

INTRO to DATA SCIENCE

LECTURE 5: MACHINE LEARNING

I. WHAT IS MACHINE LEARNING?

II. MACHINE LEARNING PROBLEMS

III. CLASSIFICATION PROBLEMS

IV. KNN CLASSIFICATION

I. WHAT IS MACHINE LEARNING?

WHAT IS MACHINE LEARNING?

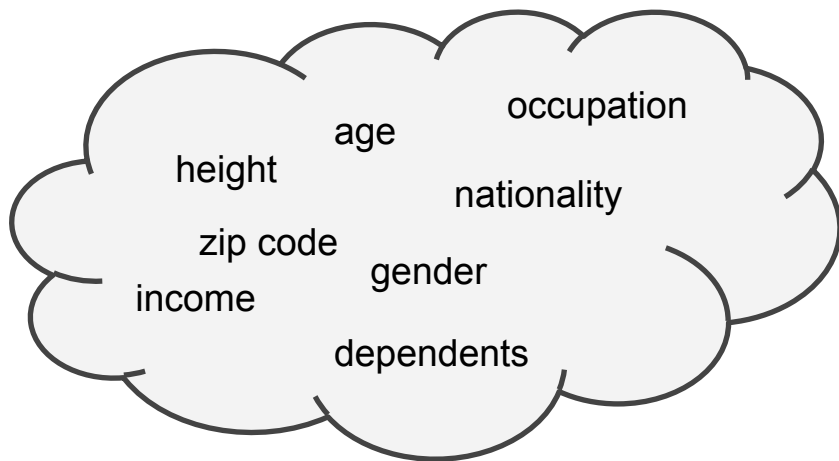
4

from Wikipedia:

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

from Wikipedia:

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

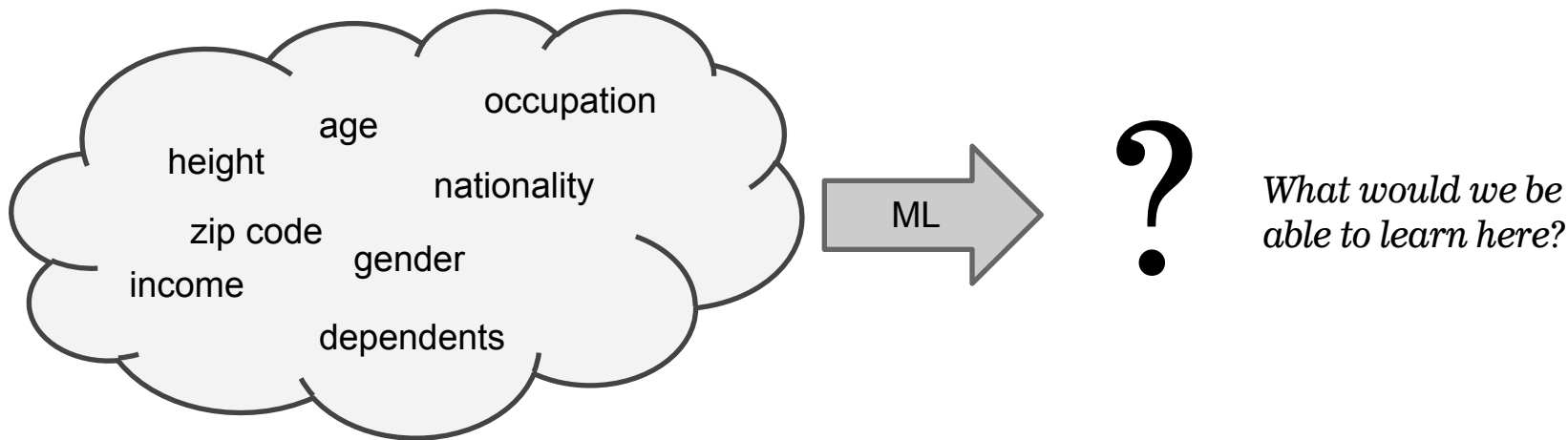


WHAT IS MACHINE LEARNING?

6

from Wikipedia:

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

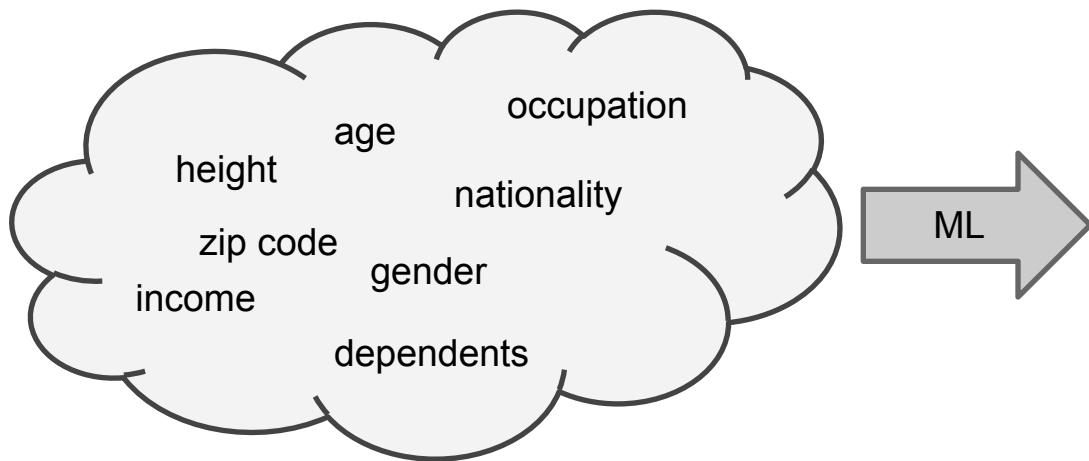


WHAT IS MACHINE LEARNING?

7

from Wikipedia:

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



Algorithm has learned to:

- predict income of person
- cluster customers in segments
- predict if person has children

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*...”

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

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

from Coursera:

“Machine learning is the science of getting computers to act without being explicitly programmed.”

from Coursera:

“Machine learning is the science of getting computers to act without being explicitly programmed.”

Instead of programming a computer how to perform a task...

from Coursera:

“Machine learning is the science of getting computers to act without being explicitly programmed.”

Instead of programming a computer how to perform a task...

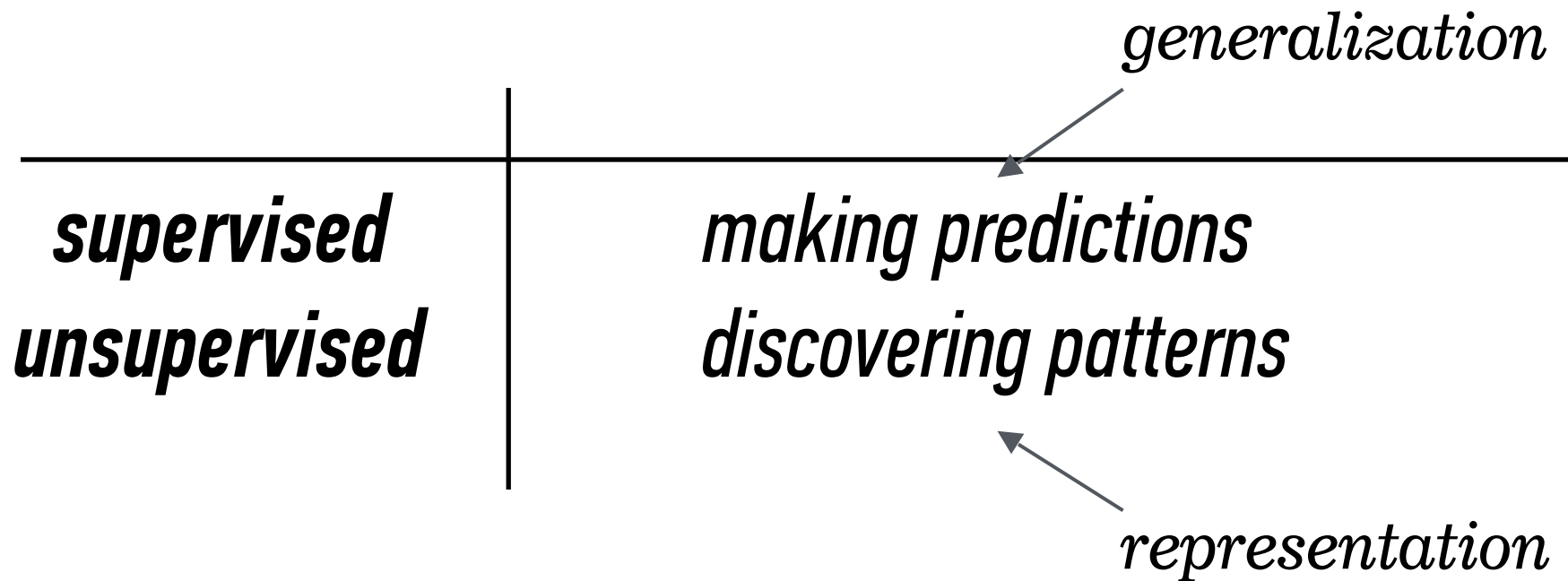
...we show the computer a history of how others performed the task,
and tell him to figure it out from there.

III. MACHINE LEARNING PROBLEMS

<i>supervised</i>	
<i>unsupervised</i>	

<i>supervised</i>	<i>labeled examples</i>
<i>unsupervised</i>	<i>no labeled examples</i>

<i>supervised</i>	<i>making predictions</i>
<i>unsupervised</i>	<i>discovering patterns</i>



continuous

categorical

--	--

continuous

categorical

quantitative

age, salary, height, etc.

qualitative

city, yes/no, vote, etc.

	<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>	Salary prediction <i>regression</i>	<i>classification</i>
<i>unsupervised</i>	<i>dimension reduction</i>	<i>clustering</i>

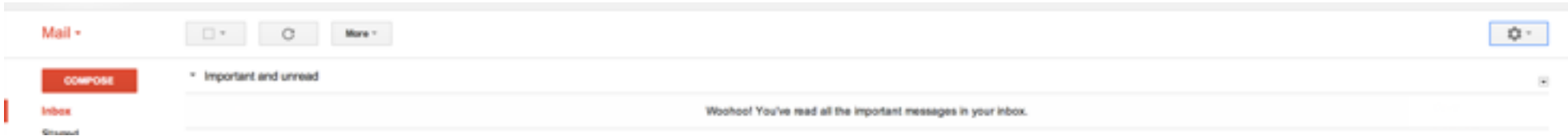
	<i>continuous</i>	<i>categorical</i>
<i>supervised</i>	Salary prediction <i>regression</i>	Vote prediction <i>classification</i>
<i>unsupervised</i>	<i>dimension reduction</i>	<i>clustering</i>

	<i>continuous</i>	<i>categorical</i>
<i>supervised</i>	Salary prediction <i>regression</i>	Vote prediction <i>classification</i>
<i>unsupervised</i>	<i>dimension reduction</i>	<i>clustering</i> customer segmentation

	<i>continuous</i>	<i>categorical</i>
<i>supervised</i>	<p>Salary prediction</p> <p><i>regression</i></p>	<p>Vote prediction</p> <p><i>classification</i></p>
<i>unsupervised</i>	<p><i>dimension reduction</i></p> <p>face recognition</p>	<p><i>clustering</i></p> <p>customer segmentation</p>

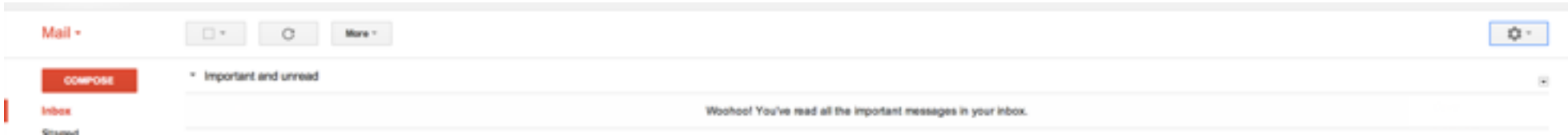
What type of problem is this?

