

INTRO to DATA SCIENCE

SESSION 3: INTRO to MACHINE LEARNING with KNN

Rob Hall

DAT13 SF // March 16, 2015

RECAP

LAST TIME:

- INTRO TO PYTHON**
- LAB: INTRO TO NUMPY & PANDAS**

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?

- *Types of machine learning problems / algorithms*
- *Generalization*
- *Types of error (training, generalization, OOS)*
- *Over-fitting and under-fitting*
- *Cross-validation*

"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.”

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

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

WHAT IS MACHINE LEARNING?

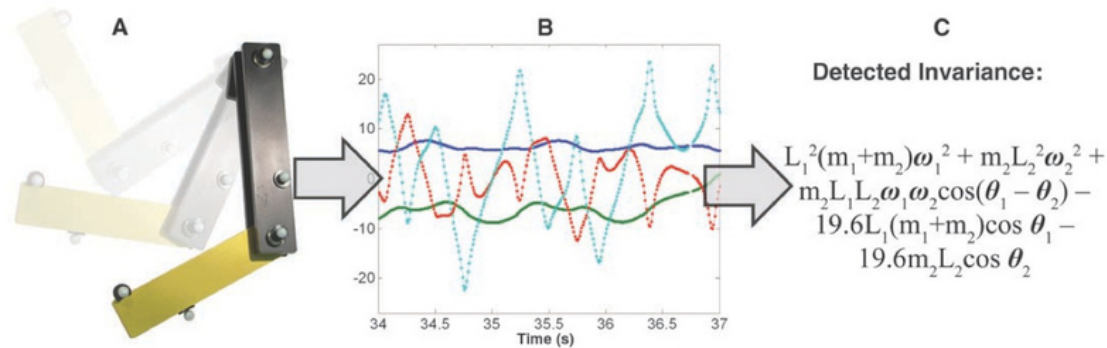
- **Machine learning** *is an area in computer science that studies and develops algorithms that can learn from data.*
- **Machine learning** *is a set of methods that can automatically detect patterns in data and use the discovered patterns to predict future data or perform other kinds of decision making*
- *Statistical learning theory, Pattern recognition*

WHEN DO WE NEED MACHINE LEARNING?

Where we need it:

- *Some observable patterns exist*
- *There no explicitly known equations or dependencies (formulas)*
- *We have data on it*

Example: Newton's second law of motion, conservation of mechanical energy, pendulum motion



From "Distilling Free-Form Natural Laws from Experimental Data." M. Schmidt and H.Lipson. Science, 2009.

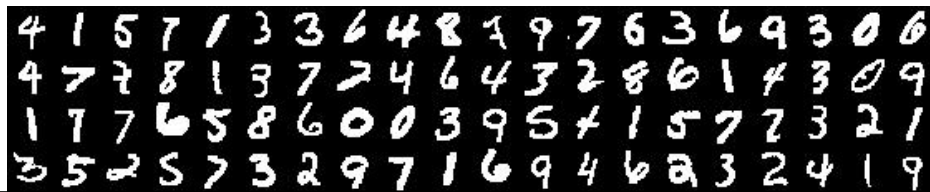
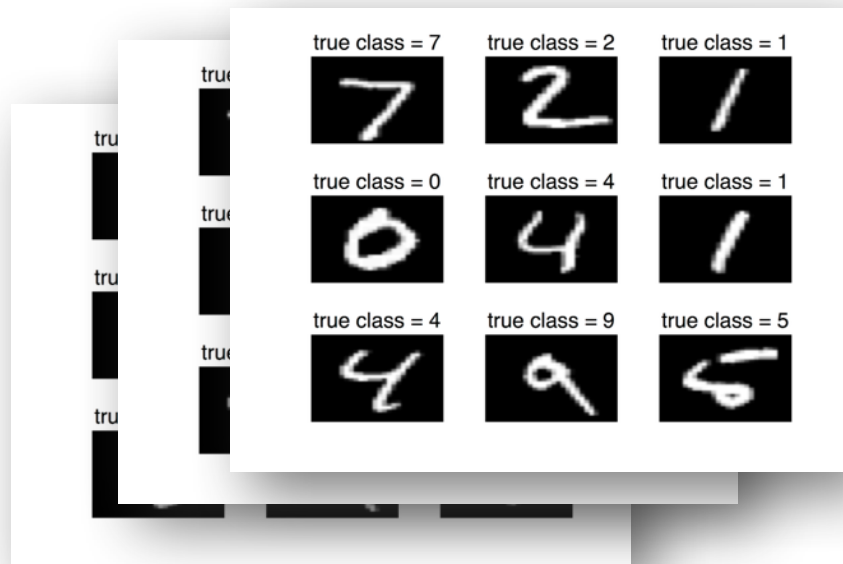
WHAT IS LEARNING?

Learning is not about memorizing and being able to recall, it is about generalizing the conclusions to previously unseen examples

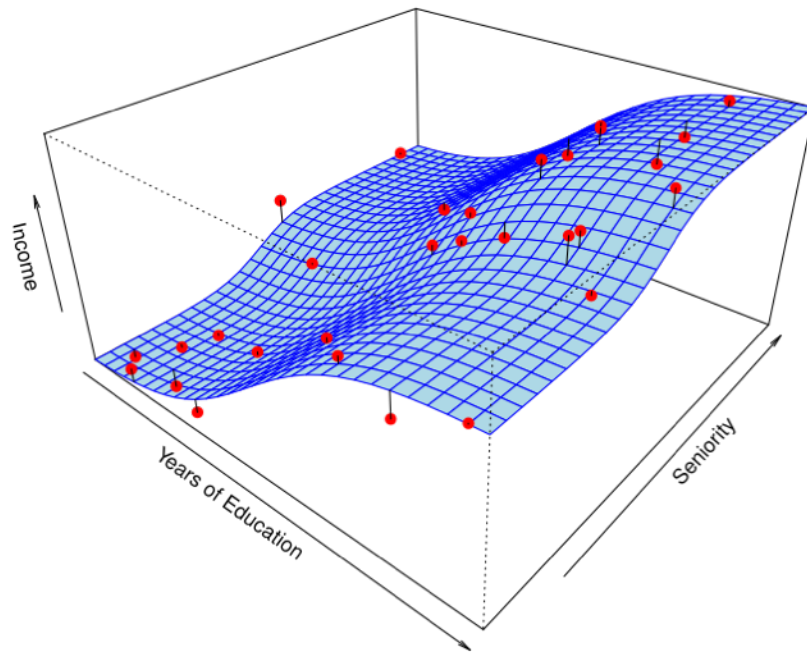
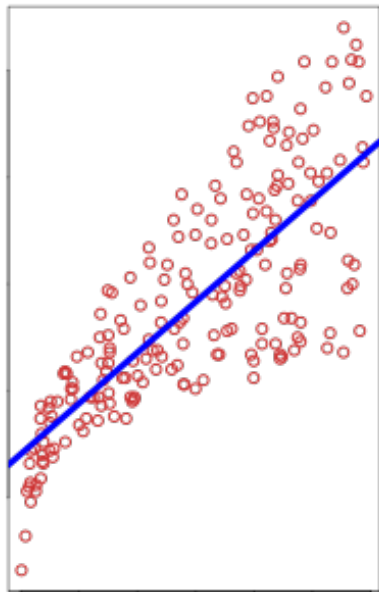
TYPES OF LEARNING?

- **Supervised learning:** the goal is to learn mapping from given inputs x to outputs y , given a labeled set of input-output pairs (when the training set contains explicit examples of what correct output should be for given input)
- **Unsupervised learning:** the goal is to learn interesting patterns and structure in data given only inputs (no output information given at all)

SUPERVISED LEARNING: CLASSIFICATION



SUPERVISED LEARNING: REGRESSION



SUPERVISED LEARNING: EXAMPLE

Credit Scoring

	<i>Client 1</i>	<i>Client 2</i>	<i>Client 3</i>
<i>Age</i>	23	30	19
<i>Gender</i>	<i>M</i>	<i>F</i>	<i>M</i>
<i>Annual salary</i>	\$30,000	\$45,000	\$15,000
<i>Years in residence</i>	3 years	1 year	3 month
<i>Years in job</i>	1 year	1 year	1 month
<i>Current debt</i>	\$5,000	\$1,000	\$10,000
<i>Paid off credit</i>	Yes	Yes	No

