

Trust in numbers

David Spiegelhalter

University of Cambridge, UK

[The address of the President delivered to The Royal Statistical Society on Wednesday, June 28th, 2017]

Summary. Those who value quantitative and scientific evidence are faced with claims both of a reproducibility crisis in scientific publication and of a post-truth society abounding in fake news and alternative facts. Both issues are of vital importance to statisticians, and both are deeply concerned with trust in expertise. By considering the ‘pipelines’ through which scientific and political evidence is propagated, I consider possible ways of improving both the trustworthiness of the statistical evidence being communicated, and the ability of audiences to assess the quality and reliability of what they are being told.

1. Introduction

This Presidential address gives me a fine opportunity to bring together two topical issues: first, the claims of a reproducibility crisis in science, which have led to concerns about the quality and reliability of at least parts of the scientific literature; second, the suggestion that we live in a ‘post-truth’ society abounding in fake news and alternative facts, in which emotional responses dominate evidence-informed judgement. These two topics have a close connection: both are associated with claims of a decrease in trust in expertise, and both concern the use of numbers and scientific evidence. They are therefore of vital importance to professional statisticians, or to any who analyse and interpret data.

A simple Internet search will reveal the daunting amount that has been written about these contested issues, and here I can only give a brief personal review of the evidence and the possible causes, focusing on the ‘filters’ that distort statistical evidence as it is passed through the information pipeline from the originators to its final consumption by the public. No single group can deal with these complex matters, but I shall argue that statisticians, and in particular the Royal Statistical Society, have an essential role both in improving the trustworthiness of statistical evidence as it flows through the pipeline, and in improving the ability of audiences to assess that trustworthiness. On statistical shoulders rests a great responsibility.

2. Reproducibility and replication

The idea of a ‘reproducibility or replication crisis’ might reasonably be said to date from John Ioannidis’s 2005 paper which notoriously proclaimed ‘Why most published research findings are false’ (Ioannidis, 2005). Although initially concerned with the biomedical literature, the idea has since been applied particularly to psychology and other social sciences. Note that, although attempts have been made to define ‘reproducibility’ and ‘replication’ precisely (Leek

Address for correspondence: David Spiegelhalter, Statistical Laboratory, Centre for Mathematical Sciences, Wilberforce Road, Cambridge, CB3 0WB, UK.
E-mail: d.spiegelhalter@statslab.cam.ac.uk

1 and Jager, 2017), I feel that we should try to avoid giving yet more technical definitions to words
 2 in routine use. (We are all familiar with misunderstandings from using ‘significance’ as a technical
 3 term, but I had always thought of ‘expected’ as fairly innocuous. That was until a journalist
 4 labelled all hospital deaths above the expected level as ‘unexpected’.) So I shall treat the terms
 5 interchangeably and distinguish when an entire study is repeated or when data is reanalysed).

6 The extent of this ‘crisis’ is contested. Ioannidis (2005) was based on modelling rather than
 7 empirical evidence: he argued that reasonable assumptions about the design of studies, biases
 8 in conduct, selection in reporting and the proportion of hypotheses investigated that were truly
 9 non-null meant a high rate of ‘false discoveries’, i.e. the proportion of published positive results
 10 that were actually null hypotheses that had been falsely rejected. In contrast, an analysis of
 11 published p -values (Jager and Leek, 2014) came up with an estimated false discovery rate of
 12 14% in the mainstream medical literature, and a recent review (Leek and Jager, 2017) concluded
 13 ‘We do not believe science is in the midst of a crisis of reproducibility, replicability, and false
 14 discoveries’.

15 So was the claim about false claims itself a false claim? This is strongly disputed by Ioannidis
 16 (2014), and a recent exercise (Szucs and Ioannidis, 2017) scraped nearly 30000 t -statistics and
 17 degrees of freedom from recent psychology and neuroscience journals, and on the basis of the
 18 observed effect sizes and low power concluded that

19 ‘Assuming a realistic range of prior probabilities for null hypotheses, false report probability is likely to
 20 exceed 50% for the whole literature’.

21
 22 Some of this apparent disagreement will be due to different literatures: Jager and Leek (2014)
 23 examined abstracts from top medical journals with many randomized controlled trials and meta-
 24 analyses, which would be expected to be much more reliable than first claims of ‘discoveries’.
 25 And even a 14% false discovery rate might be considered too high.

26 An alternative approach is purely empirical, in which the experiments behind past published
 27 claims are replicated by other teams of researchers: for example the effect of ‘power posing’,
 28 popularized in a Ted talk that has been viewed over 40 million times (Cuddy, 2012), has been
 29 subject to repeated failed replications (Ranehill *et al.*, 2015). The reproducibility project was a
 30 major exercise in which 100 psychology studies were replicated with higher power (Open Science
 31 Collaboration, 2015): whereas 97% of the original studies had statistically significant results,
 32 only 36% of the replications did. This was widely reported as meaning that the majority of the
 33 original studies were false discoveries, but Patil *et al.* (2016) pointed out that 77% of the new
 34 results lay within the 95% predictive interval from the original study, which corresponds to there
 35 not being a significant difference between the original and replication studies. This illustrates
 36 that the *difference between significant and not significant is often not significant* (Gelman and
 37 Stern, 2006). But it also means that 23% of original and replication studies had significantly
 38 different results.

39 Perhaps the main lesson is that we should stop thinking in terms of significant or not sig-
 40 nificant as determining a ‘discovery’, and instead focus on effect sizes. The reproducibility
 41 project found that replication effects were on average in the same direction as the originals but
 42 were around half their magnitude (Open Science Collaboration, 2015). This clearly displays the
 43 biased nature of published estimates in their literature, and strong evidence for what might
 44 be termed regression to the null.

45 46 47 **3. What is the cause of this ‘crisis’?**

48 It is important to note that deliberate fabrications of data do occur but appear relatively rare. A

1 review estimated that 2% of scientists admitted falsification of data (Fanelli, 2009), and the US
 2 National Science Foundation and Office of Research Integrity deal with a fairly small number
 3 of deliberately dishonest acts (Mervis, 2017), although substantial numbers of cases must go
 4 undetected as it is generally difficult to check raw material. Computational errors are more
 5 common but can be detected by repeating analyses if the original data are available.

6 Rather than deliberate dishonesty or computational incompetence, the main blame has been
 7 firmly placed on a ‘failure to adhere to good scientific practice and the desperation to publish
 8 or perish’ (Begley and Ioannidis, 2015). The crucial issue is the quality of what is submitted to
 9 journals, and the quality of what is accepted, and deficits are a product of what have become
 10 known as ‘questionable research practices’.

11 Fig. 1 shows the results of a survey of academic psychologists in the USA, which had a 36%
 12 response rate (John *et al.*, 2012). A very low proportion admitted falsification, but other practices
 13 that can severely bias outcomes were not only frequently acknowledged but also generally seen
 14 as defensible: for example the 50% who admitted selectively reporting studies gave an average
 15 score of 1.66 when asked whether this practice was defensible, where 0 \equiv no, 1 \equiv possibly and
 16 2 \equiv yes. An Italian survey found similar rates, although the respondents were more inclined to
 17 agree that the practices were not defensible (Agnoli, *et al.*, 2017).

