Generally build a classifier w/ discriminant functions $\delta_c$ for each of my classes $c = 1, \ldots, K$ so that

$$\delta_c(x) = \text{large if } x \text{ signifies that the assoc. } y \text{ is likely to be in class } c$$

So

$$\hat{y} = \arg\max_c \delta_c(x)$$

we said a Bayes' classifier

$$\delta_c(x) = P(Y=c \mid X=x)$$

- - - - - - - - -

Defn: A classifier is linear if some (increasing) function of the $\delta_c$s is linear.

$\Rightarrow$ In this case the decision boundaries btwn classes are lines.

$\underline{\text{Linear Discriminant Analysis}}$ (linear classifier)

$$\delta_c(x) = P(Y=c \mid X=x)$$

$$\propto \underbrace{P(X=x \mid Y=c)}_{N(\mu_c, \sigma^2)} \underbrace{P(Y=c)}_{\pi_c}$$

Shared last time w/ algebra

$$\delta_c(x) = x\, \frac{\mu_c}{\sigma^2} - \frac{\mu_c^2}{2\sigma^2} + \log(\pi_c)$$

$$\delta_c(x) = x \frac{\mu_c}{\sigma^2} - \frac{\mu_c^2}{2\sigma^2} + \log(\pi_c)$$

$x \in \mathbb{R}$    $\beta_c$    $\beta_{0c}$

unknown parameters $\mu_c$, $\pi_c$, $\sigma^2$

we estimate these using training data

$\hat{\mu}_c$ = mean of $x_n$s in class $c$
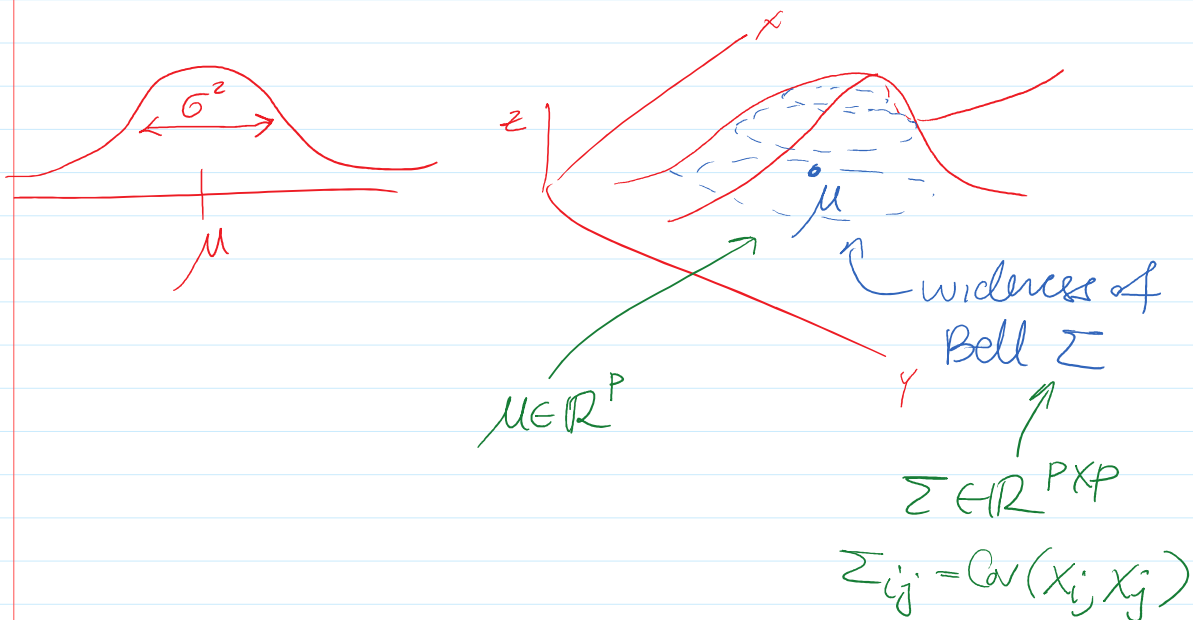
$\hat{\pi}_c$ = % of $x_n$s in class $c$

$\hat{\sigma}^2$ = pooled variance

   = calc. sample var. separately for each class and then take weighted avg.

