



National Institutes of Health
Turning Discovery Into Health

Urinalysis Unlocking the Next Phase of Alzheimer's Diagnosis

• • •

Presentation by: Andranique Green



Table Of Content

• • •

- 01. Background
- 02. Problem Statement
- 03. R&D Process
- 04. Experimental Conditions

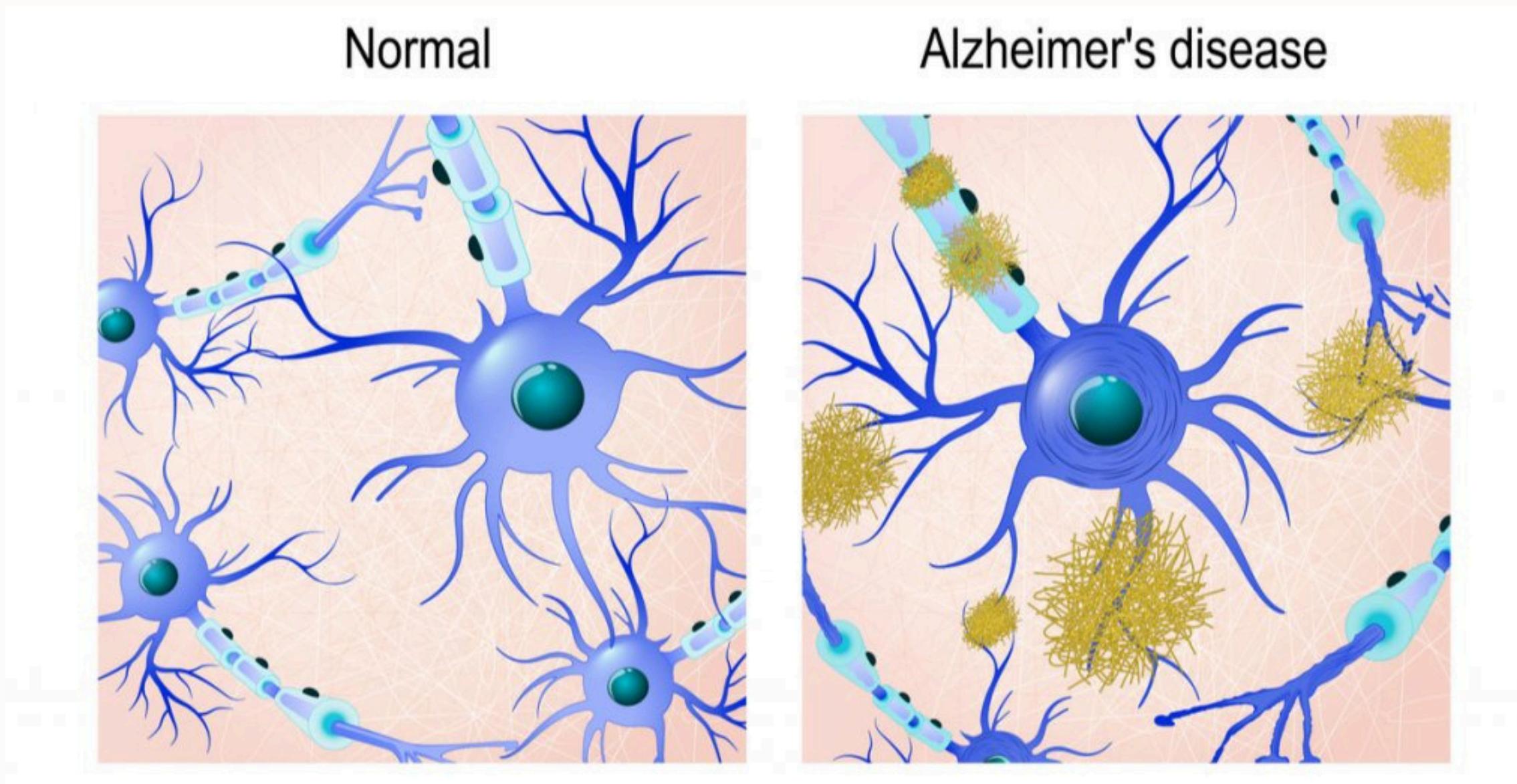
- 05. Addressing Limitations
- 06. Data Handling
- 07. Model Performance
- 08. Conclusion



