# INTRO TO DATA SCIENCE
## LECTURE 2: MACHINE LEARNING

# LAST TIME:

- FIRST LOOK AT DATA SCIENCE & THE DATA MINING WORKFLOW
- DATA EXPLORATION WITH UNIX
- DATA VISUALIZATION WITH R & GGPLOT2

# QUESTIONS?

# I. WHAT IS MACHINE LEARNING?
# II. MACHINE LEARNING PROBLEMS

# I. WHAT IS MACHINE LEARNING?

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*

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*

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

# II. MACHINE LEARNING PROBLEMS

*supervised*
*unsupervised*

| | |
|---|---|
| *supervised* | making predictions |
| *unsupervised* | discovering patterns |

| | |
|---|---|
| *supervised* | labeled examples |
| *unsupervised* | no labeled examples |

| | *continuous* | *categorical* |
|---|---|---|
| | quantitative | qualitative |

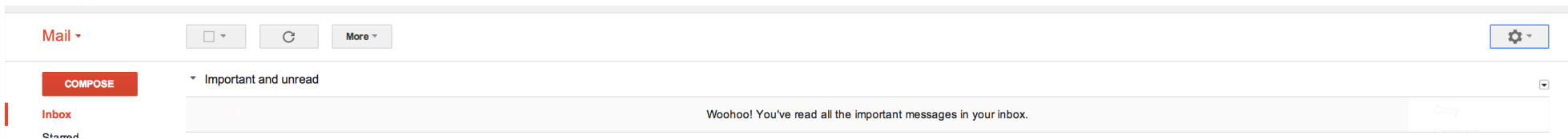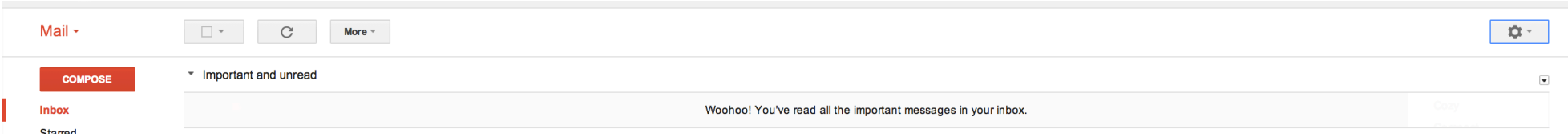|  | *continuous* | *categorical* |
| --- | --- | --- |
| *supervised* | regression | classification |
| *unsupervised* | dimension reduction | clustering |

What type of problem is this?

Priority Inbox
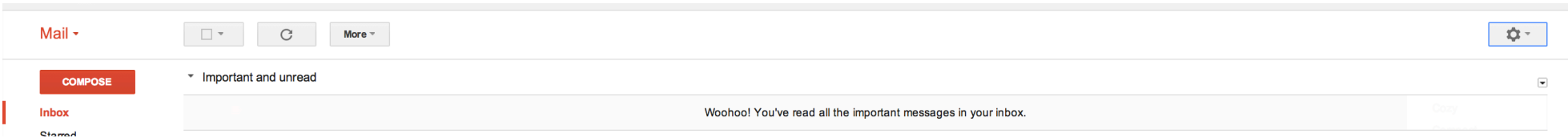
What type of problem is this?

Priority Inbox



Probably either.

**Priority Inbox: Supervised Learning**

Predict which mails users are most likely to star

**Priority Inbox: Unsupervised Learning**

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

# 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 **and the data available.**

# HOW
## DO YOU
# REPRESENT
## YOUR
# DATA?

|  | *continuous* | *categorical* |
|---|---|---|
|  | quantitative | qualitative |

|  | *continuous* | *categorical* |
|---|---|---|
| color | RGB-values | {red, blue} |
| ratings | 1 – 10 rating | 1-5 star rating |

# HOW DO YOU MEASURE OF QUALITY?

*supervised* | making predictions
*unsupervised* | extracting structure

| | |
|---|---|
| *supervised* *unsupervised* | test out your predictions … |

*supervised* | test out your predictions

# WHAT
## DO YOU
# DO
## WITH YOUR
# RESULTS?

acquire — parse — filter — mine — represent — refine — interact

**ANSWER**

Interpret them and
react accordingly.

