# Meet a New R Package: scrapeR

Joe Walsh

May 1, 2013

# What scrapeR Does

**Directly**

- downloads & preprocesses webpages
- saves as a list
- one function: scrape()

# What scrapeR Does

**Directly**

- downloads & preprocesses webpages
- saves as a list
- one function: scrape()

- Dependency package: XML

# What scrapeR Does

**Directly**

- downloads & preprocesses webpages
- saves as a list
- one function: scrape()

- Dependency package: XML

**Indirectly**

- only requires three lines of code
- can be automated
- dynamic reports
- improved replicability

# Why I Like It

- huge amounts of data
- new data being released
- can stop static reports
- eases replication

# Example: Cherry Hypothesis

# Example: Cherry Hypothesis

# Example: Cherry Hypothesis

# Example: Cherry Hypothesis

```
# install and load scrapeR package
install.packages("scrapeR", dependencies=TRUE)
library(scrapeR)
```

# Example: Cherry Hypothesis

```
# install and load scrapeR package
install.packages("scrapeR", dependencies=TRUE)
library(scrapeR)

# scrape data
URLs <- c("http://www.nhl.com/ice/standings.htm?type=lea",
          "http://www.hockeyfights.com/leaders/teams/")
pageSource <- scrape(url=URLs, headers=FALSE, parse=TRUE)
```

# Example: Cherry Hypothesis

```
# install and load scrapeR package
install.packages("scrapeR", dependencies=TRUE)
library(scrapeR)

# scrape data
URLs <- c("http://www.nhl.com/ice/standings.htm?type=lea",
          "http://www.hockeyfights.com/leaders/teams/")
pageSource <- scrape(url=URLs, headers=FALSE, parse=TRUE)

# get two tables with info we need
NHL.tables <- readHTMLTable(pageSource[[1]])
  team.records <- NHL.tables[[3]]
HockeyFights.tables <- readHTMLTable(pageSource[[2]])
  team.fights <- HockeyFights.tables[[1]]
```
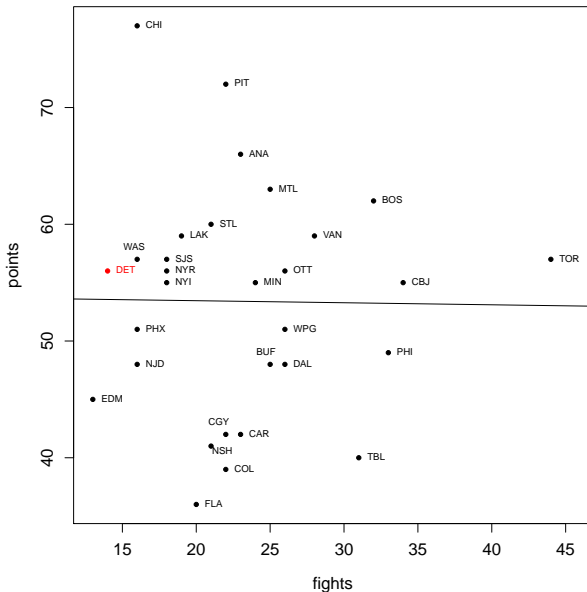
# Example: Cherry Hypothesis

# Questions?

Email: j.thomas.walsh@gmail.com
GitHub: jtwalsh0