# データマイニング

1: データマイニングの全体像

2: データセキュリティとデータサイエンスの倫理

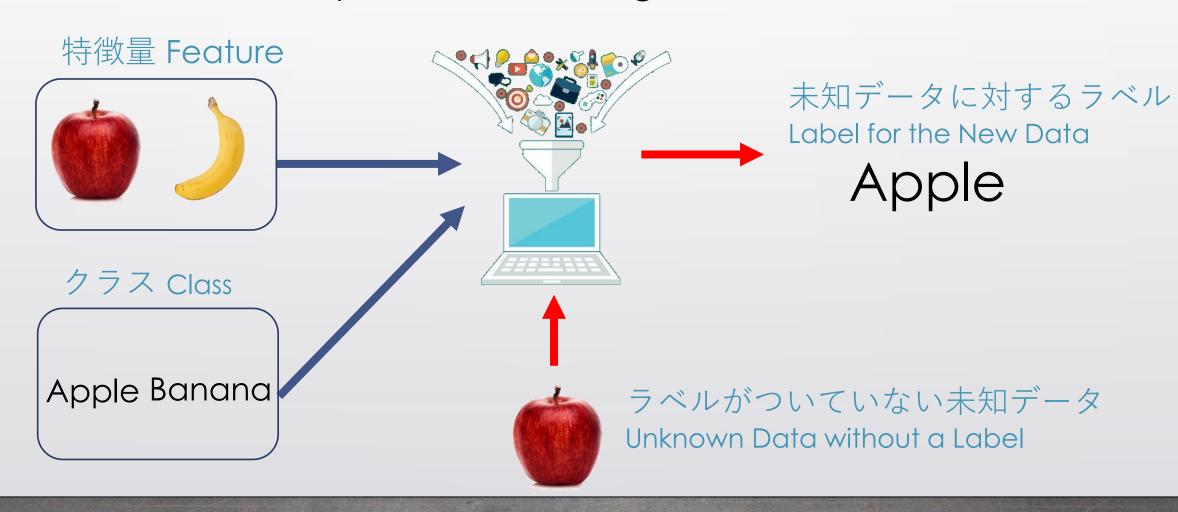
土居 裕和 Hirokazu Doi

長岡技術科学大学 Nagaoka University of Technology

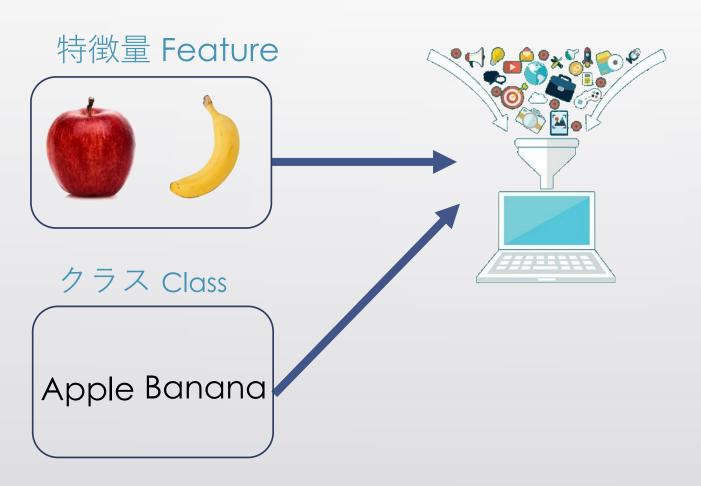
機械学習によるモデル化 Data Modelling by Machine Learning 分類 CLASSIFICATION 教師あり学習 SUPERVISED LEARNING Develop predictive model based on both input and output data 回帰 REGRESSION MACHINE LEARNING UNSUPERVISED LEARNING クラスタリング CLUSTERING Group and interpret data based only 教師なし学習 on input data

Supervised Learning versus Unsupervised Learning (Mathworks, n.d.)

