# Intro. Comp. for Data Science (FMI08)

Dr. Nono Saha

May 26, 2023

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

Spring 2023

## Course plan

1. Index objects
2. MultiIndexes
3. Reshaping data
4. Split-Apply-Combine

More `pandas` - Index objects

# Index objects: columns and indexes

When constructing a DataFrame we can specify the indexes for both the rows (`index`) and columns (`index`),

```
df = pd.DataFrame(np.random.
    randn(5, 3),
columns=['A', 'B', 'C'])
df
## A B C
## 0 -0.091 -0.37 -2.39
## 1 0.54 1.17 1.23
## 2 -0.6 1.8 0.67

df.columns
## Index(['A', 'B', 'C'], dtype
    ='object')

df.index
## RangeIndex(start=0, stop=5,
    step=1)
```

```
df = pd.DataFrame(np.random.
    randn(3, 3),
index=['x','y','z'], columns=['A
    ', 'B', 'C'])
## A B C
## x 0.61 1.12 -0.8
## y 0.90 -1.46 0.54
## z -1.95 0.75 0.65

df.columns
## Index(['A', 'B', 'C'], dtype
    ='object')

df.index
## Index(['x', 'y', 'z'], dtype
    ='object')
```

# Index objects: creating an index object

pandas' `Index` class and its subclasses provide the infrastructure necessary for lookups, data alignment, and other related tasks. You can think of them as being an immutable multiset (duplicate values are allowed).

```python
pd.Index(['A','B','C'])

## Index(['A', 'B', 'C'], dtype='object')

pd.Index(['A','B','C','A'])

## Index(['A', 'B', 'C', 'A'], dtype='object')

pd.Index(range(5))

## RangeIndex(start=0, stop=5, step=1)

pd.Index(list(range(5)))

## Int64Index([0, 1, 2, 3, 4], dtype='int64')
```

While it is not something you will need to do very often, since Indexes are "sets" the various set operations and methods are available.

```
1 a = pd.Index(['c', 'b', 'a'])
2 b = pd.Index(['c', 'e', 'd'])
3
```

```
1 a.union(b)
2 ## Index(['a', 'b', 'c', 'd', 'e
      '], dtype='object')
3
4 a.intersection(b)
5 ## Index(['c'], dtype='object')
6
7 c = pd.Index([1.0, 1.5, 2.0])
8 d = pd.Index(range(5))
9 c.union(d)
10 ## Float64Index([0.0, 1.0, 1.5,
      2.0, 3.0, 4.0], dtype='
      float64')
11
```

```
a.difference(b)
## Index(['a', 'b'], dtype='
    object')

a.symmetric_difference(b)
## Index(['a', 'b', 'd', 'e'],
    dtype='object')

e = pd.Index(["A","B","C"])
f = pd.Index(range(5))
e.union(f)
## Index(['A', 'B', 'C', 0, 1,
    2, 3, 4], dtype='object')
```

You can attach names to an index, which will then show when displaying the DataFrame or Index,

```
1 df = pd.DataFrame( np.random.
      randn(3, 3),
2 index=pd.Index(['x','y','z'],
      name="rows"),
3 columns=pd.Index(['A', 'B', 'C'
      ], name="cols")
4 )
5 df.columns
6 ## Index(['A', 'B', 'C'], dtype
      ='object', name='cols')
7
8 df.index
9 ## Index(['x', 'y', 'z'], dtype
      ='object', name='rows')
10
```

```
df.columns.rename("m")
## Index(['A', 'B', 'C'], dtype
    ='object', name='m')

df.index.set_names("n")
## Index(['x', 'y', 'z'], dtype
    ='object', name='n')

df.columns.name = "o"
df.index.rename("p", inplace=
    True)
df
```

It is possible for an index to contain missing values (e.g. `np.nan`) but this is generally a bad idea and should be avoided.

```
1 pd.Index([1,2,3,np.nan,5])
2 ## Float64Index([1.0, 2.0, 3.0, nan, 5.0], dtype='float64')
3
4 pd.Index(["A","B",np.nan,"D"])
5 ## Index(['A', 'B', nan, 'D'], dtype='object')
6
```

