

Lecture #3

F5611 Machine Learning for Astronomers
by Martin Topinka

<https://github.com/toastmaker/f5611-ML4A>



- 24 Nov 2020 Antonio D'Isanto (MPI, Heidelberg) on redshift estimation
- 15 Dec 2020 Ashish Mahabal (Caltech) on ZTF transient classification

Lecture by *Antonio D'Isanto* (MPI, Heidelberg)

24 Nov 2020

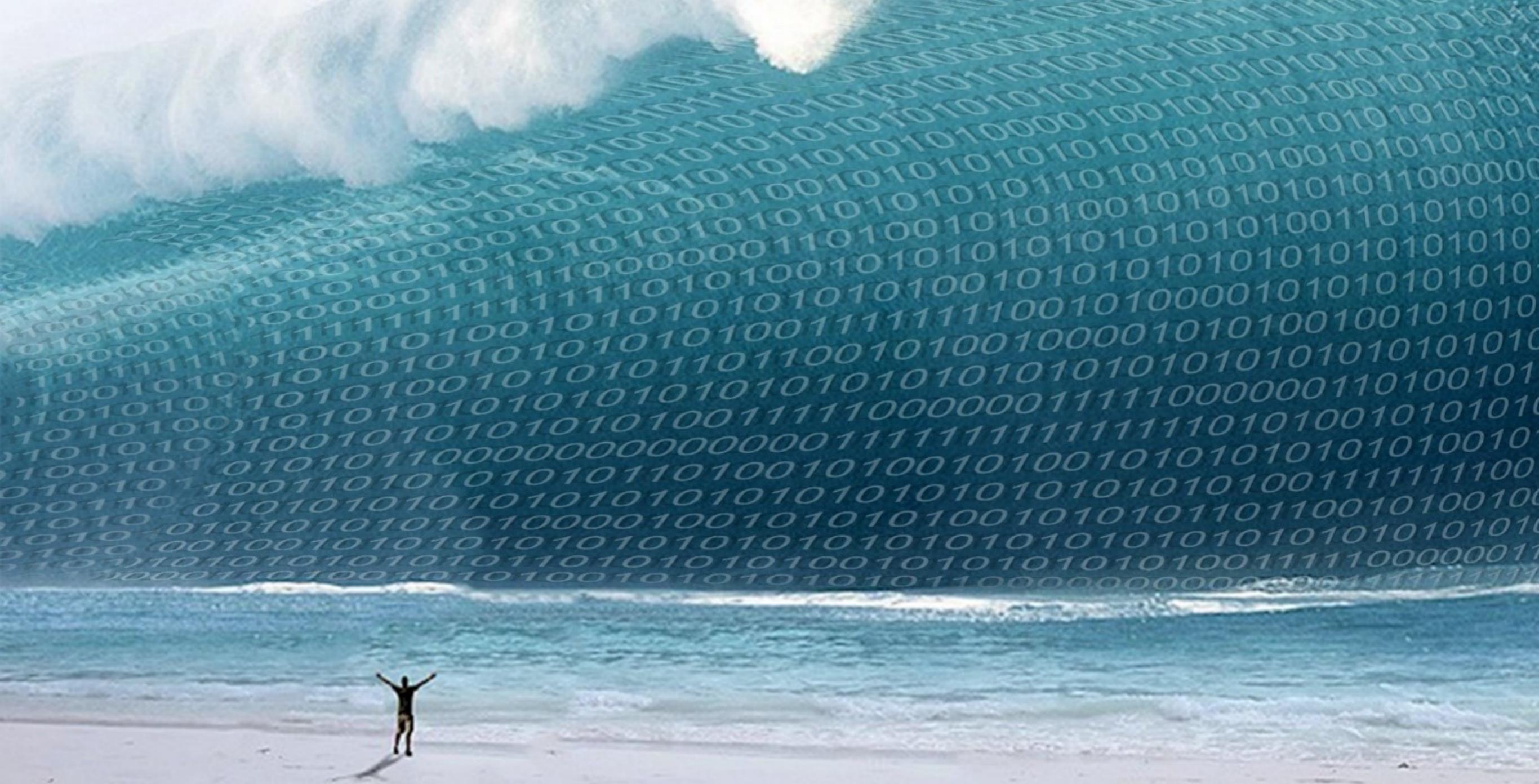
The two worlds of photometric redshift estimation via machine learning: fully automatic vs feature based

The problem of photometric redshift estimation is a major subject in astronomy, since the need of estimating distances for a huge number of sources, as required by the data deluge of the recent years. The ability to estimate redshifts through spectroscopy does not scale with this avalanche of data and photometric redshifts provide the required redshift estimates at the cost of some precision. The success of several forthcoming missions is highly dependent on the availability of photometric redshifts.

Machine learning provides a powerful and efficient solution to the problem of photometric redshift estimation.

In this lecture, I will present two models. The first is fully-automatised, based on the combination of a convolutional neural network with a mixture density network, to predict probabilistic multimodal redshifts directly from images. The second model is features-based, performing a massive combination of photometric parameters to apply a forward selection in a huge feature space. The proposed models perform very efficiently compared to some of the most common models used in the literature. Particular focus is dedicated to the correct estimation of the errors and prediction quality.

The proposed models are very general and can be applied to different topics in astronomy and beyond.



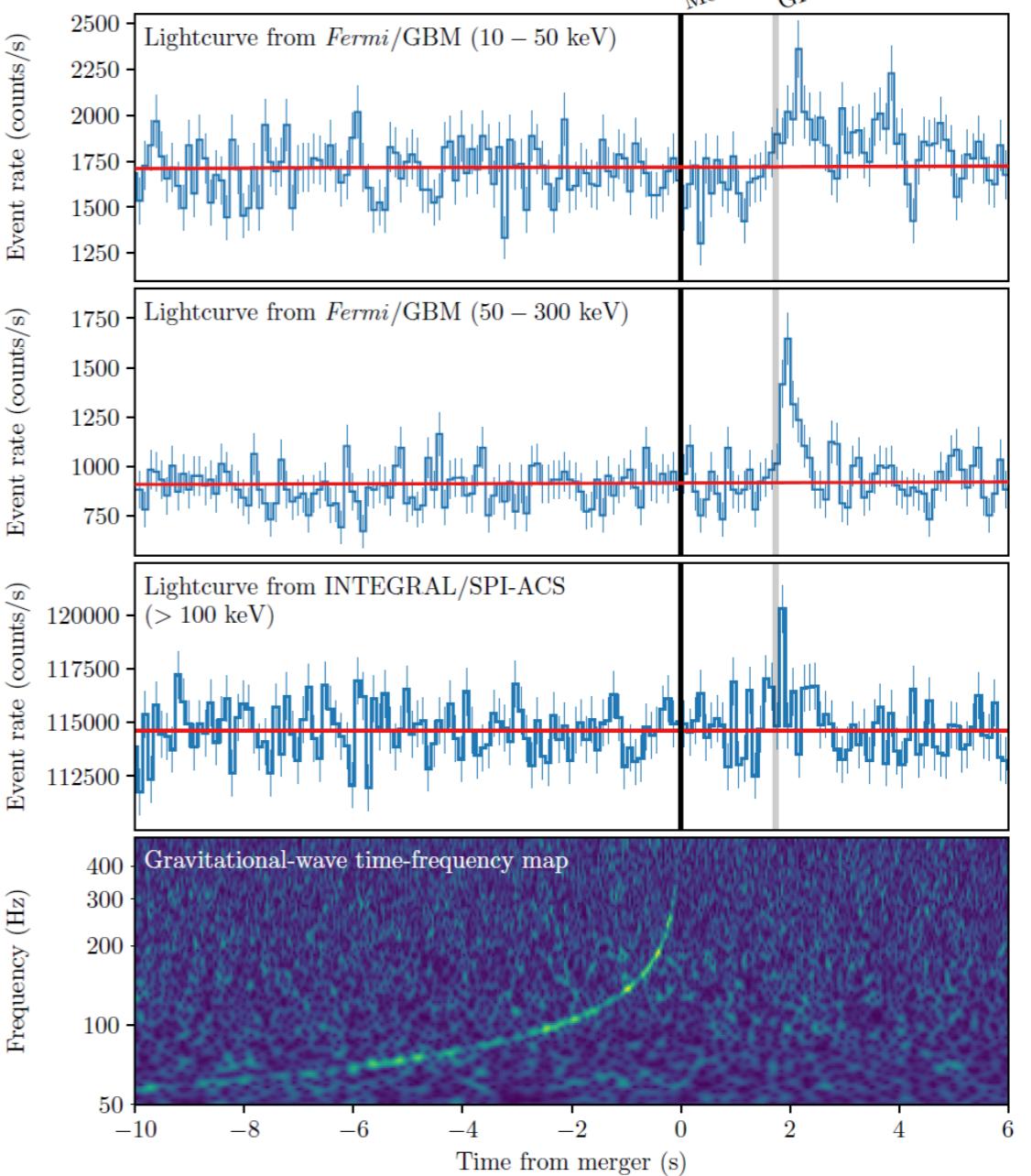
The information volumes and rates grow exponentially ⇒ Most data will never be seen by a human

Increase in the data information content ⇒ Data-driven vs hypothesis driven science

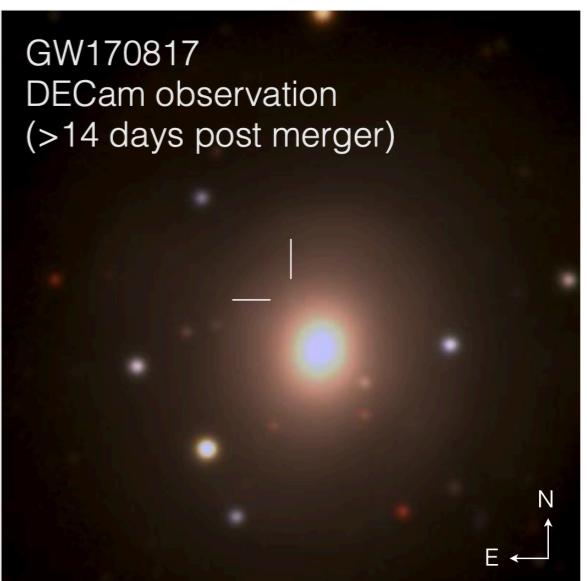
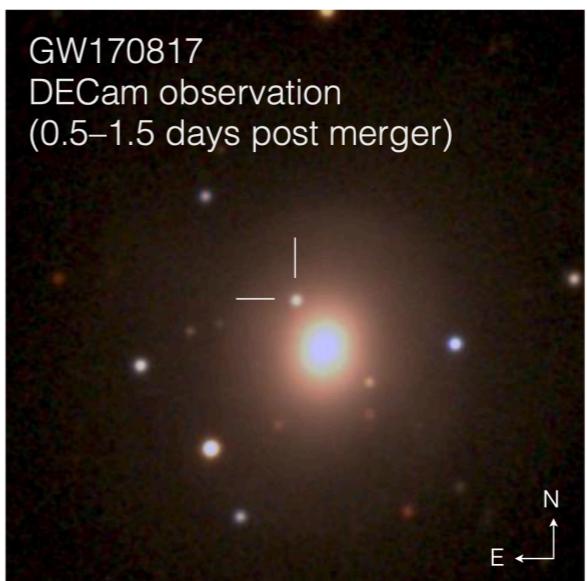
Increase in the information complexity ⇒ There are patterns in the data that cannot be comprehended by human directly

Abbott et al. 2017

Merger
GRB start



Credit: EHT Collaboration



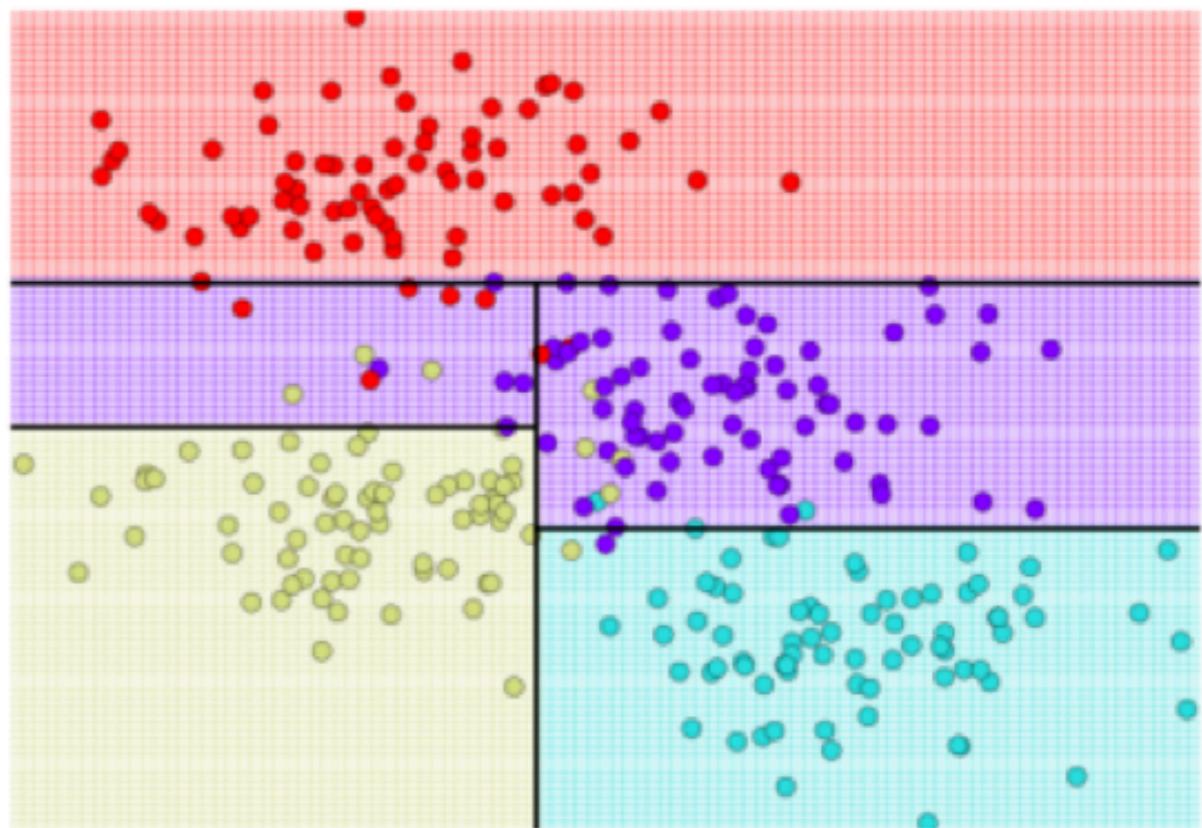
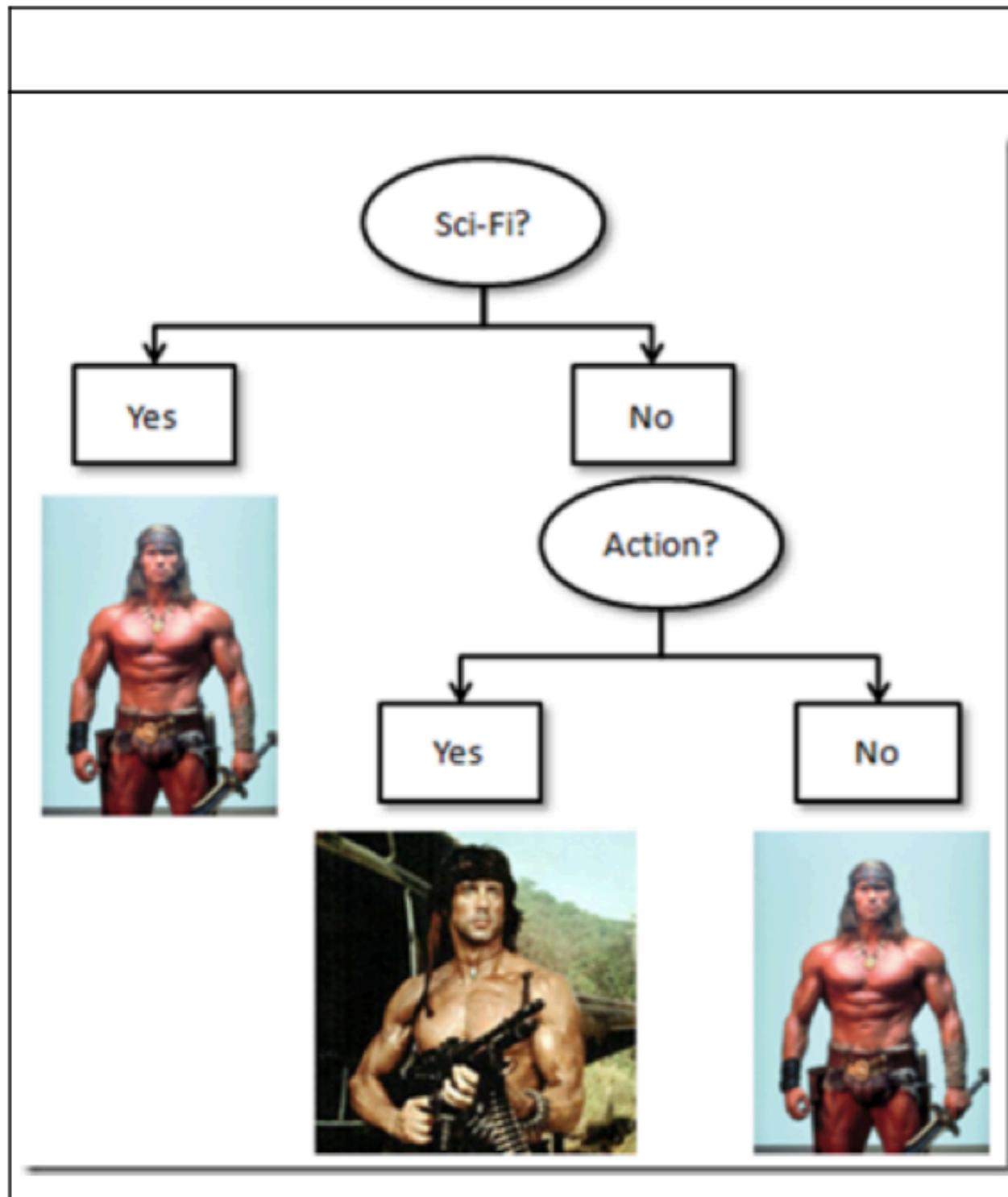
Soares-Santos et al. 2017

More advanced algorithms

- **Ensemble methods:**
 - Bagging (averaging lowers variance) "bootstrapping"
 - Voting (majority wins)
(or train a simple linear regressor of the classifiers)
 - Boosting
- **Neural networks**



Random Forest (Ensemble of Decision Trees)

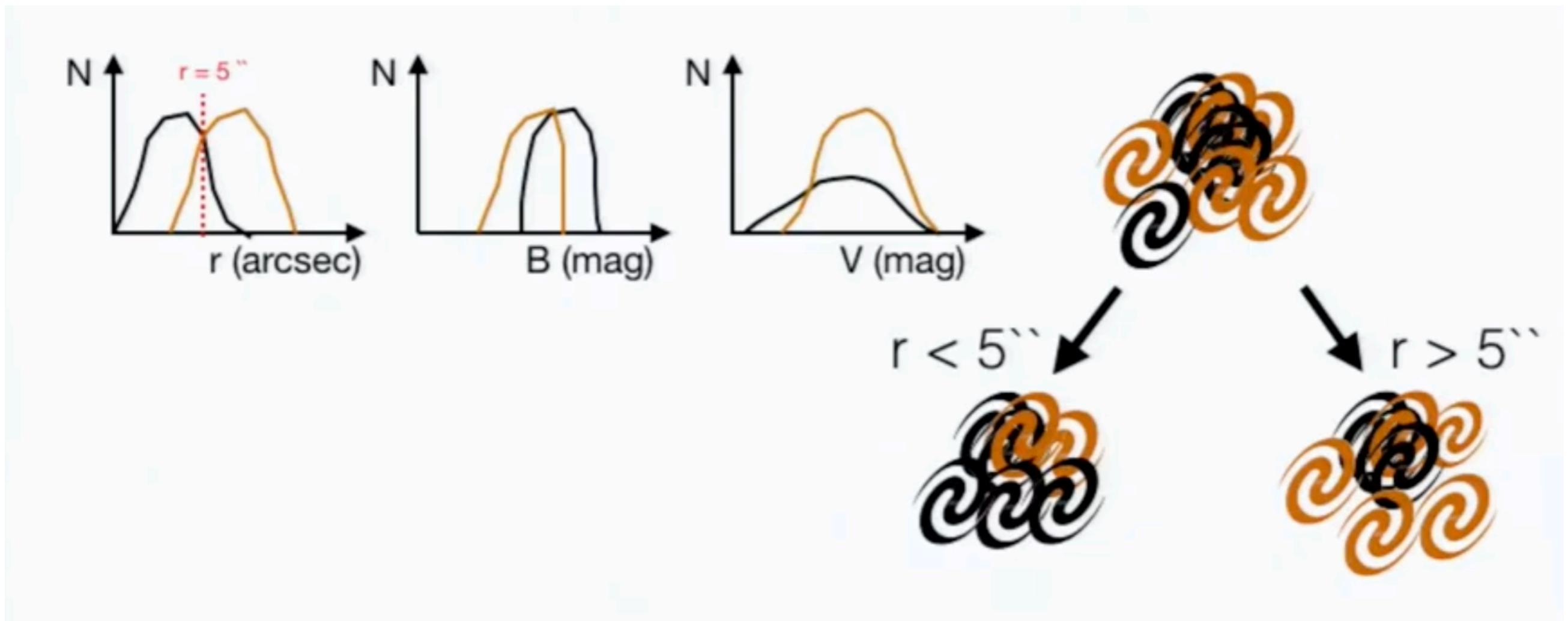


Construction of Decision Tree

Input training set: a list of galaxies

Classes: "black" galaxies and "brown" galaxies

Features: r (arcsec), B (mag), V (mag)



