











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: $y = 3 + 0.5 x$

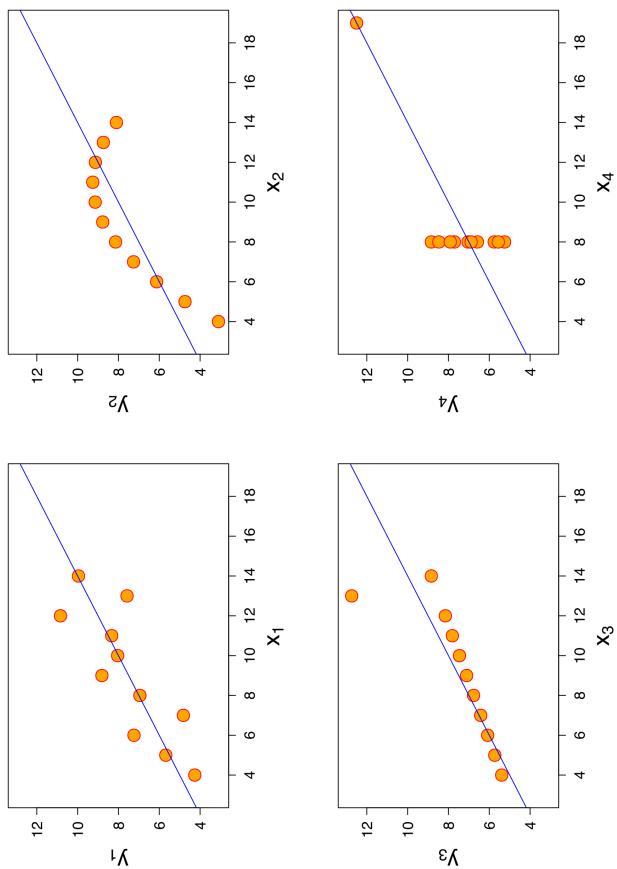
Sum of squares of $x - \bar{x} = 110.0$

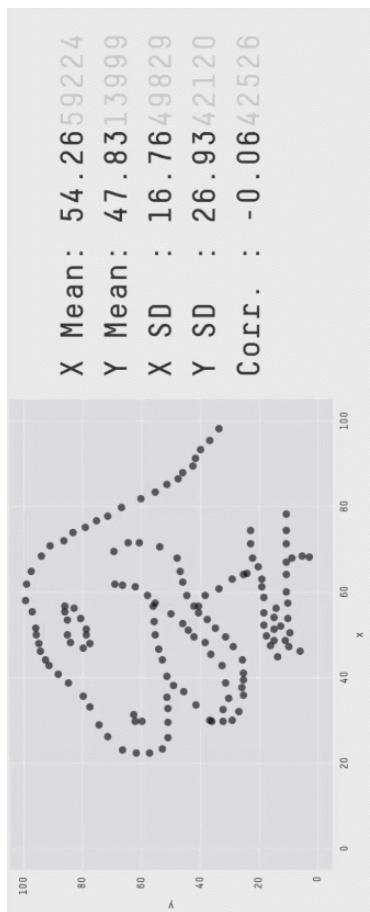
Regression sum of squares = 27.50 (1 d.f.)

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

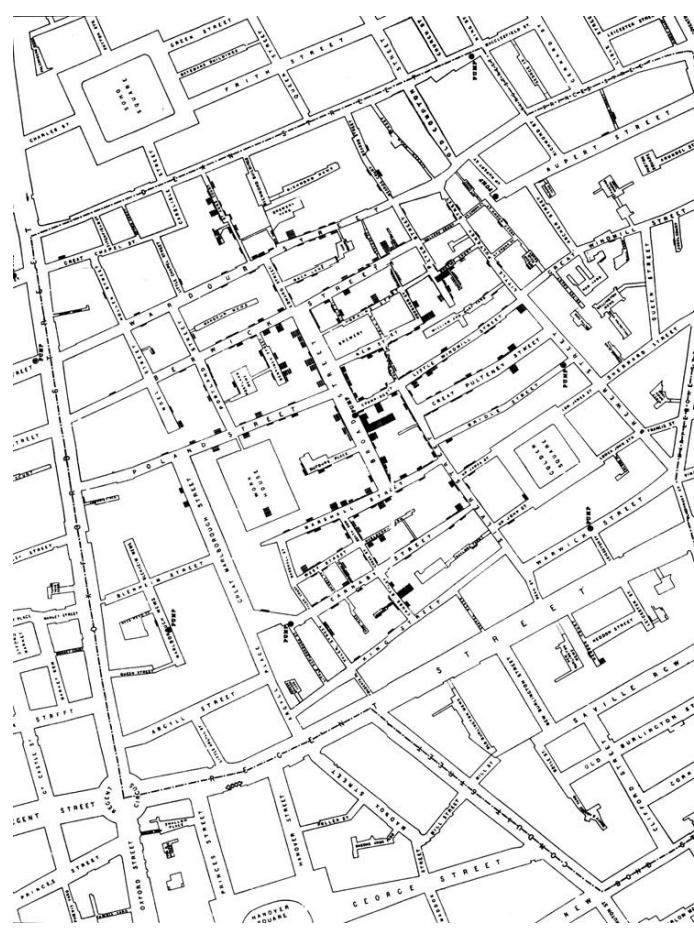
Estimated standard error of $b_1 = 0.118$

Multiple $R^2 = 0.667$

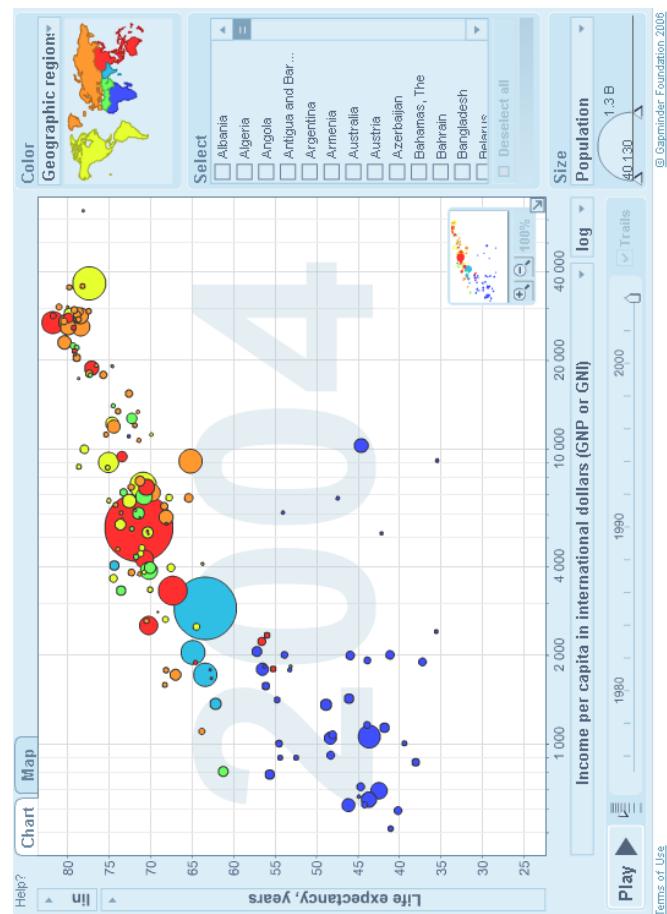


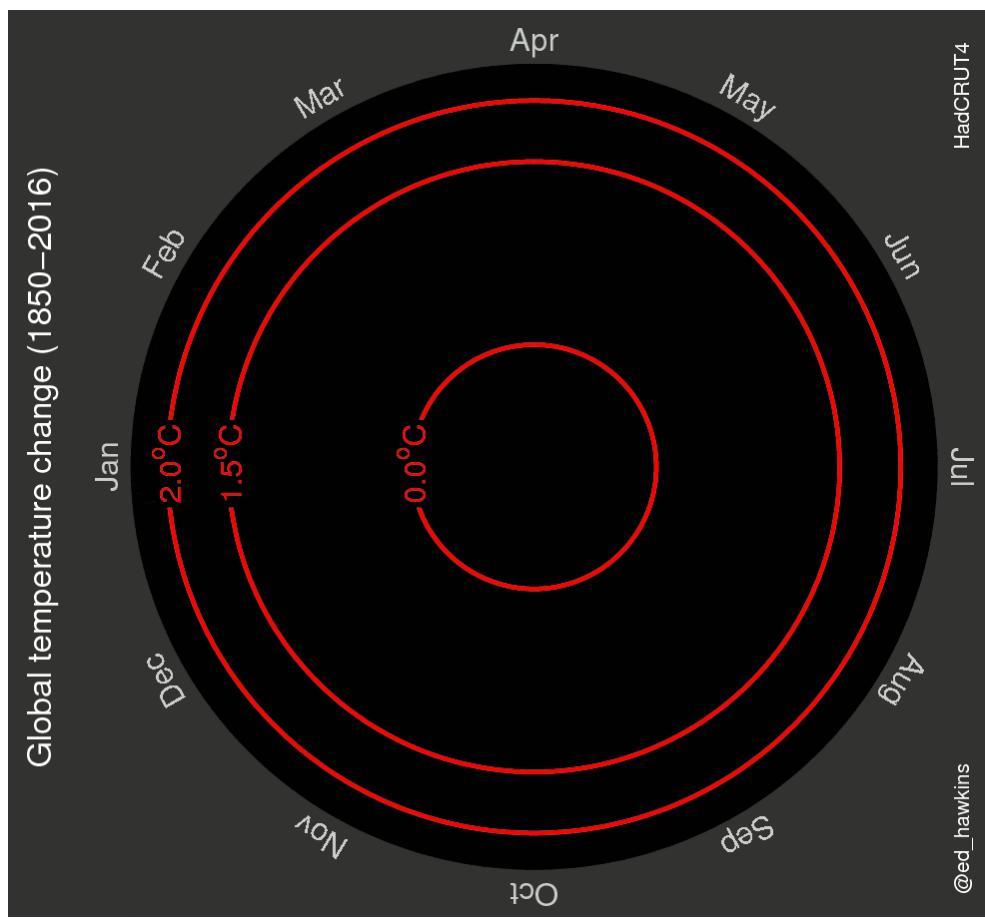




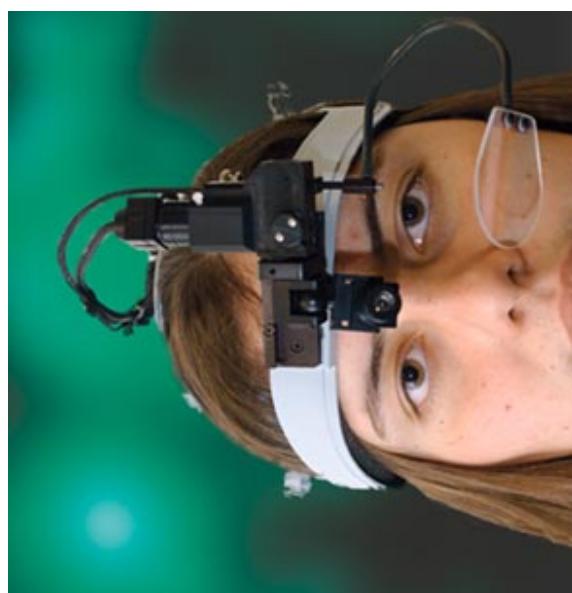












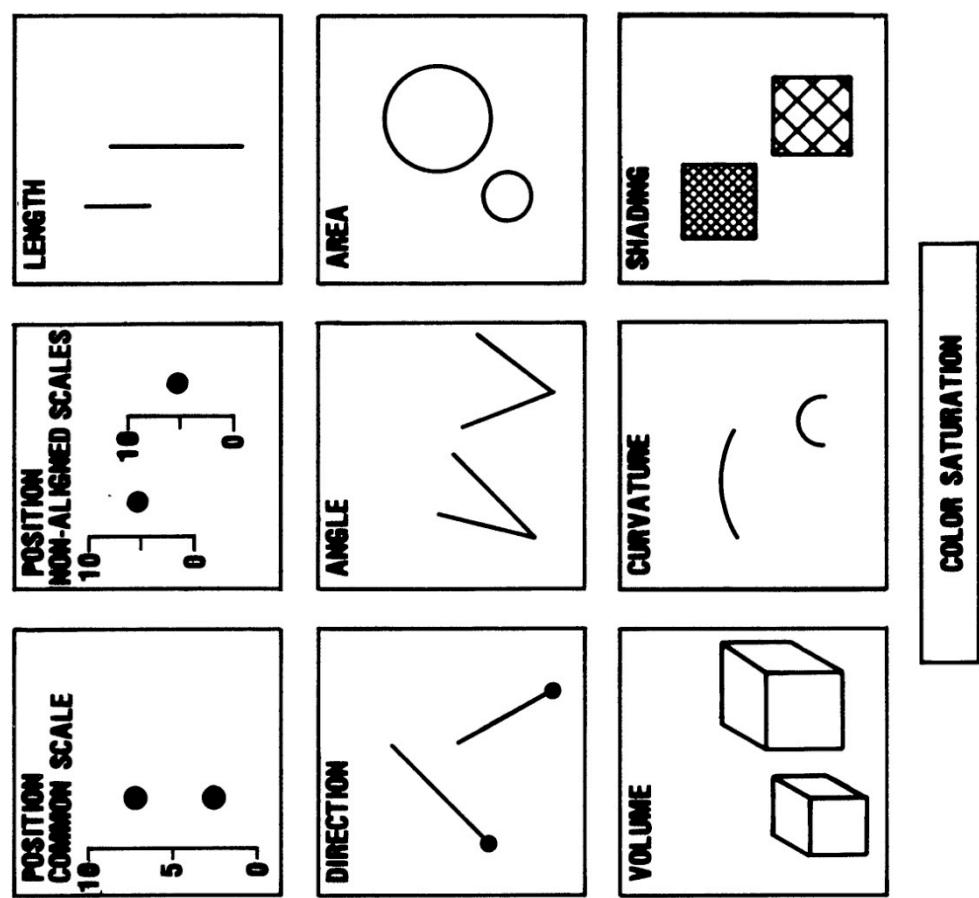


Figure 1. Elementary perceptual tasks.

