# Pandas
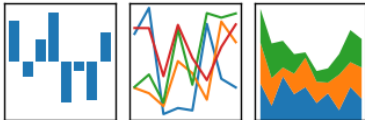
pandas

$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$

# pandas: a library for data manipulation

▶ Documentation:
   http://pandas.pydata.org/pandas-docs/stable/

# Reading in any textfile

```
1  import pandas as pd
2  FNAME = "http://www.stat.ucla.edu/projects/datasets/twins.dat"
3
4  df = pd.read_csv(FNAME, sep=',')
```

# Inspecting the DataFrame

```python
1  df.shape   # Dimensions
2  df.head()  # First 5 lines
3  df.tail()  # Last 5 lines
4  df.columns  # List of variables as list
5  df.describe()  # Summary statistics
```

# Inspecting the DataFrame

```
1  df.shape   # Dimensions
2  df.head()  # First 5 lines
3  df.tail()  # Last 5 lines
4  df.columns  # List of variables as list
5  df.describe()  # Summary statistics
```

1. How many observations do you have?
2. How many variables do you have?
3. Which variables are numeric?
4. What is the mean of variable `"DEDUC1"`?

# Understanding dtypes

```
1  df.info()
```

# Understanding dtypes

```
1  df.info()
```

| Pandas | Python | Purpose |
|--------|--------|---------|
| object | unicode | Text |
| int64 | int | Integers |
| float64 | float | Floating numbers |
| bool | bool | True and False values |
| datetime64 | | Date and time values |
| timedelta[ns] | | Differences between two datetimes |
| category | | Finite list of text values |

# Slicing the DataFrame

```python
1  df["DEDUC"]   # Column by column name
2  df.loc[0]   # Row by index name
3  df.iloc[0]   # Row by row number
4  df[["AGE", "HRWAGEH"]]   # Columns by list of names
5  df.loc["AGE":"HRWAGEH"]   # Column range by column names (will be empty)
6  df.iloc[:,5:7]   # Column range by column indices
7  df.iloc[1, 4]   # Value by row and column index
8  df.loc[18, "WHITEH"]   # Value by row and column name
9  df["WHITEH"].iloc[18]   # Same as above
```

# Slicing the DataFrame

```
1  df["DEDUC"]   # Column by column name
2  df.loc[0]   # Row by index name
3  df.iloc[0]   # Row by row number
4  df[["AGE", "HRWAGEH"]]   # Columns by list of names
5  df.loc["AGE":"HRWAGEH"]   # Column range by column names (will be empty)
6  df.iloc[:,5:7]   # Column range by column indices
7  df.iloc[1, 4]   # Value by row and column index
8  df.loc[18, "WHITEH"]   # Value by row and column name
9  df["WHITEH"].iloc[18]   # Same as above
```

1. What is the 6th entry of the 5th column?
2. What is the 5th entry of column "DTEN"?
3. What is the last entry of column "HRWAGEL"?

# Boolean indexing

```
1  df[df["AGE"] > 20]
2  df[(df["AGE"] > 20) & (df["WHITEL"] == 1)]
3  df[~df["AGE"] > 20]
4  df[df["AGE"].isin((20, 21, 22, 23))]
```

# Boolean indexing

```
1  df[df["AGE"] > 20]
2  df[(df["AGE"] > 20) & (df["WHITEL"] == 1)]
3  df[~df["AGE"] > 20]
4  df[df["AGE"].isin((20, 21, 22, 23))]
```

1. How many observations have `"WHITEL"` equal to 0?
2. How many observations have `"WHITEH"` equal to 1 and `"DEDUC1` unequal to 0?
3. What is the mean age of married twins whose L-parent is a non-white male? (Use `"DMARRIED"`, `"WHITEL"` and `"MALEL"`)
4. In how many rows do the values for `"WHITEH"` and `"WHITEL"` differ?

# Aggregate data

```
1  df["WHITEL"].value_counts()
2  pd.crosstab(df["WHITEH"], df["WHITEL"])
3  df[["DEDUC2", "EDUCL"]].corr()
```

# Aggregate data

```
1  df["WHITEL"].value_counts()
2  pd.crosstab(df["WHITEH"], df["WHITEL"])
3  df[["DEDUC2", "EDUCL"]].corr()
```

1. What is the most common value in "EDUCL"?
2. What is the most common combination of "MALEH" and "MALEL"?
3. What is the Spearman correlation between "EDUCH" and "EDUCL"? What is the Pearson correlation? (Check documentation!)

# Changing dtypes

```
1  df["DMARRIED"] = df["DMARRIED"].astype(bool)
2  df["WHITEH"] = df["WHITEH"].astype("category")
3  df["HRWAGEH"] = pd.to_numeric(df["HRWAGEH"], errors="coerce")
```

# Manipulation

```
1  df = df.sort_values(by='HRWAGEH')   # Sort by column
2  df = df[sorted(df.columns)]   # Re-order columns alphabetically
3  df['new'] = 9   # Add new column
4  df['AGETR'] = df['AGE']**3
5  df['combined'] = df['MALEH'] + df['EDUCH']
6  df["HRWAGEH_new"] = df["HRWAGEH"].fillna(0)   # Fill missings with 0
7  df = df.dropna(subset=["HRWAGEH"])   # Drop rows missing in "HRWAGEH"
```

# Grouping

```
1  grouped = df.groupby(['MALEH'])
2  print(grouped['AGE'].mean())
3  print(grouped['EDUCH'].agg(['mean', 'sum']))
4  print(grouped[['EDUCH', 'AGE']].agg(['mean', 'std']))
```

# Grouping

```
1  grouped = df.groupby(['MALEH'])
2  print(grouped['AGE'].mean())
3  print(grouped['EDUCH'].agg(['mean', 'sum']))
4  print(grouped[['EDUCH', 'AGE']].agg(['mean', 'std']))
```

$\rightarrow$ Full list at http://pandas.pydata.org/pandas-docs/stable/getting_started/basics.html#descriptive-statistics

► What is the "AGE" variance for "MALEL" == 0 individuals?

