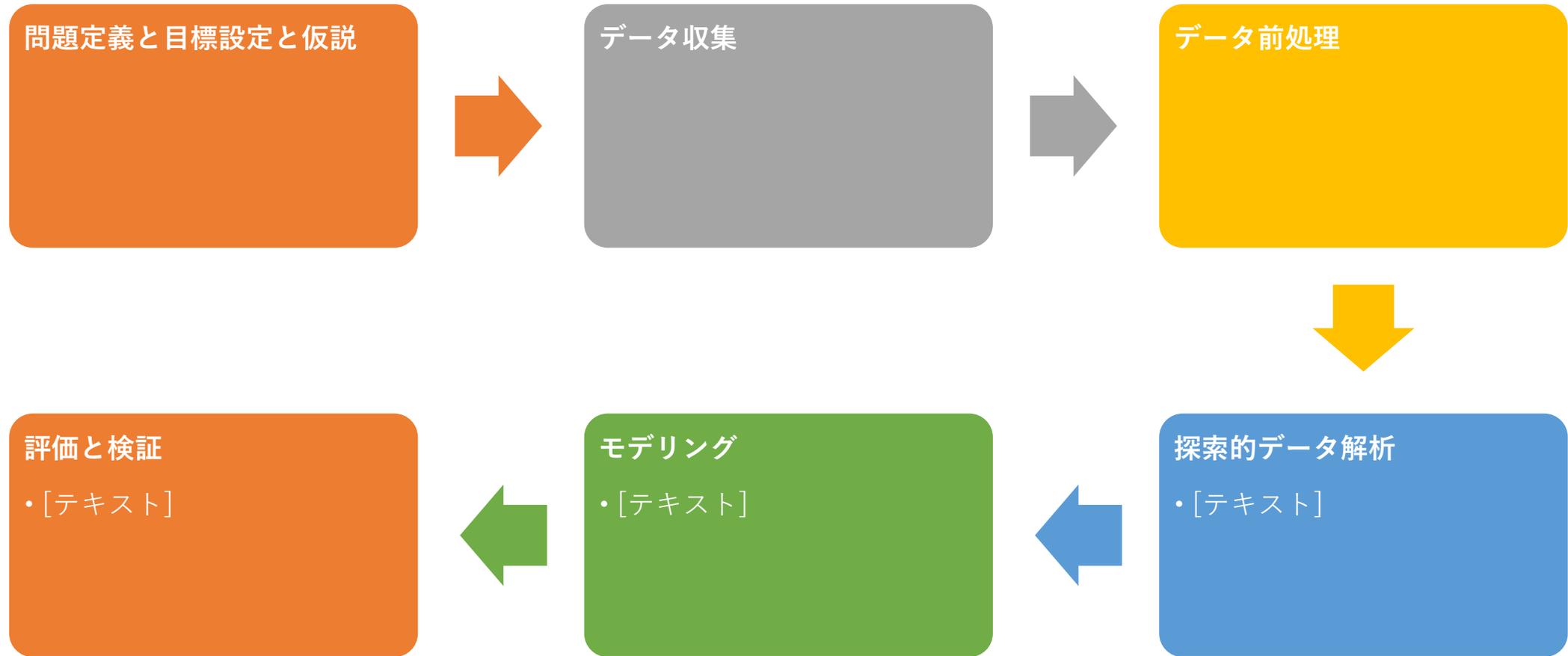


インテグレーションステップ
分析計画書

なぜ、成熟した民主国家は、投票率が
下がるのか??

データサイエンティスト育成講座
木曜クラス
佐藤 順



問題の定義と目標設定 仮説

問題の定義

- 日本の選挙の投票率が低いことが問題にされて、十数年経つがその問題は、現在でも未解決である。
- 国際的な統計でも G 7 の主要国でも選挙の投票率の低下は、指摘されている。

目標設定

- 世界の主要な民主国家の各種統計データと選挙の投票率の関連性を機械学習で分析する。
- モデルの精度や特徴量の重要度をもとに投票率と各種の統計データの関係性を考察する。
- 投票率の低下に影響する統計データについて考察する。

仮説

仮説

- 民主政治の経過年数が増えてくると国民の生活満足度が向上し政治に対する関心が希薄になり投票率が低下する。
- 各国の政治を取り巻く環境は変化するが、国民に対する教育の内容の更新が遅れ、教育レベル（教育の質）が低下し、政治への関心も低下する。

