

Bayesian Decision Theory (BDT)

- BDT is a framework for making optimal decisions on problems involving uncertainty (probabilities)
- Statistical approach to pattern classification.

Framework

1) World has states/classes, drawn from r.v. Y .

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

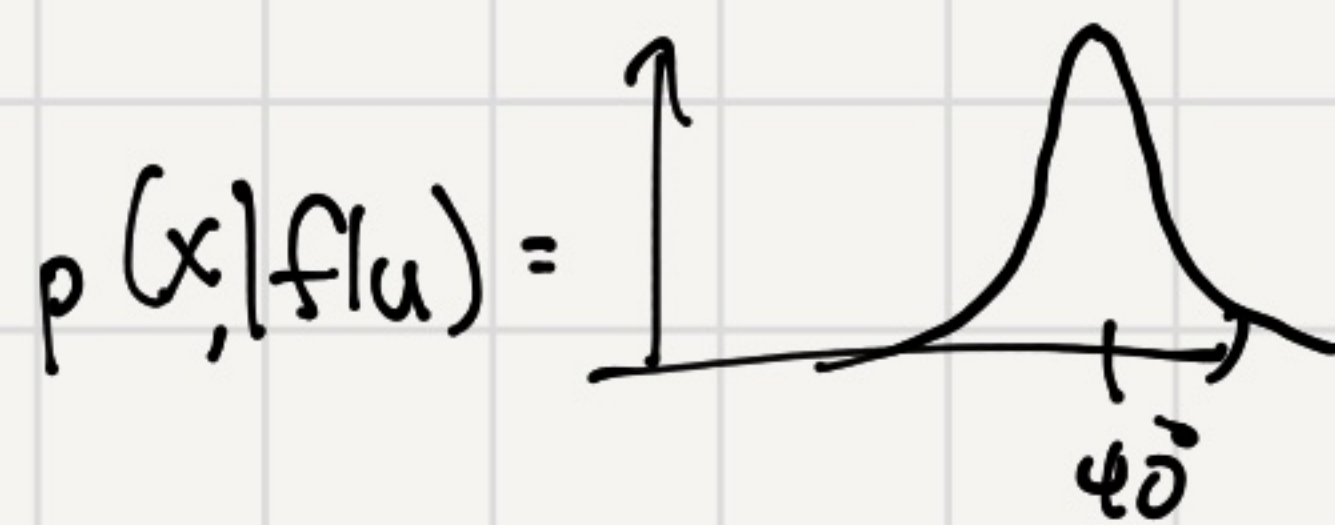
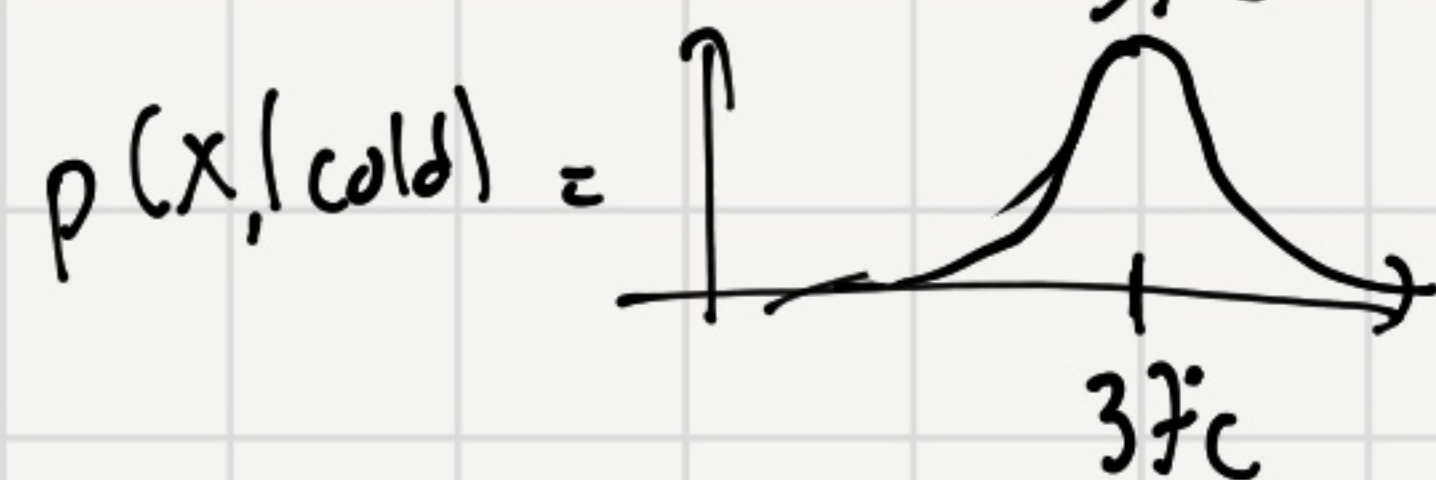
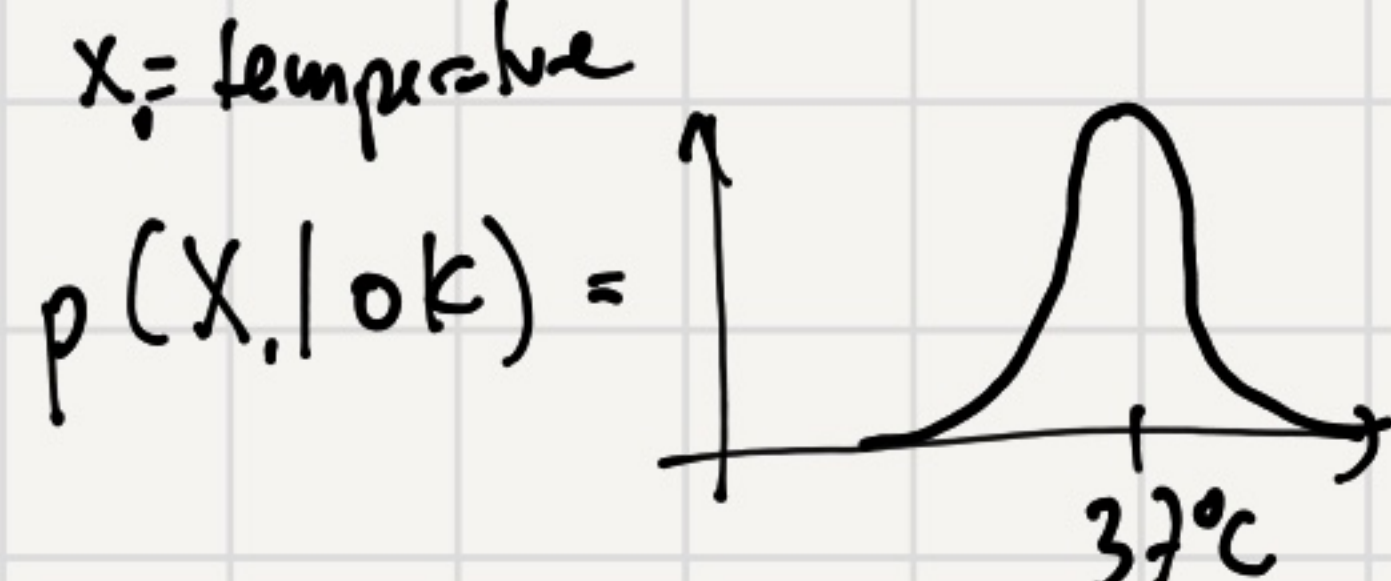
prior: $p(Y)$ = prior prob. of state occurring.

