



Statistical Machine Learning: Platelet Count and Strokes

By Emily Neuman, Marshall Roll, Kenny Nhan



Goals and Methods

1

Regression

OLS and LASSO modeling to **predict platelet levels**, which are related to stroke events

2

Classification

Logistic LASSO and a decision tree to **predict stroke events**

3

Clustering

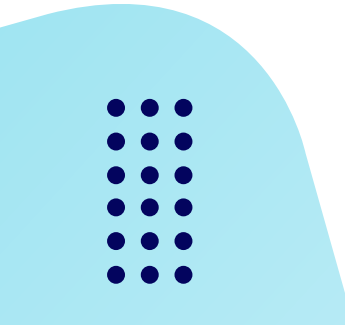
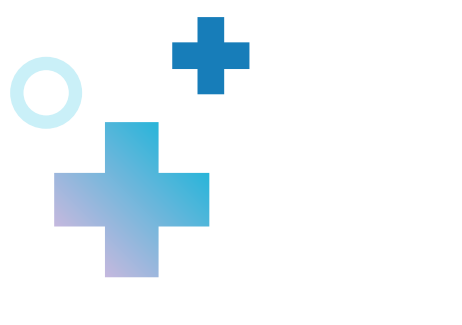
Hierarchical clustering to **uncover patterns underlying stroke events**





01

Platelets

- Age
 - Creatinine Phosphokinase
 - Ejection Fraction
 - **Platelets**
 - Serum Creatinine
 - Serum Sodium
 - Diabetes
 - High Blood Pressure
 - Sex
 - Smoking
 - Anaemia
- 
- 




Regression




Research Question

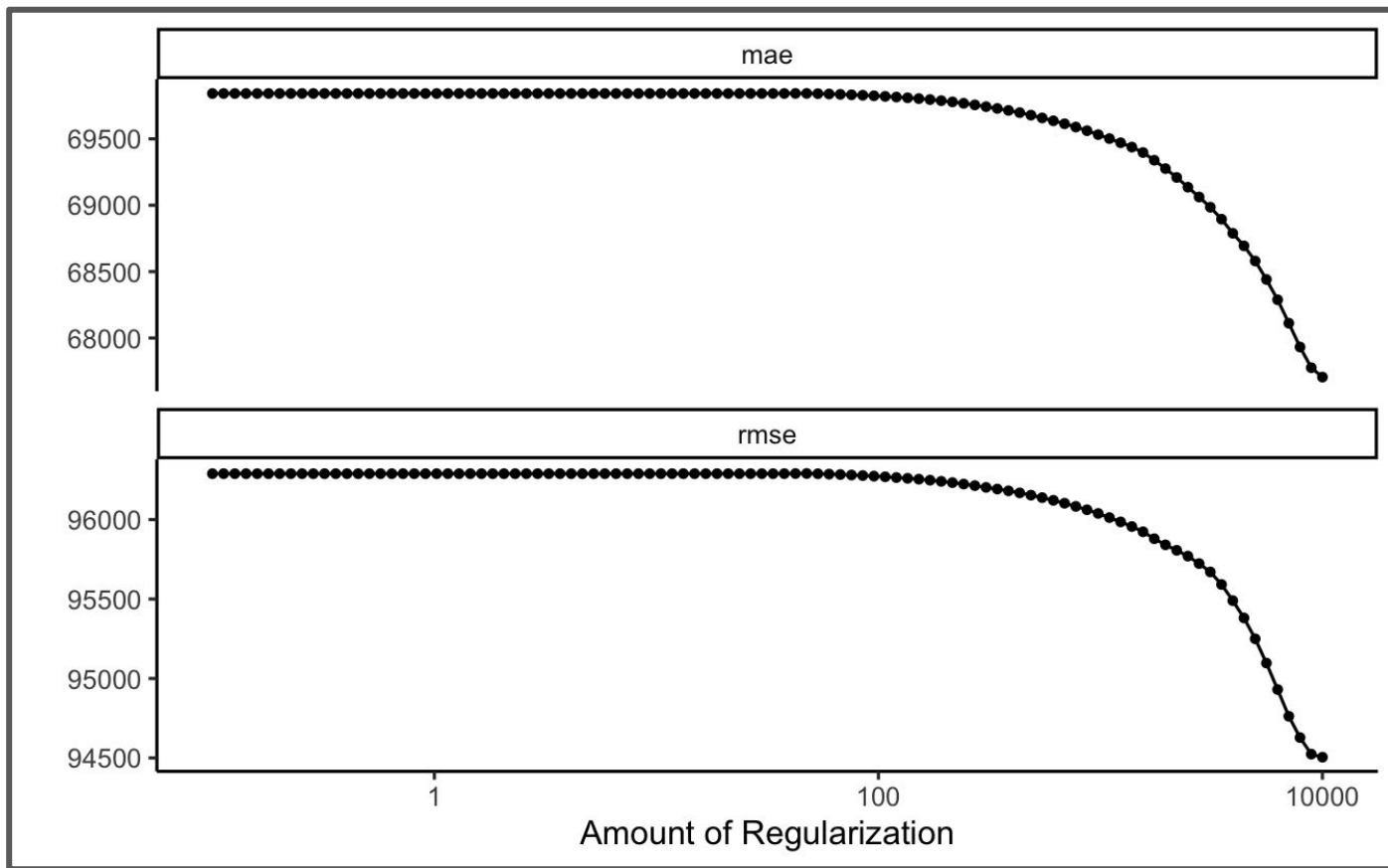
Which of the available biological factors are most accurate in **predicting platelet levels**, which are related to stroke events?





