

1 Introduction

In the early eighties, hospital administrators in the US sought to reduce hospitalization costs by changing the reimbursement system. Instead of paying hospitals per day of hospitalization, hospitals were paid a fixed rate for the type of hospitalization and the procedures performed. Following implementation of this law, the length of stay at hospitals dropped dramatically, although evidence exists that it was already in decline.

Some worry that the new system created a perverse incentive to discharge patients early, and admit them again at a later date. To ensure proper quality of care in the hospitals while keeping costs controlled, administrators have sought to establish useful quality of care metrics. Hospital readmissions have been identified as a simple metric that can establish a baseline of care; if an abnormally high number of patients from a certain hospital are quickly readmitted, it could indicate poor quality of care.

Since hospitals admit patients with varying risk of readmission, it is important to accurately estimate the effect of hospital treatment on readmission, independent of patient case-mix. Several readmission risk models have been criticized on the basis that they failed to distinguish preventable readmissions from readmissions due to chance alone. Despite evidence that clinicians cannot reliably distinguish preventable and non-preventable readmissions, pairs of diagnosis codes have been developed that are "potentially preventable".

The counterfactual model clarifies the study question: What is the difference in the proportion with an emergency readmission within 30 days if all patients had attended Hospital A vs Hospital B? The notion of a "preventable" readmission is implicit; if a patient would have been readmitted if treated at another hospital, then the readmission was preventable. Since we cannot directly observe the desired proportions, we estimate it by attempting to recreate exchangeability (control for confounding) among the populations that visited different hospitals.

Hospital administrative data is complex, but information-rich. Drug prescriptions, diagnoses, and medical procedures can all provide important information on how the effect of hospital care on readmission is confounded by patient health. Typically, however, the administrative hospital data is simplified to "comorbidity scores".

To estimate propensity of attending different hospitals as a function of hundreds of covariates relating to patient health, we used a non-parametric machine learning technique, random forest. We combined this propensity with a more traditional model of hospital readmissions as a function of hospital care and patient characteristics (penalized logistic regression) with a doubly robust technique, targeted maximum likelihood estima-

tion (TMLE).

2 Methods

2.1 Data

2.1.1 Cohort selection

We used a cohort extracted from a Canadian provincial (Quebec) administrative database of hospitalizations, obtained from the *Régie de l'assurance maladie du Québec* (RAMQ). We enrolled patients into this cohort on the month that two conditions were satisfied: 1) they had at least one diagnosis of a respiratory illness (the exact list of respiratory International Classification of Diseases, 9th Revision [ICD-9] codes is given in the Appendix) between January 1st, 1996 and March 31, 2006 (the study period), while living in the 2006 census metropolitan area of Montreal, and 2) were at least 65 years of age. We used this cohort because it represents the majority of 65-year olds who were hospitalized in the region during the study period.

From among this cohort, we selected hospital discharges for those who had accrued at least one continuous year in the cohort preceding the time of admission. We restricted our data to only the discharges from the twenty hospitals with the most discharges of patients 65 years or older within the study period; the twenty hospitals accounted for 75% of all such discharges. We only selected hospital discharges which resulted from hospital stays of at least one day. Therefore, the earliest possible hospital discharge was January 2, 1997.

2.1.2 Hospital readmissions

The unit of analysis in all models was the hospital discharge; a person could be discharged multiple times. A hospital readmission was defined as an emergency hospital admission to any Quebec hospital in the 30 days following a discharge. A person who died or had a non-emergency readmission in the 30 days following discharge was considered not readmitted.

2.1.3 Disease types

From among the identified hospital discharges, we selected only those with one of three high-volume admission diagnoses with high rates of hospital readmissions: pneumonia, acute myocardial infarction (AMI), and heart failure, the three initial conditions selected by the Centers for Medicare and Medicaid Services (CMS) to implement the Hospital Readmissions Reduction Program mandated by the Affordable Care Act. We identified each of the admission diagnoses using ICD-9 codes; for pneumonia we used codes ranging from 480-487, for heart failure we used all 428 codes, and for AMI we used all 410 codes. The following methods were applied individually to all three disease subsets.

2.2 Confounders

For each hospital discharge, we collected variables that measured states at the time of admission, or events that occurred prior to the hospital admission, and which may confound the relationship between hospital care and readmission. We used the demographic characteristics (age at time of admission (years), sex, birth year-month), the number of previous readmissions (within the study period), the admission diagnosis (as measured by the specific ICD-9 code). We also included the day of week of discharge, which has been previously shown to have an association with readmissions, and the month of discharge, because we hypothesized that readmission risk would vary by seasons in Montreal.

