# Chapter 1:
# Data and Datasets

# Data

The Latin word **data** is the plural of **datum**, "*(thing) given*"

**Data** are a set of values of *qualitative* or *quantitative* variables about one or more persons or objects, while a **datum** is a single value of a single variable.[1]

[1] "Statistical Language - What are Data?". *Australian Bureau of Statistics*. 2013-07-13

# Dataset types

- Tables
- Networks / graphs (trees)
- Text / logs
- Fields
- Static file vs. dynamic stream

# Tables

- each data **item** in a new row
- each column contains an **attribute**

| 32 | 7/16/07 | 2-High | Medium Box | 0.65 | 7/18/07 |

| A | B | C | S | T | U |
|---|---|---|---|---|---|
| Order ID | Order Date | Order Priority | Product Container | Product Base Margin | Ship Date |
| 32 | 7/16/07 | 2-High | Medium Box | 0.65 | 7/18/07 |

| Order ID | Order Date | Order Priority | Product Container | Product Base Margin | Ship Date |
|---|---|---|---|---|---|
| 3 | 10/14/06 | 5-Low | Large Box | 0.8 | 10/21/06 |
| 6 | 2/21/08 | 4-Not Specified | Small Pack | 0.55 | 2/22/08 |
| 32 | 7/16/07 | 2-High | Small Pack | 0.79 | 7/17/07 |
| 32 | 7/16/07 | 2-High | Jumbo Box | 0.72 | 7/17/07 |
| 32 | 7/16/07 | 2-High | Medium Box | 0.6 | 7/18/07 |
| 32 | 7/16/07 | 2-High | Medium Box | 0.65 | 7/18/07 |
| 35 | 10/23/07 | 4-Not Specified | Wrap Bag | 0.52 | 10/24/07 |
| 35 | 10/23/07 | 4-Not Specified | Small Box | 0.58 | 10/25/07 |
| 36 | 11/3/07 | 1-Urgent | Small Box | 0.55 | 11/3/07 |
| 65 | 3/18/07 | 1-Urgent | Small Pack | 0.49 | 3/19/07 |
| 66 | 1/20/05 | 5-Low | Wrap Bag | 0.56 | 1/20/05 |
| 69 | 6/4/05 | 4-Not Specified | Small Pack | 0.44 | 6/6/05 |
| 69 | 6/4/05 | 4-Not Specified | Wrap Bag | 0.6 | 6/6/05 |
| 70 | 12/18/06 | 5-Low | Small Box | 0.59 | 12/23/06 |
| 70 | 12/18/06 | 5-Low | Wrap Bag | 0.82 | 12/23/06 |
| 96 | 4/17/05 | 2-High | Small Box | 0.55 | 4/19/05 |
| 97 | 1/29/06 | 3-Medium | Small Box | 0.38 | 1/30/06 |
| 129 | 11/19/08 | 5-Low | Small Box | 0.37 | 11/28/08 |
| 130 | 5/8/08 | 2-High | Small Box | 0.37 | 5/9/08 |
| 130 | 5/8/08 | 2-High | Medium Box | 0.38 | 5/10/08 |
| 130 | 5/8/08 | 2-High | Small Box | 0.6 | 5/11/08 |
| 132 | 6/11/06 | 3-Medium | Medium Box | 0.6 | 6/12/06 |
| 132 | 6/11/06 | 3-Medium | Jumbo Box | 0.69 | 6/14/06 |
| 134 | 5/1/08 | 4-Not Specified | Large Box | 0.82 | 5/3/08 |
| 135 | 10/21/07 | 4-Not Specified | Small Pack | 0.64 | 10/23/07 |
| 166 | 9/12/07 | 2-High | Small Box | 0.55 | 9/14/07 |
| 193 | 8/8/06 | 1-Urgent | Medium Box | 0.57 | 8/10/06 |
| 194 | 4/5/08 | 3-Medium | Wrap Bag | 0.42 | 4/7/08 |

