

Predicting Malignant Breast Tumors – Supervised Learning Model

AI Academy Apprenticeship Capstone Project – Rachel Day

Agenda

1 | Project Goals

2 | Data Understanding

3 | Methods

4 | Results



Background and Purpose

Background

Breast imaging techniques can detect suspicious areas, but can't tell a patient whether they have cancer or not*.

An estimated 297,790 women are expected to be diagnosed with breast cancer in 2023**.



Goal

What if there was a way, through existing imaging techniques, to get tumor measurements and be able to predict if it was malignant or benign?

Goal: Create a model that accurately predicts malignant tumors given existing measurements.

*[Can a Radiologist Diagnose Breast Cancer from Imaging Tests Alone? \(healthline.com\)](https://www.healthline.com/health/cancer/can-a-radiologist-diagnose-breast-cancer-from-imaging-tests-alone)

**[Breast Cancer Statistics | How Common Is Breast Cancer?](https://www.healthline.com/health/cancer/breast-cancer-statistics)



The Data

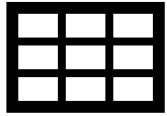
Pulled from Kaggle.com, an AirBnb for data scientists.

Kaggle is a crowd-sourced platform to attract, nurture, train and challenge data scientists.

The Dataset: Contains measurements and characteristics from breast tumors coming out of the State of Wisconsin.



The Data



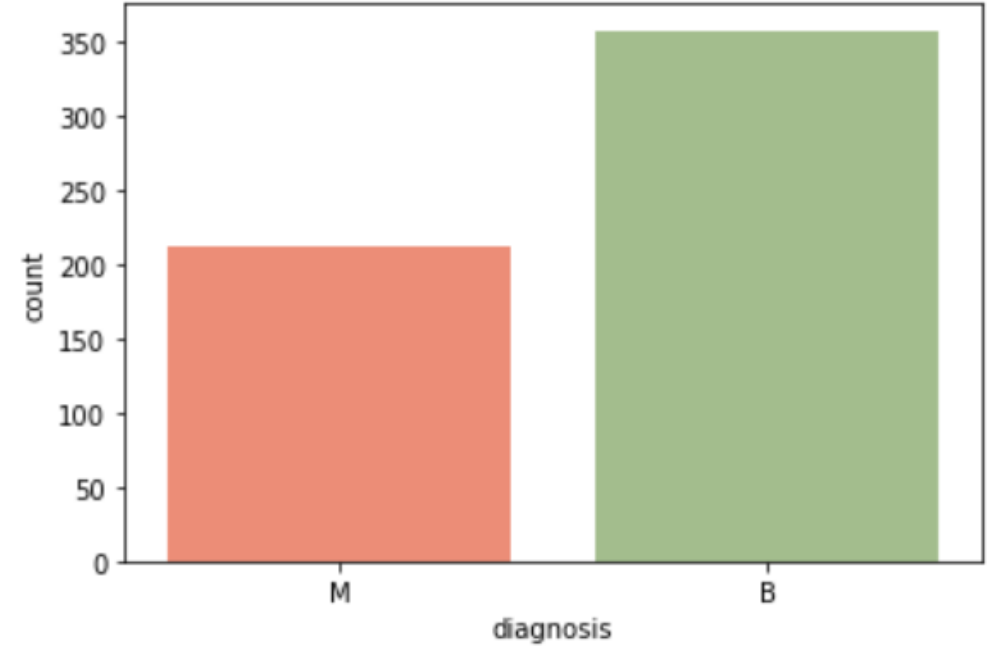
33 Columns
569 Rows

Containing
measurements of
tumors: radius,
texture, compactness,
etc.