Model Type	RMSE	Ease of Interpretability
OLS	96307	Low
LASSO	95305	High
GAM	95478	Low





LASSO Takeaways

- Model shrinks all variables to zero except for **biological sex**
- All linear regression models **perform weakly**, but the LASSO is the most interpretable
- **More biological predictors** needed to determine platelet levels



Classification



Research Question

Can we predict whether or not a patient will have a stroke based on certain characteristics?





02

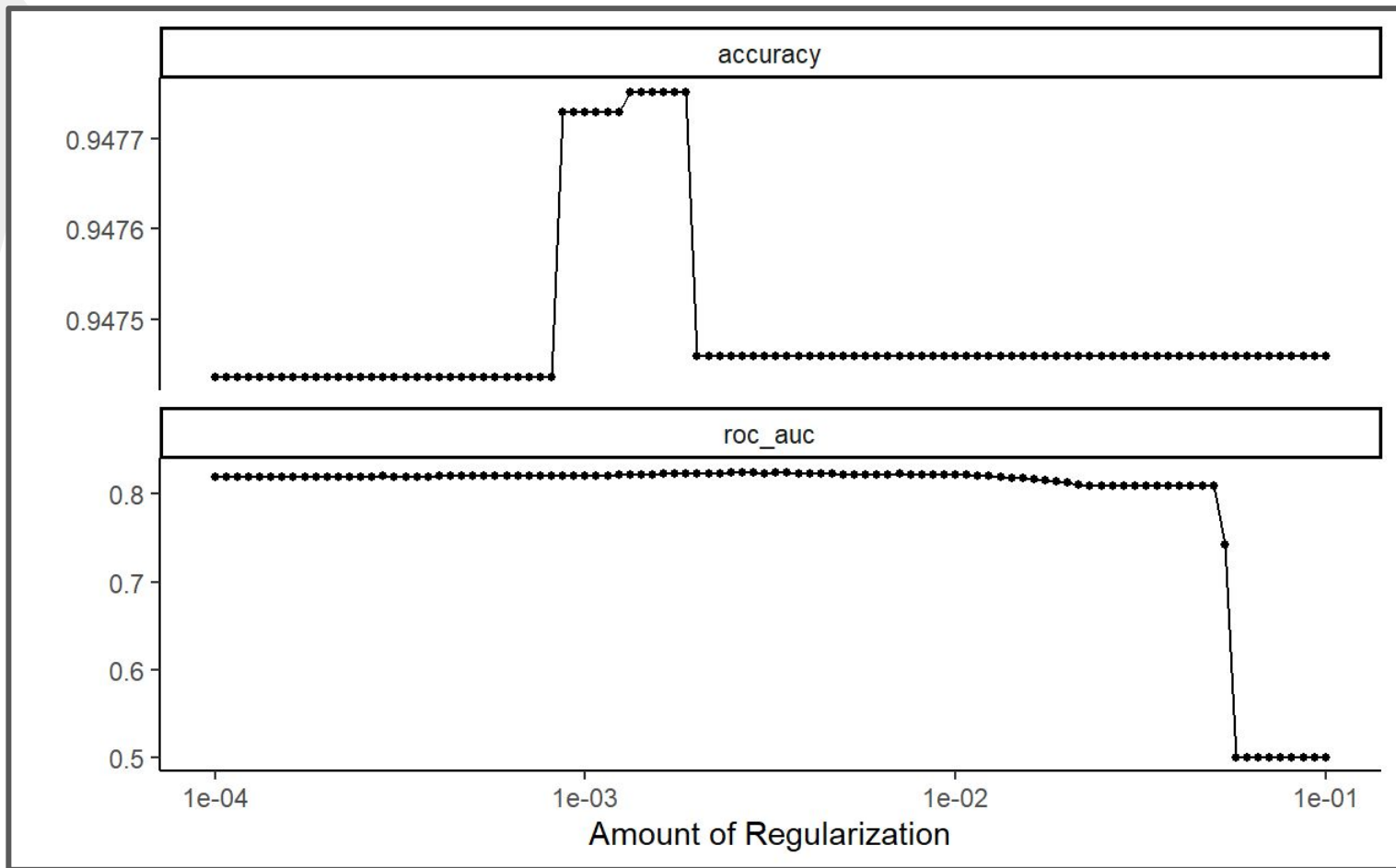
Stroke

- Age
- Average Glucose Level
- BMI
- Gender
- Hypertension
- Heart Disease
- Ever Married
- Work Type
- Residence Type
- Smoking Status
- Stroke



