# Intro. Comp. for Data Science (FMI08)

Dr. Nono Saha

May 24, 2023

Max Planck Institute for Mathematics in the Sciences
University of Leipzig/ScaDS.AI

Spring 2023

# Course plan

# Introdcution to **pandas**

## pandas?

pandas implements data frames in Python - it takes much of its inspiration from R and NumPy.

pandas aims to be the fundamental high-level building block for practical, real-world data analysis in Python. Additionally, it seeks to become the most powerful and flexible open-source data analysis/manipulation tool available in any language.

Key features:

- DataFrame object class
- Reading and writing tabular data
- Data munging (filtering, grouping, summarizing, joining, etc.)
- Data reshaping

# pandas - Series

# Series - create a `pandas` series from a `list`

The columns of a DataFrame are constructed as `Series` - a 1d array-like object containing values of the same type (similar to an `ndarray`).

```python
import pandas as pd
```

```python
pd.Series([1,2,3,4])
## 0    1
## 1    2
## 2    3
## 3    4
## dtype: int64


pd.Series(["C","B","A"])
## 0    C
## 1    B
## 2    A
## dtype: object
```

```python
pd.Series(range(5))
## 0    0
## 1    1
## 2    2
## 3    3
## 4    4
## dtype: int64

pd.Series([1,"A",True])
## 0       1
## 1       A
## 2    True
## dtype: object
```

Once constructed the components of a series can be accessed via
`array()` and `index()` methods.

```
1 s = pd.Series([4,2,1,3])
2 s.array
3 ## <PandasArray>
4 ## [4, 2, 1, 3]
5 ## Length: 4, dtype: int64
6
7 s.index
8 ## RangeIndex(start=0, stop=4, step=1)
```

An index can also be explicitly provided when constructing a Series,

```
1 t = pd.Series([4,2,1,3], index=[
    "a","b","c","d"])
2 ## a    4
3 ## b    2
4 ## c    1
5 ## d    3
6 ## dtype: int64
```

```
1 t.array
2 ## <PandasArray>
3 ## [4, 2, 1, 3]
4 ## Length: 4, dtype: int64
5
6 t.index
```

Series objects are compatible with **NumPy** like functions (vectorized)

```
1  t = pd.Series([4,2,1,3], index=["a","b","c","d"])
```

```
1  t+1
2  ## a    5
3  ## b    3
4  ## c    2
5  ## d    4
6  ## dtype: int64
7
8  np.log(t)
9  ## a    1.386294
10 ## b    0.693147
11 ## c    0.000000
12 ## d    1.098612
13 ## dtype: float64
```

```
1  np.exp(-t**2/2)
2  ## a    0.000335
3  ## b    0.135335
4  ## c    0.606531
5  ## d    0.011109
6  ## dtype: float64
```

Series can be indexed in the same way as NumPy arrays with the addition of being able to use label(s) when selecting elements.

```python
t = pd.Series([4,2,1,3], index=["a","b","c","d"])
```

```python
t[1]
## 2

t[[1,2]]
## b    2
## c    1
## dtype: int64

t["c"]
## 1

t[["a","d"]]
## a    4
## d    3
## dtype: int64
```

```python
t[t == 3]
## d    3
## dtype: int64

t[t % 2 == 0]
## a    4
## b    2
## dtype: int64

t["d"] = 6
t
## a    4
## b    2
## c    1
## d    6
```

## pandas - Series: index alignment

When performing (arithmetic) operations on Series, they will attempt to align by their index. Let us consider the following `Series`:

```
1 m = pd.Series([1,2,3,4], index = ["a","b","c","d"])
2 n = pd.Series([4,3,2,1], index = ["d","c","b","a"])
3 o = pd.Series([1,1,1,1,1], index = ["b","d","a","c","e"])
```

### Questions

For each of the following operations, what is the output?

