

# Statistics 202: Data Mining

## Linear Discriminant Analysis

Based in part on slides from textbook, slides of Susan Holmes

©Jonathan Taylor

November 9, 2012

# Discriminant analysis

Statistics 202:  
Data Mining

© Jonathan  
Taylor

## Nearest centroid rule

- Suppose we break down our data matrix as by the labels yielding  $(\mathbf{X}^j)_{1 \leq j \leq k}$  with sizes  $\text{nrow}(\mathbf{X}^j) = n_j$ .
- A simple rule for classification is:

*Assign a new observation with features  $\mathbf{x}$  to*

$$\hat{f}(\mathbf{x}) = \underset{1 \leq j \leq k}{\operatorname{argmin}} d(\mathbf{x}, \mathbf{X}^j)$$

- What do we mean by distance here?

# Discriminant analysis

Statistics 202:  
Data Mining

©Jonathan  
Taylor

## Nearest centroid rule

- If we can assign a *central point* or *centroid*  $\hat{\mu}_j$  to each  $\mathbf{X}^j$ , then we can define the distance above as distance to the centroid  $\hat{\mu}_j$ .
- This yields the nearest centroid rule

$$\hat{f}(\mathbf{x}) = \operatorname{argmin}_{1 \leq j \leq k} d(\mathbf{x}, \hat{\mu}_j)$$

# Discriminant analysis

Statistics 202:  
Data Mining

© Jonathan  
Taylor

## Nearest centroid rule

- This rule is described completely by the functions

$$h_{ij}(\mathbf{x}) = \frac{d(\mathbf{x}, \hat{\mu}_j)}{d(\mathbf{x}, \hat{\mu}_i)}$$

with  $\hat{f}(\mathbf{x})$  being any  $i$  such that

$$h_{ij}(\mathbf{x}) \geq 1 \quad \forall j.$$

# Discriminant analysis

Statistics 202:  
Data Mining

© Jonathan  
Taylor

## Nearest centroid rule

- If  $d(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2$  then the natural centroid is

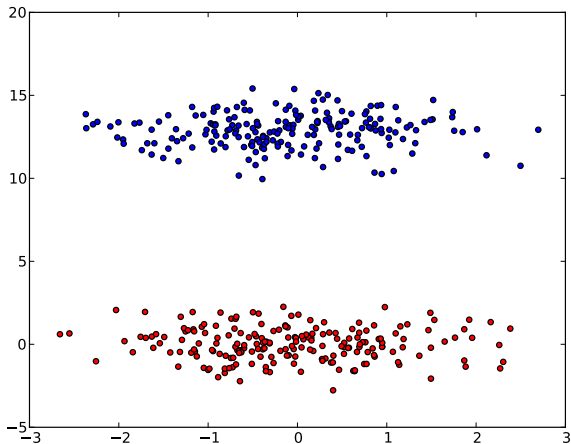
$$\hat{\mu}_j = \frac{1}{n_j} \sum_{l=1}^{n_j} \mathbf{x}_l^j$$

- This rule just classifies points to the nearest  $\hat{\mu}_j$ .
- But, if the covariance matrix of our data is not  $I$ , this rule ignores structure in the data ...

# Why we should use covariance

Statistics 202:  
Data Mining

© Jonathan  
Taylor



# Discriminant analysis

Statistics 202:  
Data Mining

© Jonathan  
Taylor

## Choice of distance

- Often, there is some background *model* for our data that is equivalent to a given procedure.
- For this nearest centroid rule, using the Euclidean distance effectively assumes that within the set of points  $\mathbf{X}^j$ , the rows are multivariate Gaussian with covariance matrix proportional to  $I$ .
- It also implicitly assumes that the  $n_j$ 's are roughly equal.
- For instance, if one  $n_j$  was very small just because a data point is close to that  $\hat{\mu}_j$  doesn't necessarily mean that we should conclude it has label  $j$  because there might be a huge number of points of label  $i$  near that  $\hat{\mu}_j$ .

