# Biomedical Data Science (GersteinLab.org/courses/452)
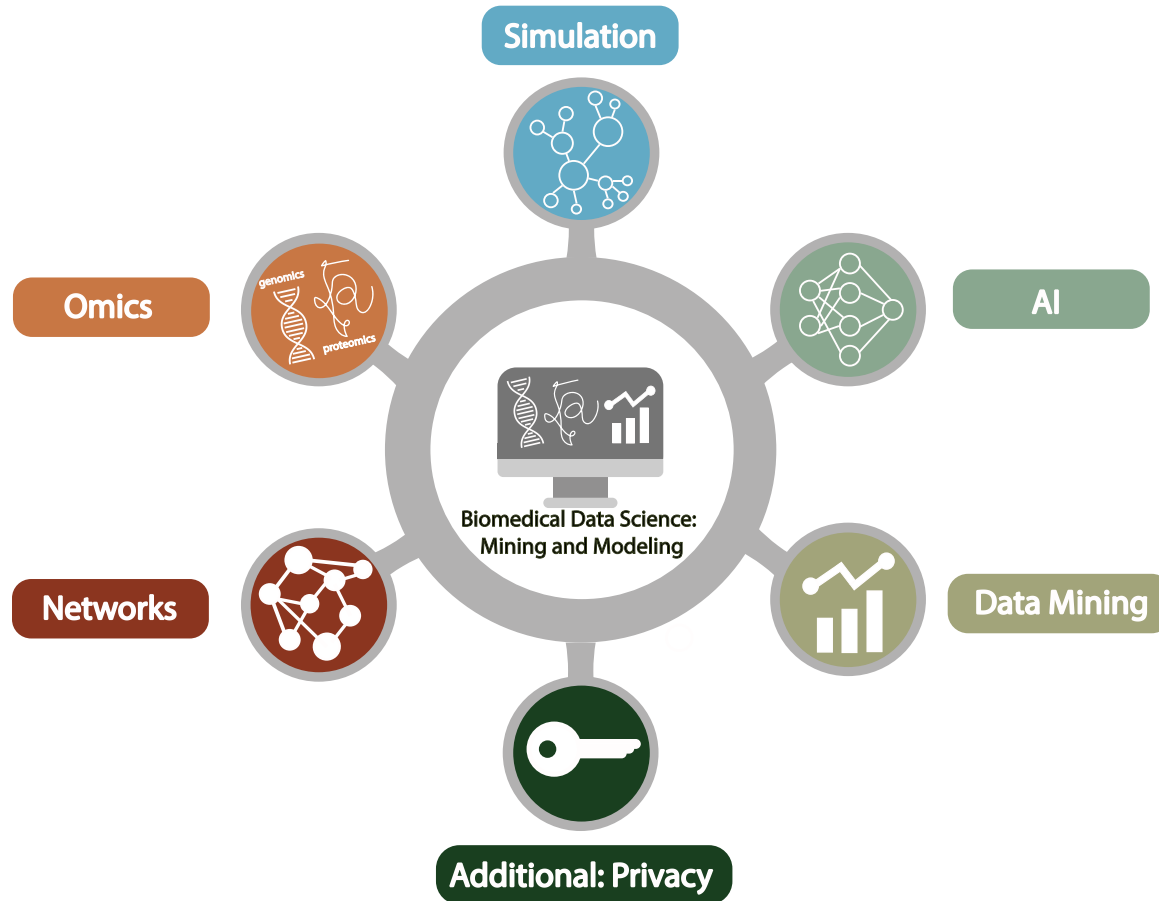
# Supervised Mining: Preliminaries + Decision Trees
### (25m8a+25m8b)
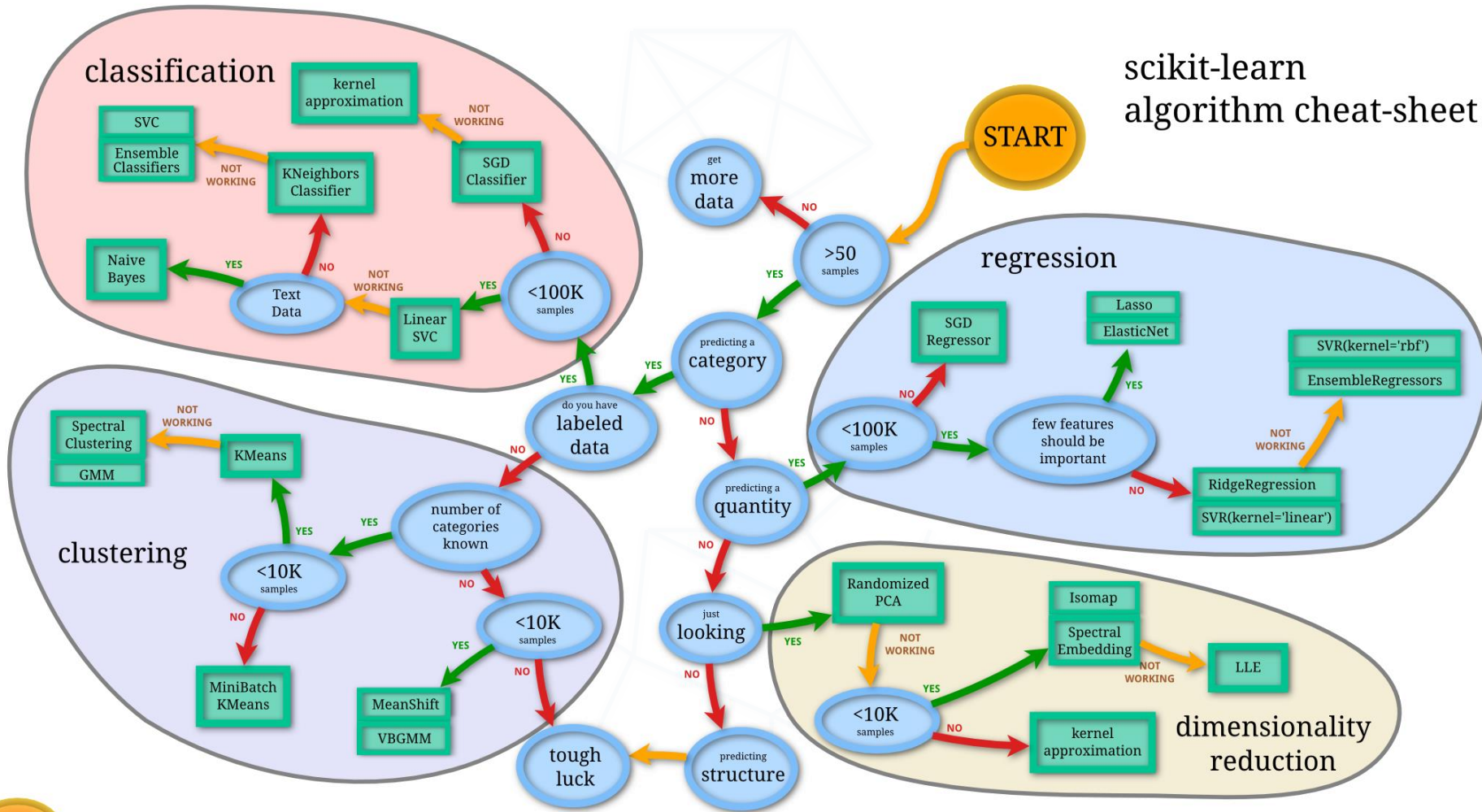
Mark Gerstein
Yale U.

