# goJitsu

## Cade Michael Lueker

## 12/29/2021

## Goal and Methods

- After taking class in R and basic data analytics I thought that it would be fun to try and find / scrape a data set of my own. Not one that was given or loaded from elsewhere. I decided that I would try to get my hands on some Jiu Jitsu data and turned to **BJJHeroes** to do so.
- as of right now my data is comprised of a
  - name
  - number of wins
  - number of losses
  - number per type of win
  - number per type of loss
- this data is relatively unintersting as it is fairly simple and doesn't allow for much prediction, ie predicting if one fighter would beat another.
- But scraping the data was a foreign task and I will try to get more detailed data in the near future to update and or re write this report.

## Data Collection

- to collect the data I decided to try **Go** as it is a language that I am very interested in and want to become more familiar with.
  - I utilized the go colly library and found that it offered a very clean and consise meta for web scraping.
  - I visit a page that lists fighters and every fighter's individual page gets visited for individual data collection
  - if this fighter has wins or losses (some fighters only have descriptions) their unique data gets entered into a CSV engendered by the program. *the go code is at the end of this report*

## Possible observations

- The only real observations I can make with the data I collected is to see the distribution of different win & loss types. I can compare the frequencies of each, I can determine what a median and or average competitor (who is well known enough for BJJ heroes) looks like in terms of wins and losses
- I will rely mainly on histograms as they are utilized to measure frequency which is what I can analyze.

## Load Required packages

```
# load libraries
library(tidyverse)
```

```
## -- Attaching packages --------------------------------------- tidyverse 1.3.1 --
```

```
## v ggplot2 3.3.5     v purrr   0.3.4
## v tibble  3.1.6     v dplyr   1.0.7
## v tidyr   1.1.4     v stringr 1.4.0
```

```
## v readr    2.1.0      v forcats 0.5.1
```

```
## -- Conflicts ----------------------------------------- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
```

```r
library(ggplot2)
library(vioplot)
```

```
## Loading required package: sm
```

```
## Package 'sm', version 2.2-5.7: type help(sm) for summary information
```

```
## Loading required package: zoo
```

```
##
## Attaching package: 'zoo'
```

```
## The following objects are masked from 'package:base':
##
##      as.Date, as.Date.numeric
```

## Read The data from the CSV

```r
goData <- read.csv("./workingFiles/winTypes.csv")
```
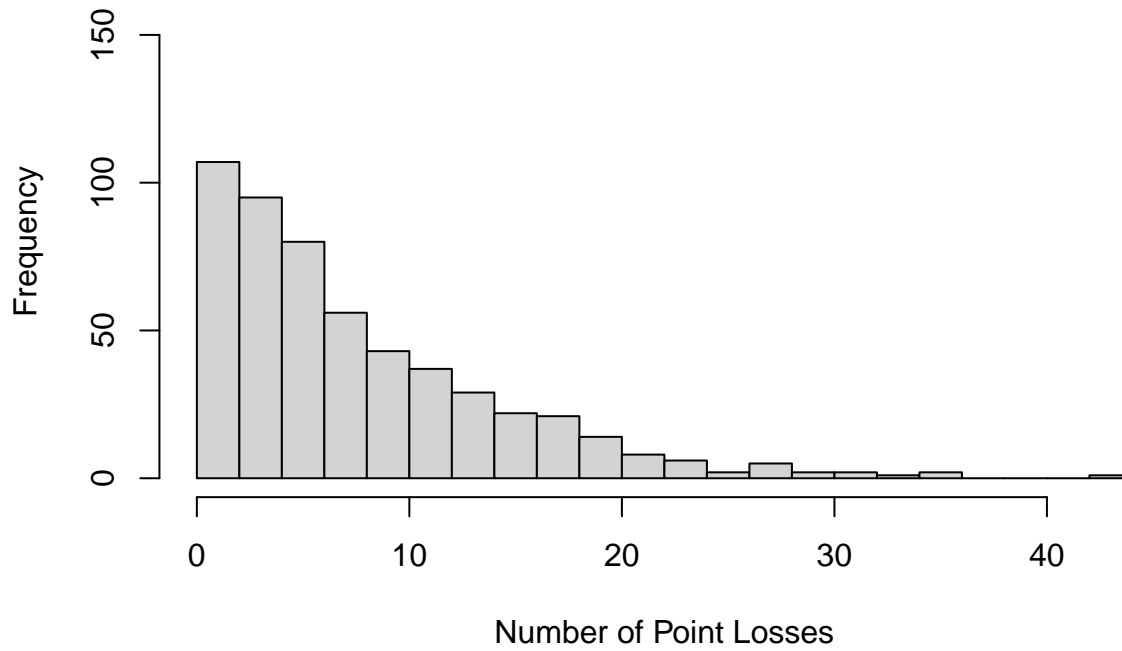
## Clean data

```r
goData$wins <- gsub("[A-Z]", "", goData$wins)
goData$losses <- gsub("[A-Z]", "", goData$losses)
nonNumer <- names(goData)[2:length(goData)]
goData[nonNumer] <- lapply(goData[nonNumer], as.numeric)
# add total number of matches
goData <- goData %>%
  mutate(total = losses + wins)
# convert to R df
winLoss_df <- as.data.frame(goData)

# add number of matches won not by submission
winLoss_df <- winLoss_df %>%
  mutate(nonSubWins = wins - Win.Submissions)
# separate losses and wins
loss_df <- winLoss_df[11:length(winLoss_df)-1]
wins_df <- winLoss_df[4:10]
```

