

Big Data Is Not Enough: People and Systems Are Needed to Benefit Health and Biomedicine

William Hersh, MD
Professor and Chair
Department of Medical Informatics & Clinical Epidemiology
Oregon Health & Science University
Portland, OR, USA
Email: hersh@ohsu.edu
Web: www.billhersh.info
Blog: <http://informaticsprofessor.blogspot.com>
Twitter: [@williamhersh](https://twitter.com/williamhersh)

References

- Amarasingham, R., Moore, B., Tabak, Y., Drazner, M., Clark, C., Zhang, S., . . . Halm, E. (2010). An automated model to identify heart failure patients at risk for 30-day readmission or death using electronic medical record data. *Medical Care*, 48, 981-988.
- Amarasingham, R., Patel, P., Toto, K., Nelson, L., Swanson, T., Moore, B., . . . Halm, E. (2013). Allocating scarce resources in real-time to reduce heart failure readmissions: a prospective, controlled study. *BMJ Quality & Safety*, 22, 998-1005.
- Amarasingham, R., Patzer, R., Huesch, M., Nguyen, N., & Xie, B. (2014). Implementing electronic health care predictive analytics: considerations and challenges. *Health Affairs*, 33, 1148-1154.
- Anonymous. (2002). Risks and benefits of estrogen plus progestin in healthy postmenopausal women - principal results from the Women's Health Initiative randomized controlled trial. *Journal of the American Medical Association*, 288, 321-333.
- Anonymous. (2014). *IDC Reveals Worldwide Big Data and Analytics Predictions for 2015*. Retrieved from Framingham, MA: <http://bit.ly/IDCBigDataFutureScape2015>
- Anonymous. (2015). Estimating the reproducibility of psychological science. *Science*, 349, aac4716.
- Anonymous. (2016a). *The Cost of Sequencing a Human Genome*. Retrieved from Bethesda, MD: <http://www.genome.gov/sequencingcosts/>
- Anonymous. (2016b). Toward fairness in data sharing. *New England Journal of Medicine*, 375, 405-407.
- Baker, M. (2016). Is there a reproducibility crisis? *Nature*, 533, 452-454.
- Barocas, S., & Selbst, A. (2015). Big data's disparate impact. *California Law Review*, 104, 2016.
- Bates, D., Saria, S., Ohno-Machado, L., Shah, A., & Escobar, G. (2014). Big data in health care: using analytics to identify and manage high-risk and high-cost patients. *Health Affairs*, 33, 1123-1131.
- Begley, C., & Ellis, L. (2012). Raise standards for preclinical cancer research. *Nature*, 483, 531-533.

- Begley, C., & Ioannidis, J. (2015). Reproducibility in science - improving the standard for basic and preclinical research. *Circulation Research*, 116, 116-126.
- Bourgeois, F., Olson, K., & Mandl, K. (2010). Patients treated at multiple acute health care facilities: quantifying information fragmentation. *Archives of Internal Medicine*, 170, 1989-1995.
- Bourne, P., Lorsch, J., & Green, E. (2015). Sustaining the big-data ecosystem. *Nature*, 527, S16-S17.
- Boyd, D., & Crawford, K. (2012). Critical Questions for Big Data. *Information, Communication & Society*, 15, 662-679.
- Broberg, C., Sklenar, J., Burchill, L., Daniels, C., Marelli, A., & Gurvitz, M. (2015). Feasibility of using electronic medical record data for tracking quality indicators in adults with congenital heart disease. *Congenital Heart Disease*, 10, E268-E277.
- Burwell, S. (2015). Setting value-based payment goals - HHS efforts to improve U.S. health care. *New England Journal of Medicine*, 372, 897-899.
- Charlson, M., Wells, M., Ullman, R., King, F., & Shmukler, C. (2014). The Charlson comorbidity index can be used prospectively to identify patients who will incur high future costs. *PLoS ONE*, 9(12), e112479.
- Cho, I., Park, I., Kim, E., Lee, E., & Bates, D. (2013). Using EHR data to predict hospital-acquired pressure ulcers: A prospective study of a Bayesian Network model. *International Journal of Medical Informatics*, 82, 1059-1067.
- Collins, F., & Varmus, H. (2015). A new initiative on precision medicine. *New England Journal of Medicine*, 372, 793-795.
- Davenport, T., & Patil, D. (2012, October, 2012). Data Scientist: The Sexiest Job of the 21st Century. *Harvard Business Review*.
- deLusignan, S., & vanWeel, C. (2005). The use of routinely collected computer data for research in primary care: opportunities and challenges. *Family Practice*, 23, 253-263.
- DesRoches, C., Painter, M., & Jha, A. (2015). *Health Information Technology in the United States 2015 - Transition to a Post-HITECH World*. Retrieved from Princeton, NJ: <http://www.rwjf.org/en/library/research/2015/09/health-information-technology-in-the-united-states-2015.html>
- Donoho, D. (2015). *50 years of Data Science*. Retrieved from Princeton NJ: <https://dl.dropboxusercontent.com/u/23421017/50YearsDataScience.pdf>
- Donzé, J., Aujesky, D., Williams, D., & Schnipper, J. (2013). Potentially avoidable 30-day hospital readmissions in medical patients: derivation and validation of a prediction model. *JAMA Internal Medicine*, 173, 632-638.
- Dwan, K., Gamble, C., Williamson, P., & Kirkham, J. (2013). Systematic review of the empirical evidence of study publication bias and outcome reporting bias - an updated review. *PLoS ONE*, 8(7), e66844.
- Dzau, V., McClellan, M., & McGinnis, J. (2016). Vital Directions for Health and Health Care: an initiative of the National Academy of Medicine. *Journal of the American Medical Association*, 316, 711-712.
- Eklund, A., Nichols, T., & Knutsson, H. (2016). Cluster failure: why fMRI inferences for spatial extent have inflated false-positive rates. *Proceedings of the National Academy of Sciences*, 113, 7900-7905.

