

K410 菊池研・斉藤研合同発表会 2017年2月4日

# 乗降履歴データの 有用性評価指標と匿名加工

明治大学 菊池研究室 4年

伊藤聡志

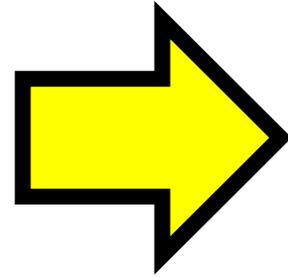
# 研究背景

**危険**

**個人情報**

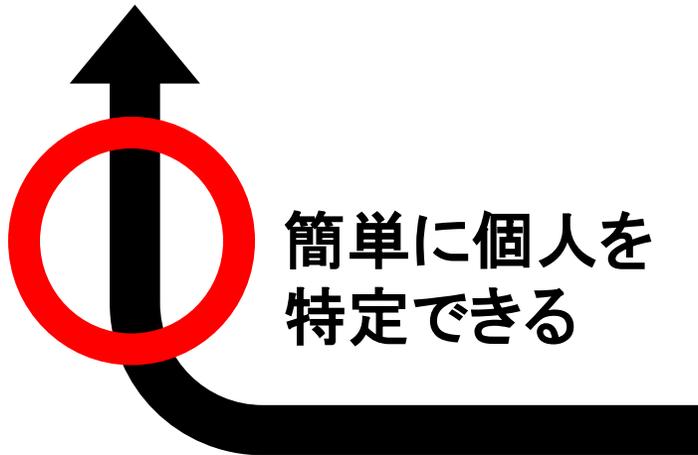
例: 駅の乗降履歴

**匿名加工**



**安全**

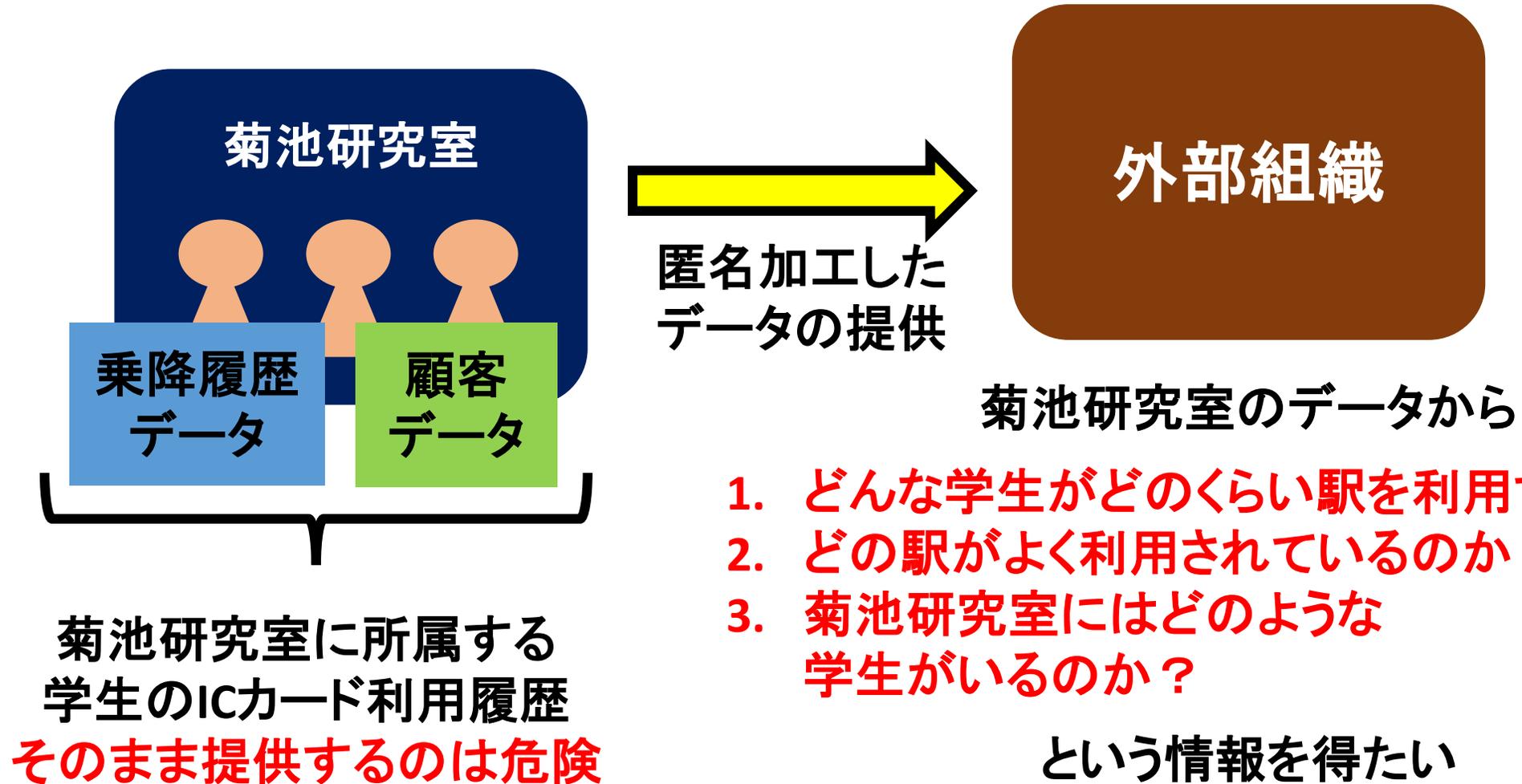
**匿名加工データ**



**攻撃者**



# 想定されるケース



## 問題点

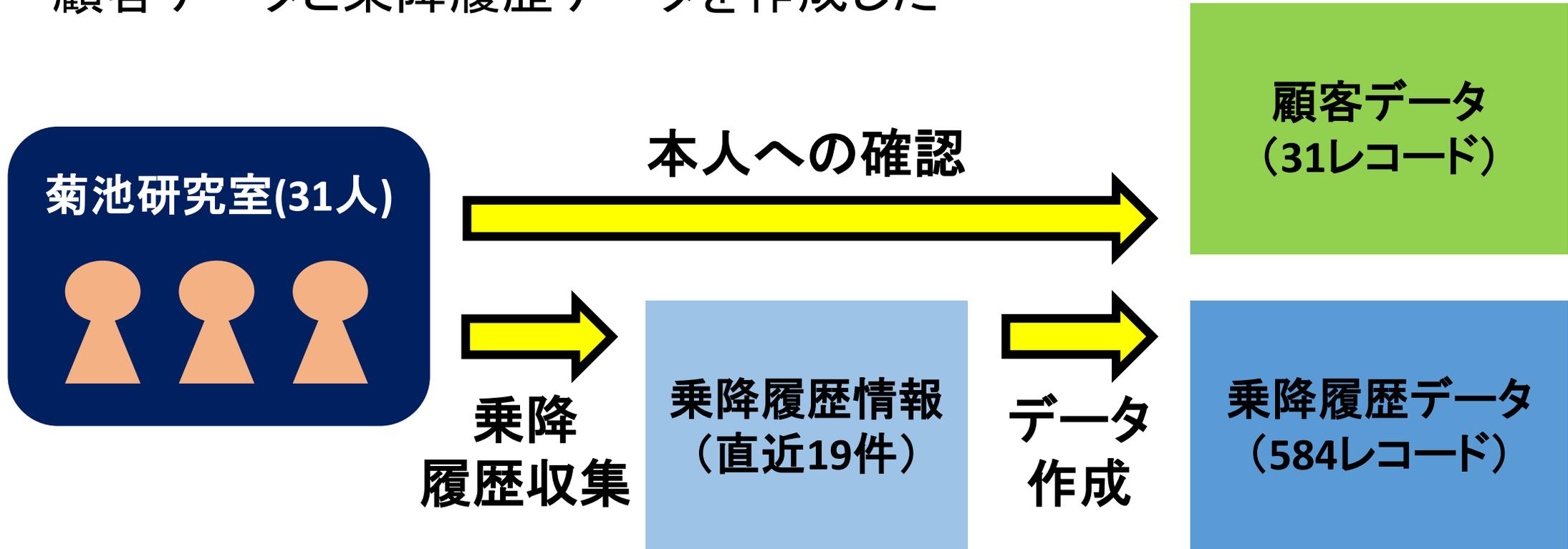
1. 企業内の個人情報にはアクセスが許されない
2. データに対してどういった匿名加工をしたらよいのか不明

