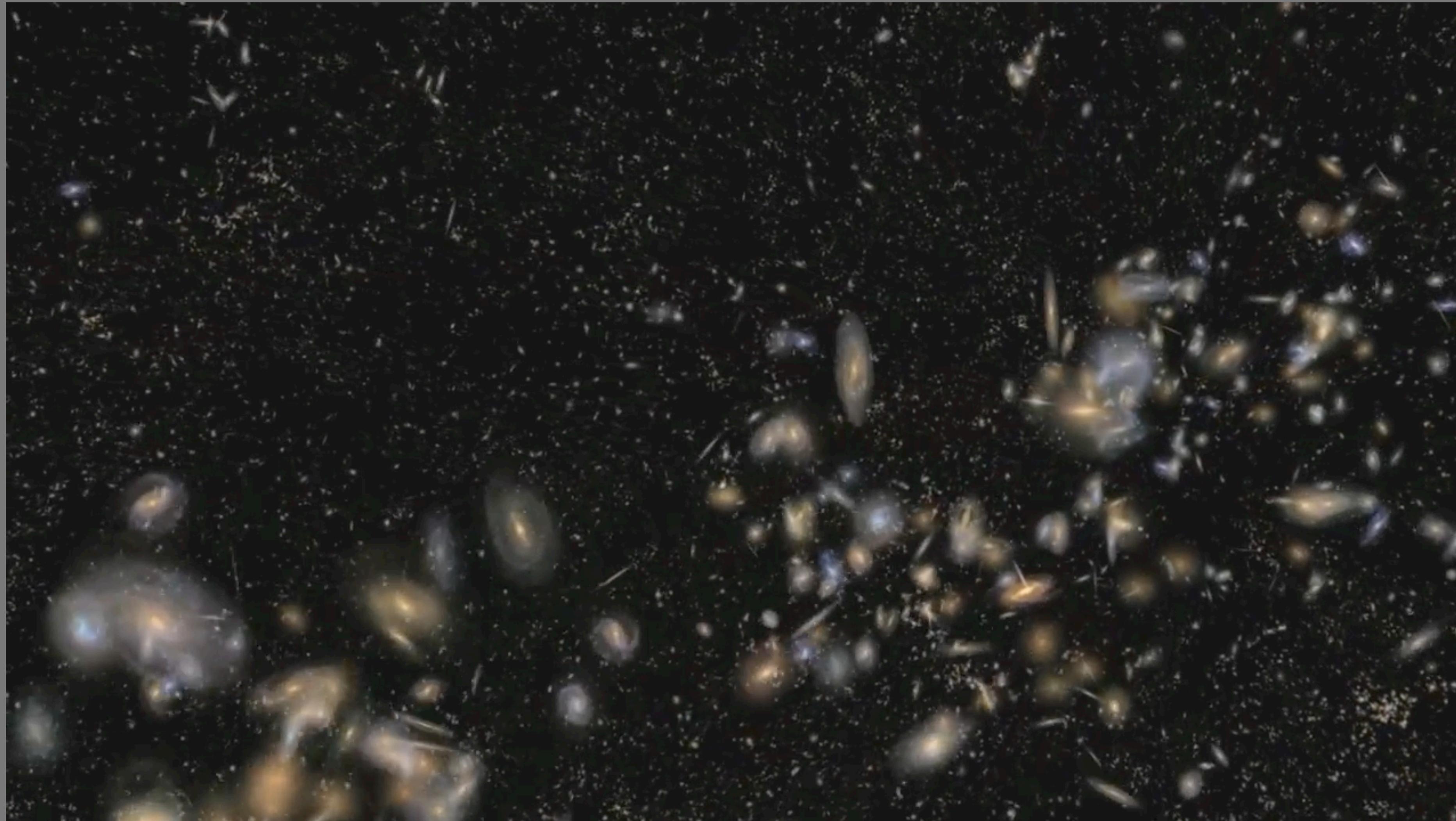


https://github.com/igc5972/comp_meth/tree/master/comp_final

Logistical Regression For Color-Mass Diagrams

Isabella G. Cox
ASTP 720 Computational Methods