- m + n
- n + m
- n + 0
- m + 0

We can construct Series from `dicts`, in which case the keys are used to form the index,

```
1 d = {"anna": "A+", "bob": "B-", "carol": "C", "dave": "D+"}
2 pd.Series(d)
3 ## anna      A+
4 ## bob       B-
5 ## carol      C
6 ## dave      D+
7 ## dtype: object
```

Index order will follow key order unless overridden by index,

```
1 pd.Series(d, index = ["dave","carol","bob","anna"])
2
3 ## dave      D+
4 ## carol      C
5 ## bob       B-
6 ## anna      A+
7 ## dtype: object
```

# pandas - Series: string series

Series containing strings can be accessed via the `str` attribute,

```
s = pd.Series(["the quick", "
    brown fox", "jumps over", "a
     lazy dog"])
s
## 0      the quick
## 1      brown fox
## 2     jumps over
## 3      a lazy dog
## dtype: object

s.str.upper()
## 0      THE QUICK
## 1      BROWN FOX
## 2     JUMPS OVER
## 3      A LAZY DOG
## dtype: object
```

```
s.str.split(" ").str[1]
## 0     quick
## 1       fox
## 2      over
## 3      lazy
## dtype: object

pd.Series([1,2,3]).str
## AttributeError: Can only use
    .str accessor with string
```

```python
pd.Series(["Mon",
    "Tue", "Wed"
    , "Thur", "
    Fri"])
## 0     Mon
## 1     Tue
## 2     Wed
## 3    Thur
## 4     Fri
## dtype: object
```

```python
pd.Series(["Mon", "Tue", "Wed", "Thur", "Fri"],
    dtype="category")
## 0     Mon
## 1     Tue
## 2     Wed
## 3    Thur
## 4     Fri
## dtype: category
## Categories (5, object): ['Fri', 'Mon', 'Thur
    ', 'Tue', 'Wed']
```

```python
pd.Series(["Mon", "Tue", "Wed", "Thur", "Fri"], dtype=pd.
    CategoricalDtype(ordered=True))
## 0     Mon
## 1     Tue
## 2     Wed
## 3    Thur
## 4     Fri
## dtype: category
## Categories (5, object): ['Fri' < 'Mon' < 'Thur' < 'Tue' < '
    Wed']
```

**pandas** - DataFrame

# pandas - DataFrame: constructing DataFrames

panda data frames can also be constructed via `DataFrame()`,
general this is done via dict of columns:

```
n = 5
d = {
"id":     np.random.randint(100, 999, n),
"weight": np.random.normal(70, 20, n),
"height": np.random.normal(170, 15, n),
"date":   pd.date_range(start='2/1/2022', periods=n, freq='D')
}
df = pd.DataFrame(d)
df

##      id      weight      height        date
## 0  168   102.915535  188.677769  2022-02-01
## 1  615    71.767364  155.907801  2022-02-02
## 2  346    76.666059  171.386839  2022-02-03
## 3  390    74.735465  173.151008  2022-02-04
## 4  556    50.538488  183.083407  2022-02-05
```

