# Programación para la Computación Científica - IA

# Data processing Cycle II
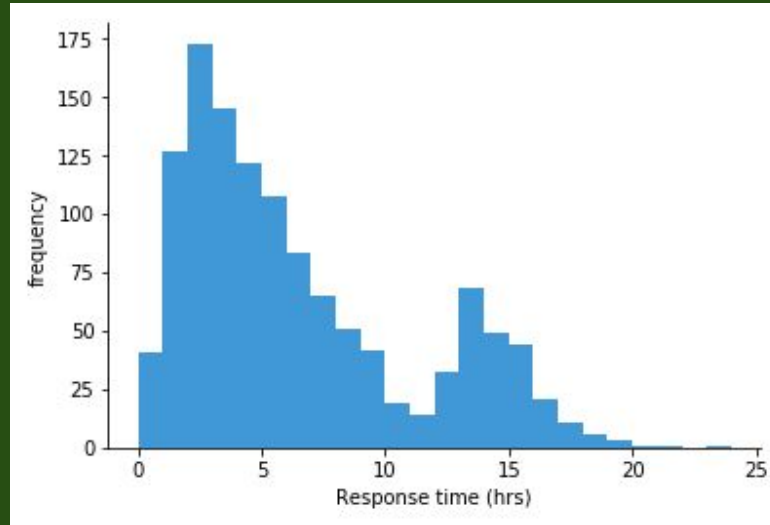
Universidad Sergio Arboleda
*Prof. John Corredor*
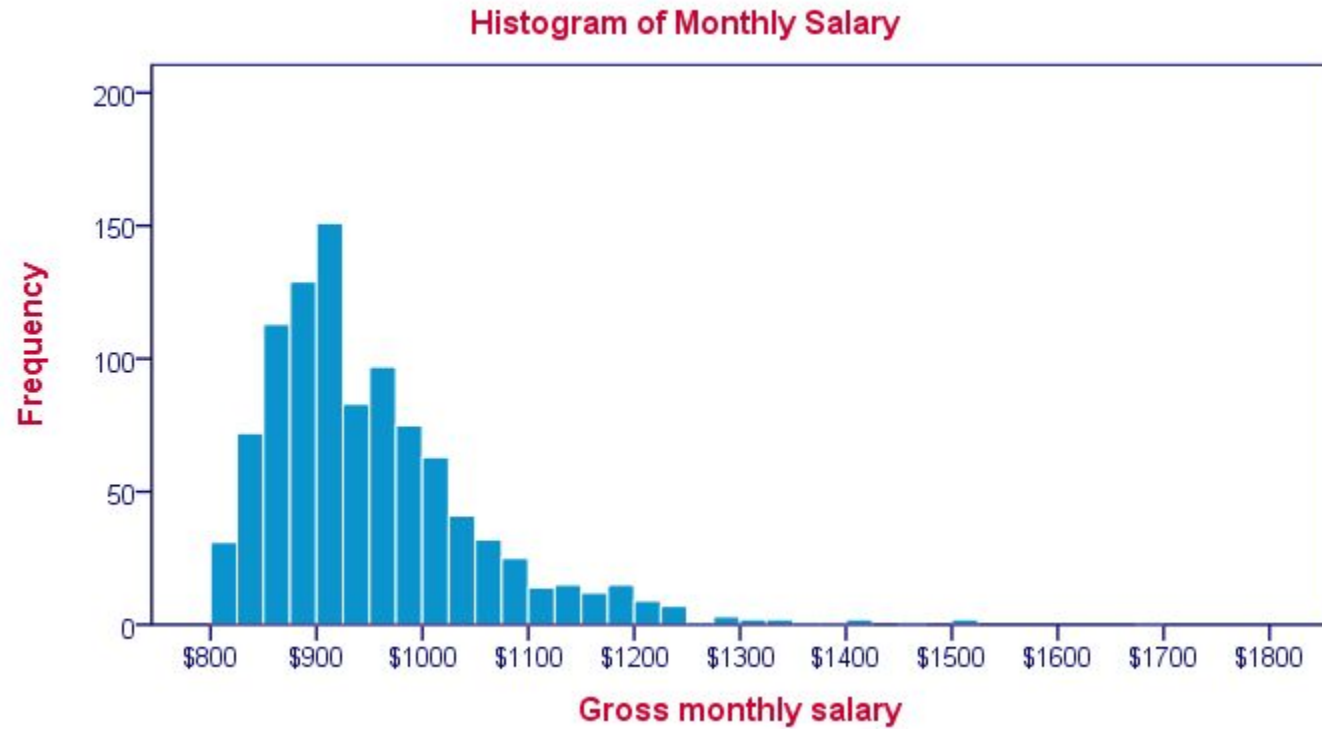
**Today goal's**

