

# Class 2: Introduction to R and categorical data



# Overview

Quick review of central statistics concepts and Quarto

Introduction to R

Categorical data

- Proportions
- Frequency tables
- Bar charts and pie plots
- If there is time: Categorical data in R

# Announcements

If you haven't done so yet, please remember to:

1. Fill out the background survey
2. Fill out the practice session survey

Also, please log into the RStudio server now

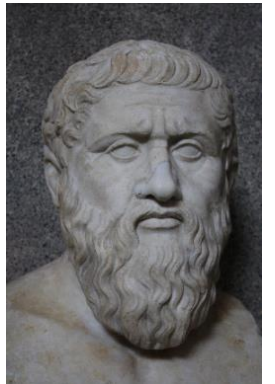
Any questions about anything?

REVIEW

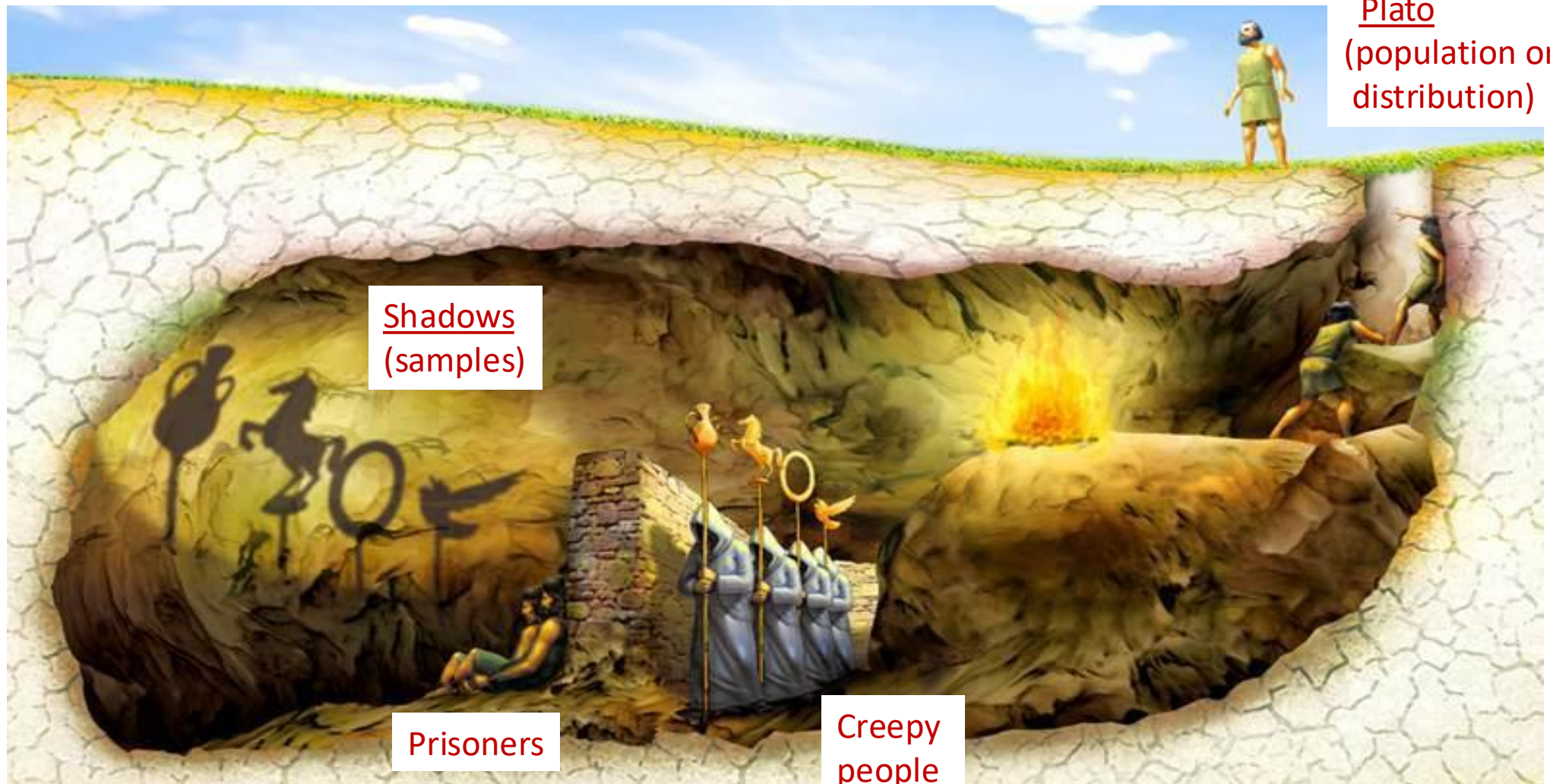
# Quiz time!

(not to be turned in)

