

Classification

Two Different Approaches

- **Generative**
 - Bayesian Classification and Naïve Bayes
 - Example: Mitochondrial Protein Prediction
- **Discriminative**
 - Support Vector Machines
 - Example: Tumor Classification

Bayesian Classification

We will pose the classification problem in **probabilistic terms**

Create **models for how features are distributed** for objects of different classes

We will use probability calculus to **make classification decisions**

Classifying Mitochondrial Proteins

Derive 7 features for all human proteins

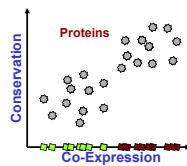
- Targeting signal
- Protein domains
- Co-expression
- Mass Spec
- Homology
- Induction
- Motifs

Predict nuclear encoded mitochondrial genes
Maestro

The figure shows a list of 7 features used for classification: Targeting signal, Protein domains, Co-expression, Mass Spec, Homology, Induction, and Motifs. To the right is a screenshot of a scientific paper titled "Systematic identification of human mitochondrial disease genes through integrative genomics" by Hwang et al. (2007). The paper abstract discusses the identification of mitochondrial disease genes using a combination of genomic, proteomic, and clinical data. A table at the bottom of the paper lists identified genes and their associated diseases.

Lets Look at Just One Feature

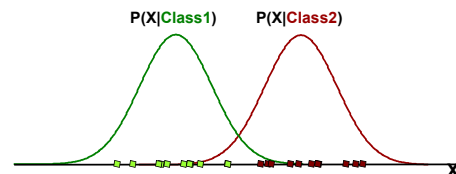
- Each object can be associated with multiple features
- We will look at the case of just one feature for now



We are going to define two key concepts....

The First Key Concept

Features for each class drawn from **class-conditional probability distributions (CCPD)**



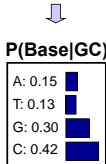
Our first goal will be to *model* these distributions

This Seems Familiar....

Remember how we modeled GC Islands when we talked about HMMs?

$P(\text{Base}|\text{GC Island})$ is a class-conditional probability distribution

We designed emissions probabilities:



Objects: Nucleotides
Class: GC Island
Feature: A,T,G,C

Makes sense: we built a simple classifier for sequences (GC vs Not)

The Second Key Concept

We model **prior probabilities** to quantify the expected a priori chance of seeing a class

$P(\text{Class2})$ & $P(\text{Class1})$

$P(\text{mito})$ = how likely is the next protein to be a mitochondrial protein before I see any features to help me decide

We expect ~1500 mitochondrial genes out of ~21000 total, so

$$P(\text{mito}) = 1500/21000$$

$$P(\sim\text{mito}) = 19500/21000$$

But How Do We Classify?

- So we have priors defining the *a priori* probability of a class

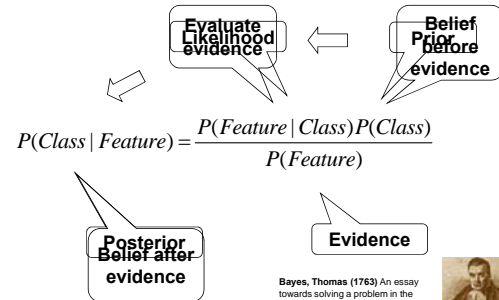
$$P(\text{Class1}), P(\text{Class2})$$

- We also have models for the probability of a feature given each class

$$P(X|\text{Class1}), P(X|\text{Class2})$$

But we want the probability of the class given a feature
How do we get $P(\text{Class1}|X)$?

Bayes Rule



Bayes Decision Rule

If we observe an object with feature X, how do decide if the object is from Class 1?

