





•

•

|



Each of the four data sets yields the same standard output from a typical regression program, namely

Number of observations (n) = 11

Mean of the x 's (\bar{x}) = 9.0

Mean of the y 's (\bar{y}) = 7.5

Regression coefficient (b_1) of y on x = 0.5

Equation of regression line: $y = 3 + 0.5x$

Sum of squares of $x - \bar{x}$ = 110.0

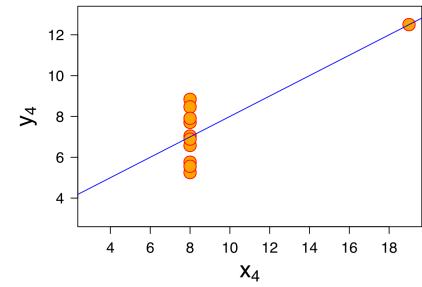
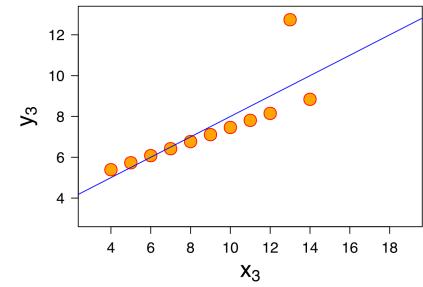
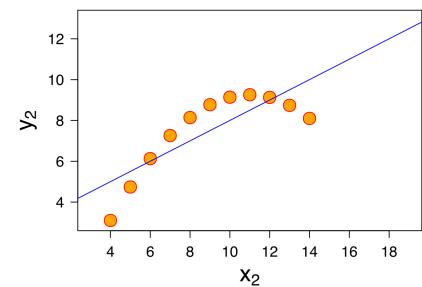
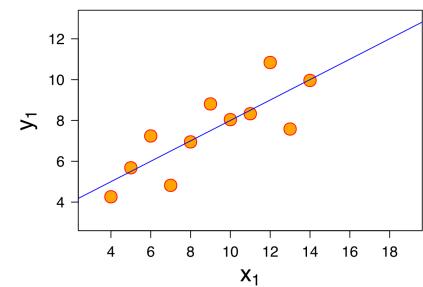
Regression sum of squares = 27.50 (1 d.f.)

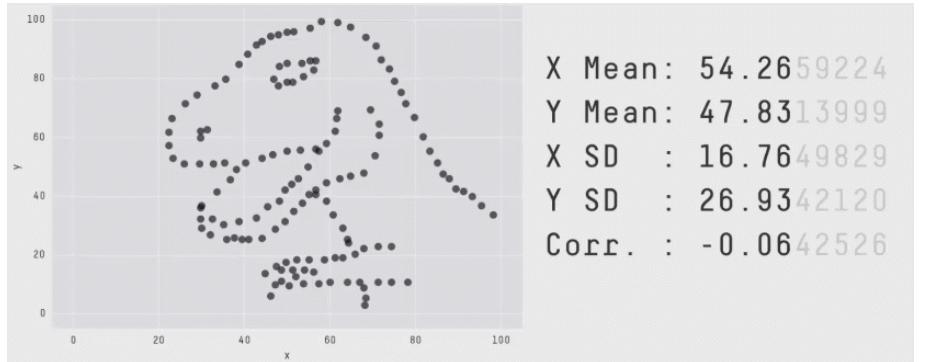
Residual sum of squares of y = 13.75 (9 d.f.)

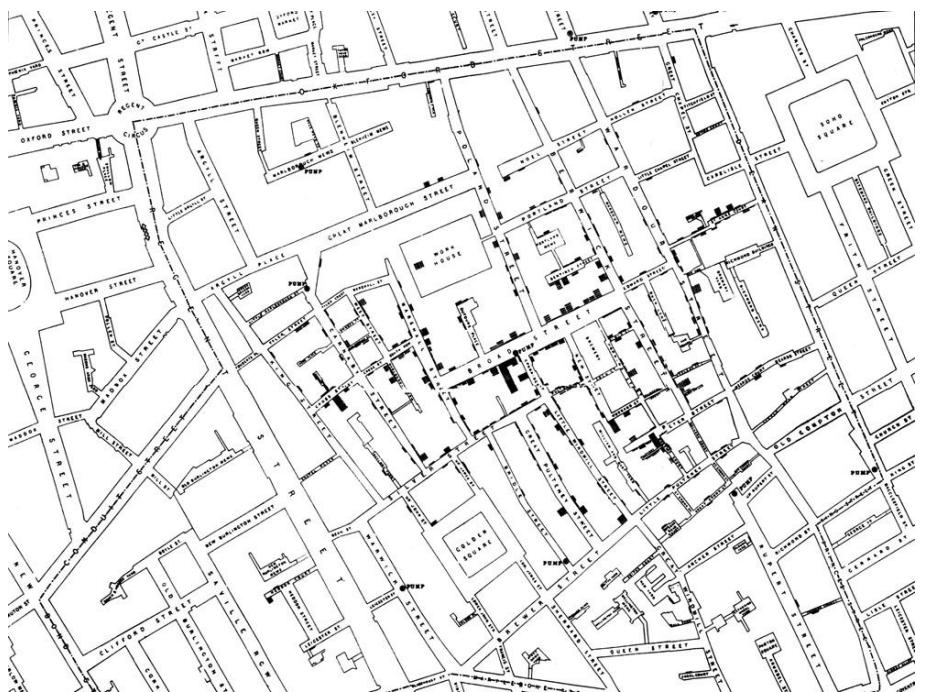
Estimated standard error of b_1 = 0.118

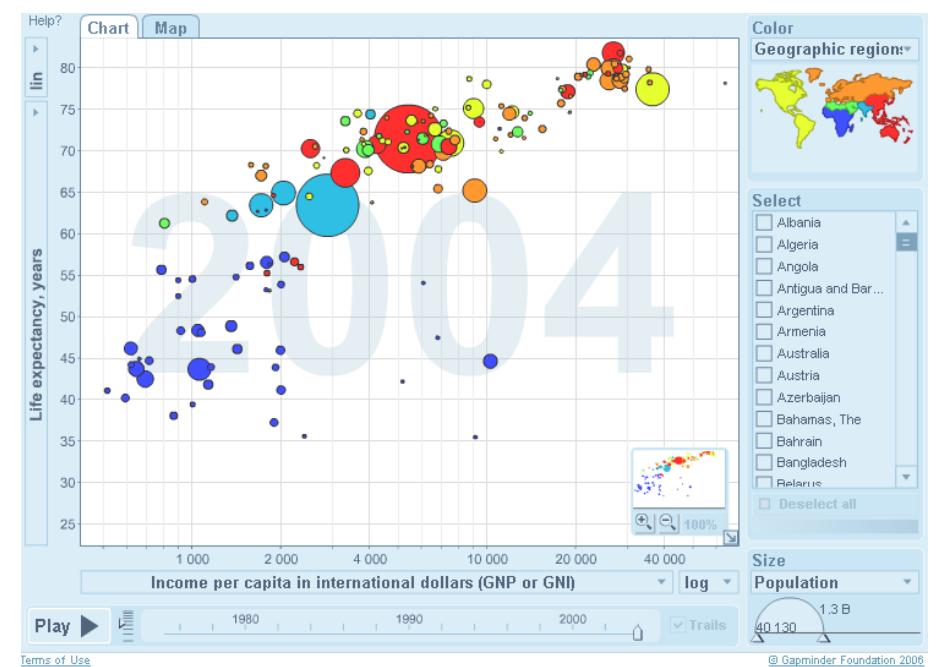
Multiple R^2 = 0.667

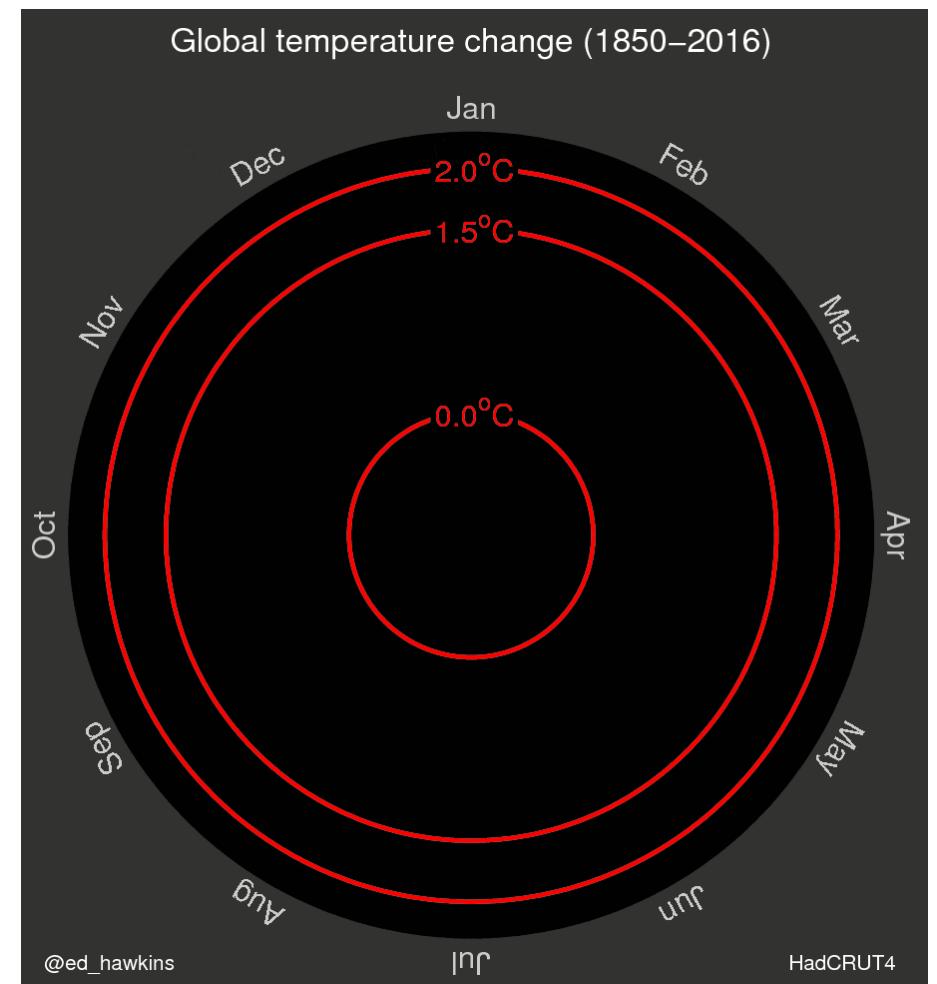


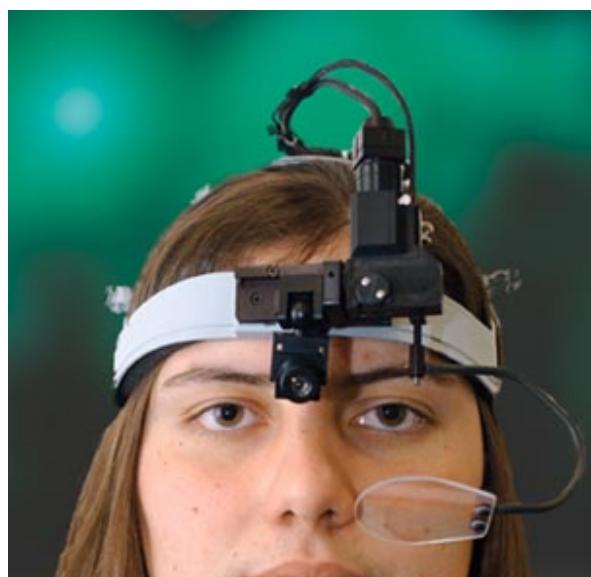












•

•



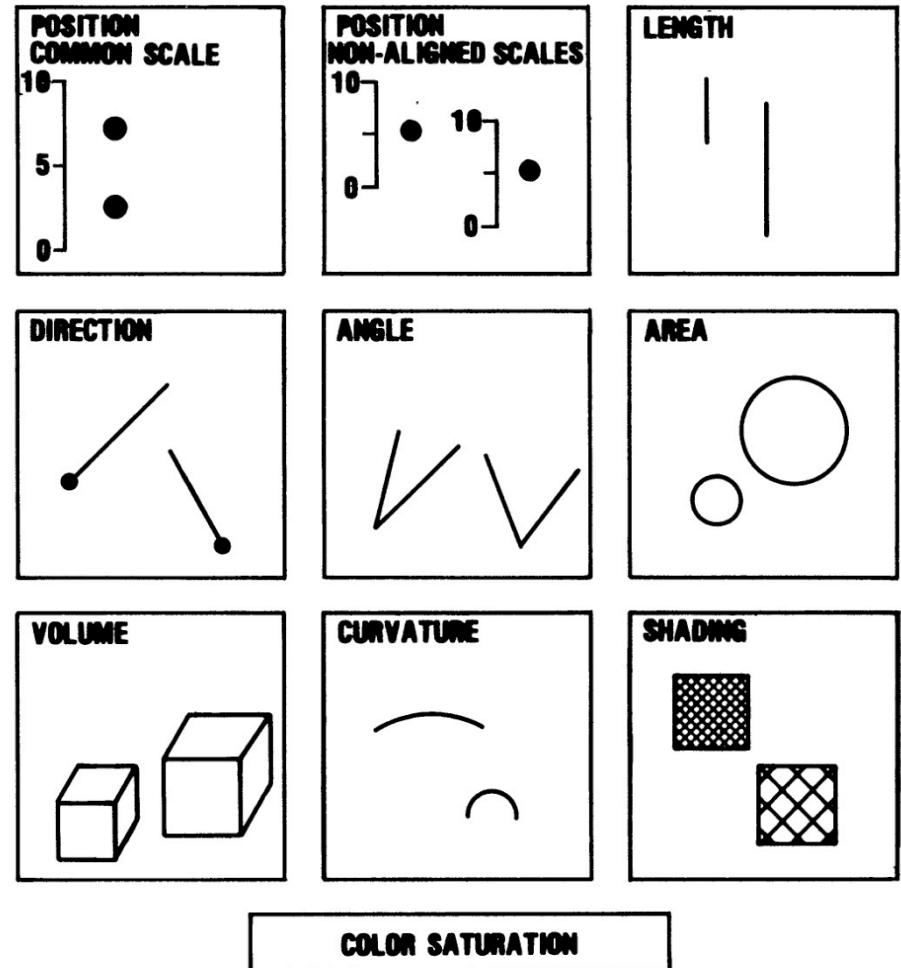


Figure 1. Elementary perceptual tasks.



