# Sailing Through Data: Discoveries and Mirages

Emmanuel Candès, Stanford University



Baidu Research, Sunnyvale, May 2019



# The replicability crisis

Begley and Ellis, Nature (2012)

Amgen could only replicate 6 of 53 studies they considered landmarks in basic cancer science

HealthCare could only replicate about 25% of 67 seminal studies

Systematic attempts to replicate widely cited priming experiments have failed



# Media coverage...



ARE HERE: LAT Home -+ Collections -+ Busine

Freeing the prisoners of NASA

#### Science has lost its way, at a big cost to humanity

Researchers are rewarded for splashy findings, not for double-checking accuracy. So many scientists looking for curves to discuss have been building on ideas that aren't even true. October 27, 2011 Michael Hiltrik

In today's world, brimful as it is with opinion and falseboods masquerading as facts, you'd think the one place you can depend on for verifiable facts is science.

You'd be wrong. Many billions of dollars' worth of wrong.

A few years ago, scientists at the Thousand Oaks biotech firm Amgen set out to double-check the results of 53 landmark papers in their fields of cancer research and blood biology.

The idea was to make sure that research on which Amgen



t few years ago, adjentiats at Amgen set out to double-che he results... (Anne Casack, Los Angeles...)



#### New Truths That Only One Can See

George Johnson RAW DATA JAN. 20, 2014



Since 1955, The Journal of Irreproducible Results has offered spoofs, parodies, whimsies, burlesques, lampoons and satires" about life in the laboratory. Among its greatest hits: "Acoustic Oscillations in Jell-O. With and Without Fruit, Subjected to Varying Levels of Stress" and "Utilizing Infinite Loops to Compute an Approximate Value of Infinity." The good-natured jibes are a backhanded celebration of science. What really goes on in the lab is, by implication, of a loftier, more serious nature.

00000

NEW YORKER

 The NEW YORKER

 The NEW YORKER

 Subscripting anywhere, everywhere.
 Subscripting anywhere of a free tore bag.

NEW YORMER

ANNALS OF SCIENCE DECEMBER IS, 2010 ISSUE

#### THE TRUTH WEARS OFF

Is there something wrong with the scientific method?

By Jonah Lehrer

f y s

On September 18, 2007, a few dozen drug-compary executive gathered in a drug-compary executive gathered in a bode conference room in Broused to hers some rarräng nears. In that ou do with class of drugs known as atypical or recordgeneration antipoyotics, which care on the market in the early sinetics. The drugs, outdouder bord manuses who as Ability Serospiel, and Zypreza, had bern tested on schizophersen in surveal large clinical trials, all de which had demonstrated a drumatic decrease in the solitics?



Many results that are rightenedy proved and accepted start shrinking in later studies. Internetion by LAIRENT GLUPTO

psychiatric symptoms. As a result, second-



The New York Times

#### Essay: The Experiments Are Fascinating. But Nobody Can Repeat Them.

Science is mired in a "replication" crisis. Fixing it will not be easy.



New York Times, Science, November 19, 2018

# The replicability problem

Early report (Kaplan, '08) 50% of Phase III FDA studies ended in failure 22 CASE STUDIES WHERE PHASE 2 AND PHASE 3 TRIALS HAD DIVERGENT RESULTS

U.S. FOOD & DRUG

anuary 2017

# Personal and societal concern

Snippets from media

"Significance chasing"

"Publication bias"

"Selective reporting"

Essay

# Why Most Published Research Findings Are False

John P.A. Ioannidis

# Personal and societal concern

Snippets from media

"Significance chasing"

"Publication bias"

"Selective reporting"

Essay Why Most Published Research Findings Are False John P.A. Joannidis

Great danger in seeing erosion of public confidence in science

Scientific community is responding

# Response: reproducibility intiatives



#### Major projects





Investigating the replicability of the 50 most impactful cancer biology studies from 2010-2012



Helping VCs, funding agencies, and others validate findings to promote high-quality research





Independently validating thousands of commercial antibodies to improve reliability



# **Reproducibility Initiative**

# http://validation. scienceexchange.com/

### Response: editorial policies

#### ANNOUNCEMENT Reducing our irreproducibility

Over the past year, Nature has published a string of articles that highlight failures in the reliability and reproducibility of published research (collected and freely available at go.nature.com/ huhbyr). The problems arise in laboratories, but journals such as this one compound them when they fail to exert sufficient scrutiny over the results that they publish, and when they do not publish enough information for other researchers to assess results properly.

From next month, Nature and the Nature research journals will introduce editorial measures to address the problem by improving the consistency and quality of reporting in life-sciences articles. To ease the interpretation and improve the reliability of published results we will more systematically ensure that key methodological details are reported, and we will give more space to methods sections. We will examine statistics more closely and encourage authors to be transparent, for example by including their raw data.

