

Data Visualization for Environmental Epidemiology with ggplot2

Mastering Presentation Grade Figures

Alexandra E. Larsen, Alison K. Krajewski, Lauren H. Wyatt

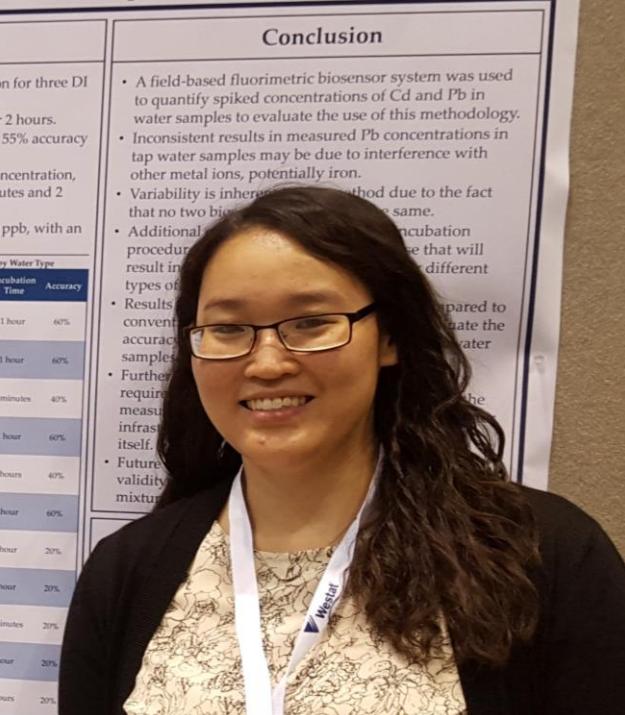
ISEE 2020 Pre-Conference Workshop
January 28, 2021

Disclaimer: The views expressed in this workshop are those of the authors and do not necessarily reflect the views or policies of the USEPA. Mention of trade names or commercial products does not constitute endorsement or recommendation for use.

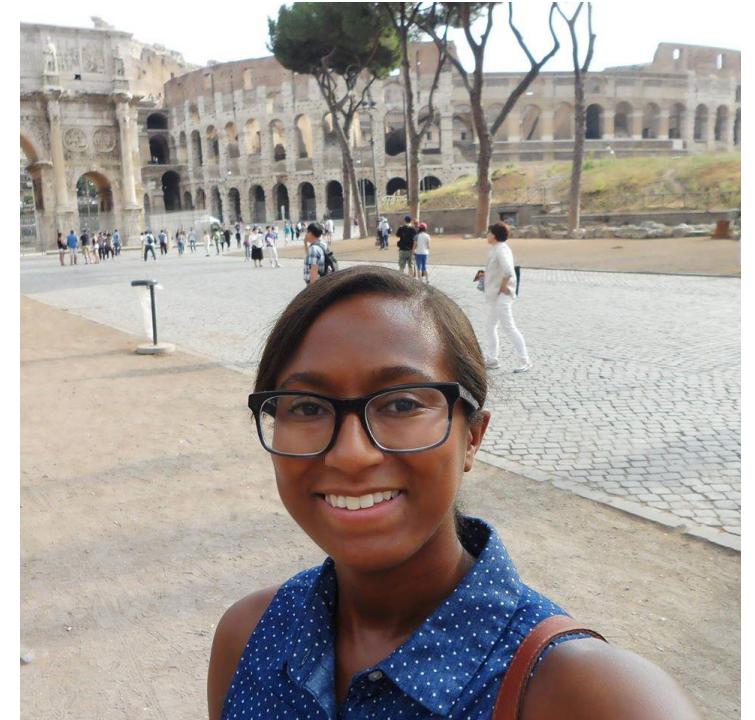
Welcome!



Alexandra Larsen
US EPA | ORD |
CPHEA| CPAD



Alison Krajewski
US EPA | ORD | CPHEA| HEEAD |
IHAB



Lauren Wyatt
US EPA | ORD | CPHEA | PHITD |
CRB

Workshop Format

- ~ Three 50-minute sessions with 10-minute breaks in between.
- Working through examples in RStudio throughout.
- Code and slides are on GitHub:
<https://github.com/USEPA/data-viz-ggplot2>
- Please feel free to ask questions!

Topics for Today

Section 0: What makes a “good” figure?

Section 1: The Basics (data formating and ggplots)

Section 2: Customization (scales, colors, theme and facets)

Section 3: Complex plots (maps, examples from the literature)

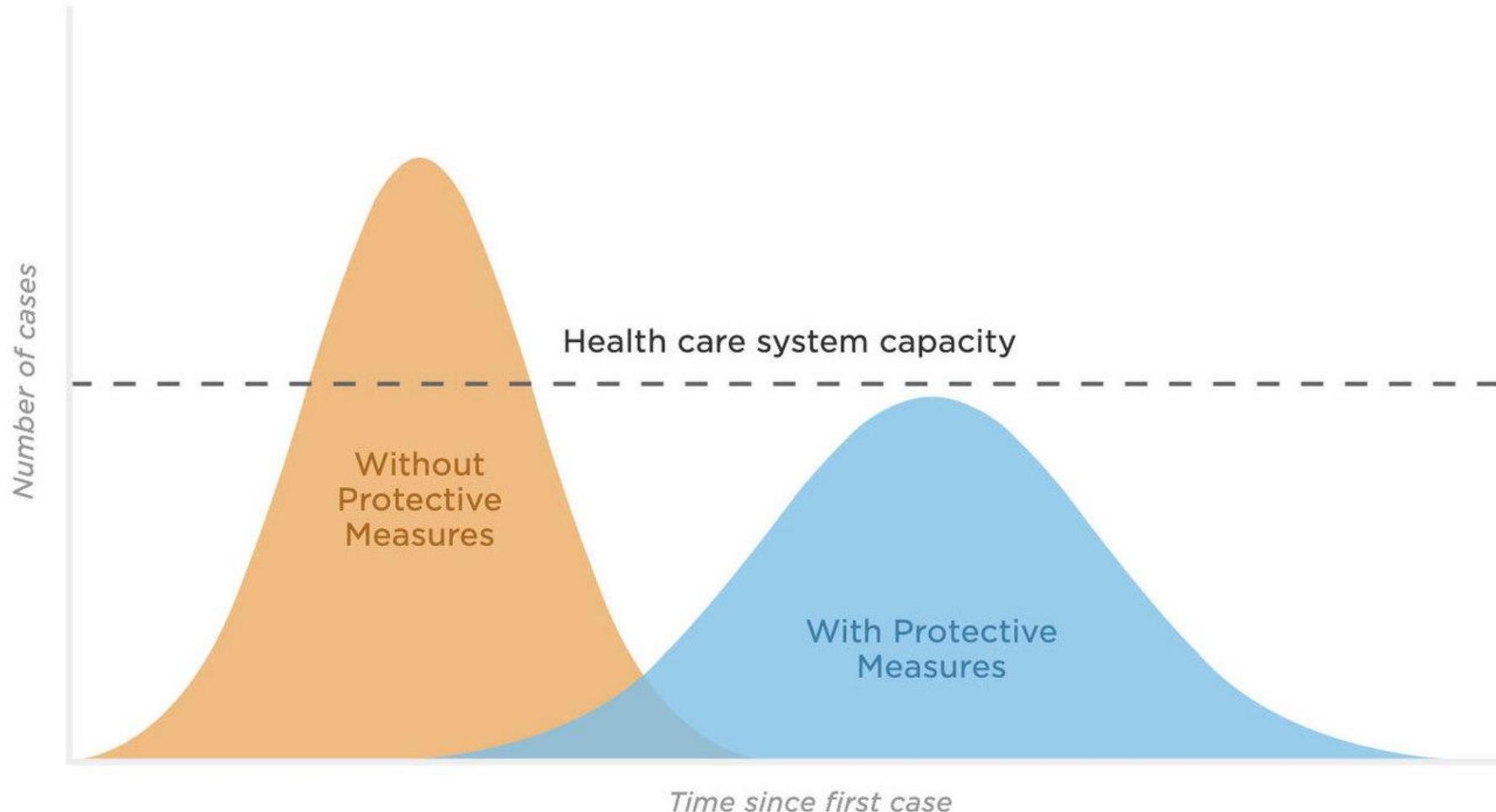
Section 0

What makes a “good” figure?

Effective Plots

- Graph type is appropriate for the data/results;
- Scales are correctly formatted;
- Message is clear enough to understand in a few minutes;
- Formatting choices help deliver the message instead;
- Facilitates informed decision-making.

'Flatten the Curve'



Example from the Literature:

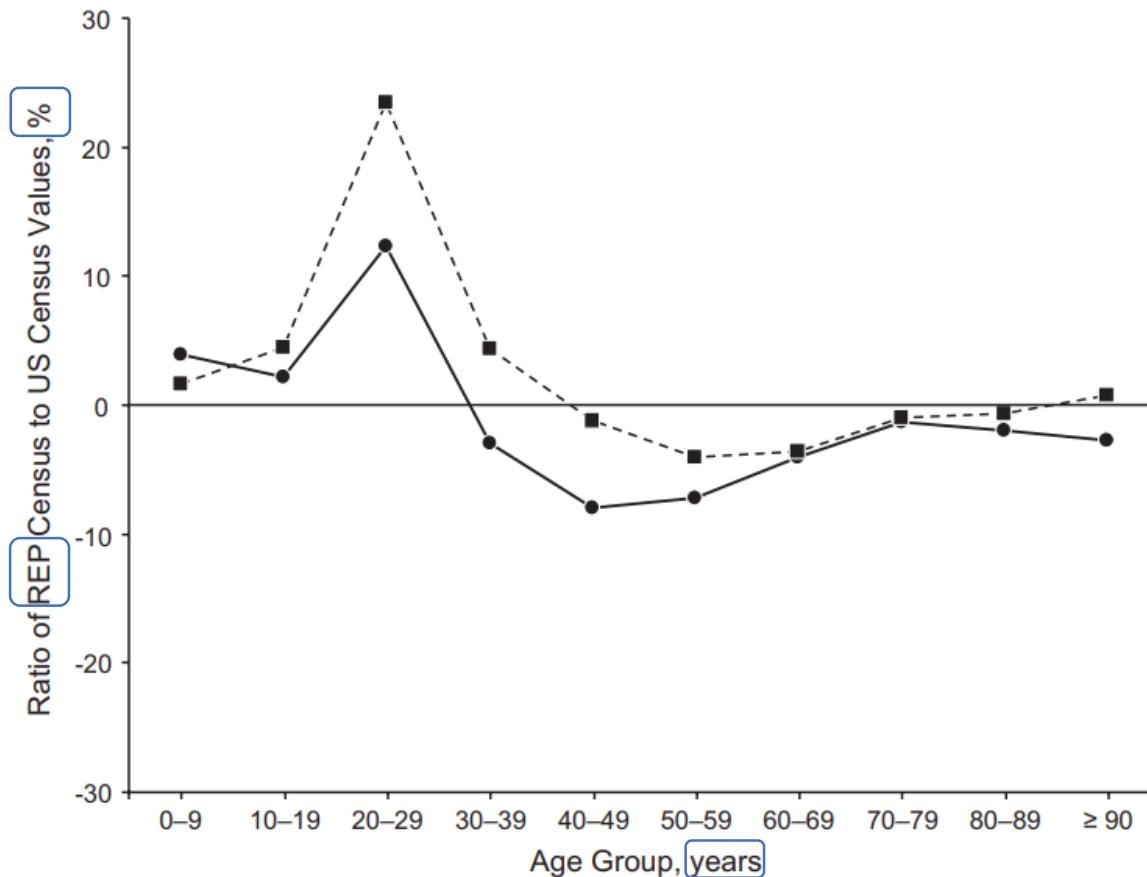
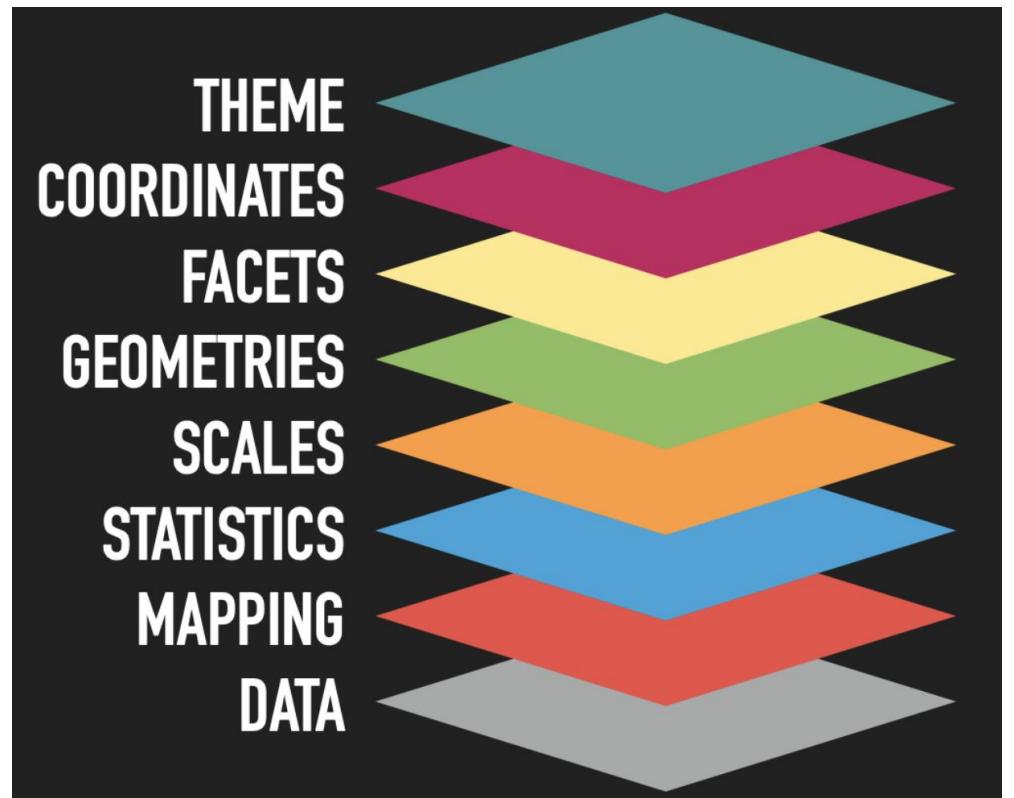


Figure 3. Age- and sex-specific capture rate by the Rochester Epidemiology Project (REP) medical records linkage system compared with US Census data (median capture rate in 1970, 1980, 1990, and 2000). Data from men (solid line, circle points) and women (dashed line, square points) are shown separately. The 0% line corresponds to perfect agreement between the system and the US Census. Values plotted above the 0% line indicate that the REP counted more persons than the US Census; values plotted below the 0% line signify that the REP counted fewer persons than the US Census.

ggplot2

- Open-source R package for making scientific graphics.
- “**ggplot**” = Grammar of Graphics (Leland Wilkinson).
- Link variables in any data set to any graphic component.
- Allows for more flexibility and customization than built-in plots.



Section 1

The Basics: Data formatting and ggplots.

Data Formatting

Reading and Writing Data Sets in R

Can read in several formats – .csv, .xlsx, .txt – but .csv is the easiest/most user-friendly.

read.csv()

- Arguments: Name and location of the csv file you want to read in, whether there are headers...
- Save to a variable

write.csv()

- Arguments: the data frame you want to write out, the name and location of the output file , whether to include headers, etc.

[Go to Example...](#)

Data Manipulation with `tidyverse`

- Tidyr is a library for data cleaning.
- Includes intuitive functions for sub-setting, pivoting, etc.
- Both ggplot2 and tidyr use the **pipe operator**:

`f(object, args) <-> object %>% f(args)`

- Allows for streamlined code that is easy to read.
- The pipe operator is “`%>%`” in tidyr and “`+`” in ggplot2

[Go to Example...](#)

Long vs. Wide Data

- Often need to convert between long and wide data.

Long (narrow, stacked):

One columns contain the values, and the other contains the description. Going from wide to long lengthens the data, increases the number of rows.

