

## Lecture 9: Classification

$$f^*(x) = \operatorname{argmax}_c P(Y=c | X=x)$$

One way to build classifier is approx these probs.  
and then

$$\hat{f}(x) = \operatorname{argmax}_c \hat{P}(Y=c | X=x)$$

↑ estimate

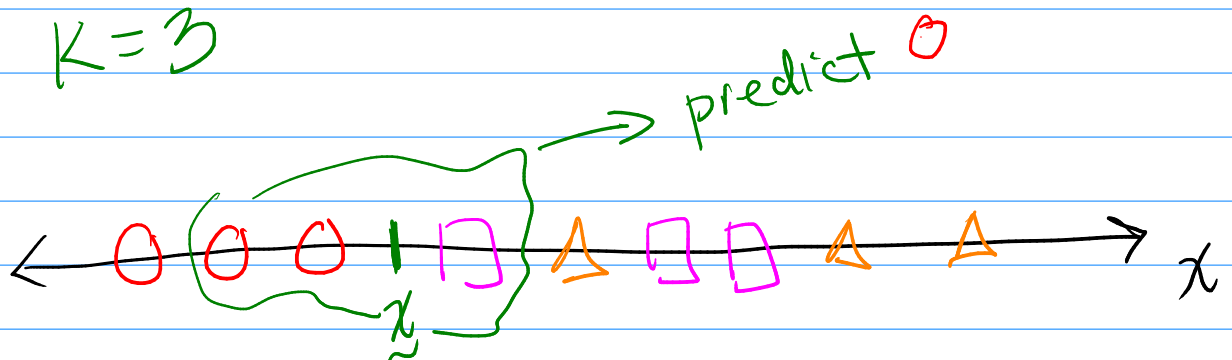
Simple way: KNN classification

$$\hat{P}(Y=c | X=x) \approx \begin{array}{l} \% \text{ of training } y_n\text{'s} \\ \text{in class } c \\ \text{w/ } x_n\text{'s near } x \end{array}$$

$$= \frac{1}{K} \sum_{x_n \in N_K(x)} \mathbb{1}(y_n=c)$$

# pts in class c

$K=3$



$$\underline{x}_1 = (x_{11}, x_{12})$$

$$\underline{x}_2 = (x_{21}, x_{22})$$

$$\|\underline{x}_1 - \underline{x}_2\| = \sqrt{(x_{11} - x_{21})^2 + (x_{12} - x_{22})^2}$$

$$\approx \sqrt{(x_{12} - x_{22})^2}$$

Correct Re-scaling

$$\hat{\mu} = \text{mean}(\text{train})$$

$$\hat{\sigma} = \text{sd}(\text{train})$$

$$z_{\text{train}} = \frac{x_{\text{train}} - \hat{\mu}}{\hat{\sigma}}$$

$$z_{\text{test}} = \frac{x_{\text{test}} - \hat{\mu}}{\hat{\sigma}}$$

Bayes' classifier says look at  $P(Y=c | \underline{X}=\underline{x})$

① discriminative models

→ model  $P(Y=c | \underline{X}=\underline{x})$  directly  
i.e. model  $Y | \underline{X}=\underline{x}$

↳ KNN, logistic regression, classification trees

## ② generative models

$$\text{Bayes' rule: } P(Y=c | \underline{X}=\underline{x}) = \frac{P(\underline{X}=\underline{x} | Y=c)P(Y=c)}{P(\underline{X}=\underline{x})}$$

$$\propto P(\underline{X}=\underline{x} | Y=c)P(Y=c)$$

model  $\underline{X} | Y=c$  and  $Y$

→ LDA/QDA, naive Bayes'

---

### Linear Classifiers

Bayes' says  $\hat{f}(\underline{x}) = \underset{c}{\operatorname{argmax}} P(Y=c | \underline{X}=\underline{x})$

more generally

$$\hat{f}(\underline{x}) = \underset{c}{\operatorname{argmax}} \delta_c(\underline{x})$$