source: http://benfry.com/phd/dissertation-110323c.pdf

# III. SUPERVISED LEARNING

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 $y$.
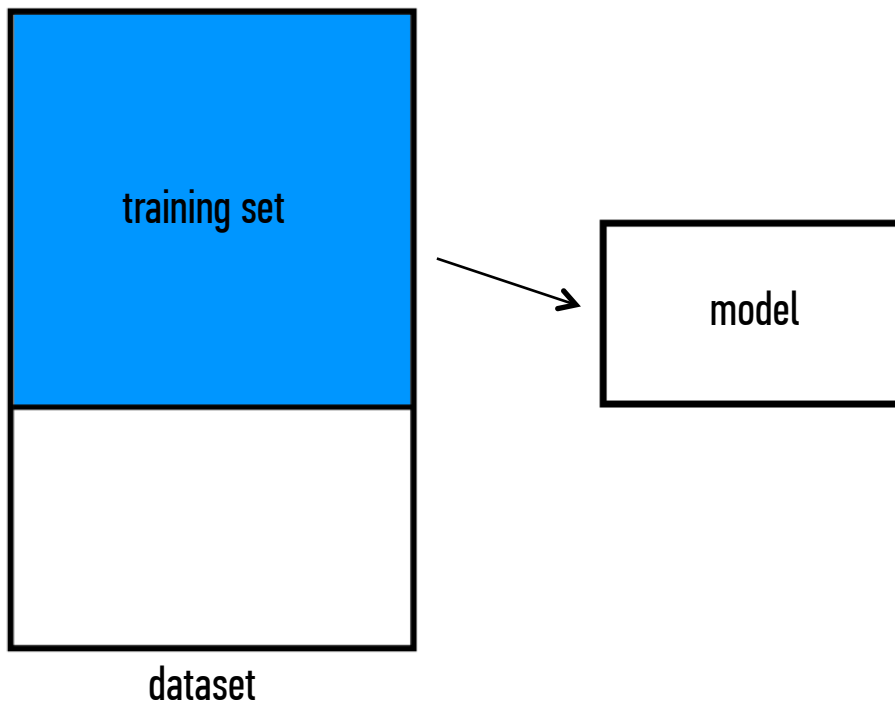
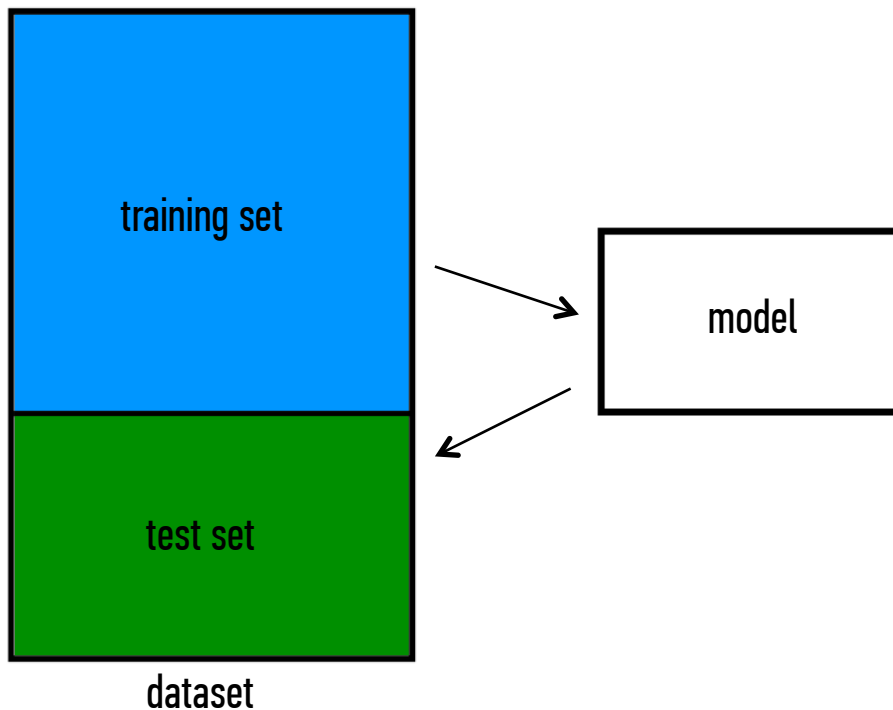Q: What steps does a classification problem require?

dataset

model

Q: What steps does a classification problem require?