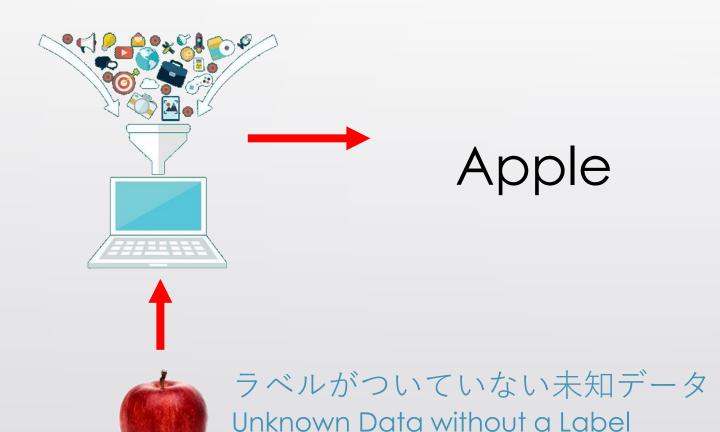
### 教師あり学習 Supervised Learning



# モデルの訓練(トレーニング) Model Training



#### 予測 Prediction



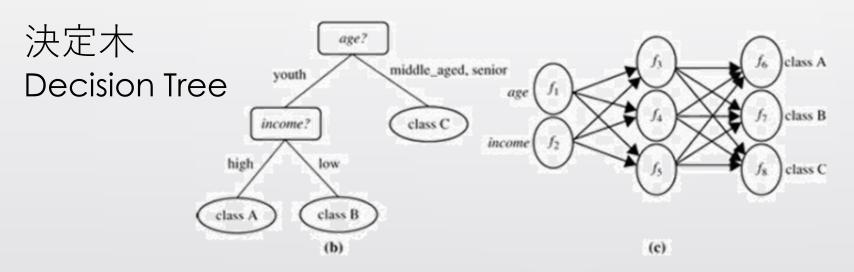
量を予測する回帰に対して、分類ではカテゴリー/クラスを予測する

Classification algorithms predict category/class while regression predicts quantity.

※量とクラスの予測の両方に対応したアルゴ リズムも多い

There are many algorithms that can make predictions about both quantity and category.

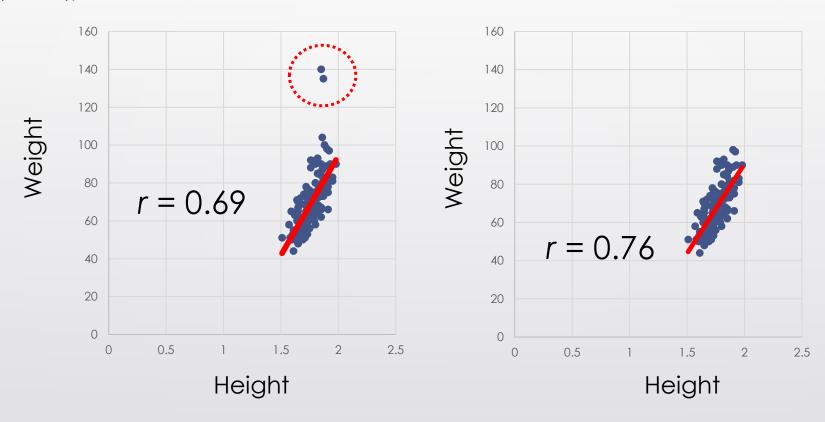
 $age(X, "youth") AND income(X, "high") \longrightarrow class(X, "A")$   $age(X, "youth") AND income(X, "low") \longrightarrow class(X, "B")$   $age(X, "middle_aged") \longrightarrow class(X, "C")$   $age(X, "senior") \longrightarrow class(X, "C")$ 



ニューラル ネットワーク Neural Network

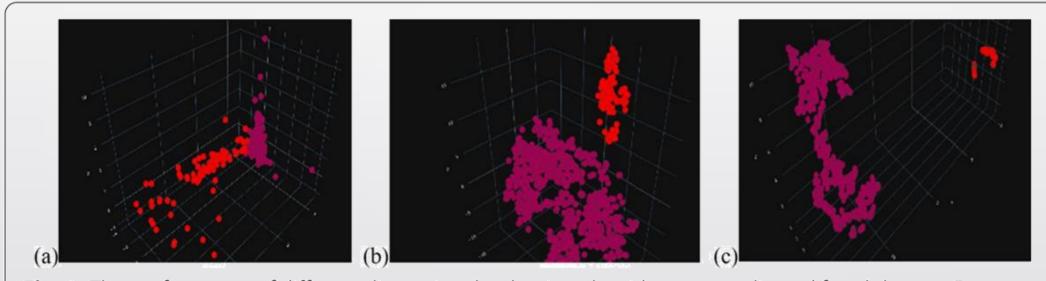
FIGURE 1.9 A classification model can be represented in various forms: (a) IF-THEN rules, (b) a decision tree, or (c) a neural network.