Central to this initiative is a checklist intended to prompt authors to disclose technical and statistical information in their submissions, and to encourage referees to consider aspects important for research reproducibility (go.nature.com/oloeip). It was developed after discussions with researchers on the problems that lead to irreproducibility, including workshops organized last year by US National Institutes of Health (NIH) institutes. It also draws on published concerns about reporting standards (or the lack of them) and the collective experience of editors at Nature journals

The checklist is not exhaustive. It focuses on a few experimental and analytical design elements that are crucial for the interpretation of research results but are often reported incompletely. For example, authors will need to describe methodological parameters that can introduce bias or influence robustness, and provide precise characterization of key reagents that may be subject to biological variability, such as cell lines and antibodies. The checklist also consolidates existing policies about data deposition and presentation. We will also demand more precise descriptions of statistics, and

we will commission statisticians as consultants on certain papers, at the editor's discretion and at the referees' suggestion

We recognize that there is no single way to conduct an experimental study. Exploratory investigations cannot be done with the same level of statistical rigour as hypothesis-testing studies. Few academic laboratories have the means to perform the level of validation required, for example, to translate a finding from the laboratory to the clinic. However, that should not stand in the way of a full report of how a study was designed, conducted and analysed that will allow reviewers and readers to adequately interpret and build on the results.

To allow authors to describe their experimental design and methods in as much detail as necessary, the participating journals, including Nature, will abolish space restrictions on the methods section.

To further increase transparency, we will encourage authors to provide tables of the data behind graphs and figures. This builds on our established data-deposition policy for specific experiments and large data sets. The source data will be made available directly from the figure legend, for easy access. We continue to encourage authors to share detailed methods and reagent descriptions by depositing protocols in Protocol Exchange (www.nature.com/ protocolexchange), an open resource linked from the primary paper.

Renewed attention to reporting and transparency is a small step. Much bigger underlying issues contribute to the problem, and are beyond the reach of journals alone. Too few biologists receive adequate training in statistics and other quantitative aspects of their subject. Mentoring of young scientists on matters of rigour and transparency is inconsistent at best. In academia, the ever increasing pressures to publish and chase funds provide little incentive to pursue studies and publish results that contradict or confirm previous papers. Those who document the validity or irreproducibility of a published piece of work seldom get a welcome from journals and funders, even as money and effort are wasted on false assumptions.

Tackling these issues is a long-term endeavour that will require the commitment of funders, institutions, researchers and publishers. It is encouraging that NIH institutes have led community discussions on this topic and are considering their own recommendations. We urge others to take note of these and of our initiatives. and do whatever they can to improve research reproducibility.



EDITORIAL CALENDAR PDF ARCHIVES ADVERTISE STATISTICIANS IN HISTORY ABOUT

Home » Additional Features, Featured, News and Announcements

#### Reproducible Research in JASA

1.88 Y 2018 1.234 VIEWS 3 COMMENTS

Montse Fuentes, Coordinating Editor of JASA and Editor of JASA ACS



Societal impact through scientific advances is predicated on discovery and new knowledge that is reliable and robust and provides a solid foundation on which further advances can be built. Unfortunately, there is evidence many published scientific results will not stand the test of time, in part due to the lack of good scientific practices for reproducibility.

Our statistical profession has a responsibility to establish publication standards that improve the transparency and robustness of what we publish and to

promote awareness within the scientific community of the need for rigor in our statistical research to ensure reproducibility of our scientific results. JASA is committed to helping lead the effort by presenting solutions that can help improve research quality and reproducibility.

398 | NATURE | VOL 496 | 25 APRIL 2013

### Response: best practices



President's Council of Advisors on Science and Technology (PCAST) Public Meeting Agenda January 31, 2014

> National Academy of Sciences (NAS) 2101 Constitution Avenue, NW Washington, DC

> > Lecture Room

- 9:00 am Welcome from PCAST Co-Chairs John Holdren, Assistant to the President for Science and Technology; Director, Office of Science and Technology Policy (OSTP); Co-Chair, PCAST Eric Lander, Co-Chair, PCAST
- 9:05 an Improving Scientific Reproducibility in an Age of International Competition and Big Dub 1: Brease-there Genne Regtor, Chief Scientific Officer and Senior Vice-President R&D, TetraLogic Pharmaceuticals Donald Berry, Professor, Department of Biostatistics, University of Texas MD Anderson Cancer Center

Daniel MacArthur, Assistant Professor, Harvard Medical School and Massachusetts General Hospital and Associate Member, Broad Institute of Harvard and MIT