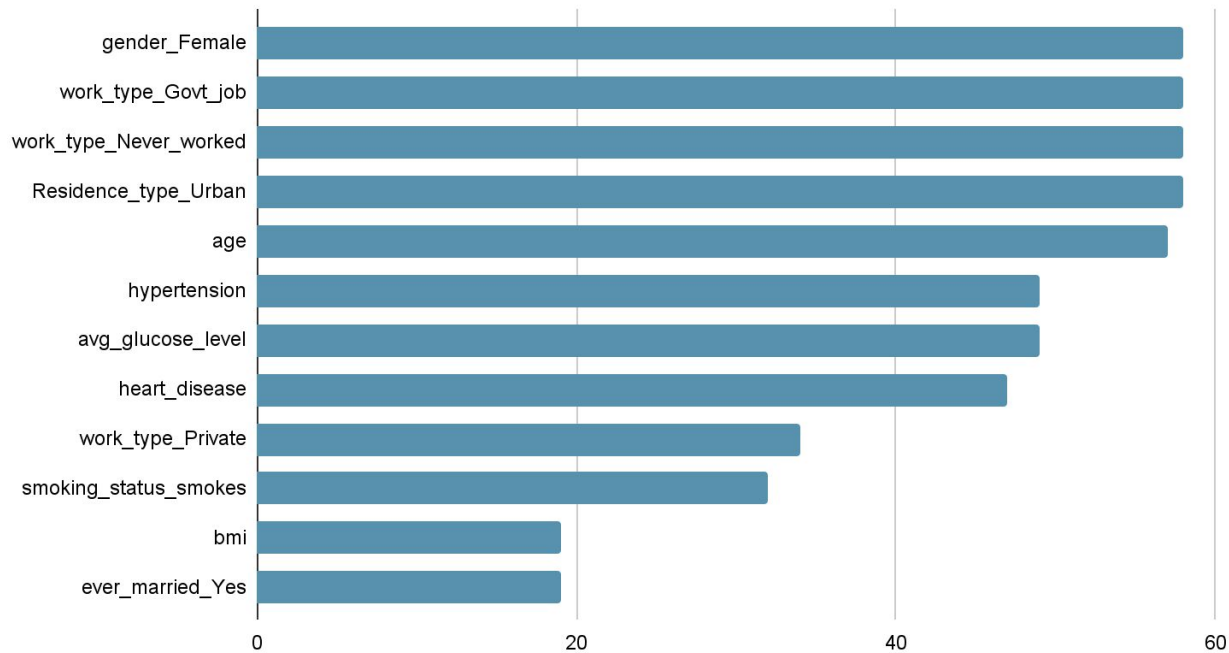


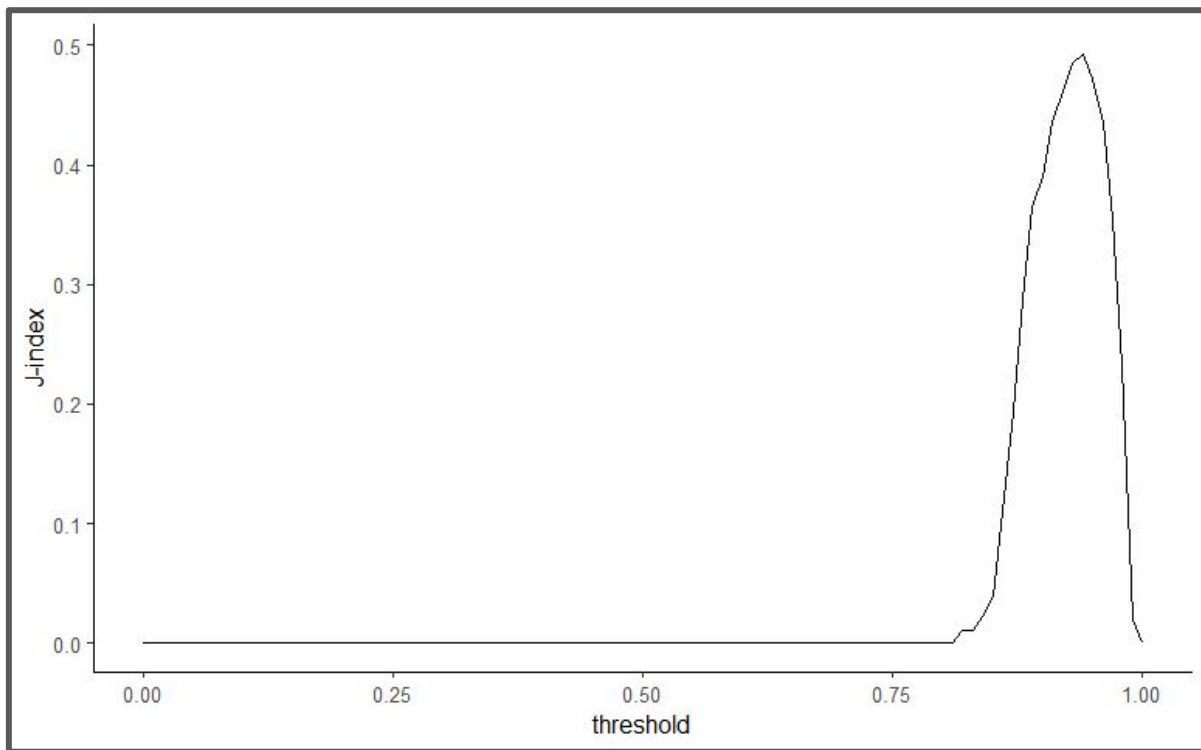
Logistic LASSO



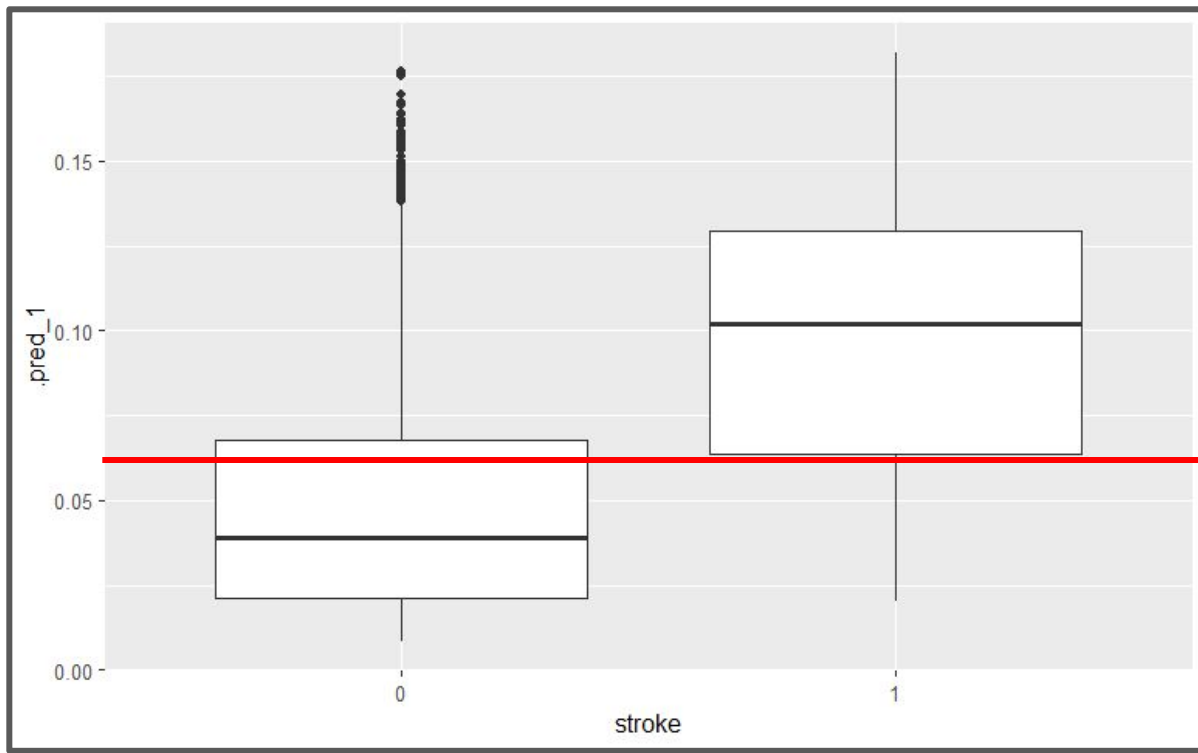


