

Project 1

Trever Yoder and Koji Takagi

Load Packages and Functions

In this section, we load all necessary libraries and our custom functions file.

```
library(tidyverse)
library(readr)
library(ggplot2)

# Load custom functions
source("functions.R")
```

Task 1: Data Processing

Question 1: Read in the dataset

We want to read in some of this Census data set, but not all of it. Here we specify which columns we want to read in and we named this data set: `df_selected`. We then `slice` the first 5 lines to display them to confirm we read the data in correctly.

```
#Read in the data while selecting specific columns
df_selected <- read_csv("https://www4.stat.ncsu.edu/~online/datasets/EDU01a.csv", show_col_types = FALSE)
df_selected <- select(df_selected, Area_name, STCOU, ends_with("D")) %>% #select specified columns
df_selected <- rename(df_selected, area_name = Area_name) #rename "Area_name" as directed

#Display the first 5 lines
df_selected %>%
slice(1:5)
```

```
# A tibble: 5 x 12
  area_name      STCOU EDU010187D EDU010188D EDU010189D EDU010190D EDU010191D
```

	<chr>	<chr>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
1	UNITED STATES	00000	40024299	39967624	40317775	40737600	41385442
2	ALABAMA	01000	733735	728234	730048	728252	725541
3	Autauga, AL	01001	6829	6900	6920	6847	7008
4	Baldwin, AL	01003	16417	16465	16799	17054	17479
5	Barbour, AL	01005	5071	5098	5068	5156	5173

```
# i 5 more variables: EDU010192D <dbl>, EDU010193D <dbl>, EDU010194D <dbl>,
#   EDU010195D <dbl>, EDU010196D <dbl>
```

Question 2

Now we want to convert the data into long format where each row has only one enrollment value for `area_name`. This converted data will be called `df_long`. We then display the first 5 rows to make sure everything looks as expected.

```
df_long <- pivot_longer(
  df_selected,
  cols = ends_with("D"),
  names_to = "Survey",
  values_to = "Enrollment Value"
)
```

```
#Display the first 5 lines
df_long %>%
  slice(1:5)
```

```
# A tibble: 5 x 4
  area_name      STCOU Survey      `Enrollment Value`
  <chr>          <chr> <chr>          <dbl>
1 UNITED STATES 00000 EDU010187D      40024299
2 UNITED STATES 00000 EDU010188D      39967624
3 UNITED STATES 00000 EDU010189D      40317775
4 UNITED STATES 00000 EDU010190D      40737600
5 UNITED STATES 00000 EDU010191D      41385442
```

Question 3

Now we need to separate some values that are currently combined in `Survey`. The first 7 digits of `Survey` are currently a measurement (public school enrollment) and the last 2 digits followed by D are the school year. We want to separate these values to create 2 corresponding variables and turn the year into a 4 digit format. Since we will not be working with any data

that was before the year 1925 or after the year 2025, we can do some simple math. The Year 1987 will be referring to the Fall 1986-1987 school year.

```
#Separate and create variables from Survey
long_updated <- df_long %>%
  mutate(
    Year = as.numeric(substr(Survey, 8, 9)),
    Year = ifelse(Year > 25, Year + 1900, Year + 2000),
    Measurement = substr(Survey, 1, 7)
  )
#Display the first 5 lines
long_updated %>%
  slice(1:5)
```

```
# A tibble: 5 x 6
  area_name      STCOU Survey      `Enrollment Value`  Year Measurement
  <chr>          <chr> <chr>                <dbl> <dbl> <chr>
1 UNITED STATES 00000 EDU010187D          40024299  1987 EDU0101
2 UNITED STATES 00000 EDU010188D          39967624  1988 EDU0101
3 UNITED STATES 00000 EDU010189D          40317775  1989 EDU0101
4 UNITED STATES 00000 EDU010190D          40737600  1990 EDU0101
5 UNITED STATES 00000 EDU010191D          41385442  1991 EDU0101
```

