Pandas

- Open source, BSD-licensed library for high-performance, easy-to-use data structures and data analysis tools for Python
- Manipulates and visualizes data in spreadsheet

Reading in Data From CSV File

I have the following data saved in the file "Grades_Short.csv":

| | F30 | ÷ 😵 | ♥ (* f: | x | | | | | |
|----|-------|--------------|---------------|------------|------------|---------------|------------|-------|------------|
| _4 | Α | В | С | D | E | F | G | Н | 1 |
| 1 | Name | Previous_Par | Participation | Mini_Exam1 | Mini_Exam2 | Participation | Mini_Exam3 | Final | Grade |
| 2 | Jake | 32 | 1 | 19.5 | 20 | 1 | 10 | 33 | A |
| 3 | Joe | 32 | 1 | 20 | 16 | 1 | 14 | 32 | A |
| 4 | Susan | 30 | 1 | 19 | 19 | 1 | 10.5 | 33 | A - |
| 5 | Sol | 31 | 1 | 22 | 13 | 1 | 13 | 34 | Α |
| 6 | Chris | 30 | 1 | 19 | 17 | 1 | 12.5 | 33.5 | Α |
| 7 | Tarik | 31 | 1 | 19 | 19 | 1 | 8 | 24 | В |
| 8 | Malik | 31.5 | 1 | 20 | 21 | 1 | 9 | 36 | Α |
| 9 | | | | | | | | | |
| 10 | | | | | | | | | |

Let's see how we read this data into pandas:

Before you use pandas you must import it. Anytime you use pandas put this line as the top of your code.

```
import pandas as pd

df_grades = pd.read_csv("Grades_Short.csv")
```

Pandas – read_csv

pd.read_csv?

```
Signature: pd.read csv(filepath or buffer, sep=',', delimiter=None, header='infer', names=None, index col=None, usecols=None, squeeze=
False, prefix=None, mangle dupe cols=True, dtype=None, engine=None, converters=None, true values=None, false values=None, skipinitials
pace=False, skiprows=None, nrows=None, na values=None, keep default na=True, na filter=True, verbose=False, skip blank lines=True, par
se_dates=False, infer_datetime_format=False, keep_date_col=False, date_parser=None, dayfirst=False, iterator=False, chunksize=None, co
mpression='infer', thousands=None, decimal=b'.', lineterminator=None, quotechar='"', quoting=0, escapechar=None, comment=None, encodin
g=None, dialect=None, tupleize cols=None, error bad lines=True, warn bad lines=True, skipfooter=0, skip footer=0, doublequote=True, de
lim whitespace=False, as recarray=None, compact ints=None, use unsigned=None, low memory=True, buffer lines=None, memory map=False, fl
oat precision=None)
Docstring:
Read CSV (comma-separated) file into DataFrame
Also supports optionally iterating or breaking of the file
into chunks.
Additional help can be found in the `online docs for IO Tools
<http://pandas.pydata.org/pandas-docs/stable/io.html>`_.
Parameters
filepath_or_buffer : str, pathlib.Path, py._path.local.LocalPath or any object with a read() method (such as a file handle or StringI
0)
    The string could be a URL. Valid URL schemes include http, ftp, s3, and
    file. For file URLs, a host is expected. For instance, a local file could
    be file ://localhost/path/to/table.csv
sep : str, default ','
    Delimiter to use. If sep is None, the C engine cannot automatically detect
    the separator, but the Python parsing engine can, meaning the latter will
    be used and automatically detect the separator by Python's builtin sniffer
    tool, `csv.Sniffer`. In addition, separators longer than 1 character and
   different from ``'\s+'`` will be interpreted as regular expressions and
   will also force the use of the Python parsing engine. Note that regex
    delimiters are prone to ignoring quoted data. Regex example: ``'\r\t'``
```

