

Generative Models for Classification

Lecture 10

Classifiers

Covered so far

K-Nearest Neighbors

Perceptron

Logistic Regression

Linear Discriminant Analysis

Quadratic Discriminant Analysis

Naïve Bayes

Along the way...

Revisiting Bayes' Rule

Projections from higher dimensions

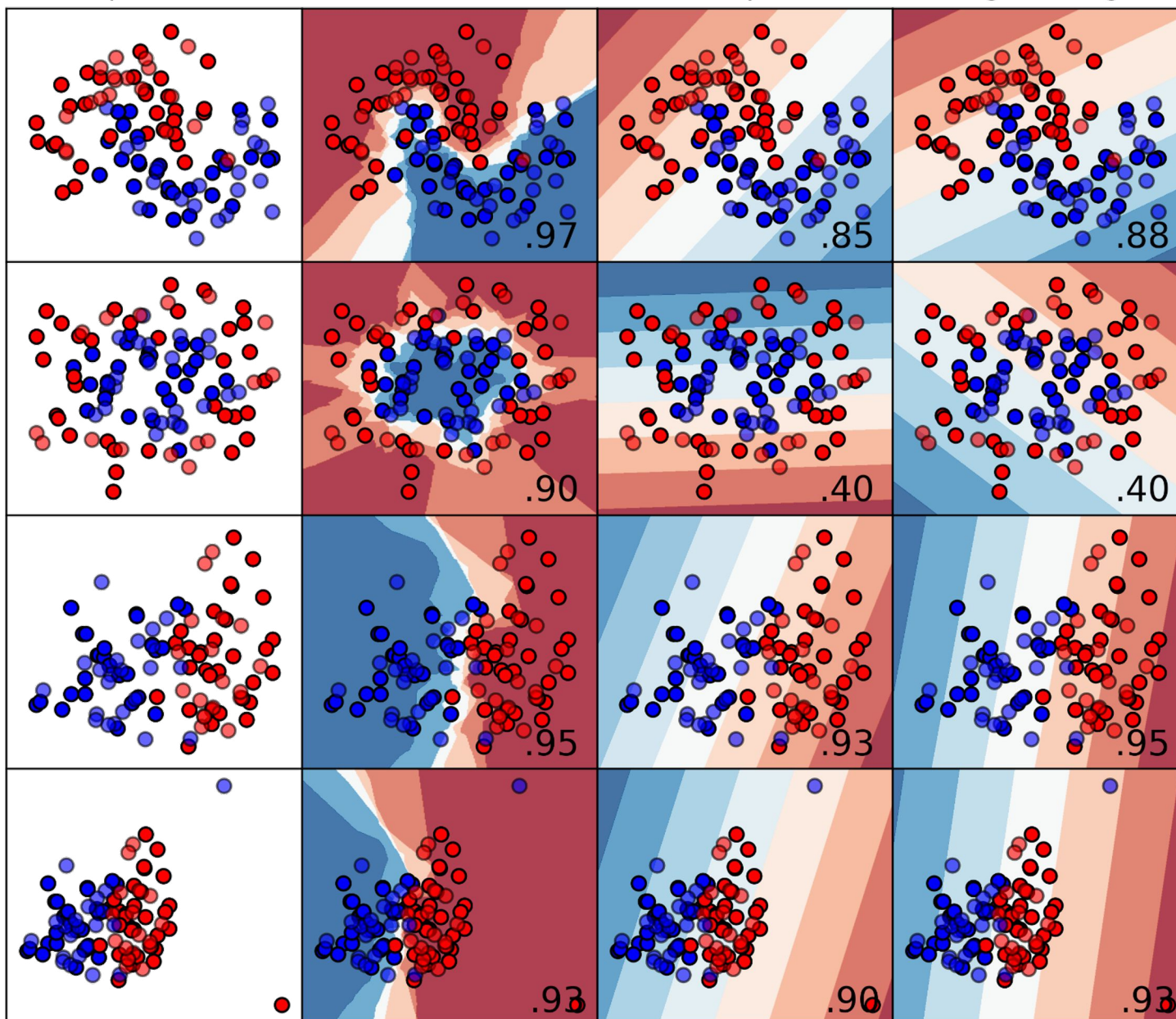
Multivariate normal distributions

Input data

KNN (k=5)

Perceptron

Logistic Reg.



Comparison of classifiers
we've seen so far

Bayes rule in the context of classification



Class 1: Light Post



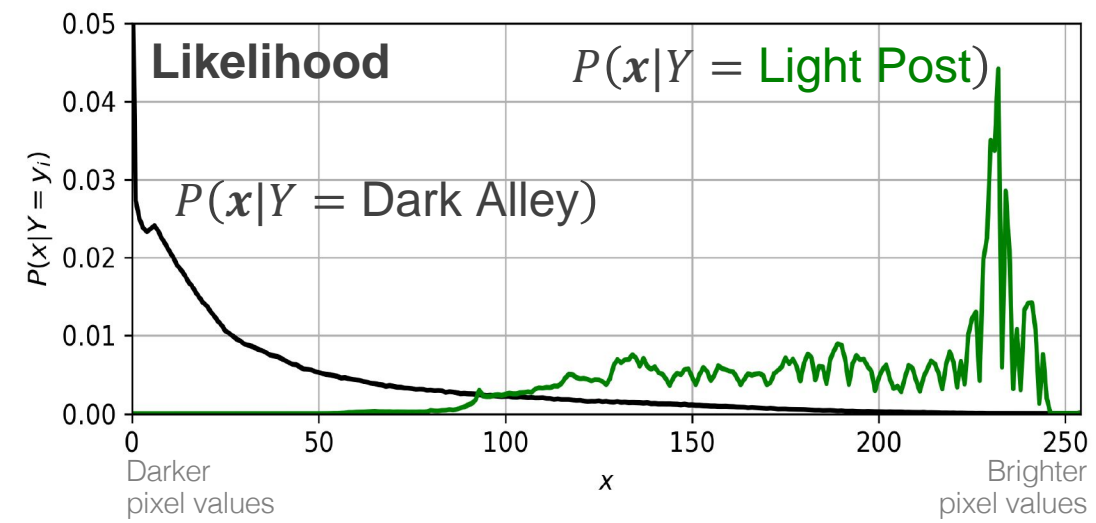
Class 0: Dark Alley

Randomly draw a pixel from either of the images:

$$x_i = 149$$

Darker pixel values are lower numbers (closer to 0), brighter pixels are higher numbers (closer to 255)

How do we determine which image it was most likely to have come from?

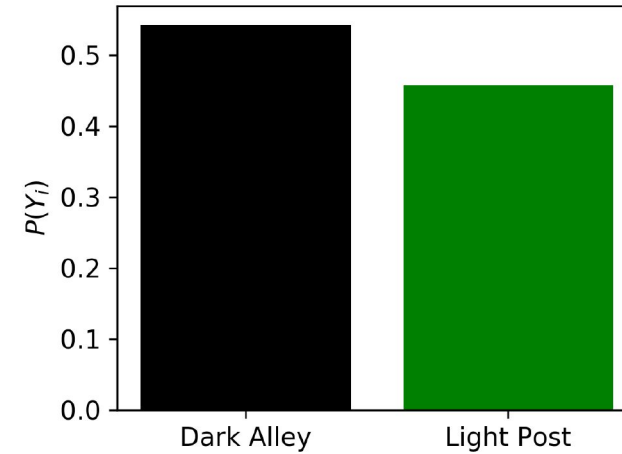


Class 1: Light Post y_1



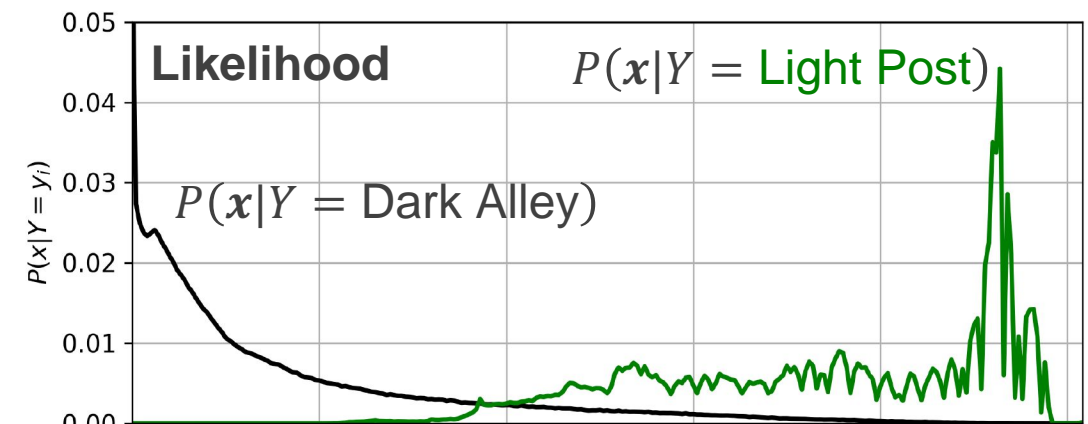
Class 0: Dark Alley y_0

Prior: $P(Y = y_i)$



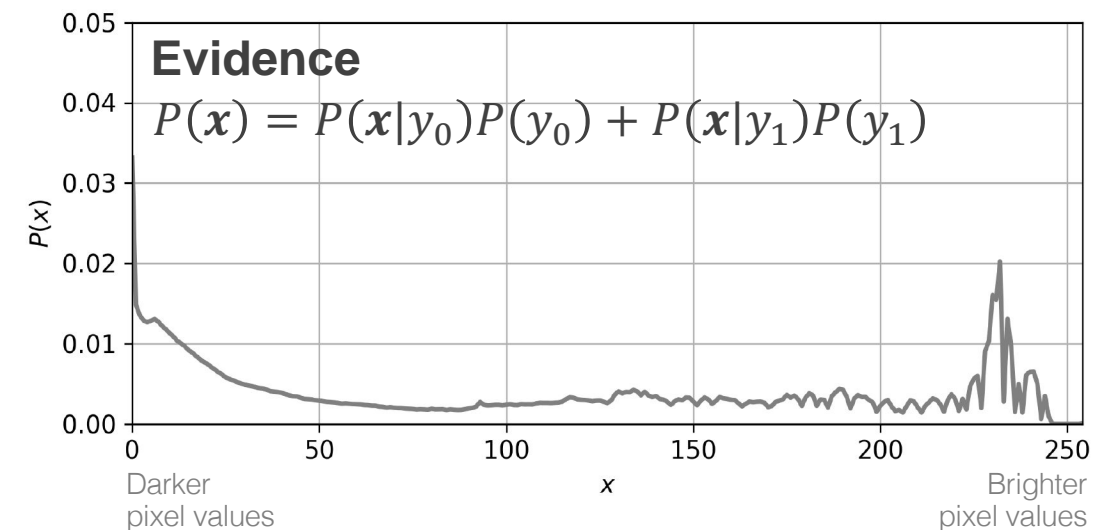
Bayes' Rule

$$\begin{array}{c} \text{Posterior} \\ P(Y = y_i | \mathbf{x}) \end{array} = \frac{\begin{array}{c} \text{Likelihood} \\ P(\mathbf{x} | Y = y_i) \end{array} \begin{array}{c} \text{Prior} \\ P(Y = y_i) \end{array}}{\begin{array}{c} P(\mathbf{x}) \\ \text{Evidence} \end{array}}$$

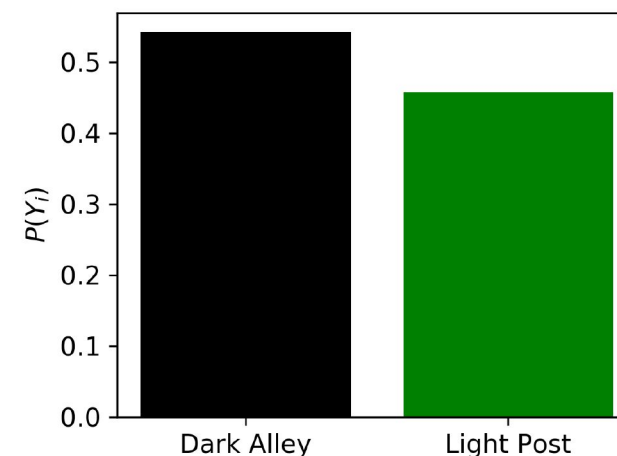


Class 1: Light Post y_1

Class 0: Dark Alley y_0

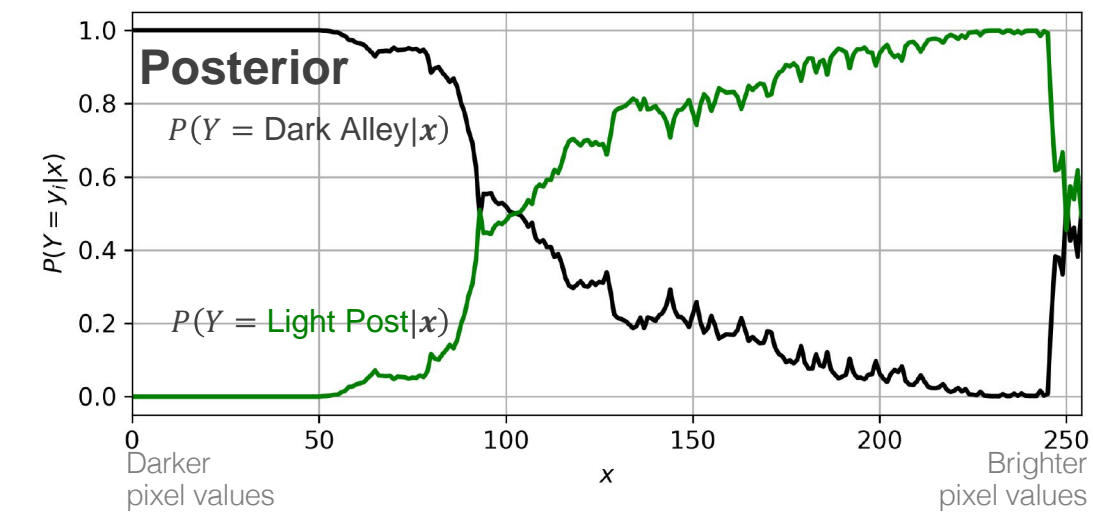
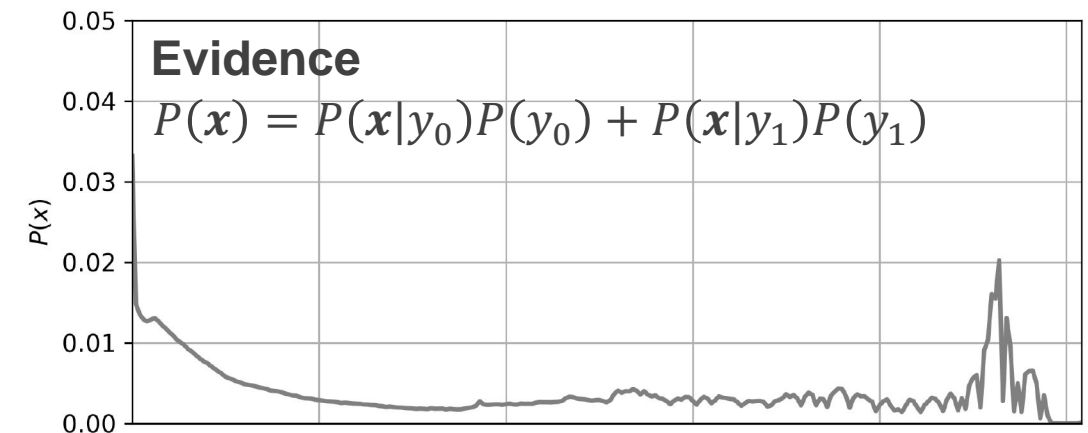
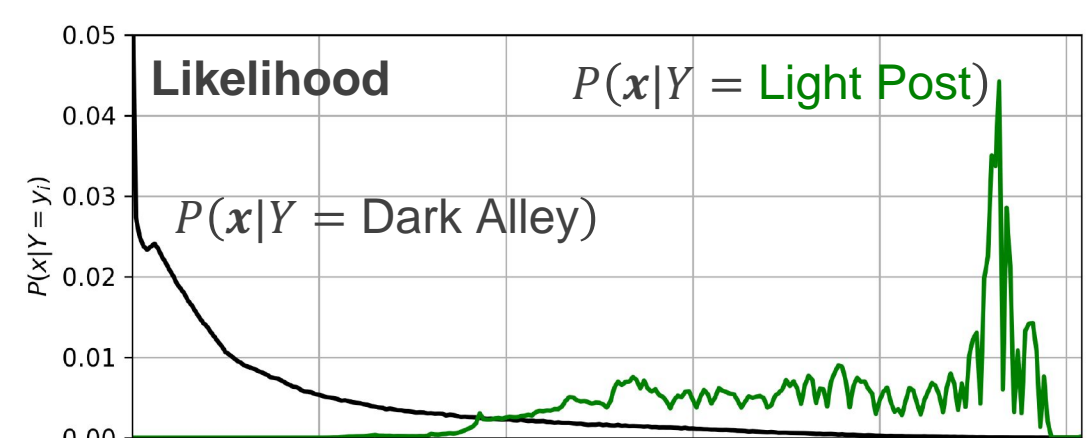