Last edit in spring '25. Condensed (with ~2 slide deletions + a little more on RF) from 22m8a, Merging in 8b before DTs. Collectively similar to 2021's M8a & M8b [which have a videos].

# Supervised Mining:

# Overview

# The World of "Classic" ML



scikit-learn algorithm cheat-sheet

**classification**

kernel approximation

SVC

Ensemble Classifiers

KNeighbors Classifier

SGD Classifier

Naive Bayes

Text Data

Linear SVC

<100K samples

START

get more data

>50 samples

predicting a category

do you have labeled data

**regression**

SGD Regressor

Lasso ElasticNet

SVR(kernel='rbf')

EnsembleRegressors

<100K samples

few features should be important

RidgeRegression

SVR(kernel='linear')

**clustering**

Spectral Clustering

GMM

KMeans

number of categories known

<10K samples

MiniBatch KMeans

MeanShift

VBGMM

<10K samples

predicting a quantity

just looking

predicting structure

tough luck

**dimensionality reduction**

Randomized PCA

Isomap

Spectral Embedding

LLE

<10K samples

kernel approximation

Back

scikit learn

SciKit learn: http://scikit-learn.org/stable/tutorial/machine_learning_map/

# Distinctions in Supervised Learning

- Regression vs Classification
  - Regression: labels are quantitative
  - Classification: labels are categorical

- Regularized vs Un-regularized
  - Regularized: penalize model complexity to avoid over-fitting
  - Un-regularized: no penalty on model complexity

- Parametric vs Non-parametric
  - Parametric: an explicit parametric model is assumed
  - Non-parametric: otherwise

- Ensemble vs Non-ensemble
  - Ensemble: combines multiple models
  - Non-ensemble: a single model

