

Survey and Analysis

INFR11158/11230 Usable Security and Privacy

Dr. Jingjie Li

29/01/2026



THE UNIVERSITY
of EDINBURGH

What is different about security

- Large **information asymmetry** between participant and researcher
 - The researcher likely understand security of their tool
 - Participant likely doesn't even know that security problem exists
- **Deception** studies are common
 - You told the participant to accomplish task A, but you are really looking to see if they do B activity

Why Johnny Can't Encrypt

Why Johnny Can't Encrypt: A Usability Evaluation of PGP 5.0

Alma Whitten
*School of Computer Science
Carnegie Mellon University
Pittsburgh, PA 15213
alma@cs.cmu.edu*

J. D. Tygar¹
*EECS and SIMS
University of California
Berkeley, CA 94720
tygar@cs.berkeley.edu*

Abstract

User errors cause or contribute to most computer security failures, yet user interfaces for security still tend to be clumsy, confusing, or near-nonexistent. Is this simply due to a failure to apply standard user interface design techniques to security? We argue that, on the contrary, effective security requires a different usability standard, and that it will not be achieved through the user interface design techniques appropriate to other types of consumer software.

To test this hypothesis, we performed a case study of a security program which does have a good user interface by general standards: PGP 5.0. Our case study used a cognitive walkthrough analysis together with a laboratory user test to evaluate whether PGP 5.0 can be successfully used by cryptography novices to achieve effective electronic mail security. The analysis found a number of user interface design flaws that may

1 Introduction

Security mechanisms are only effective when used correctly. Strong cryptography, provably correct protocols, and bug-free code will not provide security if the people who use the software forget to click on the encrypt button when they need privacy, give up on a communication protocol because they are too confused about which cryptographic keys they need to use, or accidentally configure their access control mechanisms to make their private data world-readable. Problems such as these are already quite serious: at least one researcher [2] has claimed that configuration errors are the probable cause of more than 90% of all computer security failures. Since average citizens are now increasingly encouraged to make use of networked computers for private transactions, the need to make security manageable for even untrained users has become critical [4, 9].

Why End

If an average user of email feels the need for privacy and authentication, and acquires PGP with that purpose in mind, will PGP's current design allow that person to realize what needs to be done, figure out how to do it, and avoid dangerous errors, without becoming so frustrated that he or she decides to give up on using PGP after all?

interface by general standards: PGP 5.0. Our case study used a cognitive walkthrough analysis together with a laboratory user test to evaluate whether PGP 5.0 can be successfully used by cryptography novices to achieve effective electronic mail security. The analysis found a number of user interface design flaws that may

the probable cause of more than 90% of all computer security failures. Since average citizens are now increasingly encouraged to make use of networked computers for private transactions, the need to make security manageable for even untrained users has become critical [4, 9].

ive when used
rovably correct
vide security if
to click on the
y, give up on a
re too confused
need to use, or
rol mechanisms
able. Problems
s: at least one
ration errors are

Users need to:

- understand that privacy is achieved by encryption, and figure out how to encrypt email and how to decrypt email received from other people
- understand that authentication is achieved through digital signatures, and figure out how to sign email and how to verify signatures on email from other people
- understand that in order to sign email and allow other people to send them encrypted email a key pair must be generated, and figure out how to do so
- understand that in order to allow other people to verify their signature and to send them encrypted email, they must publish their public key, and figure out some way to do so
- understand that in order to verify signatures on email from other people and send encrypted email to other people, they must acquire those people's public keys
- manage to avoid such dangerous errors as accidentally failing to encrypt, trusting the wrong public keys, failing to back up their private keys, and forgetting their pass phrases
- be able to succeed at all of the above within a few hours of reasonably motivated effort

Tested usability using two methods

- Cognitive Walkthrough
 - A set of experts review and the experts make an informed guess about what will be problematic
 - Paired with heuristics — The experts state how the user interface supports or violates common HCI principles (Heuristics)
- Lab Study
 - Ask the participant to perform a set of tasks
 - Very similar to a think aloud, but without the talking aloud part

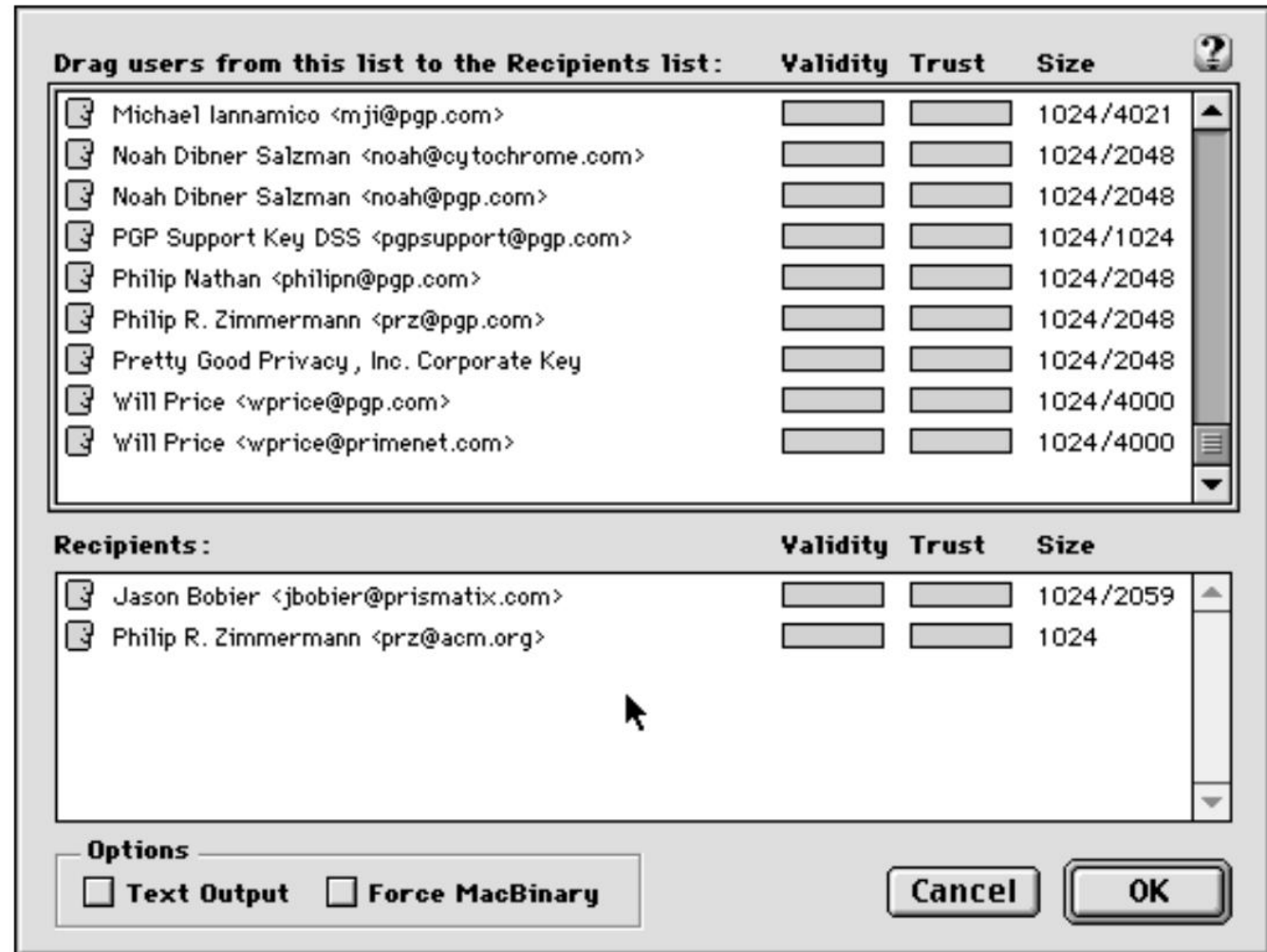
Cognitive walkthrough outcomes

- **Visual metaphors** — Do key and lock pictures make sense?
- **Different key types** — Public vs private keys, or maybe signing and encryption keys?
- **Key server** — Used for sharing keys
- **Key management policy** — Trust and validity ratings
- **Consistency** — Use of the same terms everywhere
- **Too much information** — Information like key size, hashes, and trust
- **Irreversible actions**
 - Accidentally deleting the private key
 - Accidentally publicizing a key
 - Accidentally revoking a key
 - Forgetting the pass phrase
 - Failing to back up the key rings

Lab study

- 12 participants with CS backgrounds
- Participant had to send several emails to team members (the researchers)
 - Creating a key pair
 - Sending their public key to team members
 - Getting team members' public keys
 - Sending the email
 - Decrypting response email
- 3 — emailed the private key to the team member
 - 1 never realized the error
- 1 — forgot their pass phase and had to re-generate keys
- 1 — never figured out how to encrypt
- 7 — used their public keys to encrypt
 - 1 created a separate key pair for each team member
- 3 — successfully sent an encrypted email to the whole team and were able to decrypt an response email

Whitten and Tygar evaluated PGP encryption in 1999, surely it must be more usable now.



A personal story during my PhD

Kaleido: Real-Time Privacy Control for Eye-Tracking Systems

Jingjie Li, Amrita Roy Chowdhury, Kassem Fawaz, and Younghyun Kim

University of Wisconsin–Madison

{jingjie.li, roychowdhur2, kfawaz, younghyun.kim}@wisc.edu

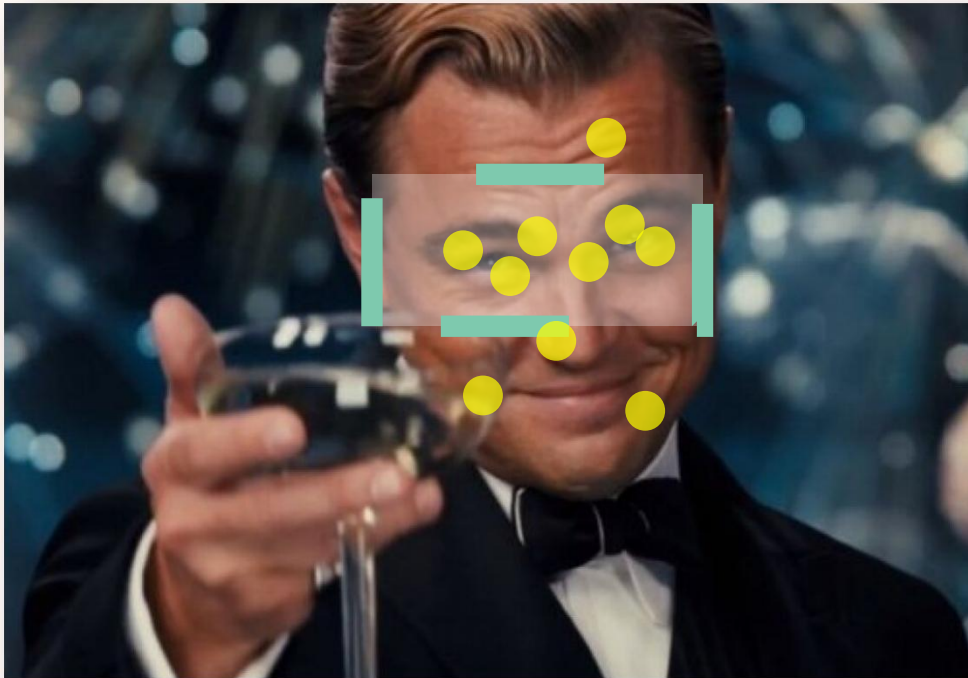
Abstract

Recent advances in sensing and computing technologies have led to the rise of eye-tracking platforms. Ranging from mobiles to high-end mixed reality headsets, a wide spectrum of interactive systems now employs eye-tracking. However, eye gaze data is a rich source of sensitive information that can reveal an individual's physiological and psychological traits. Prior approaches to protecting eye-tracking data suffer from two major drawbacks: they are either incompatible with the current eye-tracking ecosystem or provide no formal



Figure 1: Eye gaze heatmaps from an individual user with and without Kaleido's noising effect on a web page.

Privacy Implications of Eye Tracking



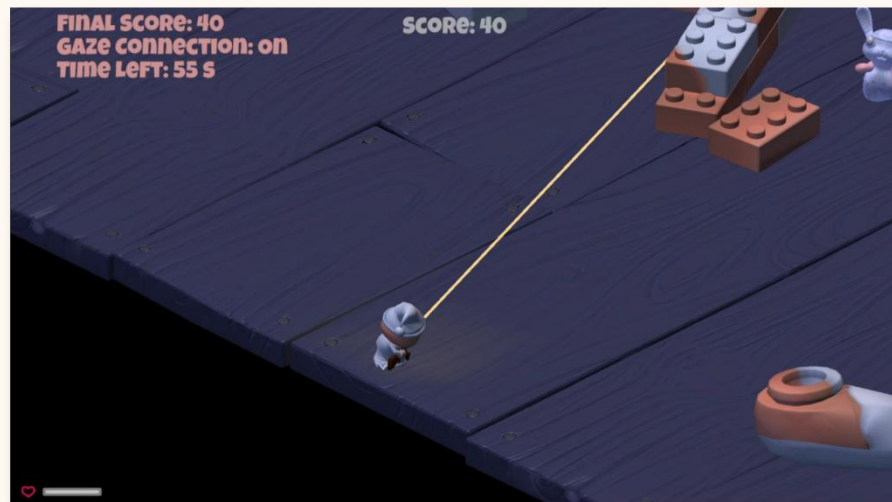
Eye gazes from people with **low social anxiety**



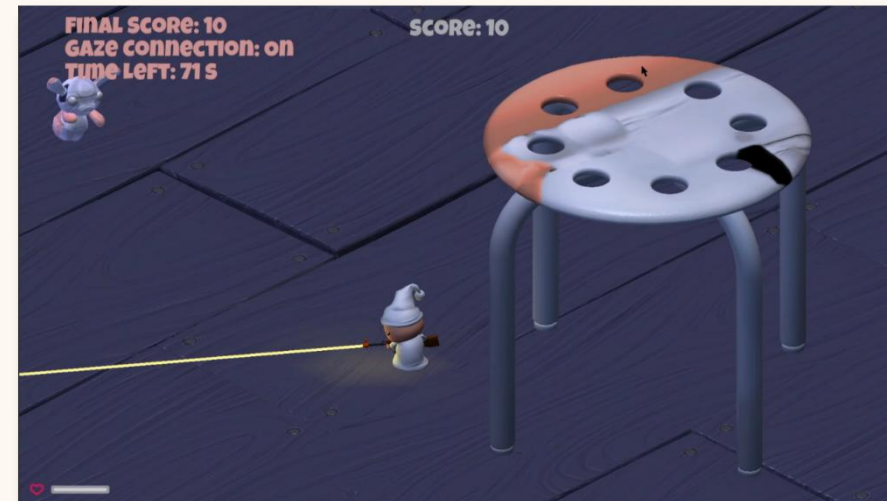
Eye gazes from people with **high social anxiety**

Some background

- Test out whether/how user experience is impacted by a privacy control we designed in an eye tracking game setting



Kalido off



Kalido on

How to do user studies?