Prior: $P(Y = y_i)$



Bayes' Rule

$$P(Y = y_i | x) = \frac{\text{Likelihood } P(x|Y = y_i) \text{ Prior } P(Y = y_i)}{\text{Evidence } P(x)}$$

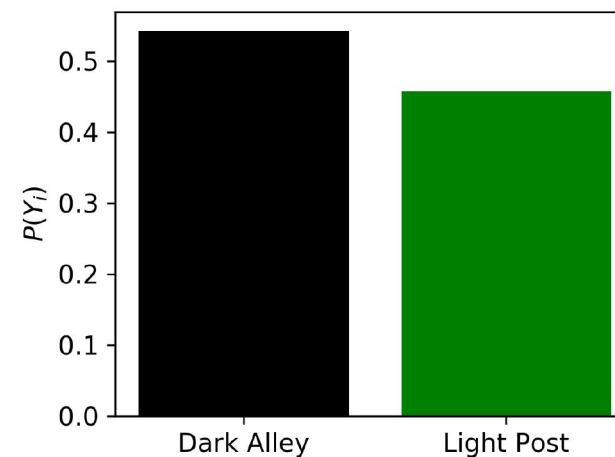


Class 1: Light Post y_1



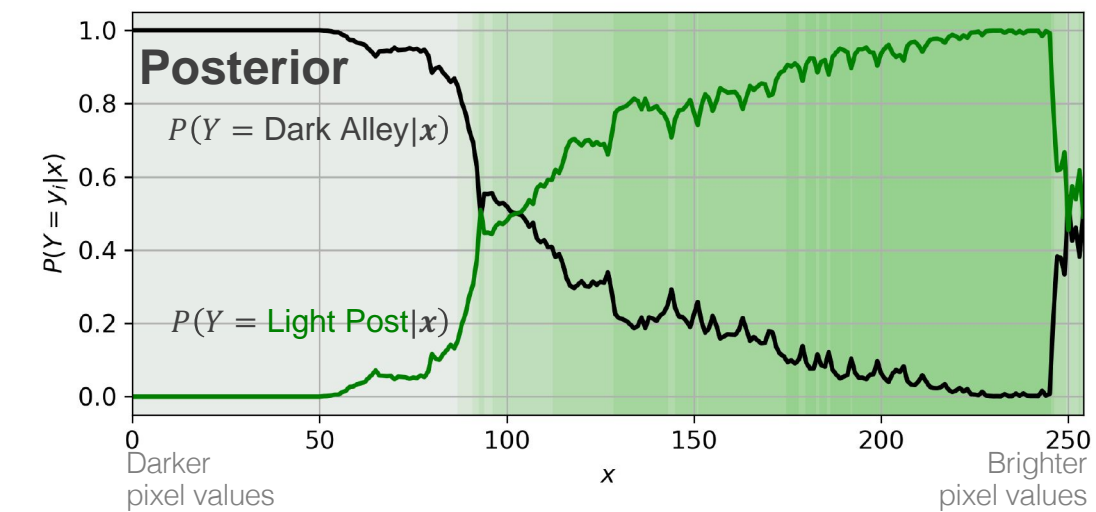
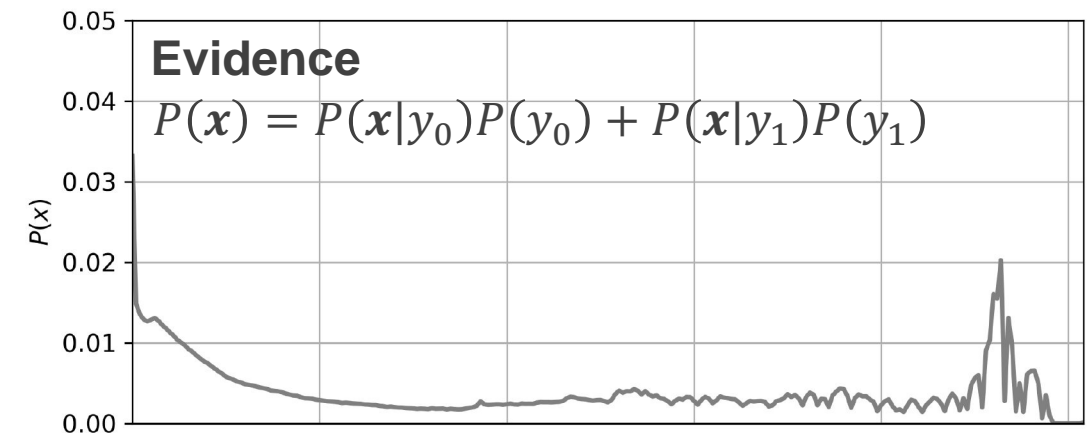
Class 0: Dark Alley y_0

Prior: $P(Y = y_i)$



Bayes' Rule

$$P(Y = y_i|x) = \frac{\text{Likelihood } P(x|Y = y_i) \text{ Prior } P(Y = y_i)}{\text{Evidence } P(x)}$$

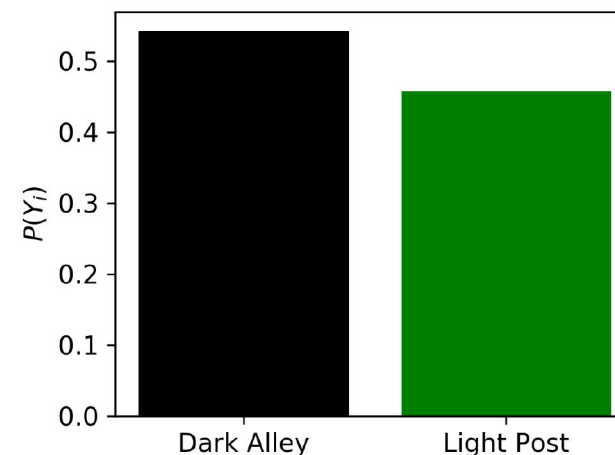


Class 1: Light Post y_1



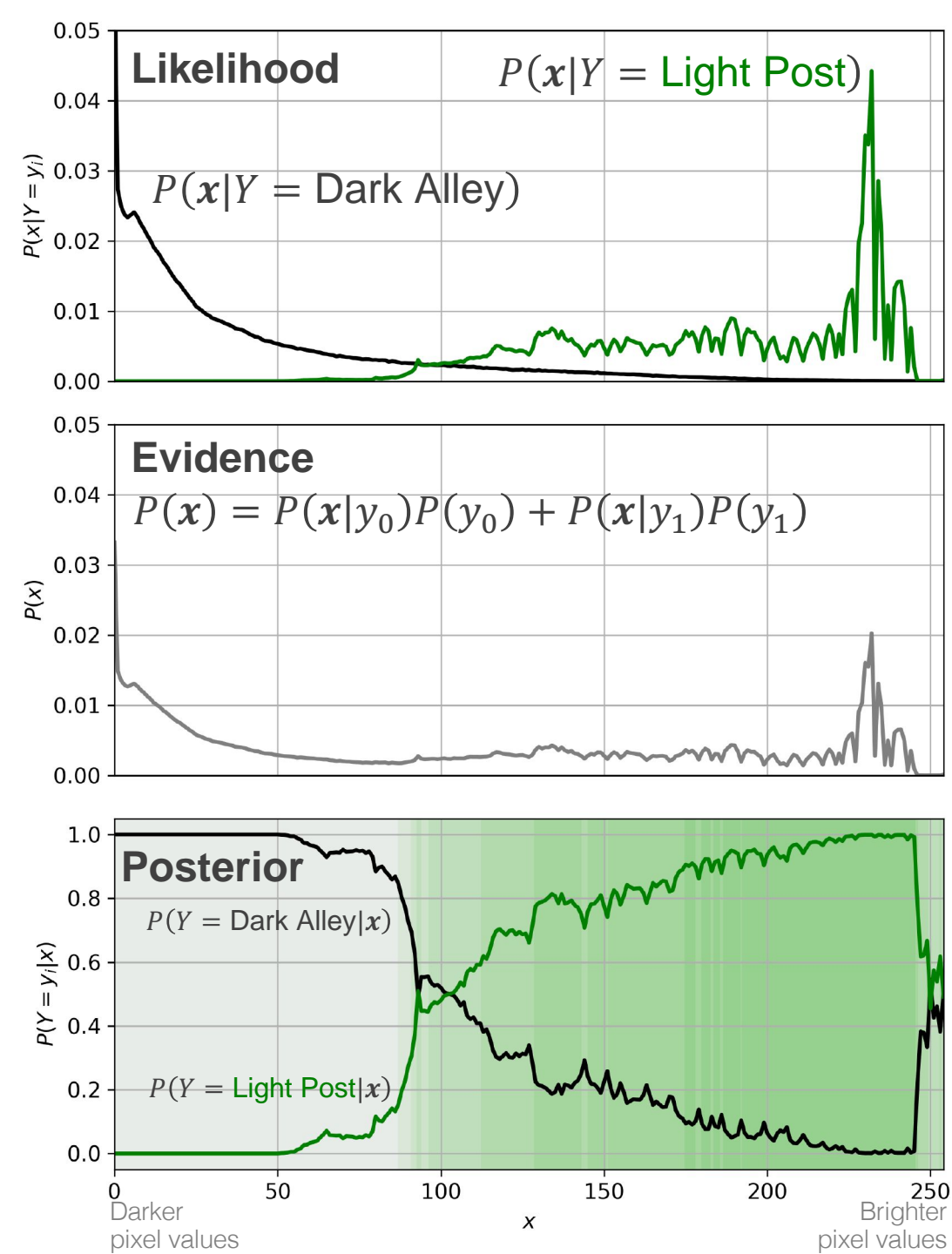
Class 0: Dark Alley y_0

Prior: $P(Y = y_i)$



Decision rule:

If $P(Y = \text{Light Post}|x) > P(Y = \text{Dark Alley}|x)$ then **Light Post**
 else **Dark Alley**



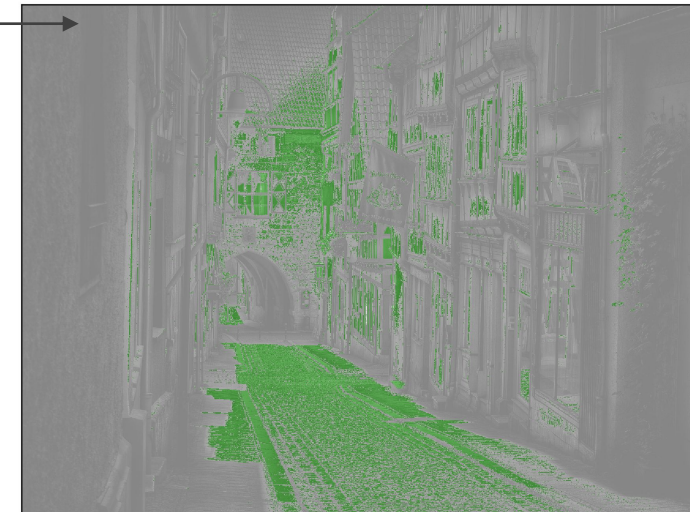
Class 1: Light Post y_1



Green = classified as from Light Post
 Grey = classified as from Dark Alley



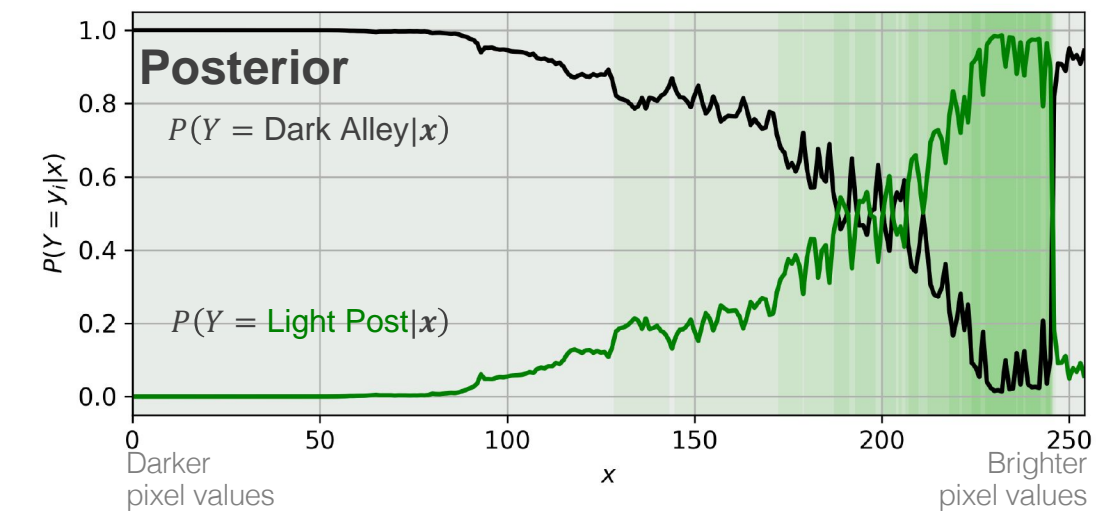
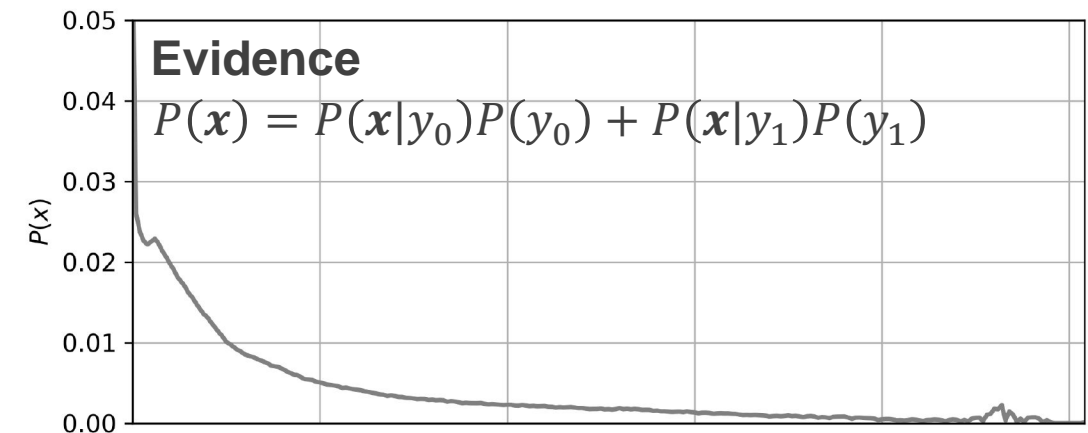
Class 0: Dark Alley y_0