Missing values can be replaced via the `fillna()` method,

```
1 pd.Index([1,2,3,np.nan,5]).fillna(0)
2 ## Float64Index([1.0, 2.0, 3.0, 0.0, 5.0], dtype='float64')
3
4 pd.Index(["A","B",np.nan,"D"]).fillna("Z")
5 ## Index(['A', 'B', 'Z', 'D'], dtype='object')
6
```

# pandas - changing a DataFrame's index

Existing columns can used as an index via `set_index()` and removed via `reset_index()`,

```
1 data
2 ##       a    b  c  d
3 ## 0  bar  one  z  1
4 ## 1  bar  two  y  2
5 ## 2  foo  one  x  3
6 ## 3  foo  two  w  4
7
```

```
1 data.set_index('a')
2
3 data.set_index('c', drop=
      False)
4
5 data.reindex(["w","x","y","z"
      ])
6
7 data.reindex(range(5,-1,-1))
8
```

```
data.set_index('a').reset_index()

data.set_index('c').reset_index(
    drop=True)

data.reindex(columns = ["a","b","c"
    ,"d","e"])

data.index = ["w","x","y","z"]
```

More **pandas** - multiIndexes

These are a hierarchical analog of standard Index objects, there are a number of methods for constructing them based on the initial object.

```
1  tuples = [('A','x'), ('A','y'),('B','x'), ('B','y'),('C','x'), (
      'C','y')]
2  pd.MultiIndex.from_tuples(tuples, names=["1st","2nd"])
3
4  pd.MultiIndex.from_product([["A","B","C"],["x","y"]], names=["1
      st","2nd"])
5
6  idx = pd.MultiIndex.from_tuples(tuples, names=["1st","2nd"])
7  pd.DataFrame(np.random.rand(6,2), index = idx, columns=["m","n"
      ])
8
9  #Column MultiIndex
10 cidx = pd.MultiIndex.from_product([["A","B"],["x","y"]], names=[
      "c1","c2"])
11 pd.DataFrame(np.random.rand(4,4), columns = cidx)
12
```

```
1 data = pd.DataFrame(np.random.rand(4,4), index= ridx, columns =
    cidx)
2 data
3 ## c1              A                       B
4 ## c2              x           y           x           y
5 ## r1 r2
6 ## m  l   0.019149  0.519056  0.924092  0.996320
7 ##    p   0.219535  0.537471  0.962619  0.968074
8 ## n  l   0.020447  0.817611  0.493241  0.632190
9 ##    p   0.432398  0.854118  0.774252  0.838321
10
```

### Examples

```
1 data["A"]
2
3 data["x"]
4
5 data["m"]
6
```

```
data["m","A"]
## KeyError: ('m', 'A')

data["A","x"]

data["A"]["x"]
```

# MultiIndex: indexing via `iloc`

```
data
## c1               A                       B
## c2               x          y            x          y
## r1 r2
## m  l    0.019149   0.519056   0.924092   0.996320
##    p    0.219535   0.537471   0.962619   0.968074
## n  l    0.020447   0.817611   0.493241   0.632190
##    p    0.432398   0.854118   0.774252   0.838321
```

```
data.iloc[0]
#Try and see the output

data.iloc[(0,1)]
## 0.519055710819791

data.iloc[[0,1]]
```

```
data.loc[("m","l")]
###

data.loc[:,("A","y")]
###
```

# MultiIndex: fancier indexing with `loc`

Index slices can also be used with combinations of indexes and index tuples,

```
1 data
2 ## c1                    A                       B
3 ## c2              x           y           x           y
4 ## r1 r2
5 ## m   l    0.019149    0.519056    0.924092    0.996320
6 ##     p    0.219535    0.537471    0.962619    0.968074
7 ## n   l    0.020447    0.817611    0.493241    0.632190
8 ##     p    0.432398    0.854118    0.774252    0.838321
9
```

```
1 data.loc["m":"n"]
2 #Try and see the output
3
4 data.loc[("m","l"):("n","l")]
5
```

```
1 data.loc[("m","p"):"n"]
2 ###
3
4 data.loc[[("m","p"),("n","l")]]
5 ###
6
```

