

# **INTRO to DATA SCIENCE**

## **LECTURE 4: INTRO TO ML & KNN CLASSIFICATION**

Francesco Mosconi  
DAT10 SF // October 15, 2014

---

**HEADER— CLASS NAME, PRESENTATION TITLE**

---

# **DATA SCIENCE IN THE NEWS**

---

## DATA SCIENCE IN THE NEWS

---

### BUILDING A RACE SIMULATOR

October 3, 2014 · by f1metrics · in Mathematical models, Predictions · 25 Comments



Source: <http://f1metrics.wordpress.com/2014/10/03/building-a-race-simulator/>

---

## DATA SCIENCE IN THE NEWS

---

# Deep Learning on Amazon EC2 GPU with Python and nolearn

by Adrian Rosebrock on October 13, 2014 in Deep Learning, Tutorials



Tweet

81



Like

9



+1

6



Share

15



Source: <http://www.pyimagesearch.com/2014/10/13/deep-learning-amazon-ec2-gpu-python-nolearn/>

---

## RECAP

---

## LAST TIME

- Cleaning data
- Dealing with missing data
- Setting up github for homework

---

**INTRO TO DATA SCIENCE**

---

**QUESTIONS?**

---

## **AGENDA**

---

**I. WHAT IS MACHINE LEARNING?**

**II. CLASSIFICATION PROBLEMS**

**III. BUILDING EFFECTIVE CLASSIFIERS**

**IV. THE KNN CLASSIFICATION MODEL**

**EXERCISES:**

**IV. LAB: KNN CLASSIFICATION IN PYTHON**

**V. BONUS LAB: VISUALIZATION WITH MATPLOTLIB (IF TIME ALLOWS)**

# **I. WHAT IS MACHINE LEARNING?**



"A field of study that gives computers the ability to learn without being explicitly programmed." (1959)



Arthur Samuel, AI pioneer  
Source: Stanford

"A computer program is said to learn from experience  $E$  with respect to some set of tasks  $T$  and performance measure  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ ". (1989)



Tom Mitchell, Professor, CMU  
(Source: CMU)

"A computer program is said to learn from experience  $E$  with respect to some set of tasks  $T$  and performance measure  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ ".

A person is said to learn from a college course  $E$  with respect to some set of readings and midterms  $T$  and grades  $P$ , if its performance at tasks in  $T$ , as measured by  $P$ , improves with  $E$ .

from Wikipedia:

“Machine learning, a branch of artificial intelligence, is about the construction and study of systems that can learn from data.”

source: [http://en.wikipedia.org/wiki/Machine\\_learning](http://en.wikipedia.org/wiki/Machine_learning)

from Wikipedia:

“Machine learning, a branch of artificial intelligence, is about the construction and study of systems that can learn from data.”

“The core of machine learning deals with representation and generalization...”

source: [http://en.wikipedia.org/wiki/Machine\\_learning](http://en.wikipedia.org/wiki/Machine_learning)

from Wikipedia:

“Machine learning, a branch of artificial intelligence, is about the construction and study of systems that can learn from data.”

“The core of machine learning deals with representation and generalization...”

- representation – extracting structure from data

source: [http://en.wikipedia.org/wiki/Machine\\_learning](http://en.wikipedia.org/wiki/Machine_learning)

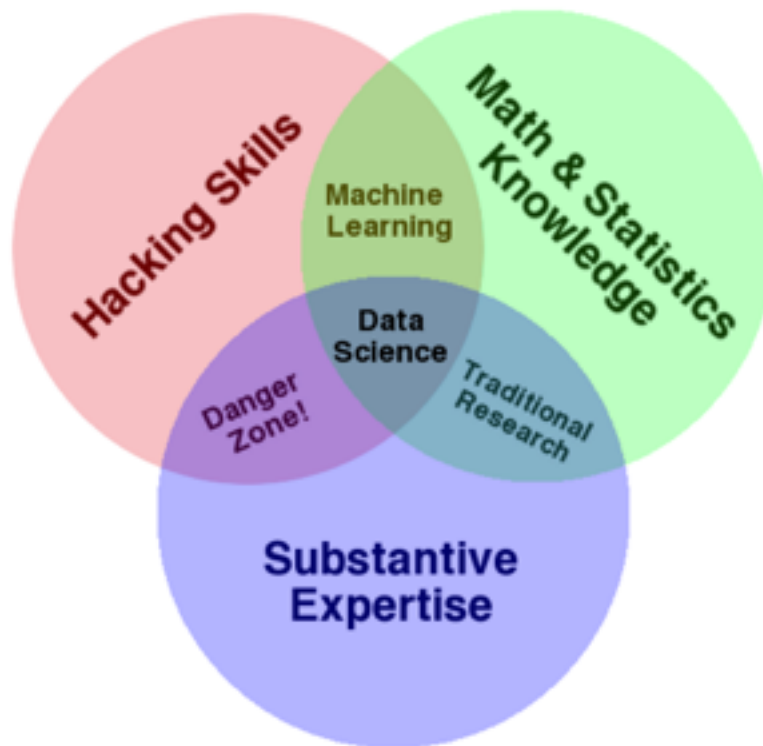
from Wikipedia:

“Machine learning, a branch of artificial intelligence, is about the construction and study of systems that can learn from data.”

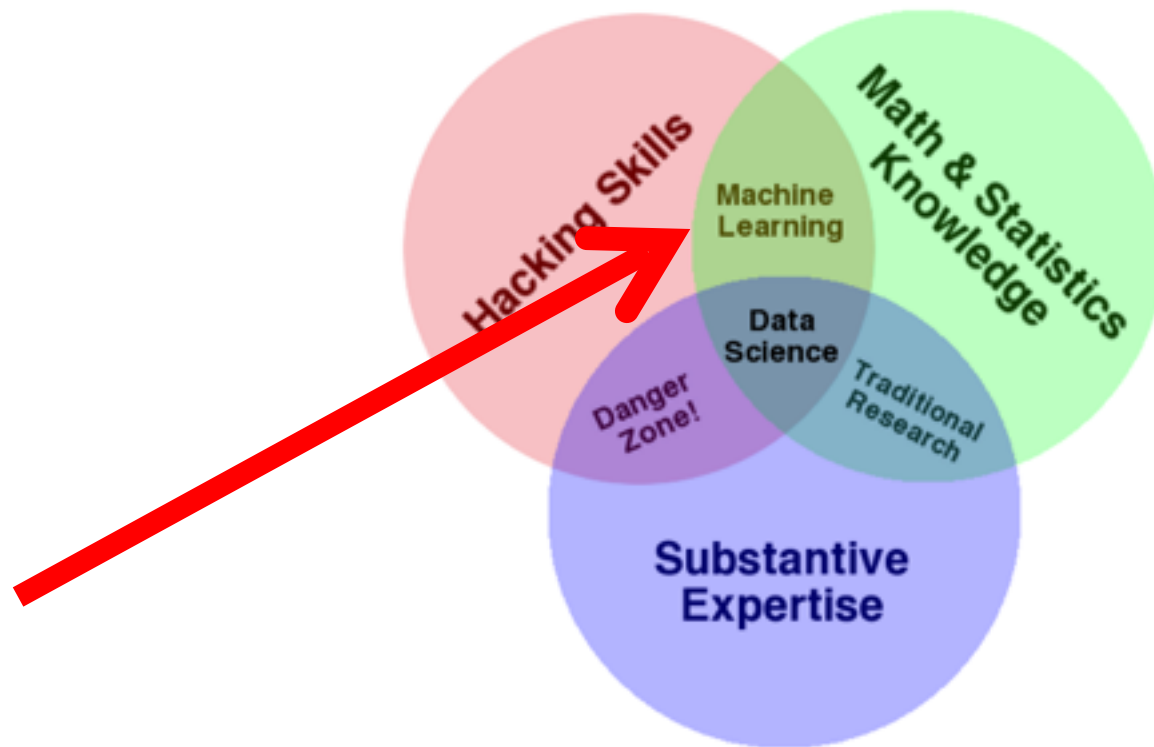
“The core of machine learning deals with representation and generalization...”

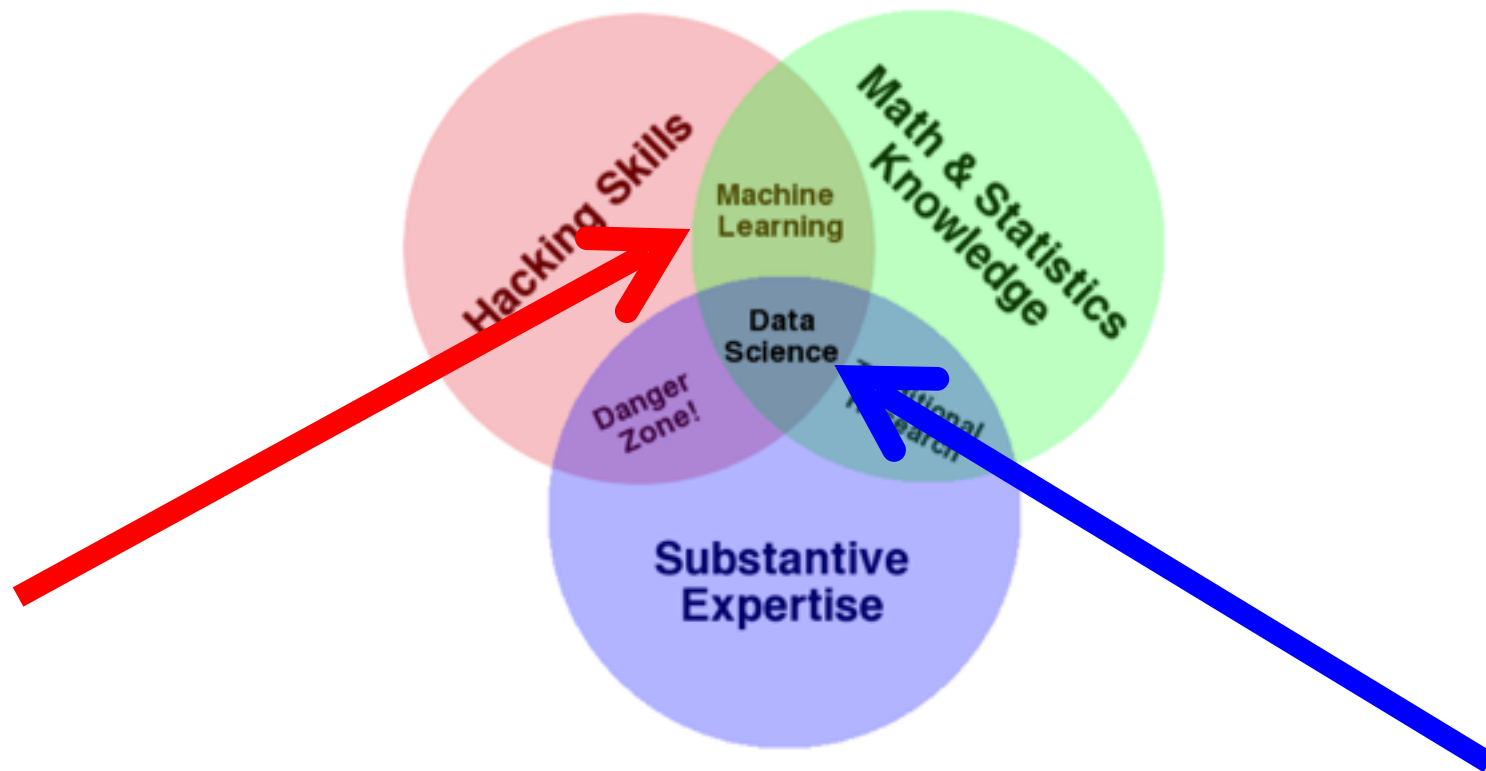
- representation – extracting structure from data
- generalization – making predictions from data