## Initial Plots on match data, Histograms
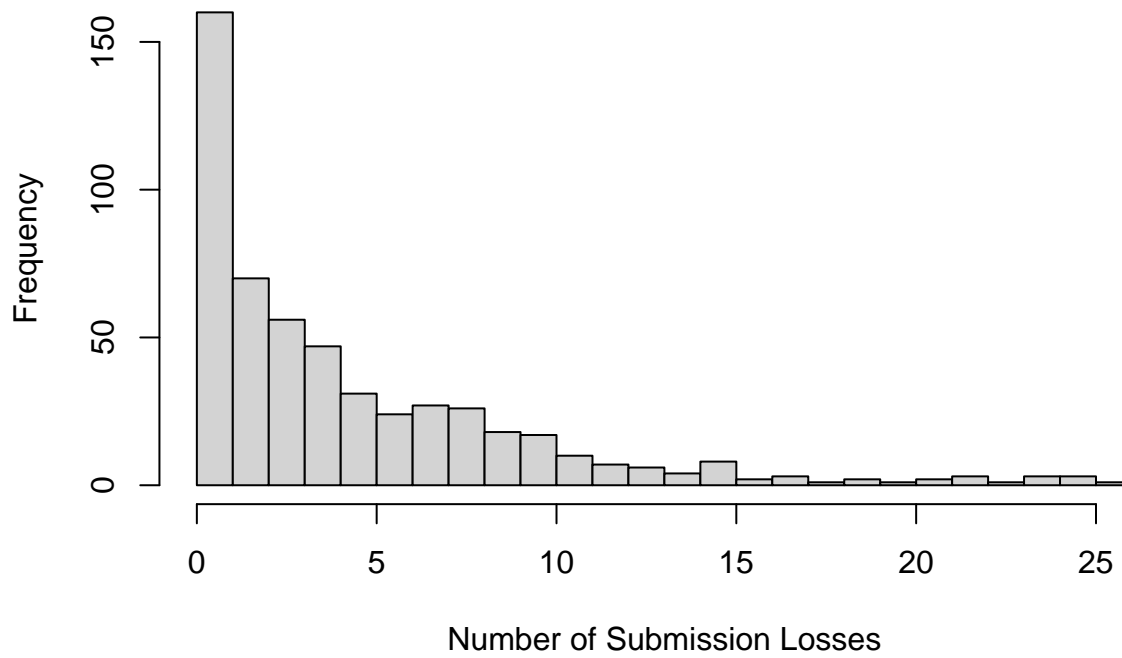
```r
# Basic Histograms
hist(winLoss_df$Losses.Points,
     ylim = range(0,160),
     xlab = "Number of Point Losses",
     main = "Frequency of Point Losses",
     breaks = 25)
```

