

# Welcome to Data Mining

Course Developer



## David Svancer

*Adjunct Professor of Business Analytics*

**George Mason University**

School of Information Systems and Operations Management

# MIS 431

## Skills You Will Develop In This Course

Fundamentals of  
Programming

The basics of R programming

Data Analysis with the  
*Tidyverse*

Data manipulation and visualization using  
the popular *tidyverse* R package

Probability and Statistics  
for Machine Learning

Applied Statistics and Probability

Machine Learning with R

Machine learning with the *tidymodels* R  
package

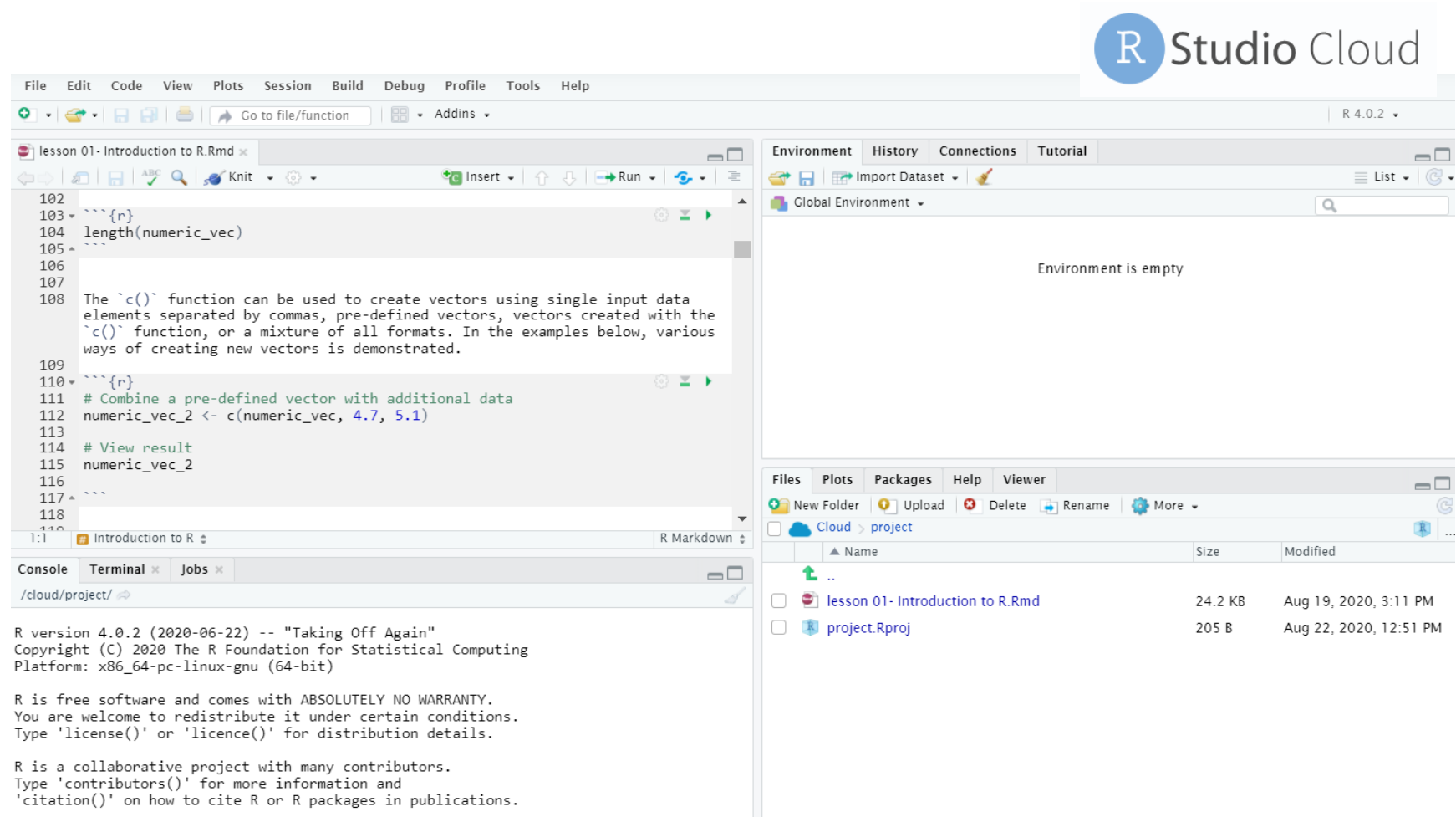
Implementing Analytics  
Projects

Communicating business value from  
analytics projects

# Course Goals

## Computer Programming Fundamentals

- Data Structures
- Writing Custom Functions
- Programming an Analytics Project



The screenshot displays the R Studio Cloud interface. The top right corner features the "R Studio Cloud" logo and the version "R 4.0.2". The main window is divided into several panes:

