# Bayes Decision Theory (BDT)

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

## Framework

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

   e.g. $Y \in \{H, T\}$, $Y \in \{A, B, C, D, F\}$, $Y \in \{ok, flu\}$

   prior: $p(Y)$ – prior probability of a state occurring in the world.

2) observer measures features/observations from r.v. $X$.
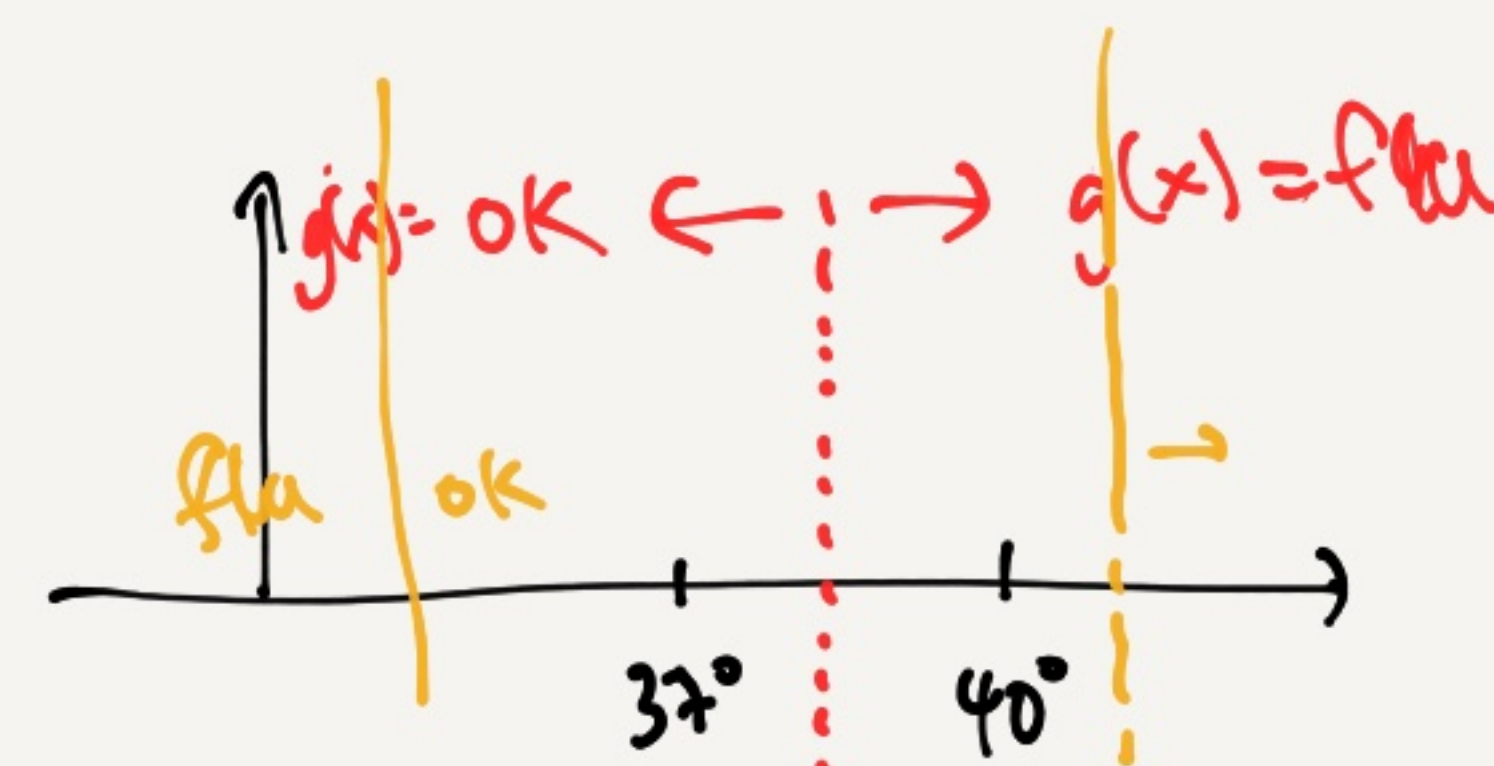
   - class conditional densities (CCDs)

     $p(X|Y)$ – observations conditioned on the state $y$.

     e.g. $Y \in \{ok, flu\}$, $X = temperature$

     $$\begin{cases} p(X|ok) = \\[2em] p(X|flu) = \end{cases}$$

     

3) decision function. – use observation to make a decision about the state of the world.