18 These questionable research practices just involve experimentation. If we consider general
 19 observational biomedical studies and surveys, then there is a vast range of additional potential
 20 source of bias: these might include

- 21 (a) sampling things that are convenient rather than appropriate,
- 22 (b) leading questions or misleading wording,
- 23 (c) inability to adjust properly for confounders and to make fair comparisons,
- 24 (d) too small a sample,
- 25 (e) inappropriate assumptions in a model and
- 26 (f) inappropriate statistical analysis.

27 And to these we might add many additional questionable practices concerned with interpretation
 28 and communication, which we shall return to later.

29 These are not just technical issues of, say, lack of adjustment of p -values for multiple testing.
 30 Many of the problems arise through more informal choices made throughout the research
 31 process in response to the data, say in selecting the measures to emphasize, choice of adjusting
 32 variables, cut points to categorize continuous quantities and so on: this has been described as
 33 the ‘garden of forking paths’ (Gelman and Loken, 2014) or ‘researcher degrees of freedom’
 34 (Simmons *et al.*, 2011) and will often take place with no awareness that these are questionable
 35 research practices.

36 There have been strong arguments that the cult of p -values is fundamental to problems of
 37 reproducibility, and recent guidance from the American Statistical Association clearly revealed
 38 their misuse (Wasserstein and Lazar, 2016). Discussants called for their replacement or at least
 39 downplaying their pivotal role in delineating ‘discoveries’ through the use of arbitrary thresh-
 40 olds. We have already seen that p -values are fragile things that need to be handled carefully in
 41 replication studies — for example a study with $p = 0.05$ would only be predicted a 50% chance
 42 of obtaining $p < 0.05$ in a precise replication.

43 This is a complex issue, and in a recent in article *Significance* (Matthews *et al.*, 2017) I
 44 confessed that I liked p -values, and that they are good and useful measures of compatibility
 45 between data and hypotheses, but insufficient distinction is made between their informal use
 46 in exploratory analysis and their more formal use in confirmatory analyses that summarize the
 47 totality of evidence—perhaps they should be distinguished as p_{exp} and p_{con} .

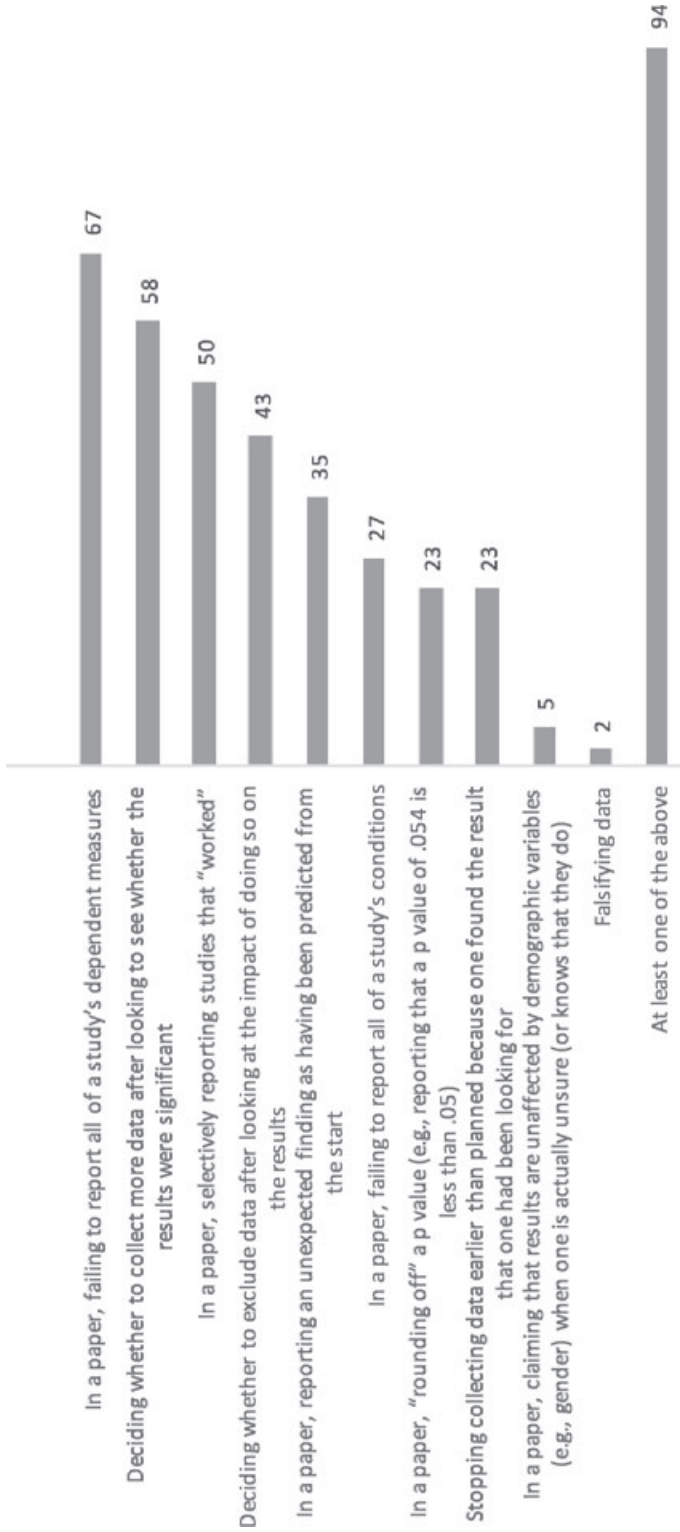


Fig. 1. Questionable research practices admitted by 2155 US academic psychologists (John et al., 2012)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

Essentially there is too strong a tendency to use p -values to jump from selected data, to a claim about the strength of evidence and to conclusions about the practical importance of the research. p -values ‘do what they say on the tin’, but people do not read the tin.

4. What gets into the scientific literature?

Questionable practices influence what is submitted to the scientific literature, and what finally appears depends on the publisher’s ability to critique and select from what is presented to them. Ideally, peer review would weed out inadequate research and reporting, and recommend publication of good science regardless of the actual results. But we know that peer review is often inadequate, and there is an urge for the leading journals, to a varying amount across different disciplines, to publish newsworthy positive ‘discoveries’ and hence they produce a skewed resource.

We should not be surprised at this, since traditionally journals were set up to report new findings rather than the totality of evidence. Now there is an explosion in the amount of research and publishing opportunities, I would agree that

‘most scientific papers have a lot more noise than is usually believed, that statistically significant results go in the wrong direction far more than 5% of the time, and that most published claims are overestimated, sometimes by a lot’

(Gelman, 2013). Although Gelman adds, more positively, that even though there are identifiable problems with individual papers, areas of science could still be moving generally in the right direction.