- Code Editor:** Shows R code for "lesson 01- Introduction to R.Rmd". The code includes comments and a function call: 

```
length(numeric_vec)
```

 and 

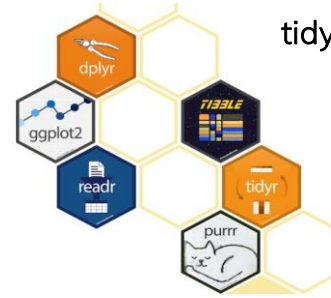
```
numeric_vec_2 <- c(numeric_vec, 4.7, 5.1)
```

. A text box explains the `c()` function.
- Environment:** Shows "Global Environment" and "Environment is empty".
- Files:** Shows a file browser with "lesson 01- Introduction to R.Rmd" (24.2 KB, Aug 19, 2020) and "project.Rproj" (205 B, Aug 22, 2020).
- Console:** Displays the R version and platform information: 

```
R version 4.0.2 (2020-06-22) -- "Taking Off Again"
Copyright (C) 2020 The R Foundation for Statistical Computing
Platform: x86_64-pc-linux-gnu (64-bit)
```

# Course Goals

## Data Analysis with the tidyverse



tidyverse.org

R packages for data science

The tidyverse is an opinionated collection of R packages designed for data science. All packages share an underlying design philosophy, grammar, and data structures.

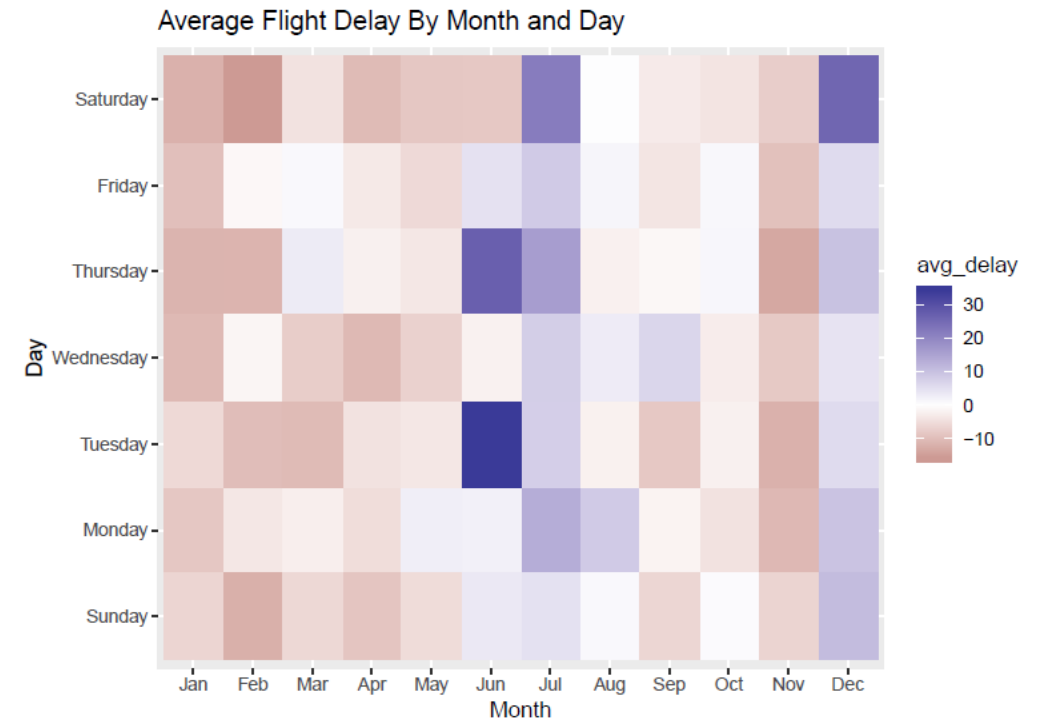
### Data Manipulation

```
heart %>% group_by(ChestPain, HeartDisease) %>%
  summarise(patients_n = n(),
            avg_chol = mean(Cholesterol),
            sd_chol = sd(Cholesterol))

# A tibble: 8 x 5
# Groups:   ChestPain [4]
  ChestPain HeartDisease patients_n avg_chol sd_chol
  <chr>      <chr>          <int>   <dbl> <dbl>
1 asymptomatic No                39     245.  48.9
2 asymptomatic Yes            103     253.  52.9
3 nonanginal No                65     247.  64.7
4 nonanginal Yes             18     239.  43.8
5 nontypical No                40     241.  45.3
```

### Data Visualization

```
ggplot(data = average_delays, mapping = aes(x = month_text, y = day_text,
                                             fill = avg_delay)) +
  geom_tile() +
  scale_fill_gradient2() +
  labs(title = "Average Flight Delay By Month and Day",
       x = "Month", y = "Day")
```



### Data Wrangling and Reshaping

Country	1999	2000
Afghanistan	745	2,666
Brazil	37,737	80,488
China	212,258	213,766

