

An Analysis of the Relationship Between Heart Disease and Osteoarthritis

Inference through Propensity Scores

Matthew Berkowitz¹, Coco Liu¹, Barinder Thind¹, & Jiahao Tian¹

1. Simon Fraser University

Abstract

The goal of this analysis was to isolate the effect of osteoarthritis on the likelihood of developing cardiovascular disease. The provided data was characterized by mild sparsity, so the analysis was carried out on a complete-cases subset and on an imputed data set. The imputation was done through a combination of Logistic and Polytomous Bayesian Regression. There were a number of covariates that had the potential to confound this effect; in order to circumvent this problem, propensity scores were used [1] in conjunction with the LASSO [3]. More specifically, the propensity scores for all three cycles were computed for each of the two types of data sets and then another analysis was done involving the ensemble of propensity scores and the LASSO. No definitive causal relationship was uncovered although there was evidence of a statistically significant effect.

Introduction

It is important to identify as many of the factors associated with any particular condition in an attempt to minimize or remove their confounding effect. The purpose of this analysis is to understand whether this kind of antibiotic relationship exists between osteoarthritis and heart disease; the approach is summarized as follows:

- We first explore the data set in an effort to understand the nature of the missing values
- An imputation is implemented on a large and relevant subset of the data
- A propensity score approach is used for inference and the results are provided
- The LASSO was used on half the data sets (for each cycle) for variable selection, after which propensity score analysis was used on the other half of the data sets to provide odds ratio estimates

Data Description

- The data is filtered out according to study eligibility criteria so that the participants who are not 20–64 years of age, and the participants who were diagnosed as either *Rheumatoid Arthritis* or *Other* are excluded from the analysis
- There are a number of variables that include information on dietary habits, age, location, general health, marital status, substance use, blood pressure, stress levels, and income; in total, 23 covariates are used
- There are separate observations for three “cycles” which reference varying time periods and had about 130,000 observations each.

Pre-Processing & Imputation

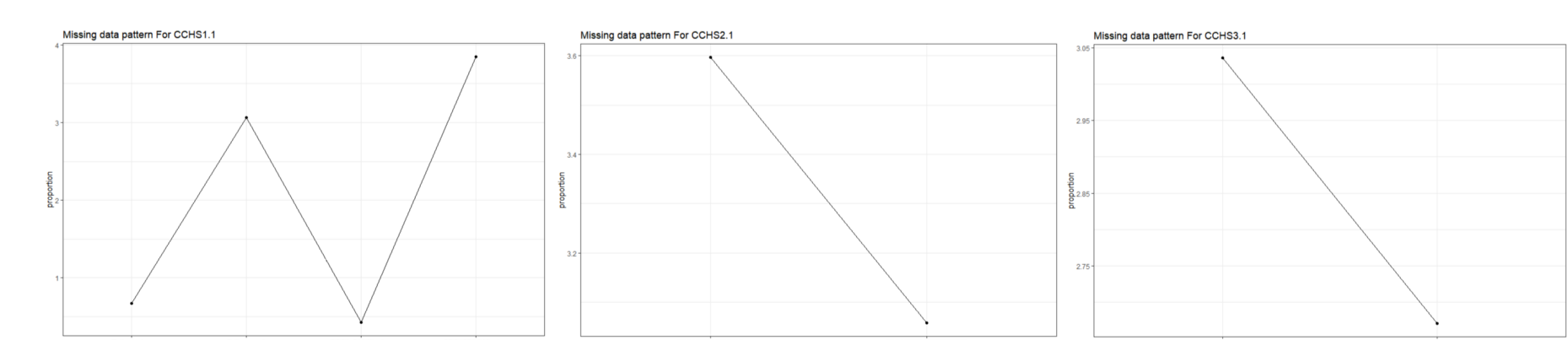
Recoding

A number of variables were recoded, for example:

- The daily consumption of vegetables was transformed from the number of servings into a categorical variable
- The geographical variable was recoded so that a distinction was only made between the territories and the provinces
- The BMI variable was recoded so that it was not a raw number, but rather categorical in nature

Polytomous Bayesian Regression for Imputation

The underlying assumption with the imputation implemented here is that the missing data is *missing at random (MAR)*. However, since the validity of this claim can easily be questioned, we decided to carry out the analysis on the complete-cases data set as well.

