

# Gaussian Discriminant Analysis

Fundamental assumption: each class comes from a normal/gaussian distribution

$$X \sim N(\mu, \sigma^2): f(\vec{x}) = \frac{1}{(\sqrt{2\pi})^d} e^{-\frac{1}{2\sigma^2} \cdot \|\vec{x} - \vec{\mu}\|^2}$$

(multivariate form, Variance is same in every direction)  
Isotropic

- For each class  $C$ , Suppose we estimate  $\mu_C, \sigma_C^2$ , prior  $\pi_C = P(Y=C)$ .

- then, given  $x$ , find Bayes decision rule  $r^*(x)$  that predicts class  $C$  that maximizes posterior probability which is the same as maximizing  $f(X=x|Y=C)\pi_C \cdot \ell(C, y)$

Want to get rid of exponential:

observe that  $\ln(\cdot)$  is monotonically increasing  $\forall w > 0$

$\Rightarrow$  Same as maximizing:

$$Q_C(x) = \ln((\sqrt{2\pi})^d f_C(x) \pi_C) = \underbrace{\frac{-\|x - \mu_C\|^2}{2\sigma_C^2}}_{\text{quadratic in } x} - d \ln \sigma_C + \ln \pi_C$$

↑  
Normal PDF, estimates  
 $f(X=x|Y=C)$ 
↑  
only need to  
calculate constants  
once

if  $S$  classes, compute all  $S$   $Q_C(x)$ , pick biggest probability (most likely class)

Quadratic Discriminant Analysis (QDAS):

Suppose only classes  $C$  &  $D$ , then

$$r^*(x) = \begin{cases} C & \text{if } Q_C(x) > Q_D(x) \\ D & \text{o/w} \end{cases}$$

quadratic decision fn.

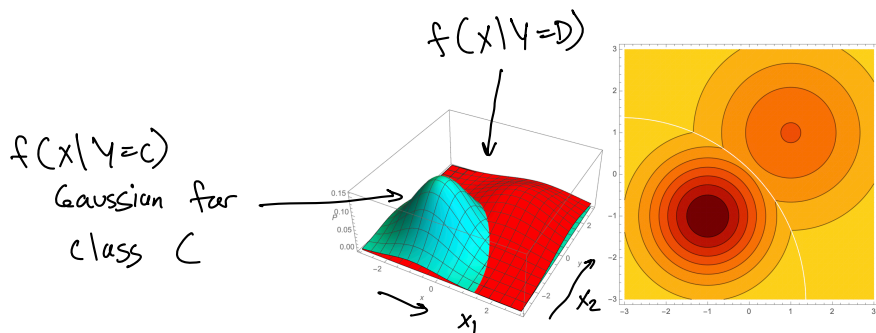
Same answer as computing posterior probabilities

Decision fn is quadratic in  $x$

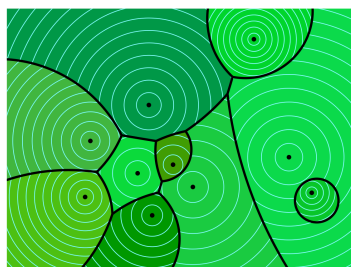
$\rightarrow$  Bayes decision boundary =  $\{x: Q_C(x) > Q_D(x)\}$

- In 1D, B.d.b may have 1 or 2 points

- In d-D, B.d.b is a quadric



multiple classes,  
each black dot is a  
mean of a gaussian  
distribution



What if we want posterior probabilities?

