# Which results can we trust?
# Using replications and prediction markets to assess the reliability of scientific results
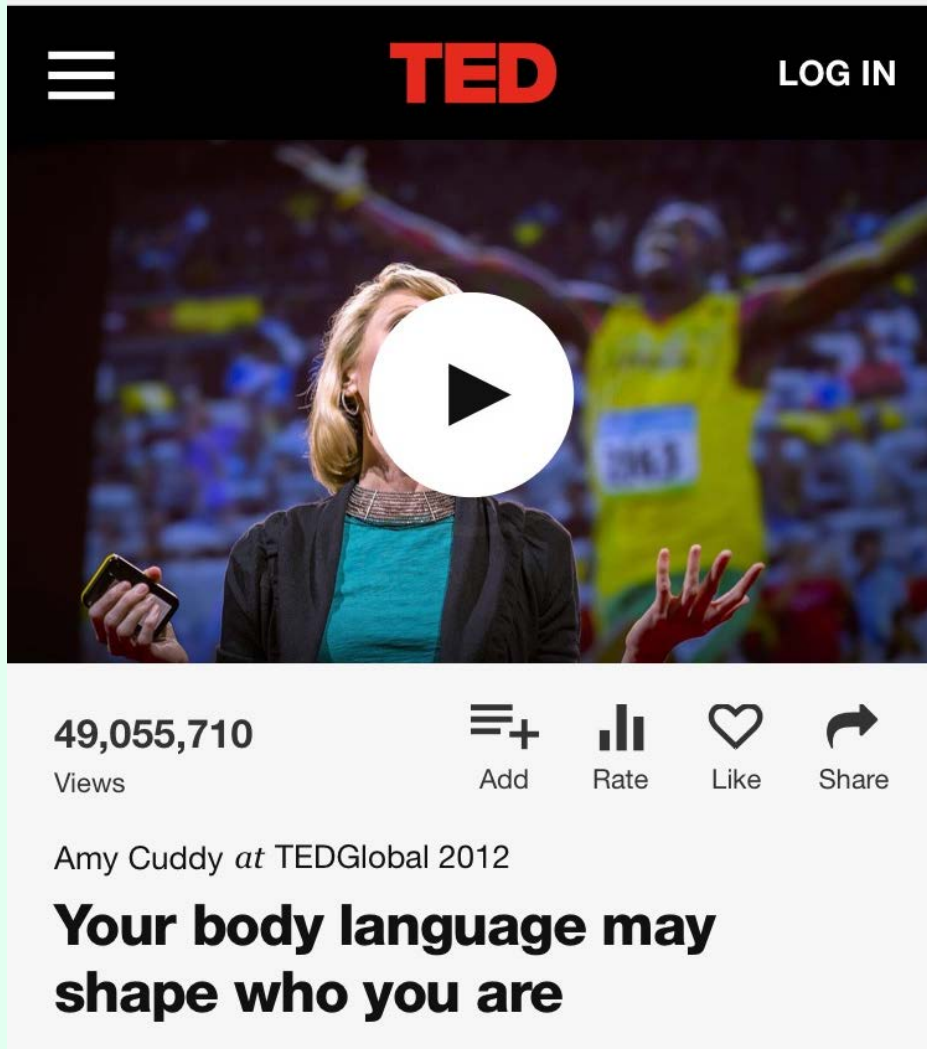
SUHF annual conference 2019 at Karolinska Institutet
Anna Dreber Almenberg
Department of Economics
Stockholm School of Economics

# Power posing



Carney et al. 2010, Ranehill et al. 2015

- How "researcher degrees of freedom" and low statistical power have lead to a replication crisis and how we should design studies and do <u>pre-analysis plans</u> to solve this problem


- Not only an experimental problem
- Not only a social science problem

# False results

# How many published claims are false?

- False positive results

- False negative results

**Most null results are never written up**
The fate of 221 social science experiments

| | |
|---|---|
| 100% | |
| 90 | |
| 80 | |
| 70 | |
| 60 | |
| 50 | |
| 40 | |
| 30 | |
| 20 | |
| 10 | |
| 0 | |

