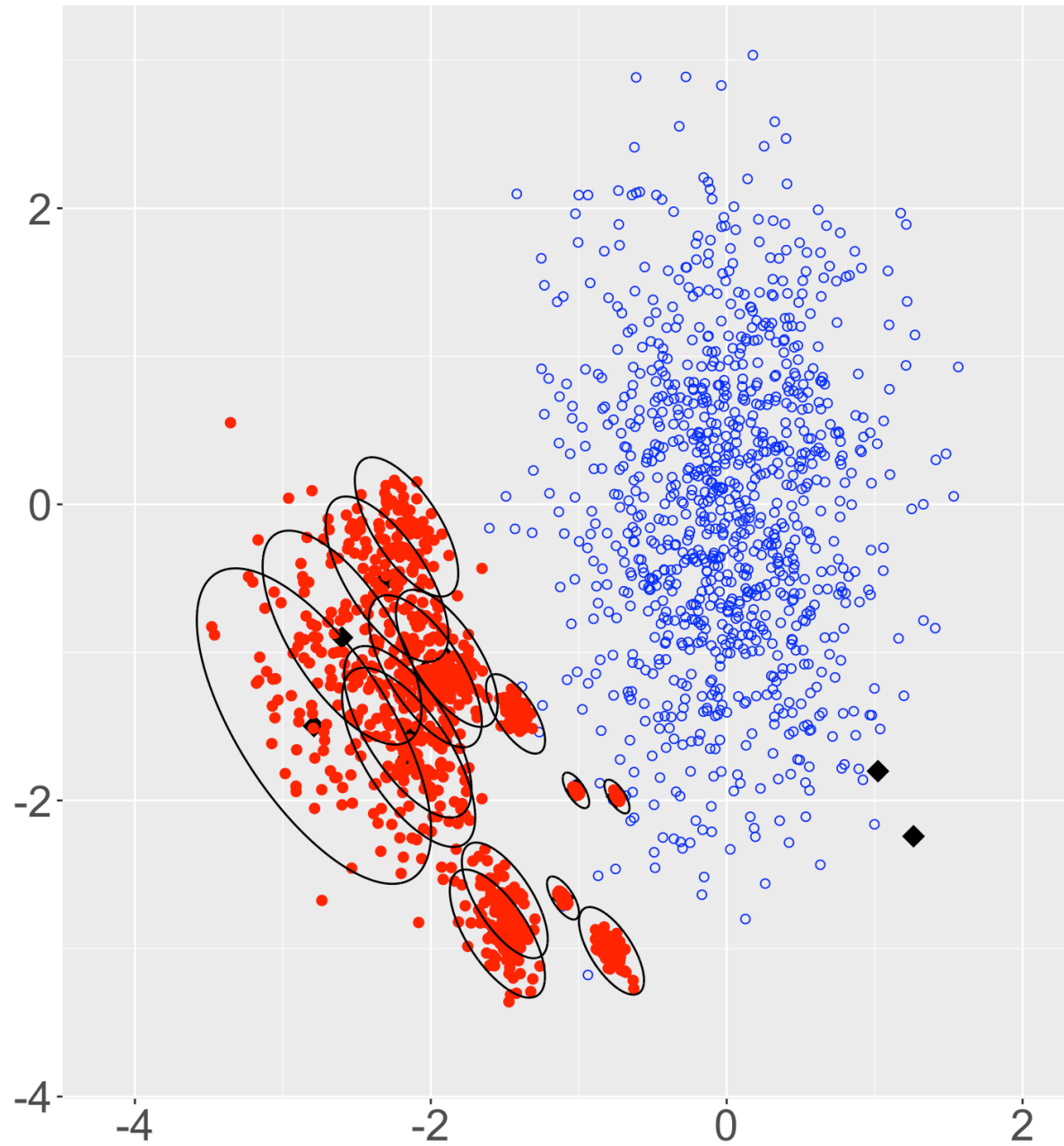
Robust ROSE

1. Select \mathbf{x}_i from the “non-outlying” minority group with probability $\propto 1 /$ (shortest distance from \mathbf{x}_i to majority case)
2. Estimate normal density distribution $\mathcal{N}(\mathbf{x}_i, H)$
 - selected observation \mathbf{x}_i as center
 - *smoothing* matrix $H =$ MCD covariance matrix estimate on “non-outlying” minority group

