

Lecture 1: Introduction to Machine Learning

Rafael Martínez-Galarza
Harvard-Smithsonian Center for Astrophysics



This lecture:

- Introduction to Neural Networks
- Feedforward Networks
- Stochastic Gradient Descent
- Back propagation algorithm

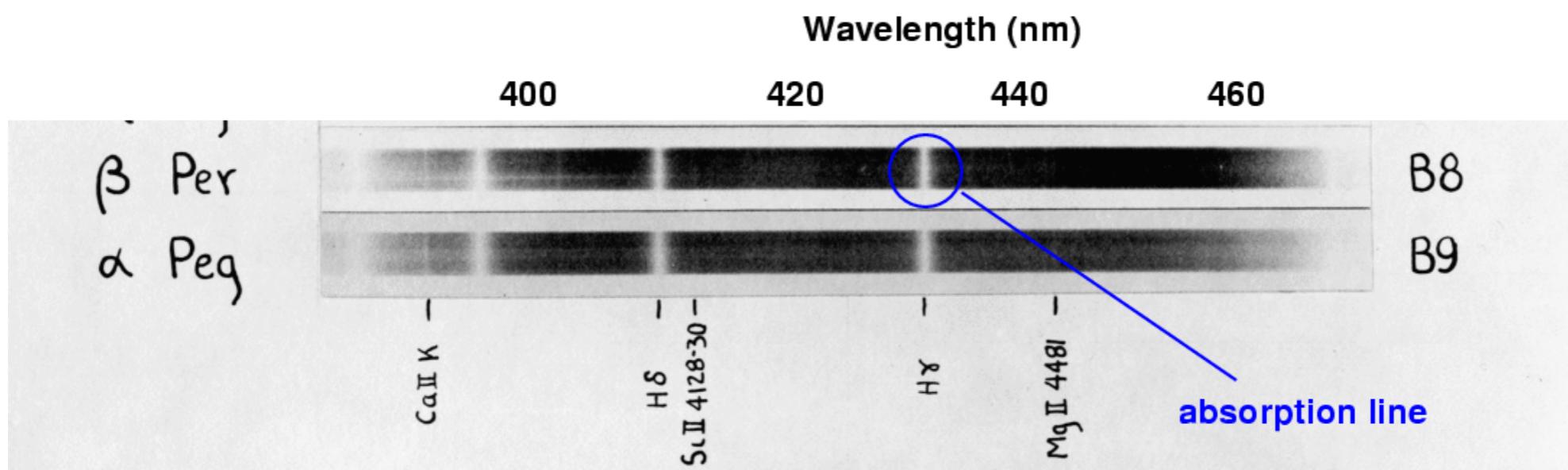
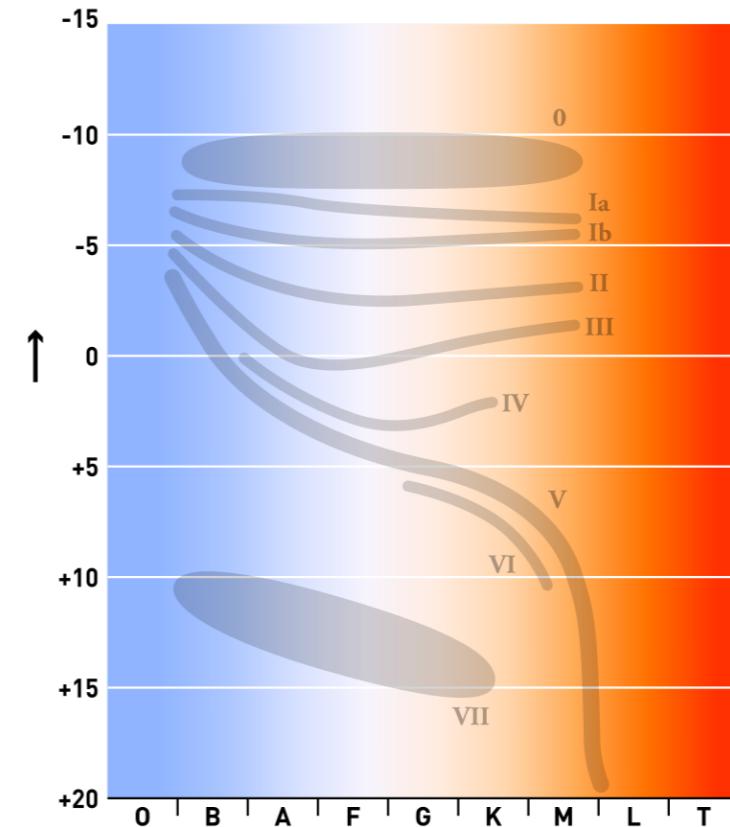
Big data in astronomy... then

Annie Jump Cannon at her Harvard desk



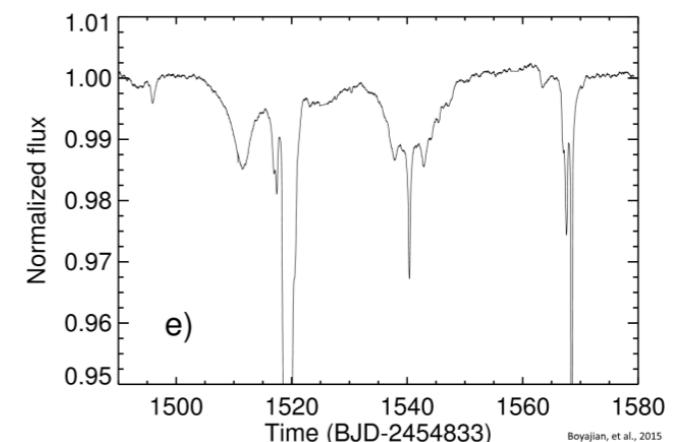
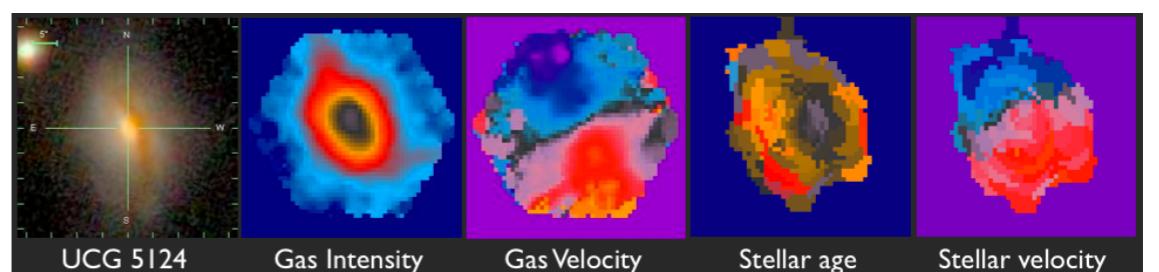
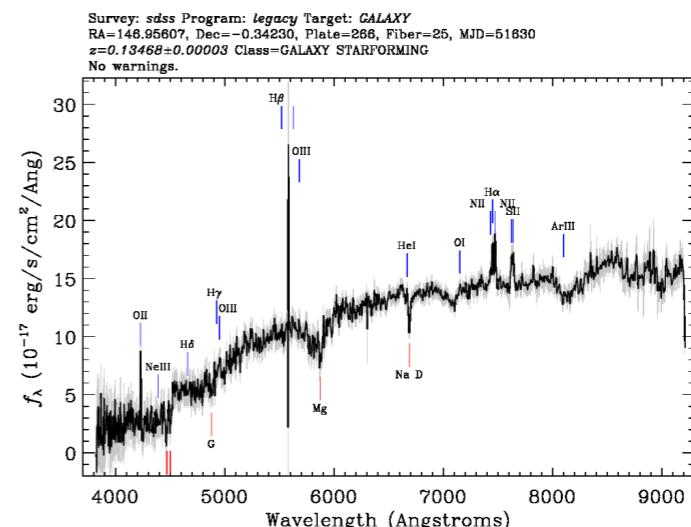
She analyzed over 300,000 stellar spectra during her lifetime.... by hand.

We owe her the stellar classification system we use today



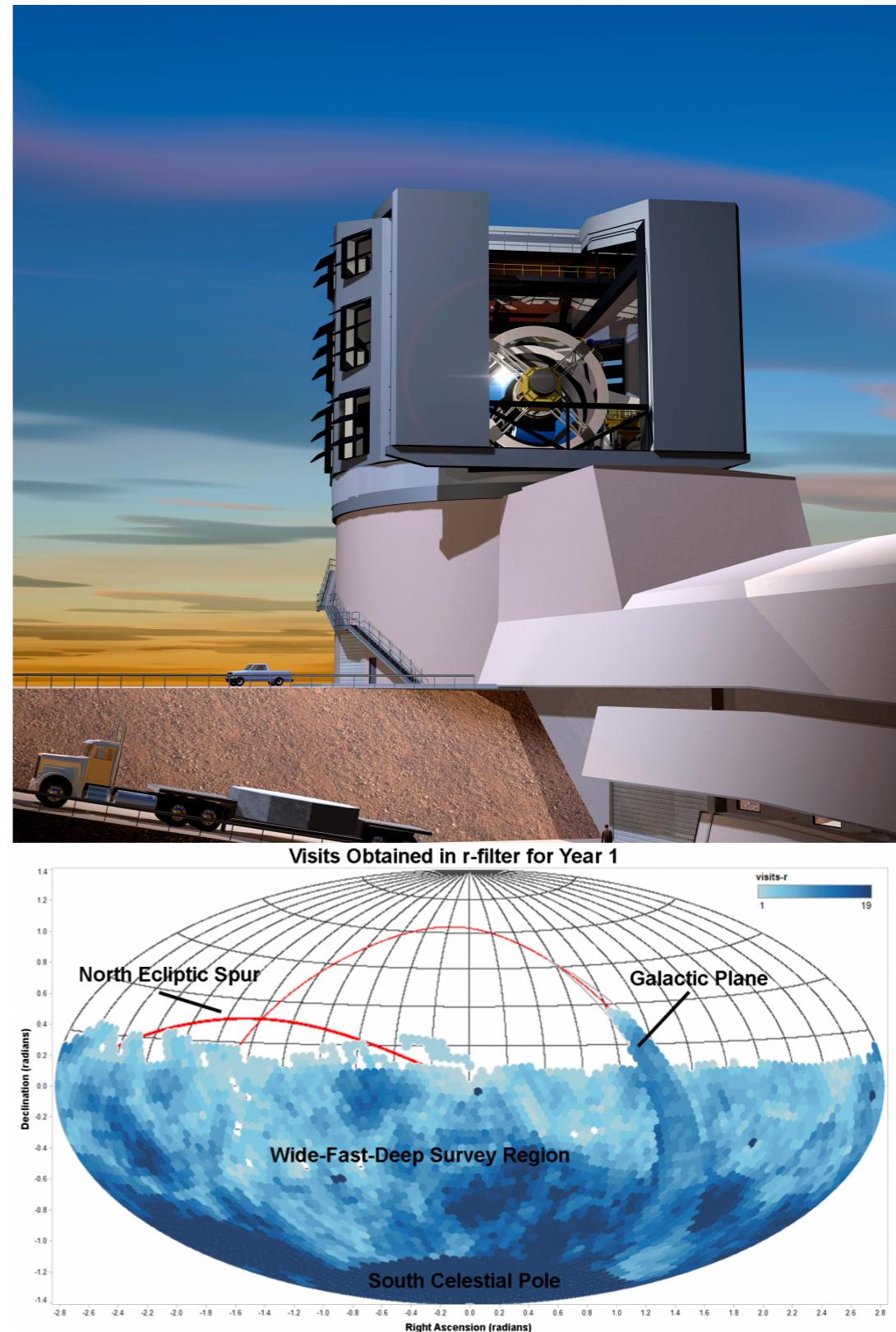
Big data in astronomy... now

- Images. Limited by spatial resolution.
Spatial resolution explosion with the next generation of telescopes.
- Spectra. Limited by wavelength coverage and spectral resolution. Huge coverage of millions of spectra in the sky possible thanks to SDSS.
- IFU spectra. 2D spectra of astronomical objects. Spatial and spectral information together in the same dataset.
- Light curves. Variation of brightness as a function of time. Future synoptic surveys will imply an explosion of light curves.



