

Classifiers!!!

BCH394P/364C Systems Biology / Bioinformatics
Edward Marcotte, Univ of Texas at Austin

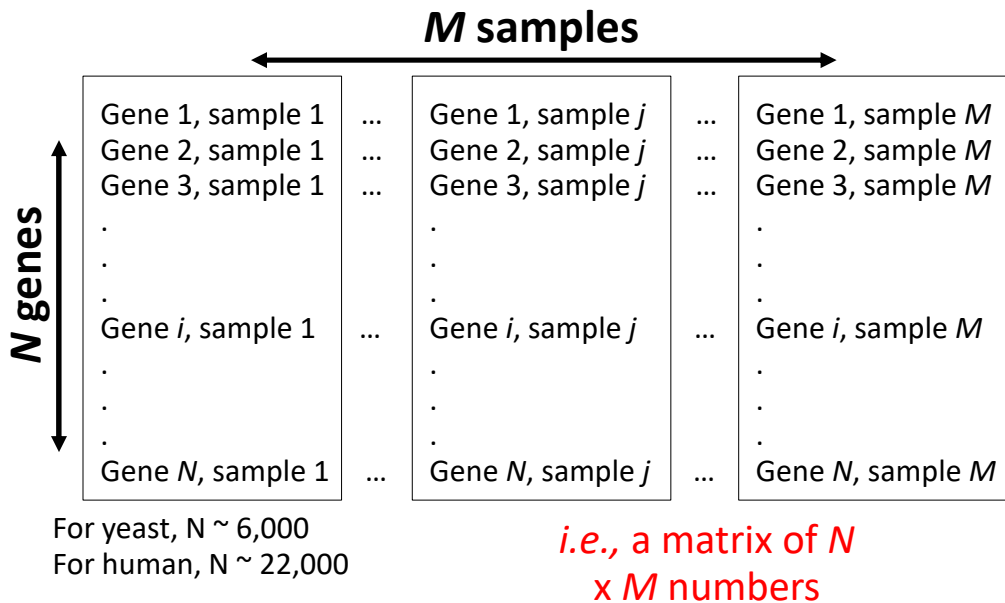
Clustering = task of grouping a set of objects in such a way that objects in the same group (a **cluster**) are more similar (in some sense) to each other than to those in other groups (clusters).

VS.

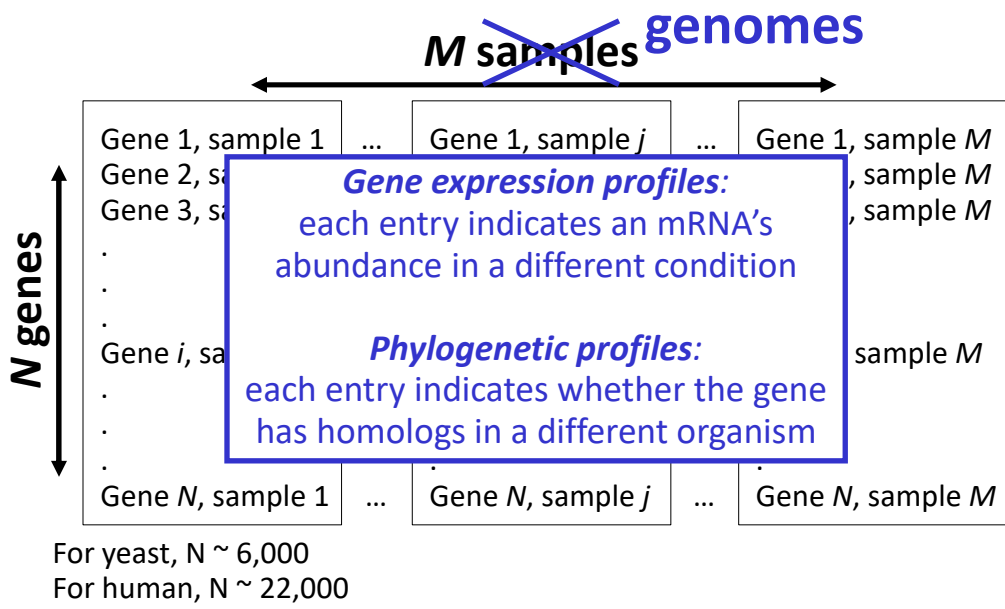
Classification = task of categorizing a new observation, on the basis of a training set of data with observations (or instances) whose categories are known

Adapted from Wikipedia

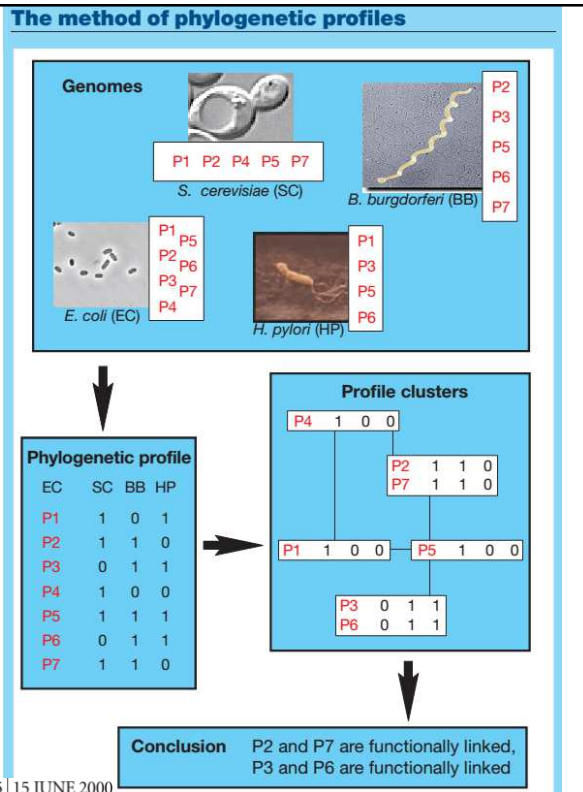
Remember, for clustering, we had a matrix of data...



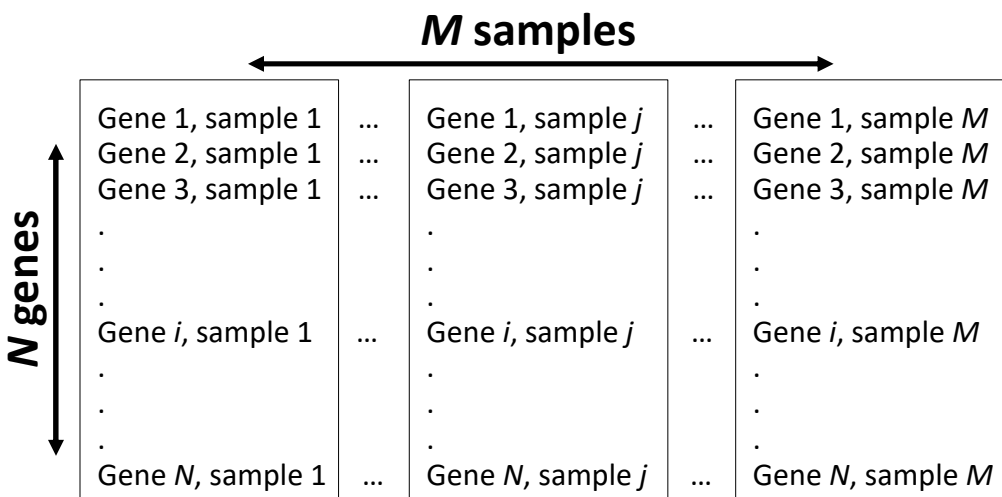
We discussed gene expression profiles. Here's another example of gene features.



