

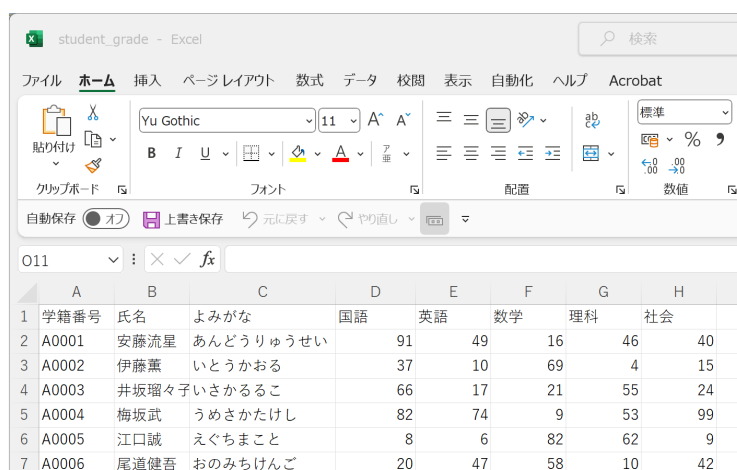
## 第9章

# 基盤モジュール: Pandas

表計算 (table calculation) ソフトウェアは、その名の通り、数値や文字列を含むデータを表にまとめて整理するためのデータ処理ツールです。Microsoft 社が開発した Excel (エクセル) は代表的な表計算ソフトウェアで、長い歴史を誇ります。その他、無料で使用することのできる Google Sheets や、Apache 財団の OpenOffice Calc などがあり、いずれも表を用いたデータ処理やグラフ作成機能を備えています。このように電子化された表計算ソフトウェアのデータは、あらかじめ整理されている分、Python にとって扱いやすいものと言えます。本章では主として Pandas パッケージを用いた Excel ファイルの扱い方を学び、簡単な統計ツールを作っていきます。

### 9.1 Excel と表

今回使用する Excel ファイルは図 9.1 に示すような、個人ごとに 5 科目 (国語, 英語, 数学, 理科, 社会) の得点が記されているものとします。



	A	B	C	D	E	F	G	H
1	学籍番号	氏名	よみがな	国語	英語	数学	理科	社会
2	A0001	安藤流星	あんどうりゅうせい	91	49	16	46	40
3	A0002	伊藤薫	いとうかおる	37	10	69	4	15
4	A0003	井坂瑠々子	いさかると	66	17	21	55	24
5	A0004	梅坂武	うめさかたけし	82	74	9	53	99
6	A0005	江口誠	えぐちまこと	8	6	82	62	9
7	A0006	尾道健吾	おのみちけんご	20	47	58	10	42

図 9.1 今回使用する Excel ファイルの内容

まず、Excelで簡単な集計を行ってみましょう。追加するのは

- ① 個人ごとの平均値 (Excel の average 関数を使用), 中央値 (median 関数), 標準偏差 (stdev.p 関数)
- ② 科目ごとの平均値, 中央値, 標準偏差をグラフ化する
- ③ 科目ごとに上位5名のリストと得点 (棒グラフも含む)

