

Recall the Bayes' classifier (optimal 0-1 loss classifier)

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

two approaches:

① discriminative: model $Y|X$

ex. KNN

$$P(Y=c|X=x) \approx \begin{array}{l} \% \text{ of nearby} \\ \text{training } y\text{'s of} \\ \text{class } c \end{array}$$

② generative: model $X|Y$ and Y

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

ex. linear discriminant analysis
(today)

↑
LDA

Defn: Linear Classifier

Bayes' says $\hat{y} = \operatorname{argmax}_c P(Y=c|X=x)$

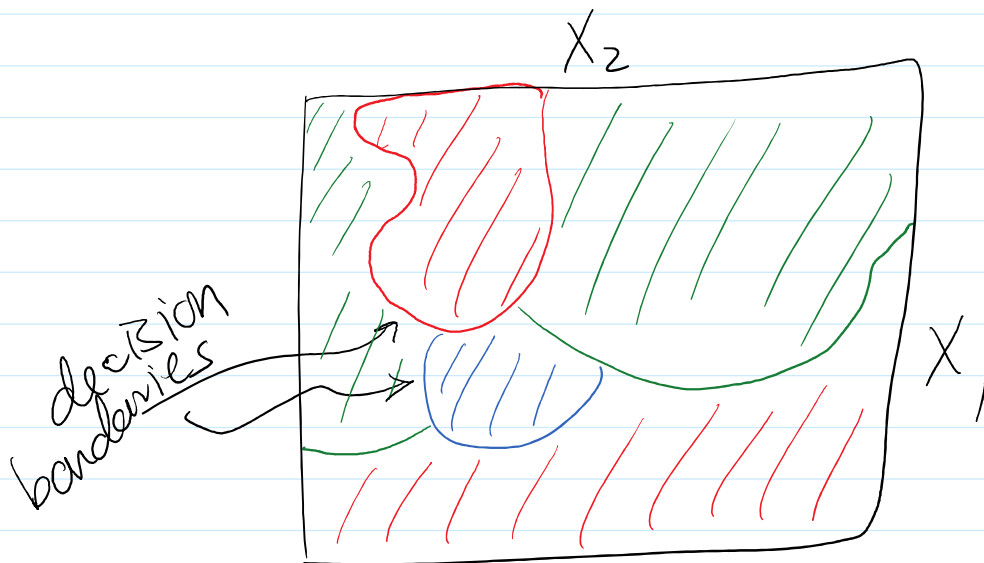
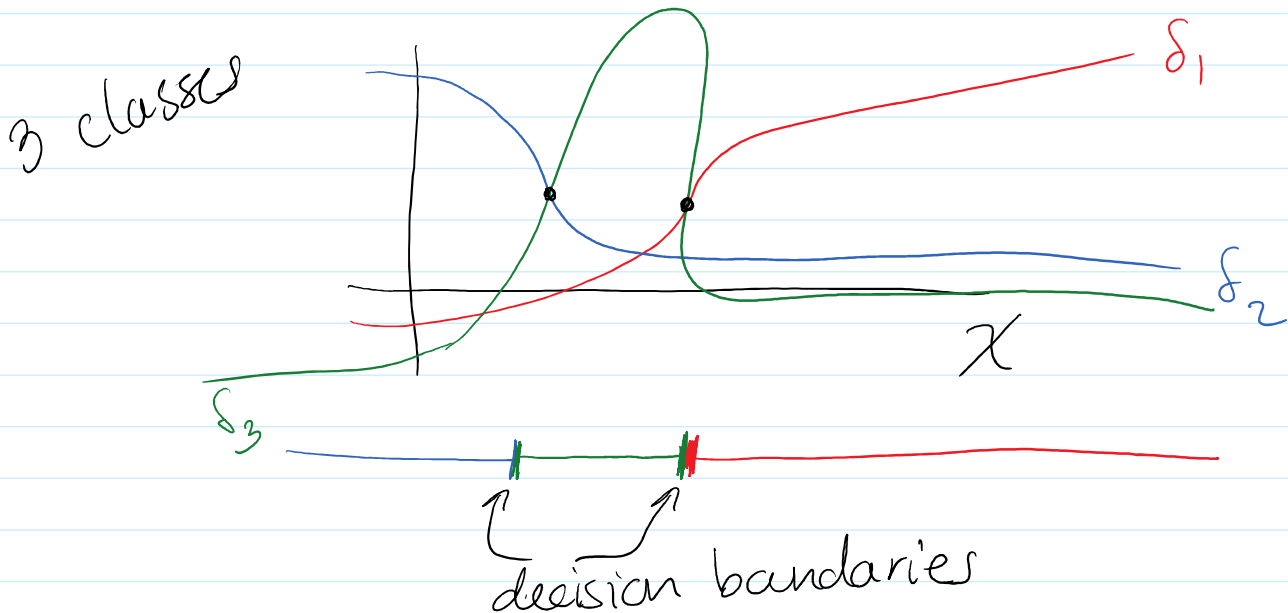
More generally $\hat{y} = \operatorname{argmax}_c \delta_c(x)$ ↑ discriminant function

more generally $y = \dots \delta_c$ ^{discriminative function}
 intuition: higher = better