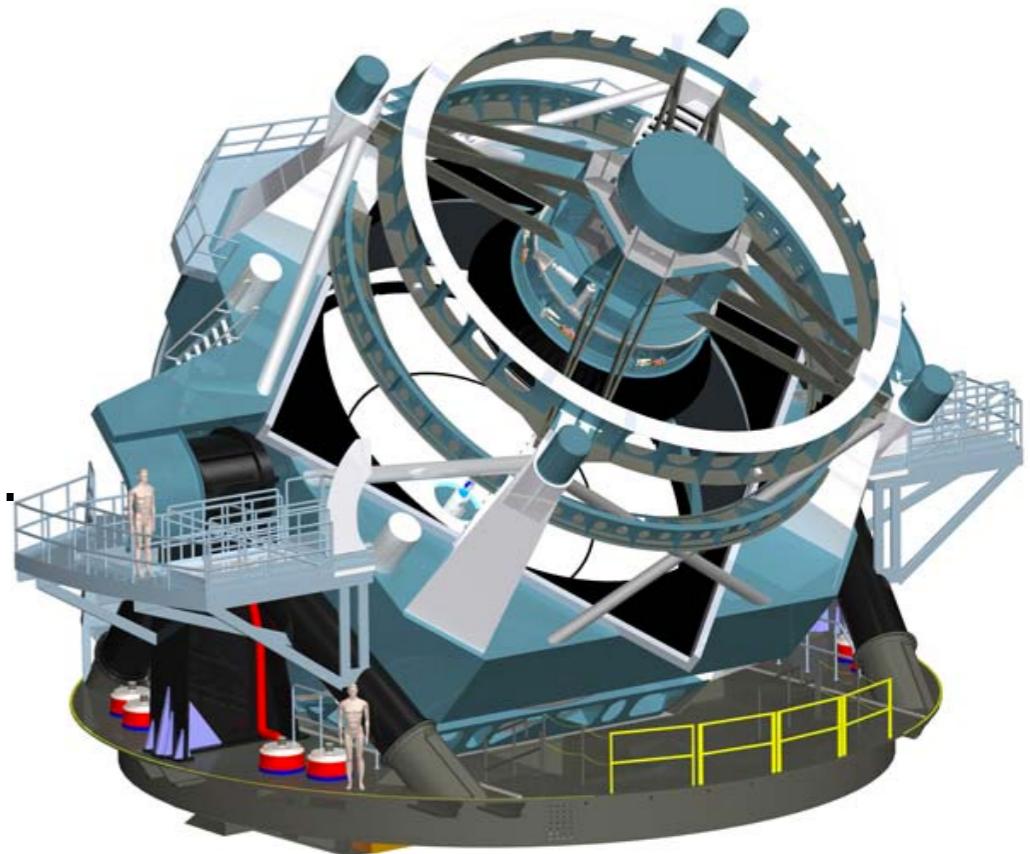
The motivation: LSST is coming

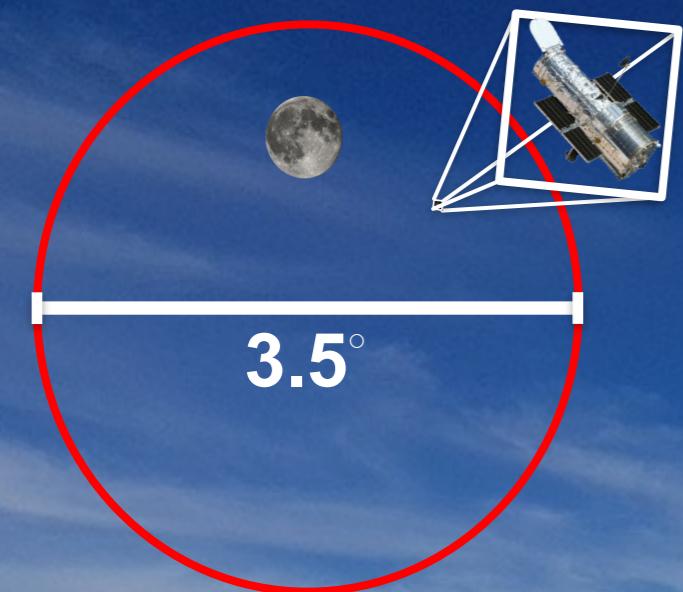
- The Large Synoptic Survey Telescope is a 8.4m reflector currently under construction in Chile (first light expected in 2021).
- Design concept: a survey that will take an image of every part of the entire visible sky every few nights, in six bands, for 10 years.
- Transients and variable stars: periodic and non-periodic variable sources will be studied in detail, and new types are expected at very short and very long timescales.



The Large Synoptic Survey Telescope

- LSST is an excellent example of what we mean by the new data-intensive astronomy
 - photometry of the entire southern sky every 3-4 nights for over 10 years.
 - ugrizy multiband data.
 - 30,000 GB per night.
 - Final catalog: 100s of petabytes.
 - ~1000 observations per field





Movie: LSST Corporation

slide: Andy Connolly for LSST Corporation

3,200,000,000
pixels



Slide by R. Hlozek

LSST, Todd Mason, Mason Productions



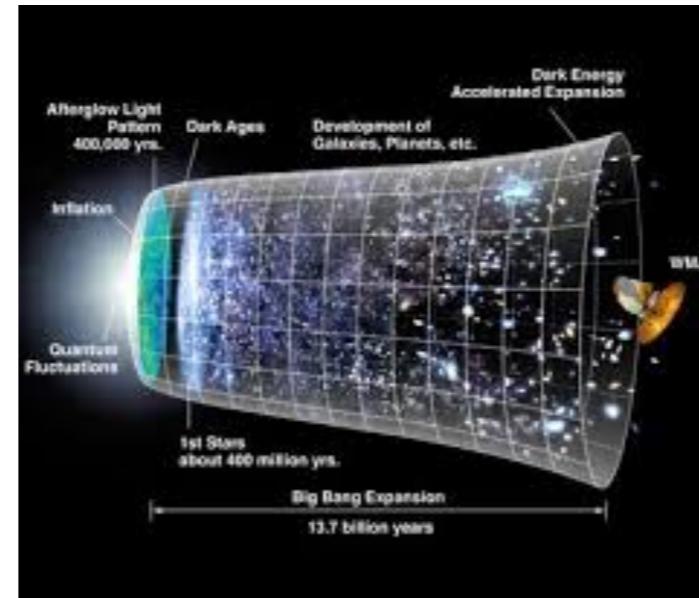
How can we use machine learning in astronomy?

Gravitational Lensing



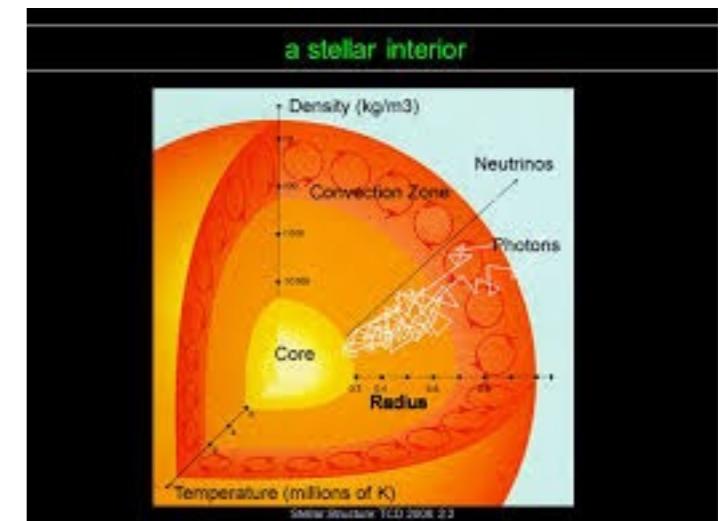
Convolutional Neural Networks can help us understand the patterns of light bended by faraway clusters of galaxies

Precision cosmology



Searches for distant supernova in large time-domain surveys will be a lot more efficient using machine learning.

Numerical models



Deep learning methods will soon allow us to replace complex computational models, such as convection, with machine learning analogs.

Machine Learning is becoming a fundamental tool for astronomers. This lecture will introduce you to some basic concepts and practical applications.

