

plotting numeric data using scales, labels, jitter, lines of best fit

Emily Malcolm-White

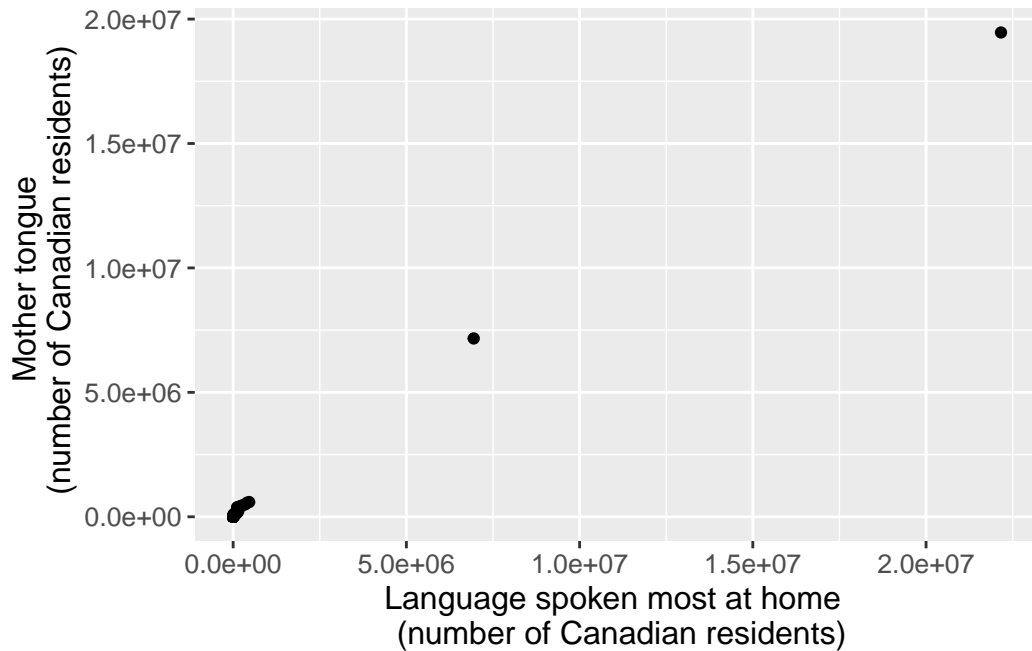
```
library(tidyverse)
```

```
#Import the can_lang dataset
```

```
can_lang <- read.csv("https://raw.githubusercontent.com/ttimbers/canlang/master/inst/extdata/can_lang.csv")
```

scatterplot of can_lang

```
ggplot(can_lang, aes(x=most_at_home, y=mother_tongue)) +  
  geom_point() +  
  xlab("Language spoken most at home \n (number of Canadian residents)") +  
  ylab("Mother tongue \n (number of Canadian residents)") +  
  theme(text = element_text(size = 12))
```

