

# Lecture 03. Data Loading and Cleaning

Instructor: Luping Yu

Mar 12, 2024

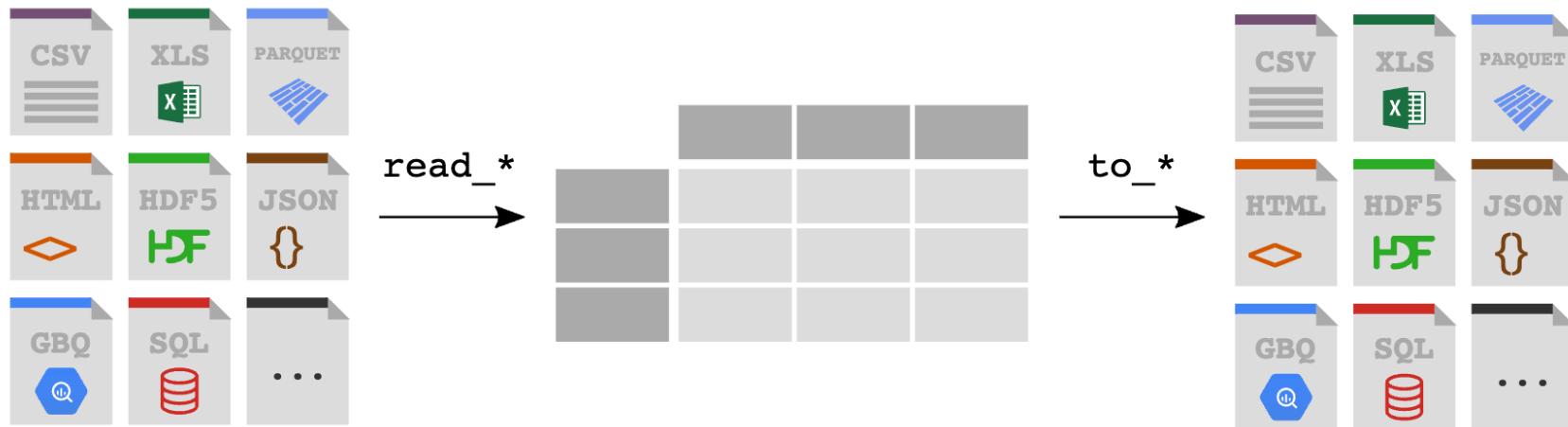
---

Accessing data is a necessary first step for using most of the tools in this course. We are going to be focused on data input and output using pandas.

---

## Reading and writing data in text format

`pandas` features a number of functions for reading **tabular data** as a `DataFrame` object.



The following table summarizes some of them, though `read_csv` is likely the ones you'll use the most.

## Function Description

**|read\_csv** | Load delimited data from a file, URL, or file-like object; use comma as default delimiter **|read\_excel** | Read tabular data from an Excel XLS or XLSX file **|read\_stata** | Read a dataset from Stata file format **|read\_sas** | Read a SAS dataset stored in one of the SAS system's custom storage formats **|read\_html** | Read all tables found in the given HTML document **|read\_json** | Read data from a JSON (JavaScript Object Notation) string representation **|read\_pickle** | Read an arbitrary object stored in Python pickle format **|read\_sql** | Read the results of a SQL query (using SQLAlchemy) as a pandas DataFrame

---

## Reading and Writing .csv (comma-separated values)

`.csv` is a delimited text file that uses a **comma** to separate values. A `.csv` file typically stores **tabular data** (numbers and text) in **plain text**.

Let's start with a small `.csv` text file: [ex1.csv](#)

```
In [1]: import pandas as pd
```

```
pd.read_csv('examples/ex1.csv') # relative path
```

```
Out[1]:
```

	a	b	c	d	message
0	1	2	3	4	hello
1	5	6	7	8	world
2	9	10	11	12	foo

```
In [2]: # absolute path (absolute path differs between Windows and Mac)
```

```
pd.read_csv('/Users/luping/desktop/teaching/2024_fdap/examples/ex1.csv')
```

```
Out [2]:
```

	a	b	c	d	message
0	1	2	3	4	hello
1	5	6	7	8	world
2	9	10	11	12	foo

`pandas.read_csv()` perform type inference. That means you don't necessarily have to specify which columns are numeric, integer, boolean, or string:

```
In [3]: df = pd.read_csv('examples/ex1.csv')
df.dtypes
```

```
Out [3]: a          int64
b          int64
c          int64
d          int64
message    object
dtype: object
```

