



The Shape of Data

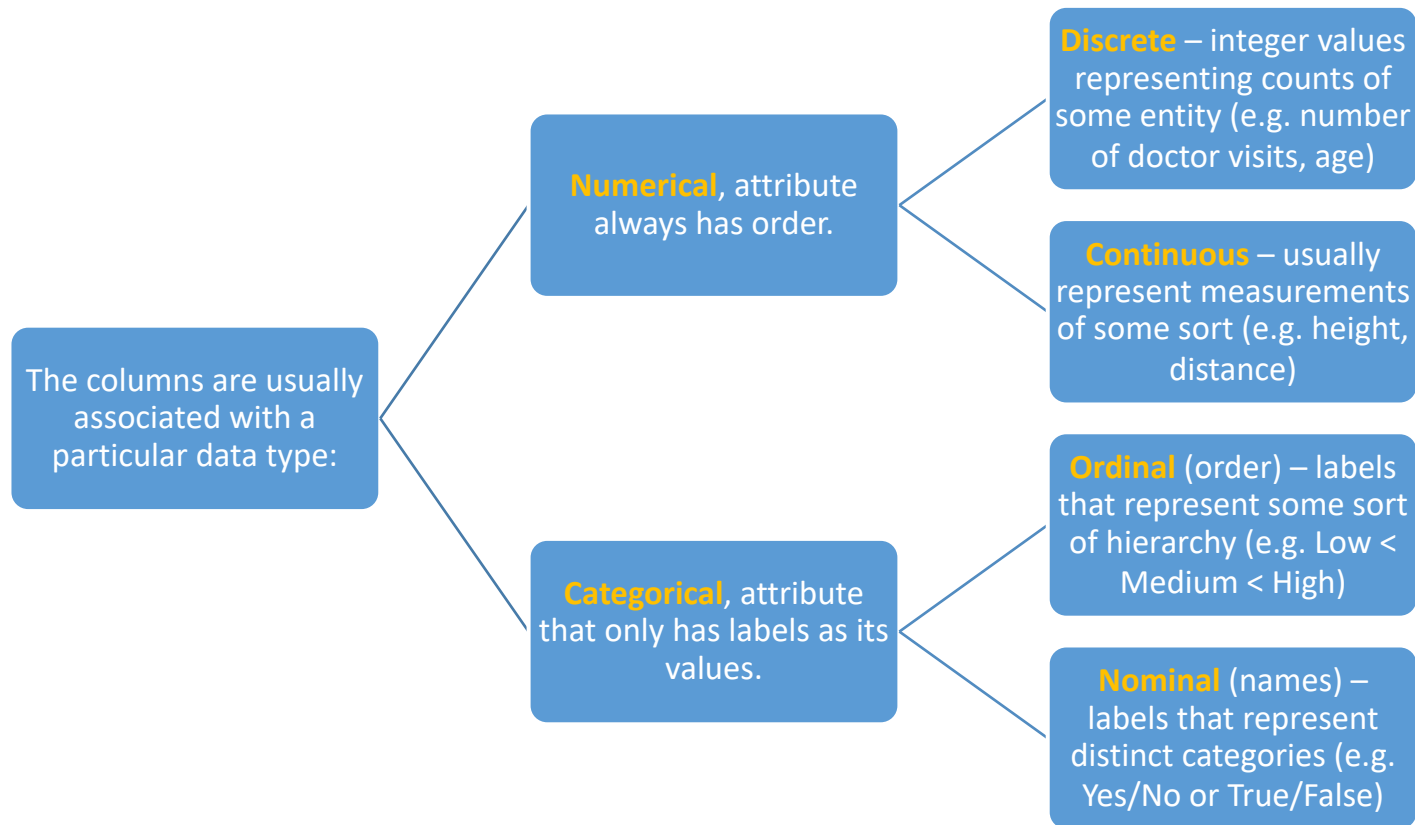
- In data science most of the data encountered is in tabular format, e.g. our tennis data set.
- We call this kind of data **structured data**.
- As opposed to **unstructured data** which usually appears in the form of text, e.g. medical reports, news articles.
- In this course we will take a look at both. We will start with structured data.

