# 2.1 Collect Datasets & Generate Graphs step-by-step
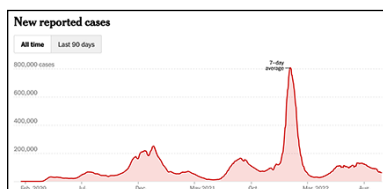
Canvas assignment page: https://sfsu.instructure.com/courses/17190/assignments/251100

NYT data repository: https://github.com/nytimes/covid-19-data
NYT coronavirus page (US): https://www.nytimes.com/interactive/2021/us/covid-cases.html

Below are the steps to download datasets from the NYT public repository of covid data (Github). Also, how to cut and paste data from NYT homepage, and to make some data visualizations from those datasets using Tableau.
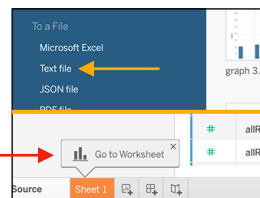
The first visualization is an area graph of covid cases in the US since Jan. 21, 2020. The NYT version looks like the image below.
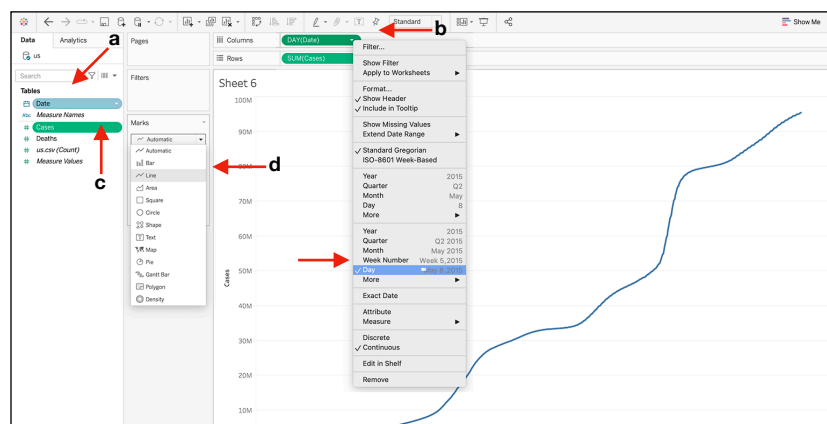


Download the datasets from NYT's Github repository:
https://github.com/nytimes/covid-19-data

Click green button (Code) —> Download ZIP.

Unzip and save the archive, it includes a file named us.csv which contains three columns: date, cases, deaths. The numbers are cumulative, so each cell is the total of the day's cases or deaths plus the tally from the day before.



In Tableau, import the us.csv text file; then start a new sheet.
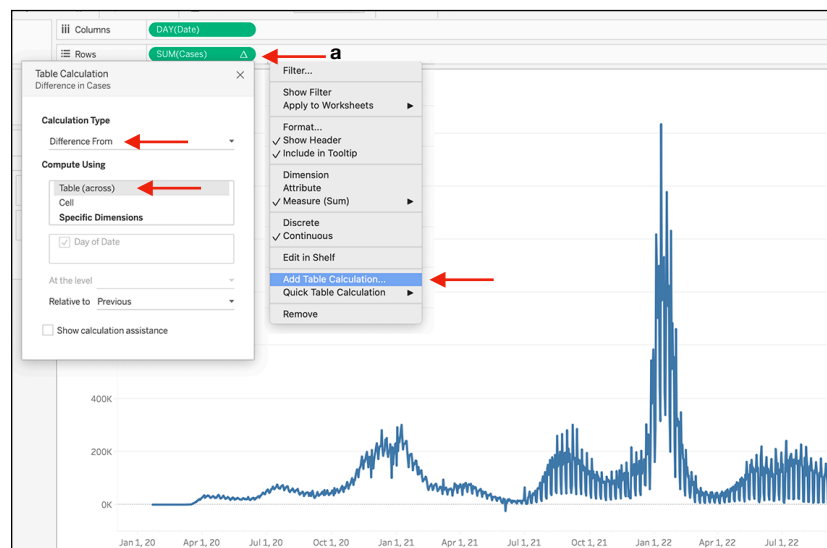


**NOTE:** screenshots are from Fall 2022. Your Spring 2023 graphs will look a bit different (for example the line will have 5 extra months with a hump on Jan. 2023.

**(1) a)** Double-click date; **b)** Right-click Date to change date type from Year to Day;
**c)** Double click cases; **d)** Change Mark type from Automatic to Line.

