# Discriminant Analysis

# Recall the Heart Data (for classification)

response variable $Y$ is Yes/No

| Age | Sex | ChestPain | RestBP | Chol | Fbs | RestECG | MaxHR | ExAng | Oldpeak | Slope | Ca | Thal | AHD |
|-----|-----|-----------|--------|------|-----|---------|-------|-------|---------|-------|-----|------|-----|
| 63 | 1 | typical | 145 | 233 | 1 | 2 | 150 | 0 | 2.3 | 3 | 0.0 | fixed | No |
| 67 | 1 | asymptomatic | 160 | 286 | 0 | 2 | 108 | 1 | 1.5 | 2 | 3.0 | normal | Yes |
| 67 | 1 | asymptomatic | 120 | 229 | 0 | 2 | 129 | 1 | 2.6 | 2 | 2.0 | reversable | Yes |
| 37 | 1 | nonanginal | 130 | 250 | 0 | 0 | 187 | 0 | 3.5 | 3 | 0.0 | normal | No |
| 41 | 0 | nontypical | 130 | 204 | 0 | 2 | 172 | 0 | 1.4 | 1 | 0.0 | normal | No |

# Linear Discriminant Analysis (LDA)

- Linear discriminant analaysis (LDA) takes a different approach to classification than logistic regression. Rather than attempting to model the conditional distribution of *Y* given *X*, $P(Y = k|X = x)$, LDA models the distribution of the predictors *X* given the different categories that *Y* takes on, $P(X = x|Y = k)$.

- In order to flip these distributions around to model *P(X = x|Y = k)* an analyst uses Bayes' theorem.

$$P(A_i|B) = \frac{P(B|A_i)P(A_i)}{\sum_{k=1}^{n} P(B|A_k)P(A_k)}$$

- In this setting with one feature (one *X*), Bayes' theorem can then be written as:

- What does this mean?

$$P(Y = k|X = x) = \frac{\pi_k P(X = x|Y = k)}{\sum_{j=1}^{K} \pi_j P(X = x|Y = j)}$$

# LDA

$$P(Y = k|X = x) = \frac{f_k(x)\pi_k}{\sum_{j=1}^{K} f_j(x)\pi_j}$$

- The left hand side, $P(Y = k|X = x)$, is called the *posterior* probability and gives the probability that the observation is in the $k^{th}$ category given the feature, *X*, takes on a specific value, *x*.  The numerator on the right is conditional distribution of the feature within category $k$, $f_k(x)$, times the *prior* probability that observation is in the $k^{th}$ category.

- The *Bayes' classifier* is then selected.  That is the observation assigned to the group for which the posterior probability is the largest.

# LDA for one predictor

- LDA has the simplest form when there is just one predictor/feature ($p = 1$). In order to estimate $f_k(x)$, we have to assume it comes from a specific distribution. If $X$ is quantitative, what distribution do you think we should use?

- One common assumption is that $f_k(x)$ comes from a Normal distribution:

$$f_k(x) = \frac{1}{\sqrt{2\pi\sigma_k^2}} \exp\left(-\frac{(x - \mu_k)^2}{2\sigma_k^2}\right)$$

- In shorthand notation, this is often written as $(X|Y = k) \sim N(\mu_k, \sigma_k^2)$, meaning, the distribution of the feature $X$ within category $k$ is Normally distributed with mean $\mu_k$ and variance $\sigma_k^2$.

# LDA for one predictor (cont.)

- An extra assumption that the variances are equal,

- $\sigma_1^2 = \sigma_2^2 = \cdots = \sigma_K^2$ will simplify are lives.

- Plugging this assumed likelihood into the Bayes' formula (to get the posterior) results in:

$$P(Y = k | X = x) = \frac{\pi_k \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu_k)^2}{2\sigma^2}\right)}{\sum_{j=1}^{K} \pi_j \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x-\mu_j)^2}{2\sigma^2}\right)}$$

- The Bayes classifier will be the one that maximizes this over all values chosen for *x*. How should we maximize?