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

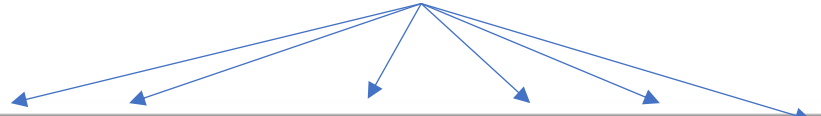
Structured Data

- Structured data consists of tables where
 - each column describes an attribute of the data objects in question. We often call the columns **variables** or **attributes**.
 - Each row describes a single **observation** or data object.
- For example, in our tennis data set each row describes a day in terms of its attributes (columns).

Data Types



Variables



Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

All of the variables in this data set are categorical variables – variables whose values only consist of labels/levels.

Outlook and PlayTennis are **nominal categorical** variables, the labels/level cannot be considered ordered, i.e. $Yes \not\leq No$ and $No \not\leq Yes$

The remaining variables are all ordinal categorical variables – the labels/levels can be considered ordered, i.e. $Cool < Mild < Hot$

ID	Gender	Age	Income	Rating
1	Male	28	\$50,000	4.5
2	Female	35	\$65,000	3.8
3	Male	22	\$40,000	4.2
4	Female	45	\$80,000	4.8
5	Male	31	\$55,000	3.5

- **ID**: Discrete numerical variable representing a unique identifier for each individual.
- **Gender**: Nominal categorical variable representing the gender of the individual (Male/Female).
- **Age**: Discrete numerical variable representing the age of the individual.
- **Income**: Continuous numerical variable representing the income of the individual.
- **Rating**: Continuous numerical variable representing a rating given by the individual.

Note: We see later that we will treat numerical ID variables like they appear in this table as nominal categorical variables because it makes no sense to use these identifiers as numerical values, we cannot order them or do mathematical transformations on them.

Real-World Data is Noisy

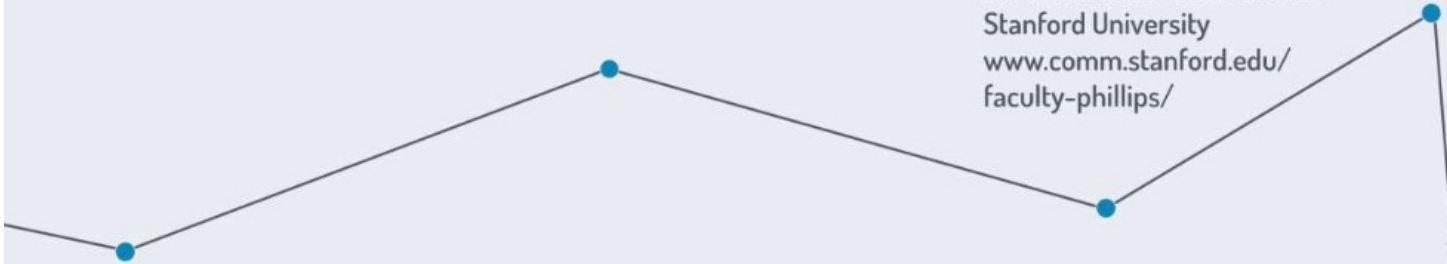
<https://www.slideshare.net/dataremixed/17-key-traits-of-data-literacy>

“ Anyone who has worked with data knows that it doesn't all come in pristine form. For this reason, a data literate person needs to learn how to handle data that needs some work, or that doesn't even exist in a data form and needs to be gathered. This is often missed, but it's one of the key points in becoming data literate.”