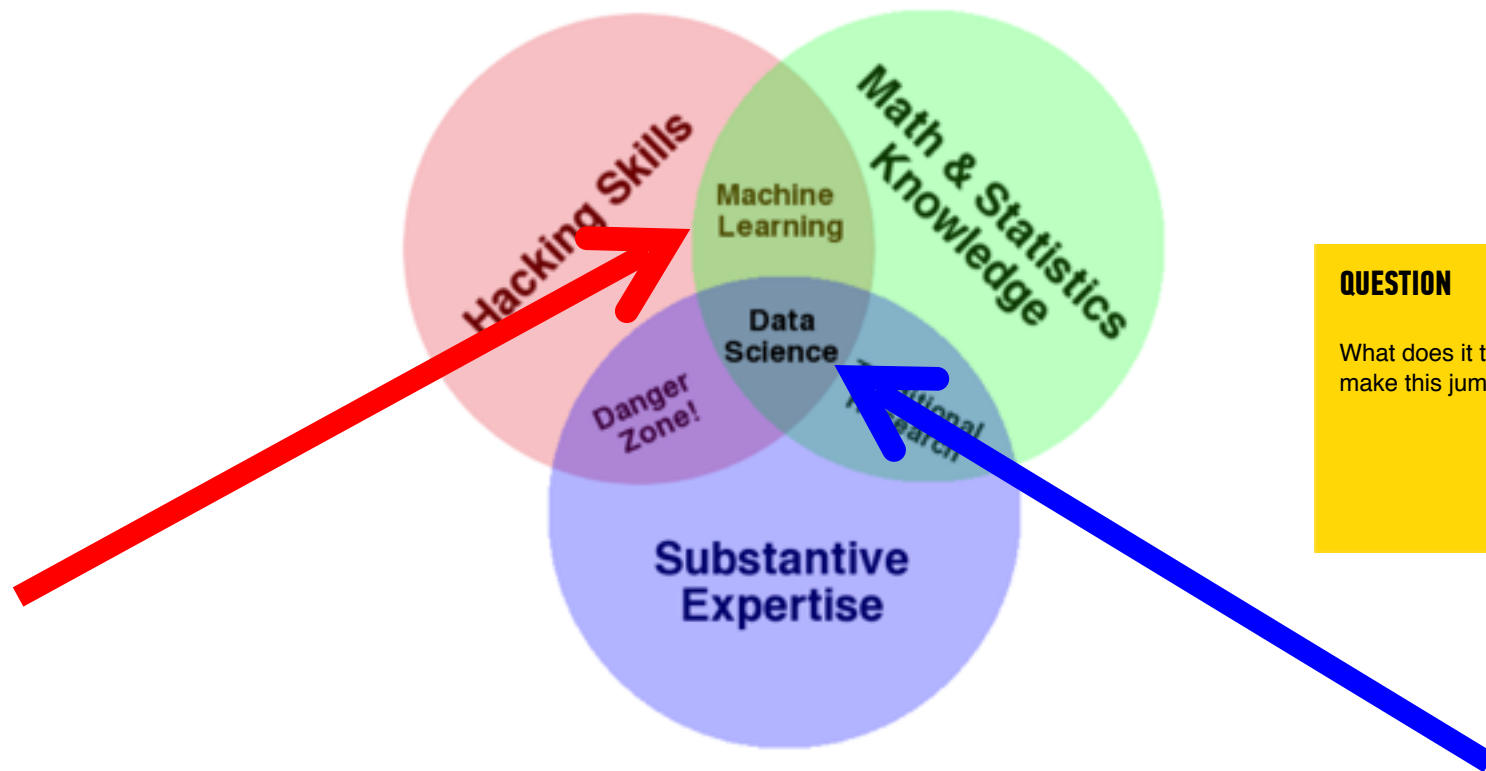
source: [http://en.wikipedia.org/wiki/Machine\\_learning](http://en.wikipedia.org/wiki/Machine_learning)











### QUESTION

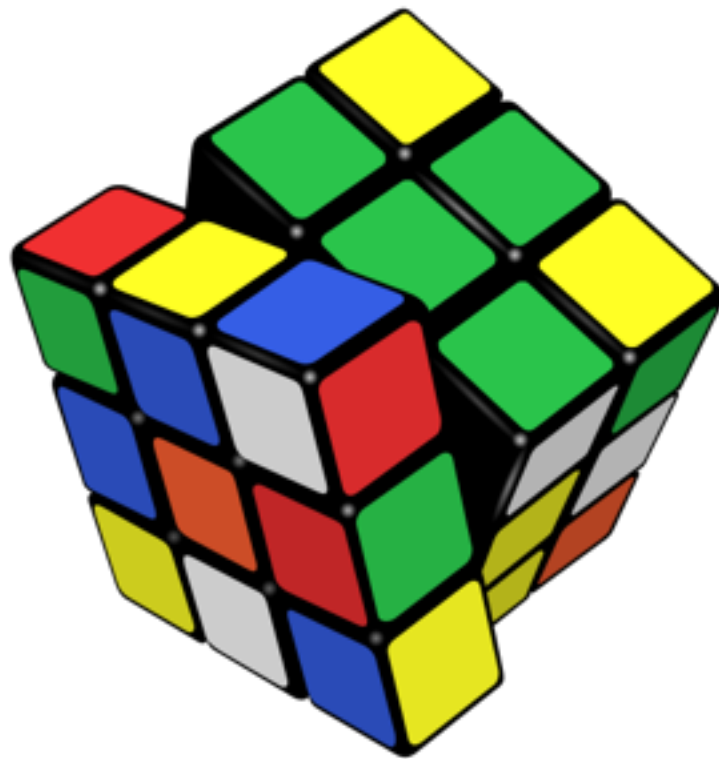
What does it take to make this jump?

**ANSWER: PROBLEM SOLVING!**

---

20



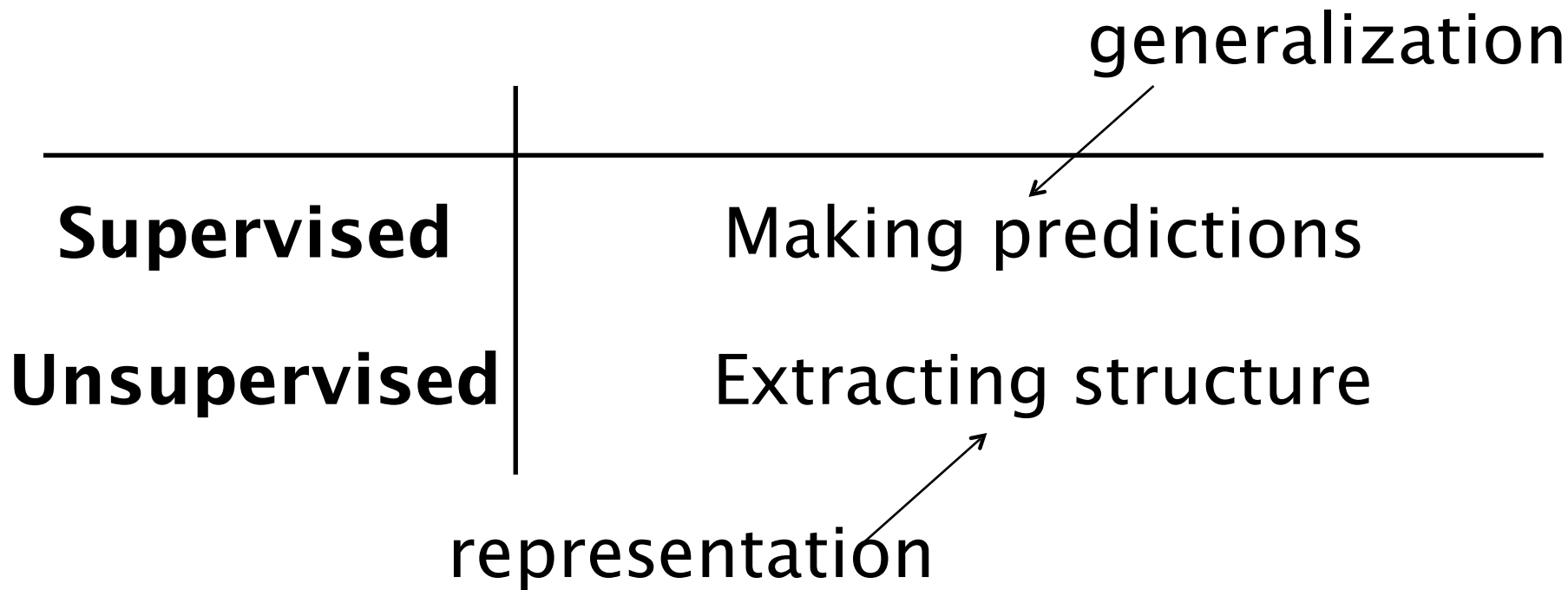


**NOTE**

Implementing solutions  
to ML problems is the  
focus of this course!

# THE STRUCTURE OF MACHINE LEARNING PROBLEMS

<b>Supervised</b>	Making predictions
<b>Unsupervised</b>	Extracting structure



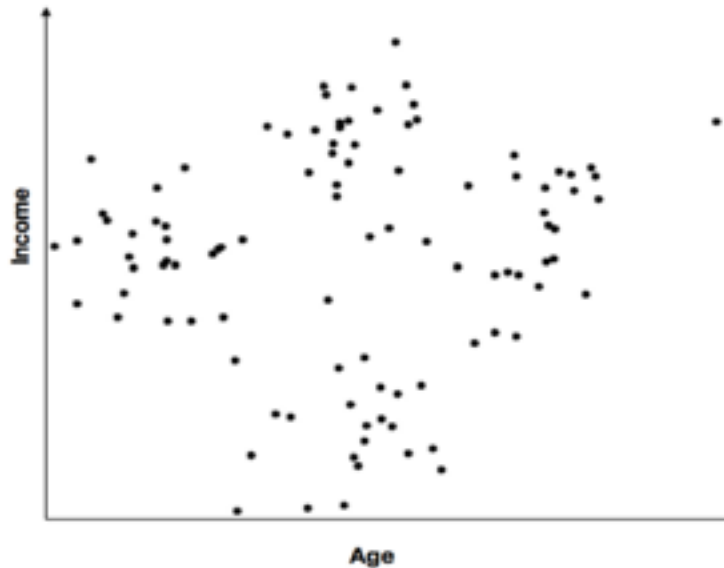


Supervised Learning - Can we create a function that predicts a value based on labeled training data?