| Order ID | Order Date | Order Priority | Product Container | Product Base Margin | Ship Date |
|---|---|---|---|---|---|
| 3 | 10/14/06 | 5-Low | Large Box | 0.8 | 10/21/06 |
| 6 | 2/21/08 | 4-Not Specified | Small Pack | 0.55 | 2/22/08 |
| 32 | 7/16/07 | 2-High | Small Pack | 0.79 | 7/17/07 |
| 32 | 7/16/07 | 2-High | Jumbo Box | 0.72 | 7/17/07 |
| 32 | 7/16/07 | 2-High | Medium Box | 0.6 | 7/18/07 |
| 32 | 7/16/07 | 2-High | Medium Box | 0.65 | 7/18/07 |
| 35 | 10/23/07 | 4-Not Specified | Wrap Bag | 0.52 | 10/24/07 |
| 35 | 10/23/07 | 4-Not Specified | Small Box | 0.58 | 10/25/07 |
| 36 | 11/3/07 | 1-Urgent | Small Box | 0.55 | 11/3/07 |
| 65 | 3/18/07 | 1-Urgent | Small Pack | 0.49 | 3/19/07 |
| 66 | 1/20/05 | 5-Low | Wrap Bag | 0.56 | 1/20/05 |
| 69 | 6/4/05 | 4-Not Specified | Small Pack | 0.44 | 6/6/05 |
| 69 | 6/4/05 | 4-Not Specified | Wrap Bag | 0.6 | 6/6/05 |
| 70 | 12/18/06 | 5-Low | Small Box | 0.59 | 12/23/06 |
| 70 | 12/18/06 | 5-Low | Wrap Bag | 0.82 | 12/23/06 |
| 96 | 4/17/05 | 2-High | Small Box | 0.55 | 4/19/05 |
| 97 | 1/29/06 | 3-Medium | Small Box | 0.38 | 1/30/06 |
| 129 | 11/19/08 | 5-Low | Small Box | 0.37 | 11/28/08 |
| 130 | 5/8/08 | 2-High | Small Box | 0.37 | 5/9/08 |
| 130 | 5/8/08 | 2-High | Medium Box | 0.38 | 5/10/08 |
| 130 | 5/8/08 | 2-High | Small Box | 0.6 | 5/11/08 |
| 132 | 6/11/06 | 3-Medium | Medium Box | 0.6 | 6/12/06 |
| 132 | 6/11/06 | 3-Medium | Jumbo Box | 0.69 | 6/14/06 |
| 134 | 5/1/08 | 4-Not Specified | Large Box | 0.82 | 5/3/08 |
| 135 | 10/21/07 | 4-Not Specified | Small Pack | 0.64 | 10/23/07 |
| 166 | 9/12/07 | 2-High | Small Box | 0.55 | 9/14/07 |
| 193 | 8/8/06 | 1-Urgent | Medium Box | 0.57 | 8/10/06 |
| 194 | 4/5/08 | 3-Medium | Wrap Bag | 0.42 | 4/7/08 |

Record / Item

| Order ID | Order Date | Order Priority | Product Container | Product Base Margin | Ship Date |
|---|---|---|---|---|---|
| 3 | 10/14/06 | 5-Low | Large Box | 0.8 | 10/21/06 |
| 6 | 2/21/08 | 4-Not Specified | Small Pack | 0.55 | 2/22/08 |
| 32 | 7/16/07 | 2-High | Small Pack | 0.79 | 7/17/07 |
| 32 | 7/16/07 | 2-High | Jumbo Box | 0.72 | 7/17/07 |
| 32 | 7/16/07 | 2-High | Medium Box | 0.6 | 7/18/07 |
| 32 | 7/16/07 | 2-High | Medium Box | 0.65 | 7/18/07 |
| 35 | 10/23/07 | 4-Not Specified | Wrap Bag | 0.52 | 10/24/07 |
| 35 | 10/23/07 | 4-Not Specified | Small Box | 0.58 | 10/25/07 |
| 36 | 11/3/07 | 1-Urgent | Small Box | 0.55 | 11/3/07 |
| 65 | 3/18/07 | 1-Urgent | Small Pack | 0.49 | 3/19/07 |
| 66 | 1/20/05 | 5-Low | Wrap Bag | 0.56 | 1/20/05 |
| 69 | 6/4/05 | 4-Not Specified | Small Pack | 0.44 | 6/6/05 |
| 69 | 6/4/05 | 4-Not Specified | Wrap Bag | 0.6 | 6/6/05 |
| 70 | 12/18/06 | 5-Low | Small Box | 0.59 | 12/23/06 |
| 70 | 12/18/06 | 5-Low | Wrap Bag | 0.82 | 12/23/06 |
| 96 | 4/17/05 | 2-High | Small Box | 0.55 | 4/19/05 |
| 97 | 1/29/06 | 3-Medium | Small Box | 0.38 | 1/30/06 |
| 129 | 11/19/08 | 5-Low | Small Box | 0.37 | 11/28/08 |
| 130 | 5/8/08 | 2-High | Small Box | 0.37 | 5/9/08 |
| 130 | 5/8/08 | 2-High | Medium Box | 0.38 | 5/10/08 |
| 130 | 5/8/08 | 2-High | Small Box | 0.6 | 5/11/08 |
| 132 | 6/11/06 | 3-Medium | Medium Box | 0.6 | 6/12/06 |
| 132 | 6/11/06 | 3-Medium | Jumbo Box | 0.69 | 6/14/06 |
| 134 | 5/1/08 | 4-Not Specified | Large Box | 0.82 | 5/3/08 |
| 135 | 10/21/07 | 4-Not Specified | Small Pack | 0.64 | 10/23/07 |
| 166 | 9/12/07 | 2-High | Small Box | 0.55 | 9/14/07 |
| 193 | 8/8/06 | 1-Urgent | Medium Box | 0.57 | 8/10/06 |
| 194 | 4/5/08 | 3-Medium | Wrap Bag | 0.42 | 4/7/08 |

