

LISBON
DATASCIENCE
ACADEMY

Reducing wrongful discharge

Predicting wrongful discharge at
the Hazel and Bazel Hospital (Report I)

Prepared for:

The Hazel and Bazel Hospital

Prepared by:

Yash Pandya

Lisbon Data Science's *Awkward Problem Solutions™*

Table Of Contents

Table Of Contents	2
Client requirements	3
1.1 Summary	3
1.2 Requirements clarifications	3
Dataset analysis	4
2.1 General analysis	4
2.2 Business questions analysis	6
2.3 Conclusions and Recommendations	7
Modelling	8
3.1 Model expected outcomes overview	8
3.2 Model specifications	8
3.3 Analysis of expected outcomes based on training set	9
3.4 Alternatives considered	10
3.5 Known issues and risks	10
Model Deployment	12
4.1 Deployment specifications	12
4.2 Known issues and risks	13
Annexes	15
5.1 Dataset technical analysis	15
5.2 Business questions technical support	20
5.3 Model technical analysis	24

1. Client requirements

1.1 Summary

The client, The Hazel and Bazel Hospital (HBH), has had multiple instances of patients being wrongfully medically discharged, resulting in severe health consequences for the patient. There have been allegations made that these wrongful discharges may be targeting specific subsets of the population based on gender, ethnicity, race and insurance status.

HBH management has deemed it necessary to carry out an investigation into these allegations, and implement an additional service into HBH's hospital patient record management system. Specifically, this additional service will be integrated into HBH's internal system, and will use available patient data to predict the likelihood of the patient being prematurely discharged. To aid with the investigation and implementation, HBH has made available $\approx 80,000$ partially anonymized patient records dating back to 2012.

As such, the primary objectives of this project are to:

1. Investigate the provided data and determine if there is evidence of discrimination based on gender, ethnicity or age when dismissing patients from care.
2. Determine if the Medical Specialty the patient is under the care of, the admission source, or their insurance status is a potential source of this discrimination
3. To identify patients who are likely to return to the hospital in less than 30 days

1.2 Requirements clarifications

Additional clarifications were requested from Dr Agnes Crumplebottom. In addition to some clarifications regarding input/imputation irregularities, identification of missing data and insurance status, Dr Crumplebottom provided the following key information:

1. A clearer target was outlined, specifically that wrongful discharges (leading to readmission within 30 days) should be minimized. This can be achieved with a high recall model, allowing for detection of most, if not all of the sick patients. However, a further requirement that at least "50% of the patients readmitted should actually be sick" has been added. This is particularly challenging as the dataset provided is quite unbalanced and machine learning models may not allow us to achieve this target of 50% precision while making a meaningful reduction in wrongful discharges. To achieve a balance between these two metrics (precision and recall), F1 and ROC-AUC scores will also be used for model optimization.
2. Discrepancies between subgroups should be minimized (10% maximum difference), although different medical specialties may still have differences (5% maximum difference) as they are linked to conditions requiring frequent hospital visits.
3. Initial correspondence stated that this data covered the years since 2012, and may have had alterations over time, however, no time stamp data was available. It is reasonable to assume that admissions are registered sequentially and while we will attempt to use this property, it will limit analysis and modeling.
4. The IT team at HBH provided an expected structure of the REST API for integration into their system, allowing for the implementation of measures to ensure it remains robust. Considering the nature of the data, and the fact that a lot of it seems to be manually entered, there will be a requirement for an extended testing phase, and we encourage the implementation of improved record keeping practices to ensure stability and reproducibility in the future.

2. Dataset analysis

2.1 General analysis

An initial analysis of the provided data was carried out using the field descriptions provided by Dr Crumplebottom to examine structure, distribution and quality of the data available. Some key aspects of the data are highlighted here. See annex 5.1 for additional description of the provided dataset, information about data types (information about processed data is in the table), and links to interactive profile reports generated based on initial and processed data.

The data contains unique admission and patient identifiers (although multiple patients appear several times throughout the data - as expected). Patient demographic data (summarized on the right), namely gender, age and race in addition to patient insurance status are included. The data shows slightly more female than male patients were admitted to the hospital during the relevant time frame (Fig. 1A), and that the hospital caters primarily to older patients (Fig. 1B). The race category has been included with redundant descriptors (such as 'white' and 'caucasian' being separate categories), although a rapid processing reveals that the patients are primarily white, with a sizable black minority and very few patients of other races (Fig. 1C). While multiple payer codes are provided in the data set (see annex 5.1 reports for additional information), our primary interest is in the 'SP' code for 'self payers' (uninsured patients), although a large number of patients have an unknown insurance status (Fig. 1D).

The dataset includes information pertaining to the medical reasons for the patient being admitted to the hospital and some information relating to their medical history and hospital stay (summarized overleaf and in annex 5.1). Unfortunately, weight data is largely missing for most patients. Although the dataset provides for the use of 8 distinct codes for the admission type and 26 for the admission source, these can be condensed and reveal that the majority of patients admitted arrive for emergency care (Fig. 2A and B).

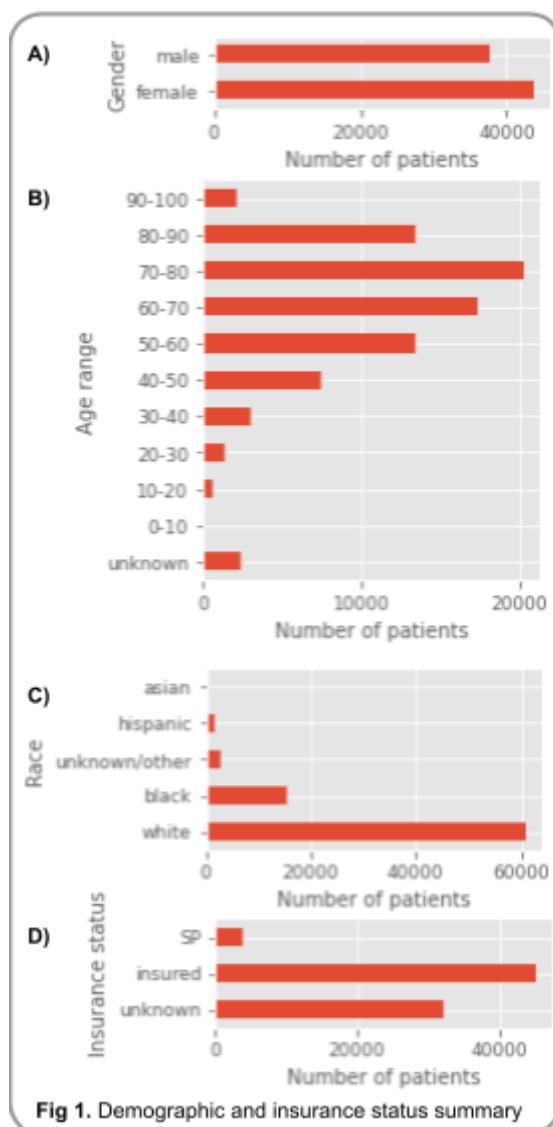


Fig 1. Demographic and insurance status summary

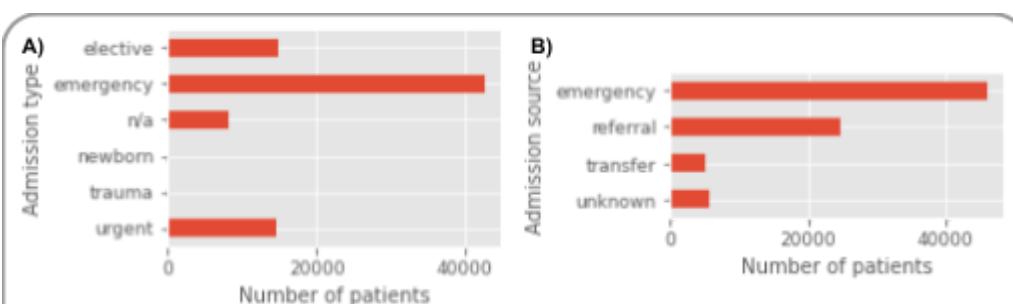


Fig 2) Admission type and source

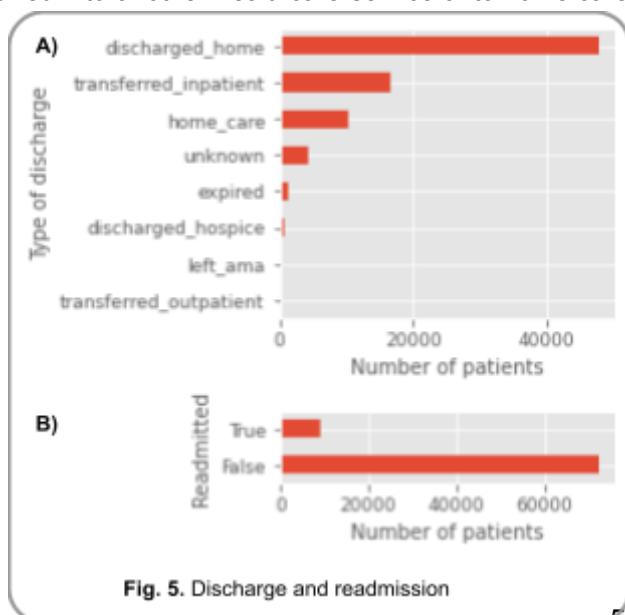
Clinical information provided (Annex 5.1 Fig 1) includes the specialty they are visiting (the top 5 specialties visited are shown - of note is that for many patients, the specialty is unknown) and the length of their hospital stay, which shows that the most common length of stay is 2-3 days, while very few patients stay for 14 days. In the dataset provided longer stays are not found, probably due to these patients being transferred to long term care facilities. Surprisingly for a hospital specializing in the treatment of diabetes, the hospital specialties with the highest number of admissions are not endocrinology, but include general practice, internal medicine and cardiology. The number of previous hospital visits by the patient (as outpatient, inpatient or emergency visits during the previous year are included, and the strong negative skews show that the vast majority of patients have very few, or no hospital visits per year, but that there are some outliers who are regularly at the hospital.

A large number of patients do not undergo medical procedures during their visit. There also appears to be a multimodal distribution regarding the number of laboratory procedures (i.e. tests) which the patients carry out, with a large number taking very few or no tests, and another peak appearing in patients who take approximately 40 tests. Finally, the amount of medications which were administered during the patients stay indicates a high degree of polypharmacy, which is explained by the fact that a large portion of the patients are elderly, and may require multiple medications on a regular basis, in addition to those being used to treat their current illness. Several clinical features were strongly linked (such as length of stay and number of procedures undergone), suggesting several features are strongly linked (see annex 5.1 Fig 2 for complete correlation data of numerical variables)

The diagnosis, as well as 2 secondary diagnoses are provided using the ICD9 codes. These fields have a high cardinality, reflecting the diversity in illnesses which patients can present with. Interestingly, the most common diagnostic codes as a primary diagnosis appear to be related primarily with circulatory, respiratory and digestive disorders, although there are endocrine disorders (including diabetes) amongst the most common secondary diagnoses.