- So we take the log of this expression and rearrange to simplify our maximization...

# LDA for one predictor (cont.)

- So we maximize the following simplified expression:

$$\delta_k(x) = x\frac{\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2} + \log \pi_k$$

- How does this simplify if we have just two classes ($K = 2$) and if we set our prior probabilities to be equal?

- This is equivalent to choosing a decision boundary for $x$ for which

$$x = \frac{\mu_1^2 - \mu_2^2}{2(\mu_1 - \mu_2)} = \frac{\mu_1 + \mu_2}{2}$$

- Intuitively, why does this expression make sense? What do we use in practice?

# LDA for one predictor (cont.)

- In practice we don't know the true mean, variance, and prior.  So we estimate them with the classical estimates, and plug-them into the expression:

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i:y_i=k} x_i$$

- and

$$\hat{\sigma}^2 = \frac{1}{n-K} \sum_{k=1}^{K} \sum_{i:y_i=k} (x_i - \hat{\mu}_k)^2$$

- where $n$ is the total sample size and $n_k$ is the sample size within class $k$ (thus, $n = \sum n_k$).

# LDA for one predictor (cont.)

- This classifier works great if the classes are about equal in proportion, but can easily be extended to unequal class sizes.

Instead of assuming all priors are equal, we instead set the priors to match the 'prevalence' in the data set:

$$\hat{\pi}_k = \hat{n}_k/n$$

- Note: we can use a prior probability from knowledge of the subject as well; for example, if we expect the test set to have a different prevalence than the training set.
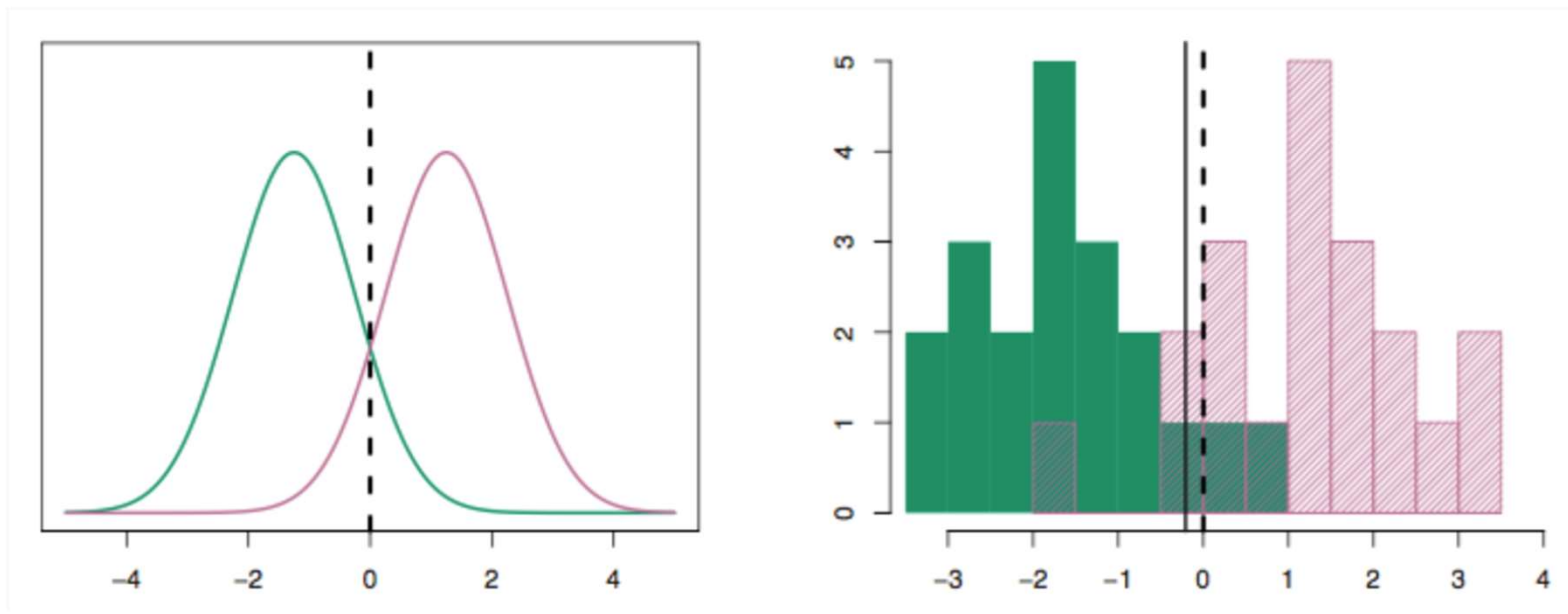
# LDA for one predictor (cont.)

