

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

Adapted from Wikipedia

Remember, for clustering, we had a matrix of data									
<i>M</i> samples									
Ngenes	Gene 1, sample 1 Gene 2, sample 1 Gene 3, sample 1 Gene <i>i</i> , sample 1 Gene <i>N</i> , sample 1	···· ····	Gene 1, sample <i>j</i> Gene 2, sample <i>j</i> Gene 3, sample <i>j</i> Gene <i>i</i> , sample <i>j</i>	···· ····	Gene 1, sample <i>M</i> Gene 2, sample <i>M</i> Gene 3, sample <i>M</i> Gene <i>i</i> , sample <i>M</i> Gene <i>N</i> , sample <i>M</i>				
For yeast, N ~ 6,000 For human, N ~ 22,000			<i>i.e.,</i> a matrix of <i>N</i> x <i>M</i> numbers						







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





















Molecular Classification of Cancer: Class Discovery and Class Prediction by Gene Expression Monitoring T. R. Golub,¹²⁴; 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. Caliguri,⁴ C. D. Bloomfield,⁴ E. S. Lander^{1,54} "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

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:

15 OCTOBER 1999 VOL 286 SCIENCE

"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."

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.

X-fold cross-validation e.g. 3-fold or 10-fold

Can also withhold 1/X (e.g. 1/3 or 1/10) of sample, build a predictor based only on the remaining samples, and predict the class of the withheld samples.

Repeat this process X times for each withheld fraction of the sample, then calculate the cumulative or average error rate.

15 OCTOBER 1999 VOL 286 SCIENCE

15 OCTOBER 1999 VOL 286 SCIENC

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</u> 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 a	nswer:						
	Negative Positive	Positive	Negative						
Algorithm		True positive	False positive						
predicts:		False negative	True negative						
	Sp Se	ecificity = T nsitivity = T)						































