

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

Framework

- 1) world has states/classes drawn from r.v. Y

eg. $Y \in \{H, L\}$, $Y \in \{ok, flu, cold\}$

prior: $p(Y)$ - prior prob. of state occurring.

- 2) observer measures features/observations from r.v. X

class - conditional density (ccd)

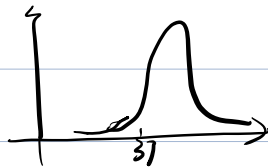
$P(X|Y)$ - observations conditioned on class/state

eg. X - temperature

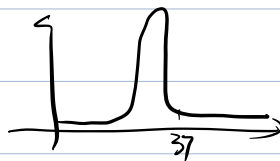
$p(X|ok)$



$p(X|flu)$



$p(X|cold)$



- 3) Decision Function - use observation to make a decision about the state.

$g(x): X \rightarrow Y$

(4) loss function - penalty for deciding the wrong Y (wrong decision)

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

eg. 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 given assumptions (loss, prior, cdd)

Bayes Decision Rule (BDR)

Risk - expected value of the loss function

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

$$= E_x [R(x)] \leftarrow \text{expectation of Conditional Risk}$$

Since $L(g(x), y) \geq 0$, then minimizing risk \Leftrightarrow
minimizing conditional risk $R(x)$, $\forall x$.

$$\begin{aligned} \forall x, \quad g^*(x) = y^* &= \underset{j \in Y}{\operatorname{argmin}} R(x) = \underset{j \in Y}{\operatorname{argmin}} \sum_y p(y|x) L(j, Y) \\ &= \underset{j}{\operatorname{argmin}} E_{Y|x} [L(j, Y)] \\ &\quad \downarrow \\ &\quad \text{function of } x \end{aligned}$$

0-1 loss function & classification

$$Y \in \{1, 2, \dots, c\}$$

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

In this case $L(g(x), y) \Rightarrow$ indicator variable

$$\text{Conditional Risk } R(x) = E_{y|x} [L(g(x), y)]$$

$$= \Pr(g(x) \neq y | x) \leftarrow \text{prob. of error, given } x$$

BDR

$$y^* = \operatorname{argmin}_j R(x) = \operatorname{argmin}_j \Pr(j \neq y | x)$$

$$= \operatorname{argmin}_j 1 - \Pr(y = j | x)$$

$$y^* = \operatorname{argmax}_j \Pr(y = j | x) \quad \text{MAP rule}$$

choose the class w/ largest posterior.

Equivalently,

$$y^* = \operatorname{argmax}_j \frac{P(x|y=j) \cdot P(y=j)}{P(x)} = \operatorname{argmax}_j P(x|y=j) P(y=j)$$

$$= \operatorname{argmax}_j \underbrace{\log P(x|y=j)}_{\text{CED}} + \underbrace{\log P(y=j)}_{\text{prior (typically estimated from data)}}$$

Example

2-class problem (0, 1)

$$p_{i2} < 0 \quad \text{if} \quad P(x|0) \cdot p_{i0} > P(x|1) \cdot p_{i1}$$

$$\Rightarrow \underbrace{\frac{P(x|0)}{P(x|1)}}_{\text{likelihood ratio}} > \frac{p_{i1}}{p_{i0}} = T \quad \leftarrow \text{threshold.}$$

Summary

for 0-1 loss function:

- BDR is MAP rule (tells the threshold)

- Risk = prob. of error.

-BDR minimizes the risk, i.e. the prob. of error
(nothing is better)

Caveat: assume the model (density is correct)

CCD & prior best you can do given them

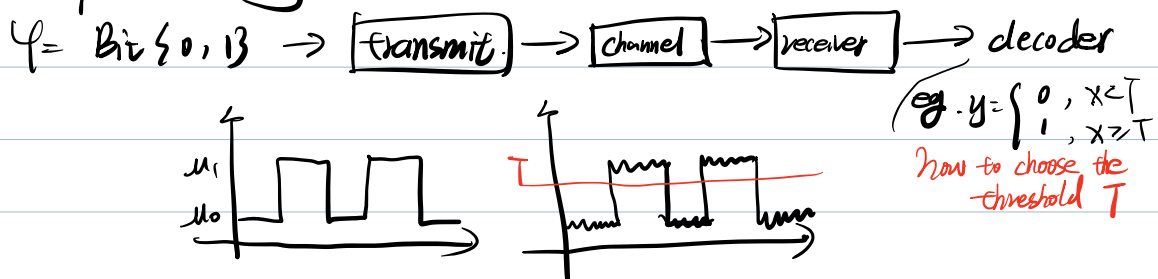
This is called a generative classification model

• model how data is generated in the world

— CCD & prior learned from data

assume they are from some existing density, i.e.,

Example: noisy channel.