Additionally, for each discharge, we collected the Quebec hospital diagnoses, Quebec hospital procedures, and drugs dispensed outside of the hospital but inside Quebec, in the year preceding the admission. The hospital procedures were recorded in the Canadian Classification of Diagnostic, Therapeutic, and Surgical Procedures (CCP) system. Hospital diagnostic codes were coded using the ICD-9 system. Finally, drugs which were prescribed and dispensed outside the hospital, and were being taken on the day of admission were also recorded for each patient in the *code commune* system, which categorizes drugs based on the chemical compound. To ease computation, before fitting any model, we removed any diagnosis, procedure or drug that occurred less than 30 times among all discharges. We chose 30 because it appeared to be a natural breakpoint; if the number of variables included is a function f of the threshold, then the first derivative of f dropped at 30 for all three disease categories.

2.3 Descriptive analysis

2.3.1 Choropleth

We believed that census tract of residence would strongly affect the probability of admission to the hospital nearest that census tract. We plotted choropleths of the rate of attendance at the twenty hospitals by census tract. The numerator was the number of live discharges at the hospital, and the denominator was the number of person-years accumulated in that census tract of residence by cohort members when their admissions would have been eligible (after the first continuous year within the cohort).

2.4 Models

2.4.1 Model g - probability of exposure

We developed a multinomial model g that predicted the probability of attending each of the twenty hospitals (the exposure) as a function of all the confounders. To fit g we used a random forest, a non-parametric model based on decision trees¹.

Our trees were decision stumps; we only used one splitting node at in each tree. We arbitrarily chose to grow 1200 trees (decision stumps), and then measured the accuracy as a function of the number of trees to ensure that growing further trees would be unlikely to significantly improve accuracy. When measuring the accuracy for each discharge, we only used trees for which the discharge was "out-of-bag", that is, we only used trees for which the bootstrap replicate did not include the discharge.

To assess the importance of the variables in predicting which hospital a patient will choose, for each tree, we calculate the increase in homogeneity of classes between the root of the tree and the leaves. To measure the homogeneity of classes, we repurpose a metric that is typically used to measure equality (homogeneity) of income, the Gini coefficient². The homogeneity is defined as the sum of squared proportions in each class (hospital), with a maximum of 1 (all discharges at the same hospital) and a minimum of $1/20$ (the discharges were evenly divided between the 20 hospitals). If the patients within the leaves of the tree made relatively homogeneous hospital choices, this variable predicts hospital choice well.

Because the model was used solely to estimate the *probability* of admission to to specific hospitals (and not to predict exactly which hospital was attended), we configured the model to favour calibration over discrimina-

tion: we weighted each of the twenty predicted hospitals by the inverse of the proportion of discharges at that hospital. We multiplied each proportion used to calculate the Gini coefficient by this weight.

Random forest traditionally classifies each item by majority vote; we converted this into a probability by taking the proportion of votes for each hospital (using only out-of-bag trees for each discharge).

2.4.2 Q model - probability of outcome

For the Q model, we developed another random forest model, very similar to the g model, except that instead of predicting the choice of hospital, we directly predicted 30-day readmission. We also included a set of 19 indicator variables for the hospital attended. We also calibrated this model based on the inverse of the proportion of those readmitted.

We also fit a regularized logistic regression model for Q , using cyclical coordinate descent to estimate the parameters efficiently in our sparse but large information matrix. We penalized the likelihood by the ℓ_2 norm, that is, the sum of the squares of the normalized regression parameters. The scale of the penalty was determined by the parameter λ . We optimized the selection of the penalty scale for the best partial likelihood using a nested a 10-fold cross-validation. Within each fold we assessed 100 λ -values spaced evenly between $\max(\lambda) \times 10^{-4}$ and $\max(\lambda)$, where $\max(\lambda)$ was the smallest λ -value that would result in a model with no non-zero coefficients.

2.4.3 Variable importance

For each tree, we calculated the decrease in Gini homogeneity when comparing the split nodes of the tree to the root. The Gini impurity at any node is, for all out-of-bag discharges.

For each hospital, the proportion of patients that chose this hospital p_h , multiplied by the probability of guessing that this patient would choose this hospital based on the distribution of patients $1 - p_h$. Summed over all hospitals, this is the Gini impurity. Gini impurity is 1 when all patients chose the same hospital.

Gini is defined as "inequity" when used in describing a society's distribution of income, or a measure of "node impurity" in tree-based classification. A low Gini (i.e. higher decrease in Gini) means that a particular predictor variable plays a greater role in partitioning the data into the defined classes.