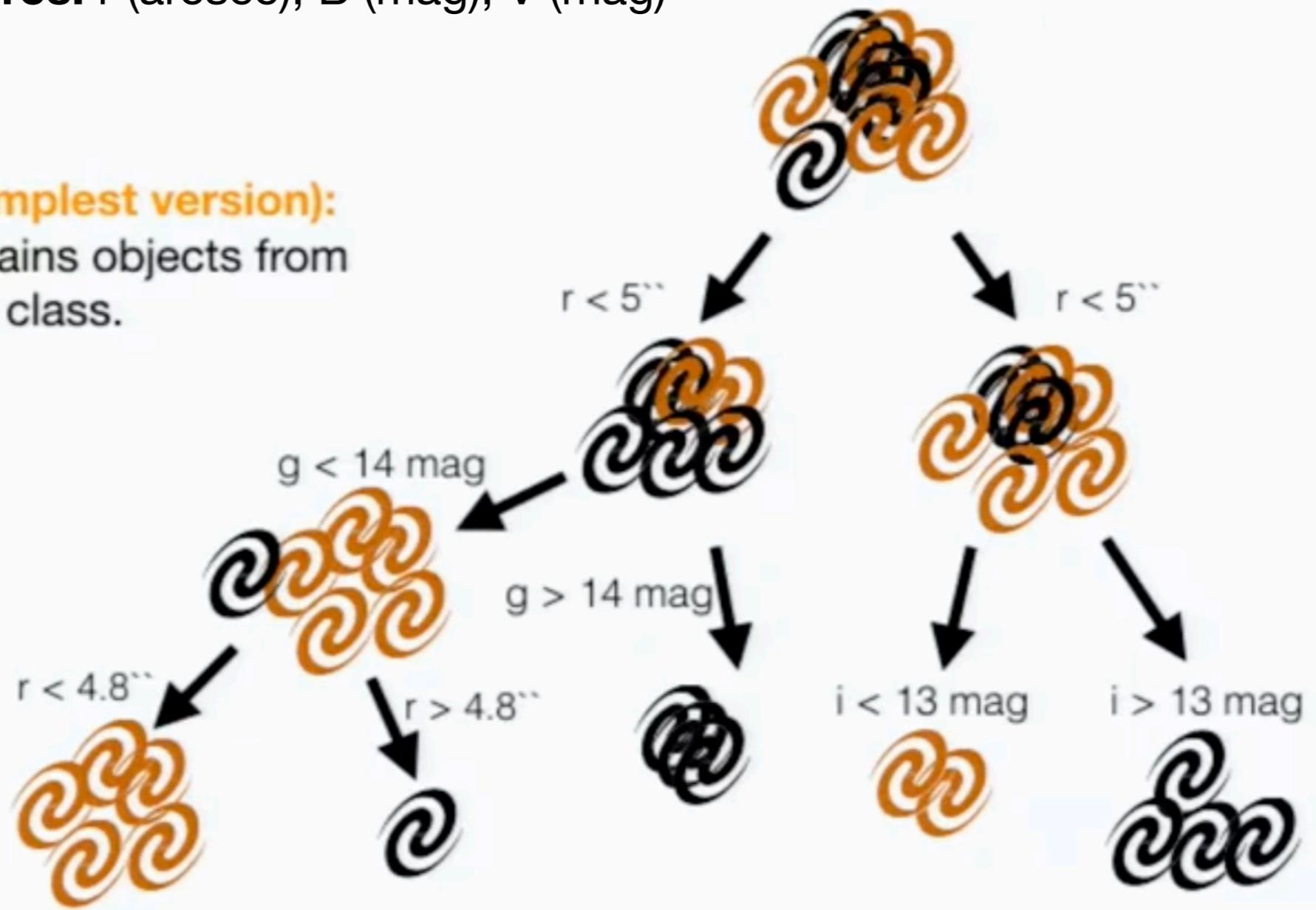
Construction of Decision Tree

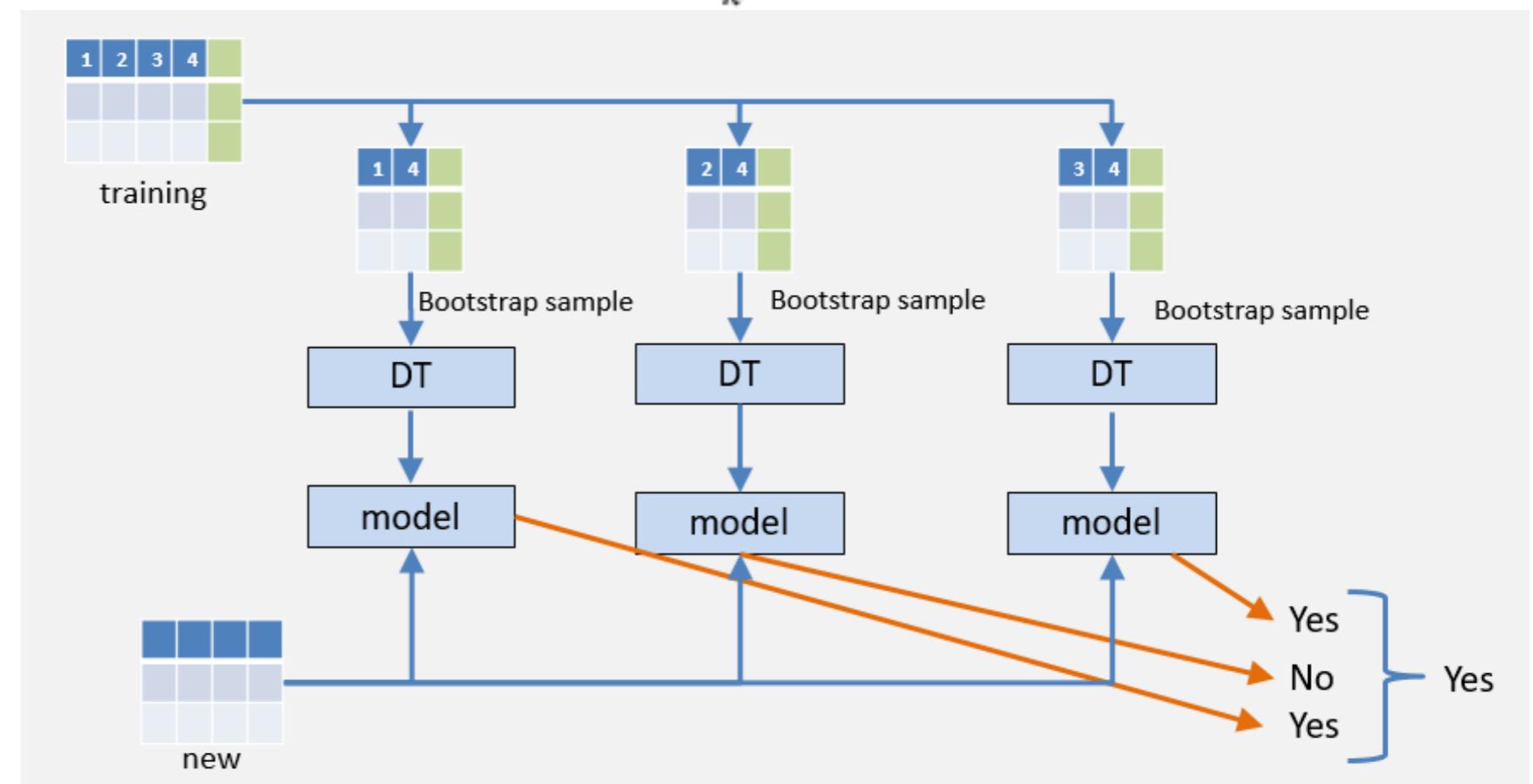
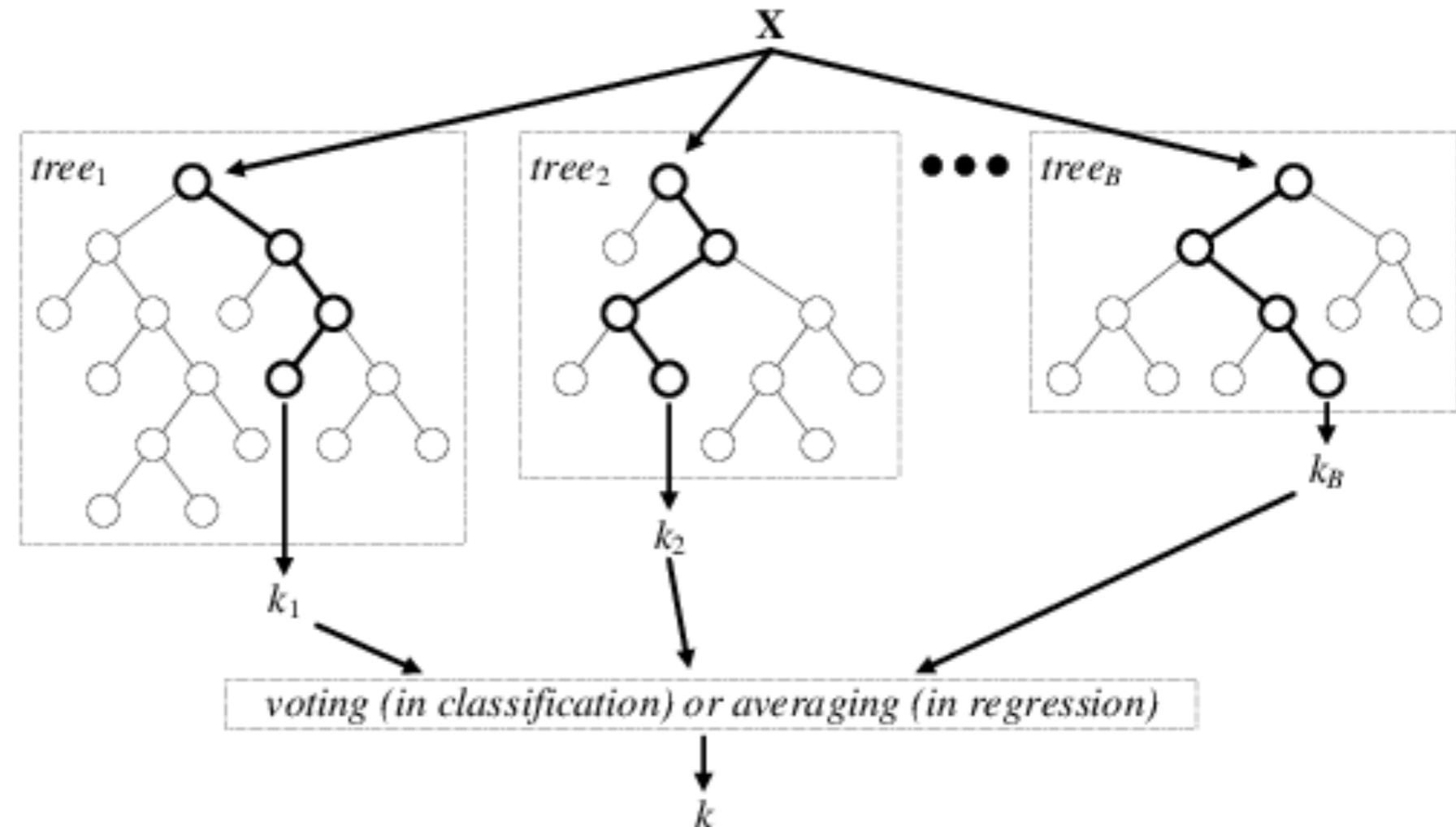
Input training set: a list of galaxies

Classes: "black" galaxies and "brown" galaxies

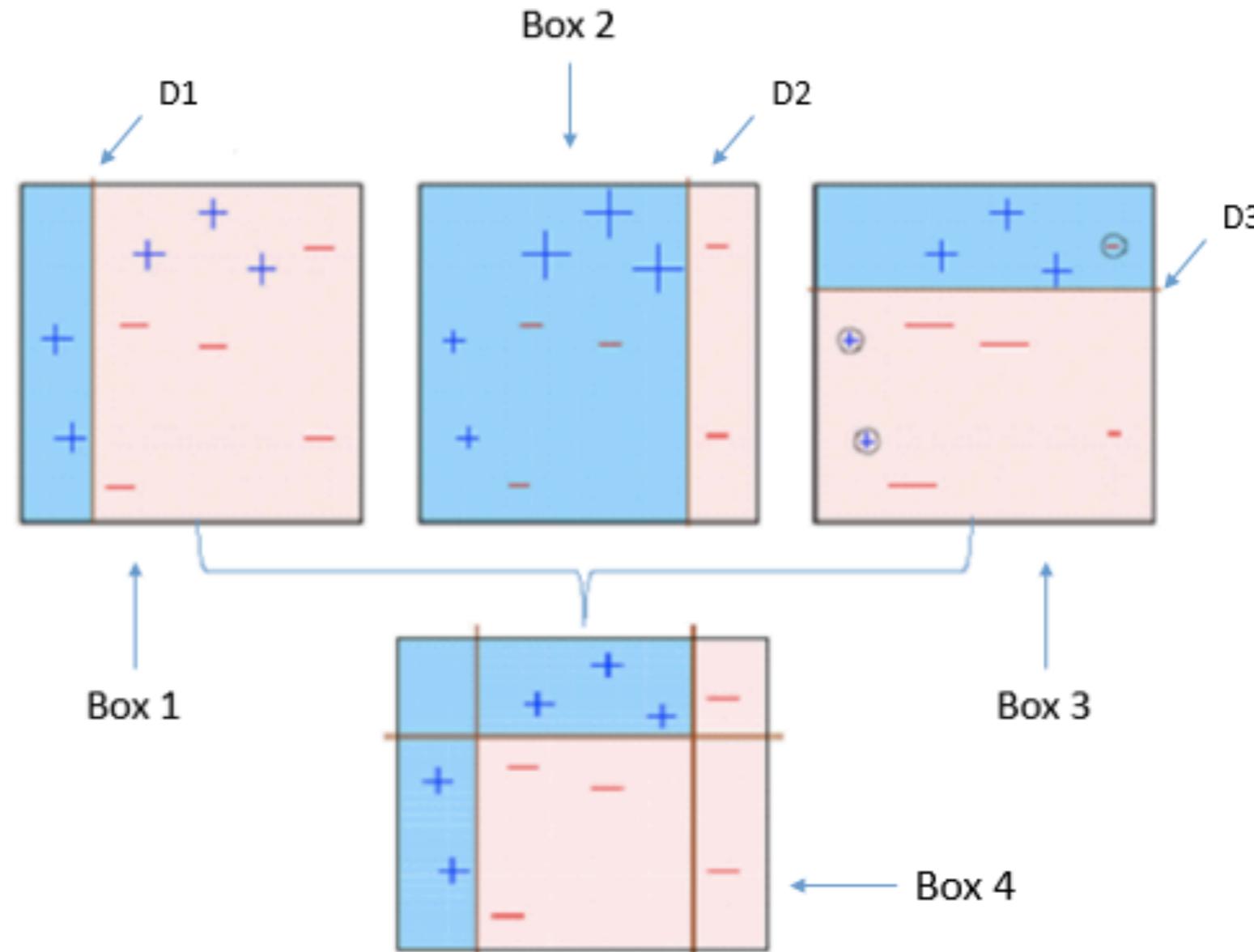
Features: r (arcsec), B (mag), V (mag)

Stop criterion (simplest version):
each terminal contains objects from
a single class.





Voting (max/average) with weights (gradient descent)

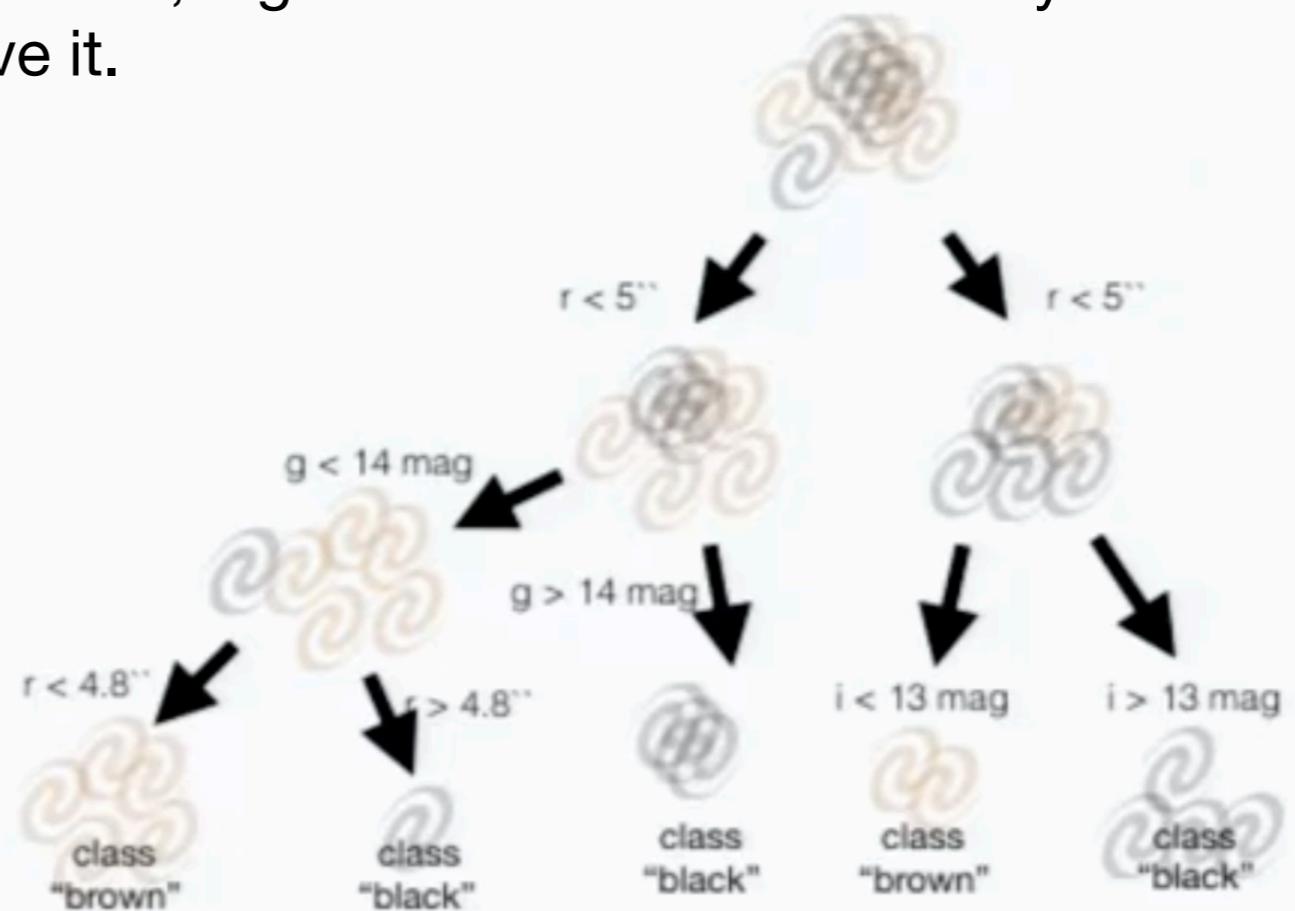


Feature importance & feature selection

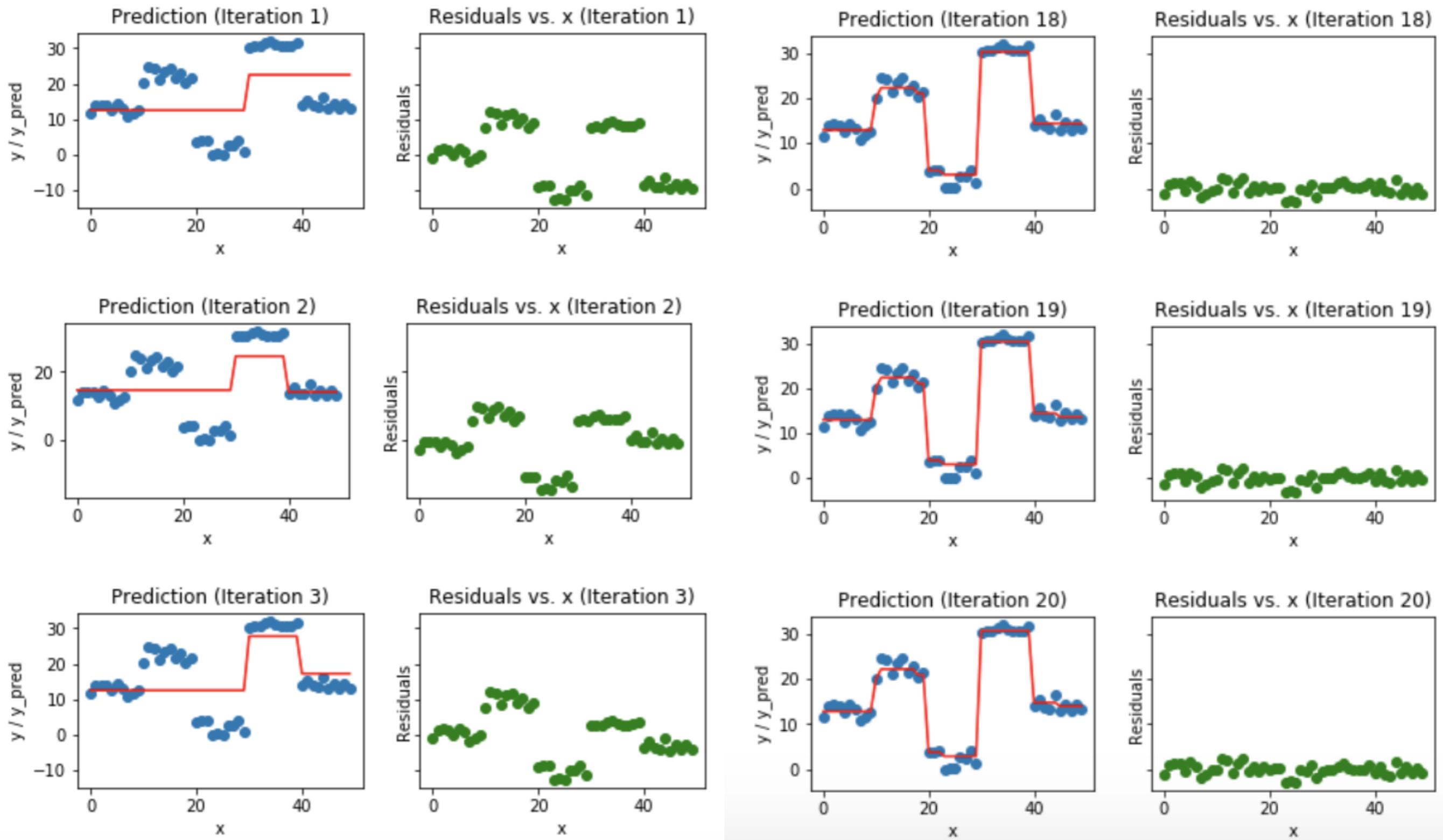
DT: The higher a feature is in the tree, the more important it is for classification.