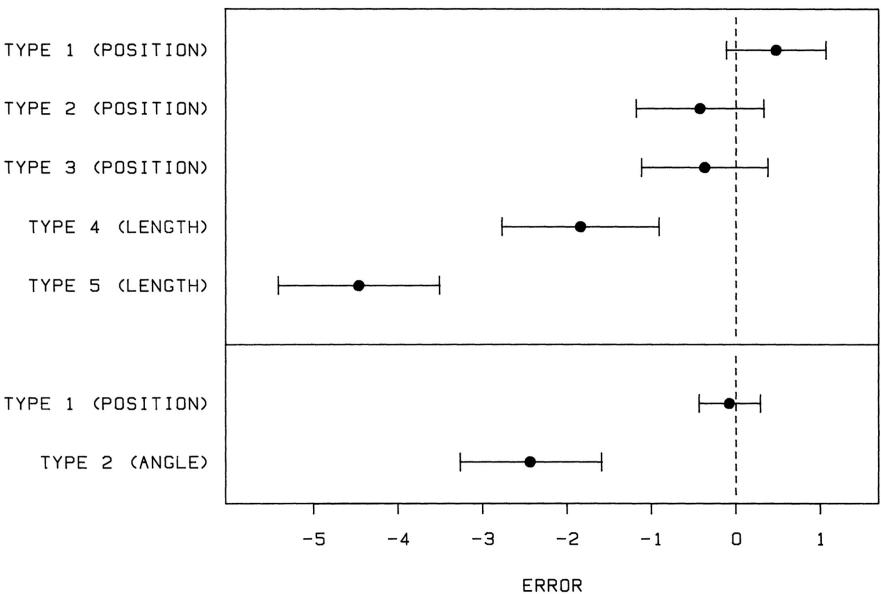


Figure 20. Error means and 95% confidence intervals for judgment types in position-length experiment (top) and position-angle experiment (bottom).



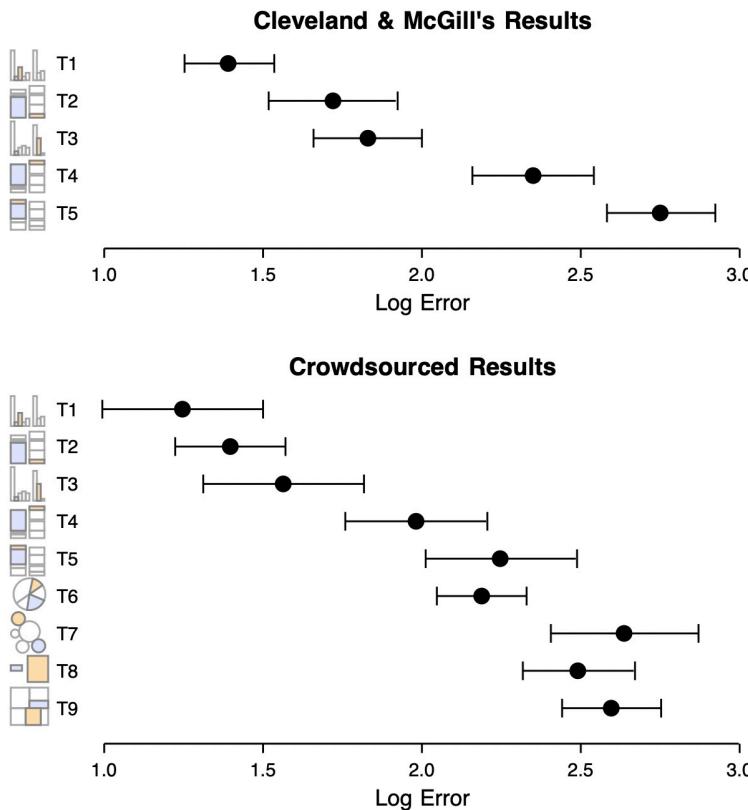
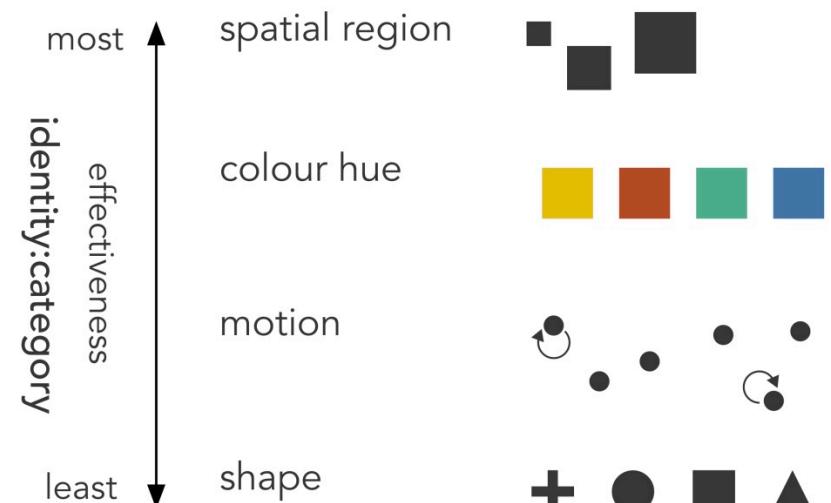
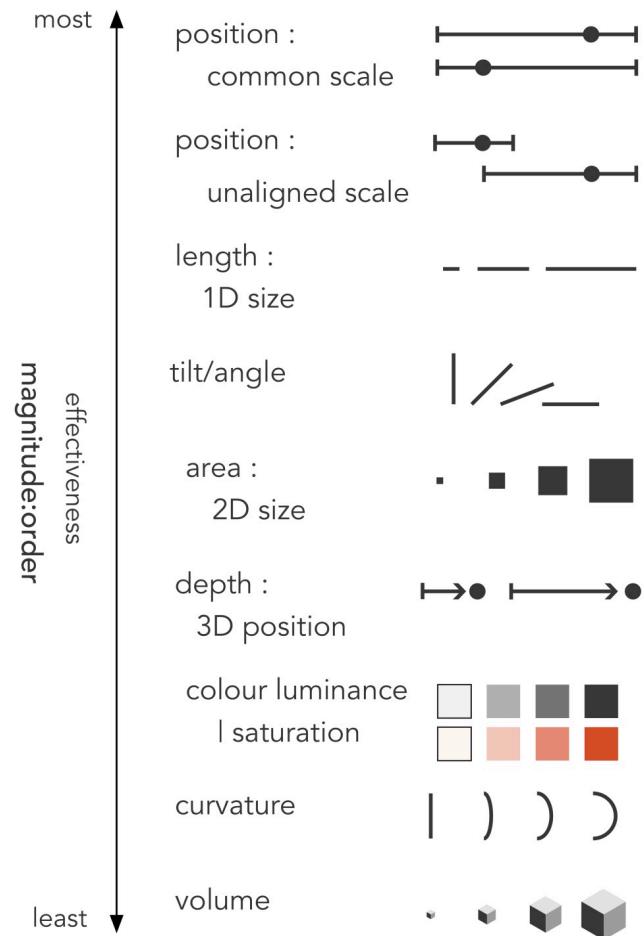


Figure 4: Proportional judgment results (Exp. 1A & B).
 Top: Cleveland & McGill's [7] lab study. Bottom: MTurk studies. Error bars indicate 95% confidence intervals.





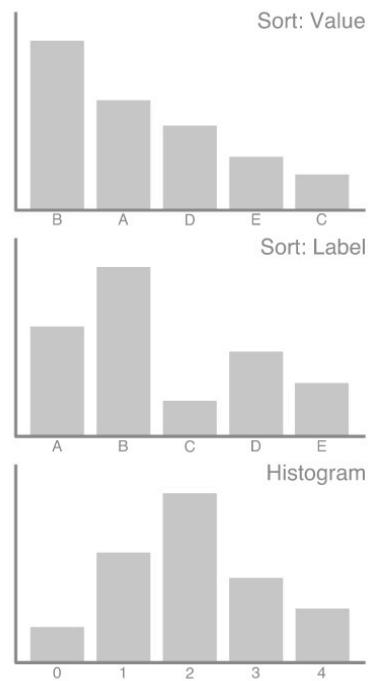
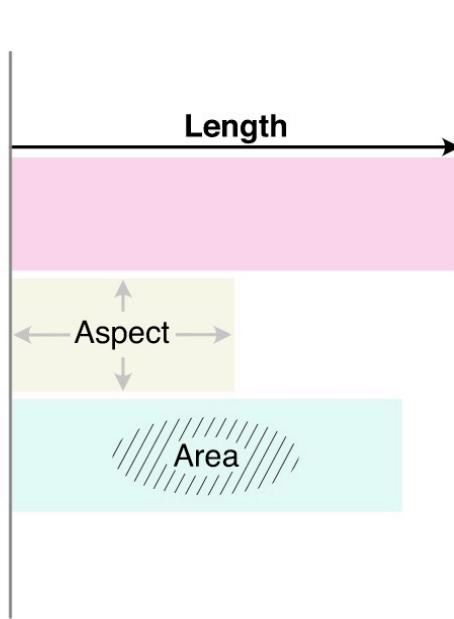


FIGURE 2. Bar charts encode values as their length, but also their area, aspect ratio, and overall shape. Sorting is often used for specific use cases.

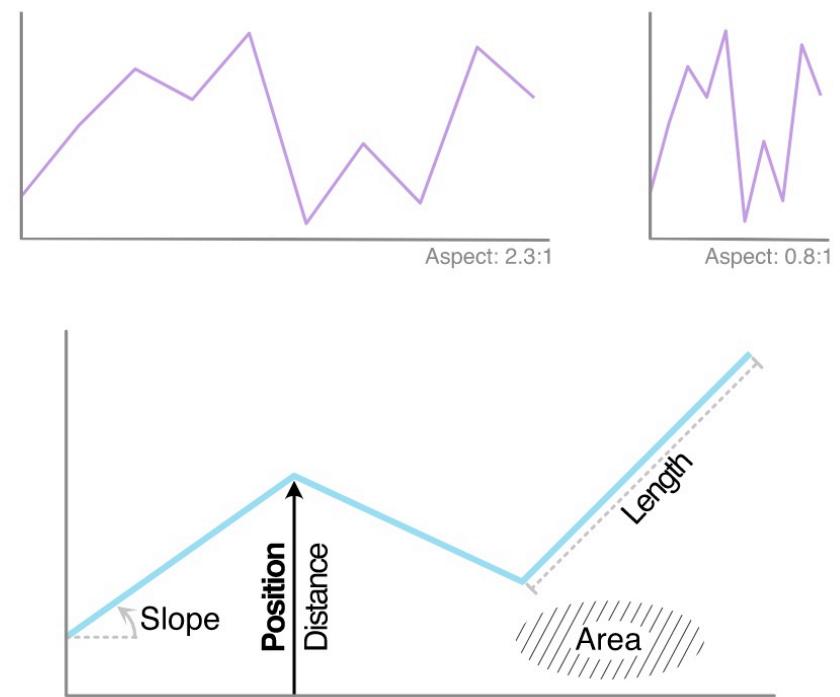


FIGURE 3. Line charts are specified by the location of the points connected by lines, but are read as slope and length, as well as area. Aspect ratio of the chart is also generally considered important.