 Improving Scientific Reproducibility in an Age of International Competition and Big Data II: Editors Marcia McNut, Editor-In-Chief, Science Philip Campbell, Editor-In-Chief, Anner and Nature Publishing Group Véronique Kierner, Executive Editor and Head of Researchers Services, Nature Public and Group

STATISTICAL CHALLENGES IN ASSESSING AND FOSTERING THE REPRODUCIBILITY OF SCIENTIFIC RESULTS

A Workshop of the Committee on Applied and Theoretical Statistics NATIONAL RESEARCH COUNCIL OF THE NATIONAL ACADEMIES

#### RESEARCH REPRODUCIBILITY, REPLICABILITY, RELIABILITY

A Speech by Ralph J. Cicerone, President National Academy of Sciences Presented at the Academy's 152<sup>nd</sup> Annual Meeting April 27, 2015

# The replicability issue

Many different components

- 1. Publishing culture
- 2. Granting agencies culture
- 3. Computational reproducibility
- 4. Statistics: how to choose a finding? Statistical methodology enhancing replicability



Can only do 3 and 4 1 and 2 above pay grade

# Why is this happening? A new scientific paradigm



#### ${\sf Collect \ data \ first} \quad \Longrightarrow \quad {\sf Ask \ questions \ later}$

Large data sets available prior to formulation of scientific hypotheses/theories

Very different from hypothesis-driven research

# Example from genomics

Historically, molecular biology was hypothesis-driven research

"Sometime in the 90's"

Eruption of high-throughput technologies

Enabled thousands of genes to be tested simultaneously for differential expression Small # of samples High # of variables

Researchers begin to look everywhere



gene expression microarray

Complete revolution: from hypothesis- to data-driven research



1000 hypotheses to test



1000 hypotheses, 100 potential discoveries



1000 hypotheses, 100 potential discoveries





$$\mathsf{FDR} = \mathbb{E}\left[\frac{\#\mathsf{false positives}}{\#\mathsf{selections}}\right]$$



Reported

$$\mathsf{FDR} = \mathbb{E}\left[\frac{\#\mathsf{false positives}}{\#\mathsf{selections}}\right]$$

#### Knockoffs: Tools for Replicable Selections

Joint with R. Barber and Y. Fan, L. Janson and J. Lv



### Some data-driven scientific problems

One response Y: phenotype; e.g. Crohn's disease status, cholesterol level Hundreds of thousands of variables X: genotype information

Ex. 1: which genetic variations affect traits, e.g. the risk of a disease?



### Some data-driven scientific problems

One response Y: phenotype; e.g. Crohn's disease status, cholesterol level Hundreds of thousands of variables X: genotype information

Ex. 1: which genetic variations affect traits, e.g. the risk of a disease?



Ex. 2: which gene expression profiles help determine severity of a tumor? Ex. 3: which factors/variables help determine whether a loan will be repaid?

### Some data-driven scientific problems

One response Y: phenotype; e.g. Crohn's disease status, cholesterol level Hundreds of thousands of variables X: genotype information

Ex. 1: which genetic variations affect traits, e.g. the risk of a disease?



Ex. 2: which gene expression profiles help determine severity of a tumor? Ex. 3: which factors/variables help determine whether a loan will be repaid?

How can we select variables without too many false positives?  $\rightsquigarrow$  do not run into problem of irreproducibility

# Formalizing the selection problem



Thousands/millions of variables XWhich ones are important? Distribution of  $Y \mid X$  depends on Xthrough which variables?

# Formalizing the selection problem



Thousands/millions of variables XWhich ones are important? Distribution of  $Y \mid X$  depends on Xthrough which variables?

Variable is a discovery if p(response | variable, others) $\neq p(\text{response} | \text{others})$ 

(Formally) j null iff  $Y \perp X_j | X_{-j}$ 

# Conditional testing

j null iff  $Y \perp X_j \mid X_{-j}$ 

Local Markov property  $\implies$  non nulls form Markov blanket of Y



# Conditional testing

j null iff  $Y \perp X_j \mid X_{-j}$ 

Local Markov property  $\implies$  non nulls form Markov blanket of Y



# Selection in the computer age

Many sophisticated tools to measure strength of dependence



Random forests



```
SVM
```





Bayes posteriors (MCMC)



#### Black-box algorithm

Deep nets











### Only one dataset: what should we report?



Modern science faces the problem of selection of promising findings from the noisy estimates of many.

Y. Benjamini and Y. Hechtlinger

# Knockoffs (Barber and Candès, 2015)

For each variable (e.g. SNP)  $X_j$ , make a knockoff version (e.g. fake SNP)  $\tilde{X}_j$ 

Run scoring procedure on features and knockoffs 'serving as controls'



Black box selects 49 original features & 24 knockoff features  $\implies$  probably  $\approx$  24 false positives among 49 original features
# How? By permutation?