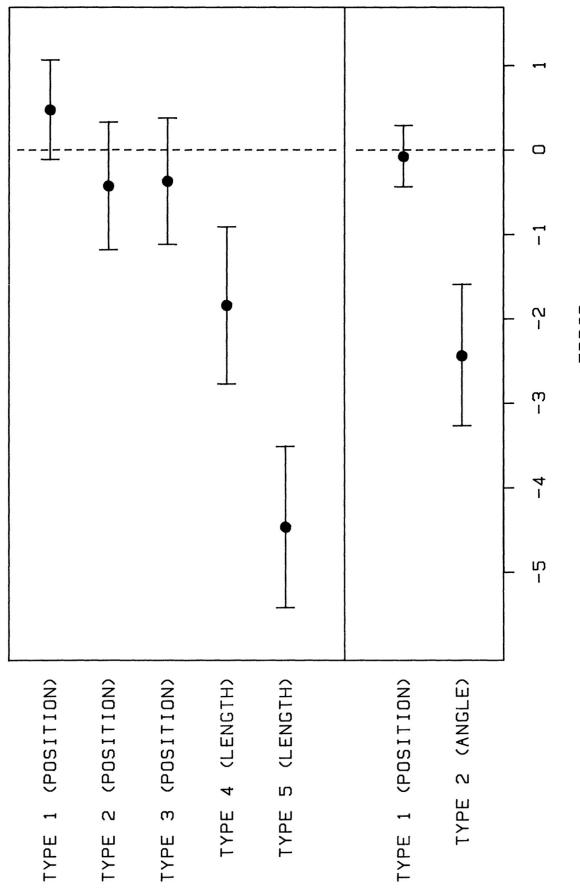


Figure 20. Error means and 95% confidence intervals for judgment types in position-length experiment (top) and position-angle experiment (bottom).

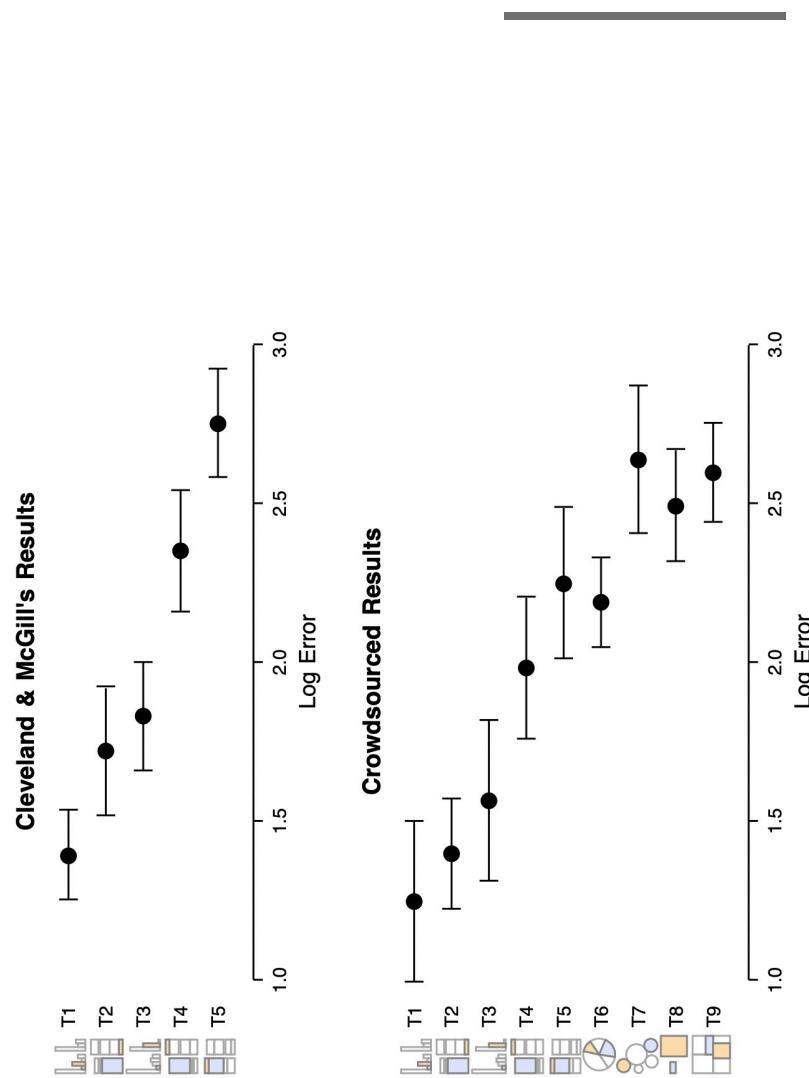
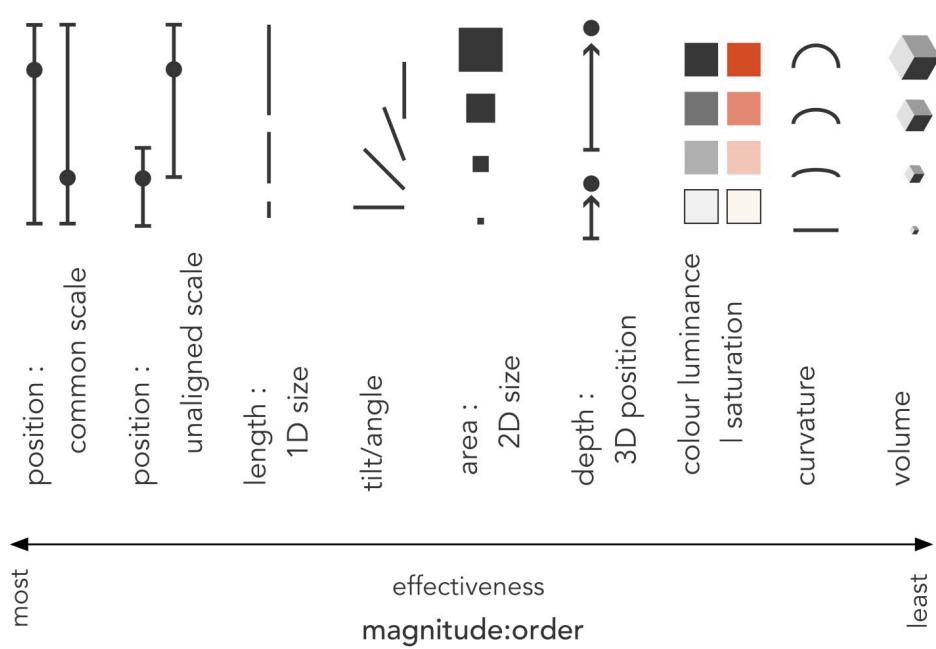
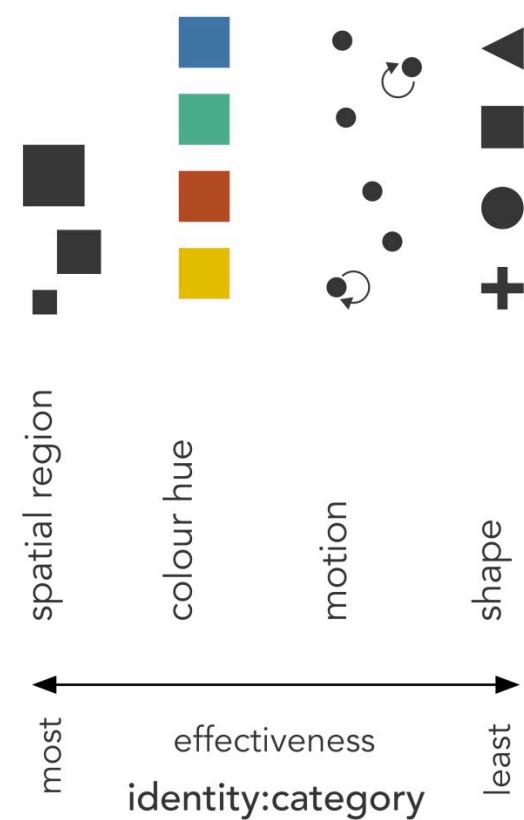


Figure 4: Proportional judgment results (Exp. 1A & B).
Top: Cleveland & McGill's [7] lab study. Bottom: MTurk studies. Error bars indicate 95% confidence intervals.



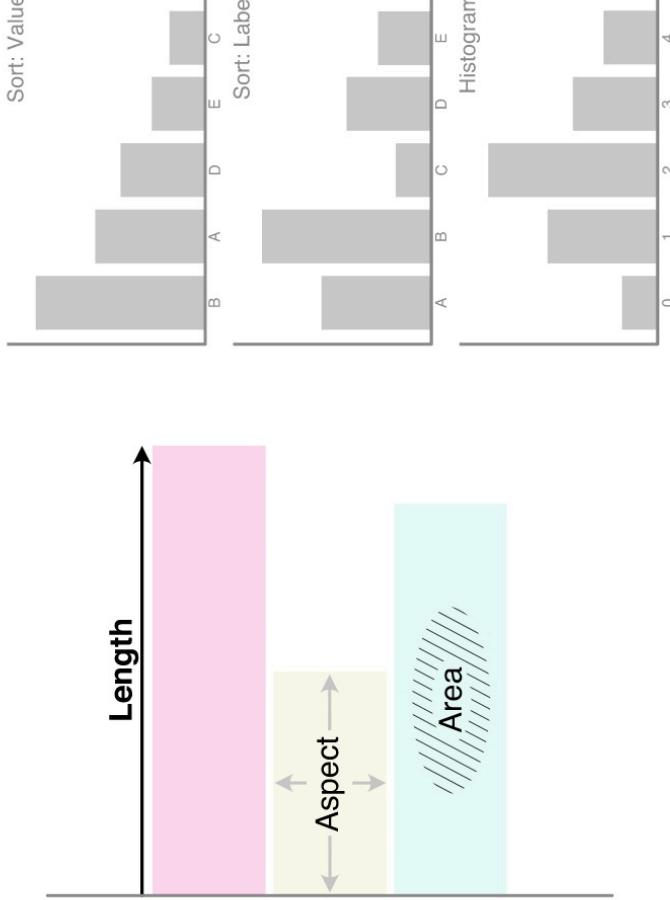


FIGURE 2. Bar charts encode values as their length, but also their area, aspect ratio, and overall shape. Sorting is often used for specific use cases.

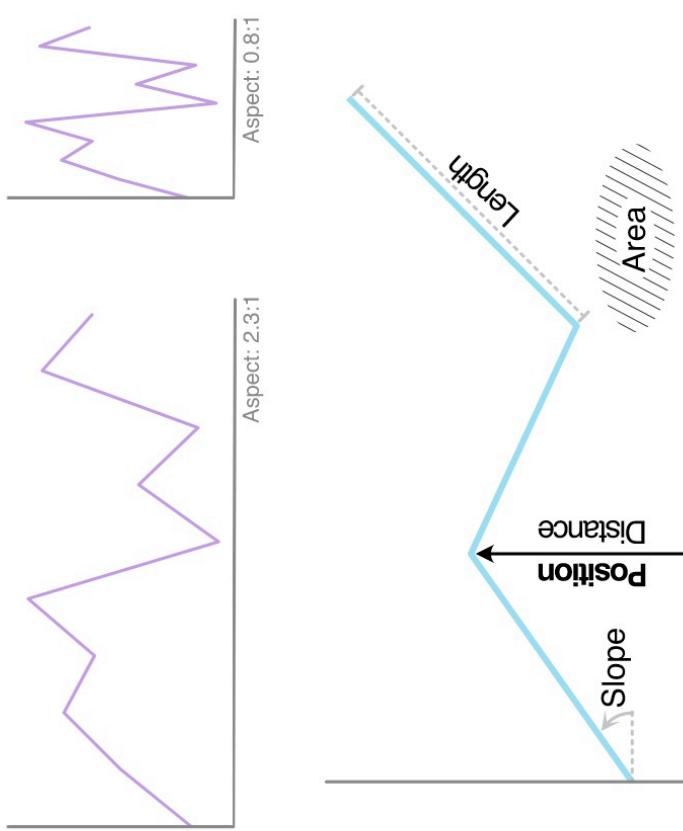


FIGURE 3. Line charts are specified by the location of the points connected by lines, but are read as slope and length, as well as area. Aspect ratio of the chart is also generally considered important.