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

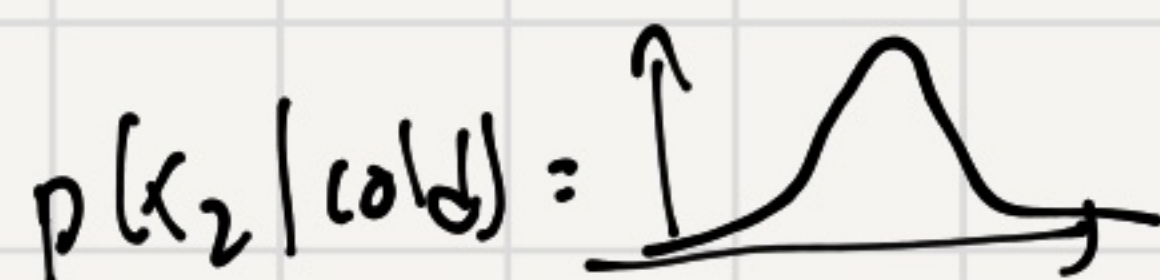
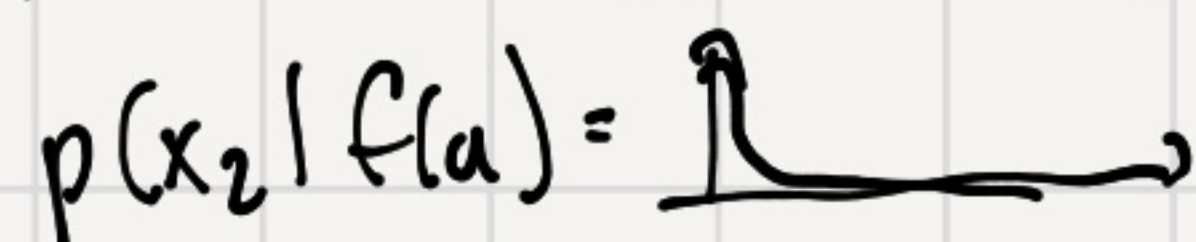
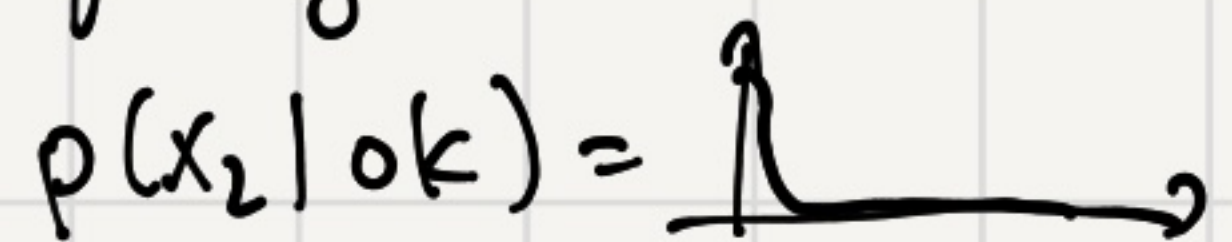
class conditional density (CCD)

$p(X|Y)$ = ^{distribution of} features for a particular class/state.