Question 4

Now we want to create a data set for non-county data and a data set for only county level data. As directed, we will add a class to the county level data tibble that's called **county** and we will create a class for the non-county data called **state**. Then we will print the first 10 rows of each tibble to make sure they look correct.

```
#Create the county and state data sets
county_idx <- grep(" ", df_long$area_name)
county_tibble <- df_long[county_idx, ]
state_tibble <- df_long[-county_idx, ]

#add class accordingly
class(county_tibble) <- c("county", class(county_tibble))
class(state_tibble) <- c("state", class(state_tibble))

#display first 10 lines of county data
county_tibble %>%
  slice(1:10)
```

```
# A tibble: 10 x 4
```

	area_name	STCOU	Survey	`Enrollment Value`
	<chr>	<chr>	<chr>	<dbl>
1	Autauga, AL	01001	EDU010187D	6829
2	Autauga, AL	01001	EDU010188D	6900
3	Autauga, AL	01001	EDU010189D	6920
4	Autauga, AL	01001	EDU010190D	6847
5	Autauga, AL	01001	EDU010191D	7008
6	Autauga, AL	01001	EDU010192D	7137
7	Autauga, AL	01001	EDU010193D	7152
8	Autauga, AL	01001	EDU010194D	7381
9	Autauga, AL	01001	EDU010195D	7568
10	Autauga, AL	01001	EDU010196D	7834

```
#display first 10 lines of noncounty data
state_tibble %>%
slice(1:10)
```

```
# A tibble: 10 x 4
```

	area_name	STCOU	Survey	`Enrollment Value`
	<chr>	<chr>	<chr>	<dbl>
1	UNITED STATES	00000	EDU010187D	40024299
2	UNITED STATES	00000	EDU010188D	39967624
3	UNITED STATES	00000	EDU010189D	40317775
4	UNITED STATES	00000	EDU010190D	40737600
5	UNITED STATES	00000	EDU010191D	41385442
6	UNITED STATES	00000	EDU010192D	42088151
7	UNITED STATES	00000	EDU010193D	42724710
8	UNITED STATES	00000	EDU010194D	43369917
9	UNITED STATES	00000	EDU010195D	43993459
10	UNITED STATES	00000	EDU010196D	44715737

Process the EDU Data Sets

We run our wrapper function on the two EDU datasets and inspect the results.

```
edu1 <- my_wrapper("https://www4.stat.ncsu.edu/~online/datasets/EDU01a.csv", value = "Enrollment")
edu2 <- my_wrapper("https://www4.stat.ncsu.edu/~online/datasets/EDU01b.csv", value = "Enrollment")

# Inspect to ensure correctness
head(edu1$county)
```

```
# A tibble: 6 x 7
  area_name STCOU Survey `Enrollment Value` Year Measurement State
  <chr>      <chr> <chr>          <dbl> <dbl> <chr>      <chr>
1 Autauga, AL 01001 EDU010187D      6829  1987 EDU0101    AL
2 Autauga, AL 01001 EDU010188D      6900  1988 EDU0101    AL
3 Autauga, AL 01001 EDU010189D      6920  1989 EDU0101    AL
4 Autauga, AL 01001 EDU010190D      6847  1990 EDU0101    AL
5 Autauga, AL 01001 EDU010191D      7008  1991 EDU0101    AL
6 Autauga, AL 01001 EDU010192D      7137  1992 EDU0101    AL
```

```
head(edu1$noncounty)
```

```
# A tibble: 6 x 7
  area_name STCOU Survey `Enrollment Value` Year Measurement Division
  <chr>      <chr> <chr>          <dbl> <dbl> <chr>      <chr>
1 UNITED STATES 00000 EDU010187D      40024299  1987 EDU0101    ERROR
2 UNITED STATES 00000 EDU010188D      39967624  1988 EDU0101    ERROR
3 UNITED STATES 00000 EDU010189D      40317775  1989 EDU0101    ERROR
4 UNITED STATES 00000 EDU010190D      40737600  1990 EDU0101    ERROR
5 UNITED STATES 00000 EDU010191D      41385442  1991 EDU0101    ERROR
6 UNITED STATES 00000 EDU010192D      42088151  1992 EDU0101    ERROR
```

