

LISBON  
DATASCIENCE  
ACADEMY

# MINIMIZING ERRONEOUS PATIENT DISCHARGES

A data science analysis on the Hazel and Bazel Hospital's unexpected patient discharge rate

**Prepared for:**  
Hazel and Bazel Hospital

**Prepared by:**  
Rafael Gil  
Data Scientist  
Awkward Problem Solutions™

# Contents

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Business Conclusions</b>                 | <b>3</b>  |
| 1.1      | Summary . . . . .                           | 3         |
| <b>2</b> | <b>Results Analysis</b>                     | <b>4</b>  |
| 2.1      | Model Performance . . . . .                 | 4         |
| 2.2      | Success on requirements . . . . .           | 5         |
| 2.3      | Population Analysis . . . . .               | 6         |
| <b>3</b> | <b>Next Steps</b>                           | <b>10</b> |
| 3.1      | Next Steps . . . . .                        | 10        |
| <b>4</b> | <b>Deployment Issues</b>                    | <b>11</b> |
| 4.1      | Re-deployment . . . . .                     | 11        |
| 4.2      | Unexpected problems . . . . .               | 12        |
| 4.3      | Learnings and Future Improvements . . . . . | 12        |

# 1 Business Conclusions

## 1.1 Summary

Awkward Problem Solutions™ provided a REST API to validate patient discharges at the Hazel and Bazel hospital. This was part of an effort to minimize the erroneous discharge rate the hospital was noticing. The service run in production for a whole week validating a total of 9694 medical encounters. Besides predicting whether a patient would be readmitted at the hospital within the following 30 days, the service also collected data for future iterations of the model, including the real outcome of each medical encounter, i.e. whether the patient was actually readmitted within a 30 day time-frame after discharge.

The service was stable throughout the whole week in production, requiring no interventions or updates. The application was able to reject malformatted requests, informing the end user of the error, and to pass valid data onto the predictive model. The predictive model itself handled all kinds of data successfully, by outputting a result at every request.

There was minor a issue regarding data gathering, causing 200 medical encounters not to receive their true label, i.e. whether that medical encounter resulted in readmission or not.

Regarding performance, the model achieved a 52 % recall score, which measures the model's capacity to minimize the erroneous discharge rate. This result means, that 52 % of the medical encounters that would result in a readmission were identified by the model, thus minimizing the patient discharge rate. However, the model failed the minimum 50 % precision score requirement, having achieved only 18 % precision. This means that only 18 % of the medical encounters classified as a readmission were actually a readmission, in other words only 18 % of the patients flagged for readmission were actually sick.

The present results indicate that the model may have overfitted to the training data, as warned in the previous report. An evidence of overfitting is that the real performance drastically dropped when compared to the expected performance. The unsatisfactory performance, could also have been related to a population shift, thus rendering the old data as obsolete, and consequently the model provided. However, a population analysis revealed no population shift, thus enhancing the overfitting probability.

Finally, according to the previous report predictions, not all discrimination requirements were met, with the exception of gender and insurance status discrimination criteria. The requirement states that the wrongful discharge rate, or readmission rate, should not vary more than 10 percentage points between sub-groups, and less than 5 percentage points between medical specialties. As mentioned, gender and insurance status meet the criteria with a maximum readmission rate variation of 0.94 % and 6.9 %, respectively. The remaining sensitive groups, age and race, did not meet the requirement by achieving a maximum readmission rate variation of 26.5 % and 13.0 %, respectively. The readmission rate among medical specialties also did not meet the criteria with a maximum readmission rate variation of 27.5 %.

Nonetheless, a data analysis revealed that, analogously to the initial training data, the collected data has multiple discrimination issues on many medical specialties across different subgroups. These issues are not related to the model itself, but to the hospital care instead.

## 2 Results Analysis

### 2.1 Model Performance

The production deployment allowed the collection of 9694 new medical encounters, and their true outcome is known for all but 200. These 200 medical encounters were discarded on the present analysis, since their readmission outcome is not known.

In order to evaluate the model’s performance, the following metrics are used: *accuracy*, *recall*, *precision*, and *F1 Score*. Each metric provides a different insight on the actual model’s performance. Table 1 lists the score of those metrics on the test set, as detailed on the previous report, and on the newly collected data. The table also includes the score difference between those datasets to ease performance evaluation.