Variable Importance





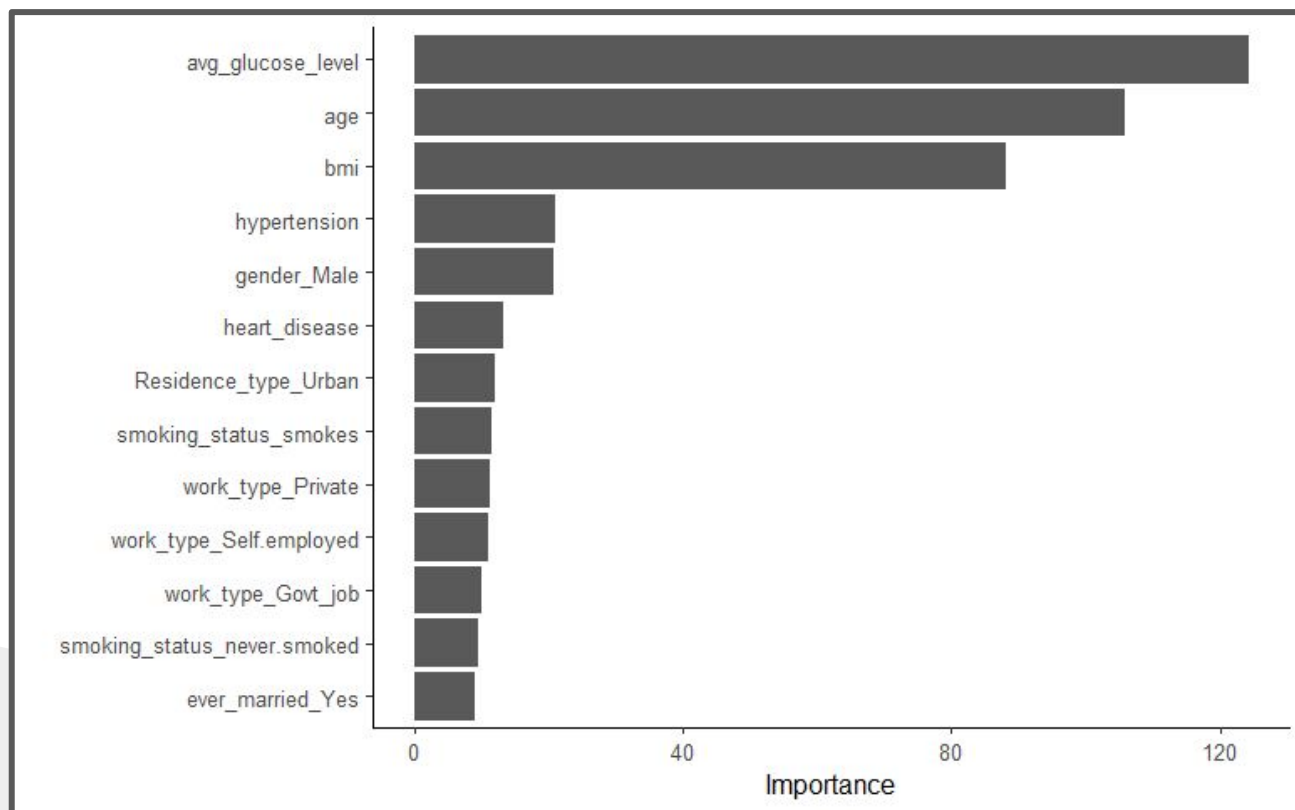
Accuracy	0.70
Sensitivity	0.78
Specificity	0.70