Classifying each of the individual pixels as being either from **Light Post** or **Dark Alley** results in classification above

Decision rule:

If $P(Y = \text{Light Post}|x) > P(Y = \text{Dark Alley}|x)$ then **Light Post**
 else **Dark Alley**

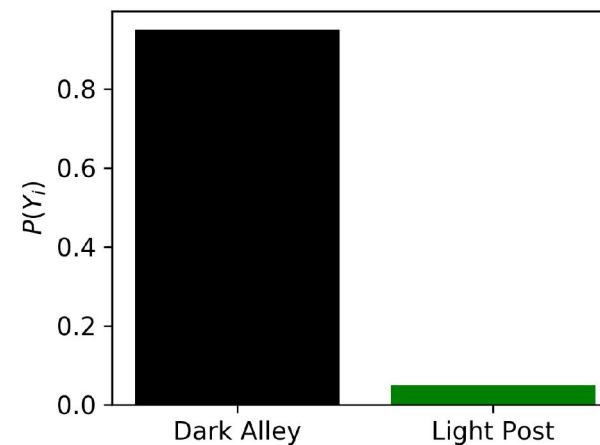


Class 1: Light Post y_1



Class 0: Dark Alley y_0

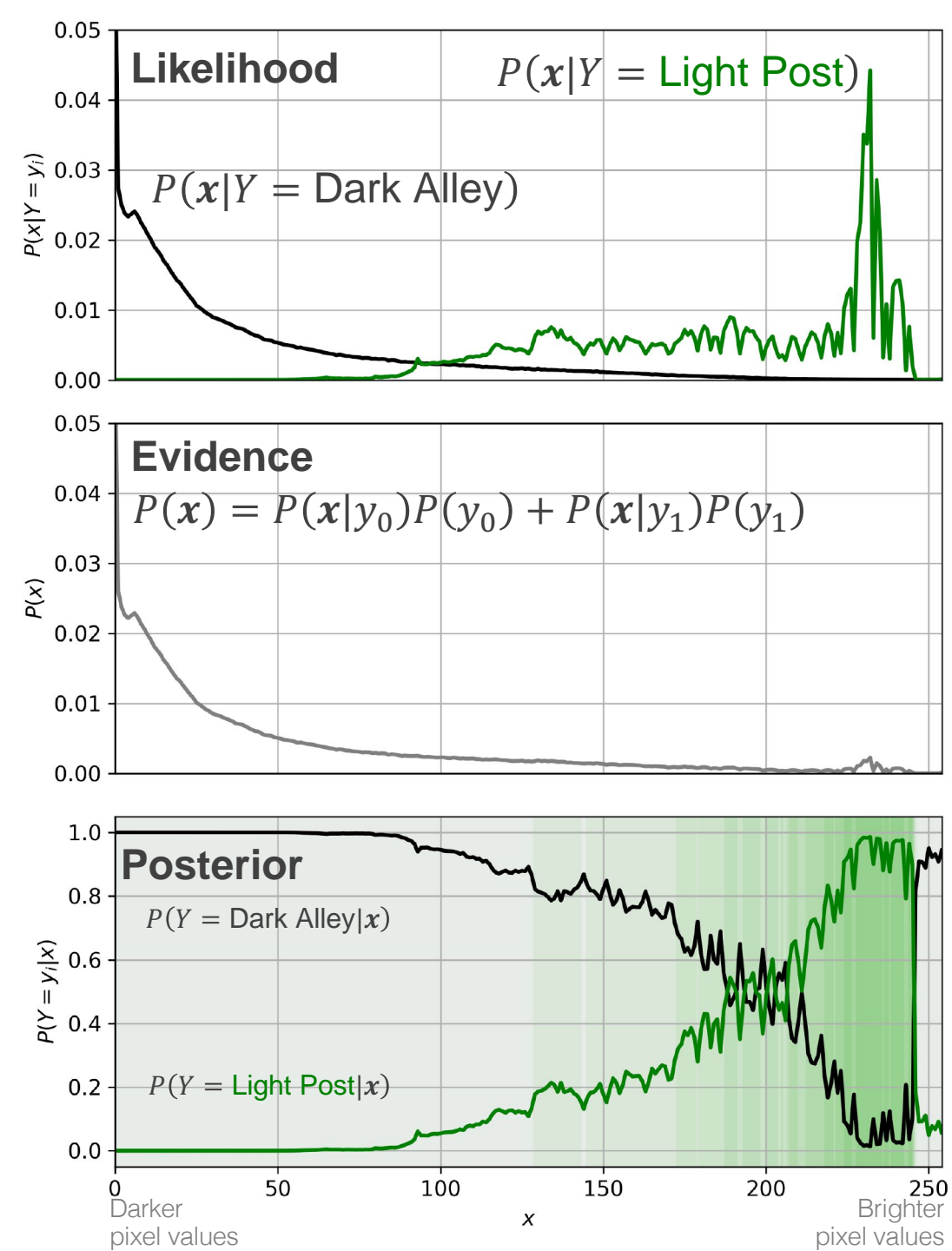
Prior: $P(Y = y_i)$



Let's assume the sampling of pixels occurred more from the **Dark Alley**

Bayes' Rule

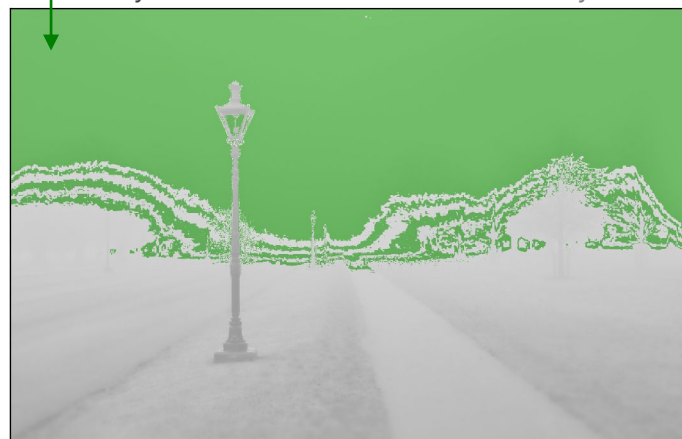
$$P(Y = y_i|x) = \frac{\text{Likelihood} \quad \text{Prior}}{\text{Evidence}} = \frac{P(x|Y = y_i)P(Y = y_i)}{P(x)}$$



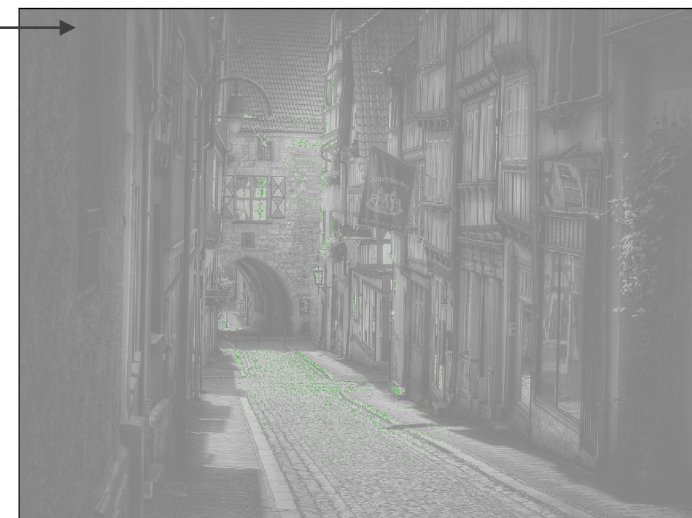
Class 1: Light Post y_1



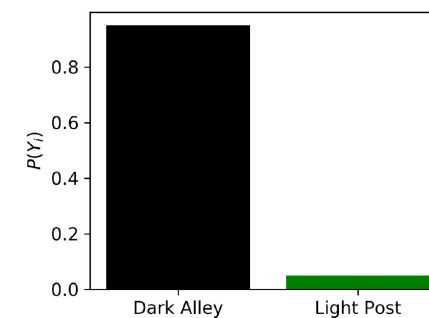
Green = classified as from Light Post
Grey = classified as from Dark Alley



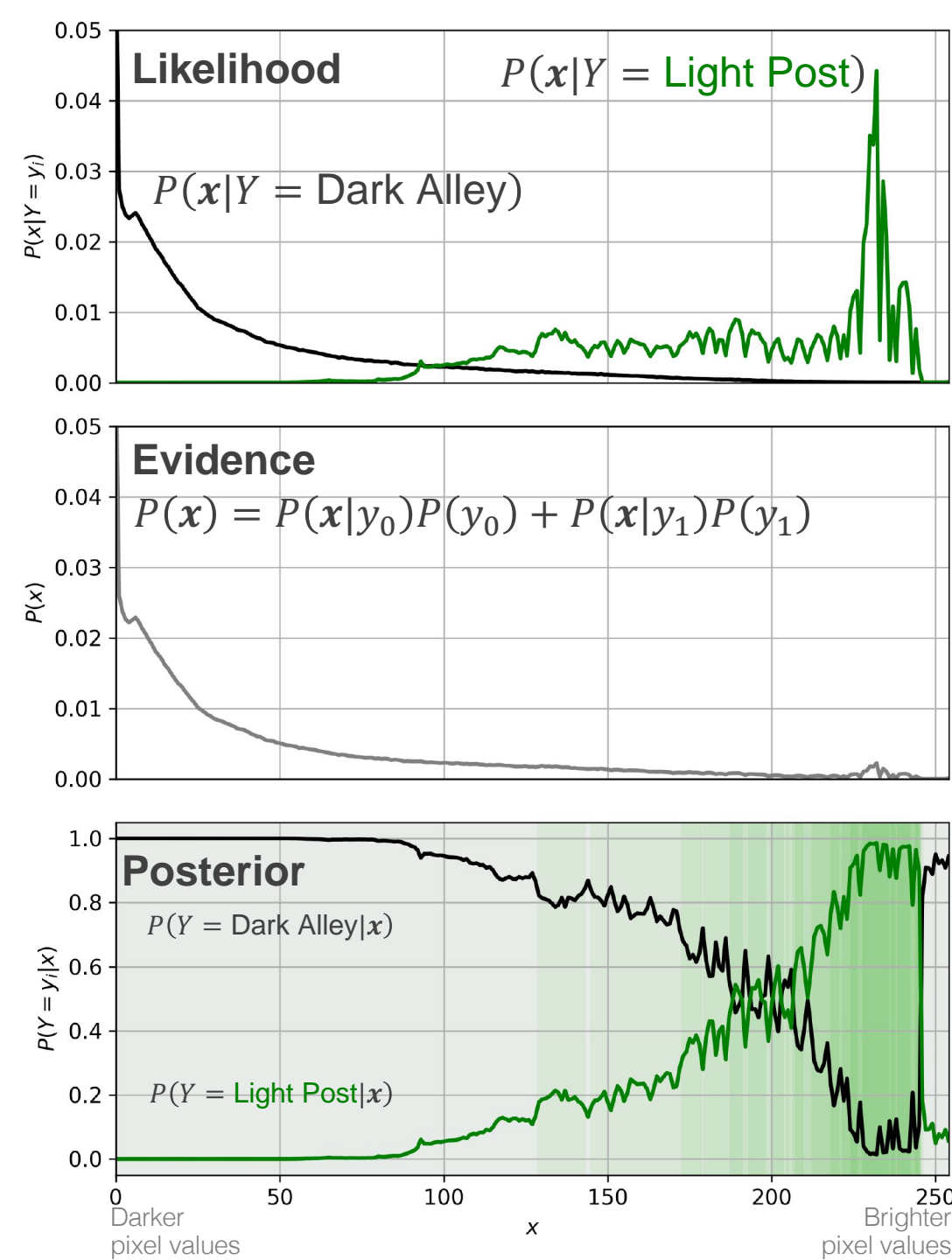
Class 0: Dark Alley y_0



Prior: $P(Y = y_i)$



Assuming we the sampling of pixels occurred more from the **Dark Alley**



Generative models model the **likelihood**

- These can also be used to generate synthetic data
- Often good performance when sample size is small

Examples: linear discriminant analysis, naïve Bayes, hidden Markov models, Gaussian mixture models, Generative Adversarial Networks

$$\begin{array}{c} \text{Posterior} \\ P(Y = y_i | x) \end{array} = \frac{\begin{array}{c} \text{Likelihood} \\ P(x|Y = y_i) \end{array} \begin{array}{c} \text{Prior} \\ P(Y = y_i) \end{array}}{\begin{array}{c} \text{Evidence} \\ P(x) \end{array}}$$

Discriminative models model the **posterior**

- Or they just directly estimate labels without a probabilistic interpretation, $f(x) \rightarrow y$
- Often better performance for large sample sizes

Examples: logistic regression, support vector machines, neural networks, k nearest neighbors

Generative modeling for classification

Assume we have c **different classes**

Define the **posterior probability** as a discriminant function: $d_i(x) = P(y = i | x)$
 $i = 1, \dots, c$

For a new sample, classify it as the class with the **largest** $d_i(x)$

Generative modeling for classification

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 classify feature \mathbf{x} to class i

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

$$= \frac{P(\mathbf{x}|y = i)P(y = i)}{\sum_{i=1}^c P(\mathbf{x}|y = i)P(y = i)} \rightarrow$$

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

Denominator is the same for all classes i , so it won't help us tell which class's posterior is higher relative to other classes, so we ignore it going forward

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

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

If we know the **true likelihood and prior** for our data, this process yields our **Bayes' classifier** (minimum misclassification error classifier)

Generative modeling for classification

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

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

We **rarely** know our **true likelihood** for our data so we need to make assumptions