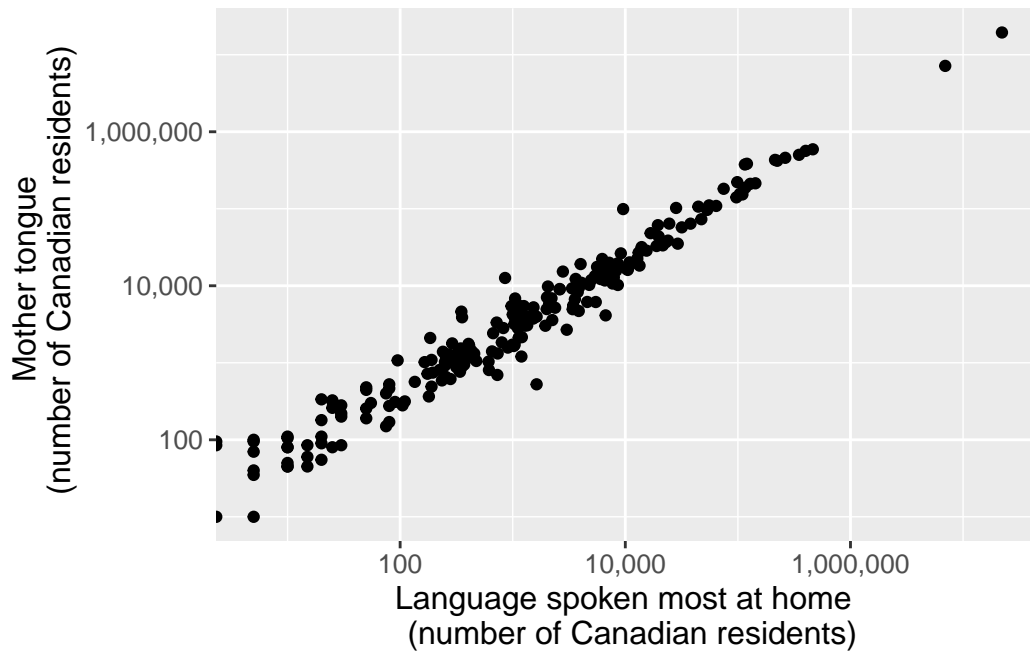


Notice anything weird about this plot?

Axis transformations using the scales package

```
library(scales)

ggplot(can_lang, aes(x = most_at_home, y = mother_tongue)) +
  geom_point() +
  xlab("Language spoken most at home \n (number of Canadian residents)") +
  ylab("Mother tongue \n (number of Canadian residents)") +
  theme(text = element_text(size = 12)) +
  scale_x_log10(labels = label_comma()) +
  scale_y_log10(labels = label_comma())
```



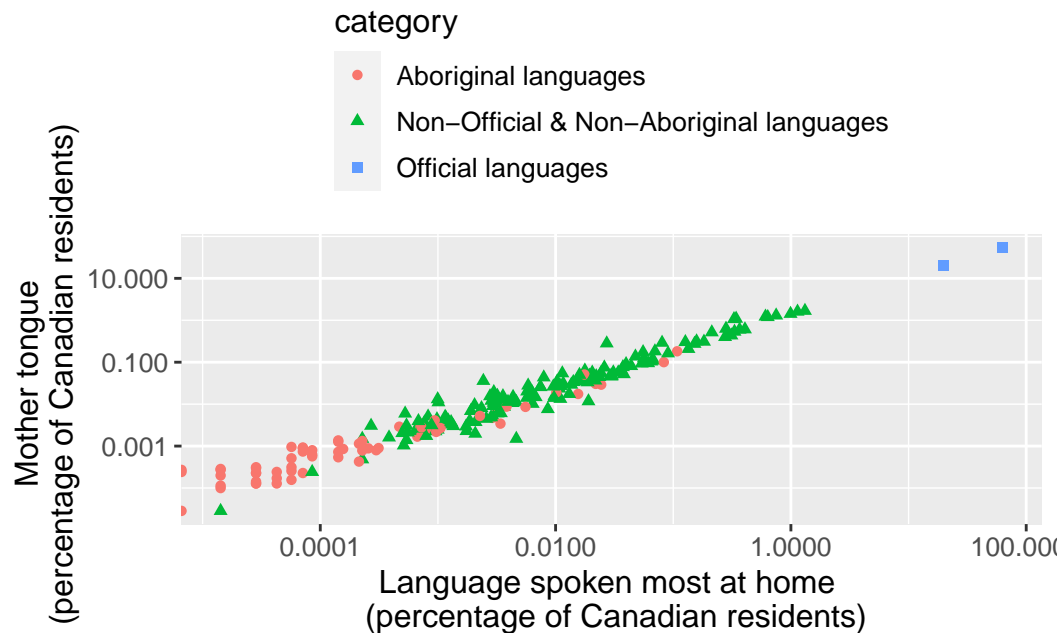
mutate to create new columns

```
can_lang <- can_lang %>%
  mutate(
    mother_tongue_percent = (mother_tongue / 35151728) * 100,
    most_at_home_percent = (most_at_home / 35151728) * 100
  )
```

Scatterplot with Percents and Colors

```
ggplot(can_lang, aes(x = most_at_home_percent,
                     y = mother_tongue_percent,
                     color = category, shape=category)) +
  geom_point() +
  xlab("Language spoken most at home \n (percentage of Canadian residents)") +
  ylab("Mother tongue \n (percentage of Canadian residents)") +
  theme(text = element_text(size = 12),
        legend.position = "top",
        legend.direction = "vertical") +
  scale_x_log10(labels = comma) +
```

```
scale_y_log10(labels = comma)
```



