

データサイエンス・オンライン講座 特別セミナー

# 今日から始められる、 Pythonによるビジネスデータ解析プログラミング講座

2024/11/09

野口 怜



# Agenda

## ■ Part 1 : 今日から始めるための Pythonプログラミングの基本 (40min)

- ✓ Google Colaboratoryの使い方 (15min)  
Google Colab使い方演習
- ✓ Pythonプログラミングの基本 (25min)  
Python基本文法演習

## ■ Part 2 : ビジネスデータ解析に 必要なスキル演習 (1hr50min)

- ✓ データ分析の進め方 (5min)
- ✓ Pythonでできる分析実体験 : 分析デモンストレーション (5min)
- ✓ 分析に必要なデータサイエンス力とデータエンジニアリング力
  - 記述統計学の基礎 (20min)
    - Pythonによる記述統計量算出演習
  - データ観察の基本 (20min)
    - Pythonによるグラフ描画演習
  - データ加工の基本 (20min)
    - Pythonによるデータ加工演習 (ダミー変数化、外れ値・欠損値の処理など)
  - データ分析 (モデル構築) の基本 (40min)
    - Pythonによるモデル構築演習 (重回帰分析、決定木分析など)

## ■ Part 3 : 実際のビジネスデータ解析 に向けた実践演習 (2.5hr)

- ✓ e-Stat, SSDSEの紹介と使い方 (15min)
- ✓ データ分析の実践演習 (2hr)
  - 決定木分析によるビジネス意思決定サポート
  - 階層的クラスタリングによるデータ層別と戦略ポイントの検討
- ✓ 今後の継続学習、実践活用のためのポイント (15min)

※青字はGoogle Colaboratory による演習

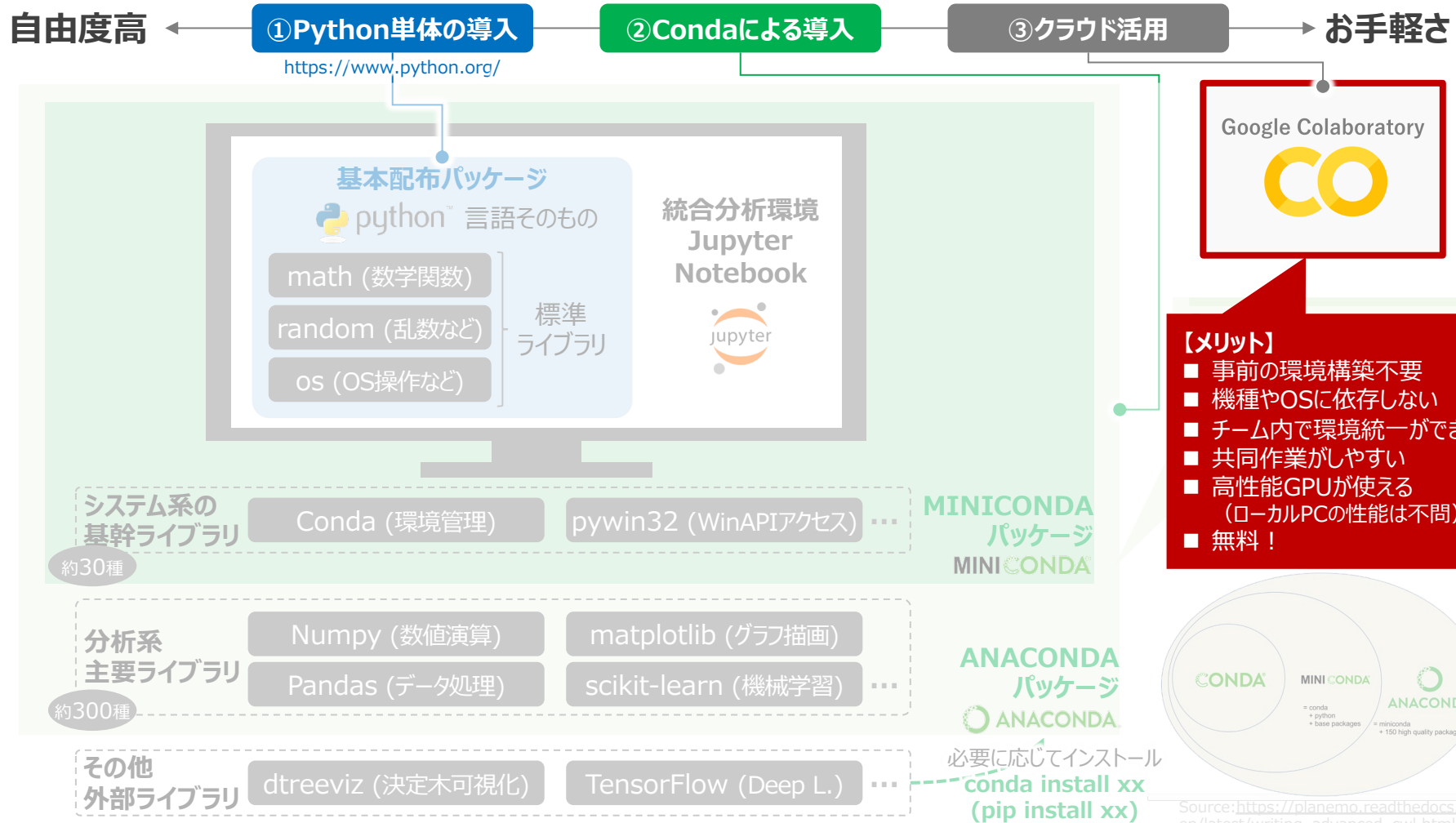
# Part 1

## 今日から始めるための Pythonプログラミングの基本

- ✓ Pythonプログラミング環境：Google Colaboratoryの使い方
- ✓ Pythonプログラミングの基本

# Python導入方法の3パターンとGoogle Colaboratory

- Pythonの初期環境構築には、①Python単体の導入、②Condaによる導入、③構築済のクラウド活用 (Google Colaboratory) の大きく3パターンが存在。環境統一や共同作業を行う場合は、③が適している



※今のところ、miniconda+conda-forgeの組合せは無償!



# Google Colaboratoryへのアクセス方法とノートブックの作成

- Google Colaboratory はインストール不要のPython開発環境であり、Googleのサーバーにインターネット接続（クラウド環境）することで、ブラウザ上でコーディング、実行ができる
- 必要なデータも随時、Googleのサーバーにアップする必要がある

1

Google Colaboratory へアクセス  
<https://colab.research.google.com>

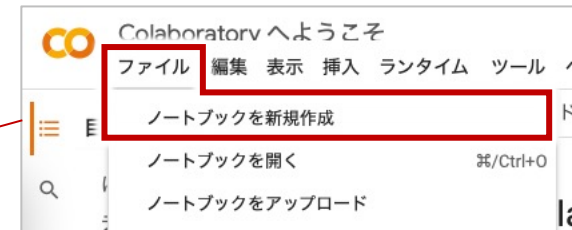
2

右上の [ログイン] ボタンより、Googleアカウントでログイン



3

左上の [ファイル] → [ノートブックを新規作成]



# Google Colaboratory画面の概要

## 画面構成

このスクリーンショットは、Google Colaboratoryのインターフェースの概要を示しています。赤い枠と線は、重要な機能や設定を指しています。

- コード目次**: 左側のメニューアイコン。
- ファイル一覧**: 左側のファイル管理アイコン。
- + コード** / **+ テキスト**: 上部のメニューからコードセルやテキストセルを追加するためのボタン。
- テキストセルを追加 (見出しなどを作りたい場合)**: テキストセル追加の注釈。
- コードセルを追加**: コードセル追加の注釈。
- コードは基本的に自動保存!**: 中央下部の重要なメッセージ。
- 共有** / **設定**: 右上の共有と設定ボタン。
- クリックでファイル名変更可能**: 共有ボタンの注釈。
- ノートブックを共有したい場合**: 共有ボタンの注釈。
- 設定画面 (以下オススメ設定)**: 設定ウィンドウの注釈。

**設定画面 (以下オススメ設定)**

- フォントサイズ**: フォントサイズを14pxに設定。
- 行番号**: 行番号を表示するオプションを有効にする。
- インデントガイド**: インデントガイドの表示を有効にする。

## コードセルの編集

このスクリーンショットは、Google Colaboratoryのコードセル編集画面の概要を示しています。

- 上部の編集メニューからも様々なセル操作が可能**: 編集メニューの注釈。
- 実行ボタン**: コードセルの左側にある実行ボタン。
- コード記載領域**: コードが記載される領域。
- セル操作**: コードセルの右側にある操作メニュー（上へ、下へ、リンク、コメント追加、設定画面で開く、削除、その他）。
- + コード** / **+ テキスト**: コードセル/テキストセルを追加するためのボタン。
- コードセル/テキストセルを追加 (マウスオーバーで表示される)**: ボタンの注釈。

# 参考：Pythonの統合開発環境 “Jupyter Notebook”

- Pythonの（ローカルな）統合開発環境として良く用いられる “Jupyter Notebook” は、**Webブラウザベース**の実行環境であり、実行履歴を残しながら、インタラクティブな分析を実現（Google Colab.もこれをベースとしている）
- コードの記述から実行、グラフの出力までを一貫して行える



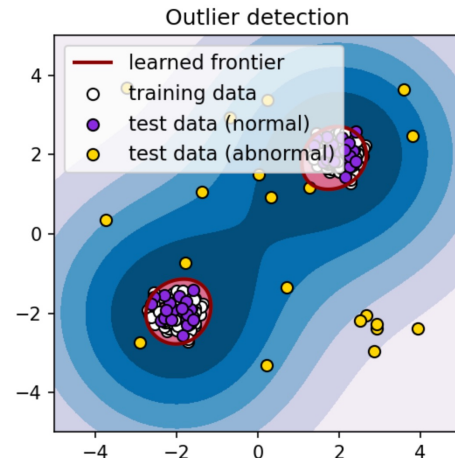
```
[12]: plt.figure(figsize=(4, 4), dpi=150)

plt.title('Outlier detection')
plt.contourf(xx, yy, Z, levels=np.linspace(Z.min(), 0, 7), cmap=plt.cm.PuBu)
a = plt.contour(xx, yy, Z, levels=[0], linewidths=2, colors='darkred')
plt.contourf(xx, yy, Z, levels=[0, Z.max()], colors='palevioletred')

s = 40
b1 = plt.scatter(X_train[:, 0], X_train[:, 1], c='white', s=s, edgecolors='k')
b2 = plt.scatter(X_test[:, 0], X_test[:, 1], c='blueviolet', s=s, edgecolors='k')
c = plt.scatter(X_outliers[:, 0], X_outliers[:, 1], c='gold', s=s, edgecolors='k')

plt.axis('tight')
plt.xlim((-5, 5))
plt.ylim((-5, 5))
plt.legend([a.collections[0], b1, b2, c],
          ['learned frontier', 'training data', 'test data (normal)', 'test data (abnormal)'],
          loc='upper left',
          prop=matplotlib.font_manager.FontProperties(size=11))
```

[12]: <matplotlib.legend.Legend at 0x173db58e0>



# 「変数」とは

- 変数とは、「データを入れておく箱」のことであり、数値に限らず、日付型や文字列型のデータも格納しておくことができる。変数同士の演算も可能で、実際の分析では変数を用いて行う
- 変数名は自由に設定できるが、使用可能な文字には制限があるため、注意が必要

## ▼プログラミングでは「変数」と呼ばれる「箱」にデータを入れる

例：数値を入れる場合

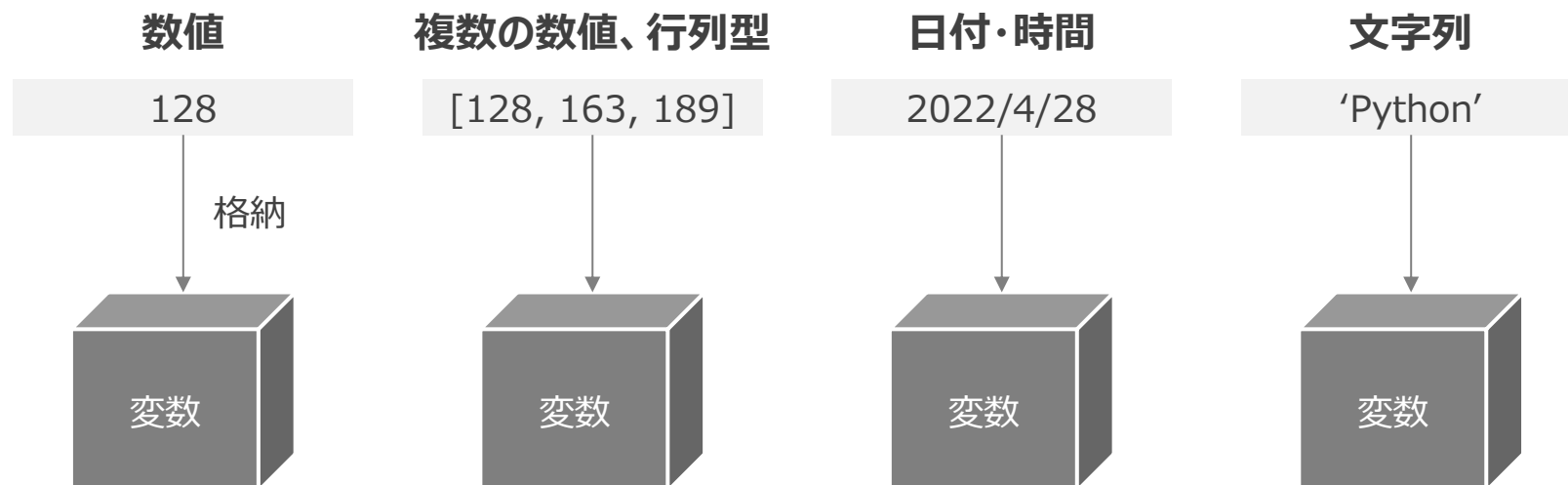
```
num = 10  
print(num)  
>> 10
```

例：文字を入れる場合

```
str = "Hello"  
print(str)  
>> Hello
```

例：演算（足し算）の処理

```
num1 = 1  
num2 = 2  
num1 + num2  
>> 3
```



## ▼変数名の制約

**OK**

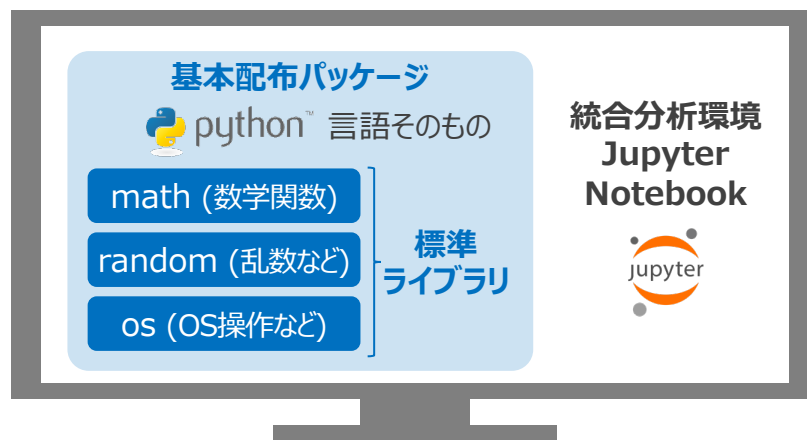
アルファベット（大文字・小文字は区別）  
ひらがな／カタカナ／漢字 ※非推奨  
数字（ただし頭文字は不可）

**NG**

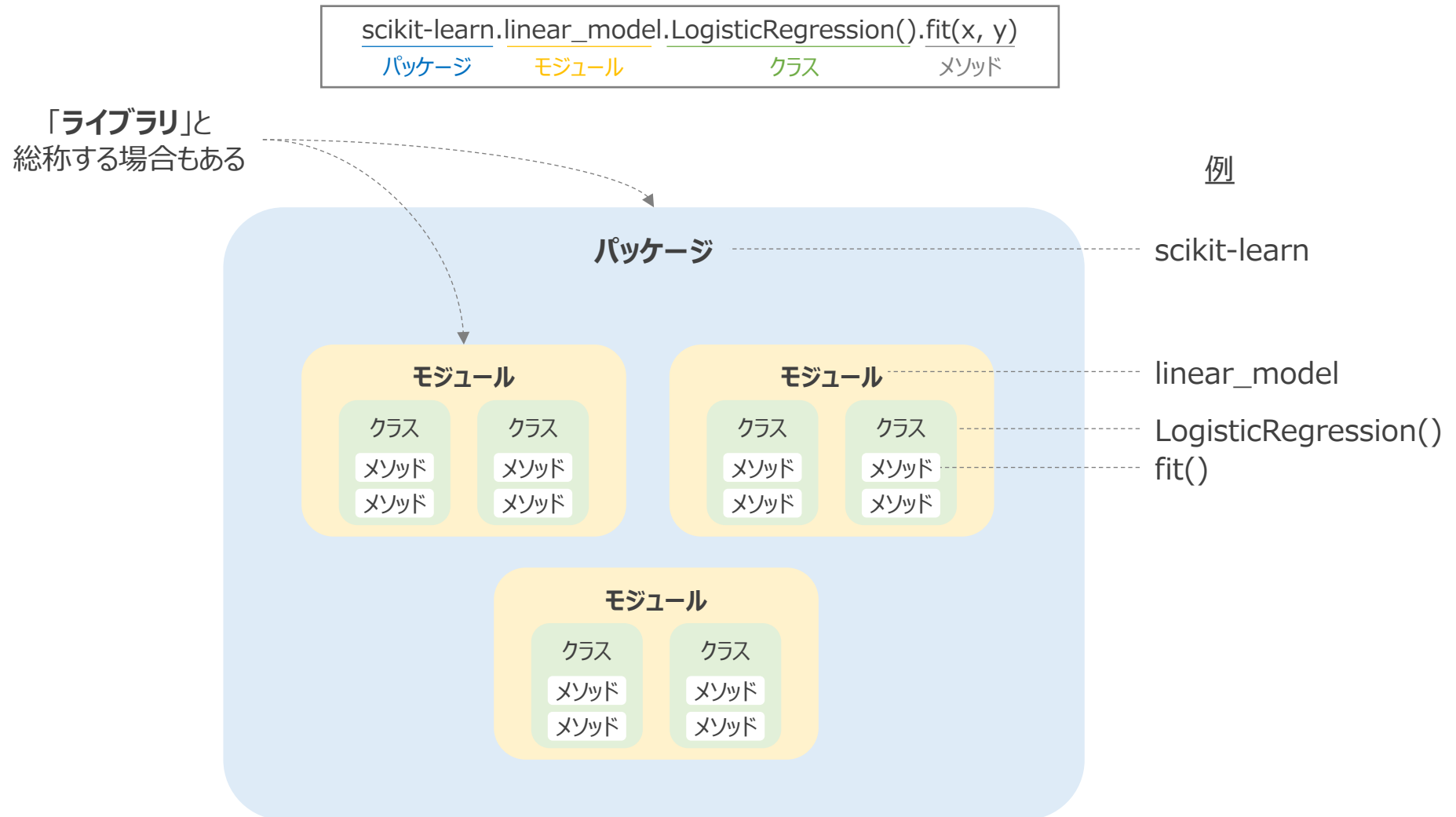
記号（@!?- など。アンダーバー \_ はOK）  
環境依存文字（① 𠬞 など）  
予約語や関数名

# 「ライブラリ」について

- Pythonでは、インストール不要な標準ライブラリが同梱されており、math, random, os, re, などがある
- まずはこれらのライブラリをimportして試してみる



# 用語の整理：「ライブラリ」「パッケージ」「モジュール」「クラス」



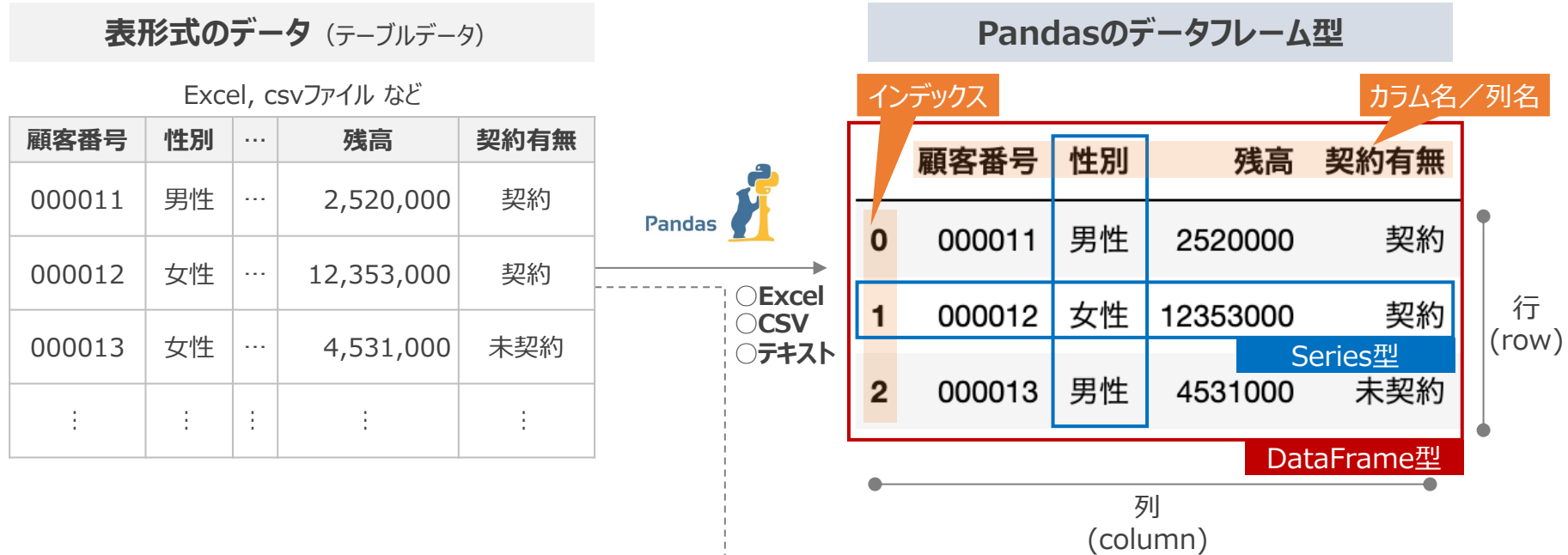
# ビッグデータ分析／機械学習でよく用いられるライブラリ

ライブラリ名 (パッケージ名)	用途
matplotlib	基本グラフ描画
Seaborn	拡張グラフ描画
Numpy	数値演算、ベクトル/行列演算
Pandas	テーブル形式データのデータ処理
scikit-learn	機械学習モデル



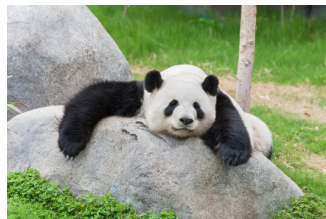
# 表データの処理ライブラリ：pandas

- pandasは、**表形式のデータ**を処理する上で必須のライブラリである
- 表データをそのまま読み込んだ**DataFrame型**と、一部の列（行）のみを抽出した**Series型**とが存在する



