



# Detecting Anomalous Patterns of Care Using Health Insurance Claims

**Sriram Somanchi**

**Mendoza College of Business**

**University of Notre Dame**

Joint work with Edward McFowland III\* and Daniel B. Neill+

\*University of Minnesota, +Carnegie Mellon University

---

# Agenda

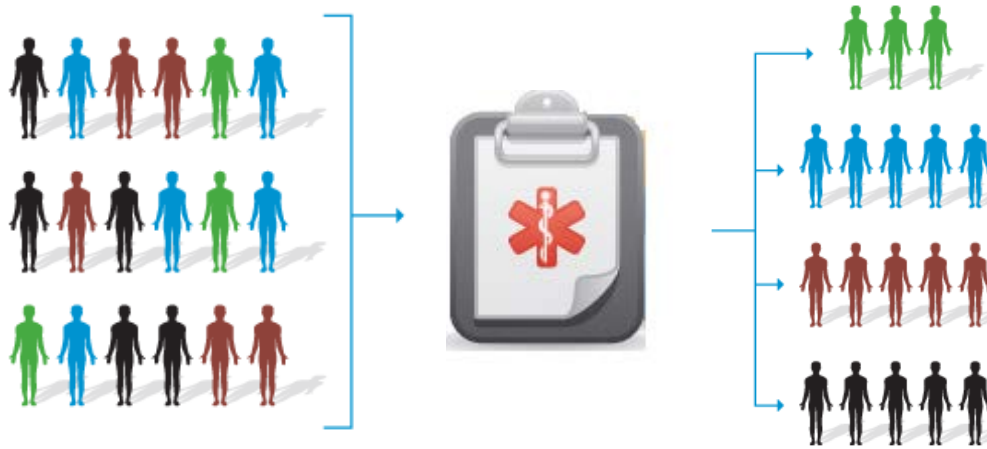
- Introduction
  - Research Question
  - Motivating Example
  - Literature and Contribution
- Methods
  - Problem Formulation
  - Algorithm
  - Modeling the scoring function
- Empirical Analysis
  - Data
  - Results
  - Validation using regression analysis

---

# Agenda

- Introduction
  - Research Question
  - Motivating Example
  - Literature and Contribution
- Methods
  - Problem Formulation
  - Algorithm
  - Modeling the scoring function
- Empirical Analysis
  - Data
  - Results
  - Validation using regression analysis

# Introduction: Healthcare Setting



- Challenges the US healthcare system faces<sup>1,2</sup>
  - Instances of over-treatment and under-treatment
  - Inconsistencies in execution of care

1. N.C. Lallemand, "Health Policy Brief: Reducing Waste in Health Care," Health Affairs, 13 Dec. 2012.  
2. L.T. Kohn, et al., To Err Is Human: Building a Safer Health System, Inst. of Medicine/Nat'l Academy Press, 1999.

---

# Introduction: Healthcare Setting

- Huge opportunity to discover novel patterns of care that are potentially effective due to availability of
  - Electronic Health Records
  - Documented care through health insurance claims
- Analyze patterns across patients and provide actionable insights

---

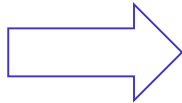
# Research Question

- Identify the treatment and the sub-population for whom that treatment corresponds to significantly better or worse outcomes
  - With multiple treatments and population characteristics varying in multiple dimensions.

# Motivating Example



Health Insurance  
Claims Data



Healthcare Analyst  
Patrick



Congestive Heart Failure Patients

1. Males
2. Age above 50
3. Similar co-morbidity (atrial fibrillation, on anticoagulant)

Taking Carvidilol correlated  
with longer stay in hospital

Can we automate the process and produce these interesting hypotheses?

