# OPPORTUNITIES AND CHALLENGES FOR CLINICAL RESEARCH WITH ELECTRONIC HEALTH RECORDS

### Benjamin A. Goldstein, PhD, MPH

ben.goldstein@duke.edu

Department of Biostatistics & Bioinformatics
School of Medicine
Duke University

July 30$^{th}$, 2017

# WHY WE WANT TO USE EHRS FOR CLINICAL RESEARCH

- Data readily available
- Often 100,000's of Patients
- Information collected over a variety of fields
- Can study just about any clinical outcome
- Representative Population

# WHAT WE CAN DO WITH ELECTRONIC HEALTH RECORDS

1. Risk Prediction
   - Near term prediction - *Risk of inhospital sepsis*
   - Long(er) term risk - *30 Day Revisit*
2. Population Health
   - Health Service Utilization - *Assessment of high utilizers*
   - Disease Epidemiology - *Experience of incident diabetes in Durham County*
3. Comparative Effectiveness Research (CER)
   - Retrospective Studies - *Assessment of community intervention for diabetics (SEDI)*
   - Prospective Studies - *Point of care randomization*
4. Association Analyses
   - Risk factors for disease - *Phenome Wide Association Studies*
   - Data mining - *Drug-Drug interactions*

# WHY WE MAY *Not* WANT TO USE EHRS FOR CLINICAL RESEARCH

## DATA ARE NOT COLLECTED FOR RESEARCH

- Data exist in disparate places
- All patients have different pieces of information
- Observational Data

1 STRUCTURE OF ELECTRONIC HEALTH RECORDS

2 RESEARCH WITH EHR DATA

3 CONCLUDING THOUGHTS

1 STRUCTURE OF ELECTRONIC HEALTH RECORDS

2 RESEARCH WITH EHR DATA

3 CONCLUDING THOUGHTS

# THE EHR FRONT END: GETTING DATA IN

# DATA MOVE FRONT END TO DATA WHAREHOUSE



Epic's
CLARITY
Data Warehouse

•Patient Demographics
•Encounters
   (Outpatient/Inpatient)
•Diagnoses
•Procedures
•Lab Results
•Medications
•Vital Signs
•Social History
•Radiological Results
•Clinician Notes
Etc.

# THE DATA STRUCTURE

It's
Complicated...

# CHECK THE BLIND SPOTS



- Data movement and curation requires decision-making.
- Decisions may not be easily accessible.
- Decisions may not be documented or documentation may not be made available.

# TURNING EHRS INTO DATA

The analysis pipeline and data platform

**Datamart**

**Core Tables**
"Building block data" close to native source data format
(often transaction level)

**Data Dictionary**
Definitions of the source data, and mappings between source and target tables

**Curation Dictionary**
Datamart-specific processing rules, logic, and algorithms used to create the derived data

**Derived Tables**

| | |
|---|---|
| Aggregations /Summary Levels (eg, summary per year) | Consistency Enforced (eg, excluding adult height of 6 inches) |
| Processing Rules (eg, patient matching and linkage between sources) | Derived Variables (eg, computable phenotypes) |

**Source Data**

Examples:
- Enterprise Data Warehouse (EDW)
- External EHR sources
- Electronic data capture systems such as REDCap and eCOS
- Auxiliary data sources such as Census data

**Analysis and Evaluation**

**Analytic Dataset Collections**
Extracted in output format compatible with statistical purposes, such as SAS. Each collection is structured for the specific analysis and its independent/dependent variables. May include limited or anonymized datasets.

**Operational Reporting**
Examples: Dashboards and other Business Intelligence (BI) platforms; includes data quality reporting

Figure 2. Datamart components and relationship to external systems and processes.

# DATA MARTS:
## STRENGTHS AND WEAKNESSES

- Strengths
  - Registry like
  - Multiple clinical subject areas for cohort
  - Regularly scheduled data refresh
    **Ideal For:** Posing variety of questions across subject area

- Soft Spots
  - More time and effort to create than data extract
  - Structure not easily adaptable
  - Data are fixed between refreshes
    **Not Ideal For:** Small, targeted analyses

# ADDING INFORMATION BACK INTO EHR

- Dashboards
- Best Practice Alerts
- Predictive Analytics
- Clinical trial recruitment (Snifters)