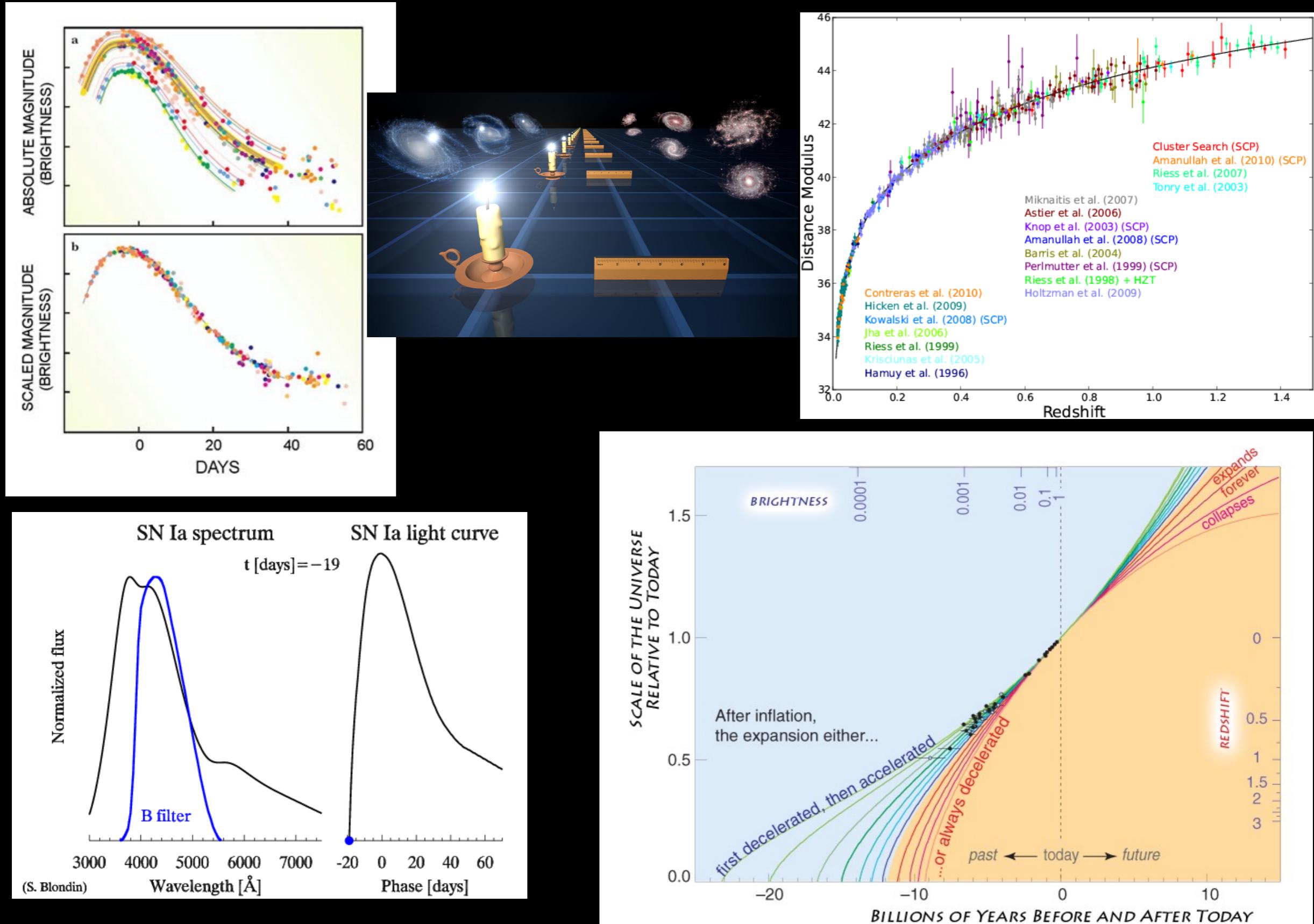


image composite: BJ Fulton



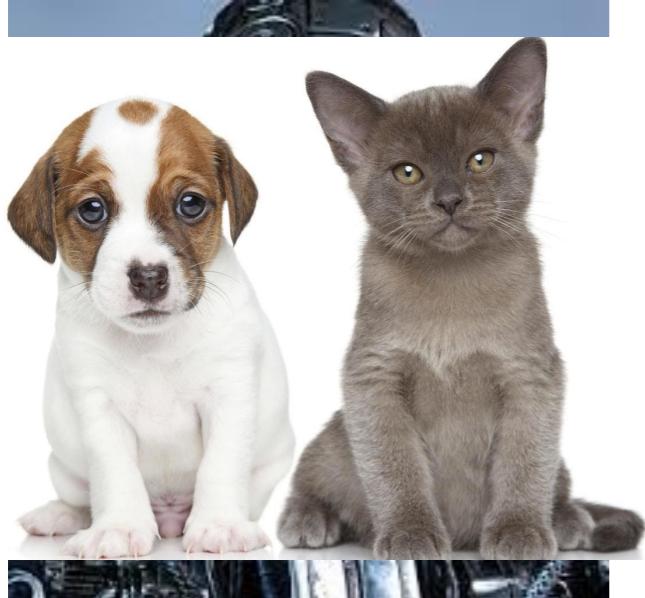
image composite: BJ Fulton

THE ACCELERATING UNIVERSE



Machine Learning

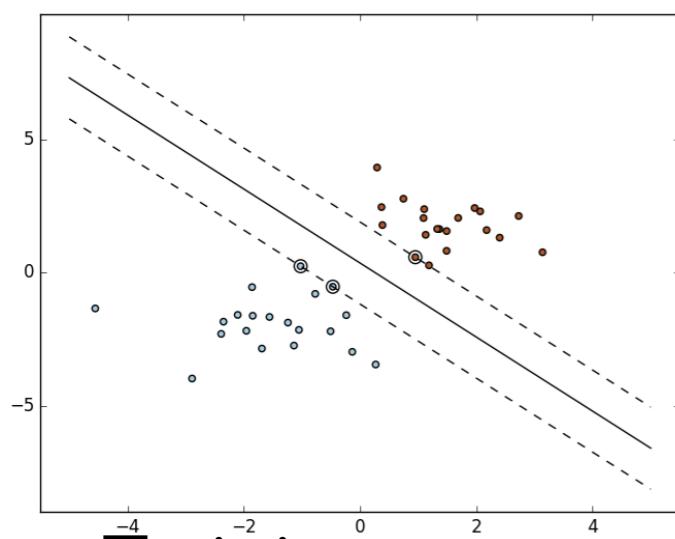
Machine Learning



The field of computer science that uses statistics and linear algebra to have computers program “learn” complex correlations between variables of a large dataset, so it achieves the skill to make predictions of new data that has not been observed.

The performance of a ML algorithm improves as more data is included in the analysis, with respect to certain metric (e.g., the accuracy at classifying “dogs vs cats”)

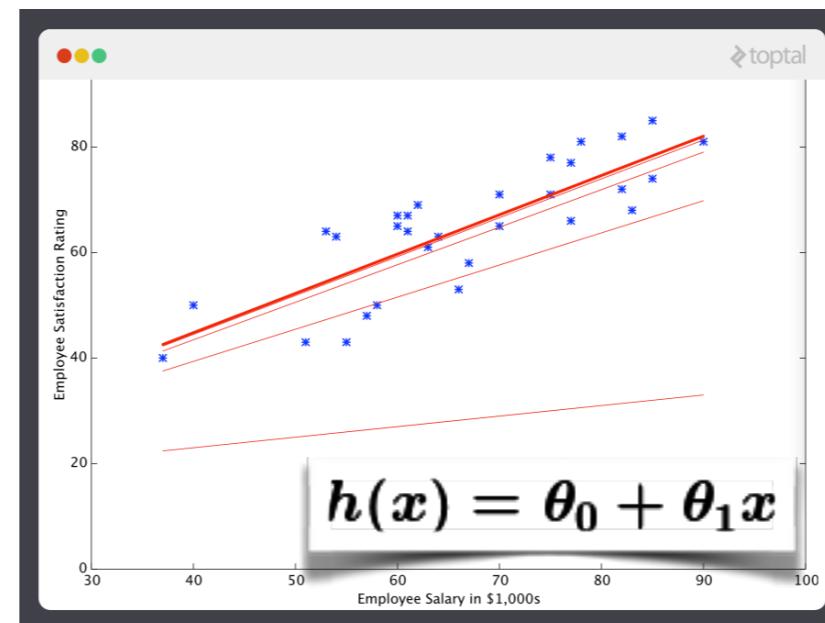
Classification



Training set:

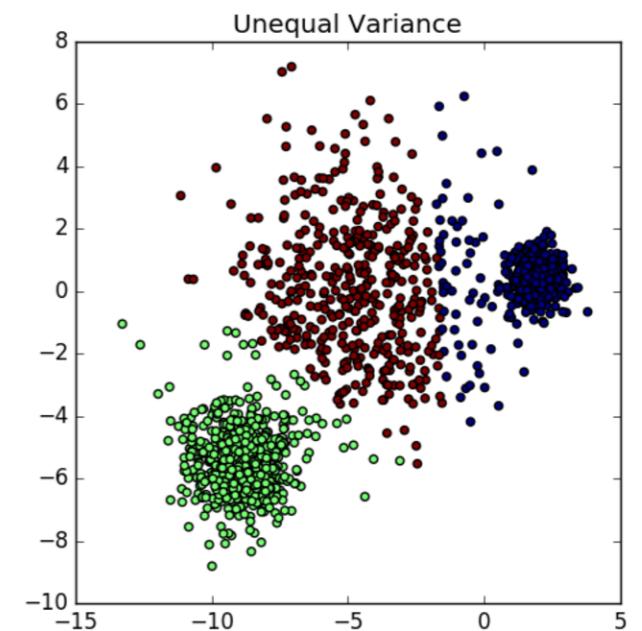
Human experience

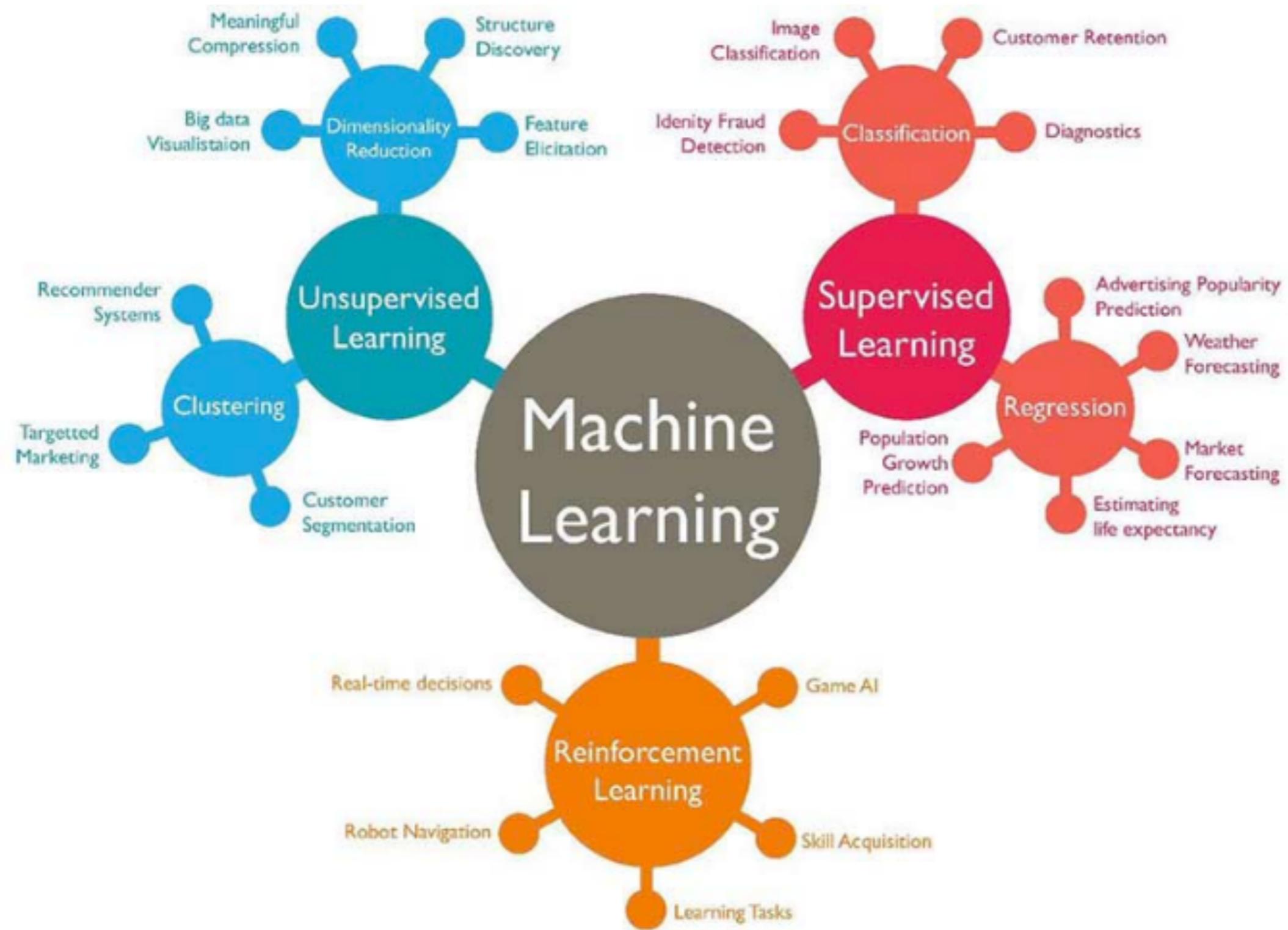
Regression



$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h(x_{t,i}) - y)^2$$

Clustering





Vincent Grainville

Types of Machine Learning

- At a glance

Supervised Learning

- ◆ Makes machine learn explicitly
- ◆ Data with clearly defined output is given
- ◆ Direct feedback is given
- ◆ Predicts outcome/ future
- ◆ Resolves classification & regression problems



