# Classifiers!!!

**BCH394P/364C Systems Biology / Bioinformatics**

**Edward Marcotte, Univ of Texas at Austin**

---

**Clustering** = task of <u>grouping</u> 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 <u>categorizing</u> a new observation, on the basis of a training set of data with observations (or instances) whose categories are known

## Remember, for clustering, we had a matrix of data…

**M samples**

| | | | | | |
|---|---|---|---|---|---|
| Gene 1, sample 1 | … | Gene 1, sample *j* | … | Gene 1, sample *M* |
| Gene 2, sample 1 | … | Gene 2, sample *j* | … | Gene 2, sample *M* |
| Gene 3, sample 1 | … | Gene 3, sample *j* | … | Gene 3, sample *M* |
| . | | . | | . |
| . | | . | | . |
| . | | . | | . |
| Gene *i*, sample 1 | … | Gene *i*, sample *j* | … | Gene *i*, sample *M* |
| . | | . | | . |
| . | | . | | . |
| . | | . | | . |
| Gene *N*, sample 1 | … | Gene *N*, sample *j* | … | Gene *N*, sample *M* |

**N genes**

For yeast, N ~ 6,000
For human, N ~ 22,000

*i.e.,* a matrix of *N* x *M* numbers

---

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

**M samples** ~~genomes~~ **genomes**

**N genes**

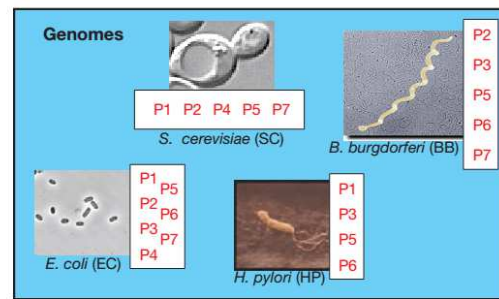| | | | | |
|---|---|---|---|---|
| Gene 1, sample 1 | … | Gene 1, sample *j* | … | Gene 1, sample *M* |
| Gene 2, sa | | | | , sample *M* |
| Gene 3, sa | | | | , sample *M* |
| . | | | | |
| . | | | | |
| . | | | | |
| Gene *i*, sa | | | | sample *M* |
| . | | | | |
| . | | | | |
| . | | | | |
| Gene *N*, sample 1 | … | Gene *N*, sample *j* | … | Gene *N*, sample *M* |

*Gene expression profiles*:
each entry indicates an mRNA's abundance in a different condition

*Phylogenetic profiles*:
each entry indicates whether the gene has homologs in a different organism

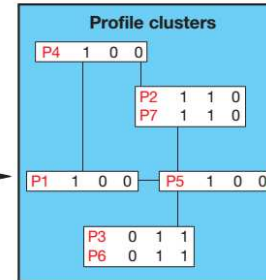For yeast, N ~ 6,000
For human, N ~ 22,000

