

CSE8803 Project: Surgeon Scorecard

Michael Shamberger

Abstract—CSE8803 Project: Surgeon Scorecard.

This project attempts to reproduce the work done by Pro Publica data scientists that made the individual surgeon complication rate for eight low risk inpatient procedures in the United States public for the first time in 2015.

A configurable open algorithm for evaluating surgical outcomes is implemented on top of the OHDSI OMOP data format. This is used to calculate the raw complication rate for individual surgeons for the eight selected procedures.

PySpark and Python Pandas are the selected technologies for the implementation and the CMS-ETL was used for converting Medicare data to the OMOP format. Various defects and scalability issues are found and corrected in the data conversion from Medicare to OMOP format and also with the Spark tool itself.

The raw complication rates for surgeons are computed on the Medicare SynPuf data and the results were published to the web. The Surgeon Scorecard application was then open sourced as a starting point for discussion on developing an open algorithm for computing surgical complication rates.

Index Terms—Big data, Health analytics, Data mining, OHDSI, Surgeon Scorecard, Pro Publica, Open source

I. INTRODUCTION AND MOTIVATION

DATA scientists at Pro Publica performed a study on the US Medicare claims data to determine the complication rate of surgeons for eight low risk inpatient surgeries. They created an **application** that allows the general public to query the complication rates of local surgeons. This was the first time that surgeons in the United States had their complication rates published to the public.

The researchers at Pro Publica published their **methodology** but did not release their source code. Having an open source algorithm is necessary to allow researchers to modify and iterate on the algorithms as there is no public consensus in the medical profession exactly how to measure a surgeon's complication rates. An example of this is an online feud between Pro Publica and the Rand Corporation where the Rand corporation **criticizes** [8] the Scorecard methods, Pro Publica **responds** [4] and Rand then made a **statement** [7] on the Pro Publica rebuttal.

There is a rigorous methodology available for evaluating surgical outcomes. It is from the American College of Surgeons, and is called NSQIP. It is the "leading nationally validated, risk-adjusted, outcomes-based program to measure and improve the quality of surgical care in the private sector." [1] The NSQIP does not publish their algorithm and the results must remain confidential to those who use it.

The goal of this work is to reproduce the **Surgeon Scorecard** and publish the application code as a scalable open source application that uses the standard OHDSI OMOP data format. It is also a starting point to defining and implementing a configurable open algorithm for evaluating surgical outcomes.

II. METHODS

A. Data

The data used for analysis is the synthetically made **Medicare 2008-2010 Data Entrepreneurs SynPuf** converted to an OMOP format. The original Pro Publica study was performed against the Medicare 100% Standard Analytic Files for years 2009-2013.

The Data Entrepreneurs SynPuf contains the medical claims information for 2.32 million synthetic patients.

The SynPuf data preserves the detailed data structure of key variables at both the beneficiary and claim levels [2]. There is no guarantee that the data will produce meaningful results at the individual surgeon level.

| OMOP Table | N |
|----------------------|---------------|
| procedure_cost | 637,854,416 |
| observation_period | 2,098,516 |
| provider | 905,493 |
| measurement | 74,128,430 |
| observation | 37,531,052 |
| location | 3,186 |
| payer_plan_period | 7,789,513 |
| device_exposure | 4,465,486 |
| care_site | 320,546 |
| death | 107,645 |
| procedure_occurrence | 264,056,962 |
| drug_cost | 111,085,970 |
| person | 2,326,857 |
| visit_occurrence | 111,637,583 |
| condition_occurrence | 276,693,282 |
| drug_exposure | 126,048,051 |
| Total | 1,657,052,988 |

TABLE I
MEDICARE SYNPUF STUDY: RECORD COUNTS

Of the tables shown in Table I only the following are used for this analysis: procedure_occurrence, condition_occurrence, person, visit_occurrence and death.