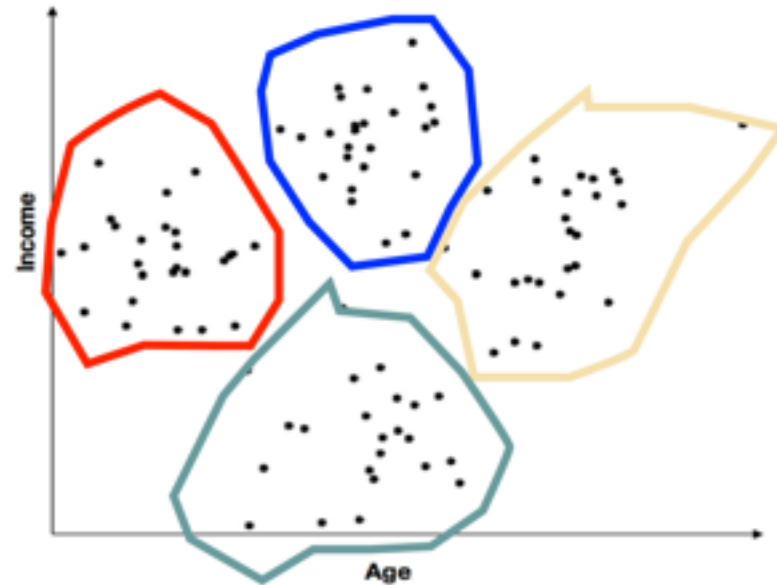
Regression example: Alan is 30 years old and can eat *four donuts an hour*. Betty is 60 years old, and can eat *two donuts an hour*. Cameron is 15 years old--how many donuts an hour eaten would be a good guess? This prediction is a regression model.

Classification example: Let's use the same data above. What is the probability that Cameron will eat eight donuts? Here, we have an answer and am now calculating the probability that an outcome has occurred.

Unsupervised Learning - Can we find structure to unlabeled data?



Unsupervised Learning - Can we find structure to unlabeled data?



	<b>Continuous</b>	<b>Categorical</b>
	Quantitative	Qualitative

Continuous	Categorical
Quantitative	Qualitative

### NOTE

The space where data live is called the *feature space*.

Each point in this space is called a *record*.

	Continuous	Categorical
Supervised	regression	classification
Unsupervised	dimension reduction	clustering

	Continuous	Categorical
Supervised	regression	classification
Unsupervised	dimension reduction	clustering

**NOTE**

We will implement solutions using *models* and *algorithms*.

Each will fall into one of these four buckets.

---

## QUESTION

---

WHAT IS THE GOAL OF  
MACHINE  
LEARNING?



<b>Supervised</b>	Making predictions
<b>Unsupervised</b>	Extracting structure

**ANSWER**

The goal is determined  
by the type of problem.

---

QUESTION

---

HOW DO YOU  
DETERMINE THE RIGHT  
APPROACH?

	Continuous	Categorical
Supervised	regression	classification
Unsupervised	dimension reduction	clustering

**ANSWER**

The right approach is determined by the desired solution.

	Continuous	Categorical
Supervised	regression	classification
Unsupervised	dimension reduction	clustering

**ANSWER**

**NOTE**

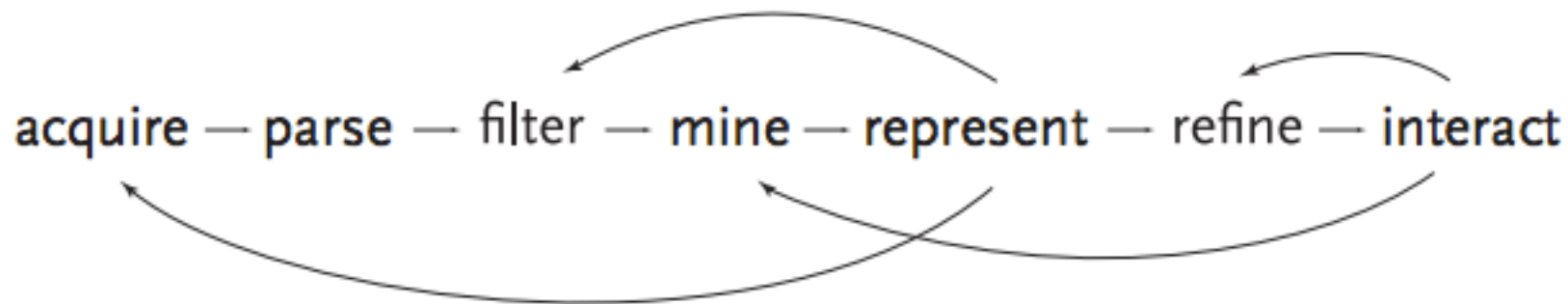
The  
det  
des  
All of this depends on  
your data!

---

QUESTION

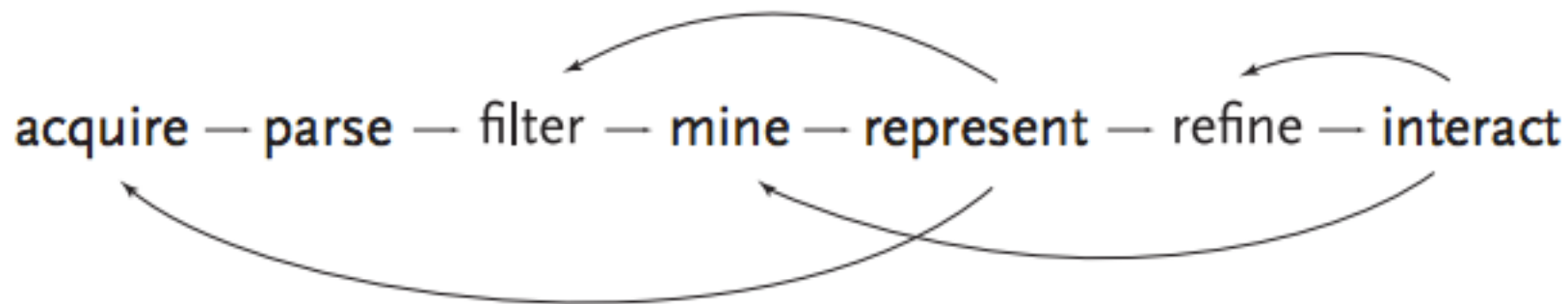
---

WHAT DO YOU DO WITH YOUR  
RESULTS?



### ANSWER

Interpret them and react accordingly.



### ANSWER

Int  
re

### NOTE

This also relies on your  
problem solving skills!

# **II. CLASSIFICATION PROBLEMS**



	Continuous	Categorical
Supervised	???	???
Unsupervised	???	???

	Continuous	Categorical
Supervised	regression	classification
Unsupervised	dimension reduction	clustering

Here's (part of) an example dataset:

Fisher's *Iris* Data

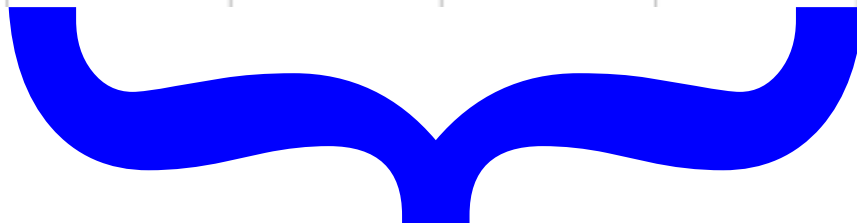
Sepal length ⇅	Sepal width ⇅	Petal length ⇅	Petal width ⇅	Species ⇅
5.1	3.5	1.4	0.2	<i>I. setosa</i>
4.9	3.0	1.4	0.2	<i>I. setosa</i>
4.7	3.2	1.3	0.2	<i>I. setosa</i>
4.6	3.1	1.5	0.2	<i>I. setosa</i>
5.0	3.6	1.4	0.2	<i>I. setosa</i>
5.4	3.9	1.7	0.4	<i>I. setosa</i>
4.6	3.4	1.4	0.3	<i>I. setosa</i>
5.0	3.4	1.5	0.2	<i>I. setosa</i>

Here's (part of) an example dataset:

Fisher's *Iris* Data

Sepal length ⇅	Sepal width ⇅	Petal length ⇅	Petal width ⇅	Species ⇅
5.1	3.5	1.4	0.2	<i>I. setosa</i>
4.9	3.0	1.4	0.2	<i>I. setosa</i>
4.7	3.2	1.3	0.2	<i>I. setosa</i>
4.6	3.1	1.5	0.2	<i>I. setosa</i>
5.0	3.6	1.4	0.2	<i>I. setosa</i>
5.4	3.9	1.7	0.4	<i>I. setosa</i>
4.6	3.4	1.4	0.3	<i>I. setosa</i>
5.0	3.4	1.5	0.2	<i>I. setosa</i>

independent  
variables



Here's (part of) an example dataset:

Fisher's Iris Data

Sepal length ⇅	Sepal width ⇅	Petal length ⇅	Petal width ⇅	Species ⇅
5.1	3.5	1.4	0.2	<i>I. setosa</i>
4.9	3.0	1.4	0.2	<i>I. setosa</i>
4.7	3.2	1.3	0.2	<i>I. setosa</i>
4.6	3.1	1.5	0.2	<i>I. setosa</i>
5.0	3.6	1.4	0.2	<i>I. setosa</i>
5.4	3.9	1.7	0.4	<i>I. setosa</i>
4.6	3.4	1.4	0.3	<i>I. setosa</i>
5.0	3.4	1.5	0.2	<i>I. setosa</i>

independent  
variables

class  
labels  
(qualitative)

Q: What does “supervised” mean?

Q: What does “supervised” mean?

A: We know the labels.

Fisher's Iris Data

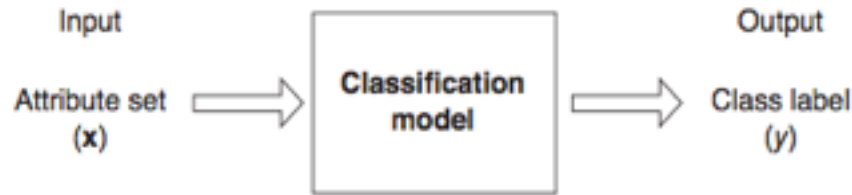
Sepal length ⇅	Sepal width ⇅	Petal length ⇅	Petal width ⇅	Species ⇅
5.1	3.5	1.4	0.2	<i>I. setosa</i>
4.9	3.0	1.4	0.2	<i>I. setosa</i>
4.7	3.2	1.3	0.2	<i>I. setosa</i>
4.6	3.1	1.5	0.2	<i>I. setosa</i>
5.0	3.6	1.4	0.2	<i>I. setosa</i>
5.4	3.9	1.7	0.4	<i>I. setosa</i>
4.6	3.4	1.4	0.3	<i>I. setosa</i>
5.0	3.4	1.5	0.2	<i>I. setosa</i>

class  
labels  
(qualitative)

Q: How does a classification problem work?

