

# Predicting Breast Cancer

A report by Jessica Montealvo





# Introduction

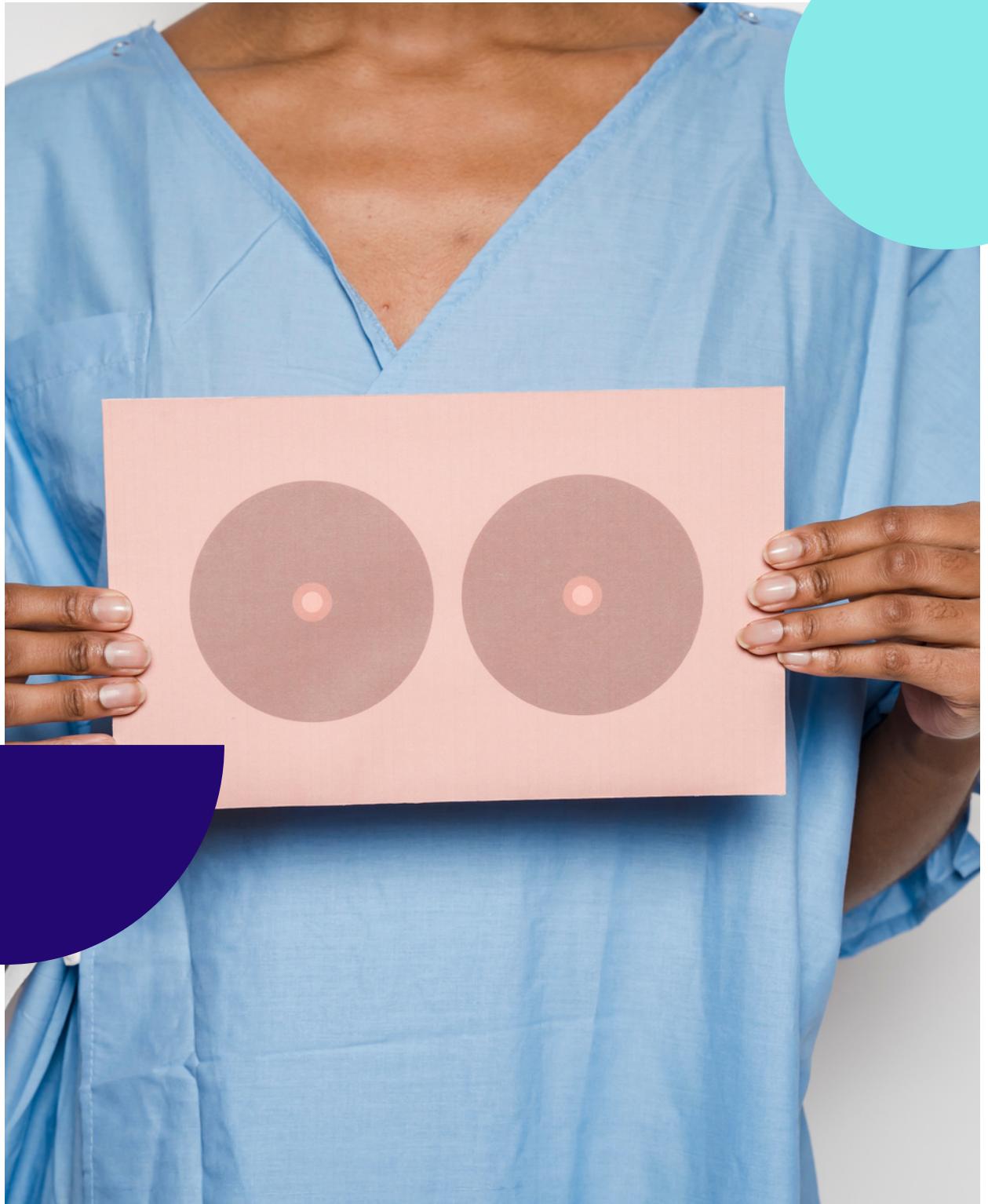
**Last year, breast cancer took the lead in the number of new cases reported, with 2.26 million diagnoses**

Although prominent, breast cancer is not as deadly as other cancers, mainly because it benefits from early detection. Early detection strategies have been implemented throughout the years and rely on both clinical and self-breast exams as well as mammogram screening. Earlier-stage cancers can be treated more effectively, which readily contributes to survival. So the goal is always to identify cancer early.

# Probelm Statement

In this project, I will use the Breast Cancer Wisconsin (Diagnostic) Data Set to train a model to predict whether a tumor is benign or malignant.





# The Data

## UCI MACHINE LEARNING REPOSITORY

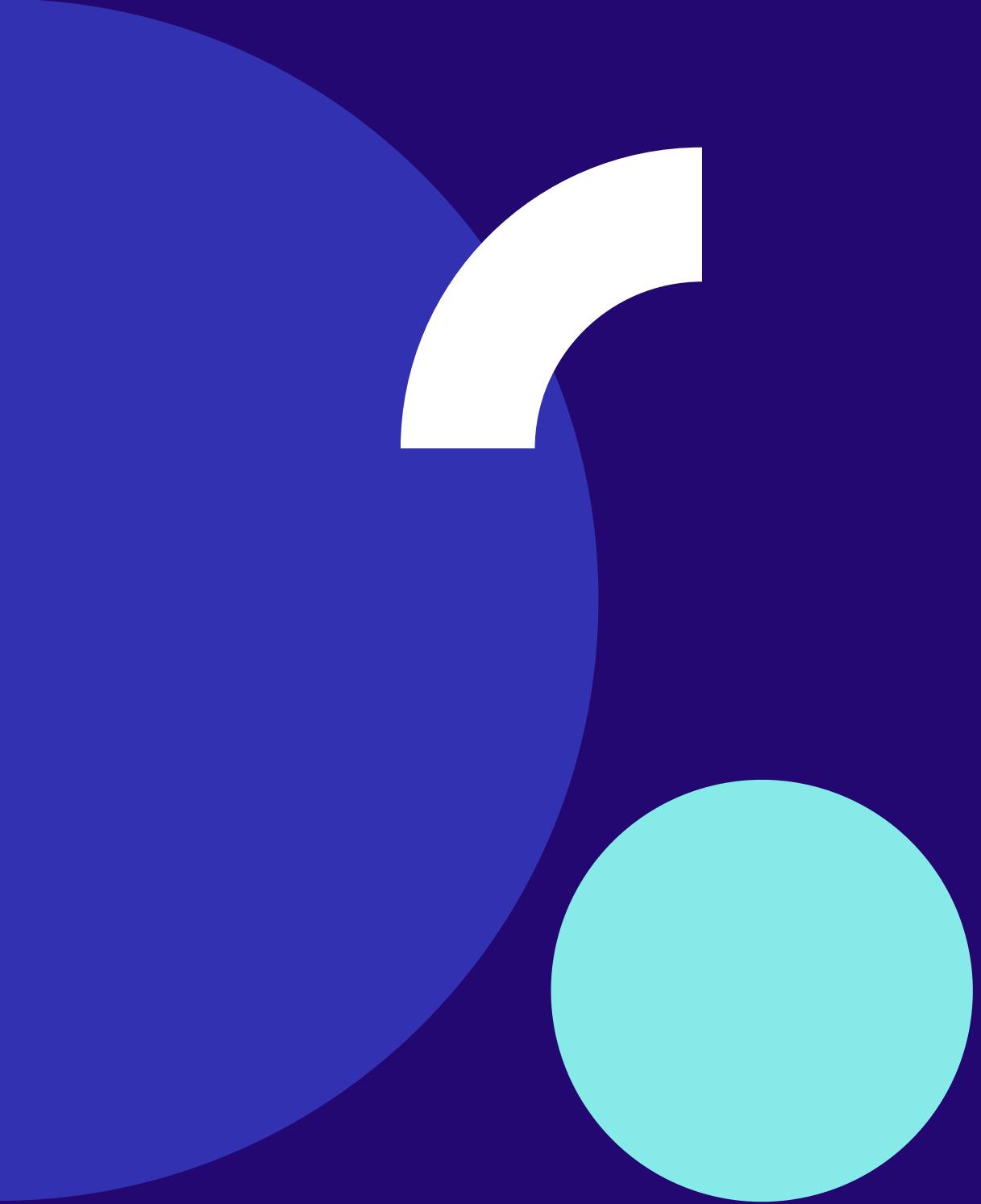
The data set used for this project can be found on the UCI Machine Learning Repository. It was gathered by Dr. William H. Wolberg and his team at the University of Wisconsin in 1995