- **Data processing - Histograms**
- **Introduction to Jupyter Notebook**
- **Principles of Pyhon**
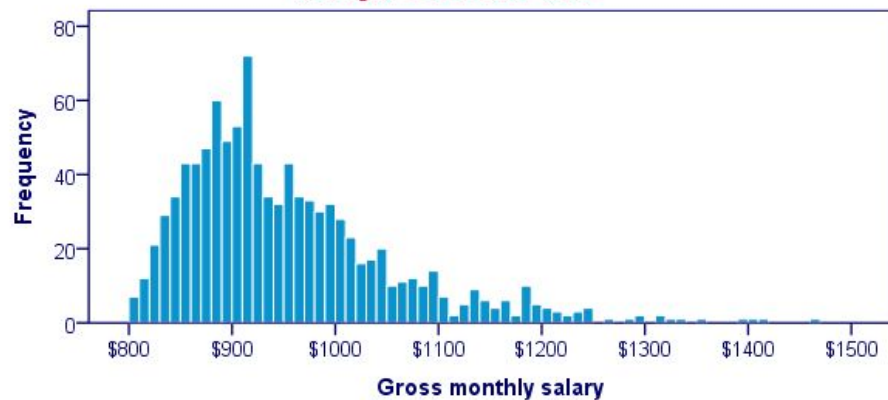
# What is a histogram?



A histogram is a chart that plots the distribution of a numeric variable's values as a series of bars. Each bar typically covers a range of numeric values called a bin or class; a bar's height indicates the frequency of data points with a value within the corresponding bin.

# What is a histogram?



Histogram of Monthly Salary

**Histogram Bin Width = $10,-**

Frequency / Gross monthly salary

**Histogram Bin Width = $25,-**

Frequency / Gross monthly salary

**Histogram Bin Width = $50,-**

Frequency / Gross monthly salary

**Histogram Bin Width = $100,-**

Frequency / Gross monthly salary

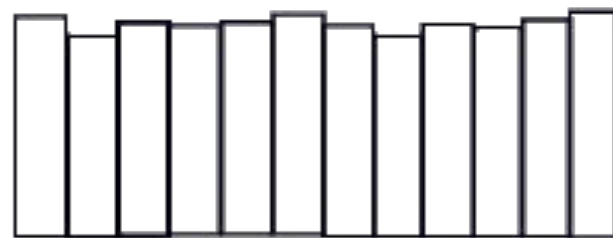| Individual Height, Measured in Inches | | | | |
|---|---|---|---|---|
| 69.9 | 68.9 | 68.2 | 66.0 | 71.9 |
| 69.0 | 70.0 | 68.5 | 66.5 | |
| 69.6 | 69.5 | 70.0 | 67.5 | |
| 68.5 | 70.4 | 66.8 | 68.3 | |
| 65.0 | 71.1 | 69.0 | 68.2 | |
| 65.9 | 71.0 | 69.3 | 69.1 | |
| 67.2 | 72.5 | 69.1 | 70.2 | |
| 67.5 | 73.1 | 69.4 | 69.5 | |
| 68.0 | 68.8 | 68.5 | 70.5 | |
| 68.6 | 71.3 | 65.5 | 70.8 | |