- Erskine, A., Karunakaran, P., Slotkin, J., & Feinberg, D. (2016). How Geisinger Health System Uses Big Data to Save Lives. *Harvard Business Review*.
- Evans, R., Benuzillo, J., Horne, J., Lloyd, J., Bradshaw, A., Budge, D., . . . Lappé, D. (2016). Automated identification and predictive tools to help identify high-risk heart failure patients: pilot evaluation. *Journal of the American Medical Informatics Association*, Epub ahead of print.
- Finnell, J., Overhage, J., & Grannis, S. (2011). *All health care is not local: an evaluation of the distribution of emergency department care delivered in Indiana*. Paper presented at the AMIA Annual Symposium Proceedings, Washington, DC.
- FitzHenry, F., Murff, H., Matheny, M., Gentry, N., Fielstein, E., Brown, S., . . . Speroff, T. (2013). Exploring the frontier of electronic health record surveillance: the case of postoperative complications. *Medical Care*, 51, 509-516.
- Fung, K. (2014, March 25, 2014). Google Flu Trends' Failure Shows Good Data > Big Data. *Harvard Business Review*.
- Geifman, N., & Butte, A. (2016). *Do cancer clinical trial populations truly represent cancer patients? A comparison of open clinical trials to the Cancer Genome Atlas*. Paper presented at the Pacific Symposium on Biocomputing, Kohala Coast, HI.
- Gilbert, P., Rutland, M., & Brockopp, D. (2013). Redesigning the work of case management: testing a predictive model for readmission. *American Journal of Managed Care*, 19(11 Spec No. 10), eS19-eSP25.
- Gildersleeve, R., & Cooper, P. (2013). Development of an automated, real time surveillance tool for predicting readmissions at a community hospital. *Applied Clinical Informatics*, 4, 153-169.
- Gold, M., & McLaughlin, C. (2016). Assessing HITECH implementation and lessons: 5 years later. *Milbank Quarterly*, 94, 654-687.
- Halamka, J. (2013, February 12, 2013). The "Post EHR" Era. Retrieved from <http://geekdoctor.blogspot.com/2013/02/the-post-ehr-era.html>
- Haug, C. (2013). The downside of open-access publishing. *New England Journal of Medicine*, 368, 791-793.
- Hebert, C., Shivade, C., Foraker, R., Wasserman, J., Roth, C., Mekhjian, H., . . . Embi, P. (2014). Diagnosis-specific readmission risk prediction using electronic health data: a retrospective cohort study. *BMC Medical Informatics & Decision Making*, 14, 65.
- Hersh, W. (2013). What is a Thinking Informatician to Think of IBM's Watson? Retrieved from <http://informaticsprofessor.blogspot.com/2013/06/what-is-thinking-informatician-to-think.html>
- Hersh, W. (2015). What is the Difference (If Any) Between Informatics and Data Science? Retrieved from <http://informaticsprofessor.blogspot.com/2015/07/what-is-difference-if-any-between.html>
- Hersh, W. (2016). Generalizability and Reproducibility of Scientific Literature and the Limits to Machine Learning.
- Hersh, W., Weiner, M., Embi, P., Logan, J., Payne, P., Bernstam, E., . . . Saltz, J. (2013). Caveats for the use of operational electronic health record data in comparative effectiveness research. *Medical Care*, 51(Suppl 3), S30-S37.
- Horner, P., & Basu, A. (2012, January/February 2012). Analytics & the future of healthcare. *Analytics*, 11-18.

