Applying Information Retrieval to the Electronic Health Record for Cohort Discovery and Rare Disease Detection

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

Distinguished Lecture Series, UCLA Biomedical Data Science Program February 13, 2020

References

Buckley, C and Voorhees, EM (2005). Retrieval System Evaluation. <u>TREC: Experiment and Evaluation in Information Retrieval</u>. E. Voorhees and D. Harman. Cambridge, MA, MIT Press: 53-75.

Chamberlin, SR, Bedrick, SD, et al. (2019). Evaluation of patient-level retrieval from electronic health record data for a cohort discovery task. *medRxiv*: 19005280. https://www.medrxiv.org/content/10.1101/19005280v2

Demner-Fushman, D, Abhyankar, S, et al. (2012). NLM at TREC 2012 Medical Records Track. *The Twenty-First Text REtrieval Conference Proceedings (TREC 2012)*, Gaithersburg, MD. National Institute for Standards and Technology

http://trec.nist.gov/pubs/trec21/papers/NLM.medical.final.pdf

Demner-Fushman, D, Abhyankar, S, et al. (2011). A knowledge-based approach to medical records retrieval. *The Twentieth Text REtrieval Conference Proceedings (TREC 2011)*, Gaithersburg, MD. National Institute for Standards and Technology

Edinger, T, Cohen, AM, et al. (2012). Barriers to retrieving patient information from electronic health record data: failure analysis from the TREC Medical Records Track. *AMIA 2012 Annual Symposium*, Chicago, IL. 180-188.

Halamka, JD (2020). A New Model for Sharing Insights While Protecting Privacy. <u>Dispatch</u> <u>from the Digital Health Frontier</u>. <u>http://geekdoctor.blogspot.com/2020/01/a-new-model-</u> <u>for-sharing-insights-while.html</u>

Hanbury, A, Müller, H, et al. (2015). Evaluation-as-a-service: overview and outlook. *arXiv.org*: arXiv:1512.07454. <u>https://arxiv.org/abs/1512.07454</u>

Harman, D (2011). <u>Information Retrieval Evaluation</u>. San Rafael, CA, Morgan & Claypool. Harman, DK (2005). The TREC Ad Hoc Experiments. <u>TREC: Experiment and Evaluation in</u> <u>Information Retrieval</u>. E. Voorhees and D. Harman. Cambridge, MA, MIT Press: 79-98. Hersh, W and Voorhees, E (2009). TREC genomics special issue overview. *Information Retrieval*. 12: 1-15.

Hersh, WR (2009). <u>Information Retrieval: A Health and Biomedical Perspective (3rd Edition)</u>. New York, NY, Springer.

Hersh, WR, Crabtree, MK, et al. (2002). Factors associated with success for searching MEDLINE and applying evidence to answer clinical questions. *Journal of the American Medical Informatics Association*. 9: 283-293.

Hersh, WR and Greenes, RA (1990). SAPHIRE: an information retrieval environment featuring concept-matching, automatic indexing, and probabilistic retrieval. *Computers and Biomedical Research*. 23: 405-420.

Hersh, WR and Hickam, DH (1995). An evaluation of interactive Boolean and natural language searching with an on-line medical textbook. *Journal of the American Society for Information Science*. 46: 478-489.

Hersh, WR and Hickam, DH (1998). How well do physicians use electronic information retrieval systems? A framework for investigation and review of the literature. *Journal of the American Medical Association*. 280: 1347-1352.

Hersh, WR, Hickam, DH, et al. (1994). A performance and failure analysis of SAPHIRE with a MEDLINE test collection. *Journal of the American Medical Informatics Association*. 1: 51-60.

Jarvelin, K and Kekalainen, J (2002). Cumulated gain-based evaluation of IR techniques. *ACM Transactions on Information Systems*. 20: 422-446.