### 外れ値 Outlier



結果に影響を与える可能性がある外れ値の存在をチェックする必要がある Data set should be checked for the existence of outliers that can influence the results

### 外れ値分析による不正検知 Fraud Detection by Outlier Analysis



**Fig. 8** The performance of different dimensional reduction algorithms on credit card fraud dataset. Feature size is reduced to dimension 3 by **a** PCA, **b** t-SNE and **c** UMAP

Benchaji et al. Journal of Big Data (2021) 8:151

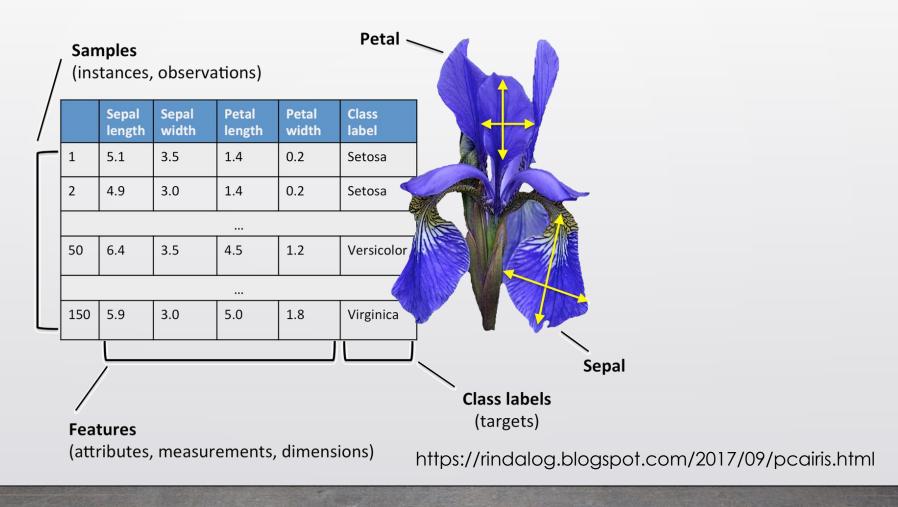
### データマイニングの流れ Steps in Data Mining

- 1. 目標設定 Goal Setting
- 2. データ収集 Data collection
- 3. 前処理 Preprocessing
- 4. 特徴量選択 Feature Selection (必要ないケースもある; can be skipped in some cases)
- 5. データ分析 Data Analysis・モデリング Modeling
- 6. 性能評価 Performance Evaluation
- 7. (ディプロイメント Deployment)

### 目標設定 Goal Setting

- 1. 調査の目的を設定する Verify the purpose of the survey
- 2. どのようなデータを集めるかを決める Decide what data to collect for our purposes
- 3. 目的を達成するには、どのようなデータ分析を行うのが適切かを、前もって考えておく Consider what type of data analysis is appropriate to achieve the goal

#### アヤメからの特徴量抽出 Feature Extraction from Iris



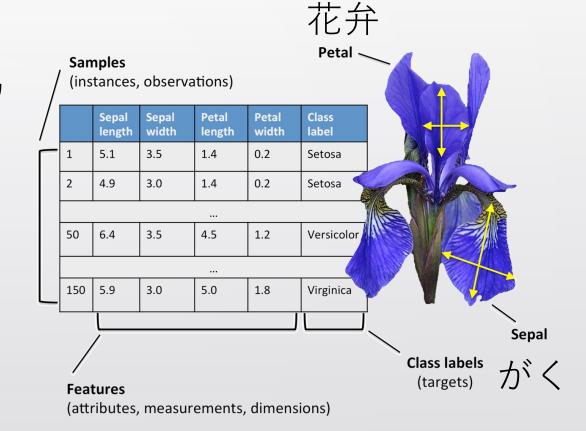
### アヤメデータセット Iris Dataset

ヒオウギアヤメ ブルーフラッグ アイリス・バージニカ









### データの読み込み Loading Data

```
Iris plants dataset
**Data Set Characteristics:**
    :Number of Instances: 150 (50 in each of three classes)
    :Number of Attributes: 4 numeric, predictive attributes and the class
    :Attribute Information:
        - sepal length in cm
        - sepal width in cm
        - petal length in cm
        - petal width in cm
        - class:
                - Iris-Setosa
                - Iris-Versicolour
                - Iris-Virginica
```