Additional characteristics include whether or not the patient has a prosthesis or is on diuretics (both of which were a very small portion of patients), patient blood type, whether or not the patient required a blood transfusion, including over 10% of patients and whether or not the patient was prescribed insulin, or new medication for diabetes management, resulting in a change in their current medication. These last features, shown in annex 5.1 Fig 3 reveal that the vast majority of patients at the Hazel and Bazel hospital do have diabetes and that they are frequently prescribed new medications, with an often changing regimen.

The final aspects of the database pertain to the patient's departure from the hospital. 29 discharge codes are available, although several are very rarely used and others are unmapped. The codes can consequently be compressed into a few categories, showing that the vast majority of patients are discharged home, while many others are transferred into another healthcare service or to home care services (Fig. 5A). The final aspect of the dataset is the one which will be key for this project: readmission. This field informs whether a patient has been readmitted within 30 days of their discharge and can act as an indicator of a wrongful discharge. In the United States, in 2018, readmission rates were approximately 14% ([AHRQ statistical brief #278, 2021](#)), placing HBH well below average with approximately 10% (Fig. 5B), although with readmissions costing a large sum, and disproportionately impacting patients with diabetes, it is important to minimize this value as much as possible.



2.2 Business questions analysis

The objectives of this investigation were to i) establish whether or not there was any discrimination based on patient gender, race, age or insurance status resulting in wrongful discharge, and to ii) determine if any admission sources or services were discriminating based on these 4 sensitive features. The technical specifications provided indicate discrepancies of up to 10% between subgroups, and up to 5% between medical specialties are acceptable for the API predictions. For this analysis however, a more stringent 5% difference was used as an indicator for areas requiring further investigation. Readmission rates of the subgroups were compared with expected readmission rates (the rate if each group was equally represented). Additionally, the difference between the lowest and highest values was extracted. In the event this difference was beyond the specified threshold of 5%, it was considered that discrimination may be occurring.

Using this approach, There do not seem to be differences between different subgroups based on gender (Fig. 6 A) or race (Fig. 6 B), as shown by the close proximity of all the readmission rates to the expected readmission rate (blue line), and the fact that the values are all within 5% of the minimum value detected (indicated by window between dashed green lines). It should be noted that, for age subgroups (Fig 6 C), while children have lower readmission rates, they are often treated by separate medical units and so the age groups up to 20 have not been used to determine minimum readmission rates. In any case, there appears to be a mild increase in readmissions in the 20-30 year old age range (less than 5% above minimum), warranting further investigation. Lastly, in the case of patient insurance status (Fig 6 D), there is no clear discrepancy between insured and uninsured patients.

The admission sources were grouped as above (emergency, referral, transfer and unknown/other) based on the field descriptions provided and total readmission rates for each source were examined (Annex 5.2 Fig 1). Subsequently records from each admission source were selected and the readmission rates pertaining to individuals based on race, gender, age and insurance status were

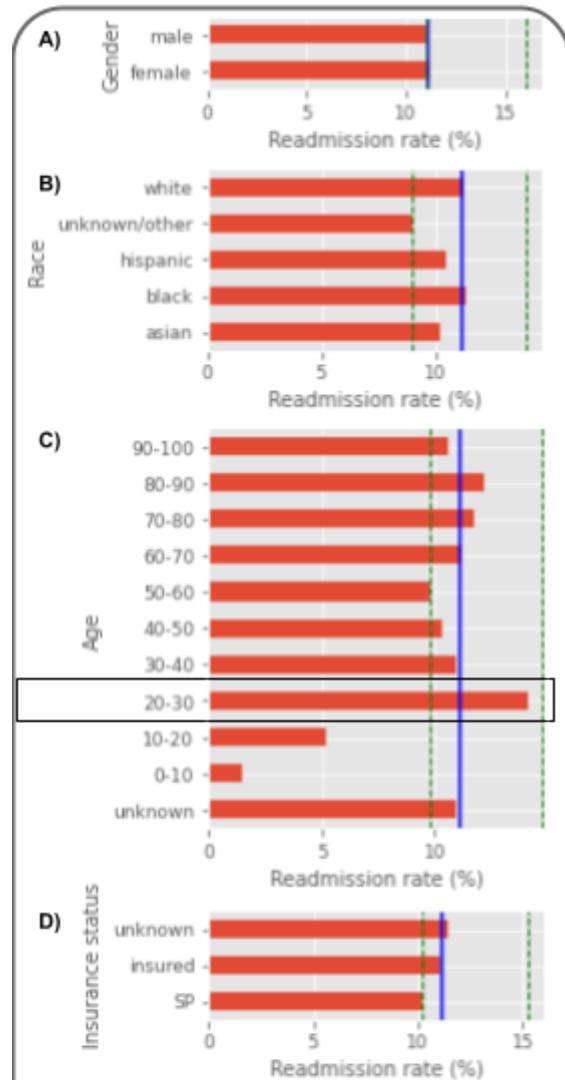


Fig 6. Readmission rates and demographic data

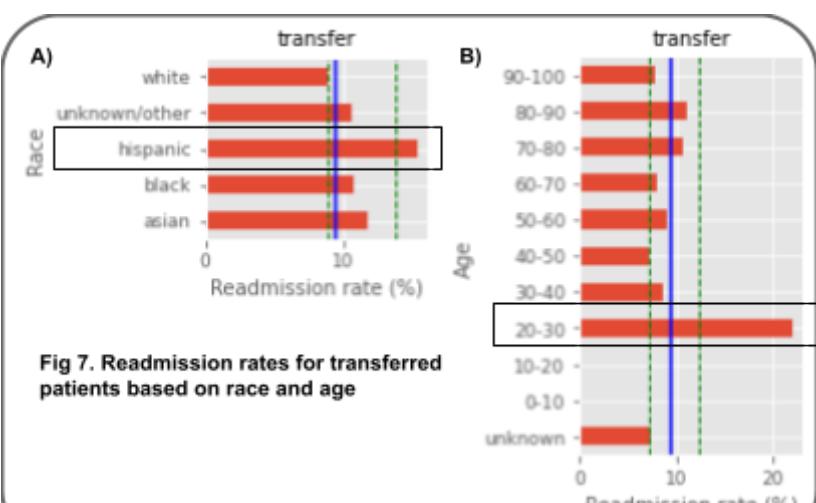


Fig 7. Readmission rates for transferred patients based on race and age

examined. Of particular interest, it appears that, among patients who are transferred from other healthcare services, hispanic patients have higher readmission rates (Fig 7 A), as do in all groups 20-30 year olds (Fig 7 B). The latter group also have higher readmission rates in all admission sources except referrals (Annex 5.2 Fig 2).

To examine the medical specialties, the data was first filtered to exclude specialties which had less than 100 entries and to exclude records where the medical specialty was unknown. Once processed, this left 24 medical specialties with a varying range of readmission rates.

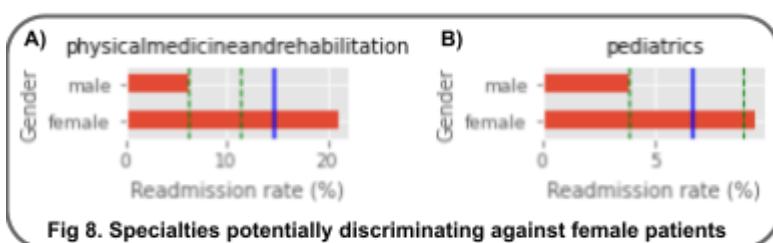


Fig 8. Specialties potentially discriminating against female patients

This analysis reveals that hematology/oncology, oncology, nephrology vascular surgery and rehabilitative medicine most commonly result in readmissions, with the first 3 having readmission rates 10% higher than the minimum in other non-age specific specialties (Annex 5.2 Fig 1). A similar approach to that for admission sources was then used with records for each medical specialty being selected and readmission rates linked to race, gender, age and insurance status was examined. Of note, both physical medicine and rehabilitation, and pediatrics are much more likely to readmit female patients (Fig 8 A and B). While there are some differences between readmission rates in specialties linked to race (Annex 5.2 Fig 4), it is difficult to draw valid conclusions here as races other than black and white have low representation (although hispanic patients frequently have higher readmission rates). Multiple specialties show increased readmission rates among young adults, reflecting the overall behavior of the dataset (Annex 5.2 Fig 5) and suggesting that this problem is not linked to specific specialties. Discrepancies which impact older populations appear in specialties linked to age related diseases, such as psychiatry and physical medicine and rehabilitation. Alarmingly, SP (self paying patients) have higher readmission rates in several specialties including include cardiothoracic surgery, gastroenterology, urology and vascular surgery (Annex 5.2 Fig 6).

2.3 Conclusions and Recommendations

HBH has a readmission rate which is slightly lower than published average values. Amongst the admission sources for the patients, there is no clear difference in readmission values, although a key role is played here by the impact of a large fraction of patients arriving for emergencies, potentially masking the effects of other admission sources.

A similar problem is presented with respect to race, as only white and black patients are present in sufficient numbers to draw conclusions, with other minorities accounting for a small fraction of patients. In any case, it is clear that HBH needs to further investigate the causes behind increased readmission of minorities in pediatrics, several surgery related specialties and hematology/oncology, where black patients are twice as likely to be readmitted as white patients.

Furthermore, it is apparent from the findings that young adults arriving at the hospital for emergency care, or being transferred from other facilities have a high risk of readmission compared to other age groups, a difference which is reflected when considering all the patients as a whole. It is important to ensure that these differences are not the result of young adults being discharged prematurely.

Lastly, it would be advisable for HBH management to engage in continuous monitoring of this data, and improve the record keeping system such that the large number of “unknown” or missing values are filled. This would allow for a clearer analysis and an easier identification of potential sources of wrongful discharge by allowing for more accurate assessment of patients prior to discharge. The missing status of patient health data is a potential indicator of quality of care, whereas the missing information linked to patient insurance status may influence accounting.

3. Modelling

3.1 Model expected outcomes overview

The initial objective was to design and implement a system which would predict patient readmission within 30 days. The target is to ensure that at least 50% of readmissions are patients who would truly require further care, with variation of no more than 10% between race, gender and age groups, and no more than 5% between different medical specialties.

The dataset provided was unbalanced (11% readmission rate) and consequently it proved to be challenging to meet the requirements desired. In particular, and as explained in section 3.3, in order to reach the 50% precision requested, the threshold at which patients would be readmitted would be set very high (at 0.67). Unfortunately, this reduces the recall score to 0.014, meaning less than 2% of wrongful discharges would be correctly identified. Consequently, it was deemed a necessary tradeoff to reduce precision and aim for a balance which would allow for the allowance of more sick patients. However, the deployed model has a precision of 18.6% and this means a significant amount of patients (81.4%) predicted wrongful discharges would be incorrectly classified as such.