٠



Can I use someone else's genetic info as control?

### Permuted dummies do not work!



# Knockoff dummies work!



#### Permuted dummies: other feature importance $Z_j$



# Knockoff dummies: other feature importance $Z_j$



## What's wrong?



# Model-X knockoffs

C., Fan, Janson and Lv ('16)

#### i.i.d. samples from ${\cal P}_{XY}$

- $P_X$  known
- $P_{Y|X}$  completely unknown

# Model-X knockoffs

C., Fan, Janson and Lv ('16)

#### i.i.d. samples from $P_{XY}$

- $P_X$  known
- $P_{Y|X}$  completely unknown

• Originals 
$$X = (X_1, \dots, X_p)$$
 • Knockoffs  $\tilde{X} = (\tilde{X}_1, \dots, \tilde{X}_p)$ 

(1) Pairwise exchangeability: for any null j

$$\tilde{X}_j, X_j, X_{-j}, \tilde{X}_{-j} \stackrel{d}{=} X_j, \tilde{X}_j, X_{-j}, \tilde{X}_{-j}$$



# Model-X knockoffs

C., Fan, Janson and Lv ('16)

- i.i.d. samples from  $P_{XY}$ 
  - $P_X$  known
  - $P_{Y|X}$  completely unknown

• Originals 
$$X = (X_1, \dots, X_p)$$
 • Knockoffs  $\tilde{X} = (\tilde{X}_1, \dots, \tilde{X}_p)$ 

(1) Pairwise exchangeability: for any null j

$$\tilde{X}_j, X_j, X_{-j}, \tilde{X}_{-j} \stackrel{d}{=} X_j, \tilde{X}_j, X_{-j}, \tilde{X}_{-j}$$

(2) Ignore Y when constructing knockoffs:  $\tilde{X} \perp Y \mid X$ 

$$(\underbrace{X, \tilde{X}}_{Swap(4)}) \stackrel{d}{=} (\underbrace{X, \tilde{X}}_{Swap(4)})$$

## Knockoffs as negative controls



All null scores are exchangeable

$$(Z_j, \tilde{Z}_j) \stackrel{d}{=} (\tilde{Z}_j, Z_j)$$

## Knockoffs with binary response

Feature importance  $Z_j$  and  $\tilde{Z}_j$  from random forests



## Knockoffs with binary response

Feature importance  $Z_j$  and  $\tilde{Z}_j$  from random forests



#### Knockoffs with binary response

Feature importance  $Z_j$  and  $\tilde{Z}_j$  from random forests



## Knockoffs-adjusted scores

#### Adjusted scores $W_j$ with flip-sign property

Combine  $Z_j$  and  $\tilde{Z}_j$  into single (knockoff) score  $W_j$ 

$$W_j = w_j(Z_j, \tilde{Z}_j)$$
  $w_j(\tilde{Z}_j, Z_j) = -w_j(Z_j, \tilde{Z}_j)$   
e.g.  $W_j = Z_j - \tilde{Z}_j$ 

#### Knockoffs-adjusted scores

#### Adjusted scores $W_j$ with flip-sign property

Combine  $Z_j$  and  $\tilde{Z}_j$  into single (knockoff) score  $W_j$ 

 $W_j = w_j(Z_j, \tilde{Z}_j)$   $w_j(\tilde{Z}_j, Z_j) = -w_j(Z_j, \tilde{Z}_j)$ 

e.g. 
$$W_j = Z_j - \tilde{Z}_j$$

Conditional on |W|, signs of null  $W_j$ 's are i.i.d. coin flips



### Knockoffs-adjusted scores

#### Adjusted scores $W_j$ with flip-sign property

Combine  $Z_j$  and  $\tilde{Z}_j$  into single (knockoff) score  $W_j$ 

 $W_j = w_j(Z_j, \tilde{Z}_j)$   $w_j(\tilde{Z}_j, Z_j) = -w_j(Z_j, \tilde{Z}_j)$ 

e.g. 
$$W_j = Z_j - \tilde{Z}_j$$

Conditional on |W|, signs of null  $W_j$ 's are i.i.d. coin flips







$$\mathsf{FDP}(t) = \frac{\#\{j \text{ null } : W_j \ge t\}}{\#\{j : W_j \ge t\} \lor 1}$$



$$\mathsf{FDP}(t) = \frac{\#\{j \text{ null } : W_j \ge t\}}{\#\{j : W_j \ge t\} \lor 1} \approx \frac{\#\{j \text{ null } : W_j \le -t\}}{\#\{j : W_j \ge t\} \lor 1}$$



$$\mathsf{FDP}(t) = \frac{\#\{j \text{ null } : W_j \ge t\}}{\#\{j : W_j \ge t\} \lor 1} \approx \frac{\#\{j \text{ null } : W_j \le -t\}}{\#\{j : W_j \ge t\} \lor 1}$$
$$\leq \frac{\#\{j : W_j \le -t\}}{\#\{j : W_j \ge t\} \lor 1} := \widehat{\mathsf{FDP}}(t)$$