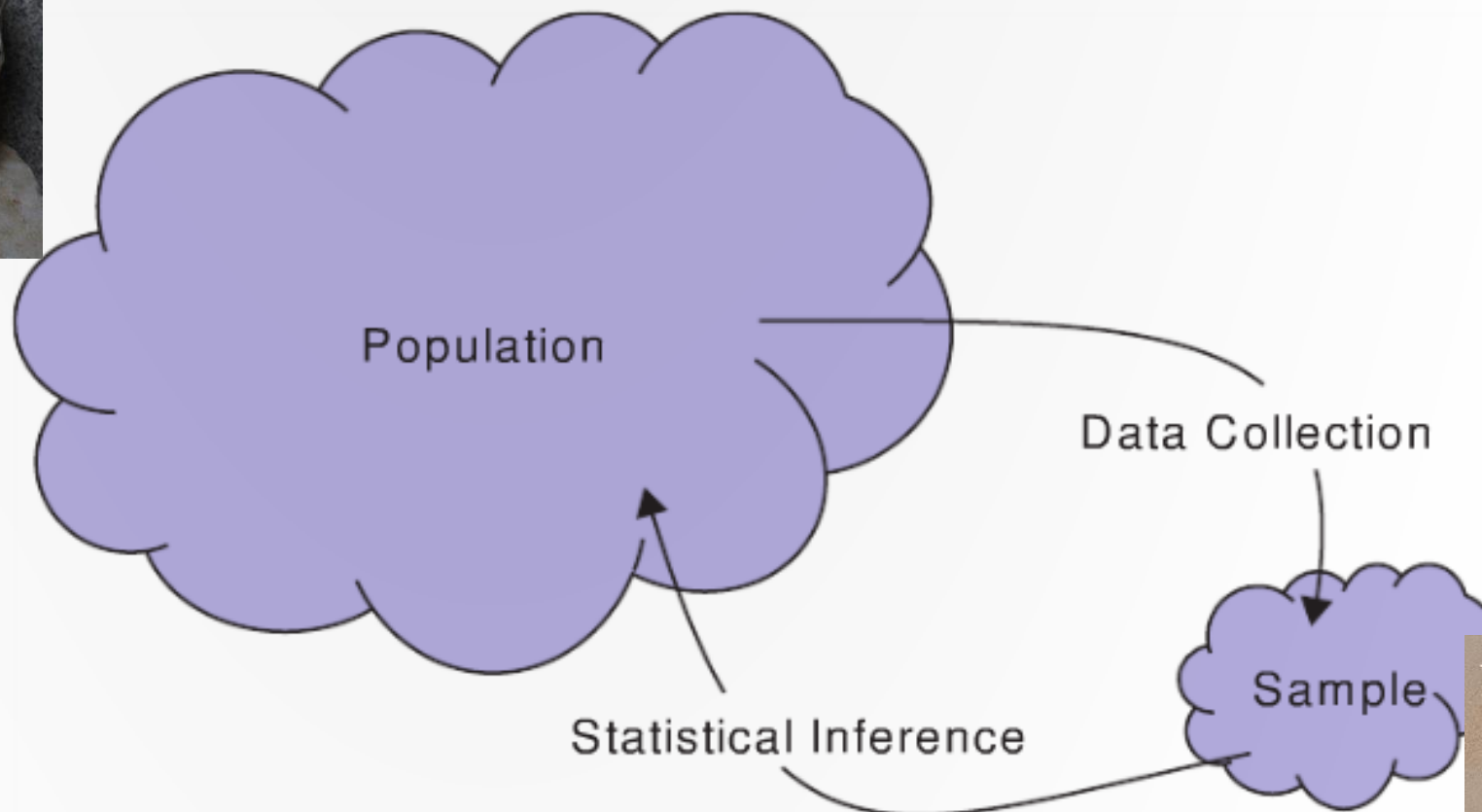
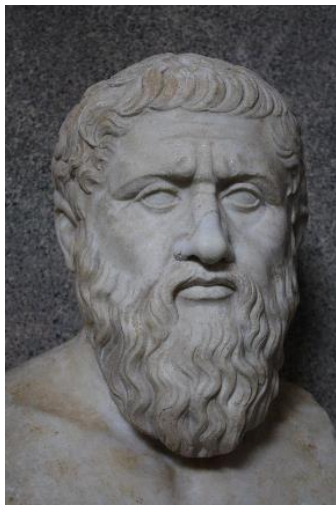
1. What is a population?
2. What is a sample?
3. What is statistical inference?
4. What are the rows of a data table called?
5. What are the columns of a data table called?
6. What is the difference between categorical and quantitative variables?
7. Who is this?



# Plato's cave



From The Republic (~ 380 BCE)



# Review Quarto

Quarto (.qmd files) allow you to embed written descriptions, R code and the output to create a reproducible research document!



Everything in R chunks is executed as code:

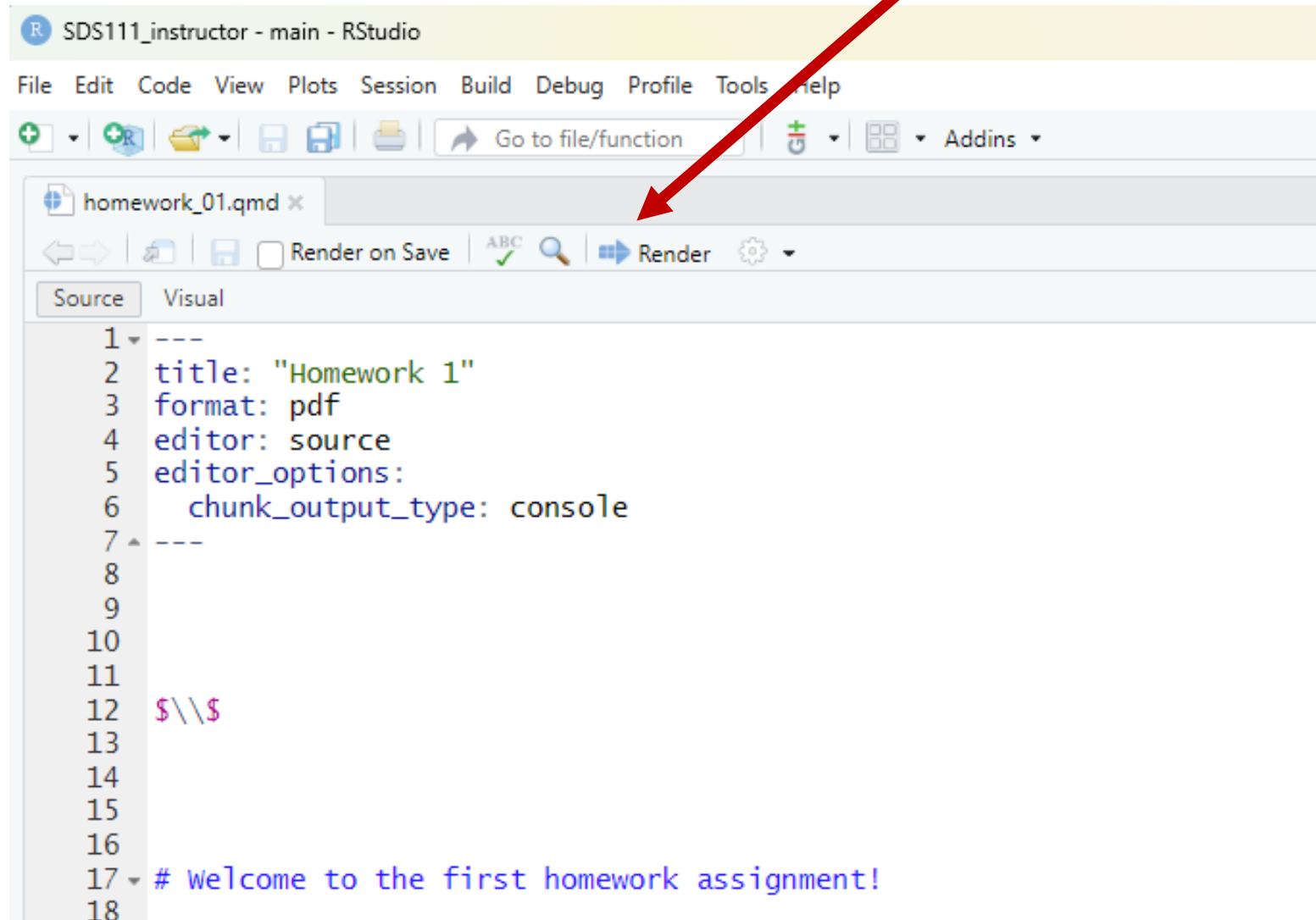
```
```${r}
  # this is a comment
  # the following code will be executed
  2 + 3
```
```