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

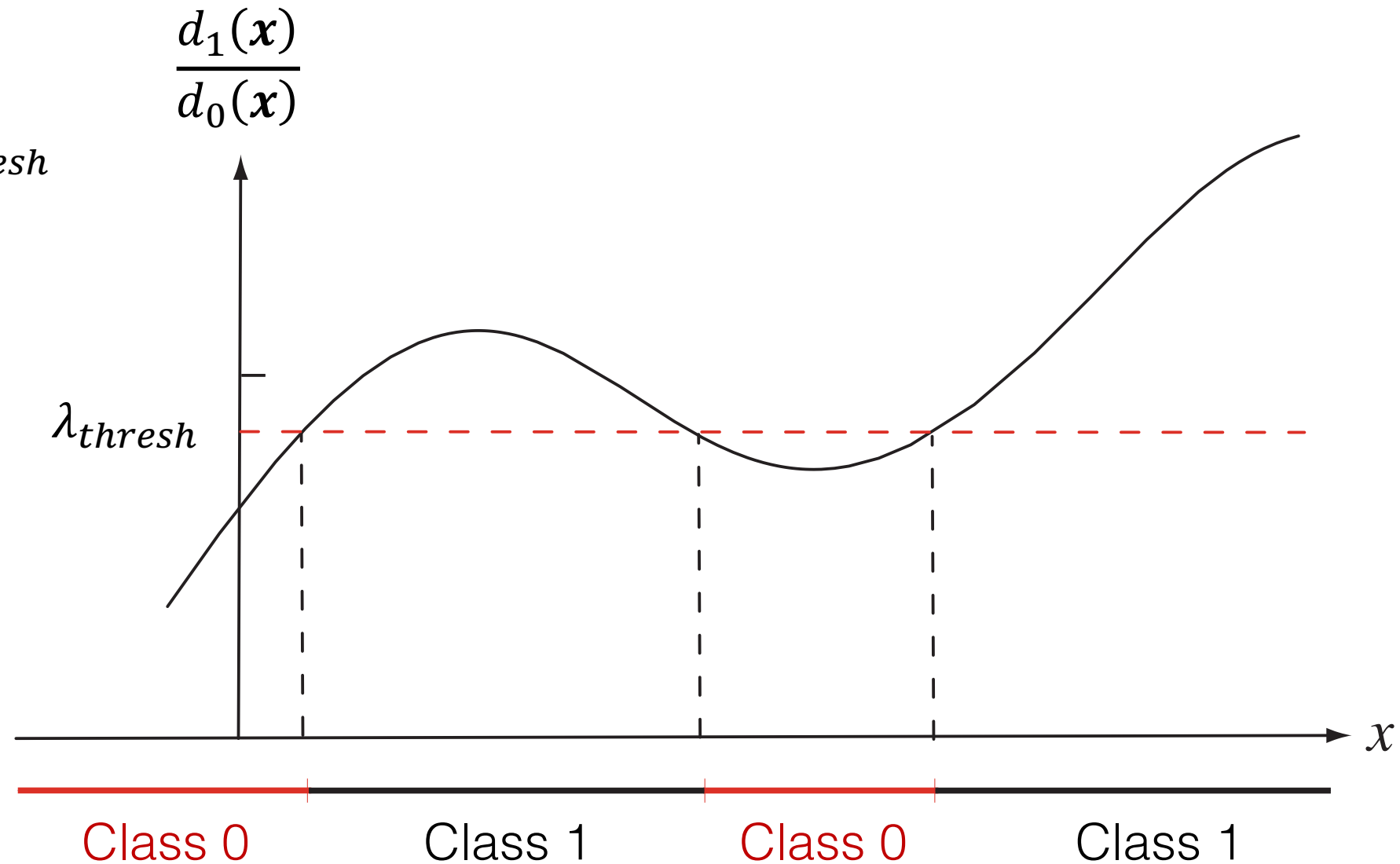
Applies to both binary and multiclass problems

Discriminant Functions: 2 classes

Decision rule:

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

Otherwise, class 0



Discriminant Function: 2 classes

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

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

Assumes these likelihoods are normal

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

The diagram illustrates the relationship between the likelihood ratio and the normal distribution parameters. It features a central inequality: $\frac{P(\mathbf{x}|y = 1)}{P(\mathbf{x}|y = 0)} > \frac{P(y = 0)}{P(y = 1)}$. From the left side of this inequality, a green line extends upwards and then turns right as an arrow pointing to the expression $N(\mu_1, \Sigma_1)$. From the right side of the inequality, a green line extends downwards and then turns right as an arrow pointing to the expression $N(\mu_0, \Sigma_0)$. Both $N(\mu_1, \Sigma_1)$ and $N(\mu_0, \Sigma_0)$ have green arrows pointing towards a text block on the right that reads: "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:

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

$N(\mu_1, \Sigma_1)$

$N(\mu_0, \Sigma_0)$

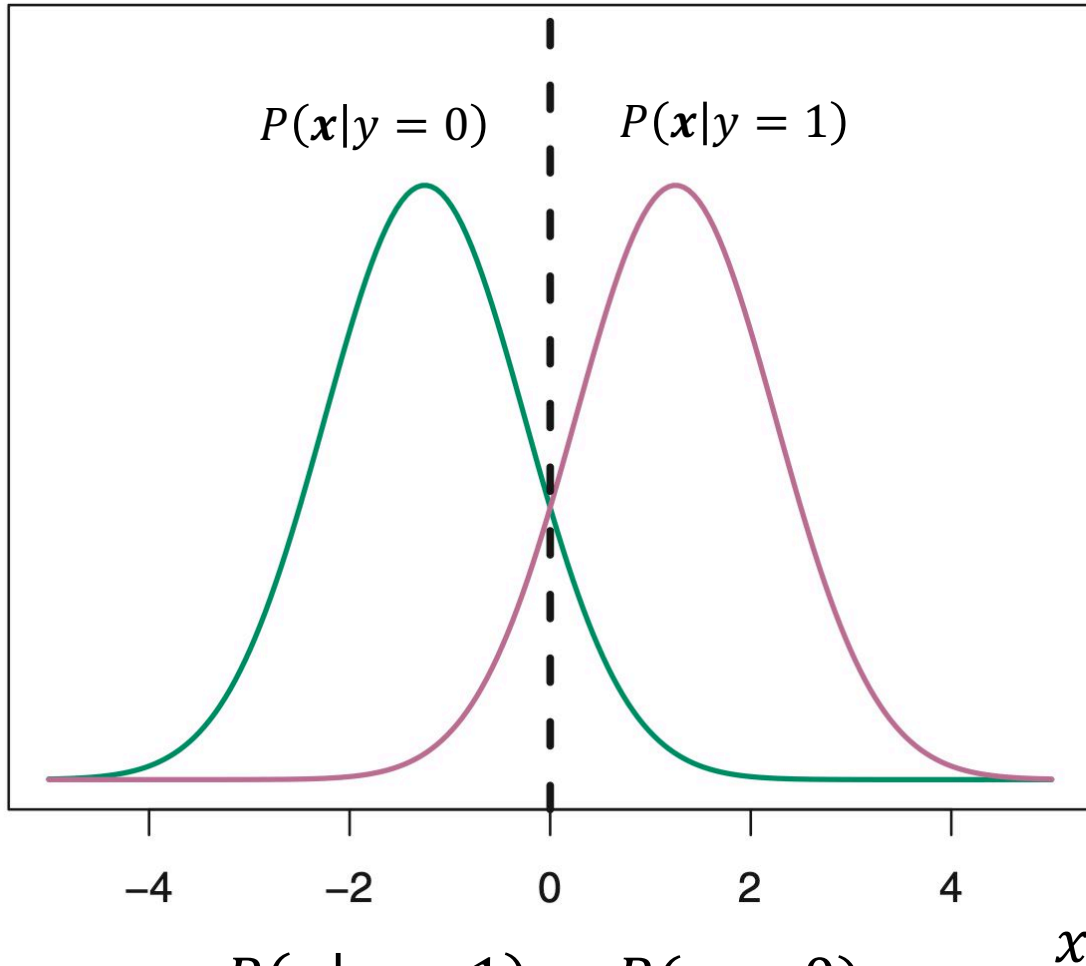
If we assume the class conditional distributions are Gaussian, this represents

Quadratic Discriminant Analysis

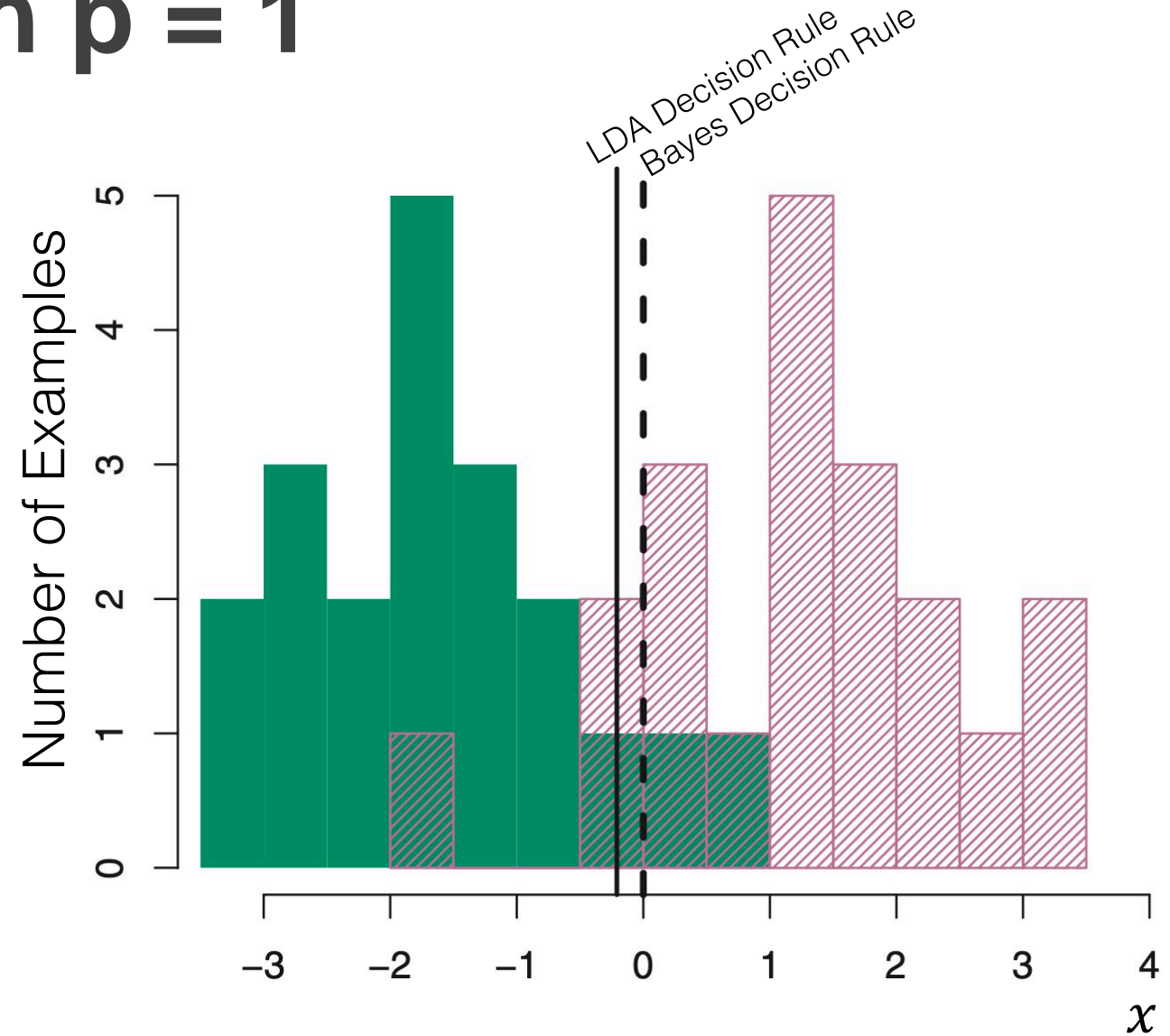
If we further assume the covariance matrices are the same, $\Sigma_0 = \Sigma_1$, this represents

Linear Discriminant Analysis

Simple example when $p = 1$



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



Multivariate Gaussian

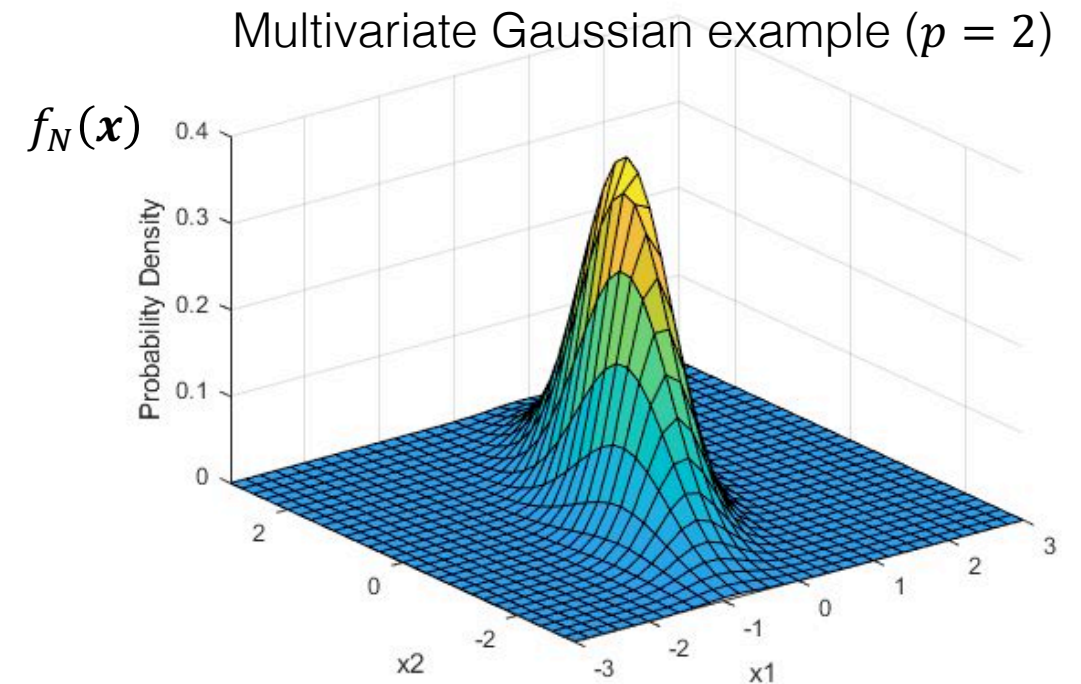
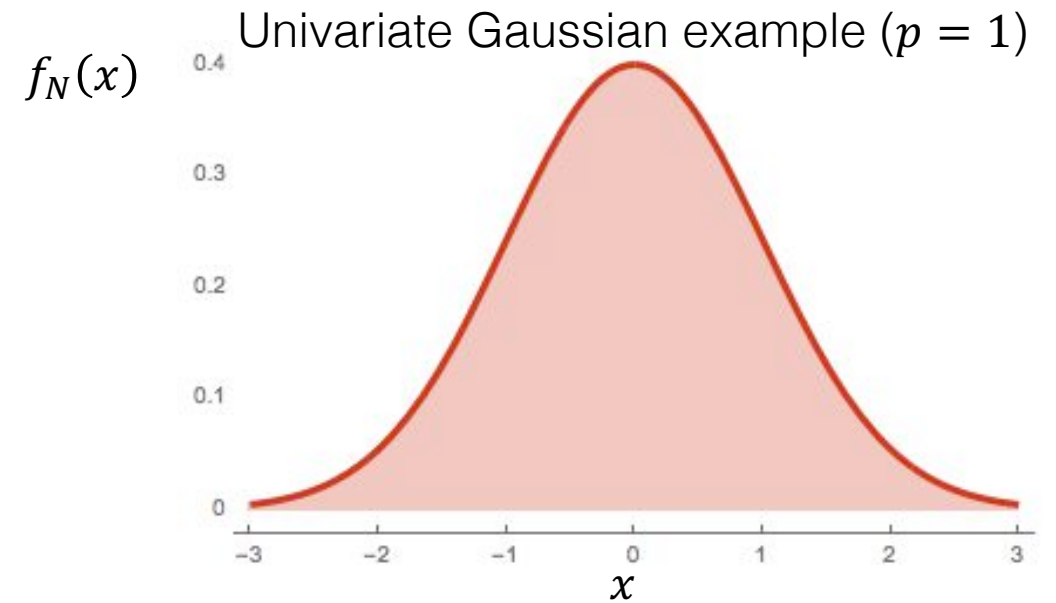
Univariate Gaussian (1 predictor)

$$f_N(x) = \frac{1}{(2\pi)^{1/2}\sigma} \exp\left[-\frac{1}{2\sigma^2}(x - \mu)^2\right]$$