My original plan in
03/2020 to do lab
studies with a VR
setup

Then...Guess what?

Questionnaires

Questionnaires

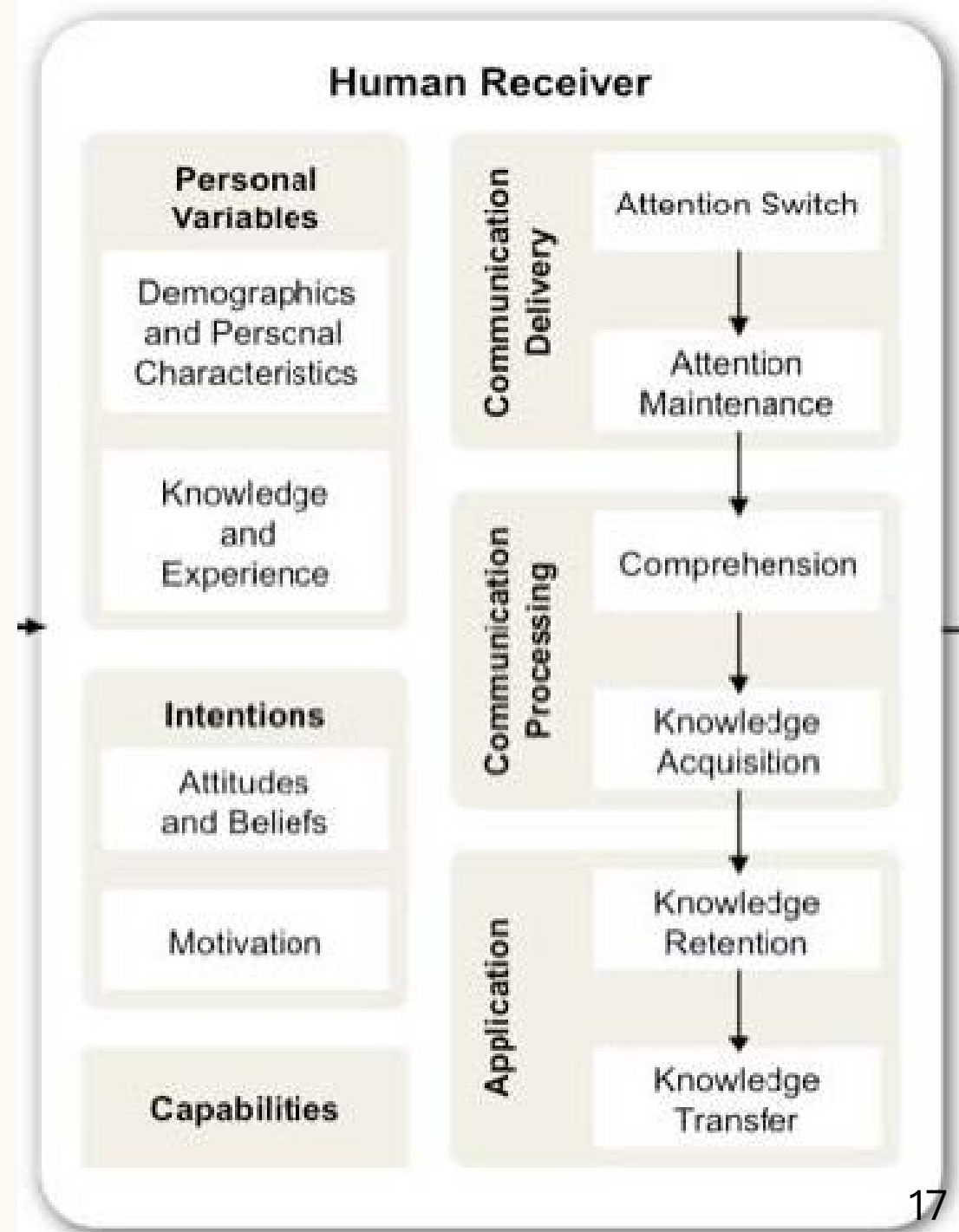
- Ask participants to answer a set of pre-defined questions.
- Pros:
 - gather data from a large number of people quickly
 - can determine how prevalent an issue or concern is
 - close-ended questions are easy to analyze
- Cons:
 - can only gather data you know about
 - careful planning is required before running a questionnaire
 - open-ended questions can take a lot of time to analyze and require careful setup

Questionnaires can be used at various points in the design process

- Understanding people
 - Understand the target population
 - Incorrect mental models
- Testing a theory
 - Are my assumptions correct?
 - Do people think that $A=B$?
- Testing a prototype design
 - How do people interpret my interface?
- Testing the final design
 - How are people actually using it?
 - What do people think after they use it?

What do you want to know?

- Attitudes
 - Do you like X?
 - Would using X work?
- Behaviors
 - How often do you use X?
 - Do you regularly do X?
- Knowledge
 - What is the best definition of X?
- Expectations
 - If the webpage did X what would you expect to happen?
- Capabilities
 - What is the result of adding 20 and 30?



Common survey elements

- Single and multiple choice checkboxes
- Matching
 - Rank the following from 1 to 5
- Rating scales
 - Likert Scales
 - 3, 5, 7 points scales
 - Semantic scales
- Open ended responses

OPEN ENDED

- Where does this URL go? What does it do?

Easier to write, harder to analyze

Harder to write, easier to analyze

CLOSE-ENDED

If you clicked on the link above, what web page would open?

- ☐ WWW3's main page
- ☐ National Geographic's main page
- ☐ World News's main page
- ☐ I will be taken to one of the sites above, but not their main page
- ☐ I will be taken to a website not listed above
- ☐ Other _____

Response Anchors

Psychologists have been working for quite some time to determine the least biased way to present a set of answers.

On the right are a set of response anchors that are known to work well.

Likert-Type Scale Response Anchors

Citation:

Vagias, Wade M. (2006). *Likert-type scale response anchors*. Clemson International Institute for Tourism & Research Development, Department of Parks, Recreation and Tourism Management. Clemson University.

Level of Acceptability

- 1 – Totally unacceptable
- 2 – Unacceptable
- 3 – Slightly unacceptable
- 4 – Neutral
- 5 – Slightly acceptable
- 6 – Acceptable
- 7 – Perfectly Acceptable

Level of Appropriateness

- 1 – Absolutely inappropriate
- 2 – Inappropriate

My beliefs

- 1 – Very untrue of what I believe
- 2 – Untrue of what I believe
- 3 – Somewhat untrue of what I believe
- 4 – Neutral
- 5 – Somewhat true of what I believe
- 6 – True of what I believe
- 7 – Very true of what I believe

Priority:

Level of Support/Opposition

- 1 – Strongly oppose
- 2 – Somewhat oppose
- 3 – neutral
- 4 – Somewhat favor
- 5 – Strongly favor

Level of Probability

- 1 – Not probable
- 2 – Somewhat improbable
- 3 – Neutral
- 4 – Somewhat probable

Level of Acceptability

- 1 – Totally unacceptable
- 2 – Unacceptable
- 3 – Slightly unacceptable
- 4 – Neutral
- 5 – Slightly acceptable
- 6 – Acceptable
- 7 – Perfectly Acceptable

• 3 – Sometimes but

Affect on X

Frequency

Fill in the blank
type question

Q2: What is your age? _____

Typical
multiple choice
question

Q8: What is the highest level of education you have achieved?

- ☐ High school or less
- ☐ Some College
- ☐ Bachelor's Degree
- ☐ Master's Degree
- ☐ Doctorate Degree

Scale where
multiple
questions are
meant to be
summed
together

Q12: To what extent do you agree or disagree with each of the following statement

Please select one answer per row

	Strongly Disagree	Disagree	Neither Agree or Disagree	Agree	Strongly Agree
I often ask others for help with the computer	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Others often ask me for help with the computer	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Likert scale
question using a
pre-defined
anchor

Q13: In terms of your Internet skills, do you consider yourself to be:

- ☐ Not at all skilled
- ☐ Not very skilled
- ☐ Fairly skilled
- ☐ Very skilled
- ☐ Expert

System Usability Scale Questionnaire

**Strongly
Disagree**

**Strongly
Agree**

1. I think that I would like to use this product frequently.

1	2	3	4	5
---	---	---	---	---

2. I found the product unnecessarily complex.

1	2	3	4	5
---	---	---	---	---

3. I thought the product was easy to use.

1	2	3	4	5
---	---	---	---	---

4. I think that I would need the support of a technical person to be able to use this product.

1	2	3	4	5
---	---	---	---	---

5. I found the various functions in the product were well integrated.

1	2	3	4	5
---	---	---	---	---

6. I thought there was too much inconsistency in this product.

1	2	3	4	5
---	---	---	---	---

7. I imagine that most people would learn to use this product very quickly.

1	2	3	4	5
---	---	---	---	---

8. I found the product very awkward to use.

1	2	3	4	5
---	---	---	---	---

9. I felt very confident using the product.

1	2	3	4	5
---	---	---	---	---

10. I needed to learn a lot of things before I could get going with this product.

1	2	3	4	5
---	---	---	---	---

System Usability Scale

#	Question	N/A	μ	σ
A1	<i>I apply software updates as soon as my computer prompts me.</i>	5 (1.0%)	3.20	1.221
A2	<i>I am happy to use an older version of a program, as long as it meets my needs.</i>	5 (1.0%)	^r 1.99	1.000
A3	<i>Whenever I step away from my computer, I lock the screen.</i>	5 (1.0%)	2.50	1.306
A4	<i>Others can access my smartphone or tablet without needing a PIN or passcode.</i>	21 (4.4%)	^r 3.34	1.545
A5	<i>When I discover a computer security problem at work, I'm likely to promptly report it to my employer.</i>	64 (13.4%)	4.08	0.995
A6	<i>It's important to use a WiFi password to prevent unauthorized people from using my home network.</i>	11 (2.3%)	4.66	0.690
A7	<i>I frequently click links in email messages to see what they are, regardless of who sent the message.</i>	5 (1.0%)	^r 4.51	0.922
A8	<i>It's important to run anti-virus software on my computer.</i>	7 (1.5%)	4.35	0.941
A9	<i>When browsing websites, I frequently mouseover links to see where they go, before clicking them.</i>	4 (0.8%)	4.13	0.977
A10	<i>When using public WiFi, I visit the same websites that I would visit when using the Internet at home.</i>	20 (4.2%)	^r 2.93	1.266
A11	<i>I usually do not pay attention to where I'm downloading software from.</i>	2 (0.4%)	^r 4.38	0.900
A12	<i>I frequently backup my computer.</i>	5 (1.0%)	3.07	1.165
A13	<i>I frequently visit websites even when my web browser warns me against it.</i>	8 (1.7%)	^r 3.98	1.028
A14	<i>I circumvent my employer's computer usage policies when they prevent me from completing a task.</i>	86 (18.0%)	^r 3.54	1.184
A15	<i>I am careful to never share confidential documents stored on my home or work computers.</i>	15 (3.1%)	4.36	0.757
A16	<i>Frequently checking the access control settings on social networking websites isn't worth the time it takes.</i>	18 (3.8%)	^r 3.56	1.165
A17	<i>I always write down my passwords to help me remember them.</i>	6 (1.3%)	^r 3.60	1.313
A18	<i>Creating strong passwords is not usually worth the effort.</i>	6 (1.3%)	^r 4.05	1.047
A19	<i>I frequently check my financial accounts for fraudulent charges.</i>	10 (2.1%)	4.11	0.914
A20	<i>If I receive a suspicious email from a company that I do business with, I'll phone the company to make sure the email is accurate.</i>	22 (4.8%)	3.52	1.236
A21	<i>I never give out passwords over the phone.</i>	7 (1.5%)	4.53	0.787
A22	<i>I frequently purchase things that I see advertised in unsolicited emails.</i>	4 (8.8%)	^r 4.51	0.840
A23	<i>I tend to ignore computer security stories in the news because they don't impact me.</i>	4 (8.8%)	^r 3.83	1.050
A24	<i>I use encryption software to secure files or email messages.</i>	10 (2.1%)	2.74	1.225
A25	<i>Once I create a password, I tend to never change it.</i>	5 (1.0%)	^r 3.30	1.182
A26	<i>I try to create a unique password for every account I have.</i>	5 (1.0%)	3.21	1.284
A27	<i>Rather than logging out of websites, I usually just navigate elsewhere or close the window when I'm done.</i>	7 (1.5%)	^r 3.06	1.299
A28	<i>I always make sure that I'm at a secure website (e.g., SSL, "https://", a lock icon) when transmitting information online.</i>	4 (0.8%)	3.80	1.173
A29	<i>I frequently use privacy software, "private browsing" or "incognito" mode when I'm online.</i>	6 (1.3%)	3.17	1.247
A30	<i>I frequently let others use my computing devices (e.g., smartphone, tablet, laptop).</i>	3 (0.7%)	^r 3.79	1.172

Table 1. Initial set of security questions evaluated on a 5-point Likert scale (from "strongly disagree" to "strongly agree") by 479 participants. Depicted are the questions, the rate of "N/A" responses, and the average responses and standard deviations after recoding negatively-phrased questions (^r).

Scaling the Security Wall: Developing a Security Behavior Intentions Scale (SeBIS)

Planning a survey

Don't panic! This is not a statistics class.

COULD BE ON THE EXAM

- Independent and dependent variables
- Correlation vs causation
- Between vs within subject design
- Study question design

WILL NOT BE ON THE EXAM

- Statistical test names
 - T-test, ANOVA, etc.
- When to use different tests
 - Chi Sq should be used with categorical dependent and independent variables
- P-values, distributions, confidence intervals or other outcomes from tests

Topics Outline

- **Descriptive questions vs testing a question**
- Correlation vs causation
- Dependent vs independent variables
- Between and within subjects testing
- Numeric vs categorical data

Planning a survey

- Surveys normally answer **multiple research questions**. With each research question tied to one or more survey questions.
- **Descriptive** — learn something about the whole population.
 - How many people have heard of the term “phishing”?
 - What words do people use to describe cookie tracking?
- **Testing for correlation or causation** — show that two things are related or one thing causes the other thing.
 - If someone has been trained on phishing in the past, are they better at differentiating phishing emails?
 - We have three training options, each user goes through one training, which training causes people to identify phishing emails the best?

Descriptive Statistics

- **Descriptive Questions** — learn something about the **whole population**.
 - How many people have heard of the term “phishing”?
 - What words do people use to describe cookie tracking?
- **Descriptive Numeric** — fancy term for all the basic measures of numeric data: **Mean, median, mode, standard deviation**
 - What % of consumers are worried about privacy?
 - What % of people know the difference between behavioral advertising and cookies?
 - On average, how long does it take to decide if an email is phishing or not?
- **Descriptive Qualitative** — use data to learn about a whole population
 - What is the most common reason people avoid using ATMs?
 - Why do some people choose to not have a Google account?

