Charlotte Hadley

1. There's a reproducibility crisis

2. Everything you can do to help reproducibility

3. Understanding Open Access and Open Data

4. Acknowledging legal and funder requirements affecting reproduciblity

Lecture 🚊 2 / 6 4

"There's a reproducibility crisis in academia"

Have you heard this b core?

There is a reproduciblity crisis in academia

B akdn 2005 a ground b eaking paper by John Ioannidis [1] exposed an unsettling truth

Why Most Published Research Findings

Are False

John P. A. Ioannidis

Summary

There is increasing concern that most current published research findings are false. The probability that a research claim is true may depend on study power and bias, the number of other studies on the same question, and, importantly, the ratio of true to no relationships among the relationships probed in each scientific field. In this framework, a research finding is less likely to be true when the studies conducted in a field are smaller; when effect sizes are smaller; when tested relationships; where there is greater number and lesser preselection of tested relationships; where there is greater financial and other interest and prejudice; and when more teams are involved in a scientific field in chase of statistical significance. Simulations show that for most study designs and settings, it is more likely for

factors that influence this problem and some corollaries thereof.

Modeling the Framework for False Positive Findings

Several methodologists have pointed out [9–11] that the high rate of nonreplication (lack of confirmation) of research discoveries is a consequence of the convenient, yet ill-founded strategy of claiming conclusive research findings solely on the basis of a single study assessed by formal statistical significance, typically for a *p*-value less than 0.05. Research is not most appropriately represented and summarized by *p*-values, but, unfortunately, there is a widespread notion that medical research articles

It can be proven that most claimed research findings are false.

being true is R/(R+1). The probability the Type II error rate). The probability circumscribed fields where either there several existing true relationships. The field targets highly likely relationships the power is similar to find any of the pre-study probability of a relationship of claiming a relationship when none is only one true relationship (among rate, \alpha. Assuming that crelationships of a study finding a true relationship vary a lot depending on whether the expected values of the 2×2 table are true relationships among thousands and millions of hypotheses that may reflects the power $1 - \beta$ (one minus truly exists reflects the Type I error be postulated. Let us also consider, many that can be hypothesized) or finding has been claimed based on are being probed in the field, the or searches for only one or a few for computational simplicity,

Lecture 🚊 4 / 6 4

There's still a crisis almost two decades later

Des pte lotsbern productive about the crisis... it's still here. And with the rise of

"m achine learning can solve anything!"

... the cris sis evolving and getting more complicated, as reported by Douglas Heaven [2].

Hundreds of AI tools have been built to catch covid. None Some have been used in hospitals, despite not being properly tested. But the pandemic ARTIFICIAL INTELLIGENCE of them helped. could help make medical AI better. By Will Douglas Heaven July 30, 2021

Reproducibility vs Replicability

meanings The American Statistical Association (ASA) provides useful advice from Broman, Çetinkaya-Quite a lot of the time you'll see these two words used interchangeably - but they have specific Rund & Nuss bma Raciorek, Peng, Turek, and Wickham [3].

Reproducible if you can take the original data and the computer code us designally to analyze the deat and reproduce all of the numerical findings from the sutdest

Your mod Leas ssesemnt must be reproducible—I mus Lobe eable to run the code after it's banted.

Reproducibility

Reproducible II ty: A suld vis reproducible if you can take the original data and the computer code us deto analyze the detachdreproduce all of the numerical findings from the study.

Reproducibility is something to think about from the start of a research project:

- Plan to record and d cument all processes in d aatcollection, wrangling and analysis
- While performing the research keep track of all the things you do particularly with data!
- When writing up your res arch ensure that all necessary materials to reproduce your findings are mad ævailab **e**.

Later in this lecture and the workshop we'll look at specific advice for achieving this.

| pcti | rp 🖰 7 / 6 /

Replicability

Replicab lity: This is the act of repeating an entire study, independently of the original investigator without the us @foriginal data (but generally using the same methods) This is a sightly hard etopic to conceptualise and much of the lack of replicability comes from what we call "q ues iotnab de res aerch practices" (QRPs)

... thes enis bhaviours lie somewhere on a continuum between scientific fraud, bias, and simple careles ses so dheir direct inclusion in the "falsification" category is debatable, although their negative impact on res arch can be dramatic

Fanelli [4]

Lecture 🚊 8 / 6 4

Questionable Research Practices

There are lots of different ways to summarise common QRPs but I quite like this table from 2012 by John, Loewens din, and Prelec [5].

Questionab & Researc IPrac it e	Self-ad msion rate (Amongst 2,000 pysyc blogists)
In a paper, failing to report all of a stud 1/5 d pend at measures	63.4%
Dec ding whether to colube thore disa tafter looking to see whether the results were significant	25.9%
In a paper, selec ively reporting stud &s that "worked"	45.8%
Dec d ing whether to exculd ed aa tafter looking at the impact of d ing so on the results	38.2%
In a paper, failing to report all of a stud 1/5 cond tilons	27.7%
In a paper, reporting an unexpec &d find mg as having ben pred ted from the start	27.0%
In a paper, "round nig off" a p value (e.g., reporting that a p value of .054 is less than .05)	22.0%
Stopping cubec ing dua barlier than planned beause one found the result that one had been looking for	15.6%
In a paper, c diming that results are unaffected by demographic variabes (e.g., gender) when one is actually unsure (or knows that they decorately decorated by d	3.0%
Falsifying d aat	%9:0

Questionable Research Practices and P-hacking

The majority of thes @RPs can be categorised as "P-hacking" or more fully, hacking the P-value.