Everything outside R chunks appears as text



# Render to a pdf

Renders to a pdf document  
(which you will submit to Gradescope)



# Formatting in Quarto

We can add formatting to text outside the code chunks

Examples:

`## Level 2 header`

**Level 2 header**

---

`**bold**`

**bold**

LaTeX {  
`$\pi$`  
`$x_{outcome}$`

$\pi$

$x_{outcome}$

# Avoid hard to debug code!

Only change a few lines at a time and then render your document to make sure everything is working!

**I.e., render your to pdf document often!**

Let's quickly try it in R...



Questions?



Introduction to



# Getting class code and the homework

To load the class functions use:

```
> library(SDS1000)
```

To get the class 2 material, on the console type:

```
> goto_class(2)
```

# R and R Studio

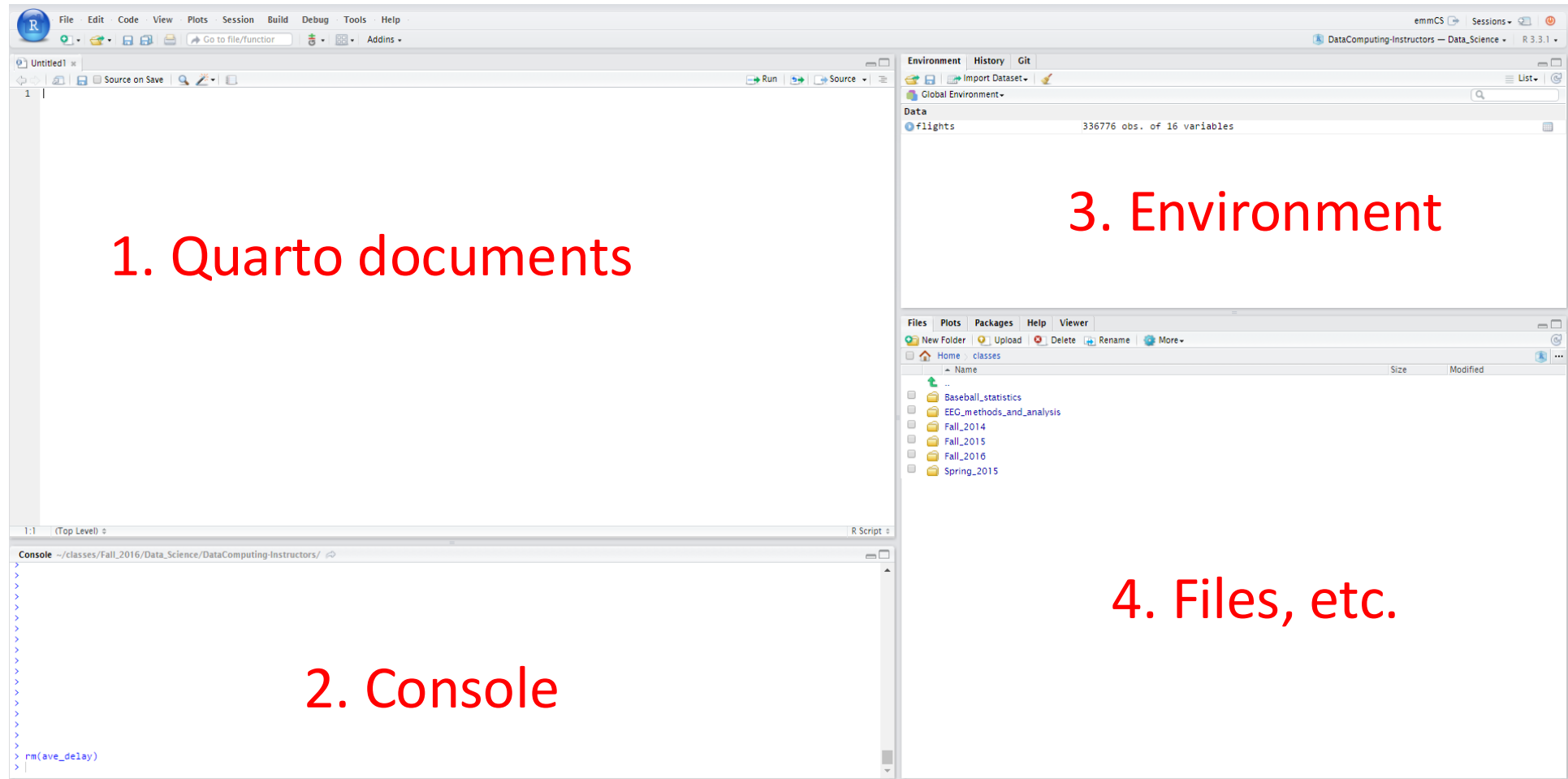
**R: Engine**



**RStudio: Dashboard**

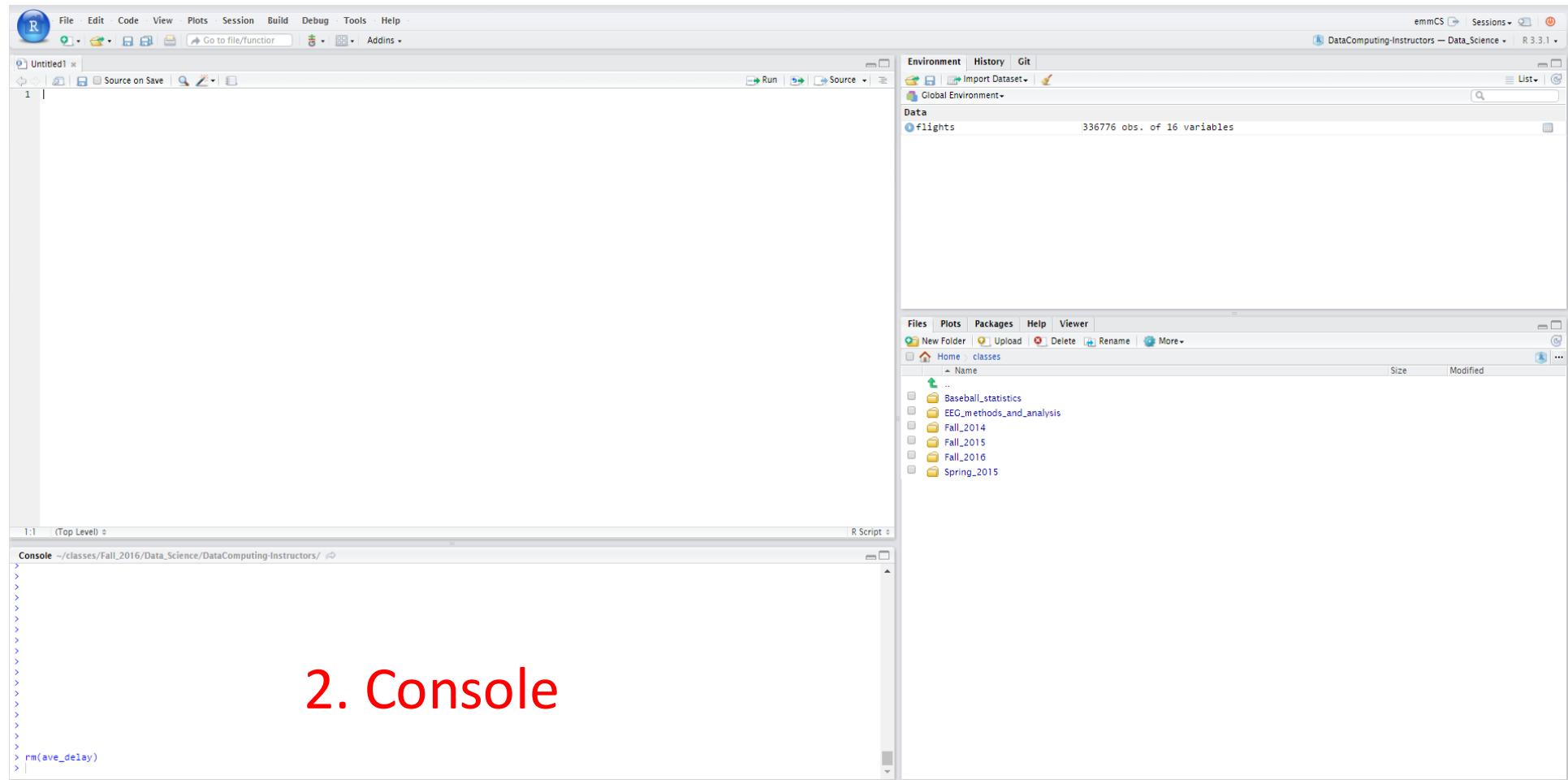


# RStudio layout





# RStudio layout



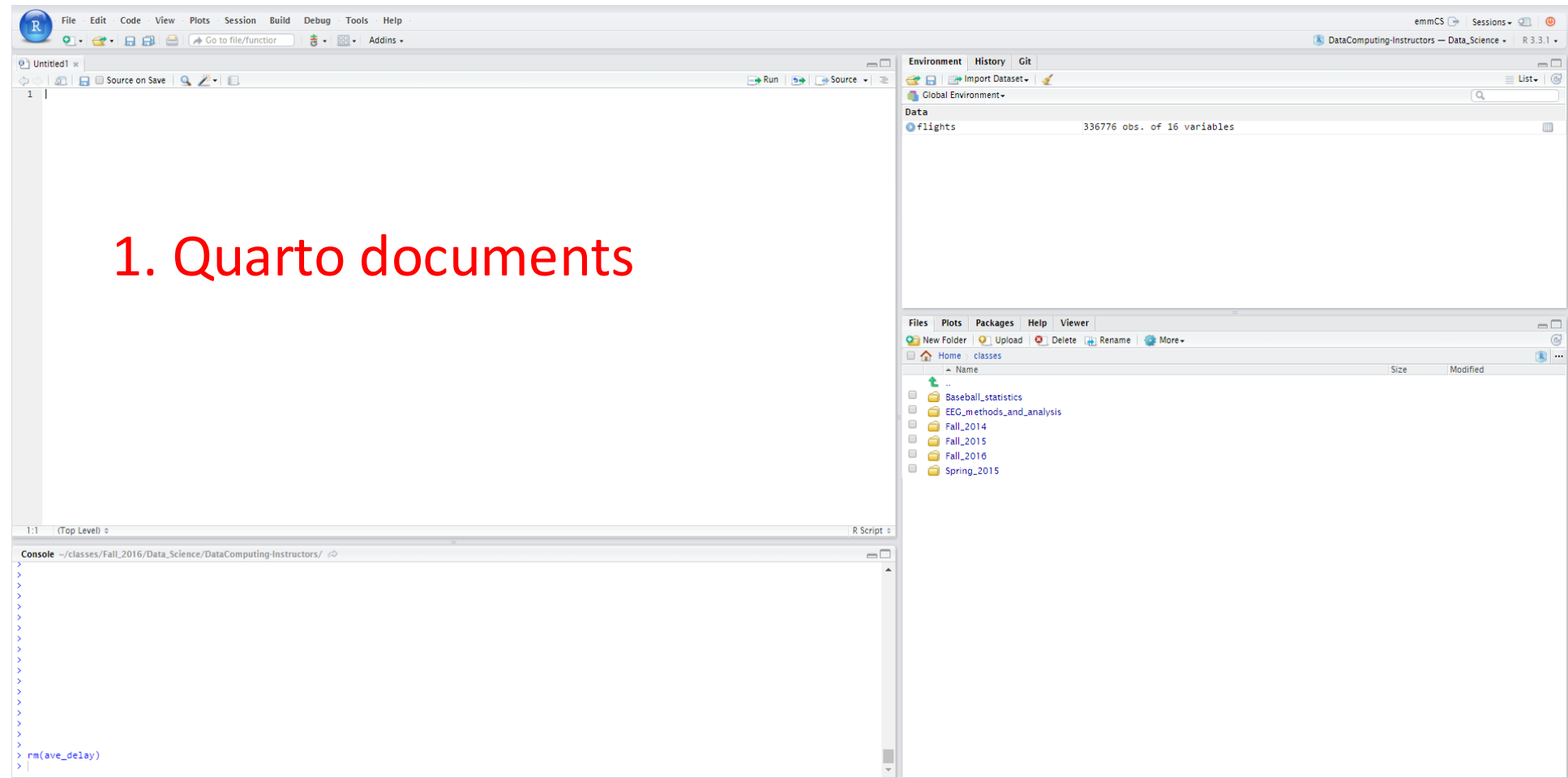
## 2. Console

### R as a calculator

`> 2 + 3`

`> 7 * 5`

# RStudio layout



# R Basics

Please follow along in RStudio!

Arithmetic:

```
> 2 + 3
```

```
> 7 * 5
```

Assignment:

```
> a <- 4
```

```
> b <- 7
```

```