Always specify the input name (order of inputs only matters if you don't)

delimiter : str, default `None``

Reading in Data From Excel

So, what is df_grades and how does it store the data?

```
import pandas as pd

df_grades = pd.read_csv("Grades_Short.csv")

df_grades
```

| | Name | Previous_Part | Participation1 | Mini_Exam1 | Mini_Exam2 | Participation2 | Mini_Exam3 | Final | Grade |
|---|-------|---------------|----------------|------------|------------|----------------|------------|-------|-------|
| 0 | Jake | 32.0 | 1 | 19.5 | 20 | 1 | 10.0 | 33.0 | Α |
| 1 | Joe | 32.0 | 1 | 20.0 | 16 | 1 | 14.0 | 32.0 | Α |
| 2 | Susan | 30.0 | 1 | 19.0 | 19 | 1 | 10.5 | 33.0 | A- |
| 3 | Sol | 31.0 | 1 | 22.0 | 13 | 1 | 13.0 | 34.0 | Α |
| 4 | Chris | 30.0 | 1 | 19.0 | 17 | 1 | 12.5 | 33.5 | Α |
| 5 | Tarik | 31.0 | 1 | 19.0 | 19 | 1 | 8.0 | 24.0 | В |
| 6 | Malik | 31.5 | 1 | 20.0 | 21 | 1 | 9.0 | 36.0 | Α |

- df_grades is a pandas dataframe.
- The data is stored in a tabular format very similar to excel.

The head() Method

Using the **head()** method

```
import pandas as pd

df_grades = pd.read_csv("Grades_Short.csv")

df_grades.head(3)
```

| | Name | Previous_Part | Participation1 | Mini_Exam1 | Mini_Exam2 | Participation2 | Mini_Exam3 | Final | Grade |
|---|-------|---------------|----------------|------------|------------|----------------|------------|-------|-------|
| 0 | Jake | 32.0 | 1 | 19.5 | 20 | 1 | 10.0 | 33.0 | Α |
| 1 | Joe | 32.0 | 1 | 20.0 | 16 | 1 | 14.0 | 32.0 | Α |
| 2 | Susan | 30.0 | 1 | 19.0 | 19 | 1 | 10.5 | 33.0 | Α- |

- If the data is really large you don't want to print out the entire dataframe to your output.
- The **head(n)** method outputs the first n rows of the data frame. If n is not supplied, the default is the first 5 rows.
- Run the head() method after reading in the dataframe to check that everything is read in correctly.
- There is also a tail(n) method that returns the last n rows of the dataframe

Basic Features

```
import pandas as pd

df_grades = pd.read_csv("Grades_Short.csv")

df_grades.head(3)
```

| | Name | Previous_Part | Participation1 | Mini_Exam1 | Mini_Exam2 | Participation2 | Mini_Exam3 | Final | Grade |
|---|-------|---------------|----------------|------------|------------|----------------|------------|-------|-------|
| 0 | Jake | 32.0 | 1 | 19.5 | 20 | 1 | 10.0 | 33.0 | Α |
| 1 | Joe | 32.0 | 1 | 20.0 | 16 | 1 | 14.0 | 32.0 | Α |
| 2 | Susan | 30.0 | 1 | 19.0 | 19 | 1 | 10.5 | 33.0 | A- |

```
#dimension of df
df_grades.shape
```

(7, 9)

#How each column is stored df_grades.dtypes

```
object
Name
Previous Part
                  float64
Participation1
                    int64
Mini Examl
                  float64
                    int64
Mini Exam2
Participation2
                    int64
Mini Exam3
                  float64
Final
                  float64
Grade
                   object
```

dtype: object

Basic Features

column names

| | _ | |
|---|---|---|
| | ı | |
| | ı | |
| ١ | ı | 1 |
| | ٧ | 1 |

| | Name | Previous_Part | Participation1 | Mini_Exam1 | Mini_Exam2 | Participation2 | Mini_Exam3 | Final | Grade |
|---|-------|---------------|----------------|------------|------------|----------------|------------|-------|-------|
| 0 | Jake | 32.0 | 1 | 19.5 | 20 | 1 | 10.0 | 33.0 | Α |
| 1 | Joe | 32.0 | 1 | 20.0 | 16 | 1 | 14.0 | 32.0 | Α |
| 2 | Susan | 30.0 | 1 | 19.0 | 19 | 1 | 10.5 | 33.0 | Α- |

row names = index

RangeIndex(start=0, stop=7, step=1)