# `pandas` - DataFrame: from nparray

For 2d ndarrays, it is also possible to construct a DataFrame - generally, providing column names and row names (indexes) is a good idea.

```
pd.DataFrame(
np.diag([1,2,3]),
columns = ["x","y","z"]
)
##    x  y  z
## 0  1  0  0
## 1  0  2  0
## 2  0  0  3

pd.DataFrame(
np.diag([1,2,3]),
columns = ["x","y","z"]
)
##    x  y  z
## 0  1  0  0
## 1  0  2  0
## 2  0  0  3
```

```
pd.DataFrame(
np.tri(5,3,-1),
columns = ["x","y","z"],
index = ["a","b","c","d","e"]
)
##      x    y    z
## a  0.0  0.0  0.0
## b  1.0  0.0  0.0
## c  1.0  1.0  0.0
## d  1.0  1.0  1.0
## e  1.0  1.0  1.0
```

```
1  df[0]
2  ## KeyError: 0
3
4  df["id"]
5  ## 0    168
6  ## 1    615
7  ## 2    346
8  ## 3    390
9  ## 4    556
10 ## Name: id, dtype: int64
11
12 df.id
13 ## 0    168
14 ## 1    615
15 ## 2    346
16 ## 3    390
17 ## 4    556
18 ## Name: id, dtype: int64
```

```
df[1:3]
##     id     weight      height
    date
## 1  615  71.767364  155.907801
    2022-02-02
## 2  346  76.666059  171.386839
    2022-02-03

df[0::2]
##     id     weight       height
    date
## 0  168  102.915535  188.677769
    2022-02-01
## 2  346   76.666059  171.386839
    2022-02-03
## 4  556   50.538488  183.083407
    2022-02-05
```

```
1  df.iloc[1] && df.iloc[[1]]
2  ## What is the difference between
       the two instructions?
3
4  df.iloc[0:2]
5  ## id  weight   height    date
6  ## 0  168  102.915535  188.677769
       2022-02-01
7  ## 1  615   71.767364  155.907801
       2022-02-02
8
9  df.iloc[lambda x: x.index % 2 != 0]
10 ##     id    weight      height
          date
11 ## 1  615  71.767364  155.907801
       2022-02-02
12 ## 3  390  74.735465  173.151008
       2022-02-04
```

```
df.iloc[1:3,1:3]
##      weight      height
## 1  71.767364  155.907801
## 2  76.666059  171.386839

df.iloc[0:3, [0,3]]
##     id       date
## 0  168 2022-02-01
## 1  615 2022-02-02
## 2  346 2022-02-03

df.iloc[0:3, [True, True,
   False, False]]
##     id     weight
## 0  168  102.915535
## 1  615   71.767364
## 2  346   76.666059
```

```
df.loc[["anna"]]
## id weight height date
## anna  168   102.915535
     188.677769 2022-02-01

df.loc["bob":"dave"]
## id weight height date
## bob 615 71.76 155.9 2022-02-02
## carol 346 76.6 171.3 2022-02-03
## dave 390 74.73 173.15 2022-02-04

df.loc[df.id < 300]
## id weight height date
## anna 168 102.91 188.67
    2022-02-01
```

```
df.loc[:, "date"]
## anna     2022-02-01
## bob      2022-02-02
## carol    2022-02-03
## dave     2022-02-04
## erin     2022-02-05
## Name: date, dtype:
    datetime64[ns]

df.loc[["bob","erin"], "
    weight":"height"]
## weight height
## bob   71.76  155.9
## erin  50.53  183.08

df.loc[0:2, "weight":"height"
    ]
## ???
```

## pandas - DataFrame: views and copies

In general, most pandas operations will generate a new object, but
some will return views, mostly the later occurring with subsetting.

```
d = pd.DataFrame(np.arange(6).reshape(3,2), columns = ["x","y"])
```

```
v = d.iloc[0:2,0:2]
v
##    x  y
## 0  0  1
## 1  2  3

d.iloc[0,1] = -1
v
##    x  y
## 0  0 -1
## 1  2  3
```

```
v.iloc[0,0] = np.pi
v
##            x  y
## 0  3.141593 -1
## 1  2.000000  3

d
##    x  y
## 0  0 -1
## 1  2  3
## 2  4  5
```

The `query()` method can be used for filtering rows. It evaluates a string expression in the context of the data frame.

```
df = pd.DataFrame(d)
df.query('date == "2022-02-01"')
##          id      weight       height        date
## anna   168   102.915535   188.677769 2022-02-01

df.query('weight > 50')
##          id       weight       height        date
## anna   168   102.915535   188.677769 2022-02-01
## bob    615    71.767364   155.907801 2022-02-02
## carol  346    76.666059   171.386839 2022-02-03
## dave   390    74.735465   173.151008 2022-02-04

df.query('weight > 50 & height < 165')
##        id      weight       height        date
## bob   615   71.767364   155.907801 2022-02-02

qid = 414
df.query('id == @qid')
```