**Absolute Frequencies Histogram of Height**

# Bi-Modal Distribution

# Unitary Distribution
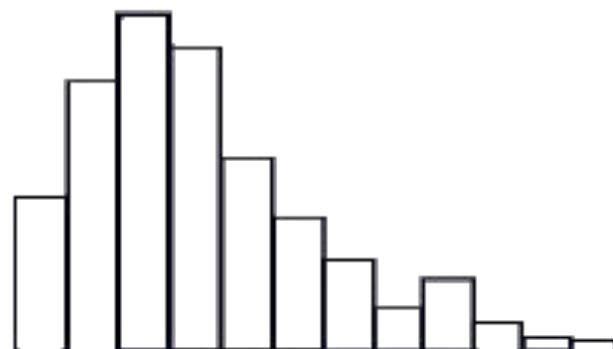
# Negatively Skewed

# - Positively Skewed

# How to start?

1. Count the number of data points (50 in our height example) .
2. Determine the range of the sample - the difference between the highest and lowest values (73.1-65, or 8.1 inches in our height example.
3. Determine the number of class intervals. (Measurement System Analysis (MSA)) You can use either of two methods as general guidelines in determining the number of intervals:

   A. Use ten intervals as a rule of thumb.

   B. Calculate the square root of the number of data points and round to the nearest whole number. In the case of our height example, the square root of 50 is 7.07, or 7 when rounded.

   You may wish to experiment with different interval numbers. If there are too many, the distribution will spread out, and the histogram will look flat. Likewise, if there are too few intervals, the distribution can look artificially tight.

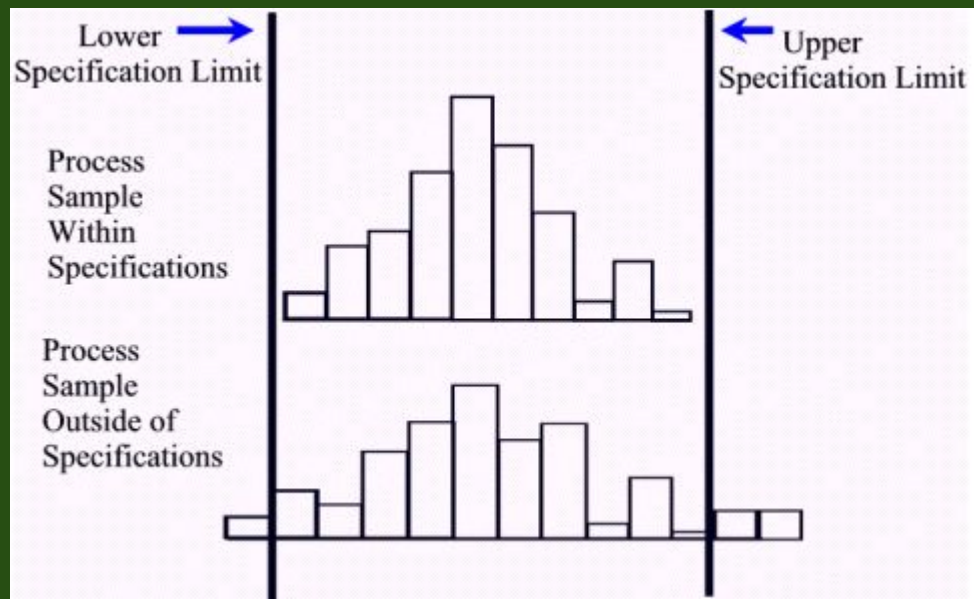4. Determine the interval class width by one of two methods:

...0 = 0.81

...ree. In this case, the height data has a Standard

**Individual Height, Measured in Inches**

| | | | | |
|---|---|---|---|---|
| 69.9 | 68.9 | 68.2 | 66.0 | 71.0 |
| 69.0 | 70.0 | 68.5 | 66.5 | 72.5 |
| 69.6 | 69.5 | 70.0 | 67.5 | 73.0 |
| 68.5 | 70.4 | 66.8 | 68.3 | 69.0 |
| 65.0 | 71.1 | 69.0 | 68.2 | 71.3 |
| 65.9 | 71.0 | 69.3 | 69.1 | 68.2 |
| 67.2 | 72.5 | 69.1 | 70.2 | 68.5 |
| 67.5 | 73.1 | 69.4 | 69.5 | 70.0 |
| 68.0 | 68.8 | 68.5 | 70.5 | 67.0 |
| 68.6 | 71.3 | 65.5 | 70.8 | 69.2 |