Threshold: 0.94



Decision Tree



Classification Takeaways

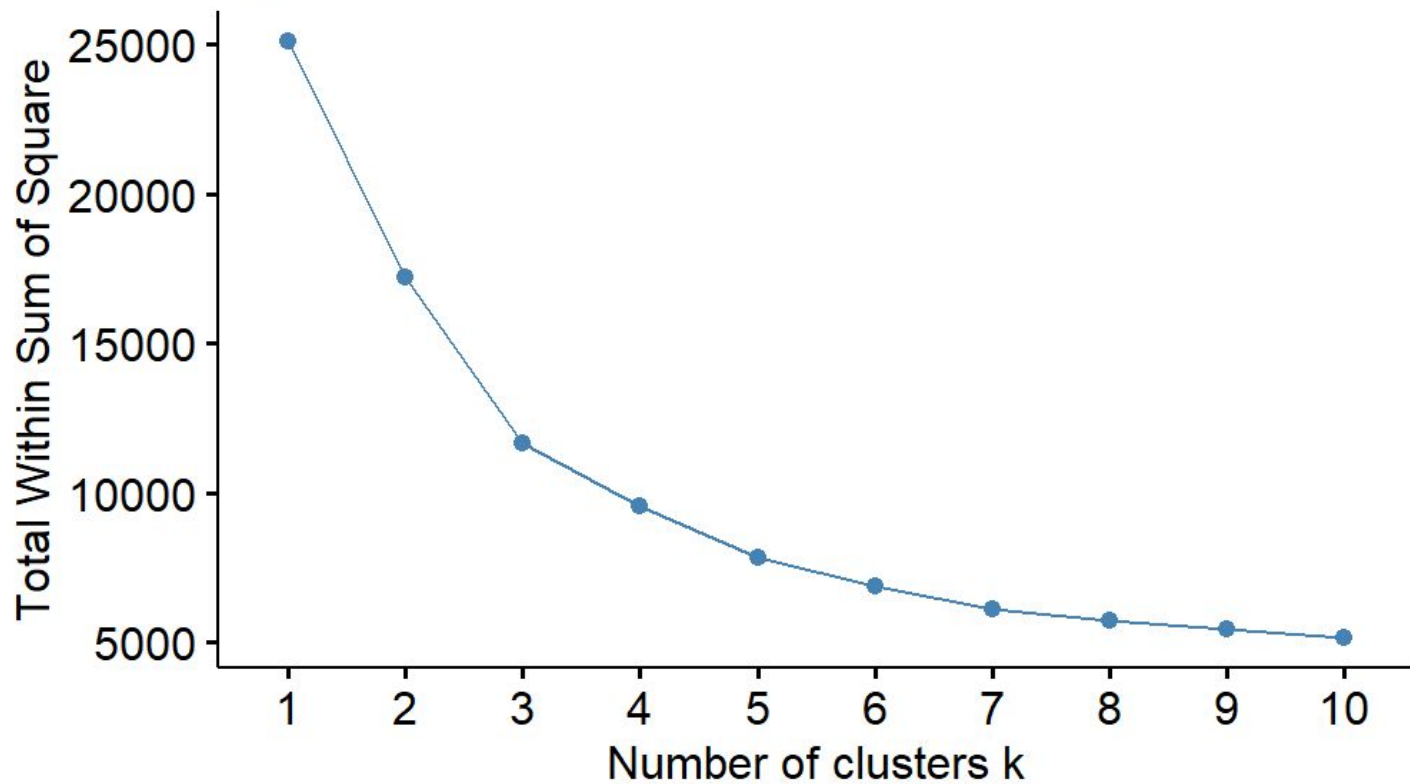
- **Logistic LASSO is a stronger model** because it is more interpretable and gives clearer metrics of variable importance
- Preferred thresholds with **higher specificity** while maintaining a reasonable sensitivity and accuracy
- Biological women, age, hypertension, working a government job, and unemployment are the strongest predictors



