Replication Crisis, Prediction Markets and Fundamental Principle of Probability

Harry Crane

Department of Statistics Rutgers

October 22, 2018

Key Ideas

- Goodhart's Law: When a measure becomes a target it ceases to be a good measure.
 - Theory: P-value measures the amount of evidence data provides for a scientific hypothesis.
 - Practice: P-value < 0.05 is the target that determines whether results are published or not.
- **Cournot's Principle**: an event of small or zero probability singled out in advance will not happen. (Shafer)
 - Gives probabilities empirical meaning.
- Fundamental Principle of Probability: If you assign a probability to an outcome, then you must accept a bet offered on the other side at the correct implied odds.
 - Gives probability statements both credibility and meaning.
 - Meaning: Probability p of A means an offer to bet against A at p/(1-p) odds.
 - *Credibility*: Someone claiming a probability *p* of *A* must believe the probability is *at least p* or else risk long-term loss.

H. Crane. (2018). The Fundamental Principle of Probability. https://www.researchers.one/article/2018-08-16

Essay

Why Most Published Research Findings Are False

John P. A. Ioannidis

Why?

- Low Power: The smaller the studies conducted in a scientific field, the less likely the research findings are to be true.
- Multiple Testing: The greater the number and the lesser the selection of tested relationships in a scientific field, the less likely the research findings are to be true.
- Researcher Degrees of Freedom: The greater the flexibility in designs, definitions, outcomes, and analytical modes in a scientific field, the less likely the research findings are to be true.
- Incentives: The greater the financial and other interests and prejudices in a scientific field, the less likely the research findings are to be true.

"Claimed Research Findings May Often Be Simply Accurate Measures of the Prevailing Bias" (Ioannidis)

What proportion of results should replicate?

	H ₀ true	H_0 false
Proportion	ϕ	$(1-\phi)$
Reject	$\alpha\phi$	$(1-eta)(1-\phi)$
Not Reject	$(1-lpha)\phi$	$\beta(1-\phi)$

Table: α is Type-I error rate; β is Type-II error rate; $(1 - \phi)/\phi$ is prior odds.

- Type-I error rate: $Pr(P < \alpha \mid H_0 \text{ true}) = \alpha$.
- Type-II error rate: $Pr(P > \alpha | H_0 \text{ false}) = \beta$.
- **Prior odds** $(1 \phi)/\phi$: ratio of false to true H_0 among all those tested.
- False positive rate (FPR): proportion of false positives among all $P < \alpha$

$$\mathsf{FPR}(\alpha,\beta,\phi) = \frac{\alpha\phi}{\alpha\phi + (1-\beta)(1-\phi)}.$$
 (1)

• Replication rate (RR): RR(α, β, ϕ) = 1 - FPR(α, β, ϕ).

Large-scale replication studies have successfully replicated

- 39% of findings (37 out of 97) in psychology.¹
- 61% (11 out of 18) in experimental economics.²
- 62% (13 out of 21) of social science articles published in Science and Nature.³

Does this contradict theoretical expectations?

α	β	$\phi/(1-\phi)$	FPR	RR
0.05	0.20	1:1	6%	94%
0.05	0.20	1:10	36%	64%
0.005	0.20	1:1	< 1%	> 99%
0.005	0.20	1:10	5%	95%

- Maybe 60% replication rate shouldn't be surprising reflects testing parameters (α, β, ϕ) .
- But published findings are/should be based on more than single hypothesis test Theoretical FPR should be upper bound.

¹Open Science Collaboration. Estimating the reproducibility of psychological science. *Science*, 2015.

²C. Camerer et al. Evaluating replicability of laboratory experiments in economics. *Science*, 2016.

³C. Camerer et al. Evaluating the replicability of social science experiments in Nature and Science between 2010 and 2015. *Nature Human Behaviour*, 2, 637–644, 2018.

What to do?

