

```
1 import numpy as np
2 import scipy.stats as stats
3 import pandas as pd
4 # import matplotlib.pyplot as plt
5 # import seaborn as sns
```

▼ モジュールpandasについて

データを集めて分析をするには、収集したデータの整理が必要になる。モジュールpandasは、データを表にして管理したり、表からさまざまな項目のデータを抽出したりするのに利用される。

次のコードを実行してから、先の内容を読み進めること。

```
1 ser1 = pd.Series([66,76,82,63,59,52,87,61,73,68,55,67,64])
2
3 data1 = {"Number":["4401","4402","4403","4404","4405","4406","4407"],
4         "Name":["Akao","Ozawa","Kurokawa","Tominaga","Negoro","Yano","Yokota"],
5         "Math":[66,76,82,63,59,52,87],
6         "English":[61,73,68,55,67,64,74],
7         "Birth_Month":["February","June","August","November","April","October"]
8
9 df1 = pd.DataFrame(data1)
```

▼ Seriesオブジェクト

seriesオブジェクトは、1列の表を扱う、配列のようなものである。次のコードによって、上で定義したser1を出力してみよ。

また、メソッド.values、.indexをser1につけると、出力がどのように変化するか、コメントアウトを外して確認せよ。

```
1 # 上のセルから、ser1というSeriesオブジェクトについて出力する
2 #ser1
3 ser1.values
4 # ser1.index

array([66, 76, 82, 63, 59, 52, 87, 61, 73, 68, 55, 67, 64])
```

▼ DataFrameオブジェクト

2次元配列を扱うには、DataFrameオブジェクトを利用する。

```
1 # 上のセルから、df1というDataFrameオブジェクトを出力する
2 df1
3 # df1.T # 行と列を入れ替える(Tは転置, Transposeの意)
```

```

4 # df1.Birth_Month # 特定の1列のみを取り出す場合の書き方
5 # df1[["Math","Birth_Month"]] # 複数の列を取り出す場合の書き方
6 # df1["English"][5] # English列の第5行を取り出す場合

```

	Number	Name	Math	English	Birth_Month
0	4401	Akao	66	61	February
1	4402	Ozawa	76	73	June
2	4403	Kurokawa	82	68	August
3	4404	Tominaga	63	55	November
4	4405	Negoro	59	67	April
5	4406	Yano	52	64	October
6	4407	Yokota	87	74	June

▼ 条件に合う行の抽出

データフレームから、条件に合う行を抽出するには、以下のセルのように条件式を書けばよい。

```

1 # df1[df1["Birth_Month"] == "June"] # Birth_Month列がJuneであるデータのみを抽出
2 df1[df1["Math"] >= 60] # Mathの成績が60以上であるデータのみを抽出

```

	Number	Name	Math	English	Birth_Month
0	4401	Akao	66	61	February
1	4402	Ozawa	76	73	June
2	4403	Kurokawa	82	68	August
3	4404	Tominaga	63	55	November
6	4407	Yokota	87	74	June

▼ 行のソート(並び替え)

データフレームの特定の列を数値の順に並び替えるには、`sort.values()` メソッドを利用して、以下のセルのように書けばよい。詳しくは、[こちら](#)などを参照。

```

1 df1.sort_values("English") # 英語の成績で昇順にソート
2 # df1.sort_values("English",ascending=False) #降順にソートしたい場合は, ascending=False

```

	Number	Name	Math	English	Birth_Month
3	4404	Tominaga	63	55	November
0	4401	Akao	66	61	February
5	4406	Yano	52	64	October
4	4405	Negoro	59	67	April



▼ 外部ファイルからデータを取得する方法

以下では、Notebookファイルと同一のフォルダに配置されたExcel, [CSV](#)のデータを読み込む方法, および, ColaboratoryのNotebookファイルと同一のフォルダに配置されたSpreadsheetデータの読み込み方について解説する。