The Cases line is now visualized, but because the values are cumulative, the line only goes up. In the next step, you apply a calculation where each value subtracts from it the previous value. This will be the exact value for each individual day and the line will be an accurate representation of the data.



**(2) a)** Right-click Cases —> Add Table Calculation.

The first, default table calculation is Type: Difference From and Compute Using: Table (across). This is actually the calculation we want (subtract the previous day value from the current value).

The line is now accurate but because there are big variations in the number of cases from day to day, and because there are more than 1,000 days in the graph (more than 2.5 years), the line is hard to read as it jumps back and forth from peak to valley.

In the next step, change the Mark type from Line to Bar. Later, you will change back to line, but bars, with a separate line on top, was how the NYT graph started, so we will replicate it. Later, the bars became unwieldy to show the full length of the pandemic, and so now NYT only uses bars for the last 90 days of data.

**Note:** If needed, Play Video symbols are linked to extra video content. Datasets and layouts might be different.

US cases bars and 7-day average line (Tableau)

US cases bars and 7-day average line (Excel)

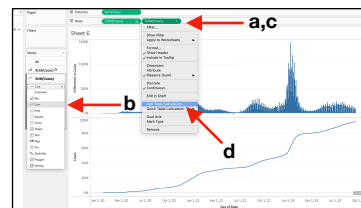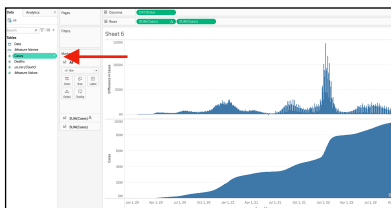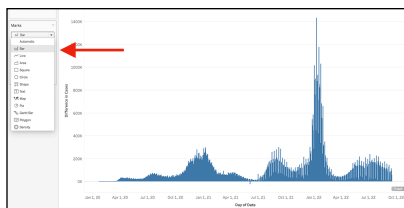California cases bars and 7-day average line (Tableau)

SF County cases bars and 7-day average line (Tableau)
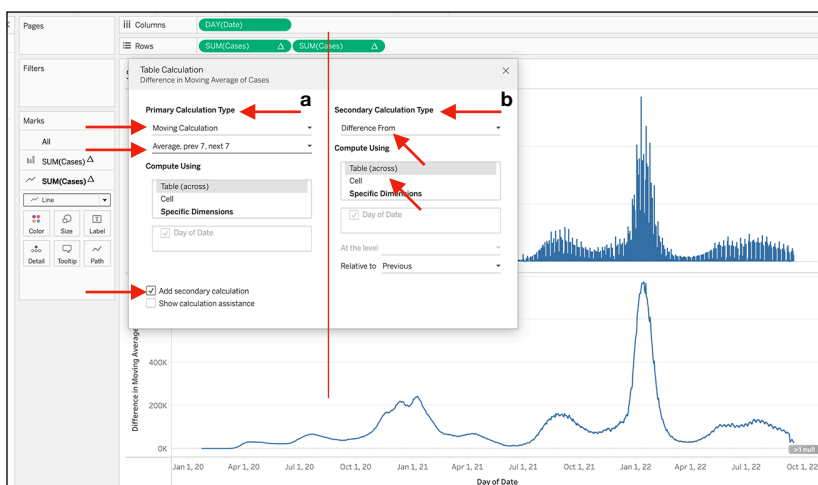
SF County deaths bars and 7-day average line (Tableau)

**(3)** Change Mark type from Line to Bar.

Because there are so many bars, one for each day for more than 1,000 days, they are now all bunched together and look like an area graph. It's OK for now. The next, multiple-step operation is to add a 7-Day Moving Average line on top of the bars. This will be a smoother line showing the weekly average of daily cases. The line will show the seven-day trend while the bars under the line will still show each daily count.

**Note: the colors in the images might not match your own colors.**

**(4)** Double-click Cases again. This will add a second bar chart below the first. Again, it looks solid because there are lots of bars. It also only goes up like in the first line graph. In the next steps you will change this second bar chart into a 7-Day Moving Average line.