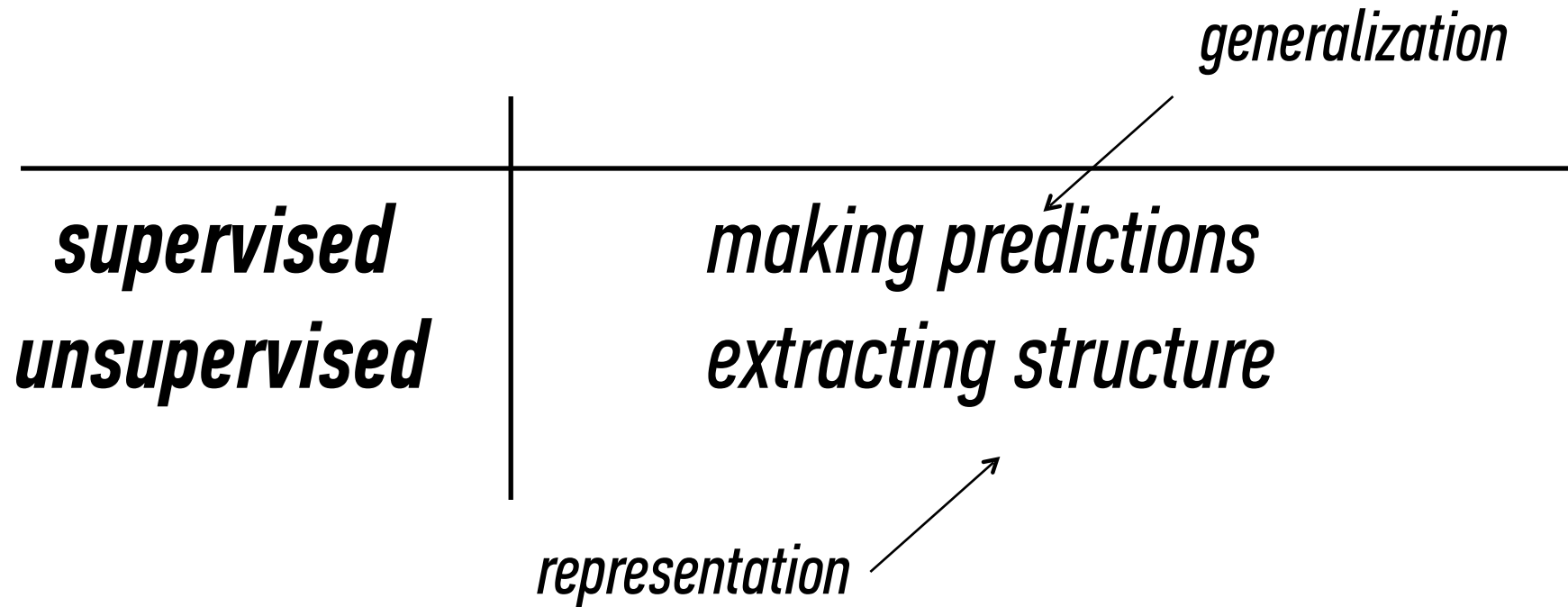
SUPERVISED LEARNING: EXAMPLE

Credit Scoring

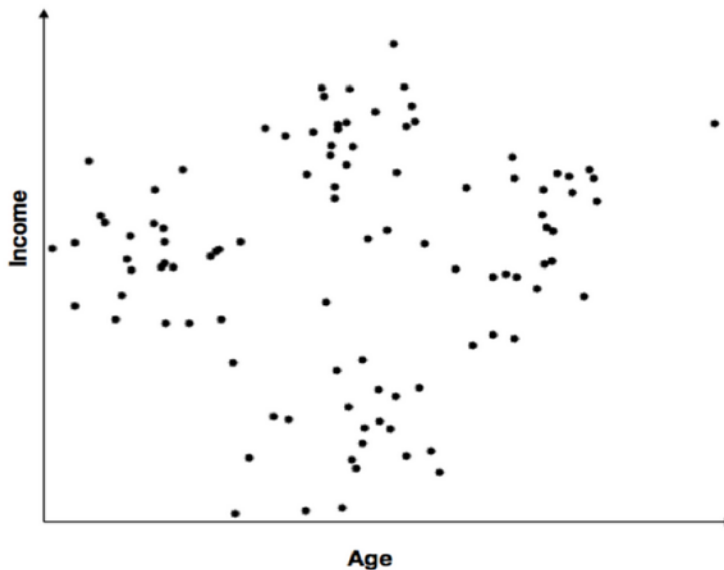
	<i>Applicant</i>
<i>Age</i>	25
<i>Gender</i>	M
<i>Annual salary</i>	\$25,000
<i>Years in residence</i>	1 year
<i>Years in job</i>	2 year3
<i>Current debt</i>	\$15,000
<i>Credit decision/ score</i>	???

THE STRUCTURE OF MACHINE LEARNING PROBLEMS

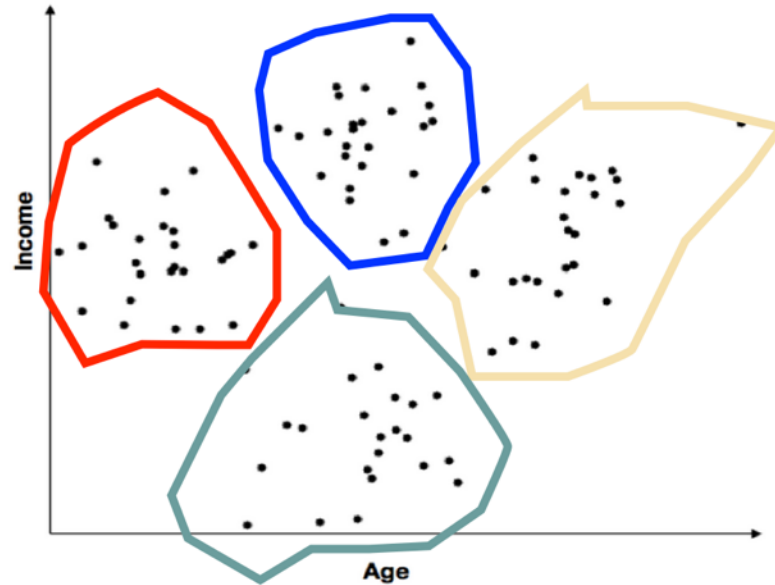
<i>supervised</i>	<i>making predictions</i>
<i>unsupervised</i>	<i>extracting structure</i>



Unsupervised Learning - Can we find structure to unlabeled data?



Unsupervised Learning - Can we find structure to unlabeled data?



continuous

categorical

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*.

	<i>continuous</i>	<i>categorical</i>
<i>supervised</i>	<i>regression</i>	<i>classification</i>
<i>unsupervised</i>	<i>dimension reduction</i>	<i>clustering</i>

	<i>continuous</i>	<i>categorical</i>
<i>supervised</i>	<i>regression</i>	<i>classification</i>
<i>unsupervised</i>	<i>dimension reduction</i>	<i>clustering</i>

NOTE

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

Each will fall into one of these four buckets.

	<i>continuous</i>	<i>categorical</i>
<i>supervised</i>	<i>regression</i>	<i>classification</i>
<i>unsupervised</i>	<i>dimension reduction</i>	<i>clustering</i>

ANSWER

The right approach is determined by the desired solution.

	<i>continuous</i>	<i>categorical</i>
<i>supervised</i>	<i>regression</i>	<i>classification</i>
<i>unsupervised</i>	<i>dimension reduction</i>	<i>clustering</i>

ANSWER**NOTE**

The is d
des
All of this depends on
your data!

II. CLASSIFICATION PROBLEMS

	<i>continuous</i>	<i>categorical</i>
<i>supervised</i>	???	???
<i>unsupervised</i>	???	???

	<i>continuous</i>	<i>categorical</i>
<i>supervised</i>	<i>regression</i>	<i>classification</i>
<i>unsupervised</i>	<i>dimension reduction</i>	<i>clustering</i>

Here's (part of) an example dataset:

Fisher's *Iris* Data

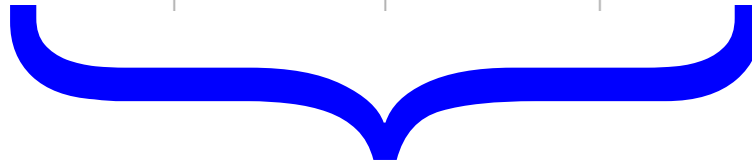
Sepal length ⇅	Sepal width ⇅	Petal length ⇅	Petal width ⇅	Species ⇅
5.1	3.5	1.4	0.2	<i>I. setosa</i>
4.9	3.0	1.4	0.2	<i>I. setosa</i>
4.7	3.2	1.3	0.2	<i>I. setosa</i>
4.6	3.1	1.5	0.2	<i>I. setosa</i>
5.0	3.6	1.4	0.2	<i>I. setosa</i>
5.4	3.9	1.7	0.4	<i>I. setosa</i>
4.6	3.4	1.4	0.3	<i>I. setosa</i>
5.0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

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.

