

# Python for Data processing

## Lecture 3: **pandas** - Part I

Gleb Ivashkevich

# What we already know

A lot about **numpy** arrays and PyTorch tensors:

- basic operations
- best practices
- optimization
- linear algebra
- very basics of how machine learning works

# Structured and unstructured data

# Unstructured data

- images
- signals (including time series)
- text

Each data element (pixel, datapoint, letter) is usually atomic and **is equal** to any other data element. You need to **perform analysis** to get the structure

# Structured data

- tabular data
- JSON
- XML

Each data element (row, DB record, XML file) has **internal structure** or **schema**

```
[{'name': 'Anny Smith', 'age': 35, 'sex': 'female'},  
 {'name': 'John Black', 'age': 62, 'sex': 'male'}, ...]
```

# Dataframe

**Tabular representation** of structured data

- well known in R world for years
- **indexed** rows and columns
- SQL-like operations<sup>(joins, filtering)</sup>, aggregations, alignment and more

# pandas

One of the most respected Python packages for data science

- started in **2008**
- **very fast** (a lot of Cython inside)
- supports **tons of operations and formats**
- extremely **flexible** and **powerful**
- It's **crazy** sometimes, but you'll **love** it

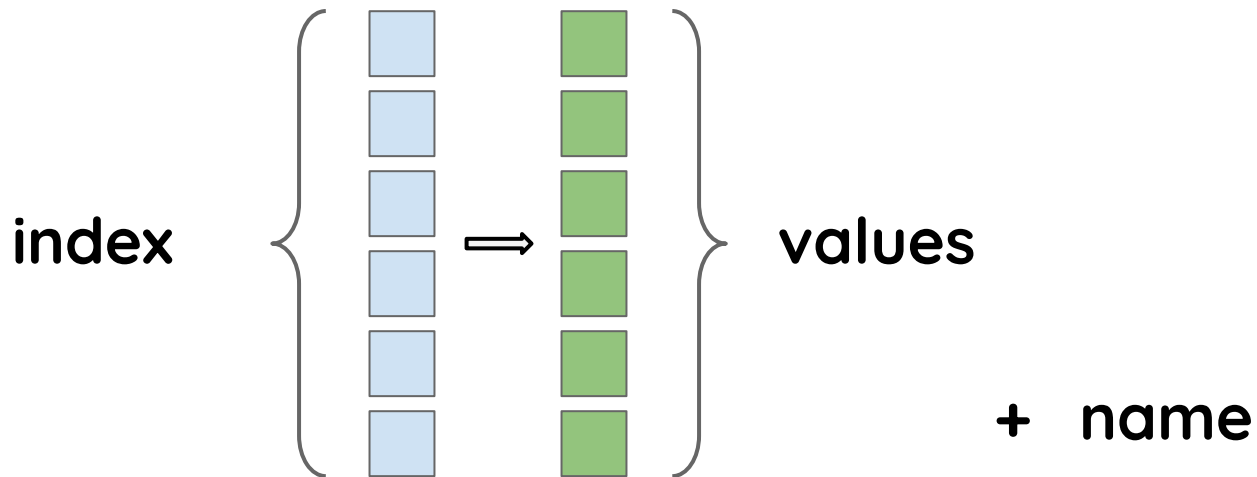
# pandas series and df's

pandas has two main data structures:

- `pd.Series` for indexed 1D data
- `pd.DataFrame` for indexed tabular data

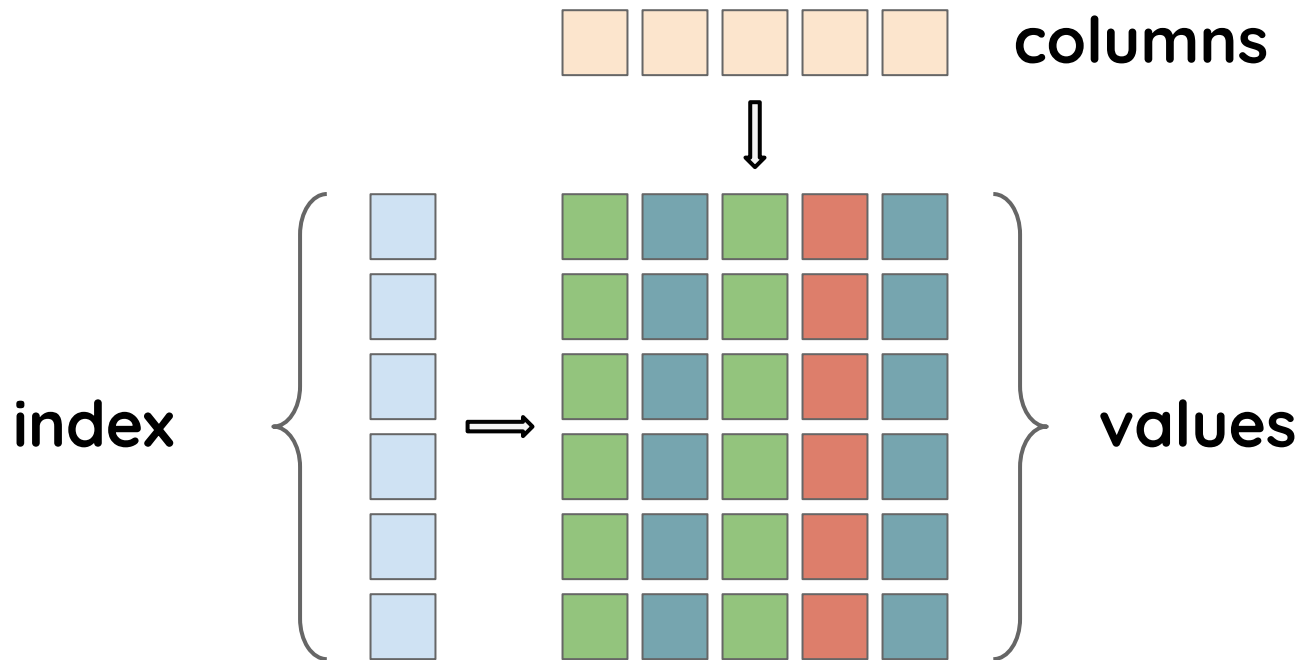


# pandas series

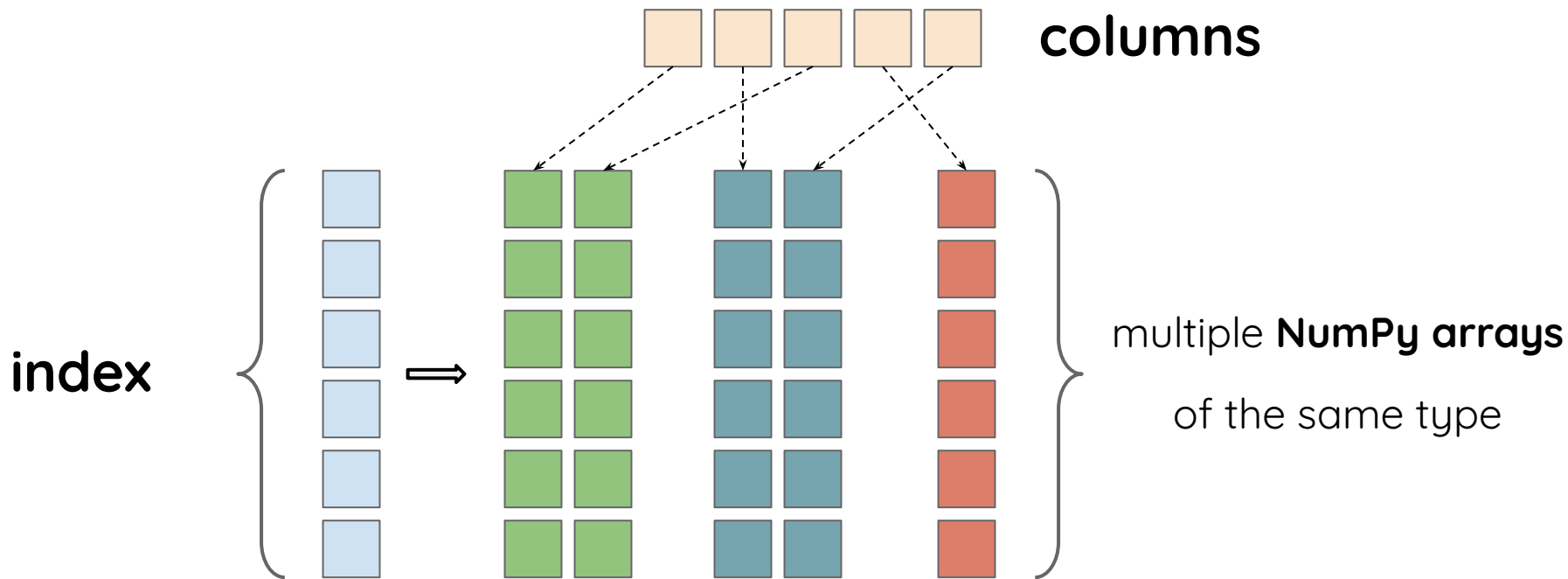


→let's try it out!

# pandas dataframe



# pandas dataframe



# Indexing series and dataframes

# Indexing series and df's

- `[]` indexing
- `.loc` - label based indexing
- `.iloc` - position based indexing

**Boolean indexing** is possible and is heavily used

→let's try it out!