## 569 CASES

The data describes 569 cases of digitized images of biopsies, in particular a fine needle aspirate (FNA) of breast tissue, when breast cancer was suspected.

## 30 FEATURES

The features in this data set are measurements derived from the images and describe the cell nuclei in the tissue.



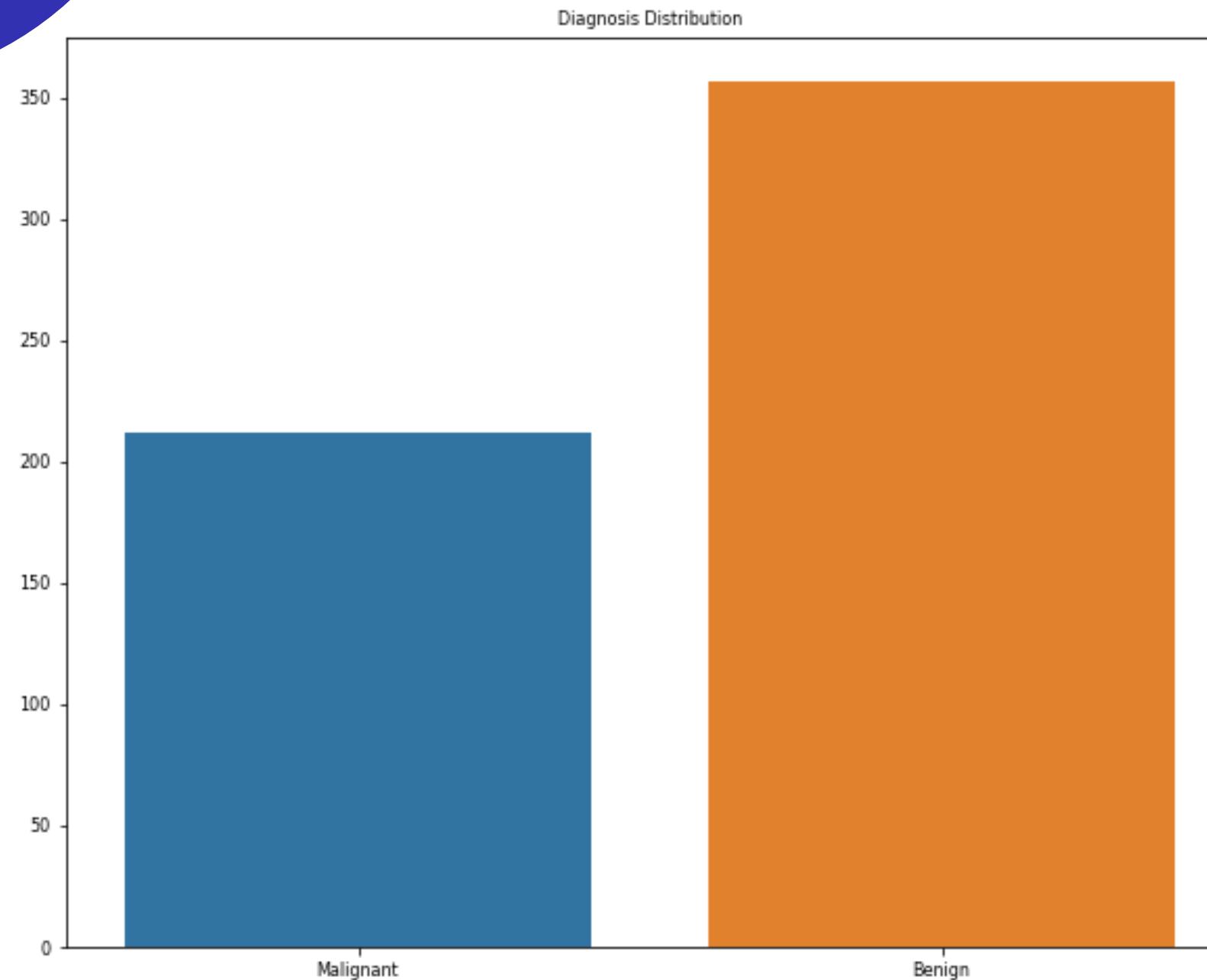
# Data Wrangling

There were no missing values, but depending on the source, the ID column can be dropped, and the diagnosis column can be translated to (malignant, benign) in preparation for the EDA.