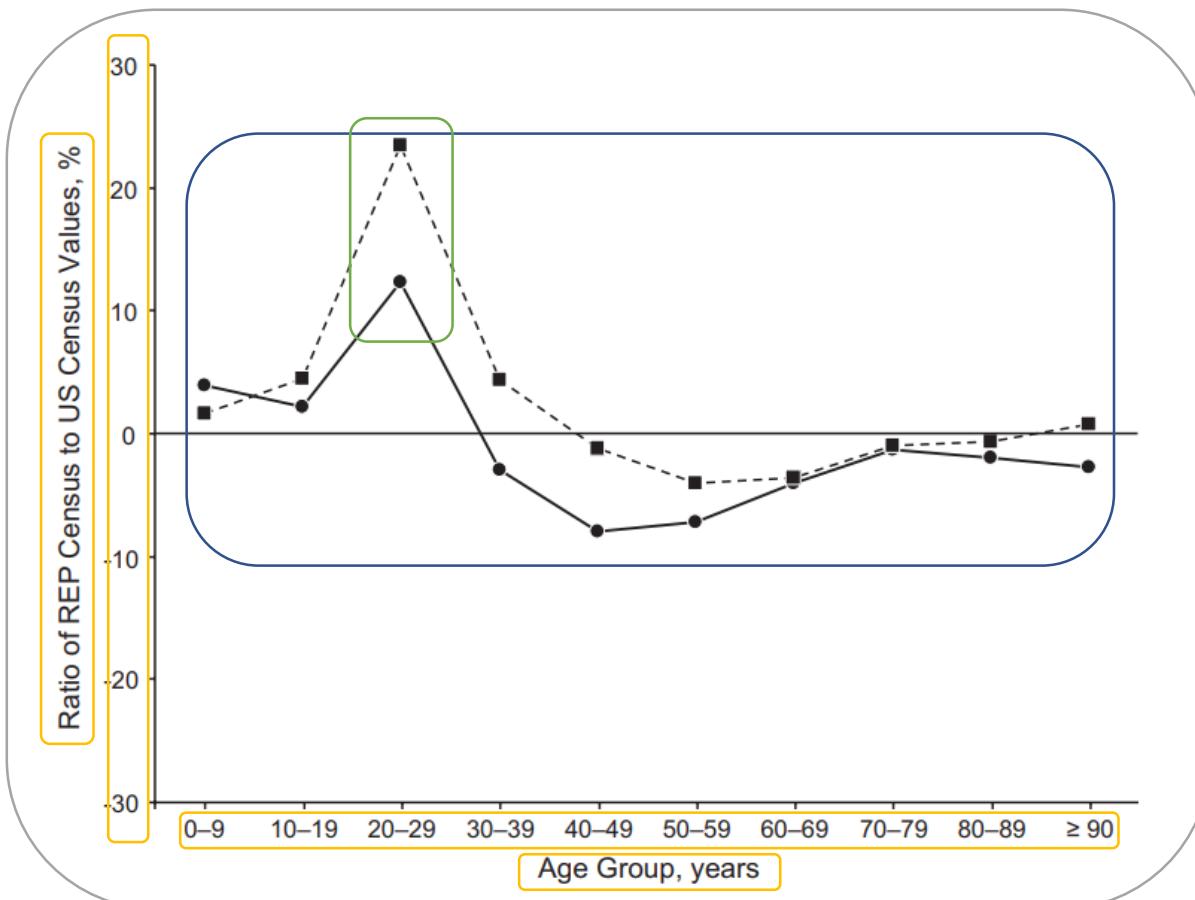
Wide:

Each variable in the data has its' own column. Going from long to wide shortens the data. *Very* useful format for plotting.

Go to Example...

ggplots

ggplot Structure



```
ggplot(data = dat.csv) +
```

```
  geom_point(aes(x, y, shape, lty)) +  
  geom_hline(y = 0) +
```

```
  theme_bw() +
```

```
  ylim(c(-30, 30)) +
```

```
  ggttitle(null)
```

```
  xlab("Age Group, years") +
```

```
  ylab("Ratio of REP Census...") +
```

```
  scale_line_manual(guide = F) +
```

```
  scale_shape_manual(guide = F)
```

Geometries
(Plot type)

Themes
(background, style)

Labels, Axes
(breaks, limits etc.)

Scales, Legends
(Controls the aesthetics)

Creating a ggplot

- Starts with a call to **ggplot()**, add graphic components in layers.
 - E.g. plot type, scales, theme, etc.
- In each component, specify data set and map aesthetics to the data.
 - Aesthetic mapping: links variables in data to graphic component (e.g. x, y, color, shape, etc.)
 - Data sets must be a `data.frame()` type, not `numeric()`, `matrix()`, etc.

Plot Types and Geometries (geoms)

`geom_{plot-type}()`

- Controls the points, lines, bars, polygons that go into making each type of plot
- Commonly used:
 - `geom_point` (scatterplot)
 - `geom_histogram` (histogram)
 - `geom_boxplot` (box plot)
 - `geom_map` (maps)
- Comprehensive list: [ggplot2 Reference page](#)

Go to Example...

Saving your ggplot Figure

ggsave()

- Call after making your figure.
- Defaults to saving the last plot created.
- Can specify figure dimensions, resolution (dpi), format.
- Several output formats (png, jpg, tiff, etc.)

grid.arrange(), multiplot(), ggarrange()

- Combine multiple plots into one figure (alternative to facet)

Go to example...

Section 2

Customization: Scales, colors, theme, and facets.

Scales

Scales

- What are they?
 - Controls the mapping from data to aesthetics
 - Everything inside the `aes()` will have scales
 - Each scale is a function from a region in data space (domain of scale) to a region in aesthetic space (range of scale)

Modifying Scales

- Scales are made up of three pieces separated by an underscore (_)
- Scale + name of aesthetic (e.g. colour, shape, or x/y) + the name of scale (e.g. discrete, continuous, brewer)
 - Examples:
 - `scale_x_continuous()`
 - `scale_color_discrete()`

Labels

- x and y axes labels will default to the variable name
- Modify the scales to change the label
 - `scale_x_continuous(name="Label Name")`
 - `ylab("Label Name")`
- Labels can include superscript, subscript, and mathematical expressions

Breaks

- Breaks control which values appear as tick marks on axes and keys on legends
- Breaks can be set on continuous or categorical scale
 - Used for labels, colors

Legends

- Legends can display multiple aesthetics (color, shape), from multiple layers
- Symbols displayed in legend varies based on the geom() used in layer
- Legends have more details that can be manipulated
 - Vertically or horizontally
 - Columns or rows
 - Size

[Go to code and examples in markdown](#)

Colors

Colors

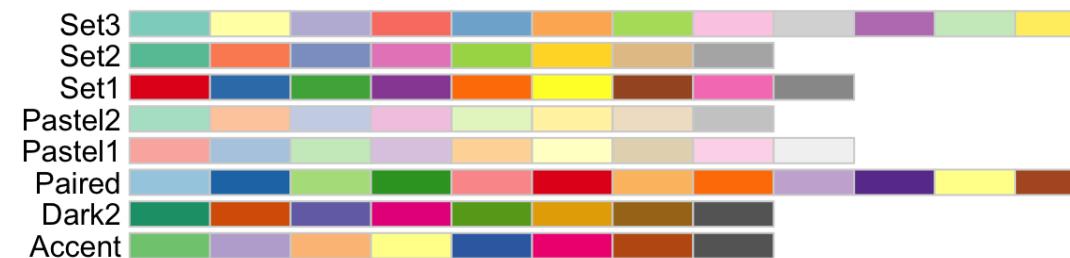
- Continuous
 - Note: for continuous color scales, keep the color scheme constant, and use a gradient scale
 - `scale_colour_gradient()` or `scale_fill_gradient()`
 - `scale_color_distiller()` or `scale_fill_distiller()`
- Discrete
 - Note: for discrete color scales, keep the color scheme color blind friendly
 - `scale_colour_brewer()`
 - Uses “ColorBrewer” colours (<https://colorbrewer2.org>)
 - `scale_colour_grey()`
 - Helpful for when needing grey-scale figures
 - `scale_colour_manual()`

ColorBrewer

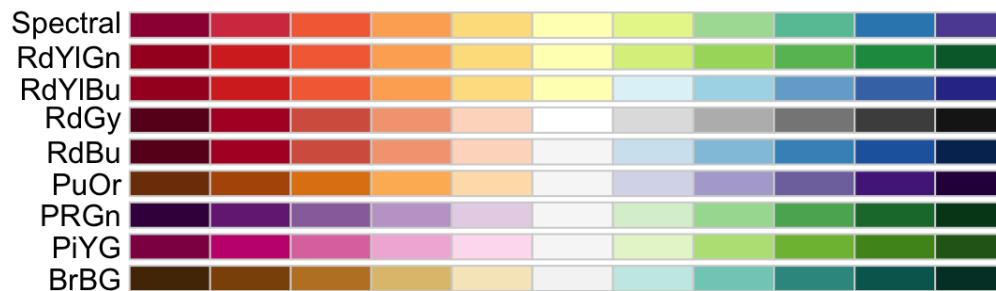
Sequential



Qualitative



Divergent



Viridis

viridis



magma



plasma



inferno



cividis



Creating Color Palettes

- Use breaks to create color palette
- Can use pre-existing color palettes
 - Specify in `scale_color_brewer()`
- Create a color palette with HEX codes
 - First, create a set of values
 - Second, specify with `scale_color_manual`

Go to code and examples in markdown

Theme

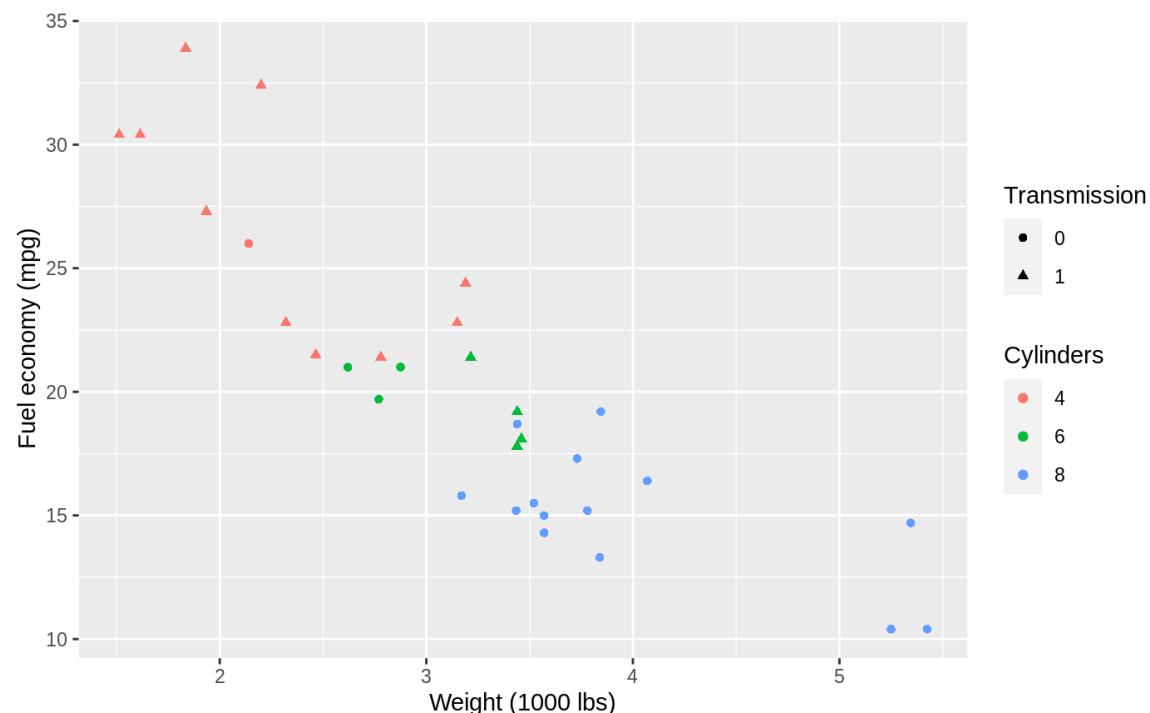
Theme

- What is a theme?
- Pre-existing themes
- Naming convention
- Customizing plot example

Theme

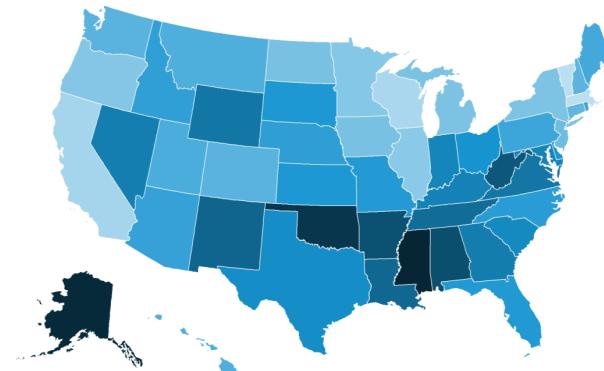
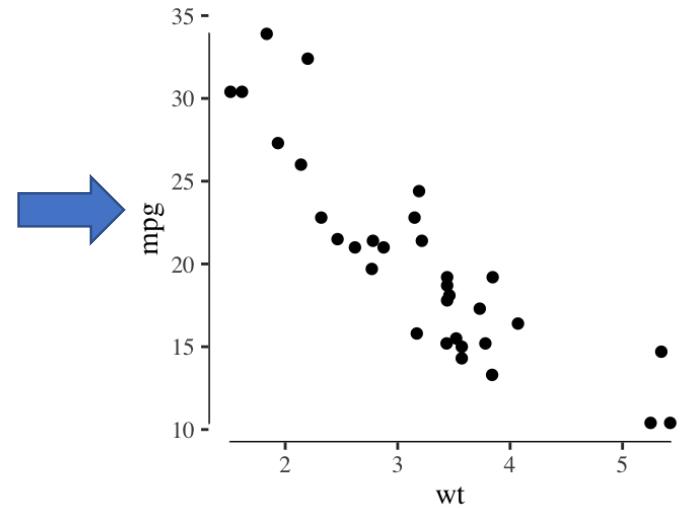
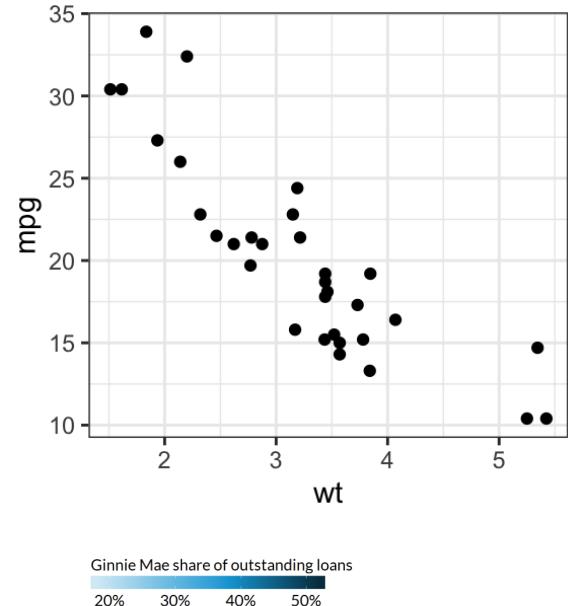
Customize non-data part of plots

- Titles, labels, fonts, background, gridlines, legends
 - Data exploration → Polished figure w/ focused message



Use theme to focus attention to data

- Choosing – decluttering
- Maps – remove plot borders



Pre-existing themes

Details

`theme_gray` The signature ggplot2 theme with a grey background and white gridlines, designed to put the data forward yet make comparisons easy.

`theme_bw` The classic dark-on-light ggplot2 theme. May work better for presentations displayed with a projector.

`theme_linedraw` A theme with only black lines of various widths on white backgrounds, reminiscent of a line drawing. Serves a purpose similar to `theme_bw`. Note that this theme has some very thin lines (<< 1 pt) which some journals may refuse.

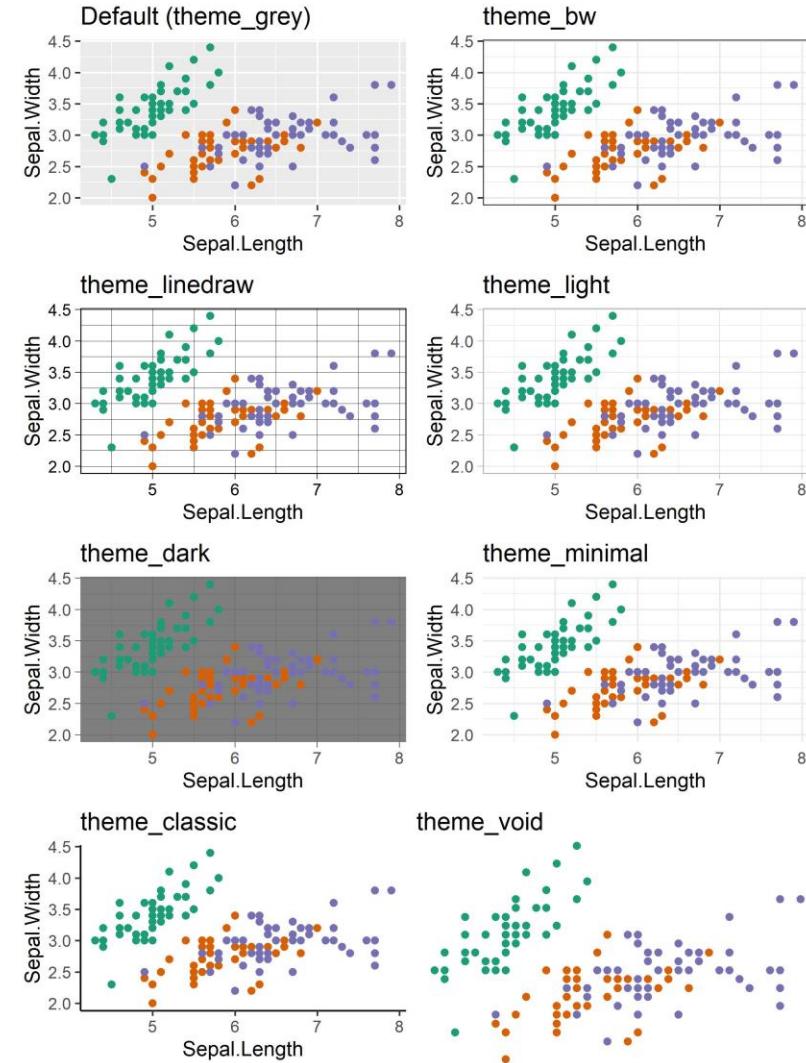
`theme_light` A theme similar to `theme_linedraw` but with light grey lines and axes, to direct more attention towards the data.