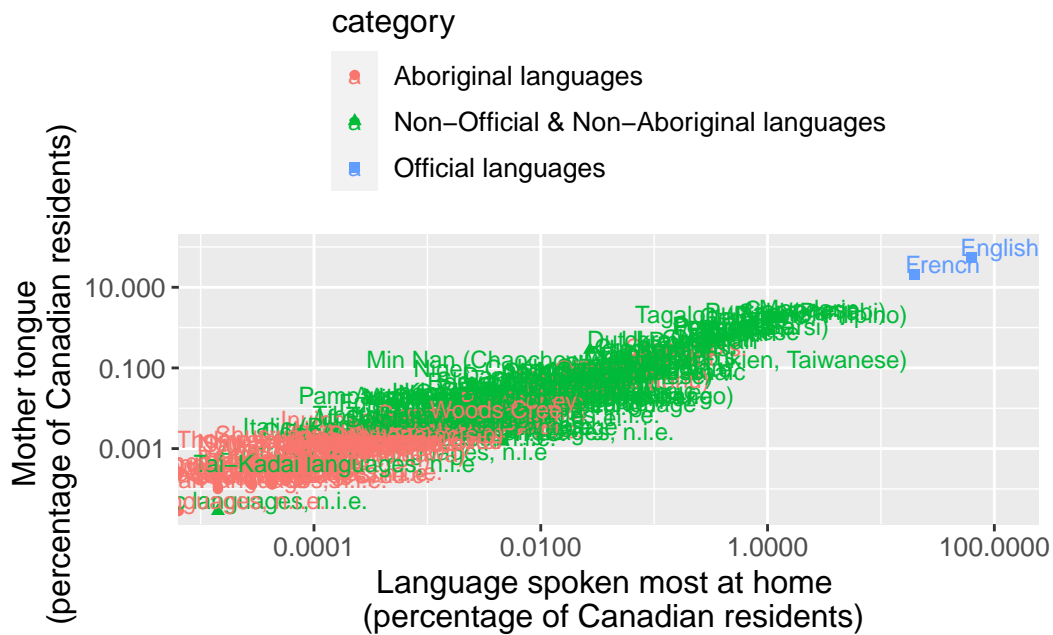
Labels

Adding text to a plot is one of the most common forms of annotation. Most plots will not benefit from adding text to every single observation on the plot, but labeling outliers and other important points is very useful.

```
ggplot(can_lang, aes(x = most_at_home_percent,
                     y = mother_tongue_percent,
                     color = category, shape=category)) +
  geom_point() +
  xlab("Language spoken most at home \n (percentage of Canadian residents)") +
  ylab("Mother tongue \n (percentage of Canadian residents)") +
  theme(text = element_text(size = 12),
        legend.position = "top",
        legend.direction = "vertical") +
  scale_x_log10(labels = comma) +
  scale_y_log10(labels = comma) +
  geom_text(aes(label=language),
            nudge_x = 0.25,
```



```
nudge_x = 0.25,
nudge_y=0.25, size=3)
```