Step-up rule: stop last time ratio between '-' and '+' below target FDR level



Select '+'s

Step-up rule: stop last time ratio between '-' and '+' below target FDR level

# FDR control

$$\hat{\mathcal{S}} = \{W_j \ge \tau\}$$
$$\tau = \min\left\{t : \widehat{\mathsf{FDP}}(t) \le q\right\}$$



## FDR control

$$\hat{\mathcal{S}} = \{W_j \ge \tau\}$$
$$\tau = \min\left\{t : \widehat{\mathsf{FDP}}(t) \le q\right\}$$



#### Theorem (Barber and C, '15)

For knockoff+

$$\mathsf{FDR} = \mathbb{E}\left[\frac{\# \text{ false positives}}{\# \text{ selections}}\right] \le q$$

## FDR control

$$\hat{\mathcal{S}} = \{W_j \ge \tau\}$$
$$\tau = \min\left\{t : \widehat{\mathsf{FDP}}(t) \le q\right\}$$



#### Theorem (Barber and C, '15)

For knockoff+

$$\mathsf{FDR} = \mathbb{E}\left[\frac{\# \text{ false positives}}{\# \text{ selections}}\right] \le q$$

Robust extension (Barber, C. and Samworth, '18):  $P_X$  not known exactly and knockoffs are exchangeable only w.r.t.  $Q_X \approx P_X$ 

# Knockoffs framework

#### Always FDR control

- $\checkmark$  Under finite sample
- ✓ Any dimension (including p > n)
- $\checkmark \ \ \, {\rm Any \ model \ for \ } Y \mid X$
- ✓ Any black-box

#### The cost?

#### How to construct knockoffs?

Need access to  $P_X$  (not  $P_{Y|X}$ )

#### Let's Make Knockoffs! (Only a Taste)

$P_X$ known:	with M. Sesia and C. Sabatti	
	with S. Bates, L. Janson and W. V	Vang

 $P_X$  unknown: with Y. Romano and M. Sesia

How to make knockoffs?

Input Features X Dist. P<sub>X</sub> How to make knockoffs?

 $Input
Features X
Dist. P_X
<math>X_1 X_2 X_p$ 

 $\underbrace{\begin{array}{c} \underline{Output}\\ Knockoffs \ \tilde{X}\\ Dist. \ P_{\widetilde{X}|X}\\\\\\ \overline{X_1 \ X_2} & \overline{X_p} \end{array}}$ 

How to make knockoffs?



Exchangeability
### Challenges



Challenges



How to make 
$$P_{\tilde{X}|X}$$
?

### The knockoff factory



Gaussian variables

C., Fan, Janson & Lv '16

- Hidden Markov models C., Sabatti & Sesia '17
- Some Bayesian networks Gimenez, Ghorbani & Zou '18
- Pretty much anything (e.g. any graphical model)

Bates, C., Janson & Wang '19

Bates, C., Janson & Wang, '19

- MCMC: Metropolis-Hastings correction
- Importance sampling
- Graphical modeling



Bates, C., Janson & Wang, '19

- MCMC: Metropolis-Hastings correction
- Importance sampling
- Graphical modeling



Bates, C., Janson & Wang, '19

- MCMC: Metropolis-Hastings correction
- Importance sampling
- Graphical modeling



Bates, C., Janson & Wang, '19

- MCMC: Metropolis-Hastings correction
- Importance sampling
- Graphical modeling



Bates, C., Janson & Wang, '19

- MCMC: Metropolis-Hastings correction
- Importance sampling
- Graphical modeling



Bates, C., Janson & Wang, '19

- MCMC: Metropolis-Hastings correction
- Importance sampling
- Graphical modeling



Bates, C., Janson & Wang, '19

- MCMC: Metropolis-Hastings correction
- Importance sampling
- Graphical modeling



Bates, C., Janson & Wang, '19

- MCMC: Metropolis-Hastings correction
- Importance sampling
- Graphical modeling



Bates, C., Janson & Wang, '19

- MCMC: Metropolis-Hastings correction
- Importance sampling
- Graphical modeling





observed variables



• Sample  $Z \sim p(Z \mid X)$  (variation on Viterbi's algorithm)



- Sample  $Z \sim p(Z \mid X)$  (variation on Viterbi's algorithm)
- Sample  $\tilde{Z}_j \sim p(Z_j \mid Z_{-j}, \tilde{Z}_{1:(j-1)})$  for  $j = 1, \dots, p$



- Sample  $Z \sim p(Z \mid X)$  (variation on Viterbi's algorithm)
- Sample  $\tilde{Z}_j \sim p(Z_j \mid Z_{-j}, \tilde{Z}_{1:(j-1)})$  for  $j = 1, \dots, p$
- Sample  $\tilde{X} \sim p(X \mid Z = \tilde{z})$  from emission probs.