`theme_dark` The dark cousin of `theme_light`, with similar line sizes but a dark background. Useful to make thin coloured lines pop out.

`theme_minimal` A minimalistic theme with no background annotations.

`theme_classic` A classic-looking theme, with x and y axis lines and no gridlines.

`theme_void` A completely empty theme.



- Great place to start, can customize further with `theme()`

Custom theme – theme()

- Main components
 - Line elements
 - axis lines, minor and major grid lines, plot panel border, axis ticks, etc.
 - Text elements
 - plot title, axis titles, legend title and text, axis tick mark labels, etc.
 - Rectangle elements
 - plot background, panel background, legend background, etc.
- Functions
 - `element_line(color, size, linetype)`
 - `element_text(face, color, size, hjust, vjust, angle)`
 - `element_rect(fill, color, size, linetype)`

Custom theme

- Naming convention

- General

- `theme(
 axis.text = element_text(size = text_size))`



Change text size for all text

- More specific

- `theme(
 axis.text.x = element_text(size = text_size))`



Change text size for only x-axis

Theme example with scatterplot

- Data
 - Daily ozone measurements in 2 cities
 - Time series
 - Source: Los Angeles Ozone Pollution Data, 1976 (package: mlbench)
- Improvements
 - Gridlines
 - Text size
 - Rotate axis labels
 - Spacing between panels
 - Legend position

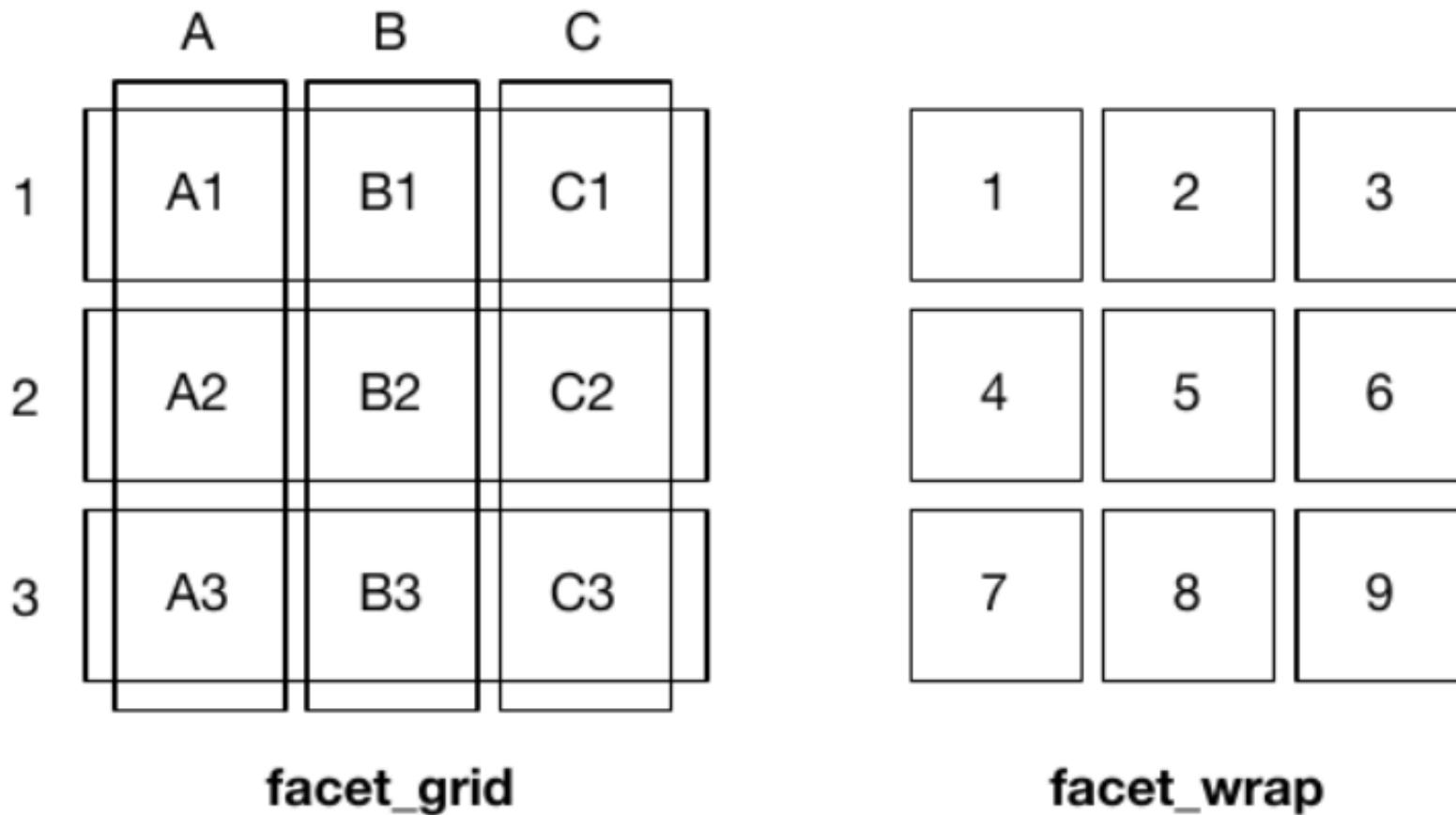


Facets

Facets

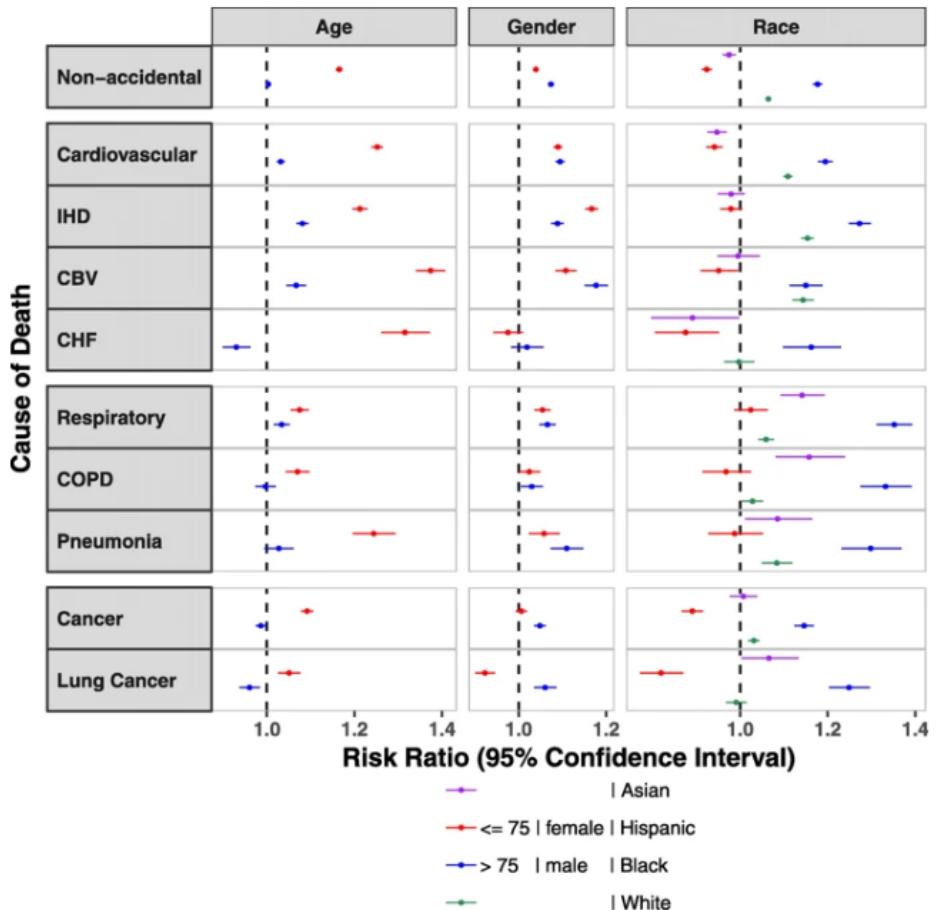
- What is it?
 - Facets generate small groupings with a different subset of the data
 - Powerful tool for exploratory data analyses
 - Readily compare patterns in different parts of data to see differences or similarities
 - Panel layout may carry meaning
 - Three types
 - `facet_null()`: a single plot
 - `facet_wrap()`: wraps a ribbon of panels
 - `facet_grid()`: produces a grid of panels defined by variables forming the rows and columns

Facet Grid vs. Facet Wrap



Example from the Literature

Fig. 2



Modification of the SES-adjusted Association between $\text{PM}_{2.5}$ and Cause-specific Mortality by Age, Sex, and Race. For each cause of death, we examined effect modification using interaction terms for age, sex and race respectively in the SES-adjusted models. Results are expressed as the risk ratio and 95% CIs per $10 \mu\text{g}/\text{m}^3$ increase in 12-month average $\text{PM}_{2.5}$. Abbreviations: IHD (Ischemic heart disease), CBV (Cerebrovascular disease), CHF (Congestive heart failure), COPD (Chronic Obstructive Pulmonary disease), SES (Socio-Economic Status), $\text{PM}_{2.5}$ (particles with aerodynamic diameters < $2.5 \mu\text{m}$). Note: Each subgroup in the death-group box follows the same order defined in the figure legend.

Wang B, Eum K, Kazzemiparkouhi F, Li C, Manjourides J, Pavlu V, and Suh H. The impact of long-term $\text{PM}_{2.5}$ exposure on specific causes of death: exposure-response curves and effect modification among 53 million US Medicare beneficiaries. Environmental Health 2020; 19: 20. doi: <https://doi.org/10.1186/s12940-020-00575-0>

Facet Grid

- Lays out figures in a grid defined by `.~` option
 - `.~a` spreads the values of variable `a` across columns, which allows for comparisons of the y-axis because the vertical scales are aligned
 - `b~.` Spreads the values of variable `b` down rows, which allows for comparisons of the x-axis because the horizontal scale are aligned
 - `a~b` spreads variable `a` across columns and variable `b` down rows
 - Usually put the variable with the highest number of levels in the columns

Facet Wrap

- Control wrap with: ncol, nrow, as.table, and dir
 - ncol and nrow controls the number of columns and rows in the arrangement
 - as.table controls the layout to be like a table, with highest values at the bottom right (as.table=TRUE) or the highest values at the top right (as.table=FALSE)
 - dir controls the direction of the wrap (horizontal or vertical)

Scales in Facets

- Position of scales can be the same in all panels (fixed) or vary between panels (free)
 - scales = “fixed”: the x and y axes are fixed across the panels
 - scales = “free_x”: the x axis is free, but the y axis is fixed
 - scales = “free_y”: the y axis is free, but the x axis is fixed
 - scales = “free”: x and y axes vary across the panel

Go to code and examples in markdown

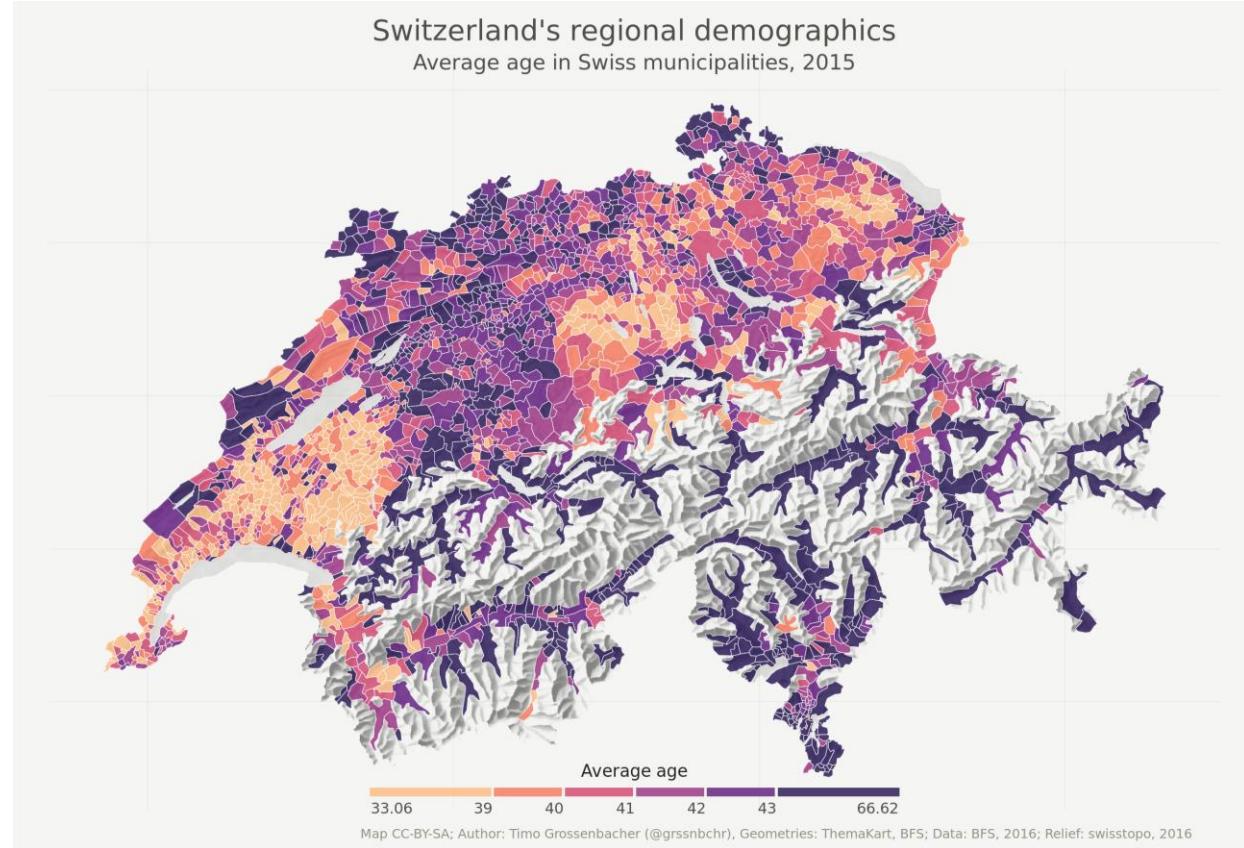
Section 3

Complex plots: Maps, examples from the literature.

Mapping

Mapping

- Fundamentals
- Layout and formatting
- Color palettes
- Create and customize examples
 - Continuous
 - Categorical

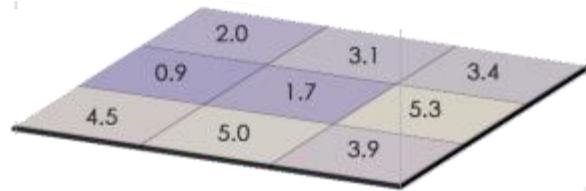


Fundamentals – spatial data types

- Vectors
 - Points (cities, landmarks)
 - Lines (roads, rivers)
 - Polygons (country borders)



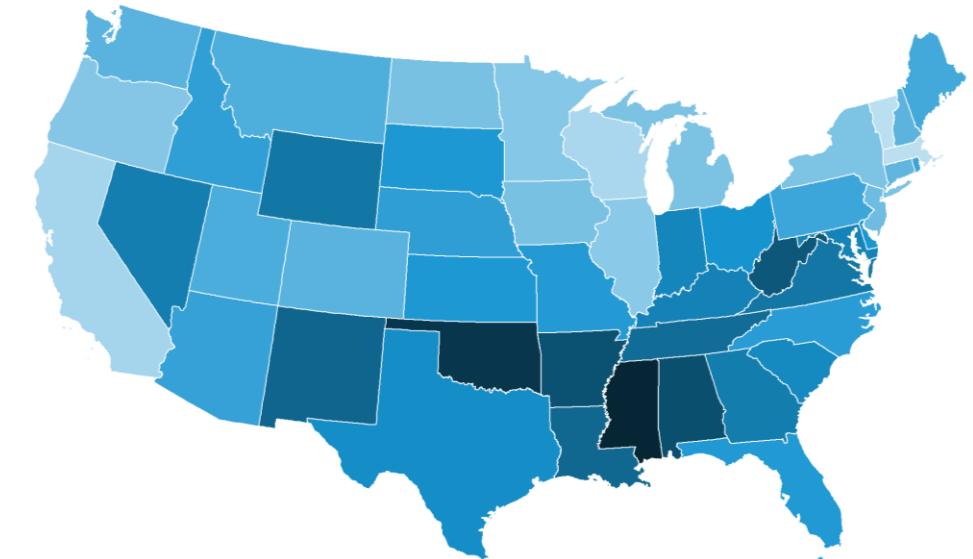
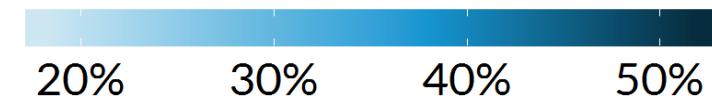
- Raster



Fundamentals – data prep → plot

sf package (spatial features)

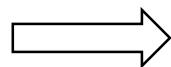
- Load data
 - Shapefiles
 - Data to display
- Join data of interest to shapefile