Questionab & Researc IPrac it: e	Self-ad rission rate (A mongst 2,000 pysyc blogists)
In a paper, failing to report all of a stud 1/5 d pend of the measures	63.4%
Decding whether to cities thore disabilities looking to see whether the results were significant	92.9%
In a paper, selec ively reporting stud és that "worked"	45.8%
Deciding whether to excilidical adalafter looking at the impact of discoon the results	38.2%
In a paper, failing to report all of a stud 1/8 cond tilions	27.7%
In a paper, reporting an unexpec &d find mg as having b en pred ted from the start	27.0%
In a paper, "round mg off" a p value (e.g., reporting that a p value of .054 is less than .05)	22.0%
Stopping cubec it gd aa tearlier than planned bue eause one found the result that one had been looking for	15.6%
In a paper, c diming that results are unaffected by denographic variabes (e.g., gender) when one is actually unsure (or knows that they decorately decorated by d	3.0%
Falsifying d aat	%9.0

Lecture 🖺 1 0'6 4

What is the p-value?

We're going to talk ab out p-values A LOT in week 10. For now I wanted to borrow a slide from Lucy D'Agos ino McGowan's talk which is well explained in this Twitter thread.

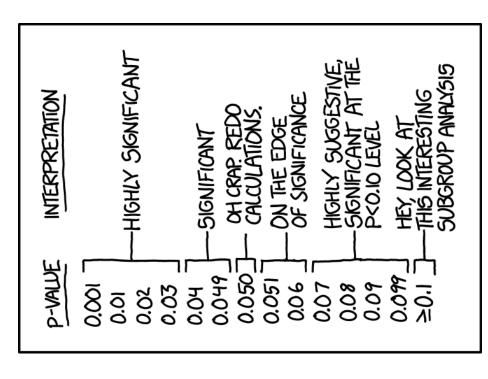
How folks treat the p-value

In some stuations" pvalues" are considered infallib evidonce of an effect or the conclusion of a sutdy

There are lots of different p-value thresholds, but the most common in healthcare datas icence is 0.05

Res ærchers who find their study results in values just above 0.05 will explore ways to get the value b æw 0.05

That's p-value hacking.



https //xkcd com/1478/

P-hacking

In 2019 As lawand @ [6] published an article in Wired.com titled "We're All 'P-Hacking' Now" which I highly recommend read rig.

It highlights an excellent study by Simmons, Nelson, and Simonsohn [7] that was able to use p-value hacking to make two increasingly absurd conclusions:

Stud yl: Lis ening to a children's song ("Hot Potato" by The Wiggles) makes people feel old

Stud 1/2: Lis ening to a s and about old age ("When I'm Sixty-Four" by The Beatles) makes people actually younger. ... I really like this paper because it makes very clear recommendations to researchers and reviewers.

Lecture 🖺 1 3/6 4

Have you ever read an academic paper? doi.org/10.1177/0956797611417632

Always expect to read a paper multiple times and make notes

We can kind of neatly s pt papers into two different types:

Clinical trials like this one from Hu, Liu, Wu, and Fang [8]

In some clinical trials the abstract includes a lot of structured information - it doesnot son the pubils br.

Some journals like B M Open even include suid ys rengths and weaknesses.

Read thes @apers in this order:

- Ab sratct
- Tab &s and or figures
- Conclus dns
- Introd ation

B ut most papers look this one from Simmons et al. [7]



Psychologial Science 22(1) 1359-1366 © The Author(s) 2011 Septimiz and permissions: suggestud-complournalisems: Suggestud-complournalisems: nav DOI: 10.1177/1958.5979.611417632 MSSAGE

Flexibility in Data Collection and Analysis Allows Presenting Anything as Significant

False-Positive Psychology: Undisclosed

General Article

Joseph P. Simmons¹, Leif D. Nelson², and Uri Simonsohn¹ The Wharton School, University of Pennsylvania, and ¹Haas School of Business, University of California, Berkeley

Abstrac

In this article, we accomplish two things. First, we show that despite empirical psychologists' nominal endorsement of a low rate of false-positive findings (≤ .05), flexibility in data collection, analysis, and reporting dramatically increases actual false-positive rates. In many cases, a researcher is more likely to falsely find evidence that an effect exists than to correctly find evidence that it does not. We present computer simulations and a pair of actual experiments that demonstrate how unacceptably easy it is to accumulate (and report) statistically significant evidence for a false hypothesis. Second, we suggest a simple, low-cost, and straightforwardly effective disclosure-based solution to rih problem. The solution involves six concrete requirements for authors and four guidelines for reviewers, all of which impose a minimal burden on the publication process.

In fact, not all med dal trials have such nice to read ab sracts.