This is useful because biological systems tend to be modular and often inherited intact across evolution. (e.g. you tend to have a flagellum or not)



Many such features are possible...



For yeast, $N \sim 6,000$
 For human, $N \sim 22,000$

i.e., a matrix of $N \times M$ numbers

We also needed a measure of the similarity between feature vectors. Here are a few (of many) common distance measures used in clustering.

Names	Formula
Euclidean distance	$\ a - b\ _2 = \sqrt{\sum_i (a_i - b_i)^2}$
Manhattan distance	$\ a - b\ _1 = \sum_i a_i - b_i $
cosine similarity	$\frac{a \cdot b}{\ a\ \ b\ }$

Wikipedia

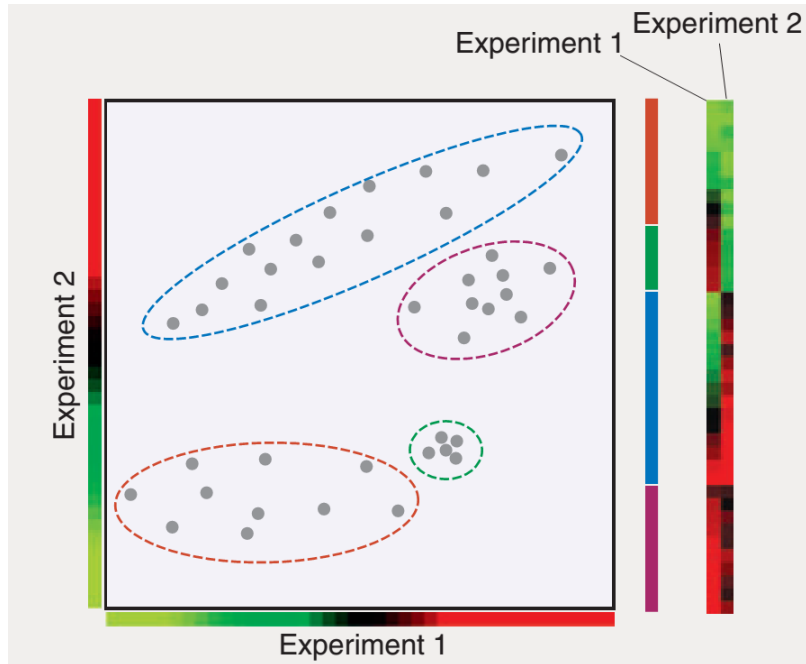
We also needed a measure of the similarity between feature vectors. Here are a few (of many) common distance measures used in clustering.

~~classifying~~

Names	Formula
Euclidean distance	$\ a - b\ _2 = \sqrt{\sum_i (a_i - b_i)^2}$
Manhattan distance	$\ a - b\ _1 = \sum_i a_i - b_i $
cosine similarity	$\frac{a \cdot b}{\ a\ \ b\ }$

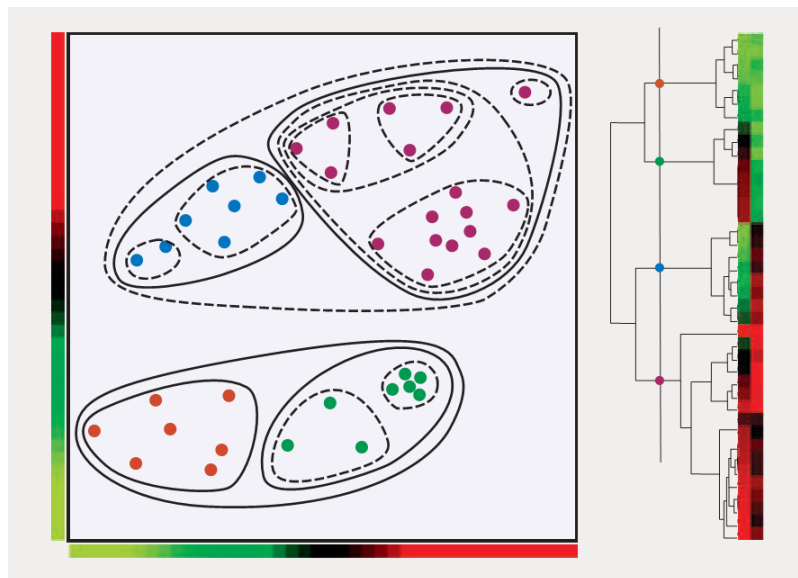
Wikipedia

Clustering refresher: 2-D example



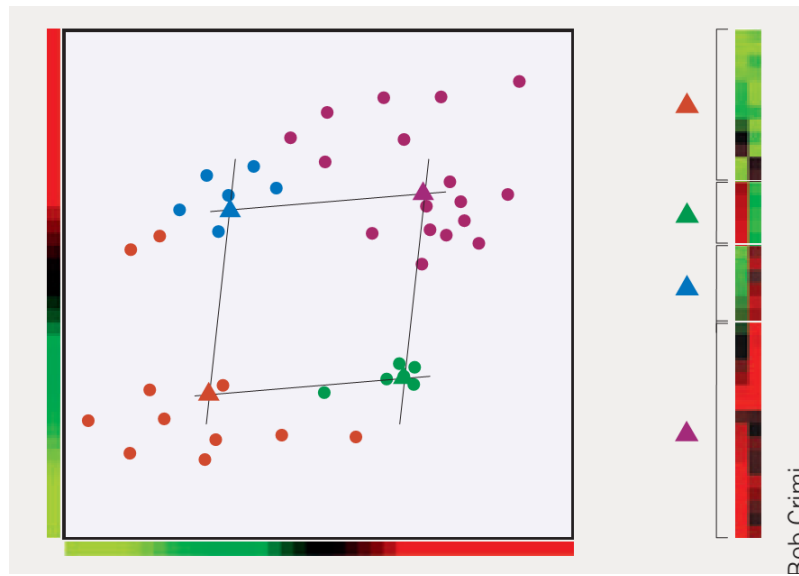
Nature Biotech 23(12):1499-1501 (2005)

Clustering refresher: hierarchical



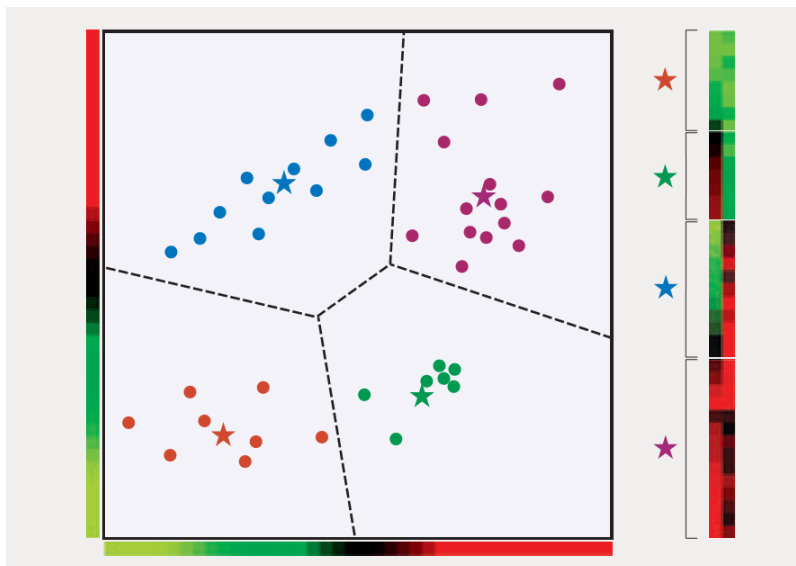
Nature Biotech 23(12):1499-1501 (2005)