1) split dataset

dataset

model

## Q: What steps does a classification problem require?

1) split dataset
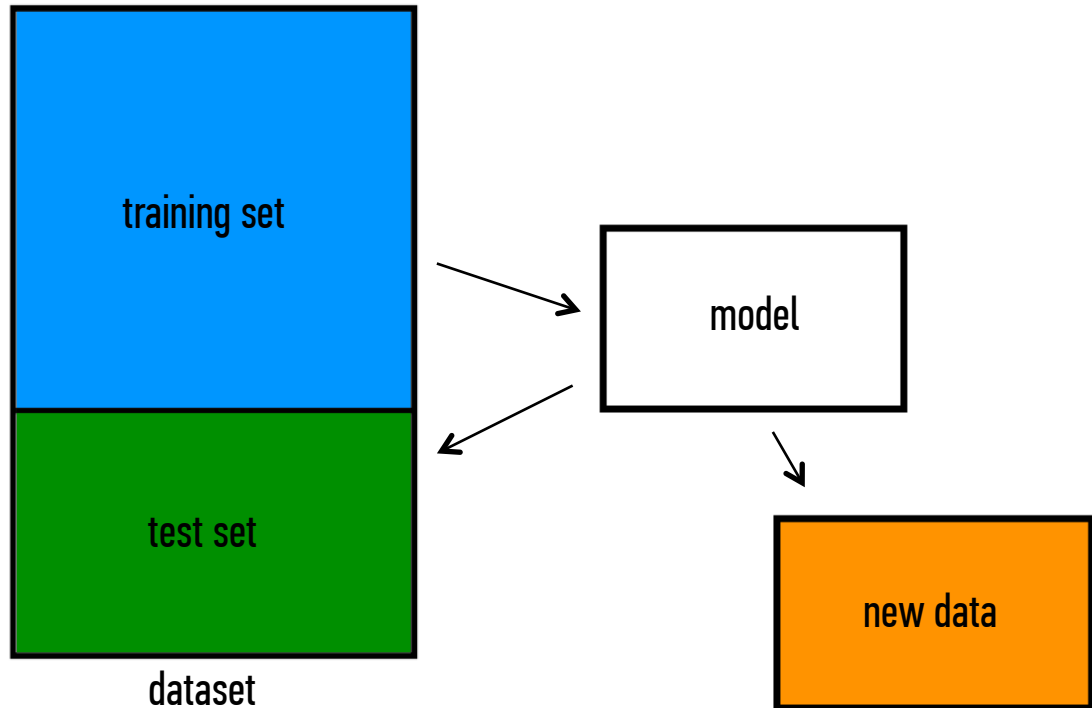2) train model

training set

model

dataset

Q: What steps does a classification problem require?