Robust ROSE

1. Select \mathbf{x}_i from the “non-outlying” minority group with probability $\propto 1 /$ (shortest distance from \mathbf{x}_i to majority case)
2. Estimate normal density distribution $\mathcal{N}(\mathbf{x}_i, H)$
 - selected observation \mathbf{x}_i as center
 - *smoothing* matrix $H =$ MCD covariance matrix estimate on “non-outlying” minority group
3. Generate a new observation from this multivariate normal density

robROSE

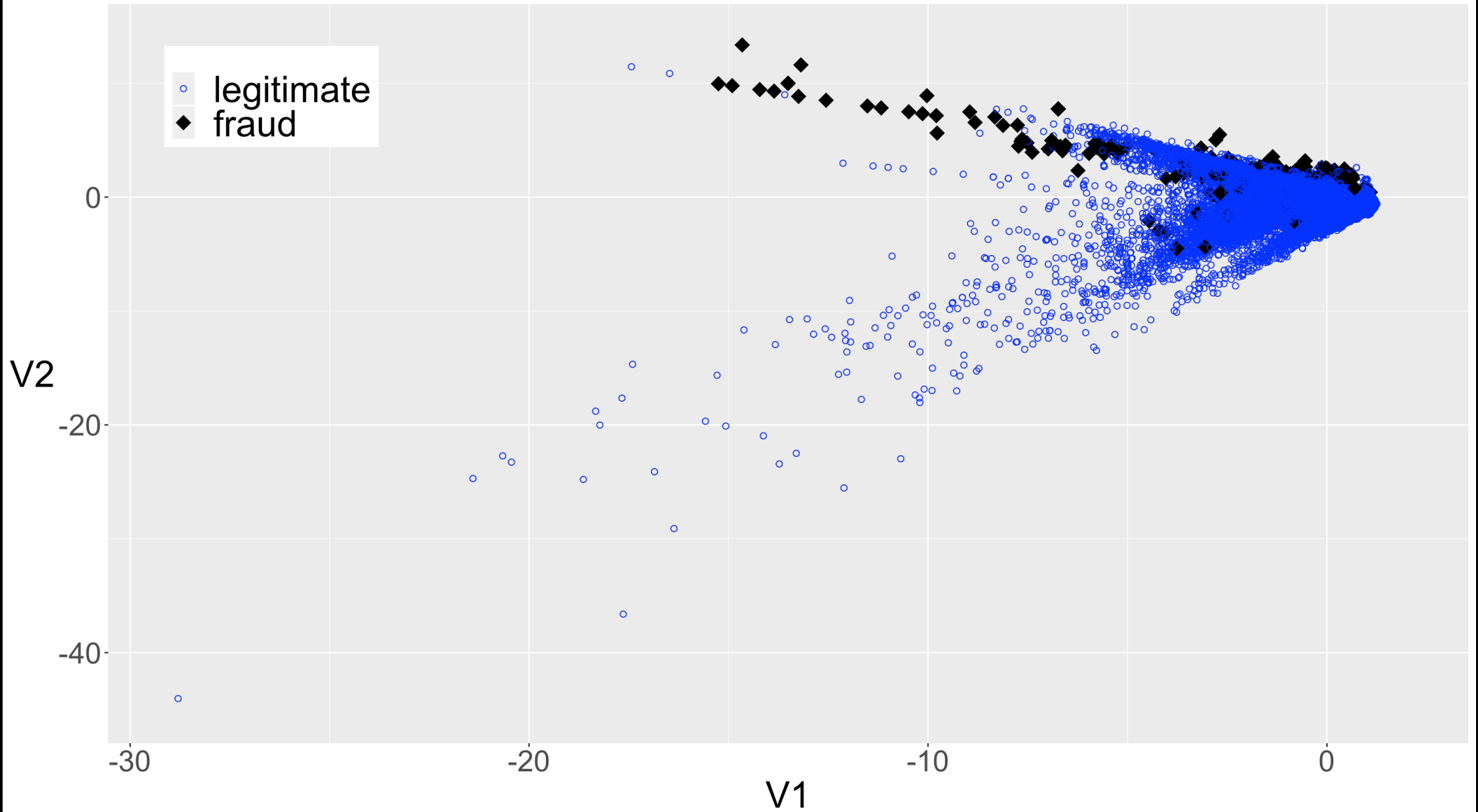


Credit Card Fraud Detection

Anonymized credit card transactions labeled as fraudulent or genuine

source: kaggle.com, made available by Andrea Dal Pozzolo et al., Calibrating Probability with Undersampling for Unbalanced Classification. In Symposium on Computational Intelligence and Data Mining (CIDM), IEEE, 2015

- Transactions made in two days by credit cards in September 2013 by European cardholders
- 497 frauds out of 284,807 transactions \Rightarrow 0.172% is fraud