**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)**
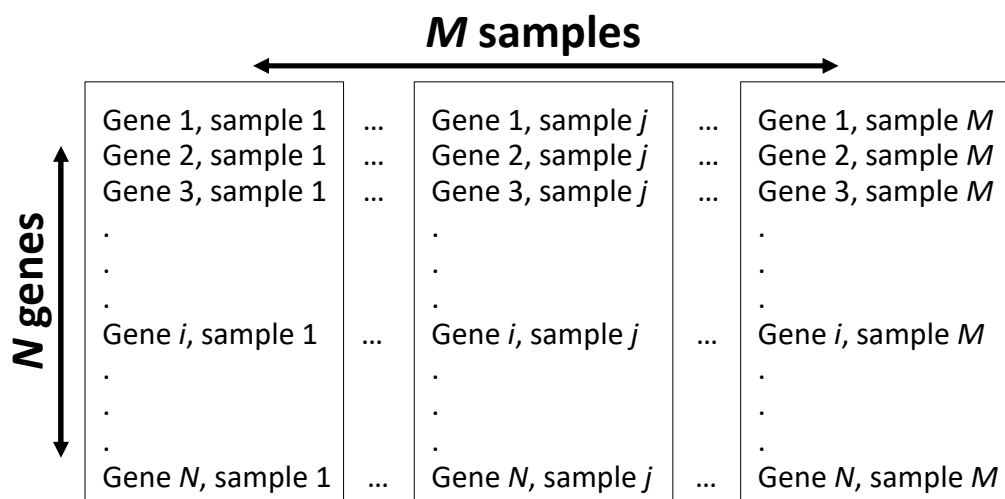


The method of phylogenetic profiles

NATURE | VOL 405 | 15 JUNE 2000

---

# Many such features are possible…

**M samples**

**N genes**

| Gene 1, sample 1 | … | Gene 1, sample *j* | … | Gene 1, sample *M* |
| Gene 2, sample 1 | … | Gene 2, sample *j* | … | Gene 2, sample *M* |
| Gene 3, sample 1 | … | Gene 3, sample *j* | … | Gene 3, sample *M* |
| . | | . | | . |
| . | | . | | . |
| . | | . | | . |
| Gene *i*, sample 1 | … | Gene *i*, sample *j* | … | Gene *i*, sample *M* |
| . | | . | | . |
| . | | . | | . |
| . | | . | | . |
| Gene *N*, sample 1 | … | Gene *N*, sample *j* | … | Gene *N*, sample *M* |

For yeast, N ~ 6,000
For human, N ~ 22,000

*i.e.,* a matrix of *N* x *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.**

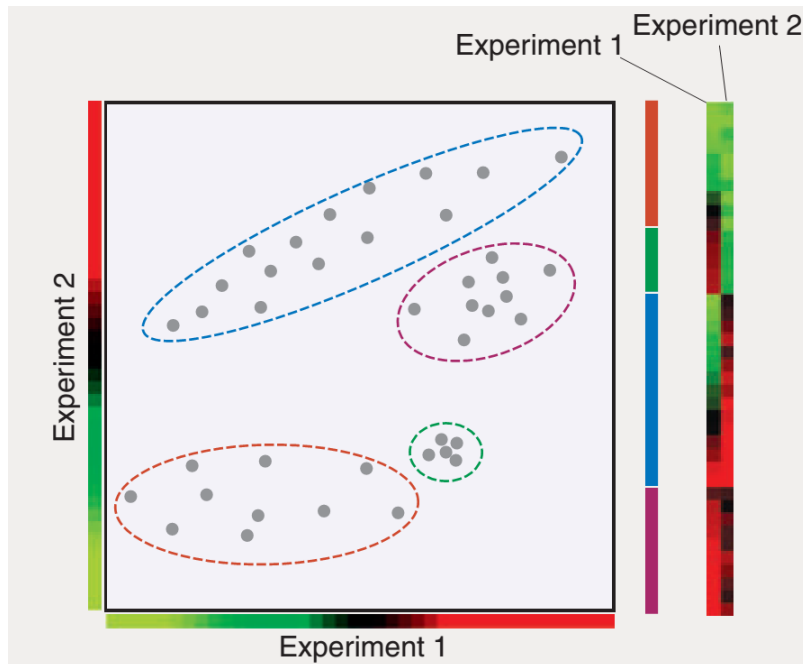| Names | Formula |
|---|---|
| Euclidean distance | $\|a - b\|_2 = \sqrt{\sum_i (a_i - b_i)^2}$ |
| Manhattan distance | $\|a - b\|_1 = \sum_i |a_i - b_i|$ |
| cosine similarity | $\dfrac{a \cdot b}{\|a\|\|b\|}$ |