**Strong results** (42% of total)  **Mixed results** (36% of total)  **Null results** (22% of total)

■ Unwritten  ■ Unpublished but written  ■ Paper in non-top journal  ■ Paper in top journal

Source: A. Franco *et al.*, *Science* (28 August)

Ioannidis 2005 PLoS Medicine: Why Most Published Research Findings Are False

# "Researcher degrees of freedom"

## Histogram of p-values



Ioannidis 2005 Why Most Published Research Findings Are False;  Simmons, Nelson and Simonsohn 2011 False-Positive Psychology: Undisclosed Flexibility in Data Collection and Analysis Allows Presenting Anything as Significant; Gelman and Loken 2013 The Garden of Forking Paths

**Table 1.** Likelihood of Obtaining a False-Positive Result

| Researcher degrees of freedom | Significance level | | |
|---|---|---|---|
| | $p < .1$ | $p < .05$ | $p < .01$ |
| Situation A: two dependent variables ($r = .50$) | 17.8% | 9.5% | 2.2% |
| Situation B: addition of 10 more observations per cell | 14.5% | 7.7% | 1.6% |
| Situation C: controlling for gender or interaction of gender with treatment | 21.6% | 11.7% | 2.7% |
| Situation D: dropping (or not dropping) one of three conditions | 23.2% | 12.6% | 2.8% |
| Combine Situations A and B | 26.0% | 14.4% | 3.3% |
| Combine Situations A, B, and C | 50.9% | 30.9% | 8.4% |
| Combine Situations A, B, C, and D | 81.5% | 60.7% | 21.5% |

Note: The table reports the percentage of 15,000 simulated samples in which at least one of a set of analyses was significant. Observations were drawn independently from a normal distribution. Baseline is a two-condition design with 20 observations per cell. Results for Situation A were obtained by conducting three $t$ tests, one on each of two dependent variables and a third on the average of these two variables. Results for Situation B were obtained by conducting one $t$ test after collecting 20 observations per cell and another after collecting an additional 10 observations per cell. Results for Situation C were obtained by conducting a $t$ test, an analysis of covariance with a gender main effect, and an analysis of covariance with a gender interaction (each observation was assigned a 50% probability of being female). We report a significant effect if the effect of condition was significant in any of these analyses or if the Gender × Condition interaction was significant. Results for Situation D were obtained by conducting $t$ tests for each of the three possible pairings of conditions and an ordinary least squares regression for the linear trend of all three conditions (coding: low $= -1$, medium $= 0$, high $= 1$).

Simmons, JP, LD Nelson, U Simonsohn, 2011, False-Positive Psychology: Undisclosed Flexibility in Data Collection and Analysis Allows Presenting Anything as Significant. Psychological Science 22(11): 1359-1366.

# Forking

- Multiple testing problem where the universe of tests is not clear
- The data decide the analysis
- Beware subgroup analyses etc
- P-values are meaningless

# Which results can we trust?

- Depends on
  - P-values and power
  - Publication bias
  - Researcher degrees of freedom
  - Priors
    - Probability of a hypothesis to be true ("prior")
    - Typically subjective and unaccessible

# How big is the problem?

# (In some of the quantitative empirical social sciences)

Open Science Collaboration (2015). "Estimating the Reproducibility of Psychological Science." *Science*.
Camerer et al. (2016) "Evaluating replicability of laboratory experiments in economics." *Science*.
Camerer et al. (2018) "Evaluating the replicability of social science experiments in *Nature* and *Science* between 2010 and 2015." *Nature Human Behaviour.*

**Fig. 1. Density plots of original and replication P values and effect sizes. (A)** P values. **(B)** Effect sizes (correlation coefficients). Lowest quantiles for P values are not visible because they are clustered near zero.