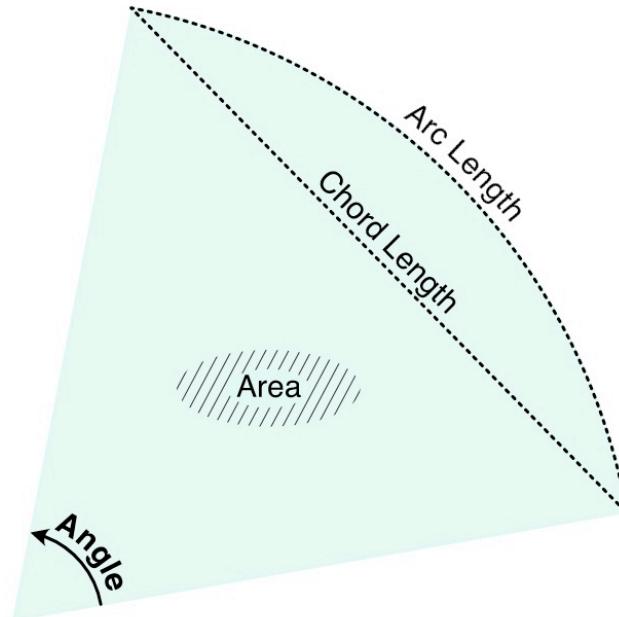


The *specified encoding* transforms a data value into a shape, ...



Angle

... which expresses that value as multiple *observable encodings*, ...



... some subset of which end up being *observed* by the user.

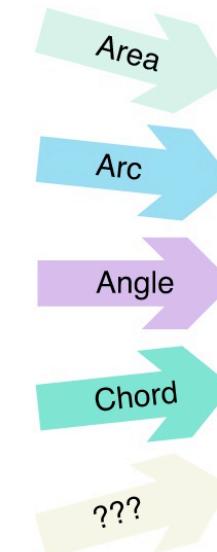
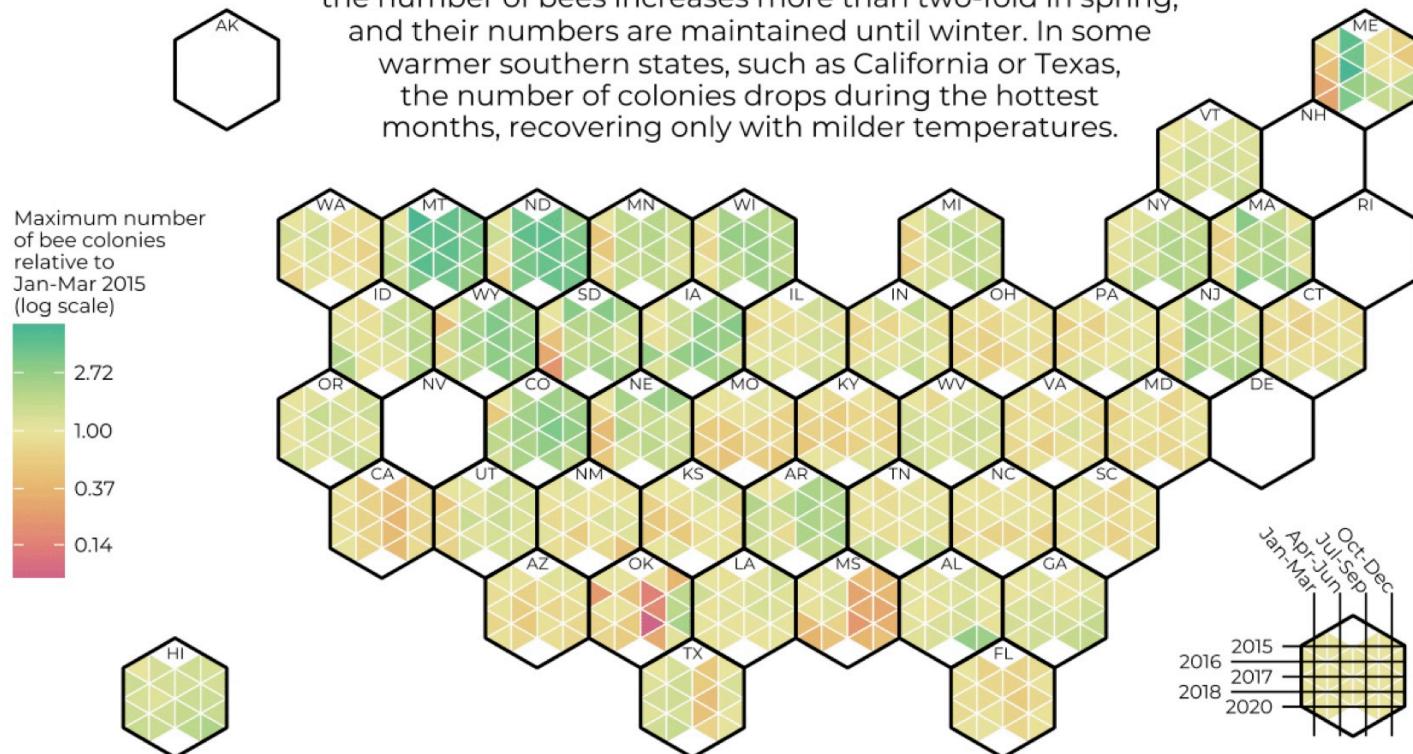


FIGURE 1. Pie charts are specified by angle, but may also be read by area, arc length, or even chord. Shape recognition is also likely for specific angles like $90^\circ/25\%$ and $180^\circ/50\%$.

Data from USDA
Graph by @irg_bio



—



To extract accu

To quantitatively

To find the largest.

To find unusual



You have a story you want to tell

The reader wants

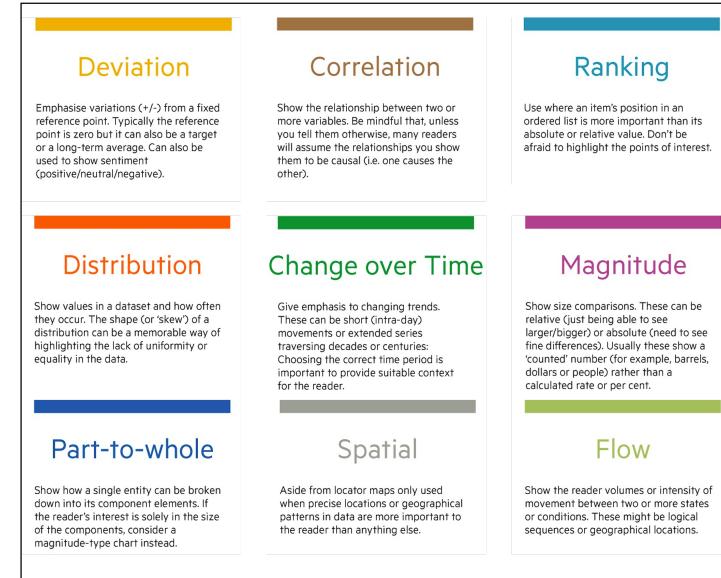
The reader has a preconception about



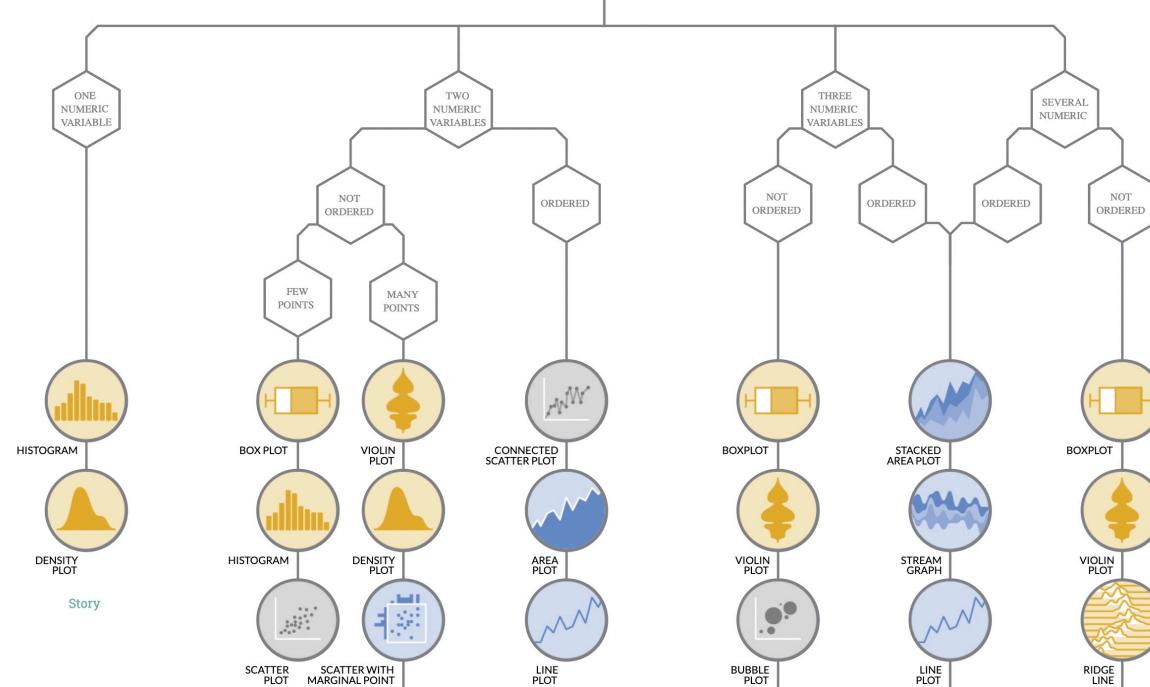

```

> msleep
# A tibble: 83 × 11
  name    genus vore order conse...¹ sleep...² sleep...³ sleep...⁴ awake brainwt bodywt
  <chr>   <chr> <chr> <chr> <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
1 Cheetah Acin... carni Carn... lc     12.1    NA     NA    11.9  NA      50
2 Owl mo... Aotus omni Prim... NA      17       1.8    NA      7    0.0155  0.48
3 Mounta... Aplo... herbi Rode... nt     14.4     2.4    NA     9.6  NA      1.35
4 Greate... Blar... omni Sori... lc     14.9     2.3   0.133   9.1  0.00029  0.019
5 Cow      Bos    herbi Arti... domest... 4       0.7    0.667   20    0.423   600
6 Three-... Brad... herbi Pil... NA     14.4     2.2   0.767   9.6  NA      3.85
7 Northe... Call... carni Carn... vu     8.7      1.4   0.383   15.3  NA      20.5
8 Vesper... Calo... NA    Rode... NA     7       NA     NA     17  NA      0.045
9 Dog      Canis carni Carn... domest... 10.1     2.9   0.333   13.9  0.07    14
10 Roe de... Capr... herbi Arti... lc     3       NA     NA     21  0.0982   14.8
# ... with 73 more rows, and abbreviated variable names `¹conservation, `²sleep_total,
#   `³sleep_rem, `⁴sleep_cycle
# i Use `print(n = ...)` to see more rows

```



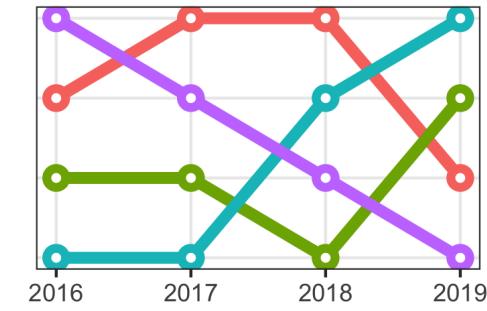
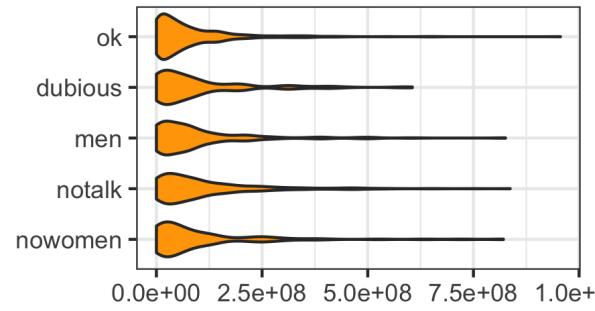
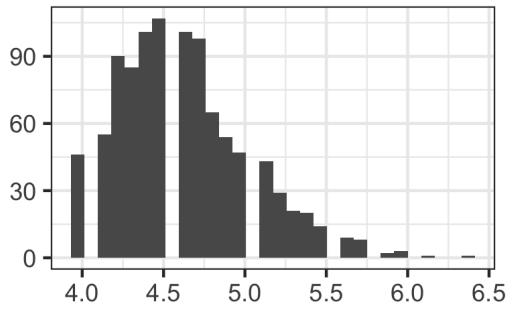
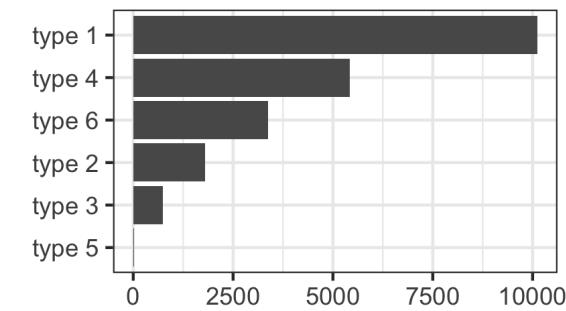
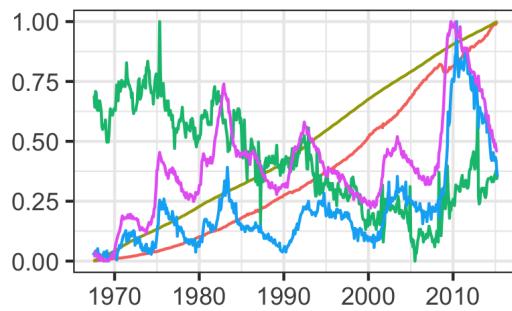
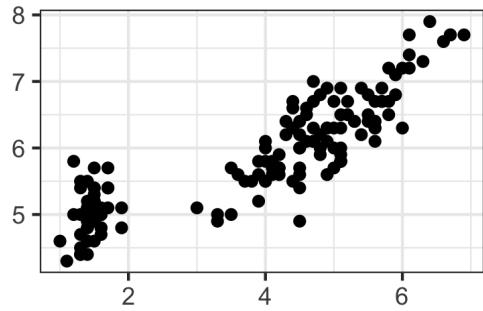
Numeric Categoric Num & Cat Maps Network Time series





|





response

- Maybe
- No
- Unsure
- Yes





Aesthetics



Geoms



Scales



Guides



Theme





```
1 msl eep %>%
2   ggpl ot() +
3   aes(
4     x = sl eep_total,
5     y = sl eep_rem
6     colour = vore
7   )
```