Unsupervised Learning

- ◆ Machine understands the data (Identifies patterns/ structures)
- ◆ Evaluation is qualitative or indirect
- ◆ Does not predict / find anything specific



Reinforcement Learning

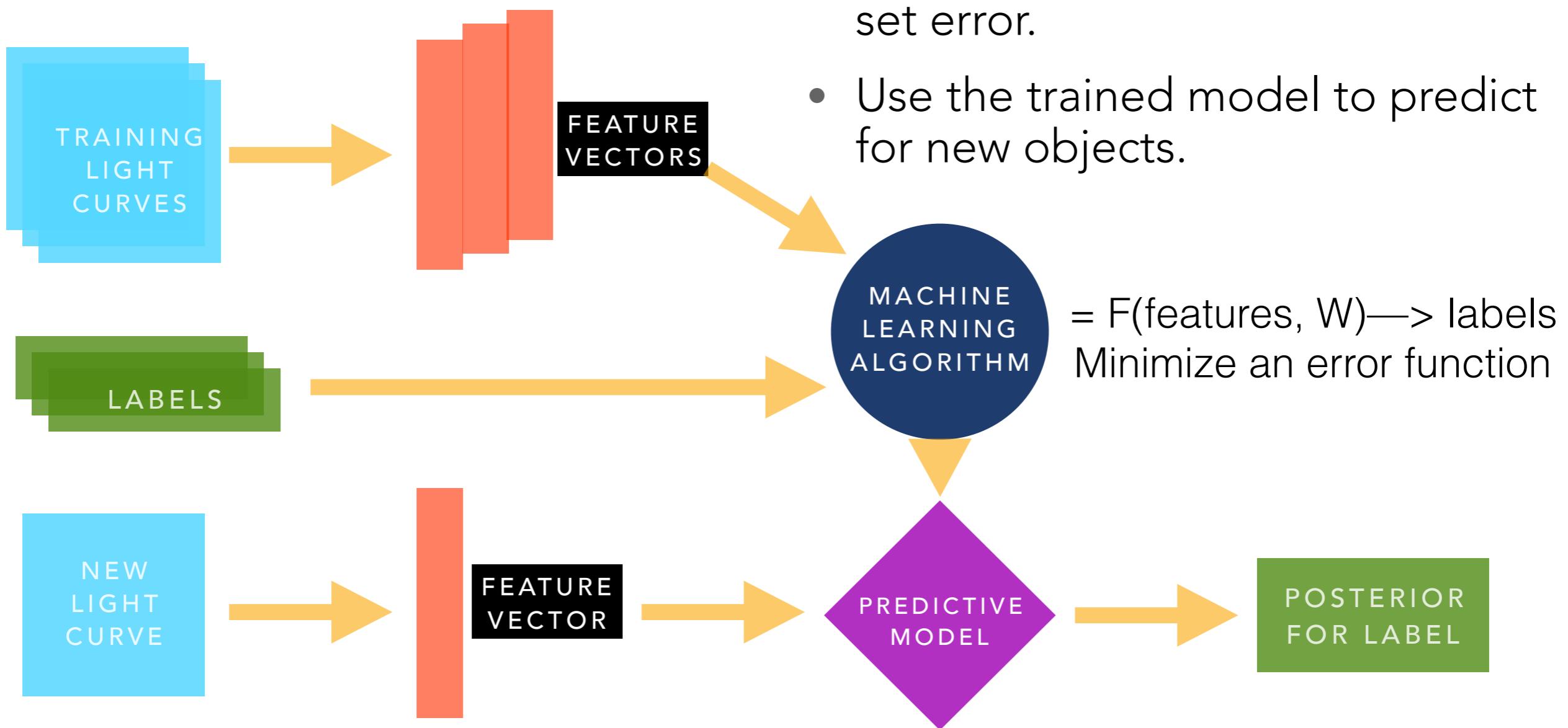
- ◆ An approach to AI
- ◆ Reward based learning
- ◆ Learning from +ve & -ve reinforcement
- ◆ Machine learns how to act in a certain environment
- ◆ To maximize rewards



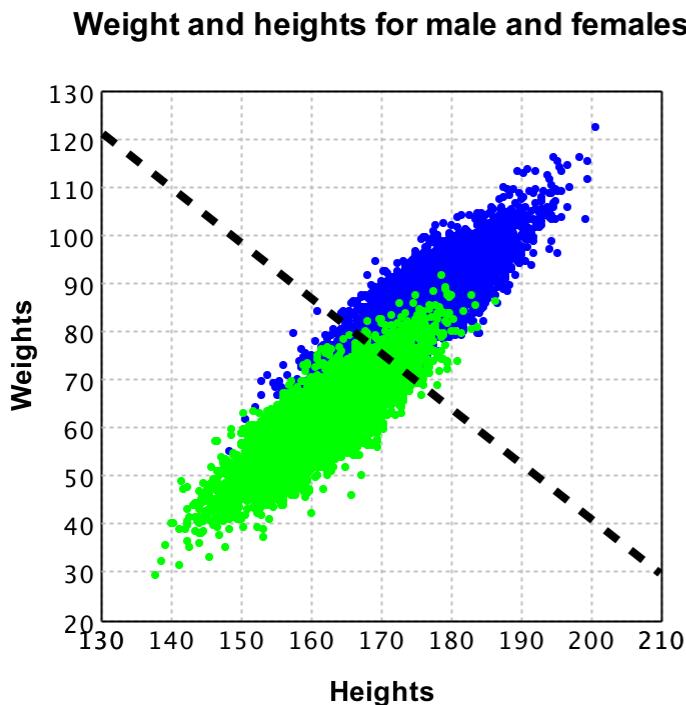
SUPERVISED LEARNING

(Supervised) ML in a nutshell

- Select a representative training set.
- Fit your model (a function of the features) minimizing the training set error.
- Use the trained model to predict for new objects.



Classification



- Consider this dataset
 - Green: females
 - Blue: males
 - How do we find the best decision boundary in the parameter space in order to decide what class to assign to a new person of whom all we know is weight and height?
-
- In ML jargon, height and weight are called the **features**, the dashed line is the **decision boundary**, or discriminator. And **female/male** are the **classes**, to which we assign labels.
 - The initial green and blue dots for which we know the classes in advance are the **training set**. The “learning” happens when the algorithm generalizes the patterns in this training set
 - Supervised classification: finding a function f such that:
$$f: f(\text{features}) \longrightarrow \text{classes}$$
 so that $\text{Error}(\text{prediction})$ is minimized

Intuition on probability: a coin flip



Head



Tail



Head



Tail

Can you predict, for a single toss, whether it will land heads or tails?

Can you predict instead the rate of tails for, say, 10, or 100 tosses?

Why?

Our intuition of probability corresponds to a model of nature

We speak of the probability $P(\text{heads})$ that the coin will land heads.

$P(\text{heads})$ is the fraction of times that the toss results in heads

If the coin is fair, then $P(\text{heads}) = 0.5$



What about a fair dice?

What is $P(2)$?

What if we have a non-fair coin or dice?

What if we toss a coin, say, 4 times?
 What is the probability of getting 4 heads ($P(4H)$)?
 Or 3 heads and 1 tail $P(3H)$?

$$P(4H) = 1/16 = 0.025$$

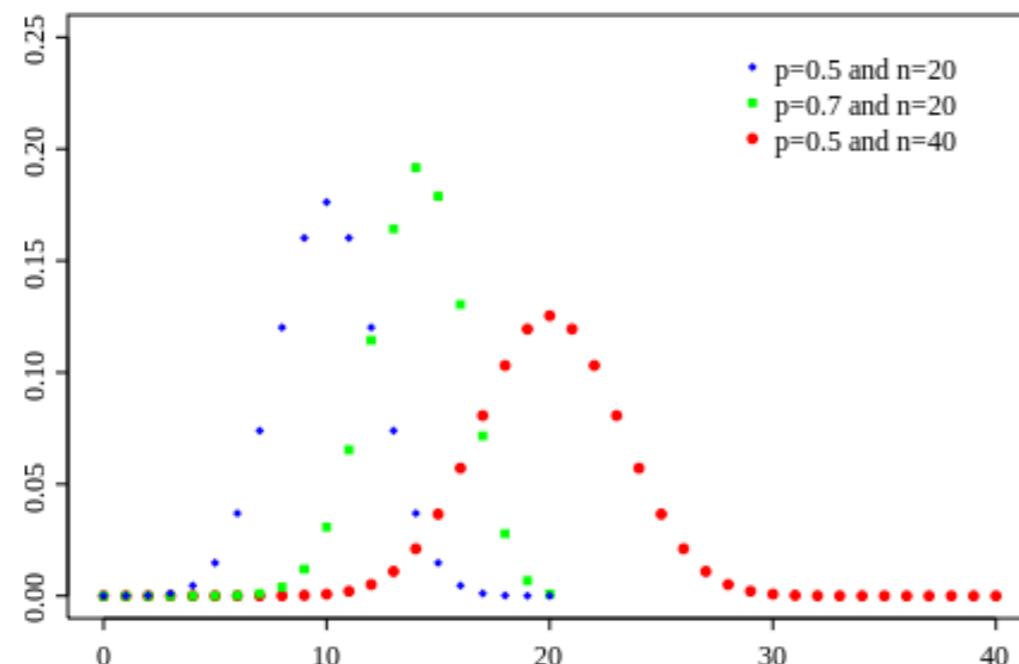
$$P(3H) = 4/16 = 0.25$$

Mathematically, we write this as:

$$\binom{n}{k} = \frac{n!}{k!(n-k)!} \quad \text{If } n=1: \text{Bernoulli}$$

There are $(n k)$ ways we can distribute k successes (heads) in a sequence on n trials
 In general (Binomial distribution):

$$f(k, n, p) = \Pr(k; n, p) = \Pr(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$