↑ discriminant functions

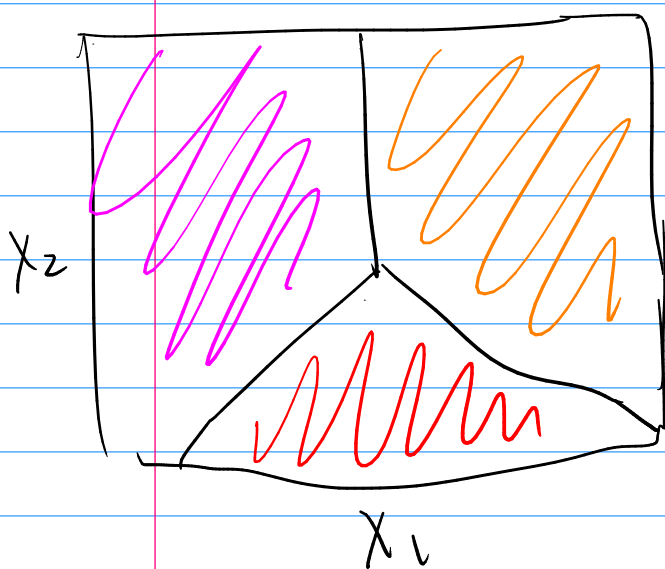
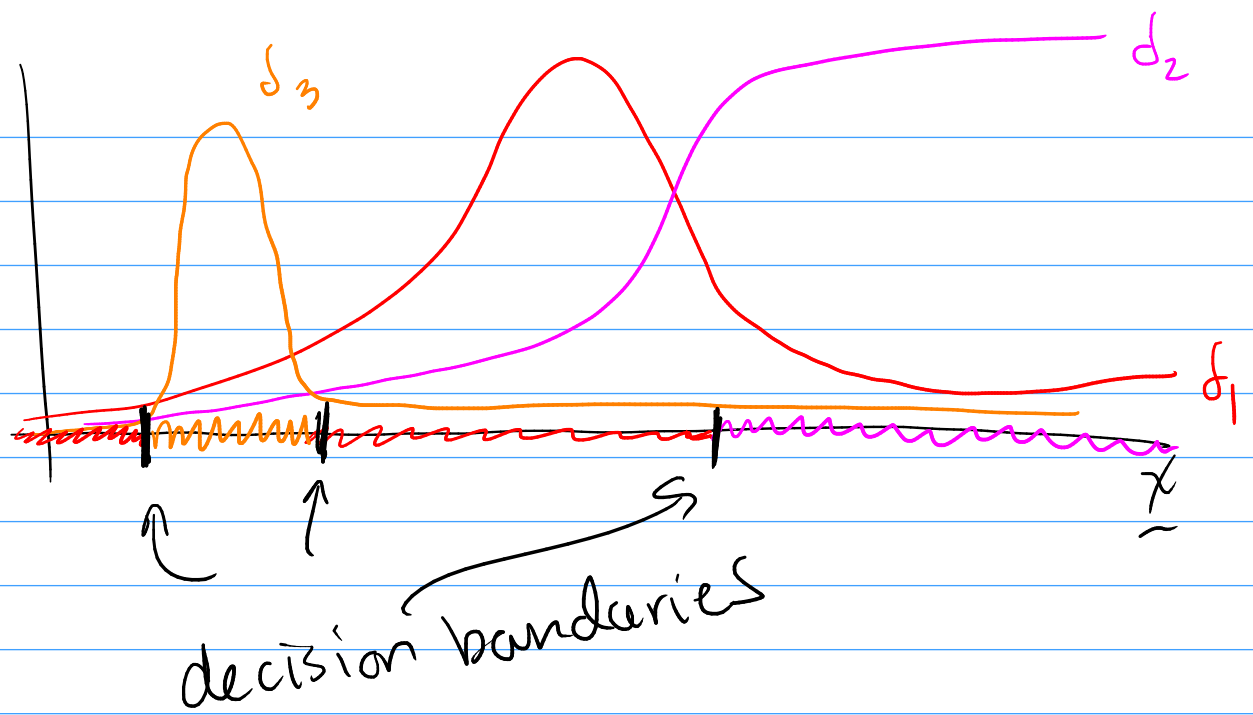
$$\underline{\text{Ex.}} \quad \delta_c(\underline{x}) = \hat{P}(Y=c | \underline{X}=\underline{x})$$

$$\underline{\text{Ex.}} \quad \delta_c(\underline{x}) = \log \hat{P}(Y=c | \underline{X}=\underline{x})$$

↑ any increasing fn

$\delta_c(\underline{x})$  = large if  $\hat{f}(\underline{x})$  likely in class  $c$ , small otherwise

$$\underline{\text{Ex.}} \quad \delta_c(\underline{x}) = \alpha_c + \beta_c^T \underline{x}$$



linear decision boundaries



non-linear

Defn: Linear classifiers have linear decision boundaries

Defn: The discriminant functions  $d_c$  of a linear classifier can be transformed to linear functions via increasing transformations (vice-versa)

Ex.  $\delta_c(\underline{x}) = \alpha_c + \beta_c^T \underline{x}$

gives a linear classifier

Ex.  $\delta_c(\underline{x}) = e^{\alpha_c + \beta_c^T \underline{x}}$   
is also linear

Reason:

$$\hat{f}(\underline{x}) = \underset{c}{\operatorname{argmax}} \delta_c(\underline{x})$$

can depend on  $\underline{x}$  but not  $c$

$$= \underset{c}{\operatorname{argmax}} T(\delta_c(\underline{x}))$$

$T$  increasing

Why do linear classifiers have linear decision boundaries?

