

lecture 8 - 22/09/2024

* WITH MULTIPLE FEATURES WE NEED MULTIVARIATE DISTRIBUTIONS FOR THE CLASS CONDITIONALS.

Multivariate Data

- IF X_1, \dots, X_p ARE REAL-VALUED RANDOM VARIABLES:
 - WE CAN TALK ABOUT A **p-DIMENSIONAL RANDOM VECTOR** $\mathbf{X} = (X_1, \dots, X_p)$

- Usual convention \rightarrow **RANDOM VECTORS = COLUMN VECTORS**:

$$- \mathbf{X} = \begin{bmatrix} X_1 \\ \vdots \\ X_p \end{bmatrix}$$

- The **Expectation / mean** of a p-dimensional random vector \mathbf{X} is:

$$- \mathbb{E} \mathbf{X} = \begin{bmatrix} \mathbb{E} X_1 \\ \vdots \\ \mathbb{E} X_p \end{bmatrix}$$

- The **sample mean** for vector \mathbf{X} is the avg:

$$\hat{\mathbf{y}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \quad \rightarrow \text{vector of coordinate-wise avg's.}$$

- The **Variance** of a p-dimensional random vector \mathbf{X} is:

$$\text{Var } \mathbf{X} = \mathbb{E} [(X - \mathbb{E} X)(X - \mathbb{E} X)^T]$$

$$= \begin{bmatrix} \text{Var } X_1 & \text{Cov}(X_1, X_2) & \dots & \text{Cov}(X_1, X_p) \\ \text{Cov}(X_2, X_1) & \text{Var } X_2 & \dots & \text{Cov}(X_2, X_p) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(X_p, X_1) & \text{Cov}(X_p, X_2) & \dots & \text{Var } X_p \end{bmatrix} \rightarrow \text{Covariance matrix}$$

Multivariate Normal Distribution

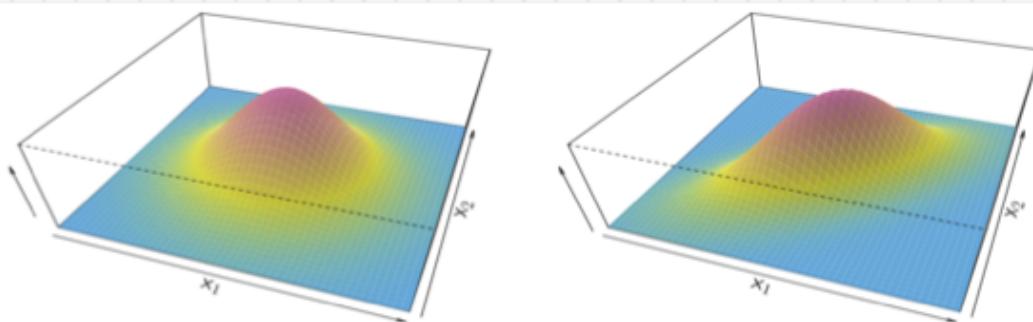


FIGURE 4.5. Two multivariate Gaussian density functions are shown, with $p = 2$. Left: The two predictors are uncorrelated. Right: The two variables have a correlation of 0.7.

- The p-dimensional multivariate Gaussian dist has prob density function:

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

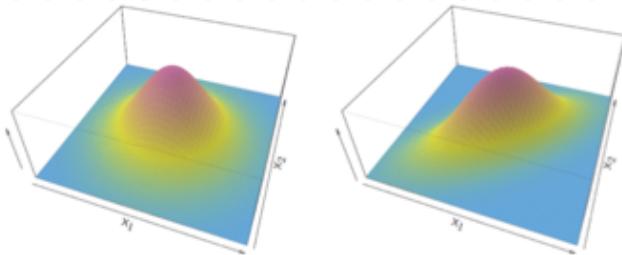
The exponent contains the (squared) Mahalanobis distance based on Σ :

$$(x - \mu)^T \Sigma^{-1} (x - \mu) = \|x - \mu\|_{\Sigma}^2$$

When $\Sigma = I$, all variables are independent and standard normal, and

$$\|x - \mu\|_I^2 = \sum_{i=1}^p (x_i - \mu_{ik})^2$$

is the squared geometric (Euclidean) distance from x to μ .



