# 03 | *pandas*

*Ivan Corneillet*

*Data Scientist*

# Learning Objectives

After this lesson, you should be able to:

‣ Define a data science problem

‣ Write a Jupyter notebook to import, format, and clean data using *pandas*

# Business Understanding

# By asking a good question and setting a clear aim:



- You set yourself up for success
    - "A problem well stated is half solved" – Charles Kettering
- You help other data scientists learn from and reproduce your work
    - You establish the basis for making your analysis reproducible
- You also help them expand on your work in the future

# The SMART Goals Framework for Data Science
([https://en.wikipedia.org/wiki/SMART_criteria](https://en.wikipedia.org/wiki/SMART_criteria))

| | |
|---|---|
| **S**PECIFIC | The dataset and key variables are clearly defined |
| **M**EASURABLE | The type of analysis and major assumptions are articulated |
| **A**TTAINABLE | The question you are asking is feasible for your dataset and is not likely to be biased |
| **R**EPRODUCIBLE | Another person (or you in 6 months!) can read your state and understand exactly how your analysis is performed |
| **T**IME-BOUND | You clearly state the time period and population for which this analysis will pertain |

Trends often change over time and vary by the population of source of your data.  It is important to clearly define who/what you included in your analysis as well as the time period for the analysis

# Models, Feature Matrix $X$, Response Vector $y$, and Tidy Data

Before **modeling**, our data needs to be **tidy** and in the form of a **feature matrix $X$** (i.e., the stimuli, e.g., *"ring bell"*) and a **response vector $y$** (i.e., the response, e.g., *"dog salivates"*)

**Feature Matrix $X$**  **Response Vector $y$**

|  | col0 | col1 | col2 | col3 |
| --- | --- | --- | --- | --- |
| row0 |  |  |  |  |
| row1 |  |  |  |  |
| row2 |  |  |  |  |
| row3 |  |  |  |  |

|  | col |
| --- | --- |
| row0 |  |
| row1 |  |
| row2 |  |
| row3 |  |

# San Francisco Dataset: a dataset we will use throughout this course
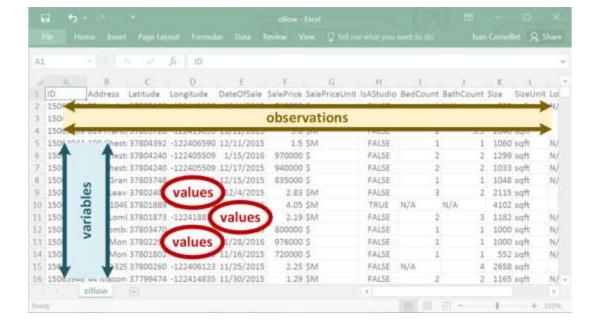
CASE STUDY

‣ **Recently Sold Homes (Source: Zillow)**

   ‣ 1,000 homes sold in San Francisco between 11/10/2015 and 2/12/2016

# What is Tidy Data?

‣ Your data is tidy if you follow these three rules:

  ‣ Each **sample** (or **observation**) in the dataset is placed in its own **row**

  ‣ Each **feature** (or **variable**) is placed in its own **column**
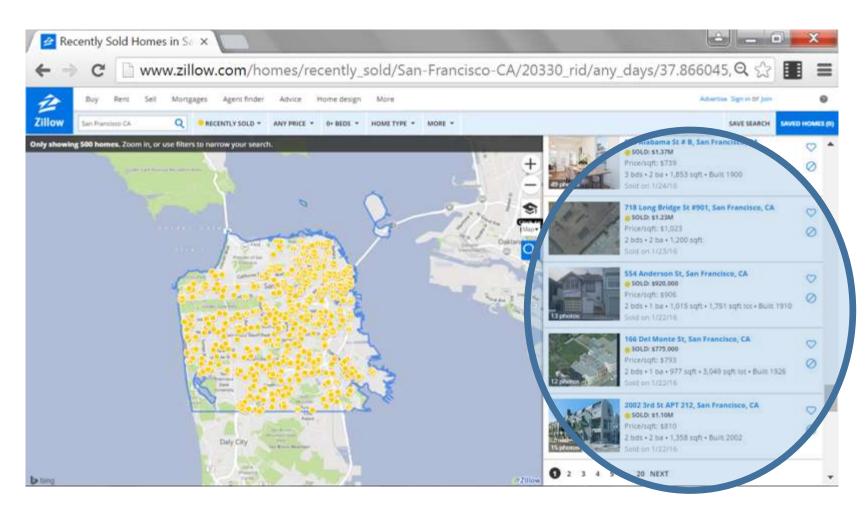
  ‣ Each **value** is placed in its own **cell**