## 研究方法

1. 実際に個人情報データを取得し、分析する
2. 取得したデータを用いて匿名加工を行う

# 菊池研究室の乗降履歴データ

菊池研究室に所属する31人の情報や交通ICカードから顧客データと乗降履歴データを作成した



# データの内容

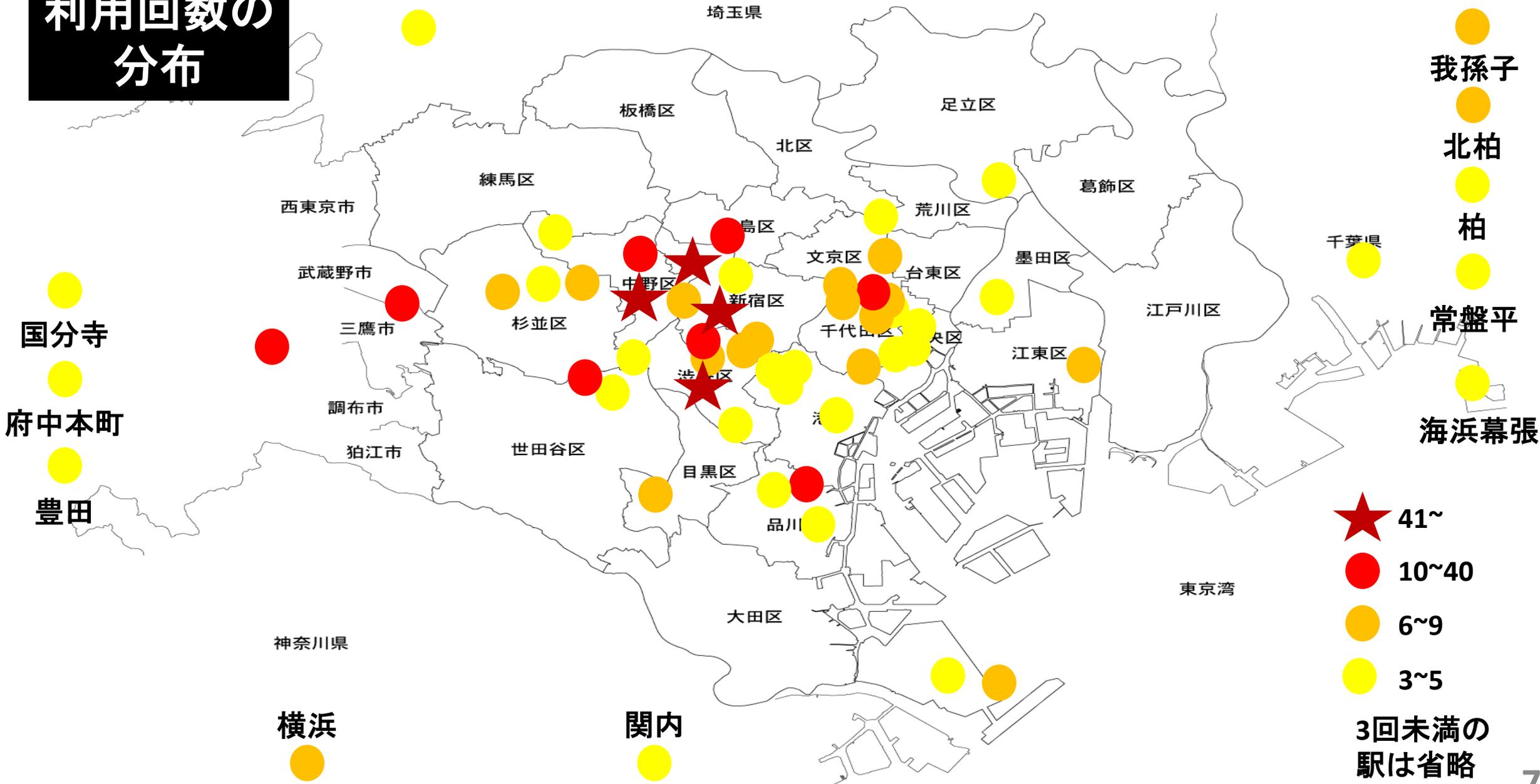
## 顧客データ例

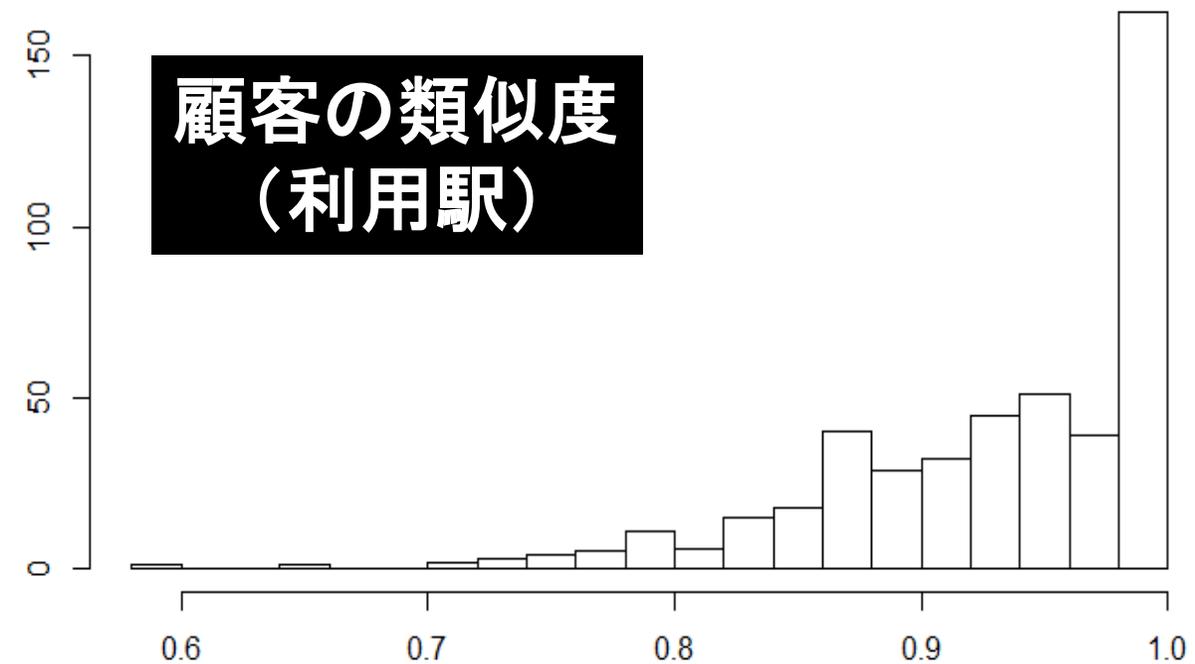