→

Country	Year	Count
Afghanistan	1999	745
Brazil	1999	37,737
China	1999	212,258
Afghanistan	2000	2,666
Brazil	2000	80,488
China	2000	213,766



# Course Goals

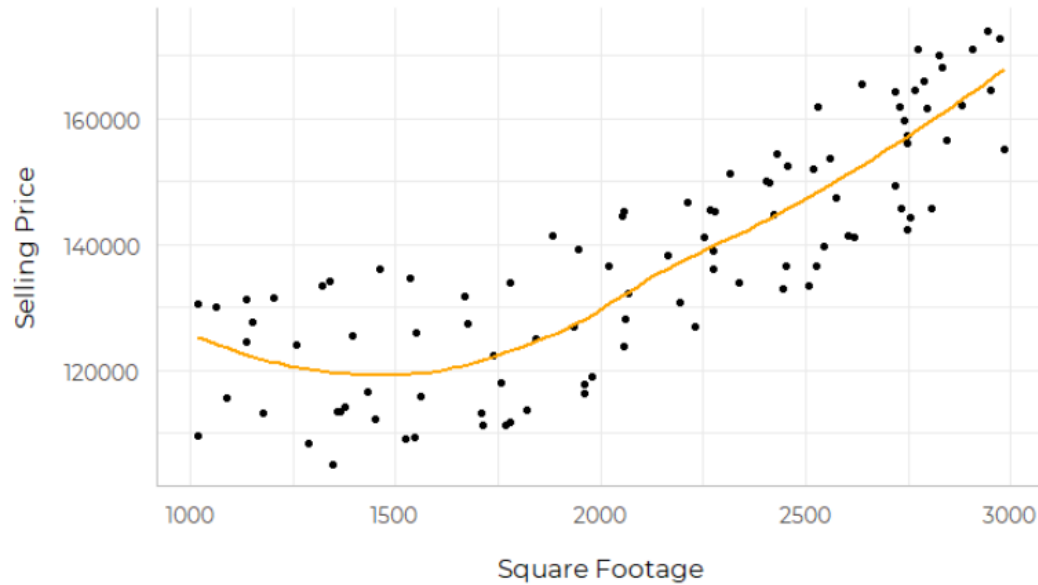
## Machine Learning with *tidymodels*

Tidymodels

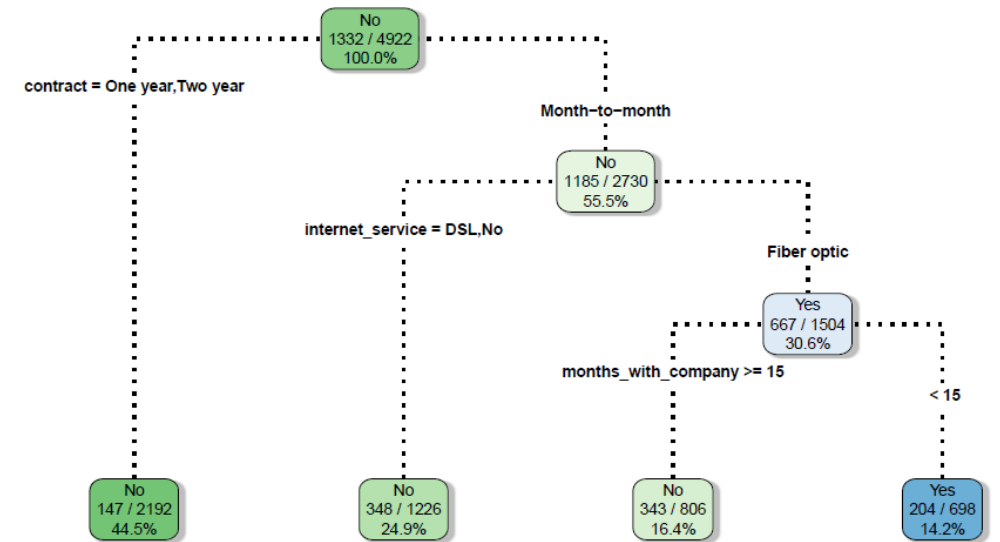


### Regression Predicting Numeric Outcomes

Predicting Home Selling Price



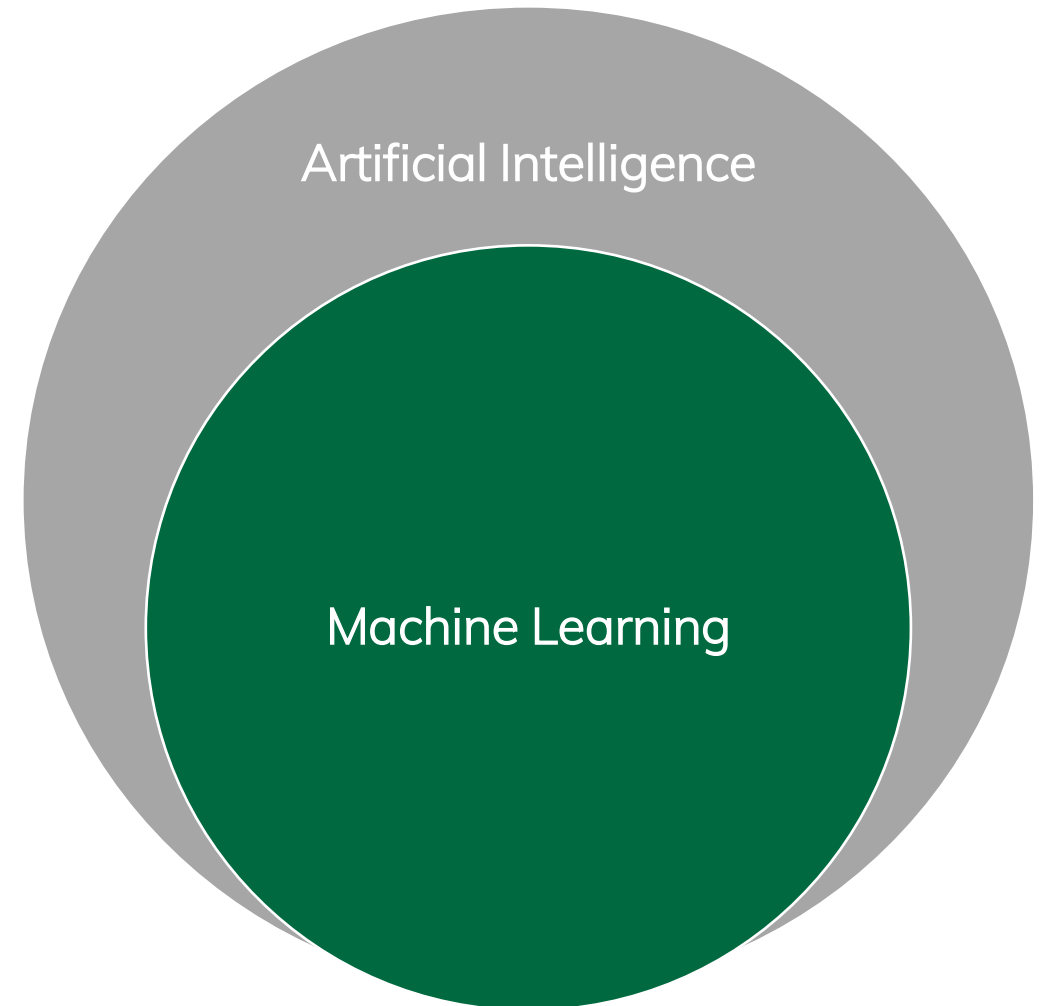
### Classification Predicting Categories



# Machine Learning

## What is Machine Learning?

A subset of *Artificial Intelligence* that gives computers the capability to learn without being *explicitly programmed*



# Machine Learning

## A New Programming Paradigm

### Before ML

Computers were explicitly programmed to achieve desired results

### Explicit Program

If input number is even → return “Yes”

If input number is odd → return “No”

