# Variation and Randomness

Understanding Political Numbers

Feb 20, 2019

## Essay 1

Overall great!

#### Feedback:

- Commit to a measure
- Direction of relationship
- Avoid "proving" and "disproving"
- Contamination of the

D۷

• Please proofread!

#### My bad:

• Latent probability vs incidence (e.g. of war)

## Essay 1

### R Exercise 1

Overall great!

Feedback:

- Commit to a measure
- Direction of relationship
- Avoid "proving" and "disproving"
- Contamination of the

D۷

• Please proofread!

My bad:

• Latent probability vs incidence (e.g. of war)

Installing vs. loading ("librarying") packages

Folder names

Data file name

# Let's talk about papers

#### Original research question

- Does *X* affect *Y*?
- What other things might affect Y (at least two other potential explainers)? Use theoretical reasoning

Collect data (at least 50 cases)

Write a 12-page paper (not counting graphics)

- What's the question
- Theory and hypotheses
- Explain your data and method

# Let's talk about papers

## What to do:

Original research question

- Does *X* affect *Y*?
- What other things might affect Y (at least two other potential explainers)? Use theoretical reasoning

Research question due March 11 (a Monday)

- Variables, hypotheses, and (proposed) data sources
- Before then: 15-minute meeting w/ Michael or me
- Extra office hours

Collect data (at least 50 cases)

Write a 12-page paper (not counting graphics)

- What's the question
- Theory and hypotheses
- Explain your data and method

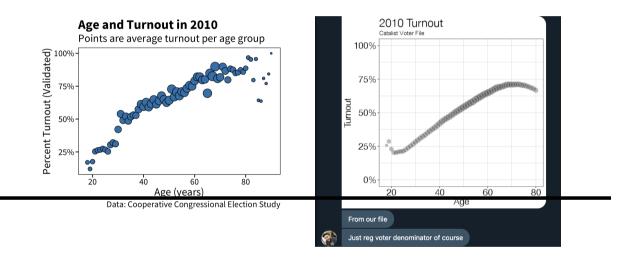
Dataset due after spring break

#### More info coming soon

- Watch for email about scheduling meetings!
   You must meet with us
- Formal assignment sheet soon!

## https://mikedecr.github.io/p 270/

## Will we ever use this in the real world?



## Where were we?

| The "real<br>world" | Your<br>data   |
|---------------------|----------------|
| "Population"        | Sample         |
| Theoretical mean    | Sample<br>mean |
| Expectation         | Estimate       |
| Parameter           | Statistic      |

 $\mu$  x

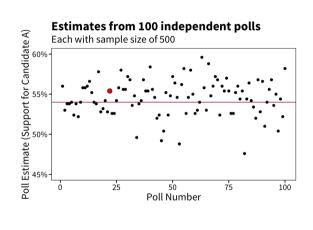
Remember:

Observed data = Truth

+ Bias + Noise

### Where were we?

| The "real<br>world" | Your<br>data   |
|---------------------|----------------|
| "Population"        | Sample         |
| Theoretical mean    | Sample<br>mean |
| Expectation         | Estimate       |
| Parameter           | Statistic      |



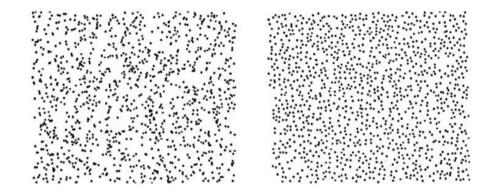
- In a given sample, I am randomly inaccurate, and that is to be expected.
  - Every time I call someone, there's a 54% chance they support A, so there's Remamberess
  - When randomness is at play, you get uncertainty

- Observed data = Truth

  It looks like it's getting wider, but it isn't.

  That's just random variation

## Which is random?



Flip a coin. Probability of heads is 0.5.

Flip a coin. Probability of heads is 0.5.

Flip a coin twice. Probability of two heads in a row is  $0.5 \times 0.5 = 0.25$ 

Flip a coin. Probability of heads is 0.5.

Flip a coin twice. Probability of two heads in a row is  $0.5 \times 0.5 = 0.25$ 

Flip a coin n times. Probability of n heads in a row is  $0.5^n$ .

Flip a coin. Probability of heads is 0.5.

Flip a coin twice. Probability of two heads in a row is  $0.5 \times 0.5 = 0.25$ 

Flip a coin n times. Probability of n heads in a row is  $0.5^n$ .

I flip a coin 1000 times. What's the longest sequences of heads I get?