# Indexing df: SettingWithCopyWarning

pandas is not like numpy:

- **It's unknown whether you get view or copy**

Why?

- **It's hard to give guarantees**

(but there are rules inside) (but you should not even try to understand them)

→let's try it out!

# pandas: quiz

For pandas series `s`, `s[1]` will return

- element at integer position 1,
- element at integer index 1,
- it depends on series index `s`.

# Operations on dataframes

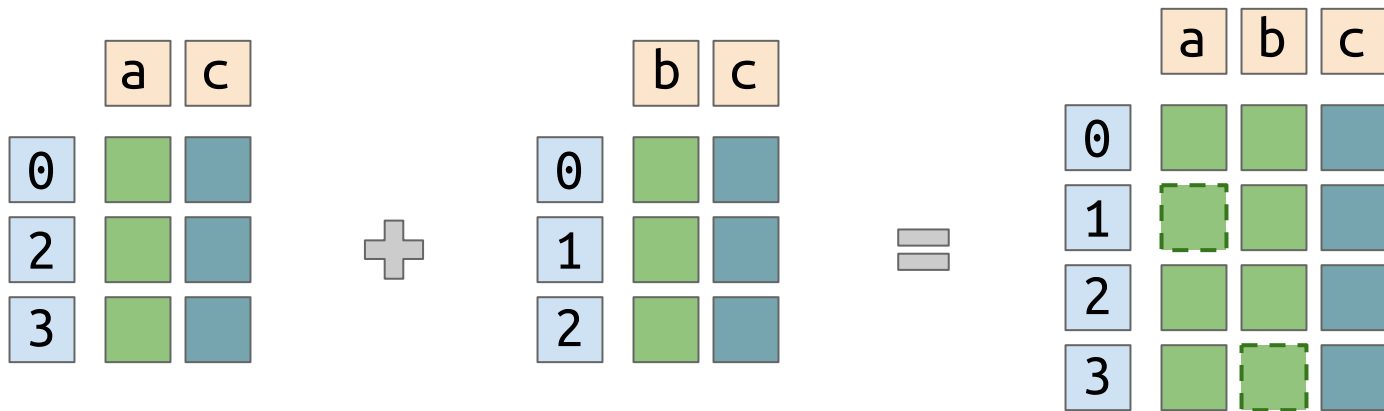


# Arithmetic: not your usual numpy

**pandas** aligns dataframes for you before performing operations by creating a **union** of row and column indexes

→let's try it out!

# Arithmetic



→let's try it out!

# Applying functions to df's

**pandas** allows you to apply custom functions across rows or columns, and elementwise.

And combines results for you appropriately.

It's usually **fast**.

→let's try it out!

# Dataframes: quiz

Given two pandas dataframes, `df1`, with index `[1,2,3]` and columns `[a, b, d]`, and `df2`, with index `[1,3,4]` and columns `[a, b, c]`, result of `df1 + df2` has

- columns `[a,b,c,d]` and index `[1,2,3,4]`,
- columns `[a,b]` and index `[1,3]`,
- neither of two.

# Reading csv files

The best tool to read CSV and other text files in Python:

```
pd.read_csv(...)
```

→let's try it out!

# Dataframe summaries

It's easy to get general information about dataframe:

`df.info()`

`df.describe()`

`df.head(), df.tail()`

→let's try it out!

# Counts and statistics

To get counts or statistics about column or row:

```
df[col].unique()
```

```
df[col].value_counts()
```

```
df.sum(axis=...), df.mean(axis=...), df.std(axis=...)
```

Powerful in combination with smart indexing.

→let's try it out!

# Groupby operations (basic)

Dataframes can be **grouped** by column or columns, or row index level. Great for:

- calculating fine grained statistics
- plotting
- applying operations group-wise

`.groupby` is an **extremely powerful tool**

→let's try it out!



# Replacing and renaming

**df.replace** allows for flexible replacement of values in dataframes:

- by value, per column

**df.rename** allows you to easily rename any label, be it column name or index label

→let's try it out!

# Missing data

**pandas** is great at handling missing data:

- infers it for you
- backward fill, forward fill and more

→let's try it out!

# Categorical data

**pandas** easily calculates one-hot encoded values for any column, adding properly named columns

→let's try it out!

# Special datatypes

pandas has very good support for:

- **strings** - great for text columns  
(split, replace and other usual string operations, vectorized)
- **datetime** - flexible indexing, handling timezones and extravagant parsing  
(great for anything time series related)

→let's try it out!

# What we've learned

Basics of **pandas**:

- creating, indexing
- operations on dataframes
- basics of grouping

# Words we know

- Jupyter
- numpy
- PyTorch
- pandas

# Assignment

- explore **pandas**
- play with Titanic dataset

questions?