

LECTURE 4 - 17/09/24

STRATEGIES FOR BUILDING A CLASSIFIER

- So far our models have approximated posterior class prob:

$$P(C_k | x) \text{ with } C_1, C_2, \dots, C_K \text{ & w/o } x$$

↳ This is because we have used **DISCRIMINATIVE** models to classify data.

↳ They model posterior class probabilities directly from data.

SO FAR THE DISCRIMINATIVE MODELS:

- Logistic regression

$$P(Y = C_1 | x) = \frac{e^{\beta_0 + \beta_1 x}}{1 + e^{\beta_0 + \beta_1 x}}$$

- K-nearest neighbours

$$P(Y = C_k | x) = \frac{1}{K} \sum_{i \in N_0} I(y_i = C_k).$$

GENERATIVE MODELS FOR CLASSIFICATION

- Another approach would be to take a step back.

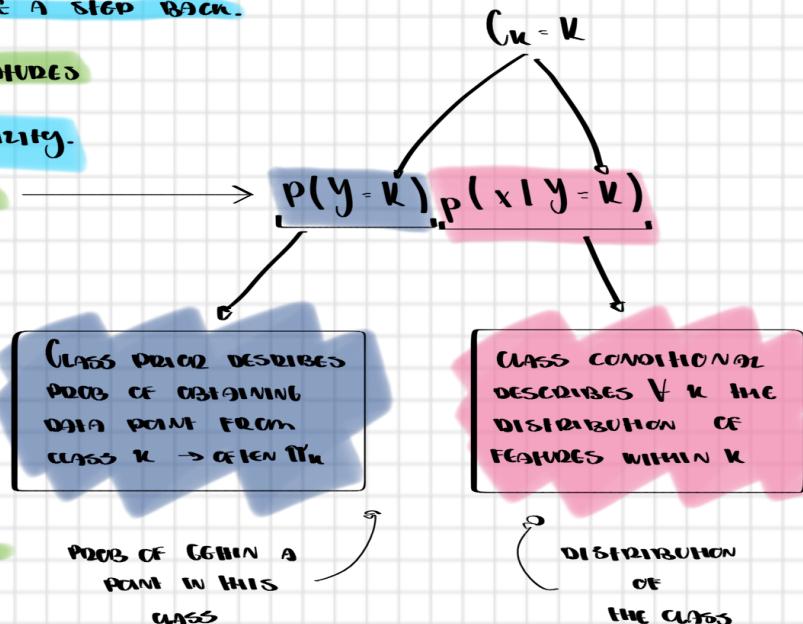
- Model joint distribution of features

↳ Gives also posterior probability.

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

- This generative model will then give us the posterior class distribution that we use for classification in discriminative models

- Generative model also gives estimates of joint dist. $\rightarrow P(x, y)$ can be used for simulating data



- These are then the strategies:

- Model the full joint distribution (generative models)

- Gives posterior probabilities

- Model class posterior probabilities (discriminative model)

- Posterior probs used as discriminant functions

- Model directly a discriminant function

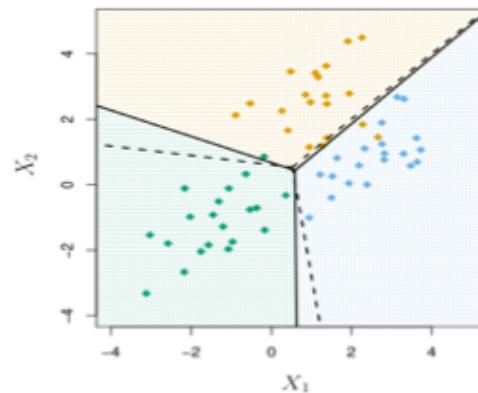
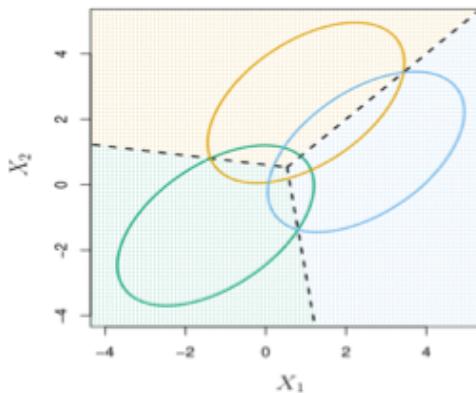
- Remember $d(x) \rightarrow$ maps input into prediction

Hierarchy!!!

Number ② refers to \rightarrow Logistic Regression, K-nearest neighbors.

Today's lecture \rightarrow ①

Linear Discriminant Analysis



So... What do we do in LDA and how do we do it?

- To build our classifier we will follow the hierarchy described before
- For linear Discriminant Analysis:
 - we model the joint
 - end up with discriminant functions, $g_k(x)$
- For LDA we assume that class conditionals are:
 - Gaussian with a class-specific mean, common var

→ FEW POINTS...