# Exploration Data Analysis

---

# Diagnosis Distribution



**Not Cancer (benign): 62.74%**  
**Cancer (malignant): 37.26%**

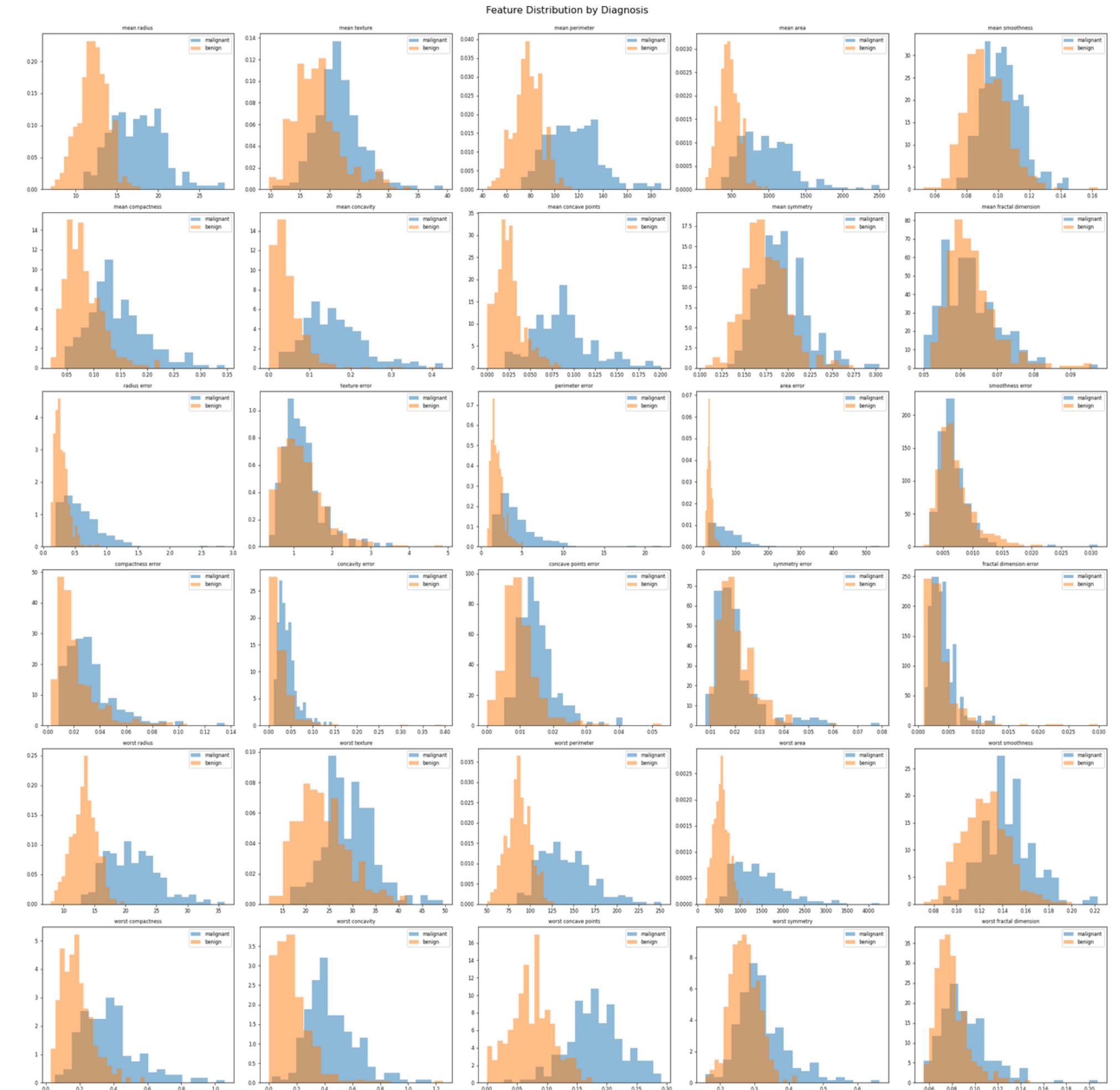
The following graph shows that our data set has two classes (cancer/ not cancer). The classes are not balanced. For this reason, I will upsample in my training data set to avoid my models predicting correctly solely based on the fact that there are more Not Cancer instances.

# Feature Distribution

Malignant in blue

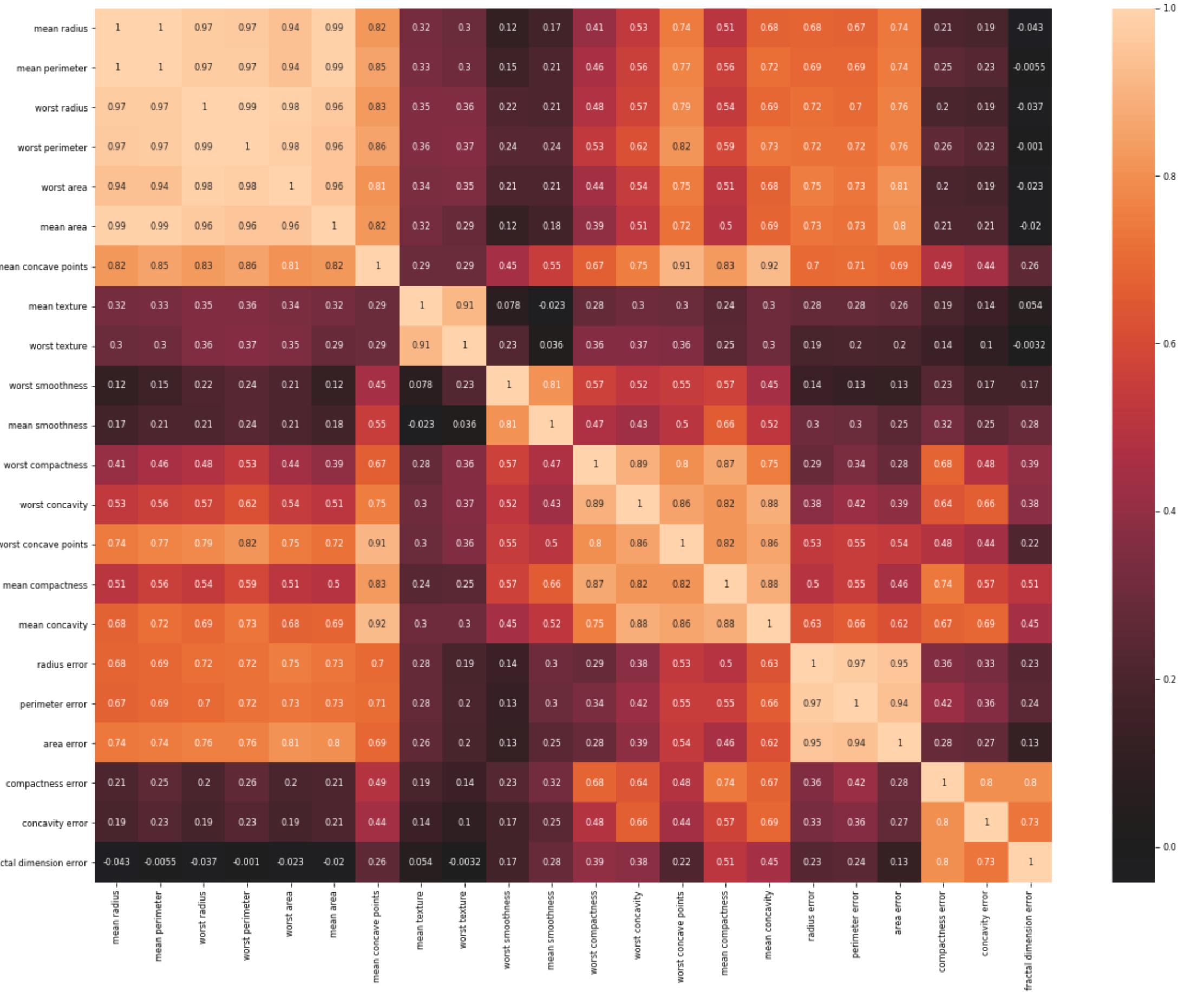
Cancer cells have less compact nuclei; accordingly, they are distributed on the larger portion of the spectrum for radius, mean texture, mean perimeter, area, concavity, and concave points.