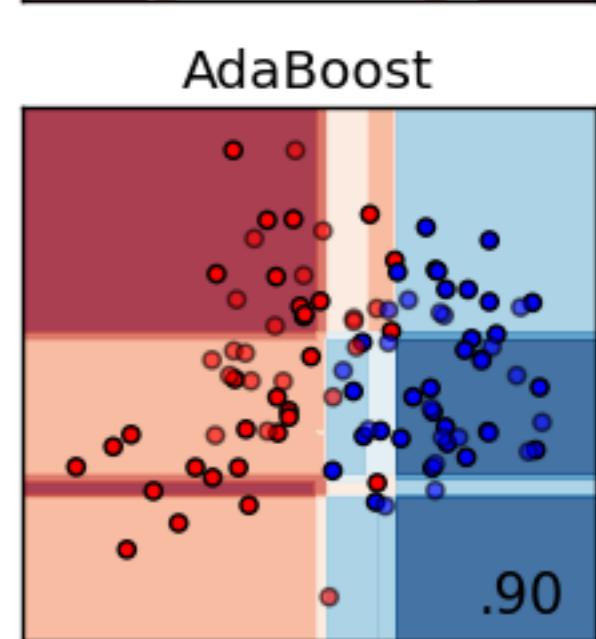
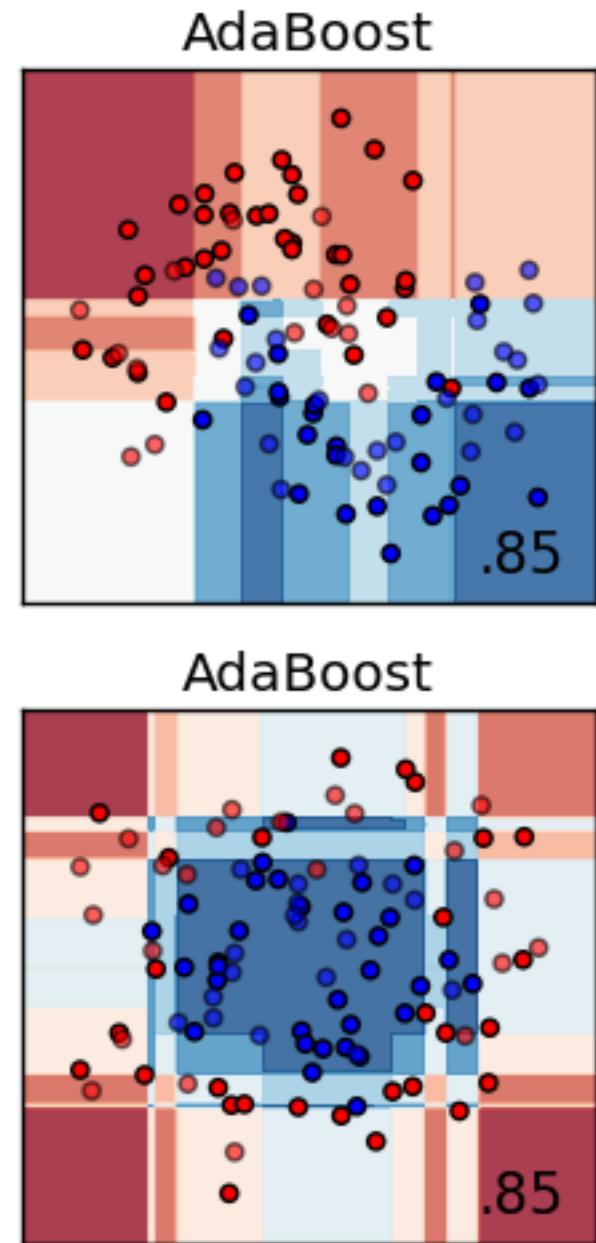
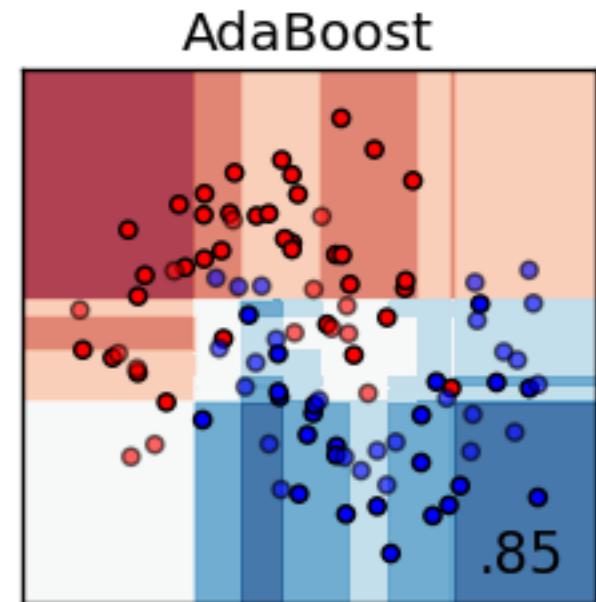
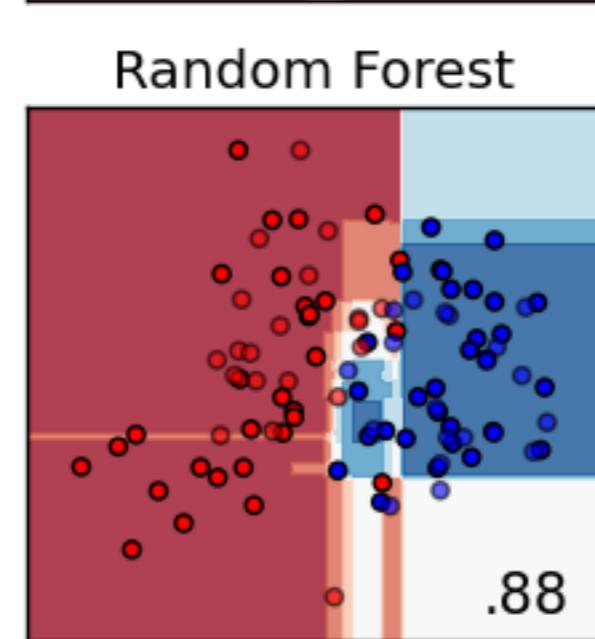
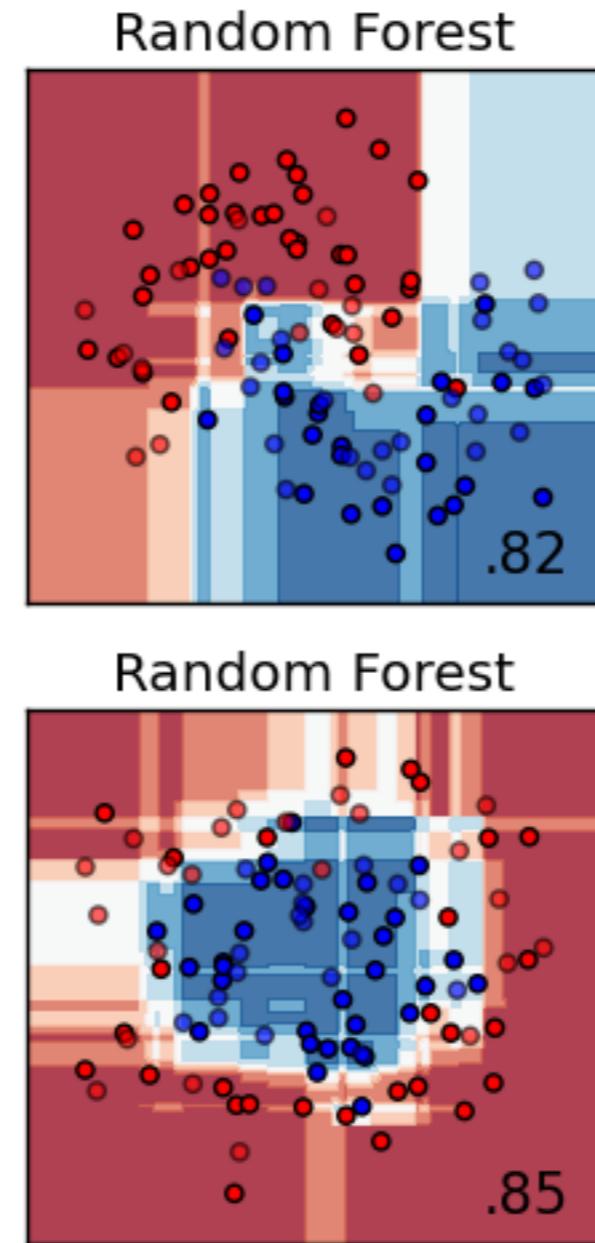
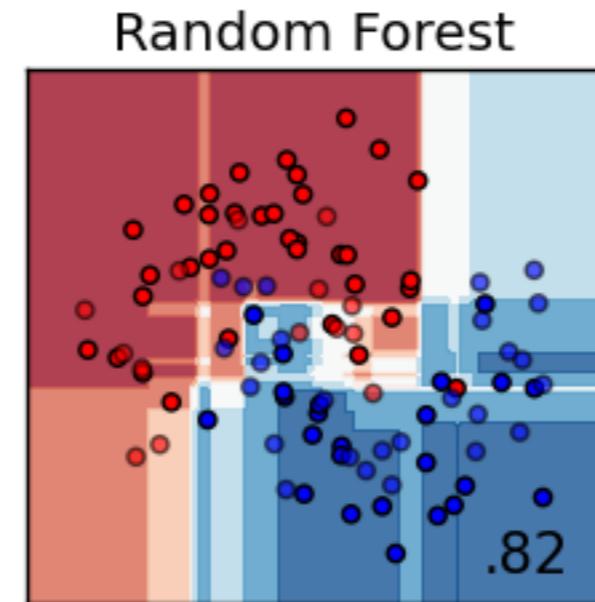
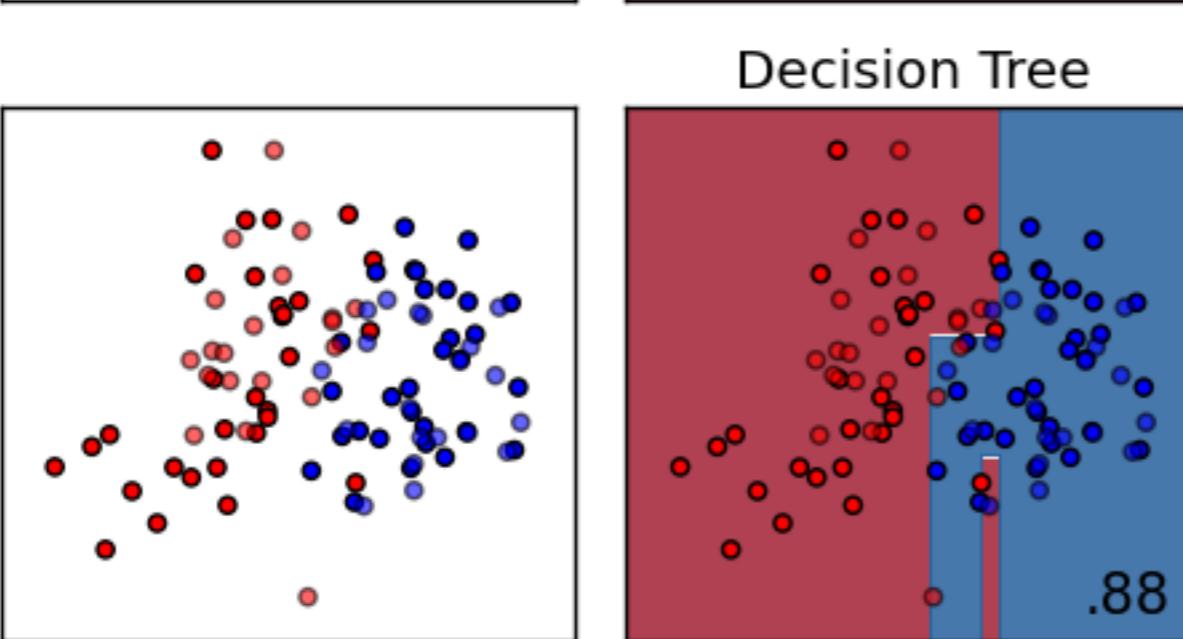
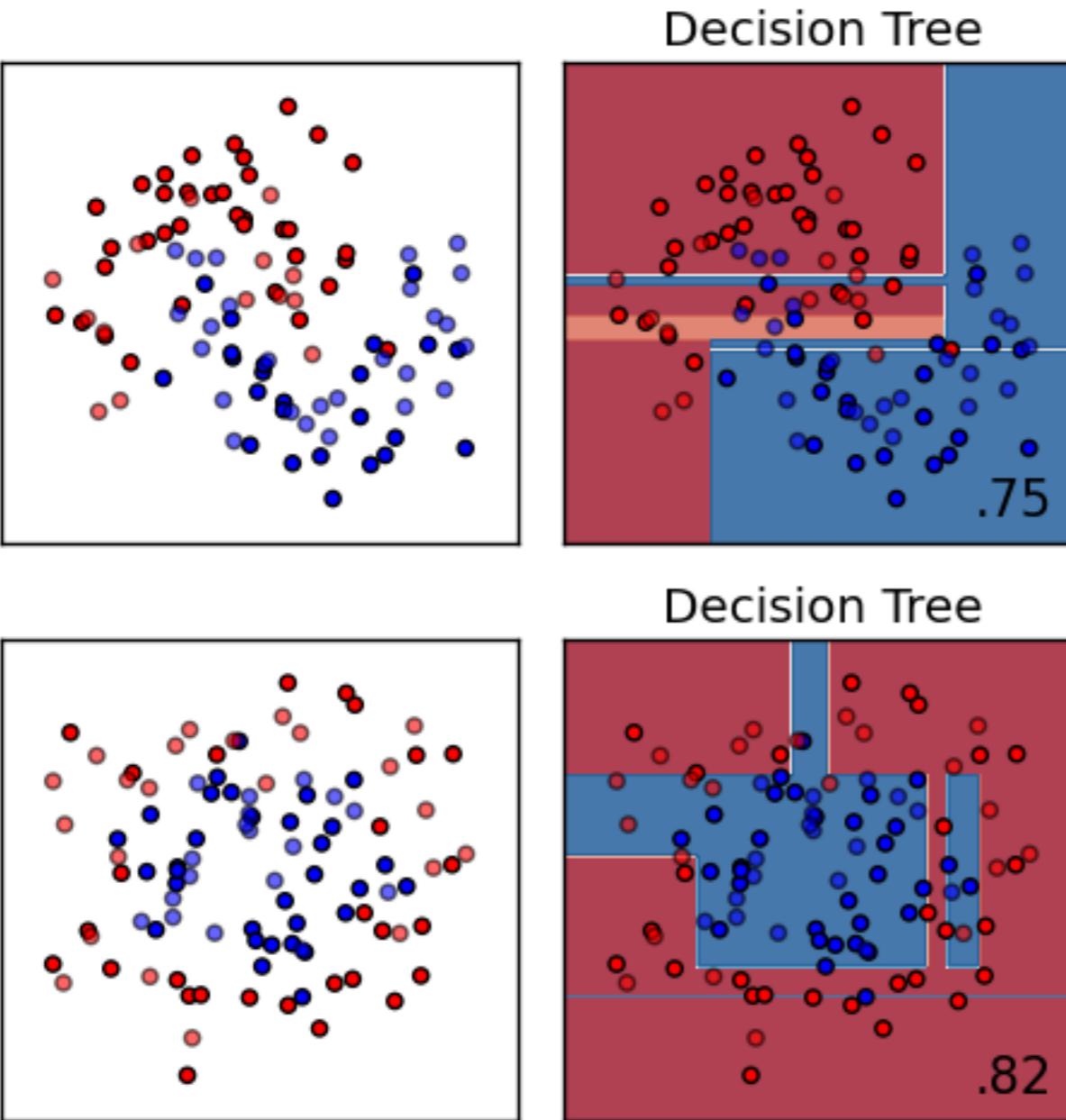
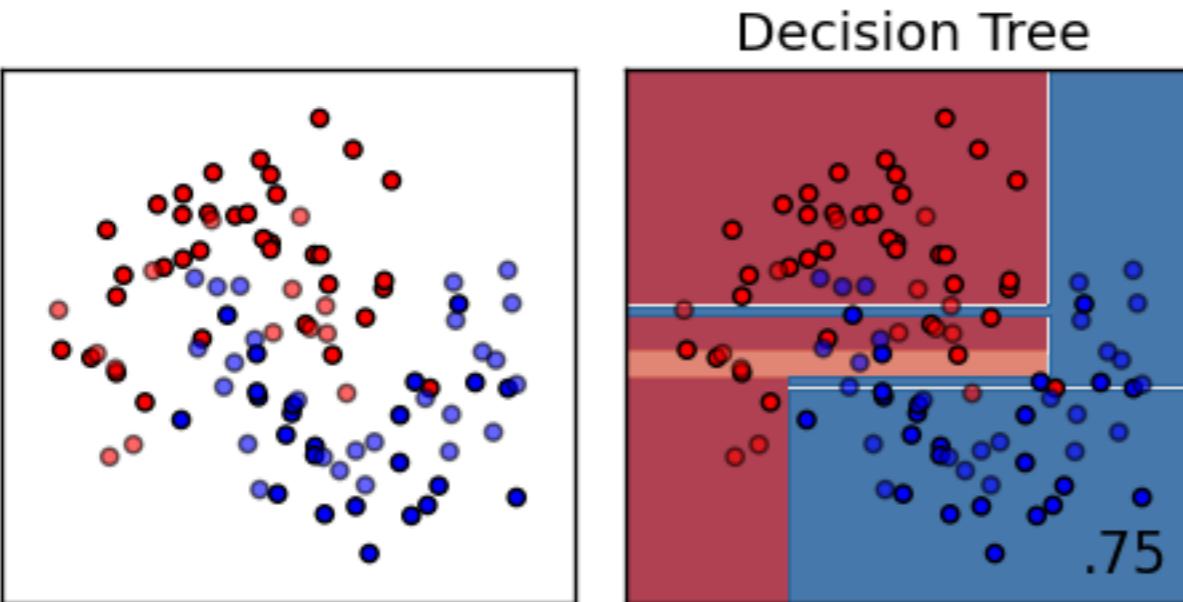
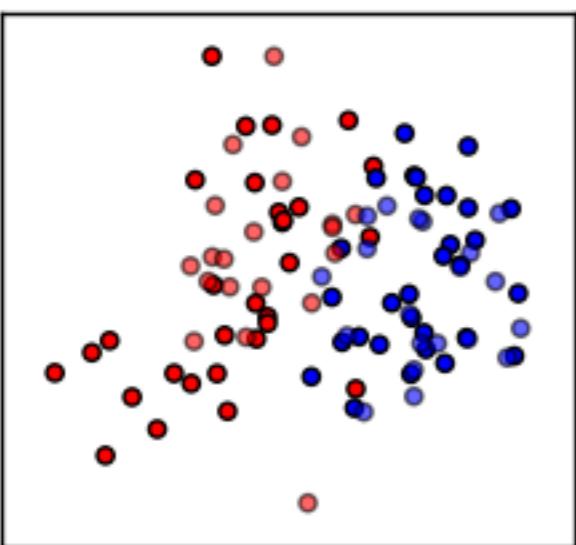
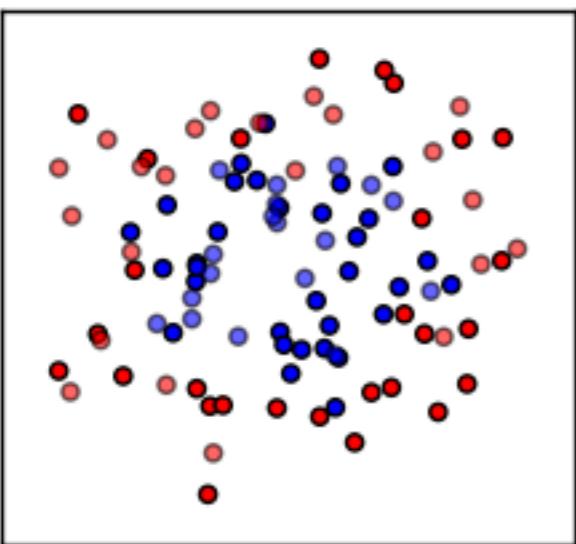
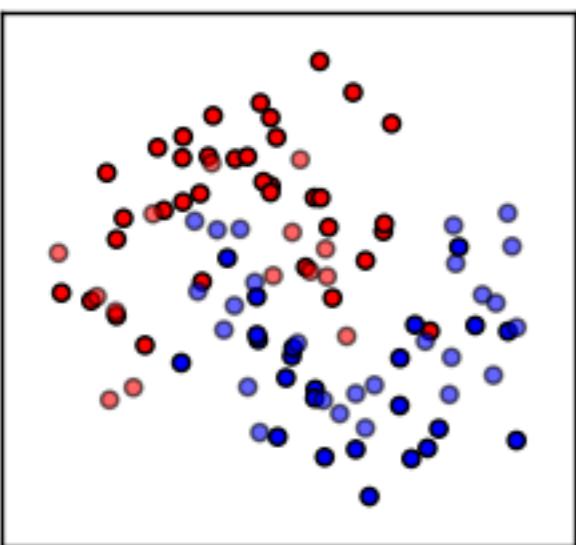
RF: Omitting a feature makes separates data faster, feature has low importance and vice versa.

Trick: add a non-informative dummy feature, e.g. random or a constant. If your physical feature is ranked lower, remove it.

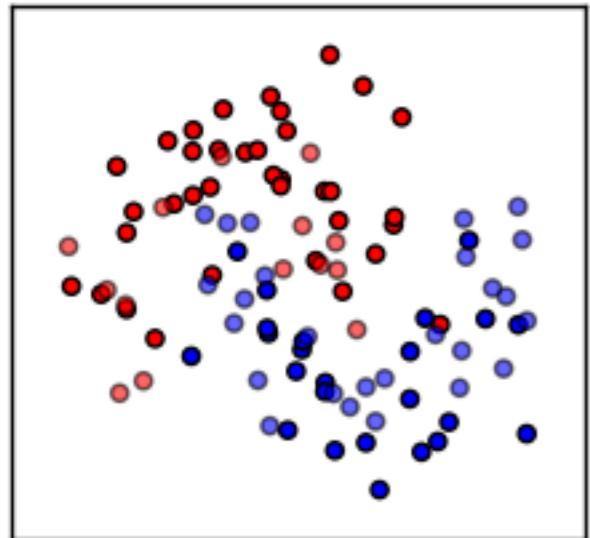


Gradient Boosted Trees

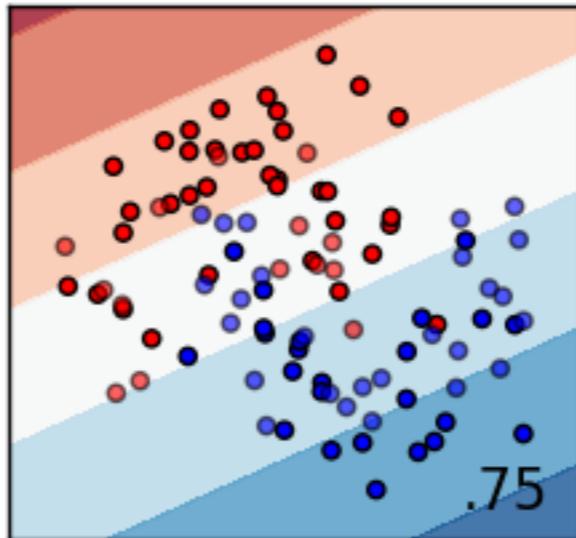




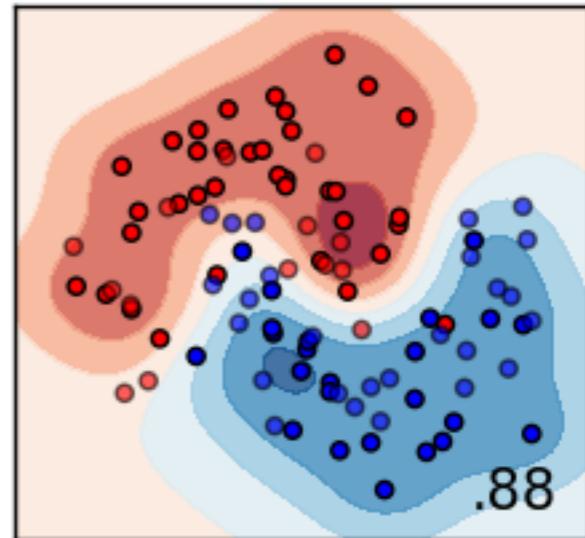
Nearest Neighbors



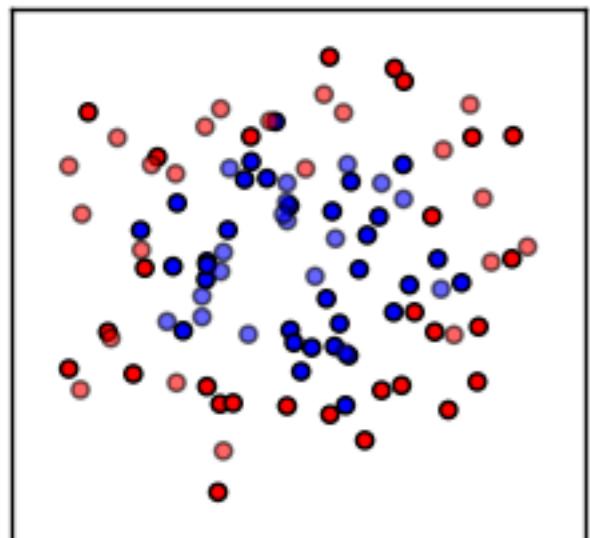
Linear SVM



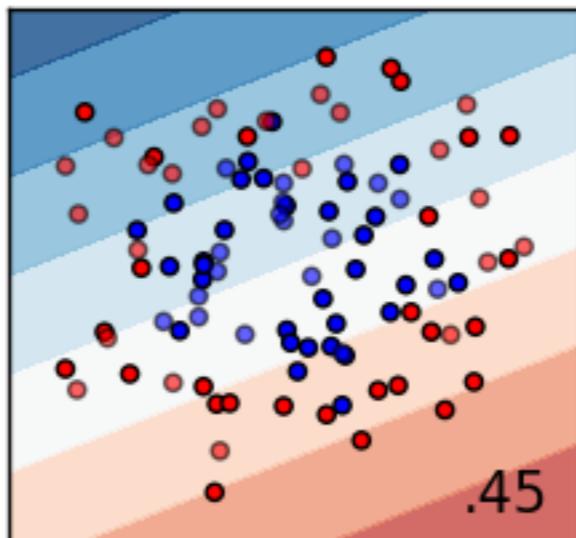
RBF SVM



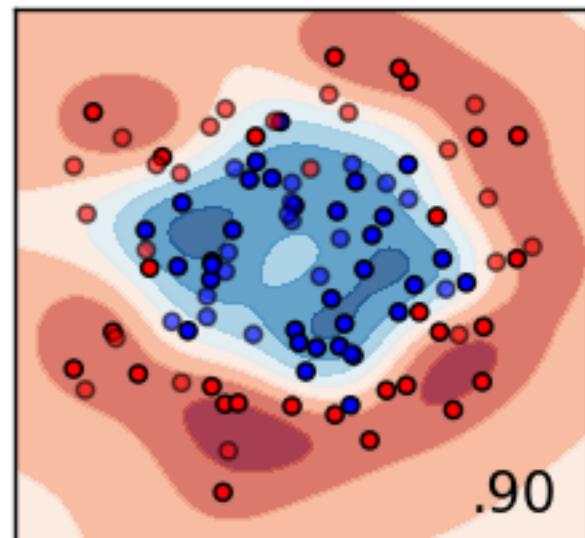
Nearest Neighbors



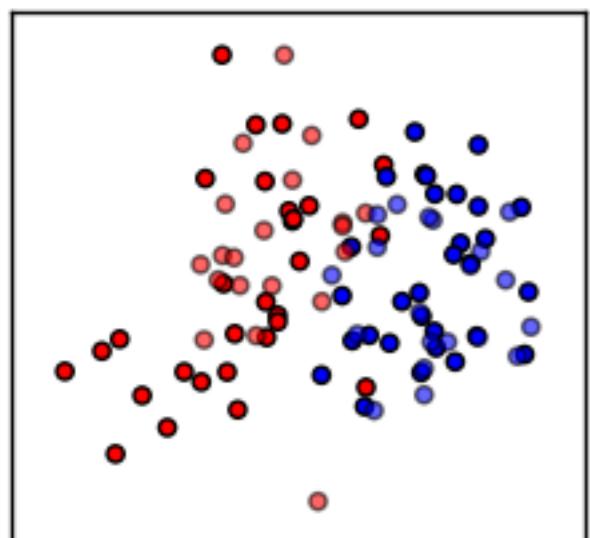
Linear SVM



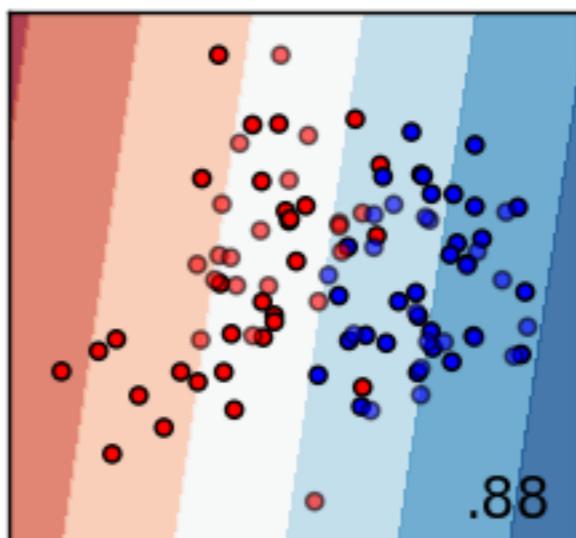
RBF SVM



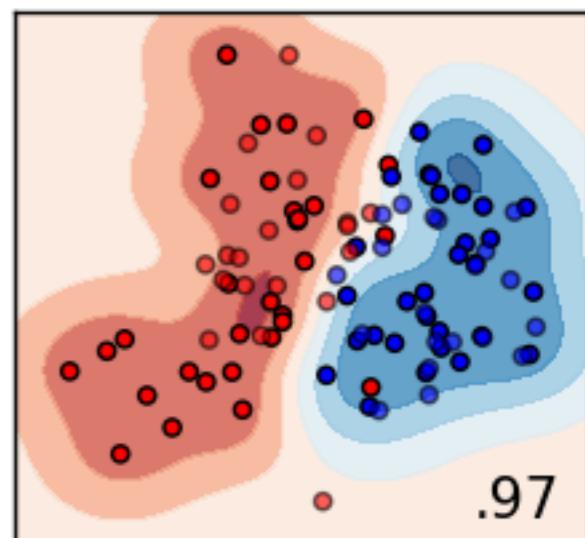
Nearest Neighbors



Linear SVM



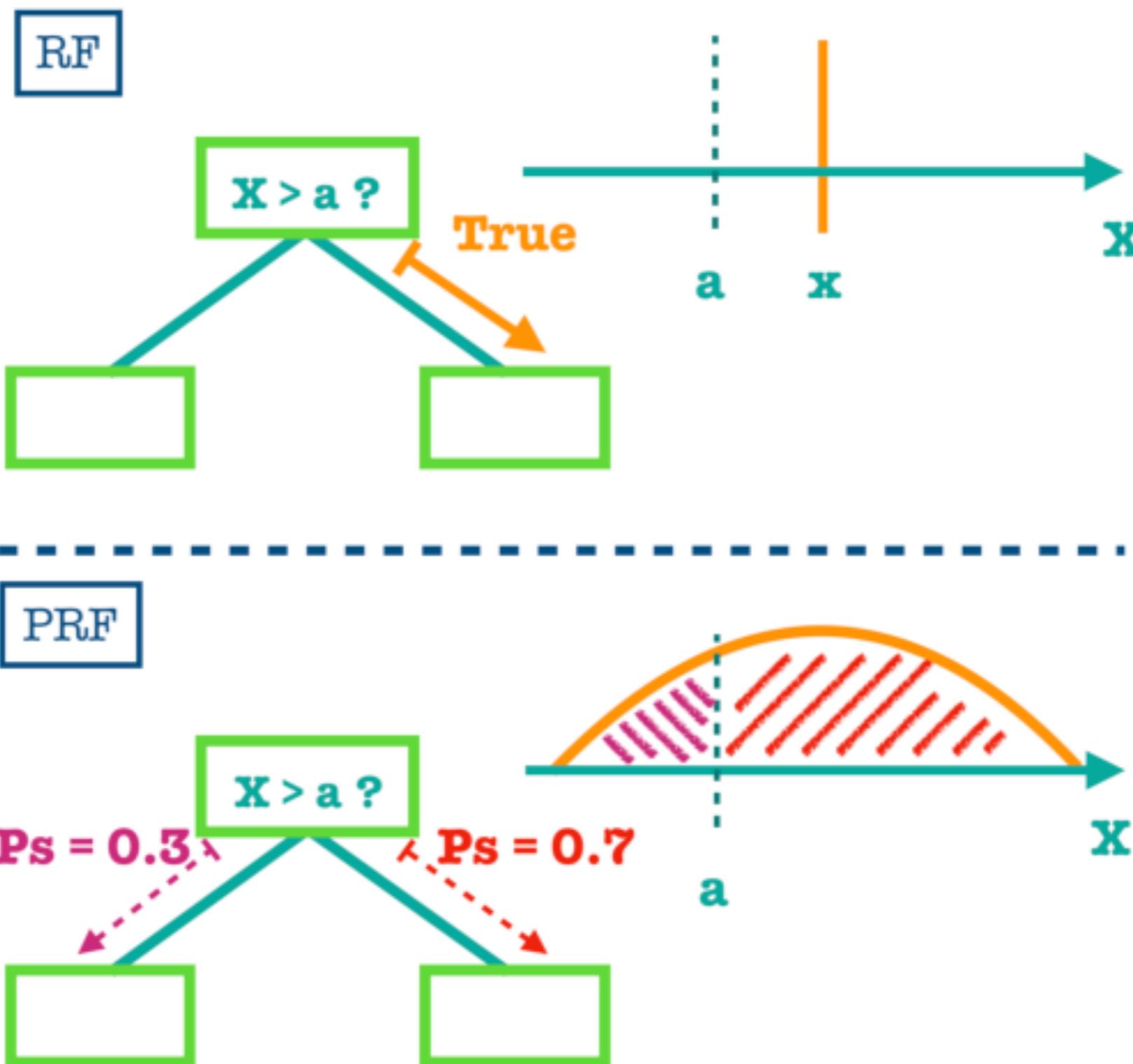
RBF SVM



Top 5 (supervised)