We can expand to when $x \in \mathbb{R}^P$ using multivariate normal distribution

## Univariate        Multivariate



$\mu \in \mathbb{R}^P$

wideness of Bell $\Sigma$

$\Sigma \in \mathbb{R}^{P \times P}$

$\Sigma_{ij} = \text{Cov}(x_i, x_j)$

## PDFs

### Univariate:

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$$

Univariate:
$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right)$$

Multivariate: $(2\pi)^{-1/2}(\sigma^2)^{-1/2} - \frac{1}{2}(x-\mu)(\sigma^2)^{-1}(x-\mu)$

$$f(x) = (2\pi)^{-P/2} \det(\Sigma)^{-1/2} \exp\left(-\frac{1}{2}\underbrace{(x-\mu)^T}_{1\times P}\underbrace{\Sigma^{-1}}_{P\times P}\underbrace{(x-\mu)}_{P\times 1}\right)$$

$$\underbrace{\phantom{(x-\mu)^T\Sigma^{-1}(x-\mu)}}_{|X|}$$

$$\delta_c(x) = \underbrace{P(X=x|Y=c)}_{N(\mu_c, \Sigma)}\underbrace{P(Y=c)}_{\pi_c}$$

$$\mu_c \in \mathbb{R}^{P\times 1} \quad \Sigma \in \mathbb{R}^{P\times P}$$

$$= (2\pi)^{-P/2} \det(\Sigma)^{-1/2} \exp\left(-\frac{1}{2}(x-\mu_c)^T \Sigma^{-1}(x-\mu_c)\right) \pi_c$$

$$\delta_c(x) \leftarrow \log \delta_c(x) = -\frac{1}{2}(x-\mu_c)^T \Sigma^{-1}(x-\mu_c) + \log \pi_c$$

$$= -\frac{1}{2}(x^T\Sigma^{-1} - \mu_c^T\Sigma^{-1})(x-\mu_c) + \log \pi_c$$

$$= -\frac{1}{2}(x^T\Sigma^{-1}x - x^T\Sigma^{-1}\mu_c + \mu_c^T\Sigma^{-1}\mu_c$$

$$- \mu_c^T\Sigma^{-1}x) + \log \pi_c$$

$$\boxed{\delta_c(x) = x^T\Sigma^{-1}\mu_c - \frac{1}{2}\mu_c^T\Sigma^{-1}\mu_c + \log \pi_c}$$

$$\underbrace{\phantom{x^T\Sigma^{-1}\mu_c}}_{\beta} \qquad \underbrace{\phantom{-\frac{1}{2}\mu_c^T\Sigma^{-1}\mu_c + \log \pi_c}}_{\beta_0}$$

$$= x^T\beta + \beta_0$$
$$\underset{\text{linear!}}{\uparrow}$$

So even in the multivariable $(P > 1)$ this

So even in the multivariable $(P > 1)$ this
is a linear classifier.



How do I estimate $\hat{\mu}_c$ and $\hat{\Sigma}$?

$$\hat{\mu}_c = \text{mean vector of } X_n\text{'s in class } c$$

$$\hat{\Sigma}_{ij} = \text{pooled covariance btwn } X_{ni}\text{'s and } X_{nj}\text{'s.}$$

Just like in regression I can always create
more dimensions w/ transformations (grow P)
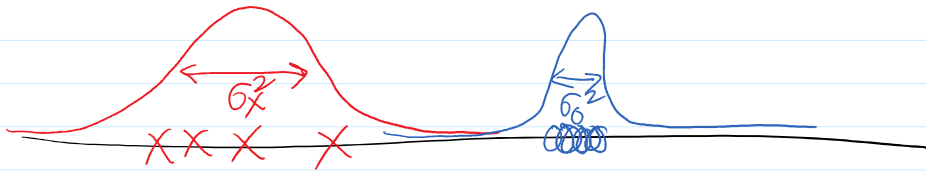
$$X = \begin{bmatrix} | & | & | \\ X & X^2 & \log X \\ | & | & | \end{bmatrix}$$