Clustering refresher: SOM



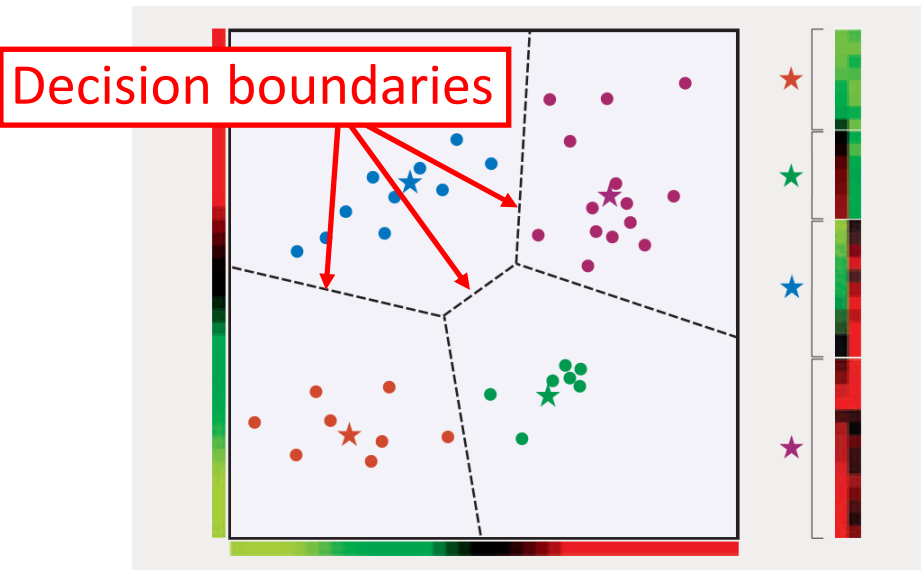
Nature Biotech 23(12):1499-1501 (2005)

Clustering refresher: *k*-means

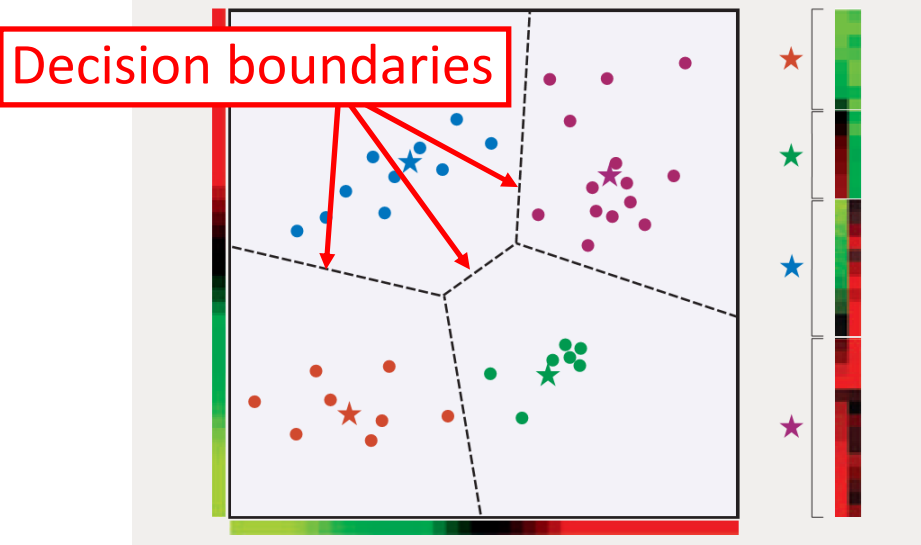


Nature Biotech 23(12):1499-1501 (2005)

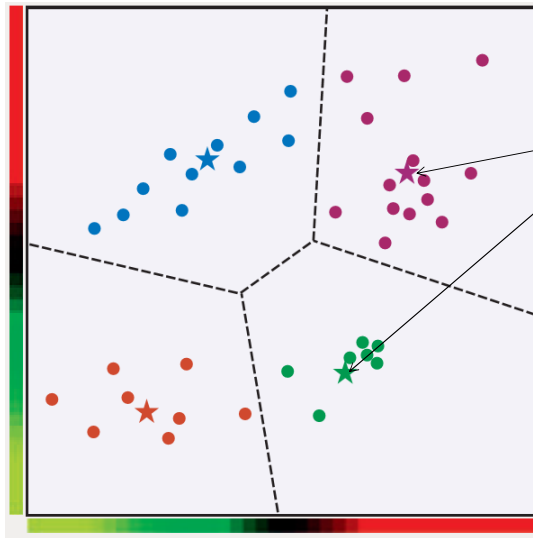
Clustering refresher: k -means



One of the simplest classifiers uses the same notion of decision boundaries.



One of the simplest classifiers uses this notion of decision boundaries.



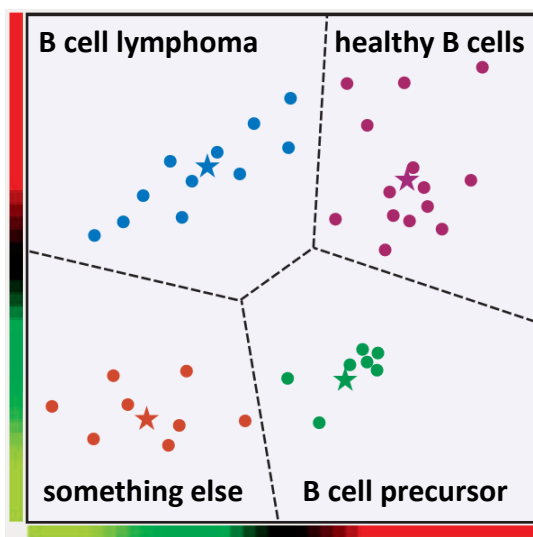
Rather than first clustering, calculate the centroid (mean) of objects with each label.

New observations are classified as belonging to the group whose mean is nearest.

= "minimum distance classifier"

Nature Biotech 23(12):1499-1501 (2005)

One of the simplest classifiers uses this notion of decision boundaries.



For example....

Nature Biotech 23(12):1499-1501 (2005)

**Molecular Classification of
Cancer: Class Discovery and
Class Prediction by Gene
Expression Monitoring**

T. R. Golub,^{1,2*} D. K. Slonim,^{1†} P. Tamayo,¹ C. Huard,¹
M. Gaasenbeek,¹ J. P. Mesirov,¹ H. Coller,¹ M. L. Loh,²
J. R. Downing,³ M. A. Caligiuri,⁴ C. D. Bloomfield,⁴
E. S. Lander^{1,5*}

Let's look at a specific
historic example:

“Enzyme-based histochemical analyses were introduced in the 1960s to demonstrate that **some leukemias were periodic acid-Schiff positive, whereas others were myeloperoxidase positive...**

This provided the first basis for classification of acute leukemias into those arising from lymphoid precursors (acute lymphoblastic leukemia, ALL), or from myeloid precursors (acute myeloid leukemia, AML).”