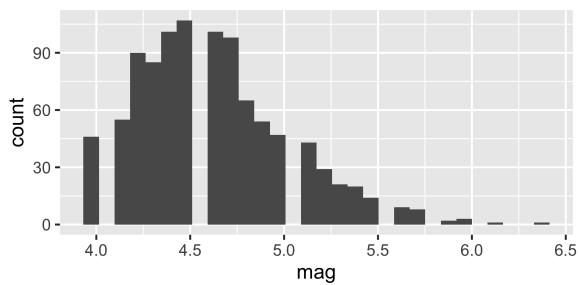


```
1 mslleep %>
2 ggplot() +
3 aes(
4   x = sleep_total,
5   y = sleep_rem,
6   colour = vore
7 ) +
8 geom_point()
```

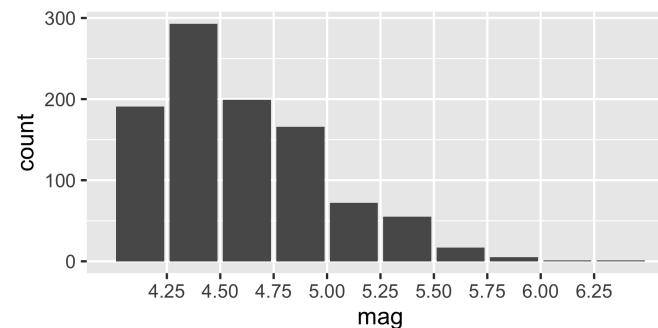
```
geom_ablignem(m)_ar,geom_bagre(m)_bi_nge(m)_bl_agre(m)_boxplot()
geom_cagle(m)_cont,geom_contour_,geom_loime(m)_ngre(m)_crossbar()
geom_curgre(m)_dens,geom_densitg(m)_density_2d_filled()
geom_densitg(m)_density_2d_gfeiorm_eddb(t)pg(m)_errorbar()
geom_error,geom_freq_pgo(m)_func_tgi_eom(m)_hegxe(m)_histogram()
geom_hli,geom_jitt,geom_labf(m)_li,geom_liner,geom_map()
geom_pa,geom_poi,geom_point,geom_point_(pol,y,g(m))_qq()
geom_qq_l,geom_quantg(m)_rass,geom_re,geom_rib,geom_bon()
geom_rgge(m)_segmgat(m)_s,geom_sf_lagbed(m)_sf_t,ext()
geom_smote(m)_spok(m)_st,geom_te,geom_t,geom_violin()
geom_vline()
```



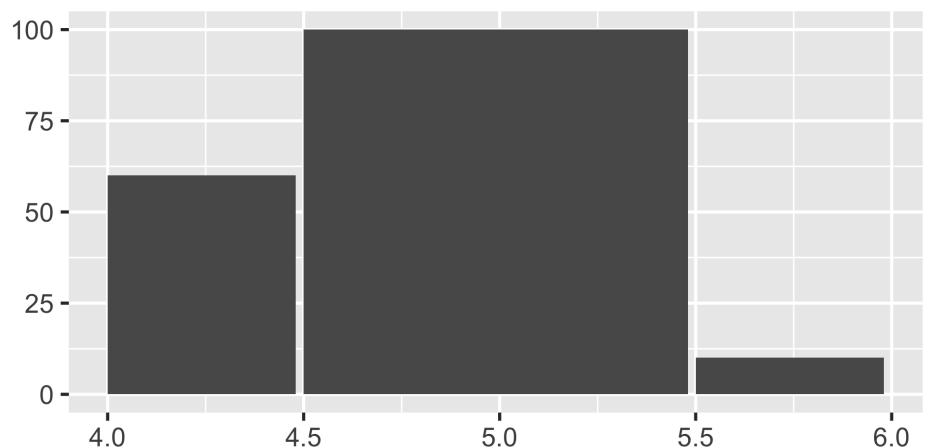
```
1 ggplot(quakes, aes(mag)) +  
2 geom_histogram()
```



```
1 ggplot(quakes, aes(mag)) +  
2 geom_bar() +  
3 scale_x_binned()
```



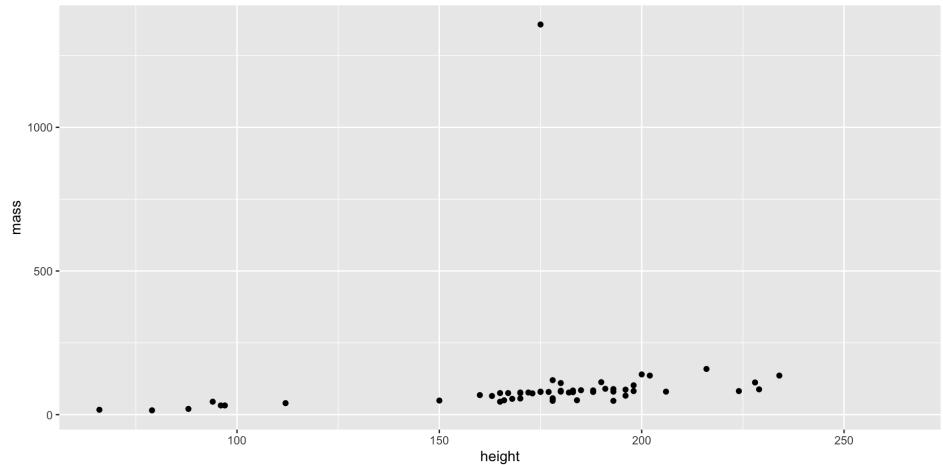
```
1 rect_data <- tribble(
2   ~x_min, ~x_max, ~y_min, ~y_max,
3   4, 4.48, 0, 60,
4   4.5, 5.48, 0, 100,
5   5.5, 5.98, 0, 10
6 )
7 rect_data %>%
8   ggplot() +
9   geom_rect(aes(xmin = x_min,
10               xmax = x_max,
11               ymin = y_min,
12               ymax = y_max)) +
13   theme_gray(base_size = 25)
```



X

y

```
1 starwars %>%
2   ggplot() +
3   aes(x = height,
4        y = mass) +
5   geom_point()
```



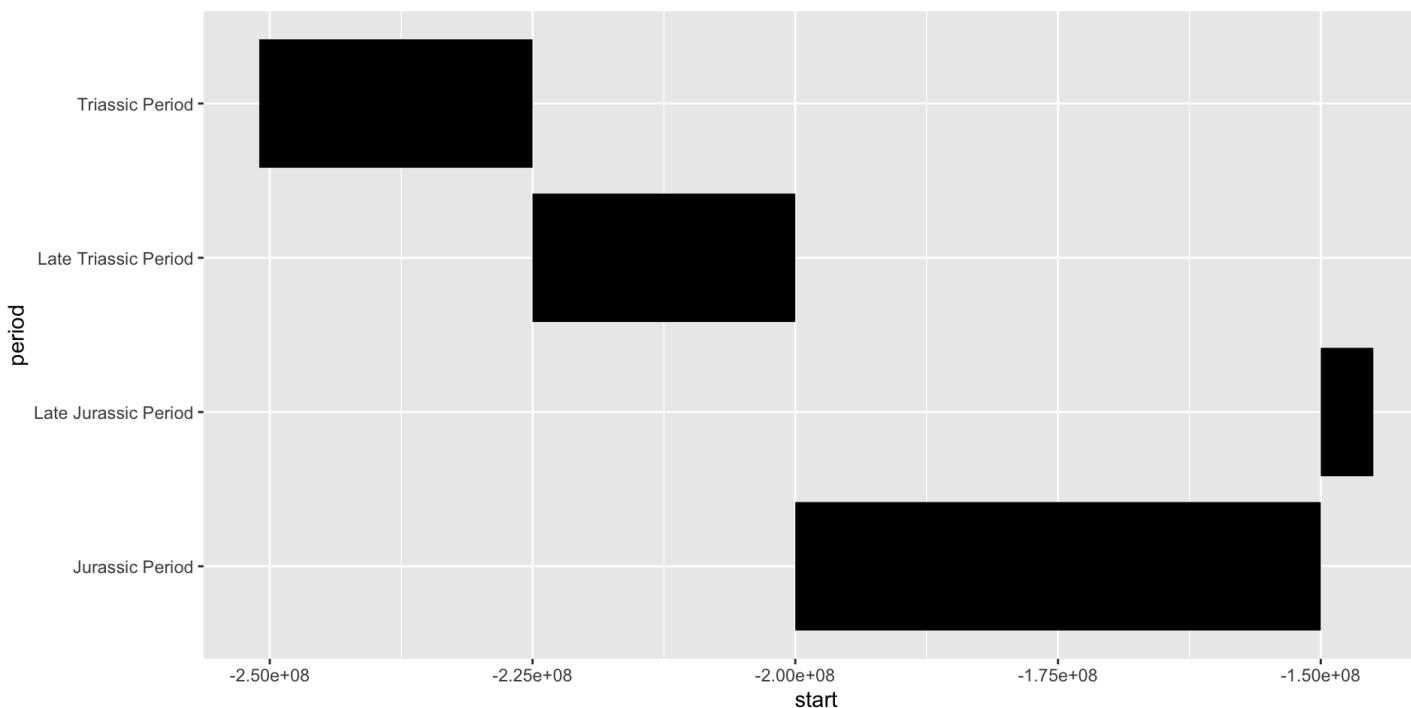
X y

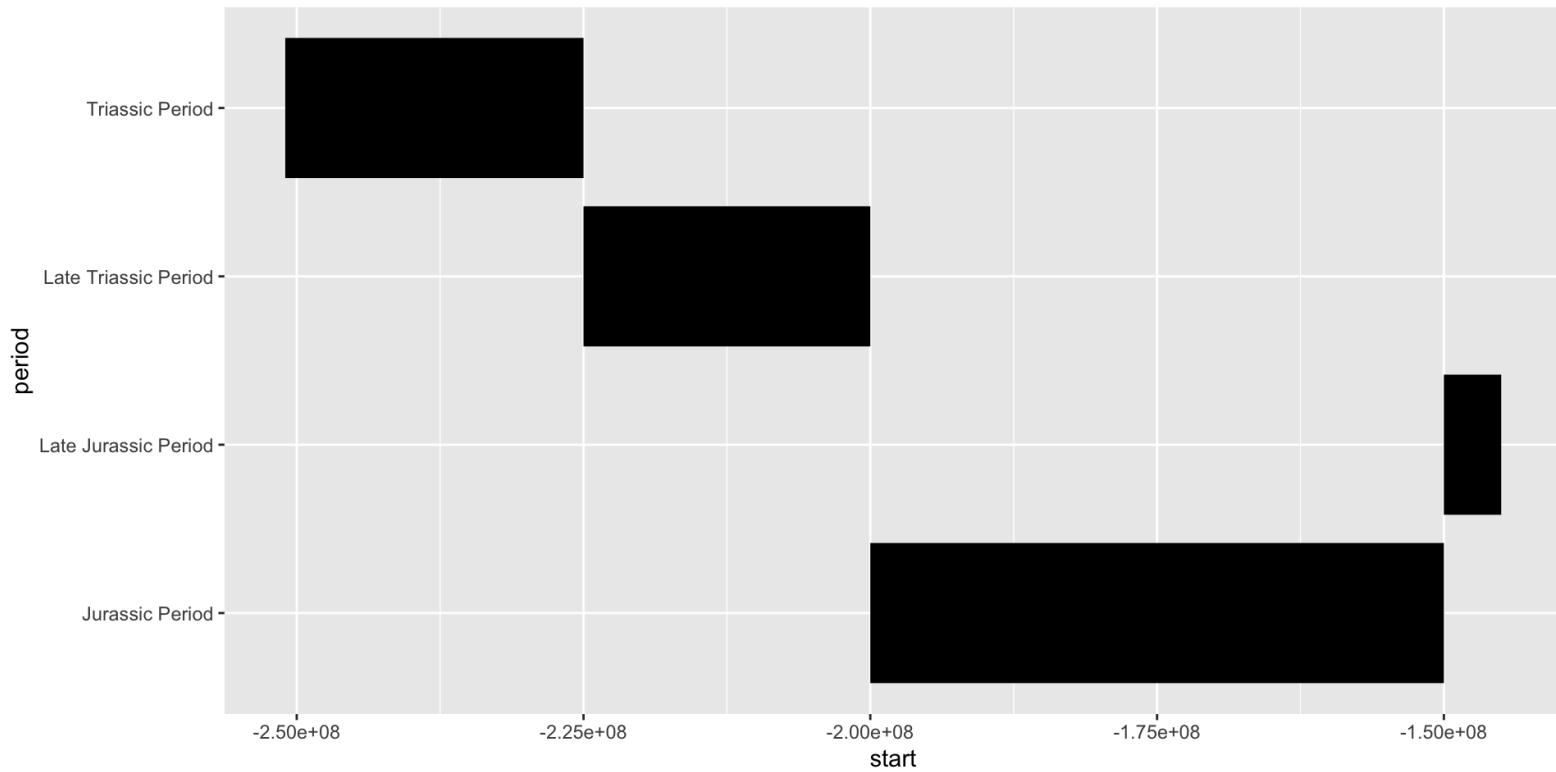
```
1 dinosaurs <- tribble(  
2   ~period, ~start, ~end,  
3   "Triassic Period", -251e6, -225e6,  
4   "Late Triassic Period", -225e6, -200e6,  
5   "Jurassic Period", -200e6, -150e6,  
6   "Late Jurassic Period", -150e6, -145e6  
7 )
```