The **Bayes Decision Rule** is simply choose Class1 if:

$$P(\text{Class1} | X) > P(\text{Class2} | X)$$

$$\frac{P(X | \text{Class1})P(L1)}{P(X)} > \frac{P(X | \text{Class2})P(L2)}{P(X)}$$

This is the same number on both sides!
 $P(X | \text{Class1})P(\text{Class1}) > P(X | \text{Class2})P(\text{Class2})$

Discriminant Function

We can create a convenient representation of the Bayes Decision Rule

$$P(X | \text{Class1})P(\text{Class1}) > P(X | \text{Class2})P(\text{Class2})$$

$$\frac{P(X | \text{Class1})P(\text{Class1})}{P(X | \text{Class2})P(\text{Class2})} > 1$$

$$G(X) = \log \frac{P(X | \text{Class1}) P(\text{Class1})}{P(X | \text{Class2}) P(\text{Class2})} > 0$$

If $G(X) > 0$, we classify as Class 1

Stepping back

What do we have so far?

We have defined the two components, **class-conditional distributions and priors**

$$P(X|Class1), P(X|Class2) \quad P(Class1), P(Class2)$$

We have used **Bayes Rule** to create a **discriminant function for classification** from these components

$$G(X) = \log \frac{P(X|Class1)P(Class1)}{P(X|Class2)P(Class2)} > 0$$

Given a new feature, X, we plug it into this equation...
...and if $G(X) > 0$ we classify as **Class1**

Flashback to Previous Lecture

We defined a scoring function for *classifying* sequences as GC islands or background DNA

Can we justify this choice in Bayesian terms?

$$Score = \log \frac{P(S|GC)}{P(S|B)} = \log \frac{P(Sequence|GCIsland)}{P(Sequence|BackgroundDNA)}$$

GC Islands



Background DNA



From the Bayesian Approach

We can write out our Bayes discriminant function:

$$G(S) = \log \frac{P(GC|S)}{P(B|S)} = \log \frac{P(S|GC)P(GC)}{P(S|B)P(B)}$$

$$= \log \frac{P(S|GC)}{P(S|B)} \quad \leftarrow \text{If } P(MB) \neq P(B) \neq 1$$

Bayes Rule

Log-likelihood ratio is Bayes classifier assuming equal priors

Unequal Priors

Bayes rule lets us correct for the fact that GC islands are *not* as common as background DNA

$$G(X) = \log \frac{P(S|MB)P(MP)}{P(S|B)P(B)} = \log \frac{P(S|MB)}{P(S|B)} + \log \frac{0.15}{0.85}$$

$$= Score_{original} - 0.75$$

We require a higher score to call a GC island
(The HMM implicitly accounted for this in the transition probabilities)

GC Islands



15%

Background DNA



85%

Back to Classification

We have two fundamental tasks

- We need to estimate the needed probability distributions
 - $P(X|Mito)$ and $P(x|\neg Mito)$
 - $P(Mito)$ and $P(\neg Mito)$
- We need to assess the accuracy of the classifier
 - How well does it classify new objects

The All Important Training Set

Building a classifier requires a set of labeled data points called the Training Set

The quality of the classifier depends on the number of training set data points

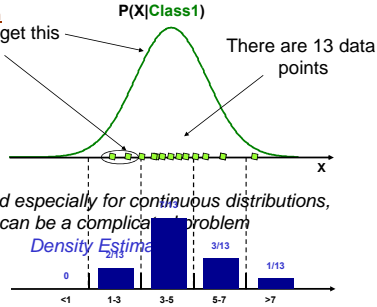


How many data points you need depends on the problem
Need to build and test your classifier

Getting P(X|Class) from Training Set

One Simple Approach

How do we get this from these?
 Divide X values into bins
 And then we simply count frequencies



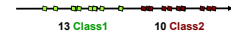
Getting Priors

Three general approaches

1. Estimate priors by counting fraction of classes in training set

$$P(\text{Class1}) = 13/23$$

$$P(\text{Class2}) = 10/23$$



2. Estimate from "expert" knowledge
 But sometimes fractions in training set are not representative of world

Example

$$P(\text{mito}) = 1500/21000$$

$$P(\text{-mito}) = 19500/21000$$

3. We have no idea – use equal (uninformative) priors

$$P(\text{Class1}) = P(\text{Class2})$$

We Are Just About There....

