

# Data Cleaning with Python and Pandas

Matt Steele

# Resources

- [Pandas Documentation](#)
- [O'Reilly Learning Platform](#)

# Agenda

- Entering your data [lists, dictionaries, series, data frames]
- Reading data sets and cleaning data
- Analyzing data

# Recap

- Functions and Arguments
- Variables
- Run Code

# Sequence Types

ordered collections of data items of the same type.

[Python Documentation - Sequence Types](#)

Tuples	A tuple is a collection which is ordered and unchan
--------	---

---

Sets	A set is a collection which is unordered, unchangea
------	---

---

Dictionaries	A dictionary is a collection which is ordered*, chang
--------------	---

# Dictionary

A dictionary is a collection which is ordered\*, changeable and does not allow duplicates.

```
1 country_roads = {"Morgantown": 30847, "Charlestown": 45879,  
2 "Reedsville": 603, "Huntington": 44934}  
3  
4 print(country_roads)
```

# Data science libraries in Python

Listed below are the major libraries that provide built-in functions, methods, and constants that are important

## Storage, M

- [Numpy](#)
  - [Pandas](#)
  - [Scipy](#)
  - [StatsMode](#)
-

# Series

Like a dictionary in the standard library, a series from Pandas allows you to store key-value pairs in python.

```
1 import pandas as pd
2
3 series_example = pd.Series([1, 2, 3, 4], index=['a', 'b', 'c', 'd'])
4 series_example
```



# Data Frame

A **Pandas DataFrame** is a two-dimensional, size-mutable, and heterogeneous tabular data structure. It is composed of:

- **Variables** (or columns), which represent the data types.
- **Observations** (or rows), which represent individual data entries.

Each variable (column) in a DataFrame typically contains data of the same type (e.g., integers, strings). However, different variables can contain different data types.

# Step 1: Create Series

- All Series must have the same size to combine
- For this example, we're creating four Series: year

```
1 import pandas as pd
2
3 # Creating individual Series for each column
4 year = pd.Series([1977, 1980, 1983], name="Year")
5
6 title = pd.Series(["Star Wars", "Empire Strikes Back", "F
7 name="Title")
8
9 length = pd.Series([121, 124, 144], name="Length")
10
11 gross = pd.Series([787, 534, 572], name="Gross")
```

# Step 2: Merge the Series into a data frame

```
1 # Creating the DataFrame
2 starwars_df = pd.DataFrame({
3     "Year": year,
4     "Title": title,
5     "Length": length,
6     "Gross": gross
7 })
8
9 # Displaying the DataFrame
10 print(starwars_df)
```

- After creating the Series, you can merge them
- Each Series becomes a column in the DataFrame

# Subsetting Variables

used to select and work with specific variables (columns) from a data frame.

```
1 starwars_df['Title']  
2  
3 mean_length = starwars_df["Length"].mean()  
4  
5 # Displaying the result  
6  
7 print("Mean Length of Movies:", mean_length)
```

# Export / Save Data

Once you are done entering your data, you can export it to your working directory. The function without built-in arguments is **write.table( )** but if are saving it as a csv, you are better using **write.csv( )**.

```
1 # Exporting the DataFrame to a CSV file
2 starwars_df.to_csv("starwars.csv", index=False)
```

# Importing/Reading Data



<code>read_csv</code>	Load delimited data from a file, URL, or database
<code>read_fwf</code>	Read data in fixed-width column format
<code>read_excel</code>	Read tabular data from an Excel XLS or XLSX file
<code>read_html</code>	Read all tables found in the given HTML document
<code>read_json</code>	Read data from a JSON (JavaScript Object Notation) file
<code>read_sas</code>	Read a SAS dataset stored in one of the SAS file formats
<code>read_spss</code>	Read a data file created by SPSS
<code>read_stata</code>	Read a data set from Stata file format
<code>read_xml</code>	Read a table of data from an XML file

# Read Non-Proprietary Data

- Non-proprietary files, like CSV, are open formats that can be used and shared across different software platforms.
- CSV (Comma-Separated Values) files store data in a plain text format where values are separated by commas. They are commonly used because they are simple and widely supported.

```
1 import pandas as pd
2
3 reviews = pd.read_csv("customer_reviews.csv")
```

# Read Proprietary Data

- Proprietary files, such as SPSS (.sav), are used by specific software systems and often require special tools to open.
- Python can handle proprietary formats like SPSS without needing the original software (SPSS), making it a versatile tool for data analysis.

```
1 import pandas as pd
2
3 demographics = pd.read_spss("demographics.sav")
4 demographics
```

# View Environment

The environment in Python is where all your variables, functions, and imported libraries are stored during your current session.

```
1 %whos
```

# Explore a Data Frame

syntax	example	description
<code>.head()</code>	<code>reviews.head()</code>	displays the first 5 rows of the DataFrame
<code>.dtypes</code>	<code>reviews.dtypes</code>	shows the data type of each column
<code>.info()</code>	<code>reviews.info()</code>	provides a summary of the DataFrame
<code>.shape</code>	<code>reviews.shape</code>	returns the dimensions of the DataFrame
<code>.columns</code>	<code>reviews.columns</code>	gives you a list of all column names
<code>.describe</code>	<code>reviews.describe()</code>	generates summary statistics for numeric columns
<code>.value_counts()</code>	<code>reviews["Class_Name"].value_counts()</code>	counts the occurrences of each unique value