   $g(x): X \to Y$

   

4) Loss function – penalty for deciding the wrong $y$ or making the wrong decision.

   $L(g(x), y) \geq 0$, e.g. 0-1 loss function

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

**Goal:** Find the optimal $g^*(x)$ for the given assumptions.

(loss function, CCD, prior)

# Bayes Decision Rule (BDR)

Risk – expected value of the loss function

$$\text{Risk} = E_{X,Y}\left[L(g(X), Y)\right]$$

$$= \sum_y \int_x p(x,y) \, L(g(x), y) \, dx$$

$\underbrace{\qquad}_{p(y|x)\,p(x)}$

$$= \int_x \sum_y p(y|x)\, p(x)\, L(g(x), y)\, dx$$

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

$\underbrace{\qquad\qquad}_{\text{conditional Risk } R(x)\text{ of a particular } x.}$

$$= E_X\left[R(x)\right] \quad \leftarrow \text{expectation of conditional risk.}$$

Minimizing the Risk can be achieved by minimizing the <u>conditional</u> risk for each x.

$$R(x) = E_{Y|x}\left[L(g(x), y)\right]$$

For an x

$$g^*(x) = y^* = \underset{j \in Y}{\arg\min}\ R(x)$$

$$= \underset{j \in Y}{\arg\min}\ \sum_y p(y|x) L(j, y)$$

↖ our decision.

$$g^*(x) = \underset{j \in Y}{\arg\min}\ E_{Y|x}\left[L(j, y)\right]$$

↖ conditional expectation of the loss function

⇑

"Bayes Decision Rule"

## 0-1 loss function for Classification

Classes: $Y \in \{1, \ldots, C\}$

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

← misclassification by $g(x)$.

## Conditional Risk:

$$R(x) = E_{Y|x}[\underbrace{L(g(x), y)}_{\text{indicator function}}]$$

⟩ expectation of indicator is the prob. of the thing happening

$$= P_r(\underbrace{g(x) \neq y \mid x}_{\substack{\text{probability of an error} \\ \text{given } x.}})$$

$$= \sum_{y \neq g(x)} p(y \mid x) = 1 - \underbrace{p(y = g(x) \mid x)}_{\text{prob. of correct}}$$

Thus minimizing $R(x)$ is equivalent to minimizing the probability of making an error.

## BDR

$$y^* = \arg\min_j 1 - p(y = j \mid x)$$ ← minimize conditional risk

$$y^* = \arg\max_j p(y = j \mid x)$$ ← choose $j$ that has highest posterior probability.

Bayes Rule (

$$y^* = \arg\max_j p(x \mid y = j)\, p(y = j)$$

$$\boxed{g^*(x) = \arg\max_j \log p(x \mid y = j) + \log p(y = j)}$$

## Simple example

2-class problem $\{0, 1\}$

$$g(x) = \arg\max_j \log p(x \mid y = j) + \log p(y = j)$$

pick class 0 if:

$$\log p(x \mid 0) + \log p(0) > \log p(x \mid 1) + \log p(1)$$

$$\log p(x \mid 0) - \log p(x \mid 1) > \log p(1) - \log p(0)$$

$$\underbrace{\frac{p(x \mid 0)}{p(x \mid 1)}}_{\text{likelihood ratio}} > \frac{p(1)}{p(0)} = T$$
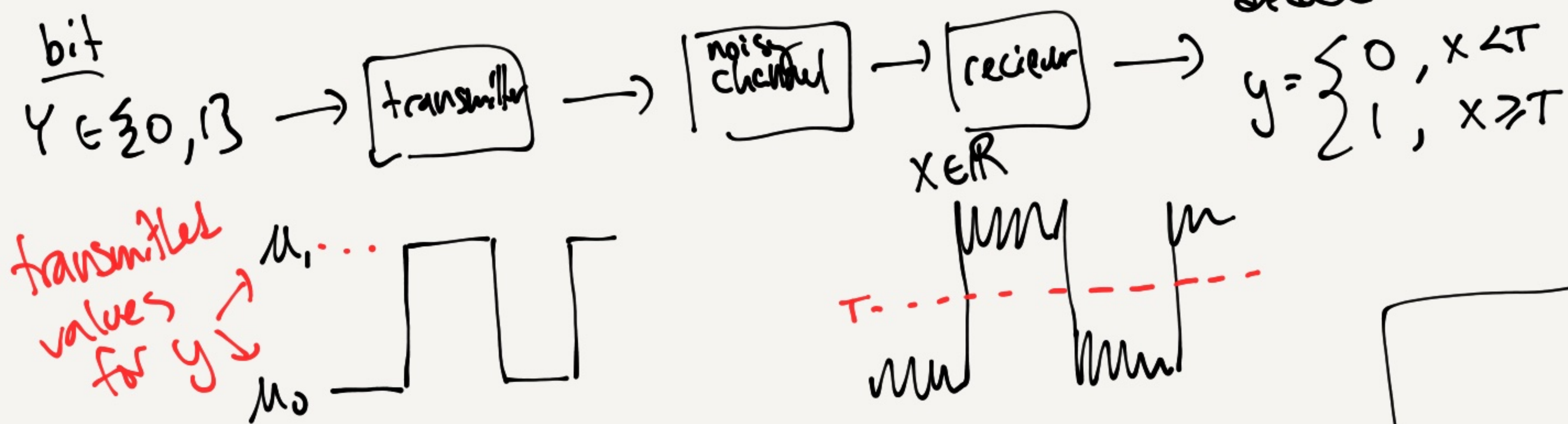
## Summary

- for 0-1 loss function
  - → BDR is MAP (pick maximum posterior)
  - → conditional risk = prob. of error for $x$
  - → Risk = prob. of error.
  - → BDR minimizes prob. of error (no other decision rule is better!)

Caveat: all the modeling assumptions are correct. (CCD & prior are correct)

This is a generative classification model

- CCD are learned from data, decision rule is computed from CCDs.

# Example: Noisy Channel

bit
$Y \in \{0,1\}$ → [transmitter] → [noisy channel] → [reciever] → decode $y = \begin{cases} 0, & x < T \\ 1, & x \geq T \end{cases}$

$X \in \mathbb{R}$

transmitted values for $y$ → $\mu_1 \cdots$ ⊓⊔ ... $\mu_0$

## What is threshold $T$?

Given measurment $X$, recover bit $Y$.