35/97 positive results replicate
Relative effect size about 50%

Open Science Collaboration (2015). "Estimating the Reproducibility of Psychological Science."
*Science*, 349(6251).

**A**

Abeler et al., AER 2011 (32)
Ambrus and Greiner, AER 2012 (33)
Bartling et al., AER 2012 (34)
Charness and Dufwenberg, AER 2011 (35)
Chen and Chen, AER 2011 (36)
de Clippel et al., AER 2014 (37)
Duffy and Puzzello, AER 2014 (38)
Dulleck et al., AER 2011 (39)
Ericson and Fuster, QJE 2011 (40)
Fehr et al., AER 2013 (41)
Friedman and Oprea, AER 2012 (42)
Fudenberg et al., AER 2012 (43)
Huck et al., AER 2011 (44)
Ifcher and Zarghamee, AER 2011 (45)
Kessler and Roth, AER 2012 (46)
Kirchler et al, AER 2012 (47)
Kogan et al., AER 2011 (48)
Kuziemko et al., QJE 2014 (49)

- Estimate
⊢—⊣ 95% CI

-1   0   1   2

0= no effect

1=same effect
as in original study

11/18 results replicate
Relative effect size about 60%

Camerer et al. 2016 *Science*

**a. Stage 1 results** **b. Stage 2 results**

Relative standardized effect size

Relative standardized effect size

Studies:
Ackerman et al. (2010)[36], Science
Aviezer et al. (2012)[37], Science
Balafoutas and Sutter (2012)[38], Science
Derex et al. (2013)[39], Nature
Duncan et al. (2012)[40], Science
Gervais and Norenzayan (2012)[41], Science
Gneezy et al. (2014)[42], Science
Hauser et al. (2014)[43], Nature
Janssen et al. (2010)[44], Science
Karpicke and Blunt (2011)[45], Science
Kidd and Castano (2013)[46], Science
Kovacs et al. (2010)[47], Science
Lee and Schwartz (2010)[48], Science
Morewedge et al. (2010)[49], Science
Nishi et al. (2015)[50], Nature
Pyc and Rawson (2010)[51], Science
Ramirez and Beilock (2011)[52], Science
Rand et al. (2012)[53], Nature
Shah et al. (2012)[54], Science
Sparrow et al. (2011)[55], Science
Wilson et al. (2014)[56], Science

Legend:
⊢—⊣ 95% confidence interval
◆ Point estimate larger than zero ($p < 0.05$)
◆ Point estimate not different from zero ($p > 0.05$)

13/21 results replicate in Stage 2

Mean relative effect size: 50%. For 13 studies that replicated: 74%, for the rest, 0%

Camerer et al. 2018 *Nature Human Behaviour*

# "Could gambling save science?"



Hanson 1995 *Social Epistemology*

# Our prediction markets on replications

- 10 days – 2 weeks
- USD 50-100
- 50-100 participants
- Central hypothesis
- Binary outcomes
- Price: predicted probability of the outcome occuring
- Participants get replication reports
- Also survey questions