Table 1: Model performance on the test set and also in production.

|            | ACCURACY | RECALL | PRECISION | F1 SCORE |
|------------|----------|--------|-----------|----------|
| PRODUCTION | 68 %     | 52 %   | 18 %      | 0.27     |
| TEST SET   | 68 %     | 79 %   | 65 %      | 0.71     |
| DIFF       | 0 %      | -27 %  | -47 %     | -0.44    |

Accuracy evaluates how many classifications the model got right, regardless of the classification label. Naturally, this metric is highly sensitive to imbalanced datasets, such as a the current medical encounter dataset. The collected data during production has 1074 readmitted medical encounters out of a total of 9494 medical encounters, thus representing a fraction of 11 % of the total data. Consequently, the accuracy score could have been achieved by mostly identifying the non-readmissions correctly. Nonetheless, the production’s accuracy score met the expectations laid out on the test set.

Recall measures the model’s ability to find all the medical encounters that were indeed readmitted. Thus, it is a proxy measure of how good the model is at minimizing the erroneous patient discharge rate. The expectations regarding recall were set at 79 % on the test set. However, that score dropped down to 52 % on the newly collected data. This result means that more than half of medical encounters that would result in a readmission were detected by the model. It is, nonetheless, a positive outcome even though there is room to improve in future iterations.

Conversely, the precision score measures the model’s ability not to mark as a potential readmission the medical encounters that won’t be readmitted. Conversely, it is a proxy measure of how good the model is at avoiding to flag patients that are not actually sick. Results regarding precision reveal a steeper drop, from 65 % to 18 %, when comparing the precision score on the test set to the precision score on the newly collected data. These result shows that, during the production week, only 18 % of the medical encounters flagged for readmission actually resulted in a readmission. The reasoning behind the steep precision drop is twofold. On one hand, the model failed to generalize the training data due to overfitting, resulting its misinterpretation of the root causes of a readmission. On the other hand, the true label, i.e. a readmitted medical encounter, is the minority class. Thus the probability of false positives, which highly impacts the precision score, increases drastically.

As mentioned in the previous report, there is a trade-off relationship between recall and precision, hence the F1 score. The F1 score is the harmonic mean of precision and recall, thus it provides a statistical measure of the balance between precision and recall. Naturally, since both precision and recall scores dropped on the production dataset, F1 score also dropped from 0.71

to 0.27. This metric provides a quick overview of the lower model performance on the newly collected data, when compared to the test set.

**Discrimination is defined as a steep variation on the readmission rate of the sensitive classes and the medical specialties. As discussed in the previous report, the readmission rate can be measured through the precision score.**

Table 2 lists the maximum precision variation, and consequently the readmission rate variation, observed across each of the sensitive classes and also across the medical specialties. Furthermore, this data can also be visualized on figures 1 and 2 as discussed in section 2.2.

Results show a readmission rate variation above 25 % across age, and medical specialties sub-groups. The readmission rate variation across race sub-groups is 13.0 %, but below 10 % for gender and insurance status. Naturally, a higher readmission rate variation means that there are higher readmission rate discrepancies across subgroups of the sensitive class, hence there are potentially higher discrimination issues.

**Table 2:** Model's discrimination results based on the production data.

| PRECISION VARIATION |        |
|---------------------|--------|
| AGE                 | 26.5 % |
| RACE                | 13.0 % |
| GENDER              | 0.94 % |
| INSURANCE STATUS    | 6.9 %  |
| MEDICAL SPECIALTY   | 27.5 % |

## 2.2 Success on requirements

The first requirement regarding the model performance is that at least 50 % of the patients identified for readmission should actually be sick. This requirement translates into a minimum 50 % precision score on the overall model, and, as mentioned in the previous section, the threshold was not met. The model achieved a precision score of 18 % according to table 1.

Nonetheless, there was one more business requirement stating that the wrongful discharge rate, or readmission rate, should not vary more than 10 percentage points between sub-groups, and less than 5 percentage points between medical specialties. As mentioned, the readmission rate can be measured using the precision score for each sub-group and medical specialty.

