# COMP 4462
# Data Visualization Tutorial

CHEN Chang
PAN Ziqi

Friday 27 September, 2024

# Visualization process

**Prepare data**

- Get data
  - Download, crawl, collect
- Load data
  - Load data into visualization software
- Transform, join and aggregate
  - Make data to the form that ready to be drawn
- Filter
  - Clean up data and remove irrelevant information

**Draw data (visualize)**

- Visual encoding design
  - It's what you learn from the lectures
  - Marks and channels
  - Position, color, size, shape, etc.
- Interactions
  - Pan and zoom, select and filter
  - Click, drag and drop, scroll, and keyboard input, etc.

# Get data

- Download data prepared by the others
  - Kaggle Dataset
  - World Bank
  - And many more
- Crawl from the web
  - Write you own program to crawl
  - From API
  - Extract from HTML or JSON
  - Python: Scrapy, Beautiful Soup
  - Nodejs: Cheerio
- Collect by yourself
  - Most costly
  - Takes time and efforts
  - But sometimes you have to when there is no existing dataset

# Load data

- Most common
  - csv, tsv: comma separated value, tab separated value
  - xlsx: Excel
  - HTTP request: Ajax, JSON, XML
- Databases
  - SQL: Oracle, MySQL, PostgreSQL, MS SQL, etc.
    - Structured, normalized
  - NoSQL: MongoDB
    - Document based
- PDF
  - Tableau supports import from PDF.
  - Import PDF to Excel: reference1, reference2.
- Other source
  - Google Cloud Public Datasets
    - Only available through Google Cloud
    - Too big to be downloaded

# Clean data

- Data is always dirty
  - Missing values
  - Typo
  - Overloaded fields (mixing continuous numbers with text)
  - Mismatch primary keys / external keys
  - Duplicated entries
  - Missing data for several days
    - Equipment failure / bugs in crawler programs / website is down
  - Non-sense error in data, e.g. integer overflow, or just not making any sense
  - Emoji / language / accent decoration
  - Identical typeface but different in unicode
- Depends on severity, it can be very nasty to deal with
- Data normalization
  - [Google Text Normalization Challenge](#)

# Transform, join and aggregate

- Manipulate data to the form for visualization
  - Wide form
  - Long form
  - Derive attributes: percentage changes, year-to-year changes
- Join
  - Linking up multiple table or data sources
  - Inner join, left join, right join, outer join
  - Commonly join on ID
    - Sometimes on date
    - Sometimes on multiple attributes
- Aggregate
  - Statistical: counting, sum, average, median, etc.
  - Grouping: binning, frequency, time slicing
  - Moving average, running sum

| Ranking | 2018 | 2017 | 2016 | 2015 |
|---------|------|------|------|------|
| CS      | 14   | 19   | 14   | 8    |
| CHEM    | 23   | 27   | 28   | 25   |

| Subject | Ranking | Year |
|---------|---------|------|
| CS      | 14      | 2018 |
| CS      | 19      | 2017 |
| CS      | 14      | 2016 |
| CS      | 8       | 2015 |
| CHEM    | 23      | 2018 |
| CHEM    | 27      | 2017 |
| CHEM    | 28      | 2016 |
| CHEM    | 25      | 2015 |

# Filter

- Reduce the number of items to show
- Focus only on relevant data, clean up irrelevant data
  - Base on user interest
  - Or users' level of authority
    - Not everyone can access all the data
  - Time relevancy
    - Outdated data are no longer relevant to real-time analysis
  - Geographic relevance
    - You don't care about restaurants outside Hong Kong (unless you're going to travel)
- Hard to show all with a limited screen size
  - Especially on mobile device
  - Reduce cluttering, more "clickable" on screen to show item details
- Zoom in to a specific small subset of data
  - Then you can show more detail of each item
  - Google Maps, zoom in to show more detailed terrain

# Tableau

Visualization with Tableau and data processing pipeline

- Install Tableau beforehand
  - Tableau student (Full version, preferred): https://www.tableau.com/academic/students
  - Or Tableau Public: https://public.tableau.com

# Tableau

- Tableau Public
  - Free
  - All saved works are public
    - Publicly viewable, downloadable
  - Must connect to the internet in order to save
  - Less data connectors
- Tableau Desktop
  - Free for students, need verification
  - Can save locally, use without connecting to the internet
  - More data connectors
- Tableau Prep
  - Prepare data for visualization
- Tableau Server
  - Standalone, dedicated server
  - Enterprise level, expensive

# Load Data

Take the dataset (global_superstore_2016.xlsx) from [GitHub](GitHub) as an example.