### Program Execution



### Benefit

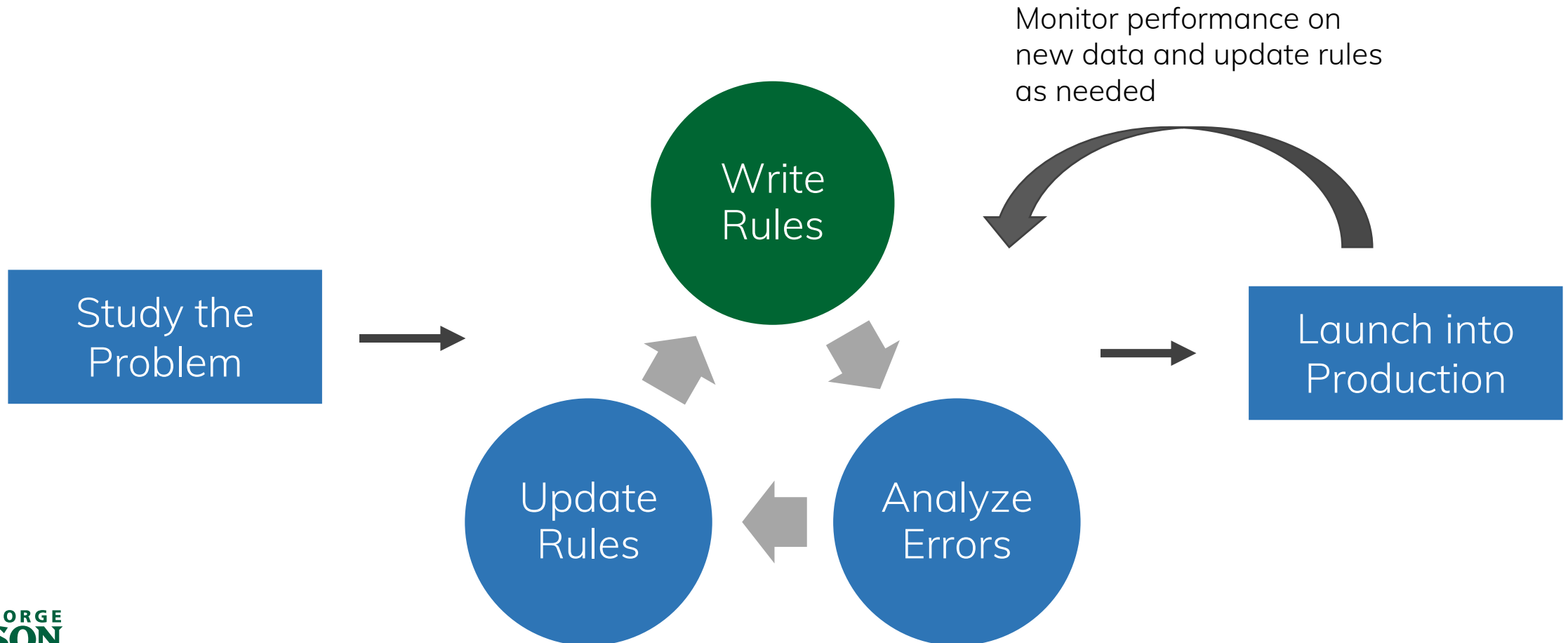
Correct output on every execution

### Challenge

All rules to accomplish task must be *known in advance*

# Machine Learning

## Explicit Programming Workflow





# Machine Learning

## Learning From Data

### Today

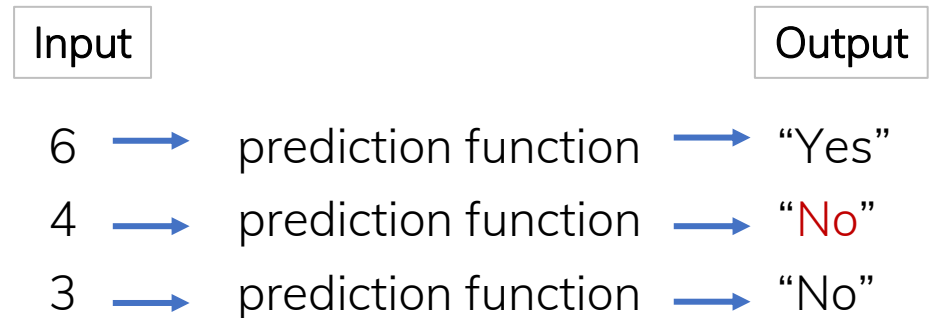
ML algorithms use vast amounts of data to discover patterns and relationships without relying on a *predetermined* equations or set of rules as a model

### ML Program

Label	Data Value
Yes	2
Yes	12
No	3
Yes	4
No	5
No	39
...	...

→ Learned prediction function

### ML Prediction Function Execution



### Benefit

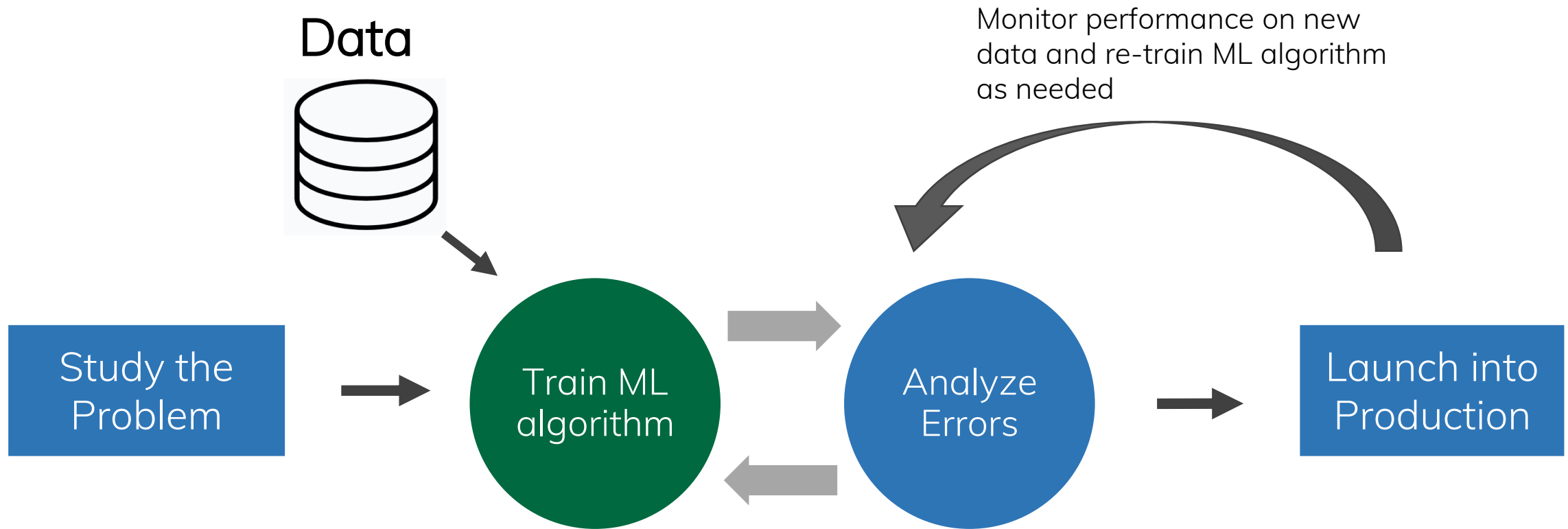
All steps/rules to accomplish **do not** have to be known or programmed explicitly