## 💡 余談：pandasの由来

諸説ありますが・・・下記が有力です  
Python and data analysis  
panel data



画像出典：Adobe Stock  
<https://stock.adobe.com/jp>

標準ライブラリのみで  
読み込むと・・・

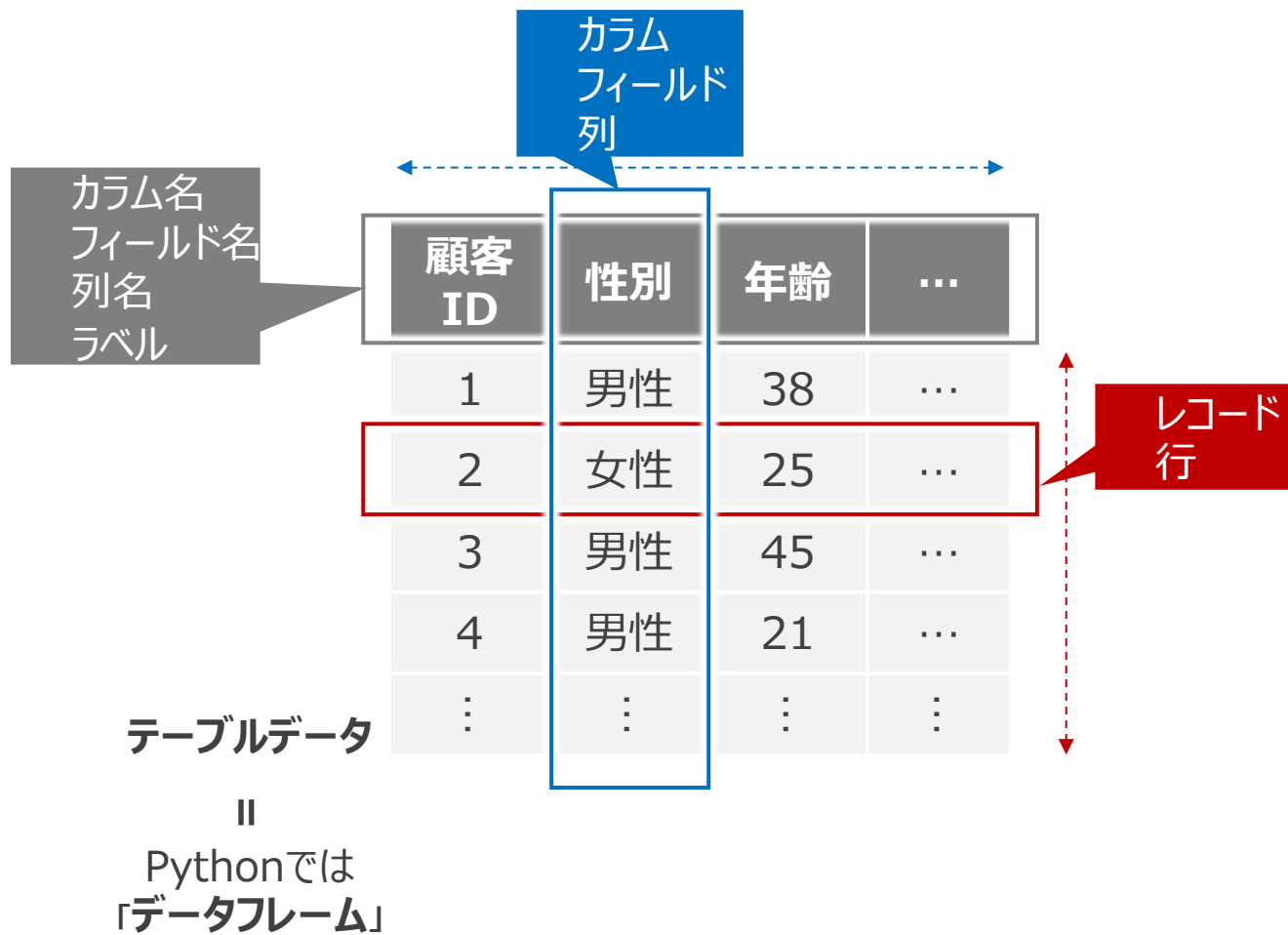
×Excel  
○CSV  
○テキスト

ネストされた2重のリスト  
(2次元配列とも呼ばれる)

```
[[ '000011', '男性', 2520000, '契約'],  
 [ '000012', '女性', 12353000, '契約'],  
 [ '000013', '男性', 4531000, '未契約']]
```



# 参考：「カラム」と「レコード」





# Google Colaboratory上での レクチャー&演習

Part 1 : Google Colab.リンク

<https://colab.research.google.com/drive/1ExQNKoZHe-oMFGUvxU23kyYEipeIxLOY?usp=sharing>

# Part 2

## ビジネスデータ解析に必要なスキルを学ぶ

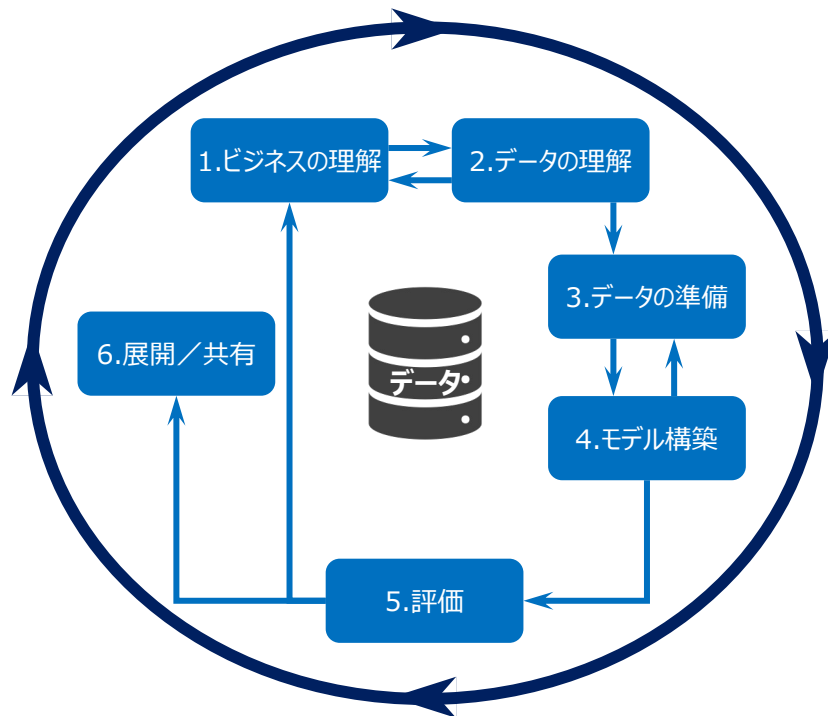
- ✓ データ分析の進め方
- ✓ 分析に必要なビジネス力
- ✓ 分析に必要なデータサイエンス力とデータエンジニアリング力

# データ分析の進め方

- データ分析の進め方に関する方法論「**CRISP-DM**」に基づいて、分析と評価を繰り返して試行錯誤しながら進めるのが一般的である

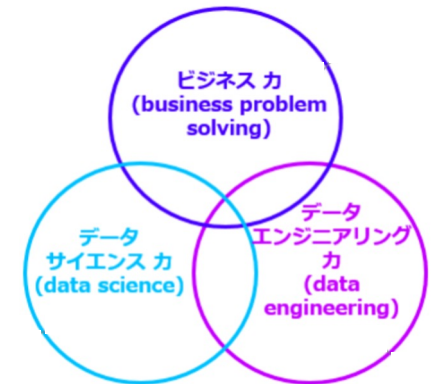
## CRISP-DM: 分析方法論

(CRoss Industry Standard Process for Data Mining)



1. ビジネスの理解	• ビジネス、分析目標の決定 • プロジェクトの立ち上げ	ビジネス
2. データの理解	• データの収集 • データの調査 • データ品質の検証	データサイエンス
3. データの準備	• データの選択や除外 • データのクリーニング • データの構築や統合	データエンジニアリング
4. モデル構築	• モデリング手法の選択 • モデルの作成 • モデルの評価	データサイエンス
5. 評価	• 分析結果の評価 • プロセスの見直し • 実行可能なアクションリストの作成	データエンジニアリング
6. 展開/共有	• 業務への導入計画 • モニタリング、メンテナンスの計画	ビジネス

一般社団法人  
データサイエンティスト協会



# データ分析プロセス（CRISP-DM）とスキルとの対応関係

- データ分析のプロセスCRISP-DMに、3つのスキルセットをマッピングすると以下ようになる
- プロセスによって求められるスキルが変わり、特に最初と最後ではビジネス力が、中盤ではデータサイエンス力とデータエンジニアリング力が求められる

分析プロセス (CRISP-DM)	#	タスク詳細	ビジネス	データ サイエンス	データエンジニア リング
ビジネスの理解	1	顧客のビジネス理解	✓		
	2	課題（分析目的）の設定	✓		
	3	ビジネス目標の設定	✓		
データの理解	4	取得可能データ確認	✓	✓	
	5	分析アプローチ検討	✓	✓	
	6	データの関係性理解/各データ項目の意味理解	✓	✓	
	7	説明変数/目的変数の設定	✓	✓	
データの準備	8	データ取込			✓
	9	データクレンジング/対象データ検討			✓
	10	データ加工（データ結合、特徴量抽出など）			✓
	11	データ傾向確認（ヒストグラム、クロス集計など）	✓	✓	
モデル構築	12	分析手法・アルゴリズムの選定		✓	
	13	モデル構築		✓	
	14	分析結果の可視化、知見抽出	✓	✓	
	15	モデルチューニング	✓	✓	
評価	16	分析結果の評価（統計観点/ビジネス観点）	✓	✓	
	17	課題抽出（ビジネス観点/システム観点）	✓		✓
展開	18	業務への導入計画策定/システム化検討	✓		✓
	19	運用定着化支援/他工程・他拠点展開検討	✓		

## 分析に必要なビジネス力

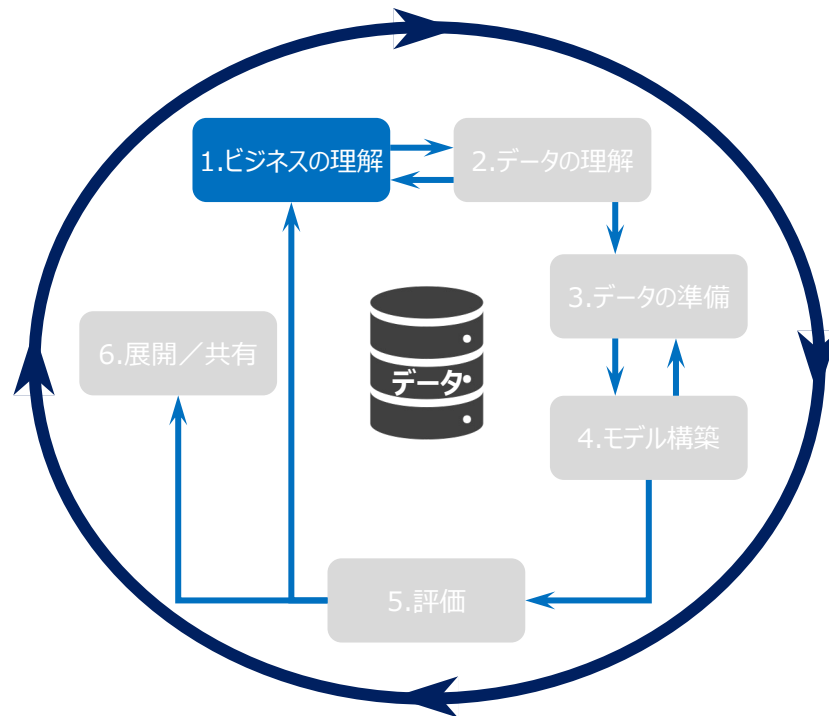
- ✓ 分析テーマの明確化
- ✓ 仮説思考力の重要性

# データ分析の進め方

- データ分析の進め方に関する方法論「CRISP-DM」に基づいて、分析と評価を繰り返して試行錯誤しながら進めるのが一般的である

## CRISP-DM: データマイニング方法論

(CRoss Industry Standard Process for Data Mining)



### 1.ビジネスの理解

- ビジネス、データマイニング目標の決定
- プロジェクトの立ち上げ

### 2.データの理解

- データの収集
- データの調査
- データ品質の検証

### 3.データの準備

- データの選択や除外
- データのクリーニング
- データの構築や統合

### 4.モデル構築

- モデリング手法の選択
- モデルの作成
- モデルの評価

### 5.評価

- データマイニングの結果の評価
- プロセスの見直し
- 実行可能なアクションリストの作成

### 6.展開/共有

- 業務への導入計画
- モニタリング、メンテナンスの計画

# 「なんかわかるでしょ？」の注意

- データ分析のきっかけは、大体ぼんやりとした依頼から始まることが多い・・・



最近ほら、AIとか流行ってるでしょ。  
うちも売上データ使ってなんかできない？  
データがあればなんかわかるよね？

売り上げを増やしたいの？  
顧客を増やしたいの？

**分析対象（目的変数）  
の明確化**

今どのくらいの貸し倒れ？  
どこまで改善目指す？

なんとか貸し倒れを減らしたいんだけど、  
データ分析して原因を見つけられる？



**現状の把握  
分析目標（ゴール）の設定**



# ビジネスの理解＝「プロセス」と「構造」の理解

- データを触る前に、分析対象ビジネスの「プロセス」と「構造」を理解しておくことが重要である
- これらと並行して、分析スコープと分析目的の検討、データの理解も行う

## 1 「ビジネスプロセス」を理解する

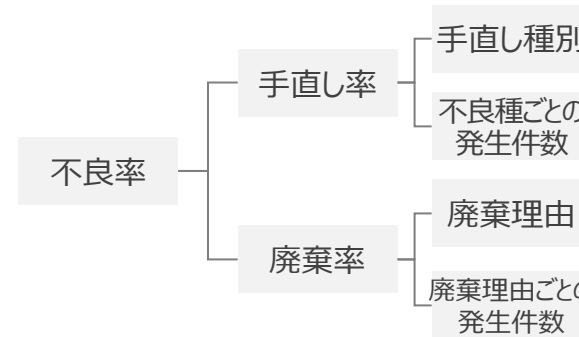
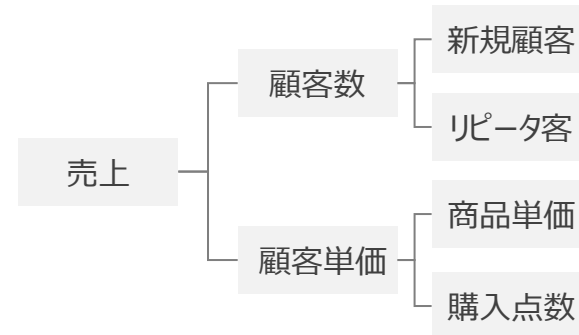
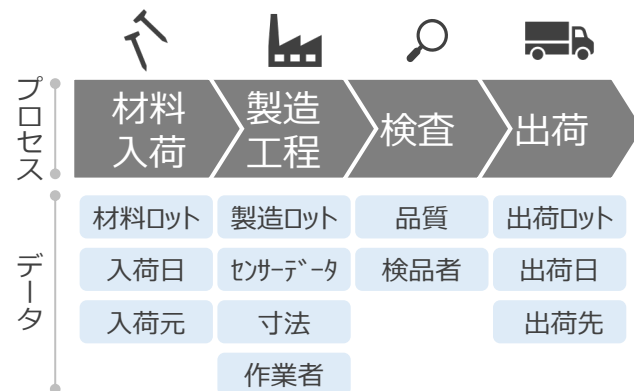
## 2 「ビジネス構造」を理解する

例

小売業

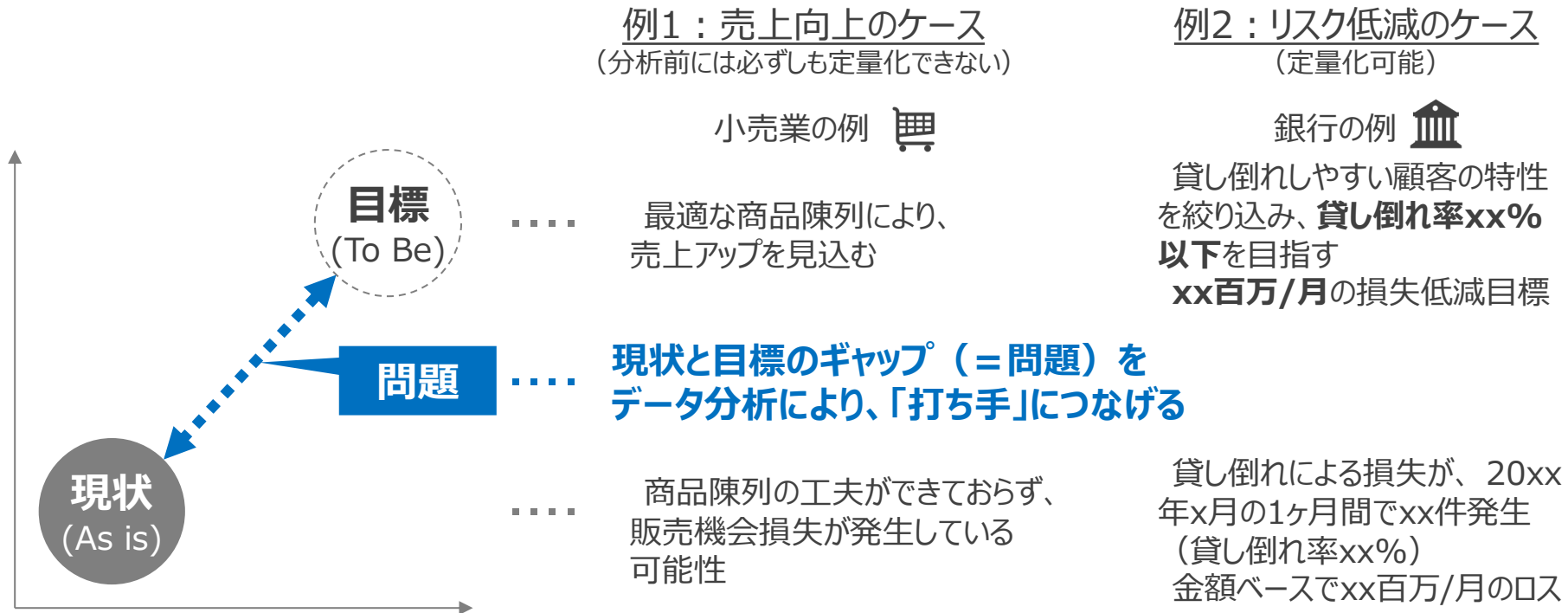


製造業



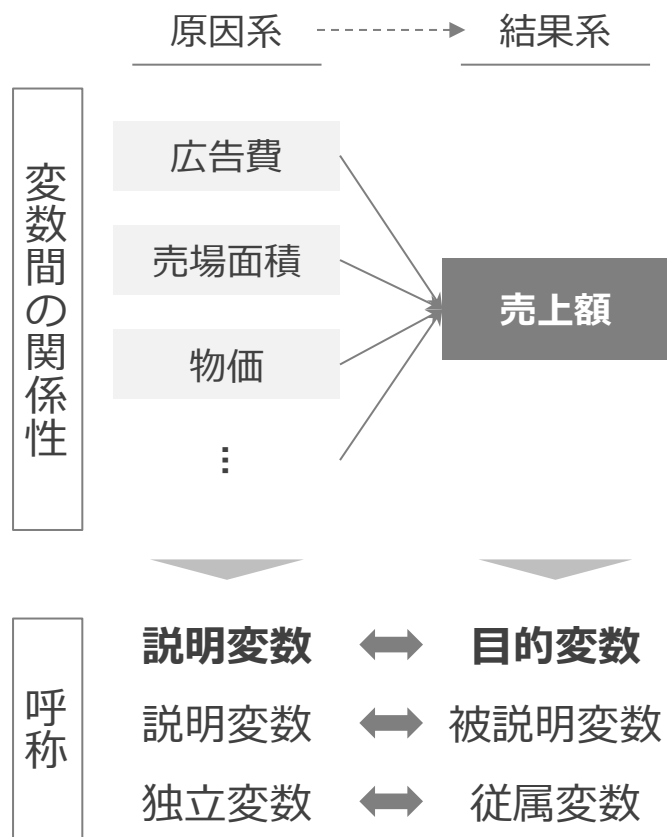
# 現状把握と目標の設定

- 理解したビジネスのプロセス・構造に基づいて、現状の把握・分析目標の設定を行ない、**データ分析により解くべき問題を明確化する**
- (特にコスト削減のケースでは) **可能な限り定量的な現状把握、目標設定を行う**



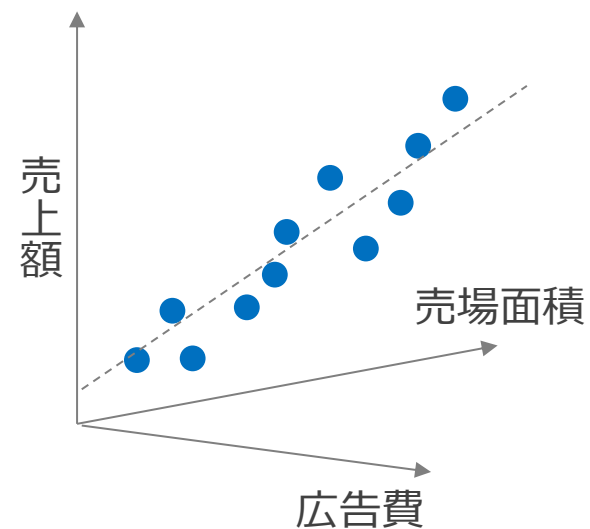
# 目的変数・説明変数とは

- 説明変数は四則演算などで必要に応じて新規変数を作る
- 目的変数は、必要に応じてフラグ値化／カテゴリ値化などの加工をする  
(分布の形を見てビンの幅や閾値を判断する)



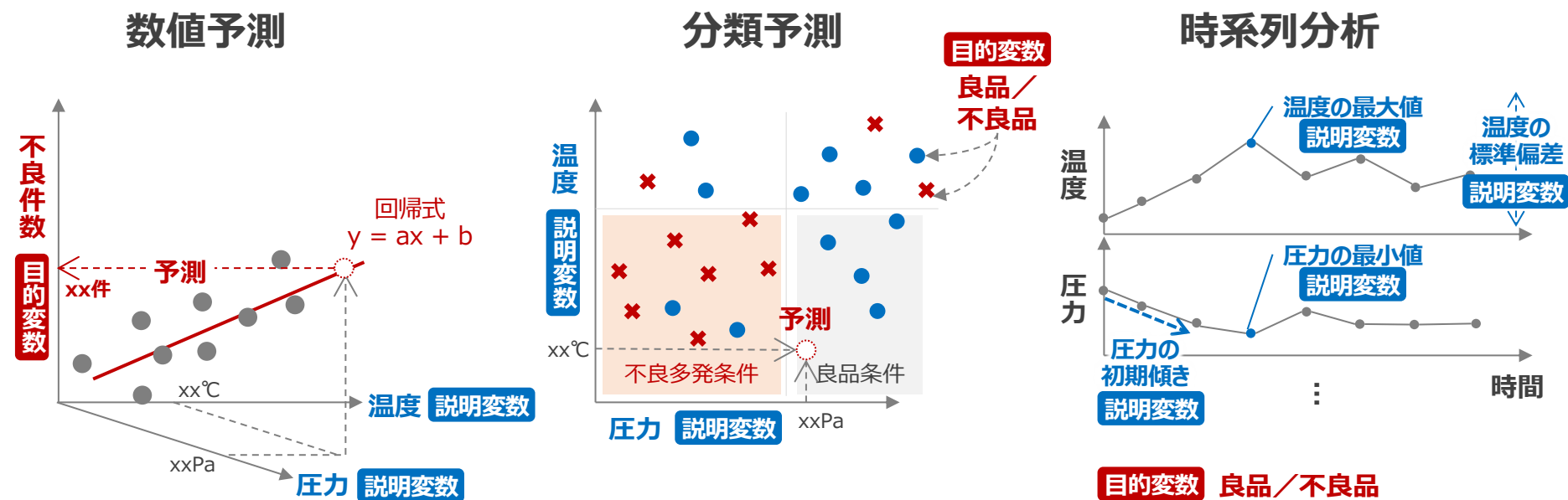
重回帰の場合：

$$\text{売上額} = a_1 \cdot \text{広告費} + a_2 \cdot \text{売場面積} + a_3 \cdot \text{物価} + \dots$$