Priority Inbox



What type of problem is this?

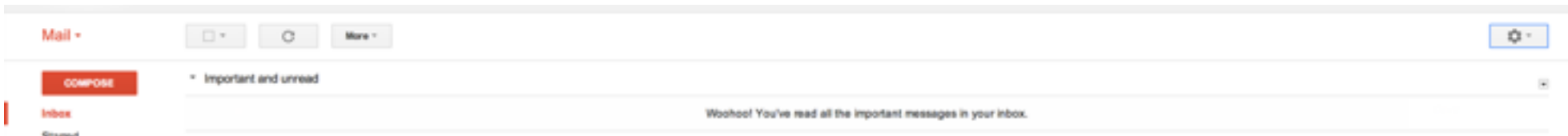
Priority Inbox



Probably either.

What type of problem is this?

Priority Inbox

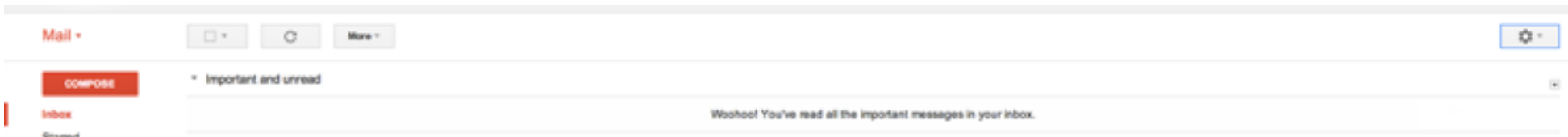


Supervised

‣ Predict which mails users are most likely to star

What type of problem is this?

Priority Inbox



Supervised

▸ Predict which mails users are most likely to star

Unsupervised

▸ Group mails into groups and decide which group represents important mails

What type of problem is this?

Music Recommendation



What type of problem is this?

Music Recommendation

Probably either.



What type of problem is this?

**Music Recommendation
as Supervised Learning**

Predict which songs a user
will 'thumbs-up'



What type of problem is this?

Music Recommendation as Unsupervised Learning

Cluster songs based on attributes
and recommend songs in the same group



QUESTION

***HOW
DO YOU
DETERMINE
THE RIGHT
APPROACH?***

	<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 and the data available.

QUESTION

HOW
DO YOU
REPRESENT
YOUR
DATA?

continuous

categorical

quantitative

qualitative

	<i>continuous</i>	<i>categorical</i>
<i>color</i>	<i>quantitative</i>	<i>or qualitative ?</i>
<i>ratings</i>		

	<i>continuous</i>	<i>categorical</i>
<i>color</i>	<i>RGB-values</i>	<i>{red, blue}</i>
<i>ratings</i>	<i>1 – 10 rating</i>	<i>1-5 star rating</i>

QUESTION

***HOW
DO YOU
MEASURE
THE
QUALITY?***

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

supervised

test out your predictions

<i>supervised</i> <i>unsupervised</i>	<i>test out your predictions</i> <i>can't really (to be continued)</i>
--	---

QUESTION

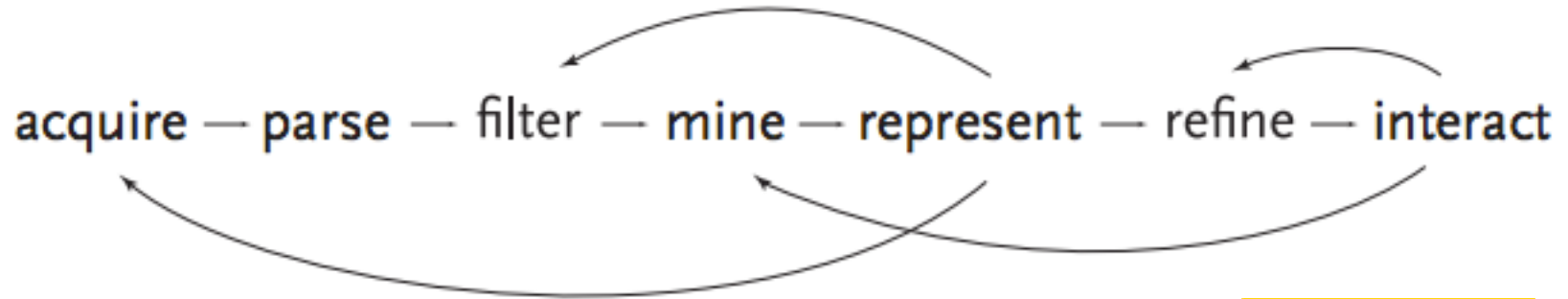
WHAT

DO YOU

DO

WITH YOUR

RESULTS?



ANSWER

Interpret them and react accordingly.

IV. CLASSIFICATION PROBLEMS

	<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 dataset (1936)

sepal_length	sepal_width	petal_length	petal_width	species
5.1	3.5	1.4	0.2	setosa
4.9	3.0	1.4	0.2	setosa
7.0	3.2	4.7	1.4	versicolor
6.4	3.2	4.5	1.5	versicolor
6.3	3.3	6.0	2.5	virginica
5.8	2.7	5.1	1.9	virginica

Here's (part of) an example dataset:

Fisher's iris dataset (1936)

sepal_length	sepal_width	petal_length	petal_width	species
5.1	3.5	1.4	0.2	setosa
4.9	3.0	1.4	0.2	setosa
7.0	3.2	4.7	1.4	versicolor
6.4	3.2	4.5	1.5	versicolor
6.3	3.3	6.0	2.5	virginica
5.8	2.7	5.1	1.9	virginica

independent variables
(also called *features*)

Here's (part of) an example dataset:

Fisher's iris dataset (1936)

sepal_length	sepal_width	petal_length	petal_width	species
5.1	3.5	1.4	0.2	setosa
4.9	3.0	1.4	0.2	setosa
7.0	3.2	4.7	1.4	versicolor
6.4	3.2	4.5	1.5	versicolor
6.3	3.3	6.0	2.5	virginica
5.8	2.7	5.1	1.9	virginica

independent variables
(also called *features*)

class labels
(qualitative)

Here's (part of) an example dataset:

Fisher's iris dataset (1936)

sepal_length	sepal_width	petal_length	petal_width	species
5.1	3.5	1.4	0.2	setosa
4.9	3.0	1.4	0.2	setosa
7.0	3.2	4.7	1.4	versicolor
6.4	3.2	4.5	1.5	versicolor
6.3	3.3	6.0	2.5	virginica
5.8	2.7	5.1	1.9	virginica

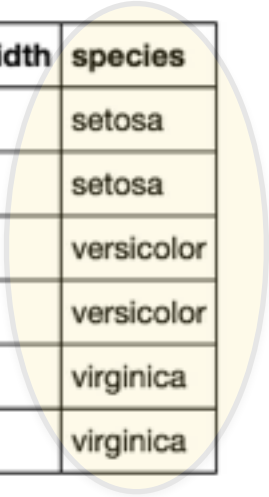
\mathbf{X} = independent variables
(also called *features*)

\mathbf{y} = class labels
(qualitative)

Q: What does “supervised” mean?

Q: What does “supervised” mean?

A: We know the labels.



sepal_length	sepal_width	petal_length	petal_width	species
5.1	3.5	1.4	0.2	setosa
4.9	3.0	1.4	0.2	setosa
7.0	3.2	4.7	1.4	versicolor
6.4	3.2	4.5	1.5	versicolor
6.3	3.3	6.0	2.5	virginica
5.8	2.7	5.1	1.9	virginica

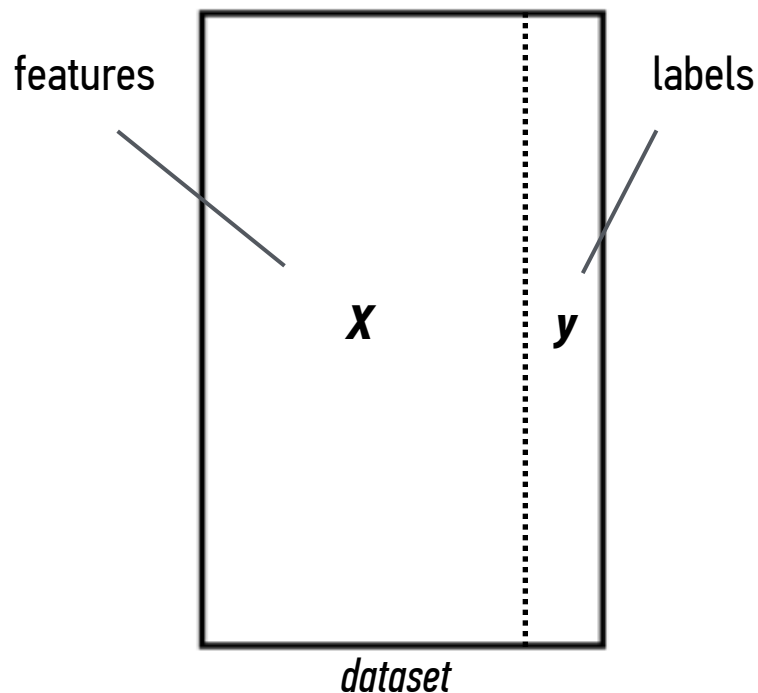
Q: How does a classification problem work?

Q: How does a classification problem work?