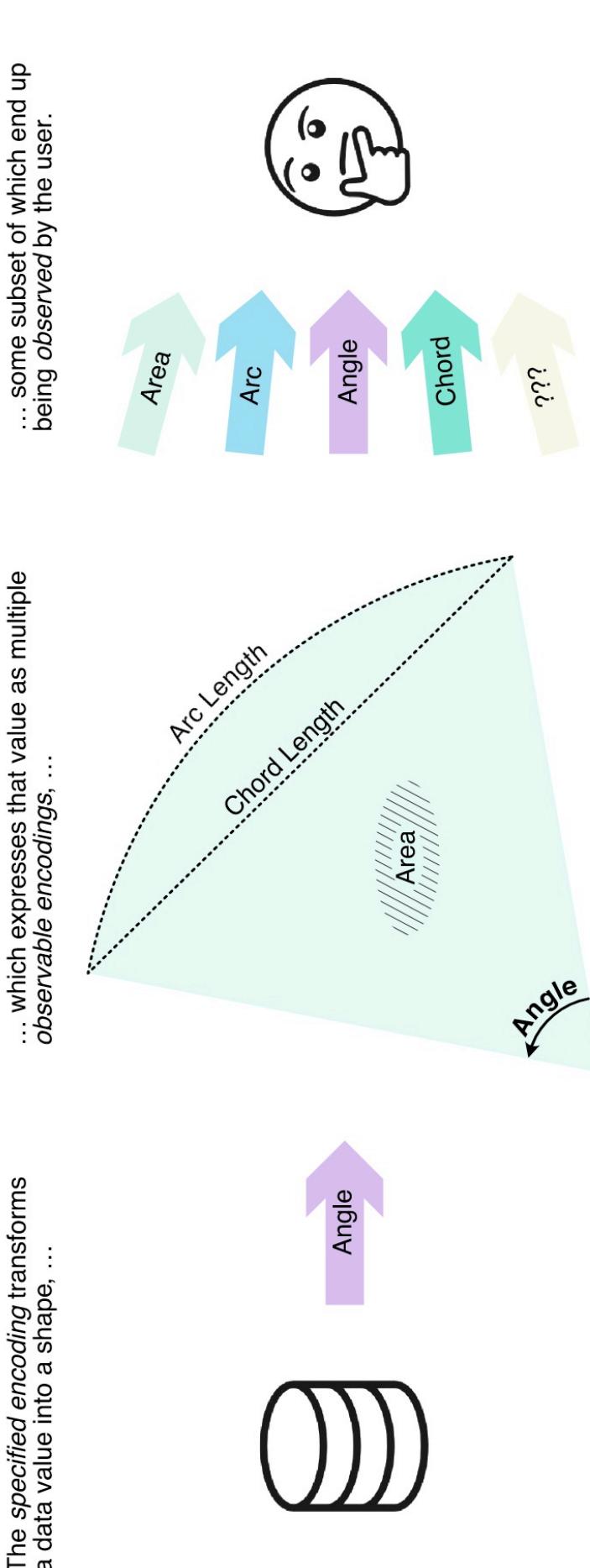
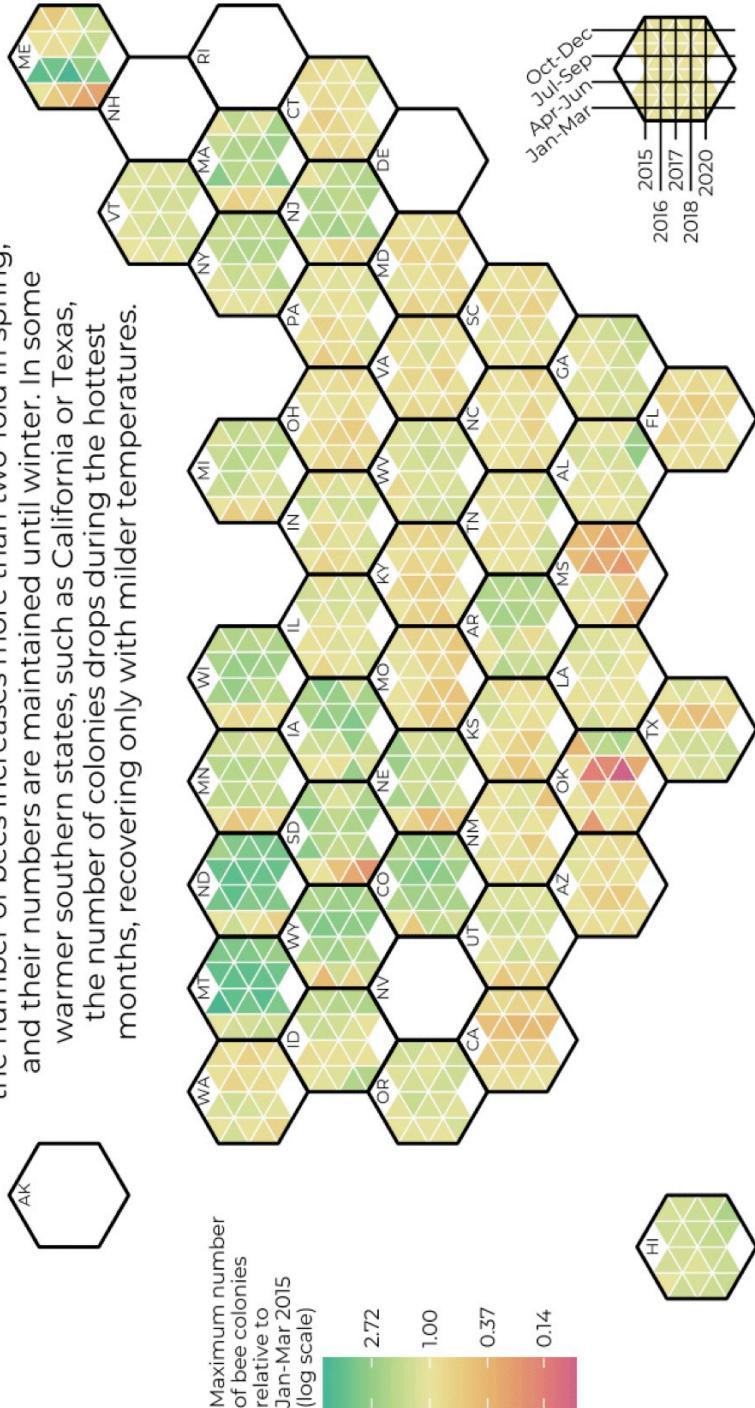


FIGURE 1. Pie charts are specified by angle, but may also be read by area, arc length, or even chord. Shape recognition is also likely for specific angles like $90^\circ/25\%$ and $180^\circ/50\%$.

Data from USDA
Graph by @irg-bio

Seasonality in the number of bee colonies

The number of bee colonies in the US changes with the seasons depending on the state. In some northern states, the number of bees increases more than two-fold in spring, and their numbers are maintained until winter. In some warmer southern states, such as California or Texas, the number of colonies drops during the hottest months, recovering only with milder temperatures.



—

To extract accruals



To quantify volatility



To find the larges t,



To find unusual



You have a story you want to tell

The reader wants

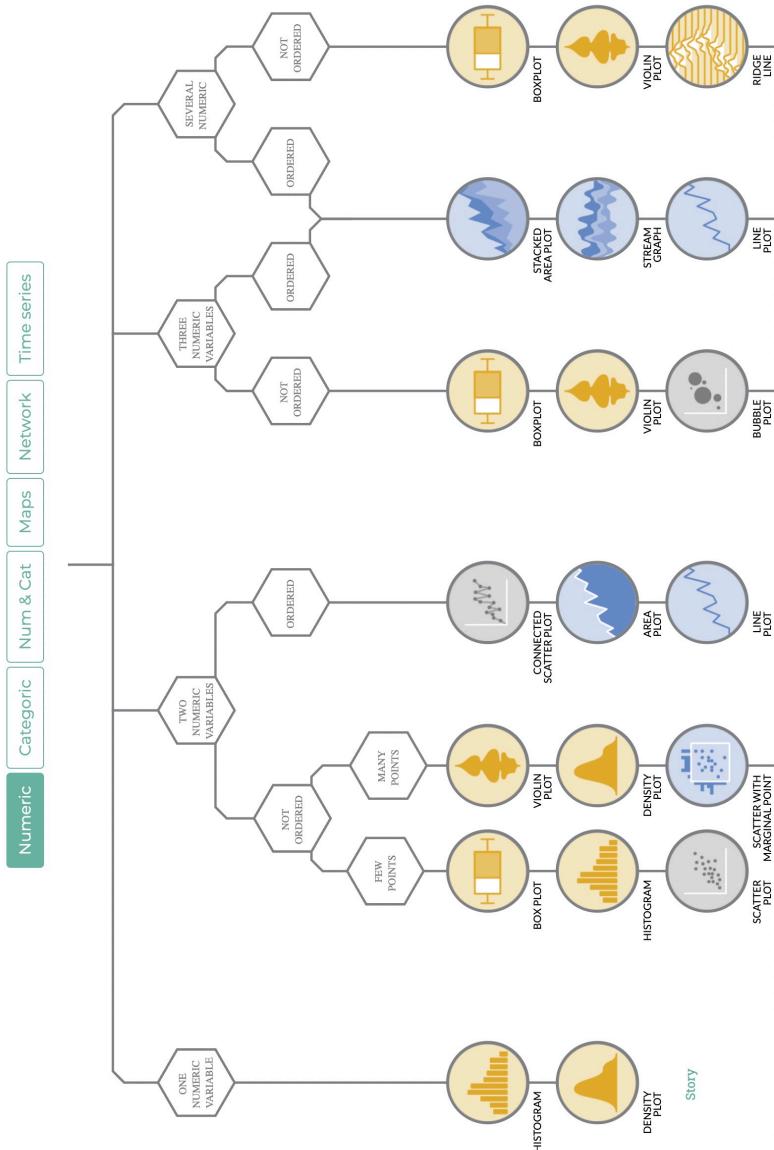
The reader has a preconception about





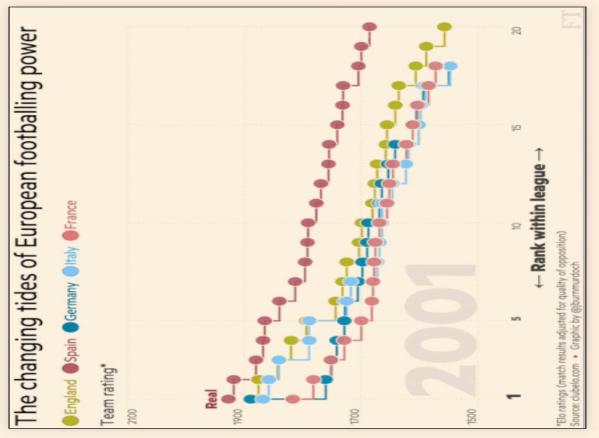
	Ranking	Magnitude	Flow							
> msleep										
# A tibble: 83 × 11										
name	genus	vore	order	conse... ¹	sleep... ²	sleep... ³	sleep... ⁴	awake	brainwt	bodywt
<chr>	<chr>	<chr>	<chr>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>	<dbl>
1 Cheetah	Acin... carn	Carni	Carni	12.1	NA	11.9	NA	50		
2 Owl	mo... Autus	omni	Prim...	17	1.8	NA	7	0.0155	0.48	
3 Mountain...	Aplo...	herbi	Rode...	nt	14.4	2.4	NA	9.6	NA	1.35
4 Great...	Bla...	omi	Sori...	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	Pilo...	NA	14.4	2.2	0.767	9.6	NA	3.85
7 Northe...	Call...	carni	Carni	Vu	8.7	1.4	0.383	15.3	NA	20.5
8 Vesper...	Calo...	NA	Rode...	NA	7	NA	17	NA	0.045	
9 Dog	Canis	carni	Carni	domest...	10.1	2.9	0.333	13.9	0.07	14
10 Roe	de...	Capr...	herbi	Arti...	3	NA	NA	21	0.0982	14.8
# ...	with 73 more rows,	and abbreviated variable names	¹ conservation, ² sleep_total,							
#	³ sleep_rem, ⁴ sleep_cycle									
# i Use `print(n = ...)` to see more rows										

```
> 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> <dbl>
1 Cheetah Acin... carn Carni 12.1 NA    11.9 NA    50
2 Owl     mo... Autus omni Prim... 17   1.8  NA    7     0.0155 0.48
3 Mountain... Aplo... herbi Rode... nt  14.4 2.4  NA   9.6  NA    1.35
4 Great... Bla... omni Sori... 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 Pilo... NA   14.4 2.2  0.767 9.6  NA    3.85
7 Northe... Call... carn... Vu   8.7   1.4  0.383 15.3 NA    20.5
8 Vesper... Calo... NA   Rode... NA   7    NA   17   NA   0.045
9 Dog      Canis carn... Carni domes... 10.1 2.9  0.333 13.9 0.07 14
10 Roe    de... Capr... herbi Arti... 3    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
```