# 目的変数・説明変数の検討

- 分析目的やデータ形式に応じて、目的変数や説明変数の形式は異なる  
「どんなデータを用いてどんな分析をしたいか」を明確にすることが重要



分析例

キーパラメータの変動予測  
機器の部品交換タイミング予測  
不良発生率の予測

良品/不良品の分類  
不良パターン (モード) の分類  
機器故障原因の分類

設備稼働の異常検知  
リアルタイム不良・不具合検出

# 統計学の基礎

- ✓ 記述統計学 と 推測統計学
- ✓ 要約統計量 (代表値 / ばらつき指標 / 順序統計量)

# 統計の落とし穴

- 以下は、日本国内でのある統計調査の結果です
- みなさん、こんな危険な食べ物、即刻禁止にすべきと思いませんか？

**心筋梗塞による死亡者の95%以上**が生前ずっとこの食べ物を食べていた

**がん患者の98%**がこの食べ物を摂取していた

**強盗や殺人などの凶悪犯の70%以上**が犯行の24時間以内にこの食べ物を口にしている

西内 啓「統計学が最強の学問である」より



**「お米」** 禁止しますか？

# 2つの「統計学」

## ■ 国勢調査と視聴率調査は何が違う？

### 国勢調査



Source: [https://www.stat.go.jp/data/kokusei/2020/kekka/pdf/summary\\_01.pdf](https://www.stat.go.jp/data/kokusei/2020/kekka/pdf/summary_01.pdf)

### 視聴率調査

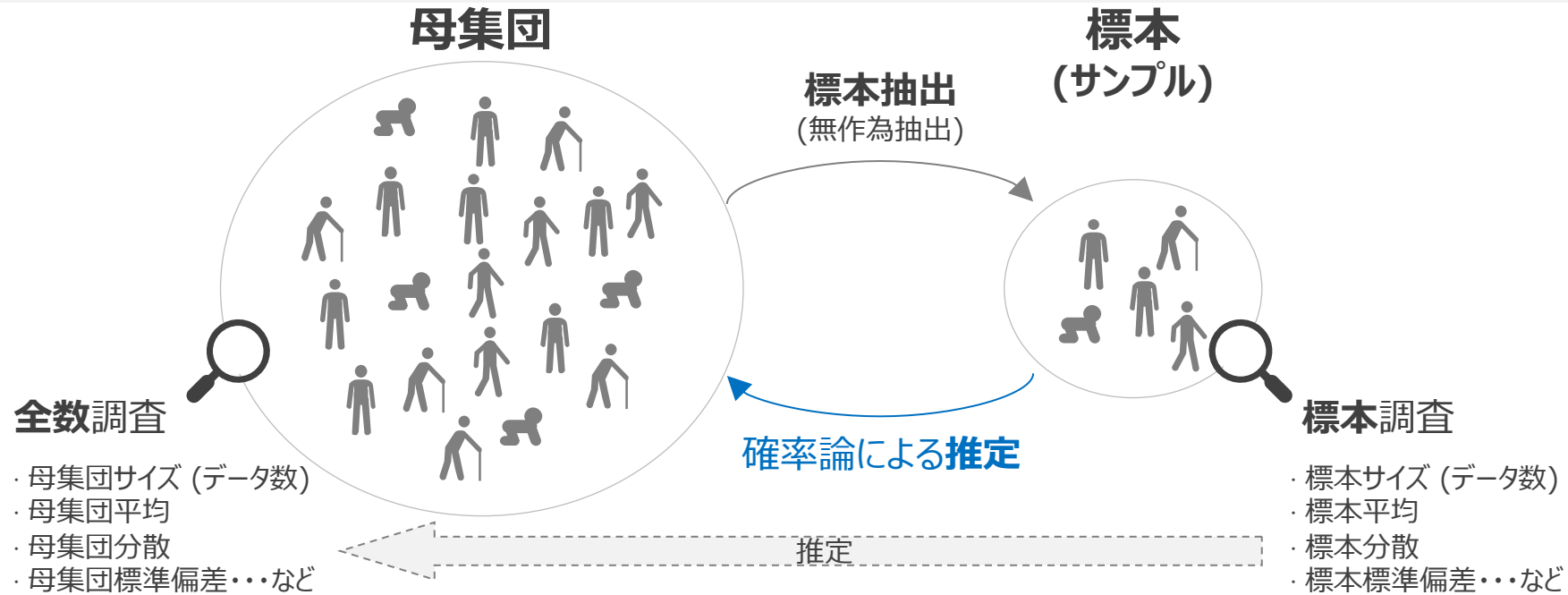
順位	番組タイトル	放送局	視聴率
1	連続テレビ小説・らんまん（7/14）	NHK総合	17.5%
2	ザワつく！金曜日（7/14）	テレビ朝日	13.1%
3	ポットンと一軒家 3時間SP（7/16）	テレビ朝日	12.3%
4	大相撲 2023年 名古屋場所・8日目（7/16 17:06～）	NHK総合	12.1%
5	日曜劇場・VIVANT（7/16）	TBS	11.5%
6	プラタモリ（7/15）	NHK総合	11.3%
7	1億人の大質問!?笑ってコラえて！ 27周年2時間SP（7/12）	日本テレビ	11.1%
8	世界まる見え！テレビ特捜部 2時間SP（7/10）	日本テレビ	10.6%
9	木曜ドラマ・ハヤブサ消防団（7/13）	テレビ朝日	10.5%
10	音楽の日2023（7/15 19:00～）	TBS	10.1%

2023年7月10日(月)～7月16日(日) 視聴率ランキング

Source: <https://www.tvguide.or.jp/tvranking/tvrank-20230728/>

# 2つの「統計学」－ 記述統計学と推測統計学

- データ全体を調べてデータの特徴や傾向を把握する「記述統計学」と、一部のサンプルを抽出して調べ、全体の特徴や傾向を推測する「推測統計学」がある



## 記述統計学

概要

全数調査により観測されたデータ全体の特徴を調べたり、データを要約するアプローチのこと

例

国勢調査、(学校の)健康診断・定期テスト など

金融では... (CRM内の) 顧客データの分析 など

## 推測統計学

データ全体 (母集団) を調べあげるのが困難な場合に、抽出したデータの一部 (標本) の特徴から、データ全体の特徴を確率的に推測するアプローチのこと

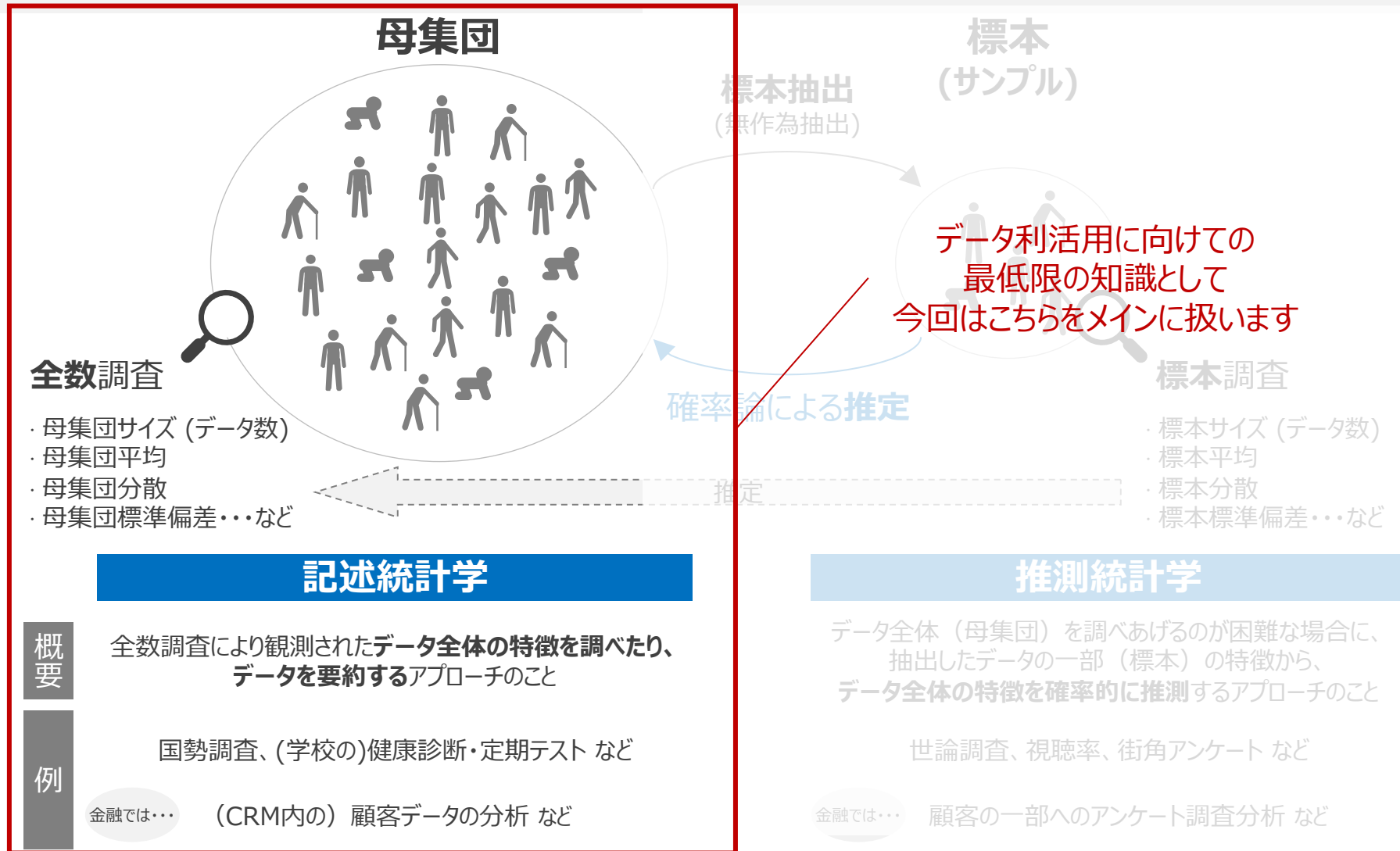
世論調査、視聴率、街角アンケート など

金融では... 顧客の一部へのアンケート調査分析 など



# 2つの「統計学」－ 記述統計学と推測統計学

- データ全体を調べてデータの特徴や傾向を把握する「記述統計学」と、一部のサンプルを抽出して調べ、全体の特徴や傾向を推測する「推測統計学」がある

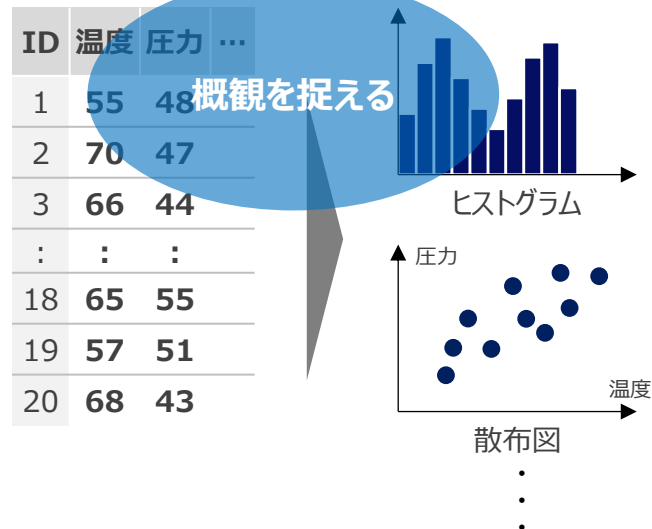


# 記述統計学の基本：データの特徴を捉える

- 個々のデータをくまなく見るのは難しいため、グラフ（ヒストグラムや散布図）や要約統計量（平均値や標準偏差）を用いて全体傾向を把握する

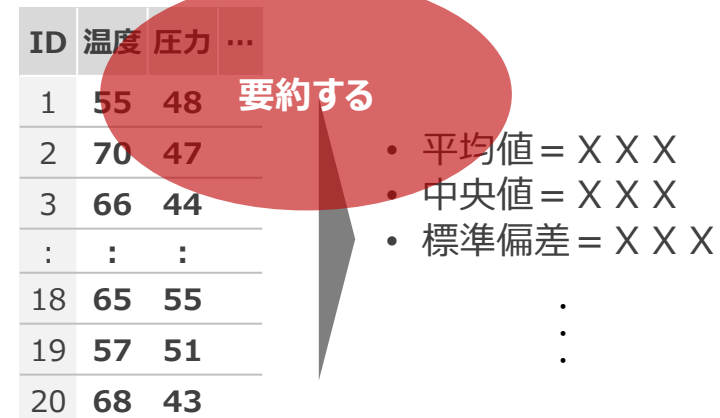
## 可視化（グラフ化）

視覚的にデータの特徴や傾向を把握



## 数値化（データ要約）

データの特徴を示す値に要約し  
比較可能な客観的傾向を掴む  
(要約統計量)

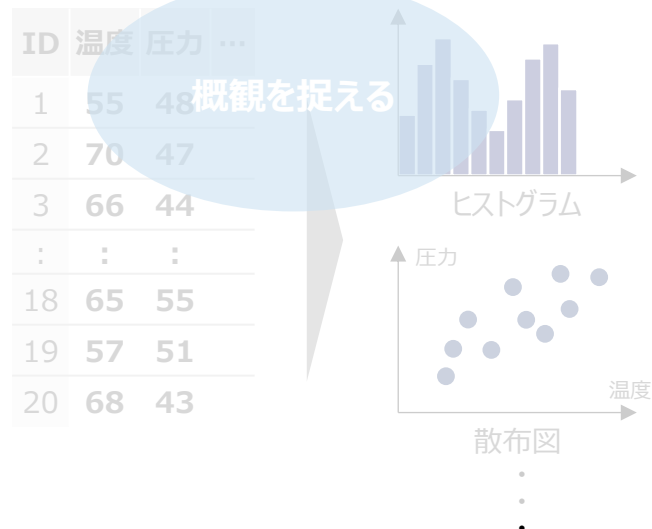


# 記述統計学の基本：データの特徴を捉える

- 個々のデータをくまなく見るのは難しいため、グラフ（ヒストグラムや散布図）や要約統計量（平均値や標準偏差）を用いて全体傾向を把握する

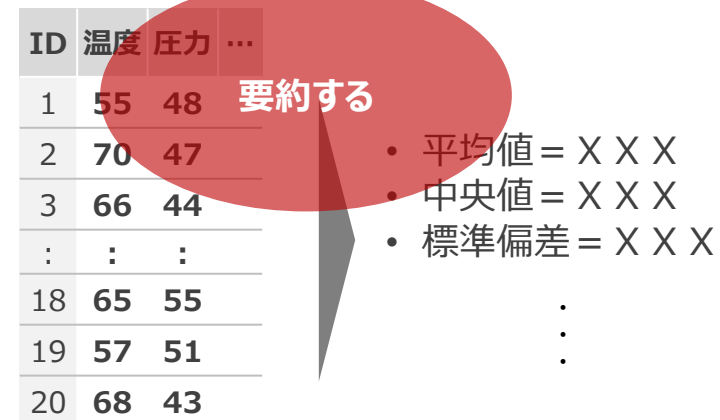
## 可視化（グラフ化）

視覚的にデータの特徴や傾向を把握



## 数値化（データ要約）

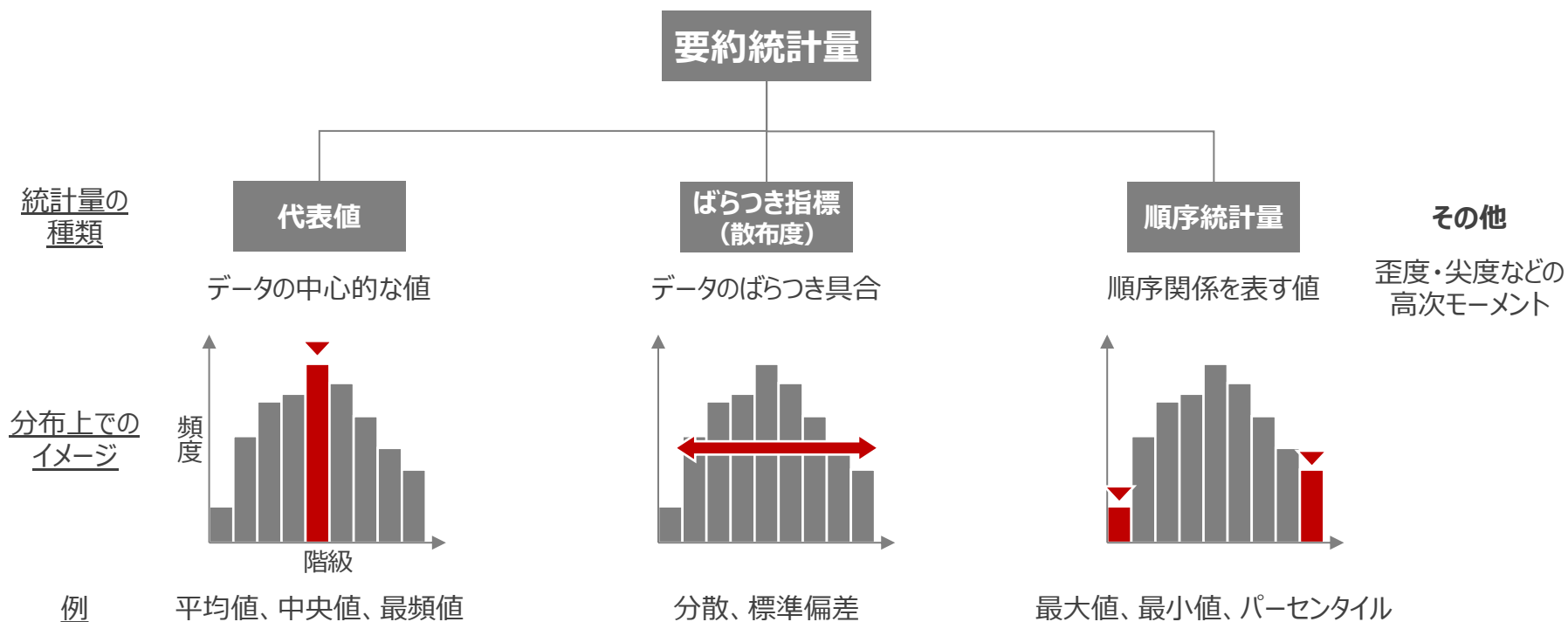
データの特徴を示す値に要約し  
比較可能な客観的傾向を掴む  
(要約統計量)



※可視化については、別途「データの理解（可視化）」のプロセスで取り扱う

# 数値化 (データ要約)

- 全データをくまなく見るのは難しいため、「要約した値 (=要約統計量)」を用いて全体傾向を把握するのが一般的



# 代表値（平均値）

- 平均値は最もよく用いられる代表値であるが、分布の形によっては思わぬ落とし穴がある

ある会社の社員9名の年収

社員	年収
A	200万円
B	200万円
C	200万円
D	400万円
E	400万円
F	500万円
G	700万円
H	900万円
I	1億円

平均値の算出

$$\begin{aligned}\text{平均値} &= \frac{200\text{万円} + 200\text{万円} + \dots + 1\text{億円}}{9\text{名}} \\ &= \mathbf{1,500\text{万円}}\end{aligned}$$

この平均値、直感と合っている？



# 代表値（中央値・最頻値）

- 中央値：データを値の小さい順に並べたときに中央に位置する値
- 最頻値：データの中で最も多く出現する値。複数存在するケースもある

ある会社の社員9名の年収

中央値と最頻値

社員	年収
A	200万円
B	200万円
C	200万円
D	400万円
E	400万円
F	500万円
G	700万円
H	900万円
I	1億円

最頻値 = 200万円

社内で最も多い年収値

中央値 = 400万円

社内の年収を小さい順に並べ、  
ちょうど真ん中に位置する年収値

**外れ値**

**中央値・最頻値は外れ値の影響を受けにくい！**

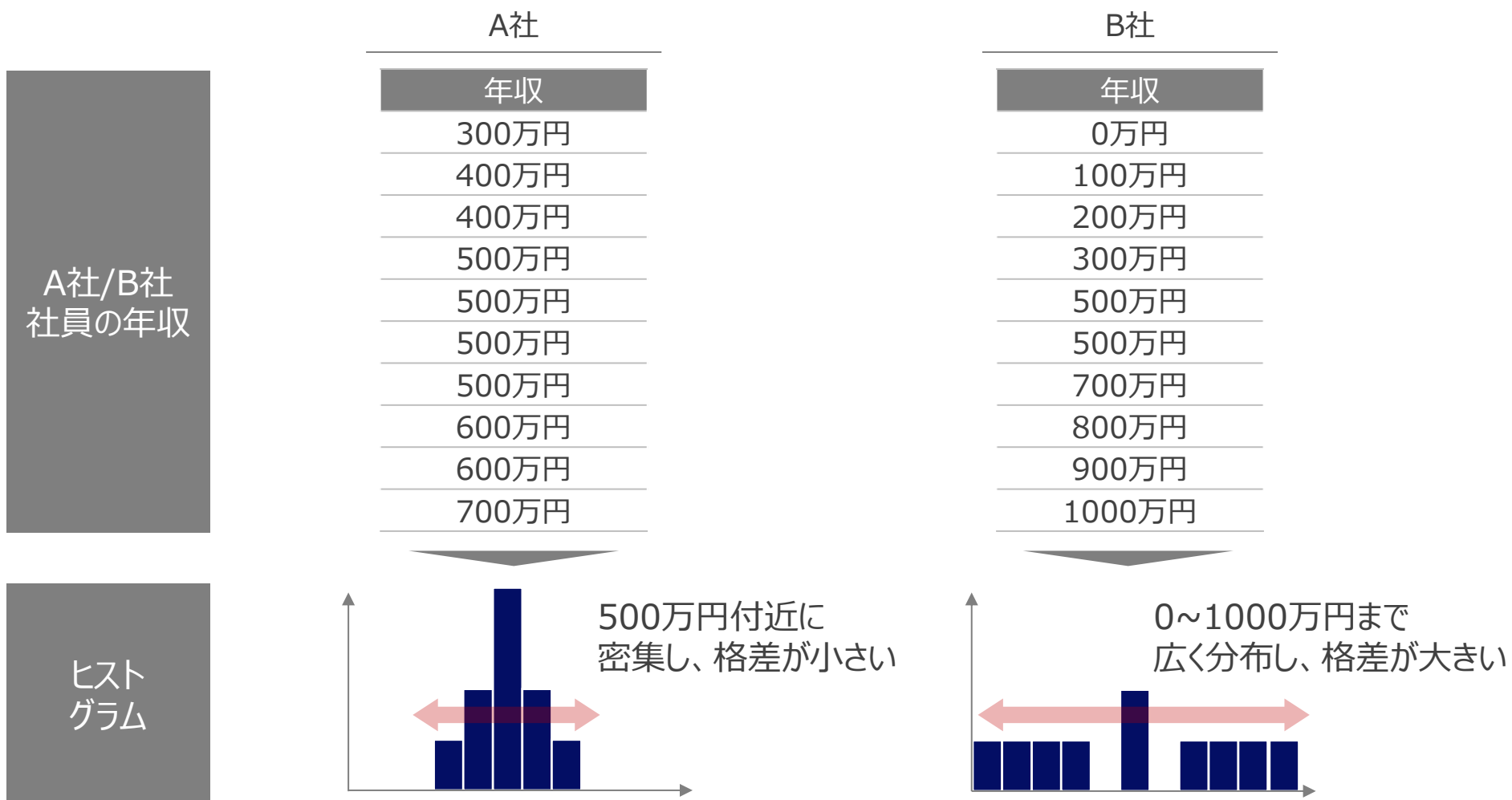
# 代表値だけではつかめない特徴

- 代表値はすべて同じ。どちらの会社に入っても待遇は同じ！と言ってよいのか？

		A社	B社
A社/B社 社員の年収		年収	年収
		300万円	0万円
		400万円	100万円
		400万円	200万円
		500万円	300万円
		500万円	500万円
		500万円	500万円
		500万円	700万円
		600万円	800万円
		600万円	900万円
		700万円	1000万円
代表値	平均値	500万円	500万円
	中央値	500万円	500万円
	最頻値	500万円	500万円

