

Classification

Lecture 16

Classifiers

Covered so far

K-Nearest Neighbors

Linear regression (can be used, but not recommended)

Perceptron

Logistic Regression

Fisher's Linear Discriminant / Linear Discriminant Analysis

Quadratic Discriminant Analysis

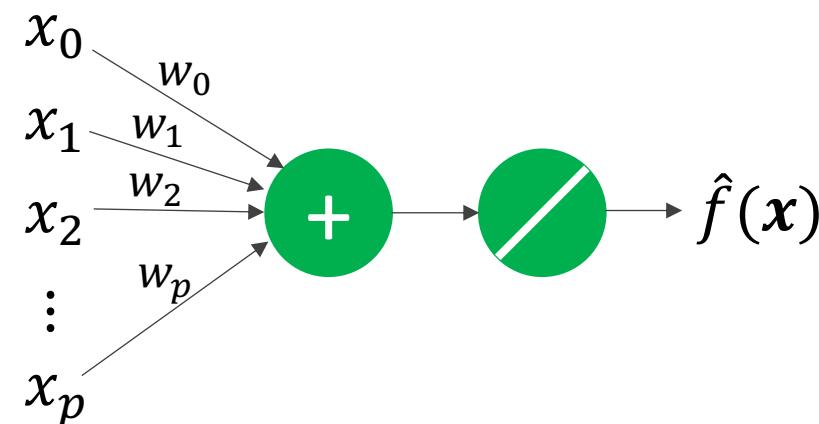
Naïve Bayes

Rely on a linear combination of weights and features: $\mathbf{w}^T \mathbf{x}$

Remember linear models?

Linear Regression

$$\hat{f}(\mathbf{x}) = \sum_{i=0}^N w_i x_i$$

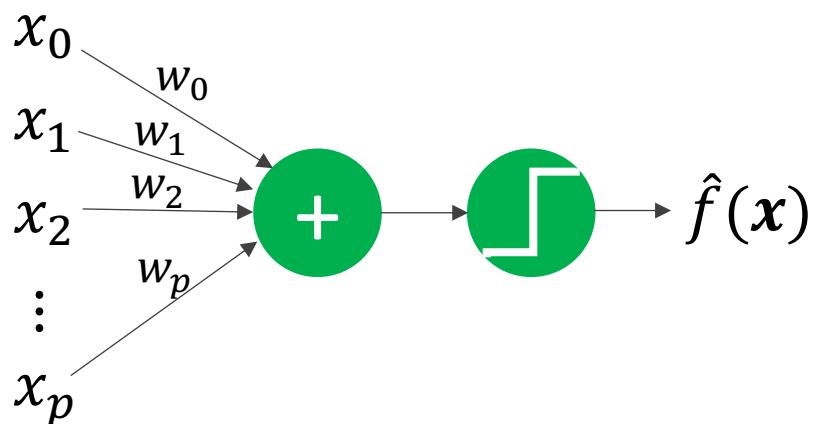


Linear Classification

Perceptron

$$\hat{f}(\mathbf{x}) = \text{sign} \left(\sum_{i=0}^N w_i x_i \right)$$

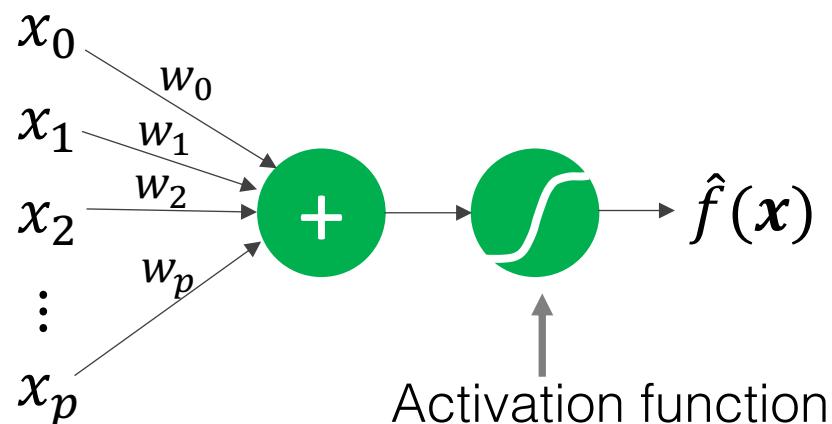
$$\text{sign}(x) = \begin{cases} 1 & x > 0 \\ -1 & \text{else} \end{cases}$$



Logistic Regression

$$\hat{f}(\mathbf{x}) = \sigma \left(\sum_{i=0}^N w_i x_i \right)$$

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$



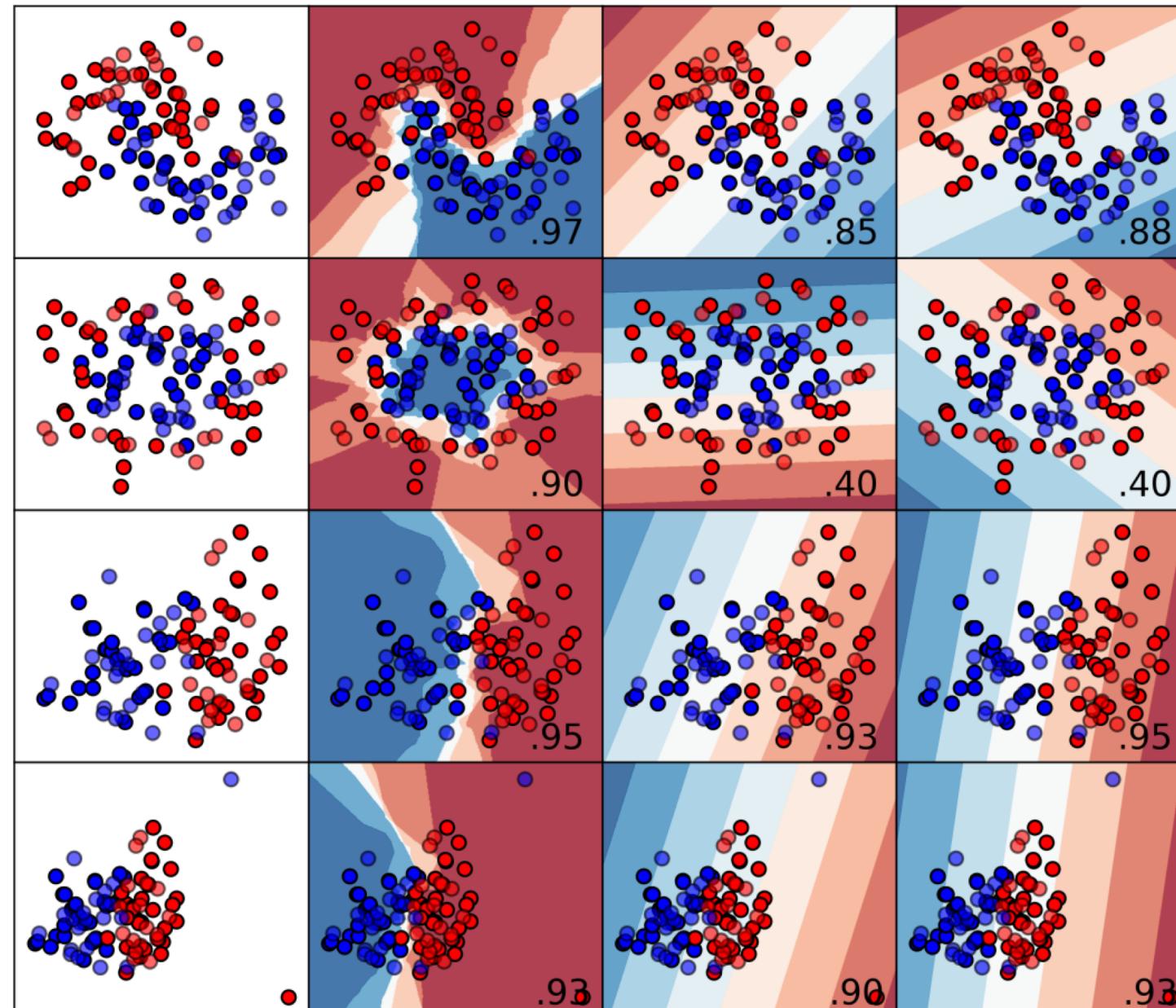
Source: Abu-Mostafa, Learning from Data, Caltech

Input data

KNN (k=5)

Perceptron

Logistic Reg.