**Frequency of Point Losses**
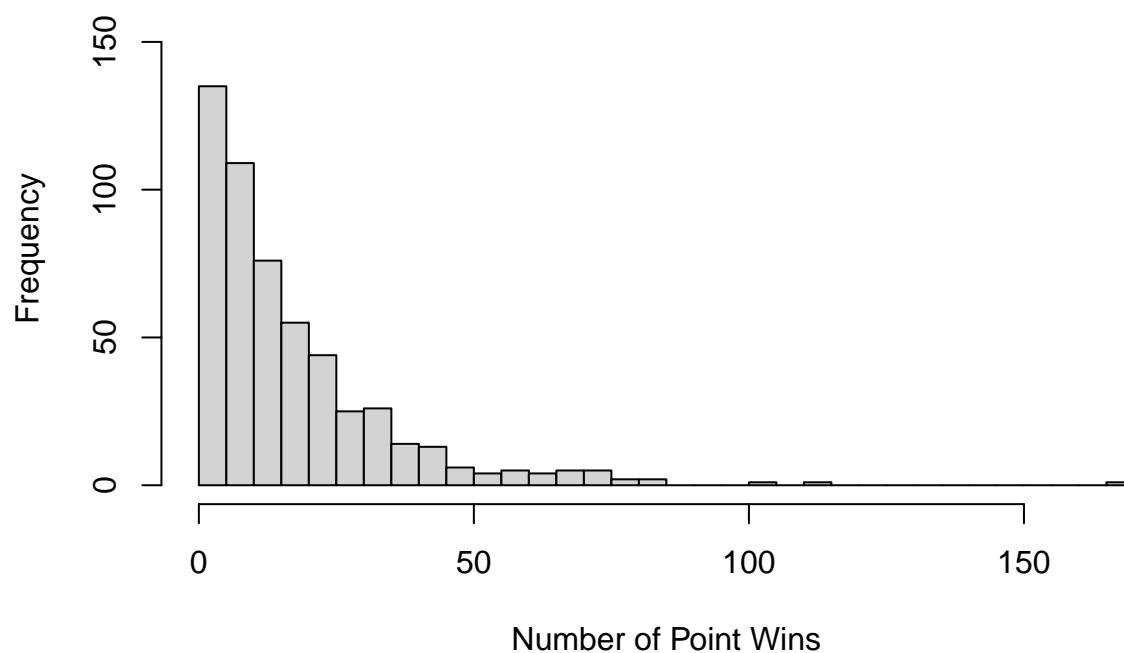


```
hist(winLoss_df$Losses.Submissions,
     ylim = range(0,160),
     xlab = "Number of Submission Losses",
     main = "Frequency of Submission Losses",
     breaks = 25)
```

## Frequency of Submission Losses
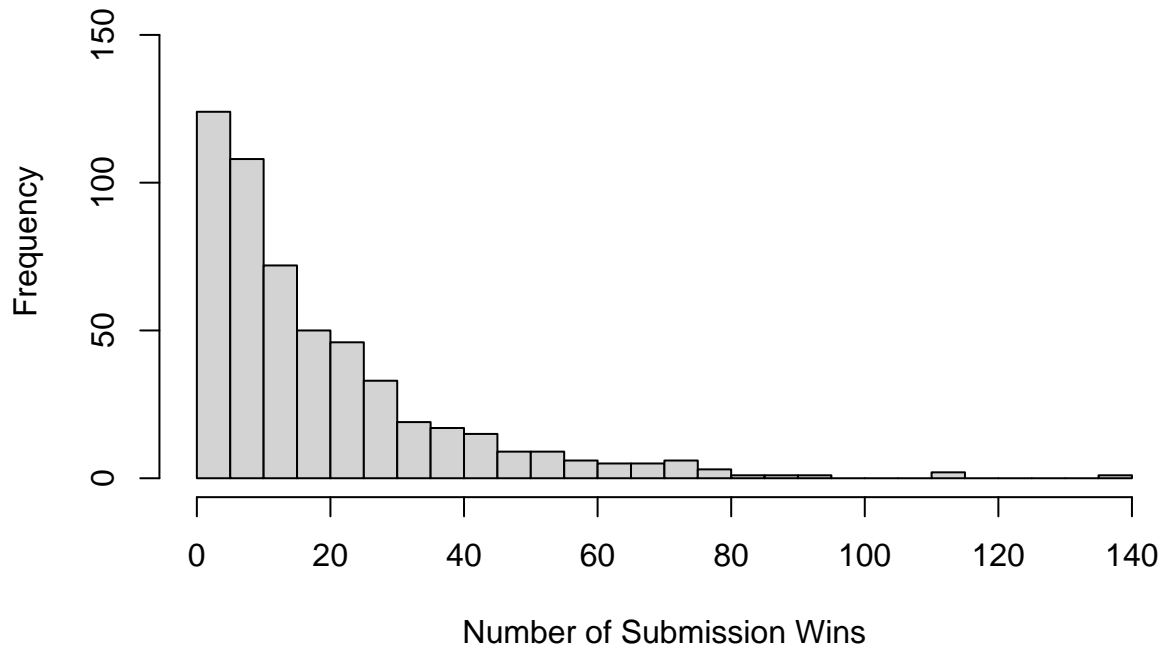


```
hist(winLoss_df$Win.Points,
    ylim = range(0,160),
    xlab = "Number of Point Wins",
    main = "Frequency of Point Wins",
    breaks = 25)
```

# Frequency of Point Wins



```
hist(winLoss_df$Win.Submissions,
     ylim = range(0,160),
     xlab = "Number of Submission Wins",
     main = "Frequency of Submission Wins",
     breaks = 25)
```