とします。元のデータは別シートを作成し, そちらにすべてのデータを移して作業してみてください。完成版は図 9.2 のようになります。

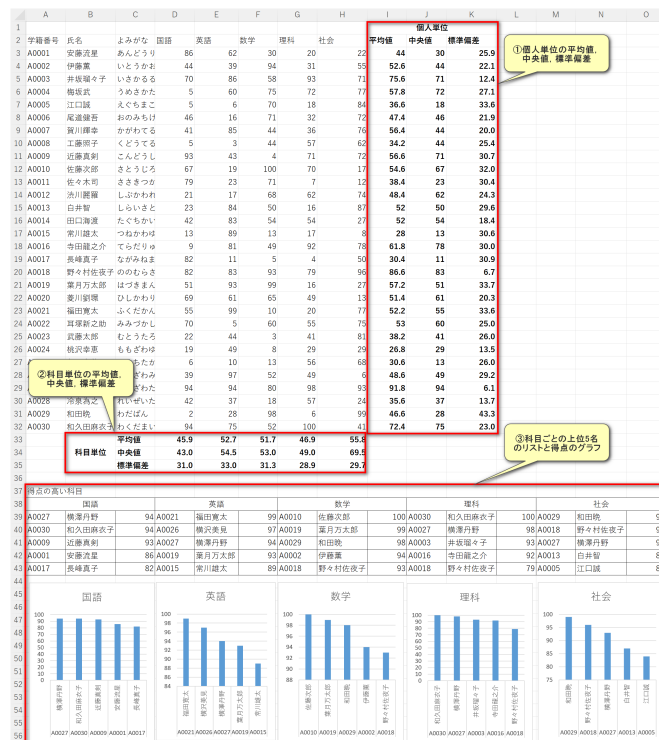


図 9.2 統計処理とグラフを付加した Excel シート

必ずしもそのようなレイアウトで作る必要はありませんので, 自分で創意工夫したレイアウトで作成してみてください。ただし, 次の点は気を付けて下さい。

1. 統計値はせいぜい小数点 1 桁で表示すること。
2. 得点のグラフは必ず最大値を 100 点に統一すること。

Excel だけでもこれらの処理は可能ですが, 例えばこれに個人単位の得点グラフを付けろとか, 個人ごとにレポート形式にせよと言われると, 手作業では手間がかかります。

ということで, ここから先は Python で作業を肩代わりしてみましょう。

### 問題 9.1

図 9.2 を作成し, 次のデータとグラフを追加せよ。

1. 科目別の最小点と最大点
2. 科目別の平均点のグラフ

## 9.2 Pandas を用いた Excel ファイルの読み出し

では図 9.1 に示した Excel ファイル `student_grade.xlsx` を Python から読み出し、図 9.2 に示すような科目別と個人別の平均値を導出してみましょう。ここで作成する Python スクリプトは「`student_grade.py`」とします。

### 9.2.1 Excel ファイルの読み出しと科目ごとの統計値算出・グラフ描画

まず `pandas` と `japanize_matplotlib` モジュールをインストールしておき、下記の部分を作成して「`student_graphde.xlsx`」が読み出せるか、確認して下さい。

ソースコード 9.1 `student_grade.py(1/3)`

```
1 # student_grade.py: 生徒の成績統計計算
2 import numpy as np # NumPy
3 import pandas as pd # Pandas
4 import matplotlib.pyplot as plt # Matplotlib
5 import japanize_matplotlib # Matplotlib 日本語対応
6
7 # Excel ファイル名
8 filename = 'student_grade.xlsx'
9 # Excel シート名
10 sheetname = '学籍番号・氏名・タイトルなし'
11 # 教科名
12 subjects = ['国語', '英語', '数学', '理科', '社会']
13
14 # Excel ファイル読み込み確認
15 try:
16     pd.ExcelFile(filename)
17 except:
18     print(filename, 'が開けませんでした。')
19
20 # Excel ファイル読み込み時のみ動作
21 with pd.ExcelFile(filename) as xls:
22     sheet = pd.read_excel(xls, sheetname)
23     print(sheet)
24
25 # 科目ごとの平均点,標準偏差,最高点,最低点
26 for i in range(len(subjects)):
27     print('{:s}の平均点,標準偏差,最高点,最低点:_{:3.1f},_{:5.3g},_{:4d},_{:4d}'.format
28         (
29         subjects[i],
30         np.average(sheet[subjects[i]]),
31         np.std(sheet[subjects[i]]),
32         np.amax(sheet[subjects[i]]),
33         np.amin(sheet[subjects[i]])
34     ))
```

この部分の動作確認ができれば、科目ごとの平均点を棒グラフ化します。

ソースコード 9.2 student\_grade.py(2)

```
55 # 平均点のグラフ化
56 x = subjects # ['国語', '英語', '数学', '理科', '社会']
57 y = [np.average(sheet[x[i]]) for i in range(5)]
58 fig, ax = plt.subplots()
59 # 棒グラフ
60 ax.set_title('科目別平均点')
61 ax.set_ylabel('得点')
62 ax.bar(x, y)
63 # グラフ表示
64 plt.show()
```