---

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

**classifying**

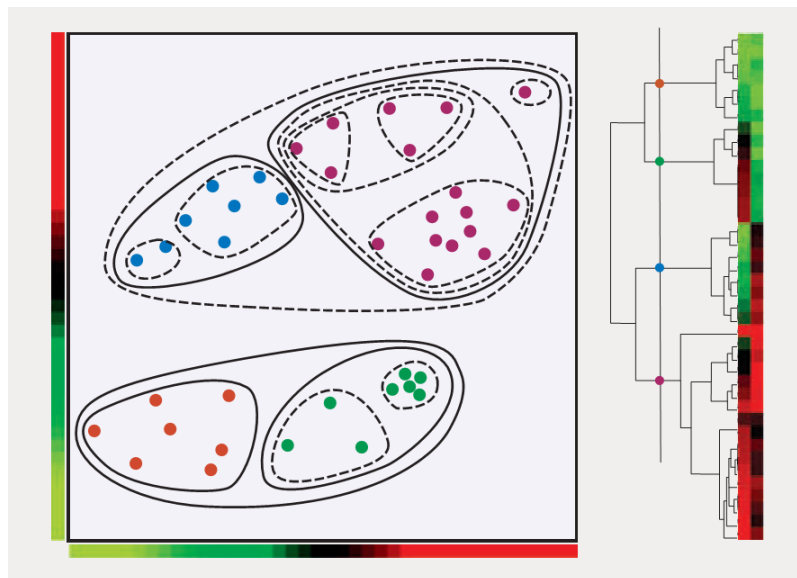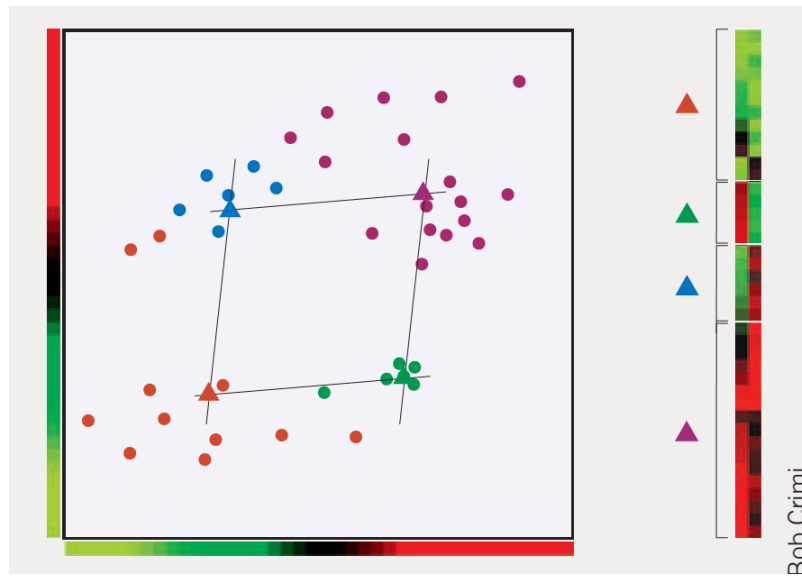| Names | Formula |
|---|---|
| Euclidean distance | $\|a - b\|_2 = \sqrt{\sum_i (a_i - b_i)^2}$ |
| Manhattan distance | $\|a - b\|_1 = \sum_i |a_i - b_i|$ |
| cosine similarity | $\dfrac{a \cdot b}{\|a\|\|b\|}$ |

# Clustering refresher: 2-D example

# Clustering refresher: hierarchical

# Clustering refresher: SOM



Bob Crimi

# Clustering refresher: *k*-means

# Clustering refresher: *k*-means



Decision boundaries

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



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.*

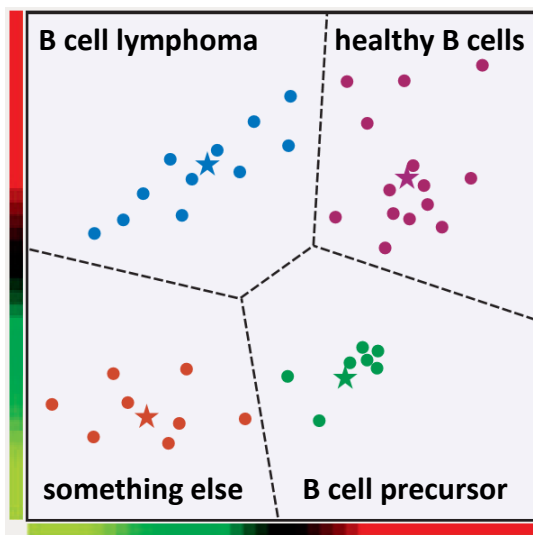**=“minimum distance classifier”**

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

**B cell lymphoma**  **healthy B cells**

**something else**  **B cell precursor**

**For example….**

**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 <u>lymphoid</u> precursors (acute lymphoblastic leukemia, ALL), or from <u>myeloid</u> precursors (acute myeloid leukemia, AML)."