Fall 2020

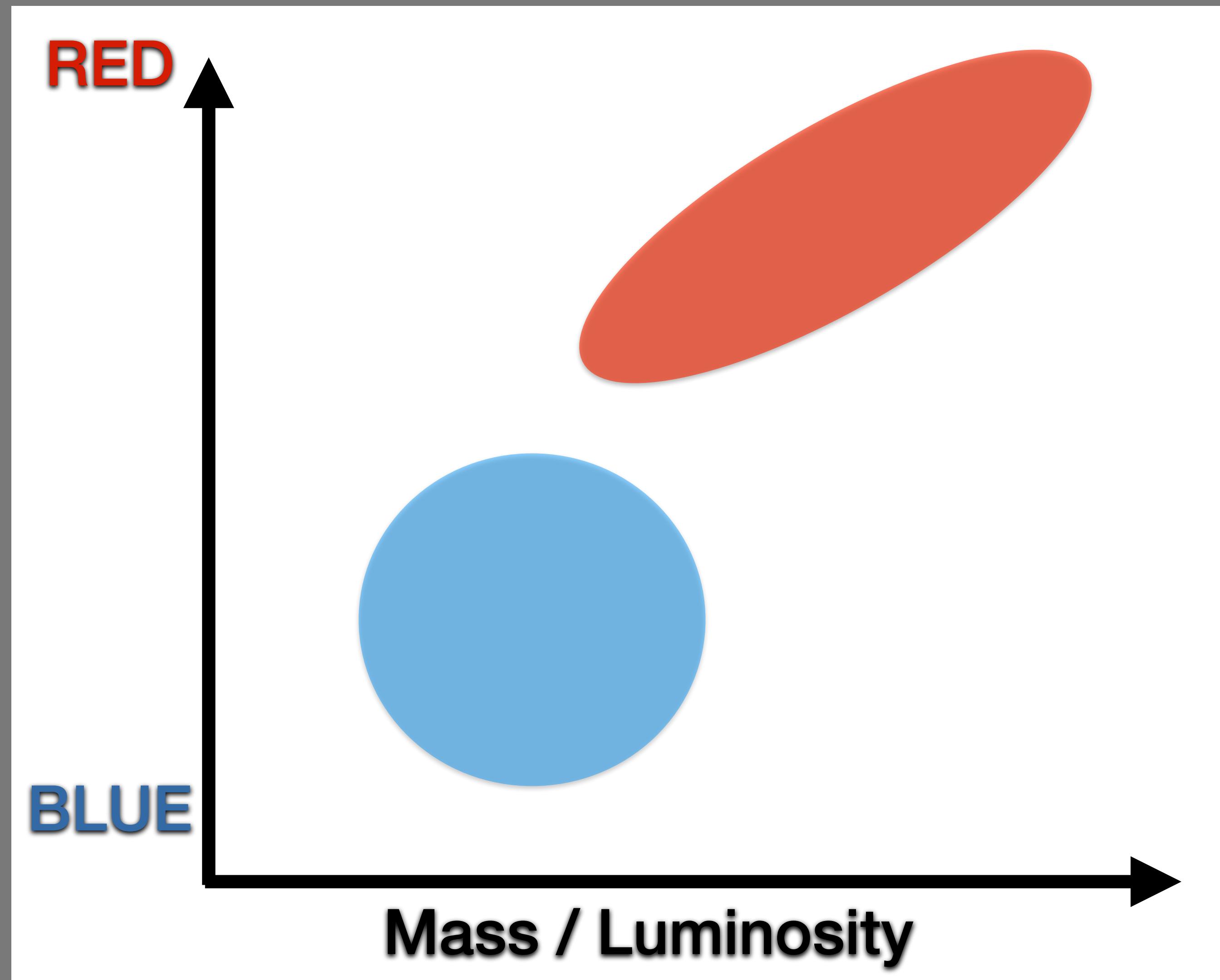
Sloan Digital Sky Survey



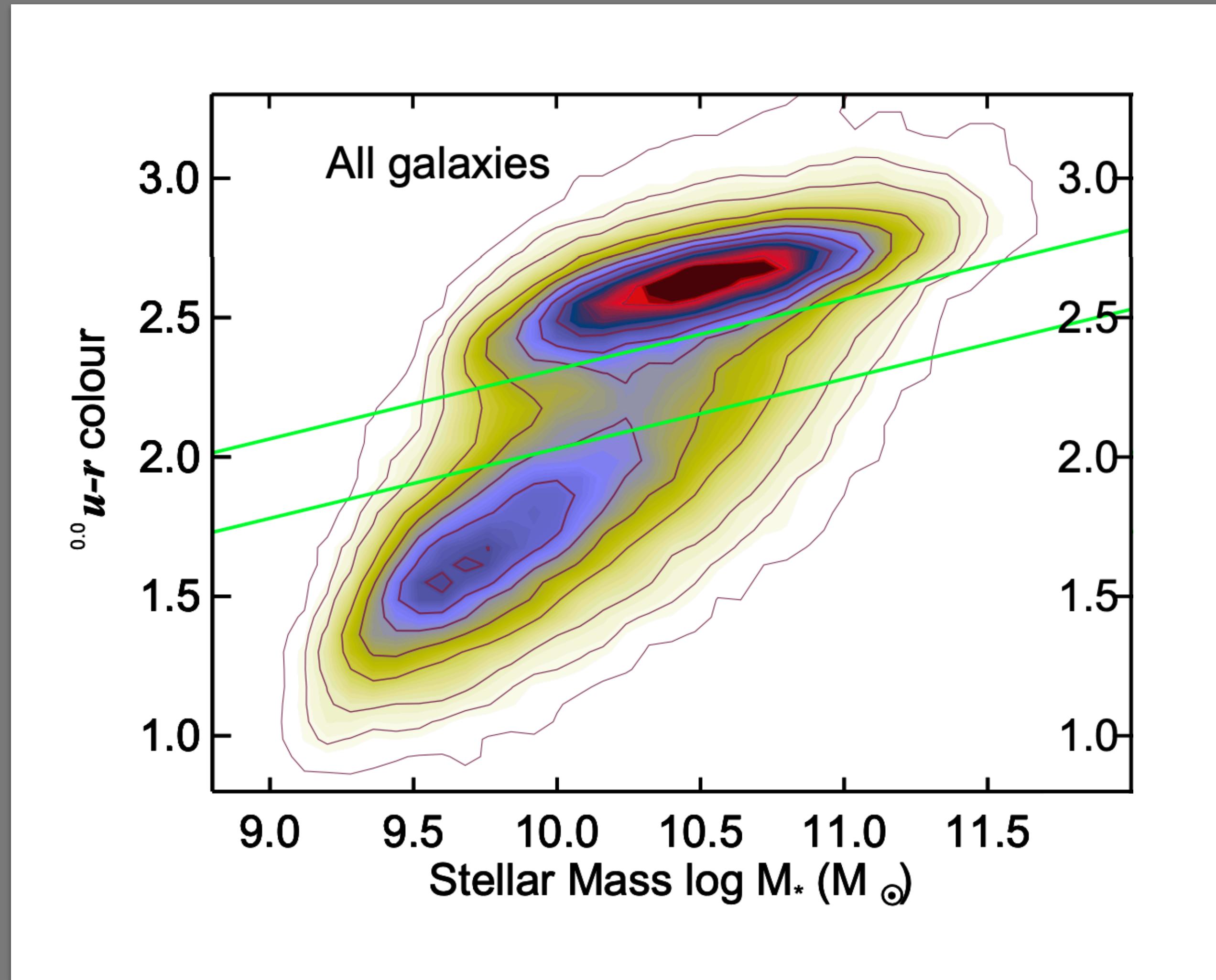
How do we sort through this vast amount of scientific data?

