

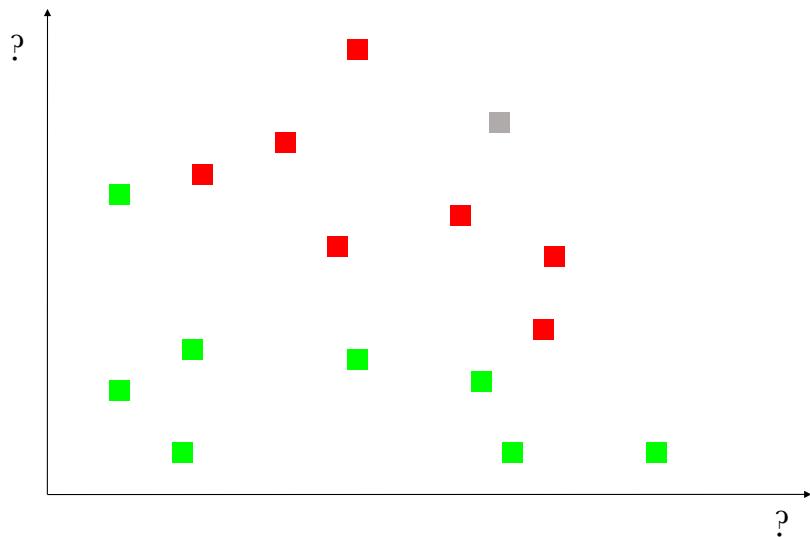
Beginner-friendly ML school

**Viviana Acquaviva
(CUNY)
vacquaviva@citytech.cuny.edu**

What is machine learning?

**THE ART OF
TEACHING A MACHINE
TO MAKE DECISIONS
(e.g. recognize objects,
similarities and differences,
patterns, signal vs noise)**

Our brain machine learns



**(of course)
ML is not the only way**

I CAN ALSO WRITE A FORMULA
(MAKE A MODEL)
TO PREDICT COLOR BASED ON COORDINATES

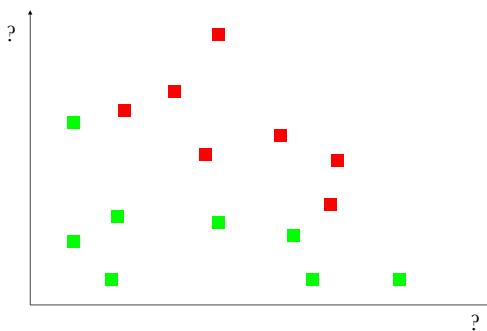
MACHINE LEARNING ALGORITHMS
PROVIDE AN “IMPLICIT MODEL” IN THE
FORM OF A PATH TO A DECISION
AND PERHAPS RESEMBLE MORE THE WAY
WE (HUMANS) SOLVE PROBLEMS

MACHINE LEARNING JARGON

Features are observable quantities known for all objects

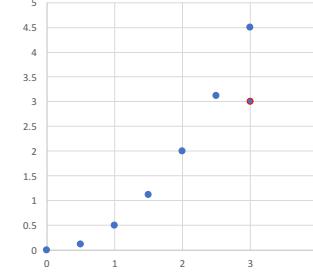
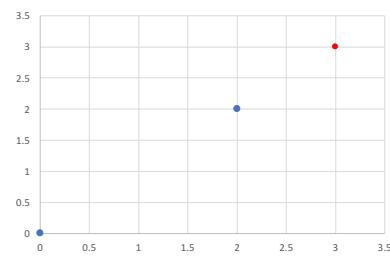
Label is the target property that we want to predict

SUPERVISED ML ASSUMES THAT WE HAVE A SET OF OBJECTS WITH KNOWN LABELS, called the **LEARNING SET**



Performance is limited by size and quality of the learning set.

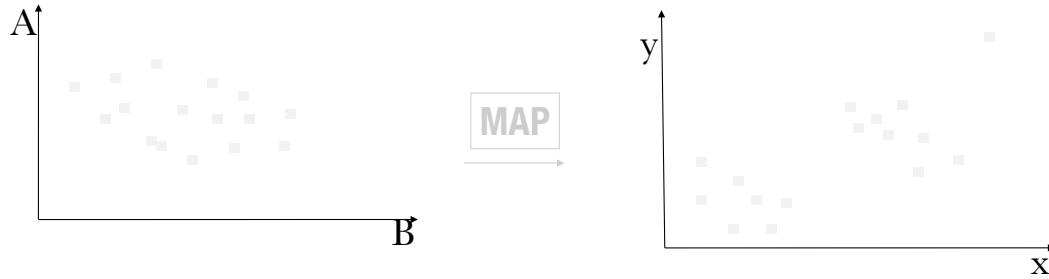
0 , 0
2 , 2
3 , 3 ?