---

**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

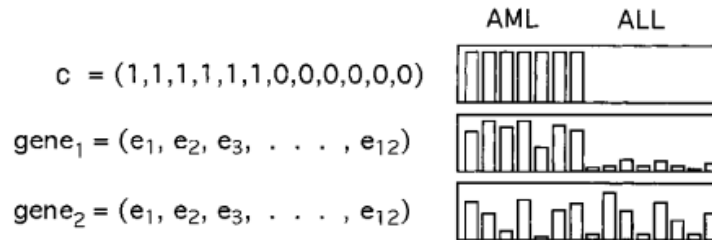most AML regimens rely on a backbone of daunorubicin and cytarabine (8).

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

**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:

$$c = (1,1,1,1,1,1,0,0,0,0,0,0)$$

$$gene_1 = (e_1, e_2, e_3, \ldots, e_{12})$$

$$gene_2 = (e_1, e_2, e_3, \ldots, e_{12})$$

AML      ALL

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

---

B

$\mu_{AML}$         $\mu_{ALL}$    AML   ALL   Weight

gene_1         $v_1$    $w_1$

gene_2         $v_2$    $w_2$

gene_3      $v_3$        $w_3$

gene_4         $v_4$    $w_4$

gene_5         $v_5$    $w_5$

$V_{AML}$   $V_{ALL}$

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

Fig. 3. (A) Prediction strengths. The scatter-plots show the prediction strengths (PSs) for the samples in cross-validation (left) and on the independent sample (right). Median PS is denoted by a horizontal line. Predictions with PS < 0.3 are considered as uncertain. (B) Genes distinguishing ALL from AML. The 50 genes most highly correlated with the ALL-AML class distinction are shown. Each row corresponds to a gene, with the columns corresponding to expression levels in different samples. Expression levels for each gene are normalized across the samples such that the mean is 0 and the SD is 1. Expression levels greater than the mean are shaded in red, and those below the mean

PS=(Vwin-Vlose)/(Vwin+Vlose), whereVwin and VLose are the vote totals for the winning and losing classes.

illustrating the value of a multigene prediction method. For a complete list of gene names, accession numbers, and raw expression values, see www.genome.wi.mit.edu/MPR.

15 OCTOBER 1999   VOL 286   SCIENCE

---



What are these?

Fig. 3. (A) Prediction strengths. The scatter-plots show the prediction strengths (PSs) for the samples in cross-validation (left) and on the independent sample (right). Median PS is denoted by a horizontal line. Predictions with PS < 0.3 are considered as uncertain. (B) Genes distinguishing ALL from AML. The 50 genes most highly correlated with the ALL-AML class distinction are shown. Each row corresponds to a gene, with the columns corresponding to expression levels in different samples. Expression levels for each gene are normalized across the samples such that the mean is 0 and the SD is 1. Expression levels greater than the mean are shaded in red, and those below the mean are shaded in blue. The scale indicates SDs above or below the mean. The top panel shows genes highly expressed in ALL, the bottom panel shows genes more highly expressed in AML. Although these genes as a group appear correlated with class, no single gene is uniformly expressed across the class,

illustrating the value of a multigene prediction method. For a complete list of gene names, accession numbers, and raw expression values, see www.genome.wi.mit.edu/MPR.