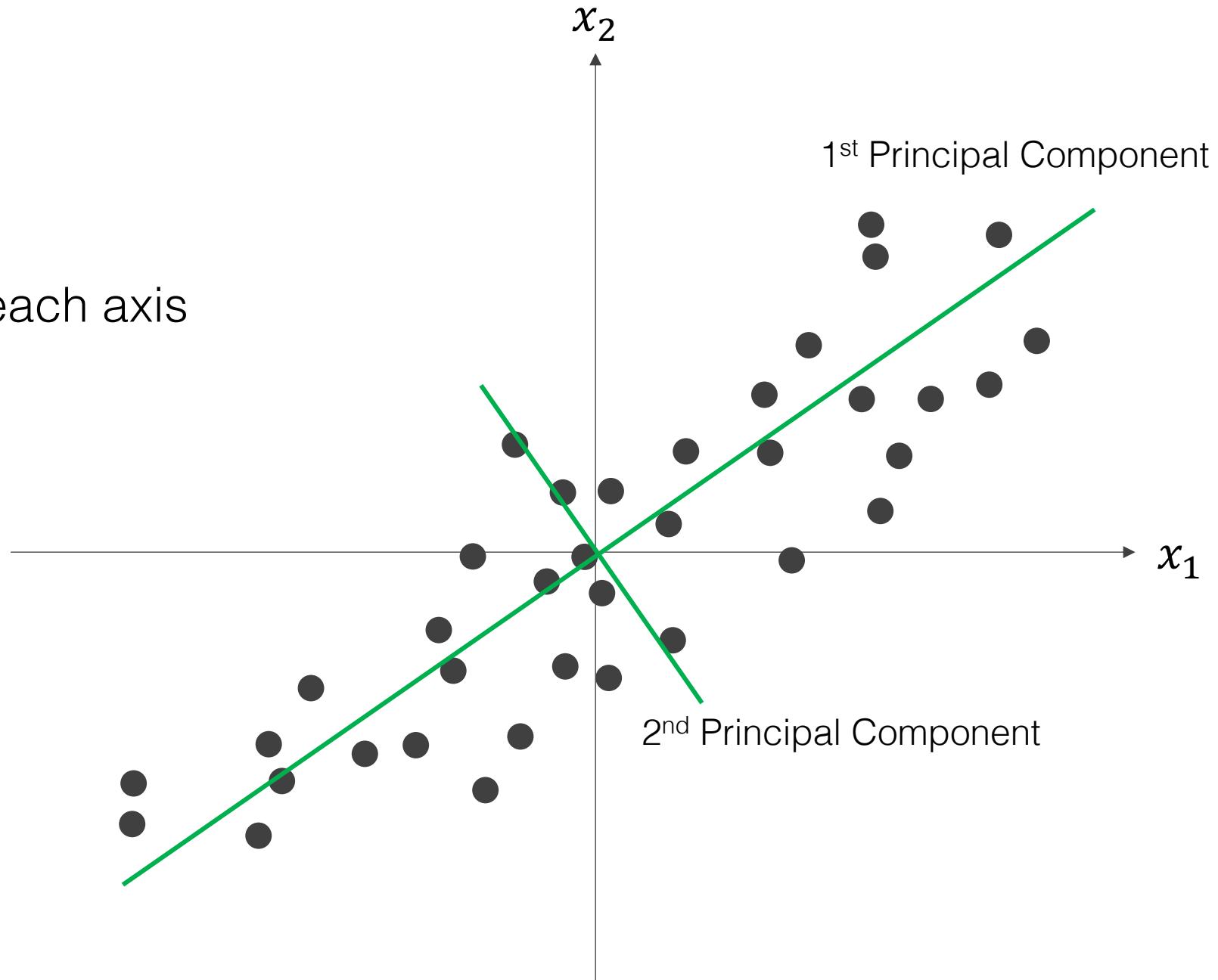


Comparison of classifiers
we've seen so far

Principal Components

Maximize the variance along each axis

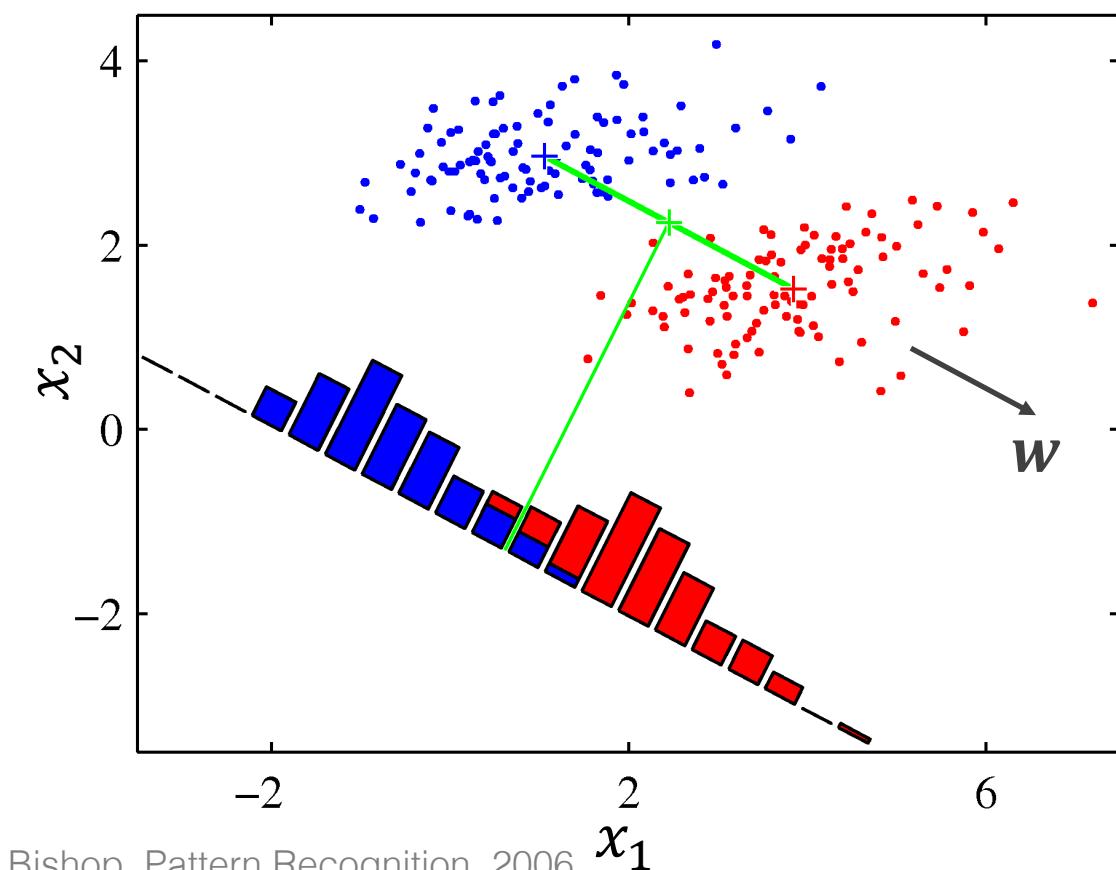
We want to maximize
the **separability** of
classes instead



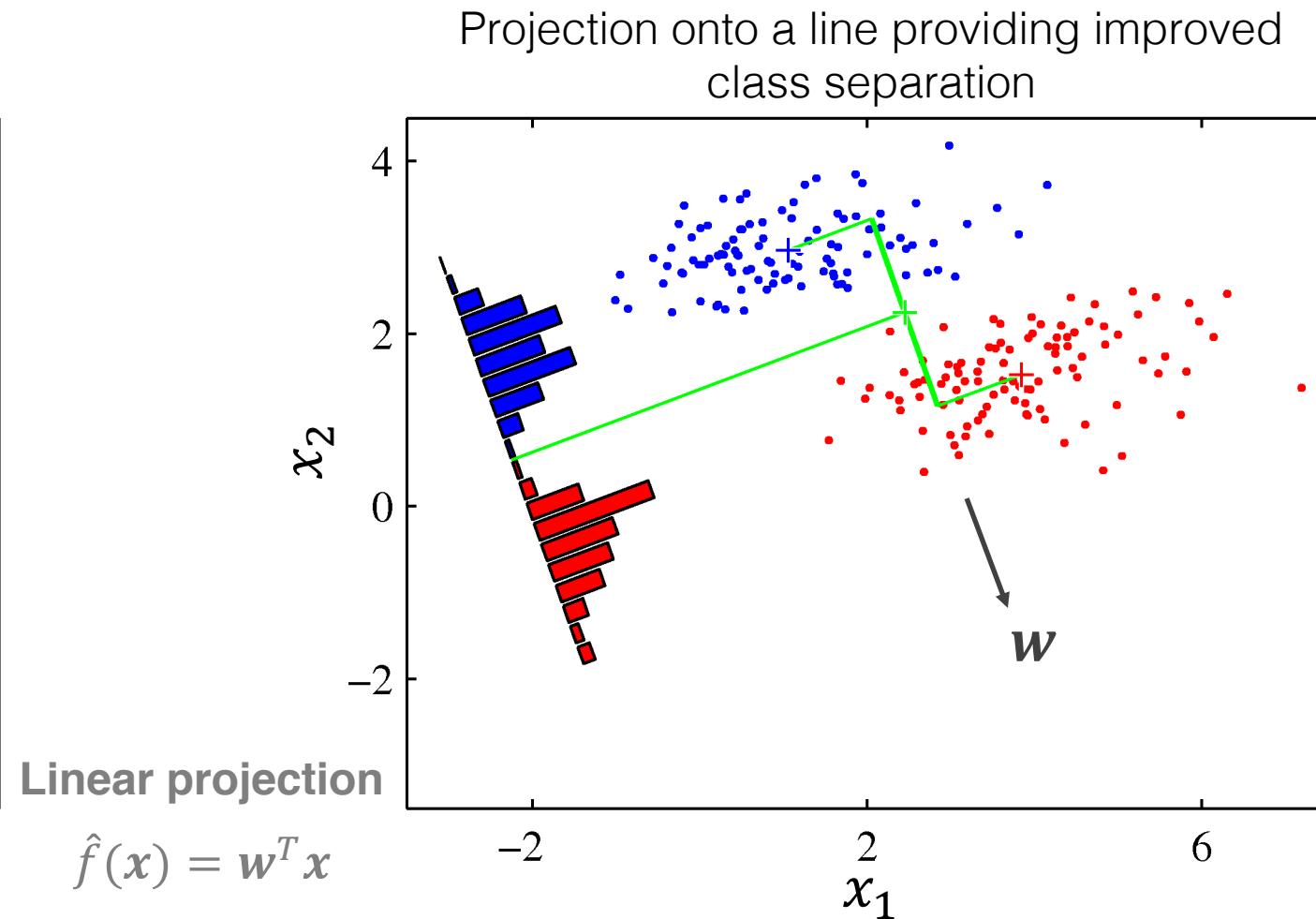
Fisher's Linear Discriminant

Looks for the projection into the one dimension that “best” separates the classes

Projection onto line connecting the means



Projection onto a line providing improved class separation



Linear projection

$$\hat{f}(x) = w^T x$$

Fisher's Linear Discriminant (FLD)