So what can be done? A group of prominent researchers recently published a ‘manifesto for reproducible science’ whose recommendations for action are summarized in Table 1, together with the relevant stakeholders (Munafò *et al.*, 2017).

The list in Table 1 demonstrates the responsibility of a wide range of stakeholders and is firm in its commitment to transparency. Statistical science has a major role in many of these proposals, in particular in methodological training, improving reporting and peer review, and the sharing of data for reanalysis. However, the one important element that seems to missing from Table 1 is the need for external commentary, critique and ‘calling out’ of poor practice, which are the responsibility of the entire scientific community, the media and the public — we shall return to this theme later.

In spite of extensive discussion about problems in the reliability of published science, these

Table 1. Proposals from the ‘manifesto for reproducible science’ (Munafò *et al.*, 2017)

<i>Theme</i>	<i>Proposal</i>	<i>Stakeholders</i>
Methods	Protecting against cognitive biases	Journals, funders
	Improving methodological training	Funders, institutions
	Independent methodological support	Funders
	Collaboration and team science	Funders, institutions
Reporting and dissemination	Promoting study pre-registration	Journals, funders
	Improving the quality of reporting	Journals
	Protecting against conflicts of interest	Journals
Reproducibility	Encouraging transparency and open science	Journals, funders, regulators
	Evaluation	Diversifying peer review
Incentives	Rewarding open and reproducible practices	Journals, funders, institutions

Table 2. Who is trusted as a source of medical research information?: responses from 1500 UK adults (Wellcome Trust, 2017)

<i>Profession</i>	<i>Trust completely or largely (%)</i>	<i>Trust very little or not at all (%)</i>
Doctors or nurses	64	6
University scientists	59	4
Medical research charities	37	11
Pharma scientists	32	16
Industry scientists	29	16
Journalists	3	59

concerns do not seem to have fed into public opinion yet. A recent survey (Wellcome Trust, 2017) revealed the trust ratings shown in Table 2.

It is ironic that pharmaceutical scientists are given low levels of trust, although they work under far greater constraints than university scientists in terms of prespecification of design and analyses for regulators, and arguably produce more trustworthy analyses (personally I would trust the opinion of pharma statisticians on medical research far more than I would many health professionals). Journalists are given very low trust ratings in spite of being a major source of information to the public.

This introduces the idea that expressions of trust are not, in general, based on careful consideration of evidence but arise as an immediate response based on our gut feelings, which brings us naturally to the way that we handle all the other numbers that deluge us as part of daily life, and in particular those that appear in the news.

5. Numbers in the news

Scientists are not the only people reporting claims based on statistical evidence. Politicians, non-governmental organizations and many other bodies are all competing for our attention, using numbers and science to provide an apparently ‘objective’ basis for their assertions. Technology has changed, encouraging an increasing diversity of sources to use on line and social media to communicate, with few controls to ensure reliable use of evidence. This has led to suggestions that we are in a time of populist political discourse in which appeals to our emotion triumph over reason.

At the extreme, there are completely fabricated, demonstrably false facts that can be honestly labelled ‘fake news’. But, as with science, I believe that deliberate fabrication is not the main issue: this will be better dealt with in the future by a combination of calling-out by fact checking organizations such as Full Fact (Full Fact, 2017), crowd sourcing on social media, automatic algorithms and possible regulation of social media sites: for example, Full Fact covered the recent election campaign in the *Evening Standard* as well as collaborating with Facebook on prominent advertising of ‘Tips for spotting false news’.

As with science, a much bigger risk is manipulation and exaggeration through inappropriate interpretation of ‘facts’ that may be technically correct but are distorted by what we might call ‘questionable interpretation and communication practices’. Fig. 2 provides a highly simplified view of the process by which we hear about statistical evidence as the end of a pipeline that starts with the originators of the data, and then goes through the ‘authorities’, then through their press and communication offices to the traditional media, and finally to us as individual members of society.

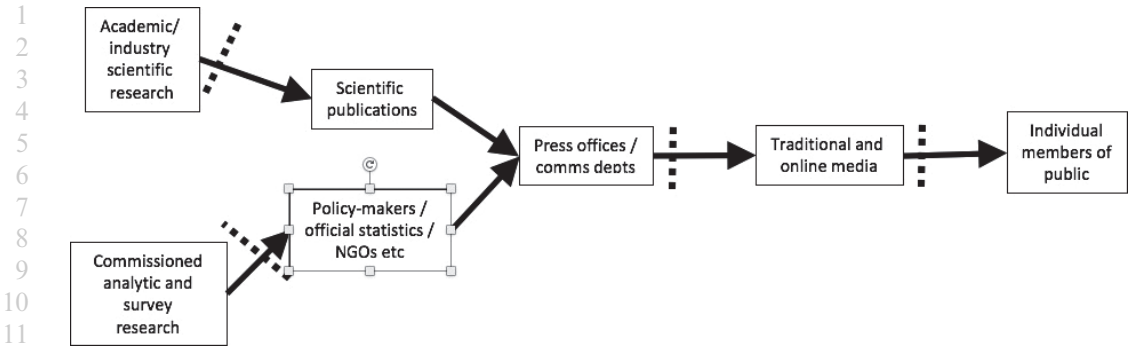


Fig. 2. Simplistic diagram of the traditional information flows from statistical sources through to the public: . . . , filters arising from questionable research, interpretation and communication practices, such as selective reporting, lack of context and exaggeration of importance

Table 3. Some highly questionable interpretation and communication practices

<ul style="list-style-type: none"> Pick stories that go against current consensus Promote stories regardless of research quality Do not report uncertainties Do not provide context or comparative perspective, such as a time series Suggest a cause when only an association is observed Exaggerate relevance and importance of findings Claim that the evidence supports a particular policy Provide only relative and not absolute risks Use positive or negative framing depending on whether the aim is to reassure or frighten Do not dwell on conflicts of interest or alternative views Use a sexy but uninformative graphic Write a headline which may have little connection to the story but will encourage clicks.

The questionable practices that have been adopted by some of the more ruthless press offices, communications teams and journalists include those in Table 3.

Many of these would be seen as defensible by those whose professional career depends on attracting readers, listeners or clicks, and it would be very interesting to conduct a survey of press officers and journalists to see how many of these they had used. But it is important to note that scientists might use these as well—in my personal experience some can, in spite of their proclaimed *caveats*, be too quick to jump to the wider implications of their work. When communicating subtle statistical evidence there seems an irresistible tendency to produce a simplifying narrative; we have seen *X*, it is because of *Y*, and so we should do *Z*.

This pipeline suggests that it is too easy to blame journalists for misreporting science. Press offices, and the journals and scientists themselves, can be to blame: a recent study found that, of 462 press releases from UK universities in 2011, 40% contained exaggerated advice, 33% contained exaggerated causal claims, 36% contained exaggerated inference to humans from animal research and the majority of exaggerations appearing in the press could be traced back to the press release (Sumner *et al.*, 2014). The same team found slightly more reassuring results in 534 press releases from major biomedical journals: causal claims or advice in a paper were exaggerated in 21% of corresponding press releases, although these exaggerations, which tended to be reported, did not produce more press coverage (Sumner *et al.*, 2016).