But the read migord eisthe same:

- Ab sratct
- Tab es and or figures
- Condus dns
- Introd ation

- On your firs bas sof a paper you are trying to understand if the paper is relevant and provides s b santial information and/or evidence for your needs.
- Often it will take another 1 or 2 passes to understand the results of the paper.
- It us ally takes even more effort to understand the methods of the paper
- Method togy information is often provided in "supplementary materials" 0
- B ut unfortunately a good portion of the time you won't find sufficient information to fully und es and the method lagy because of poor reproducibility. 0

Learning ab out replicab lity and reproducibility now will help you in understanding the existing literature and prepare you to sacced in a research career later on. Lecture 🖺 1 7/6 4

If you deid do go into a research career you'll likely be reading 10+ papers a week.

I'd highly recommend investing in learning speed reading early in your career.

There's lots of [very interesting] eye tracking and neurological research into how we read and evid enced method togies for speed reading that are nicely summarised by Clifton, Ferreira, Hend es a, Inhoff, Liversedge, Reichle, and Schotter [9]

B C gives you free access to LinkedIn Learning.

Go to linked in.com/learning-login/and login with your B & email ad dess.



COURSE Learning Speed Reading (2014)

Learn how to read faster. Improve your reading speed and comprehension with these proven speed-reading techniques.

in LinkedIn · By: Paul Nowak · Nov 2014

4.7 ★★★★ (334) · 73,108 learners · Beginner

let's get back to p-hacking

Fighting p-hacking with pre-registration

So far I've b ens paking about the academic literature at large.

But let's looks pcifically at clinical trials.

All clinical trials b gan after July 1st 2005 are explicitly required to be registered in ord eto be pub is ad in all biomedical journals overs en by the International Committee of Medical Journal E dtors (ICMJE [10].

This is ind pend at of the country in which the trial took place.

[1] Frustratingly in the clinical trials community we use the phrase "registration" but everyone else says "pre-registration".

Lecture 🖺 2 0'6 4

National Institutes of Health (NIH) and registration

From a curs by sarch for clinical trial registration information you might come to the conclusion that the ad ice is only appropriate for studies in the US.

That is not true.

In 2008 the World Med dal Association [11] updated the Declaration of Helsinki – Ethical Principles for Medical Research Involving Human Subjects to includ a paragraph ab out registration.

19. E very clinical trial must be registered in a publicly accessible database before recruitment of the firs \$ b \(\psi \)ct. Lecture 🖺 2 1/6 4

International advice on clinical trials

The NIH provid & a really useful tool for comparing the clinical research regulations from several countries :

https //clinregs niaid nih.gov/country/united-kingdom

UK trial registration

Des jate thes & eq urements" violations are still common as detailed in Bradley, Lloyd, and DeVito [12].

collab oration beween the Human Research Authority and the ISRCTN regis ny which is deailed by However, the UK has announced new infrastructure requiring 100% clinical trial registration via a B rackner [13].

Registration to "Reduce waste"

For unexplainab & reas as there is a lot written about registration of trials "reducing research waste" asithelpsred oced polication of studies.

While there are elements of truth to this, the big takeaway is pre-registration helps prevent Ques ibnab & Res arch Practices.

In factit's the reas a why GSK reached a \$2.5million settlement in 2004

GlaxoSmithKline & Clinical Trial Registration

GSK chos to state a civil case instead of engaging in an expensive legal battle over "repeated and pers si ant fraud 'concerning the use of paroxetine in treating depression in adolescents.

Asd eailed by Dyer [14]:

- Two sutdess bwed nobelefit to using paroxetine when compared with placebo.
- Three sutdes found evidence for an increase in suicidal thoughts and behaviour.

Internal company d cuments confirmed the suppression of these results:

"it would b commercially unacceptable to include a statement that efficacy had not been d enons rated as this would undermine the profile of paroxetine."

Spurgeon [15]

lecture 🖺 2 5/6 4

Pre-registration for all studies

There's a growing b **d** yof researchers both actually doing pre-registration and calling for it in all d siplines - particularly those that intersect with healthcare.

- Simmons et al. [7] have been running as ped teed orgs rice 2015 to help authors to create preregis ration reports
- The Centre for Open Science encourages pre-regis ration on OSF and has in the past run a Preregis ration Challenge with a monetary prize.
- overview of pre-regis ration and walk through multiple example of how it works in practice. I highly recommend reading Nosek, E. Bersole, DeHaven, and Mellor [16] which gives a great

At the moment there's nothing binding you to pre-register non-clinical trials but this topic is s mmering away in the b ak**g**round Lecture 🖺 2 6/6 4

Prediction markets and replication studies

outcome of the event and the market price indicates what the crowd thinks the probability of the In pred ation markets investors make predictions of future events by buying shares in the

Harvard Pres Releas e

This was firs applied to predicting replicability of research results in 2015 by Dreber, Pfeiffer, Almenb eg, Is asksnpWilson, Chen, Nosek, and Johannesson [17]. There is now evid ace for these markets being a reliable estimate of predictability Camerer, Dreber, Holzmeis &r, Ho, Hub & Johannesson, Kirchler, Nave, Nosek, Pfeiffer, Altmejd, Buttrick, Chan, Chen, Fors 🕸 Gampa, Heikens 🚓 Hummer, Imai, Isaksson, Manfredi, Rose, Wagenmakers, and Wu [18].

Prof. Anna Dreb e gives an excellent 40minute overview of replication prediction markets here https //youtu.b @5rFDKB1aZc?t=1036

cture ≙ 2 7/6 .