---

# Literature and Contribution

- Heterogeneous Treatments Effects with a given treatment
    - Randomized Control Trials
      - Imai and Ratkovic (2013)
      - McFowland, Somanchi and Neill (2017)
    - Observational Studies
      - Athey and Imbens (2016)
      - Wager and Athey (2016 arXiv)
  
  - My Contribution
    - Identify sub-populations and treatment, with multiple treatments, who have anomalous outcomes
    - Computationally efficient algorithm instead of evaluating exponentially many sub-populations
    - Observational studies

Effectively use observational data to help run future targeted control trials
-



---

# Agenda

- Introduction
  - Research Question
  - Motivating Example
  - Literature and Contribution
- Methods
  - Problem Formulation
  - Algorithm
  - Modeling the scoring function
- Empirical Analysis
  - Data
  - Results
  - Validation using regression analysis

---

# Problem Formulation

- Let  $X = (X_1, X_2, \dots, X_N)$  be the set of observed covariates for a patient
- Let  $T_1, T_2, \dots, T_M$  be the set of available treatments
- Let  $Y$  be the scalar outcome of interest

# Estimating Potential Outcome Distributions

- We want to estimate the distribution of potential outcomes for treatment assignments  $T_j = 1$ , for a given sub-population,  $S$

$$f_{j1,S} = f(y^{(1)} \mid x \in S)$$

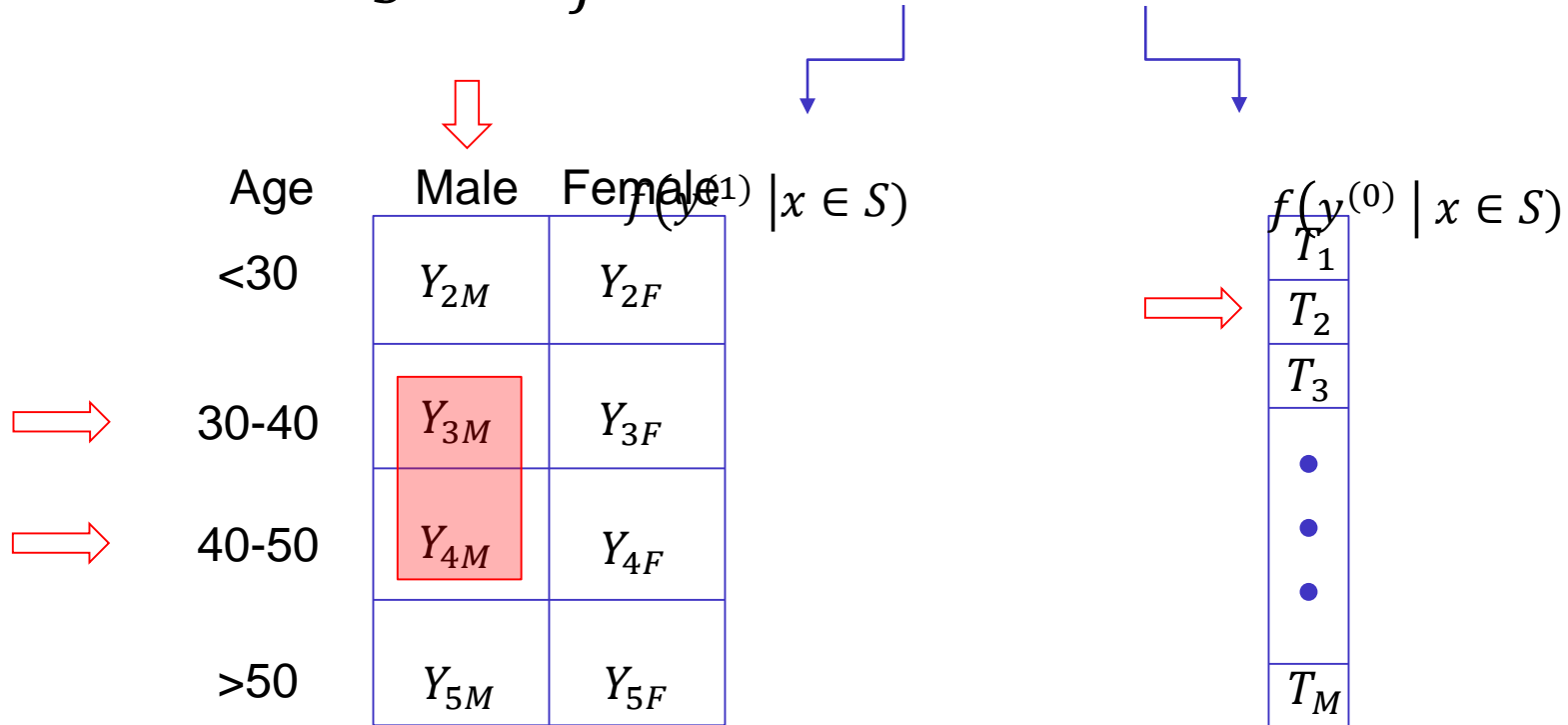
- Similarly, we want to estimate

$$f_{j0,S} = f(y^{(0)} \mid x \in S)$$

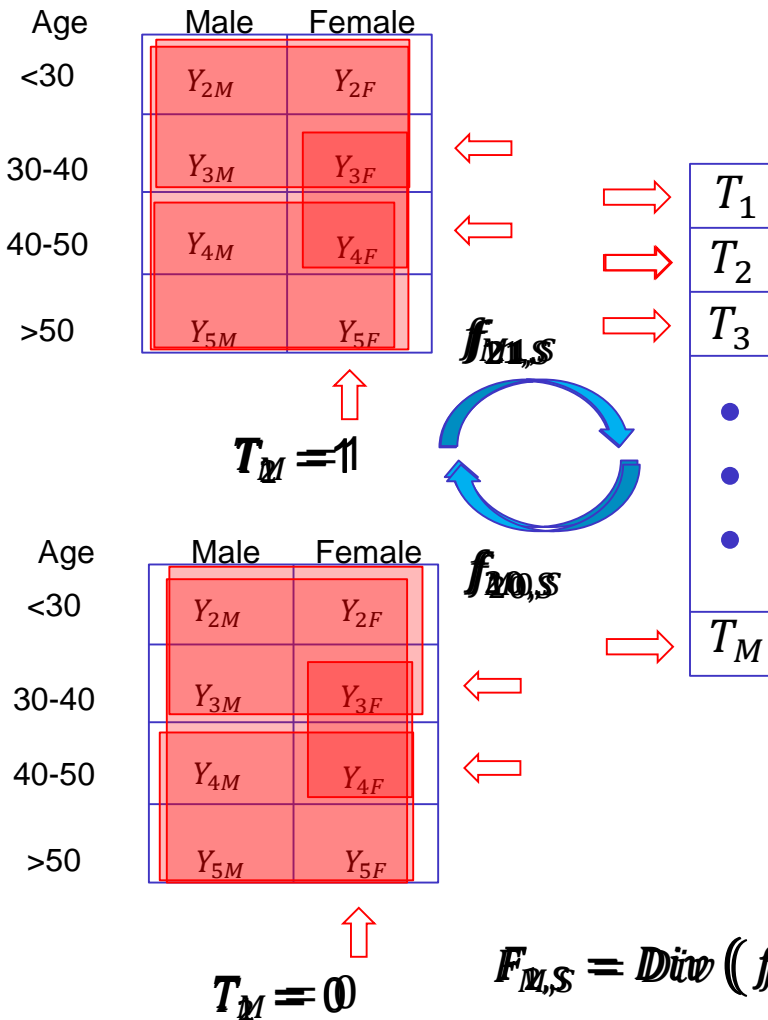
# Our Goal

- Simultaneously detect effective treatment and sub-population combination

$$\max_S \max_j \text{Div}(f_{j1,S}, f_{j0,S})$$



# Anomalous Patterns of Care Scan



1. Start with a random sub-population  $S$
2. For each  $T_j$ 
  - a. Compute the propensity scores
  - b. Reweight outcome distributions
  - c. Compute Divergence  $F_{j,S}$
3.  $j^* = \text{argmax}_j F_{j,S}$
4. Reweight entire population outcomes based on  $T_{j^*}$
5. Use MD-Scan to identify  $S^* = \text{argmax}_S F_{j^*,S}$
6. Set  $S = S^*$  and repeat steps 2 to 5 until score stops increasing
7. Repeat steps 1-6 for  $R$  times
8. Compute statistical significance by randomization testing

Iterative Ascent algorithm between sub-populations and treatments

# Inverse Propensity Score Weighting

- We use inverse propensity score weighting to estimate the outcome distribution from observational data

$$f_{j1,S} = f(y^{(1)} \mid x \in S) \approx \sum_{x \in S} \frac{f(y, T_j=1, X=x)}{P(T_j=1 \mid X=x)}$$

$$f_{j0,S} = f(y^{(0)} \mid x \in S) \approx \sum_{x \in S} \frac{f(y, T_j=0, X=x)}{P(T_j=0 \mid X=x)}$$

