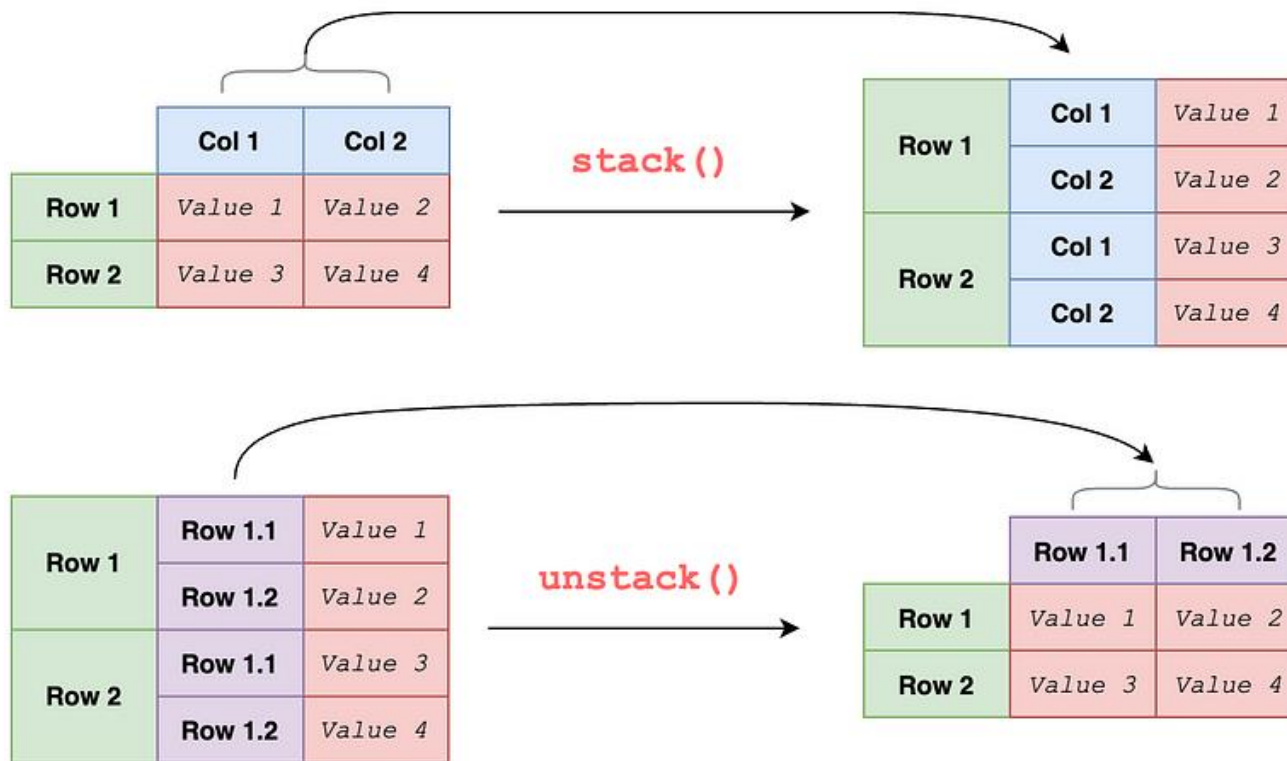# მონაცემთა ანალიტიკა Python

ლექცია 9: მონაცემების ფორმის ცვლილება. მონაცემების გახლეჩვა, დამუშავება და გაერთიანება. ჯვარედინა ტაბულაციის ცხრილები. მონაცემების ტრანსფორმაცია გრძელიდან განიერ და განიერიდან გრძელ ფორმატში.

ლიკა სვანაძე
lika.svanadze@btu.edu.ge

# Reshaping a DataFrame

- Reshaping is often needed when you work with datasets that contain variables with some kinds of sequences, say, time-series data.

- Pandas provides various built-in methods for reshaping DataFrame. Among them, stack() and unstack() are the 2 most popular methods for restructuring columns and rows (also known as index).

- **stack():** stack the prescribed level(s) from column to row.