**Look at *samples* of many galaxies,
instead of individuals.**

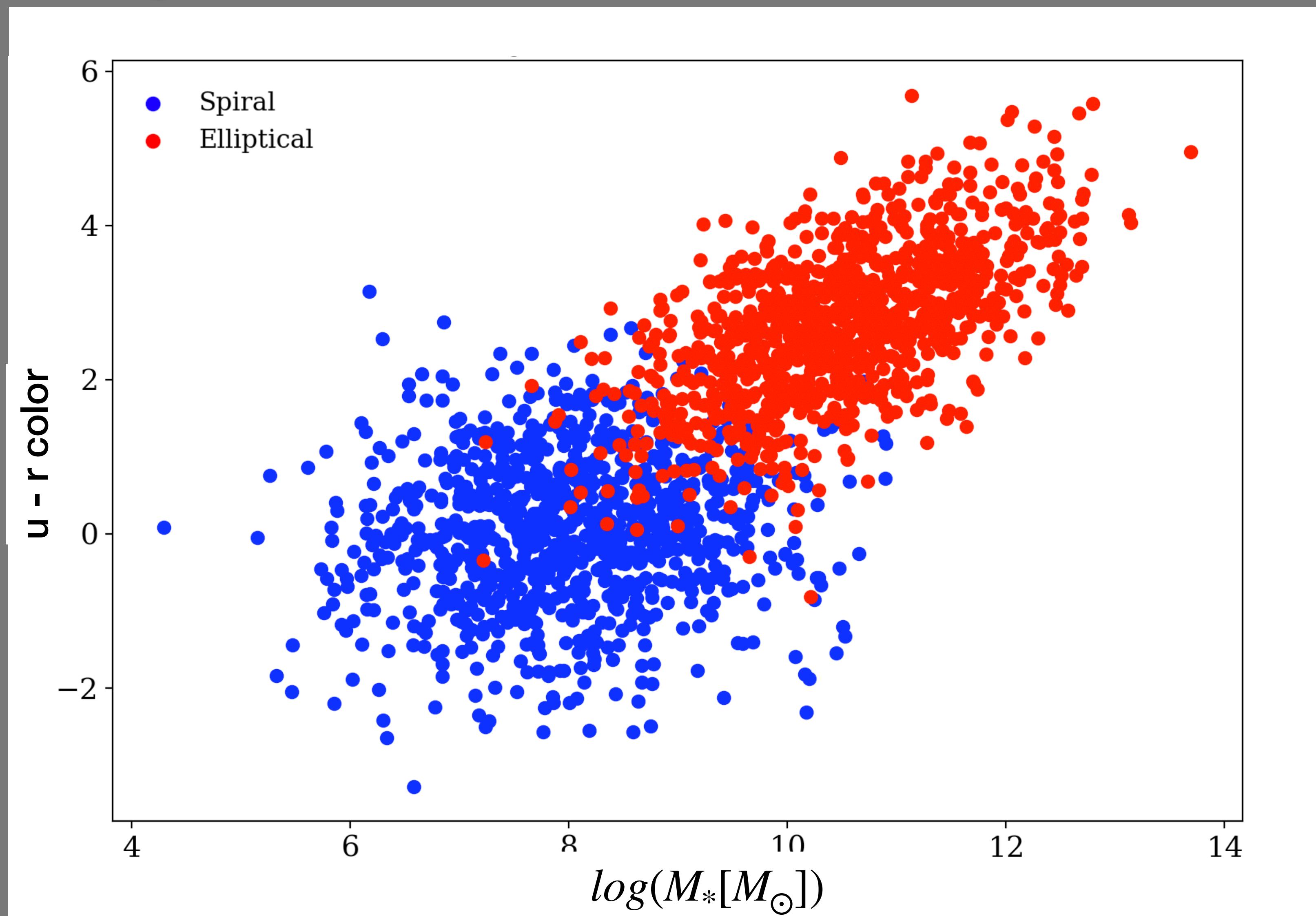
Color Bi-Modality of Galaxies



Color - Mass Diagram



My “Simulated” SDSS Data



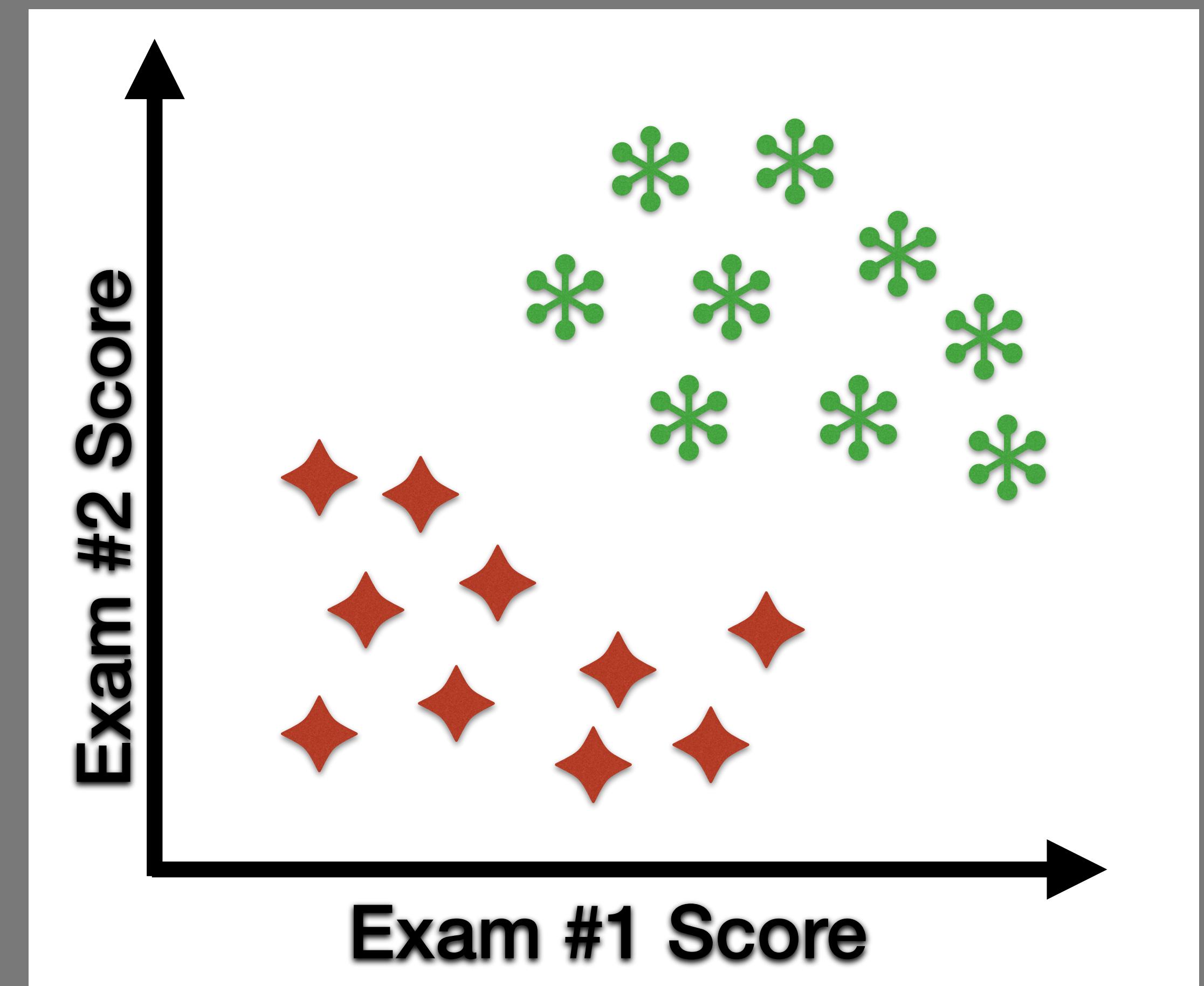
Science Question

If I assign labels to the simulated galaxies, can I correctly predict their label by drawing a line demarcating the data of a color-mass diagram into two populations?