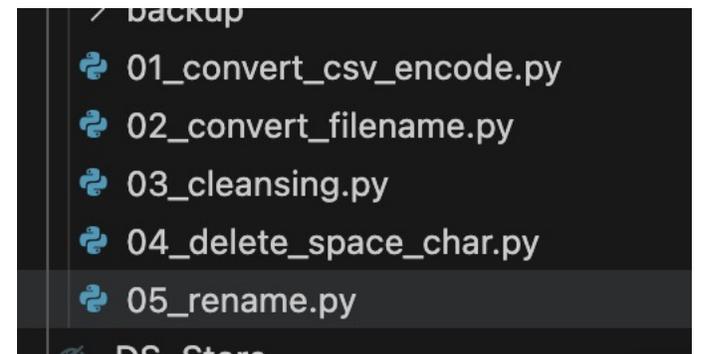
データ収集

- 各国の統計データは、Global Note を利用する。
 - <https://www.globalnote.jp/>

データ前処理

- Global Note のデータは、Shift JIS の CSV を含むテキストファイルとして、ダウンロードできる。

- 以下の処理を行うPythonコードを作成
 - 文字コードの変換 (Shift JIS → UTF-8)
 - 不要な行の削除
 - 国名をファイル名から削除し、ファイル名に付与する。
 - 不要な半角スペースを削除する。
 - ファイル名の変更。



処理<前>のCSVファイル

The screenshot shows a code editor with a file explorer on the left and a code editor on the right. The file explorer shows a directory structure with a folder named 'work3' containing a subfolder 'before_cleansing' with several CSV files. The code editor displays the content of 'data_10-Brazil-2024-01-28.csv'.

```
work3 > before_cleansing > data_10-Brazil-2024-01-28.csv > data
1  名目GDP (国連統計)
2
3  Brazil
4
5  単位: 百万US$
6  出典: 国連
7
8  年, 順位, データ, 注
9  1990, 10, 408823.1376196,
10 1991, 11, 380373.4149149,
11 1992, 12, 364362.7512962,
12 1993, 12, 408861.8974125,
13 1994, 7, 583242.23274,
14 1995, 7, 781735.8664119,
15 1996, 8, 853644.1451109,
16 1997, 8, 885685.1435997,
17 1998, 8, 857857.034833,
18 1999, 11, 596883.5234931,
19 2000, 10, 655448.1494907,
20 2001, 11, 559983.5049806,
21 2002, 13, 509795.1833208,
22 2003, 15, 558233.5953538,
23 2004, 14, 669289.2491361,
24 2005, 12, 891634.0305762,
25 2006, 10, 1107626.7114235,
26 2007, 10, 1397114.2792332,
27 2008, 8, 1695855.3432824,
28 2009, 8, 1666996.1165739,
29 2010, 7, 2208838.1085931,
30 2011, 7, 2616156.6067345,
31 2012, 7, 2465228.2938626,
32 2013, 7, 2472819.362259,
33 2014, 7, 2456043.7660634,
34 2015, 9, 1802211.9995555,
35 2016, 9, 1795693.2658242,
36 2017, 8, 2063514.6887616,
37 2018, 9, 1916933.7084041,
38 2019, 9, 1881458.564748,
39 2020, 12, 1476107.2920365,
40 2021, 13, 1649622.8862183,
41 2022, 11, 1920095.4771642,
42
43 注
44
```

The status bar at the bottom right shows 'Col 4 行 22、列 27 スペース: 4 Shift JIS CRLF CSV'. A red box highlights this area.

処理<後>のCSVファイル

The screenshot shows a CSV editor interface with a file explorer on the left and a data view on the right. The file explorer shows a directory structure with a file named '15歳未満人口比率_data_2900-Brazil-2024-01-28.csv' selected. The data view shows a table with columns for year, age, and ratio. The status bar at the bottom indicates the file is open in UTF-8 encoding with LF line endings. A red box highlights the file name in the status bar.

