

Data wrangling with Tableau and Excel

October 11 2016
JRNL 520H

What is data wrangling?

Data wrangling is the process of preparing raw data for use in a data analysis or visualization software.

What are the causes of dirty data?

- Data entry error

What are the causes of dirty data?

- Data entry error
- Incompatible tables

What are the causes of dirty data?

- Data entry error
- Incompatible tables
- Incompatible table format

What should we look out for when cleaning data?

- Table formatting

What should we look out for when cleaning data?

- Table formatting
- Variable type

What should we look out for when cleaning data?

- Table formatting
- Variable type
- Invalid character values

What should we look out for when cleaning data?

- Table formatting
- Variable type
- Invalid character values
- Invalid numeric values

What should we look out for when cleaning data?

- Table formatting
- Variable type
- Invalid character values
- Invalid numeric values
- Grouping data

What should we look out for when cleaning data?

- Table formatting
- Variable type
- Invalid character values
- Invalid numeric values
- Grouping data
- Missing values

Ideal format of data in Tableau

- Start your data in cell A1. Remove all introductory information and footnotes.
- Have the first row be the column headers/variable names
- Have every subsequent row be one observation. No cross-tabulation!

Ideal format of data in Tableau

Before				After			
City	Year	Pop	Pop Chg	City	Year	Pop	Pop Chg
Albuquerque	2000	500	0.00				
Albuquerque	2001	500	0.00				
Albuquerque	2002	500	0.00				
Albuquerque	2003	500	0.00				
Albuquerque	2004	500	0.00				
Albuquerque	2005	500	0.00				
Albuquerque	2006	500	0.00				
Albuquerque	2007	500	0.00				
Albuquerque	2008	500	0.00				
Albuquerque	2009	500	0.00				
Albuquerque	2010	500	0.00				
Albuquerque	2011	500	0.00				
Albuquerque	2012	500	0.00				
Albuquerque	2013	500	0.00				
Albuquerque	2014	500	0.00				
Albuquerque	2015	500	0.00				
Albuquerque	2016	500	0.00				
Albuquerque	2017	500	0.00				
Albuquerque	2018	500	0.00				
Albuquerque	2019	500	0.00				
Albuquerque	2020	500	0.00				
Albuquerque	2021	500	0.00				
Albuquerque	2022	500	0.00				
Albuquerque	2023	500	0.00				
Albuquerque	2024	500	0.00				
Albuquerque	2025	500	0.00				
Albuquerque	2026	500	0.00				
Albuquerque	2027	500	0.00				
Albuquerque	2028	500	0.00				
Albuquerque	2029	500	0.00				
Albuquerque	2030	500	0.00				

Ideal format of data in Tableau

Before				After			
City	Year	Pop	Pop Chg	City	Year	Pop	Pop Chg
Albuquerque	2000	500	0.00				
Albuquerque	2001	500	0.00				
Albuquerque	2002	500	0.00				
Albuquerque	2003	500	0.00				
Albuquerque	2004	500	0.00				
Albuquerque	2005	500	0.00				
Albuquerque	2006	500	0.00				
Albuquerque	2007	500	0.00				
Albuquerque	2008	500	0.00				
Albuquerque	2009	500	0.00				
Albuquerque	2010	500	0.00				
Albuquerque	2011	500	0.00				
Albuquerque	2012	500	0.00				
Albuquerque	2013	500	0.00				
Albuquerque	2014	500	0.00				
Albuquerque	2015	500	0.00				
Albuquerque	2016	500	0.00				
Albuquerque	2017	500	0.00				
Albuquerque	2018	500	0.00				
Albuquerque	2019	500	0.00				
Albuquerque	2020	500	0.00				
Albuquerque	2021	500	0.00				
Albuquerque	2022	500	0.00				
Albuquerque	2023	500	0.00				
Albuquerque	2024	500	0.00				
Albuquerque	2025	500	0.00				
Albuquerque	2026	500	0.00				
Albuquerque	2027	500	0.00				
Albuquerque	2028	500	0.00				
Albuquerque	2029	500	0.00				
Albuquerque	2030	500	0.00				

Data Interpreter

Tableau's Data Interpreter feature draws out sub-tables and removes some of that extraneous information to help prepare your data source for analysis. Note: the data interpreter only works with Microsoft Excel files, not CSV or other file types.

Data Interpreter

Tableau's Data Interpreter feature draws out sub-tables and removes some of that extraneous information to help prepare your data source for analysis. Note: the data interpreter only works with Microsoft Excel files, not CSV or other file types.

Complete Tableau exercise

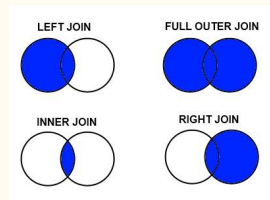
Joins

A JOIN is a means for combining columns from one or more tables by using values common to each. There are four main join types: inner, left, right and full outer.

