

Bayes Classifier and Linear Discriminant Analysis

Announcements

- All waitlisted students have been granted permission to enroll.
- HW 1 will go out today, and is due in one week.
- There will be a Python tutorial next week, stayed tuned for details.
- Please use Piazza for course related questions.
- Office hours have been posted to Canvas

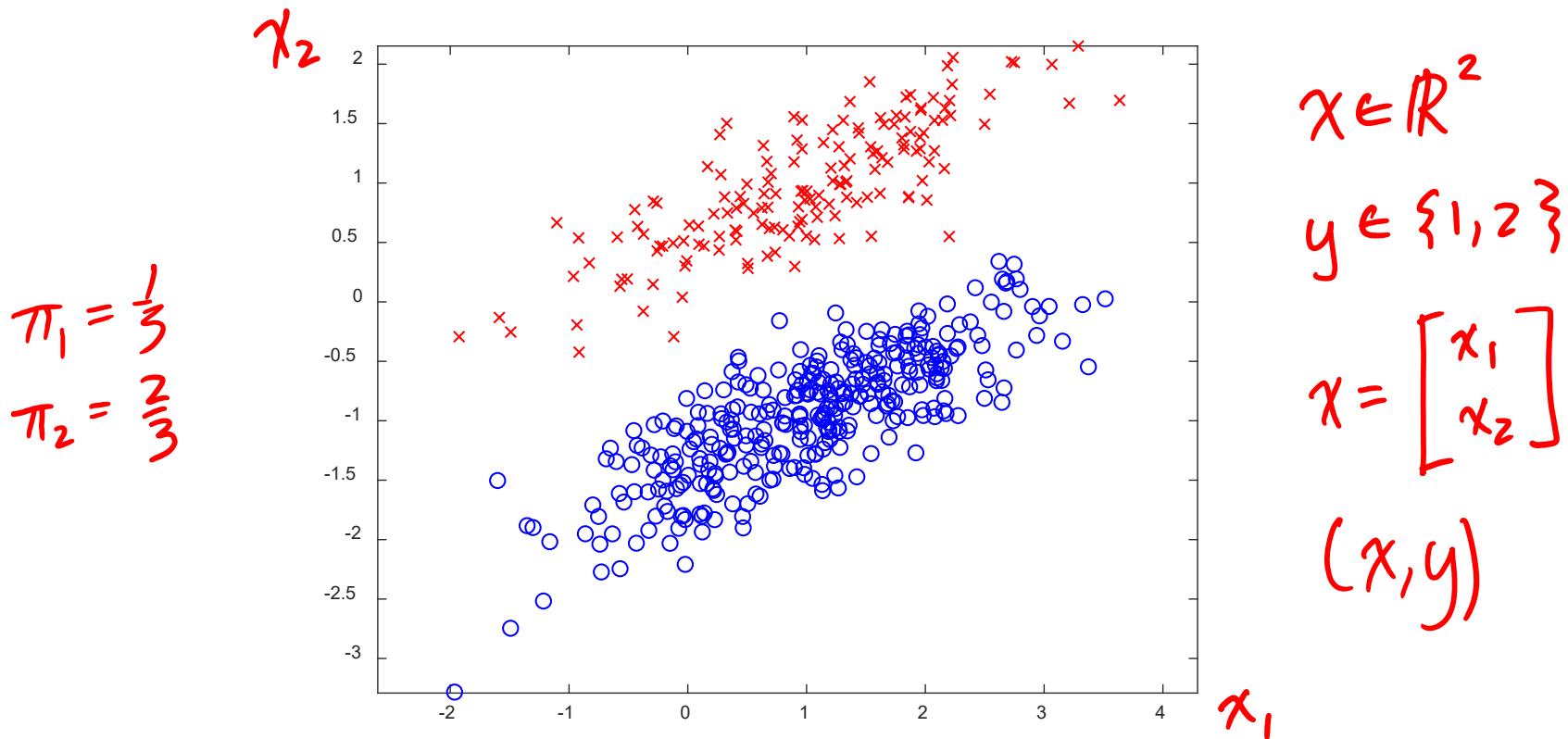
Classification: Probabilistic Setting

- Feature vector $\mathbf{X} \in \mathbb{R}^d$
- Label $Y \in \{1, \dots, K\}$
- Assume (\mathbf{X}, Y) are jointly distributed ($d + 1$ dimensional)

\uparrow Y is the label associated to X

$(x_1, y_1), \dots, (x_n, y_n)$ iid joint distrib. of (x, y)

Classification: Probabilistic Setting



- Two classes, each class has a Gaussian distribution
- Blue class is twice as probable as red class
- How can we characterize the joint distribution mathematically?