Subject Matter Background



Alzheimer's Disease is caused by the build up of abnormal proteins



Problem Statement

• • •

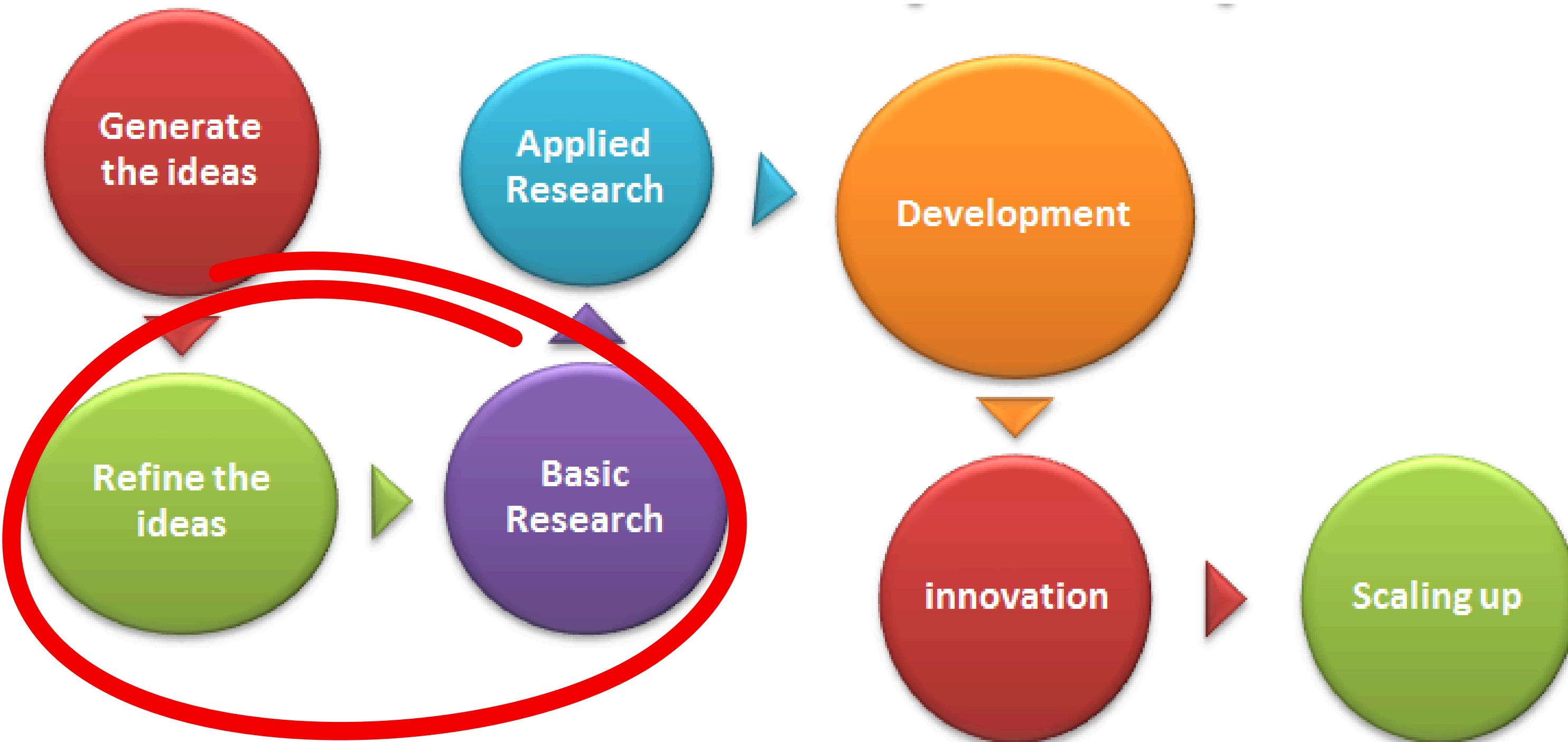
Alzheimer's Disease

It is estimated that Alzheimer's related dementia accounts for 70% of total dementia cases and number of people living with the disease is steadily rising. There are various tests that assist in the diagnosing process; unfortunately the majority of these tests tend to be inaccessible, invasive, and expensive. This project aims to change this fact.

How

The goal of this project is to identify biomarkers for detecting Alzheimer's dementia from urine and create a model that can correctly classify between individuals with Alzheimer's and those without. Which reduce pain and cost associated with current testing, and lower accessibility barriers

Research and Development Process





Experimental Data:

- 
- ✓ Niigata University Graduate School of Medical and Dental Sciences
 - ✓ Urine samples were collected from 36 individuals sex and aged matched
 - ✓ The samples were analysis by mass spectrometry
 - ✓ Results: There were changes in 109 proteins

Healthy Control

Noted as 0 in the cleaned dataset. Cognitively normal individuals



Alzheimer's Sample Group

Noted as 1 in the cleaned dataset. Individuals Experiencing Alzheimer's related demetia