Question 3: Combine EDU Data Sets

Here we use our combining function to merge the two processed data sets.

```
edu_combined <- combine_wrapper_results(edu1, edu2)
head(edu_combined$county)
```

```
# A tibble: 6 x 7
  area_name STCOU Survey `Enrollment Value` Year Measurement State
  <chr>      <chr> <chr>          <dbl> <dbl> <chr>      <chr>
1 Autauga, AL 01001 EDU010187D      6829  1987 EDU0101    AL
2 Autauga, AL 01001 EDU010188D      6900  1988 EDU0101    AL
3 Autauga, AL 01001 EDU010189D      6920  1989 EDU0101    AL
4 Autauga, AL 01001 EDU010190D      6847  1990 EDU0101    AL
5 Autauga, AL 01001 EDU010191D      7008  1991 EDU0101    AL
6 Autauga, AL 01001 EDU010192D      7137  1992 EDU0101    AL
```

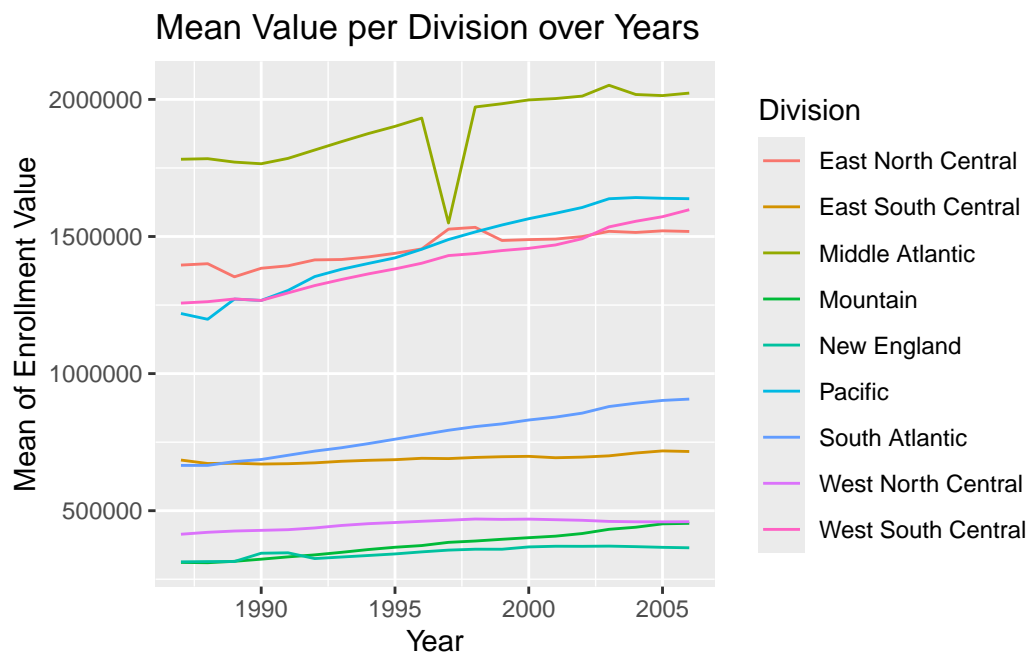
```
head(edu_combined$noncounty)
```

```
# A tibble: 6 x 7
  area_name      STCOU Survey `Enrollment Value` Year Measurement Division
  <chr>          <chr> <chr>              <dbl> <dbl> <chr>      <chr>
1 UNITED STATES 00000 EDU010187D         40024299 1987 EDU0101  ERROR
2 UNITED STATES 00000 EDU010188D         39967624 1988 EDU0101  ERROR
3 UNITED STATES 00000 EDU010189D         40317775 1989 EDU0101  ERROR
4 UNITED STATES 00000 EDU010190D         40737600 1990 EDU0101  ERROR
5 UNITED STATES 00000 EDU010191D         41385442 1991 EDU0101  ERROR
6 UNITED STATES 00000 EDU010192D         42088151 1992 EDU0101  ERROR
```

Question 4: State Plot for EDU Data

This plot shows the mean enrollment by Division across years.