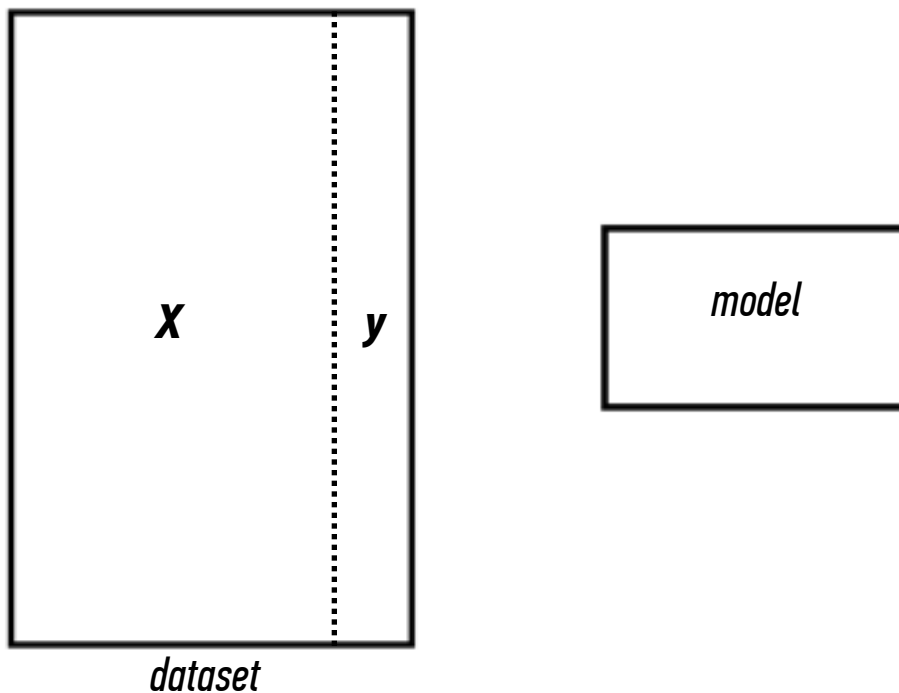


dataset

Q: How does a classification problem work?



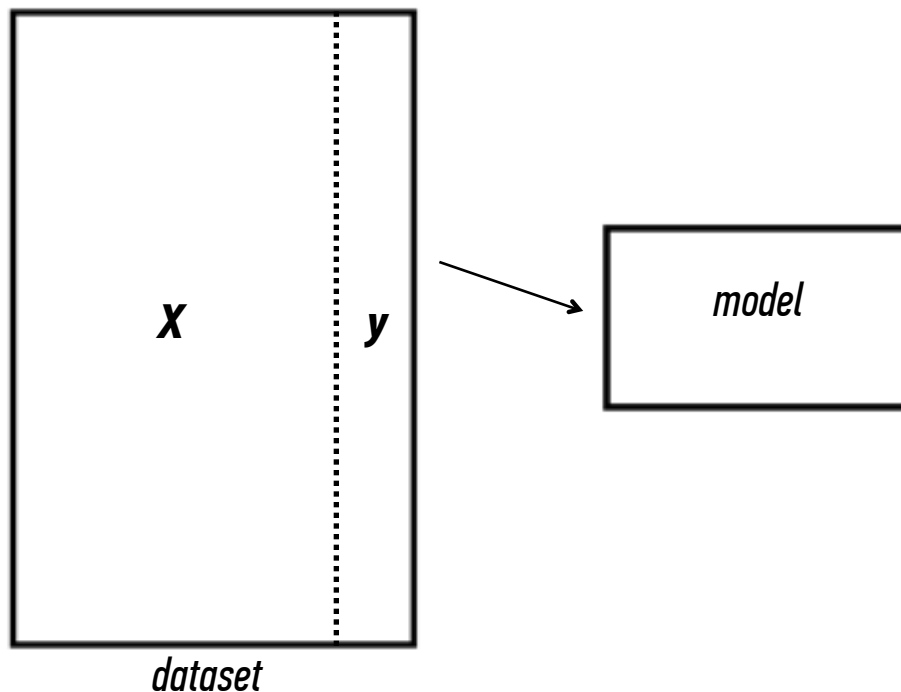
Q: How does a classification problem work?



Q: How does a classification problem work?

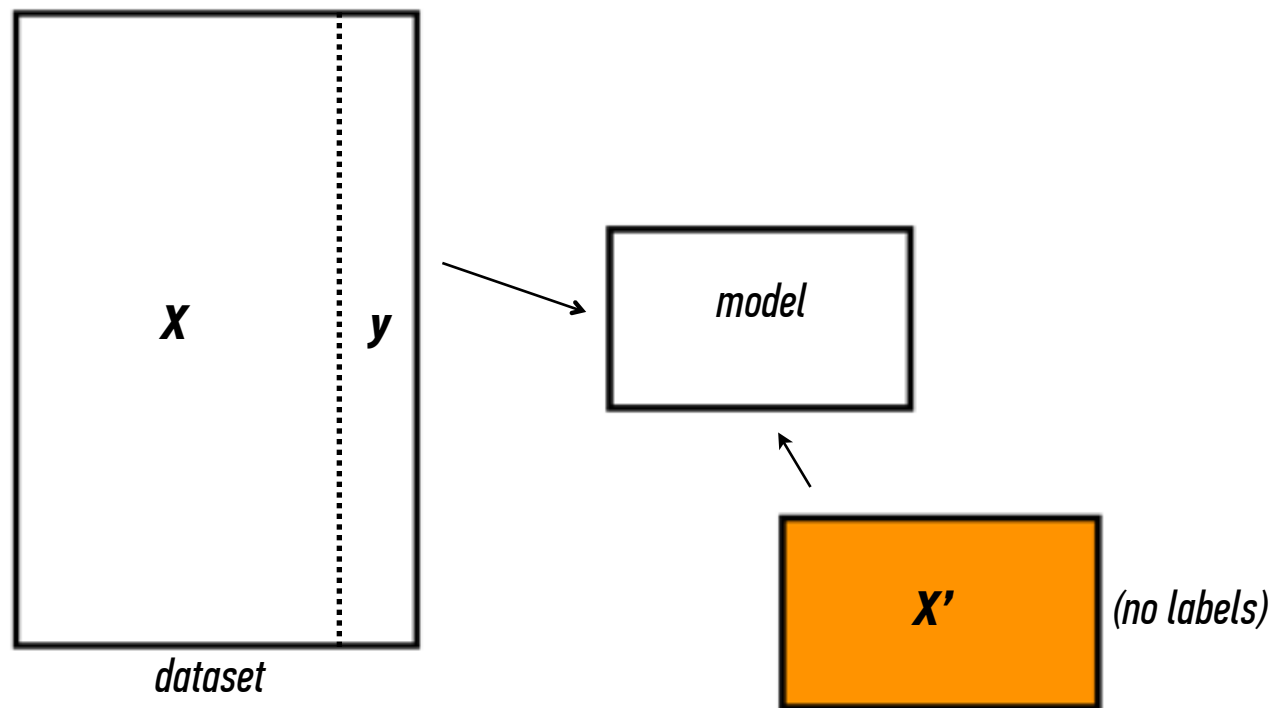
1) train model

*model 'learns' how
 \mathbf{X} and \mathbf{y} relate to
each other*



Q: How does a classification problem work?

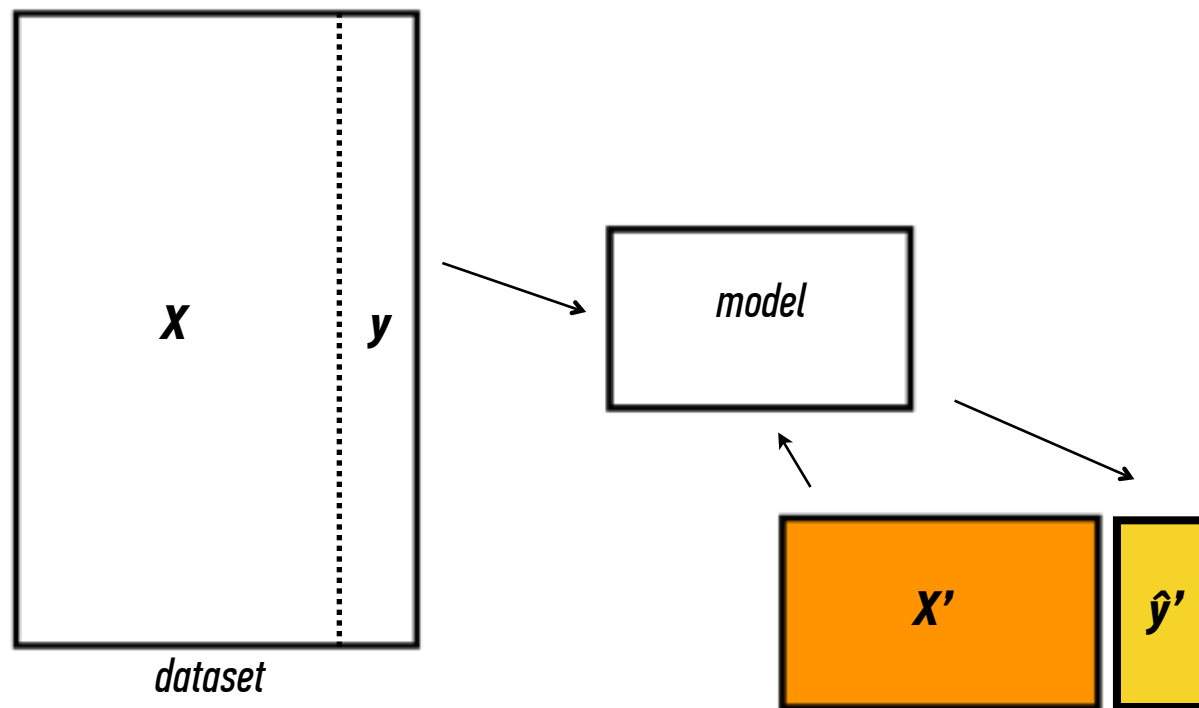
- 1) train model*
- 2) make predictions*



Q: How does a classification problem work?

- 1) *train model*
- 2) *make predictions*

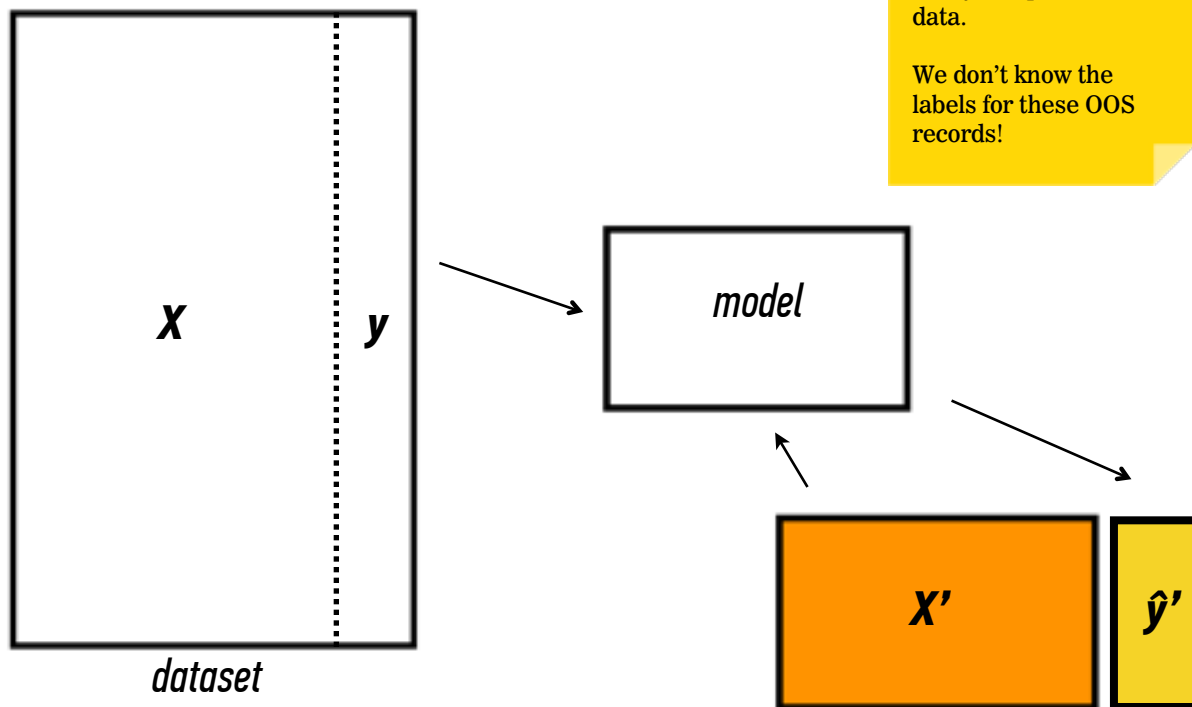
*model applies
what it learned
to new dataset X'*



Q: How does a classification problem work?