> z <- a + b
```

```
> z
```

```
[1] 11
```

# Character strings and Booleans

```
> a <- 7
```

```
> s <- "Statistics is great!"
```

```
> b <- TRUE
```

```
> class(a)
```

```
[1] numeric
```

```
> class(s)
```

```
[1] character
```

# Functions

Functions use parenthesis: functionName(x)

```
> sqrt(49)
```

```
> tolower("DATA is AWESOME!")
```

To get help

```
> ? sqrt
```

One can add comments to your code

```
> sqrt(49)  # this takes the square root of 49
```

# Vectors

Vectors are ordered sequences of numbers or letters

The `c()` function is used to create vectors

```
> v <- c(5, 232, 5, 543)
```

```
> s <- c("these", "are", "strings")
```

One can access elements of a vector using square brackets `[]`

```
> s[3]      # what will the answer be?
```

# Vectors continued

One can assign a sequence of numbers to a vector

```
> z <- 2:10
```

```
> z[3]
```

We can apply functions to vectors

```
> sqrt(z)    # returns a vector
```

```
> sum(z)     # returns a single number
```

We can see how many elements are in a vector using the `length()` function

```
> length(z)
```

# Categorical variables



# The sprinkle business

(fictional)

**ACME**  
CORPORATION



PERFECT  
Corporation