| Class | Height Intervals | Frequency | Total |
|---|---|---|---|
| 1 | 64.4 - 65.0 | X | 1 |
| 2 | 65.1 - 65.7 | X | 1 |
| 3 | 65.8 - 66.4 | XX | 2 |
| 4 | 66.5 - 67.1 | XX | 2 |
| 5 | 67.2 - 67.8 | XXXX | 4 |
| 6 | 67.9 - 68.5 | X | 1 |
| 7 | 68.6 - 69.2 | XXXXXXXXXX | 10 |
| 8 | 69.3 - 69.9 | XXXXXXXXX | 9 |
| 9 | 70.0 - 70.6 | XXXXXXX | 7 |
| 10 | 70.7 - 71.3 | XXX | 3 |
| 11 | 71.4 - 72.0 | XXXXXX | 6 |
| 12 | 72.1 - 72.7 | | 0 |
| 13 | 72.8 - 73.4 | XX | 2 |
| 14 | 73.5 - 74.1 | XX | 2 |

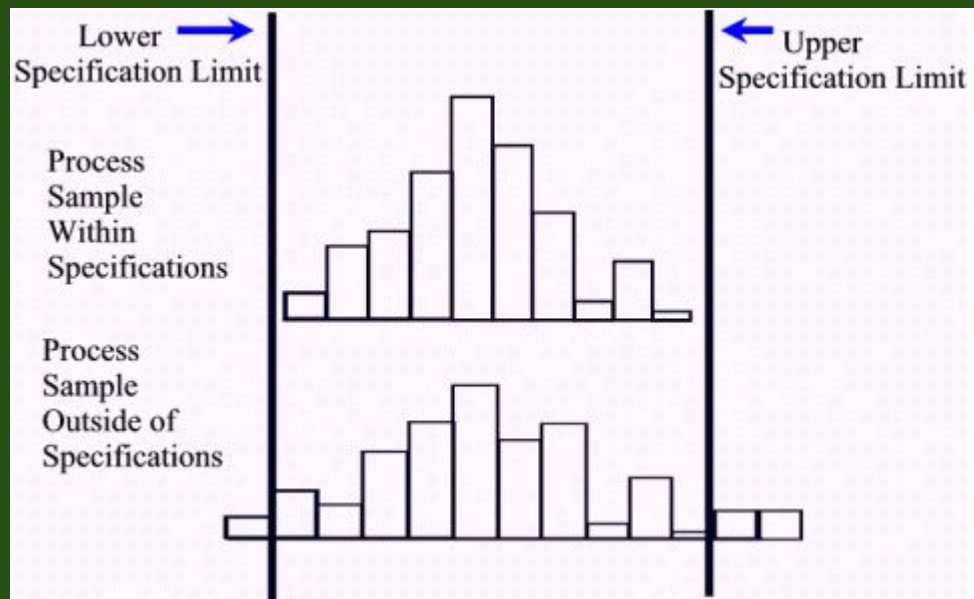Once the histogram is developed, you can analyze the data with regard to customer expectations (specifications).