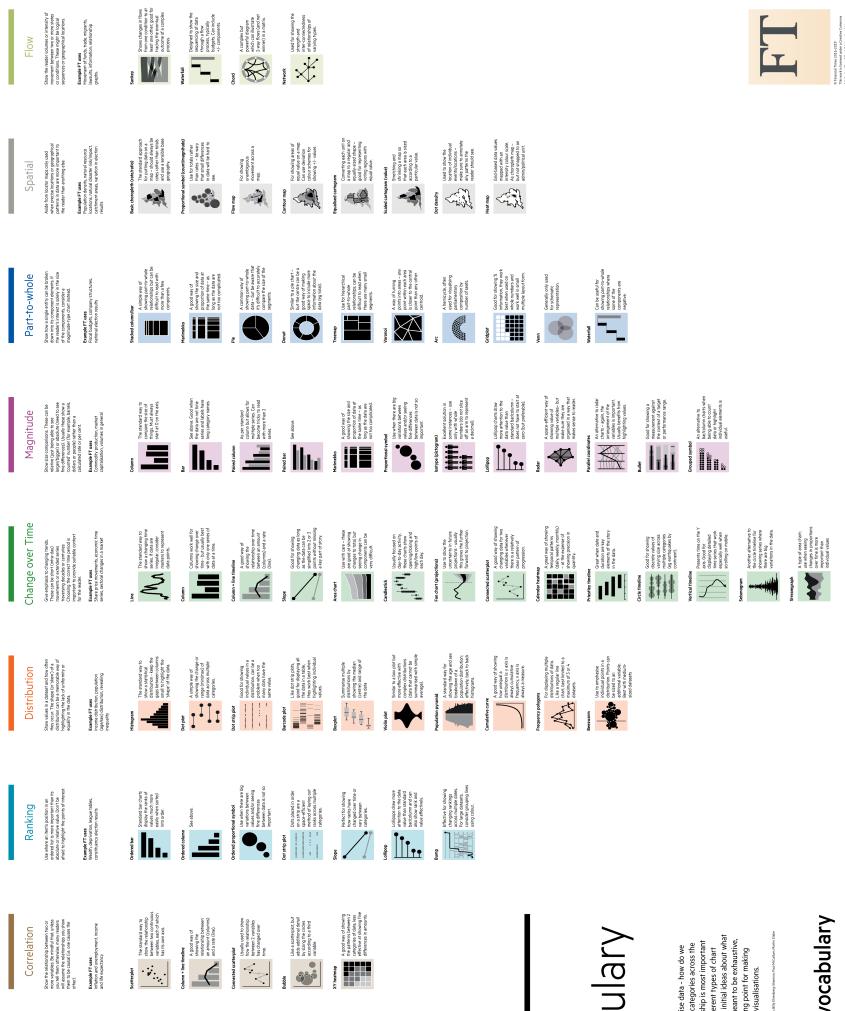


GGPLOT2 AS A CREATIVITY ENGINE

... AND OTHER WAYS R IS TRANSFORMING THE FT'S QUANTITATIVE JOURNALISM



FT



Visual vocabulary

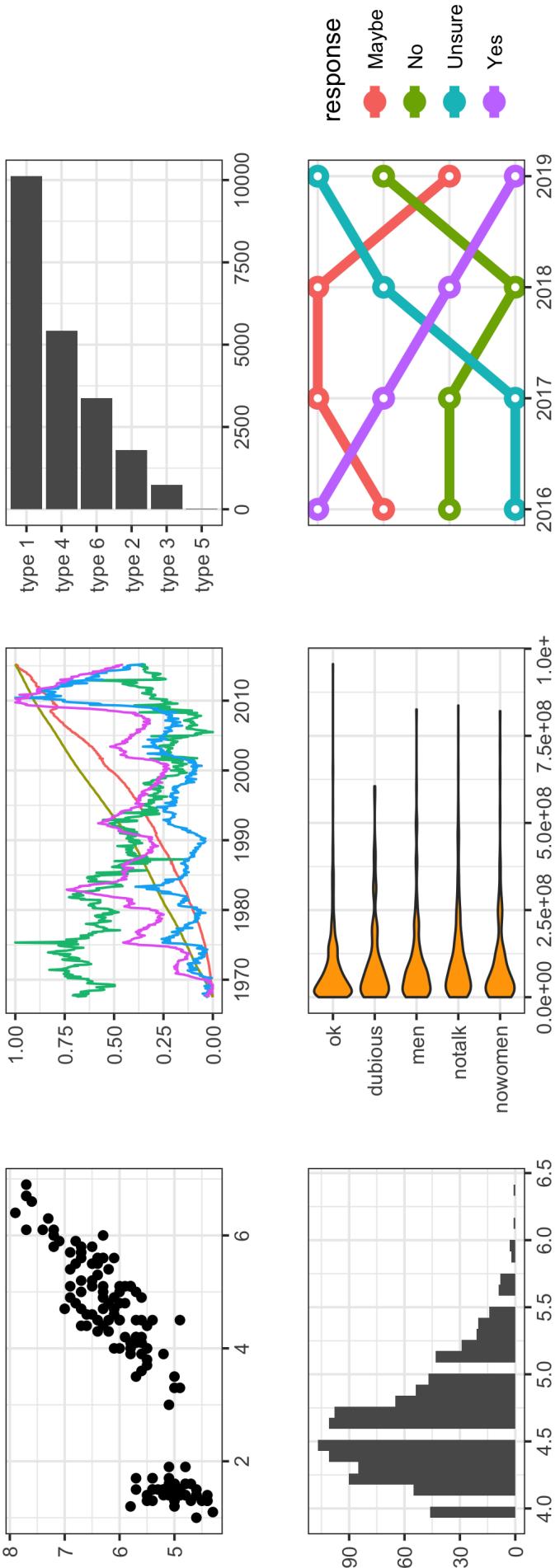
Designing with data

There are so many ways to visualize data – how do we know which one to pick? Use the categories across the top to help you decide. Each category has a chart example, a short blurb on the different types of charts within the category, and some general ideas about what might work best. This is not meant to be exhaustive, nor a guide. It is a useful starting point for reading, learning, and discussing data visualization.

ft.com/vocabulary









■ Aestheticss

▲ Geoms

◆ Scals

▷ Guidess

▼ Themes



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





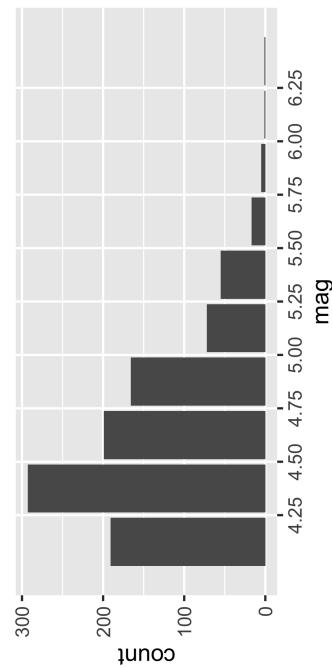
```
1 msl_eep %>%
2   ggplot() +
3     aes(
4       x = slEEP_total,
5       y = slEEP_Rem
6       colour = vore
7     ) +
8     geom_point()
```

geom_abline(a = 0, b = 1) # add identity line
geom_colege(m = contour, g = 0) # contour lines
geom_curv(g = dens, g = 0) # density contours
geom_dens(dens = density2d, g = 0) # density contours
geom_error(error = errorbar(t)) # error bars
geom_hline(g = geo(m), f = unc) # horizontal error bars
geom_point(point = point, g = 0) # points
geom_quantile(quantile = 0.5, r = 100) # quantile lines
geom_rug(rug = rug, s = 0.5) # rug lines
geom_smooth(smooth = st, g = 0) # smooth lines
geom_vline(v = 0) # vertical lines

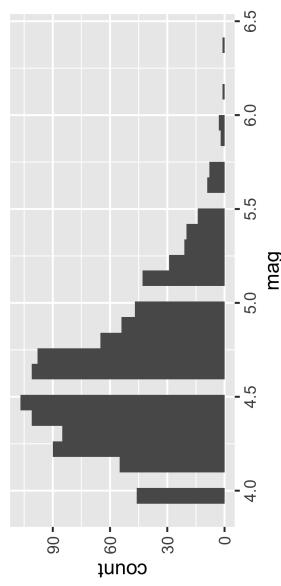




```
1 ggplot(quakes, aes(mag)) +  
2 geom_bar() +  
3 scale_x_binned()
```

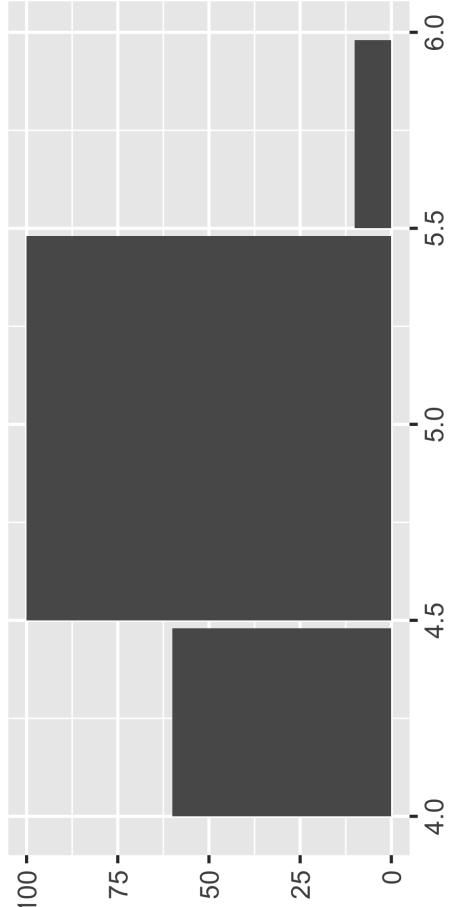


```
1 ggplot(quakes, aes(mag)) +  
2 geom_histogram()
```





```
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)
```