1 One of my favourite examples of imaginative storytelling by a communications team is when
 2 a rather dull study which found that 10% of people carried a gene which protected them against
 3 high blood pressure (Newton-Cheh *et al.*, 2009), which was reframed negatively as nine in
 4 10 people carrying a gene which *increases* the risk of high blood pressure: this duly received
 5 international press coverage (Devlin, 2009).

6 Another recent classic was a careful Swedish population study whose published abstract
 7 (Khanolkar *et al.*, 2016) said ‘We observed consistent associations between higher socio-
 8 economic position and higher risk of glioma’; the press release headlined with ‘High levels
 9 of education linked to heightened brain tumour risk’ (Medical Xpress, 2016) and the subeditor
 10 of the *Daily Mirror* finally turned it into ‘Why going to university increases risk of getting a
 11 brain tumour’ (Gregory, 2017).

12 The use of relative risks without absolute risks is a standard complaint and is explicitly
 13 warned against in British Broadcasting Corporation statistical guidelines (British Broadcasting
 14 Corporation, 2017). It is known that relative risks, which are often referred to by the media
 15 as simply an ‘increased risk’ regardless of magnitude, are an effective way of making a story
 16 look more exciting, and this is not helped by the fact that odds, rate and hazard ratios are the
 17 standard output from most biomedical studies. The gripping headline ‘Why binge watching your
 18 TV box sets could kill you’ (Donnelly, 2016) arose from an epidemiological study that estimated
 19 an adjusted hazard ratio of 2.5 for a fatal pulmonary embolism associated with watching more
 20 than 5 h of television a night compared with less than 2.5 hours (Shirakawa *et al.*, 2016). But
 21 careful scrutiny of the absolute rate in the high risk group (13 in 158000 person-years) could be
 22 translated as meaning that you can expect to watch more than 5 h of television a night for 12000
 23 years before experiencing the event, which somewhat lessens the impact. The newspaper article
 24 was, as usual, much better than the headline, but whose fault is it not to insist on including this
 25 *perspective*: the journalist, the press office, the publication or the scientists?

26 I am unsure whether this misuse of statistical evidence is becoming worse. There are certainly
 27 more outlets promoting partial news, but mainstream Web sites, newspapers and television are
 28 perhaps under more scrutiny. What would fact checkers have found were they active say 30 years
 29 ago? My unsupported feeling is that it would have been even worse than now.

30 But in my darkest moods I follow what could be called the ‘Groucho principle’: because
 31 stories have gone through so many filters that encourage distortion and selection, the very fact
 32 that I am hearing a claim based on statistics is reason to disbelieve it.

33 6. Trust in expertise

34
 35
 36 When faced with stories about how the natural world or society works, we can rarely check them
 37 for ourselves. So *trust* is an inevitable element in dealing with statistical evidence, and therefore
 38 recent claims that there has been a decrease in trust in expertise is worth serious attention.

39 This claim is often associated with Michael Gove in the British exist campaign saying that
 40 people having had enough of experts, but it is important to quote him in full:

41 ‘people have had enough of experts from organisations with acronyms saying that they know what is
 42 best and getting it consistently wrong’
 43

44 (Youtube, 2016). This sounds a little more reasonable and reflects recent high-profile failures
 45 to predict and control financial markets, and the lamentable quality of many political forecasts
 46 identified by the ‘Good judgement project’ (Tetlock and Gardner, 2015).

47 The evidence for such a decline is mixed. The Edelman trust barometer claims that trust is ‘in
 48 crisis’, and their poll shows that a ‘person like yourself’ is now as credible as a technical expert,

1 and yet their data show an overall increase in trust in government, the media, business and
 2 non-governmental organizations since 2012 (Edelman, 2017). A recent YouGov poll showed
 3 scientists trusted by 71%, although this is 63% versus 83% depending on whether voting to leave
 4 or remain in the European Union, and scientists come fourth in a UK trust league table behind
 5 nurses, doctors, and your own general practitioner (YouGov, 2017). Levels of trust in official
 6 statistics remain high and have increased (National Centre for Social Research, 2017): of those
 7 able to give an opinion in 2016,

- 8 (a) 90% trust the Office for National Statistics,
- 9 (b) 85% trust the statistics produced by the Office for National Statistics,
- 10 (c) 78% agree that official figures are accurate,
- 11 (d) 26% agree that the UK Government presents official figures honestly and
- 12 (e) 18% agree that newspapers present official figures honestly.

13 These figures look reassuring to official statistics, although not to government or the media, but
 14 might improve further were prerelease access of politicians and their advisors to official statistics
 15 abolished, which is an objective of a continuing campaign by the Royal Statistical Society(RSS).
 16 (A letter on this topic was published in *The Times* during the recent election campaign, with
 17 114 signatories (including Baroness Onora O’Neill), reflecting the 114 people with prerelease
 18 access to labour market statistics.)

19 In her Reith lectures, philosopher Onora O’Neill pointed out the undifferentiated nature
 20 of these questions about whom we trust, and perhaps reflect a more general mood of suspi-
 21 cion (O’Neill, 2002). More important is *active* trust, judged by our actions, which display that
 22 we commonly put our trust in institutions that we profess not to trust. Crucially, she went
 23 on to say that nobody can just expect to be trusted—they need to demonstrate *trustworthi-*
 24 *ness*.

25
 26
 27 **7. Improving trustworthiness**

28 In a remarkable Ted talk (O’Neill, 2013), Onora O’Neill argued that, rather than aiming to
 29 build trust, the task should be to become trustworthy—this means demonstrating competence,
 30 honesty and reliability. But you also have to provide usable evidence that allows others to check
 31 whether you are trustworthy, which necessitates making yourself vulnerable to the other party.
 32 Although identifying *deception* as being a key breaker of trust, she emphasized the danger of
 33 too much focus on policing deliberate fraud—this is too low a bar. This reinforces the need to
 34 avoid excessive attention to deliberate dishonesty, say through data fabrication or demonstrably
 35 fake news, because it could distract attention from the more pressing problem of misleading,
 36 incompetent and unreliable use of evidence.

37 There seem to be three main ways of building more trustworthiness in the statistical evi-
 38 dence pipeline shown in Fig. 2: change the communication structure, improve the filters for the
 39 information being passed and improve the ability of audiences to check trustworthiness.

40
 41
 42 **7.1. Changing the communication structure**

43 There are increasing possibilities to bypass potentially distorting filters. These include direct-
 44 to-public communication through social media by scientists, agencies, statistical ‘experts’, and
 45 even US Presidents. Although these innovations open up exciting opportunities for direct com-
 46 munication, there is also the risk of bypassing filters that have a positive role in weeding out
 47 poor science and statistics, and this emphasizes even more the need for audiences to be able to
 48 appraise the reliability of what is being claimed.