A final word on replicability studies

Mos (finotall) of suid (esabout replicability are deeply technical and rely on statistical methods we d other time to cover in this course. This includ at the found aidhal paper "Why Most Published Research Findings Are False" by Ioannidis

Your as ssesemtd esnot require you to understand or replicate any of the methodologies behind replicab lity sutdes

Reproducibility

Let's get b akdo the ASA reproducibility recommendations from Broman, et al. [3].

Reproduciblity: A suld ys reproducible if you can take the original data and the computer code us **d** to analyze the d aatand reproduce all of the numerical findings from the study. In ord efor reproduiblity to be possible we need the original daatto be accessible - we need Open Data.

Open Data

Most papers don't provide their data

More often than papers simply do not provide the daatthat they:

- Us &o create charts and tables
- Us &o performs atis ital tests
- Us &o generate their conclusions

Thes papers are the antithesis of reproducible.

"Data available on request"

It's common to set he phrase "data available on request" but that's frequently meaningless:

" Data req ues atto authors are successful in 27-59% of cases, whereas the request is ignored in 14-41% cas **s** " Ted es o, Küngas Oras Köster, E enmaa, Leijen, Pedaste, Raju, Astapova, Lukner, Kogermann, and Sepp [19]

E en in cas swhere d aatsreturned it's often insufficient for reproducibility, Roche, Kruuk, Lanfear, and B inning [20]:

- Lata" might actually becreenshots of charts stored in an Excel workbook
- Lata" might b & dired in other "non-machine-readable" formats like PDF or images
- Lata" might only b eprovided in the post analysis form

Lecture 🖺 3 2/6 4

Researchers have confusing opinions

In one of the earlies , to agest and most cited studies in data sharing from 2011 by Tenopir, Allard, Douglas Ayd moglu, Wu, Read, Manoff, and Frame [21] surveyed 1,329 scientists found that " Mas tes pnd ats (at least 60% across disciplines) agree that lack of access to data generated by other res **a**rchers or ins i**t**utions is a major impediment to progress in science."

₽ 9 :: "A majority of all res pndents indicate they are not willing to place all of their data in central repos tories with no res rictions"

Thes find higs have ben replicated again, and again.

Lecture 🖺 3 3/6 4

Open Data helps everyone

Researchers

- Open Data helps reprod @e previous s utd &s
- Open Data gives an ad dtibnal way res archers can b exited
- There is clear evid eace Open Data is linked with higher citation rates eg Colavizza, Hrynas kiewicz, Stad ea, Whitaker, and McGillivray [22].

Society

- The coord mated glob adOVID-19 response b enefited significantly from Open Data portals [23].
- Red aing energy cons aption of buildings
 b a sia occupancy from public datasets
 Roth, Lim, Jain, and Grueneich [24].
- E nærgency planning via use of the
 B æavioral Ris Æactor Surveillance System
 (B ÆSS) [25]

You c a find some healthc as spec flic examples on the c arse website eng7218.netlify.app/resourc s/open-data. Please do researc Iyour own - and c asider sharing them with the group.

Lecture 🖺 3 4/6 4

Open Data vs Big Data

B & Data is great. It's the driver behind the Internet of Things and much of modern healthcare technology. But in mostircumstances it is not Open Data.

The Open Definition

The Open Definition sets out principles that define "openness" in relation to data and content.

It makes precise the meaning of "open" in the terms "open data" and "open content" and thereby ensures quality and encourages **compatibility** between different pools of open material.

It can be summed up in the statement that:

"Open means anyone can freely access, use, modify, and share for any purpose (subject, at most, to requirements that preserve provenance and openness)."

Put most succinctly:

"Open data and content can be freely used, modified, and shared by anyone for any purpose"

Screens bt of [26]

Open Data short definition

Let's minimis thes a finitions:

Opendaatmus to e

- Legally open. The daatmust be subject to an open data license
- Technically open.
- Data files mus **b** emachine-readable and non-proprietary, which often means plain text. 0
- Acces & From a pub its server without password protection 0

Open daatmus bapen to humans and computers.

Lecture 🖺 3 6/6 4

Open Data licenses

In mos tas sthe Creative Commons "Choose a license tool" is the b sthot eas disthoice if you have a dataset you want to make into Open Data.

The CCOlicens & the most permissive license.

There are special Open Data licenses used by Governmental/Charity organisations designed to waive liab lity for us æg the Open Government License from the UK Government [27]

In general it's b & to us @ata licenses for data and software licenses for software.

Open Data examples

There are SO MANY differents arces (or publishers) of Open Data, for a good sample checkout the Open Data Es atials page from the World Bank [28].

