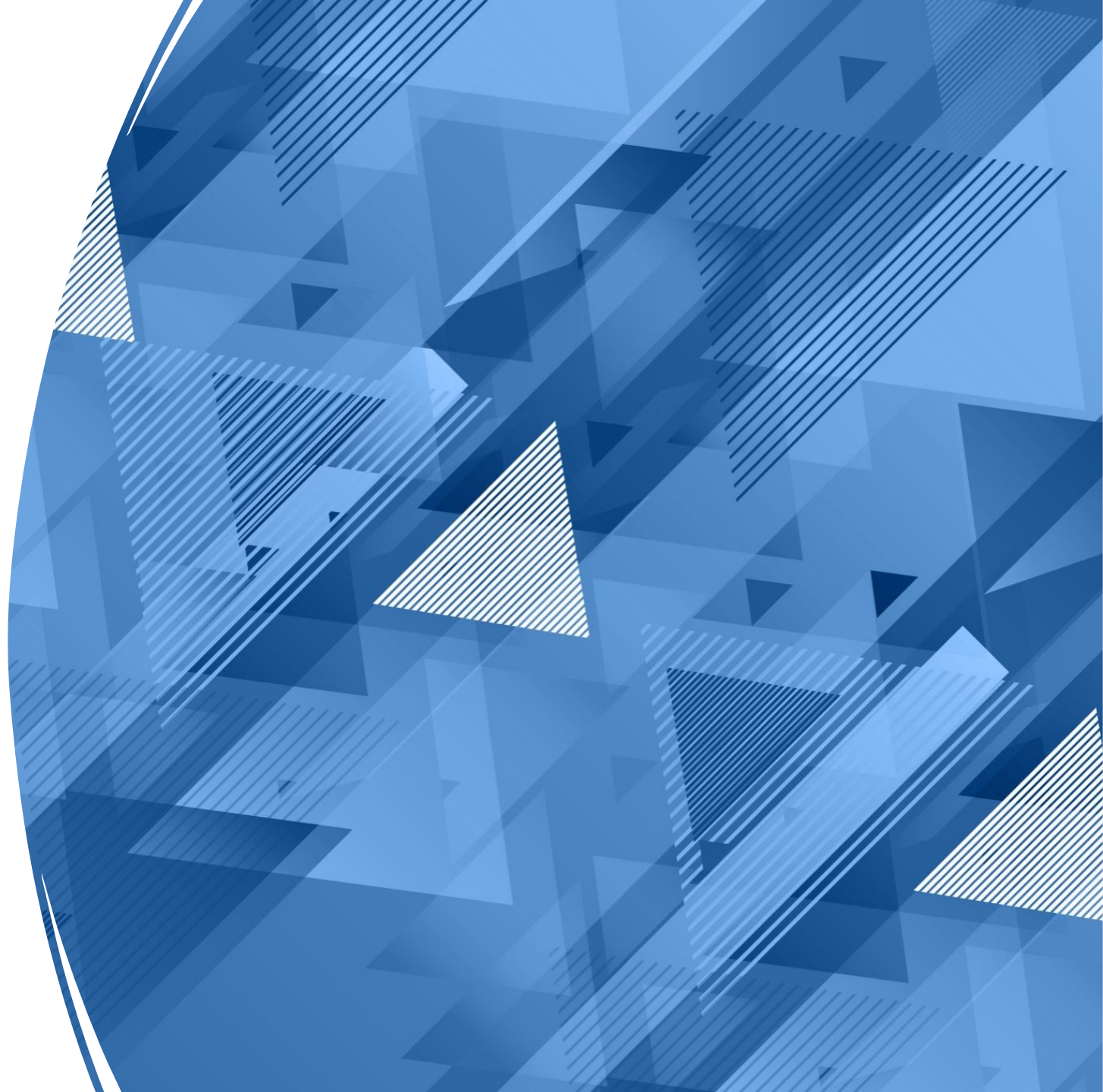


# Healthcare Dataset Data Science Specialisation

---

Exploratory Data  
Analysis  
And Modelling



# Data Description

---

- ABC company (United States) delegated the investigation of persistence of an NTM (non-tuberculous mycobacteria) drug (antibiotic) as a function of various factors
- 3242 participants with data from 67 features including:
  - Patient demographics: Age, gender, ethnicity, etc.
  - Risk factors: Smoking status, obesity, family history, etc.
  - Underlying conditions and comorbidities: Diabetes, hypertension, cardiovascular diseases, etc.
  - Physician information: Specialist or general practitioner, experience level, etc.
- Data mainly categorical (mostly binary variables) and two ordinal variables - frequency of DEXA (bone density) scans and number of pre-existing or ongoing risks for each participant