B. Identifying Patient Cohort

The Cohort is made up of the 1.88 million patients that had a hospital visit during the period 2008-2010. No screening for existing comorbidities was applied to the Cohort.

| Patients | N |
|---------------------|-----------|
| All Patients | 2,326,857 |
| Patients with visit | 1,885,277 |

TABLE II
COHORT INFORMATION

C. Identifying Procedures for Analysis

Similar to the Pro Publica study, eight elective inpatient procedures were selected. These surgeries include hip replace-

ment, knee replacement, 3 types of spinal fusion, gall bladder removal, prostate removal, and prostate resections.

| ICD9 CODE | Procedure | N |
|-----------|--|--------|
| 51.23 | Laparoscopic cholecystectomy | 11,263 |
| 60.5 | Radical prostatectomy | 2,585 |
| 60.29 | Transurethral prostatectomy (TURP) | 2,746 |
| 81.02 | Cervical fusion of the anterior column, anterior technique | 2,634 |
| 81.07 | Lumbar and lumbosacral fusion of the posterior column, posterior technique | 11,122 |
| 81.08 | Lumbar and lumbosacral fusion of the anterior column, posterior technique | 10,622 |
| 81.51 | Total hip replacement | 10,903 |
| 81.54 | Total knee replacement | 29,013 |
| Total | | 80,888 |

TABLE III

MEDICARE SYNPUF STUDY: PROCEDURES INCLUDED IN ANALYSIS AND COUNT OF INDEX ADMISSIONS IN PATIENT COHORT

D. Identifying Index Admissions

An index procedure is defined as one of the individual surgeries for which the outcome was analyzed.

- Each procedure in Table III was identified by a set of qualifying concept primary admission icd9 codes.
- The procedure is checked to make sure it occurs during an inpatient visit. This allows the algorithm to avoid billing codes that might occur weeks after the surgery resulting in a duplicate index procedure being identified.

The following were part of the Pro Publica study but were not implemented:

- Excludes patients if the surgery occurred during an ER visit
- Excludes patients transferred from a nursing home or correctional facility.

E. Complications

Table IV shows examples of complications found from the patient cohort where the complication icd9 code is the primary admission code.

For each patient that was found to have an index procedure, the data was searched for two negative outcomes.

- Death - Patient died during hospital stay or within 30 days of being discharged.
- Complication - Patient was discharged alive but was admitted to a hospital within 30 days[11] of discharge with a principal diagnosis indicating a negative surgical outcome.

The index readmissions were detected by searching a set of icd9 primary diagnosis codes for codes that were indicating a complication based on the index procedure and within the time window.

For the full list of surgical codes and their counts that were determined to indicate surgical complications for the procedures of interest, see [Readmission Count Primary](#) online.

| Complication Type | N | Example |
|-------------------|--------|---|
| Infection | 2,834 | 998.59-Postoperative infection |
| Clot | 118 | 415.11-Iatrogenic pulmonary embolism |
| Reaction | 242 | 996.69 - Infection and inflammatory reaction due to internal joint prosthesis |
| Mechanical | 513 | 996.47 - Mechanical complication of prosthetic joint implant |
| Sepsis | 15,399 | 03.89 - Septicemia |
| Bone | 406 | 996.44 - Peri-prosthetic fracture around prosthetic joint |
| Hematoma | 574 | 998.12 - Hematoma complicating a procedure |
| Wound | 154 | 998.2 - Accidental puncture or laceration during a procedure |
| Hemorrhage | 567 | 998.11 - Hemorrhage complicating a procedure |
| Pain | 230 | 338.18 - Acute postoperative pain |
| Digestive | 0 | 997.49 - Digestive system complications |
| C.diff | 1,521 | 00.845 - Intestinal infection due to Clostridium difficile |
| Misc. Comp. | 35 | 787.22 - Dysphagia, oropharyngeal phase |
| Seroma | 113 | 998.13 - Seroma complicating a procedure |
| Fever | 49 | 780.62 - Postprocedural fever |
| Urinary | 192 | 997.5 - Surgical complications of the urinary tract |
| Total | 22,947 | |