15 OCTOBER 1999 VOL 286 SCIENCE

**Molecular Classification of
Cancer: Class Discovery and
Class Prediction by Gene
Expression Monitoring**

T. R. Golub,^{1,2*} D. K. Slonim,^{1†} P. Tamayo,¹ C. Huard,¹
M. Gaasenbeek,¹ J. P. Mesirov,¹ H. Coller,¹ M. L. Loh,²
J. R. Downing,³ M. A. Caligiuri,⁴ C. D. Bloomfield,⁴
E. S. Lander^{1,5*}

Let's look at a specific
historic example:

“**Distinguishing ALL from AML is critical for successful treatment...**

chemotherapy regimens for ALL generally contain corticosteroids, vincristine, methotrexate, and L-asparaginase, whereas

most AML regimens rely on a backbone of daunorubicin and cytarabine (8).

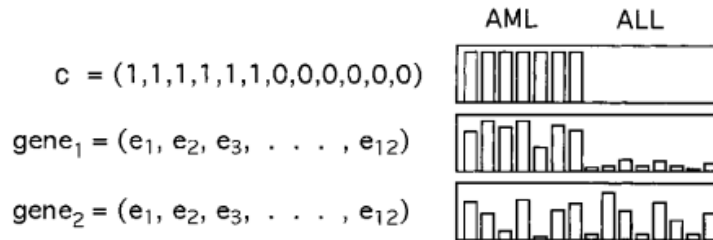
Although remissions can be achieved using ALL therapy for AML (and vice versa), **cure rates are markedly diminished**, and unwarranted toxicities are encountered.”

15 OCTOBER 1999 VOL 286 SCIENCE

Molecular Classification of Cancer: Class Discovery and Class Prediction by Gene Expression Monitoring

T. R. Golub,^{1,2*} D. K. Slonim,^{1†} P. Tamayo,¹ C. Huard,¹
 M. Gaasenbeek,¹ J. P. Mesirov,¹ H. Coller,¹ M. L. Loh,²
 J. R. Downing,³ M. A. Caligiuri,⁴ C. D. Bloomfield,⁴
 E. S. Lander^{1,5*}

Let's look at a specific historic example:



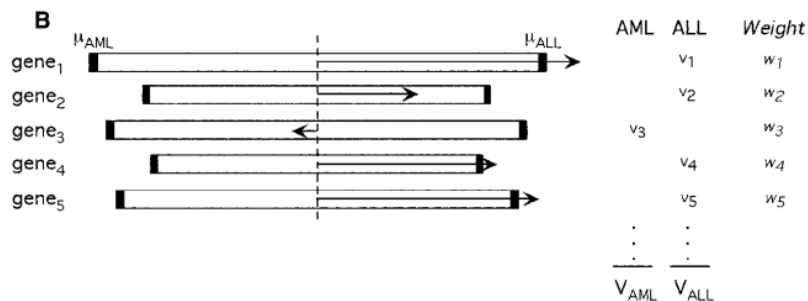
Take labeled samples, find genes whose abundances separate the samples...

15 OCTOBER 1999 VOL 286 SCIENCE

Molecular Classification of Cancer: Class Discovery and Class Prediction by Gene Expression Monitoring

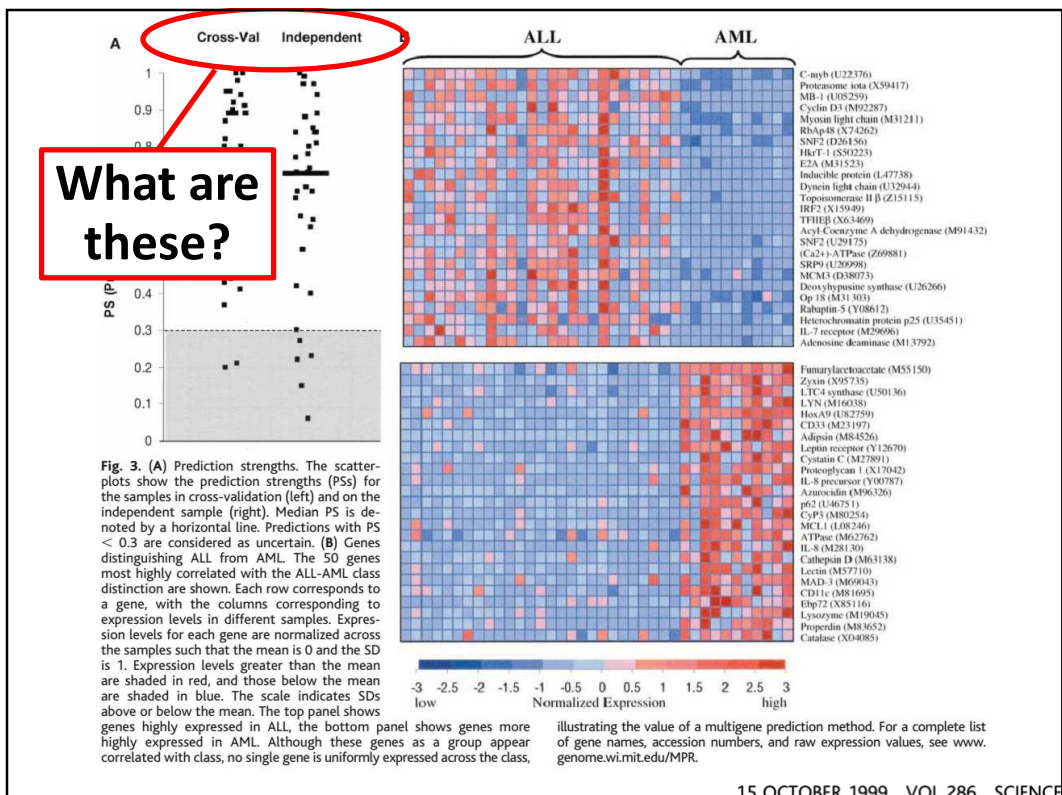
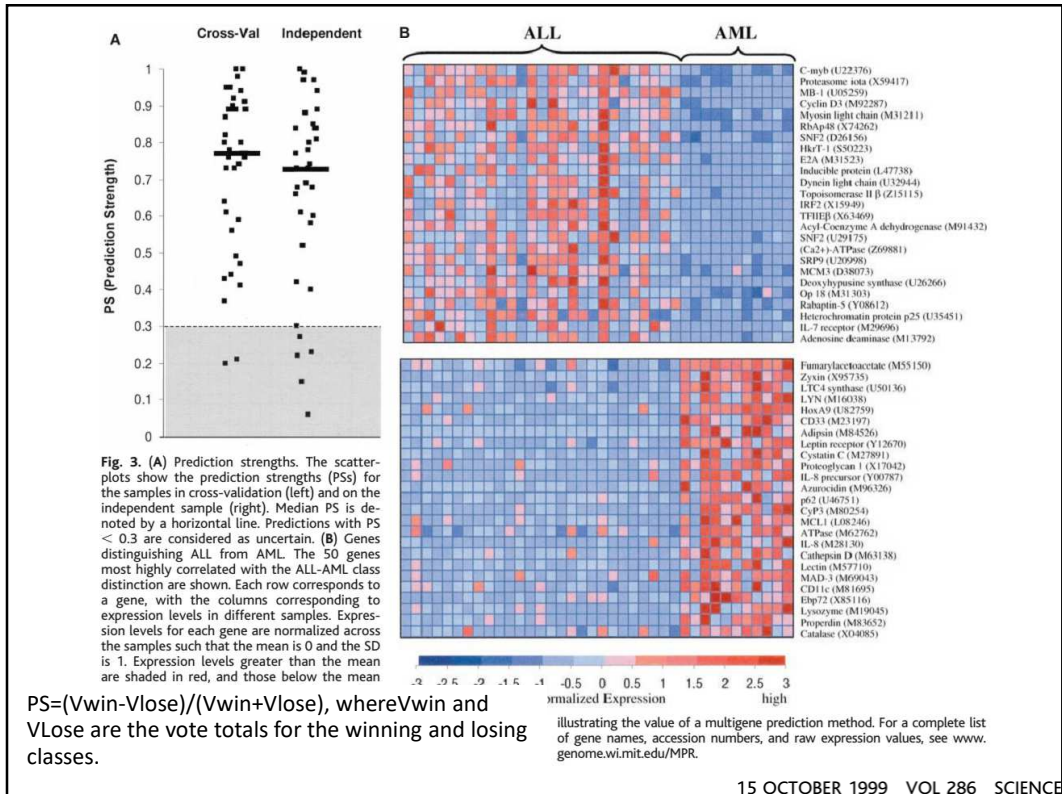
T. R. Golub,^{1,2*} D. K. Slonim,^{1†} P. Tamayo,¹ C. Huard,¹
 M. Gaasenbeek,¹ J. P. Mesirov,¹ H. Coller,¹ M. L. Loh,²
 J. R. Downing,³ M. A. Caligiuri,⁴ C. D. Bloomfield,⁴
 E. S. Lander^{1,5*}