A file will not always have a **header row**. Consider this file: [ex2.csv](#)

```
In [4]: pd.read_csv('examples/ex2.csv')
```

```
Out [4]:
```

	1	2	3	4	hello
0	5	6	7	8	world
1	9	10	11	12	foo

To read this file, you have a couple of options. You can allow pandas to assign default column names, or you can specify names yourself:

```
In [5]: pd.read_csv('examples/ex2.csv', names=['a', 'b', 'c', 'd', 'message'])
```

```
Out [5]:
```

	a	b	c	d	message
0	1	2	3	4	hello
1	5	6	7	8	world
2	9	10	11	12	foo

Suppose you wanted the message column to be the index of the returned `DataFrame`. You can use the `index_col` argument:

```
In [6]: pd.read_csv('examples/ex2.csv', names=['a', 'b', 'c', 'd', 'message'], index_col='message')
```

```
Out [6]:
```

message	a	b	c	d
hello	1	2	3	4
world	5	6	7	8
foo	9	10	11	12

Handling missing values is an important and frequently nuanced part of the file parsing process. Consider this file: [ex3.csv](#)

```
In [7]: pd.read_csv('examples/ex3.csv')
```

```
Out [7]:
```

	something	a	b	c	d	message
0	one	1	2	3.0	4	NaN
1	two	5	6	NaN	8	world
2	three	9	10	11.0	12	foo

Missing data is usually either not present (empty string) or marked by some **sentinel** value, such as **NA** and **NULL**.

```
In [8]: df = pd.read_csv('examples/ex3.csv')
```

```
pd.notnull(df)
```

```
Out[8]:
```

	something	a	b	c	d	message
0	True	True	True	True	True	False
1	True	True	True	False	True	True
2	True	True	True	True	True	True

Using `to_csv()` method, we can write the data out to a comma-separated file:

```
In [9]: df.to_csv('examples/out1.csv')
```

With no other options specified, both the row and column labels are written. Both of these can be disabled:

```
In [10]: df.to_csv('examples/out2.csv', index=False)
```

---

## Parameters of data loading functions

Because of how messy data in the real world can be, data loading functions (especially `read_csv()`) have grown very complex in their options over time. The **online pandas documentation** has many examples about how each of them works.

API reference (pandas documentation) of `read_csv()`: [https://pandas.pydata.org/docs/reference/api/pandas.read\\_csv.html](https://pandas.pydata.org/docs/reference/api/pandas.read_csv.html)

---

## Reading Microsoft excel files

`pandas` also supports reading tabular data stored in Excel 2003 (and higher) files using `pandas.read_excel()` function:

Internally these tools use the add-on packages **xlrd** and **openpyxl** to read XLS and XLSX files, respectively. You may need to install these manually with pip.

```
In [ ]: df = pd.read_excel('examples/ex1.xlsx', 'Sheet1')
```

```
df
```

To write pandas data to Excel format, you can pass a file path to `to_excel()` :

```
In [ ]: df.to_excel('examples/out1.xlsx')
```

---

## Data cleaning and preparation

During the course of doing data analysis and modeling, a significant amount of time is spent on data preparation: **loading**, **cleaning**, **transforming**, and **rearranging**. Such tasks are often reported to take up 80% or more of an analyst's time.

Sometimes the way that data is stored in files or databases is not in the right format for a particular task. Fortunately, `pandas` provides you with a high-level, flexible, and fast set of tools to enable you to manipulate data into the right form.

---

## Handling Missing Data

Missing data (**NA**, which stands for **not available**) occurs commonly in many data analysis applications. For numeric data, pandas uses the floating-point value **NaN** (not a number) to represent missing data.

With `DataFrame` objects, you may want to drop rows or columns that are all NA or only those containing any NAs.