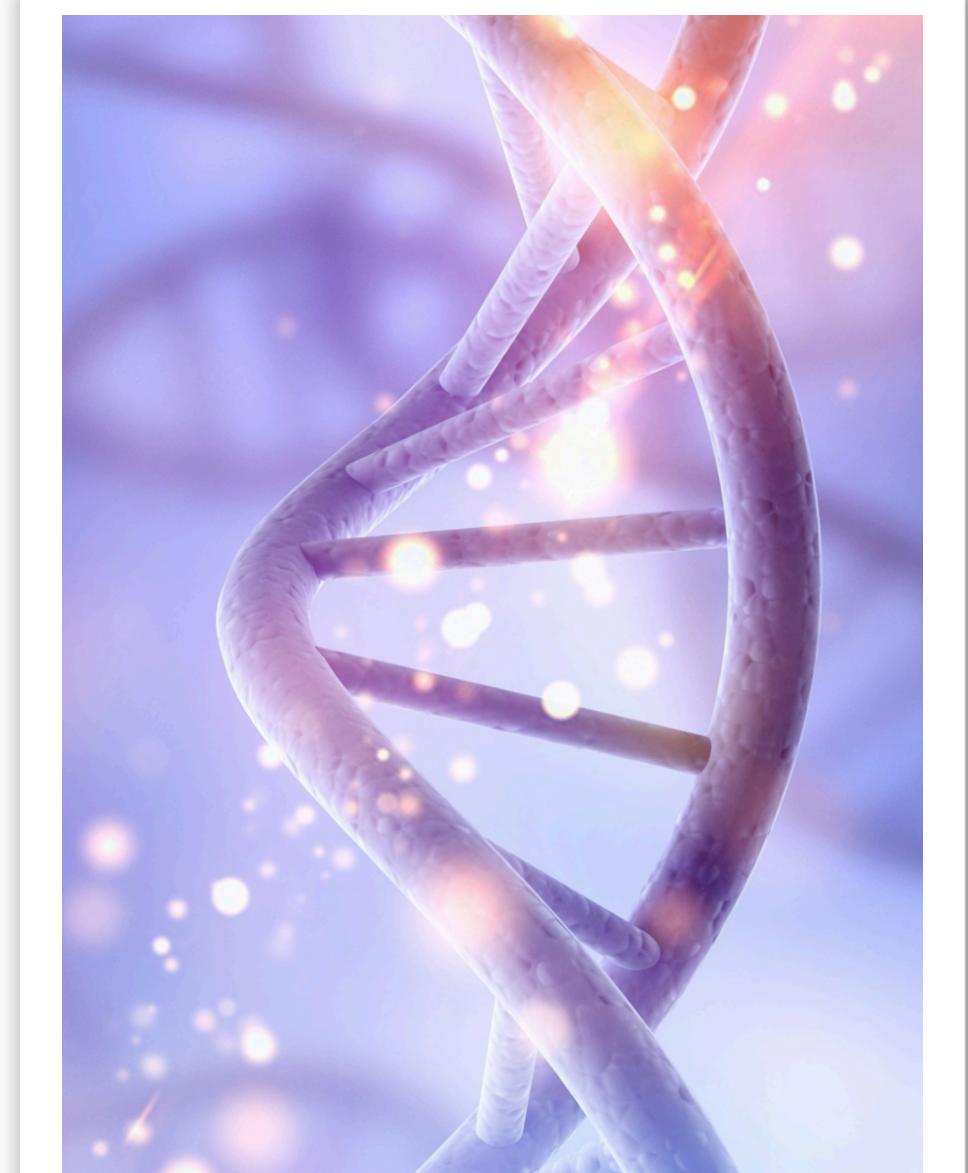
Data Limitations

• • •

The experimental data I have was only partially completed.

While I had very comprehensive data on a singular patient there are still 36 total samples.

There isn't much we can do about the dataset but being that it is early in test development we can use the results as a guiding path for research advancement



Data Handling



The individual sample data was very extensive and the reports formatting made it very hard to extract.

For this reason I relied specifically on emPAI values or estimated protein ratios based on peptide presence.

102	32979	103	12840	45660
2132400	32979	103	681	6
132313	1400	920	150054	481
15	13843	627	478	401358
135	355	51940	4703	146
13	462	4764981	4613433	1384
3	51781	4311894	1321	355
	4764816	1845	13678	462
	4311729	5468778	190	51781
	1680	655014	297	1781
	5468613	1682	51616	32192
	654849	156453	1616	1973
	170	1682	32027	32192
	56288	681	1808	1781
	1517	31862	1808	1176
	516	1643	1616	157149
	97	1643	1011	973
	8	1451	156984	5258
		846	808	806
		350	500	100

Data Included :

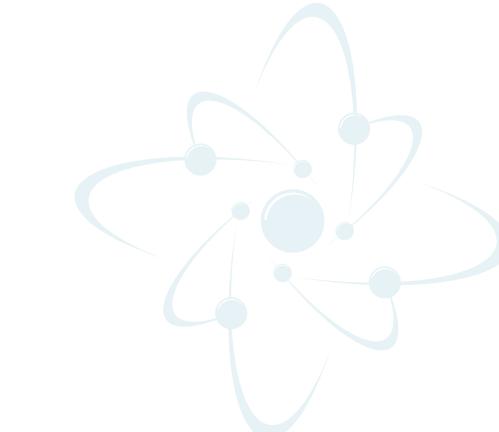
protein_desc- the names of proteins present

emPAI value

Alzheimer's- Target column



Binary Classification Models



• • •

Keeping the data limitation in mind I used a simple machine learning model to prevent overfitting. I ended up comparing models based on their accuracy score solely to identify how well they could differentiate based on protein values without preference for TP and TN. With later models metrics like sensitivity would be better suited to meet the problem's needs.

I chose logistic regression for model and I built two versions. One without added penalties and another using ElasticNet. The models performed the same regardless of the penalties