Study

Prediction market and survey beliefs

The power of prediction markets
Nature.com

*Atlantic*
Online Bettors Can Sniff Out Weak Psychology Studies
2 months ago

Psychologists' betting market hints at most reliable research findings
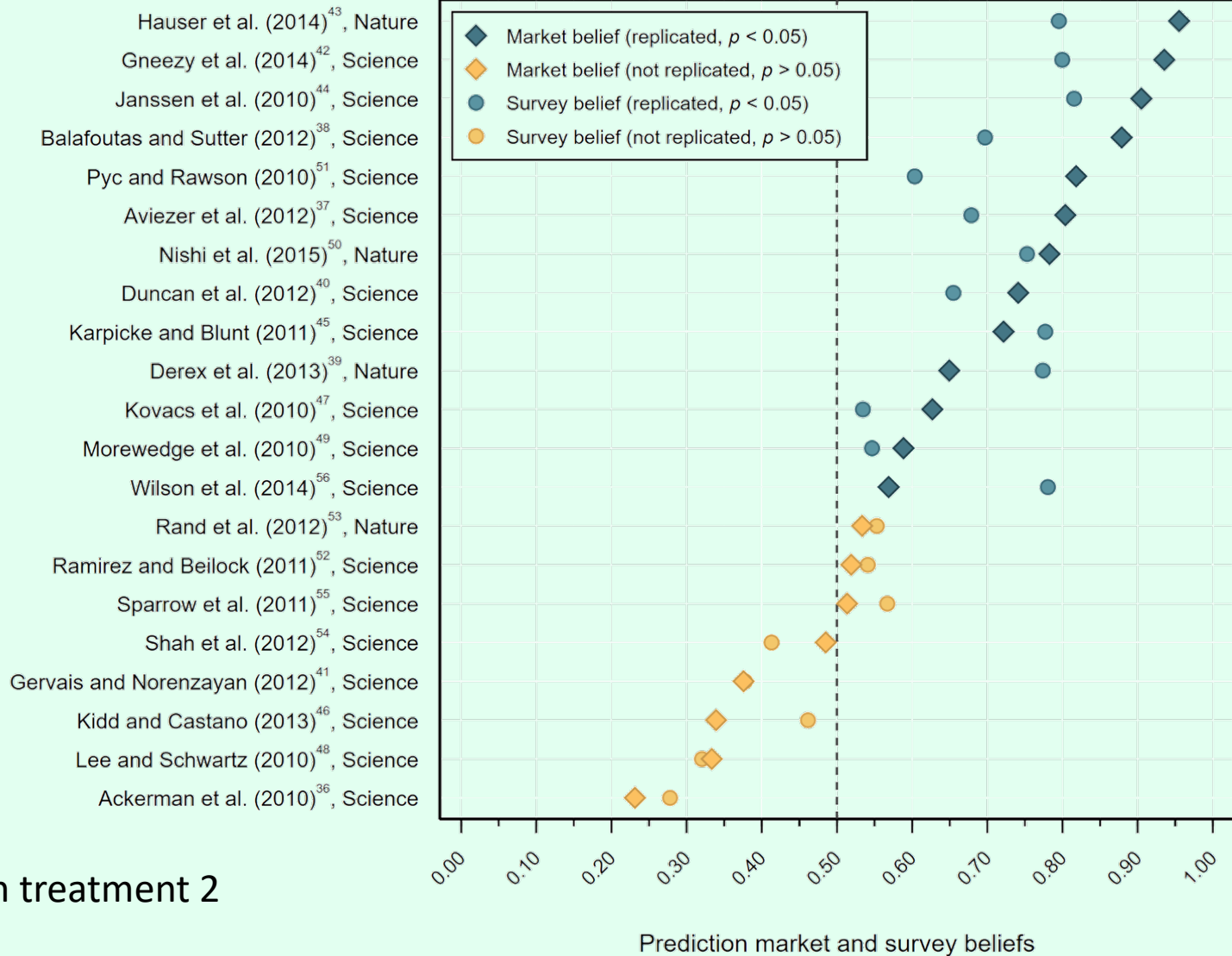Nature.com