`dropna()` by default drops any row containing a missing value:

```
In [13]: df = pd.DataFrame([[1., 6.5, 3.],  
                           [1., None, None],  
                           [None, None, None],  
                           [None, 6.5, 3.]])
```

```
df
```

```
Out[13]:
```

	0	1	2
0	1.0	6.5	3.0
1	1.0	NaN	NaN
2	NaN	NaN	NaN
3	NaN	6.5	3.0

```
In [14]: df.dropna()
```

```
Out[14]:
```

	0	1	2
0	1.0	6.5	3.0

Passing `how='all'` will only drop rows that are all NA:

```
In [15]: df.dropna(how='all')
```

```
Out[15]:
```

	0	1	2
0	1.0	6.5	3.0
1	1.0	NaN	NaN
3	NaN	6.5	3.0

To drop columns in the same way, pass `axis=1` :

```
In [16]: df[3] = None
```

```
df
```

```
Out[16]:
```

	0	1	2	3
0	1.0	6.5	3.0	None
1	1.0	NaN	NaN	None
2	NaN	NaN	NaN	None
3	NaN	6.5	3.0	None

```
In [17]: df.dropna(axis=1, how='all')
```

```
Out[17]:
```

	0	1	2
0	1.0	6.5	3.0
1	1.0	NaN	NaN
2	NaN	NaN	NaN
3	NaN	6.5	3.0

---

## Filling In Missing Data

Rather than filtering out missing data, you may want to fill in the "holes" in any number of ways. The `fillna` method is the function to use.

Calling `fillna` with a constant replaces missing values with that value:

```
In [18]: df = pd.DataFrame([[10, 30, 20, 40],
                           [8, 25, 15, 35],
                           [6, 20, 10, None],
                           [None, None, None, None],
                           [None, None, 10, 30]],
                           columns=['class participation', 'homework', 'midterm', 'final'])

df
```



```
Out[18]:
```

	class participation	homework	midterm	final
0	10.0	30.0	20.0	40.0
1	8.0	25.0	15.0	35.0
2	6.0	20.0	10.0	NaN
3	NaN	NaN	NaN	NaN
4	NaN	NaN	10.0	30.0

```
In [19]: df.fillna(5)
```

```
Out[19]:
```

	class participation	homework	midterm	final
0	10.0	30.0	20.0	40.0
1	8.0	25.0	15.0	35.0
2	6.0	20.0	10.0	5.0
3	5.0	5.0	5.0	5.0
4	5.0	5.0	10.0	30.0

Calling `fillna()` with a dict, you can use a different fill value for each column:

```
In [20]: df.fillna({'class participation': 5, 'final': 30})
```

```
Out[20]:
```

	class participation	homework	midterm	final
0	10.0	30.0	20.0	40.0
1	8.0	25.0	15.0	35.0
2	6.0	20.0	10.0	30.0
3	5.0	NaN	NaN	30.0
4	5.0	NaN	10.0	30.0

The **interpolation methods** can be used with fillna:

```
In [21]: df.fillna(method='ffill')
```

```
Out[21]:
```

	class participation	homework	midterm	final
0	10.0	30.0	20.0	40.0
1	8.0	25.0	15.0	35.0
2	6.0	20.0	10.0	35.0
3	6.0	20.0	10.0	35.0
4	6.0	20.0	10.0	30.0

With fillna you can do lots of other things with a little creativity. For example, you might pass the mean or median value of a Series:

```
In [22]: df.describe() #summary statistics
```

```
Out[22]:
```

	class participation	homework	midterm	final
<b>count</b>	3.0	3.0	4.000000	3.0
<b>mean</b>	8.0	25.0	13.750000	35.0
<b>std</b>	2.0	5.0	4.787136	5.0
<b>min</b>	6.0	20.0	10.000000	30.0
<b>25%</b>	7.0	22.5	10.000000	32.5
<b>50%</b>	8.0	25.0	12.500000	35.0
<b>75%</b>	9.0	27.5	16.250000	37.5
<b>max</b>	10.0	30.0	20.000000	40.0

```
In [23]: df.fillna(df.mean())
```

```
Out[23]:
```

	class participation	homework	midterm	final
0	10.0	30.0	20.00	40.0
1	8.0	25.0	15.00	35.0
2	6.0	20.0	10.00	35.0
3	8.0	25.0	13.75	35.0
4	8.0	25.0	10.00	30.0