- 1 Finds a projection into a lower dimension that “best” separates the classes

$$\hat{f}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$

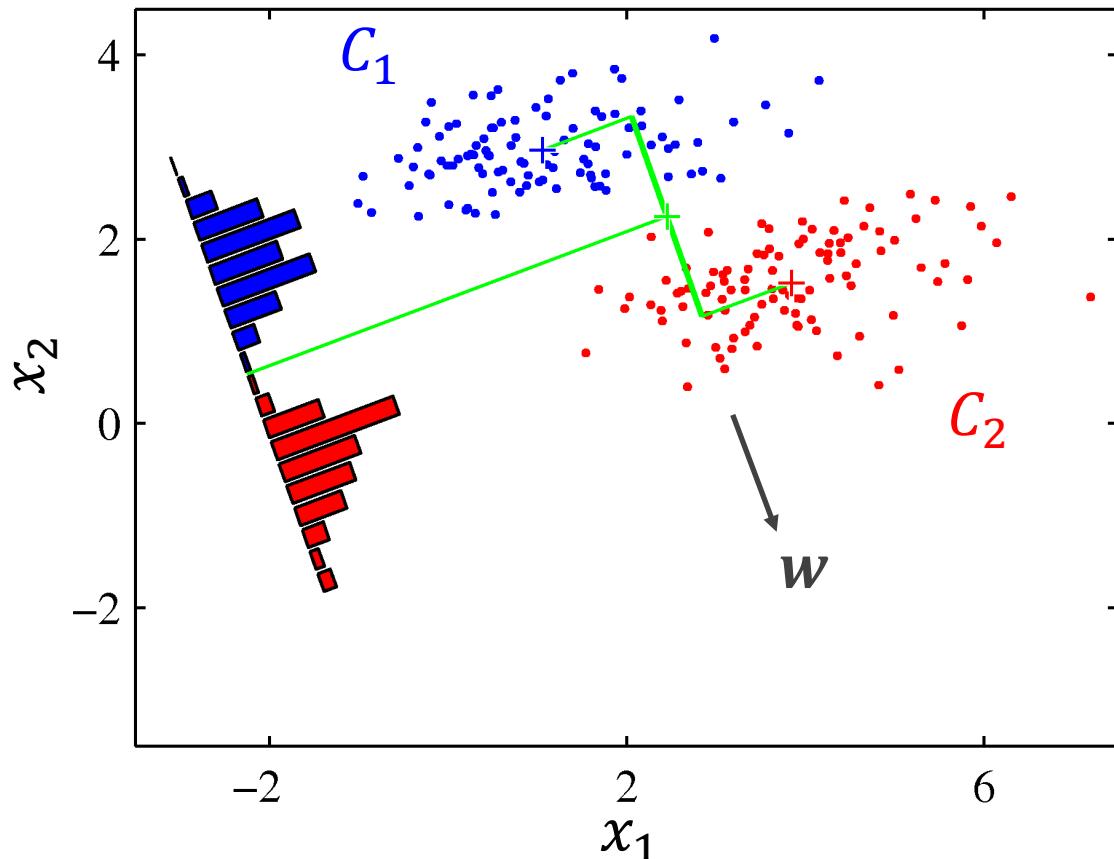
- 2 We then classify the data in this space

Similar to PCA, but accounts for class separability

Our decision rule becomes:

if	$\hat{f}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} > \lambda_{thresh}$	Class 1
else		Class 2

FLD: how do we choose the vector w ?



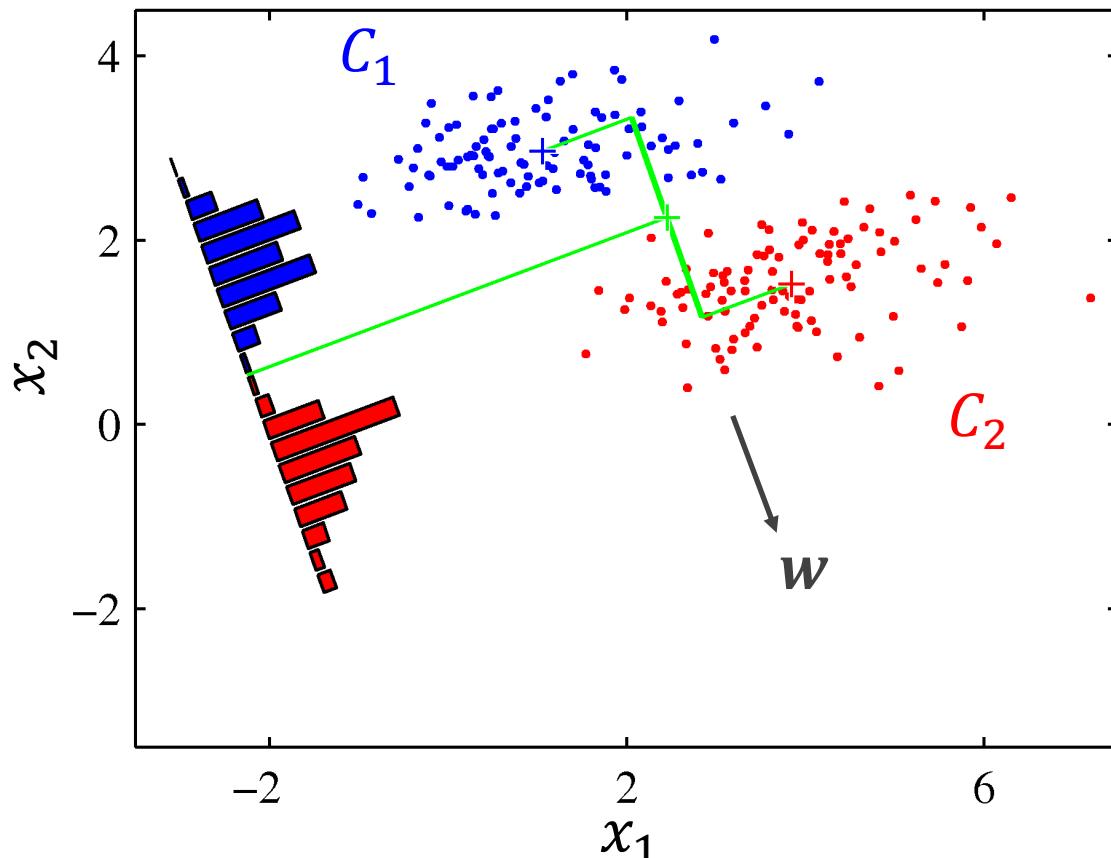
Increase the distance between the **means**

Decrease the **variance** within each class

$$y = \hat{f}(x) = w^T x$$

FLD: how do we choose the vector w ?

Increase the distance between the **means**



$$y = \hat{f}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$

$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{i \in C_1} \mathbf{x}_i$$

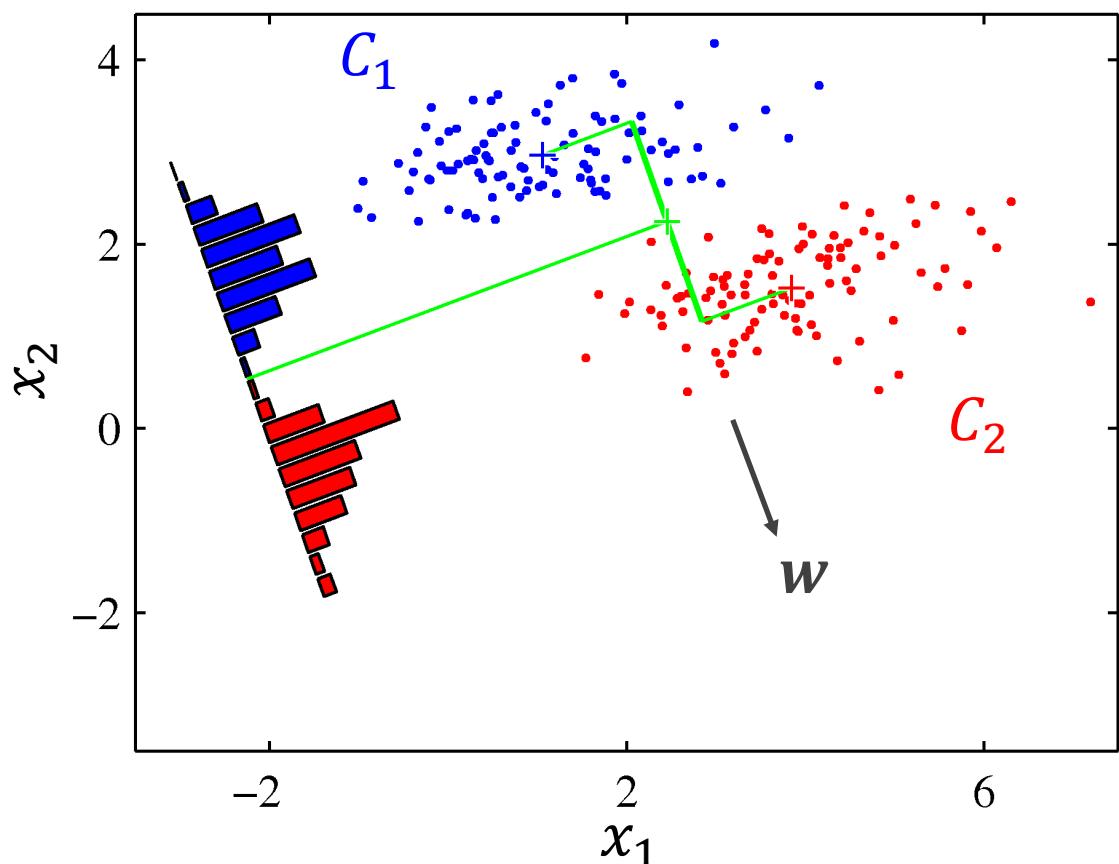
$$\mathbf{m}_2 = \frac{1}{N_2} \sum_{i \in C_2} \mathbf{x}_i$$

The means projected onto w : $m_k = w^T \mathbf{m}_k$

The distance between the means:

$$m_2 - m_1 = w^T(\mathbf{m}_2 - \mathbf{m}_1)$$

FLD: how do we choose the vector w ?