7.2. Improving the filters

We have already seen the proposals in Table 1 for improving the reproducibility, and hence the trustworthiness, of published science. Many of these are concerned with transparency, but O’Neill has observed that transparency does not necessarily avoid deception (O’Neill, 2002). She recommended ‘intelligent transparency’, which requires information to be ‘accessible, intelligible, assessable, and useable’ (Royal Society, 2012). The crucial element is that audiences need to be able to inquire and not just to accept assurances on trust.

Many of the ideas that are listed in Table 1 would also serve to improve trustworthiness in evidence in general, such as training of professionals, improved reporting standards, openness and protection against conflicts of interest. Other measures that could enhance the reputation of scientific and statistical expertise might include clear demonstration of the following factors.

- (a) *Uncertainty*: many have recommended a greater willingness to embrace uncertainty (Makri, 2017) and to display humility (Shafik, 2017). I strongly concur, but I would add that this does not mean a reluctance to speak out confidently when faced with clear false statements or beliefs: perhaps we need a form of *muscular uncertainty*.
- (b) *Engagement*: it seems essential to have empathy with audiences, and in particular understanding their beliefs and concerns. As we shall see below, this can also allow some pre-emption of misunderstandings.
- (c) *Impartiality*: trustworthiness can be demonstrated by meticulous avoidance of broader agendas, so that there is a clear demarcation between the description of the evidence and any potential policy recommendations. If scientists and statisticians are seen as advocates, then they must expect their objectivity to be questioned.

This final point is especially relevant to the history of the Royal Statistical Society, whose founding principles in 1834 included the pious assertion that

‘The Statistical Society will consider it to be the first and most essential rule of its conduct to exclude carefully all opinions from its transactions and publications – to confine its attention rigorously to facts — and, as far as it may be found possible, to facts which can be stated numerically and arranged in tables.’

This ‘essential rule’ was immediately ignored by Fellows who made bold recommendations on the basis of flimsy data. Even contemporary commentators commented on the ambiguity of the term ‘facts’ (McConway, 2016), with its implication of unquestionable veracity and authority, whereas data do not exist in a vacuum and only acquire meaning and value in a context.

This means acknowledging that numbers do not speak for themselves and so entails a responsibility to provide interpretation and potential implications of data, but without slipping into advocacy or suggesting that the evidence mandates a particular decision without taking account more general societal values. The current RSS strapline — ‘Data, evidence, decisions’ — explicitly recognizes the role of statistical science at all stages of this path: for example the RSS’s data manifesto encourages the publication of the evidence behind policies (Royal Statistical Society, 2016) and so to ‘show the working’. Interpretation can be provided in a clearly separate section of a report, as practised by the Office for National Statistics.

When it comes to the specific outputs from press offices and the media, the primary aim should be to avoid the sort of questionable communication practices that are listed in Table 3. This might be helped by the following strategies:

- (a) Construction and adoption of reporting guidelines, such as the simple list commissioned by the Levesen Inquiry (Fox, 2012). (the RSS made a major contribution to revised BBC

- 1 guidelines on reporting statistics (British Broadcasting Corporation, 2017), and the BBC
 2 is also investing in data journalism—other broadcasters might follow their example);
 3 (b) establishing close links between statisticians and journalists, although this is not without
 4 problems (McConway, 2016), and journalism training, such as carried out by the RSS;
 5 (c) working with dedicated organizations such as the Science Media Centre (Science Media
 6 Centre, 2017) and Sense about Science (Sense about Science, 2017);
 7 (d) encouraging good storytelling with data, with appropriate and attractive narratives and
 8 visualization.

9
 10 Although there are many exhortations to turn numbers into stories, the process does carry
 11 risks. Stories need an arc and a well-rounded conclusion, which science rarely provides, and
 12 so it is tempting to oversimplify and overclaim. We need to encourage stories that are true to
 13 the evidence: its strengths, weaknesses and its uncertainties. We need, for example, to be able
 14 to say that a drug or another medical intervention is neither good nor bad, it has benefits and
 15 harms, that people might weigh them up in different ways and quite reasonably come to different
 16 conclusions. Journalists seem to shy away from such nuanced narratives but, say by including
 17 testimony from people with differing views, a good communicator should be able to make these
 18 stories gripping.

19 As an apparently rare example of such a story, Christie Aschwanden from FiveThirtyEight
 20 discussed the statistics about breast screening, and then said that she had decided to avoid the
 21 procedure, whereas her smart friend, provided with the same evidence, had made the oppo-
 22 site decision (Aschwanden, 2015). This neatly asserts the importance of personal values and
 23 concerns, while still respecting the statistical evidence.

24 But it is not enough simply to invent lists of things that could be done to improve communi-
 25 cation—we need active research into how best to do them. For example, how can we best
 26 communicate uncertainty about facts and the future without jeopardizing trust and credibility,
 27 and how can our techniques be tailored to audiences with different attitudes and knowledge? In
 28 addition, there seems a remarkable lack of research into different ways of communicating how
 29 policy decisions are expected to impact society.

30
 31 *7.3. Trust as feeling and trust as analysis*

32 The concept of a ‘dual process’ in psychology has been popularized by Kahneman’s image of
 33 thinking fast or slow: a rapid automatic non-conscious system 1, and a more considered con-
 34 scious system 2 (Kahneman, 2011). This idea has proved useful in examining different attitudes
 35 to risk: Slovic *et al.* (2004) distinguished ‘risk as feelings’, our immediate personal reactions to
 36 perceived threats, from ‘risk as analysis’, the analytic approach that is more familiar to statisti-
 37 cians and actuaries.

38 Trust might be approached similarly. When we are the recipient of a claim, trust is generally
 39 viewed as a ‘feeling’, a product as much of whether we warm to the ‘expert’ than by careful
 40 consideration of what is said: we have seen how pharma scientists suffer mistrust through broad
 41 suspicion of the industry. This is often a good heuristic, but like all heuristics it can be gamed
 42 by manipulative persuaders. In the spirit of Kahneman, we might distinguish ‘fast trust’ and
 43 ‘slow trust’, with fast trust dominated by our feelings about the topic and whether we feel that
 44 the source shares our values and has our interests at heart. Slow trust is based on the type of
 45 considered inquiry and analysis that was encouraged by O’Neill.

46 But is it possible to move people from ‘trust as feeling’ to ‘trust as analysis’? Can people be
 47 ‘reasoned’ out of gut feelings, when they have an emotional investment in an opinion and their
 48 ‘motivated reasoning’ means that they are not shifted by evidence? This is not a new debate: in

1 1682 the English poet John Dryden optimistically claimed ‘A Man is to be cheated into Passion,
 2 but to be reason’d into Truth’ (Dryden, 1682), but in 1721 Jonathan Swift presented a directly
 3 opposing view: ‘Reasoning will never make a Man correct an ill Opinion, which by Reasoning
 4 he never acquired’ (Swift, 1843).