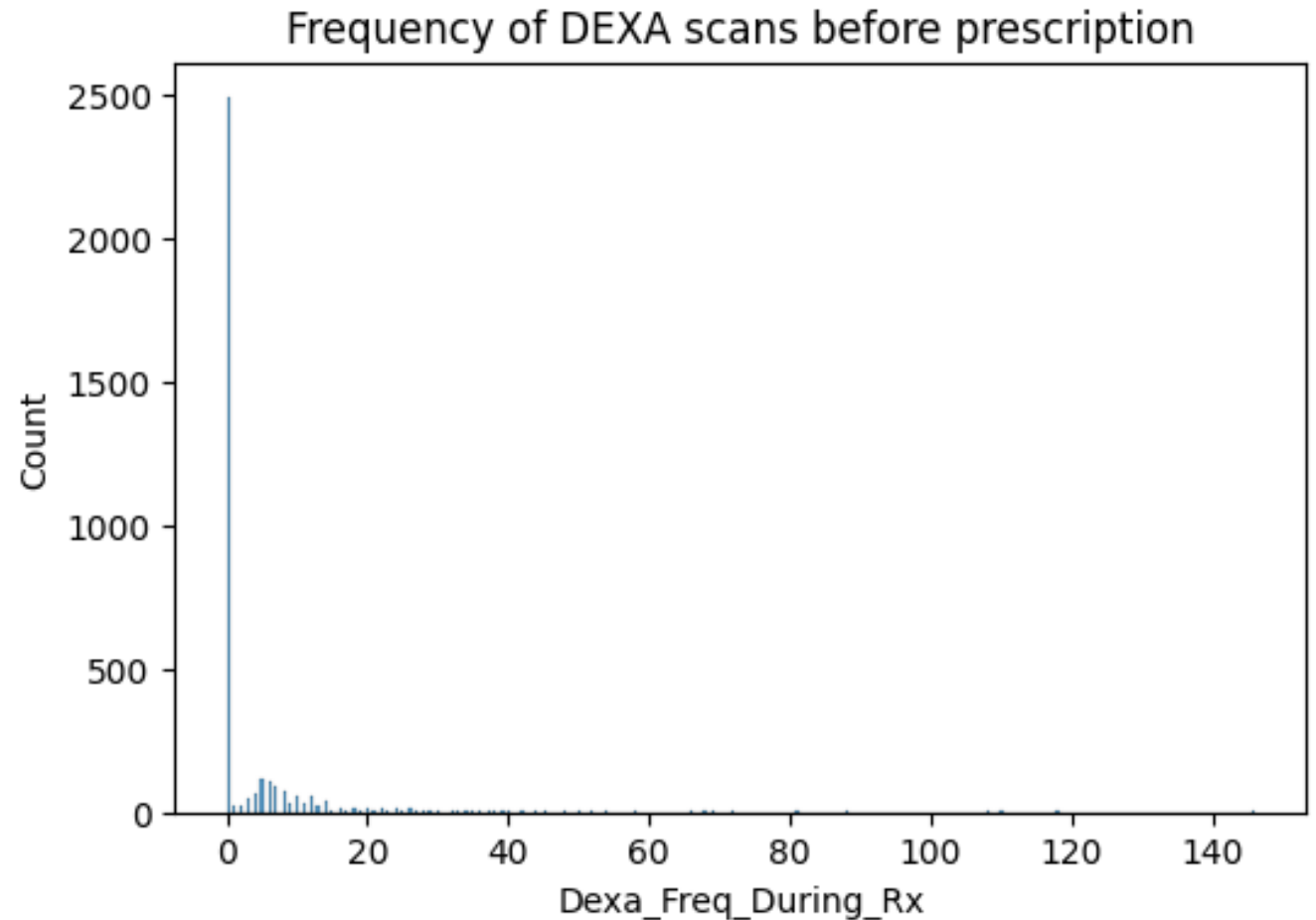
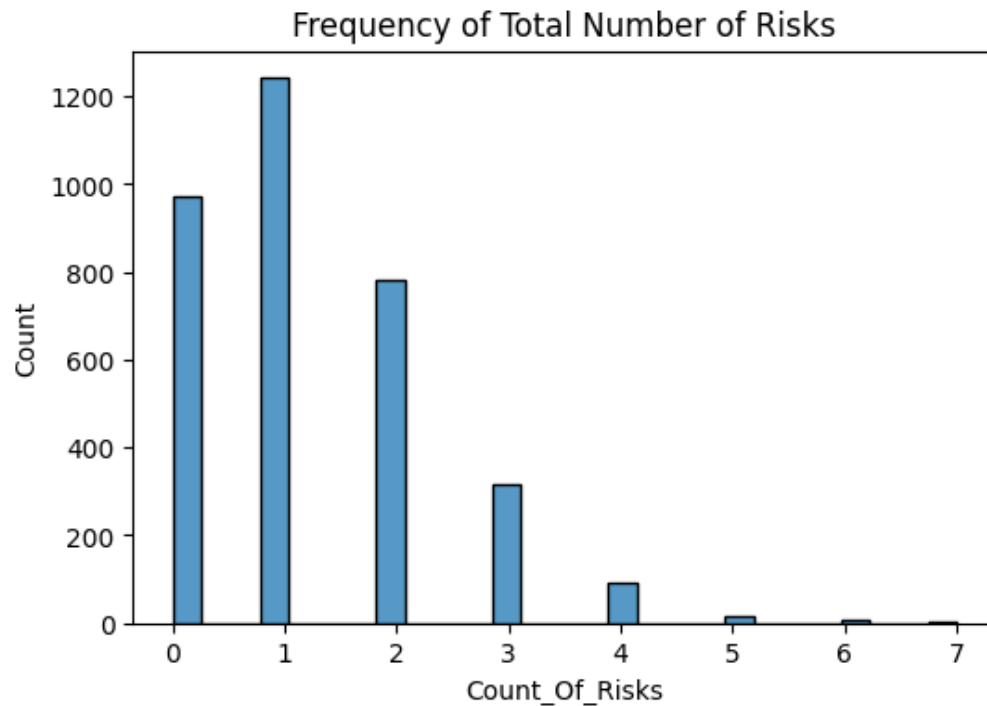
# Data Cleaning

---

- Data do not contain any NaN/null variables, or any duplicates
- Most categorical variables are imbalanced: some categories of a variable containing hundreds more data points than others
- For such data I removed the variables which presented at least one category that had fewer data points than a specific threshold
- The threshold was determined based on one of the variables used to measure general health: change in risk segment (from high risk to low risk and reverse)
- The risk segment measures the health severity of a person, and in this case also considers the number of serious medical conditions
- The 'Change in Risk Segment' variable contains four levels: Unknown (most patients, 2229), No change (1052 patients), Worsened (121) and Improved (22)
- The patients who reported a certain improvement in their risk segment are very few, but sufficient to have their data analysed even without oversampling it or using further alteration methods; thus the number of patients contained in this category represents the threshold for which I eliminated binary variables with categories containing fewer data points than this limit
- Seven such variables were removed as a result
- For ordinal variables I used common statistical elimination methods depending on their skewness: for severely skewed variables such as DEXA scan frequency, I used the flooring and capping method (replacing everything outside of the 25% and 75% quantiles with their corresponding values, give or take the value of the interquartile range, respectively); for the less skewed Count of Risks I replaced everything above the 99% quantile with this value (patients presenting 5, 6 and 7 risks - 23 in total - were newly noted as only presenting 4 risks)

# Skewness of Count of Risks and DEXA scan frequency as shown in histograms

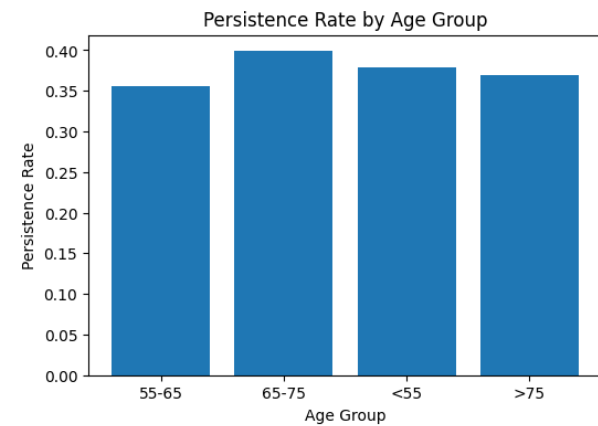
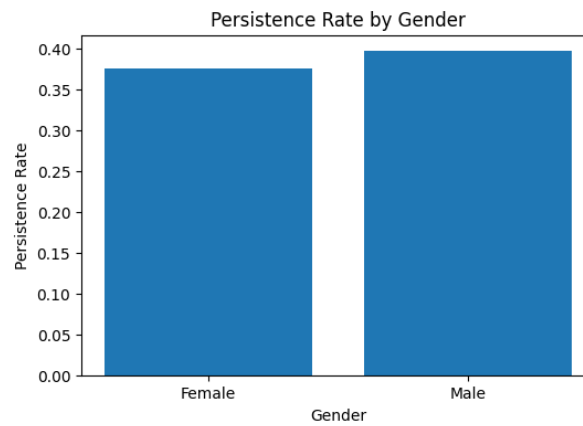
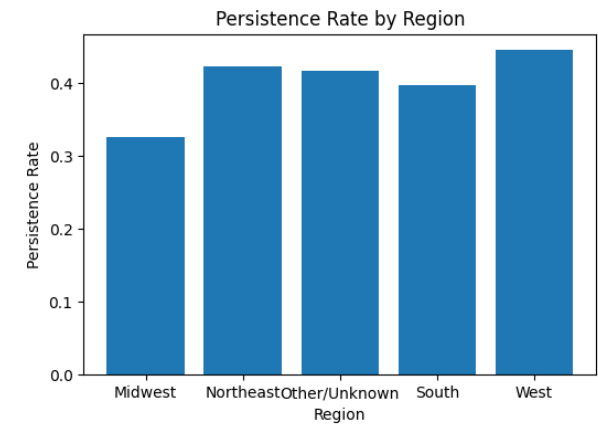
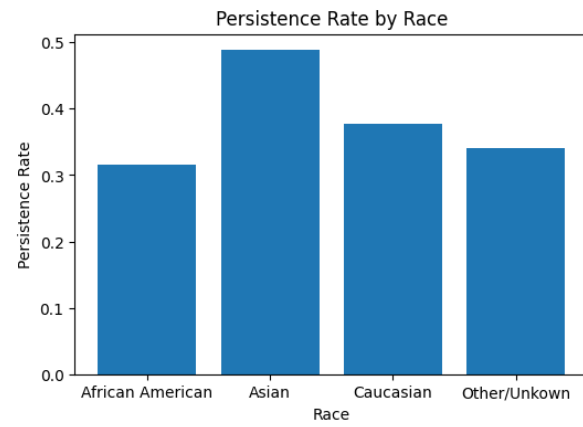
---