Logistic Regression Example

Classification Algorithm. The dependent variable is predicted label for each datapoint, determined with respect to demarcation line separating populations

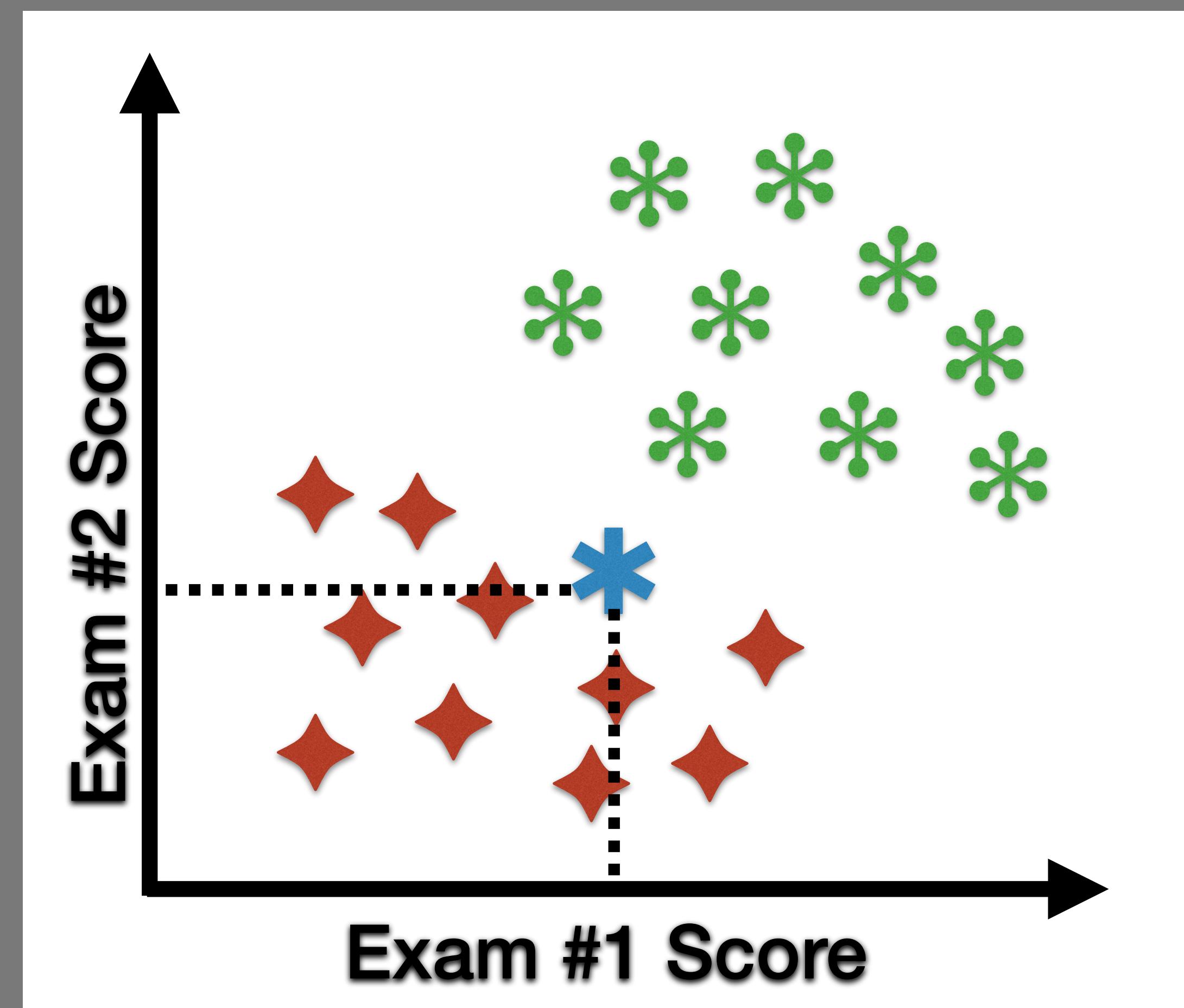
- Admitted Students
- Rejected Students



Logistic Regression Example

Classification Algorithm. The dependent variable is predicted label for each datapoint, determined with respect to demarcation line separating populations

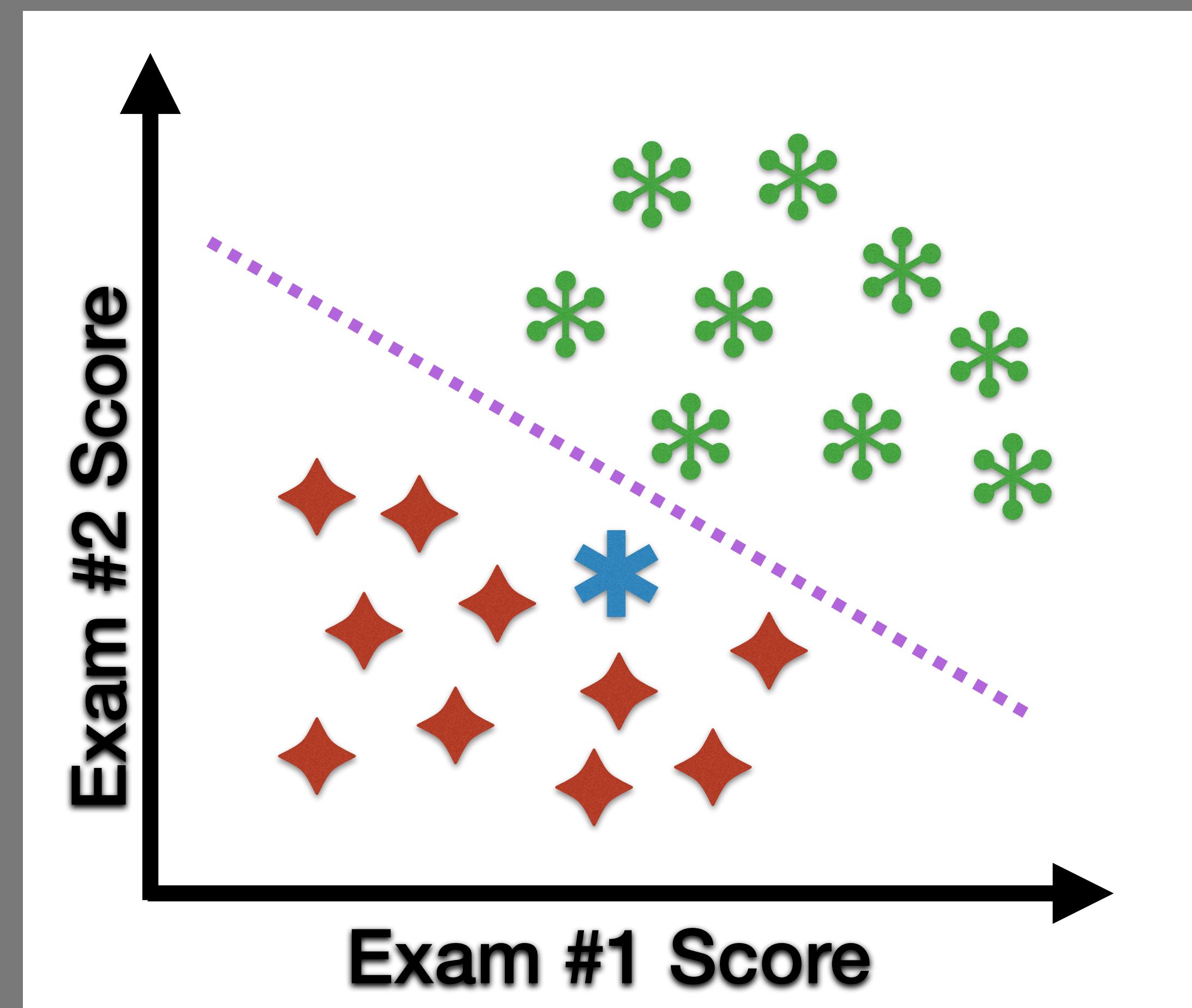
- Admitted Students
- Rejected Students
- New Student Data