e.g. x_1 = temperature

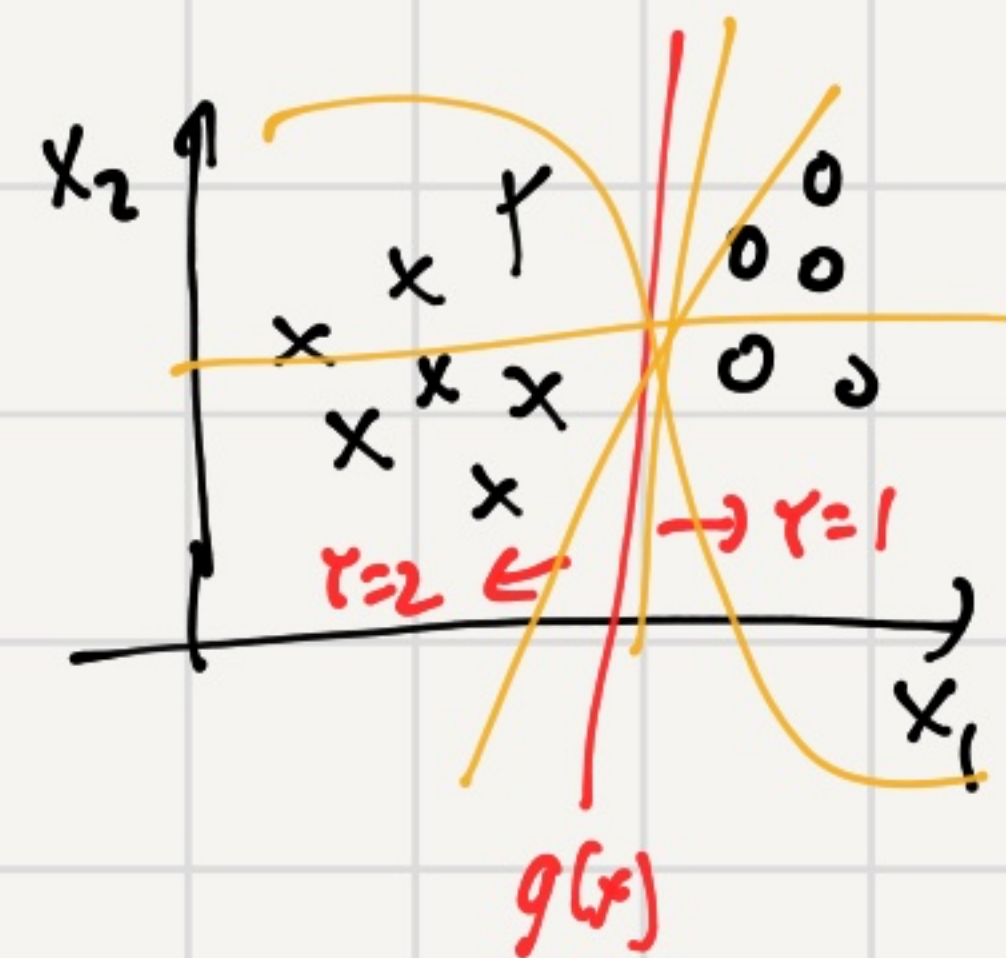


x_2 = quantity of snot



3) Decision Function - use features to infer the state.

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



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

$$L(g(x), y) = \begin{cases} 0, & g(x) = y \\ 1, & g(x) \neq y \end{cases} \quad \left. \begin{array}{l} \uparrow \text{prediction} \\ \uparrow \text{true} \end{array} \right\} \text{0-1 loss functions.}$$

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

Goal: Find an optimal decision function $g^*(x)$ for the above assumptions (loss, CCD, prior)

Bayes Decision Rule (BDR)

Risk - expected value of the loss

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

$$= \sum_{y \in Y} \int_X \underbrace{p(x, y)}_{p(y|x)p(x)} L(g(x), y) dx$$

$$= \int_X \sum_y \underbrace{p(y|x)p(x)}_{\uparrow} L(g(x), y) dx$$

$$= \int_X p(x) \left[\sum_y p(y|x) L(g(x), y) \right] dx = E_x [R(x)]$$

conditional risk: $R(x)$ (given observation x) expectation of cond. risk.

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

Given an x ,

$$g^*(x) = g^* = \argmin_{j \in Y} R(x) = \argmin_{j \in Y} \sum_y p(y|x) L(j, y)$$

value of $g(x)$ ↓

$$= \argmin_{j \in Y} E_{y|x} [L(j, y)]$$

conditional exp. of loss.