TABLE IV

EXAMPLES OF 20 TYPES OF COMPLICATION ASSOCIATED WITH SURGERY

| ICD9 CODE | Procedure | N |
|-----------|--|-----|
| 51.23 | Laparoscopic cholecystectomy | 136 |
| 60.5 | Radical prostatectomy | 14 |
| 60.29 | Transurethral prostatectomy (TURP) | 12 |
| 81.02 | Cervical fusion of the anterior column, anterior technique | 21 |
| 81.07 | Lumbar and lumbosacral fusion of the posterior column, posterior technique | 97 |
| 81.08 | Lumbar and lumbosacral fusion of the anterior column, posterior technique | 91 |
| 81.51 | Total hip replacement | 139 |
| 81.54 | Total knee replacement | 383 |
| Total | | 893 |

TABLE V

MEDICARE SYNPUF STUDY: READMISSIONS DUE TO SURGICAL COMPLICATION FOR EACH OF THE PROCEDURES INCLUDED IN ANALYSIS, BROKEN DOWN BY PRINCIPAL DIAGNOSIS ON READMISSION.

F. Surgeon Complication Rate

Table VI shows the counts of index complications that followed an index procedure.

The complication rate for each surgeon was calculated by taking the total amount of index procedures and dividing those by the sum of instances of death and readmission for that procedure. While the Pro Publica study applied a risk adjustment based on comorbidities to each complication rate this study only produces the raw surgeon complication rate.

G. Further Information

More detailed information about the medical justifications for selecting these procedures and the icd9 codes that indicate them in the Medicare data can be found in the Pro Publica whitepaper [Surgeon Level Risk Methodology](#)[6].

| Procedure | N |
|-----------|--------|
| 51.23 | 7,389 |
| 60.5 | 2,103 |
| 60.29 | 2,248 |
| 81.02 | 2,189 |
| 81.07 | 6,964 |
| 81.08 | 6,719 |
| 81.51 | 6,648 |
| 81.54 | 14,041 |

TABLE VI
MEDICARE SYNPUF STUDY: NUMBER OF SURGEONS BY PROCEDURE

III. IMPLEMENTATION

A. Analytic Infrastructure

A dedicated server with 64 GB of ram and 8 cpu cores was used for the scalability testing of the PySpark application.

B. Technologies

The application needs to be capable of handling large data sets. The OHDSI OMOP format already has had 682 million patient records converted to it[9] while the The Medicare Fee-For-Service program included around 52.5 million patients in 2013. The smaller Synpuf dataset has 2.32 million patient records and when converted to the OMOP format had a total of 1.6 billion records as shown in Table I. Overall, it is estimated that only 1% of surgical outcomes are being measured[13] and this percentage will increase in the future due to advancements in electronic data capturing technology. Due to the scale of the data involved with the analysis it was decided to use big data technologies[12][14] for the implementation.

The following technologies were used to implement the Cohort and Readmission components as well as the Surgeon Scorecard application:

- Python - The OHDSI communities existing tools show that they are most comfortable with scripting languages such as R and sql logic but these existing tools have limitations when dealing with large data sizes[10]. To develop an application that would be based on similar technologies, python in PySpark with Spark Sql was used.
- Python Pandas - The Surgeon Scorecard sometimes sends processed data to a pandas dataframe. This is only done when it is known that the data can fit into the memory of a single server. There are more API's implemented for processing Dataframes in pandas then are currently available in PySpark allowing for rapid development.
- Spark - Spark is a framework that allows operations to utilize all cores and memory on a server. This is in contrast to a relational database such as Mysql which can only run 1 query on 1 cpu core. [3] In addition, Data load and tuning times of clusters are lower in a map reduce system.[15]

C. Software Components

- Cohort- A PySpark application class Cohort was developed. Based on a properties files, OMAP format data is read in from CSV files and a cohort of patients is selected.

The data can be read uncompressed or compressed in various formats. After the cohort is selected, the remaining data is filtered to remove data that is not related to this cohort. The option exists to write this cohort and filtered data back out to OMAP format CSV files. This smaller dataset can be loaded into a database in order to be able to use the full set of OHDSI tools.