# Exploratory Data Analysis

---

- The challenges in performing exploratory data analysis on the present dataset was the presence of numerous variables, many of them referring to very specific health conditions; knowledge about the interaction of this drug's effects with each of them required high expertise – and it may have been inefficient to investigate further
- Some visualisations are presented for the usual factors: demographics, general information about the prescriber's medical specialisation, and general health information
- I am reporting rates (%) of persistence rather than total numbers of patients showing persistence; to generalise, stats were conducted on this relative persistence; this is because of the number imbalance between the levels of the reported variables. No significant differences were observed at such relative rates – whether this is true of the general population remains unknown, but significant statistical values would have been observed when looking at individual numbers.

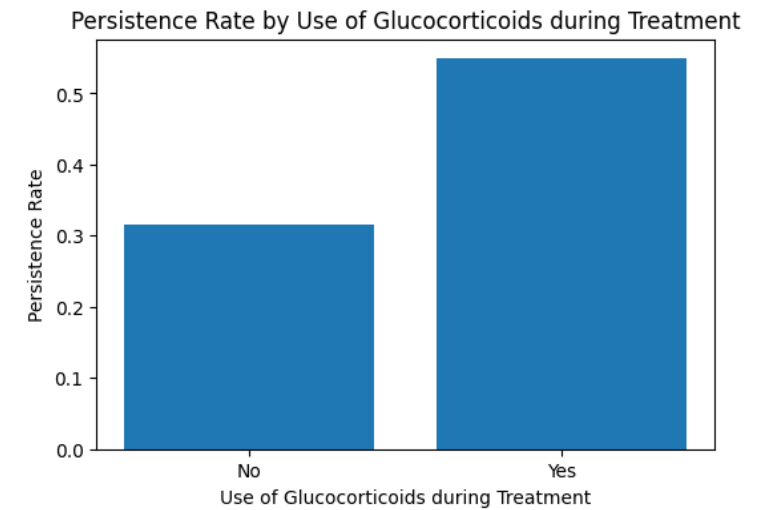
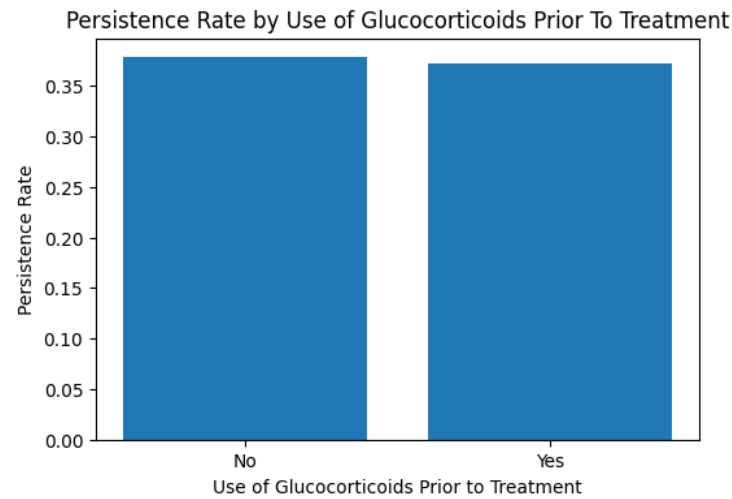
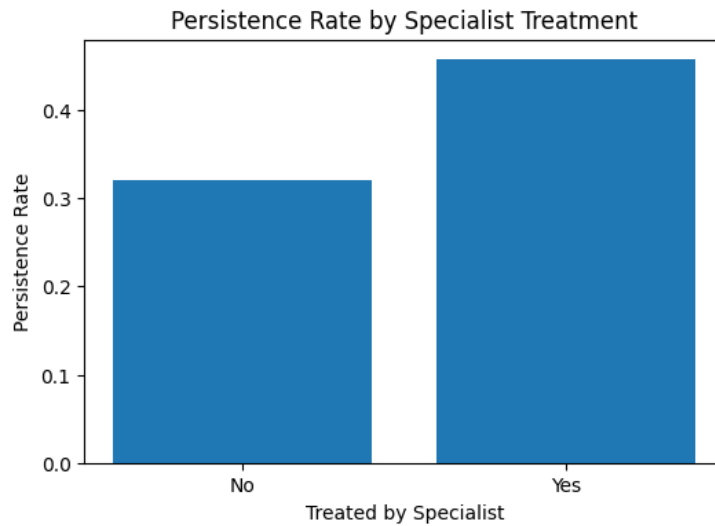


# Demographics

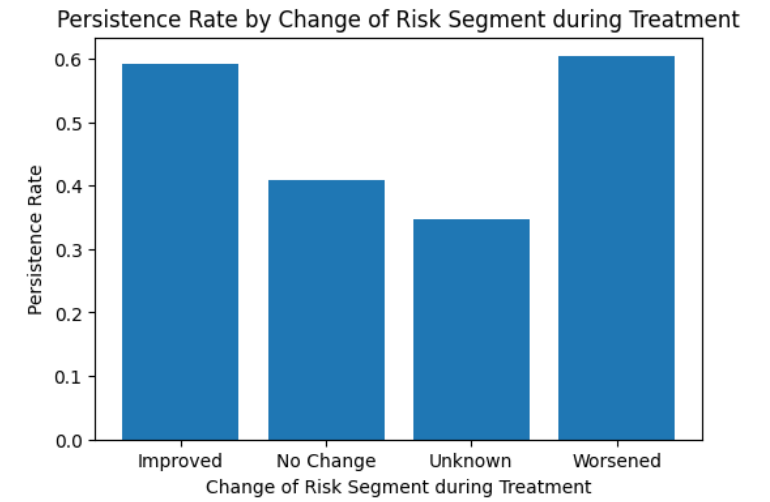
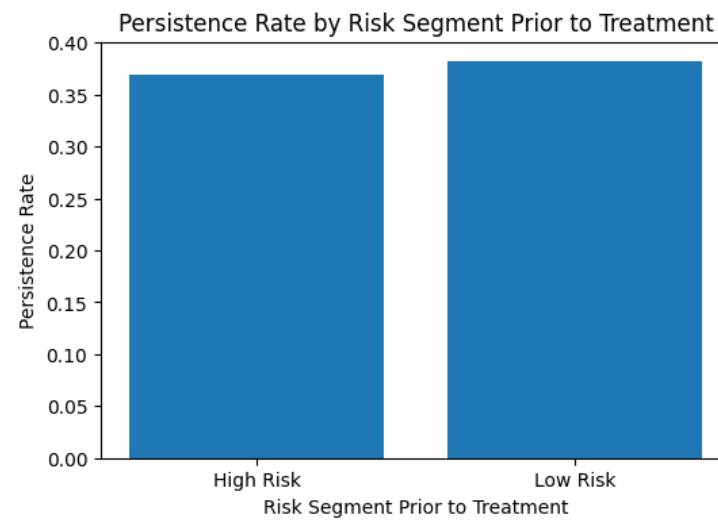
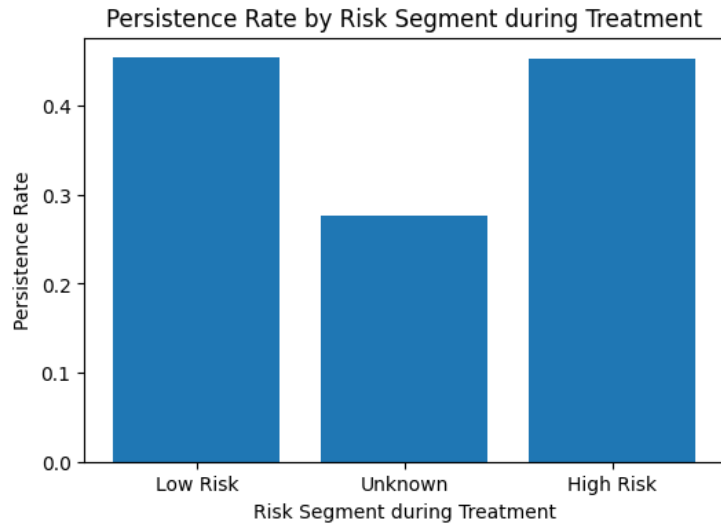
---