$$y = \hat{f}(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$$

Decrease the **variance** within each class

The “scatter” of the projected data:

$$s_k^2 = \sum_{i \in C_k} (y_i - m_k)^2$$

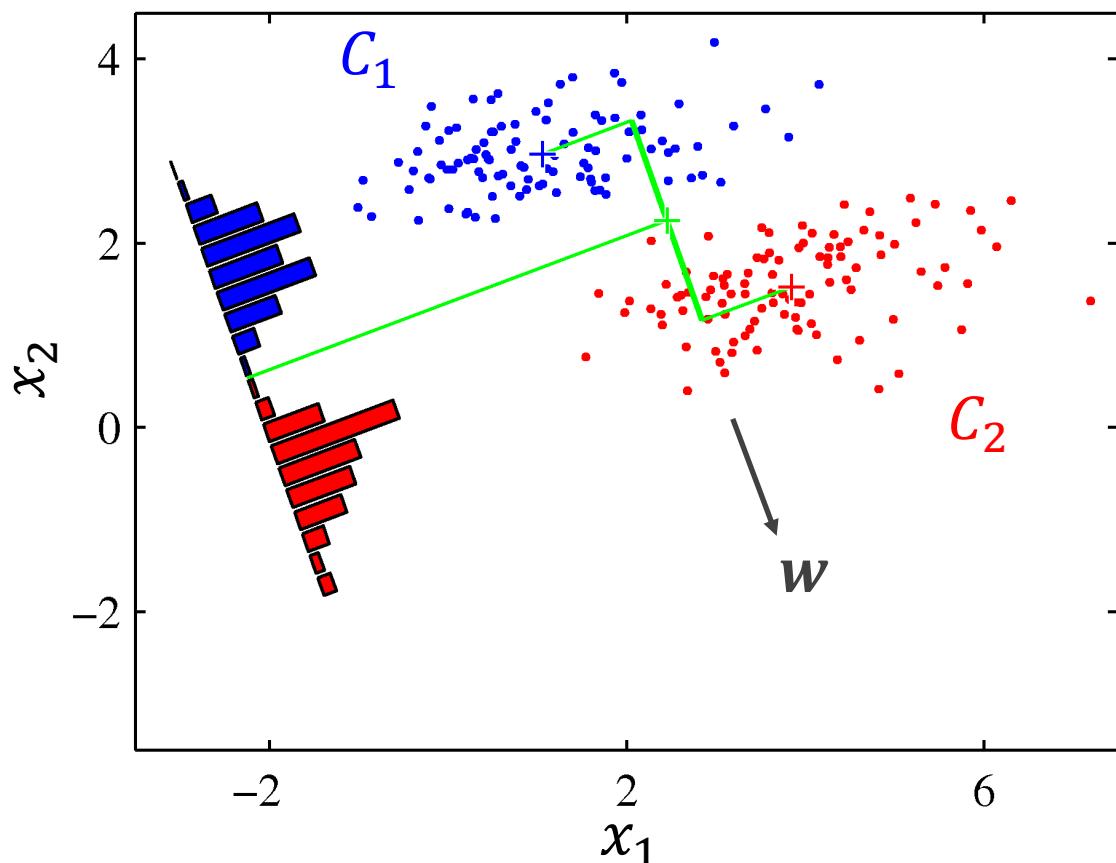
where $m_k = \mathbf{w}^T \mathbf{m}_k$

$$y_i = \mathbf{w}^T \mathbf{x}_i$$

Therefore the total within-class scatter:

$$S = s_1^2 + s_2^2$$

FLD: how do we choose the vector w ?



$$y = \hat{f}(x) = w^T x$$

Increase the distance between the **means**

$$m_2 - m_1 = w^T (m_2 - m_1)$$

Decrease the **variance** within each class

$$S = s_1^2 + s_2^2$$

The Fisher criterion is then:

$$J(w) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}$$

We want to maximize this and solve for w

FLD: how do we choose the vector w ?

We want to maximize this and solve for w

$$J(w) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}$$
$$= \frac{w^T S_B w}{w^T S_W w} \quad (\text{see appendix slides for full derivation})$$

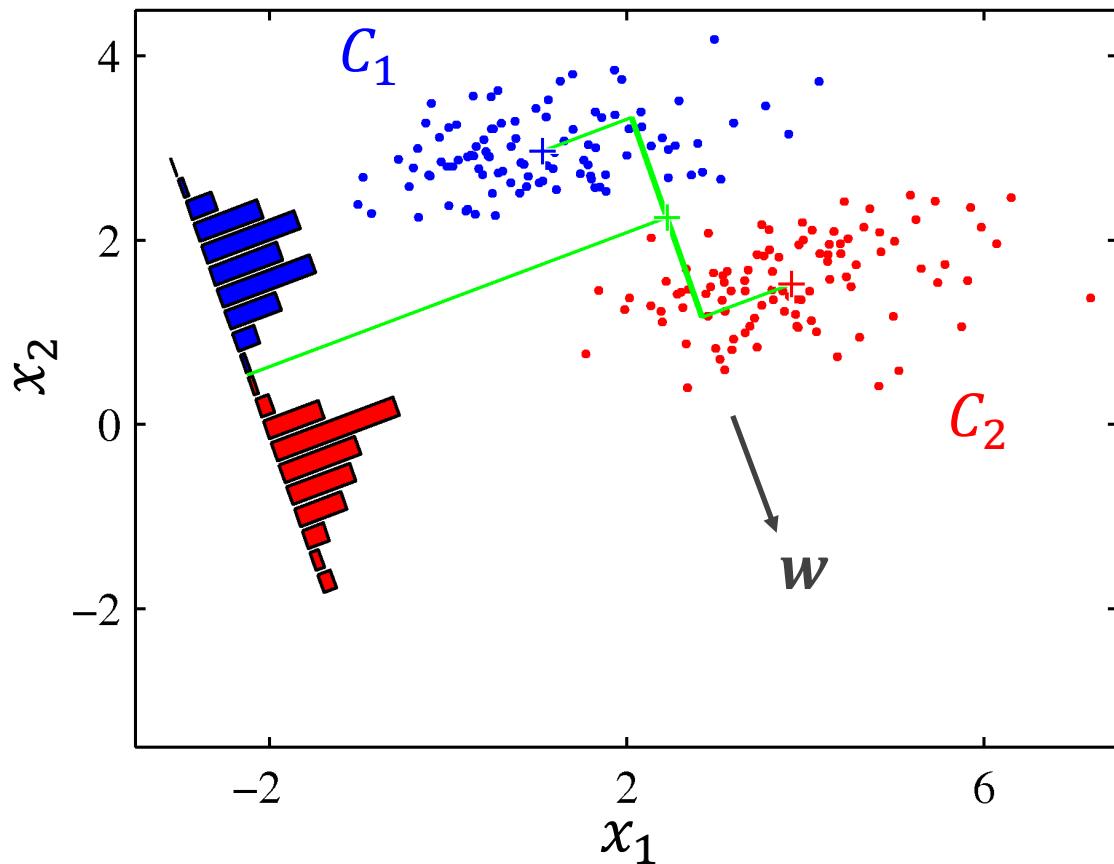
Take the derivative (gradient), set it equal to zero, solve for w

(see appendix slides for full derivation)

$$w \propto S_W^{-1} (m_2 - m_1)$$

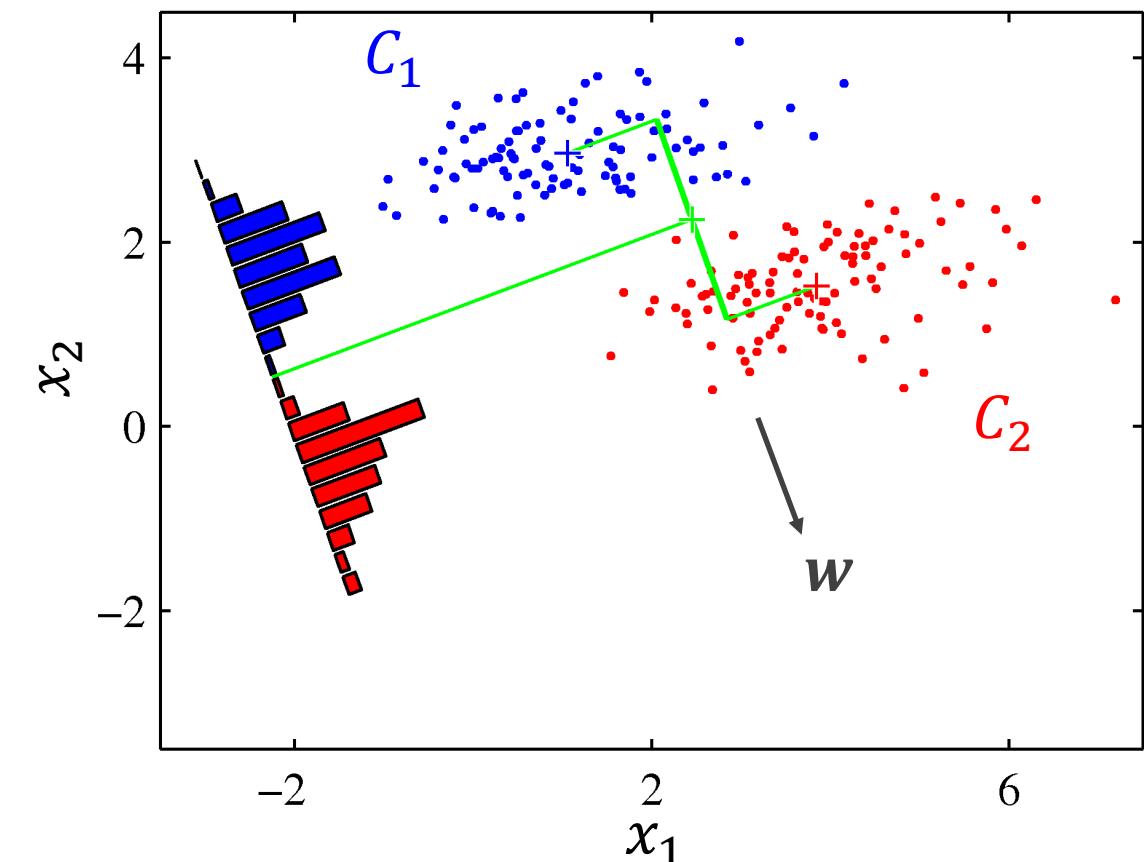
$$w \propto (\Sigma_1 + \Sigma_2)^{-1} (m_2 - m_1)$$