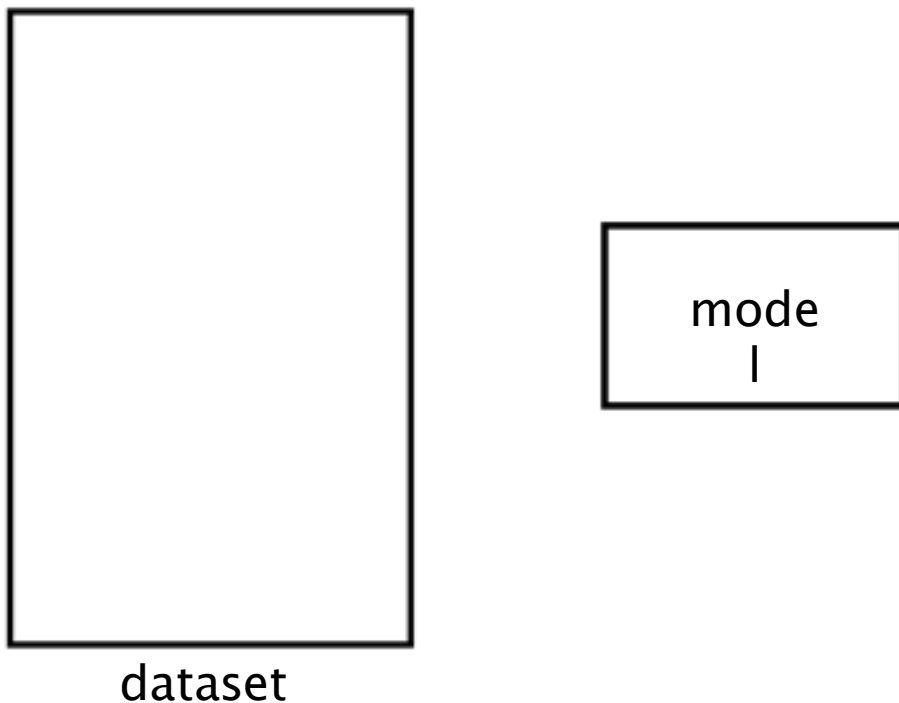


Q: How does a classification problem work?  
A: Data in, predicted labels out.



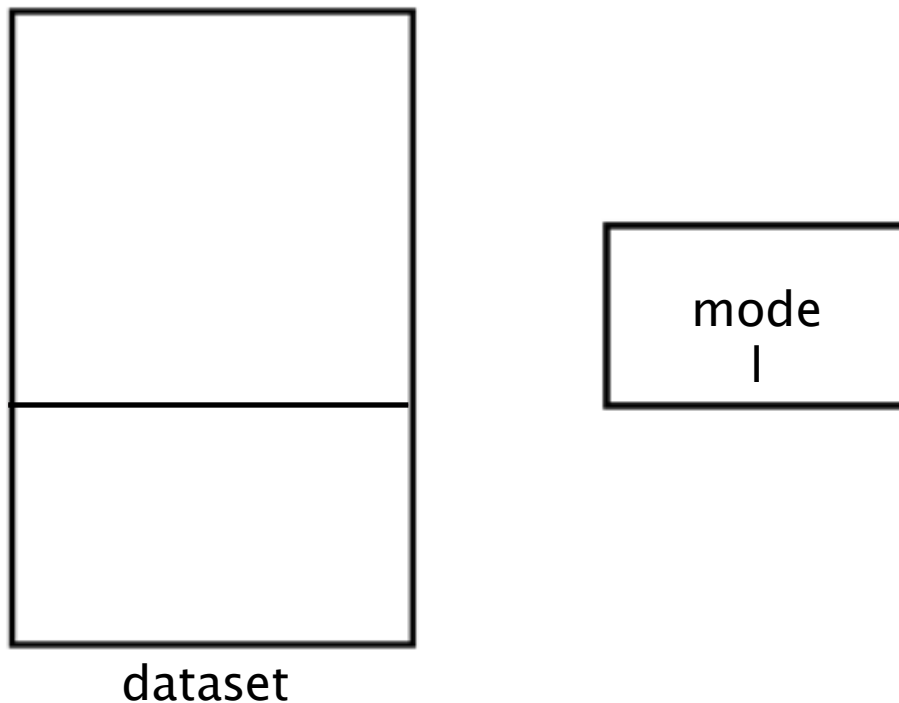
**Figure 4.2.** Classification as the task of mapping an input attribute set  $x$  into its class label  $y$ .

Q: What steps does a classification problem require?



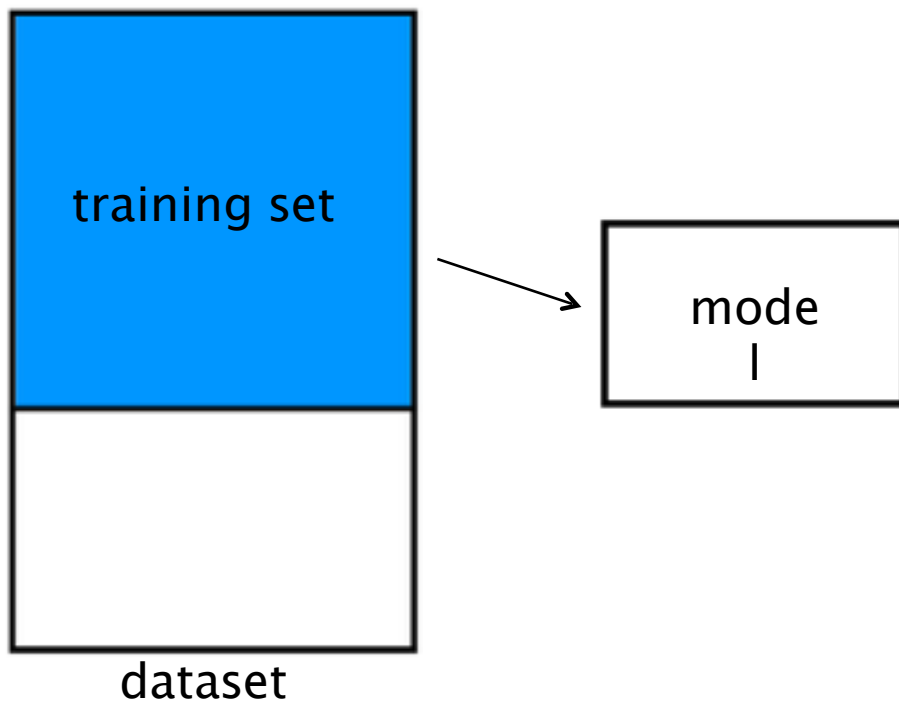
Q: What steps does a classification problem require?

1) split dataset



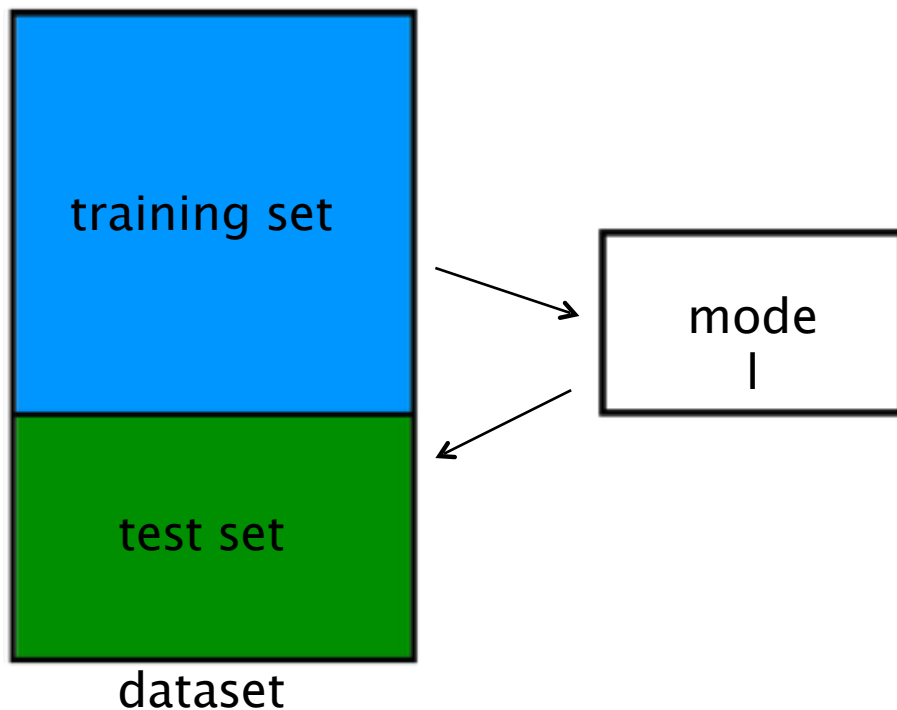
Q: What steps does a classification problem require?

- 1) split dataset
- 2) train model



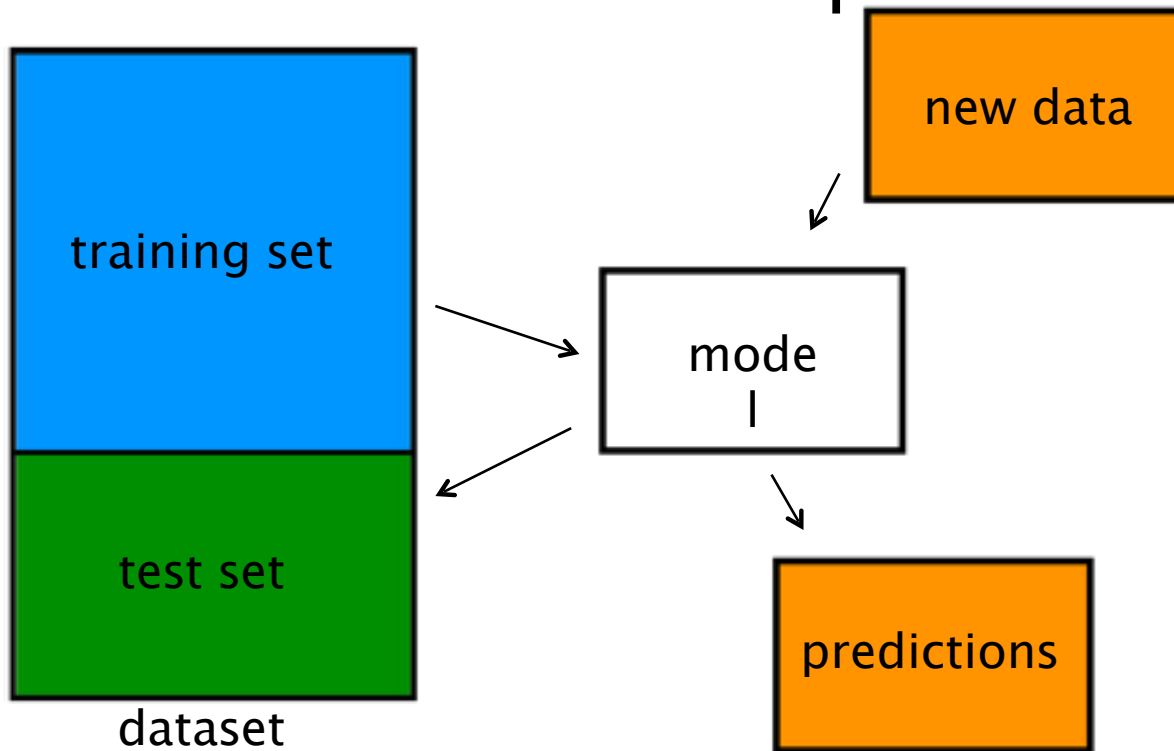
Q: What steps does a classification problem require?