Country-Level Open Data

- <u>Australia</u>
- Brazil
- · Costa Rica
 - Ghana Chile
- India
 - <u>Italy.</u> <u>Kenya</u>
- Moldova
- Morocco
- Russian Federation Philippines
- United Kingdom
- United States of America

City- & Subnational-Level Open Data

- · Buenos Aires, Argentina
- Edmonton, Canada Chicago, U.S.A.
 - Edo State, Nigeria London, U.K.
- Nantes, France Rennes, France
- San Francisco, U.S.A.
- Vienna, Austria
- Vancouver, Canada

Sector	Website
Agriculture	The USDA National Farmers Market Directory
Agriculture	U.K. Department of Agriculture and Rural Development
Budgets & Public Finance	WB Open Budgets
Budgets & Public Finance	<u>OpenSpending</u>
Budgets & Public Finance	International Budget Partnership
Budgets & Public Finance	The International Aid Transparency Initiative (IATI)
Budgets & Public Finance	U.S. IRS Tax Statistics
Education	Ed Data Inventory.
Education	MyData Office of Educational Technology.
Education	<u>CheckMySchool</u>
Energy & Extractive Industries	Extractive Industries Transparency Initiative
Energy & Extractive Industries	U.S. Department of Energy.
Energy & Extractive Industries	Enel Open Data - Largest power company in Italy
Environment	Open Climate Data
Environment	Fuel Economy Data, U.S. Environmental Protection Agen
Environment	New York City Environment Open Data
Geospatial	<u>OpenStreetMap</u>
Geospatial	Haiti Data geospatial information
Health	The U.S. Department of Health & Human Services
Health	Agency for Healthcare Research & Quality (AHRQ) Databases on healthcare cost & utilization in the U.S.
Health	WB Health Data
Information & Communication Technologies (ICT)	Australian ICT Open Datasets
Transport	<u>OpenPlans</u>
Transport	European Public Sector Information Platform: Transport
Water	Global water database

Open Data is awes one. But if we're responsible for data that can identify individuals or groups we have a [legal] d thy of care to protect that data. In the UK we have the Data Protection Act (legis ation.gov.uk [29]) which is the UK's implementation of the General Data Protection Regulation (GDPR).

In week 6 we will d stors &DPR and the Data Protection Act in the context of anonmyising data.

This is It a cours an the law school so we won't go hard into the definitions But there are some things we need to d s as .s

In the DPA (legis ation.gov.uk [29]) there are 6 different types of sensitive data - I'm highlighting the ones that cover d aatthat you might reasonably collect and consider "health data".

- (a) the procesing of personal data revealing rac al or ethnic origin, politic bopinions, religious or philosophic bbeliefs or trade union membership;
- (b) the procesing of genetic data for the purpose of uniquely identifying an individual;
- (c) the processing of biometric data for the purpose of uniquely identifying an individual;
- (d) the procesing of datacoc ening health;
- (e) the processing of data concentration; an individual's sex life or sexual orientation;
- (f) the processing of personal data as to [commission or alleged commission of an offence]

There's some recurs on, so let's pull out the definitions of "health data" from Part 7 Section 205 and lis everything together / AV / M

Part 7 Section 205 d fines the following

- "b ometricd aa" meanspersonal data resulting from spec flictec hior laproc osing relating to the physical, physiological or behavioural c barac etristic sofan individual, whic hallows or c of irms the unique identification of that individual, such as fac al images or dac ylosc pic data;
- "d da concerning health" means personal data relating to the physic bor mental health of an individual, including the provision of health c are servic s, whic Ineveals information about his or her health status;
- "geneticd ata" means personal data relating to the inherited or ac qired genetic characteristics of an individual which gives unique information about the physiology or the health of that individual and which results, in particular, from an analysis of a biologic asample from the individual in question;

But we s buld als as ittinclude:

- (a) the procesing of personal data revealing rac all or ethnic origin,
- (e) the processing of data concentration; an individual's sex life or sexual orientation;

Lecture 🖺 4 1/6 4

There's more read efriend yd cumentation about health data and DPA 2018 from the Information Commis sovier's Office [30].

We'll talk ab out this more later.

Upen Data mandates and UKRI

UKRI is res pns bi & for the 6 UK research councils who fund most university-based research in the UK.

Some res arch councils have their own "data sharing policies", but others depend on the "common principles on on res arch data" from UKRI [31].

b emad expenly availab & with as few restrictions as possible in a timely and responsible manner. Pub icly fund deres arch data are a public good and produced in the public interest. They should

Open Data mandates and UKRI

- Arts and Humanities Researc Nounc II(AHRC [32]): Relies on the UKRI [31] common principles.
- Biotec hology and Biologic bScienc & Researc (Counc Ii(BBSRC [33]): "BBSRC expects that all data (with accompanying metadata) should be shared in a timely fashion as soon as it is verified".
- Engineering and Physic aSc anc Researc (Counc Ii(EPSRC [34]): Has the most explicit data policy, including a requirement
- Ec nomicand Soc al Researc (Counc Ii(ESRC [35]): Explic tirequirement that "data will be made available [...] as Open Data".
- Medic bResearc Counc Ii(MRC): There's an "expec ation" from MRC [36] that data must be made open, they helpfully provide lots of advic about patient and population data.
- Natural Environment Researc (Counc Ii(NERC [37]): Flubs it by saying all research "must include a statement on how the supporting data and any other relevant researc Imaterials c a be accessed"
- Scence and Technologies Fac lities Council (STFC [38])": "STFC expects that published data should be made publicly available within six months of public tion unless justified otherwise"

Understanding the whole process

- 1. Sub rit a Data Management Plan (DMP) to a fund e, includ ing a data bring plan
- 2. Pre-regis dr your res arch (often smultaneous ywith s dp 1)
- Res eve a DOI to s dre your research data at a d aa tepos tiory
- 4. Do the resaerch
- Keep rawd aatpris ithe. Do not modify your rawd aa.t
- Keep track of how you wrangle the data (much eas dr when you write R code!)
- Craft an anonymis d s breable version of your d aas tesd os db ed in your DMP.