Algorithm	Comments
Neural Networks	<ul style="list-style-type: none">• Take long to train - lot of CPU• Overfits• Requires lot of data
Gradient Boosted Trees	<ul style="list-style-type: none">• Fast• Overfit danger
Random Forest	<ul style="list-style-type: none">• Robust to overfitting
SVM w/non-linear kernel	<ul style="list-style-type: none">• Pretty good
Gaussian Processes	<ul style="list-style-type: none">• non-parametric fitting

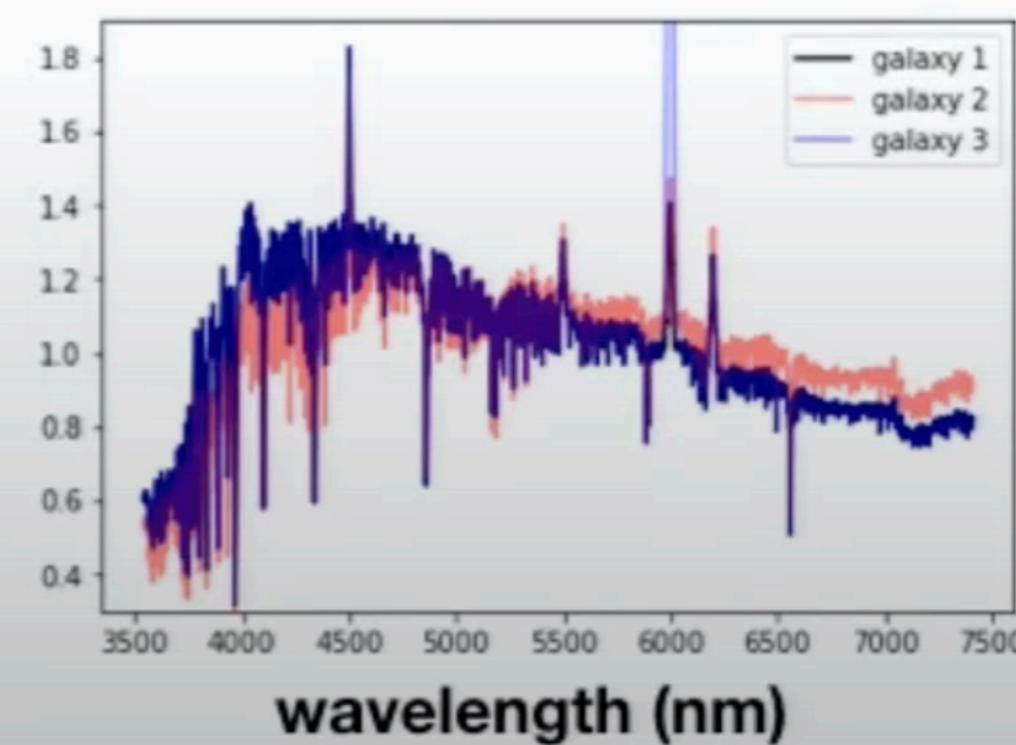
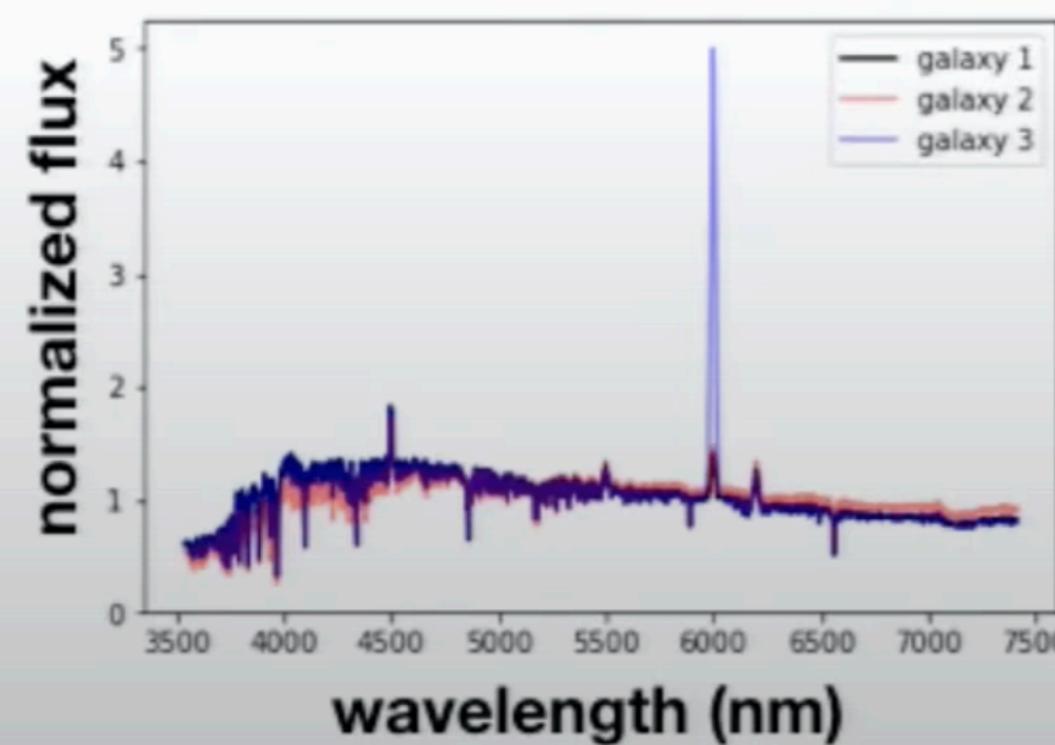
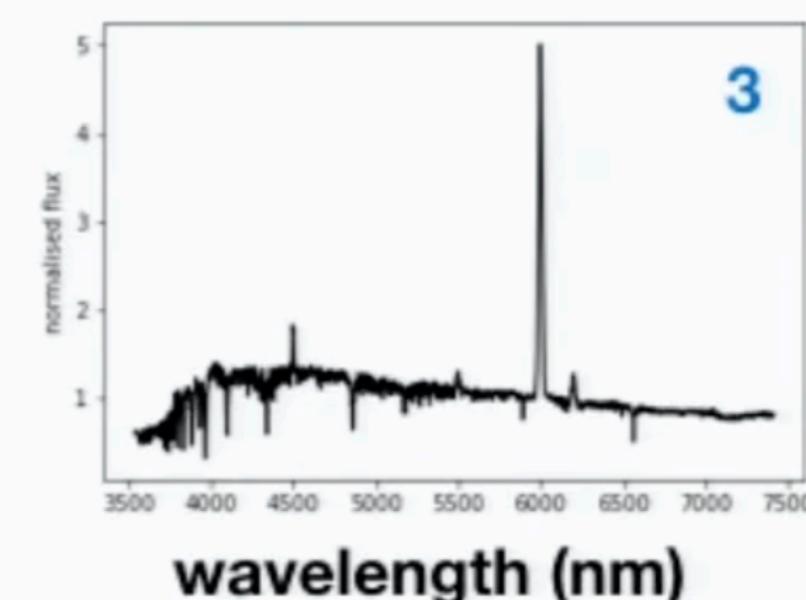
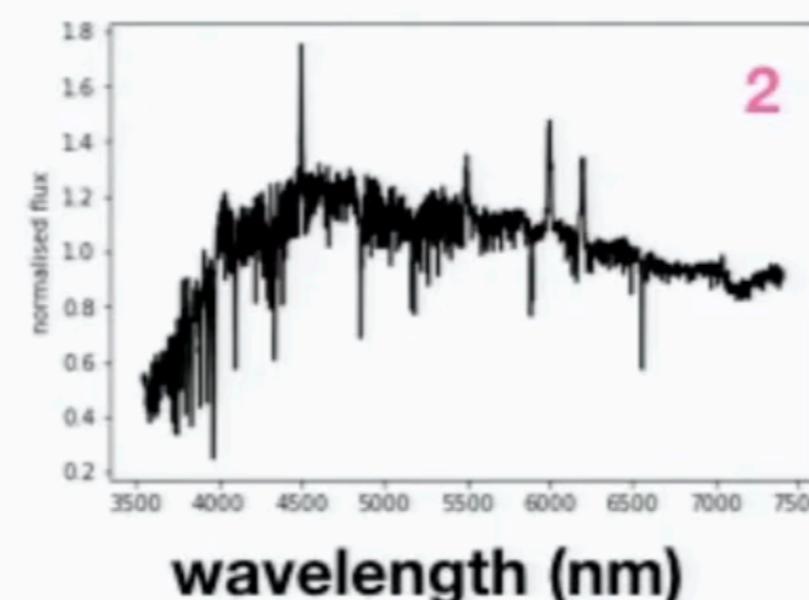
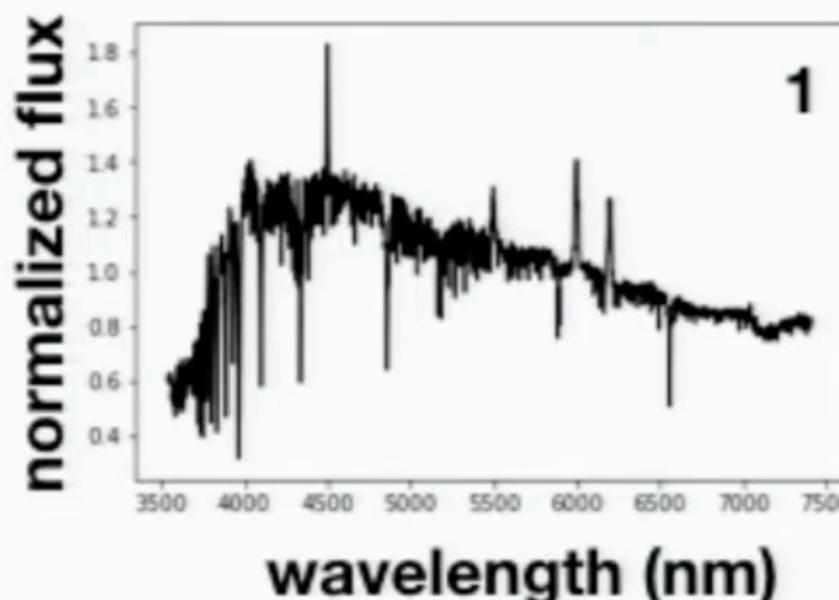
Probabilistic Random Forest



Unsupervised Random Forest

Random Forest can be used as an unsupervised algorithm, to produce pair-wise similarity for the objects in our sample.

Why do we need to measure distances between objects?

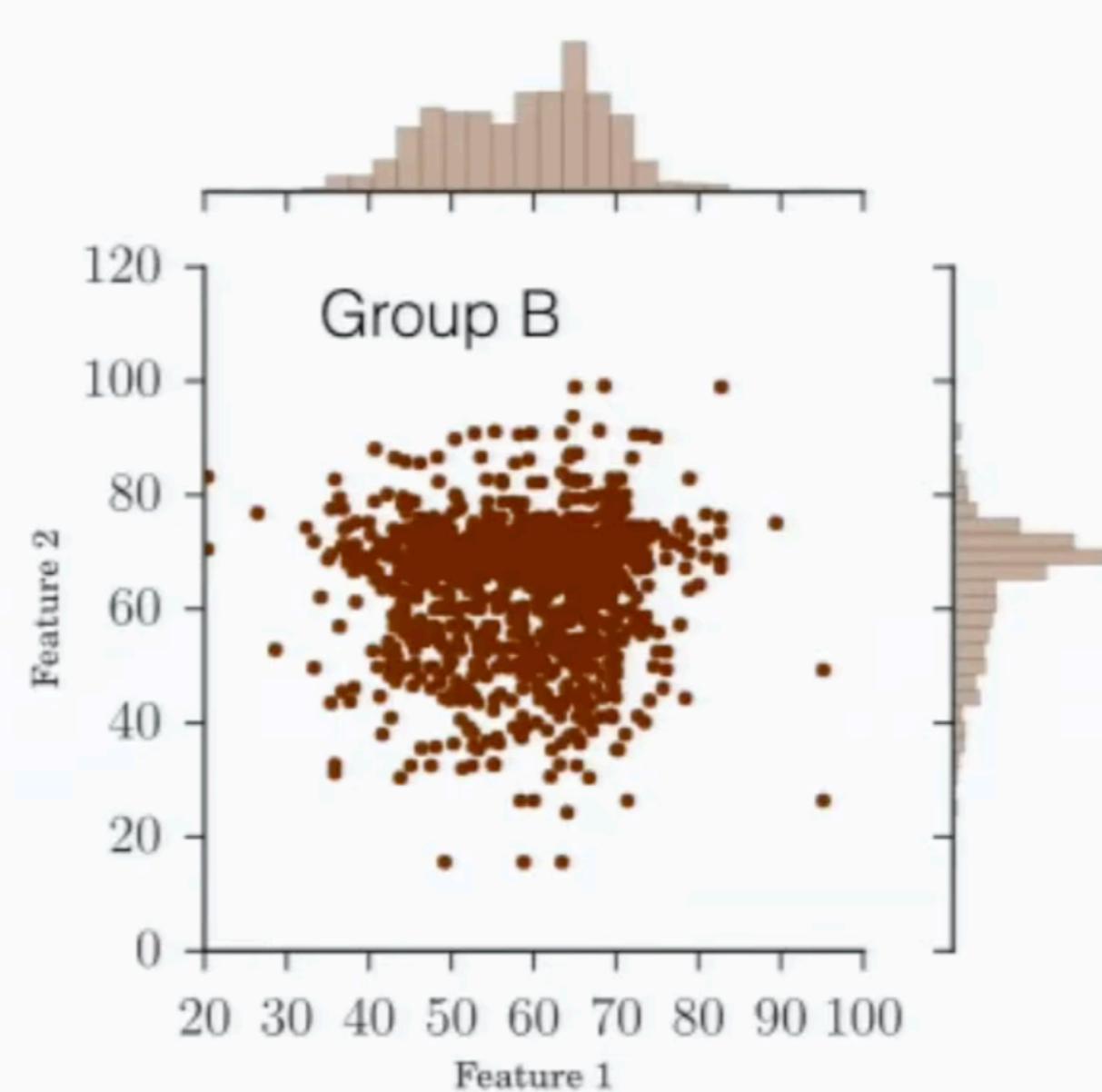
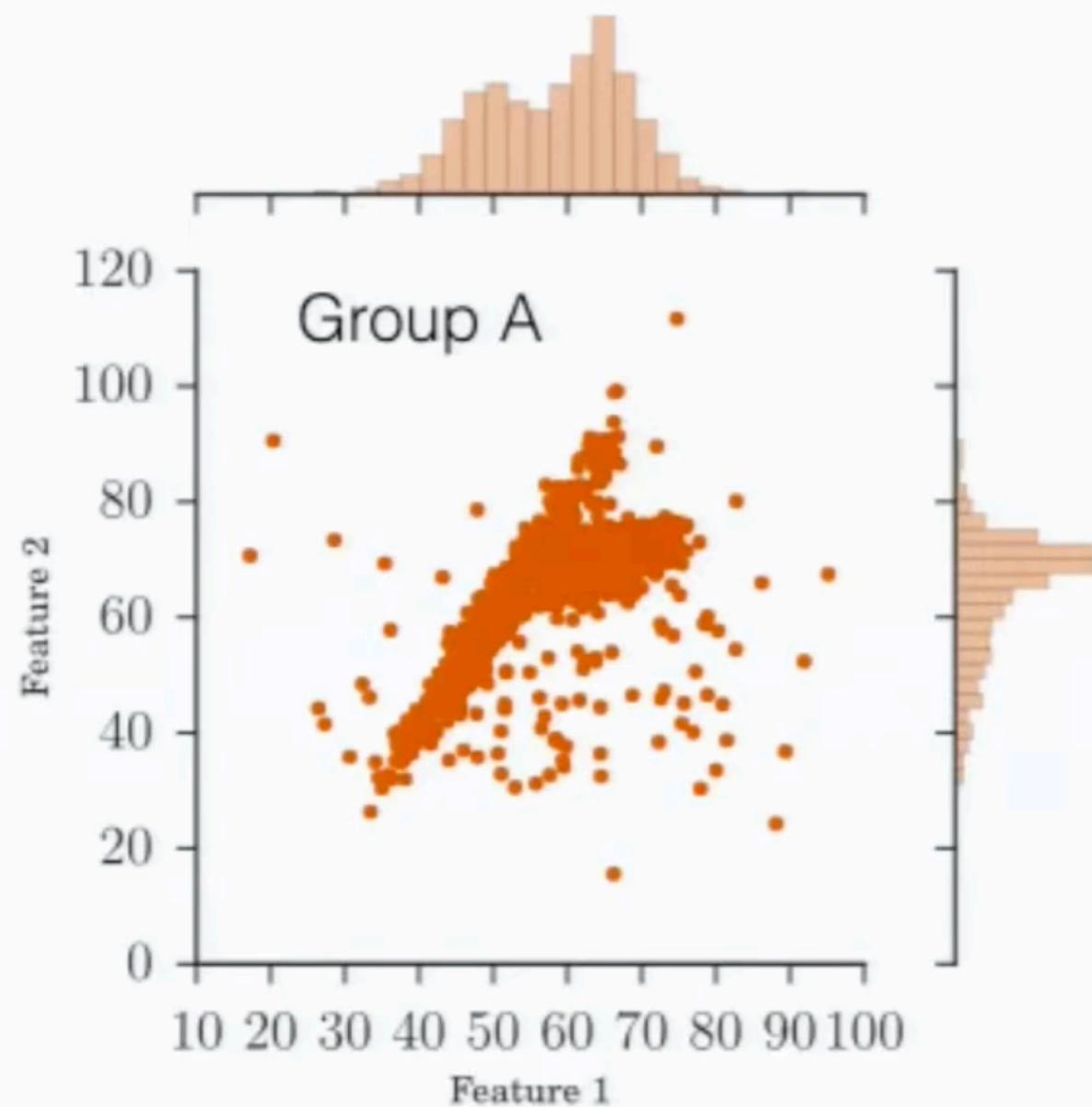


Unsupervised Random Forest

Random Forest can be used as an unsupervised algorithm, to produce pair-wise similarity for the objects in our sample.

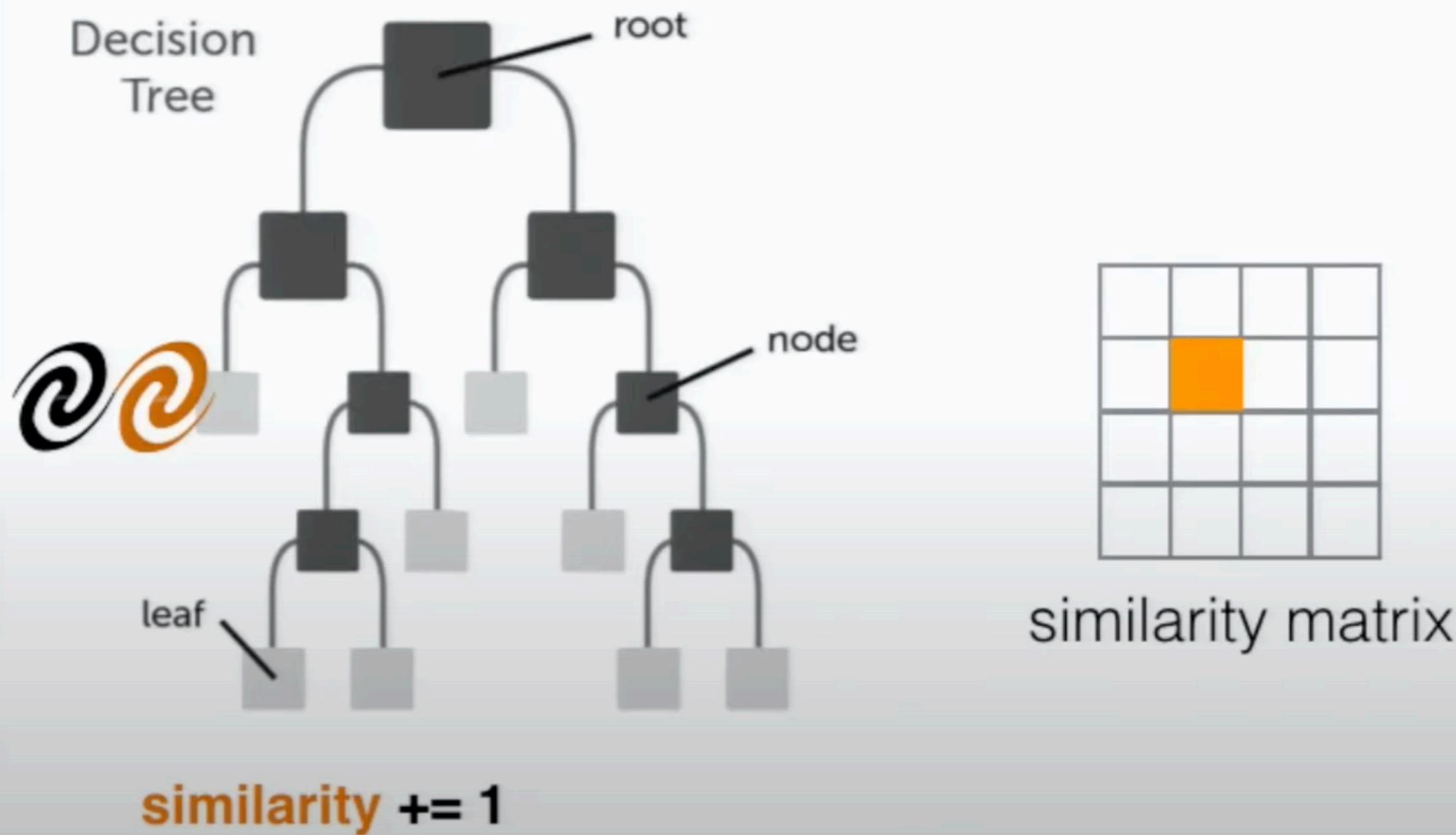
Input dataset: a list of objects with measured features, but no labels!

Random Forest is trained to distinguish between real and synthetic datasets.