- Plugging all of these estimates back into the original logged maximization formula we get:

$$\hat{\delta}_k(x) = x\frac{\hat{\mu}_k}{\hat{\sigma}^2} - \frac{\hat{\mu}_k^2}{2\hat{\sigma}^2} + \log \hat{\pi}_k$$

- Thus this classifier is called the linear discriminant classifier: this discriminant function is a linear function of *x*.
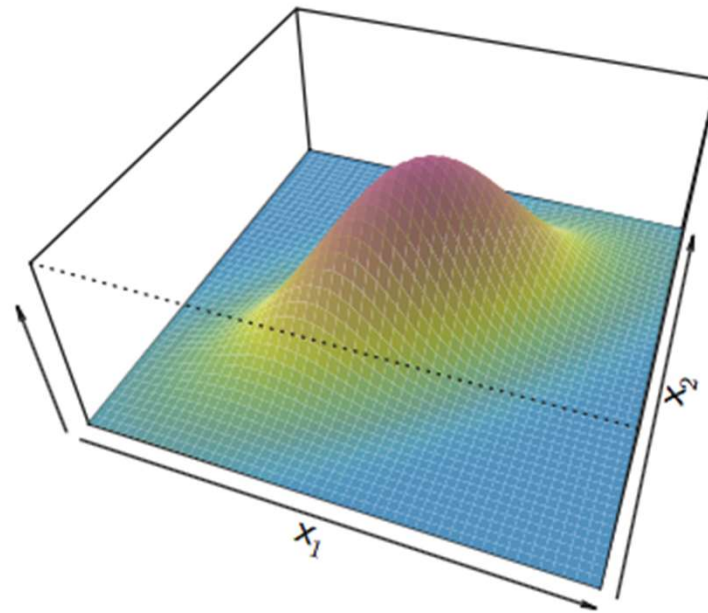
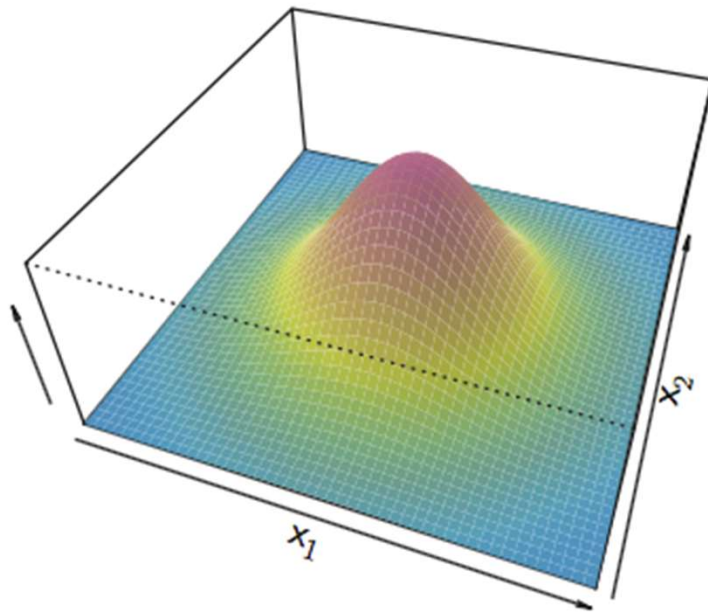# Illustration of LDA when $p = 1$

# LDA when $p > 1$

- LDA generalizes 'nicely' to the case when there is more than one predictor.