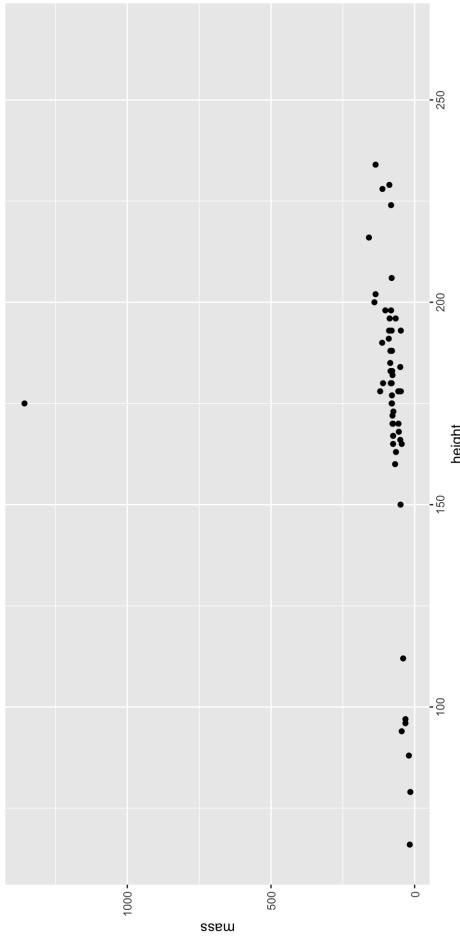




X

Y

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





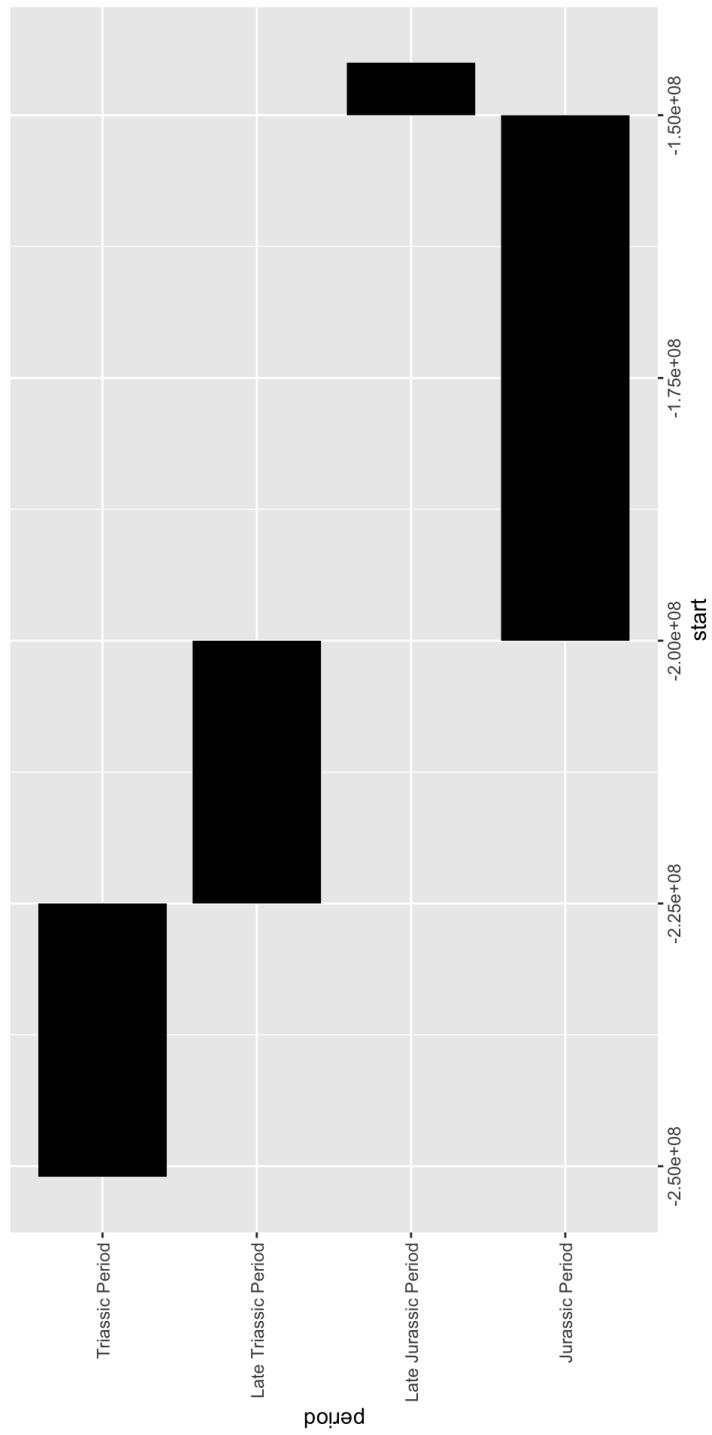
y

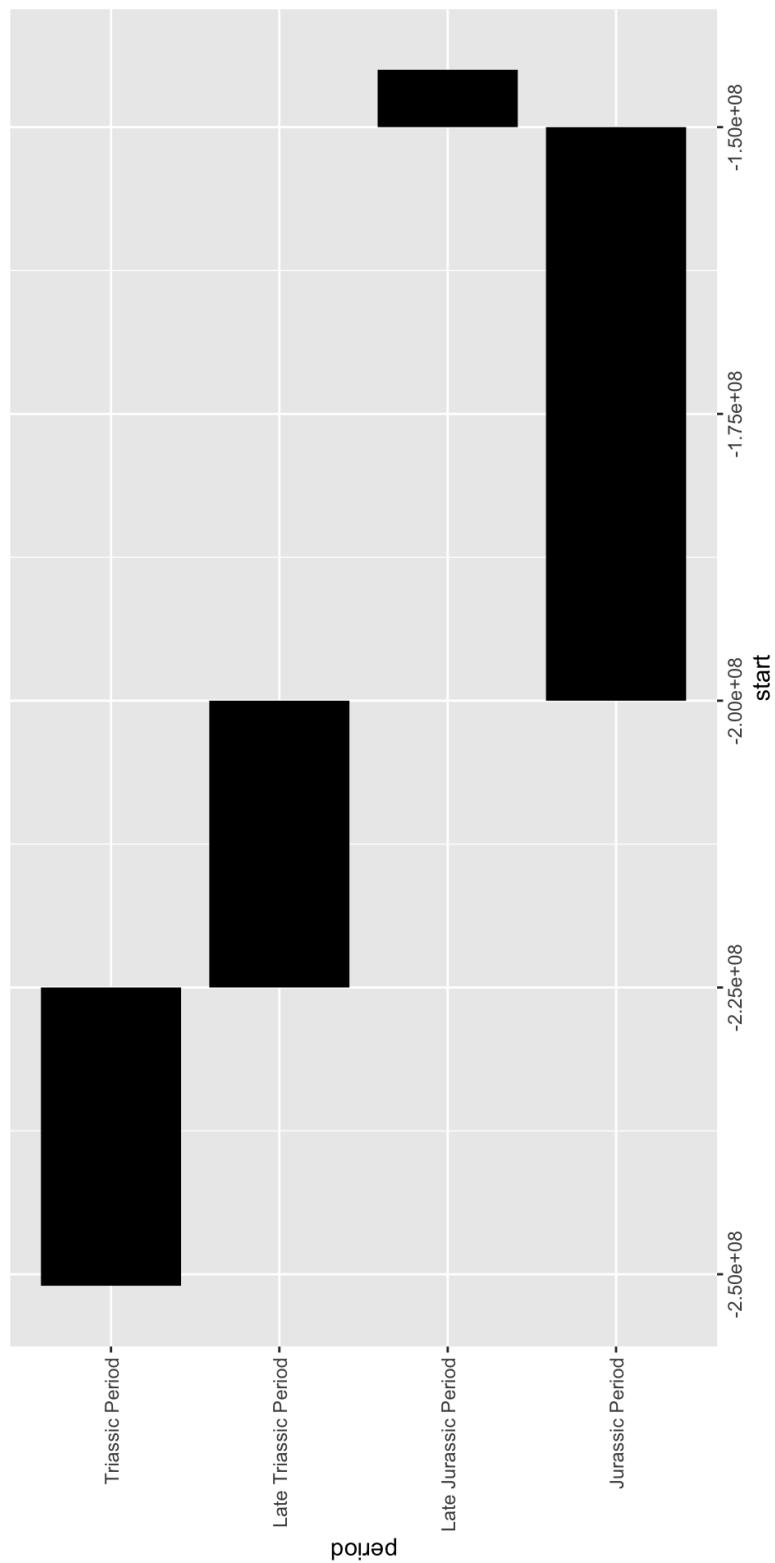
x

```
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 )
```

Size

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

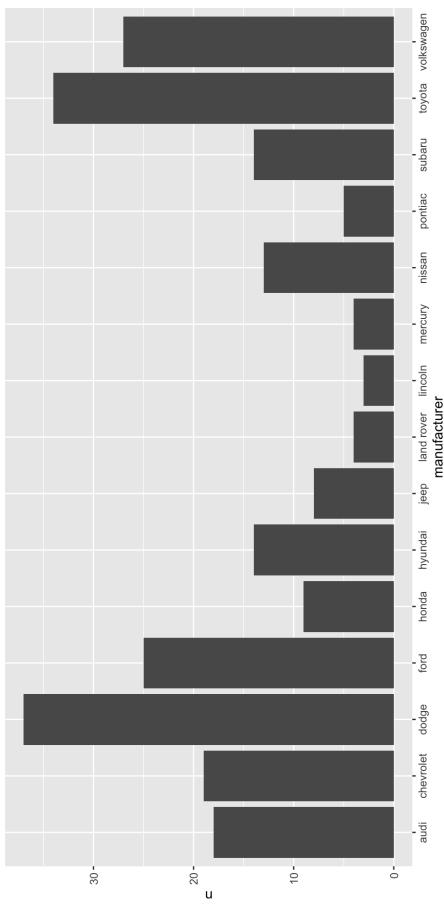




```

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

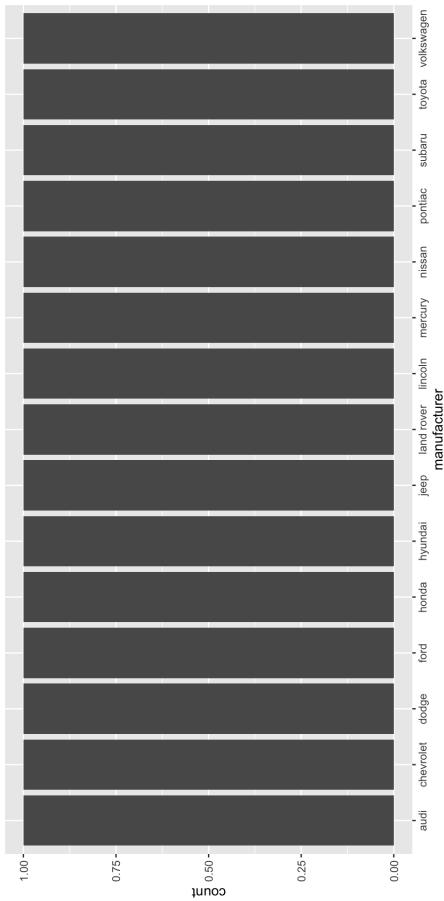
```



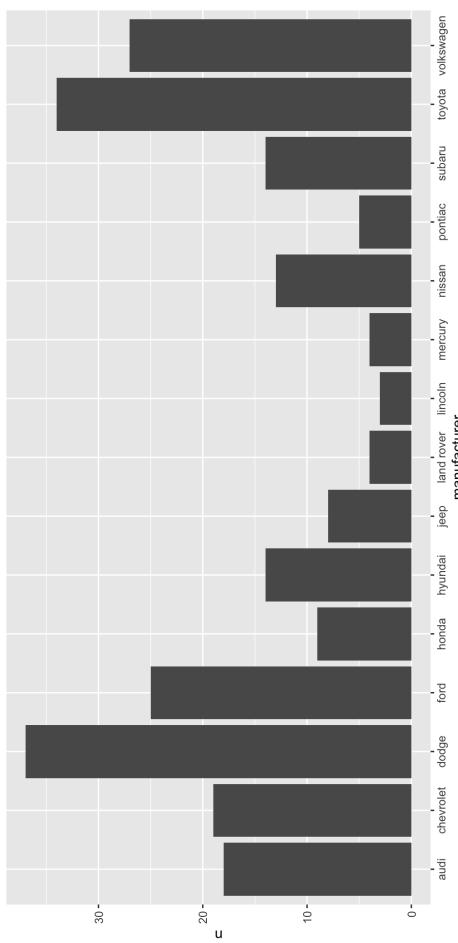
```

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

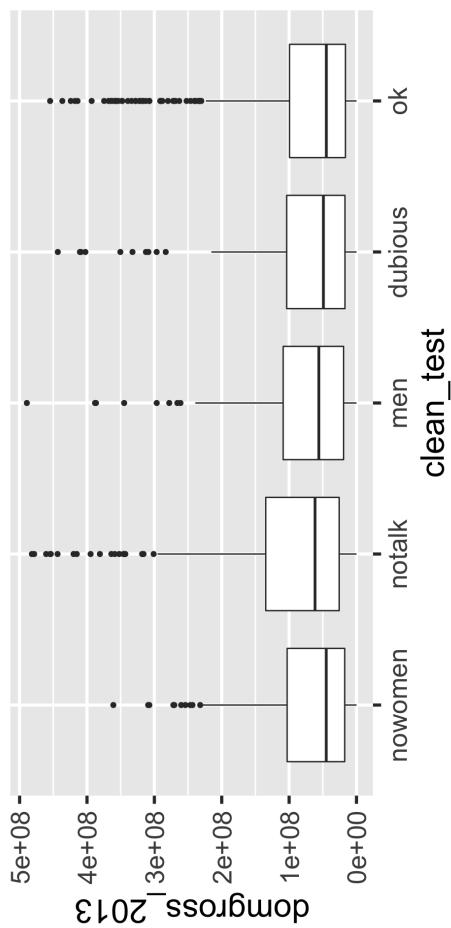
```



Stat



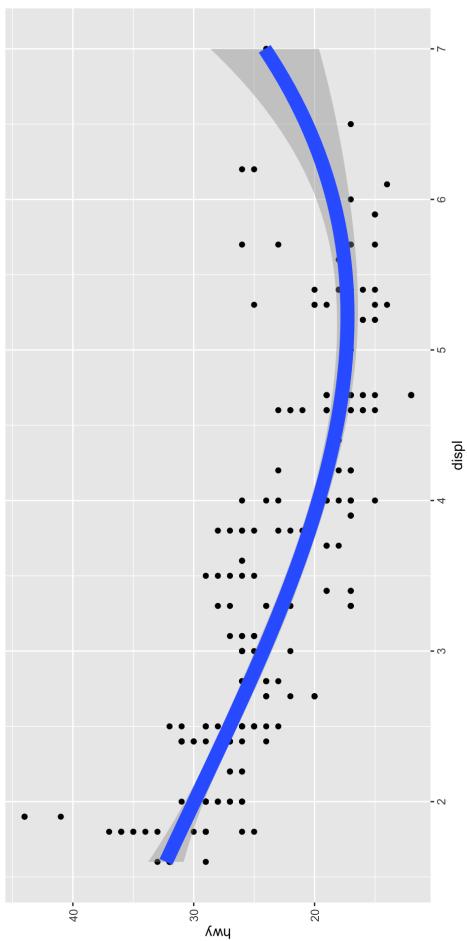
```
1 mpg %>%  
2 count(manufacturer) %>%  
3 ggplot() +  
4 geom_bar(aes(x = manufacturer,  
y = n),  
5 stat = "identity")  
6
```



```
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  Never married  4
#> 2 No answer  Separated    22
#> 3 No answer  Separated    3
#> 4 No answer  Divorced    13
#> 5 No answer  Widowed     7
#> 6 No answer  Married     44
#> 7 Don't know Never married  6
#> 8 Don't know Separated    3
#> 9 Don't know Divorced    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,
6   size = 5)
```

```
1 msl_eep %>%
2   ggplot() +
3     aes(
4       x = sleep_total,
5       y = sleep_rem
6       colour = vore
7     ) +
8     geom_point() +
9     scale_colour_manual(
10       values = c("car ni" = "#c03728",
11         "omni" = "#ff d8f 24",
12         "insecti" = "#f 5c04a",
13         "her bi" = "#919c4c",
14         "NA" = "#e68c7c")
15     )
```