- No noticeable difference in drug persistence between males and females, or between age groups
- Asians show the most drug persistence, followed by the other groups
- As for region, Midwest shows the lowest drug persistence

- Persistence rates are higher if antibiotic prescribed by specialist than by GP; in terms of treatment effects, persistence rates higher if antibiotic taken together with glucocorticoids (anti-inflammatory), but not if these were taken before the commencement of the NTM treatment

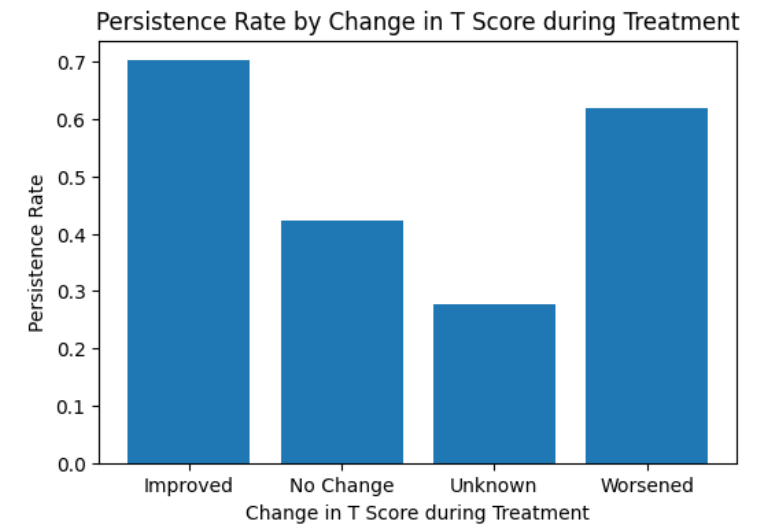
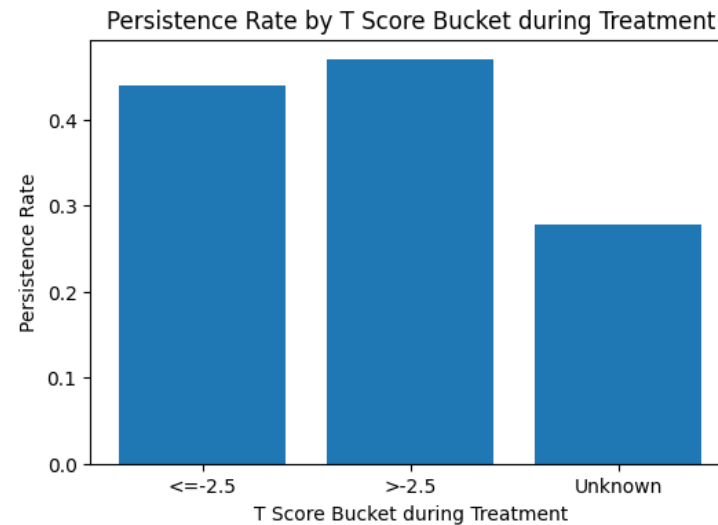
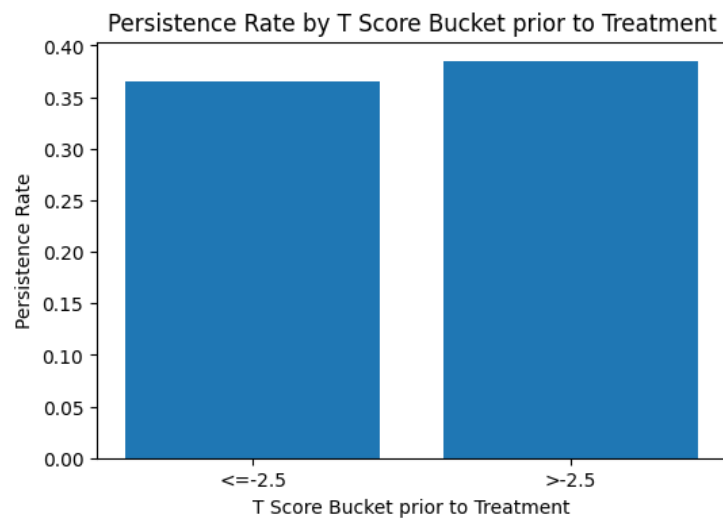