- Instead of assuming the one predictor is Normally distributed, it assumes that the set of predictors for each class is 'multivariate normal distributed' (shorthand: MVN).  What does that mean?

- This means that the vector of $X$ for an observation has a multidimensional normal distribution with a mean vector, $\mu$, and a covariance matrix, $\Sigma$.

# Multivariate Normal Distribution

- Here is a visualization of the Multivariate Normal distribution with 2 variables:

# MVN Distribution

• The joint PDF of the Multivariate Normal distribution, $\vec{X} \sim MVN(\vec{\mu}, \Sigma)$ , is:

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

• where $\vec{X}$ is a $p$ dimensional vector and $|\Sigma|$ is the determinant of the $p$ x $p$ covariance matrix.

• Let's do a quick dimension analysis sanity check...

• What do $\vec{X}$ and $\Sigma$ look like?

# LDA when $p > 1$

- Discriminant analysis in the multiple predictor case assumes the set of predictors for each class is then multivari $\vec{X} \sim MVN(\vec{\mu}, \Sigma)$

- Just like with LDA for one predictor, we make an extra assumption that the covariances are equal in each group, $\Sigma_1 = \Sigma_2 = \cdots = \Sigma_K$. in order to simplify our lives.

- Now plugging this assumed likelihood into the Bayes' formula (to get the posterior) results in:

$$P(Y = k | \vec{X} = \vec{x}) = \frac{\pi_k \frac{1}{2\pi^{p/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\vec{x} - \vec{\mu}_k)^T \Sigma^{-1} (\vec{x} - \vec{\mu}_k)\right)}{\sum_{j=1}^{K} \frac{1}{2\pi^{p/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\vec{x} - \vec{\mu}_j)^T \Sigma^{-1} (\vec{x} - \vec{\mu}_j)\right)}$$

# LDA when $p > 1$ (cont.)

- Then doing the same steps as before (taking log and maximizing), we see that the classification will for an observation based on its predictors, $\vec{x}$, will be the one that maximizes (maximum of $K$ of these $\delta_k(\vec{x})$):

$$\delta_k(\vec{x}) = \vec{x}^T \Sigma^{-1} \vec{\mu}_k - \frac{1}{2}\vec{\mu}_k^T \Sigma^{-1}\vec{\mu}_k + \log \pi_k$$
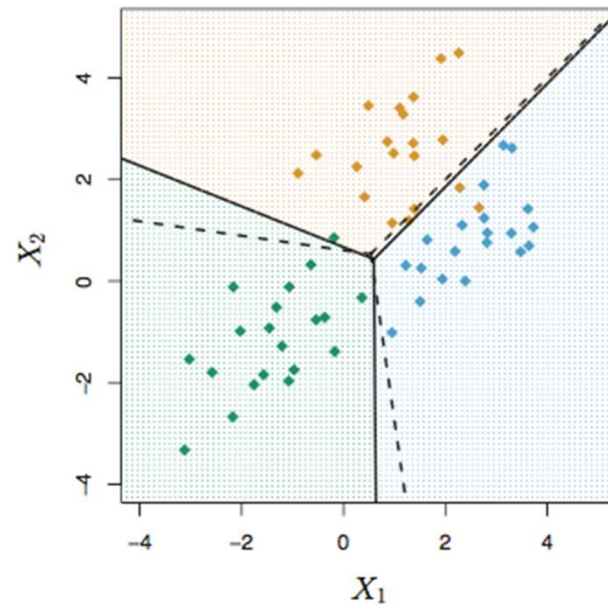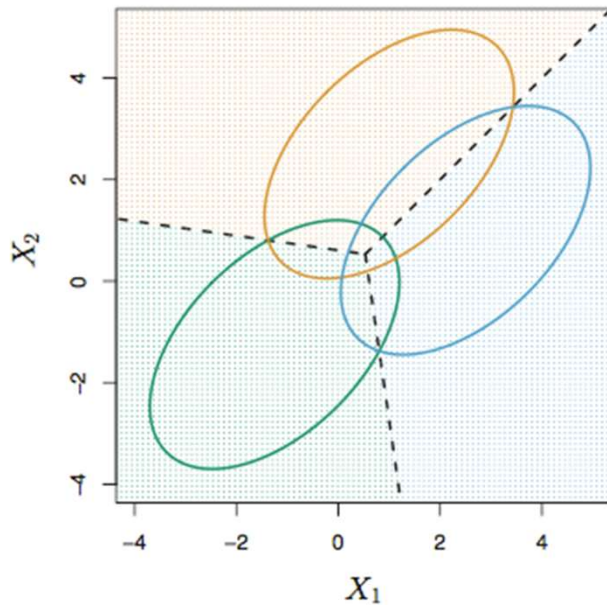
- Note: this is just the vector-matrix version of the formula we saw earlier in lecture:

$$\delta_k(x) = x\frac{\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2} + \log \pi_k$$

- What do we have to estimate now with the vector-matrix version?  How many parameters are there?

- There are $pK$ means, $pK$ variances, $K$ prior proportions, and $\binom{p}{2} = \frac{p(p-1)}{2}$ covariances to estimate.

# LDA when $K > 2$

- The linear discriminant nature of LDA still holds not only when $p > 1$, but also when $K > 2$ for that matter as well. A picture can be very illustrative:
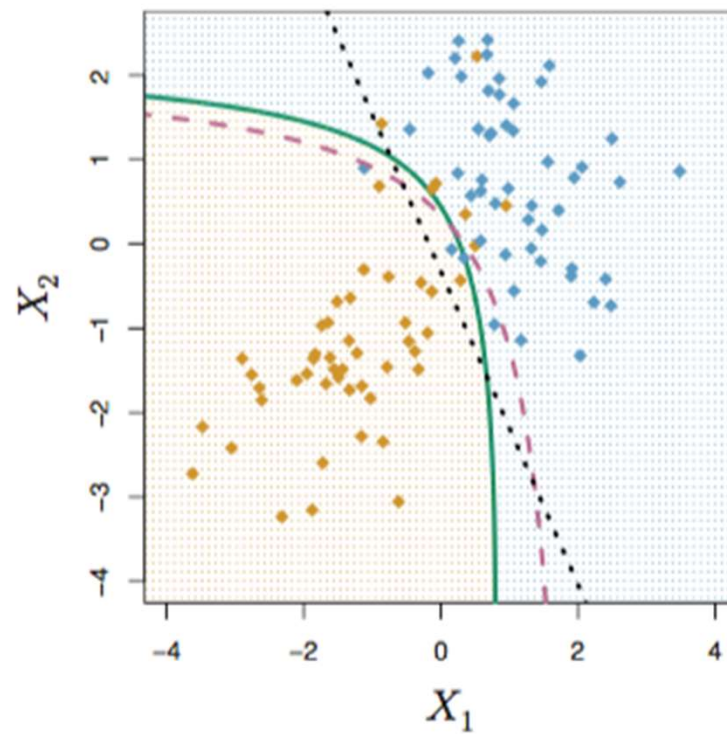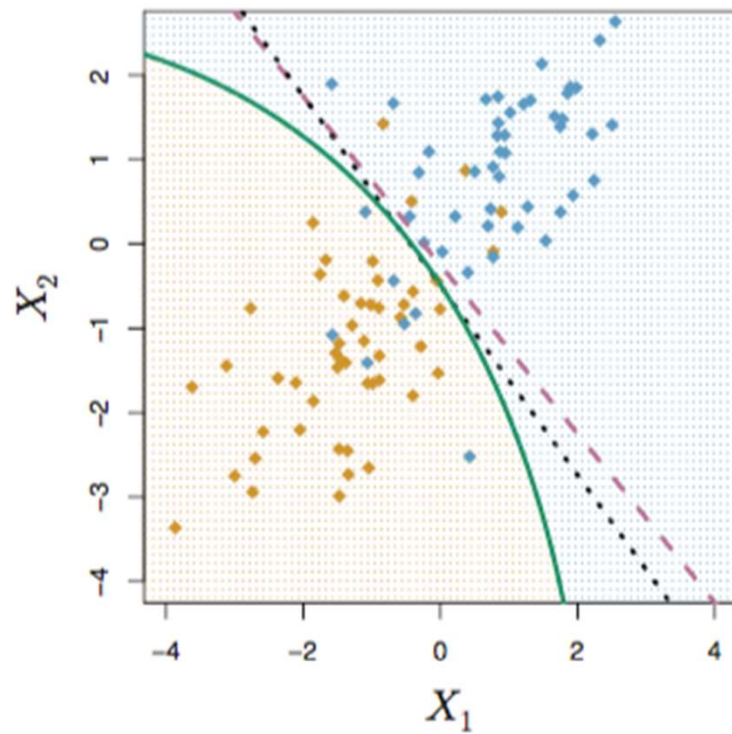
# Quadratic Discriminant Analysis (QDA)