Selecting a Single Column

| | Name | Previous_Part | Participation1 | Mini_Exam1 | Mini_Exam2 | Participation2 | Mini_Exam3 | Final | Grade |
|---|-------|---------------|----------------|------------|------------|----------------|------------|-------|-------|
| 0 | Jake | 32.0 | 1 | 19.5 | 20 | 1 | 10.0 | 33.0 | Α |
| 1 | Joe | 32.0 | 1 | 20.0 | 16 | 1 | 14.0 | 32.0 | Α |
| 2 | Susan | 30.0 | 1 | 19.0 | 19 | 1 | 10.5 | 33.0 | A- |
| 3 | Sol | 31.0 | 1 | 22.0 | 13 | 1 | 13.0 | 34.0 | Α |
| 4 | Chris | 30.0 | 1 | 19.0 | 17 | 1 | 12.5 | 33.5 | Α |
| 5 | Tarik | 31.0 | 1 | 19.0 | 19 | 1 | 8.0 | 24.0 | В |
| 6 | Malik | 31.5 | 1 | 20.0 | 21 | 1 | 9.0 | 36.0 | Α |

```
#Get Name column
df_grades['Name']
```

```
0    Jake
1    Joe
2    Susan
3    Sol
4    Chris
5    Tarik
6    Malik
Name: Name, dtype: object
```

- Between square brackets, the column must be given as a string
- Outputs column as a series
 - A series is a one-dimensional dataframe

Selecting a Single Column

| | Name | Previous_Part | Participation1 | Mini_Exam1 | Mini_Exam2 | Participation2 | Mini_Exam3 | Final | Grade |
|---|-------|---------------|----------------|------------|------------|----------------|------------|-------|------------|
| 0 | Jake | 32.0 | 1 | 19.5 | 20 | 1 | 10.0 | 33.0 | Α |
| 1 | Joe | 32.0 | 1 | 20.0 | 16 | 1 | 14.0 | 32.0 | Α |
| 2 | Susan | 30.0 | 1 | 19.0 | 19 | 1 | 10.5 | 33.0 | A - |
| 3 | Sol | 31.0 | 1 | 22.0 | 13 | 1 | 13.0 | 34.0 | Α |
| 4 | Chris | 30.0 | 1 | 19.0 | 17 | 1 | 12.5 | 33.5 | Α |
| 5 | Tarik | 31.0 | 1 | 19.0 | 19 | 1 | 8.0 | 24.0 | В |
| 6 | Malik | 31.5 | 1 | 20.0 | 21 | 1 | 9.0 | 36.0 | Α |

```
#Get Name column df_grades.Name
```

```
0    Jake
1    Joe
2    Susan
3    Sol
4    Chris
5    Tarik
6    Malik
Name: Name, dtype: object
```

- Exactly equivalent way to get Name column
 - +: don't have to type brackets or quotes
 - -: won't generalize to selecting multiple columns,, won't work if column names have spaces, can't create new columns this way

Selecting Multiple Columns

| | Name | Previous_Part | Participation1 | Mini_Exam1 | Mini_Exam2 | Participation2 | Mini_Exam3 | Final | Grade |
|---|-------|---------------|----------------|------------|------------|----------------|------------|-------|-------|
| 0 | Jake | 32.0 | 1 | 19.5 | 20 | 1 | 10.0 | 33.0 | Α |
| 1 | Joe | 32.0 | 1 | 20.0 | 16 | 1 | 14.0 | 32.0 | Α |
| 2 | Susan | 30.0 | 1 | 19.0 | 19 | 1 | 10.5 | 33.0 | A- |
| 3 | Sol | 31.0 | 1 | 22.0 | 13 | 1 | 13.0 | 34.0 | Α |
| 4 | Chris | 30.0 | 1 | 19.0 | 17 | 1 | 12.5 | 33.5 | Α |
| 5 | Tarik | 31.0 | 1 | 19.0 | 19 | 1 | 8.0 | 24.0 | В |
| 6 | Malik | 31.5 | 1 | 20.0 | 21 | 1 | 9.0 | 36.0 | Α |