- Readmissions- The PySpark application class Readmission takes a set of OMOP data and detects hospital inpatient readmissions following a procedure of interest. It can be used in combination with the Cohort class to find readmissions on a filtered cohort of users. The codes of procedures of interest and their possible complication codes are defined in a properties files diagnosis.properties and readmission.properties. There is no limit to the amount of procedures that can be investigated during a single run. This component can work with other codes besides icd9 such as icd10 and SNOMED as long as consistent codes are used in both the diagnosis.properties and readmission.properties as well as the data being analyzed. There is no functionality for conversion of codes.

D. Library Gaps

There are no existing open source Python libraries available for conversion of medical codes such as between icd9 and icd10 or making relationships between these codes and comorbidities. An R library exists for icd conversion and relation to comorbidities. It would be a good model for a future python library but converting it will be a significant development effort.

Due to this library gap, it was not possible to apply risk adjustment due to comorbidities which was done in the Pro Publica study.

E. Data Issues

There are some data conversion issues with the existing OHDSI CMS ETL that impacted the study. The Pro Publica study used primary admission diagnosis for identifying a qualifying procedure. Examination of the logic in the ETL showed that while the primary admission diagnosis was mismatched for outpatient to a first position diagnosis, it was completely dropped when doing an inpatient record migration. The mapping to the CONDITION_TYPE_CONCEPT_ID that holds the information on whether the diagnosis was primary or not was incomplete in both inpatient and outpatient record conversions.

| Concept Type ID | Description | N |
|-----------------|------------------------------|------------|
| 38000200 | Inpat Condition 1st Position | 8,317,475 |
| 38000251 | Inpat Procedure 1st Position | 3,592,580 |
| Total | | 11,910,055 |

TABLE VII
INITIAL INPATIENT COUNTS OF MISMATCHED CONCEPT_TYPE_ID FROM CMS-ETL

The source data defines a primary diagnosis code and several other sets of up to 10 other icd procedure codes. To

fix this defect, it was decided that the actual position in the Medicare source data was not considered medically important. Only the difference between primary and not primary was significant. The python ETL was modified and github [pull request](#) was sent to the OHDSI CMS-ETL project. The ETL was re-run and the scorecard was recalculated based on the correctly mapped data.

| Concept Type ID | Description | N |
|-----------------|----------------------------------|------------|
| 38000199 | Inpat Condition Primary Position | 959,691 |
| 38000200-209 | Inpat Condition Not Primary | 7,932,634 |
| DNE | Inpat Procedure Primary Position | 291,287 |
| 38000251-260 | Inpat Procedure Not Primary | 3,383,288 |
| Total | | 12,566,900 |

TABLE VIII
FINAL COUNTS OF CONCEPT_TYPE_ID FROM CMS-ETL

This made a difference as the code was now considering only 1.3 million condition and procedure occurrences instead of 11.8 million. The results would have been off by an order of magnitude.

F. Scalability Issues In Spark

Spark has some 2G limits in its 2.x codebase. These are due to 32 bit types used in various places in the code. The first of these happens when reading a data block that has been stored to disk. The code uses an instance of MappedByteBuffer which cannot exceed a size of 2G. The second is when using kryo serialized data. The serialized data is stored in a byte[] the size of which cannot exceed 2G. The third situation occurs when RPC writes data to be sent to a channel. The code uses a ByteBuffer which means that it cannot transfer data over 2G in memory. The fourth and final issue occurs when an RPC message is received where again a ByteBuffer is used which cannot exceed 2G. The result to the end user is the same for all of the above issues. Their big data application crashes with a “Size exceeds Integer.MAX_VALUE at sun.nio.ch.FileChannelImpl.map”.

The recommended workaround for this limitation is to set additional partitions beyond the default 200 on the dataframes. The properties file for the application allows the partition size to be set for all dataframes created. The workaround did not help with the Surgeon Scorecard application. Tests were done with up to 4000 partitions but instability was still seen as the study parameters were modified.

A high priority [defect](#) is open against the Spark application for this scalability issue. There is currently a [patch](#) pending to fix this issue but it has not yet been incorporated to a released version of spark.