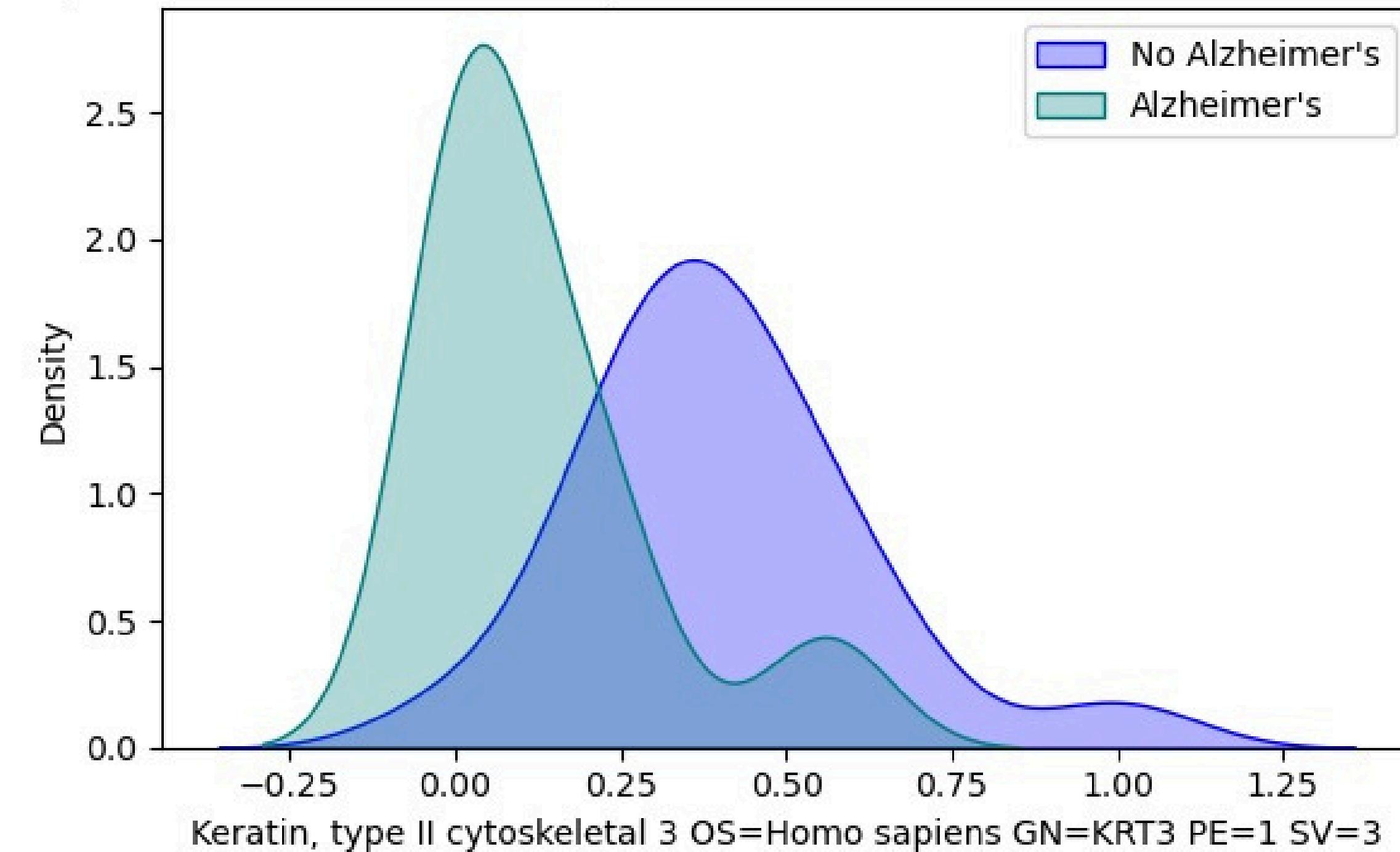
Target Protein Group

Extracted through calculating mean variance between groups

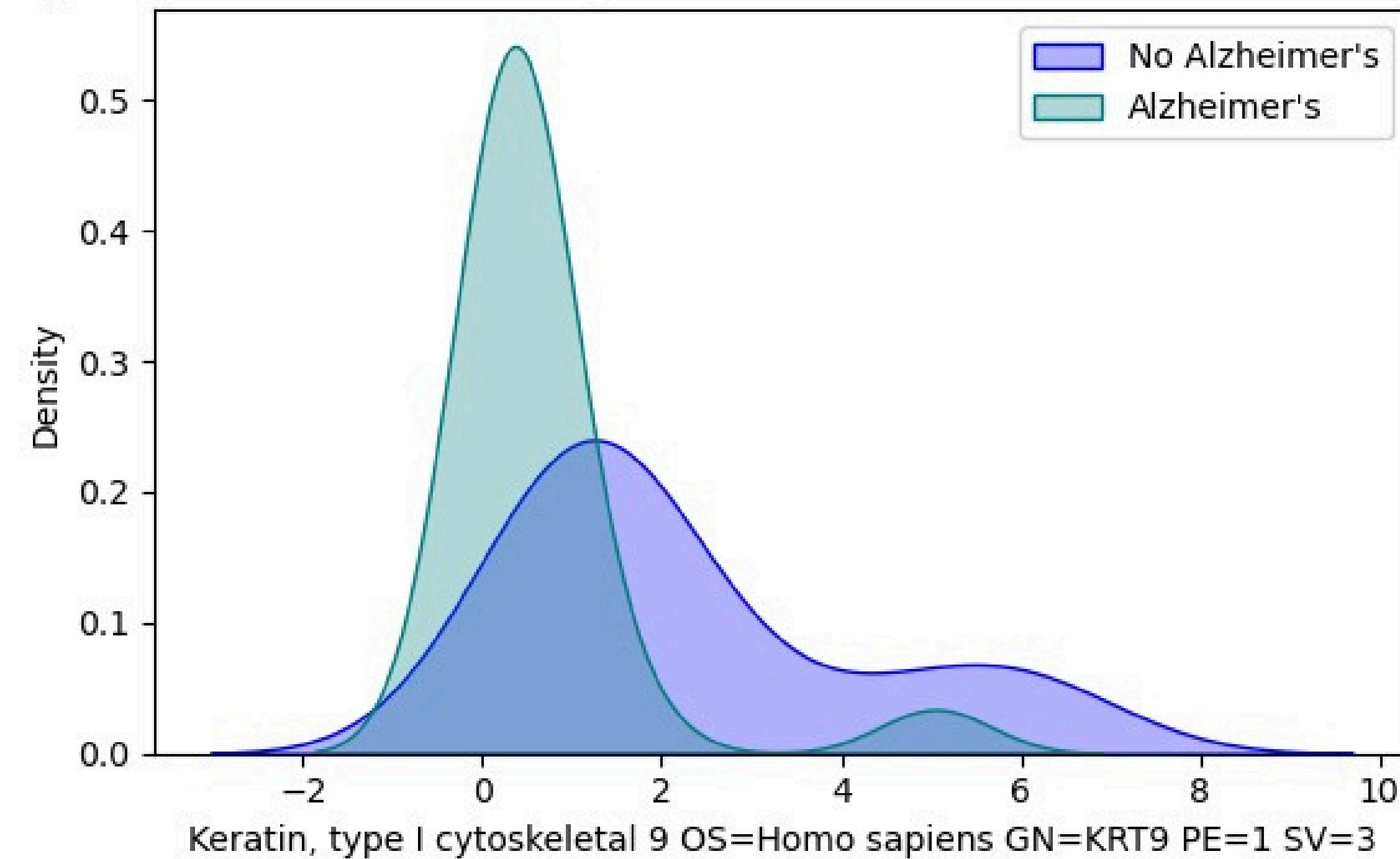
Important Features Group

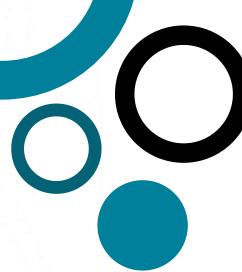