Correlations are evident.



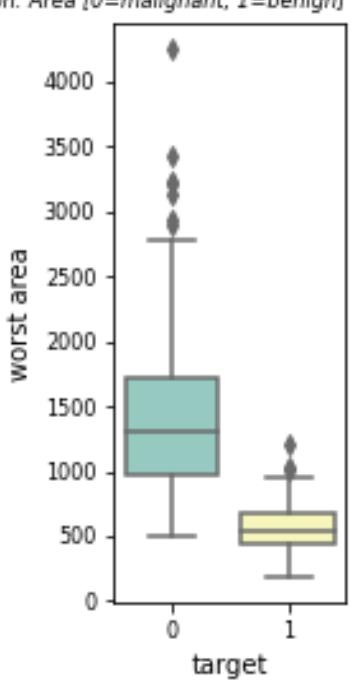
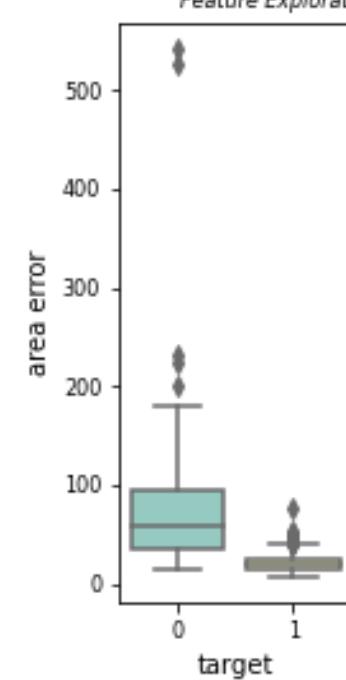
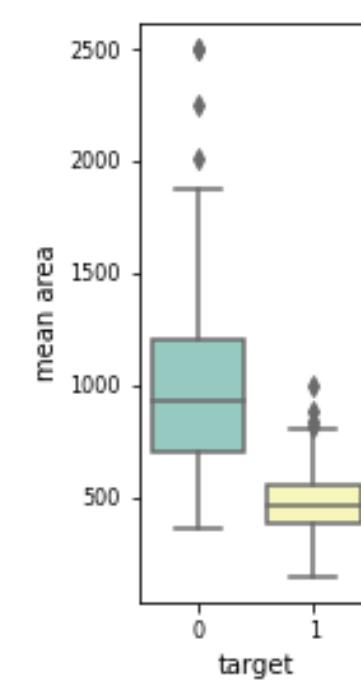
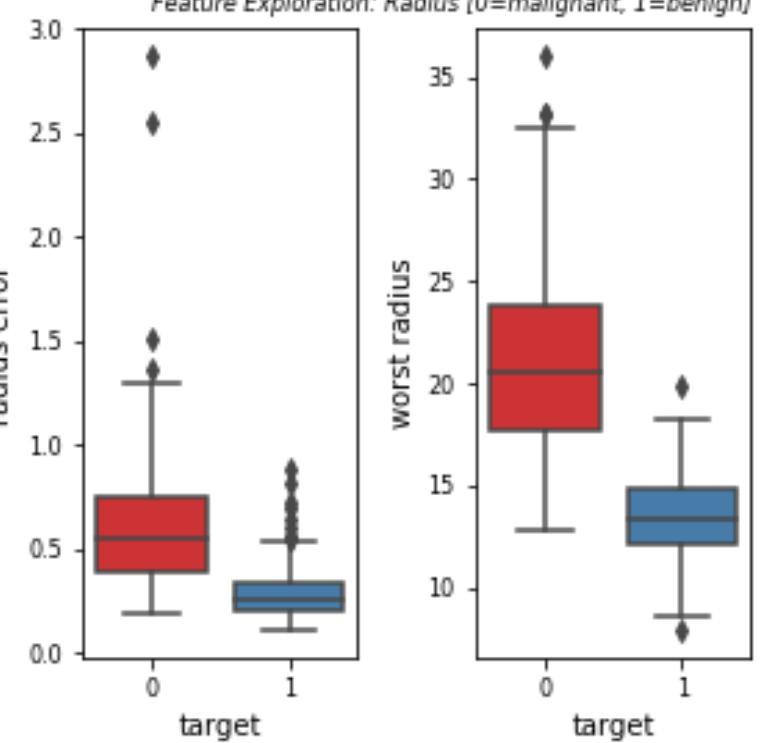
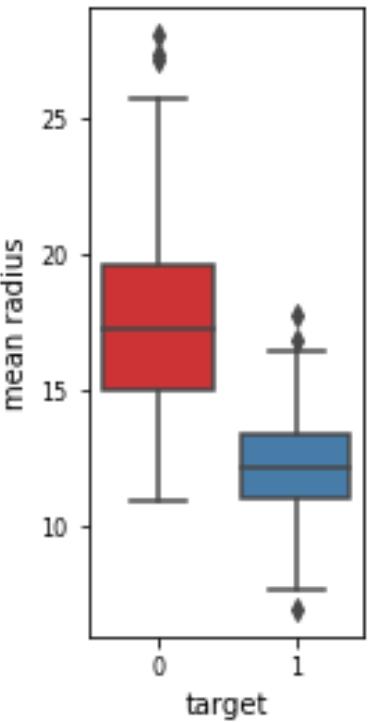
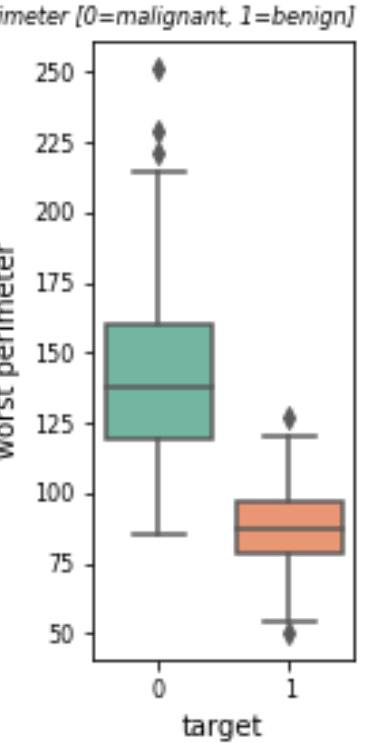
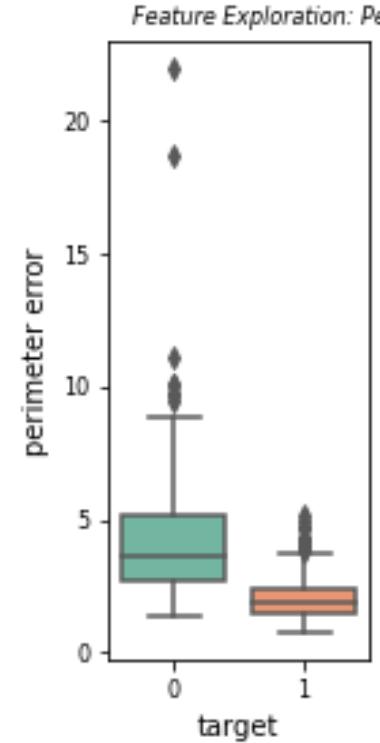
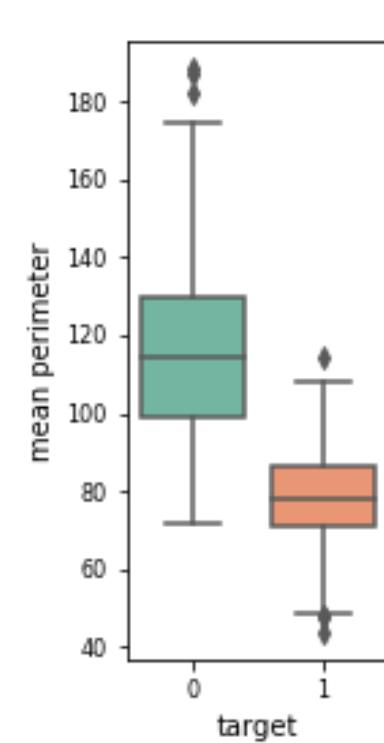
# Correlation Matrix Snapshot