# Frequency of Submission Wins



```
# number of losses for Subs and Points
sum(winLoss_df$Losses.Submissions)
```

```
## [1] 2517
```

```
sum(winLoss_df$Losses.Points)
```
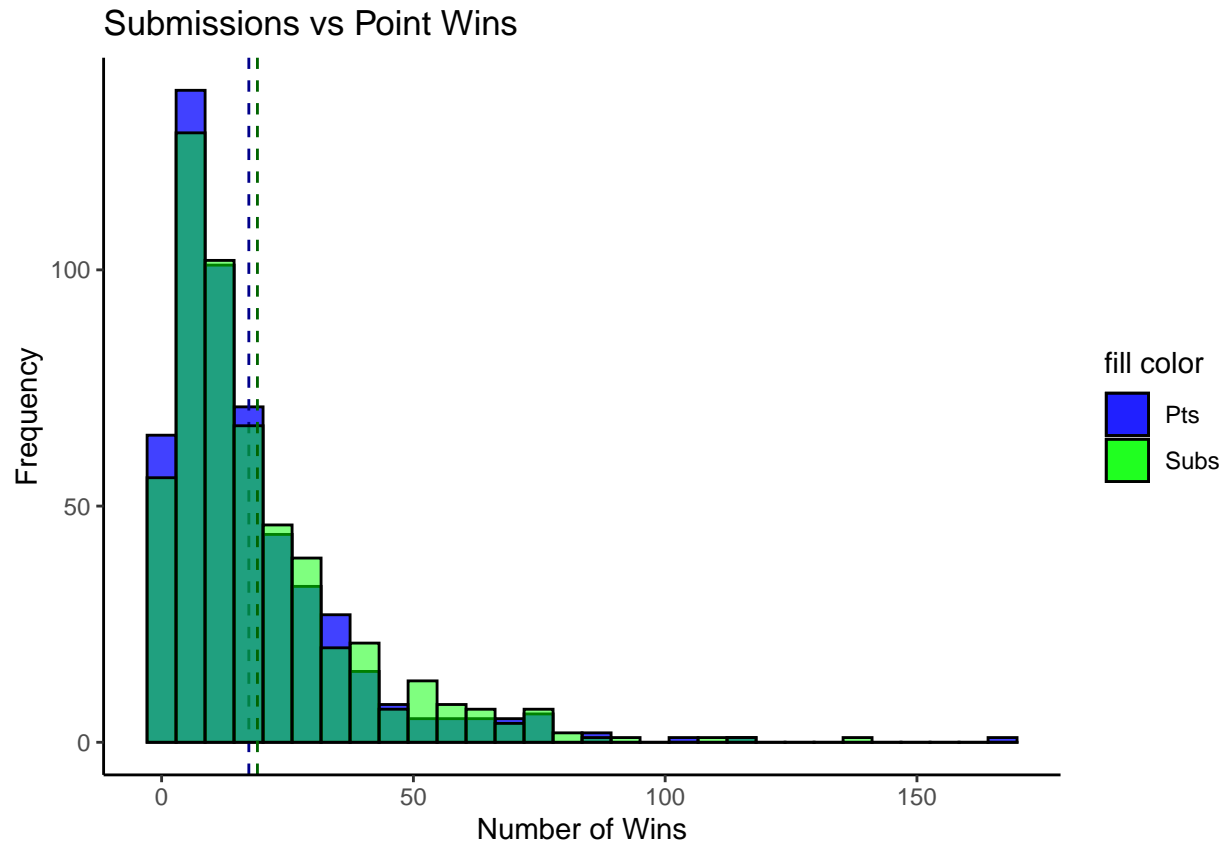
```
## [1] 4253
```