Data representation (sometimes called feature engineering) and determining whether you have enough data to create a good model are crucial.

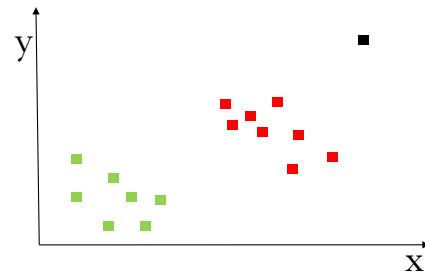
UNSUPERVISED MACHINE LEARNING

No labeled examples



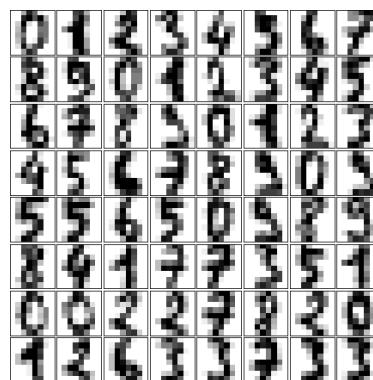
Useful to group together similar objects, find outliers, or find more efficient representations of data

Can be combined with human input or limited labels to understand the groups



Regression vs. classification

Usually we talk about classification when the target is a discrete variable (or class). For example in this image recognition problem:



There are a finite (10) numbers of possible outcomes.

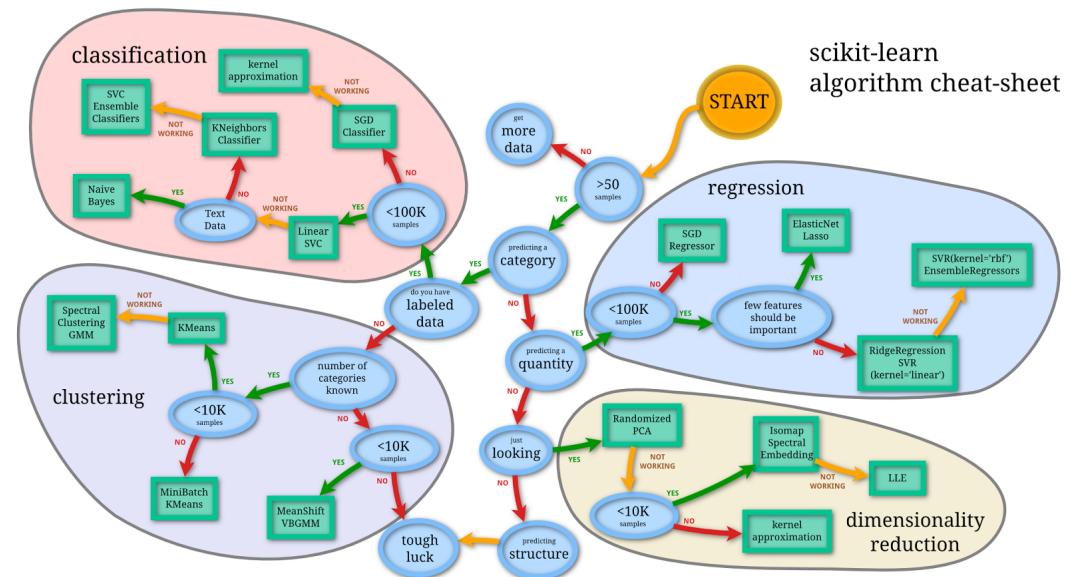
Regression vs. classification

- Vice versa, if we are trying to predict, say, the probability that it will rain in half a hour based on the current weather conditions, the outcome (target) is a continuous variable that can have all values between 0 and 1.



What if I was trying to decide whether or not I should bring an umbrella?

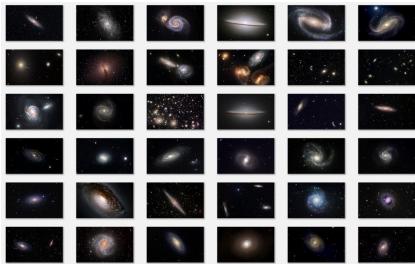
ALGORITHMS ABOUND



What can we do with them?

1. Save time.

Galaxy morphology



Trained humans are the best classifiers.
But what to do when
you have millions of objects?