Bayes' Decision Rule

0-1 loss function & classification

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

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

Conditional Risk:

$$R(x) = E_{y|x} [L(g(x), y)] = \underbrace{P_r(g(x) \neq y | x)}_{\substack{\text{indicator variable} \\ \text{prob. of misclassifying} \\ x \cdot (\text{prob. of error})}} = 1 - \underbrace{P(g(x) = y | x)}_{\substack{\text{prob. correct} \\ \text{classification} \\ \text{of } x. \\ p(y = g(x) | x)}}$$

$$\text{BDR: } g^*(x) = g^* = \argmin_{j \in Y} \underbrace{1 - p(y=j | x)}_{R(x)}$$

$$\boxed{g^*(x) = \argmax_{j \in Y} p(y=j | x)}$$

MAP rule: choose j w/ largest posterior probability.

Equivalently, $g^*(x) = \argmax_j \frac{p(x|y=j)p(y=j)}{p(x)}$

$$\boxed{g^*(x) = \argmax_j p(x|y=j)p(y=j) = \argmax_j \log p(x|y=j) + \log p(y=j)}$$

Example: 2-class problem $y \in \{0, 1\}$

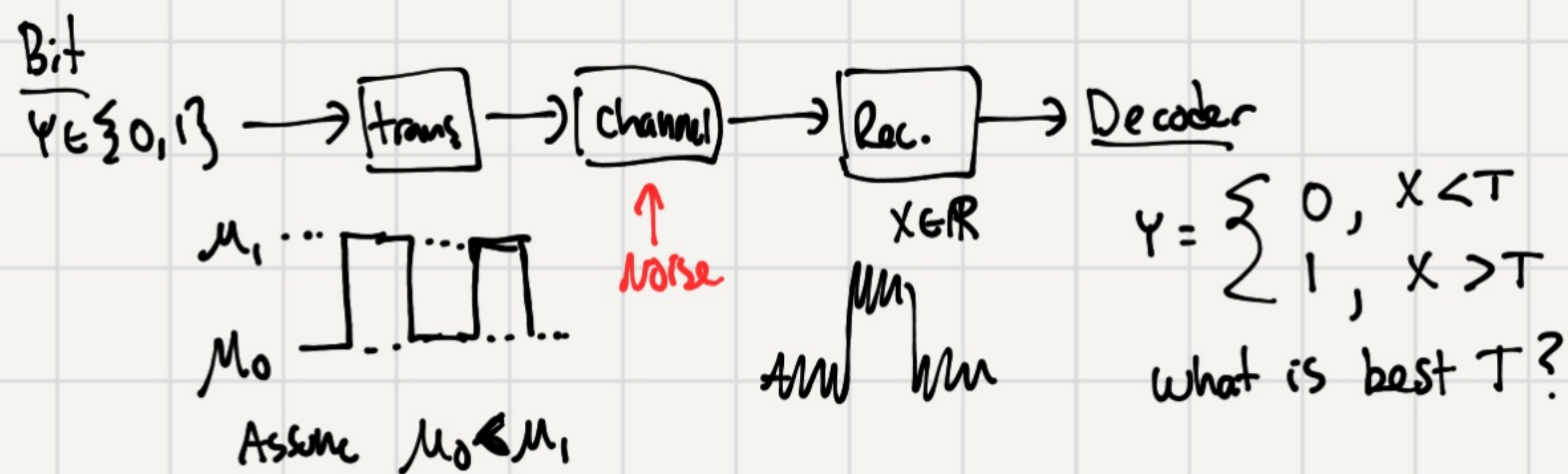
$$\text{pick 0 if } p(x|0)p(0) > p(x|1)p(1) \Rightarrow \underbrace{\frac{p(x|0)}{p(x|1)}}_{\text{likelihood ratio test}} > \underbrace{\frac{p(1)}{p(0)}}_{\text{threshold}} = T$$

Summary: for 0-1 loss:

- BDR is MAP rule
- Risk = prob. of error
- BDR minimizes risk (best you can do)
- caveat: assume densities are correct.