# DIFFERENT TYPES OF CLINIC ENVIRONMENTS

- Clinic Based System (e.g. Practice Fusion, Flatiron)
  - Capture Routine Care
  - Local Population
  - Misses inpatient activity
- Hospital Based System
  - Observe inpatient procedures and events
  - Only observe when sick
  - Referral hospitals may not represent local or stable population
- Comprehensive Medical System (e.g. VA, Kaiser)
  - Observe all types of patient encounters
  - May represent artificial population

Collaborative Clinical Research

# FOUR WAYS EHR DATA DIFFER FROM TRADITIONAL CLINICAL DATA

1. We don't have everything we want
2. Outcomes are not defined - need to phenotype data
3. Data irregularly and potentially densely observed
4. Data not observed randomly - Informed Presence

# MOST EHRS ARE INCOMPLETE

- Patients seek care at multiple facilities
- Missing information on when individuals are healthy
- EHRs don't always contain all the data you want

# LINKING EHR DATA

- Data from other facilities (PCORNet)
- Claims: Center for Medicaire & Medicaid Services (CMS)
- Mortality: National Death Index (NDI) & Social Security Death Index (SSDI)
- Genetic Data
- GeoCode Information: American Community Survey (ACS)
- Personal Tracking Data: FitBit, sensors

# ISSUES OF DATA DEFINITION: WHAT IS A DIABETIC?



Research and applications

## A comparison of phenotype definitions for diabetes mellitus
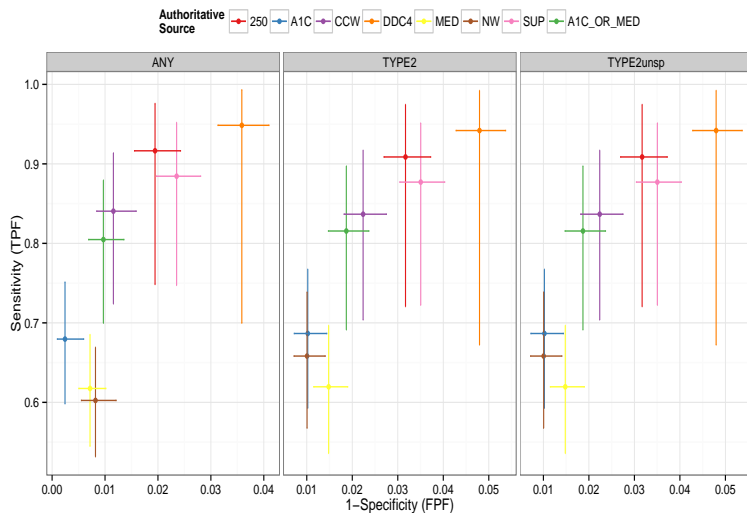
Rachel L Richesson,[1] Shelley A Rusincovitch,[2] Douglas Wixted,[3] Bryan C Batch,[4] Mark N Feinglos,[4] Marie Lynn Miranda,[5] W Ed Hammond,[2,6] Robert M Califf,[3,7] Susan E Spratt[4]

**ABSTRACT**
**Objective** This study compares the yield and characteristics of diabetes cohorts identified using heterogeneous phenotype definitions.
**Materials and methods** Inclusion criteria from seven diabetes phenotype definitions were translated into query algorithms and applied to a population (n=173 503) of adult patients from Duke University Health System. The numbers of patients meeting criteria for each definition and component (diagnosis, diabetes-associated medications, and laboratory results) were compared.
**Results** Three phenotype definitions based heavily on ICD-9-CM codes identified 9–11% of the patient population. A broad definition for the Durham Diabetes Coalition included additional criteria, and identified 13%. The electronic medical records and genomics, NYC A1c Registry, and diabetes-associated medications definitions, which have restricted or no ICD-9-CM criteria, identified the smallest proportions of patients (7%). The demographic characteristics for all seven phenotype definitions were similar (56–57% women, mean age ...)



**Figure 1** Overlap of diabetes cohorts identified from different categories of phenotype eligibility criteria; n=24 520 patients identified by criteria from any of the three categories.