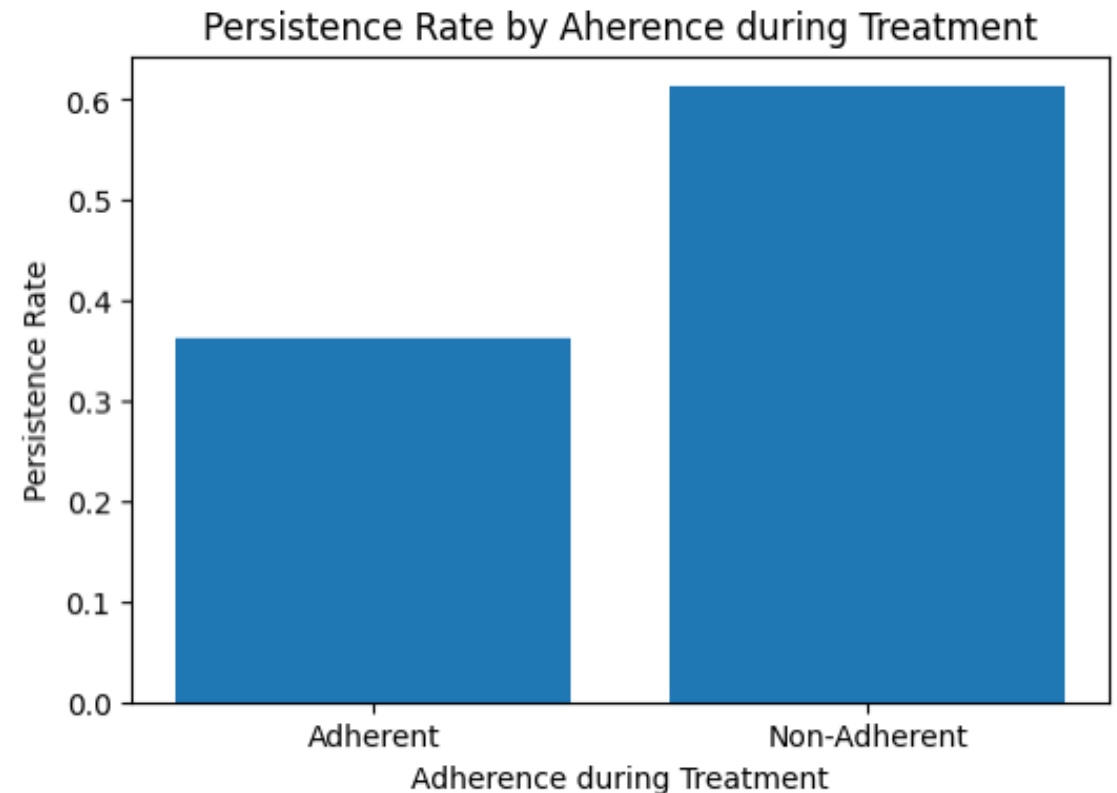
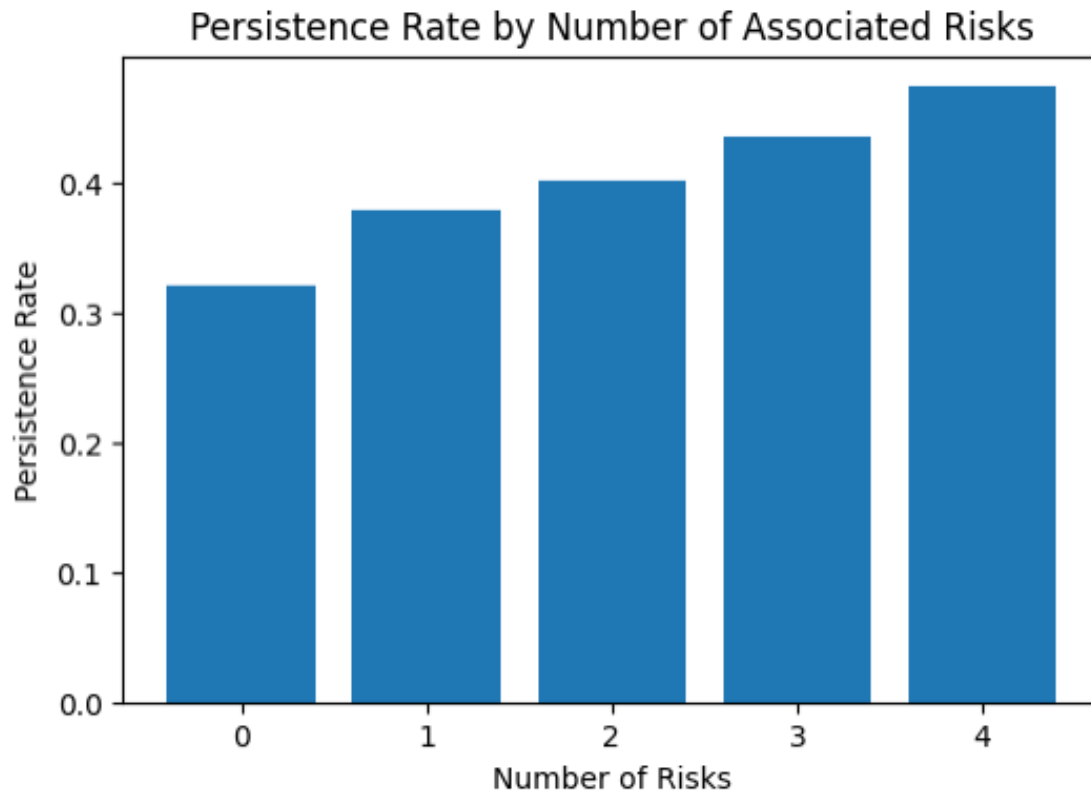


It seems persistence may increase if the risk segment changes – whether this means improving or worsening; but there is no significant difference and the numbers in these categories are quite small

# Similar results for change in T Score as for change in risk segment



Persistence rates are higher for patients who present risks for higher numbers of conditions, but are lower in patients who show less adherence to treatment (compliance to medical recommendations)