Figure 2 plots the precision score of the deployed predictive model for the sensitive sub-groups. Two sub-groups meet the criteria: gender, and insurance status with precision scores ranging a maximum of 0.94 % and 6.9 %, respectively. However, the criteria is not met for age, and race with with precision scores ranging a maximum of 26.5 % and 13.0 %, respectively. **Nonetheless, the observed precision variation is lower than the 36.7 % and 19.1 % variation expected on report 1, i.e. an improvement of 10.1 % and 6.1 % respectively.**

Regarding, the medical specialty the criteria is not met, since the precision scores ranges a maximum of 27.5 %. Once again, the observed precision variation, for the medical specialties, is lower than the 47.3 % variation expected on report 1, i.e an improvement of 19.8 %.

However, as also discussed on report 1, it is worth noting that the sub-groups where the criteria failed are also the sub-groups with a higher number of categories. Empirically, spreading the existing medical encounters throughout a higher number of categories, increases the preponderance of those medical encounters. Thus, even though the criteria failed, immediate conclusions should not be drawn because a higher volume of data could prove otherwise.

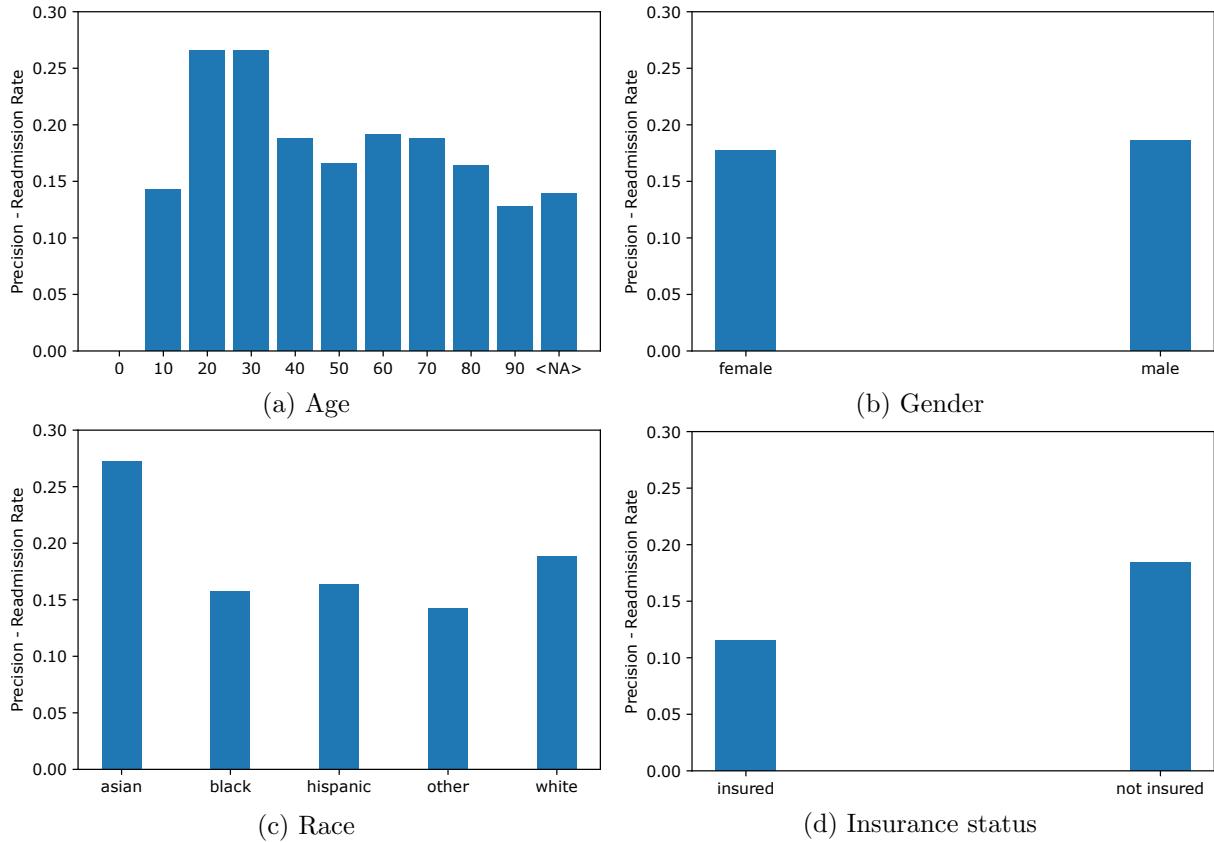


Figure 1: Model readmission rate for sensitive sub-groups