# Structure of Genomic Features Matrix

# Represent predictors in abstract high dimensional space

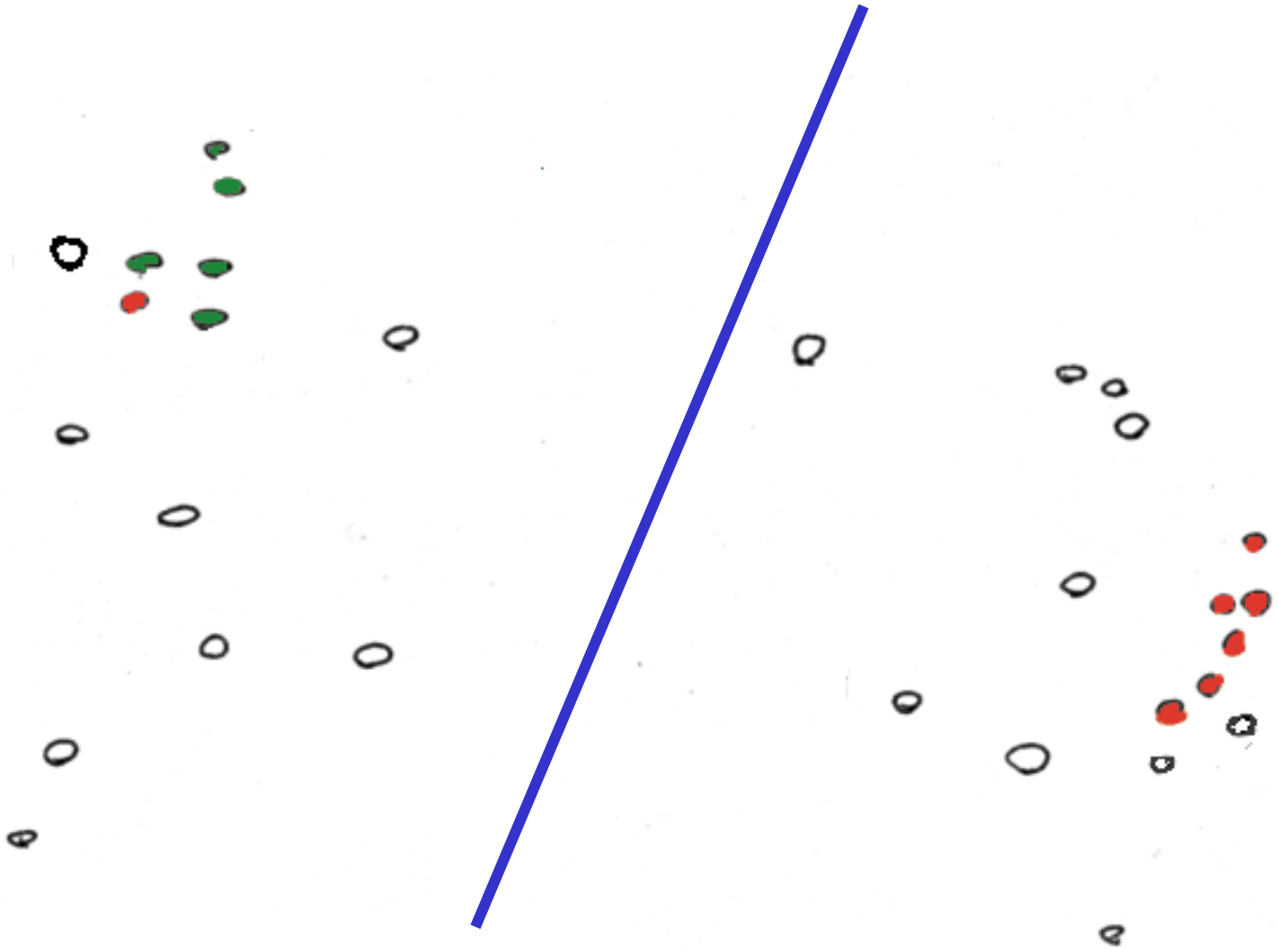# "Label" Certain Points

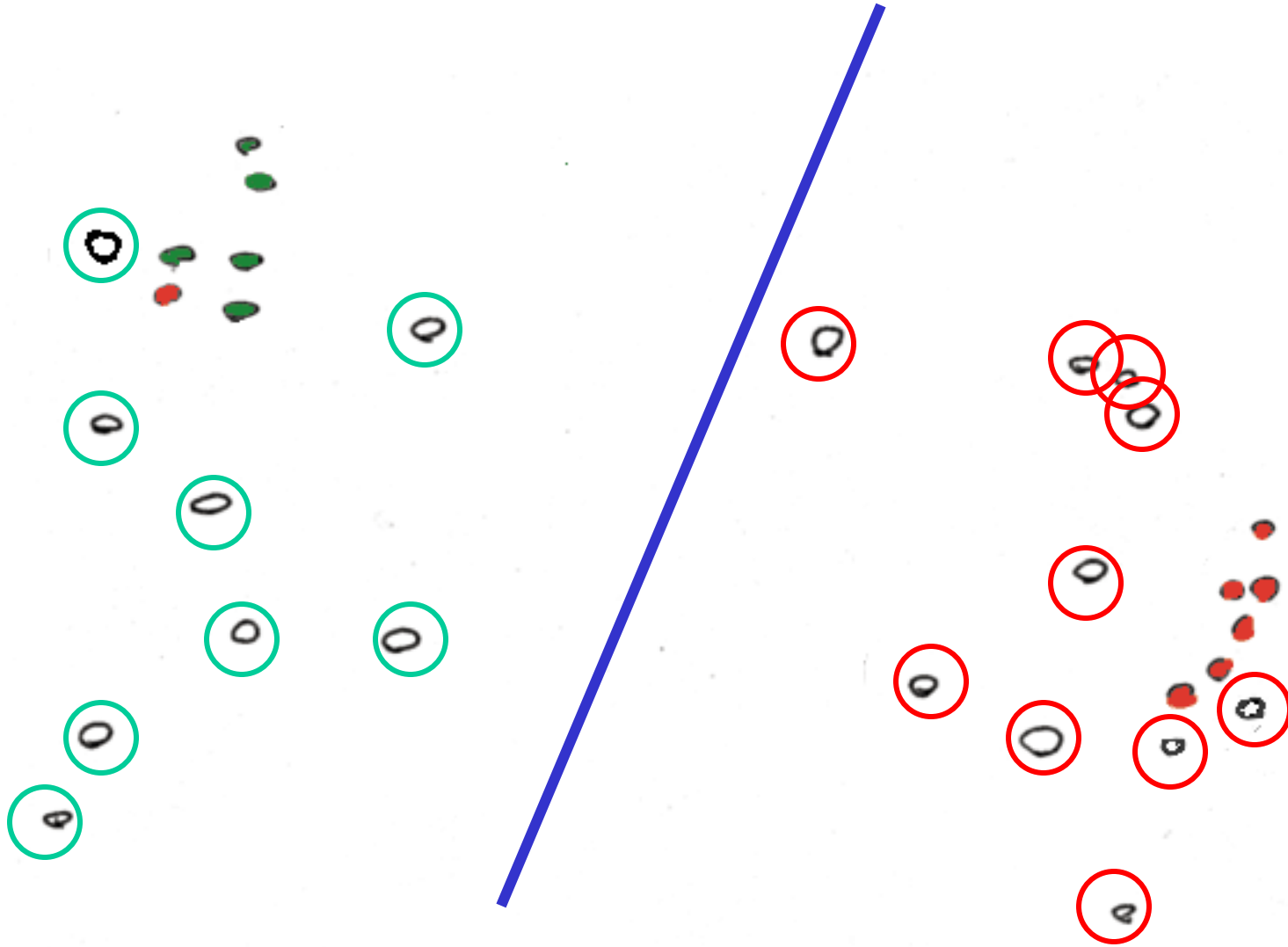# "Cluster" predictors (Unsupervised)

# Use Clusters to predict Response
## (Unsupervised, guilt-by-association)

# Find a Division to Separate Tagged Points

# Extrapolate to Untagged Points

# Find a Division to Separate Tagged Points

# Supervised Mining:

# Assessment, Cross-Validation & ROC Curves
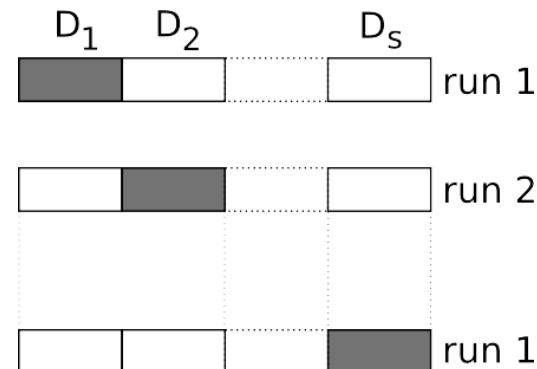
# Evaluating performance: What? How?