357 Benign
212 Malignant

Classified as
B-Benign, or
M-Malignant



symmetry_worst	fractal_dimension_worst
0.4601	0.11890
0.2750	0.08902

diagnosis	radius_mean	texture_mean	perimeter_mean
M	17.99	10.38	122.80
M	20.57	17.77	132.90

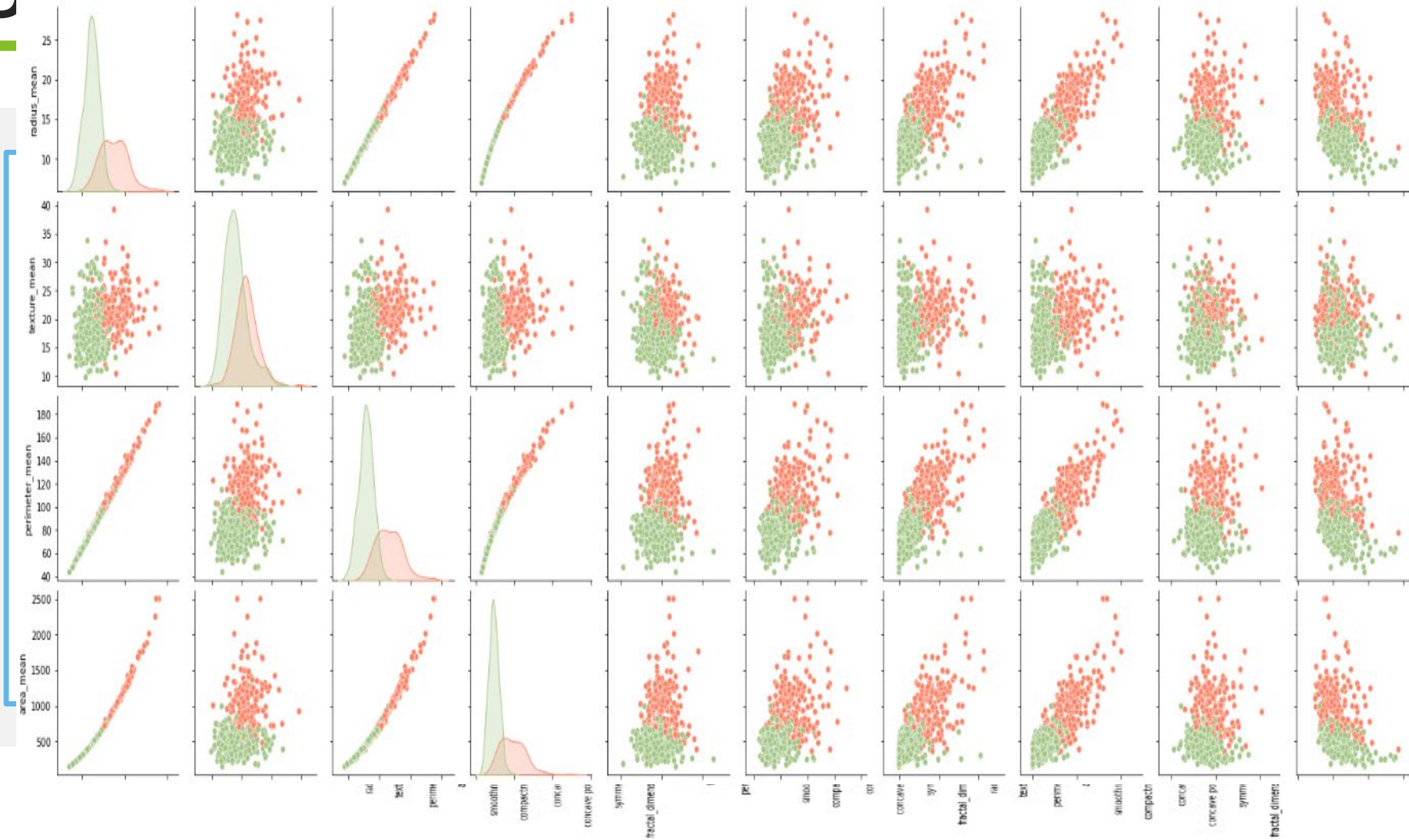
area_mean	smoothness_mean	compactness_mean	concavity_mean
1001.0	0.11840	0.27760	0.3001
1326.0	0.08474	0.07864	0.0869



D

diagnosis	1	0.73	0.42	0.74	0.71	0.36	0.6	0.7	0.78	0.35	0.013	0.57	0.008	0.5	0.55	0.067	0.29	0.25	0.41	0.0063	0.78	0.78	0.46	0.78	0.73	0.42	0.59	0.66	0.79	0.42	0.32
radius_mean	0.73	1	0.32	1	0.99	0.17	0.51	0.68	0.02	0.15	-0.31	0.68	0.097	0.67	0.74	-0.22	0.21	0.19	0.38	-0.1	-0.043	0.97	0.3	0.97	0.94	0.12	0.41	0.53	0.74	0.16	0.0071