As evidenced by the correlation matrix, namely by the peach-colored squares in the matrix, in many instances, the mean, worst, and error measurements for each category correlate to each other.



# Presence of Collinear Features

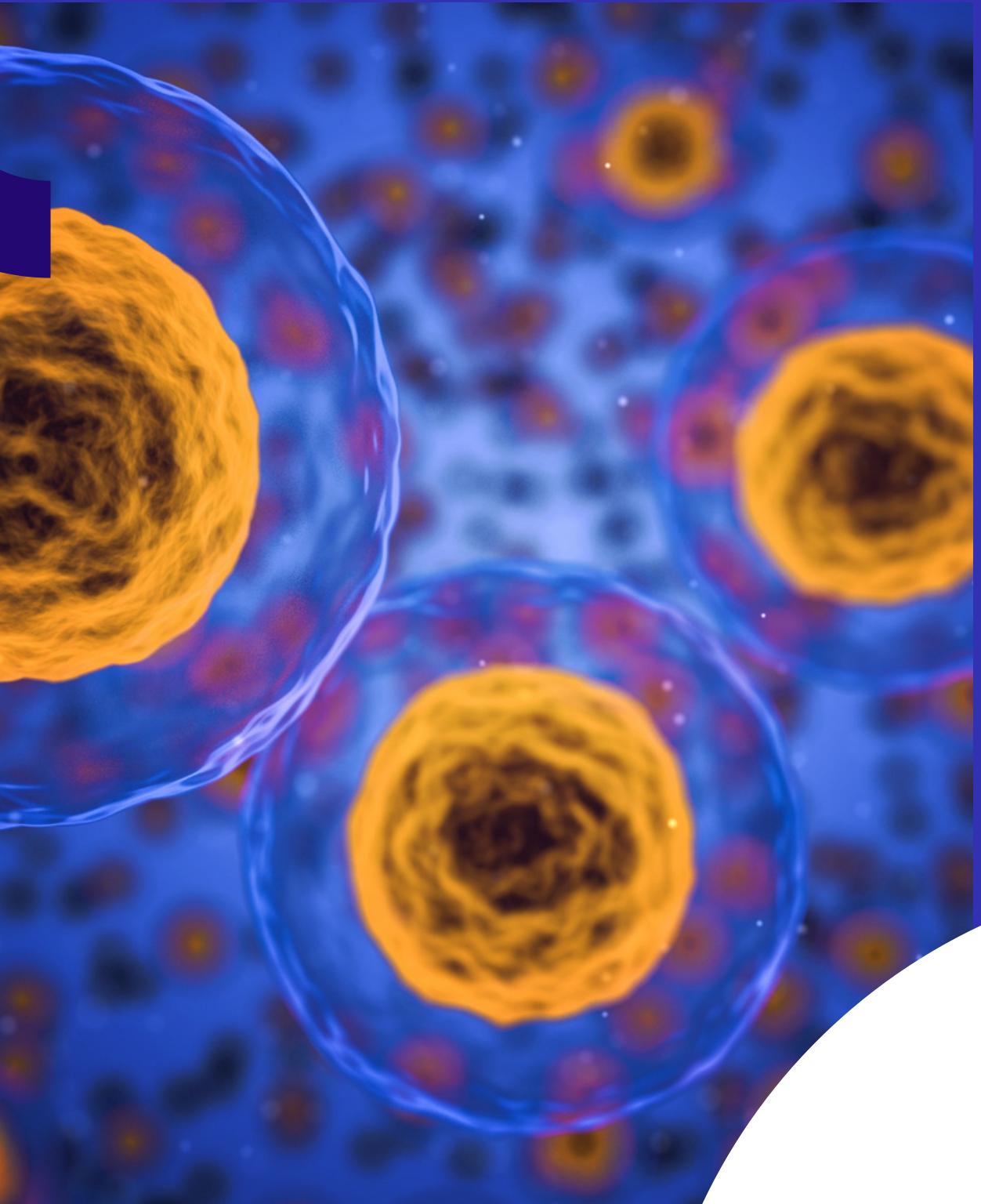
Features that describe the same phenomenon, which may confuse or mislead a model. A robust feature selection will need to tease these apart.



