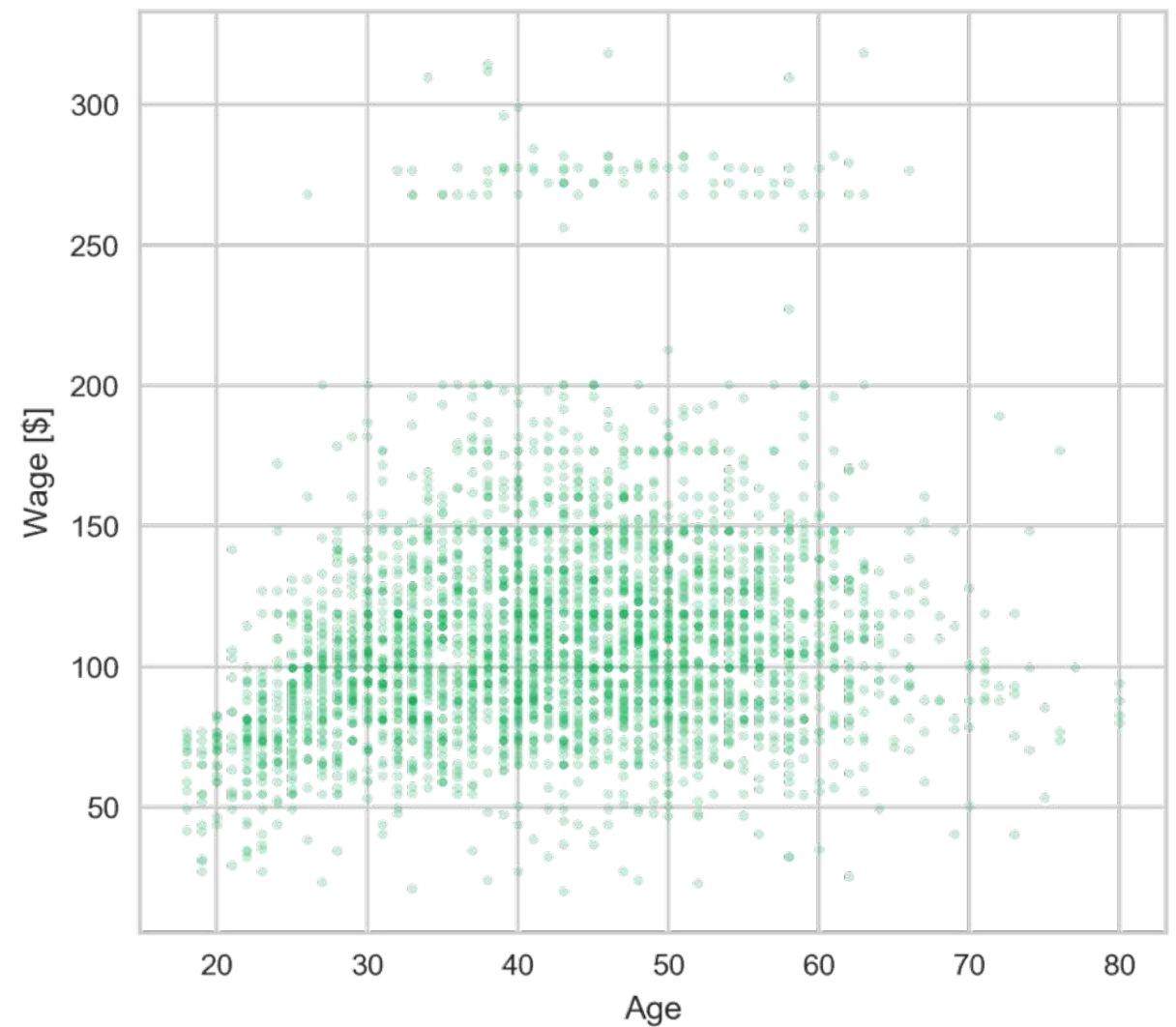
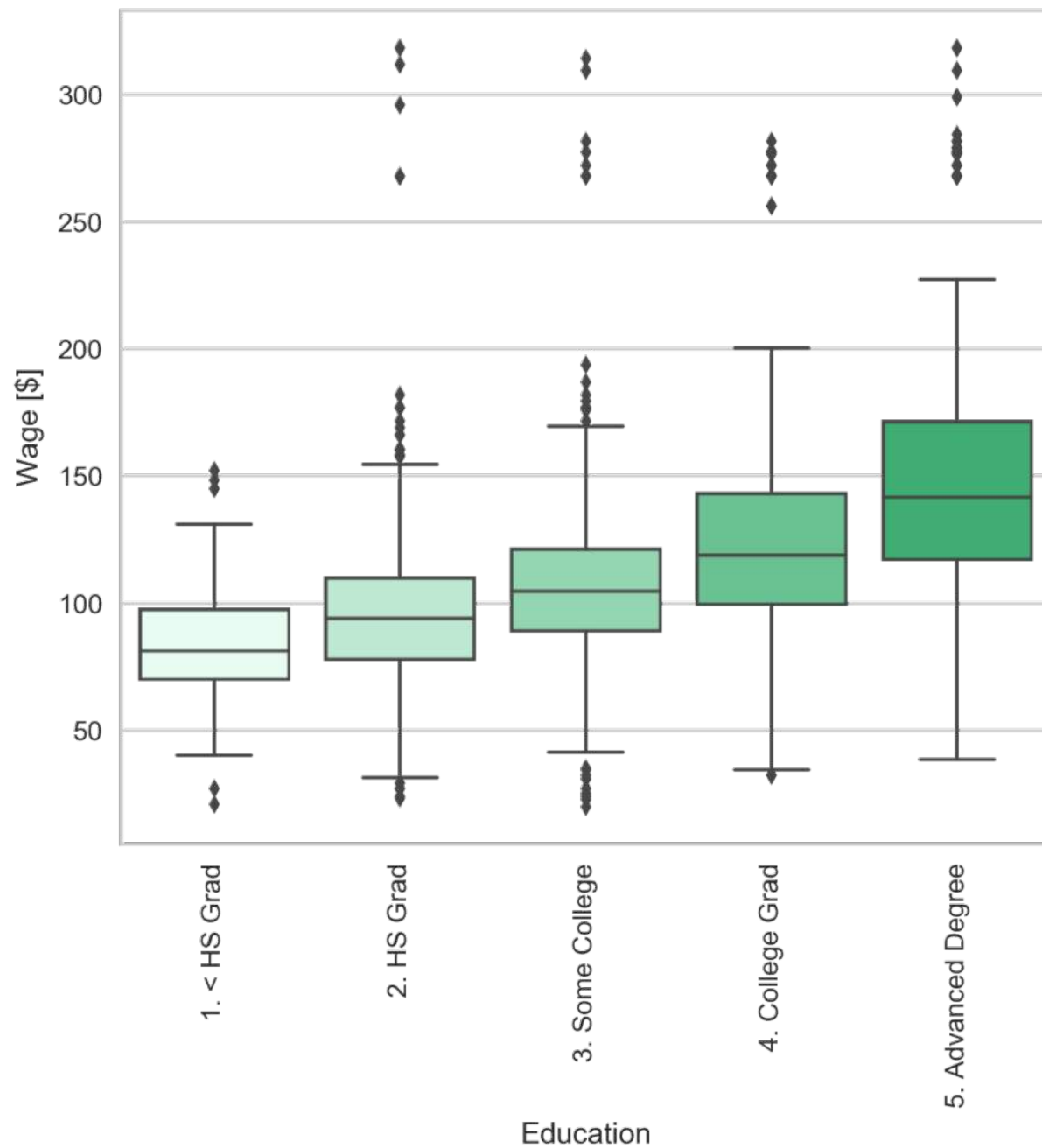


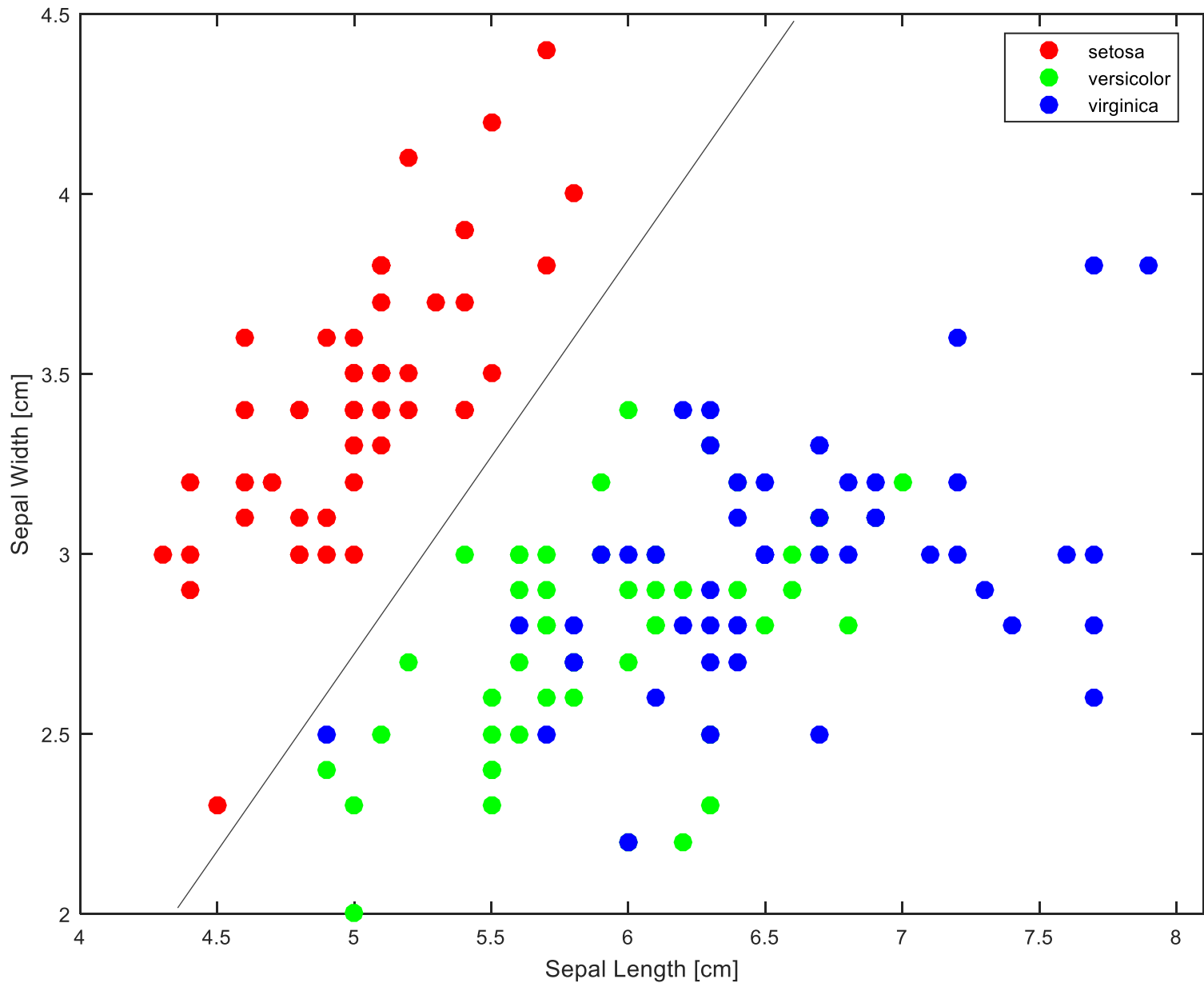
# What is machine learning?

