

From Data To Discovery

We will start momentarily. In the meantime:

Class Starter (respond on pollev.com/tiasondjaja)

Suppose that a data frame called `results` stores the results of a 5k (3.1 mile) race as well as the gender and age group of the participating runners:

Gender	AgeGroup	RaceTime
Female	40-59	28
Male	20-39	35
Male	40-59	19
Other	20-39	25
⋮	⋮	⋮
⋮	⋮	⋮

Suppose that you would like to create a new column called `Pace`, which consists of the number of minutes it took each runner to run one mile during this race.

Which R code would accomplish this task?

- A. `results$Pace <- 3.1 / RaceTime`
- B. `Pace <- results$RaceTime / 3.1`
- C. `results$Pace <- results$RaceTime / 3.1`
- D. None of the above

Open the Lesson 04 Jupyter Notebook

- ▶ Our NYU Classes site > Lessons > 4. Exploring Data (part 2); Click link
- ▶ No Jupyter Notebook for Lesson 5



Update your name to
First name + Last Init.
(+ pronouns, optional)



Mute your microphone
Turn your camera on



Make sure you can
raise your hand and
give nonverbal responses



Open the
chat
window

Quick Concept Check

(respond on pollev.com/tiasondjaja)

Suppose that a data frame called `results` store the results of a 5k (3.1 mile) race as well as the gender and age group of the participating runners:

Gender	AgeGroup	RaceTime (minutes)
Female	40-59	28
Male	20-39	35
Male	40-59	19
Other	20-39	25
.	.	.
.	.	.

Suppose that you would like to find the average race time among runners of each gender.

Which of the following R commands would give you the answer?

- A. `mean(filter(results, Gender == 'Female'))`
- B. `groupedresults <- filter(results, Gender == 'Female')`
`mean(groupedresults$RaceTime)`
- C. `groupedresults <- group_by(results)`
`summarize(groupedresults, aveRaceTime = mean(RaceTime))`
- D. None of the above works

Today's Plan

- ▶ Lesson02: Introduction to Jupyter Notebook and R
- ▶ Lesson03: Exploring Data
- ▶ Finish Lesson04: Exploring Data
- ▶ Start Lesson05: Visualizing data

Lesson 5: Visualizing Data

Goals and Key Ideas

1. Types of data visualizations

- ▶ How to determine the right type of data visualization for a given type of variable.

2. Good and bad data visualizations

- ▶ Bad data visualizations could mislead

Why visualize data?

Why visualize data?

- ▶ Help create a **visual summary** of data
- ▶ Help **identify patterns** in data
- ▶ Help **identify relationships** between variables

Why visualize data?

- ▶ Help create a **visual summary** of data
- ▶ Help **identify patterns** in data
- ▶ Help **identify relationships** between variables
- ▶ Help **communicate** or **describe** results of data analysis

Why visualize data?

- ▶ Help create a **visual summary** of data
- ▶ Help **identify patterns** in data
- ▶ Help **identify relationships** between variables
- ▶ Help **communicate** or **describe** results of data analysis

Bad data visualizations can **miscommunicate** or **misrepresent** information.

Breakout Activity: Types of Data Visualizations

Respond here: <https://pollev.com/tiasondjaja>

Types of Data Visualizations

Types of Data Visualizations

1. Bar Graphs

Types of Data Visualizations

1. Bar Graphs

Two types:

- ▶ (the most important one) Describe the **distribution** of a **categorical** variable
- ▶ Describe the relationship between a **categorical** variable and a **numerical** variable

Types of Data Visualizations

1. Bar Graphs

Two types:

- ▶ (the most important one) Describe the **distribution** of a **categorical** variable
- ▶ Describe the relationship between a **categorical** variable and a **numerical** variable

2. Histograms

Types of Data Visualizations

1. Bar Graphs

Two types:

- ▶ (the most important one) Describe the **distribution** of a **categorical** variable
- ▶ Describe the relationship between a **categorical** variable and a **numerical** variable

2. Histograms

Describe the **distribution** of a **numerical** variable.