**(5) a)** Select the second Cases button in the Rows field; **b)** Change Marks type from Bar to Line; **c)** Right-click Cases button —> **d)** Add Table Calculation.
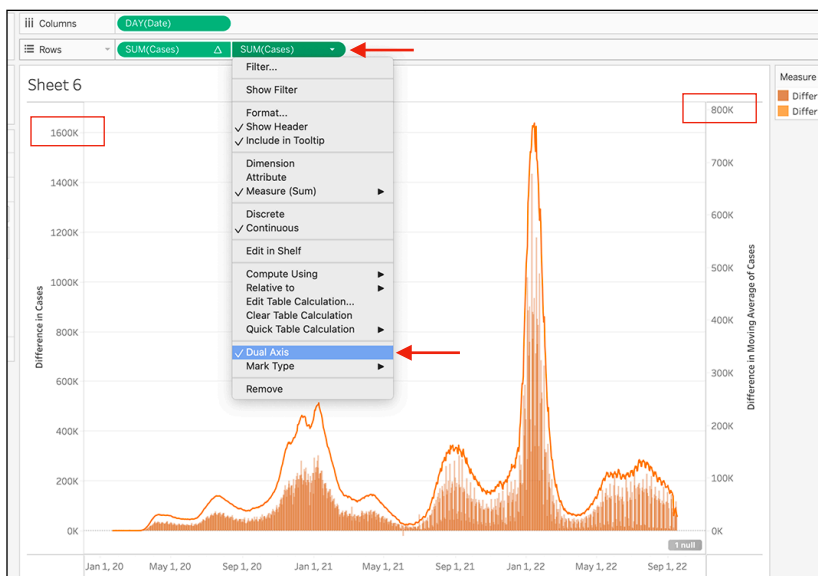


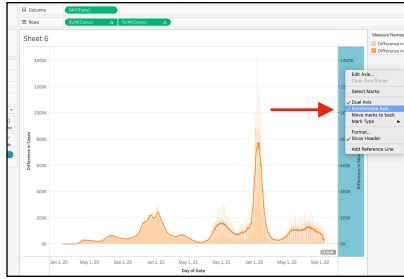**(6)** In this step the average calculation is set up first:

**a)** Under Primary Calculation Type, select Moving Calculation from the Calculation Type drop-down menu (menu not shown), then select Average from the Summarize Values drop-down menu (menu not shown); then Previous Values: 7; Next Value: 7. Then, check the Add secondary calculation box.

**b):** The default secondary calculation type is what you want: Difference From. The average line is rendered in the lower pane.
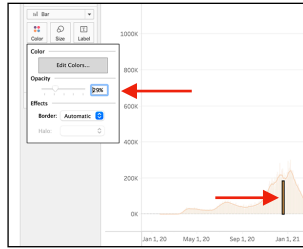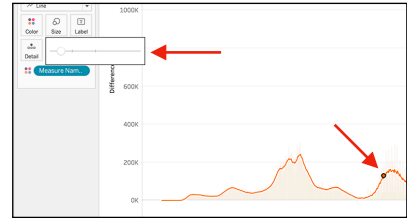


**(7)** Right-click Cases —> Dual Axis.

Line and bars now overlap (the colors may change and also be different for you). Note that the line does not appear to be an average somewhere between the tips of the tall and the short bars. That's because the two vertical axes (bars and line axes) are now on two different scales: bars at about 1600K and line at about 800K. So the two need to be "synchronized". (in next step).
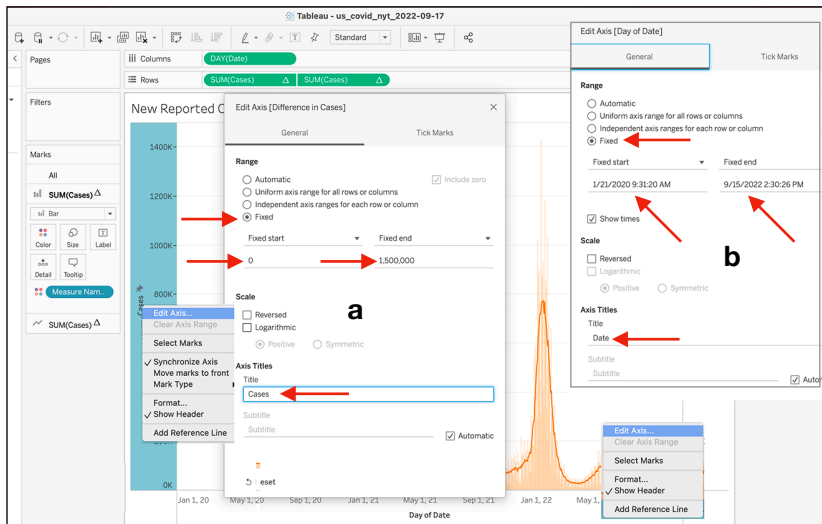
**(8)** Right-click right axis —> Synchronize Axis. Now the 7-Day Average line overlaps the bars.