Machine Learning Day



## Wage data from workers in the mid-Atlantic region

Data source: James et al., 2013



setosa



versicolor



virginica



# Challenges



What  
is  
this?

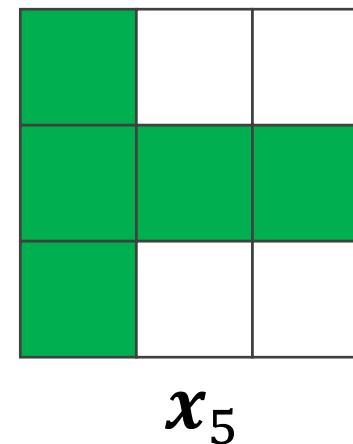
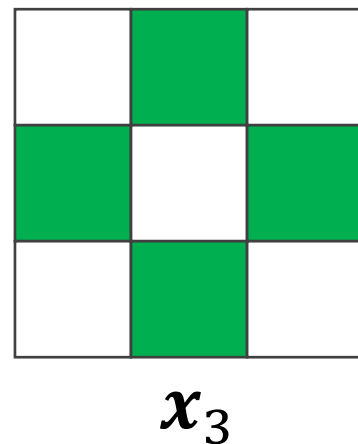
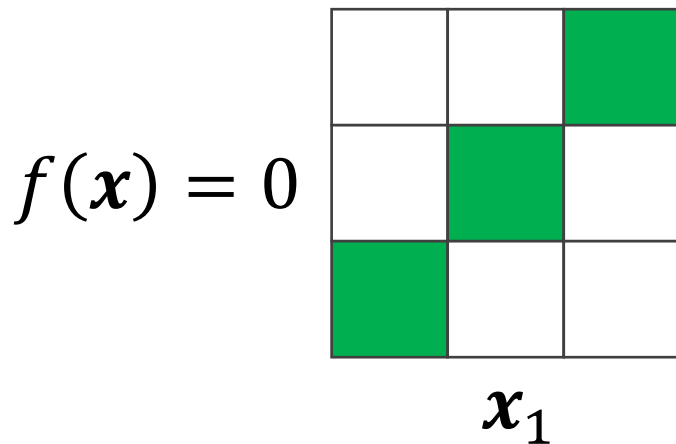
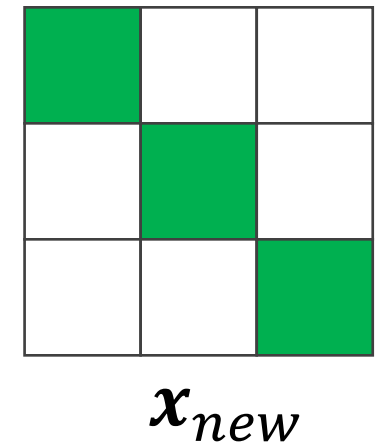
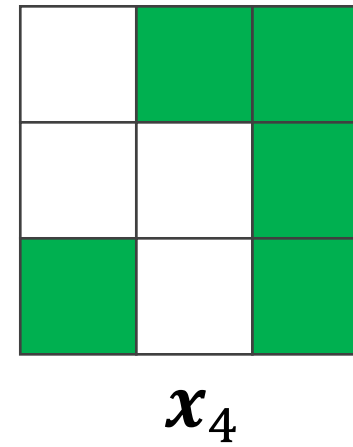
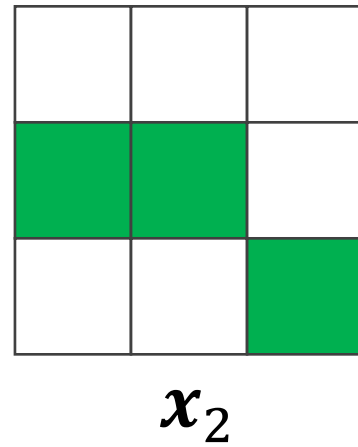
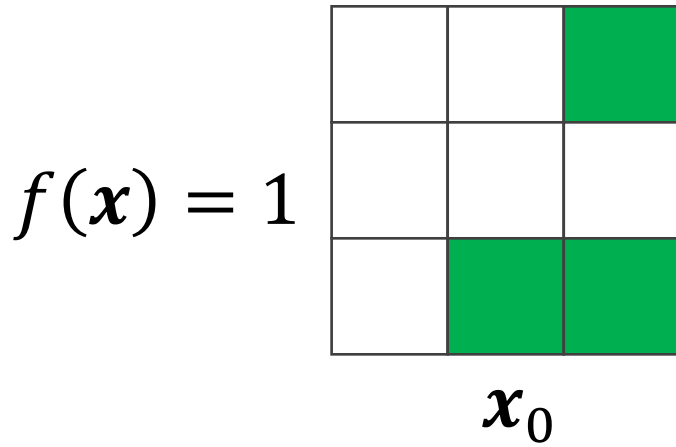


# We **generalize** from past experiences



Image: "It's not what it seems" by artist Hikaru Cho

# Predict which class $x_{\text{new}}$ belongs to...



$f(x_{\text{new}}) = ?$

# 1

# Machine learning is an **ill-posed problem**

There are often **many** models that will fit your data – how do we choose which to use?