# Efficiently Optimizing for Divergence

- Parametric form
  - Compute the sufficient statistic
  - Expectation-based Subset Scan framework
- Non-parametric form
  - Compute p-values for outcomes
  - Non-parametric Subset Scan framework
- In order to efficiently optimize, the divergence score needs to satisfy **Linear Time Subset Scanning (LTSS)** property

# Multi-Dimensional Scan (MD-Scan)

$$S^* = \operatorname{argmax}_S F_{j,S}$$

2<sup>4</sup> combinations of age groups

↓

Instead we just need to evaluate 4 combinations

		↓	↓
	Age	Male	Female
→	<30	$Y_{2M}$	$Y_{2F}$
→	30-40	$Y_{3M}$	$Y_{3F}$
→	40-50	$Y_{4M}$	$Y_{4F}$
→	>50	$Y_{5M}$	$Y_{5F}$

Each step is computationally efficient if divergence function satisfies LTSS property



# Modeling the Scoring Function

- We model the scoring function as generalized log-likelihood ratio statistic
- We assume a parametric distribution for the outcome and compute the sufficient statistics of the expected distribution from the control ( $T_j = 0$ )
  - Expectation Based Poisson
  - Expectation Based Gaussian
  - Exponential family distributions

# Expectation Based Poisson statistic for potential outcomes


$$\mathbf{H}_0 \quad : \quad Y_i^{(1)} \mid X_i \in X_S \sim \text{Poisson}(\lambda_S) \quad \forall X_S$$

$$\lambda_S = E[Y^{(0)} \mid X \in X_S]$$

$$\mathbf{H}_1(\mathbf{S}, q) : \quad Y_i^{(1)} \mid X_i \in X_S \sim \text{Poisson}(q * \lambda_S) \quad X_S \in \mathbf{S}$$

$$\mathbf{H}_1(\mathbf{S}, q) : \quad Y_i^{(1)} \mid X_i \in X_S \sim \text{Poisson}(\lambda_S) \quad X_S \notin \mathbf{S}$$

$$F(S|q) = \log \frac{P(\text{Data} \mid H_1(S, q))}{P(\text{Data} \mid H_0)}$$

$$F(S) = \max_q F(S|q)$$

$$S^* = \max_S F(S)$$

# Traditional Causal Estimands

- Average Treatment Effect

$$\tau_{ATE} = E[Y(1) - Y(0)]$$

- Conditional Average Treatment Effect

$$\tau_{CATE}(x) = E[Y(1) - Y(0) | X = x]$$

- Marginal Conditional Average Treatment Effect (Grimmer, Messing, Westwood 2017)

$$\tau_{MCATE}(x^s) = \int E \left[ Y(1) - Y(0) \mid X^1, X^2, \dots, X^s = x^s, \dots, X^d \right] dF_{X^{-s} | X^s = x^s}$$

# Our General Causal Estimand

- Distributional Average Treatment Effect

$$\tau_{DATE}(S) = E_{x \in S} \left[ \text{Div} \left( F_{Y(1)|X=x}, F_{Y(0)|X=x} \right) \right]$$

- Parametric Distributional Average Treatment Effect

$$\tau_{PDATE}(S) = \max_q \tau_{PATE}^q(S) = \max_q F(S | q)$$

# Scoring function with observational data

- **Theorem 1:** If we assume unconfoundedness and we have balancing propensity score ( $e(X)$ ) weights, then  $F_{obs}(S|q)$  using weighted observed outcomes ( $Y$ ) is unbiased estimator of  $F(S|q)$  using potential outcomes ( $Y^{(0)}, Y^{(1)}$ )

# Statistical Properties

- Subpopulation Exactness
  - If the signal is  $\alpha$  – strong (**Lemma 1**)
    - $S^* \subseteq S^T$
  - If the signal is  $\beta$  – *homogeneous* (**Lemma 2**)
    - $S^* \supseteq S^T$
  - We show that  $\alpha$  and  $\beta$  are bounded