Can a Futures Market Save Science?
The Atlantic

Replication:
- Did not replicate
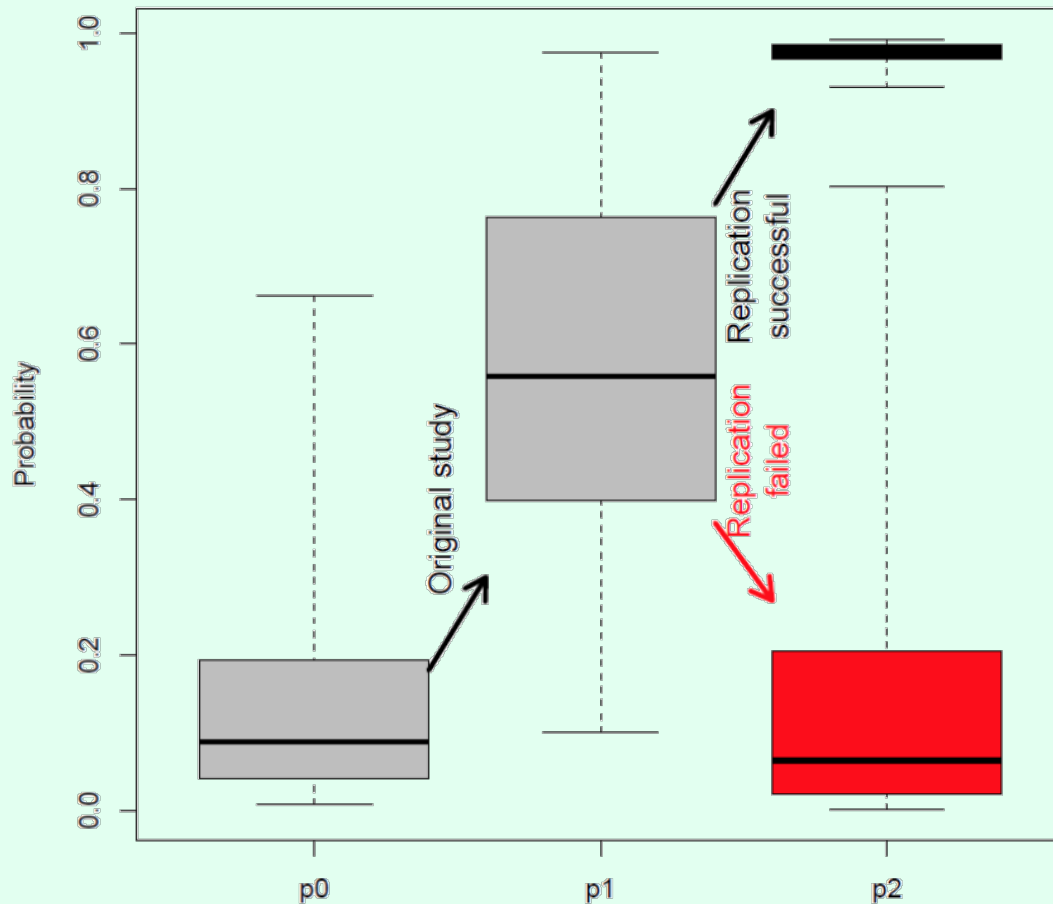- Replicated

Beliefs:
- Market beliefs
- Survey beliefs

Pooling 4 prediction market studies: 73% (76/104) correct prediction rate
Pooling 4 surveys: 66% (68/103) correct prediction rate

Work in progress

# Prediction markets results *Nature* and *Science*



From treatment 2

# Probability of hypothesis being true at 3 stages of testing for RPP



- Initial priors are low (median 8.8%)
- Positive result in initial publication moves prior to intermediate level (median 56%)
- If successful replication, probability moves up (median 98%)
- If failed replication, probabiliby close to initial prior (median 6.3%)

Dreber et al. 2015 PNAS

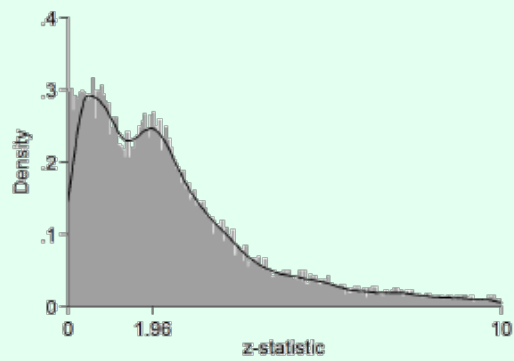Whiskers: range
Boxes: 1st to 3rd quartiles
Thick lines: medians