## Box Plots For wins

```
# sub wins vs points wins
winBar <- winLoss_df %>%
  ggplot() +
    geom_histogram(aes(x = Win.Points,fill= "blue"), color = "black", alpha = 0.75) +
    geom_vline(aes(xintercept = mean(Win.Points)), color = "dark blue", linetype = "dashed") +
    geom_histogram(aes(x = Win.Submissions, fill= "green"), color = "black", alpha = 0.5) +
    geom_vline(aes(xintercept = mean(Win.Submissions)), color = "dark green", linetype = "dashed") +
    labs(x = ("Number of Wins"), y = ("Frequency"), title = "Submissions vs Point Wins") +
    scale_fill_identity(name = "fill color", guide = "legend", labels = c("Pts", "Subs")) +
    theme_classic()
# plot
winBar
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```
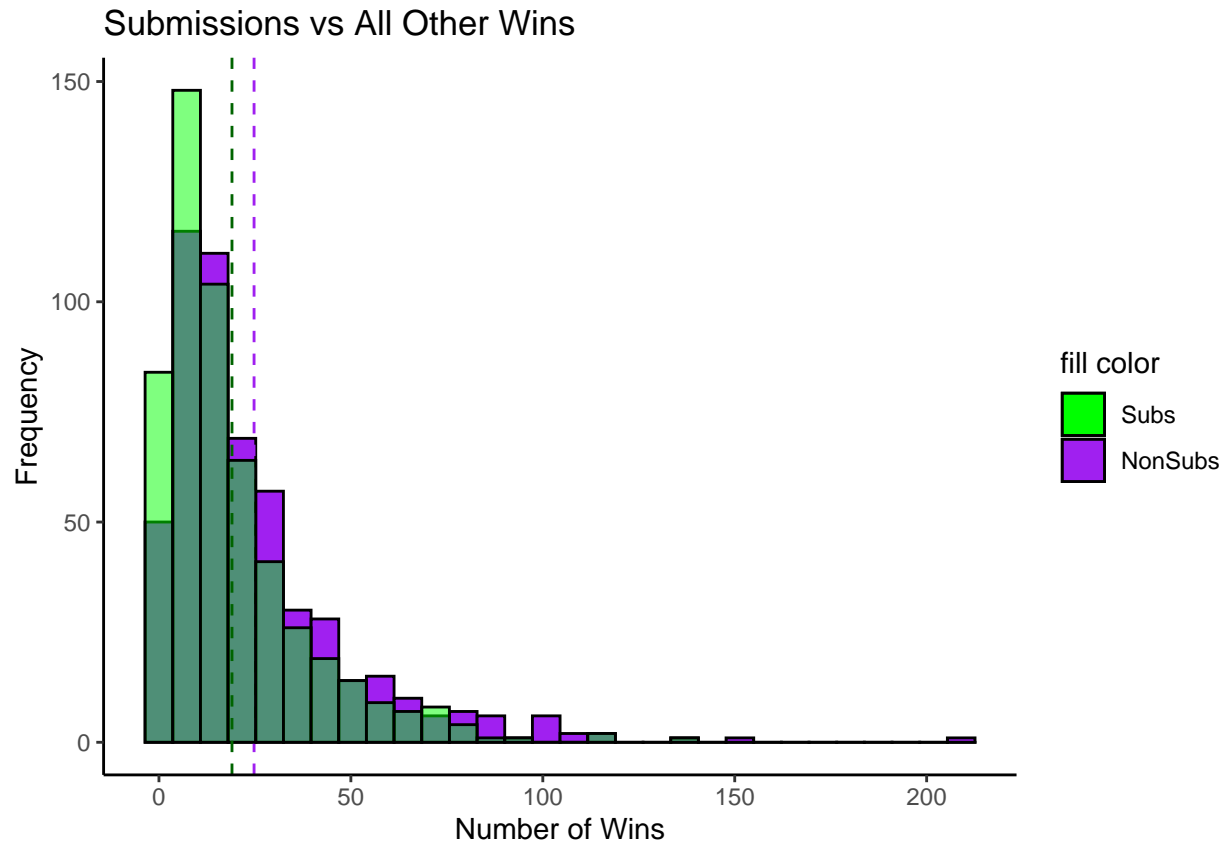
## Submissions vs Point Wins



```r
# sub wins vs all other types of wins
subVsNon <- winLoss_df %>%
  ggplot() +
    geom_histogram(aes(x = nonSubWins, fill= "purple"), color = "black", alpha = 1) +
    geom_vline(aes(xintercept = mean(nonSubWins)), color = "purple", linetype = "dashed") +
    geom_histogram(aes(x = Win.Submissions, fill= "green"), color = "black", alpha = 0.5) +
    geom_vline(aes(xintercept = mean(Win.Submissions)), color = "dark green", linetype = "dashed") +
    labs(x = ("Number of Wins"), y = ("Frequency"), title = "Submissions vs All Other Wins") +
    scale_fill_identity(name = "fill color", guide = "legend", labels = c("Subs", "NonSubs")) +
    theme_classic()