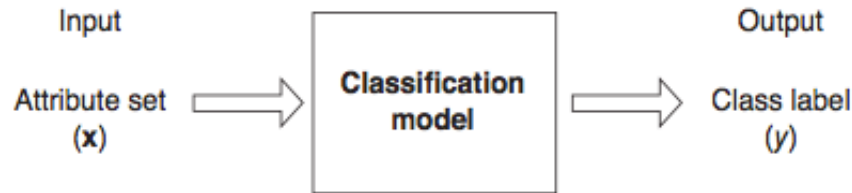
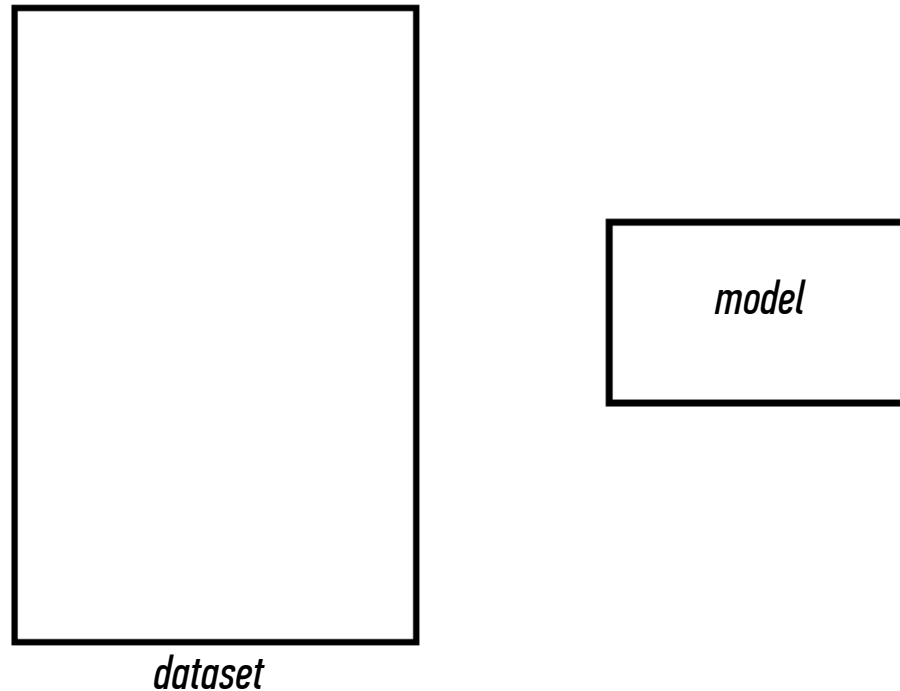


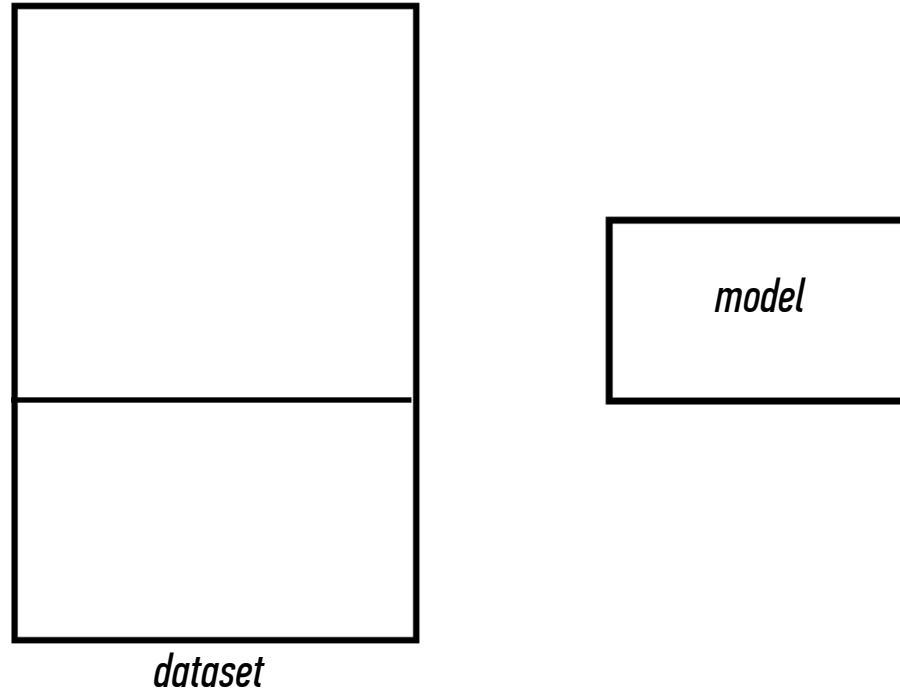
Figure 4.2. Classification as the task of mapping an input attribute set x into its class label y .

Q: What steps does a classification problem require?



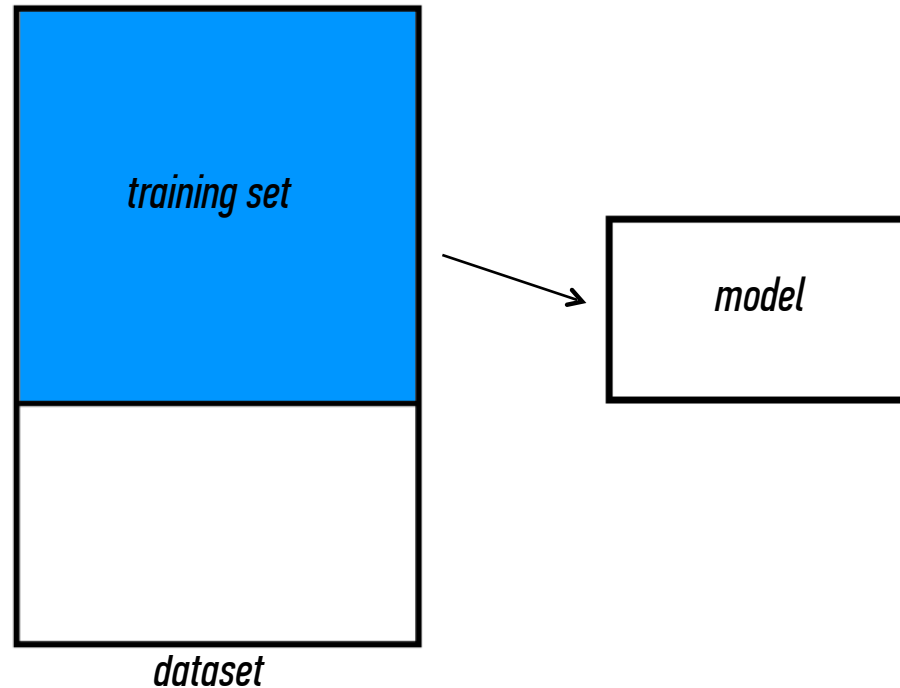
Q: What steps does a classification problem require?

1) split dataset



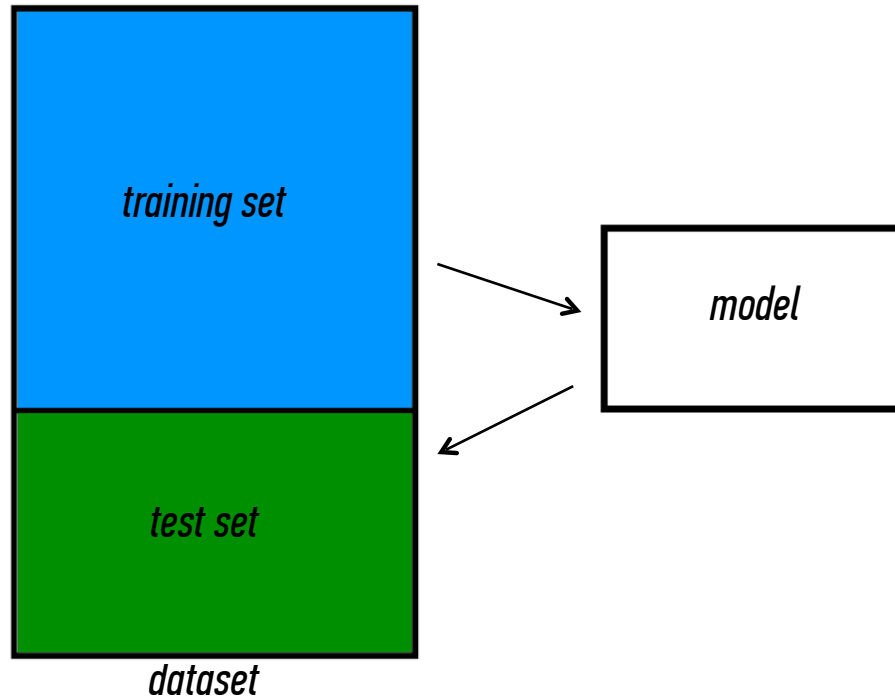
Q: What steps does a classification problem require?