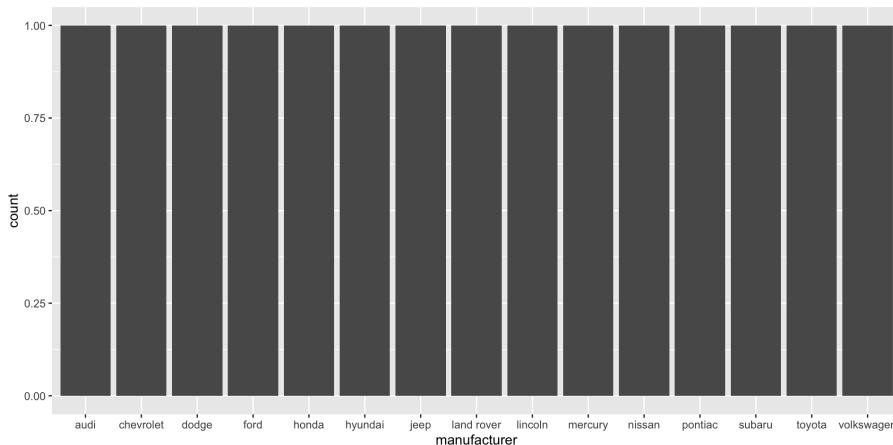
S i z e

```
1 dinosaurs %>%
2   ggplot() +
3   aes(x = start, xend = end,
4        y = period, yend = period) +
5   geom_segment(size = 30)
```

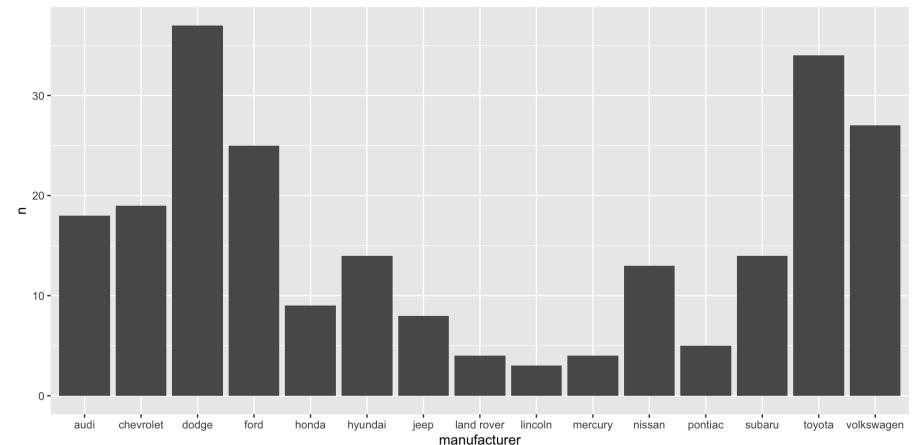




```
1 mpg %>%
2   count(manufacturer) %>%
3   ggplot() +
4   geom_bar(aes(manufacturer))
```

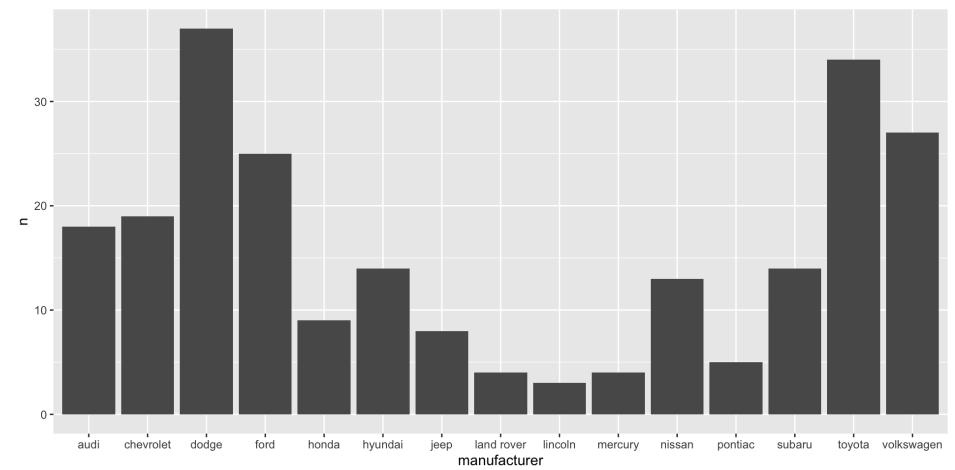


```
1 mpg %>%
2   count(manufacturer) %>%
3   ggplot() +
4   geom_col(aes(x = manufacturer, y = n))
```

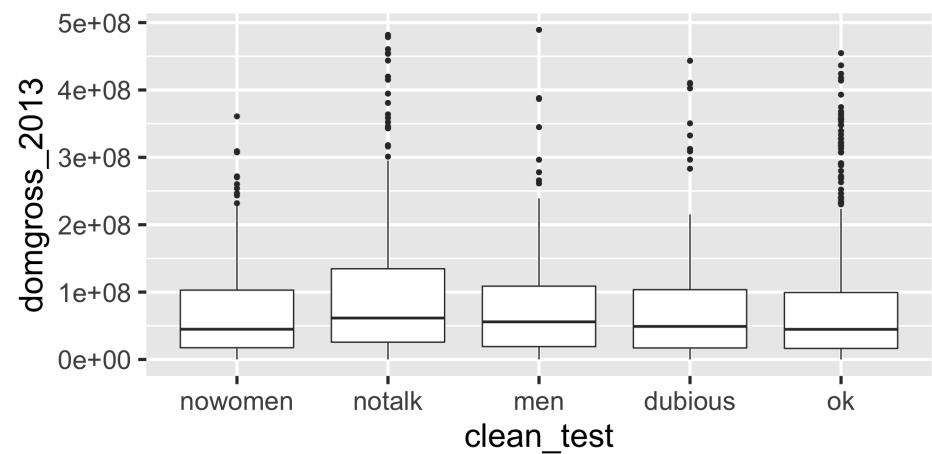


s t a t

```
1 mpg %>%
2   count(manufacturer) %>%
3   ggplot() +
4   geom_bar(aes(x = manufacturer,
5               y = n),
6             stat = "identity")
```



```
1 bechdel %>%
2   filter(complete.cases(.),
3         domgross_2013 < 0.5e9) %>%
4   ggplot(aes(clean_test,
5             domgross_2013)) +
6   geom_boxplot() +
7   theme_gray(base_size = 25)
```



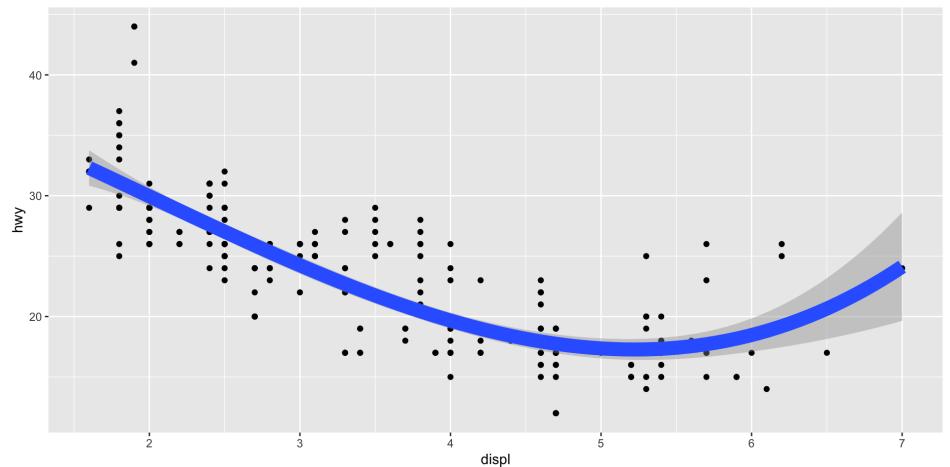
-
-
-

```
1 gss_cat %>%
2   count(relig, marital)

# A tibble: 78 × 3
  relig      marital     n
  <fct>     <fct>     <int>
1 No answer  No answer    4
2 No answer  Never married 22
3 No answer  Separated     3
4 No answer  Di vorced    13
5 No answer  W dowed       7
6 No answer  Married       44
7 Don't know Never married  6
8 Don't know Separated      3
9 Don't know Di vorced     1
10 Don't know Married       5
# ... with 68 more rows
```



```
1 ggplot(mpg,
2       aes(displ, hwy)) +
3     geom_point() +
4     geom_smooth(method = lm,
5                 formula = y ~ splines::bs(x, 3),
6                 size = 5)
```





-
-

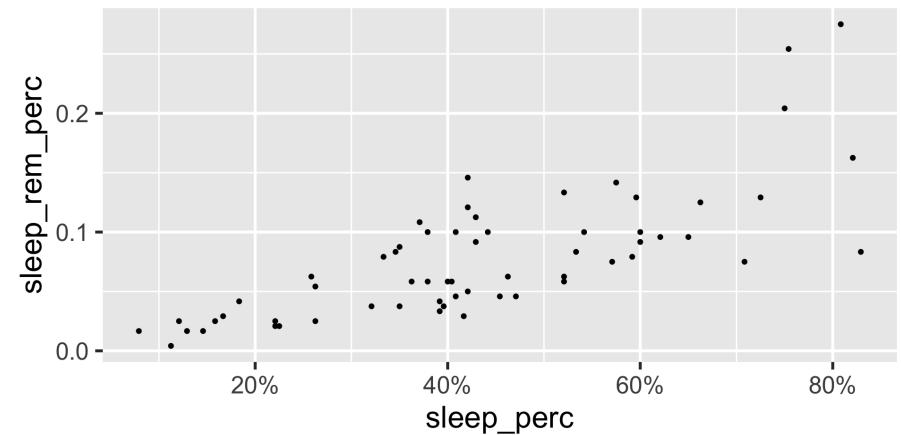
```
1 msl eep %>%
2   ggpl ot() +
3   aes(
4     x = sl eep_total ,
5     y = sl eep_rem
6     colour = vore
7   ) +
8   geom_poi nt() +
9   scale_col our_manual (
10    values = c("carni " = "#c03728",
11                  "omni " = "#fd8f24",
12                  "insecti " = "#f5c04a",
13                  "herbi " = "#919c4c",
14                  "NA" = "#e68c7c")
15  )
```




```
1 percent(c(0.3, 0.5, 0.6))
```

```
[1] "30%" "50%" "60%"
```

```
1 mslleep %>%
 2   mutate(sleep_perc = sleep_total / 24,
 3         sleep_rem_perc = sleep_rem / 24) %>%
 4   ggplot() +
 5   aes(x = sleep_perc,
 6        y = sleep_rem_perc) +
 7   geom_point() +
 8   scale_x_continuous(label = percent_format()) +
 9   theme_gray(base_size = 24)
```