2.4.4 Model Q* - updated targeted maximum likelihood estimation

To avoid convergence problems for this one parameter model, we fit the parameter using a quasi-Newton method simultaneously discovered by Broyden³, Fletcher⁴, Goldfarb⁵ and Shanno⁶ (BFGS).

2.5 Software

The data were cleaned and prepared for statistical analysis using the Postgres relational database (version 9.2.6). We implemented our models using the R statistical package (version 3.1.0)⁷. We implemented the random forest using the "bigrf" package (version 0.1.9)⁸. We implemented the GLM fitting using coordinate descent using the "glmnet" package (version 1.9.5)⁹. We plotted our figures using the "ggplot2" package (version 1.0.0)¹⁰.

3 Results

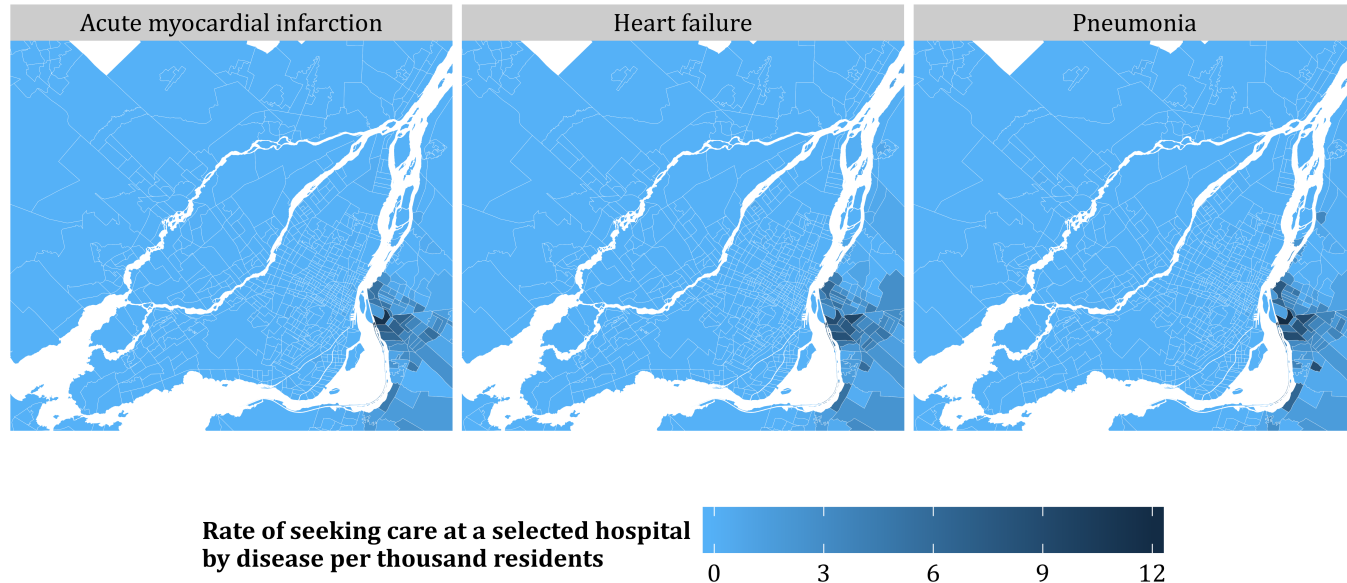


Figure 1: Rates of seeking care at a selected hospital by admission diagnosis. A descriptive sentence