► What are the second and the third quartile of "EUDCL" for "MALEL" == 0 individuals?

► What is the average age for individuals with "MALEL" == 0 and "MALEH" = 0?

# Creating DataFrames from other objects

## Creating Pandas DataFrames from Python Lists and Dictionaries

| | Dictionary | List |
|---|---|---|

**Row Oriented**

```
sales = [{'account': 'Jones LLC', 'Jan': 150, 'Feb': 200, 'Mar': 140},
         {'account': 'Alpha Co', 'Jan': 200, 'Feb': 210, 'Mar': 215},
         {'account': 'Blue Inc', 'Jan': 50, 'Feb': 90, 'Mar': 95 }]
df = pd.DataFrame(sales)
```

```
sales = [('Jones LLC', 150, 200, 50),
         ('Alpha Co', 200, 210, 90),
         ('Blue Inc', 140, 215, 95)]
labels = ['account', 'Jan', 'Feb', 'Mar']
df = pd.DataFrame.from_records(sales, columns=labels)
```

default

from_records

| | account | Jan | Feb | Mar |
|---|---|---|---|---|
| 0 | Jones LLC | 150 | 200 | 140 |
| 1 | Alpha Co | 200 | 210 | 215 |
| 2 | Blue Inc | 50 | 90 | 95 |

**Column Oriented**

```
sales = {'account': ['Jones LLC', 'Alpha Co', 'Blue Inc'],
         'Jan': [150, 200, 50],
         'Feb': [200, 210, 90],
         'Mar': [140, 215, 95]}
df = pd.DataFrame.from_dict(sales)
```

```
sales = [('account', ['Jones LLC', 'Alpha Co', 'Blue Inc']),
         ('Jan', [150, 200, 50]),
         ('Feb', [200, 210, 90]),
         ('Mar', [140, 215, 95])]
df = pd.DataFrame.from_items(sales)
```

from_dict

from_items

When using a dictionary, column order is not preserved.
Explicitly order them:
```
df = df[['account', 'Jan', 'Feb', 'Mar']]
```

Practical Business Python - pbpython.com

# Creating DataFrames from other objects, cont.

```
1  d = {'employee': ['Hannes', 'Fabiana', 'George', 'Olga'],
2       'group': ['Accounting', 'Engineering', 'Engineering', 'HR']}
3  df1 = pd.DataFrame.from_dict(d)
4  t = [('Hannes', 2004), ('Fabiana', 2008), ('George', 2012), ('Olga', 2014)]
5  df2 = pd.DataFrame.from_records(t, columns=["employee", "hire_date"])
```

# Appending, Concatening and Merging

```
1  df3 = df1.append(df2)
2  df4 = pd.concat([df1, df2])
3
4  df5 = pd.concat([df1, df2], axis=1)
5  df6 = df1.merge(df2, left_on="employee", right_on="employee")
```

# Appending, Concatening and Merging

```
1  df3 = df1.append(df2)
2  df4 = pd.concat([df1, df2])
3
4  df5 = pd.concat([df1, df2], axis=1)
5  df6 = df1.merge(df2, left_on="employee", right_on="employee")
```

- ► How do objects df3 and df4 differ?
- ► How do objects df5 and df6 differ?

# What is pivot?

## Pivot

df

```
df.pivot(index='foo',
         columns='bar',
         values='baz')
```

| | foo | bar | baz | zoo |
|---|---|---|---|---|
| 0 | one | A | 1 | x |
| 1 | one | B | 2 | y |
| 2 | one | C | 3 | z |
| 3 | two | A | 4 | q |
| 4 | two | B | 5 | w |
| 5 | two | C | 6 | t |

| bar | A | B | C |
|---|---|---|---|
| **foo** | | | |
| one | 1 | 2 | 3 |
| two | 4 | 5 | 6 |

from: "Reshaping and Pivot Tables"

# Pivoting and melting

```
1  pivoted = df6.pivot(index='employee', columns='group', values='hire_date')
2  reverse = (pivoted.reset_index()
3                     .melt(id_vars="employee", value_name="hire_date")
```

# Pivoting and melting

```
1  pivoted = df6.pivot(index='employee', columns='group', values='hire_date')
2  reverse = (pivoted.reset_index()
3                       .melt(id_vars="employee", value_name="hire_date")
```

- How do you make `reverse` look like `df6` again?

# Output

- https://pandas.pydata.org/pandas-docs/stable/user_guide/io.html

# Output

- https://pandas.pydata.org/pandas-docs/stable/user_guide/io.html

---

```
1  df.to_csv(FNAME, sep=";")
2  df.to_html(FNAME, decimal=",", justify="center")
3  df.to_stata(FNAME, write_index=False)
```

---