- 1) split dataset*
- 2) train model*



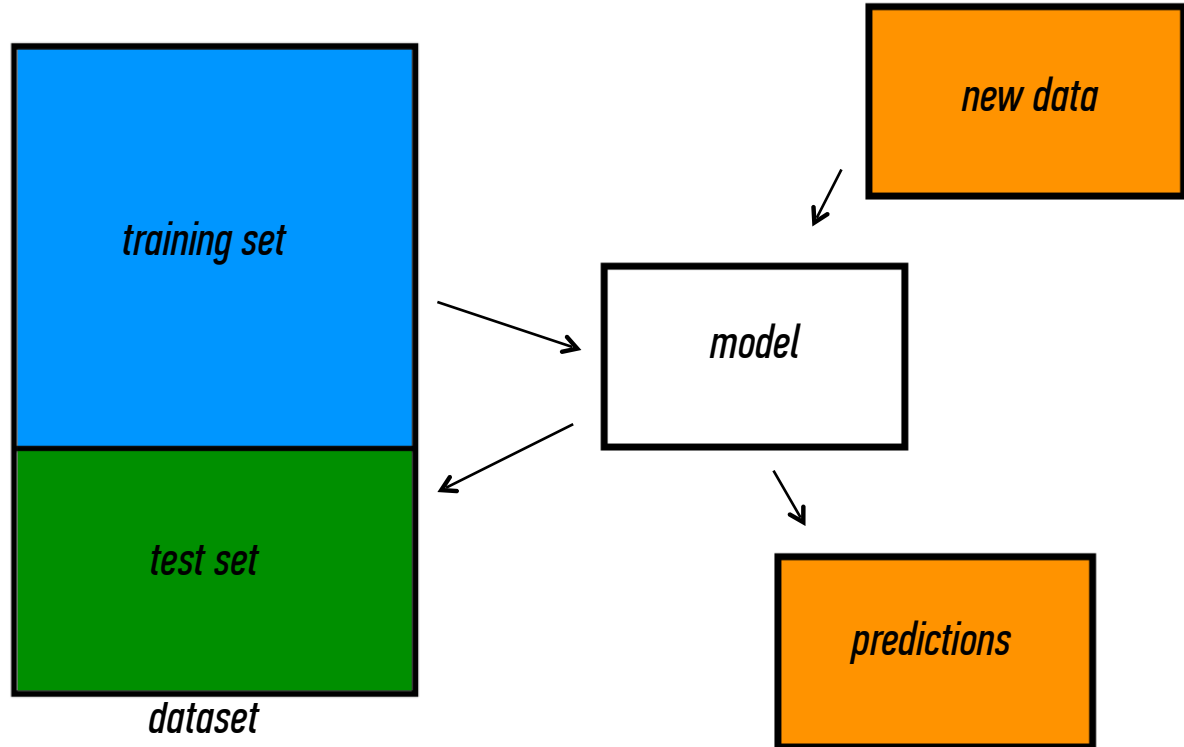
Q: What steps does a classification problem require?

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



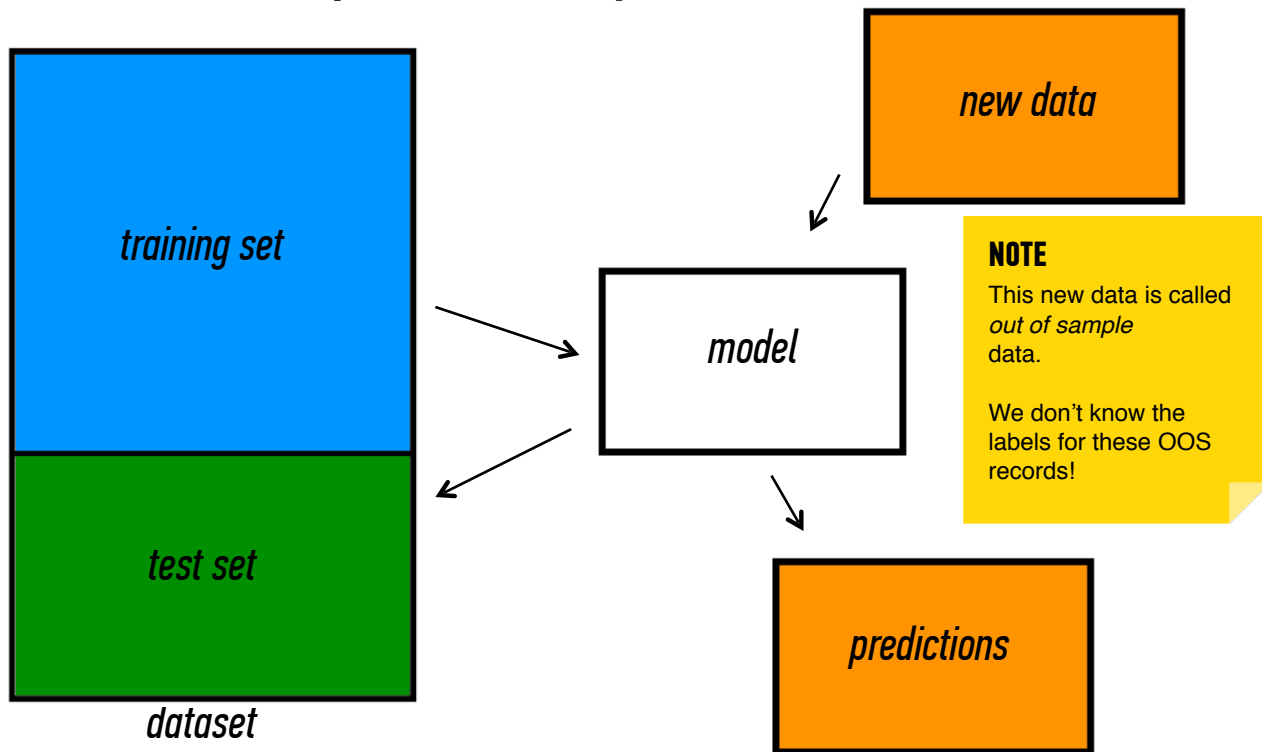
Q: What steps does a classification problem require?

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



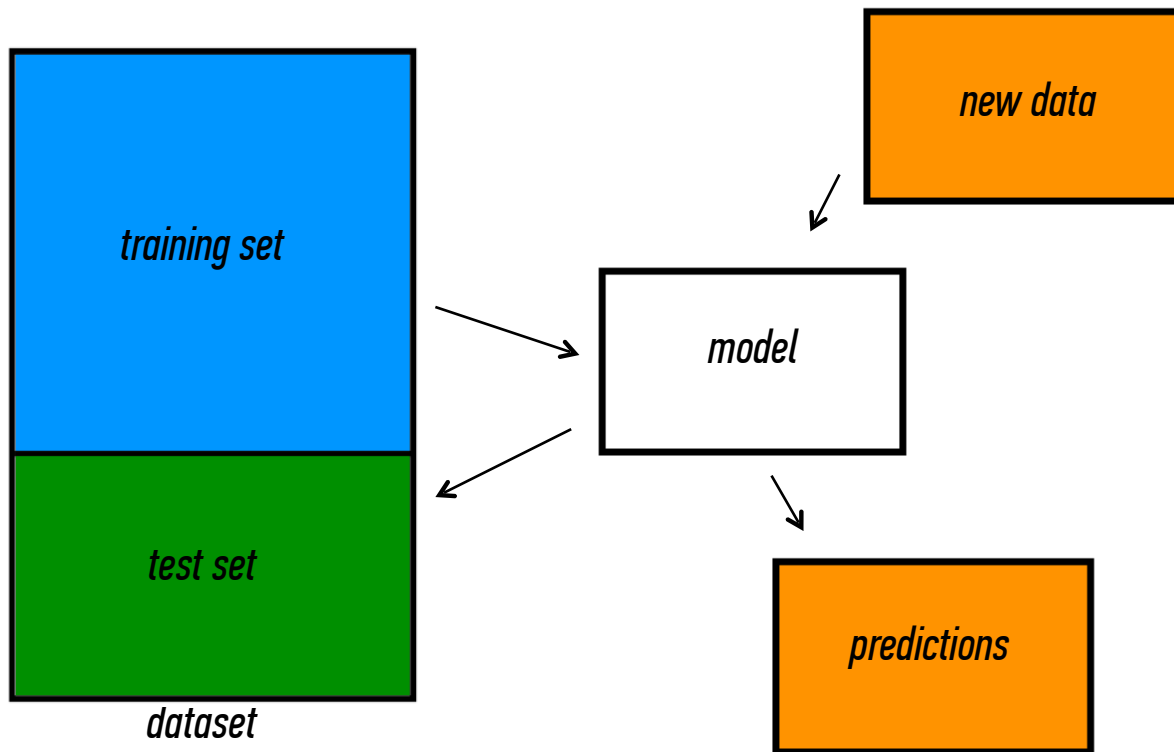
Q: What steps does a classification problem require?