Citizen science

The screenshot shows the Galaxy Zoo interface. At the top, there are navigation links: CLASSIFY, STORY, SCIENCE, GALAXY ZOO (which is highlighted in yellow), DISCUSS, PROFILE, and LANGUAGE. Below this, a main heading reads "Few have witnessed what you're about to see" with a subtext "Experience a privileged glimpse of the distant universe as observed by the SDSS, CTIO and VST.". A section titled "Classify Galaxies" contains a brief description: "To understand how galaxies formed we need your help to classify them. Every classification you make, no matter how simple or complex, you may even be the first person to see the galaxies you're asked to classify." A yellow "Begin Classifying" button is at the bottom of this section. The central part of the page features a large, detailed image of a spiral galaxy.

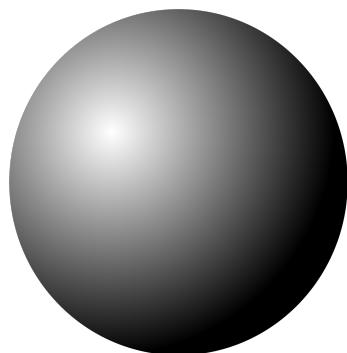
Automated Classification via Machine Learning

(supervised/unsupervised
approach, see e.g. Hocking et al 2017)

2. Provide an alternative to simplistic models.

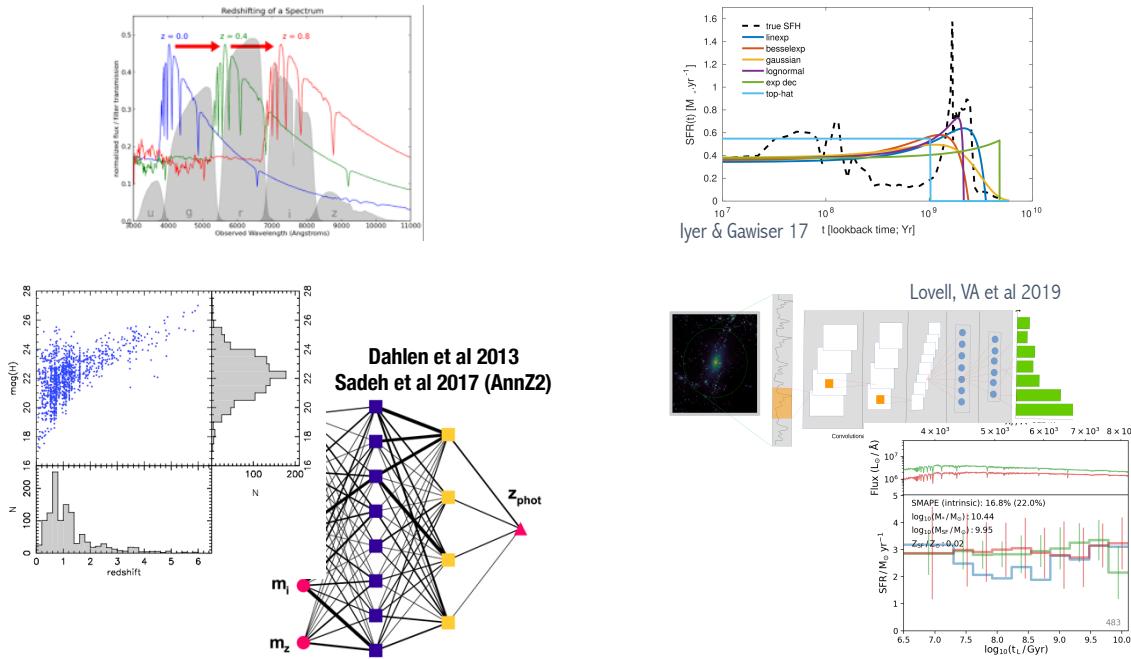


≠



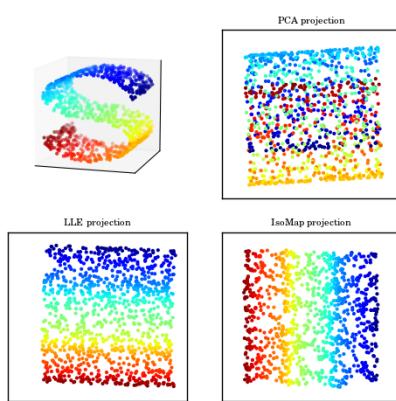
2. Provide an alternative to simplistic models.

Example: Galaxy Photometric Redshifts or SFH.