King, B, Wang, L, et al. (2011). Cengage Learning at TREC 2011 Medical Track. *The Twentieth Text REtrieval Conference Proceedings (TREC 2011)*, Gaithersburg, MD. National Institute for Standards and Technology

Matheny, M, Israni, ST, et al., Eds. (2019). <u>Artificial Intelligence in Health Care: The Hope,</u> <u>the Hype, the Promise, the Peril</u>. Washington, DC, National Academy of Medicine.

Roberts, K, Simpson, M, et al. (2016). State-of-the-art in biomedical literature retrieval for clinical cases: a survey of the TREC 2014 CDS track. *Information Retrieval Journal*. 19: 113-148.

Roegiest, A and Cormack, GV (2016). An architecture for privacy-preserving and replicable high-recall retrieval experiments. *Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval*, Pisa, Italy. 1085-1088.

Safran, C, Bloomrosen, M, et al. (2007). Toward a national framework for the secondary use of health data: an American Medical Informatics Association white paper. *Journal of the American Medical Informatics Association*. 14: 1-9.

Sardh, E, Harper, P, et al. (2019). Phase 1 trial of an RNA interference therapy for acute intermittent porphyria. *New England Journal of Medicine*. 380: 549-558.

Voorhees, E and Hersh, W (2012). Overview of the TREC 2012 Medical Records Track. *The Twenty-First Text REtrieval Conference Proceedings (TREC 2012)*, Gaithersburg, MD.

National Institute of Standards and Technology

http://trec.nist.gov/pubs/trec21/papers/MED12OVERVIEW.pdf

Voorhees, EM (2013). The TREC Medical Records Track. *Proceedings of the International Conference on Bioinformatics, Computational Biology and Biomedical Informatics,* Washington, DC. 239-246.

Voorhees, EM and Tong, RM (2011). Overview of the TREC 2011 Medical Records Track. *The Twentieth Text REtrieval Conference Proceedings (TREC 2011)*, Gaithersburg, MD. National Institute of Standards and Technology

Wang, Y, Wen, A, et al. (2019). Test collections for electronic health record-based clinical information retrieval. *JAMIA Open.* 2: 360-368.

https://academic.oup.com/jamiaopen/advance-

article/doi/10.1093/jamiaopen/ooz016/5510566

Wu, S, Liu, S, et al. (2017). Intra-institutional EHR collections for patient-level information retrieval. *Journal of the American Society for Information Science & Technology*. 68: 2636-2648.

Zhu, D, Wu, ST, et al. (2014). Using large clinical corpora for query expansion in textbased cohort identification. *Journal of Biomedical Informatics*. 49: 275-281.































Failure analysis for 2011 topics (Edinger, 2012)

| | Number | Number |
|--|-----------|-----------|
| Reasons for Incorrect Retrieval | of Visits | of Topics |
| Visits Judged Not Relevant | | |
| Topic terms mentioned as future possibility | 16 | 9 |
| Topic symptom/condition/procedure done in the past | 22 | 9 |
| All topic criteria present but not in the time/sequence specified by the topic description | 19 | 6 |
| Most, but not all, required topic criteria present | 17 | 8 |
| Topic terms denied or ruled out | 19 | 10 |
| Notes contain very similar term confused with topic term | 13 | 11 |
| Non-relevant reference in record to topic terms | 37 | 18 |
| Topic terms not present-unclear why record was ranked highly | 14 | 8 |
| Topic present-record is relevant-disagree with expert judgment | 25 | 11 |
| Visits Judged Relevant | | |
| Topic not present-record is not relevant-disagree with expert judgment | 44 | 21 |
| Topic present in record but overlooked in search | 103 | 27 |
| Visit notes used a synonym or lexical variant for topic terms | 22 | 10 |
| Topic terms not named in notes and must be inferred | 3 | 2 |
| Topic terms present in diagnosis list but not visit notes | 5 | 5 |