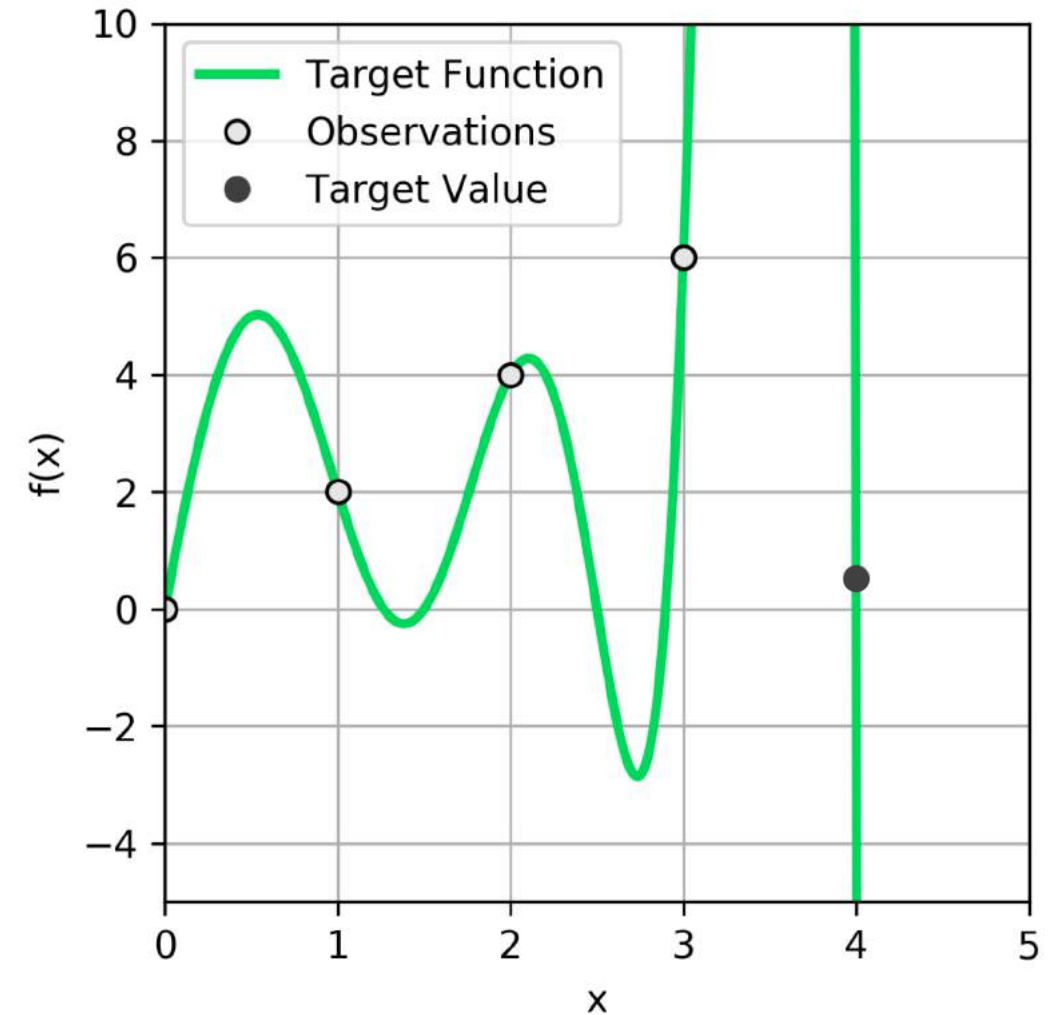


the best models  
**generalize** well

# Predict the next value in the sequence...

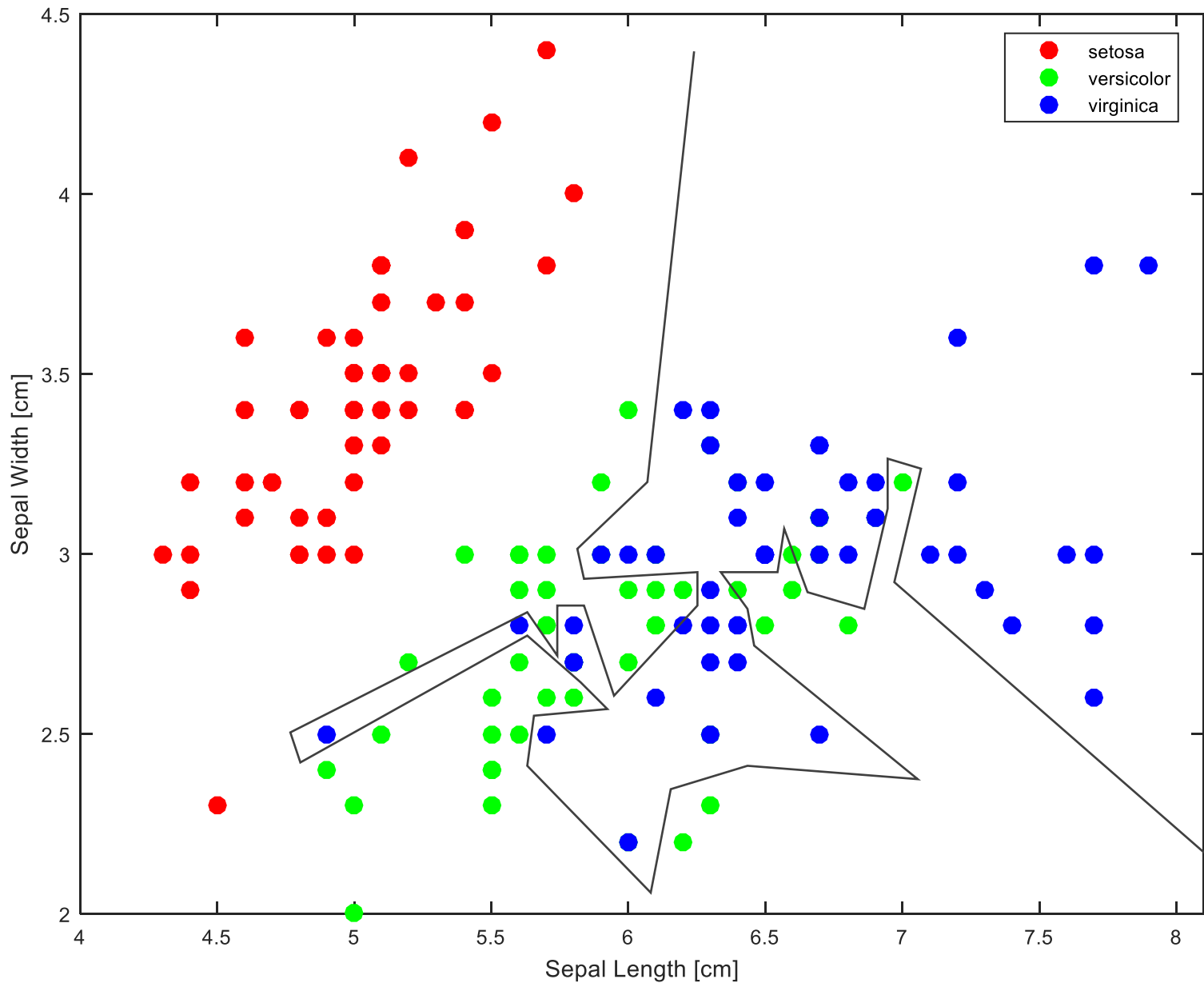
$x$	0	1	2	3	4
$f(x)$	0	2	4	6	?

$$f(4) = \boxed{0.530}$$



Our guess:

$$f(x) = 16.2x - 6.36x^2 - 11.9x^3 - 4.77x^4 + 7.03x^5 + 8.32x^6 - 9.01x^7 + 2.75x^8 - 0.275x^9$$



setosa



versicolor



virginica



# 2

## Overfit works against generalization

We can use **regularization** to prevent overfit

# What is machine learning?

A class of techniques where the **goal** is to **describe**, **predict**, and **strategize**...

...**based on** data, past experiences, and/or direct instruction...

...and do so **automatically**, with minimal human intervention.

# Types of machine learning tools

## Types of learning

**Unsupervised** learning

**Supervised** learning

**Reinforcement** learning

## Common use case

**Describe**

**Predict**

**Strategize**



# Types of machine learning

	Unsupervised Learning	Supervised Learning	Reinforcement Learning
Goal	<b>Describe</b> ...structure in data	<b>Predict</b> ...from examples	<b>Strategize</b> Learn through interaction
Data available	predictors, $x$	predictor and response pairs, $(x, y)$	actions and delayed responses (called rewards)
Examples	<ul style="list-style-type: none"><li>• Density estimation</li><li>• Clustering</li><li>• Dimensionality reduction</li><li>• Anomaly detection</li></ul>	<ul style="list-style-type: none"><li>• Classification</li><li>• Regression</li></ul>	<ul style="list-style-type: none"><li>• Model-free learning</li><li>• Model-based learning</li></ul>

# Sale Price Prediction

## Input Data:

Home characteristics  
(Numerical & Categorical)

## Target Data:

Price estimate (numerical)

## Learning Category:

Supervised Learning  
Regression



Kyle Bradbury

## 27708 Real Estate

1 home for sale

[Homes for You](#) [Newest](#) [Cheapest](#) [More](#)

**HOUSE FOR SALE**  
**\$599,900** 5 bds · 4 ba · 3,264 sqft  
1640 Marion Ave, Durham, NC

### 1640 Marion Ave, Durham, NC 27705

5 beds · 4 baths · 3,264 sqft

SPACIOUS RANCH W FINISHED LL WALKOUT! 5 BEDROOMS AND 4 BRAND NEW BATHS! RENOVATED WITH CUSTOM FEATURES THRUOUT! CONTEMPORARY HOME WITH MANY HANDICAP ACCESSIBLE REQUIREMENTS ALREADY IN PLACE! VAULTED CEILINGS! SECLUDED TREED LOT! GREAT HOME FOR LIVING AND ENTERTAINING WITH LARGE REAR DECK! WONDERFUL CONTEMPORARY FEEL THAT LIVES LARGE WITH EASY ACCESS TO DUKE UNIVERSITY: SHOPPING; HEALTH CARE; PARKS; R SHOPPING; AND EASY HIGHWAY AC

**FOR SALE**  
**\$599,900**  
Price cut: -\$79,100 (6/17)  
Zestimate®: \$619,585

EST. MORTGAGE  
\$2,284/mo   
[Get pre-qualified](#)

**Zestimate®: \$619,585**

# Video Recommendations



## Sherlock

97% Match 2017 TV-14 4 Series

97% Match



Season 3's episode "The Abominable Bride," which originally aired as a TV movie, won two Emmys.



MY LIST



### Input Data:

User video ratings  
(numerical and categorical)

### Target Data:

User rating of video  
(numerical)

### Learning Category:

Recommender Systems  
~Supervised & Unsupervised



OVERVIEW

EPISODES

MORE LIKE THIS

DETAILS

NETFLIX



# Spam Filters

**From:** Internal Revenue Service  
[mailto:yourtaxrefund@InternalRevenueService.com]

**Sent:** Tuesday, July 22, 2008 9:47 AM

**Subject:** Get your tax refund now

**Importance:** High

After the last annual calculations of your account activity we have determined that you are eligible to receive a tax refund of \$479.30 .

Please submit the tax refund request and allow us 2-6 days in order to process it.

A refund can be delayed for a variety of reasons. For example submitting invalid records or applying after the deadline.

To access the form for your tax refund, please click here (<http://e-dlogs.rta.mi.th:84/www.irs.gov/>)

Note: Deliberate wrong inputs will be prosecuted by law.

Regards,

Internal Revenue Service

**Input Data:**

Email text (text)

**Target Data :**

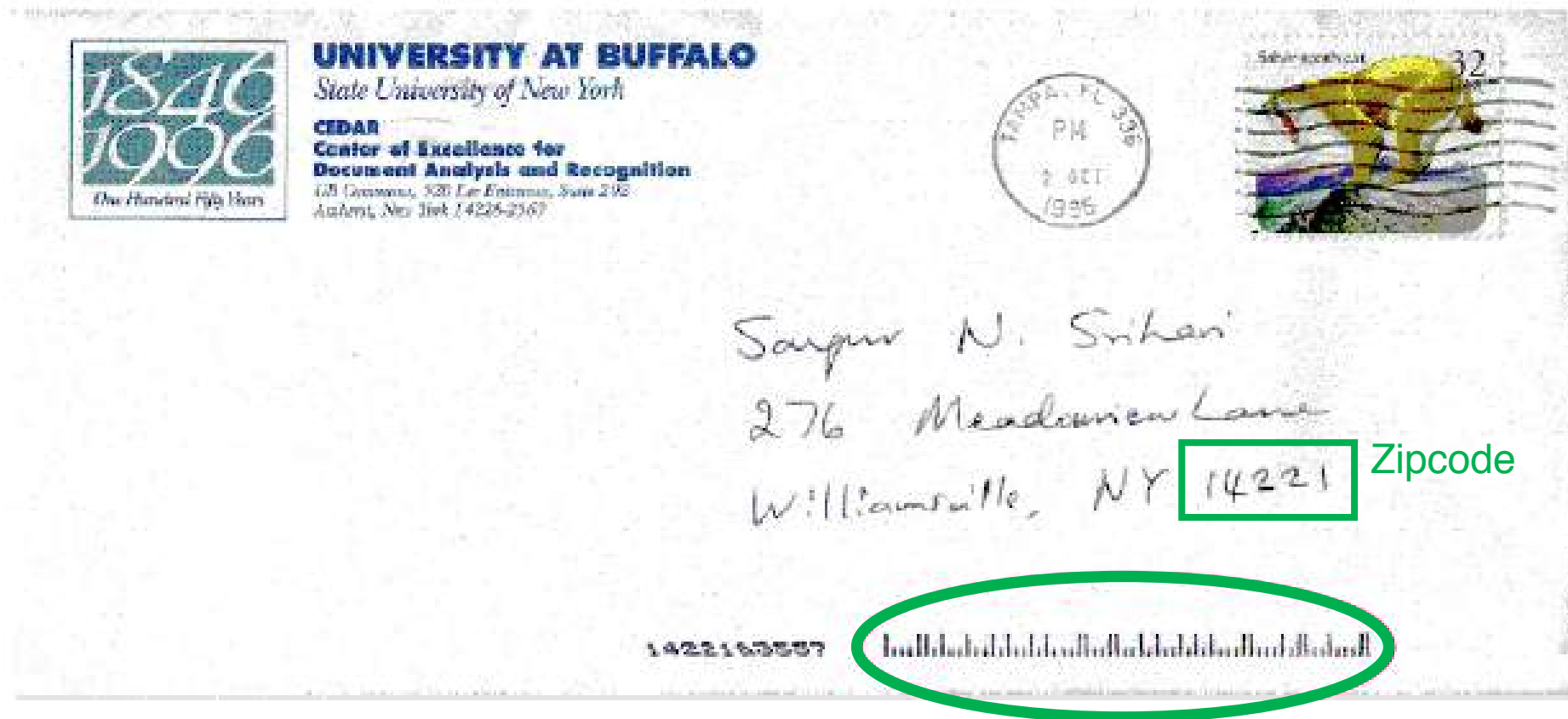
Spam/not spam  
(category)

**Learning Category:**

Supervised Learning  
Classification (binary)

Spam example source: itservices.uchicago.edu

# Handwriting and Optical Character Recognition



**Input Data:**

Imagery

**Target Data:**

Text Characters

**Learning Category:**

Supervised Learning

Classification (multiclass)