For Desktop

# Load Data

Take the dataset (global_superstore_2016.xlsx) from [GitHub](#) as an example.

For Web Authoring

# Load Data



Drag "Orders"

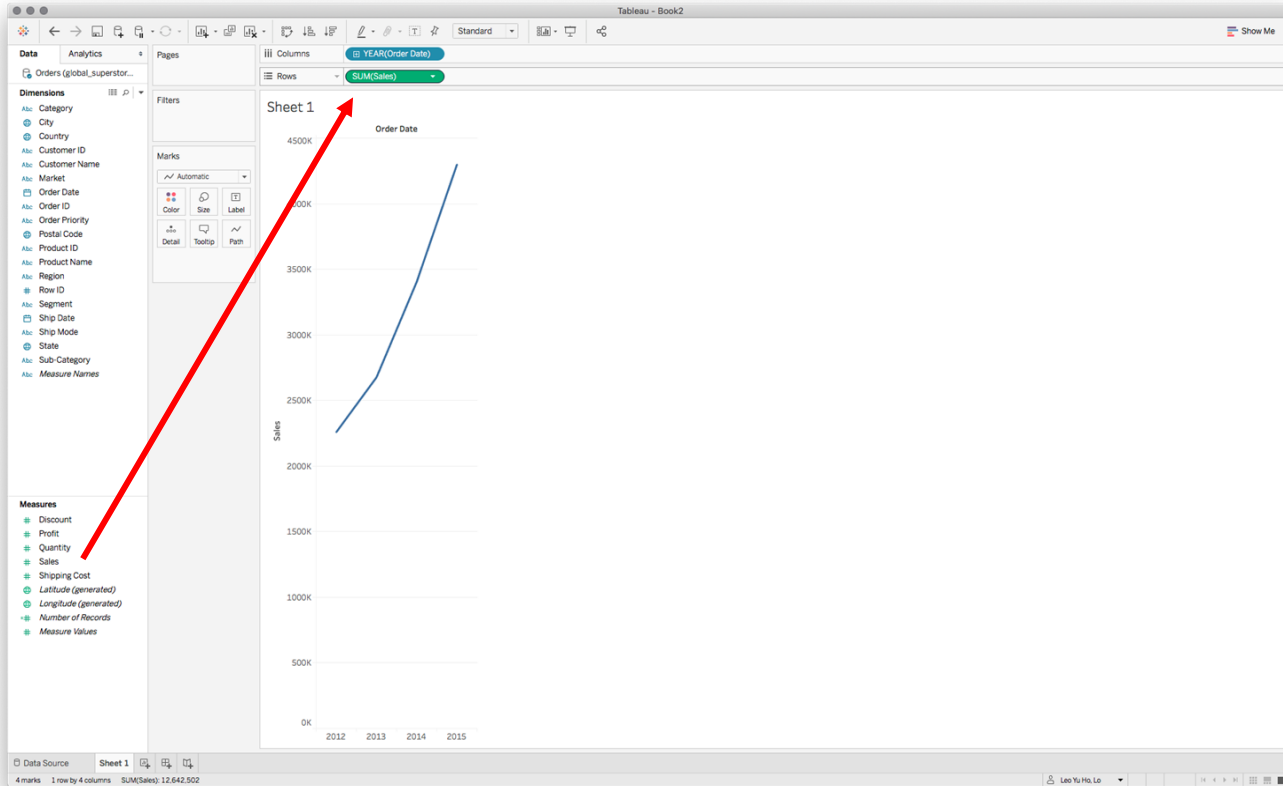If you page is like this, click **update now** to reach the