# plot
subVsNon
```

```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

## Submissions vs All Other Wins
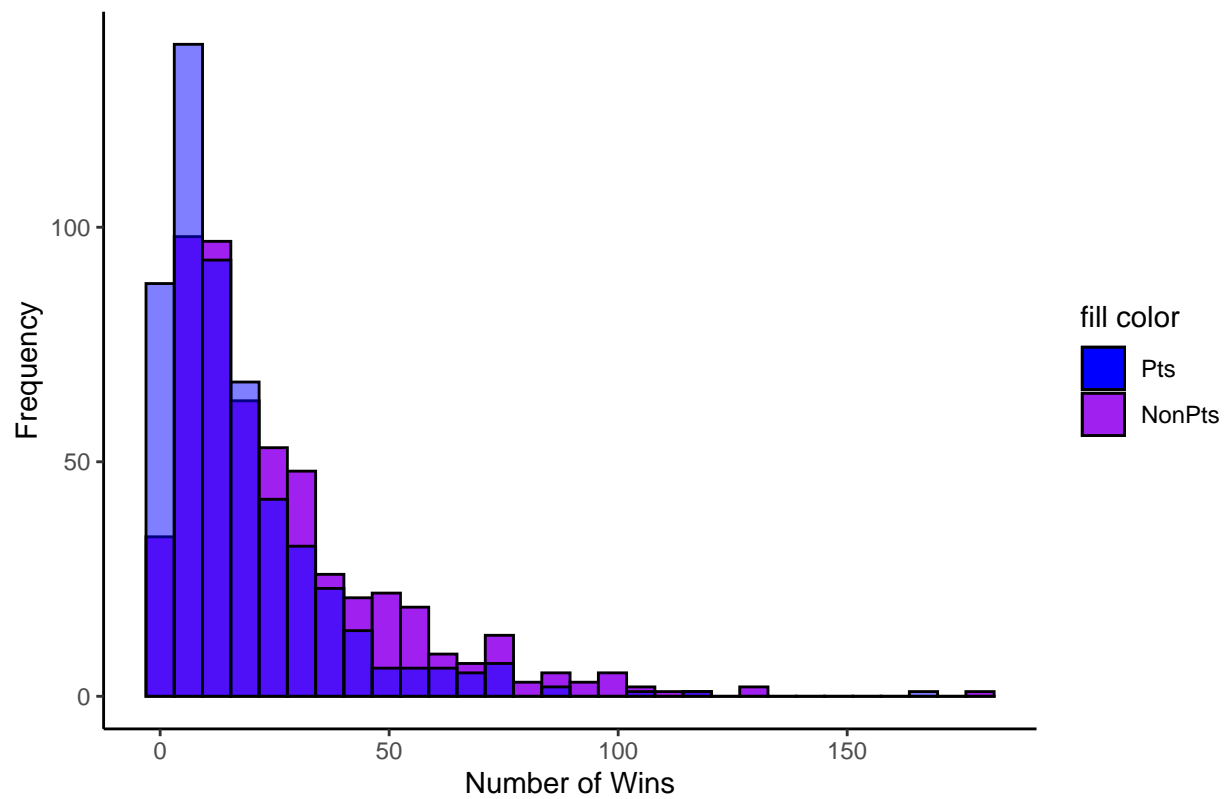


```
# wins by points vs all other types of wins
PtsVsNon <- winLoss_df %>%
  ggplot() +
    geom_histogram(aes(x = wins - Win.Points, fill= "purple"), color = "black", alpha = 1) +
    geom_histogram(aes(x = Win.Points, fill= "blue"), color = "black", alpha = 0.5) +
    labs(x = ("Number of Wins"), y = ("Frequency"), title = "Points vs All Other Wins") +
    scale_fill_identity(name = "fill color", guide = "legend", labels = c("Pts", "NonPts")) +
    theme_classic()

# plot
PtsVsNon
```
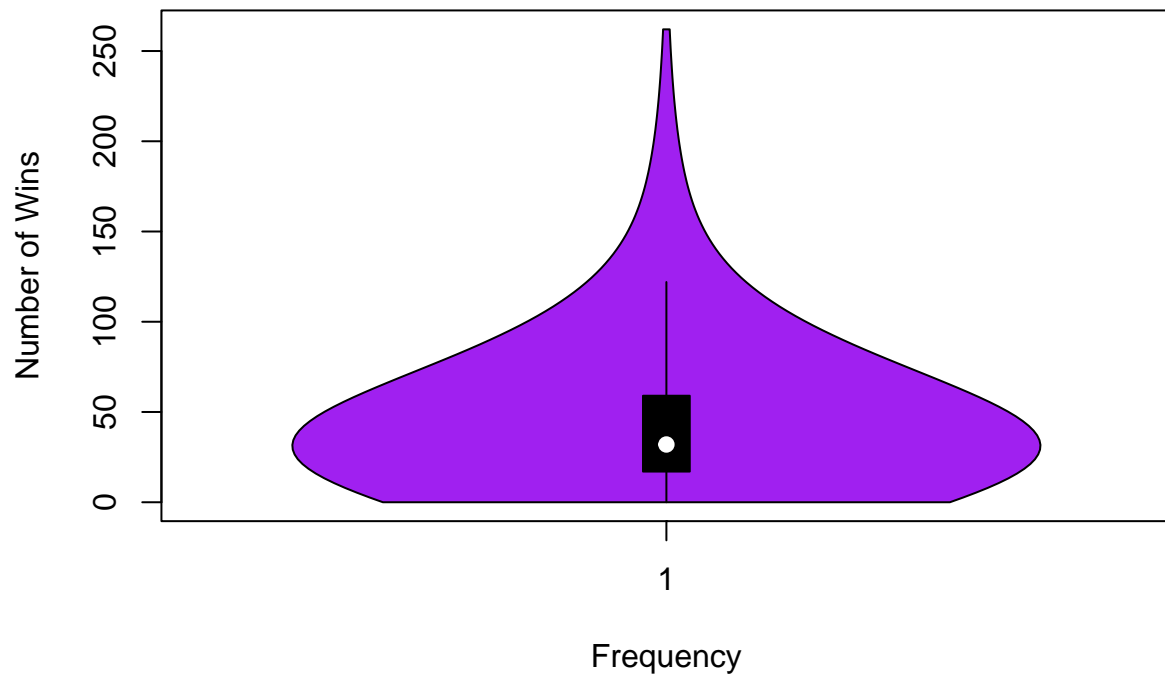
```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

## Points vs All Other Wins



## Violin Plots