ACME corporation believes that if they had the correct ratio (proportion) of red sprinkles that PERFECT corporation uses, their sales will increase

# Where do samples/data come from?

To assess the proportion of sprinkles that PERFECT corporation uses, AMCE sampled 100 of PERFECT corporation's sprinkles

- The **sample size** is 100 ( $n = 100$ )



|   |        |
|---|--------|
| 1 | orange |
| 2 | red    |
| 3 | green  |
| 4 | white  |
| 5 | white  |
| 6 | white  |
| 7 | white  |
| 8 | white  |
| 9 | red    |

# Sampling example



## Questions:

1. What are the observational units (cases)?
2. What is the variable?
3. Is the variable categorical or quantitative?
4. What is the population?
5. Do you think the samples we are getting are representative of the population?

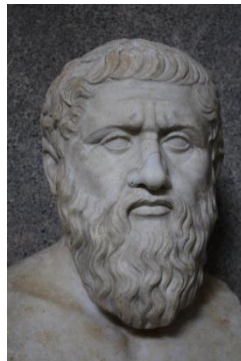
|   |        |
|---|--------|
| 1 | orange |
| 2 | red    |
| 3 | green  |
| 4 | white  |
| 5 | white  |
| 6 | white  |
| 7 | white  |
| 8 | white  |
| 9 | red    |

# Population parameters vs. sample statistics

A **statistic** is a number that is computed from ***data in a sample***

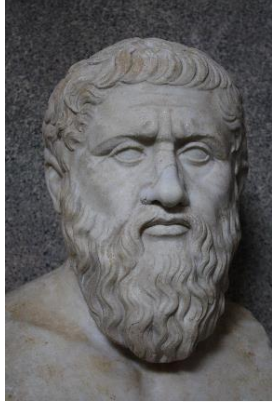
- Not to be confused with Statistics, which is a field of study