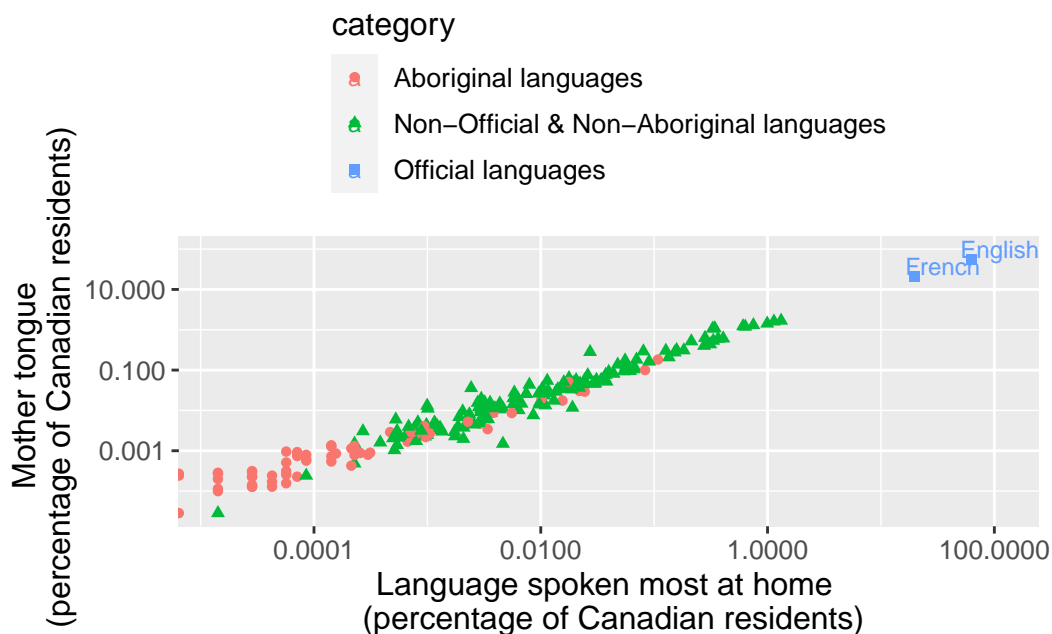
Subset the labels

Create a new column for the labels. Use `case_when` (or `ifelse`) to only use the official language names and not to put a label for other language categories.

```
can_lang <- can_lang %>%
  mutate(official_languages = case_when(category == "Official languages" ~ language, TRUE

ggplot(can_lang, aes(x = most_at_home_percent,
  y = mother_tongue_percent,
  color = category, shape=category)) +
  geom_point() +
  xlab("Language spoken most at home \n (percentage of Canadian residents)") +
  ylab("Mother tongue \n (percentage of Canadian residents)") +
  theme(text = element_text(size = 12),
    legend.position = "top",
    legend.direction = "vertical") +
  scale_x_log10(labels = comma) +
```

```
scale_y_log10(labels = comma) +
geom_text(aes(label=official_languages),
          nudge_x = 0.25,
          nudge_y=0.25, size = 3)
```



Using ggrepel

```
library(ggrepel)

ggplot(can_lang, aes(x = most_at_home_percent,
                    y = mother_tongue_percent,
                    color = category, shape=category)) +
  geom_point() +
  xlab("Language spoken most at home \n (percentage of Canadian residents)") +
  ylab("Mother tongue \n (percentage of Canadian residents)") +
  theme(text = element_text(size = 12),
        legend.position = "top",
        legend.direction = "vertical") +
  scale_x_log10(labels = comma) +
  scale_y_log10(labels = comma) +
  geom_text_repel(aes(label=official_languages), min.segment.length=0, box.padding=1)
```

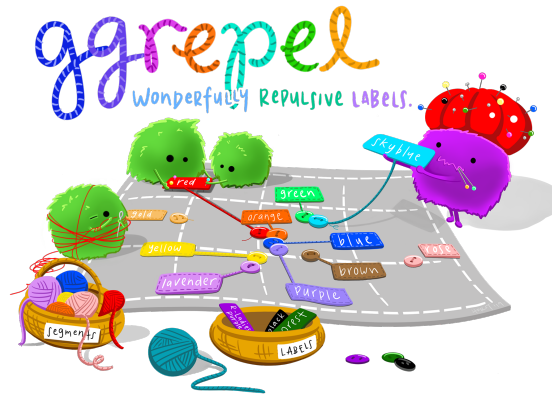
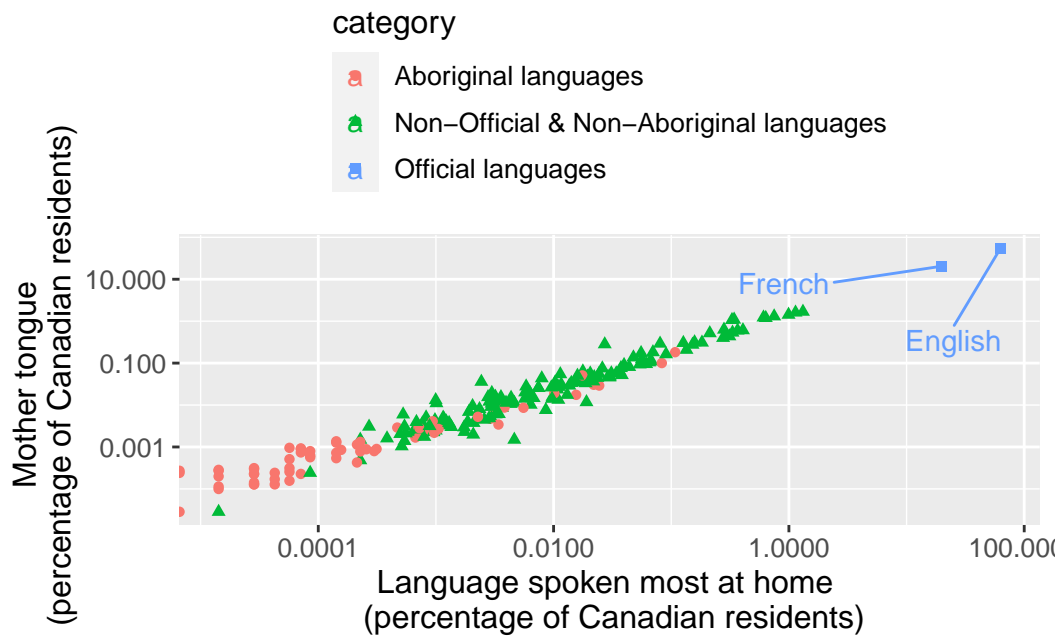