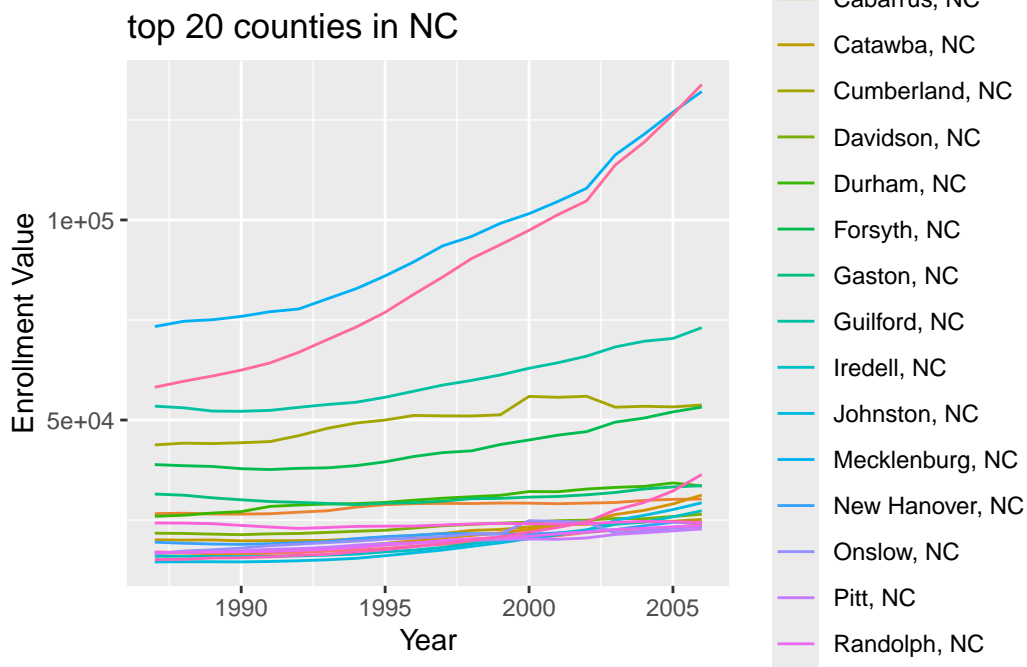
```
plot(edu_combined$noncounty, var_name = "Enrollment Value")
```



Question 5: County Plots for EDU Data

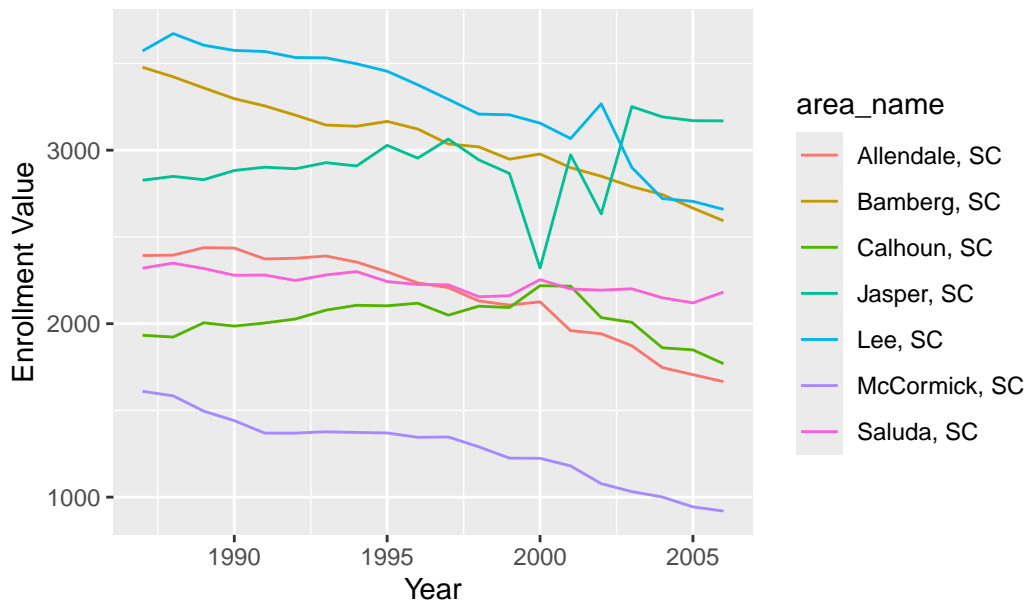
Below are various plots for county data, demonstrating flexibility in selecting state, top/bottom, and count.

```
# NC, top 20  
plot(edu_combined$county, var_name = "Enrollment Value", state = "NC", top_or_bottom = "top")
```