Types of Data Visualizations

1. Bar Graphs

Two types:

- ▶ (the most important one) Describe the **distribution** of a **categorical** variable
- ▶ Describe the relationship between a **categorical** variable and a **numerical** variable

2. Histograms

Describe the **distribution** of a **numerical** variable.

3. Scatterplots

Types of Data Visualizations

1. Bar Graphs

Two types:

- ▶ (the most important one) Describe the **distribution** of a **categorical** variable
- ▶ Describe the relationship between a **categorical** variable and a **numerical** variable

2. Histograms

Describe the **distribution** of a **numerical** variable.

3. Scatterplots

Describe the relationship between a **numerical** variable and another **numerical** variable.

Types of Data Visualizations

1. Bar Graphs

Two types:

- ▶ (the most important one) Describe the **distribution** of a **categorical** variable
- ▶ Describe the relationship between a **categorical** variable and a **numerical** variable

2. Histograms

Describe the **distribution** of a **numerical** variable.

3. Scatterplots

Describe the relationship between a **numerical** variable and another **numerical** variable.

4. Time Series

Types of Data Visualizations

1. Bar Graphs

Two types:

- ▶ (the most important one) Describe the **distribution** of a **categorical** variable
- ▶ Describe the relationship between a **categorical** variable and a **numerical** variable

2. Histograms

Describe the **distribution** of a **numerical** variable.

3. Scatterplots

Describe the relationship between a **numerical** variable and another **numerical** variable.

4. Time Series

Describes how quantities change over time.

Types of Data Visualizations

1. Bar Graphs

Two types:

- ▶ (the most important one) Describe the **distribution** of a **categorical** variable
- ▶ Describe the relationship between a **categorical** variable and a **numerical** variable

2. Histograms

Describe the **distribution** of a **numerical** variable.

3. Scatterplots

Describe the relationship between a **numerical** variable and another **numerical** variable.

4. Time Series

Describes how quantities change over time.

5. Pie Charts

Types of Data Visualizations

1. Bar Graphs

Two types:

- ▶ (the most important one) Describe the **distribution** of a **categorical** variable
- ▶ Describe the relationship between a **categorical** variable and a **numerical** variable

2. Histograms

Describe the **distribution** of a **numerical** variable.

3. Scatterplots

Describe the relationship between a **numerical** variable and another **numerical** variable.

4. Time Series

Describes how quantities change over time.

5. Pie Charts

Describe the proportion of observations that belong to each category.

Types of Data Visualizations

1. Bar Graphs

Two types:

- ▶ (the most important one) Describe the **distribution** of a **categorical** variable
- ▶ Describe the relationship between a **categorical** variable and a **numerical** variable

2. Histograms

Describe the **distribution** of a **numerical** variable.

3. Scatterplots

Describe the relationship between a **numerical** variable and another **numerical** variable.

4. Time Series

Describes how quantities change over time.

5. Pie Charts

Describe the proportion of observations that belong to each category.

6. (and others!)