**(9)** Edit the colors and the size of the bars and the line. Select a bar; change the color by clicking Edit Colors and choosing a different color in the color palette (not shown) or you can change the opacity of the current color. Click Size if you want to change the thickness of the bars, however here they are already the smallest possible size.



**(10)** Edit the color and size (thickness) of the line. Select the line. Use the Color and Size tools to modify the line. In this design based on the NYT graph, you want the bars a very light orange, and the line a bit thick and in a dark orange color.
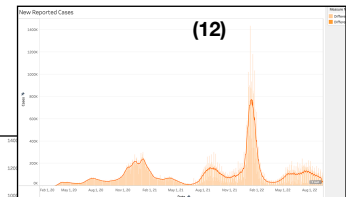


**(11)** Edit the vertical and horizontal axes.

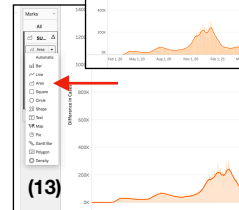**a)** Right-click the vertical axis (Cases): Range: Fixed —> Start: 0; End: 1,500,000; Axis Title: Cases.

Setting the origin at zero will also remove any "data anomaly", if for some reason there are any negative values that result in bars pointing downward.
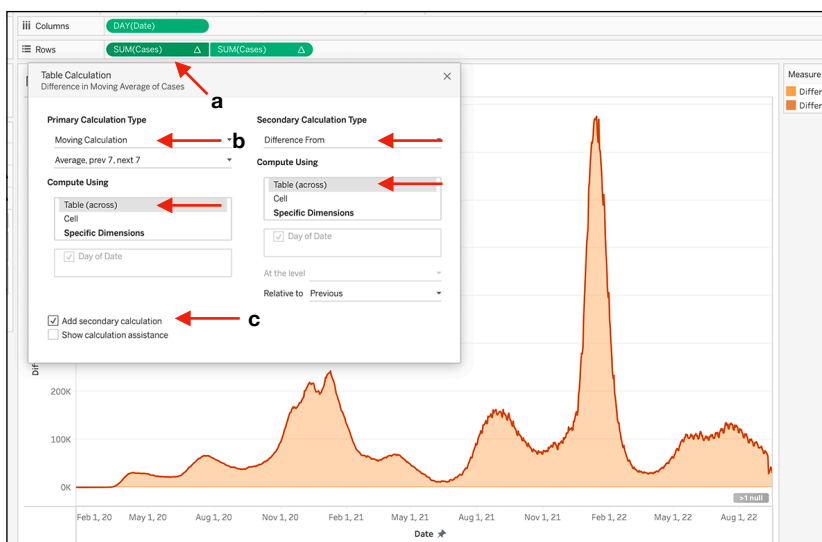
**b)** Right-click the horizontal axis (Dates): Range: Fixed —> Start: 1/21/2020; End: 2/27/2023; Axis Title: Date.
Note: your end date may be different later.



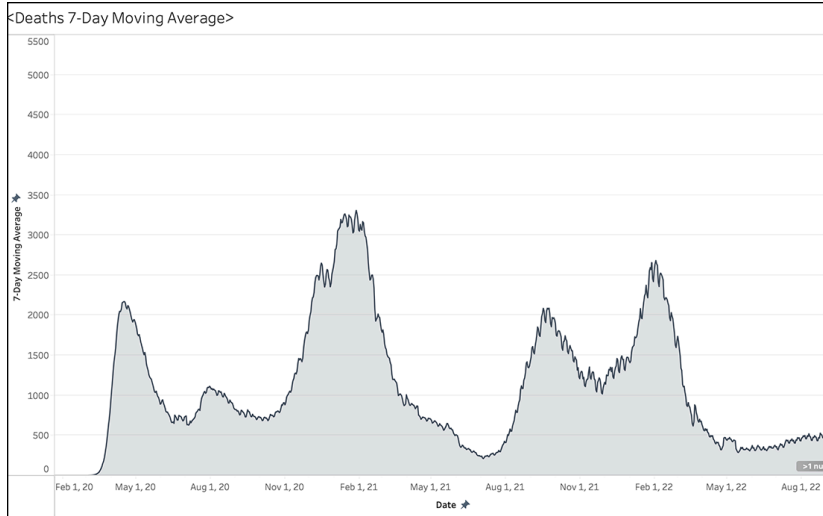**(12)** The final with both bars and line.