status shown on the left.
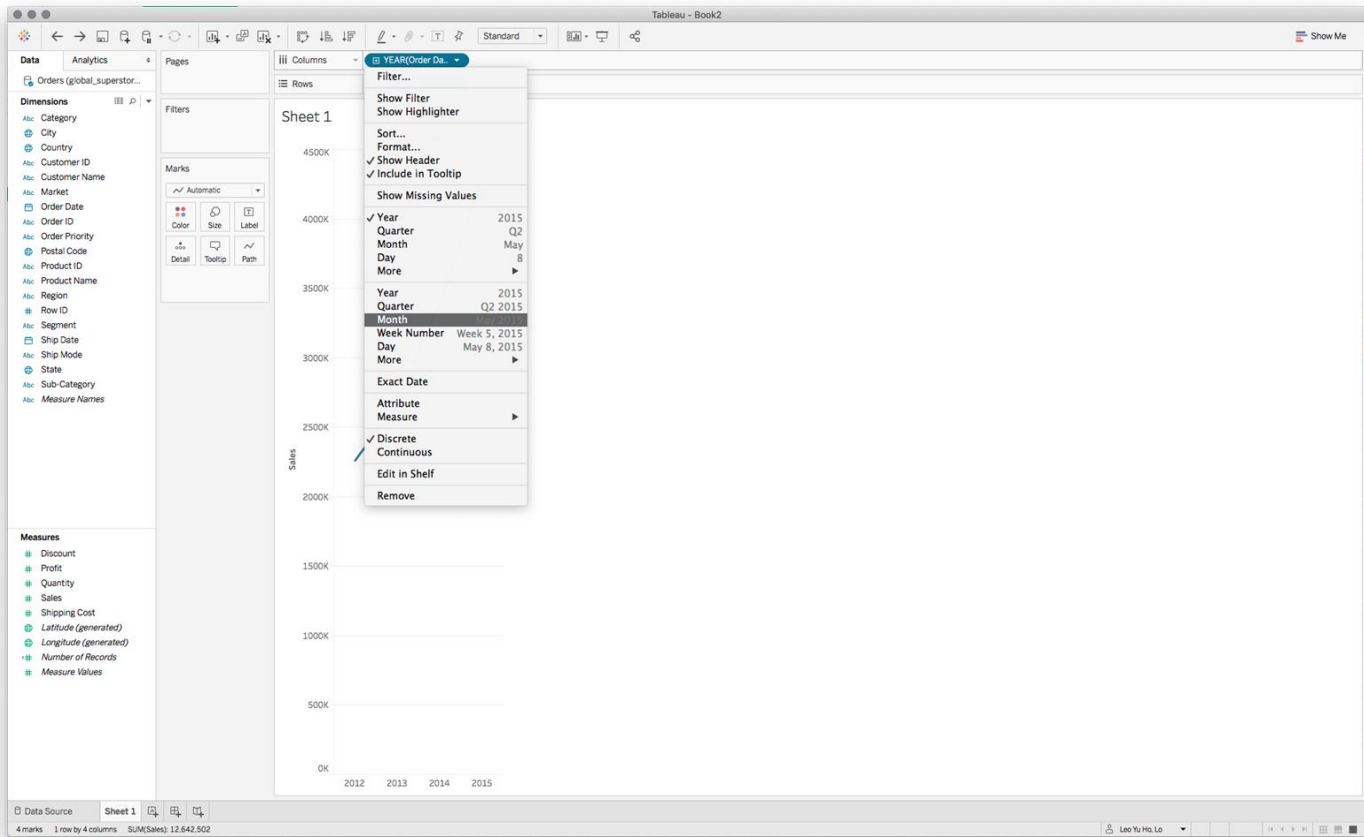
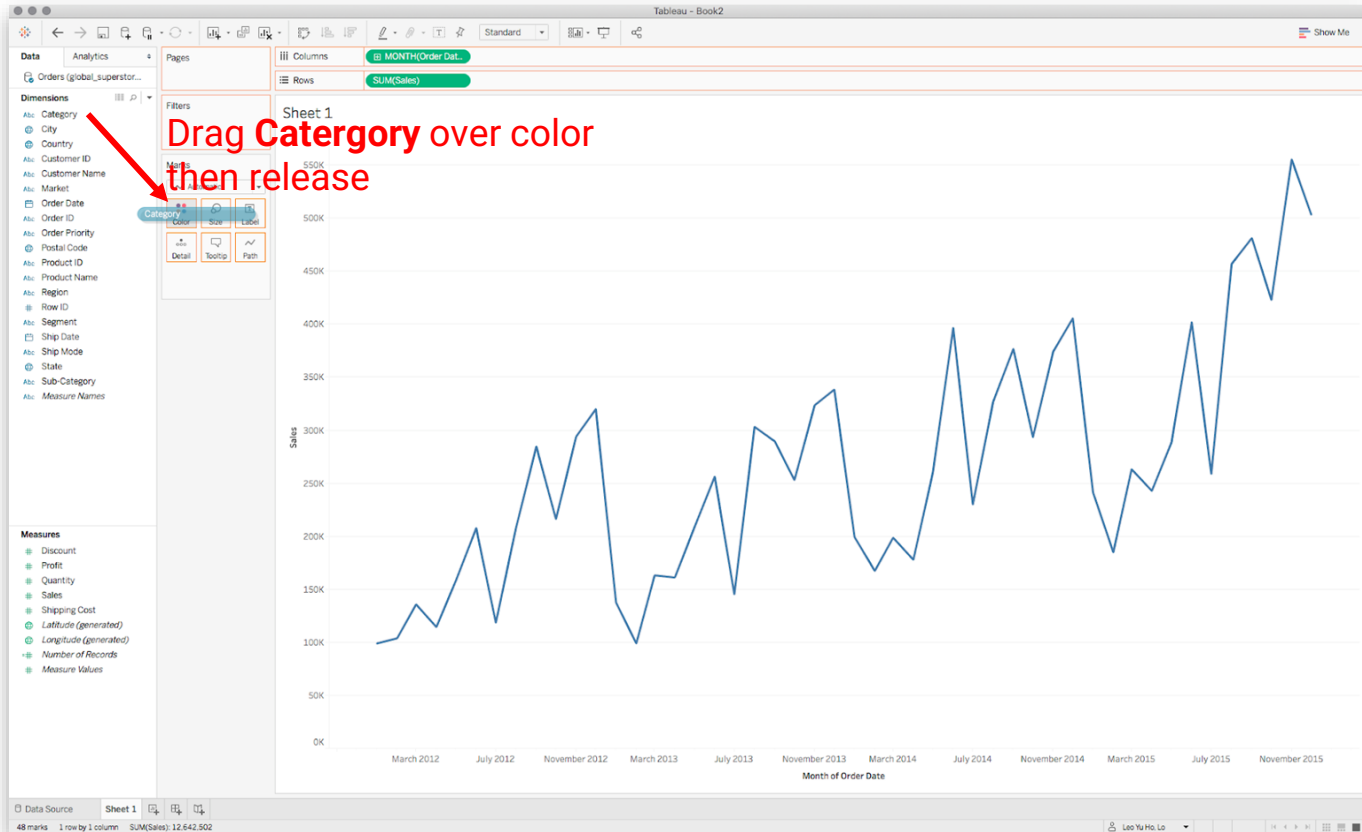# Change to Sheet 1