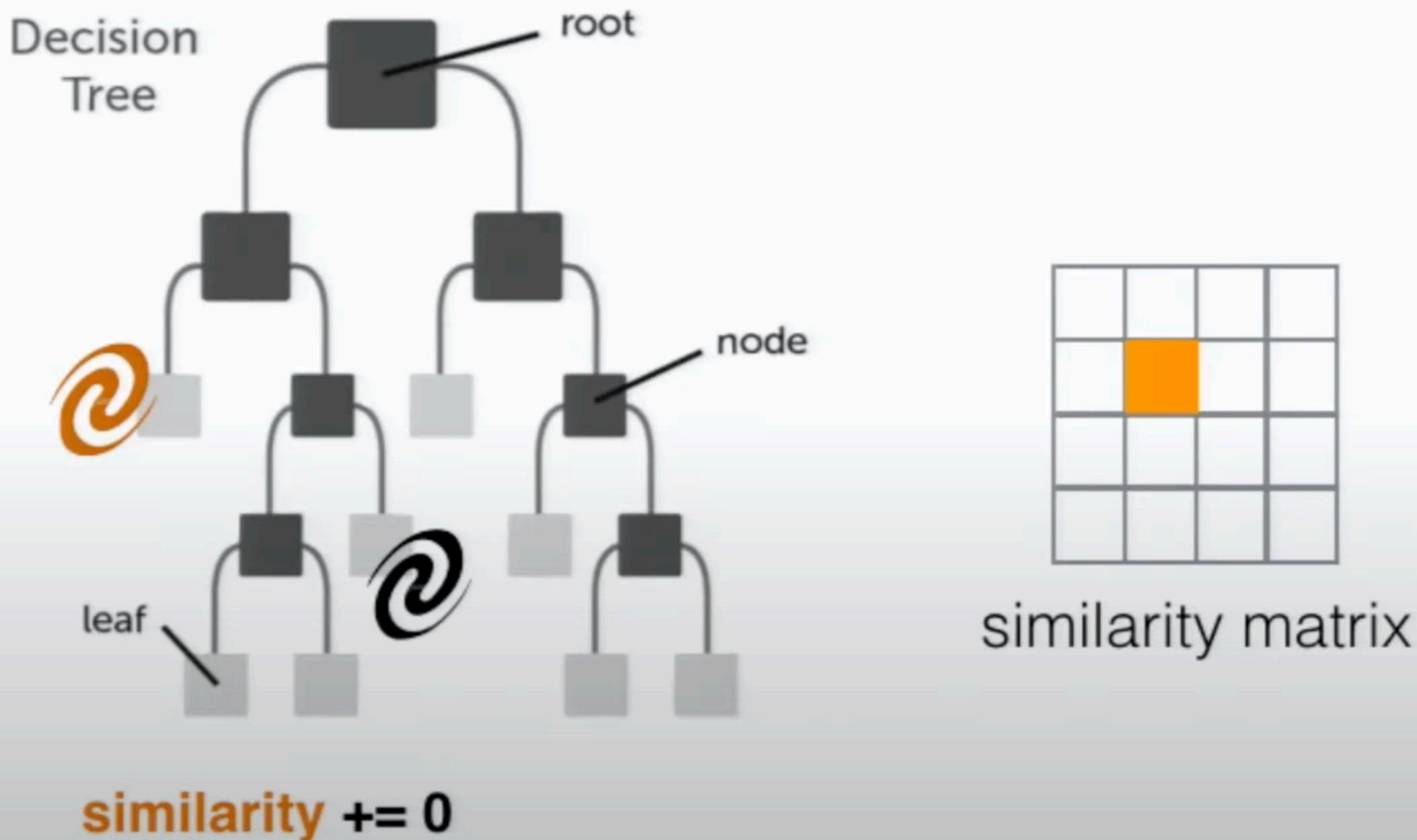
Unsupervised Random Forest

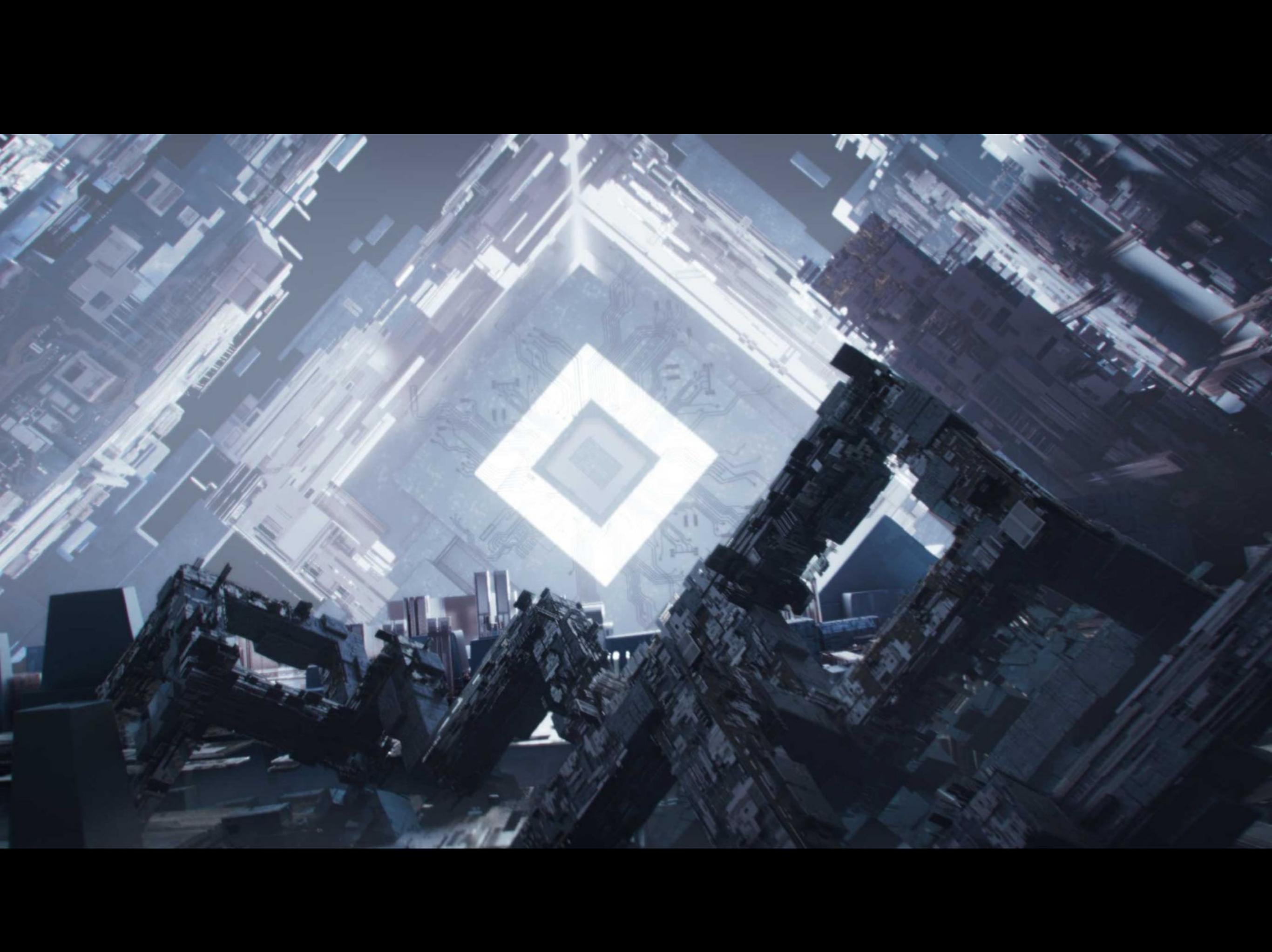
We train the Random Forest to distinguish between groups A and B.
For group A (real data), we propagate the objects and obtain a similarity matrix.



Unsupervised Random Forest

We train the Random Forest to distinguish between groups A and B.
For group A (real data), we propagate the objects and obtain a similarity matrix.

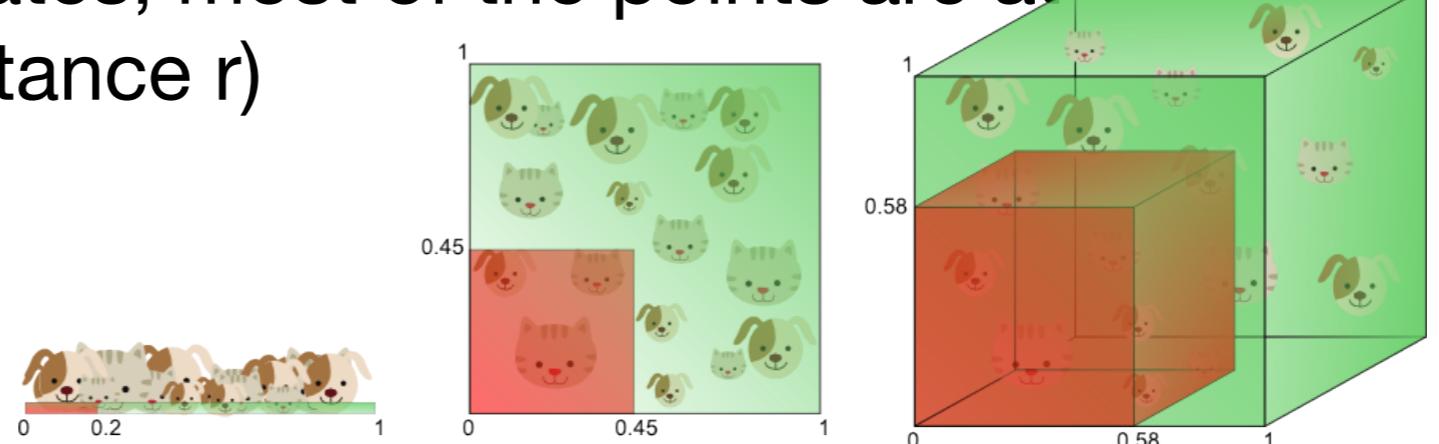




Curse of Dimensionality

- Too many features
 - Expensive to store
 - Slowing down computation
 - Subject to *Dimensionality curse*
 - Sample space gets harder and harder to fill as dimensions grow
 - A reason why too many features lead to overfitting as data become **sparse**
-
- “*If people can see in multi-dimensions we would not need machine learning*”
 - More and more data needed to fill the same % of space
 - Distance measure degenerates, most of the points are at the surface of a sphere (distance r)

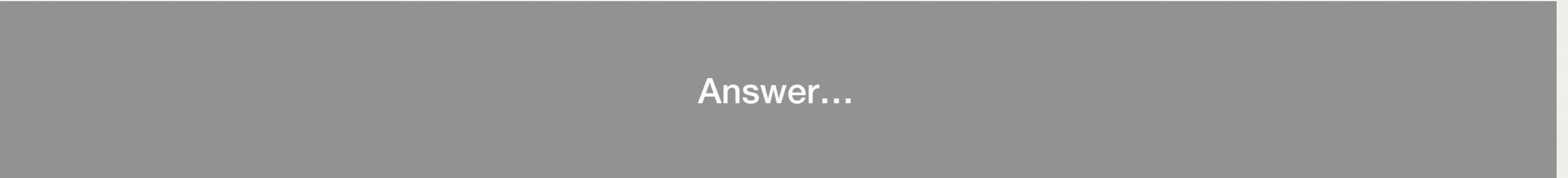
$$\frac{V_{\text{hypersphere}}}{V_{\text{hypercube}}} = \frac{\pi^{d/2}}{d2^{d-1}\Gamma(d/2)} \rightarrow 0 \text{ as } d \rightarrow \infty.$$



How Many Shades of Gray Can you Distinguish?



Answer...

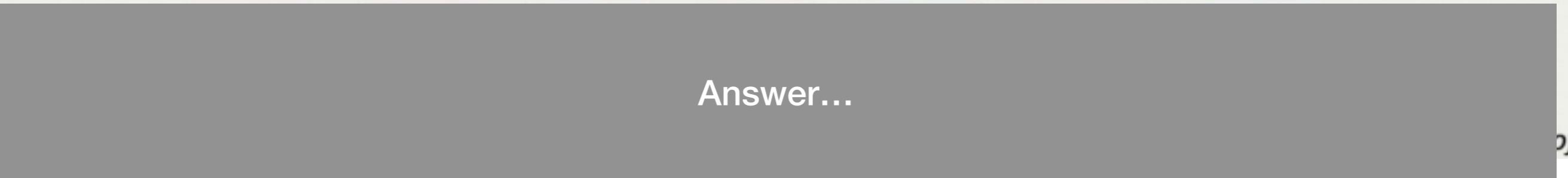


Value encodes continuous variables (less well)

How Many Colors?



Hue encodes nominal variables



Answer...

(off)

How Many Shades of Gray Can you Distinguish?



Value easily encodes ordinal variables

How Many Colors?

Hue encodes nominal variables

Answer...

(off)

How Many Shades of Gray Can you Distinguish?

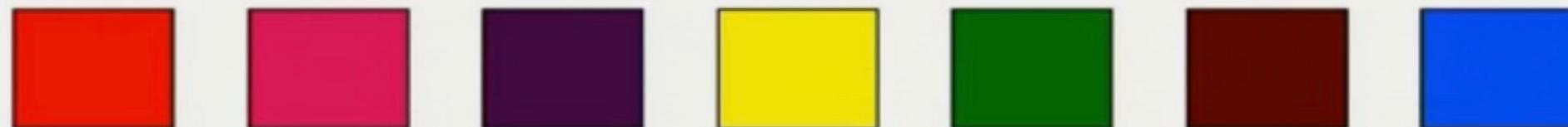
Value easily encodes ordinal variables



Value encodes continuous variables (less well)

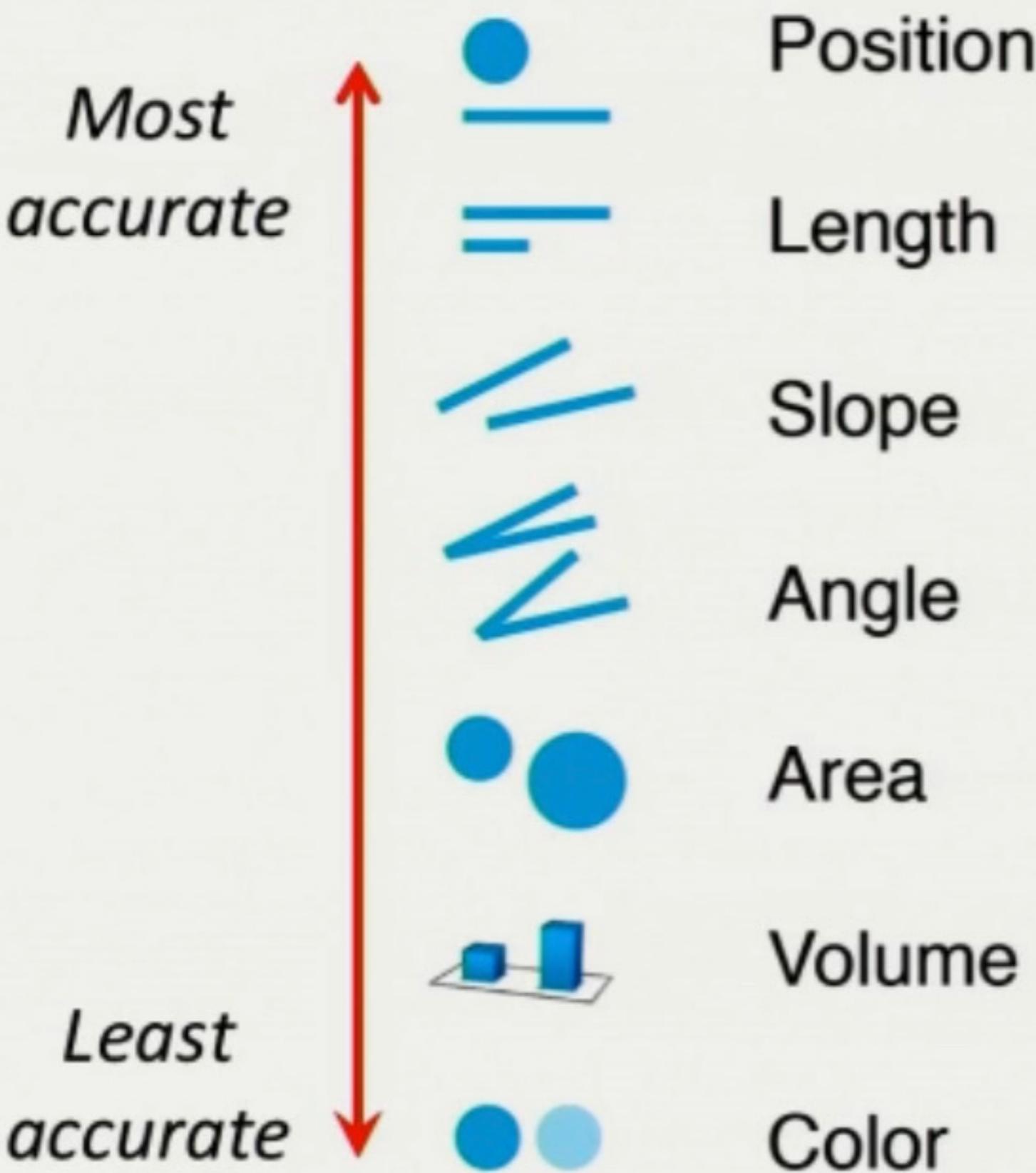
How Many Colors?

Hue encodes nominal variables

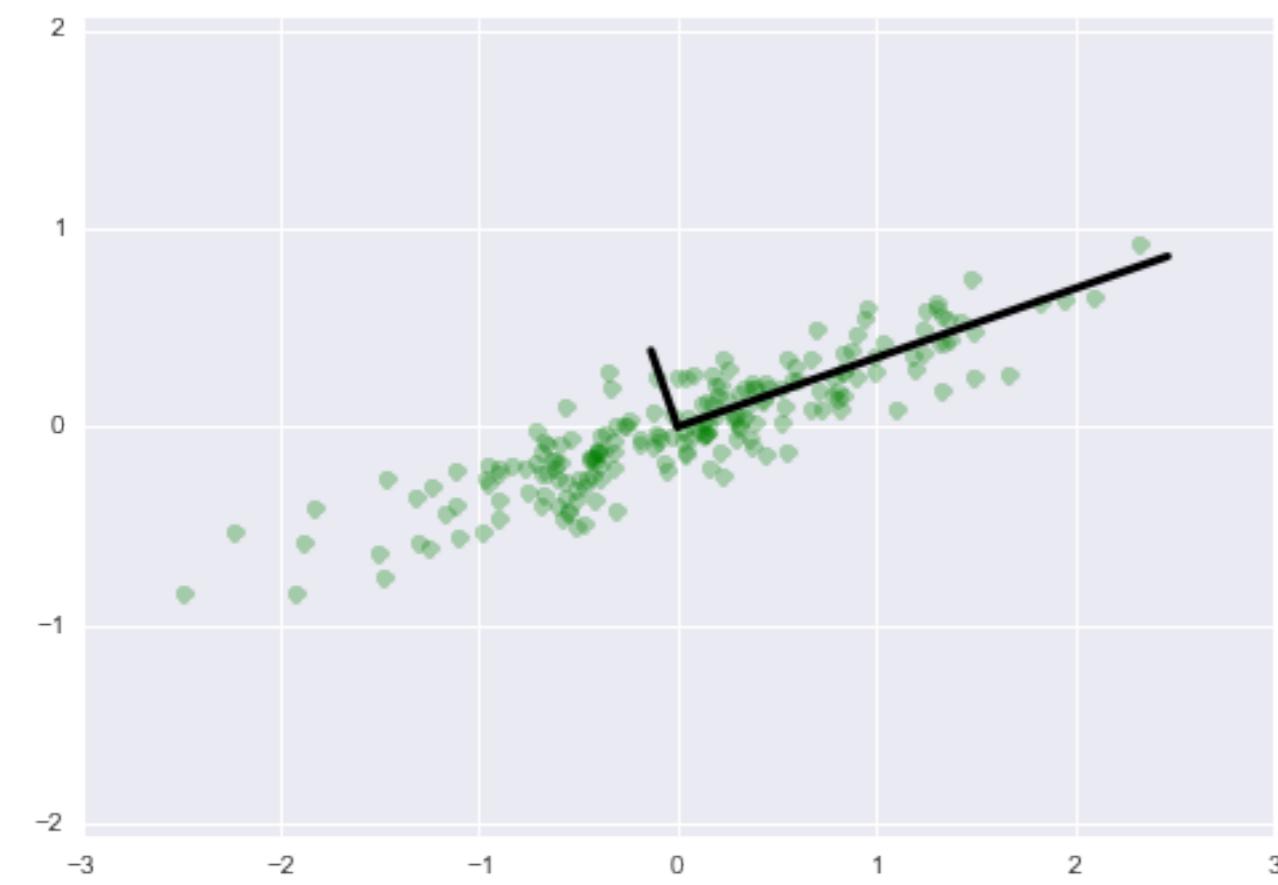


(from S. Davidoff)

Relative accuracy of the visualisation space axes:

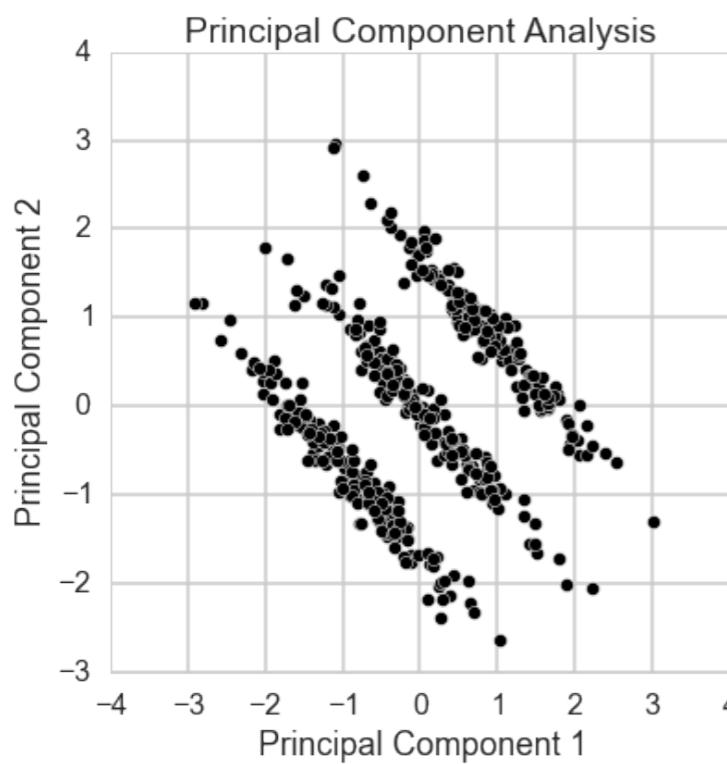


PCA (unsupervised)

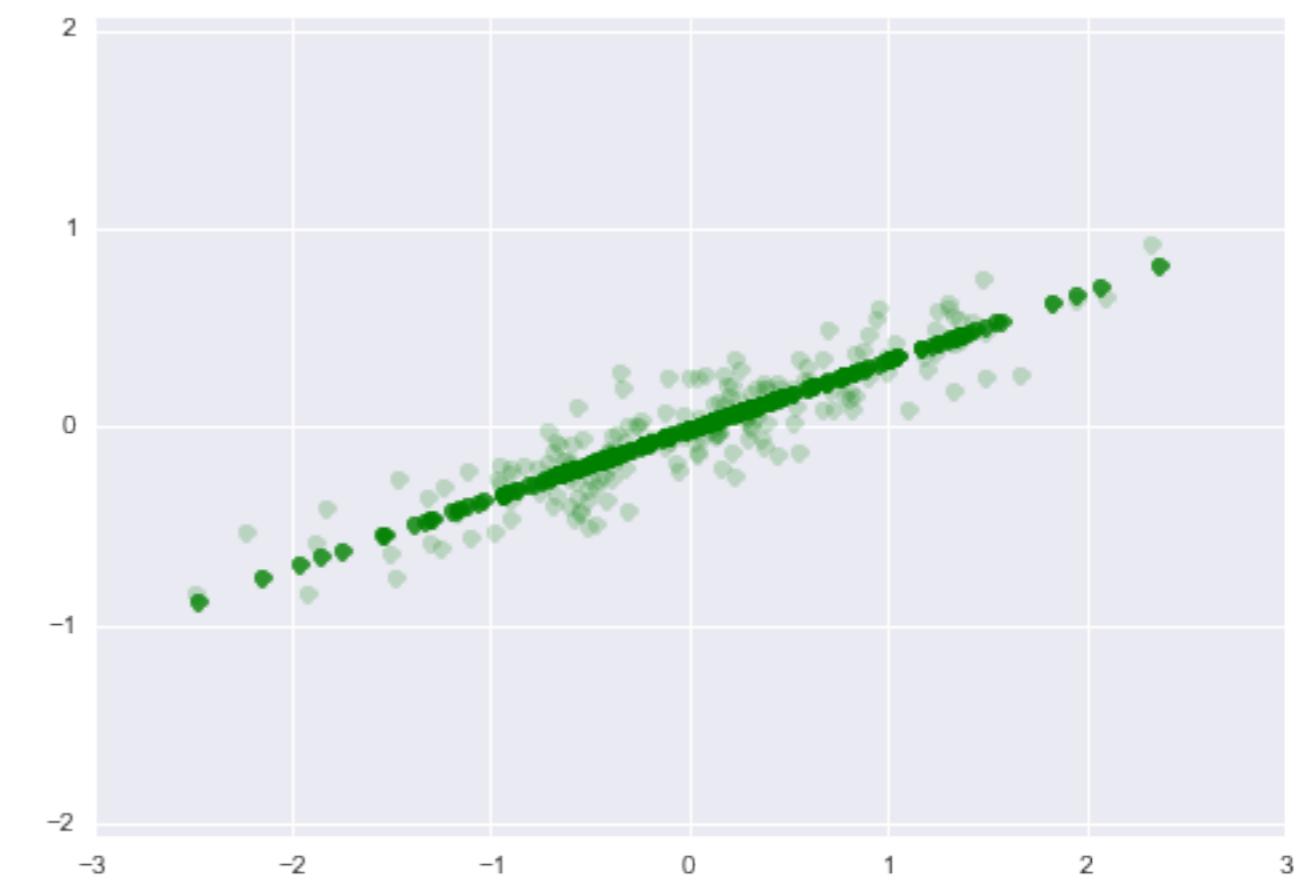
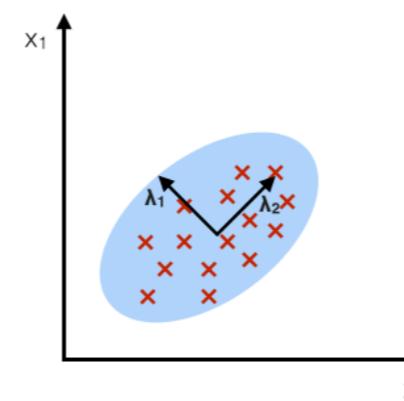


$$var(x) = \frac{\sum(x_i - \bar{x})^2}{N}$$

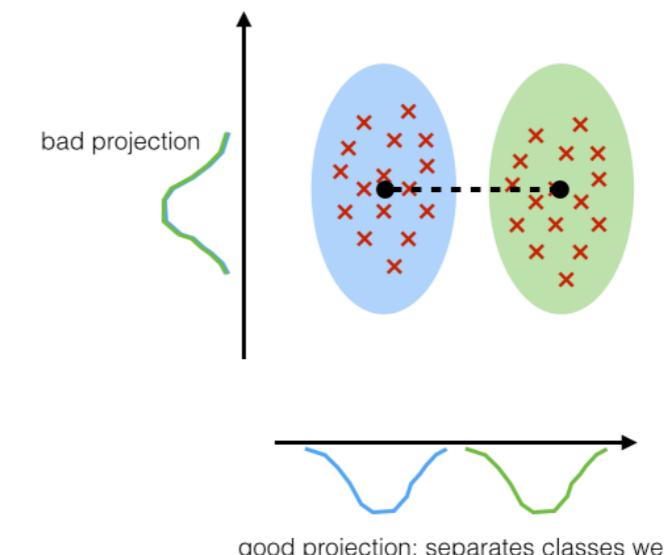
$$cov(x, y) = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{N}$$