Clustering

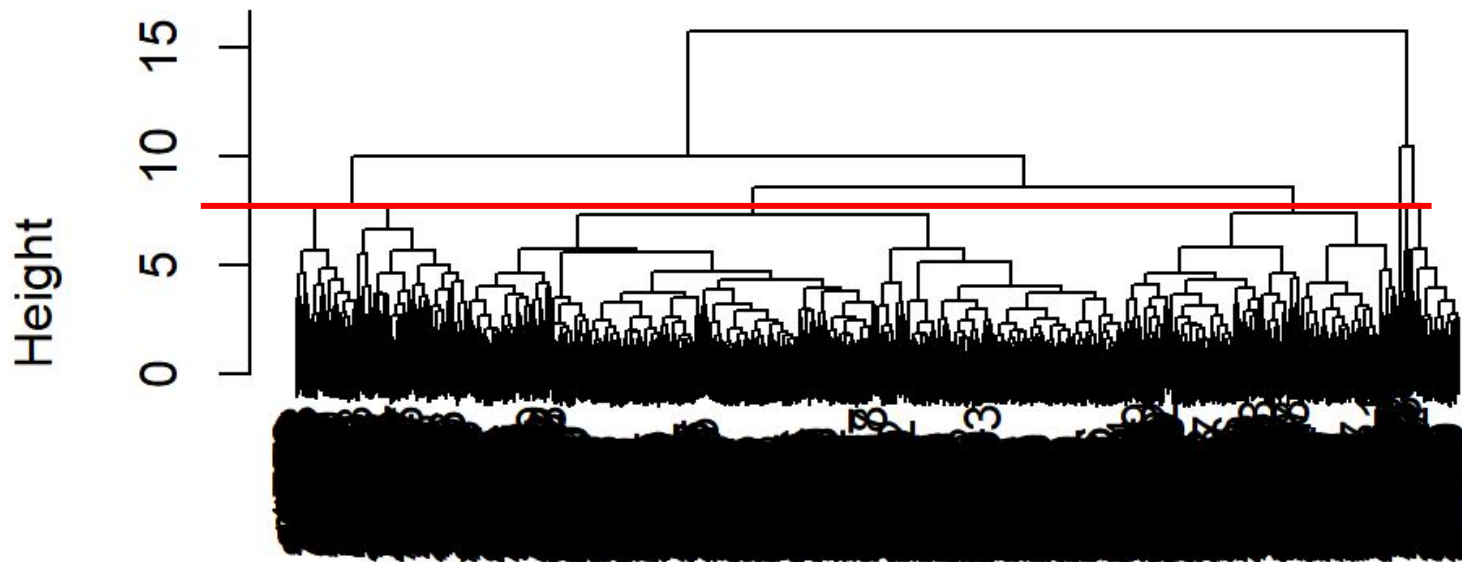
- Uncover **underlying patterns** determining whether a patient is at risk for a stroke
- Hierarchical clustering using complete linkage
- Created **6 clusters** cutting at a height of ~12.5