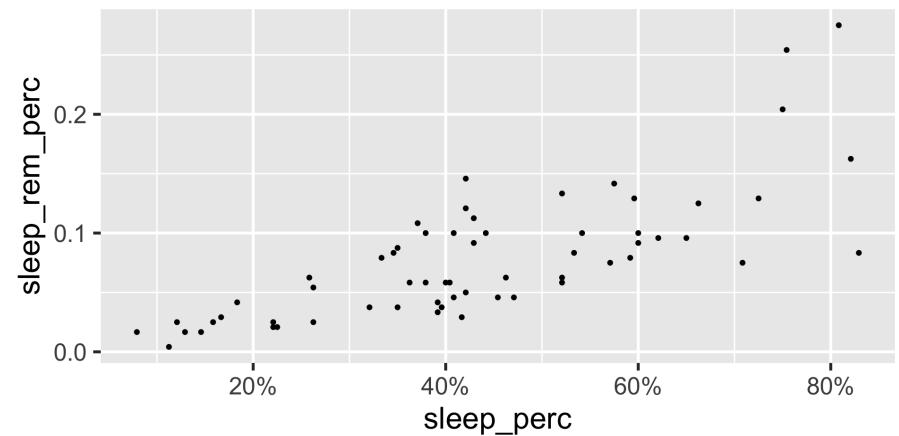
|

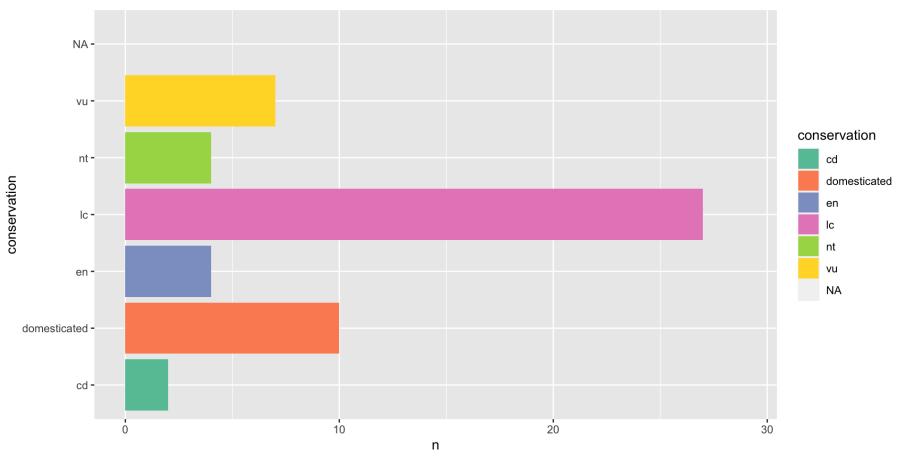


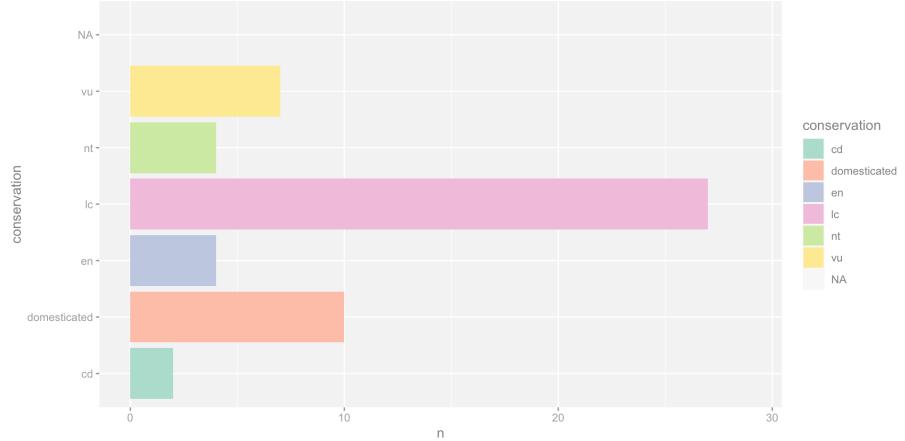
```
1 label_percent()(c(0.3, 0.5, 0.6))
```

```
[1] "30%" "50%" "60%"
```

```
1 mslleep %>%
2   mutate(sleep_perc = sleep_total / 24,
3         sleep_rem_perc = sleep_rem / 24) %>%
4   ggplot() +
5   aes(x = sleep_perc,
6        y = sleep_rem_perc) +
7   geom_point() +
8   scale_x_continuous(label = label_percent()) +
9   theme_gray(base_size = 24)
```



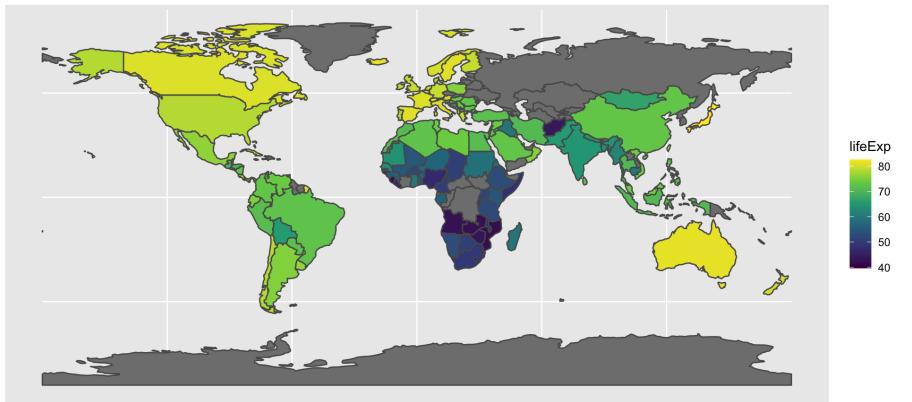




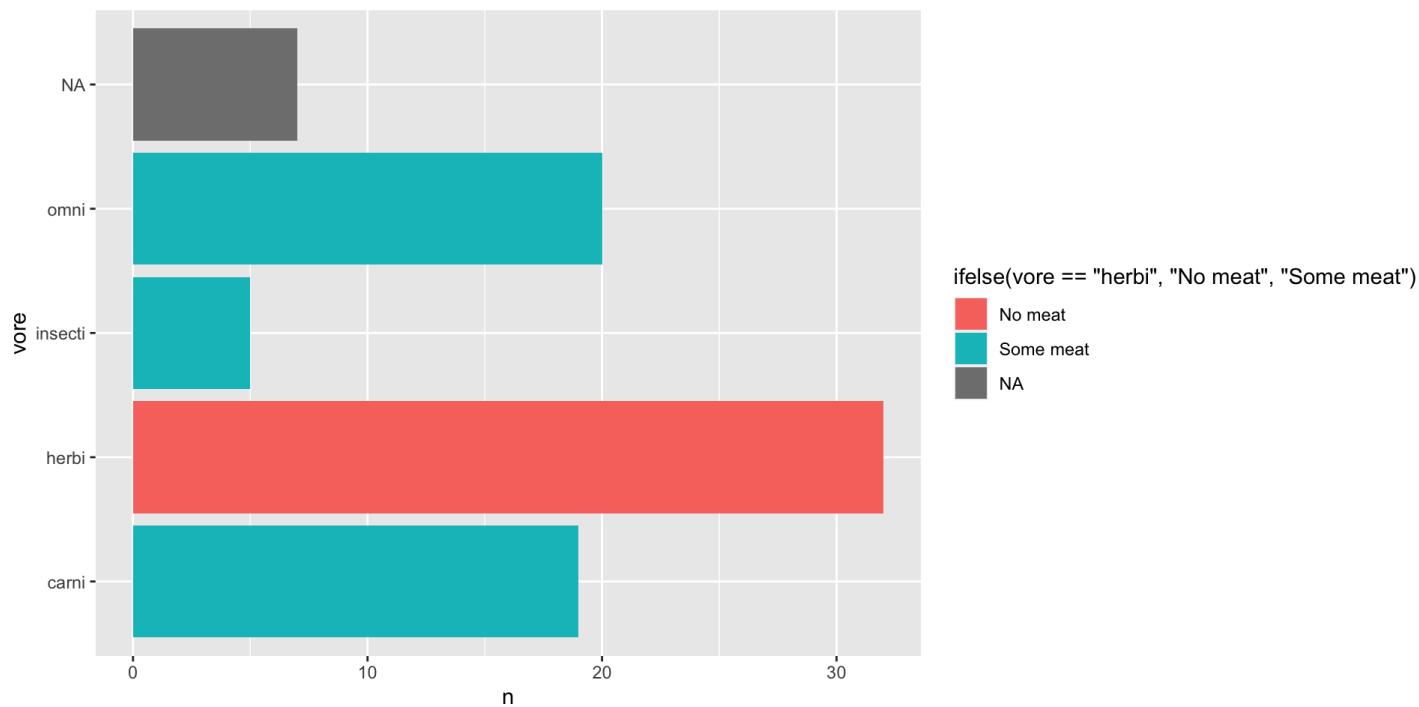
```

1 countries110 %>%
2   st_as_sf() %>%
3   left_join(filter(gapminder, year == 2007),
4             by = c("name" = "country")) %>%
5   ggplot() +
6   geom_sf(aes(fill = lifeExp)) +
7   scale_fill_viridis_c()

```



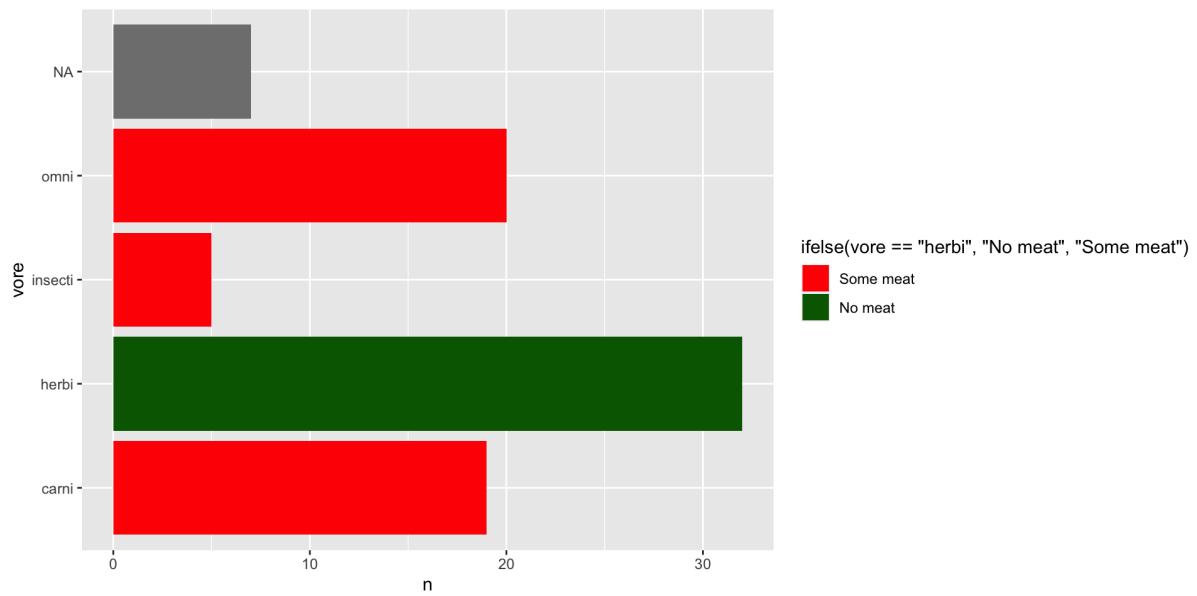
```
1 mslleep %>%
2   count(vore) %>%
3   ggplot() +
4   aes(x = n,
5       y = vore,
6       fill = ifelse(vore == "herbi", "No meat", "Some meat")) +
7   geom_col()
```



```

1 msl %>%
2   count(vore) %>%
3   ggplot() +
4   aes(x = n,
5     y = vore,
6     fill = ifelse(vore == "herbi", "No meat", "Some meat")) +
7   geom_col() +
8   scale_fill_manual(values = c("Some meat" = "red",
9                           "No meat" = "darkgreen"))

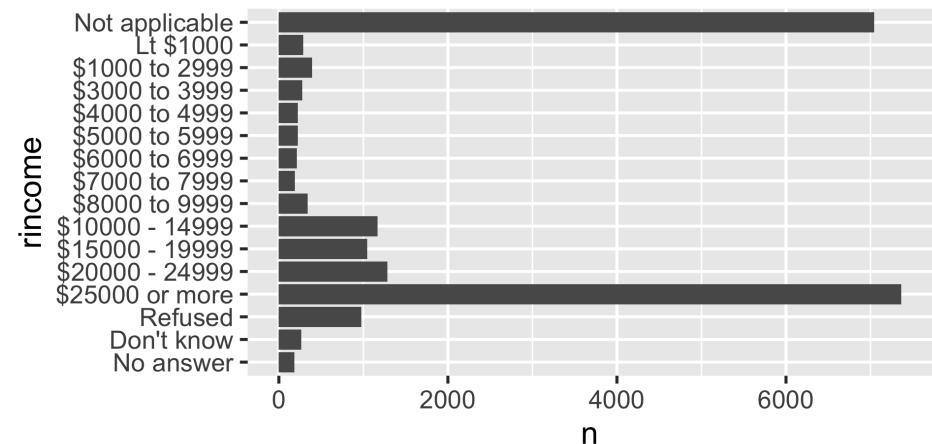
```



```
1 gss_cat %>%
2   head() %>%
3   pull(rincome)
```

```
[1] $8000 to 9999 $8000 to 9999 Not applicable Not applicable Not applicable
[6] $20000 - 24999
16 Levels: No answer Don't know Refused $25000 or more ... Not applicable
```

```
1 gss_cat %>%
2   count(rincome) %>%
3   ggplot() +
4   aes(x = n,
5       y = rincome) +
6   geom_col() +
7   theme_gray(base_size = 24)
```