- 1) *train model*
- 2) *make predictions*

*model applies
what it learned
to new dataset X'*



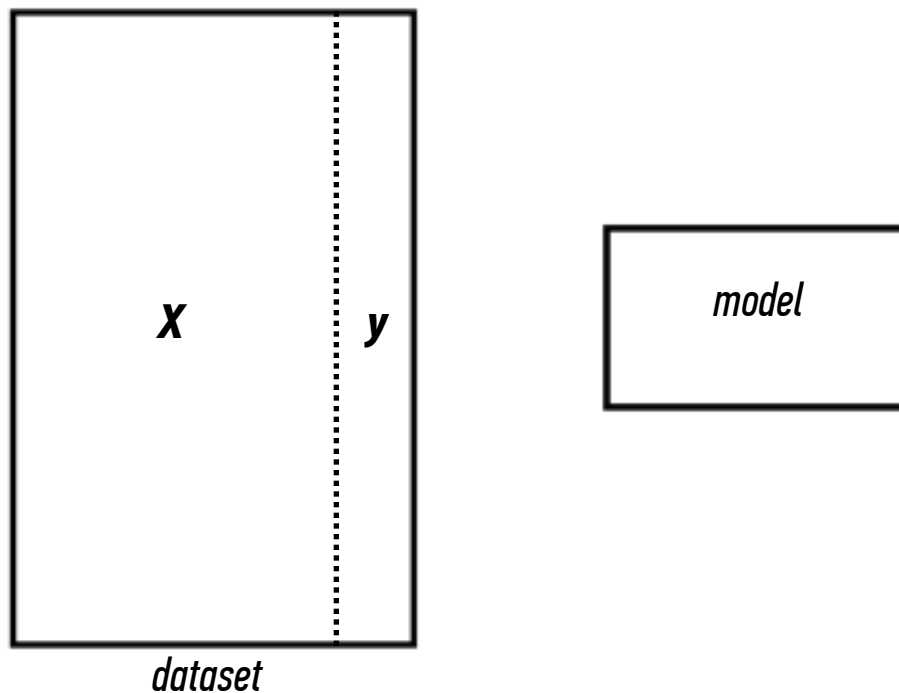
QUESTION

***HOW
DO YOU
MEASURE
THE
QUALITY?***

supervised

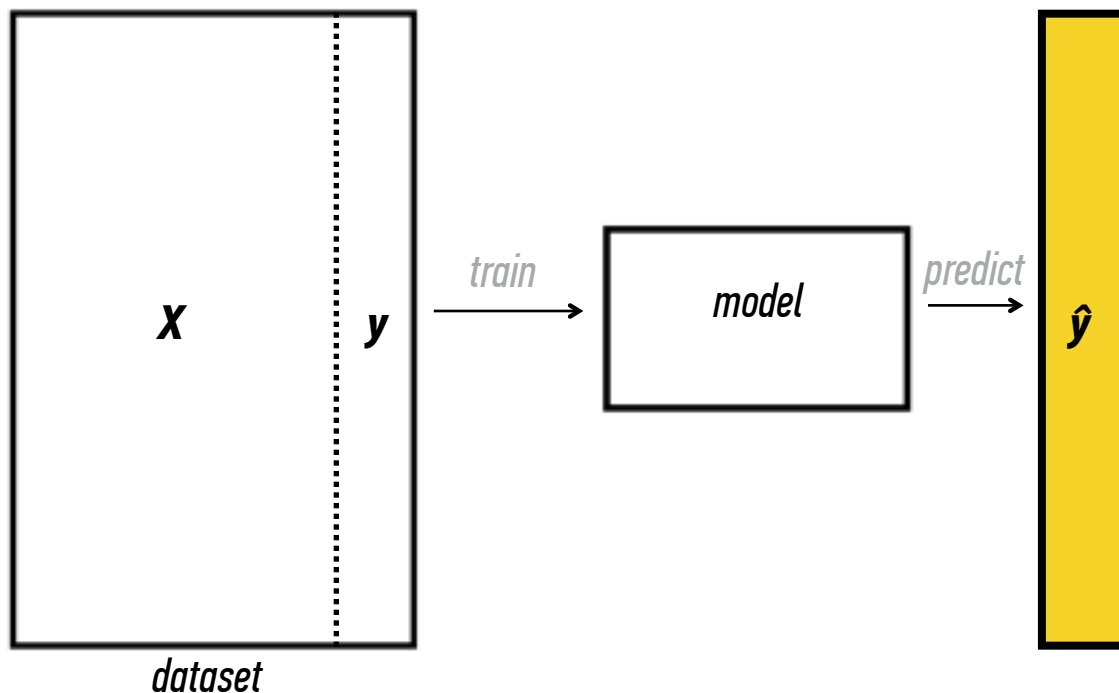
test out your predictions

Q: How do we test the model's predictions?



Q: How do we test the model's predictions?

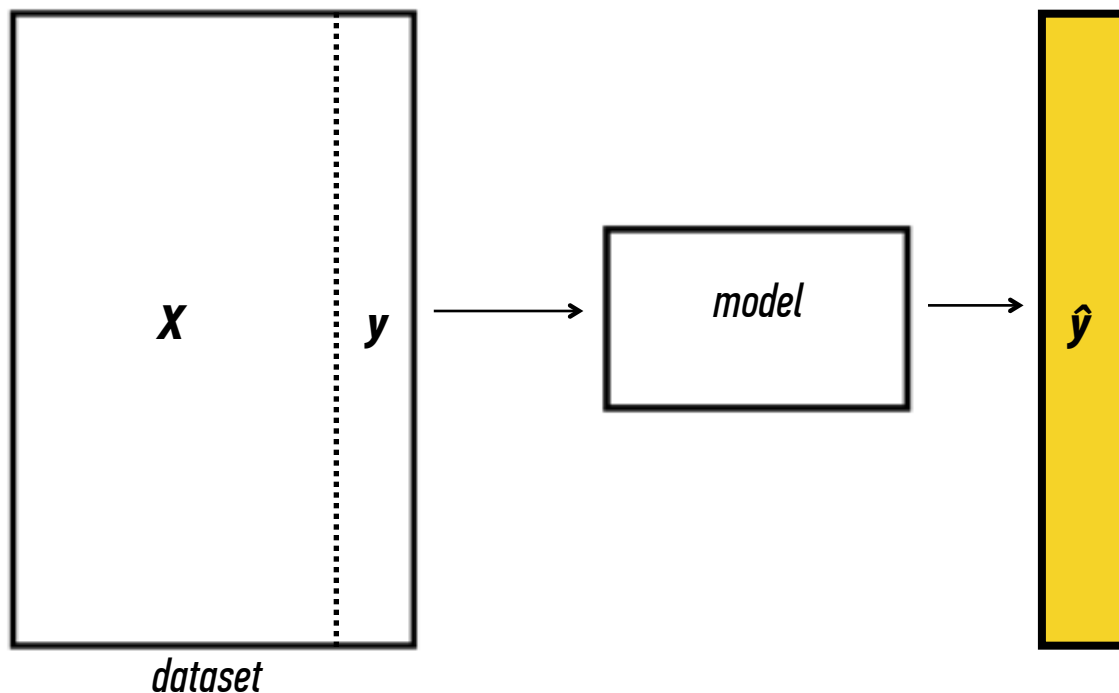
We could apply the model on the given dataset \mathbf{X} and test predictions \mathbf{y}



Q: How do we test the model's predictions?

We could apply the model on the given dataset X and test predictions y

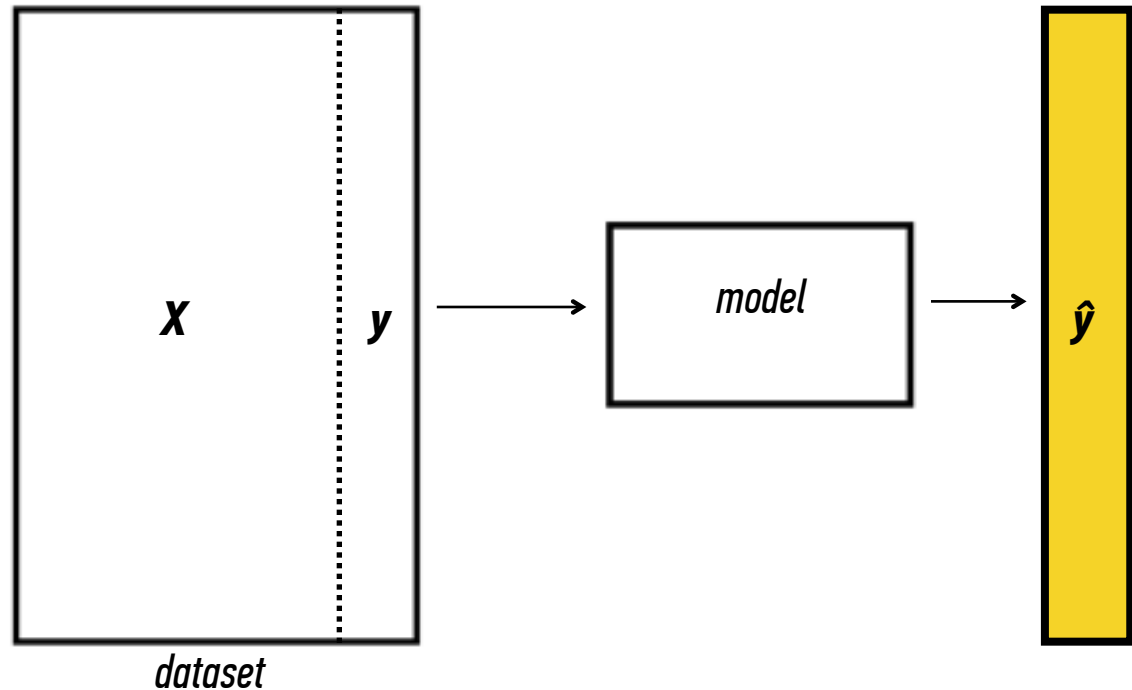
What could possibly go wrong here?



Q: How do we test the model's predictions?