Testing for correlation or causation

- Testing **for correlation or causation** — show that two things are related, or one thing causes the other thing.
 - If someone has been trained on phishing in the past, are they better at differentiating phishing emails?
 - We have three training options, each user goes through one training, which training **causes** people to identify phishing emails the best?
- These tests require more complex statistics, such as:
 - T-test
 - ANOVA
 - Linear Models
 - CHI Squared

Topics Outline

- Descriptive questions vs testing a question
- **Correlation vs causation**
- Dependent vs independent variables
- Between and within subjects testing
- Numeric vs categorical data

Correlation vs. Causation

- Correlation
 - Two things tend to behave in a way that seems inter-related, where if one thing changes the other thing will also change in a related way.
 - For example, if the price of rice goes up at the same time as the price for beans.
- Causation
 - When one thing changes it causes the other thing to change.
 - For example, when the weather gets cold more people wear coats.
Cold weather causes more people to wear coats.

Does consuming
chocolate
increase the
number of Nobel
Laureates?

This is a
correlation, not
necessarily a
causation.

Chocolate Consumption, Cognitive
Function,
and Nobel Laureates
Franz H. Messerli, M.D.

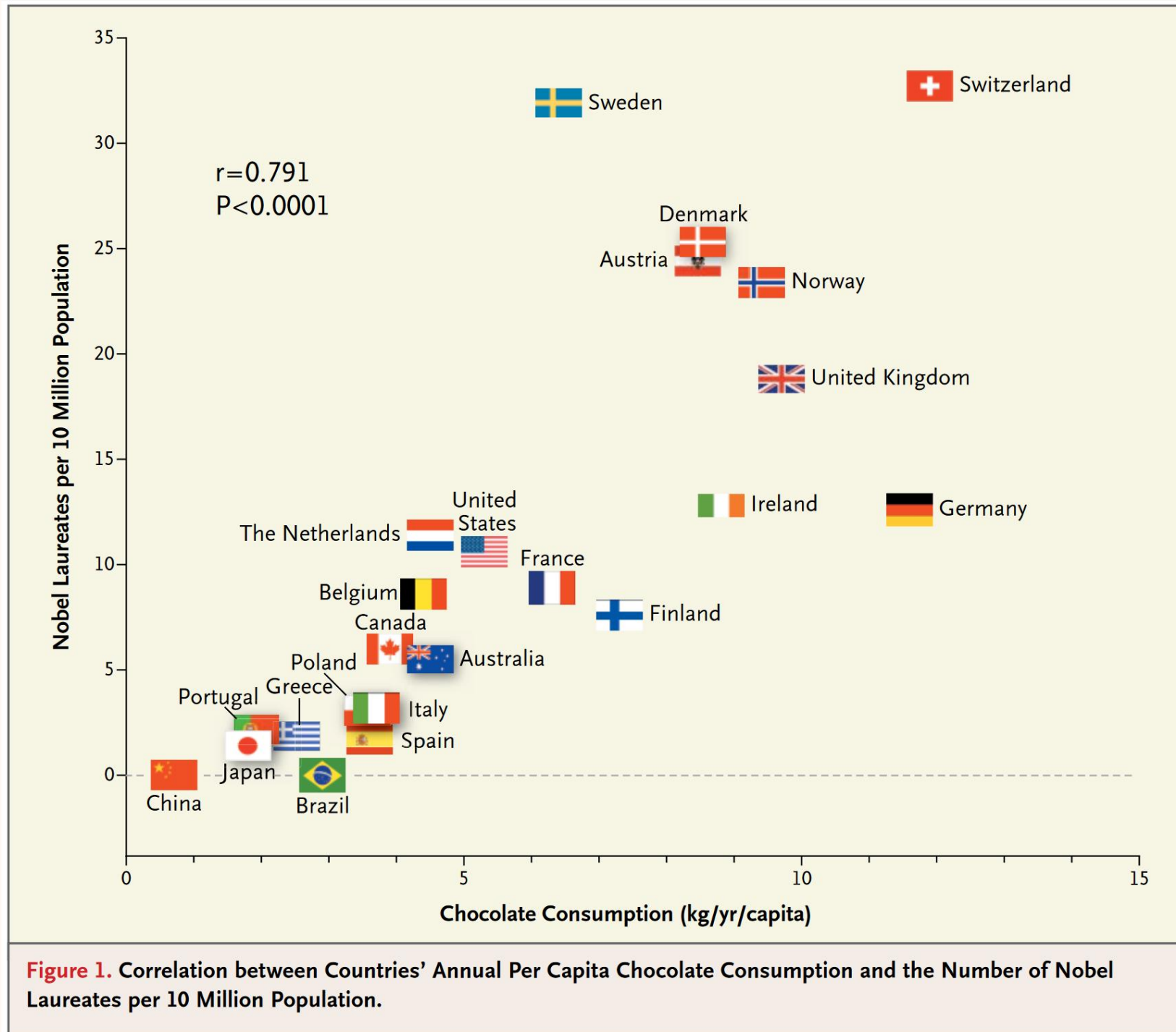
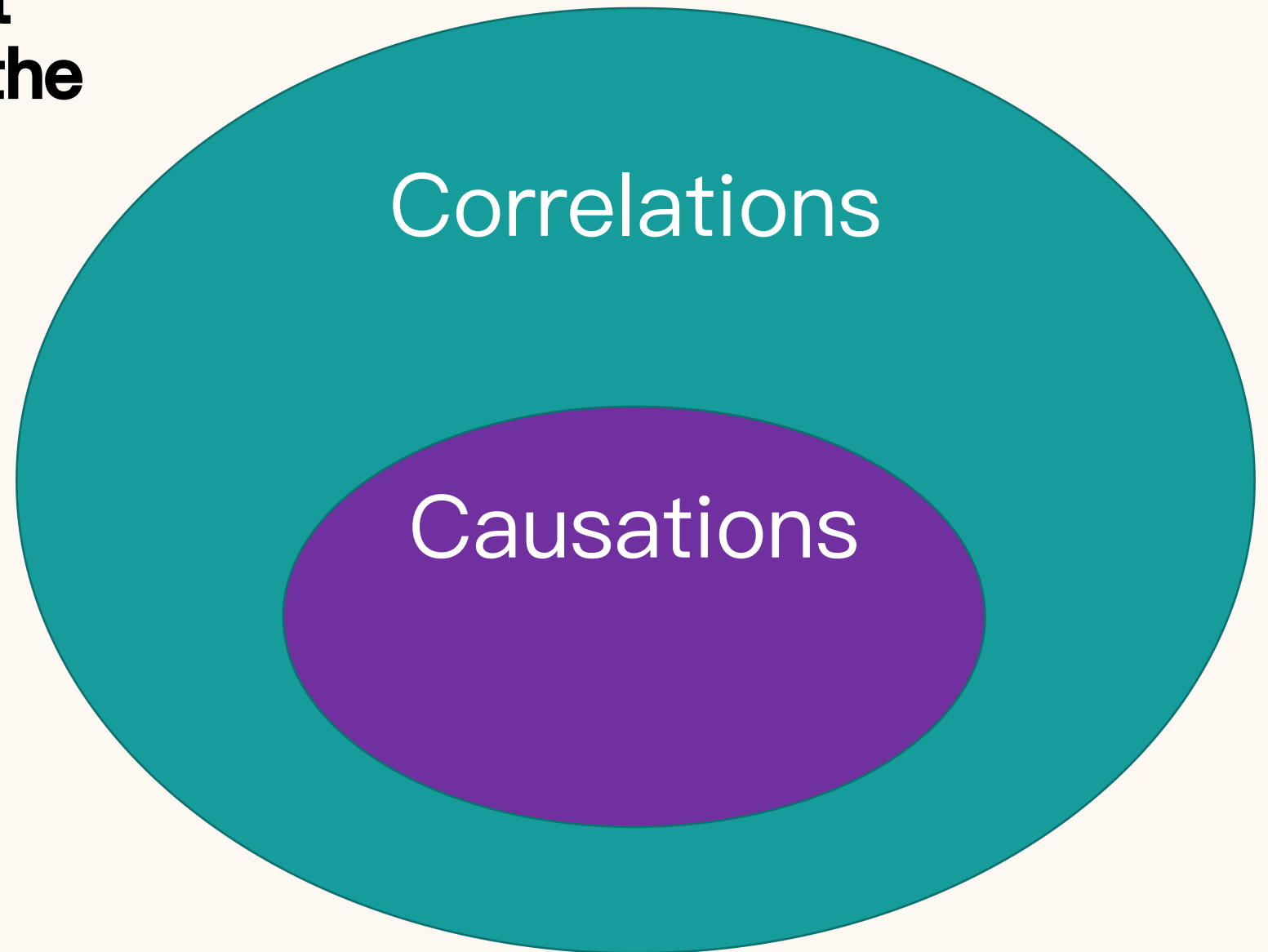


Figure 1. Correlation between Countries' Annual Per Capita Chocolate Consumption and the Number of Nobel Laureates per 10 Million Population.

**Causations are
Correlations, but
not necessarily the
other way round**



History +
CTR is a
correlation.

How might
you test if
it is really
a causation?

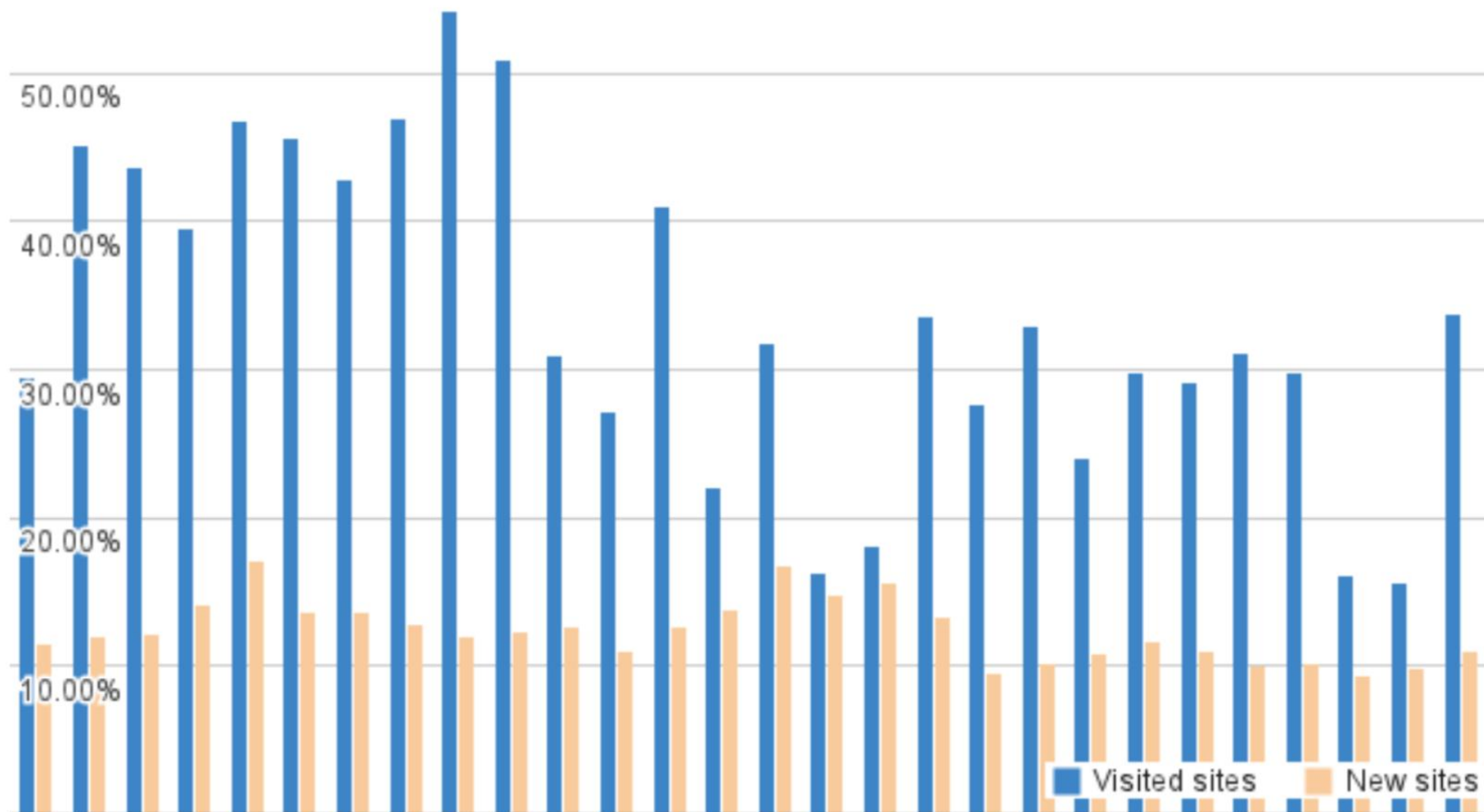


Figure 3: Daily CTR, separated by whether the website was already in the user's browsing history. For 28 days in January-February 2014.

Topics Outline

- Descriptive questions vs testing a question
- Correlation vs causation
- **Dependent vs independent variables**
- Between and within subjects testing
- Numeric vs categorical data

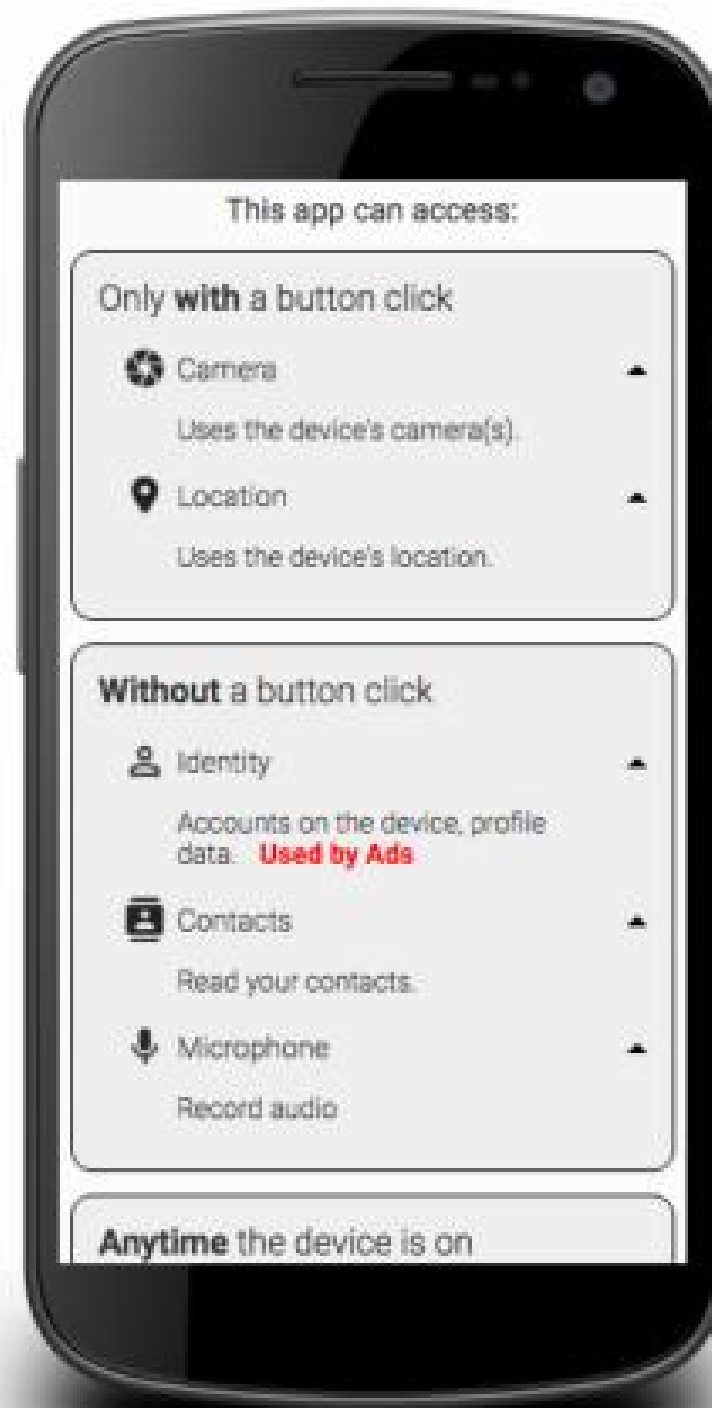
What are you going to measure?

