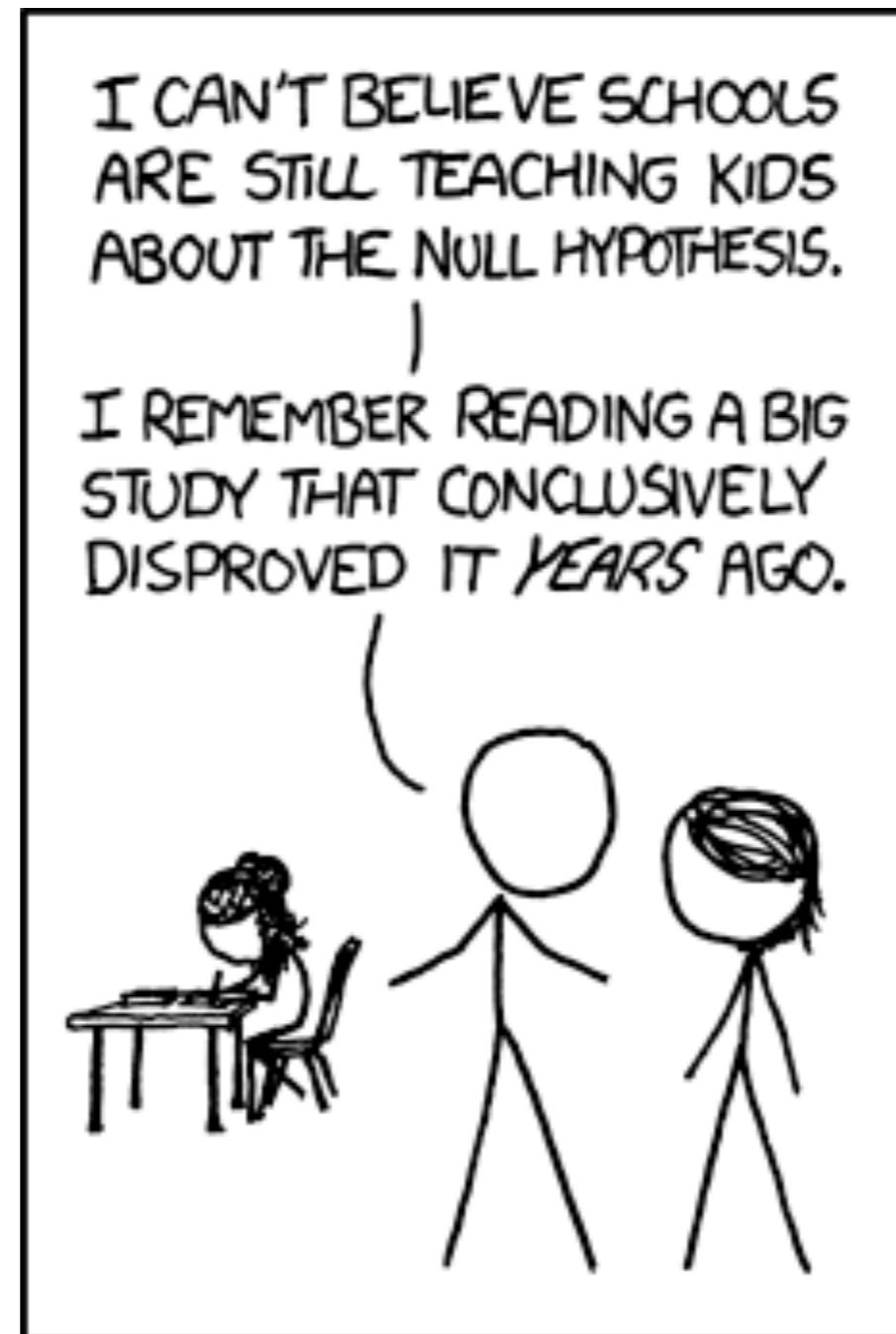


Modeling data



01/28/2019

Logistics

Homework 2

Homework 2

```
# 21, Jamie
# 22, Gale
# 23, Robbie
# 24, Tracy
# 25, Merrill
# 26, Noel
# 27, Dee
# 28, Sunny
# 29, Paris
# 30, Ariel
# 31, Rene
# 32, Johnnie
# 33, Jan
# 34, Layne
# 35, Devon
#
# If you can't quite figure out how to compute the most unisex names, then filter your data based on th
#####
# Per year - calculate proportions and unisex from counts; also filter to data to the relevant time
df.data.cleaned <- df.data %>%
  select(-prop) %>%
  spread(sex, n) %>%
  mutate(year_total = M+F,
    male = M/year_total,
    female = F/year_total,
    unisex = abs(female - 0.5)) %>%
  filter(year >= 1930 & year <= 2012)

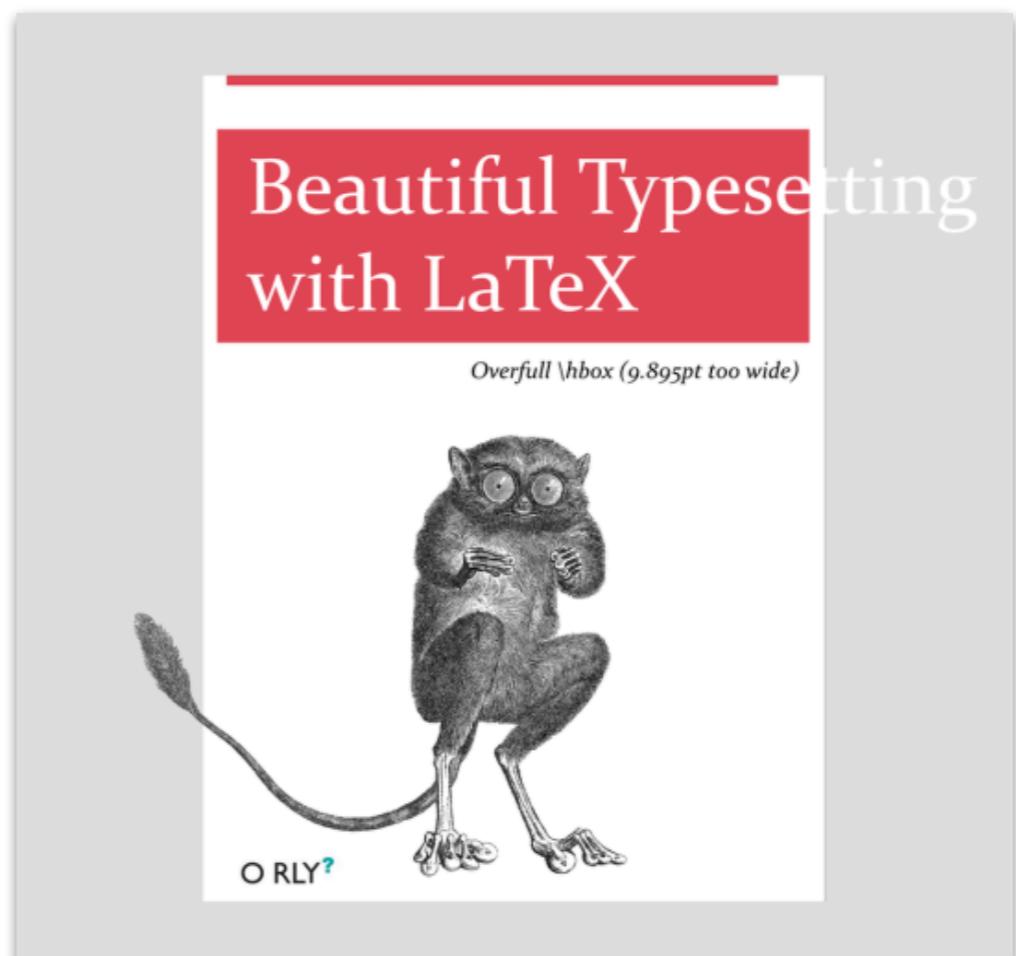
#####
# Per name across years - filter to names that were given at least 9000 times (overall) and occur at
df.data.35.names <- df.data.cleaned %>%
  group_by(name) %>%
  summarize(occurrences = n_distinct(year),
    overall_total = sum(year_total),
    mse = mean((female - 0.5)^2)) %>%
  filter(occurrences >=75 &
    overall_total >= 9000) %>%
  arrange(mse) %>%
  top_n(-35, mse)

#####
# Filter cleaned data to just data about the 35 names
df.data.35 <- df.data.cleaned %>%
  filter(name %in% df.data.35.names$name)

#####
# Per name - find the most unisex year and value
df.data.35.most <- df.data.35 %>%
  group_by(name) %>%
  top_n(-1, unisex)

#####
# Tidy cleaned data for visualization
df.data.35 <- df.data.35 %>%
  select(-F, -M, year_total) %>%
  gather(gender, prop, male:female)
df.data.35.most <- df.data.35.most %>%
```

6



margin column

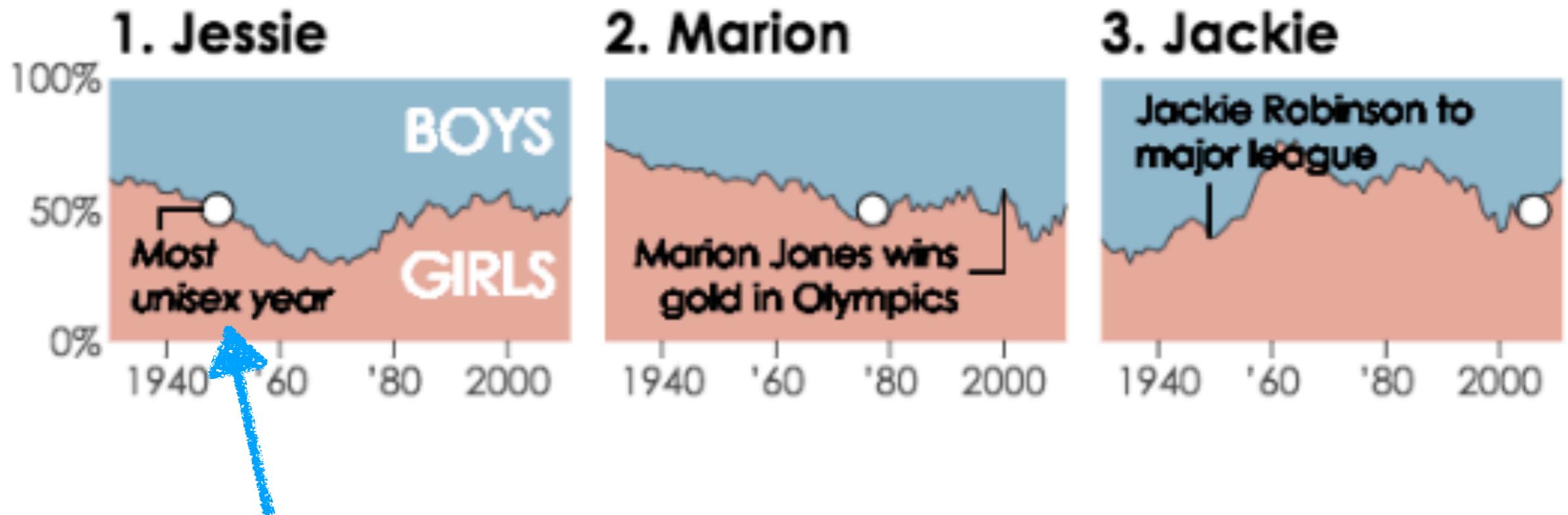
```
visualization2.Rmd x
1 --- 
2   title: "Class 3"
3   author: "Tobias Gerstenberg"
4   date: "January 11th, 2019"
5   output:
6     bookdown::html_document2:
7       toc: true
8       toc_depth: 4
9       theme: cosmo
10      highlight: tango
11 
12 
13 {r setup, include=FALSE}
14 # these options here change the formatting of how comments are rendered
15 knitr::opts_chunk$set(
16   collapse = TRUE,
17   comment = "#>")
18 ...
19 
20 # Visualization 2
21 
22 In this lecture, we will lift our `ggplot2` skills to the next level!
23 
24 ## Learning objectives
25 
26 - Deciding what plot is appropriate for what kind of data.
27 - Customizing plots: Take a sad plot and make it better.
28 - Saving plots.
29 - Making figure panels.
30 - Debugging.
31 - Making animations.
32 - Defining snippets.
33 

12:1 (Top Level) ▾
```

R Markdown ▾

- Visualization 2
- Learning objectives
- Install and load pack...
- Overview of different...
- Proportions
- Stacked bar charts
- Pie charts
- Comparisons
- Boxplots
- Violin plots
- Joy plots
- Practice plot 1
- Relationships
- Scatter plots
- Raster plots
- Temporal data
- Customizing plots
- Changing the order...
- Dealing with legends
- Choosing good colors
- Customizing themes
- Saving plots
- Creating figure panels
- Peeking behind the ...
- Making animations
- Shiny apps
- Defining snippets
- Additional resources
- Cheatsheets
- Data camp courses
- Books and chapters
- Misc
- Session info

Homework 2



if text looks pixelated, it's likely that there are many layers of text on top of each other

`geom_text()` needs a separate data frame with one entry per facet

Homework 2

I learned something new!

```
1 data = c(1, 3, 4, 2, 5)
2 prediction = c(1, 2, 2, 1, 4)
3
4 # calculate root mean squared error the pipe way
5 rmse = (prediction - data) ^ 2 %>%
6   mean() %>%
7   sqrt() %>%
8   print()
```



can we pipe this even more?

```
1 rmse = prediction %>%
2   subtract(data) %>%
3   raise_to_power(2) %>%
4   mean() %>%
5   sqrt() %>%
6   print()
```

Homework 2

I learned something new!



```
library("magrittr")
```

extract	`[`
extract2	`[[`
inset	`[<-`
inset2	`[[<-`
use_series	`\$`
add	`+`
subtract	`-`
multiply_by	`*`
raise_to_power	`^`
multiply_by_matrix	`%*%`
divide_by	`/`
divide_by_int	`%/`
mod	`%%`
is_in	`%in%`
and	`&`
or	` `
equals	`==`
is_greater_than	`>`
is_weakly_greater_than	`>=`
is_less_than	`<`
is_weakly_less_than	`<=`
not(`n'est pas`)	`!`

Your feedback

Your feedback

Central limit theorem was a little bit confusing.....

Good explanation. I haven't really understood the CLT before

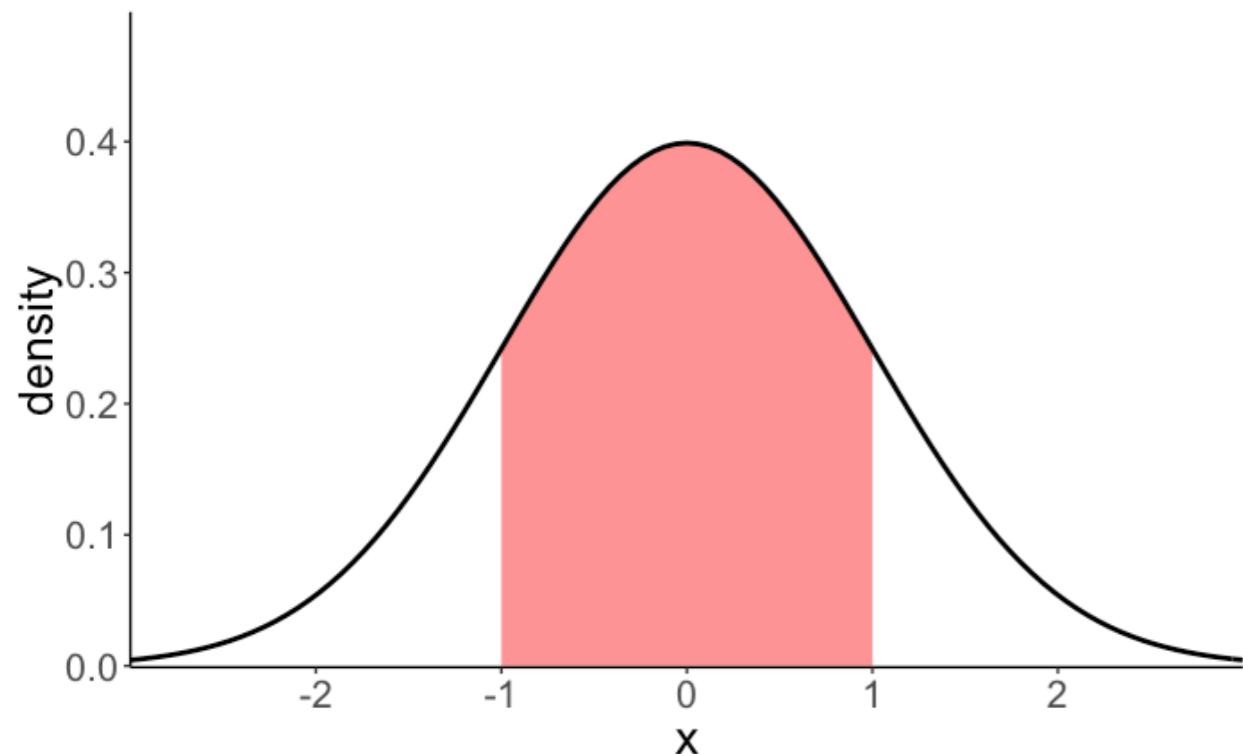
I think you spent too long on CLT which isn't intuitively difficult and would like more time on the harder topics near the end