- Hripcsak, G., & Albers, D. (2012). Next-generation phenotyping of electronic health records. *Journal of the American Medical Informatics Association*, 20, 117-121.
- Hripcsak, G., Ryan, P., Duke, J., Shah, N., Park, R., Huser, V., . . . Madigan, D. (2016). Characterizing treatment pathways at scale using the OHDSI network. *Proceedings of the National Academy of Sciences*, 113, 7329-7336.
- Ioannidis, J. (2005a). Contradicted and initially stronger effects in highly cited clinical research. *Journal of the American Medical Association*, 294, 218-228.
- Ioannidis, J. (2005b). Why most published research findings are false. *PLoS Medicine*, 2(8), e124.
- Joppa, L., McInerney, G., Harper, R., Salido, L., Takeda, K., O'Hara, K., . . . Emmott, S. (2013). Troubling trends in scientific software use. *Science*, 340, 814-815.
- Kesselheim, A., & Avorn, J. (2017). New "21st Century Cures" legislation: speed and ease vs science. *Journal of the American Medical Association*, Epub ahead of print.
- Khurana, H., Groves, R., Simons, M., Martin, M., Stoffer, B., Kou, S., . . . Parthasarathy, S. (2016). Real-time automated sampling of electronic medical records predicts hospital mortality. *American Journal of Medicine*, 129, 688-698.
- Kim, C., & Prasad, V. (2015). Strength of validation for surrogate end points used in the US Food and Drug Administration's approval of oncology drugs. *Mayo Clinic Proceedings*, Epub ahead of print.
- Krumholz, H. (2014). Big data and new knowledge in medicine: the thinking, training, and tools needed for a learning health system. *Health Affairs*, 33, 1163-1170.
- Kush, R., & Goldman, M. (2014). Fostering responsible data sharing through standards. *New England Journal of Medicine*, 370, 2163-2165.
- Longo, D., & Drazen, J. (2016). Data sharing. *New England Journal of Medicine*, 374, 276-277.
- Lowes, L., Noritz, G., Newmeyer, A., Embi, P., Yin, H., & Smoyer, W. (2016). 'Learn From Every Patient': implementation and early results of a learning health system. *Developmental Medicine & Child Neurology*, 59, 183-191.
- Manor-Shulman, O., Beyene, J., Frndova, H., & Parshuram, C. (2008). Quantifying the volume of documented clinical information in critical illness. *Journal of Critical Care*, 23, 245-250.
- Manyika, J., Chui, M., Brown, B., Bughin, J., Dobbs, R., Roxburgh, C., & Byers, A. (2011). *Big data: The next frontier for innovation, competition, and productivity*. Retrieved from http://www.mckinsey.com/insights/business_technology/big_data_the_next_frontier_for_innovation
- Merali, Z. (2010). Computational science: ...Error. *Nature*, 467, 775-777.
- Moher, D., & Moher, E. (2016). Stop predatory publishers now: act collaboratively. *Annals of Internal Medicine*, 164, 616-617.
- Murphy, D., Laxmisan, A., Reis, B., Thomas, E., Esquivel, A., Forjuoh, S., . . . Singh, H. (2014). Electronic health record-based triggers to detect potential delays in cancer diagnosis. *BMJ Quality & Safety*, 23, 8-16.
- Murphy, D., Wu, L., Thomas, E., Forjuoh, S., Meyer, A., & Singh, H. (2015). Electronic trigger-based intervention to reduce delays in diagnostic evaluation for cancer: a cluster randomized controlled trial. *Journal of Clinical Oncology*, 33, 3560-3567.
- Prasad, V., & Cifu, A. (2015). *Ending Medical Reversal: Improving Outcomes, Saving Lives*. Baltimore, MD: Johns Hopkins University Press.

