

Health Informatics

Lecture 8

Samantha Kleinberg
samantha.kleinberg@stevens.edu

Project proposals

- Success \neq hypothesis was true
- Success = result valid
- A proposal marked “complete” doesn’t mean it’s acceptable if I raised serious objections (just received on time for grading purposes)

- How do we know which EHR interface is better?
- Does an NLP method get all mentions of uncertainty or diabetes in a dataset?

Evaluation methods

- Quantitative
 - Compare to Gold standard
 - Classification/prediction
 - Evaluate cost, time
- Qualitative
 - Compare to existing literature
 - Usability
- Replication

Gold standard

- Create labels for data
- Evaluate based on how many items found
- Example: have MDs mark Framingham symptoms in text, evaluate what percentage NLP method can find
- Challenges?
 - Inter-rater reliability
 - Cost
 - Evaluating precision

Inter-rater reliability

How well do annotators agree?

E.g. do clinicians agree on which records mention diabetes?

What's difference delay between two people annotating activities? Does it differ significantly?

How to annotate?

- Domain experts
 - + Expertise
 - + Can help develop annotation scheme
 - \$\$\$, time
 - Smaller pool of annotators
- Crowdsourcing
 - + Large pool of potential annotators
 - + Cheaper
 - Need training
 - May try to game system
 - Privacy issues

Crowdsourcing

Journal List > J Med Internet Res > v.15(4); Apr 2013 > PMC3636329



Journal of Medical Internet Research
The leading peer-reviewed eHealth journal

[Current Issue](#) [Submit](#) [Membership](#) [Editorial Board](#)

J Med Internet Res. 2013 April; 15(4): e73.

PMCID: PMC3636329

Published online 2013 April 2. doi: [10.2196/jmir.2426](https://doi.org/10.2196/jmir.2426)

Web 2.0-Based Crowdsourcing for High-Quality Gold Standard Development in Clinical Natural Language Processing

Monitoring Editor: Gunther Eysenbach

Reviewed by Miguel Luengo-Oroz, Hua Xu, and Lynette Hirschman

[Haijun Zhai](#), PhD,^{#1} [Todd Lingren](#), MA,^{#1} [Louise Deleger](#), PhD,¹ [Qi Li](#), PhD,¹ [Megan Kaiser](#), BA,¹ [Laura Stoutenborough](#), BSN,¹ and [Imre Solti](#), MD, PhD^{✉1}

[Author information](#) ▶ [Article notes](#) ▶ [Copyright and License information](#) ▶

This article has been [cited by](#) other articles in PMC.

Classification, prediction

- Classification: train model to determine who has CHF. Apply to new data to separate CHF/non-CHF.
- Prediction:
 - Use past data, predict who will develop CHF, and follow cohort.
 - Predict next value for measurement
- Challenges?
 - Labels aren't certain
 - Getting labels can be difficult, expensive
 - May have to follow cohort over a long time
 - Predictions are indirect measure, doesn't ensure causality

Cost/time

- Does new method reduce readmissions to hospital?
- Did expenditures go down?
- Are orders sped up by new decision support?
- For algorithms: memory needs, computational time

Replication

- Compare against prior method
- Compare findings in two populations

- Recall earlier discussion paper!

Literature review

- Do findings conflict with prior knowledge?
- Are there known mechanisms that support a hypothesis?
- Problems:
 - Can find supporting results for almost any finding
 - What if it's truly new?

Usability

- Metrics from HCI
 - Surveys, questionnaires, interviews
 - Lab vs in actual usage environment
 - Focus groups?
 - Screen capture
 - Eye tracking

Evaluation strategy depends on..

- What you're evaluating
 - NLP, decision support, new interface?
- Goal
 - Evaluate usability vs. annotation accuracy
- Constraints
 - Time, money, people
 - In-depth study of small group vs. superficial test with large population

Today's paper

Study methods

- **Observational**
 - Cross-sectional
 - Longitudinal cohort study
 - Case-control
- **Intervention**
 - Clinical trials
 - n-of-1 studies

Is Health IT improving outcomes

Say we observe...

- Reduced readmissions
- Fewer errors
- Less time spent on notes

among hospitals using electronic records

How do we know if it's a result of EHR use?

What causes intervention vs what causes effect?

- More skilled doctors may choose a particular procedure
- Healthier patients may use an app

Reasons to randomize

- Subjective assessment
- Delayed effects (side effects)
- Intervention indirectly causing effect

RCTs

- Two groups identical in all respects except for treatment
- Differences between them should then be due to treatment

James Lind and Scurvy

- Divided 12 sailors into groups of two
- Each pair received an addition to their diet
- Citrus led to improved symptoms

But didn't actually randomly assign the treatments

- Individual
- Group-level/cluster
 - Consider true sample size

Comparing two fitness trackers.