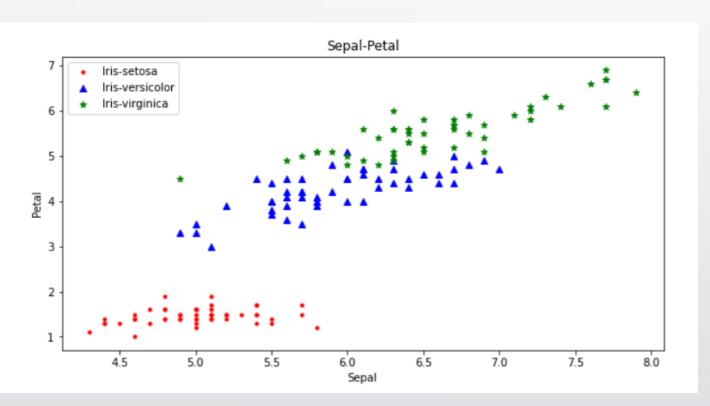
### 記述統計 Descriptive Statistics

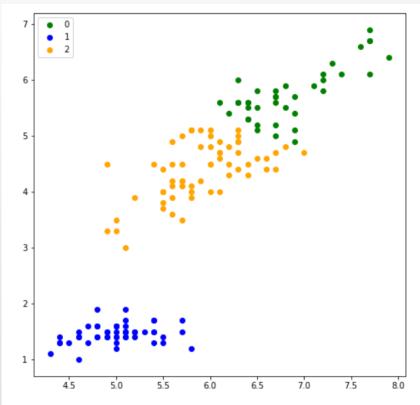
In [15]: df.describe()

Out[15]:

	const longth (om)	const width (om)	notal longth (om)	notal width (am)
	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)
count	150.000000	150.000000	150.000000	150.000000
mean	5.843333	3.057333	3.758000	1.199333
std	0.828066	0.435866	1.765298	0.762238
min	4.300000	2.000000	1.000000	0.100000
25%	5.100000	2.800000	1.600000	0.300000
50%	5.800000	3.000000	4.350000	1.300000
75%	6.400000	3.300000	5.100000	1.800000
max	7.900000	4.400000	6.900000	2.500000







	copui iongin (om)	sepai width (cm)	petal length (cm)	petal width (cm)	target
0	5.1	3.5	1.4	0.2	setosa
1	4.9	3.0	1.4	0.2	setosa
2	4.7	3.2	1.3	0.2	setosa
3	4.6	3.1	1.5	0.2	setosa
4	5.0	3.6	1.4	0.2	setosa
					/ \.
145	6.7	3.0	5.2	2.3	virginica
146	6.3	2.5	5.0	1.9	virginica
147	6.5	3.0	5.2	2.0	vrginica
148	6.2	3.4	5.4	2.3	virginica
149	5.9	3.0	5.1	1.8	virginica



	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)
0	5.1	3.5	1.4	0.2
1	4.9	3.0	1.4	0.2
2	4.7	3.2	1.3	0.2
3	4.6	3.1	1.5	0.2
4	5.0	3.6	1.4	0.2
145	6.7	3.0	5.2	2.3
146	6.3	2.5	5.0	1.9
147	6.5	3.0	5.2	2.0
148	6.2	3.4	5.4	2.3
149	5.9	3.0	5.1	1.8