A. *__What__ do we want to evaluate?*

## GENERALIZATION

Therefore, it is mandatory to divide your dataset:

| TRAIN | VALIDATION | TEST |
|-------|------------|------|

Alternatively, use Cross Validation:

# B. _How do we evaluate performance?_

## 1. Classification problems



_Accuracy_

TP+TN/(TP+FP+FN+TN)

_Sensitivity (or TPR)_

TP/P=TP/(TP+FN)

_Specificity_

TN/N=TN/(TN+FP)

_Positive predictive value (PPV)_

TP/(TP+FP)

_False positive rate (FPR)_

FP/N=FP/(FP+TN)

_False discovery rate (FDR)_

FP/(FP+TP)

## 2. Regression problems  Sum of squares error

Root Mean Square error

_ROC analysis is good for comparing binary classifiers_

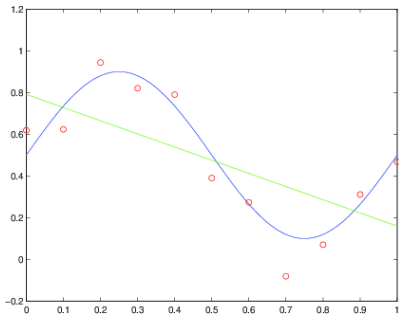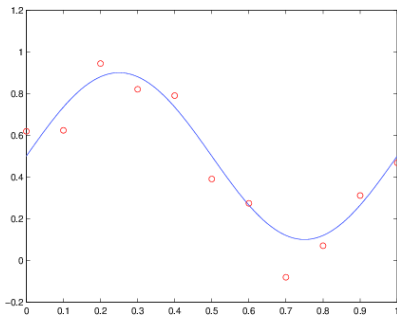https://en.wikipedia.org/wiki/Sensitivity_and_specificity

# Model dimensionality and overfitting
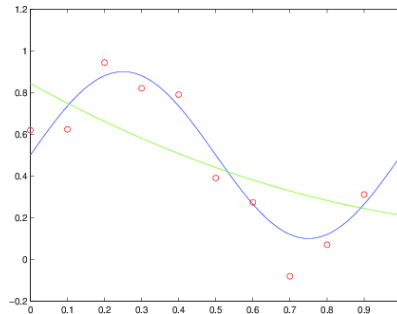


We are given the red dots.

We assume that they are noisy samples from a signal/(function) – the blue curve – which we do not have (we only have the red dots).

We want to predict new points, i.e. the $y$ coordinates for other values of $x$ (e.g. $x > 1$)
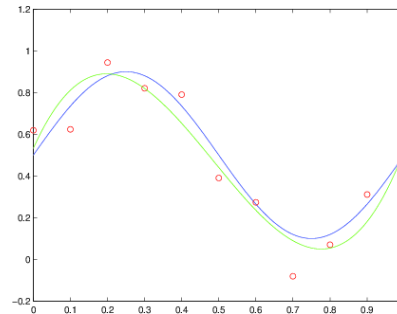
Our model needs to approximate the blue function. We decide to do it with polynomials.
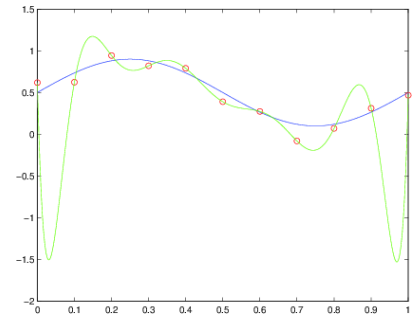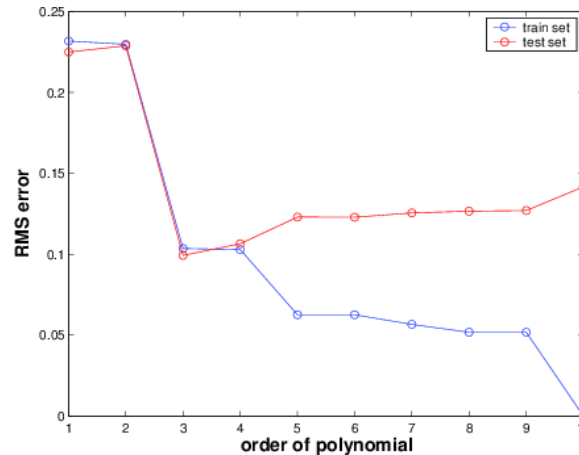


Degree 1 polynomial

Degree 2 polynomial

Degree 3 polynomial

Degree 10 polynomial

Which one is best? And why?

How does the GENERALIZATION performance vary, as we increase the complexity of the polynomial?



- **Occam's razor** *(William of Occam, ~1300):* **Accept the simplest explanation that fits the data.**

We should prefer simpler models to more complex models, and this preference should be traded off against the extent to which the model fits the data.
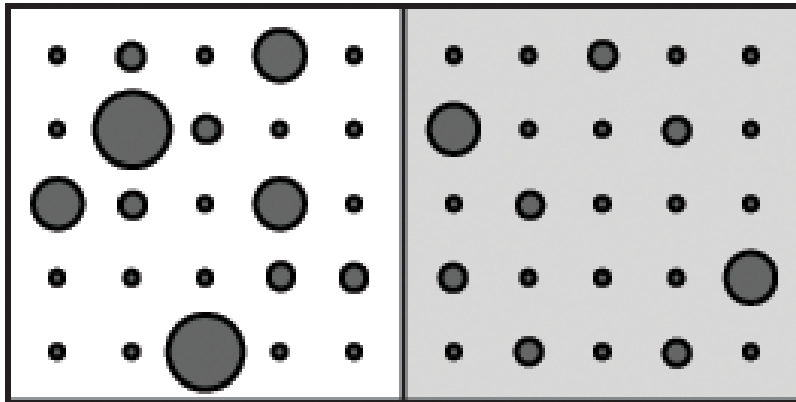
Related to "Bias-Variance" tradeoff.

- IMPORTANT: increasing the number of features may lead to a reduction in performance if the number of datapoints is not increased. Why?
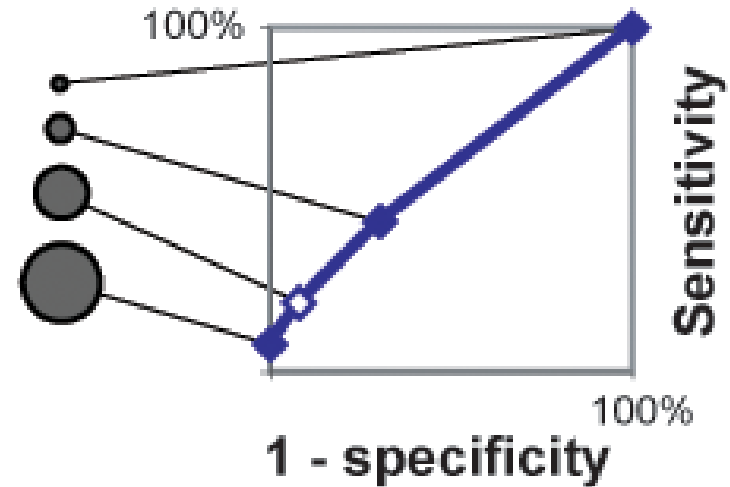


This is related to the "Curse of Dimensionality" Bellman, 1961.