To get posterior probabilities in 2-class case,

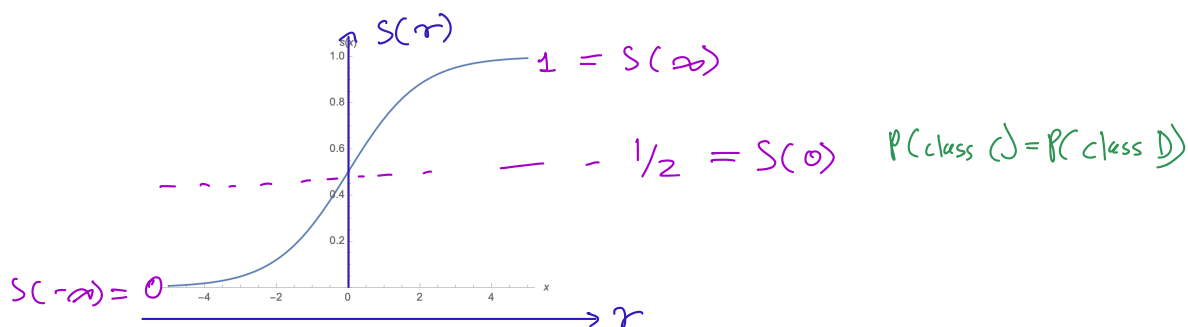
- apply Bayes thm to find  $P(Y=c|X) = \frac{f(X|Y=c) \pi_c}{f(X|Y=c) \pi_c + f(X|Y=D) \pi_D}$

need to relate  $f(X|Y=c) \pi_c$  back to  $Q_c(x)$ .

- recall  $e^{Q_c(x)} = (\sqrt{2\pi})^d f_c(x) \pi_c$

$$\Rightarrow P(Y=c|X) = \frac{e^{Q_c(x)}}{e^{Q_c(x)} + e^{Q_D(x)}} = \frac{1}{1 + e^{Q_D(x) - Q_c(x)}}$$

SIGMOID fn:  $= S(Q_c(x) - Q_D(x))$ , where  $S(r) = \frac{1}{1 + e^{-r}}$



Linear Discriminant Analysis (LDA):

- every class has same Variance/Covariance  $\leftarrow$  key assumption  $\sigma_c = \sigma_D$
- $\hookrightarrow$  always linear decision boundaries
- $\hookrightarrow$  less likely to overfit

$\rightarrow Q_c(x) - Q_D(x) = \underbrace{\frac{(\mu_c - \mu_D) \cdot x}{\sigma^2}}_{W \cdot x} - \underbrace{\frac{\|\mu_c\|^2 - \|\mu_D\|^2}{2\sigma^2} + \ln \pi_c + \ln \pi_D}_{d}$

*only dependence on x, linear term*

If multiple classes, compute LD fn for each class then choose C that maximizes those LD fns.

LD fn:  $\frac{\mu_c \cdot x}{\sigma^2} - \frac{\|\mu_c\|^2}{2\sigma^2} + \ln \pi_c$

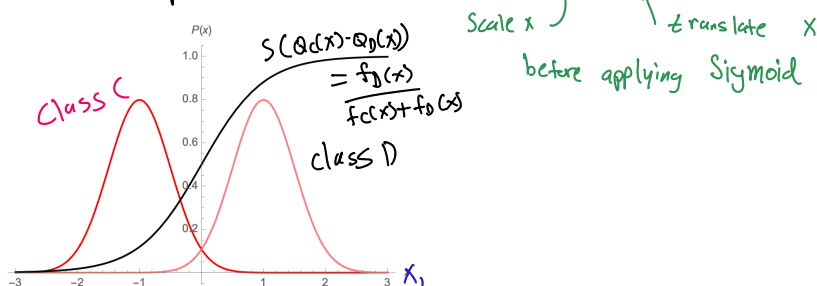
$\swarrow$  Same as doing  $Q_c(x) - Q_D(x)$  for every possible pair

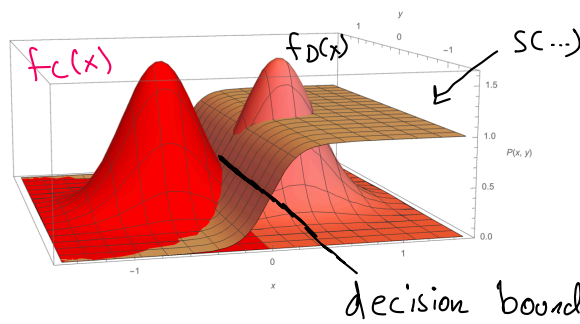
e.g. MNIST, calculate  $Q_0(x) \dots Q_9(x)$  and pick biggest