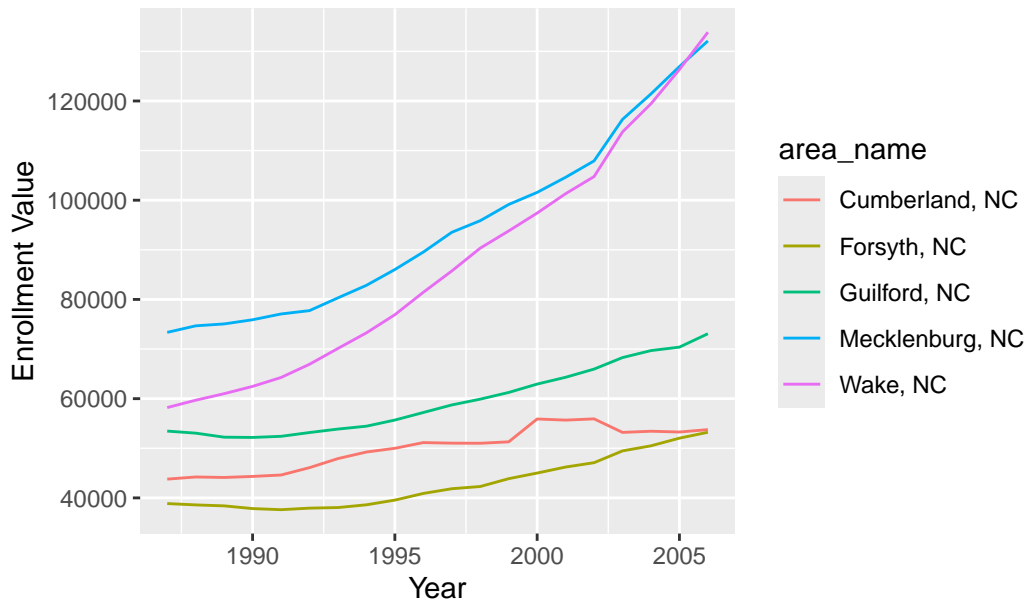
```
# SC, bottom 7  
plot(edu_combined$county, var_name = "Enrollment Value", state = "SC", top_or_bottom = "bottom")
```