- 1. Write up the suld y
- 2. Choos a journal to s b rit to
- That's an entire proces in of itself!
- 3. Make the daatd postipublic when your resarch is mad public
- Applying an emb ago can help automate this proces s

You might want to create follow up studies or news utd es where you ad to the existing data depos ti

This works nicely thanks to DOI versioning!

Lecture 🖺 4 5/6 4

Understanding the whole process (more!)

Itweeted as kng for ad ice about things I missde

- Meta d aartecord sinclud ing data
 d tionaries are really us ful tools for
 s utd to Thes to an b to book to in your articles.
- I haven't mentioned ethics forms and cons at forms Thes are very important. Refer to your Res arch Support teams for help in d & gining thes a auments.

... what about the code?

So far I've only s pken about the data element of reproducibility, we'll get onto the code in the works bp.

Digital Object Identifiers (DOI)

DOI are extremely important to ensure res arch availab lity into the future. Academic journal links are fragile and could change at any time:

s ienced riect.com/s ience/article/pii/S2665927122000879

DOI are persistent and extremely long-term id etifiers that look like this:

10.1016/j.crfs 2022.05.015

The is then res pns bi & for directing you to the res orce by creating a URL like this:

d dorg/10.1016/j.crfs 2022.05.015

Lecture 🖺 4 8/6 4

DOI for more than just publications

Initially DOI were only issued by academic publishers to resolve journal articles.

Data Repos toriess arted issuing DOI so we could resolve links to data, code and more.

Specialist Data Repositories

Sometimes you need a repos tory with spiralis features eg:

- Genomes **q** wences
- Proteins q ences
- Climated aat

Nature [39] provid san excellent overview of thes cols.

General P urpose Repositories

Thes cools all have sighty different advantages and ds a dnatages

- Figs are
- Open Science Framework
- Zenod o

Lecture 🖺 4 9/6 4

Identifiers for researchers?

DOI are great for resaing our research outputs, what about uniquely identifying researchers?

Res archers often s are the same names or change their names throughout their career.

The only open resarcher idatifier is ORCID.

Thankfully-itals oworksreally well!

It keeps track of all pub ications and deposits on datepos tories.

Here's mine: https //orcid org/0000-0002-3039-6849



I want you to regis & for an ORCID now so you can use it for everything - including in your CV

Lecture 🚊 5 0'6

Exercise: Setup a Collection on Figshare

E erything we've s pken about has been very theoretical. I want you to go through the steps of creating a collection on Figshare. We're creating a Collection because it can contain multiple Figshare items. At the beginning of your res arch you likely d of t know exactly how many data files you'll end up with.

- Sign up for Figs are with your ORCID
- Go to " Md aatin Figs bre
- Go to "Collections "
- Create a Collection
- Res eve a DOI

- Go to " Md aat
- Ad da newitem
- Res ave a DOI
- Go to your collection and add this data item.

Figshare for talks and more

Figs are is very much a general purpose repository.

If you creates mething you want to make available for the future the best thing you can do is get a

- Cons d e us rig Figs bre for presentations
- Cons d e us rig Figs bre for posters

We need to talk about Open Access

It us de to be extremely hard to access the research that UKRI funds - despite it being funded by public money. The Open Acces snovement began in the 90s and is ever growing.

It's now a requrement of UKRI [40] funding that

d winload via an online publication platform, publishers' website, or institutional or subject the final Vers on of Record or the Author's Accepted Manuscript must be free to view and repos tory within a maximum of 12 months of publication

... however, this often means that someone is paying an Article Processing Charge (APC).

Lecture 🖺 5 3/6 4

Different routes to Open Access

Someone pays an AP C

APCs can be but between publishers and authors in different ways.

UKRI is us ally res pns bi & for paying the author's portion of the bili.

- Gold Open Acces sAll articles in a journal are Open Acces s.
- Hyrb id Open Acces Specific articles in a journal are mad @pen Access through APCs Journals receive money through both s b sriptions and APCs.

No AP Gspaid

The mos Important category here is Green Open Access.

In Green Open Acces the author self-archives their article in a pub icly available repository.

This gets complicated.

- Some pub is bars req irre that only pre-prints are s barchived
- Some pub is allow post-print publishing.

If you're interes &d read Gadd and Troll Covey [41].

Lecture 🖺 5 4/6 4

Negative Results

Positive Publication Bias

There is a significant and very clear bias to researchers publishing "positive results" (Mlinari, Horvat, and Šupak Smoli [42]) - which you can even see in article titles.

This pos as gnificant is sessin the literature.

- It's really us faul to know X doesn't work! It means others d on't need to repeat the result.
- Negative effects can be under reported, as per the GSK laws it in 2004.
- A lack of negative res Lts introduces bias to meta-analys s.

The MRC [36] Open Res arch Data policy explicitly req ures bath positive and negative res utsofs utdesbandished within 24 months of the trial end.

