

Lecture 5 - 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. Y

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

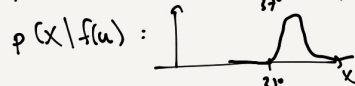
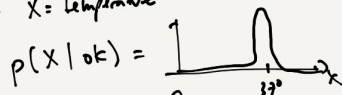
prior: $p(Y)$ - prior prob. of state occurring.

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

class-conditional density (CCD):

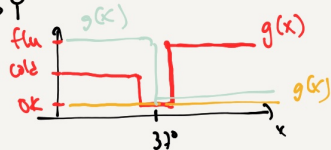
$p(X|Y)$ - observations conditioned on the class/state.

e.g. X = temperature



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

$$g(x): X \rightarrow Y$$



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

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

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, CCD, ...)

Bayes Decision Rule (BDR)

Risk - expected value of the loss function

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

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

For an x ,

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

"Bayes Decision Rule" - will give the minimum risk!

0-1 loss function & classification

$$y \in \{1, 2, \dots, C\}$$

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

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

$$\text{Conditional Risk: } R(x) = E_{\text{r.v.}} [L(g(x), y)]$$

indicator variable

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

BDA

$$\begin{aligned}y^* &= \underset{j \in Y}{\operatorname{argmin}} R(x) = \underset{j \in Y}{\operatorname{argmin}} \Pr(j \neq y | x) \\ &= \underset{j}{\operatorname{argmin}} 1 - \Pr(y = j | x)\end{aligned}$$

$$y^* = \underset{j}{\operatorname{argmax}} \Pr(y = j | x)$$

MAP rule - choose the class w/ largest posterior.

Equivalently

$$y^* = \underset{j}{\operatorname{argmax}} \frac{p(x|y=j)p(y=j)}{p(x)} = \underset{j}{\operatorname{argmax}} p(x|y=j)p(y=j)$$

$$y^* = \underset{j}{\operatorname{argmax}} \underbrace{\log p(x|y=j)}_{\text{CCD}} + \underbrace{\log p(y=j)}_{\text{prior}}$$

Example: 2-class problem (0,1)

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

likelihood ratio test threshold

Summary

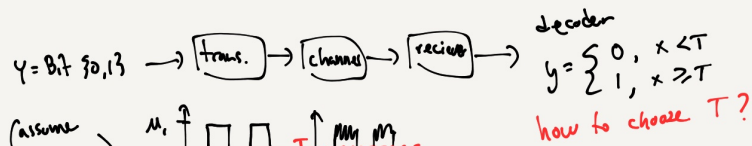
for 0-1 loss function:

- BDR is MAP rule (tells us the threshold)
- Risk = probability of error
- BDR minimizes the risk, i.e. the prob. of error (nothing is better)
- caveat: assuming the model (densities are correct) (CCD, prior)

This is called a generative classification model.

- model how data is generated in the world.
- CCD & prior learned from data.

Example: Noisy channel



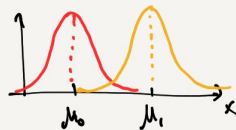
Given measurement x , recover bit y

class prob. $p(y=0) = p(y=1) = \frac{1}{2}$

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

$$p(x|y=0) = N(x|\mu_0, \sigma^2)$$

$$p(x|y=1) = N(x|\mu_1, \sigma^2)$$



Assume 0-1 loss, the BDR is:

$$\begin{aligned}
 y^* &= \underset{j}{\operatorname{argmax}} \log p(x|j) + \log p(j) \\
 &= \underset{j}{\operatorname{argmax}} \log N(x|\mu_j, \sigma^2) + \log \frac{1}{2} \\
 &= \underset{j}{\operatorname{argmax}} \underbrace{\left[-\frac{1}{2\sigma^2}(x-\mu_j)^2 - \frac{1}{2}\log \sigma^2 - \frac{1}{2}\log 2\pi + \log \frac{1}{2} \right]}_{\text{constant w.r.t } j} \\
 &= \underset{j}{\operatorname{argmax}} \underbrace{-\frac{1}{2\sigma^2}}_{\text{constant}} \underbrace{(x^2 - 2x\mu_j + \mu_j^2)}_{\text{w.r.t } j} \\
 &= \underset{j}{\operatorname{argmin}} -2x\mu_j + \mu_j^2
 \end{aligned}$$