The first histogram of a process sample falls within the specifications, while the second has a portion of the histogram outside of the specifications.
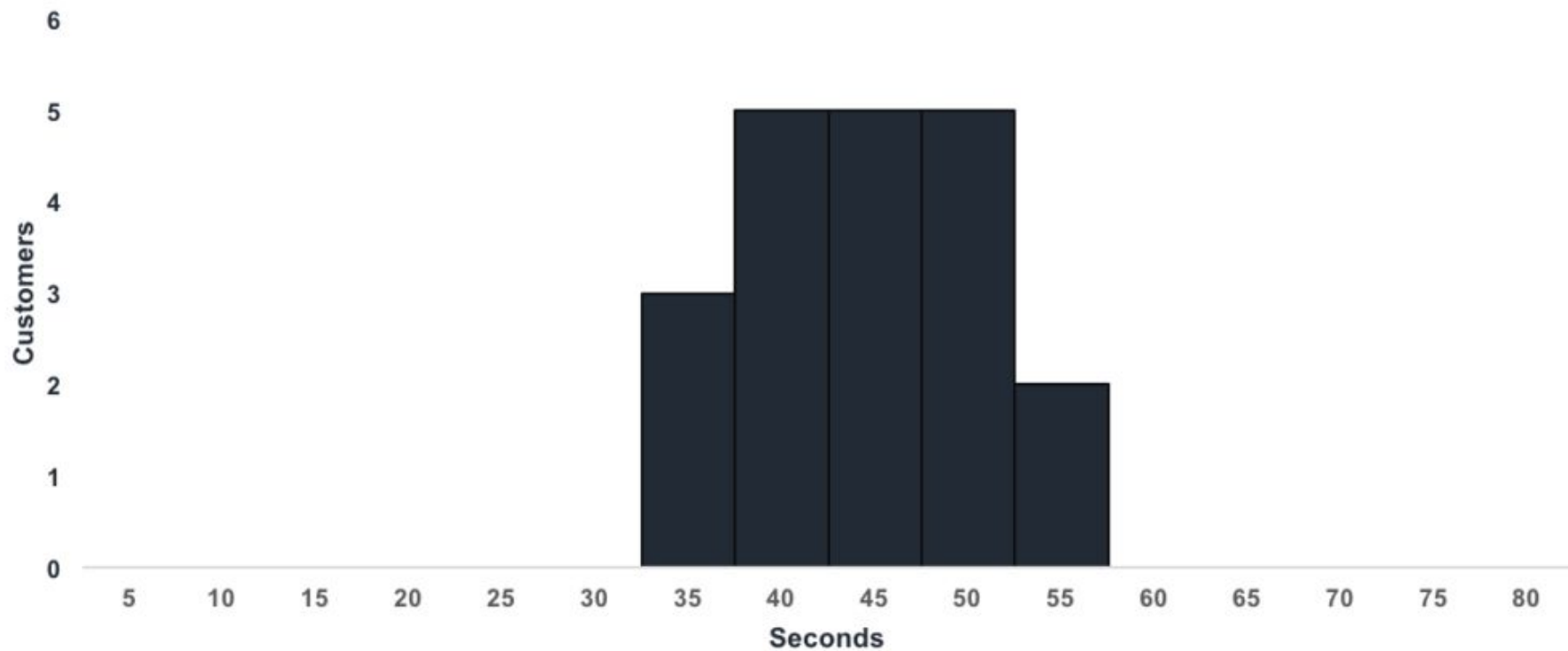
The second histogram has too much dispersion, or variability, to meet customer expectations. The indication is that action must be taken to make the output more consistent, or some number of defects will be produced. **(Statistical Process Control)**

The first histogram of a process sample falls within the specifications, while the second has a portion of the histogram outside of the specifications.

The second histogram has too much dispersion, or variability, to meet customer expectations. The indication is that action must be taken to make the output more consistent, or some number of defects will be produced. (Statistical Process Control)

Customer Wait Time

| Customer Waiting Time (in mins) |
| --- |
| 2.30 |
| 5.00 |
| 3.55 |
| 2.50 |
| 5.10 |
| 4.21 |
| 3.33 |
| 4.10 |
| 2.55 |
| 5.07 |
| 3.45 |
| 4.10 |
| 5.12 |