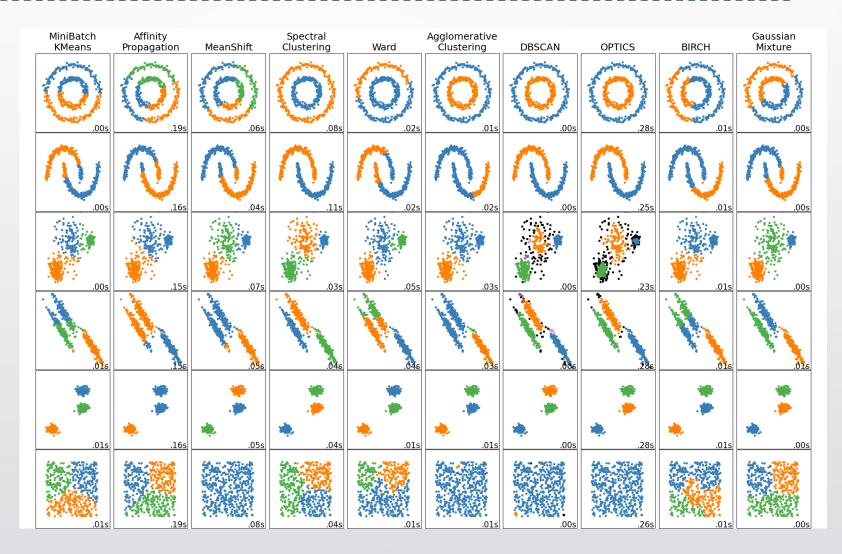
### クラスタリング

Final result of clustering depends on

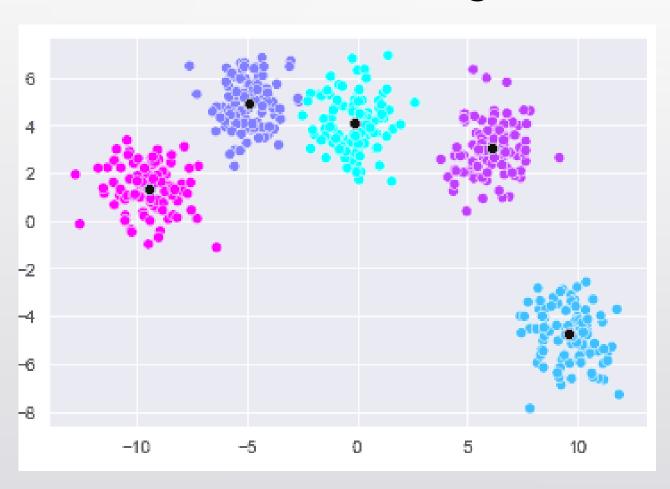
Type of algorithm アルゴリズムの種類

Parameter Setting パラメータ設定

> https://scikitlearn.org/stable/ modules/clusterin g.html



### K-Means Clustering



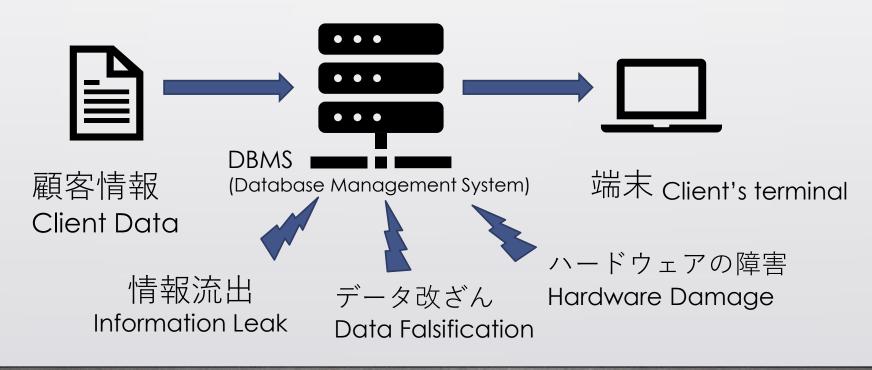
クラスターの数を指定しなくては いけない

You have to specify the number of clusters, k.

データセキュリティ Data Security

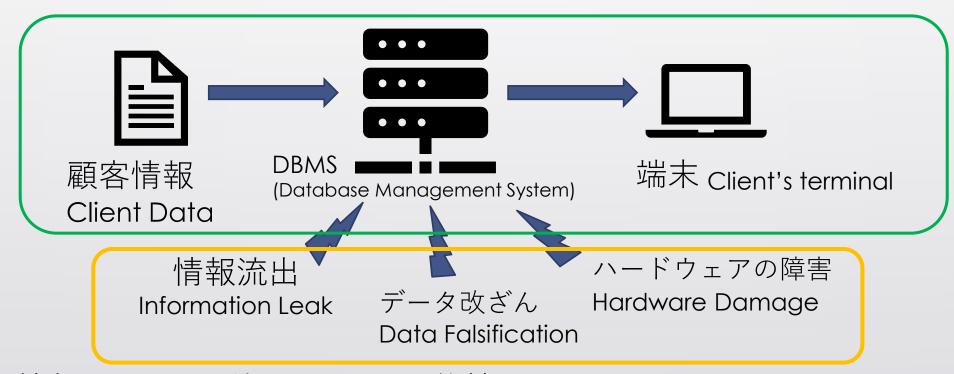
#### 情報資產 Information Asset

個人や企業にとって、価値のある情報や情報システム
Information or the integrity of information system valuable for an individual or corporation.