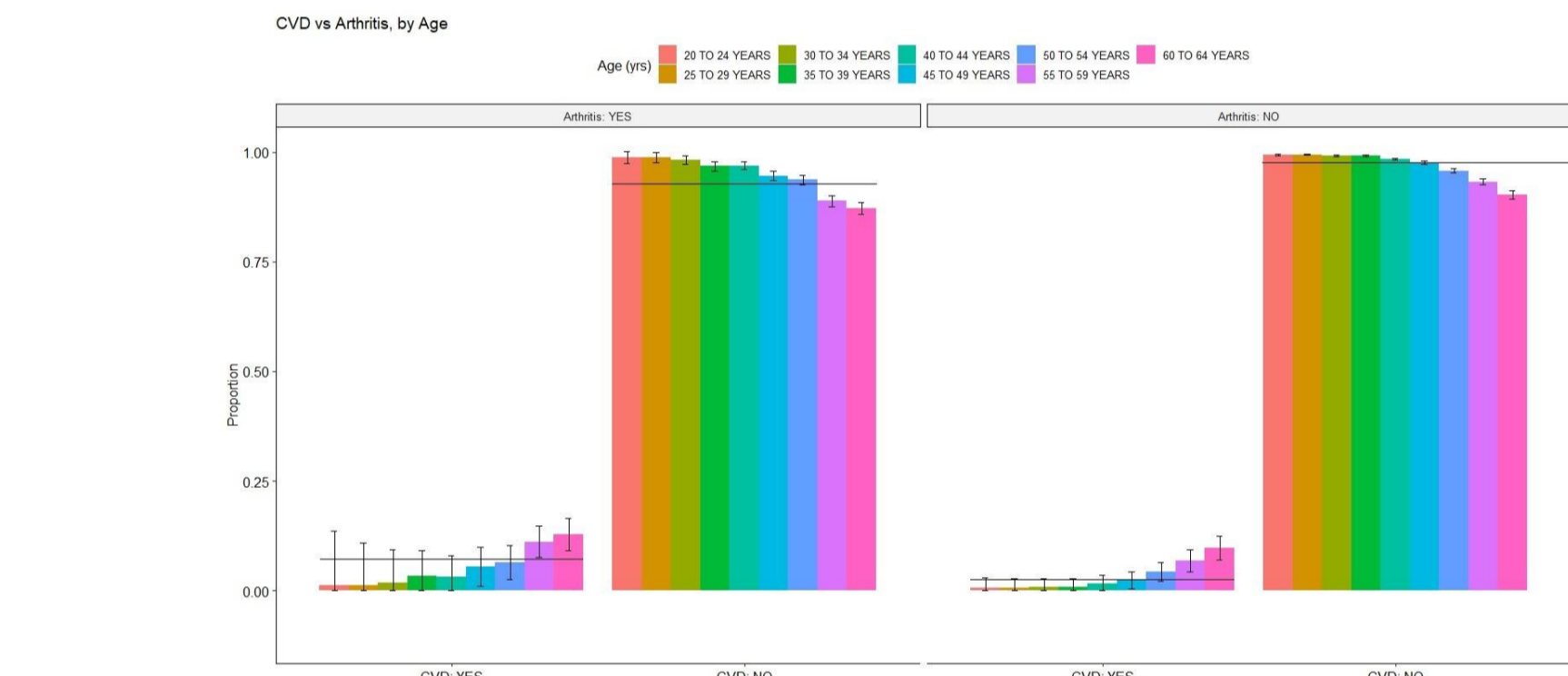


- The proportions of people who have osteoarthritis was calculated for each group. If the MCAR assumption holds, we would not expect to see any pattern of proportions associated with each group in the figure
- The plots shown indicate that the MCAR assumption does not hold in our data set, thus making analyzing the complete cases potentially inappropriate