Publishing negative results

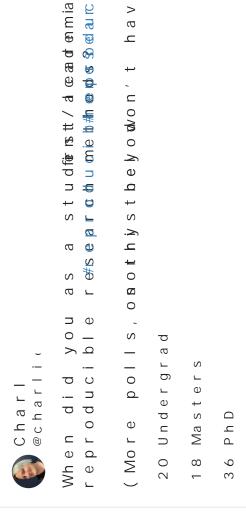
Thisisan open prob ∉m.

includ ing the Journal of Negative Results in Biomedicine The remaining journals have very low impact There was a pus flor new negative result journals in the early 2010s but several of these folded, factors as they are pub ished by smaller publishers. In general the journals with the highest impact factors are getting better at publishing negative results.

Why are we learning all of this now?

tweetrmd::tweeth_templsed/(/twitter.com/charliejhadley/status/1559 = 20&t = MOfOBvpbqLi KpDxBHi Jy7w"

2



6 Later

44 voffiersal

 \sim 202 16, A u g ∑ d 2:34

other Charlotwtee'ts See Lecture 🖺 5 8/6 4

There's a reproducibility crisis

2. Everything you can do to help reproducibility

3. Understanding Open Access and Open Data

4. Acknowledging legal and funder requirements affecting reproduciblity

Lecture 🖺 5 9/6 4

- [1] J. P. A. Ioannidis. "Why Most Published Research Findings Are False". In: PLoS Medicine 2.8 (Aug. 2005), p. e124. ISSN: 1549-1676. DOI: 10.1371 Journal.pmed.0020124.
- https://www.technologyreview.com/2021/07/30/1030329/machine-learning-ai-failed-covid-hospital-diagnosis-pandemic/.Jul. 2021. [2] W. Douglas Heaven. Hundreds of Al Tool Have Been Buil to Catch Covid. None of Them Helped.
- [3] K Broman et al. Recommendations to Funding Agencies for Supporting Reproducible Research. Tech. rep. American Statistial Association, Jan. 2017,
- Ω [4] D. Fanelli. "How Many Scientists Fabricate and Falsify Research? A Systematic Review and Meta-Analysis of Survey Data". In: PLO SO N EA. (May. 2009), p. e5738. ISSN: 1932-6203. DOI: 10.1371/journal.pone.0005738.
- [5] L. K. John et al. "Measuring the Prevalence of Questionable Research Practices With Incentives for Truth Telling". In: Psychol. dical Science 23.5 (May. 2012), pp. 524-532. ISSN: 0956-7976, 1467-9280. DOI: 10.1177/0956797611430953.
- [6] C. Aschwanden. "We're All 'P-Hacking' Now'. In: Wired (Nov. 2019).
- [7] J. P. Simmons et al. "False-Positive Psychology: Undisclosed Flexibility in Data Collection and Analysis Allows Presenting Anything as Significant". In: Psychol gical Science 22.11 (Nov. 2011), pp. 1359-1366. ISSN: 0956-7976. DOI: 10.1177/0956797611417632
- [8] Y. Hu et al. "Factors Influencing Self-Care Behaviours of Patients with Type 2 Diabetes in China Based on the Health Belief Model: A Cross-Sectional Study". In: BMJ O pre12.8 (Aug. 2022), p. e044369. ISSN: 2044-6055, 2044-6055. DOI: 10.1136/bmjopen-2020-044369.
- [9] C. Clifton et al. "Eye Movements in Reading and Information Processing: Keith Rayner's 40year Legacy". In: Journal of Memory and Language 86 (Jan. 2016), pp. 1-19. ISSN: 0749-596X. DOI: 10.1016/j.jml.2015.07.004.
- [10] Recommendations for the Conduct, Reporting, Editing, and Publication of Scholarly Work in Medical Journals (ICMJE). https://www.icmje.org/icmjerecommendations.pdf. 2022.

Lecture 🚊 6 0'6 4

https://www.wma.net/policies-post/wma-declaration-of-helsinki-ethical-principles-for-medical-research-involving-human-subjects/. 2013. [11] World Medical Association. D et ration of Heliski, Ethical Principles for Medical Research Involving Human Subjects

[12] S. H. Bradley et al. "Automatic Registration for UK Trials". In: BMJ 376 (Jan. 2022), p. 041. ISSN: 1756-1833. DOI: 10.1136/bmj.041.

[13] T. Bruckner. UK Launches New System to Achieve 100% Clinical Trial Registration (Updated). https://www.transparimed.org/single-post/ukclinical-trial-registration. Oct. 2021 [14] O. Dyer. "GlaxoSmithK Ine Faces US Lawsuit over Concealment of Trial Results". In: BMJ: British Medical Journal 328.7453 (Jun. 2004), p. 1395. ISSN: 0959-8138. [15] D. Spurgeon. "GlaxoSmithK Ine Staff Told Not to Publicise Ineffectiveness of Its Drug". In: BMJ: British Medical Journal 328.7437 (Feb. 2004), p. 422. ISSN: 0959-8138. [16] B. A. Nosek et al. "The Preregistration Revolution". In: Proceedings of the National Academy of Sciences 115.11 (Mar. 2018), pp. 2600-2606. DOI: 10.1073/pnas.1708274114.

[17] A. Dreber et al. "Using Prediction Markets to Estimate the Reproducibility of Scientific Research". In: Proceedings of the National Academy of Sciences 112.50 (Dec. 2015), pp. 15343-15347. DOI: 10.1073/pnas.1516179112.