# ROC plots & Comparison of Predictions against a Positive & Negative Gold Standard
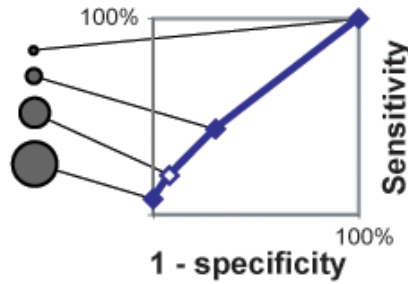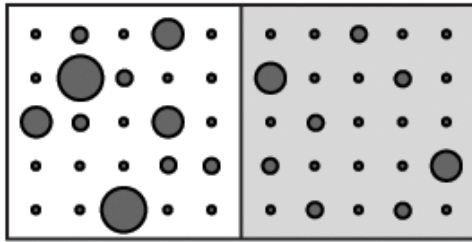
Threshold "predictions" (strength of positive score is represented by circle size) at different levels and compare to + and - gold standards (represented by white & gray squares). 50 total instances, half + and half -. A concrete example would be doing cancer prediction 50 individuals with known cancer status.



"Error Rate" (FP/N)
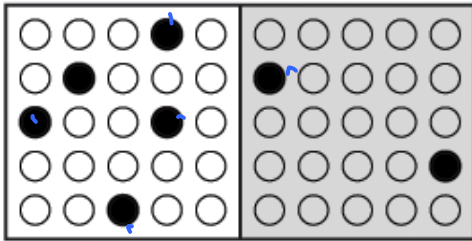
ROC plot
(cross validated)

Effect on Predictions of Large Number of Negatives

(e.g. terrorist identification or breast cancer screening)

**Sensitivity**



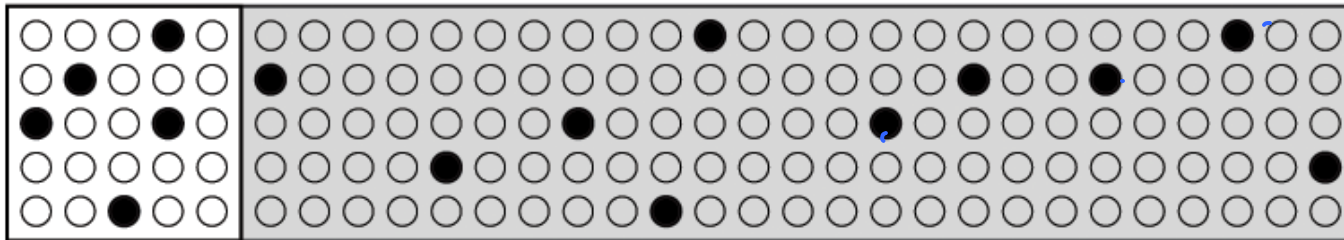$$\frac{5}{25} = 20\%$$

**1- specificity**

$$\frac{2}{25} = 8\%$$

**Positive predictive value**
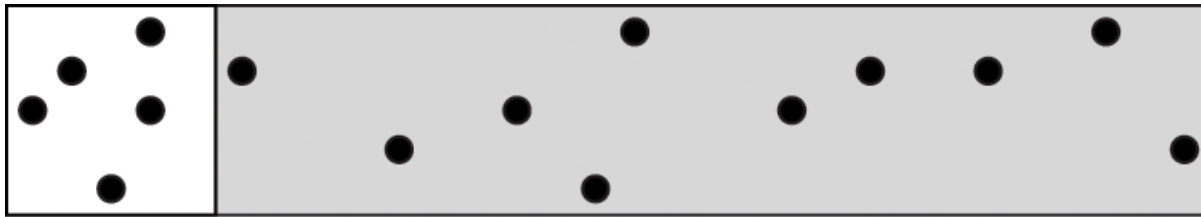
$$\frac{5}{5+2} \approx 71\%$$



$$\frac{5}{25} = 20\%$$

$$\frac{10}{125} = 8\%$$
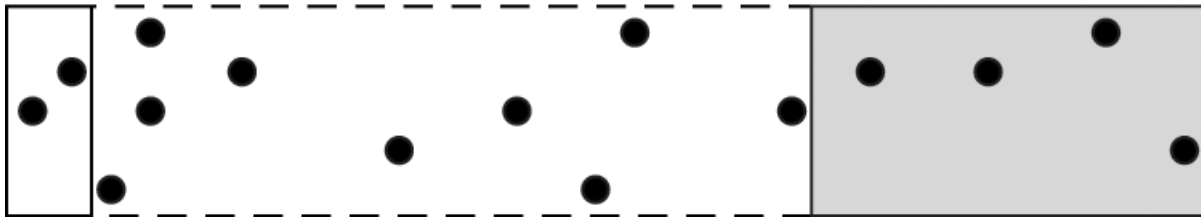
$$\frac{5}{5+10} \approx 33\%$$