Richesson RL, Rusincovitch SA, Wixted D, Batch BC, Feinglos MN, Miranda ML, Hammond WE, Califf RM, Spratt SE. A Comparison of Phenotype Definitions for Diabetes Mellitus. J Am Med Inf Assoc 2013 (epub ahead of print). http://www.ncbi.nlm.nih.gov/pubmed/24026307

# ISSUES OF DATA DEFINITION: WHAT IS A DIABETIC?

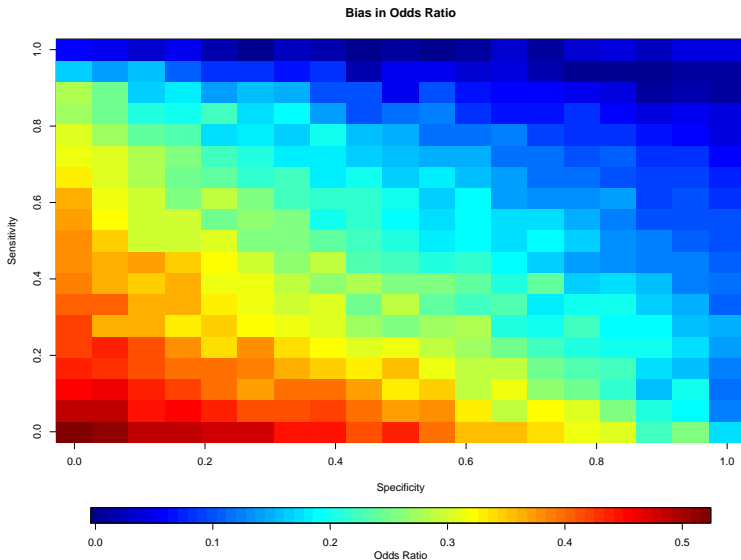| | ICD-9 250.xx | ICD-9 250.x0 & 250.x2 (exclude type I) | Expand. ICD-9 (249.xx, 357.2, 362.0x, 366.41) | HbA1c | Glucose | Abnormal OGTT | Diabetes Meds |
|---|---|---|---|---|---|---|---|
| ICD-9 250.xx | X | | | | | | |
| CMS CCW | X* | | X* | | | | |
| NYC A1c Registry | | | | X | | | |
| Meds | | | | | | | X |
| DDC | | X | X | X | X | X | X |
| SUPREME-DM | X* | | X* | X | X | X | X |
| eMERGE | | X* | | X | X | | X |

\* Distinction between Inpatient and Outpatient Visits

# DEFINITION DIFFERENCES



Diabetes Validation Results faceted by Endpoint

# IMPACT OF POORER DEFINITIONS



Bias in Odds Ratio

## ADDITIONAL PHENOTYPING CHALLENGES

- **Death:** Internal work estimates 20% capture of deaths
- **Disease Incidence:** Need to apply 'burn-in' periods
- **Censoring:** Need to apply 'burn-out' periods
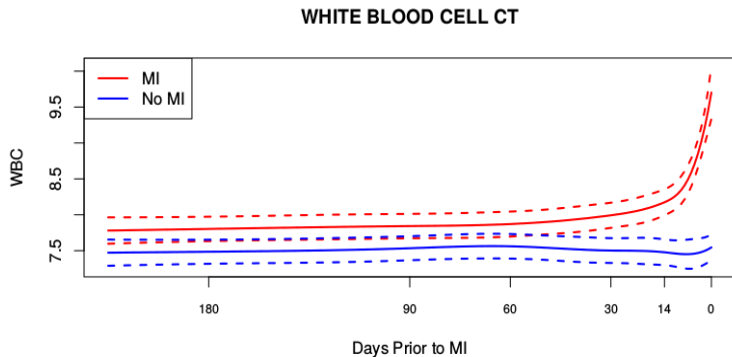
# MULTIPLE MEASUREMENTS PER PERSON

### OPPORTUNITIES

- Get to observe patient's evolving health status
- More frequent visits than a typical longitudinal study
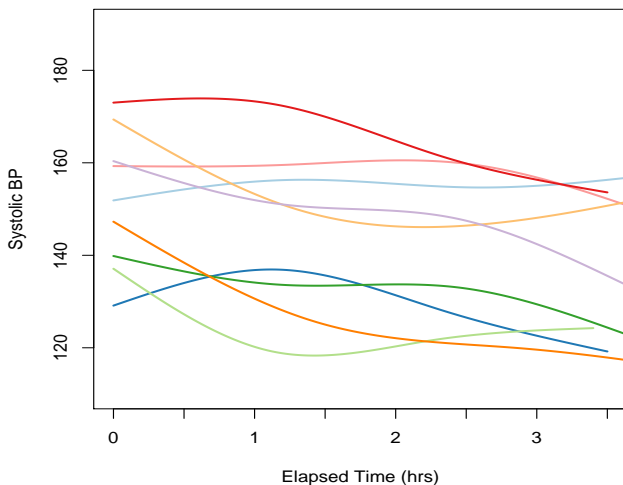- Denser visit information