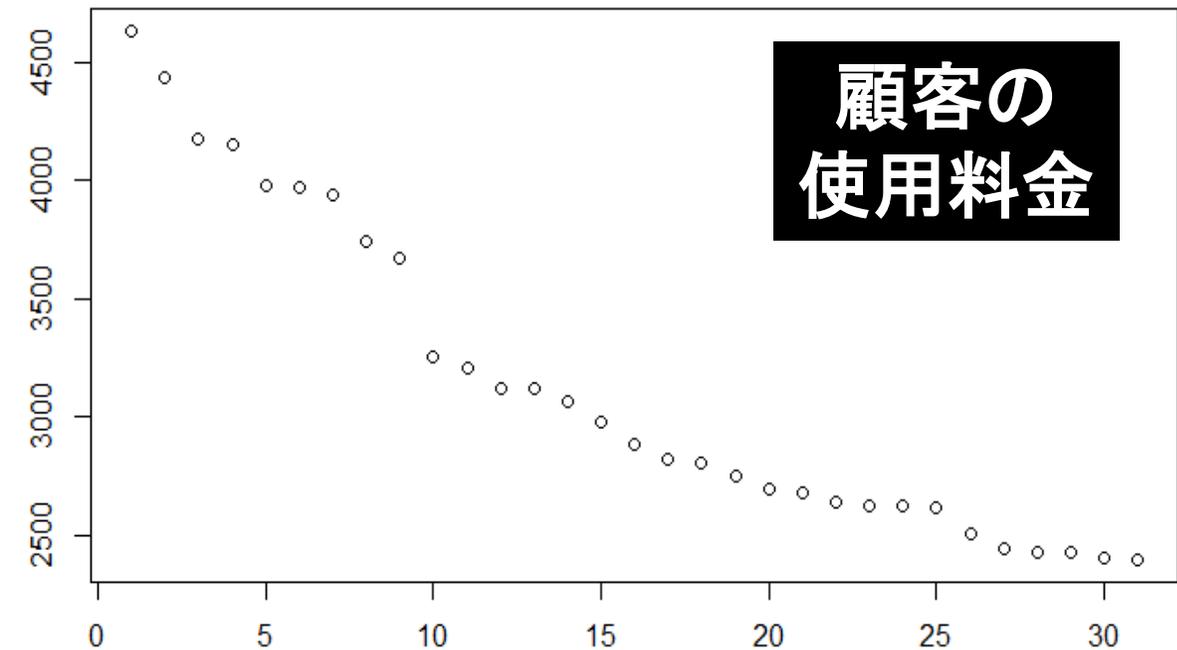
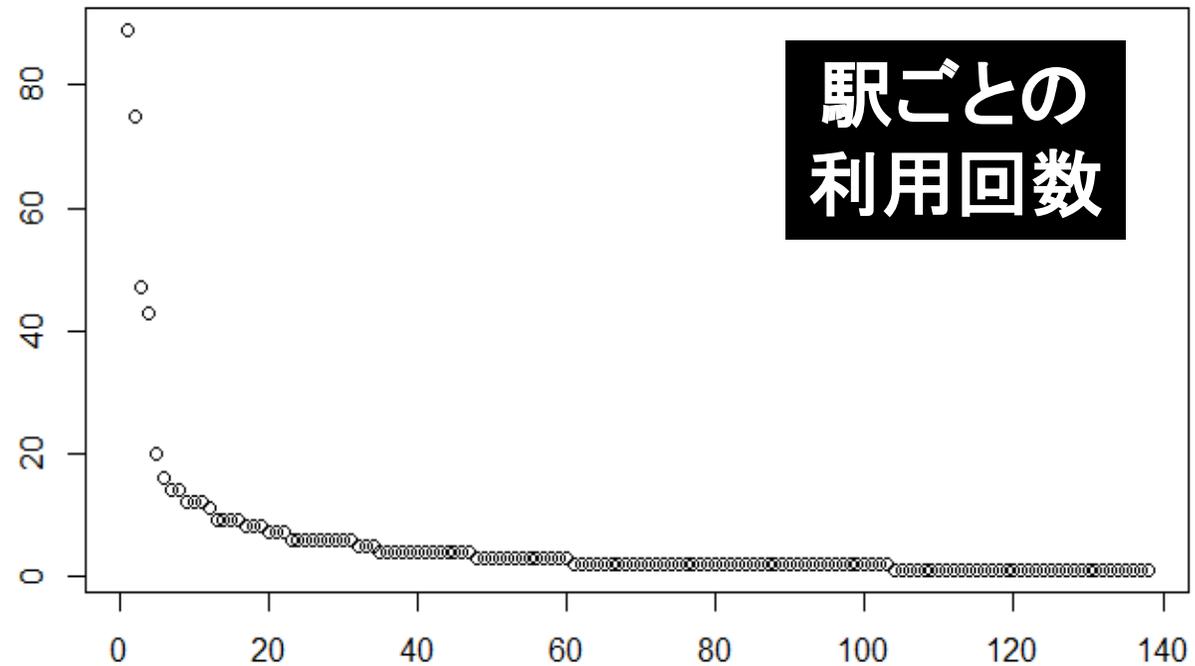
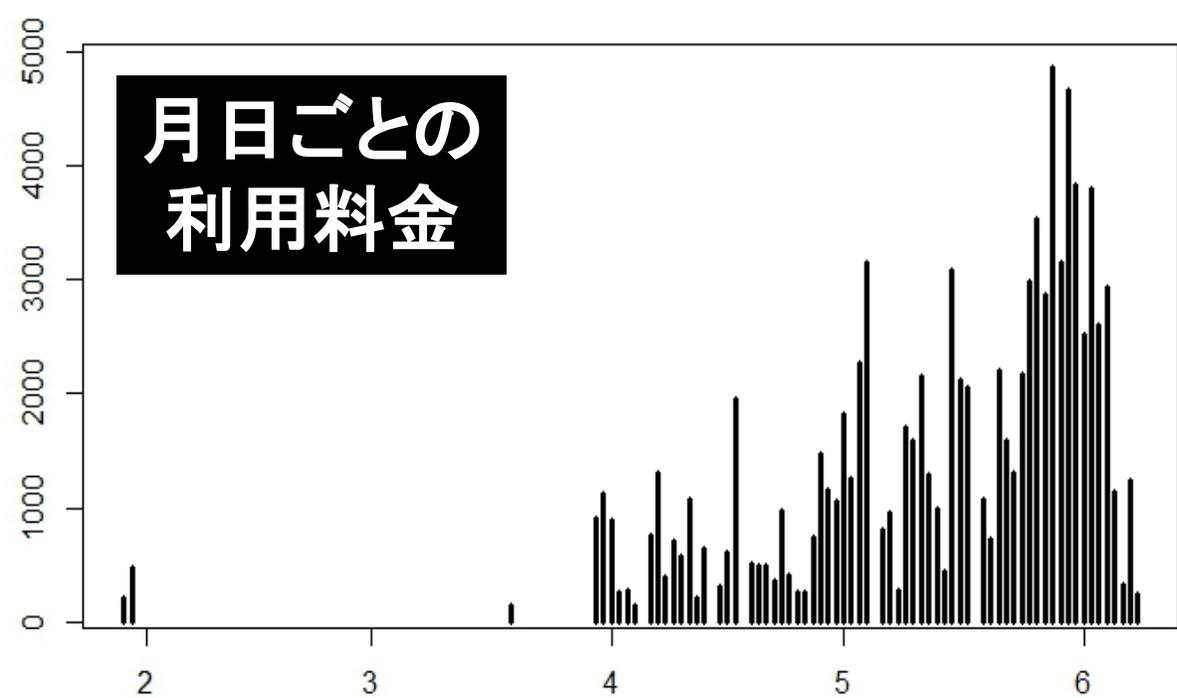
顧客ID	性別	学年	住所	定期券範囲1	定期券範囲2
1	男	1	千葉県	NA	NA
2	女	3	東京都	中野	新宿