### Challenge

Prediction error

# Machine Learning

## Machine Learning Workflow



# Machine Learning

## Example - Image Recognition

### Task

Identify handwritten digits

### For a Human

Easy

### For a Computer

Extremely difficult

### MNIST Database of Handwritten Digits



# Machine Learning

## Without ML – Explicit Program

### Explicit Program to Identify Digits

Imagine having to develop explicit instructions for a program to correctly identify handwritten digits

- You must identify **every possible variation** of how digits appear and instruct a computer to label them correctly
- Practically impossible – your program would be millions of lines long!

### MNIST Database of Handwritten Digits

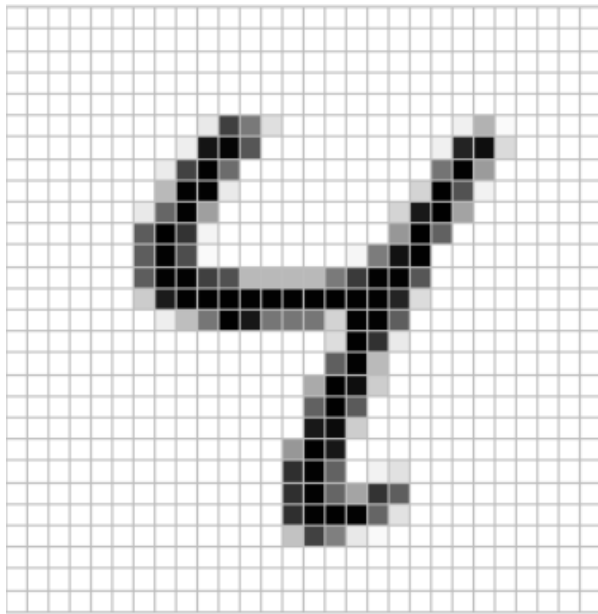


# Machine Learning

## A Machine Learning Approach

Encode Color Intensities and Apply ML Algorithms to Learn Patterns

28 x 28 image grid



Color intensities (0 – 255)

Number	Region_1	...	Region_467	Region_468	...	Region_783	Region_784
4	0		158	242		0	0
5	85		0	63		16	66
1	32		0	92		0	93
9	10		95	0		55	73
3	0		60	25		92	139
...	...	...	...	...	...	...	...

