

Lecture 6 Bayesian Decision Theory

BDT

- BDT is a framework for making optimal decisions on problems involving uncertainty.

Framework

1) world has states / classes drawn from r.v. \bar{Y} .

e.g. $\bar{Y} \in \{H, T\}$, $\bar{Y} \in \{ok, flu, cold\}$

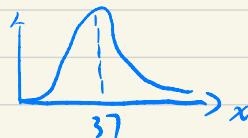
Prior: $P(\cdot | \bar{Y})$ - prior probability of state occurring.

2) Observer measures features / observations from r.v. X
class-conditional density (CCD)

$P(X | \bar{Y})$ - observations conditioned on the class/state.

e.g. $X = \text{temperature}$.

$$P(X | \text{ok}) =$$



$$P(X | \text{flu}) =$$



$$P(X | \text{cold}) =$$



3) Decision Function - uses observation to make a decision about the state $g(x): X \rightarrow \bar{Y}$

4) loss function - penalty for deciding the wrong \hat{Y} (wrong decision)

$$L(g(x), y) \geq 0$$

e.g. 0-1 loss function $L(g(x), y) = \begin{cases} 0, & g(x) = y \\ 1, & \text{otherwise} \end{cases}$

Goal: Find the optimal decision function $g^*(x)$ for the given assumptions (loss, prior, (c) ...)

Bayes Decision Rule (BDR)

Risk - expected value of the loss function

$$\begin{aligned} \text{Risk} &= E_{x,y}[L(g(x), Y)] = \sum_y \int_x \frac{p(x|y)}{p(y|x)p(x)} L(g(x), Y) dx \\ &= \int_x p(x) \left[\sum_y p(y|x) L(g(x), Y) \right] dx \\ &\quad \text{conditional risk } R(x) \\ &\quad \text{function of } x \\ &= E_x[R(x)] \\ &\quad \text{expectation of conditional Risk} \end{aligned}$$

Since $L(g(x), y) \geq 0$, then minimizing the risk is equivalent to minimizing the conditional risk $R(x)$ for each x

0-1 loss function and classification.

$$Y \subseteq \{1, 2, \dots, C\}$$

$$g(x) \in \{1, 2, \dots, C\}$$

$$L(g(x), y) = \begin{cases} 1 & g(x) \neq y \text{ (misclassification)} \\ 0 & \text{otherwise} \end{cases}$$

Conditional Risk : $R(x) = \mathbb{E}_{y|x} [L(g(x), y)]$

indicator variable
 $= \Pr(g(x) \neq y | x)$
 probability of error given the x .

BDR.

$$\begin{aligned} y^* &= \operatorname{argmin}_{j \in Y} R(x) = \operatorname{argmin}_{j \in Y} \Pr_{y=j} L(y|x) \\ &= \operatorname{argmin}_j 1 - \Pr_{y=j}(y|x) \\ y^* &= \operatorname{argmax}_j \Pr_{y=j}(y|x) \end{aligned}$$

MAP rule - choose the class with largest posterior.

$$\begin{aligned} \text{Equivalent} \\ y^* &= \operatorname{argmax}_j \frac{\Pr_{y=j} \Pr_{x|y=j}}{\Pr_x} = \operatorname{argmax}_j \Pr_{x|y=j} \Pr_{y=j} \end{aligned}$$

$$y^* = \operatorname{argmax}_j \underbrace{\log \Pr_{x|y=j}}_{\text{CD}} + \underbrace{\log \Pr_{y=j}}_{\text{Prior}}$$

Example : 2-class problem (0, 1)

$$\text{pick } 0 \text{ if } \underbrace{p(x|0)}_{\text{CD}} \underbrace{p(0)}_{\text{prior}} > p(x|1) p(1) \Rightarrow \underbrace{\frac{p(x|0)}{p(x|1)}}_{\text{likelihood ratio test}} > \frac{p(1)}{p(0)} = T$$

Summary

for 0-1 loss function.

- BDR is MAP rule (tell us the threshold)
- Risk = probability of error
- BDR minimize the risk . i.e. the Probability of an error.
(nothing is better)
- Caveat : assuming the model (densities are correct)
 CD Prior

This is called a generative classification model.

model how data is generated in the world.

- CD and prior learned from data

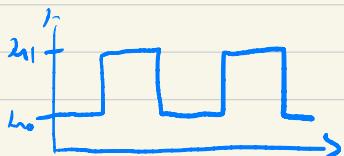
Example: Noisy channel.