- In statistics there are classically two types of measurements (variables): dependent and independent
- Dependent
 - Also known as the **outcome variable**
 - “Dependent” on the study
 - Measures the usability **goal**
- Independent
 - Anything **you are directly manipulating**
 - An element of the study which is under your control
 - A pre-existing feature of your participant

Some research questions:

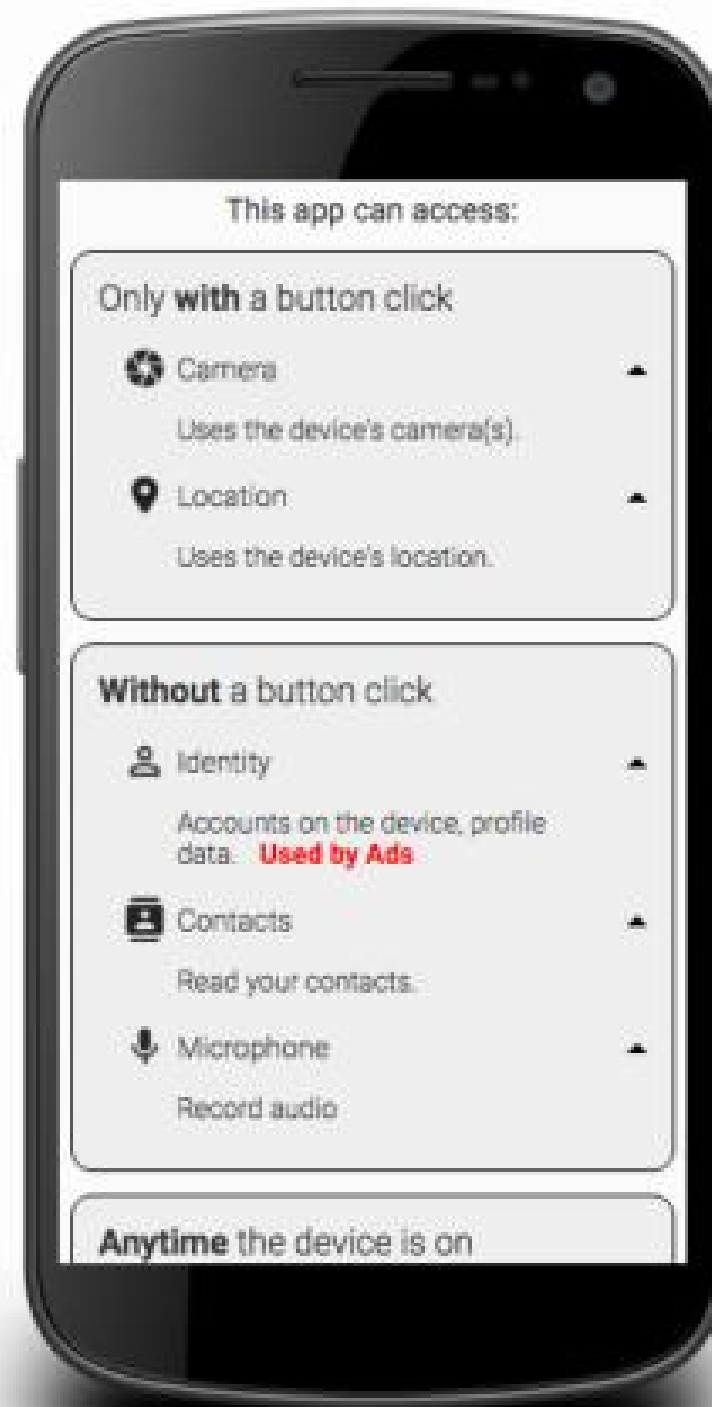
- Can people differentiate between a subdomain and a domain when reading a URL?
- Does [a new system] help people differentiate between malicious URLs and safe ones?
- Can users use [a new password manager] faster and with less errors than [the old password manager]?
- Does knowing how an app will use its permissions impact app installation decisions?
- Using [website], can users successfully opt-out of cookie tracking without forming inaccurate mental models?

**Lets use this study as
an example**



Research Question:

Can users reliably identify if an app can or cannot perform an action directly tied to a permission.





Awesome App

can access

- Location
Uses the device's location
- Camera
Uses the device's camera(s)

Dependent variable:

Count of the number of questions the participant answered correctly



Awesome App

can access

Without a button click

- Microphone
Record audio
- Camera
Uses the device's camera(s).
- Location
Uses the device's location. **Used by Ads**

What can this app do?

Independent variable:

Which of the two interfaces the participant was shown

Charge purchases to your credit card at any time.
Get your location.
Allow ads to know your location.
Load ads.
Write on the SD card

Absolutely
Possible

<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Variables that would make sense

- Research Question: Can users reliably identify if an app can or cannot perform an action directly tied to a permission?
- Dependent
 - Which permissions correctly/incorrectly read
 - Count of permissions correctly/incorrectly read
 - Time spent reading each permission screen
- Independent
 - Study group (which screen was shown)
 - If the permission was privacy sensitive or not
 - Order of the tasks
 - Time of day
 - Type of most used device (laptop, mobile, PC)
 - Demographics of the participants (gender, age, native language, ...)

Common dependent things to measure

- Number of dangerous errors made
- Time to complete task
- Percent of task completed
- Percent of task completed per unit of time
- Ratio of successes to failures
- Time spent in errors
- Percent or number of errors
- Percent or number of competitors better than it
- Frequency of help and documentation use

Topics Outline

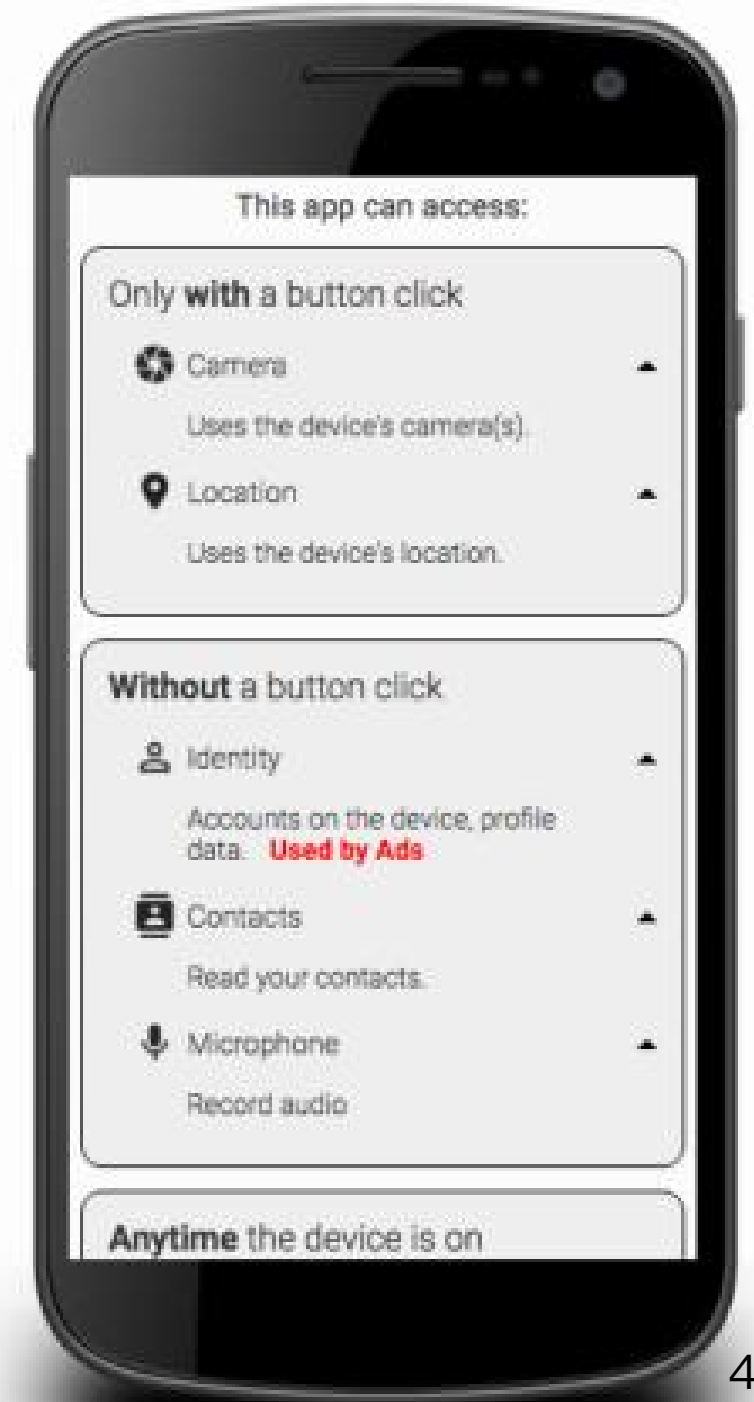
- Descriptive questions vs testing a question
- Correlation vs causation
- Dependent vs independent variables
- **Between and within subjects testing**
- Numeric vs categorical data

Between vs. Within subjects

- Between subjects
 - Your study only shows one interface to one person
 - You are measuring how well the people randomly assigned to the A interface did compared to the people randomly assigned to the B interface
 - **Lots of variability with this method**
- Within subjects
 - Your study shows all interfaces to all people
 - You are measuring the difference in how they do on the two interfaces
 - **Less variability (same person) but more learning effects and priming**

Study design

- RQ: Does [my new interface] enable people to accurately determine what permissions an app will use?
- A/B test between the existing and new interface
- Between subjects
- 10 Tasks shown in the same order to all participants
- Dependent variables
 - Accuracy on task
- Independent variables
 - Which interface (A or B)



Topics Outline

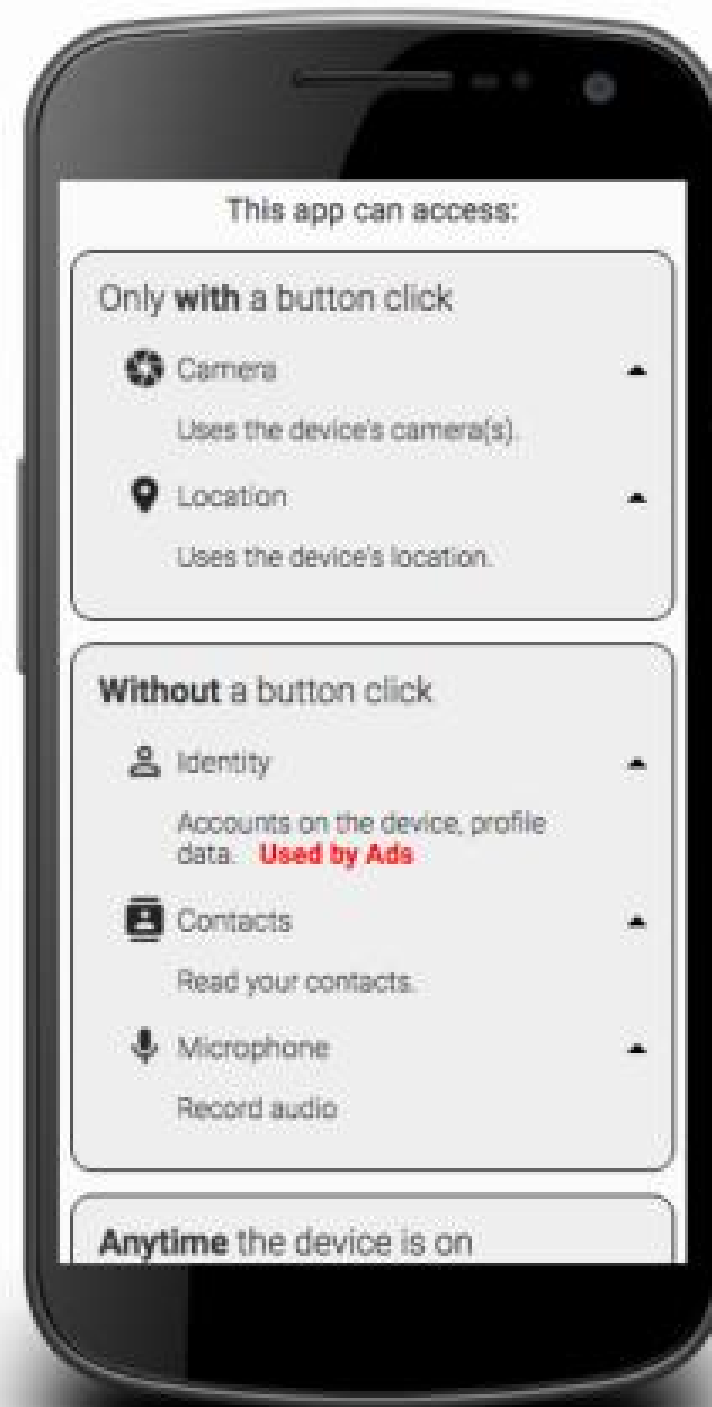
- Descriptive questions vs testing a question
- Correlation vs causation
- Dependent vs independent variables
- Between and within subjects testing
- **Numeric vs categorical data**

Types of data

- Numeric
 - **Continuous** — Any value on the range is possible including decimal (1–5)
 - **Discrete** — Only certain values on the range are possible (1,2,3,4,5)
 - **Interval** — Only certain values on the range are possible and each has equal distance from its neighboring values (strongly agree, agree, neutral, disagree, strongly disagree)
- Categorical
 - **Binary** — Only two possibilities (true, false)
 - **Ordinal** — The values have an ordering (slow, medium, fast)
 - **Nominal** — The values have no ordering (apple, pear, kiwi, banana)