Decision Tree importance assigned proteins

Keratin, type II cytoskeletal 3 OS=Homo sapiens GN=KRT3 PE=1 SV=3 - Protein Abundance Distribution

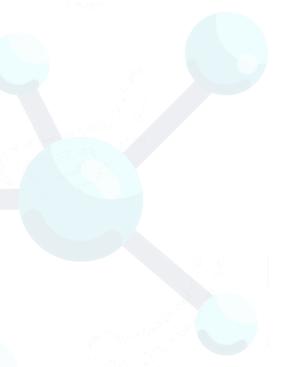
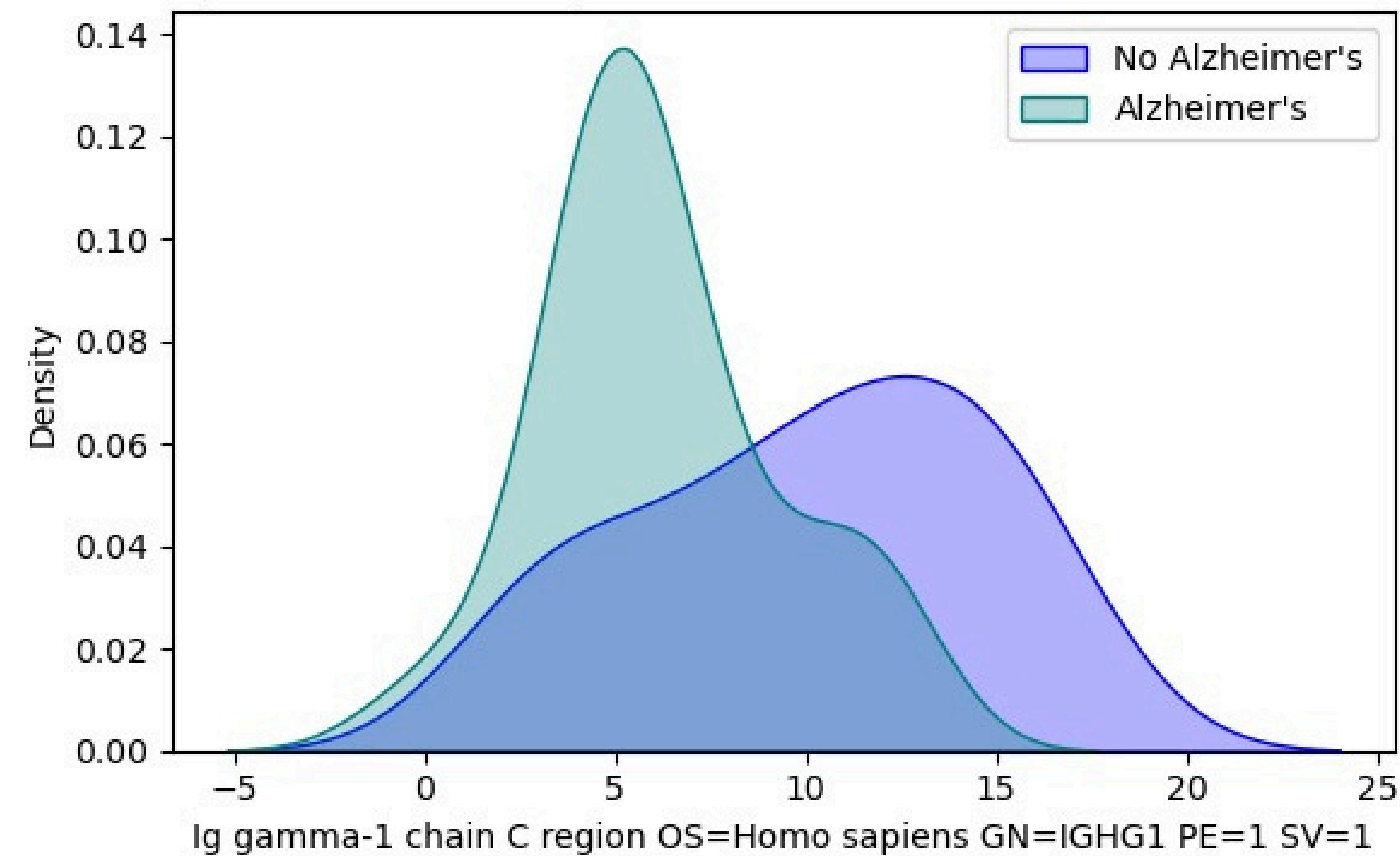