sometimes it's tricky
to get it right

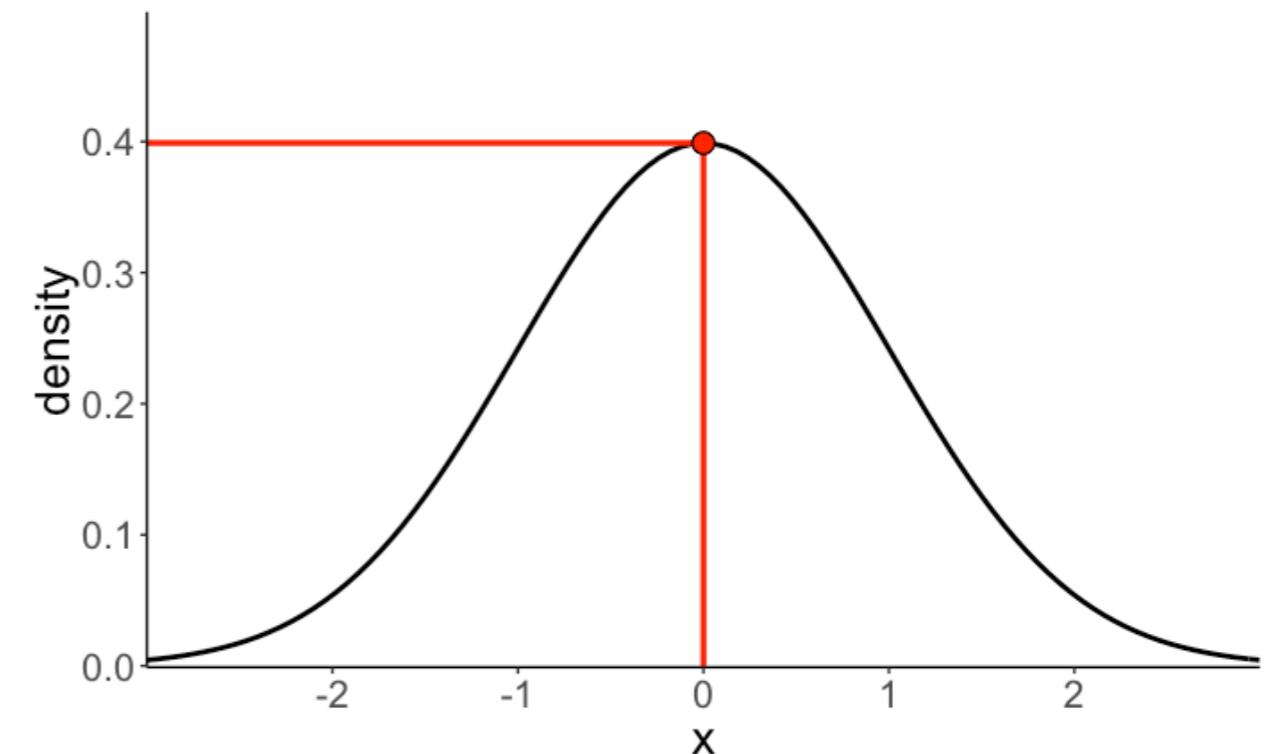
Probability vs. likelihood

Probability vs. likelihood

Probability

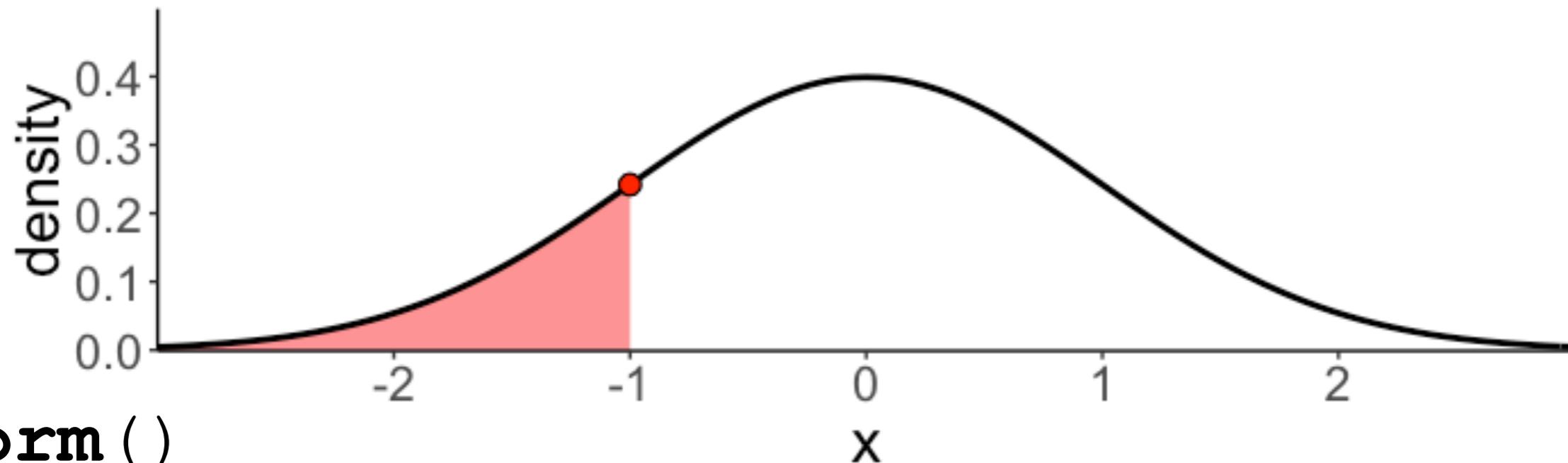


Likelihood

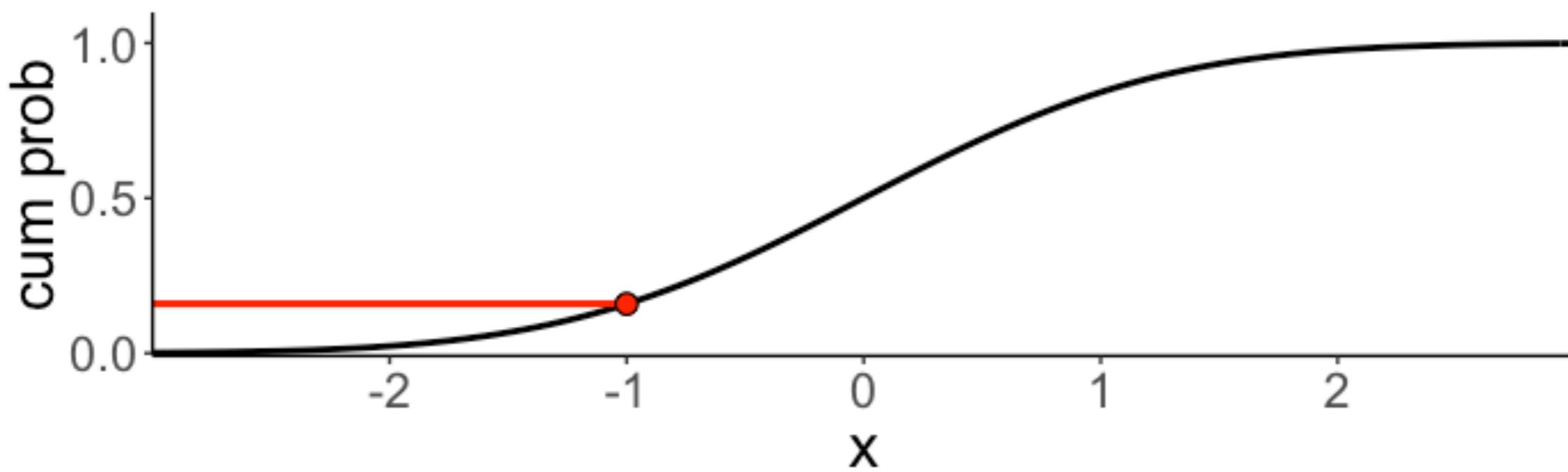


Probability vs. likelihood

dnorm ()

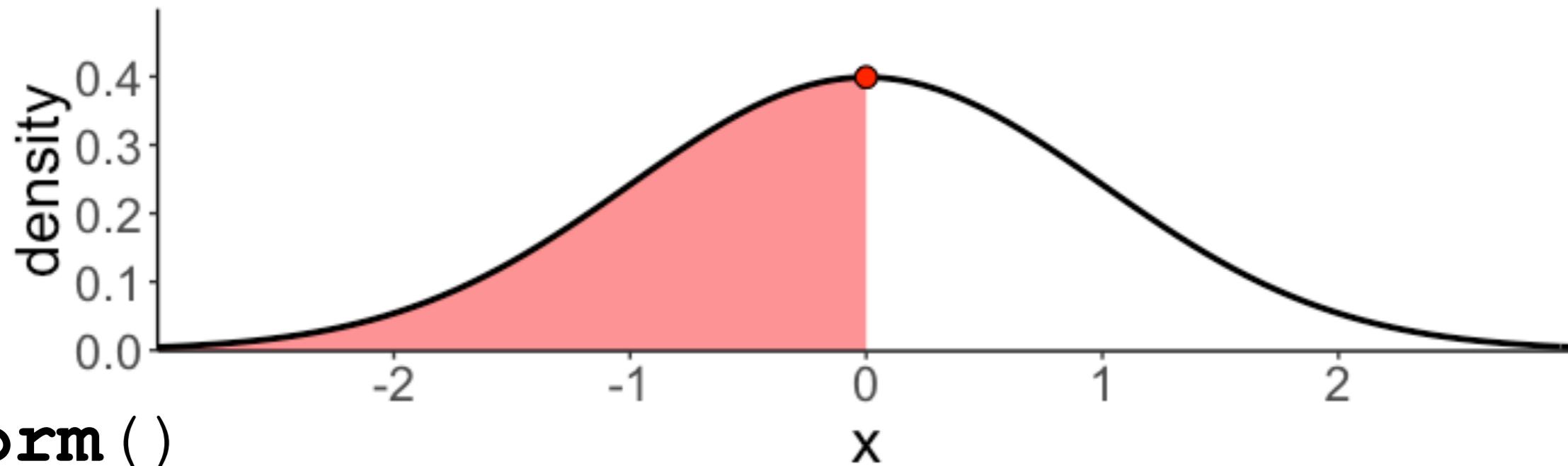


pnorm ()

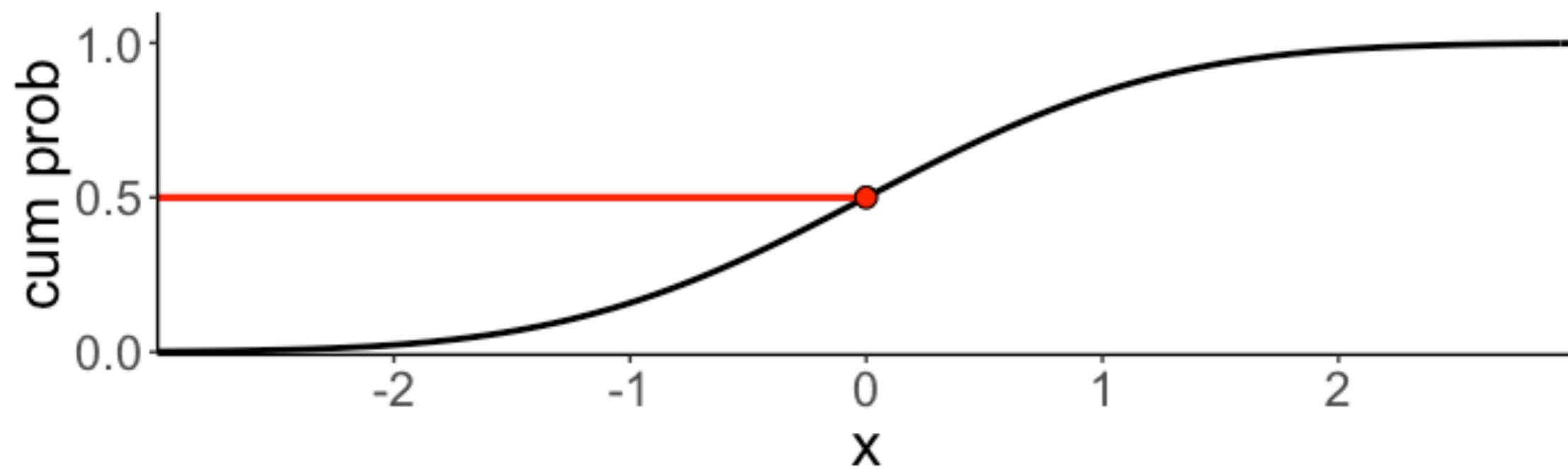


Probability vs. likelihood

dnorm ()

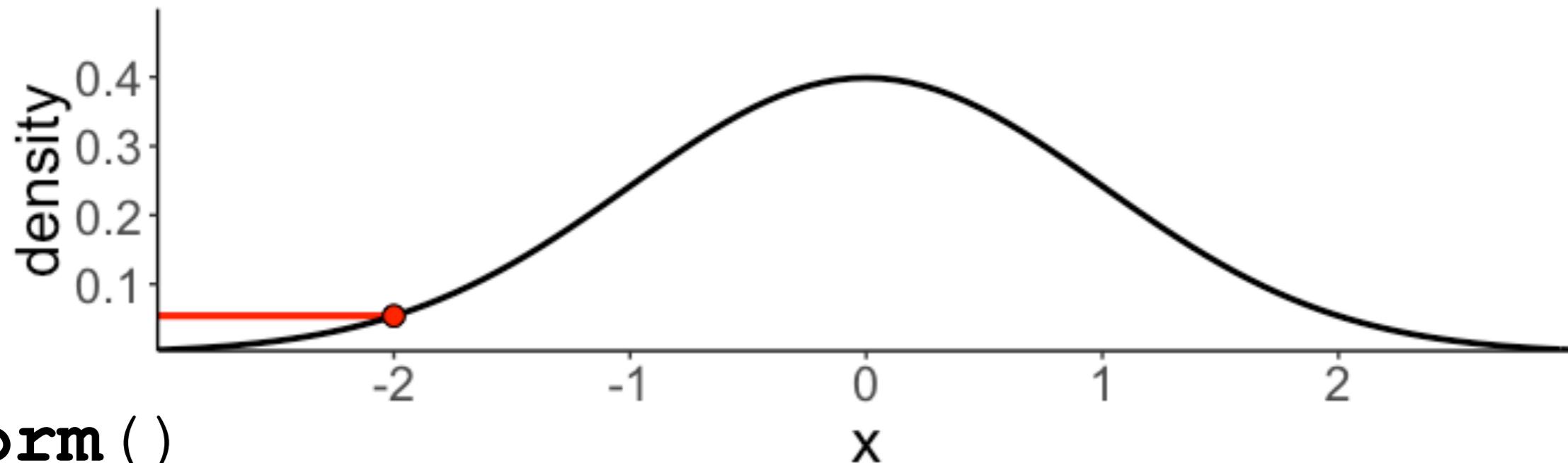


pnorm ()

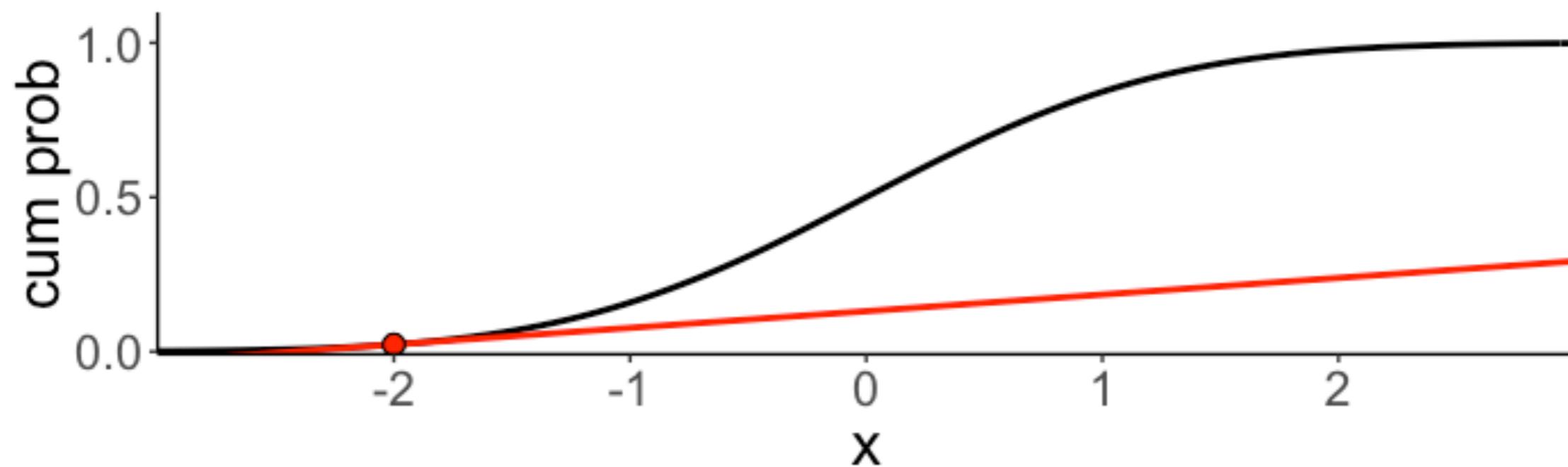


Probability vs. likelihood

dnorm ()



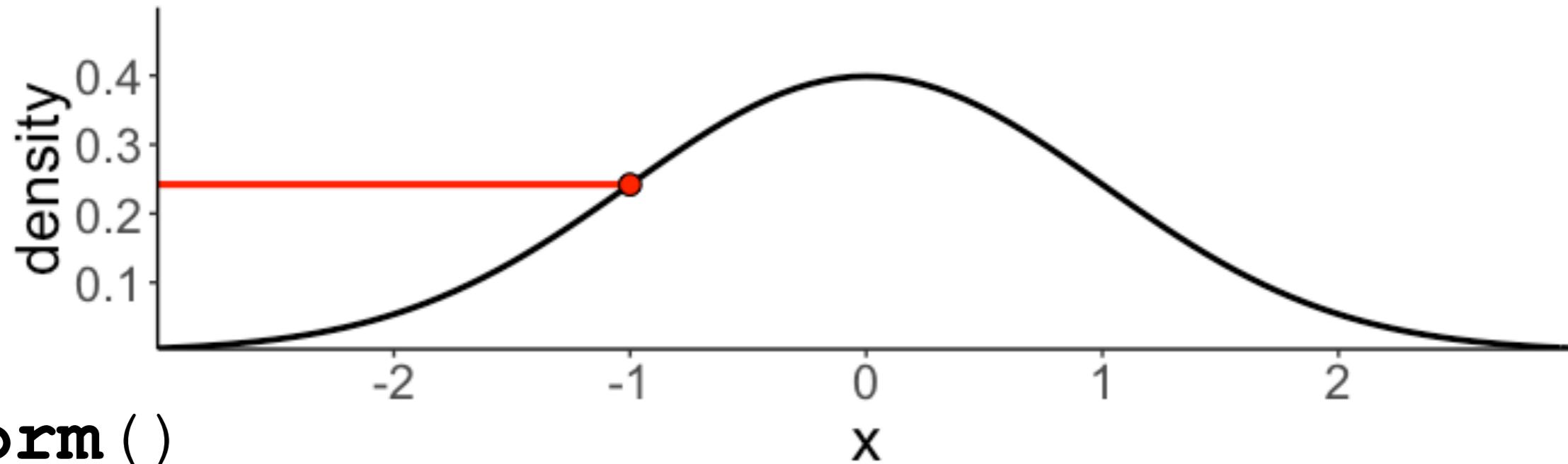
pnorm ()



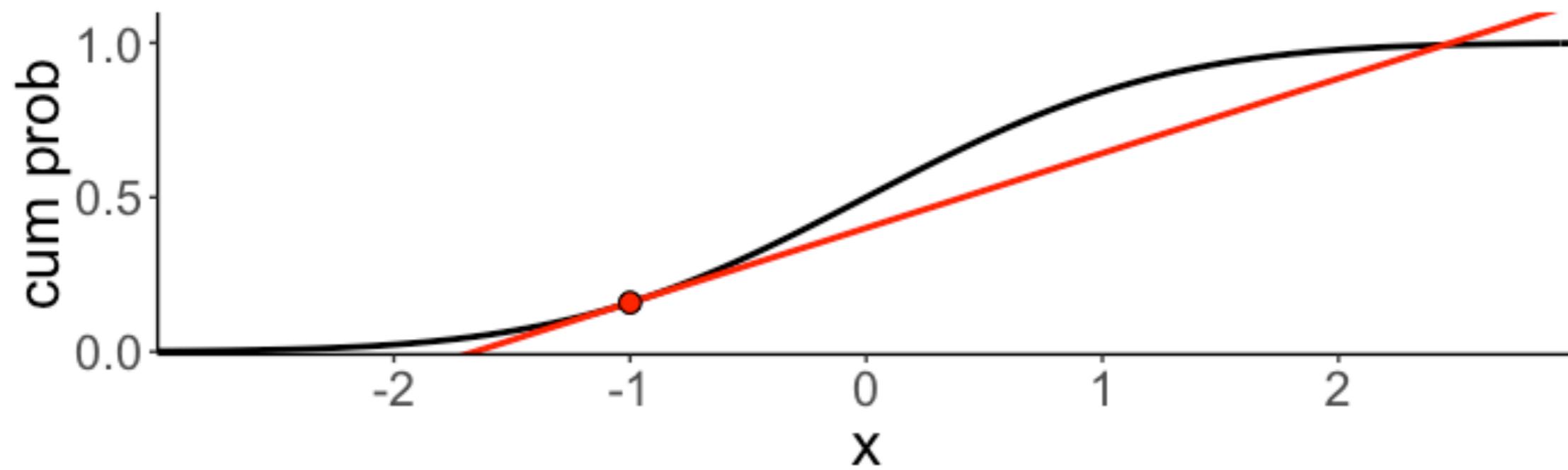
dnorm () is the first derivative of **pnorm ()**