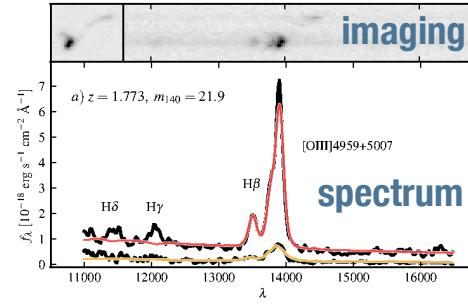
3. Make problems more tractable, e.g.:

Via dimensionality reduction



3D -> 2D
With various degrees
of information loss

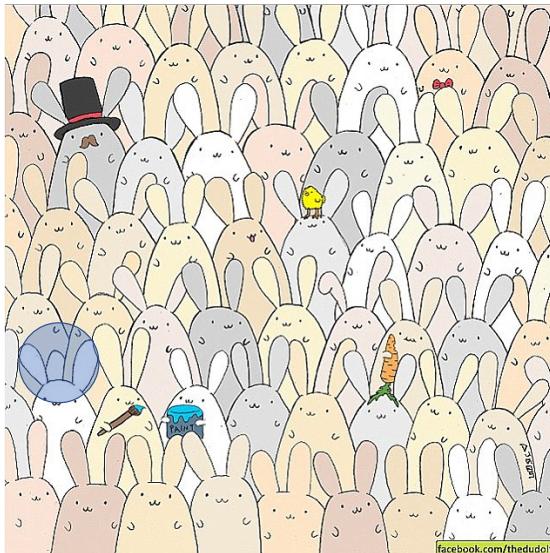
By making mix-and-match easier



Parameter estimation:
 $P(\text{model} \mid \text{combined data}) = \text{mess} \odot$

In machine learning:
Higher tolerance to mixed features ☺

4. Allow serendipitous and data-driven discoveries



Outliers might be very interesting objects. In model fitting they are often discarded.

Data mining techniques can tell you about new categories of objects.

Machine Learning vs Model Fitting

ML

- Data-driven (only as good as the data)
- Usually generalizes poorly (model derived using some data can't be applied blindly to different data)
- Answers questions; interpretation is possible but might be non-trivial
- Fast(er)
- More robust/accommodating of mixed and missing data
- Allows serendipitous discoveries

MF

- Intuition or model-driven (only as good as the scientist :))
- Generalizes well if model (physics) is well understood
- Builds subject matter knowledge; easier to interpret
- Might be computationally intensive
- Dealing with heterogenous data often a pain in the neck
- Leads to loss of information if models are too simplistic

Synergy is often the best strategy

Questions?

Let's play a fun game 😊

Let's answer these questions

Is this supervised or unsupervised ML?

Is this a classification or regression problem?

What could be useful features (data) to collect?

Let us focus on supervised learning for now.

Important:
you shouldn't use all of your learning set
to build your model.

It is customary to split the learning set into
a **training set** and **test set**



By building a model on the training set and applying it
to the test set, you “mimic” what happens when your model sees new data
for which the labels are not known. **Note:**
Otherwise you would be too optimistic! **You still are.**

For example...

Math Quiz #1 - Teacher's Answer Key

$$1) \ 2 \ 4 \ 5 = 3$$

$$2) \ 5 \ 2 \ 8 = 2$$

$$3) \ 2 \ 2 \ 1 =$$

$$4) \ 4 \ 2 \ 2 =$$

TEST

**Diagnosing and Improving
Machine Learning
algorithm performance:
cross validation,
performance metrics,
diagnostics**

The goal of the training set and test set split is to be able to evaluate performance on unseen examples.

The test set “mimics” new data.

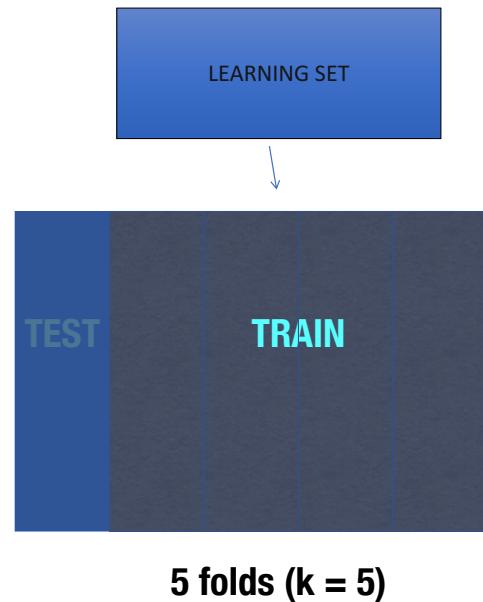
However, it might be better to pick more than one test split.

Why would this be a good idea?