5 An active area of research focuses on whether people’s demonstrably inaccurate opinions can
 6 be corrected through provision of evidence. There are many studies of the ‘backfire’ effect, a form
 7 of confirmation bias, which says that simply correcting ‘fake news’ can end up reinforcing the
 8 very belief that is being countered. However, there is increasing evidence that misconceptions
 9 can to some extent be overcome by persuasive information (Spinney, 2017), studies show it
 10 is possible to protect pre-emptively (‘inoculate’) public attitudes about climate change against
 11 real world misinformation (van der Linden *et al.*, 2017; Cook, *et al.*, 2017) and that good
 12 visualizations can improve “immunity to misleading anecdote” (Fagerlin *et al.*, 2005). People
 13 do not like to be deceived.

14 7.4. Improving the assessment of trustworthiness

15 There appear to be two main ways of ensuring that trustworthiness can be properly assessed:
 16 training audiences in critical appraisal, and encouraging platforms dedicated to response and
 17 ‘calling-out’. Possible routes are shown in Fig. 3.

18 Although each of the groups of ‘assessors’ will have different capacities and interests, rather
 19 similar principles should apply to whoever is considering the trustworthiness of statistical ev-
 20 idence, whether it is policy professionals critiquing the impact assessments provided by their
 21 analysts, or patients confronted by information leaflets. For example, Stempra’s (2017) ‘Guide
 22 to being a press officer’ emphasizes the need to be clear about the limitations of the study, Sense
 23 about Science (2017) have aimed directly at the public with their ‘Ask for Evidence’ campaign
 24 and recent randomized trials in Africa have shown that families can be taught, by using comic
 25 books and audio lessons, to question claims made about medical treatments (Informed Health
 26 Choices, 2017).

27 Training can involve the development of teaching and assessment material in critical ap-
 28 praisal, provision of checklists and awareness of tricks illustrated with gripping examples that
 29 are relevant to the specific audience. I have already mentioned Facebook and Full Fact’s check-
 30 list for detecting ‘false news’, and I am pleased that the RSS is active in creating a more general
 31 list of questions that can be adapted to specific circumstances. Three aspects of a story can be
 32 critiqued:

- 33 (a) questions about the *research*—i.e. the trustworthiness of the number itself (‘internal
 34 validity’);
- 35 (b) questions about the *interpretation*—i.e. the trustworthiness of the conclusions drawn
 36 (external validity);
- 37 (c) questions about the *communication*—i.e. the trustworthiness of the source and what
 38 we are being told (‘spin’).

39 Fig. 3 shows that fact checkers, blogs and official watchdogs such as the UK Statistics Au-
 40 thority can all publicly ‘name and shame’ bad practice in the use of statistics. In contrast, the
 41 corresponding opportunity for the scientific community to comment on publications is mainly
 42 limited to a myriad of personal blogs, because of the rather dysfunctional publication model
 43 that does not encourage even a moderated on-line discussion forum, even though the RSS
 44 has had published commentaries on papers for nearly two centuries. Short of retraction, there
 45 still seem to be few penalties for scientists indulging in questionable practices or slipping into
 46 advocacy.

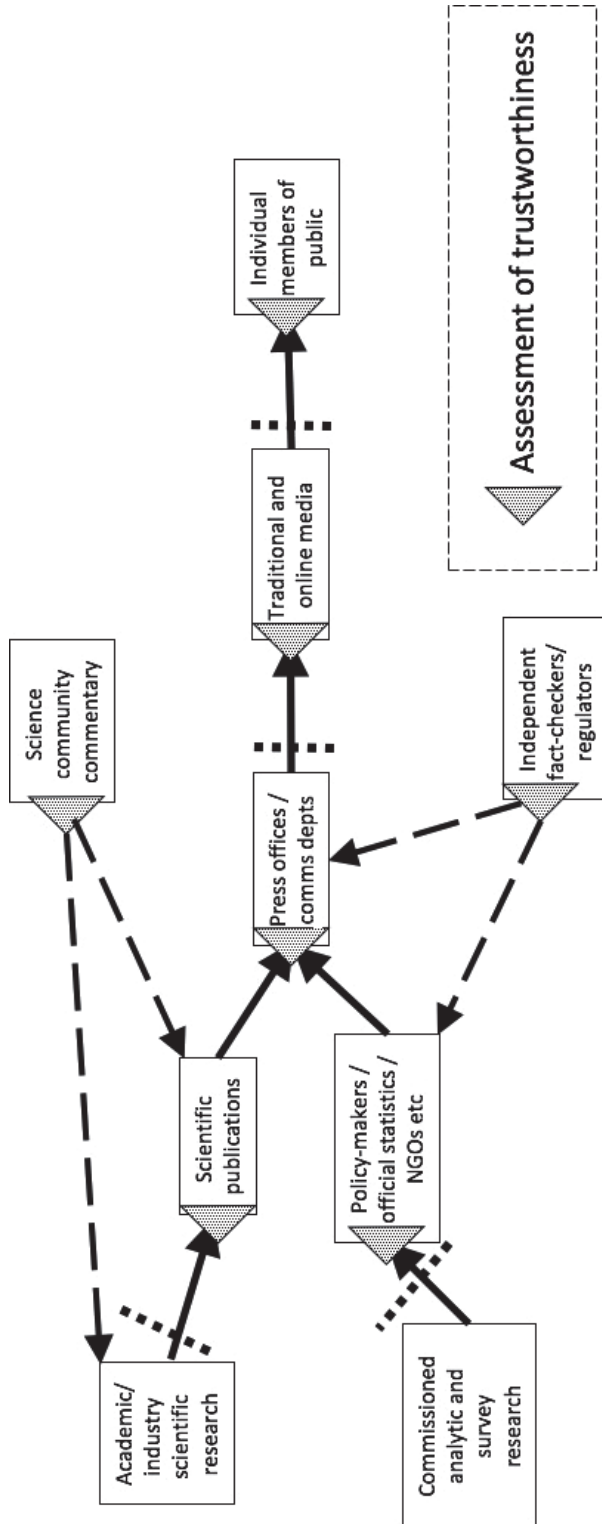


Fig. 3. Potential for the assessment of trustworthiness of statistical evidence

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

1 The first step in good communication is to shut up and to listen, and the flow of trustworthy
 2 evidence will only improve when providers are aware that at least part of their audience is
 3 carefully monitoring the quality of what is delivered and will publicly call them out if they
 4 deviate too much from competence, honesty and reliability.

6 8. Conclusions

8 I have tried to bring together two related issues—lack of scientific reproducibility and dubious
 9 numbers in the news—by framing them both as threats to trust, which should be countered by
 10 improving trustworthiness. The pipeline through which we receive information, complete with a
 11 series of filters, applies to both contexts, and the possible measures to improve trustworthiness of
 12 what is published in both the scientific and the general media have much in common. Audiences
 13 need to be able to assess trustworthiness, and again the measures to improve their ability to do
 14 so are similar in both scientific and general media.