It is expected that the model will adhere to the requirements for avoiding discrimination with respect to gender (maximum difference = 1.0%), admission sources (maximum difference = 2%) or patient insurance status (uninsured vs insured difference = 4%). With respect to race, the current model may have moderate discrimination against some groups (maximum difference = 10.0%), and asians and hispanic patients, although poorly represented in the dataset may be discriminated against. This model will need to be improved if these patient groups increase in proportion. Additionally, the data reflects a clear bias for young adults being prematurely discharged, and our system also predicts that they are more likely to be (maximum difference of 37.5%, with 20-30 year olds much higher than next closest age group). Lastly, certain specialties are more likely to result in readmission, however, these specialties are linked to frequent patient visits such as rehabilitative medicine and oncology, consequently, this aspect is not likely to be particularly problematic (A summary of these results is in Annex 5.3).

3.2 Model specifications

The implemented model consists of a pipeline culminating with a Random Forest Classifier, with several preprocessing steps on the data, summarized as follows:

1. The target for prediction was a Boolean value built based on the 'readmitted' value, with values of 'Yes' being replaced by True, and 'No' by False.
2. Data was sorted by admission_id (assumed to be sequential) and a test/train split, keeping 10% of entries for the test set was done. Random oversampling of the true class was used to generate a balanced training dataset.
3. Numerical variables ('time_in_hospital', 'num_lab_procedures', 'num_procedures', 'num_medications', 'number_outpatient', 'number_emergency', 'number_inpatient', 'number_diagnoses', 'hemoglobin_level'):
 - a. Missing values are imputed with the median of the column
 - b. Values are scaled using a robust scaling method to reduce the impact of outliers (minimum-maximum scaling gave similar results)
4. Binary variables ('has_prosthesis', 'blood_transfusion', 'diuretics', 'insulin', 'change', 'diabetesMed'):
 - a. Are encoded as True/False boolean values
5. Ordinal categorical variables ('max_glu_serum', 'A1Cresult', 'age', 'weight')

- a. These were converted to categories based on those available, and arranged in their natural order
 - b. The order was then used to encode a value based on their position in the order (e.g. for age 0-10, the value was 1, for 10-20, it was 2 etc).
 - c. Values which were not included in the ordered categories (such as missing data '?' values) were encoded as -1.
6. A feature selection step was added at this stage to allow for easy iteration and optimization - can be excluded if needed (see annex 5.3 for details of selected features and justification).
7. Categorical variables were each processed individually and one hot encoded prior to fitting the classifier (see annex 5.3 for details):
8. A Random Forest Classifier, selected to minimize risk of overfitting, was fitted after a grid search was used to select hyperparameters. The options tested in the grid search are listed in the table to the right. Optimal hyperparameters are in the final column and were selected based on F1 score. Unlisted parameters were used at their default values to economize on the time required to train the models (min_samples_leaf=1, min_weight_fraction_leaf=1, max_leaf_nodes=None, min_impurity_decrease=0.0, oob_score=False, verbose=0, ccp_alpha=0, max_samples=None, for definitions of hyperparamters, see documentation [here](#)).
9. The threshold used to decide predictions was maintained at the default value (0.5) as calculation of geometric means yielded an optimal threshold that was extremely close to this value (0.49, see section 3.3 for further details).

Hyperparameter	Options	Optimal value
Sample bootstrapping	True, False	False
class_weight	balanced, balanced_subsample	balanced
criterion	gini, entropy	gini
max_depth	2, 5, 10	10
max_features	auto (same as sqrt), log2	auto
min_samples_split	2, 5, 10	2
n_estimators	10 , 100, 1000	1000

The model generated was evaluated based on a combination of F1 score, precision scores and ROC-AUC scores and the differences in precision score between categories within the identified target variables. Specialties with at least 100 patients admitted were selected - although low, some of these had very high readmission rates, perhaps indicating a link to treatment of severe/rare diseases and consequently acting as good indicators for patients who may require repeated visits.

3.3 Analysis of expected outcomes based on training set

Based on the implemented model and supplied training set, the model has a precision rate of just under 20% (20% of patients predicted for readmission are actually readmitted) and a recall rate of slightly over 60% (correctly identifies more than 60% of the cases of potential wrongful discharge) (See annex 5.3 Fig 2). Calculation of geometric means to optimize the threshold for true positive rates did not provide a difference from the default value (0.49 was the indicator received, whereas default values are 0.5), and so the value was left as is (5.3 Fig 2).

The model also has an accuracy (the proportion of correct predictions) of slightly under 70%, and the ROC curve indicates a significant tradeoff between precision and recall (annex 5.3 Fig 3). In order to achieve the target of 50% of true predictions being correct, the threshold for a patient being predicted as readmitted would be so high as to mean less than 5% of cases of wrongful discharge would be identified. This would not solve the problem facing management, consequently, selecting a higher recall rate allows for more patients to be identified. Considering most patients are insured, identifying

them for additional care would not be a significant financial burden, whereas wrongful discharge lawsuits can result in heavy financial and reputational damage.

Comparison of precision scores generated for separate subgroups reveals that this model does not result in a large discrepancy between different genders (Male: 19%, Female: 18%), insurance status (Insured: 18%, Uninsured: 16%), admission sources (all values ranging between 18-21%) or races (White: 18%, Black: 20%, other races had higher scores but were present in low numbers), but it does still disproportionately affect a few specialities (The maximum difference was of 28%). Most concerning is the impact on 20-30 year olds who, as expected based on exploratory data analysis, had significantly higher predicted readmission rates which we were unable to mitigate without significantly compromising performance in other areas (A summary of methodology and results is available in annex 5.3).

3.4 Alternatives considered

Alternative approaches were explored during the initial stages of the project. The options considered and the reasons for them eventually being discarded are briefly explained below:

1. Similar results were achieved with random undersampling, and with alternative classifiers such as BalancedRandomForest and RUSBoost which implement random undersampling methodology. These were discarded in favor of random oversampling in order to ensure that information was not lost from the majority class.
2. Decision tree classifier:
This model was examined to allow for a high degree of interpretability, potentially granting additional insights, however, it had poorer performance when implemented as part of the pipeline used, quite possibly due to a high degree of overfitting to the training set.
3. KNN classifier:
An easily implemented, rapid algorithm which also allows for a high degree of interpretability. While initial testing with a small subset of the data was promising, when transitioning to the full data set, the primary flaw of this algorithm - struggles with larger datasets became very apparent and results suffered.
4. Gradient boosting classifiers:
The results were poor in the absence of additional approaches to handle unbalanced datasets. It also proved easier to tune the parameters for a random forest classifier, which gave similar results independent of being trained on over/undersampled data or the true dataset, consequently the gradient boosting options, including those with random undersampling at each iteration, were discarded.

3.5 Known issues and risks

The implementation applied here has several potential issues and risks, affecting the way features are processed, and the modeling approach selected.

Regarding feature engineering and selection steps, the compression of various categories into smaller groups (i.e, turning many categories into a few), while helping to reduce dimensionality of the data, may have resulted in the loss of resolution and, possibly, a loss of signal linked to specific categorical variables. Furthermore, there are over 10,000 ICD-9 diagnostic codes, and over 700 of these were present in the data. The approach employed means that there is room for any diagnostic code to be included in the broader categories, this approach has resulted in a loss of signal and resolution. A potential option to explore would be the imputation of ICD9 diagnoses using the codes, and treatment of the text using vectorization approaches and principal component analysis to minimize dimensionality.

It is important to note that sensitive features specified for investigation, gender, race, age, insurance status and admission source were included within the training of the model. While this may allow for groups facing unfair treatment to be identified, ideally, these would not be factors involved in determining patient suitability for discharge.

Using the Random Forest Classifier approach, minimizes overfitting, at the cost of interpretability. The random selection of features at each step makes it challenging to clearly identify important features and the reduction in variance afforded by averaging multiple predictions also comes at the cost of increased bias. This approach also takes significantly longer to train, a potential issue if the dataset expands much more, or if it becomes necessary to retrain frequently.

Lastly, changes in the data over time (data drift), such as alterations in the demographics of patients coming to the hospital, or changes in the quality of care due to changes in methods, staff, or facilities for example may result in impaired model performance. Additionally, as the definition of wrongful discharge is refined (concept drift), potentially accounting for different specialties or links between the medical reasons for initial admissions and readmissions, it is likely that the model will no longer be as useful.

4. Model Deployment

4.1 Deployment specifications

The developed model has been deployed behind an HTTP server on heroku following the steps outlined below:

1. Model serialization and deserialization

The developed pipeline was serialized using the python pickle library (part of python core) to serialize columns and datatypes, and the joblib library to serialize the pipeline. In the app, the model is deserialized for use using the same libraries (see code for details)

2. HTTP framework selection and database creation

Flask (<http://flask.pocoo.org/>) is the HTTP micro-framework selected for this deployment primarily because this project requires only two endpoints. Heroku postgresql was the database used for data storage in this implementation via the peewee library. The database allows for storage of the admission ID (must be a unique field), as well as the data received, the prediction returned and the true result (can be null).

3. Endpoints

1. **/predict:** The predict endpoint allows for the user to input data regarding the admission in question (including patient ID, diagnosis, discharge disposition - see the schematic overleaf for structure of the expected payload), and generate a prediction for whether that patient will be readmitted to the hospital within the next 30 days. This endpoint goes through the following steps:
 - a. Convert the request into a dictionary for further processing
 - b. Call validation functions that inspect the resulting dictionary to ensure:
 - i. The presence of all the required fields (even if they are null or empty)
 - ii. The presence of numerical admission and patient IDs
 - iii. That datatypes of columns are, or can be converted to, the expected datatype for use in the model (particularly important for numerical fields such as 'num_medication', 'num_procedures' etc)
 - iv. Checking that strictly defined categorical/string fields ('blood type', 'change', 'insulin', 'diuretics') contain an expected value.
 - c. Failure to pass any of these tests results in the return of a 422 error, and a message explaining the source.
 - c. Return the prediction ("Yes" or "No") and store the data. If the admission ID already exists, it will return an error informing that the data was not stored.
2. **/update:** This endpoint allows for true labels (the actual readmission status) to be added to the database, alongside previously made predictions and details about the admission. This can be used for analysis and potential retraining of the model if required. This endpoint also checks inputs for validity, to ensure true labels correspond to an admission for which a prediction has been generated, and stored in the database.