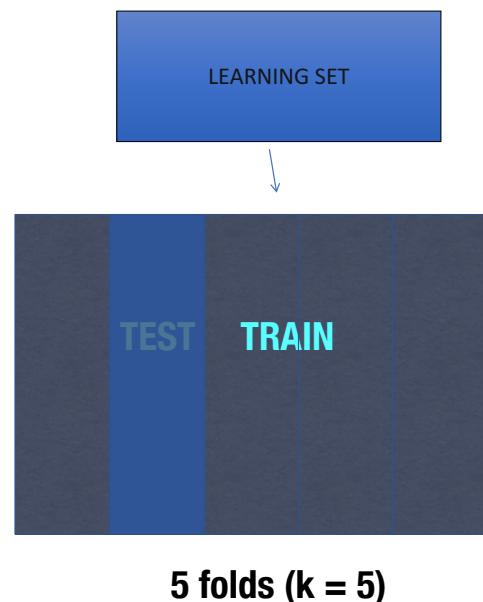
1. We use all the training data for training!

2. We avoid the risk of under/overestimating performance because of a “weird” pick of train/test split.

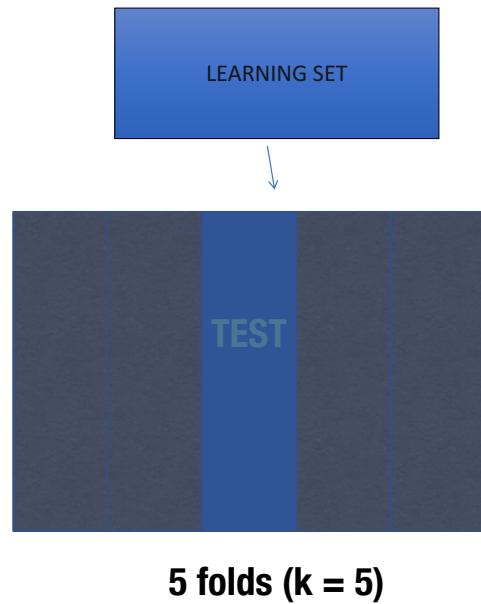
k-FOLD CROSS VALIDATION



k-FOLD CROSS VALIDATION

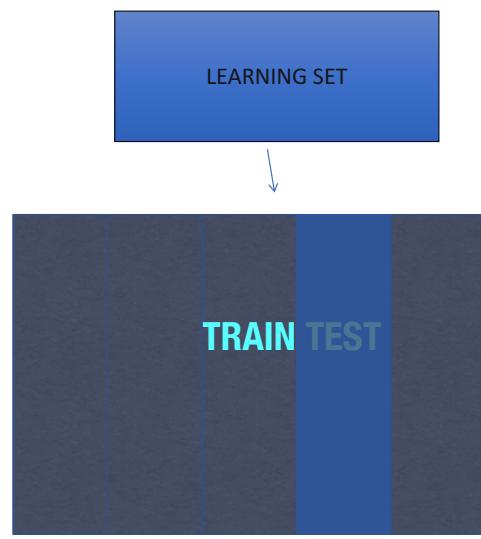


k-FOLD CROSS VALIDATION



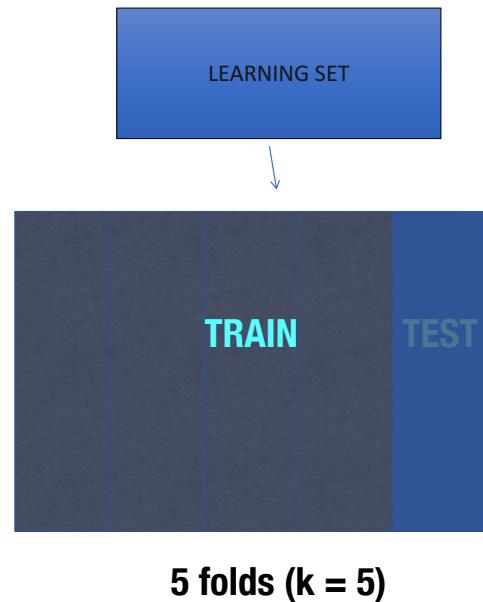
5 folds ($k = 5$)