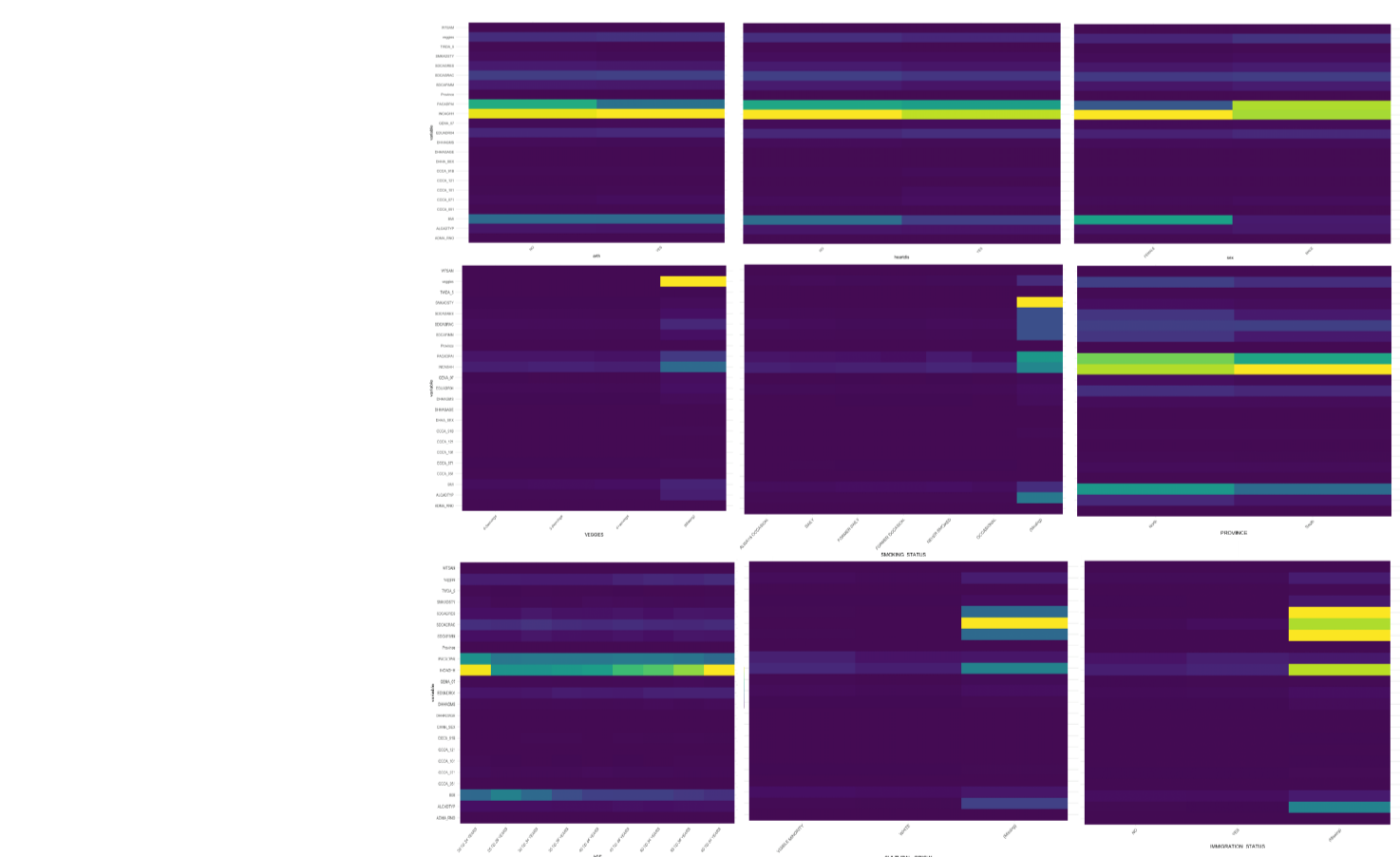
The imputation technique used here was Polytomous (Multinomial) Logistic Regression. This models how multinomial response variable Y depends on a set of m explanatory variables, $X=(X_1, X_2, \dots, X_m)$; in essence, this is a generalized linear model where the random component assumes that the distribution of Y is multinomial(n, π), where π is a vector with probabilities of “success” for each category [2]. Note that the responses here are all of the case where ordinality is not of particular importance.

Exploratory Data Analysis

In this section, some plots of interest are presented which helped guide the decisions made later in the analysis.



- The above plot shows the relationship between osteoarthritis and cardiovascular disease (CVD), broken down by age. As can clearly be seen visually, age appears to be a significant predictor of CVD, while osteoarthritis status (yes / no) appears to matter much less.