4. Application structure

The application consists of the pipeline, data types and column files outputted from pickling (explained above), in addition to the following files and components:

1. Heroku.yml and dockerfiles which allow for the implementation on the heroku platform
2. A requirements.txt file which allows for replication on other systems if required
3. 2 python files:

- a. The app.py file which contains all the code pertaining to database setup, model unpickling and implementation of endpoints
 - b. A unit test file which contains a variety of functions designed to examine the request input validity - these functions are called at various points in the app.py file
4. A custom_transformer package containing pipeline elements necessary for the model. These elements are all based on elements available from the sci-kit learn library:
- a. A preprocessor file (preprocessor.py) containing functions which clean and treat the data, as well as some minor processing of categorical variables
 - b. A combination imputer/scaler for numerical features (custom_impute_scale.py)
 - c. A one hot encoder (custom_onehot_encoder.py)
 - d. An ordinal encoder to process features which have a natural order in their categories such as age, containing the mapping for the relevant features from this dataset (custom_ordinal_encoder.py)
 - e. A feature selection element (featureselector.py) which takes a list of features to be included.

4.2 Known issues and risks

There are several potential weaknesses in the current deployment, affecting various areas, these are briefly outlined below:

1. The preprocessing methodology makes use of unknown and new data, by assigning many categorical values to the ‘unknown’ group. Consequently, the system is not protected from unusual entries, such as “race”:“blue”. This could be addressed by implementing improved recording systems - such as the use of drop down lists with well defined categories, or by increasing stringency on data validity. This can be altered as required.
2. The functions to look for data validity have been tested thoroughly, however, it is impossible to explore all possibilities and this means that the following issues may arise:
 - a. The system is currently set up to accept unique admission id values. While it will generate a prediction for duplicate values, it will not store these values, so imputation errors can result in data loss
 - b. Inability to generate a prediction with what would be considered a valid request
 - c. Predictions being generated for what should be considered an invalid request
3. The implemented transformers would require significant reworking if new features were to become available, and in their current state are dependent on being run in a very specific order.
4. The API has been deployed using Heroku, which may bring several limitations in the future:
 - a. Running a high traffic API on Heroku can become expensive very quickly, as well as potentially impacting performance
 - b. As Heroku provides the majority of infrastructure, there is the possibility of becoming “locked in”, as transitioning to alternative platforms becomes difficult and costly
 - c. The heroku database currently being used has a limit of 10,000 entries - so data will periodically need to be exported and the database reset. Failure to do so may result in data being lost or the app crashing.
5. Lastly, considering the sensitive nature of the data being stored, the deployed API has a major security flaw: weak access control (anybody with the URL can make requests in this

```
Prediction payload:
{
  "admission_id": 0,
  "patient_id": 0,
  "race": "string",
  "gender": "string",
  "age": "string",
  "weight": "string",
  "admission_type_code": 0,
  "discharge_disposition_code": 0,
  "admission_source_code": 0,
  "time_in_hospital": 0,
  "payer_code": "string",
  "medical_specialty": "string",
  "has_prosthesis": true,
  "complete_vaccination_status": "string",
  "num_lab_procedures": 0,
  "num_procedures": 0,
  "num_medications": 0,
  "number_outpatient": 0,
  "number_emergency": 0,
  "number_inpatient": 0,
  "diag_1": "string",
  "diag_2": "string",
  "diag_3": "string",
  "number_diagnoses": 0,
  "blood_type": "string",
  "hemoglobin_level": 0,
  "blood_transfusion": true,
  "max_glu_serum": "string",
  "A1Cresult": "string",
  "diuretics": "string",
  "insulin": "string",
  "change": "string",
  "diabetesMed": "string"
}

Update payload:
{
  "admission_id": 0,
  "readmitted": "string"
}
```

case). This leaves the API vulnerable to attacks such as SQL Injection, resulting in access to sensitive data and Distributed Denial of Service (rendering the endpoints unusable with excess traffic). The integration of this API with the hospital IT systems could therefore introduce new security weaknesses into these systems.

5. Annexes

5.1 Dataset technical analysis

Prior to any manipulation of the dataset, a profile report was generated from the provided file. The software details shown below will allow for a rapid replication.

Reproduction

Analysis started	2022-02-01 17:22:29.832680
Analysis finished	2022-02-01 17:23:27.469446
Duration	57.64 seconds
Software version	pandas-profiling v3.1.0
Download configuration	config.json

An interactive html report is available and can be downloaded at:

<https://drive.google.com/file/d/1SSookY72Dl4wggDs2saegGhA5VKpO1le/view?usp=sharing>

The file should be saved and opened using any web browser and all files will be available until at least April 30th 2022.

Subsequently, minor cleaning and processing of the data was performed namely:

1. Assignment of data types to columns
 - a. In categorical features containing missing or '?', an additional category was created titled 'unknown'. This decision was made considering that failure to collect patient data within certain specialties etc may reflect poor quality of care and assist with prediction of premature discharge.
2. Assignment (and compression) of categories a rapid technical summary of the data was generated (results below for each field provided)
 - a. Redundant categories within race were combined
 - b. Admission and discharge codes were combined to form larger categories - see model specifications and code for details on exact combination methodology.
 - c. Diagnoses were grouped by their primary category according to CDC classifications - see model specifications for these categories. This resulted in over 700 unique diagnoses being compressed into 17 categories.

A profile report of the processed data was also generated and is available at:

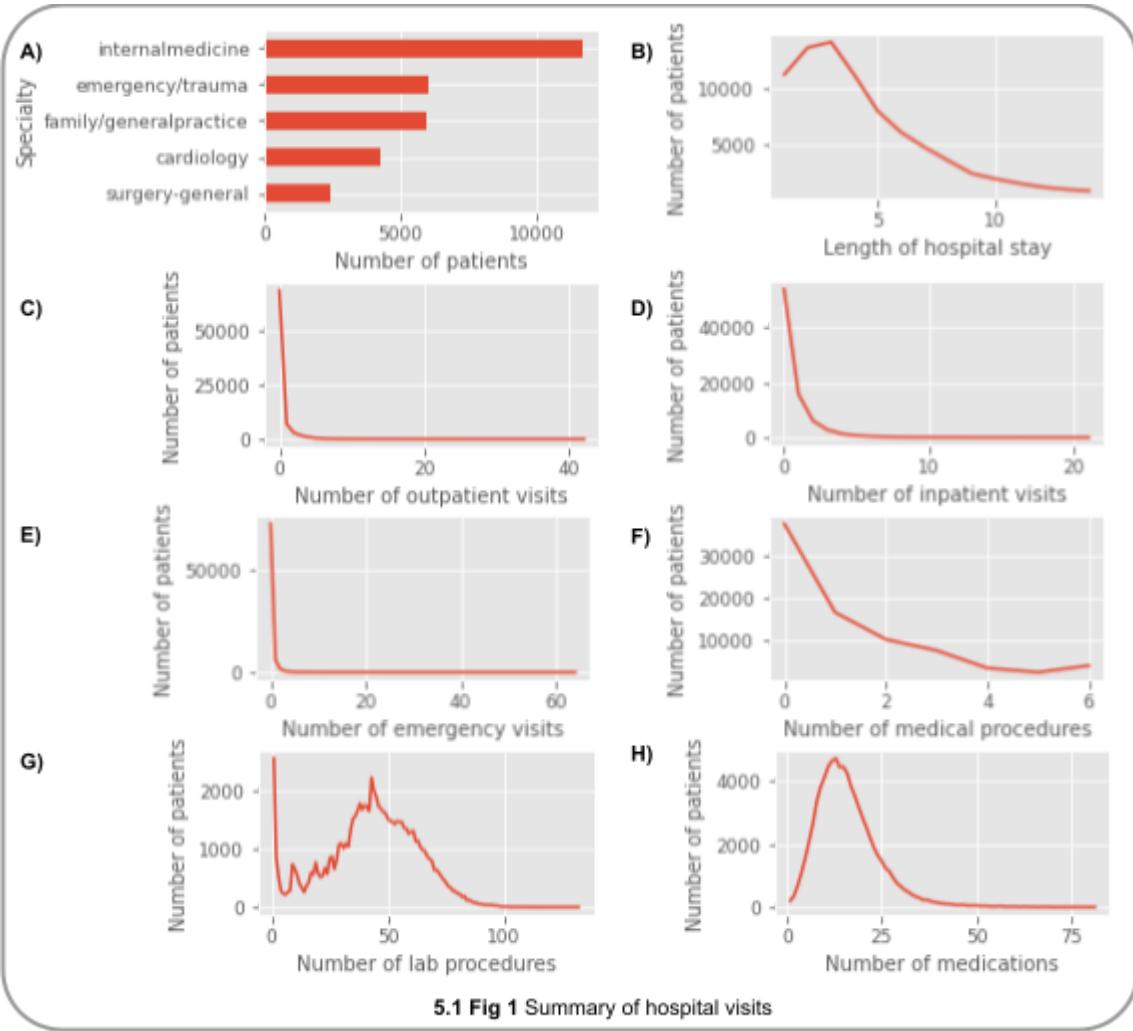
<https://drive.google.com/file/d/1JmzQmHlj3aktfVJbZ57gahxF85aKnImv/view?usp=sharing>

The table below shows a brief description of each of the fields provided after preprocessing (N.B, patient_id was used as an identifier and is only included here for completeness)

Name: patient_id, dtype: float64 count 8.141200e+04 mean 1.086395e+08 std 7.732453e+07 min 1.980000e+02 25% 4.683906e+07 50% 9.083437e+07	Name: race, dtype: object count 81412 unique 5 top white freq 60873	Name: gender, dtype: object count 81412 unique 3 top female freq 43719
---	---	--

75% 1.751117e+08 max 3.790052e+08		
Name: age, dtype: object count 81412 unique 11 top 70-80 freq 20261	Name: weight, dtype: object count 81412 unique 10 top unknown freq 78913	Name: admission_type_code, dtype: object count 80250 unique 6 top emergency freq 42562
Name: discharge_disposition_code, dtype: object count 81412 unique 7 top discharged_home freq 47854	Name: admission_source_code, dtype: object count 81412 unique 4 top emergency freq 45942	Name: time_in_hospital, dtype: float64 count 81412.000000 mean 4.395924 std 2.975844 min 1.000000 25% 2.000000 50% 4.000000 75% 6.000000 max 14.000000
Name: payer_code, dtype: object count 81412 unique 3 top insured freq 45131	Name: medical_specialty, dtype: object count 81412 unique 25 top unknown freq 40020	Name: has_prosthesis, dtype: object count 81412 unique 2 top False freq 80550
Name: complete_vaccination_status, dtype: object count 81412 unique 2 top True freq 67434	Name: num_lab_procedures, dtype: float64 count 79919.000000 mean 43.071197 std 19.630405 min 1.000000 25% 32.000000 50% 44.000000 75% 57.000000 max 132.000000	Name: num_procedures, dtype: float64 count 81412.000000 mean 1.341768 std 1.708465 min 0.000000 25% 0.000000 50% 1.000000 75% 2.000000 max 6.000000
Name: num_medications, dtype: float64 count 78734.000000 mean 16.024424 std 8.107235 min 1.000000 25% 10.000000 50% 15.000000 75% 20.000000 max 81.000000	Name: number_outpatient, dtype: float64 count 81412.000000 mean 0.370953 std 1.278538 min 0.000000 25% 0.000000 50% 0.000000 75% 0.000000 max 42.000000	Name: number_emergency, dtype: float64 count 81412.000000 mean 0.197588 std 0.881290 min 0.000000 25% 0.000000 50% 0.000000 75% 0.000000 max 64.000000
Name: number_inpatient, dtype: float64 count 81412.000000 mean 0.637793	Name: diag_1, dtype: object count 81412 unique 17 top circulatory	Name: diag_2, dtype: object count 81412 unique 17 top circulatory