In 2-class case: - decision boundary =  $W \cdot x + d = 0$

$\rightarrow W = \frac{\mu_c - \mu_D}{\sigma^2}$ , d above

- posterior prob  $P(Y=c|X=x) = S(W \cdot x + d)$





If  $\pi_C = \pi_D = \frac{1}{2}$ , then  $\underbrace{(\mu_C - \mu_D) \cdot x - (\mu_C - \mu_D) \cdot \frac{(\mu_C + \mu_D)}{2}}_{\text{decision boundary}} = 0$

This is the **centroid method**

## MAXIMUM LIKELIHOOD estimation of parameters

(Ronald Fisher, 1912)

Let's flip a coin  $\sim \text{Bernoulli}(p)$

10 flips, 8 heads, 2 tails

What value of  $p$  gives this outcome?

Let  $X = \# \text{ heads}$ ,  $X \sim \text{Binom}(n, p)$

$$P(X=x) = \binom{n}{x} p^x (1-p)^{n-x}$$

$$\Rightarrow P(X=8) = \binom{10}{8} p^8 (1-p)^2 = 45 p^8 (1-p)^2 = \mathcal{L}(p) \quad (\text{likelihood fn})$$

find  $p$  that maximizes this

Probability of 8 heads in 10 flips:

Written as  $\mathcal{L}(p)$  of distr. params, this is the **likelihood func.**

**Maximum Likelihood estimation (MLE):** method of estimating params of a statistical distribution by picking params that maximize  $\mathcal{L}$

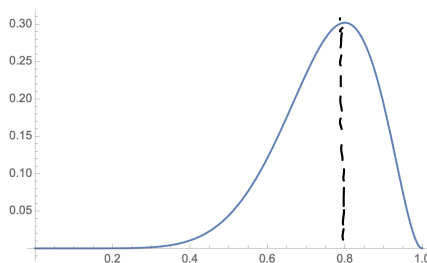
— one method of **density estimation**: estimating PDF from data

→ Find  $p$  that maximizes  $\mathcal{L}(p)$ :

Solve by setting  $\frac{d}{dp}$  to 0:

$$\frac{d}{dp} \mathcal{L}(p) = 360 p^7 (1-p)^2 - 90 p^8 (1-p) = 0$$

$$\Rightarrow 4(1-p) \cdot p = 0 \Rightarrow p = 0.8$$



$\hat{\pi}_C = 0.8$  = chance a newly sampled point is in class C

MLE used on a Gaussian:

given  $x_1 \dots x_n$ , find best-fit Gaussian

Since this is continuous, the likelihood is no longer a probability

regardless, likelihood of generating these points is:

$$\mathcal{L}(\mu, \sigma; x_1, \dots, x_n) = f(x_1) \cdot f(x_2) \cdot \dots \cdot f(x_n) \quad (\text{intersection})$$

choose  $\mu$  &  $\sigma$  that maximize  $\mathcal{L}$

take log to get a sum:

log likelihood:  $\ell(\mu, \sigma; x_1, \dots, x_n)$  is natural log of  $\mathcal{L}(\cdot)$

maximizing likelihood  $\Leftrightarrow$  maximizing log likelihood

$$\begin{aligned} \rightarrow \ell(\cdot) &= \ln f(x_1) + \dots + \ln f(x_n) \\ &= \sum_{i=1}^n \left( \underbrace{\frac{-\|x_i - \mu\|^2}{2\sigma^2} - \frac{1}{2} \ln \sqrt{2\pi} - \frac{1}{2} \ln \sigma}_{\text{ln of normal PDF}} \right) \end{aligned}$$

Want to set  $\nabla_{\mu} \ell = 0$ ,  $\frac{\partial \ell}{\partial \sigma} = 0$  to find critical point

$$\nabla_{\mu} \ell = \sum_{i=1}^n \frac{x_i - \mu}{\sigma^2} = 0 \Rightarrow \hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i \quad (\text{average of samples is good est. for } \mu \text{ of distribution})$$

$$\frac{\partial \ell}{\partial \sigma} = \sum_{i=1}^n \frac{\|x_i - \mu\|^2 - \sigma^2}{\sigma^3} \Rightarrow \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n \|x_i - \hat{\mu}\|^2$$

mean estimate of true  $\mu$

Sample Variance

don't know exact  $\mu$