Ex. Bayes classifier
 is a particular example
 $\delta_c(x) = P(Y=c/X=x)$

Ex. $\delta_c(x) = \beta_{0c} + \beta_c^T x$
 intercept β_{0c} slope β_c^T

Ex. $\delta_c(x) = \exp(\tanh(-\log(\alpha_c x^2)))$



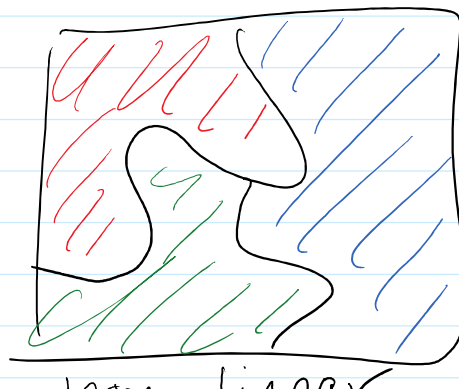
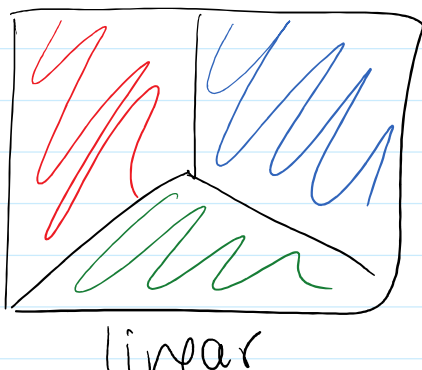
A linear classifier is a classifier whose discriminant functions can be transformed to linear functions (of x) by a increasing (monotone) transformation.

Ex. If the $\delta_c(x)$ is linear in x
 i.e. $\delta_c(x) = \beta_{0c} + \beta_c^T x$
 then the classifier is linear.

Ex. there is a monotone non-decreasing transformation T so that
 $T(\delta_c(x)) = \beta_{0c} + \beta_c^T x$

Ex. $T(x) = \log(x)$ or $T(x) = \exp(x)$

(*) The decision boundaries of linear classifiers are linear. (*)



linear

non-linear

Two classes: $\{x | \delta_1(x) = \delta_2(x)\}$

decision boundary

Condition:

$$\delta_1(x) = \delta_2(x)$$

if δ_1 and δ_2 are linear in x then so is the

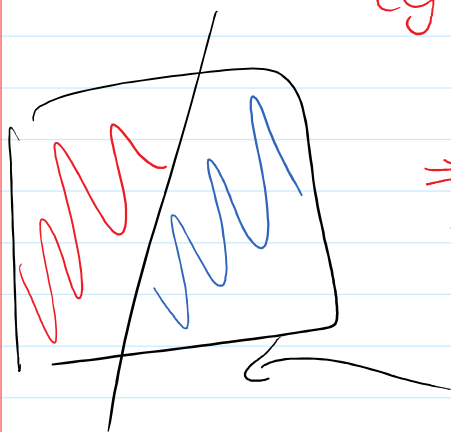
eg. $\delta_1(x)$ solution set $\delta_2(x)$

$$\beta_{01} + \beta_1^T x = \beta_{02} + \beta_2^T x$$

$$\Rightarrow (\beta_{01} - \beta_{02}) = (\beta_2 - \beta_1)^T x$$

linear algebra

says solution set is a subspace



Notice if T is increasing then

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

adding T gives

EXACT same classifier

→ decision boundaries are the same

x	$T(P(Y=1 x))$	$T(P(Y=2 x))$	$T(P(Y=3 x))$
x^2	5^2	3^2	2^2

$$T(x) = x^2 \quad \begin{array}{c} 5 \\ \uparrow \\ y=1 \end{array} \quad | \quad \begin{array}{c} 3 \\ \uparrow \\ y=1 \end{array} \quad | \quad \begin{array}{c} 2 \\ \uparrow \\ y=1 \end{array}$$