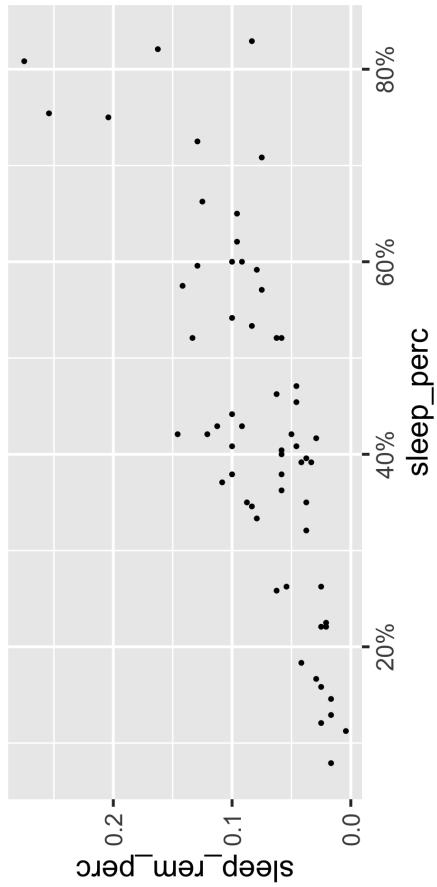






```
1 percent(c(0.3, 0.5, 0.6))  
[1] "30%" "50%" "60%"
```

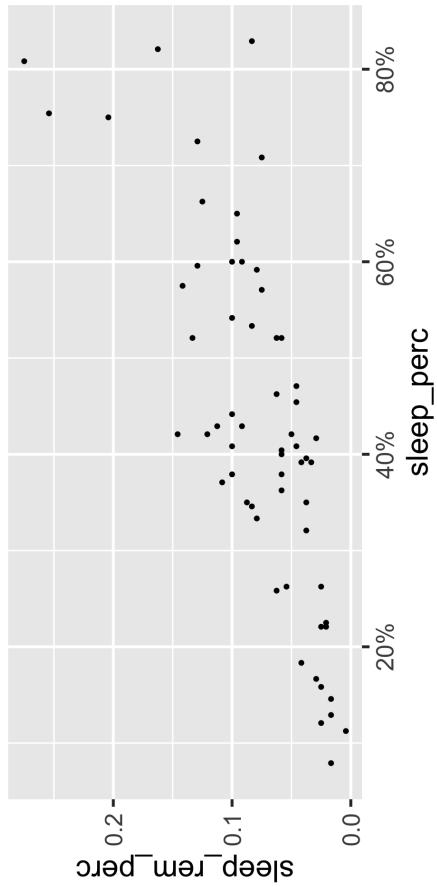
```
1 mSleep %>%  
2   mutate(sleep_perc = sleep_total / 24,  
3         sleep_remperc = sleep_rem / 24) %>%  
4   ggplot() +  
5   aes(x = sleep_perc,  
6        y = sleep_remperc) +  
7   geom_point() +  
8   scale_x_continuous(label = percent_format()) +  
9   theme_gray(base_size = 24)
```



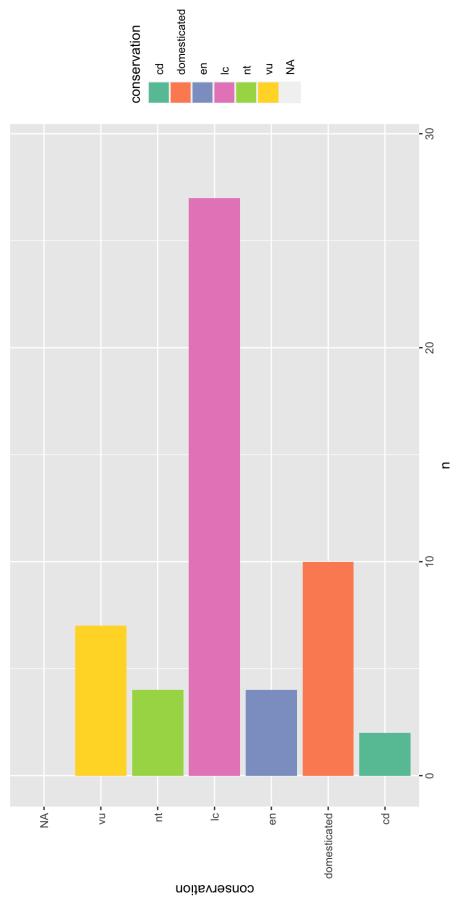


```
1 label_percent() (c(0.3, 0.5, 0.6))  
[1] "30%" "50%" "60%"
```

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

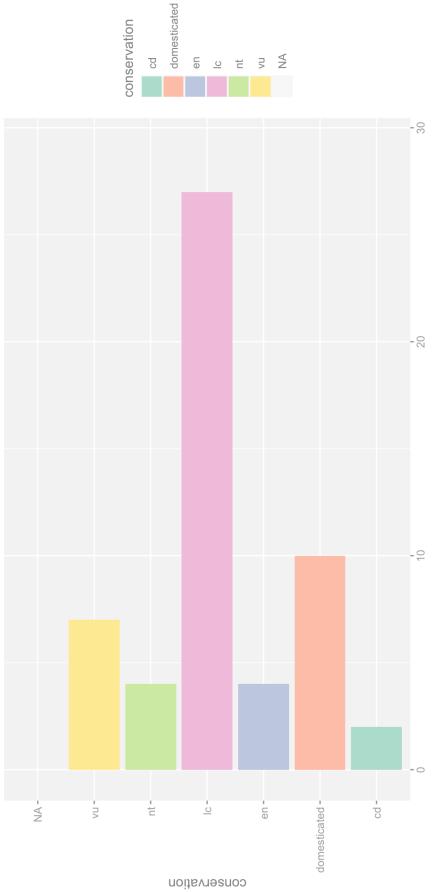
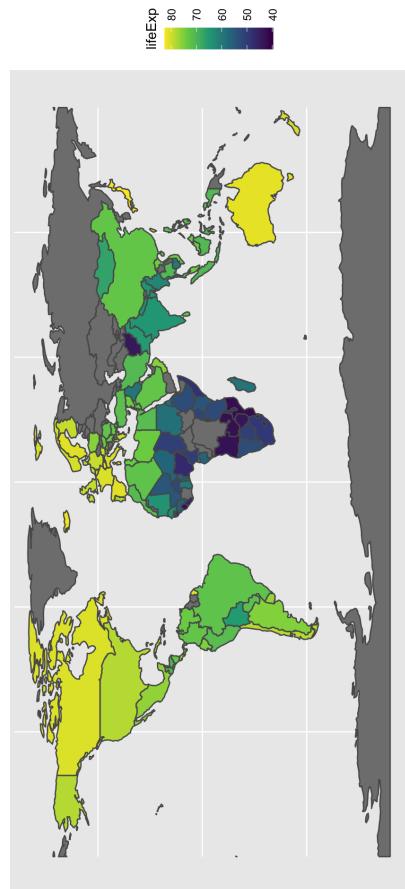






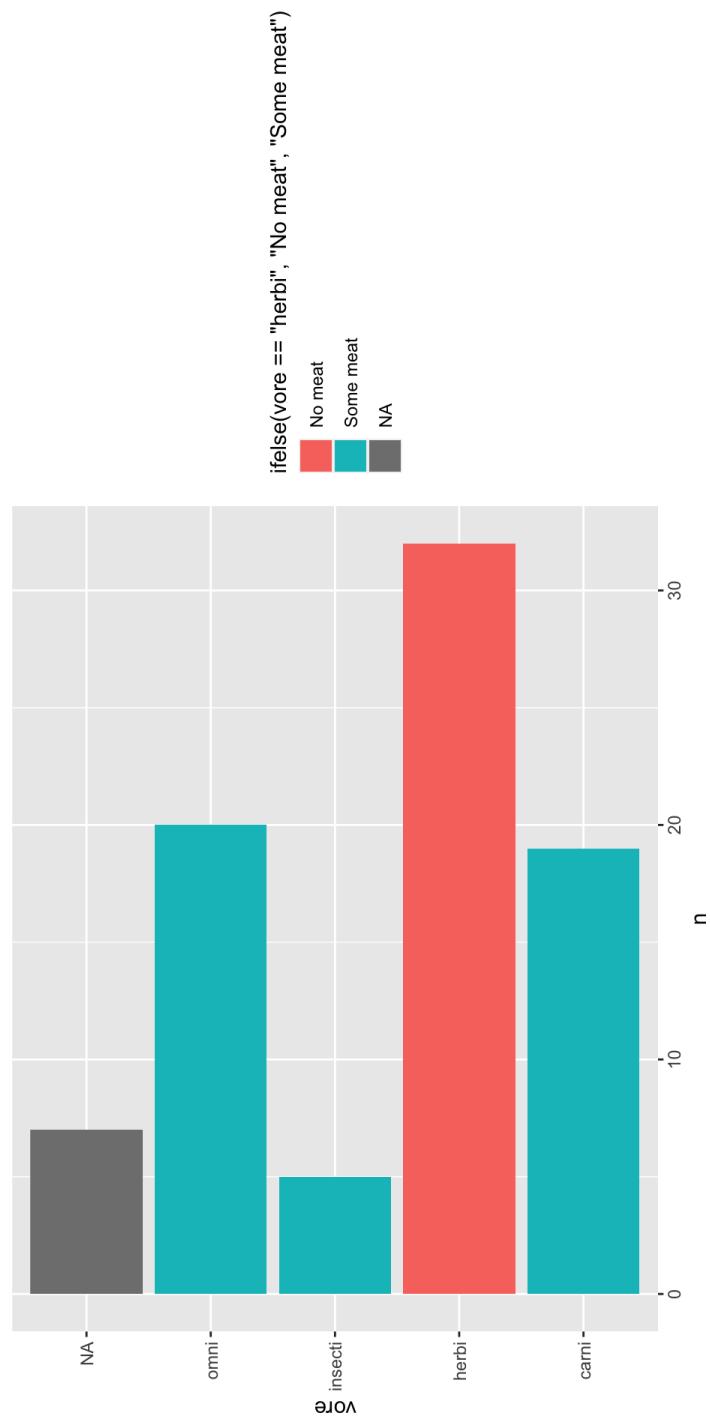


```
1 countries110 %>%
2   st_as_sf() %>%
3   left_join(filter(gapminder, year == 2007),
4             by = c("name" = "country")) %>%
5   ggplot() +
6   geom_sf(aes(fill = lifeExp)) +
7   scale_fill_viridis()
```



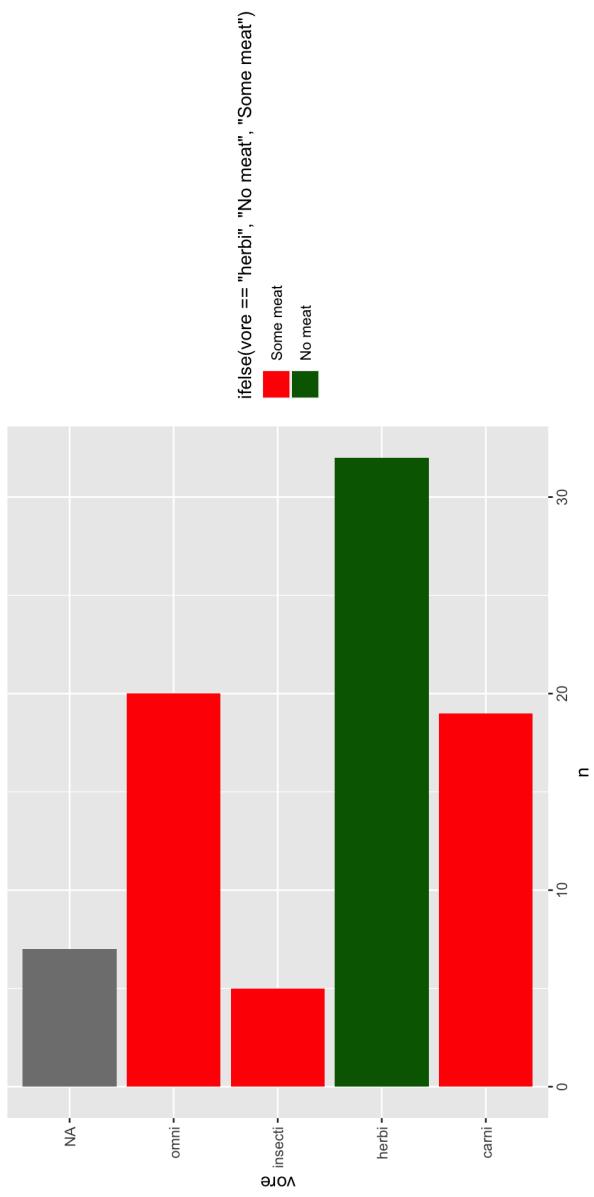


```
1 nseleep %>%  
2   count(vore) %>%  
3     ggplot() +  
4       aes(x = n,  
5               y = vore,  
6               fill = ifelse(vore == "herbi", "No meat", "Some meat")) +  
7     geom_col()
```





```
1 n$leep %>%
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)