Logistic Regression Example

Classification Algorithm. The dependent variable is predicted label for each datapoint, determined with respect to demarcation line separating populations

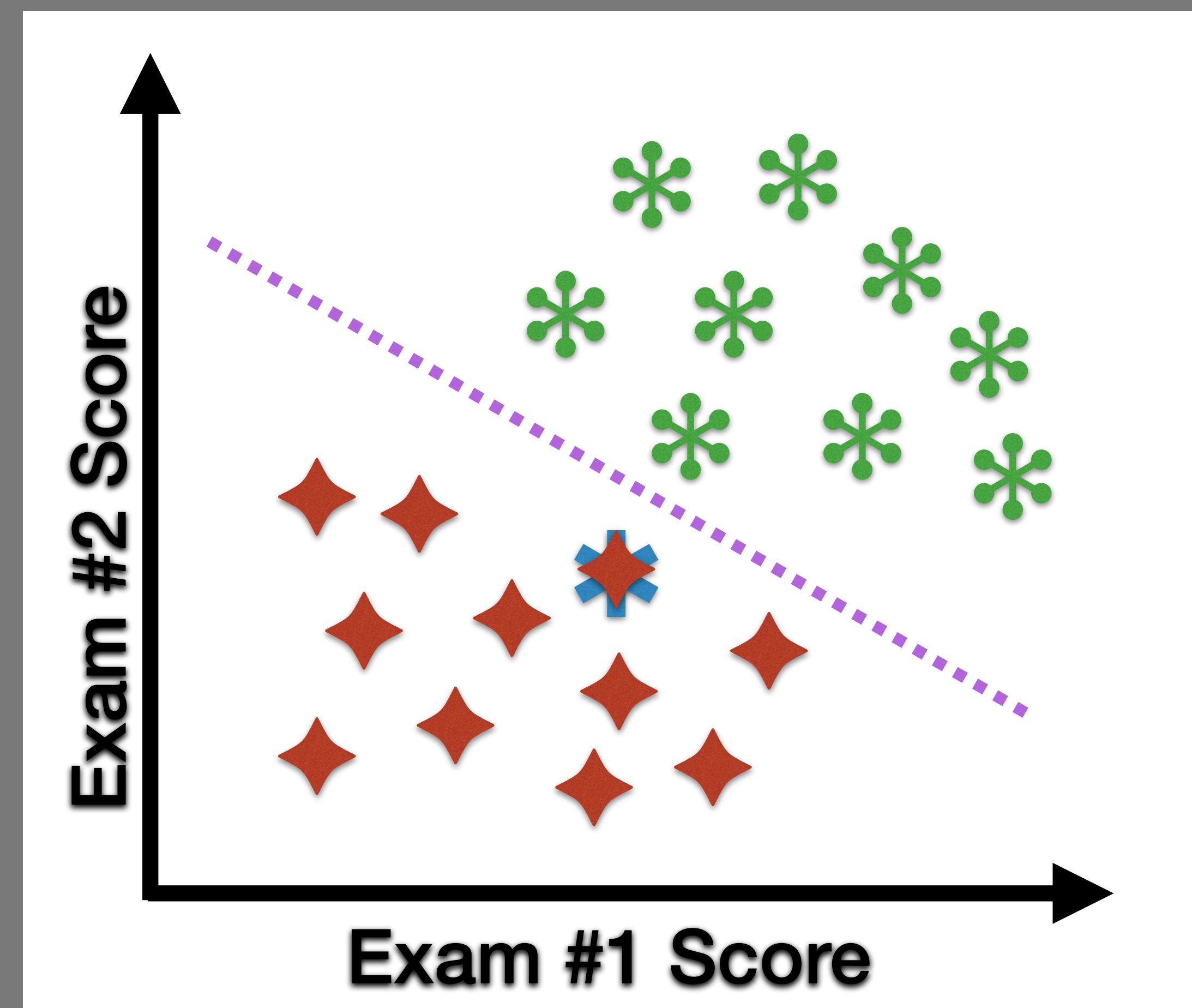
- Admitted Students
- Rejected Students
- New Student Data



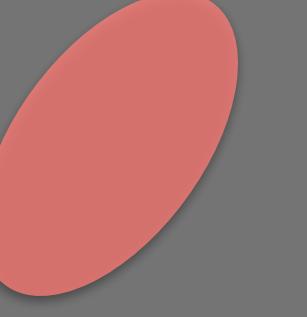
Logistic Regression Example

Classification Algorithm. The dependent variable is predicted label for each datapoint, determined with respect to demarcation line separating populations

- Admitted Students
- Rejected Students
- New Student Data



Terminology

θ = parameters	y = labels	X = features
$\theta_1, \theta_2, \theta_3$ Will be used to construct eq. of line for decision boundary.	Elliptical [1]  Spiral [0] 	X_1 : Color [float] X_2 : Stellar Mass [float]

The Plan

Go through a large number of iterations, and for each step.....

The Plan

Go through a large number of iterations, and for each step.....

(1) Calculate the sigmoid functional value based on current θ values

The Plan

Go through a large number of iterations, and for each step.....

- (1) Calculate the sigmoid functional value based on current θ values
- (2) Calculate the corresponding gradient rate to minimize the cost function

The Plan

Go through a large number of iterations, and for each step.....

- (1) Calculate the sigmoid functional value based on current θ values
- (2) Calculate the corresponding gradient rate to minimize the cost function
- (3) Decrease θ by the gradient rate times the learning rate