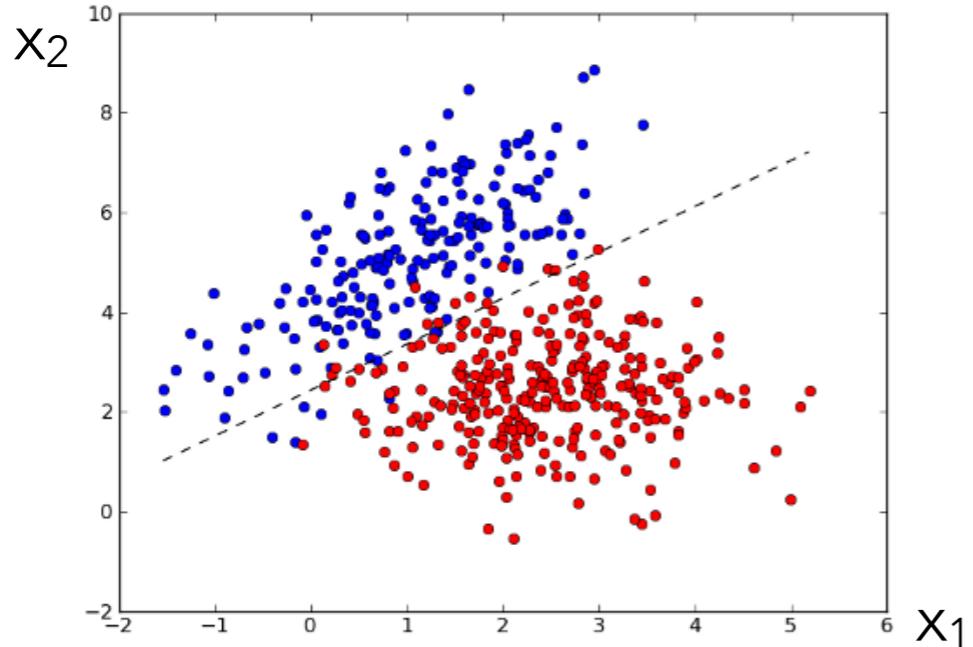
4 coins

Number of Heads	Results of Flips	Number of Ways
0	0 0 0 0	1
1	0 0 0 1 0 0 1 0 0 1 0 0 1 0 0 0	4
2	0 0 1 1 0 1 0 1 1 0 0 1 0 1 1 0 1 0 1 0 1 1 0 0	6
3	0 1 1 1 1 0 1 1 1 1 0 1 1 1 1 0	4
4	1 1 1 1	1

Logistic regression as a classifier



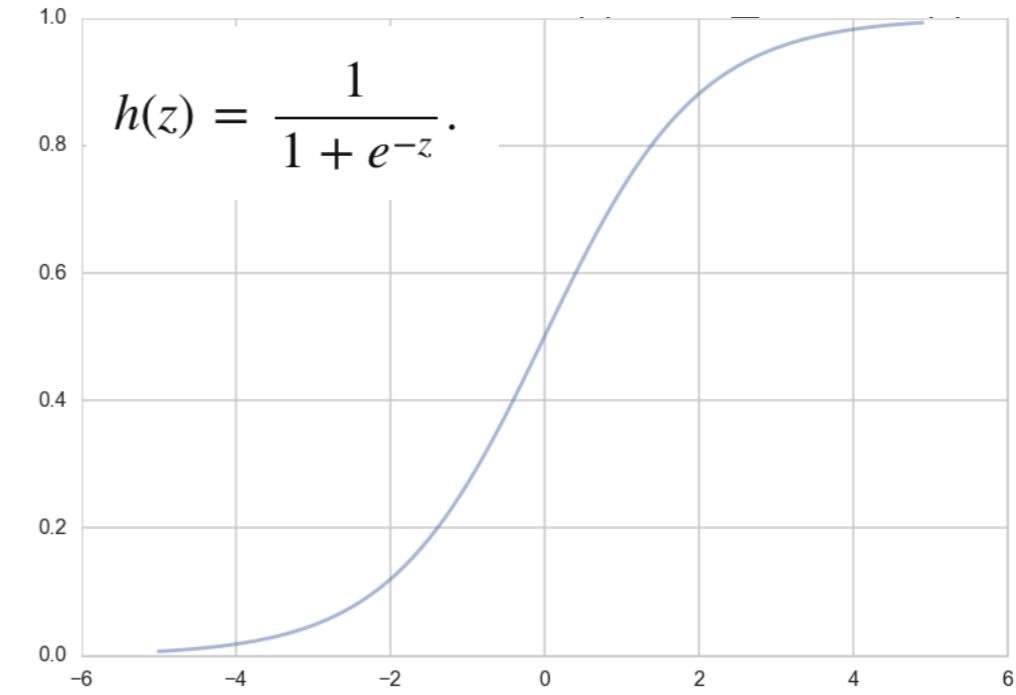
- Simple idea: find a line (surface) that separates (two) classes in the feature space.



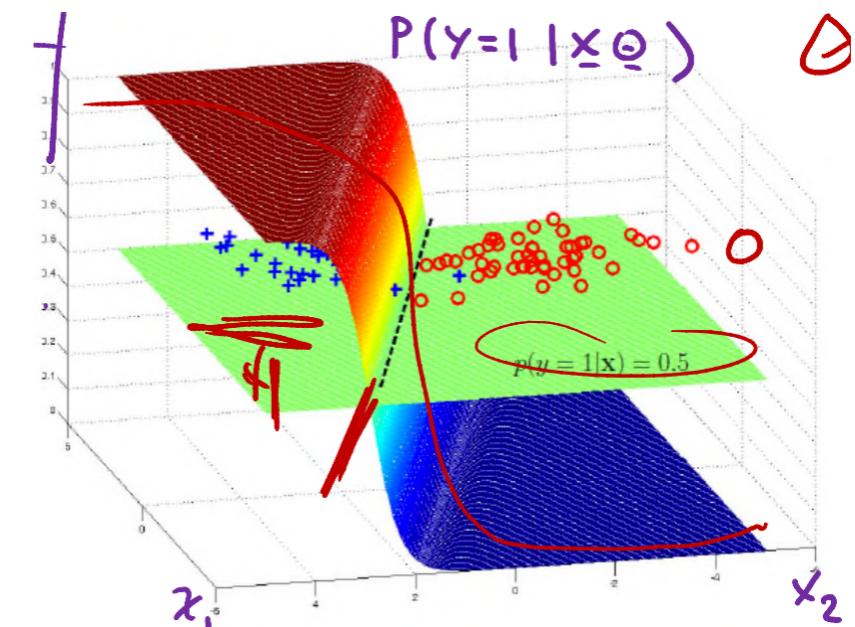
$$P(y|\mathbf{x}, \mathbf{w}) = h(\mathbf{w} \cdot \mathbf{x})^y (1 - h(\mathbf{w} \cdot \mathbf{x}))^{(1-y)}$$

$$P(y|\mathbf{x}, \mathbf{w}) = \prod_{y_i \in \mathcal{D}} h(\mathbf{w} \cdot \mathbf{x}_i)^{y_i} (1 - h(\mathbf{w} \cdot \mathbf{x}_i))^{(1-y_i)}$$

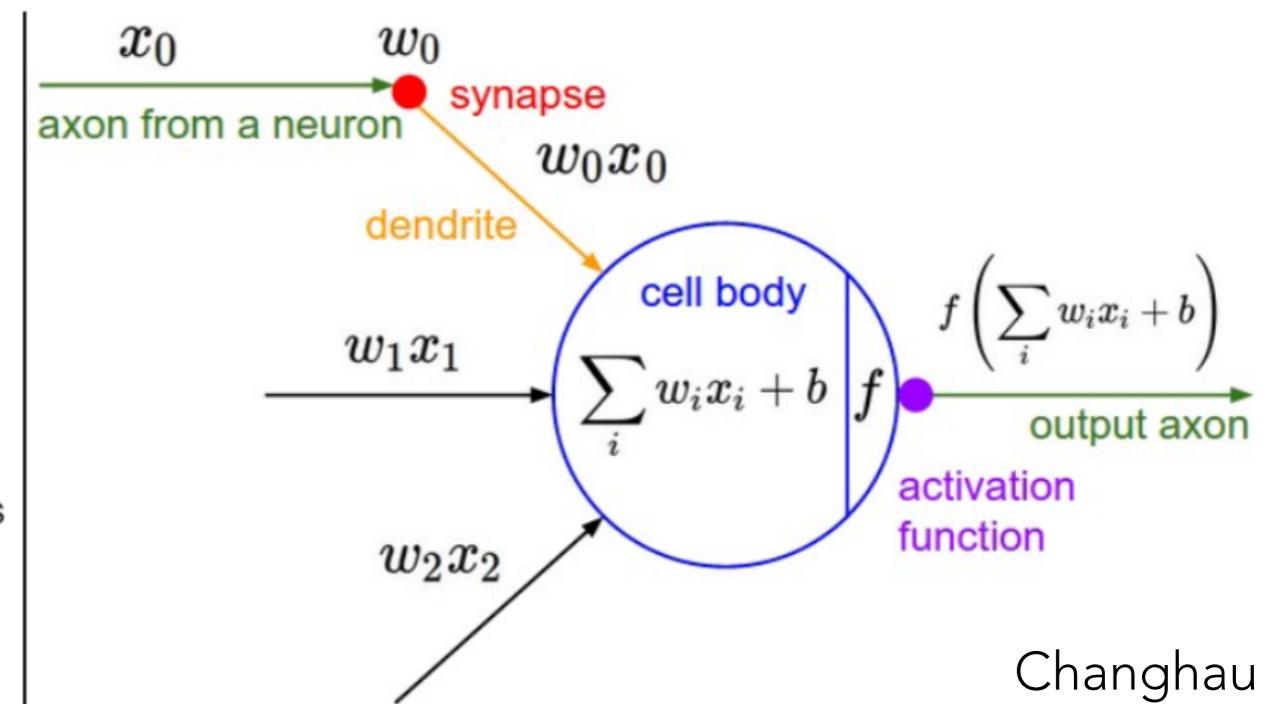
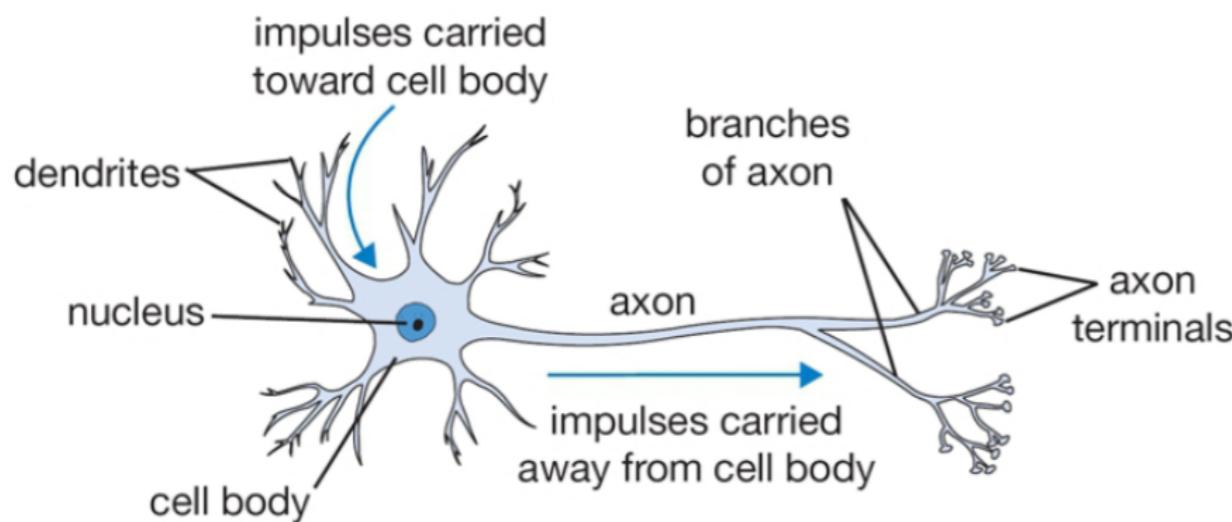
Use training set to learn the values of w that give $P(y|\mathbf{x}, w)$ large when the object is blue and $P(y|\mathbf{x}, w)$ small when the object is red.