In order for the scorecard to run reliably, the patch was compiled into the main branch and resulting binary was deployed to the implementation server. The patch solved the issue and this was reported back to the pull request in github. The patched source code is available in a [github repository](#) and the binary for ubuntu platform can be downloaded from [here](#).

IV. RESULTS

A. Surgeon Scorecard

A set of reusable usable software components were developed that can be used with the OHDSI OMOP data format. These components were used to reproduce the Surgeon Scorecard raw readmission rates.

The library successfully calculates the complication rate for each of the 8 inpatient surgical procedures of interest. The scorecard results for the 1000 patient dataset can be seen [here](#). The full data set scorecard results are available [here](#).

Similar to the Pro Publica [Surgeon Level Risk Appendix](#)[5], index procedure counts using icd9 codes of interest were generated and can be seen at this [location](#). The counts of index complication codes were calculated and available [here](#).

The goal of scalability was achieved as the tool ran the scorecard for the eight procedures for the 2.32 million patients in 100 minutes. It is expected that it would complete the scorecard for the full Medicare dataset of 52.5 million users in less than a day on a single server.

B. Comparison To Pro Publica Results

The Pro Publica study used the Medicare 100% Standard Analytic Files for years 2009-2013 and included those patients that were enrolled in the Medicare Fee-For-Service program. The Medicare Fee-For-Service program included around [52.5](#) million patients in 2013.

The analysis on the SynPuf dataset compares well to that done on the real Medicare Standard analytic files as seen in Table IX. The SynPuf adjusted figures are found by multiplying by a 22.6 scaling factor for data size and also 1.33 for the length of the study since the real Medicare data covered 4 years while the Synpuf data covered 3. Adjusted for data size the number of index procedures matches closely with the Pro Publica study while the index complications are off by factor of 2. The number and type of inpatient procedures was probably a key variable when generating the SynPuf data but a more complicated measure such as 30 day readmission rates after a procedure was not considered.

| Description | N Medicare | N SynPuf | N SynPuf Adjusted |
|---------------------|------------|----------|-------------------|
| Total Procedures | 2,376,851 | 80,888 | 2,437,425 |
| Total Complications | 64,367 | 893 | 26,909 |

TABLE IX
COMPARISON OF FULL MEDICARE DATA RESULTS VS SYNPUF DATASET

Without access to the real Medicare data and the source code written by Pro Publica for their Surgeon Scorecard it is not possible to determine the accuracy of their implementation. Pro Publica was also prevented from publishing their raw complications rates by the CMS organization that administers the data since reidentification may be possible when the amount of procedures and complications for a surgeon are low.

Pro Publica applied a complication risk adjustment for each surgeon based on patient comorbidities and age but this was not implemented due to the lack of open source software libraries in the selected implementation language.

C. Cohort Selection

Another benefit of this work is to allow the OHDSI tools to be used even when the data is initially too massive to be loaded into a relational database without requiring special hardware and skills. Researchers are able to build a Cohort on a subset of the data based on their specific criteria and generate csv format files that can be loaded into a relational database to make use of the rich set of OHDSI tools already in place.

D. Big Data Tool Improvements

The application encountered data size limitations in Spark and published a method to bypass them by patching the main branch of code.

The CMS-ETL produced the first publicly available massive health care dataset in the OMOP format. Defects in the ETL related to conversion of the primary admission diagnosis were uncovered and a fix provided to the OHDSI.

V. FINAL COMMENTS

An open algorithm for evaluating surgeon performance could have credibility across the healthcare industry if it was backed by an organization such as the OHDSI. It would be beneficial to healthcare data scientists if they would form a working group and publish an open algorithm for different hospital and physician evaluation metrics.

More comprehensive open source medical libraries need to be built into the OHDSI tools. The current OHDSI tools are built as applications but would benefit from having a set of reusable components implemented in a scripting language such as python. Python API's for medical code conversion (For example: icd9/10 to SNOMED) and relationships to comorbidities are also needed. The benefits of having a common schema to combine datasets are limited if it is not possible to have common medical codes within it.