Examples of (good and bad) data visualizations

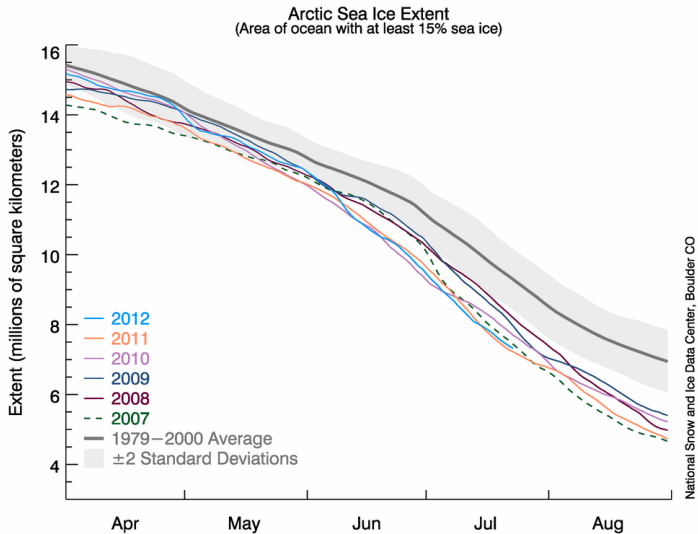
Handout

https://drive.google.com/file/d/1KLjLuK_5IzW2FyYR7Z10derewmsXL-LU/view?usp=sharing

Question & Task: Which of these examples are good data visualizations? Which are bad? Discuss!

Breakout Activity

<https://docs.google.com/document/d/10nDovYYmeyFXPAgySrBQQgYpbnoWosXDr1Dy7Q8sMCs/edit?usp=sharing>



23 Jul 2012

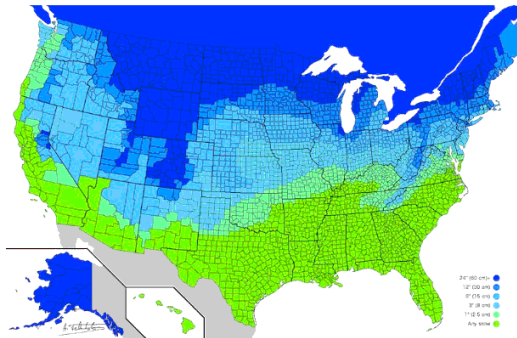
How Much Snow Before America Cancels School?

A map shows how many inches it takes before various regions call it off.

by **ERIC RANDALL** • 2/3/2014, 11:10 a.m.



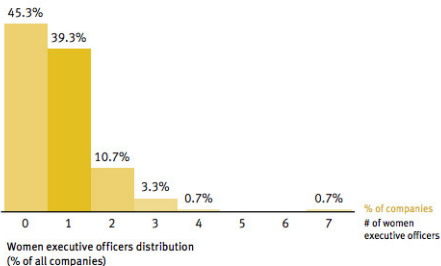
Get a compelling long read and must-have lifestyle tips in your inbox every Sunday morning — great with coffee!



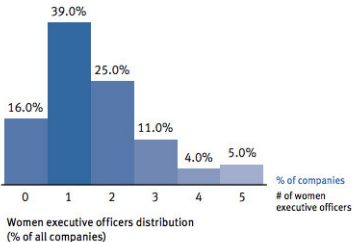
WOMEN EXECUTIVE OFFICERS DISTRIBUTION — 2013 PROXY SEASON

SV 150
2013

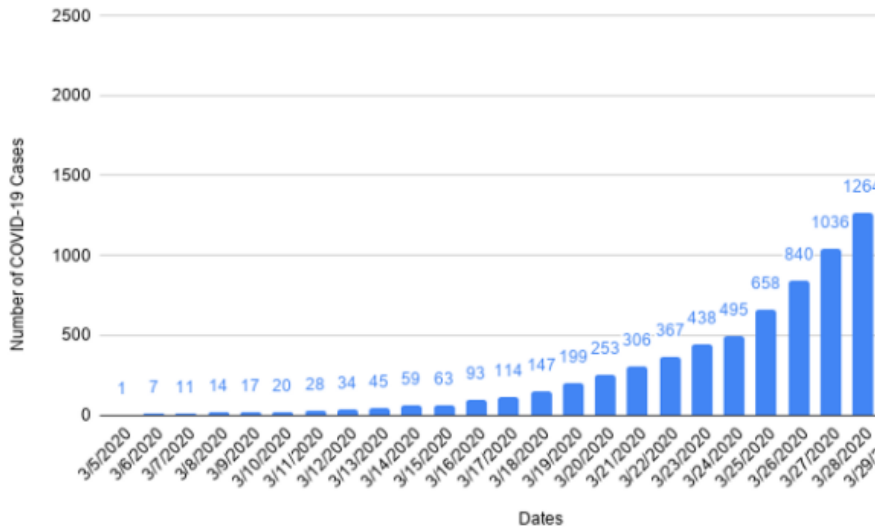
% of companies
with at least
1 woman
executive officer