- We can define $z = \mathbf{w}^* \mathbf{x}$ to squash its value in the range [0,1] and interpret $h(z)$ as the probability that one particular sample is blue.



ARTIFICIAL NEURONS



Changhau

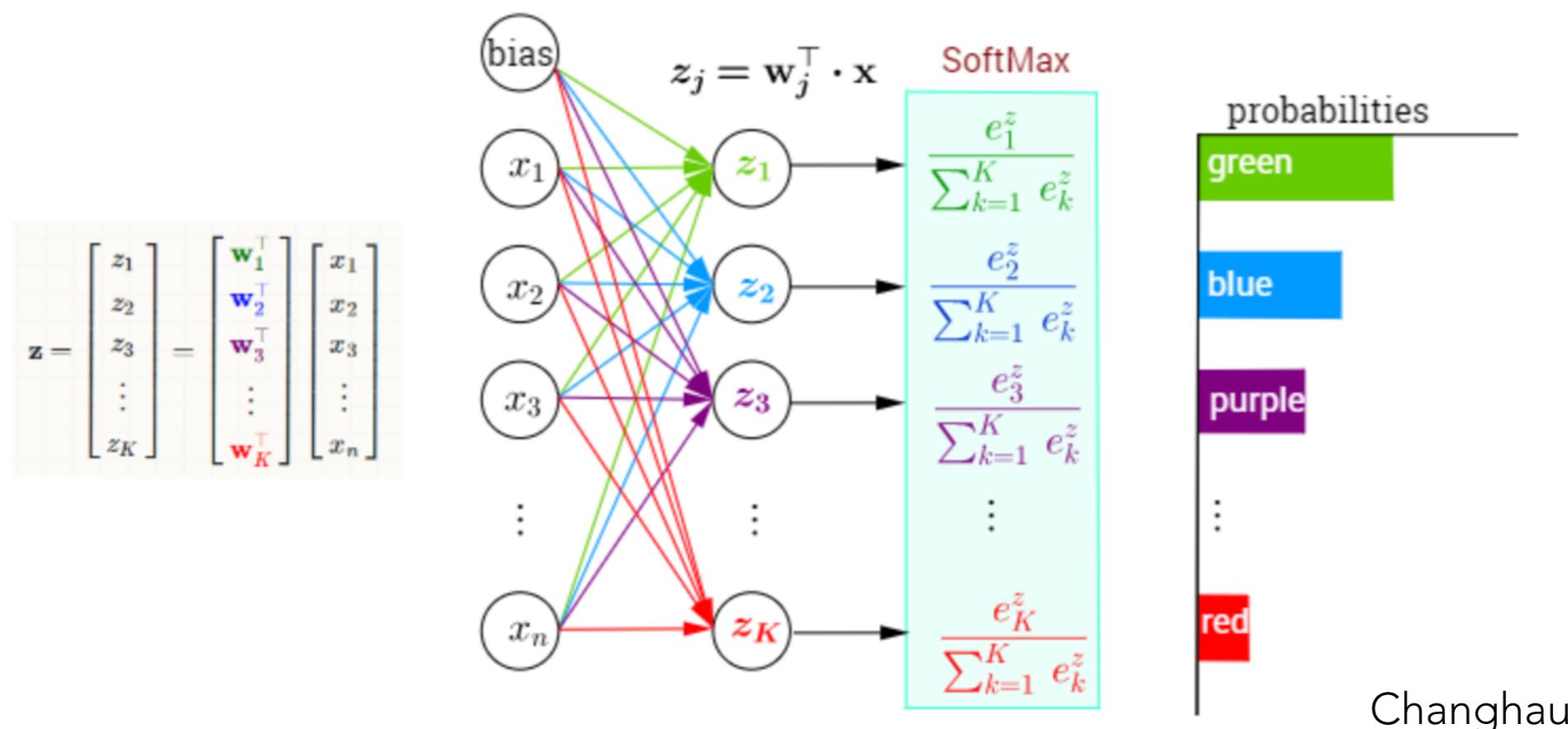
- A non-linear activation function is applied to the linear combination of the features.
- The loss (error) function that we want to minimize is minus the logarithm of the joint probability of the labels observed in the training set:

$$J(\theta) = - \sum_i (y^{(i)} \log(h_\theta(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)})))$$

THE SOFTMAX FORMULATION

- Now suppose that we want to deal with a multi-class problem. The binary nature of the sigmoid function does not (quite) apply anymore.
- We now want to estimate $P(y=k|x)$ for each value of the label k .
- We can use the softmax formulation instead (generalization of the logistic), that squashes the z vector to a K -dimensional vector with values in $[0,1]$

Multi-Class Classification with NN and SoftMax Function



$$J(\theta) = - \left[\sum_{i=1}^m \sum_{k=1}^K 1 \{y^{(i)} = k\} \log \frac{\exp(\theta^{(k)\top} x^{(i)})}{\sum_{j=1}^K \exp(\theta^{(j)\top} x^{(i)})} \right]$$

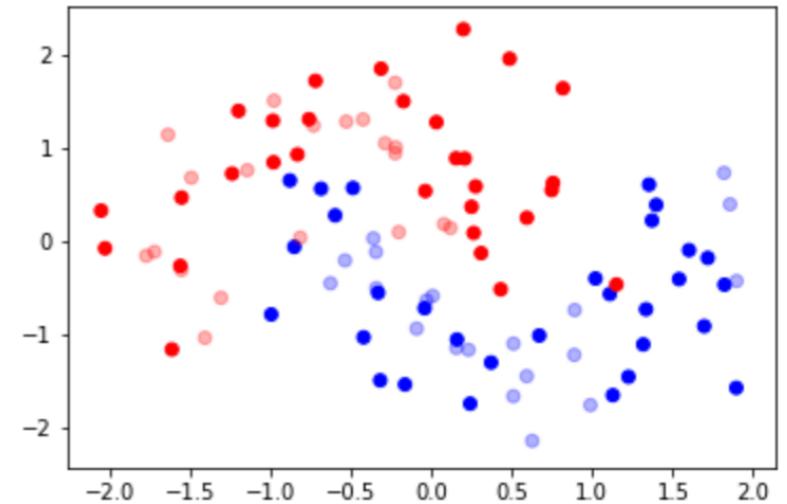
Training, validating, testing

Machine Learning is about minimizing the loss function, which is $-\log(P(y|x,w))$:

$$J(\theta) = - \sum_i (y^{(i)} \log(h_\theta(x^{(i)})) + (1 - y^{(i)}) \log(1 - h_\theta(x^{(i)})))$$

You can think of it as the error in assigning correct classes to the objects in the training set, given their features.

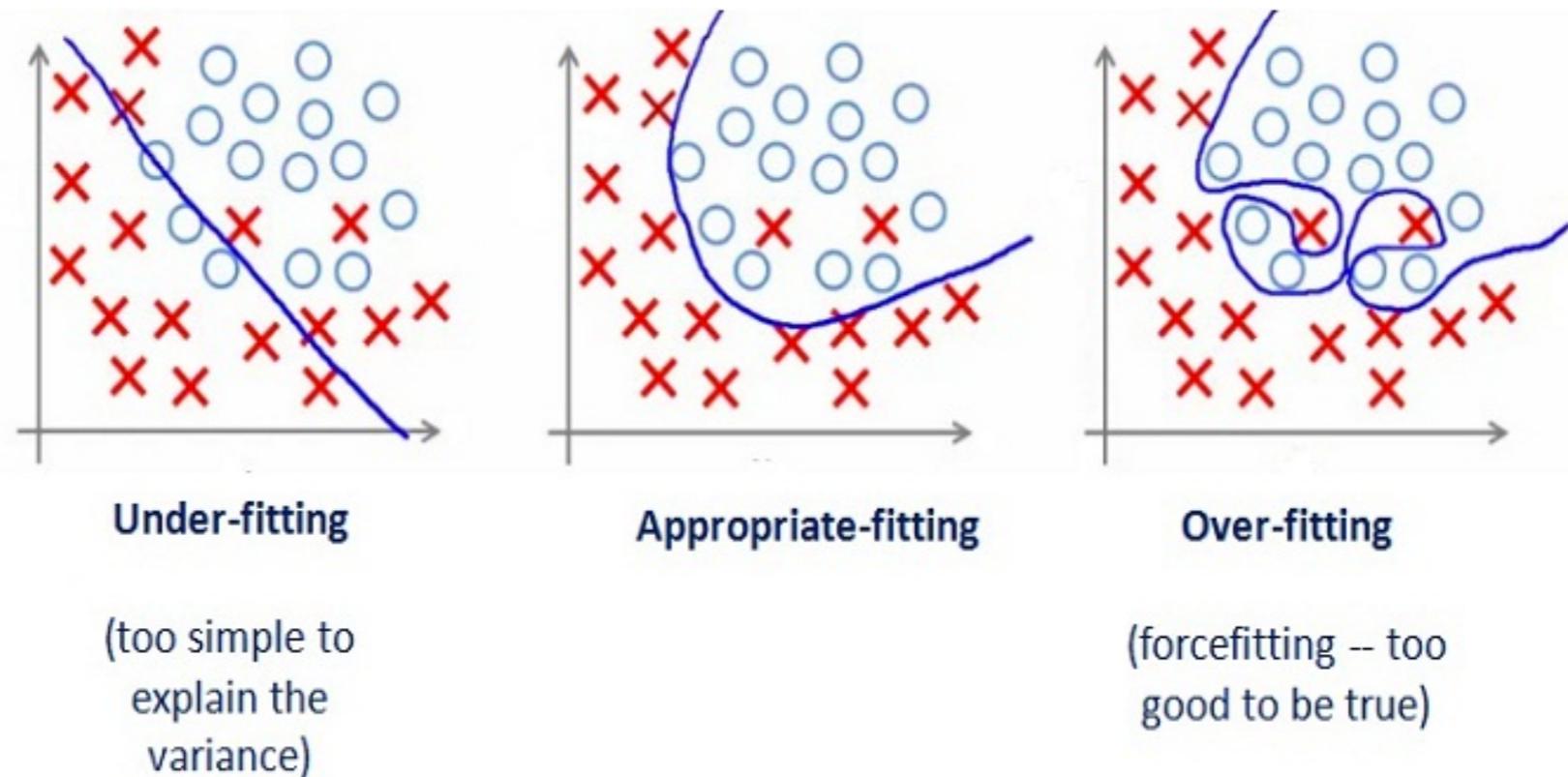
But remember, you want to be able to predict. How do you know that you are doing well in a set of data that the algorithm has not seen before?