### CHALLENGES

- Visits are irregularly spaced
- Different ways to aggregate
- Don't know what you are not seeing

# LOOK AT CHANGES OVER LONG PERIODS OF TIME...



**WHITE BLOOD CELL CT**

# ...OR SHORT PERIODS OF TIME



**Individual Blood Pressure Curves**

## ANALYZING REPEATED MEASURES

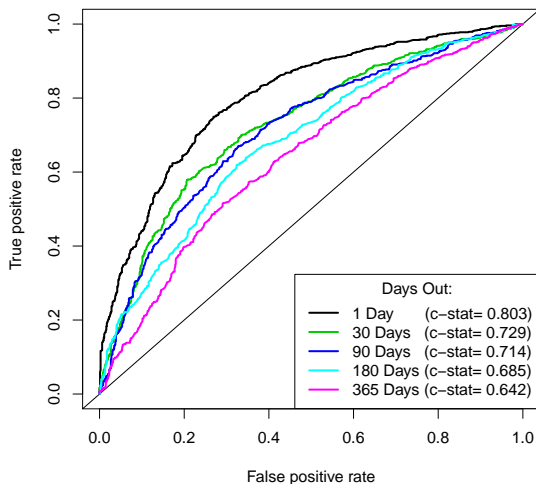| **Summarizing Data** | **Modelling Progression** |
| --- | --- |
| Mean/Median Values | Regression Splines |
| Extreme Values | Functional Data Analysis |
| Variability | Joint Models |
| Number of Measurements | |

# SIMPLER METHODS OFTEN WORK BEST

# EHR DATA OPTIMIZED FOR NEARER TERM PREDICTION



**ROC Curves for Forecasting SCD**

Days Out:
— 1 Day     (c–stat= 0.803)
— 30 Days   (c–stat= 0.729)
— 90 Days   (c–stat= 0.714)
— 180 Days  (c–stat= 0.685)
— 365 Days  (c–stat= 0.642)

True positive rate

False positive rate

# TOP PREDICTORS

|   | 1 Day | 7 Days | 30 Days |
|---|---|---|---|
| 1 | LabValue: Albumin | LabValue: Albumin | LabValue: Albumin |
| 2 | Pre Systolic BP | Pre Systolic BP | Pre Systolic BP |
| 3 | Pre MAP | Pre MAP | Lowest Systolic BP |
| 4 | Pre Pulse Pressure | LabValue: WBC | LabValue: Creatinine |
| 5 | LabValue: Hemoglobin | Medication Dose: Epogen | Pre MAP |
| 6 | Lowest Systolic BP | LabValue: Creatinine | Post MAP |

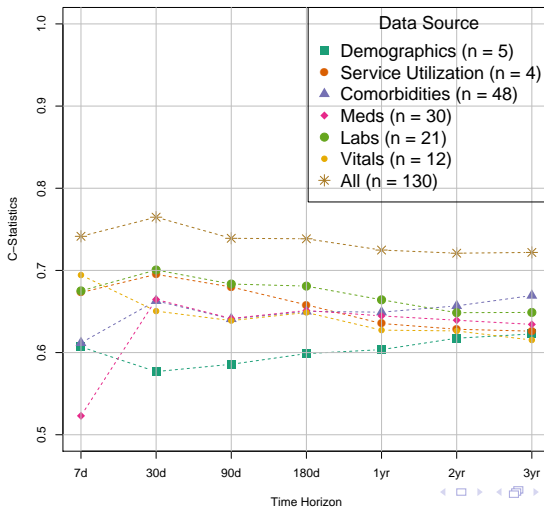|   | 90 Days | 180 Days | 365 Days |
|---|---|---|---|
| 1 | LabValue: Albumin | LabValue: Albumin | LabValue: Albumin |
| 2 | Pre Weight | Pre Weight | Medication Dose: Epogen |
| 3 | Pre Systolic BP | Pre Map | Age |
| 4 | Pre Pulse Pressure | Post Weight | LabValue: Creatinine |
| 5 | Medication Dose: Epogen | Medication Dose: Epogen | Pre Systolic BP |
| 6 | Post Weight | Pre Systolic BP | Pre Pulse Pressure |