Can assign half participants to one and half to the other, but is there a better approach?

What if we randomize two schools (one to trackers, one not)?

Who to randomize?

Testing heartburn medication...

Should population be...?

- Everyone
- People with history of heartburn
- People w/heartburn and not taking drugs that may interact w/proposed one
- People w/o heartburn due to another condition

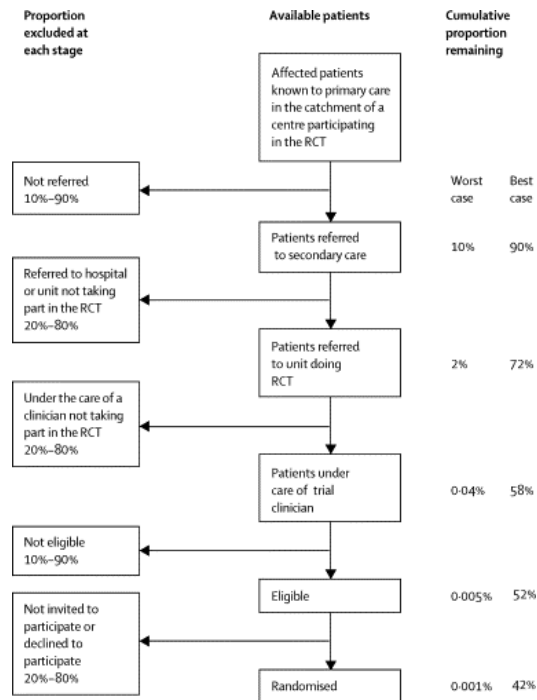


Figure 4. Schematic diagram illustrating effect of multiple stages of selection inherent in clinical practice on proportion of patients in catchment of participating centre entered into an RCT done in secondary care Worst case assumes proportion of patients exc...

Peter M Rothwell

External validity of randomised controlled trials: “To whom do the results of this trial apply?”

null, Volume 365, Issue 9453, 2005, 82–93

[http://dx.doi.org/10.1016/S0140-6736\(04\)17670-8](http://dx.doi.org/10.1016/S0140-6736(04)17670-8)

Randomizing to avoid selection bias but...

Who to analyze?

- People not available for all data collection on outcomes = lost to follow-up
- Ignore subjects with missing data?
- Use all data?

Survival bias

Example: analyze habits of people who've lost weight and kept it off for 2 years.

What's wrong with this?

How to randomize

To assign treatments we could...

- Have researcher alternate between treatment/control for each person
- Flip a coin each time person joins study
- Enclose assignment in set of envelopes that researcher only opens after each subject is enrolled

What are we controlling?

1st RCT compared bed rest (then current standard of care) and streptomycin for treating tuberculosis

Why not compare streptomycin against nothing?

Why placebo

Act of treatment (even with no active ingredient) can lead to real change in outcomes

Just because something changes after treatment, doesn't mean it's because of the treatment!

Decisions aren't treat/don't treat, but WHICH treatment

Control ethics

- If effective treatment exists, should compare against that, not placebo
- Principle: we should in theory not know whether one treatment is better than the other.
- If we know a treatment exists and we withhold it, that's unethical

Control and bias

- Ever see commercials for diet programs or supplements?
 - Just shows A diet better than none, not that this particular one is better than any other
- Remember, choosing population and HOW the other treatment is given (dosage, etc)

Choosing a control

Mimic the process as much as possible

Examples

- sugar pill
- sham acupuncture
- fake surgery
- active placebos

Placebo effect

- Dosage
- Injection vs pill
- Color of pills
- Packaging/branding
- Surgeries

THE LANCET

Volume 299, Issue 7763, 10 June 1972, Pages 1279–1282

Originally published as Volume 1, Issue 7763



Medical Education

DEMONSTRATION TO MEDICAL STUDENTS OF PLACEBO RESPONSES AND NON-DRUG FACTORS

Barry Blackwell, Saul S. Bloomfield, C. Ralph Buncher

[+](#) **Show more**

Abstract

A class experiment for medical students **Summary** was devised to demonstrate the influence of the placebo effect and non-drug factors on response to drugs. The subjects were conditioned to expect sedative or stimulant effects, but all received placebo in one or two blue or pink capsules. Predictions about the size and nature of the placebo response and influence of the non-drug factors were made before the experiment and discussed afterwards. Four of six predictions were fully confirmed. Drug-associated changes were reported by 30% of the subjects and were severe in 1 or 2 individuals. Two capsules produced more noticeable changes than one, and blue capsules were associated with more sedative effects than pink capsules. Students rated the experiment highly both as a learning experience and for its relevance to their future practice of medicine.