Field / Dimension / Attribute

| Order ID | Order Date | Order Priority | Product Container | Product Base Margin | Ship Date |
|---|---|---|---|---|---|
| 3 | 10/14/06 | 5-Low | Large Box | 0.8 | 10/21/06 |
| 6 | 2/21/08 | 4-Not Specified | Small Pack | 0.55 | 2/22/08 |
| 32 | 7/16/07 | 2-High | Small Pack | 0.79 | 7/17/07 |
| 32 | 7/16/07 | 2-High | Jumbo Box | 0.72 | 7/17/07 |
| 32 | 7/16/07 | 2-High | Medium Box | 0.6 | 7/18/07 |
| 32 | 7/16/07 | 2-High | Medium Box | 0.65 | 7/18/07 |
| 35 | 10/23/07 | 4-Not Specified | Wrap Bag | 0.52 | 10/24/07 |
| 35 | 10/23/07 | 4-Not Specified | Small Box | 0.58 | 10/25/07 |
| 36 | 11/3/07 | 1-Urgent | Small Box | 0.55 | 11/3/07 |
| 65 | 3/18/07 | 1-Urgent | Small Pack | 0.49 | 3/19/07 |
| 66 | 1/20/05 | 5-Low | Wrap Bag | 0.56 | 1/20/05 |
| 69 | 6/4/05 | 4-Not Specified | Small Pack | 0.44 | 6/6/05 |
| 69 | 6/4/05 | 4-Not Specified | | 0.6 | 6/6/05 |
| 70 | 12/18/06 | 5-Low | | 0.59 | 12/23/06 |
| 70 | 12/18/06 | 5-Low | | 0.82 | 12/23/06 |
| 96 | 4/17/05 | 2-High | | 0.55 | 4/19/05 |
| 97 | 1/29/06 | 3-Medium | | 0.38 | 1/30/06 |
| 129 | 11/19/08 | 5-Low | Small Box | 0.37 | 11/28/08 |
| 130 | 5/8/08 | 2-High | Small Box | 0.37 | 5/9/08 |
| 130 | 5/8/08 | 2-High | Medium Box | 0.38 | 5/10/08 |
| 130 | 5/8/08 | 2-High | Small Box | 0.6 | 5/11/08 |
| 132 | 6/11/06 | 3-Medium | Medium Box | 0.6 | 6/12/06 |
| 132 | 6/11/06 | 3-Medium | Jumbo Box | 0.69 | 6/14/06 |
| 134 | 5/1/08 | 4-Not Specified | Large Box | 0.82 | 5/3/08 |
| 135 | 10/21/07 | 4-Not Specified | Small Pack | 0.64 | 10/23/07 |
| 166 | 9/12/07 | 2-High | Small Box | 0.55 | 9/14/07 |
| 193 | 8/8/06 | 1-Urgent | Medium Box | 0.57 | 8/10/06 |
| 194 | 4/5/08 | 3-Medium | Wrap Bag | 0.42 | 4/7/08 |