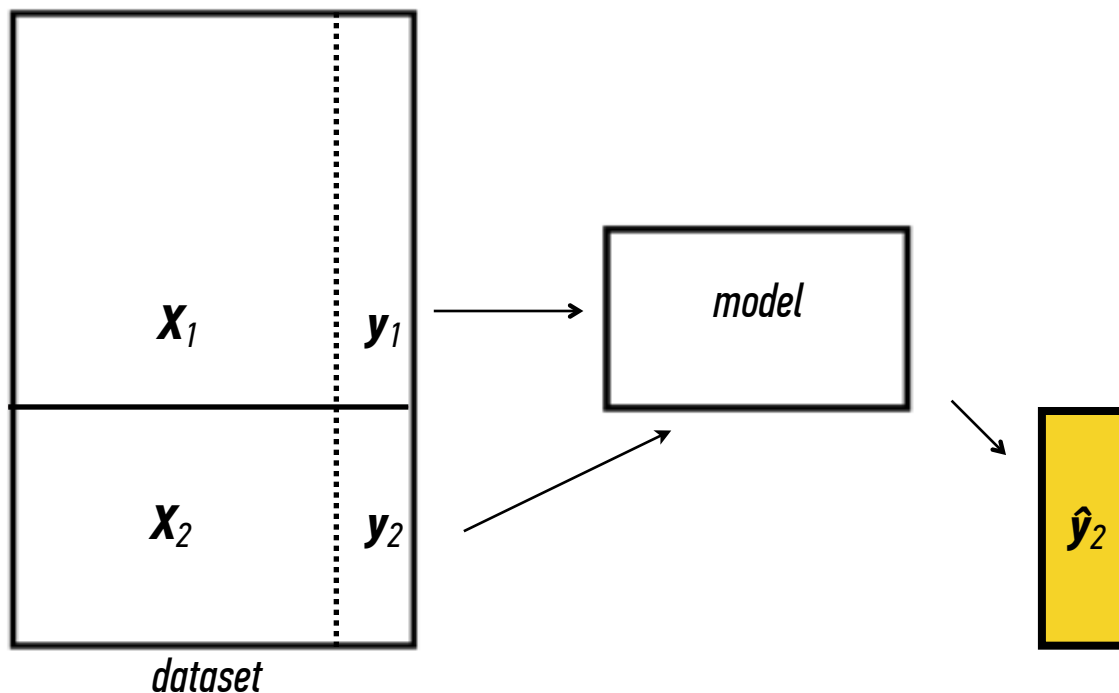
We could apply the model on the given dataset X and test predictions y

***Model could just have memorized all labels
(like a cheating student)***



Q: How do we test the model's predictions?

*Train model on a part
of \mathbf{X} , and test the results
on the rest of the data*



Q: What steps does a classification problem require?

Q: What steps does a classification problem require?



dataset

Q: What steps does a classification problem require?

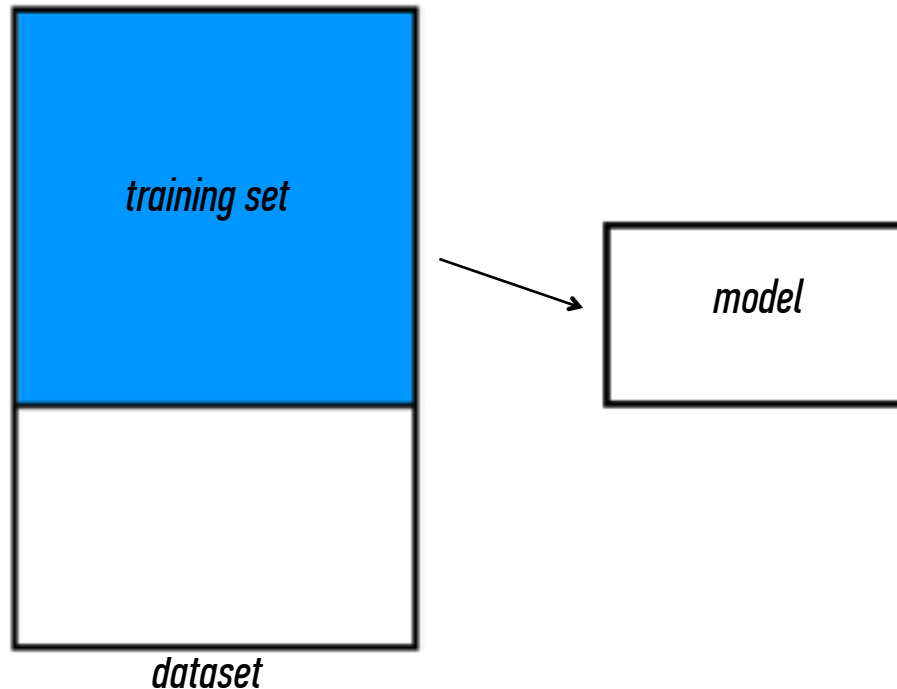
1) split dataset



dataset

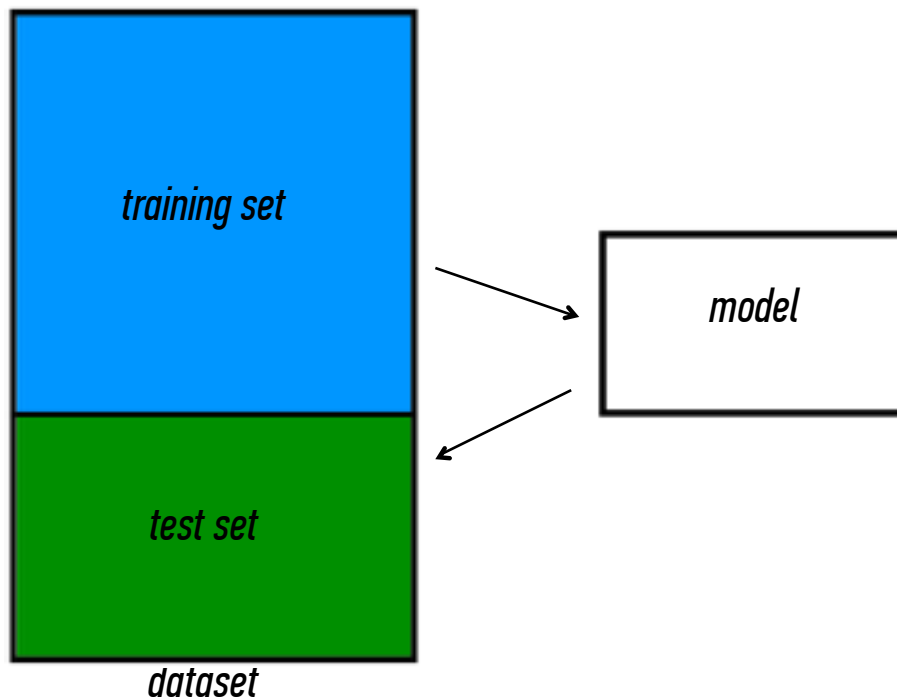
Q: What steps does a classification problem require?

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



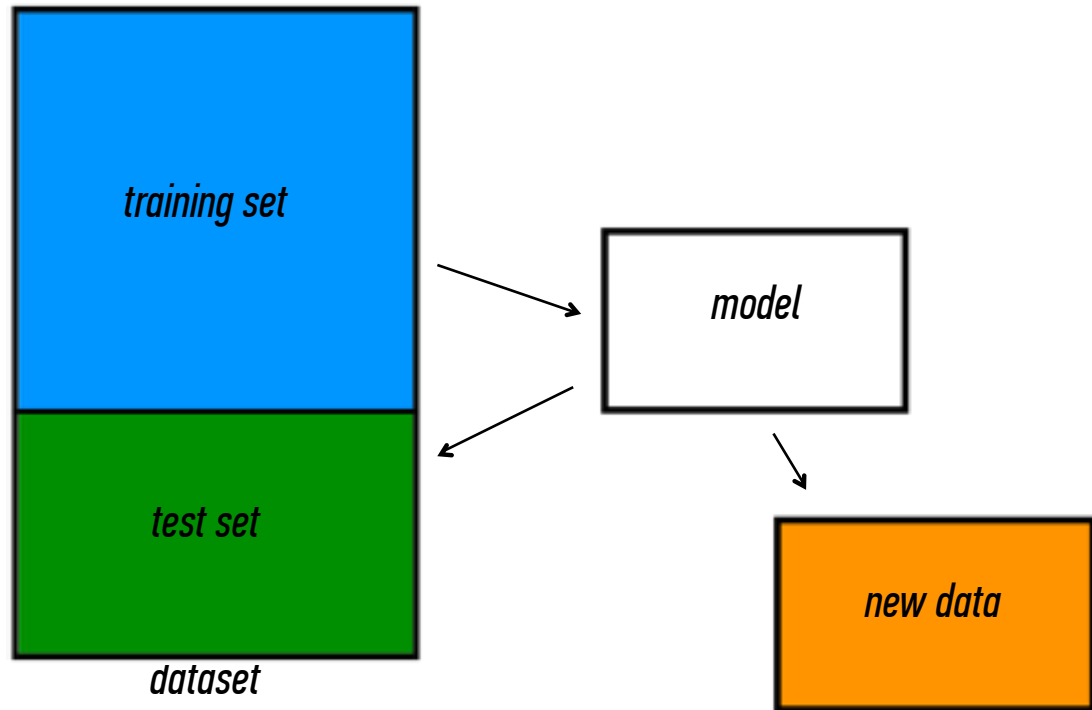
Q: What steps does a classification problem require?