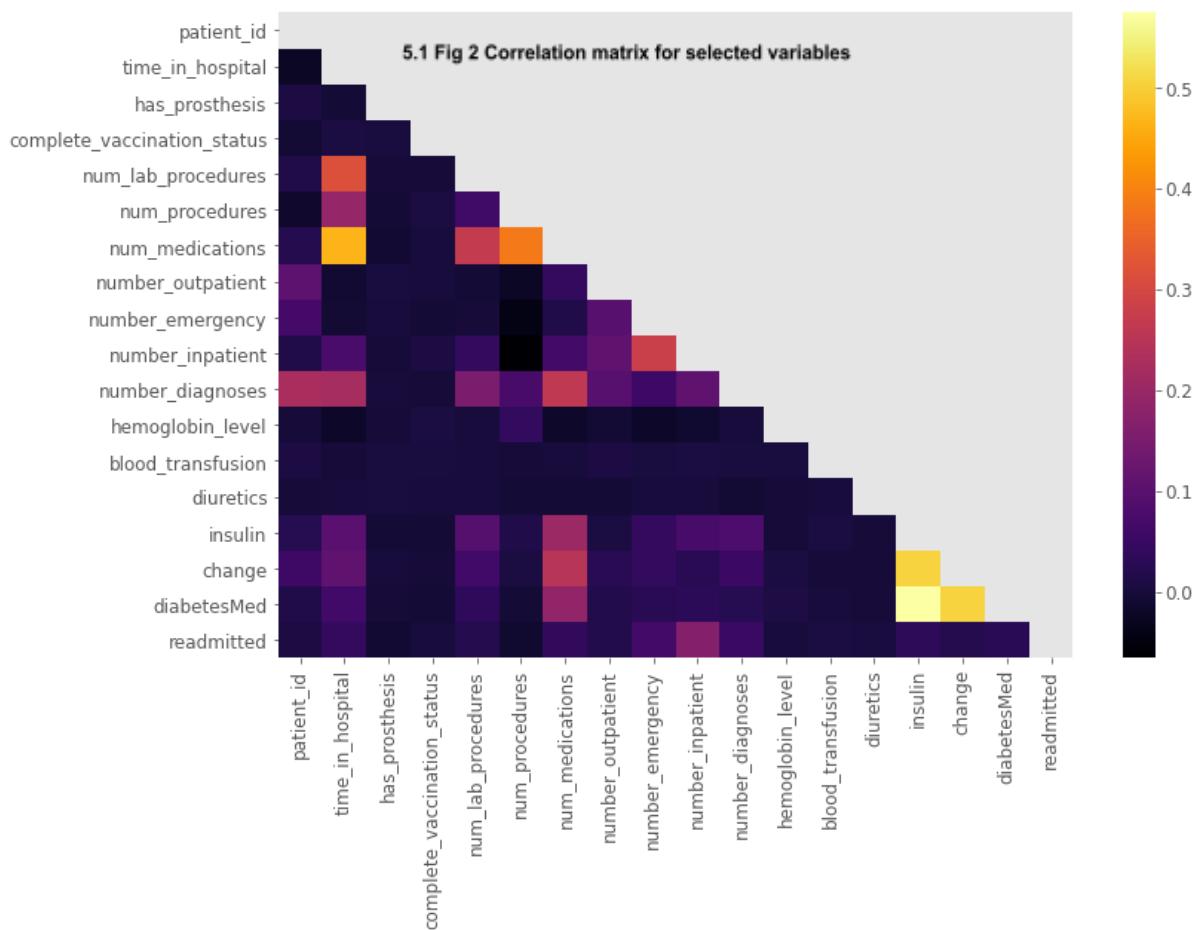
std min 25% 50% 75% max	1.265472 0.000000 0.000000 0.000000 1.000000 21.000000	freq 	24193	freq 	24666
Name: diag_3, dtype: object count unique top freq	81412 17 circulatory 23979	Name: number_diagnoses, dtype: float64 count mean std min 25% 50% 75% max	81412.000000 7.421965 1.931480 1.000000 6.000000 8.000000 9.000000 16.000000	Name: blood_type, dtype: object count unique top freq	81412 8 O+ 32053
Name: hemoglobin_level, dtype: float64 count mean std min 25% 50% 75% max	81412.000000 14.192328 1.060000 10.500000 13.400000 14.100000 15.000000 18.600000	Name: blood_transfusion, dtype: object count unique top freq	81412 2 False 71697	Name: max_glu_serum, dtype: object count unique top freq	81412 4 unknown 77159
Name: A1Cresult, dtype: object count unique top freq	81412 4 unknown 67807	Name: diuretics, dtype: object count unique top freq	81412 2 False 79893	Name: insulin, dtype: object count unique top freq	81412 2 True 44360
Name: change, dtype: object count unique top freq	81412 2 False 43772	Name: diabetesMed, dtype: object count unique top freq	81412 2 True 62718	Name: readmitted, dtype: object count unique top freq	81412 2 False 72340



5.1 Fig 1 Summary of hospital visits

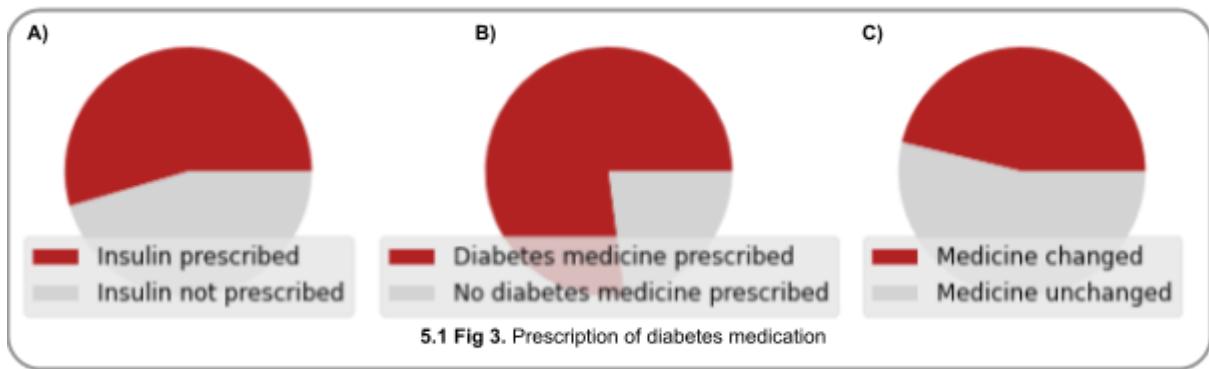
Numerical features were briefly explored and, a correlation matrix (overleaf) was produced to determine if any variable were strongly correlated to one another, or if they were duplicating information. Below are the top and bottom 10 correlation pairs and correlation coefficients:

insulin	diabetesMed	0.574090	num_procedures	number_inpatient	-0.065887	
change	diabetesMed	0.506269		number_emergency	-0.041881	
insulin	change	0.505491		number_outpatient	-0.024351	
time_in_hospital	num_procedures	0.464316	patient_id	time_in_hospital	-0.023890	
num_procedures	num_medications	0.386588		number_emergency	hemoglobin_level	-0.020732
time_in_hospital	num_lab_procedures	0.315046		time_in_hospital	hemoglobin_level	-0.018590
number_emergency	number_inpatient	0.281446		num_medications	hemoglobin_level	-0.018230
num_lab_procedures	num_medications	0.268071		patient_id	num_procedures	-0.014419
num_medications	number_diagnoses	0.261044		number_inpatient	hemoglobin_level	-0.011542
	change	0.247635		num_procedures	readmitted	-0.011276



Finally, considering the objective of this project being to predict readmission, correlation coefficients of the features with readmission were extracted (shown below)

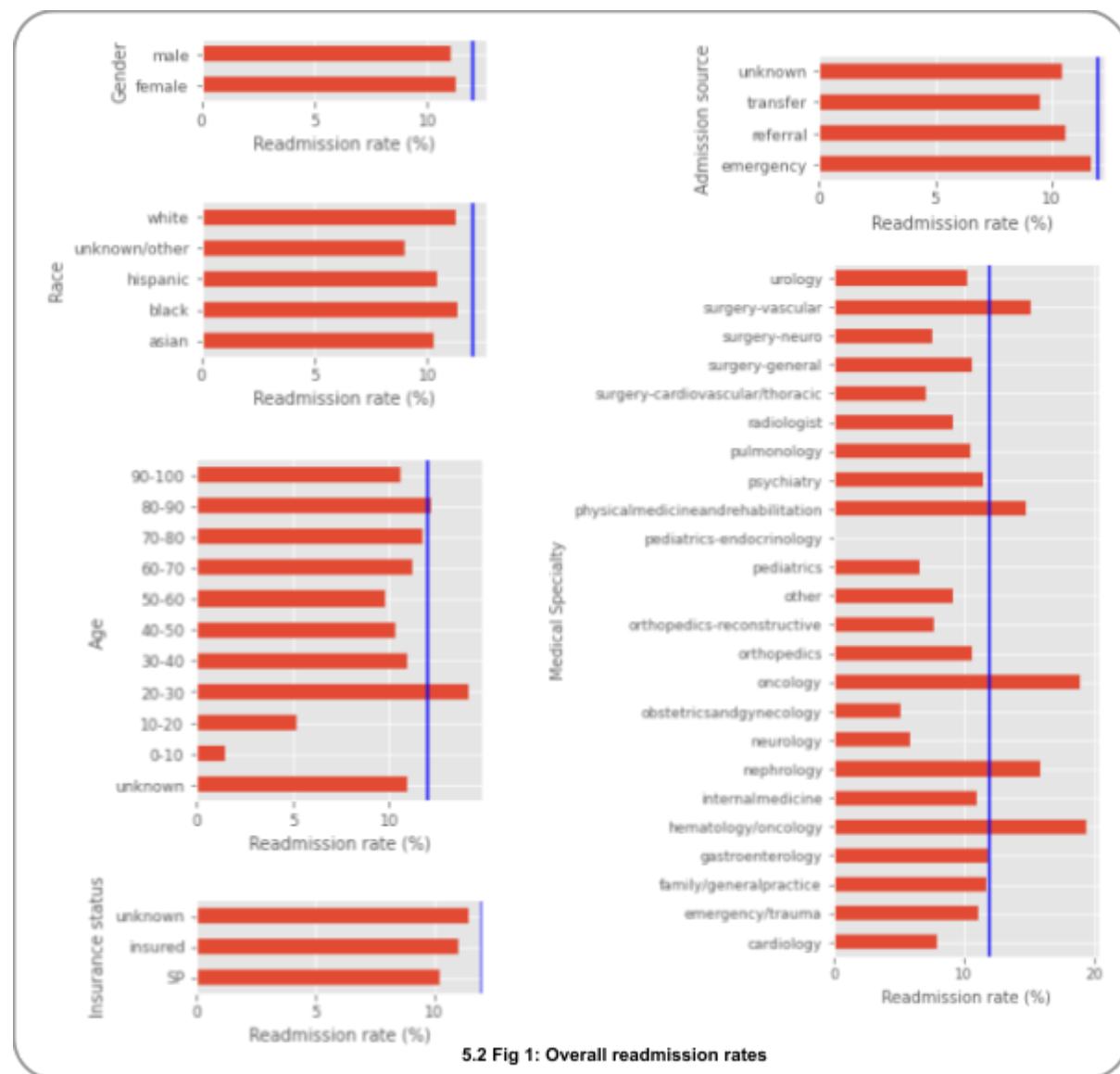
number_inpatient	0.164781
number_emergency	0.065709
number_diagnoses	0.049817
time_in_hospital	0.042647
num_medications	0.039105
insulin	0.033144
diabetesMed	0.027206
change	0.019311
num_lab_procedures	0.019160
number_outpatient	0.017730
patient_id	0.008108
blood_transfusion	0.005350
diuretics	0.001942
hemoglobin_level	-0.000169
complete_vaccination_status	-0.000558
has_prosthesis	-0.008412
num_procedures	-0.011276