S&P 100
2013

% of companies
with at least
1 woman
executive officer



Number of COVID-19 Cases in Russia from March 5 to March 31

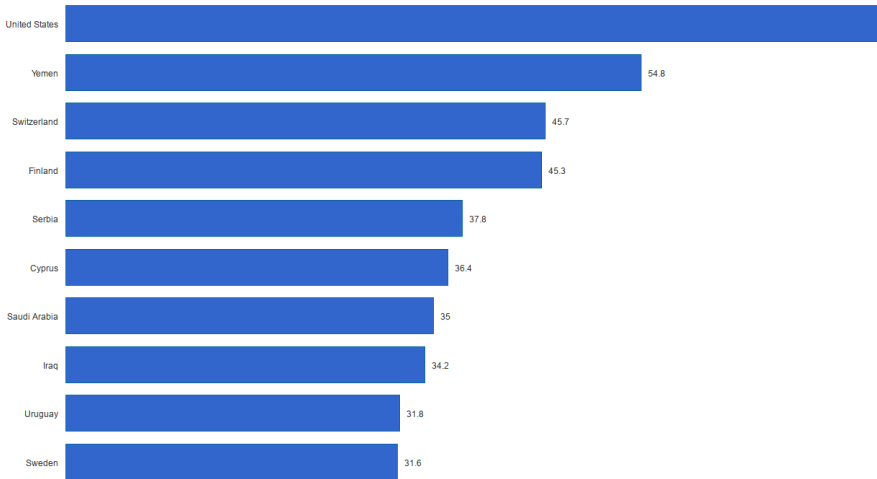


Gun ownership by country

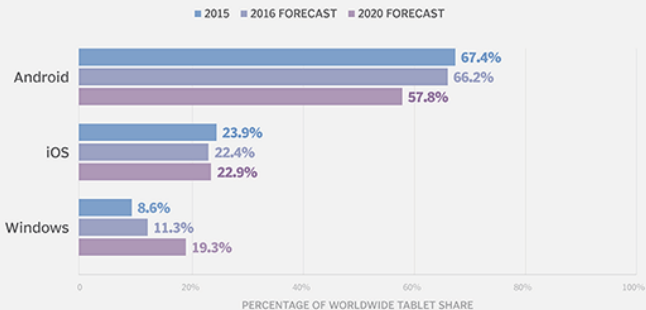
Click to switch view

Average firearms per 100 people

Homicide by firearm rate per 100,000 pop



IDC WORLDWIDE TABLET SHARE FORECAST BY OS 2015-2020

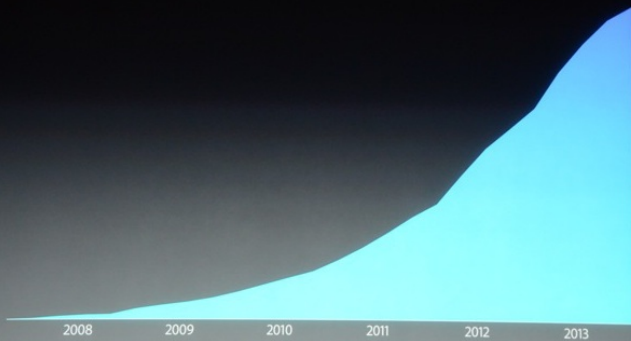


SOURCE: IDC WORLDWIDE QUARTERLY TABLET TRACKER, AUGUST 30, 2016