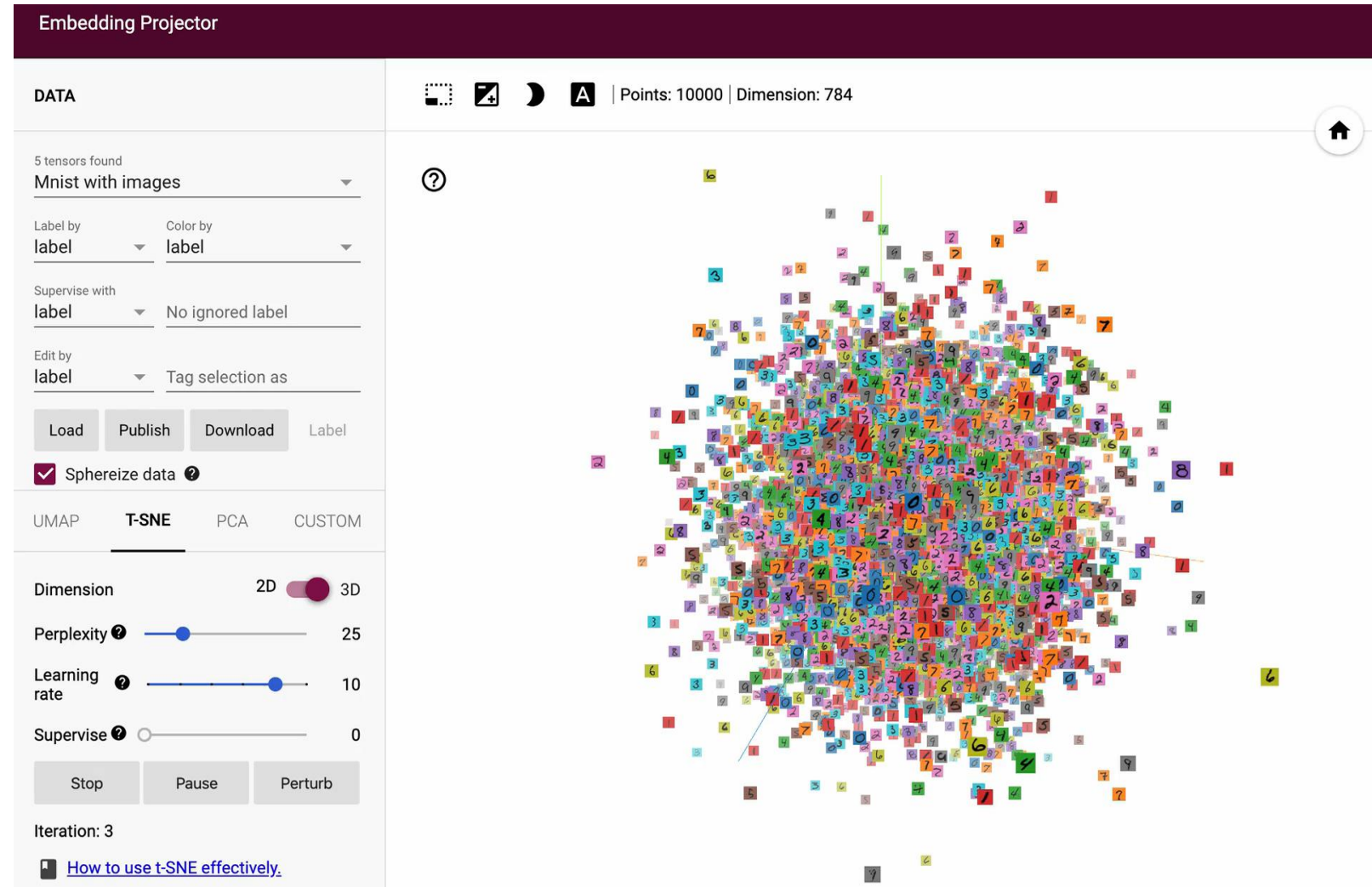
# Machine Learning

## Demonstration of ML Algorithm

TensorFlow projection tool

<https://projector.tensorflow.org/>

- o **Goal** – find the optimal way to compress digit image data to 3 dimensions so that the same digits are grouped together
- o Once this model is discovered, we can use it to predict new images based on where they fall in this 3-dimensional space



# Machine Learning Methods

## Supervised Learning

Supervised learning algorithms learn prediction functions from *labeled training data*.

Labeled data set from a hospital

- Each row represents a patient who eventually did or did not develop heart disease (*the response variable – Heart Disease*)
- Our goal might be to predict whether a new patient will develop heart disease using the predictor variables
  - For each set of predictor values, we have a known outcome
  - We also have a set of predictor values for each known outcome

Response (Target, Dependent) variable

Heart Disease	Age	Chest Pain	Resting BP	Cholesterol
No	63	typical	145	233
Yes	67	asymptomatic	160	286
Yes	67	asymptomatic	120	229
No	37	nonanginal	130	250
No	41	nontypical	130	204

Predictor (Feature, Independent) variables

# Machine Learning Methods

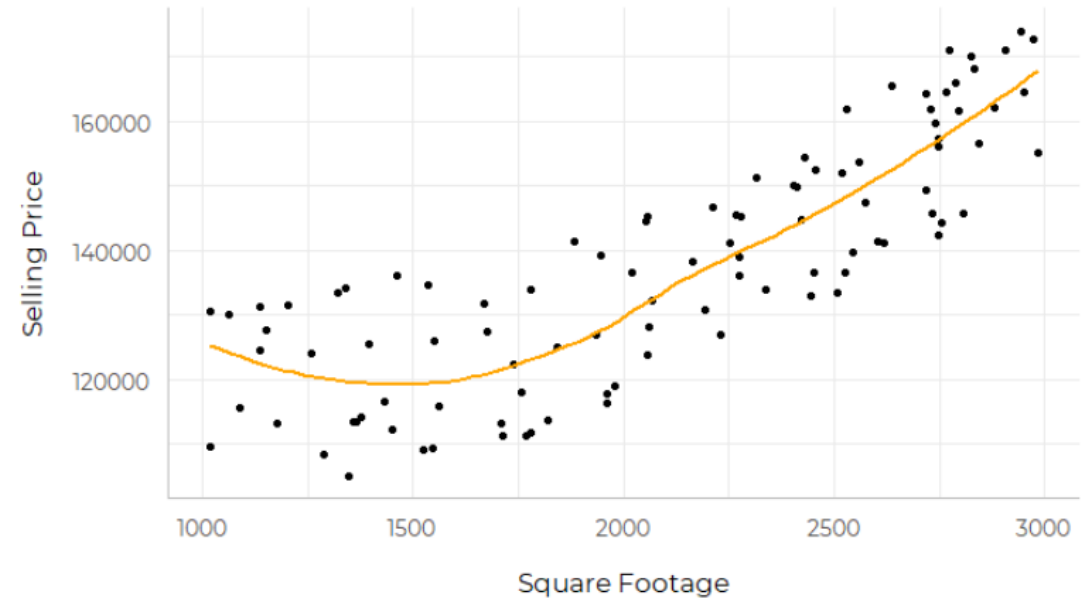
## Supervised Learning - Regression

### Regression

- Supervised learning methods are used to predict **quantitative** response variables
- Example**
  - Predict the selling price of homes using features such as square footage, age, location