1 = Quantitative
2 = Nominal
3 = Ordinal

| Order ID | Order Date | Order Priority | Product Container | Product Base Margin | Ship Date |
|---|---|---|---|---|---|
| 3 | 10/14/06 | 5-Low | Large Box | 0.8 | 10/21/06 |
| 6 | 2/21/08 | 4-Not Specified | Small Pack | 0.55 | 2/22/08 |
| 32 | 7/16/07 | 2-High | Small Pack | 0.79 | 7/17/07 |
| 32 | 7/16/07 | 2-High | Jumbo Box | 0.72 | 7/17/07 |
| 32 | 7/16/07 | 2-High | Medium Box | 0.6 | 7/18/07 |
| 32 | 7/16/07 | 2-High | Medium Box | 0.65 | 7/18/07 |
| 35 | 10/23/07 | 4-Not Specified | Wrap Bag | 0.52 | 10/24/07 |
| 35 | 10/23/07 | 4-Not Specified | Small Box | 0.58 | 10/25/07 |
| 36 | 11/3/07 | 1-Urgent | Small Box | 0.55 | 11/3/07 |
| 65 | 3/18/07 | 1-Urgent | Small Pack | 0.49 | 3/19/07 |
| 66 | 1/20/05 | 5-Low | Wrap Bag | 0.56 | 1/20/05 |
| 69 | 6/4/05 | 4-Not Specified | Small Pack | 0.44 | 6/6/05 |
| 69 | 6/4/05 | 4-Not Speci | | 0.6 | 6/6/05 |
| 70 | 12/18/06 | 5-Low | | 0.59 | 12/23/06 |
| 70 | 12/18/06 | 5-Low | | 0.82 | 12/23/06 |
| 96 | 4/17/05 | 2-High | | 0.55 | 4/19/05 |
| 97 | 1/29/06 | 3-Medium | | 0.38 | 1/30/06 |
| 129 | 11/19/08 | 5-Low | Small Box | 0.37 | 11/28/08 |
| 130 | 5/8/08 | 2-High | Small Box | 0.37 | 5/9/08 |
| 130 | 5/8/08 | 2-High | Medium Box | 0.38 | 5/10/08 |
| 130 | 5/8/08 | 2-High | Small Box | 0.6 | 5/11/08 |
| 132 | 6/11/06 | 3-Medium | Medium Box | 0.6 | 6/12/06 |
| 132 | 6/11/06 | 3-Medium | Jumbo Box | 0.69 | 6/14/06 |
| 134 | 5/1/08 | 4-Not Specified | Large Box | 0.82 | 5/3/08 |
| 135 | 10/21/07 | 4-Not Specified | Small Pack | 0.64 | 10/23/07 |
| 166 | 9/12/07 | 2-High | Small Box | 0.55 | 9/14/07 |
| 193 | 8/8/06 | 1-Urgent | Medium Box | 0.57 | 8/10/06 |
| 194 | 4/5/08 | 3-Medium | Wrap Bag | 0.42 | 4/7/08 |
| 194 | 4/5/08 | 3-Medium | Wrap Bag | 0.84 | 4/7/08 |

1 = Quantitative
2 = Nominal
3 = Ordinal

# Dataset for this Tutorial

- **Happiness** X **World Factbook** X **HDI** X **IHDI**

- built for this tutorial, data from 2017

- Download Link: https://github.com/Christoph/Happiness-Teaching-Dataset/archive/master.zip

- Data-Description: https://github.com/Christoph/Happiness-Teaching-Dataset

# Some Dataset Values Explained

- ## Happiness Score
  Calculated by asking sampled people in a country the question: *"How would you rate your happiness on a scale of 0 to 10 where 10 is the happiest."*

- ## Biggest Official Language
  Official language with most speakers in a country

- ## Human Development Index
  A composite index measuring average achievement in three basic dimensions of human development—a long and healthy life, knowledge and a decent standard of living.[2]

- ## Gini Coefficient
  This index measures the degree of inequality in the distribution of income in a country. Measured from 0 to 100 where 0 means perfect equality and 100 means perfect inequality.

[2] See *Technical note 1 at* http://hdr.undp.org/sites/default/files/hdr2018_technical_notes.pdf for details on how the HDI is calculated