k-FOLD CROSS VALIDATION



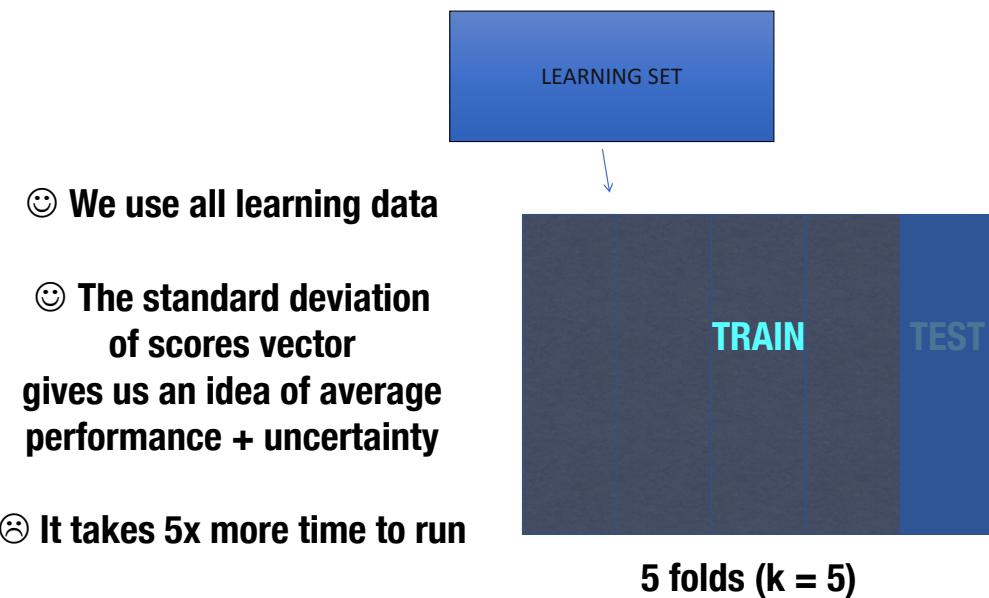
5 folds ($k = 5$)

k-FOLD CROSS VALIDATION



5 folds ($k = 5$)

k-FOLD CROSS VALIDATION



Diagnosing a ML algorithm

BIAS

Algorithm
can't capture
complexity
of rule
connecting
input and output

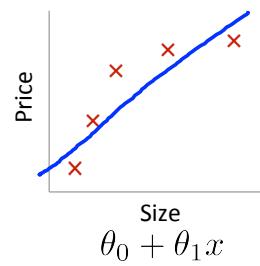
VARIANCE

Algorithm
is excessively
tailored
to training set
and generalizes
poorly

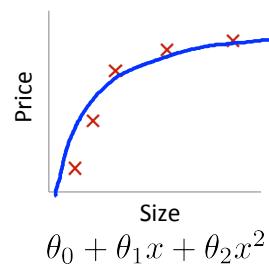
UNDERFITTING

OVERFITTING

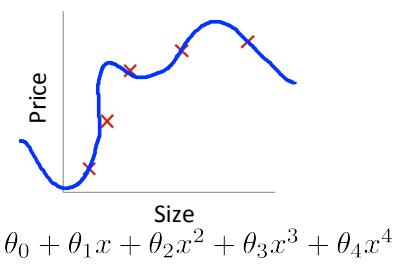
Bias/variance



High bias
(underfit)
 $d=1$



"Just right"
 $d=2$



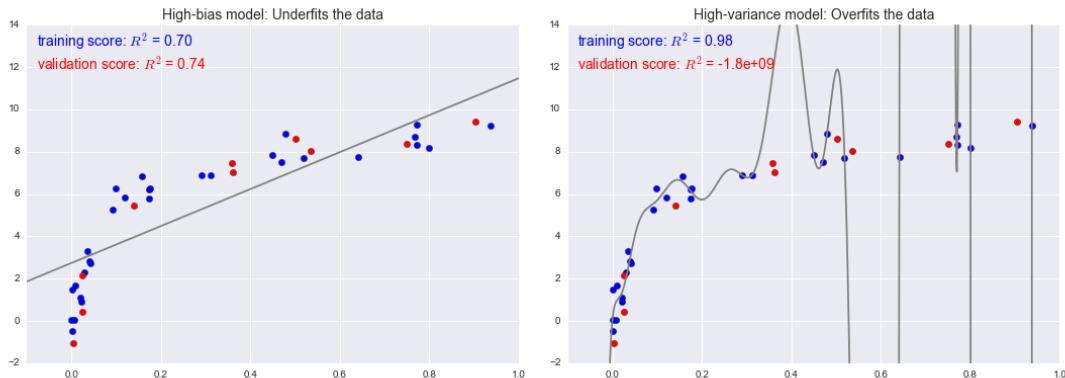
High variance
(overfit)
 $d=4$

Andrew Ng

slide from Andrew Ng's Coursera ML class

How can we diagnose high variance vs high bias?

figure from Jake VanderPlas' book



High bias: test and train error are similar but high

High variance: there is a gap between test and train error because algorithm does not generalize well

Improving high bias

1. Try using different features.
2. Try engineering new features.
3. Try a more complex algorithm.

Improving high variance

1. Try reducing the number of features.
2. Try a less complex algorithm/change parameters.

Also: Check if you need more training data

Useful diagnostics: learning curves

plot performance of algorithm for train and test set
as a function of size of training set



Note: which algorithm is the best?