Study design

- Accuracy on all tasks
 - Discrete
- Which interface
 - Categorical binary



Statistical tests

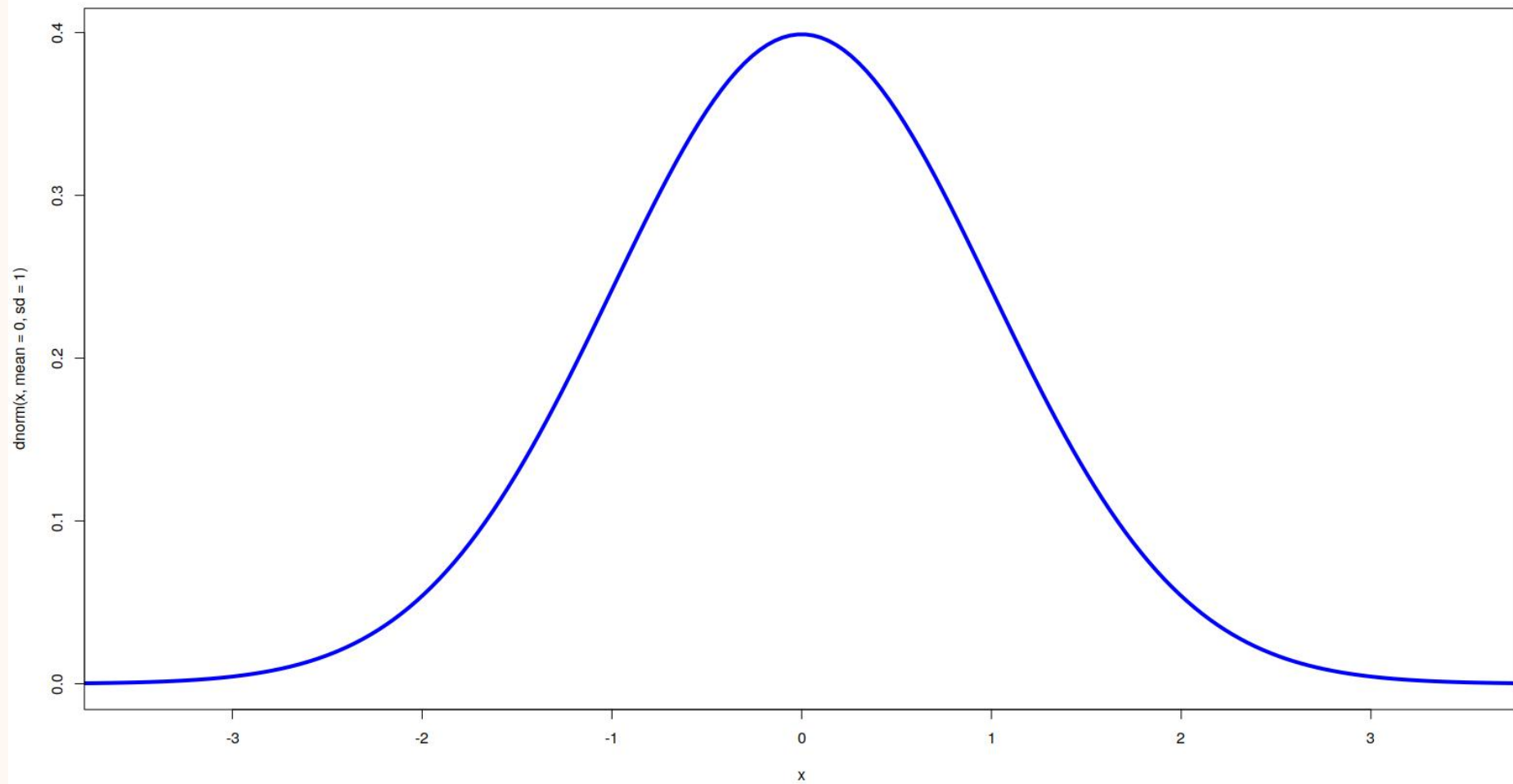
Comparing	Dependent	Independent	Parametric (Dependent variable is mostly normally distributed)	Non-parametric
The means of two independent groups	Continuous / scale	Categorical / nominal	Independent t-test	Mann-Whitney test
The means of 2 paired (matched) samples	Continuous / scale	Time variable (before/after)	Paired t-test	Wilcoxon signed rank test
The means of 3+ independent groups	Continuous / scale	Categorical / nominal	One-way ANOVA	Kruskal-Wallis test
3+ measurements on the same subject	Continuous / scale	Time variable	Repeated measures ANOVA	Friedman test
Relationship between 2 continuous variables	Continuous / scale	Continuous / scale	Pearson's Correlation Coefficient	Spearman's Correlation Co-efficient
Predicting the value of one variable from the value of a predictor variable	Continuous / scale	Any	Simple Linear Regression	
Assessing the relationship between two categorical variables	Categorical / nominal	Categorical / nominal		Chi-squared test

**t-test: Test if two groups have the same mean
(average)**

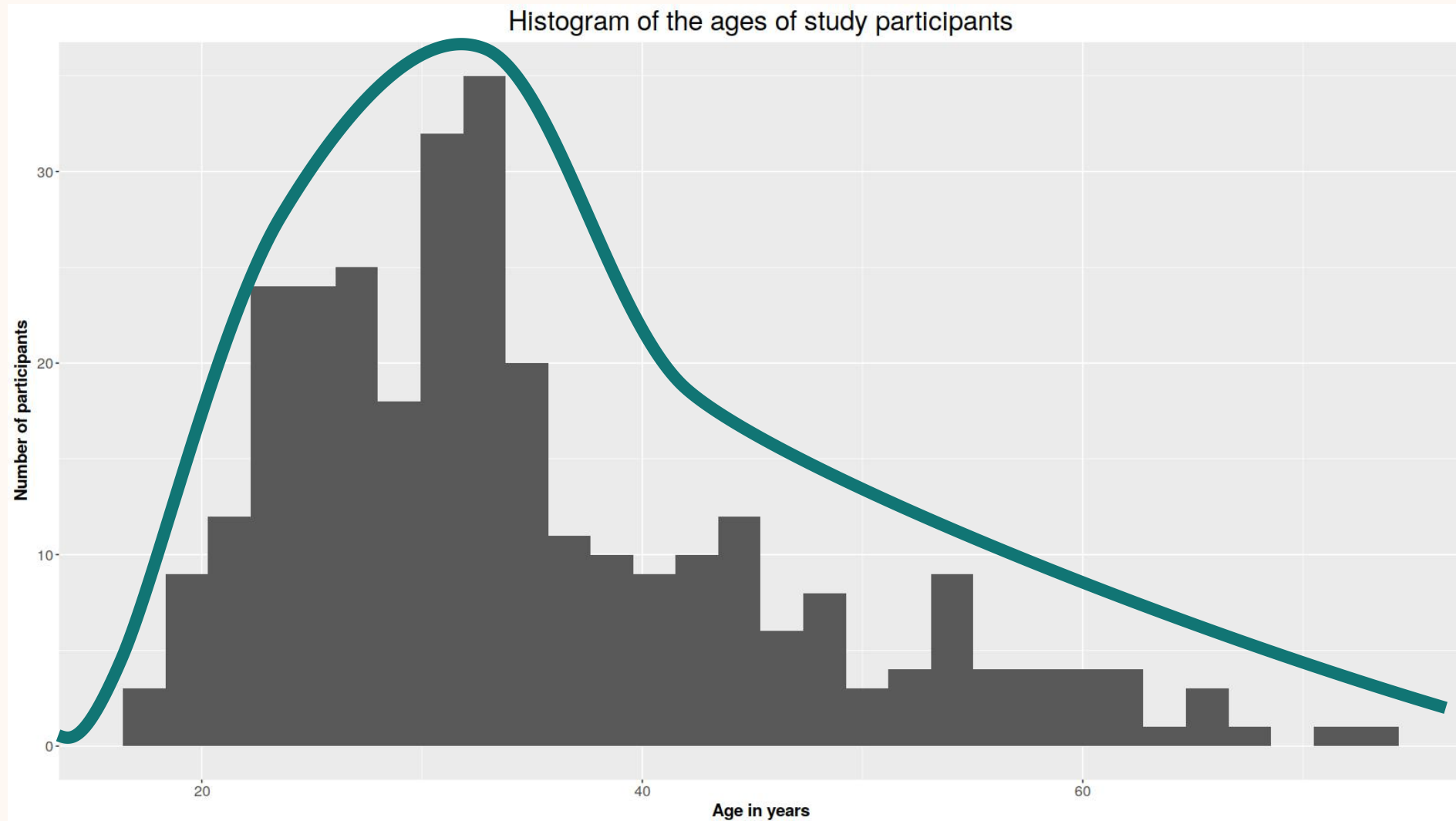
T-test requires:

- Independent variable: categorical binary
- Dependent variable: numeric (continuous or discrete)
- Data must be **normally distributed**