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



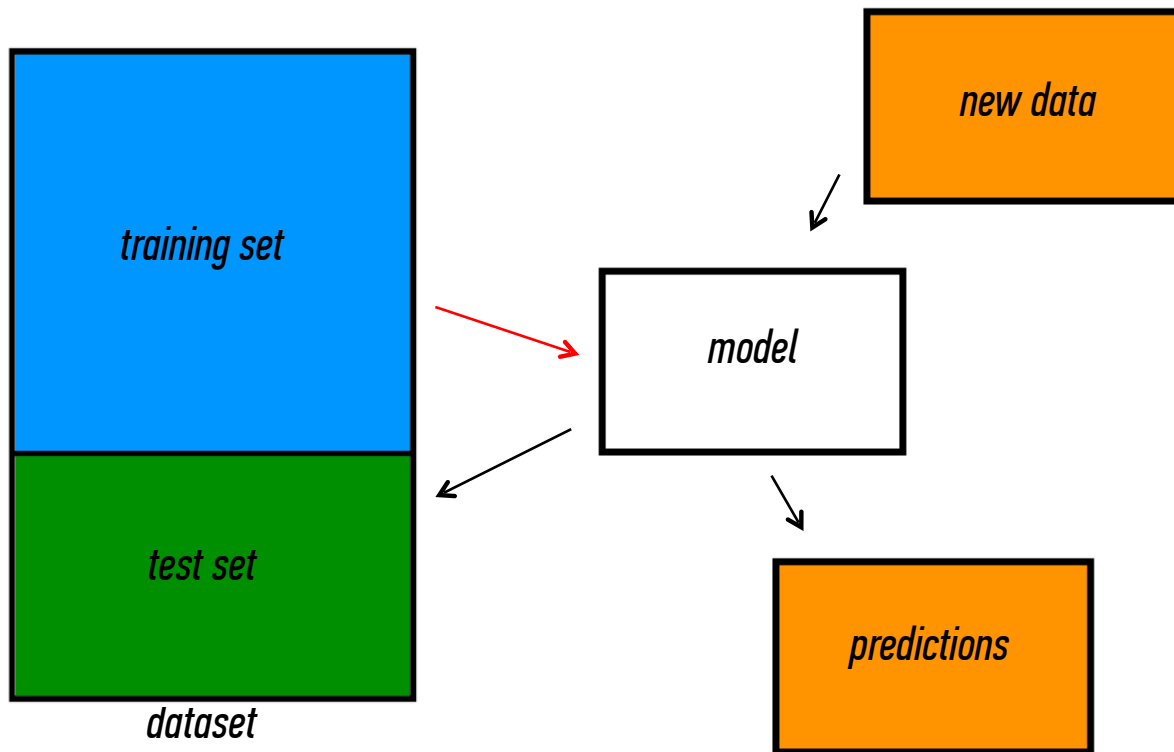
III. BUILDING EFFECTIVE CLASSIFIERS

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



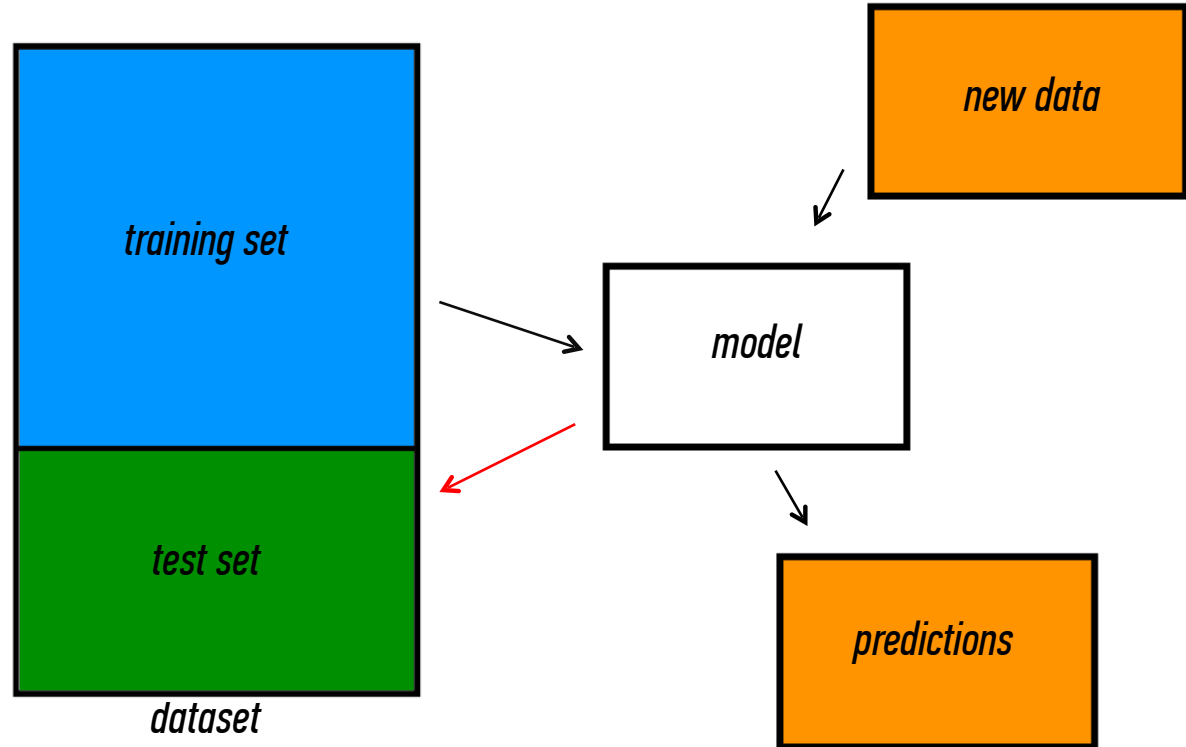
Q: What types of prediction error will we run into?

1) training error



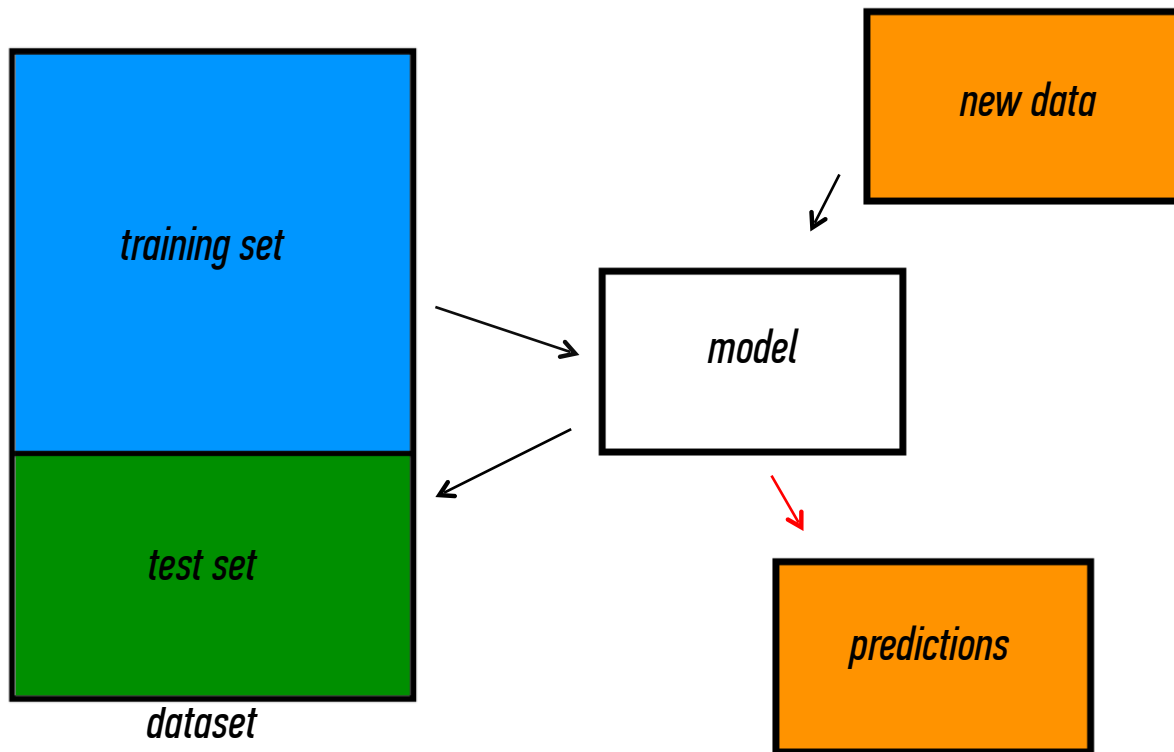
Q: What types of prediction error will we run into?