# Importance of Balanced Positive and Negative Examples
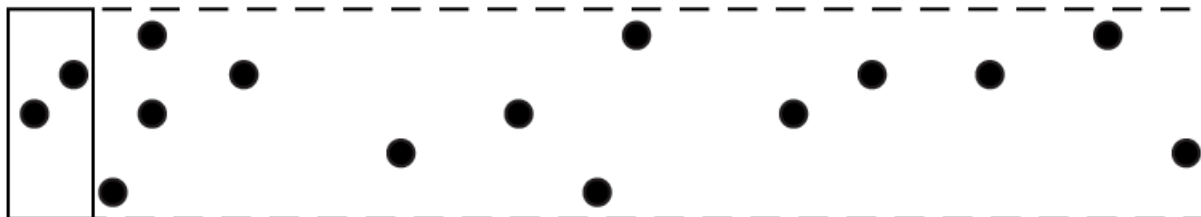
$$\frac{5}{?} = ? \qquad \frac{10}{?} = ? \qquad \frac{5}{5+10} \approx 33\%$$

$$\frac{2}{?} = ? \qquad \frac{4}{?} = ? \qquad \frac{2}{2+4} \approx 33\% \text{ (estimate)}$$

$$\frac{2}{?} = ? \qquad \frac{?}{?} = ? \qquad \frac{2}{2+?} = ?$$
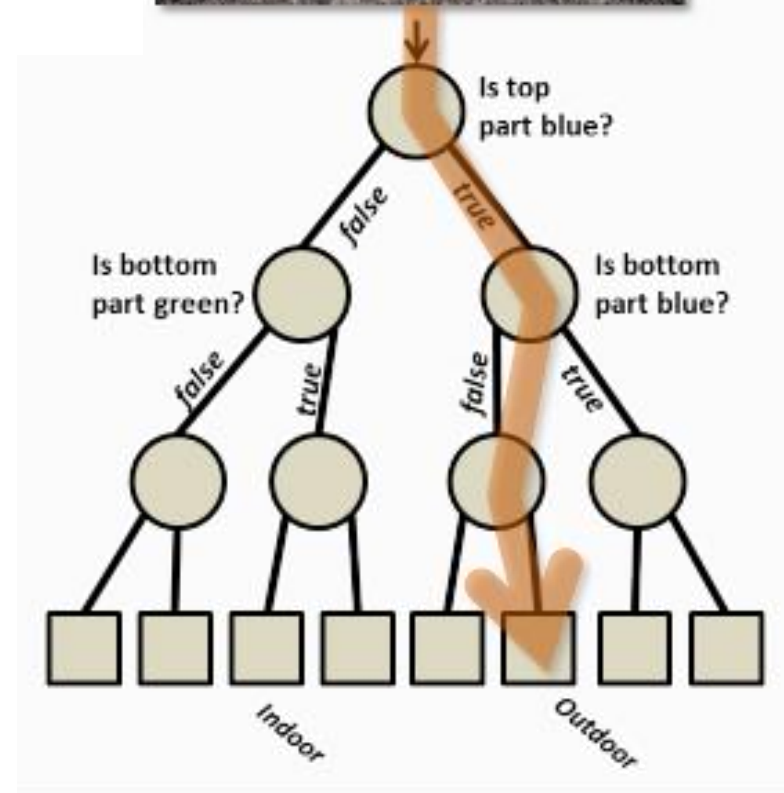
# Supervised Mining:

# Decision Trees

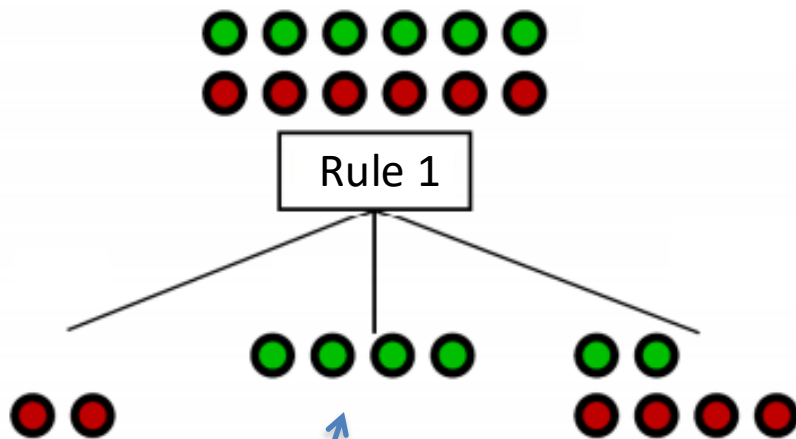**("Jumping to first method")**

# Decision Trees

- **Classify data by asking questions** that divide data in subgroups
- Keep asking questions until subgroups become homogenous
- Use **tree** of questions to make predictions



- Example: Is a picture taken inside or outside?
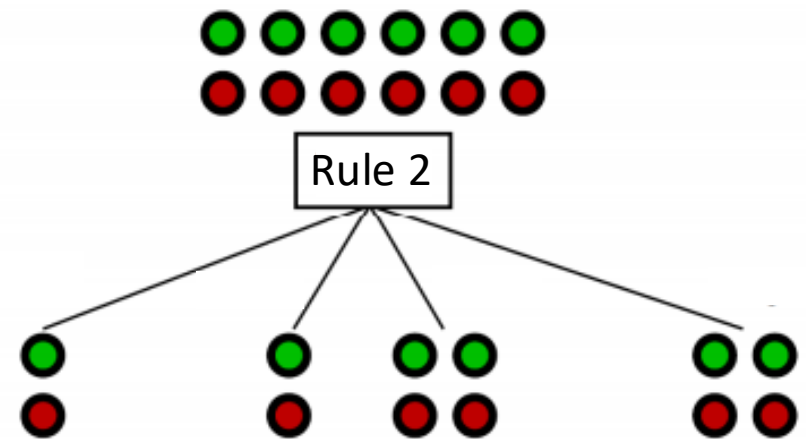
Criminisi, Shotton, and Konukoglu *Microsoft Technical Report* 2011