For two class problem  $Y = 1$  or  $2$

decision boundary:  $\{\underline{x} : \delta_1(\underline{x}) = \delta_2(\underline{x})\}$

$$\delta_c(\underline{x}) = \alpha_c + \beta_c^T \underline{x}$$

then  $\delta_1(\underline{x}) = \alpha_1 + \beta_1^T \underline{x} = \alpha_2 + \beta_2^T \underline{x} = \delta_2(\underline{x})$

$$\Rightarrow (\alpha_1 - \alpha_2) + (\beta_1 - \beta_2)^T \underline{x} = 0$$

↑ linear system

Generally,

linearizable  
by  $T$

$$\delta_1(\underline{x}) = \delta_2(\underline{x})$$

$$T(\delta_1(\underline{x})) = T(\delta_2(\underline{x}))$$

linear system ...

## LDA: Linear Discriminant Analysis

Model  $\underline{X} | Y = c$  and  $Y$

$$\delta_c(\underline{x}) \equiv P(Y=c | \underline{X}=\underline{x}) \equiv P(\underline{X}=\underline{x} | Y=c) \cdot P(Y=c)$$

$P=1$

up to some  
inc. transf

deps on  $c$

LDA:

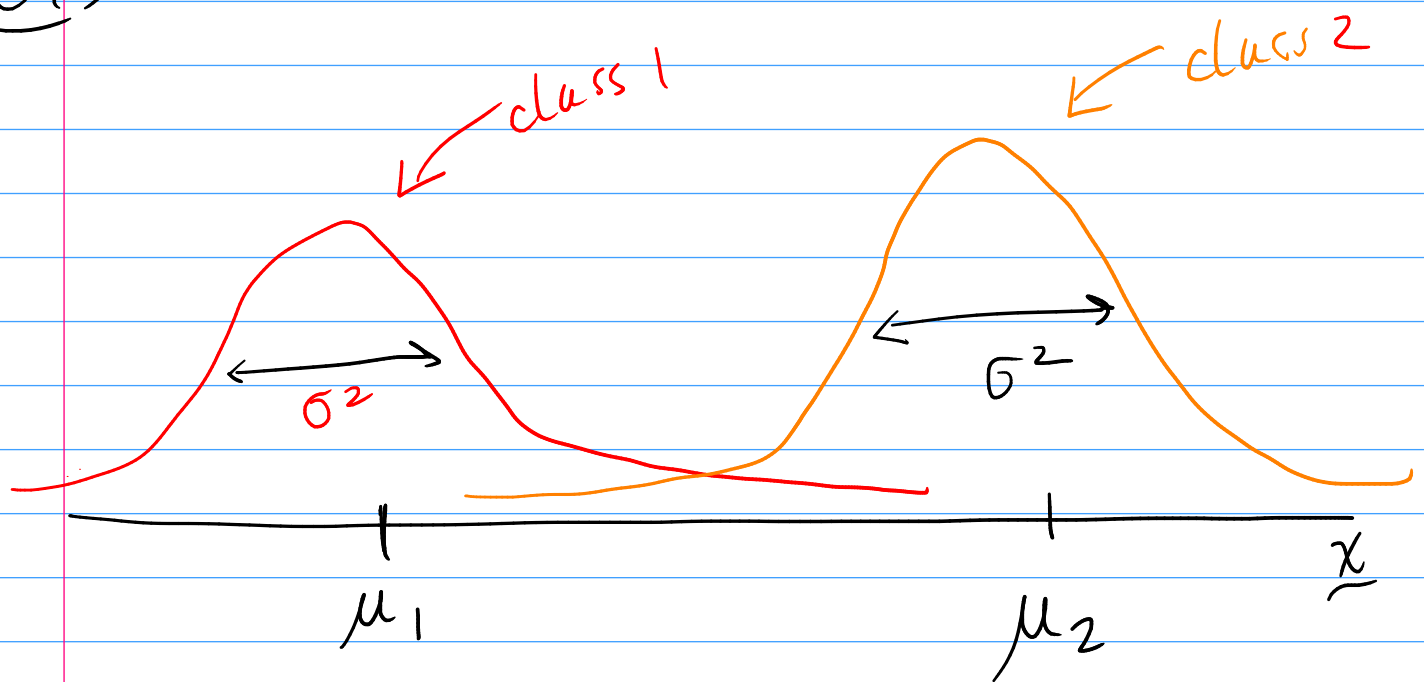
$$\left\{ \begin{array}{l} \underline{X} | Y = c \sim N(\mu_c, \sigma^2) \\ Y \text{ is discrete} \end{array} \right.$$

shared across  $c$

$Y$	1	2	...	$K$
$P(Y=c)$	$\pi_1$	$\pi_2$	...	$\pi_K$

$$\pi_c \geq 0 \text{ and } \sum_c \pi_c = 1$$

$\mathcal{X}$ ,  $K=2$  and  $\underline{x}$  is 1-dim'l



Learning LDA: reduces to learning  
 $\pi_1, \dots, \pi_K, \mu_1, \dots, \mu_K, \sigma^2$

---