Flip a coin. Probability of heads is 0.5.

Flip a coin twice. Probability of two heads in a row is  $0.5 \times 0.5 = 0.25$ 

Flip a coin n times. Probability of n heads in a row is  $0.5^n$ .

I flip a coin 1000 times. What's the longest sequences of heads I get?

```
# I run this simulation in R
max(one_thousand_flips$heads_in_a_row)
```

## [1] 10

Wow! The probability of getting that many heads in a row is  $0.5^{10} = 0.0009766!$ 

Flip a coin. Probability of heads is 0.5.

Flip a coin twice. Probability of two heads in a row is  $0.5 \times 0.5 = 0.25$ 

Flip a coin n times. Probability of n heads in a row is  $0.5^n$ .

I flip a coin 1000 times. What's the longest sequences of heads I get?

```
# I ran this simulation in R
max(one_thousand_flips$heads_in_a_row)
```

## [1] 10

Wow! The probability of getting that many heads in a row is  $0.5^{10} = 0.0009766!$ 

Unlikely things aren't always unusual.

When you have lots of data, unlikely things are common.

#### Who will win the presidency?

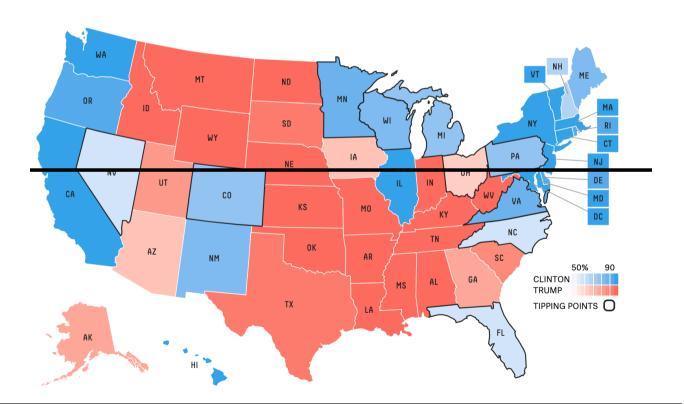


#### Chance of winning









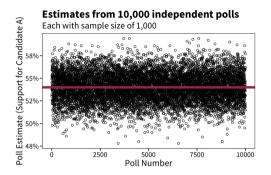
# Time for something amazing

## 10,000 Independent Polls

#### Why 10,000?

- reality is clumpy
- but it's clumpy when you look up close
- when you zoom out really far, clumpiness cancels out (randomly)

## 10,000 Independent Polls



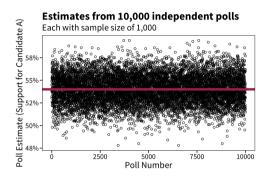
#### Why 10,000?

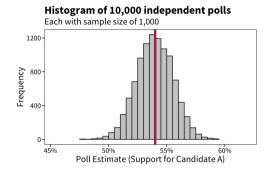
- reality is clumpy
- but it's clumpy when you look up close
- when you zoom out really far, clumpiness cancels out (randomly)

#### bee swarm

- if we were to zoom in, we would see that every poll is different
- Every poll is off a little bit
- but overall, they're close to the truth
- It looks like the farther you get from the truth, the fewer polls are out there

## 10,000 Independent Polls





#### Why 10,000?

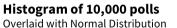
- reality is clumpy
- but it's clumpy when you look up close
- when you zoom out really far, clumpiness cancels out (randomly)

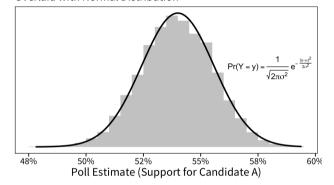
#### bee swarm

- if we were to zoom in, we would see that every poll is different
- Every poll is off a little bit
- but overall, they're close to the truth
- It looks like the farther you get from the truth, the fewer polls are out there

#### The normal distribution

If some variable Y is affected (at least in part) by an accumulation of random fluctuations, then *the distribution of Y* will be approx. normal