We have created the class-conditional distributions and priors

$$P(X|\text{Class1}), P(X|\text{Class2}) \quad P(\text{Class1}), P(\text{Class2})$$

And we are ready to plug these into our discriminant function

$$G(X) = \log \frac{P(X|\text{Class1}) P(\text{Class1})}{P(X|\text{Class2}) P(\text{Class2})} > 0$$

But there is one more little complication....

But What About Multiple Features?

- We have focused on a single feature for an object
- But mitochondrial protein prediction (for example) has 7 features

Targeting signal
Protein domains
Co-expression
Mass Spec
Homology
Induction
Motifs

So $P(X|\text{Class})$ become $P(X1, X2, X3, \dots, X8|\text{Class})$ and our discriminant function becomes

$$G(X) = \log \frac{P(X_1, X_2, \dots, X_7|\text{Class1}) P(\text{Class1})}{P(X_1, X_2, \dots, X_7|\text{Class2}) P(\text{Class2})} > 0$$

Distributions Over Many Features

Estimating $P(X1, X2, X3, \dots, X8|\text{Class1})$ can be difficult

- Assume each feature binned into 5 possible values
- We have 5^8 combinations of values we need to count the frequency for
- Generally will not have enough data
 - We will have lots of nasty zeros

Naïve Bayes Classifier

We are going to make the following assumption:

All features are independent given the class

$$P(X_1, X_2, \dots, X_n | \text{Class}) = P(X_1 | \text{Class}) P(X_2 | \text{Class}) \dots P(X_n | \text{Class})$$

$$= \prod_{i=1}^n P(X_i | \text{Class})$$

We can thus estimate individual distributions for each feature and just multiply them together!

Naïve Bayes Discriminant Function

Thus, with the Naïve Bayes assumption, we can now rewrite, this:

$$G(X_1, \dots, X_7) = \log \frac{P(X_1, X_2, \dots, X_7 | \text{Class1}) P(\text{Class1})}{P(X_1, X_2, \dots, X_7 | \text{Class2}) P(\text{Class2})} > 0$$

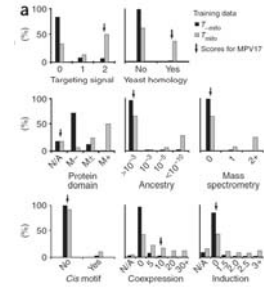
As this:

$$G(X_1, \dots, X_7) = \log \frac{\prod P(X_i | \text{Class1}) P(\text{Class1})}{\prod P(X_i | \text{Class2}) P(\text{Class2})} > 0$$

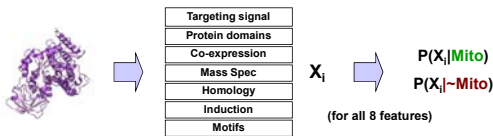
Individual Feature Distributions

Instead of a single big distribution, we have a smaller one for each feature (and class)

$P(\text{Target}|\text{Mito})$ $P(\text{Target}|\sim\text{Mito})$
 $P(\text{Domain}|\text{Mito})$ $P(\text{Domain}|\sim\text{Mito})$
 $P(\text{CE}|\text{Mito})$ $P(\text{CE}|\sim\text{Mito})$
 $P(\text{Mass}|\text{Mito})$ $P(\text{Mass}|\sim\text{Mito})$
 $P(\text{Homology}|\text{Mito})$ $P(\text{Homology}|\sim\text{Mito})$
 $P(\text{Induc}|\text{Mito})$ $P(\text{Induc}|\sim\text{Mito})$
 $P(\text{Motif}|\text{Mito})$ $P(\text{Motif}|\sim\text{Mito})$



Classifying A New Protein



Plug these and priors into the discriminant function

$$G(X_1, \dots, X_7) = \log \frac{\prod P(X_i | \text{Mito}) P(\text{Mito})}{\prod P(X_i | \sim \text{Mito}) P(\sim \text{Mito})} > 0$$

IF $G > 0$, we predict that the protein is from class Mito

Maestro Results

Apply Maestro to Human Proteome

Total predictions: 1,451 proteins
490 novel predictions

Slide Credit: S. Calvo

How Good is the Classifier?