Probability vs. likelihood

dnorm ()



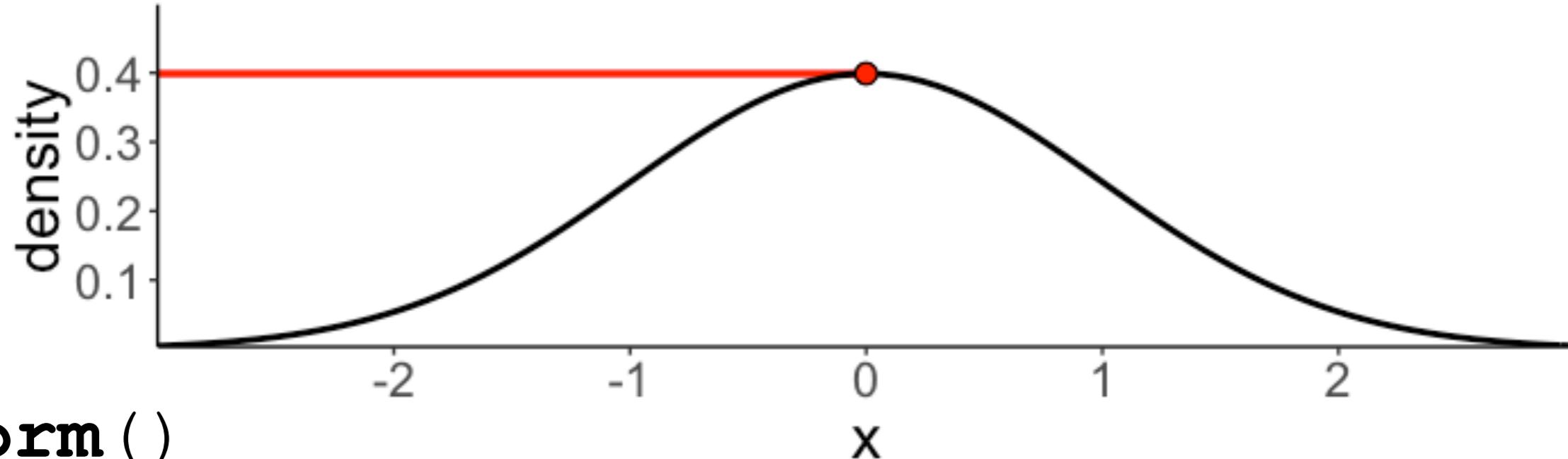
pnorm ()



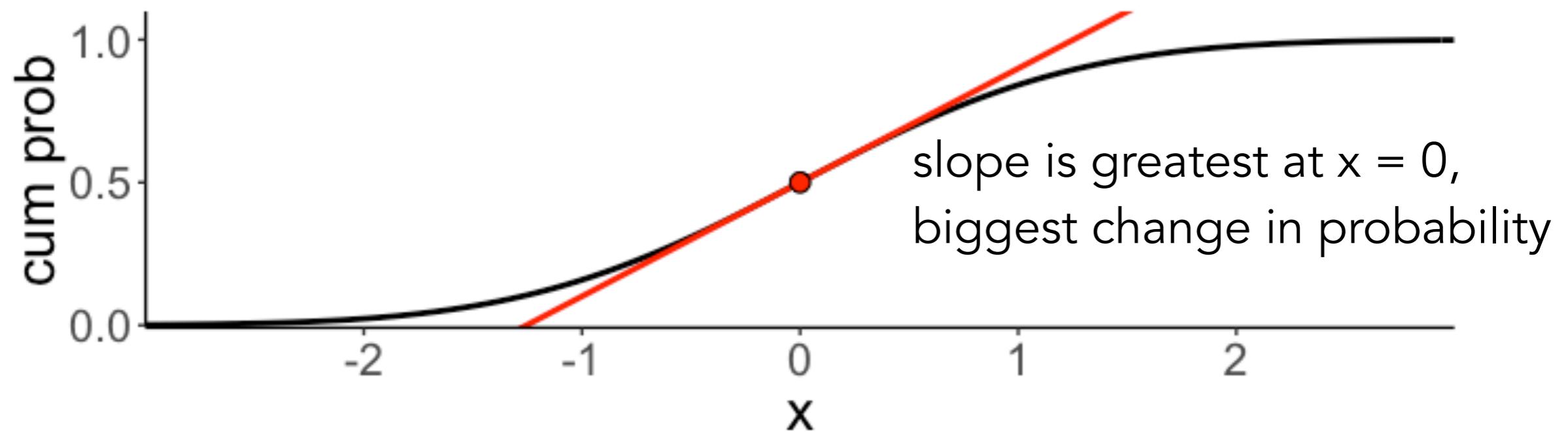
dnorm () is the first derivative of **pnorm ()**

Probability vs. likelihood

dnorm ()



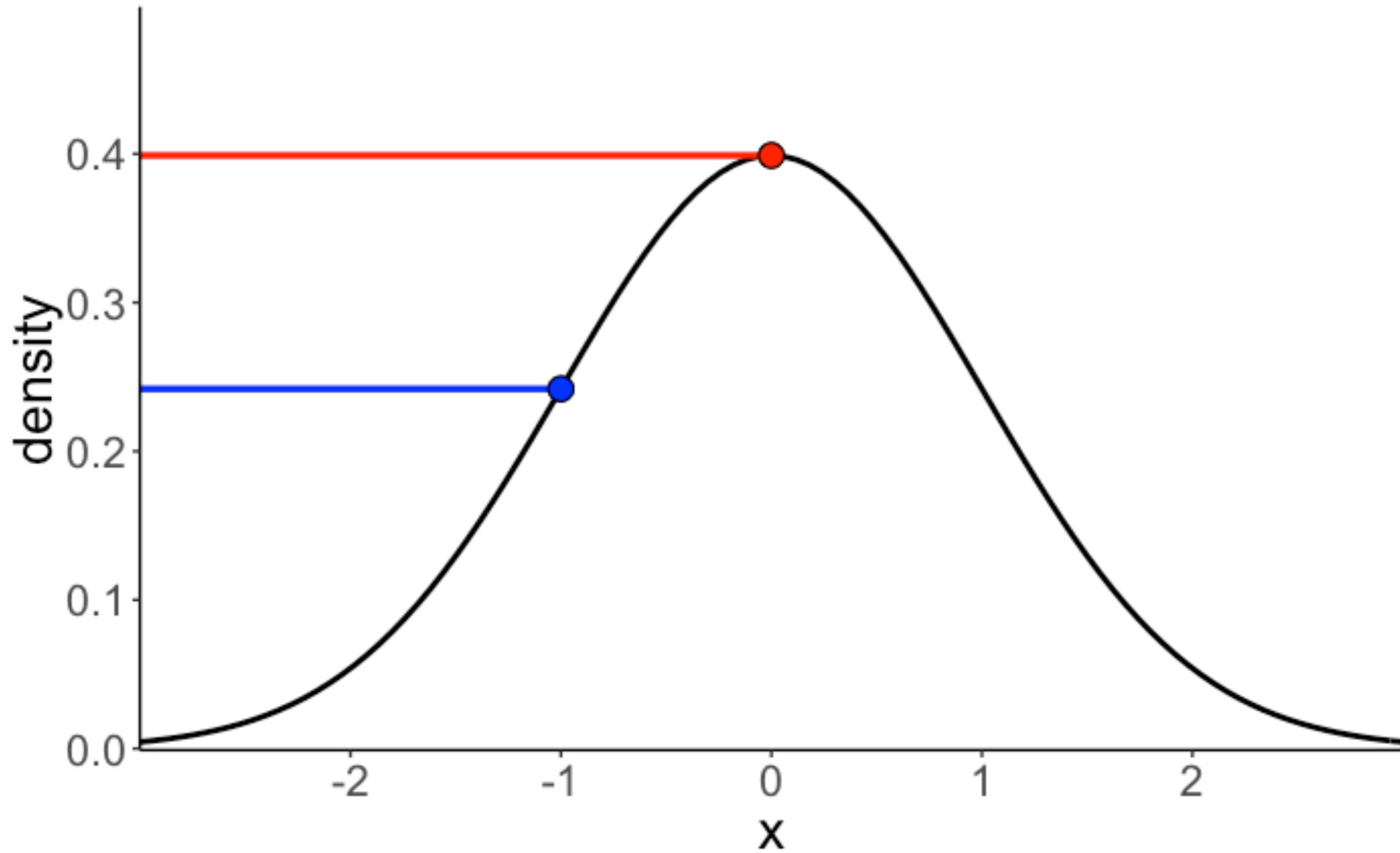
pnorm ()



dnorm () is the first derivative of **pnorm ()**

Probability vs. likelihood

$$\text{dnorm}(0) / \text{dnorm}(-1) = 1.6487$$

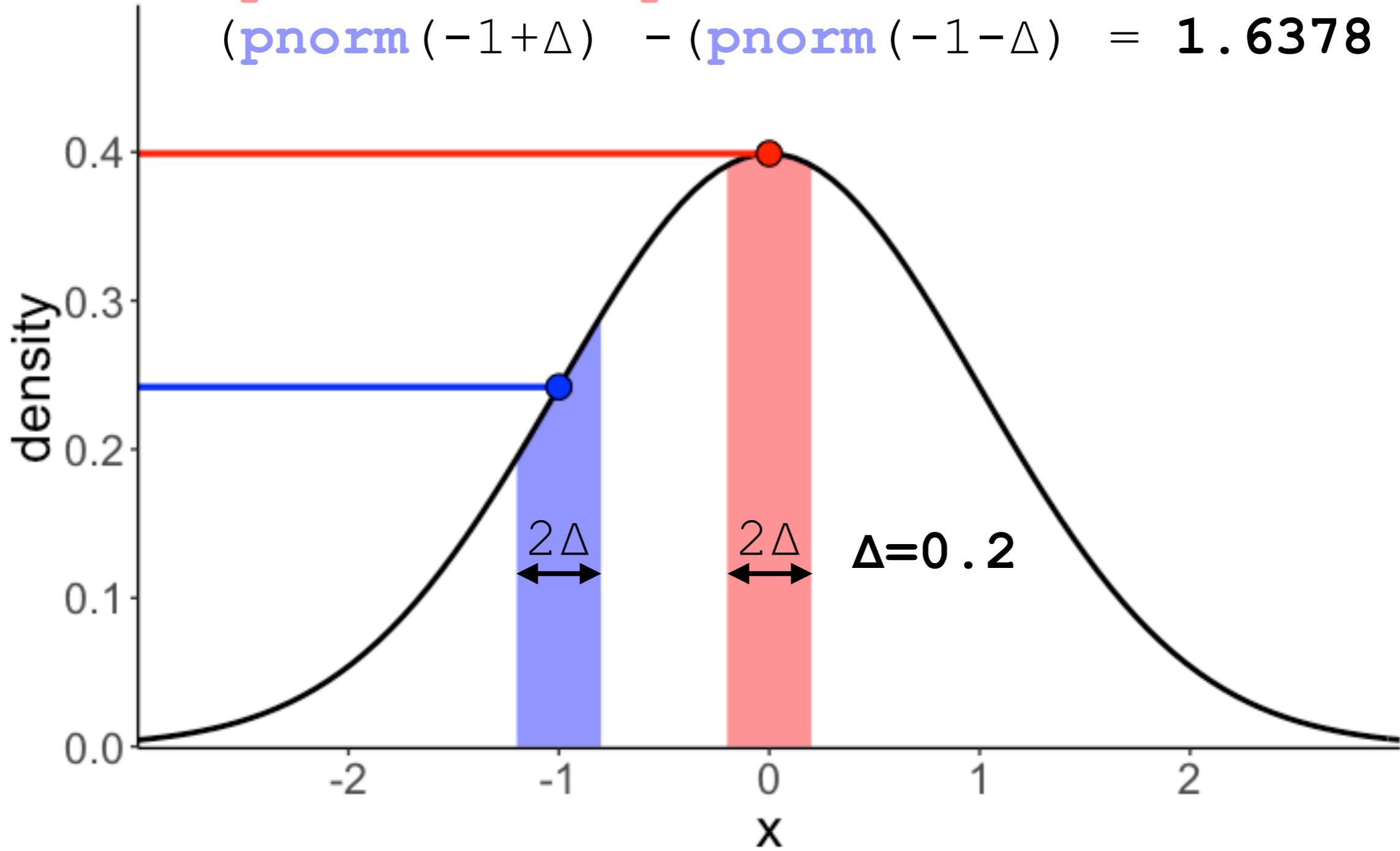


relative probability of one value vs. another

Probability vs. likelihood

$$\text{dnorm}(0) / \text{dnorm}(-1) = 1.6487$$

$$\frac{(\text{pnorm}(0+\Delta) - \text{pnorm}(0-\Delta))}{(\text{pnorm}(-1+\Delta) - \text{pnorm}(-1-\Delta))} = 1.6378$$

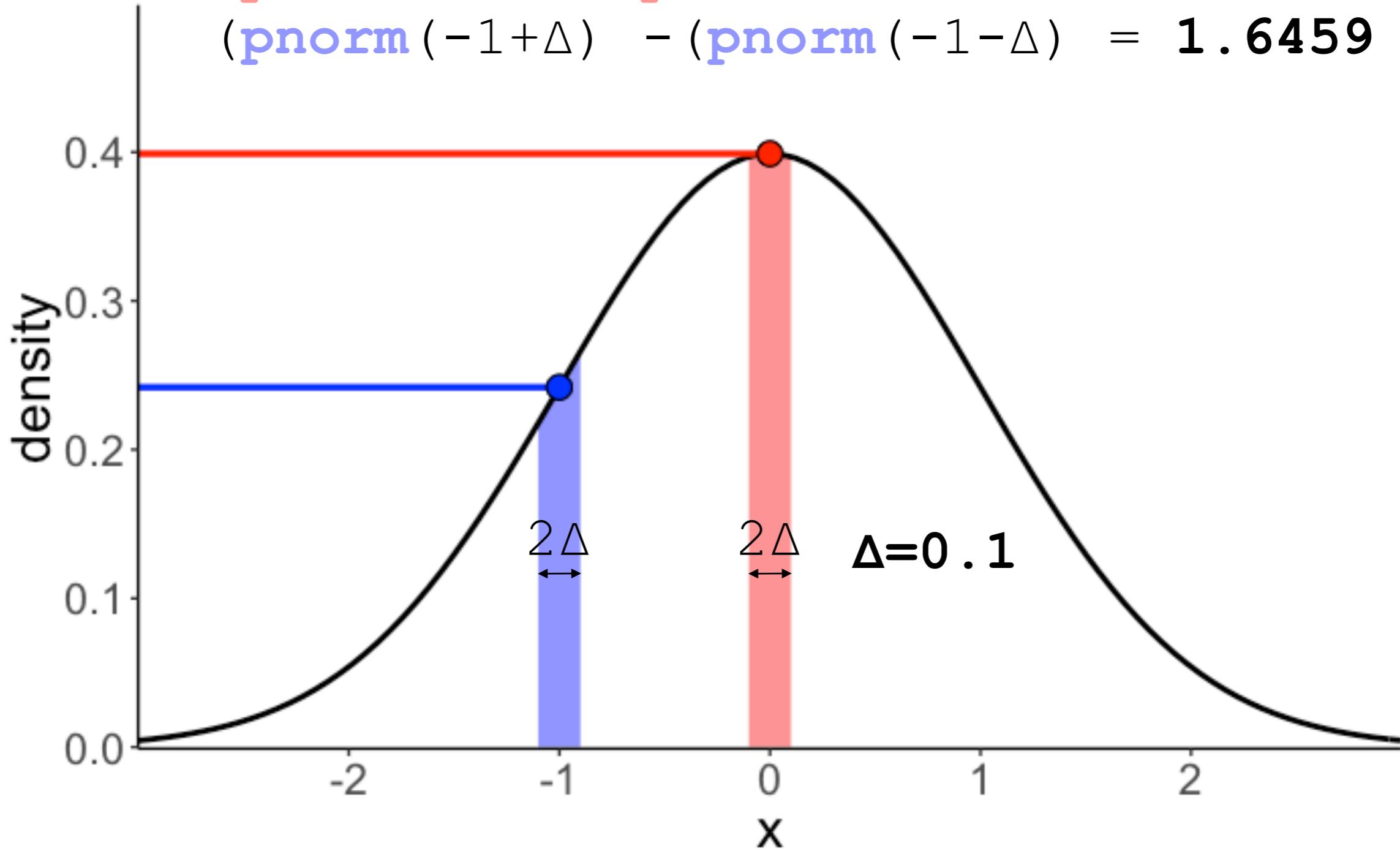


relative probability of one value vs. another

Probability vs. likelihood

$$\text{dnorm}(0) / \text{dnorm}(-1) = 1.6487$$

$$\frac{(\text{pnorm}(0+\Delta) - \text{pnorm}(0-\Delta))}{(\text{pnorm}(-1+\Delta) - \text{pnorm}(-1-\Delta))} = 1.6459$$