# Discriminant analysis

Statistics 202:  
Data Mining

© Jonathan  
Taylor

## Gaussian discriminant functions

- Suppose each group with label  $j$  had its own mean  $\mu_j$  and covariance matrix  $\Sigma_j$ , as well as proportion  $\pi_j$ .
- The Gaussian discriminant functions are defined as

$$h_{ij}(\mathbf{x}) = h_i(\mathbf{x}) - h_j(\mathbf{x})$$

$$h_i(\mathbf{x}) = \log \pi_i - \log |\Sigma_i|/2 - (\mathbf{x} - \mu_i)^T \Sigma_i^{-1} (\mathbf{x} - \mu_i)/2$$

- The first term weights the prior probability, the second two terms are  $\log \phi_{\mu_i, \Sigma_i}(\mathbf{x}) = \log L(\mu_i, \Sigma_i | \mathbf{x})$ .
- Ignoring the first two terms, the  $h_i$ 's are essentially within-group Mahalanobis distances . . .



# Discriminant analysis

Statistics 202:  
Data Mining

©Jonathan  
Taylor

## Gaussian discriminant functions

- The classifier assigns  $\mathbf{x}$  label  $i$  if  $h_i(\mathbf{x}) \geq h_j(\mathbf{x}) \forall j$ .
- Or,

$$f(\mathbf{x}) = \operatorname{argmax}_{1 \leq i \leq k} h_i(\mathbf{x})$$

- This is equivalent to a Bayesian rule. We'll see more Bayesian rules when we talk about naïve Bayes ...
- When all  $\Sigma_i$  and  $\pi_i$ 's are identical, the classifier is just nearest centroid using Mahalanobis distance instead of Euclidean distance.

# Discriminant analysis

Statistics 202:  
Data Mining

©Jonathan  
Taylor

## Estimating discriminant functions

- In practice, we will have to estimate  $\pi_j, \mu_j, \Sigma_j$ .
- Obvious estimates:

$$\hat{\pi}_j = \frac{n_j}{\sum_{j=1}^k n_j}$$

$$\hat{\mu}_j = \frac{1}{n_j} \sum_{l=1}^{n_j} \mathbf{x}_l^j$$

$$\hat{\Sigma}_j = \frac{1}{n_j - 1} (\mathbf{X}^j - \hat{\mu}_j \mathbf{1})^T (\mathbf{X}^j - \hat{\mu}_j \mathbf{1})$$

# Discriminant analysis

Statistics 202:  
Data Mining

© Jonathan  
Taylor

## Estimating discriminant functions

- If we assume that the covariance matrix is the same within groups, then we might also form the pooled estimate

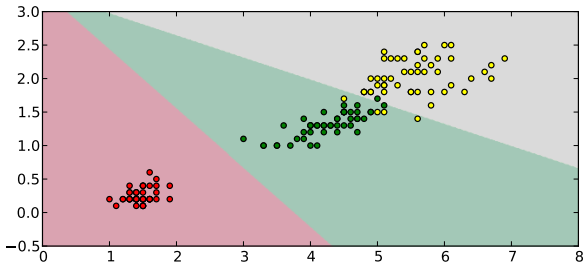
$$\hat{\Sigma}_P = \frac{\sum_{j=1}^k (n_j - 1) \hat{\Sigma}_j}{\sum_{j=1}^k n_j - 1}$$

- If we use the pooled estimate  $\Sigma_j = \hat{\Sigma}_P$  and plug these into the Gaussian discriminants, the functions  $h_{ij}(\mathbf{x})$  are linear (or affine) functions of  $\mathbf{x}$ .
- This is called Linear Discriminant Analysis (LDA).
- Not to be confused with the other LDA (Latent Dirichlet Allocation) ...