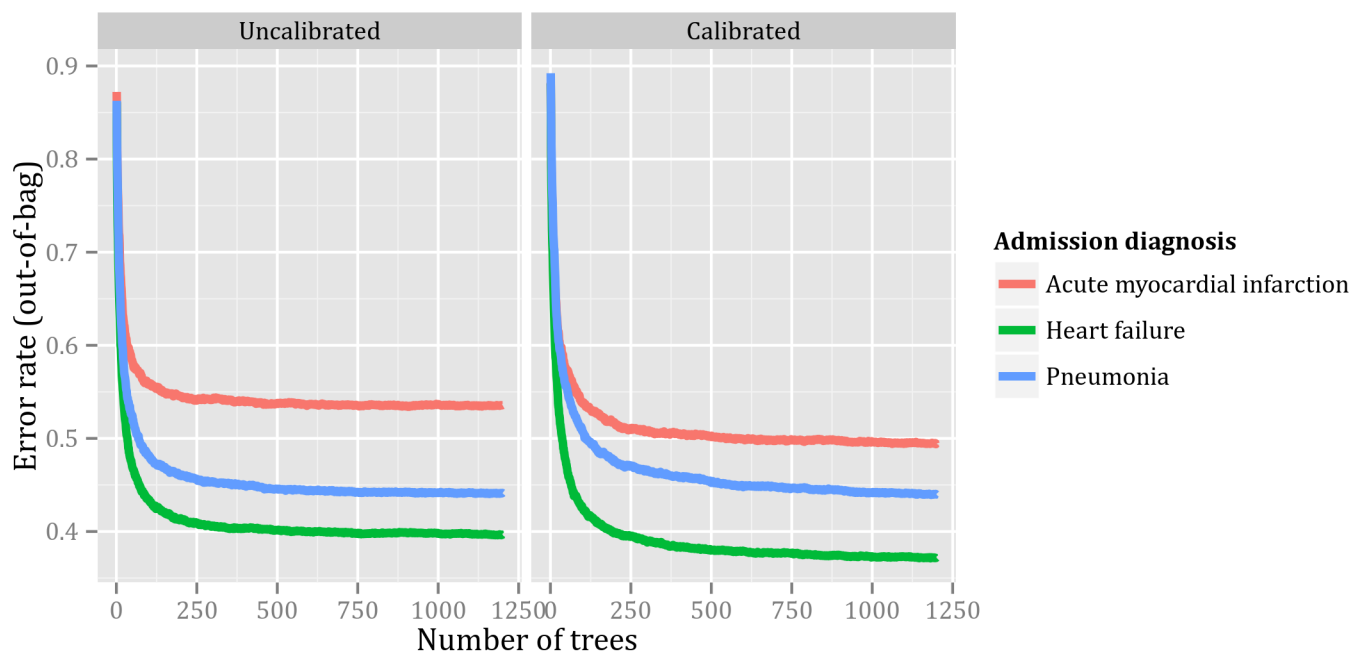


Figure 2: Error rate for random forest model of hospital choice. A descriptive sentence.

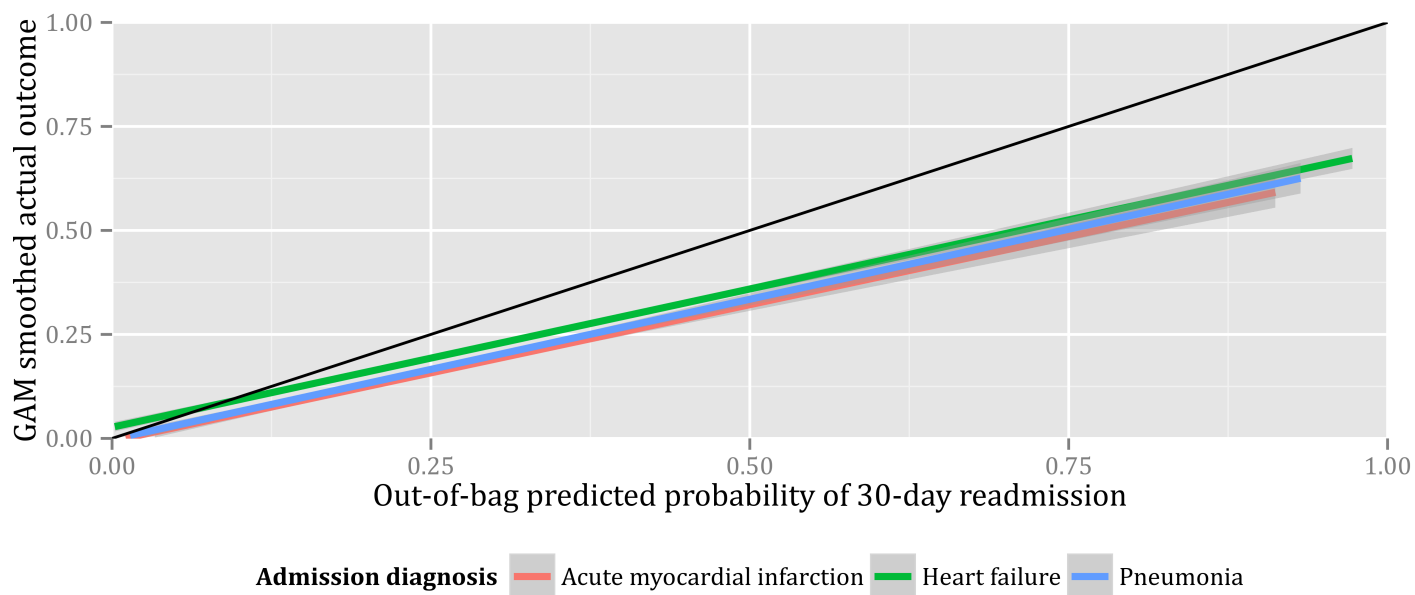


Figure 3: Calibration for random forest model of hospital readmission. A descriptive sentence.