Among the first handwritten addresses sorted automatically in October 1996

Image source: Sargur Srihari, SUNY



# Where's Waldo = Computer Vision Problem



**Input Data:**  
Color Imagery (Image)

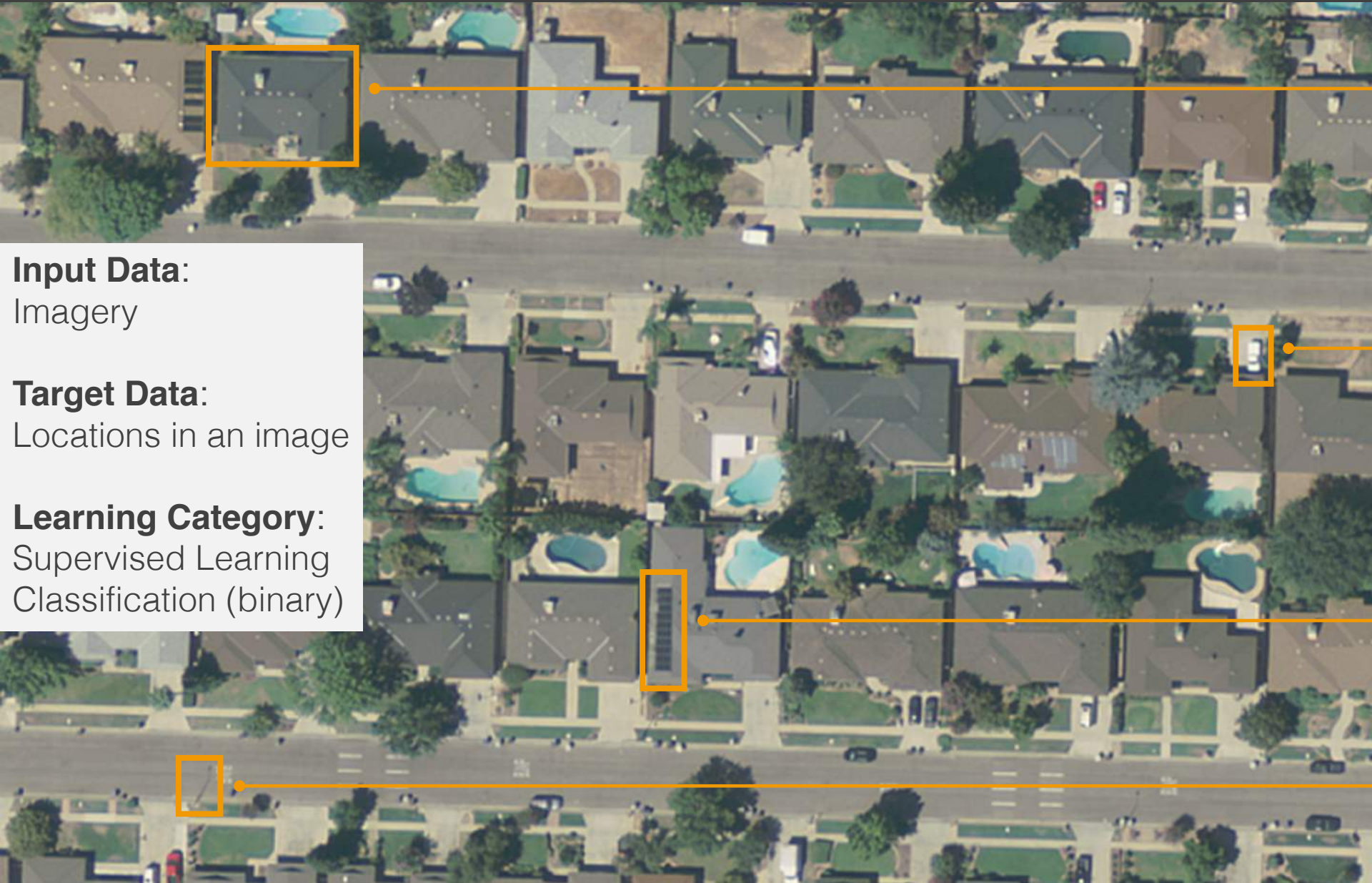
**Target Data:**  
Locations in an image  
(label for each pixel)

**Learning Category:**  
Supervised Learning  
Classification (binary)

Image source: [www.whereswaldo.com/](http://www.whereswaldo.com/)



# Object Recognition: Energy Systems



**Input Data:**

Imagery

**Target Data:**

Locations in an image

**Learning Category:**

Supervised Learning  
Classification (binary)

**Building**

behind-the-meter  
energy consumption

**Car**

transportation  
energy consumption

**Solar Array**

distributed energy  
resources

**Light Pole**

access to electricity

# Types of machine learning

	Unsupervised Learning	Supervised Learning	Reinforcement Learning
Goal	<b>Describe</b> ...structure in data	<b>Predict</b> ...from examples	<b>Strategize</b> Learn through interaction
Data available	predictors, $x$	predictor and response pairs, $(x, y)$	actions and delayed responses (called rewards)
Examples	<ul style="list-style-type: none"><li>• Density estimation</li><li>• Clustering</li><li>• Dimensionality reduction</li><li>• Anomaly detection</li></ul>	<ul style="list-style-type: none"><li>• Classification</li><li>• Regression</li></ul>	<ul style="list-style-type: none"><li>• Model-free learning</li><li>• Model-based learning</li></ul>

# Credit Fraud

## Input Data:

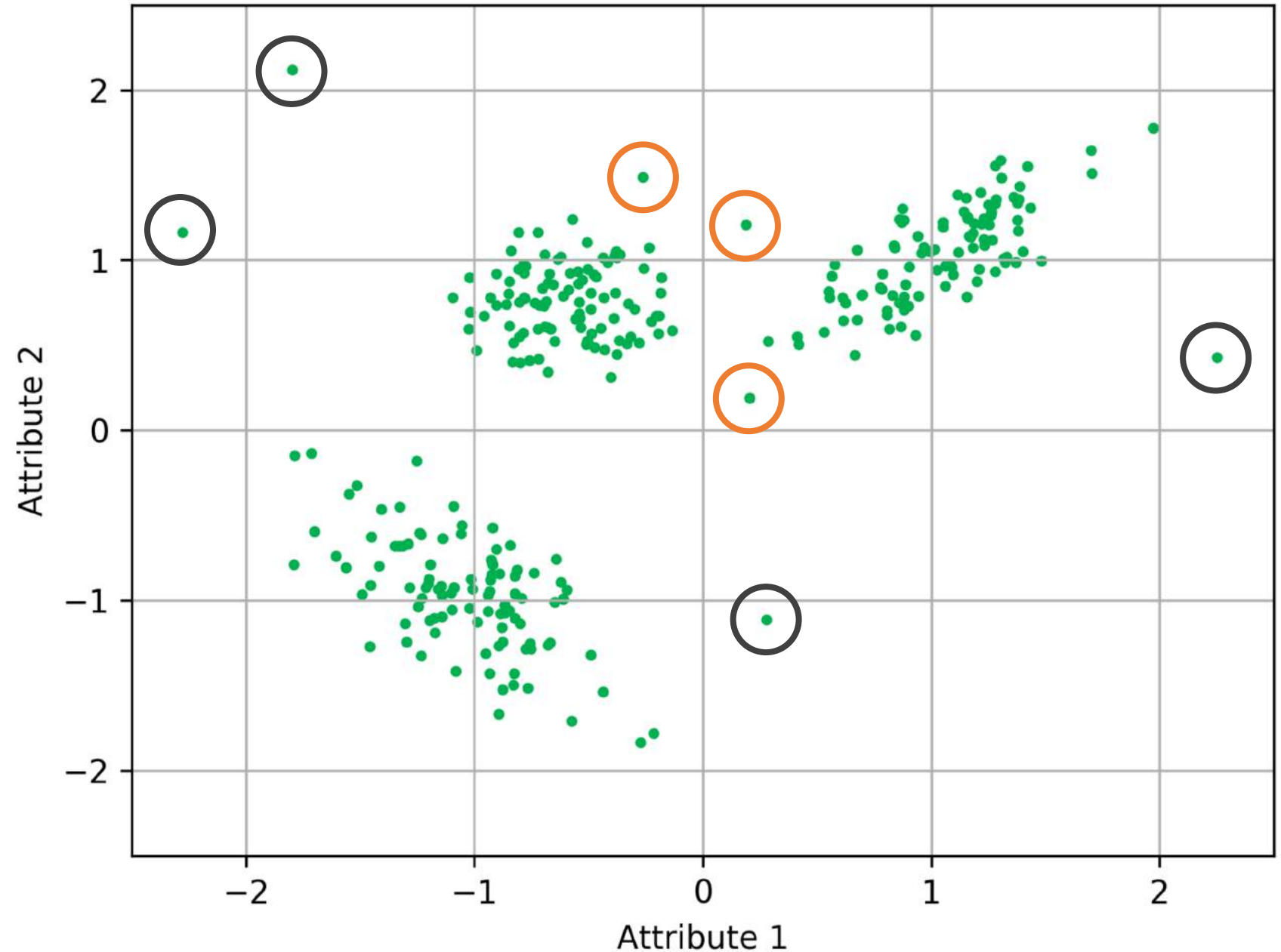
Account transactions, dates,  
locations, demographic  
information  
(Numerical and categorical)