# Linear Discriminant Analysis using (petal.width, petal.length)

Statistics 202:  
Data Mining

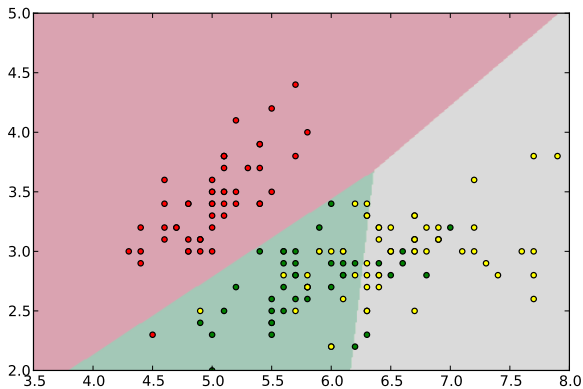
© Jonathan  
Taylor



# Linear Discriminant Analysis using (sepal.width, sepal.length)

Statistics 202:  
Data Mining

© Jonathan  
Taylor



# Discriminant analysis

Statistics 202:  
Data Mining

© Jonathan  
Taylor

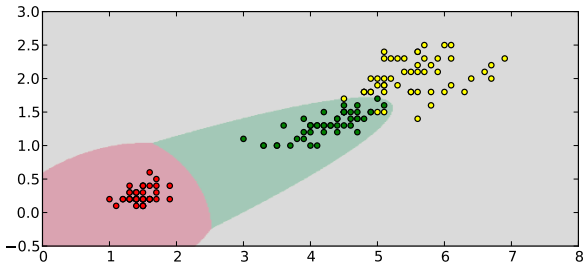
## Quadratic Discriminant Analysis

- If we don't use pooled estimate  $\Sigma_j = \hat{\Sigma}_j$  and plug these into the Gaussian discriminants, the functions  $h_{ij}(\mathbf{x})$  are *quadratic* functions of  $\mathbf{x}$ .
- This is called Quadratic Discriminant Analysis (QDA).

# Quadratic Discriminant Analysis using (petal.width, petal.length)

Statistics 202:  
Data Mining

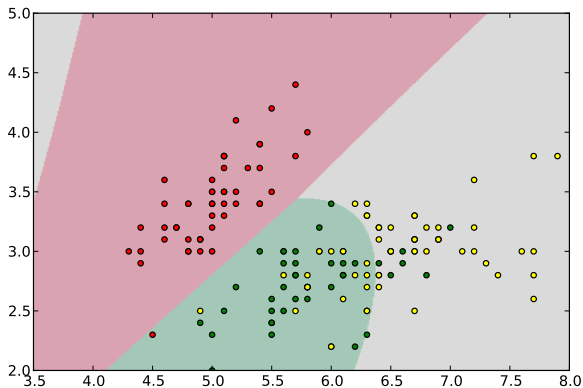
© Jonathan  
Taylor



# Quadratic Discriminant Analysis using (sepal.width, sepal.length)

Statistics 202:  
Data Mining

© Jonathan  
Taylor

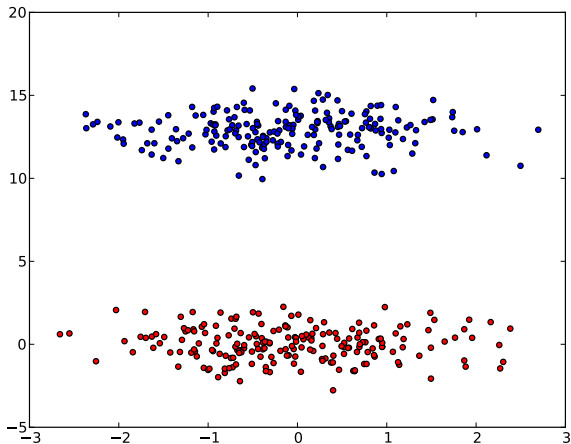




# Motivation for Fisher's rule

Statistics 202:  
Data Mining

© Jonathan  
Taylor



# Discriminant analysis

Statistics 202:  
Data Mining

© Jonathan  
Taylor

## Fisher's discriminant function

- Fisher proposed to classify using a linear rule.
- He first **decomposed**

$$(\mathbf{X} - \hat{\mu}\mathbf{1})^T(\mathbf{X} - \hat{\mu}\mathbf{1}) = \widehat{SS}_B + \widehat{SS}_W$$