### Application to genetic data

C., Sabatti and Sesia ('17)



### Haplotypes and genotypes well modeled by HMMs

Scheet ('06), Marchini ('07, '11), Li ('10), Browning ('10)

### Application to genetic data

C., Sabatti and Sesia ('17)



Haplotypes and genotypes well modeled by HMMs

Scheet ('06), Marchini ('07, '11), Li ('10), Browning ('10)

Wellcome Trust Case Control Consortium  $\approx 5,000$  subjects and 400,000 SNPs Response: Crohn's disease (CD) Northern Finland 1966 Birth Cohort

pprox 4,700 subjects and 330,000 SNPs

Response: lipid levels

Datacat	Number of discoveries				
Dataset	Original study	Knockoffs (average)			
CD	9	22.8			
HDL	5	8			
LDL	6	9.8			

Nominal FDR level at 10%

Selection frequency	SNP (cluster size)	Chr.	Position range (Mb)	Franke et al. '10	WTCCC '07
100%	rs11209026 (2)	1	67.31-67.42	yes	yes
99%	rs6431654 (20)	2	233.94-234.11	yes	yes
98%	rs6688532 (33)	1	169.4-169.65		yes
97%	rs17234657 (1)	5	40.44-40.44	yes	yes
95%	rs11805303 (16)	1	67.31–67.46	yes	yes
91%	rs7095491 (18)	10	101.26-101.32	yes	yes
91%	rs3135503 (16)	16	49.28-49.36	yes	yes
81%	rs7768538 (1145)	6	25.19-32.91	yes	yes
80%	rs6601764 (1)	10	3.85-3.85		yes
75%	rs7655059 (5)	4	89.5-89.53		
73%	rs6500315 (4)	16	49.03-49.07	yes	yes
72%	rs2738758 (5)	20	61.71–61.82	yes	
70%	rs7726744 (46)	5	40.35-40.71	yes	yes
68%	rs11627513 (7)	14	96.61-96.63		
66%	rs4246045 (46)	5	150.07-150.41	yes	yes
62%	rs9783122 (234)	10	106.43-107.61		
61%	rs6825958 (3)	4	55.73-55.77		

Table: SNP clusters found to be important for CD over 100 repetitions of knockoffs.

Selection frequency	SNP (cluster size)	Chr.	Position range (Mb)	Confirmed in Willer et al. '13	Found in Sabatti et al. '09
100%	rs1532085 (4)	15	58.68-58.7	yes	yes
100%	rs7499892 (1)	16	57.01-57.01	yes	yes
100%	rs1800961 (1)	20	43.04-43.04	yes	
99%	rs1532624 (2)	16	56.99-57.01	yes	yes
95%	rs255049 (142)	16	66.41-69.41	yes	yes

Table: SNP clusters found to be important for HDL over 100 repetitions of knockoffs.

Selection frequency	SNP (cluster size)	Chr.	Position range (Mb)	Confirmed in Willer et al. '13	Found in Sabatti et al. '09
99%	rs4844614 (34)	1	207.3-207.88		yes
97%	rs646776 (5)	1	109.8-109.82	yes	yes
97%	rs2228671 (2)	19	11.2-11.21	yes	yes
94%	rs157580 (4)	19	45.4-45.41	yes	yes
92%	rs557435 (21)	1	55.52-55.72	yes	
80%	rs10198175 (1)	2	21.13-21.13	yes	yes
76%	rs10953541 (58)	7	106.48-107.3		
62%	rs6575501 (1)	14	95.64-95.64		

Table: SNP clusters found to be important for LDL over 100 repetitions of knockoffs.