Let's look at a specific historic example:



Calculate weighted average of indicator genes to assign class of an unknown

15 OCTOBER 1999 VOL 286 SCIENCE



Cross-validation

Withhold a sample, build a predictor based only on the remaining samples, and predict the class of the withheld sample.

Repeat this process for each sample, then calculate the cumulative or average error rate.

15 OCTOBER 1999 VOL 286 SCIENCE

X-fold cross-validation e.g. here, 5-fold:

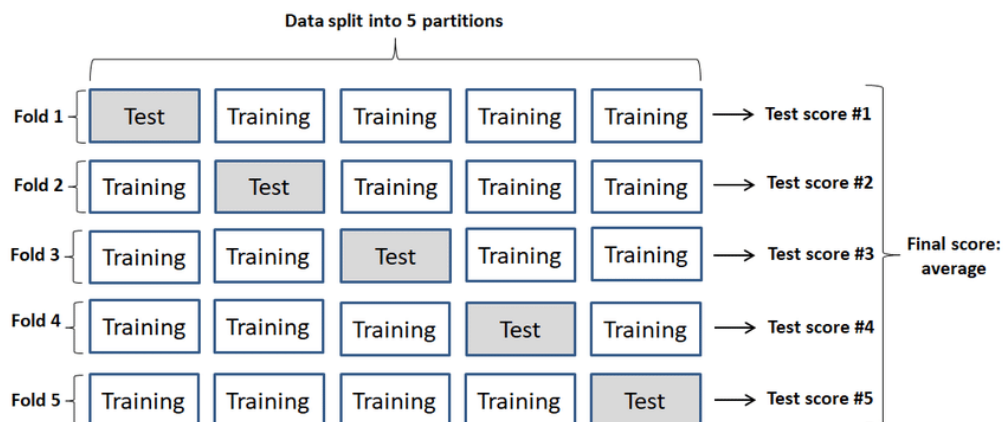


Image CC-BY, from DOI:10.3390/app9214500

Independent data

Withhold an entire dataset, build a predictor based only on the remaining samples (**the training data**).

Test the trained classifier on the independent **test data** to give a fully independent measure of performance.

15 OCTOBER 1999 VOL 286 SCIENCE

You already know how to measure how well these algorithms work (way back in our discussion of gene finding!)

True answer:

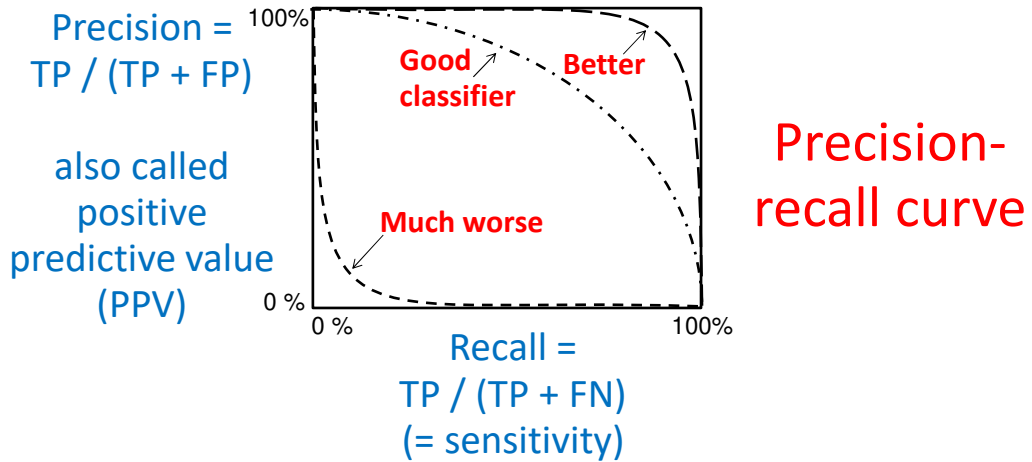
		Positive	Negative
Algorithm predicts:	Positive	True positive	False positive
	Negative	False negative	True negative

$$\text{Specificity} = \text{TP} / (\text{TP} + \text{FP})$$

$$\text{Sensitivity} = \text{TP} / (\text{TP} + \text{FN})$$

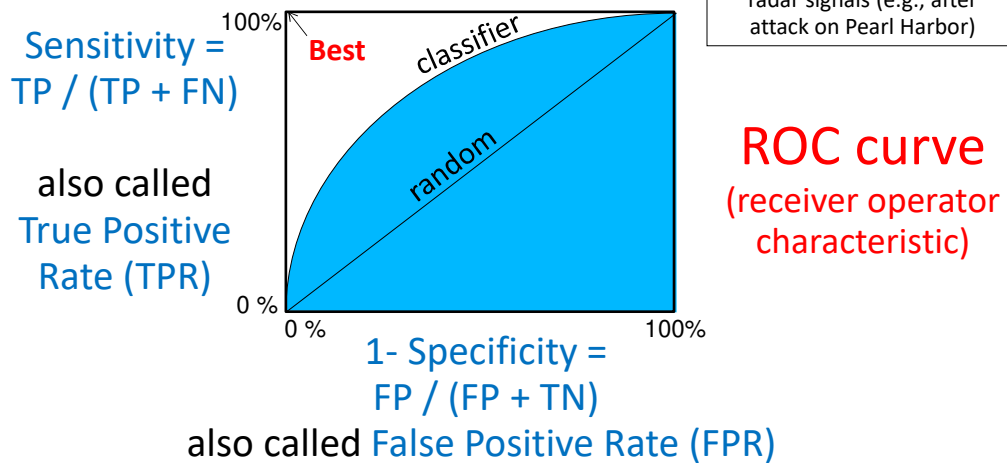
You already know how to measure how well these algorithms work (way back in our discussion of gene finding!)