PCA:
component axes that
maximize the variance

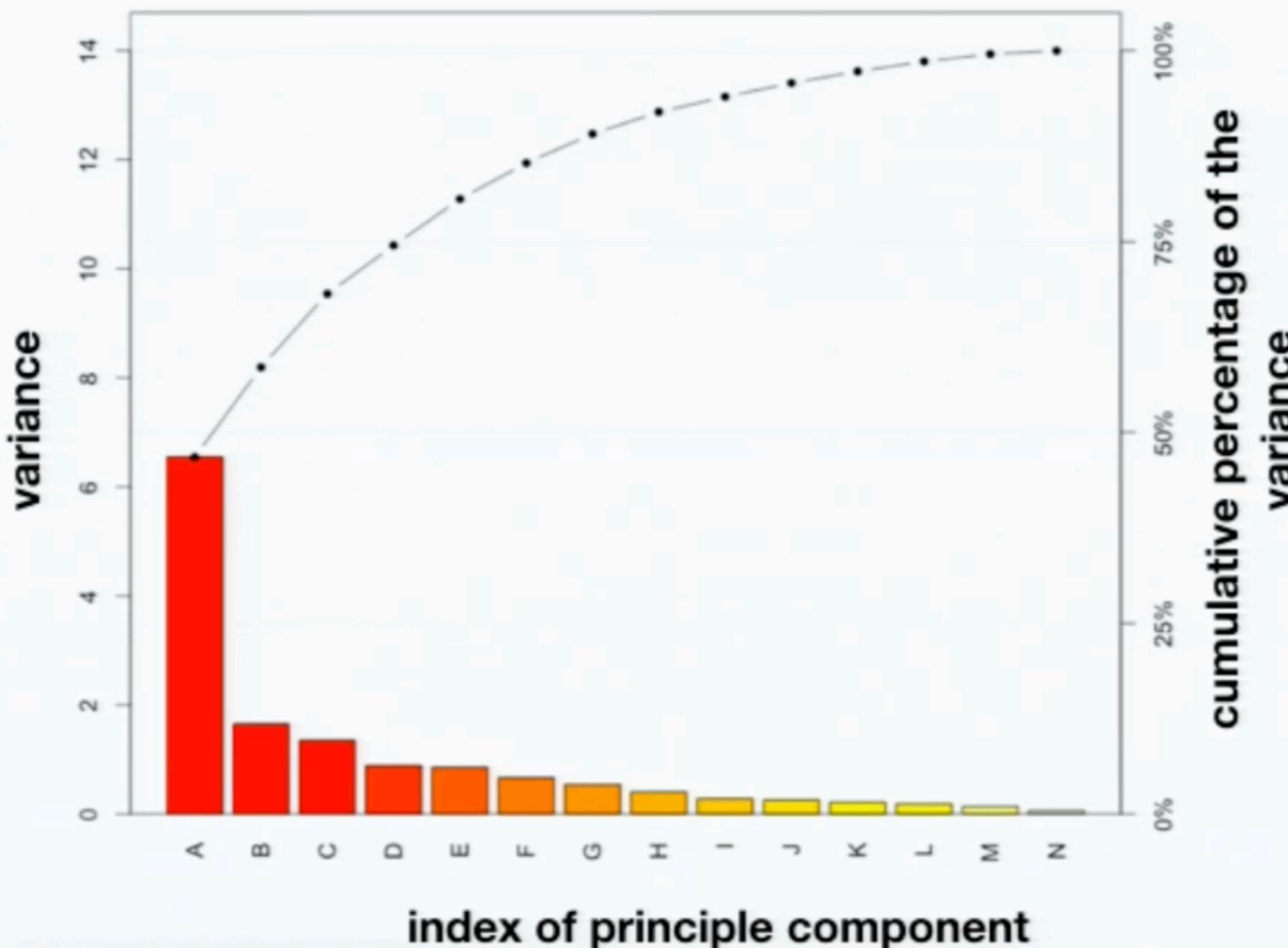


LDA:
maximizing the component
axes for class-separation



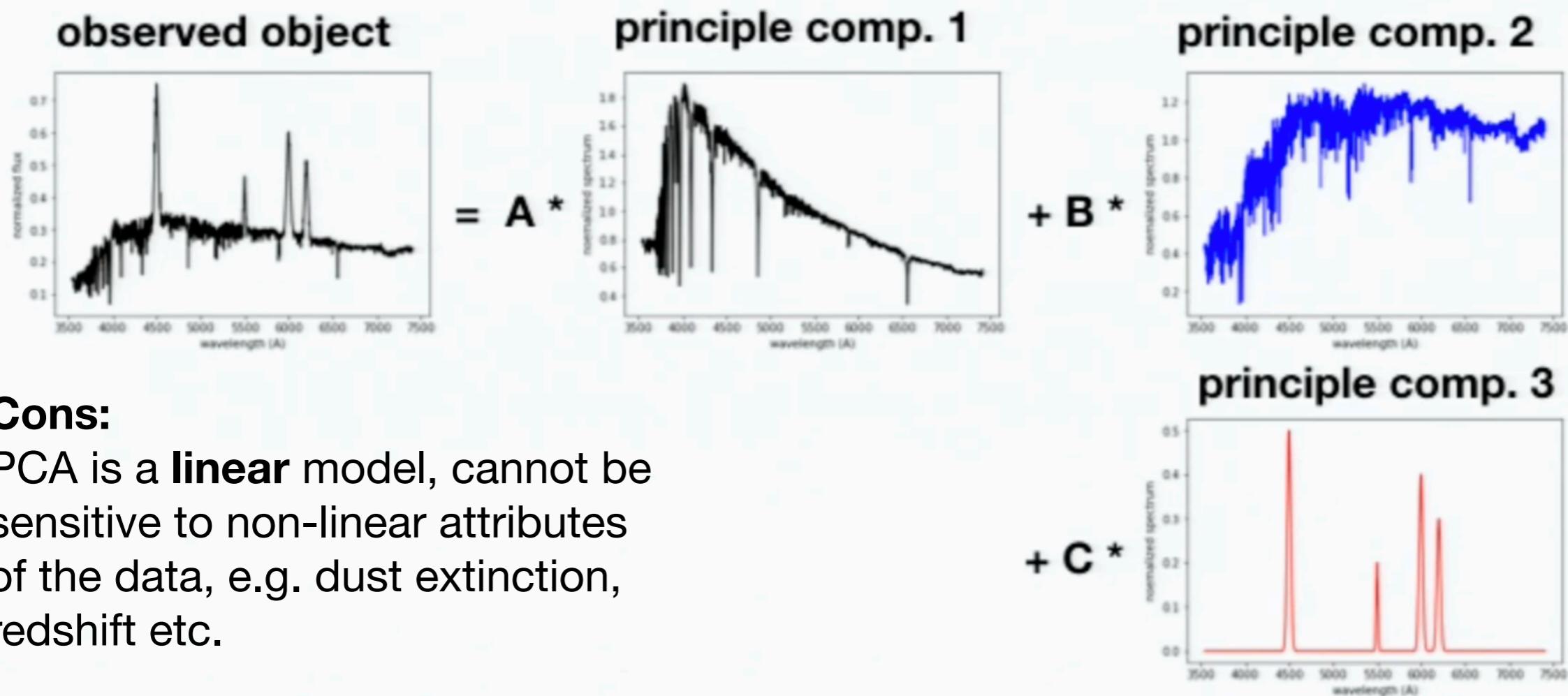
Principle Component Analysis (PCA)

PCA allows us to compress the data, by representing each object as a projection on the first principle components.



Principle Component Analysis (PCA)

The principle components **may** represent the true **building blocks** of the objects in our dataset.



Cons:

PCA is a **linear** model, cannot be sensitive to non-linear attributes of the data, e.g. dust extinction, redshift etc.

Why do we need dimensionality reduction?

- Improve performance of supervised learning algorithms:
 - features can be correlated and redundant
 - most algorithms cannot handle 1000's of features
- Compressing data (SKA, LSST...)
- Data visualisation and interpretation
- Uncover complex trends
- Look for outliers and "unknown unknowns"



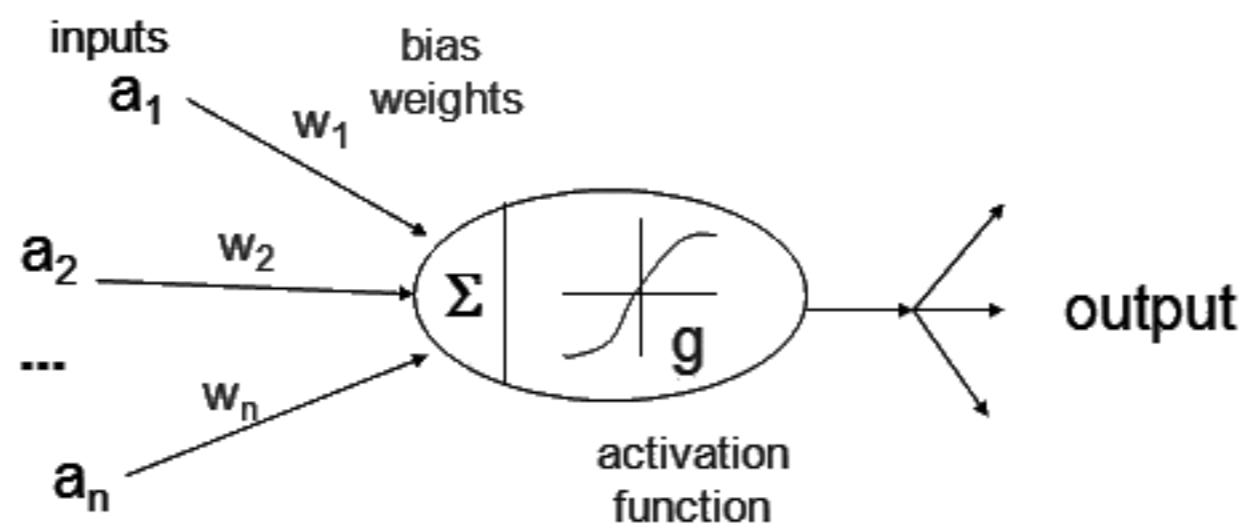
“Deep learning (neural networks) is like teenage sex: everyone talks about it, nobody really knows how to do it, everyone thinks everyone else is doing it, so everyone claims they are doing it....”

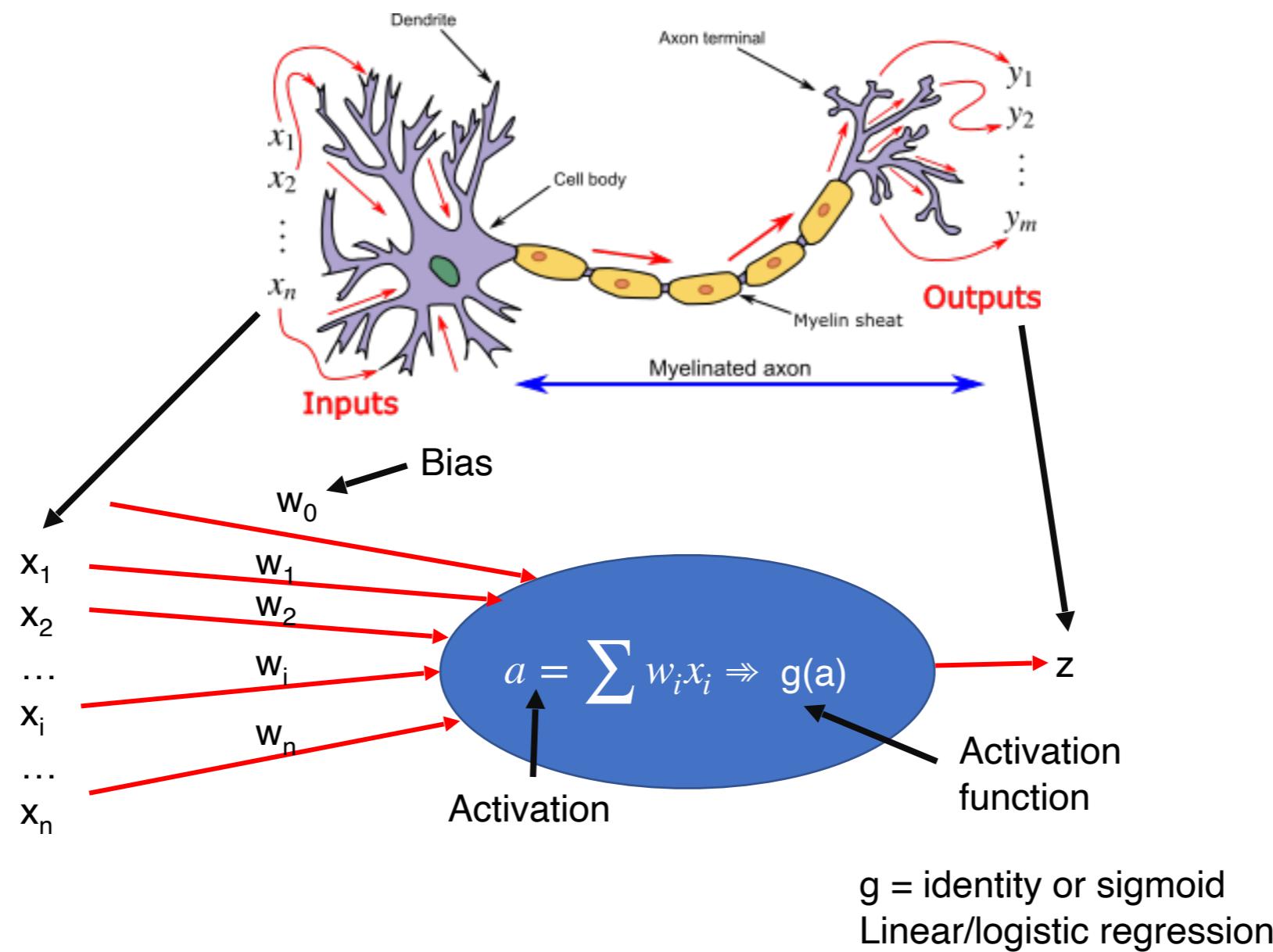
– a ML expert

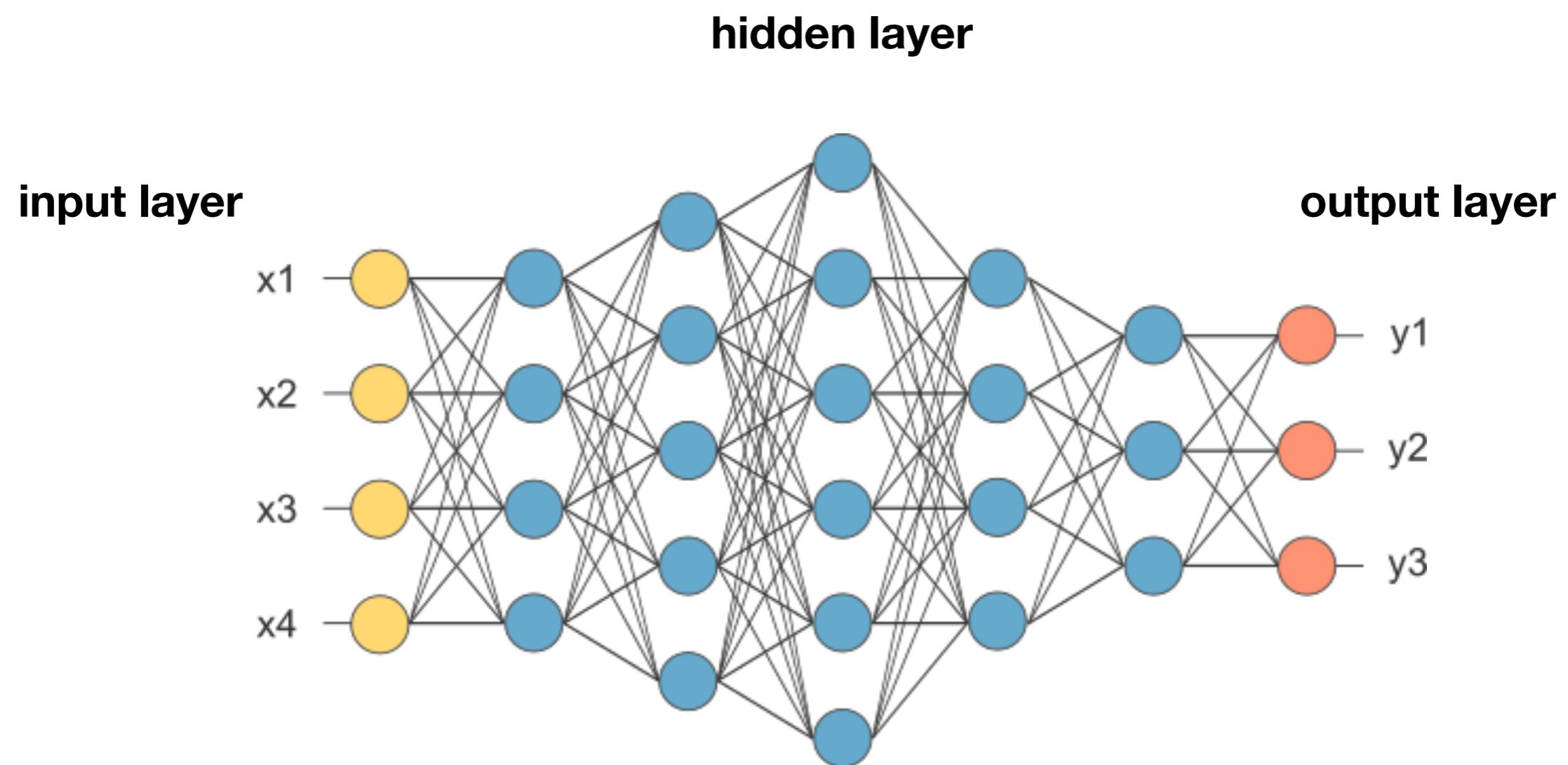
Buzz words: machine learning, deep learning, AI, neural network, Bayesian analysis, gravitational waves, exoplanets, dark matter...

According to a report by MMC Ventures, close to 40% of European companies claiming to use AI don't actually use it

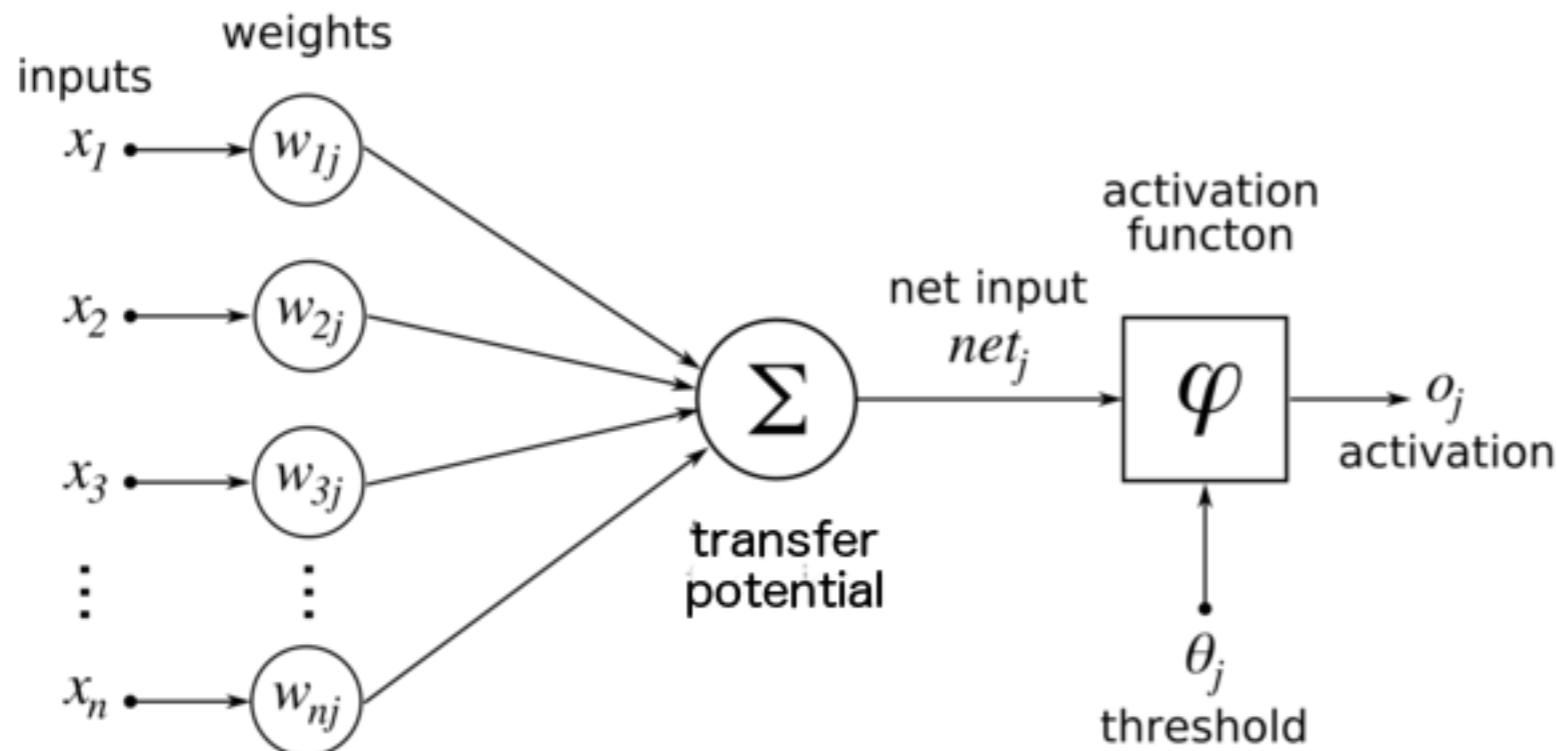






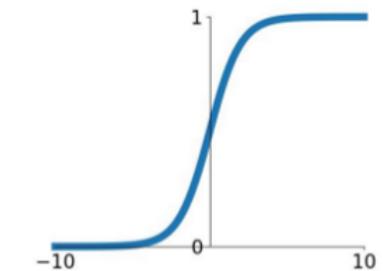


Activation Functions



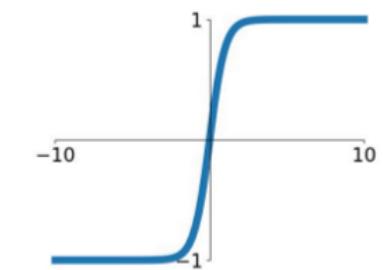
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



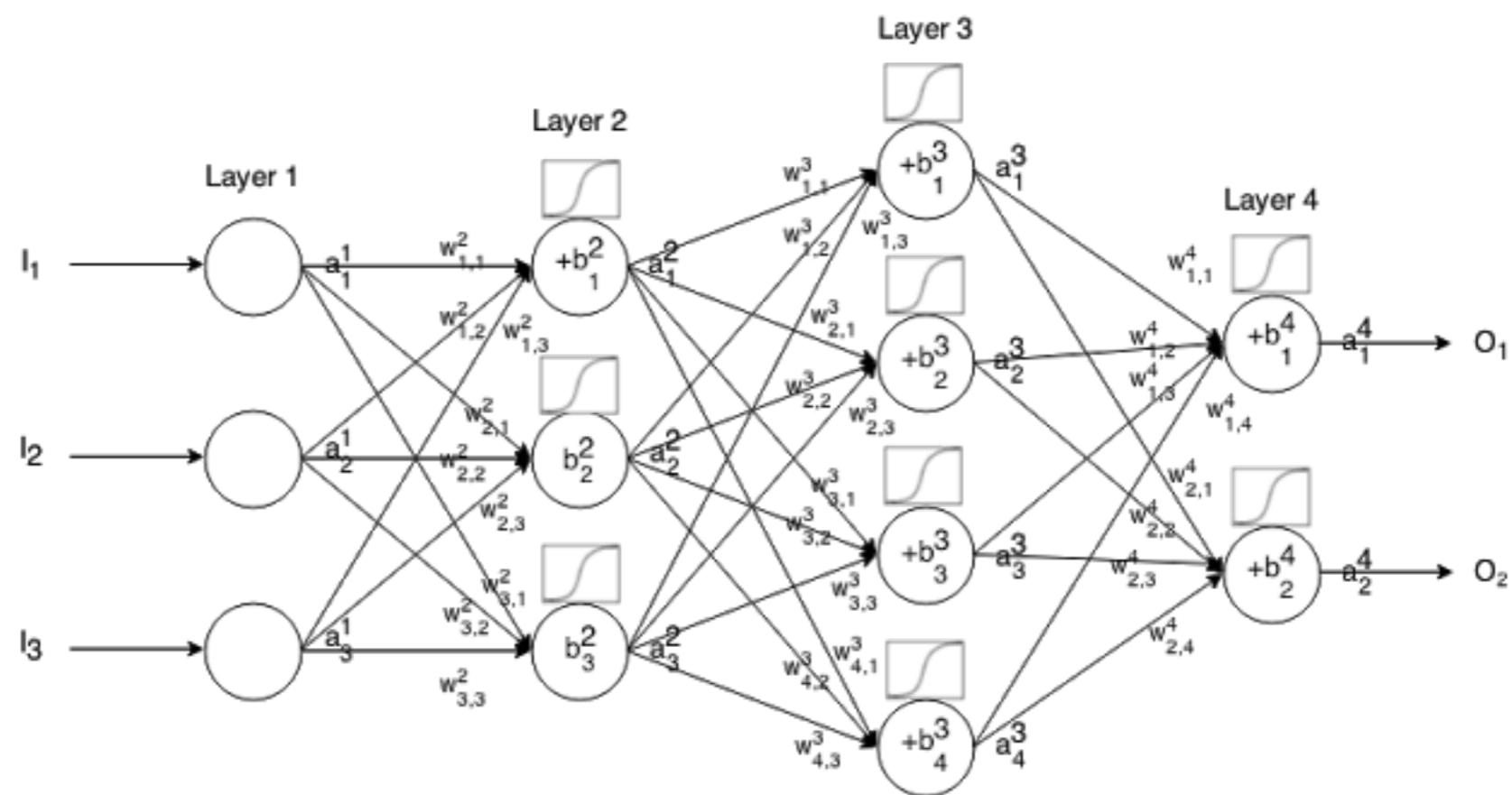
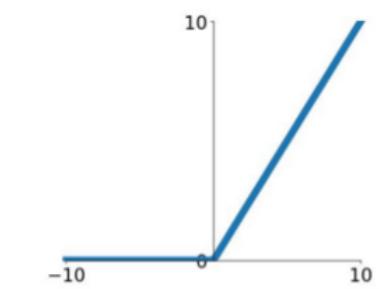
tanh