# What makes a good rule?

- Want resulting groups to be as homogenous as possible

2/3 Groups homogenous
→Good rule

All groups still 50/50
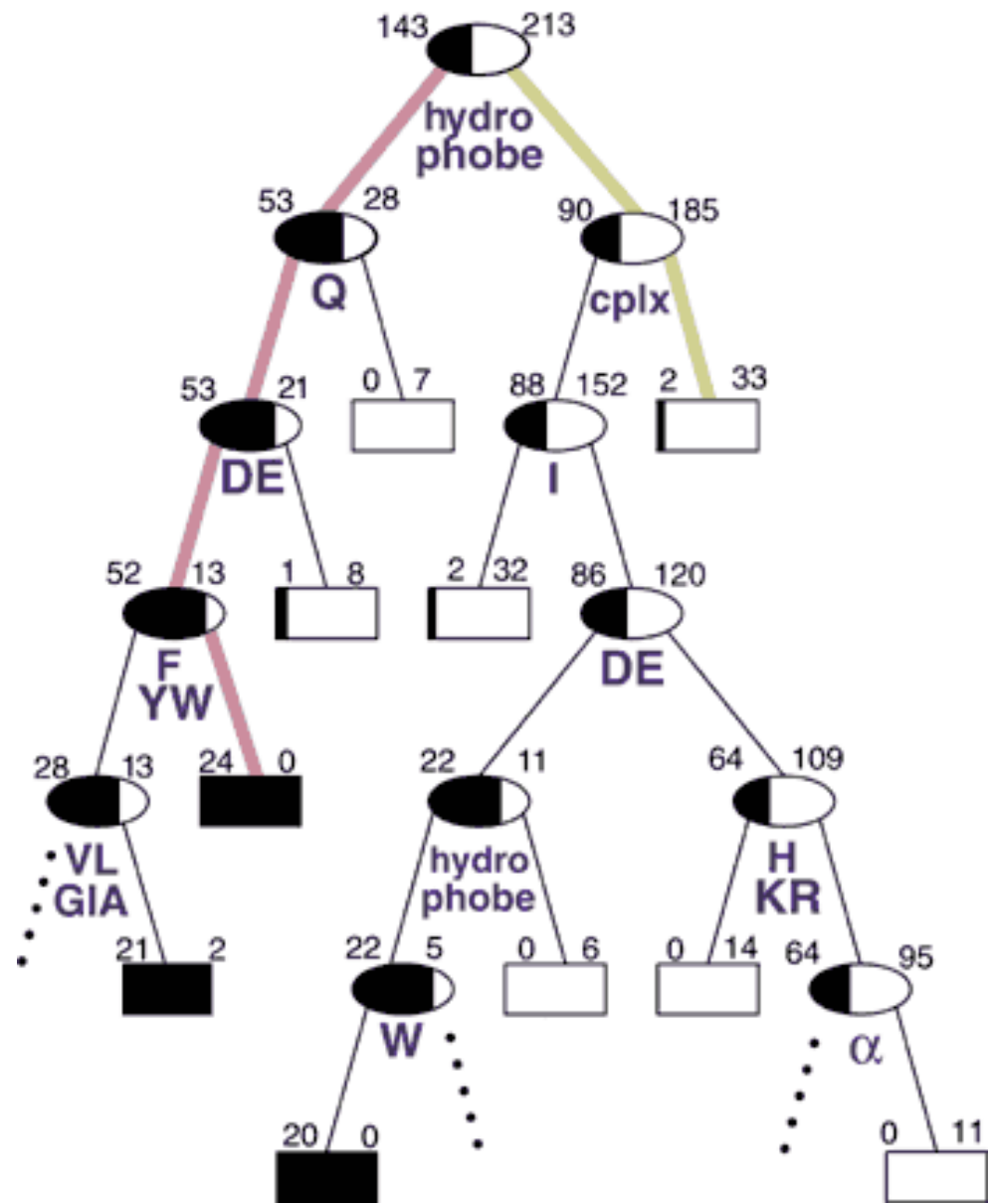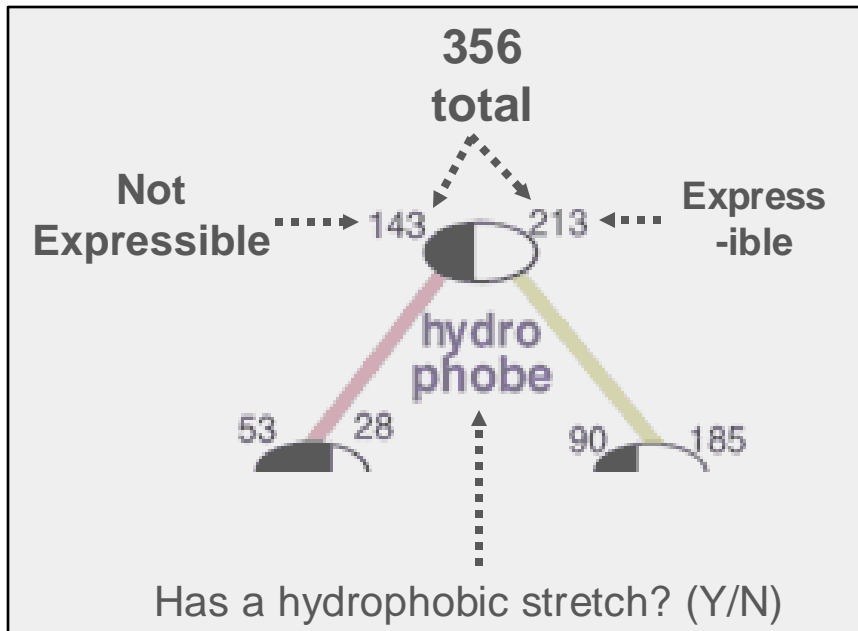→ Unhelpful rule

# Quantifying the value of rules

- ## Decrease in inhomogeneity (or increase in homogeneity)
  - Most popular metric: Information theoretic entropy

    $$S = -\sum_{i=1}^{m} p_i \log p_i$$

  - Use frequency of classifier characteristic within group as probability
  - Minimize entropy to achieve homogenous group