Sort the data by their classifier score, then step from best to worst and plot the performance:

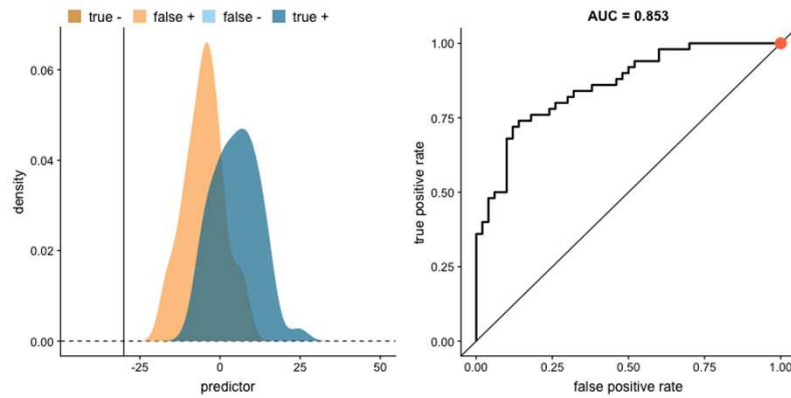


Another good option:

Sort the data by their classifier score, then step from best to worst and plot the performance:

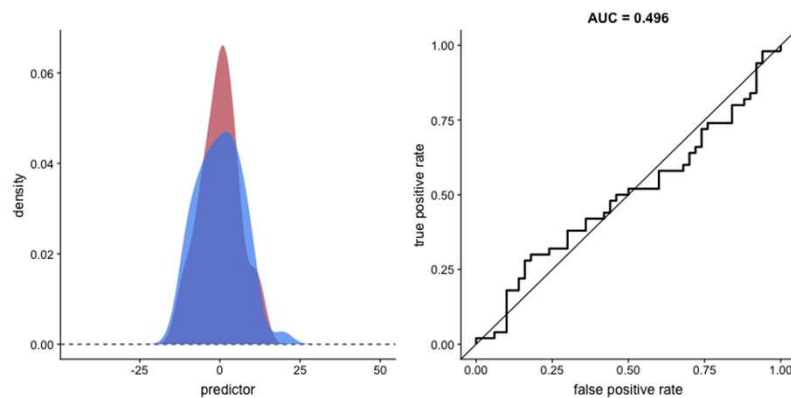


ROC curve, as you go from stronger to weaker predictions



Thanks to Dariya Sydykova (UT Austin), for her excellent visualizations, available here: https://github.com/dariyasdykova/open_projects/tree/master/ROC_animation

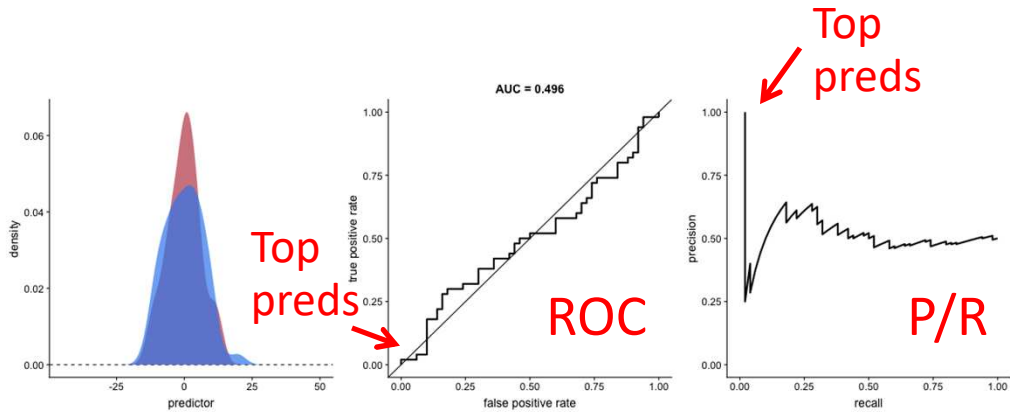
ROC curve, as you go from stronger to weaker classifiers



Thanks to Dariya Sydykova (UT Austin), for her excellent visualizations, available here: https://github.com/dariyasdykova/open_projects/tree/master/ROC_animation

ROC versus Recall/Precision

The 2 measures are related and both useful. They differ strongly in performance as proportions of positive and negative classes change.

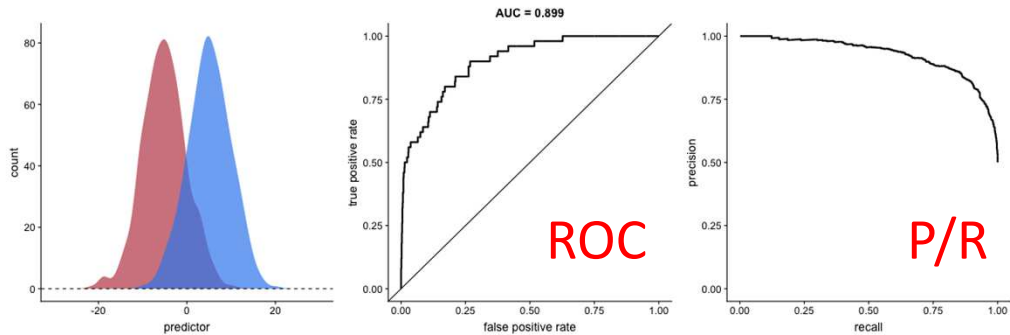


Thanks to Dariya Sydykova (UT Austin), for her excellent visualizations, available here: https://github.com/dariyasdykova/open_projects/tree/master/ROC_animation

ROC versus Recall/Precision

- R/P depends strongly on relative rates of the 2 classes
- ROC performance is independent of their relative rates

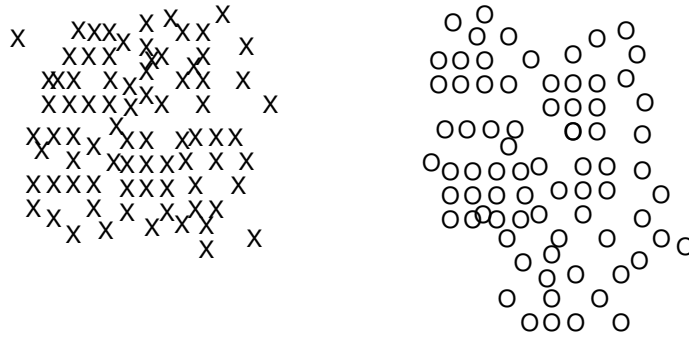
(It may be important or not for your particular problem...)