Figure 1: Artwork by @allisonhorst



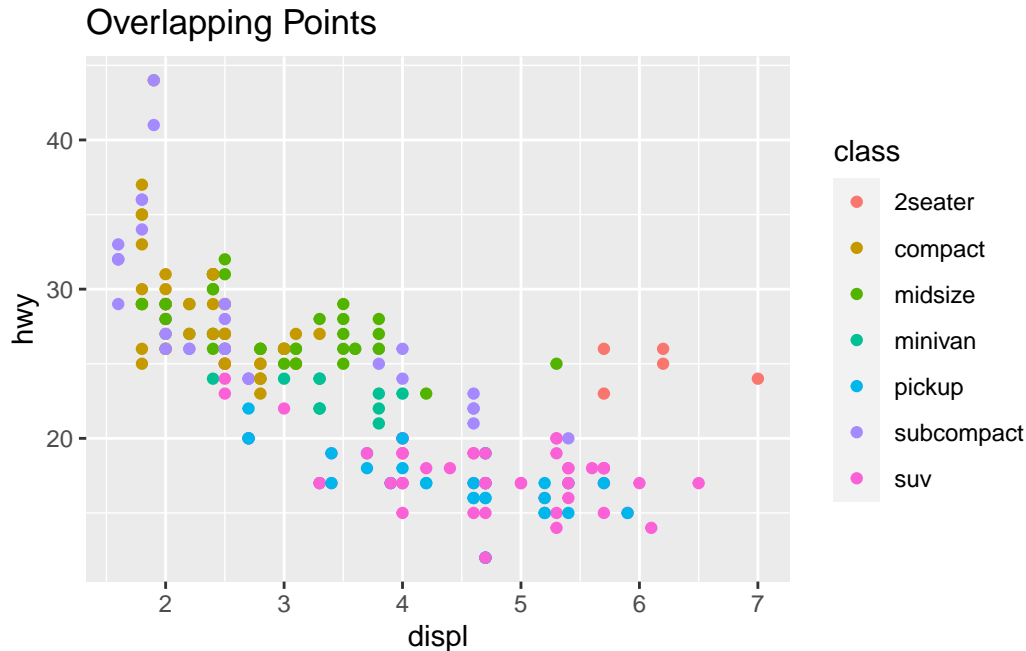
(Optional) Jitter

This dataset contains a subset of the fuel economy data that the EPA makes available on <https://fuelconomy.gov/>. It contains only models which had a new release every year between 1999 and 2008 - this was used as a proxy for the popularity of the car.

```
data("mpg")
```



```
ggplot(data = mpg, aes(x = displ, y = hwy)) +
  geom_point(aes(color = class)) +
  ggtitle("Overlapping Points")
```