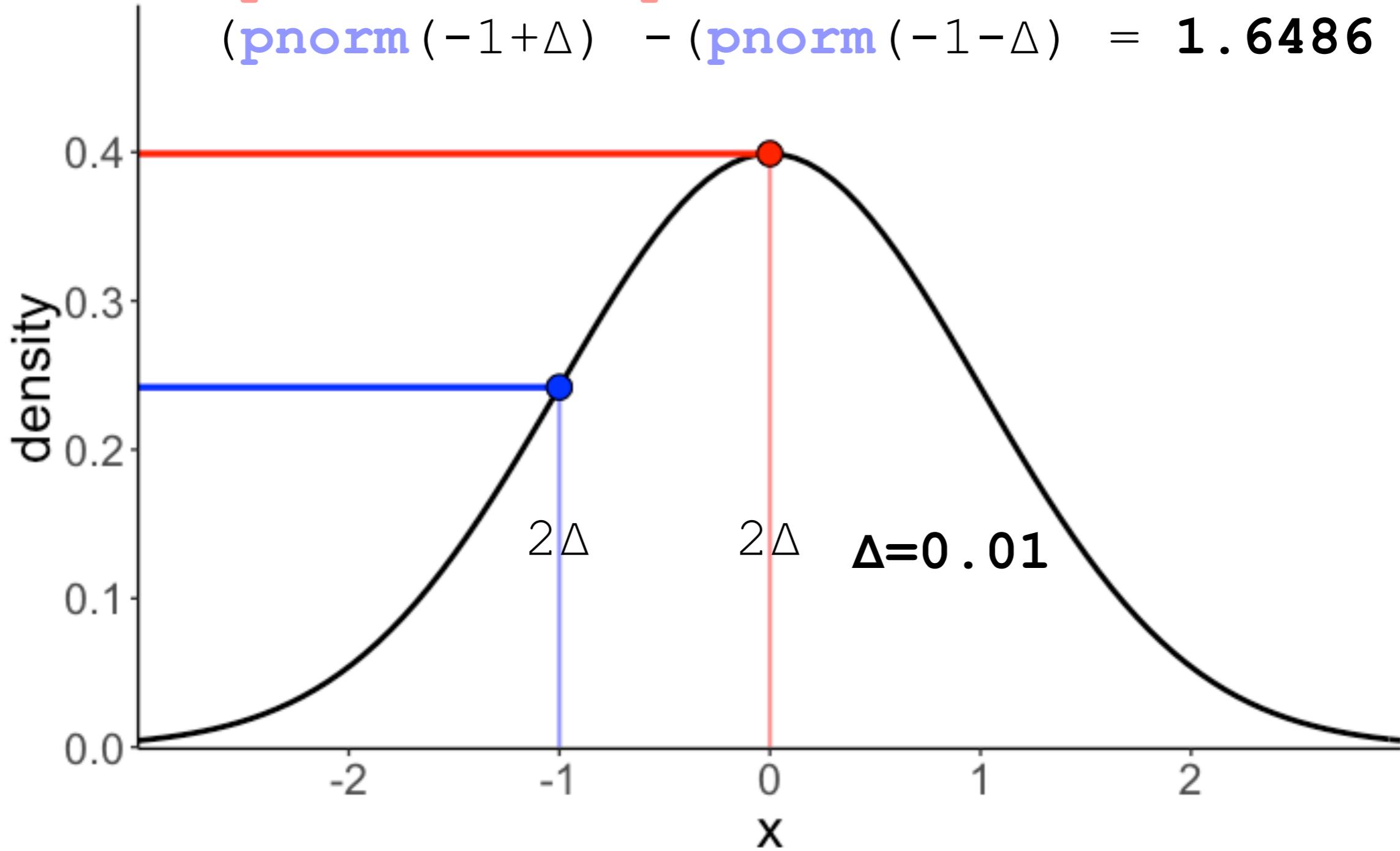


relative probability of one value vs. another

Probability vs. likelihood

$$\text{dnorm}(0) / \text{dnorm}(-1) = 1.6487$$

$$\frac{(\text{pnorm}(0+\Delta) - \text{pnorm}(0-\Delta))}{(\text{pnorm}(-1+\Delta) - (\text{pnorm}(-1-\Delta))} = 1.6486$$



relative probability of one value vs. another

Probability vs. likelihood

StatQuest makes me feel so happy....

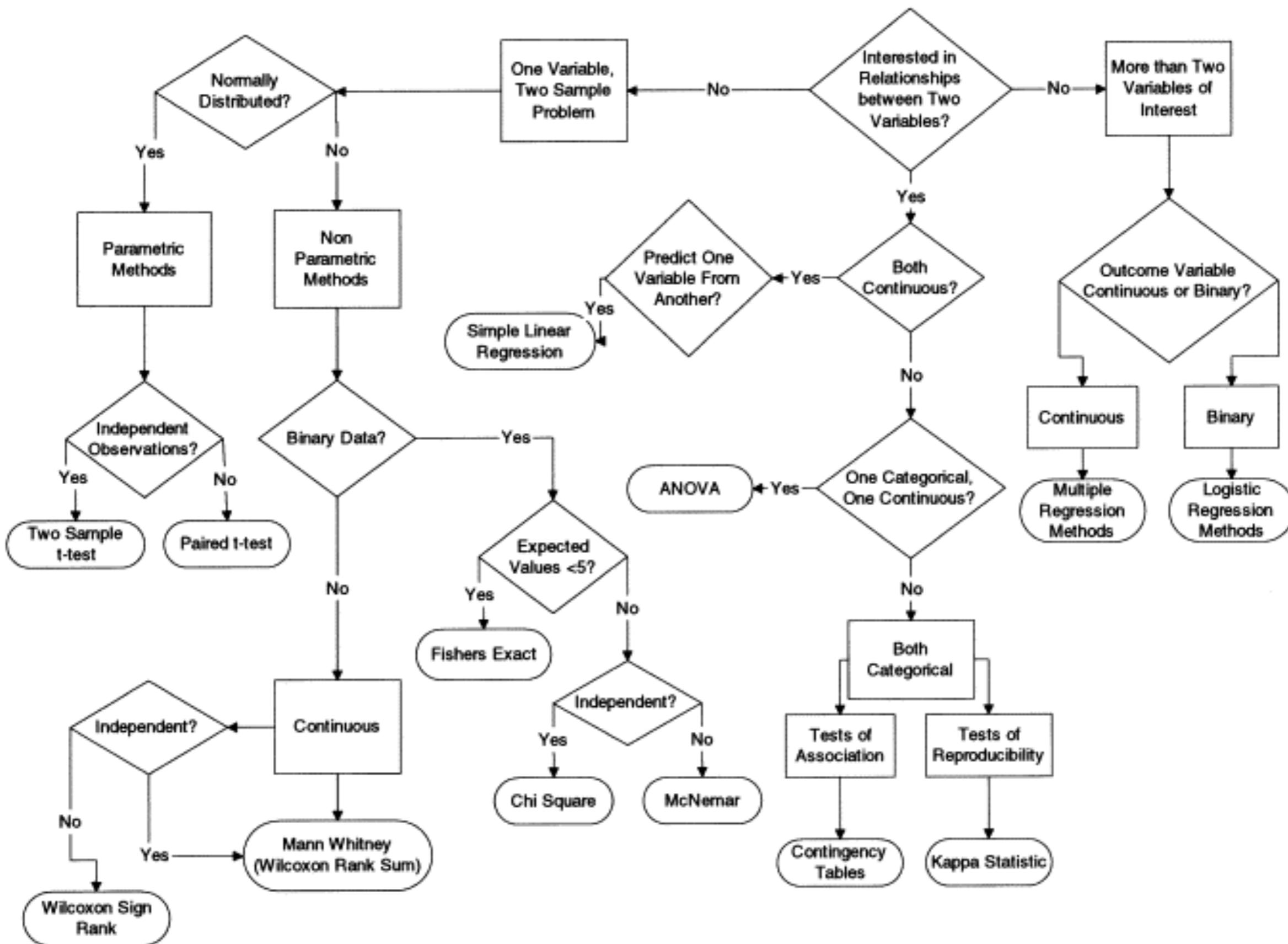
<https://www.youtube.com/watch?v=pYxNSUDSFH4>

Plan for today

- Cookbook vs. Model Comparison
- Modeling data
- Definitions of error and parameter estimates
- Models of error
- Statistical inferences about parameter values

Cookbook vs. Model Comparison

The cookbook approach

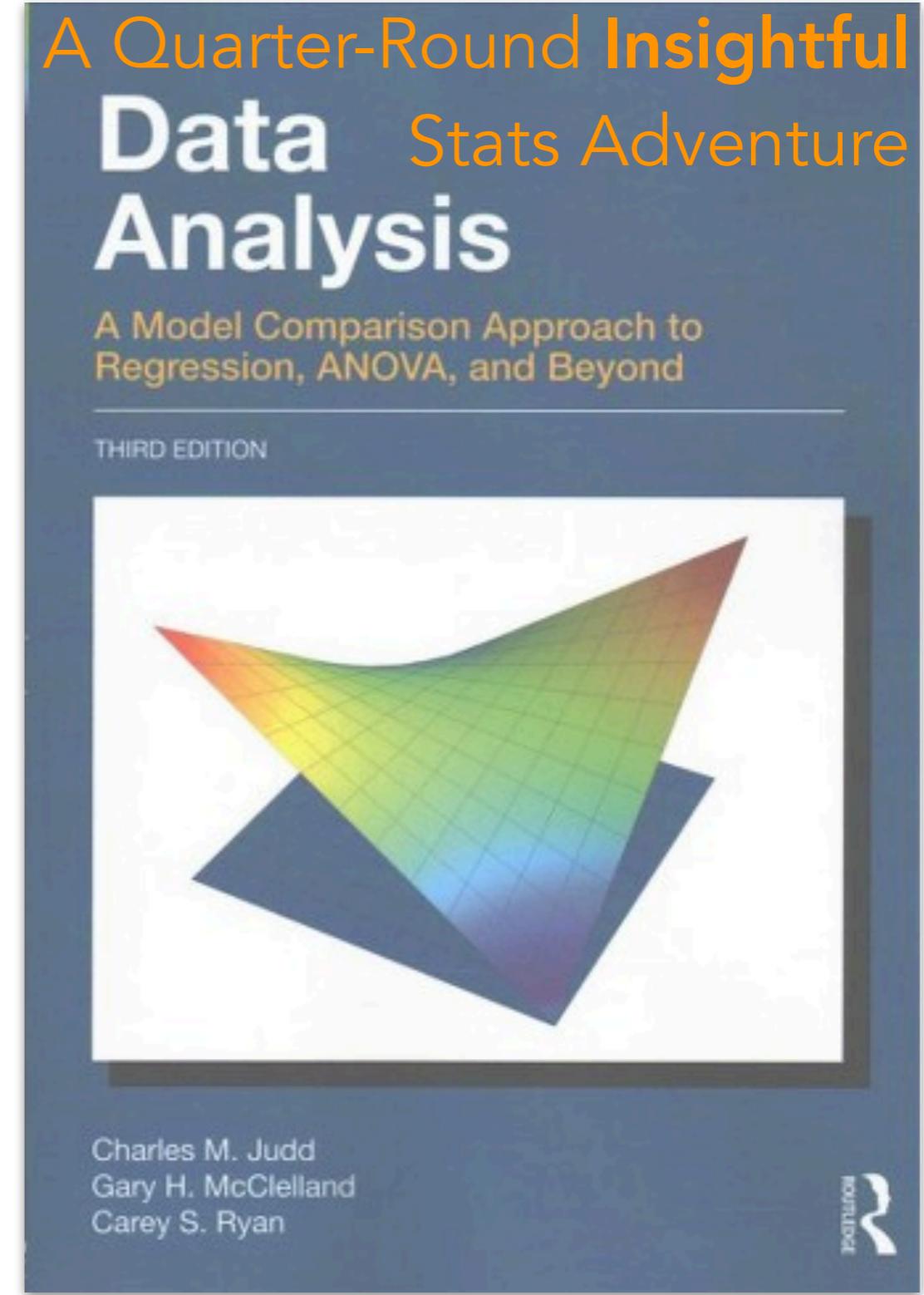


The cookbook approach



- many statistics textbooks are organized in this way
- works reasonably well if what we want to cook is in the book
- leaves us with no idea what to do if we can't find a recipe

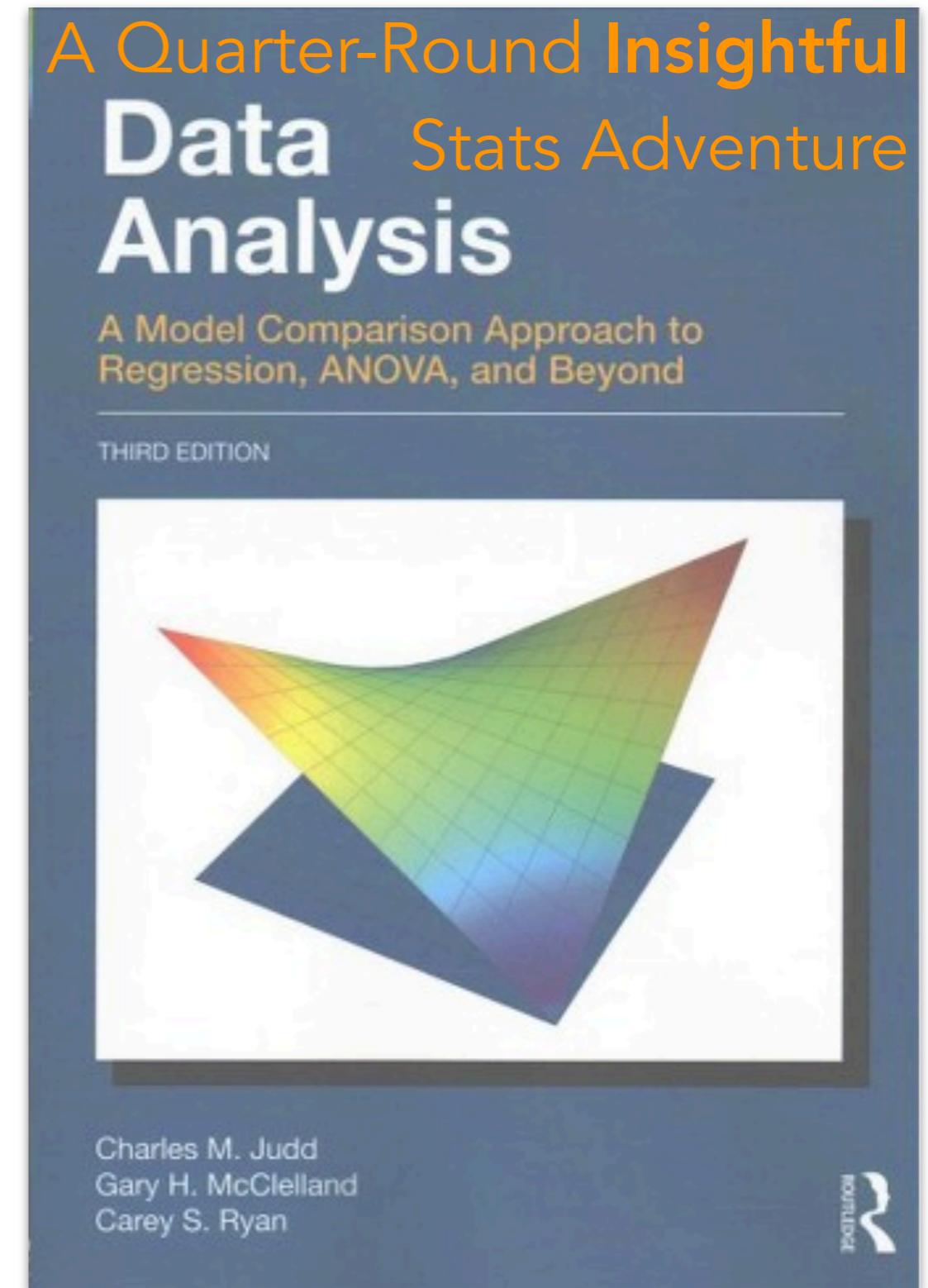
Model comparison approach



Judd, C. M., McClelland, G. H., & Ryan, C. S. (2011). *Data analysis: A model comparison approach*. Routledge.

Model comparison approach

- more flexible approach
- hopefully generates better insight
- thinking of statistical analysis as modeling
- allows for a smoother transition into Bayesian data analysis, and probabilistic modeling more generally



Judd, C. M., McClelland, G. H., & Ryan, C. S. (2011). *Data analysis: A model comparison approach*. Routledge.

Modeling data

Data = Model + Error



what's a good
model?



how shall we
define this?

= residual: the part that's left over after we have used the model to predict/explain the data



$$\text{Error} = \text{Data} - \text{Model}$$

to reduce error we can:

improve the quality of the data

e.g. run good experiments



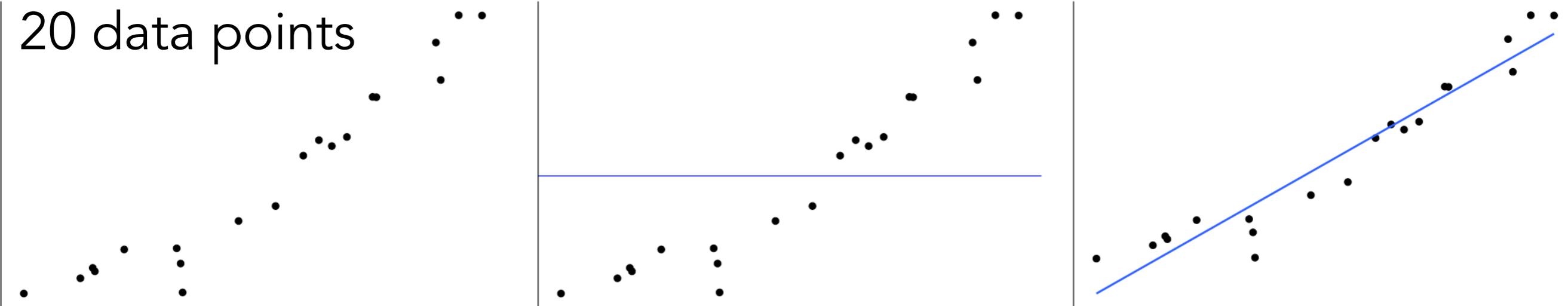
improve the model

e.g. make predictions conditional on additional information

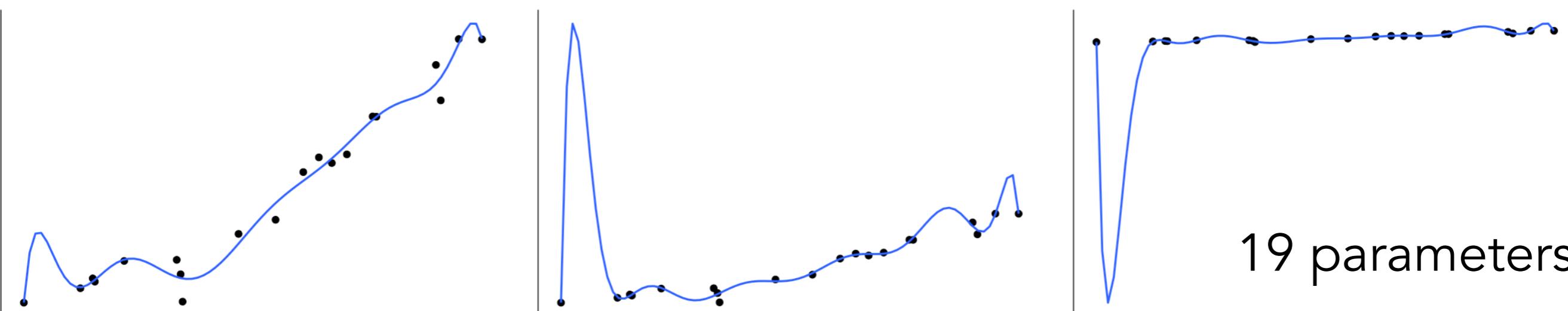
$$\text{Error} = \text{Data} - \text{Model}$$

- we build models with parameters, and fit those parameters to minimize error
- adding additional parameters to the model will always improve the model fit and reduce error
- fundamental trade-off between **simplicity** and **accuracy**

20 data points



Which model describes the data best?