- A generalization to linear discriminant analysis is quadratic discriminant analysis (QDA).

- Why do you suppose the choice in name?

- The implementation is just a slight variation on LDA. Instead of assuming the covariances of the MVN distributions within classes are equal, we instead allow them to be different.

- This relaxation of an assumption completely changes the picture...

# QDA in a picture

- A picture can be very illustrative:

# QDA (cont.)

- When performing QDA, performing classification for an observation based on its predictors $\vec{x}$ is equivalent to maximizing the following over the $K$ classes:

$$\delta_k(\vec{x}) = -\frac{1}{2}\vec{x}^T \Sigma_k^{-1} \vec{x} + \vec{x}^T \Sigma_k^{-1} \vec{\mu}_k - \frac{1}{2}\vec{\mu}_k^T \Sigma_k^{-1} \vec{\mu}_k - \frac{1}{2}\log|\Sigma_k| + \log \pi_k$$

- Notice the `quadratic form' of this expression. Hence the name QDA.

- Now how many parameters are there to be estimated?

- There are $pK$ means, $pK$ variances, $K$ prior proportions, and $\binom{p}{2}K = \left(\frac{p(p-1)}{2}\right)K$ covariances to estimate. This could slow us down very much if $K$ is large...

# Discriminant Analysis in Python

- LDA is already implemented in Python via the `sklearn.discriminant_analysis` package through the `LinearDiscriminantAnalysis` function.

- QDA is in the same package and is the `QuadraticDiscriminantAnalysis` function.

- It's very easy to use.  Let's see how this works

# QDA vs. LDA

- So both QDA and LDA take a similar approach to solving this classification problem: they use Bayes' rule to flip the conditional probability statement and assume observations within each class are multivariate Normal (MVN) distributed.

- QDA differs in that it does not assume a common covariance across classes for these MVNs.  What advantage does this have?  What disadvantage does this have?

# QDA vs. LDA (cont.)

- So generally speaking, when should QDA be used over LDA?  LDA over QDA?

- The extra covariance parameters that need to be estimated in QDA not only slow us down, but also allow for another opportunity for overfitting.  Thus if your training set is small, LDA should perform better for **out-of-sample prediction**, aka, predicting future observations.

# Quadratic Discriminant Analysis (QDA)

- We have seen 3 major methods for doing classification:

  - Logistic Regression

  - $k$-NN

  - Discriminant Analysis (LDA and QDA)

- For a specific problem, which approach should be used?


- Well of course, it depends on the nature of the data.  So how should we decide?