- 1) training error*
- 2) generalization error*



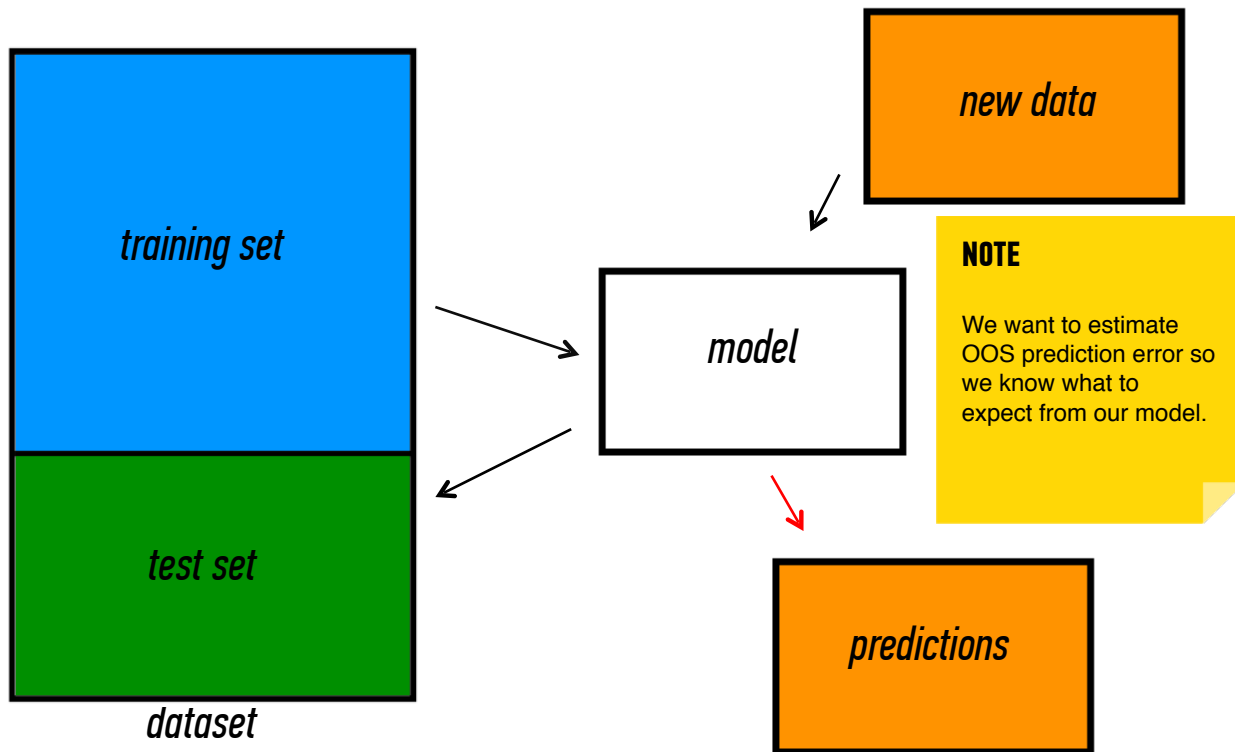
Q: What types of prediction error will we run into?

- 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*.

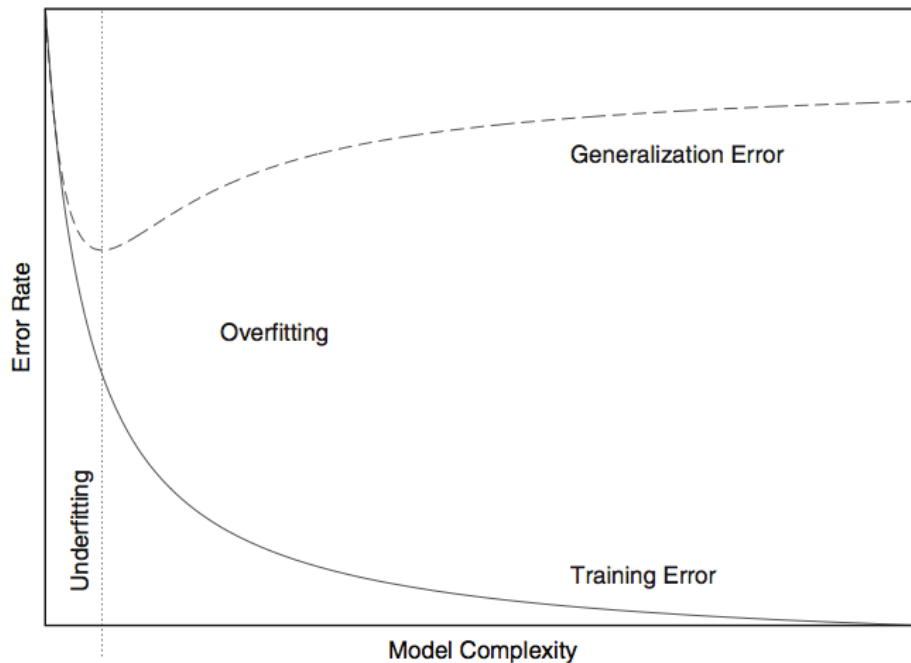
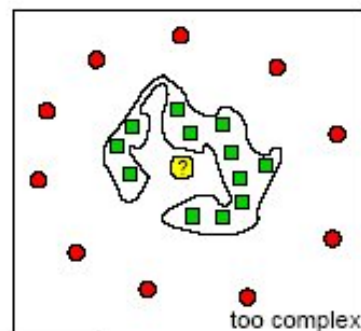
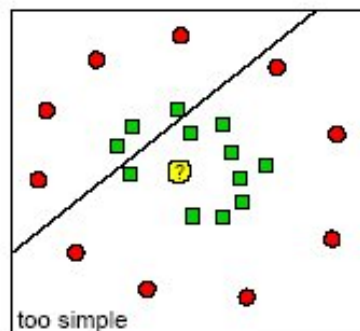
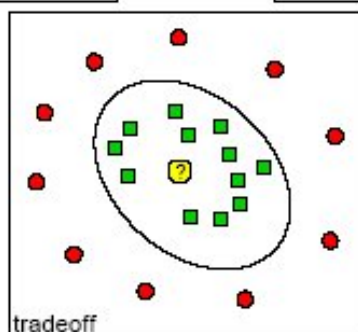


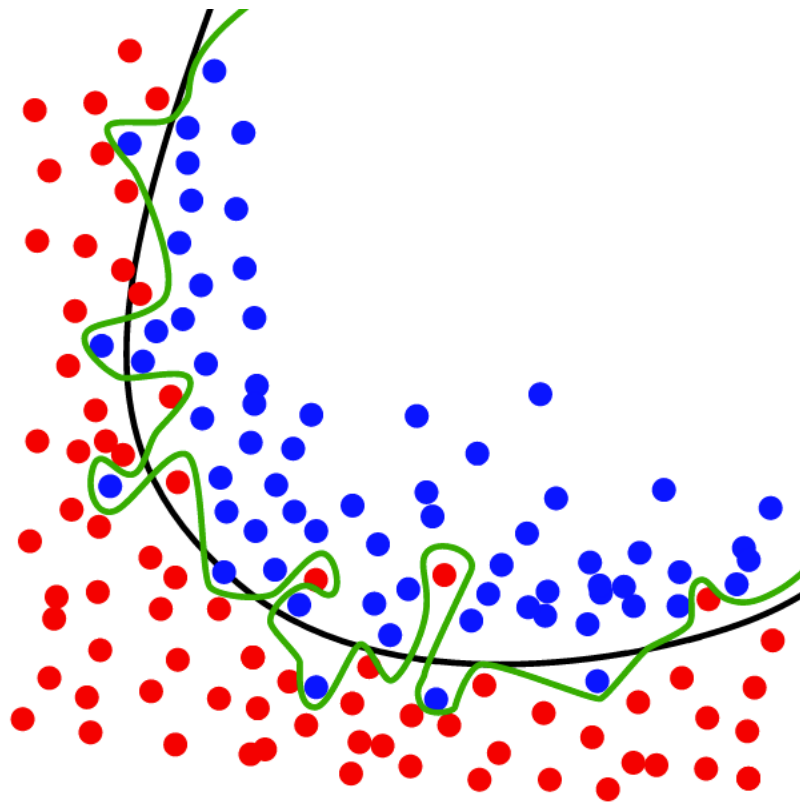
FIGURE 18-1. *Overfitting: as a model becomes more complex, it becomes increasingly able to represent the training data. However, such a model is overfitted and will not generalize well to data that was not used during training.*

Underfitting and Overfitting



- negative example
- positive example
- ⊙ new patient





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?

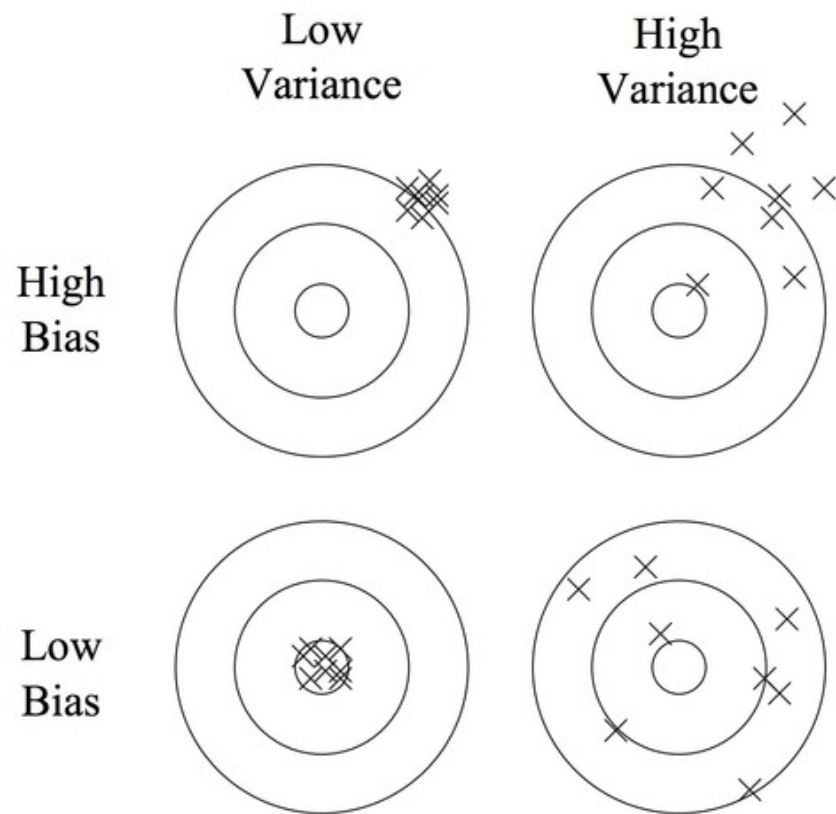
A: Of course not!