The Rule

We *must* test our classifier on a different set from the training set: the **labeled test set**

The Task

We will classify each object in the test set and count the **number of each type of error**

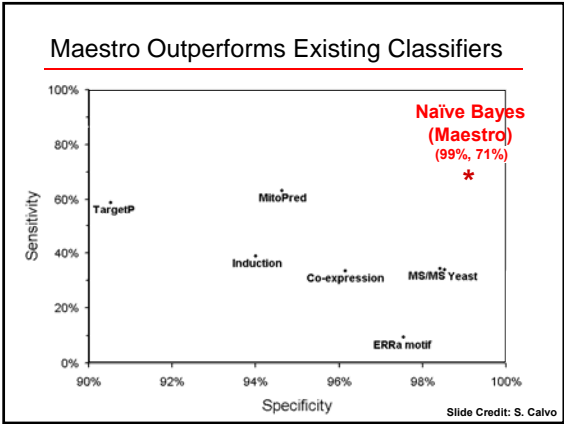
Binary Classification Errors

	True (Mito)	False (~Mito)
Predicted True	TP	FP
Predicted False	FN	TN

$$\text{Sensitivity} = \text{TP}/(\text{TP}+\text{FN}) \quad \text{Specificity} = \text{TN}/(\text{TN}+\text{FP})$$

- Sensitivity**
 - Fraction of all Class 1 (True) that we correctly predicted at Class 1
 - How good are we at finding what we are looking for
- Specificity**
 - Fraction of all Class 2 (False) called Class 2
 - How many of the Class 2 do we filter out of our Class 1 predictions

In both cases, the higher the better



Support Vector Machines

Discriminative Classification

Support Vector Machines (SVMs)

Easy to select a line

But many lines will separate these training data

What line should we choose?

Support Vector Machines (SVMs)

A sensible choice is to select a line that maximizes the *margin* between classes

SVM Formulation

We define a vector w parallel to the separating line

The distance from the line to the origin is b

$x_i \cdot w - b \geq 1$ for $y_i = +1$
 $x_i \cdot w - b \leq -1$ for $y_i = -1$

$y_i (x_i \cdot w - b) \geq 1$

Labels
 ● $Y=+1$
 ■ $Y=-1$

An Optimization Problem

For full derivation, see Burger

Only need dot product of input data!

Minimize $L_D = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j x_i \cdot x_j$ Quadratic Programming

subject to $\sum_i \alpha_i y_i = 0$ and $\alpha_j > 0$

Solve for α

$\alpha_i (y_i (x_i \cdot w - b) - 1) = 0$

$w = \sum_i \alpha_i y_i x_i$ Only some α_i are non-zero

x_i with $\alpha_i > 0$ are the *support vectors*
 w is *determined by these data points!*

Using an SVM

Given a new data point we simply assign it the label:

$$y_i = \text{sign}(\mathbf{w} \cdot \mathbf{x}_{\text{new}})$$

$$= \text{sign}\left(\sum_i \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x}_{\text{new}}\right)$$

Again, only dot product of input data!

Labels
 ● $Y=+1$
 ■ $Y=-1$

Non-linear Classifier

- Some data not linearly separable in low dimensions
- What if we transform it to a higher dimension?

1 dimensional data → Kernel function → 2 dimensional data

Noble, 2006. NATURE BIOTECHNOLOGY 24:1565.

Kernel Mapping

Want a **mapping** from input space to other euclidean space

$$\Phi(x): \mathbb{R}^d \rightarrow H$$

But $\Phi(X)$ can be a mapping to an infinite dimensional space i.e. d points become an infinite number of points

$$X=(x_1, x_2) \rightarrow \Phi(X)=(\phi_1, \phi_2, \phi_3, \dots, \phi_\infty)$$

Rather difficult to work with!