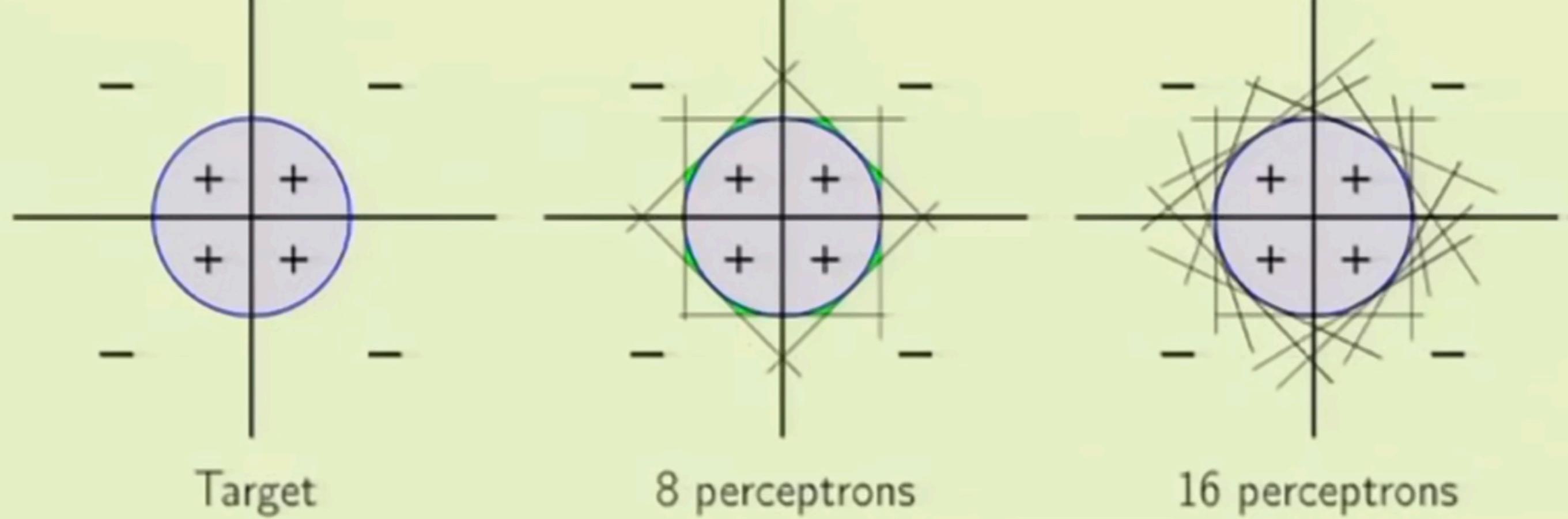
$$\tanh(x)$$

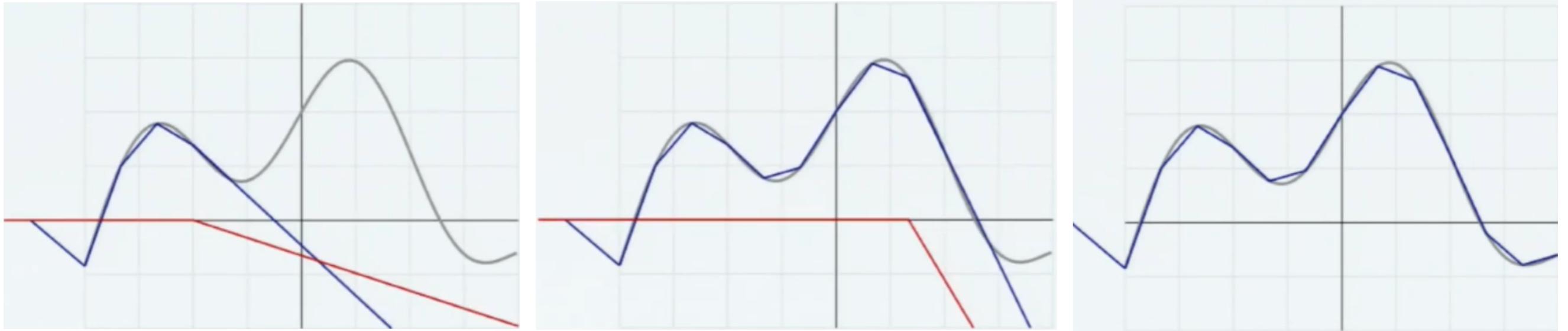


ReLU

$$\max(0, x)$$





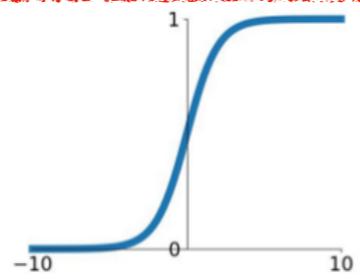


Activation Functions

[0,1]
"probability"
final layer for
classification

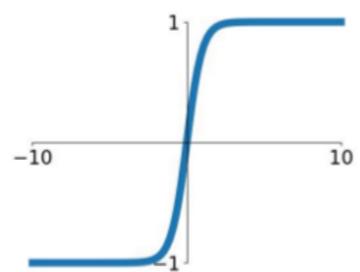
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



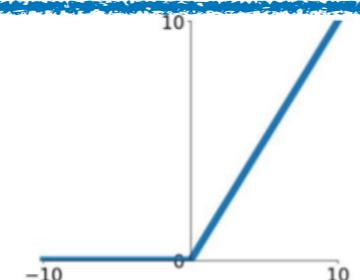
tanh

$$\tanh(x)$$



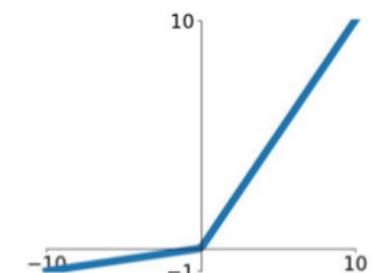
ReLU

$$\max(0, x)$$



robust
against
"vanishing
gradient"

Leaky ReLU
 $\max(0.1x, x)$

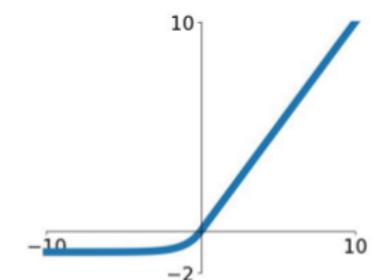


Maxout

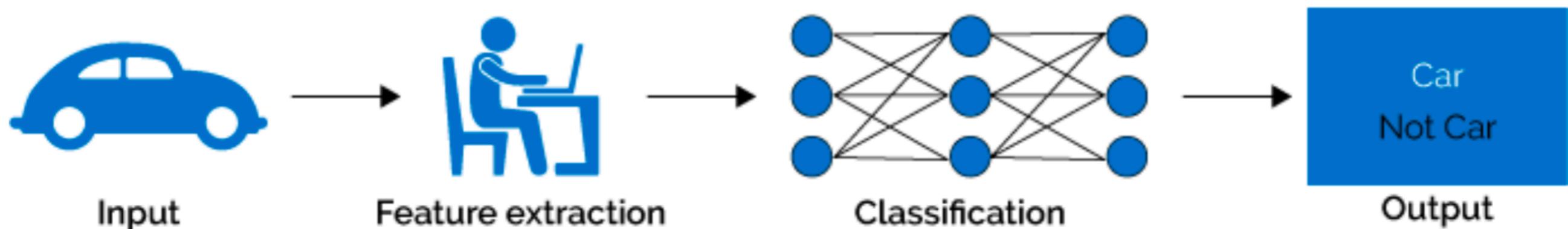
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

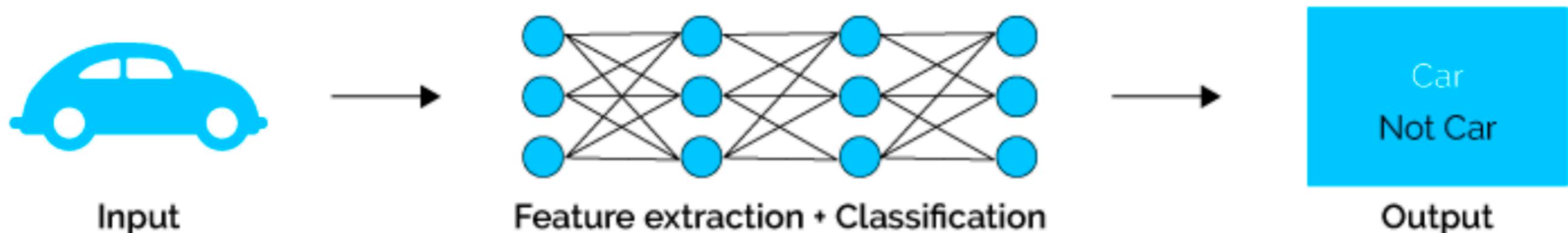
$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



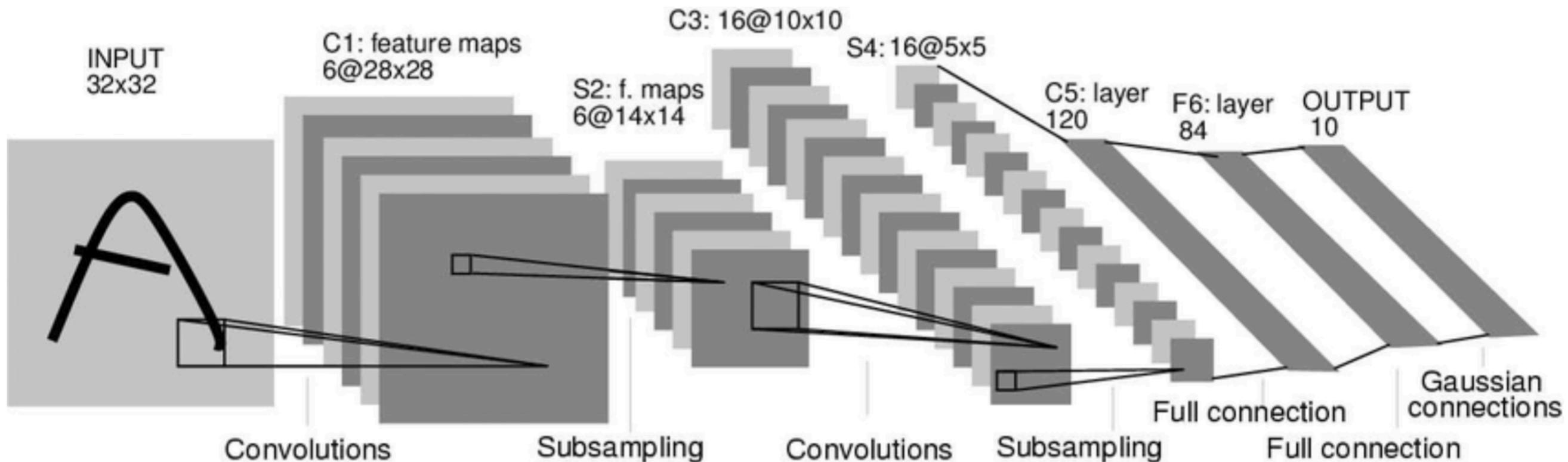
Machine Learning

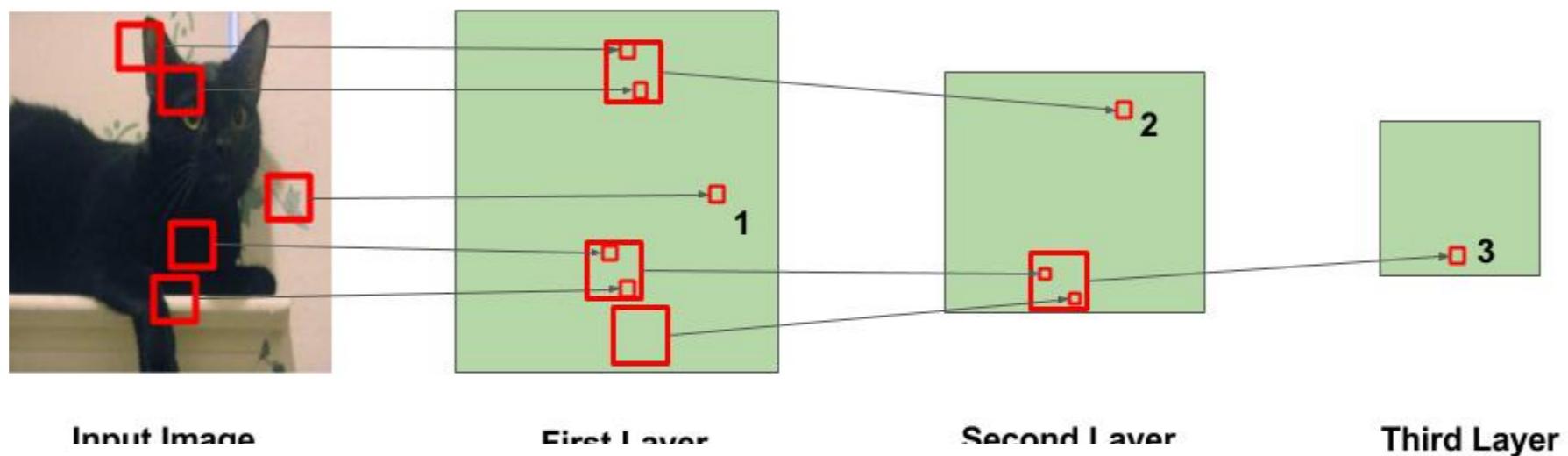


Deep Learning

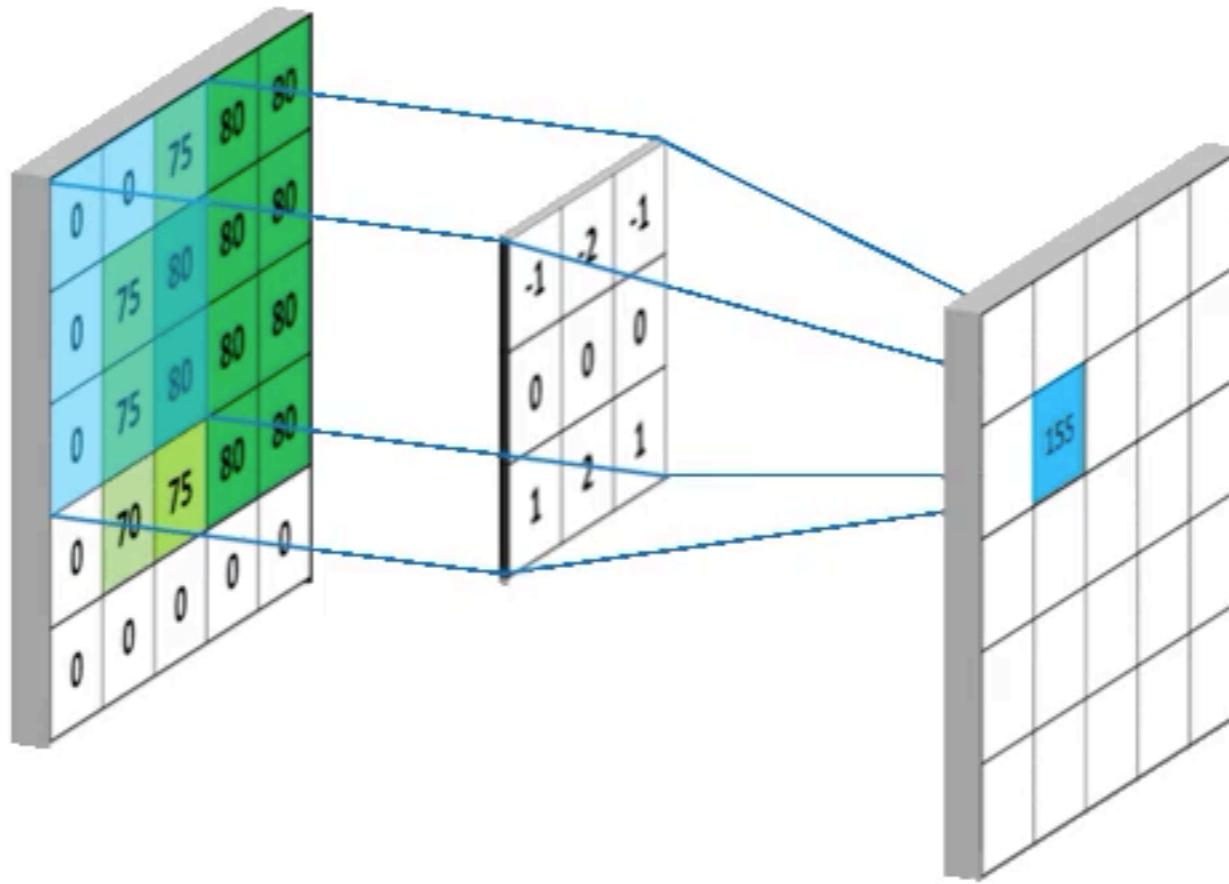


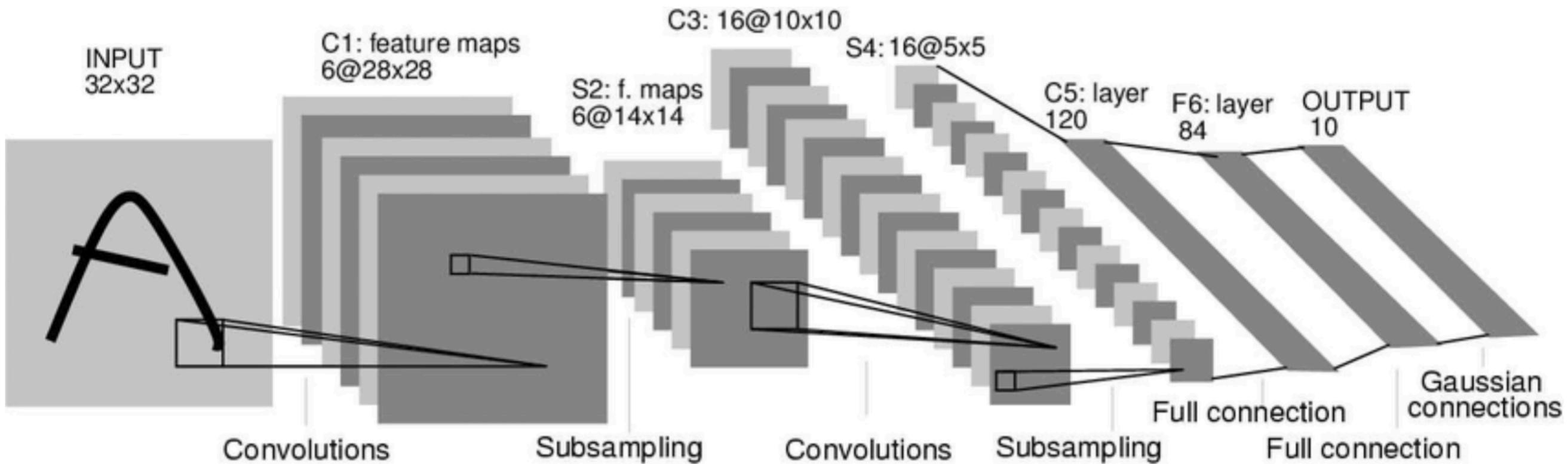
Convolutional Neural Networks



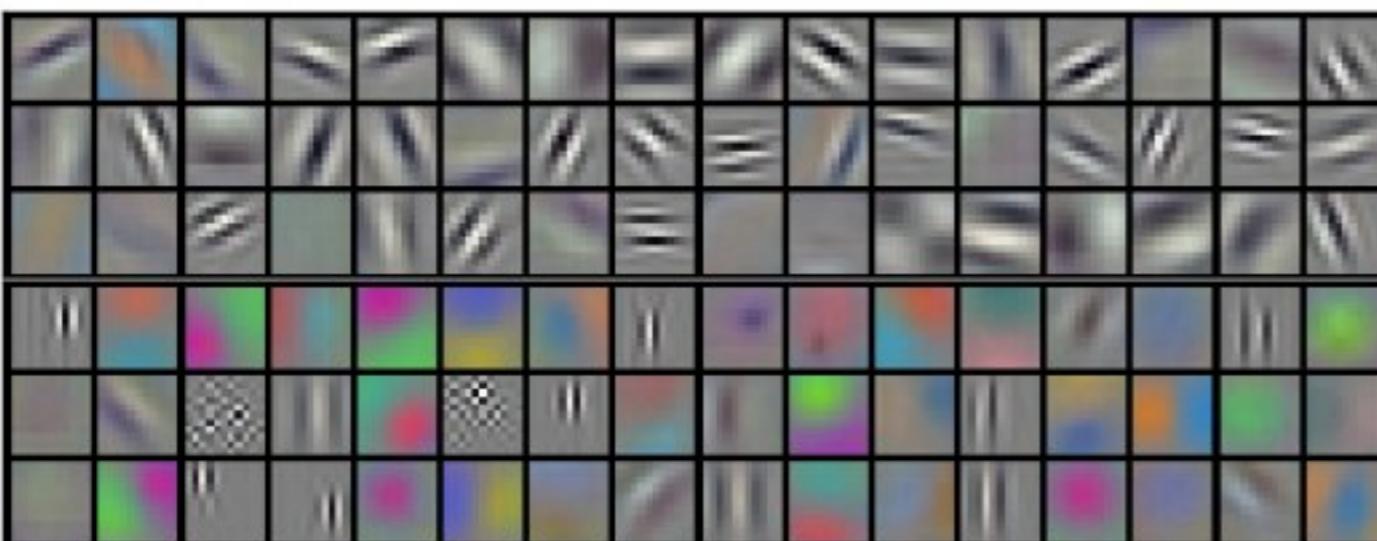


Convolution (with a kernel that is a parameter to be fit)

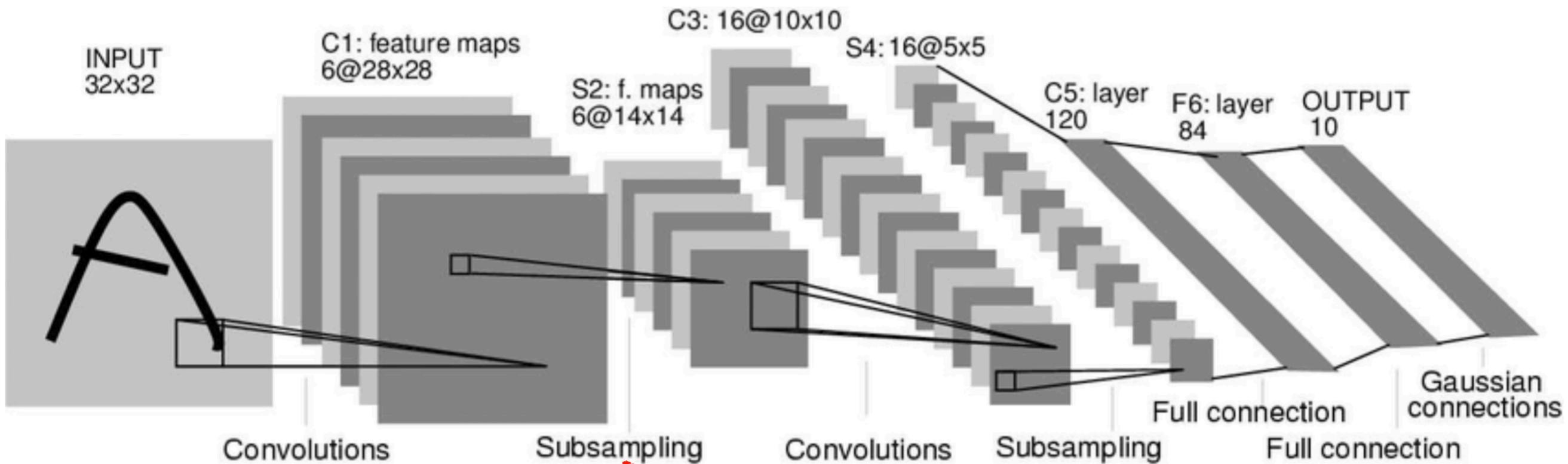




feature kernels



**flattening
and
fully connected
neural network**



pooling (max, min, mean...)