generative classification model

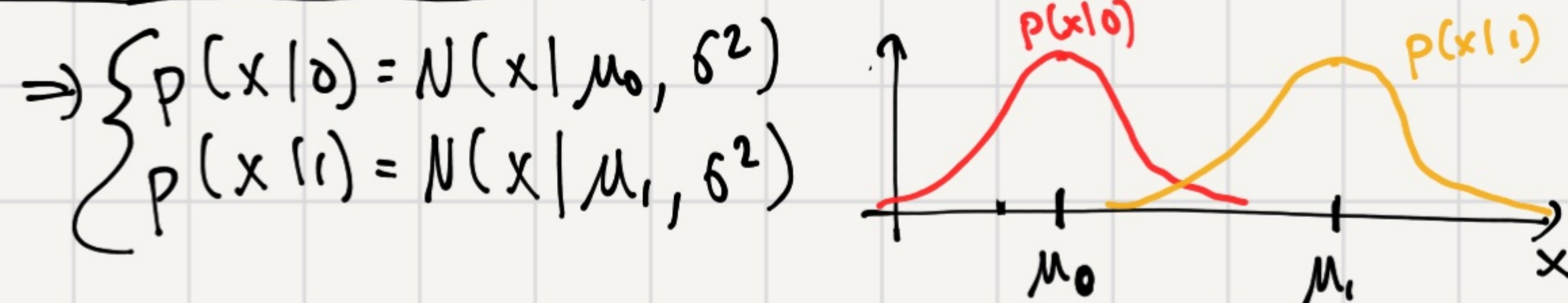
Example: Noisy Channel



Goal: given X , recover bit Y .

Assume: $p(Y=0) = p(Y=1) = \frac{1}{2}$

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



BDR w/o 0-1 loss

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

$$= \arg\max_j \underbrace{-\frac{1}{2\sigma^2} (x - \mu_j)^2}_{\text{scalar}} - \underbrace{\frac{1}{2} \log 2\pi - \frac{1}{2} \log \sigma^2 + \log \frac{1}{2}}_{\text{constant}}$$

$$= \arg\min_j (x - \mu_j)^2 = \arg\min_j \underbrace{x^2 - 2x\mu_j + \mu_j^2}_{\text{constant}}$$

$$y^* = \arg\min_j \mu_j^2 - 2x\mu_j$$

Hence: pick 0 when $\mu_0^2 - 2x\mu_0 < \mu_1^2 - 2x\mu_1$

$$2x(\mu_1 - \mu_0) < \mu_1^2 - \mu_0^2 \iff a^2 - b^2 = (a-b)(a+b)$$

Assumptions

- 1) 0-1 loss (MPE)
- 2) uniform prior
- 3) Gaussian iid noise, additive

$$X < \frac{\mu_0 + \mu_1}{2}$$

intuitive threshold \rightarrow halfway btwn means.

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

e.g. Channel coding: $7 \rightarrow 1111110$

\downarrow PS6-3

BDR: pick 0 if:

$$X < \underbrace{\frac{\mu_1 + \mu_0}{2}}_{\text{same as before}} + \underbrace{\frac{\sigma^2}{\mu_1 - \mu_0} \log \frac{p(Y=0)}{p(Y=1)}}_{\text{same as before}}$$