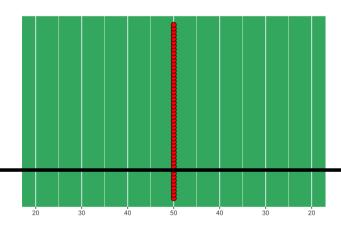


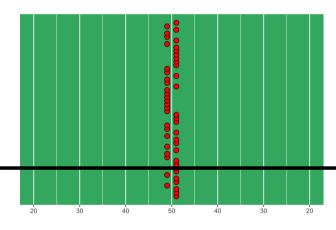


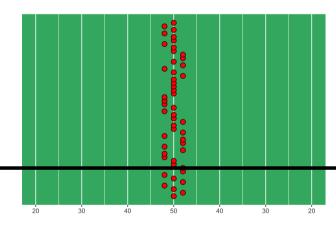
#### A "draw" from a

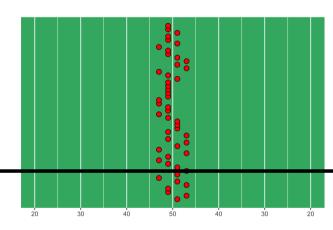
#### normal distribution:

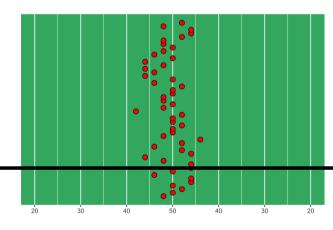
Any *individual*observation of *y* is just one number, but the probability of observing that value (relative to the mean) is given by the normal distribution

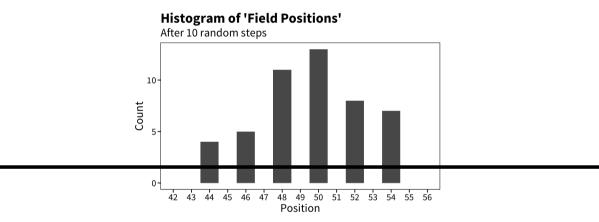












# What do "accumulating fluctuations" have to do with the mean

or, "Why are means normally distributed?"

$$\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

# What do "accumulating fluctuations" have to do with the mean

or, "Why are means normally distributed?"

$$\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

Summing the data is "adding up fluctuations." Cool, huh!

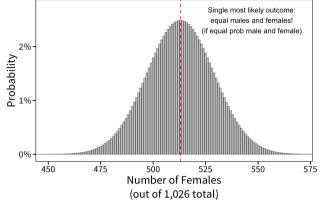
## "Order arising from disorder" or just *pure chaos*?

## Most likely babies



## Most likely babies





# How we use the normal distribution in statistics

## "The sampling distribution of $\bar{x}$ is normal"

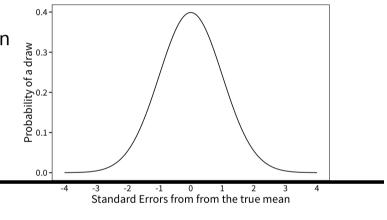
```
This means "the
                                   Data x where \bar{x} = 10
    sample mean is a
                                   ## # A tibble: 100 x 3
    normal draw from the
                                                  x \cdot x - xbar \cdot (x - xbar)^2
    true mean".
                                                             <dbl>
                                            <int>
                                                                                   <dbl>
                                   ##
    How accurate is the
                                                                                        81
    mean?
                                   ##
                                                 10
                                                 14
• nowhaltdrawmendispegrsion single number, bugit's more likely go appear in the
                                                                                        81
   center of the diagraphic and less likely to appear in the tail _4
                                                                                        16
• We don't know where the true mean is but we know enough about the properties of the Normal distribution that we can build an informed guess 10
    deviation with datas at:
                                                                                        36
                                                                                        81
   s(x) = \sqrt{\frac{\sum_{i=0}^{n} (x_i - \bar{x})^2}{n-1}} ## # ... with 90 more rows
```

## "The sampling distribution of $\bar{x}$ is normal"

#### Standard error is

"average error between  $\bar{x}$  and the true mean  $\mu$ "

Standard error of  $\bar{x}$ :



#### Rules for averages:

 $se(\bar{x}) =$ 

 ~ 68% of estimates are within ± 1 std. error from "true" value

## Statistics is about dealing with uncertainty in real data

Here is how most people do statistics

Every sample is an *imperfect* representation of the Lots of statistics depend on assumptions about the "distribution" of our data underlying population

- a lot of these assumptions are like, assuming that our estimate of the mean is normal(Valistributed sample estimate (such as a mean) and its
  or assuming that our data are normally distributed
  Standard error

It often feels like this is a tenuous assumption, that could easily be broken Most of the time, your estimate is within ± 2 standard errors

it is actually hard to break, because it isn't fragile

true value"

The normal distribution is just what happens when you add lots of fluctuations

- together
- and lots of things are influenced by lots of small fluctuations