5.2 Business questions technical support

1) Analysis of overall discrimination

Categories within each of the sensitive characteristics (age, gender, insurance status and race were grouped and the readmission rates were plotted. The same approach was used for admission sources and medical specialties, as shown below). The specialties represented are those which had at least 100 admissions in the time period covered, all others were grouped in a category titled “other” and a significant portion were unknown (data for the latter two is not shown as it is not informative). Readmission rates for the whole data set were approximately 11%. To allow for a slight variation the threshold of 12% was selected to identify factors resulting in above average readmission rates. In the future, it would be beneficial to carry out a deeper analysis and examine odds ratios and risk, as well as address the massive quantity of missing data. However, the simple analysis shown here allows us to rapidly identify problematic situations. Of note, patients between 20 and 30 years of age, and patients visiting hematology/oncology, oncology, vascular surgery, physical rehabilitation and nephrology had higher readmission rates.

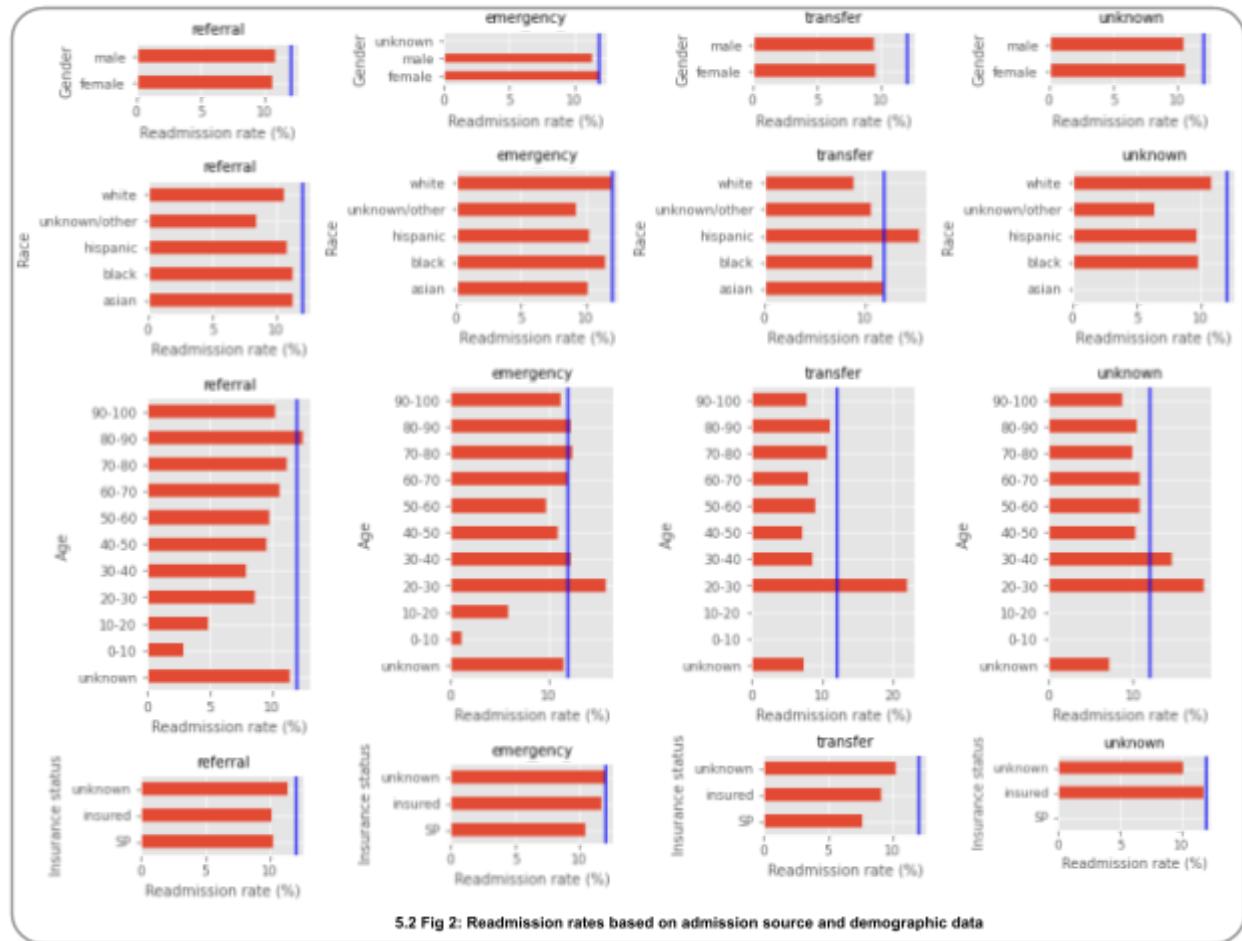


5.2 Fig 1: Overall readmission rates

2) Analysis of discrimination originating from admission sources

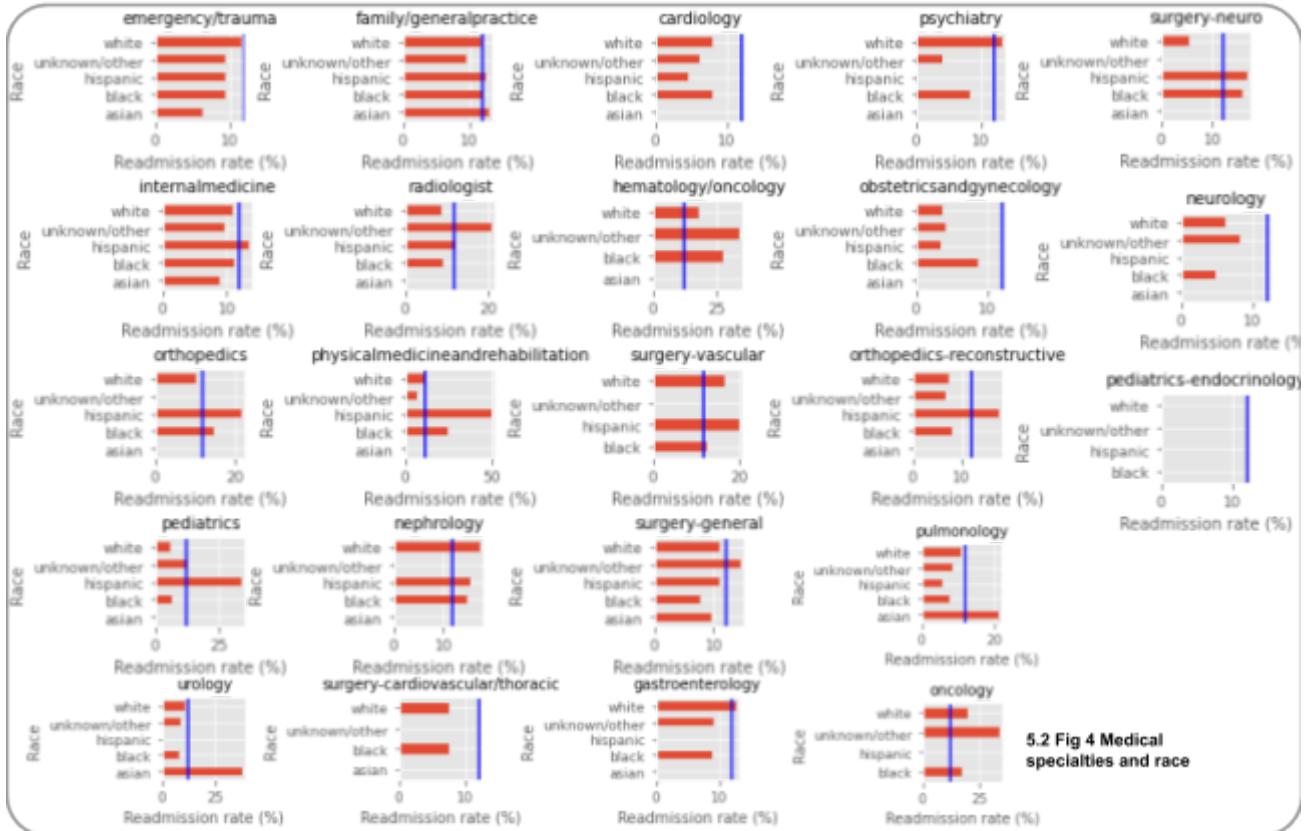
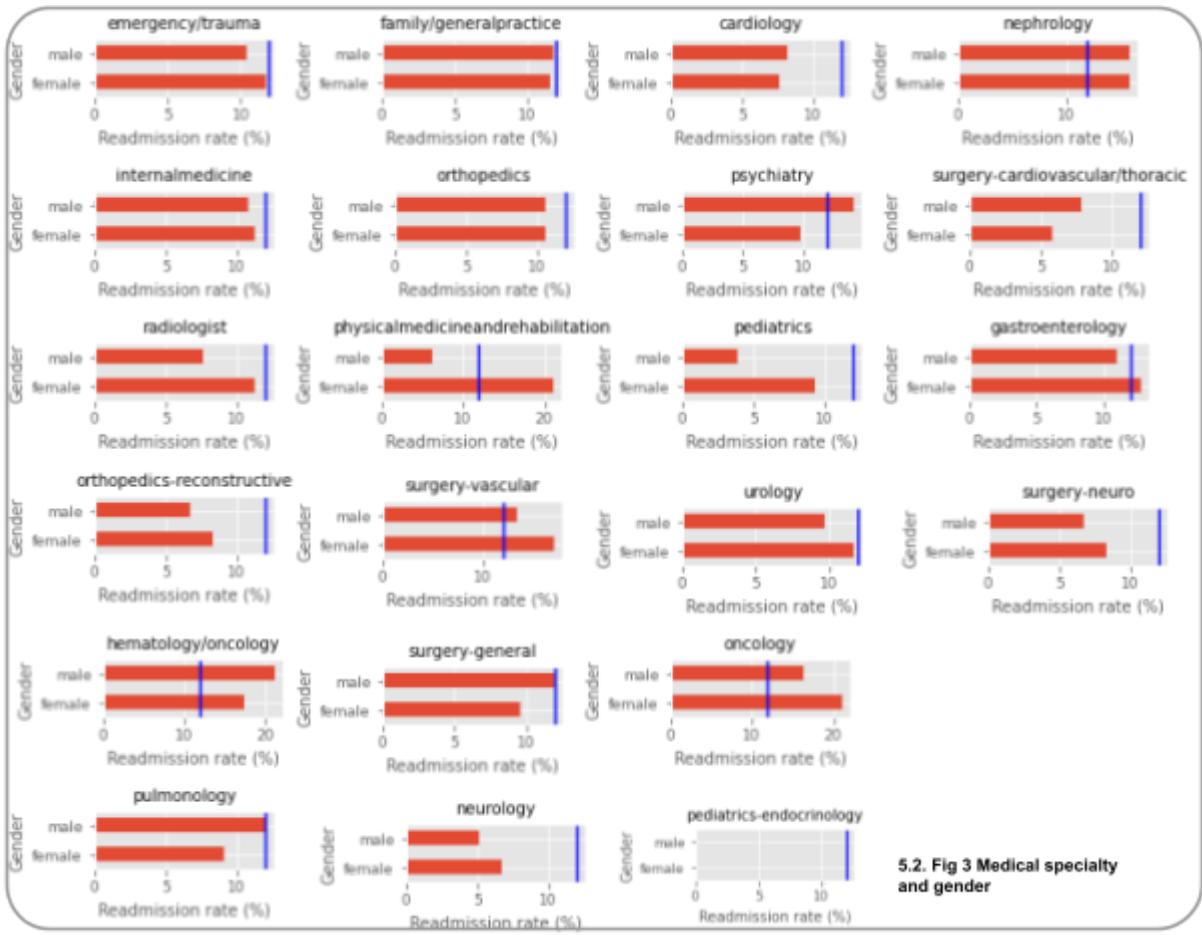
Admission sources were grouped into referrals, emergencies, transfers and unknown and the data was filtered to select each of these. Subsequently, divided into categories within the sensitive