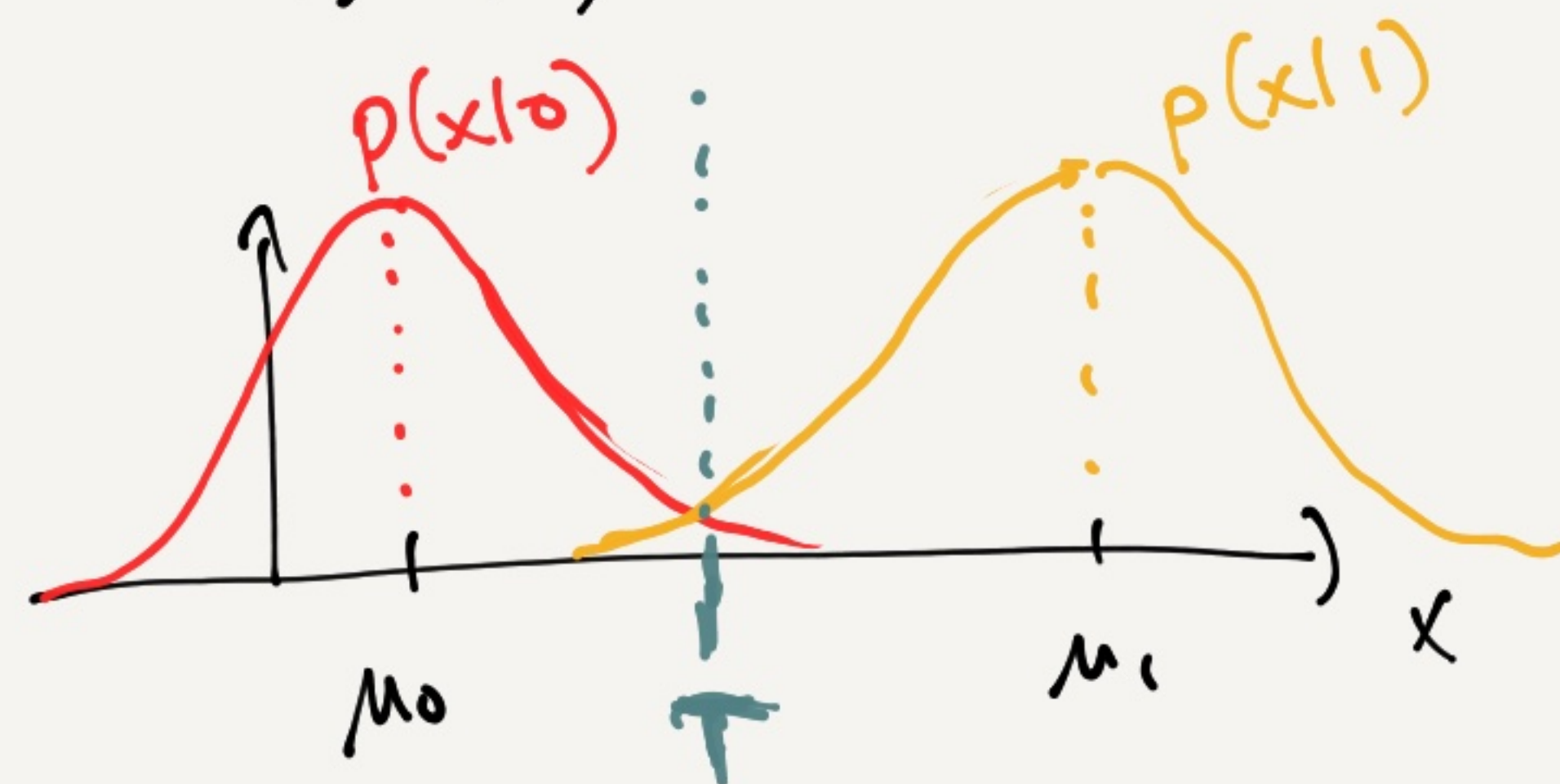- prior: $p(y=0) = p(y=1) = \frac{1}{2}$

- CCD: Assume Gaussian additive noise

$$X = \mu_y + \epsilon \quad , \quad \epsilon \sim N(0, \sigma^2)$$

meas. ↑ xmitted value ↑ channel noise

$$\begin{bmatrix} p(x|y=0) = N(x \mid \mu_0, \sigma^2) \\ p(x|y=1) = N(x \mid \mu_1, \sigma^2) \end{bmatrix}$$

Assume $\mu_0 < \mu_1$

$p(x|0)$    $p(x|1)$

$\mu_0$    $T$    $\mu_1$    $x$

# BDR w/ 0-1 loss

$$y^* = \underset{j}{\arg\max} \; \log p(x|j) + \log p(j)$$

$$= \underset{j}{\arg\max} \; -\frac{1}{2\sigma^2}(x-\mu_j)^2 - \frac{1}{2}\log 2\pi - \frac{1}{2}\log \sigma^2 + \log \frac{1}{2}$$

constant     constant

$$= \underset{j}{\arg\max} \; -(x-\mu_j)^2$$

$$= \underset{j}{\arg\max} \; -(x^2 - 2x\mu_j + \mu_j^2)$$

constant

$$= \underset{j}{\arg\min} \; \mu_j^2 - 2x\mu_j$$

choose $y^* = 0$ when

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

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

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

$$\Rightarrow \quad x < \frac{\mu_1^2 - \mu_0^2}{2(\mu_1 - \mu_0)} = \frac{\mu_1 + \mu_0}{2}$$

threshold is the midpoint between $\mu_1$ & $\mu_0$

Assumptions are explicit
1) 0-1 cost BDR
2) uniform classpriors
3) Gaussian additive noise

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

## BDR:
pick 0 if:

$$X < \frac{\mu_1 + \mu_0}{2} + \frac{\sigma^2}{\mu_1 - \mu_0} \log \frac{p(y=0)}{p(y=1)}$$

$\underbrace{\phantom{\frac{\mu_1 + \mu_0}{2}}}$ same as before

increase threshold if $p(0) > p(1)$, i.e. 0 is more frequent.