- Prasad, V., Kim, C., Burotto, M., & Vandross, A. (2015). The strength of association between surrogate end points and survival in oncology: a systematic review of trial-level meta-analyses. *JAMA Internal Medicine*, 175, 1389-1398.
- Prasad, V., Vandross, A., Toomey, C., Cheung, M., Rho, J., Quinn, S., . . . Cifu, A. (2013). A decade of reversal: an analysis of 146 contradicted medical practices. *Mayo Clinic Proceedings*, 88, 790-798.
- Prieto-Centurion, V., Rolle, A., Au, D., Carson, S., Henderson, A., Lee, T., . . . Krishnan, J. (2014). Multicenter study comparing case definitions used to identify patients with chronic obstructive pulmonary disease. *American Journal of Respiratory and Critical Care Medicine*, 190, 989-995.
- Rajkomar, A., Yim, J., Grumbach, K., & Parekh, A. (2016). Weighting primary care patient panel size: a novel electronic health record-derived measure using machine learning. *JMIR Medical Informatics*, 4(4), e29.
- Randhawa, A., Babalola, O., Henney, Z., Miller, M., Nelson, T., Oza, M., . . . So, S. (2016). A collaborative assessment among 11 pharmaceutical companies of misinformation in commonly used online drug information compendia. *Annals of Pharmacotherapy*, 50, 352-359.
- Richards, N., & King, J. (2014). Big data ethics. *Wake Forest Law Review*, 49, 393-432.
- Rothman, M., Rimar, J., Coonan, S., Allegretto, S., & Balcezak, T. (2015). Mortality reduction associated with proactive use of EMR-based acuity score by an RN team at an urban hospital. *BMJ Quality & Safety*, 24, 734-735.
- Sainani, K. (2011, September 1, 2011). Error! – What Biomedical Computing Can Learn From Its Mistakes. *Biomedical Computation Review*.
- Schank, R. (2016). The fraudulent claims made by IBM about Watson and AI. They are not doing "cognitive computing" no matter how many times they say they are. Retrieved from <http://www.rogerschank.com/fraudulent-claims-made-by-IBM-about-Watson-and-AI>
- Schoenfeld, J., & Ioannidis, J. (2013). Is everything we eat associated with cancer? A systematic cookbook review. *American Journal of Clinical Nutrition*, 97, 127-134.
- Shadmi, E., Flaks-Manov, N., Hoshen, M., Goldman, O., Bitterman, H., & Balicer, R. (2015). Predicting 30-day readmissions with preadmission electronic health record data. *Medical Care*, 53, 283-289.
- Singer, D., Jacks, T., & Jaffee, E. (2016). A U.S. "Cancer Moonshot" to accelerate cancer research. *Science*, 353, 1105-1106.
- Stead, W., Searle, J., Fessler, H., Smith, J., & Shortliffe, E. (2011). Biomedical informatics: changing what physicians need to know and how they learn. *Academic Medicine*, 86, 429-434.
- Strom, B., Buyse, M., Hughes, J., & Knoppers, B. (2016). Data sharing — is the juice worth the squeeze? *New England Journal of Medicine*, 375, 1608-1609.
- Taichman, D., Backus, J., Baethge, C., Bauchner, H., deLeeuw, P., Drazen, J., . . . Wu, S. (2016). Sharing clinical trial data: a proposal from the International Committee of Medical Journal Editors. *New England Journal of Medicine*, 374, 384-386.
- Tenenbaum, J., Avillach, P., Benham-Hutchins, M., Breitenstein, M., Crowgey, E., Hoffman, M., . . . Freimuth, R. (2016). An informatics research agenda to support precision medicine: seven key areas. *Journal of the American Medical Informatics Association*, 23, 791-795.