## Target Data:

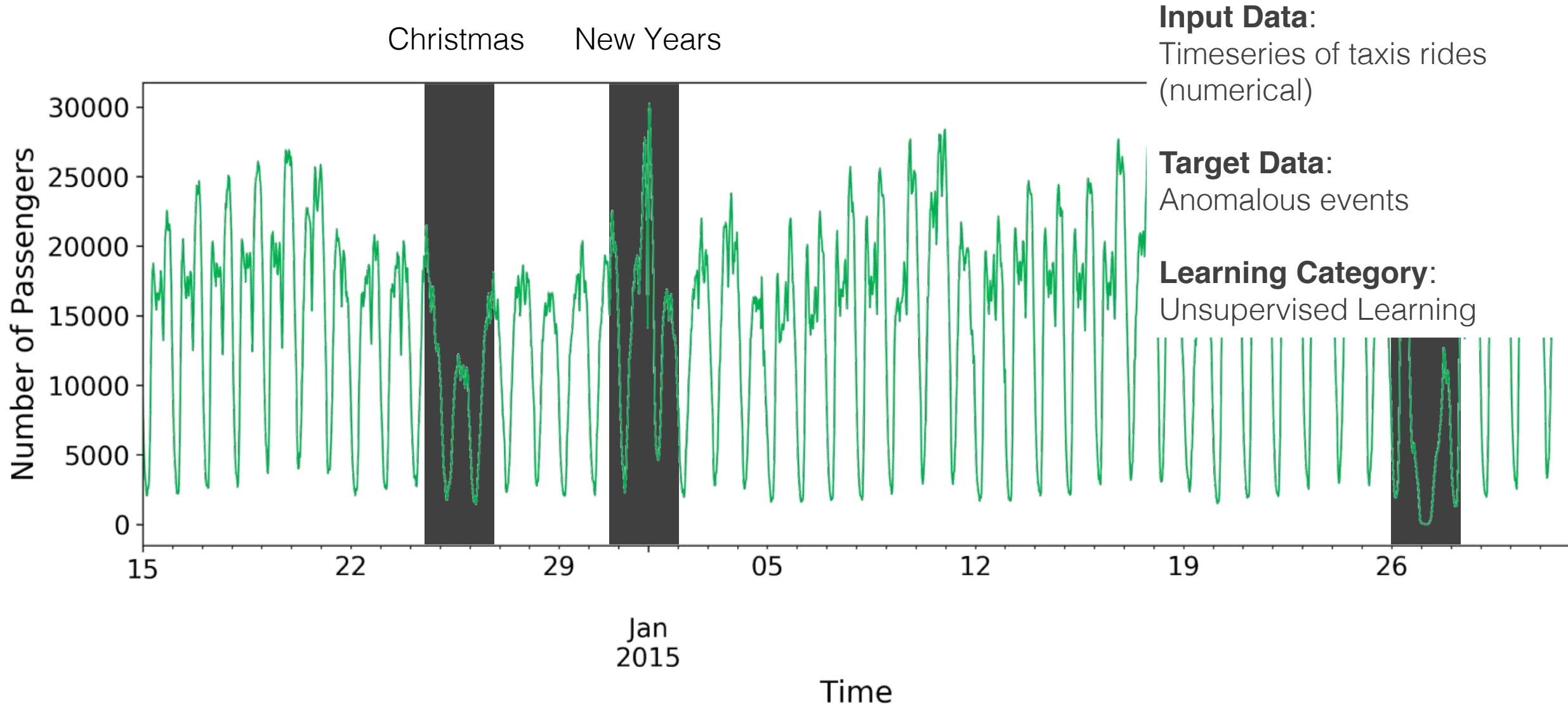
Anomalous transactions

## Learning Category:

Unsupervised Learning  
Clustering, Density  
Estimation



# Anomalous Event Detection: NYC Taxis



Data source: Numenta Anomaly Benchmark (NAB), from kaggle.com



# Types of machine learning

	Unsupervised Learning	Supervised Learning	Reinforcement Learning
Goal	<b>Describe</b> ...structure in data	<b>Predict</b> ...from examples	<b>Strategize</b> Learn through interaction
Data available	predictors, $x$	predictor and response pairs, $(x, y)$	actions and delayed responses (called rewards)
Examples	<ul style="list-style-type: none"><li>• Density estimation</li><li>• Clustering</li><li>• Dimensionality reduction</li><li>• Anomaly detection</li></ul>	<ul style="list-style-type: none"><li>• Classification</li><li>• Regression</li></ul>	<ul style="list-style-type: none"><li>• Model-free learning</li><li>• Model-based learning</li></ul>

# Learning a strategy to master games

## Input Data:

Moves taken and occasional feedback on win/loss  
(Numerical and categorical)

## Target Data:

Win/loss (Maximizing rewards)

## Learning Category:

Reinforcement Learning



## THE ULTIMATE GO CHALLENGE

GAME 3 OF 3

27 MAY 2017



**RESULT B + Res**

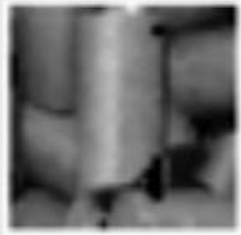


Google DeepMind



# Manufacturing – learn to pick up iron cylinders

Success Failure



## Input Data:

Actions taken and occasional feedback on success/failure (Numerical and categorical)

## Target Data:

Success/failure (Maximizing rewards)

## Learning Category:

Reinforcement Learning



Source: MIT Technology Review; Company: **FANUC**

# Types of machine learning

	Unsupervised Learning	Supervised Learning	Reinforcement Learning
Goal	<b>Describe</b> ...structure in data	<b>Predict</b> ...from examples	<b>Strategize</b> Learn through interaction
Data available	predictors, $x$	predictor and response pairs, $(x, y)$	actions and delayed responses (called rewards)
Examples	<ul style="list-style-type: none"><li>• Density estimation</li><li>• Clustering</li><li>• Dimensionality reduction</li><li>• Anomaly detection</li></ul>	<ul style="list-style-type: none"><li>• Classification</li><li>• Regression</li></ul>	<ul style="list-style-type: none"><li>• Model-free learning</li><li>• Model-based learning</li></ul>

THIS IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG  
PILE OF LINEAR ALGEBRA, THEN COLLECT  
THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

JUST STIR THE PILE UNTIL  
THEY START LOOKING RIGHT.