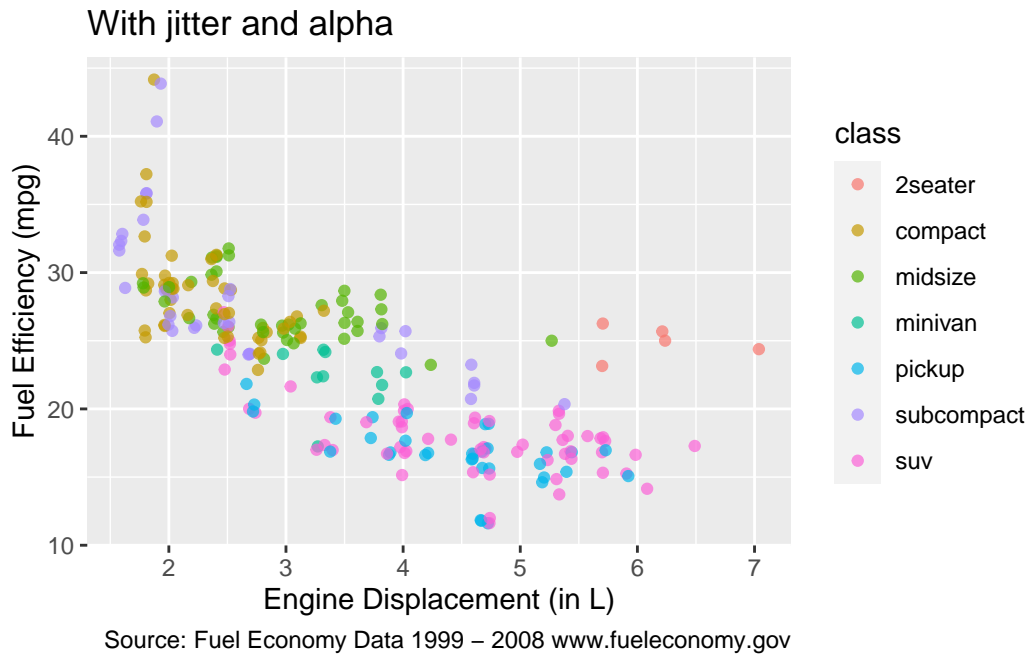


Note how there is only a fixed number of values for the x-axis and a fixed number of values for the y-axis. It's possible that some points might be overlapping others!

Jittering is a technique for adding random noise to data points that have identical values in a plot. It is a useful method for avoiding overplotting and making it easier to visualize the density of data points that are otherwise hidden behind each other.

We add `position="jitter"` inside the `geom_point` layer:

```
ggplot(data = mpg, aes(x = displ, y = hwy)) +
  geom_point(aes(color = class), position = "jitter", alpha=0.7) +
  ggtitle("With jitter and alpha") +
  xlab("Engine Displacement (in L)") +
  ylab("Fuel Efficiency (mpg)") +
  labs(caption = "Source: Fuel Economy Data 1999 - 2008 www.fueleconomy.gov")
```



```
library("DAAG")
```

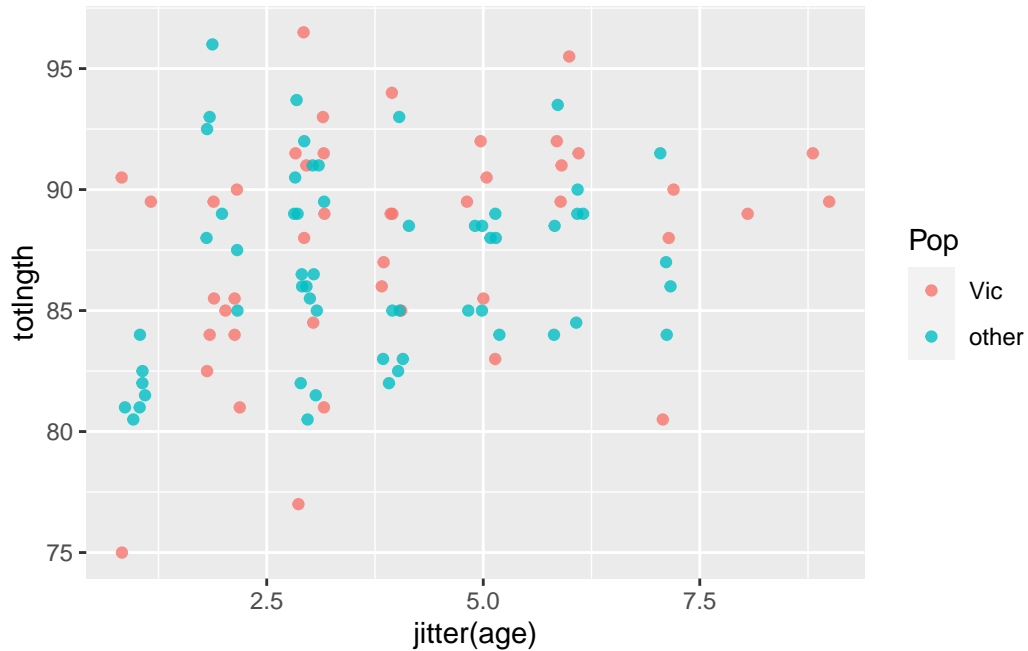
Jitter

Jittering is a technique for adding random noise to data points that have identical values in a plot. It is a useful method for avoiding overplotting and making it easier to visualize the density of data points that are otherwise hidden behind each other.