15 These are complex issues with many actors, of which the RSS is just one player, and we
 16 are fortunate that the UK has an active and collaborative ecology of organizations who are
 17 trying to improve the reliability of our science publications and the ‘factfulness’ of the media.
 18 Again, the RSS strapline: ‘*Data, evidence, decisions*’, has never been so pertinent, and it is
 19 a noble role to negotiate the delicate steps along that process. My personal heuristic is that
 20 statisticians are a trustworthy bunch, good and conscientious at their job, if a little nerdy. I believe
 21 that they should have a higher profile in promoting impartial evidence, and this means that at
 22 least some need to become better at converting their insights into accessible (and trustworthy)
 23 stories.

24 Of course this endeavour is not restricted to those who would label themselves professional
 25 statisticians and join the RSS. Hans Rosling, master statistical storyteller, was a public health
 26 physician, and there is a growing and vibrant community of people who analyse and communi-
 27 cate data. Those who use, and misuse, statistics come from a wide variety of backgrounds, and
 28 so the aim must be to promote the trustworthiness of statistically based claims and decisions,
 29 not just the trustworthiness of statisticians. Nevertheless, when making such an assessment it
 30 may be reasonable to take into account the professionalism of the source.

31 I hope that the RSS will continue to be at the forefront of both improving the trustworthiness
 32 of the numbers that are used in society, and the ability of audiences to assess that trustworthiness.

34 Acknowledgements

36 I am indebted to many for comments and encouragement on this diatribe, in particular Kevin
 37 McConway, Alex Freeman, Michael Blastland, Theresa Marteau and Iain Wilton. And I could
 38 not be working in this area at all without the unquestioning support of David Harding of Winton
 39 Capital Management.

41 References

- 43 Agnoli, F., Wicherts, J. M., Veldkamp, C. L. S., Albiero, P. and Cubelli, R. (2017) Questionable research practices
 44 among Italian research psychologists. *PLOS ONE*, 12, no. 3, article e0172792.
 45 Aschwanden, C. (2015). Science won't settle the mammogram debate. *Five Thirty Eight*, Oct. 20th.
 46 Begley, C. G. and Ioannidis, J. P. A. (2015) Reproducibility in science: improving the standard for basic and
 47 preclinical research. *Circlin Res.*, **116**, 116–126.
 48 British Broadcasting Corporation (2017) Reporting statistics. British Broadcasting Corporation, London. (Available from <http://downloads.bbc.co.uk/rmhttp/guidelines/editorialguidelines/pdfs/ReportingStatistics.pdf>.)

- 1 Cook, J., Lewandowsky, S. and Ecker, U. K. H. (2017) Neutralizing misinformation through inoculation: Exposing
 2 misleading argumentation techniques reduces their influence. *PLOS ONE*, **12**, no. 5, article e0175799.
- 3 Cuddy, A. (2012) Your body language shapes who you are. (Available from [http://www.ted.com/talks/
 4 amy_cuddy_your_body_language_shapes_who_you_are](http://www.ted.com/talks/amy_cuddy_your_body_language_shapes_who_you_are).)
- 5 Devlin, K. (2009) Nine in 10 people carry gene which increases chance of high blood pressure. *Telegraph*, Feb.
 6 15th.
- 7 Donnelly, L. (2016) Why binge watching your TV box-sets could kill you. *Telegraph*, July 25th.
- 8 Dryden, J. (1682) Religio laici, or, A laymans faith a poem. (Available [http://name.umdl.umich.edu/
 9 A36673.0001.001](http://name.umdl.umich.edu/A36673.0001.001).)
- 10 Edelman, (2017) Trust barometer. (Available from <http://www.edelman.com/trust2017/>.)
- 11 Fagerlin, A., Wang, C. and Ubel, P. A. (2005) Reducing the influence of anecdotal reasoning on people's health
 12 care decisions: is a picture worth a thousand statistics? *Med. Decsion Makng*, **25**, 398–405.
- 13 Fanelli, D. (2009) How many scientists fabricate and falsify research?: a stematic review and meta-analysis of
 14 survey data. *PLOS ONE*, **4**, no. 5, artical e5738.
- 15 Fox, F. (2012) 10 best practice guidelines for reporting science & health stories. (Available [http://webar
 16 chive.nationalarchives.gov.uk/20140122145147/http://www.levesoninquiry.org.uk/
 17 wp-content/uploads/2012/07/Second-Submission-to-inquiry-Guidelines-for-Science
 18 and-Health-Reporting.pdf](http://webarhive.nationalarchives.gov.uk/20140122145147/http://www.levesoninquiry.org.uk/wp-content/uploads/2012/07/Second-Submission-to-inquiry-Guidelines-for-Science-and-Health-Reporting.pdf).)
- 19 Full Fact (2017) Full Fact is the UK's independent factchecking organisation. Full Fact, London.
- 20 Gelman, A. (2013) Difficulties in making inferences about scientific truth from distributions of published p-values.
 21 (Available from [http://andrewgelman.com/2013/09/26/difficulties-in-making-inferen
 22 ces-about-scientific-truth-from-distributions-of-published-p-values/](http://andrewgelman.com/2013/09/26/difficulties-in-making-inferences-about-scientific-truth-from-distributions-of-published-p-values/).)
- 23 Gelman, A., and Loken, E. (2014) The statistical crisis in science. *Am. Scient.*, **102**, 460.
- 24 Gelman, A. and Stern, H. (2006) The difference between “significant” and “not significant” is not itself statistically
 25 significant. *Am. Statistn*, **60**, 328–331.
- 26 Gregory, A. (2017) Why going to university increases risk of getting a brain tumour. *Mirror Online*, June 20th.
- 27 Informed Health Choices (2017) Using evidence to change the world. (Available from [http://www.
 28 informedhealthchoices.org/](http://www.informedhealthchoices.org/).)
- 29 Ioannidis, J. (2014). Discussion: Why ‘An estimate of the science-wise false discovery rate and application to the
 30 top medical literature’ is false. *Biostatistics*, **15**, 28–36.
- 31 Ioannidis, J. P. A. (2005) Why most published research findings are false. *PLOS Med*, **2**, no. 8, article e124.
- 32 Jager, L. R. and Leek, J. T. (2014) An estimate of the science-wise false discovery rate and application to the top
 33 medical literature. *Biostatistics*, **15**, 1–12.
- 34 John, L. K., Loewenstein, G. and Prelec, D. (2012) Measuring the prevalence of questionable research practices
 35 with incentives for truth telling. *Psychol. Sci.*, **23**, 524–532.
- 36 Kahneman, D. (2011) *Thinking, Fast and Slow*. New york: Farrar, Straus and Giroux.
- 37 Khanolkar, A. R., Ljung, R., Talbäck, M., Brooke, H. L., Carlsson, S., Mathiesen, T. and Feychting, M. (2016)
 38 Socioeconomic position and the risk of brain tumour: a Swedish national population-based cohort study.
 39 *J. Epidem. Commty Hlth*, **70**, 1222–1228.
- 40 Leek, J. T. and Jager, L. R. (2017) Is most published research really false? *A. Rev. Statist. Appl.*, **4**, 109–122.
- 41 van der Linden, S., Leiserowitz, A., Rosenthal, S. and Maibach, E. (2017) Inoculating the public against mis-
 42 information about climate change. *Globl Chall.*, **1**, no. 2, article 1600008.
- 43 Makri, A. (2017) Give the public the tools to trust scientists. *Nat. News*, **541**, 261.
- 44 Matthews, R., Wasserstein, R. and Spiegelhalter, D. (2017) The ASA's *p*-value statement, one year on. *Significance*,
 45 **14**, 38–41.
- 46 McConway, K. (2016). Statistics and the media: a statistician's view. *Journalism*, **17**, 49–65.
- 47 Medical Xpress (2016) High levels of education linked to heightened brain tumor risk. (Available
 48 from [https://medicalxpress.com/news/2016-06-high-linked-heightened-brain-tumor.
 html](https://medicalxpress.com/news/2016-06-high-linked-heightened-brain-tumor.html).)
- Mervis, J. (2017) Data check: NSF sends Congress a garbled message on misconduct numbers. *Science*, Mar. 24th.
- Munafò, M. R., Nosek, B. A., Bishop, D. V. M., Button, K. S., Chambers, C. D., Percie du Sert, N., ... Ioannidis,
 J. P. A. (2017) A manifesto for reproducible science. *Nat. Hum. Behav.*, **1**, no.1, article 0021.
- National Centre for Social Research (2017) Public confidence in official statistics. National Centre for Social
 Research, London. (Available from [http://natcen.ac.uk/our-research/research/public-con
 fidence-in-official-statistics/](http://natcen.ac.uk/our-research/research/public-confidence-in-official-statistics/).)
- Newton-Cheh, C., Johnson, T., Gateva, V., Tobin, M. D., Bochud, M., Coin, L. ... , Munroe, P. B. (2009) Genome-
 wide association study identifies eight loci associated with blood pressure. *Nat. Genet.*, **41**, 666–676.
- O'Neill, O. (2002) *A Question of Trust: the BBC Reith Lectures 2002*. Cambridge: Cambridge University Press.
- O'Neill, O. (2013) What we don't understand about trust. (Available from [https://www.ted.com/talks/
 onora_o_neill_what_we_dont_understand_about_trust/transcript?language=en](https://www.ted.com/talks/onora_o_neill_what_we_dont_understand_about_trust/transcript?language=en).)
- Open Science Collaboration (2015) Estimating the reproducibility of psychological science. *Science*, **349**, article.
 4716.
- Patil, P., Peng, R. D. and Leek, J. T. (2016) What should researchers expect when they replicate studies?: a statistical
 view of replicability in psychological science. *Perspect. Psychol. Sci.*, **11**, 539–544.