1.0



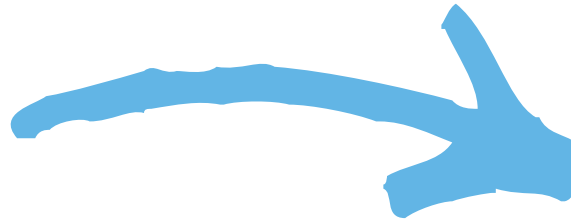
2



Pre-processing



**Split Dataset
into Training and
Testing Subsets**



Scale the Data

Modeling

0	83	5
1	3	52
	0	1

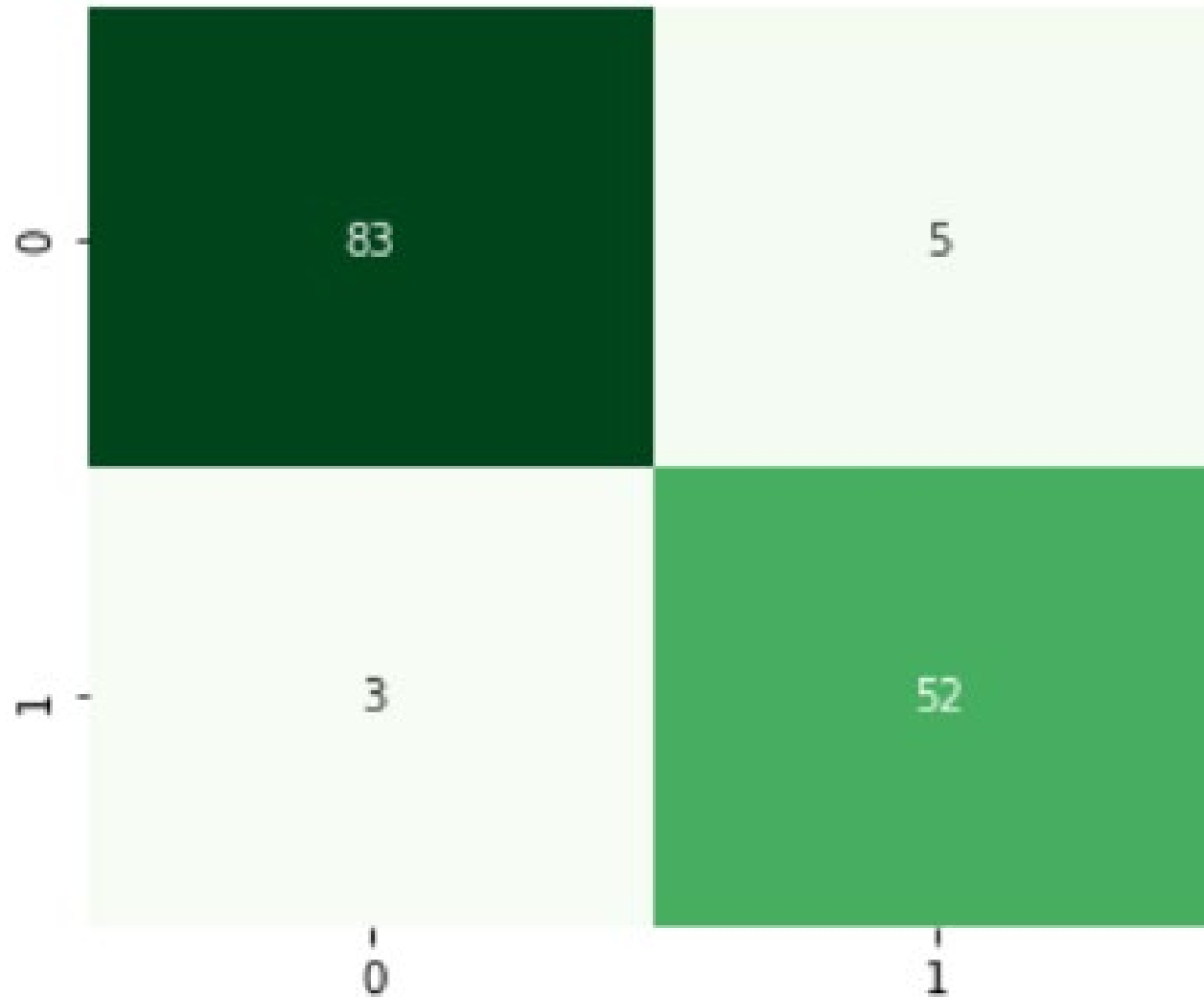
Baseline Model

Results:

95% Recall

94% Accuracy

Modeling



0	83	5
1	3	52
	0	1

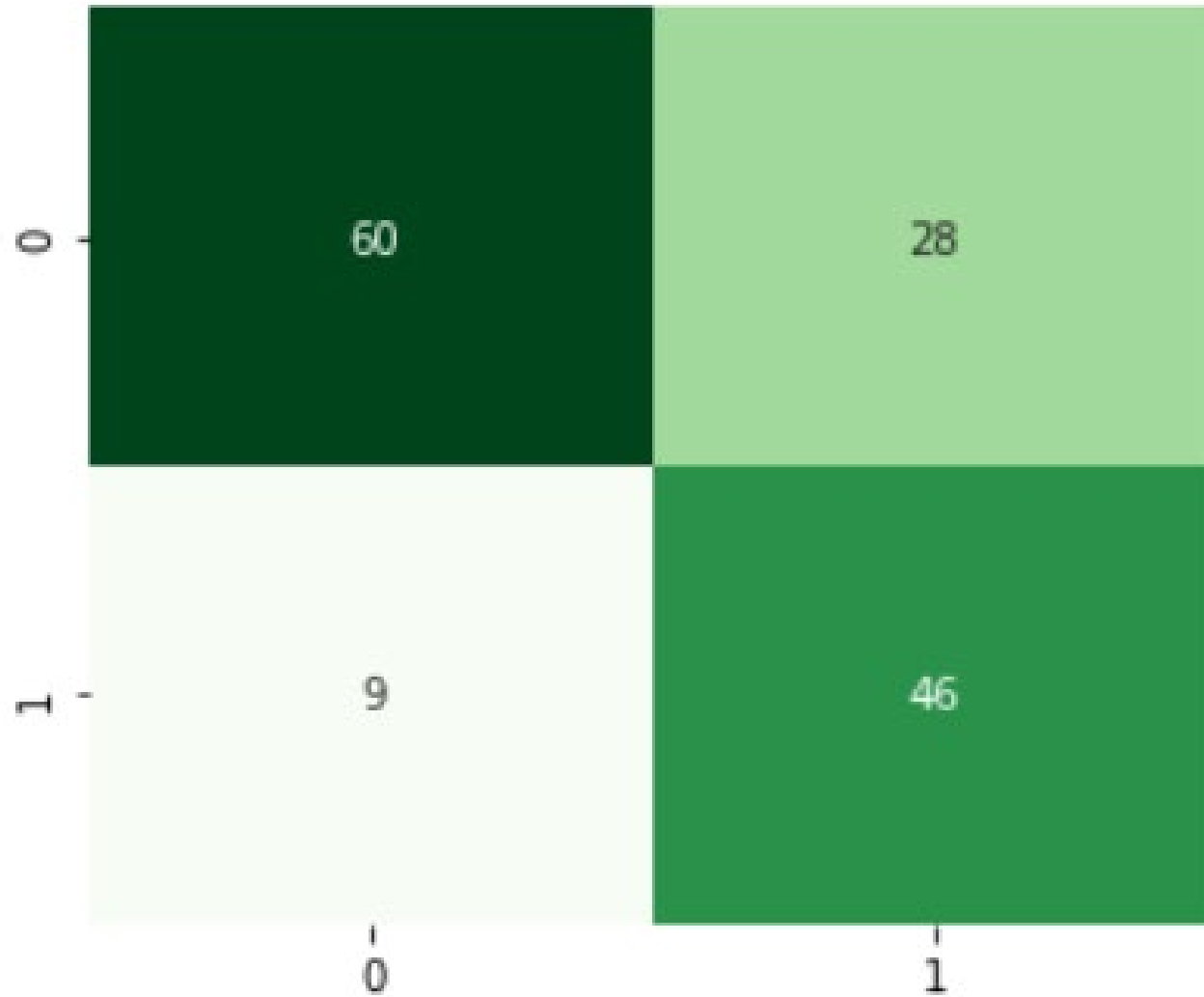
Tuned Logistic Regression Model

Tuned with a GridSearch

Results:
95% Recall

94% Accuracy

Modeling



0	60	28
1	9	46
	0	1

Tuned Decision Tree Model

Tuned with a GridSearch