| #Select multiple columns |
|---|
| <pre>df_grades[["Name", "Grade"]]</pre> |

| Name | Grade |
|-------|--------------------------------|
| Jake | Α |
| Joe | Α |
| Susan | A - |
| Sol | Α |
| Chris | Α |
| Tarik | В |
| Malik | Α |
| | Jake Joe Susan Sol Chris Tarik |

- List of strings, which correspond to column names.
- You can select as many column as you want.
- Columns don't have to be contiguous.

Storing Result

```
#Print the column
 df_grades["Name"]
     Jake
      Joe
2
    Susan
3
      Sol
    Chris
    Tarik
    Malik
Name: Name, dtype: object
 #Store the column
 names= df grades["Name"]
 names k
     Jake
      Joe
                The variable name stores a
2
    Susan
                series
      Sol
    Chris
    Tarik
    Malik
Name: Name, dtype: object
```

Why store a slice?

- We might want/have to do our analysis is steps.
 - Less error prone
 - More readable

Slicing a Series

```
Slice/index through
the index, which is
usually numbers
```

```
names= df_grades["Name"]
names

Jake
Joe
Susan
Sol
Chris
Tarik
Malik
Name: Name, dtype: object
```

Picking out single element

Contiguous slice

names[1:4]

Arbitrary slice

```
names[0]
```

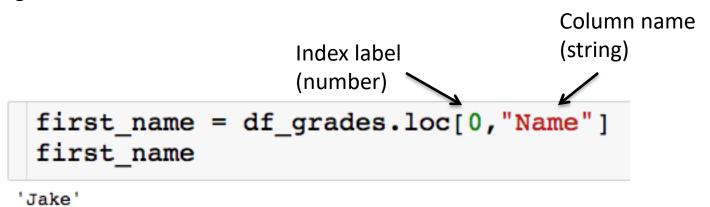
```
1 Joe
2 Susan
3 Sol
Name: Name, dtype: object
```

names[[1,2,4]]

1 Joe
2 Susan
4 Chris
Name: Name, dtype: object

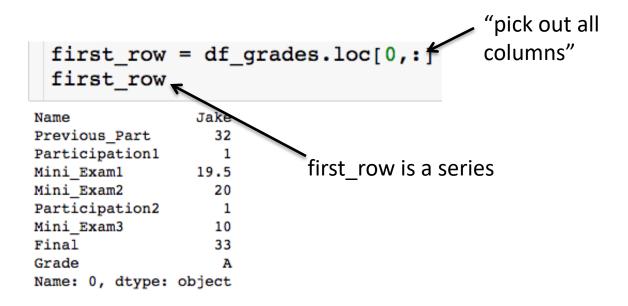
| | Name | Previous_Part | Participation1 | Mini_Exam1 | Mini_Exam2 | Participation2 | Mini_Exam3 | Final | Grade |
|---|-------|---------------|----------------|------------|------------|----------------|------------|-------|-------|
| 0 | Jake | 32.0 | 1 | 19.5 | 20 | 1 | 10.0 | 33.0 | Α |
| 1 | Joe | 32.0 | 1 | 20.0 | 16 | 1 | 14.0 | 32.0 | Α |
| 2 | Susan | 30.0 | 1 | 19.0 | 19 | 1 | 10.5 | 33.0 | A- |
| 3 | Sol | 31.0 | 1 | 22.0 | 13 | 1 | 13.0 | 34.0 | Α |
| 4 | Chris | 30.0 | 1 | 19.0 | 17 | 1 | 12.5 | 33.5 | Α |
| 5 | Tarik | 31.0 | 1 | 19.0 | 19 | 1 | 8.0 | 24.0 | В |
| 6 | Malik | 31.5 | 1 | 20.0 | 21 | 1 | 9.0 | 36.0 | Α |