- 1) split dataset
- 2) train model
- 3) test model



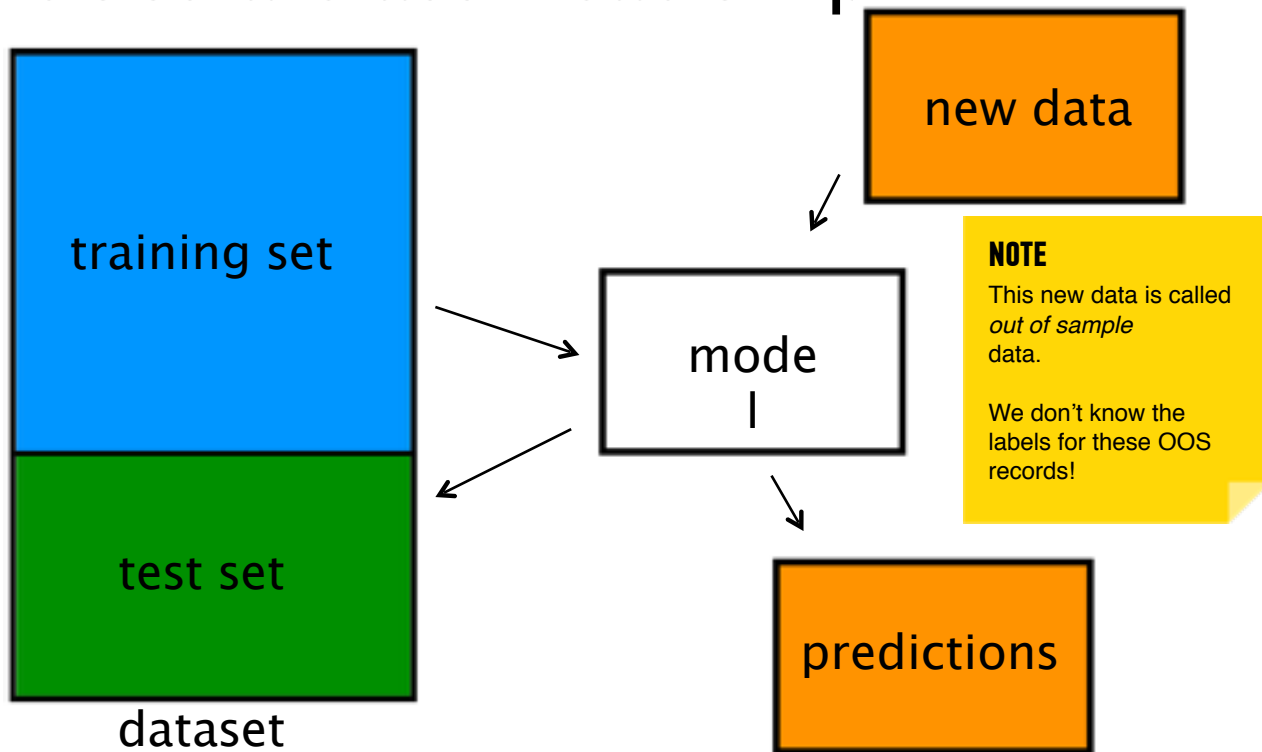
Q: What steps does a classification problem require?

- 1) split dataset
- 2) train model
- 3) test model
- 4) make predictions



Q: What steps does a classification problem require?

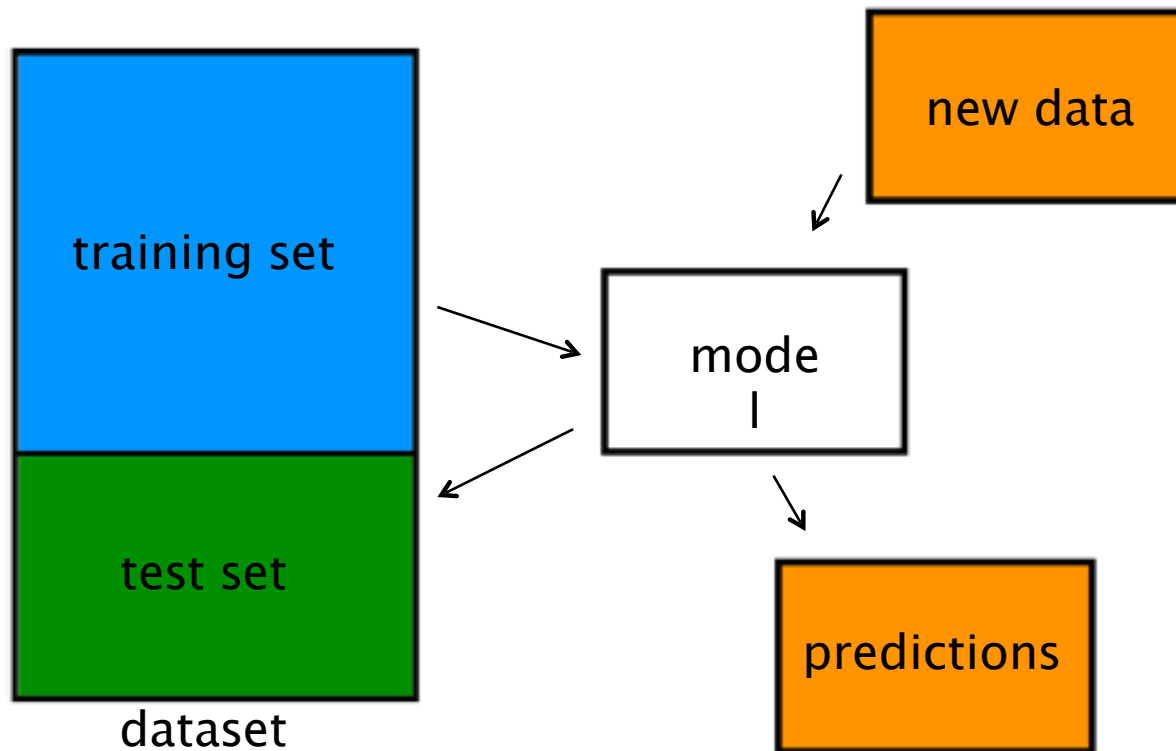
- 1) split dataset
- 2) train model
- 3) test model
- 4) make predictions



# **III. BUILDING EFFECTIVE CLASSIFIERS**

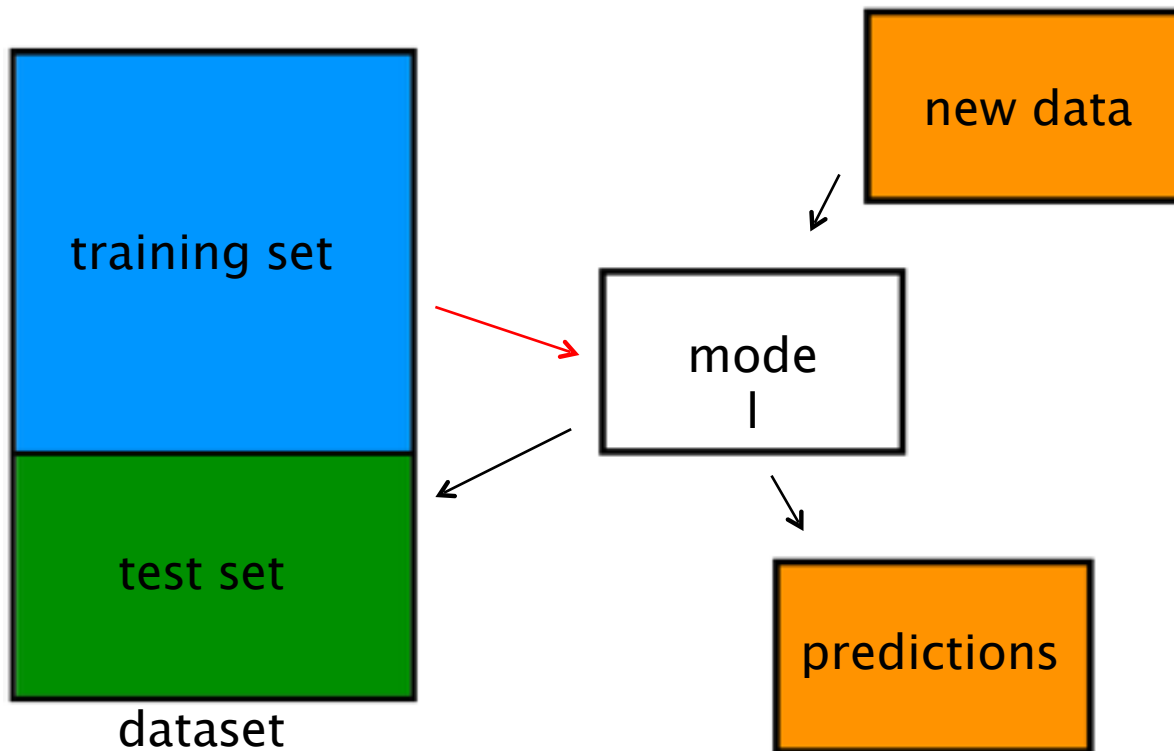


Q: What types of prediction error will we run into?