That's why you want to break your *labelled* dataset into a training and a testing set.

Training, validating, testing

But what about validation? Well, consider this:

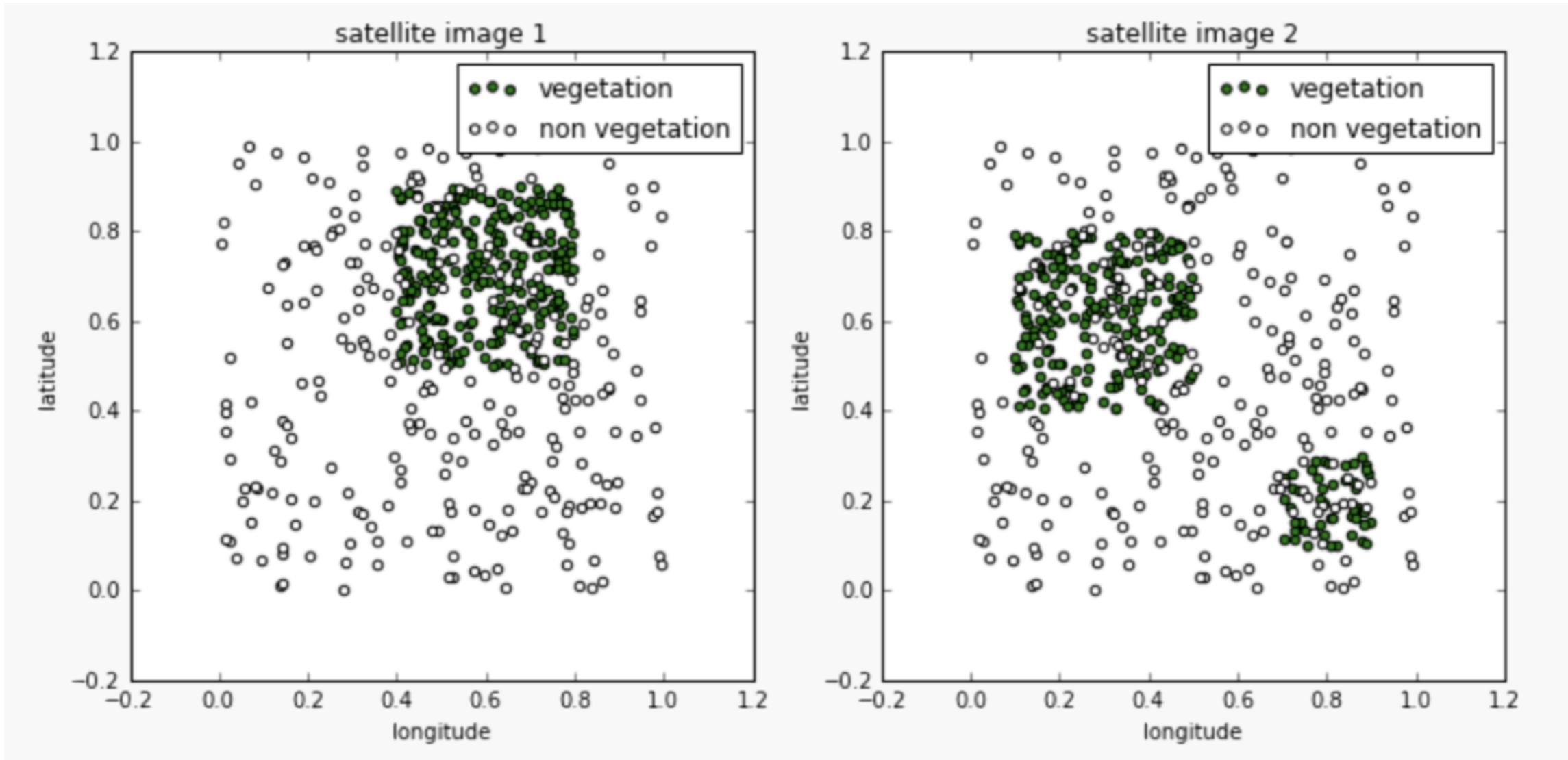


How well should you reproduce the training set without loosing your prediction power for other unseen datasets?

Validation is a way to make sure that you are not **overfitting**

See k-fold cross validation later in the lectures.

CONSIDER GEOMETRY OF DATA



Classes are still well separated in the feature space, but the decision boundaries cannot be described by a simple, analytical expression.

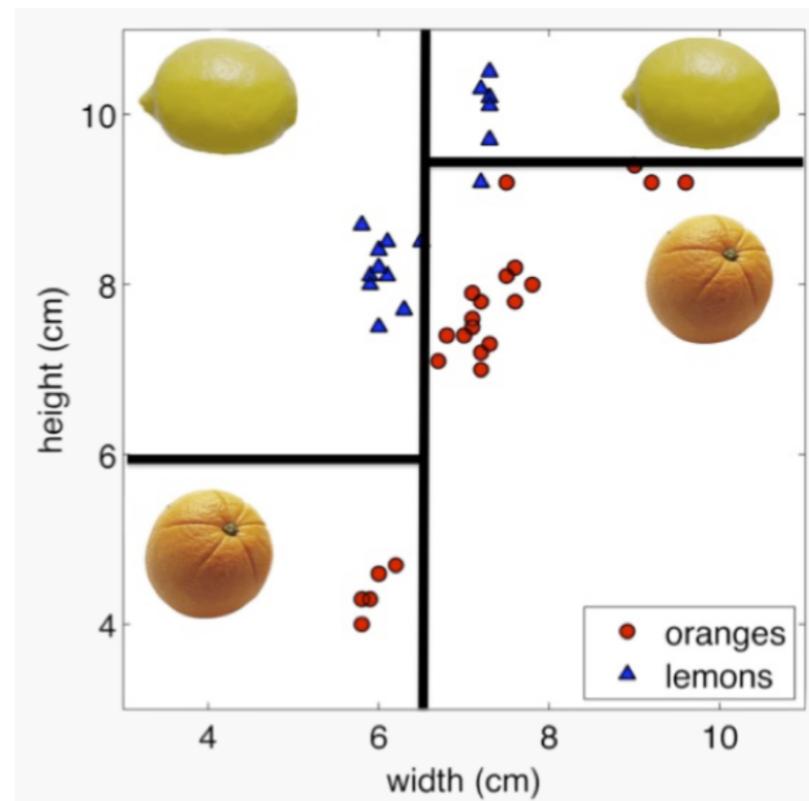
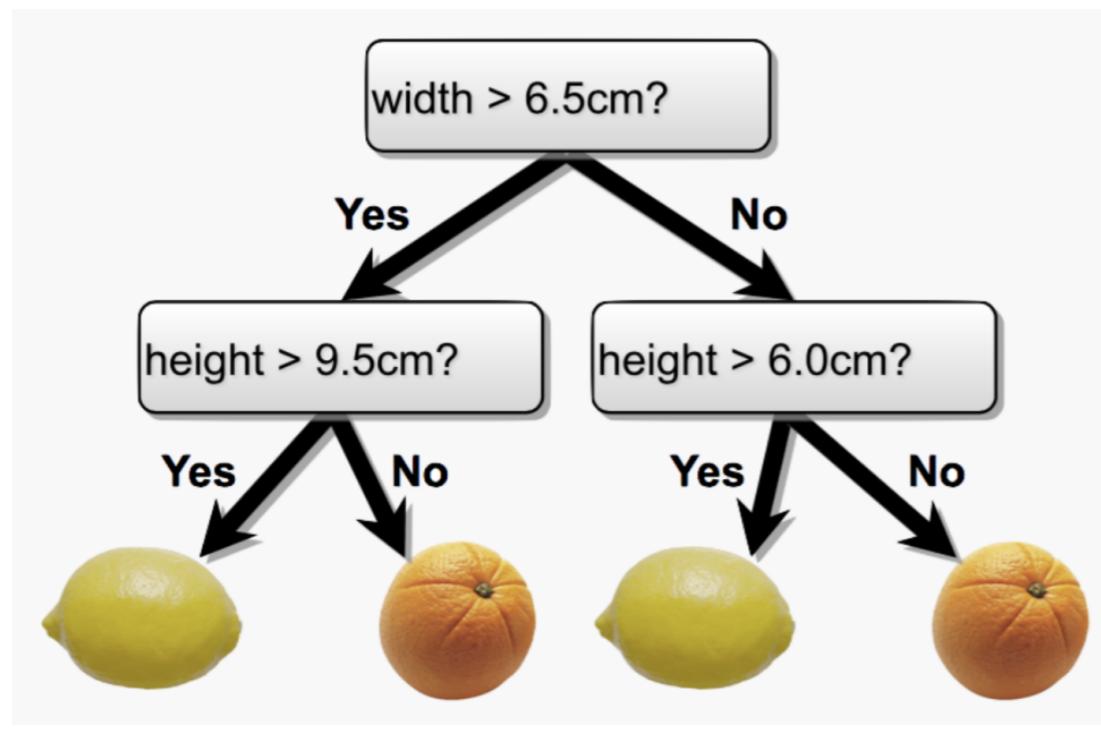
Logistic regression not very useful here.

We need an interpretable model with complex boundaries.

Our brain does that all the time!