Keratin, type I cytoskeletal 9 OS=Homo sapiens GN=KRT9 PE=1 SV=3 - Protein Abundance Distribution

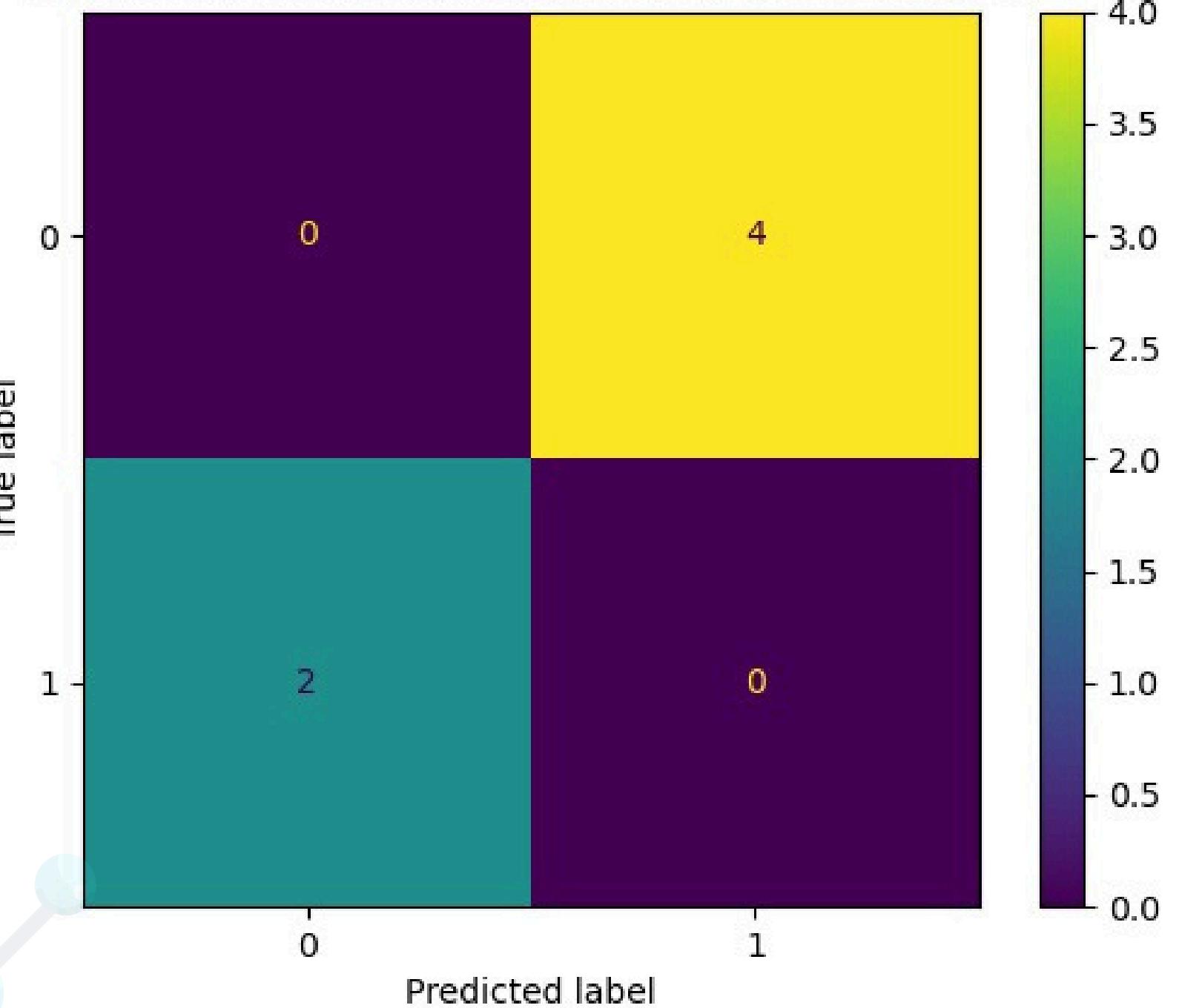




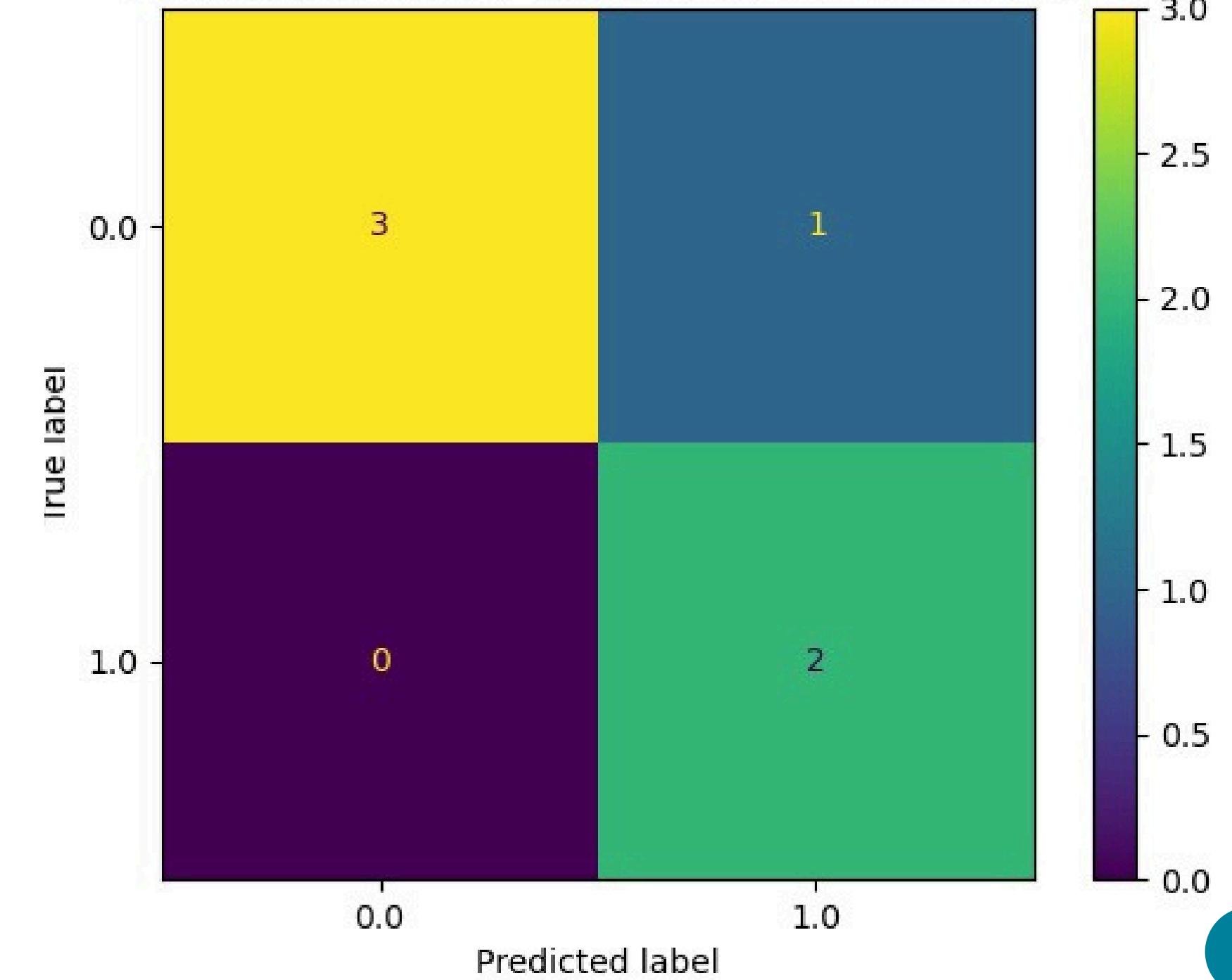
Ig gamma-1 chain C region OS=Homo sapiens GN=IGHG1 PE=1 SV=1 - Protein Abundance Distribution



Confusion Matrix for Important Features (Complex)



Confusion Matrix for Target Predictors (Complex)



Conclusion

• • •

Move forward with additional research on the 3 proteins identified by decision trees

Metrics:

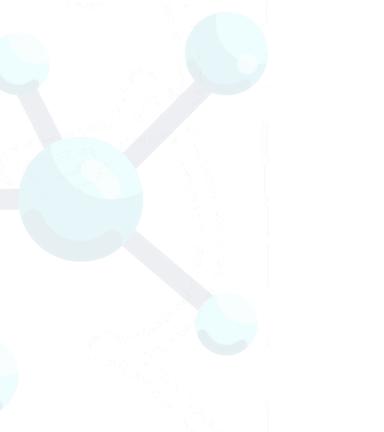
The Logarithmic model trained on the Decision Tree features had a 100% accuracy rate with testing and 87.5% with the training data