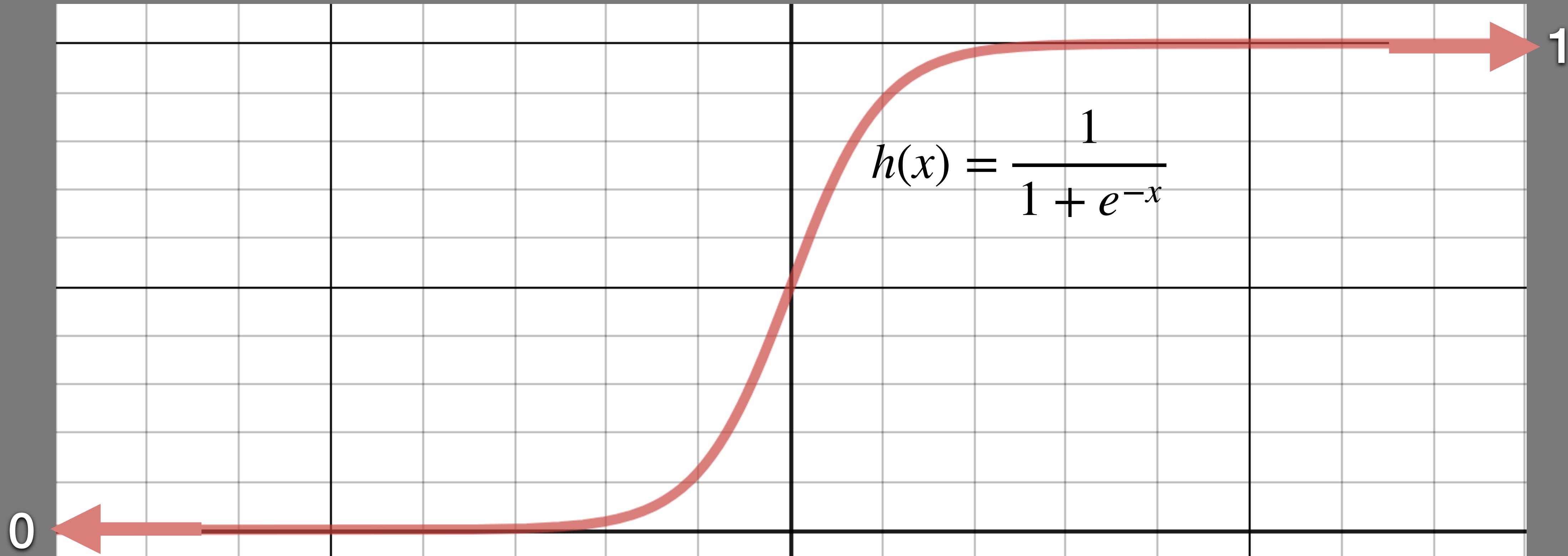
The Plan

Go through a large number of iterations, and for each step.....

- (1) Calculate the sigmoid functional value based on current θ values
- (2) Calculate the corresponding gradient rate to minimize the cost function
- (3) Decrease θ by the gradient rate times the learning rate

Once we have final θ values, we can make predictions by calculating sigmoid functional value for these final θ s and then use that to cast each entry to 0 or 1 (the binary label).

Function: Hypothesis Function

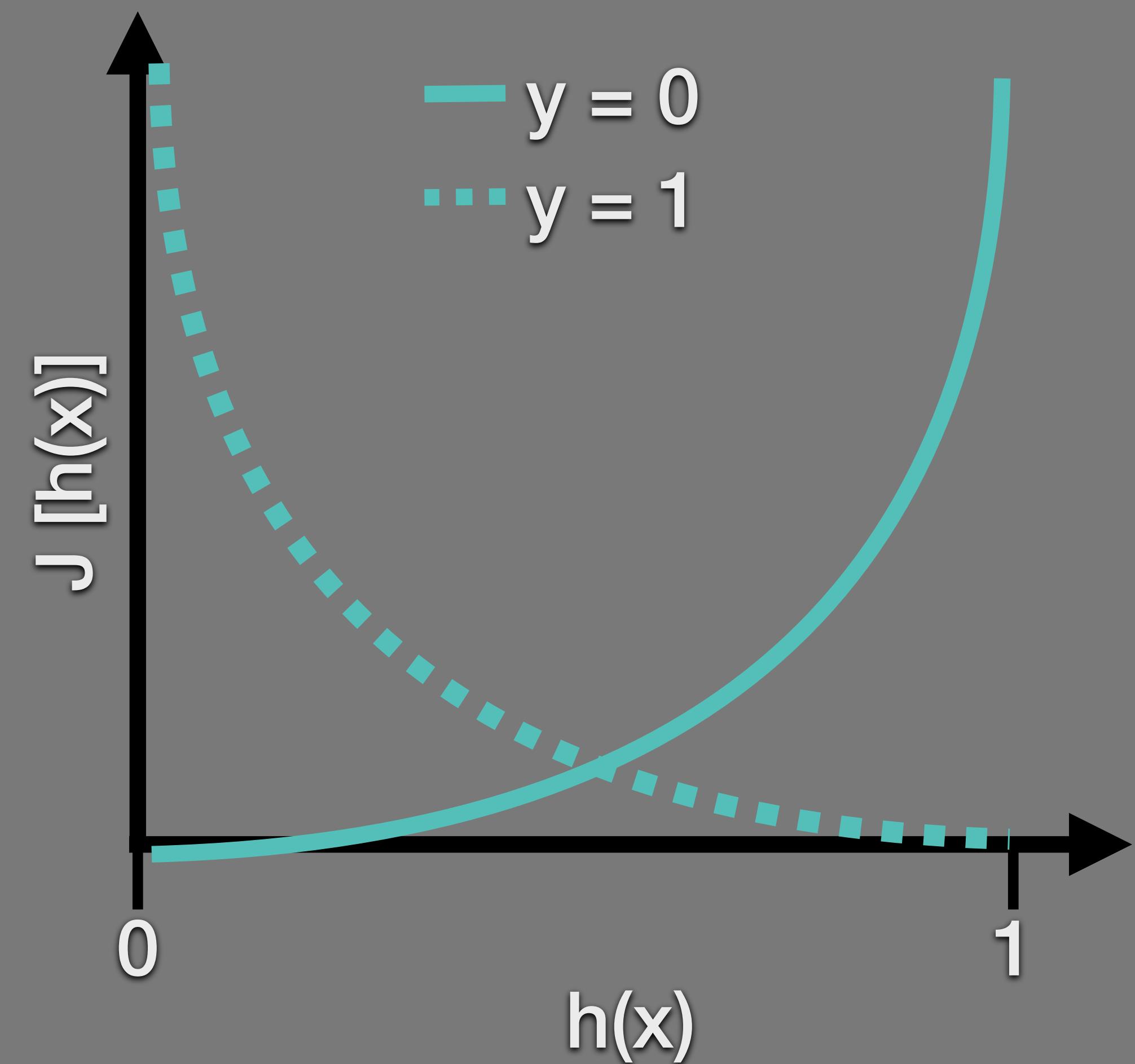


Function: Cost Function

$$J(\theta) = \begin{cases} -\log[h(x)], & \text{if } y = 1 \\ -\log[1 - h(x)], & \text{if } y = 0 \end{cases}$$

Function: Cost Function (another view)

$$J(\theta) = \begin{cases} -\log[h(x)], & \text{if } y = 1 \\ -\log[1 - h(x)], & \text{if } y = 0 \end{cases}$$



Function: Gradient

Hypothesis

$$\frac{\partial J}{\partial \theta} = \frac{[h(z) - y]X}{m}$$

Array of features

Number of data points

Putting it all together

```
for i in range(no_iter):
    z = XθT
    h = sigmoid(z)
    grad = gradient(X, h, y)
    θ -= rate * grad

# OUTPUT: θ has been populated
```