steps in convolution can be also used

12	20	30	0
8	12	2	0
34	70	37	4
112	100	25	12

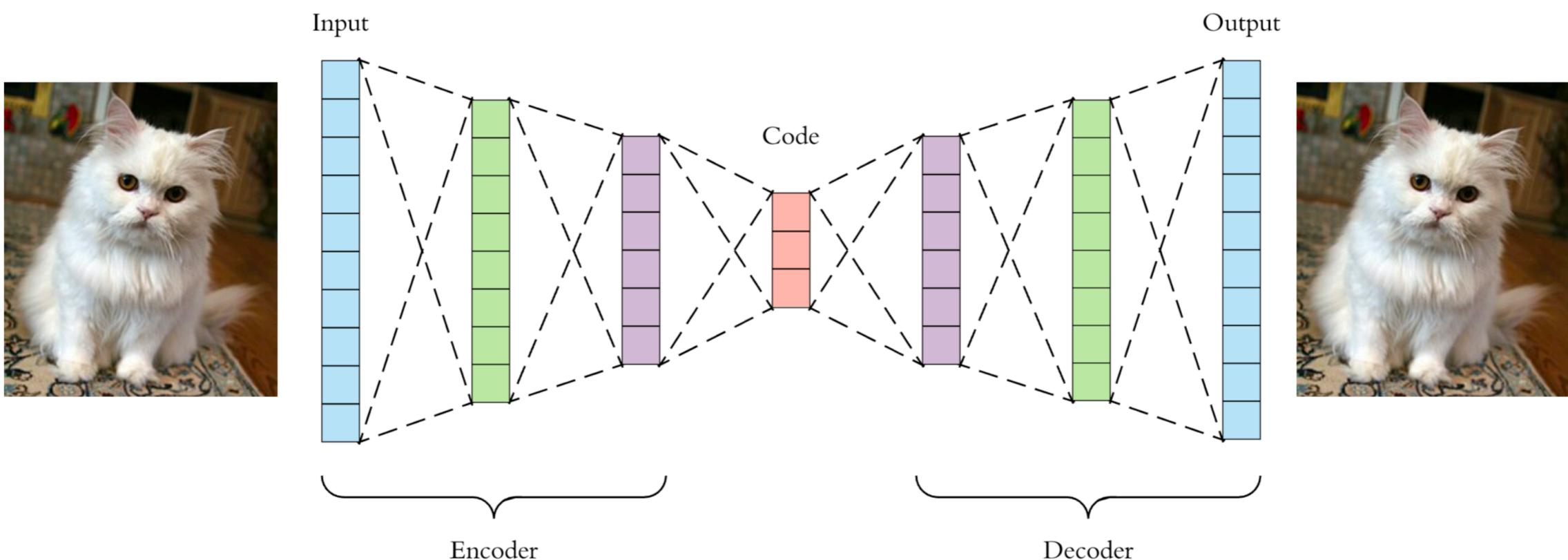
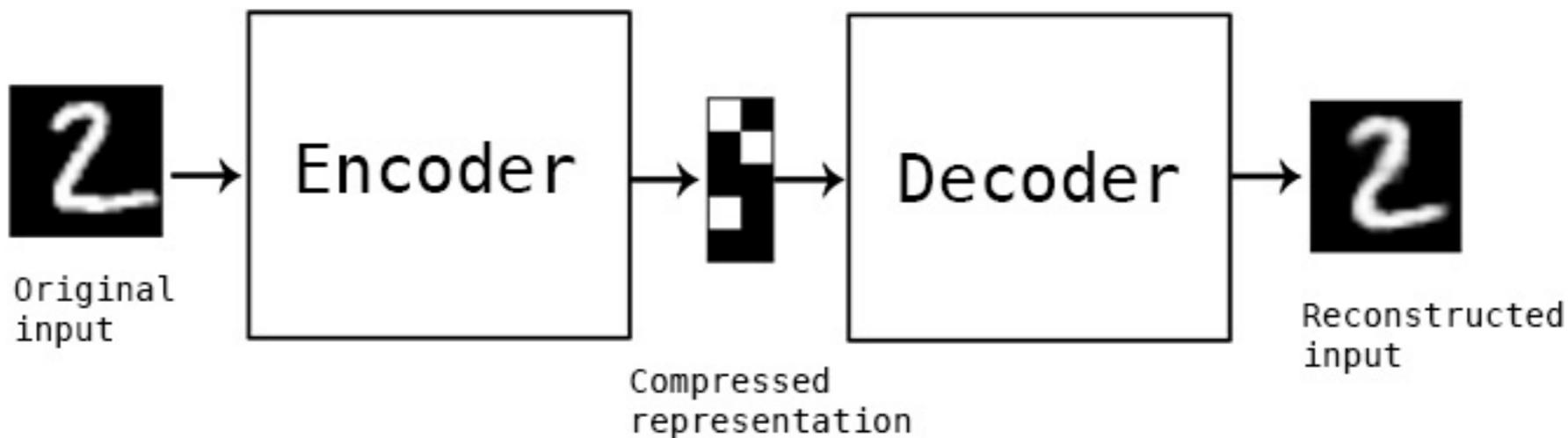
2×2 Max-Pool

20	30
112	37



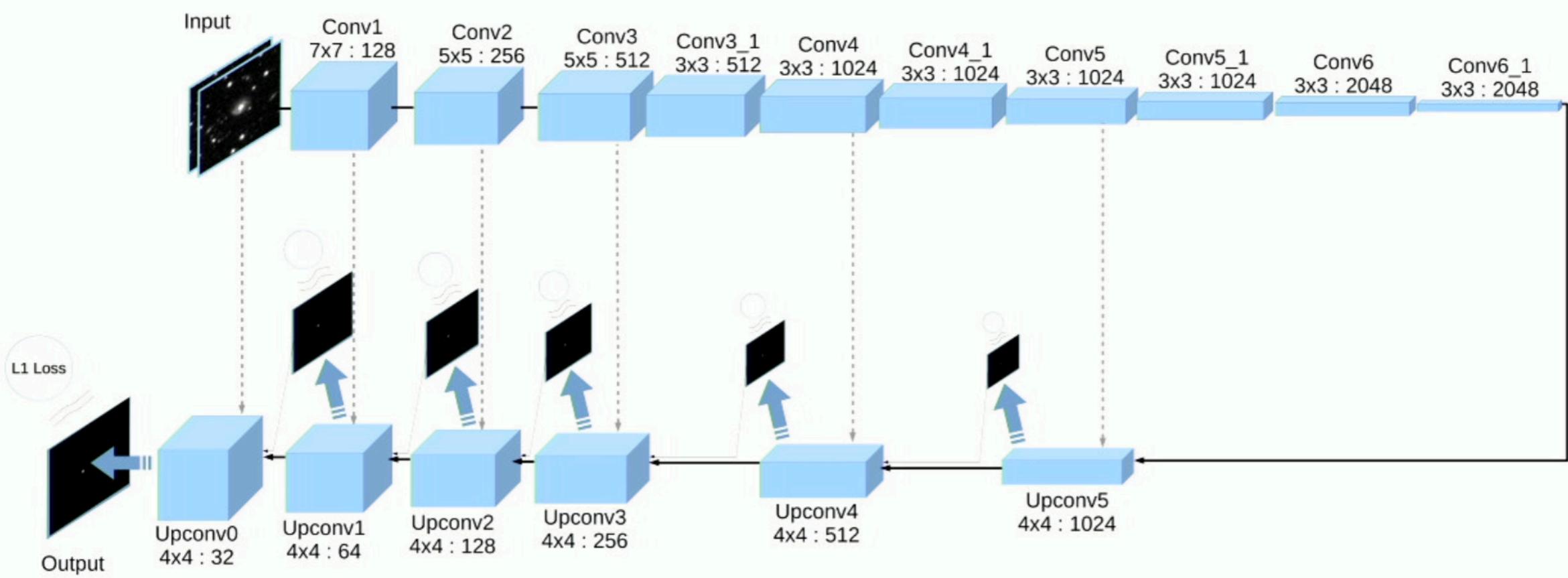
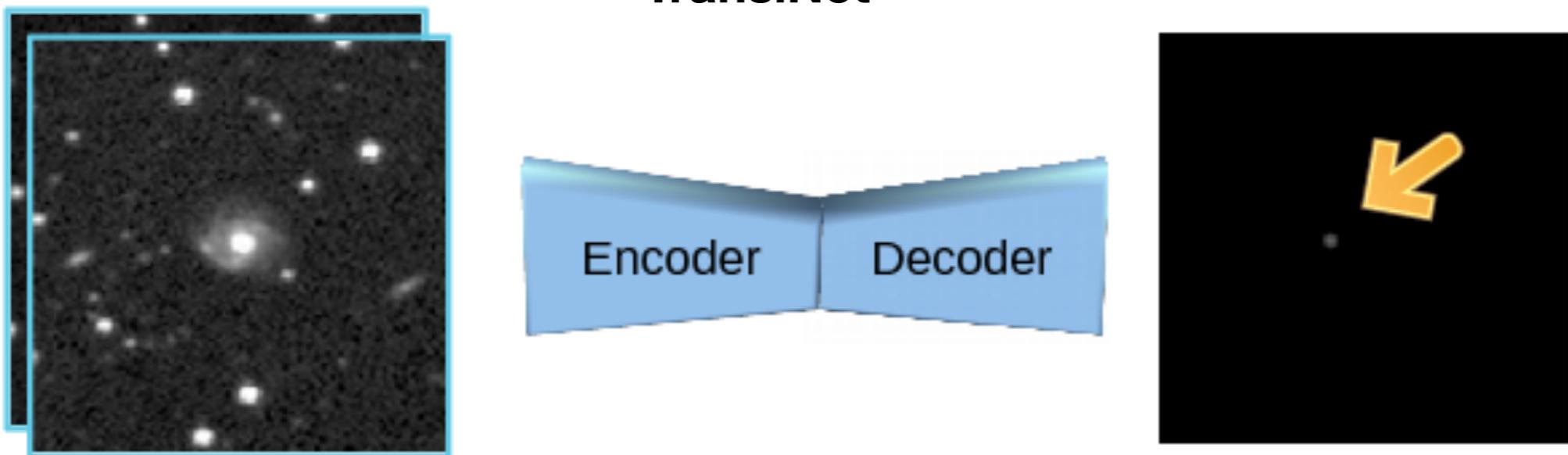
**Deep learning can identify the
patients gender with 95% accuracy!**

Encoder-decoders (Autoencoders)

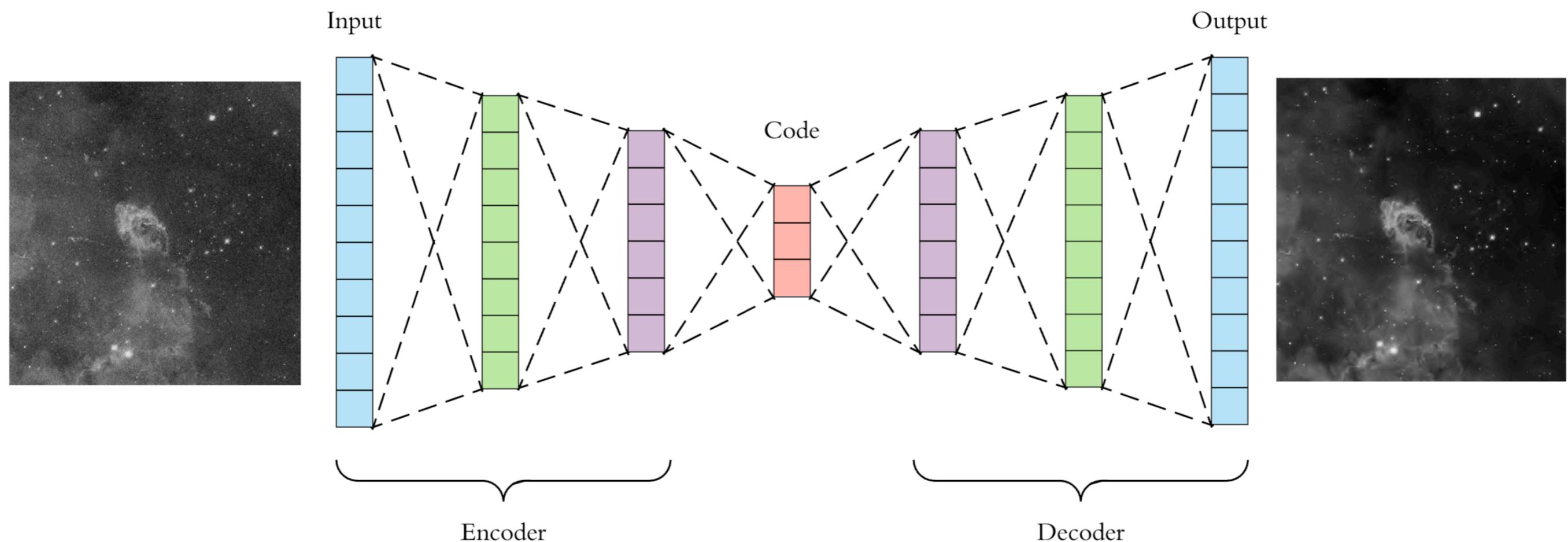


$$\text{Loss} = \sum (\text{original} - \text{product})^2$$

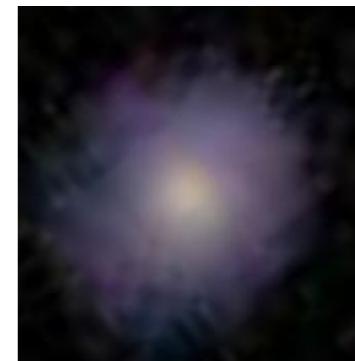
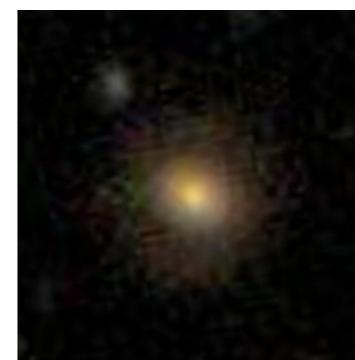
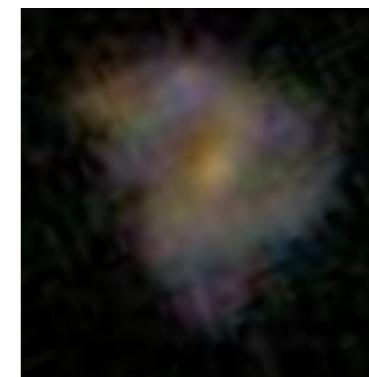
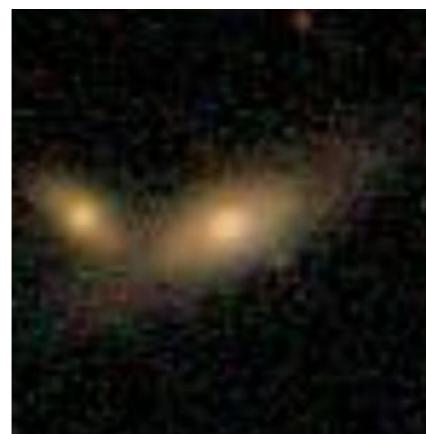
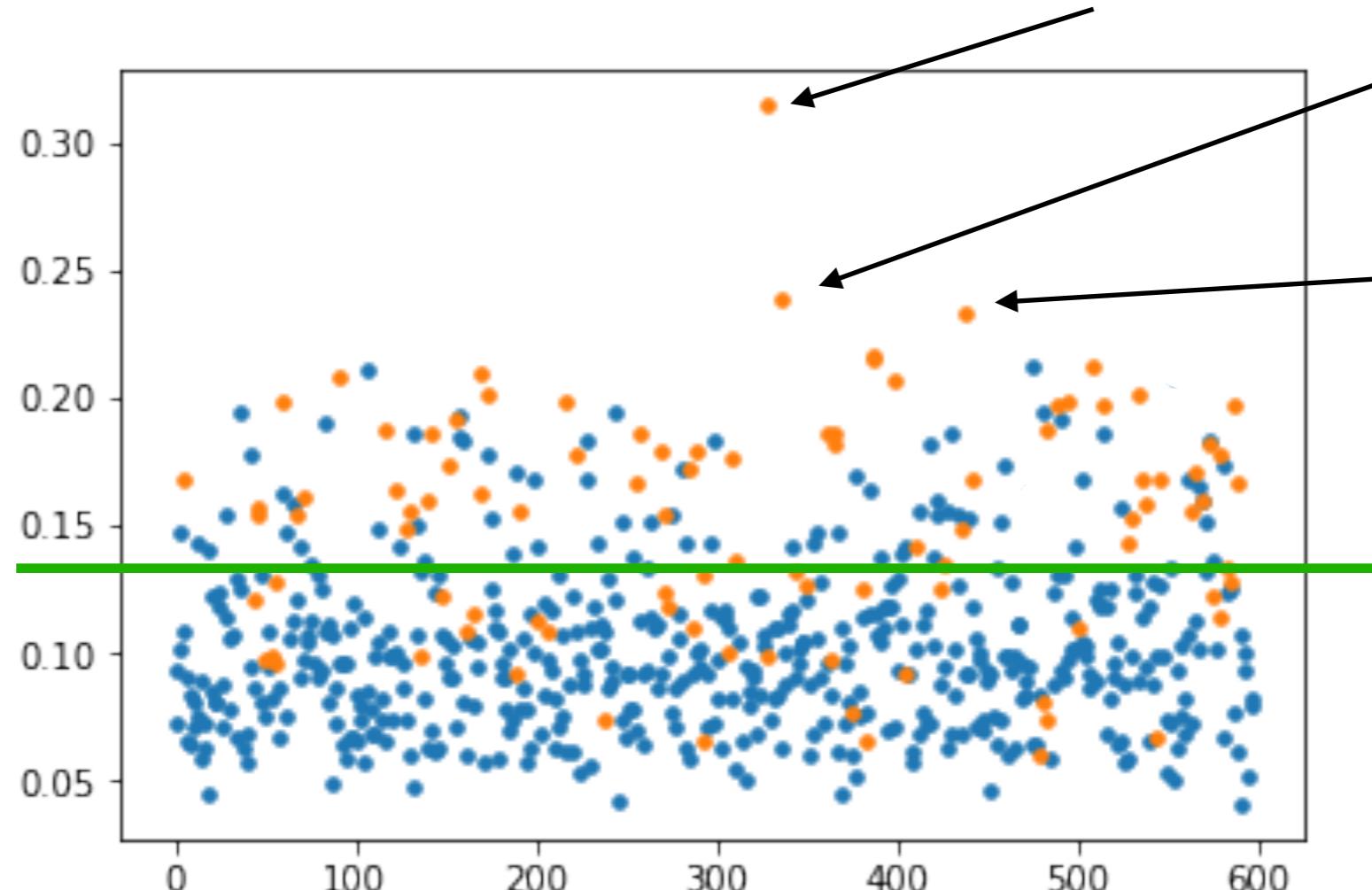
TransiNet



E.g. De-noising



Anomaly detection



Encoder/decoder usage

- Dimensionality reduction
- Denoising
- Outlier detection (measure of reconstruction error)

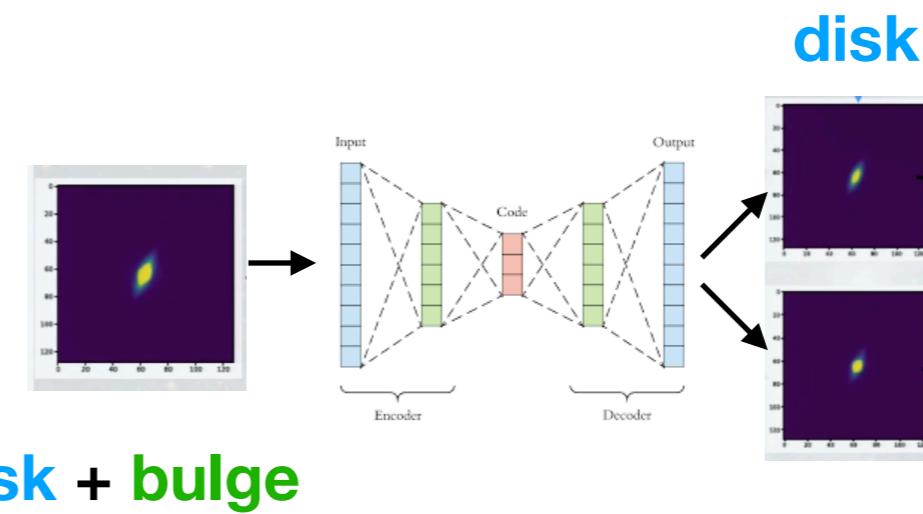
- Image segmentation

- Deblending

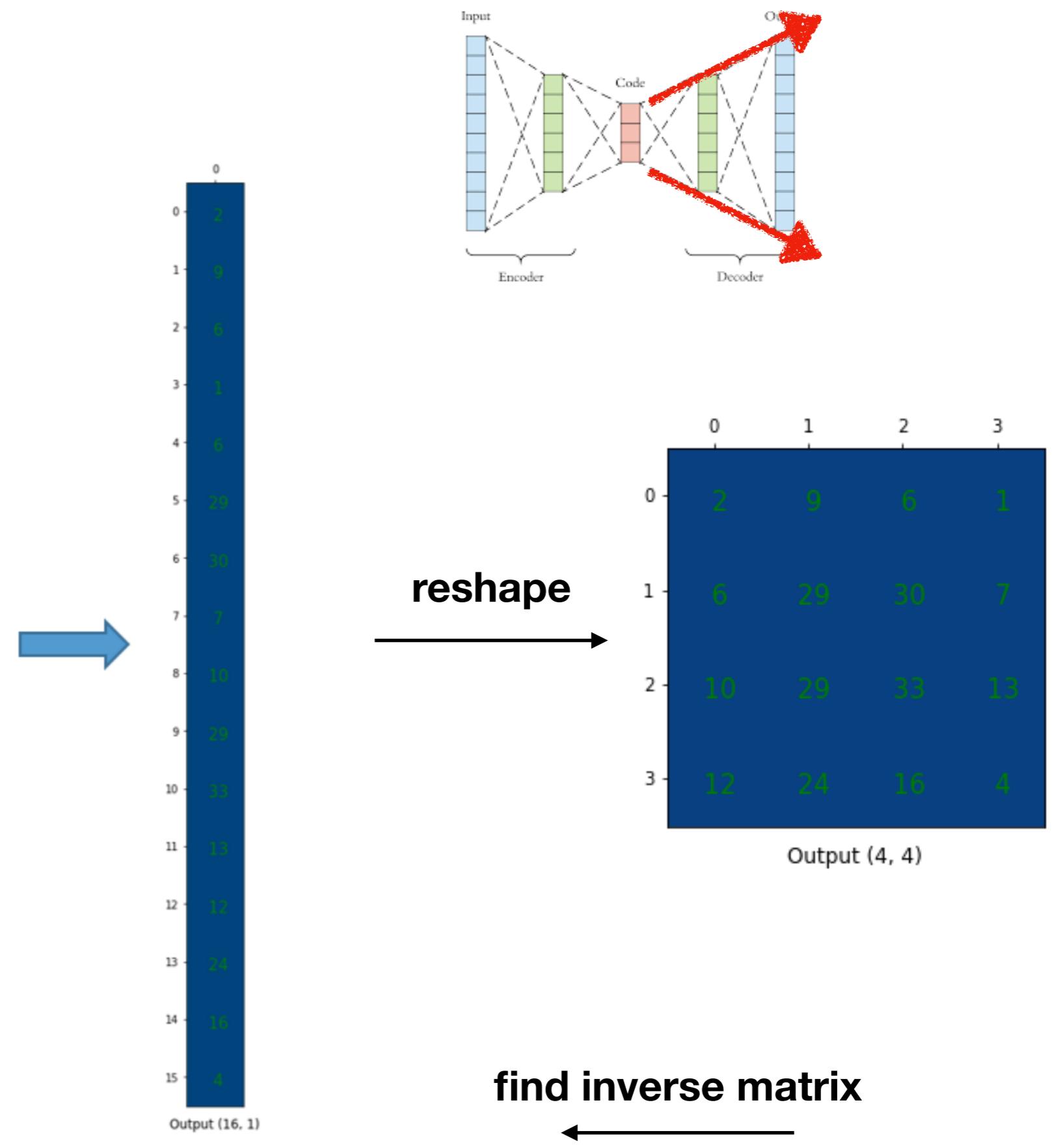
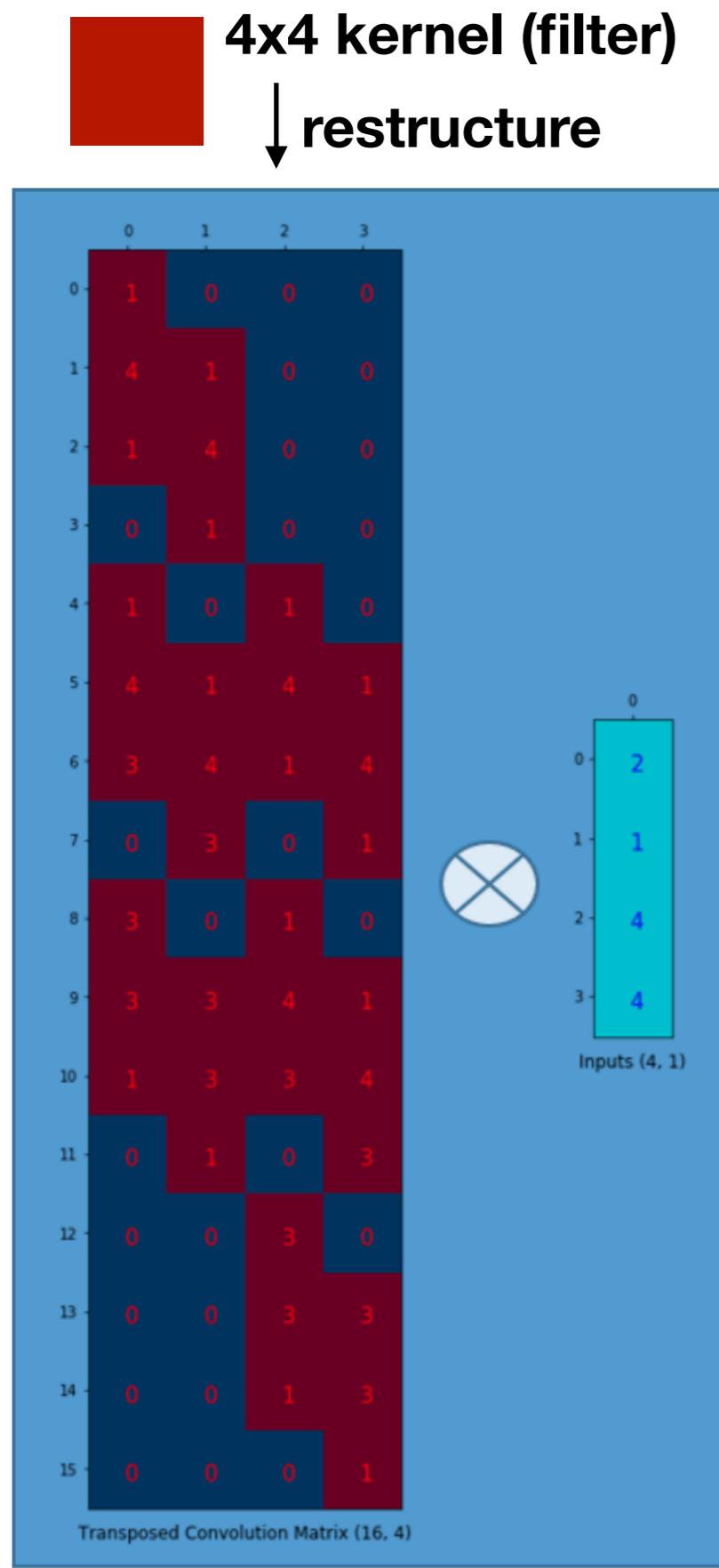
- Predicting next move

- Many others...

- (Can be considered as an unsupervised learning)



Transposed convolution (up or inverse convolution)



THIS IS YOUR MACHINE LEARNING SYSTEM?

YUP! YOU POUR THE DATA INTO THIS BIG
PILE OF LINEAR ALGEBRA, THEN COLLECT
THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

JUST STIR THE PILE UNTIL
THEY START LOOKING RIGHT.