**THE BEST WAY TO
EXPLAIN OVERFITTING**

Example

Percentage of households that had internet access in the year 2013 by US state

i	internet	state	college	auto	density
1	79.0	AK	28.0	1.2	1.2
2	63.5	AL	23.5	1.3	94.4
3	60.9	AR	20.6	1.7	56.0
4	73.9	AZ	27.4	1.3	56.3
5	77.9	CA	31.0	0.8	239.1
6	79.4	CO	37.8	1.0	48.5
7	77.5	CT	37.2	1.0	738.1
8	74.5	DE	29.8	1.1	460.8
9	74.3	FL	27.2	1.2	350.6
10	72.2	GA	28.3	1.1	168.4

set before looking at the data

$$\text{model}_1: Y_i = 75 + \text{ERROR}$$

worth it?

fit to the data

$$\text{model}_2: Y_i = \beta_0 + \text{ERROR}$$

worth it?

additional predictor

$$\text{model}_3: Y_i = \beta_0 + \beta_1 X_i + \text{ERROR}$$

fit to the data

Compact model

model_C: $Y_i = \beta_0 + \text{ERROR}$

Augmented model

model_A: $Y_i = \beta_0 + \beta_1 X_i + \text{ERROR}$

$$\text{ERROR}(C) \geq \text{ERROR}(A)$$

Proportional reduction in error (PRE)

$$\text{PRE} = \frac{\text{ERROR}(C) - \text{ERROR}(A)}{\text{ERROR}(C)}$$

Compact model

$$\text{model}_C: Y_i = \beta_0 + \text{ERROR}$$

$$\text{ERROR}(C) = 50$$

Augmented model

$$\text{model}_A: Y_i = \beta_0 + \beta_1 X_i + \text{ERROR}$$

$$\text{ERROR}(A) = 30$$

Proportional reduction in error (PRE)

$$\begin{aligned} \text{PRE} &= 1 - \frac{\text{ERROR}(A)}{\text{ERROR}(C)} \\ &= 1 - \frac{30}{50} = .40 \end{aligned}$$

Increasing the complexity of the model reduced the error by 40%. **worth it?**

worth it?

Compact model

model_C: $Y_i = \beta_0 + \text{ERROR}$

Augmented model

model_A: $Y_i = \beta_0 + \beta_1 X_i + \text{ERROR}$

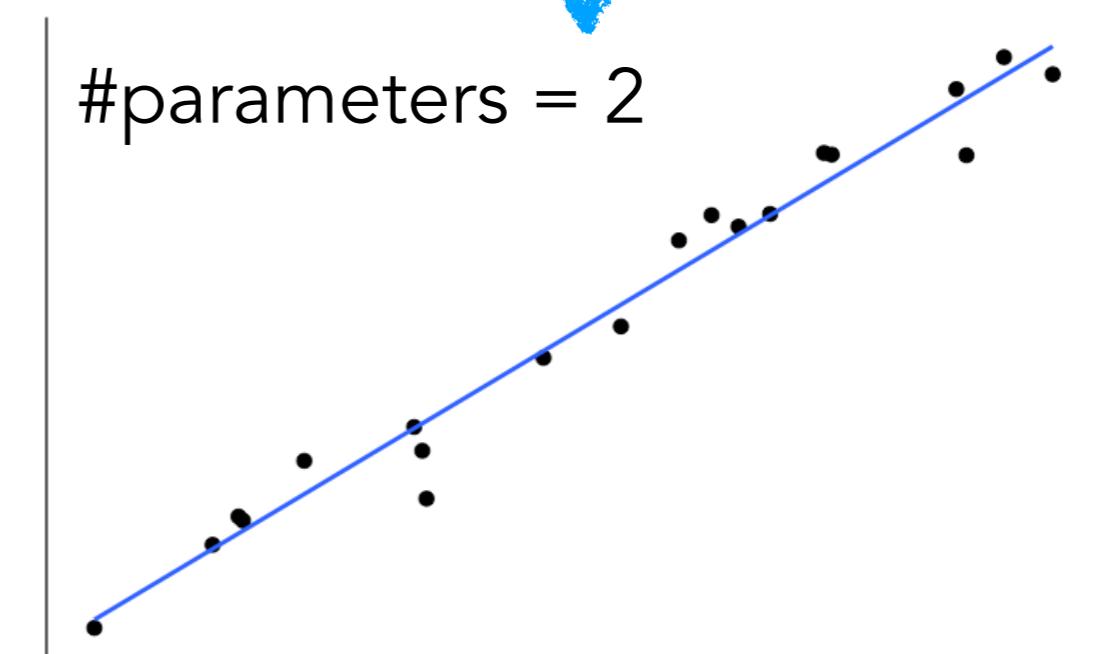
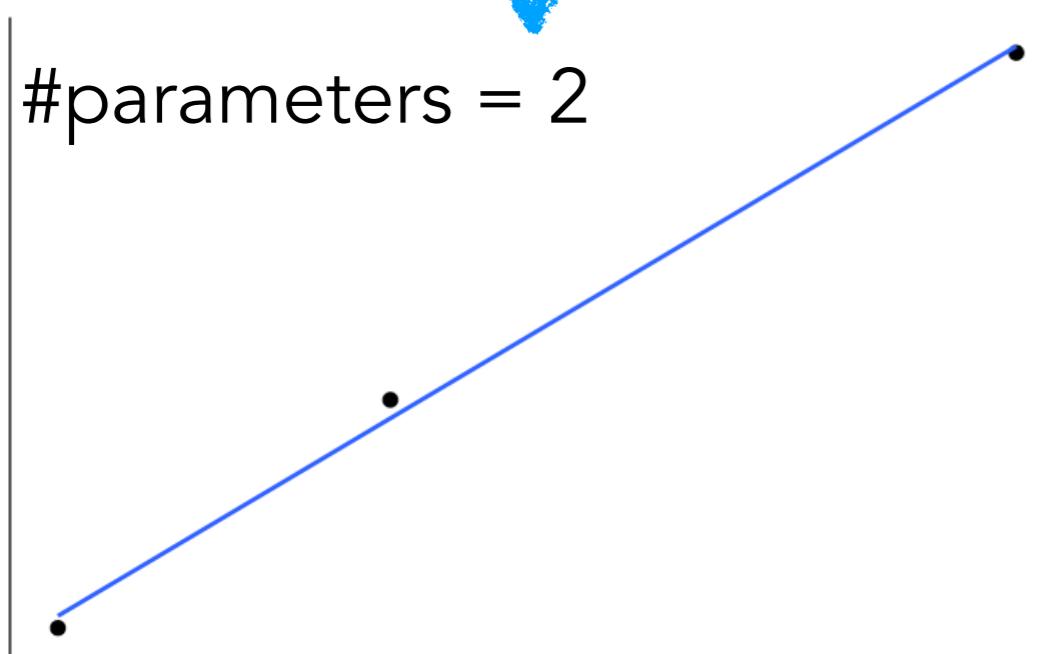
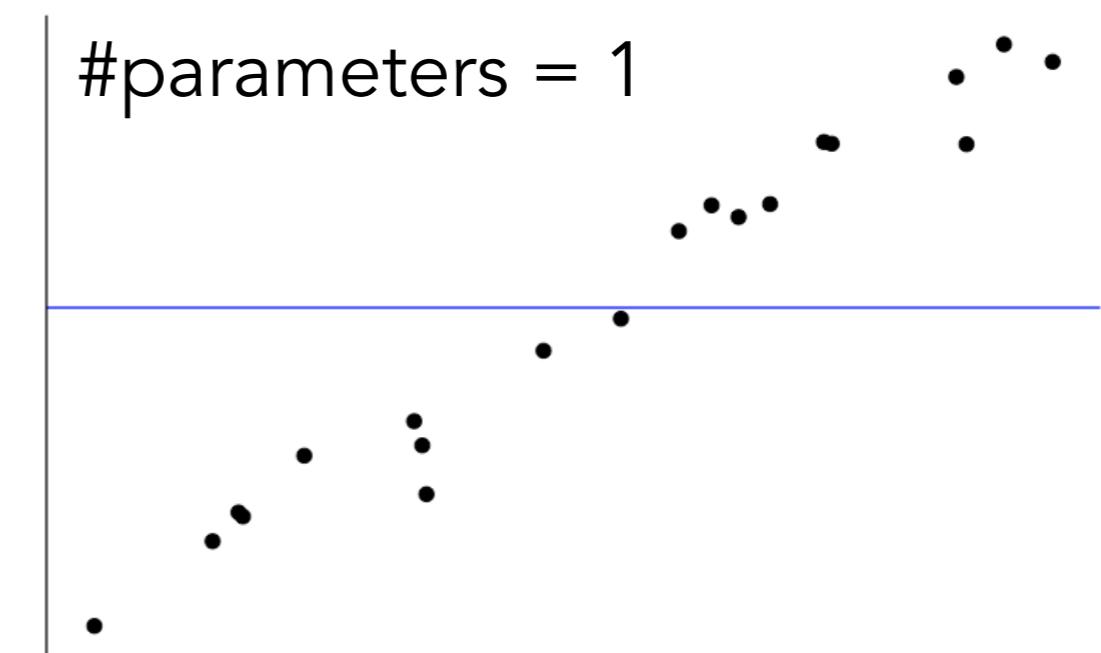
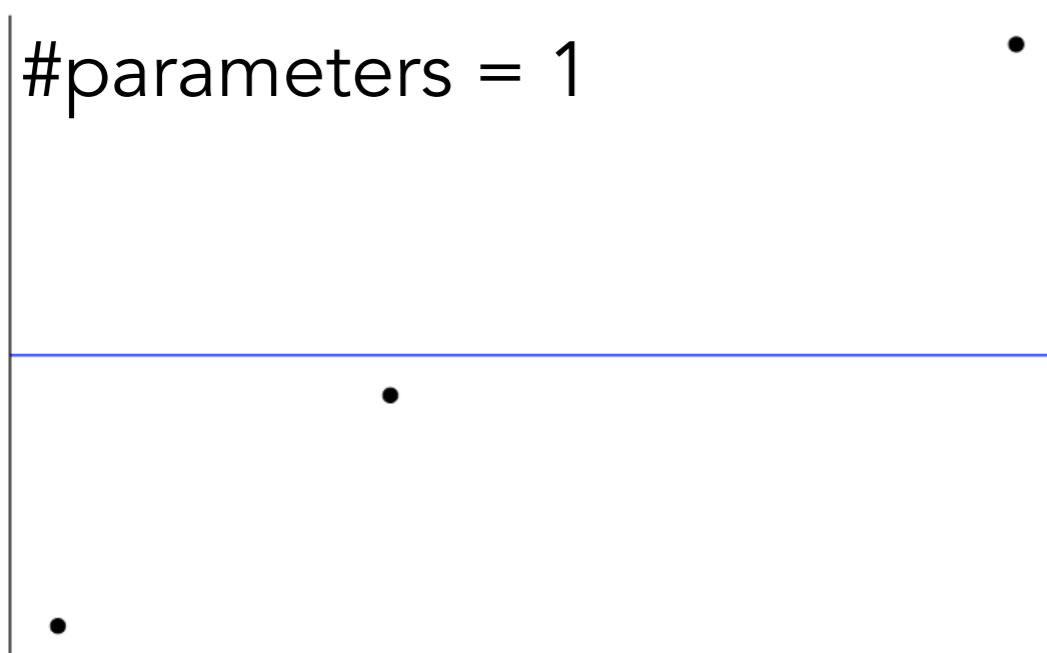
Proportional reduction in error (PRE)

$$\begin{aligned}\text{PRE} &= 1 - \frac{\text{ERROR}(A)}{\text{ERROR}(C)} \\ &= 1 - \frac{30}{50} = .40\end{aligned}$$

- to answer the **worth it?** question we need inferential statistics (define ERROR, sampling distributions, etc.)
- more likely to be **worth it** if:
 1. **PRE** is high
 2. the number of additional parameters in A compared to C is low (**PRE per additional parameter**)
 3. the number of parameters that could have been added to model_C to create model_A but were not is low

more impressed if the number of observations n is much greater than the number of parameters

PRE per parameter for different n



neato!



impressive!

General procedure

- for any question we want to ask about our DATA
 - we define model_C and model_A
 - compare the models using PRE
 - determine whether PRE is **worth it**
 - in standard frequentist lingo:
 - model_C = H_0 (null hypothesis) 
 - model_A = H_1 (alternative hypothesis) 
 - hypothesis test:
 - H_0 : **all** the parameters that are included in model_A but not in model_C are 0
 - H_1 : **not all** the parameters that are included in model_A but not in model_C are 0
- model comparison**

Notation

DATA = MODEL + ERROR

$$Y_i = \beta_0 + \epsilon_i \text{ simple model (true parameters)}$$

$$Y_i = b_0 + e_i \text{ simple model (estimated parameters)}$$

$$\hat{Y}_i = b_0$$

college

$$Y_i = b_0 + b_1 X_{i1} + e_i \text{ more complex model}$$

Percentage of households that had internet access in the year 2013 by US state

i	internet	state	college	auto	density
1	79.0	AK	28.0	1.2	1.2
2	63.5	AL	23.5	1.3	94.4
3	60.9	AR	20.6	1.7	56.0
4	73.9	AZ	27.4	1.3	56.3
5	77.9	CA	31.0	0.8	239.1
6	79.4	CO	37.8	1.0	48.5
7	77.5	CT	37.2	1.0	738.1
8	74.5	DE	29.8	1.1	460.8
9	74.3	FL	27.2	1.2	350.6
10	72.2	GA	28.3	1.1	168.4



Greek letters β or ϵ represent the true but unknowable parameters in the population.

Roman letters b or e represent estimates of these parameters using our DATA.

Definitions of error and parameter estimates

Definitions of error and parameter estimates

1. How should individual errors be aggregated into a summary index ERROR?
 - sum of absolute errors
 - sum of squared errors
 - count of errors
2. What's the best estimator of the data for each kind of error?
3. Which error shall we choose?