- "Redefine Statistical Significance"
 - Change default cutoff for statistical significance from 0.05 to 0.005.
 - Why? FPR under 5% level is too high even in absence of misconduct or reporting issues.
- Pre-registration and results-blind peer review
 - Decide whether to publish results based on methods, not results.
 - Why? Incentivize sound methods over results.
- Prediction Markets
 - Use betting markets to make better publication decisions.
 - Why? Incentivize more accurate assessment of replication chances by scientific community prior to publication.
- Fundamental Principle of Probability
 - Authors set their own replication criteria and probability and back up claims with their own money.
 - Why? Incentivize more accurate and honest assessment of replication by authors. Penalize inaccurate and dishonest claims (i.e., P-hacking).

Observation: Current paradigm promotes skewed incentives, reporting bias, and replication crisis.

Solution: Change the paradigm to eliminate reporting bias (publish everything), align incentives (reward accuracy, punish inaccuracy), and directly target replication.

Harry Crane (Rutgers)

Fundamental Principle of Probability

Title: Redefine Statistical Significance

Authors: Daniel J. Benjamin^{1*}, James O. Berger², Magnus Johannesson^{3*}, Brian A. Nosek^{4,5}, E.-J. Wagenmakers⁶, Richard Berk^{7, 10}, Kenneth A. Bollen⁸, Björn Brembs⁹, Lawrence Brown¹⁰, Colin Camerer¹¹, David Cesarini^{12, 13}, Christopher D. Chambers¹⁴, Merlise Clyde², Thomas D. Cook^{15,16}, Paul De Boeck¹⁷, Zoltan Dienes¹⁸, Anna Dreber³, Kenny Easwaran¹⁹, Charles Efferson²⁰, Ernst Fehr²¹, Fiona Fidler²², Andy P. Field¹⁸, Malcolm Forster²³, Edward I. George¹⁰, Richard Gonzalez²⁴, Steven Goodman²⁵, Edwin Green²⁶, Donald P. Green²⁷, Anthony Greenwald²⁸, Jarrod D. Hadfield²⁹, Larry V. Hedges³⁰, Leonhard Held³¹, Teck Hua Ho³², Herbert Hoijtink³³, James Holland Jones^{39,40}, Daniel J. Hruschka³⁴, Kosuke Imai³⁵, Guido Imbens³⁶, John P.A. Ioannidis³⁷, Minieong Jeon³⁸, Michael Kirchler⁴¹, David Laibson⁴², John List⁴³, Roderick Little⁴⁴, Arthur Lupia⁴⁵, Edouard Machery⁴⁶, Scott E, Maxwell⁴⁷, Michael McCarthy⁴⁸, Don Moore⁴⁹, Stephen L. Morgan⁵⁰, Marcus Munafó^{51, 52}, Shinichi Nakagawa⁵³, Brendan Nyhan⁵⁴, Timothy H. Parker⁵⁵, Luis Pericchi⁵⁶, Marco Perugini⁵⁷, Jeff Rouder⁵⁸, Judith Rousseau⁵⁹, Victoria Savalei⁶⁰, Felix D, Schönbrodt⁶¹, Thomas Sellke⁶², Betsy Sinclair⁶³, Dustin Tingley⁶⁴, Trisha Van Zandt⁶⁵, Simine Vazire⁶⁶, Duncan J. Watts⁶⁷, Christopher Winship⁶⁸, Robert L. Wolpert², Yu Xie⁶⁹, Cristopal Young⁷⁰, Jonathan Zinman⁷¹, Valen E, Johnson⁷²*

One Sentence Summary: We propose to change the default *P*-value threshold for statistical significance for claims of new discoveries from 0.05 to 0.005.

- Benjamin, et al (2017, *Nature Human Behavior*) propose to "redefine statistical significance" from *P* < 0.05 to *P* < 0.005.
- Claimed benefits of RSS:
 - Reproducibility would "immediately improve".
 - False positive rates will fall as low as 5% and improve by factors greater than 2.
 - Replication rate will approximately double.
- Other claims:
 - Lower cutoff supported by "critical mass of researchers".
 - 0.05 cutoff is "leading cause of non-reproducibility".