# What have we learned?

- Common false interpretation of p<0.05: 95% probability of hypothesis being true
- For this to be the case, a p<0.05 finding needs to supported in a high-powered replication
- Meta-analyses will also have inflated effect sizes – we need replications
- Are the incentives for replications appropriate?
- There is something systematic about results that fail to replicate – and experts "know" this
  - So why are so many false results published?
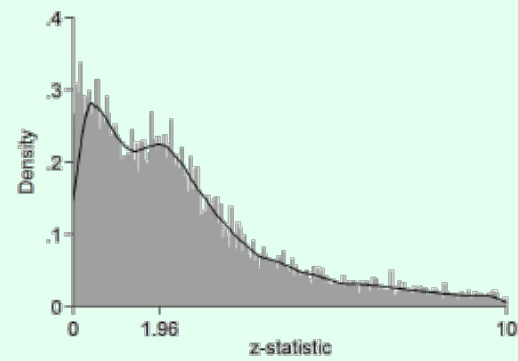
See, e.g., 2019 book chapter by Camerer, Dreber and Johannesson for more
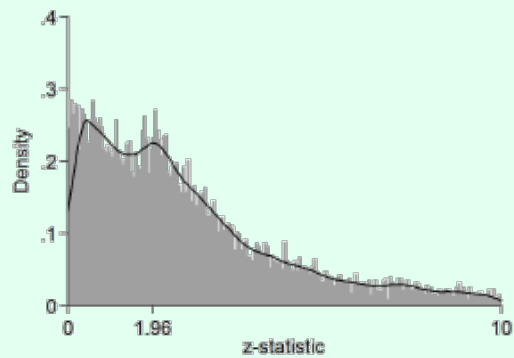
# Other thoughts

- <u>Pre-analysis plans</u>
- Problems probably worse for non-experimental work
- Higher power and team science
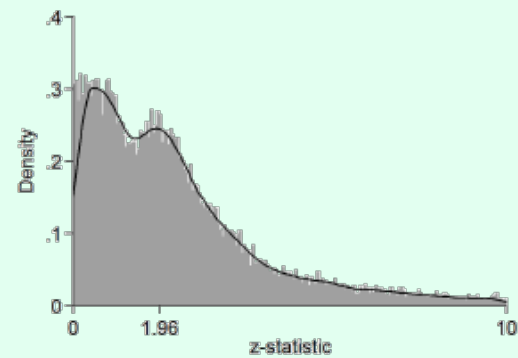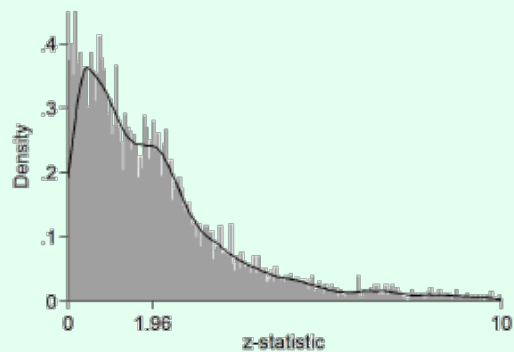  - Munafo et al. 2017 Nature Human Behaviour
- $p<0.005$

(a) Eye-catchers.

(b) No eye-catchers.

(c) Model.
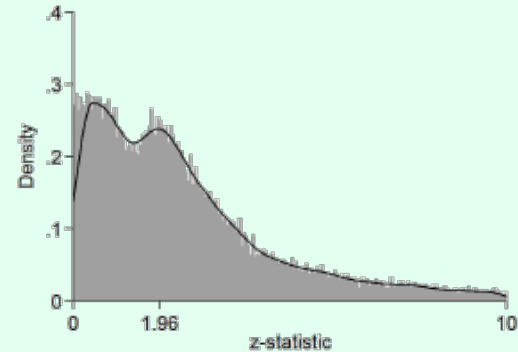
(d) No model.

(e) Lab. experiments or RCT data.

(f) Other data.