Here, pick 0 when $-2x\mu_0 + \mu_0^2 < -2x\mu_1 + \mu_1^2$

$$-2x\mu_0 + \mu_0^2 < \mu_1^2 - \mu_0^2$$

$$x \cdot 2 \cdot (\mu_1 - \mu_0) < \mu_1^2 - \mu_0^2$$

$$x < \frac{\mu_1^2 - \mu_0^2}{2(\mu_1 - \mu_0)} \leftarrow (\mu_1 - \mu_0)(\mu_1 + \mu_0)$$

assume $\mu_0 < \mu_1$

$$\Rightarrow \boxed{\text{pick 0 when } x < \frac{\mu_1 + \mu_0}{2}}$$

intuitive threshold \rightarrow halfway between μ_0 & μ_1

Assumptions are explicit:

1) 0-1 loss, BDR

2) uniform class prior ($p(y=0) = \frac{1}{2}$)

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

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

PS 6-3

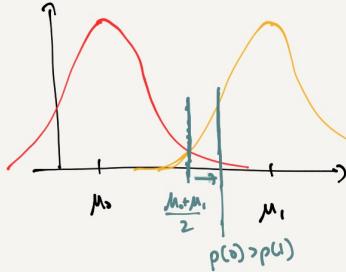
e.g. channel coding: $7 \Rightarrow 1111110$

BDR: pick 0 if $x < \frac{\mu_0 + \mu_1}{2} + \frac{\sigma^2}{\mu_1 - \mu_0} \log \frac{p(y=0)}{p(y=1)}$

Same as before

$p(y=0) > p(y=1) \Rightarrow \log(\frac{0}{1}) \Rightarrow > 0$

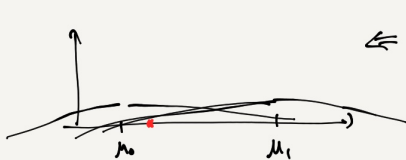
• increase the threshold if 0 is more frequent (predict more 0), vice versa



"normalized distance from means"



⇐ normalized distance is large
⇒ ignore the priors



⇐ normalized distance is small
⇒ use the priors.

Gaussian classifier

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

$p(y=j) = \pi_j$

$x \in \mathbb{R}^d$, CDS are m.v. Gaussians

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

BDR: $g(x) = \arg \max_j \log p(x|j) + (\log p(j))$
 $= \arg \max_j \underbrace{-\frac{1}{2} \|x - \mu_j\|_{\Sigma_j}^2 - \frac{1}{2} \log |\Sigma_j| + \log \pi_j}_{g_j^*(x) = \text{discriminant function for class } j.}$

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

∴ (linear)

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

← discriminant is linear function.

where $w_j = \frac{1}{\sigma^2} \mu_j$

$b_j = -\frac{1}{2\sigma^2} \mu_j^T \mu_j + \log \pi_j$

Geometric meaning

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

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

hyperplane

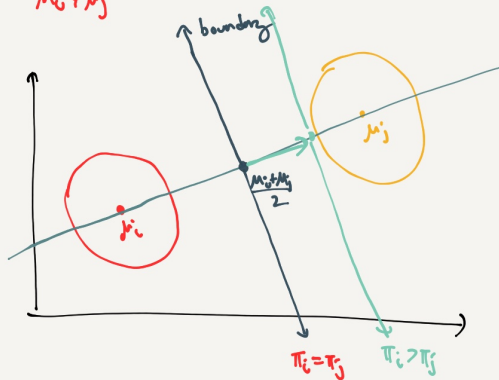


$$\Rightarrow w^T(x - x_0) = 0 \quad \leftarrow \text{hyperplane normal to } w, \text{ goes through } x_0$$

$$w = \frac{1}{\delta^2} (\mu_i - \mu_j) \quad \leftarrow \text{vector from } \mu_j \text{ to } \mu_i$$

$$x_0 = \underbrace{\left(\frac{\mu_i + \mu_j}{2} \right)}_{\text{midpoint b/w } \mu_i \text{ \& } \mu_j} + \underbrace{(\mu_j - \mu_i)}_{\text{from } \mu_i \text{ to } \mu_j} \underbrace{\left[\frac{\delta^2}{\|\mu_i - \mu_j\|^2} \log \frac{\pi_i}{\pi_j} \right]}_{\pi_i > \pi_j \Rightarrow \text{positive}}$$

moves the boundary based on priors.



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