以下では、各自利用している環境に合わせて、実行するコードを選ぶこと。

▼ Jupyter Notebookのローカル環境下での作業用

CSV形式のファイルからデータを読み込む場合は、このまま下のセルを実行する。

Excel形式のファイルからデータを読み込む場合は、下のセルの8 ~ 9行目をコメントアウト(行頭に#を加える)し、2 ~ 5のコメントアウトを外す。

```

1 # Excelファイルから読み込む場合
2 # excel_file_name = "DogOrCat_Test.xlsx" # 読み込むファイル名を文字列として入力(拡張子を
3 # xl_data = pd.ExcelFile(excel_file_name) # データをxl_dataに格納
4 # sheets = xl_data.sheet_names # Excelのシートのリストを作成
5 # df = xl_data.parse(sheets[0]) # 0番目のシートのデータをDataFrame形式にしてdfに格納
6
7 # CSV形式のファイルから読み込む場合
8 csv_file_name = "DogOrCatOr_4-3.csv" # 読み込むファイル名を文字列として, csv_file_nameに
9 df = pd.read_csv(csv_file_name) # データをDataFrame形式にしてdfに格納

```

▼ Google Colab内での作業用

Google Colaboratoryを使用する場合は、Google Drive内に置いたスプレッドシートからデータを読み込む必要がある。

基本的には、スプレッドシートのURLを下のコードの所定の場所に貼り付け、順にコードを実行すればよい。詳細は、この[Qiitaの記事](#)などを参考にされたい。

```

1 # 必要なモジュールのインストール
2 !pip install --upgrade -q gspread

```

```

1 # 必要なモジュールのインポート
2 from google.colab import auth

```

```

3 from oauth2client.client import GoogleCredentials
4 import gspread
5
6 # GoogleDriveをColaboratoryにロードする(認証作業あり)
7 auth.authenticate_user()
8 gc = gspread.authorize(GoogleCredentials.get_application_default())
9
10 # 文字列urlに読み込むスプレッドシートのURLを入力
11 url = "https://docs.google.com/spreadsheets/d/1ym4x88o5221PyvCLuv-SmB5D4o3xa1afi
12
13 sheet = gc.open_by_url(url).get_worksheet(0) # urlに指定したURLから、スプレッドシートのC
14 df = pd.DataFrame(sheet.get_all_values()) # sheetのデータをDataFrame形式にしてdfに格
15 df

```

	0	1
0	RandNum	DogOrCatOr
1	0.002127418695	「犬派」
2	0.0116918293	中立
3	0.02330211606	中立
4	0.03087775439	「猫派」
...
110	0.9854067944	中立
111	0.9866553689	「猫派」
112	0.9869737949	「犬派」
113	0.9903752154	「犬派」
114	0.9928164818	「犬派」

115 rows × 2 columns



▼ 模擬アンケート結果の分析

先週採ったアンケートの結果の分析を行おう。

問題設定

「模擬アンケート①」に回答した10回生を母集団とし、回答者から無作為に抽出した20人の回答結果から、母集団における「犬派」と回答した人数を、信頼度95%の信頼区間を求めることにより、推測する。

1. 抽出するサンプル数 n を20に設定する。
2. 信頼度を設定する。ここでは α を、信頼区間が母集団における「犬派」と回答した人数を推定し損なう確率として、信頼度 $(100(1 - \alpha)\%)$ を間接的に設定する。

分析

1. 回答結果を収めたデータフレーム(`df` , 上のセルで定義したもの)の行数(総回答数)を L とする.
2. リスト $[1, L]$ (1以上 L 以下の整数値)から20個(n に格納した値)の一様乱数(どの整数も等確率で選ばれる, 一様分布に従う乱数)を非復元抽出(`randlist`)し, 抽出した乱数に対応する回答のみを含むデータフレームを構成する(`df.iloc[randlist]`).
3. `df.iloc` に含まれる「犬派」の回答件数を計上し, `n_dog` に格納する. このとき, 抽出した標本における「犬派」の割合 `p_dog` は

$$p_dog = n_dog / n$$

と計算される.

4. 母集団における「犬派」の割合(母比率)を, 標本における「犬派」の割合 `p_dog` から, 信頼度を95%として, 以下の考え方に従って区間推定せよ.