Placebos without Deception: A Randomized Controlled Trial in Irritable Bowel Syndrome

Ted J. Kaptchuk^{1,2*}, Elizabeth Friedlander¹, John M. Kelley^{3,4}, M. Norma Sanchez¹, Efi Kokkotou¹, Joyce P. Singer², Magda Kowalczykowski¹, Franklin G. Miller⁵, Irving Kirsch⁶, Anthony J. Lembo¹

1 Beth Israel Deaconess Medical Center, Harvard Medical School, Boston, Massachusetts, United States of America, **2** Osher Research Center, Harvard Medical School, Boston, Massachusetts, United States of America, **3** Psychology Department, Endicott College, Beverly, Massachusetts, United States of America, **4** Massachusetts General Hospital, Harvard Medical School, Boston, Massachusetts, United States of America, **5** Department of Bioethics, National Institutes of Health, Bethesda, Maryland, United States of America, **6** Department of Psychology, University of Hull, Hull, United Kingdom

Abstract

Background: Placebo treatment can significantly influence subjective symptoms. However, it is widely believed that response to placebo requires concealment or deception. We tested whether open-label placebo (non-deceptive and non-concealed administration) is superior to a no-treatment control with matched patient-provider interactions in the treatment of irritable bowel syndrome (IBS).

Methods: Two-group, randomized, controlled three week trial (August 2009-April 2010) conducted at a single academic center, involving 80 primarily female (70%) patients, mean age 47 ± 18 with IBS diagnosed by Rome III criteria and with a score ≥ 150 on the IBS Symptom Severity Scale (IBS-SSS). Patients were randomized to either open-label placebo pills presented as "placebo pills made of an inert substance, like sugar pills, that have been shown in clinical studies to produce significant improvement in IBS symptoms through mind-body self-healing processes" or no-treatment controls with the same quality of interaction with providers. The primary outcome was IBS Global Improvement Scale (IBS-GIS). Secondary measures were IBS Symptom Severity Scale (IBS-SSS), IBS Adequate Relief (IBS-AR) and IBS Quality of Life (IBS-QoL).

Findings: Open-label placebo produced significantly higher mean (\pm SD) global improvement scores (IBS-GIS) at both 11-day midpoint (5.2 ± 1.0 vs. 4.0 ± 1.1 , $p < .001$) and at 21-day endpoint (5.0 ± 1.5 vs. 3.9 ± 1.3 , $p = .002$). Significant results were also observed at both time points for reduced symptom severity (IBS-SSS, $p = .008$ and $p = .03$) and adequate relief (IBS-AR, $p = .02$ and $p = .03$); and a trend favoring open-label placebo was observed for quality of life (IBS-QoL) at the 21-day endpoint ($p = .08$).

Blinding

- Single: Patient unaware of what treatment being received
- Double: Patient + clinician unaware of what treatment being administered/received
- Triple: Double + person analyzing results not aware of which group is which

Example

RCT of two treatments + placebo for multiple sclerosis

Examination by blinded and unblinded neurologists

Unblinded Neurologists showed one treatment beneficial at 6, 12, 24 months, $p\text{-value} < 0.05$

Noseworthy, J. H., Ebers, G. C., Vandervoort, M. K., Farquhar, R. E., Yetisir, E., & Roberts, R. (1994). The impact of blinding on the results of a randomized, placebo-controlled multiple sclerosis clinical trial. *Neurology*, 44(1), 16-16.

Using RCT results

- Tension between study adequately testing hypothesis... and broader applicability
- Study tests one intervention, but in reality do we only change one thing?
- Homogeneous study population isolates cause, but what happens when we move to more varied setting?
- RCT tells us treatment CAN cause something, not that it WILL when we try to use it

Key questions

- Context
 - Are needed features present?
- Will treatment work the same way in another place?
 - Ex: incentives, mosquito nets/fishing
- Side effects
- Interactions
 - Factors that will undermine the cause?

- Intervention: office chairs
- Outcome: weight loss



Weight loss in first test!

Failure in new population. How?

Reviews and Overviews

Why Olanzapine Beats Risperidone, Risperidone Beats Quetiapine, and Quetiapine Beats Olanzapine: An Exploratory Analysis of Head-to-Head Comparison Studies of Second-Generation Antipsychotics

Stephan Heres, M.D.

John Davis, M.D.

Katja Maino, M.D.

Elisabeth Jetzinger, M.D.

Werner Kissling, M.D.

Stefan Leucht, M.D.