# Feature Selection

**Six feature selection methods used**

- Select K Best
- Recursive Feature Elimination
- LassoCV
- Random Forest Feature Importance + Hierarchical clustering
- Only Means
- <80% correlated features



# Feature Selection Results

Feature Set	Selected based on	Number of Features	Baseline (LogisticRegression) ROC AUC Score
Random Forest Feature Importance	Hierarchy based on Correlation and node distance	14	0.9975
Recursive Feature Elimination with Cross-Validation	RFECV	18	0.9974
Original	Kept all features	30	0.9972
Lasso CV	LassoCV/LogReg	22	0.9961
SelectKBest	SelectKBest	15	0.9906
Less Correlated Features	Removed correlated features (threshold=0.8)	16	0.9896
Mean Features	Removed error and worst feature sets	10	0.9894

# Machine Learning

---

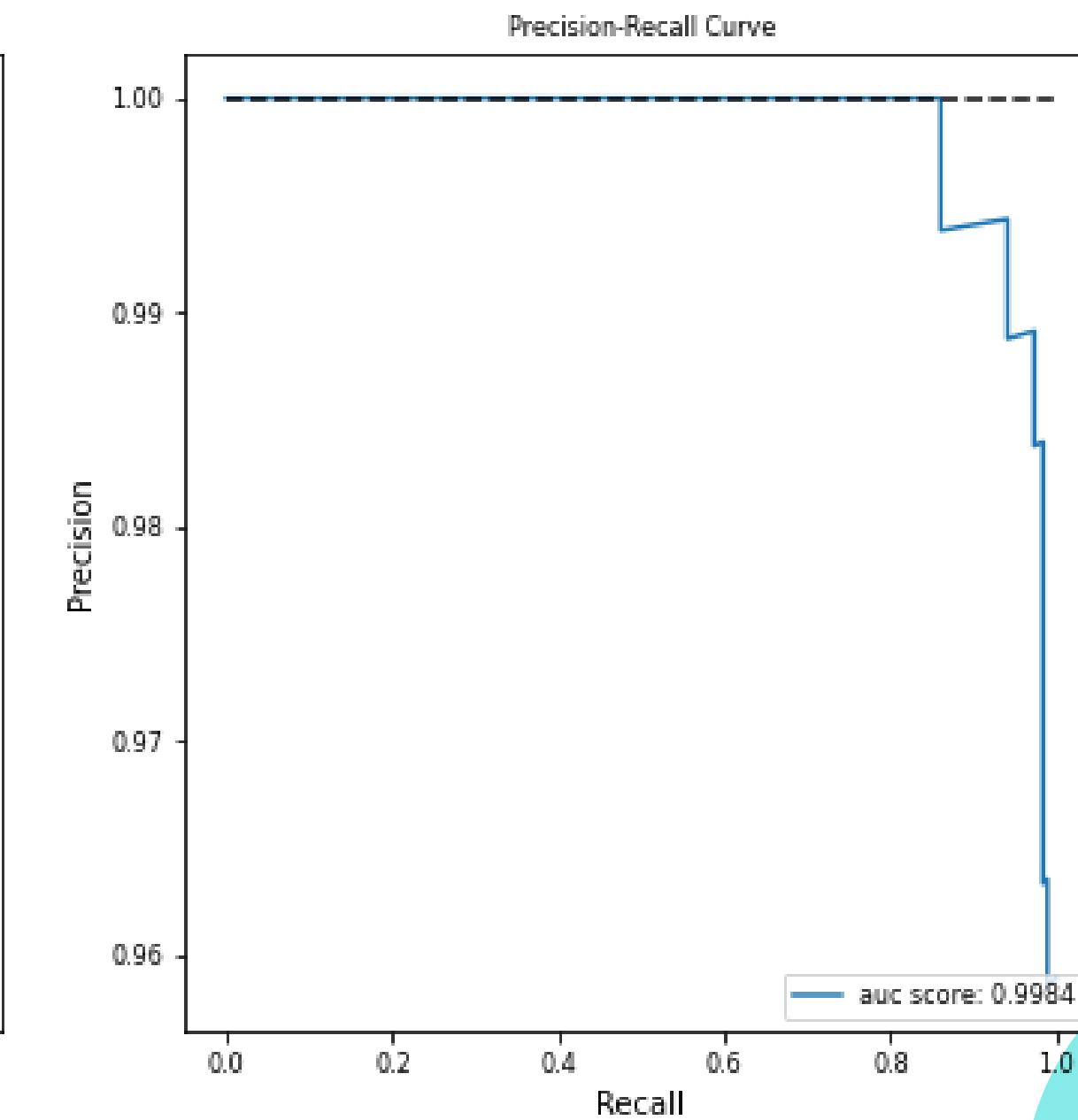
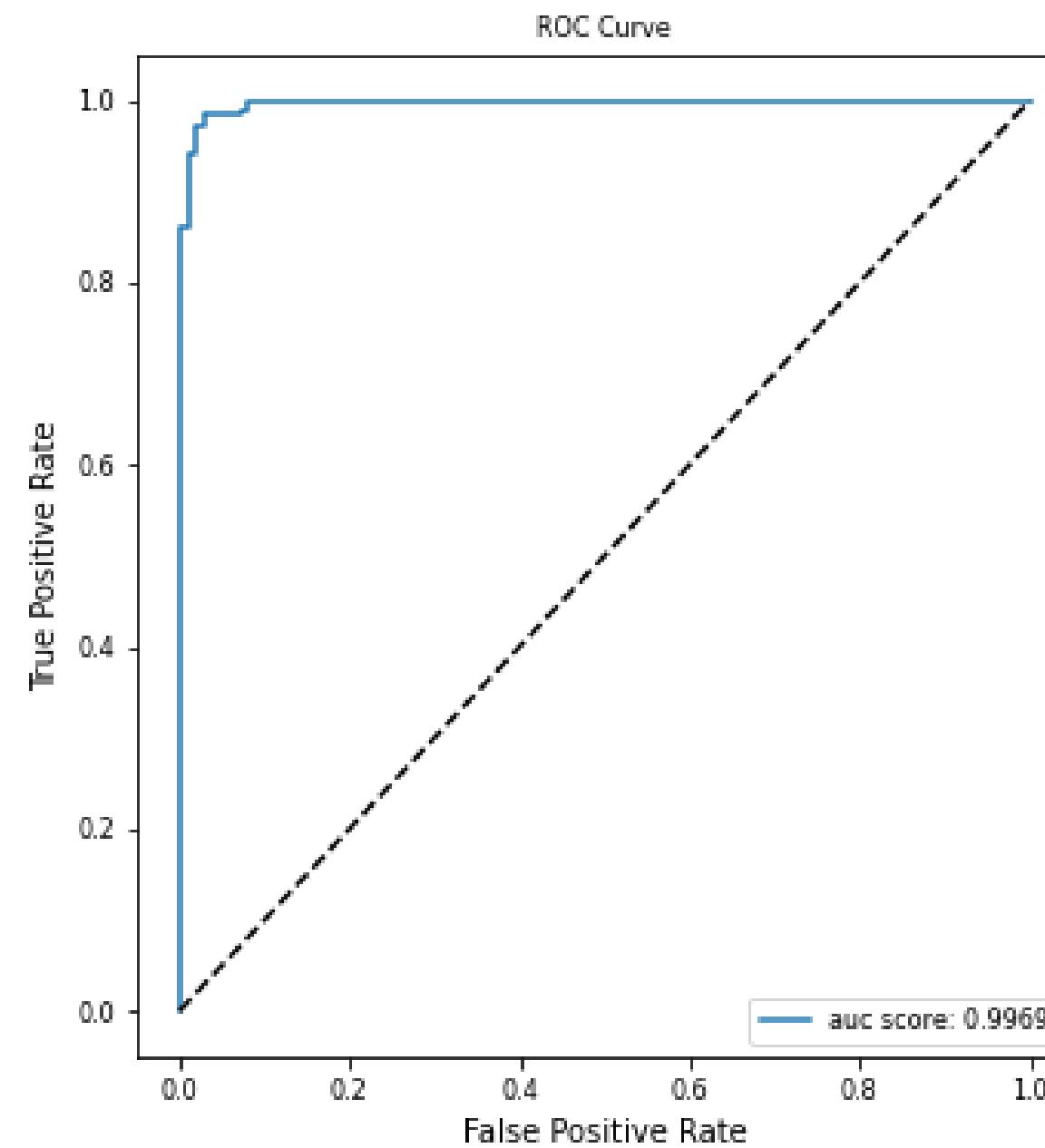
# Modeling