$\frac{\mu_1 - \mu_0}{\sigma^2}$ = normalized distance b/twn means

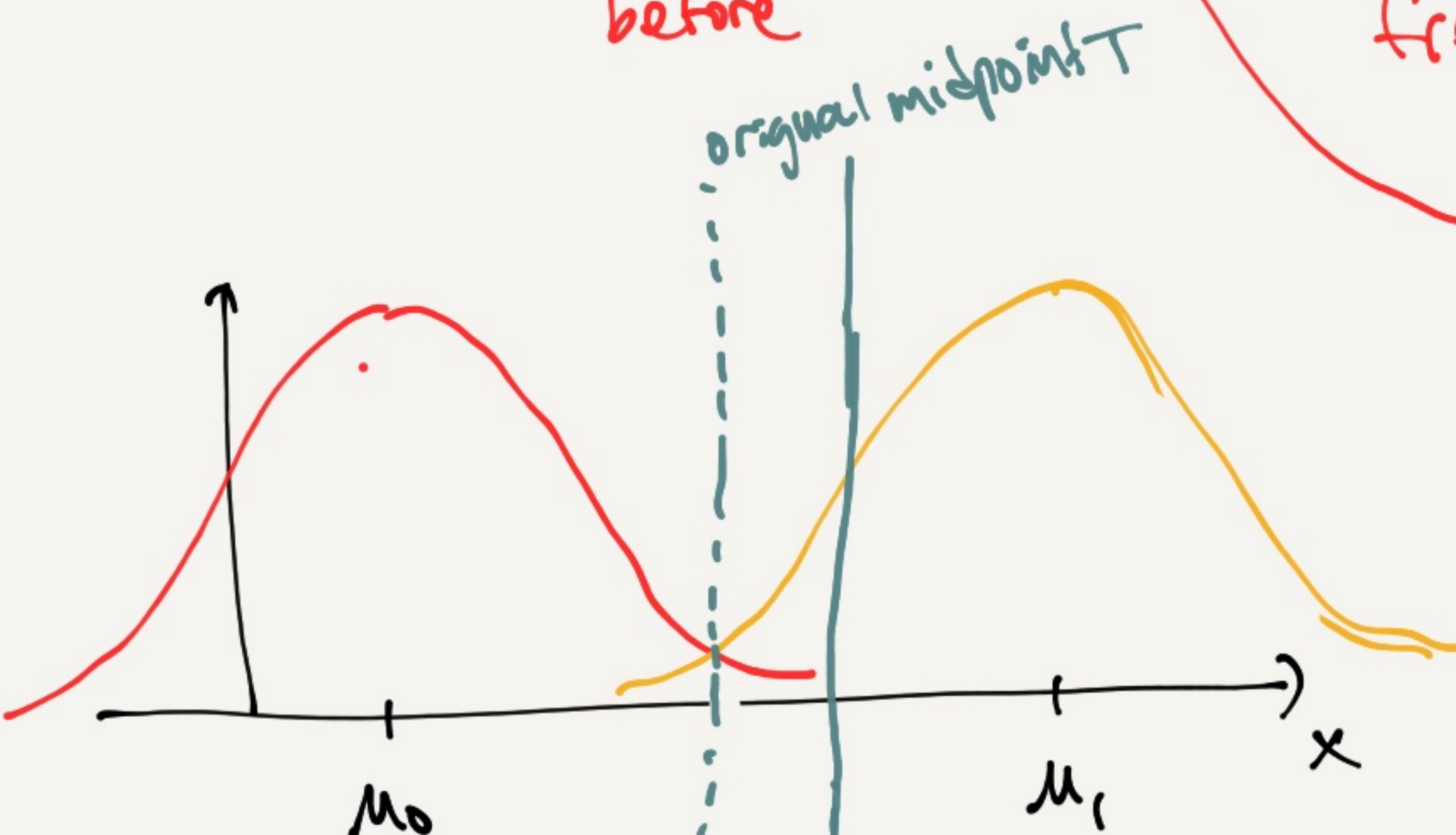$\Rightarrow \frac{\sigma^2}{\mu_1 - \mu_0} = \frac{1}{\text{normalized distance}}$