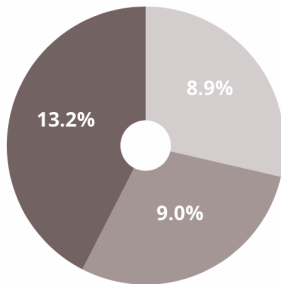
©2016 TECHTARGET. ALL RIGHTS RESERVED



Cumulative iPhone sales



PRETERM BIRTH BY RACE & ETHNICITY



■ Non-Hispanic White

■ Hispanic

■ Non-Hispanic Black

Top 20 Tourist Generating Countries To UNITED STATES OF AMERICA

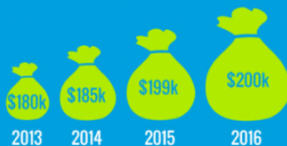


SALES ARE STRONG & PRICES ARE GOING UP

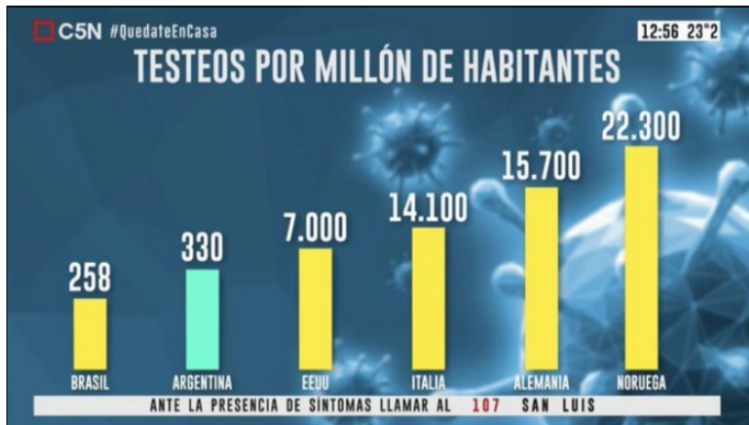
Number of Sold Businesses
Are Going Up



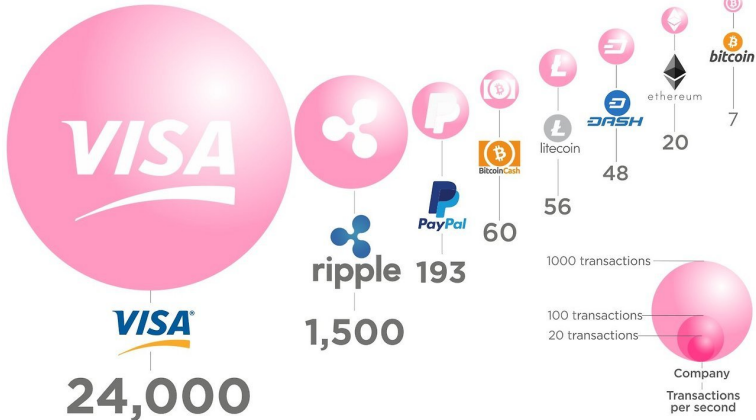
Average Sale Price of Businesses
Are Going Up