# Basic Plotting: Select Row and Column

# Basic Plotting: Select Row and Column

# Basic Plotting: Adjust Date

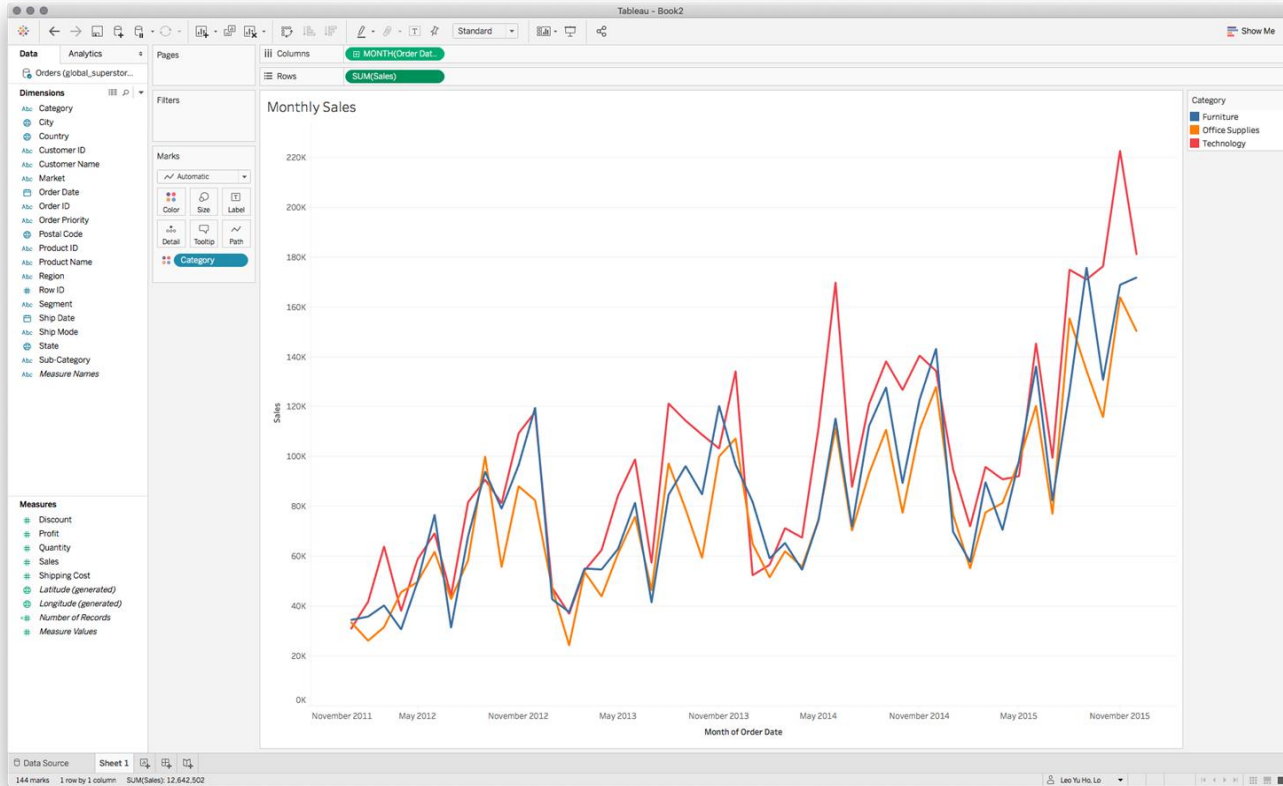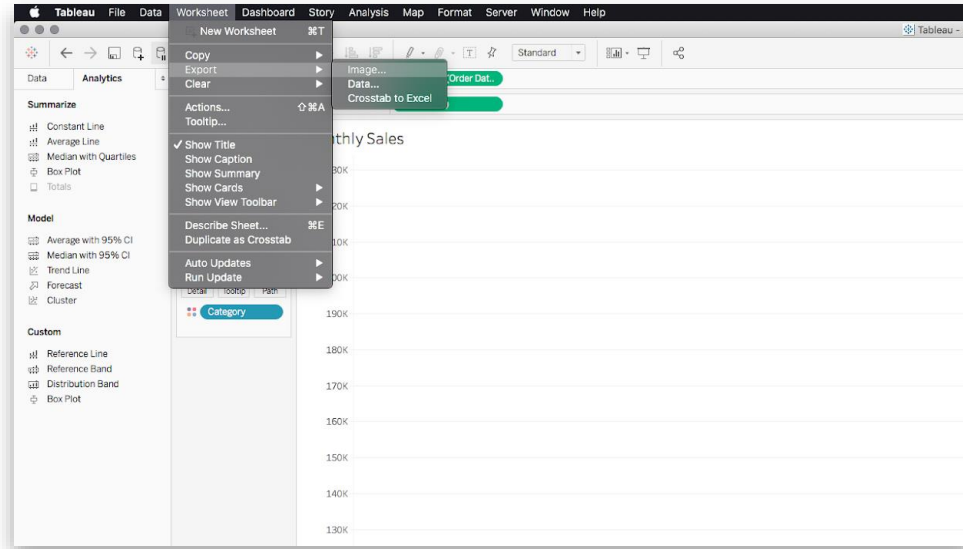# Basic Plotting: Marks with Color