**(13)** To match the area/line chart from NYT (see picture on first page) you need to first change the Bar mark to Area.



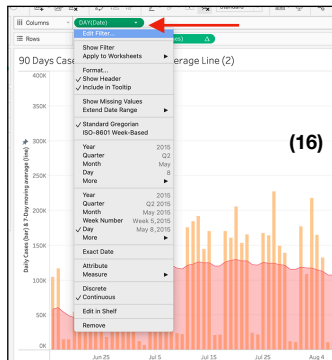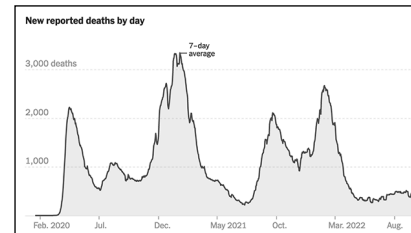**(14)** Then, smooth the edges of the area by applying the same average calculation as the line.
**a)** Right-click the left Cases button.
**b)** Change the existing Difference calculation (not shown) to the 7-Day average as for the line. See Fig. 6a and text for detailed steps. (the area will temporarily go only up again).
**c)** Add secondary Calculation —> Difference from Table across. The edge of the light area will now match the darker line above it.

<Deaths 7-Day Moving Average>

**(15)** You can try to create the same area chart as #14 but using the Deaths variable instead of Cases. Here I used a light gray for the area and a dark gray for the line, trying to match the same chart on the NYT website, reprinted below:
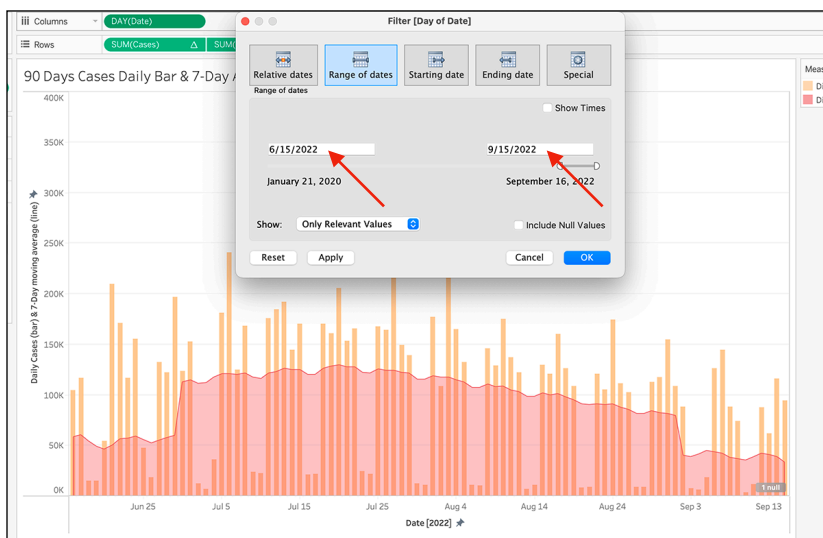


New reported deaths by day

**Cases in last 90 days**
In this variation, showing cases for only the last 90 days, follow the same steps from #1 to #12, but before #3 … apply the date filter below:

**(16)** To apply the filter, right-click Day/Date in the Columns field.

Apply a filter to Day/Date to restrict the time frame to start at Dec. 1, 2022 and end at Feb. 28, 2023.
(Screenshot may look slightly different).



**(17)** Filter the Day/Date. Use the slider or type the desired dates in the left and right fields: here I used June 15, 2022 and September 15, 2022.

Change the dates to a later date when you actually do this graph, using an updated dataset. Note: the pics show the Filter steps, however these were done before all the other steps that resulted in the final chart shown here.

**California counties choropleth map (colored map)** showing average daily covid cases per 100,000 in last seven days (exact dates on day of web scrape).



Click on link below:
https://www.nytimes.com/interactive/2021/us/california-covid-cases.html

On webpage, click Show All button to expand the table and include all counties. **(18)** Select all table information including column headers. Copy and paste the data into a new Excel file.

**NOTE: for this type of cut-and-paste and in general, I highly recommend that you use an external mouse as an input device, instead of the trackpad on your laptop.**



**(19)** Save the file after pasting the data into Excel. Note the following items that need work:
1. Header names are in two rows; combine info into just one row, taking care to clearly label each column.
2. Delete the California row.
3. Note empty rows between data rows. Delete all empty rows.
4. The minus signs are dashes but should be hyphens. Find/replace all in text editor.
5. In text-editor, delete all percentage symbols (%).