**Challenge 01**

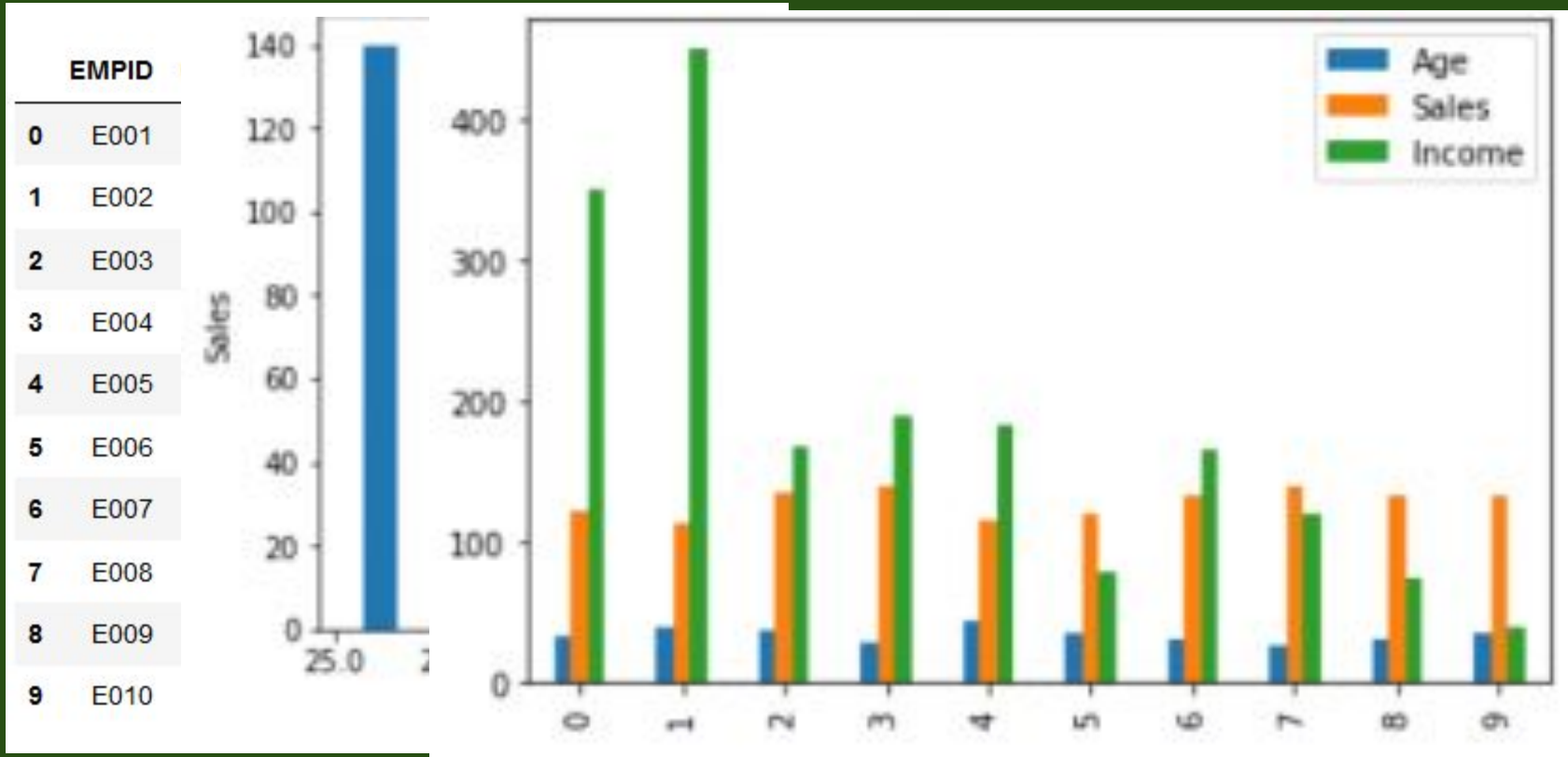| Sr No | Height (in cms) |
|-------|-----------------|
| 1     | 141             |
| 2     | 143             |
| 3     | 145             |
| 4     | 145             |
| 5     | 147             |
| 6     | 152             |
| 7     | 143             |
| 8     | 144             |
| 9     | 149             |
| 10    | 141             |
| 11    | 138             |
| 12    | 143             |
| 13    | 145             |
| 14    | 148             |
| 15    | 145             |

**Challenge 02**

**El Sr. Larry, un médico famoso, está realizando una investigación sobre la altura de los estudiantes que estudian en el octavo estándar. Ha reunido una muestra de 15 estudiantes, pero quiere saber cuál es la categoría máxima a la que pertenecen.**

# Data visualization with different Charts

# Data visualization with different Charts

# Data visualization with different Charts

# Data visualization with different Charts

| | EMPID | Gender | Age | Sales | BMI | Income |
|---|---|---|---|---|---|---|
| 0 | E001 | M | 34 | 123 | Normal | 350 |
| 1 | E002 | F | 40 | 114 | Overweight | 450 |
| 2 | E003 | F | 37 | 135 | Obesity | 169 |
| 3 | E004 | M | 30 | 139 | Underweight | 189 |
| 4 | E005 | F | 44 | 117 | Underweight | 183 |
| 5 | E006 | M | 36 | 121 | Normal | 80 |
| 6 | E007 | M | 32 | 133 | Obesity | 166 |
| 7 | E008 | F | 26 | 140 | Normal | 120 |
| 8 | E009 | M | 32 | 133 | Normal | 75 |
| 9 | E010 | M | 36 | 133 | Underweight | 40 |

