Title: Pet Data Insights

Part 1
Pet Data CSV
2022

Part 2
Reddit Data Crawling
2023. 4 — 2024. 4

Part 1:
Shiny App:
- Main steps:
1. Create separate shiny apps ... & merge into one major app. Then merge w/ part 2 together to a final integrated app
2. Data Wrangling:
  clean & merge cvs files to a single csv
3. The initial variables we wanted to investigate are vaccination rate, expenditure, adoption rate, microchip rate.
  BUT During data wrangling, we decided not to simply visualize these two variables. Instead we divided return ÷ adoption → adoption rate
1. expenditure → distribution of expenditure
      which visualization to use?

? possible options: stacked bar ✓
                    grouped bar → too many states
                    tree map (×) → effect not good

5. Main Plots:
① choropleth geo plots
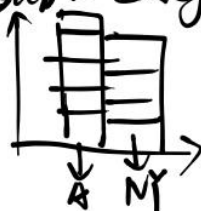        (USA)         ← vaccine
                        microchip
                        return rates
                        total score



⬇ leaflet → interactivity

② stacked bar plot for expenditure
  In order to display sub - cost of each sub - category | interactivity rendered by shiny.
                                The user can select state data
                                to display ( >1) or specific "All"

# Challenges (Solved) (Hawaii & Alaska are excluded)

1. Since the data set only has 48 state data "map" function cannot be used directly. ⟹ Manually added geometry, lat, long data to the dataset.

st_

states_data + state-st = states_data_st
                         (geometry)
                          polygon

---

Leaflet choropleth:
- Challenges: adding the ABBR on the static geomap. solved by manually adding abbrs to the dataset

(Solved)

- tweak the legend for expenditure ⟹ Fail to make it shiny because of the legend
  ↓
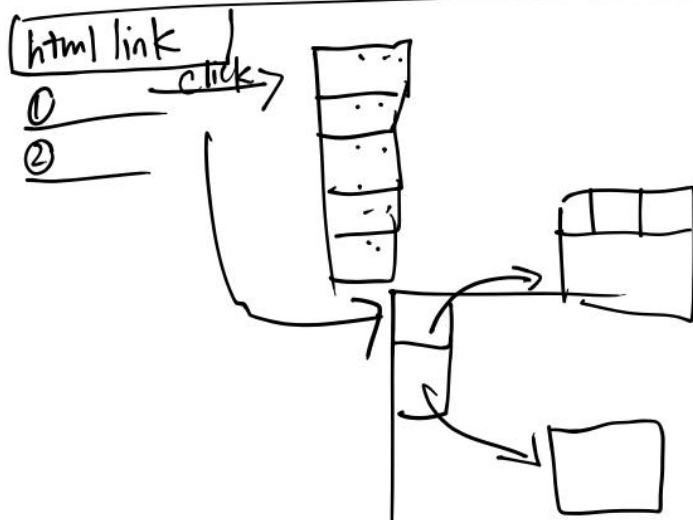  simple leaflet html plot for displaying dog vs. cat population ratio in 2022
- change the Pal multiple times in order to keep the palette consistent and obvious & * express sentiment & emotions.
- eg. adoption - return rates is ∧ just wanna make the state w/ highest return rates obvious

---

Adjust the UI height too 600, put the legend at the bottom right

---

if we finally merge 2 shiny apps to one, how to embed the html leaflet? (Dog vs Cat population)
                                       ↓
embed html to the main dashboard title?

---

(html link)
                click→
0 _____
② _____

# Part 2: Reddit Data Analysis & DV

## Step 1: How to get the data (So hard !!!!)

① = Registered for Reddit API to directly fetch data

⇒ **Failed**, due to connection issues. ☹

② Searched online for existing code

⇒ Even worse, Recent updates to Reddit's terms imposed stricter data retrieval limits and robust anti-Scraping measures

② Found a script for the real time data scraping

⇒ Success, but not access historical

④ Further efforts: Combined Several online code example.

⇒ Get some past data, limits 1000 record, within a year.

Final Strategy = Created three Separate loops to extract data by year, month, week
This is the maximum data collection. within given restrictions.

## Step 2: Data cleaning

① remove duplicates, special characters, space

② remove basic stop word & name

③ Extracted three detailed articles from a pet-focused website.
{ Clean this three. txt, just like what I did in ②
Select Top 100 most frequency words.
generic pet specific stop-word list

④ clean data again.

## Step4 = Data Analysis

/\ Topic Model ⇒ Using Data 2023/04 − 2024/04

① Creat DTM

② Using LDA to get topic. Words. (5)
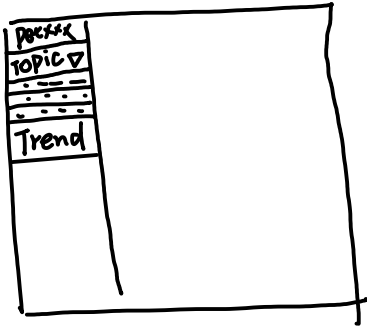
③ Calculate topic word Frequency

④ Calculate. Coherent Score.

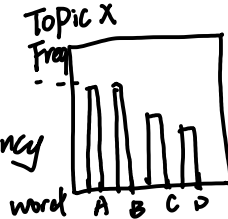2. Topic Trend Analysis ⇒ using data segmenting into Season
Just do what 1. Topic Model did for each Season.

Step 5 = ✗ Data Visualization ✗ ⇒ Shiny !



Docket
TOPIC ▽
Trend
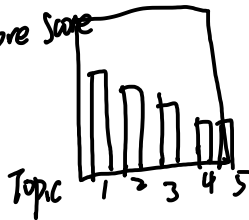
TOPIC X
Freq
word A B C D

Topic name & word frequency

coherent Score Score
Topic 1 2 3 4 5

Topic Model

word cloud

Trend ⇒

Season
24S5
24S1
23S4
23S3
23S2
word A B C D

heat