# TOP PREDICTORS

| | 1 Day | 7 Days | 30 Days |
|---|---|---|---|
| 1 | LabValue: Albumin | LabValue: Albumin | LabValue: Albumin |
| 2 | Pre Systolic BP | Pre Systolic BP | Pre Systolic BP |
| 3 | Pre MAP | Pre MAP | Lowest Systolic BP |
| 4 | Pre Pulse Pressure | LabValue: WBC | LabValue: Creatinine |
| 5 | LabValue: Hemoglobin | Medication Dose: Epogen | Pre MAP |
| 6 | Lowest Systolic BP | LabValue: Creatinine | Post MAP |

| | 90 Days | 180 Days | 365 Days |
|---|---|---|---|
| 1 | LabValue: Albumin | LabValue: Albumin | LabValue: Albumin |
| 2 | Pre Weight | Pre Weight | Medication Dose: Epogen |
| 3 | Pre Systolic BP | Pre Map | Age |
| 4 | Pre Pulse Pressure | Post Weight | LabValue: Creatinine |
| 5 | Medication Dose: Epogen | Medication Dose: Epogen | Pre Systolic BP |
| 6 | Post Weight | Pre Systolic BP | Pre Pulse Pressure |

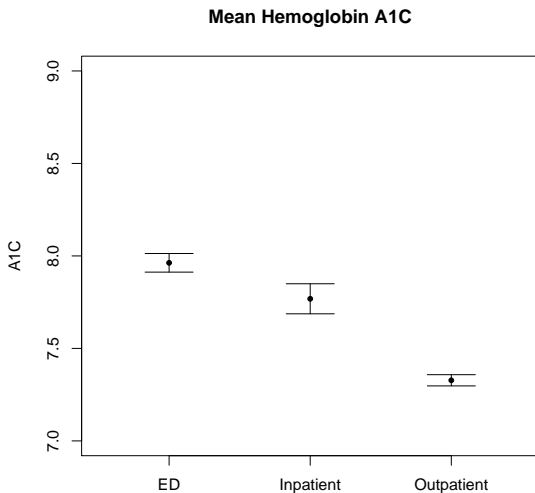# DIFFERENT DATA ELEMENTS HAVE DIFFERENT PREDICTABILITY



Predicting Death at Different Horizons
With Different Data Sources

# BIASES IN EHRS:
# INFORMED PRESENCE

- We only see patients when they are sick
- We only see information that is deemed important
- Different environments have different policies

# INFORMED PRESENCE I:
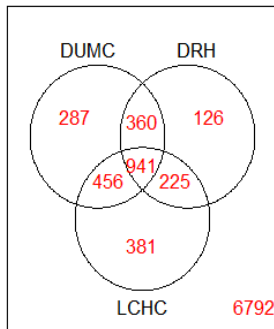## WHERE A PERSON SEEKS CARE IS INFORMATIVE

**Mean Hemoglobin A1C**

## LOCATION IMPACTS INFERENCE

- Hazard Ratio for HgB A1C for time to Myocardial Infarction

| Type | Hazard Ratio | P-value |
|---|---|---|
| Unadjusted | 1.06 (1.01, 1.11) | 0.026 |
| Adjusted for Location | 0.97 (0.92, 1.02) | 0.178 |
| OP Only | 1.07 (1.00, 1.14) | 0.044 |
| ED Only | 0.94 (0.89, 0.99) | 0.022 |