- Then, he proposed,

$$\hat{v} = \operatorname{argmax}_{v: v^T \widehat{SS}_W v = 1} v^T \widehat{SS}_B v$$

- Having found  $\hat{v}$ , form  $\mathbf{X}^j \hat{v}$  and the centroid  $\eta_j = \operatorname{mean}(\mathbf{X}^j \hat{v})$
- In the two-class problem  $k = 2$ , this is the same as LDA.

# Discriminant analysis

Statistics 202:  
Data Mining

©Jonathan  
Taylor

## Fisher's discriminant functions

- The direction  $\hat{v}_1$  is an eigenvector of some matrix. There are others, up to  $k - 2$  more.
- Suppose we find all  $k - 1$  vectors and form  $\mathbf{X}_j V^T$ , each one an  $n_j \times (k - 1)$  matrix with centroid  $\eta_j \in \mathbb{R}^{k-1}$ .
- The matrix  $V$  determines a map from  $\mathbb{R}^p$  to  $\mathbb{R}^{k-1}$ , so given a new data point we can compute  $V\mathbf{x} \in \mathbb{R}^{k-1}$ .
- This gives rise to a classifier

$$\hat{f}(\mathbf{x}) = \text{nearest centroid}(V\mathbf{x})$$

- This is LDA (assuming  $\pi_j = \frac{1}{k}$ ) ...

# Discriminant analysis

Statistics 202:  
Data Mining

© Jonathan  
Taylor

## Discriminant models in general

- A discriminant model is generally a model that estimates

$$P(Y = j|\mathbf{x}), 1 \leq j \leq k$$

- That is, given that the features I observe are  $\mathbf{x}$ , the probability I think this label is  $j$  ...
- LDA and QDA are actually *generative models* since they specify

$$P(X = \mathbf{x}|Y = j).$$

- There are lots of discriminant models ...

# Discriminant analysis

Statistics 202:  
Data Mining

©Jonathan  
Taylor

## Logistic regression

- The logistic regression model is ubiquitous in binary classification (two-class) problems
- Model:

$$P(Y = 1|\mathbf{x}) = \frac{\alpha + e^{\mathbf{x}^T \beta}}{1 + e^{\alpha + \mathbf{x}^T \beta}} = \pi(\alpha, \beta, \mathbf{x})$$

# Discriminant analysis

Statistics 202:  
Data Mining

© Jonathan  
Taylor

## Logistic regression

- Software that fits a logistic regression model produces an estimate of  $\beta$  based on a data matrix  $\mathbf{X}_{n \times p}$  and binary labels  $\mathbf{Y}_{n \times 1} \in \{0, 1\}^n$
- It fits the model minimizing what we call the *deviance*

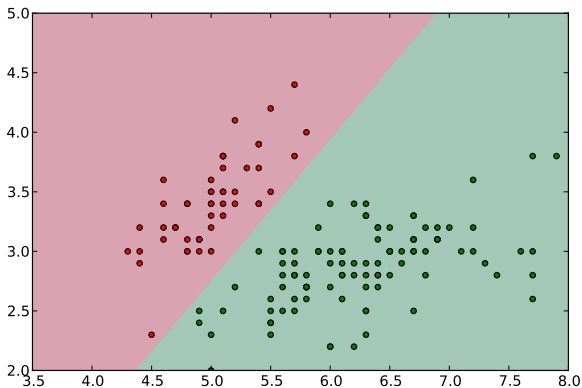
$$\text{DEV}(\beta) = -2 \sum_{i=1}^n (\mathbf{Y}_i \log \pi(\beta, \mathbf{X}_i) + (1 - \mathbf{Y}_i) \log(1 - \pi(\beta, \mathbf{X}_i)))$$

- While not immediately obvious, this is a convex minimization problem, hence is fairly easy to solve.
- Unlike trees, the convexity yields a globally optimal solution.