• Flaws in argument:

• Numerous — see Amrhein–Greenland; Lakens, et al; Trafimow, et al; Crane.

Main criticism:

- Argument based on hidden/misleading assumptions.
- Leads to exaggerated conclusions.

• Irony:

- RSS article in *Nature Human Behavior* ignores an obvious part of human behavior (i.e. P-hacking and incentives of publication).
- Makes non-replicable claims in argument claiming to improve replication crisis.

H. Crane. (2018). The impact of P-hacking on "Redefine Statistical Significance". Basic and Applied Social Psychology.

Harry Crane (Rutgers)

https://projects.fivethirtyeight.com/2018-midterm-election-forecast/senate/new-jersey/

ELECTION 2018						
FiveThirtyEight	House forecast	Senate	Governor	Midterms coverage	More politics v	ODO NEWS

9 in 10

Chance the Democrat wins (90.4%)

1 in 10

£у

Chance the Republican wins (9.6%)

Image: Image:

Candidate	Forecasted vote share	Chance of winning
Robert Menendez (D)	53.5	9 in 10 (90.4%)
Bob Hugin (R)	43.2	1 in 10 (9.6%)
Other candidates	3.3	<1 in 100 (<0.1%)
	10% 20 30 40 50 60 70	80 90 100

FiveThirtyEight Estimate:

Pr(Menendez re-elected) = 0.904

Pr(Menendez not re-elected) = 0.096

Prediction Markets (from PredictIt.com)

Will Bob Menendez be re-elected to the U.S. Senate in New Jersey in 2018?

Latest Price: 78¢ + 1¢



If this prediction comes true, Predictit will redeem all Yes shares at \$1. Shares in No will have zero value. If this prediction does not come true, Predictit will redeem all No shares at \$1. Shares in Yes will have zero value.



Bob Menendez shall be the winner of the 2018 general election for U.S. Senator from New Jersey. Predictit may determine how and when to settle the market based on all information available to Predictit at the relevant time.

Predictit reserves the right to wait for further official, party, judicial or other relevant announcements, reports or decisions to resolve any ambiguity or uncentarity before the market is settled. Markets may stop open or incur a delay in settlement well parts due of the contest in certain circumstances. If there is any change to an event, or any situation arises, that is not in Predictit's wie addressed adequately by the market rules. Predictit will exide the integer and noncorrotate course of action.

Predictit's decisions and determinations under this rule shall be at Predictit's sole discretion and shall be final.

Prediction Market Estimate:

Pr(Menendez re-elected) = 0.78

Pr(Menendez not re-elected) = 0.23

・ロト ・同ト ・ヨト ・ヨ

	https://projects.fivethirtyeight.com/2018-midterm-election-forecast/senate/new-jersey/				
Will Bob Menendez be re-elected to the U.S. Senate in New Jersey in 2018? Latest Price: 786 + 16	FiveThirtyEight	House forecast	Senate Governor	Midterms coverag	n More politics ~ @NEWS
Buy Yes Dia to match Offen starting at 786, or to make year own, lower Offen. Dito to make year own, lower Offen.	New Jerse	UKELY D			f 9
*** pretiction zone an use. Preficit will indeen all Yas shares at \$1.5 uma in the will have zon whee, if this predictor does not come true, Predict will indeen all to shares at \$1.5 uma in Yas will have zon whee.	9 in 10 Chance the Democrat w	ins (90.4%)		Chance	1 in 10 the Republican wins (9.0%)
Data Rules Prices	Candidate	Fores	asted vote share		Chance of winning
Bob Menendez shall be the winner of the 2018 general election for U.S. Senator from New Jersey. Predictit may determine how and when to settle the market based on all information available to Predict) at the relevant time.	Robert Menendez (D) Incembert			53.5	9 in 10 (90.4%)
dists reserves the right to wait for further official, party, judicial or other relevant announcements, reports or decisions to resolve any his situ or providenty before the market is written Markets may stay open or invert a delay to writtenent well not the date of the createst	Bob Hugin (R)		43.2		1 in 10 (0.6%)
v certain discurstances. If there is any charge to an event, or any situation arises, that is not in Predictit's view addressed adequately by the narket rules. Predictit will decide the fairest and most appropriate course of action.	Other candidates	3.3			<1 in 100 (40.00)
redictifs decisions and determinations under this rule shall be at Predictifs sole discretion and shall be final.			0% 20 30 40	50 63 70 83	90 100