A contour curve for the PDF:

- * Shows all points with a given density val.
- * Is a quadratic form
- * Shows all points with distance from μ
- * No changing shape when we take dot.

Looking only at factors involving x the pdf is

$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{p/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

$$\propto e^{-\frac{1}{2} \|x - \mu\|_{\Sigma}^2}$$

(1xn) (nxn) (nx1)

$\|x\|_{\Sigma}$

SCALAR

of random vars

If X is p-dimensional multivariate normal, $MVN(\mu, \Sigma)$, then:

$$Ax + b \sim MVN(A\mu + b, A\Sigma A^T)$$

where A is any $q \times p$ -matrix

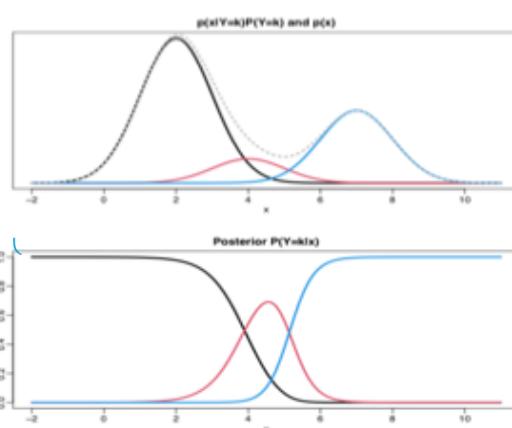
JDA & QDA with multiple features!

* Remember...

- JDA is a generative model where class conditionals $p(x | Y=k)$ are assumed Gaussian with individual class means, but equal covariance matrices.
- QDA is the same apart from the fact that they have class-specific covariance matrix.

JDA EXAMPLE

- In 1 dimension we have what we've seen till now:



In 1 dimension: Conditionally on the class we assumed the feature x to have a univariate Gaussian distribution as

$$p(x | Y = \text{black}) = \mathcal{N}(2, 1)$$

$$p(x | Y = \text{red}) = \mathcal{N}(4, 1)$$

$$p(x | Y = \text{blue}) = \mathcal{N}(7, 1)$$

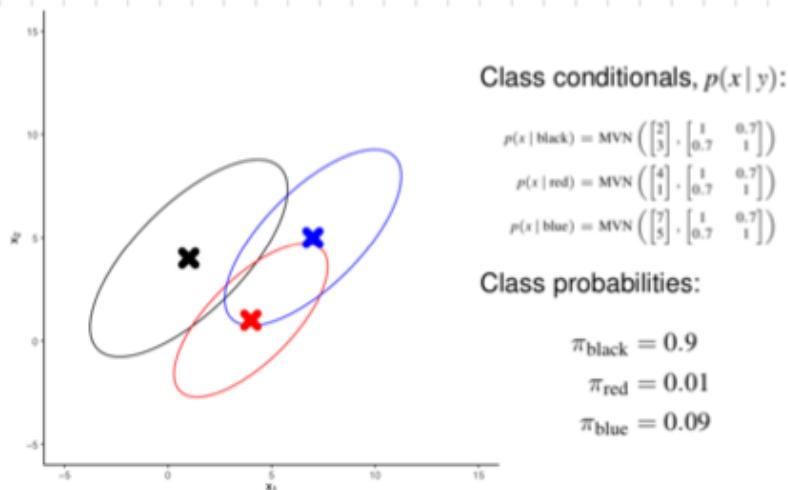
The class probabilities we took to be

$$\pi_{\text{black}} = 0.6$$

$$\pi_{\text{red}} = 0.1$$

$$\pi_{\text{blue}} = 0.3$$

- Now by using multiple features:



Plot of $p(x,y)$ contours.