Classification: Probabilistic Setting

- How should we think about the joint distribution of (\mathbf{X}, Y) ?
- First way: Marginal distribution of Y , and conditional distribution of \mathbf{X} given $Y = y$, for each y
- Second way: Marginal distribution of \mathbf{X} , and conditional distribution of Y given $\mathbf{X} = \mathbf{x}$, for each \mathbf{x}

Classification: Probabilistic Setting

- Notation:

- 1st way*
- $\pi_k := \Pr(Y = k)$ class prior
 - $g_k(x) := \text{pdf/pmf of } X \text{ given } Y = k$ class-conditional pdf/pmf
- 2nd way*
- $\eta_k(x) := \Pr(Y = k | X = x)$ class posterior
 - $g(x) := \text{pdf/pmf of } X$ marginal distrib. of X

$$g(x) = \sum_{k=1}^K \pi_k g_k(x)$$

$$\eta_k(x) = \frac{\pi_k g_k(x)}{g(x)}$$

Classification: Probabilistic Setting

- First way: Marginal distribution of Y , and conditional distribution of \mathbf{X} given $Y = y$, for each y

red class : $y=1$

$$\mathbf{X} | Y=1 \sim N(\mu_1, \Sigma)$$

blue " : $y=2$

$$\mathbf{X} | Y=2 \sim N(\mu_2, \Sigma)$$

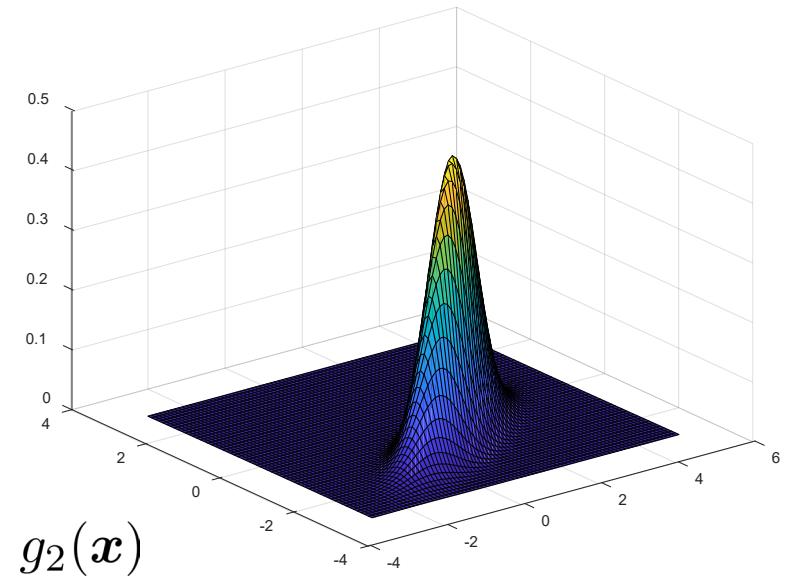
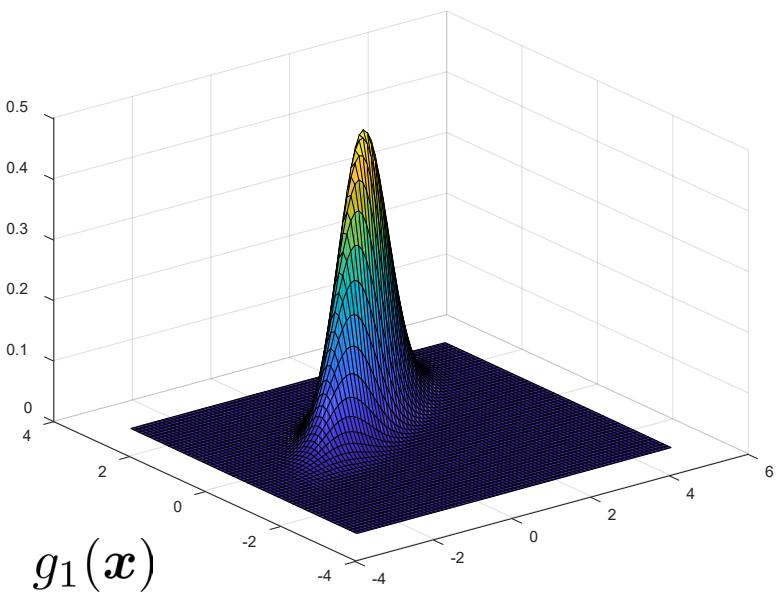
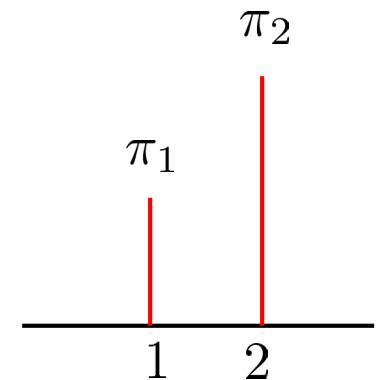
$$\pi_1 = \frac{1}{3}, \pi_2 = \frac{2}{3}$$