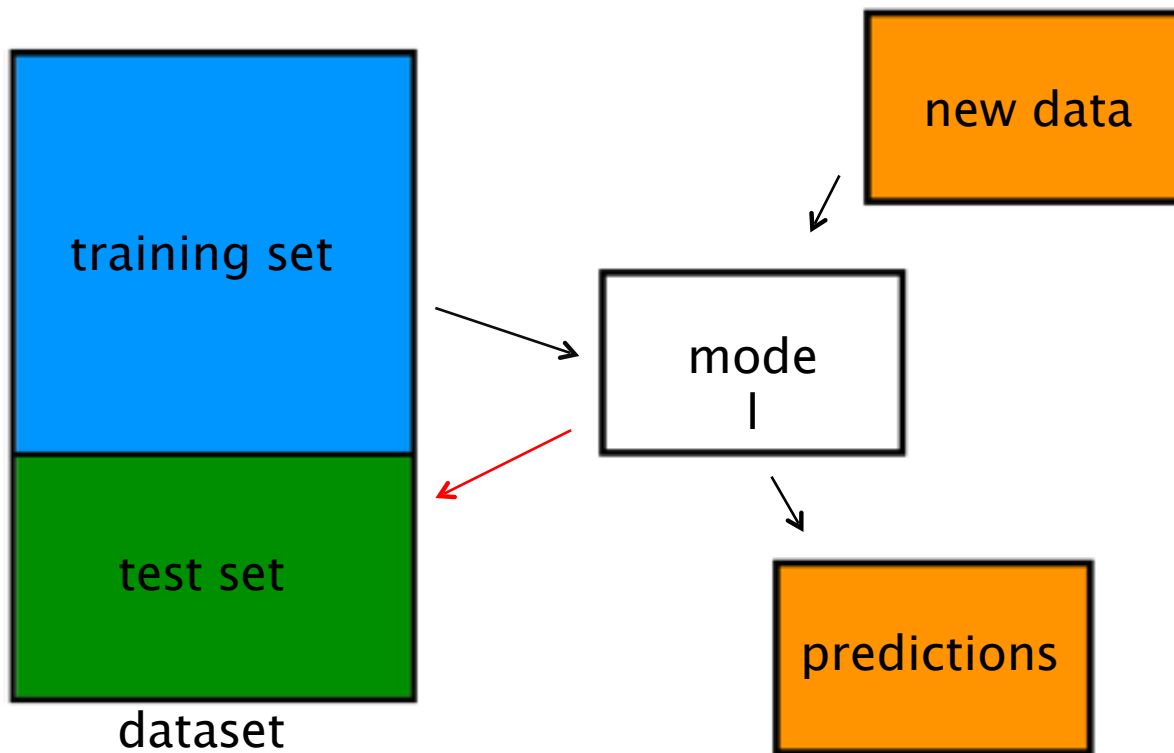
Q: What types of prediction error will we run into?

1) training error



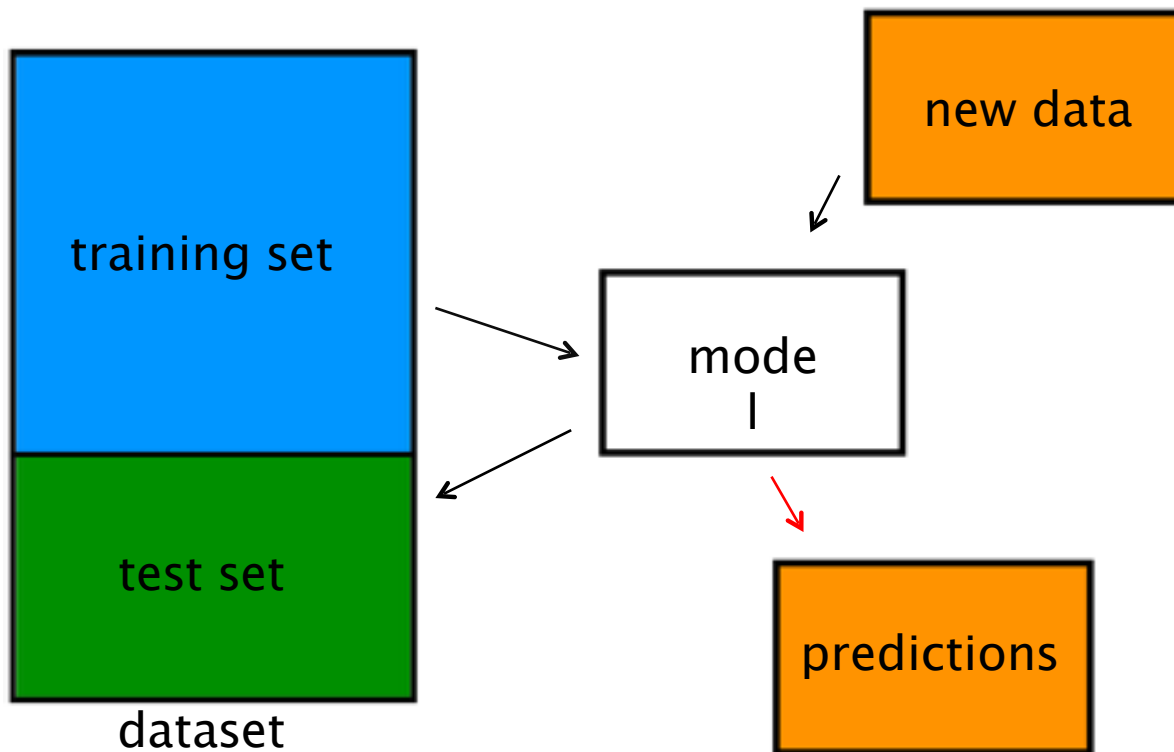
Q: What types of prediction error will we run into?

- 1) training error
- 2) generalization error



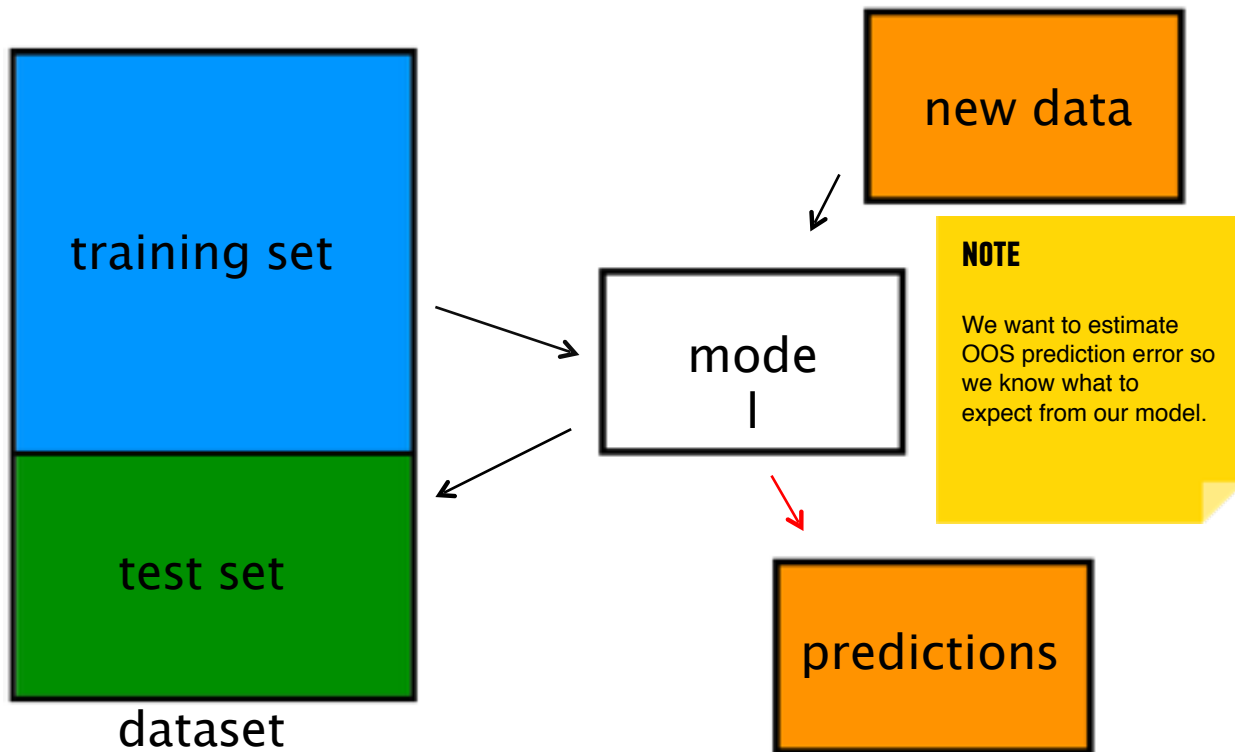
Q: What types of prediction error will we run into?

- 1) training error
- 2) generalization error
- 3) OOS error



Q: What types of prediction error will we run into?

- 1) training error
- 2) generalization error
- 3) OOS error



Q: Why should we use training & test sets?

Q: Why should we use training & test sets?

Thought experiment:

Suppose instead, we train our model using the entire dataset.

Q: Why should we use training & test sets?

Thought experiment:

Suppose instead, we train our model using the entire dataset.

Q: How low can we push the training error?



Q: Why should we use training & test sets?

Thought experiment:

Suppose instead, we train our model using the entire dataset.

Q: How low can we push the training error?

- We can make the model arbitrarily complex (effectively “memorizing” the entire training set).

Q: Why should we use training & test sets?

Thought experiment:

Suppose instead, we train our model using the entire dataset.

Q: How low can we push the training error?

- We can make the model arbitrarily complex (effectively “memorizing” the entire training set).

A: Down to zero!

Q: Why should we use training & test sets?

Thought experiment:

Suppose instead, we train our model using the entire dataset.

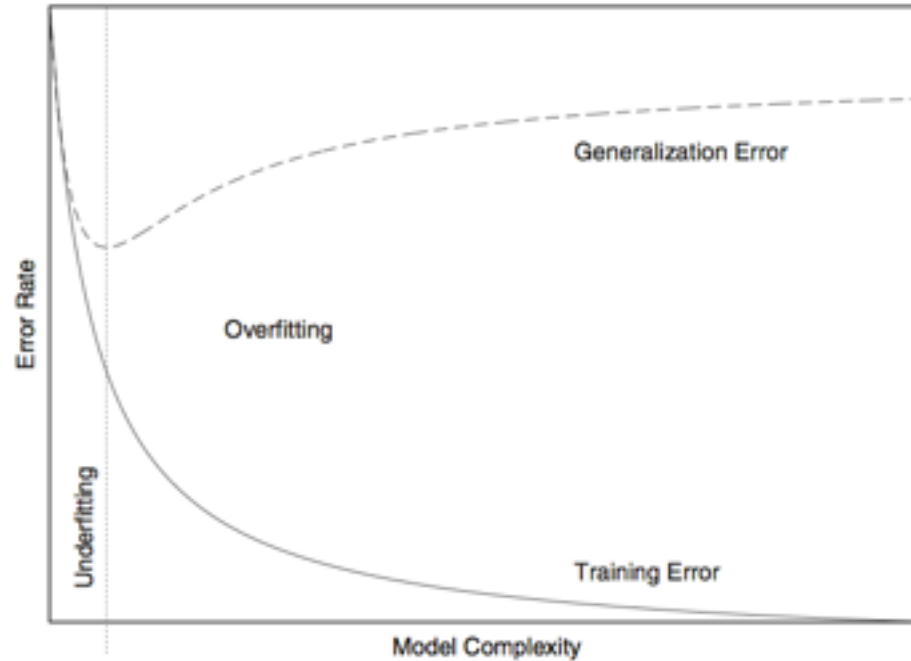
Q: How low can we push the training error?