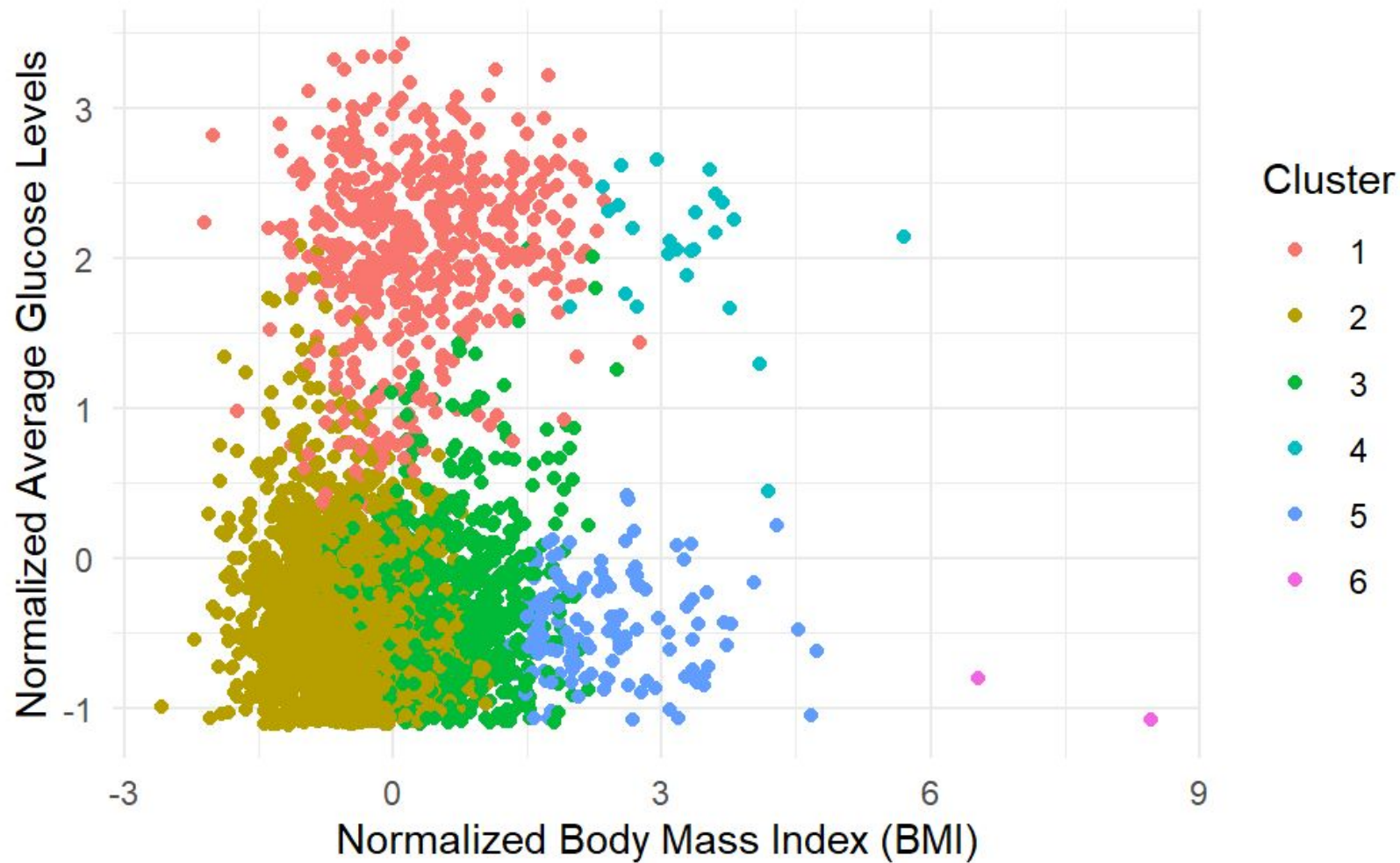
Optimal number of clusters

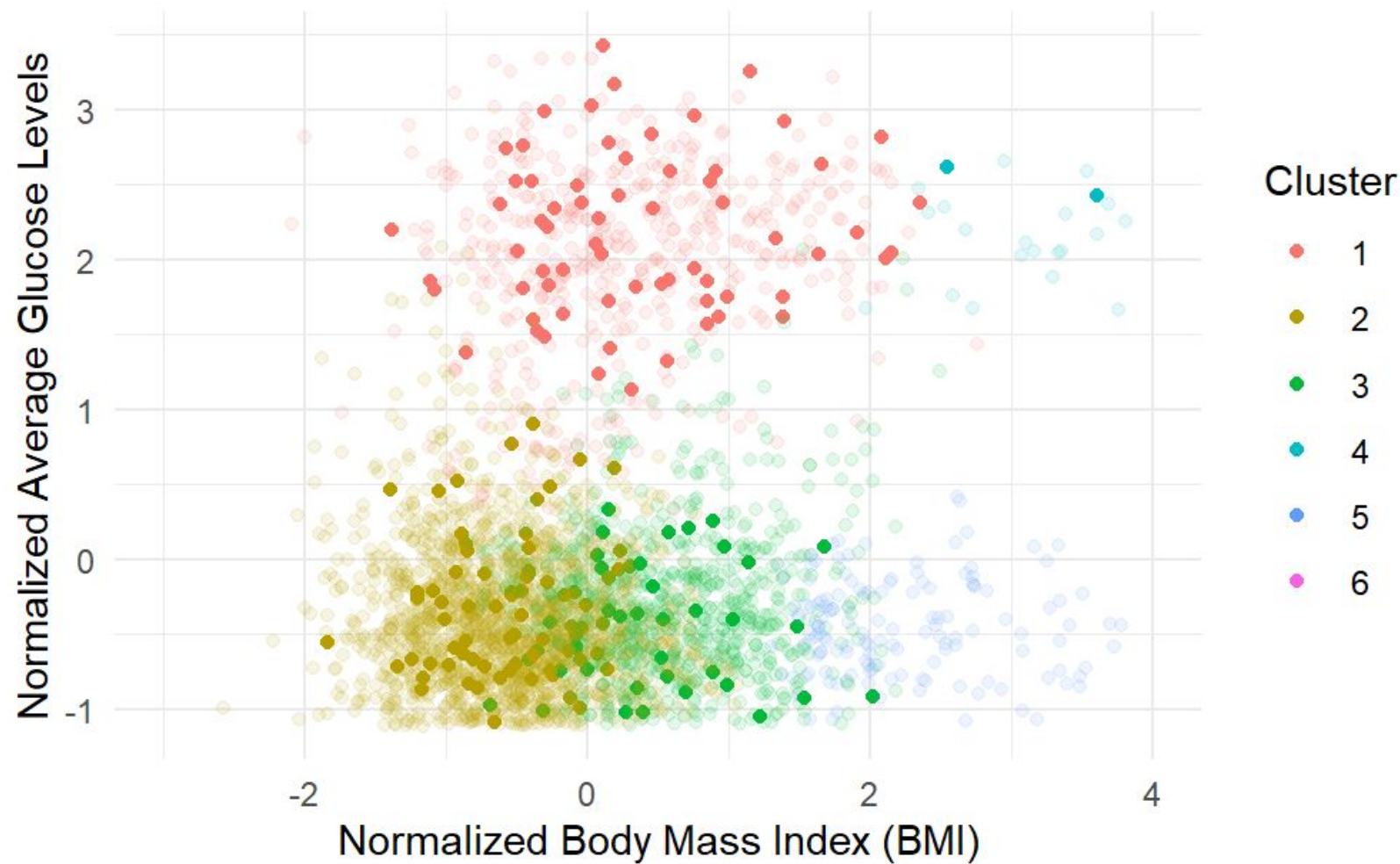




Cluster Dendrogram





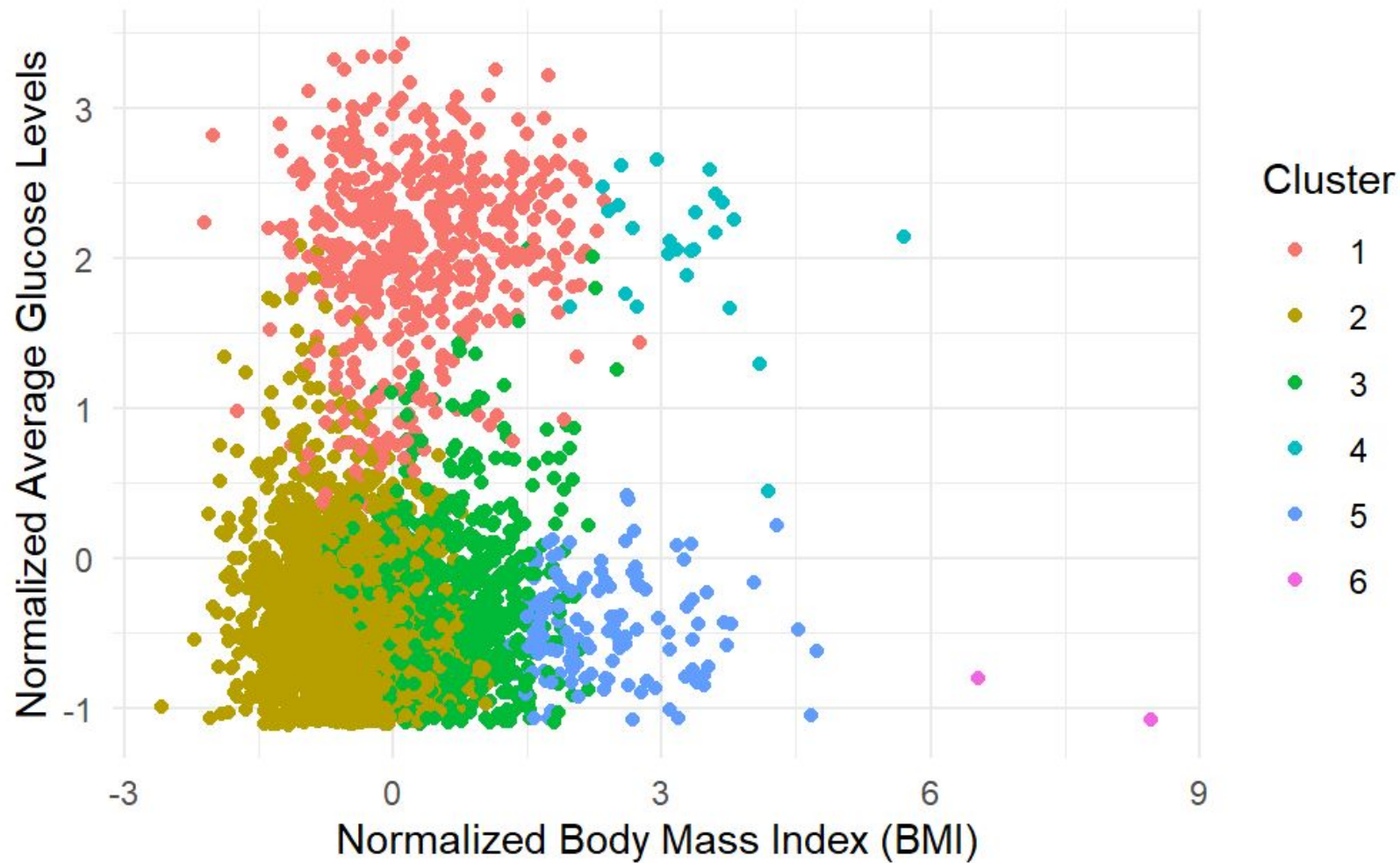


Relative Frequency of Stroke Event by Cluster



Cluster	Total Cases	Number of Strokes	Proportion of Strokes
1	508	66	0.13
4	25	2	0.08
3	816	39	0.05
2	1929	73	0.04
5	146	0	0
6	2	0	0





Clustering Takeaways

- Hierarchical clustering primarily groups cases by **average glucose level and BMI**
- Clusters with **higher average glucose levels tend to have higher instances of stroke**, whereas BMI does not tend to be as important in the dataset