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

$$g_1(\mathbf{x}) = \phi(\mathbf{x}; \mu_1, \Sigma)$$

$$g_2(\mathbf{x}) = \phi(\mathbf{x}; \mu_2, \Sigma)$$

Classification: Probabilistic Setting



Classification: Probabilistic Setting

- Second way: Marginal distribution of \mathbf{X} , and conditional distribution of Y given $\mathbf{X} = \mathbf{x}$, for each \mathbf{x}

$$g(\mathbf{x}) = \pi_1 g_1(\mathbf{x}) + \pi_2 g_2(\mathbf{x})$$

$$= \frac{1}{3} \phi(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}) + \frac{2}{3} \phi(\mathbf{x}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma})$$

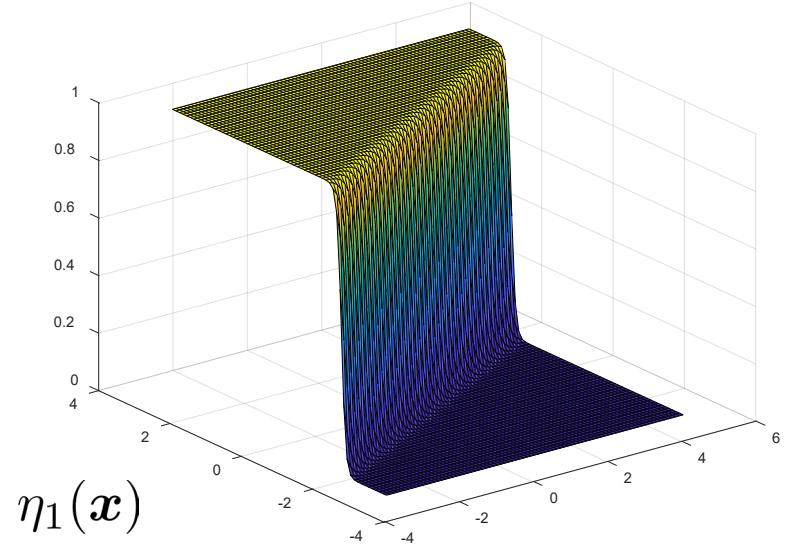
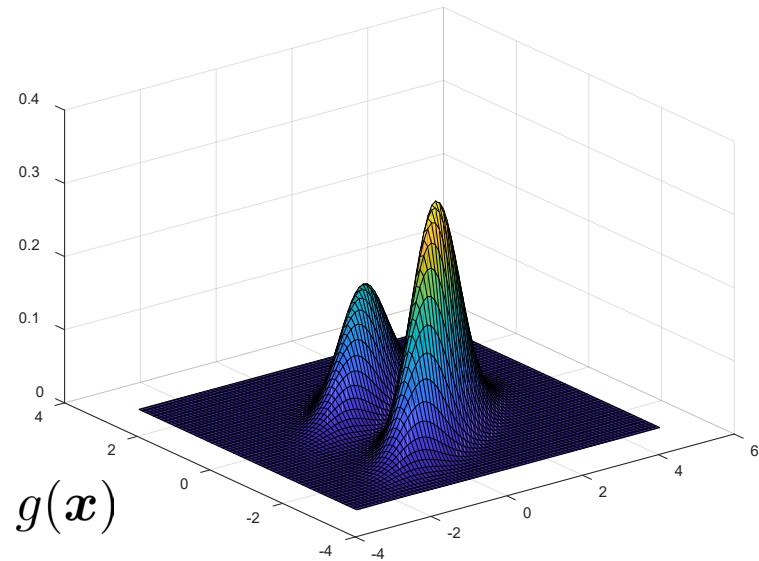
$$\gamma_1(\mathbf{x}) = \frac{\pi_1 g_1(\mathbf{x})}{\pi_1 g_1(\mathbf{x}) + \pi_2 g_2(\mathbf{x})}$$