We use this to project the features into one dimension for classification, $w^T x$



$$y = \hat{f}(x) = w^T x$$

Fisher's Linear Discriminant



$$y = \hat{f}(x) = \mathbf{w}^T \mathbf{x}$$

Makes no assumptions about the distribution of the data, and allows for different covariance matrices

This is a **projection** into one dimension that can be used to construct a discriminant

$$\mathbf{w} \propto \mathbf{S}_W^{-1} (\mathbf{m}_2 - \mathbf{m}_1)$$

$$\mathbf{w} \propto (\Sigma_1 + \Sigma_2)^{-1} (\mathbf{m}_2 - \mathbf{m}_1)$$

Discriminant Functions

If we have c different classes, we define a discriminant function, $d_i(\mathbf{x})$ for $i = 1, \dots, c$

If $d_i(\mathbf{x}) > d_j(\mathbf{x})$ for $i \neq j$, then we assign feature \mathbf{x} to class i

$$\begin{aligned}d_i(x) &= P(y = i | \mathbf{x}) = \frac{P(\mathbf{x}|y = i)P(y = i)}{P(\mathbf{x})} \\&= \frac{P(\mathbf{x}|y = i)P(y = i)}{\sum_{i=1}^c P(\mathbf{x}|y = i)P(y = i)}\end{aligned}$$

Bayes' Rule: $P(Y|X) = \frac{\text{Likelihood} \cdot \text{Prior}}{\text{Posterior} \cdot \text{Evidence}}$

$$P(Y|X) = \frac{P(X|Y)P(Y)}{P(X)}$$

We can simply write $d_i(x) = P(\mathbf{x}|y = i)P(y = i)$

Or in log form: $\ln d_i(x) = \ln P(\mathbf{x}|y = i) + \ln P(y = i)$

Discriminant Functions

$$d_i(x) = P(\mathbf{x}|y = i)P(y = i)$$

$$\ln d_i(x) = \ln P(\mathbf{x}|y = i) + \ln P(y = i)$$

Equivalent discriminant functions since log is monotonic

- 1 Assume a form for $P(\mathbf{x}|y = i)$

Gaussian for Linear and Quadratic Discriminant Analysis

Gaussian mixture models

Nonparametric density estimates

Naïve Bayes models

- 2 Assign the class, i , for which $d_i(x)$ is largest

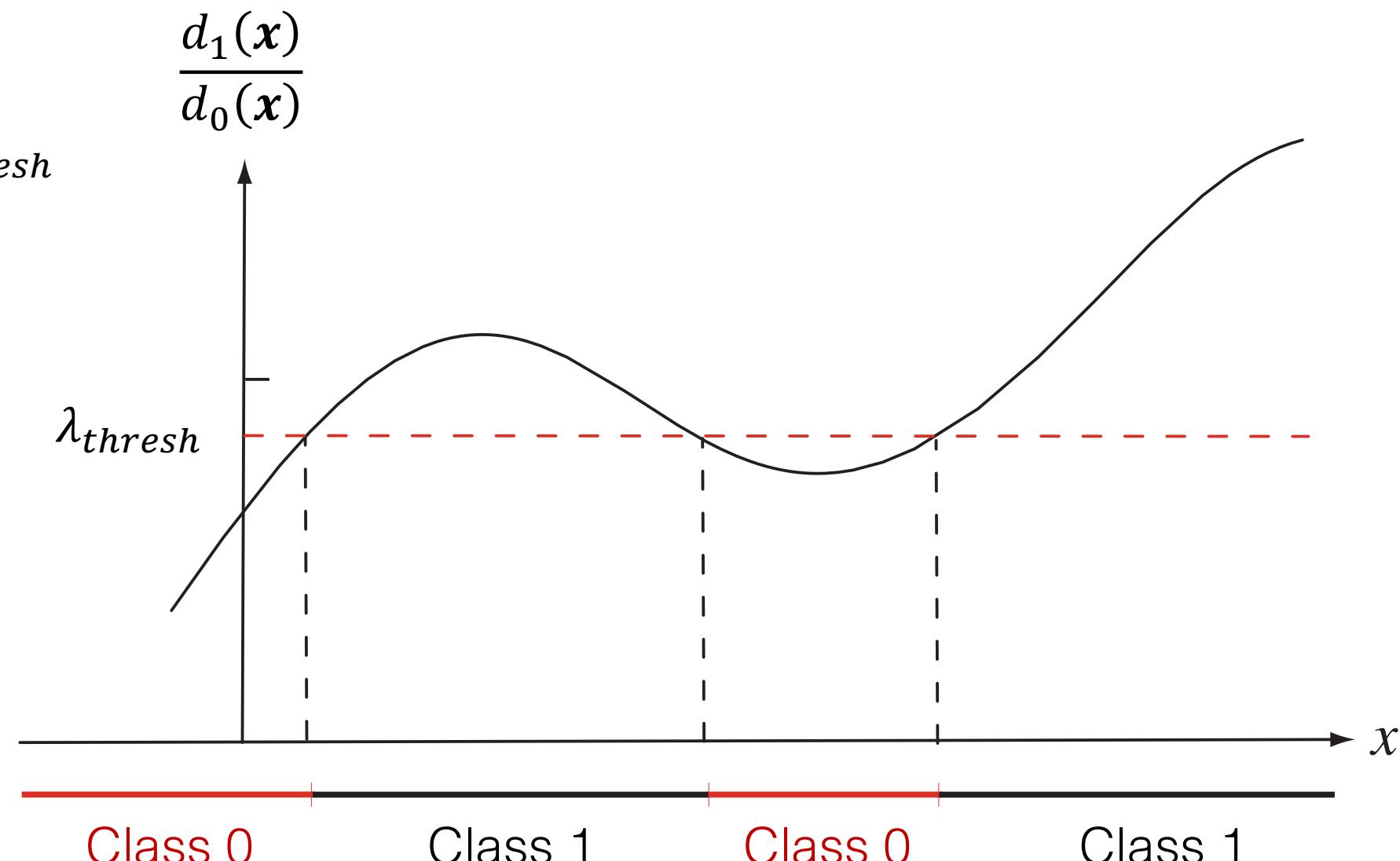
Applies to both binary and multiclass problems

Discriminant Functions: 2 classes

Decision rule:

Class 1 if: $\frac{d_1(x)}{d_0(x)} > \lambda_{thresh}$

Otherwise, class 0

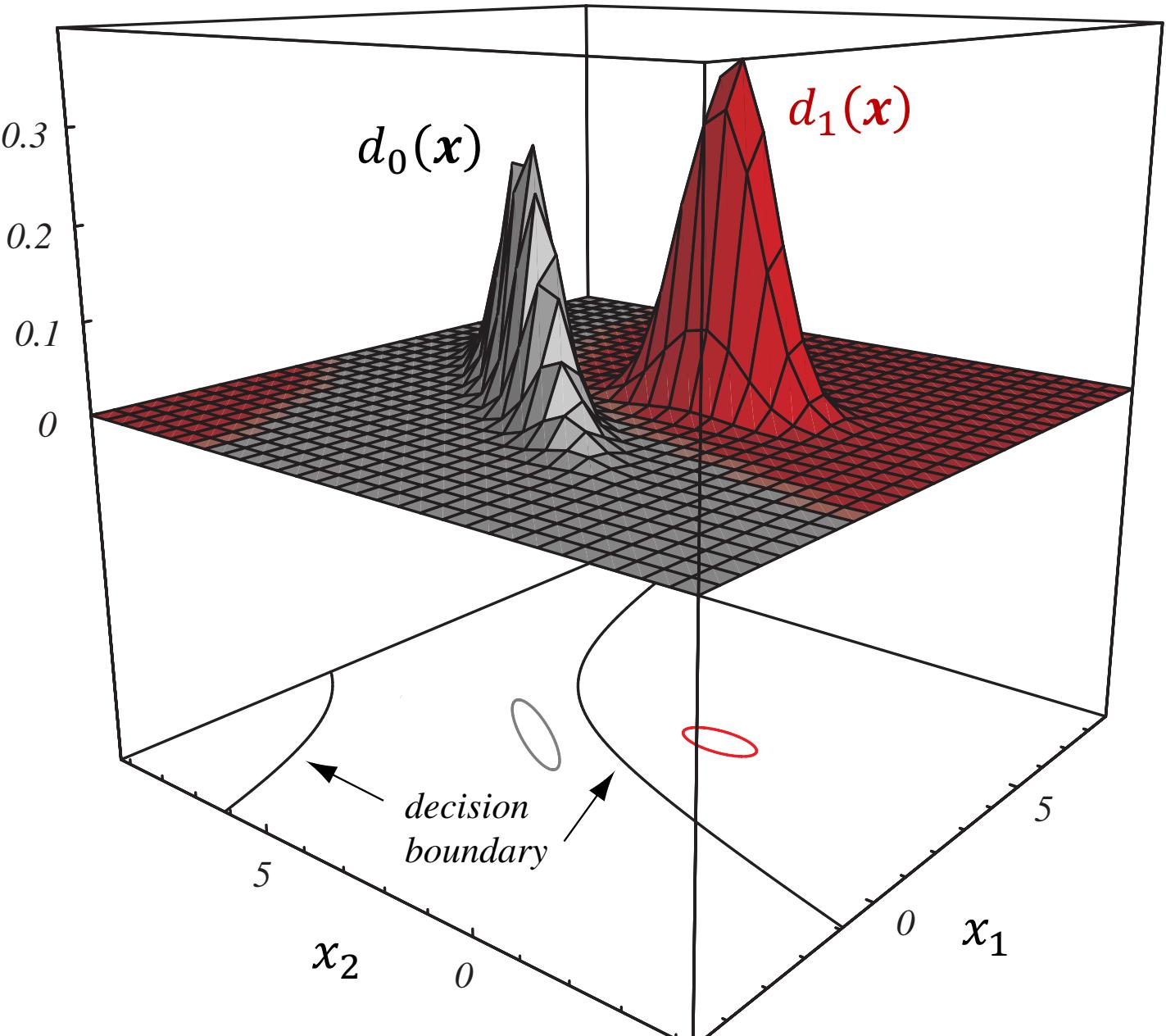


Discriminant Functions: 2 classes, 2 dimensions

Decision rule:

Class 1 if: $\frac{d_1(x)}{d_0(x)} > \lambda_{thresh}$

Otherwise, class 0



Discriminant Function: 2 classes

We build a classifier that assigns the class with the higher posterior probability:

If $\frac{P(y = 1|\mathbf{x})}{P(y = 0|\mathbf{x})} = \frac{P(\mathbf{x}|y = 1)P(y = 1)}{P(\mathbf{x}|y = 0)P(y = 0)} > 1$ Assign class 1, else class 0

Assumes these likelihoods are normal

$$= \frac{P(\mathbf{x}|y = 1)}{P(\mathbf{x}|y = 0)} > \frac{P(y = 0)}{P(y = 1)}$$

$N(\mu_1, \Sigma_1)$ $N(\mu_0, \Sigma_0)$

Estimate the class-conditional mean and covariance matrix from the data

Discriminant Function: 2 classes

We build a classifier that assigns the class with the higher posterior probability:

Likelihood ratio: $\frac{P(\mathbf{x}|y=1)}{P(\mathbf{x}|y=0)} > \frac{P(y=0)}{P(y=1)}$

The diagram consists of two green arrows originating from the terms $P(\mathbf{x}|y=1)$ and $P(\mathbf{x}|y=0)$ in the likelihood ratio equation. The top arrow points to the text $N(\mu_1, \Sigma_1)$. The bottom arrow points to the text $N(\mu_0, \Sigma_0)$.

If we assume the class conditional distributions are Gaussian, this represents

Quadratic Discriminant Analysis

If we further assume the covariance matrices are the same, $\Sigma_0 = \Sigma_1$, this represents

Linear Discriminant Analysis

Comparison

	Fisher's Linear Discriminant (FLD)	Quadratic Discriminant Analysis (QDA)	Linear Discriminant Analysis (LDA)
Assumes Gaussian Likelihood $P(\mathbf{x} y)$ (class conditional density)	No	Yes	Yes
Assumes equivalent covariance matrices $\Sigma_1 = \Sigma_0$	No	No	Yes

FLD and LDA reduce the dimensionality of the data to make them more separable

FLD and LDA are often presented as the same, but there are differences

Linear Discriminant Analysis ($\Sigma_0 = \Sigma_1$)

We build a classifier that assigns the class with the higher posterior probability:

$$\frac{P(\mathbf{x}|y=1)}{P(\mathbf{x}|y=0)} = \frac{\frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_1)\right]}{\frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_0)\right]}$$

$$\ln \left| \frac{P(\mathbf{x}|y=1)}{P(\mathbf{x}|y=0)} \right| = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_0)$$

$$\ln \left| \frac{P(\mathbf{x}|y=1)}{P(\mathbf{x}|y=0)} \right| = \mathbf{x}^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) - \frac{1}{2}(\boldsymbol{\mu}_1 + \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)$$

(see appendix slides
for full derivation)

Linear Discriminant Analysis ($\Sigma_0 = \Sigma_1$)

$$\ln \left| \frac{P(\mathbf{x}|y=1)}{P(\mathbf{x}|y=0)} \right| = \mathbf{x}^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) - \frac{1}{2} (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_0) \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)$$

Since our decision rule is to classify as class 1 if the following is true:

$$\ln \left| \frac{P(\mathbf{x}|y=1)}{P(\mathbf{x}|y=0)} \right| > \ln \left| \frac{P(y=0)}{P(y=1)} \right|$$

We can rewrite our decision rule as:

$$\mathbf{x}^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) > \frac{1}{2} (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_0) \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) + \ln \left| \frac{P(y=0)}{P(y=1)} \right|$$

Or simply as:

$$\mathbf{x}^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) > \lambda_{thresh}$$

If we define $\mathbf{w} = \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)$, then this becomes $\mathbf{x}^T \mathbf{w} > \lambda_{thresh}$

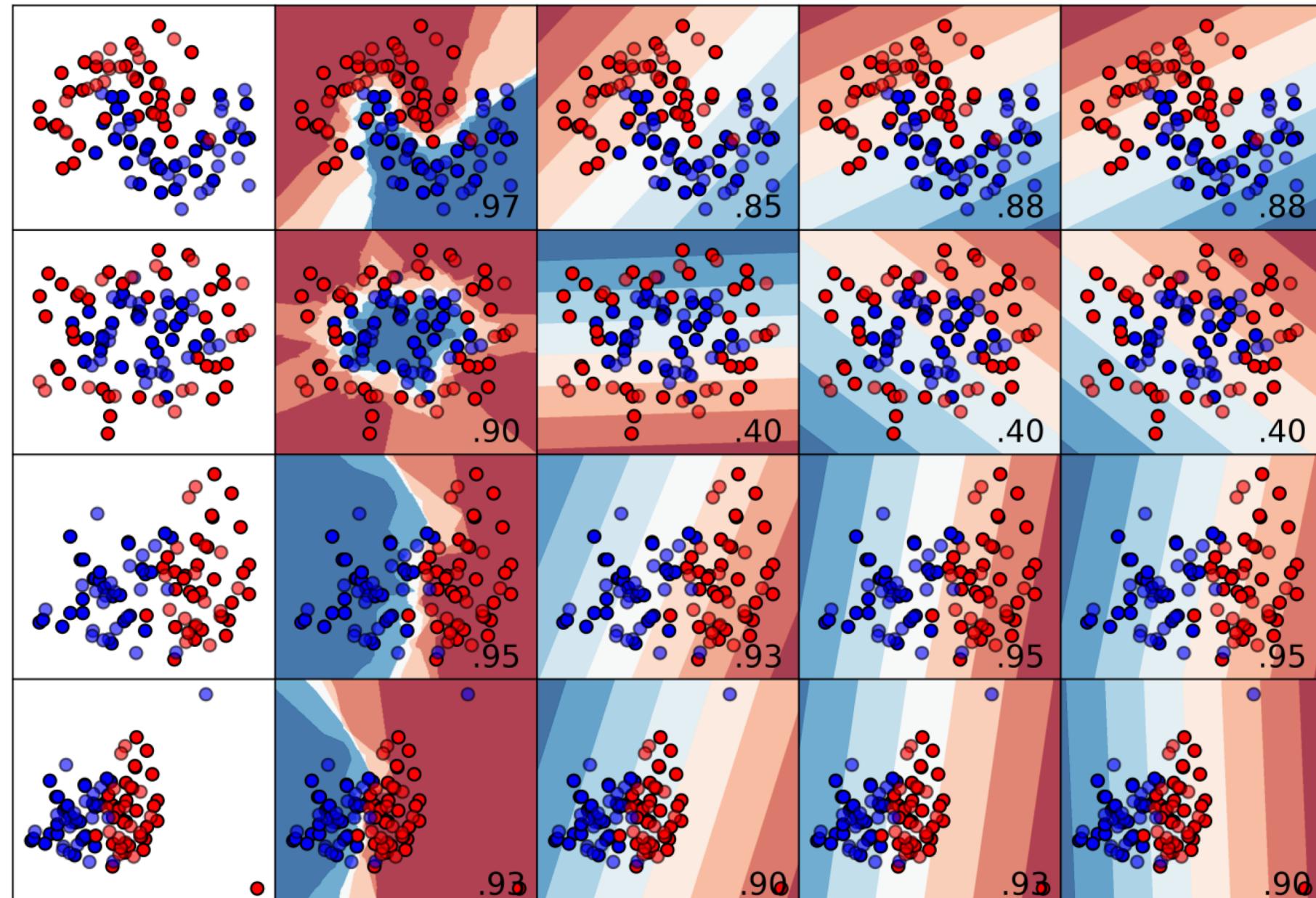
Input data

KNN (k=5)

Perceptron

Logistic Reg.

LDA



Quadratic Discriminant Analysis ($\Sigma_0 \neq \Sigma_1$)

We build a classifier that assigns the class with the higher posterior probability:

$$\frac{P(\mathbf{x}|y=1)}{P(\mathbf{x}|y=0)} = \frac{\frac{1}{(2\pi)^{d/2}|\boldsymbol{\Sigma}_1|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}_1^{-1} (\mathbf{x} - \boldsymbol{\mu}_1)\right]}{\frac{1}{(2\pi)^{d/2}|\boldsymbol{\Sigma}_0|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}_0^{-1} (\mathbf{x} - \boldsymbol{\mu}_0)\right]}$$