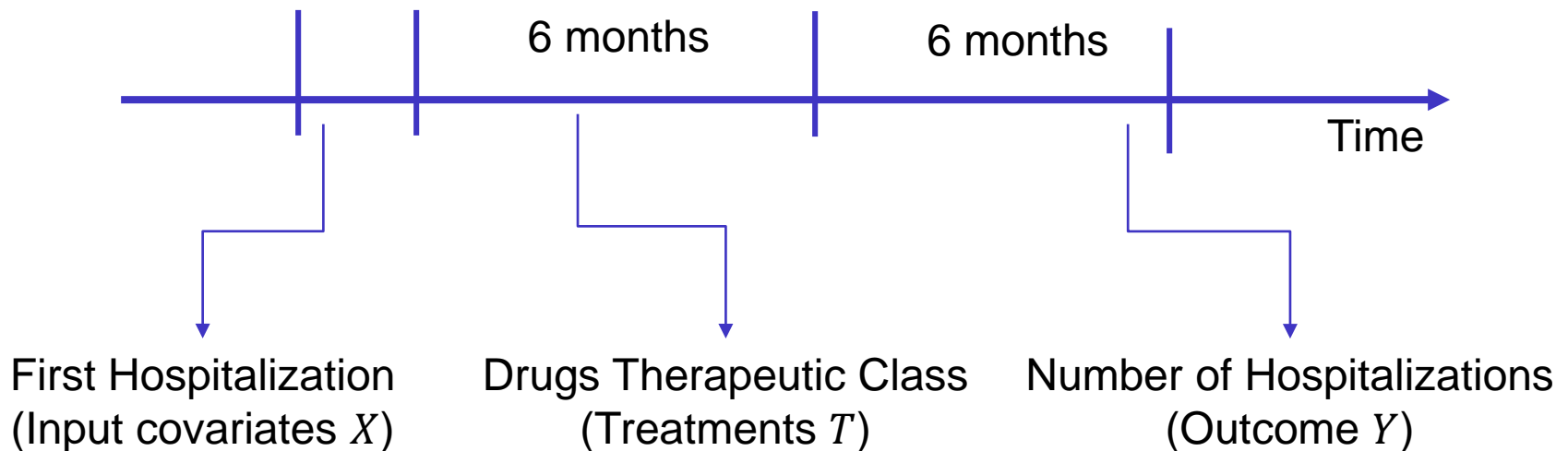
---

# Agenda

- Introduction
  - Research Question
  - Motivating Example
  - Literature and Contribution
- Methods
  - Problem Formulation
  - Algorithm
  - Modeling the scoring function
- Empirical Analysis
  - Data
  - Results
  - Validation using regression analysis

# Highmark Claims Data

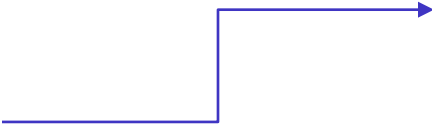
- Patients with primary or admission diagnosis as 'diseases of the circulatory system' from the year 2008 to 2014
  - ~125K patients





# Highmark Claims Data

- Covariates ( $X$ ) were built based on
  - Demographics
  - Median income at patient's zip code level
  - Diagnosis (primary and secondary)
  - Charlson Comorbidity Index<sup>1</sup>
  - Length of current stay
  - Previous outpatient visits
- Treatments ( $T_j$ )
  - Drug Therapeutic Class
- Outcome ( $Y$ )
  - Number of hospitalizations, Total length of stay



Bronchial Dilators  
Glucocorticoids  
Thyroid Preparations  
Diabetic Therapy  
Lipotropics  
Hypotensives  
Vasodilators  
Digitalis Preparations  
Cardiovascular  
Preparations  
Anticoagulants  
Diuretics

1. Quan et al (2005). Coding algorithms for defining comorbidities in icd-9-cm and icd-10 administrative data. Medical Care, 43(11):1130–1139

# Descriptive Statistics

Characteristics	Values	Percentage of Patients
Entire Population		100% (124,146)
Gender	Male	53.0%
	Female	47.0%
Age	Below40	2.8%
	40to60	19.8%
	60to80	43.5%
	Above80	33.9%
Hypertensive	Yes	53.9%
	No	46.1%
Diabetic	Yes	29.2%
	No	70.8%
Obese	Yes	11.1%
	No	88.9%
Primary Diagnosis	Rheumatic (390-398)	0.5%
	Hypertensive (401-405)	3.5%
	Ischemic (410-414)	24.5%
	Pulmonary (415-417)	3.7%
	Heart Failure (420-429)	33.0%
	Cerebrovascular (430-438)	16.6%
	Arteries (440-448)	5.0%
Veins and lymphatics (451-459)	13.2%	