# Algorithm

- For each characteristic:
  - Split into subgroups based on each possible value of characteristic
- Choose rule from characteristic that maximizes decrease in inhomogeneity
- For each subgroup:
  - if (inhomogeneity < threshold):
    - Stop
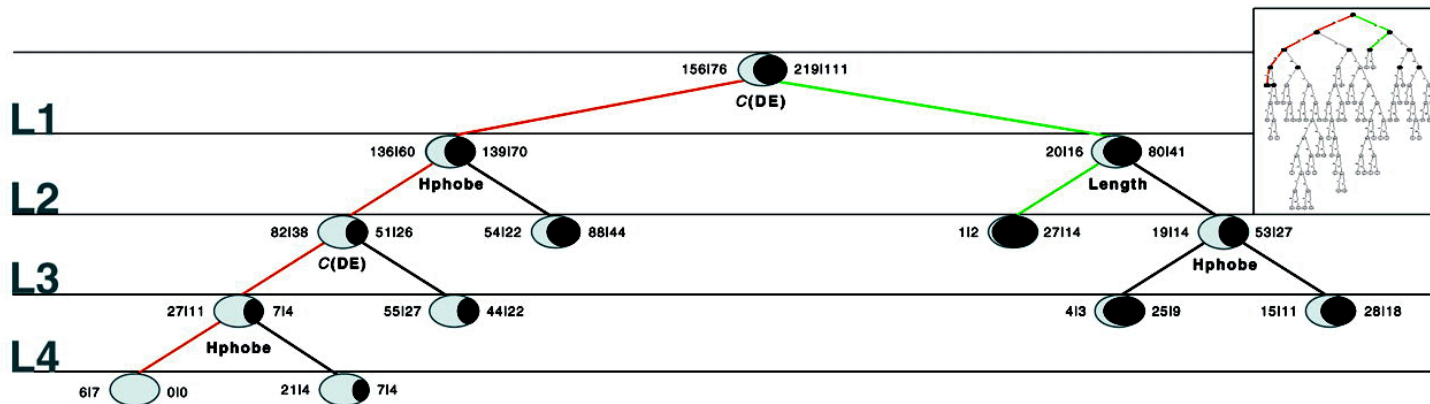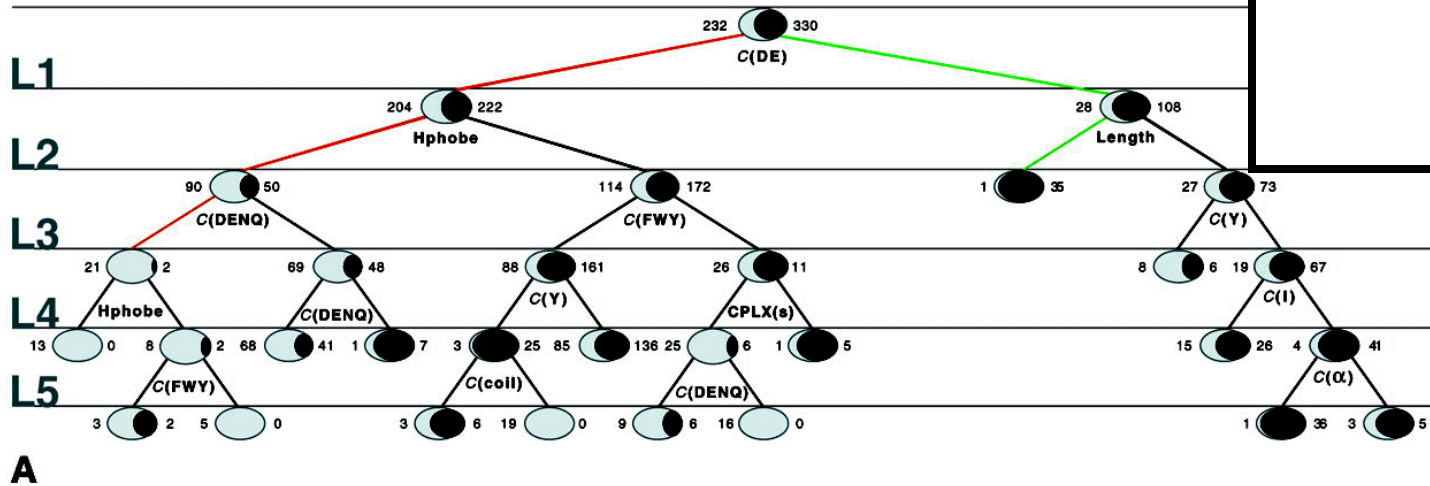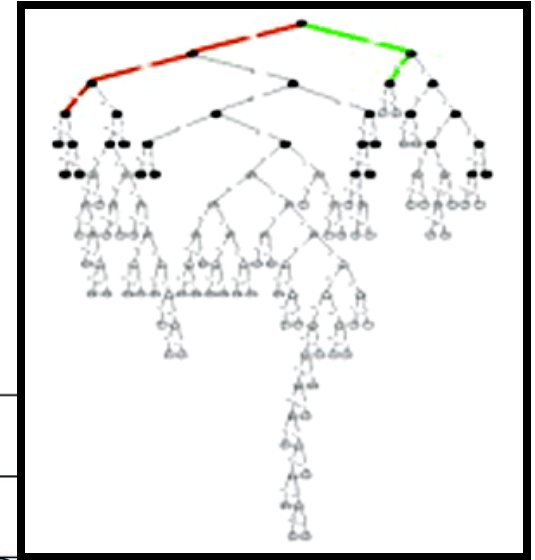  - else:
    - Restart rule search (recursion)

# Retrospective Decision Trees



Analysis of the Suitability of 500 M. thermo. proteins
to find optimal sequences purification

[Bertone et al. NAR ('01)]

# Overfitting, Cross Validation, and Pruning

# Random Forest (RF)

- Basic decision tree (DT) method is very sensitive to dataset selection & noise in the data

- RFs are ensemble of DTs; address this issue
  - Build many DTs on bootstrapped training samples. (Reduces sensitivity to noise.)
  - Each time a split in a tree is considered, a random sample of m predictors is chosen as split candidates from the full set of predictors. (Decorrelates "bagged" DTs.)
  - Finally, we average or vote amongst the trees

# References

- James, Gareth, Witten, Daniela, Hastie, Trevor, Tibshirani, Robert
  An Introduction to Statistical Learning: with Applications in R
  [ ISLR (2nd edition) ]
  https://www.amazon.com/Introduction-Statistical-Learning-Applications-Statistics/dp/1071614177/  +  https://www.statlearning.com
  (Chap 2 gives a nice overview on key concepts in ML.
  Chapter 8 to 8.2.2 gives background on DTs.)

- Greener, J. G., Kandathil, S. M., Moffat, L., & Jones, D. T. (2021).
  A guide to machine learning for biologists.
  Nature Reviews Molecular Cell Biology, 23(1), 40–55.
  https://doi.org/10.1038/s41580-021-00407-0
  (Good reference, but for this pack just go to up to section on "key concepts.")

- Agarwal, R. (2024, March 29).
  ROC Curves and AUC: The Ultimate guide. Built In.
  https://builtin.com/data-science/roc-curves-auc
  (Optional extra background on ROC)