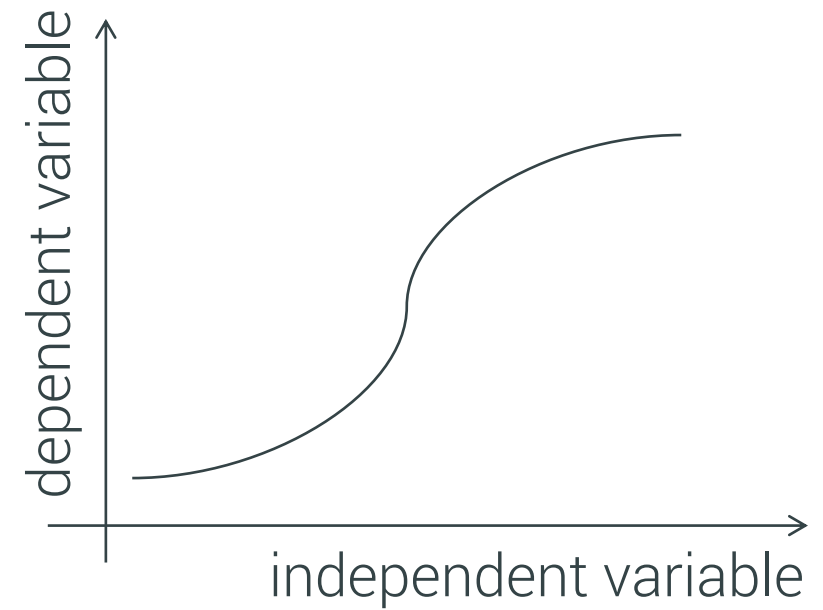


Image: xkcd.com

# A Taste of Supervised Learning

Classification and Regression

# Common language



**independent** variable

input

predictor

feature

$x$

**dependent** variable

output

response

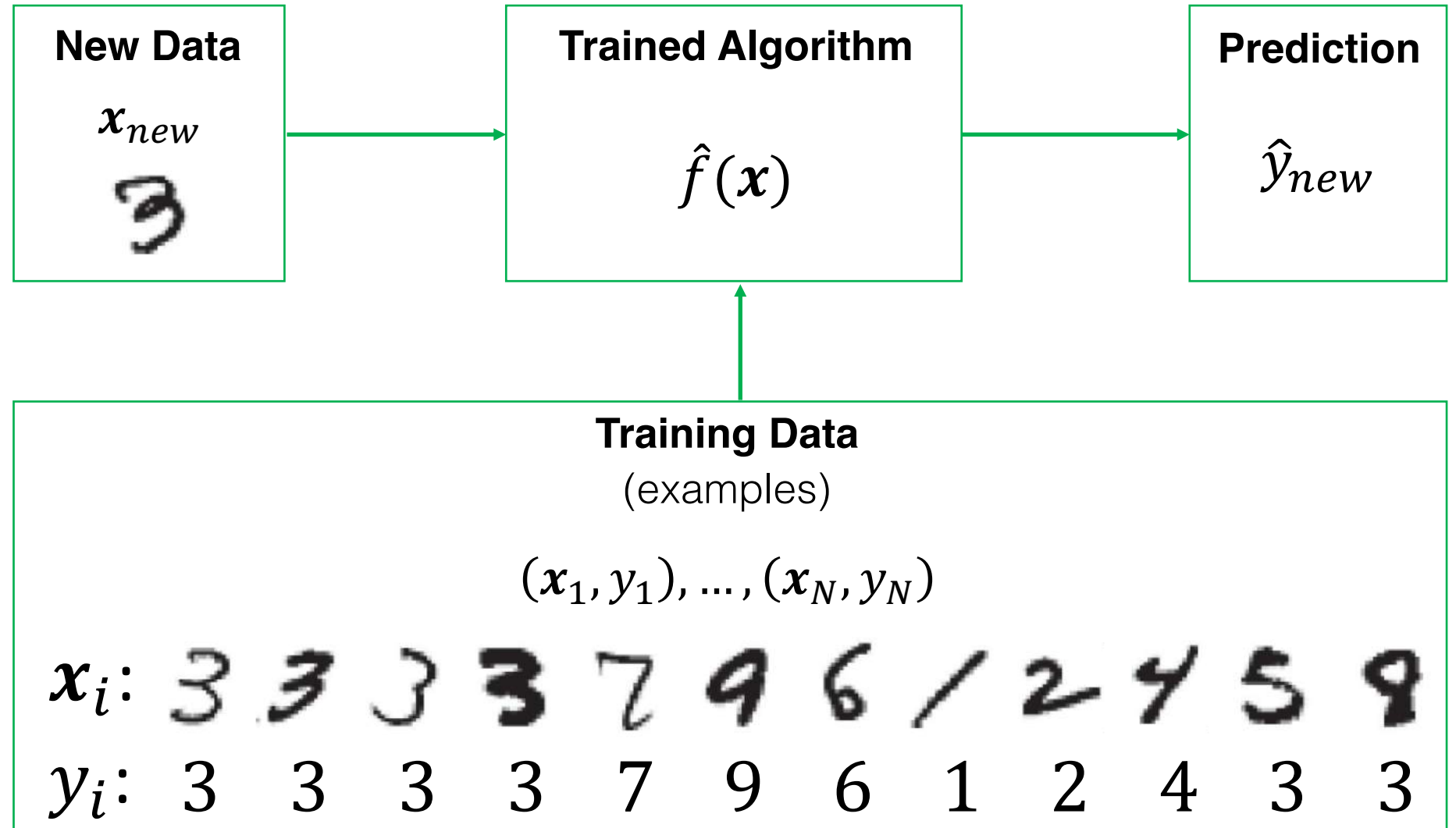
target

$y$

# Supervised learning

Objective: create an algorithm that predicts well

Example:  
**Digits classification**

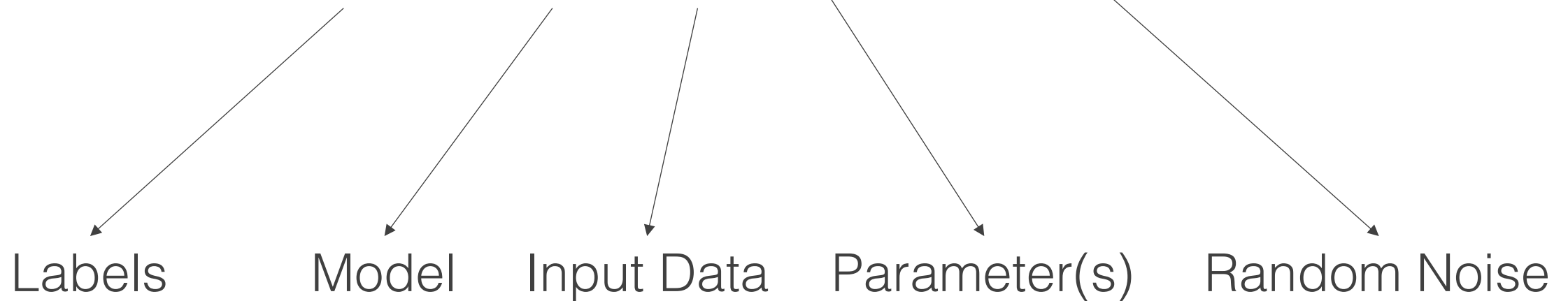




# Supervised machine learning model

We search for  
the model that  
best fits our  
data

$$y = f(x, w) + \epsilon$$



Hypothesis

Typically prevents  
perfect performance

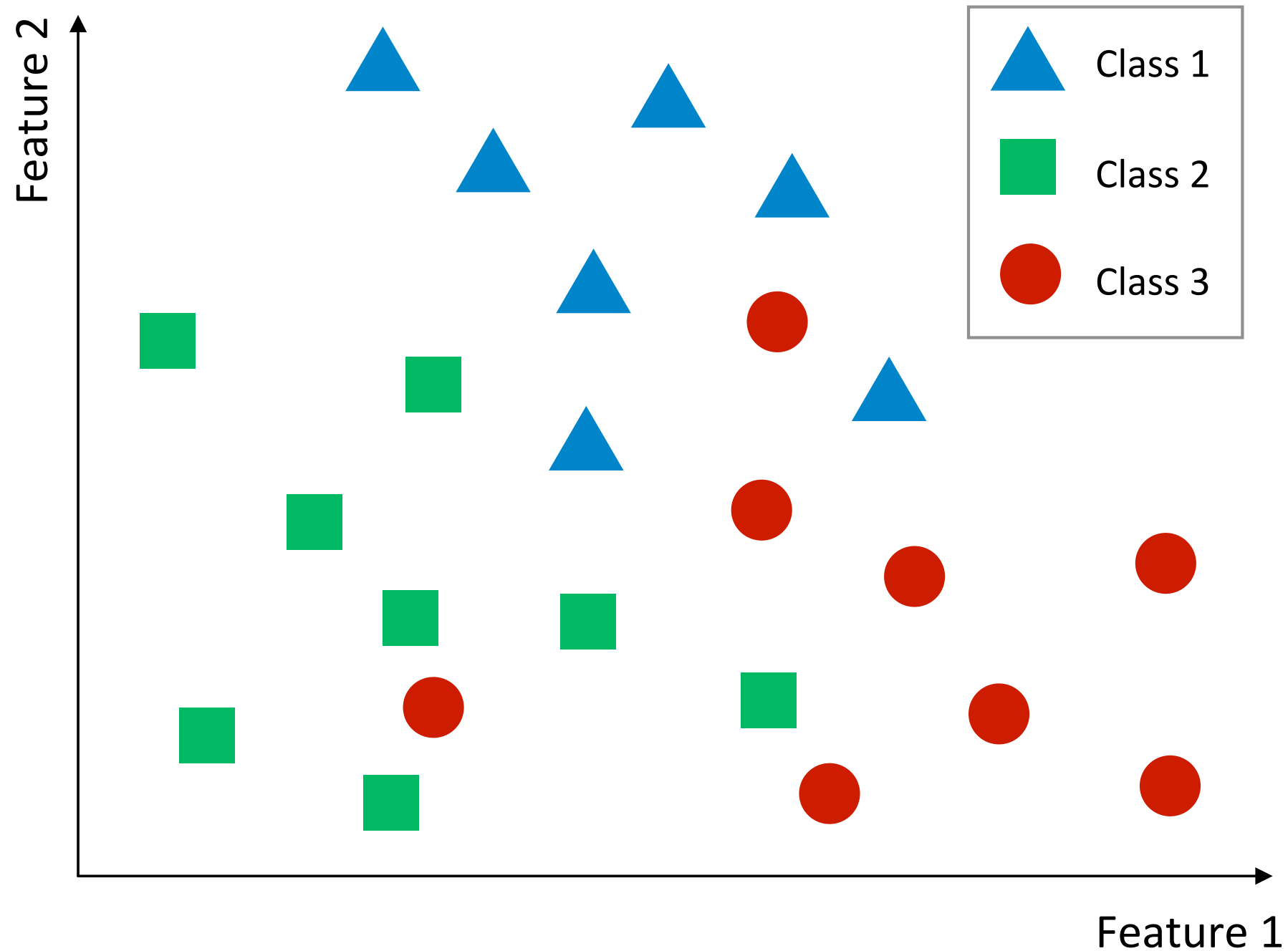
