

コンピュータリテラシ発展 ～Pythonを学ぶ～

第11回： 計算やデータ 析をやってみよう

情 学 情 学科 情 ディア 攻
清水 也 (shimizu@info.shonan-it.ac.jp)

今回の授業内容

今回の授業内容

- 前回の課題解
- データの分析
- データの可視化
- 課題

前回の課題解説

前回の課題解説

- 前回の課題の解法を示します
- 解法について質問があればご連絡ください

解答

https://colab.research.google.com/drive/1BN8_-3cvnWOEKBMtSmVQ-gVyd1qjLbTL?usp=sharing

データの分析

データの準備

- これから行う分析のためにデータを準備します
- Moodleにある「**data_analysis.zip**」をダウンロードして解凍してください
- Google Driveの作業場所に「**data_analysis**」フォルダごとアップロードします
- 「**data_analysis**」の身は次のとおりです
 - 「customer.csv」
 - 「item.csv」
 - 「transaction_1.csv」
 - 「transaction_2.csv」

データの読み込み

目標：商品の購買データと顧客データを組み合わせ、購買分析を行います

- pandasを使ってデータを^タ読み込み、顧客ID上位5名の顧客情^タを示します
- `read_csv()` 関数：CSVファイルをpandasにデータを^タ読み込みます
- 読み込んだデータはDataFrame オブジェクトとして扱うことができます
- DataFrame オブジェクトの `head()` メソッドでデータの先頭5行を表示します

```
import pandas as pd # pandasモジュールをインポート

# CSVファイルのパスを指定
path = '/content/drive/MyDrive/???/data_analysiss/customer.csv'

# 指定されたCSVファイルを読み込み、データフレームに格納
customer = pd.read_csv(path)

# データフレームの最初の5行を表示
customer.head()
```


データの読み込み

- データセットの読み込み関数
- pandas は各種データセットを読み込む関数が用意されています

データ形式	関数	解説
CSV	<code>read_csv()</code>	区切り文字で区切られたデータを読み込む
excel	<code>read_excel()</code>	Excel 形式のデータ（.xls, .xlsx）を読み込む
json	<code>read_json()</code>	JSONデータを読み込む
html	<code>read_html()</code>	HTMLファイル内のデータを読み込む

データの読み込み

- データの基礎情報を確認します
- pandasには、データの基礎情報を簡単に示すプロパティも提供されています
- DataFrame オブジェクトのプロパティを下の表にまとめます

オプション	解説
<code>shape</code>	行数と列数を示する（例: 15行6列の場合, (15,6)と出力）
<code>columns</code>	列名を示す
<code>dtypes</code>	列名と、各列のデータ型を示す

データの結合

- 下の2つの結合方法について解説します
 - データを縦方向に結合する処理
 - 列の数を一定にして横方向にデータを結合する処理

データの結合

データを縦 向に結合

- `concat()` 関数を使ってデー を縦に 結します
 - **transaction_1.csv**と**transaction_2.csv**は同じ種類のデー を扱っています
 - デー 分析をする際に1つのデー となっている^トが扱いやすい
 - この2つのデー を縦に 結します

データの結合

- `concat()` 関数 データを縦に 結合する
 - `ignore_index` を `True` に 設定 すること , 結合 後のデータ インデックス を 視し, 新しくデータ インデックス を作成します
 - `ignore_index` を `True` に 設定 しない場合 (省略), 元データ のデータ インデックス をそのまま使用します

```
path = '/content/drive/MyDrive/???.data_analysi s/'
```

```
# 指定されたCSVファイルを読み込み、データフレームに格納
```

```
transaction_1 = pd.read_csv(path + 'transaction_1.csv')
```

```
transaction_2 = pd.read_csv(path + 'transaction_2.csv')
```

```
# 2つのデータフレームを結合し、新しいデータフレームに格納
```

```
# ignore_index=True は、結合後のデータフレームのインデックスを再設定することを示す
```

```
transaction = pd.concat([transaction_1, transaction_2], ignore_index=True)
```

```
transaction
```

データの結合

特定の列をキーにして横方向にデータを結合

- `merge()` 関数 の 1 をもとにデータを横に結合する
 - transactionデータ は顧客IDがある
 - customerデータ は顧客IDにくく顧客情報（顧客名, 性別, 年齢）がある
 - 顧客IDをキーとして結合処理を行う

データの結合

- `merge()` 関数 の ーをもとにデー を横に結合します
 - 引数 `on` に ー（顧客ID）を指定します
 - `concat()` 関数とデー の指定方法は異なるので注意

```
# transactionデータフレームとcustomerデータフレームを '顧客ID' 列でマージ
# '顧客ID' 列の値が一致する行を結合して、新しいデータフレームを作成
sales_data = pd.merge(transaction, customer, on='顧客ID')

# マージされたデータフレームの内容を表示
sales_data
```

データの集計

- 合計を求める
 - データの集計作業を行う
 - 顧客IDごとの合計購入金額と合計購入数を求める
 - `merge` した `sales_data` データを使う
 - `groupby()` を使ってグループを指定する
 - 合計購入金額 → 「購入金額」
 - 合計購入数 → 「購入数」
 - 合計 → `.sum()`