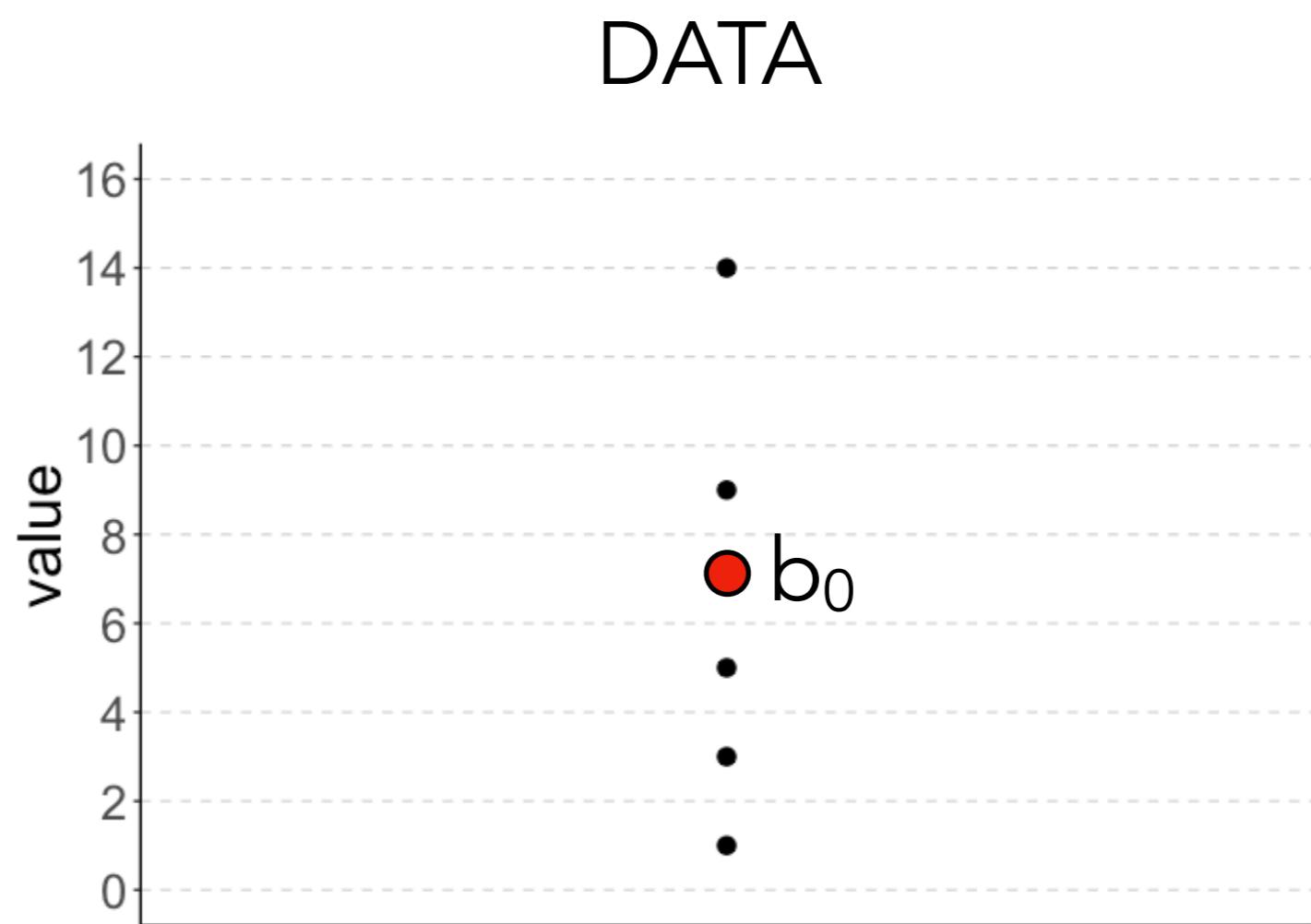
Simple model

$$Y_i = \beta_0 + \epsilon_i$$

$$Y_i = b_0 + e_i$$

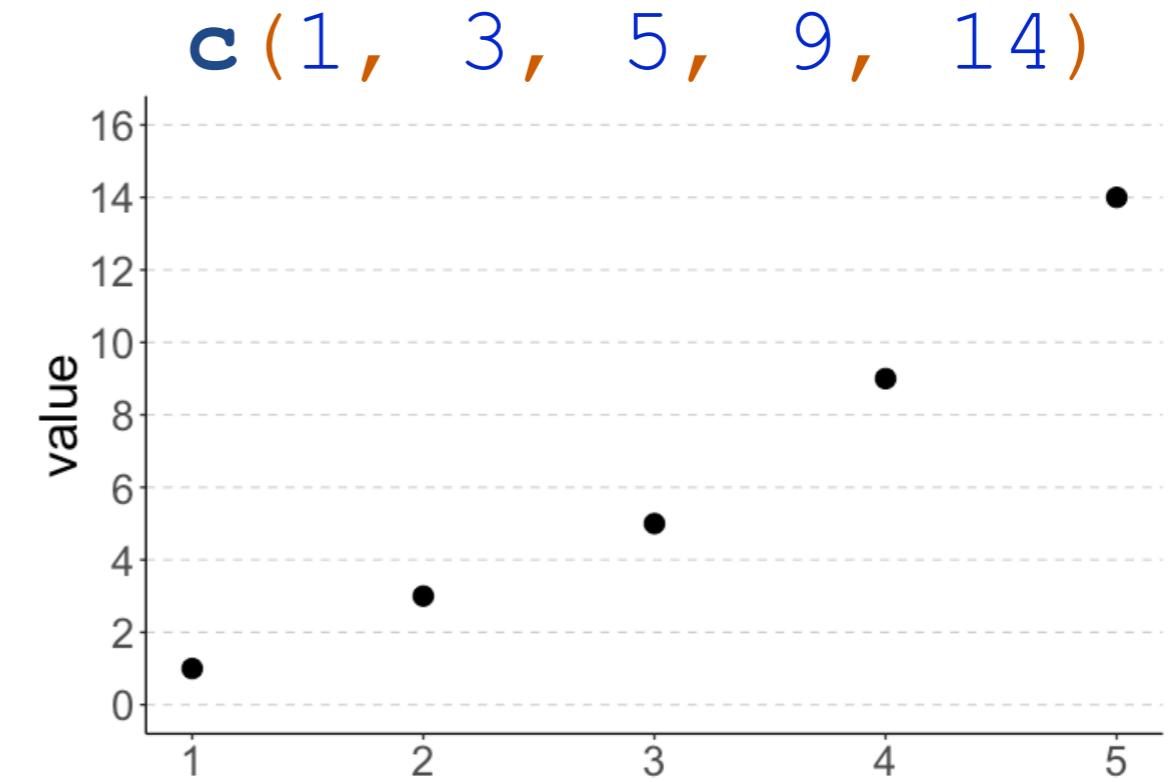
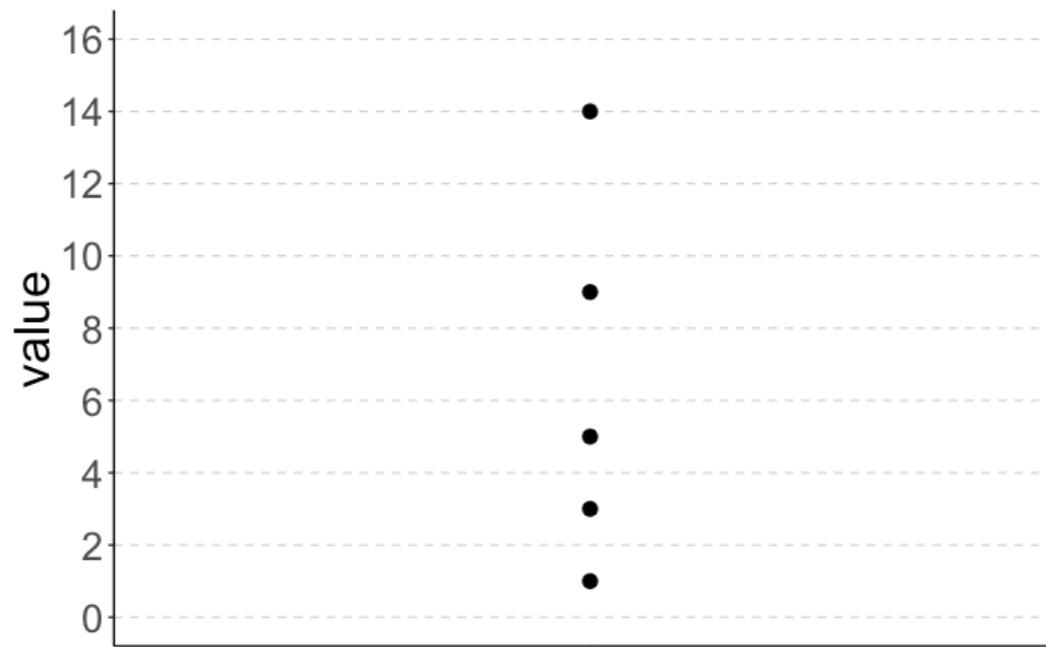
depends on how the
error is defined!

what value is the best
model of the data?

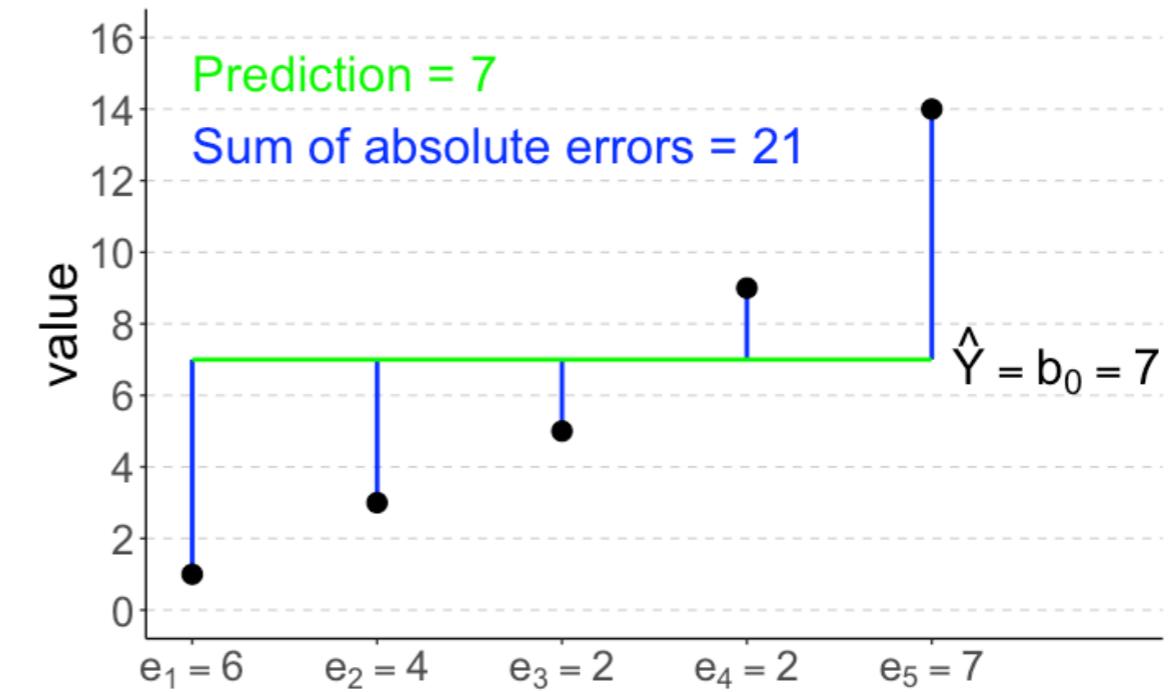
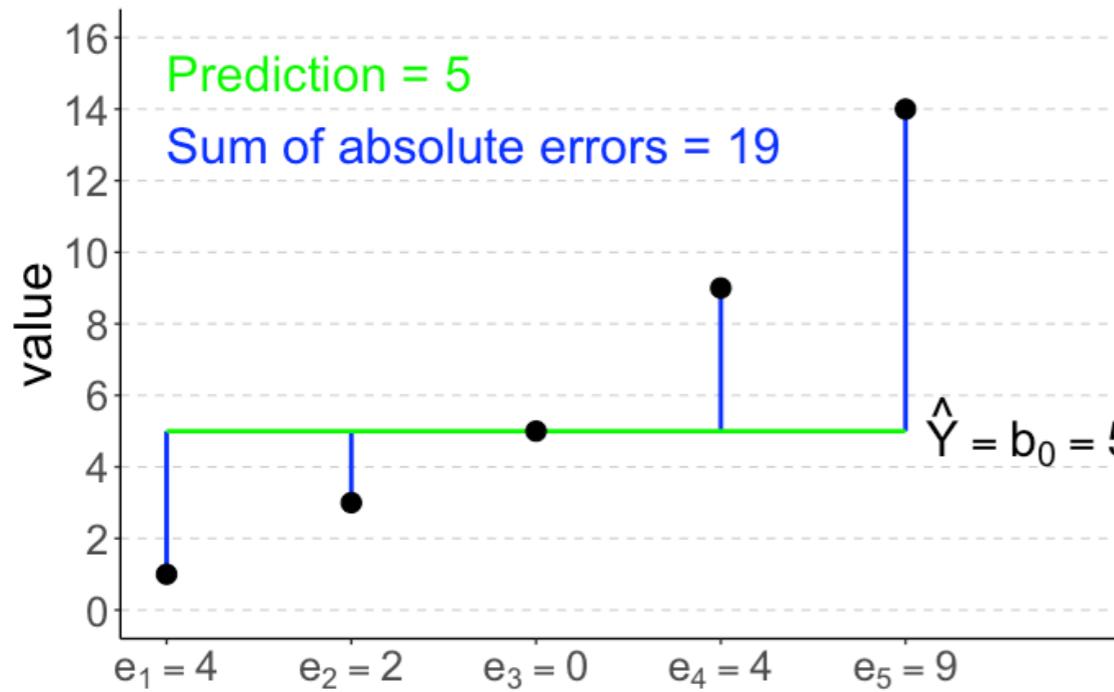


Error = sum of absolute errors

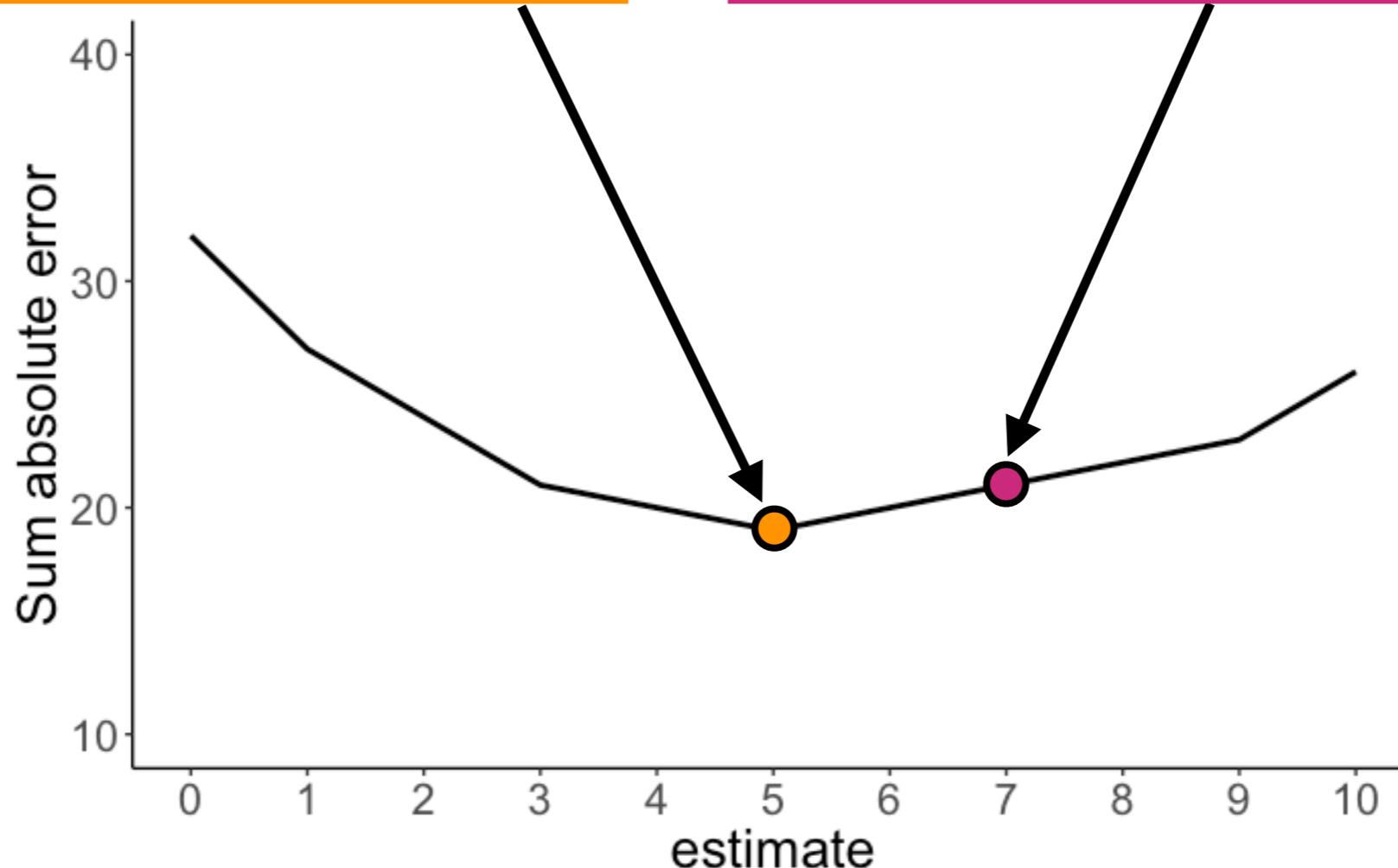
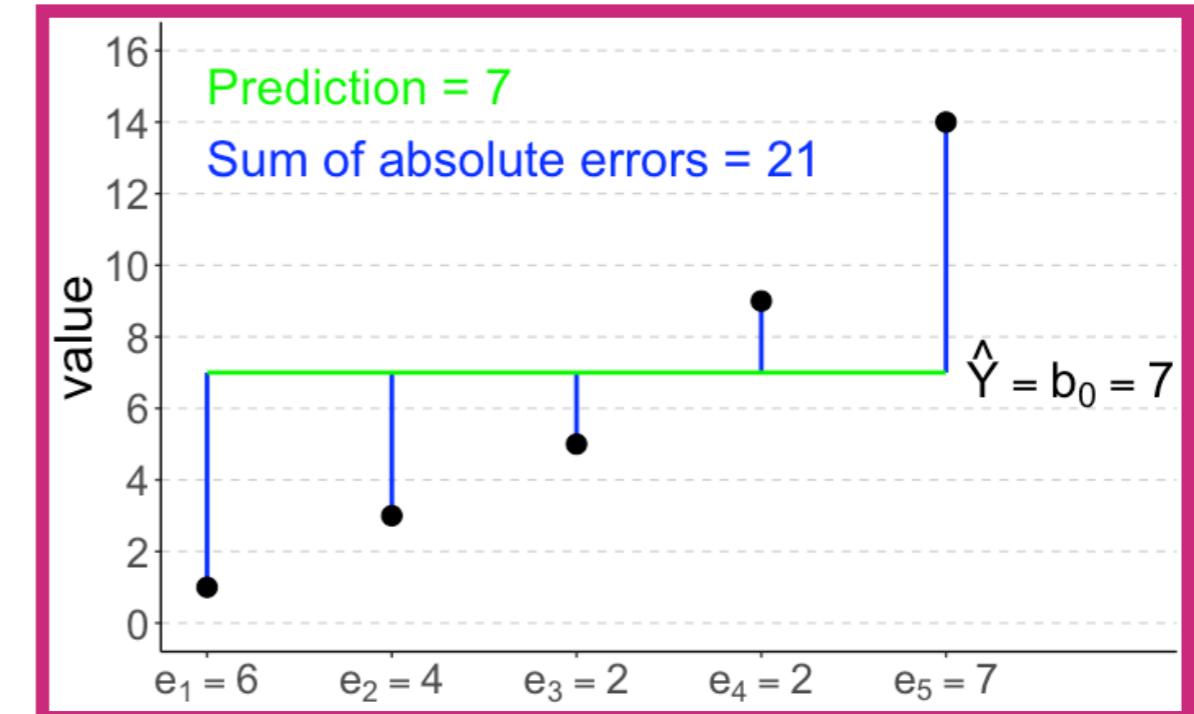
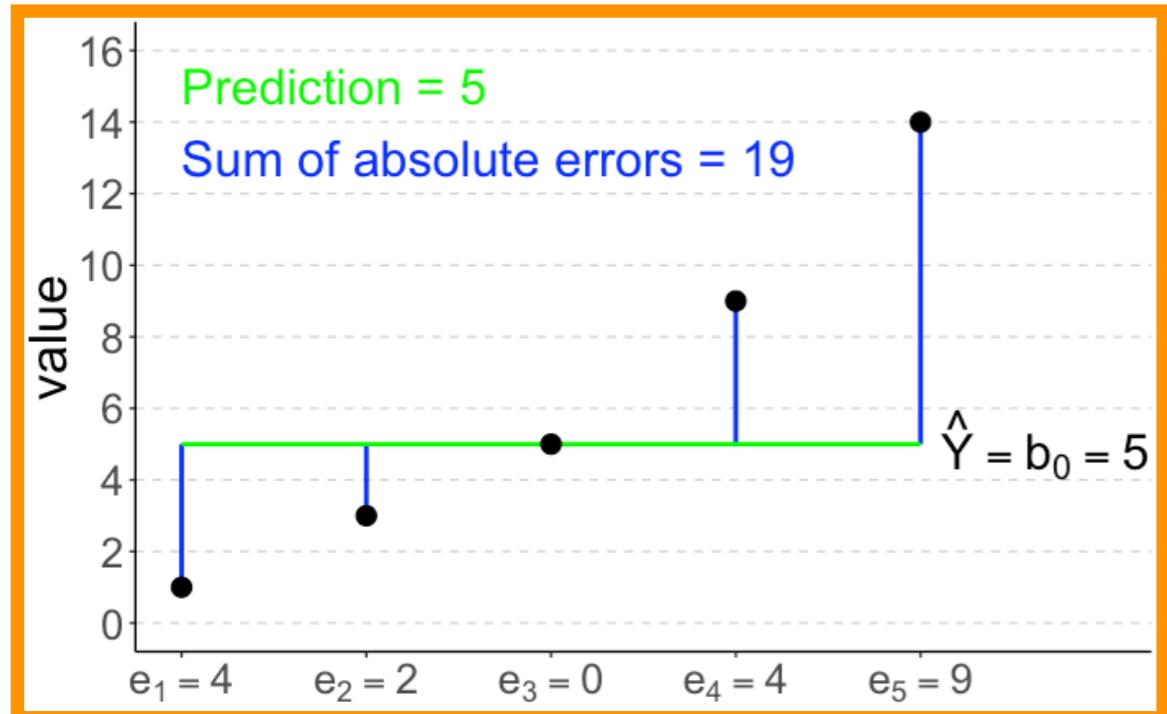
Model



Error as the sum of **line lengths**



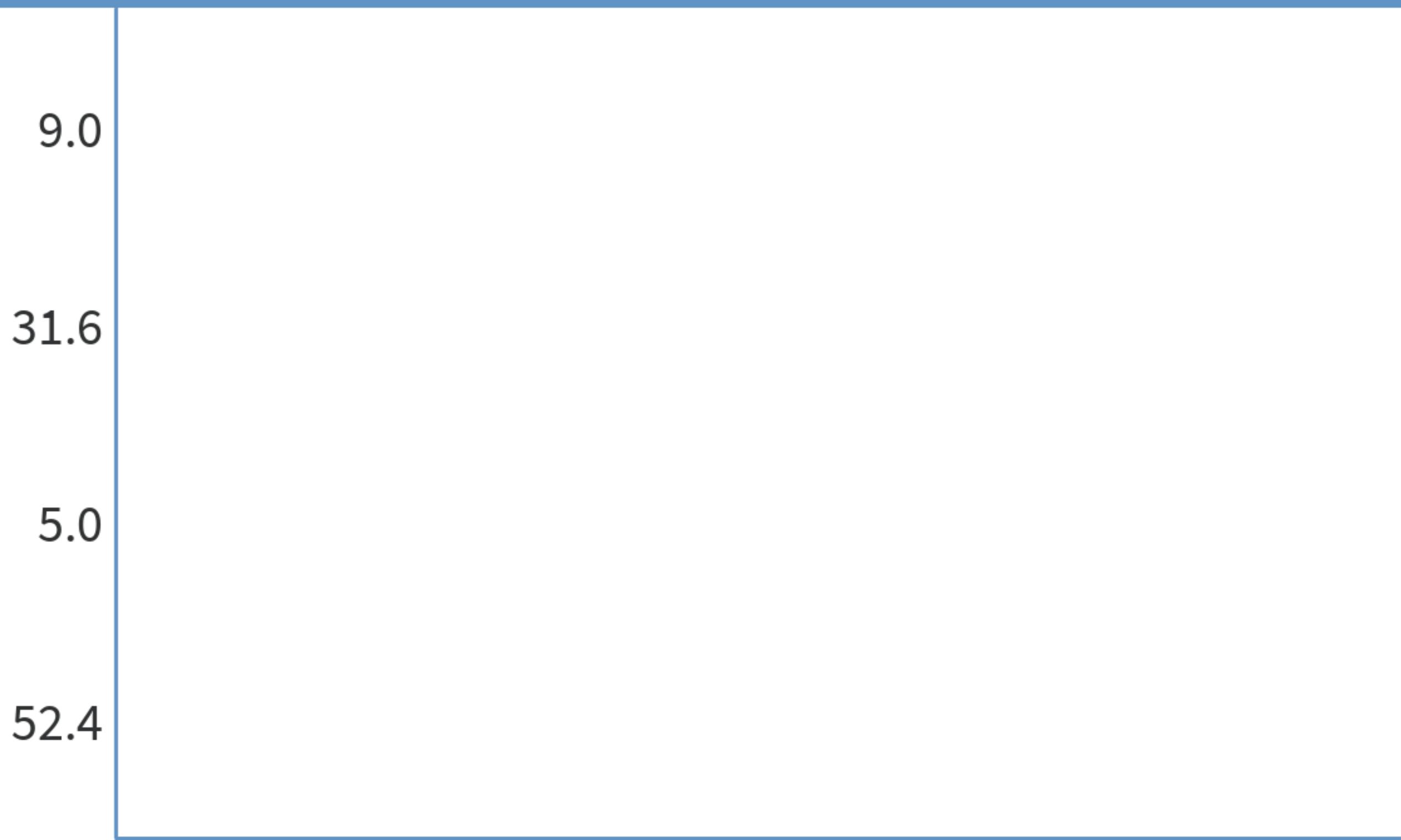
What's the best simple model b_0 for this error measure?