Response	Predictor
Selling Price	Square Footage
\$105,667	1,100
\$118,659	1,490
\$134,268	1,850
\$165,000	2,300

Predicting Home Selling Price





# Machine Learning Methods

## Supervised Learning - Classification

### Classification

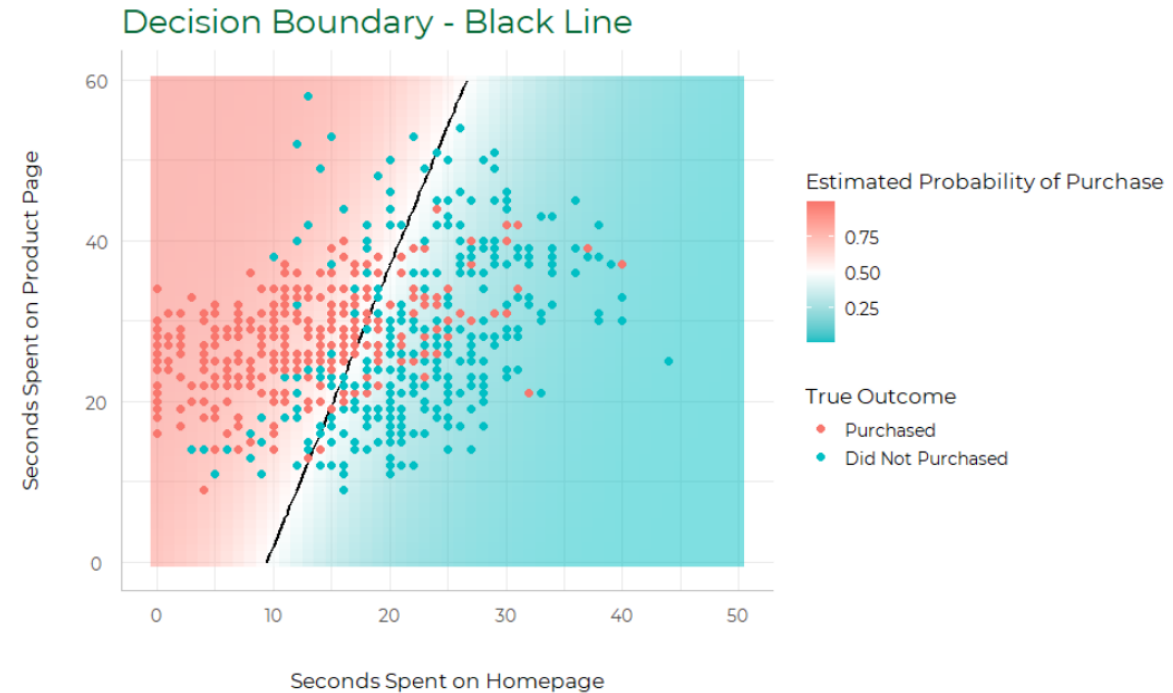
Supervised learning methods used to predict **categorical** response variables

#### Example

- Predict whether a customer will purchase a product based on the seconds they have spent browsing a company's homepage and product page

Response	Predictors	
	Seconds Homepage	Seconds Product Page
Did Not Purchase	4	30
Purchased	32	43
Did Not Purchase	2	22
Purchased	24	36

Segmenting the predictor values into distinct, non-overlapping regions to predict a category



# Machine Learning Methods

## Unsupervised Learning

In **unsupervised learning**, there are *feature* or *input* variables, but no labeled outcome variable

- no “correct” prediction

In this setting, it is typically of interest to learn the **structure** and **relationships** present in the unlabeled input data

- Methods include Clustering and Principal Components (PCA)

**Marketing Example:** Are there customer segments based on purchasing behavior?

Are there different types or species of plants present in the data below?

```
# A tibble: 150 x 4
  Sepal.Length Sepal.Width Petal.Length Petal.Width
      <dbl>      <dbl>      <dbl>      <dbl>
1         5.1         3.5         1.4         0.2
2         4.9         3         1.4         0.2
3         4.7         3.2         1.3         0.2
4         4.6         3.1         1.5         0.2
5         5         3.6         1.4         0.2
6         5.4         3.9         1.7         0.4
7         4.6         3.4         1.4         0.3
8         5         3.4         1.5         0.2
9         4.4         2.9         1.4         0.2
10        4.9         3.1         1.5         0.1
# ... with 140 more rows
```

### K-means Clustering

Finding observations that group together based on their proximity in the input data space