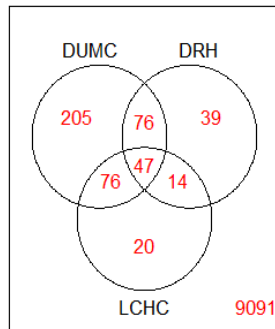- Interaction between A1C and location

# INFORMED PRESENCE II:
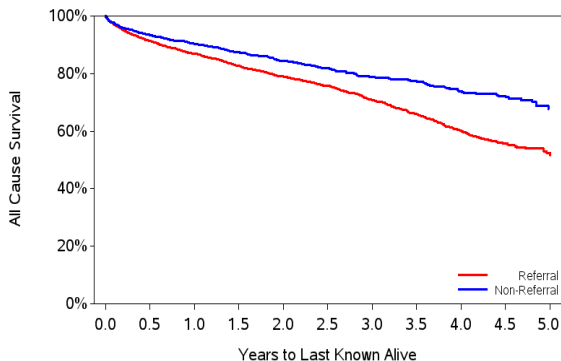## WHICH HOSPITAL A PATIENT USES IS INFORMATIVE



**Diabetes**
**N=2,783**

DUMC    DRH

287   360   126

941

456   225

381

LCHC    6792

**Cancer**
**N=477**

DUMC    DRH

205   76   39

47

76   14

20

LCHC    9091

## FACILITY IMPACTS INFERENCE

- Odds Ratio for Cancer Status on Diabetes

| Location | Odds Ratio | 95% CI |
|---|---|---|
| All Facilities | 1.69 | (1.36, 2.10) |
| DUMC Only | 1.46 | (1.15, 1.87) |
| DRH Only | 0.89 | (0.63, 1.26) |
| LCHC Only | 1.08 | (0.74, 1.56) |

# INFORMED PRESENCE III:
## REFERAL HOSPITALS ARE AN *Admixed* POPULATION



| Number at risk | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Referral | 5522 | 3307 | 2690 | 2159 | 1748 | 1360 | 995 | 697 | 474 | 282 | 65 |
| Non-Referral | 2114 | 1532 | 1318 | 1110 | 882 | 697 | 519 | 387 | 262 | 171 | 64 |

## ADMIXTURE BIAS

- Comparison of Local and Referal Patients at Cardiac Catheterization Lab

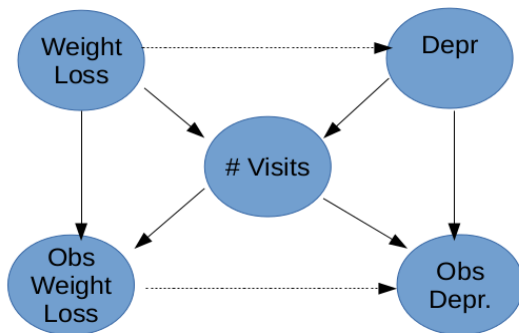| **Local Patients** | **Referal Patients** |
| --- | --- |
| Older | Younger |
| More Comorbidities | More severe valve disease |
| Disease due to ageing | Disease due systematic factors |
| Better outcomes | More follow-up procedures |

# INFORMED PRESENCE IV:
## NEED TO ACCOUNT FOR NUMBER OF ENCOUNTERS

Regression of Depression on Weight Loss

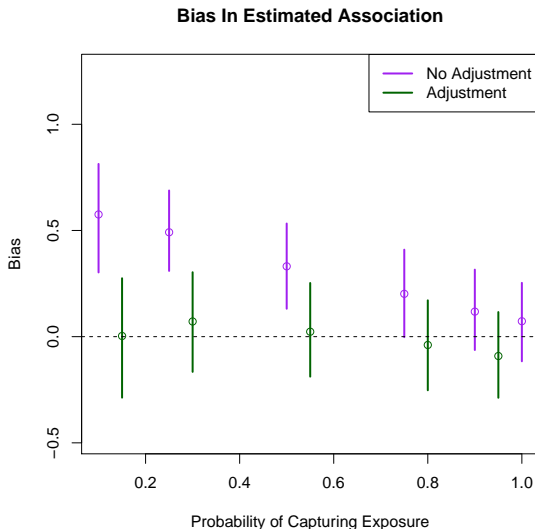|  | Odds Ratio | Δ log(OR) | Δ OR |
|---|---|---|---|
| Minimally Adjusted | 3.98 (3.81, 4.17) | — | — |
| + No. Encounters | 2.37 (2.26, 2.50) | -0.52 | -1.61 |
| + Comorbidities | 2.82 (2.69, 2.96) | -0.35 | -1.16 |
| + No. Encounters & Comorb | 2.30 (2.18, 2.42) | -0.55 | -1.68 |

# NUMBER OF ENCOUNTERS POTENTIAL CONFOUNDER

## NEED TO ACCOUNT FOR NUMBER OF ENCOUNTERS

|             |             | Median Number of Encounters | |
|-------------|-------------|-------------------|----------------|
|             | Sensitivity | Without Condition | With Condition |
| Depression  | 56.3%       | 6                 | 38             |
| Weight Loss | 9.3%        | 7                 | 45             |