<u>Quadratic Discriminant Analysis</u> (not linear)

<u>LDA</u>: assumes equal variances $(\hat{\mu}_c, \hat{\pi}_c, \Sigma)$
<u>QDA</u>: relax this, classes can have diffent vars.
$$(\hat{\mu}_c, \hat{\pi}_c, \Sigma_c)$$

## Discr. fn for QDA:

$$\delta_c(x) = \underbrace{P(X=x \mid Y=c)}_{N(\mu_c, \Sigma_c)} \underbrace{P(Y=c)}_{\pi_c}$$

### For $p=1$

quadratic

$$\delta_c(x) = -\log \sigma_c - \frac{(x-\mu_c)^2}{2\sigma_c^2} + \log(\pi_c)$$

### For $p>1$

$$\delta_c(x) = -\log \det \Sigma_c - \frac{1}{2}(x-\mu_c)^T \Sigma^{-1}(x-\mu_c) + \log \pi_c$$

### LDA v. QDA

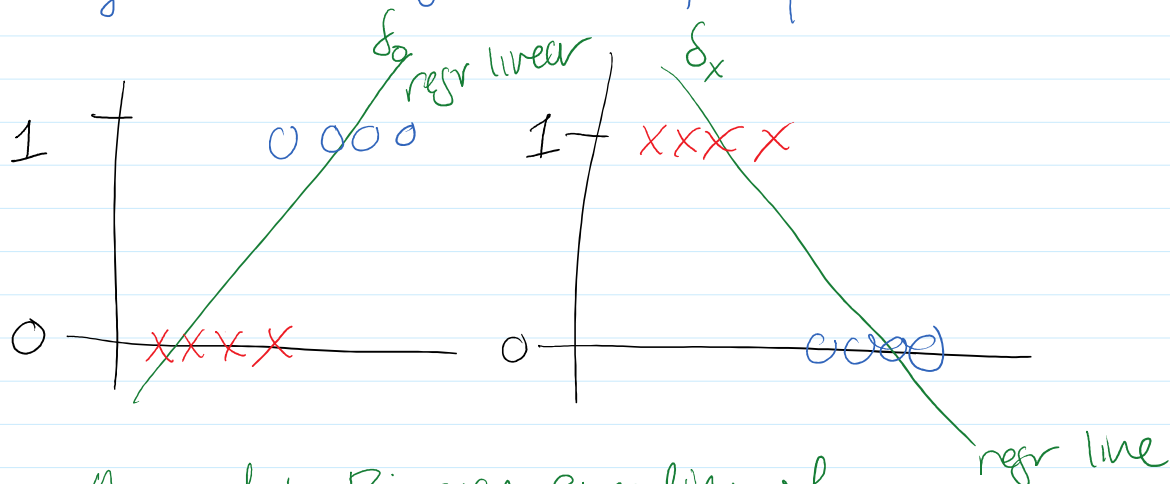| # parameters | $(K-1)(p+1)$ | $(K-1)\left(\frac{p(p+3)}{2}+1\right)$ |
|---|---|---|
| | $\underset{\approx}{Kp}$ | $\underset{\approx}{Kp^2}$ |

much more flex.

Linear classifier: very simply is a linear classifier

$$\delta_c(x) = \beta_{0c} + \beta_c' x$$

LDA estimates $\hat{\beta}_{0c}, \hat{\beta}_c$ using traing data
and an assumption of Normality.

Why not use regression to find $\hat{\beta}s$?



Approach: Binary encoding of
$Y_s$ and fit regr line to this
and use as descri. fns.

$$\delta_c(x) = \hat{\beta}_{0c} + \hat{\beta}_c^T x$$

obtained from LS regress.

Punchline: → this isn't horrible if # classes
(K) is small

→ If $K = 2$ (w/ some caveats)
this will give me LDA.

→ When K is large we have
masking problems

→ when k is large we have masking problems



$\delta_X$  $\delta_0$  $\delta_\square$

Never predict X

predict 0

predict $\square$