Prediction Market:Pr(Menendez re-elected) = 0.78538 Forecast:Pr(Menendez re-elected) = 0.90

Which is more accurate?

- If 538 estimate is right, the market is offering a 12% edge to bet on Menendez.
- Either (i) markets very inefficient or (ii) 538 estimates are unreliable. (Maybe both.)

Who has more incentive to be accurate?

• Who is more likely to suffer from slight deviations from accuracy? Markets.

Prediction Markets and Replication

Evaluating the replicability of social science experiments in *Nature* and *Science* between 2010 and 2015

Colin F. Camerer¹⁰, Anna Dreber¹³, Felix Holzmeister¹³⁴⁸, Teck-Hua Ho¹⁰, Jirgen Huber¹⁴⁶, Magnus Johannesson¹³²⁴, Michael Kirchler³³⁵, Gideon Nave⁶¹⁰, Brian A. Nosek⁰³³⁶⁴, Thomas Pfelfer¹³²⁹, Adam Altmejd¹³, Nick Butrick³³, Taizan Chan⁵, Yiling Cher¹, Eskil Forsell³², Anup Gampa³⁴, Erma Heikensteri, Lily Hummer¹, Taisuke Imal¹⁰¹⁰, Siri Isaksson², Dylan Manfredi⁴, Julia Rose¹, Erica Jan Wagenmaker¹³⁴⁸ and Hang Wu³



Idea: Use Prediction Markets to get more accurate measure of scientific community's assessment of replication.

Submit paper to journal \rightarrow betting market opens up.

Conditional on a replication attempt, the claimed results will replicate:

YES: \$0.67 NO: \$0.33

- If replication attempt successful, YES worth \$1.00, NO worth \$0.00.
- If replication unsuccessful, YES worth \$0.00, NO worth \$1.00.
- If no replication attempt within (say) 1 year, then no action.
- After 2 months (or so), journal decides whether to accept/reject paper. Can use prediction market price in its decision.
- If accepted, the prediction market price after 2 months is published with article.

R. Hanson. Could Gambling Save Science? Encouraging an Honest Consensus. R. Hanson. Shall We Vote on Values, But Bet on Beliefs? Recall loannidis's reasons for low replicability:

- Low Power: Can and should be priced in by the prediction market.
- Multiple Testing: Can be priced in by bettors, but information asymmetry between authors and market.
- Researcher Degrees of Freedom: Limited because authors specify replication criteria prior to submission.
- Incentives: Author incentive to pump up the price of replication.
- Claimed Research Findings May Often Be Simply Accurate Measures of the Prevailing Bias": High prediction market prices may just be more accurate measure of prevailing bias.

Main issue: Authors still benefit from information asymmetry.

- Authors benefit from inaccurate high market prices \Rightarrow incentive to mislead.
- Authors punished by inaccurate low market prices ⇒ paradigm-shifting and unpopular ideas suppressed.
- P-hacking, QRPs, etc. still viable career strategy.
- Market probability is a more accurate reflection of the community opinion, but can't account for unreported insider (author) information and authors still have incentive to withhold information.

So what do we expect to change?

Fundamental Principle of Probability (FPP): If you assign a probability to an outcome, then you must accept a bet on the other side at the correct implied odds.