```

```

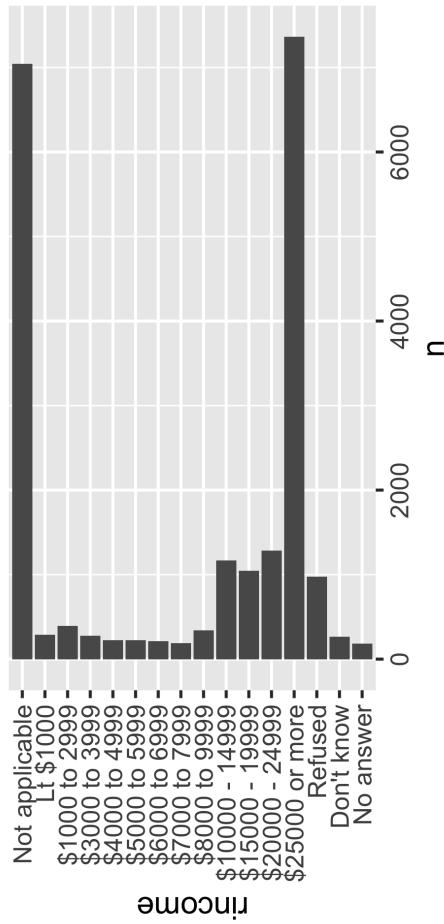
[1] $8000 to 9999 $8000 to 9999 Not appl i cable Not appl i cable
[6] $20000 - 24999
16 Levels: No answer Don't know Refused $25000 or more ... Not appl i cable

```

```

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

```



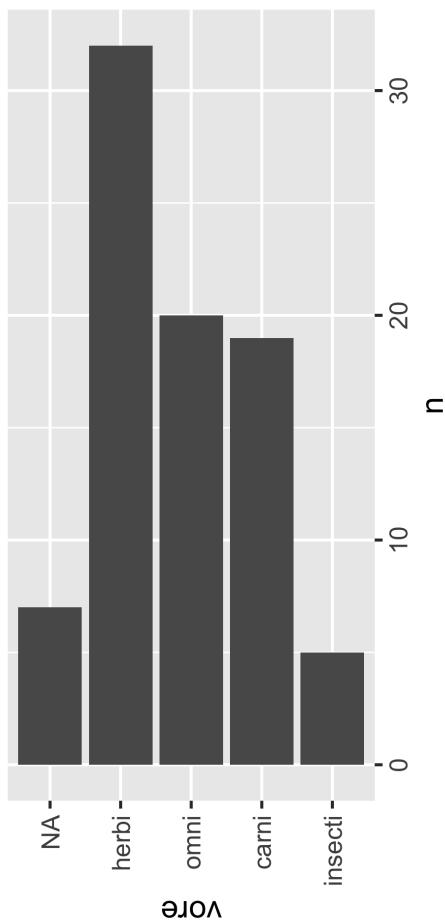




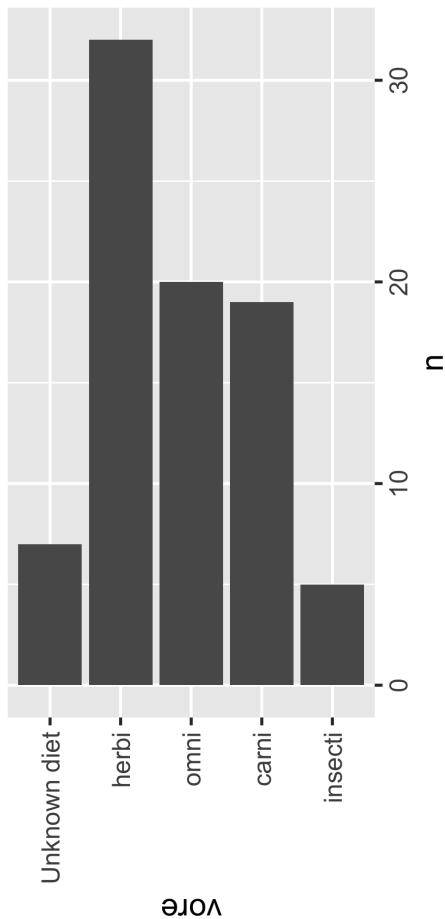
```
1 msleep %>%
2   count(vore)
```

A tibble: 5 × 2

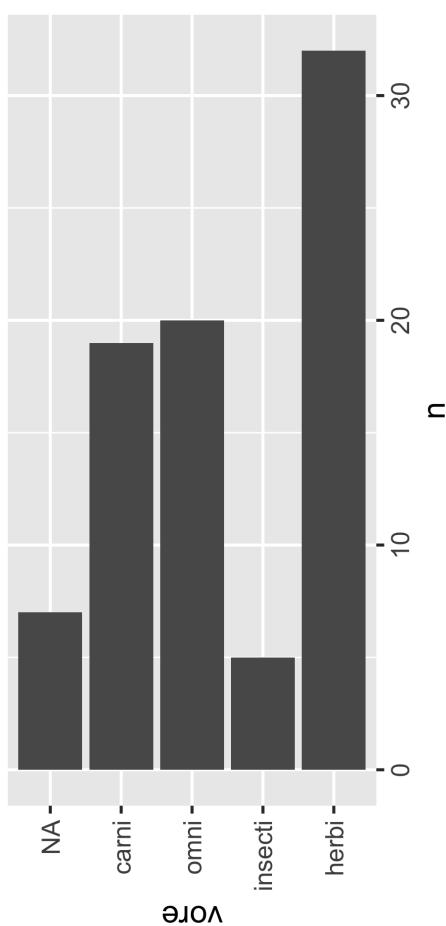
vore	n
<chr>	<int>
1 carni	19
2 herbi	32
3 insecti	5
4 omni	20
5 <NA>	7



```
1 msl_eep %>%
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)
```



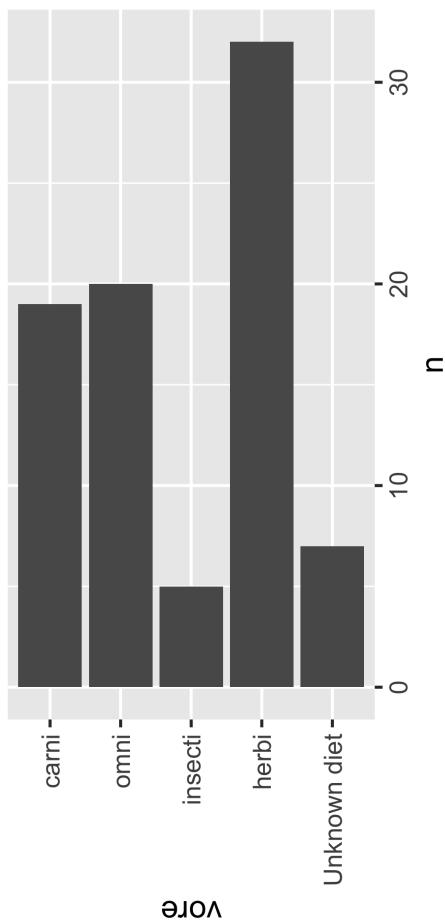
```
1 msl_eep %>%  
2   count(vore) %>%  
3     mutate(vore = fct_reorder(vore, n),  
4           vore = fct_expl(i_ci_t_na(vore, "Unknown diet")))  
5 ggplot() +  
6   aes(x = n,  
7       y = vore) +  
8   geom_col() +  
9   theme_gray(base_size = 24)
```



```

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

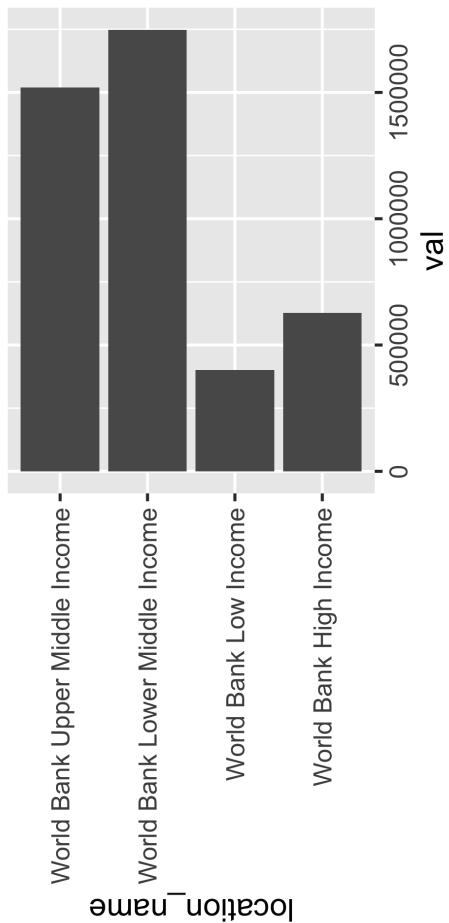


```
1 download_file("https://raw.githubusercontent.com/charliedey/eng7218_data-science-for-healthcare-applications_bcu-  
2 destfile = "data/global-burden-of-disease-data.csv")
```





• • • •

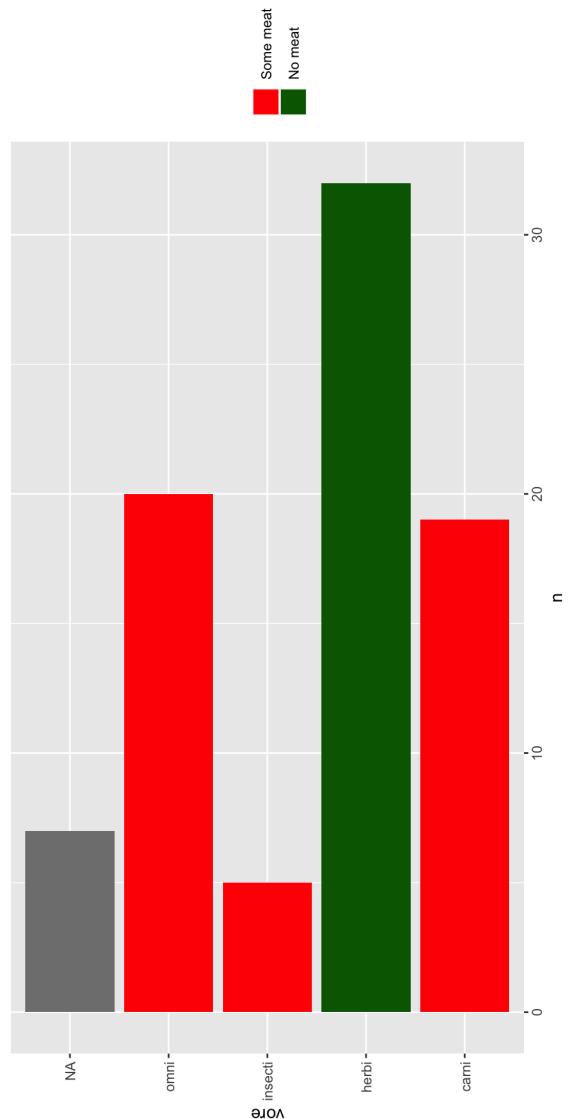


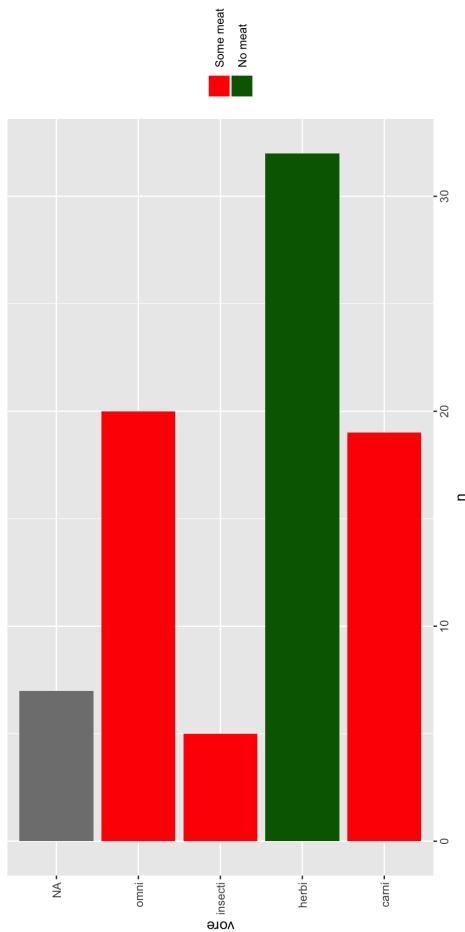
```
1 gdb_injuries %>%
2   ggplot() +
3     aes(x = val,
4       y = location_name) +
5     geom_col() +
6     theme_gray(base_size = 24)
```





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