# Basic Plotting: Ta-Da

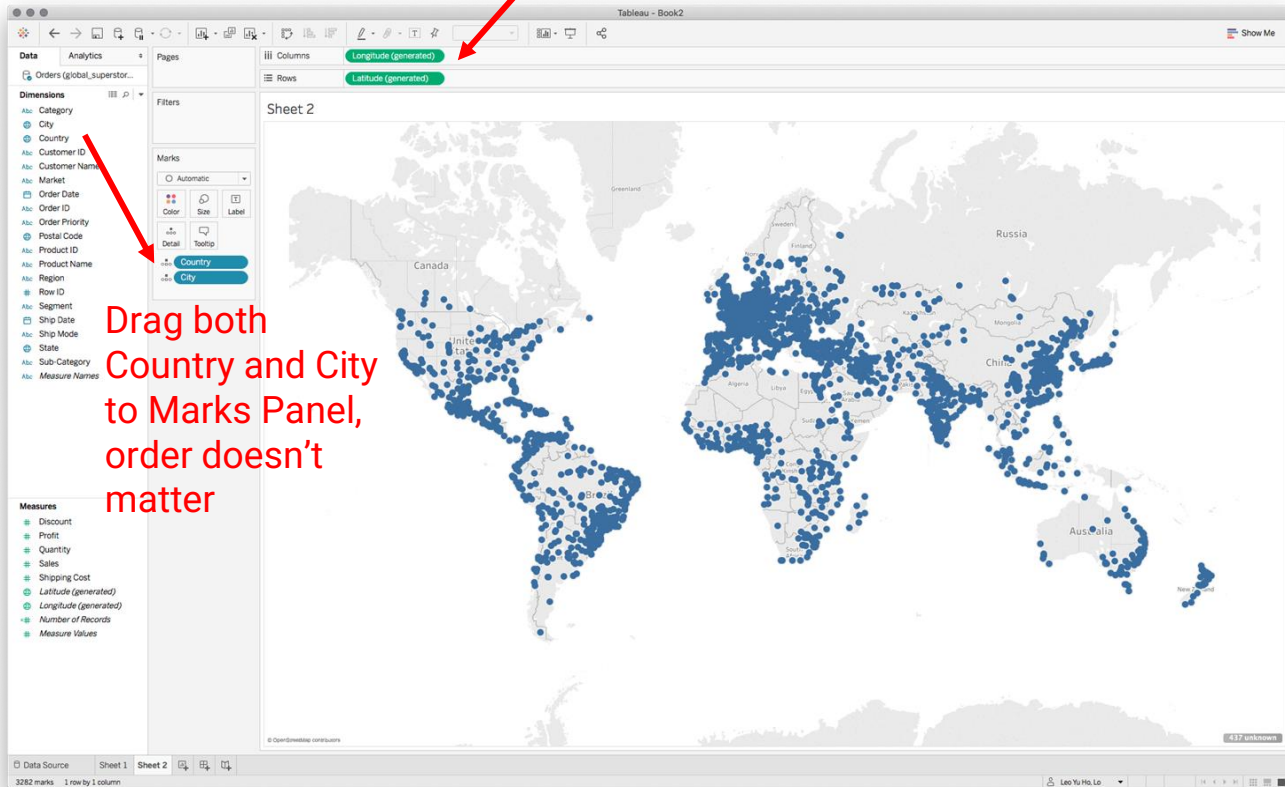# Export Image (Not available in Tableau Public)

# Plotting with Map

Create a new sheet first.