Multivariate Gaussian (p predictors)

$$f_N(\mathbf{x}) = \frac{1}{(2\pi)^{p/2}|\mathbf{\Sigma}|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right]$$

Estimate our means and variances and we can use this as our likelihood



Images from Brilliant.org (top) and the Mathworks (bottom)

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

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

$$\begin{aligned}\frac{P(\mathbf{x}|y = 1)}{P(\mathbf{x}|y = 0)} &= \frac{\frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \mu_1)^T \Sigma^{-1}(\mathbf{x} - \mu_1)\right]}{\frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \mu_0)^T \Sigma^{-1}(\mathbf{x} - \mu_0)\right]} \\ &= \frac{\exp\left[-\frac{1}{2}(\mathbf{x} - \mu_1)^T \Sigma^{-1}(\mathbf{x} - \mu_1)\right]}{\exp\left[-\frac{1}{2}(\mathbf{x} - \mu_0)^T \Sigma^{-1}(\mathbf{x} - \mu_0)\right]}\end{aligned}$$

$$\ln \left| \frac{P(\mathbf{x}|y = 1)}{P(\mathbf{x}|y = 0)} \right| = -\frac{1}{2}(\mathbf{x} - \mu_1)^T \Sigma^{-1}(\mathbf{x} - \mu_1) + \frac{1}{2}(\mathbf{x} - \mu_0)^T \Sigma^{-1}(\mathbf{x} - \mu_0)$$

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

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

Linear Discriminant Analysis: 2 Class

Σ is the same for both classes

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

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:

(see appendix slides
for full derivation)

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

Project the data into a
1-dimensional space

Set the threshold for classifying as
class 1 or 0

Linear Discriminant Analysis: 2 Class

Σ is the same for both classes

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

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:

(see appendix slides
for full derivation)

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

Or simply as:

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

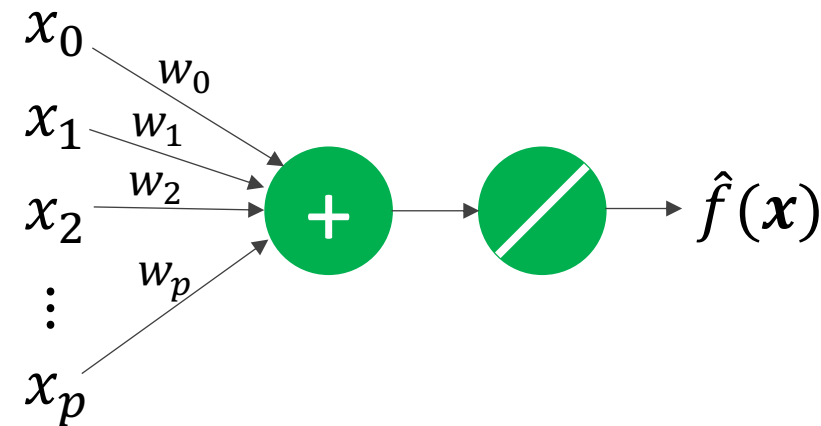
If we define $\mathbf{w} = \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)$, then this becomes $\mathbf{x}^T \mathbf{w} = \mathbf{w}^T \mathbf{x} > \lambda_{thresh}$

Remember linear models?

Linear Regression

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

$\mathbf{w}^T \mathbf{x}$

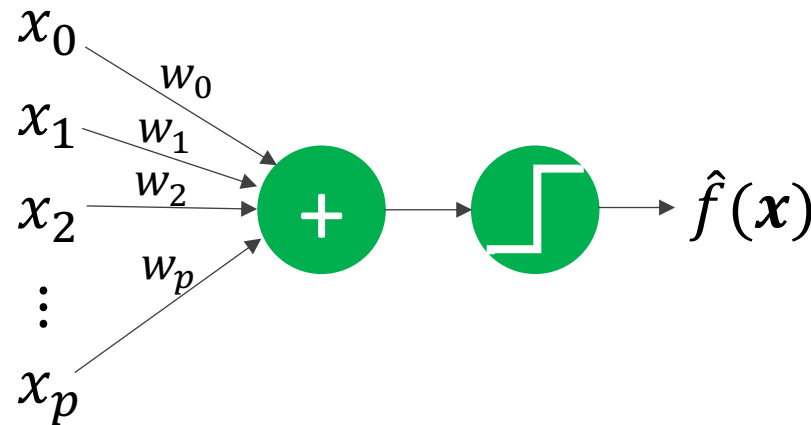


Linear Classification

Perceptron

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

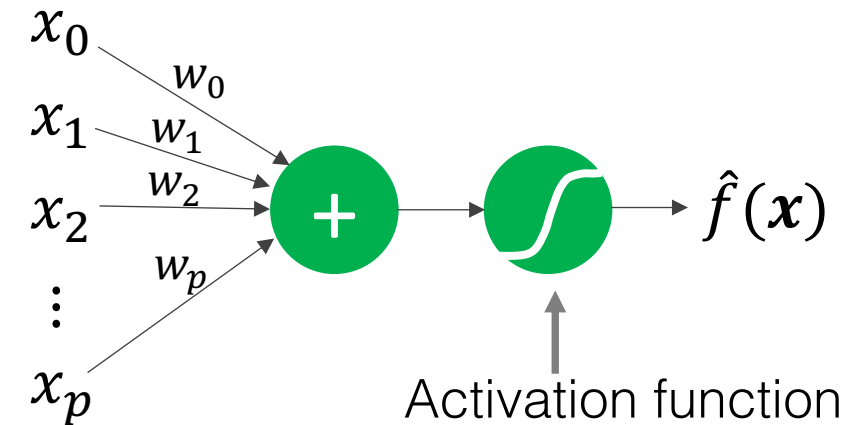
$$\text{sign}(x) = \begin{cases} 1 & x > 0 \\ -1 & \text{else} \end{cases}$$



Logistic Regression

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

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



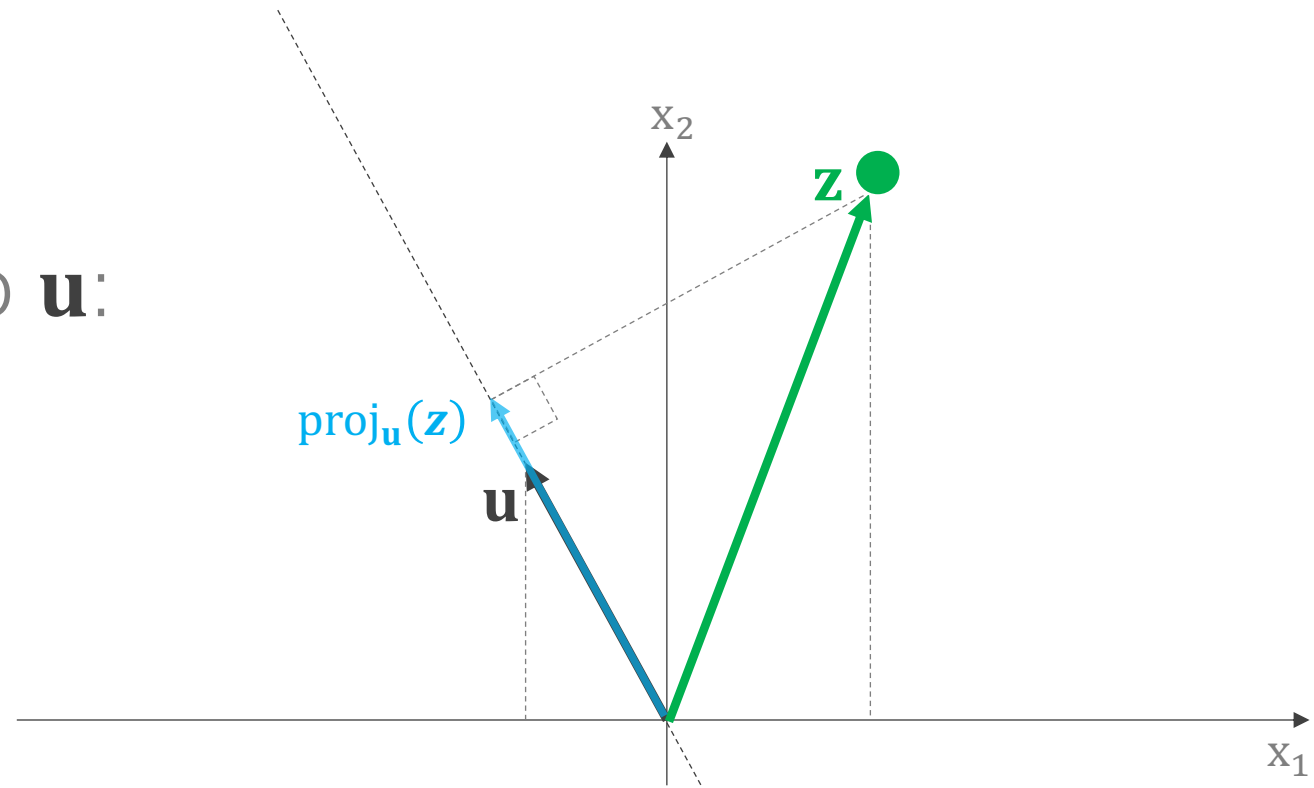
Source: Abu-Mostafa, Learning from Data, Caltech

Intuitively interpreting LDA in terms of projections

Projections

The vector projection of \mathbf{z} onto \mathbf{u} :

$$\text{proj}_{\mathbf{u}}(\mathbf{z}) = \left(\frac{\mathbf{u}^T \mathbf{z}}{\|\mathbf{u}\|^2} \right) \mathbf{u}$$



The scalar length (Euclidean or L_2 norm) of the vector \mathbf{u} is $\|\mathbf{u}\|$

If we assume \mathbf{u} is a unit vector then $\|\mathbf{u}\| = 1$ $\text{proj}_{\mathbf{u}}(\mathbf{z}) = (\mathbf{u}^T \mathbf{z}) \mathbf{u}$

↑
Magnitude of projection onto direction of \mathbf{u}

Projections

$$\mathbf{u}^T \mathbf{z} = [u_1 \quad u_2] \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}$$

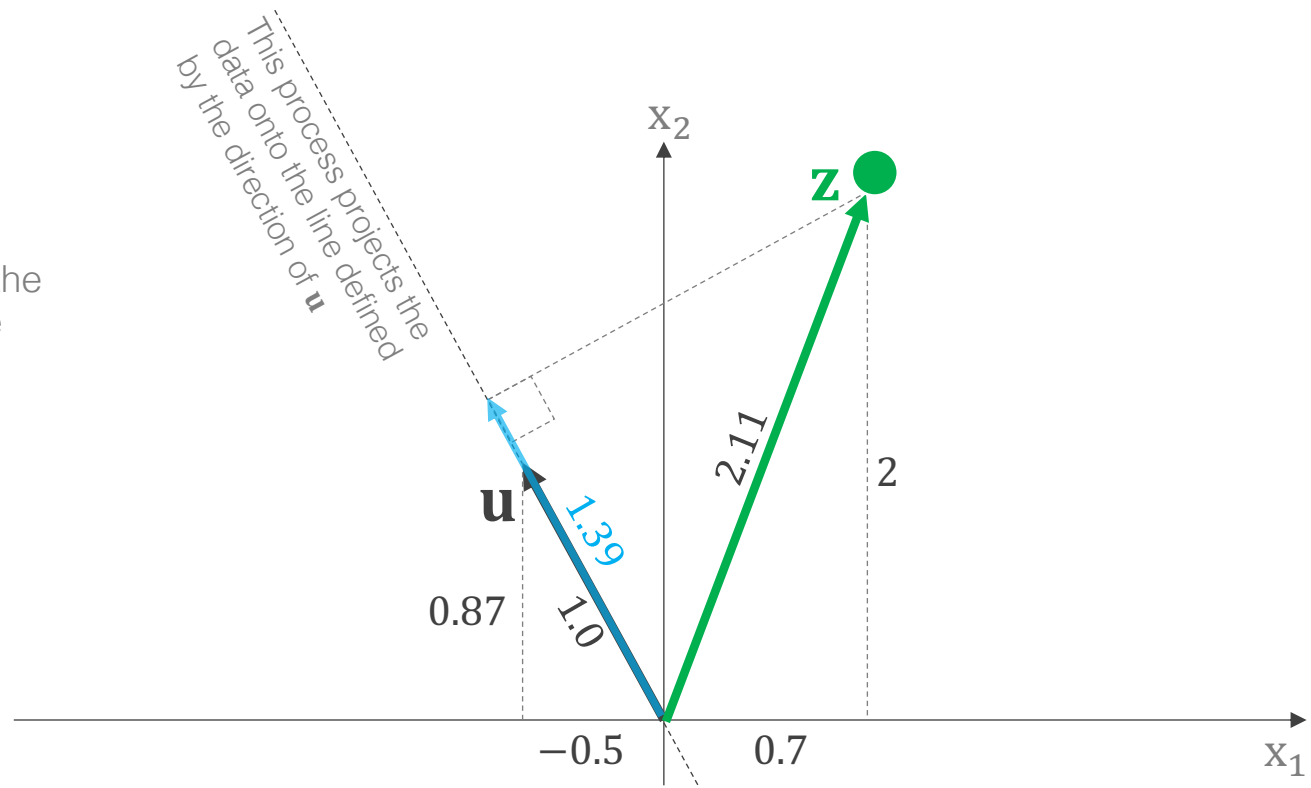
$$= u_1 z_1 + u_2 z_2$$

$$= (-0.5)(0.7) + (0.87)(2)$$

$$= 1.39$$

Length (magnitude) of the projection of \mathbf{z} onto \mathbf{u}

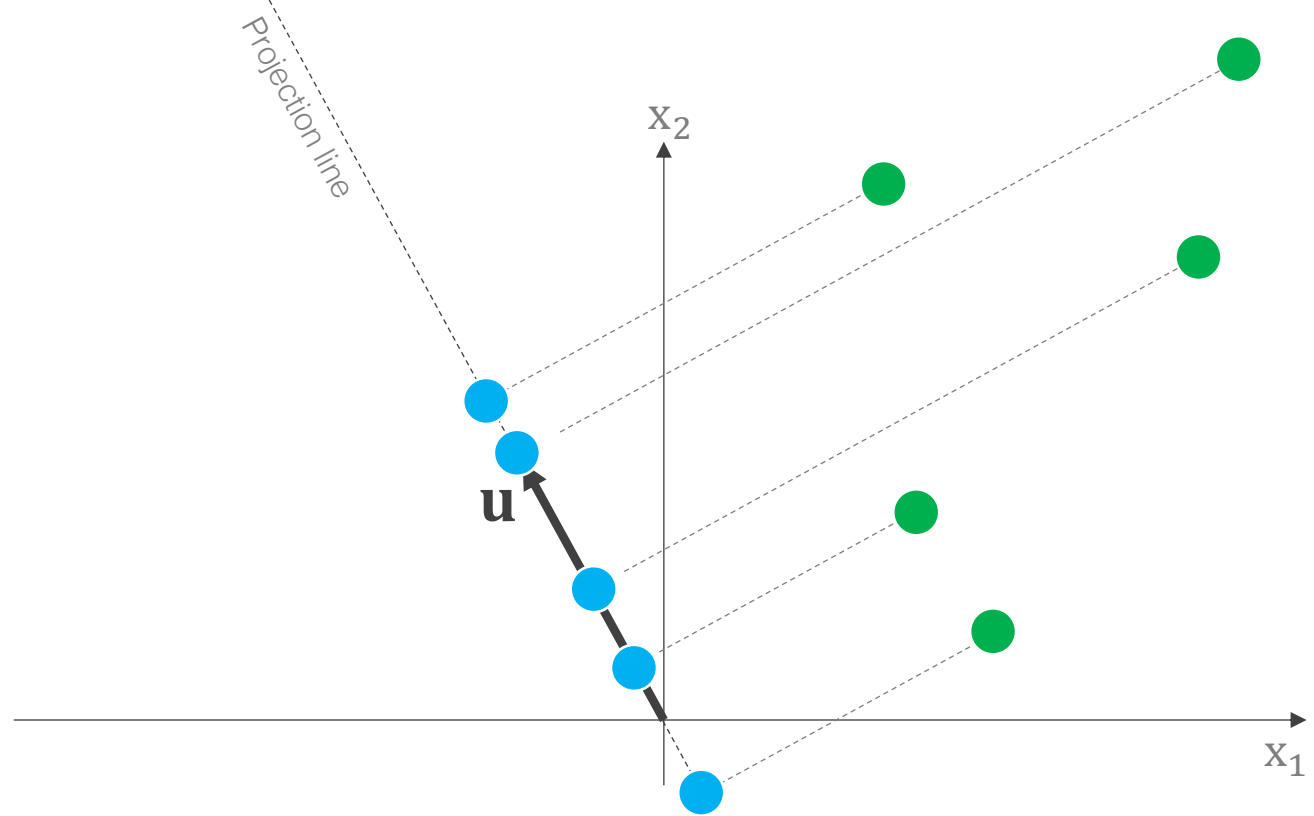
This is an inner product, but assuming \mathbf{u} is a unit vector computes this as the magnitude (length) of the projection of \mathbf{z} onto \mathbf{u}