# 代表値だけではつかめない特徴

- ヒストグラムで見ると、社員によって年収が大きく異なっていることがわかる





# データの「ばらつき」の表し方

- データのばらつきは、「平均からのズレ」= 偏差 として捉えるのが基本的な考え方

	A社		B社	
A社/B社 社員の年収	年収	平均値との差	年収	平均値との差
	300万円	-200万円	0万円	-500万円
	400万円	-100万円	100万円	-400万円
	400万円	-100万円	200万円	-300万円
	500万円	0万円	300万円	-200万円
	500万円	0万円	500万円	0万円
	500万円	0万円	500万円	0万円
	500万円	0万円	700万円	+200万円
	600万円	+100万円	800万円	+300万円
	600万円	+100万円	900万円	+400万円
	700万円	+200万円	1000万円	+500万円
	平均年収 500万円		平均年収 500万円	

# データの「ばらつき」の計算

- 分散：平均からのズレを、（正負の符号を消すために）2乗して平均した値
- 標準偏差：元のデータと単位を揃えるために、分散の値を変換した値

		A社		B社	
A社/B社 社員の数	年収	300万円	平均値との差 -200万円	0万円	平均値との差 -500万円
		400万円	-100万円	100万円	-400万円
		400万円	-100万円	200万円	-300万円
		500万円	0万円	300万円	-200万円
		500万円	0万円	500万円	0万円
		500万円	0万円	500万円	0万円
		500万円	0万円	500万円	0万円

A社年収の**分散** = 
$$\frac{(-200万円)^2 + (-100万円)^2 + \dots + (100万円)^2 + (200万円)^2}{10名}$$

= A社の年収  
標準偏差 × A社の年収  
標準偏差

実際に標準偏差を計算すると、A社は109.5万円、B社は328.6万円

# 代表的な要約統計量まとめ

- どれか一つを算出すれば良いわけではなく、様々な観点でデータの概観を把握することが重要
- また、時系列データの分析ではこれら統計量自体を「特徴量」として扱うことも多い

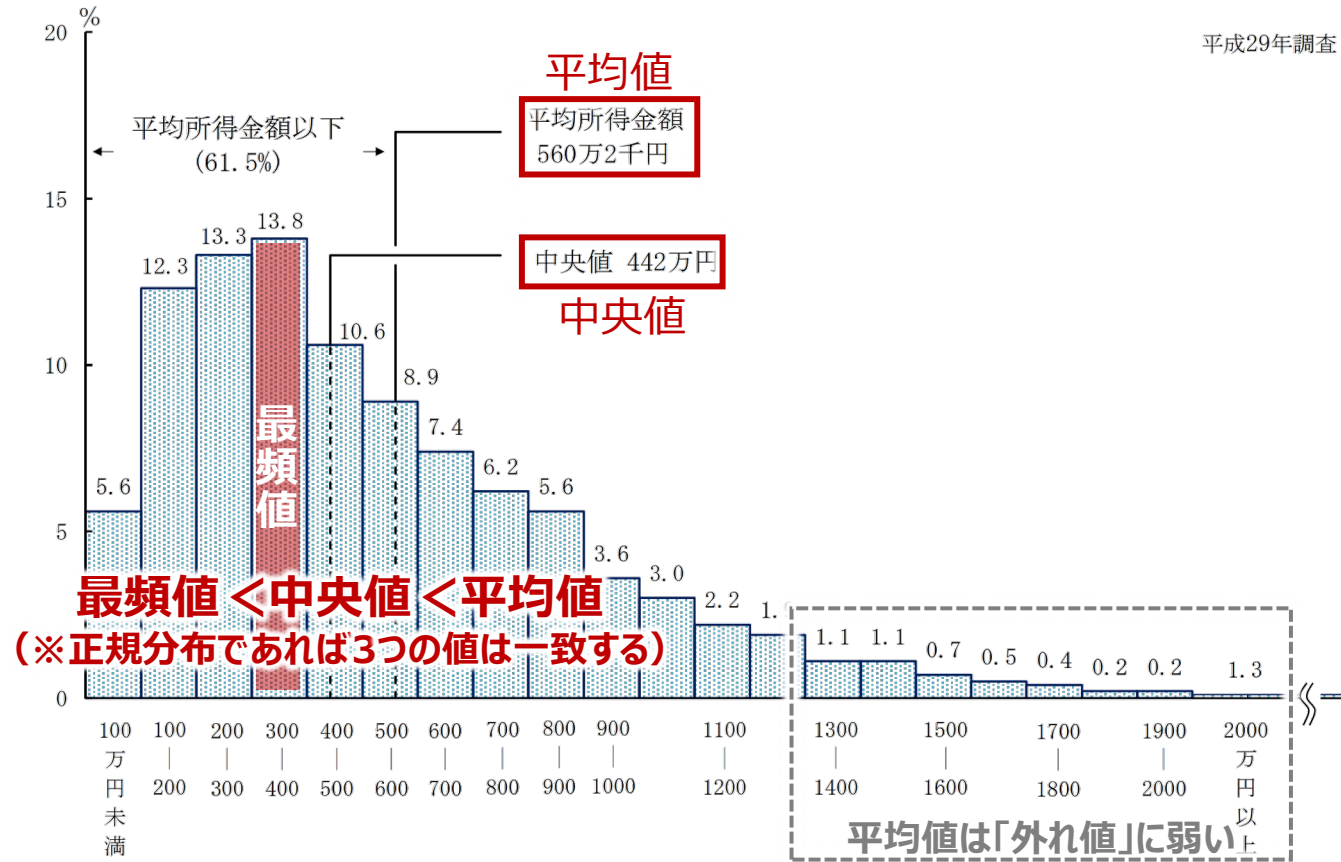
分類	統計量	説明
代表する値 (代表値)	平均	データの中心位置を表す値
	中央値 (メディアン)	データを小さい順に並べたとき、ちょうど中央の位置にくる値
	最頻値 (モード)	データの中で最頻出する値
ばらつき具合 を示す値 (散布度)	分散	データのばらつき具合を表す値
	標準偏差	データのばらつき具合を表す値 元データと <b>同じ単位</b> のため直感的なばらつき把握が可能
順序を表す値 (順序統計量)	最大値 / 最小値	データの中で最も大きい / 小さい値
	第1 / 第3四分位数	データを小さい順に並べたとき、 それぞれ25% / 75% の位置にくる値



# 参考：年収分布における代表値の差異

- 正規分布ではない分布では、平均値だけを採用するとデータの特徴を見誤る恐れがある

図9 所得金額階級別世帯数の相対度数分布



出典：厚生労働省 平成 29 年 国民生活基礎調査の概況



# Google Colaboratory上での レクチャー&演習

Part 2 : Google Colab.リンク

[https://colab.research.google.com/drive/1uTxPMJjfDWqGMbEJZsOYM\\_W3T4eJgfJh?usp=sharing](https://colab.research.google.com/drive/1uTxPMJjfDWqGMbEJZsOYM_W3T4eJgfJh?usp=sharing)

## データ観察の基本

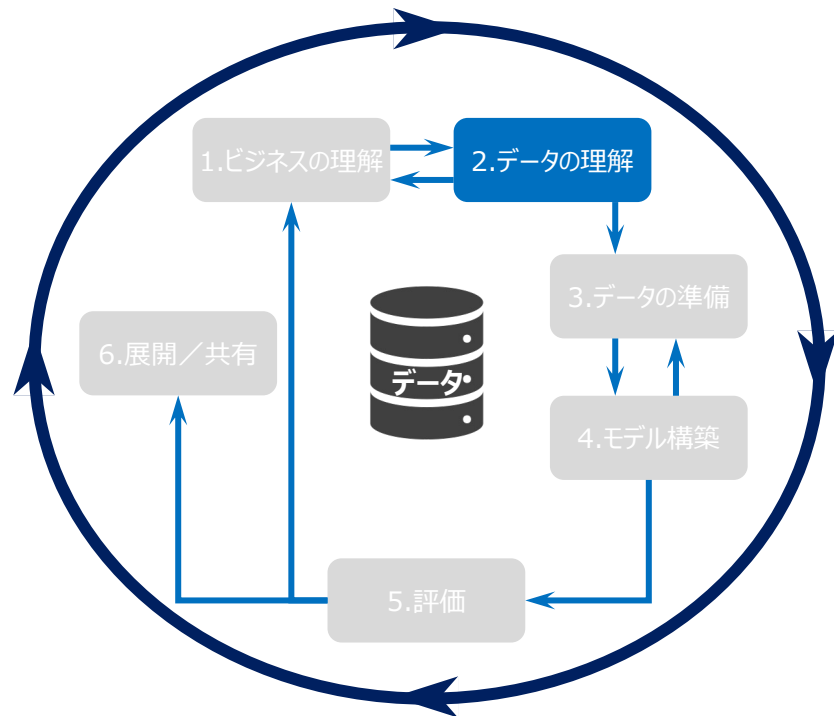
- ✓ グラフと要約統計量によるデータの傾向把握
- ✓ 相関分析によるデータ間の関係性把握
- ✓ 相関と因果の違い

# データ分析の進め方

- データ分析の進め方に関する方法論「CRISP-DM」に基づいて、分析と評価を繰り返して試行錯誤しながら進めるのが一般的である

## CRISP-DM: データマイニング方法論

(CRoss Industry Standard Process for Data Mining)



### 1. ビジネスの理解

- ビジネス、データマイニング目標の決定
- プロジェクトの立ち上げ

### 2. データの理解

- データの収集
- データの調査
- データ品質の検証

### 3. データの準備

- データの選択や除外
- データのクリーニング
- データの構築や統合

### 4. モデル構築

- モデリング手法の選択
- モデルの作成
- モデルの評価

### 5. 評価

- データマイニングの結果の評価
- プロセスの見直し
- 実行可能なアクションリストの作成

### 6. 展開/共有

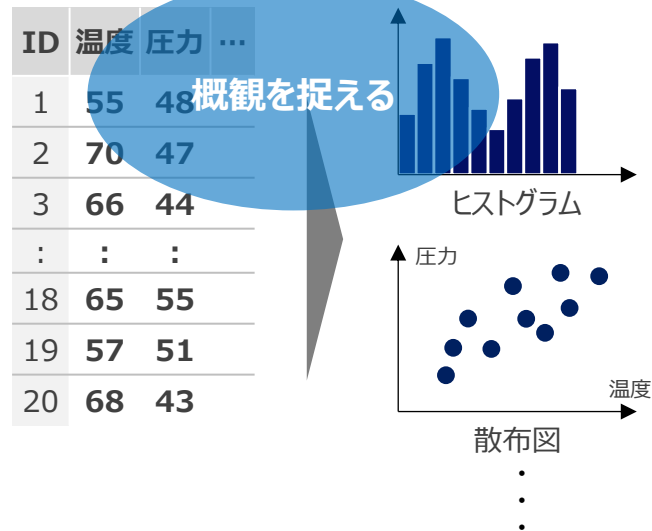
- 業務への導入計画
- モニタリング、メンテナンスの計画

# 記述統計学の基本：データの特徴を捉える

- 個々のデータをくまなく見るのは難しいため、グラフ（ヒストグラムや散布図）や要約統計量（平均値や標準偏差）を用いて全体傾向を把握する

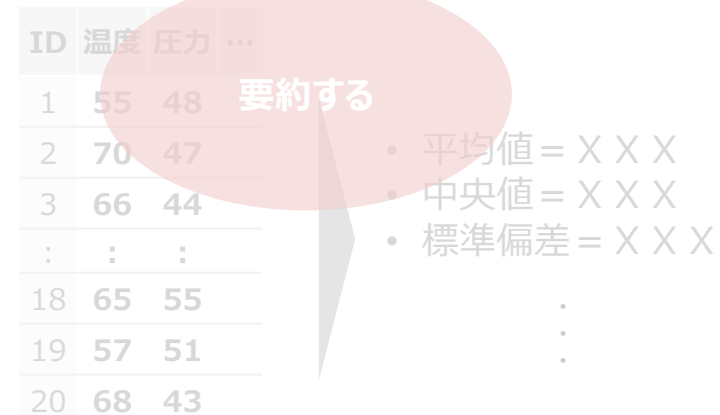
## 可視化（グラフ化）

視覚的にデータの特徴や傾向を把握



## 数値化（データ要約）

データの特徴を示す値に要約し  
比較可能な客観的傾向を掴む  
(要約統計量)





# 集団の比較

- あなたは、「飲むログ」でデートで使うお店を選んでいきます。候補となったお店は2つ。でもレビュー点数が全く同じです。
- 特徴を捉えるにはどうすれば良い？

## レストランA

★★★★☆ 3.45

#	評点	#	評点
1	3.5	11	3.0
2	3.5	12	3.5
3	4.5	13	3.0
4	3.5	14	4.0
5	3.5	15	4.0
6	2.5	16	4.0
7	3.0	17	2.5
8	3.5	18	4.5
9	4.0	19	1.0
10	3.0	20	5.0



どちらも同じ  
評価だから  
どちらでも  
いいのかな...

## レストランB

★★★★☆ 3.45

#	評点	#	評点
1	3.0	11	1.0
2	4.5	12	2.0
3	5.0	13	4.5
4	4.5	14	2.0
5	2.5	15	5.0
6	2.5	16	4.5
7	4.0	17	4.5
8	2.0	18	4.0
9	5.0	19	2.5
10	2.0	20	4.0

# ヒストグラム（度数分布図）

- ヒストグラムは、データを一定間隔（**階級**）ごとに頻度集計（**度数**）をとり、横軸に階級、縦軸に度数をとって柱状（**ビン**）の集合で表したグラフ
- データの分布状況を可視化して、直感的にデータの特徴を捉えやすくする目的がある

## 元データ

レストランAの  
レビューアごとの評点データ

#	評点	#	評点
1	3.5	11	3.0
2	3.5	12	3.5
3	4.5	13	3.0
4	3.5	14	4.0
5	3.5	15	4.0
6	2.5	16	4.0
7	3.0	17	2.5
8	3.5	18	4.5
9	4.0	19	1.0
10	3.0	20	5.0

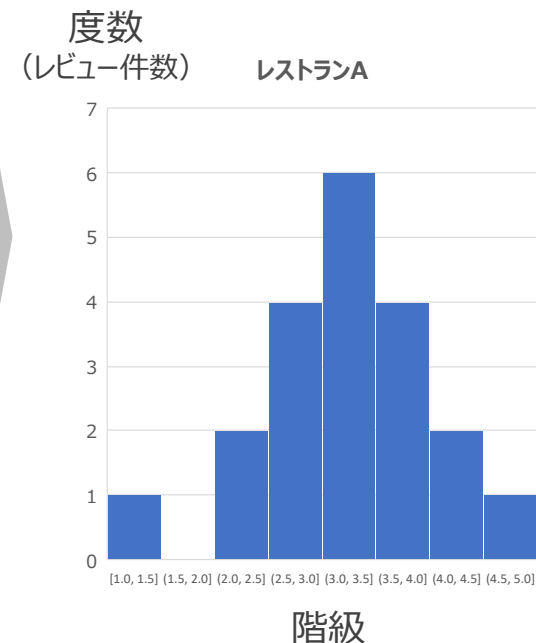
## 度数分布表

一定間隔（階級）ごとに  
レビュー件数（度数）を集計

階級	度数
0.0~0.5	0
0.5~1.0	0
1.0~1.5	1
1.5~2.0	0
2.0~2.5	2
2.5~3.0	4
3.0~3.5	6
3.5~4.0	4
4.0~4.5	2
4.5~5.0	1

## ヒストグラム（度数分布図）

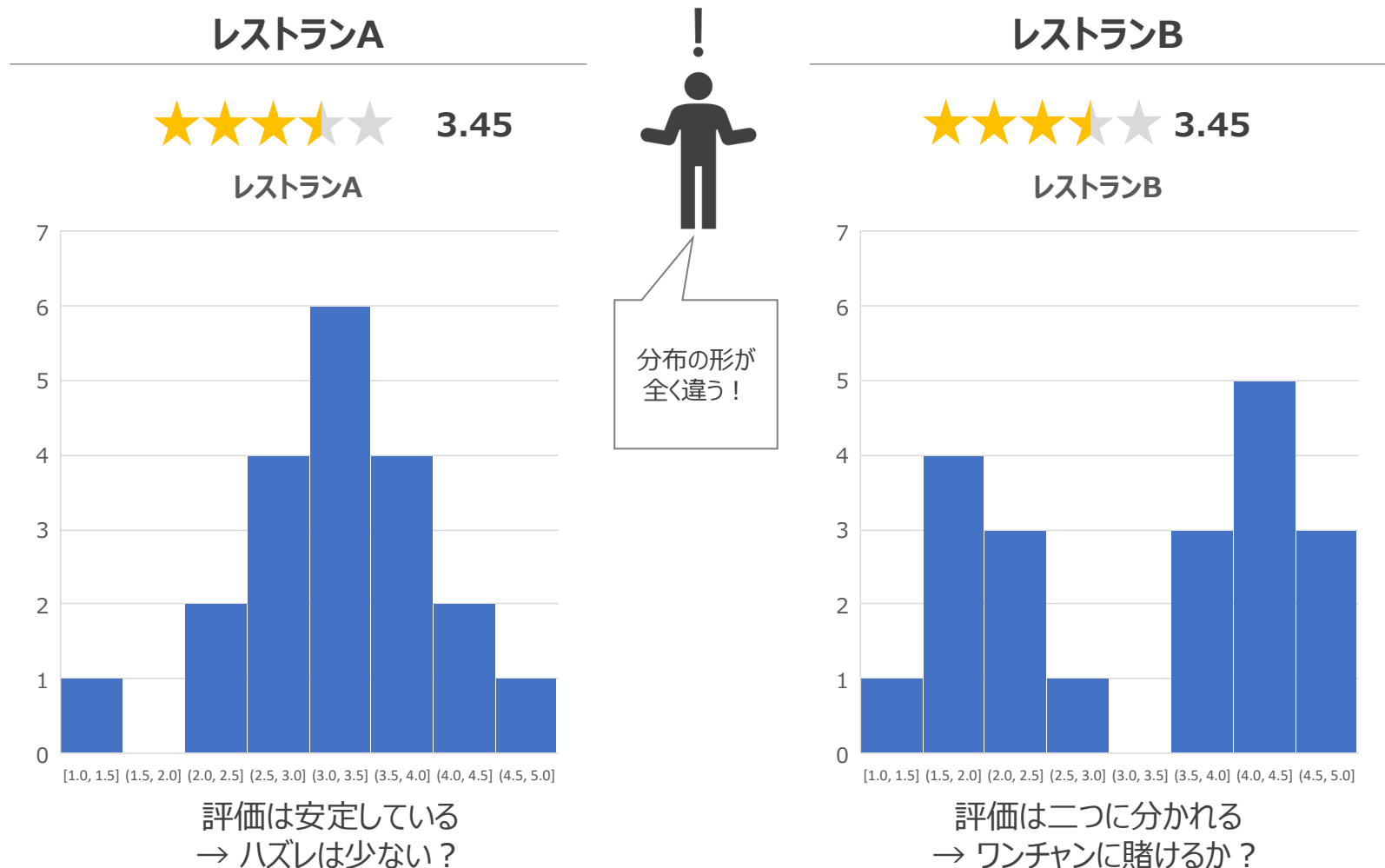
横軸に階級、縦軸に度数を取り、  
柱（ビン）の集合で可視化



※境界値は小さい側の階級に含める  
(ただし、1.0は便宜上1.0~1.5に含める)

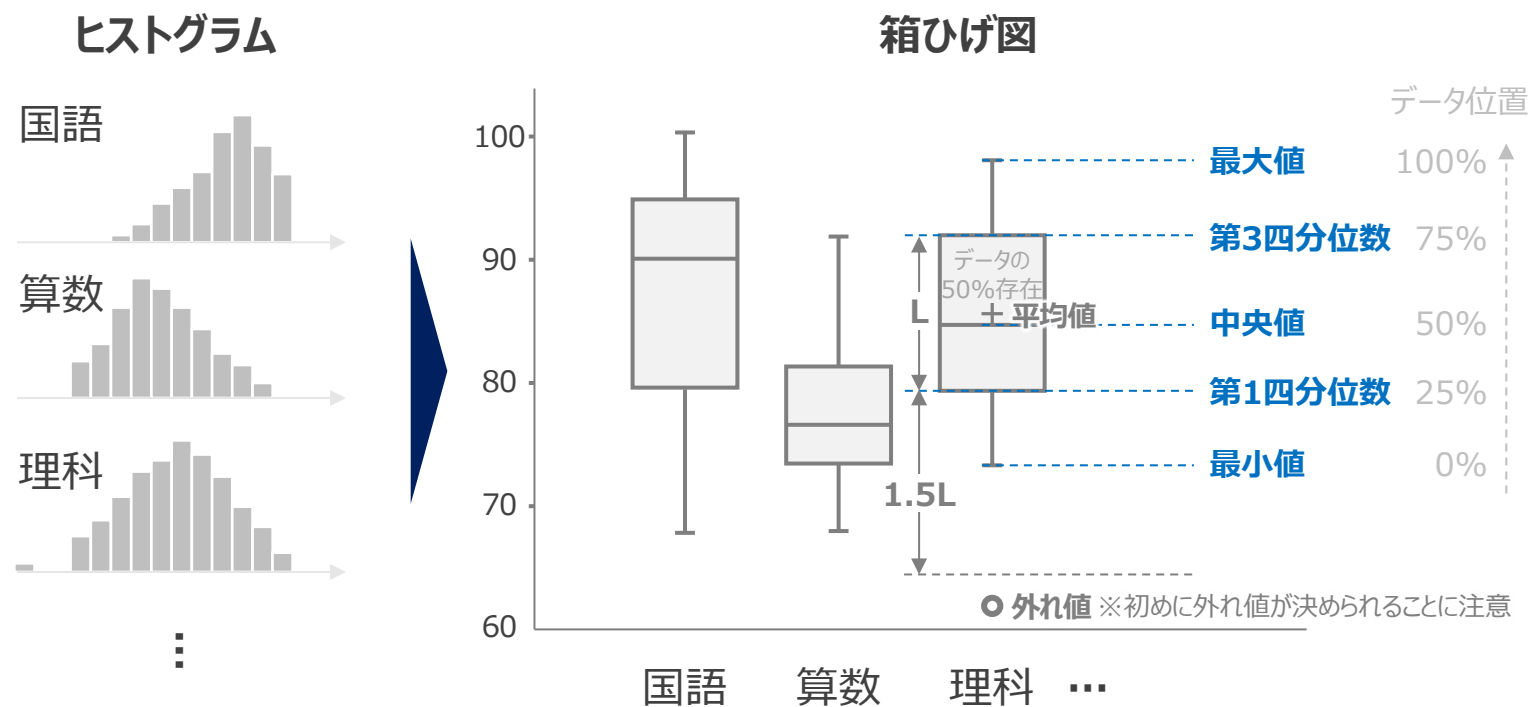
# ヒストグラムの描画と比較

- ヒストグラムを描くと、「分布の形」（≡データの散らばり具合）が可視化される



# 箱ひげ図

- 分布形状の把握は極めて重要であるが、主にヒストグラムは1変数、散布図は2変数での分布可視化に適している。一方、箱ひげ図は、分布の概形を**多変数間で比較**するのに適している

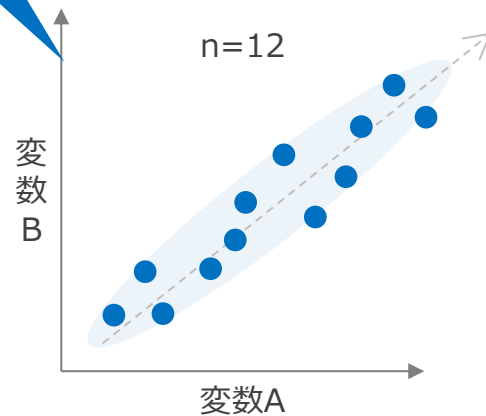


# 散布図と相関

- 二つの値（**変数**）間において、一方が上がれば他方も上がる（or 下がる）ような関係性のことを「**相関**」と呼ぶ
- 各変数を各軸にとってグラフ化した「**散布図**」を描くことで、相関関係の視覚的な把握が可能