[18] C. F. Camerer et al. "Evaluating the Replicability of Social Science Experiments in Nature and Science between 2010 and 2015". In: Nature Human Behaviour 2.9 (Sep. 2018), pp. 637-644. ISSN: 2397-3374. DOI: 10.1038/s41562-018-0399-z. [19] L. Tedersoo et al. "Data Sharing Practices and Data Availability upon Request Differ across Scientific Disciplines". In: Scientific Data 8.1 (Jul. 2021), p. 192. ISSN: 2052-4463. DOI: 10.1038/s41597-021-00981-0. [20] D. G. Roche et al. "Public Data Archiving in Ecology and Evolution: How Well Are We Doing?" In: PLO SBiol 9y 13.11 (Nov. 2015), p. e1002295. ISSN: 1545-7885. DOI: 10.1371/journal.pbio.1002295.

Lecture 🖺 6 1/6 4

[21] C. Tenopir et al. "Data Sharing by Scientists: Practices and Perceptions". In: PLO SO N IS.6 (Jun. 2011), p. e21101. ISSN: 1932-6203. DOI: 10.1371 Journal.pone.0021101 [22] G. Colavizza et al. "The Citation Advantage of Linking Publications to Research Data". In: PLO SO N B 5.4 (Apr. 2020), p. e0230416. ISSN: 1932-6203. DOI: 10.1371/journal.pone.0230416.

[23] DATA in the Time of COV $\, \mathbb{D} \,$ - 91 Nov. 2020.

[24] J. Roth et al. "Examining the Feasibility of Using Open Data to Benchmark Building Energy Usage in Cities: A Data Science and Policy Perspective". In: Energy Pol cy 1 39 (Apr. 2020), p. 111327. ISSN: 0301-4215. DOI: 10.1016/j.enpol.2020.111327.

[25] CD C- BRFSS Annual Survey D & https://www.cdc.gov/brfss/annual_data/annual_data.htm. Aug. 2021

[26] The O pred finition - O pred finition - D finition - D finition.org/. Aug. 2022.

[27] UK Government. O preGovernment Licence. https://www.nationalarchives.gov.uk/doc/open-government-licence/version/3/. 2022.

[28] World Bank. O p reD at Essential D at http://opendatatoolkit.worldbank.org/en/essentials.html. May. 2021

[29] legislation.gov.uk. D a Protection Act 2018. https://www.legislation.gov.uk/ukpga/2018/12/contents/enacted. 2022.

[30] Information Commissioner's Office. Heal htD a.https://ico.org.uk/for-organisations/guide-to-data-protection/guide-to-the-general-dataprotection-regulation-gdpr/right-of-access/health-data/. Jan. 2021. Lecture 🖺 6 2/6 4

[31] UK R. UKRI Common Principl son Research D aa. https://www.ukri.org/manage-your-award/publishing-your-research-findings/making-yourresearch-data-open/. Aug. 2022.

[32] AHRC. AHRC Research Funding Guide. https://www.ukri.org/wp-content/uploads/2022/07/AHRC-210722-ResearchFundingGuide.pdf.Aug.

[33] BBSRC. BBSRC D & Sharing Pol cy. https://www.ukri.org/wp-content/uploads/2021/07/data-sharing-policy-v1.22.pdf. Aug. 2022.

[34] EPSRC. ESPRC Pol cy Framework on Research D & Aa. https://www.ukri.org/about-us/epsrc/our-policies-and-standards/policy-framework-onresearch-data/# contents-list. Aug. 2022.

[35] ESRC. ESRC Research D & Pol ciy. https://www.ukri.org/wp-content/uploads/2021/07/ESRC-200721-ResearchDataPolicy.pdf. Aug. 2022.

[36] MRC. MRC O preResearch D at Advice. https://www.ukri.org/about-us/mrc/our-policies-and-standards/research/open-research-data-clinicaltrials-and-public-health-interventions/. Aug. 2022.

[37] NERC. NERC D & Pol cy. https://www.ukri.org/wp-content/uploads/2022/03/NERC-080322-policy-data-021219.pdf. Aug. 2022.

[38] STFC. STFC O preD at Pol cy. https://www.ukri.org/councils/stfc/guidance-for-applicants/what-to-include-in-your-proposal/datamanagement-plan. Aug. 2021.

[39] Nature. D da Repository Guidance | Scientific D da. https://www.nature.com/sdata/policies/repositories. 2022.

[40] UK R. UKRI O preAccess Pol cy. https://www.ukri.org/wp-content/uploads/2022/07/UK R-28072022-Final_UK R-0 pen-Access-Policy_ Version-1.5_July-2022.pdf. Jul. 2022. Lecture 🖺 6 3/6 4

[41] E. Gadd et al. "What Does`Green' Open Access Mean? Tracking Twelve Years of Changes to Journal Publisher Self-Archiving Policies". In: Journal of Librarianship and Information Science 51.1 (Mar. 2019), pp. 106-122. ISSN: 0961-0006. DOI: 10.1177/0961000616657406. [42] A. Mlinari et al. "Dealing with the Positive Publication Bias: Why You Should Really Publish Your Negative Results". In: Biochemia Medica 27.3 (Oct. 2017), p. 030201. ISSN: 1330-0962. DOI: 10.11613/BM.2017.030201.