# Unfortunately, data will usually come to you **raw**, i.e., **unstructured…**

**CASE STUDY**

```
<div class="property-info"
id="yui_3_18_1_1_1456167242885_71870"><strong
id="yui_3_18_1_1_1456167242885_71869"><dt class="property-address"
id="yui_3_18_1_1_1456167242885_71868"><a href="/homedetails/149-
Shipley-St-San-Francisco-CA-94107/15147894_zpid/" class="hdp-link
routable" title="149 Shipley St, San Francisco, CA Real Estate"
id="yui_3_18_1_1_1456167242885_71873">149 Shipley St, San
Francisco, CA</a></dt></strong><dt class="listing-type zsg-
content_collapsed" id="yui_3_18_1_1_1456167242885_71875"><span
class="zsg-icon-recently-sold type-icon"></span>Sold:
$1.18M</dt><dt class="zsg-fineprint"
id="yui_3_18_1_1_1456167242885_71877">Price/sqft: $1,116</dt><dt
class="property-data" id="yui_3_18_1_1_1456167242885_71880"><span
class="beds-baths-sqft">3 bds • 2 ba • 1,057 sqft</span><span
class="built-year" id="yui_3_18_1_1_1456167242885_71879"> • Built
1992</span></dt><dt class="sold-date zsg-fineprint"
id="yui_3_18_1_1_1456167242885_71975">Sold on 2/22/16</dt></div>
```

# (E.g., raw/unstructured scrapped data)

# ... and/or **messy**...

EXAMPLE

‣ Trouble tickets inspect and maintain manholes in New Year City

‣ "Service box," a common piece of infrastructure, had at least 38 variants, including `SB, S, S/B, S.B, S?B, S.B.`, `SBX, S/BX, SB/X, S/XB, /SBX, S.BX, S &BX, S?BX, S BX, S/B/X, S BOX, SVBX, SERV BX, SERV-BOX, SERV/BOX,` and `SERVICE BOX`

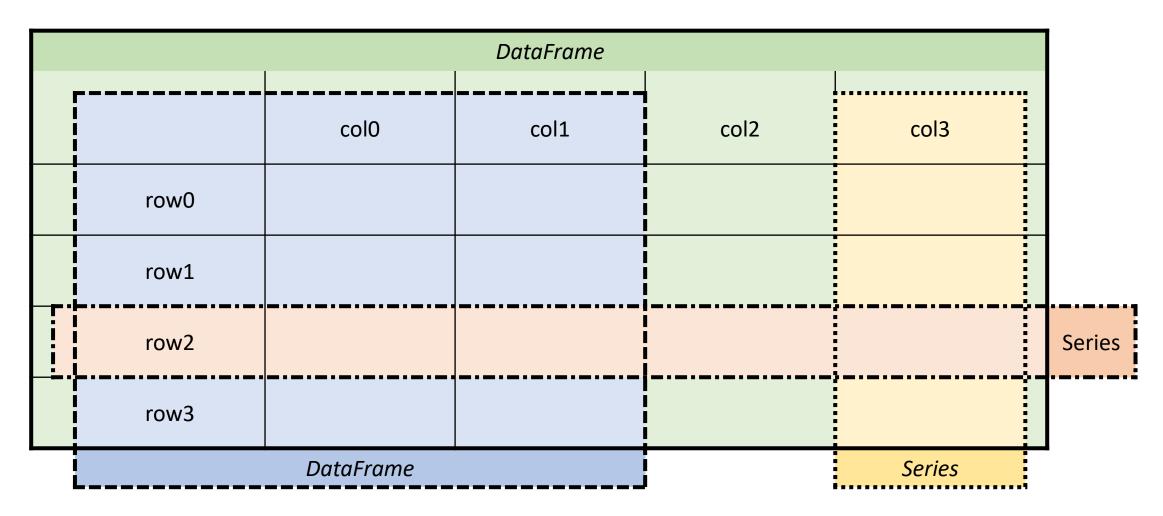Source: Big Data: A Revolution That Will Transform How We Live, Work, and Think

# Question: What tool can we use to wrangle **raw data** into **tidy data**?

- Answer: *pandas*

  - *pandas* is a Python library that provides the ability to *index*, *retrieve*, *tidy*, *reshape*, *combine*, *slice*, tabular and other multidimensional datasets

  - *pandas* also provides facilities to perform statistical and mathematical analysis which will come handy for exploratory data analysis

- Wrangling data is the most fruitful skill you can learn as a data scientist. It will save you hours of time and make your data much easier to visualize, manipulate, and model

- Today, we will use *pandas* to explore and manipulate the San Francisco housing dataset

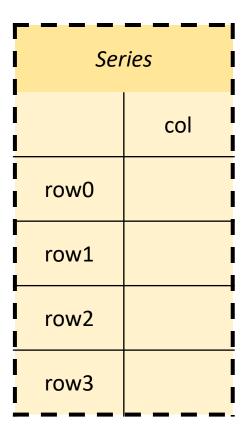# *pandas*

# *pandas.DataFrame* and *pandas.Series* (cont.)

When reaching the **modeling** step, our **feature matrix $X$** will be modeled as a DataFrame and the **response vector $y$** as a Series

**Feature Matrix $X$**

| DataFrame | | | | |
|---|---|---|---|---|
| | col0 | col1 | col2 | col3 |
| row0 | | | | |
| row1 | | | | |
| row2 | | | | |
| row3 | | | | |

**Response Vector $y$**

| Series | |
|---|---|
| | col |
| row0 | |
| row1 | |
| row2 | |
| row3 | |

Slides © 2017 Ivan Corneillet Where Applicable
Do Not Reproduce Without Permission