- Tien, M., Kashyap, R., Wilson, G., Hernandez-Torres, V., Jacob, A., Schroeder, D., & Mantilla, C. (2015). Retrospective derivation and validation of an automated electronic search algorithm to identify post operative cardiovascular and thromboembolic complications. *Applied Clinical Informatics*, 6, 565-576.
- Voorhees, E., & Hersh, W. (2012). *Overview of the TREC 2012 Medical Records Track*. Paper presented at the The Twenty-First Text REtrieval Conference Proceedings (TREC 2012), Gaithersburg, MD.
- Weng, C., Li, Y., Ryan, P., Zhang, Y., Liu, F., Gao, J., . . . Hripcsak, G. (2014). A distribution-based method for assessing the differences between clinical trial target populations and patient populations in electronic health records. *Applied Clinical Informatics*, 5, 463-479.
- White, J., Briggs, J., & Mandel, J. (2016). NIH and ONC Launch the Sync for Science (S4S) Pilot: Enabling Individual Health Data Access and Donation. Retrieved from <https://www.healthit.gov/buzz-blog/health-innovation/nih-and-onc-launch-the-sync-for-science-pilot/>
- Young, N., Ioannidis, J., & Al-Ubaydli, O. (2008). Why current publication practices may distort science. *PLoS Medicine*, 5(10), e201.
- Zheng, L., Wang, Y., Hao, S., Shin, A., Jin, B., Ngo, A., . . . Ling, X. (2016). Web-based real-time case finding for the population health management of patients with diabetes mellitus: a prospective validation of the natural language processing-based algorithm with statewide electronic medical records. *JMIR Medical Informatics*, 4(4), e37.

Big Data Is Not Enough: People and Systems Are Needed to Benefit Health and Biomedicine

William Hersh, MD
Professor and Chair
Department of Medical Informatics & Clinical Epidemiology
Oregon Health & Science University
Portland, OR, USA
Email: hersh@ohsu.edu
Web: www.billhersh.info
Blog: <http://informaticsprofessor.blogspot.com>
Twitter: [@williamhersh](https://twitter.com/williamhersh)

1



Big Data is not enough

- Many use cases for Big Data
- Growing quantity of data available at decreasing cost
- Much demonstration of predictive ability; less so of value
- Many caveats for different types of biomedical data
- Effective solutions require people and systems

2



Many use cases for Big Data in medicine (Bates, 2014)

- High-cost patients – looking for ways to intervene early
- Readmissions – preventing
- Triage – appropriate level of care
- Decompensation – when patient's condition worsens
- Adverse events – awareness
- Treatment optimization – especially for diseases affecting multiple organ systems

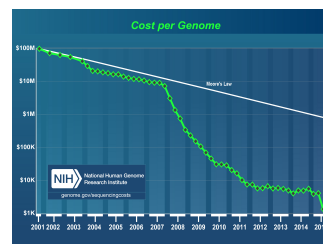
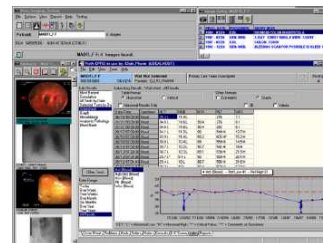


3



Growing quantity at increasingly lower cost of data

- Last half-decade has seen dramatic growth in adoption of electronic health record (EHR) by hospitals (96%) and physicians (83%) (DesRoches, 2015; Gold, 2016)
- Cost of genome sequencing has fallen faster than Moore's Law (NHGRI, 2016)
- Proliferation of other data sources
 - Imaging
 - Wearables
 - Web and social media



4

Important data-related initiatives from US government

- Big Data to Knowledge (BD2K) (Bourne, 2015) – <https://datascience.nih.gov>
- Sync for Science (White, 2016) – <http://syncfor.science>
- Vital Directions for Health and Health Care (Dzau, 2016)
- Precision Medicine Initiative (Collins, 2015) – <https://www.nih.gov/research-training/allofus-research-program>
- Cancer Moonshot (Singer, 2016) – <https://www.cancer.gov/research/key-initiatives/moonshot-cancer-initiative>
- 21st Century Cures (Kesselheim, 2017)

5



Rationale

- Growing quantity and complexity of healthcare data through EHR capture, genomics, and other sources require more decision support (Stead, 2011)
- With shift of payment from “volume to value,” healthcare organizations will need to manage information better to deliver better care (Horner, 2012; Burwell, 2015)
- New care delivery models (e.g., accountable care organizations) will require better access to data (e.g., health information exchange, HIE)
 - Halamka (2013): ACO = HIE + analytics

6



Ever-growing number of studies demonstrating predictive ability

- Using EHR data to predict patients at risk for readmission (Amarasingham, 2010; Donzé, 2013; Gildersleeve, 2013; Hebert, 2014; Shadmi, 2015)
- Identifying patients who might be eligible for participation in clinical studies (Voorhees, 2012)
- Detecting postoperative complications (FitzHenry, 2013; Tien, 2015)
- Detecting potential delays in cancer diagnosis (Murphy, 2014)
- Predicting future patient costs (Charlson, 2014)

7



Predictive studies (cont.)