Idea: Apply FPP to get more accurate measure of author's assessment of claims and directly tie accuracy of author assessment to outcomes of replication.

- Publish paper along with:
 - (I) **Replication Criteria**: run [experimental protocol] for sample size $N \ge 100$, compute test statistic *T*. If $|T| > t_c$, then declare successful replication; otherwise, not.
 - (II) **Replication Probability**: the above procedure will replicate with probability *p*.
 - (III) Exposure Limit: Authors put some amount of money (e.g., \$10,000) in escrow for pre-determined period of scrutiny (say, 2 years).
- **(a)** During 2 year period, anyone else can put up A to gain $p/(1-p) \times A$ in event of failed replication.
 - If replication attempt successful, Authors gain \$A.
 - If replication unsuccessful, Authors lose \$A × p/(1 − p).
 - If no replication attempt, then no action.
- In the second second
- All papers (even failed replications) are published.

・ロト ・同ト ・ヨト ・ヨト

Recall loannidis's reasons for low replicability:

- Low Power
- Multiple Testing

Researcher Degrees of Freedom

All 3 can and should be priced in by authors. Failure to do so directly exposes authors to loss.

- Incentives: Author incentive to deflate the probability of replication. Think about how casinos set odds.
- Claimed Research Findings May Often Be Simply Accurate Measures of the Prevailing Bias": Claimed probabilities reflect only the authors bias/naivety/ misguidedness/expertise. If authors are wrong, opportunity for other researchers to gain. If community is wrong, authors gain.

- Incentivize accurate reporting: Authors benefit from accurate assessment of their studies and are penalized for inaccurate assessment.
- Organic funding mechanism: Research funds (gain/loss from betting) move in the direction of smarter scientists and move away from bad/dishonest scientists.
- Information alignment: Information asymmetry between authors and community (P-hacking, multiple testing, etc.) priced into replication probabilities.
- Reverse P-hacking: Stated probability should be a conservative.

Possible objection:

• Publishing everything will make science less reliable.

No, it will make it more reliable.

- Authors can't filter out their own bad findings.
- Editors/reviewers can't suppress/censor "bad" ideas. Community must demonstrate that a claim is false by betting the other side and holding up during replication.

H. Crane and R. Martin. (2018). In peer review we (don't) trust: How peer review's filtering poses a systemic risk to science. https://researchers.one/article/2018-09-17

Comparison: Prediction Markets vs. FPP

		Prediction Market	FPP
	Replication Probability	Market-driven	Set by authors
	Role of Probability	Publication Decision, Signaling	Price of real bets Skin in the game
	Consequence: Prob too high	Type I error: Publication Author benefits	Author loses money/funds Challenger gains
_	Consequence: Prob too low	Type II error: Rejection Author penalized	Conservative estimate No harm, no foul
	Author incentive	Overstate probability Deception	Understate probability Conservative
	Information asymmetry	Between Market & Author	Price in or risk ruin
	Challenges to paradigm	Upper hand for "normal science" New ideas suppressed	Equal footing with <i>status quo</i> ১ ৰ 🖉 ১ ৰ ই ১ ৰ ই – পএকে
	Harry Crane (Rutgers)	Fundamental Principle of Probability	Rutgers: October 22, 2018 19 / 20

I'm talking about a specific, extra type of integrity that is [beyond] not lying, but bending over backwards to show how you're maybe wrong, that you ought to have when acting as a scientist. (Feynman)

- Foundations of Probability (through FPP) can make science more reliable without assuming trust or integrity.
- Scientists who overstate importance/reliability of results go "out of business".
- Remove all barriers to publication: Publication no longer prestigious and therefore no longer the "target" of scientific research. (Goodhart's Law)

H. Crane. (2018). The Fundamental Principle of Probability. https://www.researchers.one/article/2018-08-16

H. Crane and R. Martin. (2018). In peer review we (don't) trust: How peer review's filtering poses a systemic risk to science. https://researchers.one/article/2018-09-17