- 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*



	<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>

All supervised machine learning problems require using a training and test set

INTRO TO DATA SCIENCE

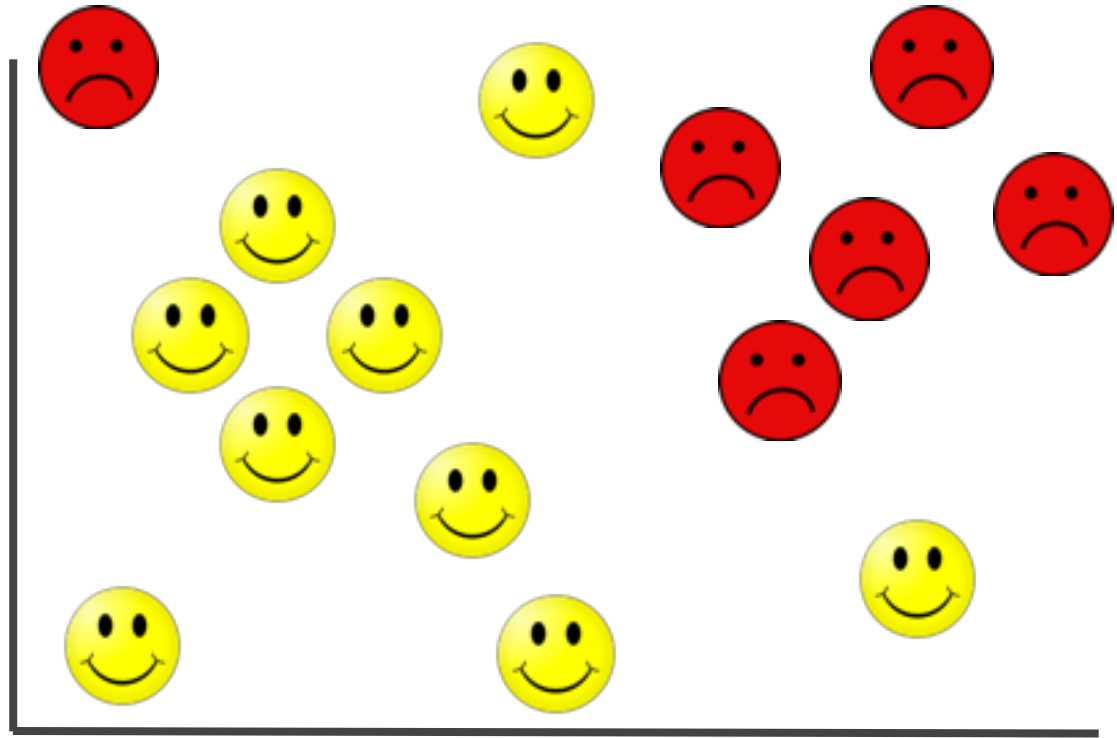
KNN CLASSIFICATION

	<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>

kNN

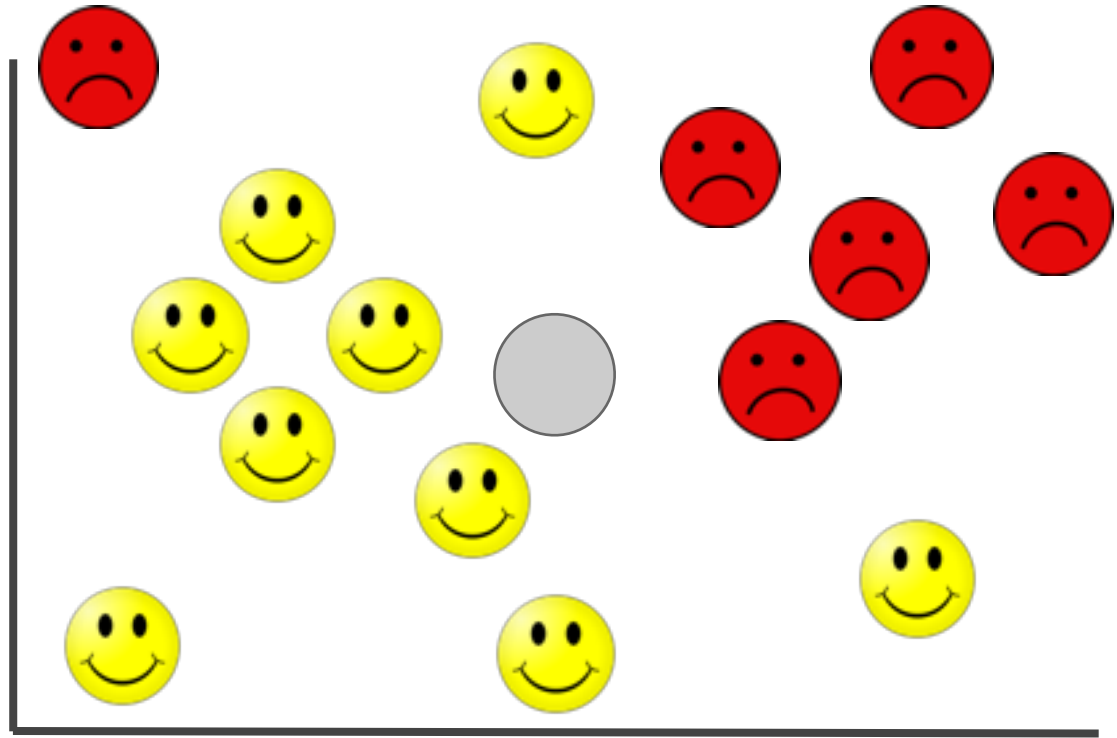
*Supervised problem
(labeled data)*

*Categorical data
(happy vs. sad)*



Want to predict:

is the grey face happy?

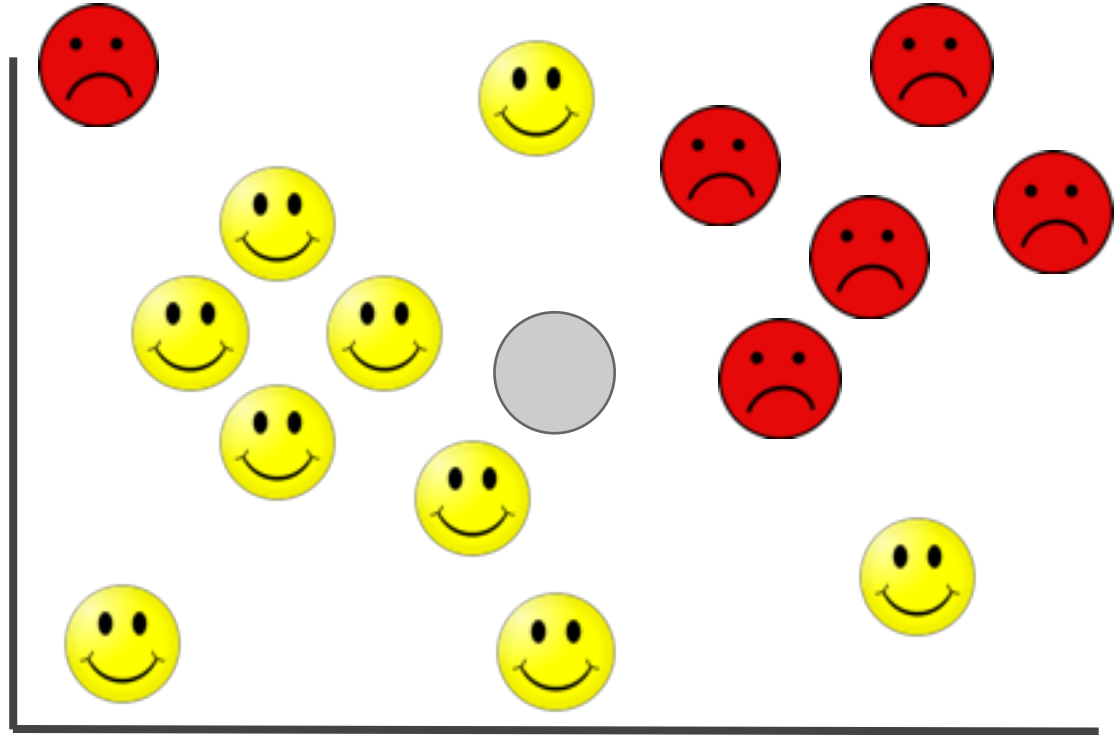


Want to predict:

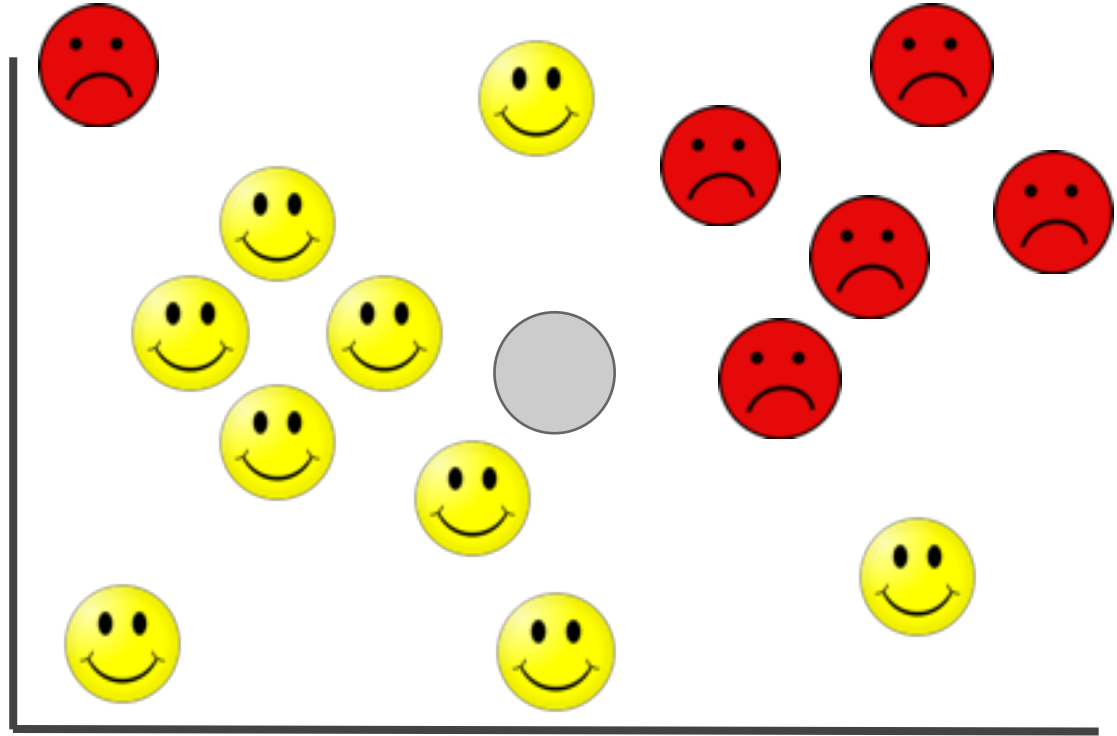
is the grey face happy?

*what do **you** think?*

why?