# Logistic regression, *setosa* vs *virginica*, *versicolor* using (sepal.width, sepal.length)

Statistics 202:  
Data Mining

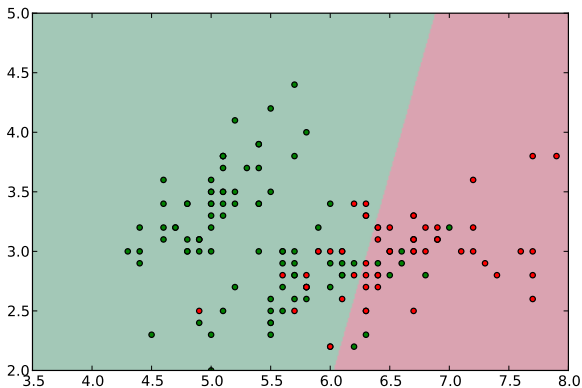
© Jonathan  
Taylor



# Logistic regression, *virginica* vs *setosa*, *versicolor* using (sepal.width, sepal.length)

Statistics 202:  
Data Mining

© Jonathan  
Taylor

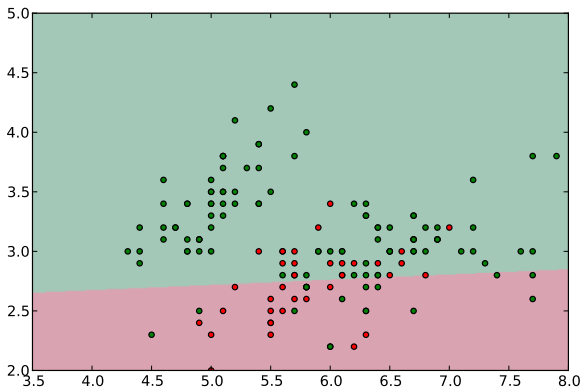




# Logistic regression, *versicolor* vs *setosa*, *virginica* using (sepal.width, sepal.length)

Statistics 202:  
Data Mining

© Jonathan  
Taylor



# Discriminant analysis

Statistics 202:  
Data Mining

©Jonathan  
Taylor

## Logistic regression

- Logistic regression produces an estimate of

$$P(\mathbf{y} = 1|\mathbf{x}) = \hat{\pi}(\mathbf{x})$$

- Typically, we classify as 1 if  $\hat{\pi}(\mathbf{x}) > 0.5$ .
- This yields a  $2 \times 2$  confusion matrix

	Predicted: 0	Predicted: 1
Actual : 0	<i>TN</i>	<i>FP</i>
Actual : 1	<i>FN</i>	<i>TP</i>

- From the  $2 \times 2$  confusion matrix, we can compute Sensitivity, Specificity, etc.

# Discriminant analysis

Statistics 202:  
Data Mining

© Jonathan  
Taylor

## Logistic regression

- However, we could choose the threshold differently, perhaps related to estimates of the prior probabilities of 0's and 1's.
- Now, each threshold  $0 \leq t \leq 1$  yields a new confusion matrix
- This yields a  $2 \times 2$  confusion matrix

	Predicted: 0	Predicted: 1
Actual : 0	$TN(t)$	$FP(t)$
Actual : 1	$FN(t)$	$TP(t)$

# Discriminant analysis

Statistics 202:  
Data Mining

© Jonathan  
Taylor

## ROC curve

- Generally speaking, we prefer classifiers that are both highly sensitive and highly specific.
- These confusion matrices can be summarized using an ROC (Receiver Operating Characteristic) curve.
- This is a plot of the curve