Drag Longitude to Columns, and Latitude to Rows

Drag both Country and City to Marks Panel, order doesn't matter

# Plotting with Map: Encode with Size and Color



Drag **Profit** over color, and **Sales** over size

# Adjust Color and Size

# Calculated Field

Create another new sheet.



For Tableau desktop, click "Create Calculated Field" by right clicking on "Dimensions" or "Measures"

For Tableau public, click "Create--Calculated Field" by right clicking on "Dimensions" or "Measures"

# Calculated Field



Enter a calculation formula by dragging fields

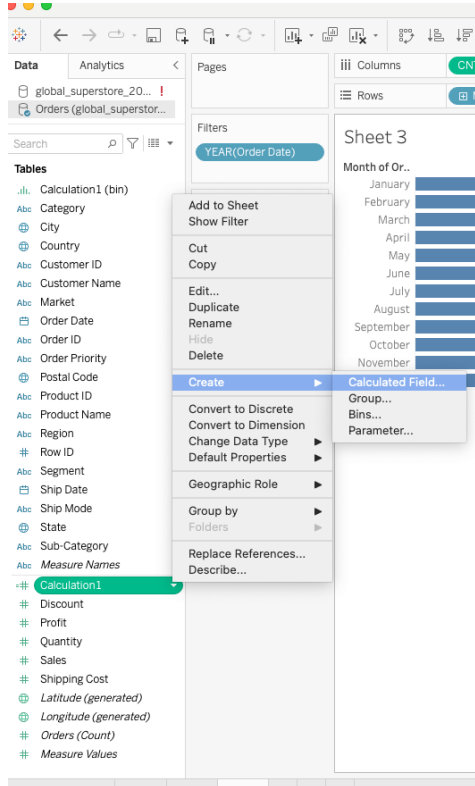# Aggregate



Add to rows and then pull
out drop-down menu

# Plotting is easy using "Show Me"



Gray-out means not suitable for the current selected data types.

Try different combinations for different plots!

# Filter

# Filter

# Interactive filtering



Check "Show Filter", a filter panel will show up on the right edge

# Lab exercise

- Tasks
  - Download dataset (world_data.csv) from GitHub
  - Import data from **Text** file
  - Watch the video Hans Rosling: 200 years in 4 minutes - BBC News
  - Recreate the bubble chart in the video with data of **2020** （An example is shown in the next page）
    - Column: Gdppc Cppp
    - Row: Life Expectancy
    - Size: Population
    - Color: World 4Region
    - Label: Country
    - Remember to set a filter to get Year 2020 data
  - Take a screenshot and upload to Canvas in .png format
    - Mac: cmd+shift+4
    - Windows: Snipping Tool

# Lab exercise

- Optional
    - Try using "Show me" to create different charts
    - Plotting data on map, adjust color and size
    - Create "Calculated Field", e.g. total GDP
    - Add an interactive filter

# Lab exercise



An example of submitted image.

# More topics on Tableau

- Coursera course
  - https://www.coursera.org/learn/analytics-tableau
- Tableau training videos
  - https://www.tableau.com/learn/training
- Tableau Viz Gallery
  - https://www.tableau.com/solutions/gallery
- Other notable features of Tableau
  - Dashboard, Storyboard
  - Parameters
  - Grouping
  - Table join
  - Features in "Analytics" tab, e.g. Trend Line, Cluster
  - Quick table calculation (e.g. running sum)
  - Tableau Prep
  - Import data from pdf

# Next Tutorial

Data scientist toolbox: Python, Jupyter Notebook and Pandas

- Prepare your Google account beforehand
  - For using [Google Colab](#)
  - Jupyter notebook environment
  - Free!
  - No setup
- Alternatively, you can use Jupyter notebook on your computer, but that is cumbersome
- Learn more about Pandas