Original EHR data – 100K OHSU patients having ≥3 visits

| Туре | Patients | Encounters | Records | Average | Median | Max |
|-----------------------|---------------------|------------------------|------------|---------|--------|--------|
| Administered Meds | 47,208 | 125,831 | 6,497,157 | 51.634 | 6 | - |
| Ambulatory Encounters | 99,965 | 3,760,205 | 3,760,205 | - | - | - |
| Current Meds | 92,783 | - | 31,997,402 | 344.863 | 64 | 20,102 |
| Demographics | 99,965 | - | | | | |
| Encounter Attributes | <mark>99,965</mark> | <mark>6,273,137</mark> | 6,273,137 | | | |
| Encounter Diagnoses | 99,938 | 3,725,603 | 18,170,896 | 4.877 | 4 | 107 |
| Notes | 99,868 | 3,491,659 | 10,111,930 | | | |
| Hospital Encounters | 73,303 | 466,252 | 466,252 | | | |
| Lab Results | 83,435 | 733,461 | 20,186,748 | 27.523 | 12 | 19488 |
| Microbiology Results | 27,515 | 65,373 | 296548 | 4.536 | 1 | 268 |
| Medications Ordered | 94,089 | 1,388,086 | 5,336,506 | 3.845 | 1 | 1551 |
| Procedures Ordered | 98,514 | 1,880,309 | 7,229,854 | 3.845 | 1 | 6681 |
| Problem List | 90,722 | - | 761,260 | 8.391 | 6 | 182 |
| Result Comments | 72,716 | 468,814 | 916,554 | 1.955 | 1 | 691 |
| Surgeries | 18,640 | 29,895 | 31,889 | 1.067 | 1 | 41 |
| Vitals | 99,098 | 1,362,431 | 6,647,115 | 4.879 | 2 | 6387 |



Judgments from Patient Relevance Assessment Interface (PRAI)

| ic Description: Women who had a pregnar | ncy during which they had a 3rd trimester ou | tpatient visit, didn't smok | e, and didn't have in | ntellectual d | isability, mood disorder, schizophrenia, autism, or ADHD. | | |
|---|--|-----------------------------|-----------------------|---------------|--|----------------|-------------------|
| Pool 2 / Topic 1 / | Basic Info | | | | | | |
| atient | ion | | | | | | |
| Encounters | Demographics | | | | | | |
| Ambulatory Encounters | Filter Results | | | | | | |
| Hospital Encounters | Judge | OHSU_MRN CU | RRENT_AGE_YRS | BIRTH_ | DATE GENDER PATIENT_ALIVE DEATH_DATE ADD | RESS_STATE ADD | RESS_COUNTY GEN |
| Encounter Diagnoses | it⊘ Pro IIQ Con | | | | OR | WAS | HINGTON N |
| Vitals | 1.0/1 | | | | | | |
| Lab Results | | | | | | | - |
| Result Comments | Problems | | | | | | |
| Microbiology Results | Filter Results | | | | | | |
| Administered Medications | Judge | DX_START_DATE | DX_END_DATE | DX_ICD | DX_NAME | PROBL | EM_LIST_DX_STATUS |
| Ordered Medications | iô Pro 🛛 🖗 Con | | 9999-12-31 | 314.00 | ATTENTION DEFICIT DISORDER WITHOUT MENTION OF HYPERACTIVITY | ACTIVE | |
| Notes | IO Pro IIQ Con | | 9999-12-31 | 250.01 | DIABETES MELLITUS TYPE I | ACTIVE | |
| Ordered Procedures | IO Pro IIQ Con | | 9999-12-31 | 250.01 | TYPE 1 DIABETES MELLITUS | ACTIVE | |
| Surgeries | itô Pro II⊋ Con | | 9999-12-31 | 251.2 | HYPOGLYCEMIA | ACTIVE | |
| | | | | | | | |
| | | | | | | | |
| | | | | | | | |



























Questions?

William Hersh, MD Professor and Chair Department of Medical Informatics & Clinical Epidemiology School of Medicine Oregon Health & Science University Portland, OR, USA http://www.ohsu.edu/informatics

Email: <u>hersh@ohsu.edu</u> Web: <u>http://www.billhersh.info</u> Blog: <u>http://informaticsprofessor.blogspot.com</u> Twitter: <u>@williamhersh</u>