bottom 7 counties in SC

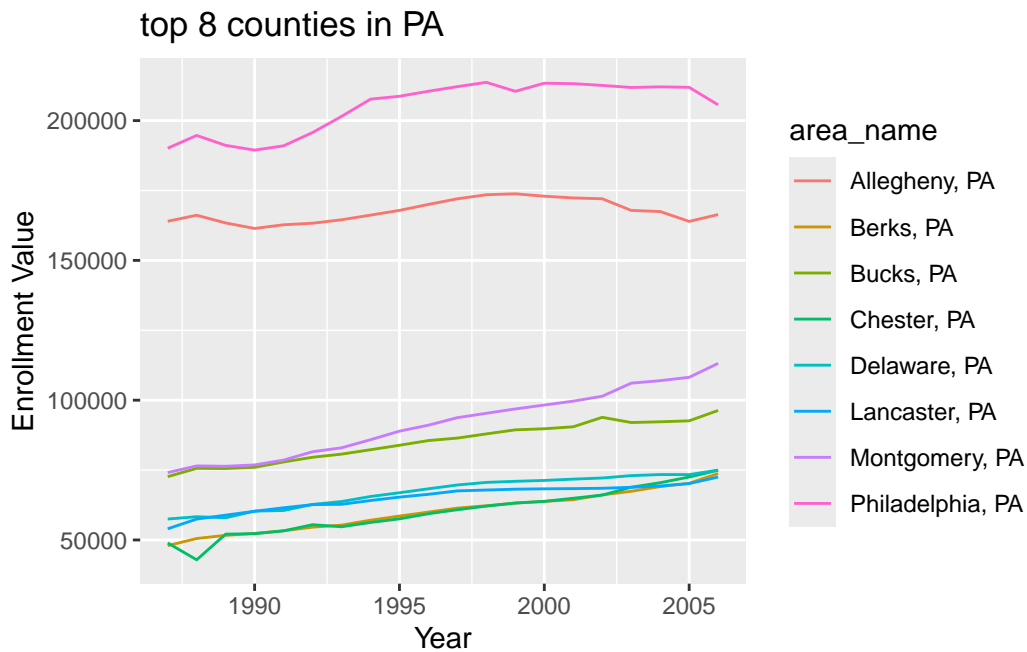


```
# Default (uses NC top 5)
plot(edu_combined$county, var_name = "Enrollment Value")
```

top 5 counties in NC




```
# PA, top 8
plot(edu_combined$county, var_name = "Enrollment Value", state = "PA", top_or_bottom = "top")
```



Question 6: Process PST Data Sets

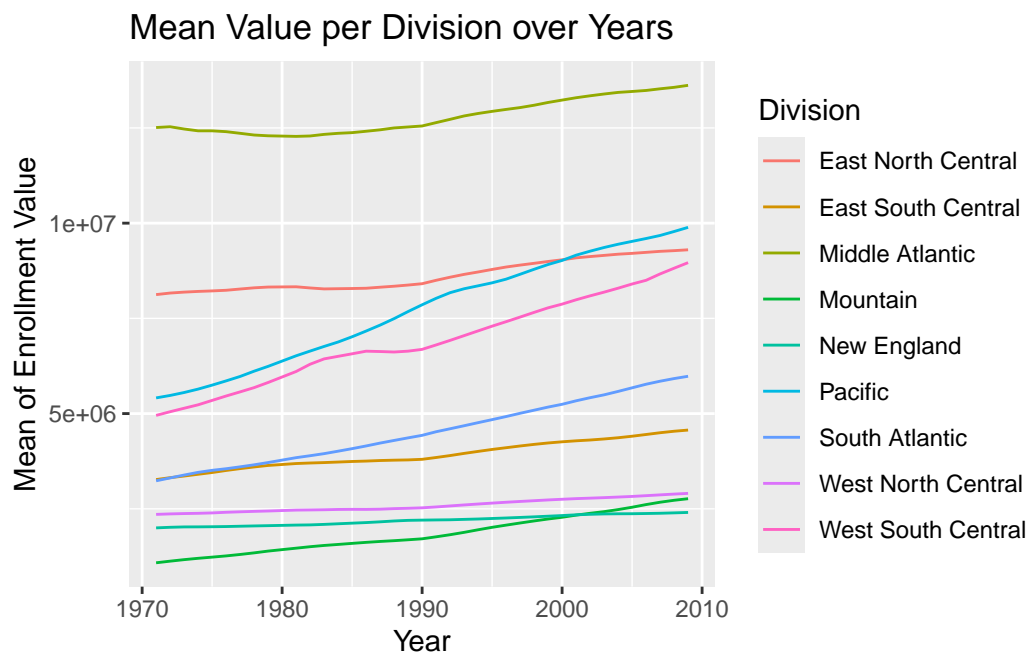
We repeat the same workflow for the four PST datasets.

```
pst1 <- my_wrapper("https://www4.stat.ncsu.edu/~online/datasets/PST01a.csv", value = "Enrollm")
pst2 <- my_wrapper("https://www4.stat.ncsu.edu/~online/datasets/PST01b.csv", value = "Enrollm")
pst3 <- my_wrapper("https://www4.stat.ncsu.edu/~online/datasets/PST01c.csv", value = "Enrollm")
pst4 <- my_wrapper("https://www4.stat.ncsu.edu/~online/datasets/PST01d.csv", value = "Enrollm")

# Combine step by step
pst12 <- combine_wrapper_results(pst1, pst2)
pst34 <- combine_wrapper_results(pst3, pst4)
pst_combined <- combine_wrapper_results(pst12, pst34)
```