15 OCTOBER 1999   VOL 286   SCIENCE

11

## 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.
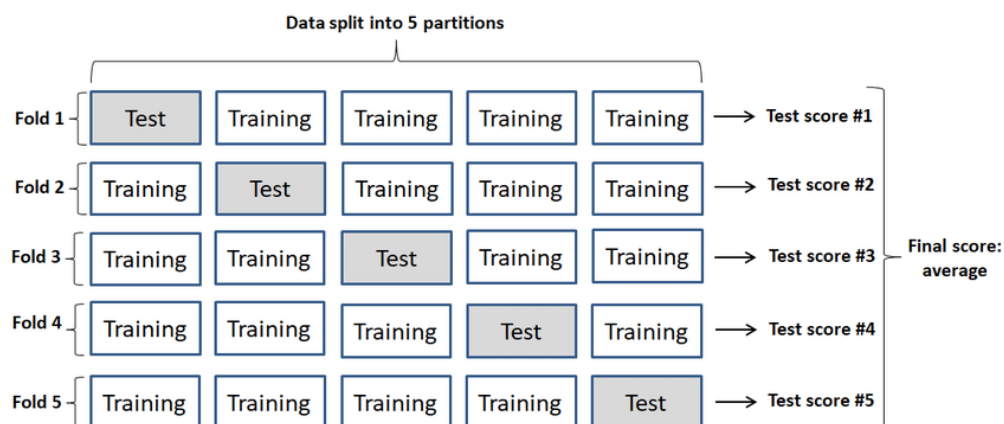
---

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

**Data split into 5 partitions**

| | | | | | |
|---|---|---|---|---|---|
| Fold 1 | Test | Training | Training | Training | Training | → Test score #1 |
| Fold 2 | Training | Test | Training | Training | Training | → Test score #2 |
| Fold 3 | Training | Training | Test | Training | Training | → Test score #3 |
| Fold 4 | Training | Training | Training | Test | Training | → Test score #4 |
| Fold 5 | Training | Training | Training | Training | Test | → Test score #5 |

Final score: average

## Independent data

Withhold <u>an entire dataset</u>, build a predictor based only on the remaining samples
(the training data).

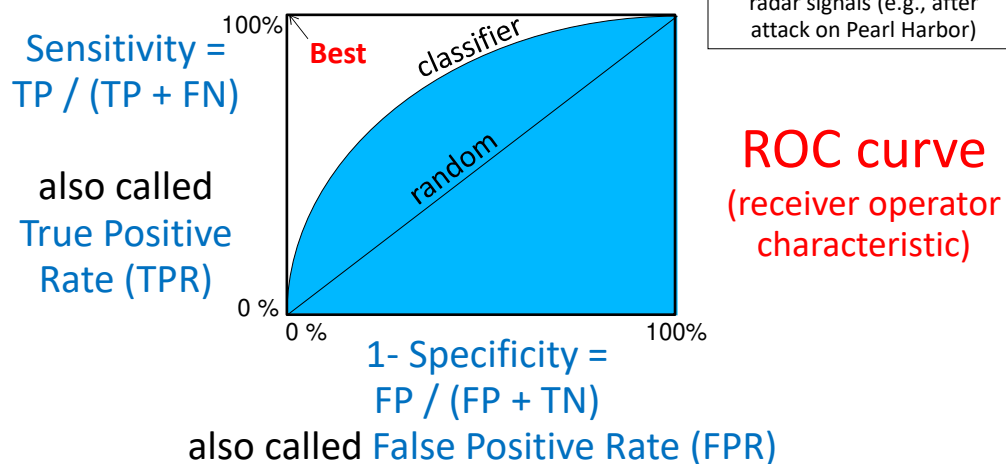Test the trained classifier on the independent test data to give <u>a fully independent measure of performance</u>.

---

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

**True answer:**

|  | Positive | Negative |
|---|---|---|
| **Positive** | True positive | False positive |
| **Negative** | False negative | True negative |

**Algorithm predicts:**

Specificity = TP / (TP + FP)

Sensitivity = TP / (TP + 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:

Precision =
TP / (TP + FP)

also called
positive
predictive value
(PPV)



Precision-recall curve

Recall =
TP / (TP + FN)
(= sensitivity)

---

Another good option:

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

First used in WWII to analyze radar signals (e.g., after attack on Pearl Harbor)

Sensitivity =
TP / (TP + FN)

also called
True Positive
Rate (TPR)



ROC curve
(receiver operator characteristic)

1- Specificity =
FP / (FP + TN)
also called False Positive Rate (FPR)

# ROC curve, as you go from stronger to weaker <u>predictions</u>

# ROC curve, as you go from stronger to weaker <u>classifiers</u>

# 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.
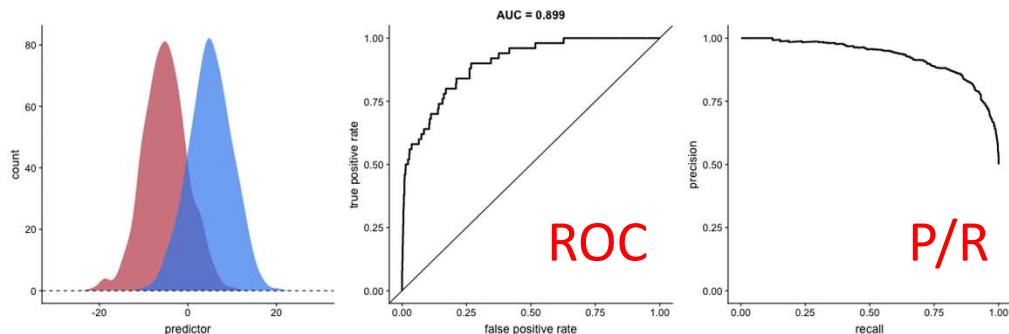


Thanks to Dariya Sydykova (UT Austin), for her excellent visualizations, available here:
https://github.com/dariyasydykova/open_projects/tree/master/ROC_animation