- Optimizing primary care physician panel size (Rajkomar, 2016)
- Real-time alerting of mortality risk and prolonged hospitalization from EHR data (Khurana, 2016)
- Elucidating treatment pathways for common diseases (Hripcsak, 2016)
- NLP-based case-finding algorithm of HIE data increased detection of diabetes cases (Zheng, 2016)
- The list goes on and on ...

8



BUT, studies demonstrating improved patient outcomes are fewer

- Readmission tool applied with case management reduced readmissions (Gilbert, 2013)
- Bayesian network model embedded in EHR to predict hospital-acquired pressure ulcers led to tenfold reduction in ulcers and one-third reduction in intensive care unit length of stay (Cho, 2013)
- Readmission risk tool intervention reduced risk of readmission for patients with congestive heart failure but not those with acute myocardial infarction or pneumonia (Amarasingham, 2013)
- Use of EHR-based acuity score allowed intervention that reduced in-hospital mortality from 1.9% to 1.3% (Rothman, 2015)
- Tool to reduce delay in cancer diagnosis led to earlier diagnosis for colorectal and prostate cancer (Murphy, 2015)

9



Newer studies of outcomes

- Use of predictive report based on NLP tool reduced time in discharge planning meetings and 30-day all-cause mortality although not cost or readmissions (Evans, 2016)
- Development and use of a universal data architecture at Geisinger has led to successes in (Erskine, 2016)
 - Closing loop on appropriate treatment and lack of follow-up
 - Early detection and treatment of sepsis
 - Monitoring and control of surgery costs and outcomes
- In cohort of children with cerebral palsy, implementation of a learning health system led to (Lowes, 2016)
 - 43% reduced hospital days
 - 30% reduction in emergency department visits
 - 210% reduction in healthcare costs

10



Some challenges for analytical use of clinical (EHR) data

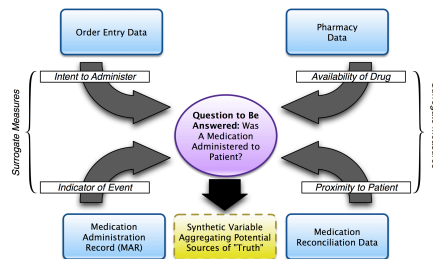
- Data quality and accuracy is not a top priority for busy clinicians (de Lusignan, 2005)
- Data quantity can be overwhelming – average pediatric ICU patient generates 1348 information items per 24 hours (Manor-Shulman, 2008)
- Patients get care at different institutions (Bourgeois, 2010; Finnell, 2011)
- Much data is “locked” in text (Hripcsak, 2012)
- EHRs of academic medical centers not easy to combine for aggregation (Broberg, 2015)

11



Caveats for use of operational EHR data (Hersh, 2013) – may be

- Inaccurate
- Incomplete
- Transformed in ways that undermine meaning
- Unrecoverable
- Of unknown provenance
- Of insufficient granularity
- Incompatible with research protocols



12



Many “idiosyncrasies” of clinical data (Hersh, 2013)

- “Left censoring” – First instance of disease in record may not be when first manifested
- “Right censoring” – Data source may not cover long enough time interval
- Data might not be captured from other clinical (other hospitals or health systems) or non-clinical (OTC drugs) settings
- Bias in testing or treatment
- Institutional or personal variation in practice or documentation styles
- Inconsistent use of coding or standards

13



Information from scientific publications can also be problematic

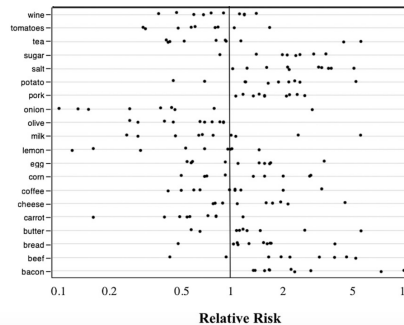
- Science, driven by experimentation, is the best source of truth, but just because something is written in a journal article does not mean it is true
 - Winner’s curse (Ioannidis, 2005; Young, 2008) leads to publication bias (Dwan, 2013)
 - Reproducibility (Begley, 2012; Science, 2015; Begley, 2015; Baker, 2016)
 - Clinical trials may not be representative of patient populations (Weng, 2014; Prieto-Centurion, 2014; Geifman, 2016)
 - Use of surrogate endpoints may distort efficacy (Kim, 2015)
 - Reversal (Ioannidis, 2005; Prasad, 2013; Prasad, 2015)
 - Erroneous information in reference materials (Randhawa, 2015)
 - Outright fraud not infrequent (RetractionWatch.com), may be driven by predatory publishing (Haug, 2013; Moher, 2016)