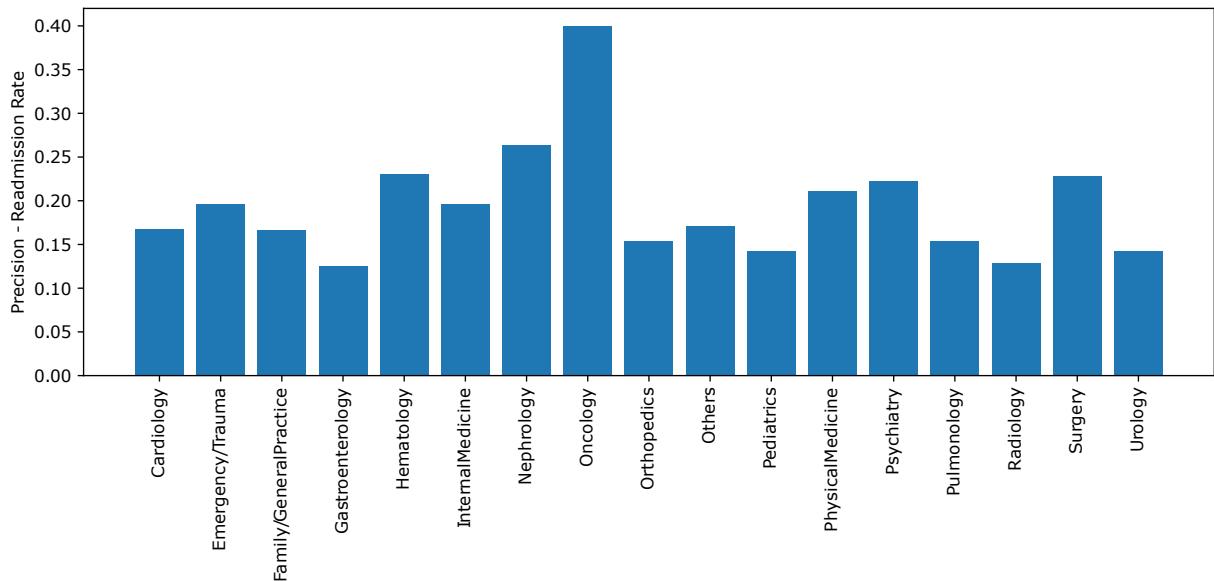


Figure 2: Model readmission rate for medical specialties

### 2.3 Population Analysis

The production deployment allowed the collection of 9694 new medical encounters, and their true outcome is known for all but 200. These 200 medical encounters were discarded on the present analysis, since their readmission outcome is not known.

This new dataset has 50 different medical specialties. Analogously to the first dataset, roughly 50 % of the medical encounters are missing a medical specialty. As discussed in the previous report, `medical_specialty` had a high volume of unique classes, and to mitigate that the classes were merged into the main medical specialty category. Out of the 50 different medical specialties, only did not exist previously: *Surgery-PlasticwithinHeadandNeck*. Consequently, it was mapped into the *Others* category, even though it belong to the *Surgery* class. Thus, this new class was added to the pre-processing logic.

After this update to the medical specialties mapping, figure 3 was plotted with the distribution of medical encounters per medical specialty. The figure shows that the medical specialty distribution is consistent with the initial dataset. The top-7 medical specialties are exactly the same both in the initial dataset and in the current collected data from the REST API.