### リスクと脆弱性 Risk and Vulnerability

脆弱性:情報システムの弱点 Weak point of an information system



リスク:情報システムが損なわれる可能性 Possibility that an information system is damaged

### 機密性 Confidentiality

許可された人のみが情報にアクセスできる Only persons granted permission can access information



アクセス制限 Access Control

サーバールームの施錠 Locked Server Room

暗号化 Encryption

Confidentiality

https://www.simplilearn.com/what-is-information-security-article

### 完全性 Integrity

改ざんなどされることなく、情報が完全に保たれている Information is being kept intact without falsification etc



アクセスログ・操作 Logging Access and Operation 直感的なユーザー・インタフェース Intuitive User Interface

Integrity

https://www.simplilearn.com/what-is-information-security-article

### 可用性 Availability

許可された人が、必要な時に、いつでも情報にアクセスできる Any persons grated permission can access information whenever necessary



複数箇所におけるデータ保存 Data storage in multiple sites

定期的なバックアップ Periodical BackUp

**Availability** 

https://www.simplilearn.com/what-is-information-security-article

AUTO & TRUCK MANUFACTURERS 2011年4月28日 / 6:00 午前 / 12年前更新

# UPDATE 1-Sony may face global legal scrutiny over breach

Tom Hals, Leigh Jones

5 分で読む



- \* U.S. lawyers focused on data breaches eye legal action
- \* British government watchdog launches investigation
- \* U.S. state attorneys general also discussing incident (Recasts with comments from U.S. legislators)

WILMINGTON, Delaware/NEW YORK, April 27 (Reuters) - Sony Corp 6758.T could face legal action across the globe after it belatedly revealed one of the biggest online data breaches ever.

In the United States, several members of Congress seized on the breach, in which hackers stole names, addresses and possibly credit card details from users of Sony's PlayStation Network, to push for tougher laws protecting personal information.



AP Photo/Shizuo Kambayashi

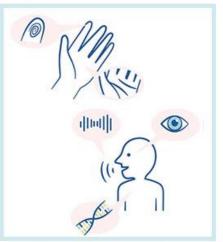
#### 個人情報 Personal Information

#### 個人情報とは



特定の個人を 識別できるもの

Information that identifies particular individual



個人の身体の データ

Data extracted from individual's body



個人に割り振られる 公的な番号

Official number assigned to each 政府広報オンライン individual

「個人情報保護法」

Personal Information Protection Law

生存する個人に関する情報で、氏名、生年月日、 住所、顔写真などにより特定の個人を識別でき る情報

Information concerning living individual... Information by which a particular individual is identifiable

データはあらかじめ決まった目的でしか利用しない Use information only for pre-defined purposes

流出を防止するため、データを安全に管理する Should be managed securely to avoid data leakage



第三者とデータを共有する場合は、関係者の同意を得る Consent must be obtained from relevant individuals when their data is shared with the third party. 請求があった場合には、関連情報を開示する Should disclose any relevant information whenever requested

#### Samarati & Sweeney, 1998

マサチューセッツ州のGIC (Group Insurance Commission)が、氏名を削除

したうえで、医療保険に関連する情報を民間企業に販売

GIC in Massachusetts sold to private companies the data relevant to medical insurance after "anonymyzation", i.e. deleting individuals' names

Medical Data Released as Anonymous

SSN	Name	Ethnicity	Date Of Birth	Sex	ZIP	Marital Status	Problem
		asian	09/27/64	female	02139	divorced	hypertension
		asian	09/30/64	female	02139	divorced	obesity
		asian	04/18/64	$_{ m male}$	02139	married	chest pain
		asian	04/15/64	$_{ m male}$	02139	married	obesity
		black	03/13/63	$_{ m male}$	02138	married	hypertension
		black	03/18/63	$_{ m male}$	02138	married	shortness of breath
		black	09/13/64	$_{ m female}$	02141	married	shortness of breath
		black	09/07/64	female	02141	married	obesity
		white	05/14/61	male	02138	single	chest pain
		white	05/08/61	$_{ m male}$	02138	$_{ m single}$	obesity
		white	09/15/61	$_{ m female}$	02142	widow	shortness of breath

### Samarati & Sweeney, 1998

マサチューセッツ州ケンブリッジの選挙人名簿は20ドルで購入可能だった Voters' list in Cambridge, Massachusetts could be purchased for 20 dollars

	Voter List									
	Name	Address	$\mathbf{City}$	ZIP	DOB	Sex	Party			
)	Sue J. Carlson	1459 Main St.	Cambridge	02142	9/15/61	female	democrat			