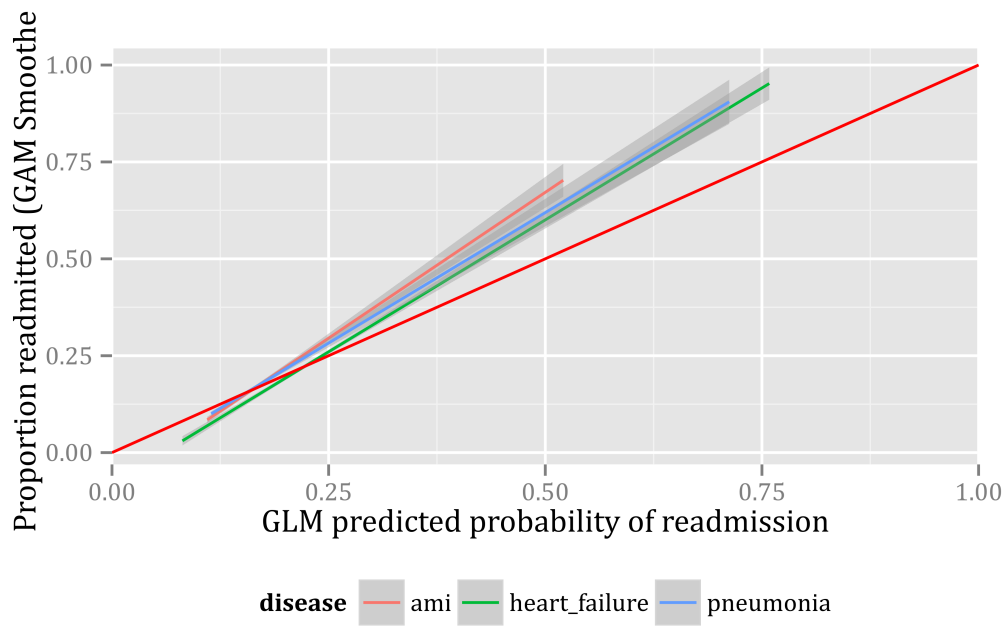


Figure 4: Calibration for GLMnet model of hospital readmission. A descriptive sentence.



Figure 5: 10 most important variables for the G model, by disease. A descriptive sentence.