# ROC versus Recall/Precision

- R/P depends <u>strongly</u> on relative rates of the 2 classes
- ROC performance is <u>independent</u> 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/dariyasydykova/open_projects/tree/master/ROC_animation

# For example, the FDA evaluates diagnostic tests (≈classifiers) using ROC-like measures

**BinaxNOW™ COVID-19 Ag Card Performance within 7 days of symptom onset against the Comparator Method**

| BinaxNOW™ COVID-19 Ag Card | Comparator Method | | |
|---|---|---|---|
| | Positive | Negative | Total |
| Positive | 99 | 5 | 104 |
| Negative | 18 | 338 | 356 |
| Total | 117 | 343 | 460 |
| Positive Agreement: 99/117 | 84.6% (95% CI: 76.8% - 90.6%) | | |
| Negative Agreement: 338/343 | 98.5% (95% CI: 96.6% - 99.5%) | | |

← **Sensitivity (TPR)**

← **1- Specificity (FPR)**

**BinaxNOW™**
**COVID-19 Ag**
Abbott CARD

**Any guesses why?**

**Hint: How would COVID test performance change for ROC vs Precision/Recall as the infection rate in the population changes?**

https://www.fda.gov/media/141570/download

---

# 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?

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.**



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.**



**a**     kNN algorithm     **b**     Effect of *k* on kNN boundaries

$k = 1$     $k = 3$

$k = 5$     $k = 7$

$k = 3$     $k = 7$

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

19

Back to leukemias. There was a follow-up study in 2010:

- **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"**

# 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] |

In practice, if you want to explore classifiers, I also <u>strongly</u> recommend always testing these classifiers:
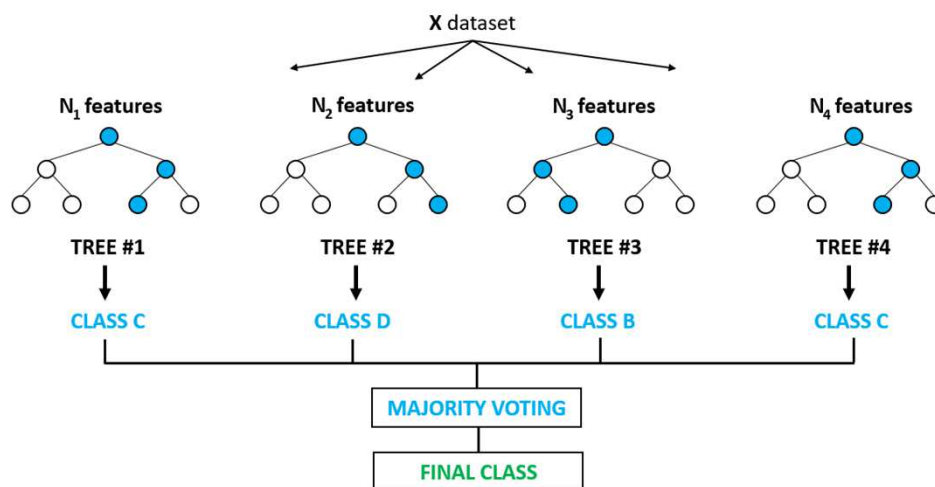
## Random forests
## Support vector machines (SVM)

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 usually 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 when you have data to support it.
→ We'll talk about NNs in the next 2 lectures, including large language models.