Objective: In many parts of the world, second-generation antipsychotics have largely replaced typical antipsychotics as the treatment of choice for schizophrenia. Consequently, trials comparing two drugs of this class—so-called head-to-head studies—are gaining in relevance. The authors reviewed results of head-to-head studies of second-generation antipsychotics funded by pharmaceutical companies to determine if a relationship existed between the sponsor of the trial and the drug favored in the study's overall outcome.

Method: The authors identified head-to-head comparison studies of second-generation antipsychotics through a MEDLINE

sources of bias that could have affected the results in favor of the sponsor's drug.

Results: Of the 42 reports identified by the authors, 33 were sponsored by a pharmaceutical company. In 90.0% of the studies, the reported overall outcome was in favor of the sponsor's drug. This pattern resulted in contradictory conclusions across studies when the findings of studies of the same drugs but with different sponsors were compared. Potential sources of bias occurred in the areas of doses and dose escalation, study entry criteria and study populations, statistics and methods, and reporting of results and wording of findings.

Is randomization always the answer?

Effects of remote, retroactive intercessory prayer on outcomes in patients with bloodstream infection: randomised controlled trial

Leonard Leibovici

Department of
Medicine, Beilinson
Campus, Rabin
Medical Center,
Petah-Tiqva 49100,
Israel
Leonard Leibovici
professor
leibovic@post.
tau.ac.il

Bmj 2001;323:1450-1

Abstract

Objective To determine whether remote, retroactive intercessory prayer, said for a group of patients with a bloodstream infection, has an effect on outcomes.

Design Double blind, parallel group, randomised controlled trial of a retroactive intervention.

Setting University hospital.

Subjects All 3393 adult patients whose bloodstream infection was detected at the hospital in 1990-6.

Intervention In July 2000 patients were randomised to a control group and an intervention group. A remote, retroactive intercessory prayer was said for the well being and full recovery of the intervention group.

Main outcome measures Mortality in hospital, length of stay in hospital, and duration of fever.

Results Mortality was 28.1% (475/1691) in the intervention group and 30.2% (514/1702) in the control group (P for difference = 0.4). Length of stay in hospital and duration of fever were significantly shorter in the intervention group than in the control group (P = 0.01 and P = 0.04, respectively).

Conclusions Remote, retroactive intercessory prayer said for a group is associated with a shorter stay in hospital and shorter duration of fever in patients with a bloodstream infection and should be considered for use in clinical practice.

were included in the study. Bloodstream infection was defined as a positive blood culture (not resulting from contamination) in the presence of sepsis.

In July 2000 a random number generator (Proc Uniform, SAS, Cary, NC, USA) was used to randomise the patients into two groups. A coin was tossed to designate the intervention group. A list of the first names of the patients in the intervention group was given to a person who said a short prayer for the well being and full recovery of the group as a whole. There was no sham intervention.

Three primary outcomes were compared: the number of deaths in hospital, length of stay in hospital from the day of the first positive blood culture to discharge or death, and duration of fever. Patients were defined as having fever on a specific day if one of three temperature measurements taken on that day showed a temperature of > 37.5°C.

The χ^2 test was used to test for the significance of the results shown in the tables. As most of the continuous variables did not have a normal distribution, the Wilcoxon rank sum test was used for comparisons.

Results

Of 3393 patients with a bloodstream infection, 1691 patients were randomised to the intervention group and 1702 to the control group. No patients were lost to

- Inappropriate surrogate outcome measures, so that what is said to have been measured is in fact not (see Gøtzsche et al. 1996; Jaeschke and Sackett 1989).
- Incomplete analysis (Feinberg and Working Group 1998).
To argue from the fact that RCTs have certain advantages, other things being equal, to the claim that the RCT is a gold standard, is like arguing that since being tall makes for a good high-jumper, it follows that a 6' elderly drunkard with a spinal injury is bound to be a better high-jumper than a 5'11" Olympic athlete.
- Problems with inadequate controls (Sackett 1989);
- Inability to detect differences or failures (Sackett 1989);
- False negatives, resulting from small numbers or insensitive outcome measures (Jaeschke and Sackett 1989).

accurate statistical
Medicine

ulting in inade-
on 1998);

ambiguous objec-
Jaeschke and

JASON GROSSMAN AND FIONA J. MACKENZIE

What do we need to know before we can use the results of an RCT?

- Study is internally valid
 - i.e. it can answer the question it aims to answer
- Factors affecting external validity
 - Characteristics of setting
 - Selection of patients
 - Characteristics of patients
 - Follow-up
- Control?
- Blinding – single, double, triple

See also: Rothwell, P. M. (2006). Factors that can affect the external validity of randomised controlled trials. *PLoS Clin Trials*, 1(1), e9.

Next week

- See syllabus for reading!
- Email any requests for reading