## MultiIndex: selecting nested levels

The previous methods don't give easy access to indexing on nested index levels. This is possible via the cross-section method `xs()`,

```
data
## c1                 A                       B
## c2          x           y           x           y
## r1 r2
## m   l   0.019149   0.519056   0.924092   0.996320
##      p   0.219535   0.537471   0.962619   0.968074
## n   l   0.020447   0.817611   0.493241   0.632190
##      p   0.432398   0.854118   0.774252   0.838321

```

```
data.xs("p", level="r2")
#Try and see the output

data.xs("m", level="r1")

```

```
data.xs("y", level="c2", axis=1)
###

data.xs("B", level="c1", axis=1)
###

```

## MultiIndex: setting multiIndexes

It is also possible to construct a MultiIndex or modify an existing one
using `set_index()` and `reset_index()`,

```
1  data
2
3  ##       a    b   c   d
4  ## 0   bar  one   z   1
5  ## 1   bar  two   y   2
6  ## 2   foo  one   x   3
7  ## 3   foo  two   w   4
8
```

```
1  data.set_index(['a','b'])
2  #Try and see the output
3
4  data.set_index('c', append=True)
5
```

```
1  data.set_index(['a','b']).
       reset_index()
2  ###
3
4  data.set_index(['a','b']).
       reset_index(level=1)
5  ###
6
```

# pandas - Reshaping data

# Reshaping data: long to wide (pivot) and wide to long (melt)

```python
df = pd.read_csv("../data/reshaping.csv", index_col=0)

df_wide = df.pivot(index=["country","year"],
columns="type",  values="count")

```

```python
df_wide.index
## MultiIndex([('A', 1999)
        ,
## ('A', 2000),
## ('B', 1999),
## ('B', 2000),
## ('C', 1999),
## ('C', 2000)],
##   names=['country', '
    year'])

df_wide.columns
## Index(['cases', 'pop'],
       dtype='object', name
    ='type')

```

```python
df_wide.reset_index().rename_axis(
    columns=None)
##    country  year cases    pop
## 0        A  1999  0.7K    19M
## 1        A  2000    2K    20M
## 2        B  1999   37K   172M
## 3        B  2000   80K   174M
## 4        C  1999  212K     1T
## 5        C  2000  213K     1T

df_long = df.melt(
id_vars="country",
var_name="year"
)
df_long
```

14

```
1 df = pd.read_csv("../data/rate.csv", index_col=0)
2 df
3 ##    country  year       rate
4 ## 0        A  1999   0.7K/19M
5 ## 1        A  2000     2K/20M
6 ## 2        B  1999   37K/172M
7 ## 3        B  2000   80K/174M
8 ## 4        C  1999    212K/1T
9
```

```
1 df.assign(rate = lambda d:
        d.rate.str.split("/")
        ).explode("rate")
2 .assign(type = lambda d: [
        "cases", "pop"] * int(
        d.shape[0]/2))
3
```

```
df.assign(rate = lambda d: d.rate.str.
    split("/"))

df.assign(rate = lambda d: d.rate.str.
    split("/"))
.explode("rate").assign(type = lambda
    d: ["cases", "pop"] * int(d.shape
    [0]/2))
.pivot(index=["country","year"],
    columns="type", values="rate")
.reset_index()
```

15

# Reshaping data: separate example - a better way

```
1 df.assign(
2 counts = lambda d: d.rate.str.split("/").str[0],
3 pop    = lambda d: d.rate.str.split("/").str[1])
4
```

If you dont want to repeat the split,

```
1 df.assign(
2 rate = lambda d: d.rate.str.split("/"),
3 counts = lambda d: d.rate.str[0],
4 pop    = lambda d: d.rate.str[1]
5 ).drop("rate", axis=1)
6
7 df.assign(
8 counts = lambda d: d.rate.str.split("/").str[0],
9 pop    = lambda d: d.rate.str.split("/").str[1]
10 )
11
```

Create a DataFrame from the data available at `../data/rent.csv` using `pd.read_csv()`.

These data come from the 2017 American Community Survey and reflect the following values:

- `name` - name of state
- `variable` - Variable name: income = median yearly income, rent = median monthly rent
- `estimate` - Estimated value
- `moe` - 90% margin of error

Using these data, find the state(s) with the lowest income-to-rent ratio.