This is valid because \mathbf{u} is a unit vector (length is 1: $\|\mathbf{u}\|_2 = \sqrt{u_1^2 + u_2^2} = \sqrt{(-0.5)^2 + (0.87)^2} \cong 1$)

Projections

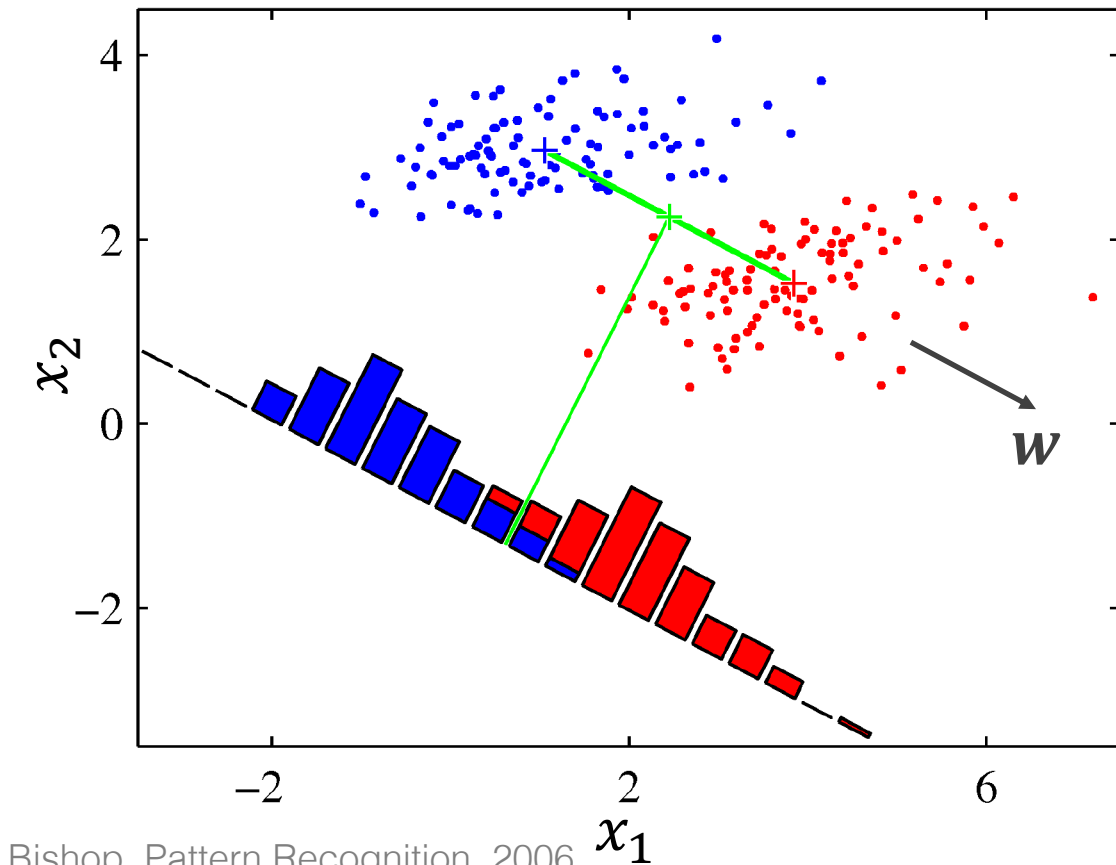
We could project any points in this space onto the line defined by the direction of unit vector \mathbf{u}



Linear Discriminant Analysis

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

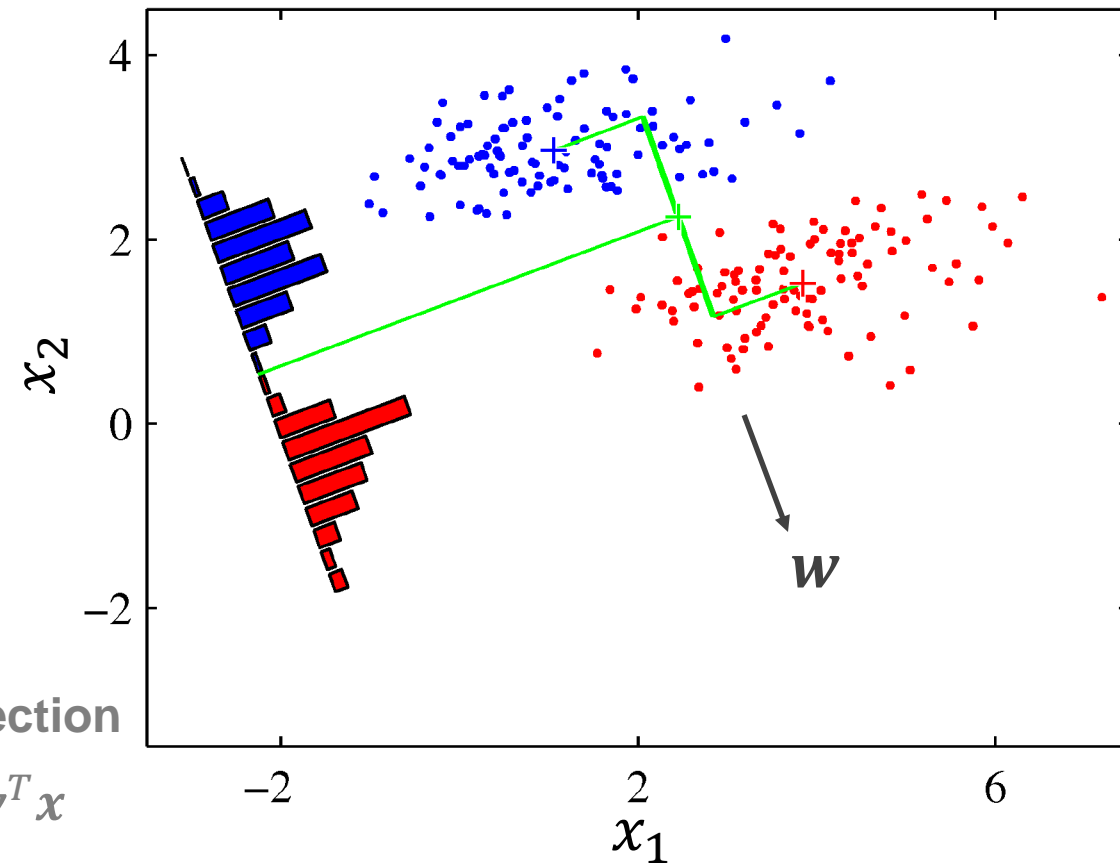
Projection onto line connecting the means



Linear projection

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

Projection onto a line providing improved class separation



Linear Discriminant Analysis (LDA)

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

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

Consider \mathbf{w} is a unit vector of parameters

Similar to PCA, but accounts for class separability

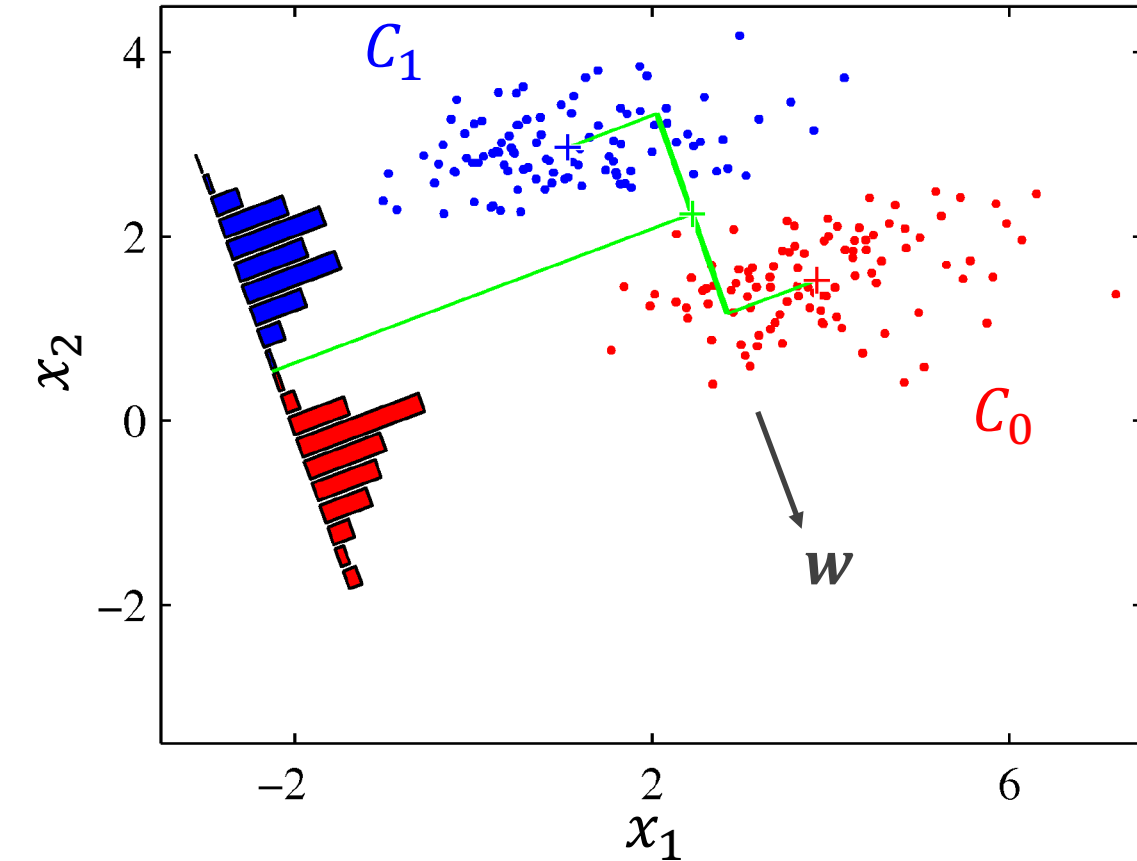
- 2 We then classify the data in this space linearly

Our decision rule becomes:

if	$\hat{f}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} > \lambda_{thresh}$	Class 1
	else	Class 0

2-Class LDA: what does vector w accomplish?

Bishop, Pattern Recognition, 2006



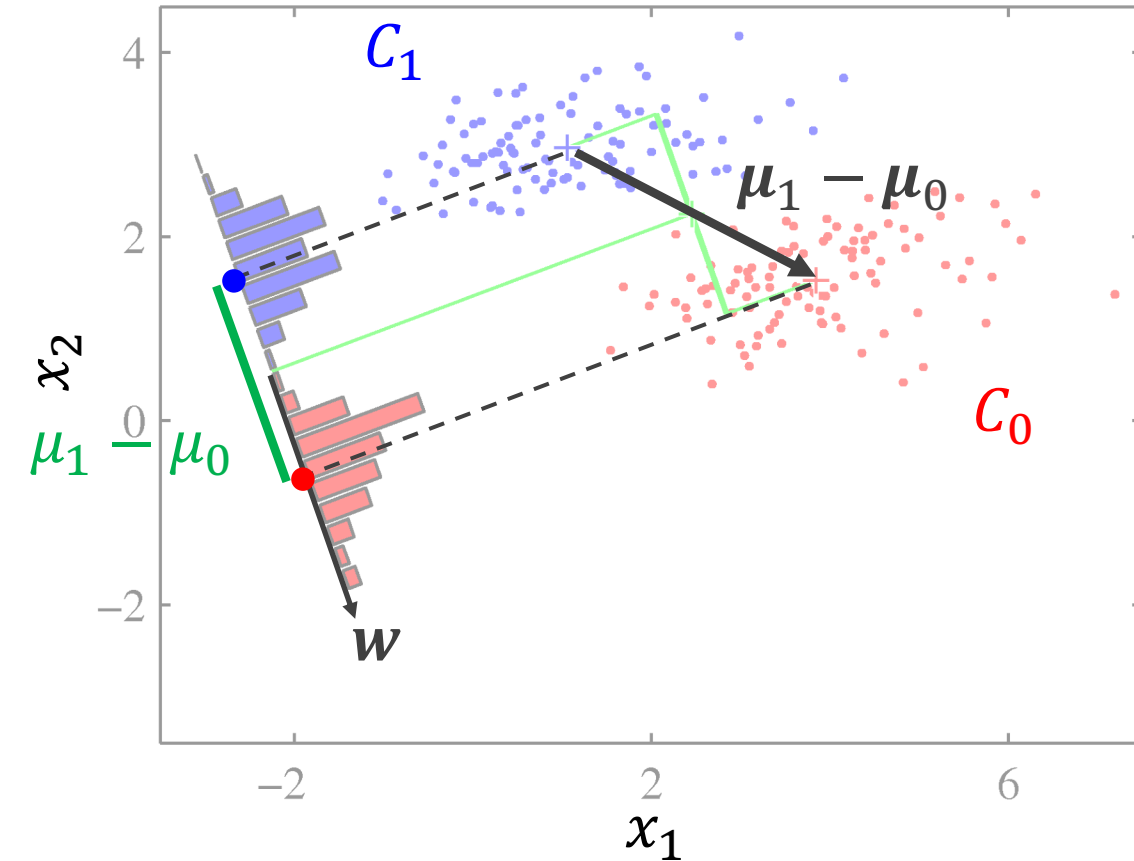
$$y = \hat{f}(x) = w^T x$$

Increase the distance between class **means**

Decrease the **variance** within the classes

2-Class LDA: what does vector w accomplish?