---

# Results

- We ran our methodology on this dataset to identify patterns of interest
- We have ranked order of the highest scoring combination of subpopulation and treatments
- We discuss the details of the highest scoring subpopulation and treatment pair

# Highest Scoring Subpopulation-Treatment Combination

## ■ Subpopulation Characteristics Identified

### □ Gender

- Male

### □ Medical condition

- Hypertension
- Obese or Overweight

### □ Age

- 40 to 80

### □ Primary diagnosis

- Ischemic Heart disease (ICD9 410 – 414)
- Heart Failure (ICD9 420 – 429)
- Cerebrovascular heart disease (ICD9 430 – 439)

### □ Secondary diagnosis

- No respiratory (ICD9 460 – 519)
- Endocrine and Immunity disorders (ICD9 240 – 279)

## ■ Drug therapeutic class

### □ Glucocorticoids

## ■ Outcome

### □ More number of hospitalizations

	Glucocorticoids	
	Yes	No
Number of Patients	264	1713
Mean Number of Hospitalizations	0.606 (0.069)	0.280 (0.016)

---

# Validation of our results

- There is huge literature in the medical community on Glucocorticoids and Cardiovascular issues:
  - Association using 10 years of observational data (Heart, 2004)
  - Metabolic and tissue level effects in heart (European Journal of Endocrinology, 2007)
  - Experiments at micro level analysis of glucocorticoids signaling certain receptors in heart for mice (J of Biochemical and Molecular Biology, 2015)

# Understanding the results using regression analysis

- In order to understand the results we split the data into
  - 60% for running our APC Scan
  - 40% for running the regression analysis
- Regression with outcome  $Y$  as number of hospitalizations with Glucocorticoids as one of independent variable  $X$ , for
  - The entire population
  - The entire population with a dummy for subpopulation identified by APC Scan
  - The subpopulation identified by APC Scan
  - The complementary subpopulation

# Regression analysis (Poisson) on a Hold-Out set

	Number of Hospitalizations		Number of Hospitalizations	
	(1)	(2)	(3)	(4)
Glucocorticoids	0.101*** (0.007)	0.099*** (0.007)	0.410*** (0.089)	0.099*** (0.007)
Glucocorticoids* Subpopulation		0.265*** (0.088)		
Subpopulation		-0.313*** (0.068)		
Age	0.079*** (0.004)	0.079*** (0.004)	-0.040 (0.079)	0.080*** (0.004)
Females	0.116*** (0.008)	0.113*** (0.008)		0.113*** (0.008)
Hypertensive	-0.163*** (0.008)	-0.161*** (0.008)		-0.161*** (0.008)
Diabetic	0.286*** (0.008)	0.286*** (0.008)	0.193*** (0.089)	0.287*** (0.008)
Obesity	0.007 (0.013)	0.020 (0.013)		0.020 (0.013)
...	...	...	...	...
Constant	-0.773*** (0.044)	-0.772*** (0.044)	-1.634*** (0.120)	-0.772*** (0.044)
Observations	49,658	49,658	796	48,862