散布図

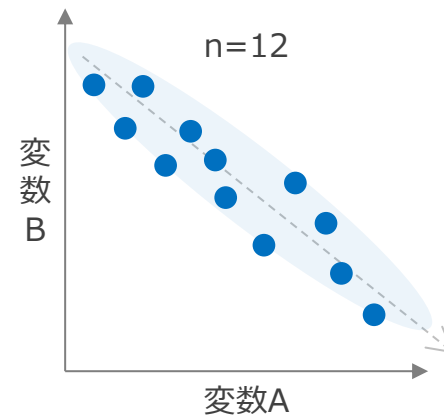
## 正の相関



一方 (A) が増加すると  
他方 (B) も増加  
(右上がりの楕円状)

例) A・・・(夏の)平均気温  
B・・・エアコンの販売台数

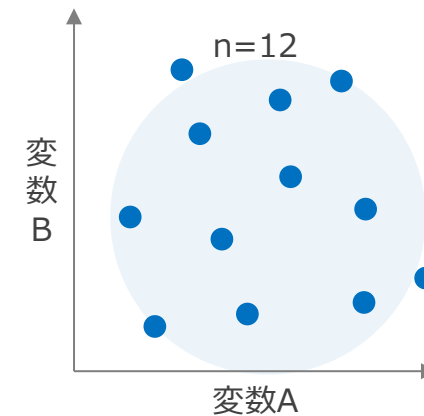
## 負の相関



一方 (A) が増加すると  
他方 (B) は減少  
(右下がりの楕円状)

例) A・・・(冬の)平均気温  
B・・・ストーブの販売台数

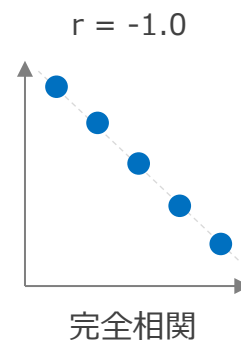
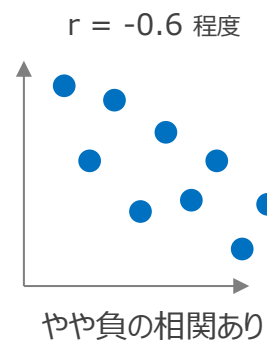
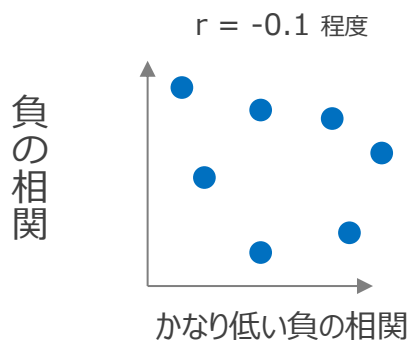
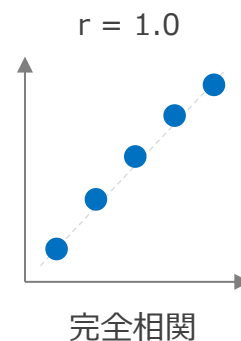
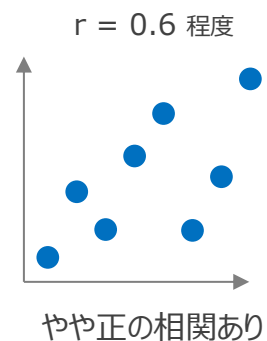
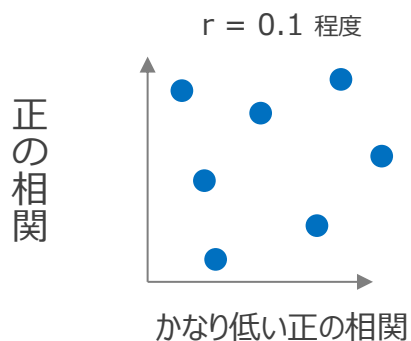
## 無相関



双方の変数間に  
関連性なし  
(まばらな円状)

# 相関係数

- 相関係数 $r$  (correlation coefficient) とは、2つの変数間の相関の度合いを表す指標
- $-1 \leq r \leq 1$  の値を取り、正の場合は正相関、負の場合は負相関、0の場合は無相関



(参考) 相関係数の目安

相関係数 (絶対値)	相関の強弱
1.0	強い相関
0.7	やや相関あり
0.4	弱い相関
0.2	ほぼ相関なし
0	

# 相関行列による網羅的な変数間の関係性把握

- 事前に各変数同士の相関係数を総当たりで調べておくと、後々の結果解釈に役立つ（**相関行列**）
- また、**共線性が高い変数**（相関の高い）が複数混ざっていると、その変数の影響を強く受け、偏った分析結果になることがある。この場合、共線性が高い変数は除外することが有効

	打率	試合	打席数	打数	安打	本塁打	打点	盗
打率	1.00							
試合	-0.22	1.00						
打席数	0.02	0.85	1.00					
打数	-0.03	0.83	0.95	1.00				
安打	0.56	0.58	0.80	0.81	1.00			
本塁打	0.00	0.34	0.48	0.38	0.30	1.00		
打点	0.07	0.40	0.48	0.39	0.36	0.92	1.00	
盗塁	0.13	0.28	0.44	0.38	0.42	0.12	-0.17	1.00
四球							0.60	
死球	0.20	0.14			0.11	-0.16	-0.02	
三振	-0.34	0.44	0.47	0.41	0.13	0.63	0.57	
犠打	-0.33	0.12	0.05	0.09	-0.11	-0.49	-0.54	
併殺打	-0.06	0.08	0.10	0.12	0.05	0.32	0.46	
出塁率	0.73	-0.10	0.10	-0.14	0.32	0.31	0.34	
長打率	0.38	0.11	0.32	0.19	0.38	0.90	0.86	
OPS	0.55	0.04	0.27	0.09	0.40	0.77	0.76	
RC27	0.63	0.04	0.28	0.08	0.43	0.65	0.64	
XR27	0.59	0.07	0.30	0.09	0.42	0.70	0.67	

各変数同士の総当たり形式で相関係数を算出  
(相関行列)

# 散布図行列と色分け可視化

- 各変数を個別に見るだけでなく、「**散布図行列**」として網羅的に各変数の関係性を俯瞰したり、**目的変数に応じて色分け**することで、影響因子を仮説立てできることもある

目的変数で色分けした  
散布図行列の例







Google Colaboratory上での  
レクチャー&演習

# 変数の尺度 (名義尺度・順序尺度・間隔尺度・比例尺度)

- 変数の種類は大きく「**質的データ (カテゴリーデータ)**」と「**量的データ (数量データ)**」に分けられ、それぞれの特性に合わせて扱う必要がある
- 数量データは先述の記述統計量、グラフが有効であるが、カテゴリーデータはクロス集計が有効

種類	変数の尺度	概要	データの例	扱い方		
				大小 (A<B)	差分 (A-B)	比率 (A/B)
質的データ (カテゴリーデータ)	名義尺度	単にデータを区別するための分類ラベル。 演算不可で、順序も意味をなさない	性別、血液型、顧客ID 作業員、個品ID、 良品/不良品	-	-	-
	順序尺度	<b>順序</b> (大小関係) にのみ意味がある尺度。 したがって、平均値は意味を持たないが、順序統計量 (最大・最小など) は算出可能	顧客満足度、震度 書道の「段」や検定の「級」 不良レベル、工程順序	●	-	-
量的データ (数量データ)	間隔尺度	数値演算可能だが、 <b>値の差</b> のみに意味がある尺度。 0はあくまで相対的な位置関係でしかない	年齢、西暦、偏差値 温度 (°C)、製造日時	●	●	-
	比例尺度	数値演算可能で、値の差に加え、 <b>値の比</b> にも意味がある尺度。 0が「何もない」という絶対的な意味を持つ	身長、売上金額 寸法、圧力、作業時間、 絶対温度	●	●	●

クロス集計 (次頁) が有効!

※集計によるカウントのみ可能

平均値に意味あり

# クロス集計とは

- カテゴリーデータは、平均値や標準偏差などの記述統計量を算出してもあまり意味がないため、当該データ項目に登場する**カテゴリごとの頻度集計を算出**することが多い
- 特に、複数のデータ項目を掛合せて算出（例えば年代×性別）した頻度集計表を**クロス集計**と呼ぶ

## ▼ある血液検査の結果

年代	性別	血清ナトリウム	...
70	男性	136	...
50	男性	136	...
60	男性	129	...
50	男性	137	...
60	女性	116	...
90	男性	132	...
70	男性	137	...
60	男性	131	...
60	女性	138	...
80	男性	133	...
70	男性	131	...
60	男性	140	...
40	男性	137	...
50	男性	137	...
40	女性	138	...
80	男性	136	...
80	男性	140	...
...	...	...	...

カテゴリーデータのため  
平均値などの統計量  
に意味を持たない

## ▼性別ごとの頻度集計

性別	人数
男性	194
女性	105
合計	299

## ▼年代ごとの頻度集計

年代	人数
40	47
50	82
60	93
70	52
80	19
90	6
合計	299

複数データ項目を  
掛け合わせる

## ▼性別×年代のクロス集計表

年代	性別		人数 合計
	男性	女性	
40	27	20	47
50	56	26	82
60	58	35	93
70	34	18	52
80	15	4	19
90	4	2	6
合計	194	105	299

## ▼応用編：他項目の「平均値」を表示

年代	性別		血清ナトリウム 平均値
	男性	女性	
40	137.1	135.6	136.4
50	137.3	136.9	137.1
60	136.0	137.1	136.4
70	136.7	136.6	136.7
80	135.5	138.0	136.1
90	133.0	141.0	135.7
血清ナトリウム 平均値	136.5	136.8	136.6

Excelの  
ピボットテーブル  
のようなイメージ



Google Colaboratory上での  
レクチャー&演習

# 参考：「データの理解」でよく用いる関数

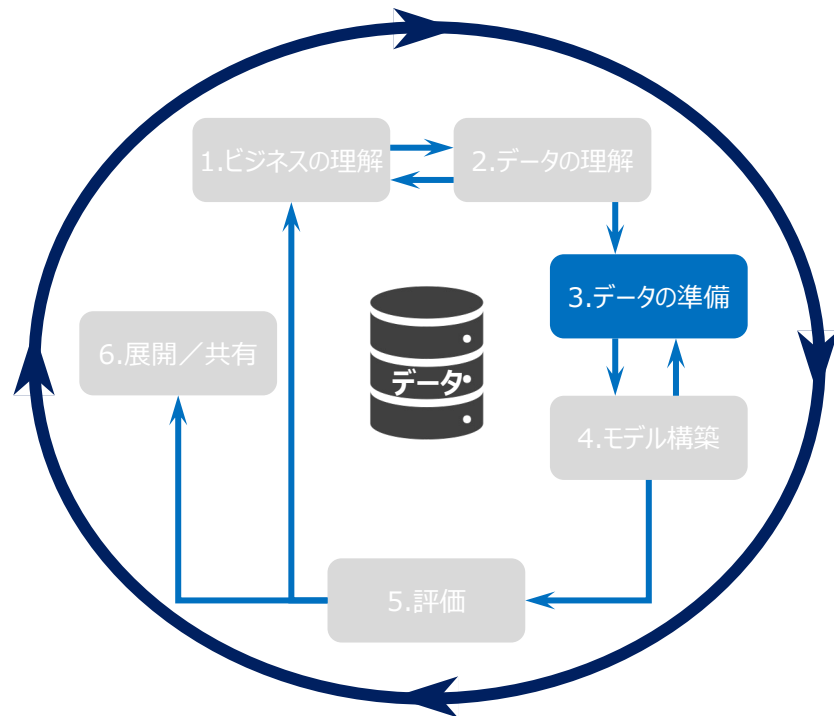
	確認内容	関数	使用例
基本情報	データの先頭/末尾確認	データフレーム.head(行数) / tail(行数)	df.head(10) / df.tail(10)
	データ件数確認	データフレーム.shape	df.shape
	各列のデータ型などの確認	データフレーム.info()	df.info()
	各列の欠損値確認	データフレーム.isnull().sum()	df.isnull().sum()
統計値	要約統計量の算出	データフレーム.describe()	df.describe()
	相関行列	データフレーム.corr()	df.corr(numeric_only=True) #文字列の列がある場合"numeric_only=True"の指定が必要
集計値	クロス集計	pd.crosstab(行見出し, 列見出し)	pd.crosstab(df['役割'], df[['成約]契約まで'])
グラフ	散布図行列	sns.pairplot(df, hue='色分け対象列名') ※seabornのインポート必要	sns.pairplot(df, hue='[成約]契約まで')
	ヒストグラム	データフレーム.hist()	df.hist()
	箱ひげ図	データフレーム.boxplot()	df.boxplot()
	棒グラフ (縦棒/横棒/積み上げ)	縦棒：データフレーム.plot.bar() 横棒：データフレーム.plot.barh()	df.plot.bar(), df.plot.barh() ※積上棒グラフは df.plot.bar( <b>stacked=True</b> )
	その他	過去資料を参考に	

# データ分析の進め方

- データ分析の進め方に関する方法論「CRISP-DM」に基づいて、分析と評価を繰り返して試行錯誤しながら進めるのが一般的である

## CRISP-DM: データマイニング方法論

(CRoss Industry Standard Process for Data Mining)



### 1. ビジネスの理解

- ビジネス、データマイニング目標の決定
- プロジェクトの立ち上げ

### 2. データの理解

- データの収集
- データの調査
- データ品質の検証

### 3. データの準備

- データの選択や除外
- データのクリーニング
- データの構築や統合

### 4. モデル構築

- モデリング手法の選択
- モデルの作成
- モデルの評価

### 5. 評価

- データマイニングの結果の評価
- プロセスの見直し
- 実行可能なアクションリストの作成

### 6. 展開/共有

- 業務への導入計画
- モニタリング、メンテナンスの計画

# 「データの準備」で主に行うこと

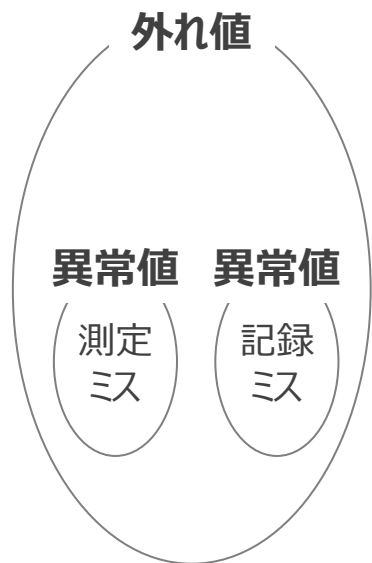
- 前段のデータの理解や後段のモデル構築と並行しながら、必要に応じてデータ加工を行う

レコード方向の増減	クレンジング	<p><b>今回の演習のターゲット</b></p> <p>欠損値を含むレコードの除外／補完 異常値を含むレコードの除外／修正</p>	<table border="1"> <thead> <tr> <th>顧客ID</th> <th>氏名</th> <th>性別</th> <th>...</th> </tr> </thead> <tbody> <tr> <td>0001</td> <td>山田 太郎</td> <td>男性</td> <td>...</td> </tr> <tr> <td>0002</td> <td>鈴木 花子</td> <td>...</td> <td>...</td> </tr> <tr> <td>0003</td> <td>佐藤 次郎</td> <td>男性</td> <td>...</td> </tr> <tr> <td>0004</td> <td>田中 太郎</td> <td>男性</td> <td>...</td> </tr> <tr> <td>...</td> <td>...</td> <td>...</td> <td>...</td> </tr> </tbody> </table>	顧客ID	氏名	性別	...	0001	山田 太郎	男性	...	0002	鈴木 花子	...	...	0003	佐藤 次郎	男性	...	0004	田中 太郎	男性	...	...	...	...	...	「データ観察」 で行うことが多い																																																					
	顧客ID	氏名	性別	...																																																																													
0001	山田 太郎	男性	...																																																																														
0002	鈴木 花子	...	...																																																																														
0003	佐藤 次郎	男性	...																																																																														
0004	田中 太郎	男性	...																																																																														
...	...	...	...																																																																														
	サンプリング	<p>学習データ／テストデータの分割 不均衡データのバランシング (アンダー／オーバーサンプリング)</p>	<table border="1"> <thead> <tr> <th>患者ID</th> <th>容態</th> <th>...</th> </tr> </thead> <tbody> <tr> <td>0001</td> <td>健常</td> <td>...</td> </tr> <tr> <td>0002</td> <td>健常</td> <td>...</td> </tr> <tr> <td>0003</td> <td>罹患</td> <td>...</td> </tr> <tr> <td>0004</td> <td>健常</td> <td>...</td> </tr> </tbody> </table> <p>健康40:罹患10    健康10:罹患10</p> <p>アンダーサンプリング</p>	患者ID	容態	...	0001	健常	...	0002	健常	...	0003	罹患	...	0004	健常	...	「モデル構築」 で行うことが多い																																																														
患者ID	容態	...																																																																															
0001	健常	...																																																																															
0002	健常	...																																																																															
0003	罹患	...																																																																															
0004	健常	...																																																																															
カラム方向の増減	データの統合	<p>複数データファイルの結合</p>	<table border="1"> <thead> <tr> <th colspan="4">トランザクション</th> <th colspan="4">結合</th> <th colspan="4">マスタ</th> </tr> <tr> <th>取引ID</th> <th>レジ</th> <th>顧客ID</th> <th>購入日</th> <th>...</th> <th>顧客ID</th> <th>氏名</th> <th>性別</th> <th>...</th> <th>顧客ID</th> <th>氏名</th> <th>性別</th> <th>...</th> </tr> </thead> <tbody> <tr> <td>2019070100001</td> <td>1</td> <td>0001</td> <td>20190701</td> <td>...</td> <td>0001</td> <td>山田 太郎</td> <td>男性</td> <td>...</td> <td>0001</td> <td>山田 太郎</td> <td>男性</td> <td>...</td> </tr> <tr> <td>2019070100001</td> <td>1</td> <td>0001</td> <td>20190701</td> <td>...</td> <td>0002</td> <td>鈴木 花子</td> <td>女性</td> <td>...</td> <td>0002</td> <td>鈴木 花子</td> <td>女性</td> <td>...</td> </tr> <tr> <td>2019070100001</td> <td>1</td> <td>0001</td> <td>20190701</td> <td>...</td> <td>0003</td> <td>佐藤 次郎</td> <td>男性</td> <td>...</td> <td>0003</td> <td>佐藤 次郎</td> <td>男性</td> <td>...</td> </tr> <tr> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> <td>...</td> </tr> </tbody> </table>	トランザクション				結合				マスタ				取引ID	レジ	顧客ID	購入日	...	顧客ID	氏名	性別	...	顧客ID	氏名	性別	...	2019070100001	1	0001	20190701	...	0001	山田 太郎	男性	...	0001	山田 太郎	男性	...	2019070100001	1	0001	20190701	...	0002	鈴木 花子	女性	...	0002	鈴木 花子	女性	...	2019070100001	1	0001	20190701	...	0003	佐藤 次郎	男性	...	0003	佐藤 次郎	男性	...	...	...	...	...	...	...	...	...	...	...	...	...	...	
トランザクション				結合				マスタ																																																																									
取引ID	レジ	顧客ID	購入日	...	顧客ID	氏名	性別	...	顧客ID	氏名	性別	...																																																																					
2019070100001	1	0001	20190701	...	0001	山田 太郎	男性	...	0001	山田 太郎	男性	...																																																																					
2019070100001	1	0001	20190701	...	0002	鈴木 花子	女性	...	0002	鈴木 花子	女性	...																																																																					
2019070100001	1	0001	20190701	...	0003	佐藤 次郎	男性	...	0003	佐藤 次郎	男性	...																																																																					
...	...	...	...	...	...	...	...	...	...	...	...	...																																																																					
変数の加工 (説明変数・目的変数)	カテゴリーデータのダミー変数化	<p>数値データのダミー変数化</p>	<table border="1"> <thead> <tr> <th>顧客ID</th> <th>氏名</th> <th>性別</th> <th>年齢</th> <th>...</th> <th>顧客ID</th> <th>氏名</th> <th>性別</th> <th>年代</th> <th>...</th> </tr> </thead> <tbody> <tr> <td>0001</td> <td>山田 太郎</td> <td>男性</td> <td>52</td> <td>...</td> <td>0001</td> <td>山田 太郎</td> <td>0</td> <td>50</td> <td>...</td> </tr> <tr> <td>0002</td> <td>鈴木 花子</td> <td>女性</td> <td>23</td> <td>...</td> <td>0002</td> <td>鈴木 花子</td> <td>1</td> <td>20</td> <td>...</td> </tr> <tr> <td>0003</td> <td>佐藤 次郎</td> <td>男性</td> <td>41</td> <td>...</td> <td>0003</td> <td>佐藤 次郎</td> <td>0</td> <td>40</td> <td>...</td> </tr> <tr> <td>0004</td> <td>田中 太郎</td> <td>男性</td> <td>34</td> <td>...</td> <td>0004</td> <td>田中 太郎</td> <td>0</td> <td>30</td> <td>...</td> </tr> </tbody> </table>	顧客ID	氏名	性別	年齢	...	顧客ID	氏名	性別	年代	...	0001	山田 太郎	男性	52	...	0001	山田 太郎	0	50	...	0002	鈴木 花子	女性	23	...	0002	鈴木 花子	1	20	...	0003	佐藤 次郎	男性	41	...	0003	佐藤 次郎	0	40	...	0004	田中 太郎	男性	34	...	0004	田中 太郎	0	30	...																												
	顧客ID	氏名	性別	年齢	...	顧客ID	氏名	性別	年代	...																																																																							
	0001	山田 太郎	男性	52	...	0001	山田 太郎	0	50	...																																																																							
0002	鈴木 花子	女性	23	...	0002	鈴木 花子	1	20	...																																																																								
0003	佐藤 次郎	男性	41	...	0003	佐藤 次郎	0	40	...																																																																								
0004	田中 太郎	男性	34	...	0004	田中 太郎	0	30	...																																																																								
	数値データのダミー変数化	<p>数値データの標準化 (スケーリング)</p>	<table border="1"> <thead> <tr> <th>顧客ID</th> <th>氏名</th> <th>性別</th> <th>年齢</th> <th>...</th> <th>顧客ID</th> <th>氏名</th> <th>性別</th> <th>年代</th> <th>...</th> </tr> </thead> <tbody> <tr> <td>0001</td> <td>山田 太郎</td> <td>男性</td> <td>52</td> <td>...</td> <td>0001</td> <td>山田 太郎</td> <td>0</td> <td>50</td> <td>...</td> </tr> <tr> <td>0002</td> <td>鈴木 花子</td> <td>女性</td> <td>23</td> <td>...</td> <td>0002</td> <td>鈴木 花子</td> <td>1</td> <td>20</td> <td>...</td> </tr> <tr> <td>0003</td> <td>佐藤 次郎</td> <td>男性</td> <td>41</td> <td>...</td> <td>0003</td> <td>佐藤 次郎</td> <td>0</td> <td>40</td> <td>...</td> </tr> <tr> <td>0004</td> <td>田中 太郎</td> <td>男性</td> <td>34</td> <td>...</td> <td>0004</td> <td>田中 太郎</td> <td>0</td> <td>30</td> <td>...</td> </tr> </tbody> </table>	顧客ID	氏名	性別	年齢	...	顧客ID	氏名	性別	年代	...	0001	山田 太郎	男性	52	...	0001	山田 太郎	0	50	...	0002	鈴木 花子	女性	23	...	0002	鈴木 花子	1	20	...	0003	佐藤 次郎	男性	41	...	0003	佐藤 次郎	0	40	...	0004	田中 太郎	男性	34	...	0004	田中 太郎	0	30	...																												
顧客ID	氏名	性別	年齢	...	顧客ID	氏名	性別	年代	...																																																																								
0001	山田 太郎	男性	52	...	0001	山田 太郎	0	50	...																																																																								
0002	鈴木 花子	女性	23	...	0002	鈴木 花子	1	20	...																																																																								
0003	佐藤 次郎	男性	41	...	0003	佐藤 次郎	0	40	...																																																																								
0004	田中 太郎	男性	34	...	0004	田中 太郎	0	30	...																																																																								
	四則演算による変数加工	<p>四則演算による変数加工</p>																																																																															