→ decision boundaries are the same

So if $\exists T$ that makes δ_c linear then

$$\hat{y} = \operatorname{argmax}_c \delta_c(x) = \operatorname{argmax}_c \underbrace{T(\delta_c(x))}_{\text{linear}}$$

another way ex

decision boundary defined $T(\delta_1(x)) = T(\delta_2(x))$

Linear Discriminant Analysis (LDA)

$$\delta_c(x) = P(Y=c|X=x) = \frac{P(X=x|Y=c)P(Y=c)}{P(X=x)}$$

Need: model $X|Y$ and Y

$P(X=x)$

denom doesn't depend on c

LDA: $\begin{cases} X|Y & \text{as Gaussian (Normal)} \\ Y & \text{as discrete} \end{cases}$

$$\operatorname{argmax}_y f(y)$$

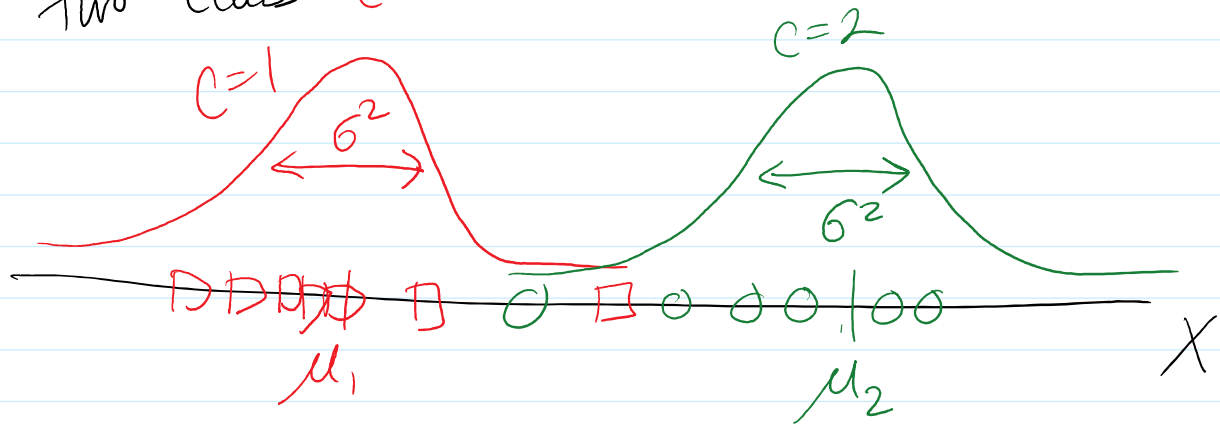
$$\operatorname{argmax}_y K \cdot f(y)$$

Ex let $x \in \mathbb{R}$ i.e. $P=1$

$$\text{LDA: } X|Y=c \sim N(\mu_c, \sigma^2)$$

$$P(Y=c) = \pi_c \text{ where } \sum_c \pi_c = 1, \pi_c \geq 0$$

Ex. two-class $c=1$ or $c=2$



$$\boxed{P=1}$$

$$X|Y=c \sim N(\mu_c, \sigma^2)$$

$$P(Y=c) = \pi_c$$

learn them
(estimate)
from training data

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

really
use

↓

$$\delta_c(x) \leftarrow P(X=x | Y=c) P(Y=c)$$

$$= \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu_c)^2\right) \pi_c \leftarrow$$

$$\delta_c(x) \leftarrow \log(\text{wavy line})$$

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

$$= \underbrace{-\frac{1}{2} \log(2\pi\sigma^2)}_{\substack{\text{ignore} \\ \text{b/c doesn't depend} \\ \text{on } c}} - \frac{1}{2\sigma^2}(x - \mu_c)^2 + \log \pi_c$$

$$\delta_c(x) \leftarrow -\frac{1}{2\sigma^2}(\underbrace{x^2 - 2x\mu_c + \mu_c^2}_{\text{ignore}}) + \log \pi_c$$

$$\boxed{\delta_c(x) = -\frac{\hat{\mu}_c^2}{2\hat{\sigma}^2} + \frac{x\hat{\mu}_c}{\hat{\sigma}^2} + \log \hat{\pi}_c}$$

reality: need to estimate μ_c 's, σ^2 , π_c 's

How: $\hat{\mu}_c$ = mean of training x_n 's in class c

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

$\hat{\pi}_c$ = % of training data in class c

Note: For LDA:

$$\delta_c(x) = \underbrace{-\frac{\hat{\mu}_c^2}{2\hat{\sigma}^2}}_{\hat{\beta}_{0c}} + \underbrace{\frac{\hat{\mu}_c}{\hat{\sigma}^2}}_{\hat{\beta}_c} x + \underbrace{\log \hat{\pi}_c}_{\hat{\beta}_{0c}}$$

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

So LDA is linear.