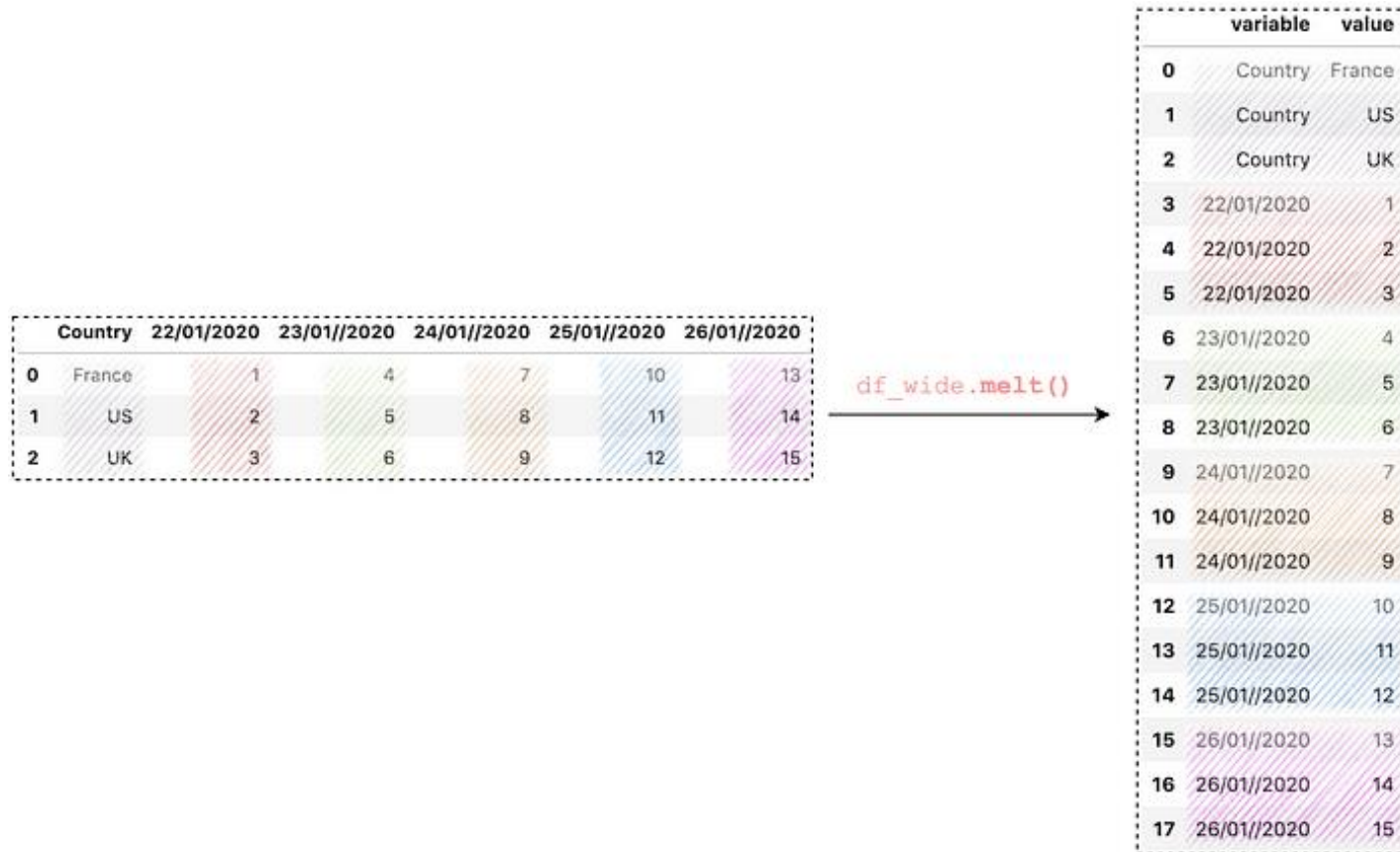- **unstack():** unstack the prescribed level(s) from row to column. The inverse operation from stack.

# Stack() and unstack()

lika.svanadze@btu.edu.ge

# Reshaping a DataFrame using Pandas melt()

- **melt()** function is used to reshape the data from a wide format to a long format to facilitate further data analysis

- The simplest melt() doesn't require any argument and it will turn all columns into rows (shown as a column variable) and list all associated values in a new column value.

- However, this output often doesn't make much sense, so the general use case at least specifies the id_vars argument. For example, id_vars='Country' will tell pandas to keep Country as a column, and turn all the other columns into rows.

lika.svanadze@btu.edu.ge

# Reshaping a DataFrame using Pandas melt()

# Reshaping a DataFrame using Pandas pivot()

- Reshaping a DataFrame from long to wide format using

- In practice, we often need the complete opposite operation as well — to reshape the data from a long to a wide format. That's where the Panda pivot() comes to help. In short, the Pandas pivot() is the complete opposite of melt().

# Reshaping a DataFrame using Pandas explode ()

In the step of data pre-processing, we often need to prepare our data in specific ways before feeding it into a machine learning model. One of the examples is to transform list-like columns into rows. Pandas provides various methods for that, among them apply() and explode() are the two most popular methods.

lika.svanadze@btu.edu.ge

# Frequency, Relative Frequency and CRF

| Data Value | Frequency | Relative Frequency | Cumulative Relative Frequency |
|:---:|:---:|:---:|:---:|
| 2 | 3 | 3/20 = 0.15 | 0.15 |
| 3 | 5 | 5/20 = 0.25 | 0.15 + 0.25 = 0.40 |
| 4 | 3 | 3/20 = 0.15 | 0.40 + 0.15 = 0.55 |
| 5 | 6 | 6/20 = 0.30 | 0.55 + 0.30 = 0.85 |
| 6 | 2 | 2/20 = 0.10 | 0.85 + 0.10 = 0.95 |
| 7 | 1 | 1/20 = 0.05 | 0.95 + 0.05 = 1.00 |

# Contingency Table - crosstab

✳ A contingency table, sometimes called a two-way frequency table, is a tabular mechanism with at least two rows and two columns used in statistics to present categorical data in terms of frequency counts.

✳ To know the relationship between two ordinal or nominal variables then look for contingency table which displays this relationship.

| | | Sport Preference | | |
|---|---|---|---|---|
| | | Archery | Boxing | Cycling |
| **Gender** | Female | 35 | 15 | 50 | 100 |
| | Male | 10 | 30 | 60 | 100 |
| | | **45** | **45** | **110** | **200** |

BTU