Data set selected from the Random Forest Feature Importance ranking and selection. Three different models assessed via Grid Search Cross-Validation.

Model	Best Score	Best Parameters	ROC AUC Score
LogisticRegression	0.9676	{'C': 1.0, 'max_iter': 100, 'penalty': 'l1', 'solver': 'liblinear'}	0.9969
Gradient Boosting	0.9529	{'learning_rate': 0.5, 'max_depth': 2, 'n_estimators': 500}	0.9915
Random Forest	0.9471	{'criterion': 'entropy', 'max_depth': 5, 'max_features': 'log2', 'min_samples_leaf': 1, 'min_samples_split': 2, 'n_estimators': 15}	0.9906

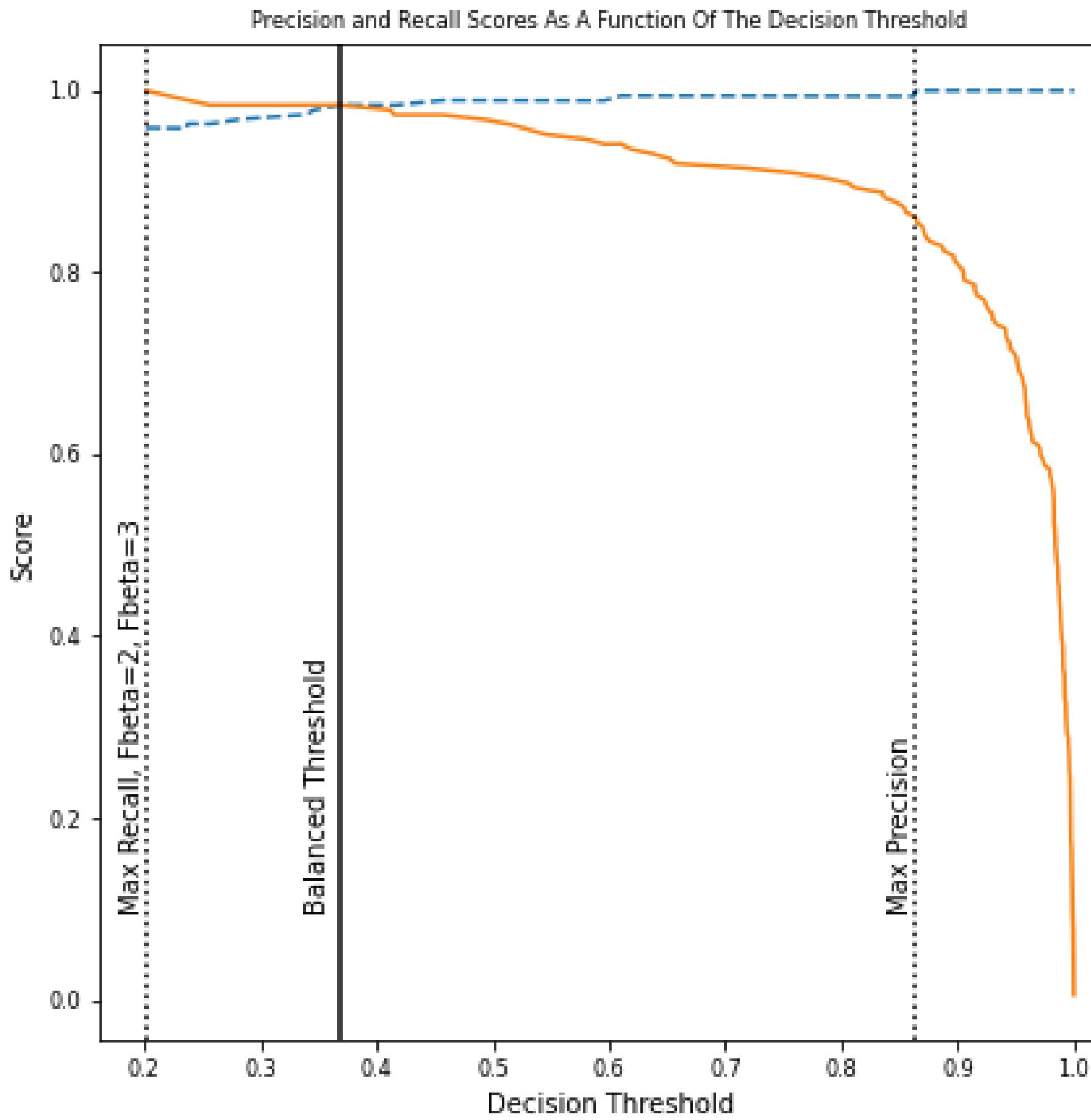
# Best: Logistic Regression

ROC AUC = 0.9969 Precision/Recall = 0.9984



Near perfect model,  
but we need a useful  
model that we can  
apply to the field.