## 乗降履歴データ例

顧客ID	日付	回数	乗車駅	降車駅	乗車路線	降車路線	用途	使用場所	料金
1	2016/10/30	2	上野	高田馬場	JR東北本線	JR山手線	交通	NA	-194
1	2016/10/30	1	高田馬場	上野	JR山手線	JR東北本線	交通	NA	-194
1	2016/10/8	1	NA	NA	NA	NA	チャージ	券売機	2000

# 利用回数 の 分布





# 顧客/乗降履歴データに対する匿名加工

今回の場合、提供先の外部組織が知りたいのは

1. どんな学生がどのくらい駅を利用するのか？
2. どの駅がよく利用されているのか？
3. 菊池研究室にはどのような学生がいるのか？

であるため、これらの情報を保持した匿名加工をする

情報を保持できているかどうかは**有用性指標**を用いて評価し、  
データが安全かどうかは**安全性指標**を用いて評価する

# 顧客/乗降履歴データの有用性/安全性指標

データの有用性を評価する指標を3つ用意した(U1, U2, U3)

有用性指標(U1)  
顧客属性ごとの駅利用回数

有用性指標2(U2)  
駅利用回数の順位上位

有用性指標3(U3)  
顧客属性のクロス集計の人数

データの安全性を評価する指標を2つ用意した(S1, S2)

安全性指標(S1)  
特殊な利用駅がない

安全性指標(S2)  
特殊な顧客がない

これらの指標の評価をもとに、  
有用性を保ちつつ安全性を上げる匿名加工をする

# 顧客/乗降履歴データの匿名加工 前

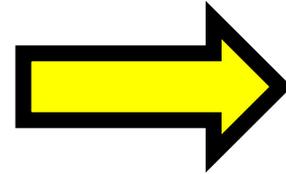
簡易顧客データ M

顧客ID	性別	学年	group
1	男	1	A
2	男	1	A
3	男	2	B
4	男	2	B
5	女	4	C

加工された簡易顧客データ M\*

顧客ID	性別	学年	group
1	男	1	A
2	男	1	A
3	男	2	B
4	男	2	B
5	男	2	B

匿名加工



特殊な顧客の  
属性を書き換える

# 顧客/乗降履歴データの匿名加工 後

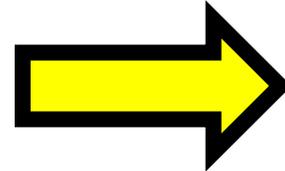
簡易乗降履歴データ

利用駅をグループ内でシャッフルする

加工された簡易乗降履歴データ T\*

顧客ID	乗車駅	降車駅	group	顧客ID	乗車駅	降車駅	group
1	新宿	品川	A	1	新宿	新宿	A
1	品川	新宿	A	1	高田馬場	品川	A
2	高田馬場	新宿	A	2	品川	中野	A
2	新宿	中野	A	2	新宿	新宿	A
3	中野	新宿	B	3	品川	新宿	B
3	新宿	中野	B	3	中野	品川	B
4	高田馬場	品川	B	4	高田馬場	中野	B
4	品川	熱海	B	4	新宿	新宿	B
5	中野	東京	B	5	中野	東京	B
5	東京	中野	B	5	東京	中野	B