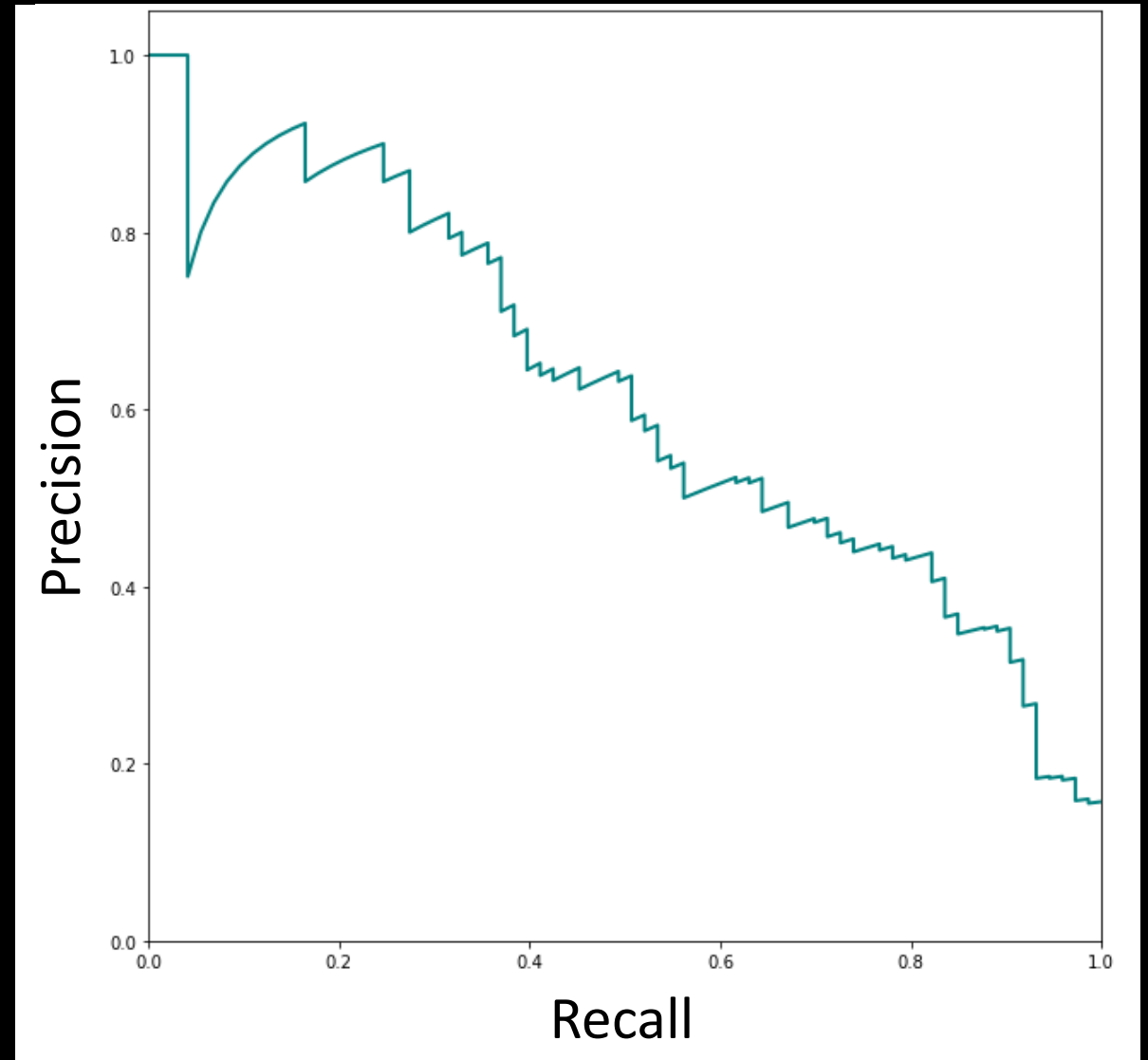
- Precision = $\frac{TP}{TP+FP}$

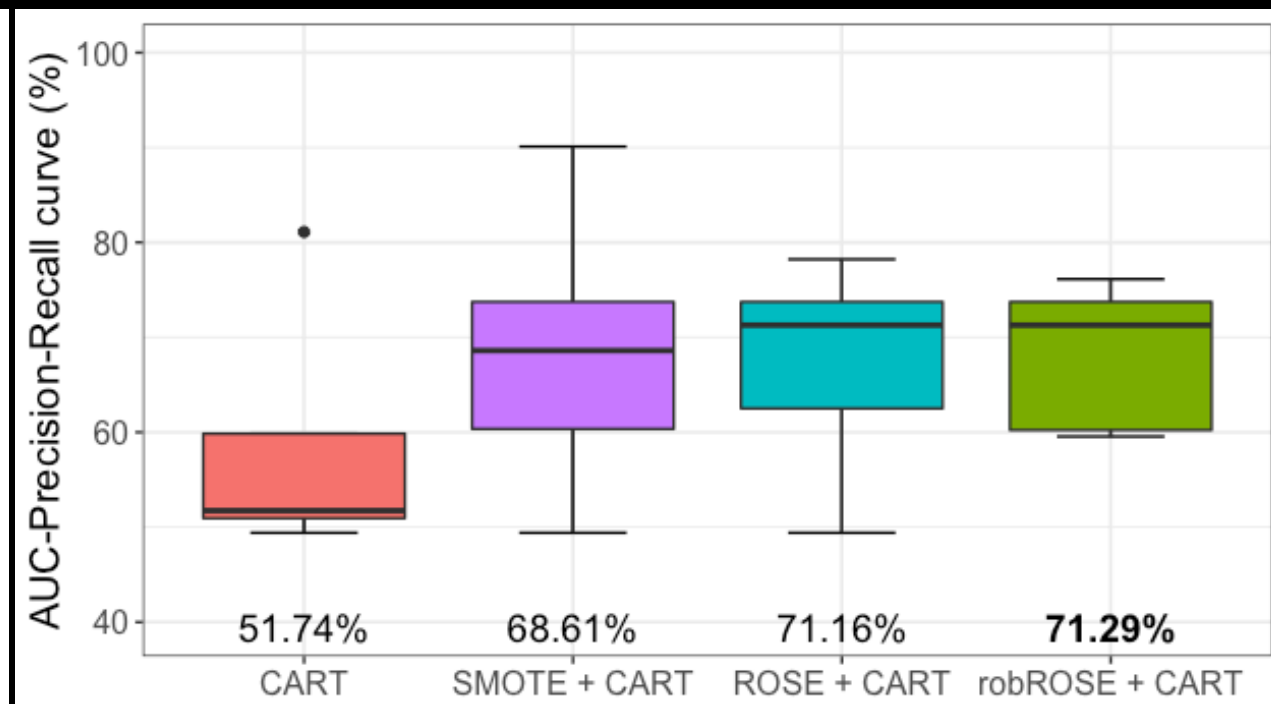
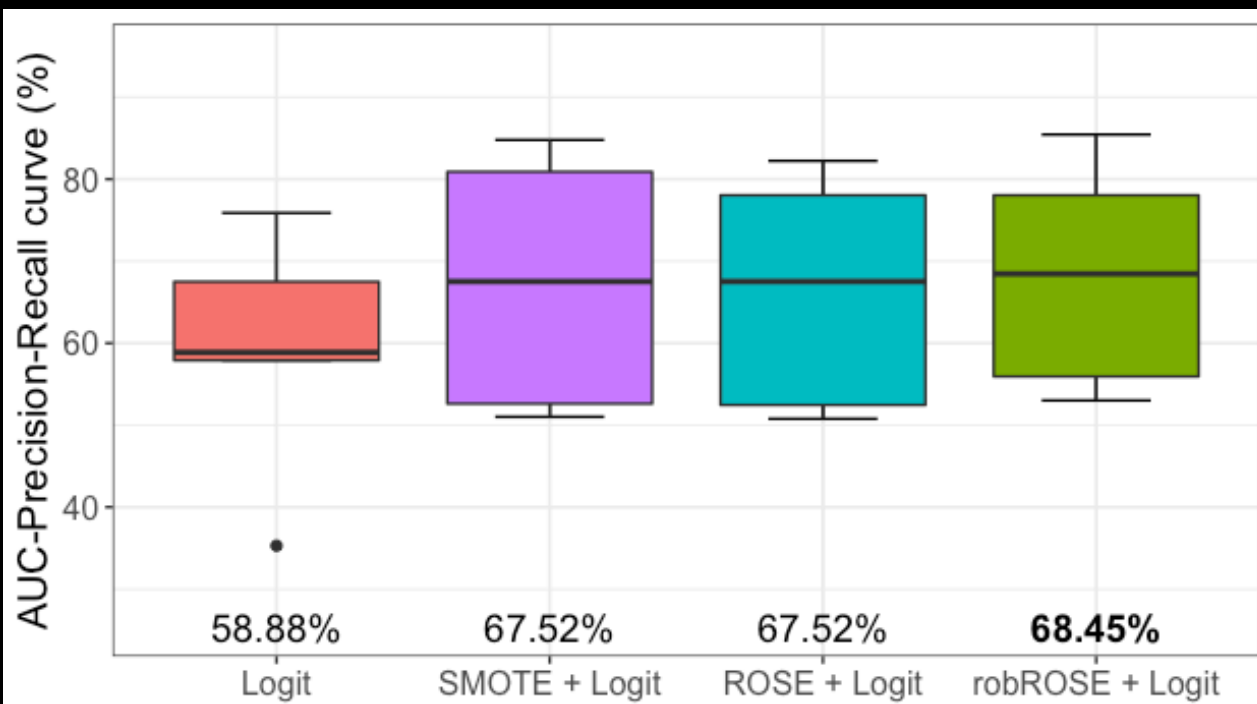
Out of all cases classified as fraud, how many are actually fraud?

- Recall = $\frac{TP}{TP+FN}$

Out of all fraud cases, how many are detected?

- Evaluation measure:
area under precision-recall curve





Thank you

SMOTE

Nitesh V Chawla, Kevin W Bowyer, Lawrence O Hall, and W Philip Kegelmeyer. Smote: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–357, 2002.

ROSE

Giovanna Menardi and Nicola Torelli. Rose: random over-sampling examples. *Data Mining and Knowledge Discovery*, 28(1):92–122, 2014.