$Y = \text{Bit } \{0, 1\} \rightarrow \boxed{\text{trans}} \rightarrow \boxed{\text{channels}} \rightarrow \boxed{\text{receive}} \rightarrow$

$$\text{decoder } y = \begin{cases} 0, & x < T \\ 1, & x \geq T \end{cases}$$

How to choose T ?

assume $l_{in} < l_{in}'$



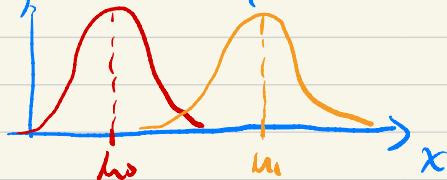
Given measurement X , recover bit \hat{Y} .

$$\text{Class Probability: } P(\hat{Y}=0) = P(Y=1) = \frac{1}{2}$$

CCD: Gaussian additive noise: $X = \mu_0 + \epsilon$, $\epsilon \sim N(0, \sigma^2)$

$$P(X|Y=0) = N(X|\mu_0, \sigma^2)$$

$$P(X|Y=1) = N(X|\mu_1, \sigma^2)$$



Assume 0-1 loss, the BDR is:

$$y^* = \arg \max_j (\log p(X|j)) + \log \frac{1}{2}$$

$$= \arg \max_j (\log N(X|\mu_j, \sigma^2)) + \log \frac{1}{2}$$

$$= \arg \max_j -\frac{1}{2\sigma^2} (X - \mu_j)^2 - \frac{1}{2} \log 2\pi - \frac{1}{2} \log \sigma^2 + \log \frac{1}{2}$$

$$= \arg \max_j -\frac{1}{2\sigma^2} (X^2 - 2X\mu_j + \mu_j^2)$$

$$= \arg \min_j -2X\mu_j + \mu_j^2$$

Hence, pick 0 when $-2X\mu_0 + \mu_0^2 < -2X\mu_1 + \mu_1^2$

$$\Rightarrow 2X(\mu_1 - \mu_0) < (\mu_1 - \mu_0)(\mu_1 + \mu_0)$$

$$\Rightarrow X < \frac{1}{2}(\mu_1 + \mu_0)$$

intuitive threshold \rightarrow halfway between μ_0 and μ_1

Assumption are explicit:

1) 0-1 loss. BDR.

2) uniform class prior ($P(Y=0) = \frac{1}{2}$)

3) Gaussian noise (i.i.d), additive, same for each bit.

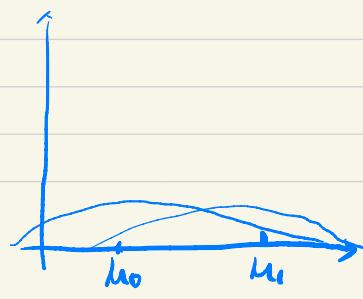
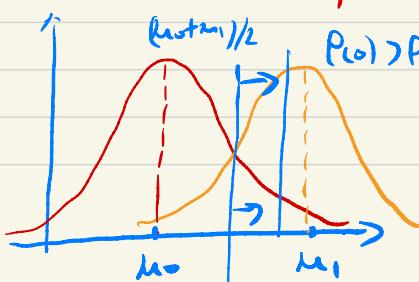
What if $p(Y)$ is not uniform?

e.g. coding: $7 \Rightarrow 1111110$

Pick 0 if

$$\log p(x|\mu_0, \sigma^2) + \log p(0) > \log p(x|\mu_1, \sigma^2) + \log p(1)$$
$$-\frac{1}{2} \log 2\pi - \frac{1}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (x - \mu_0)^2 + \log p(0) >$$
$$-\frac{1}{2} \log 2\pi - \frac{1}{2} \log \sigma^2 - \frac{1}{2\sigma^2} (x - \mu_1)^2 + \log p(1)$$
$$\Rightarrow \frac{1}{2\sigma^2} [x^2 - 2x\mu_0 + \mu_0^2 - x^2 + 2x\mu_1 - \mu_1^2] + \log \frac{p_0}{p_{(1)}} > 0$$
$$\Rightarrow \frac{1}{2\sigma^2} [(\mu_1 - \mu_0)(\mu_1 + \mu_0) - 2x(\mu_1 - \mu_0)] + \log \frac{p_0}{p_{(1)}} > 0$$
$$\Rightarrow \frac{2\sigma^2}{\mu_1 - \mu_0} \log \frac{p_0}{p_{(1)}} + (\mu_1 + \mu_0) - 2x =$$
$$\Rightarrow x < \underbrace{\frac{1}{2}(\mu_1 + \mu_0)}_{\text{Same as before}} + \underbrace{\frac{\sigma^2}{(\mu_1 - \mu_0)} \log \frac{p_0}{p_{(1)}}}_{\text{Same as before}}$$

1 ↓
"normalized distance between means"
 $\left| \begin{array}{l} p(y=0) > p(y=1) \Rightarrow \log \frac{p_0}{p_{(1)}} > 1 \Rightarrow > 0 \\ \text{increase the threshold if 0 is more frequent (predict more 0) and vice versa} \end{array} \right.$



↓
normalized distance is
(large \Rightarrow ignore the priors.) | normalized distance
is small \Rightarrow use priors.