17

```
1 df
2 ##          id       weight      height       date
3 ## anna    168   102.915535  188.677769  2022-02-01
4 ## bob     615    71.767364  155.907801  2022-02-02
5 ## carol   346    76.666059  171.386839  2022-02-03
6 ## dave    390    74.735465  173.151008  2022-02-04
7 ## erin    556    50.538488  183.083407  2022-02-05
```

```
1 df[0,0]
2 ## KeyError: (0, 0)
3
4 df.iat[0,0]
5 ## 168
6
7 df.id[0]
8 ## 168
9
10 df[0:1].id[0]
11 ## 168
```

```
1 df["anna", "id"]
2 ## KeyError: ('anna', 'id')
3
4 df.at["anna", "id"]
5 ## 168
6
7 df["id"]["anna"]
8 ## 168
9
10 df["id"][0]
11 ## 168
```

```
1 df
2 ##          id       weight      height        date
3 ## anna    168  102.915535  188.677769  2022-02-01
4 ## bob     615   71.767364  155.907801  2022-02-02
5 ## carol   346   76.666059  171.386839  2022-02-03
6 ## dave    390   74.735465  173.151008  2022-02-04
7 ## erin    556   50.538488  183.083407  2022-02-05
```

```
1 df.size
2 ## 20
3
4 df.shape
5 ## (5, 4)
6
7 df.info()
8 ## Try it on your computer and
      analyze the output.
```

```
1 df.dtypes
2 ## id                   int64
3 ## weight             float64
4 ## height             float64
5 ## date        datetime64[ns]
6 ## dtype: object
7
8 df.describe()
```

Beyond the use of `loc()` and `iloc()`, there is also the `filter()` method which can
be used to select columns (or indices) by name with pattern-matching.

```
1 df.filter(items=["id","weight"])
2 ##            id      weight
3 ## anna     168   102.915535
4 ## bob      615    71.767364
5 ## carol    346    76.666059
6 ## dave     390    74.735465
7 ## erin     556    50.538488
8
9 df.filter(like = "i")
10 ## id weight height
11 ## anna     168   102.91   188.67
12 ## bob      615    71.76   155.90
13 ## carol    346    76.66   171.38
14 ## dave     390    74.73   173.15
15 ## erin     556    50.53   183.08
```

```
1 df.filter(regex="ght$")
2 ## weight  height
3 ## anna   102.91  188.67
4 ## bob     71.76  155.90
5 ## carol   76.66  171.38
6 ## dave    74.73  173.15
7 ## erin    50.53  183.08
8
9 df.filter(like="o", axis=0)
10 ##  id  weight  height  date
11 ## bob  615   71.76   155.90
       2022-02-02
12 ## carol  346   76.66   171.38
       2022-02-03
```

# pandas - DataFrame: adding columns

Indexing with assignment allows for in-place modification of a
DataFrame, while `assign()` creates a new object (but is chainable).

```
df['student'] = [True, True, True,
    False, None]
df['age'] = [19, 22, 25, None, None
    ]
df

## id weight height date student
    age
## anna   168  102.91  188.67
    2022-02-01      True   19.0
## dave   390   74.7  173.15
    2022-02-04     False    NaN
## erin   556   50.5  183.08
    2022-02-05      None    NaN
```

```
df.assign(student = lambda x:
    np.where(x.student, "yes
    ", "no"),
rand = np.random.rand(5)
)

## id weight   height   date
    student  age   and
## anna   168  102.91  188.67
    2022-02-01     yes  19.0
    0.60
## bob    615   71.76  155.90
    2022-02-02     yes  22.0
    0.54
```