$$\gamma_2(\mathbf{x}) = 1 - \gamma_1(\mathbf{x})$$

Classification: Probabilistic Setting

$$g(\mathbf{x}) = \pi_1 g_1(\mathbf{x}) + \pi_2 g_2(\mathbf{x})$$

$$\eta_1(\mathbf{x}) = \frac{\pi_1 g_1(\mathbf{x})}{\pi_1 g_1(\mathbf{x}) + \pi_2 g_2(\mathbf{x})}$$



Bayes Classifier

- A *classifier* is a function $f : \mathbb{R}^d \rightarrow \{1, \dots, K\}$
- Given a joint distribution of (X, Y) , what is the best possible classifier?
- The best classifier depends on the performance measure. The most common performance measure is the probability of error, or *risk*, defined by:

$$R(f) = \Pr(f(x) \neq y) \in [0, 1]$$

i.e., the probability of the event

$$\{(x, y) \in \mathbb{R}^d \times \{1, \dots, K\} \mid f(x) \neq y\}$$

sample space

- The smallest risk of any classifier is called the *Bayes risk* and is denoted R^*
- If $R(f) = R^*$, f is called a *Bayes classifier*

Bayes Classifier

$$\Pr(Y=k | X=x)$$

- **Theorem:** The classifier

$$\begin{aligned} f^*(x) &= \arg \max_{k \in \{1, \dots, K\}} \eta_k(x) \\ &= \arg \max_{k \in \{1, \dots, K\}} \pi_k g_k(x) \end{aligned}$$

is a Bayes classifier.

The pmf of $Y | X=x$ is $\eta_k(+)$

$$\eta_k(x) = \frac{\pi_k g_k(x)}{g(x)}$$

Poll

True or false: A Bayes classifier always predicts the correct label

- (A) True
- (B) False

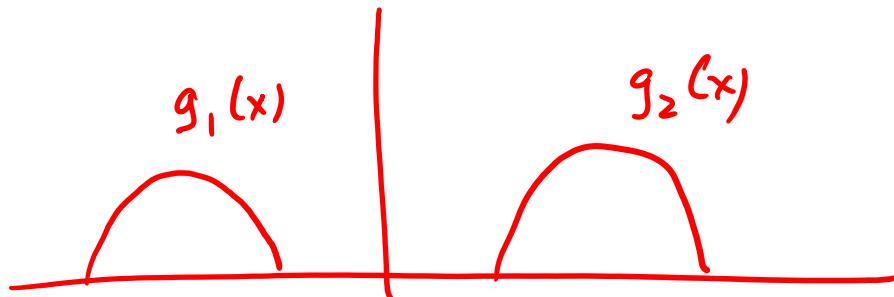
Questions

1. Give an example of a classification problem where $R^* = 0$.
2. Give an example of a classification problem where $R^* = \frac{1}{2}$.
3. What is R^* when $f(\mathbf{x})$ always predicts class k ?

1. $\pi_1 = 1, \pi_2 = 0$

or

$g_1(x)$ and $g_2(x)$ don't overlap



2. $g_1(x) = g_2(x) \forall x$

$\eta_1(x) = \frac{1}{2} \forall x$

3. $1 - \pi_k$

Plug-In Classifiers

- In machine learning, the quantities $\pi_k, g_k(\mathbf{x}), \eta_k(\mathbf{x})$, are not known, so we can't know the Bayes' classifier.
- However, the formula for the Bayes' classifier is still useful. We can estimate the quantities in the formula from training data

$$(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n),$$

and plug those estimates in to the formula to get a classifier.

- Linear discriminant analysis and Naïve Bayes have the form

$$f(\mathbf{x}) := \arg \max_k \widehat{\pi}_k \widehat{g}_k(\mathbf{x})$$

- Logistic regression has the form

$$f(\mathbf{x}) := \arg \max_k \widehat{\eta}_k(\mathbf{x})$$

- Other estimators are possible.