# Why should I prepare my data?

- Garbage in, garbage out
- Reduce errors
- Remove duplicate records
- Fix missing values
- Correct range values
- Fix formatting (i.e. date, text, number)

# Experience Check

- How many people have experience with Python?
- What types of data formats do you use in your organizations?
    - CSV, Excel, PDF, JSON, XML, SQL databases, etc
- What types of tools do you use?
    - Hoja de Calculo, ACL, IDEA, SQL Server, Python, R, SAS, Cognos, etc

# What types of data formats might I encounter?

- Comma Separated Value (CSV)
- Hoja de calculo
- JavaScript Object Notation (JSON)
- Structured Query Language (SQL)
- And more!

# CVS example

- SFO Airport Survey Results

```
RESPNUM,CCGID,RUNID,INTDATE,AIRLINE,FLIGHT,DESTINATION,DESTGEO,DESTMARK,GATE,BAREA,STRATA,PEAK,DEPTIME,ARRTIME,HOWLONG,METHOD,Q2PURP1,Q2PURP2,Q2PURP3,Q3GETTO1,Q3GETTO2,Q3GETTO3,Q3PARK,Q4BAGS,Q4STORE,Q4FOOD,Q4WIFI,Q5TIMESFLOWN,Q5FIRSTTIME,Q6LONGUSE,SAQ,Q7ART,Q7FOOD,Q7STORE,Q7SIGN,Q7WALKWAYS,Q7SCREENS,Q7INFODOWN,Q7INFOUP,Q7WIFI,Q7ROADS,Q7PARK,Q7AIRTRAIN,Q7LTPARKING,Q7RENTAL,Q7ALL,Q8COM1,Q8COM2,Q8COM3,Q9BOARDING,Q9AIRTRAIN,Q9RENTAL,Q9FOOD,Q9RESTROOM,Q9ALL,Q9COM1,Q9COM2,Q9COM3,Q10SAFE,Q10COM1,Q10COM2,Q10COM3,Q11TSAPRE,Q12PRECHEKCRATE,Q12COM1,Q12COM2,Q12COM3,Q13COUNTY,Q13GETRATE,Q14FIND,Q14FASSTHRU,Q15PROBLEM,Q15COM1,Q15COM2,Q15COM3,Q16LIVE,HOME,Q17CITY,Q17STATE,Q17ZIP,Q17COUNTRY,Q18AGE,Q19GENDER,Q20INCOME,Q21FLY,Q22SJC,Q22OAK,LANG,WEIGHT
```

(Remaining rows are CSV survey data — dense numeric/text records such as:)

```
3,460,15076,13,8,242,18,3,4,58,D,1,2,7:20 AM,6:00 AM,80,1,2,,,2,,,1,2,1,5,1,2,3,1,4,1,5,5,5,4,3,3,3,6,4,5,3,6,4,5,3,6,4,4,,,5,6,6,4,4,,,5,8,,,1,5,2,,,1,3,5,5,2,,,,1,3,OAKLAND,CA,94619,US,6,0,1,2,2,2,1,0.720538
...
56,521,15079,12,8,5988,50,1,4,58,D,1,1,10:53 AM,9:30 AM,83,1,3,,,2,,,2,3,1,2,2,4,4,4,4,4,5,5,6,5,6,4,5,,,4,5,0,4,4,4,,,4,,1,5,,,5,4,4,5,1,3,,,1,1,SAN FRANCISC,CA,,US,7,2,0,2,2,2,1,0.720538
```