匿名加工



特殊な駅を最頻値に書き換える

# 加工手法の評価

	加工前	加工後
U1	0	2.67
U2	0	0
U3	0	0.2
S1	×	○
S2	×	○

データを加工したことにより  
有用性が少し下がった  
(数値が低いほうが有用性が高い)

データを加工したことにより  
安全性が上がった

# まとめ

- 31人の交通ICカードから顧客データと乗降履歴データを取得し、これらのデータの分析を行った
- 取得したデータに対し、想定したケースに対応する評価指標と加工手法を検討した

# 質疑応答用スライド

# 乗降履歴データの匿名加工 1

簡易顧客データ M

顧客ID	性別	学年	group
1	男	1	A
2	男	1	A
3	男	2	B
4	男	2	B
5	女	4	C

簡易乗降履歴データ T

顧客ID	乗車駅	降車駅	group
1	新宿	品川	A
1	品川	新宿	A
2	高田馬場	新宿	A
2	新宿	中野	A
3	中野	新宿	B
3	新宿	中野	B
4	高田馬場	品川	B
4	品川	熱海	B
5	中野	東京	C
5	東京	中野	C

危険な場合1  
利用駅をそのまま  
データを提供する

# 乗降履歴データの匿名加工 2

簡易顧客データ M

顧客ID	性別	学年	group
1	男	1	A
2	男	1	A
3	男	2	B
4	男	2	B
5	女	4	C

顧客属性が同じグループ内で  
利用駅をシャッフルする  
(U1~U3は損なわない)

簡易乗降履歴データ T\*

顧客ID	乗車駅	降車駅	group
1	新宿*	新宿*	A
1	高田馬場*	品川*	A
2	品川*	中野*	A
2	新宿*	新宿*	A
3	品川*	新宿*	B
3			
4			
4			
5			
5			

危険な場合2  
顧客属性の  
組み合わせが  
独特のものがある

# 乗降履歴データの匿名加工 3

簡易顧客データ M\*

顧客ID	性別		
1	男		
2	男		
3	男		
4	男		
5	男*	2*	B*

危険な場合3  
特殊な利用駅  
(利用回数が少ない・  
場所が遠い)がある

簡易乗降履歴データ T\*

	乗車駅	降車駅	group
	新宿	新宿	A
	高田馬場	品川	A
	品川	中野	A
2	新宿	新宿	A
3	品川	新宿	B
3	中野	品川	B
4	高田馬場	中野	B
4	新宿	熱海	B
5	中野	東京	B*
5	東京	中野	B*

顧客属性の組み合わせが独特な  
ユーザーを別グループに変える  
(U1, U3を少し損なう)

# 乗降履歴データの匿名加工 4

簡易顧客データ M\*

顧客ID	性別	学年	group
1	男	1	A
2	男	1	A
3	男	2	B
4	男	2	B
5	男	2	B

簡易乗降履歴データ T\*\*

顧客ID	乗車駅	降車駅	group
1	新宿	新宿	A
1	高田馬場	品川	A
2	品川	中野	A
2	新宿	新宿	A
3	品川	新宿	B
3	中野	品川	B
4	高田馬場	中野	B
4	新宿	新宿*	B
5	中野	東京	B
5	東京	中野	B

特殊な駅(利用回数が少ない・  
場所が遠い)を利用回数1位の  
駅に書き換える  
(U1~U3は損なわない)

# 有用性指標の詳細

$$U_1(M, T, M^*, T^*) = \frac{\sum_{g=1}^g |T_{\text{station}}(X_i) - T^*_{\text{station}}(X_i)|}{g}$$

$$U_2(M, T, M^*, T^*) = 5 - |\{s \in S_5(\text{rank}(T, s) = \text{rank}(T^*, s))\}|$$

$$U_3(M, T, M^*, T^*) = \frac{\sum_{\text{num}(\text{sex})}^{i=1} \sum_{\text{num}(\text{grade})}^{j=1} |\text{Cross}_{\text{sex,grade}}(i, j) - \text{Cross}^*_{\text{sex,grade}}(i, j)|}{\text{num}(\text{sex}) * \text{num}(\text{grade})}$$

