

同様に、「item\_price」についても確認してみましょう。

■図2-4：データの揺れ(商品金額)

```
In [4]: uriage_data["item_price"].head()

Out[4]: 0    100.0
        1     NaN
        2     NaN
        3   2600.0
        4     NaN
        Name: item_price, dtype: float64
```

欠損値「NaN」がデータとして確認できます。このような状態をデータ欠損といい、欠損値をどのように補完するかが今後のデータ分析に影響します。

このように集計対象のデータに揺れや欠損値が存在していると正しい集計が得られません。

試しに、このまま集計を行ってみましょう。

### ⑩ ノック13： データに揺れがあるまま集計してみよう

データの揺れがどれくらい集計に影響するかを確認する事で、いかにデータの整合性が重要か分かります。

まずは「売上履歴」から商品ごとの月売上合計を集計してみましょう。

```
uriage_data["purchase_date"] = pd.to_datetime(uriage_data["purchase_date"])

uriage_data["purchase_month"] = uriage_data["purchase_date"].dt.strftime("%Y%m")

res = uriage_data.pivot_table(index="purchase_month", columns="item_name",
aggfunc="size", fill_value=0)

res
```

ノック13：データに揺れがあるまま集計してみよう

■図2-5：データ補正前の集計結果(商品毎)

ノック13：データに揺れがあるまま集計しよう

```
In [0]: uriage_data["purchase_date"] = pd.to_datetime(uriage_data["purchase_date"])
uriage_data["purchase_month"] = uriage_data["purchase_date"].dt.strftime("%Y%m")
res = uriage_data.pivot_table(index="purchase_month", columns="item_name", aggfunc="size", fill_value=0)

Out[0]:
```

item_name	商品 W	商品 n	商品 E	商品 M	商品 P	商品 S	商品 W	商品 X	商品 O	商品 Q	...	商品 k	商品 l	商品 o	商品 p	商品 r	商品 s	商品 t	商品 v	商品 x	商品 y
purchase_month																					
201901	0	1	0	0	0	0	0	0	0	0	...	1	1	1	0	0	0	0	0	0	0
201902	0	0	0	0	0	0	0	1	0	0	...	0	0	0	0	0	1	1	1	0	0
201903	0	0	1	1	1	0	0	0	0	0	...	0	0	0	0	0	0	0	0	0	0
201904	1	0	0	0	0	0	0	0	1	...		0	0	0	0	0	1	0	0	0	0
201905	0	0	0	0	0	1	0	0	0	0	...	0	1	0	0	0	0	0	0	0	1
201906	0	0	0	0	0	0	1	0	0	0	...	0	0	0	1	0	0	0	0	1	0
201907	0	0	0	0	0	0	0	0	1	0	...	0	0	1	0	2	0	0	0	0	0

7 rows x 99 columns

1～2行目は日付型の定義と、日付を年月の形に変換を行っています。これは第1章でも実施していますので復習しておきましょう。集計単位に合わせて日付を変換する処理は実際かなり行われます。

3行目で縦軸に購入年月、横軸に商品として件数を集計しています。

4行目は画面上に集計結果を表示しています。

データの揺れを補正せずに集計をしてみました。結果の表を見ると、一部省略されていますが、「商品S」や「商品s」等、本来同じ商品が別の商品として集計されている事が確認できます。

また、表の最後の「7 rows x 99 columns」に注目してください。

今回横軸(columns)は「商品」として集計しました。つまり、データの揺れがあるため、本来26個の商品が99商品に増えてしまっている事が分かります。

同様に横軸に「item\_price」を設定し集計してみましょう。

```
res = uriage_data.pivot_table(index="purchase_month", columns="item_name",
values="item_price", aggfunc="sum", fill_value=0)

res
```