Thanks to Dariya Sydykova (UT Austin), for her excellent visualizations, available here: https://github.com/dariyasdykova/open_projects/tree/master/ROC_animation

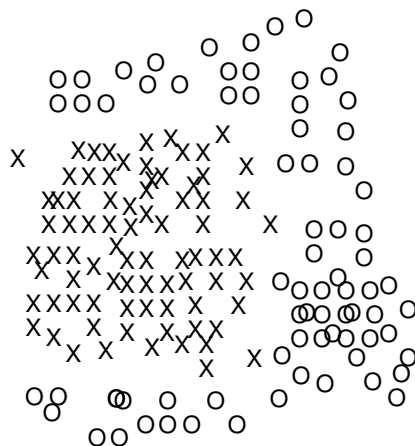
Back to our minimum distance classifier...

Would it work well for this data?



Back to our minimum distance classifier...

How about this data? What might?



Back to our minimum distance classifier...

How about this data? What might?

```
XXXXO O O O XXXXO O O O
XXXXO O O O XXXXO O O O
XXXXO O O O XXXXO O O O
XXXXO O O O XXXXO O O O
O O O O X X X X O O O O X X X X
O O O O X X X X O O O O X X X X
O O O O X X X X O O O O X X X X
O O O O X X X X O O O O X X X X
X X X X O O O O X X X X O O O O
X X X X O O O O X X X X O O O O
X X X X O O O O X X X X O O O O
X X X X O O O O X X X X O O O O
O O O O X X X X O O O O X X X X
O O O O X X X X O O O O X X X X
O O O O X X X X O O O O X X X X
O O O O X X X X O O O O X X X X
```

This is a great case for something called
a ***k-nearest neighbors classifier***:

**For each new object, calculate the k closest data points.
Let them vote on the label of the new object.**

```
XXXXO O O O XXXXO O O O
XXXXO O O O XXXXO O O O
XXXXO O O O XXXXO O O O
XXXXO O O O XXXXO O O O
O O O O X X X X O O O O X X X X
O O O O X X X X O O O O X X X X
O O O O X X X X O O O O X X X X
O O O O X X X X O O O O X X X X
X X X X O O O O X X X X O O O O
X X X X O O O O X X X X O O O O
X X X X O O O O X X X X O O O O
X X X X O O O O X X X X O O O O
O O O O X X X X O O O O X X X X
O O O O X X X X O O O O X X X X
O O O O X X X X O O O O X X X X
O O O O X X X X O O O O X X X X
```

This is surrounded by O's
and will probably be voted
to be an O.

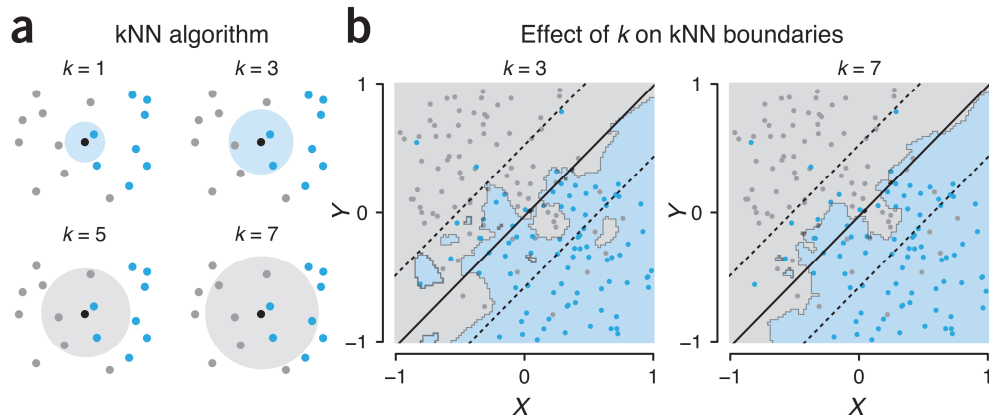
This one is surrounded by
X's and will probably be
voted to be an X.

This is a great case for something called

a ***k*-nearest neighbors classifier:**

For each new object, calculate the *k* closest data points.

Let them vote on the label of the new object.



kNN can (and often will) have complex, non-linear decision boundaries

NATURE METHODS | VOL.15 NO.1 | JANUARY 2018

Back to leukemias.

There was a follow-up study in 2010:

Clinical Utility of Microarray-Based Gene Expression Profiling in the Diagnosis and Subclassification of Leukemia: Report From the International Microarray Innovations in Leukemia Study Group

Torsten Haferlach, Alexander Kohlmann, Lothar Wiczorek, Giuseppe Basso, Geertruy Te Kronnie, Marie-Christine Béné, John De Vos, Jesus M. Hernández, Wolf-Karsten Hofmann, Ken I. Mills, Amanda Gilkes, Sabina Chiaretti, Sheila A. Shurtleff, Thomas J. Kipps, Laura Z. Rassenti, Allen E. Yeoh, Peter R. Papenhausen, Wei-min Liu, P. Mickey Williams, and Robin Foà

- **Tested clinical use of expression profiling to subtype leukemias**
- **Meta-analysis of 11 labs, 3 continents, 3,334 patients**
- **Stage 1 (2,096 patients):**
92.2% classification accuracy for 18 leukemia classes (99.7% median specificity)
- **Stage 2 (1,152 patients):**
95.6% median sensitivity and 99.8% median specificity for 14 subtypes of acute leukemia
- **Microarrays outperformed routine diagnostics in 29 (57%) of 51 discrepant cases**

Conclusion: "Gene expression profiling is a robust technology for the diagnosis of hematologic malignancies with high accuracy"

J Clin Oncol 28:2529-2537. © 2010

Current commercial breast cancer gene expression panels use this same strategy

Summary of breast cancer commercially available gene expression signatures.

Gene Signature	Biomarker Sources	Analysis Type	Clinical Outcome	No. Genes	Reference
Oncotype DX Breast	Breast tumor tissue	mRNA	Survival, benefit of chemotherapy	21	2004 Paik [82]
MammaPrint	Breast tumor tissue	mRNA	Survival	70	2002 van't Veer [83]
Endopredict	Breast tumor tissue	mRNA	Survival	12	2017 Warf [84]
Prosigna/PAM50	Breast tumor tissue	mRNA	Survival	50	2009 Parker [85]
Breast Cancer Index	Breast tumor tissue	mRNA	Survival, benefit of hormone therapy after 5 years	7	2008 Ma, 2013 Sgroi [86,87]

Prognostic Cancer Gene Expression Signatures: Current Status and Challenges (2021) *Cells* 10(3): 648

In practice, if you want to explore classifiers, I also strongly recommend always testing these classifiers:

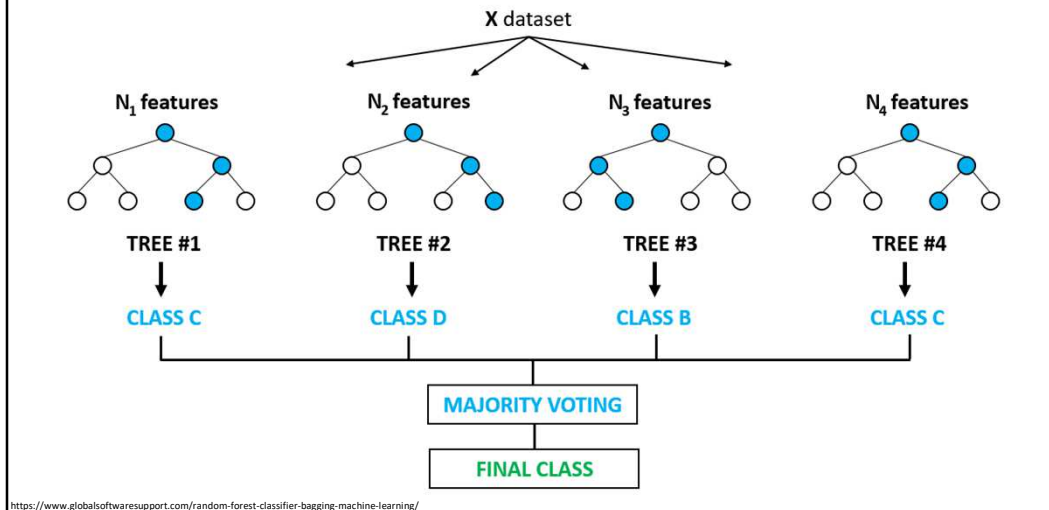
Random forests Support vector machines (SVM)

These two are surprisingly often the best for many biological classification problems. Weka can do both of them.

→ Note that I didn't say neural networks. Deep neural networks can be extremely powerful (e.g. AlphaFold) but are significantly more expert level and require extensive training examples. In general, you'll often be better off starting off with the above classifiers for many problems, only moving to deep neural networks if you really need to and only when you have data to support it.

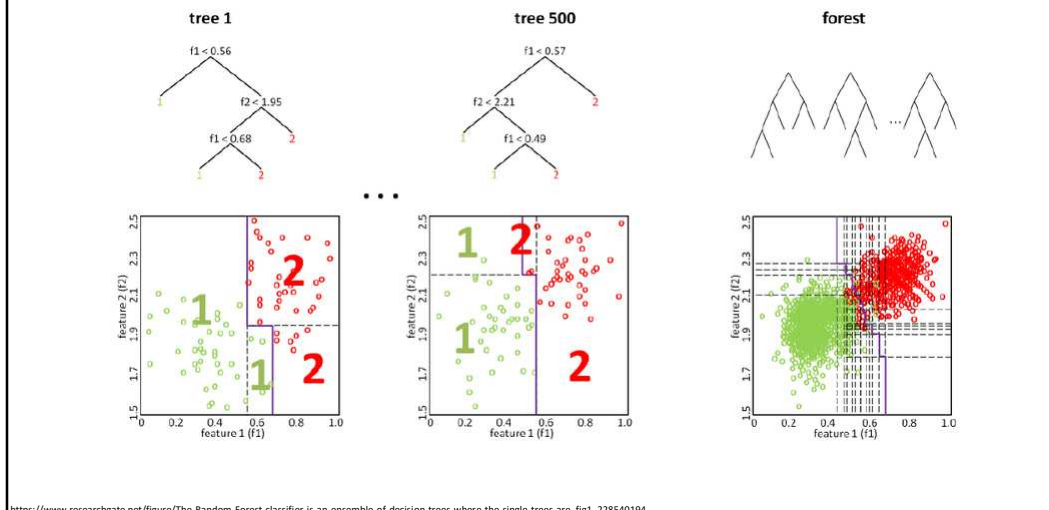
The two-slide overview of **Random forest classifiers:**

- (1) Construct many decision trees from random subsets of your features. Because the features vary across trees, trees tend to be weak but uncorrelated
- (2) All the trees “vote” on the answer, majority wins.



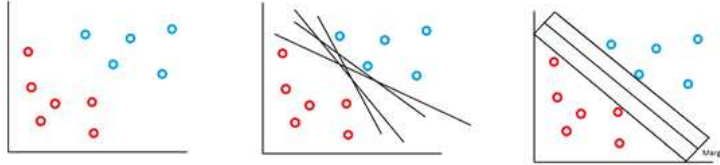
The two-slide overview of **Random forest classifiers:**

- (1) Construct many decision trees from random subsets of your features. Because the features vary across trees, trees tend to be weak but uncorrelated
- (2) All the trees “vote” on the answer, majority wins.

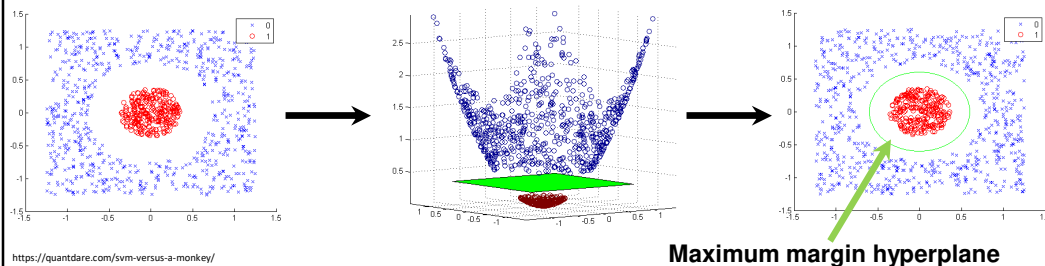


The one-slide overview of **Support vector machines:**

(1) Goal: make a linear classifier, choosing a decision boundary that *maximizes the distance margin* between classes



(2) But what if the boundary is non-linear? Use **kernels** to implicitly map the data to higher dimension where a linear decision can be made



<https://quantdare.com/svm-versus-a-monkey/>

Maximum margin hyperplane

In practice, if you want to explore classifiers, I strongly recommend the Weka package:

<http://www.cs.waikato.ac.nz/ml/weka/>



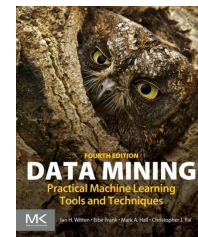
It's free, and easy to install, use, & troubleshoot. It lets you quickly test many alternative (well-vetted) classifiers, all in a proper cross-validated/precision-recall framework.

Here's a nice step-by-step intro for biologists :

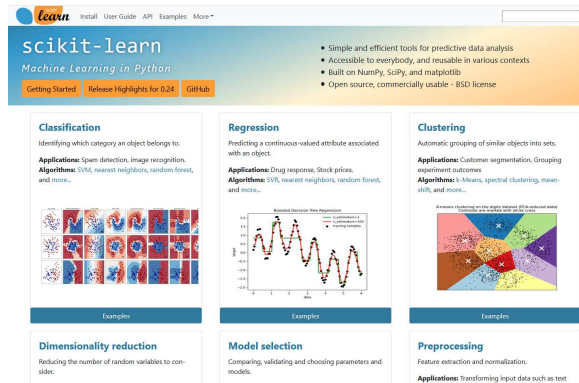
Introducing Machine Learning Concepts with WEKA, in *Statistical Genomics, Methods in Molecular Biology*, v. 1418, p. 353-378, 24 March 2016

http://link.springer.com/content/pdf/10.1007%2F978-1-4939-3578-9_17.pdf

There's also a great book to walk you through the entire process. Highly recommended!!!



In Python, you can also use the scikit-learn library:
<https://scikit-learn.org/stable/>
Like Weka, it's free, easy to install & use, and very powerful



I recommend combining it with the Pandas library for data analysis to make it easy to work with big, tabular datasets:
<https://pandas.pydata.org/>