M, T : 元データ

M\*, T\* : 加工されたデータ

T<sub>station</sub>(X<sub>i</sub>) : TについてのグループX<sub>i</sub>の駅利用総回数

g : Tのグループ数

S<sub>5</sub> : 上位5駅の集合

rank(T, s) : 駅sのTにおける利用回数順位

Cross<sub>sex,grade</sub> : Mの(性別, 学年)属性についてのクロス集計値

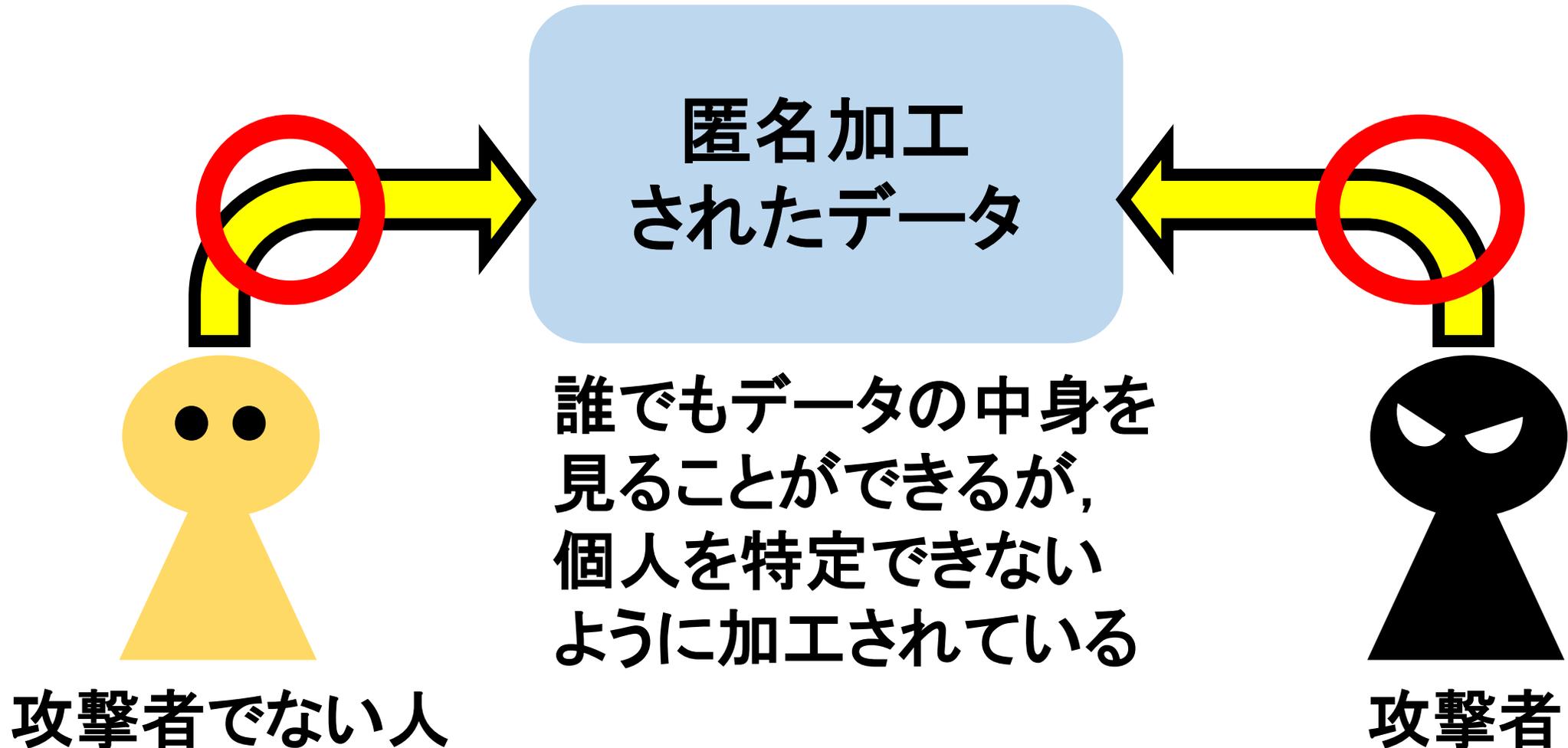
num(属性名) : 属性の種類数

これらの有用性指標の値が0に近いほど、データ(T\*, M\*)の有用性は高い

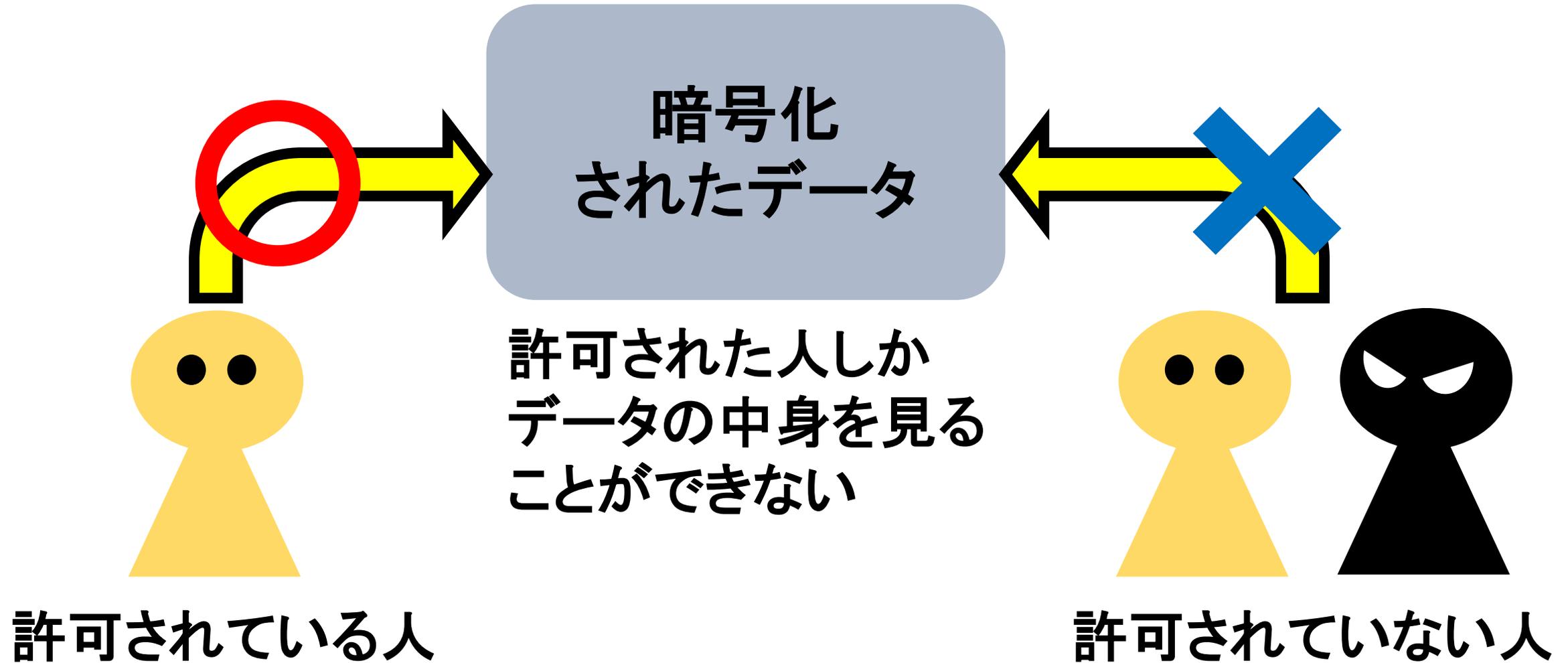
# 有用性指標の計算例(U1)

$$\begin{aligned}U_1(\mathbf{M}, \mathbf{T}, \mathbf{M}^*, \mathbf{T}^*) &= \frac{\sum_{g=1}^g |\mathbf{T}_{\text{station}}(\mathbf{X}_i) - \mathbf{T}_{\text{station}}^*(\mathbf{X}_i)|}{g} \\&= \frac{|\mathbf{T}_{\text{station}}(A) - \mathbf{T}_{\text{station}}^*(A)| + |\mathbf{T}_{\text{station}}(B) - \mathbf{T}_{\text{station}}^*(B)| + |\mathbf{T}_{\text{station}}(C) - \mathbf{T}_{\text{station}}^*(C)|}{3} \\&= \frac{|8 - 8| + |8 - 12| + |4 - 0|}{3} \\&= \frac{0 + 4 + 4}{3} \\&= \frac{8}{3} \\&= 2.66 \dots \approx 2.67\end{aligned}$$