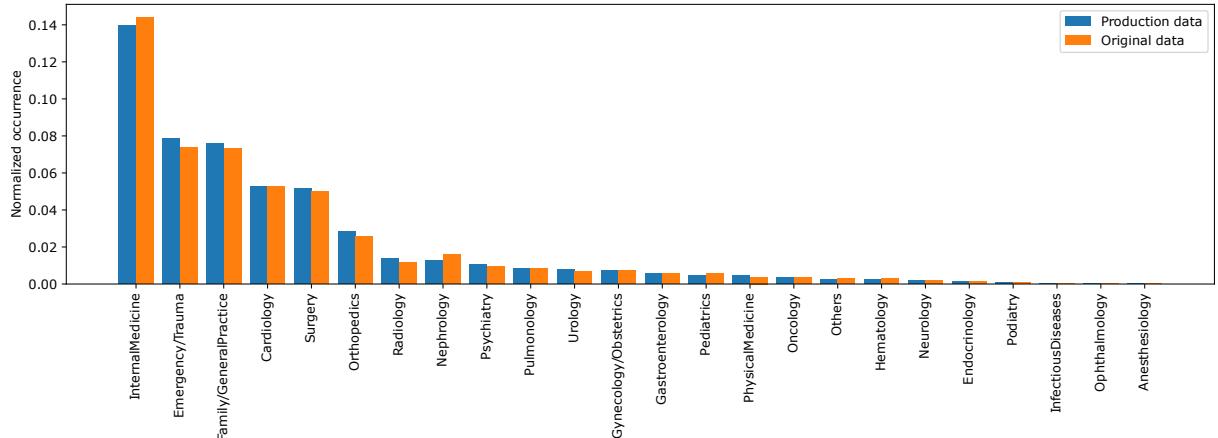


Figure 3: Medical specialties' occurrence on both datasets after processing

Analogously to the previous report, figure 4 plots the occurrence of each sensitive class on the new dataset, including age, gender, race, and insurance status. In blue its plotted the occurrence of each sensitive class on the whole dataset, whilst in orange its plotted the occurrence of the same sensitive class across the readmission-only subset.

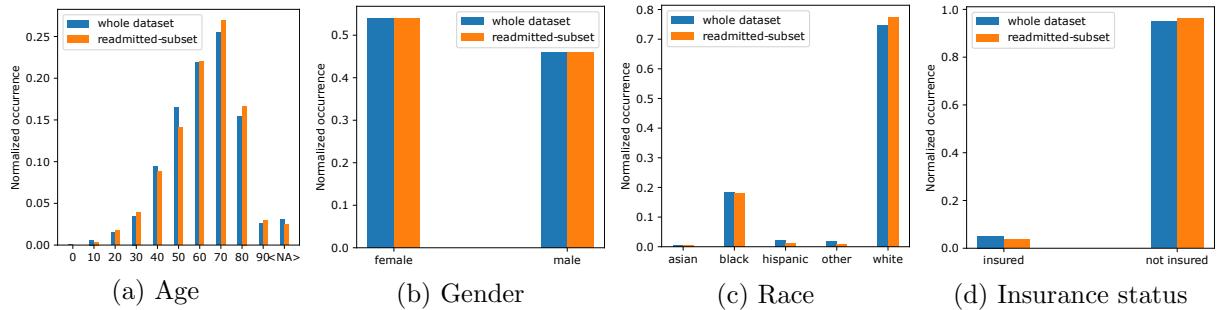


Figure 4: Sensitive classes distribution across the whole dataset and a readmitted-only subset

The graphs on figure 4 are highly similar to the graphs from the previous reports. There are only slight variations, for instance, on figure 4c, the **black** race on the new data set have a slightly higher occurrence rate on the whole dataset when compared to the readmitted subset. Nonetheless, these slight variations are neglectable when looking at the overall distribution. The new data, extracted during the REST API production week, looks like a careful replication of the population from the initial dataset.

Furthermore, as suggested by table 3, there are still evidences of discrimination on several

medical specialties against different subgroups, such as age, gender, and race. Nonetheless, there are no evidences of discrimination based on the insurance status of the patient on any of the medical specialties.

Proceeding with the corruption analysis, figure 5 reveals the analysis to the new collected dataset. Recall that a medical encounter is considered corrupted if there is a second one where the same patient has a different `blood_type`. Figure 5a shows that 7.0 % of the data set is corrupted, i.e. 7.0 % of the medical encounters have at least another medical encounter where `blood_type` does not match for the same patient.