- The variable with the most missing values is the total household income from all sources, and the second is physical activity index.

Methodology

Logistic Regression

The response variable here is dichotomous and hence, a logistic regression is the simple and obvious first approach to such a problem. The generalized linear model in this problem is defined as:

$$\Pr(\text{CCCA}=121 = 1 \mid \theta) = \frac{\exp(\beta_0 + \beta_1 \text{CCCA}-051 + \dots + \beta_{21} \text{Province})}{1 + \exp(\beta_0 + \beta_1 \text{CCCA}-051 + \dots + \beta_{21} \text{Province})} \quad (1)$$

where θ is the set of parameters defining the model. This model is then expanded on with propensity scores.

LASSO

The least absolute shrinkage and selection operator (LASSO) is a least squares technique that has the effect of minimizing the coefficients for par-

Contact Information:

Matthew Berkowitz Barinder Thind
matthew_berkowitz@sfu.ca bthind@sfu.ca
Coco Liu Jiahao Tian
sla214@sfu.ca jtian_3@sfu.ca

ticular covariates down to 0. It is mathematically defined as:

$$\hat{\beta}_{\text{lasso}} = \underset{\beta \in R^m}{\operatorname{argmin}} \|y - \mathbf{X}\beta\|_2^2 + \lambda \sum_{i=1}^m |\beta_i| \quad (2)$$

$$= \underset{\beta \in R^n}{\operatorname{argmin}} \underbrace{\|y - \mathbf{X}\beta\|_2^2}_{\text{Loss}} + \underbrace{\lambda \sum_{i=1}^m |\beta_i|}_{\text{Penalty}} \quad (3)$$

Here, this was used to lower the dimensionality of our covariates so that we could compare the results between the non-LASSO and LASSO propensity score approaches.

Propensity Scores

A propensity score, $p(x_i = 1|\theta)$ is the conditional probability, for subject i , ($i = 1, \dots, n$), of being assigned to some particular treatment (in this case, 1) given some covariates, θ . In other words:

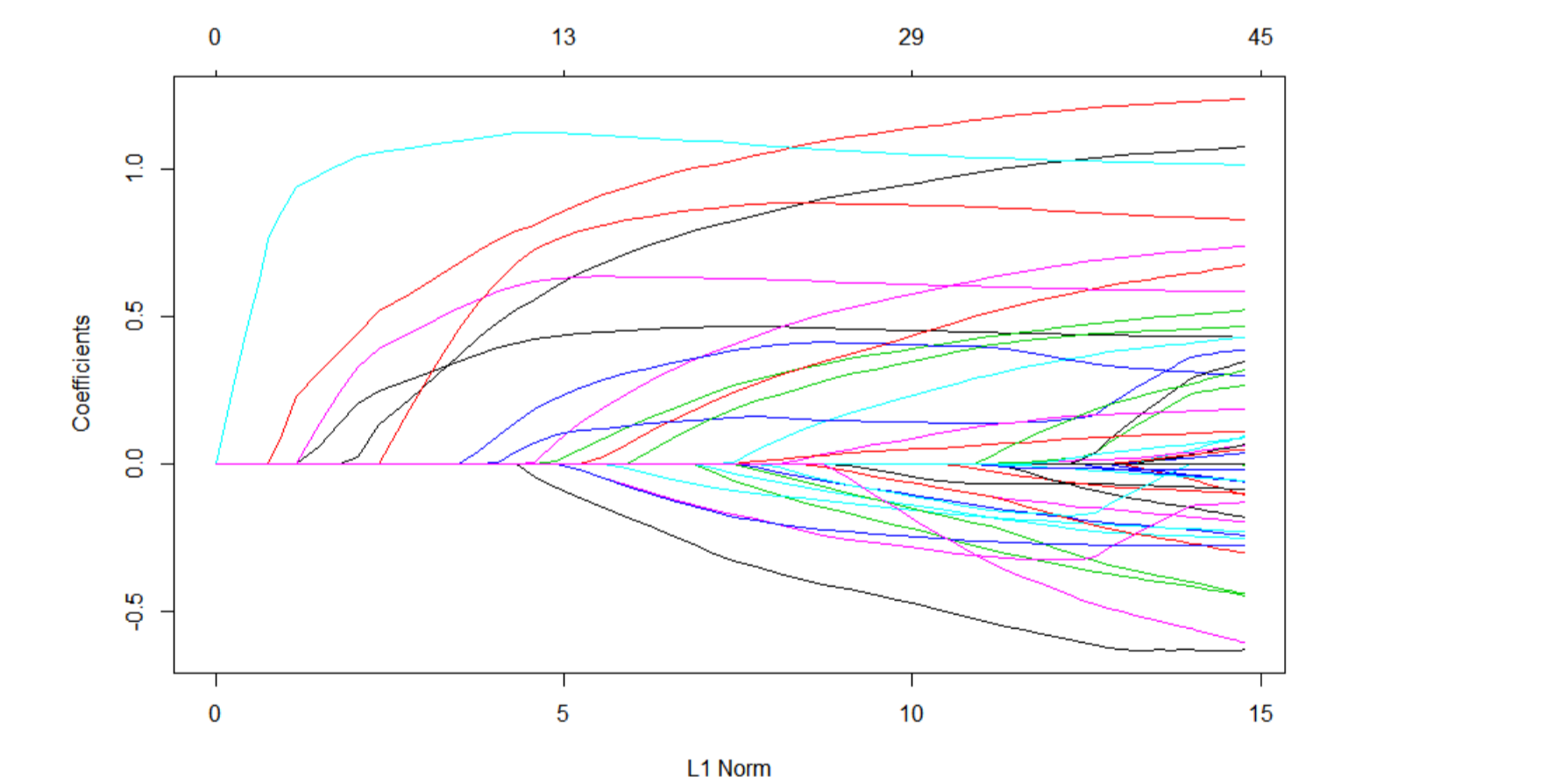
$$P(l) = P(x_i = 1|\theta) \quad (4)$$

Naturally, it follows that if a randomized design were used, we would expect the conditional probability here to be 0.5. For this case study, we applied as follows:

- We used Propensity Score Matching to estimate the effect of having osteoarthritis by accounting for an array of covariates: Propensity score, $P(L)=P(O=1|L)$, where O is having osteoarthritis, and L are the covariates. This attempts to deconfound the effect of osteoarthritis on cardiovascular disease.
- The final model uses the propensity score matched sample and incorporates survey weights into a quasibinomial logistic regression. The matched sample sizes from each cycle, the OR estimates and associated CIs are summarized in the results tables.