Gaussian classifier

$$Y \in \{1, \dots, C\} \quad C \text{ classes} \quad P(y=j) = \pi_j$$

$x \in \mathbb{R}^d$, C CDs are m.u. Gaussians.

$$P(x|y=j) = N(x|\mu_j, \Sigma_j)$$

BDR:

$$g(x) = \arg \max y \log P(x|y=j) + (1-y) \log P(x)$$

$$= \arg \max_j -\frac{1}{2} \|x - \mu_j\|_{\Sigma_j}^2 - \frac{1}{2} [\log |\Sigma_j| + \log \pi_j]$$

$g_j(x)$ = discriminant function for class j

special case: $\Sigma_j = \sigma^2 I$ (shared isotropic covariances,

($\nabla g_j(\bar{x}) = 0$)

$$g_j(x) = w_j^T x + b_j$$

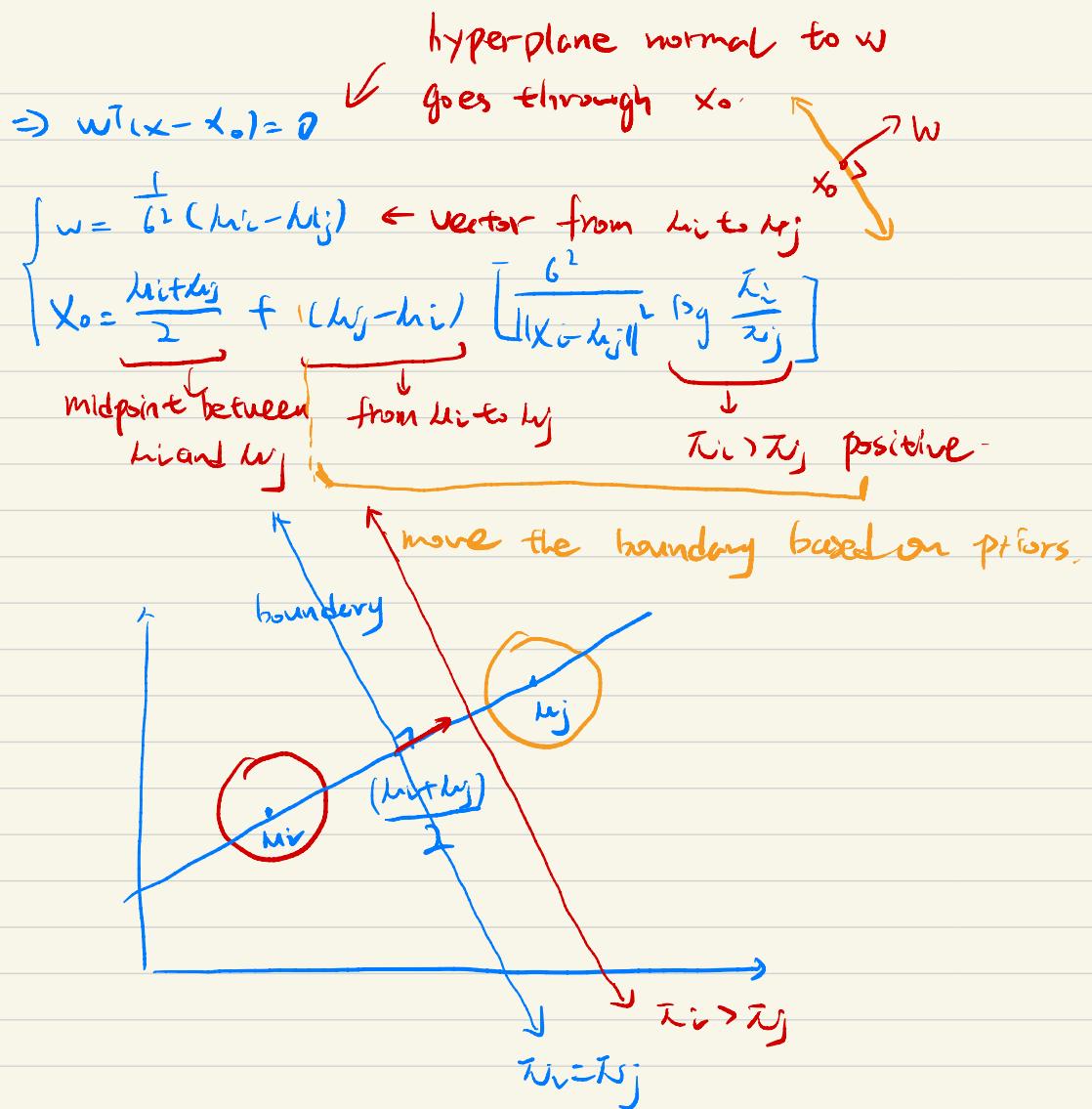
$$\text{where } \begin{cases} w_j = \frac{1}{\sigma^2} \mu_j \\ b_j = \frac{1}{2\sigma^2} \mu_j^T \mu_j + \log \pi_j \end{cases}$$

Geometric meaning

classes i and j share a boundary if $g_i(x) = g_j(x)$

$$w_i^T x + b_i = w_j^T x + b_j$$

: tutorial
↓



analogous to the 1-D version (noisy channel)
 \Rightarrow hyperplane \sim high-dim threshold.