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



dataset
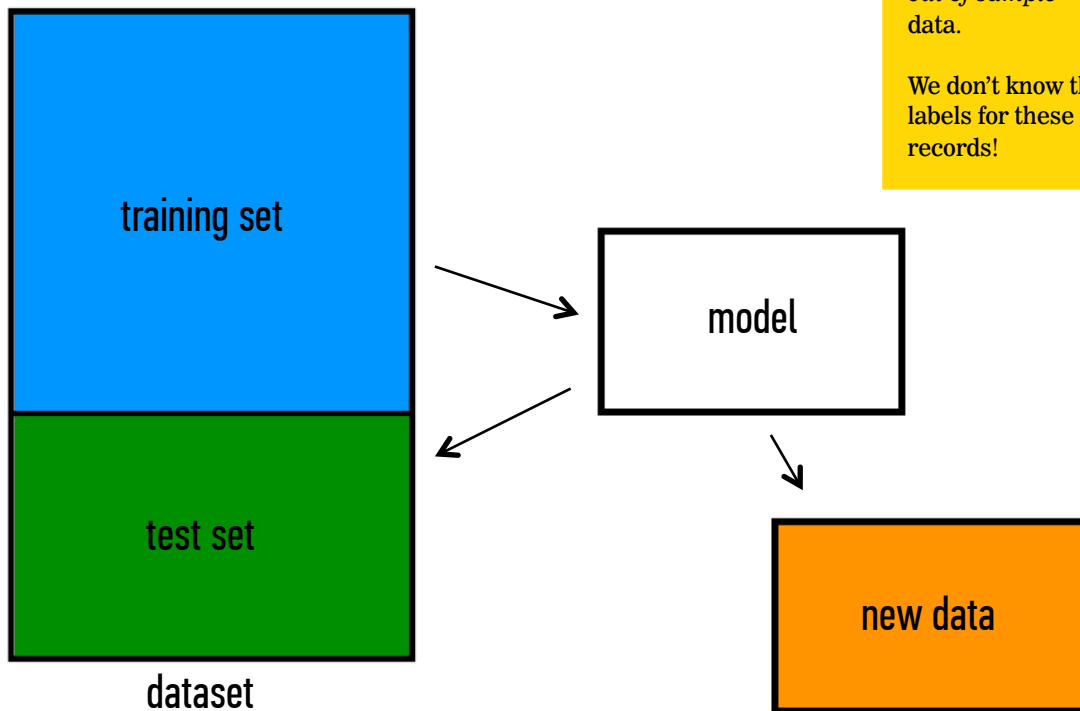
Q: What steps does a classification problem require?