• if $p(Y=0) = p(Y=1) \Rightarrow$ same as before

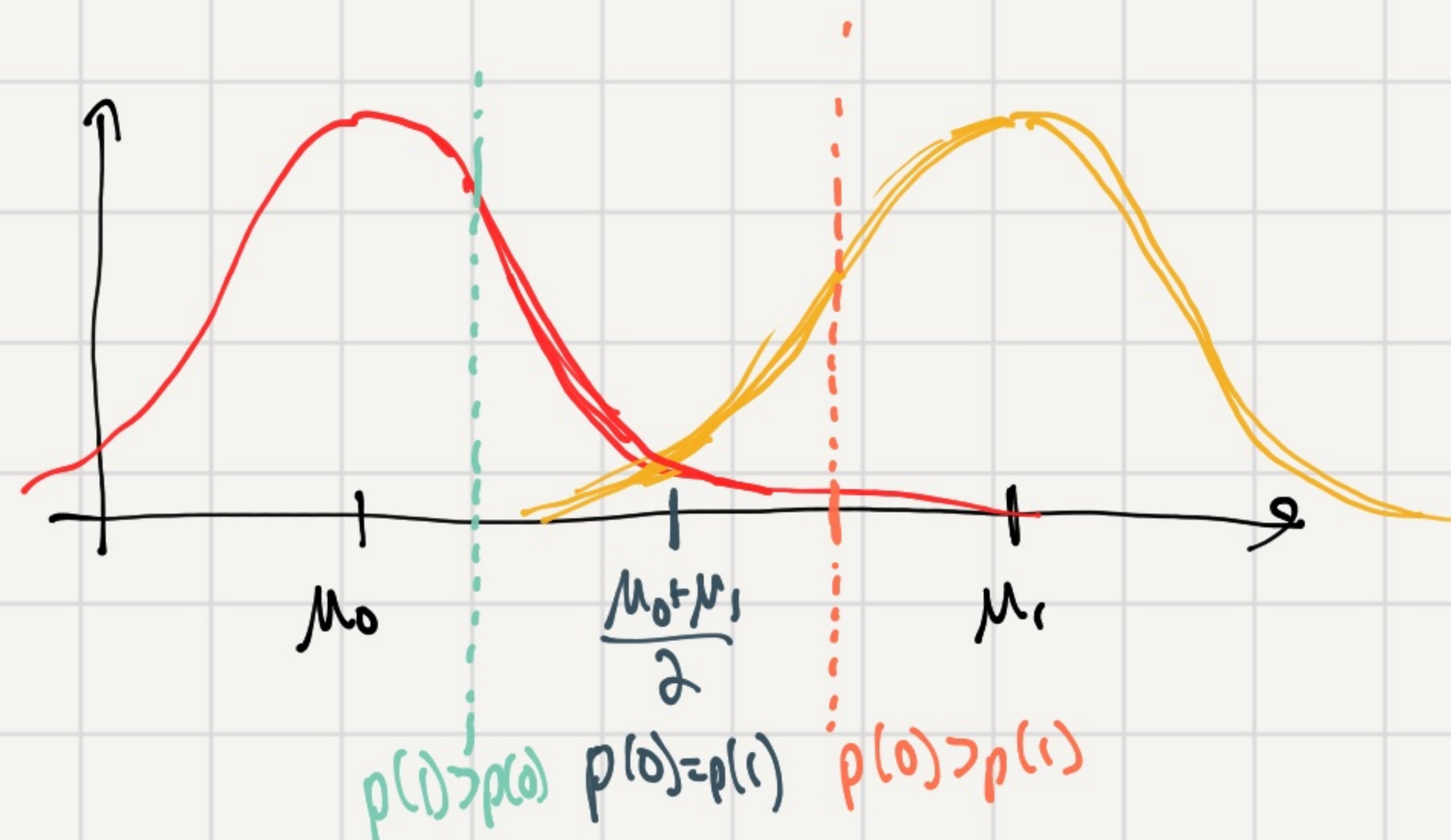
• if $p(Y=0) > p(Y=1) \Rightarrow \log \frac{p(0)}{p(1)} > 0 \Rightarrow$ increase threshold (capture more space for bit 0)

• if $p(1) > p(0) \Rightarrow \log \frac{p(0)}{p(1)} < 0 \Rightarrow$ decrease threshold

"normalized distance btwn means"

• if means are far apart, then ignore priors

• if means are close, use prior.