Kernel Mapping

Want a **mapping** from input space to other euclidean space

From previous slide, SVMs only depend on **dot product**

$$\Phi(x): \mathbb{R}^d \rightarrow H \quad \mathbf{x}_i \cdot \mathbf{x}_j \text{ becomes } \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)$$

Here is **trick**: if we have a kernel function such that

$$K(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)$$

We can just use **K** and never know $\Phi(x)$ explicitly!
 $\Phi(X)$ is high dimensional
 K is a scalar

Kernels

So the key step is to take your input data and transform it into a **kernel matrix**

$\Phi(x_i) \cdot \Phi(x_j) = \text{scalar!}$

$K(\mathbf{X}_i, \mathbf{X}_j)$

We have then done two very useful things:

- Transformed X into a high (possibly infinite) dimensional space (where we hope are data are separable)
- Taken dot products in this space to create **scalars**

Example Kernels

$K(x, x_j) = x_i^T x_j$ **Linear**

$K(x, x_j) = (\gamma x_i^T x_j + r)^d$ **Polynomial**

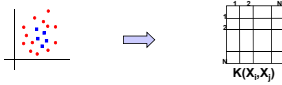
$K(x, x_j) = \exp(-\gamma \|x_i - x_j\|^2)$ **Radial Basis Function**

$K(x, x_j) = \tanh(\gamma x_i^T x_j + r)$ **Sigmoid**

What $K(\mathbf{x}_i, \mathbf{x}_j)$ are valid kernels?
 Answer given by **Mercer's Condition** (see Burgess 1998)

Using (Non-Linear) SVMs

Step 1 – Transform data to **Kernel Matrix K**



Step 2 – **Train SVM on transformed data** – get support vectors

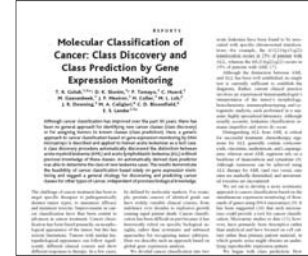
$$\text{Minimize } L_D = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \bullet \mathbf{x}_j = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

Step 2 – **Test/Classify on new samples**

$$y_{\text{new}} = \text{sign}(\mathbf{w} \bullet \mathbf{x}_{\text{new}}) = \text{sign}\left(\sum_i \alpha_i y_i \mathbf{x}_i \bullet \mathbf{x}_{\text{new}}\right) = \text{sign}\left(\sum_i \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}_{\text{new}})\right)$$

Classifying Tumors with Array Data

- **Primary samples:**
 - 38 bone marrow samples
 - 27 ALL, 11 AML
 - obtained from acute leukemia patients at the time of diagnosis;
- **Independent samples:**
 - 34 leukemia samples
 - 24 bone marrow
 - 10 peripheral blood samples
- **Assay ~6800 Genes**



Weighted Voting Classification

General approach of Golub et al (1999) paper:

- Choosing a set of **informative genes** based on their correlation with the class distinction
- Each informative gene casts a **weighted vote** for one of the classes
- Summing up the votes to determine the winning class and the **prediction strength**

Results

Initial Samples

- 36 of the 38 samples as either ALL or ALL.
- All 36 samples **agree** with clinical diagnosis
- 2 not predicted

Independent Samples

- 29 of 34 samples are strongly predicted with **100% accuracy**.
- 5 not predicted

Training Set

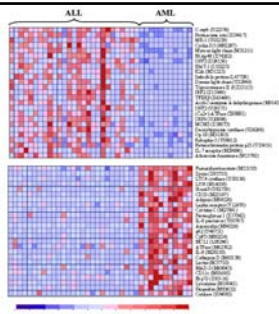
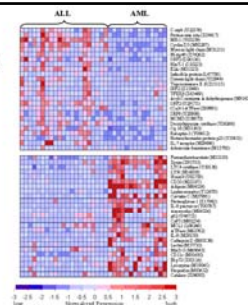