```r
win_violin <- vioplot(winLoss_df$wins,
                      col = "purple",
                      xlab = "Frequency",
                      ylab = "Number of Wins")
```

```
win_violin
```

```
## $upper
## [1] 122
##
## $lower
## [1] 0
##
## $median
## [1] 32
##
## $q1
## [1] 17
##
## $q3
## [1] 59
```

```
sub_violin <- vioplot(winLoss_df$Win.Submissions,
                      col = "green",
                      xlab = "Frequency",
                      ylab = "Number of Subs")
```
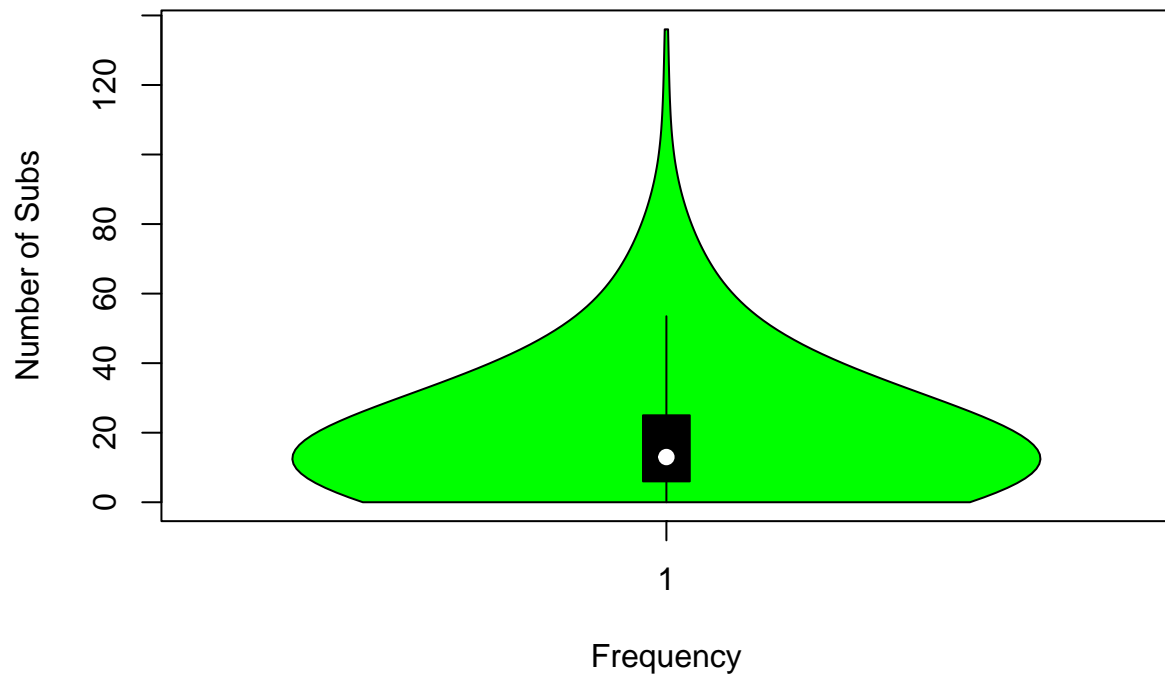
```
sub_violin
```

```
## $upper
## [1] 53.5
##
## $lower
## [1] 0
##
## $median
## [1] 13
##
## $q1
## [1] 6
##
## $q3
## [1] 25
```

## Go code for scraping

```go
package main

import (
    "encoding/csv"
    "fmt"
    "log"
    "os"
```

```go
        "github.com/gocolly/colly"
)

// main function
func main() {
    createList()
    visitAll()
    indiv("https://www.bjjheroes.com/?p=7556")
}

// visitAll function
func visitAll() {
    /* instatiate colly */
    c := colly.NewCollector(
        colly.AllowedDomains(
            "www.bjjheroes.com/",
            "bjjheroes.com/",
            "https://bjjheroes.com/",
            "www.bjjheroes.com",
            "bjjheroes.com",
            "https://bjjheroes.com",
        ),
    )
    // if a link is found in a table
    c.OnHTML("td.column-1 a[href]", func(e *colly.HTMLElement) {
        url := e.Request.AbsoluteURL(e.Attr("href")) // get the link as an absolute
        indiv(url)                                   // run indiv on the absolute link
    })
    // start scraping
    c.Visit("https://www.bjjheroes.com/a-z-bjj-fighters-list")
}