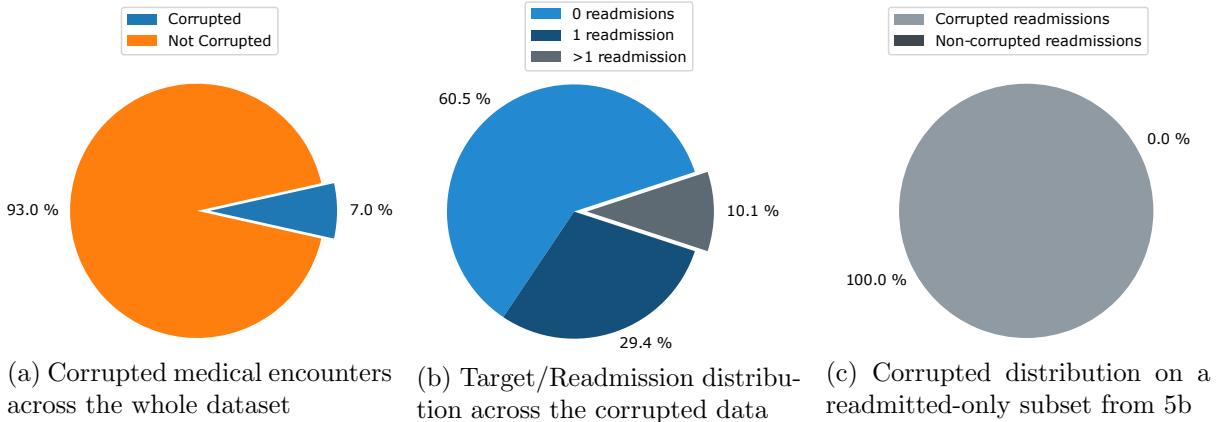


Figure 5: Data integrity analysis of based on patients who visited the hospital at least twice

Figure 5b drills down into the corrupted entries and focus on the target variable: `readmitted`. It shows that 60.5 % of those corrupted medical entries belong to patients that were never readmitted into the hospital. 29.4 % belongs to patients that were readmitted once. The remaining, 10.1 % are entries of patients that were readmitted into the hospital at least twice.

Finally, figure 5c focus on the 10.1 % of corrupted medical entries that belong to patients readmitted more than once. Assuming that a medical encounter is handled with extra care when a readmission occurs, and knowing that these patients were readmitted more than once, this chart checks if the all their readmitted medical encounters have a matching `blood_type`, i.e. if the re-admissions are corrupted when only compared to re-admissions. It turns out that all of those medical encounters have mismatched `blood_type` observations for the same patient.

Table 3: Detailed per-specialty discrimination analysis on the new collected data. All values are normalized to the overall medical specialty readmission occurrence on the last column. Values above the 25 % threshold are marked in **bold**, values below the threshold are underlined. Dashes mean there were no occurrences for that class, whilst zeros mean there were no re-admissions for that class.

## 3 Next Steps

### 3.1 Next Steps

Looking at the results drawn from the collected data during production, there are areas of improvement regarding the predictive model. The model failed some of the key business requirements and those issues should receive most of the attention.

The model seems to have overfitted to the training data, which was one of the dangers raised on the first report. The main reason for this issue is probably related to the unbalanced dataset that forced an undersampling of the training data available. Perhaps, alternatives to the undersampling should be considered, such as oversampling. Oversampling could provide better results, given the current scenario, since no data is discarded, thus helping the model to better fit the problem.

Nonetheless, an investigation should be conducted in order to identify measures that would prevent the *Gradient Boosting Classifier* from overfitting. The previous report ignored hyper parameters such as `min_samples_leaf` and `min_samples_split` that seem beneficial to prevent overfitting. Increasing these hyper-parameters makes it harder for the model to memorize single observations, and potentially preventing this issue.

The discrimination requirements, as mentioned above, are hard to achieve on features with a high number of categories. This is related to the spread of observations across a high number of categories leading to a scarce amount of observations per category. This phenomenon will skew metrics to extremes, either very high or very low, making it harder to meet criteria.

Furthermore, the Hazel and Bazel hospital seems to be dealing with an issue of discrimination on different medical specialties, as warned on the first report and repeated on the new data on the present report. The hospital should focus on analysing the root causes of such issues on a per medical specialty basis. Higher resources should be applied trying to solve the issue at the root, i.e. on the patient care itself, rather than on the readmission validation, since these criteria lower the model performance. Naturally, lowering the model performance results in a model that is less accurate compared to what it could possibly be. In the end it, the Hazel and Bazel hospital should evaluate if the discrimination issues outweigh a wrongful discharge.