Note: add "county" label in top-left cell. Also, add an extra column (not shown) titled "state" and write "California" in every cell.



**(20)** Scroll down in the file to note additional problematic characters:
6. Change data with "<" symbols; write "0.5"
7. A different dash that should be an hyphen; find/replace all.
8. Delete "Unknown" row.

Save the file as CSV, then reopen it in Excel to delete the rows, but reopen it in text-editor to do find/replace. Save as CSV.



**(21)** This view in BBEdit (text editor) shows: New State column (arrow)
a. Simplified header names (no spaces)
b. All percentage symbols (%) have been removed.
c. Note quotation marks if data point is 1,000 or higher. This is OK here as Tableau will interpret the data correctly.

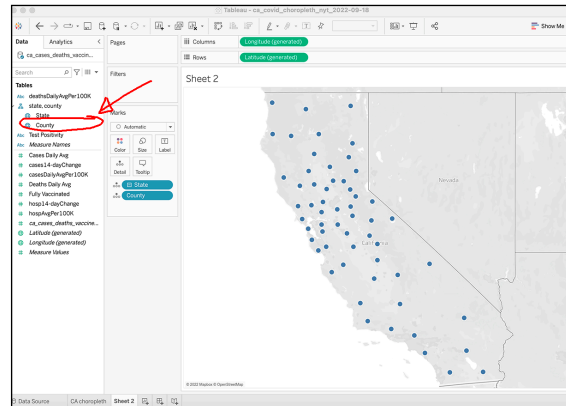Saving the file as Excel or CSV should not matter here. Tableau will work with both.

**(22)** The same file opened in Excel, showing the County header label, the additional State column, the other simplified header names. The percentage symbols (%) have been removed. Save the file and connect to it from Tableau.
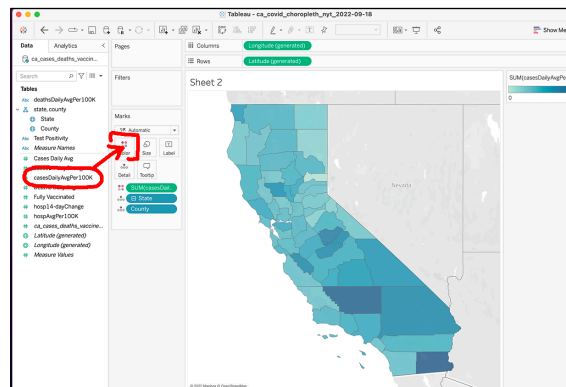
US cases per 100K choropleth map (Tableau)

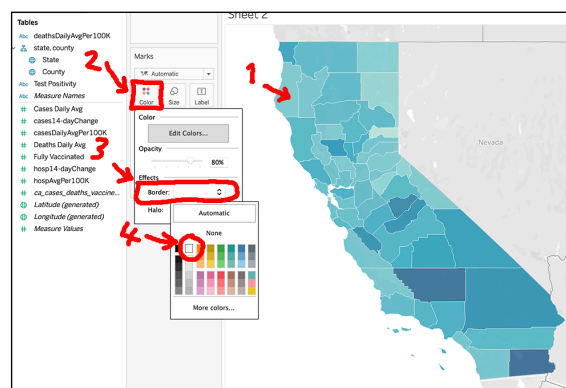California deaths per 100K choropleth map (Tableau)



**(23)** After connecting to the file in Tableau, start a new sheet and double-click County in the Tables pane. Sometimes not every county gets visualized by a dot and some counties might be missing. If this happens (you will notice it later if there are unfilled color areas), double-click State as well.



**(24)** Drag casesDailyAvgPer100K (or the corresponding name in your file) to the color symbol. Each county now is colored with a shade corresponding to its number of cases (darker color = more cases).
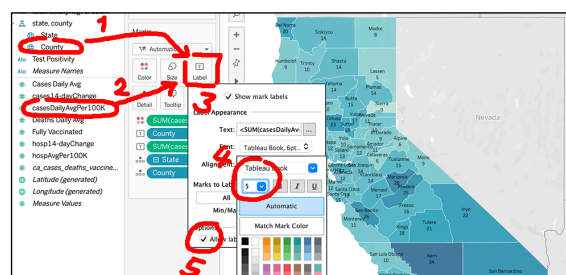
**Note: Map colors will reflect most recent data used.**