- want to classify x with highest discriminant $g_k(x)$

→ For Bayes:

- choose posterior probs
- joint dist.
- log scale

Let's build Bayes classifier:

- We can derive it under the following assumption:

- we know the true distribution of data.

- Gaussian class conditionals:

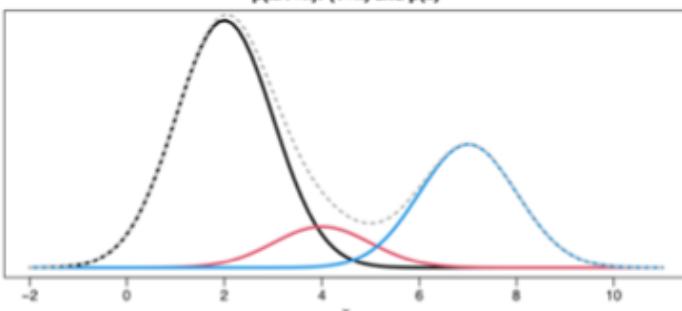
- we assume feature x to have a univariate Gaussian distribution:

$$p(x|Y=\text{BLACK}) = N(2, 1) \quad \rightarrow \pi_{\text{BLACK}} = 0.6$$

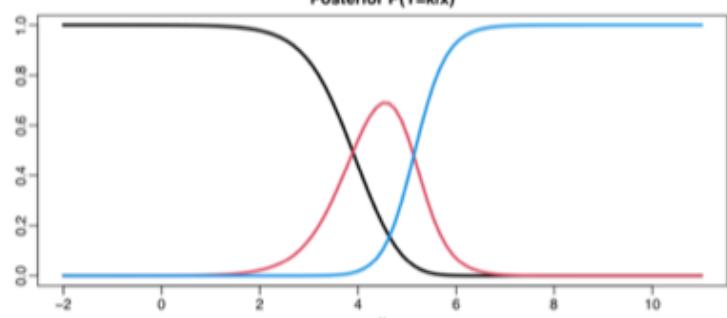
$$p(x|Y=\text{RED}) = N(4, 1) \quad \rightarrow \pi_{\text{RED}} = 0.1$$

$$p(x|Y=\text{BLUE}) = N(7, 1) \quad \rightarrow \pi_{\text{BLUE}} = 0.3$$

$p(x|Y=k)p(Y=k)$ and $p(x)$



Posterior $P(Y=k|x)$



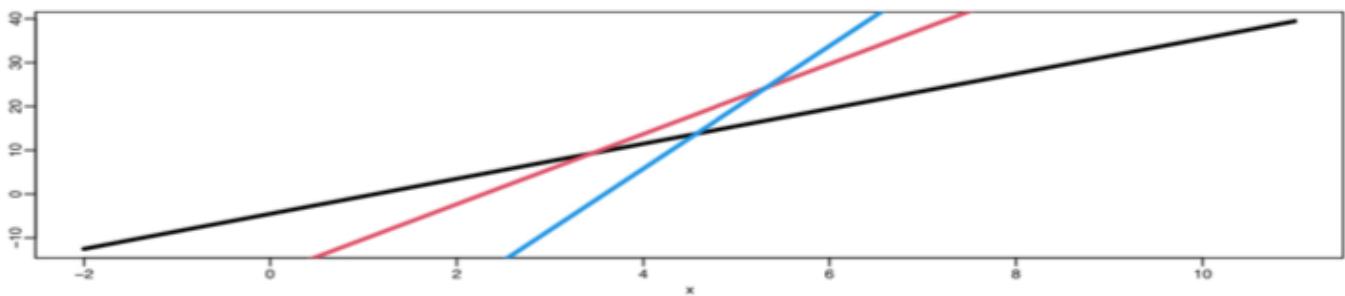
DISCRIMINANT FOR JDA WITH ONE FEATURE:

- Look at the log of the joint distribution,

$$\begin{aligned}
 \log p(x, y) &= \log P(y=k) + \log p(x|y=k) \\
 &= \log \pi_k + \log \left\{ \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu_k)^2}{2\sigma^2}} \right\} \text{ (Our gaussian assumption)} \\
 &= \log \pi_k + \log \frac{1}{\sqrt{2\pi\sigma^2}} + \log e^{-\frac{(x-\mu_k)^2}{2\sigma^2}} \\
 &= \log \pi_k + \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{(x-\mu_k)^2}{2\sigma^2} \\
 &= \log \pi_k + \log \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{x^2 - 2\mu_k x + \mu_k^2}{2\sigma^2} \\
 &= \log \pi_k + \underbrace{\log \frac{1}{\sqrt{2\pi\sigma^2}}}_{\text{same for all } k} - \frac{x^2}{2\sigma^2} + \frac{2\mu_k x}{2\sigma^2} - \frac{\mu_k^2}{2\sigma^2}
 \end{aligned}$$

- This means that the Bayes classifier chooses the class with the highest

$$g_k(x) = \frac{\mu_k}{\sigma^2} x - \frac{\mu_k^2}{2\sigma^2} + \log(\pi_k) \rightarrow \text{LINEAR IN } x$$



(Note that we modelled the joint distribution as Gaussian (class-specific mean, common variance). This makes this a generative model).