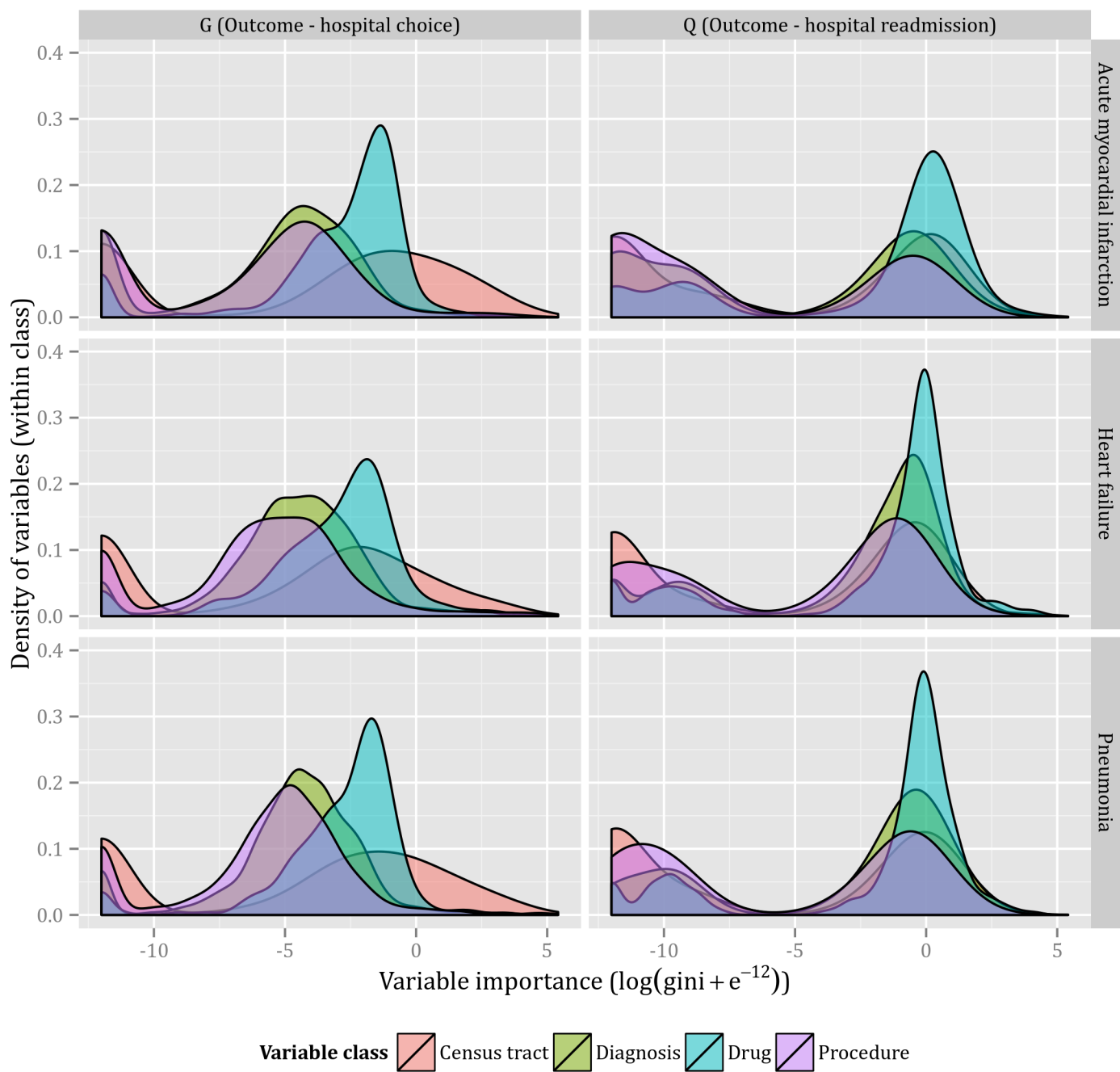


Figure 6: Variable importance by model and variable class. A descriptive sentence.

Hospital					Length of stay (mean days)					Random Forest			GLMnet		
	Admitted	Died	(%)	Discharged	Admitted	Died	Discharged	Readmitted	(%)	Q	ϵ	Q*	Q	ϵ	Q*
1	763	112	0.15	651	15	10	16	105	0.16	0.25	-0.09	0.22	0.16	0.00	0.16
2	1557	148	0.10	1409	13	12	14	191	0.14	0.25	-0.09	0.22	0.16	0.00	0.16
3	606	83	0.14	523	14	12	14	84	0.16	0.25	-0.06	0.23	0.16	0.00	0.16
4	1022	125	0.12	897	11	7	12	136	0.15	0.25	-0.08	0.22	0.16	0.00	0.16
5	729	150	0.21	579	14	12	15	98	0.17	0.26	-0.05	0.24	0.16	0.00	0.16
6	826	119	0.14	707	11	8	12	106	0.15	0.26	-0.06	0.23	0.16	0.00	0.16
7	1491	241	0.16	1250	15	14	16	216	0.17	0.25	-0.04	0.24	0.16	0.00	0.16
8	1270	198	0.16	1072	14	12	15	138	0.13	0.24	-0.11	0.20	0.16	-0.01	0.15
9	780	152	0.19	628	13	12	14	130	0.21	0.27	-0.09	0.23	0.16	0.00	0.16
10	778	124	0.16	654	13	9	14	123	0.19	0.28	-0.04	0.26	0.16	0.00	0.16
11	705	125	0.18	580	12	12	12	97	0.17	0.26	-0.06	0.24	0.16	0.01	0.17
12	1284	266	0.21	1018	15	14	15	166	0.16	0.25	-0.12	0.21	0.16	0.00	0.16
13	739	86	0.12	653	15	13	15	110	0.17	0.26	-0.04	0.24	0.16	0.00	0.16
14	1307	184	0.14	1123	12	8	13	210	0.19	0.26	-0.07	0.23	0.16	0.00	0.16
15	1152	168	0.15	984	16	14	17	129	0.13	0.25	-0.11	0.21	0.16	-0.01	0.15
16	408	70	0.17	338	9	11	9	43	0.13	0.25	-0.06	0.23	0.16	0.00	0.16
17	807	123	0.15	684	11	15	11	134	0.20	0.27	-0.04	0.25	0.16	0.00	0.16
18	894	144	0.16	750	13	17	12	116	0.15	0.25	-0.04	0.24	0.16	0.00	0.16
19	499	94	0.19	405	9	12	8	50	0.12	0.24	-0.08	0.21	0.16	-0.01	0.15
20	1025	184	0.18	841	13	9	14	143	0.17	0.26	-0.14	0.21	0.16	0.00	0.16

Table 1: Acute myocardial infarction (AMI).