- We can make the model arbitrarily complex (effectively “memorizing” the entire training set).

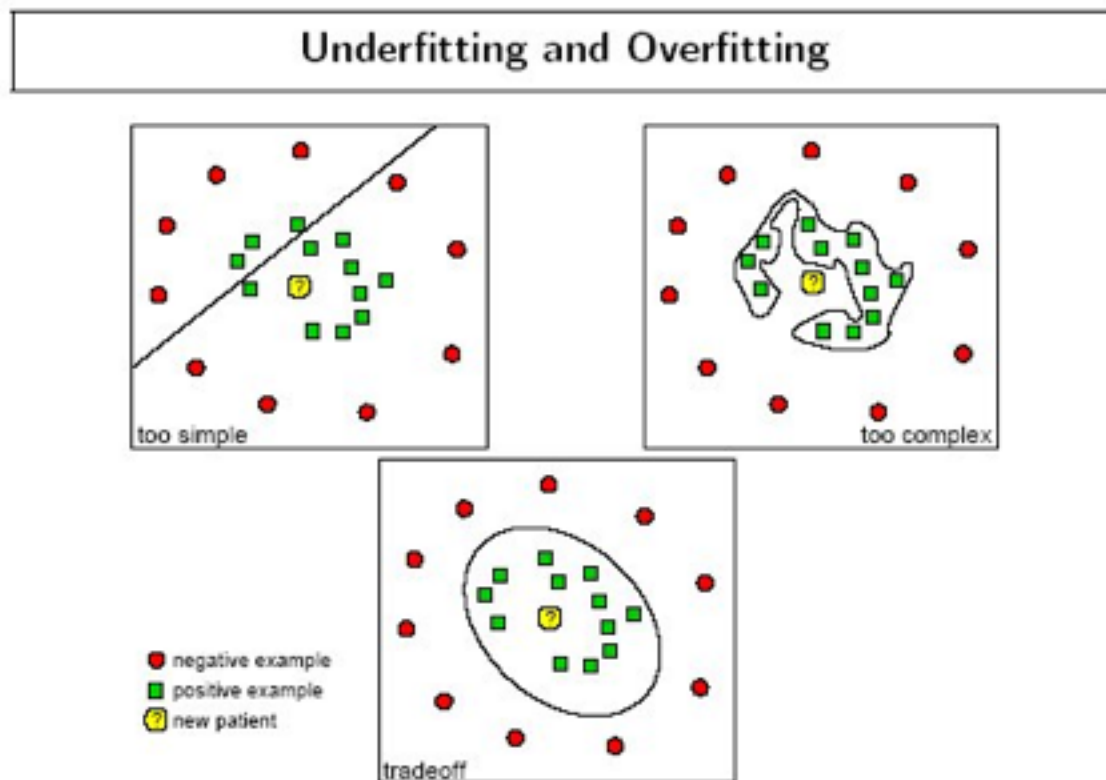
A: Down to zero!

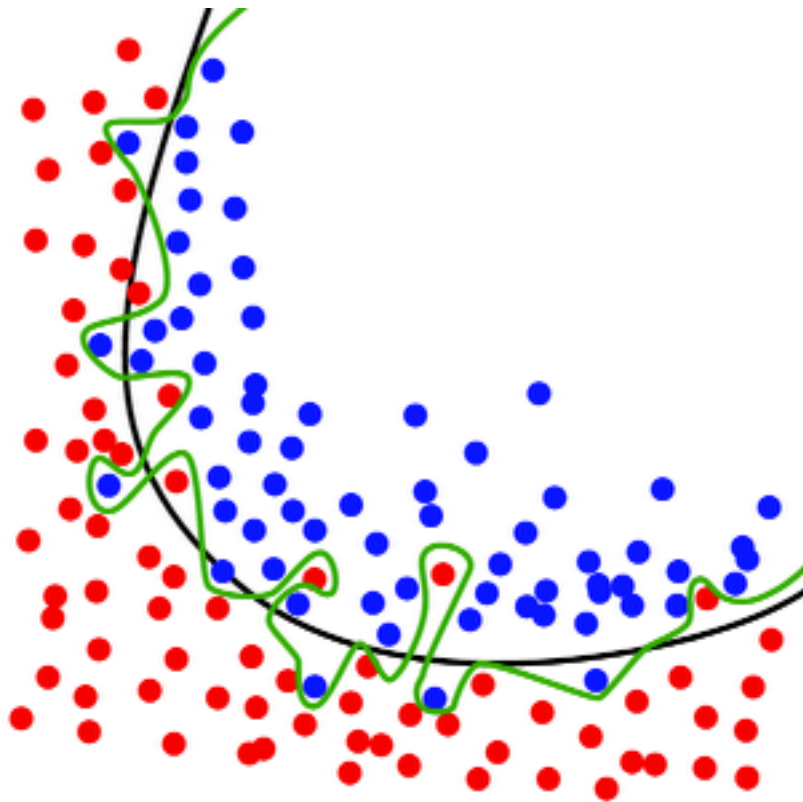
**NOTE**

This phenomenon  
is called  
*overfitting*.



**FIGURE 18-1.** Overfitting: as a model becomes more complex, it becomes increasingly able to represent the training data. However, such a model is overfitted and will not generalize well to data that was not used during training.





Q: Why should we use training & test sets?

Thought experiment:

Suppose instead, we train our model using the entire dataset.

Q: How low can we push the training error?

- We can make the model arbitrarily complex (effectively “memorizing” the entire training set).

A: Down to zero!

A: Training error is not a good estimate of OOS accuracy.

**NOTE**

This phenomenon is called *overfitting*.

Suppose we do the train/test split.



Suppose we do the train/test split.

Q: How well does generalization error predict OOS accuracy?

Suppose we do the train/test split.

Q: How well does generalization error predict OOS accuracy?

Thought experiment:

Suppose we had done a different train/test split.

Suppose we do the train/test split.

Q: How well does generalization error predict OOS accuracy?

Thought experiment:

Suppose we had done a different train/test split.

Q: Would the generalization error remain the same?

Suppose we do the train/test split.

Q: How well does generalization error predict OOS accuracy?

Thought experiment:

Suppose we had done a different train/test split.

Q: Would the generalization error remain the same?

A: Of course not!

Suppose we do the train/test split.

Q: How well does generalization error predict OOS accuracy?

Thought experiment:

Suppose we had done a different train/test split.

Q: Would the generalization error remain the same?

A: Of course not!

A: On its own, not very well.

Suppose we do the train/test split.

Q: How well does generalization error predict OOS accuracy?

Thought experiment:

Suppose we had done a different train/test split.

Q: Would the generalization error remain the same?

A: Of course not!

A: On its own, not very well.

**NOTE**

The generalization error gives a *high-variance estimate* of OOS accuracy.

Something is still missing!

Something is still missing!

Q: How can we do better?



Something is still missing!

Q: How can we do better?

Thought experiment:

Different train/test splits will give us different generalization errors.

Something is still missing!

Q: How can we do better?

Thought experiment:

Different train/test splits will give us different generalization errors.

Q: What if we did a bunch of these and took the average?

Something is still missing!

Q: How can we do better?

Thought experiment:

Different train/test splits will give us different generalization errors.

Q: What if we did a bunch of these and took the average?

A: Now you're talking!

Something is still missing!

Q: How can we do better?

Thought experiment:

Different train/test splits will give us different generalization errors.

Q: What if we did a bunch of these and took the average?

A: Now you're talking!

A: Cross-validation.

Steps for n-fold cross-validation:

## Steps for n-fold cross-validation:

- 1) Randomly split the dataset into  $n$  equal partitions.

## Steps for n-fold cross-validation:

- 1) Randomly split the dataset into  $n$  equal partitions.
- 2) Use partition 1 as test set & union of other partitions as training set.

## Steps for n-fold cross-validation:

- 1) Randomly split the dataset into  $n$  equal partitions.
- 2) Use partition 1 as test set & union of other partitions as training set.
- 3) Find generalization error.



## Steps for n-fold cross-validation:

- 1) Randomly split the dataset into  $n$  equal partitions.
- 2) Use partition 1 as test set & union of other partitions as training set.
- 3) Find generalization error.
- 4) Repeat steps 2–3 using a different partition as the test set at each iteration.

## Steps for n-fold cross-validation:

- 1) Randomly split the dataset into  $n$  equal partitions.
- 2) Use partition 1 as test set & union of other partitions as training set.
- 3) Find generalization error.
- 4) Repeat steps 2–3 using a different partition as the test set at each iteration.
- 5) Take the average generalization error as the estimate of OOS accuracy.

## Features of n-fold cross-validation:

## Features of n-fold cross-validation:

- 1) More accurate estimate of OOS prediction error.

## Features of n-fold cross-validation:

- 1) More accurate estimate of OOS prediction error.
- 2) More efficient use of data than single train/test split.
  - Each record in our dataset is used for both training and testing.

## Features of n-fold cross-validation:

- 1) More accurate estimate of OOS prediction error.
- 2) More efficient use of data than single train/test split.
  - Each record in our dataset is used for both training and testing.
- 3) Presents tradeoff between efficiency and computational expense.
  - 10-fold CV is 10x more expensive than a single train/test split

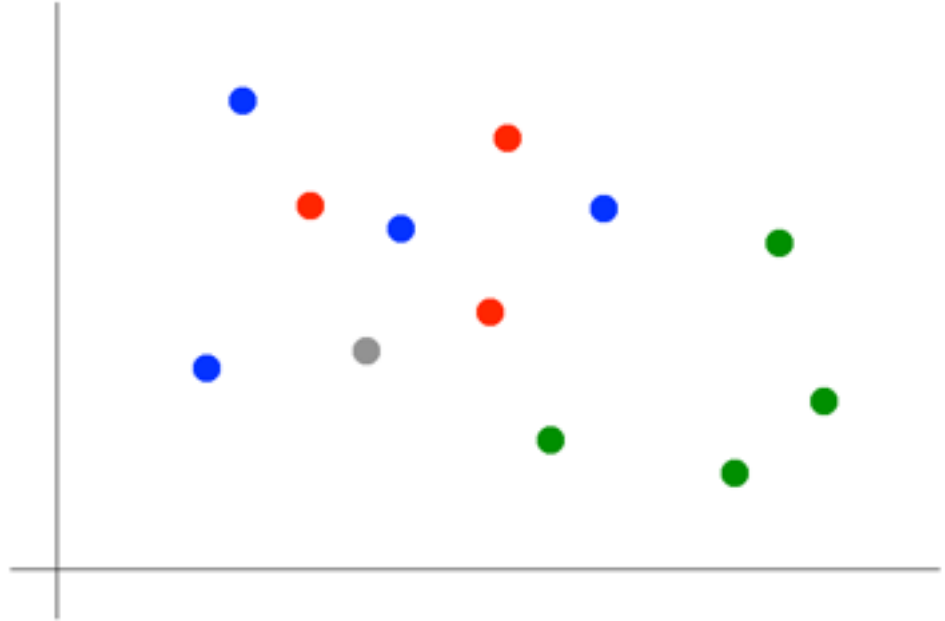
## Features of n-fold cross-validation:

- 1) More accurate estimate of OOS prediction error.
- 2) More efficient use of data than single train/test split.
  - Each record in our dataset is used for both training and testing.
- 3) Presents tradeoff between efficiency and computational expense.
  - 10-fold CV is 10x more expensive than a single train/test split
- 4) Can be used for model selection.