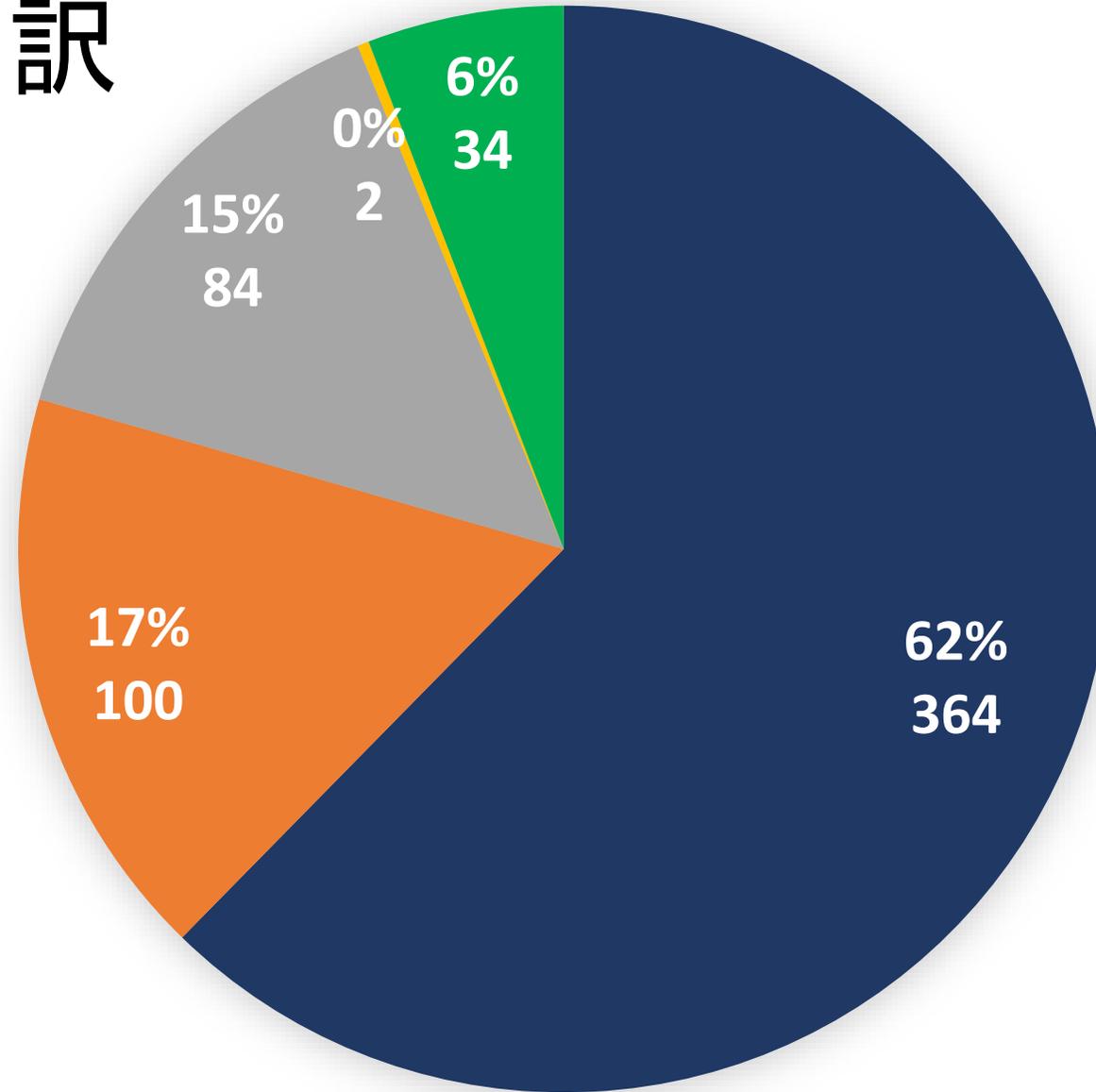
# 匿名加工と暗号化の違い 1



# 匿名加工と暗号化の違い 2

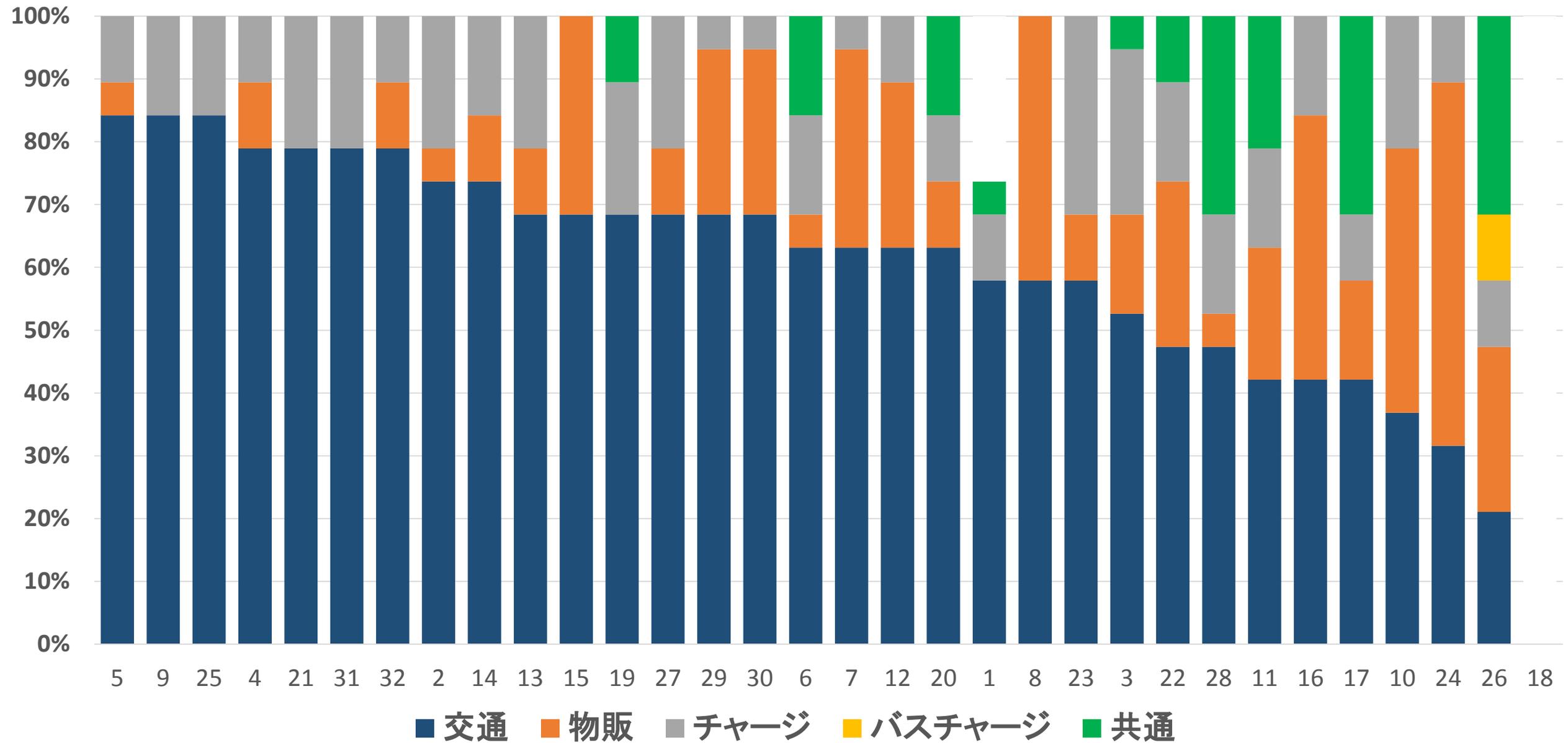


# suicaデータの内訳

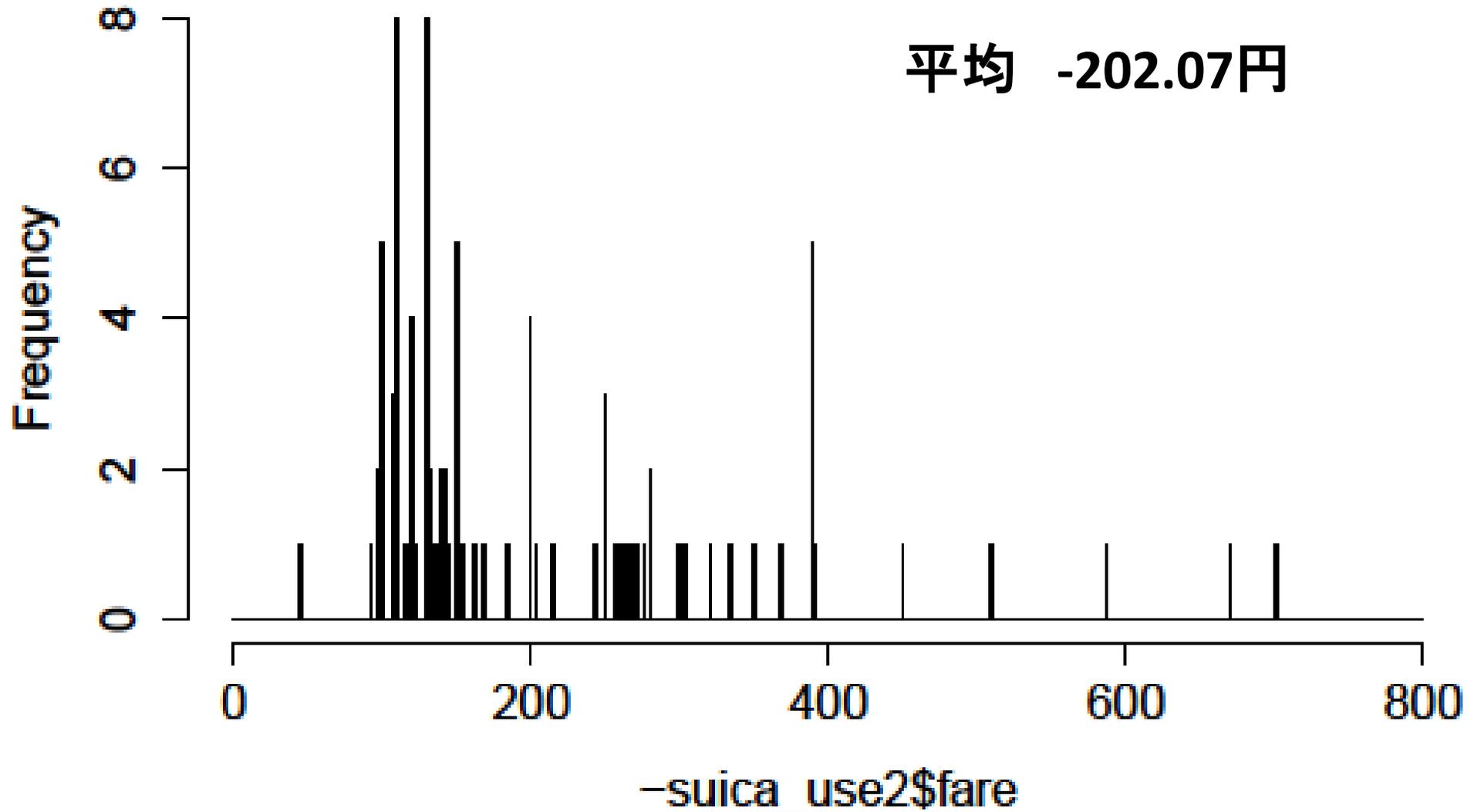


■ 交通 ■ 物販 ■ チャージ ■ バスチャージ ■ 共通

# ユーザごとの用途の内訳



# USE=2(物販) Histogram of -suica\_use2\$fare



# USE=3(チャージ)Histogram of suica\_use3\$fare