# NUMBER OF ENCOUNTERS POTENTIAL CONFOUNDER



**Bias In Estimated Association**

1 STRUCTURE OF ELECTRONIC HEALTH RECORDS

2 RESEARCH WITH EHR DATA

3 CONCLUDING THOUGHTS

# EXTRA CARE NEEDED

- Need to be mindful from where the data come
- There is not always one way to turn raw data into analytic data
- Which data to *cut* is more important than how you analyze it
- New analytic techniques may be useful/necessary

# QUESTIONS TO ASK WHEN DESIGNING EHR BASED STUDIES

- Where in the health system are the data collected?
- What is the coverage/catchment area of your health system?
- Is the patient population receiving care across multiple institutions/centers?
- Do the data constitute different catchments? (Admixture)
- How are you defining exposures and outcomes? (Phenotyping)
- How are you defining person-time?
    - What is an appropriate burn-in period to define a cohort?
    - Is a burn-out period necessary to define censoring?
- Do different populations produce more information (i.e. sicker patients have more encounters)?

# ADDITIONAL FRONTIERS

- Micro-randomized trials
- Integration of external data
- Real time risk assessment

# IS IT ALL BAD?

## A LOT OF OPPORTUNITIES WITH EHRS

- More studies
- Cheaper studies
- Faster studies
- (Perhaps) More representative studies

# REFERENCES

- Phelan, M., Bhavsar, N.A., Goldstein, B.A. Illustrating Informed Presence Bias in Electronic Health Records Data: How Patient Interactions with a Health System Can Impact Inference. *eGEMS*, In Press.

- Goldstein B.A., Pomann, GM, Winkelmayer, WC., and Pencina, MJ. Comparison of risk prediction methods using repeated observations with application to Electronic Health Record. *Statistics in Medicine*, 2017, 36(17) 2750–2763.

- Goldstein B.A. Pencina M. J., Montez-Rath, M.E. and Winkelmayer W.C. Value of different categories of information available in electronic health records for prediction of mortality in patients on dialysis. *Journal of the American Medical Informatics Association*, 2017, 24(1): 176-181.

- Goldstein B.A., Navar, A.M., Pencina, M.J., and Ioannidis, J.P. Opportunities and challenges in developing risk prediction models with electronic health records data: A systematic review. *Journal of the American Medical Informatics Association*, 2017, 24(1): 198-208.

- Spratt, SE. Pereira, K., Granger, BB.,Batch, BC., Phelan, M., Pencina, M., Miranda, ML., Boulware, E., Lucas, JE., Nelson, CL., Neely, B., Goldstein, BA., Barth, P., Richesson, RL., Riley, IL., Corsino, L., McPeek Hinz, ER., Rusincovitch, S., Green, J., Barton, AB., and the DDC Phenotype Group., Assessing Electronic Health Record Phenotypes against Gold-Standard Diagnostic Criteria for Diabetes Mellitus. *Journal of the American Medical Informatics Association*, 2017, 24, e121-e128.

- Goldstein B.A., Bhavsar, N.A., Phelan, M., and Pencina, M.J. Controlling for informed presence bias due to the number of health encounters in an Electronic Health Record. *American Journal of Epidemiology*, 2016, 184(11): 847-855.

- Goldstein BA, Cheng TI, and Winkelmayer, WC. Classifying Individuals Based on a Densely Captured Sequence of Vital Signs: An Example using Repeated Blood Pressure Measurements during Hemodialysis Treatment. *Journal of Biomedical Informatics*, 2015, 57: 219-224.

- Goldstein, B.A., Assimes TL., Winkelmayer WC, and Hastie T. Detecting clinically meaningful biomarkers with repeated measurements in an Electronic Health Record. *Biometrics*, 2015, 71, 478-486.

- Goldstein, B.A., Chang, TI, Mitani, A.A., Assimes TL., and Winkelmayer WC. Near-term prediction of sudden cardiac death in older patients using the electronic health records. *Clinical Journal of the American Society of Nephrology*, 9, 2014, 82-91.

- Richesson RL, Rusincovitch SA, Wixted D, Batch BC, Feinglos MN, Miranda ML, Hammond WE, Califf RM, and Spratt SE.A comparison of phenotype definitions for diabetes mellitus. *Journal of the American Medical Informatics Association*, 2013, 20, e319-326.