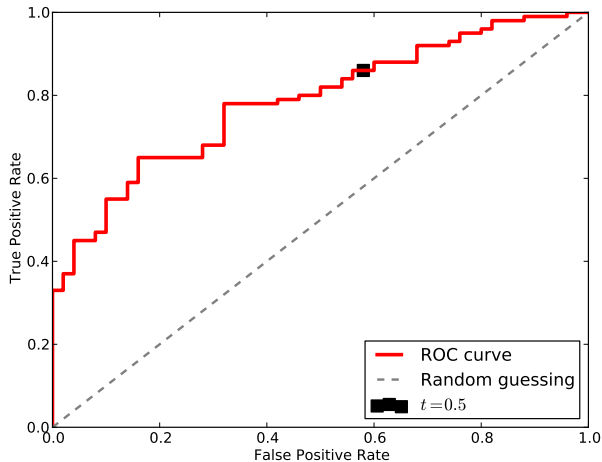
$$(1 - \text{Specificity}(t), \text{Sensitivity}(t))_{0 \leq t \leq 1}$$

- Often, Specificity is referred to as *TNR* (True Negative Rate), and  $(1 - \text{Specificity}(t))$  as *FPR*.
- Often, Sensitivity is referred to as *TPR* (True Positive Rate), and  $(1 - \text{Sensitivity}(t))$  as *FNR*.

# AUC: Area under ROC curve

Statistics 202:  
Data Mining

© Jonathan  
Taylor



# Discriminant analysis

Statistics 202:  
Data Mining

©Jonathan  
Taylor

## ROC curve

- Points in the upper left of the ROC curve are good.
- Any point on the diagonal line represents a classifier that is guessing randomly.
- A tree classifier as we've discussed seems to correspond to only one point in the ROC plot.
- *But* one can estimate probabilities based on frequencies in terminal nodes.

# Discriminant analysis

Statistics 202:  
Data Mining

©Jonathan  
Taylor

## ROC curve

- A common numeric summary of the ROC curve is

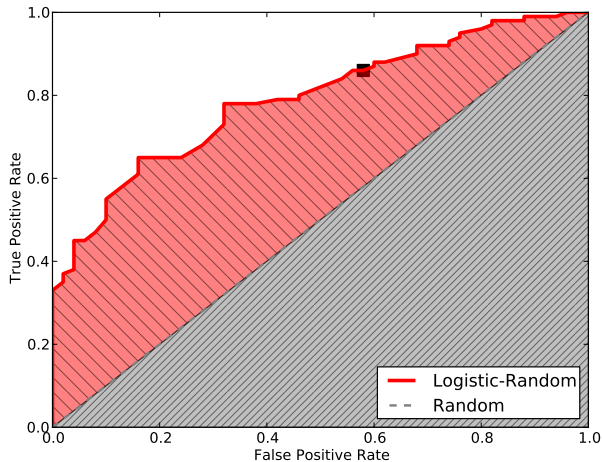
$$AUC(ROC \text{ curve}) = \text{Area under ROC curve.}$$

- Can be interpreted as an estimate of the probability that the classifier will give a random positive instance a higher score than a random negative instance.
- Maximum value is 1.
- For a random guesser, AUC is 0.5

# AUC: Area under ROC curve

Statistics 202:  
Data Mining

© Jonathan  
Taylor

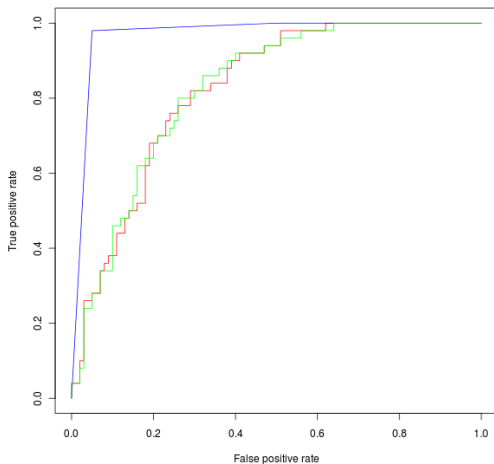




# ROC curve: logistic, rpart, lda

Statistics 202:  
Data Mining

© Jonathan  
Taylor



Statistics 202:  
Data Mining

© Jonathan  
Taylor