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



training set

test set

dataset

model

new data

## Q: What steps does a classification problem require?

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

**NOTE**
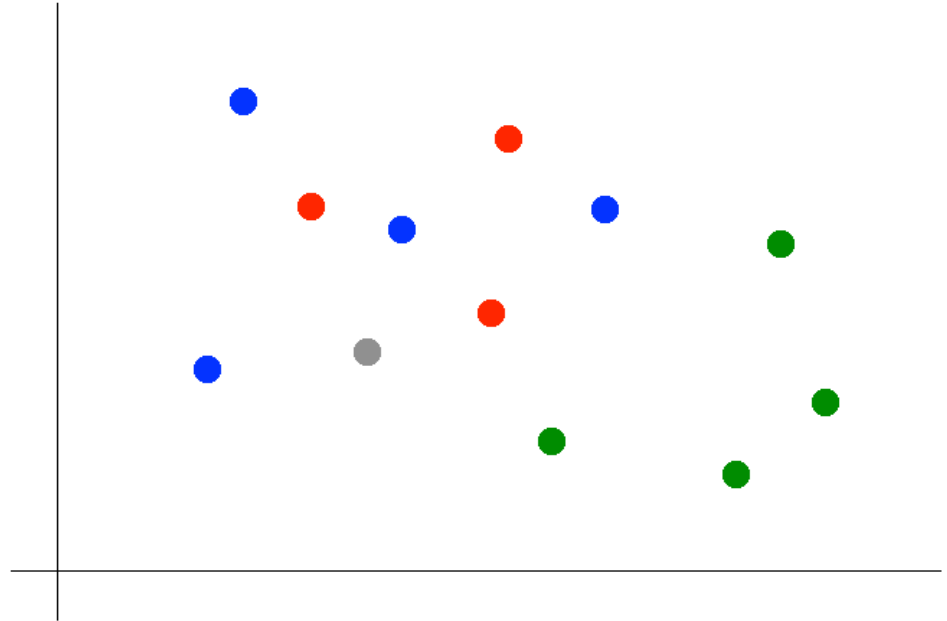This new data is called *out of sample* data.

We don't know the labels for these OOS records!
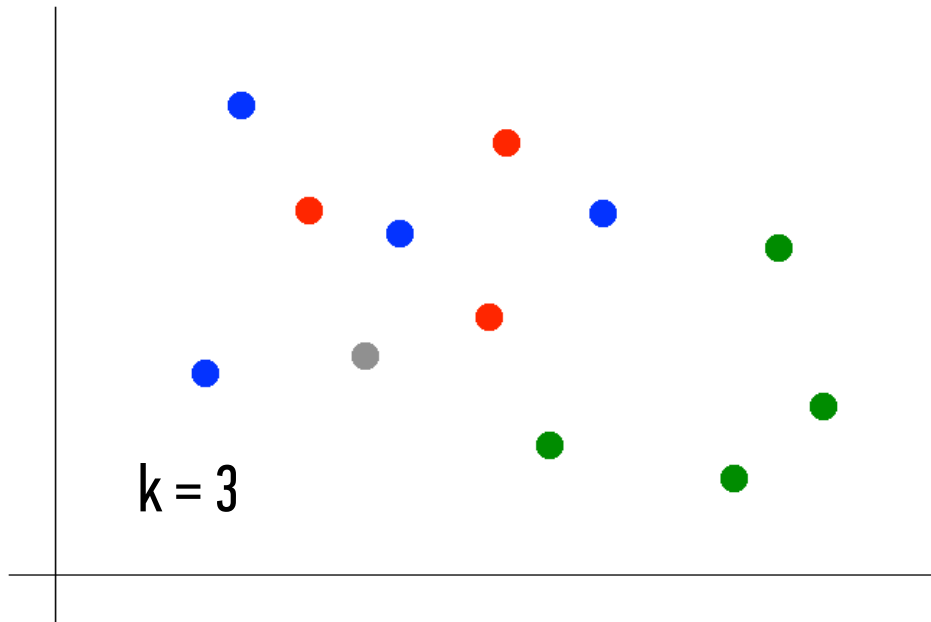
training set

test set

dataset

model

new data

# III. KNN CLASSIFICATION

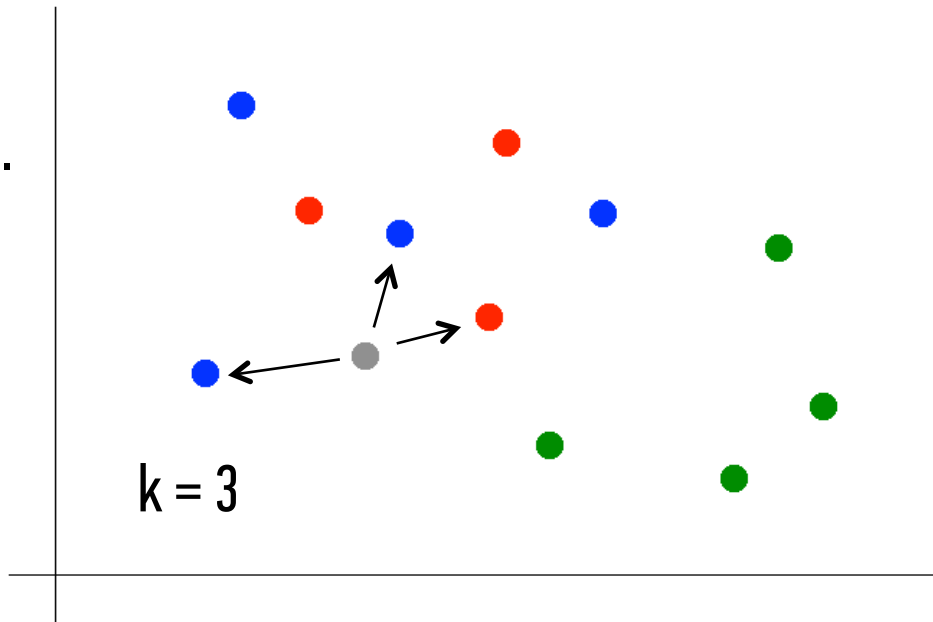Suppose we want to predict the color of the grey dot.