---

## Removing Duplicates

Duplicate rows may be found in a `DataFrame` for any number of reasons. Here is an example:

```
In [24]: df = pd.DataFrame({'k1': ['one', 'two'] * 3 + ['two'],  
                           'k2': [1, 1, 2, 3, 3, 4, 4]})  
  
df
```

```
Out[24]:
```

	k1	k2
0	one	1
1	two	1
2	one	2
3	two	3
4	one	3
5	two	4
6	two	4

The `DataFrame` method `duplicated()` returns a boolean `Series` indicating whether each row is a duplicate (has been observed in a previous row) or not:

```
In [25]: df.duplicated()
```

```
Out[25]: 0    False
         1    False
         2    False
         3    False
         4    False
         5    False
         6     True
         dtype: bool
```

Relatedly, `drop_duplicates()` returns a `DataFrame` where the duplicated array is False:

```
In [26]: df.drop_duplicates()
```

```
Out[26]:
```

	k1	k2
0	one	1
1	two	1
2	one	2
3	two	3
4	one	3
5	two	4

`drop_duplicates()` considers all of the columns; alternatively, you can specify any **subset** of them to detect duplicates.

Suppose we had an additional column of values and wanted to filter duplicates only based on the 'k1' column:

```
In [27]: df['k3'] = range(7)
         df
```

```
Out[27]:
```

	k1	k2	k3
0	one	1	0
1	two	1	1
2	one	2	2
3	two	3	3
4	one	3	4
5	two	4	5
6	two	4	6

```
In [28]: df.drop_duplicates(['k1'])
```

```
Out[28]:
```

	k1	k2	k3
0	one	1	0
1	two	1	1

`drop_duplicates()` and `drop_duplicates()` by default keep the **first** observed value combination. Passing `keep='last'` will return the last one:

```
In [29]: df.drop_duplicates(['k1', 'k2'], keep='last')
```

```
Out[29]:
```

	k1	k2	k3
0	one	1	0
1	two	1	1
2	one	2	2
3	two	3	3
4	one	3	4
6	two	4	6

---

## Vectorized string functions in pandas

Cleaning up a messy dataset for analysis often requires a lot of string **regularization**. For example, a column containing strings will sometimes have missing data:

```
In [30]: df = pd.Series({'Dave': 'dave@google.com',  
                        'Jack': 'jack@xmu.edu.cn',  
                        'Steve': 'steve@gmail.com',  
                        'Rose': 'rose@xmu.edu.cn',  
                        'Tony': None})  
  
df
```

```
Out[30]: Dave      dave@google.com  
Jack      jack@xmu.edu.cn  
Steve     steve@gmail.com  
Rose      rose@xmu.edu.cn  
Tony              None  
dtype: object
```

`Series` has array-oriented methods for string operations that skip NA values. These are accessed through Series's `str` attribute.

For example, we could check whether each email address has 'xmu.edu' in it with `str.contains` :

```
In [31]: df.str.contains('xmu.edu')
```

```
Out[31]: Dave      False
         Jack      True
         Steve     False
         Rose      True
         Tony      None
         dtype: object
```

You can similarly **slice** strings using this syntax:

```
In [32]: df.str[:5]
```

```
Out[32]: Dave      dave@
         Jack      jack@
         Steve     steve
         Rose      rose@
         Tony      None
         dtype: object
```

```
In [33]: df.str.split('@')
```

```
Out[33]: Dave      [dave, google.com]
         Jack      [jack, xmu.edu.cn]
         Steve     [steve, gmail.com]
         Rose      [rose, xmu.edu.cn]
         Tony      None
         dtype: object
```

```
In [34]: df.str.split('@').str.get(0)
```

```
Out[34]: Dave      dave
         Jack      jack
         Steve     steve
         Rose      rose
         Tony      None
         dtype: object
```

Partial listing of vectorized string methods.

## Method Description

|cat|Concatenate strings element-wise with optional delimiter |contains|Return boolean array if each string contains pattern/regex  
|count|Count occurrences of pattern |extract|Use a regular expression with groups to extract one or more strings from a Series of strings  
|endswith|Equivalent to x.endswith(pattern) for each element |startswith|Equivalent to x.startswith(pattern) for each element  
|findall|Compute list of all occurrences of pattern/regex for each string |get|Index into each element (retrieve i-th element) |join|Join  
strings in each element of the Series with passed separator |len|Compute length of each string |lower,upper|Convert cases;equivalent to  
x.lower() or x.upper() for each element |match|Use re.match with the passed regular expression on each element |replace|Replace  
occurrences of pattern/regex with some other string |slice|Slice each string in the Series |split|Split strings on delimiter or regular  
expression |strip|Trim whitespace from both sides, including newlines