## 4 Deployment Issues

### 4.1 Re-deployment

The deployment went smoothly from start to finish with no major warnings during the production week. Naturally, not having access to the complete logs of the application forced constant monitoring to catch unexpected errors. Nonetheless, there is the risk of having missed some error messages, but taking into account the amount of data collected and the expected volume of traffic, it seems unlikely.

Given the model's under-performance based on the collected data analysis, a new version aiming towards a performance improvement was indispensable to face the second week in production. The existing model's under-performance is the main driver of a redeployment, because, according to section 2.3, there is not a noticeable data drift on the input data received during the first week in production. The distribution of medical encounters across sensitive variable's groups, and also the medical specialties distribution, both remain relatively similar to the initial training dataset.

Thus, the new collected data was merged with the old data in order to create a new bigger dataset. The final merged data set has a total of 90 906 observations, where only 10 146 are readmitted medical encounters. In order to balance the data, and avoid data loss, the under-represented class was re-sampled 8 times, resulting in a balanced dataset with a total of 161 928 medical encounter where 50.13 % are readmitted.

Afterwards, a grid search was conducted in order to find the best hyper-parameters for the *Gradient Boosting Classifier*. This time, as suggested in the previous section, the search also included `min_samples_leaf` and `min_samples_split` parameters in an effort to prevent overfitting. The grid search result is listed in table 4.

Table 4: Retrained model's hyper-parameters.

| HYPER-PARAMETER                |      |
|--------------------------------|------|
| <code>random_state</code>      | 42   |
| <code>n_estimators</code>      | 100  |
| <code>learning_rate</code>     | 30 % |
| <code>max_depth</code>         | 10   |
| <code>min_samples_leaf</code>  | 100  |
| <code>min_samples_split</code> | 125  |

The retrained model expected performance improved when compared to the previous model. Its new performance metrics, based on the test set, are available on table 5.

Table 5: Retrained model's performance on the test set (20 % of the whole merged dataset).

|          | ACCURACY | RECALL | PRECISION | F1 SCORE |
|----------|----------|--------|-----------|----------|
| TEST SET | 81 %     | 92 %   | 75 %      | 0.83     |

Finally, after a discrimination analysis for the final model similar to the one on section 2.2, the discrimination requirements also improved. The new model is now expected to fulfill discrimination threshold for the race, gender, and insurance status sub-groups with a maximum variance of 7.48 %, 1.90 %, and 0.56 % respectively. It almost meets the age discrimination threshold with a maximum expected variance of 11.12 %. Regarding, the medical specialty the

maximum discrimination threshold is still overrun, but the variance has drastically improved to a maximum 18.09 % when compared to the previous model.

## 4.2 Unexpected problems

As mentioned, the deployment went smoothly from start to finish during the production week. There were no errors and both the app and predictive model were able to provide an answer to all the requests, either with a prediction or a detailed explanation of the error.

Nonetheless, it is known that the production deployment allowed the collection of 9694 new medical encounters, but their true outcome is known for all except 200. The issue that led to these 200 missing outcomes is suspected to be related to extra fields on the request that were not explicitly removed, such as `id` and `index`. These fields were indeed removed for the `/predict` endpoint, but not for the `/update` endpoint, thus leading to this issue. This minor issue was known during production week, but the risk of re-deploying at the time outweighed the benefits, since a re-deployment could impose permanent data loss and also cause a system unavailability.