# **IV. KNN CLASSIFICATION**

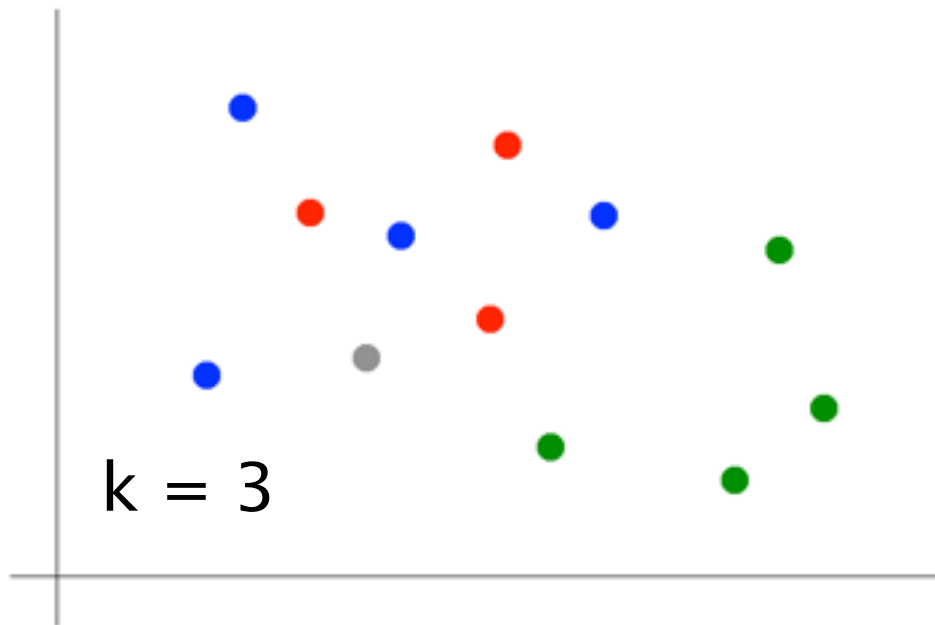


Suppose we want to predict the color of the grey dot.



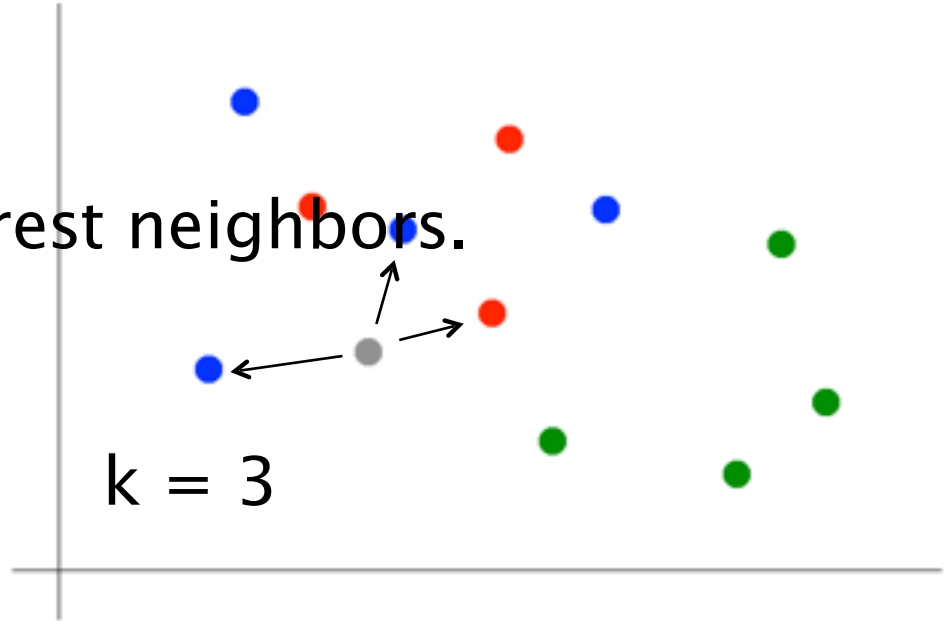
Suppose we want to predict the color of the grey dot.

1) Pick a value for  $k$ .



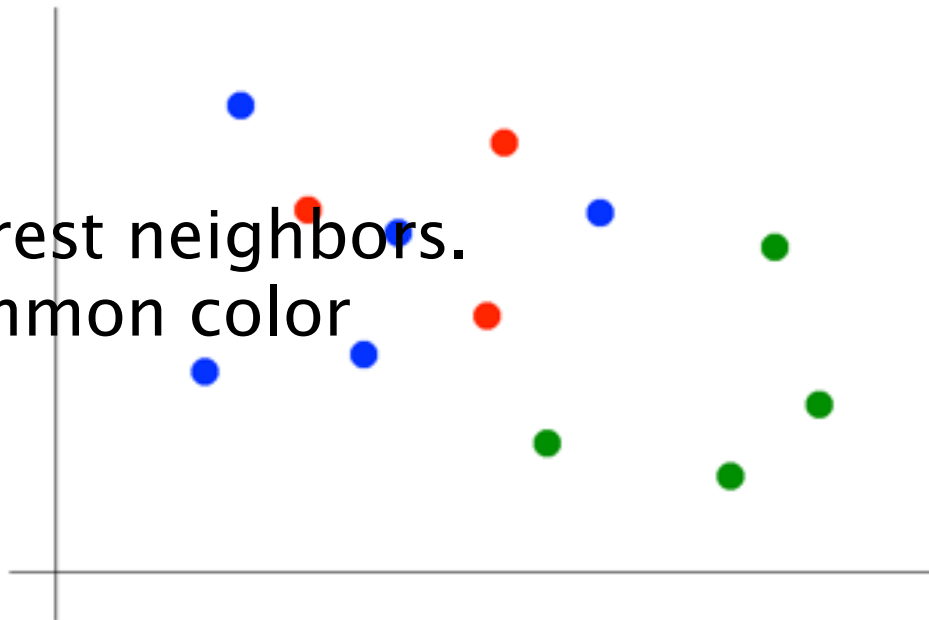
Suppose we want to predict the color of the grey dot.

- 1) Pick a value for  $k$ .
- 2) Find colors of  $k$  nearest neighbors.



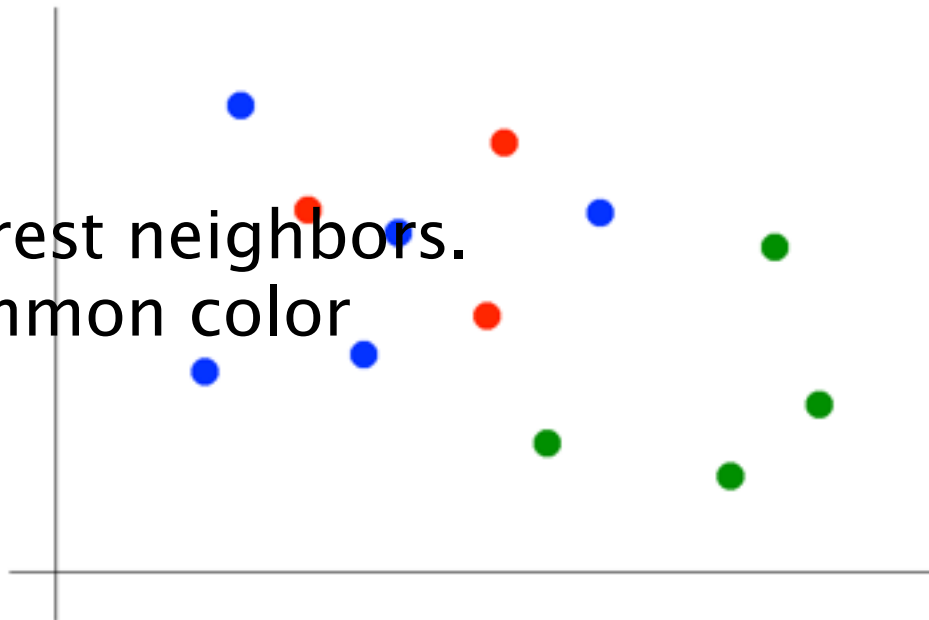
Suppose we want to predict the color of the grey dot.

- 1) Pick a value for  $k$ .
- 2) Find colors of  $k$  nearest neighbors.
- 3) Assign the most common color to the grey dot.



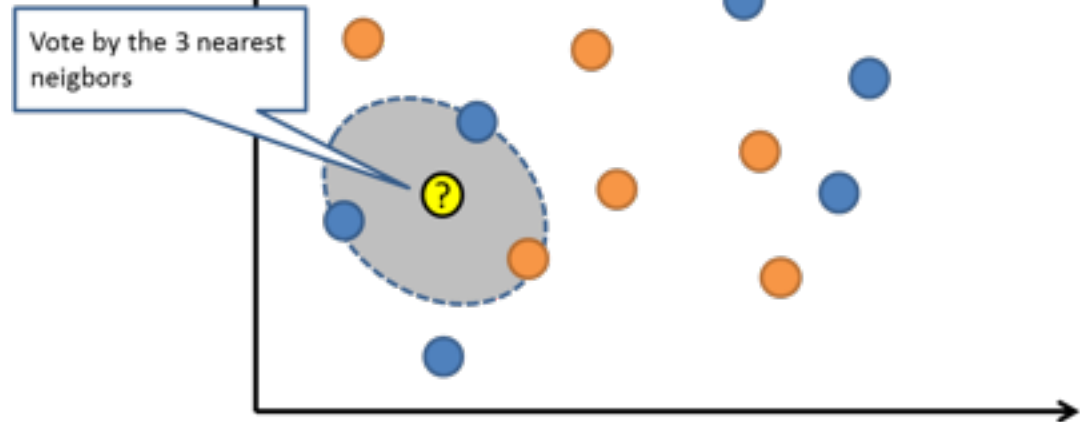
Suppose we want to predict the color of the grey dot.

- 1) Pick a value for  $k$ .
- 2) Find colors of  $k$  nearest neighbors.
- 3) Assign the most common color to the grey dot.

**OPTIONAL NOTE**

Our definition of "nearest" implicitly uses the *Euclidean distance function*.

Another example with  $k = 3$   
Will our new example be  
blue or orange?



---

**INTRO TO DATA SCIENCE**

---

**LABS**