The synthetic data provided by CMS is a good start and the open source Surgeon Scorecard application could not have been developed without it. There is a need for a refresh of the data since the Medicare coding has moved to icd10 codes in 2015 from icd9. An application will need to be able to handle both data with icd9 and icd10 codes as well as convert the icd9 codes to icd10 in order to perform multi-year studies.

VI. SUPPLEMENTAL MATERIAL

Source code - <https://github.com/opme/SurgeonScorecard/>

Video Presentation - <https://youtu.be/vRWRjVtP1Xo>

REFERENCES

- [1] Acs national surgical quality improvement program. <https://www.facs.org/quality-programs/acs-nsqip>. Accessed: 2016-11-23.
- [2] De 1.0 data users document appendix a. https://www.cms.gov/Research-Statistics-Data-and-Systems/Downloadable-Public-Use-Files/SynPUFs/Downloads/SynPUF_DUG.pdf. Accessed: 2016-11-23.
- [3] Using apache spark and mysql for data analysis. <https://www.percona.com/blog/2015/10/07/using-apache-spark-mysql-data-analysis/>. Accessed: 2016-10-24.
- [4] S. Engelberg and P. O. Our rebuttal to rands critique of surgeon scorecard. <https://www.propublica.org/article/our-rebuttal-to-rands-critique-of-surgeon-scorecard>.
- [5] S. Engelberg and P. O. Surgeon level risk appendix. <https://static.propublica.org/projects/patient-safety/methodology/surgeon-level-risk-appendices.pdf>.
- [6] S. Engelberg and P. O. Surgeon level risk methodology. <https://static.propublica.org/projects/patient-safety/methodology/surgeon-level-risk-methodology.pdf>.
- [7] M. W. Friedberg, K. Y. Bilimoria, P. J. Pronovost, S. D. M., C. L. Damberg, and A. M. Zaslavsky. Response to propublica's rebuttal of our critique of the surgeon scorecard. <http://www.rand.org/pubs/perspectives/PE170z1.html> Santa Monica, CA: RAND Corporation, 2015.
- [8] M. W. Friedberg, P. J. Pronovost, P. J. Bilimoria, S. D. M., C. L. Damberg, and A. M. Zaslavsky. A methodological critique of the propublica surgeon scorecard. http://www.rand.org/pubs/perspectives/PE100/PE170/RAND_PE170.pdf Santa Monica, CA: RAND Corporation, 2015.
- [9] H. G. R. PB, and D. JD. Characterizing treatment pathways at scale using the ohdsi network. *Proceedings of the National Academy of Sciences of the United States of America*. 2016;113(27):7329-7336.
- [10] M. Kane, J. W. Emerson, and S. Weston. Scalable strategies for computing with massive data. *Journal of Statistical Software*, 55(1):1-19, 2013.
- [11] H. M. Krumholz, Z. Lin, E. E. Drye, M. M. Desai, L. F. Han, M. T. Rapp, and S. L. T. Normand. An administrative claims measure suitable for profiling hospital performance based on 30-day all-cause readmission rates among patients with acute myocardial infarction. *Cardiovascular Quality and Outcomes* 4(2), 243-252.
- [12] S. Landset, T. M. Khoshgoftaar, A. N. Richter, and T. Hasanin. A survey of open source tools for machine learning with big data in the hadoop ecosystem. *Journal of Big Data* 2015 2:24.
- [13] H. Lyu, M. Cooper, K. Patel, M. Daniel, and M. Makary. Prevalence and data transparency of national clinical registries in the united states. *Journal For Healthcare Quality*.
- [14] C. Ordonez. Can we analyze big data inside a dbms? *Proceedings of the sixteenth international workshop on Data warehousing and OLAP*, pp. 85-92, October 2013.
- [15] E. R. e. a. Pavlo, A. Paulson. Comparison of approaches to large-scale data analysis. In *Proceedings of the 2009 ACM SIGMOD International Conference on Management of data*.