# K-Nearest Neighbors

Classification and Regression

# K Nearest Neighbor Classifier

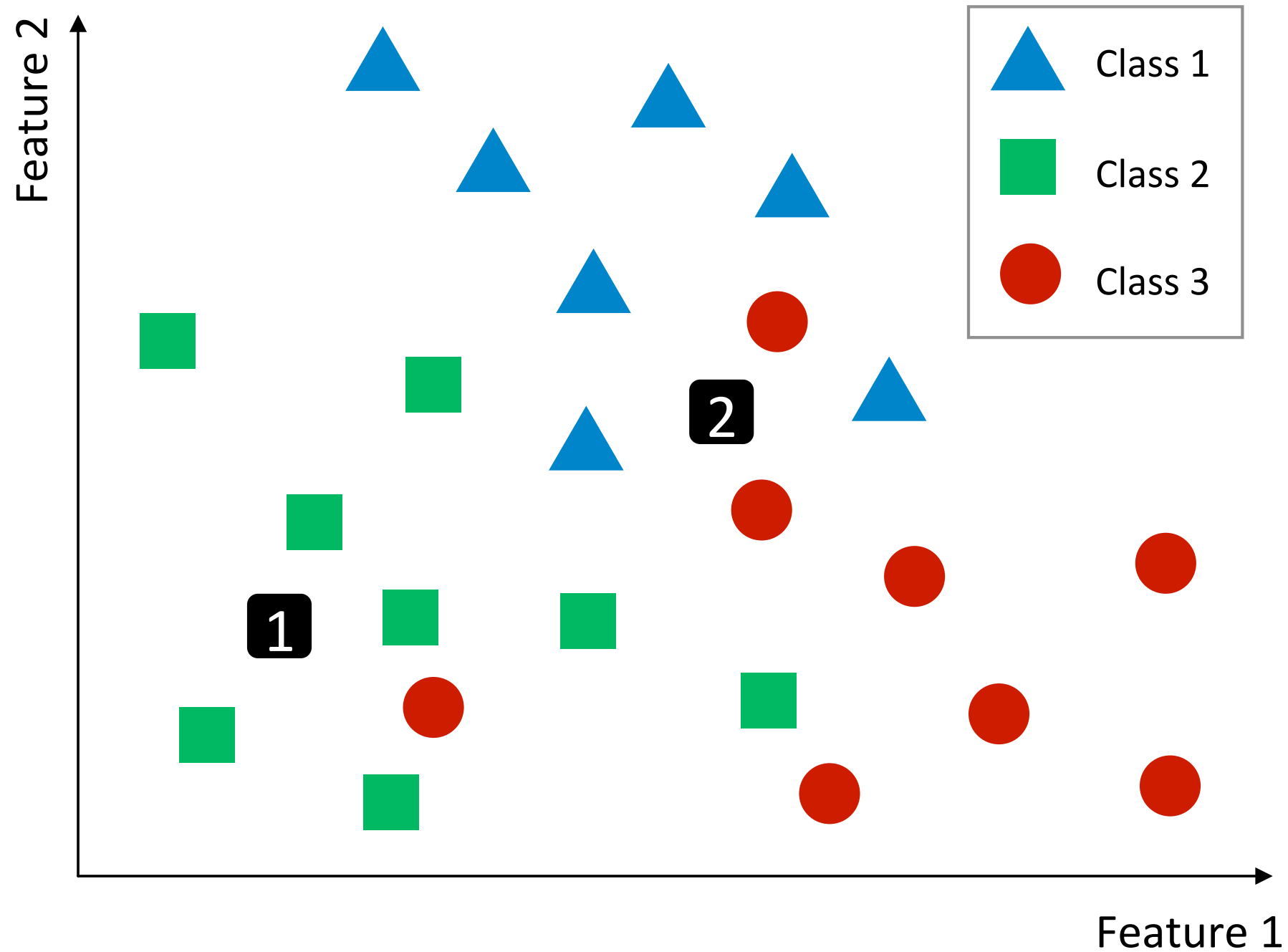
## Step 1: Training

Every new data point  
is a model parameter



# K Nearest Neighbor Classifier

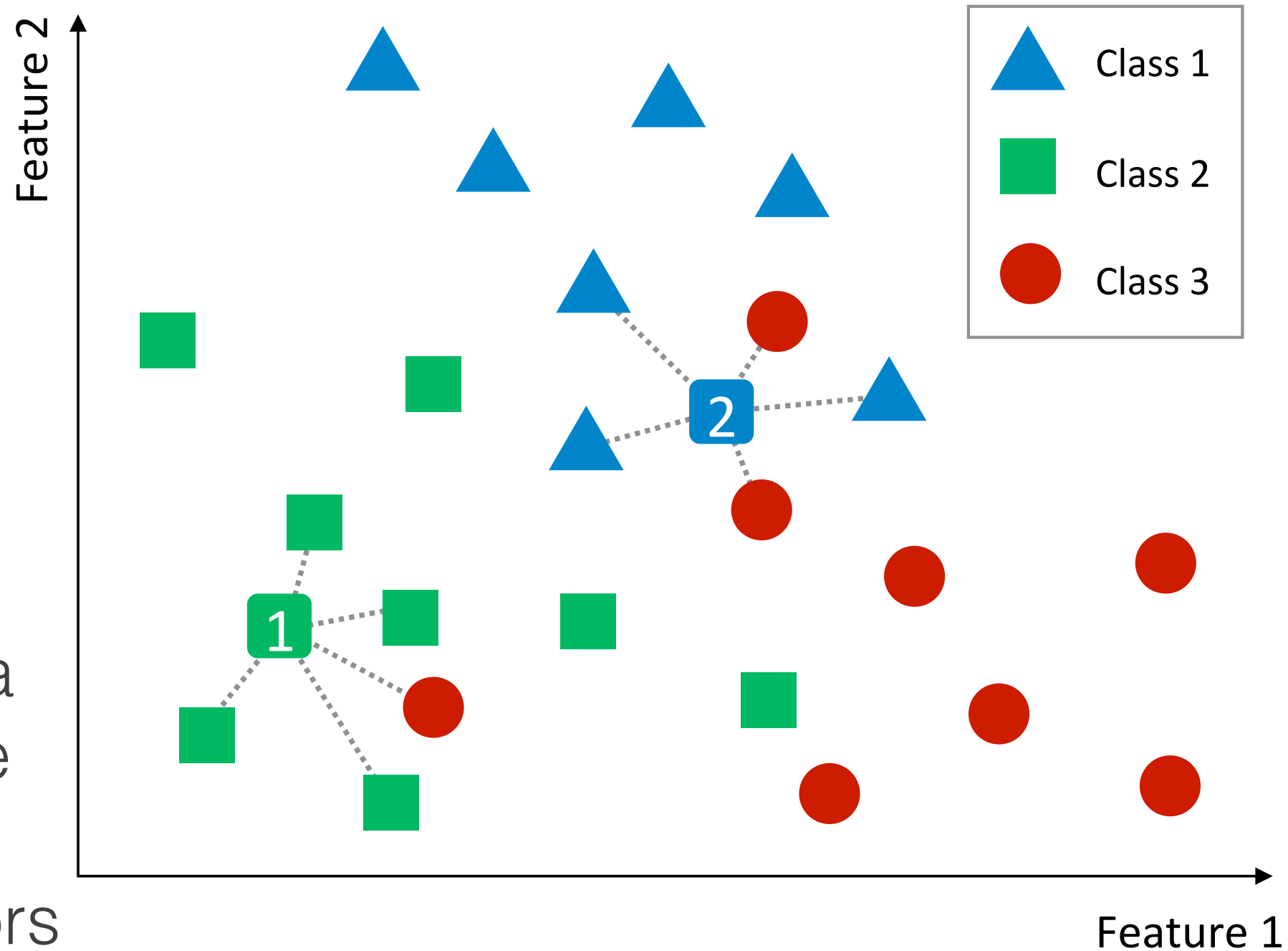
**Step 2:**  
Place new  
(unseen)  
examples in the  
feature space



# K Nearest Neighbor Classifier

## Step 3:

Classify the data by assigning the class of the  $k$  nearest neighbors



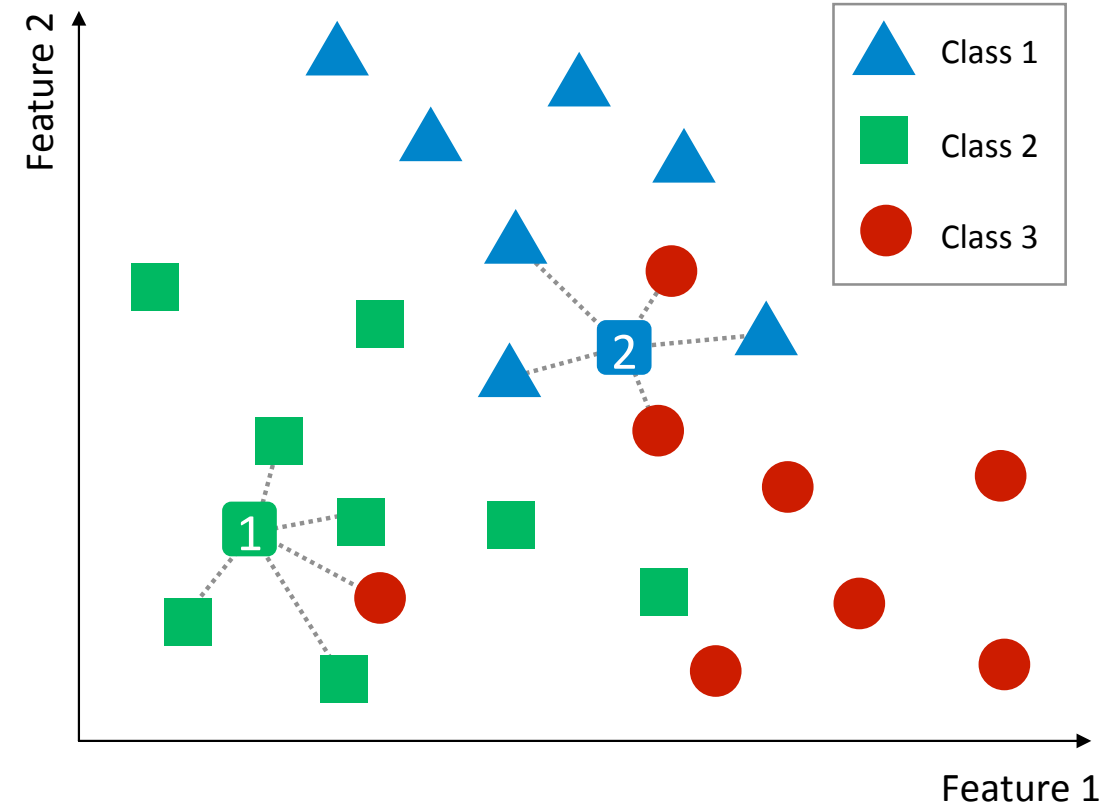
# K Nearest Neighbor Classifier

## Score vs Decision :

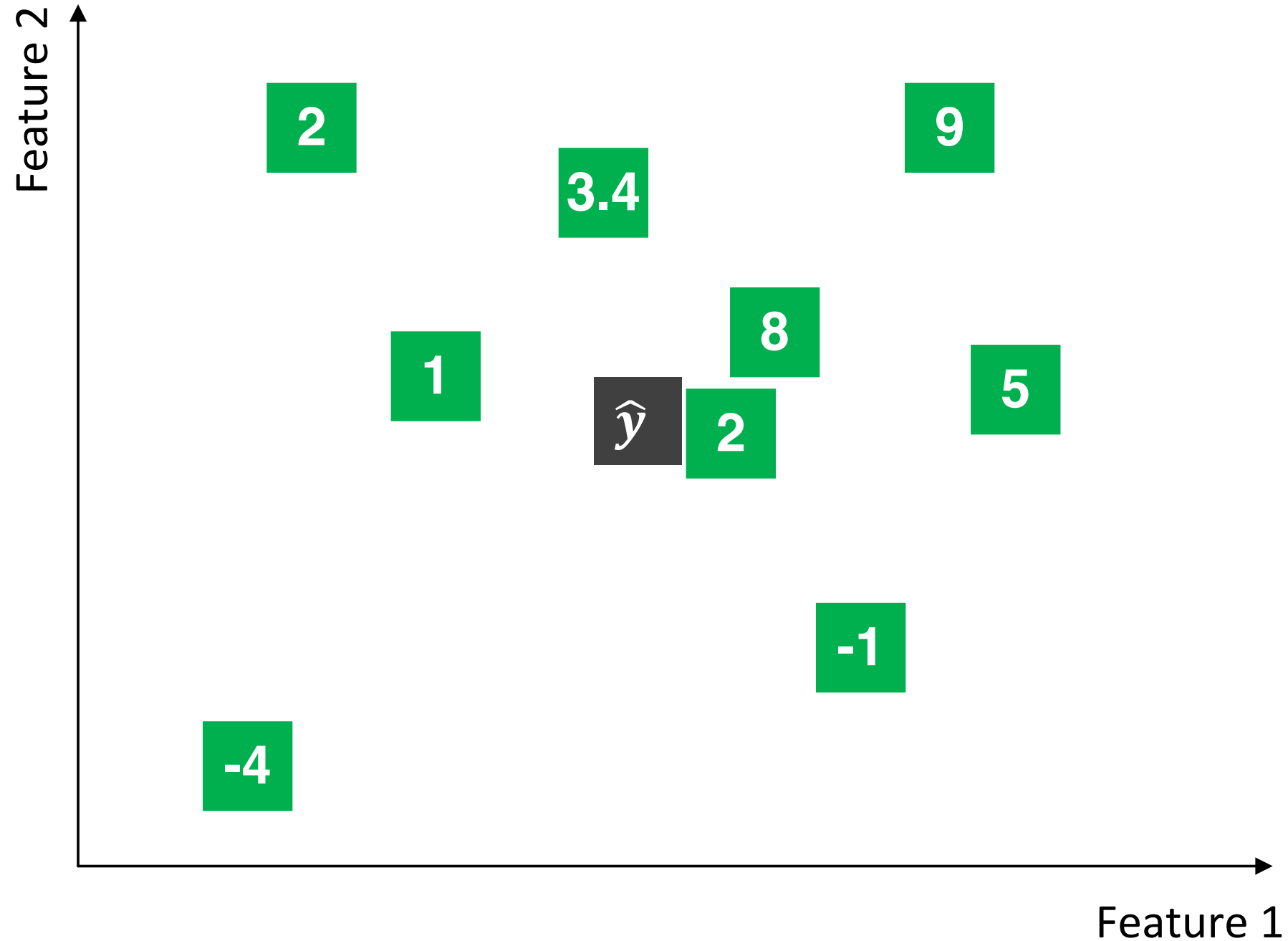
For 5-NN, the confidence score that a sample belongs to a class could be:  $\{0, 1/5, 2/5, 3/5, 4/5, 1\}$