What if the blue value was 140 instead of 14?
What would be the best estimate of b_0 then?



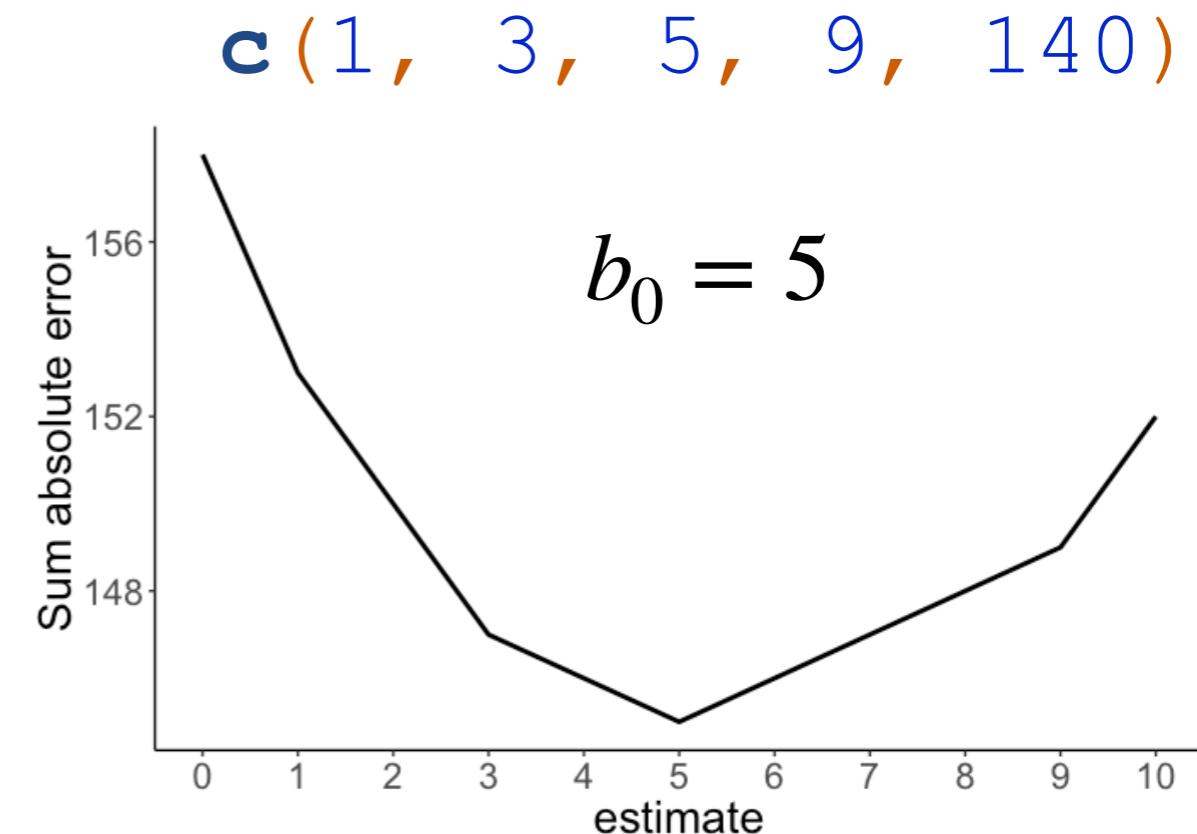
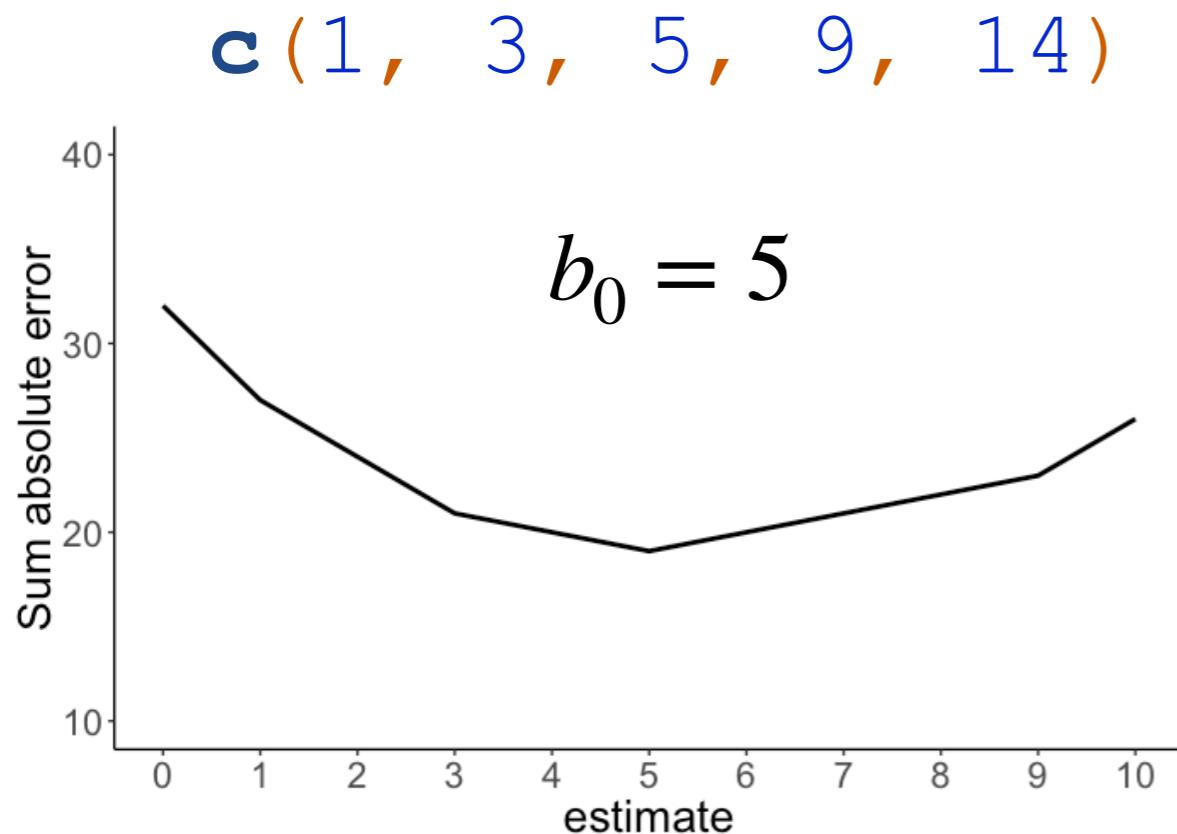
**What's the estimate that minimizes the sum of absolute errors
for the values 1, 3, 5, 9, 140?**



Sum of absolute errors

$$Y_i = b_0 + e_i$$

$$\text{ERROR} = \sum_{i=1}^n |e_i| = \sum_{i=1}^n |Y_i - b_0|$$



the **median** minimizes the sum of absolute errors

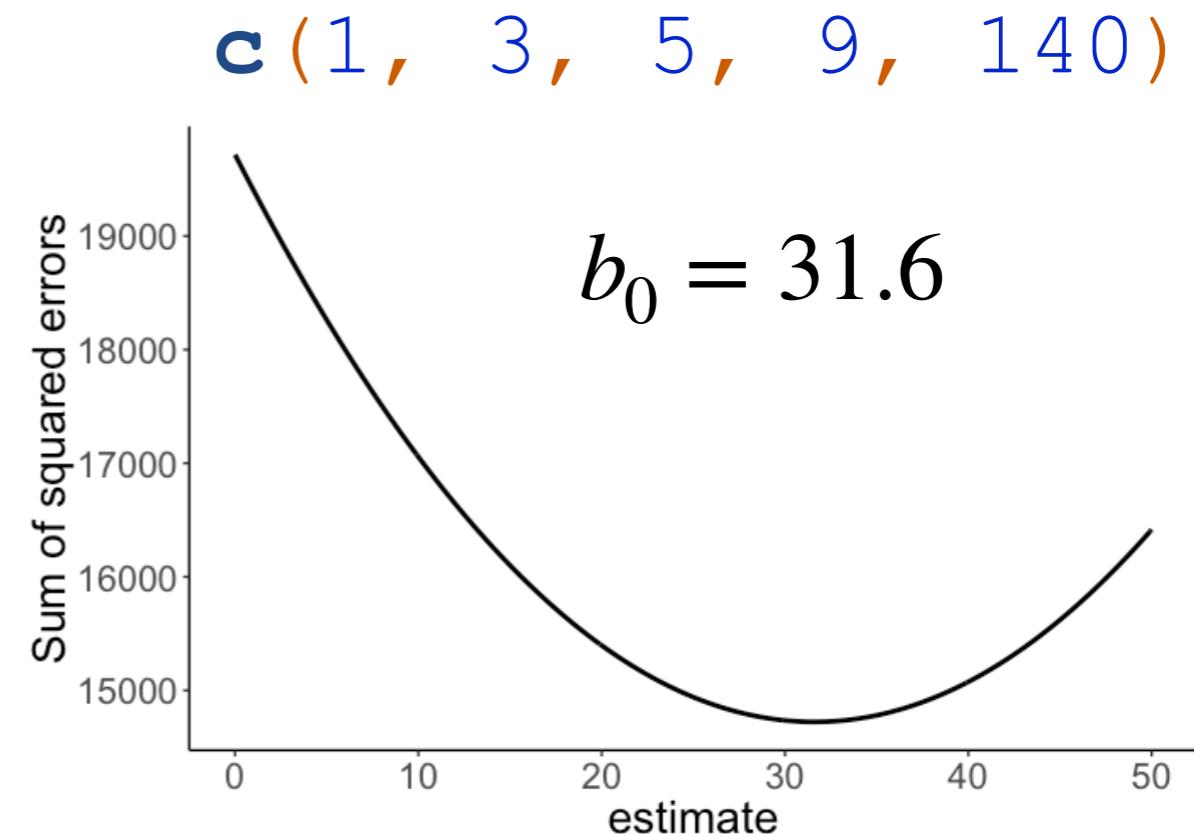
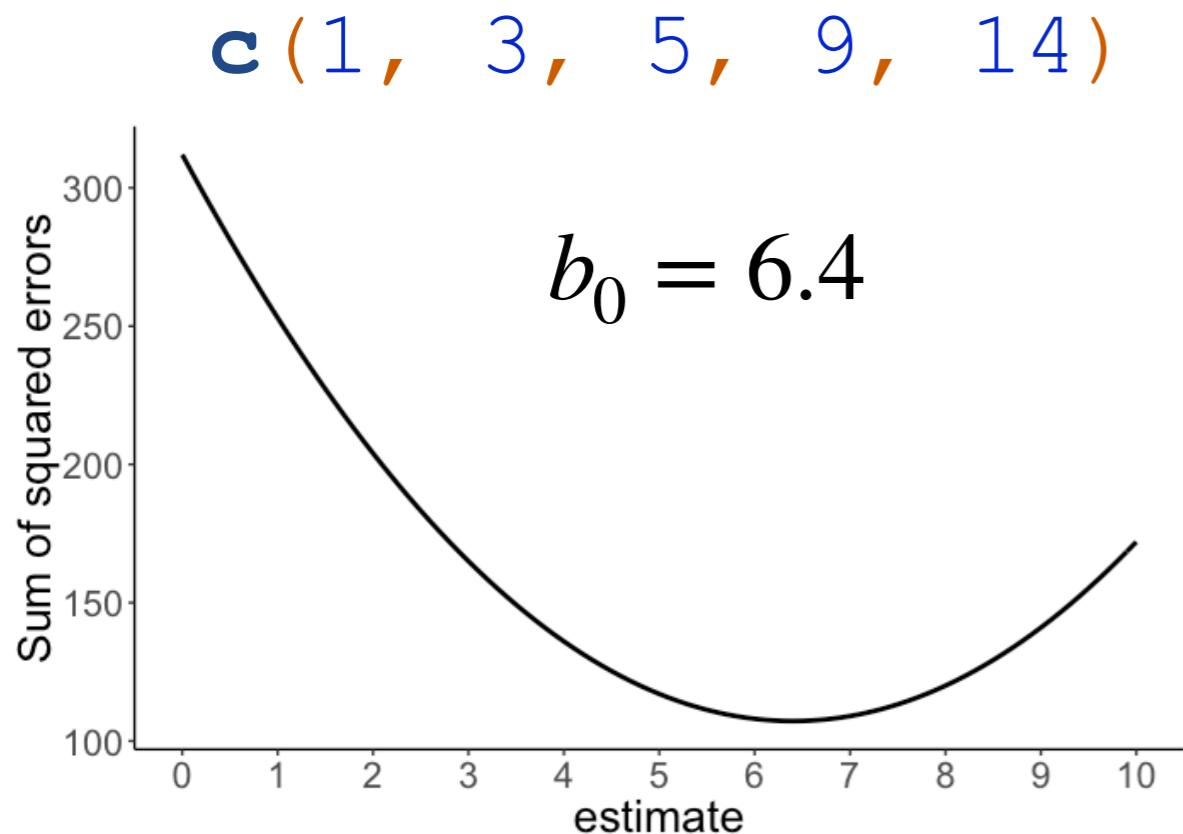
is robust to outliers!

Error = sum of squared errors

Sum of squared errors

$$Y_i = b_0 + e_i$$

$$\text{ERROR} = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (Y_i - b_0)^2$$



the **mean** minimizes the sum of squared errors

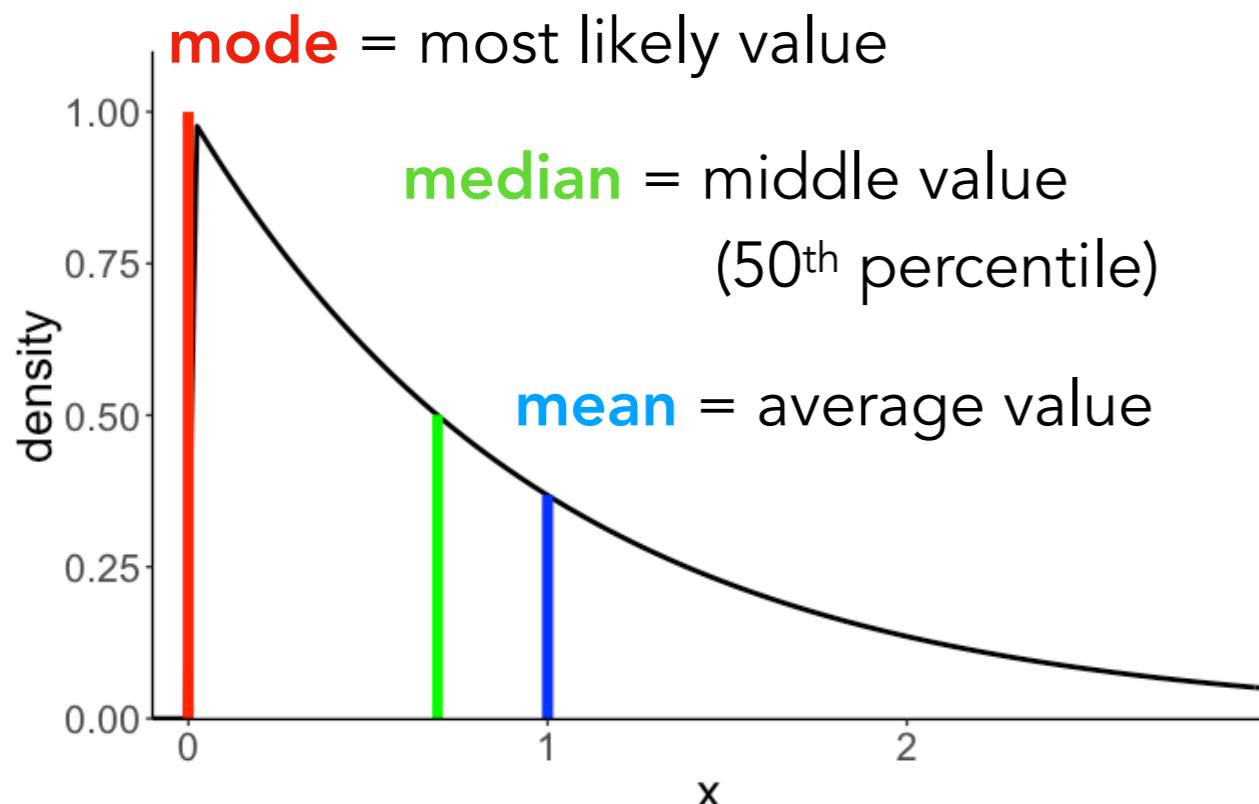
is strongly affected by outliers!