データの集計

- 合計を求める
 - データの集計作業を行う
 - 顧客IDごとの合計購買金額と合計購入数を求める

```
# sales_data データフレームを '顧客ID' 列でグループ化し、  
# '購買金額' と '購入数' 列の合計を計算  
sales_per_customer = sales_data.groupby('顧客ID')[['購買金額', '購入数']].sum()  
  
# 顧客ごとの合計購買金額と合計購入数を含むデータフレームの内容を表示  
sales_per_customer
```

データの集計

- 合計を求める - データの集計作業を行う
 - `groupby()` データをグループ化された結果を `GroupBy` オブジェクトと呼ぶ。
 - `GroupBy` オブジェクトの各グループを以下にまとめる

メソッド名	解説
<code>count()</code>	グループ化されたデータの個数を示す
<code>mean()</code>	グループ化されたデータの平均を示す
<code>sum()</code>	グループ化されたデータの総和を示す
<code>describe()</code>	グループ化されたデータの統計情を示す（：最大、標準偏差）

データの集計

- 平均を求める
 - データの集計作業を行う
 - `mean()` を使って「購入日」ごとの平均売上を求める

```
# sales_data データフレームを '購入日' 列でグループ化し、  
# '購入金額' 列の平均を計算  
sales_per_day = sales_data.groupby('購入日').購入金額.mean()  
  
# 日ごとの平均購入金額を含むデータフレームの内容を表示  
sales_per_day
```

データの可視化

データを可視化するためのライブラリ

- 集計データを使ってグラフを作成します
- 使用するライブラリ
 - [Matplotlib](#)
 - Pythonでグラフを描くための基本的なライブラリ
 - グラフの基本的な設定を行う
 - [seaborn](#)
 - Matplotlibをベースにしたより複雑な可視化を簡単に行うことができるライブラリ

データを可視化するためのライブラリ

- 集計データを使ってグラフを作成します
- **Matplotlib**も**seaborn**もColabにはインストール済みです
- インポートする必要がある場合があります

```
# matplotlibモジュールのpyplotサブモジュールをインポート（グラフ描画のためのツール）  
from matplotlib import pyplot as plt
```

```
# seabornモジュールをインポート（データの可視化を行うための高レベルなインターフェース）  
import seaborn as sns
```

データを可視化するためのライブラリ

- 集計データを使ってグラフを作成します
- 日本語フォントの使用について
 - Matplotlibは日本語に対応していないので、日本語を表示するためには日本語フォントを指定する必要があります
 - `japanize_matplotlib` ライブラリをインストールしてインポートします

```
# pipコマンドを使って japanize_matplotlib パッケージをインストール
!pip install japanize_matplotlib
```

```
# japanize_matplotlib パッケージをインポート
import japanize_matplotlib
```

データを可視化するためのライブラリ

- 集計データを使ってグラフを作成する
- `sales_per_customer`データを棒グラフで表示します
- `barplot()`関数を使用します
- 引数「**data**」に対象のDataFrameオブジェクトを指定する

```
# 顧客ごとの合計購買金額を棒グラフで表示  
# x軸に顧客ID、y軸に購買金額を指定  
ax = sns.barplot(x=sales_per_customer.index, y='購買金額', data=sales_per_customer)
```


データを可視化するためのライブラリ

- 集計データを使ってグラフを作成する
- `sales_per_day`データを折れ線グラフで表示する
- `lineplot()`関数を使用する
- 引数「**data**」に対象のDataFrameを指定する

```
# 日ごとの平均購買金額を折れ線グラフで表示  
# data引数にsales_per_dayデータフレームを指定  
ax = sns.lineplot(data=sales_per_day)
```

データを可視化するためのライブラリ

- ラフのサイズを整する
 - x軸のデータが重なってしまい見づらいの → ラフを横に広くします
 - **Matplotlib**の設定 修正可能
 - `figure()` 関数の `figsize` 引数を改 する

```
# グラフのサイズを指定 (幅18、高さ4)
plt.figure(figsize=(18, 4))

# 日ごとの平均購買金額を折れ線グラフで表示
# data引数にsales_per_dayデータフレームを指定
ax = sns.lineplot(data=sales_per_day)
```

データを可視化するためのライブラリ

- ラフの 示を設
 - `set_title()` ト: ラフ イト を
 - `set()` ト: 引数 `xlabel` , `ylabel` を してx軸, y軸のラ を作成

```
# グラフのサイズを指定 (幅18、高さ4)
plt.figure(figsize=(18, 4))

# 日ごとの平均購買金額を折れ線グラフで表示
# data引数にsales_per_dayデータフレームを指定
ax = sns.lineplot(data=sales_per_day)

# グラフのタイトルを設定
ax.set_title(' 購買日ごとの平均売上 ')

# x軸とy軸のラベルを設定
ax.set(xlabel=' 購買日', ylabel=' 平均売上 (円) ')
```

課題

課題

- Moodleにある「SCfCL-11-prac.ipynb」ファイルをアップロードしてColabにアップロードしてください
- 課題が完了したら「File」>「Download」>「Download .ipynb」→「.ipynb」形式をアップロードしてください
- アップロードした**.ipynb**ファイルをMoodleに提出してください
- 出期日は**7月4日(木) 20時まで**です