Linear Discriminant Analysis

- Training data $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$.

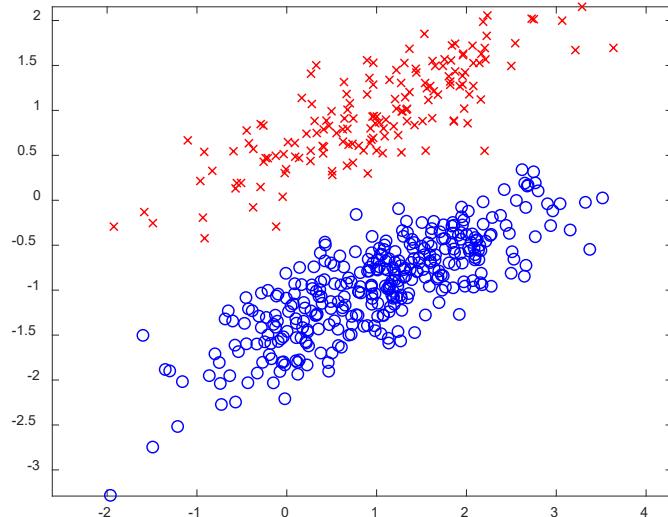
- *LDA assumption:*

$$\begin{aligned} y_1, \dots, y_n &\stackrel{iid}{\sim} \text{discrete}(\pi_1, \dots, \pi_K) \\ \mathbf{x} \mid y = k &\sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}), \quad k = 1, \dots, K \end{aligned}$$

for some unknown $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K$ and $\boldsymbol{\Sigma}$. Equivalently

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

- LDA is the plug-in rule based on this model. We use training data to estimate $\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}$.



LDA Estimates

- LDA is the classifier obtained by plugging the following into the Bayes classifier formula:

$$\hat{\pi}_k = \frac{n_k}{n}, \quad n_k = |\{i : y_i = k\}|$$

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

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \hat{\mu}_{y_i})(\mathbf{x}_i - \hat{\mu}_{y_i})^T$$

$$\begin{aligned} f(\mathbf{x}) &= \underset{k}{\operatorname{argmax}} \quad \hat{\pi}_k \cdot \hat{g}_k(\mathbf{x}) \\ &= \underset{k}{\operatorname{argmax}} \quad \hat{\pi}_k \cdot \phi(\mathbf{x}; \hat{\mu}_k, \hat{\Sigma}) \end{aligned}$$

- These estimates are all *maximum likelihood estimates*
- $\hat{\mu}_k$ is the *sample mean* for each class
- Alternatively, $\hat{\Sigma}$ can be the *pooled sample covariance*,

$$\hat{\Sigma} = \frac{1}{n-d} \sum_{i=1}^n (\mathbf{x}_i - \hat{\mu}_{y_i})(\mathbf{x}_i - \hat{\mu}_{y_i})^T$$