Hospital					Length of stay (mean days)					Random Forest			GLMnet		
	Admitted	Died	(%)	Discharged	Admitted	Died	Discharged	Readmitted	(%)	Q	ε	Q*	Q	ε	Q*
1	1229	141	0.11	1088	13	15	13	248	0.23	0.29	-0.17	0.23	0.22	-0.02	0.21
2	2071	166	0.08	1905	13	19	12	441	0.23	0.29	-0.07	0.27	0.22	0.00	0.22
3	1243	134	0.11	1109	14	18	13	285	0.26	0.30	-0.19	0.23	0.22	-0.01	0.21
4	1076	122	0.11	954	12	15	12	214	0.22	0.29	-0.09	0.26	0.22	-0.01	0.22
5	1550	181	0.12	1369	12	17	11	288	0.21	0.29	-0.14	0.24	0.22	-0.01	0.21
6	827	107	0.13	720	11	13	11	128	0.18	0.29	-0.12	0.24	0.22	-0.02	0.21
7	2917	386	0.13	2531	13	17	12	666	0.26	0.30	-0.03	0.29	0.23	0.00	0.24
8	1456	197	0.14	1259	12	19	11	232	0.18	0.28	-0.14	0.23	0.22	-0.02	0.21
9	881	111	0.13	770	13	21	12	157	0.20	0.29	-0.10	0.25	0.22	0.00	0.22
10	1410	149	0.11	1261	10	18	9	311	0.25	0.29	-0.12	0.25	0.22	-0.01	0.22
11	1297	153	0.12	1144	17	24	16	258	0.23	0.29	-0.11	0.25	0.22	-0.02	0.21
12	1323	162	0.12	1161	15	20	15	192	0.17	0.28	-0.07	0.25	0.20	0.00	0.20
13	1231	102	0.08	1129	16	21	15	262	0.23	0.29	-0.16	0.23	0.22	-0.01	0.21
14	2110	234	0.11	1876	11	15	11	424	0.23	0.29	-0.10	0.25	0.22	-0.01	0.22
15	1389	190	0.14	1199	18	20	17	203	0.17	0.28	-0.07	0.26	0.22	0.00	0.22
16	681	94	0.14	587	11	15	10	111	0.19	0.29	-0.12	0.24	0.22	-0.03	0.20
17	1438	139	0.10	1299	12	26	10	328	0.25	0.30	-0.09	0.27	0.22	0.01	0.23
18	1984	212	0.11	1772	14	24	13	438	0.25	0.29	-0.10	0.26	0.22	-0.01	0.22
19	932	99	0.11	833	9	15	8	163	0.20	0.29	-0.12	0.25	0.22	-0.02	0.21
20	1048	167	0.16	881	11	12	11	171	0.19	0.29	-0.08	0.26	0.22	0.00	0.22

Table 2: Heart failure

Hospital					Length of stay (mean days)					Random Forest			GLMnet		
	Admitted	Died	(%)	Discharged	Admitted	Died	Discharged	Readmitted	(%)	Q	ε	Q*	Q	ε	Q*
1	1184	176	0.15	1008	13	14	13	159	0.16	0.23	-0.16	0.18	0.16	-0.01	0.15
2	199	11	0.06	188	13	15	12	31	0.16	0.24	-0.16	0.19	0.16	-0.01	0.15
3	1085	132	0.12	953	13	15	13	160	0.17	0.23	-0.12	0.19	0.16	0.00	0.15
4	863	91	0.11	772	13	14	12	113	0.15	0.23	-0.10	0.20	0.16	0.00	0.16
5	923	147	0.16	776	11	13	11	143	0.18	0.25	-0.17	0.19	0.16	0.00	0.16
6	788	136	0.17	652	12	15	11	89	0.14	0.23	-0.13	0.19	0.16	0.00	0.15
7	2194	228	0.10	1966	15	16	14	328	0.17	0.24	-0.08	0.21	0.16	0.00	0.16
8	1485	243	0.16	1242	11	14	11	173	0.14	0.23	-0.12	0.19	0.16	-0.01	0.15
9	990	166	0.17	824	15	16	15	158	0.19	0.25	-0.14	0.20	0.16	0.00	0.16
10	1214	139	0.11	1075	11	14	11	181	0.17	0.23	-0.12	0.19	0.16	0.00	0.15
11	892	147	0.16	745	14	20	14	119	0.16	0.24	-0.10	0.20	0.16	0.00	0.16
12	1102	185	0.17	917	14	13	15	91	0.10	0.22	-0.18	0.16	0.15	-0.04	0.12
13	1914	204	0.11	1710	14	14	14	325	0.19	0.24	-0.07	0.22	0.16	0.00	0.16
14	1980	278	0.14	1702	11	12	11	263	0.15	0.23	-0.12	0.19	0.16	-0.01	0.15
15	1365	179	0.13	1186	11	14	11	163	0.14	0.23	-0.12	0.19	0.16	-0.01	0.15
16	541	77	0.14	464	13	17	12	46	0.10	0.22	-0.11	0.19	0.16	-0.01	0.15
17	1338	163	0.12	1175	10	18	9	193	0.16	0.24	-0.10	0.20	0.16	0.00	0.16
18	1356	168	0.12	1188	15	22	14	200	0.17	0.23	-0.09	0.20	0.16	0.00	0.16
19	1020	123	0.12	897	10	19	9	122	0.14	0.23	-0.13	0.18	0.16	-0.02	0.14
20	1152	171	0.15	981	15	16	14	126	0.13	0.23	-0.10	0.20	0.16	0.00	0.15