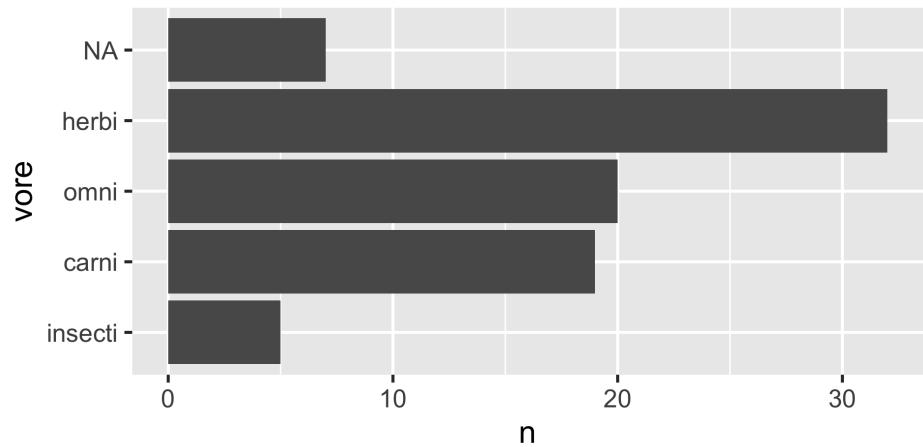



```
1 msl eep %>%
2   count(vore)
```

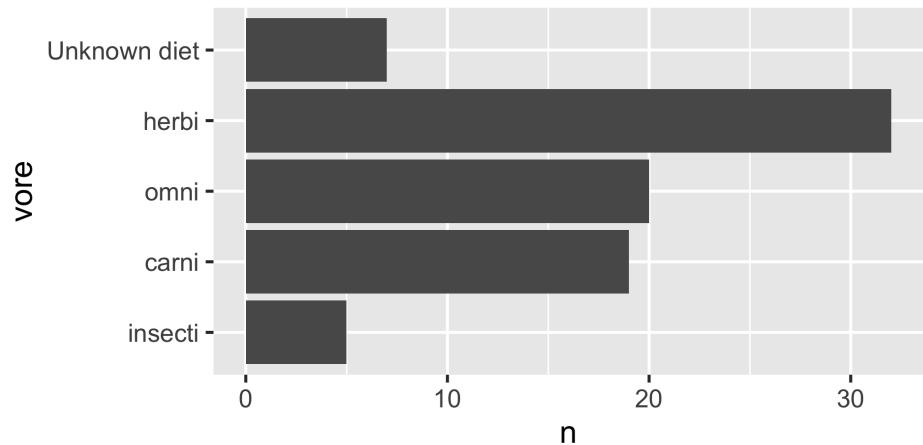
```
# A tibble: 5 × 2
  vore      n
  <chr>    <int>
1 carni     19
2 herbi     32
3 insecti    5
4 omni      20
5 <NA>       7
```



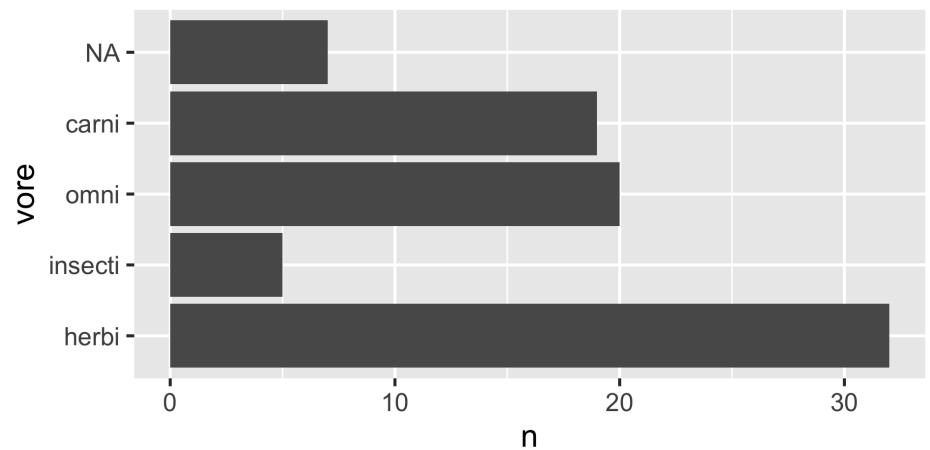
```
1 msl %>%  
2   count(vore) %>%  
3   mutate(vore = fct_reorder(vore, n)) %>%  
4   ggplot() +  
5   aes(x = n,  
6        y = vore) +  
7   geom_col() +  
8   theme_gray(base_size = 24)
```



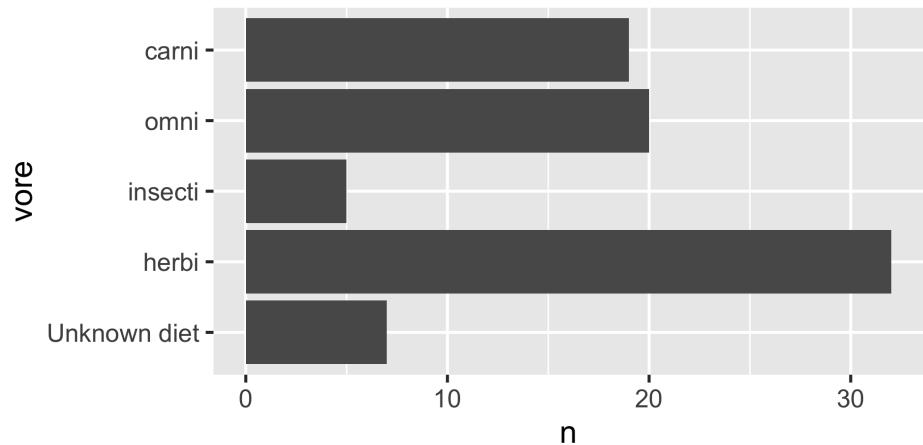
```
1 msl %>%  
2   count(vore) %>%  
3   mutate(vore = fct_reorder(vore, n),  
4         vore = fct_explicit_na(vore, "Unknown diet"))  
5 ggplot() +  
6   aes(x = n,  
7        y = vore) +  
8   geom_col() +  
9   theme_gray(base_size = 24)
```



```
1 order_vore <- c("carni", "omni", "insecti", "herbi")
2
3 msl %>%
4   count(vore) %>%
5   mutate(vore = fct_relevel(vore, order_vore),
6         vore = fct_rev(vore)) %>%
7   ggplot() +
8   aes(x = n,
9        y = vore) +
10  geom_col() +
11  theme_gray(base_size = 24)
```



```
1 msl %>%  
2   count(vore) %>%  
3   mutate(vore = fct_relevel(vore, order_vore),  
4         vore = fct_rev(vore),  
5         vore = fct_expli ci t_na(vore, "Unknown diet"),  
6         vore = fct_relevel(vore, "Unknown diet", after = 1),  
7   ggplot() +  
8   aes(x = n,  
9        y = vore) +  
10  geom_col() +  
11  theme_gray(base_size = 24)
```





|

```
1 download.file("https://raw.githubusercontent.com/charliedaley/eng7218_data-science-for-healthcare-applications_bcu-  
2           destfile = "data/global-burden-of-disease-data.csv")
```





•

•

•

•

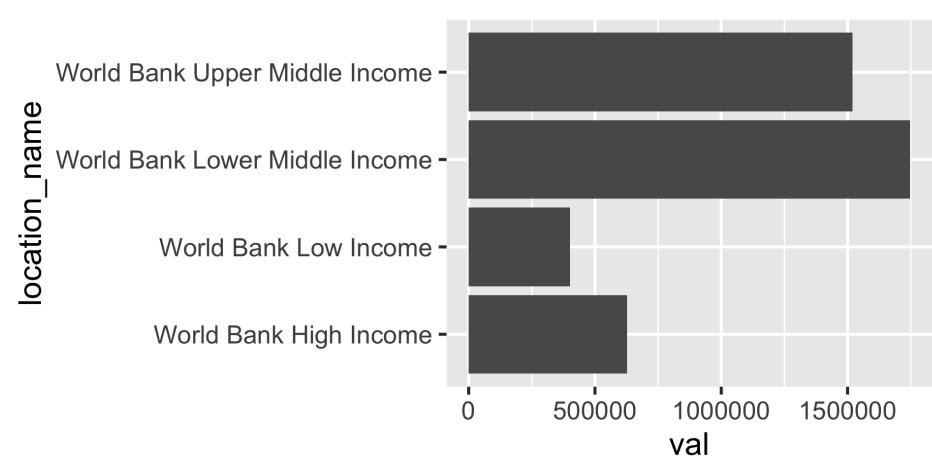
•



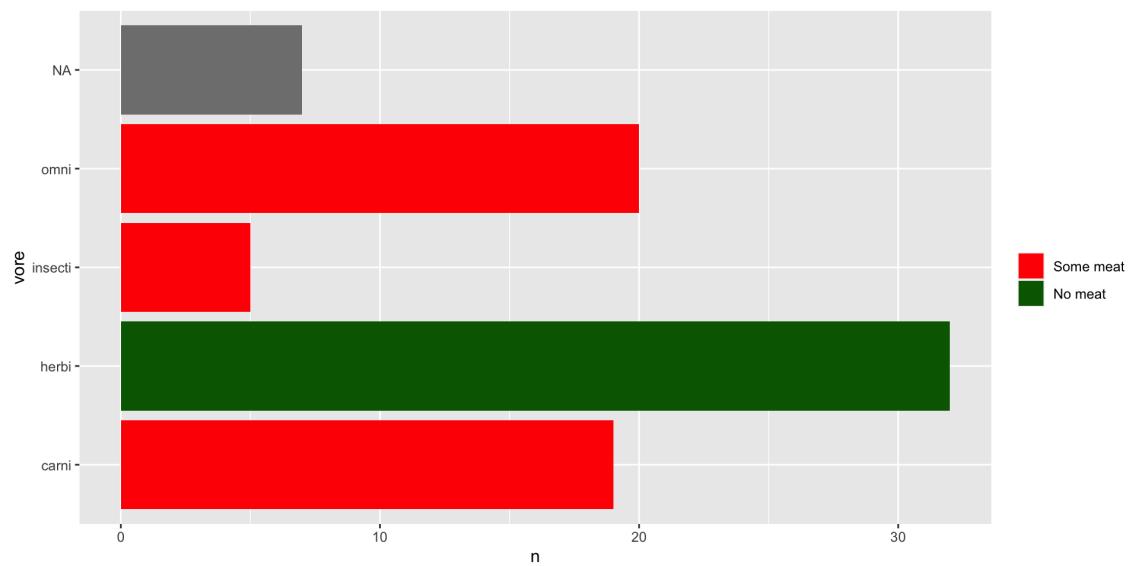


-
-

```
1 gdb_inj_uris %>%
2   ggplot() +
3   aes(x = val,
4       y = location_name) +
5   geom_col() +
6   theme_gray(base_size = 24)
```