# 外れ値／異常値／欠損値の確認と対応

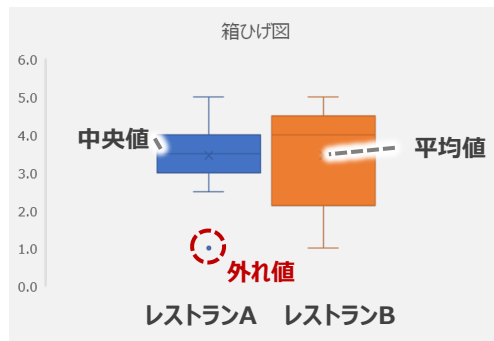
- データに外れ値や異常値、欠損値が存在した場合、**分析結果に悪影響を与える可能性があるため、事前に存在を確かめ、対応方法を検討しておくことが重要**

## 外れ値と異常値



..... データの大部分から大きく外れた値

..... 外れ値のうち、原因がハッキリしているもの (測定ミス、記録ミスなど)



箱ひげ図の描画により  
外れ値の検出が可能

## 欠損値

	A	B	C
欠損値	121	15	187
		21	175
	120	12	177
	115	25	169
	⋮	⋮	⋮

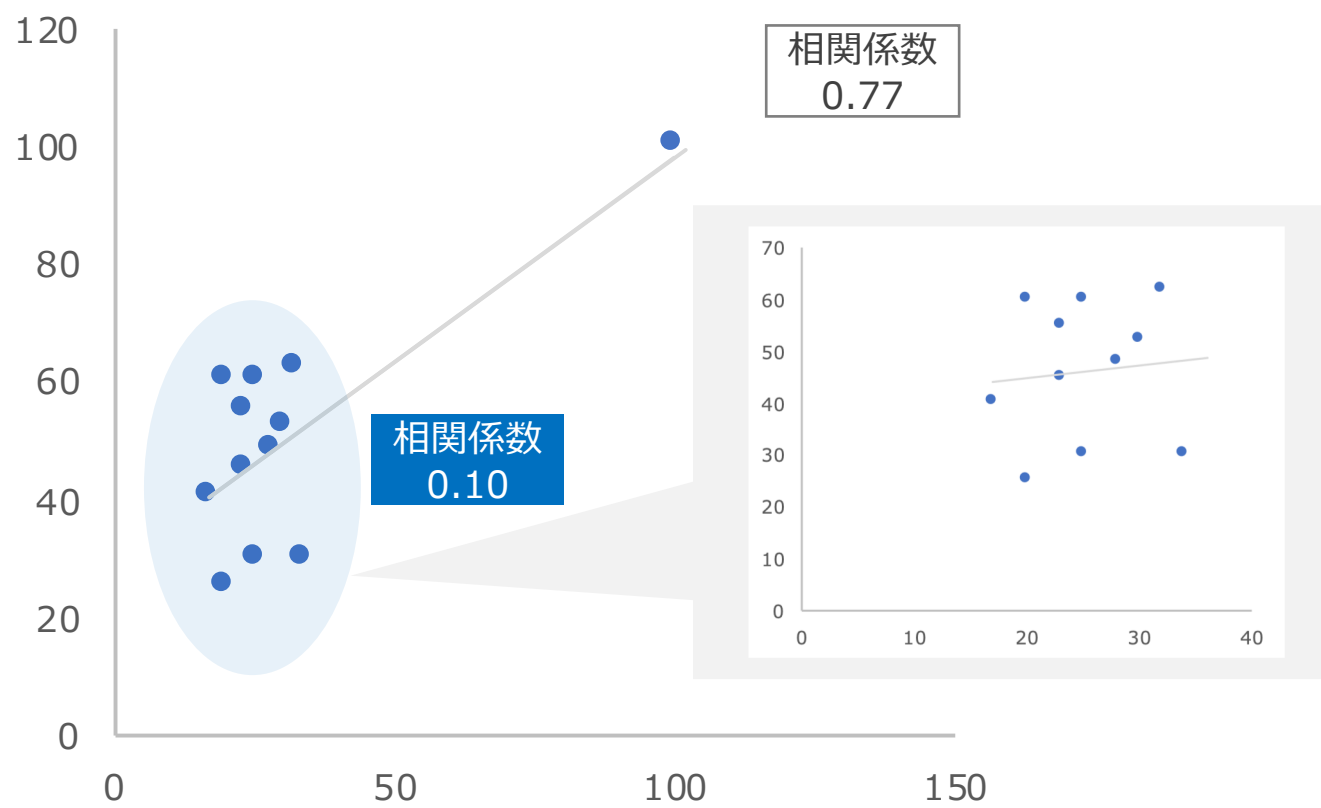
① **行ごと除外**  
※あまりに欠損の多い  
列は列ごと除外

② **補完**  
(平均値／中央値、予測値など)



# 外れ値の影響例

- 相関係数は外れ値の影響を大きく受けるため、数字だけに惑わされぬよう、散布図の確認も併せて行うことが重要である

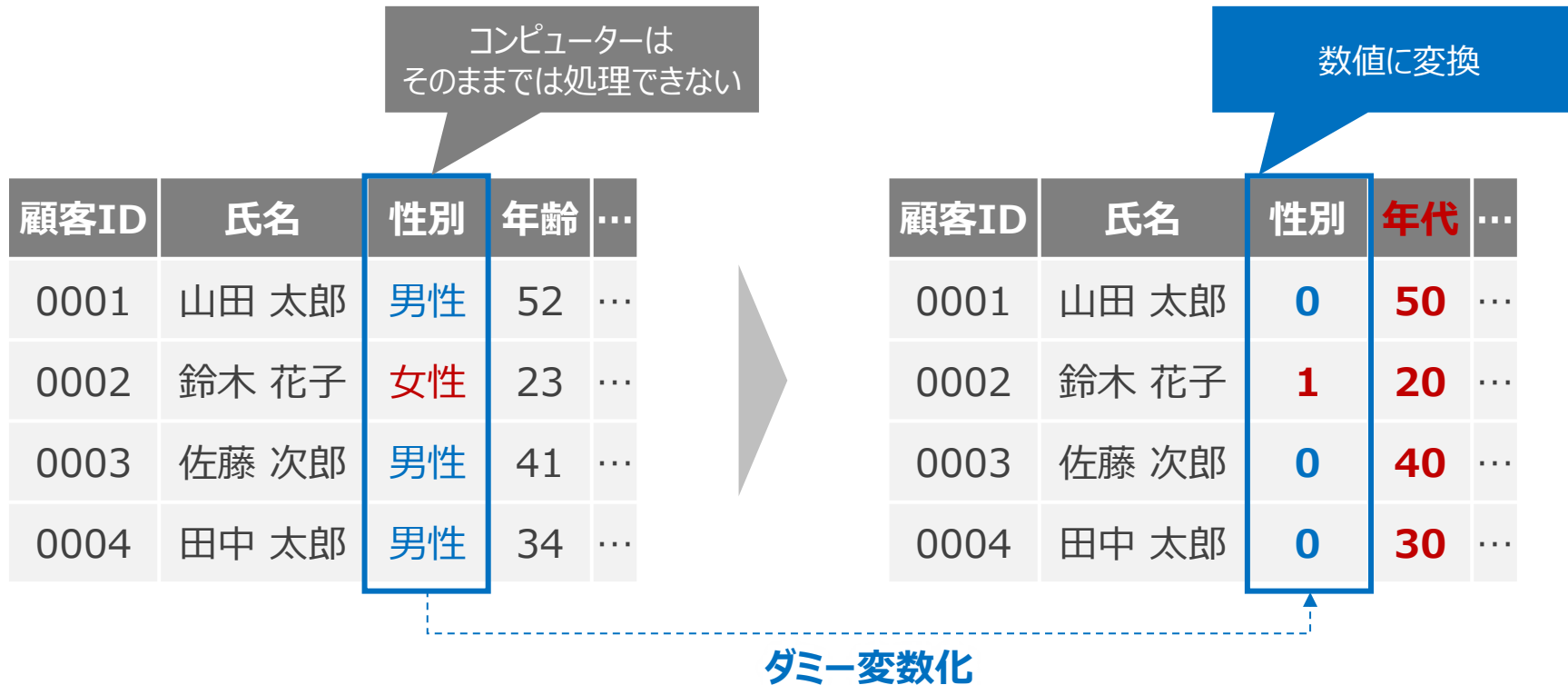




Google Colaboratory上での  
レクチャー&演習

# カテゴリデータ（文字列データ）とダミー変数化

- 通常の数値データに対して、血液型（A, B, O, AB）や部署（経営企画部、DX推進部、…）、性別など、分類上の区別に用いられる変数を「**カテゴリデータ**」と呼ぶ
- カテゴリデータは通常文字列で表されるため、擬似的に数値変換する「**ダミー変数化**」が必要



# ダミー変数化の種類

- ダミー変数化のメジャーな考え方は2つあり、男性／女性、契約／未契約のような**2カテゴリのデータは連続値化**をし、職種や部署のような**カテゴリ間に連続性がない場合には横持ち化**を行う

## ① 連続値化 (Label/Ordinal Encoding)

2カテゴリ (男/女) やアンケートのような連続的な尺度に有効

元データ		ダミー変数化
男性	----->	<b>0</b>
女性	----->	<b>1</b>
不満である	----->	<b>0</b>
どちらでもない	----->	<b>1</b>
満足している	----->	<b>2</b>

## ② 横持ち化 (One Hot Encoding)

カテゴリ間に直接的な連続性がない場合に有効

元データ		ダミー変数化
...	職種	...
...	会社員	...
...	自営業	...
...	会社員	...
...	会社員	...
...	会社員	...
...	会社員	...
...	医師	...
...	⋮	...

...	自営業	医師	会社員
...			1
...	1		
...			1
...			1
...			1
...		1	
...	⋮		

自営業	----->	<b>0</b>
会社員	----->	<b>1</b>
医師	----->	<b>2</b>

この場合、連続値化は不適



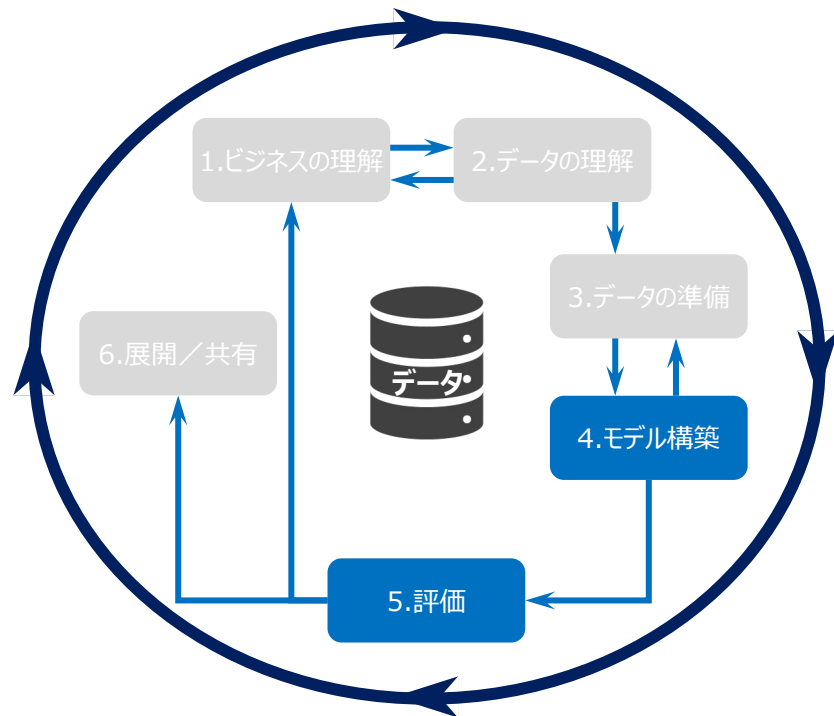
Google Colaboratory上での  
レクチャー&演習

# データ分析の進め方

- データ分析の進め方に関する方法論「CRISP-DM」に基づいて、分析と評価を繰り返して試行錯誤しながら進めるのが一般的である

## CRISP-DM: データマイニング方法論

(CRoss Industry Standard Process for Data Mining)



### 1. ビジネスの理解

- ビジネス、データマイニング目標の決定
- プロジェクトの立ち上げ

### 2. データの理解

- データの収集
- データの調査
- データ品質の検証

### 3. データの準備

- データの選択や除外
- データのクリーニング
- データの構築や統合

### 4. モデル構築

- モデリング手法の選択
- モデルの作成
- モデルの評価

### 5. 評価

- データマイニングの結果の評価
- プロセスの見直し
- 実行可能なアクションリストの作成

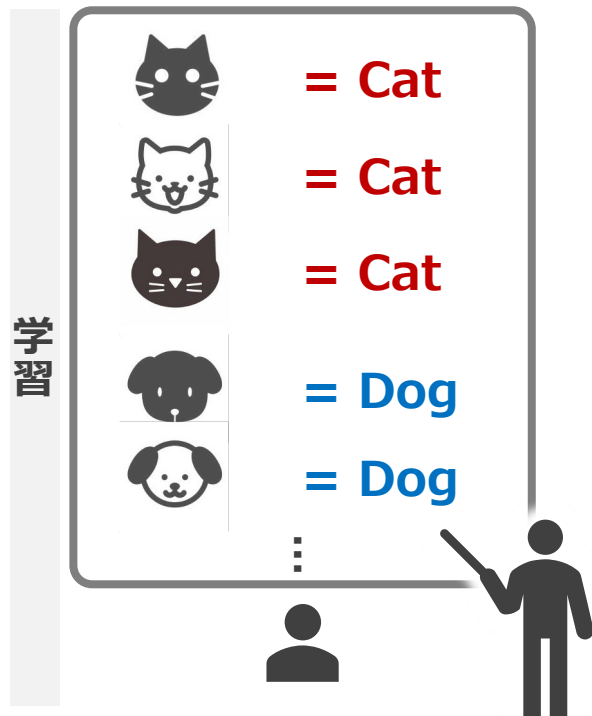
### 6. 展開/共有

- 業務への導入計画
- モニタリング、メンテナンスの計画

# 機械学習の分類

## 教師あり学習

あらかじめ「正解」を与えて  
各データと正解の関係を学習させる



推論



## 教師なし学習

「正解」を与えずに、各データ  
のパターンから自分で学習する



推論



※分類されたグループの意味づけは人が行う

# 教師あり学習のイメージ（数値予測とクラス分類）

- 各顧客レコードに対して数値もしくはカテゴリ値（クラス）の解答を与え関係性を学習

## 数値予測

(回帰)

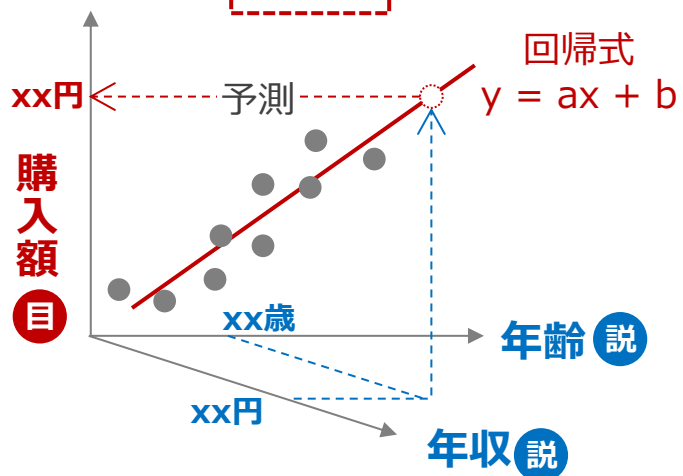
説明変数 目的変数

顧客ID	名前	年齢	年収	購入額	購入有無	...
0001	xx	25	300万	35,000	購入	...
0002	xx	35	600万	68,000	購入	...
0003	xx	18	120万	0	非購入	...
0004	xx	42	820万	85,000	購入	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮

データ例

イメージ

分析例



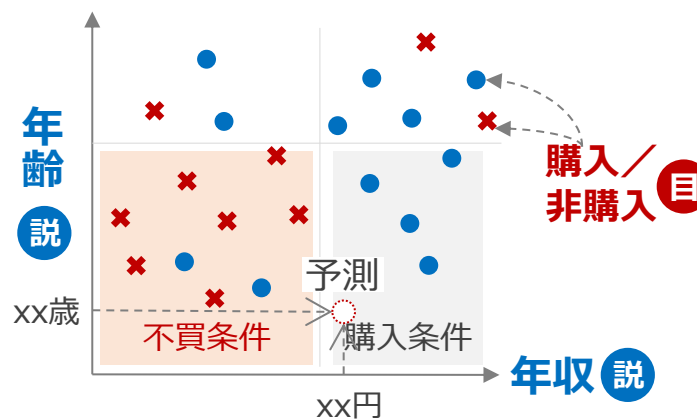
- 小売：売上の予測
- 小売：需要予測
- 製造：不良発生率の予測

## クラス分類

(2クラス or 多クラス分類)

説明変数 目的変数

顧客ID	名前	年齢	年収	購入額	購入有無	...
0001	xx	25	300万	35,000	購入	...
0002	xx	35	600万	68,000	購入	...
0003	xx	18	120万	0	非購入	...
0004	xx	42	820万	85,000	購入	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮



- 小売：購入 / 非購入顧客の分類
- 医療：生存条件の分析
- 製造：故障種類の分類



# 教師なし学習のイメージ (クラスタリング)

- 各データ間の距離に基づき、近接データ (=類似度が高いデータ) 同士のグループ (クラスタ) を作り、データを分類する手法。学習データなしでデータを大きく層別したい場合に有効

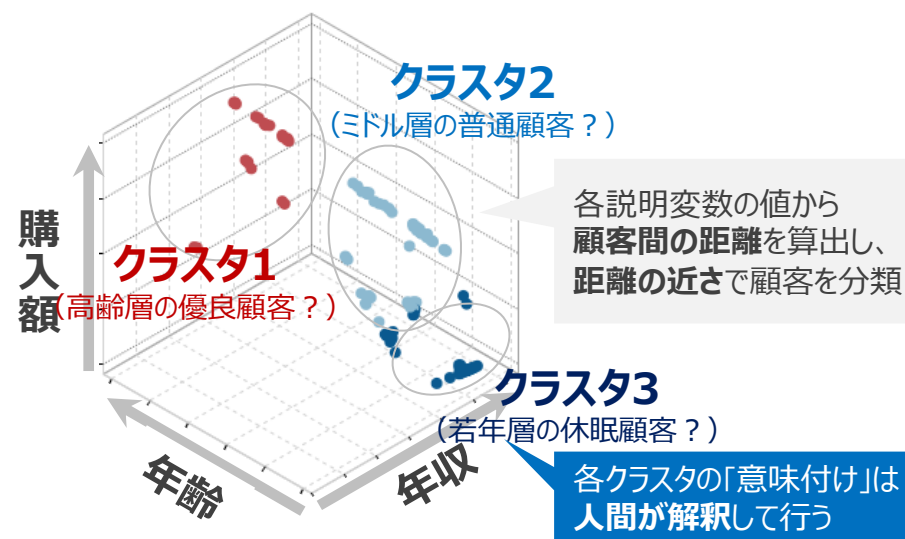
## データ例

顧客ID	名前	年齢	年収	購入額	購入有無	...
0001	xx	25	300万	35,000	購入	...
0002	xx	35	600万	68,000	購入	...
0003	xx	18	120万	0	非購入	...
0004	xx	42	820万	85,000	購入	...
⋮	⋮	⋮	⋮	⋮	⋮	...

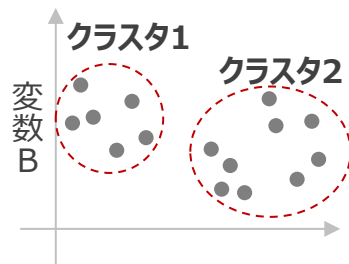
説明変数

※目的変数は無し

クラスタリング



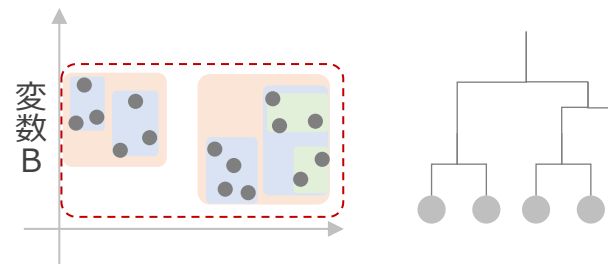
## 非階層的クラスタリング



主な手法

- k-means法 (k平均法)
- 混合ガウス

## 階層的クラスタリング



主な手法

- 最短距離法
- 最長距離法
- 群平均法
- ウォード法

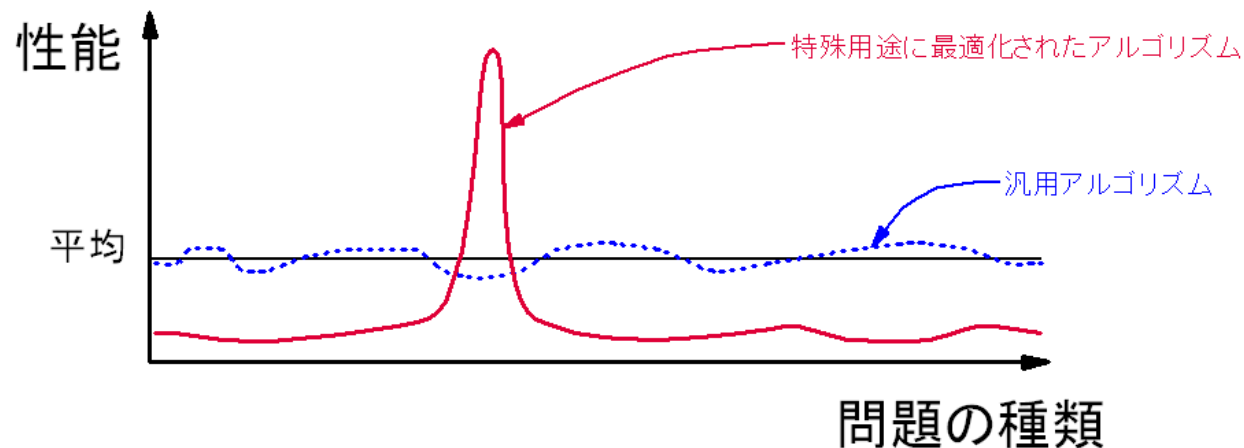
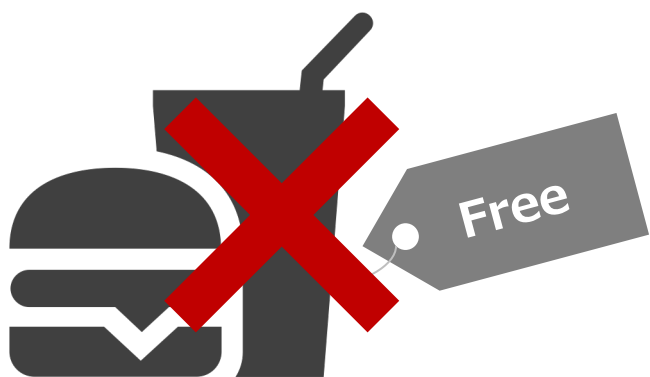
# 代表的な機械学習手法