1. 母集団(大きさ $L =$ 全回答数)のうち「犬派」が k_{dog} 件であるとする. 母比率は $\frac{k_{\text{dog}}}{L}$ である.
2. そこから無作為に大きさ $n = 20$ の標本を(非復元)抽出する. 標本に含まれる「犬派」の割合 $p_{\text{dog}} = \frac{n_{\text{dog}}}{n}$ は, 母比率の推定量となる.
3. 母集団から大きさ $n = 20$ の標本を無作為抽出するとき, 標本の中に含まれる「犬派」の割合は, 二項分布 $B(1, p_{\text{dog}})$ に従うとみなせる. この二項分布を近似する正規分布 $N(1 \cdot p_{\text{dog}}, 1 \cdot p_{\text{dog}} \cdot (1 - p_{\text{dog}}))$ から大きさ $n = 20$ の標本を抽出すると考えて, 上下2.5%点の値($z_{0.025} = 1.96$)を求めることで, 信頼度95%の信頼区間

$$\left[p_{\text{dog}} - z_{0.025} \cdot \sqrt{\frac{p_{\text{dog}} \cdot (1 - p_{\text{dog}})}{20}}, \quad p_{\text{dog}} + z_{0.025} \cdot \sqrt{\frac{p_{\text{dog}} \cdot (1 - p_{\text{dog}})}{20}} \right]$$

の計算式に代入する.

分析結果の集約

1. クラスルームに配布されている, フォーム「114_Estimation_Interval_test_4-3」に, 上下側信頼限界を入力する.

コードのあとに 母比率の推定に関する解説を載せている. 必要に応じて参照してほしい.

```
1 # 抽出するサンプル数を指定する
2 n = 20
3
4 # 信頼度の設定(1-aが信頼度, つまり, aは閾値を超える辺境の確率を定めている)
5 a = 0.05
6
7 # 標準正規分布の累積分布関数の値がa/2となるxの値をパーセント点から求めている(累積分布関数の逆関数)
8 z_ppv = stats.norm.ppf(1 - a / 2, loc=0, scale=1)
9 t_ppv = sp.stats.t.ppf(1 - a / 2, n - 1)
10
```

```

11 print("標準正規分布に基づく",100*(1-a/2),"%点の値は:",z_ppv)
12 # print("自由度",n-1,"のティー分布に基づく",100*(1-a/2),"%点の値は:",t_ppv)
13
14 L = len(df) # DataFrameの行数の取得
15 ser = np.arange(1,L+1) # リスト[1,L]
16 randlist = np.random.choice(ser,n,replace=False) # N個の乱数のリストを生成
17
18 # 回答結果のDataFrameから上の乱数リストのindexを抜き出し、標本の中で「犬派」と回答している件数をn
19 print("抽出された標本の回答結果は以下の通り:\n",df.iloc[randlist,1])
20 n_dog = len(df.iloc[randlist,1][df.iloc[randlist,1] == "「犬派」"])
21 p_dog = n_dog / n
22 print("「犬派」の回答件数は, ",n_dog,"件であり、その標本に占める割合は, ",p_dog)
23
24
25 # 標本平均の信頼度(1-a)の上側信頼限界と下側信頼限界を求めている
26 I_dogL = p_dog - z_ppv * np.sqrt(p_dog * (1 - p_dog) / n)
27 I_dogU = p_dog + z_ppv * np.sqrt(p_dog * (1 - p_dog) / n)
28 print(100*(1-a),"%信頼区間は",[I_dogL,I_dogU])
29
30 # I_dog_tL = p_dog - t_ppv * np.sqrt(p_dog * (1 - p_dog) / n)
31 # I_dog_tU = p_dog + t_ppv * np.sqrt(p_dog * (1 - p_dog) / n)
32 # print("ティー分布に基づいて計算した",100*(1-a),"%信頼区間は",[I_dog_tL,I_dog_tU])

--NORMAL--

標準正規分布に基づく 97.5 %点の値は: 1.959963984540054
抽出された標本の回答結果は以下の通り:
66      「猫派」
70      「犬派」
90      「犬派」
10      「犬派」
60      「犬派」
43      「猫派」
26      「犬派」
84      「犬派」
11      「犬派」
99      「犬派」
41      「犬派」
28      「犬派」
68      「猫派」
36      「犬派」
85      「犬派」
38      「犬派」
51      「犬派」
95      中立
25      「犬派」
92      中立
Name: 1, dtype: object
「犬派」の回答件数は, 15 件であり、その標本に占める割合は, 0.75
95.0 %信頼区間は [0.5602273032177509, 0.9397726967822491]