# Hoja de Calculo

- SFO Airport Survey Results

# JSON Examples

- Trip Advisor JSON File      * Yelp JSON file

{"Reviews": [{"Ratings": {"Service": "4", "Cleanliness": "5", "Overall": "5.0", "Value": "4", "Sleep Quality": "4", "Rooms": "5", "Location": "5"}, "AuthorLocation": "Boston", "Title": "\u201cExcellent Hotel & Location\u201d", "Author": "gowharr32", "ReviewID": "UR126946257", "Content": "We enjoyed the Best Western Pioneer Square. My husband and I had a room with a king bed and it was clean, quiet, and attractive. Our sons were in a room with twin beds. Their room was in the corner on the main street and they said it was a little noisier and the neon light shone in. But later hotels on the trip made them appreciate this one more. We loved the old wood center staircase. Breakfast was included and everyone was happy with waffles, toast, cereal, and an egg meal. Location was great. We could walk to shops and restaurants as well as transportation. Pike Market was a reasonable walk. We enjoyed the nearby Gold Rush Museum. Very, very happy with our stay. Staff was helpful and knowledgeable.", "Date": "March 29, 2012"}, {"Ratings": {"Overall": "5.0"}, "AuthorLocation": "Madison, Wisconsin", "Title": "\u201cGreat Visit to Seattle!\u201d", "Author": "Nancy W", "ReviewID": "UR126795011", "Content": "Great visit to Seattle thanks to our stay at the Best Western Pioneer Square! The hotel was reasonably priced and close to everything we wanted to see - ferry ride, Underground Tour, Klondike Museum, short walk to Pike Market and other shopping. The staff was amazingly helpful and accommodating. Our room was very clean and had everything we needed. Breakfast was plentiful and very good. Before we booked, I read about some potential issues with the area. I can honestly say that the area was just fine! In fact, if you enjoy historic and quaint parts of town, this is definitely where you want to stay. I will be recommending this hotel to anyone who is headed to Seattle.", "Date": "March 27, 2012"}, {"Ratings": {"Service": "5", "Cleanliness": "5", "Overall": "5.0", "Value": "5", "Sleep Quality": "5", "Rooms": "5", "Location": "5"}, "AuthorLocation": "Ketchikan, Alaska", "Title": "\u201cExcellent in Everyway\u201d", "Author": "Janet H",

{
    "business_id": "PK6aSizckHFWk8i0oxt5DA",
    "full_address": "400 Waterfront Dr E\nHomestead\nHomestead, PA 15120",
    "hours": {},
    "open": true,
    "categories": [
        "Burgers",
        "Fast Food",
        "Restaurants"
    ],
    "city": "Homestead",
    "review_count": 5,
    "name": "McDonald's",
    "neighborhoods": [
        "Homestead"
    ],
    "longitude": -79.910032,
    "state": "PA",
    "stars": 2,
    "latitude": 40.412086,
    "attributes": {
        "Take-out": true,
        "Wi-Fi": "free",
        "Drive-Thru": true,
        "Good For": {
            "dessert": false,
            "latenight": false,
            "lunch": false,
            "dinner": false,
            "breakfast": false,
            "brunch": false
        },
        "Caters": false,
        "Noise Level": "average",
        "Takes Reservations": false,
        "Delivery": false

# SQL Examples

- Sanoke customer data

# Python

- Object-oriented, high-level programming  language
- Used as a scripting or glue language to connect existing components together
- Simple, easy to learn syntax emphasizes readability
- Supports modules and packages
- Python interpreter and the extensive standard library are FREE!

# Python

Key Python Package:
- Pandas
    - Open source library that allows you to work with CSV, Excel, JSON, and SQL database files, pull them into tables (called dataframes), and perform various data analysis techniques.
-

# References

★ **Kernighan, Brian W., and Dennis M. Ritchie. The C Programming Language. Vol. 2. Englewood Cliffs: prentice-Hall, 1988.**

★ **Silberschatz, Abraham, Peter B. Galvin, and Greg Gagne. Operating System Concepts. Vol. 8. Wiley, 2013.**

★ https://planningtank.com/computer-applications/data-processing-cycle

★ https://www.talend.com/resources/what-is-data-processing/

★ https://peda.net/kenya/ass/subjects2/computer-studies/form-3/data-processing/dpc2

★ https://www.academia.edu/38210518/What_is_Data_Processing

★ https://www.studocu.com/en/document/polytechnic-university-of-the-philippines/information-and-communication-technology/lecture-notes/data-processing-lectures-in-data-processing/3167716/view

★ http://download.nos.org/srsec330/330L2.pdf