# Model recommendations

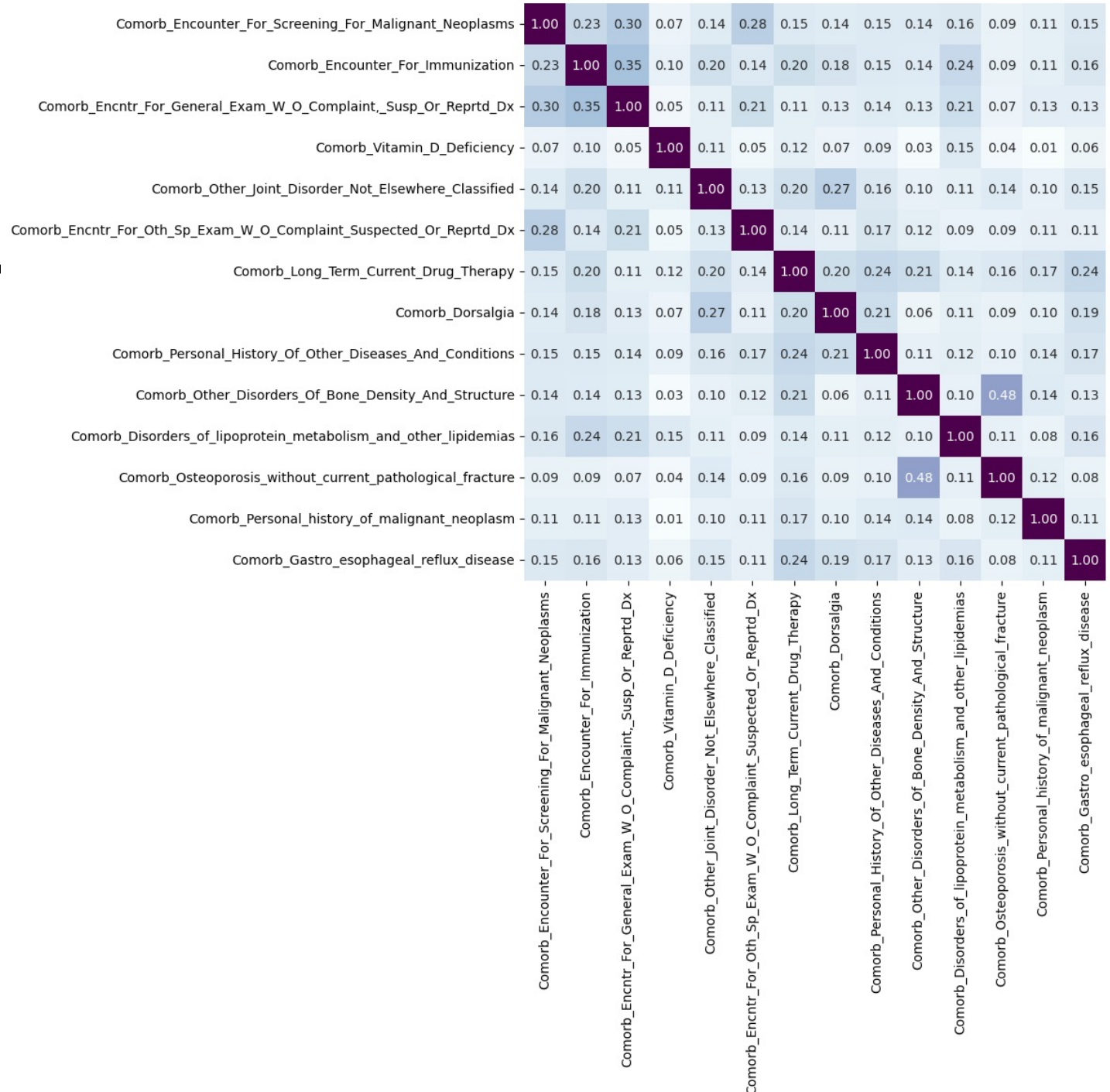
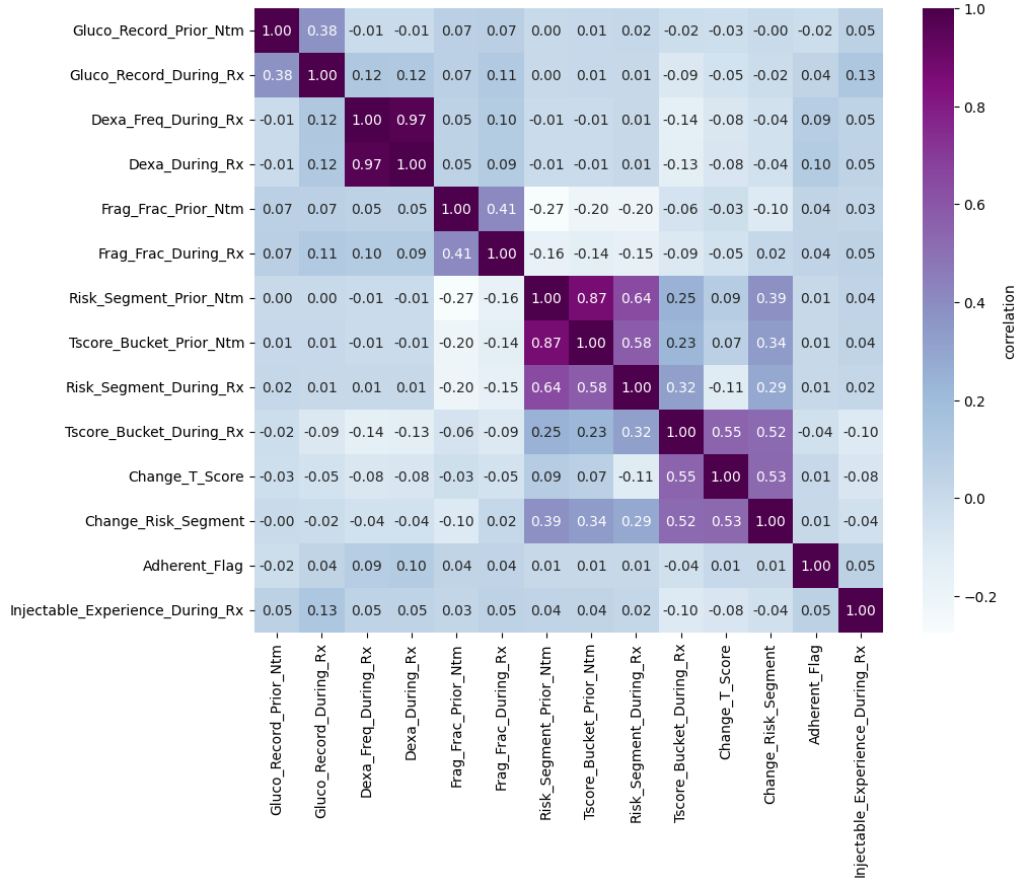
---

- Dimensionality reduction
- One way to start is by performing correlations on the different categories of variables
- Principal components analysis (PCA) can be used to remove correlated variables
- Recursive feature elimination can be used to remove features one by one and then check best model performance
- Can try oversampling for imbalanced data
- As a general model, logistic regression is recommended for a binary categorical outcome

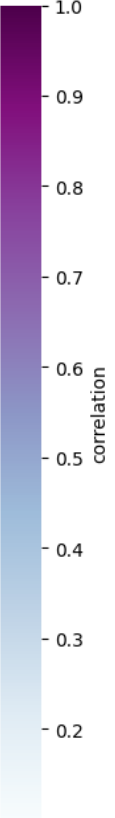
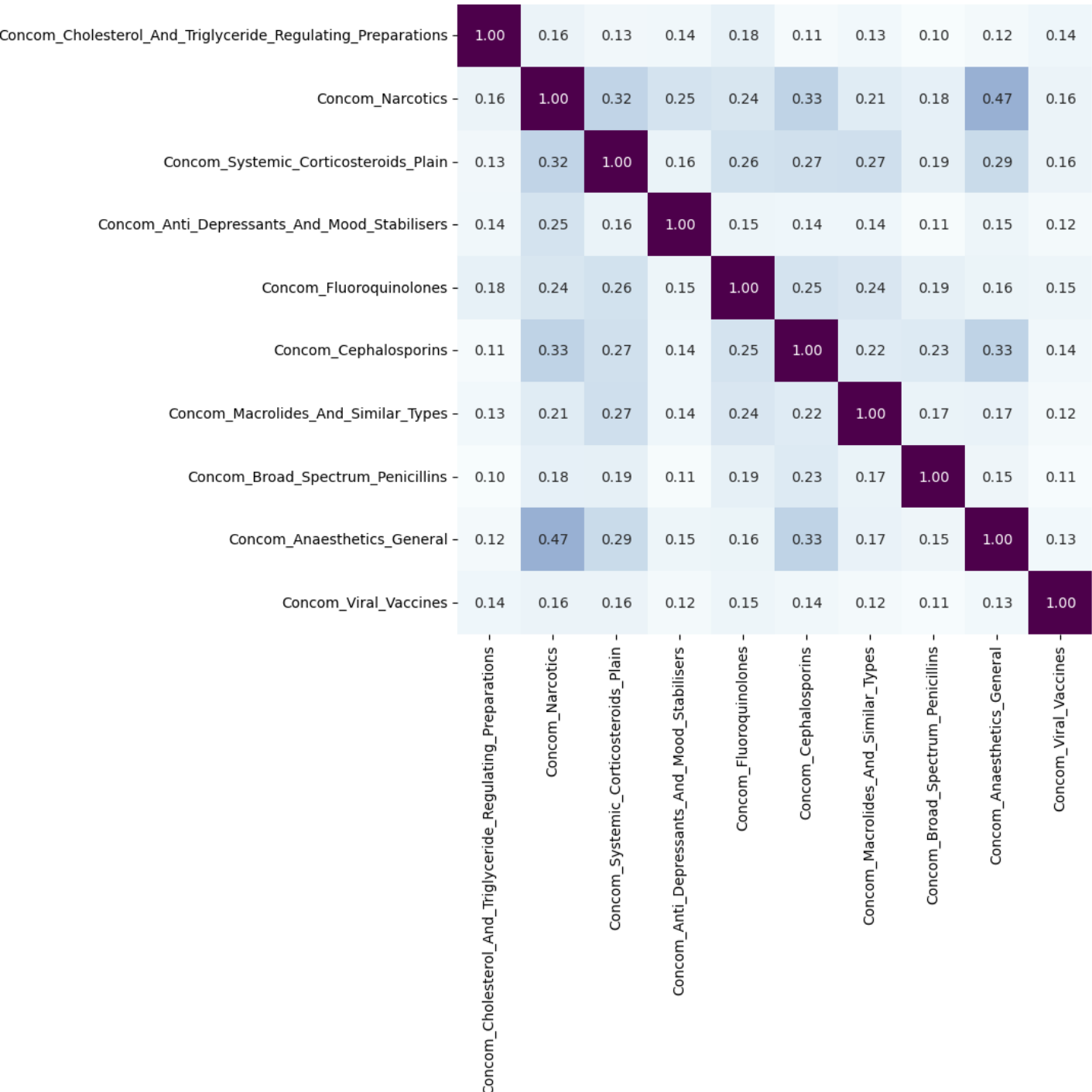
# Correlations