# Precision and Recall Decision Threshold



## FBETA CALCULATED

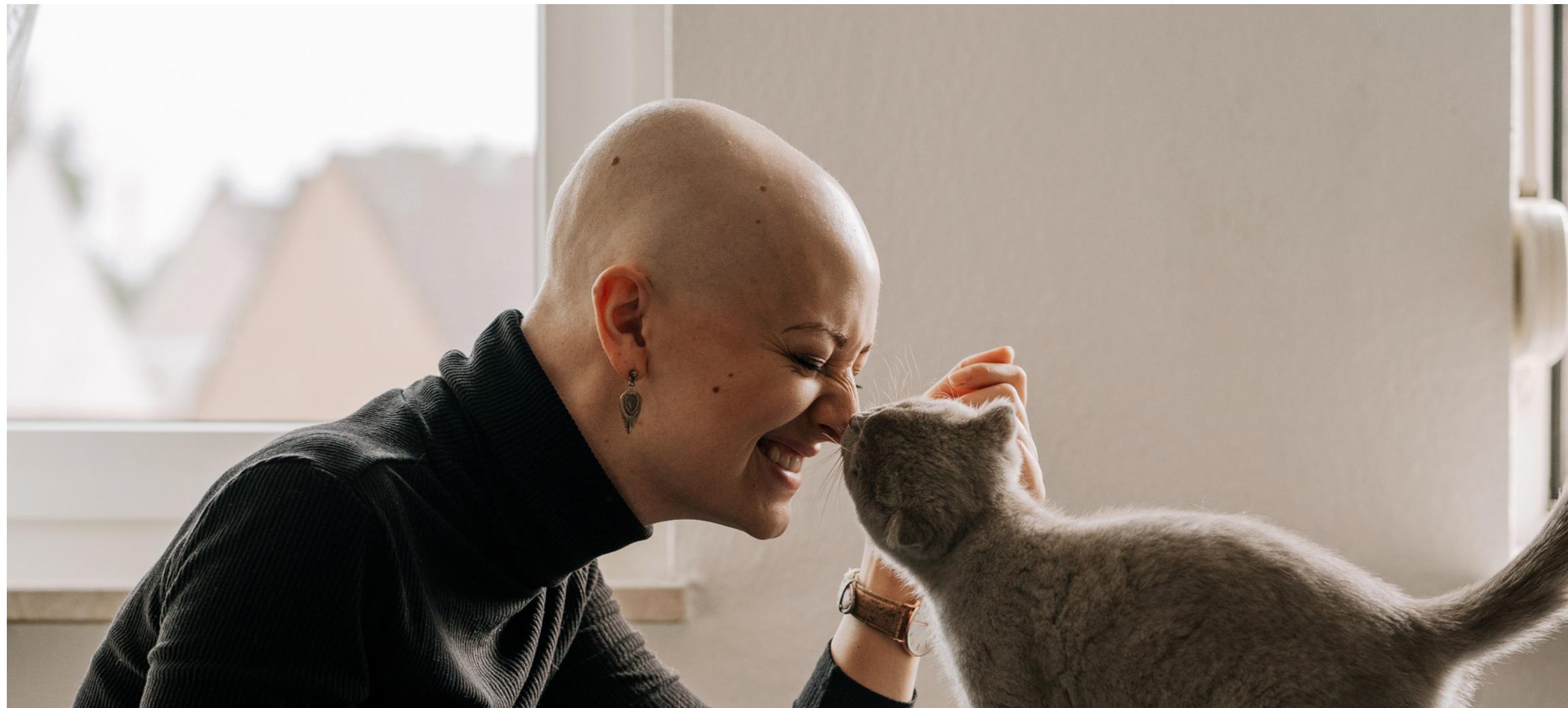
The Fbeta measure is calculated using the Precision and the Recall of a model. However, it uses a coefficient, beta, to add weight to either parameter. used  $f\beta = 2$  and  $f\beta = 3$ .

## PRIORITIZING RECALL

These would be the instances where a person walks away with a “not cancer” diagnosis, all the while harvesting a malignant tumor.

## PRIORITIZING PRECISION

Minimize false-positive errors, which in this project are less harmful since diagnosing a human with cancer while not having cancer will not be detrimental.



# Max Recall

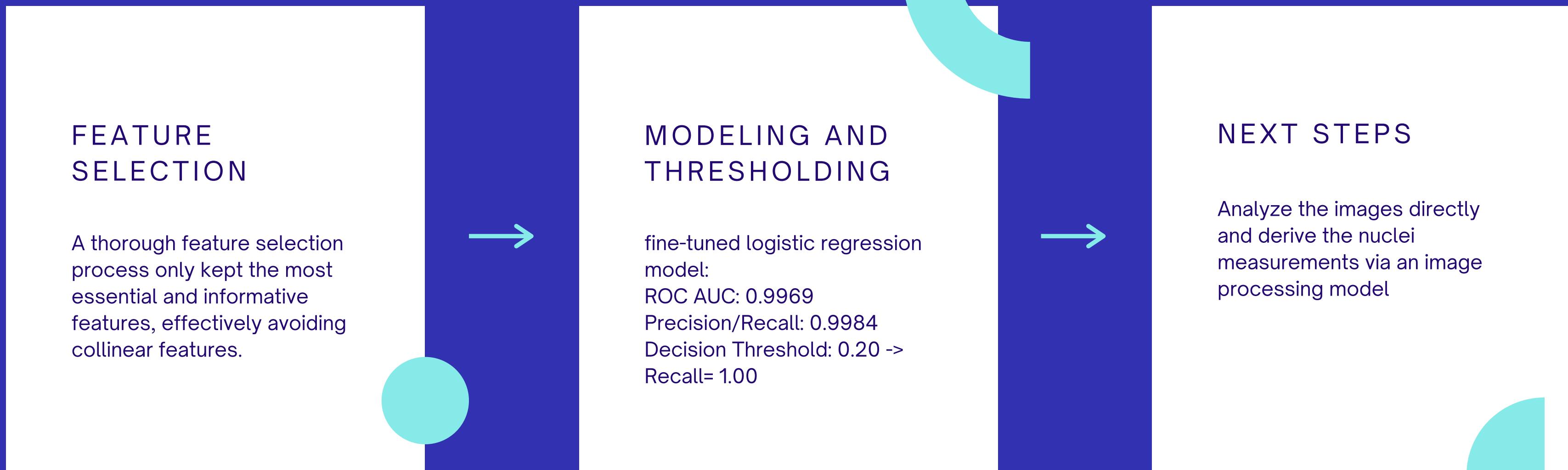
Threshold: 0.201174

Classification Report - Max Recall

	precision	recall	f1-score	support
malignant	1.00	0.92	0.96	98
benign	0.96	1.00	0.98	187
accuracy			0.97	285
macro avg	0.98	0.96	0.97	285
weighted avg	0.97	0.97	0.97	285

Setting a threshold of 0.20 will prevent any false negatives predictions, thereby giving the patient a chance at early detection and treatment.

# Streamlining Platforms



# References

WORLD HEALTH ORGANIZATION

<https://www.who.int/news-room/fact-sheets/detail/cancer>

SCIKIT-LEARN DOCUMENTATION

<https://scikit-learn.org/stable/index.html>

