

# Understanding Trends in Influenza Strains Over Time

Ashkan Nikfarjam

## Objectives

I wanted to understand in my project how influenza trends have changed over time in California.

## Dataset

The dataset comes from the California Government Influenza Surveillance data. It is updated annually. The current data from public health laboratories runs from 2009-2020.

The raw dataset contains 30073 entries, with each row containing the following columns:

- `season`
  - string
  - formatted as `{YEAR\_1}-{YEAR\_2}` (e.g. 2009-2010)
  - a season runs from the 40th week to the 39th week
- `date\_code`
  - number
  - formatted as `{YEAR}{WEEK}` (e.g. 200952)
- `weekending`
  - date
  - formatted as `MM/DD/YYYY` (e.g. 10/04/09)
  - this is the last day in the date\_code week
- `region`
  - string
  - one of: "Bay Area", "Lower Southern", "Central", "Northern", "Upper Southern", "California"
  - each region is defined as a set of counties
- `Influenza\_Category`
  - text
  - formatted as `Influenza\_{TYPE}` (e.g. Influenza\_B)
- `Count`
  - number
  - number of specimens meeting the criteria for the category lineage

## Results

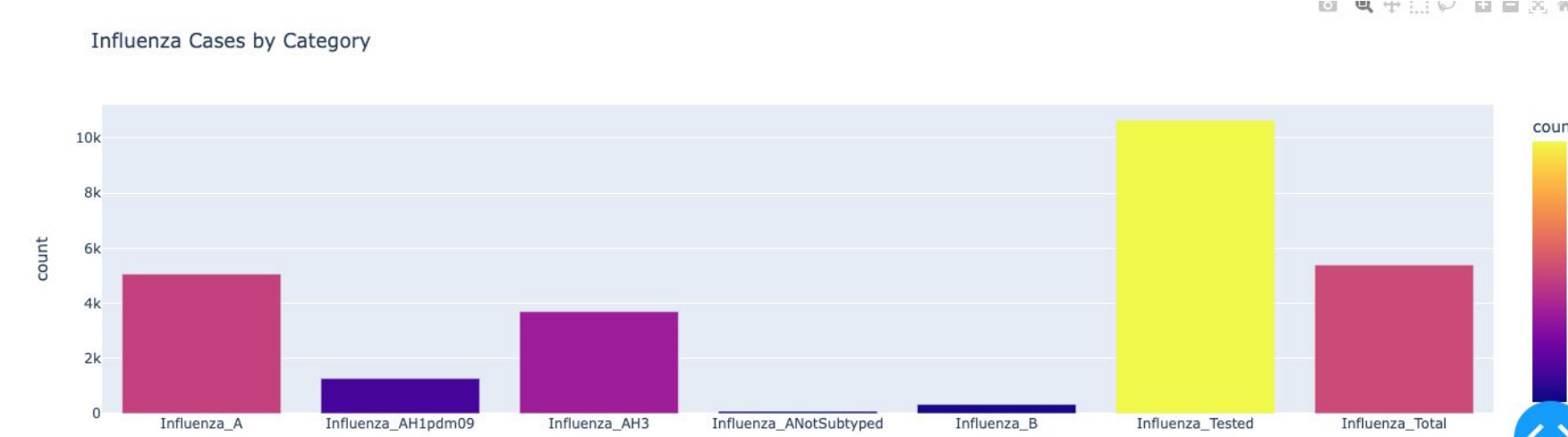
I was interested in finding out if i could use my visualization to be able to assess each region total performance in management of this virus. I used bash script, employing awk, to group and aggregate the data. I used python libraries and packages to create a heatmap that represented the death rate based on region and the type of Influenza.

Based on the result of my visualization I found out, the Bay Area and the Lower Southern california region had the highest death rate. But it also because the density of the population in these regions.

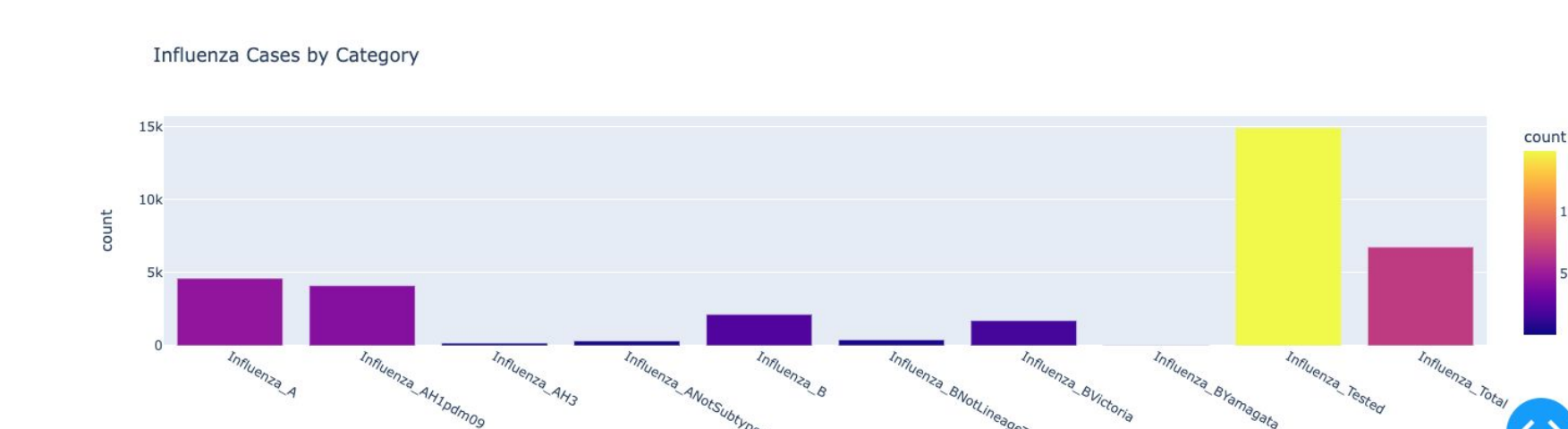


I also wanted to analyze the year-over-year trends for each influenza strain. To do this, I grouped data by year and flu strain. This helped me find that while strain A is overall more prevalent than strain B, strain B has been increasing in ratio with A.

2012



2020



## Tools and Commands

I used bash scripts to preap the data for visualization

- from the original data set i groups and aggregate data by year. Using awk, and sort in my yearly.sh bash script. The result data set is acting as the master data set that i aggregate most of the rest from it.
- yrly\_data.sh creates different data set one for each year from the master dataset. It includes region, type, count. I dropped the year column because the name of the file was sufficient indication of year.
- type\_year.sh create several datasets from the yearly one that i created in previous task, and it shows the total death count for each individual type for each year.
- total\_death.sh and total\_type.sh creates two data sets that shows the total death based on region and on based on influanzatype. The output are stored in total\_Deat.csv and type\_data.csv in the root file.

For visual representation I created a python Flask application.

I employed Pandas and Plotly libraries.

## References

- [1] <https://data.ca.gov/dataset/influenza-surveillance>