---

- For this I split the variables into their identified categories – mainly for visualisation purposes
- Correlations were observed, at each category:
  - Patient demographics: Not performed.
  - Risk factors: Medium correlations between the total number of risks and being at risk of Vitamin D insufficiency, the risk of smoking tobacco, and being at risk of chronic malnutrition or malabsorption, respectively – possibly indicating that being at risk for these three increases other health risks, but these relationships are outside the current scope
  - Concomitant conditions: Moderate correlation between the concomitant use of narcotics and anaesthetics
  - Comorbidities: Moderate correlation between disorders of bone density and structure and osteoporosis without current pathological fracture
  - Clinical Factors: High correlations between the presence of DEXA scans during treatment and the number of scans; the risk and T scores segments before treatment (but only small during); a relatively high correlation between risks segments before and during treatment; moderate correlation between T score prior to treatment and risk segment during treatment; and medium relationships between the T score, the change in T score and the change in risk segment during treatment.

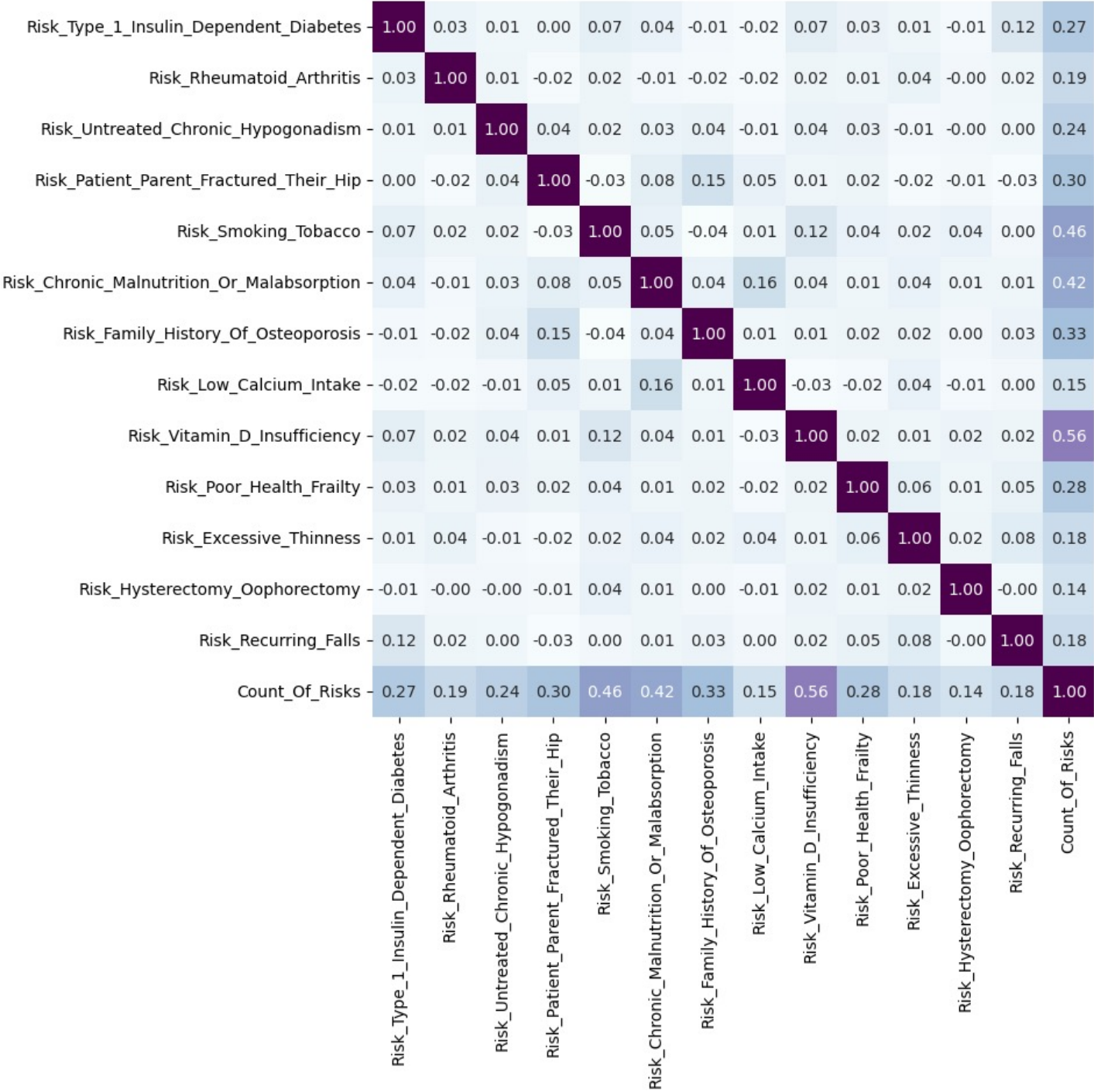


\_\_\_\_\_



\_\_\_\_\_

\_\_\_\_\_