年	順位	データ	注
1990	121	35.2394638	
1991	121	34.7938483	
1992	122	34.3081434	
1993	122	33.7897202	
1994	121	33.2502261	
1995	121	32.6958761	
1996	121	32.1285771	
1997	122	31.5509646	
1998	123	30.9750765	
1999	123	30.4027486	
2000	126	29.8340145	
2001	125	29.2743124	
2002	126	28.7201692	
2003	125	28.1693617	
2004	124	27.6372709	
2005	124	27.1345482	
2006	124	26.6507789	
2007	124	26.1783587	
2008	121	25.7060475	
2009	122	25.2305634	
2010	122	24.7547629	
2011	124	24.2822404	
2012	124	23.8132776	
2013	128	23.3430589	
2014	129	22.8980912	
2015	129	22.4977463	
2016	130	22.1115988	
2017	130	21.7477747	
2018	129	21.4398249	
2019	129	21.1451272	
2020	132	20.8347709	
2021	133	20.5406675	
2022	133	20.2713094	

探索的データ解析 (EDA)

- Jupyter Notebook での<前>処理
 -

モデリング

- 以下のモデルを作成して分析した。
 - LinearRegression (線形回帰分析)
 - RandomForestRegressor (ランダムフォレスト)
 - LGBMRegressor (Light GBM)

- 国については、以下の 22 カ国を用いた。
 - アルゼンチン、オーストラリア、ブラジル、カナダ、コロンビア
 - コスタリカ、エジプト、フィンランド、ギリシャ、アイスランド
 - インド、インドネシア、アイルランド、日本、ケニヤ、韓国
 - クウェート、ニュージーランド、ノルウェー、スウェーデン
 - イギリス、アメリカ
- 主に民主的に政治が行われていると思える国を選んだ。
 - 中国、北朝鮮、ロシアは、含めなかった。

- 特徴量として、以下の統計量を用いた。
 - 15歳未満人口比率、65歳以上人口比率（高齢化率）
 - ひとり当たり名目GDP（国連統計）、科学論文数（全分野合計）
 - 公的教育費の対GDP比、政治の民主化度
 - 大学進学率（短期大学含む）、名目GDP（国連統計）

- 以下の特徴量は、国によって、統計が存在しないものがあり、使用できなかった。
 - 経常収支、輸入総額、輸出総額、トマトの生産量
 - 教育制度への市民満足度、医療制度への市民満足度
 - 大卒人口比率（25歳-64歳）、15歳-64歳人口比率
 - 高齢者扶養率、年間労働時間（全就業者）

- 以下の特徴量については、対数変換を行なった。
 - 科学論文数（全分野合計）
 - 名目GDP（国連統計）
 - ひとり当たり名目GDP（国連統計）

- 最終的に 96 件のデータを準備することができた。

評価と検証

	LinearRegression	RandomForestRegressor	LightGBMRegressor
RMSE	115.64	73.54	8.89

- LinearRegression

順位	特徴量	Coefficient (影響度)
1位	政治の民主化度	18.27
2位	名目GDP (国連統計)	4.94
3位	公的教育費の対GDP比	0.62
4位	15歳未満人口比率	0.44
5位	大学進学率 (短期大学含む)	0.13
-3位	65歳以上人口比率 (高齢化率)	-0.46
-2位	ひとり当たり名目GDP (国連統計)	-2.30
-1位	科学論文数 (全分野合計)	-5.50

- RandomForestRegressor

順位	特徴量	Importance (重要度)
1位	政治の民主化度	0.42
2位	公的教育費の対GDP比	0.14
3位	大学進学率 (短期大学含む)	0.098
4位	ひとり当たり名目GDP (国連統計)	0.095
5位	科学論文数 (全分野合計)	0.08
6位	15歳未満人口比率	0.06
7位	65歳以上人口比率 (高齢化率)	0.04
8位	名目GDP (国連統計)	0.04

- LightGBMRegressor

順位	特徴量	Importance (重要度)
1位	名目GDP (国連統計)	47
2位	公的教育費の対GDP比	27
3位	ひとり当たり名目GDP (国連統計)	20
4位	政治の民主化度	17
5位	科学論文数 (全分野合計)	16
6位	65歳以上人口比率 (高齢化率)	15
7位	15歳未満人口比率	10
8位	大学進学率 (短期大学含む)	5

結論

- 三つのモデルで共に「公的教育費の対GDP比」が比較的に重要度が高いことがわかった。
- 今回の分析では、公的教育費を高めることで、日本の選挙の投票率を上げることができる。という結論とする。

今後の展望と課題

- より深度を深めるためになぜ公的教育費を高めることが選挙の投票率に影響するのか??という研究が今後の課題かと思う。
- 過去の関連研究に対する文献調査を実施したい。
- もともとの主題と仮説に対する調査が不十分だった。

データミックス教官の評価

- 用いたデータの件数が少なく、対象の国の範囲も狭いため学術的に価値を認めるには、まだ、不十分である。