Bishop, Pattern Recognition, 2006



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

Increase the distance between the **means**

$$\mu_1 = \frac{1}{N_1} \sum_{i \in C_1} x_i$$

mean of class 1

$$\mu_0 = \frac{1}{N_0} \sum_{i \in C_0} x_i$$

mean of class 0

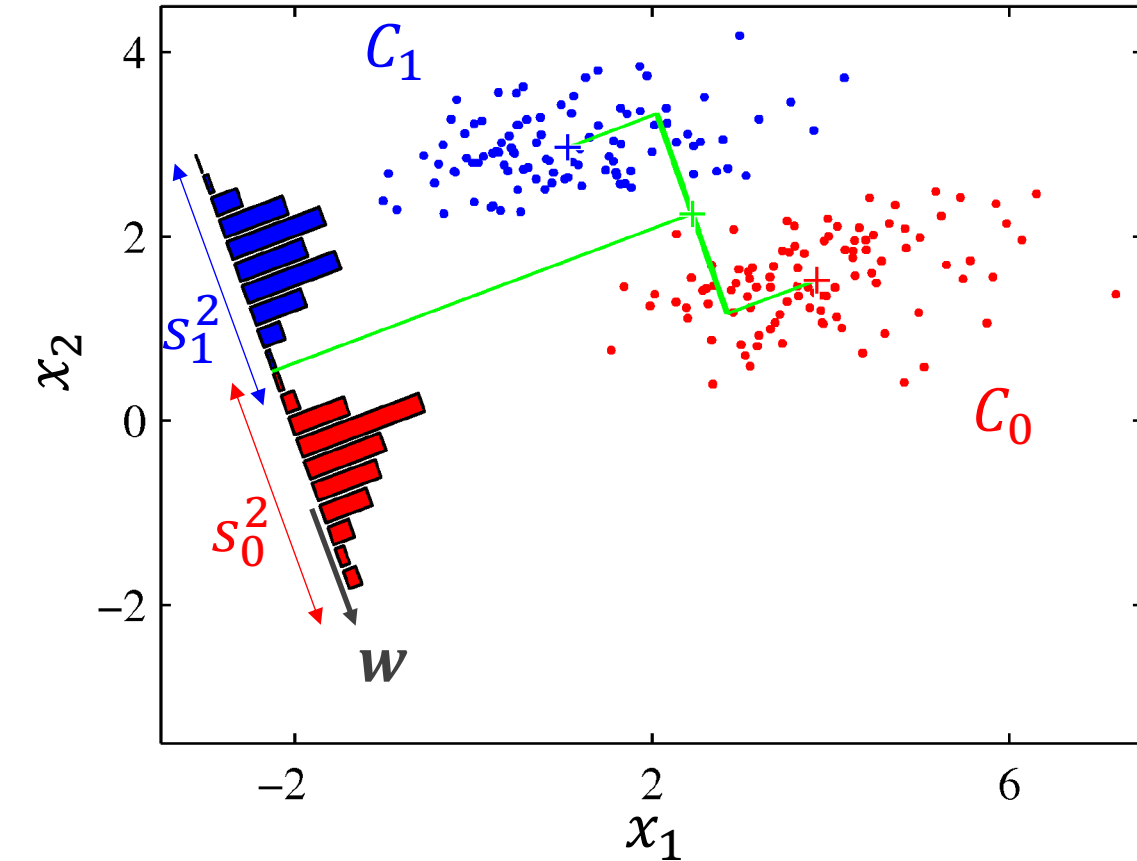
The means projected onto w : $\mu_k = w^T \mu_k$

The distance between the means:

$$\mu_1 - \mu_0 = w^T (\mu_1 - \mu_0)$$

2-Class LDA: what does vector w accomplish?

Bishop, Pattern Recognition, 2006



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

Decrease the **variance** within the classes

The “scatter” of the **projected** data:

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

where $\mu_k = w^T \mu_k$

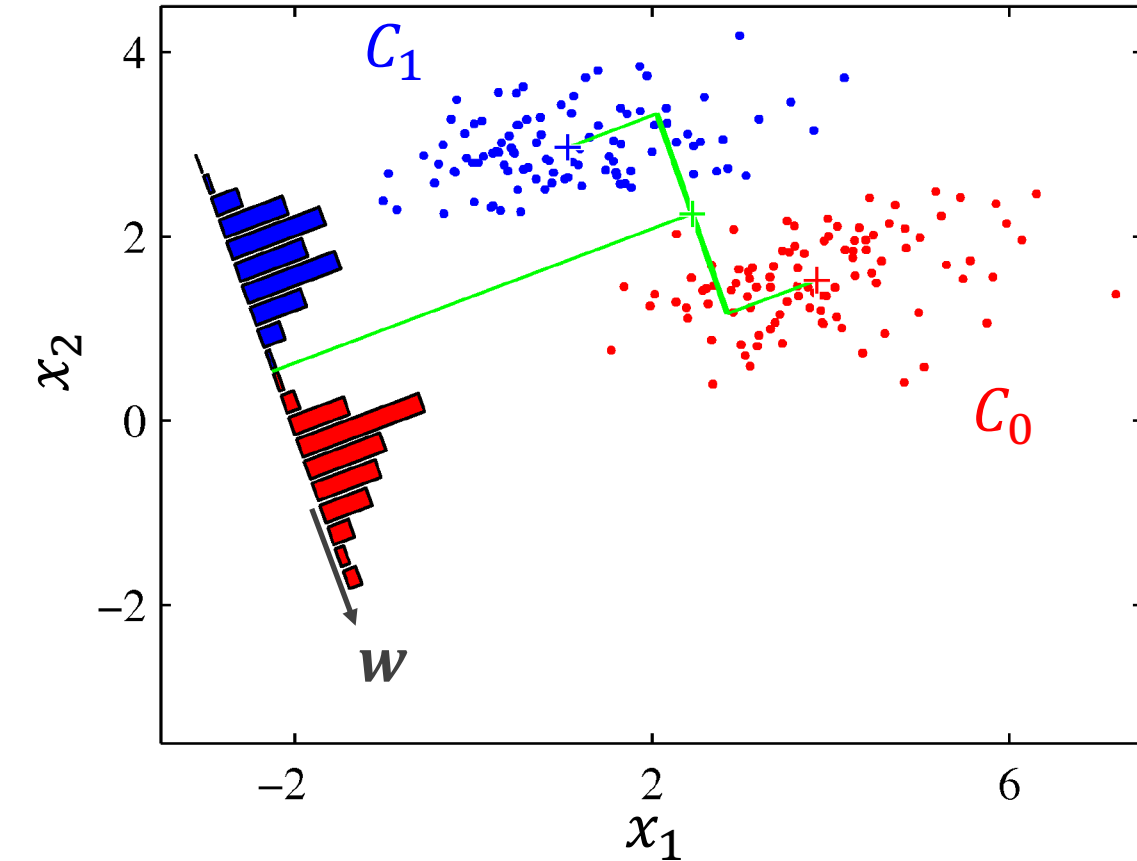
$$y_i = w^T x_i$$

We assume these are the same:

$$S = s_1^2 = s_0^2$$

2-Class LDA: what does vector w accomplish?

Bishop, Pattern Recognition, 2006



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

Increase the distance between the **means**

$$\mu_1 - \mu_0 = w^T (\mu_1 - \mu_0)$$

Decrease the **variance** within each class

$$S = s_1^2 = s_0^2$$

$$w = \Sigma^{-1} (\mu_1 - \mu_0)$$

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

**This approach is a supervised
dimensionality reduction technique
that we use for classification**

Linear Discriminant Analysis: Multiclass

Σ is the same for all classes

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

$$\delta_k(\mathbf{x}) = P(\mathbf{x}|y = k)P(y = k) = \frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)\right] \pi_k \quad P(y = k) \triangleq \pi_k$$

$$\begin{aligned} \ln|\delta_k(\mathbf{x})| &= -\frac{p}{2}\ln|2\pi| - \frac{p}{2}\ln|\Sigma| - \underbrace{\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)}_{\substack{\text{expand} \\ \text{bracket}}} + \ln|\pi_k| \\ &\quad - \frac{1}{2}[\mathbf{x}^T \Sigma^{-1} \mathbf{x} - 2\mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_k + \boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k] \end{aligned}$$

$$\ln|\delta_k(\mathbf{x})| = \cancel{-\frac{p}{2}\ln|2\pi|} - \cancel{\frac{p}{2}\ln|\Sigma|} - \cancel{\frac{1}{2}\mathbf{x}^T \Sigma^{-1} \mathbf{x}} - \mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_k + \frac{1}{2}\boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \ln|\pi_k|$$

$$\ln|\delta_k(\mathbf{x})| = -\mathbf{x}^T \Sigma^{-1} \boldsymbol{\mu}_k + \frac{1}{2}\boldsymbol{\mu}_k^T \Sigma^{-1} \boldsymbol{\mu}_k + \ln|\pi_k|$$

Since we'll be looking for the choice of class k that maximizes the value of $\delta_k(\mathbf{x})$, we can ignore terms that are independent of class

We compute this for each k and assign the class with the largest discriminant $\delta_k(\mathbf{x})$

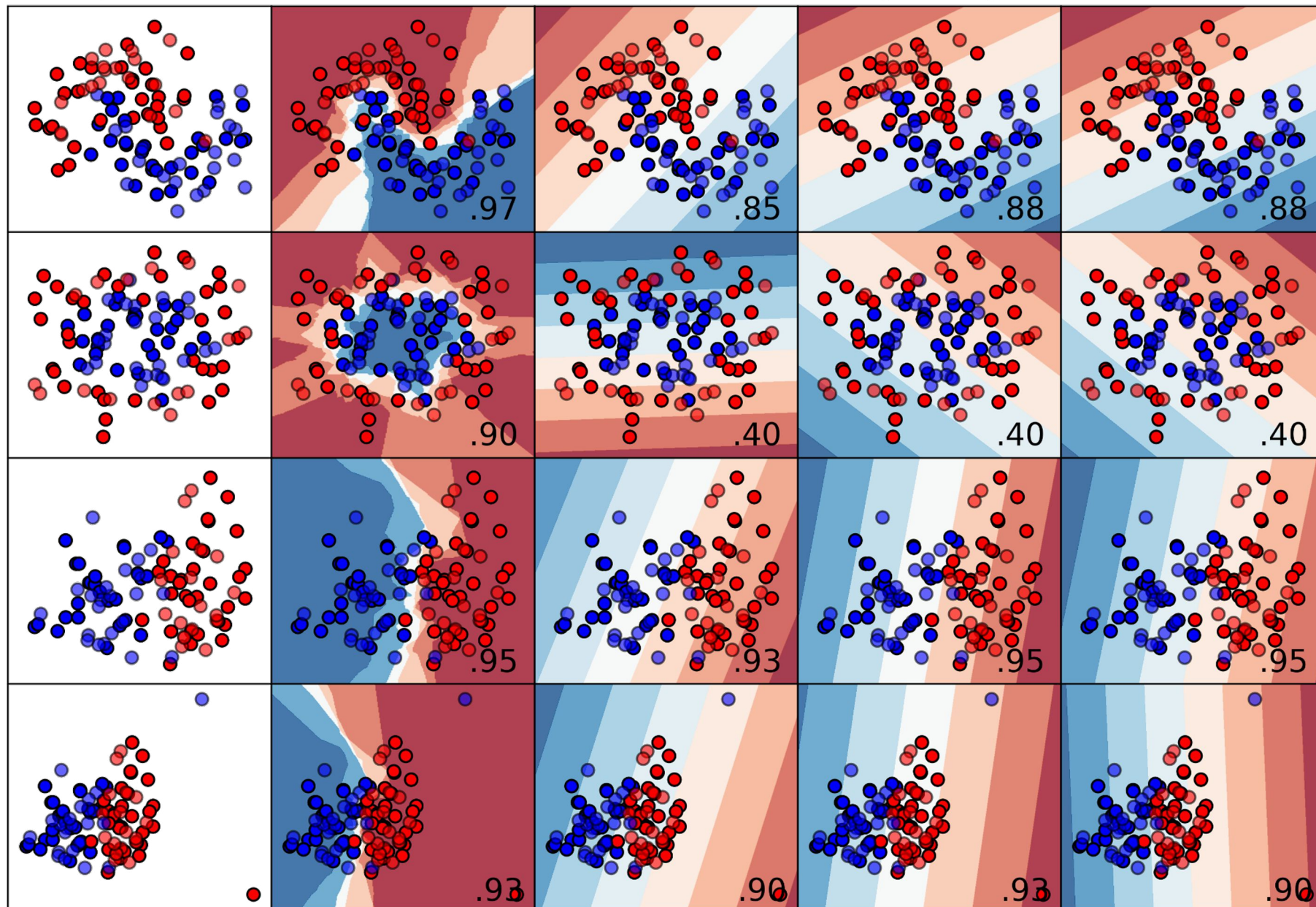
Input data

KNN (k=5)

Perceptron

Logistic Reg.

LDA



Quadratic Discriminant Analysis

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

$$d_k(\mathbf{x}) = \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_k|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right] \pi_k$$

We assume a normal distribution, but **different covariance matrices**

Produces a quadric decision boundary

Summary Comparison

	Linear Discriminant Analysis (LDA)	Quadratic Discriminant Analysis (QDA)
Assumes Gaussian Likelihood $P(\mathbf{x} \mathbf{y})$ (class conditional density)	Yes	Yes
Assumes equivalent covariance matrices $\Sigma_i = \Sigma_j$	Yes	No

Input data

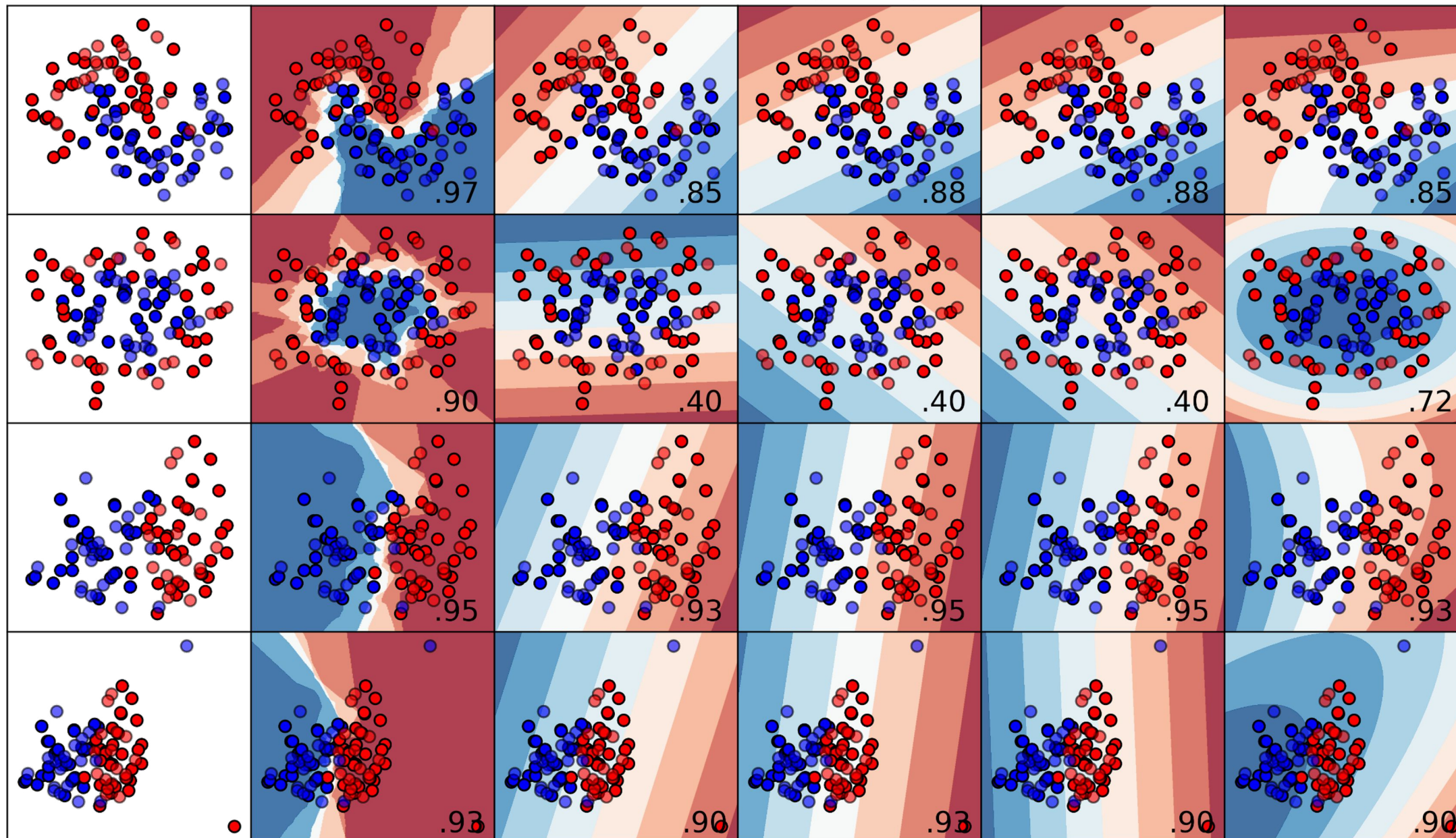
KNN (k=5)

Perceptron

Logistic Reg.