\_\_\_\_\_

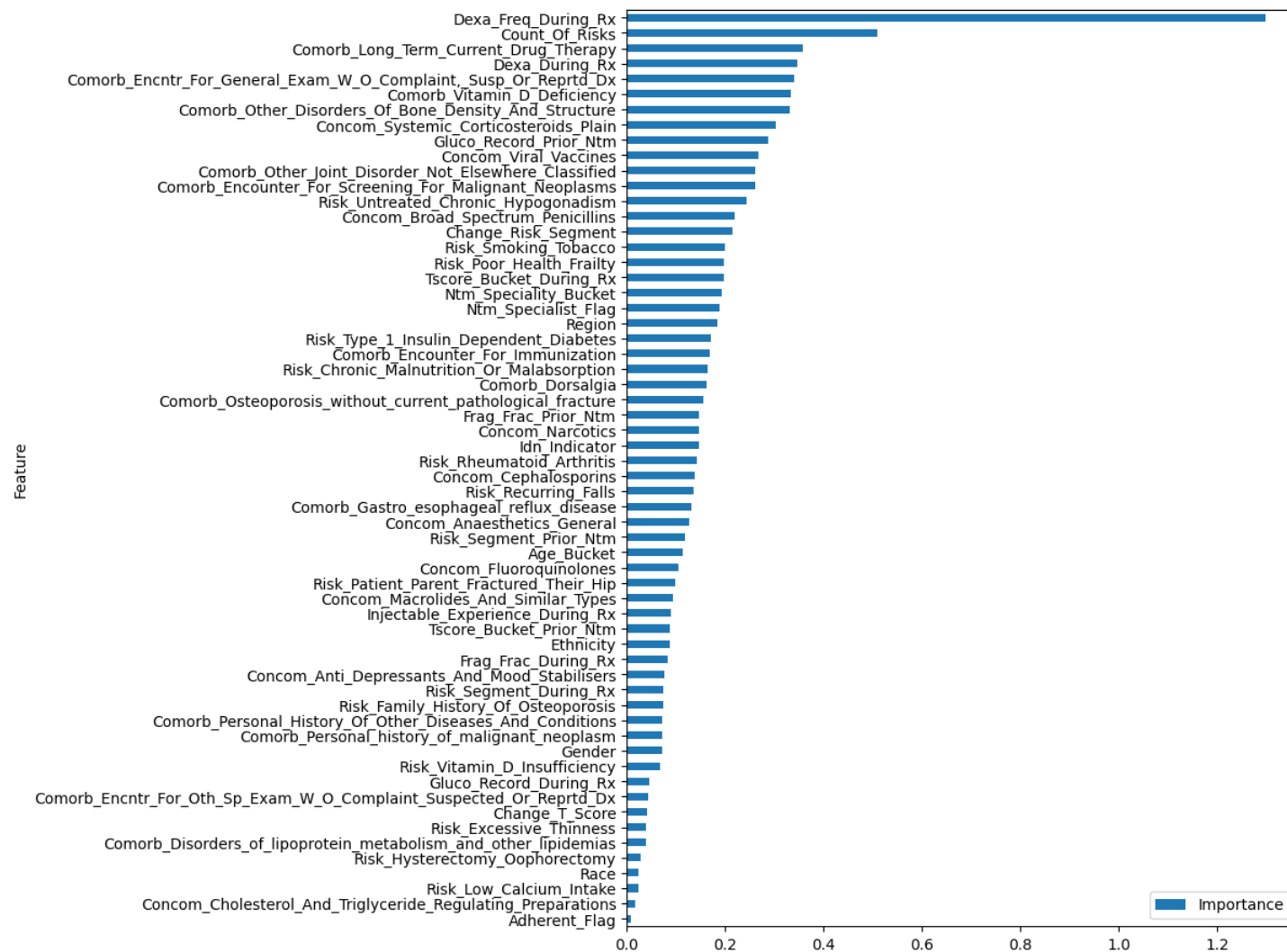


## Modelling – logistic regression

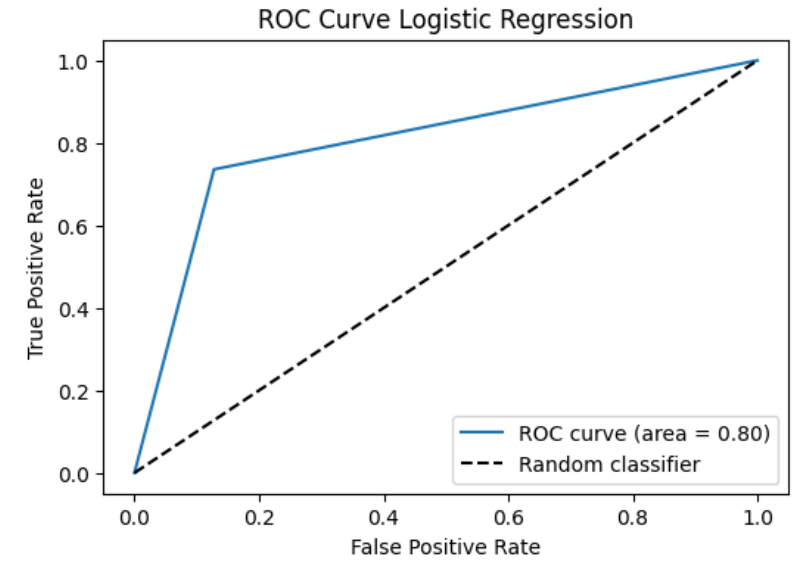
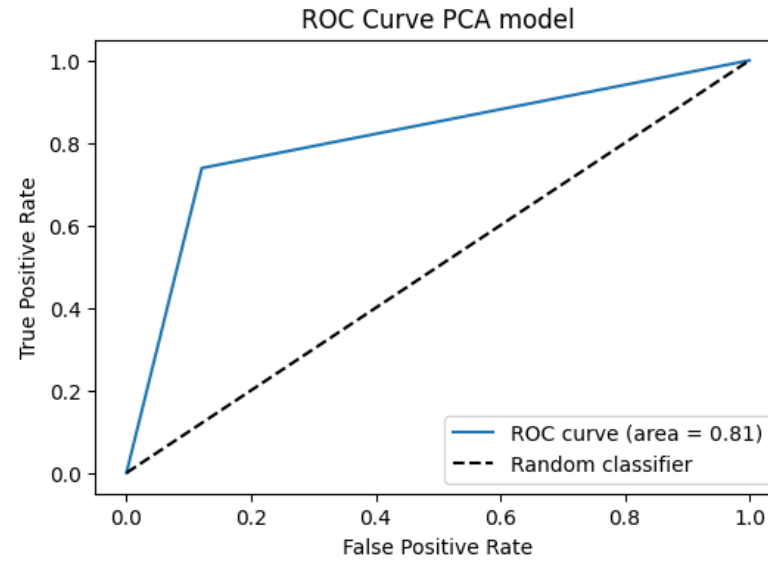
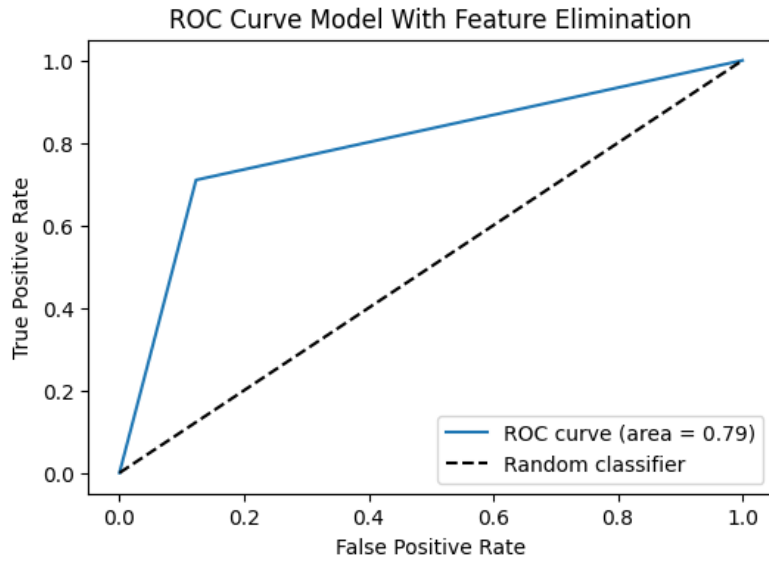
---

- Logistic regression performed on the cleaned dataset showed between 79 and 84 % accuracy with each run
- So did models on subsequently PCA-ed data (with optimal number of components extracted) or data further reduced using recursive feature elimination
- Oversampling did not lead to a more balanced dataset so was not considered for further modelling

- Model feature importance (normalised data)
- From logistic regression without automated feature elimination



# ROC curves for model variations



# Conclusion

- Logistic regression was a reasonable chosen model
- All model variations performed equally well, suggesting that data reduction for this particular data set/model was not needed

# Future recommendations

---

- Try same model with even less features ( $<10$ )
- Try a different model
- Perform different models for different kinds of data instead of trying to merge them into one (different ones for demographics, risks, etc.)
- Try different ways to identify meaningful and important variables - perhaps by looking at existing scientific results (this could be done either manually or automatically depending on resources)

End

---