グラフは図 9.3 のように表示されます。

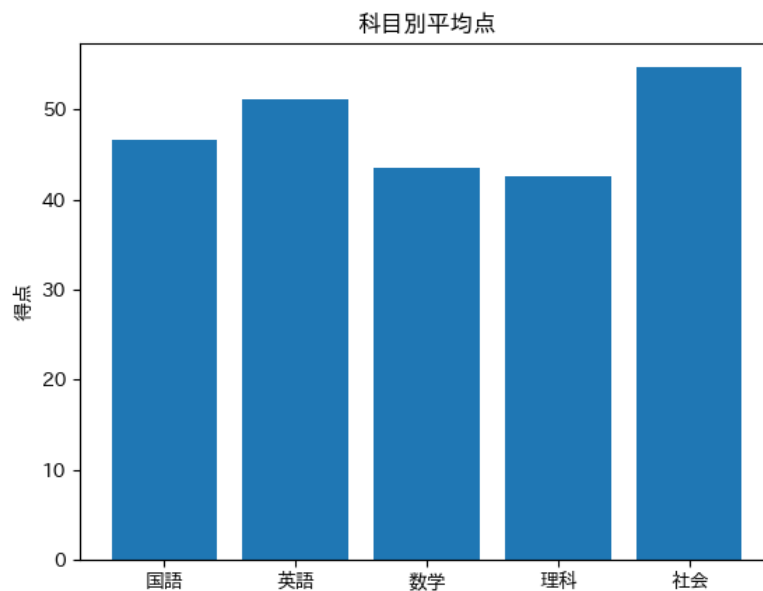


図 9.3 科目別平均点のグラフ

## 問題 9.2

中央値、標準偏差も導出しグラフ化せよ。

### 9.2.2 生徒ごとの統計値の導出

次に、生徒ごとの統計値の導出と、グラフ描画を行います。

ソースコード 9.3 student\_grade.py(2)

```
63 # 生徒ごとの統計値
64 student_name = input('生徒の氏名:')
65 find_row = sheet.loc[sheet['氏名'] == student_name]
66
67 # 生徒が見つかったときのみ動作
```

```

68     if(find_row.empty == False):
69         print(student_name, '→\n', find_row)
70         #print(subjects[0], '::', find_row[subjects[0]].to_numpy()[0])
71         #print(student_name, ' → \n', find_row.to_numpy())
72         #for i in range(len(subjects)):
73         # print(subjects[i], ' → ', find_row[subjects[i]])
74
75         #find_row_scores = find_row.loc[student_name, subjects[0]]
76         find_row_scores = [find_row[subjects[i]].to_numpy()[0] for i in range(len(
77             subjects))]
78         print(find_row_scores)
79
80         # 個人成績のグラフ化
81         x = subjects # ['国語', '英語', '数学', '理科', '社会']
82         y = find_row_scores # [find_row[x[i]] for i in range(5)]
83         #print(y)
84         fig, ax = plt.subplots()
85         # 棒グラフ
86         ax.set_title(student_name + 'の得点')
87         ax.set_ylabel('得点')
88         ax.bar(x, y)
89         # グラフ表示
90         plt.show()
91
92         print(student_name, 'の平均点,標準偏差,最高点,最低点:',
93             np.average(find_row_scores),
94             np.std(find_row_scores),
95             np.amax(find_row_scores),
96             np.amin(find_row_scores)
97         )
98     else:
99         print('名前:', student_name, 'が見つかりません。')

```

個人別の得点グラフは図 9.4 ようになります。

### 問題 9.3

次の処理を追加せよ。

1. 科目ごとの上位 5 名のリストアップと棒グラフ化
2. 科目ごとの下位 5 名のリストアップと棒グラフ化

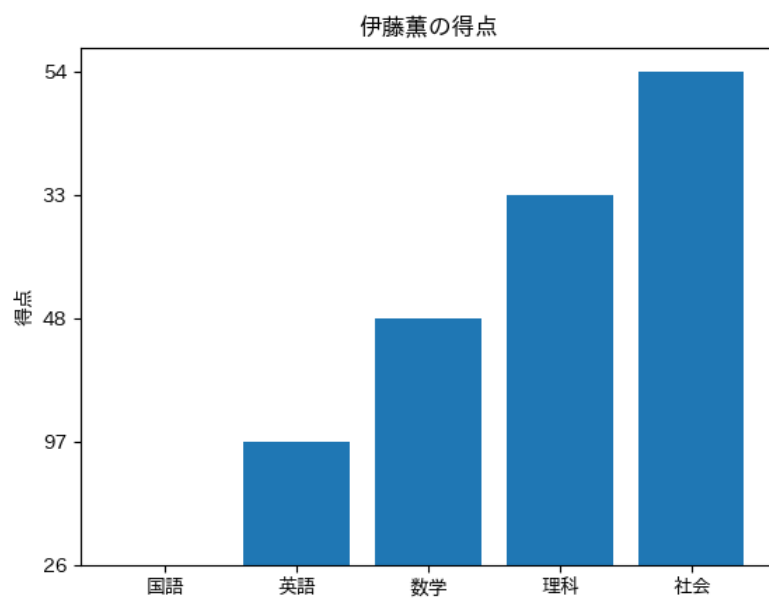


図 9.4 個人別得点グラフ