Columns can be dropped via the `drop()` method,

```
1 df.drop(['student'])
2 ## KeyError: "['student'] not
    found in axis."
3
4 df.drop(['student'], axis=1)
5 ## id weight height  date   age
6 ## anna   168  102.91  188.67
    2022-02-01  19.0
7 ## bob    615   71.76  155.90
    2022-02-02  22.0
8 ## carol  346   76.66  171.38
    2022-02-03  25.0
9 ## dave   390   74.73  173.15
    2022-02-04   NaN
10 ## erin   556   50.53  183.08
    2022-02-05   NaN
```

```
df.drop(columns = df.columns ==
    "age")
## KeyError: '[False, False,
    False, False, False, True]
    not found in axis'

df.drop(columns = df.columns[df.
    columns == "age"])
## id  weight height  date
    student
## anna   168  102.9  188.67
    2022-02-01    True
## bob    615   71.7  155.9
    2022-02-02    True
## carol  346   76.6  171.3
    2022-02-03    True
```

## `pandas` - DataFrame: row binds

DataFrames can have their rows joined via the the `concat()` function (`append()` is also available but deprecated),

```
1  df1 = pd.DataFrame(np.arange(6).
      reshape(3,2), columns=list("
      xy"))
2  df1
3  ##    x  y
4  ## 0  0  1
5  ## 1  2  3
6  ## 2  4  5
7
8  pd.concat([df1,df2])
9  ##    x   y
10 ## 0  0   1
11 ## 1  2   3
12 ## 2  4   5
13 ## 0  12  11
14 ## 1  10  9
15 ## 2  8   7
```

```
df2 = pd.DataFrame(np.arange
    (12,6,-1).reshape(3,2),
    columns=list("xy"))
df2
##    x   y
## 0  12  11
## 1  10  9
## 2  8   7

pd.concat([df1.loc[:,["y","x"]],
    df2])
##    y   x
## 0  1   0
## 1  3   2
## 2  5   4
## 0  11  12
## 1  9   10
```

23

# Missing value problems and native Nas

pandas encodes missing values using NaN (mostly),

```
1 s = pd.Series({"anna": "A+", "
    bob": "B-", "carol": "C", "
    dave": "D+"},
2 index = ["erin","dave","carol","
    bob","anna"])
3 s
4 ## erin      NaN
5 ## dave       D+
6 ## carol       C
7 ## bob        B-
8 ## anna       A+
9 ## dtype: object
10
11 pd.isna(s) = ?
12
```

```
1 s = pd.Series({"anna": 97, "bob"
    : 82, "carol": 75, "dave":
    68},
2 index = ["erin","dave","carol","
    bob","anna"],dtype='int64')
3 s
4 ## erin       NaN
5 ## dave      68.0
6 ## carol     75.0
7 ## bob       82.0
8 ## anna      97.0
9 ## dtype: float64
10
11 pd.isna(s) = ?
12
```

```
1 s = pd.Series([1,2,3,None])
2 s
3 ## 0    1.0
4 ## 1    2.0
5 ## 2    3.0
6 ## 3    NaN
7 ## dtype: float64
```

```
1 pd.isna(s)
2 ## 0    False
3 ## 1    False
4 ## 2    False
5 ## 3     True
6 ## dtype: bool
7
8 s == np.nan
9 ## 0    False
10 ## 1    False
11 ## 2    False
12 ## 3    False
13 ## dtype: bool
14
```

```
1 np.nan == np.nan
2 ## What about this?
3
4 np.nan != np.nan
5 ## And this?
6
```

25

## pandas - Series: native NAs

Recent versions of pandas have attempted to adopt a more native missing value, particularly for integer and boolean types,

```
1 pd.Series([1,2,3,None])
2 ## 0    1.0
3 ## 1    2.0
4 ## 2    3.0
5 ## 3    NaN
6 ## dtype: float64
7
8 pd.Series([True,False,None])
9 ## 0     True
10 ## 1    False
11 ## 2     None
12 ## dtype: object
13
```

```
1 pd.isna( pd.Series([1,2,3,None])
     )
2 ## 0    False
3 ## 1    False
4 ## 2    False
5 ## 3     True
6 ## dtype: bool
7
8 pd.isna( pd.Series([True,False,
     None]) )
9 ## 0    False
10 ## 1    False
11 ## 2     True
12 ## dtype: bool
13
```