Making Predictions

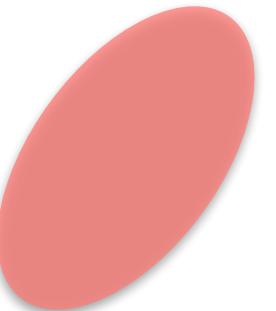
$$z = X\theta^T$$

$$s = \text{sigmoid}(z)$$

```
for i in range(len(X)):
```

```
    if s[i] >= 0.5:
```

```
        prediction = 1
```

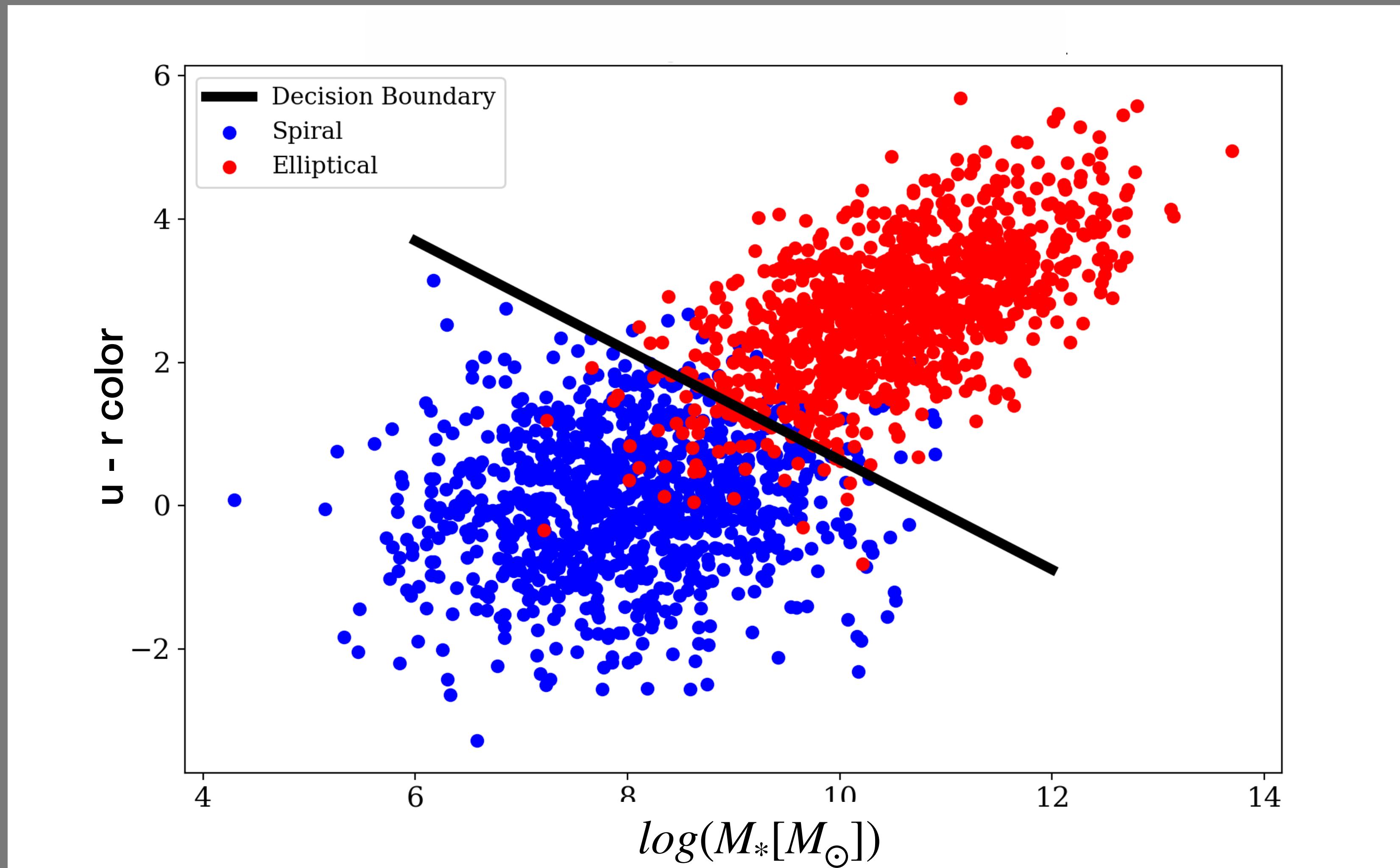


```
    elif s[i] < 0:
```

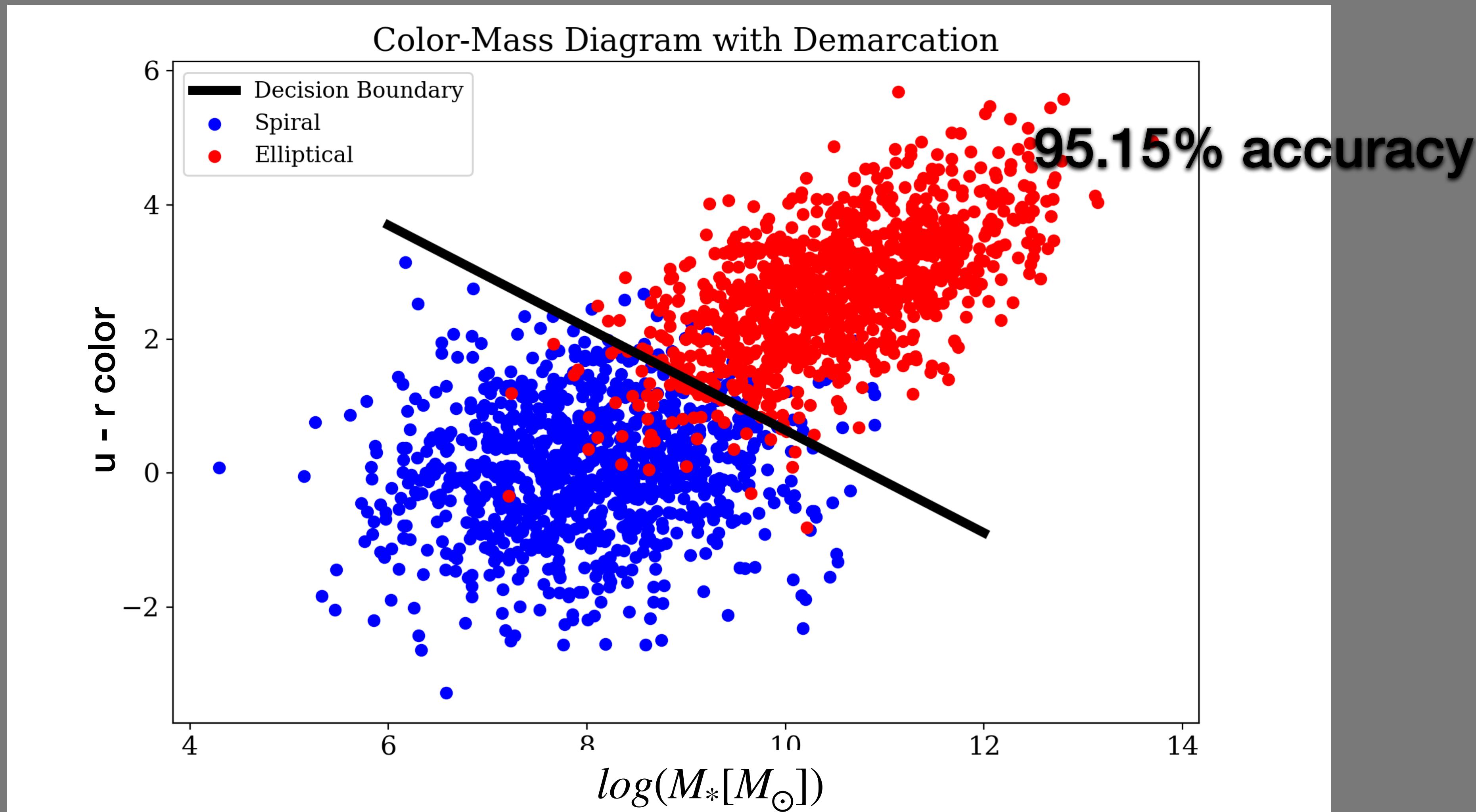
```
        prediction = 0
```



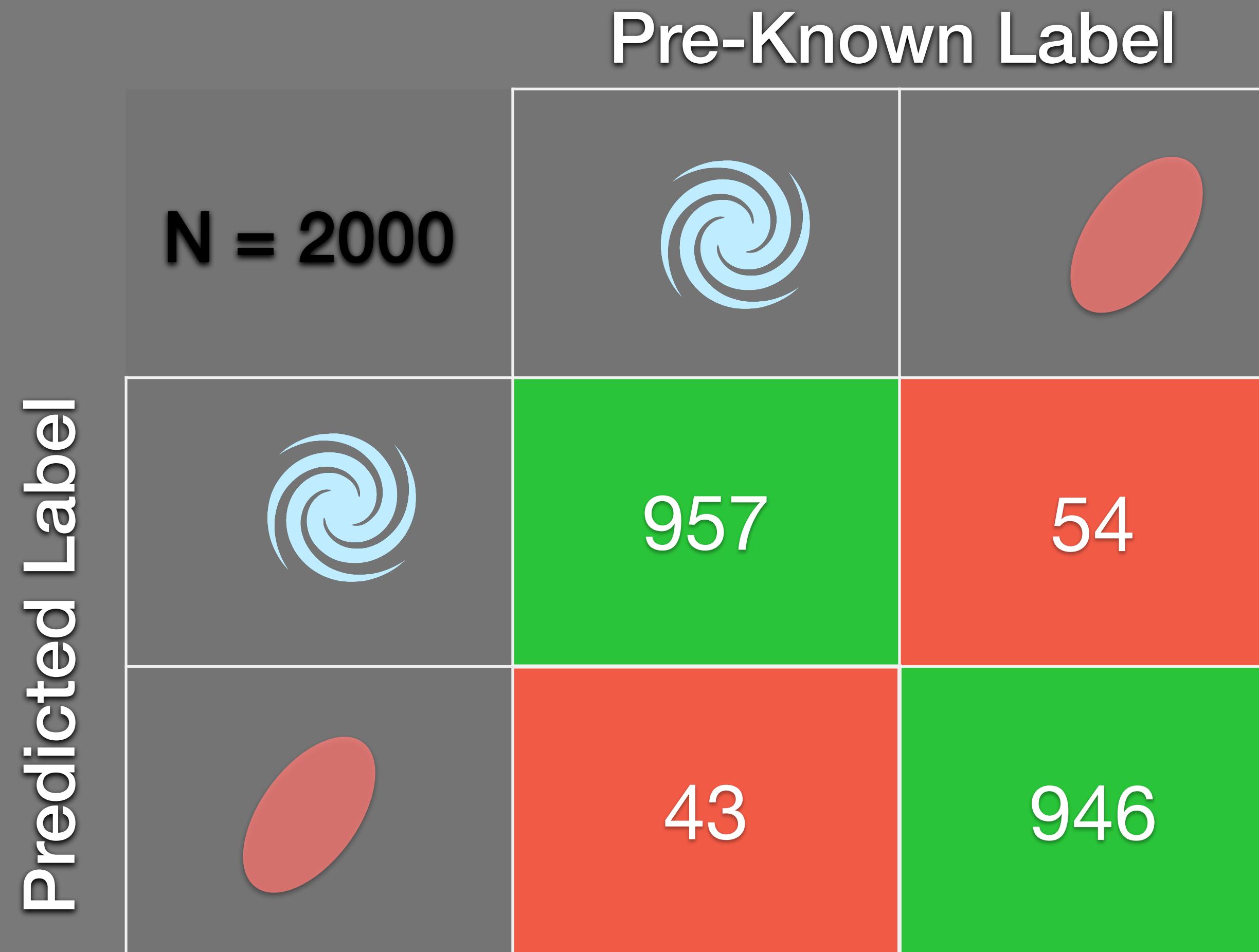
Results: Plot



Results: Plot



Results: Confusion Matrix



Future Direction

- Extend to multiple feature (from two features —> N features)
 - Find an equation that predicts label based on many values of many features

Future Direction

- Extend to multiple feature (from two features —> N features)
 - Find an equation that predicts label based on many values of many features
- Use actual observational data from SDSS and Galaxy Zoo

Questions?

References / Sources Consulted

Chandrasekaran, Dinesh. "Logistic Regression from Scratch Using Python." (2019).

Jurafsky, Daniel. "Logistic Regression." (2019).

Schawinski, Kevin et al. "The green valley is a red herring: Galaxy Zoo reveals two evolutionary pathways towards quenching of star formation in early- and late-type galaxies". Monthly Notices of the Royal Astronomical Society 440. 1(2014): 889–907.

Ungar, Lyle. "Logistic Regression." (2020).

Function: Cost Function

$$J(\theta) = \frac{-y\log(h(z)) - (1 - y)\log(1 - h(z))}{m}$$

Array of labels

Hypothesis

Number of data points

The diagram illustrates the cost function formula for a logistic regression model. The formula is:

$$J(\theta) = \frac{-y\log(h(z)) - (1 - y)\log(1 - h(z))}{m}$$

Annotations with cyan arrows point to specific components:

- An arrow points from the label "Array of labels" to the term $-y$.
- An arrow points from the label "Hypothesis" to the term $h(z)$.
- An arrow points from the label "Number of data points" to the variable m .