REMEMBER:

- $\text{MVN}(\mu, \Sigma)$

$\mu = \text{mean}$

$$\Sigma = \begin{bmatrix} \text{Var}_{xx} & \text{Cov}_{xy} \\ \text{Cov}_{yx} & \text{Var}_{yy} \end{bmatrix}$$

- This function is a bit **very** so we can get rid of some params:
 - to do this we can take dots & multiply by 2

$$g_k(x) = 2\log \pi_k - (x - \mu_k)^T \Sigma^{-1} (x - \mu_k)$$

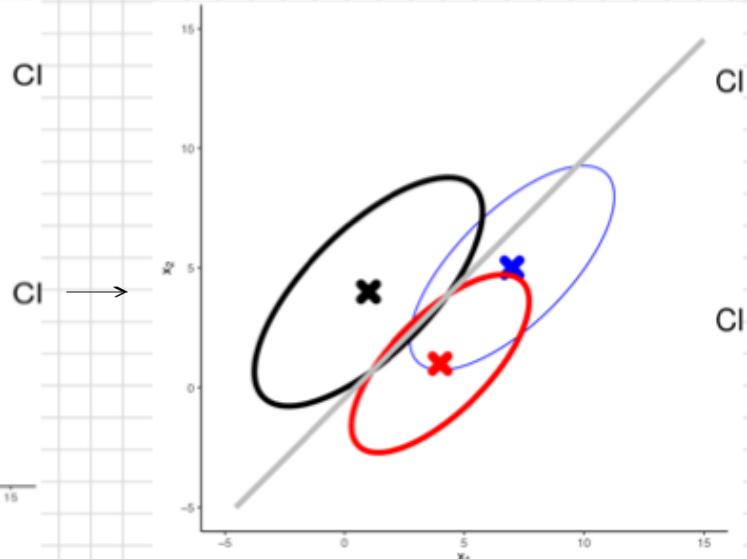
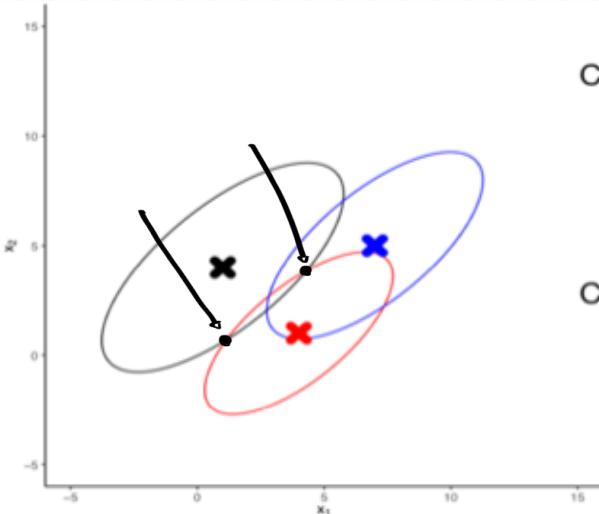
- A point is classified to the closest mean in terms of **Mahalanobis distance**.

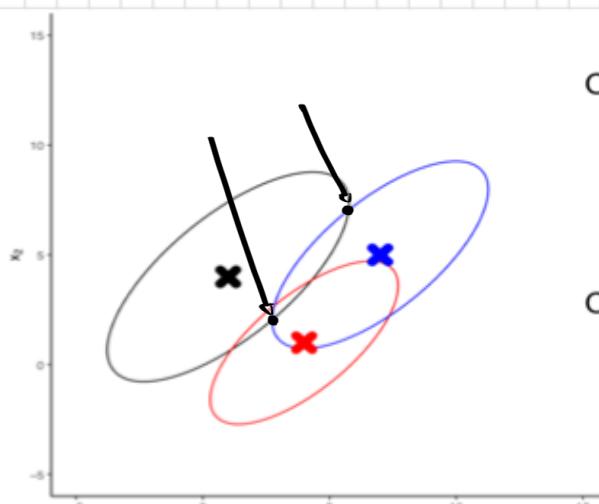
Constructing decision boundary

- We wanna **classify** x by choosing k with **HIGHEST DISCRIMINANT**
- The **decision boundary** between j & k consists of **all points where**:

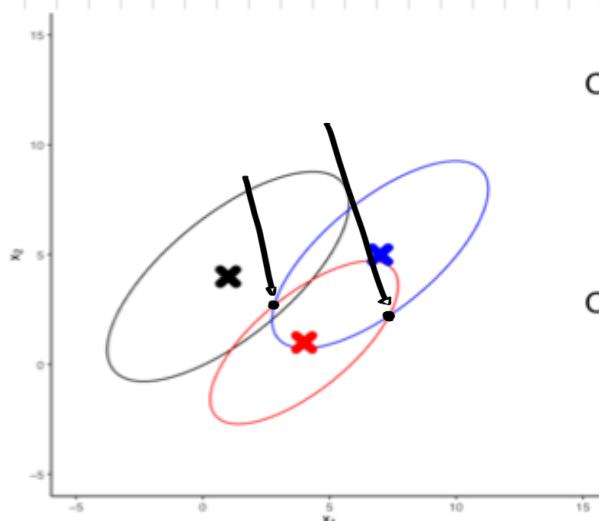
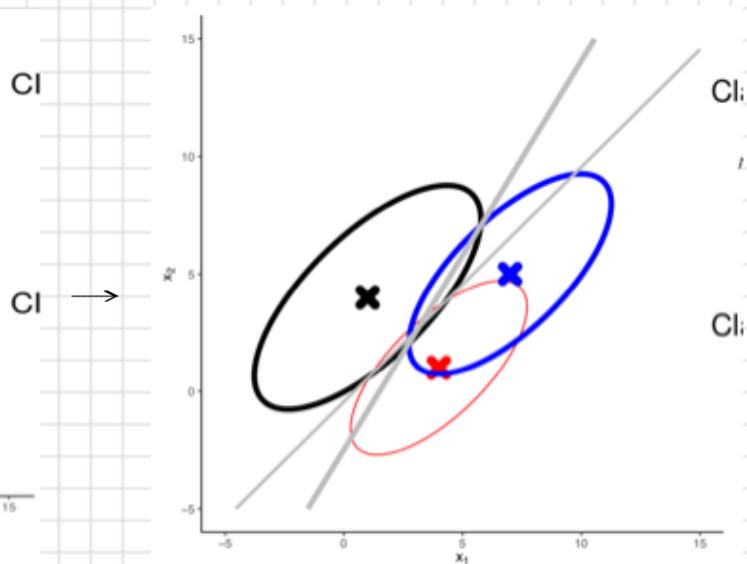
$$g_j(x) = g_k(x) \rightarrow \text{INTERSECTION}$$
- The **intersection** points between contour curves for a value g are where:

$$g_j(x) = g_k(x) = g$$
- We call the intersection between the curves a **HYPERPLANE**:
 - For **ONE FEATURE** a **POINT**
 - For **TWO FEATURES** a **LINE**
 - For **THREE FEATURES** a **PLANE**.