Table 3: Pneumonia

4 Discussion

We did not attempt to measure whether *individual* readmissions were preventable, we instead opted to place readmissions within a counterfactual framework, and only identify the difference in readmission risk between hospitals, after controlling for differences in case-mix.

References

- [1] Leo Breiman. "Random Forests". English. In: *Machine Learning* 45.1 (2001), pp. 5–32. ISSN: 0885-6125. DOI: 10.1023/A:1010933404324. URL: <http://dx.doi.org/10.1023/A:1010933404324>.
- [2] Corrado Gini. *Variabilità e Mutabilità: Contributo allo studio delle distribuzioni e delle relazioni statistiche*. ita. Bologna: C. Cuppini, 1912.
- [3] C. G. Broyden. "The Convergence of a Class of Double-rank Minimization Algorithms 1. General Considerations". en. In: *IMA Journal of Applied Mathematics* 6.1 (1970), pp. 76–90. ISSN: 0272-4960, 1464-3634. DOI: 10.1093/imamat/6.1.76. URL: <http://imamat.oxfordjournals.org/cgi/doi/10.1093/imamat/6.1.76> (visited on 06/02/2014).
- [4] R. Fletcher. "A new approach to variable metric algorithms". en. In: *The Computer Journal* 13.3 (Mar. 1970), pp. 317–322. ISSN: 0010-4620, 1460-2067. DOI: 10.1093/comjnl/13.3.317. URL: <http://comjnl.oupjournals.org/cgi/doi/10.1093/comjnl/13.3.317> (visited on 06/02/2014).
- [5] Donald Goldfarb. "A family of variable-metric methods derived by variational means". en. In: *Mathematics of Computation* 24.109 (Jan. 1970), pp. 23–23. ISSN: 0025-5718. DOI: 10.1090/S0025-5718-1970-0258249-6. URL: <http://www.ams.org/jourcgi/jour-getitem?pii=S0025-5718-1970-0258249-6> (visited on 06/02/2014).
- [6] D. F. Shanno. "Conditioning of quasi-Newton methods for function minimization". en. In: *Mathematics of Computation* 24.111 (Sept. 1970), pp. 647–647. ISSN: 0025-5718. DOI: 10.1090/S0025-5718-1970-0274029-X. URL: <http://www.ams.org/jourcgi/jour-getitem?pii=S0025-5718-1970-0274029-X> (visited on 06/02/2014).
- [7] R. Core Team. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing, 2014. URL: <http://www.R-project.org/>.
- [8] Aloysius Lim, Leo Breiman, and Adele Cutler. *bigrf: Big Random Forests: Classification and Regression Forests for Large Data Sets*. R package version 0.1-9. 2014. URL: <https://github.com/alloysius-lim/bigrf>.

- [9] Jerome H. Friedman, Trevor Hastie, and Rob Tibshirani. "Regularization Paths for Generalized Linear Models via Coordinate Descent". In: *Journal of Statistical Software* 33.1 (2010), pp. 1–22. ISSN: 1548-7660. URL: <http://www.jstatsoft.org/v33/i01>.
- [10] Hadley Wickham. *ggplot2: elegant graphics for data analysis*. Springer New York, 2009. ISBN: 978-0-387-98140-6. URL: <http://had.co.nz/ggplot2/book>.