// createList function
func createList() {
    // instantiate csv
    file, err := os.OpenFile("list.csv", os.O_APPEND|os.O_CREATE|os.O_WRONLY, 0644)
    if err != nil { // deal with errors
        log.Fatalf("cannot create file due to %s", err)
        return
    }
    defer file.Close()
    writer := csv.NewWriter(file) // create a writer to modify csv
    // write column names to csv
    writer.Write([]string{
        "Name",
        "wins",
        "losses",
        "Win Points",
        "Win Advantages",
        "Win Submissions",
        "Win Decision",
        "Win Penalties",
        "Win EBI / OT",
```

```go
        "Win DQ",
        "Losses Points",
        "Losses Advantages",
        "Losses Submissions",
        "Losses Decision",
        "Losses Penalties",
        "Losses EBI / OT",
        "Losses DQ",
    })
    writer.Flush() // ensure data is written
}

// indiv function
func indiv(url string) {
    // instantiate or open csv
    file, err := os.OpenFile("list.csv", os.O_APPEND|os.O_CREATE|os.O_WRONLY, 0644)
    if err != nil {
        log.Fatalf("cannot create file due to %s", err)
        return
    }
    defer file.Close()
    writer := csv.NewWriter(file) // create a writer to modify csv
    defer writer.Flush()          // defer pushes to end of function
    /* instatiate colly */
    c := colly.NewCollector(
        colly.AllowedDomains(
            "www.bjjheroes.com/",
            "bjjheroes.com/",
            "https://bjjheroes.com/",
            "www.bjjheroes.com",
            "bjjheroes.com",
            "https://bjjheroes.com",
        ),
    )
    // create empty column variables
    // using standard vs short declaration as short is a bit messy in this case
    fName := "0"
    wins := "0"
    lose := "0"
    wPoints := "0"
    wAdvantage := "0"
    wSubs := "0"
    wDec := "0"
    wPen := "0"
    wEBI := "0"
    wDQ := "0"
    lPoints := "0"
    lAdvantage := "0"
    lSubs := "0"
    lDec := "0"
    lPen := "0"
    lEBI := "0"
    lDQ := "0"
```

```go
// find name
  c.OnHTML("h1", func(a *colly.HTMLElement) {
      fName = a.Text
  })
//find wins
  c.OnHTML("div.Win_title", func(b *colly.HTMLElement) {
      wins = b.ChildText("em")
  })
// find win types
  c.OnHTML("div.wrapper_canvas li", func(d *colly.HTMLElement) {
      switch d.ChildText("span.by_points") {
      case "BY POINTS":
          wPoints = d.ChildText("span.per_no")
      case "BY ADVANTAGES":
          wAdvantage = d.ChildText("span.per_no")
      case "BY SUBMISSION":
          wSubs = d.ChildText("span.per_no")
      case "BY DECISION":
          wDec = d.ChildText("span.per_no")
      case "BY PENALTIES":
          wPen = d.ChildText("span.per_no")
      case "BY EBI/OT":
          wEBI = d.ChildText("span.per_no")
      case "BY DQ":
          wDQ = d.ChildText("span.per_no")
      }
  })
// find losses
c.OnHTML("div.Win_title_lose", func(t *colly.HTMLElement) {
  lose = t.ChildText("em")
})
// find loss types
  c.OnHTML("div.wrapper_canvas_lose li", func(d *colly.HTMLElement) {
      switch d.ChildText("span.by_points") {
      case "BY POINTS":
          lPoints = d.ChildText("span.per_no_lose")
      case "BY ADVANTAGES":
          lAdvantage = d.ChildText("span.per_no_lose")
      case "BY SUBMISSION":
          lSubs = d.ChildText("span.per_no_lose")
      case "BY DECISION":
          lDec = d.ChildText("span.per_no_lose")
      case "BY PENALTIES":
          lPen = d.ChildText("span.per_no_lose")
      case "BY EBI/OT":
          lEBI = d.ChildText("span.per_no_lose")
      case "BY DQ":
          lDQ = d.ChildText("span.per_no_lose")
      }
  })


// start the scraper
  c.Visit(url)
```

```go
    // write to csv
    if wins != "0" || lose != "0" {
        writer.Write([]string{
            fName,
            wins,
            lose,
            wPoints,
            wAdvantage,
            wSubs,
            wDec,
            wPen,
            wEBI,
            wDQ,
            lPoints,
            lAdvantage,
            lSubs,
            lDec,
            lPen,
            lEBI,
            lDQ,
        })
    }
    // for every individual fighter print done
    fmt.Printf("!!!Done!!!")
}
```