Given measurement X , recover bit Y

class prob. $P(Y=0) = P(Y=1) = \frac{1}{2}$

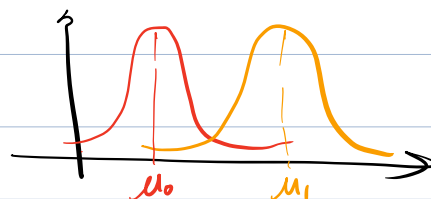
CCD: Gaussian additive noise

(assume $\mu_1 > \mu_0$)

$$X = \mu_Y + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2)$$

$$p(X|Y=0) = \mathcal{N}(X|\mu_0, \sigma^2)$$

$$p(X|Y=1) = \mathcal{N}(X|\mu_1, \sigma^2)$$



Assume 0-1 loss, the BDR is

$$y^* = \arg \max_j \log p(X|j) + \log p(j)$$

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

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

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

pick 0 when $-2X\mu_0 + \mu_0^2 \leq -2X\mu_1 + \mu_1^2$

$$X < \frac{\mu_1^2 - \mu_0^2}{2(\mu_1 - \mu_0)} = \frac{\mu_0 + \mu_1}{2}$$

intuitively halfway between the two

assumptions are explicit:

1) 0-1 loss, BDR

2) uniform class prior

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

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

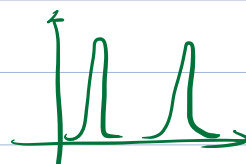
eg. channel coding: 7 \Rightarrow 1111110

BDR: pick 0 if $X < \underbrace{\frac{\mu_1 + \mu_0}{2}}_{\text{midpoint}} + \underbrace{\frac{\sigma^2}{\mu_1 - \mu_0} \log \frac{p(Y=0)}{p(Y=1)}}_{\text{eg. } p(Y=0) > p(Y=1) \Rightarrow \log(1.71) \Rightarrow > 0}$

increase \uparrow if 0 is more frequent.
make sense since expecting to see more 0

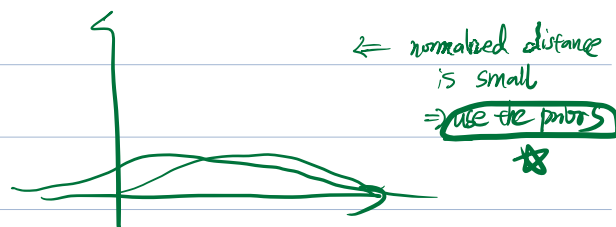
normalized
distance
between means

eg:



\Leftarrow normalized distance
is large
 \Rightarrow ignore priors

Unlikely we have
confusion since
they center around
the mean



Gaussian Classifier

$Y \in \{1, \dots, C\}$, C classes

$$p(Y=j) = \pi_j$$

$X \in \mathbb{R}^d$, c.c.d.s are m.v. Gaussian.

$$p(X|Y=j) = \mathcal{N}(X|\mu_j, \Sigma_j)$$

BDR
$$g(x) = \arg \max_j \log p(X|j) + \log p(j)$$

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

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

Special case: $\Sigma_j = \sigma^2 I$ (shared isotropic covariances)

(tutorial 1)

$$g_j(x) = w_j^T x + b_j \quad \text{discriminant is linear.}$$

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 & j share a boundary if $g_i(x) = g_j(x)$

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

tutorial

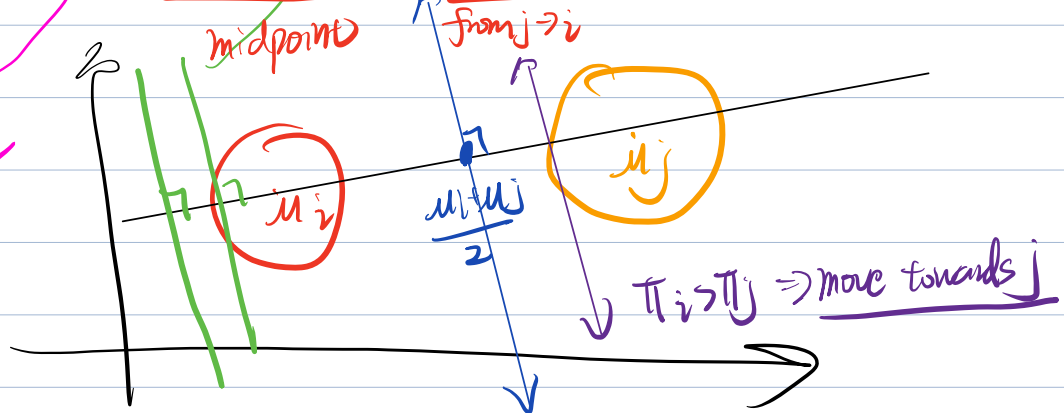
$$\Rightarrow w^T(x - x_0) = 0$$

hyperplane normal to w , goes through x_0

$$w = \frac{1}{\sigma^2} (\mu_i - \mu_j)$$

$$x_0 = \left(\frac{\mu_i + \mu_j}{2} \right) + (\mu_j - \mu_i) \left[\frac{\sigma^2}{\|\mu_i - \mu_j\|^2} \log \frac{\pi_i}{\pi_j} \right]$$

parameters
determine
boundary



analogous to 1-D version (noisy channel)
 \Rightarrow hyper plane \sim high-dim.