### Samarati & Sweeney, 1998

Medical Data Released as Anonymous

SSN	Name	Ethnicity	Date Of Birth	Sex	ZIP	Marital Status	Problem
		asian	09/27/64	female	02139	divorced	hypertension
		asian	09/30/64	female	02139	divorced	obesity
		asian	04/18/64	$_{ m male}$	02139	married	chest pain
		asian	04/15/64	$_{ m male}$	02139	married	obesity
		black	03/13/63	$_{ m male}$	02138	married	hypertension
		black	03/18/63	$_{ m male}$	02138	married	shortness of breath
		black	09/13/64	$_{ m female}$	02141	married	shortness of breath
		black	09/07/64	female	02141	married	obesity
		white	05/14/61	$_{ m male}$	02138	single	chest pain
		white	05/08/61	$_{ m male}$	02138	single	obesity
		white	09/15/61	female	02142	widow	shortness of breath

#### Voter List

Name	Address	City	ZIP	DOB	Sex	Party	
Sue J. Carlson	1459 Main St.	Cambridge	02142	9/15/61	female	democrat	

### 個人情報の例:表1

## Example of Personal Information: Table 1

Identification Number	User ID	Name	Gender	Age	Income
識別ナンバー	ユーザーID	氏名	性別	年齢	年収
339829Q	sanapon	真田昌幸	男	26	411万円
905473R	oggi1985	荻野吟子	女	33	536万円
099878L	murachan	紫式部	女	39	681万円
013214H	shozan.s	佐久間象山	男	23	309万円

### 個人情報の例:表2

#### Example of Personal Information: Table 2

User ID	Product	Price	Date	Store
ユーザーID	購買物品	購買価格	購買日時	購買店舗
murachan	人参	100円	2021/2/3 18:09	Cマート 代田2丁目店
oggi1985	バナナ	150円	2021/2/3 21:13	Cマート 小石川店
oggi1985	ダイエットコーク	160円	2021/2/4 21:15	Cマート 小石川店
murachan	粉ミルク	980円	2021/2/4 21:16	Cマート 代田2丁目店
murachan	紙おむつ	1700円	2021/2/4 21:16	Cマート 代田2丁目店

名前が削除されていても、表に含まれた情報がら、個人を特定で きる場合がある

Even if name is deleted, sometimes, an individual can be identified from information contained in a table.

### 仮名化 Pseudonymization 匿名化 Anonymization

識別ナンバー	ユーザーID	氏名	性別	年齢	年収	住所
339829Q	sanapon	真田昌幸	男	26	411万円	. \ ./
905473R	oggi1985	荻野吹子	女	33	536万円	$\cdot \cdot \mathbf{X}$
0998/8L	murachan	紫式部	女	39	681万円	• / •
0/3214H	shozan.s	佐久間象山	男	23	309万円	/ · · \

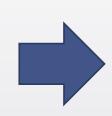


ユーザーID	性別	年齢	年収
sanapon	男	26	411万円
oggi1985	女	33	536万円
murachan	女	39	681万円
shozan.s	男	23	309万円

個人識別情報を削除する Deleting person identifiable information

### *k*-匿名化 *k*-Anonymization

ユーザーID	性別	年齢	年収
sanapon	男	26	411万円
oggi1985	女	33	536万円
murachan	女	39	681万円
shozan.s	男	23	309万円



性別	年齢	年収	
男	[20-29]	[300-499]	万円
女	[30-39]	[500-699]	万円
女	[30-39]	[500-699]	万円
男	[20-29]	[300-499]	万円

データの中に、同じ属性を持つ人が、少なくともk (>=2)人含まれる

At least k (>= 2) persons with the same attributes are included in a table.

データサイエンスの倫理 Ethics in Data Science

### 一般データ保護規則 General Data Protection Rule (GDPR)

EU/UKに住む人々の 個人情報が対象 What?

STA

Data Protection regulation that applies to processing personal data of **EU/UK residents** 

Which?



Any information relating to EU/UK citizens whether they can be identified directly or indirectly

個人情報を保護し、 プライバシーを守る Why?



To protect personal data from misuse and to ensure data privacy

How?

More **obligations** on Data Controller & provide **rights** to data owners to control their data 個人が、自分の個人 情報をコントロール できるようにする



Global

Applies **globally** to any organization processing information on EU/UK residents

EU/UKに住む人々の データを扱う、**世界中 の機関**に適用される



Penalty

Penalties up to 4% (or €20m whichever is higher) for major breaches

Privacy Preservation in Federated Learning: Insights from the GDPR Perspective

### 忘れられる権利 The Right to be Forgotten

13 May 2014 - The court's decision comes by appeal of Mario Costeja González, a Spanish man who sought to remove evidence of his home's repossession ...