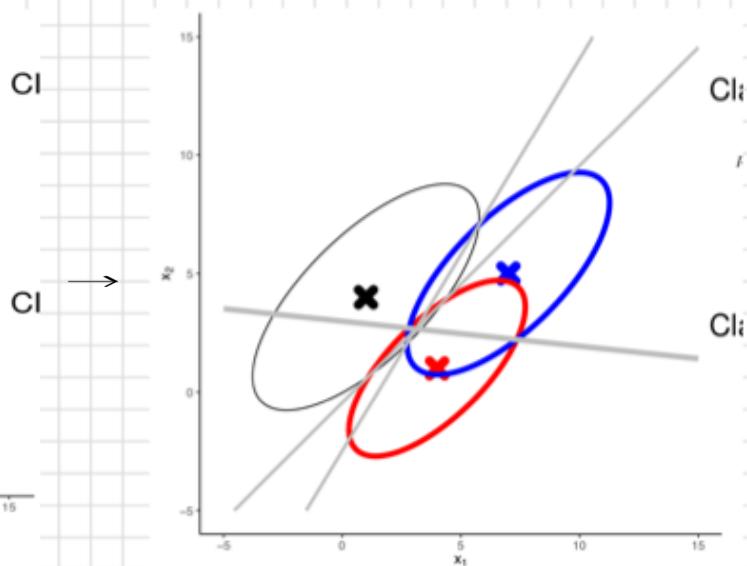




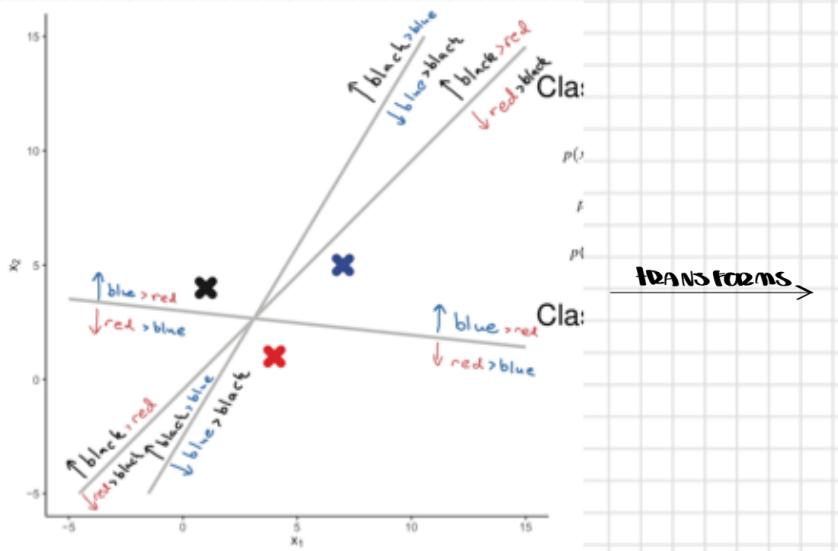
Plot of $p(x, y)$ contours.



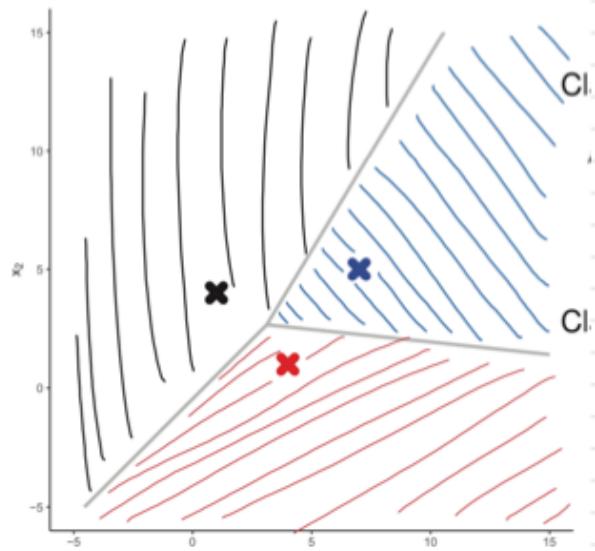
Plot of $p(x, y)$ contours.