Joins

Market rental pricing for July 2016					Market rental pricing for June 2016				
1 Bedroom		2 Bedrooms			1 Bedroom		2 Bedrooms		
City	Price	MM %	Price	MM %	City	Price	MM %	Price	MM %
Vancouver	1740	0.024	2750	-0.011	Vancouver	\$1,700	-2.30%	\$2,780	-0.70%
Toronto	1300	0.023	1720	0.042	Toronto	\$1,300	1.90%	\$1,850	-1.20%
Calgary	1110	0.037	1340	0.031	Calgary	\$1,100	1.90%	\$1,400	2.90%
Victoria	1075	-0.023	1370	-0.021	Victoria	\$1,070	-3.70%	\$1,300	-3.70%
Ottawa	1040	0.03	1280	-0.015	Ottawa	\$1,010	-1.90%	\$1,300	-1.90%
Edmonton	1010	0.031	1200	0.016	Edmonton	\$980	3.20%	\$1,200	2.00%
Regina	980	0.021	1180	0	Regina	\$960	-4.00%	\$1,180	0.00%
Montreal	960	0.011	1200	0.042	Montreal	\$950	-2.20%	\$1,200	-4.00%
Kingston	940	0.011	1070	-0.037	Kingston	\$900	0.00%	\$1,100	4.80%
Kelowna	920	0.022	1400	0.029	Kelowna	\$910	-4.20%	\$1,130	1.80%
Barrie	910	0.034	1300	0.036	Barrie	\$900	-2.20%	\$1,200	1.70%
Halifax	910	0.022	1100	0.019	Halifax	\$900	1.10%	\$1,360	1.00%

Joins



Joins

Market rental pricing for July 2016					Market rental pricing for June 2016				
1 Bedroom		2 Bedrooms			1 Bedroom		2 Bedrooms		
City	Price	MM %	Price	MM %	City	Price	MM %	Price	MM %
Vancouver	1740	0.024	2750	-0.011	Vancouver	\$1,700	-2.30%	\$2,780	-0.70%
Toronto	1300	0.023	1720	0.042	Toronto	\$1,300	1.90%	\$1,850	-1.20%
Calgary	1110	0.037	1340	0.031	Calgary	\$1,100	1.90%	\$1,400	2.90%
Victoria	1075	-0.023	1370	-0.021	Victoria	\$1,070	-3.70%	\$1,300	-3.70%
Ottawa	1040	0.03	1280	-0.015	Ottawa	\$1,010	-1.90%	\$1,300	-1.90%
Edmonton	1010	0.031	1200	0.016	Edmonton	\$980	3.20%	\$1,200	2.00%
Regina	980	0.021	1180	0	Regina	\$960	-4.00%	\$1,180	0.00%
Montreal	960	0.011	1200	0.042	Montreal	\$950	-2.20%	\$1,200	-4.00%
Kingston	940	0.011	1070	-0.037	Kingston	\$900	0.00%	\$1,100	4.80%
Kelowna	920	0.022	1400	0.029	Kelowna	\$910	-4.20%	\$1,130	1.80%
Barrie	910	0.034	1300	0.036	Barrie	\$900	-2.20%	\$1,200	1.70%
Halifax	910	0.022	1100	0.019	Halifax	\$900	1.10%	\$1,360	1.00%

Joins

Market rental pricing for July 2016					Market rental pricing for June 2016				
1 Bedroom		2 Bedrooms			1 Bedroom		2 Bedrooms		
City	Price	MM %	Price	MM %	City	Price	MM %	Price	MM %
Vancouver	1740	0.024	2750	-0.011	Vancouver	\$1,700	-2.30%	\$2,780	-0.70%
Toronto	1300	0.023	1720	0.042	Toronto	\$1,300	1.90%	\$1,850	-1.20%
Calgary	1110	0.037	1340	0.031	Calgary	\$1,100	1.90%	\$1,400	2.90%
Victoria	1075	-0.023	1370	-0.021	Victoria	\$1,070	-3.70%	\$1,300	-3.70%
Ottawa	1040	0.03	1280	-0.015	Ottawa	\$1,010	-1.90%	\$1,300	-1.90%
Edmonton	1010	0.031	1200	0.016	Edmonton	\$980	3.20%	\$1,200	2.00%
Regina	980	0.021	1180	0	Regina	\$960	-4.00%	\$1,180	0.00%
Montreal	960	0.011	1200	0.042	Montreal	\$950	-2.20%	\$1,200	-4.00%
Kingston	940	0.011	1070	-0.037	Kingston	\$900	0.00%	\$1,100	4.80%
Kelowna	920	0.022	1400	0.029	Kelowna	\$910	-4.20%	\$1,130	1.80%
Barrie	910	0.034	1300	0.036	Barrie	\$900	-2.20%	\$1,200	1.70%
Halifax	910	0.022	1100	0.019	Halifax	\$900	1.10%	\$1,360	1.00%

Complete Tableau exercise

Wrangling in Excel

Sometimes the data interpreter in Tableau isn't able to detect all of the errors in the dataset. In cases like this, you will need to manually clean the data in Excel.

Complete Tableau exercise

Pivot

Tabular format		Columnar format	
ID	Year	Revenue	%
Canada2	1992	568	100%
Provinces	1992	568	100%
Newfoundland and L.	1992	537	94%
Prince Edward Island	1992	503	88%
Alberta	1992	543	95%
British Columbia	1992	562	99%
New Scotia	1992	562	99%
Ontario	1992	572	101%
Quebec	1992	473	83%
New Brunswick	1992	461	81%
Manitoba	1992	514	90%
Saskatchewan	1992	454	80%
Alberta	1992	553	97%
British Columbia	1992	682	120%
Metropolitan Areas	1992	568	100%

Complete Tableau exercise

Pivot

Tabular format		Columnar format	
ID	Year	Revenue	%
Canada2	1992	568	100%
Provinces	1992	568	100%
Newfoundland and L.	1992	537	94%
Prince Edward Island	1992	503	88%
Alberta	1992	543	95%
British Columbia	1992	562	99%
New Scotia	1992	562	99%
Ontario	1992	572	101%
Quebec	1992	473	83%
New Brunswick	1992	461	81%
Manitoba	1992	514	90%
Saskatchewan	1992	454	80%
Alberta	1992	553	97%
British Columbia	1992	682	120%
Metropolitan Areas	1992	568	100%

Complete Tableau exercise