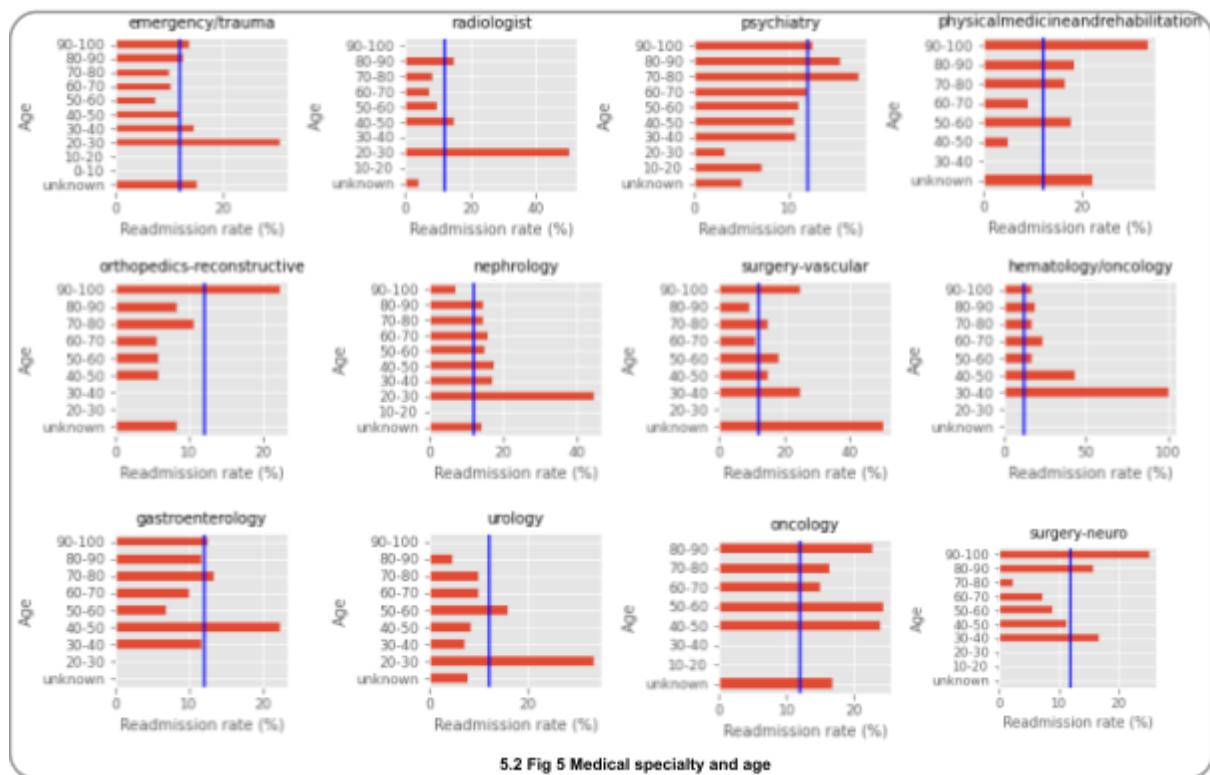
characteristics (age, gender, insurance status and race). Again, readmission rates above the threshold of 12% were deemed potentially problematic and further examined. In particular, although 20-30 year olds are a small portion of this dataset, the group consistently has higher readmission rates when having unreferrals visits to the hospital. A similar tendency is observed for hispanic patients who are transferred to the hospital from other health care services, and they also show a significantly raised readmission rate.



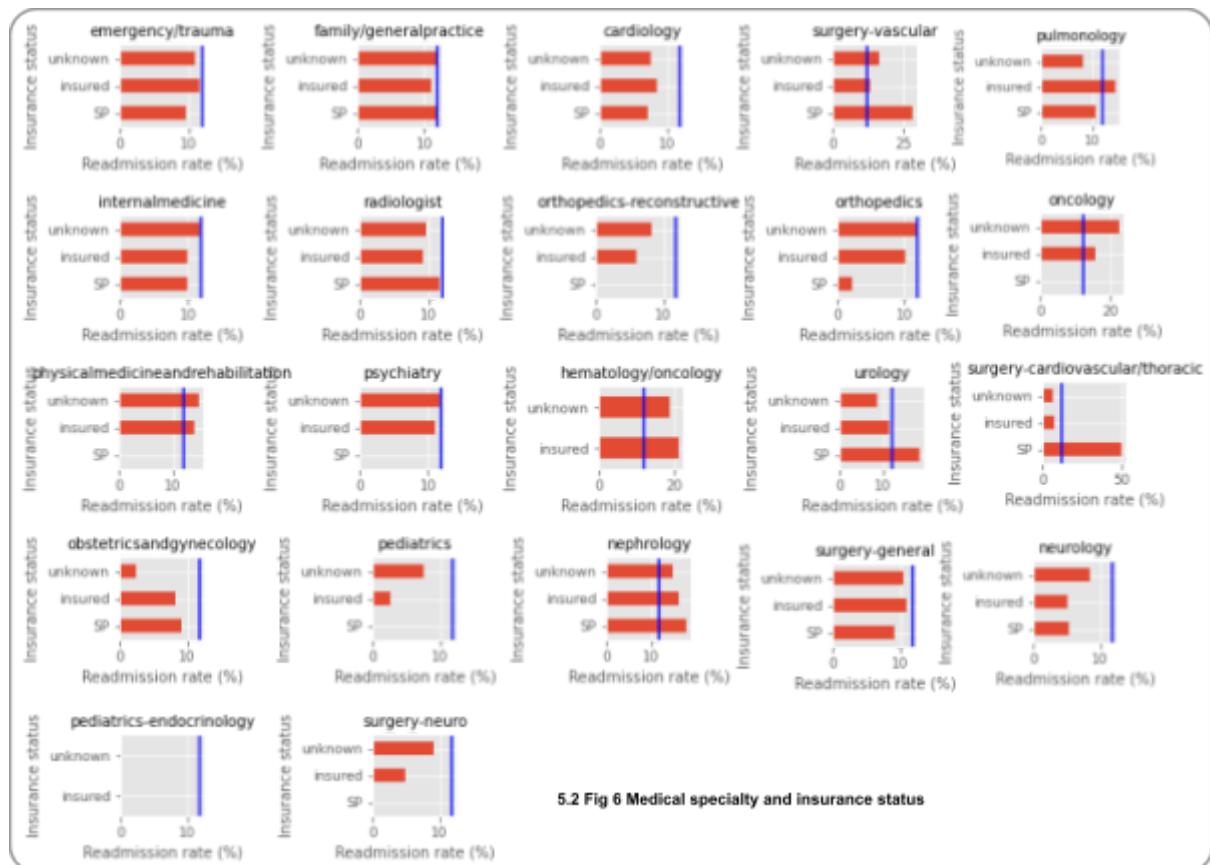
3) Analysis of discrimination originating from different specialties

Specialties with at least 100 admissions present were analysed. The data was filtered to select each of these and divided into categories within the sensitive characteristics (age, gender, insurance status and race) as above. The 5 specialties identified as having higher readmission rates had similar results here, although physical medicine and rehabilitation had a large discrepancy between men and women. Data for race is summarized, but it is difficult to draw conclusion as races other than black and white have very low representation in the dataset. Data for age shows that certain specialties have differences in readmission based on age, possibly due to the nature of the illnesses they treat, selected data is shown for this field. Lastly, there are discrepancies linked to patient insurance status in several specialties which may require further investigation (Data shown on following pages).