LDA

QDA



Joint vs Marginal Densities

The marginal densities don't factor in relationships between features

What if the joint density is too hard to estimate?

Note: The plot in the middle is actually a scatterplot, but could be used to estimate $P(x_1, x_2|Y)$

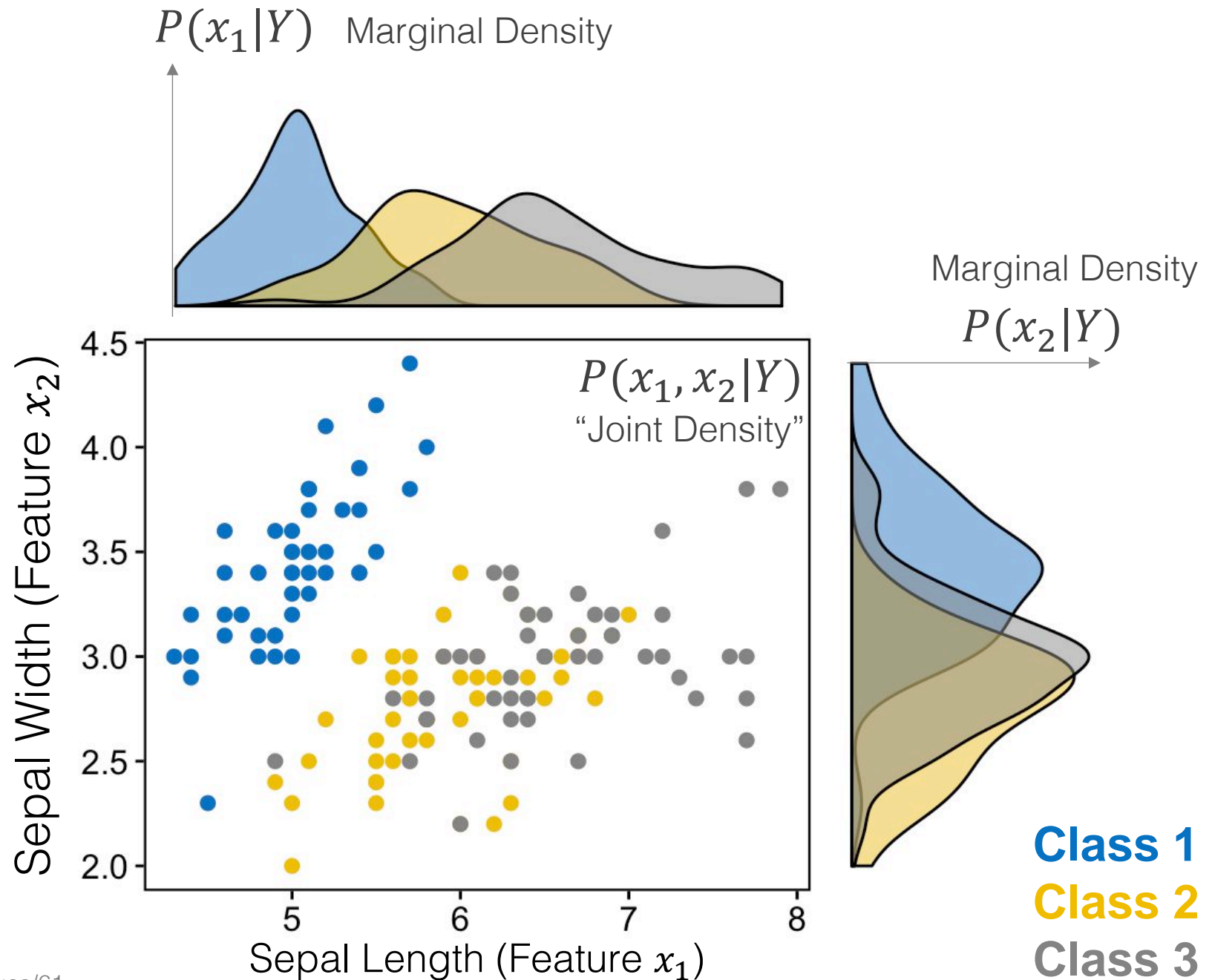


Image adapted from: <https://github.com/daattali/ggExtra/issues/61>

Naïve Bayes

Sometimes called “Idiot’s Bayes”

For independent events: A, B, and C
 $P(A \text{ and } B \text{ and } C) = P(A)P(B)P(C)$

Start with our original expression for our posterior distribution

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

Write out the full expression with all the terms in \mathbf{x}
(assume p predictors/features)

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

Assumption: Given the class, the features are independent

Note: The denominator (evidence) is a constant if we know the values of the predictor variables

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

Predict the class with the highest posterior probability

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

Naïve Bayes

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

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

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, $P(y = i | x_1, x_2, \dots, x_p)$

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

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 p features and for each class. We multiply **univariate** distributions together

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

Naïve Bayes: Parameters

p predictors, c classes

For each predictor, x_i , and class, y_j :
 $(\mu_{ij}, \sigma_{ij}^2)$

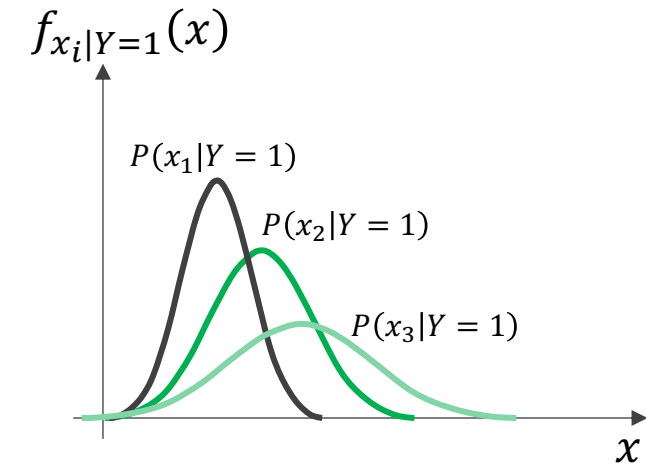
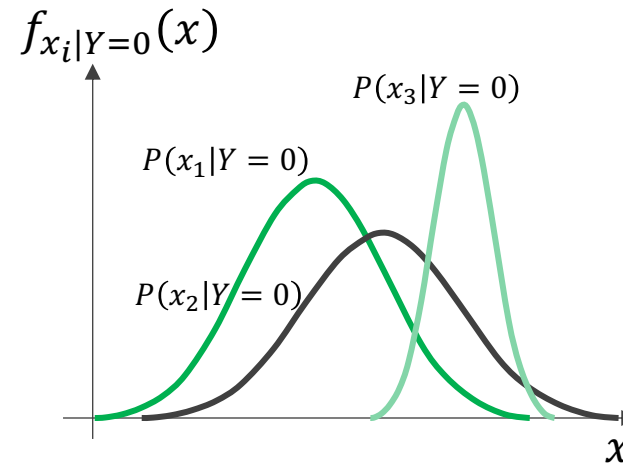
Total parameters = $2cp$

Without the Naïve Bayes assumption,
each class would be a multivariate
Gaussian with (μ_j, Σ_j)

$$\mu_j = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{bmatrix}, \Sigma_j = \begin{bmatrix} \sigma_{11}^2 & \sigma_{12}^2 & \cdots & \sigma_{1p}^2 \\ \sigma_{21}^2 & \sigma_{22}^2 & \cdots & \sigma_{2p}^2 \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{p1}^2 & \sigma_{p2}^2 & \cdots & \sigma_{pp}^2 \end{bmatrix}$$

Total parameters = $c(p + p^2)$

Naïve Bayes ($p = 3$)



Multivariate Gaussian (example shown for $p = 2$)

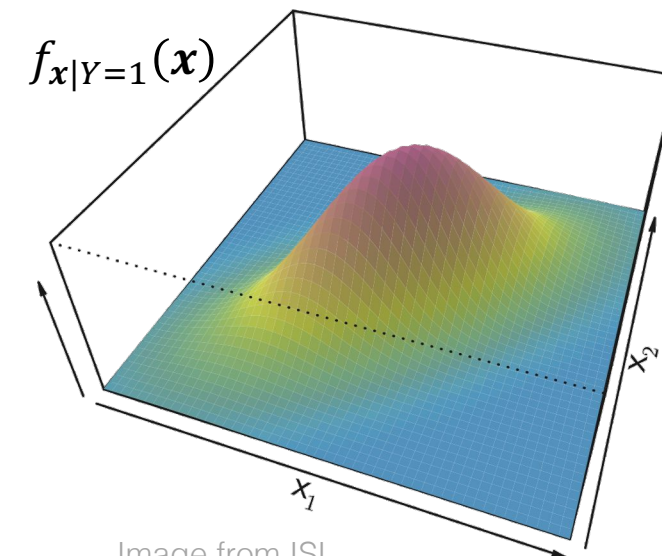
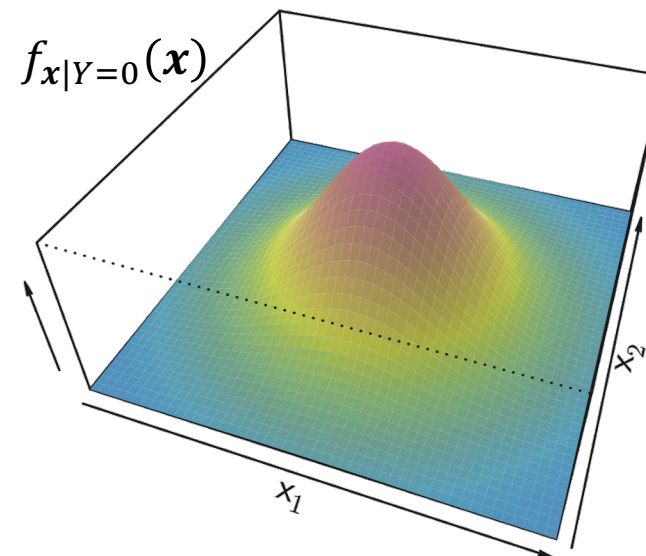


Image from ISL

Input data

KNN (k=5)

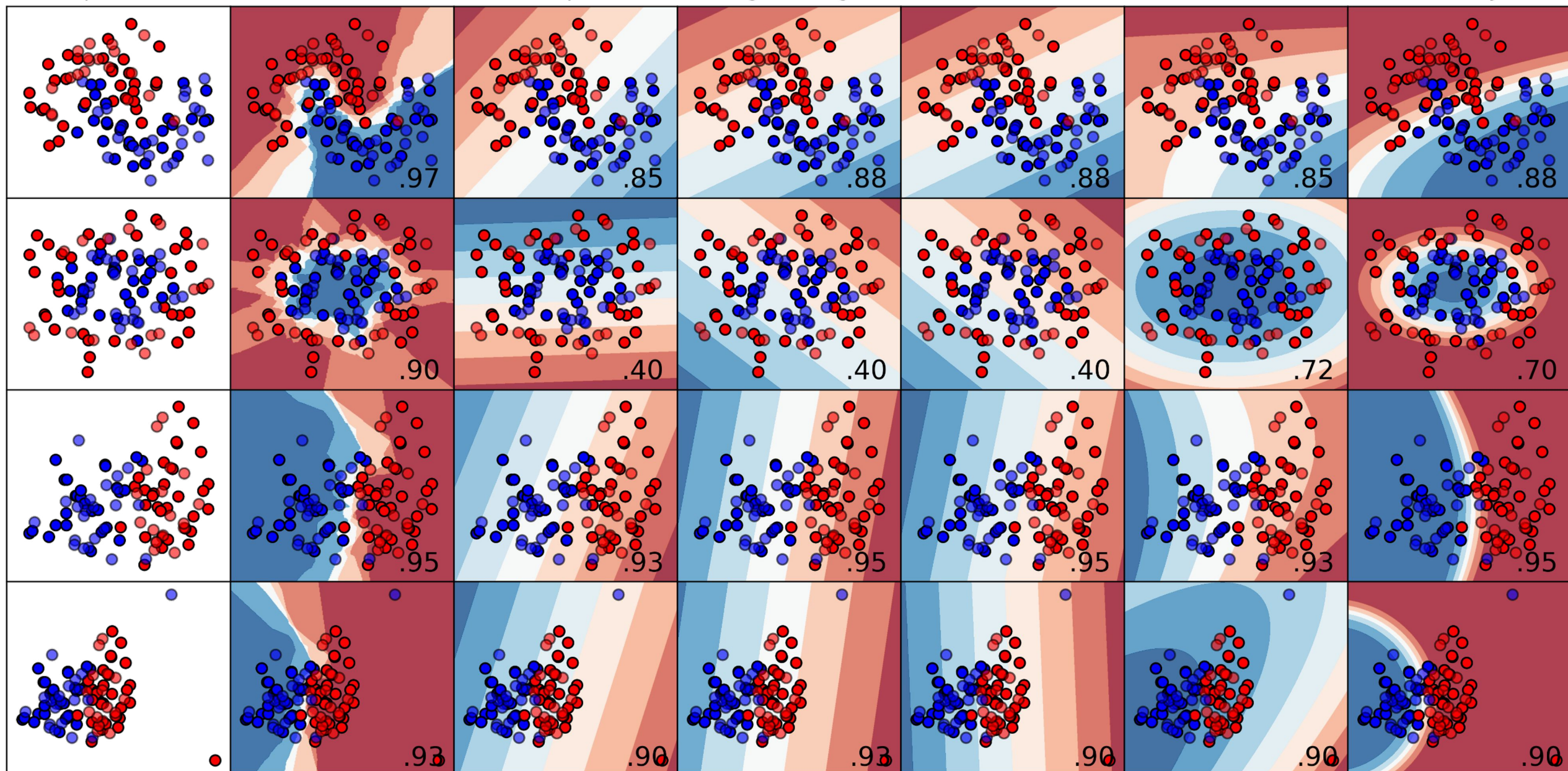
Perceptron

Logistic Reg.

LDA

QDA

Naive Bayes



Classifiers

Covered so far

K-Nearest Neighbors

Perceptron

Logistic Regression

Linear Discriminant Analysis

Quadratic Discriminant Analysis

Naïve Bayes

Have closed-form solutions
Apply to multiclass problems
Have no hyperparameters
Fast to train

Requires small amounts of training data
Only model choice is the form of $P(X|Y)$
Fast to train