---
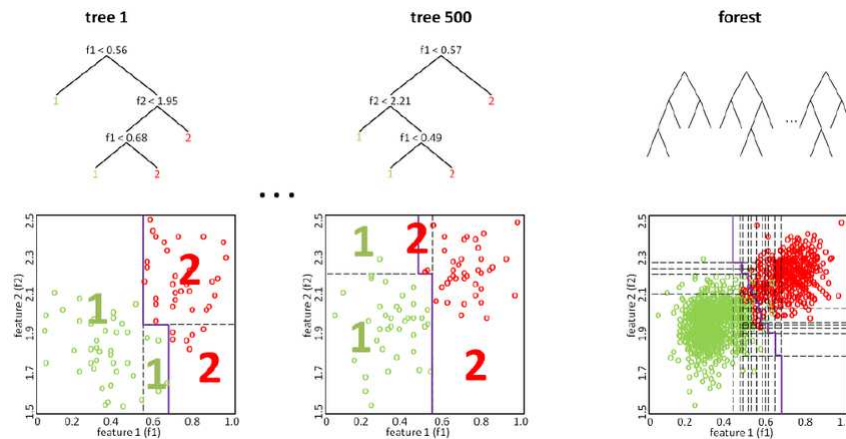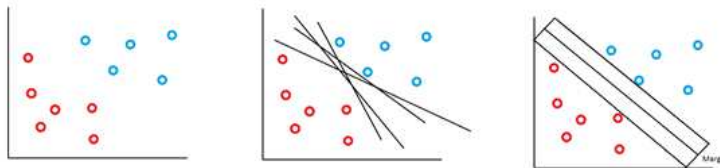
The two-slide overview of  **Random forest classifiers:**
(1) Construct many decision trees from random subsets of your features. Because the features vary across trees, trees tend to be weak but uncorrelated
(2) All the trees "vote" on the answer, majority wins.



https://www.globalsoftwaresupport.com/random-forest-classifier-bagging-machine-learning/

21

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
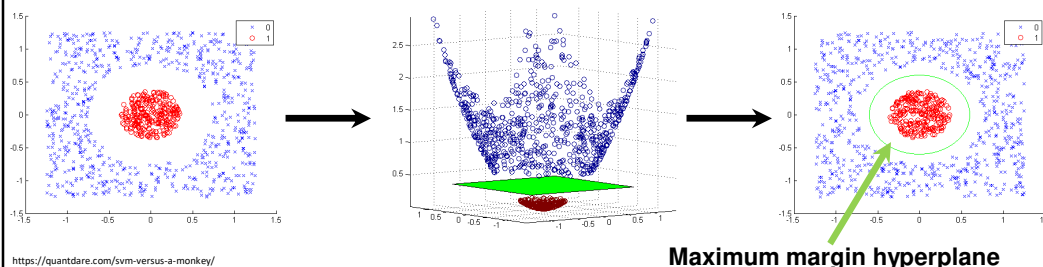(2) All the trees "vote" on the answer, majority wins.

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

**Maximum margin hyperplane**

In practice, if you want to explore classifiers, I <u>strongly</u> recommend the Weka package:
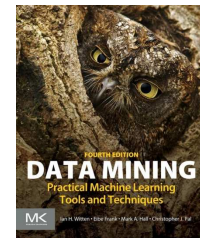
http://www.cs.waikato.ac.nz/ml/weka/

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.

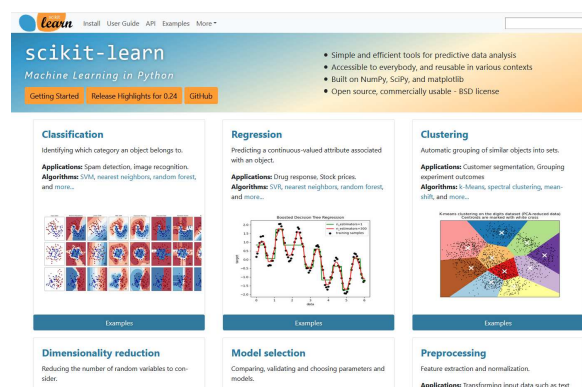Here's a nice step-by-step intro for biologists :
**Introducing Machine Learning Concepts with WEKA**, in *Statistical Genomics, Methods in Molecular Biology,* v. 1418, p. 353-378, 24 March 2016

**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/

Coming up:

The next two lectures will be guest lectures covering the basics of deep neural networks, with talks on

**Protein 3D structural modeling and prediction w/ AlphaFold/ChimeraX**

**DeepNNs, Large Language Models, & ESM**