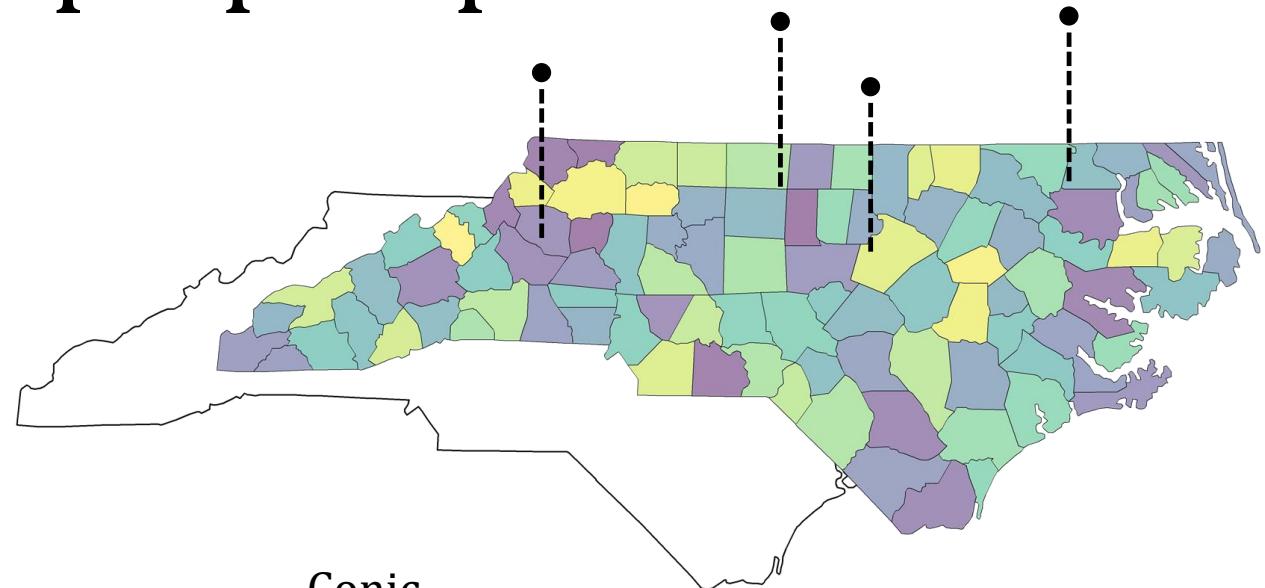
Fundamentals – data prep → plot

- Geom order is important!
- Consider map projection

Cylindric



Conic





Choropleth example

US counties

- Data
 - Spatial wildfire occurrence data for the United States
 - County-level counts of fires (2008 and 2009)
 - Source: USDA Forest Service
- Prepare data
 - Load shapefiles for state and county boundaries
 - Some packages have shapefiles for certain geographies (maps)
 - Manually load, `st_read()`
 - Combine with data of interest
- Improvements
 - Adjust focus to map, background adjustment
 - Projection
 - Color palette
 - Line colors and thickness
 - Facet

Example from the Literature

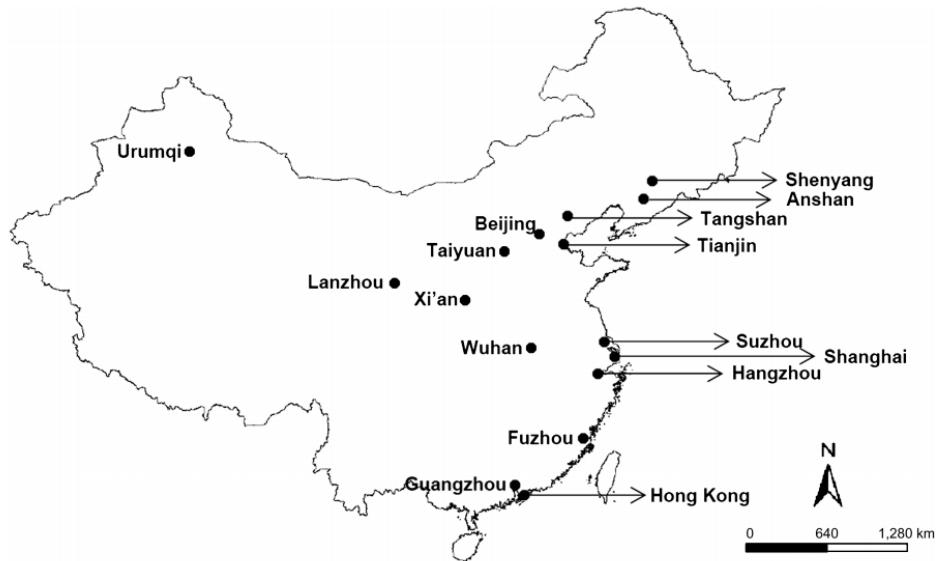


Figure 1. Location of the CAPES cities, China, 1996–2008. CAPES, China Air Pollution and Health Effects Study.

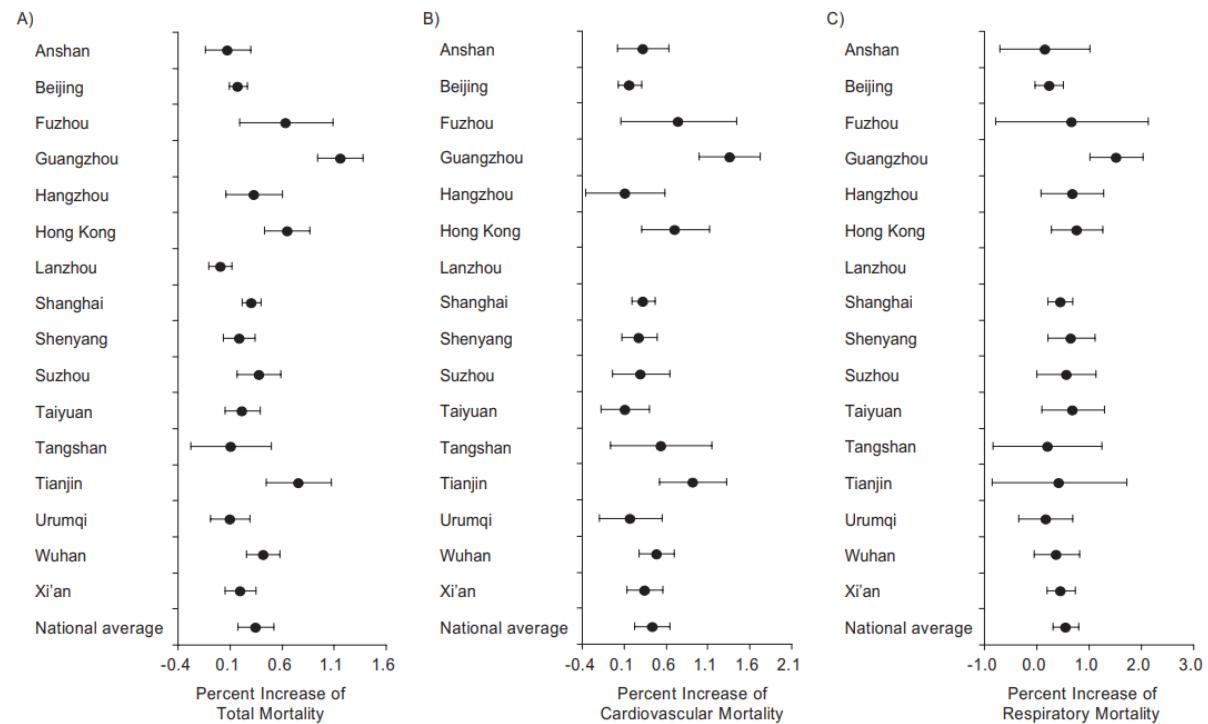


Figure 2. Percentage increase of mortality associated with a $10\text{-}\mu\text{g}/\text{m}^3$ increase of 2-day moving average PM_{10} concentrations in the CAPES cities, China, 1996–2008. Effect estimates of individual cities (mean and 95% confidence interval) and national average values (mean and 95% posterior intervals) are shown. A, total mortality; B, cardiovascular mortality; C, respiratory mortality (cause-specific mortality data were not available in Lanzhou). CAPES, China Air Pollution and Health Effects Study; PM_{10} , particulate matter with an aerodynamic diameter of less than $10\text{ }\mu\text{m}$.

Examples

Example 1: Point estimates with error bars

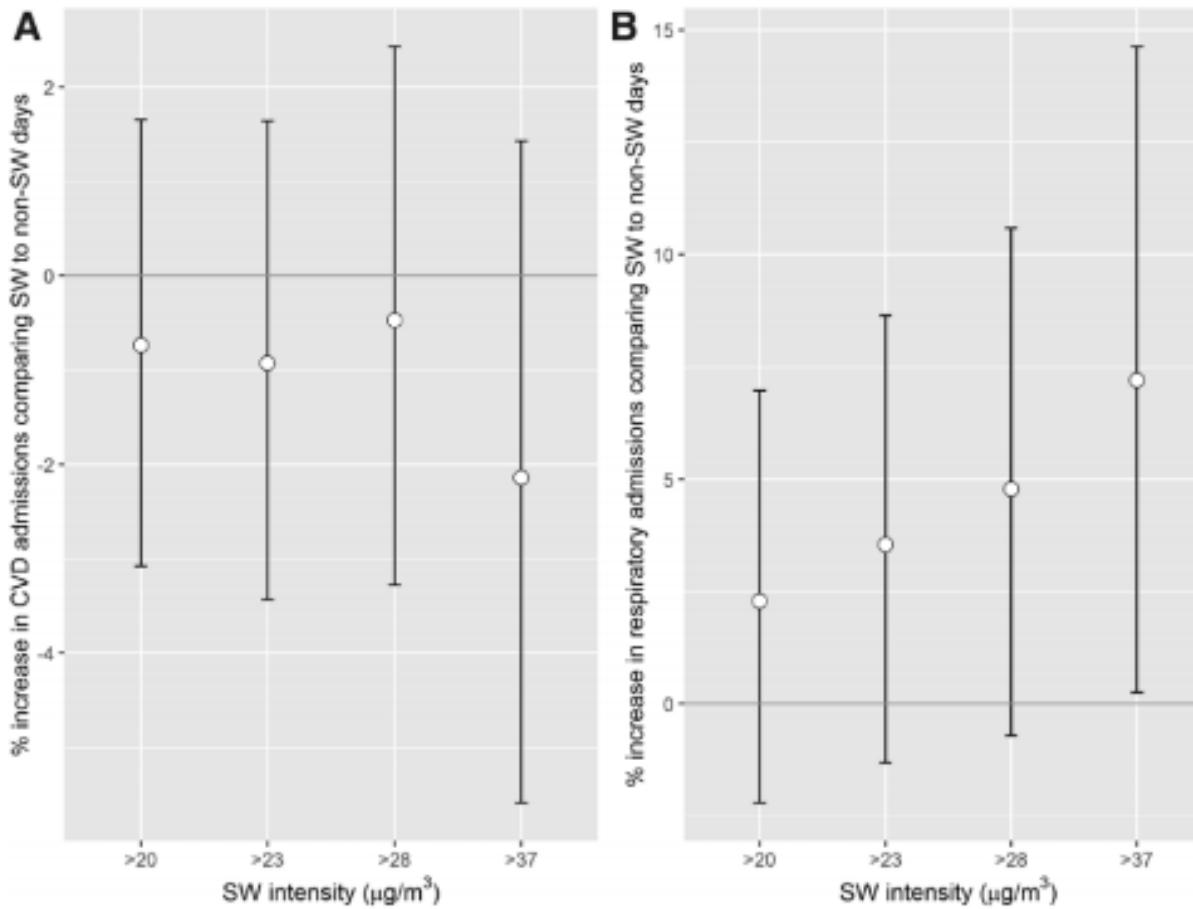
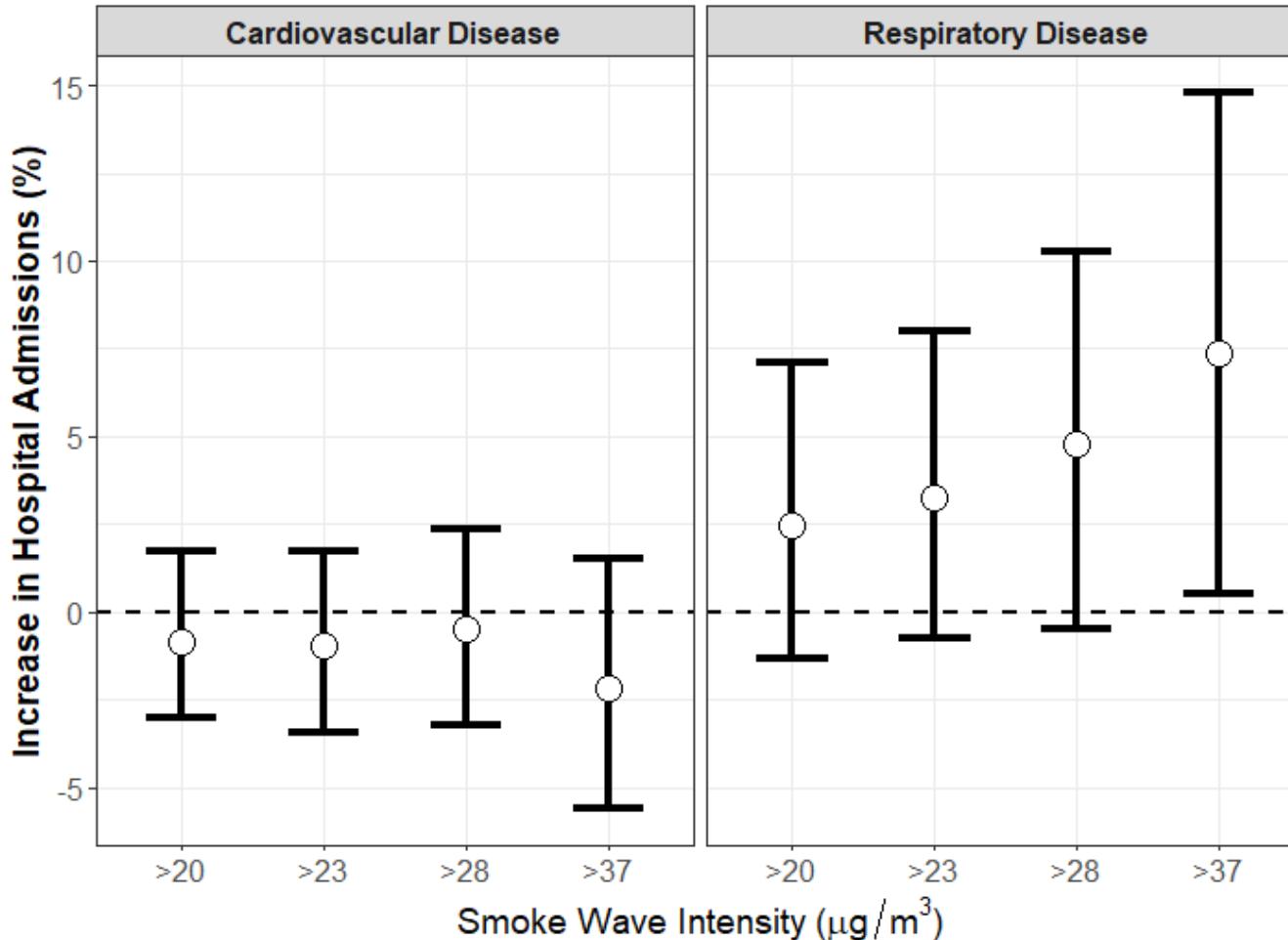


FIGURE 2. Associations between hospital admissions and exposure to smoke wave days (compared with non-smoke-wave days) for (A) cardiovascular disease and (B) respiratory disease, by different intensity (level of wildfire-specific $\text{PM}_{2.5}$) definitions of a smoke wave. SW indicates smoke wave.