SOURCE: RealPage, Fourth Quarter 2016 Insight Report

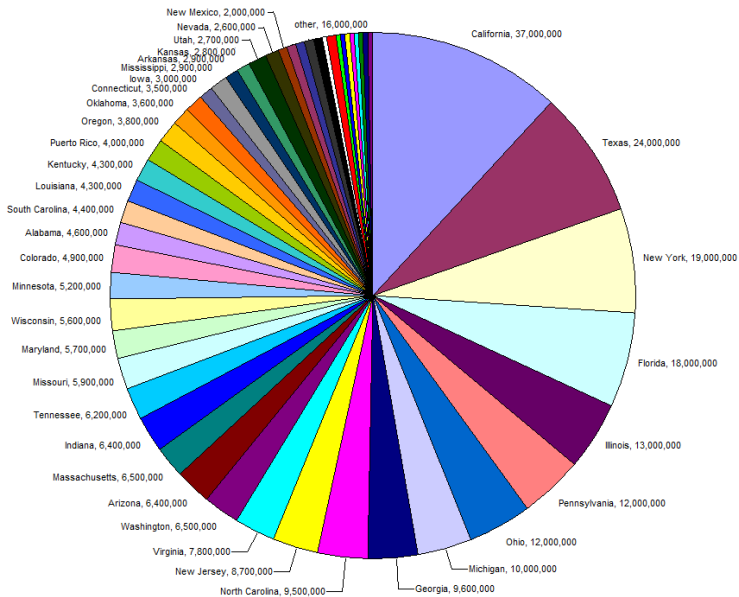


Cryptocurrencies Transaction Speeds Compared to Visa & Paypal

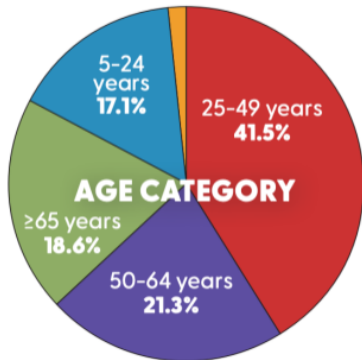


Article & Sources:
<https://howmuch.net/articles/crypto-transaction-speeds-compared>
<https://howmuch.net/sources/crypto-transaction-speeds-compared>