- Now that we were able to construct the decision boundaries
 - We need to actually cover the divided area with the classes



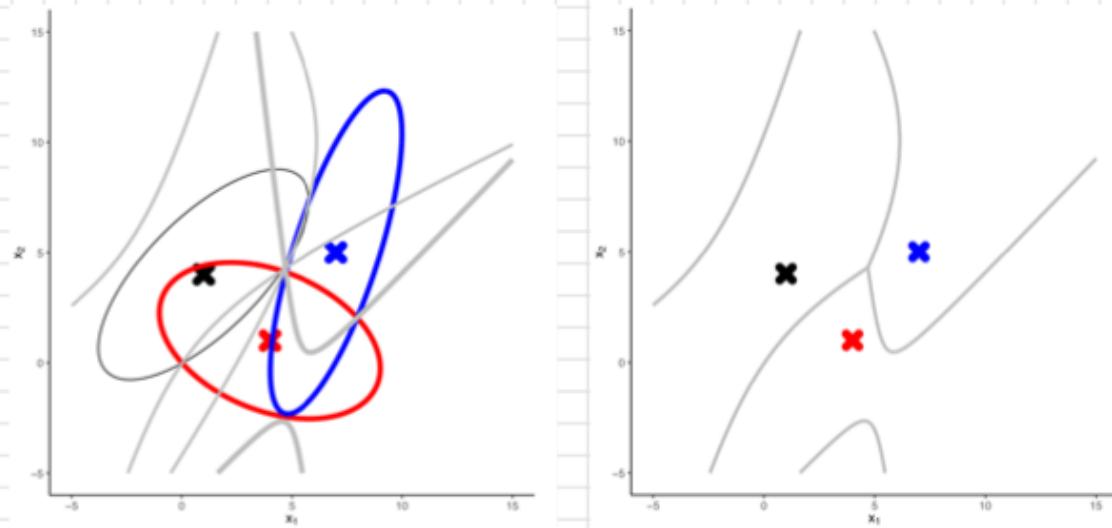
TRANSFORMS



- By changing priors the decision boundary moves away:
 - It moves away from the mean of K with highest prior.

Constructing Decision Boundary for QDA!

- IT IS DEFINITELY MORE COMPLEX TO CONSTRUCT FOR QDA!
 - THIS IS DUE TO ITS FLEXIBILITY (CURVE SHAPED)



Issues:

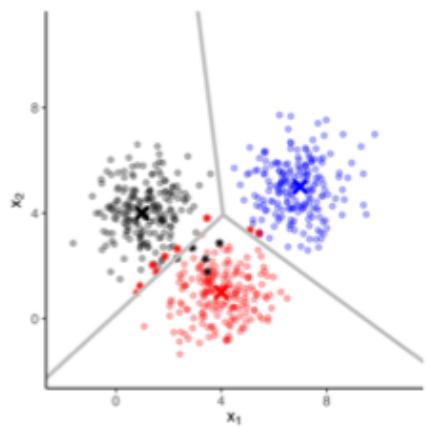
- THE MODEL THAT BEST CAPTURES VARIABILITY IN DATA MAY HAVE TOO MANY PARAMS
 - A COVARIANCE MATRIX NEEDS $p(p+1)/2$ PARAMS.
 - LUCKILY → A PRACTICAL REDUCTION
 - * QDA → ≠ COVARIANCE MATRIX Σ
 - * JDA → SAME COVARIANCE MATRIX FOR ALL CLASSES
- GAUSSIAN Naïve Bayes:
 - DIAGONAL COVARIANCE MATRIX ($\neq \Sigma$ OR $\neq \Sigma_k$)

NOTE!

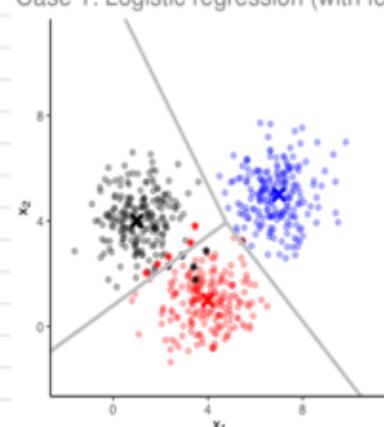
Naïve Bayes classifiers are generative models with a simplifying assumption that all features are independent, when specifying class conditionals.

Case Studies for LDA, QDA, Logistic Regression!