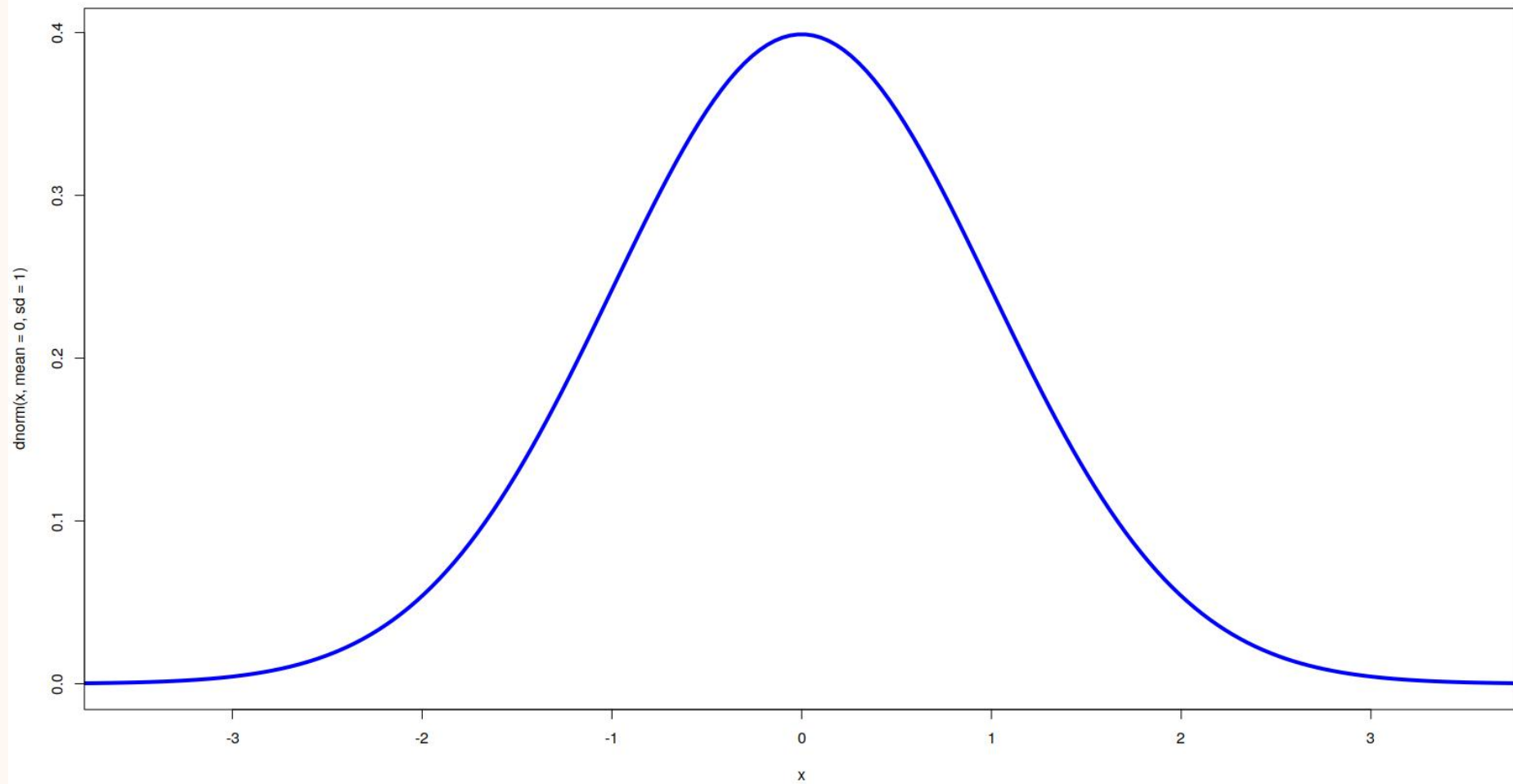
Normal distribution



Real data is messy

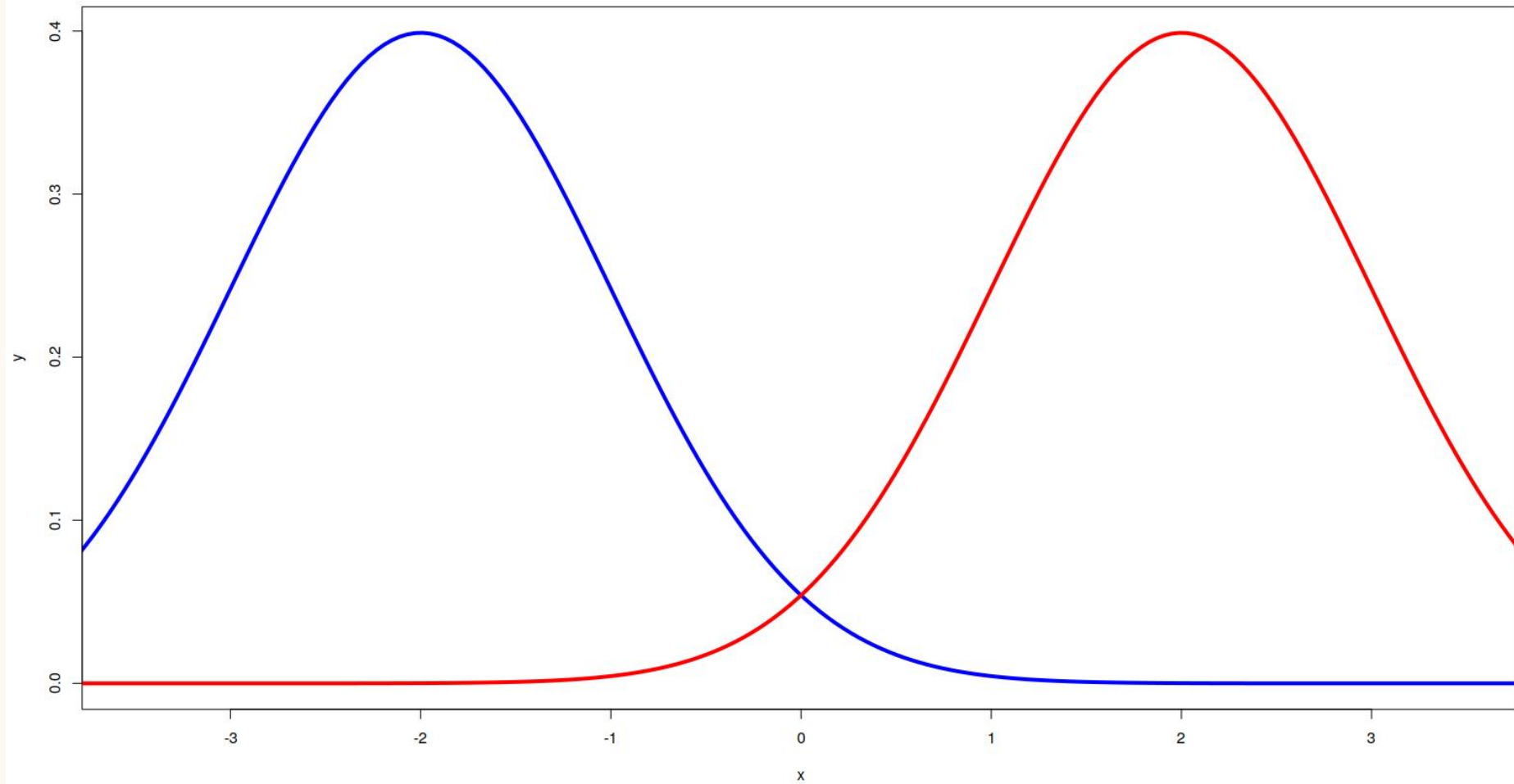


Normal distribution

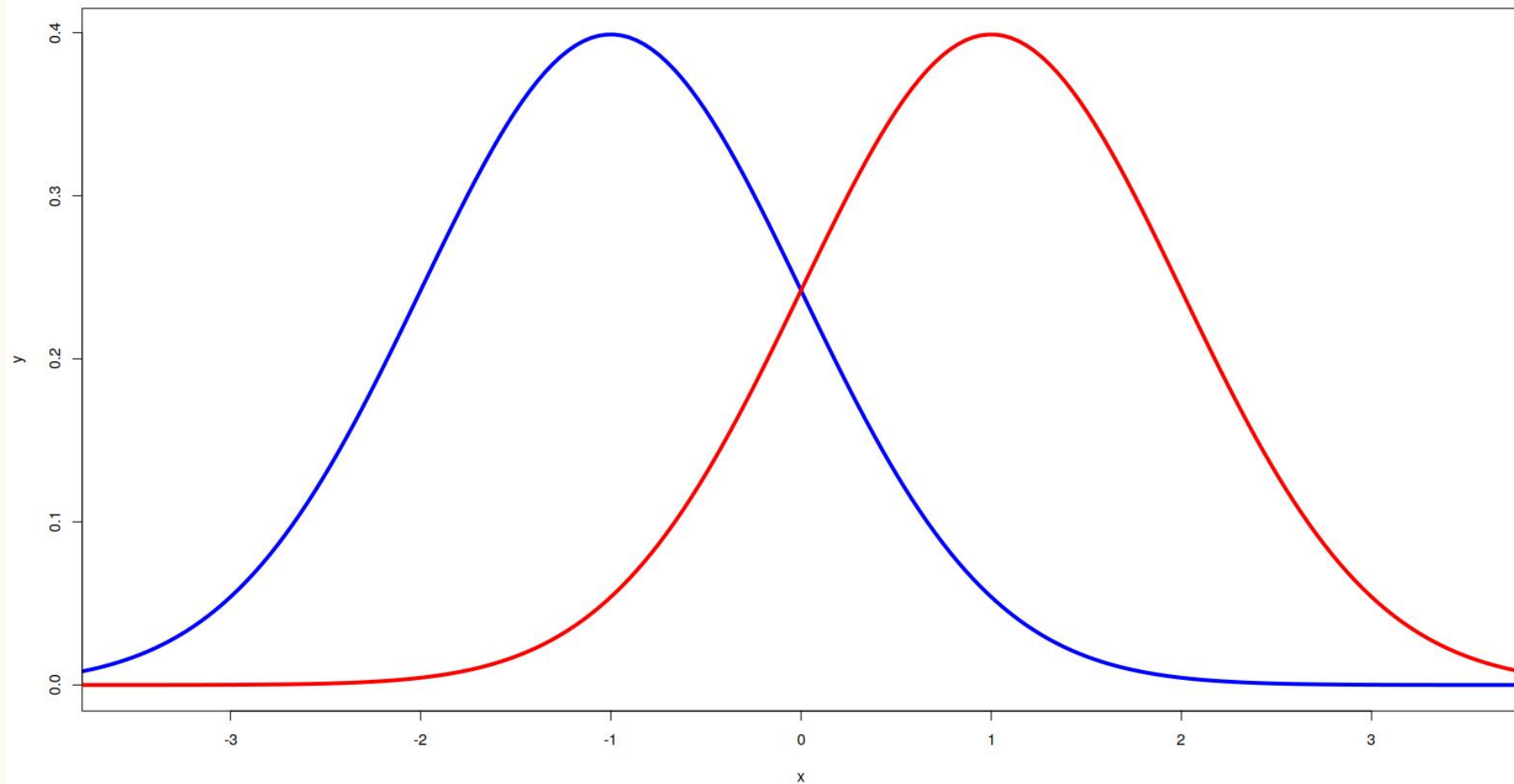


T-test: Do two populations have the same mean?

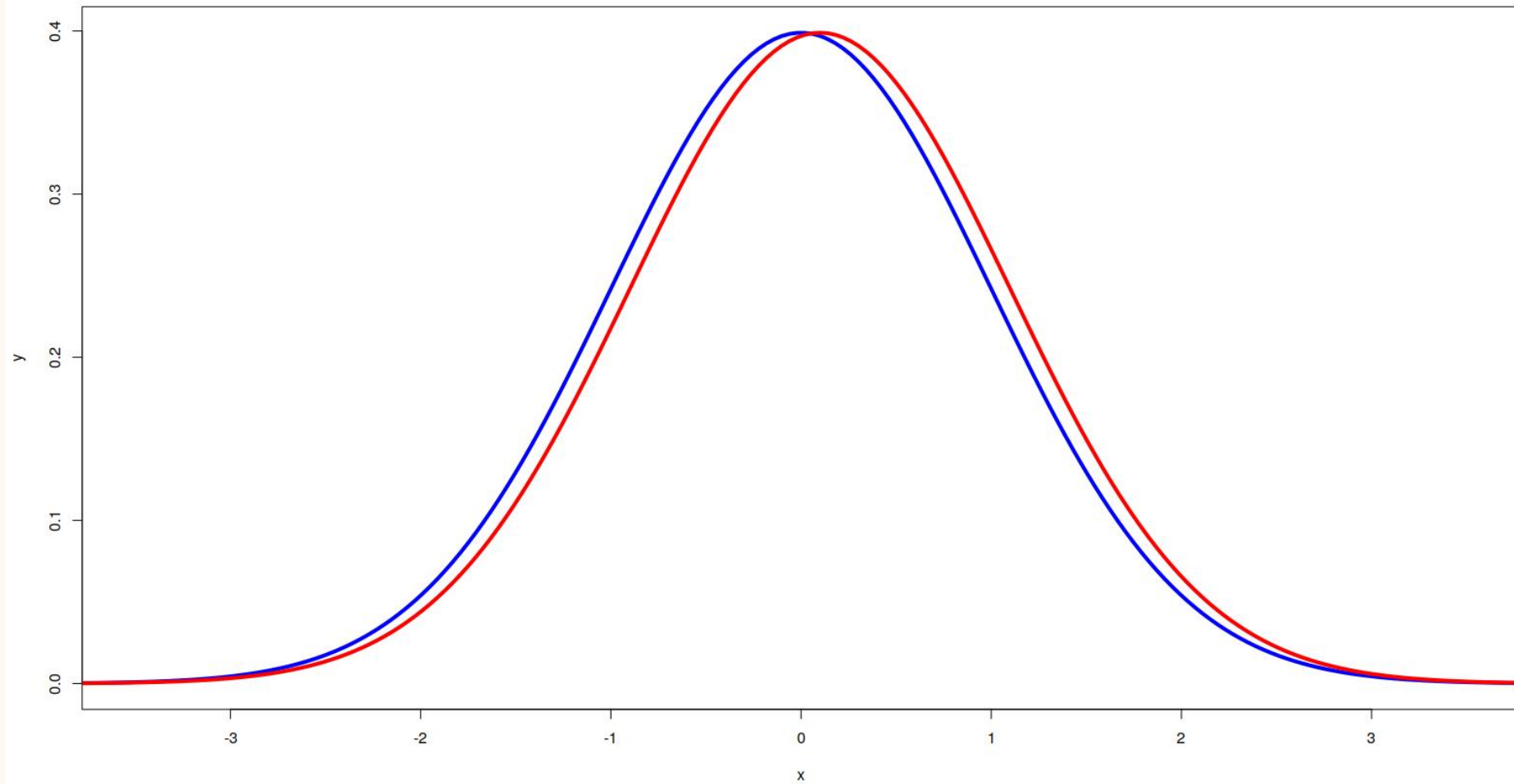
Different means



Maybe? different means



Likely not different means



I showed participants 4 code samples and asked them what the code would do. I then asked them how confident they were in their answer.

Research Question: Does the code sample shown impact confidence in their answer?

Research Question:

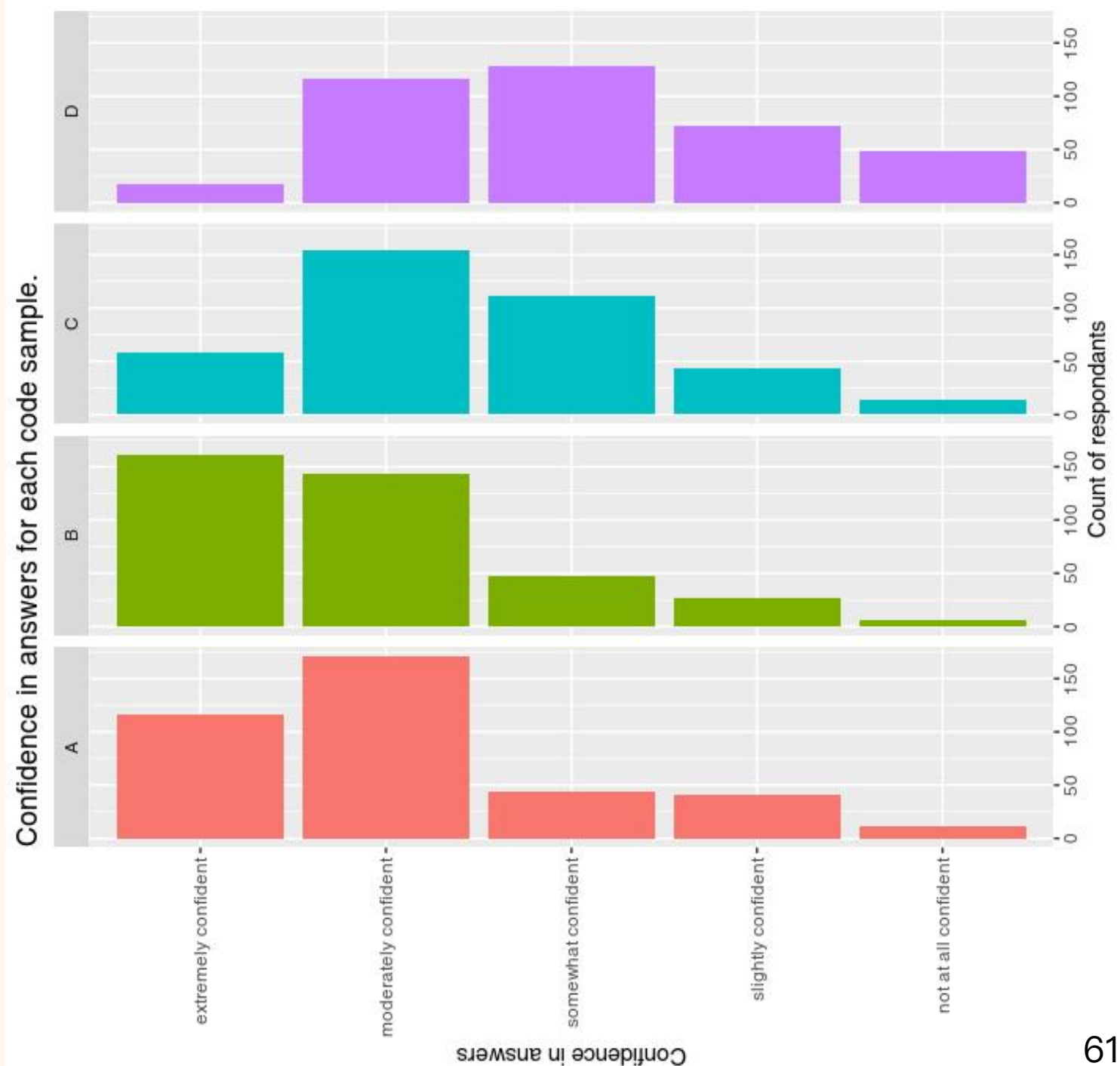
Does the code sample shown impact confidence in their answer?

Within-subjects

Independent:

Which code sample shown

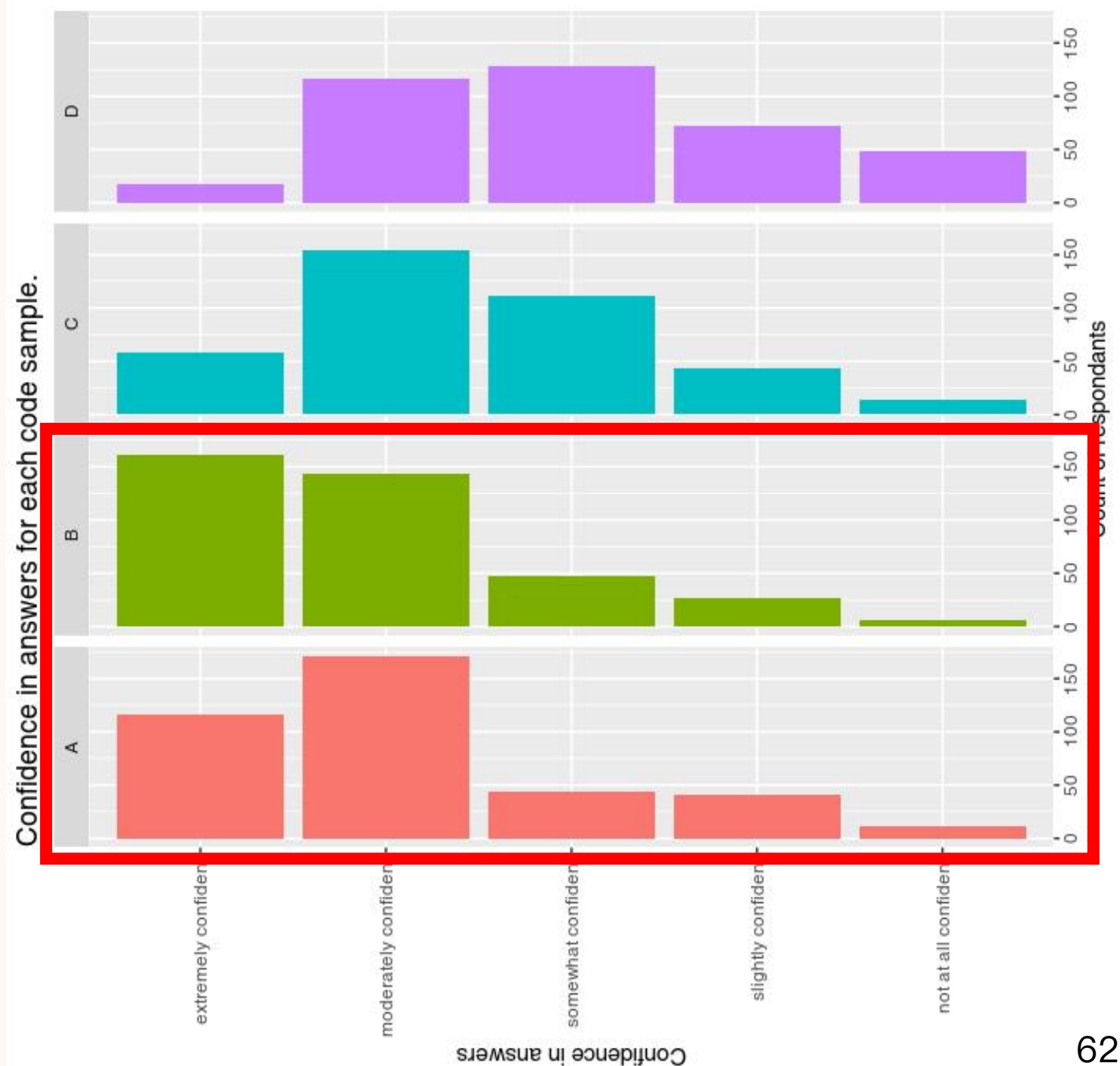
Dependent:
Confidence



Problem: My categorical variable (code sample) is not binary, there are 4 levels.