CHERYL PHILLIPS

Lorry I. Lokey Visiting Professor
in Professional Journalism at
Stanford University
[www.comm.stanford.edu/
faculty-phillips/](http://www.comm.stanford.edu/faculty-phillips/)



The Effects of Noisy Data

- A marketing firm was tasked to determine how Volkswagen ranks among the top brands
- They found a lot of misspellings of the name Volkswagen

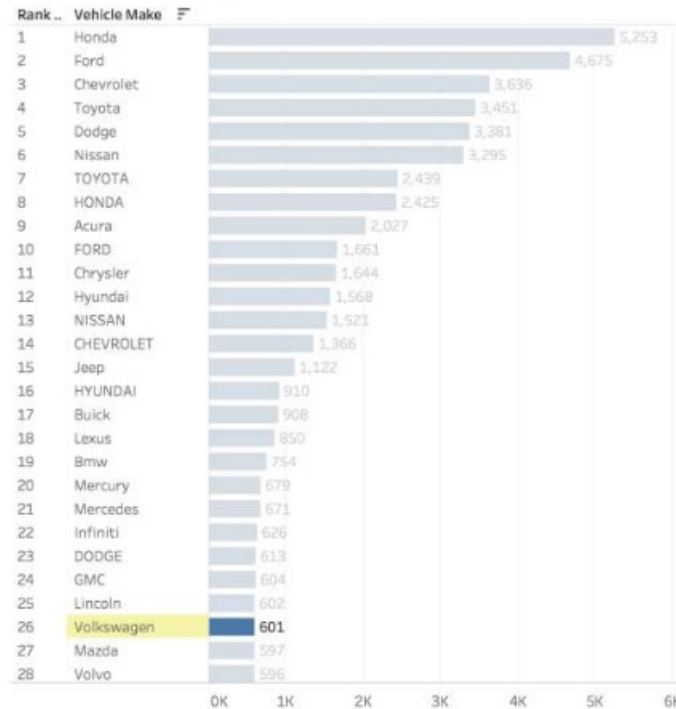


The Effects of Noisy Data

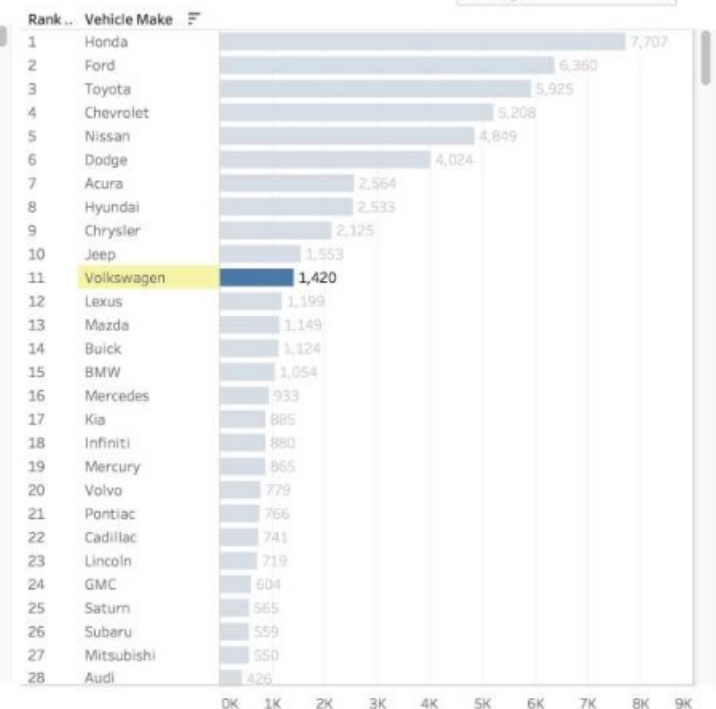


Noisy vs. Clean Data

Vehicle Makes - Original



Vehicle Makes - Cleaned



- With the noisy data the Volkswagen brand ranked at position #26.
- Once cleaned, that is, once the spelling error of the name were removed the brand moved to position #11.
- A huge difference!