### Deep knockoffs

Romano, Sesia & C. '18



#### $P_X$ unknown? Repurpose deep generative models

#### • KnockoffsGAN

Jordon, Yoon & van der Schaar '19

#### • Knockoffs via VAE

Liu & Zheng '18

### Deep knockoffs: overall view



$$J_{\theta}(\mathbf{X}, \mathbf{\tilde{X}}) = \sum_{j=1}^{p} \mathcal{D}\left((\mathbf{X}, \mathbf{\tilde{X}}), (\mathbf{X}, \mathbf{\tilde{X}})_{\mathrm{Swap}(j)}\right) + \delta \sum_{j=1}^{p} (\mathbf{X}_{j}^{T} \mathbf{\tilde{X}}_{j})^{2}$$

Maximum Mean Discrepancy (MMD) [Gretton et al. ('12)]

- Classic two-sample problem: given **U** and **V**, test whether  $P_U = P_V$
- Discrepancy measure ( $\mathcal{H}$  is RKHS)

$$\mathcal{D}_{\boldsymbol{\phi}} = \|\mathbb{E}_{\boldsymbol{U}}[\boldsymbol{\phi}(\boldsymbol{U})] - \mathbb{E}_{\boldsymbol{V}}[\boldsymbol{\phi}(\boldsymbol{V})]\|_{\mathcal{H}}^2$$

• E.g.  $U, V \in \mathbb{R}$ 

$$\phi(U) = U$$
  $\longrightarrow$  Distance between means  
 $\phi(U) = (U, U^2)$   $\implies$  Error in first two moments

• How to compare higher-order moments?

MMD & the 'kernel-trick' [Gretton et al. ('12)]

$$\mathcal{D}_{\phi} = \|\mathbb{E}_{\boldsymbol{U}}[\phi(\boldsymbol{U})] - \mathbb{E}_{\boldsymbol{V}}[\phi(\boldsymbol{V})]\|_{\mathcal{H}}^2$$

Expand quadratic and replace inner products with kernel operations

 $MMD(P_{U}, P_{V}) = \mathbb{E}_{UU'}[\kappa(U, U')] - 2\mathbb{E}_{UV}[\kappa(U, V)] + \mathbb{E}_{VV'}[\kappa(V, V')]$ 

Characteristic kernel, e.g. Gaussian, implies MMD = 0 iff  $P_U = P_V$ 

MMD & the 'kernel-trick' [Gretton et al. ('12)]

$$\mathcal{D}_{\boldsymbol{\phi}} = \|\mathbb{E}_{\boldsymbol{U}}[\boldsymbol{\phi}(\boldsymbol{U})] - \mathbb{E}_{\boldsymbol{V}}[\boldsymbol{\phi}(\boldsymbol{V})]\|_{\mathcal{H}}^2$$

Expand quadratic and replace inner products with kernel operations

 $MMD(P_{U}, P_{V}) = \mathbb{E}_{UU'}[\kappa(U, U')] - 2\mathbb{E}_{UV}[\kappa(U, V)] + \mathbb{E}_{VV'}[\kappa(V, V')]$ 

Characteristic kernel, e.g. Gaussian, implies MMD = 0 iff 
$$P_U = P_V$$

Unbiased estimate

$$\widehat{\mathsf{MMD}}(\mathbf{U},\mathbf{V}) = \frac{1}{n(n-1)} \sum_{i=1}^{n} \sum_{j \neq i} \kappa \left( U^{i}, U^{j} \right) - \frac{2}{n^{2}} \sum_{i=1}^{n} \sum_{j=1}^{n} \kappa \left( U^{i}, V^{j} \right) + \frac{1}{n(n-1)} \sum_{i=1}^{n} \sum_{j \neq i} \kappa \left( V^{i}, V^{j} \right)$$

## Optimization: stochastic gradient descent

# Generate knockoffs

Evaluate 
$$\tilde{X}^{i} = f_{\theta_{t}}(X^{i}, Z^{i})$$
 for each example. The network  $f_{\theta_{t}}$  is fixed  
Evaluate loss  

$$J_{\theta_{t}}(\mathbf{X}_{t}, \tilde{\mathbf{X}}_{t}) = \sum_{j=1}^{p} \widehat{MMD}((\mathbf{X}, \tilde{\mathbf{X}}), (\mathbf{X}, \tilde{\mathbf{X}})_{\mathrm{Swap}(j)}) + \delta \sum_{j=1}^{p} (\mathbf{X}_{j}^{T} \tilde{\mathbf{X}}_{j})^{2}$$
Update parameters  
 $\theta_{t+1} \leftarrow \theta_{t} - \mu \nabla_{\theta_{t}} J_{\theta_{t}}(\mathbf{X}_{t}, \tilde{\mathbf{X}}_{t})$ 

- Mini-batch SGD
- Random swaps
- Evaluate MMD on disjoint subsets of samples

### Software tools

$\bullet \bullet \bullet < > \bullet$		web.stanford.edu	(
Deep Knockoffs Back	Deep Knockoffs Approximate knockoffs for model-f	ree variable selection.	_
Home	DeepKnockoffs is a software packa	ge for sampling approximate model-X knockoffs using deep generative models	
Tutorial 1	The methods described in the pape	r below are implemented in Python with the help of the PyTorch library.	
Tutorial 2 Tutorial 3	The code is publicly available from <b>Reference</b>	Bitbucket: https://bitbucket.org/msesia/deepknockoffs.	
Tutorial 4	<i>"Deep Knockoffs",</i> Yaniv Romano, Matteo Sesia and	Emmanuel Candès. arXiv:1811.06687 (2018). Link to the paper.	
	License: GPLv3		
	Main features		
	<ul> <li>Generation of approximate ki</li> <li>Goodness-of-fit diagnostics fo</li> <li>Knockoff filter for variable set</li> </ul>	iockoff copies. ir knockoffs. ection.	