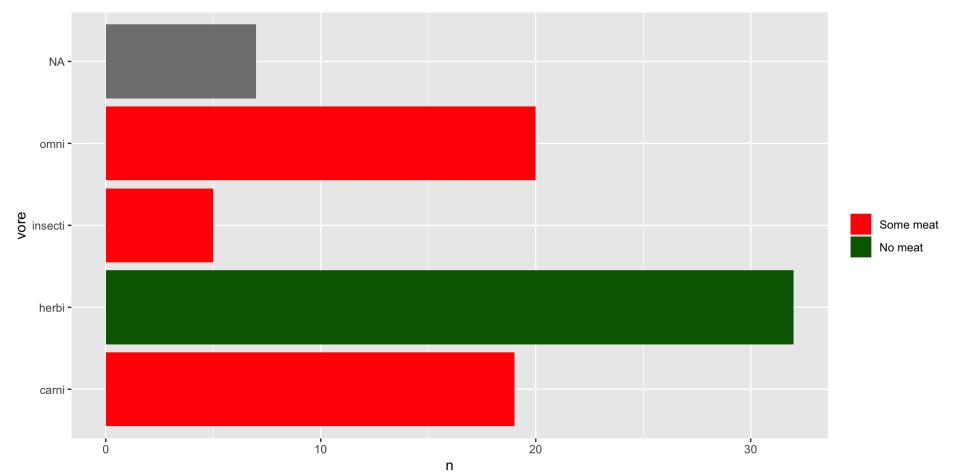
```
1 ggplot() +  
2 geom_line(showLegend = FALSE) +  
3 guides(alpha = guide_legend()) +  
4 theme(legend.position = "bottom")
```



```

1 msl %>%
2   count(vore) %>%
3   ggplot() +
4   aes(x = n,
5     y = vore,
6     fill = ifelse(vore == "herbi",
7                   "No meat",
8                   "Some meat")) +
9   geom_col() +
10  scale_fill_manual(values = c("Some meat" = "red",
11                      "No meat" = "darkgreen"))
12  name = "")

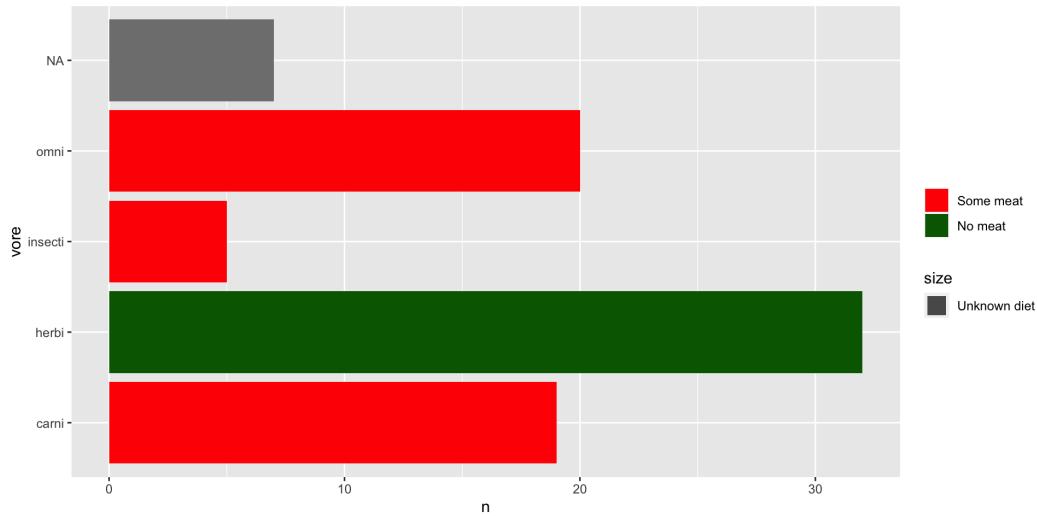
```



```

1 msl %>%
2   count(vore) %>%
3   ggplot() +
4   aes(x = n,
5     y = vore,
6     fill = ifelse(vore == "herbi", "No meat", "Some meat")) +
7   geom_col(aes(size = "Unknown diet")) +
8   scale_fill_manual(values = c("Some meat" = "red",
9                             "No meat" = "darkgreen"),
10                      name = "")

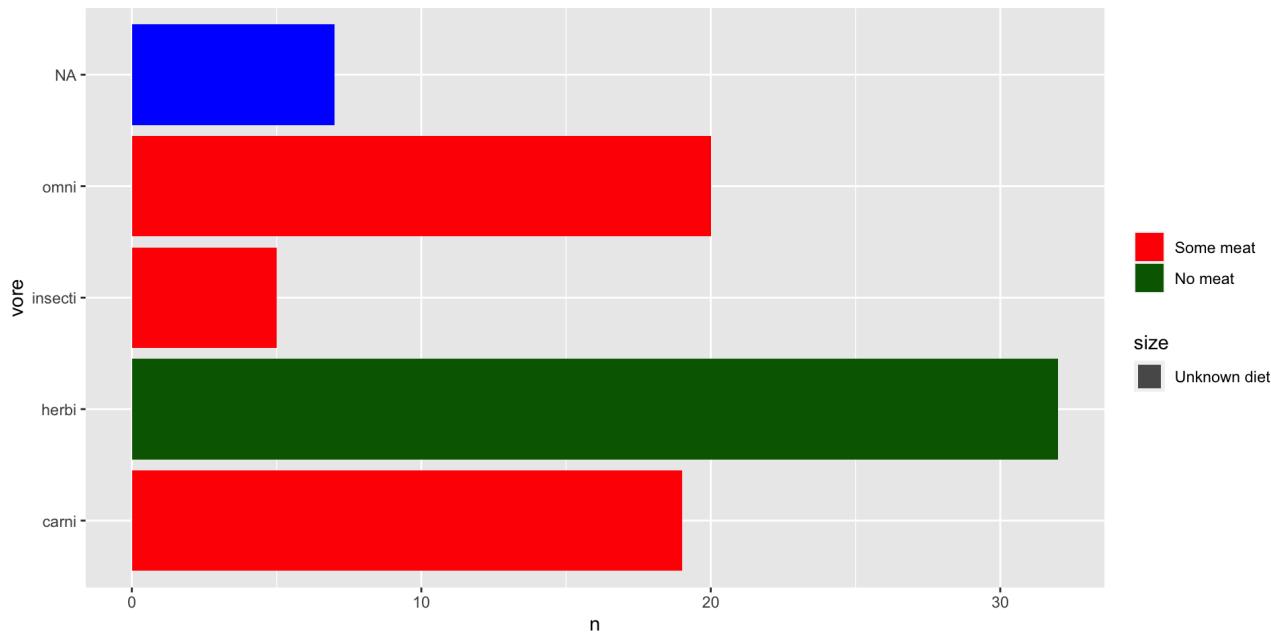
```



```

1 msl %>%
2   count(vore) %>%
3   ggplot() +
4   aes(x = n,
5     y = vore,
6     fill = ifelse(vore == "herbi", "No meat", "Some meat")) +
7   geom_col(aes(size = "Unknown diet")) +
8   scale_fill_manual(values = c("Some meat" = "red",
9                             "No meat" = "darkgreen"),
10                      name = "",
11                      na.value = "blue")

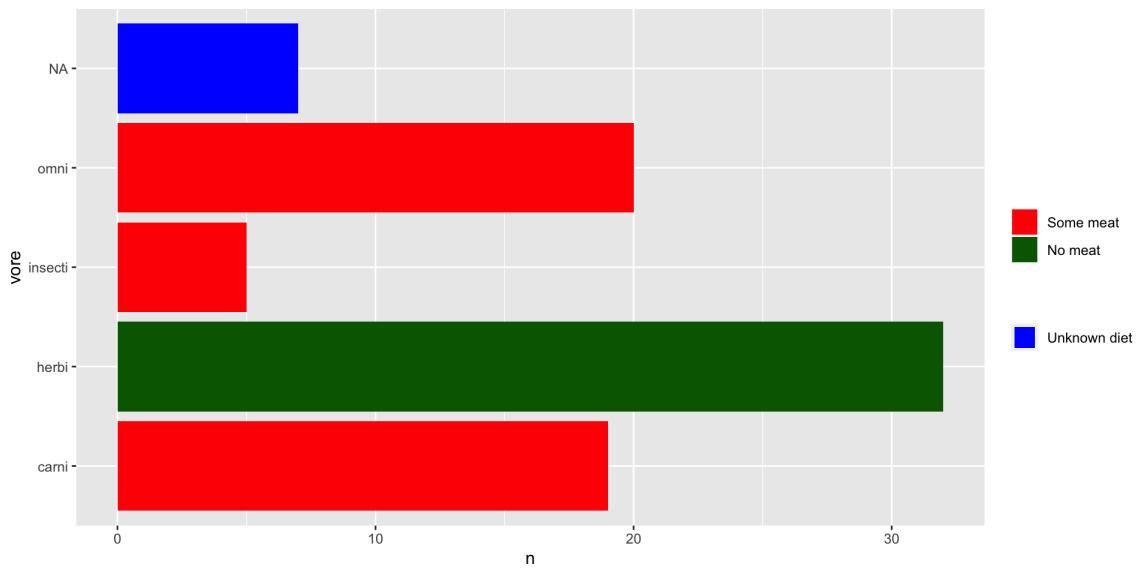
```



```

1 msl %>%
2   count(vore) %>%
3   ggplot() +
4   aes(x = n,
5     y = vore,
6     fill = ifelse(vore == "herbi", "No meat", "Some meat")) +
7   geom_col(aes(size = "Unknown diet")) +
8   scale_fill_manual(values = c("Some meat" = "red",
9                           "No meat" = "darkgreen"),
10                      name = "",
11                      na.value = "blue") +
12   guides(size = guide_legend(title = "",
13                     override.aes = list(fill = "blue")))

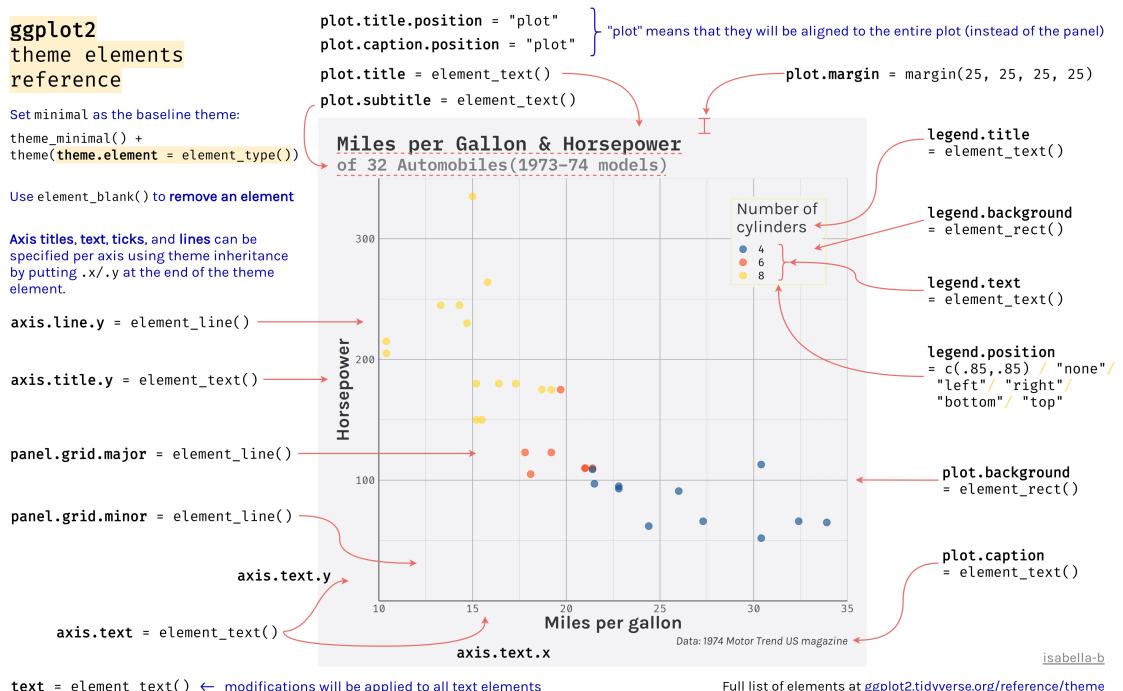
```



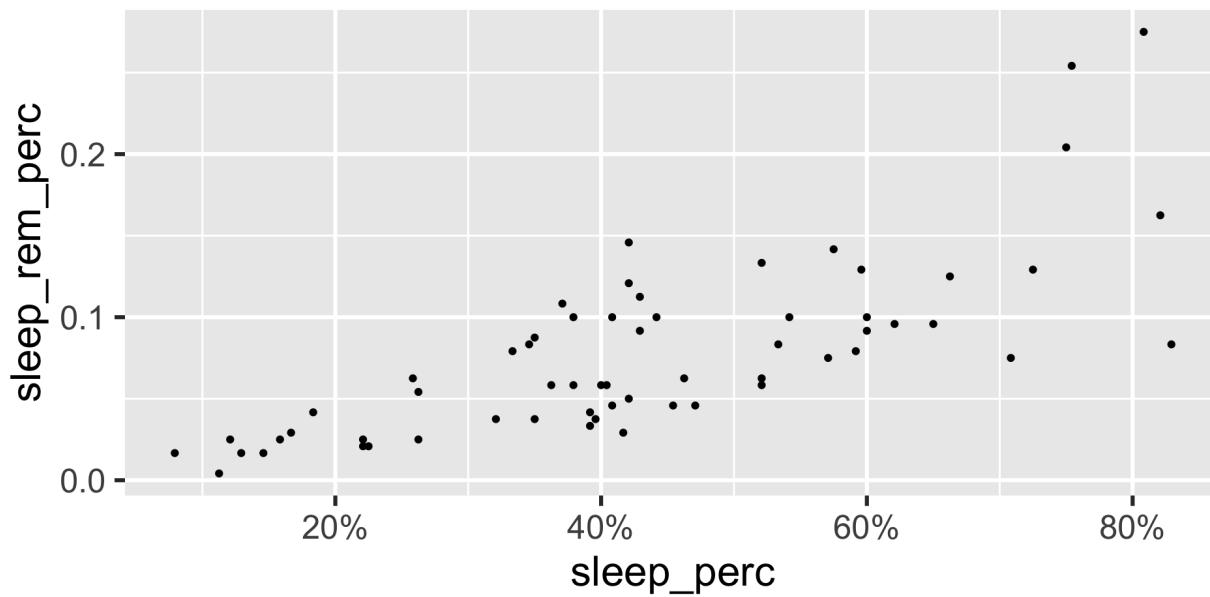
•

•





```
1 mslleep %>%
2   mutate(sleep_perc = sleep_total / 24,
3         sleep_rem_perc = sleep_rem / 24) %>%
4   ggplot() +
5   aes(x = sleep_perc,
6        y = sleep_rem_perc) +
7   geom_point() +
8   scale_x_continuous(label = label_percent()) +
9   theme_gray(base_size = 24)
```



```
1 theme_fivethirtyeight() +  
2   theme(panel.grid.major = element_line(colour = "red"))
```



•

•

•



References

Review Quarterly

Management

The American Statistician 27

Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems

On the mode of communication of cholera

Notes on Matters Affecting the Health Efficiency and Hospital Administration of the British Army

The best stats you've ever seen

EuroVis 2016 - Short Papers

Journal of the American Statistical Association 107(4)

Proceedings of the 28th international conference on Human factors in computing systems - CHI '00

ISPRS International Journal of Geo-Information 10(1)

Applications 42

IEEE Computer Graphics and



Twitter