A: On its own, not very well.

NOTE

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

BIAS-VARIANCE



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.

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:

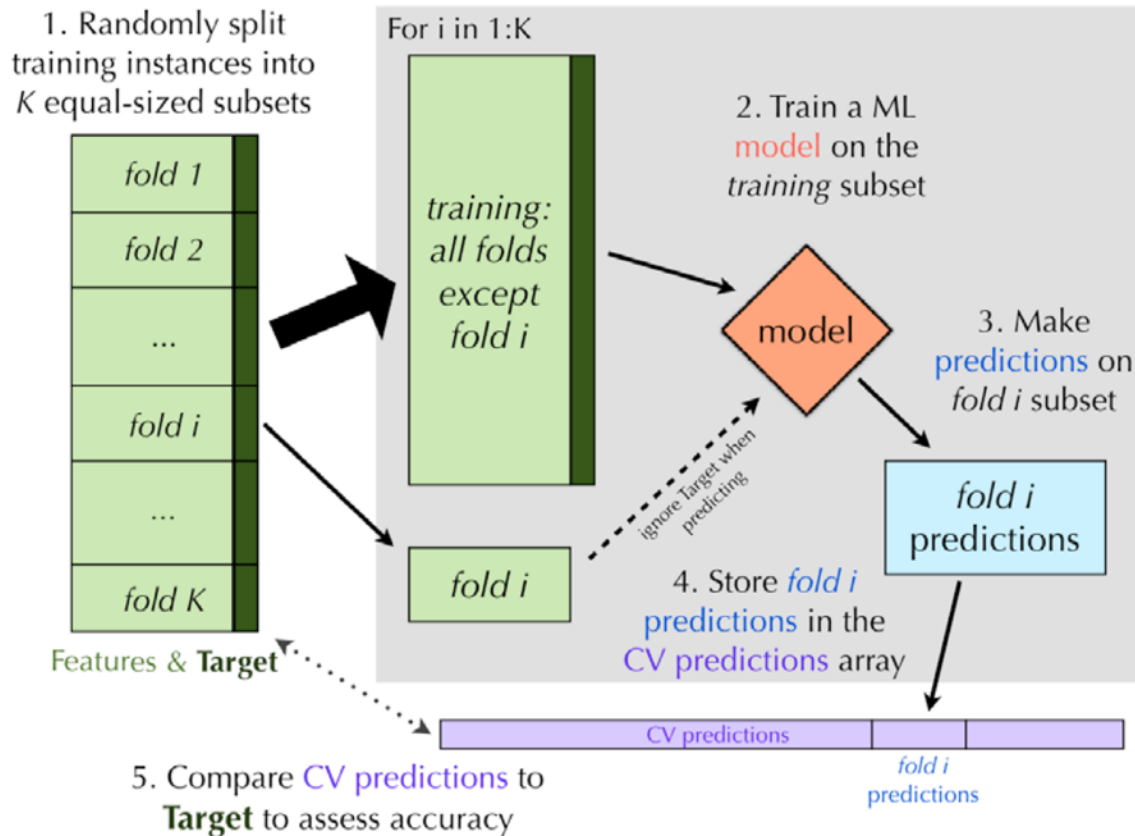
- 1) Randomly split the dataset into n equal partitions.*
- 2) Use partition 1 as test set & union of other partitions as training set.*
- 3) Find generalization error.*
- 4) Repeat steps 2-3 using a different partition as the test set at each iteration.*

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.*

Dataset	Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	<u>Accuracy</u>
1	Test	Train	Train	Train	Train	$k_1 \%$
2	Train	Test	Train	Train	Train	$k_2 \%$
3	Train	Train	Test	Train	Train	$k_3 \%$
4	Train	Train	Train	Test	Train	$k_4 \%$
5	Train	Train	Train	Train	Test	$k_5 \%$

$$5\text{-Fold Generalization Error} = (k_1 + k_2 + k_3 + k_4 + k_5) / 5$$



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:

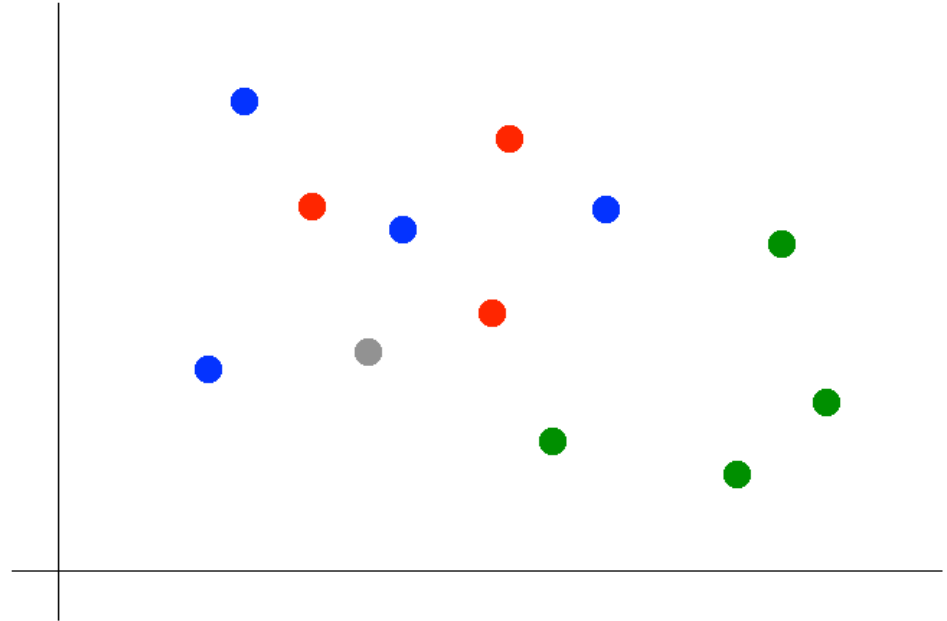
- 1) More accurate estimate of OOS prediction error.*
- 2) More efficient use of data than single train/test split.*
 - Each record in our dataset is used for both training and testing.*
- 3) Presents tradeoff between efficiency and computational expense.*
 - 10-fold CV is 10x more expensive than a single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