Example 1: Point estimates with error bars



Modifications:

- Removed grey background
- Increased font size
- Increased size of point estimates
- Increased width of lines for error bars
- Created a facet

Example 2: Trend line with confidence band

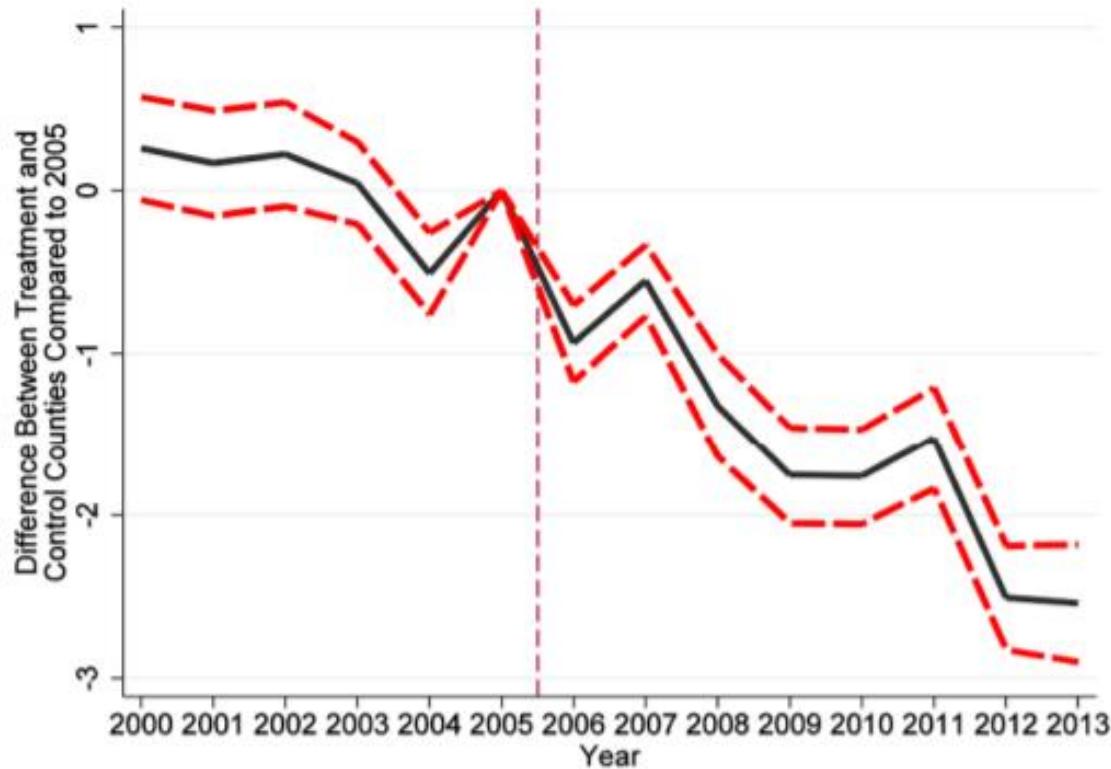
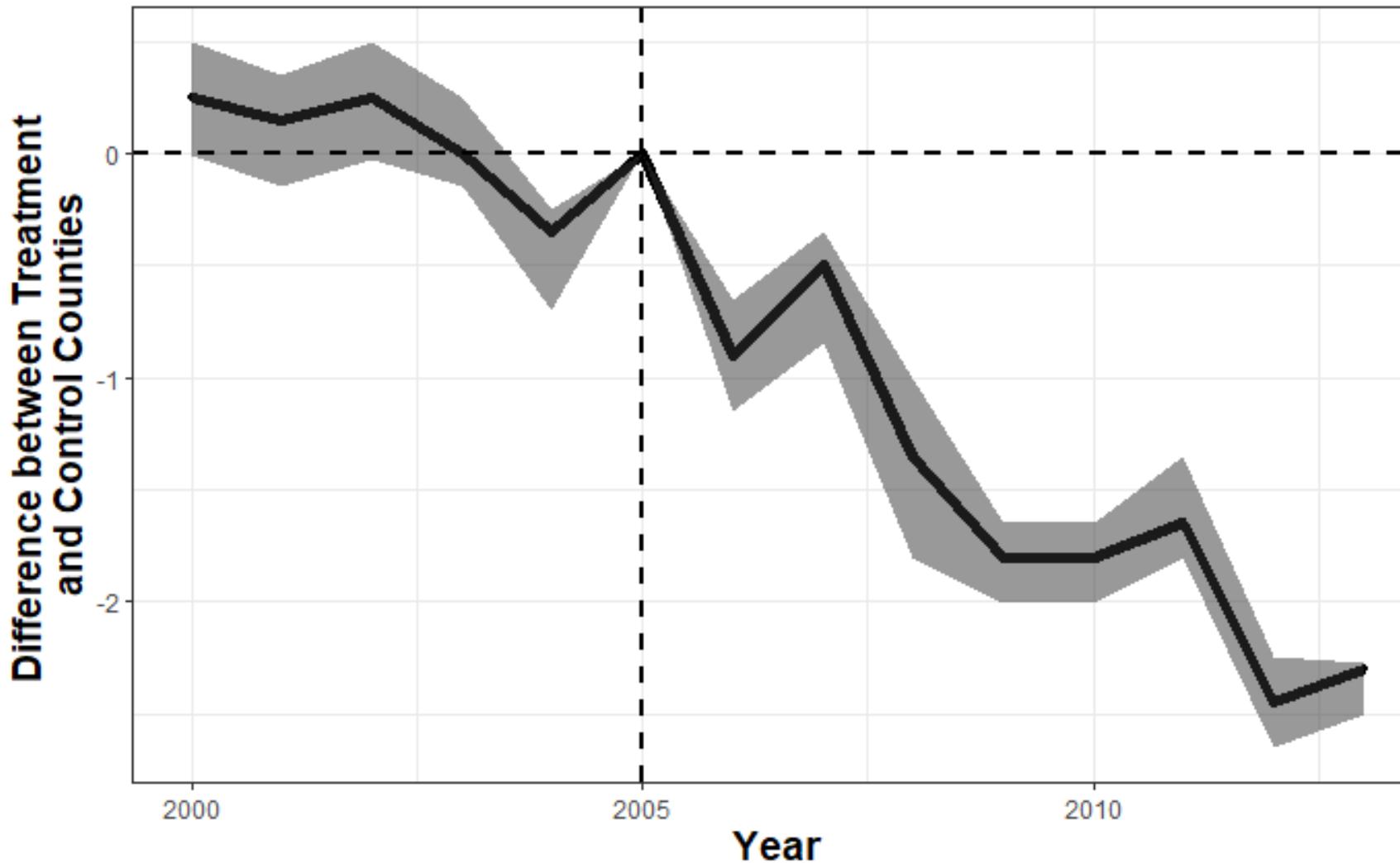


FIGURE 2. Event study estimates for PM2.5. Estimates represent the difference in PM2.5 concentrations for nonattainment counties versus attainment counties. Dotted lines represent 95% confidence interval based on standard errors clustered on county. All regressions include year dummy variables, county fixed effects, income per capita, share employed, mean temperature, maximum temperature, and precipitation. The vertical line represents when the air quality standard for PM2.5 began. 2005 is the reference category. Figure is available in color online.

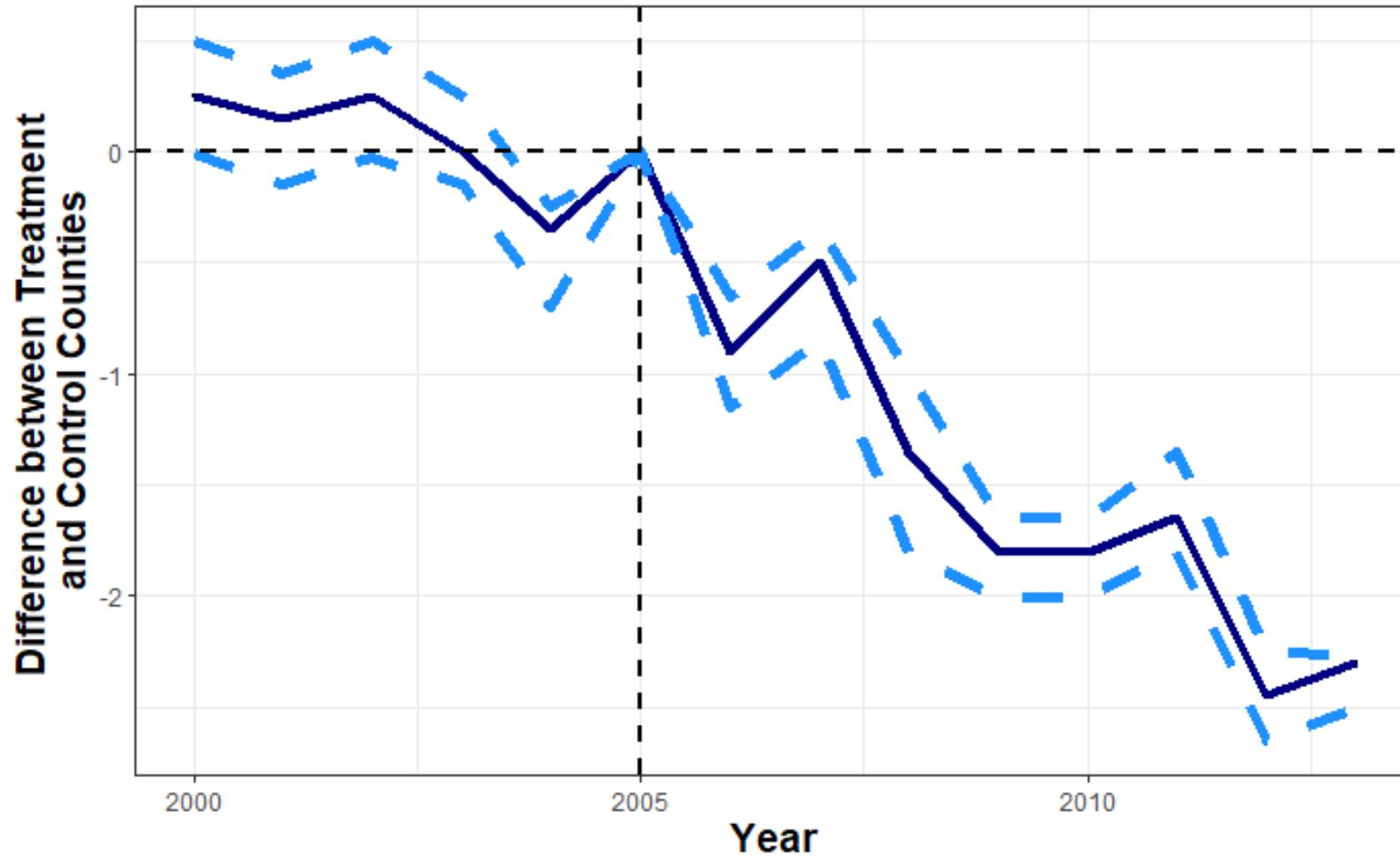
Example 2: Trend line with confidence band



Modifications:

- Increased axis label font size
- Added a confidence band
- Changed color for confidence band
- Changed the color of the reference line

Example 2: Trend line with confidence band



Modifications:

- Increased axis label font size
- Increased the size of the trend line and confidence bands
- Changed the color of confidence bands
- Change the color of reference line

Example 3: Facets with Legend Placement

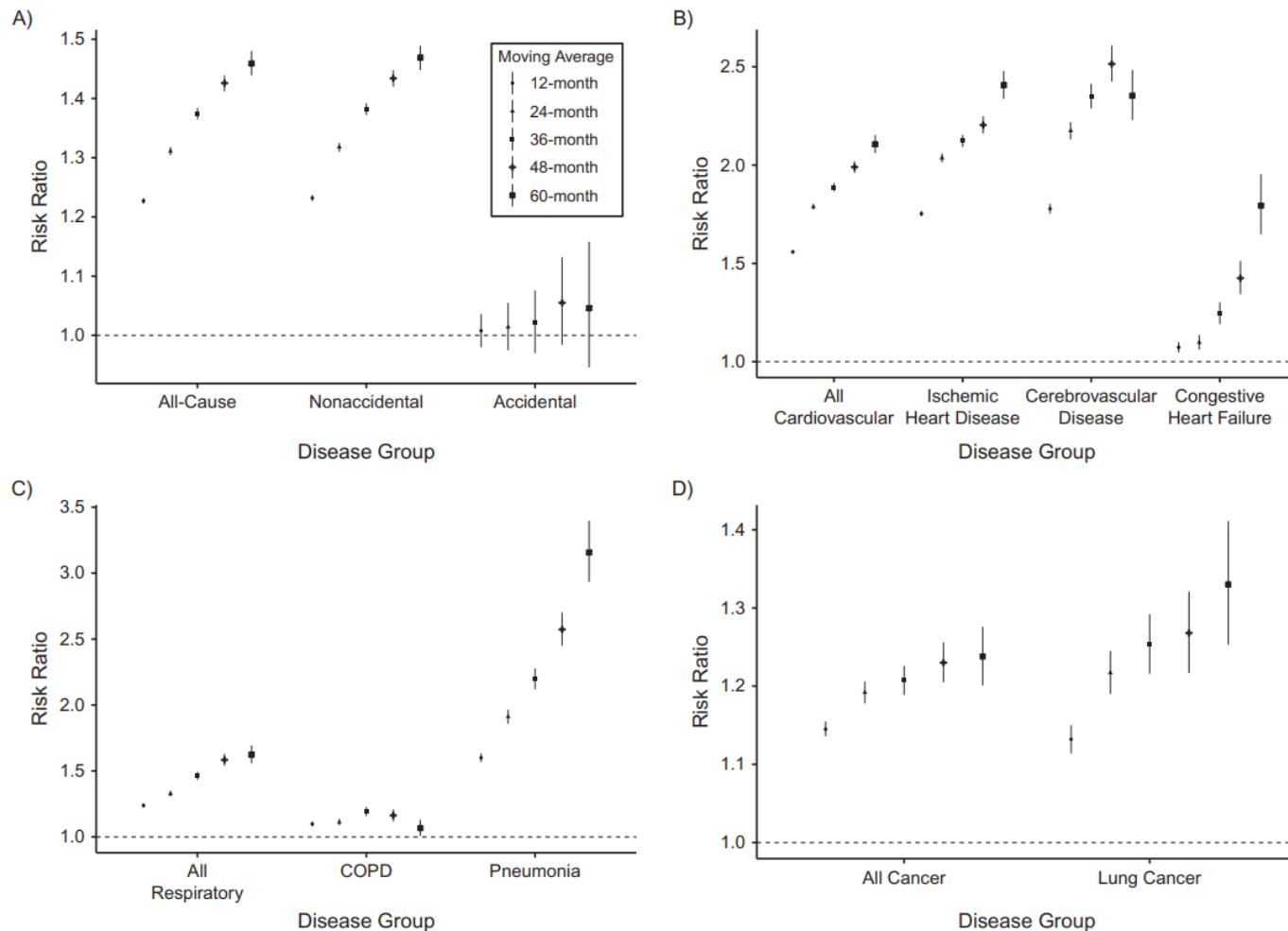
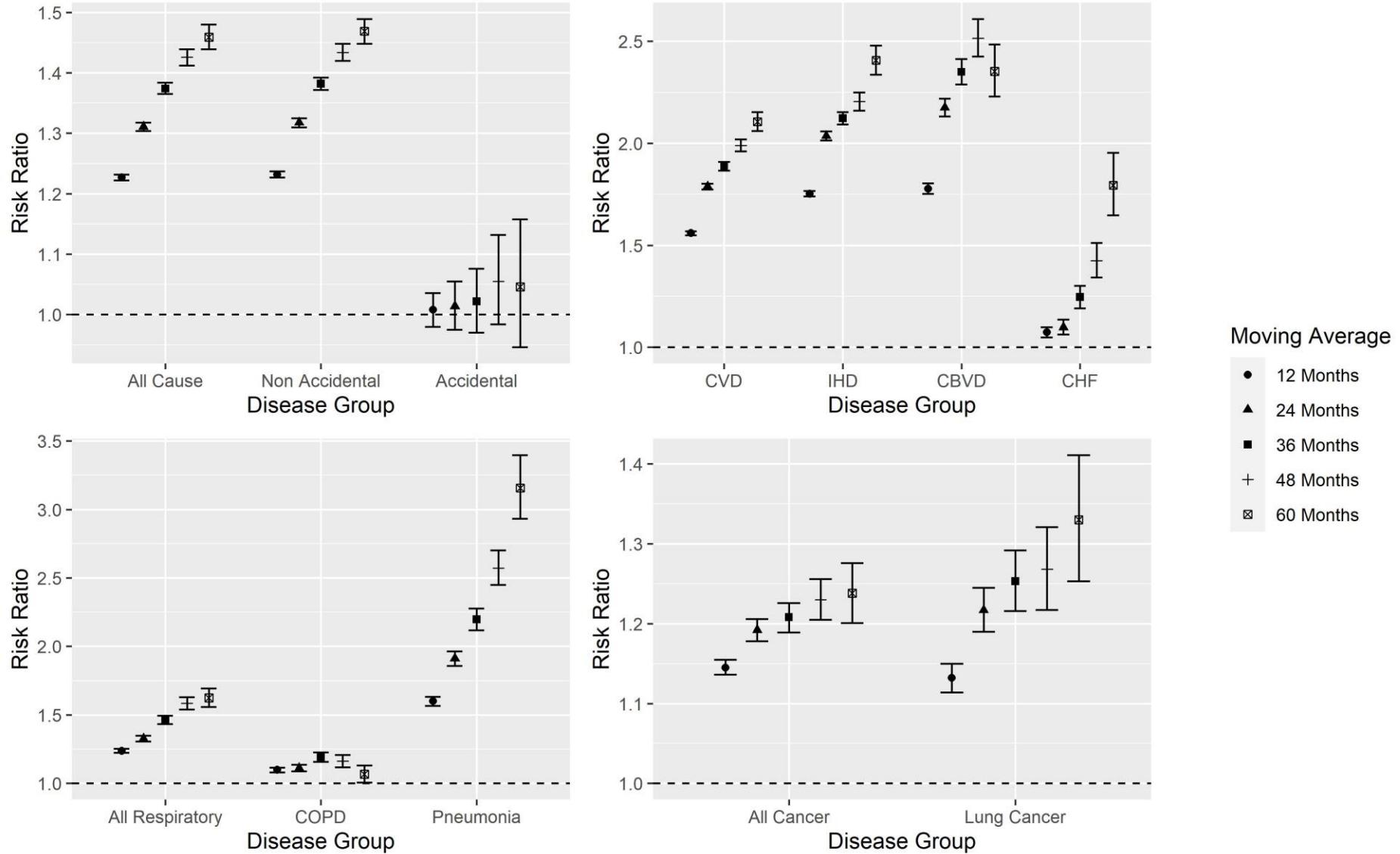


Figure 1. Risk ratios for all-cause (A), cardiovascular disease (B), respiratory disease (C), and cancer (D) mortality associated with $10\text{-}\mu\text{g}/\text{m}^3$ increases in 12- to 60-month moving average exposure to particulate matter less than or equal to $2.5\text{ }\mu\text{m}$ in aerodynamic diameter ($\text{PM}_{2.5}$) nationwide, United States, 2000–2008. Bars, 95% confidence intervals. COPD, chronic obstructive pulmonary disease.