83.6% Recall

74% Accuracy



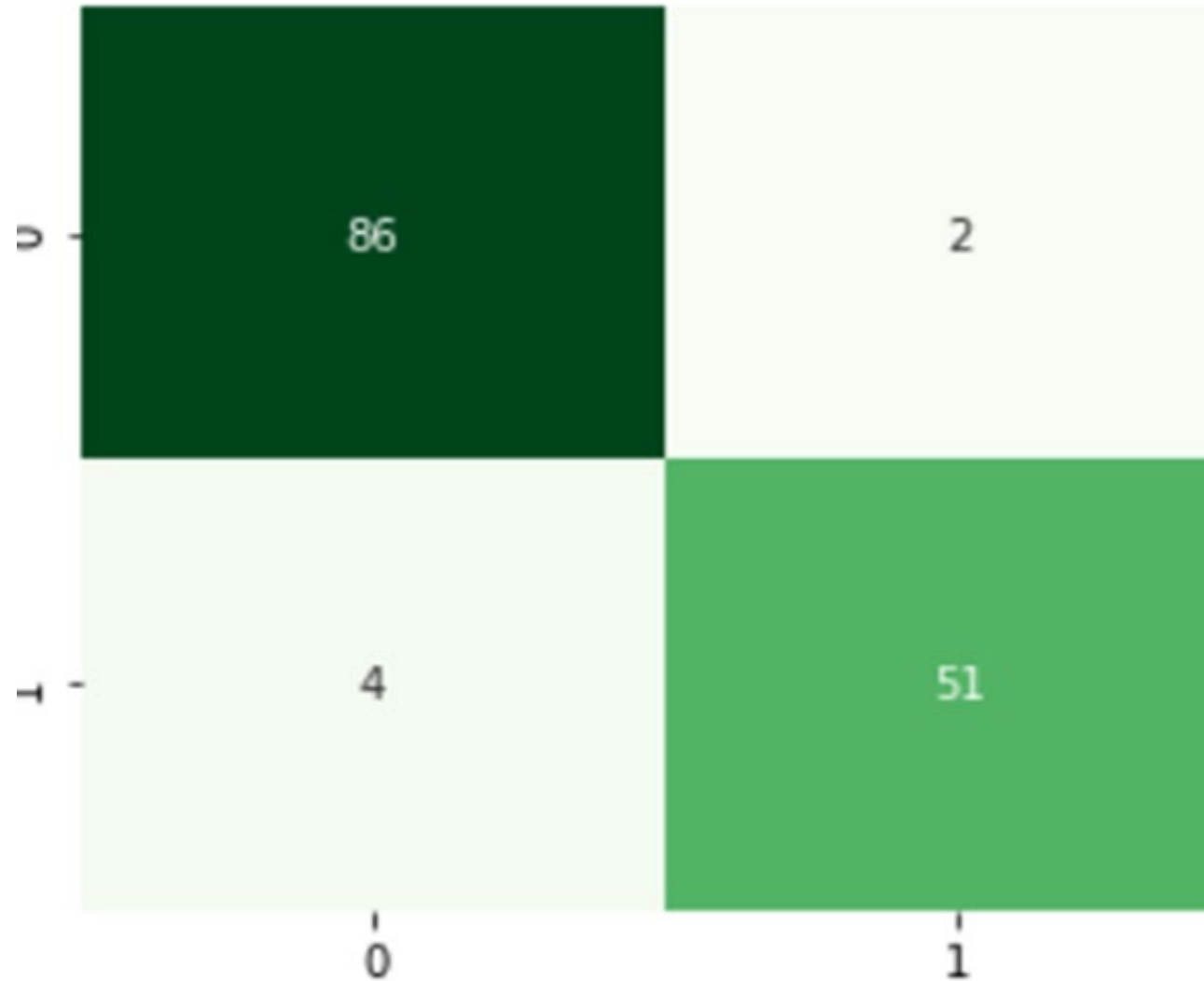
Modeling

82	6
4	51
0	1

Decision Tree Model

Results:
92.7% Recall
93% Accuracy

Modeling

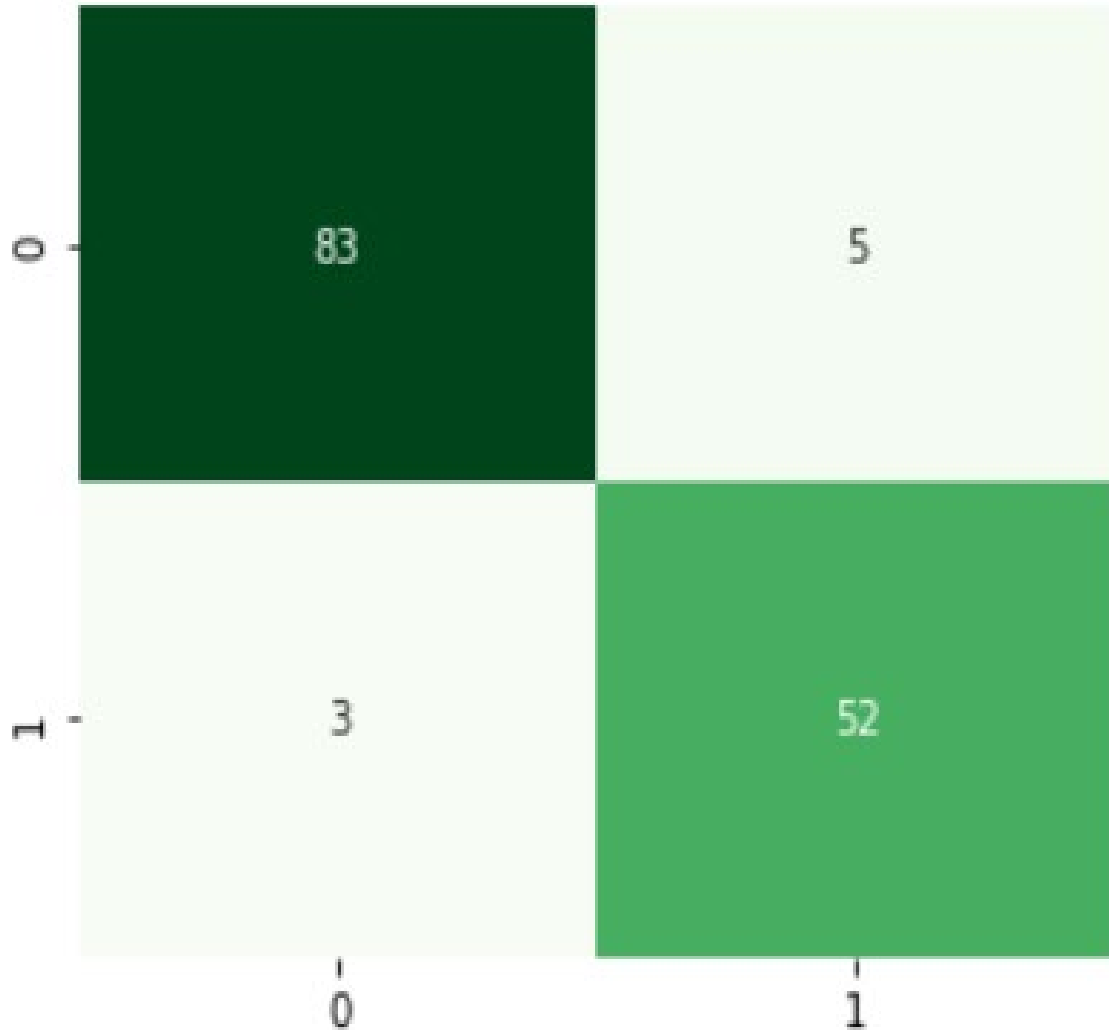


0	86	2
1	4	51
	0	1

Random Forest Model

Results:
92.7% Recall
93% Accuracy

Optimal Model - Logistic Regression Model



0	83	5
1	3	52
	0	1

Results

Out of 100 tumors observed, approximately 2 malignant tumors will be incorrectly classified as benign.

95% recall, 94% accuracy

Next Steps: tune so that false negatives are minimized.



Thank You!

Questions?