- 1 Ranhill, E., Dreber, A., Johannesson, M., Leiberg, S., Sul, S. and Weber, R. A. (2015) Assessing the robustness
 2 of power posing: no effect on hormones and risk tolerance in a large sample of men and women. *Psychol. Sci.*,
 3 **26**, 653–656.
- 4 Royal Society (2012) Science as an open enterprise. *Report*. Royal Society, London. (Available from <https://royalsociety.org/topics-policy/projects/science-public-enterprise/report/>.)
- 5 Royal Statistical Society (2016) Data manifesto. Royal Statistical Society, London. (Available from http://www.rss.org.uk/RSS/Influencing_Change/Data_manifesto/RSS/Influencing_Change/Data_democracy_sub/Data_manifesto.aspx?hkey=5dd70207-82e7-4166-93fd-bcf9a2a1e496.)
- 6 Science Media Centre (2017). Where science meets the headlines. Science Media Centre, London.
- 7 Sense about Science (2017) Because evidence matters. Sense about Science, London.
- 8 Shafik, M. (2017) In experts we trust? Bank of England, London, (Available from <http://www.bankofengland.co.uk/publications/Documents/speeches/2017/speech964.pdf>.)
- 9 Shirakawa, T., Iso, H., Yamagishi, K., Yatsuya, H., Tanabe, N., Ikehara, S., ... Tamakoshi, A. (2016) Watching
 10 television and risk of mortality from pulmonary embolism among Japanese men and women: the JACC study
 11 (Japan Collaborative Cohort). *Circulation*, **134**, 355–357.
- 12 Simmons, J. P., Nelson, L. D. and Simonsohn, U. (2011) False-positive psychology: undisclosed flexibility in data
 13 collection and analysis allows presenting anything as significant. *Psychol. Sci.*, **22**, 1359–1366.
- 14 Slovic, P., Finucane, M. L., Peters, E. and MacGregor, D. G. (2004) Risk as analysis and risk as feelings: some
 15 thoughts about affect, reason, risk, and rationality. *Risk Anal.*, **24**, 311–322.
- 16 Spinney, L. (2017) How Facebook, fake news and friends are warping your memory. *Nat. News*, **543**, 168.
- 17 Stempa (2017) Guide to being a press officer. Stempa. (Available from https://stempa.org.uk/wp-content/themes/stempa/downloads/2017_stempa_guide_to_being_a_media_officer.pdf.)
- 18 Sumner, P., Vivian-Griffiths, S., Boivin, J., Williams, A., Bott, L., Adams, R., ... Chambers, C. D. (2016) Exagger-
 19 ations and caveats in press releases and health-related science news. *PLOS ONE*, **11**, no 12, article e0168217.
- 20 Sumner, P., Vivian-Griffiths, S., Boivin, J., Williams, A., Venetis, C.A., Davies, A., ... Chambers, C. D. (2014)
 21 The association between exaggeration in health related science news and academic press releases: retrospective
 22 observational study. *Br. Med. J.*, **349**, article 7015.
- 23 Swift, J. (1843) *The Works of Jonathan Swift ...: Containing Interesting and Valuable Papers, not hitherto Published*
 24 *... with Memoir of the Author*. London: Bohn.
- 25 Szucs, D. and Ioannidis, J. P. A. (2017) Empirical assessment of published effect sizes and power in the recent
 26 cognitive neuroscience and psychology literature. *PLOS Biol.*, **15**, no.3, article e2000797.
- 27 Tetlock, P. E. and Gardner, D. (2015) *Superforecasting: the Art and Science of Prediction*. Toronto: McClelland
 28 and Stewart.
- 29 Wasserstein, R. L. and Lazar, N. A. (2016) The ASA's statement on *p*-values: context, process, and purpose. *Am.*
 30 *Statistn*, **70**, 129–133.
- 31 Wellcome Trust (2017) Public views on medical research. Wellcome Trust, London. (Available from <https://wellcome.ac.uk/what-we-do/our-work/public-views-medical-research>.)
- 32 YouGov (2017). Leave voters are less likely to trust any experts—even weather forecasters. You Gov, London.
 33 (Available from <http://yougov.co.uk/news/2017/02/17/leave-voters-are-less-likely-trust-any-experts-eve/>.)
- 34 Youtube (2016) Gove: Britons “Have had enough of experts”. (Available from <http://www.youtube.com/watch?v=GGgiGtJk7MA>.)
- 35
36
37
38
39
40
41
42
43
44
45
46
47
48