

Data distributions

Cumulative distribution functions



Data by percentile rank

- PMFs are handy exploratory tools, but as with histograms, the binwidth can strongly influence what your plot looks like
- We can overcome this problem if we convert the data into a sorted list of percentile ranks
- **Advantages**
 - Don't need to select a binsize
 - Easier to compare similarities and differences of different data distributions
 - Different classes of data distributions have distinct shapes
- The **cumulative distribution function** (CDF) lets us map between percentile rank and each value in a data column

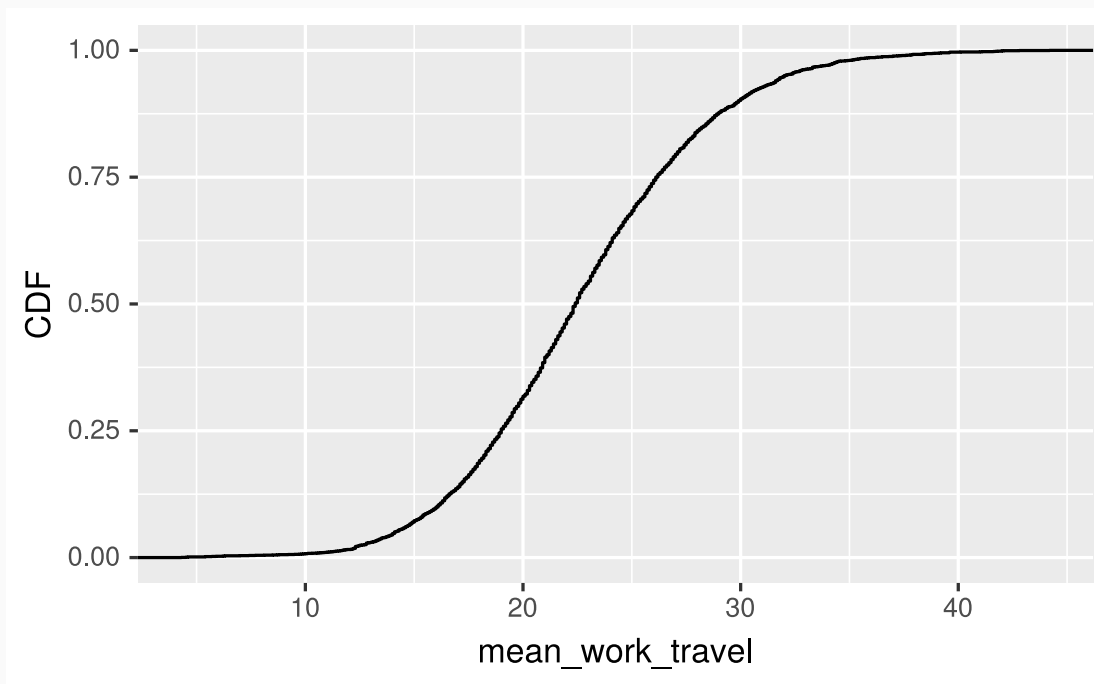
Creating CDFs in R

`ggplot2` comes with a handy convenience function `stat_ecdf()`, which lets you create CDF functions from your data

Creating CDFs in R

`ggplot2` comes with a handy convenience function `stat_ecdf()`, which lets you create CDF functions from your data

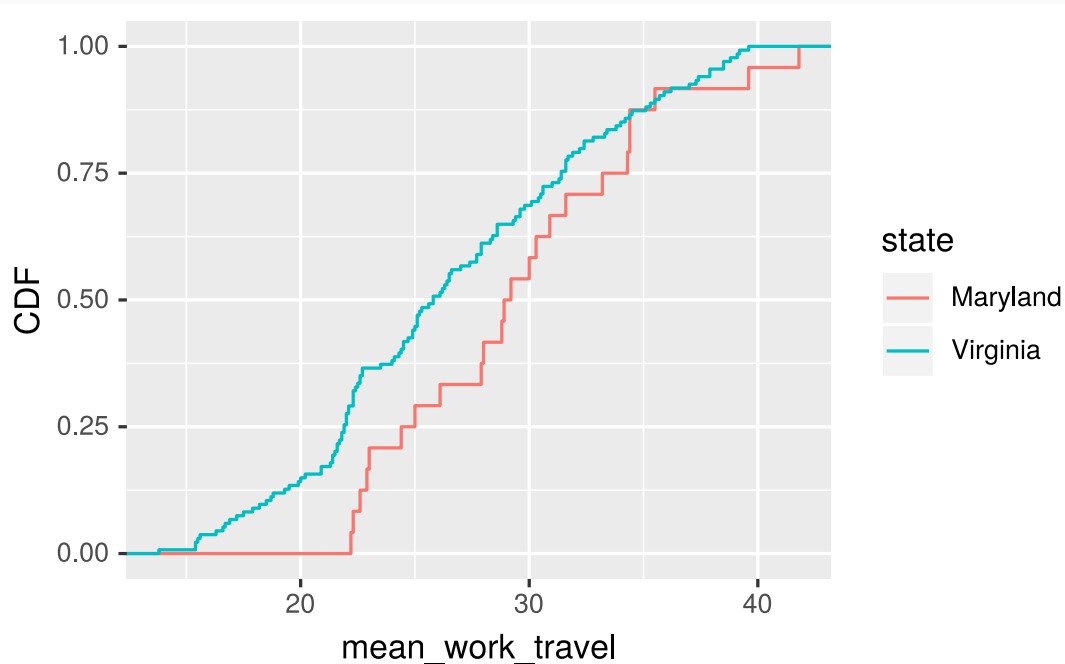
```
county %>%  
  ggplot() +  
  stat_ecdf(mapping = aes(x = mean_work_travel)) +  
  labs(y = "CDF")
```



Creating CDFs in R

We can do all the usual operations, such as grouping by state

```
county %>%  
  filter(state == "Virginia" | state == "Maryland") %>%  
  ggplot() +  
  stat_ecdf(mapping = aes(x = mean_work_travel, color = state)) +  
  labs(y = "CDF")
```



Computing the CDF

To compute the CDF, we use the `cume_dist()` function along with `filter()`, `group_by()`, and `mutate()`:

```
va_md_cdf_df <- county %>%  
  filter(state == "Virginia" | state == "Maryland") %>%  
  group_by(state) %>%  
  mutate(cdf = cume_dist(mean_work_travel)) %>%  
  select(state, mean_work_travel, cdf)
```

Get CDF data out of plot

state	mean_work_travel	cdf
Virginia	13.8	0.0074627
Virginia	15.4	0.0223881
Virginia	15.4	0.0223881
Virginia	15.5	0.0298507
Virginia	15.6	0.0373134
Virginia	16.3	0.0447761
Virginia	16.6	0.0522388
Virginia	16.7	0.0597015
Virginia	16.9	0.0671642
Virginia	17.2	0.0746269

Credits

License

Creative Commons Attribution-NonCommerical-ShareAlike 4.0 International