## Decision Rule:

If the confidence score for a class  $>$  threshold, predict that class

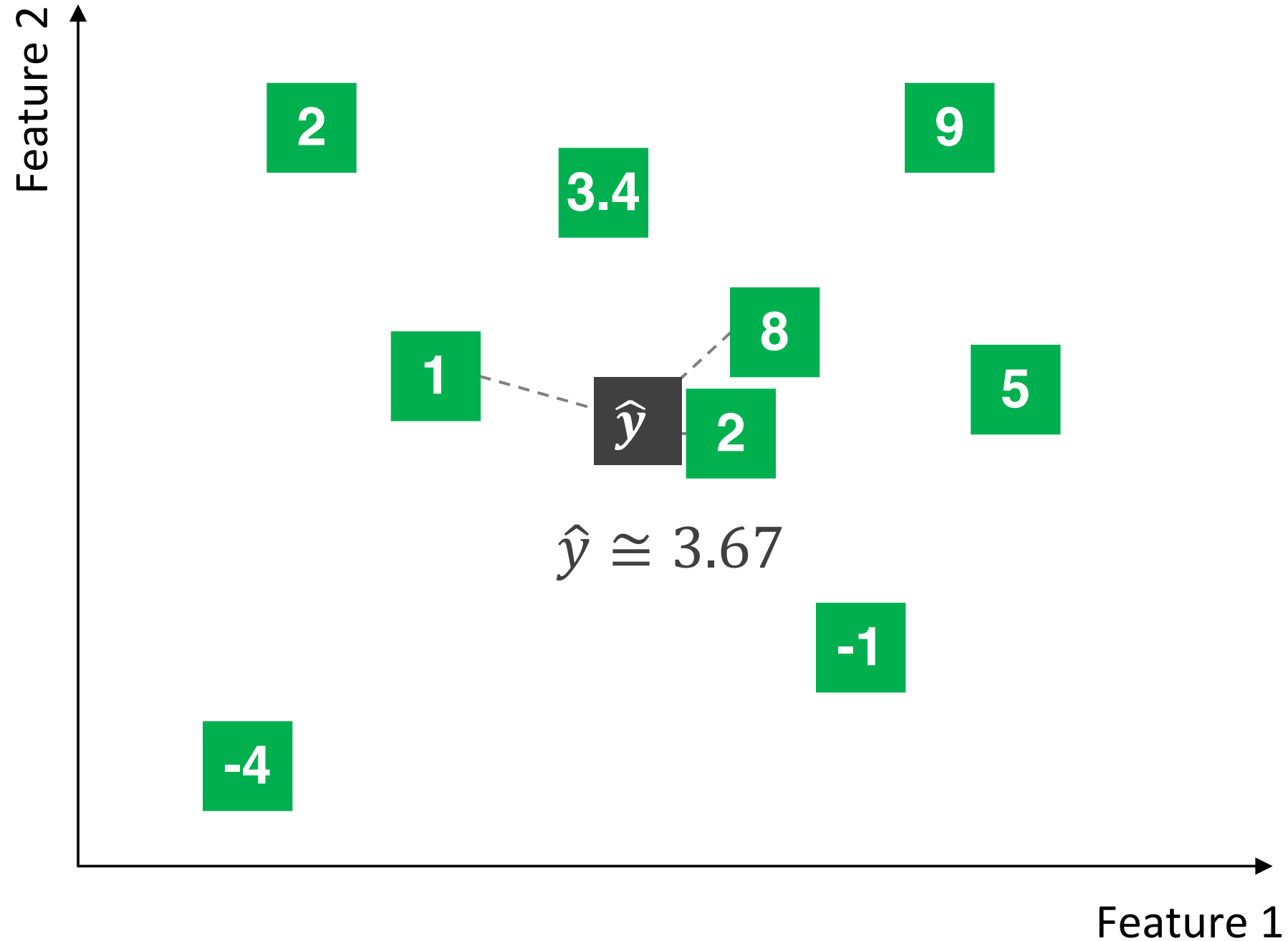


# K Nearest Neighbor Regression



# K Nearest Neighbor Regression

$$\hat{y} = \frac{1}{k} \sum_{y_i \in \{\text{k nearest}\}} y_i$$





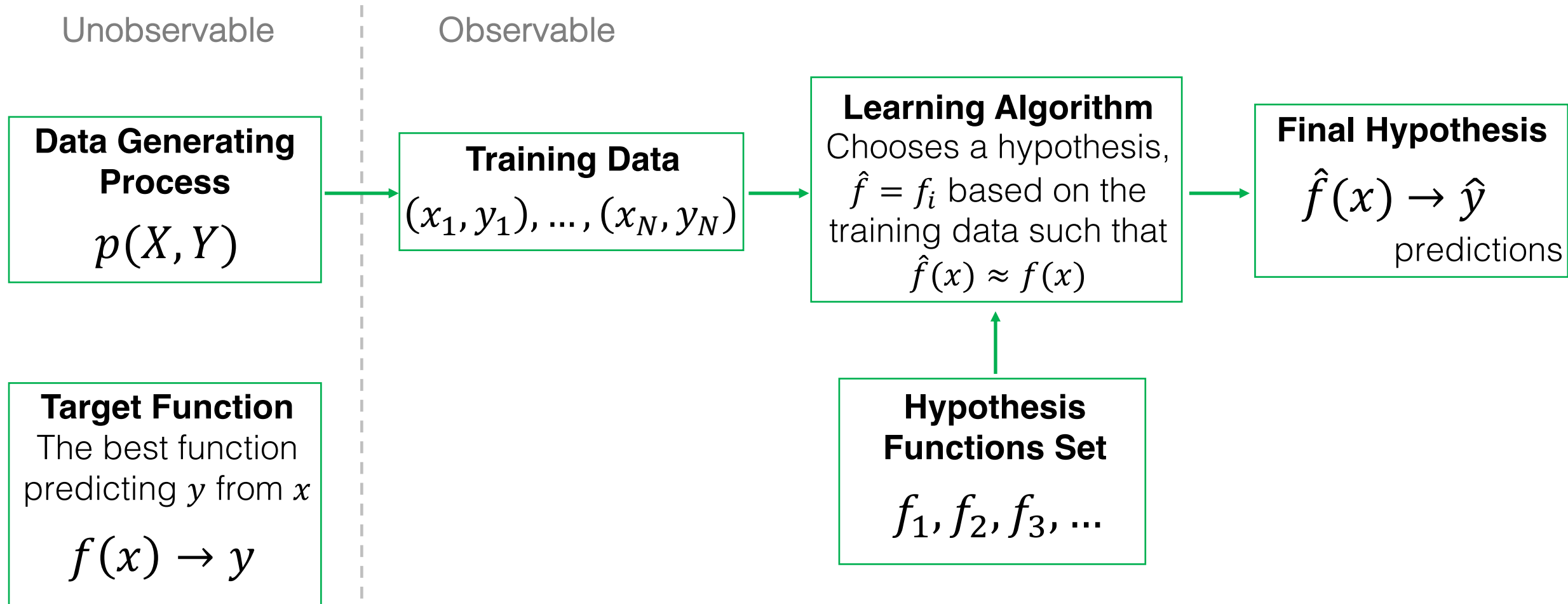
# Summarizing supervised learning

Let's review and bring it all together

# Components of supervised learning

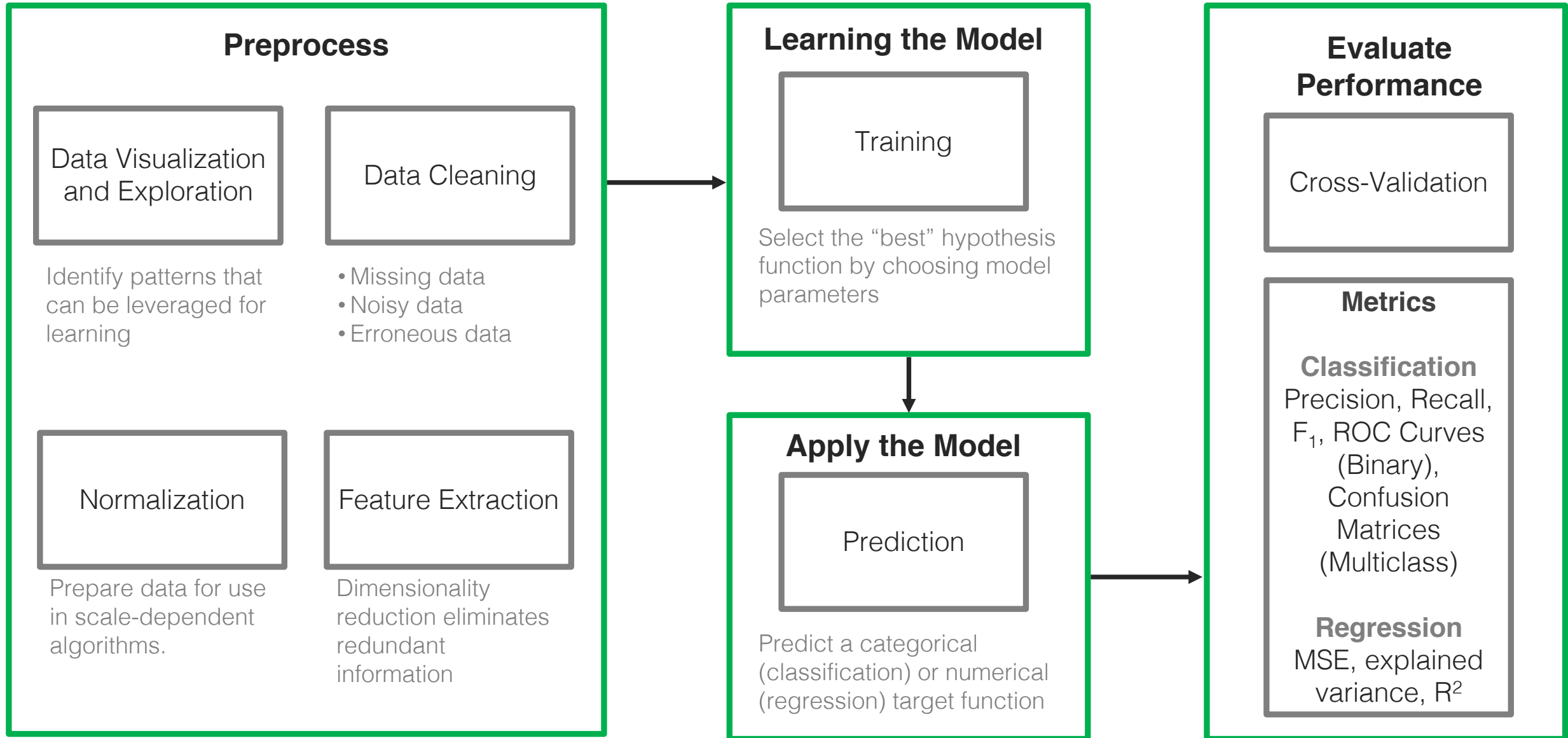
<b>Input</b>	$\mathbf{x}$	
<b>Output</b>	$y$	
<b>Training Data</b>	$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)$	
<b>Target function</b>	$f(\mathbf{x}) \rightarrow y$	This is unknown, but the best you could ever do
<b>Hypothesis set</b>	$f_i(\mathbf{x}) \rightarrow \hat{y}$	Functions to consider in trying to approximate $f(\mathbf{x})$
<b>Learning algorithm</b>	Optimization technique that searches the hypothesis set for the function $f_i$ that best approximates $f$ (typically by choosing parameters in a model)	

# Supervised Learning



- Need to select the hypothesis functions (models to train)
- Need to select the learning algorithm (for fitting the models to the data)

# Supervised learning in practice



# Want to learn more?

List of additional resources:

<http://www.kylebradbury.org/datascience.html>

# ENERGY

data analytics lab