10.6%

50.6%

(1) Entire Population

(2) Entire Population with dummy for the subpopulation

(3) Only with the subpopulation identified by APC-Scan

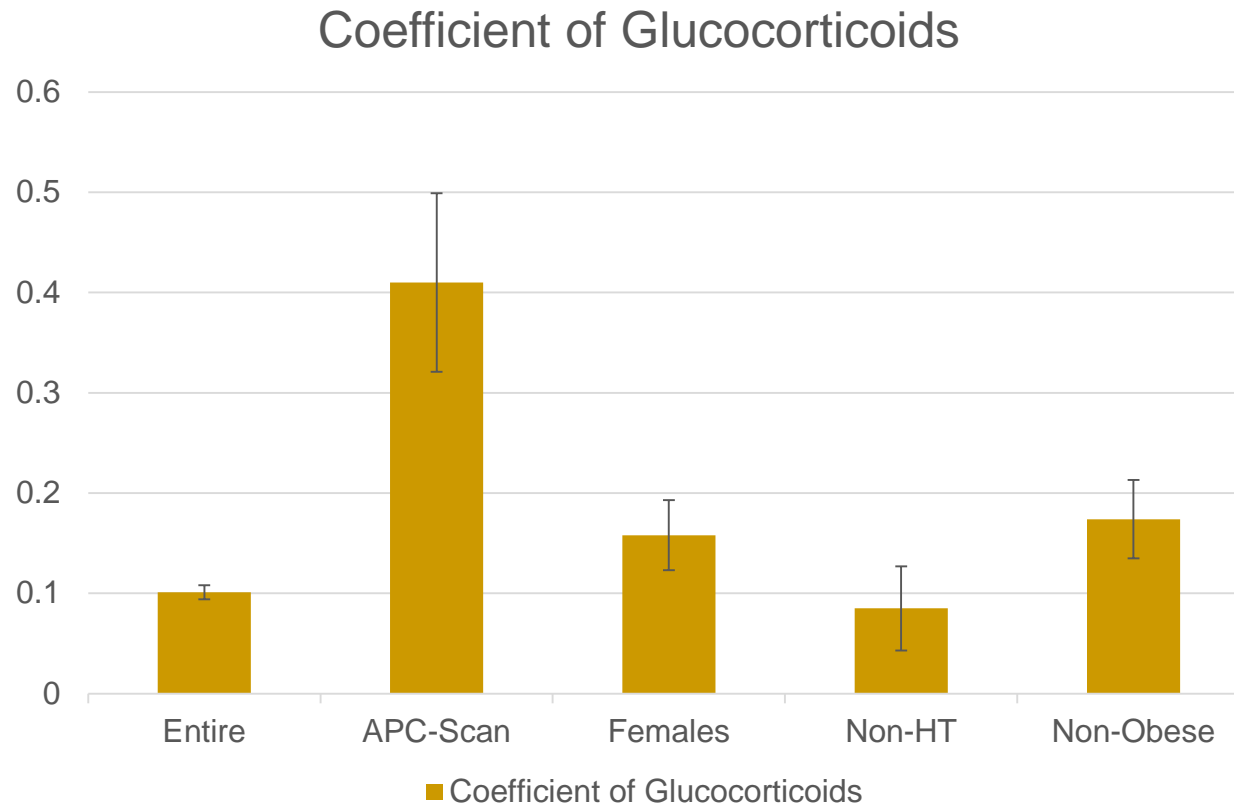
(4) Only with the complementary subpopulation

We have included all input characteristics  $X$  for our regression

Note: \* $p < 0.1$ ; \*\* $p < 0.05$ ; \*\*\* $p < 0.01$

# Sensitivity analysis

- Subpopulation identified was slightly modified





---

# Sensitivity analysis

- Typical diseases treated using Glucocorticoids
  - Rheumatic Arthritis
  - Chronic Obstructive Pulmonary Disease
  - Cushing's syndrome
- Alternative drugs to Glucocorticoids
- Ruled out hospital level biases in treating with Glucocorticoids
  - Overlap coefficient between two groups is 0.78

# Summary of our contributions

- Developed a general framework for detecting subpopulations and treatment combinations that have large deviations in their observed outcomes
- Used multidimensional constraints to scan a large number of subpopulation and treatment combinations in a computationally efficient manner
- Theoretical analysis:
  - Showed that our scoring functions with propensity reweighted outcomes removes the bias from the observed characteristics
  - Showed statistical properties of false positive rate and subpopulation exactness
- Empirical evaluation:
  - Generated interesting hypothesis related to heart disease by analyzing large, complex and observational health care claims data

# Future Research

Inputs

Categorical

Continuous

Treatments

Binary

Groups

Sequences

Outcomes

Scalar

Vector  
(Panel Outcomes)

Parametric

Non-parametric

# Future Research

Inputs

Categorical

Continuous

Treatments

Binary

Groups

Sequences

Outcomes

Scalar

Vector  
(Panel Outcomes)

Parametric

Non-parametric

# Future Research