#### Authors

Matteo Sesia and Yaniv Romano.

Page generated 2018-11-24 13:08:56 PST, by jemdoc+MathJax.

#### 

	]			iii web.stanford.edu	Ċ	0	00+
Knockoffs	Deep Knock	<b>coffs</b> Is for model-free variable sel	ection.				- 1
Home Tutoral 1 Tutoral 2 Tutoral 3 Tutorial 3 Tutorial 4	ר א ג ג ג ג ג ג ג ג ג ג ג ג ג ג ג ג ג ג	Vumerical experi lotebook written by Ma tanford University, Departmen ast updated on: November 19, e-update of the networks of the graphical graphical processing coad the required lil	ments I (f iteo Sesia and t of Statistics 2018 o allow the numeri anit.	training) J Yaniv Romano cal experimenta described in the pep	er to be reproduced easily.	Punning this code may take a few hours	n
	In [1];	mport numpy as np rom DeepKnockoffs import rom DeepKnockoffs import mport data mport parameters Data generating m We model $X \in \mathbb{R}^p$ as a multiv	KnockoffMach GaussianKnoc	# Train the machin	<i>ne</i> e knockoff mac	thine")	
	In [2]: P P d d d	<pre>orbabon parameter for this c     Number of features     = 100     Load the built-in mu:     The currently availat     - gaugatan : Multivat     - gaugatan : Multivat     - gaugatan : Multivat     - aparase : Multivat     - aparase : Multivat     istribution_params = 1     Initialize the data </pre>		<pre>machine.train(X_t) Fitting the knocke [ 1/ 100], Loss; [ 2/ 100], Loss; [ 3/ 100], Loss; [ 3/ 100], Loss; [ 6/ 100], Loss; [ 6/ 100], Loss; [ 7/ 100], Loss; [ 8/ 100], Loss; [ 9/ 100], Loss; [ 10/ 100], Loss; [ 100],</pre>	rain) off machine : 0.1876, MMD : 0.1454, MMD : 0.1332, MMD : 0.1294, MMD : 0.1260, MMD : 0.1260, MMD : 0.1253, MMD : 0.1238, MMD : 0.1230, MMD	0.1726, Cov: 1.289, 0.1366, Cov: 0.927, 0.1265, Cov: 0.786, 0.1236, Cov: 0.730, 0.1220, Cov: 0.730, 0.1221, Cov: 0.671, 0.1207, Cov: 0.694, 0.1204, Cov: 0.680, 0.1195, Cov: 0.650, 0.1195, Cov: 0.650,	Decorr: 0.2 Decorr: 0.4 Decorr: 0.5 Decorr: 0.5 Decorr: 0.5 Decorr: 0.5 Decorr: 0.5 Decorr: 0.5 Decorr: 0.5 Decorr: 0.5

### HIV Drug Resistance

- Detect mutations in HIV associated with drug resistance (to protease inhibitors)
- Y: log-fold-increase of lab-tested drug resistance
- $X_j$ : presence or absence of mutation #j
- n = 1431, p = 150



### Real X with simulated Y: FDR and Power



X fixed Resample YResample  $\tilde{X}$ 

### Real data example



Method selects variables that mostly correspond to real (replicable) effects

### Summary and challenges



'Wrapper' around black-box algorithm rigorously addresses reproducibility issue

### Summary and challenges



This is not just about not being wrong (irreproducibility)

Technology

# Liking curly fries on Facebook reveals your high IQ

By PHILIPPA WARR

Tuesday 12 March 2013

What you Like on Facebook could reveal your race, age, IQ, sexuality and other personal data, even if you've set that information to "private". Robustness?

Would want predictions to be valid in different samples collected in different circumstances

"Constant conjunction" is a property of causal effects (Hume)

### Fairness: can computer programs be racist and sexist?



Blind application of machine learning runs risk of amplifying biases and prejudices

Identifying variables  $\rightsquigarrow$  chance to scrutinize model built from one sample:

Do we believe these variables are "structurally" important, or are they just reflecting a spurious association in this sample?

Are we learning something about the world or reifying our prejudices?

Guido Rosa/Getty Images/Ikon Images