Results & Inference

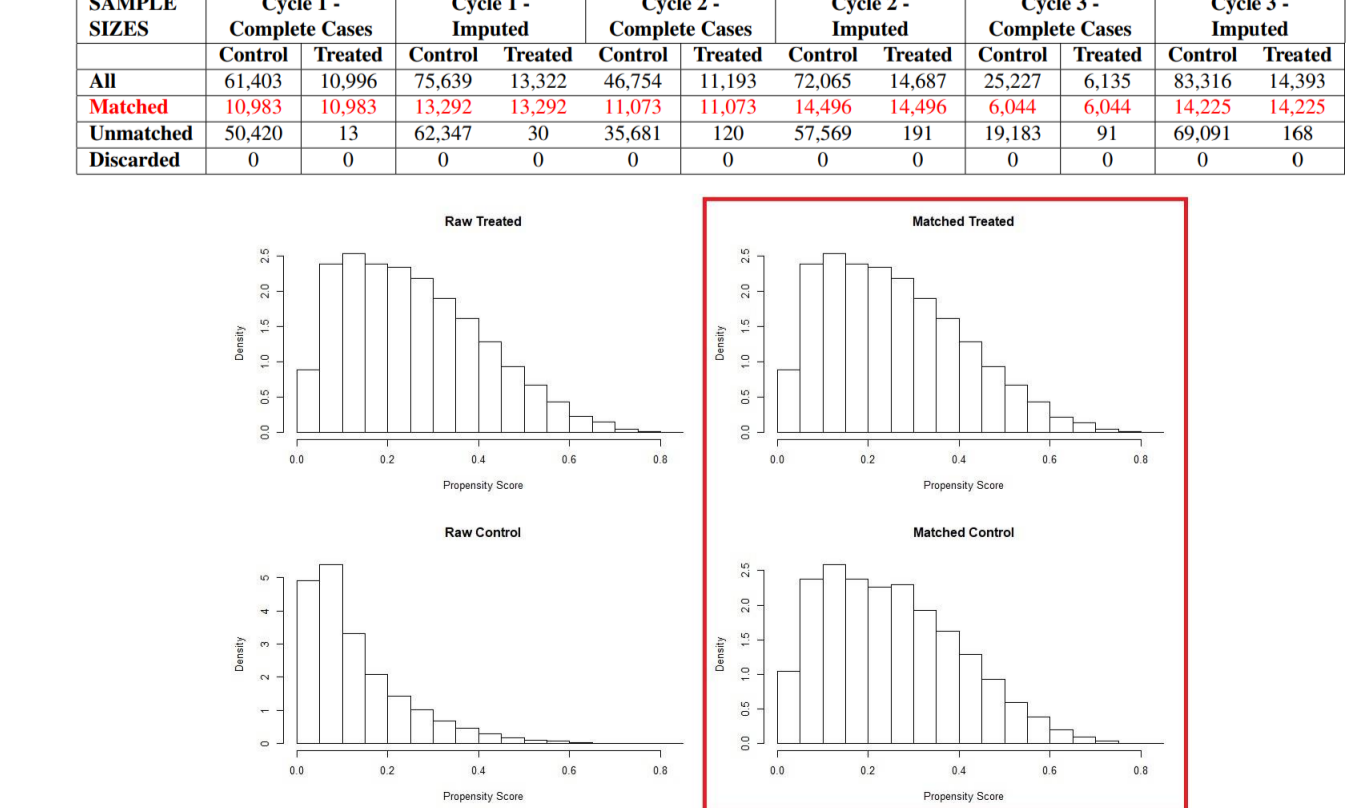
LASSO Path



- The variables selected by this regularization were: age, sex, income, type of smoker, blood pressure status, diabetic status, emphysema or chronic obstructive pulmonary disease (COPD) status, self-perceived stress.

Propensity Scores

The following tables/plots provide information on our results.



Table/Figure: Covariate balance in matched sample

SUMMARY ESTIMATES (PSA)	Cycle 1 - Complete	Cycle 1 - Imputed	Cycle 2 - Complete	Cycle 2 - Imputed	Cycle 3 - Complete	Cycle 3 - Imputed
OR Estimate	1.39	1.45	1.39	1.55	1.55	1.48
95% Confidence Interval - Lower Bound	1.18	1.25	1.34	1.33	1.23	1.26
95% Confidence Interval - Upper Bound	1.64	1.68	1.90	1.80	1.95	1.73

Table 2: Odds ratio results using propensity scores

SUMMARY ESTIMATES (PSA / LASSO)	Cycle 1 - Lasso	Cycle 2 - Lasso	Cycle 3 - Lasso
OR Estimate	1.49	1.70	1.67
95% Confidence Interval - Lower Bound	1.18	1.40	1.43
95% Confidence Interval - Upper Bound	1.88	2.07	1.96