A **parameter** is a number that describes some aspect of a ***population***

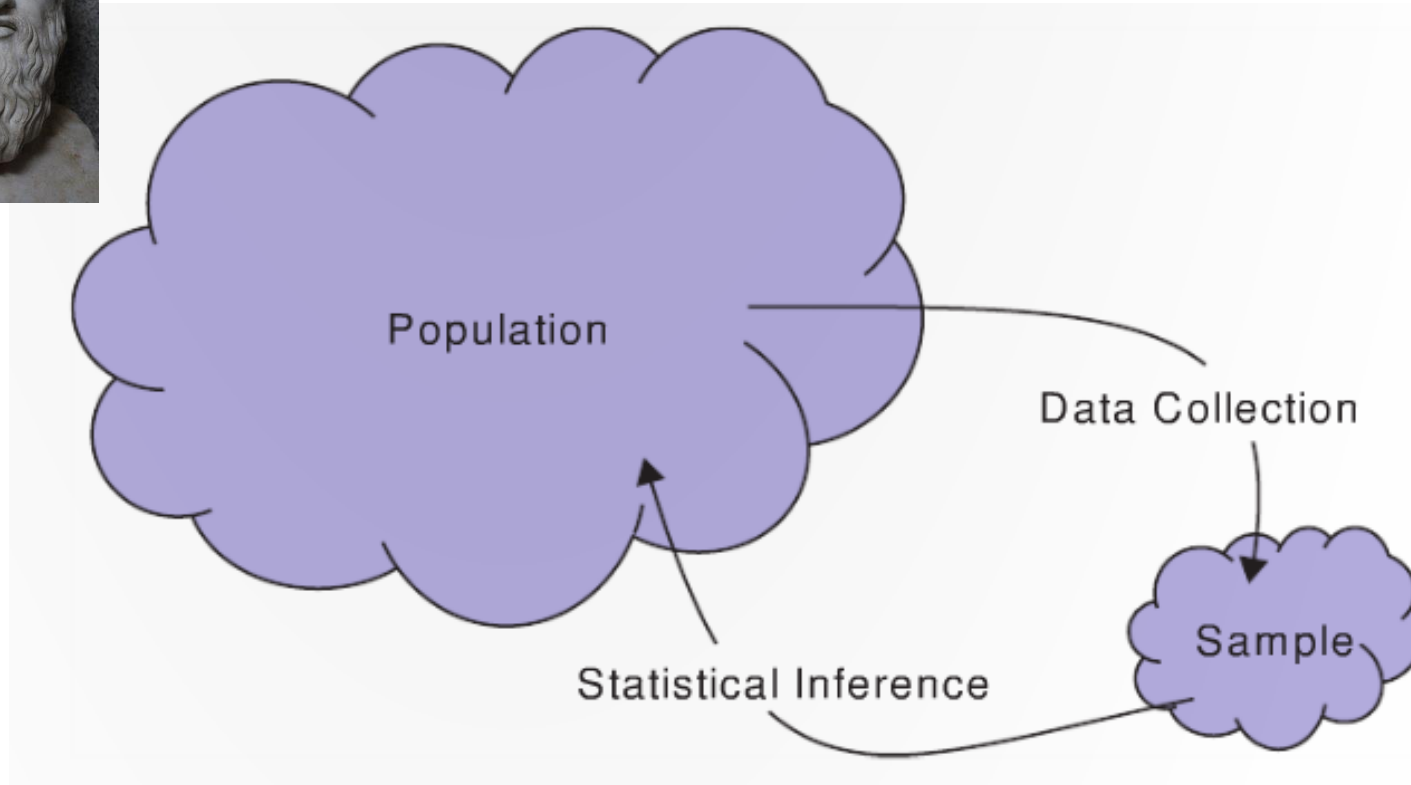


?

# Parameters and statistics



Parameters



statistics



# Proportions

For a *single **categorical variable***, the main ***statistic*** of interest is the *proportion* in each category

- E.g., the proportion of red sprinkles

$$\text{Proportion in a category} = \frac{\text{number in that category}}{\text{total number}}$$

# Example proportion of red sprinkles

The sample

- orange, red, green, white, white, white, ..., pink

The proportion for a **sample** is denoted  $\hat{p}$  (pronounced “p-hat”)

- $\hat{p}_{\text{red}} = 13/100 = 0.13$

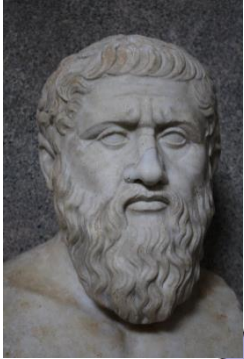
The proportion for a **population** is denoted  $\pi$  (the book uses  $p$ )

- $\pi_{\text{red}}$  proportion if we had measured all sprinkles in the population

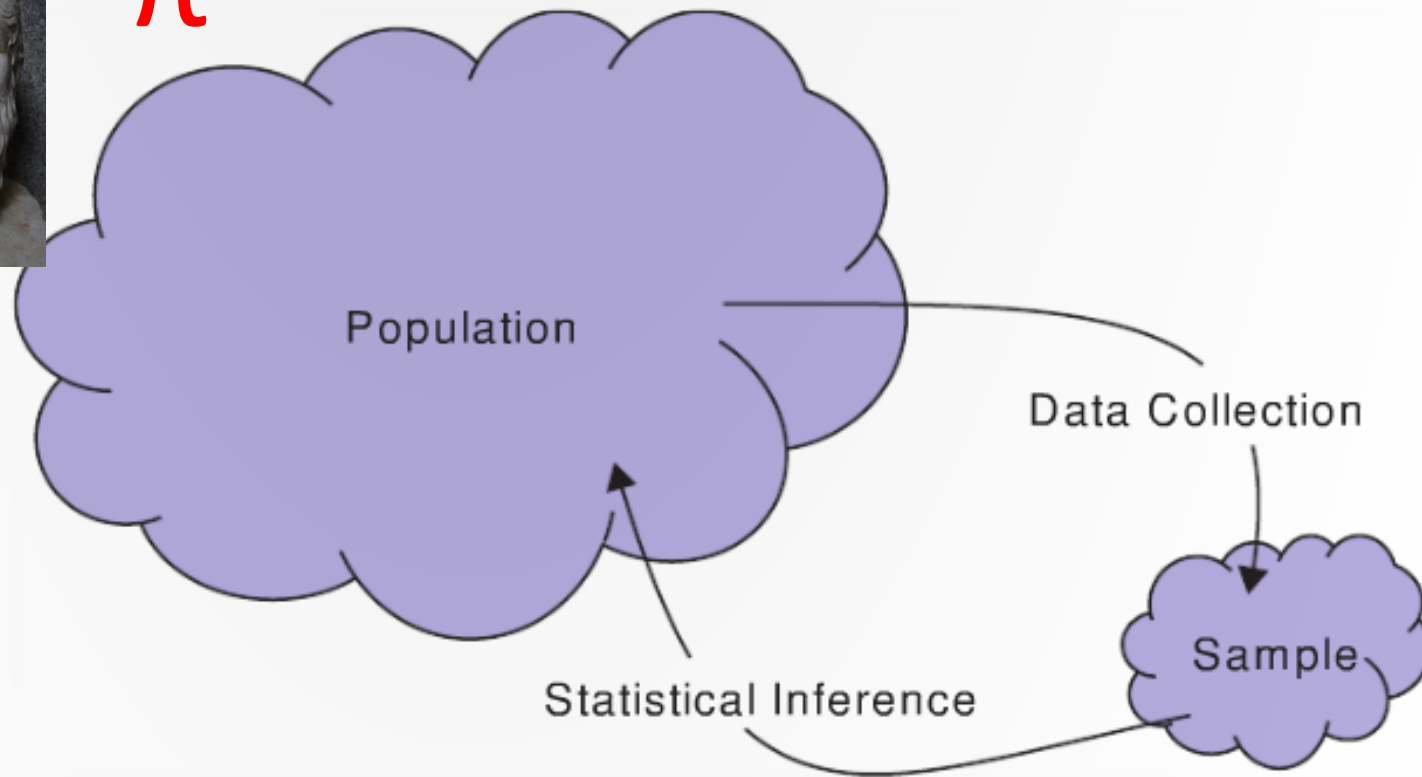
$\hat{p}$  is a **point estimate** of  $\pi$