Solution: Run the t-test on each pair. So test A vs B, A vs C, C vs D.

Real solution: Use an ANOVA (not covered in this class)



Running the t-test

- This is a “**within** subjects” test where one person gave a confidence answer for **both** Code Sample A and Code Sample B
 - So we use a **Paired t-test**
- Create two arrays (or Excel columns) one with Code Sample A confidence, the other with Code Sample B confidence
- Two-sided (tailed)
 - For now, just do this. I don't have time to explain.
- Alpha of 0.05
 - p-value needs to be less than 0.05 to show that the two code samples produce different levels of confidence
 - Means that 5% of the time we will get the wrong answer from the statistical test

```
> t.test(a.confidence,b.confide
```

Paired t-test

data: a.confidence and b.confidence

t = -5.2699, df = 383, p-value = 2.285e-07

alternative hypothesis: true difference in means is not equal to 0

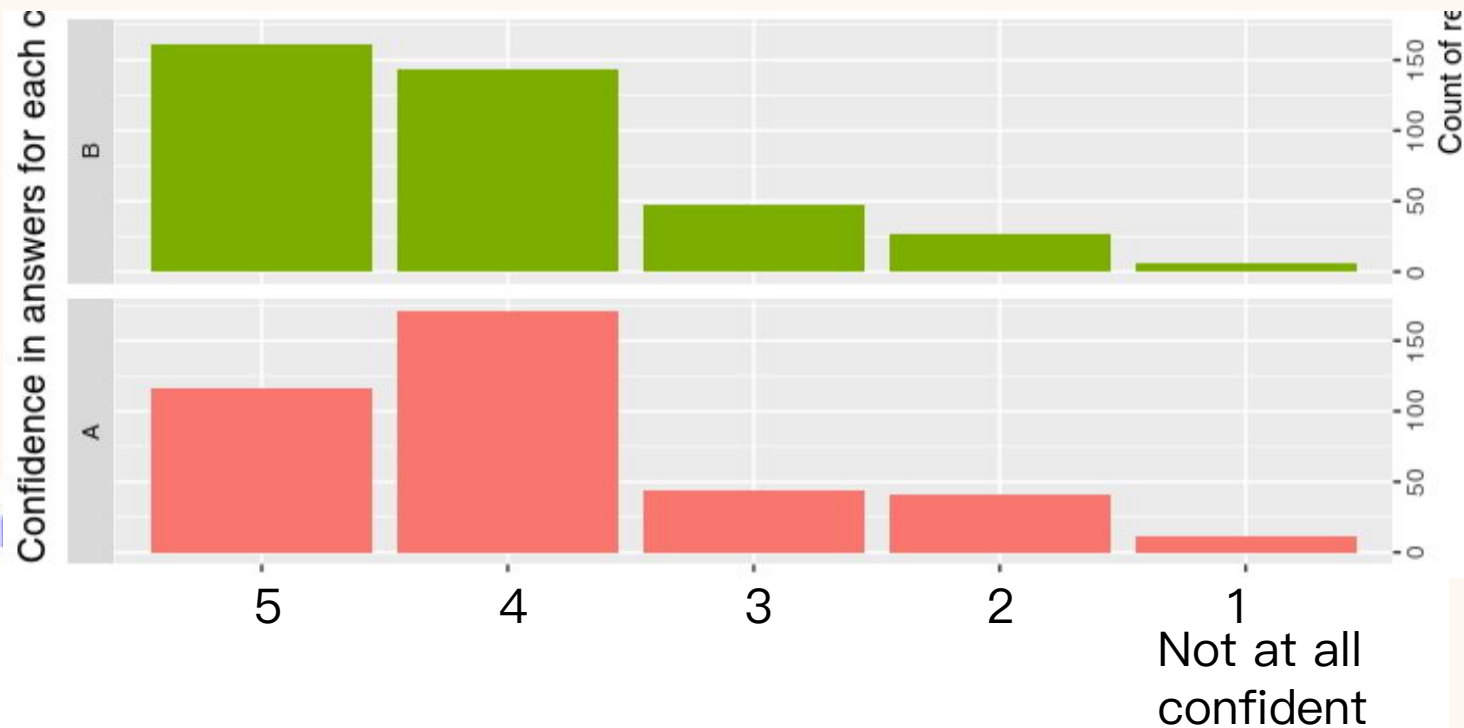
95 percent confidence interval:

-0.3218198 -0.1469302

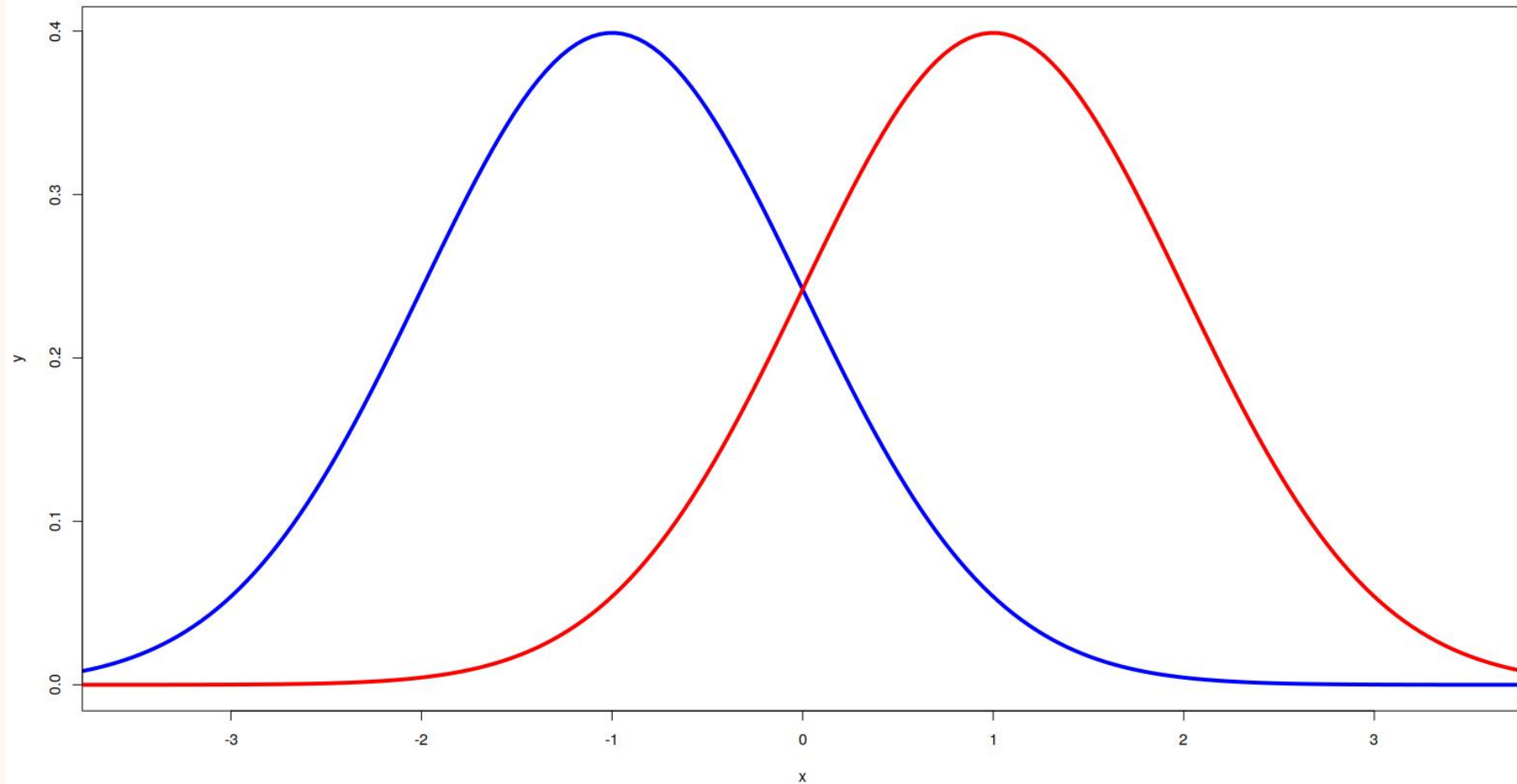
sample estimates:

mean of the differences

-0.234375



Different means, small difference



I ran a survey to learn about software update behaviors.

Research Question: Do women and men feel like they ask others for technical help with different frequency?

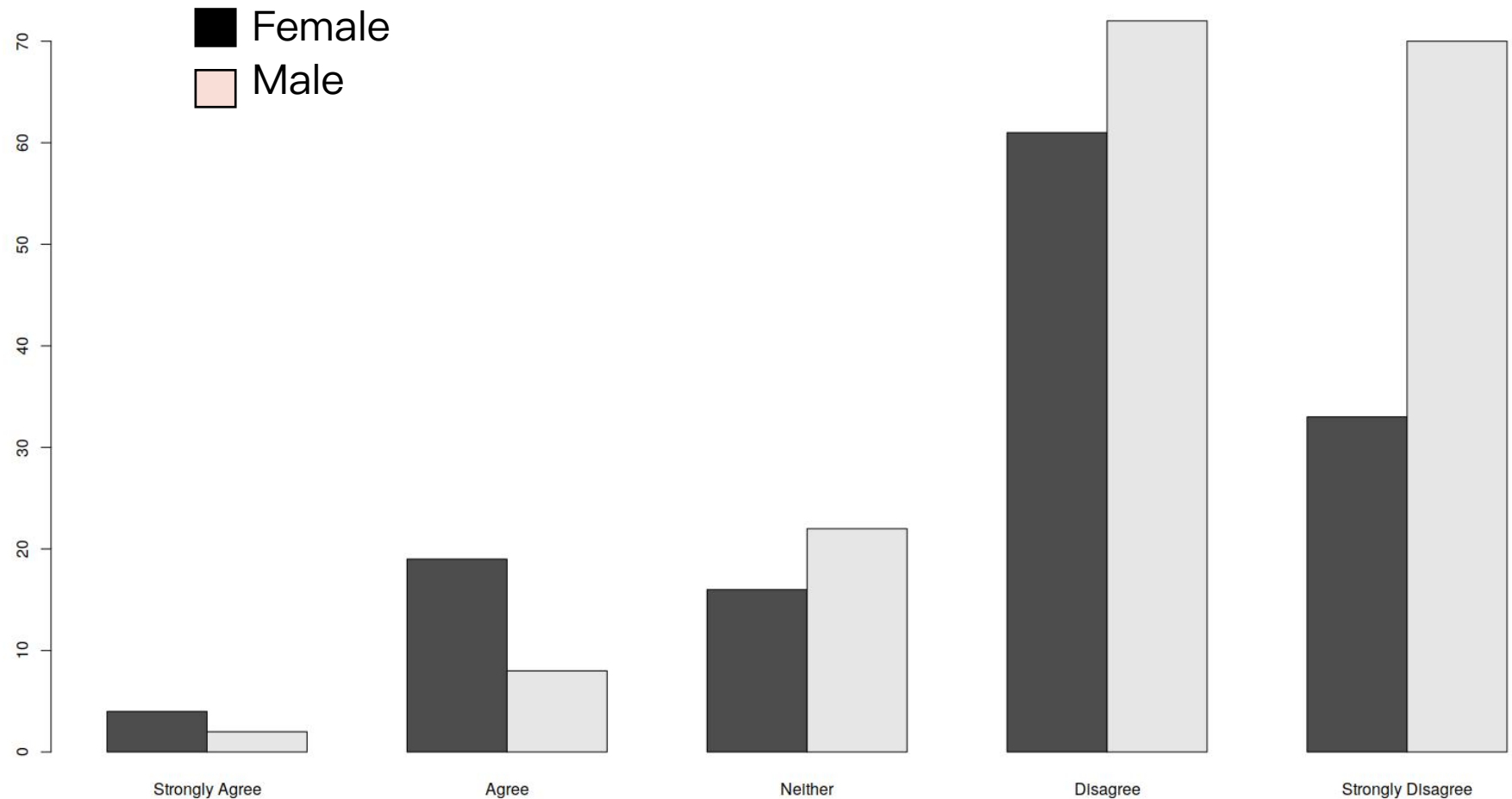
Research

Question: Do women and men feel like they ask others for help with different frequency?