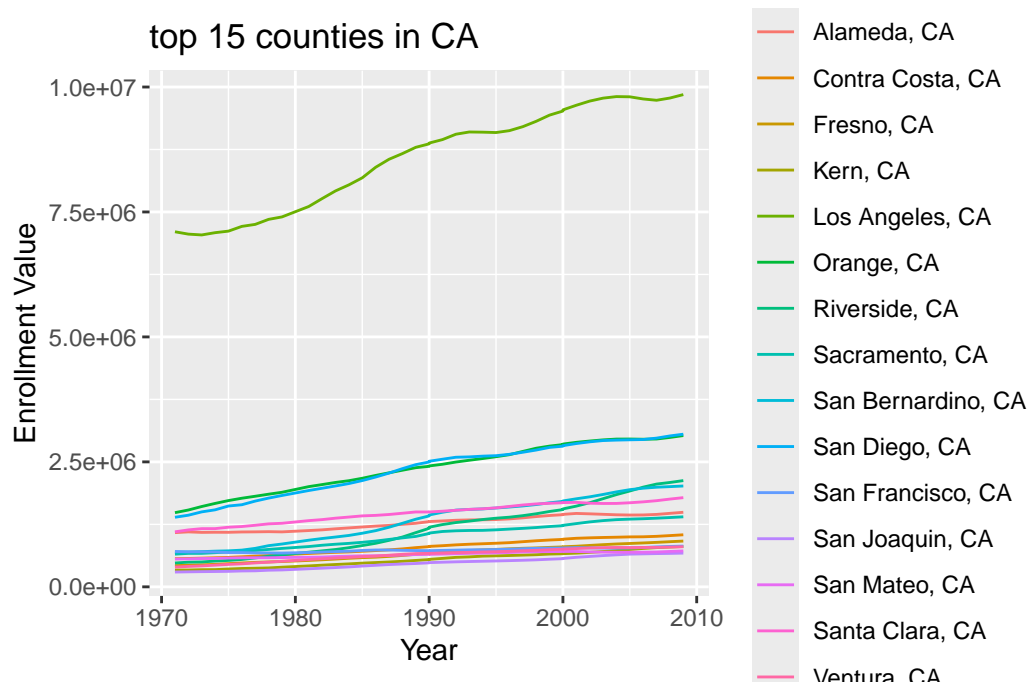
Question 7: State Plot for PST Data

```
plot(pst_combined$noncounty, var_name = "Enrollment Value")
```