- i.e.,  $\hat{p}$  our best guess of what  $\pi$  is

# Sample vs. Population proportion



$\pi$



Different samples yield different values for the statistic

$$\hat{p}_{s1\text{-red}} = 0.13$$

$$\hat{p}_{s2\text{-red}} = 0.11$$

$$\hat{p}_{s3\text{-red}} = 0.15$$

$\hat{p}$





# Calculating counts on a categorical variable

The count of how many items are in each category can be summarized in a ***frequency table***

| Color | green | orange | pink | red | white | yellow |  | Total |
|-------|-------|--------|------|-----|-------|--------|--|-------|
| Count | 20    | 11     | 9    | 13  | 36    | 11     |  | 100   |

In R: `my_table <- table(v)`

# Calculating proportions (relative frequencies)

We can convert a frequency table into a ***relative frequency table*** by dividing each cell by the total number of items

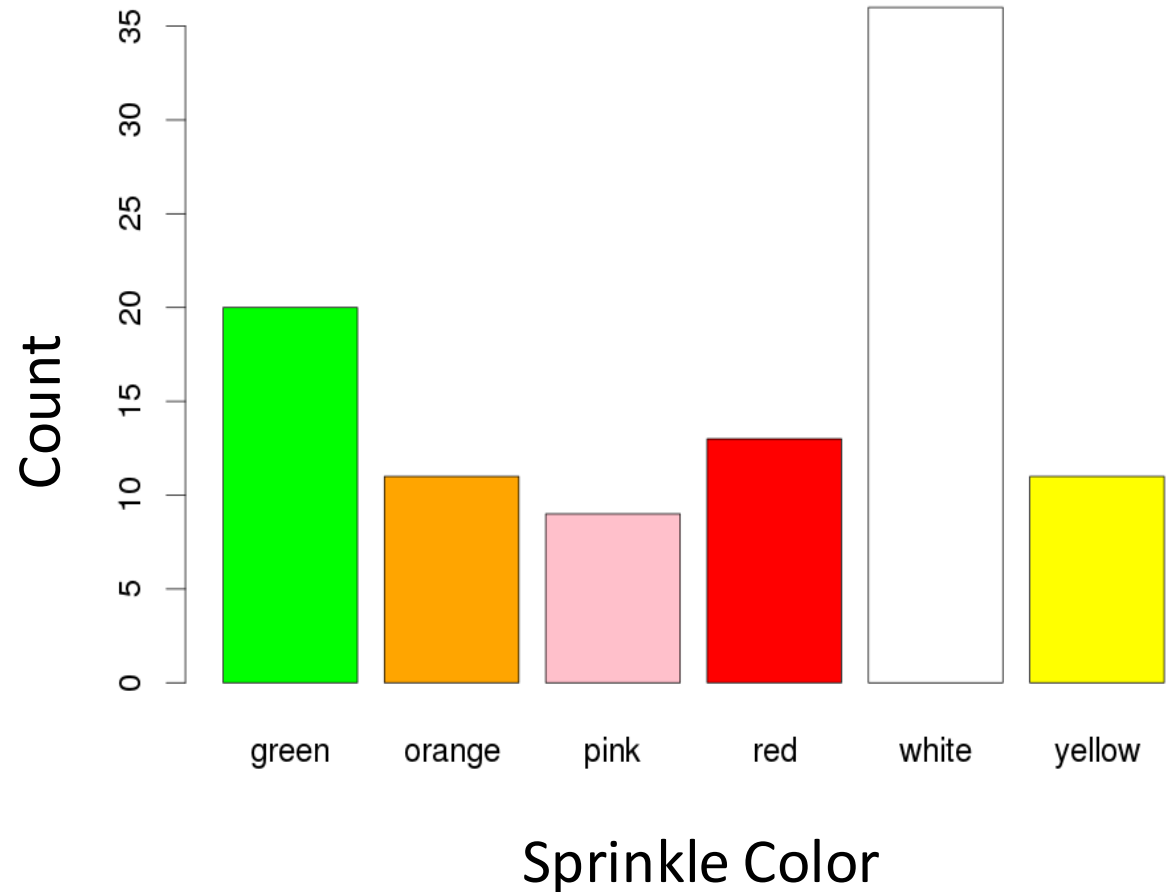
| Color | green | orange | pink | red | white | yellow |  | Total |
|-------|-------|--------|------|-----|-------|--------|--|-------|
| Count | .20   | .11    | .09  | .13 | .36   | .11    |  | 1     |

In R: `prop.table(my_table)`

# Visualizing categorical data: The bar plot

A bar plot shows the number of items in each category

The height of each bar corresponds to the number of items in a given category



In R: `barplot(my_table)`

# Visualizing categorical data: The pie chart

A pie chart plots the proportion of items in each category

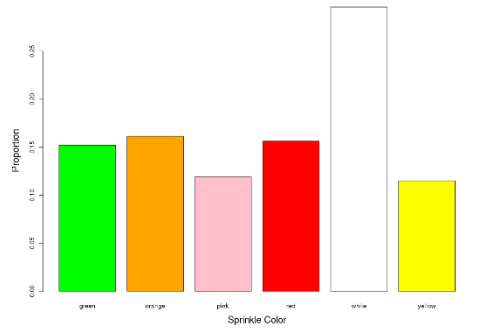
The area of each segment corresponds to the proportion of items in that segment