# Descriptive Statistics.

**syntax**

**example**

`.count()`

`reviews.count()`

`.first(), .last()`

`reviews['Age'].first(), reviews['Age'].last()`

`.mean(), .median()`

`reviews['Age'].mean(), reviews['Age'].median()`

`.min(), .max()`

`reviews['Age'].min(), reviews['Age'].max()`

`.std(), .var()`

`reviews['Age'].std(), reviews['Age'].var()`

# Cleaning Data

- Data in columns and rows are not ordered in the correct way
  - Creating values or ignoring missing data
  - Units are not correct or are wrong in some way
  - Order of magnitude is off
  - Outliers and skewing of the data
- Dplyr function: filter



# Filter

allows you to *select rows* in your data frame that meet specific conditions or criteria in a variable

```
1 #find the mean of rating that people who bought from the General Department
2
3 reviews_filter = reviews[reviews["Division_Name"]== "General"]
```

# Boolean Operators

boolean operators allow you to build criteria in your code

==	EQUAL
----	-------

!=	NOT EQUAL
----	-----------

<	LESS THAN
---	-----------

>	GREATER THAN
---	--------------

<=	LESS THAN OR EQUAL
----	--------------------

>=	GREATER THAN OR EQUAL
----	-----------------------

# Filter with Boolean

let's filter the data frame for characters who have blue eyes and were born after 50 BBY

```
1 reviews_filter_2 = reviews[(reviews["Division_Name"] == "General") &  
2 (reviews["Age"] < 40)]
```

# Select

allows you to *keep* or *discard* variables

```
1 # view
2
3 reviews.columns
4
5 # select variables
6
7 reviews_select = reviews[['Age', 'Title', 'Division_Name',
8                             'Department_Name', 'Class_Name']]
9
10 reviews_select
```

# Assign

*creates* new variables in your data or *change* existing variables by performing calculations or transformations.

```
1 demographics
2
3 demographics = demographics.assign(income = (demographics["income"]/1000))
4 demographics
```

NOTE: if you name your variable as an *existing variable*, it will *overwrite* the existing variable. If you give it a *new name*, it will create a *new variable*

# Recode

# Transform the values of a variable

```
1 # view observation categories
2
3 reviews["Recommended_IND"].value_counts()
4
5 # create a new variable
6
7 reviews = reviews.assign(Recommended_IND=
8 reviews["Recommended_IND"].map({0: "Not Recommended",
9
```

# Recode Data Values

The `.astype` function will allow you to change the data type of a variable.

```
1 # view types of data values in the dataframe
2
3 reviews.dtypes
4
5 # Clothing ID should be a string
6 reviews['Clothing_ID'] = reviews['Clothing_ID'].astype("str")
7
8 # Recommended_IND, Division_Name, Department_Name, and Class_Name should be
9 reviews[['Recommended_IND', 'Division_Name', 'Department_Name', 'Class_Name
10 reviews[['Recommended_IND', 'Division_Name', 'Department_Name', 'Class_Name
11 astype("category")
12
13 # view the values
14
15 reviews.dtypes
```



# Rename

Rename the column

```
1 reviews =  
2 reviews.rename(columns={"Recommended_IND": "Recommended_num",  
3 "Recommended_recode": "Recommended_label"})  
4  
5 reviews
```

# Relocate

Move the column location in the data frame

```
1 reviews.columns
2
3 new_order = ['Clothing_ID', 'Age', 'Title',
4             'Review_Text', 'Rating', 'Recommended_num', 'Recommended_label',
5             'Positive_Feedback_Count', 'Division_Name',
6             'Department_Name', 'Class_Name']
7
8 reviews = reviews[new_order]
9 reviews
```

# Sort

allows you to *sort* variables

```
1 reviews.sort_values(by="Rating", ascending=False)
```

# Group\_by & aggregate

the **group\_by** function allows you to group

**aggregate** function allows you to get descr

```
1  #Groupby
2
3  class_rating = reviews[["Rating", "Class_Name"]].g
4  class_rating
5
6
7  #class_rating = reviews[["Rating", "Class_Name"]].
8  #class_rating
9
```

# Missing Data

missing data in numeric fields can cause an issue when trying to calculate descriptive statistics

```
1 # are there missing NA values
2
3 number_missing = reviews.isna().sum()
4 number_missing
```

# dropna

removes **all** missing data from data frames or variables

```
1 reviews.shape
2
3 #remove all observations with na
4
5 reviews_na = reviews.dropna()
6
7 number_missing = reviews_na.isna().sum()
8 number_missing
```

we can also just drop NAs from a **variable**

```
1 reviews_na = reviews[reviews['Title'].notna()]
2 reviews_na
3
4 number_missing = reviews_na.isna().sum()
5 number_missing
```

# fillna

you can also recode the NA values for observations with *fillna*

```
1 reviews["Title"] = reviews["Title"].fillna("None Given")  
2 reviews["Title"]
```



# Export: to\_csv

the `to_csv` function allows us to export data frames to a csv file once we are done cleaning it up or when we have done some analysis that we want to export

```
1 # now that we have this date frame cleaned let's save it
2
3 # let's export the file
4
5 reviews.to_csv("cleaned_reviews.csv")
```