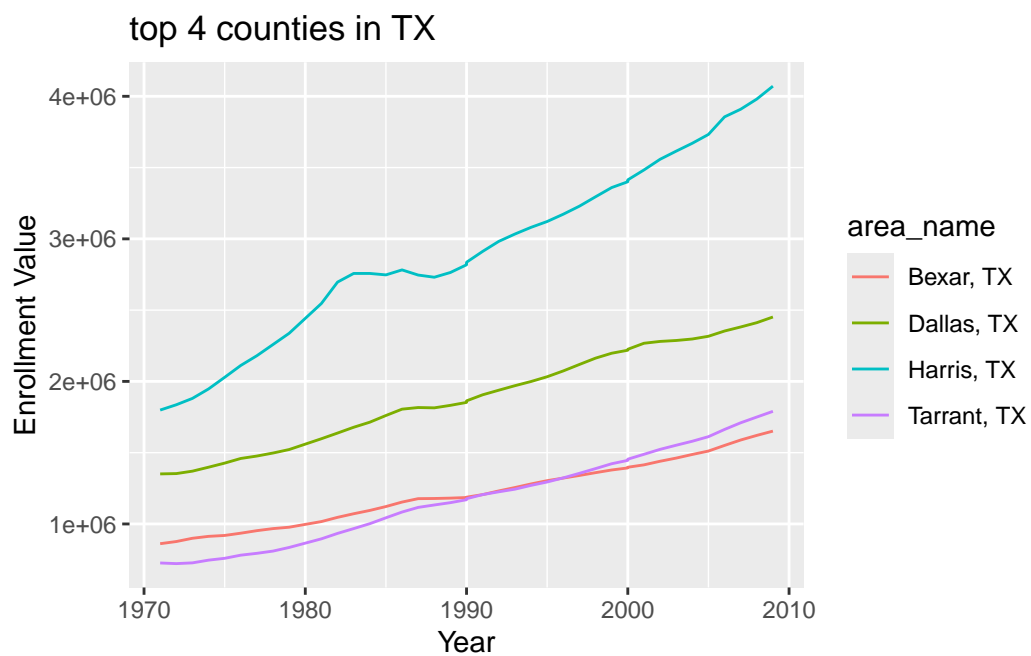
Question 8: County Plots for PST Data

```
# CA, top 15
plot(pst_combined$county, var_name = "Enrollment Value", state = "CA", top_or_bottom = "top")
```

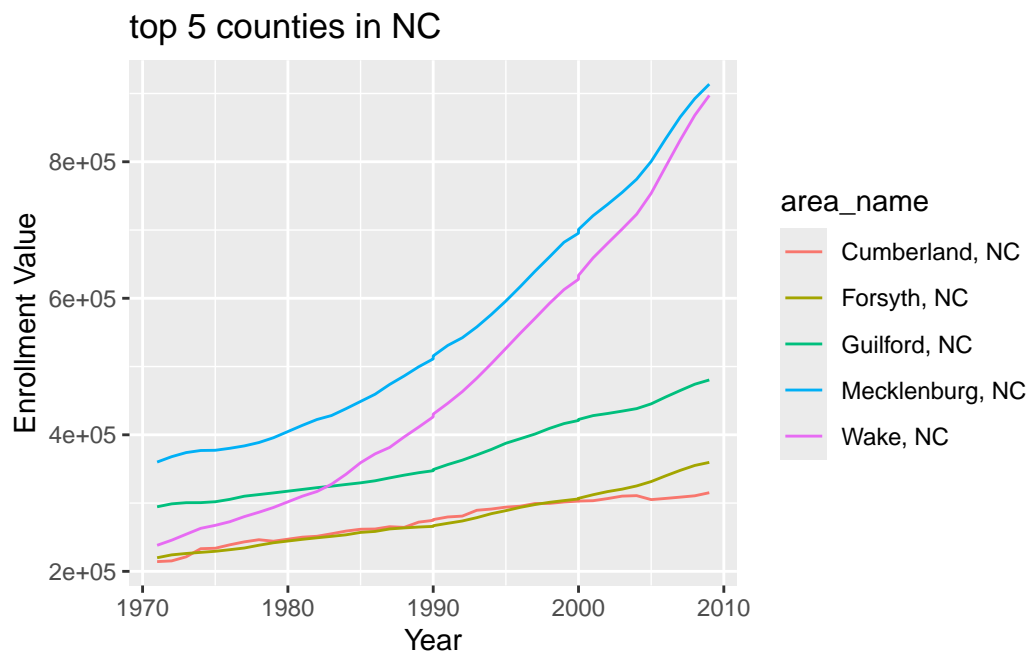


TX, top 4

```
plot(pst_combined$county, var_name = "Enrollment Value", state = "TX", top_or_bottom = "top")
```



```
# Default
plot(pst_combined$county, var_name = "Enrollment Value")
```



```
# NY, top 10
plot(pst_combined$county, var_name = "Enrollment Value", state = "NY", top_or_bottom = "top")
```