Brodeur et al 2016

Sources: AER, JPE, and QJE (2005-2011). Distributions are unweighted and plotted using de-rounded statistics. Lines correspond to kernel density estimates.
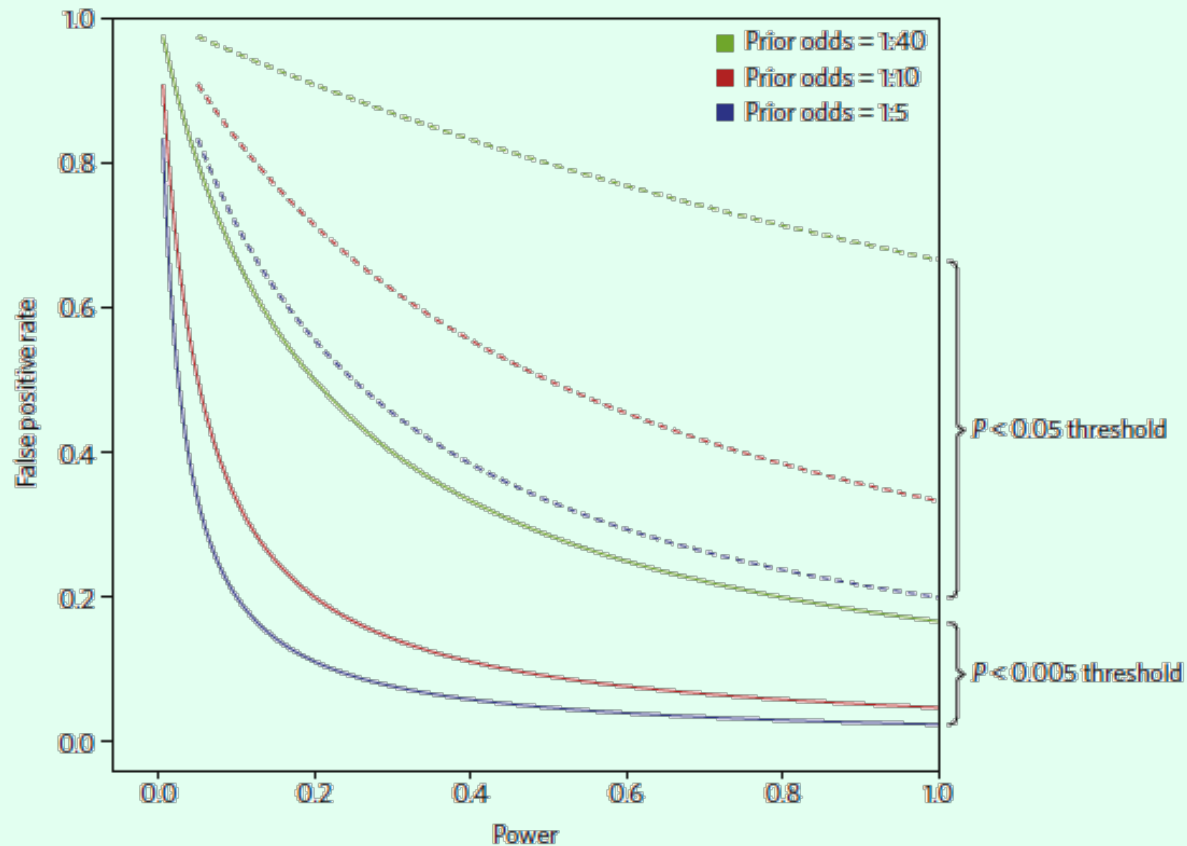
# p<0.005



**Fig. 2 | Relationship between the P value threshold, power, and the false positive rate.**
Calculated according to equation (2), with prior odds defined as $\frac{1-\phi}{\phi} = \frac{Pr(H_1)}{Pr(H_0)}$. For more details, see the Supplementary Information.

Benjamin et al. 2018 "Redefine Statistical Significance" Nature Human Behavior

# Thanks!

[anna.dreber@hhs.se](mailto:anna.dreber@hhs.se)

[www.replicationmarkets.com](http://www.replicationmarkets.com)