- 目的やデータの特徴に合わせて手法を選ぶが、適宜、**手法の試行錯誤**が必要となることも多い

		手法	概要	主なアルゴリズム
教師あり (答え合わせあり)	数値予測	数値予測	連続的な数値に対して、いくつかの影響因子を用いてある関数モデル ( $y=f(x)$ ) を当てはめる	線形／非線形回帰 (2変数：単回帰、多変数：重回帰) 回帰木 ニューラルネットワーク
		時系列予測	継時的な変動に対して、過去の履歴やトレンドから将来の動きを予測	ARモデル (自己回帰モデル) ARIMAモデル (自己回帰移動平均モデル)
	分類	クラス分類 (多クラス/ 二クラス)	あらかじめ分類のラベル付けがされたデータに対して、いくつかの影響因子を用いて、分類境界を抽出する	ロジスティック回帰 決定木分析 判別分析 ニューラルネットワーク サポートベクタマシン (SVM)
教師なし (答え合わせなし)		クラスタリング (グルーピング)	データに含まれる各変数の特徴から、類似データをグルーピングする	K-means 階層的クラスタリング
	異常値検出	データの分布に基づき、分布から外れる異常な値を検出する	k-近傍法 (k-NN) Local Outlier Factor (LOF)	
教師なし (答え合わせなし)	パターン抽出 (アソシエーションルール分析)	頻出するデータのパターン (データの組み合わせ) を抽出する (例：オムツとビールの併売)	アプリアリ シーケンス	
	データ要約 (次元削減)	高次元 (多変数) のデータに対して、より少ない変数で説明できる変数や軸を抽出する	主成分分析 因子分析	
強化学習 (報酬)	最適化	試行錯誤を通じてある目的達成のための最適な行動を学習 (例：自動運転、アルファ碁)	Q-Learning SARSA (State-action-reward-state-action) モンテカルロ法	

# No Free Lunch定理：タダ飯は無い！

- あらゆる問題で高性能を示すモデルは理論上、構築不可能であることが示されている
- 問題設定や分析目的に応じて、最適なアルゴリズムの検討が必要である



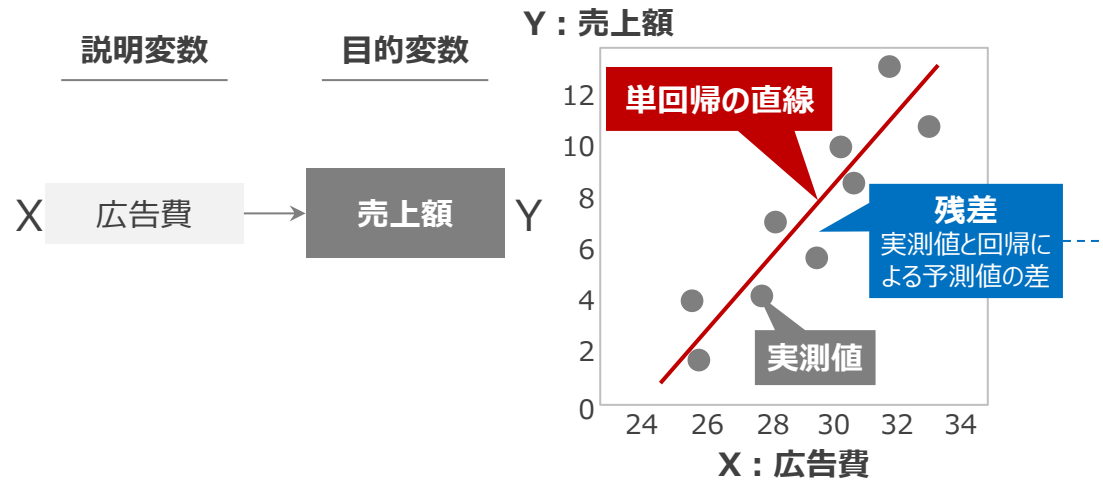
\*画像出典：Wikipedia

**No-free-lunch theorem:** コスト関数の極値を探索するあらゆるアルゴリズムは、全ての可能なコスト関数に適用した結果を平均すると同じ性能となる (Wolpert and Macready, 1995年)

# 線形回帰分析による数値予測

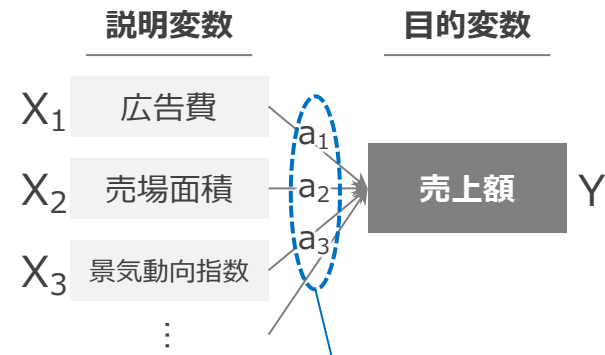
- 目的変数Yとその影響因子である説明変数Xとで関数 ( $Y=f(X)$ ) に当てはめることを**回帰**という
- 特に $f(X)$ が線形関数の場合、**線形回帰**と呼び、1変数では**単回帰**、多変数では**重回帰**という

## 単回帰



残差の2乗和が最小になるように  
係数 $a$ , 定数項 $b$ を決める (最小二乗法)

## 重回帰



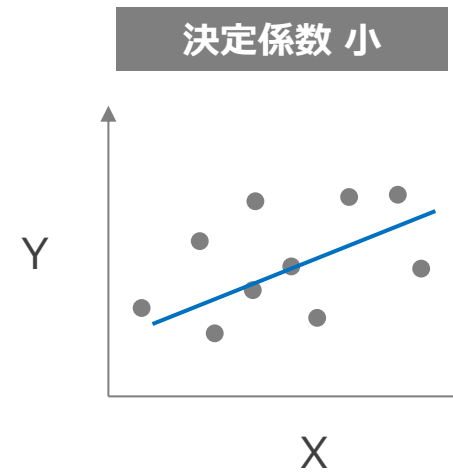
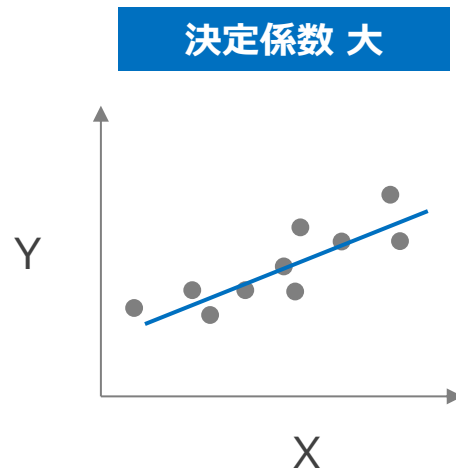
各回帰係数は、  
各説明変数の影響の大きさを表す

$$\text{売上額} = a \times \text{広告費} + b$$

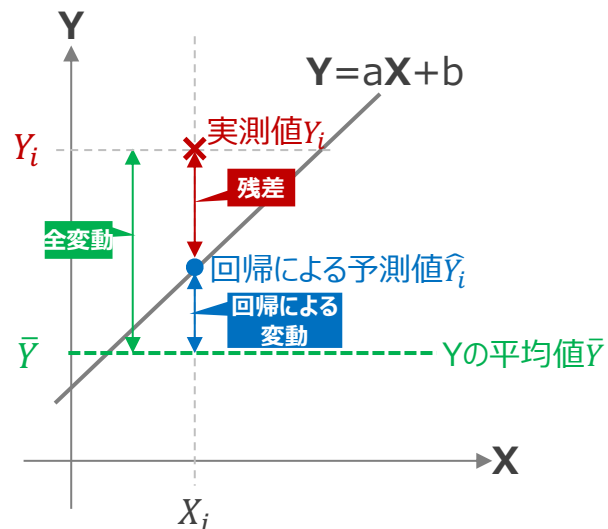
<b>目的変数</b> (Y)	<b>回帰係数</b> (傾き)	<b>説明変数</b> (X)	<b>定数項</b> (Y切片)
結果となる 変数	Xの影響 の大きさ	原因となる 変数	広告費=0の ときの売上額

# 決定係数によるモデルの当てはまり評価

- モデルの当てはまりの良さは、「決定係数」によって表すことが多い



## ▼参考：決定係数の定義イメージ



$$\text{決定係数 } R^2 = \frac{\text{回帰による変動 (平方和)}}{\text{全変動 (平方和)}}$$

(全変動のうち、回帰式で説明できる変動の割合)

- ・ **0~1**の値をとる
- ・ (一般的な定義では) **相関係数の2乗に等しい**

(参考) 相関係数の目安

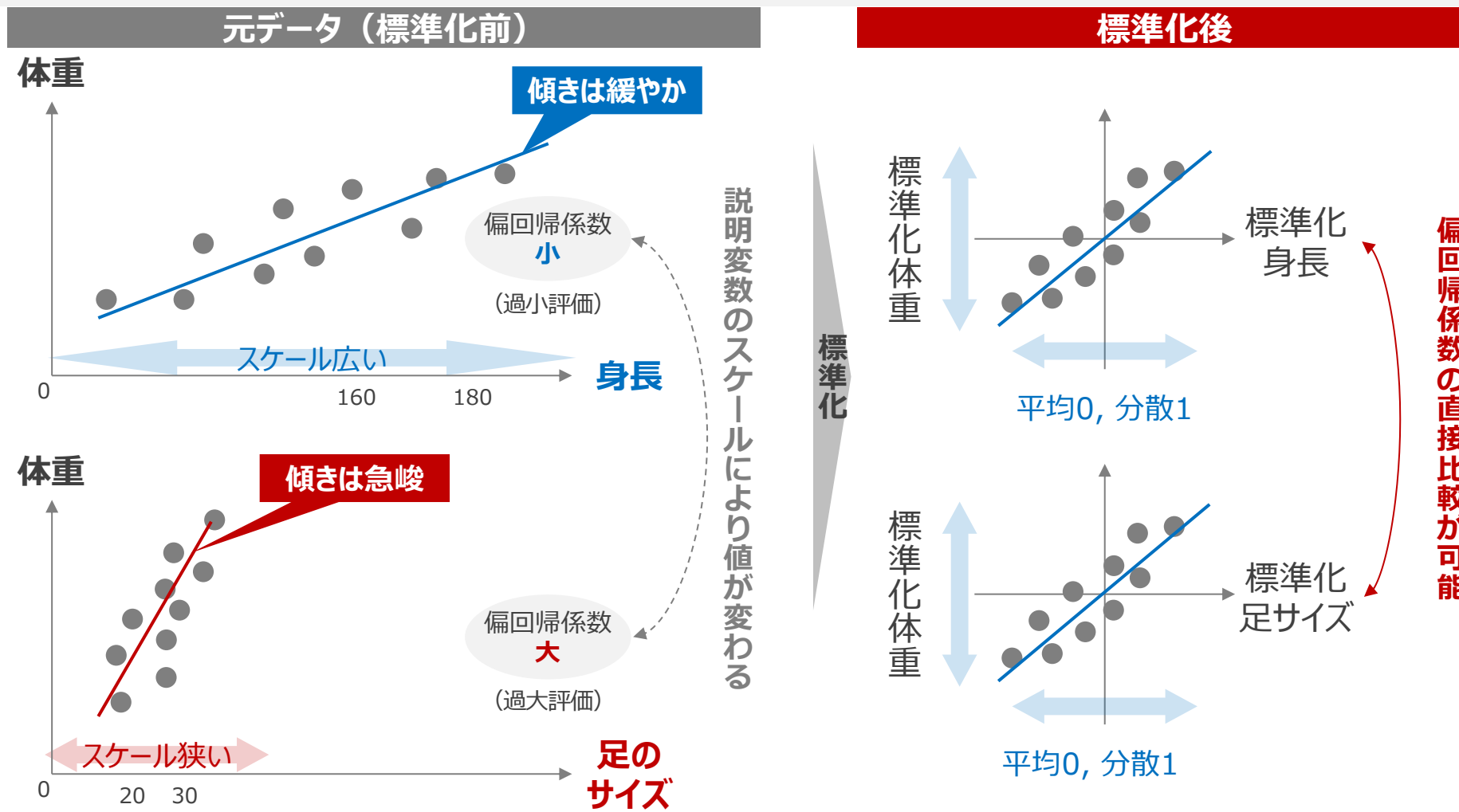
相関係数 (絶対値)	相関の強弱
1.0	強い相関
0.7	やや相関あり
0.4	弱い相関
0.2	ほぼ相関なし
0	



Google Colaboratory上での  
レクチャー&演習

# 回帰の注意点：データの標準化

- 重回帰分析の偏回帰係数は、目的変数に対する各説明変数の「傾き」を表しており、**その大きさは各変数のスケールに依存するため、そのままでは直接比較が困難**
- **予めデータを標準化しておくことで、偏回帰係数の直接比較が可能となる**

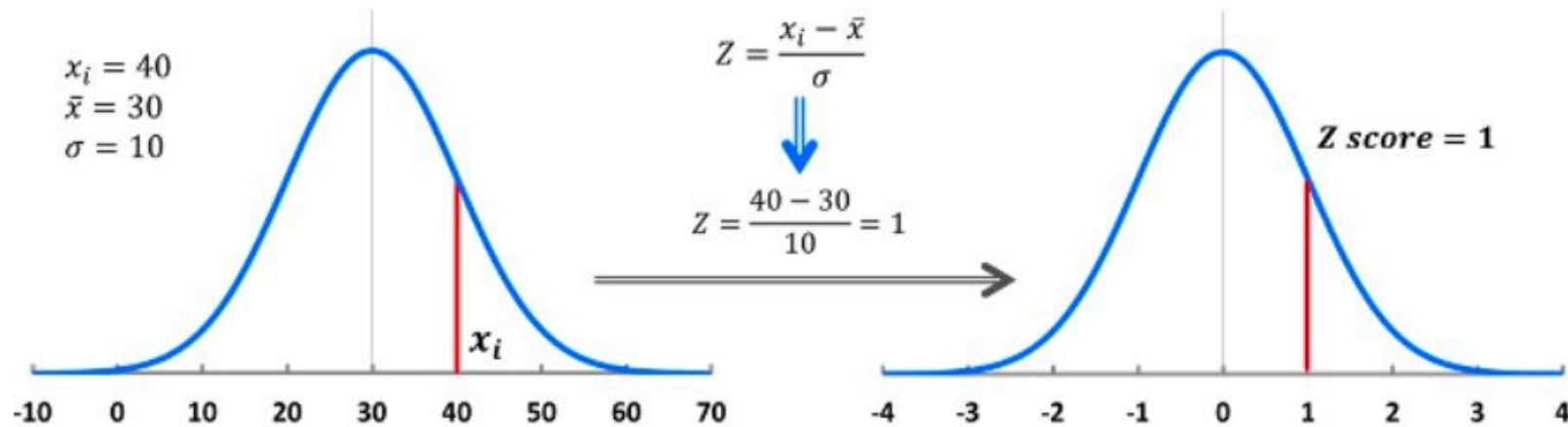


# 標準化の方法 (Zスコア)

- 重回帰分析の偏回帰係数は、各変数のスケールの大小に依存するため、直接比較が困難
- **Zスコア**を算出することでデータの標準化が行え、偏回帰係数の直接比較が可能となる

$$Z \text{ score} = \frac{x_i - \mu}{\sigma}$$

$x_i$  = 母集団を構成する要素iの値  
 $\mu$  = 母平均  
 $\sigma$  = 標準偏差



出典 : <https://www.monodukuri.com/gihou/article/703>

(参考) 偏差値

偏差値は、「平均が50、標準偏差が10」になるように変換したもの :  $50 + 10 \times Z \text{スコア}$





Google Colaboratory上での  
レクチャー&演習

# 決定木分析によるクラス分類

- 決定木分析は、与えられた分類情報に合致するように分割条件を導いていく手法
- 分析結果の解釈性が非常に高く技術的な解釈が容易なため、製造業の分析でも頻用される

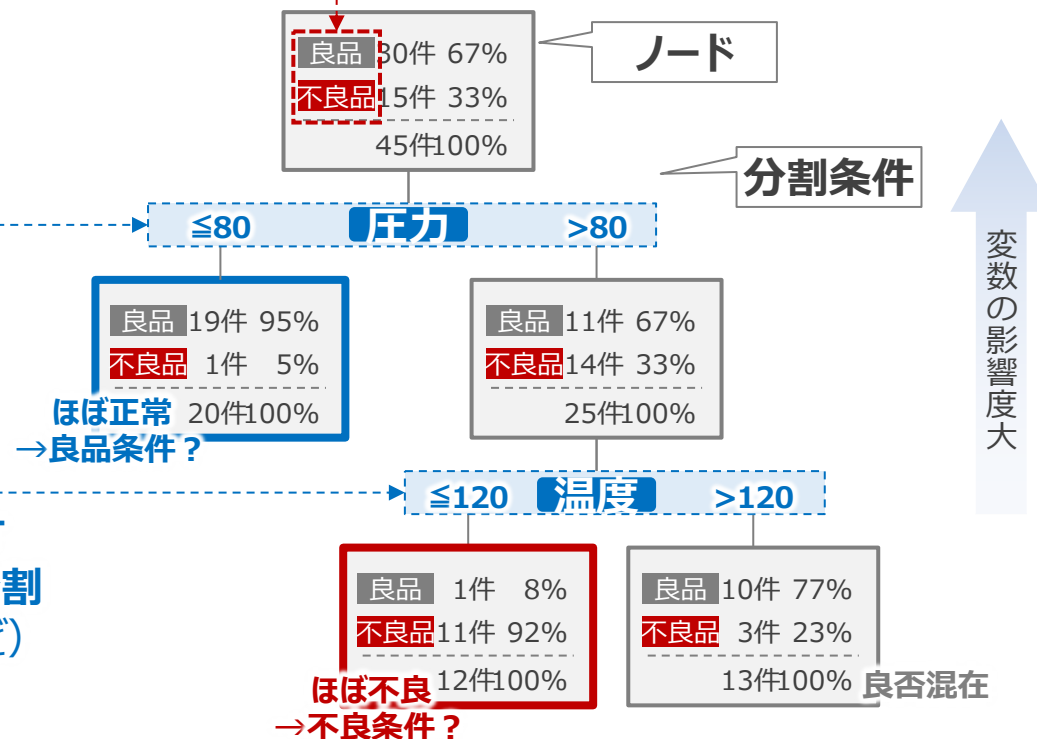
データ例

ID	品質	圧力	温度	流量	事象A	作業者	...
0001	良品	77	121	25.1	あり	A	...
0002	不良品	82	115	35.2	なし	B	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮

目的変数

説明変数

決定木モデル



各クラス（目的変数）の混在度を表す「不純度」が最も減少するような条件で分割  
（不純度の例：ジニ係数、エントロピーなど）

※分岐の縦方向は「AND」条件、横方向は「OR」条件



Google Colaboratory上での  
レクチャー&演習

# モデルの評価指標：Precision と Recall

- モデル評価は、**分析対象や目的に応じて最適な指標を選択したり**、複数の指標を用いて行う
- 例えばマーケティングではDM配布先を限定してコスト抑制するためにPrecisionを、医療では重大所見を見逃すと命に関わるためRecallを重視するなど、目的によって重視する指標は異なる

**混同行列**  
※予測対象が「不良品」の場合

		予測結果		
		不良品	良品	
実データ	不良品	<b>1</b> 正答 (True Positive)	<b>3</b> 誤り：見逃し (False Negative)	再現率
	良品	<b>2</b> 誤り：誤報 (False Positive)	<b>4</b> 正答 (True Negative)	特異度



**評価指標**

**正解率** (Accuracy) =  $\frac{1+4}{1+2+3+4}$

**適合率**, 精度, 陽性的中率 (Precision, PPV\*) =  $\frac{1}{1+2+3+4}$  誤報の少なさ  
\*Positive Predictive Value

陰性的中率 (NPV\*) =  $\frac{4}{1+2+3+4}$  負例の誤報の少なさ  
 \*Negative Predictive Value

**再現率**, 感度, 真陽性率 (Recall, Sensitivity) =  $\frac{1}{1+2+3+4}$  見逃し (取りこぼし) の少なさ

特異度, 真陰性率 (Specificity) =  $\frac{4}{1+2+3+4}$  負例の見逃しの少なさ

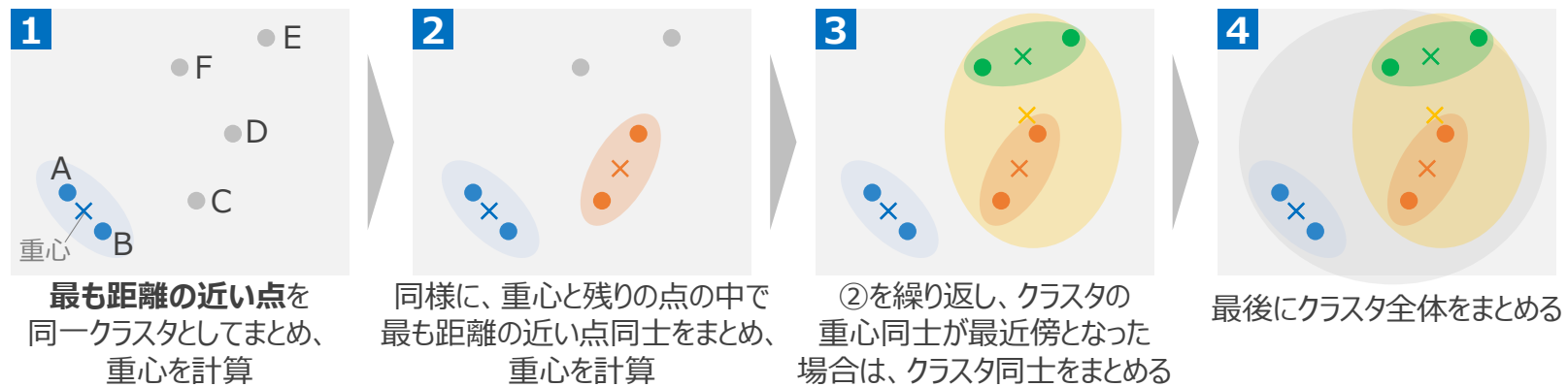
**F値**, F1スコア (F-measure) =  $\frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$  誤報と見逃しの少なさのバランス  
Precision と Recall の調和平均



Google Colaboratory上での  
レクチャー&演習

# 階層的クラスタリング (凝集型階層クラスタリング)

- 凝集型階層クラスタリングは、距離に応じて小さいクラスタを束ねて階層的に分類する手法
- クラスタ数は自動的に決定してくれる他、分類過程を可視化した**樹形図 (デンドログラム)** も同時に出力されるので、結果の解釈やクラスタ数の決定に役立つ



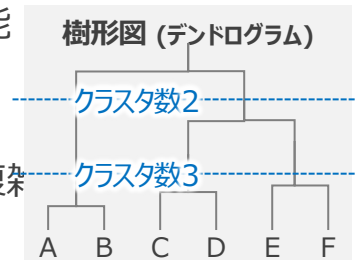
## 凝集型階層的クラスタリング (agglomerative hierarchical clustering)

メリット

クラスタ数は自動決定  
樹形図により**分類過程が可視化**されることで、  
妥当なクラスター数を人が判断可能

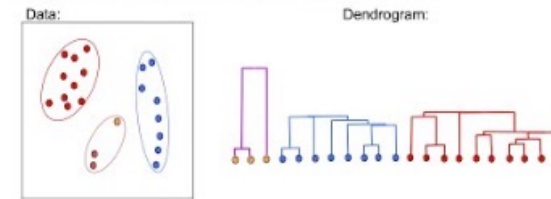
デメリット

計算量が膨大  
データ量が多い場合、樹形図が複雑  
となり、解釈が困難になる



### Iteration m-3

Builds up a sequence of clusters ("hierarchical")



In matlab: "linkage" function (stats toolbox)