Furthermore, it is worth noting a small issue imposed by the `heroku` platform. At the middle production week, `heroku` relocated the server where the application was running as detailed by the logs on figure 6.

```
2022-02-09T13:45:18.192202+00:00 heroku[router]: at=info method=POST path="/predict" host=e6281d438ef244b3b42f775aae.herokuapp.com request_id=84b040fa-c60e-4b69-b86d-5d4a7d63b5b3 fwd="52.208.22.146" dyno=web.1 connect=0ms service=104ms status=200 bytes=172 protocol=https
2022-02-09T13:46:02.197842+00:00 heroku[router]: at=info method=POST path="/predict" host=e6281d438ef244b3b42f775aae.herokuapp.com request_id=50cd5650-c750-4d78-9faf-6add43d88c41 fwd="52.208.22.146" dyno=web.1 connect=0ms service=211ms status=200 bytes=172 protocol=https
2022-02-09T13:46:24.125698+00:00 heroku[web.1]: Relocating dyno to a new server
2022-02-09T13:46:24.125698+00:00 heroku[web.1]: State changed from up to down
2022-02-09T13:46:24.229539+00:00 heroku[web.1]: State changed from down to starting
2022-02-09T13:46:46.634709+00:00 heroku[web.1]: Starting process with command `bin/sh -c /bin/bash\ -c\ \'unicorn\ app:app\'`
2022-02-09T13:46:48.504742+00:00 app[web.1]: [2022-02-09 13:46:48 +0000] [5] [INFO] Starting gunicorn 20.0.0
2022-02-09T13:46:48.505281+00:00 app[web.1]: [2022-02-09 13:46:48 +0000] [5] [INFO] Listening at: http://0.0.0.0:54413 (5)
2022-02-09T13:46:48.505439+00:00 app[web.1]: [2022-02-09 13:46:48 +0000] [5] [INFO] Using worker: sync
2022-02-09T13:46:48.526524+00:00 app[web.1]: [2022-02-09 13:46:48 +0000] [8] [INFO] Booting worker with pid: 8
2022-02-09T13:46:48.967890+00:00 heroku[web.1]: State changed from starting to up
2022-02-09T13:46:59.058159+00:00 heroku[web.1]: Stopping all processes with SIGTERM
2022-02-09T13:47:00.434356+00:00 app[web.1]: [2022-02-09 13:47:00 +0000] [5] [INFO] Handling signal: term
2022-02-09T13:47:00.434360+00:00 app[web.1]: [2022-02-09 13:47:00 +0000] [8] [INFO] Worker exiting (pid: 8)
2022-02-09T13:47:00.735883+00:00 heroku[web.1]: Process exited with status 143
2022-02-09T13:47:27.044795+00:00 heroku[router]: at=info method=POST path="/predict" host=e6281d438ef244b3b42f775aae.herokuapp.com request_id=6c29a823-3047-4c0e-8642-2177e18b4518 fwd="52.208.22.146" dyno=web.1 connect=0ms service=107ms status=200 bytes=173 protocol=https
2022-02-09T13:48:10.100357+00:00 heroku[router]: at=info method=POST path="/predict" host=e6281d438ef244b3b42f775aae.herokuapp.com request_id=01583019-5b94-489e-b120-b1fed91e4706 fwd="52.208.22.146" dyno=web.1 connect=0ms service=296ms status=200 bytes=173 protocol=https
```

Figure 6: Logs of heroku dyno relocation to a new server

This issue does not seem to have impacted the application, however one can note that it took 1 min and 30 s to complete the whole application relocation. Knowing the rate of requests was roughly one per minute, it is possible that some data was lost in the process.

This issue was imposed by the `heroku` platform, and as highlighted on the first report this kind of issues were out of the Awkward Problem Solutions™ control.

Other than that, no unexpected error occurred during the whole week in production.

## 4.3 Learnings and Future Improvements

The application deployment was perhaps the part of the current project that was better executed. Some of the key takeaways that revealed important are related to the overall app's robustness.

First of all, it is highly important to develop a broad model pipeline capable of receiving all types of data, specially `NaN`, and still output a result.

Then, the app itself has to be able to interpret the data that the REST API user provides. During this project an effort was made to interpret the data, discarding requests only as the last resort. For instance, when a field type is an integer, but the REST API user provides the string

*"11"* the app should be able to convert into the number 11, without discarding the request as invalid. Relying on the database to do this kind of heavy lifting was an important learning from this project. Since the database is, by definition, able to do this type of conversions leveraging this functionality improved the app robustness and increased the development speed.

Moreover, an important feature that was overlooked during the preparation for the production week is access to logs. This was a known limitation from `heroku`, but the information it provides is highly valuable and should never be missed in a production application.

Given the limitation, one alternative would be adding a second table with the triggered exceptions. However, this alternative was not viable for the current project, since it imposed a data loss risk due to the 10 000 rows limitation from `heroku`. A different alternative would be creating a local log file to be accessed later.

Finally, `DBeaver` revealed to be an important tool during the deployment and production stage. It allowed accessing the remote database hosted by `heroku` and proved to be a valuable monitoring tool. On top of that, it made the extraction of the database into a `csv` file seamless and quick.