Often when we “jitter”, we might also adjust the opacity of the points to be able to better see ALL the points.

A plot of possum age vs. total length with age “jittered”:

```
ggplot(possum, aes(x=jitter(age), y=totlngth)) +  
  geom_point(aes(color=Pop), alpha=0.8)
```

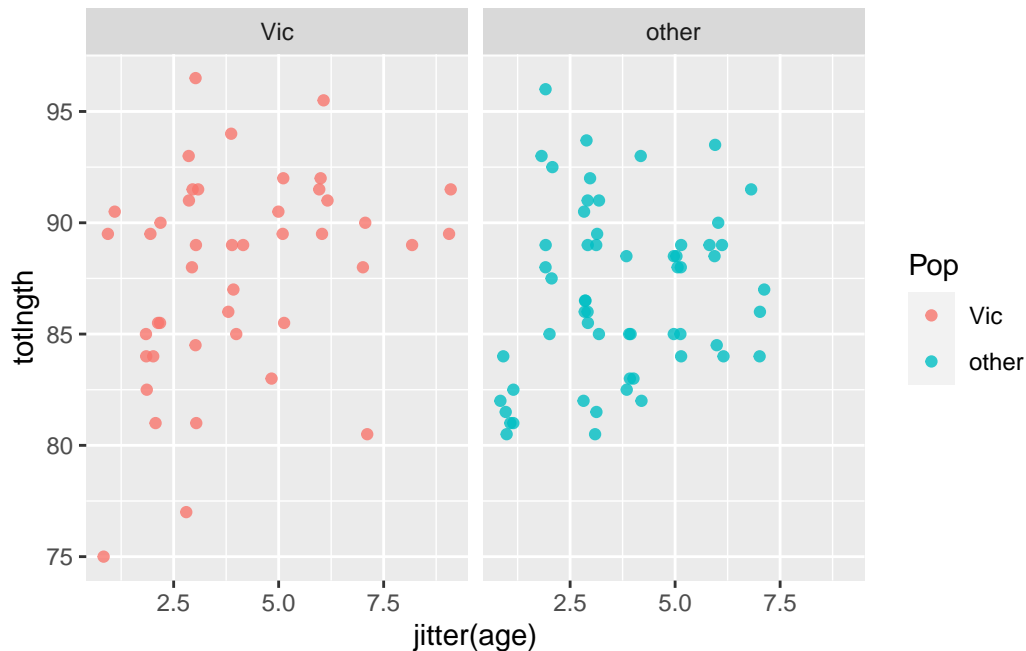


Facet Wrap

`facet_wrap()` is a function in the `ggplot2` package that allows you to create a multi-panel plot showing a similar plot over different subsets of the data, usually different values of a categorical variable.

A scatter plot of age vs. total length, creating separate side-by-side plots for each different population.

```
ggplot(possum, aes(x=jitter(age), y=totlength)) +  
  geom_point(aes(color=Pop), alpha=0.8) +  
  facet_wrap(~Pop)
```



Labels

Adding text to a plot is one of the most common forms of annotation. Most plots will not benefit from adding text to every single observation on the plot, but labeling outliers and other important points is very useful.

You can use `geom_text` to add labels in base R, but it can get quite messy when you have overlapping labels. I prefer using the `ggrepel` package:

```
# required packages
library(ggrepel)
```

A scatterplot of age vs. total length, with labels for each possum number:

```
ggplot(possum, aes(x=jitter(age), y=totlength)) +
  geom_point(aes(color=Pop), alpha=0.8) +
  #geom_text(aes(label=case), nudge_x = 0.5, nudge_y = 0.5, size=2) #<1>
  geom_text_repel(aes(label=case), size=2) #<2>
```

- ① This is the code we would have used if we were using the classic ggplot labeling
- ② Adds a layer with labels using the `ggrepel` package