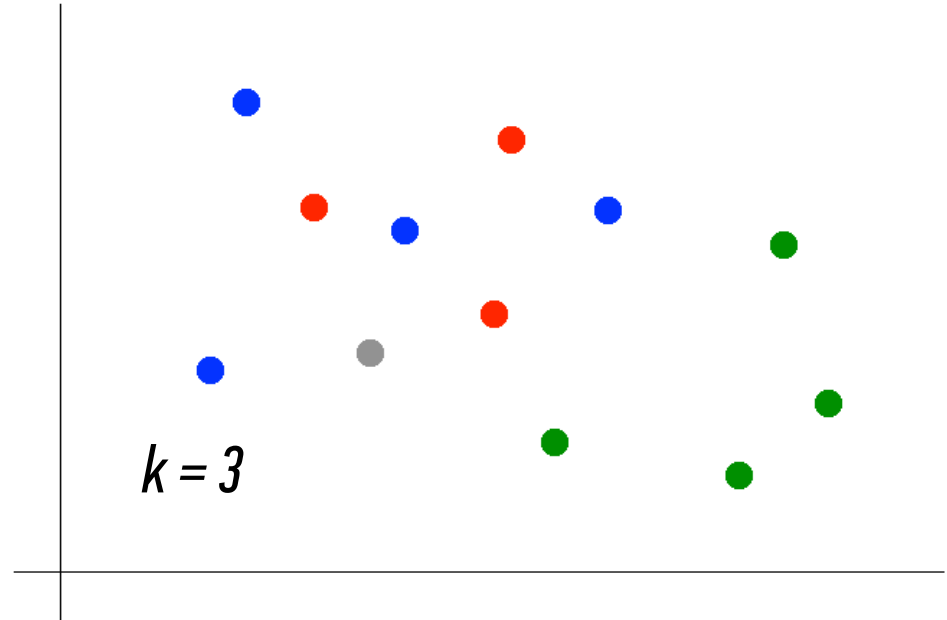
IV. KNN CLASSIFICATION

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



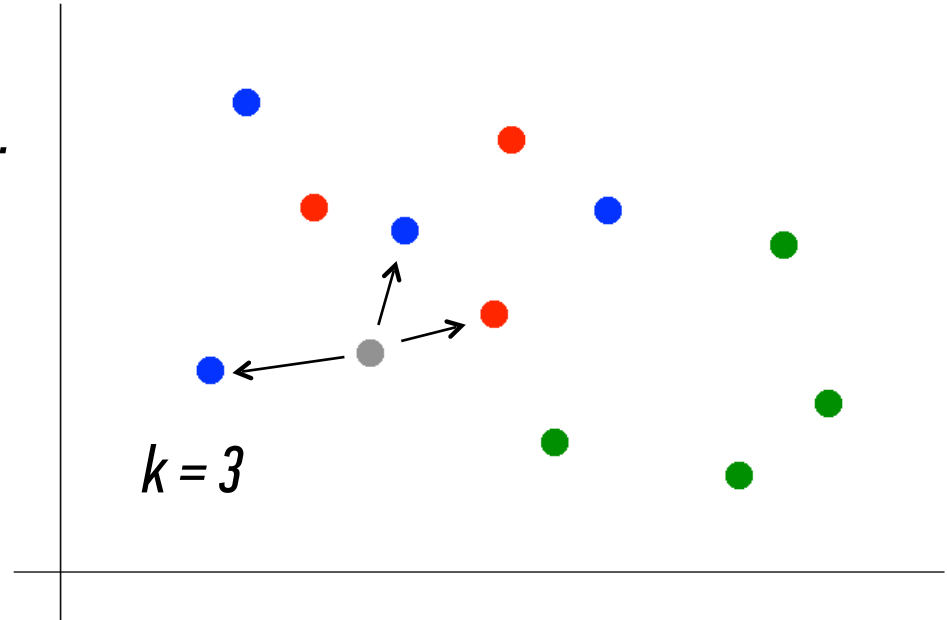
Suppose we want to predict the color of the grey dot.

1) Pick a value for k .



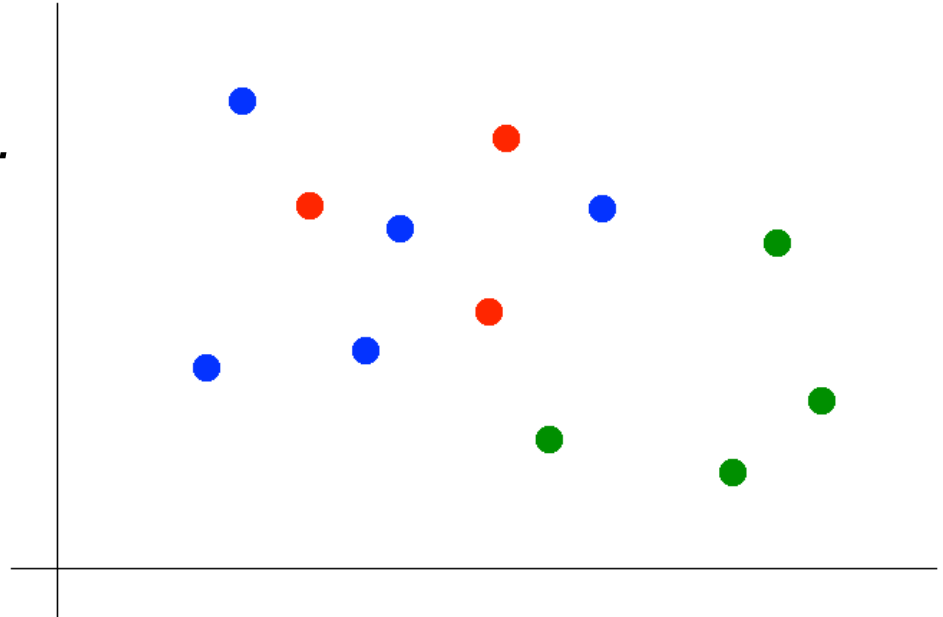
Suppose we want to predict the color of the grey dot.

- 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.*

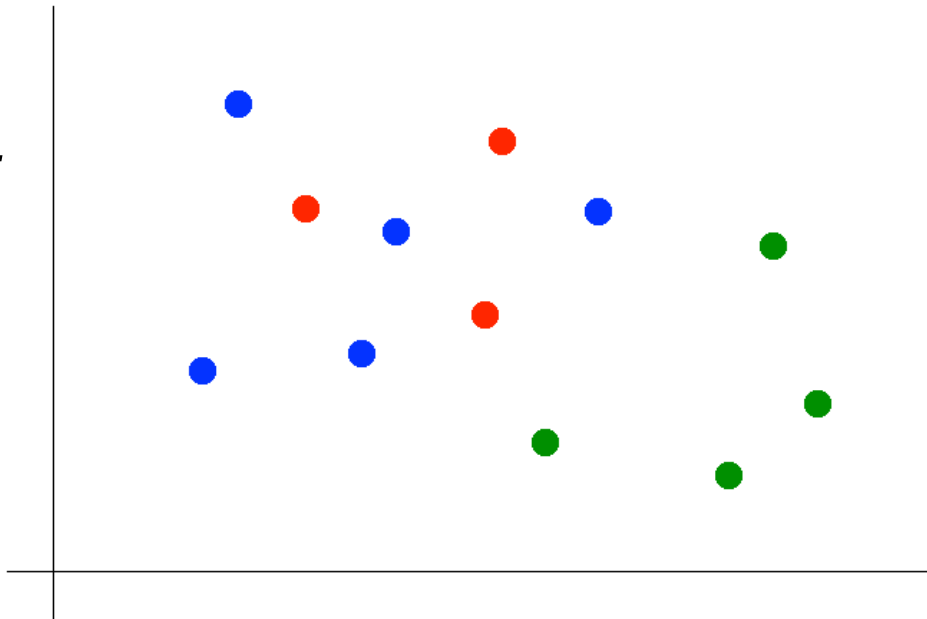


Suppose we want to predict the color of the grey dot.

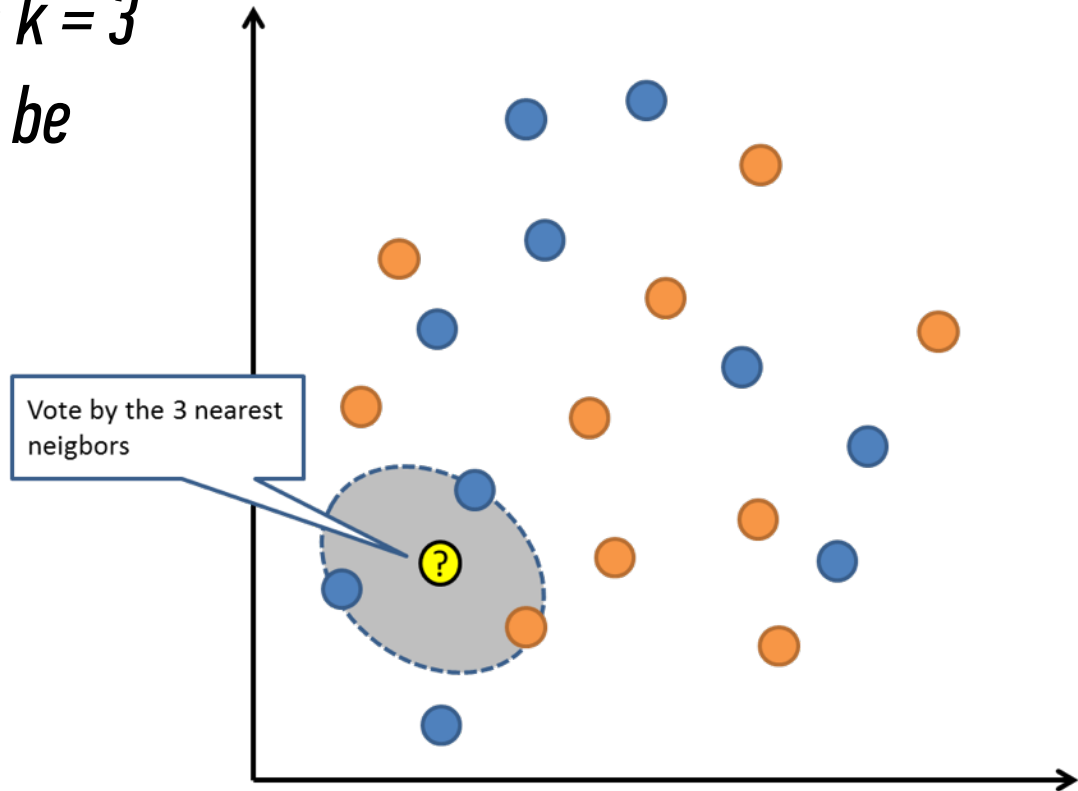
- 1) Pick a value for k .*
- 2) Find colors of k nearest neighbors.*
- 3) Assign the most common color to the grey dot.*

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?*



- *Types of machine learning problems / algorithms*
- *Generalization*
- *Types of error (training, generalization, OOS)*
- *Over-fitting and under-fitting*
- *Cross-validation*

INTRO TO DATA SCIENCE

LABS