We assume a normal distribution, but different covariance matrices

Produces a quadric decision boundary

Input data

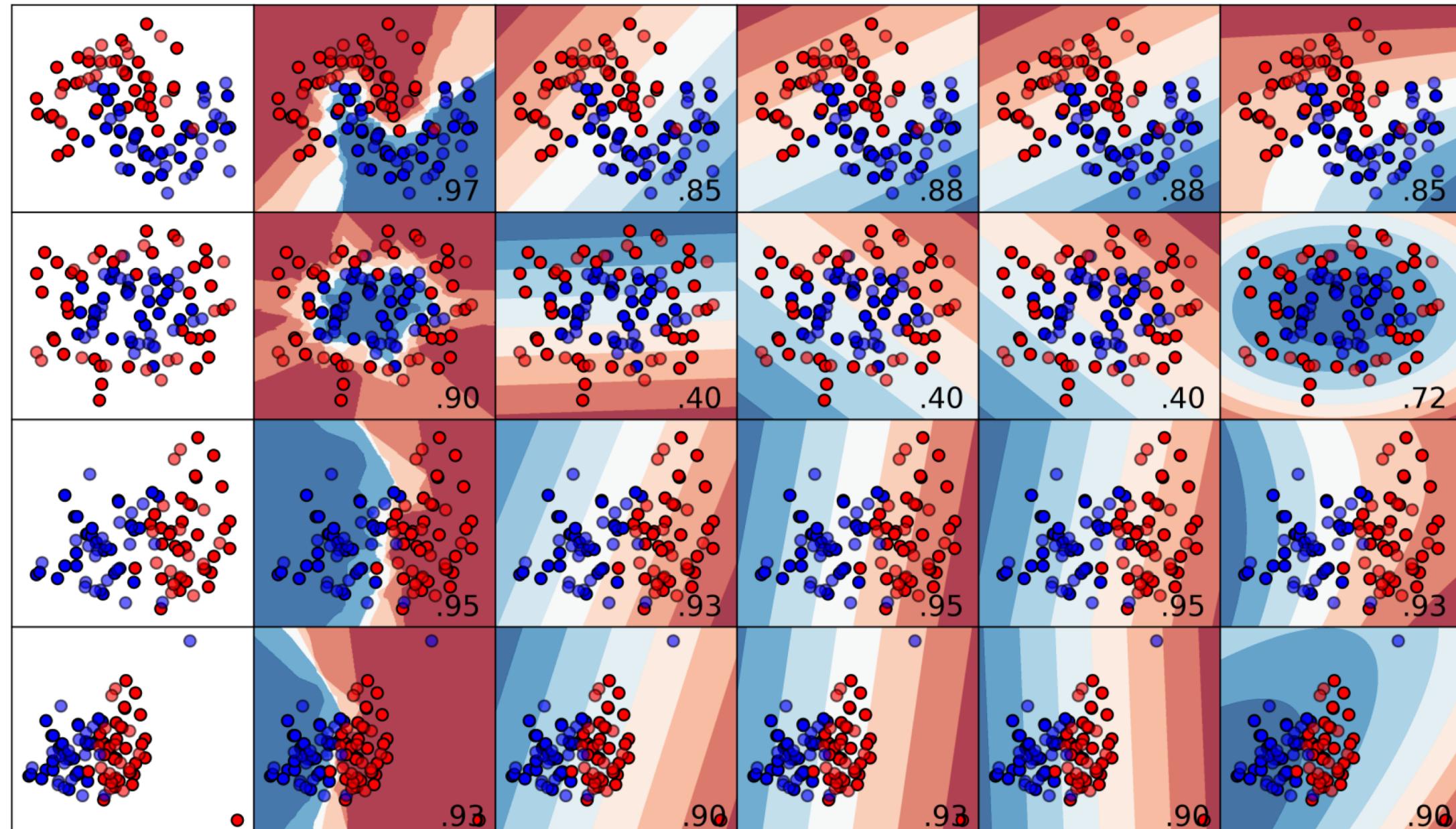
KNN (k=5)

Perceptron

Logistic Reg.

LDA

QDA



Naïve Bayes

Sometimes called “Idiot’s Bayes”

Start with our original expression for our posterior [recall this was our discriminant function, $d_i(\mathbf{x})$]

Write out the full expression with all the terms in \mathbf{x}

$$P(y = i|\mathbf{x}) = \frac{P(\mathbf{x}|y = i)P(y = i)}{P(\mathbf{x})}$$

$$P(y = i|x_1, x_2, \dots, x_D) = \frac{P(x_1, x_2, \dots, x_D|y = i)P(y = i)}{P(x_1, x_2, \dots, x_D)}$$

Assumption: Given the class, the features are independent

$$P(y = i|x_1, x_2, \dots, x_D) = \frac{P(y = i) \prod_{j=1}^D P(x_j|y = i)}{P(x_1, x_2, \dots, x_D)}$$

$$P(y = i|x_1, x_2, \dots, x_D) \propto P(y = i) \prod_{j=1}^D P(x_j|y = i)$$

Since the denominator is a proportionality constant

Naïve Bayes

We assign the class that has the largest posterior, $P(y = i|x_1, x_2, \dots, x_D)$

$$P(y = i|x_1, x_2, \dots, x_D) \propto P(y = i) \prod_{j=1}^D P(x_j|y = i)$$

This implies we estimate the density of each feature **separately**

This independence assumption is a strong assumption that is rarely valid

Considerably simplifies computation and data needs

Is flexible to allow for different distributional forms (i.e. Gaussian) or nonparametric techniques

Naïve Bayes: Gaussian example

We assign the class that has the largest posterior, $P(y = i|x_1, x_2, \dots, x_D)$

$$P(y = i|x_1, x_2, \dots, x_D) \propto P(y = i) \prod_{j=1}^D P(x_j|y = i)$$

This implies we estimate the density of each feature **separately**

If $P(x_j|y = i)$ is $N(\mu_{ji}, \sigma_{ji}^2)$, so for each class we estimate one mean and variance for each of the D features. We multiply **univariate** distributions together

$$P(y = i|x_1, x_2, \dots, x_D) \propto P(y = i) \prod_{j=1}^D N(\mu_{ji}, \sigma_{ji}^2)$$

If we didn't have the independence assumption we would have to estimate full covariance matrices for each of the c classes requiring cD^2 parameters instead of cD

Input data

KNN (k=5)

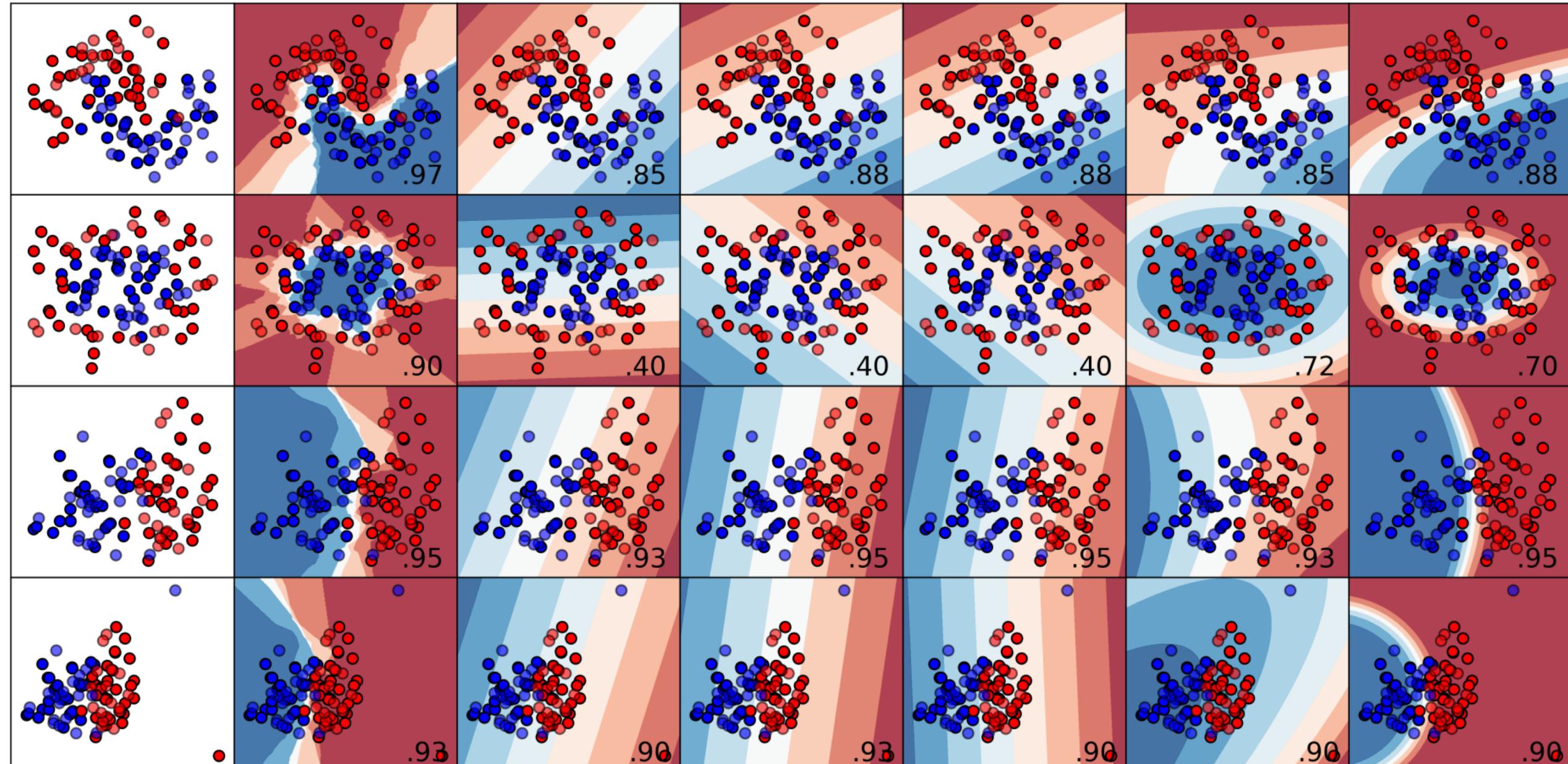
Perceptron

Logistic Reg.