Gaussian Classifier

$Y \in \{1, \dots, C\}$, C classes, $p(Y=j) = \pi_j$
 $x \in \mathbb{R}^d$

CCD are Gaussian: $p(x|y=j) = N(x|\mu_j, \Sigma_j)$

different mean,
cov, & prior for
each class

BDR & 0-1 loss: $g^*(x) = \arg \max_j \log p(x|y=j) + \log p(j)$

$$g^*(x) = \arg \max_j \underbrace{-\frac{1}{2} \|x - \mu_j\|_{\Sigma_j}^2 - \frac{1}{2} \log |\Sigma_j| + \log \pi_j}_{g_j(x)}$$

Special case: Assume: $\Sigma_j = \sigma^2 I$ (shared isotropic cov.)

$$\begin{aligned} g_j(x) &= -\frac{1}{2\sigma^2} \|x - \mu_j\|^2 - \frac{1}{2} \log |\sigma^2 I| + \log \pi_j \\ &= -\frac{1}{2\sigma^2} (x^T x - 2x^T \mu_j + \mu_j^T \mu_j) + \log \pi_j + \text{const.} \\ &= -\frac{1}{2\sigma^2} (-2x^T \mu_j + \mu_j^T \mu_j) + \log \pi_j + \text{const.} \end{aligned}$$

$$\Rightarrow g_j(x) = \underbrace{w_j^T x + b_j}_{\text{linear discriminant function}} \quad \text{where} \quad \begin{cases} w_j = \frac{1}{\sigma^2} \mu_j \\ b_j = -\frac{1}{2\sigma^2} \mu_j^T \mu_j + \log \pi_j \end{cases}$$

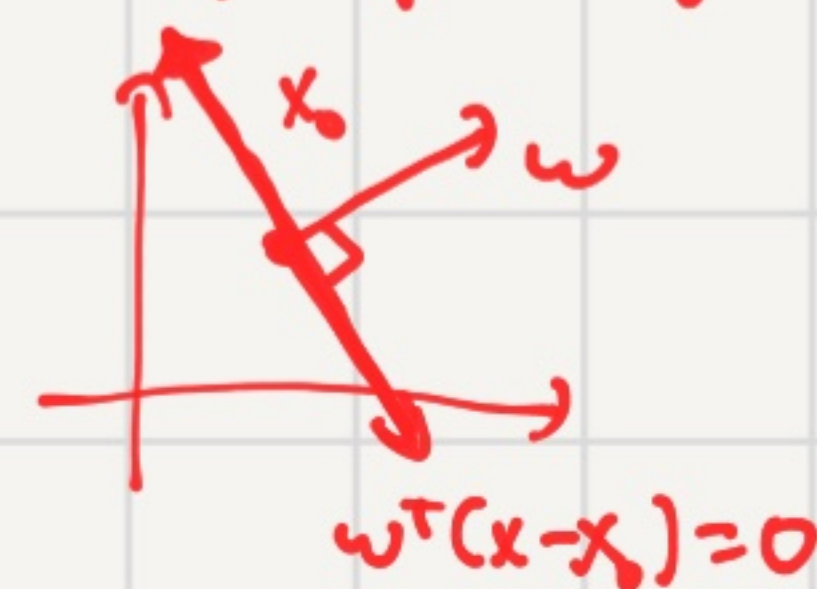
Note: Similar to 1-D
case \Rightarrow hyperplane as
high-dim version of
threshold.

Geometric meaning

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

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

$\Rightarrow w^T(x - x_0) = 0$ \leftarrow hyperplane normal to w , & passing through x_0 .



$$\begin{cases} w = \frac{1}{\sigma^2} (\mu_i - \mu_j) \\ x_0 = \underbrace{\frac{\mu_i + \mu_j}{2}}_{\text{midpoint btwn means}} + \underbrace{(\mu_j - \mu_i)}_{\text{vector from } \mu_i \rightarrow \mu_j} \left[\underbrace{\frac{\sigma^2}{\|\mu_i - \mu_j\|^2}}_{\text{normalized distance}} \log \frac{\pi_i}{\pi_j} \right] \end{cases}$$

$\pi_i > \pi_j \Rightarrow \log(\cdot) > 0$
 $\pi_i < \pi_j \Rightarrow \log(\cdot) < 0$
 move $x_0 \Rightarrow$ boundary