```

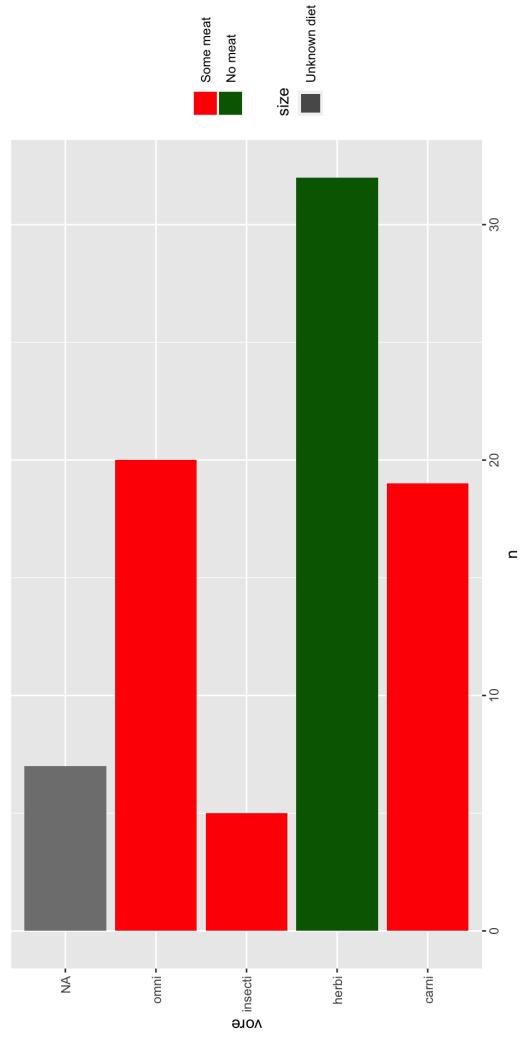
1 msl_eep %>%
2   count(vore) %>%
3     ggplot() +
4       aes(x = n,
5         y = vore,
6         fill = ifelse(vore == "herbivore",
7                         "No meat",
8                         "Some meat")) +
9       geom_col() +
10      scale_fill_manual(values = c("Some meat" = "red",
11     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

```

```

1 msl_eep %>%
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 = "")

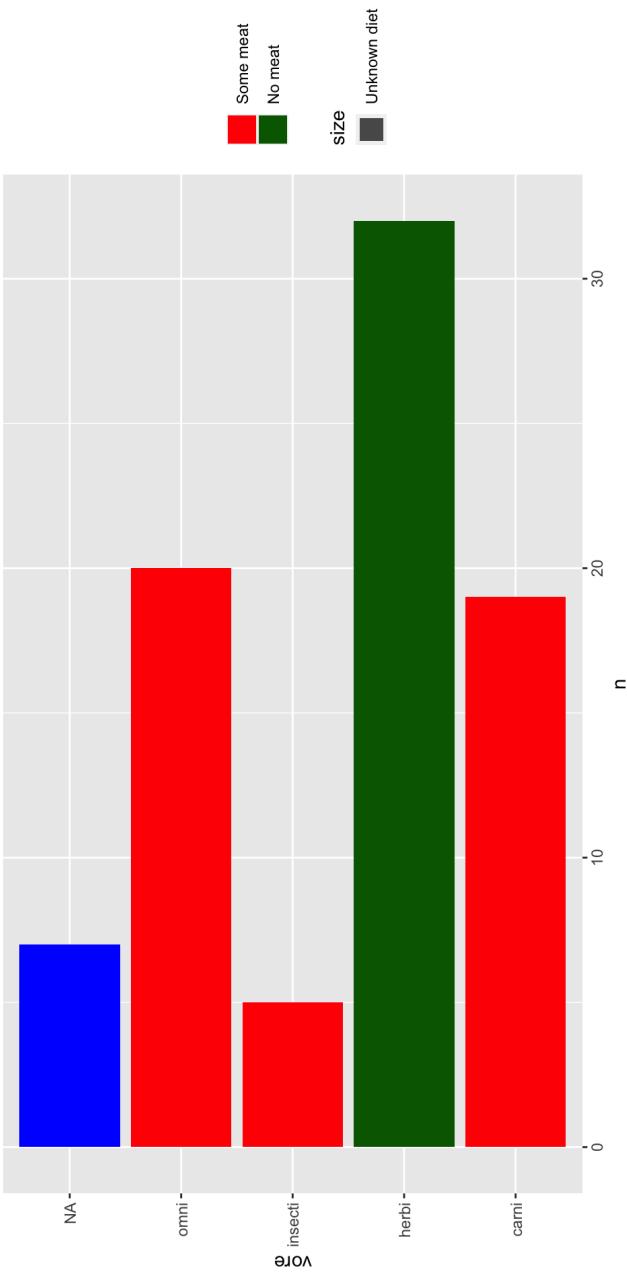
```



```

1 nseep %>%
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")

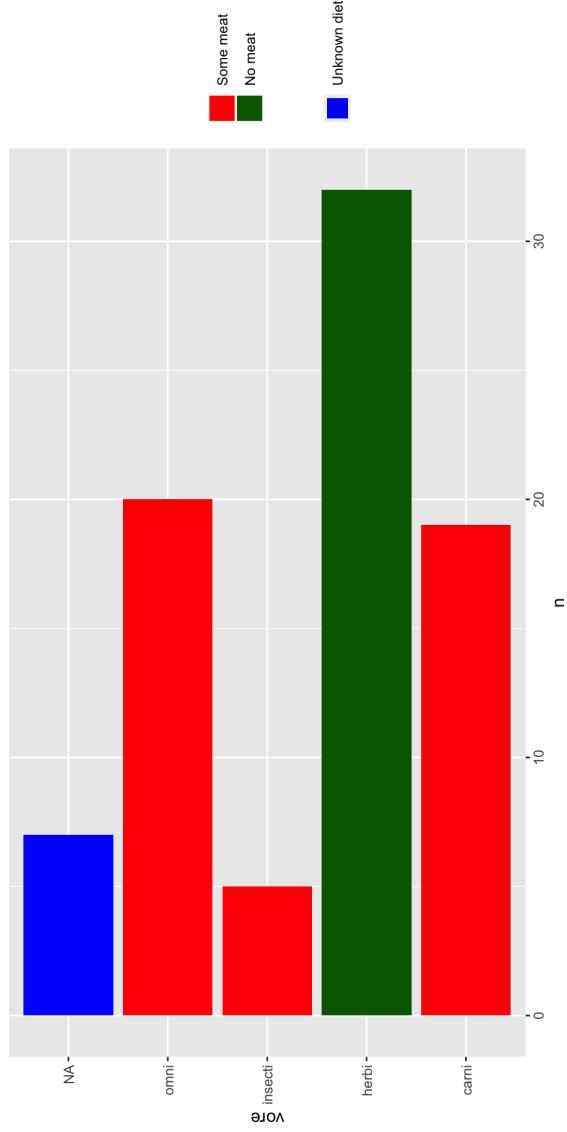
```



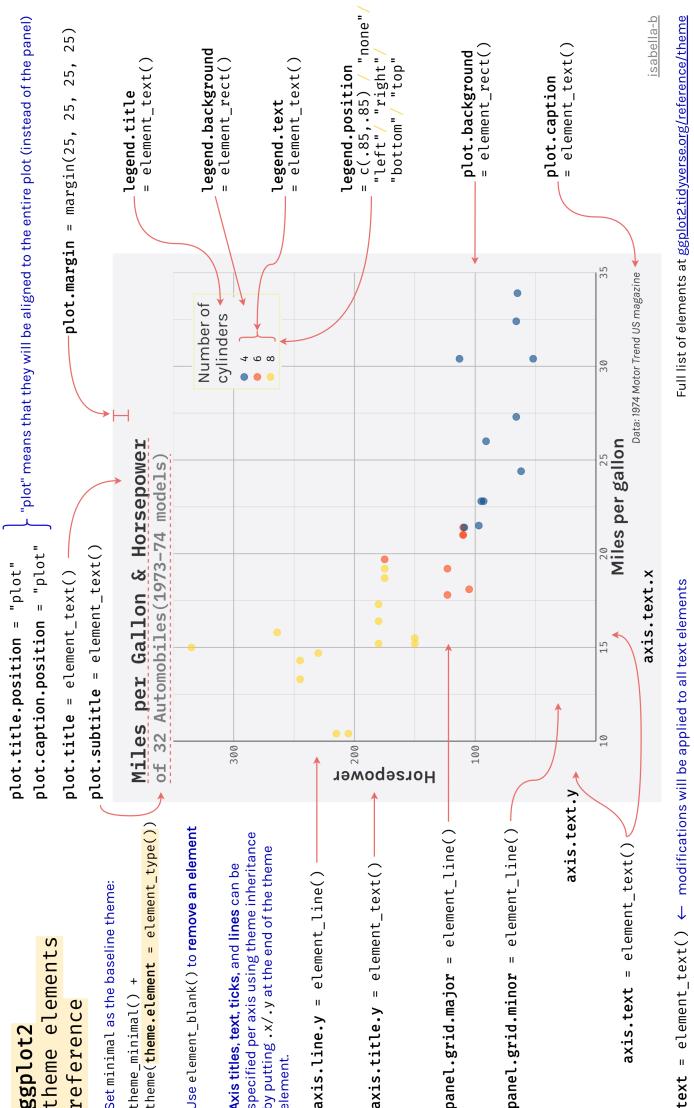
```

1 nseep %>%
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")))

```







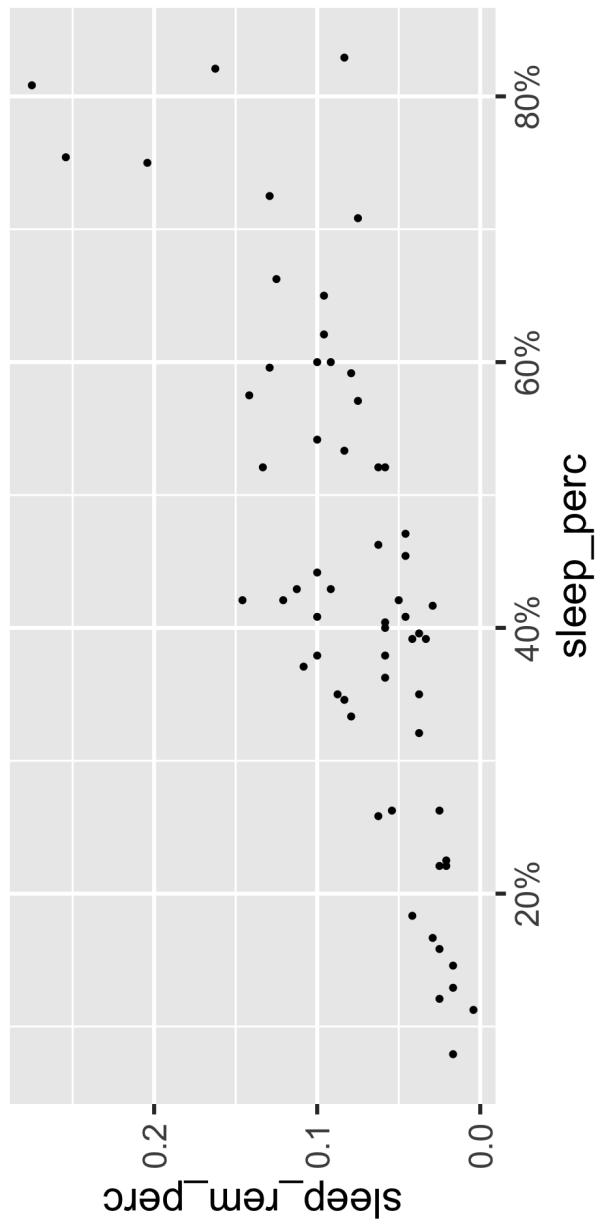
isabella:b

Full list of elements at ggplot2.tidyverse.org/reference/theme.html

```

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

```





```
1 theme_fi vethir tyei ght() +
2 theme(panel . grid. major = element_l ine(colour = "red"))
```


References

Revi & Quarter y

Management

The Ameri an Stat i ts ic@ 27

Proceedi gs of the 2017 CHI Conference on Human Factors i Computing
Systems

On the mode of communi ati nofchol r@

Notes on Matters Affecti gthe Heal ht E ffici ency and Hospi at Admini strati on of the Bri tish Army

The best stats you' ever seen

EuroVi 2016 - Short Papers

Journal of the Ameri an Stat i stic Associ ati on 1979

Proceedi gs of the 28th Internati onal conference on Human factors i ncomputi ngsystems - CHI ' 0

ISPRS Internati onal Journal of Geo-Information 1010

Appl i cati ons 42

IEEE Computer Graphi cs and



Twi ter