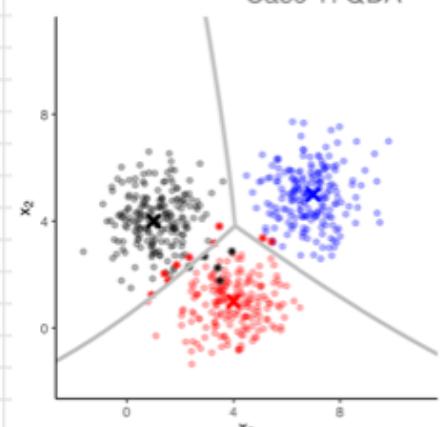
Case 1: LDA



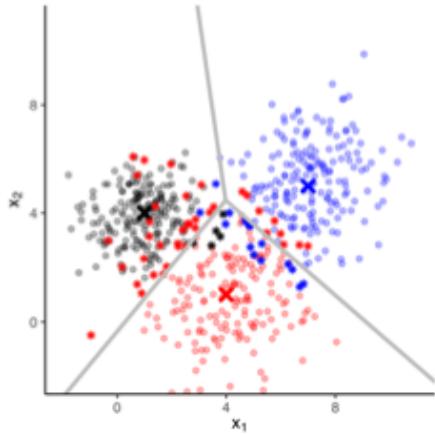
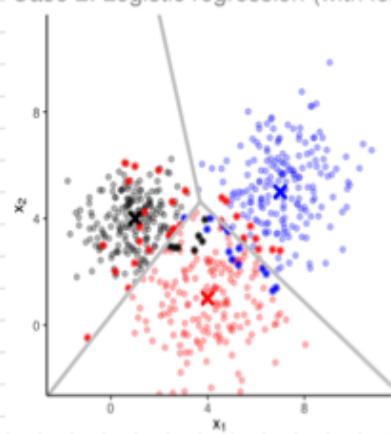
Case 1: Logistic regression (with features x_1 and x_2)



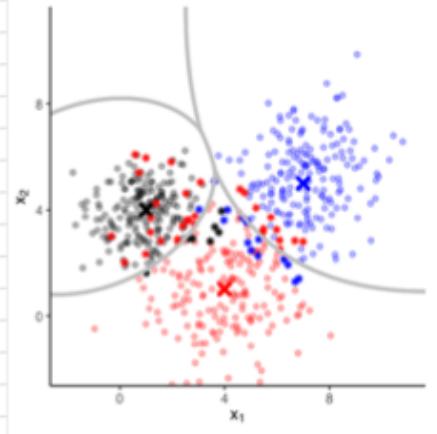
Case 1: QDA



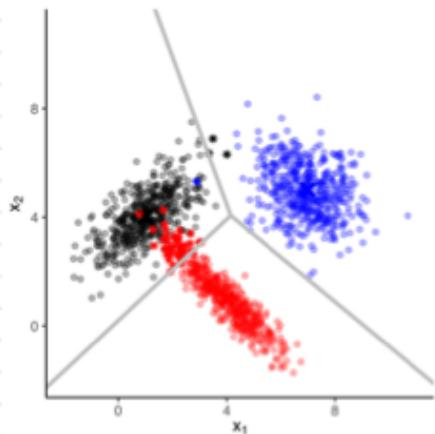
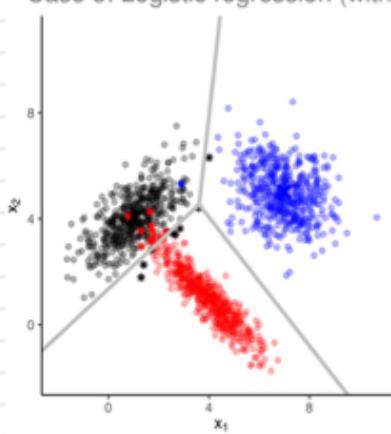
Case 2: LDA

Case 2: Logistic regression (with features x_1 and x_2)