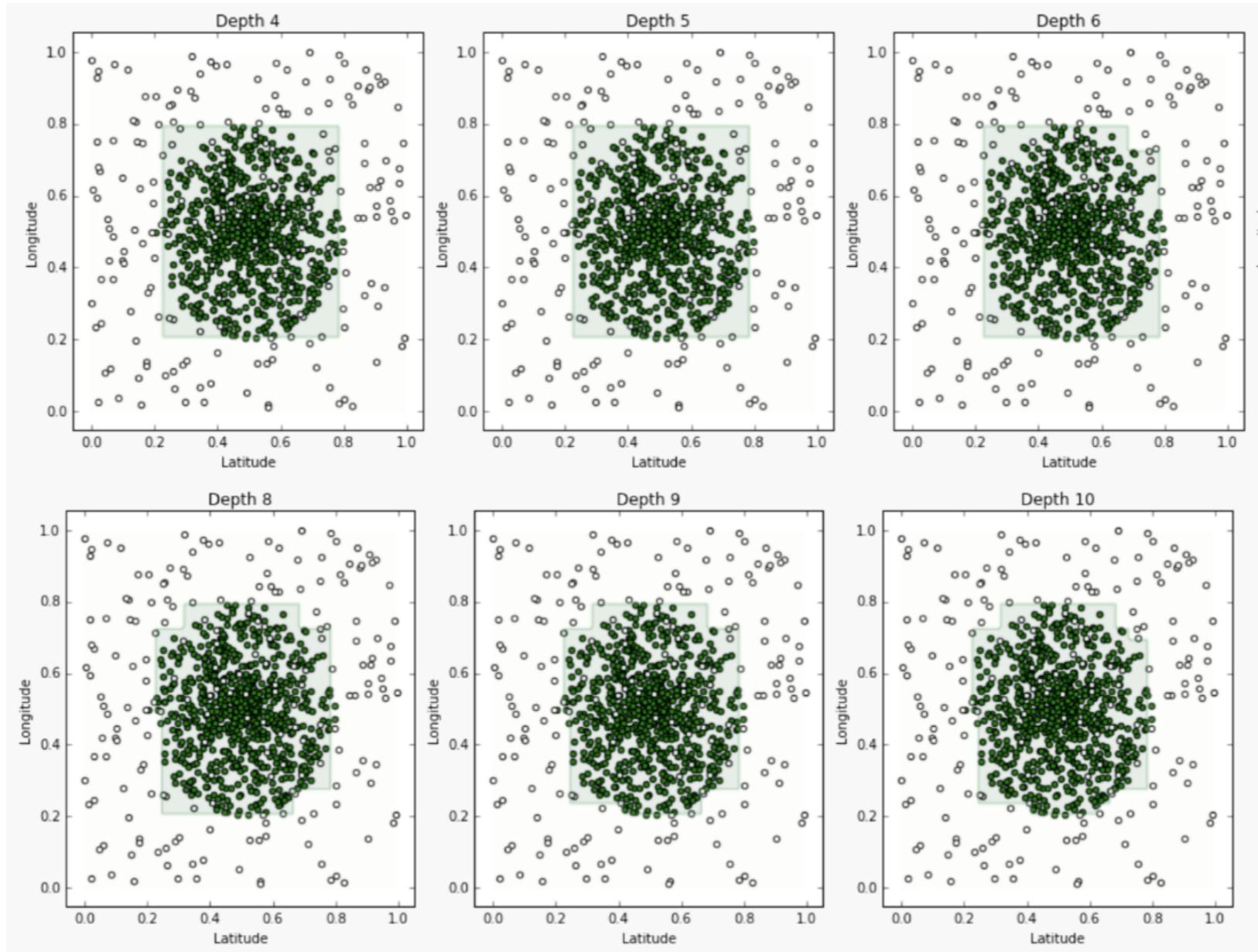
DECISION TREES

- Flow charts whose graph is a tree (connected and no cycles) represents a model called a **decision tree**.
- A decision tree model is one in which the final outcome of the model is based on a series of comparisons of the values of the predictors (features) against threshold values.



- Every tree corresponds to a partition of the feature space by axis-aligned lines or (hyper) planes.
- Given a training set, learning a decision tree means to produce an optimal partition of the features space where each region is given a class label based on the largest class in that region.

BIAS VS VARIANCE



- Bias error: difference between your model's prediction and true values. Algorithms with high bias do not capture complex signals from dataset (under-fitting)
- Variance error: sensitivity to specific sets of training data. Algorithms with high variance will produce drastically different models depending on the training set (over-fitting).

RANDOM FORESTS

- A ***random forest*** is an ensemble of independent decision trees designed to reduce the variance error of individual trees.
- The RF predicts the class of an object as the average prediction of all the trees in the ensemble.
- To de-correlate the trees, we:
 - Train each tree on a separate sub-sample of the training set.
 - For each tree, at each split, we randomly select a subset of predictors from the full set of predictors. From this subset we select the optimal predictor and the optimal threshold.

Scikit-learn & AstroML

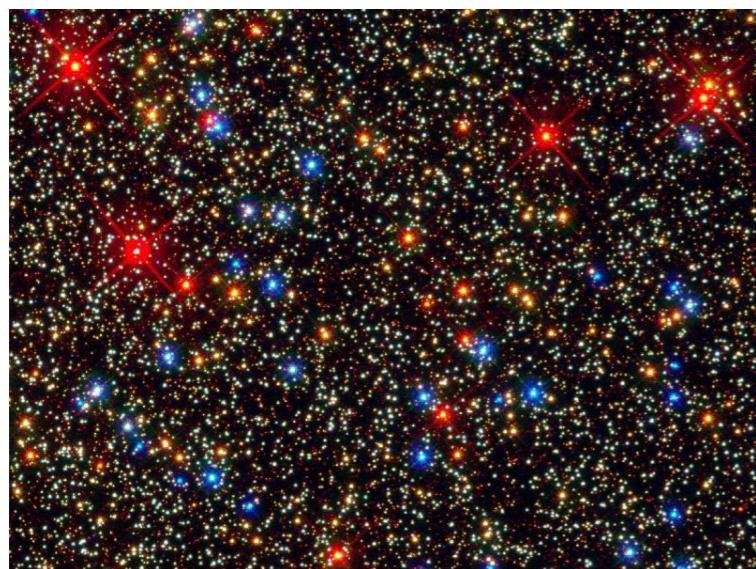
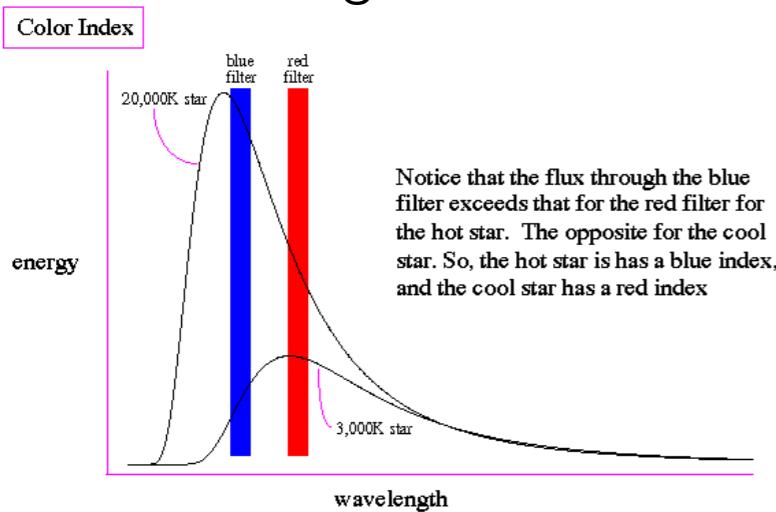
- These are two packages for ML that are very useful in astronomy.
- But you can't just use them if you don't really understand them

Classifying stars based on their color

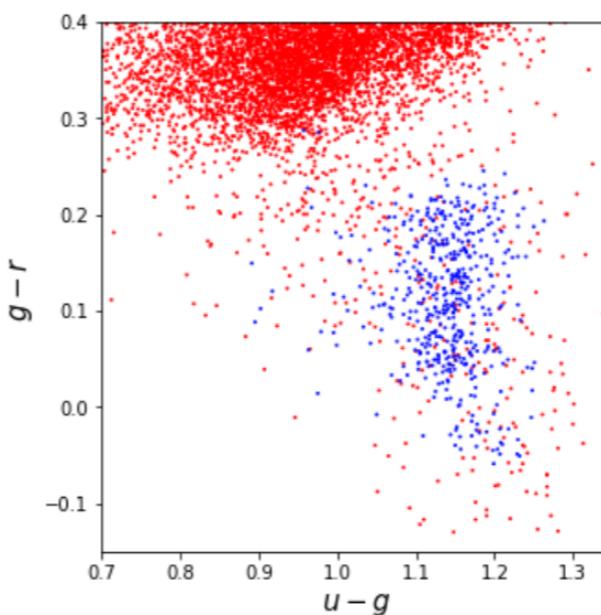
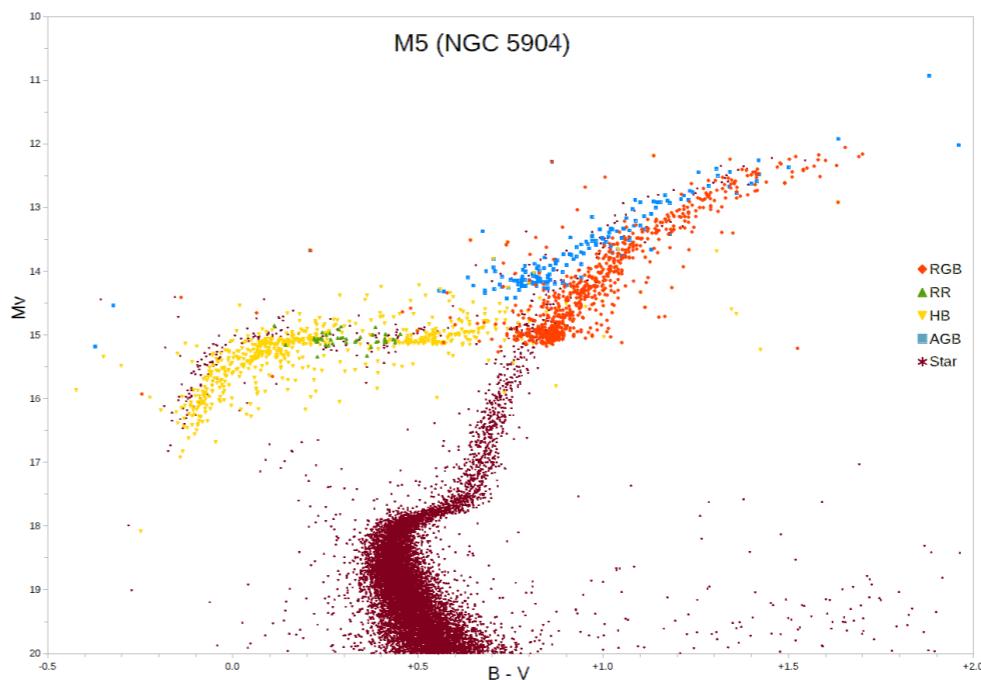
Magnitude: logarithmic measure of brightness

$$m_1 - m_{\text{ref}} = -2.5 \log_{10} \left(\frac{I_1}{I_{\text{ref}}} \right).$$

Color: difference between two magnitudes



RR Lyrae stars are pulsating stars used as standard candles. They have just finished their red giant phase



We can use their optical colors to separate them from other stars. But the diagnostic is not perfect. We need a probabilistic approach: ML!