14



Results can be misleading, conflicting, or hyped

- Observational studies can mislead us, e.g., Women's Health Initiative (JAMA, 2002)
- Observational studies do not discern cause and effect, e.g., diet and cancer (Schoenfeld, 2013)
- Hype about new technologies not yet fully assessed, e.g., IBM Watson – much promise but much hype (Hersh, 2013; Hersh, 2016; Schank, 2016)



15



Biomedical researchers are not necessarily good software engineers

- Many scientific researchers write code but are not always well-versed in best practices of testing and error detection (Merali, 2010)
- Scientists have history of relying on incorrect data or models (Sainani, 2011)
- They may also not be good about selection of best software packages for their work (Joppa, 2013)
- 3000 of 40,000 studies using fMRI may have false-positive results due to faulty algorithms and bugs (Eklund, 2016)

```
>
>
...SCIENTISTS AND THEIR
SOFTWARE
A survey of nearly 2,000
researchers showed how coding
has become an important part of
the research toolkit, but it
also revealed some potential
problems.
> 45% said scientists spend
more time today developing
software than five years ago."
> 38% of scientists spend at
least one fifth of their time
developing software.
> Only 47% of scientists
have a good understanding of
software testing.
> Only 34% of scientists
think that formal training
in developing software is
important.
```

16



Should there be more sharing of scientific data? Yes, but ...

- Came to fore with ICMJE guidelines (Taichman, 2016) and NEJM “research parasites” editorial (Longo, 2016)
 - Pro: fairness to funders (taxpayers) and subjects (patients)
 - Con: researchers who carried out the heavy work need period of embargo and protection from misuse of their data (ICFTDS, 2016); costs of curating and organizing 27K clinical trials per year; amount of actual use modest (Strom, 2016)
- Informatics issues: need for attention to standards (Kush, 2014); workflows, patient engagement (Tennenbaum, 2016)

17



Other concerns

- Boyd (2012) – critical questions for Big Data
 - Big Data changes the definition of knowledge
 - Claims to objectivity and accuracy are misleading
 - Bigger data are not always better data
 - Taken out of context, Big Data loses its meaning
 - Just because it is accessible does not make it ethical
 - Limited access to Big Data creates new digital divides
- Fung (2014) – Big Data is OCCAM
 - **O**bservational
 - **L**acking **C**ontrols
 - **S**eemingly **C**omplete
 - **A**dapted
 - **M**erged
- Big Data not neutral; reflects our values and priorities (Richards, 2014; Barocas, 2015)

18



Big Data requires more than the data; also takes people

- Data scientists – the “sexiest profession of the 21st century” (Davenport, 2012)
- McKinsey (Manyika, 2011) – need in US in all industries (not just healthcare) for
 - 140,000-190,000 individuals who have “deep analytical talent”
 - 1.5 million “data-savvy managers needed to take full advantage of big data”
- Similar analysis by IDC (2014) of need for 180,000 with “deep” talent and 5-fold around with skills in data management and interpretation

19



Big Data also requires systems

- Infrastructure (Amarasingham, 2014)
 - Stakeholder engagement
 - Human subjects research protection
 - Protection of patient privacy
 - Data assurance and quality
 - Interoperability of health information systems
 - Transparency
 - Sustainability
- New models of thinking and training users of data (Krumholz, 2014)

20



Some axes to grind

- Is data science really new or different?
 - Statisticians (Donoho, 2016) and informaticians (Hersh, 2015) have been doing some of this for a long time
- Will Big Data transform medicine?
 - In some areas, but need more demonstration of value than ability to predict
- How can we optimize its use?
 - Research focused on its applications and their outcomes
 - Don't oversell it, especially to clinicians



21



Much promise for Big Data in Health and Biomedicine, but need

- Other aspects of informatics
 - Robust EHRs and other clinical data sources
 - Standards and interoperability
 - Health information exchange
 - Usability of clinical systems
- Improved completeness and quality of data
- Research demonstrating how best applied to improve health and outcomes
- Human expertise and systems to apply and disseminate

22