5.2 Fig 5 Medical specialty and age



5.2 Fig 6 Medical specialty and insurance status

5.3 Model technical analysis

Feature selection

A Random Forest Plot with no features excluded was carried out on transformed data. Examining feature importance revealed particular importance for discharge disposition, various patient health indicators (number of hospital visits, medications etc, some diagnosis subgroups).

Some variables, such as medical specialties, whether or not the patient has a prosthesis or is on diuretics did not appear to have a role.

Furthermore, features with a large amount of missing data (specifically ‘weight’ were excluded). In categorical analysis, missing data was imputed as ‘unknown’ in its own category.

While patient ID appears to contribute, it is likely that this is due to a few patients who are frequently at the hospital.

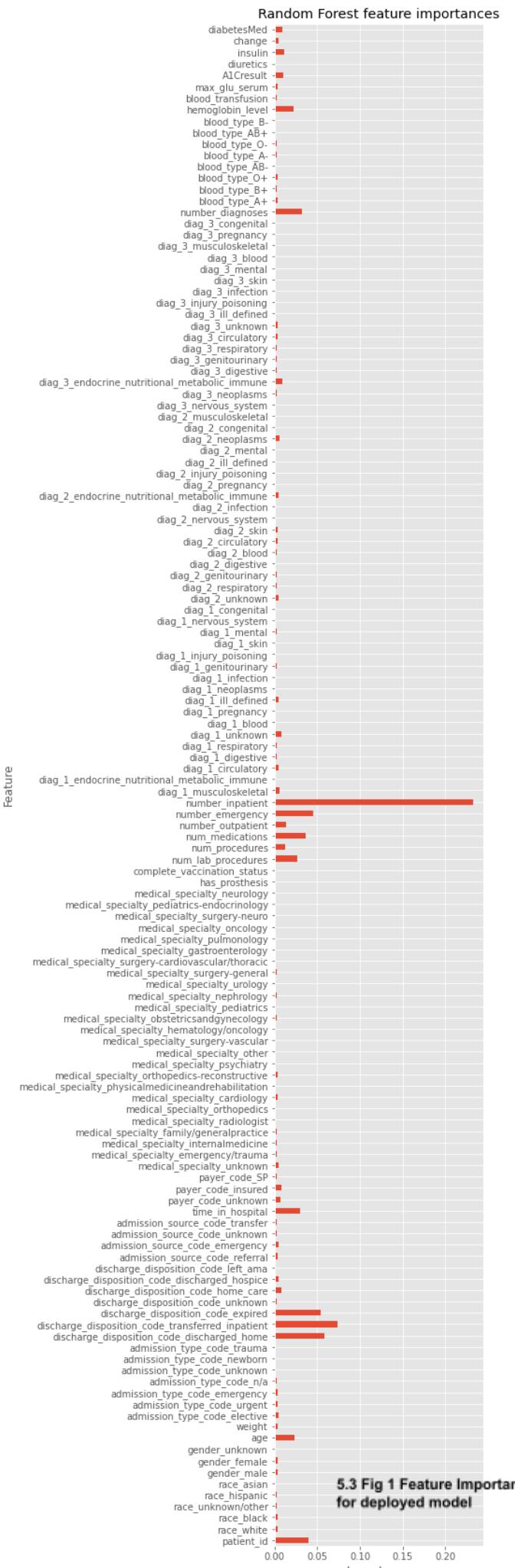
In conclusion, all features were used for model deployment except the following: weight, diuretics status, vaccination status, prosthesis status, medical specialty and patient ID.

Sensitive features were also included in the model - if they are potential reasons for wrongful discharge, it is reasonable to “correct” for that in a predictive system looking to prevent these incidents. Once root causes for this potential discrimination are solved, these features can be excluded.

Categorical encoding

Categorical variables were encoded as follows:

10. 'race': string entries were processed to automatically assign ‘white’, ‘black’, ‘hispanic’, ‘asian’ and ‘other’
11. 'gender': left as is
12. 'admission_type_code': codes were grouped to assign either ‘referral’, ‘transfer’ or ‘emergency’, missing data was treated as ‘unknown’
13. 'discharge_disposition_code': codes were grouped to assign either ‘home’, ‘transfer’ or ‘home_care’, ‘hospice’, ‘left_ama’ or ‘expired’, missing data was treated as ‘unknown’
14. 'admission_source_code': codes were grouped to assign either ‘emergency’,

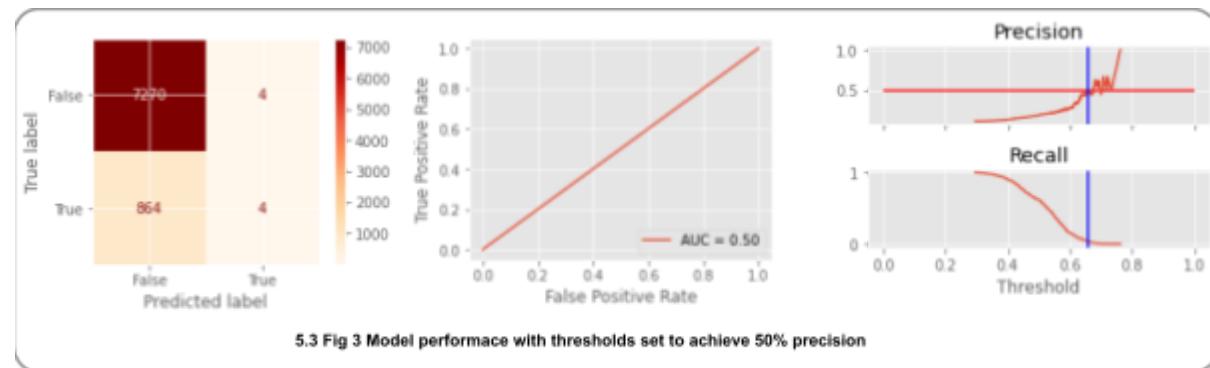
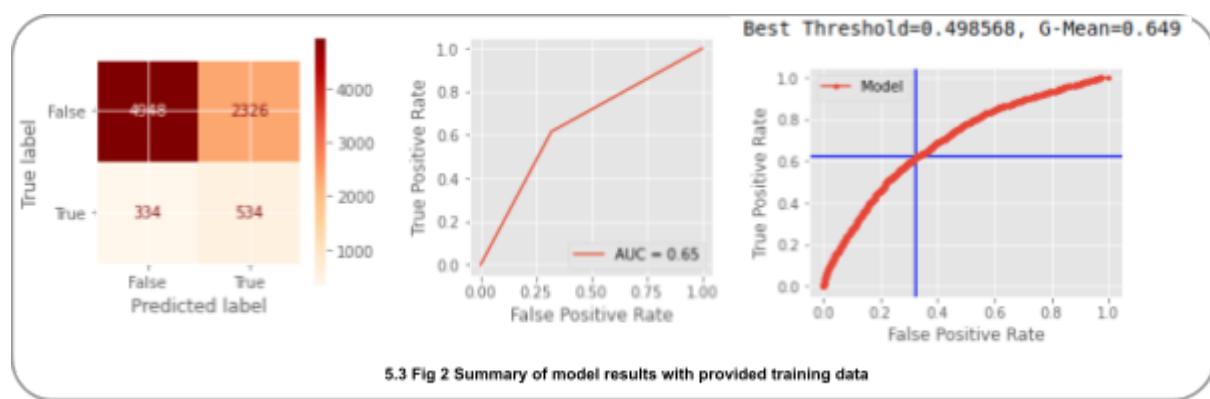


- 'elective', 'newborn', 'trauma' or 'urgent', missing data was treated as 'unknown'
15. 'payer_code': patients with codes other than SP were considered insured, missing data was classified as 'unknown'
 16. 'medical_specialty': specialties with at least 100 admissions were retained, others were grouped into 'other', missing data was treated as 'unknown'
 17. 'blood_type': left as is (imputation errors were found, with patient blood types changing - further clarification and correct values were not provided)
 18. 'diag_1', 'diag_2' and 'diag_3': were grouped and assigned into one of 19 categories based on CDC classification of diagnostic codes. This data is parsed and available at:

<http://icd9.chrisendres.com/index.php?action=contents>

Model performance

Model performance evaluated using ROC AUC and F1 scores, using current deployed system (top) and thresholds to attain 50% true positive rate (bottom). The latter figure clearly shows that the current mode, with thresholds set to attain a 50% true positive rate would not result in a reduction in wrongful discharge and consequently, this option was discarded. A precision-recall curve was generated but provided no additional insight, so is not included.



Discrimination assessment

To assess model discrimination, precision scores were calculated for different subgroups with at least 50 representatives in the test data set. The maximum difference of the precision for tests within each feature was used as an indicator for discrimination, according with the technical requirements provided (10% difference between sub-groups, and less than 5% between medical specialties). Results are shown with feature in bold, followed by a list of category:precision score pairs and a final statement of whether or not the requirement was met.

'medical_specialty':
'cardiology': 0.1111111111111111,
'unknown': 0.1912162162162162,
'internalmedicine': 0.15037593984962405,
'emergency/trauma': 0.17843866171003717,
'family/generalpractice': 0.23383084577114427,
'nephrology': 0.2222222222222222,
'psychiatry': 0.13043478260869565,
'surgery-cardiovascular/thoracic': 0.0,
'pulmonology': 0.2857142857142857,
'other': 0.18518518518518517,
'surgery-general': 0.2261904761904762,
'orthopedics': 0.16666666666666666666,
'urology': 0.125,
'radiologist': 0.15384615384615385,
'orthopedics-reconstructive': 0.14285714285714285,
'gastroenterology': 0.2},
Maximum difference = 0.29, Requirement not met

'payer_code':
'insured': 0.17593123209169054,
'unknown': 0.2059732234809475,
'SP': 0.16037735849056603},
Maximum difference = 0.04, Requirement met

'gender':
'male': 0.19096671949286848,
'female': 0.1814102564102564},
Maximum difference = 0.01, Requirement met

'Race':
'white': 0.18140589569160998,
'unknown/other': 0.19753086419753085,
'black': 0.1958333333333333,
'hispanic': 0.23809523809523808,
'asian': 0.2857142857142857},
Maximum difference = 0.10, Requirement met (value is rounded down to 10%, asians have low representation (only 55 in test set), but meet threshold of 50)

'age':
'40-50': 0.1836734693877551,
'80-90': 0.18681318681318682,
'70-80': 0.17209908735332463,
'90-100': 0.1320754716981132,
'60-70': 0.19047619047619047,
'50-60': 0.1972972972972973,
'30-40': 0.25925925925925924,
'20-30': 0.375,
'unknown': 0.17142857142857143,
'10-20': 0.0},
Maximum difference = 0.375, Requirement not met

'admission_source_code':
'emergency': 0.18147549811523964,

'referral': 0.19175911251980982,
'transfer': 0.20253164556962025,
'unknown': 0.19318181818181818},
Maximum difference = 0.02, Requirement met