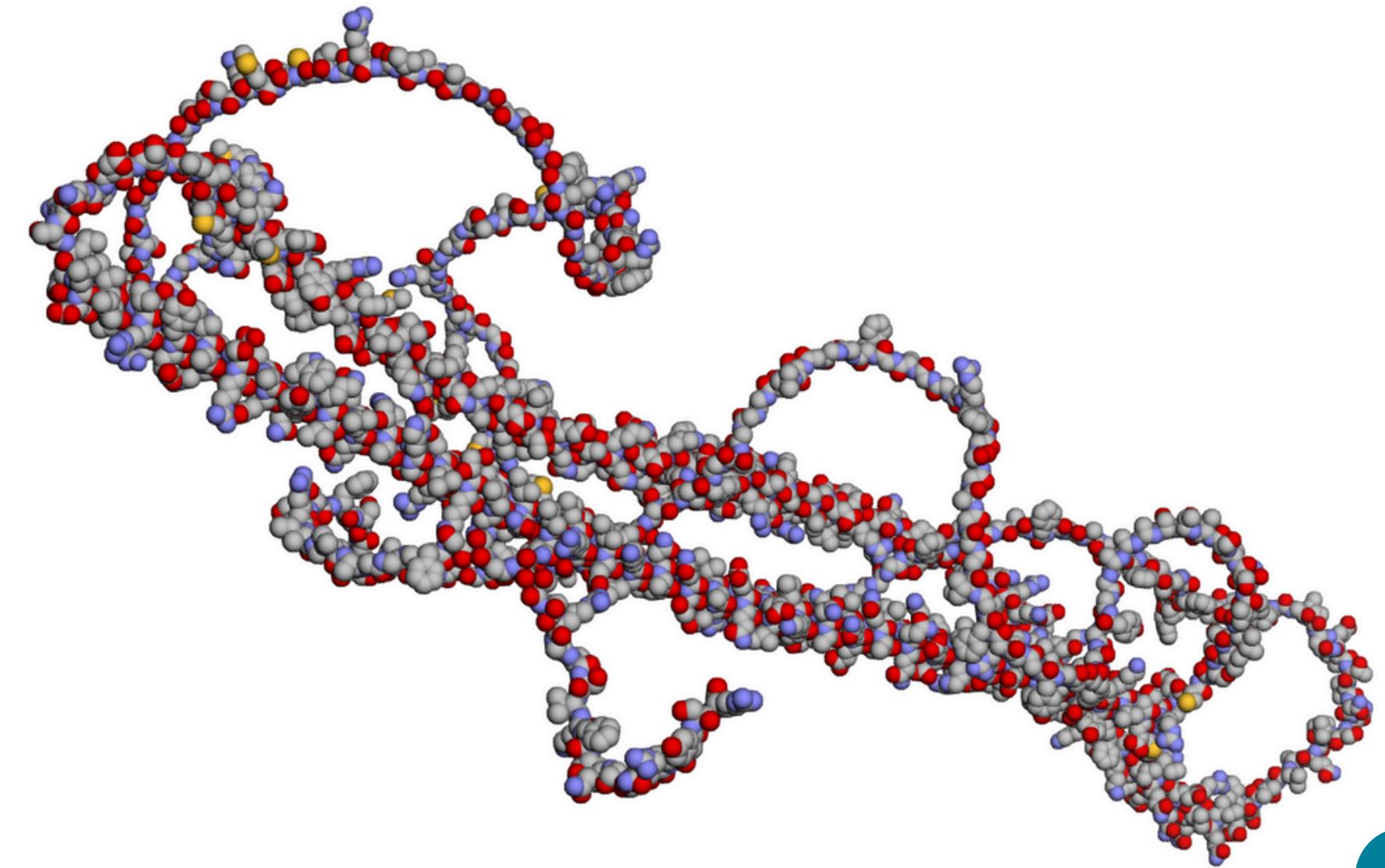
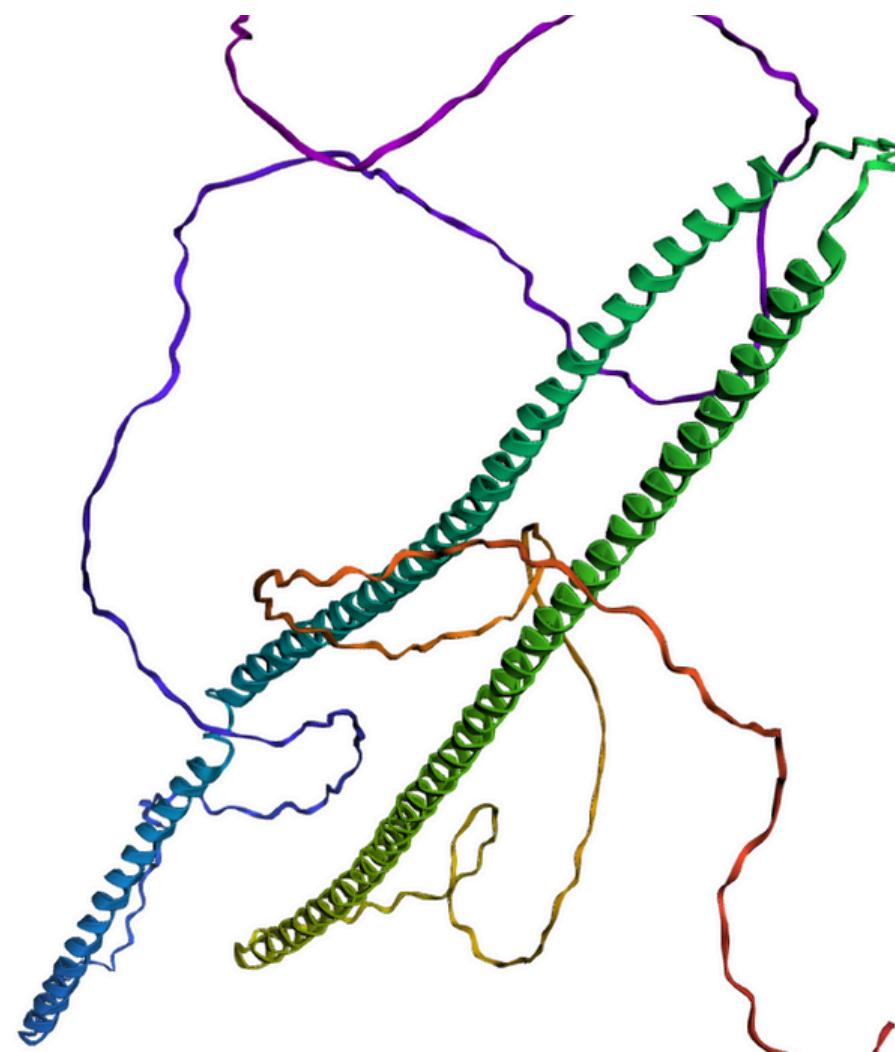
R&D

Next Steps

• • •

- Expand Dataset:
 - Collect Repeated Samples over the Disease Progression:
 - Incorporate More Information Available from CSV Report:
 - Employ Decision Curve Analysis and Patient Stratification:
- 
- 

Protein Visualizations: Keratin, type II



Thank You for Listening