Example 3: Facets with Legend Placement

Modifications:

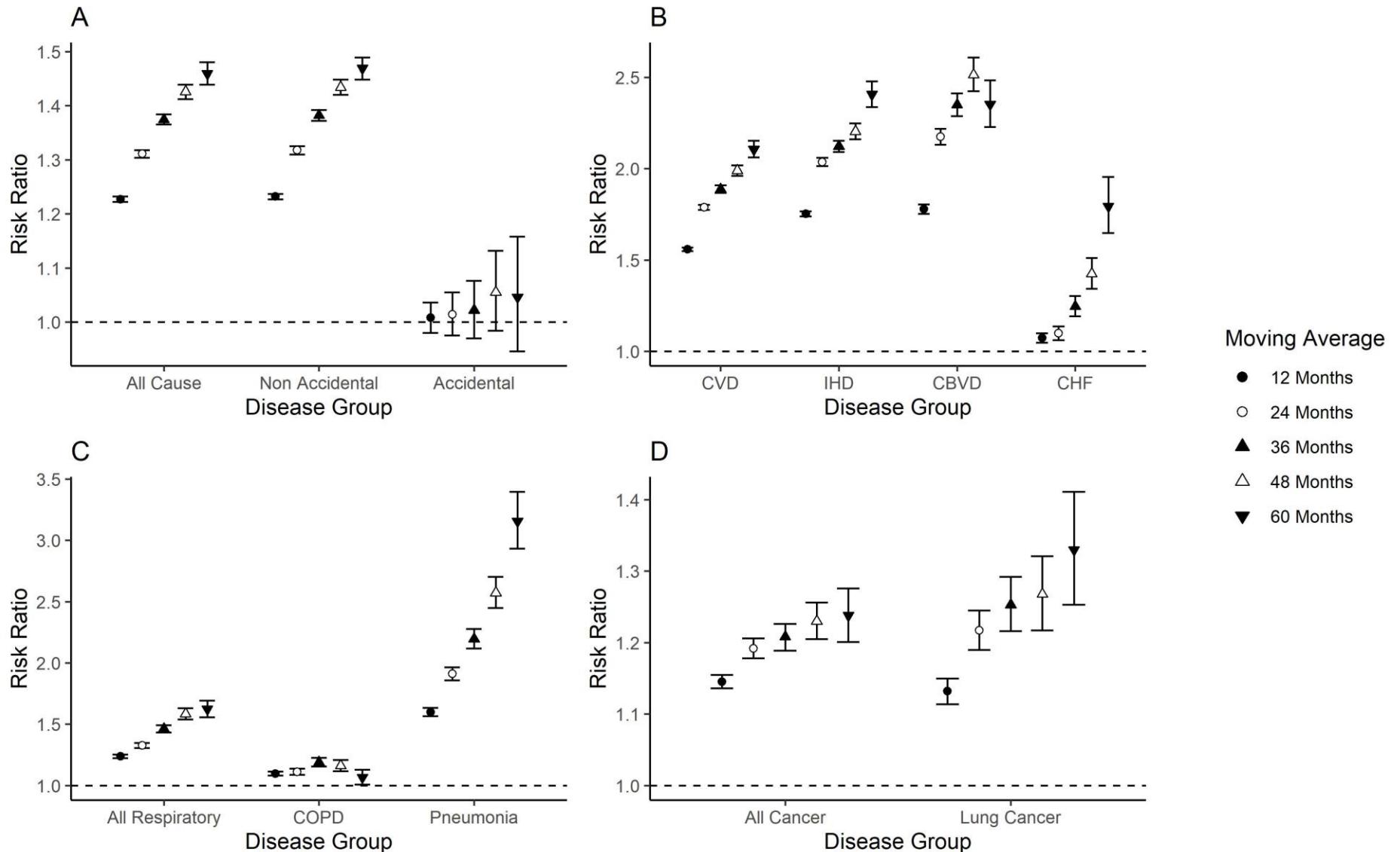
- Moved the legend position
- Increased size of the shapes of the estimates
- Increased the axis labels
- Added in background



Example 3: Facets with Legend Placement

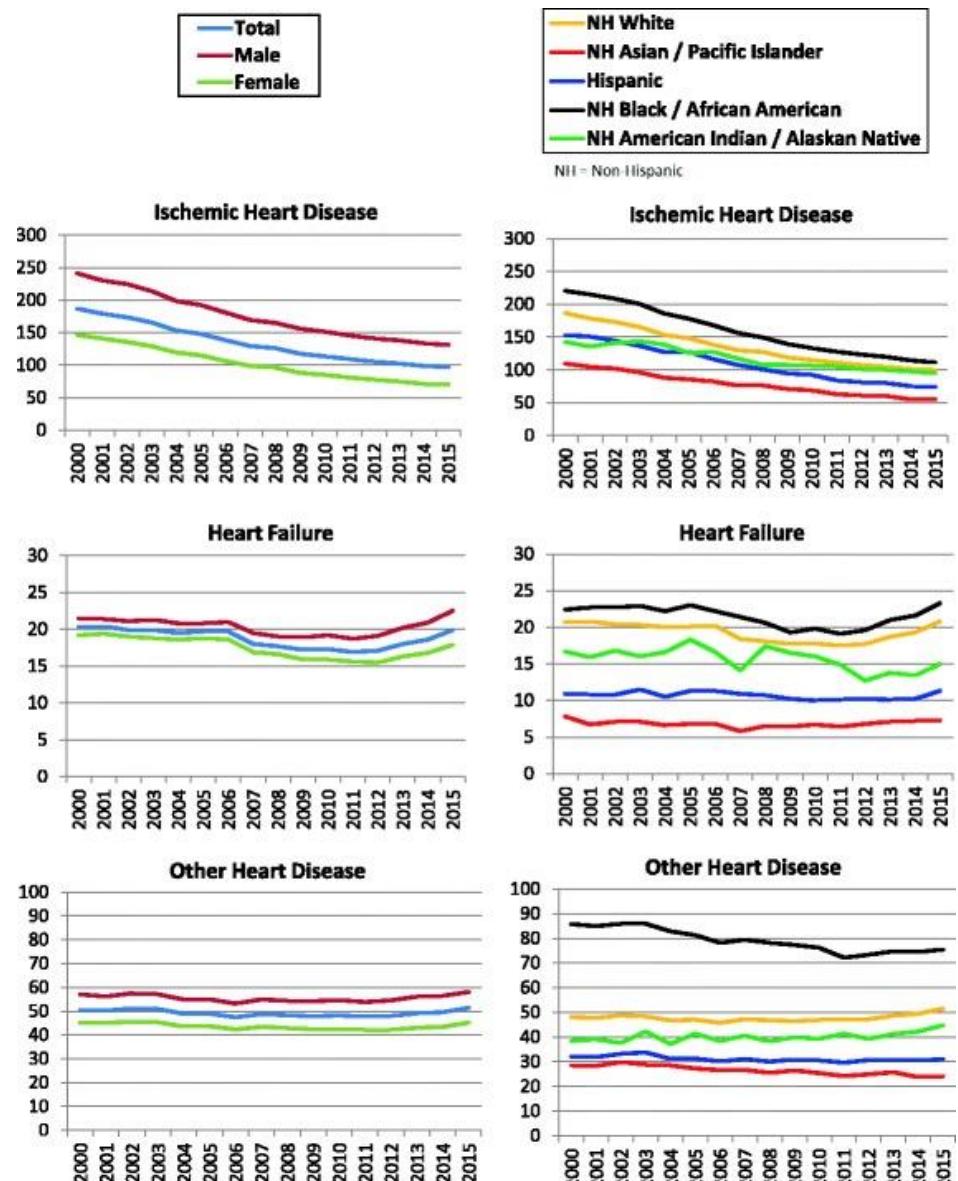
Modifications:

- Moved the legend position
- Increased size of the shapes of the estimates
- Increased the axis labels



Example 4: Facets and Colors

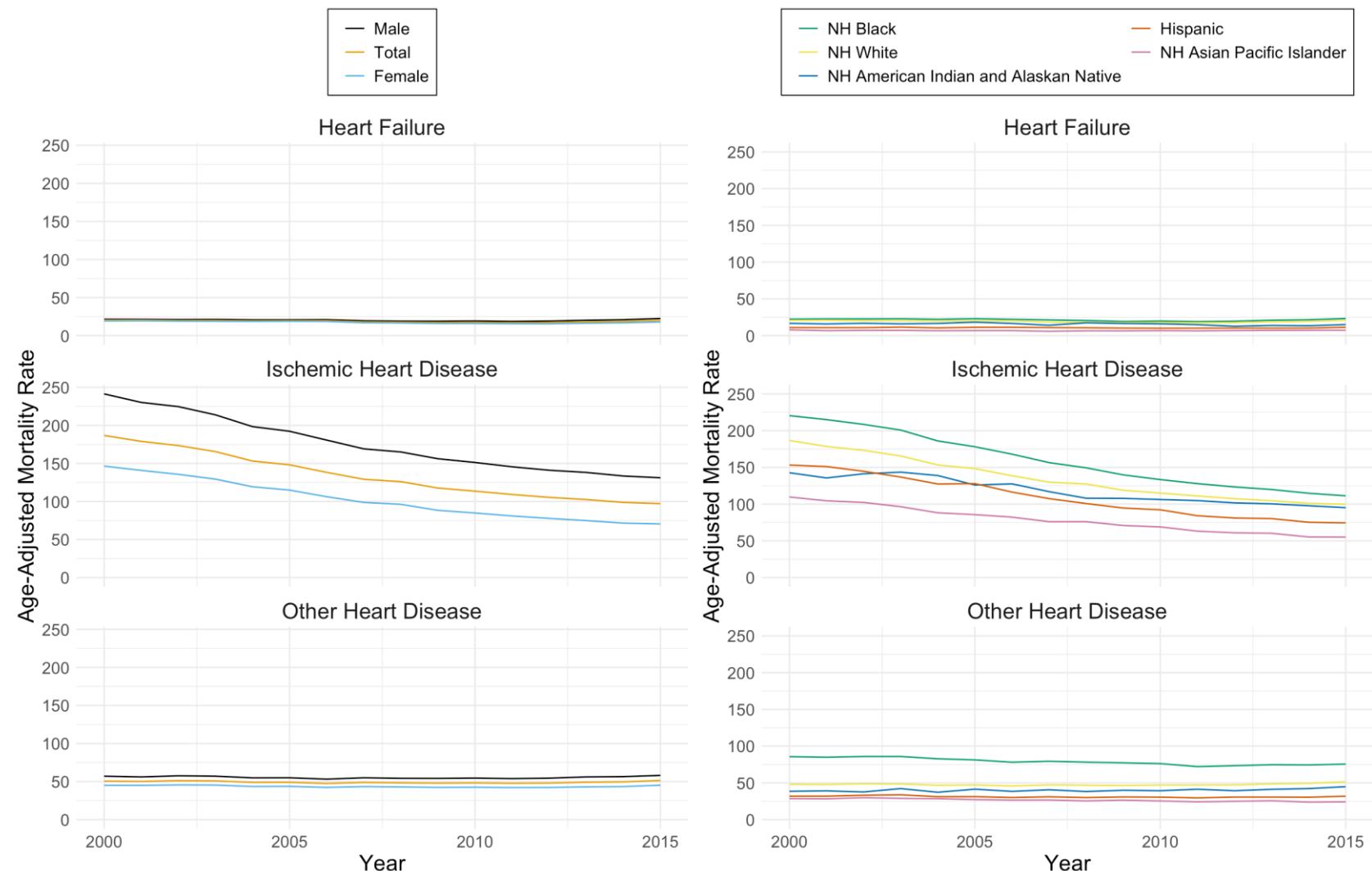
Figure 1: Age-adjusted mortality in US, 2000-2015 by sex and race-ethnicity. Legend: Total, Male, Female, NH White, NH/Asian/Pacific Islander, Hispanic, NH Black, NH American Indian/Alaskan Native (NH=Non-Hispanic)



Example 4: Facets and Colors

Modifications:

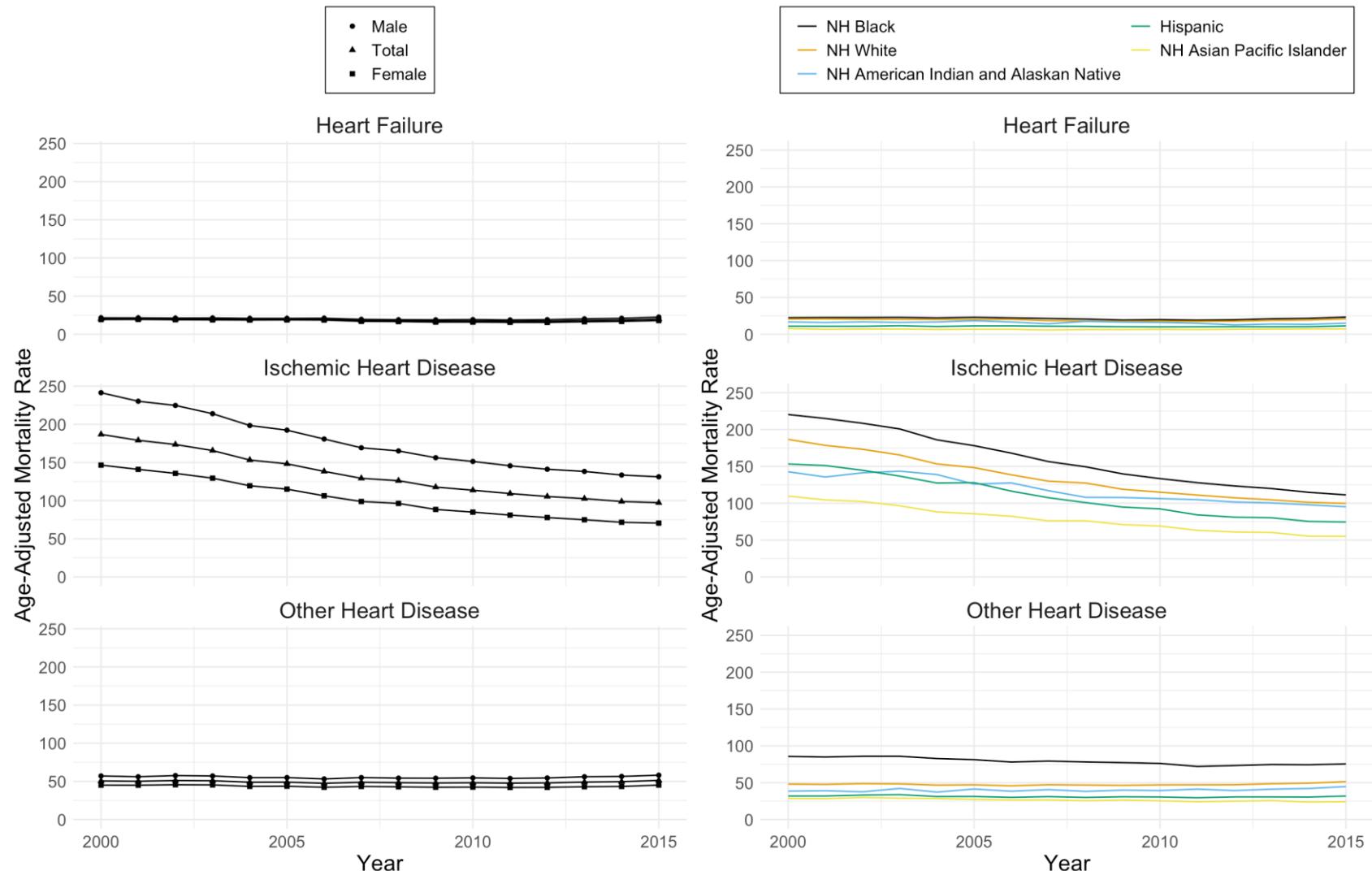
- Used unique colors from the same colorblind-friendly palette.
- Re-ordered the variable names to match the order of the lines.
- Used facet for disease type.



Example 4: Facets and Colors

Modifications:

- To better distinguish between the variable types, use shapes for sex and colors for race.