**(25)**
1.   Click on a county.
2.   Click Color.
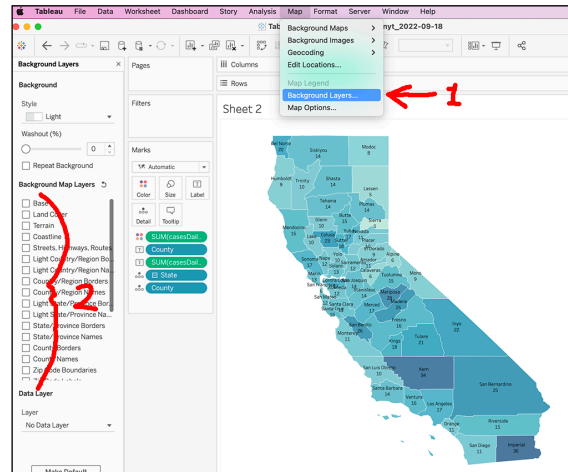3-4.  Border —> White

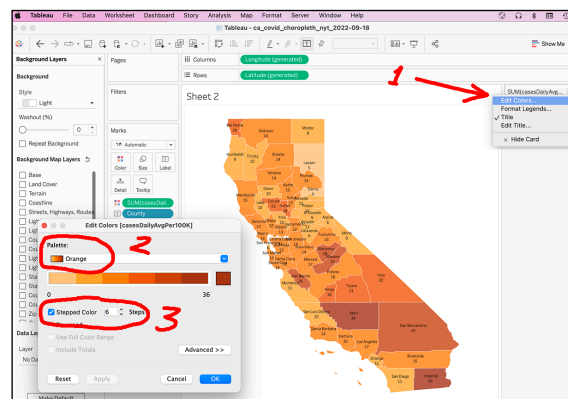This will make the county borders white, helping to separate them from each other.



**(26)**
1.   Drag County to Label.
2.   Drag casesDaily… to Label.
3.   Click Label
4.   Label type size: 5 pts
5.   Check "Allow labels to overlap…"

This will add the county names and the cases number to the map. Note: if "California" gets added and repeated with the county names, delete "State" under the marks pane and tools.

**(27)** Remove map background elements:
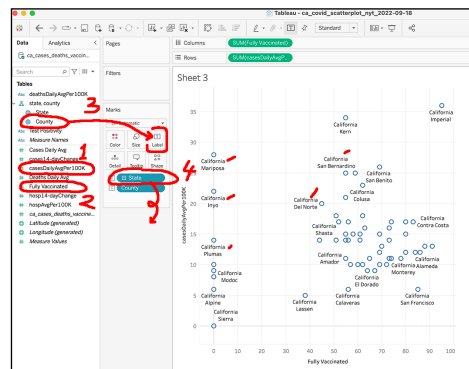1. Map —> Background Layers
2. Uncheck all Background Map Layers.

The map is much cleaner now.



**(28)** Change the color to orange and add steps:
1. In color key at right, click Edit Colors…
2. In pop-up menu, select a different color scale.
3. Check Stepped Color; 6 steps.
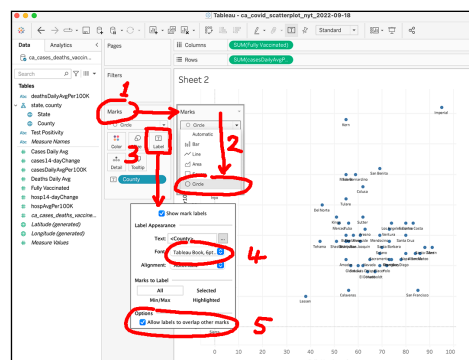
Save the Tableau file and Print/Save the graph as PDF.

California deaths per 100K choropleth map (Illustrator edits)