In R: `pie(my_table)`

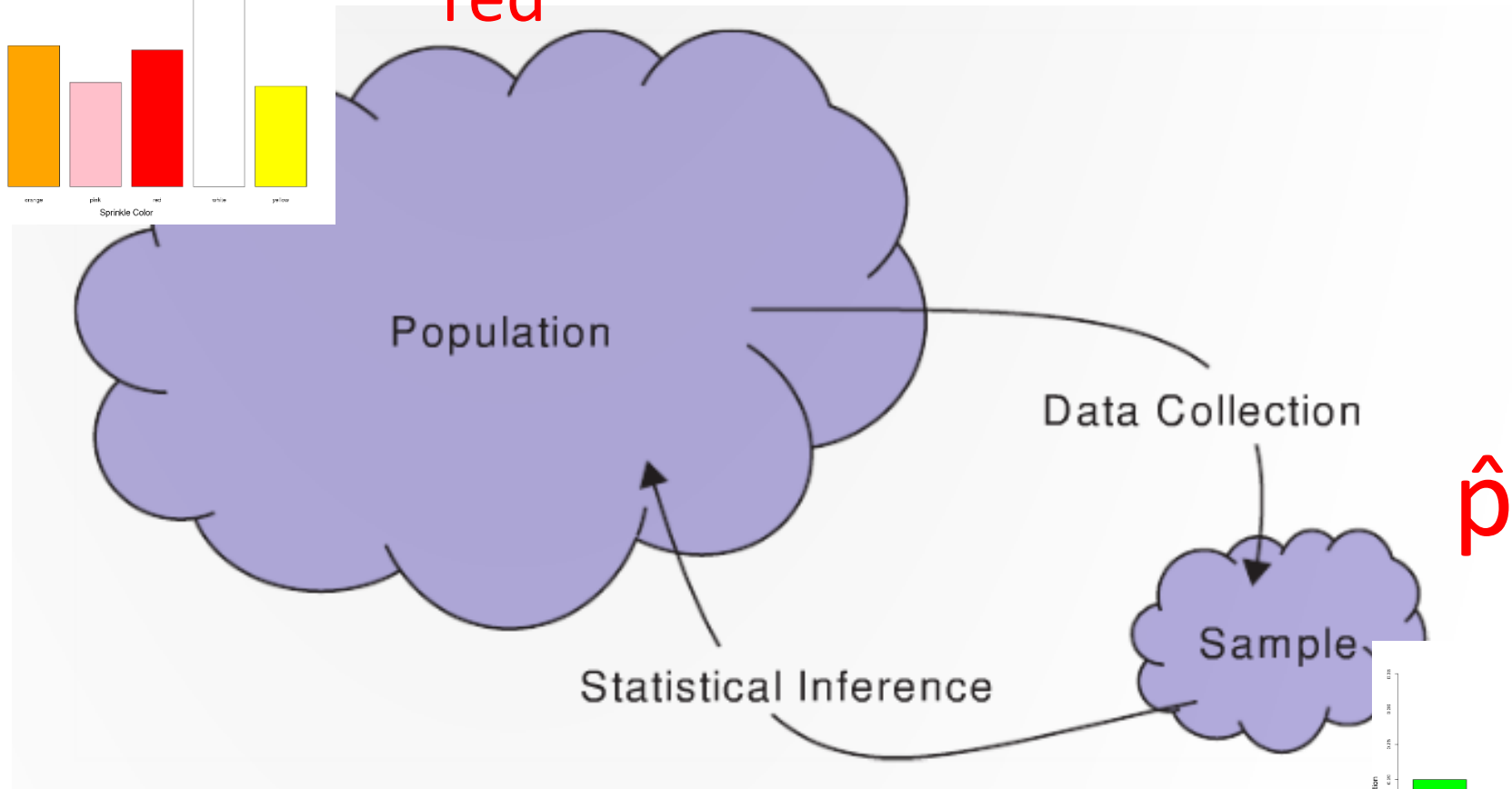


# Summary: Sample and Population proportion

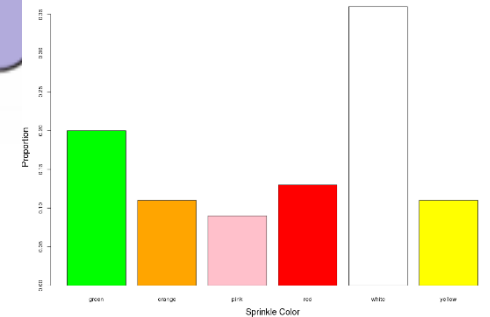
Categorical  
distribution



$\pi_{\text{red}}$



Bar chart



Let's sample virtual sprinkles in R...



# Sampling virtual sprinkles

```
library(SDS100)
```

```
sprinkle_sample <- get_sprinkle_sample(100)
```

```
sprinkle_count_table <- table(sprinkle_sample)
```

```
sprinkle_prop_table <- prop.table(sprinkle_count_table)
```

```
prop_red <- get_proportion(sprinkle_sample, "red")    # SDS1000 function
```

```
barplot(sprinkle_count_table)
```

```
pie(sprinkle_count_table)
```

# Summary of concepts

1. A **statistic** is a number that is computed from ***data in a sample***
  - The number of items in a sample is called the ***sample size*** and is usually denoted with the symbol  $n$
2. A **parameter** is a number that describes some aspect of a ***population***
3. A **point estimate** is using a value of a statistic as a guess for the value of a parameter
4. **When calculating proportions:**
  - The proportion statistic is denoted  $\hat{p}$
  - The population proportion is denoted  $\pi$
  - Thus  $\hat{p}$  is a ***point estimate*** of  $\pi$
5. Proportions can be summarized in a **relative frequency table** and can be visualized using **bar plots** and **pie charts**



# Summary of R

# a vector of character strings (or factors)

```
my_sample <- c("orange", "red", "green", "white", " white", ... )
```

# creating a table using the table() function

```
my_table <- table(my_sample)
```

# creating a frequency table using the prop.table() function

```
prop.table(my_table)
```

# get proportion in a category (SDS1000 function)

```
get_proportion(my_table, "category_name")
```

# creating bar and pie charts

```
barplot(my_table)
```

```
pie(my_table)
```