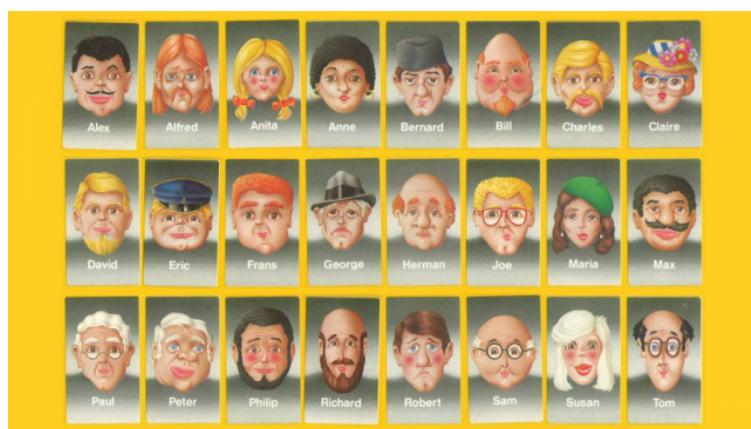
- **High bias low variance**
- **High variance low bias**
- **Lowest gap between train/test**
- **Highest test scores**

Note: which algorithm is the best?

- High bias low variance
- High variance low bias
- Lowest gap between train/test
- **Highest test scores**

DECISION TREES

- Work by splitting data on different values of features
- Simplest trees are binary trees
- If categorical features, the split would be on yes/no
- If numerical, the split would be on a certain value (e.g. $x > 100$ or $x < 100$)



Example: Look at this 2-feature data set. How should we split?

Figure credit:
Gilles Louppes

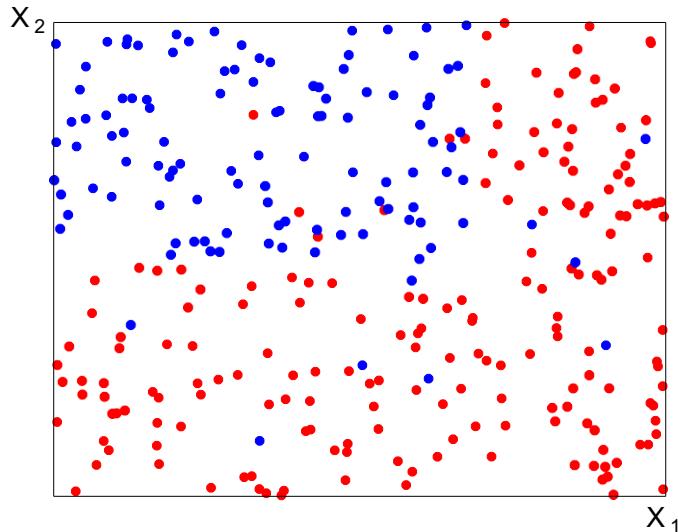
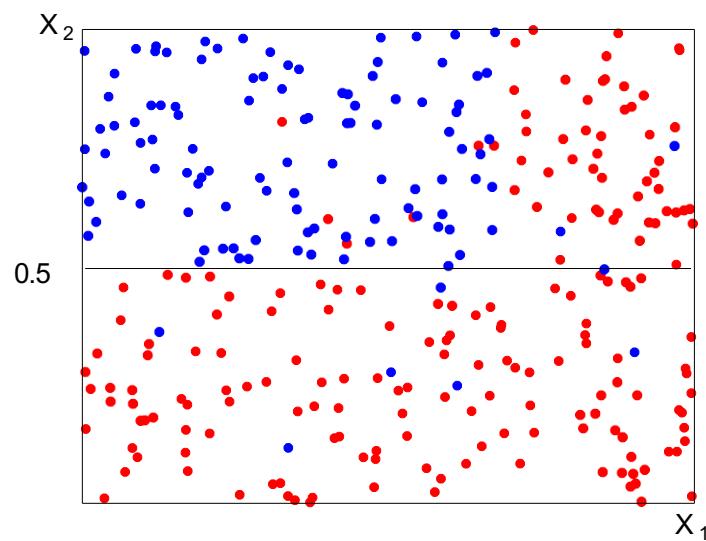
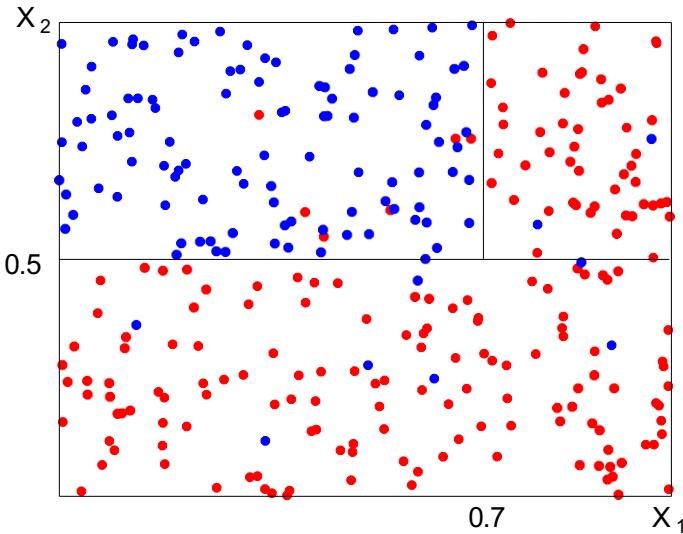