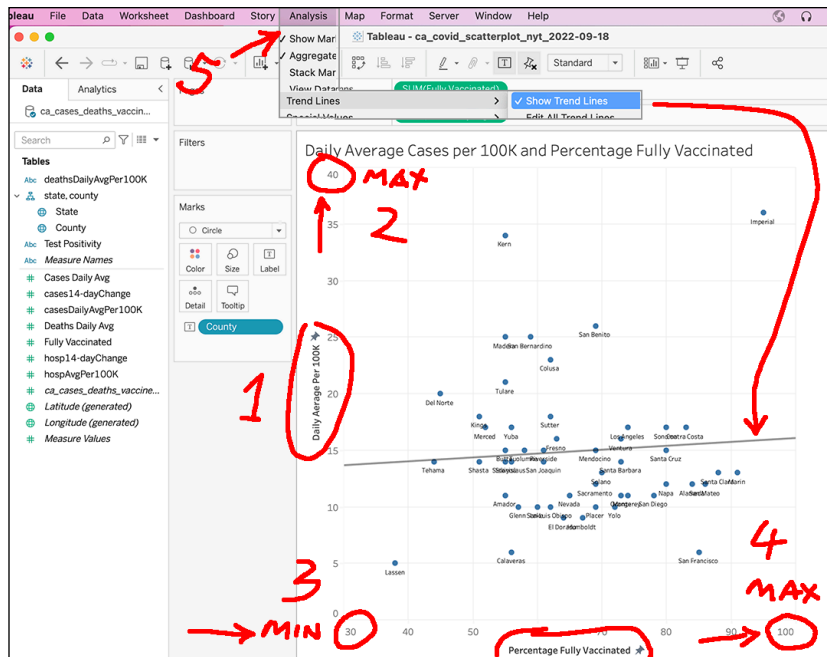


**CA cases and vaccinations scatterplot**
Start a new Tableau file and connect to the same dataset used for the California choropleth map (#s 18–28).

**(29)** Steps:
1. Double-click casesDaily…
2. Double-click Fully Vaccinated
3. Drag County to Label
4. Delete State item under tools if "California" label is being repeated next to each county name on the scatterplot (Right-click to delete).
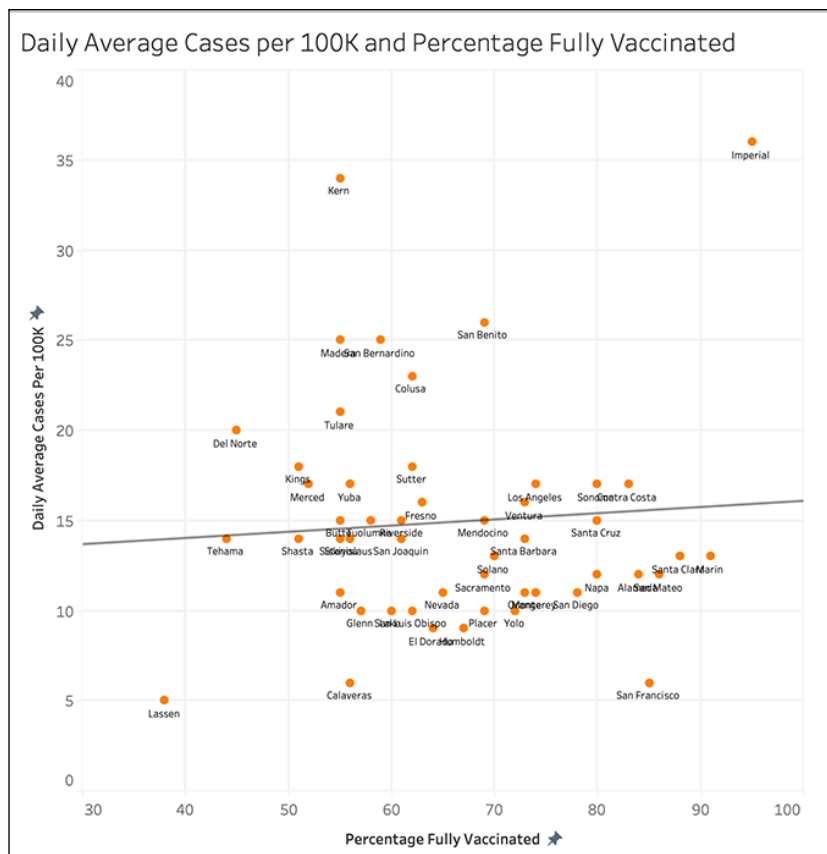


**(30)** Steps:
1. Right-click Marks and change from automatic to…
2. Circle (solid dot)
3. Edit Label
4. Type size —> 6 pts
5. Check "Allow label to overlap…"

California counties scatterplot cases & vaccines (Tableau)

**(31)** Steps:
1. Right-click and edit vertical axis: Simplify axis label wording
2. Fixed axis: Max: 40
3. Right-click and edit horizontal axis: Min: 30 (percent)
4. Max: 100 (percent) and add "Percentage" to axis label
5. Add a trend line: Analysis —> Show trend lines



Daily Average Cases per 100K and Percentage Fully Vaccinated

**(32)** Change colors if desired. Right-click color and edit colors (tools not shown).

Save Tableau file and Print/Save as PDF.