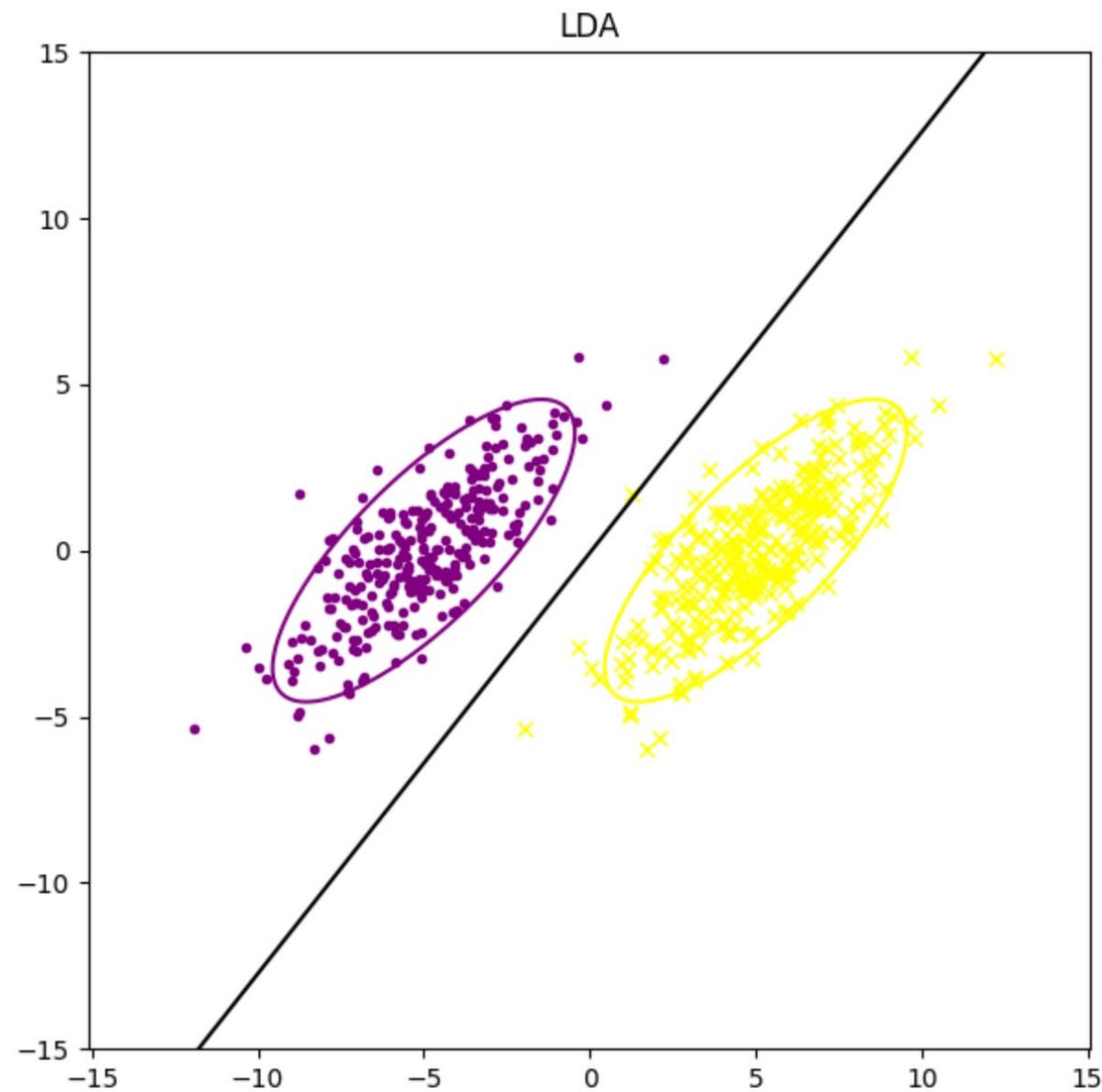
LDA is a Linear Classifier

- Binary setting, $Y \in \{-1, 1\}$.
- The LDA classifier is

$$\hat{f}(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

where

$$\begin{aligned}\mathbf{w} &= \widehat{\boldsymbol{\Sigma}}^{-1}(\widehat{\boldsymbol{\mu}}_1 - \widehat{\boldsymbol{\mu}}_{-1}) \\ b &= \log(\widehat{\pi}_1/\widehat{\pi}_{-1}) + \frac{1}{2}(\widehat{\boldsymbol{\mu}}_{-1}^T \widehat{\boldsymbol{\Sigma}}^{-1} \widehat{\boldsymbol{\mu}}_{-1} - \widehat{\boldsymbol{\mu}}_1^T \widehat{\boldsymbol{\Sigma}}^{-1} \widehat{\boldsymbol{\mu}}_1).\end{aligned}$$



Generative vs Discriminative

- A machine learning algorithm is *generative* if it is based on a full probabilistic model for the data, in other words, a model that could be used to generate data (without needing to have seen data previously).
- A machine learning algorithm is *discriminative* if it is not generative.
- LDA is *generative*

Parametric vs. Nonparametric

A learning algorithm is *nonparametric* if the amount of space needed to store the learned model (the case of LDA, a classifier) grows with n , otherwise it is *parametric*.

Poll: Is LDA parametric or nonparametric?

- A. parametric
- B. nonparametric

LDA Summary

- One of the very earliest ML algorithms invented (c. 1930s)
- Many applications throughout history, although better methods now available
- Requires $n > d$ although extensions are possible
- Logistic regression almost always preferable (makes a weaker assumption)
- Optional: read about Mahalanobis distance in lecture notes