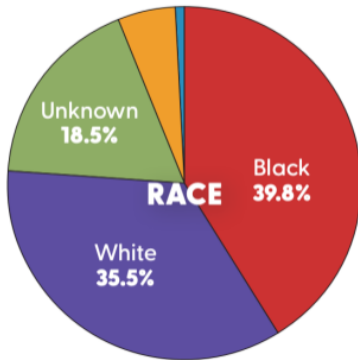
howmuch.net



DEMOGRAPHIC CHARACTERISTICS

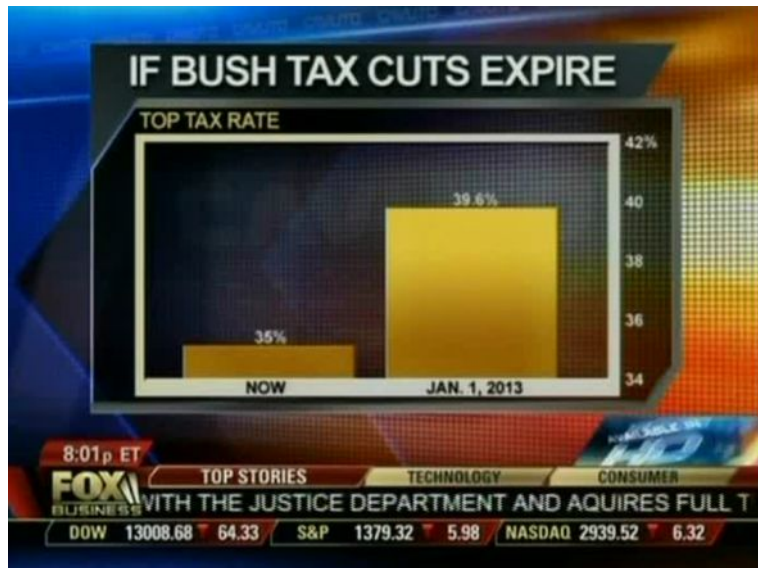


0-4 Years **1.6%**

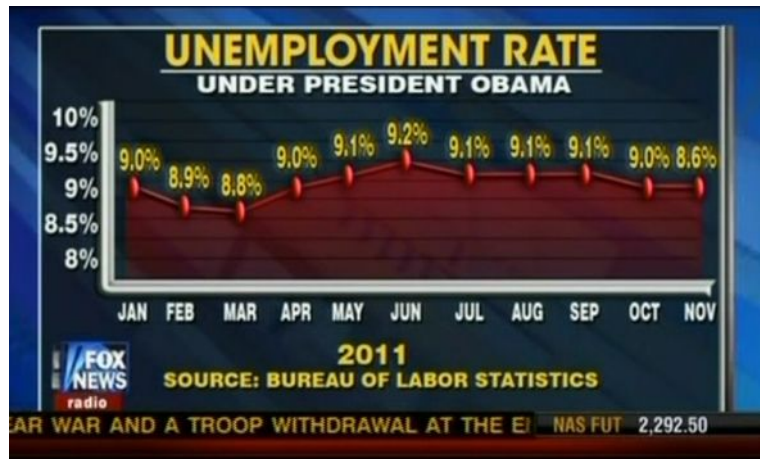


Other **5.8%** Asian **0.4%**

SCREENSHOT/ALABAMA DEPARTMENT OF HEALTH



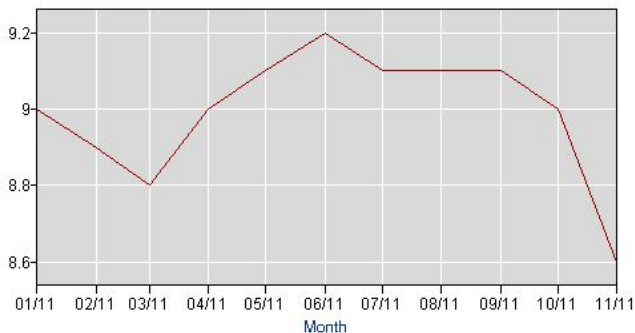




Data extracted on: December 12, 2011 (9:50:59 AM)

Labor Force Statistics from the Current Population Survey

Series Id: LNS14000000
Seasonally Adjusted
Series title: (Seas) Unemployment Rate
Labor force status: Unemployment rate
Type of data: Percent or rate
Age: 16 years and over



Partial list of image sources:

- ▶ [https://qz.com/1872980/
how-bad-covid-19-data-visualizations-mislead-the-public/](https://qz.com/1872980/how-bad-covid-19-data-visualizations-mislead-the-public/)
- ▶ [https://towardsdatascience.com/
stopping-covid-19-with-misleading-graphs-6812a61a57c9](https://towardsdatascience.com/stopping-covid-19-with-misleading-graphs-6812a61a57c9)

Some principles for good data visualizations

- ▶ Use the appropriate type of graphs/charts (consistent with type of data)

Some principles for good data visualizations

- ▶ Use the appropriate type of graphs/charts (consistent with type of data)
- ▶ Clearly label and explain the axes

Some principles for good data visualizations

- ▶ Use the appropriate type of graphs/charts (consistent with type of data)
- ▶ Clearly label and explain the axes
- ▶ Axes should generally start from 0 and consistent in scale

Some principles for good data visualizations

- ▶ Use the appropriate type of graphs/charts (consistent with type of data)
- ▶ Clearly label and explain the axes
- ▶ Axes should generally start from 0 and consistent in scale
- ▶ Avoid pie charts or other shapes; use bar graphs

Some principles for good data visualizations

- ▶ Use the appropriate type of graphs/charts (consistent with type of data)
- ▶ Clearly label and explain the axes
- ▶ Axes should generally start from 0 and consistent in scale
- ▶ Avoid pie charts or other shapes; use bar graphs
- ▶ In bar charts and histograms, the area of each bar should be proportional to the quantity represented.

Some principles for good data visualizations

- ▶ Use the appropriate type of graphs/charts (consistent with type of data)
- ▶ Clearly label and explain the axes
- ▶ Axes should generally start from 0 and consistent in scale
- ▶ Avoid pie charts or other shapes; use bar graphs
- ▶ In bar charts and histograms, the area of each bar should be proportional to the quantity represented.
- ▶ When visualizing proportions or percentages, clearly state what the population is.