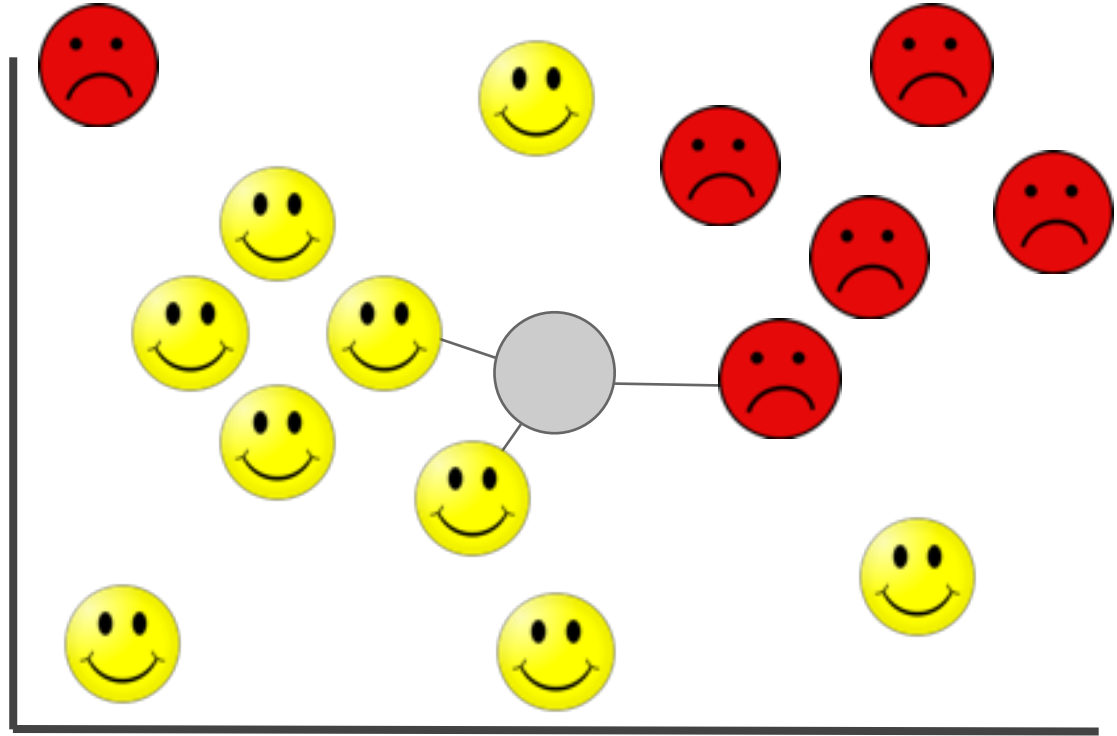


Choose k
e.g., $k = 3$



Choose k
e.g., $k = 3$

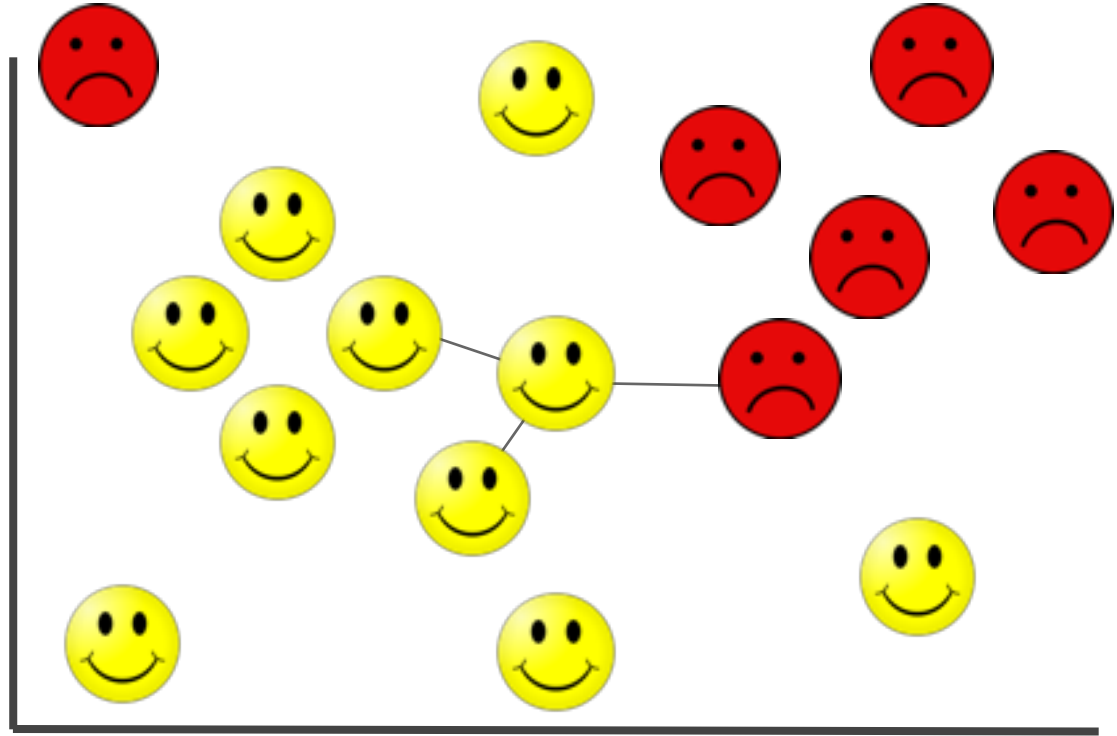
Find k nearest neighbors



Choose k
e.g., $k = 3$

Find k nearest neighbors

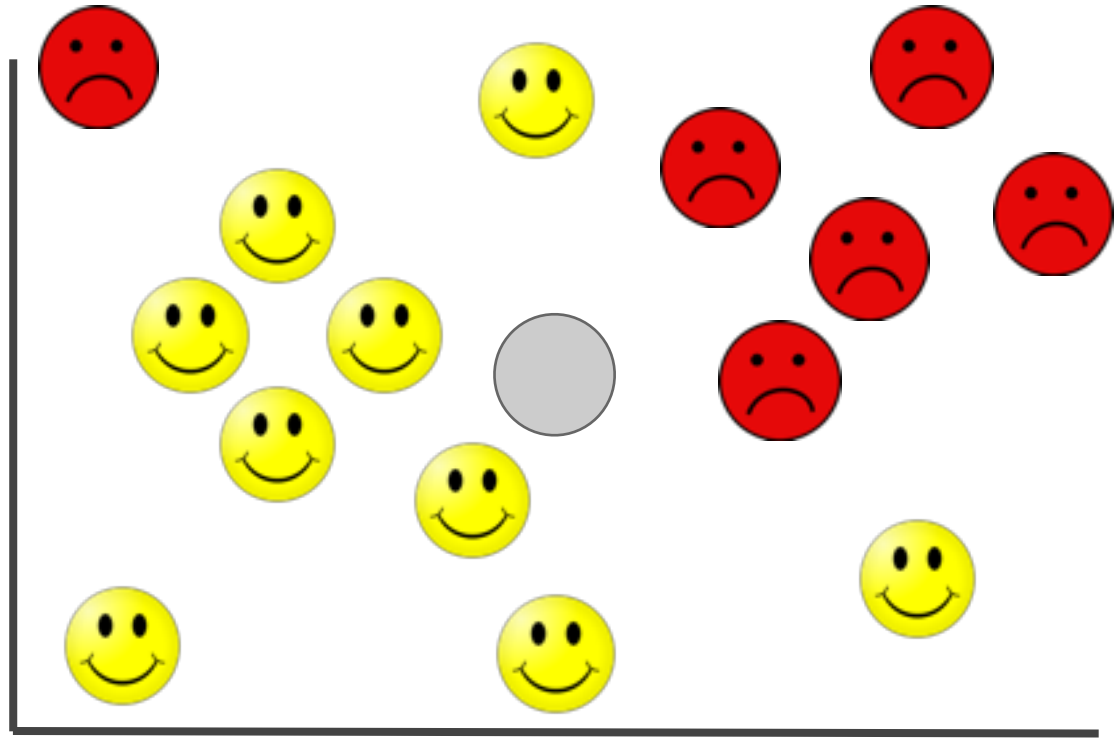
Take majority vote



QUESTION

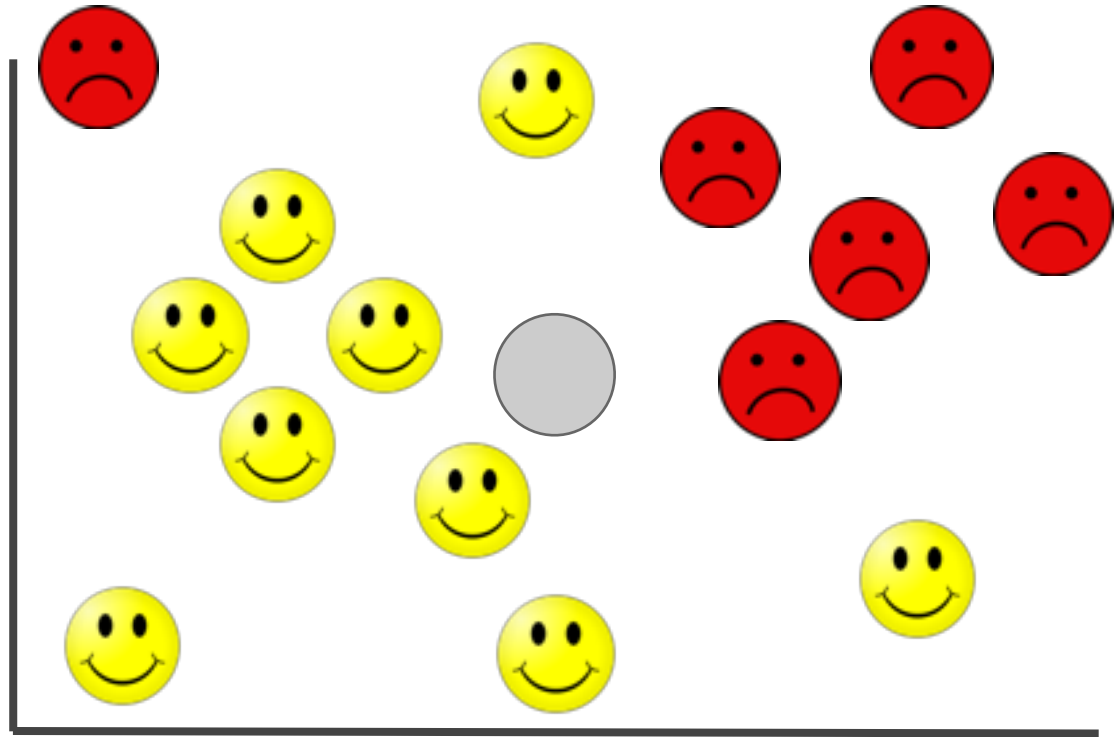
CAVEATS OF KNN

Q: What could possibly go wrong here?



Q: What could possibly go wrong here?

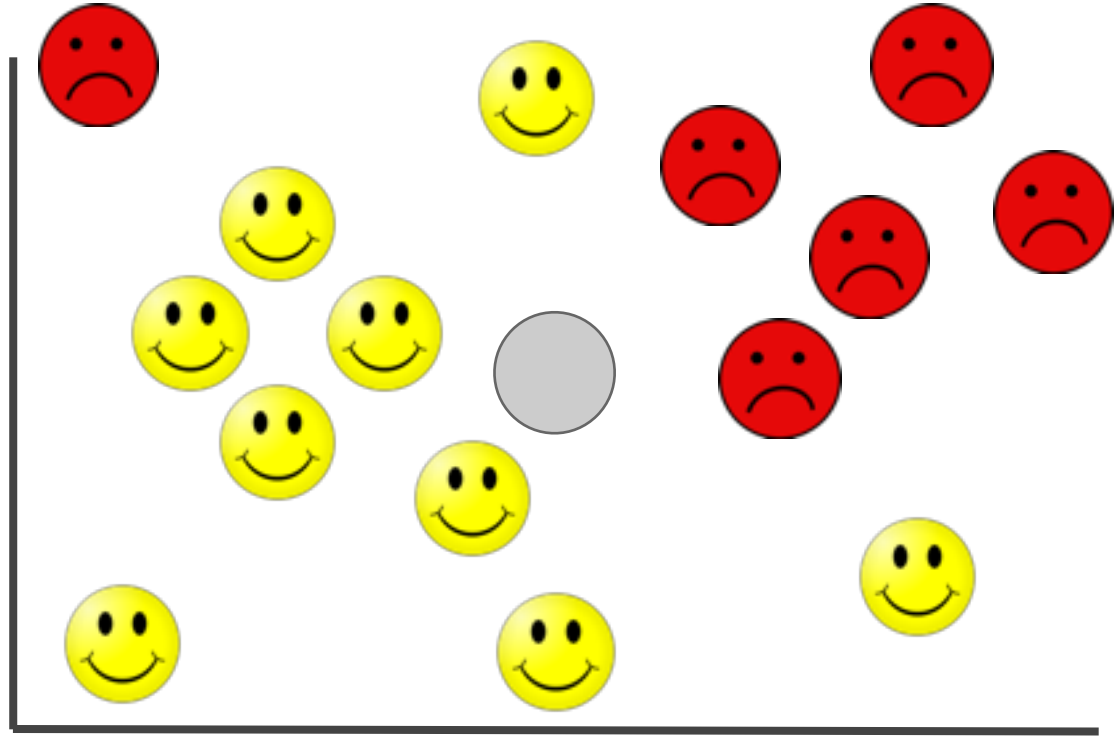
What k ?