```

▼ 一般論(母比率の推定)

1. 母集団 U における、性質 A をもつ個体の含まれる割合を調べたい。大きさ n の標本を無作為抽出すると、性質 A をもつ個体が k 個あった。このことから、母集団における性質 A をもつ個

体の比率(母比率)の信頼度 $100(1 - \alpha)\%$ の信頼区間は

$$\left[\frac{k}{n} - z_{\frac{\alpha}{2}} \cdot \sqrt{\frac{\frac{k}{n} \cdot \left(1 - \frac{k}{n}\right)}{n}}, \quad \frac{k}{n} + z_{\frac{\alpha}{2}} \cdot \sqrt{\frac{\frac{k}{n} \cdot \left(1 - \frac{k}{n}\right)}{n}} \right]$$

で与えられる. 教科書では, $p_0 = \frac{k}{n}$ として表記している.

2. 厳密には, 母分散が未知であるから, ここで推定するためのパーセント点として正規分布のものではなく, 母分散を標本不偏分散に置き換えた, 自由度 $n - 1$ のティー分布のパーセント点を利用することになる.

$$\left[\frac{k}{n} - t_{\frac{\alpha}{2}}(n - 1) \cdot \sqrt{\frac{\frac{k}{n} \cdot \left(1 - \frac{k}{n}\right)}{n}}, \quad \frac{k}{n} + t_{\frac{\alpha}{2}}(n - 1) \cdot \sqrt{\frac{\frac{k}{n} \cdot \left(1 - \frac{k}{n}\right)}{n}} \right]$$

ティー分布については教科書で触れられていないが, KP等で統計的処理を行い, 推定を行う際には知っておくべき存在である.

【解説】母比率の推定

十分な大きさの母集団 U において, 性質 A をもつ個体の含まれる割合が p_0 であるとする. このとき, 母集団 U から無作為に1個の個体を抽出して, 性質 A をもつ確率は p_0 である. 確率変数 X を, 抽出した個体が性質 A をもつとき $X = 1$, もたないとき $X = 0$ で定義すると, X は二項分布 $B(1, p_0)$ に従う: $X \sim B(1, p_0)$.

母集団 U から個体を n 個抽出し, 大きさ n の標本 $[s_1, s_2, \dots, s_n]$ を作る. s_i の性質の有無を確率変数 X_i で表すことにする. $X_i \in B(1, p_0)$ であり, (二項分布は再生性をもつので)

$$\sum_{i=1}^n X_i = X_1 + X_2 + \dots + X_n \sim B(n, p_0).$$

$B(n, p_0)$ を, (平均と分散を等しくする)正規分布 $N(np_0, np_0(1 - p_0))$ に置き換える:

$$\sum_{i=1}^n X_i \sim N(np_0, np_0(1 - p_0))$$

この結果から, $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ の平均と分散は

$$\begin{aligned} E(\bar{X}) &= E\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n} E\left(\sum_{i=1}^n X_i\right) = \frac{1}{n} \cdot np_0 \\ &= p_0, \end{aligned}$$

$$\begin{aligned} V(\bar{X}) &= V\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} V\left(\sum_{i=1}^n X_i\right) = \frac{1}{n^2} \cdot np_0(1 - p_0) \\ &= \frac{p_0(1 - p_0)}{n}. \end{aligned}$$

ゆえに, $\bar{X} \sim N\left(p_0, \frac{p_0(1 - p_0)}{n}\right)$ である.

あとの流れは, 区間推定の求め方に従う.

▼ <本授業の学び>

本授業で学んだことを, 下のテキストボックスに記入して下さい.

(ここに本授業の学びを記入する)