Algorithmic Complexity:  $O(m^2 \log m) + (m-3) \cdot O(m \log m) +$

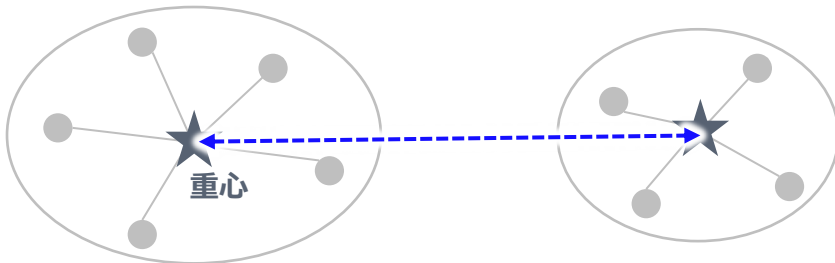
出典: <https://youtu.be/OcoE7JlbXvY>

# 参考：クラスタ間「近さ」の評価尺度バリエーション

- クラスタ間の「近さ」を測る指標には様々あるが、一概にどれが良いとは言えないため、**複数試して比較**するのが一般的である。ただし、一般には、群平均法やWard法（次頁）が頻用される
- 最短距離／最長距離法は、計算量が少なくて済む反面、1点の影響を大きく受けやすい

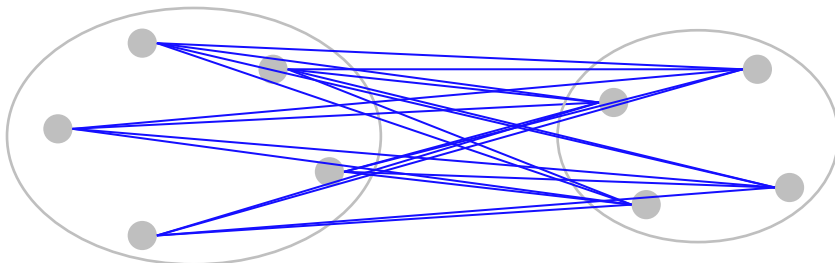
重心法

重心間の距離が近いクラスタを結合



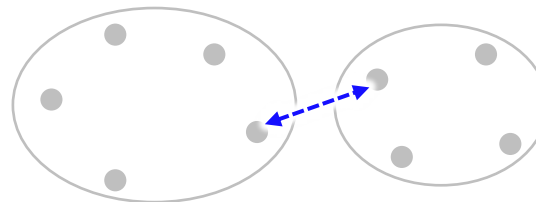
群平均法

クラスタ間で全データ間の距離を算出し、その**平均値**が近いクラスタを結合



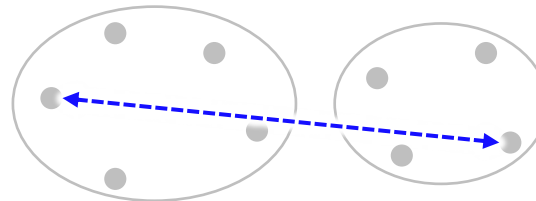
最短距離法

2つのクラスタ間で**最近傍**のデータをクラスタ間距離として採用



最長距離法

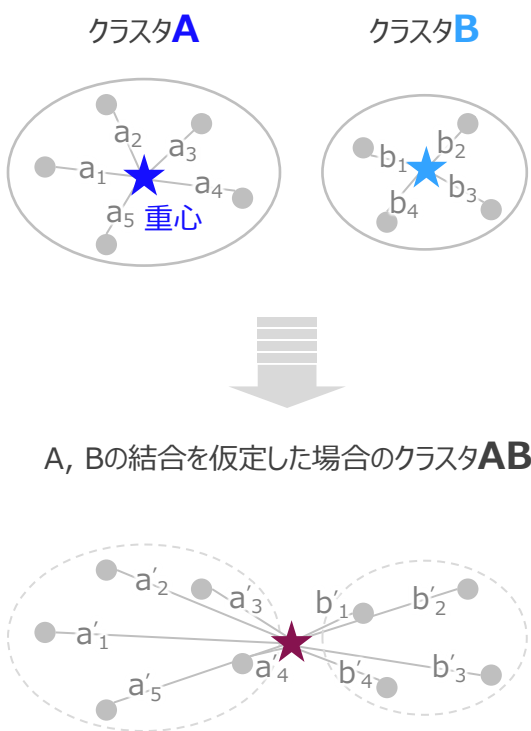
2つのクラスタ間で**最遠方**のデータをクラスタ間距離として採用



# 参考：Ward法の考え方

- Ward法\*は最もよく用いられる手法であり、計算量が多いが、各データ点とクラスタ重心との関係性まで評価しているため、他手法に比べ、**分類感度が高い**とされる

\*米国の統計学者Joe H. Ward, Jr.が1963年に発表した論文にちなむ



- 1 「クラスタ重心」と、「当該クラスタ内の各データ」との距離の総和（二乗和）をクラスタごとに算出

クラスタAの場合

$$A = a_1^2 + a_2^2 + a_3^2 + a_4^2 + a_5^2$$

クラスタBの場合

$$B = b_1^2 + b_2^2 + b_3^2 + b_4^2$$

- 2 注目する2つのクラスタを結合した場合を仮定し、「結合後のクラスタ重心」と「当該クラスタ内の各データ」との距離の総和（二乗和）を算出

$$AB = a_1'^2 + a_2'^2 + a_3'^2 + a_4'^2 + a_5'^2 + b_1'^2 + b_2'^2 + b_3'^2 + b_4'^2$$

- 3 **1**と**2**の差、つまり、 $AB - (A+B)$  が**最小**となるクラスタ結合を採用（結合前後でクラスタ内のばらつきに変化なし→統合してもOKと判定）

※近くにあり、ばらつきの小さいクラスタ同士が結合しやすい





Google Colaboratory上での  
レクチャー&演習

# Part 3

## 実際のビジネスデータ解析 に向けた実践演習

- ✓ e-Stat, SSDSEの紹介と使い方
- ✓ データ分析の実践演習
- ✓ 今後の継続学習、実践活用のためのポイント

# 政府統計のオープンデータサイト “e-Stat”

- オープンデータサイト “e-Stat” では、様々な政府統計を参照・ダウンロードすることができる  
\*各府省の統計データを一元的にみれるポータルサイト。総務省統計局が整備し、独立行政法人統計センターが運用管理している
- 国民経済の状況を把握する上で不可欠かつ、唯一無二なデータも多い

## ■ 政府統計の総合窓口 e-Stat

<https://www.e-stat.go.jp/>

**e-Stat** 統計で見る日本  
政府統計の総合窓口

e-Statは、日本の統計が閲覧できる政府統計ポータルサイトです

お問い合わせ | ヘルプ | English

ログイン 新規登録

統計データを探す 統計データの活用 統計データの高度利用 統計関連情報 リンク集

- 統計データを探す (政府統計の調査結果を探します) その他の絞り込み
  - すべて (政府統計一覧の中から探します)
  - 分野 (17の統計分野から探します)
  - 組織 (統計を作成した府省等から探します)
- 統計データを活用する
  - グラフ (主要指標をグラフで表示 (統計ダッシュボード))
  - 時系列表 (主要指標を時系列表で表示 (統計ダッシュボード))
  - 地図 (地図上に統計データを表示)
  - 地域 (都道府県、市区町村の主要データを表示)

キーワード検索:  検索

利用ガイド

- 統計データの高度利用
  - マイクロデータの利用 (公的統計のマイクロデータの利用案内)
  - 開発者向け (API、LODで統計データを取得)
- 統計関連情報
  - 統計分類・調査計画等

# 教育用標準データセット (SSDSE)

- 教育用標準データセット (SSDSE; Standarzed Statistical Data Set for Education) は、独立行政法人統計センターが作成・公開している、データ分析教育用の統計データ
- 主要な公的統計を地域別に一覧できる表形式のデータセットで、直ちにデータ分析に利用可能



## ▼ SSDSEデータセット一覧

名称	内容
<u>SSDSE-市区町村 (SSDSE-A)</u>	<b>1741市区町村</b> ×多分野125項目 全国の全市区町村の、人口、経済、教育、労働、医療、福祉など、様々な分野の統計データを収録
<u>SSDSE-県別推移 (SSDSE-B)</u>	<b>47都道府県</b> ×12年次×多分野109項目 人口、経済、教育、労働、医療、福祉など、様々な分野の統計データを、12年分の時系列で収録
<u>SSDSE-家計消費 (SSDSE-C)</u>	<b>全国・47都道府県庁所在市</b> ×家計消費226項目 1世帯当たりの食料の年間支出金額（消費額）を、魚介、肉、野菜、果物、菓子、飲料などに分類し、それぞれ詳細な品目別にデータを収録
<u>SSDSE-社会生活 (SSDSE-D)</u>	<b>全国・47都道府県</b> ×男女別×社会生活121項目 男女別に、スポーツ、趣味・娯楽、ボランティアなどの詳細な活動データや、1日の睡眠、食事、学業、家事、仕事、趣味・娯楽などの時間配分データを収録
<u>SSDSE-基本素材 (SSDSE-E)</u>	<b>全国・47都道府県</b> ×多分野90項目 人口、経済、教育、文化、医療、福祉など、様々な分野の統計データを収録し、初学者にも扱いやすいデータセット
<u>SSDSE-気候値 (SSDSE-F)</u>	<b>47都道府県庁所在市</b> ×月・年×気象42項目 気温、気圧、風速、日照、降水、降雪など、様々な気象データについて、月・年別の平年値を収録



# Google Colaboratory上での レクチャー&演習

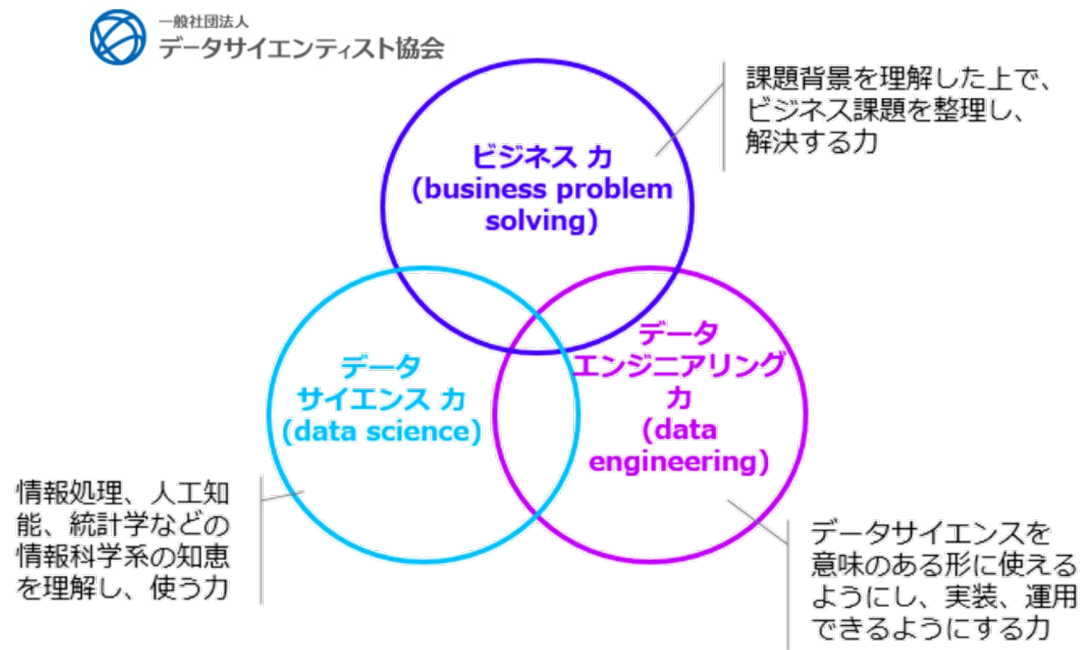
Part 3 : Google Colab.リンク

<https://colab.research.google.com/drive/1Po8q0UYq2f5kF0BSAsfqak-GlCh15Di9?usp=sharing>

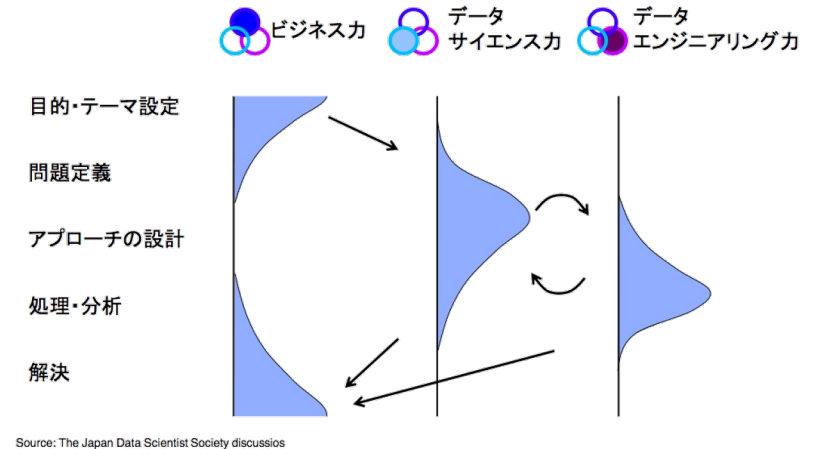
今後の継続学習、  
実践活用のためのポイント

# データサイエンティストに求められる知識・スキルセット

- データサイエンティスト協会の定義：「データサイエンティストとは、**データサイエンス力、データエンジニアリング力**をベースにデータから価値を創出し、**ビジネス**課題に答えを出すプロフェッショナル」
- これら3スキルはどれも不可欠で、分析フェーズによって中心となるスキルが変化する、としている



## 課題解決の各フェーズで要求されるスキルセットのイメージ



出展：データサイエンティスト協会資料  
<https://www.datascientist.or.jp/news/n-news/post-255/>

# それぞれのスキルの伸ばし方 – ビジネス力

## ■自業務における課題について、問題点を整理してみる

## ■仮説思考力の高め方のコツ

✓起こった事象や自分の思考結果に対して、

- 「**だから何なのか？**」(so what) と
- 「**なぜそうなるのか？**」(how) を繰り返す

参考：なぜなぜ分析（トヨタ生産方式から生まれたフレームワーク）

✓日常生活の中で将来予測をする癖をつける

- 新聞記事・ニュース
- 職場での会話
- 日常会話

✓以下はバイブルとして知られる本です

（これ以外にも今は色々出ていると思うので、本屋などで自身に合うものを探してみてください）

- 仮説思考 BCG流 問題発見・解決の発想法 | 内田 和成 <https://www.amazon.co.jp/dp/4492555552>
- イシューからはじめよ——知的生産の「シンプルな本質」 | 安宅和人 <https://www.amazon.co.jp/dp/4862760856>
- ロジカル・シンキング Best solution | 照屋 華子, 岡田恵子 <https://www.amazon.co.jp/dp/B00978ZQOG>
- 入門 考える技術・書く技術 | 山崎 康司 <https://www.amazon.co.jp/dp/B0081WMO4W>



# それぞれのスキルの伸ばし方 - データサイエンスカ / データエンジニアリングカ

## ■実際に自身の課題感を持って実践的な分析を行うことが最も効果的です

✓ 自業務におけるデータを活用する

✓ 自業務における切迫した課題を分析テーマとする

✓ オープンデータを用いて、自分なりに気になる社会課題などについて分析する

○ e-Stat (政府統計データサイト) : <https://www.e-stat.go.jp/>

✓ コンペティションに参加する

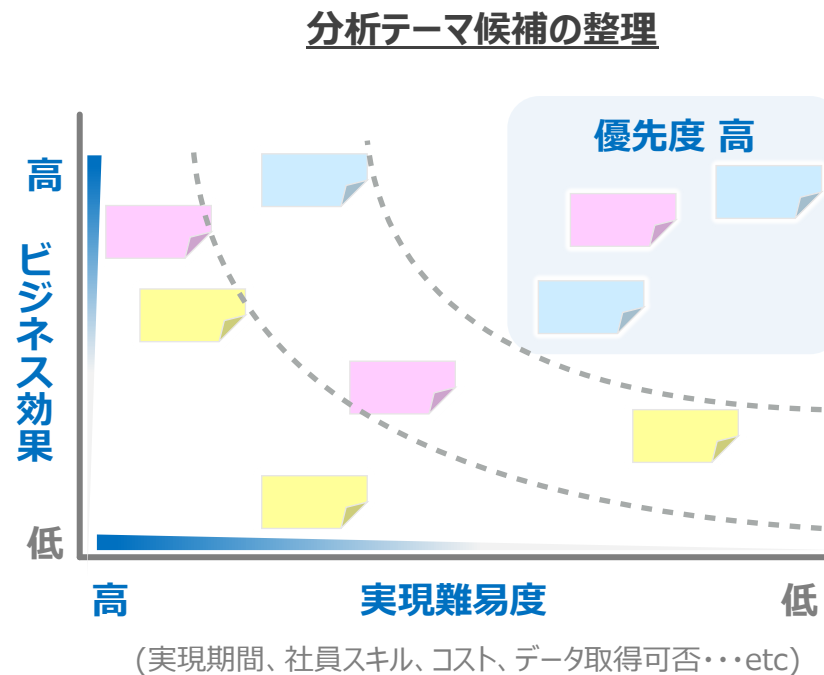
○ Kaggle (米国) : <https://www.kaggle.com/>

○ SIGNATE (日本) : <https://signate.jp>

## ■モチベーション維持や相互学習のために、勉強会などで誰かと一緒に取り組むことも極めて有効です

# 分析課題の選定・具体化

- 社内課題を洗い出し、「ビジネス効果」と「実現難易度」の観点で整理して優先度の高いテーマを選定・具体化する



## 分析テーマの具体化 (例)

分析テーマ	...
テーマ概要	...
使用データ	...
目的変数	...
説明変数	...
ビジネス効果	...
分析ゴール	...

Sample

# 分野を超えた手法の適用

- 同じデータを扱うという意味では分野の壁はなく、様々な分野の最新事例、動向にアンテナを張りながら、積極的に自分の分野に取り入れていくことが重要
- 異分野との交流会は極めて貴重な機会であり、積極的に参加していくべきである

マーケティング領域で用いられる「アソシエーション分析」(教師なし学習手法)

## 個客の購買データ



購買1件  
あたり  
1レコード

## データ構造化

購入ID	牛肉	ワイン	リンゴ	みかん	...
0001	✓	✓	✓		...
0002			✓	✓	...
0003		✓			...
0004	✓	✓	✓	✓	...
0005	✓	✓	✓		...
⋮	⋮	⋮	⋮	⋮	⋮

牛肉とワインが**同時購入**  
される割合が高い

## ルール抽出

大量データの中から  
共起する頻出パターン  
(=ルール)を抽出

条件      帰結  
牛肉 ⇒ ワイン

## 活用例

商品陳列の工夫 (同時購入されやすい製品を隣接させる)

レコメンデーション (「条件」側の製品がカゴに入ったら、「帰結」側の製品も推薦)

## 製造・医療への応用

### 製造

製品ID	割れ・ 欠け	シワ	キズ	異物 混入	...
0001	✓	✓	✓		...
0002		✓	✓	✓	...
0003		✓	✓		...
0004	✓		✓	✓	...
0005		✓	✓		...
⋮	⋮	⋮	⋮	⋮	⋮

シワとキズが**同時発生**  
する割合が高い

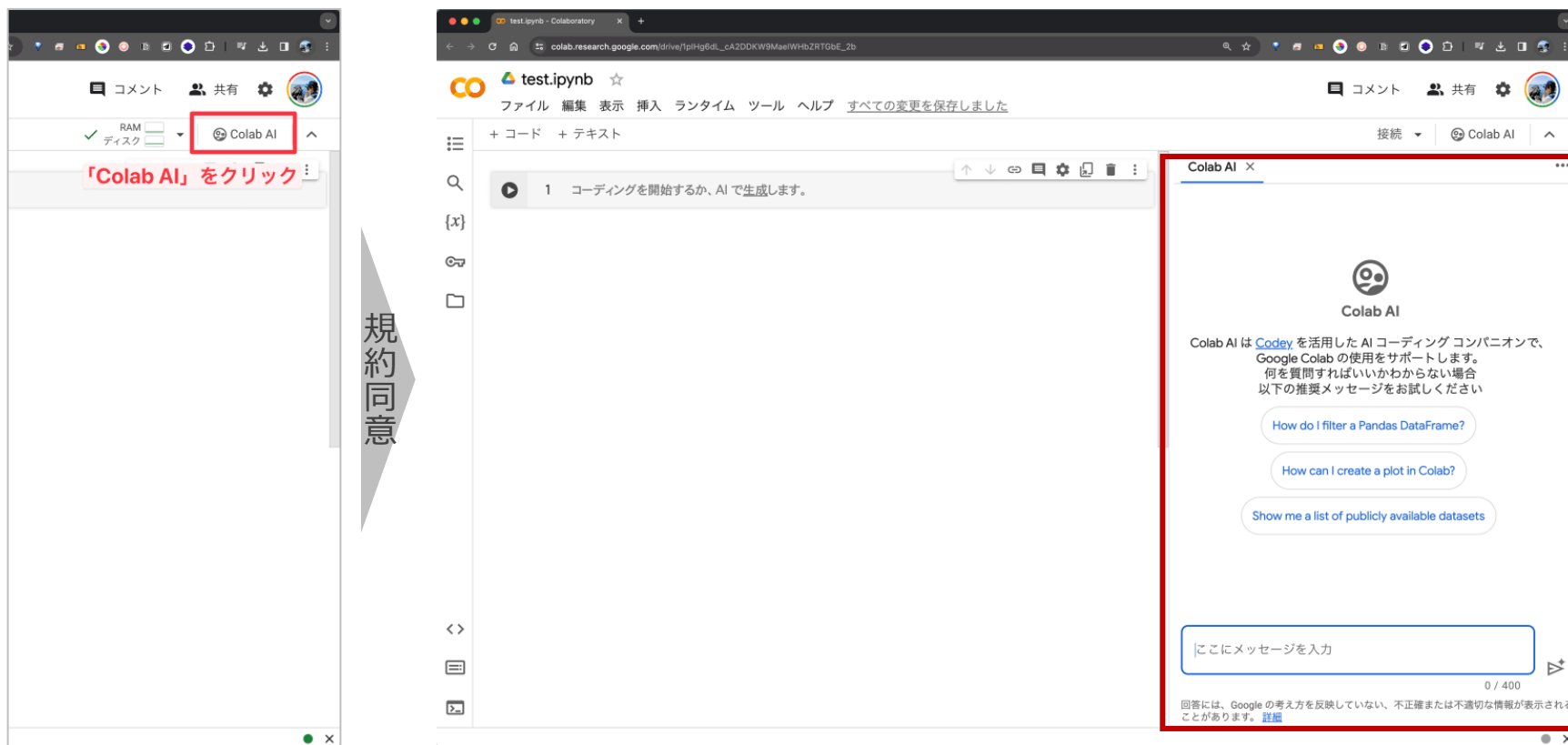
### 医療

症例ID	虚血性 心疾患	糖尿病	高血圧	脂質異 常症	...
0001	✓	✓	✓		...
0002			✓	✓	...
0003		✓			...
0004	✓	✓	✓	✓	...
0005	✓	✓	✓		...
⋮	⋮	⋮	⋮	⋮	⋮

虚血性心疾患と糖尿病が  
**同時発症**する割合が高い

# 参考：Colab AI（生成AIによる支援機能）の活用

- Google Colaboratory では、生成AIによるコード作成支援機能として、“Colab AI” が無料で活用できる
- 簡単なコード作成からデバッグなど、様々な活用でき、初心者でも効率的にコード作成ができる  
ただし、スキルアップ／定着の観点では、初めのうちこそ、できる限り、自力で作成することが望ましい



## ▼参考：プログラミングでの生成AI活用 （前回講演資料抜粋）

活用例①：簡単なコード作成

活用例②：コードのデバッグ（エラー修正）

活用例③：コードレビュー

活用例④：コードの説明

活用例⑤：コードを別の言語に変換

活用例⑥：変数名や関数名を作成

活用例⑦：ダミーデータを作成

**End of File**