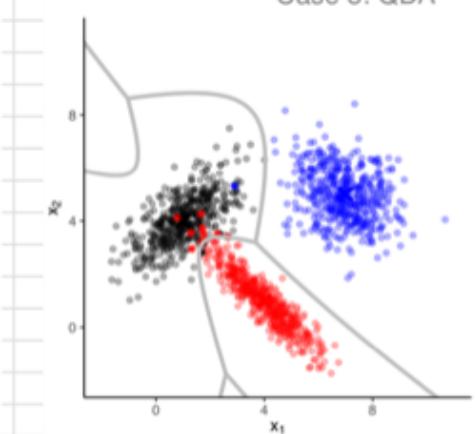
Case 2: QDA



Case 3: LDA

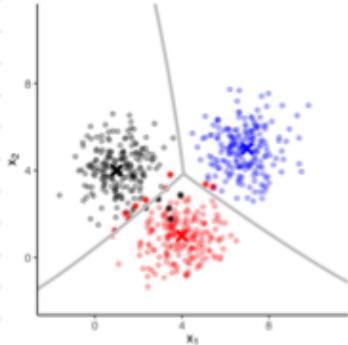
Case 3: Logistic regression (with features x_1, x_2)

Case 3: QDA

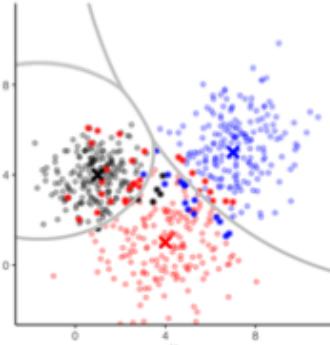


CASE Study for Naive Bayes

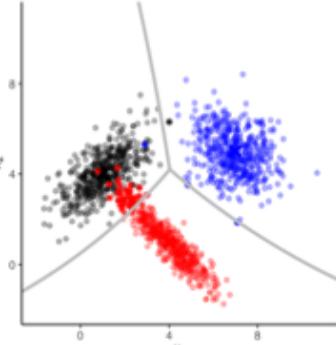
Case 1: Naive Bayes (different, but diagonal Cov)



Case 2: Naive Bayes (different, but diagonal Cov)



Case 3: Naive Bayes (different, but diagonal Cov)



Test Error for the three cases

Case 1 → Equal covariance matrices

Case 2 → Unequal covariance matrices, no correlation

Case 3 → Unequal covariance matrices with correlation

	LDA	QDA	NB	LR	LR(sq.)
Case 1	1.20	1.30	1.33	1.40	1.33
Case 2	8.00	7.27	7.43	7.60	7.57
Case 3	2.27	0.63	2.20	0.87	0.60

JDA vs Logistic Regression

SMALL K
BIG K

- WHAT HAPPENS IS:

- WE CLASSIFY TO K OVER K WHENEVER:

$$P(y=k|x) > P(y=K|x) \xrightarrow{\text{THIS IS}} \log\left(\frac{P(y=k|x)}{P(y=K|x)}\right) > 0$$

- WHAT LOGISTIC REGRESSION DOES:

- DIRECTLY MODELS LOG OF POSTERIOR ODDS BETWEEN K & K:

It models it as a LINEAR COMB → $\log\left(\frac{P(y=k|x)}{P(y=K|x)}\right) = \mathbf{a}_k + \mathbf{b}_k^T \mathbf{x}$

- WHAT JDA DOES:

- JDA ASSUMPTIONS IMPLY A LINEAR MODEL FOR THE LOG OF POSTERIOR ODDS

- * QDA IMPLIES A QUADRATIC MODEL.

- JDA ASSUMES THAT CLASS CONDITIONALS ARE GAUSSIAN

↳ IF TRUE, THEN JDA IS A MORE EFFICIENT CLASSIFIER THAN LR

JR OFFERS MORE FLEXIBILITY IN CREATING DECISION BOUNDARIES (X_1, X_2 , 1 Predictor, etc...)

BUT MODEL GIVE LINEAR DECISION BOUNDARIES!!!