Q: What could possibly go wrong here?

What k ?

What if $k = 1000$?

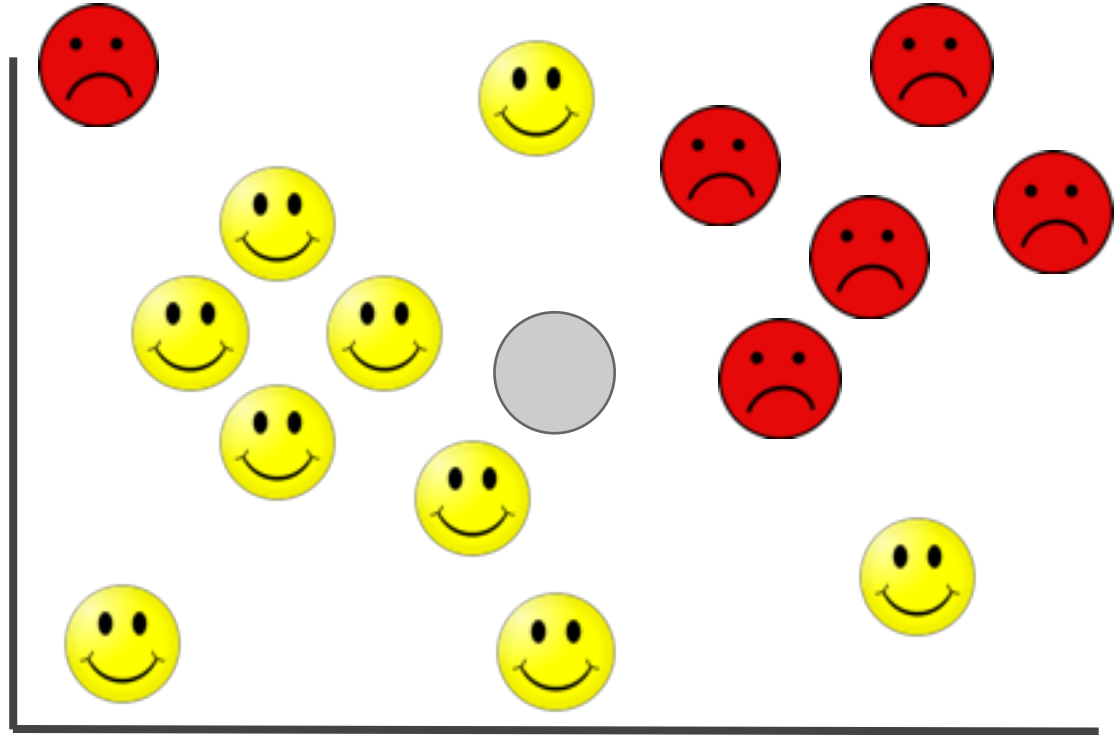


Q: What could possibly go wrong here?

What k ?

What if $k = 1000$?

What is 'nearest'?

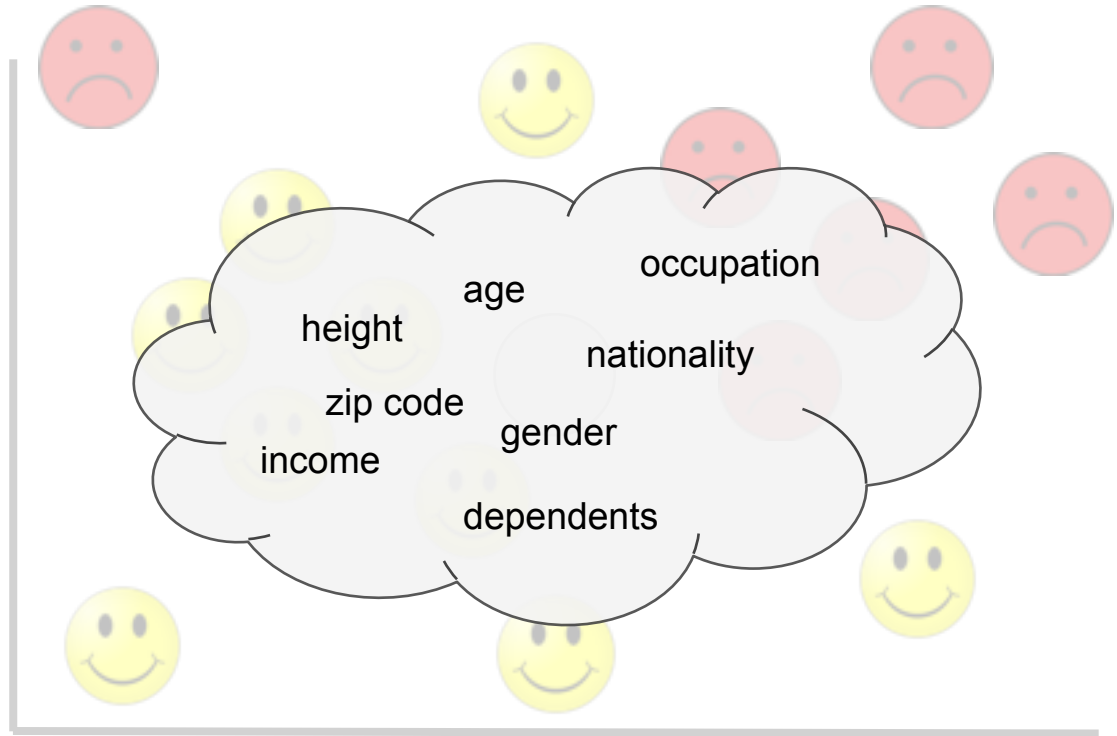


Q: What could possibly go wrong here?

What k ?

What if $k = 1000$?

What is 'nearest'?

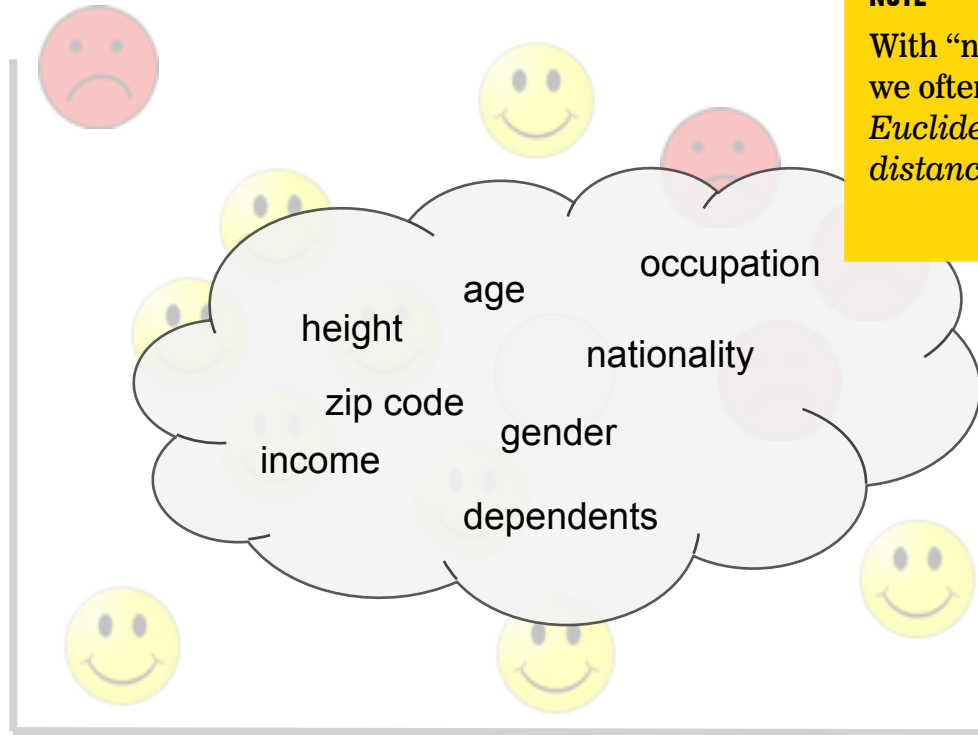


Q: What could possibly go wrong here?

What k ?

What if $k = 1000$?

What is 'nearest'?



NOTE

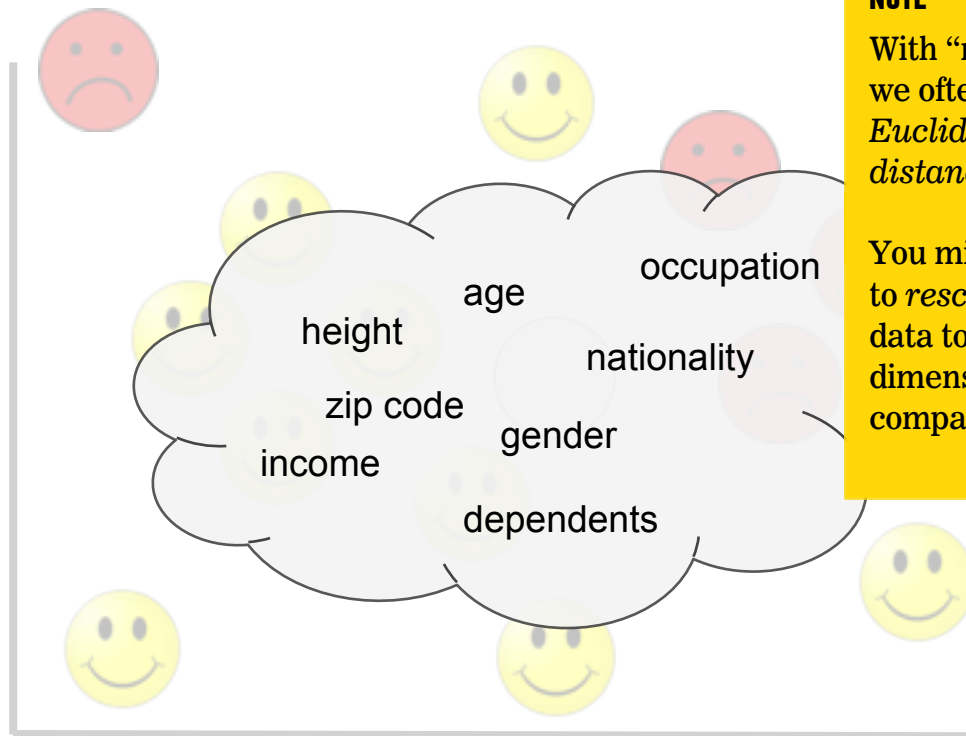
With “nearest” we often mean *Euclidean distance*

Q: What could possibly go wrong here?

What k ?

What if $k = 1000$?

What is 'nearest'?



NOTE

With “nearest” we often mean *Euclidean distance*

You might want to *rescale* your data to make the dimensions comparable

INTRO TO DATA SCIENCE

DISCUSSION