Table 3: Propensity score OR results after the LASSO regularization

Stratified by CCCA_051	NO	YES	OR
20 TO 24 YEARS	234 (1.9)	358 (1.7)	
25 TO 29 YEARS	352 (3.2)	527 (3.8)	
30 TO 34 YEARS	577 (5.3)	1073 (7.8)	
35 TO 39 YEARS	856 (7.8)	884 (8.8)	
40 TO 44 YEARS	1286 (11.7)	1252 (12.8)	
45 TO 49 YEARS	1636 (14.9)	1583 (14.4)	
50 TO 54 YEARS	1999 (18.2)	2063 (18.2)	
55 TO 59 YEARS	2830 (18.5)	2895 (18.1)	
60 TO 64 YEARS	3815 (28.5)	3927 (28.3)	
DUAL_SEX = MALE (%)	4184 (38.1)	4187 (37.4)	
DUAL_SEX = FEMALE (%)	776 (7.2)	742 (6.8)	
COWBOY-LAW	6156 (56.8)	6162 (55.8)	
HANDLED	1451 (13.2)	1453 (13.0)	
SINGLE	2065 (21.7)	2077 (24.2)	
WIDOW/SEVERELY	1451 (13.2)	1453 (13.0)	
SOCAGRA = MATE (%)	2830 (18.5)	2895 (18.1)	
SOCAGRA = YES (%)	1217 (11.3)	1242 (11.4)	
SOCAGRA = NO (%)	88 (6.8)	88 (8.7)	
10 YEARS OR MORE	1151 (14.5)	1112 (14.1)	
NOT APPLICABLE	9746 (88.7)	9775 (100.0)	
EDUCATION	3113 (28.3)	3114 (28.7)	
OTHER POST-SEC.	840 (7.6)	828 (7.5)	
POST-SEC. GRAD.	5047 (46.8)	5038 (45.8)	
SECONDARY GRAD.	1585 (14.1)	1581 (14.6)	
DICOM (%)	2117 (19.3)	2124 (19.3)	
\$15,000-\$29,999	2788 (24.6)	2848 (24.2)	
\$30,000-\$49,999	2628 (23.9)	2566 (23.4)	
\$50,000-\$79,999	1822 (16.6)	1786 (16.3)	
\$80,000 OR MORE	1659 (15.1)	1786 (16.3)	
LESS THAN \$15,000	57 (6.5)	62 (6.4)	
NO INCOME	2351 (28.5)	2387 (28.1)	
PACAPAZ (%)	6174 (56.2)	6227 (56.7)	
INACTIVE	2074 (21.3)	2048 (21.2)	
MODERATE	18867 (91.7)	18581 (92.8)	
TOTAL = YES (%)	115 (1.8)	112 (1.2)	
DAILY	3024 (31.2)	3448 (31.8)	
FORMER DAILY	3279 (28.9)	3325 (30.3)	
FORMER OCCASIONALLY	1375 (14.5)	1415 (13.1)	
NEVER SMOKE	2367 (21.8)	2317 (21.3)	
NEVER SMOKE	319 (1.8)	319 (2.9)	
ALCATYP (%)	1842 (14.8)	1828 (15.1)	
POWER DRINKER	484 (4.4)	483 (4.4)	
NEVER DRINK	275 (2.4)	268 (2.4)	
REGULAR DRINKER	6879 (55.3)	6972 (54.4)	
CCCA_051 = YES (%)	2482 (22.4)	2438 (22.1)	
CCCA_051 = YES (%)	751 (6.9)	889 (7.4)	
CCCA_051 = NO (%)	18026 (93.3)	18026 (92.5)	
NOT APPLICABLE	348 (1.3)	327 (1.7)	
YES	348 (1.3)	288 (1.8)	
GENA_051 (%)	4260 (38.8)	4262 (39.2)	
ATTEMPT	818 (7.4)	821 (7.5)	
NOT AT ALL	1853 (16.6)	971 (8.8)	
NOT SURE	2854 (26.1)	2854 (26.1)	
QUITE A BIT	2782 (25.3)	2828 (25.7)	
WE (%)	4854 (36.9)	3933 (35.8)	
overweight	4720 (41.3)	4885 (42.5)	
underweight	288 (1.8)	345 (1.7)	
3-servings	3429 (31.2)	3487 (31.8)	
2-servings	5150 (47.8)	5242 (47.7)	
1-servings	2395 (21.8)	2314 (21.3)	
Province = South (%)	18882 (98.4)	18818 (98.5)	

Figure: Regression results using propensity scores

Conclusions & Future Considerations

The goal of this analysis was to isolate the effect of osteoarthritis on cardiovascular disease. We did the analysis on two sets of data and used an ensemble of the LASSO and propensity scores. Through this, our OR estimates of the effect of osteoarthritis on cardiovascular disease are all similar, with the LASSO-derived estimates producing greater uncertainty, as seen with the wider confidence intervals. There appears to be a statistically significant impact of osteoarthritis on the risk of cardiovascular disease; however, it is doubtful that the effect is causative. Further research is needed to study the mediating variables responsible for this relationship.

References

- [1] Eshan Karim. Ps-survey workshop (participant).
- [2] Penn State. 8.1 - polytomous (multinomial) logistic regression.
- [3] Robert Tibshirani. Regression shrinkage and selection via the lasso. *J. Royal Stat. Soc.*, 58:267–288, 1996.

Acknowledgements

Thanks to Joan Hu and Rachel Altman for their helpful comments. Thanks to SFU's department of Statistics & Actuarial Science for their support.