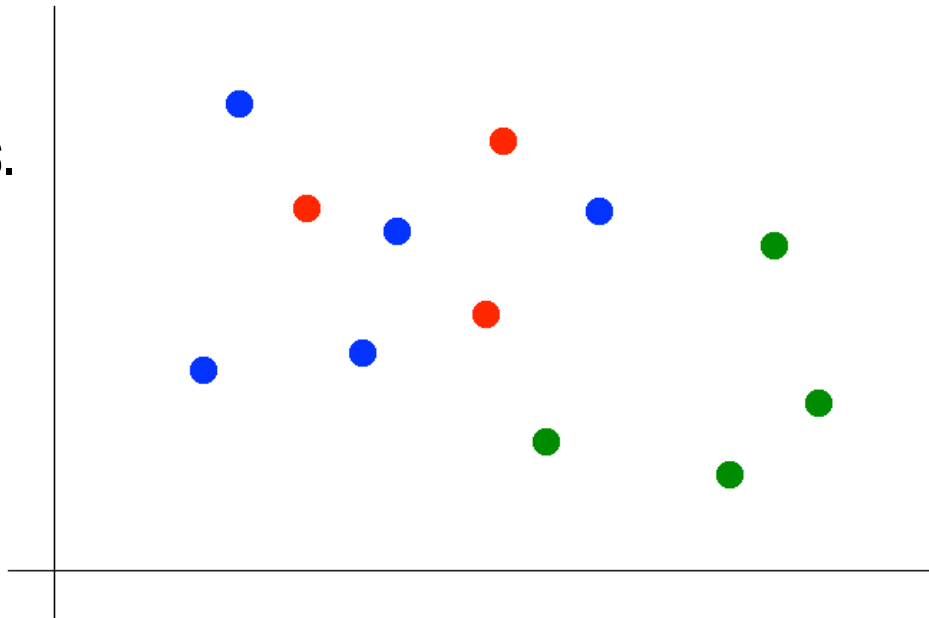Suppose we want to predict the color of the grey dot.

1) Pick a value for k.



k = 3

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

1) Pick a value for k.
2) Find colors of k nearest neighbors.
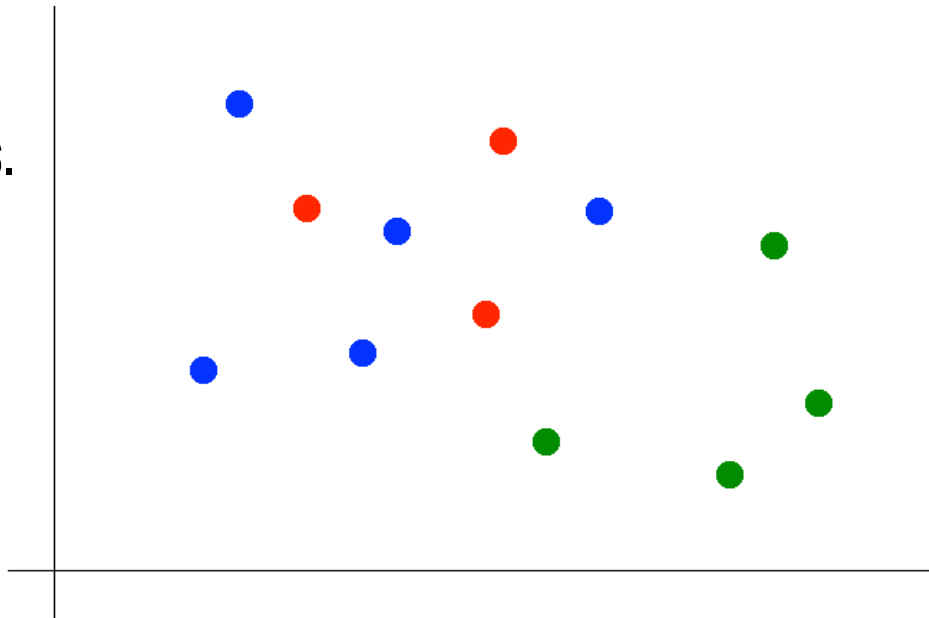
k = 3

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.


Q: What does nearest mean?

# DISCUSSION