Figure credit:
Gilles Louppes

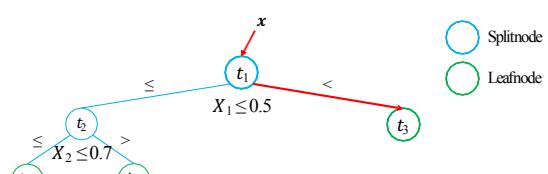
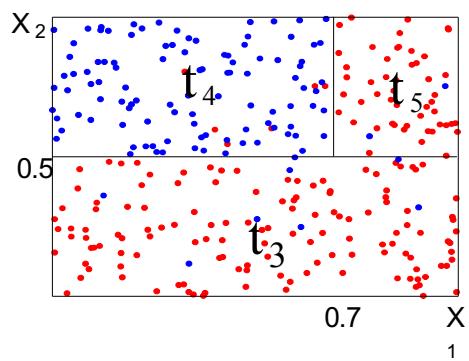


Should we stop?

Figure credit:
Gilles Louppes



Decision trees: defined by splits and leaves



Figures credit:
Gilles Louppes

How many splits in this tree?

How many leaves?

How do we decide whether we should keep splitting?

Pseudo code for decision trees

```

function BuildDecisionTree(L)
    Create node  $t$  from the learning sample  $L_t = L$ ;
    calculate (im)purity
    if the stopping criterion is met for  $t$  then
         $y^t = \text{some constant value}/\text{class}$  (MAKE PREDICTION)
    else
        Find the split on  $L_t$  that maximizes impurity
        decrease
         $s^* = \arg \max_{s \in Q} \Delta_i(s, t)$ 
        Partition  $L_t$  into  $L_{t_L} \cup L_{t_R}$  according to  $s^*$ 
         $t_L = \text{BuildDecisionTree}(L_{t_L})$ 
         $t_R = \text{BuildDecisionTree}(L_{t_R})$ 
    end if
    return  $t$ 
end function

```

stopping criterion
Gini (im)purity = 0

Gini (node L) =
$$1 - \sum f(i)^2$$

where $f(i)$ is the frequency of
the i-th class

Gini (splits Lt and Lr) =
$$L_L/L * (1 - \sum f(i)^2) +$$

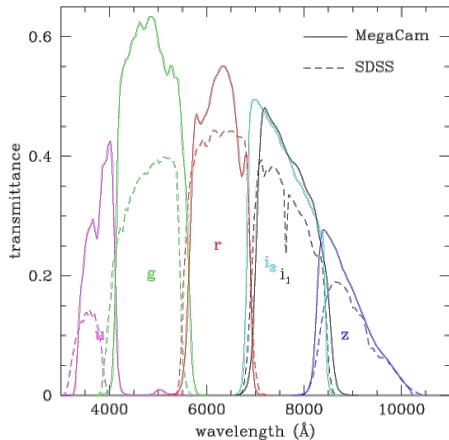
$$L_R/L * (1 - \sum f(i)^2)$$

where $f(i)$ is the frequency of
the i-th class

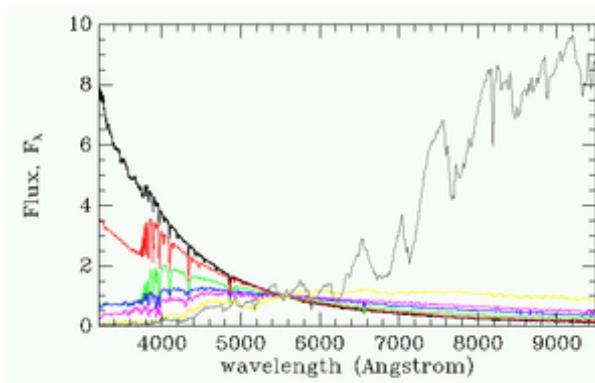
Note: splits
happen along features!

Code adapted from Gilles Louppes

Our first example will be a supervised classification problem, in which we are trying to decide if a star is a variable star (RR Lyrae), based on imaging data in 5 bands (u, g, r, i, z) and four colors (u-g, g-r, r-i, i-z)



range of observed brightness



spectra of different stellar types