Figure 3b. Genes distinguishing ALL from AML. The 50 genes most highly correlated with the ALL/AML class distinction are shown. Each row corresponds to a gene, with the columns corresponding to expression levels in different samples. Expression levels for each gene are normalized across the samples such that the mean is 0 and the standard deviation is 1. Expression levels greater than the mean are shaded in red, and those below the mean are shaded in blue. The scale indicates standard deviations above or below the mean. The top panel shows genes highly expressed in ALL, the bottom panel shows genes more highly expressed in AML. Note that while these genes as a group appear correlated with class, no single gene is uniformly expressed across the class, illustrating the value of a **multi-gene prediction method**.

Independent Set



Supplementary fig. 2. Expression levels of predictive genes in independent dataset. The expression levels of the 50 genes most highly correlated with the ALL-AML distinction in the initial dataset were determined in the independent dataset. Each row corresponds to a gene, with the columns corresponding to expression levels in different samples. The expression level of each gene in the independent dataset is shown relative to the mean of expression levels for that gene in the initial dataset. Expression levels greater than the mean are shaded in red, and those below the mean are shaded in blue. The scale indicates standard deviations above or below the mean. The top panel shows genes highly expressed in ALL, the bottom panel shows genes more highly expressed in AML.

SVM Approach

Support Vector Machine Classification of Microarray Data

Sayan Mukherjee & Ryan Rifkin

Artificial Intelligence Laboratory and
The Center for Biological and Computational Learning
Massachusetts Institute of Technology
Cambridge, Massachusetts 02139

<http://www.ai.mit.edu>

The Problem: Use the learning from examples paradigm to make class predictions and infer genes involved in these predictions from DNA microarray expression data. Specifically, we use a Support Vector Machine (SVM) classifier [6] to predict cancer morphologies and treatment success and determine the relevant genes in the inference.

Motivation:

Previous Work: A genetic approach to classifying two types of acute leukemia was introduced in Golub et al. [8]. SVMs have been applied to this problem [5] and also to the problem of predicting functional roles of uncharacterized yeast ORFs [1].

Approach: We used a SVM classifier to discriminate between two types of leukemia. The output of classical SVMs is a class designation ± 1 . In this particular application it is important to be able to report points for which the classifier is not confident enough. We introduced a confidence interval on the output of the SVM that allows us to reject genes with low confidence values. It is also important in this application to infer which genes are important for the classification. We have preliminary results for a feature selection algorithm for SVM classifiers.



Methods

- Generate 4 classifiers using different numbers of genes
 - 7129, 999, 99, 49 most informative
- Linear SVM
- Distance from hyperplane (i.e. margin) provides **confidence level**

Results

genes	rejects	errors	confidence level	$ d $
7129	3	0	~ 93%	.1
40	0	0	~ 93%	.1
5	3	0	~ 92%	.1

Results

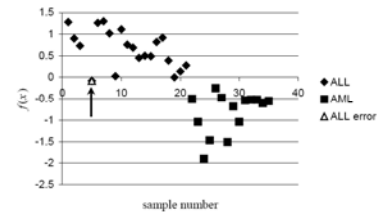
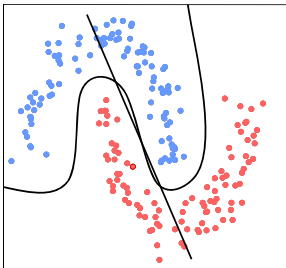


Figure 9.6 The signed distance, $f(x)$, from the optimal separating hyperplane for the test samples. The diamonds are the correctly labeled ALL samples. The squares indicate the correctly labeled AML samples. The triangle marks the misclassified ALL case (see arrow).

Bringing Clustering and Classification Together

Semi-Supervised Learning



Common Scenario

- Few labeled
- Many unlabeled
- Structured data

What if we cluster first?

Then clusters can help us classify