Pick a single value out.



| | Name | Previous_Part | Participation1 | Mini_Exam1 | Mini_Exam2 | Participation2 | Mini_Exam3 | Final | Grade |
|---|-------|---------------|----------------|------------|------------|----------------|------------|-------|-------|
| 0 | Jake | 32.0 | 1 | 19.5 | 20 | 1 | 10.0 | 33.0 | Α |
| 1 | Joe | 32.0 | 1 | 20.0 | 16 | 1 | 14.0 | 32.0 | Α |
| 2 | Susan | 30.0 | 1 | 19.0 | 19 | 1 | 10.5 | 33.0 | A- |
| 3 | Sol | 31.0 | 1 | 22.0 | 13 | 1 | 13.0 | 34.0 | Α |
| 4 | Chris | 30.0 | 1 | 19.0 | 17 | 1 | 12.5 | 33.5 | Α |
| 5 | Tarik | 31.0 | 1 | 19.0 | 19 | 1 | 8.0 | 24.0 | В |
| 6 | Malik | 31.5 | 1 | 20.0 | 21 | 1 | 9.0 | 36.0 | Α |

Pick out entire row:



| | Name | Previous_Part | Participation1 | Mini_Exam1 | Mini_Exam2 | Participation2 | Mini_Exam3 | Final | Grade |
|---|-------|---------------|----------------|------------|------------|----------------|------------|-------|-------|
| 0 | Jake | 32.0 | 1 | 19.5 | 20 | 1 | 10.0 | 33.0 | Α |
| 1 | Joe | 32.0 | 1 | 20.0 | 16 | 1 | 14.0 | 32.0 | Α |
| 2 | Susan | 30.0 | 1 | 19.0 | 19 | 1 | 10.5 | 33.0 | A- |
| 3 | Sol | 31.0 | 1 | 22.0 | 13 | 1 | 13.0 | 34.0 | Α |
| 4 | Chris | 30.0 | 1 | 19.0 | 17 | 1 | 12.5 | 33.5 | Α |
| 5 | Tarik | 31.0 | 1 | 19.0 | 19 | 1 | 8.0 | 24.0 | В |
| 6 | Malik | 31.5 | 1 | 20.0 | 21 | 1 | 9.0 | 36.0 | Α |

Pick out contiguous chunk:

Endpoints are inclusive!

| | Name | Previous_Part | Participation1 | Mini_Exam1 | Mini_Exam2 |
|---|-------|---------------|----------------|------------|------------|
| 0 | Jake | 32.0 | 1 | 19.5 | 20 |
| 1 | Joe | 32.0 | 1 | 20.0 | 16 |
| 2 | Susan | 30.0 | 1 | 19.0 | 19 |

| | Name | Previous_Part | Participation1 | Mini_Exam1 | Mini_Exam2 | Participation2 | Mini_Exam3 | Final | Grade |
|---|-------|---------------|----------------|------------|------------|----------------|------------|-------|-------|
| 0 | Jake | 32.0 | 1 | 19.5 | 20 | 1 | 10.0 | 33.0 | Α |
| 1 | Joe | 32.0 | 1 | 20.0 | 16 | 1 | 14.0 | 32.0 | Α |
| 2 | Susan | 30.0 | 1 | 19.0 | 19 | 1 | 10.5 | 33.0 | A- |
| 3 | Sol | 31.0 | 1 | 22.0 | 13 | 1 | 13.0 | 34.0 | Α |
| 4 | Chris | 30.0 | 1 | 19.0 | 17 | 1 | 12.5 | 33.5 | Α |
| 5 | Tarik | 31.0 | 1 | 19.0 | 19 | 1 | 8.0 | 24.0 | В |
| 6 | Malik | 31.5 | 1 | 20.0 | 21 | 1 | 9.0 | 36.0 | Α |

Pick out arbitrary chunk:

| | Name | Grade | | |
|---|-------|------------|--|--|
| 0 | Jake | Α | | |
| 2 | Susan | A - | | |
| 3 | Sol | Α | | |