Between-subjects

Independent:
Gender

Dependent:
Agreement



I often ask others for help with technical questions

Running the t-test

- This is a “**between** subjects” test where each person gave only one answer
 - So we use a **normal t-test** (not paired)
- Create two arrays one with women’s responses, one with men’s
- Two-sided (tailed)
 - For now, just do this. I don’t have time to explain.
- Alpha of 0.05
 - p-value needs to be less than 0.05 to show that the two genders produce different levels of confidence
 - This choice means that 5% of the time we will get the wrong answer from the statistical test

```
> t.test(as.numeric(d$i_ask_others_for_help),
```

Welch Two Sample t-test

```
data: as.numeric(d[i_ask_others_for_help[d$gender == "Female"] and  
t = -3.4481, df = 253.99 p-value = 0.0006606
```

```
alternative hypothesis: true difference in means is not equal to 0
```

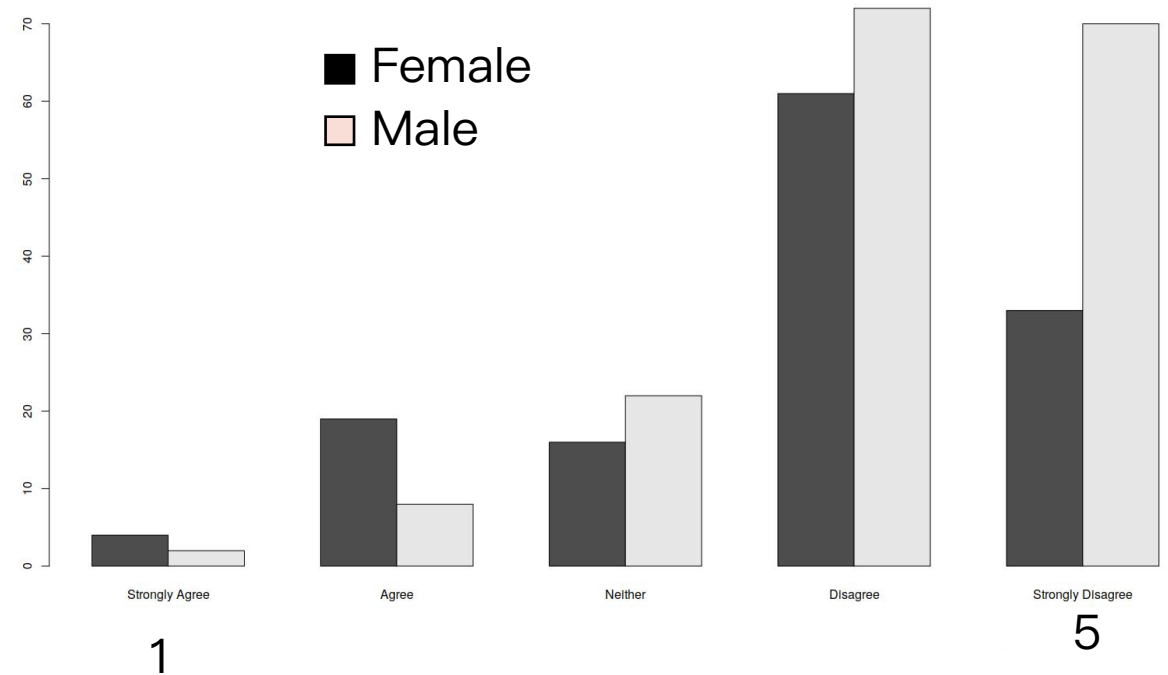
```
95 percent confidence interval:
```

```
-0.6245978 -0.1704934
```

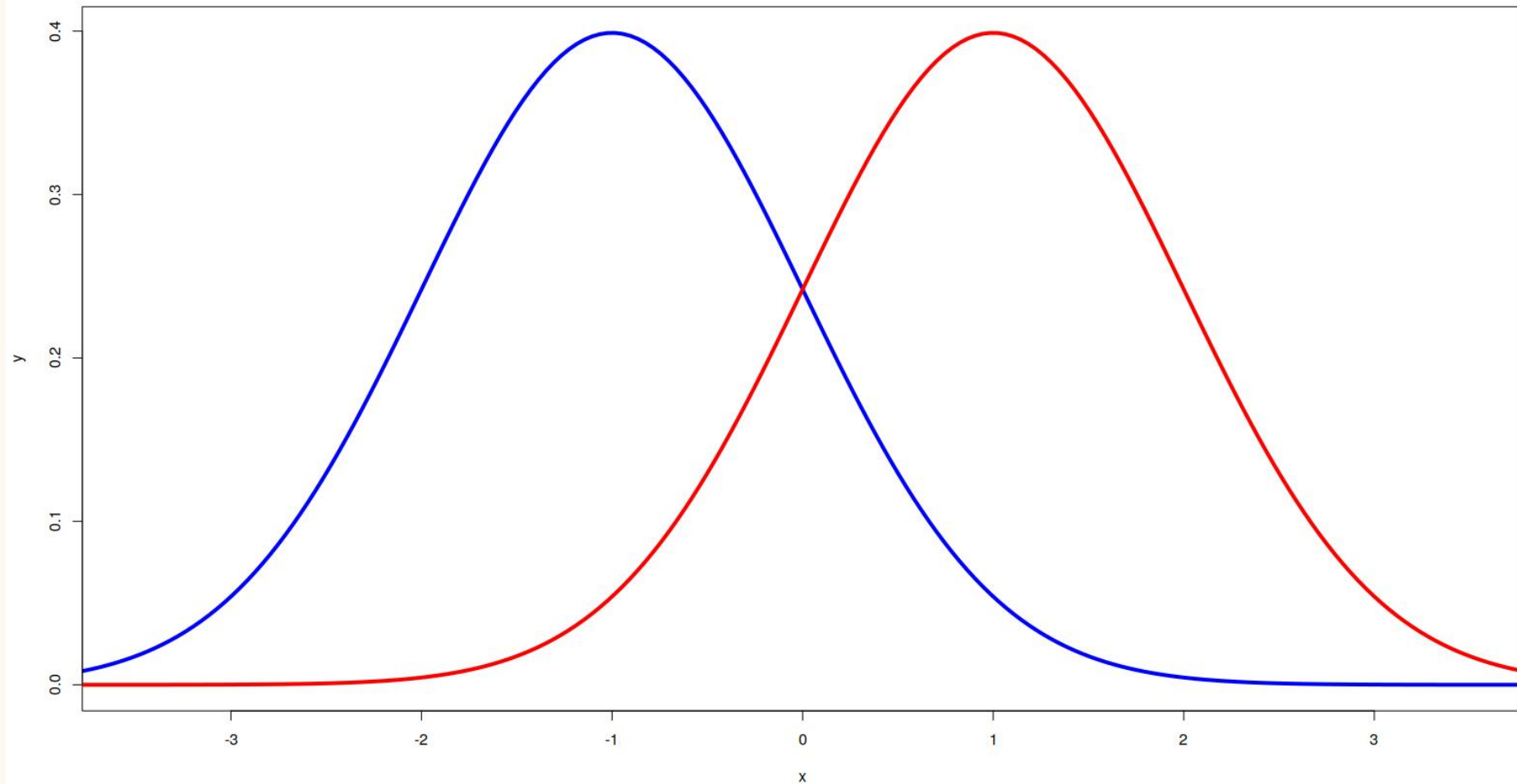
```
sample estimates:
```

```
mean of x mean of y
```

```
3.751880 4.149425
```



Maybe? different means



I asked participants to tell me a story about a prior software update.

Research Question: Are people who relate positive stories older or younger?

Research

Question: Are people who relate positive stories older or younger?

Between-subjects

Dependent:

- Age

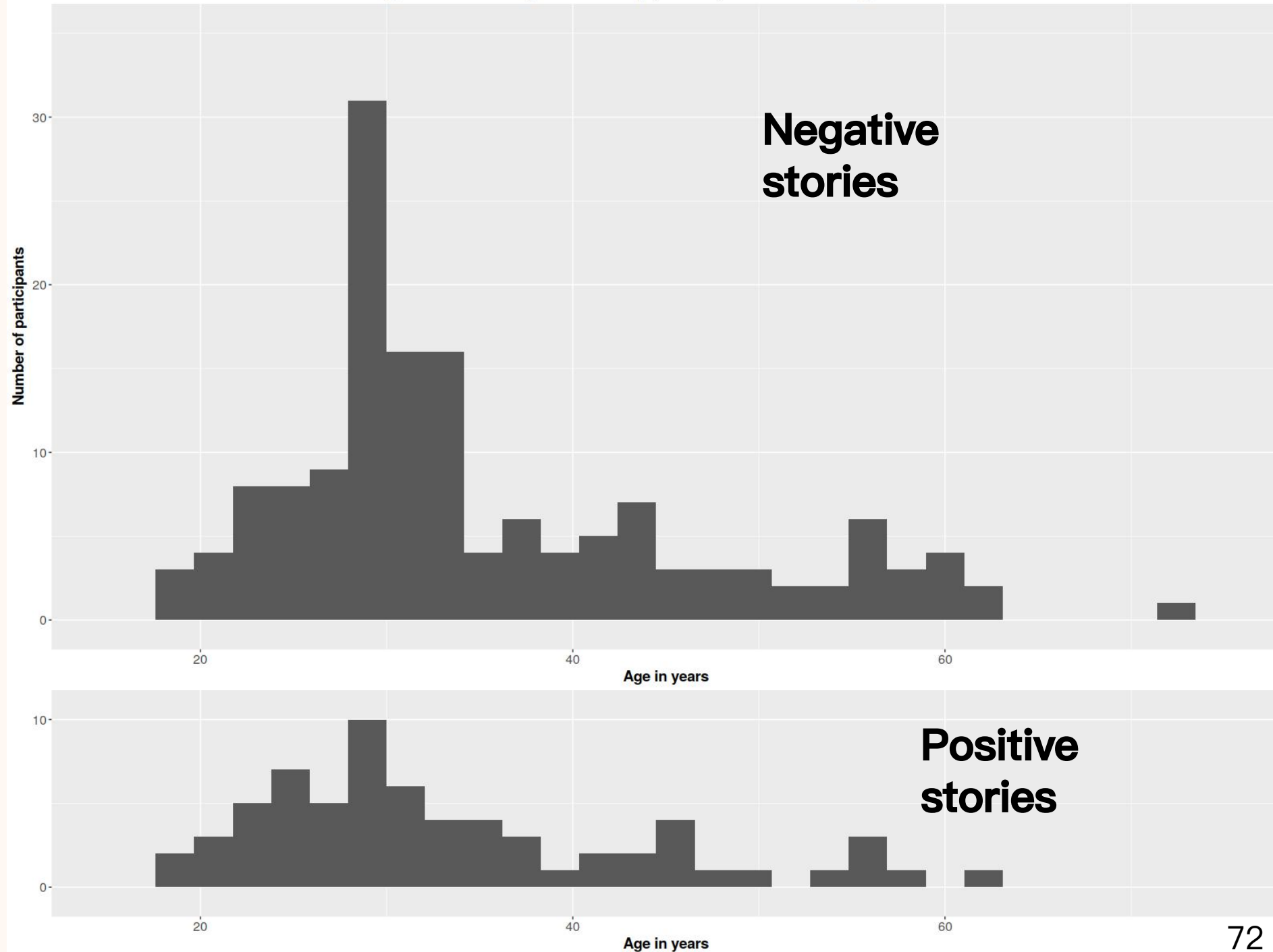
- Numerical

Independent:

- Negative or

- Positive

- Binary



```
> t.test(s_neg$age, s_pos$age)
```

Welch Two Sample t-test

data: s_neg\$age and s_pos\$age

t = 0.75677, df = 123.07, p-value = 0.4506

alternative hypothesis: true difference in means is not equal to 0

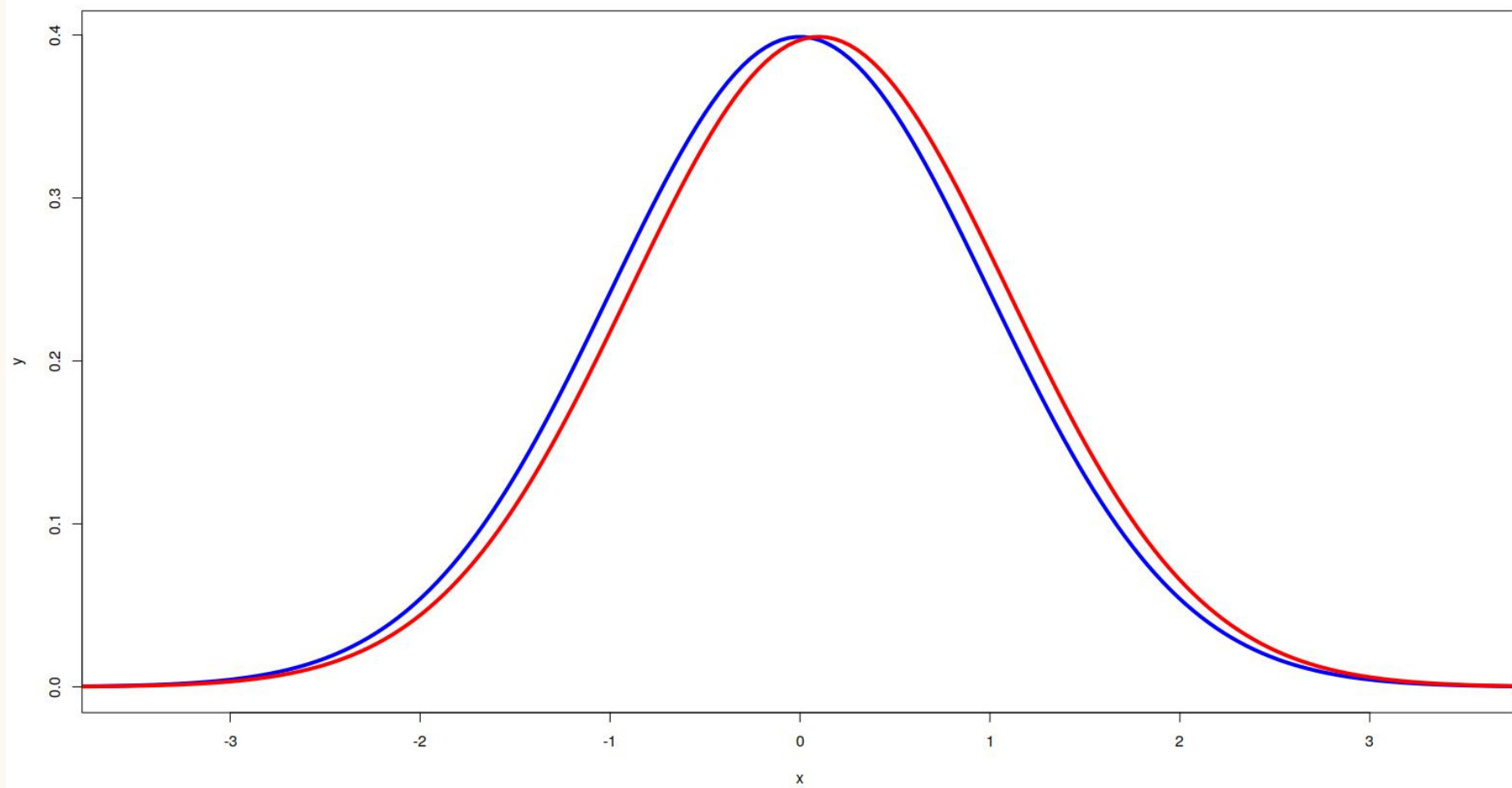
95 percent confidence interval:

-2.063833 4.618658

sample estimates:

mean of x mean of y

35.42667 34.14925



Questions

Take-home

- Frank, J., Herbert, F., Ricker, J., Schönherr, L., Eisenhofer, T., Fischer, A., Dürmuth, M. and Holz, T., 2024, May. A representative study on human detection of artificially generated media across countries. In *2024 IEEE Symposium on Security and Privacy (SP)* (pp. 55-73). IEEE.
- Li, J., Sun, K., Huff, B.S., Bierley, A.M., Kim, Y., Schaub, F. and Fawaz, K., 2023, May. “It’s up to the Consumer to be Smart”: Understanding the Security and Privacy Attitudes of Smart Home Users on Reddit. In *2023 IEEE Symposium on Security and Privacy (SP)* (pp. 2850-2866). IEEE.