Some results may have been removed under data protection law in Europe. Learn more



Unknown - Use precise location - Learn more

Help Send feedback

Privacy & Terms

Retrieved on 2023/01/28 from https://www.cima.ned.org/publication/right-to-be-forgotten-threat-press-freedom-digital-age/

### 透明性とトラスト Transparency and Trust

#### 透明性 Transparency:

データマイニングシステムやAIが、結論を導くプロセスが、人々に説明可能になっていること

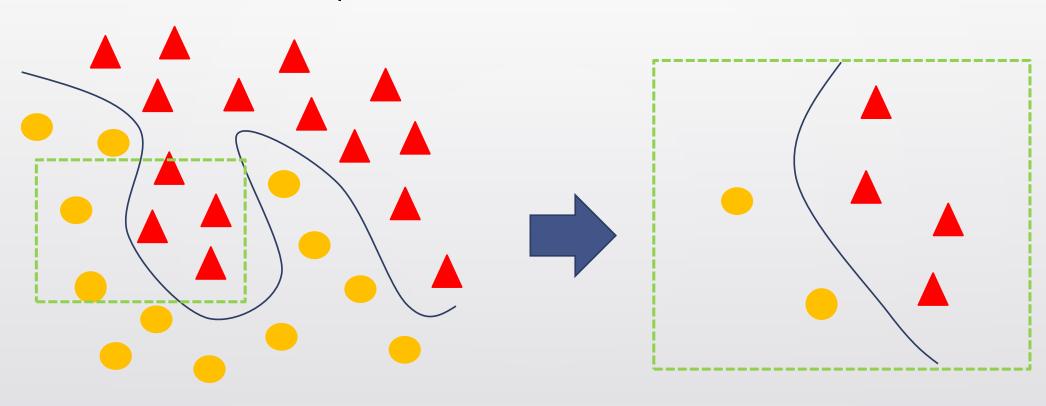
To make the process of reaching conclusions by data-mining system and Al understandable and accessible to ordinary people

#### トラスト Trust:

人々が、そのデータ処理プロセスを理解したうえで、データマイニングシステムやAIを信頼すること

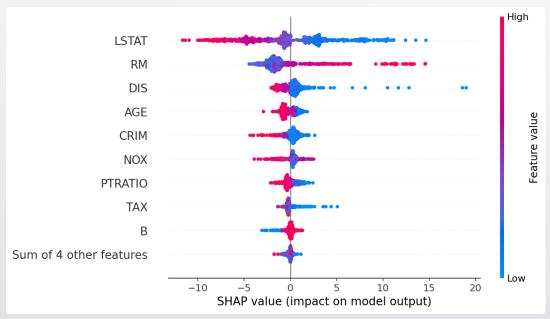
People trust data-mining system and AI after understanding their process of data processing

# 説明可能なAl Explainable Al



### 説明可能なAl Explainable Al

#### SHAP value



https://github.com/slundberg/shap

#### Grad CAM (Gradient-weighted Class Activation Mapping)

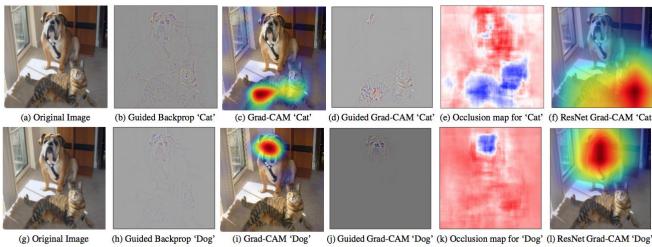


Figure 1: (a) Original image with a cat and a dog. (b-f) Support for the cat category according to various visualizations for VGG and ResNet. (b) Guided Backpropagation [46]: highlights all contributing features. (c, f) Grad-CAM (Ours): localizes class-discriminative regions, (d) Combining (b) and (c) gives Guided Grad-CAM, which gives high-resolution class-discriminative visualizations. Interestingly, the localizations achieved by our Grad-CAM technique, (c) are very similar to results from occlusion sensitivity (e), while being orders of magnitude cheaper to compute. (f, 1) are Grad-CAM visualizations for ResNet-18 layer. Note that in (d, f, i, 1), red regions corresponds to high score for class, while in (e, k), blue corresponds to evidence for the class. Figure best viewed in color.