LDA

QDA

Naive Bayes



Classifiers

Covered so far

K-Nearest Neighbors

Linear regression (can be used, but not recommended)

Perceptron

Logistic Regression

Fisher's Linear Discriminant / Linear Discriminant Analysis

Quadratic Discriminant Analysis

Naïve Bayes

Decision Trees

Ensemble methods (random forests, gradient boosting)

Rely on a linear combination of weights and features: $\mathbf{w}^T \mathbf{x}$

Appendix (Derivations)

FLD: Fisher Criterion Maximization

$$J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}$$

$$= \frac{\mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1) (\mathbf{m}_2 - \mathbf{m}_1)^T \mathbf{w}}{\sum_{i \in C_1} (y_n - m_k)^2 + \sum_{i \in C_2} (y_n - m_k)^2}$$

$$= \frac{\mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1) (\mathbf{m}_2 - \mathbf{m}_1)^T \mathbf{w}}{\sum_{i \in C_1} (\mathbf{w}^T \mathbf{x}_i - \mathbf{w}^T \mathbf{m}_1)^2 + \sum_{i \in C_2} (\mathbf{w}^T \mathbf{x}_i - \mathbf{w}^T \mathbf{m}_2)^2}$$

$$= \frac{\mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1) (\mathbf{m}_2 - \mathbf{m}_1)^T \mathbf{w}}{\mathbf{w}^T [\sum_{i \in C_1} (\mathbf{x}_i - \mathbf{m}_1)(\mathbf{x}_i - \mathbf{m}_1)^T + \sum_{i \in C_2} (\mathbf{x}_i - \mathbf{m}_2)(\mathbf{x}_i - \mathbf{m}_2)^T] \mathbf{w}}$$

$$m_2 - m_1 = \mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1)$$

$$s_k^2 = \sum_{i \in C_k} (y_n - m_k)^2$$

$$m_k = \mathbf{w}^T \mathbf{m}_k$$

$$y_k = \mathbf{w}^T \mathbf{x}_k$$

Factoring out the
 \mathbf{w} in denominator

FLD: Fisher Criterion Maximization

$$J(\mathbf{w}) = \frac{\mathbf{w}^T (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T \mathbf{w}}{\mathbf{w}^T [\sum_{i \in C_1} (\mathbf{x}_i - \mathbf{m}_1)(\mathbf{x}_i - \mathbf{m}_1)^T + \sum_{i \in C_2} (\mathbf{x}_i - \mathbf{m}_2)(\mathbf{x}_i - \mathbf{m}_2)^T] \mathbf{w}}$$
$$\mathcal{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T$$
$$\mathcal{S}_W = \sum_{i \in C_1} (\mathbf{x}_i - \mathbf{m}_1)(\mathbf{x}_i - \mathbf{m}_1)^T + \sum_{i \in C_2} (\mathbf{x}_i - \mathbf{m}_2)(\mathbf{x}_i - \mathbf{m}_2)^T$$
$$= \Sigma_1 + \Sigma_2$$

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathcal{S}_B \mathbf{w}}{\mathbf{w}^T \mathcal{S}_W \mathbf{w}}$$

Generalized
Raleigh Quotient

We want to maximize this and solve for \mathbf{w}

FLD: Fisher Criterion Maximization

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

Take the derivative (gradient), set it equal to zero, solve for \mathbf{w}

Recall the quotient rule for differentiation:

$$f(x) = \frac{u(x)}{v(x)} \quad \frac{df}{dx} = \frac{u'v - uv'}{v^2}$$

Matrix derivatives of the form $\mathbf{x}^T \mathbf{A} \mathbf{x}$ with respect to \mathbf{x} are:

$$\frac{d\mathbf{x}^T \mathbf{A} \mathbf{x}}{d\mathbf{x}} = \mathbf{x}^T (\mathbf{A} + \mathbf{A}^T)$$

If \mathbf{A} is symmetric (as it is for our scatter matrices), then $\mathbf{A} = \mathbf{A}^T$, therefore:

$$\mathbf{x}^T (\mathbf{A} + \mathbf{A}^T) = 2\mathbf{x}^T \mathbf{A}$$

Therefore, we can write:

$$\frac{dJ(\mathbf{w})}{d\mathbf{w}} = \frac{(2\mathbf{w}^T \mathbf{S}_B)(\mathbf{w}^T \mathbf{S}_W \mathbf{w}) - (\mathbf{w}^T \mathbf{S}_B \mathbf{w})(2\mathbf{w}^T \mathbf{S}_W)}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})^2} = 0$$

We want to solve this for \mathbf{w}

FLD: Fisher Criterion Maximization

$$\frac{dJ(\mathbf{w})}{d\mathbf{w}} = \frac{(2\mathbf{w}^T \mathbf{S}_B)(\mathbf{w}^T \mathbf{S}_W \mathbf{w}) - (\mathbf{w}^T \mathbf{S}_B \mathbf{w})(2\mathbf{w}^T \mathbf{S}_W)}{(\mathbf{w}^T \mathbf{S}_W \mathbf{w})^2} = 0 \quad \text{We want to solve this for } \mathbf{w}$$

Since the denominator will not approach infinity, only the numerator matters

$$(2\mathbf{w}^T \mathbf{S}_B)(\mathbf{w}^T \mathbf{S}_W \mathbf{w}) - (\mathbf{w}^T \mathbf{S}_B \mathbf{w})(2\mathbf{w}^T \mathbf{S}_W) = 0$$

$$(\underbrace{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}_{\alpha})(\mathbf{w}^T \mathbf{S}_B) = (\underbrace{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}_{\beta})(\mathbf{w}^T \mathbf{S}_W)$$

α β $[1 \times D][D \times D][D \times 1] \rightarrow \text{scalar}$

These will only affect magnitude. We assume that \mathbf{w} is of unit length, so we replace these with variables α and β .

$$\alpha \mathbf{w}^T \mathbf{S}_B = \beta \mathbf{w}^T \mathbf{S}_W$$

FLD: Fisher Criterion Maximization

$$\alpha \mathbf{w}^T \mathbf{S}_B = \beta \mathbf{w}^T \mathbf{S}_W$$

$$\alpha \mathbf{S}_B^T \mathbf{w} = \beta \mathbf{S}_W^T \mathbf{w}$$

$$\alpha \mathbf{S}_B \mathbf{w} = \beta \mathbf{S}_W \mathbf{w}$$

$$\alpha(\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T \mathbf{w} = \beta \mathbf{S}_W \mathbf{w}$$

 scalar $\mathbf{m}_2 - \mathbf{m}_1$, call this γ

$$\alpha\gamma(\mathbf{m}_2 - \mathbf{m}_1) = \beta \mathbf{S}_W \mathbf{w}$$

Property of matrix transposition:

$$(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T$$

The scatter matrices are symmetric:

$$\mathbf{S}_B = \mathbf{S}_B^T$$

$$\mathbf{S}_W = \mathbf{S}_W^T$$

Between-class scatter matrix:

$$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^T$$

Aside: dimensionality reduction

Rearranging, this is an eigenvalue problem

$$\mathbf{S}_W^{-1} \mathbf{S}_B \mathbf{w} = \lambda \mathbf{w}$$

For multiclass problems, we can use the eigenvalue of $\mathbf{S}_W^{-1} \mathbf{S}_B$, much like PCA to get projections into lower dimensional subspaces where the classes are well-separated

FLD: Fisher Criterion Maximization

$$\alpha\gamma(\mathbf{m}_2 - \mathbf{m}_1) = \beta S_W \mathbf{w}$$

Solving for \mathbf{w} :

$$\mathbf{w} = \frac{\alpha\gamma}{\beta} S_W^{-1} (\mathbf{m}_2 - \mathbf{m}_1) \quad \text{We only care about the direction of } \mathbf{w}$$

$$\mathbf{w} \propto S_W^{-1} (\mathbf{m}_2 - \mathbf{m}_1)$$

Note: if S_w is isotropic
(proportional to the identity matrix, i.e. if $S_w = aI$),
then this is just the difference between the means

$$\mathbf{w} \propto (\Sigma_1 + \Sigma_2)^{-1} (\mathbf{m}_2 - \mathbf{m}_1)$$

Linear Discriminant Analysis ($\Sigma_0 = \Sigma_1$)

We build a classifier that assigns the class with the higher posterior probability:

$$\frac{P(\mathbf{x}|y=1)}{P(\mathbf{x}|y=0)} = \frac{\frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_1)\right]}{\frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_0)\right]}$$

$$= \frac{\exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_1)\right]}{\exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_0)\right]}$$

$$\ln \left| \frac{P(\mathbf{x}|y=1)}{P(\mathbf{x}|y=0)} \right| = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_0)$$

Linear Discriminant Analysis ($\Sigma_0 = \Sigma_1$)

$$\ln \left| \frac{P(\mathbf{x}|y=1)}{P(\mathbf{x}|y=0)} \right| = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_0)^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}_0)$$

Expanding this expression yields:

These combine since $\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_i = \boldsymbol{\mu}_i^T \boldsymbol{\Sigma}^{-1} \mathbf{x}$ for symmetric matrices

$$\begin{aligned} &= -\frac{1}{2} [\cancel{\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}} - \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 - \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 \\ &\quad - \cancel{\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}} + \mathbf{x}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_0 + \boldsymbol{\mu}_0^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - \boldsymbol{\mu}_0^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_0] \\ &= \frac{1}{2} [2\mathbf{x}^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) - (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_0) \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)] \\ &= \mathbf{x}^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) - \frac{1}{2} (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_0) \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) \end{aligned}$$

Linear Discriminant Analysis ($\Sigma_0 = \Sigma_1$)

$$\ln \left| \frac{P(\mathbf{x}|y=1)}{P(\mathbf{x}|y=0)} \right| = \mathbf{x}^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) - \frac{1}{2} (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_0) \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)$$

Since our decision rule is to classify as class 1 if the following is true:

$$\ln \left| \frac{P(\mathbf{x}|y=1)}{P(\mathbf{x}|y=0)} \right| > \ln \left| \frac{P(y=0)}{P(y=1)} \right|$$

We can rewrite our decision rule as:

$$\mathbf{x}^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) > \frac{1}{2} (\boldsymbol{\mu}_1 + \boldsymbol{\mu}_0) \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) + \ln \left| \frac{P(y=0)}{P(y=1)} \right|$$

Or simply as:

$$\mathbf{x}^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) > C$$

If we define $\mathbf{w} = \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)$, then this becomes $\mathbf{x}^T \mathbf{w} > C$