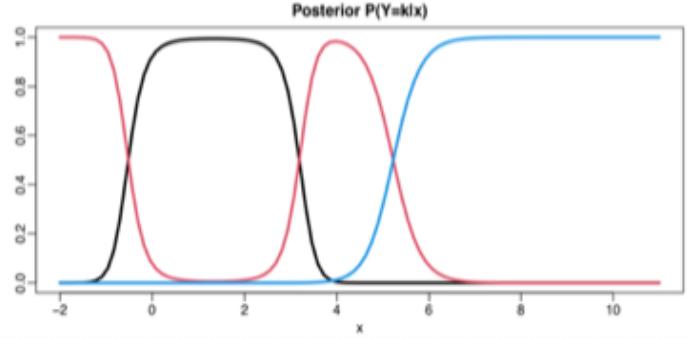
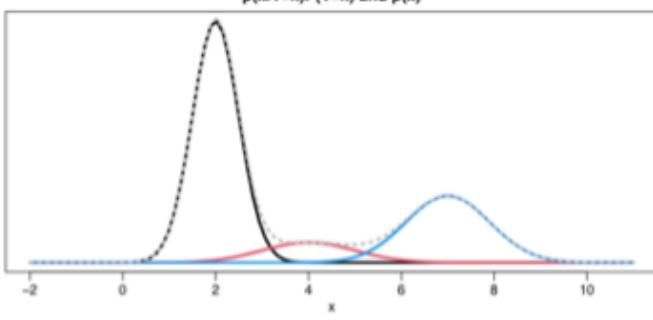
Quadratic Discriminant Analysis (QDA)

- The assumption of a common variance in JDA may be too strict
- In QDA we assume class conditionals are:
 - Gaussian with both class-specific means & class-specific variance
- QDA introducing class-specific variances:

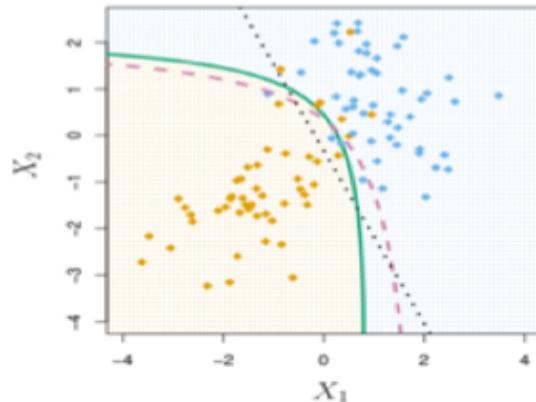
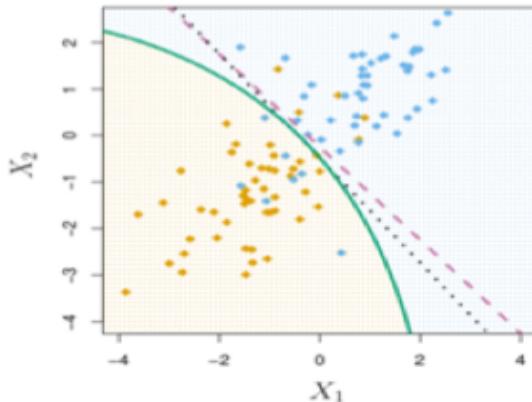
$$p(x | y=\text{black}) = N(2, 0.25) \rightarrow \sigma_{\text{black}}^2 = 0.6$$

$$p(x | y=\text{red}) = N(4, 1) \rightarrow \sigma_{\text{red}}^2 = 0.1$$

$$p(x | y=\text{blue}) = N(7, 0.81) \rightarrow \sigma_{\text{blue}}^2 = 0.3$$



LDA vs QDA



Purple: Bayes boundary, black: LDA, green: QDA.

QDA gives more flexibility (variance!). LDA vs QDA is a bias-variance tradeoff.

- LDA & QDA ARE PLUG-IN CLASSIFIERS
 - THE "GOAL STANDARD" Bayes classifier is derived:
 - By assuming class conditionals are GAUSSIAN
 - By assuming we know all their PARAMS (MEAN, VARIANCE)
 - In practice though:
 - CLASS CONDITIONALS ARE NOT TRULY GAUSSIAN
 - VARIANCE & MEAN ARE UNKNOWN.
 - TRUE CLASS PRIORS ARE UNKNOWN.
- LDA & QDA APPROXIMATE THE Bayes classifier

Usually we do not know model parameters:

- For class priors $\hat{P}(k)$ we use MLE, the empirical frequencies in training set:

$$\hat{P}(k) = \frac{n_k}{n} \quad \text{# of observations in } k$$

- For means we typically use MAX LIKELIHOOD ESTIMATES:

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

- For variances we have different methods for LDA & QDA:

LDA:

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

QDA:

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