Error = count of errors

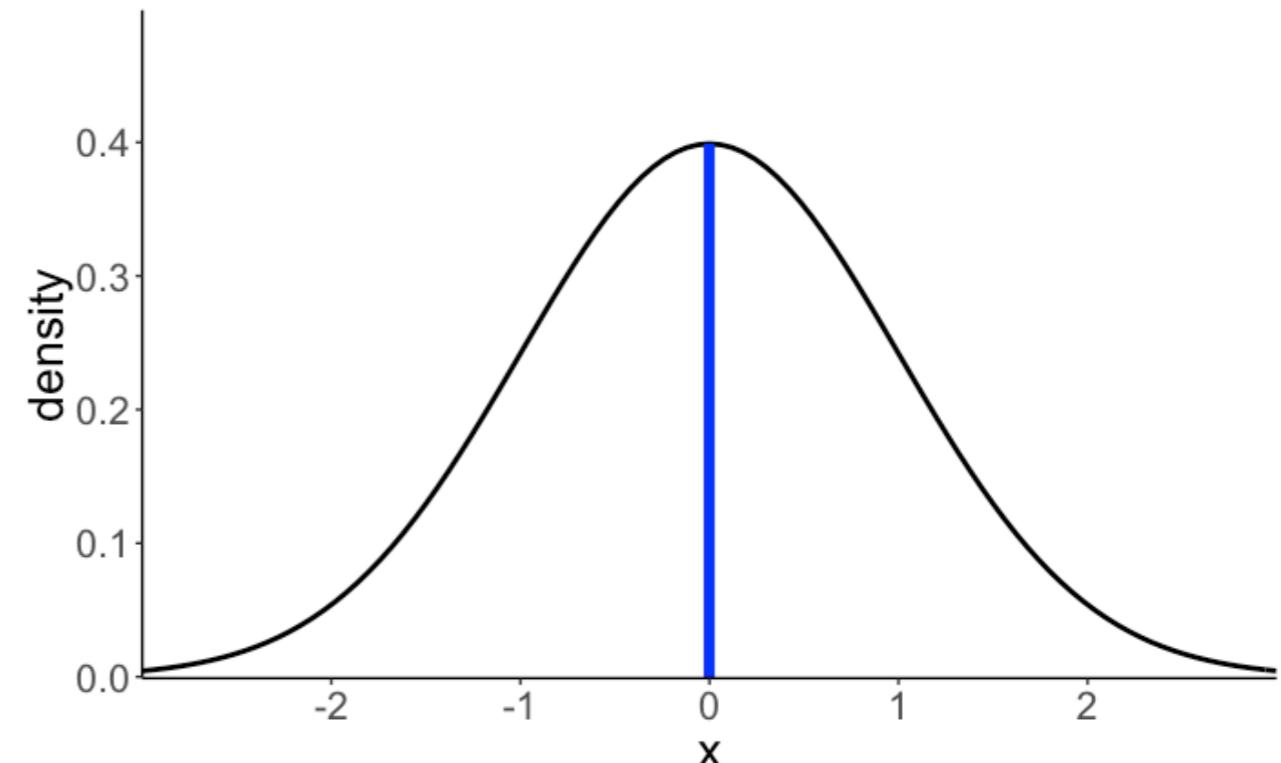
$$Y_i = b_0 + e_i \quad \text{ERROR} = \sum_{i=1}^n I(e_i) = \sum_{i=1}^n I(Y_i - b_0)$$

the **mode** minimizes the count of errors

Quick recap



exponential distribution



normal distribution

Error definition	Best estimator
Count of errors	Mode = most frequent value
Sum of absolute errors	Median = middle observation of all values
Sum of squared errors	Mean = average of all values

Models of error

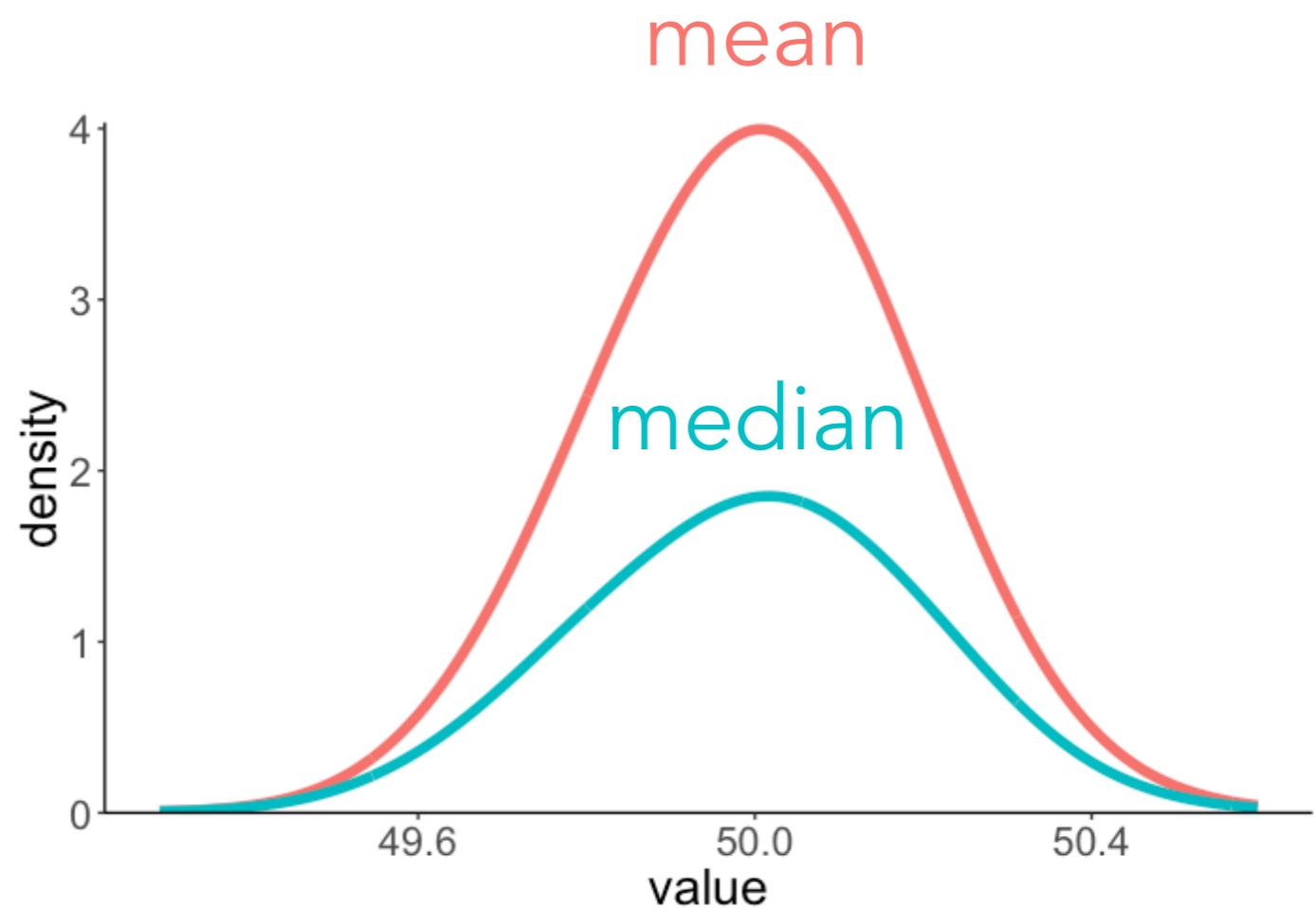
which model for error
shall we choose?

Sampling distributions

$$Y_i = 50 + \epsilon \text{ the true model}$$

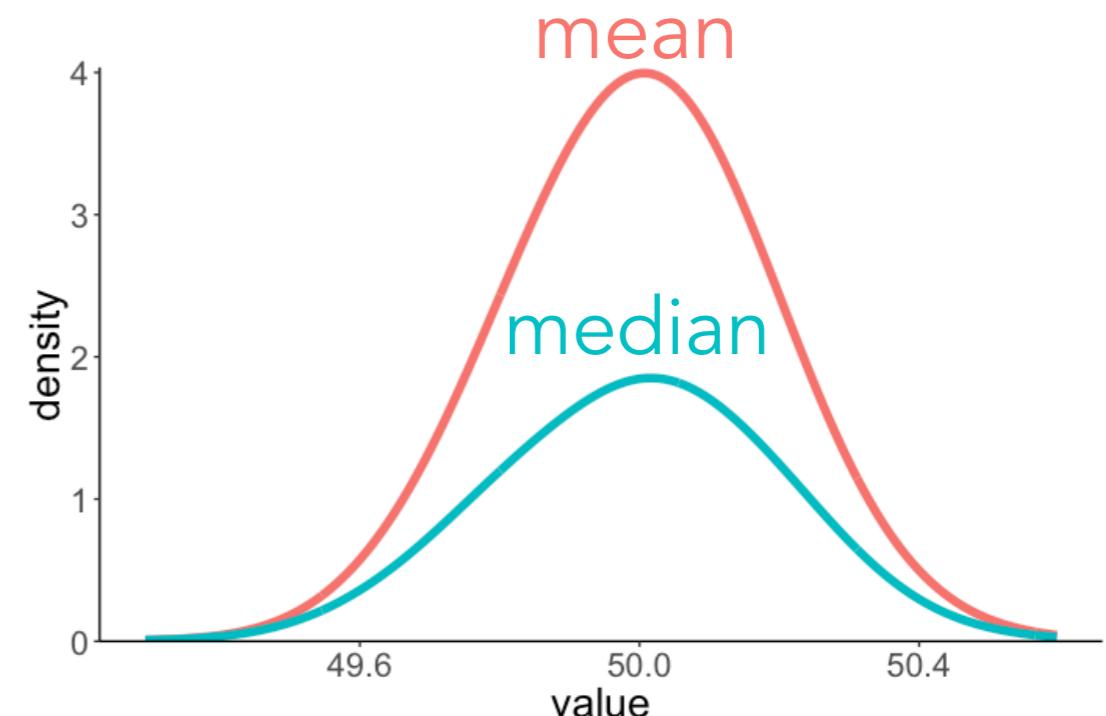
Recipe

- take m samples of size n
- for each sample, calculate the mean () and median ()
- plot the distribution (histogram, or density)



Properties of estimators

- **Unbiasedness**
 - does the average value of distribution match the true value?
- **Efficiency**
 - how precisely does the estimator capture the true value for a given sample size?
- **Consistency**
 - how does the estimators precision change as the sample size increases?



$$Y_i = 50 + \epsilon$$

$$\epsilon \sim \mathcal{N}(\mu = 0, \sigma)$$

but how was the error generated in the true model?

I assumed normally distributed errors!

justification?



The central limit theorem!

the distribution of the sum of individual error components will approximate a normal distribution

Quick recap

$$Y_i = \beta_0 + \epsilon_i$$

$$Y_i = b_0 + e_i$$

mean

sum of squared
errors

- Central limit theorem suggests that (very often) errors are normally distributed
- the mean is the *most efficient* (and unbiased) estimator when errors are normally distributed
- the mean minimizes the *sum of squared errors*

Inferring the error variance



$$Y_i = \beta_0 + \epsilon_i \quad \epsilon \sim \mathcal{N}(\mu = 0, \sigma)$$

Mean squared error (MSE) is an unbiased estimator of the population error variance.

variance

$$\text{MSE} = \frac{\text{SSE}}{n - 1} = \frac{\sum_{i=1}^n (Y_i - \bar{Y})^2}{n - 1}$$

standard deviation

$$\sqrt{\text{MSE}}$$

For more complex models with p parameters:

$$\text{MSE} = \frac{\text{SSE}}{n - p} = \frac{\sum_{i=1}^n (Y_i - \bar{Y})^2}{n - p}$$

Statistical inferences about parameter values

Procedure

1. Start with research question
2. Formulate hypothesis as a comparison between compact and augmented model
3. Fit parameters in each model
4. Calculate the proportional reduction of error (PRE)
5. Decide whether PRE is **worth it**

Research question and hypotheses

Is the average percentage of internet users per state significantly different from 75%?

Model_C: $Y_i = B_0 + \epsilon_i$

0 parameters

$$Y_i = 75 + e_i$$

Model_A: $Y_i = \beta_0 + \epsilon_i$

1 parameter

$$Y_i = b_0 + e_i$$

$$= \bar{Y} + e_i$$

i	internet	state	college	auto	density
1	79.0	AK	28.0	1.2	1.2
2	63.5	AL	23.5	1.3	94.4
3	60.9	AR	20.6	1.7	56.0
4	73.9	AZ	27.4	1.3	56.3
5	77.9	CA	31.0	0.8	239.1
6	79.4	CO	37.8	1.0	48.5
7	77.5	CT	37.2	1.0	738.1
8	74.5	DE	29.8	1.1	460.8
9	74.3	FL	27.2	1.2	350.6
10	72.2	GA	28.3	1.1	168.4

Fit parameters and calculate PRE

$$C: Y_i = 75 + e_i \quad A: Y_i = \bar{Y} + e_i$$

i	state	internet	compact_b	compact_se	augmented_b	augmented_se
1	AK	79.0	75	16.00	72.81	38.37
2	AL	63.5	75	132.25	72.81	86.60
3	AR	60.9	75	198.81	72.81	141.75
4	AZ	73.9	75	1.21	72.81	1.20
5	CA	77.9	75	8.41	72.81	25.95
6	CO	79.4	75	19.36	72.81	43.48
7	CT	77.5	75	6.25	72.81	22.03
8	DE	74.5	75	0.25	72.81	2.87
9	FL	74.3	75	0.49	72.81	2.23
10	GA	72.2	75	7.84	72.81	0.37

$$\begin{aligned} \text{PRE} &= 1 - \frac{\text{ERROR}(A)}{\text{ERROR}(C)} \\ &= 1 - \frac{1355}{1595} \approx .15 \end{aligned}$$

Model A has
15% less error
than Model C.

$$\text{SSE}(C) = 1595 \quad \text{SSE}(A) = 1355$$

Decide whether it's **worth it**

- PRE is the estimate of an unknown true reduction of error η^2
- we need a sampling distribution of PRE
 - a distribution of what PRE would look like if Model C (our H_0) were true
 - we could just simulate such a sampling distribution ...
- PRE is closely related to the F statistic!

Decide whether it's **worth it**

- To compute the F statistic, we need:

- PRE
- number of parameters in Model C (PC) and Model A (PA)
- number of observations n

- more likely to be **worth it** if:
 1. PRE is high
 2. the number of additional parameters in A compared to C is low (**PRE per additional parameter**)
 3. the number of parameters that could have been added to model_C to create model_A but were not

**difference in parameters
between models A and C**

$$F = \frac{\text{PRE}/(\text{PA} - \text{PC})}{(1 - \text{PRE})/(n - \text{PA})}$$

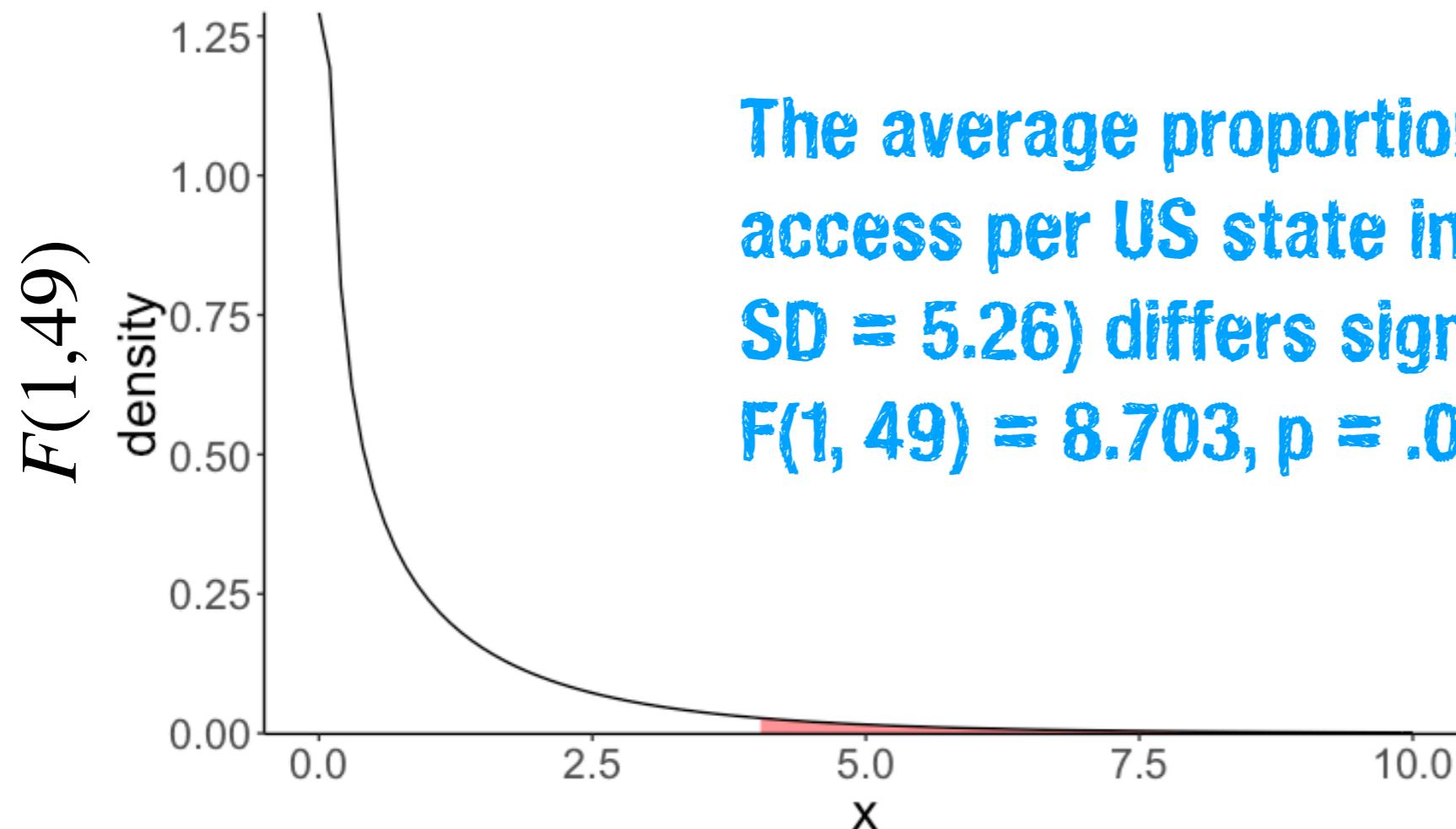
**number of observations
vs. parameters in Model A**

Decide whether it's **worth it**

$$F = \frac{\text{PRE}/(\text{PA} - \text{PC})}{(1 - \text{PRE})/(n - \text{PA})}$$

$$= \frac{.15/(1 - 0)}{(1 - .15)/(50 - 1)} = 8.703 \quad p = .00486$$

Note: I've used the approximated PRE here (PRE = 0.15), but I'm reporting the F value for the non-approximated PRE value.



we just performed a one sample t-test ...

```
t.test(df.internet$internet, mu = 75)
```

One Sample t-test

```
data: df.internet$internet  
t = -2.9502, df = 49, p-value = 0.00486  
alternative hypothesis: true mean is not equal to 75
```

but ...

- we will apply the general procedure we went through to day to a large range of different situations
- thinking of hypothesis testing as model comparison is (hopefully more) intuitive
- this approach still works even if we can't find a recipe in our favorite stats cook book

Thank you!