Example 5: Map Projection

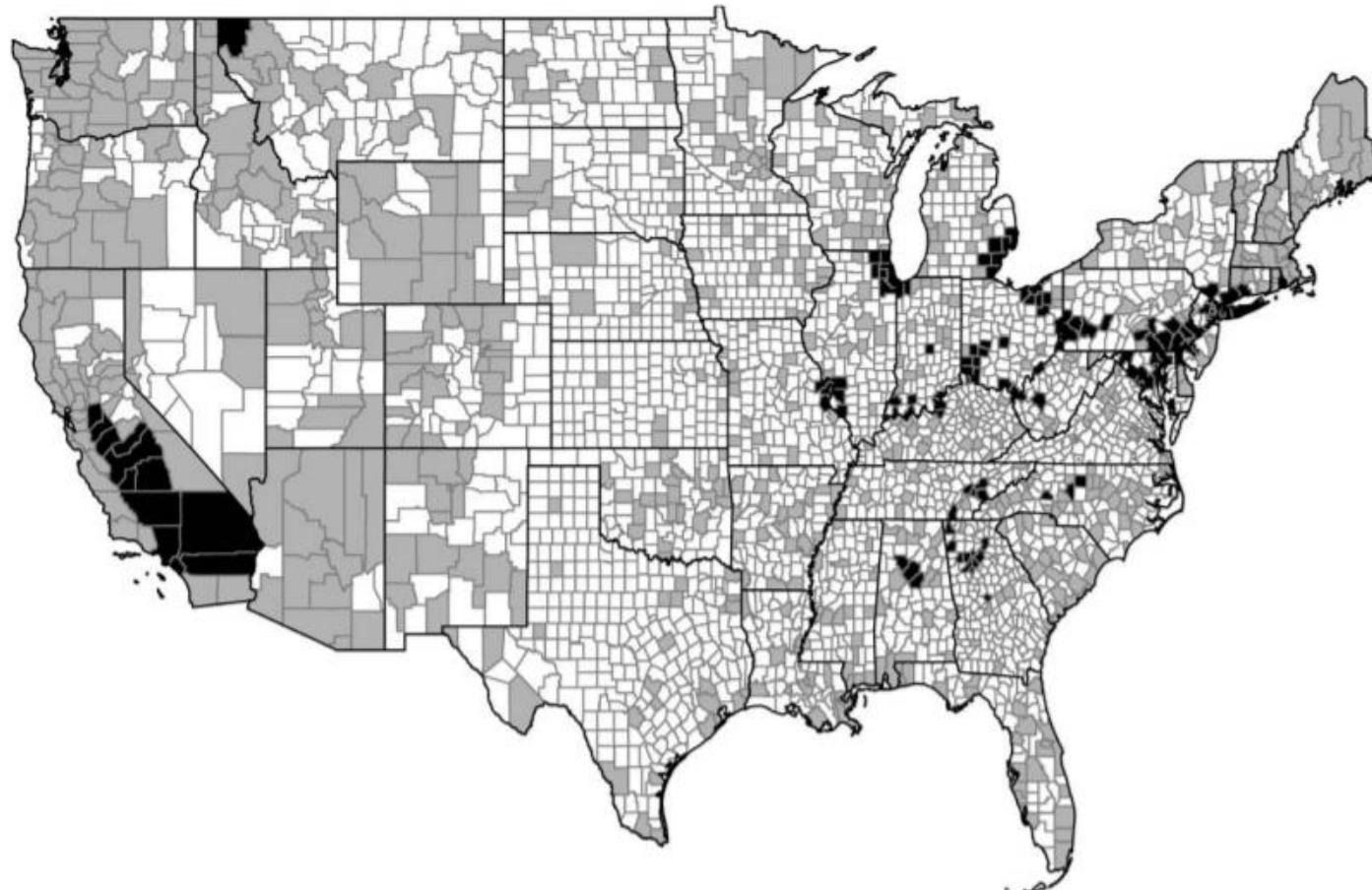
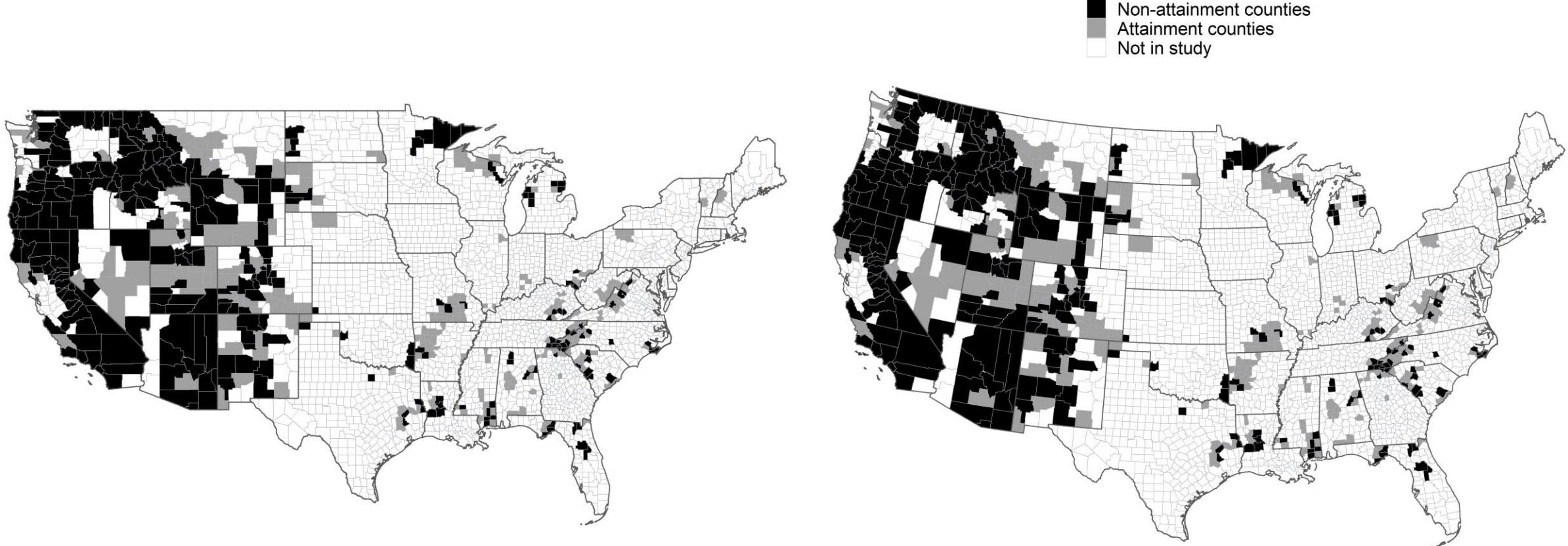


FIGURE 1. Map of nonattainment counties. Black and gray shaded areas represent counties in our sample. The black shade are the nonattainment counties as of 2005, while the gray are the attainment counties.

Sanders NJ, Barreca AI, Neidell MJ. Estimating Causal Effects of Particulate Matter Regulation on Mortality. *Epidemiology*. 2020 Mar;31(2):160-167. doi: 10.1097/EDE.0000000000001153.

Example 5: Map Projections

- Modifications:
- Changed projection to conic
 - Added legend



Example 6: Maps Boundaries

Correia AW, Pope CA 3rd, Dockery DW, Wang Y, Ezzati M, Dominici F. Effect of air pollution control on life expectancy in the United States: an analysis of 545 U.S. counties for the period from 2000 to 2007. *Epidemiology*. 2013 Jan;24(1):23-31. doi: 10.1097/EDE.0b013e3182770237.

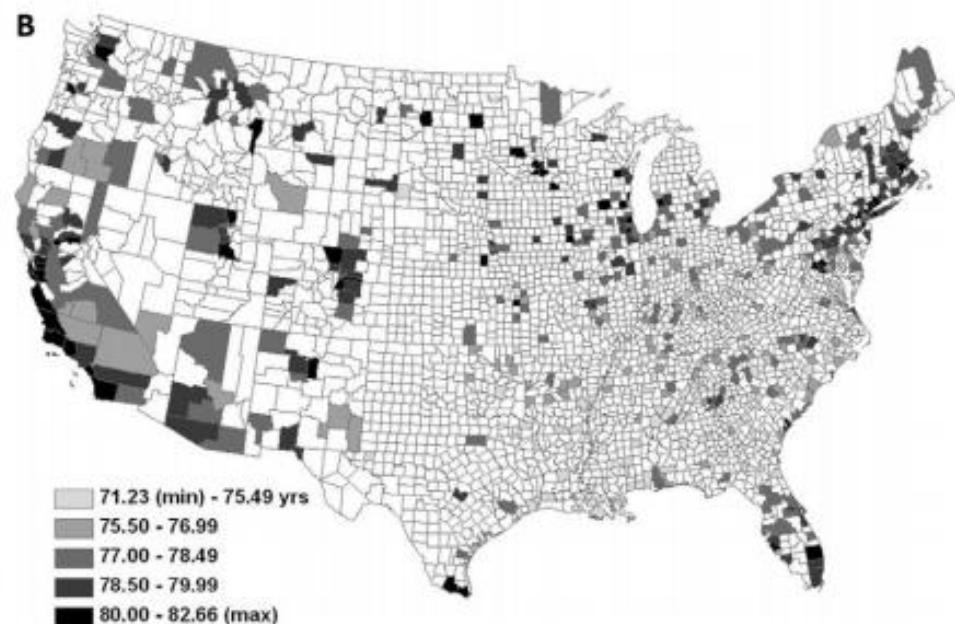
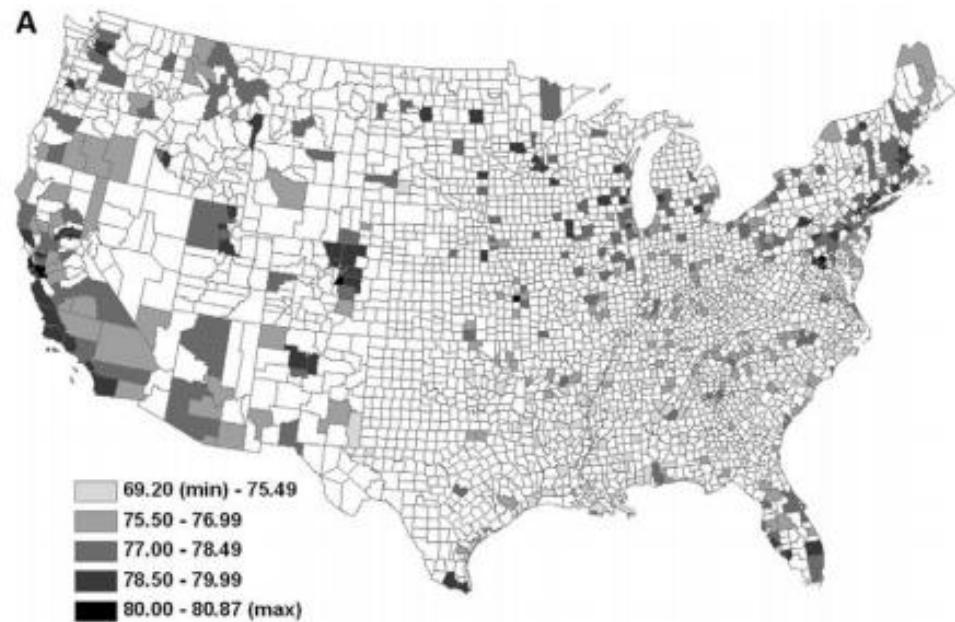
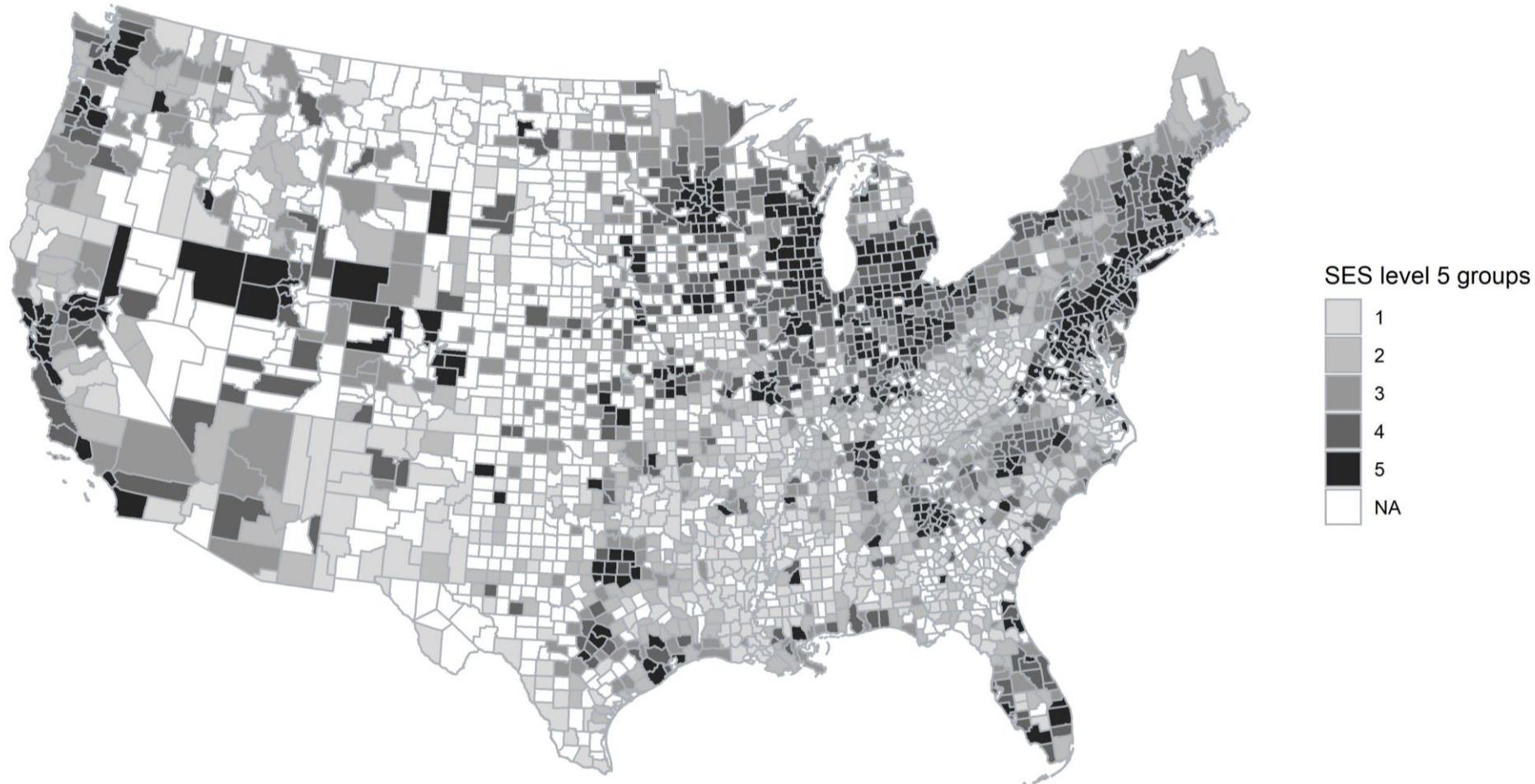


FIGURE 1. Map of United States with the 545 counties from data set 1 shaded according to (A) year 2000 and (B) year 2007 life expectancies.

Example 6: Map Boundaries

Similar map

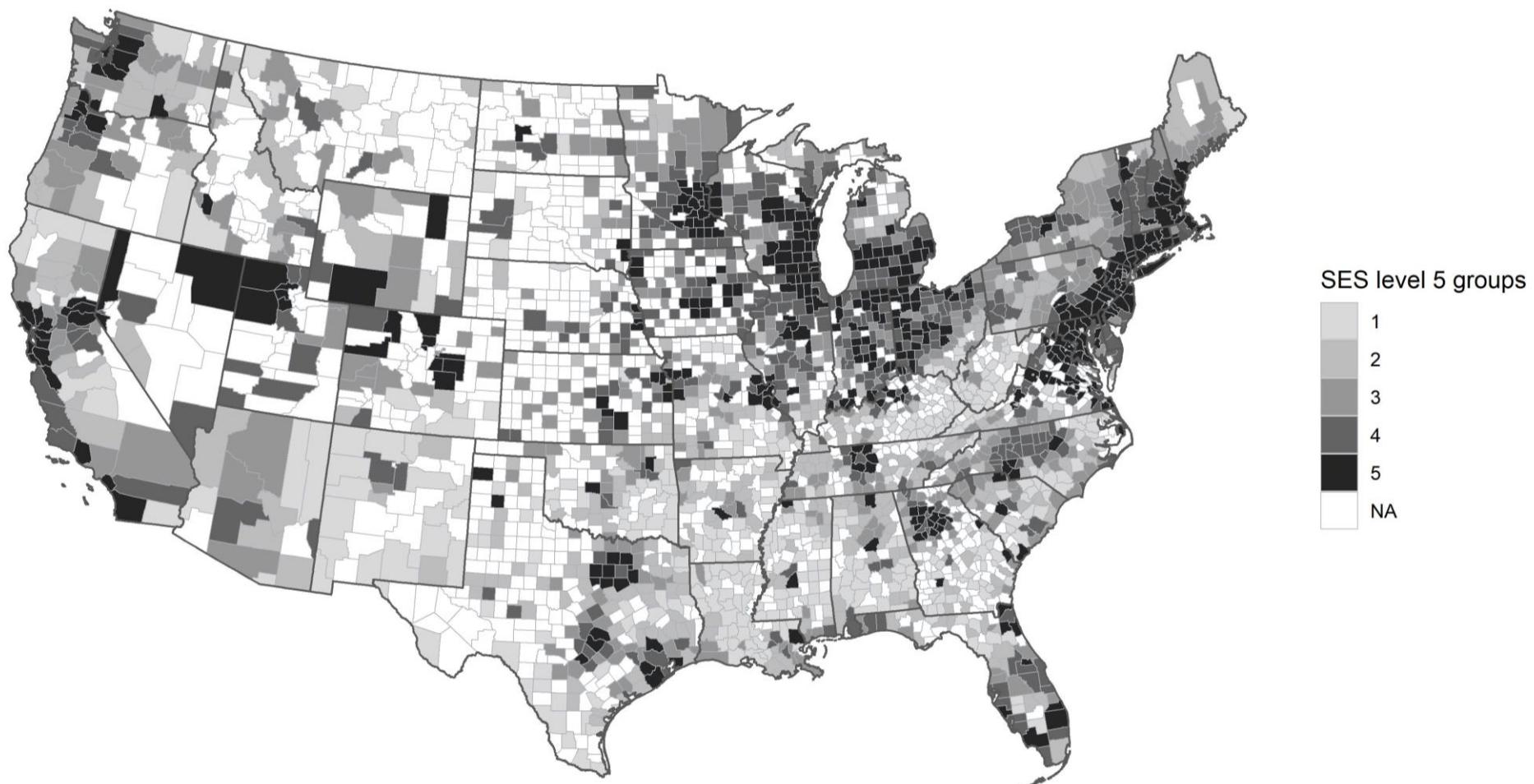
- All boundaries same size



Example 6: Map Boundaries

Modifications:

- Reduced/changed the county boundary outlines
- Altered state boundary lines



Example 7: Map Colors and Labels

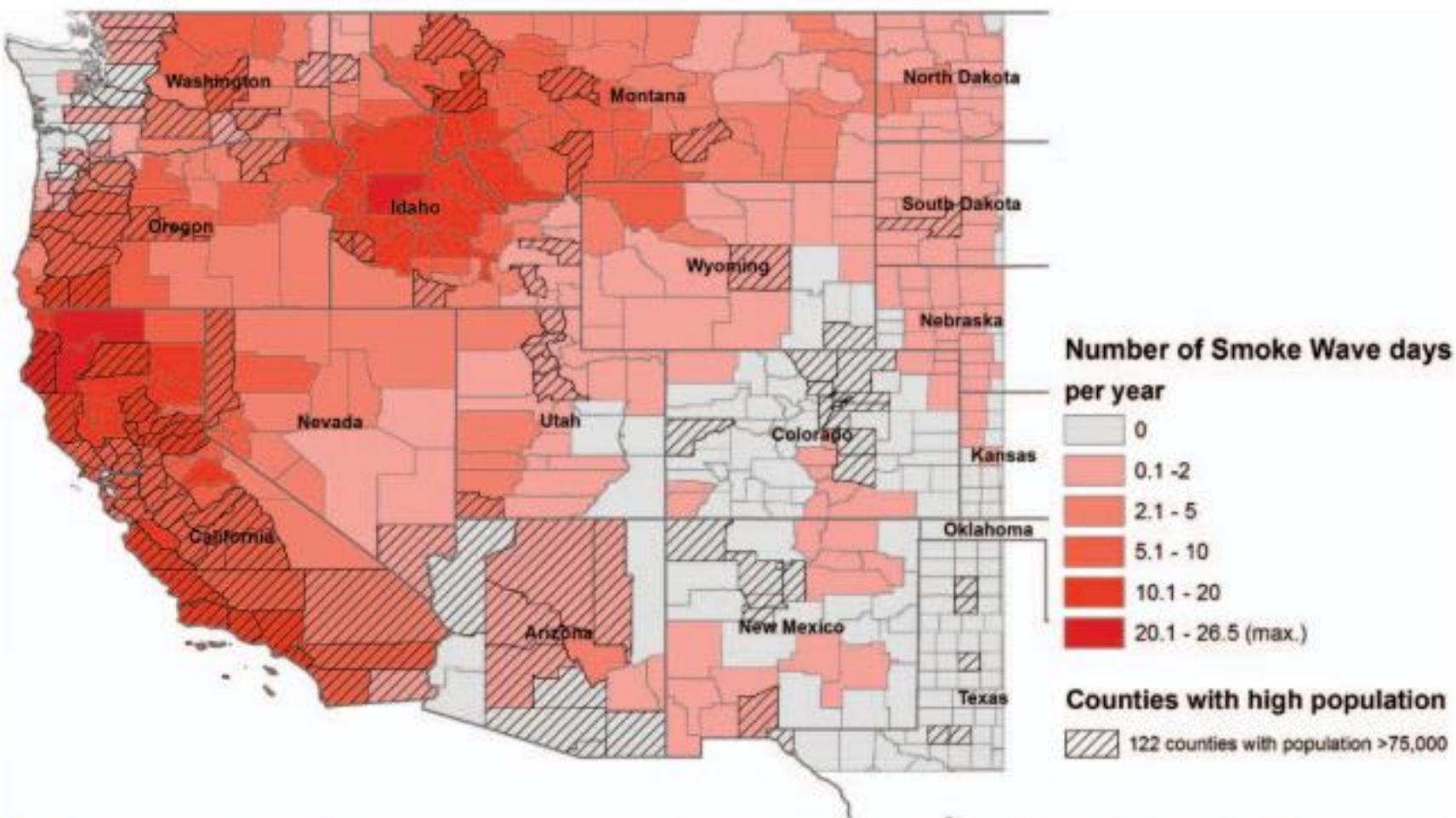
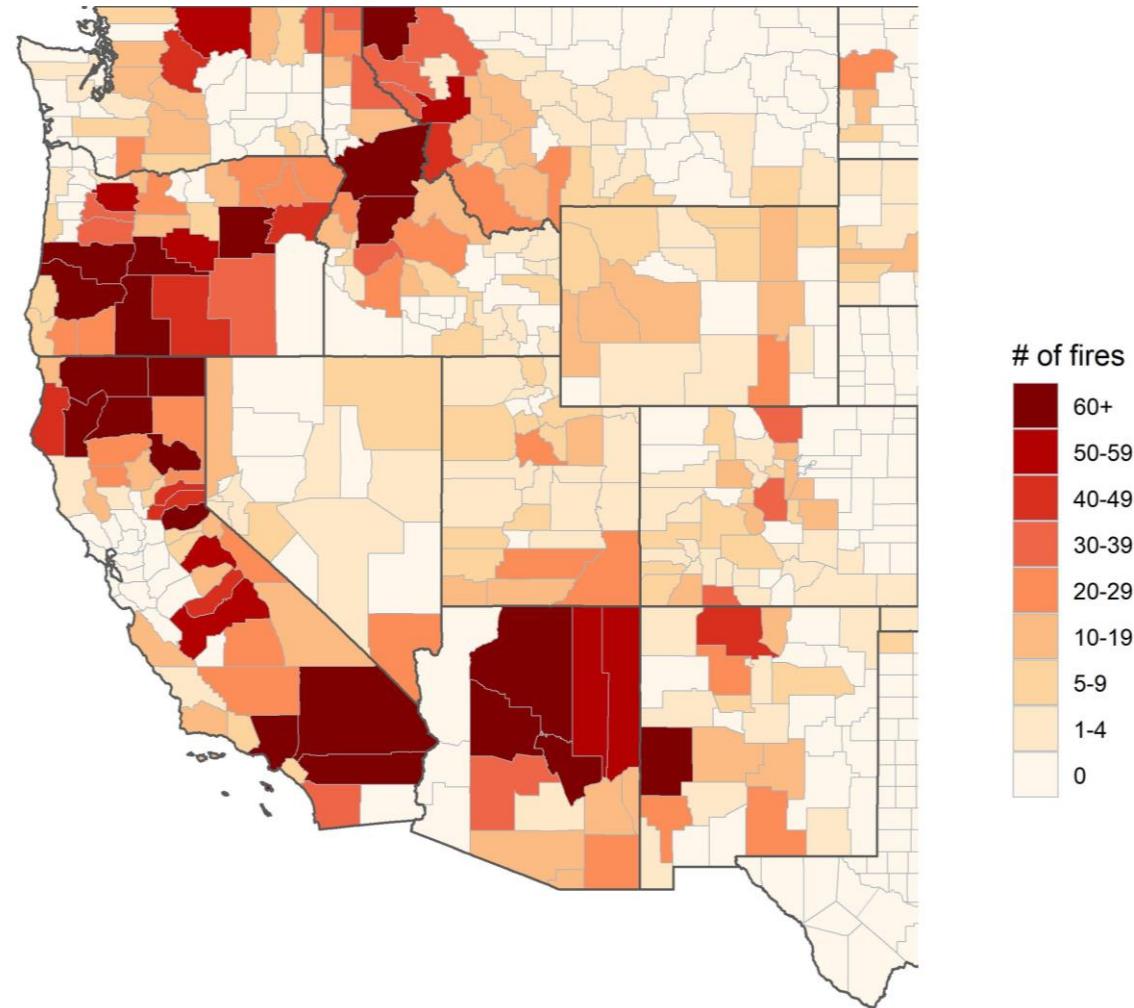


FIGURE 1. Average number of smoke wave days/year for 561 Western United States counties during 2004–2009. Hashed counties have population >75,000 in the 2010 Census.

Example 7: Map Colors and Labels

Modifications:

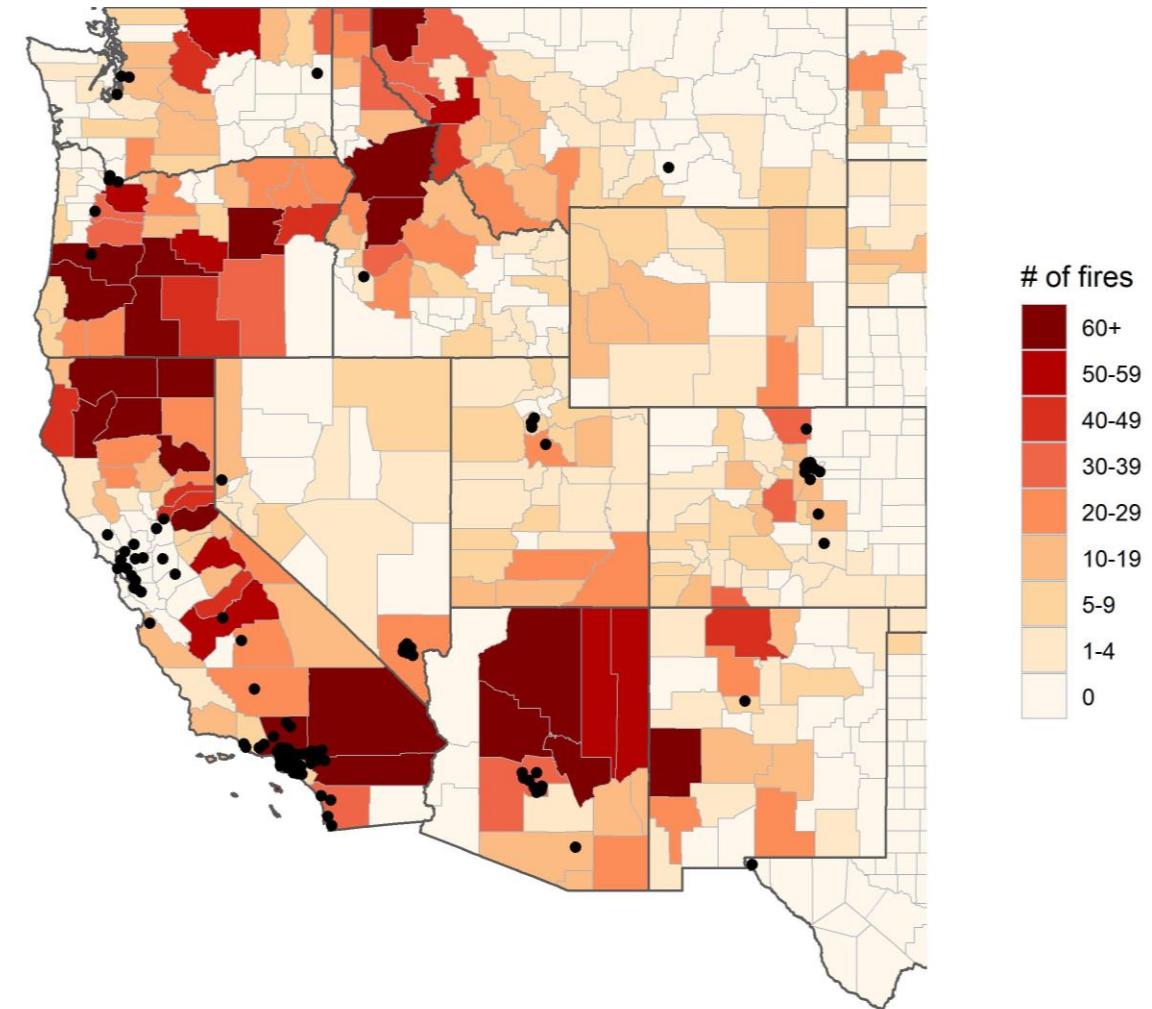
- Removed state labels
- Created state boundaries
- Changed the county boundary outlines



Example 7: Map Colors and Labels

Modifications:

- Removed state labels
- Created state boundaries
- Changed the county boundary outlines
- Added in points for populous cities



References

References

- Pederson TL. ggplot2 Workshop [Internet]. Copenhagen: GitHub; 2020 [updated 2020 March 25; cited 2020]. Available from: https://github.com/thomasp85/ggplot2_workshop.
- Peng RD and Dominici F. Statistical Methods for Environmental Epidemiology with R: A Case Study in Air Pollution and Health. New York: Springer-Verlag; 2008. Available from: <https://www.springer.com/gp/book/9780387781662>
- R Studio. Data Visualizations with ggplot2 Cheat Sheet [Internet]. Boston: R Studio; 2020 [updated 2015 March; cited 2020]. Available from: <https://rstudio.com/wp-content/uploads/2015/03/ggplot2-cheatsheet.pdf>
- Scherer C. A ggplot2 Tutorial for Beautiful Plotting in R [Internet]. 2019 [updated 2019 November 1; cited 2020]. Available from: <https://cedricscherer.netlify.app/2019/08/05/a-ggplot2-tutorial-for-beautiful-plotting-in-r/>
- Wickham H. ggplot2: Elegant Graphic for Data Analysis. New York: Springer-Verlag; 2009. Available from: <https://ggplot2-book.org/>
- Wilkinson L. The Grammar of Graphics. 2nd ed. New York: Springer-Verlag; 2005. Available from: <https://www.springer.com/gp/book/9780387245447>
- ZevRoss. Beautiful plotting in R: A ggplot2 cheatsheet [Internet]. Ithaca: ZevRoss; 2014 [updated 2016 January 20; cited 2020]. Available from: <http://zevross.com/blog/2014/08/04/beautiful-plotting-in-r-a-ggplot2-cheatsheet-3/>

References for Data

- United States Environmental Protection Agency (US EPA). National Air Quality: Status and Trends of Key Air Pollutants, Air Quality Trends. Available from: <https://www.epa.gov/air-trends>. [Accessed 2020].
- Institute for Health Metrics and Evaluation (IHME). United States Hypertension Estimates by County 2001-2009. 2013. Seattle, United States of America: Institute for Health Metrics and Evaluation (IHME). Available from: <http://ghdx.healthdata.org/record/ihme-data/united-states-hypertension-estimates-county-2001-2009>. [Accessed 2020].
- U.S. Department of Agriculture, Forest Service, Rocky Mountain Research Station. Spatial wildfire occurrence data for the United States, 1992-2011 [FPA_FOD_20130422] (1st Edition) Data publication contains GIS data Author(s): Short, Karen C. Publication Year:2013. Available from: <https://www.fs.usda.gov/rds/archive/Catalog/RDS-2013-0009>. [Accessed 2020].

References for Examples

- St Sauver JL, Grossardt BR, Yawn BP, Melton LJ 3rd, Rocca WA. Use of a medical records linkage system to enumerate a dynamic population over time: the Rochester epidemiology project. *Am J Epidemiol.* 2011 May 1;173(9):1059-68. doi: 10.1093/aje/kwq482. Epub 2011 Mar 23.
- Liu JC, Wilson A, Mickley LJ, Dominici F, Ebisu K, Wang Y, Sulprizio MP, Peng RD, Yue X, Son JY, Anderson GB, Bell ML. Wildfire-specific Fine Particulate Matter and Risk of Hospital Admissions in Urban and Rural Counties. *Epidemiology.* 2017 Jan;28(1):77-85. doi: 10.1097/EDE.0000000000000556.
- Sanders NJ, Barreca AI, Neidell MJ. Estimating Causal Effects of Particulate Matter Regulation on Mortality. *Epidemiology.* 2020 Mar;31(2):160-167. doi: 10.1097/EDE.0000000000001153.
- Pun VC, Kazemiparkouhi F, Manjourides J, Suh HH. Long-Term PM2.5 Exposure and Respiratory, Cancer, and Cardiovascular Mortality in Older US Adults. *Am J Epidemiol.* 2017 Oct 15;186(8):961-969. doi: 10.1093/aje/kwx166.
- Sidney S, Quesenberry CP Jr, Jaffe MG, Sorel M, Go AS, Rana JS. Heterogeneity in national U.S. mortality trends within heart disease subgroups, 2000-2015. *BMC Cardiovasc Disord.* 2017 Jul 18;17(1):192. doi: 10.1186/s12872-017-0630-2.
- Wang B, Eum KD, Kazemiparkouhi F, Li C, Manjourides J, Pavlu V, Suh H. The impact of long-term PM2.5 exposure on specific causes of death: exposure-response curves and effect modification among 53 million U.S. Medicare beneficiaries. *Environ Health.* 2020 Feb 17;19(1):20. doi: 10.1186/s12940-020-00575-0.
- Correia AW, Pope CA 3rd, Dockery DW, Wang Y, Ezzati M, Dominici F. Effect of air pollution control on life expectancy in the United States: an analysis of 545 U.S. counties for the period from 2000 to 2007. *Epidemiology.* 2013 Jan;24(1):23-31. doi: 10.1097/EDE.0b013e3182770237.