

# Lab 1

Student Name

YYYY-MM-DD

## Data

We'll work with the #tidytuesday data for 2019, specifically the #rstats dataset, containing nearly 500,000 tweets over a little more than a decade using that hashtag.

The data is in under Dataset tab of Week 3 module on Canvas.

You can import the dataset using the code below.

```
d <- rio::import(here::here("data", "rstats_tweets.rds"),
                 setclass = "tbl_df")
```

If you need help with processing text data, please revisit the notebook introduced in Week 1.

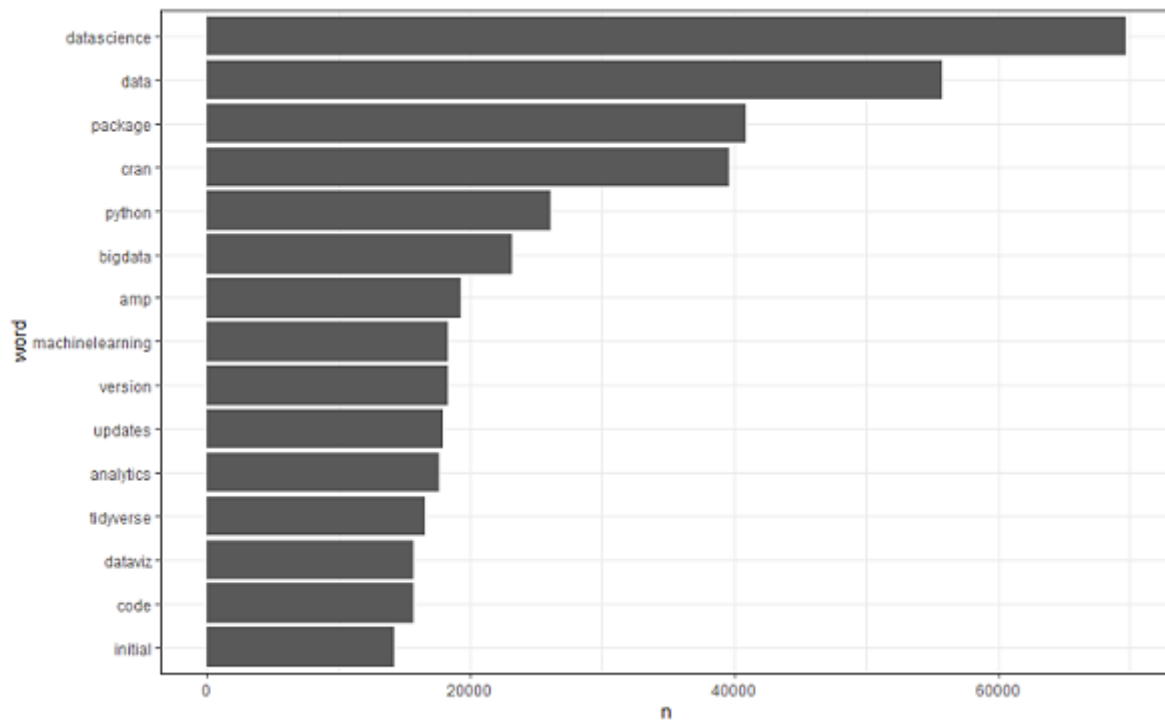
<https://www.kaggle.com/code/uocoeeds/introduction-to-textual-data>

## Histogram and Density plots

1. Create a histogram the column `display_text_width` using the `ggplot2` package and `geom_histogram()` function. Try at least four different numbers of bins (e.g., 20, 30, 40, 50) by manipulating the `bins=` argument. Select what you think best represents the data for each. Provide a brief justification for your decision. For all plots you created, change the default background color from grayish to white.
2. Create a density plot for the column `display_text_width` using the `ggplot2` package and `geom_density()` function. Fill the inside of density plot with a color using the `fill=` argument. Try at least four different numbers of smoothing `bw` (e.g., 0.2, 1.5, 3, 5) by manipulating the `bw=` argument. Select what you think best represents the data for each. Provide a brief justification for your decision.

## Barplot

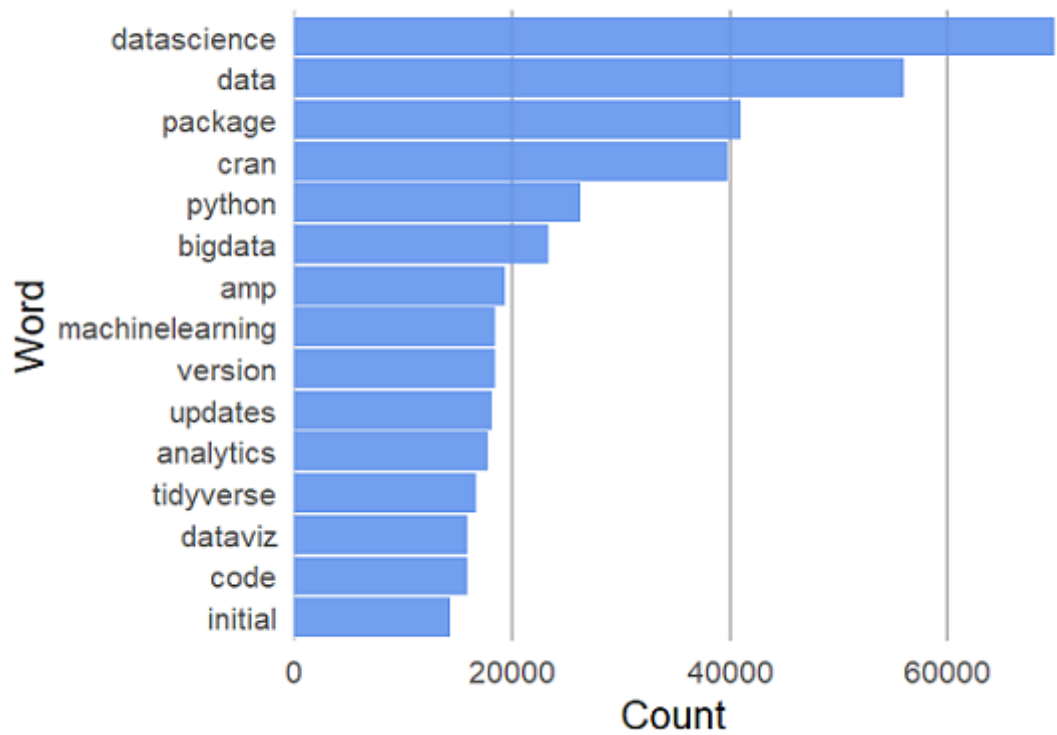
3. Using the information `text` column, create the following figure of the 15 most common words represented in these posts by using the `ggplot2()` package and `geom_col()` function. Remove the stop words, and also exclude the words such as 't.co', 'https', 'http', 'rt', 'rstats'.



4. Style the plot so it (mostly) matches the below. It does not need to be exact, but it should be close.

## Word frequencies in posts

Top 15 words displayed



Data from Mike Kearny, distributed via #tidytuesday