$\Rightarrow$ if means are far apart, then move T a little (ignore prior)

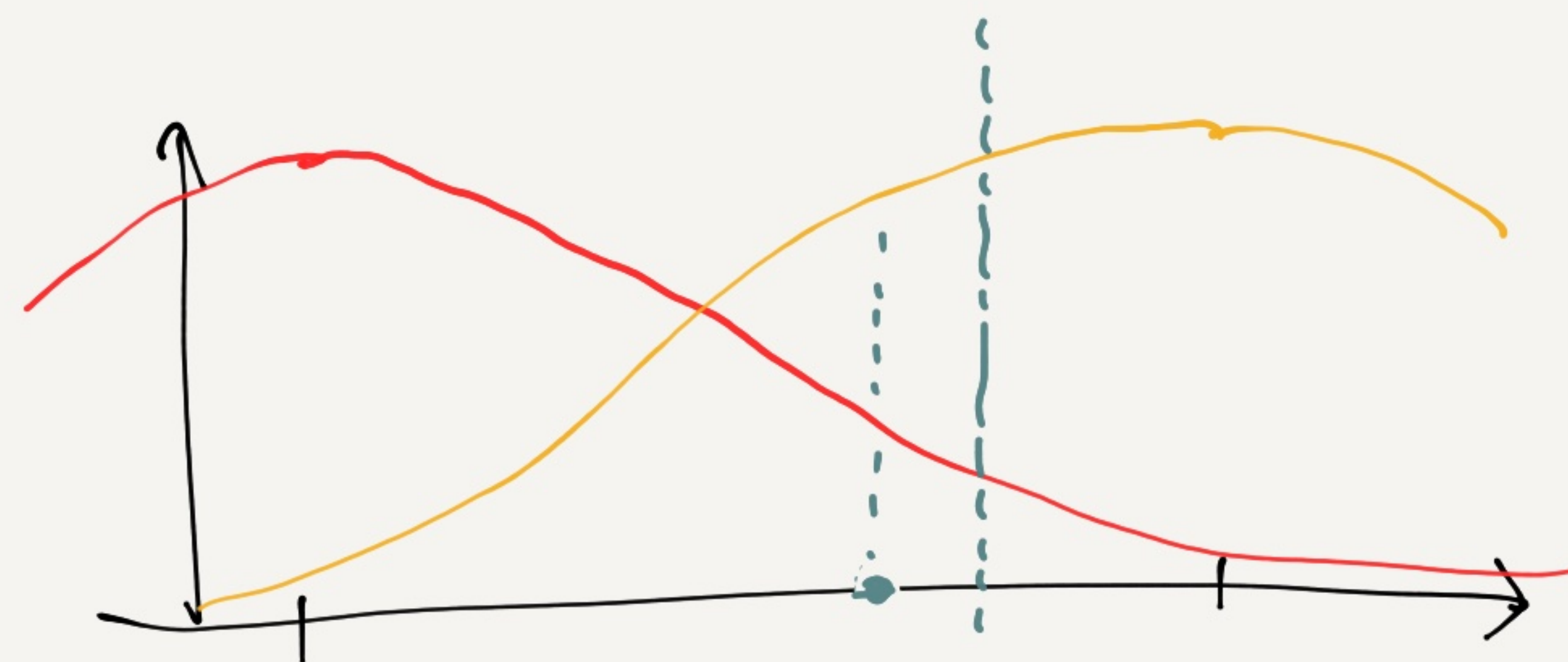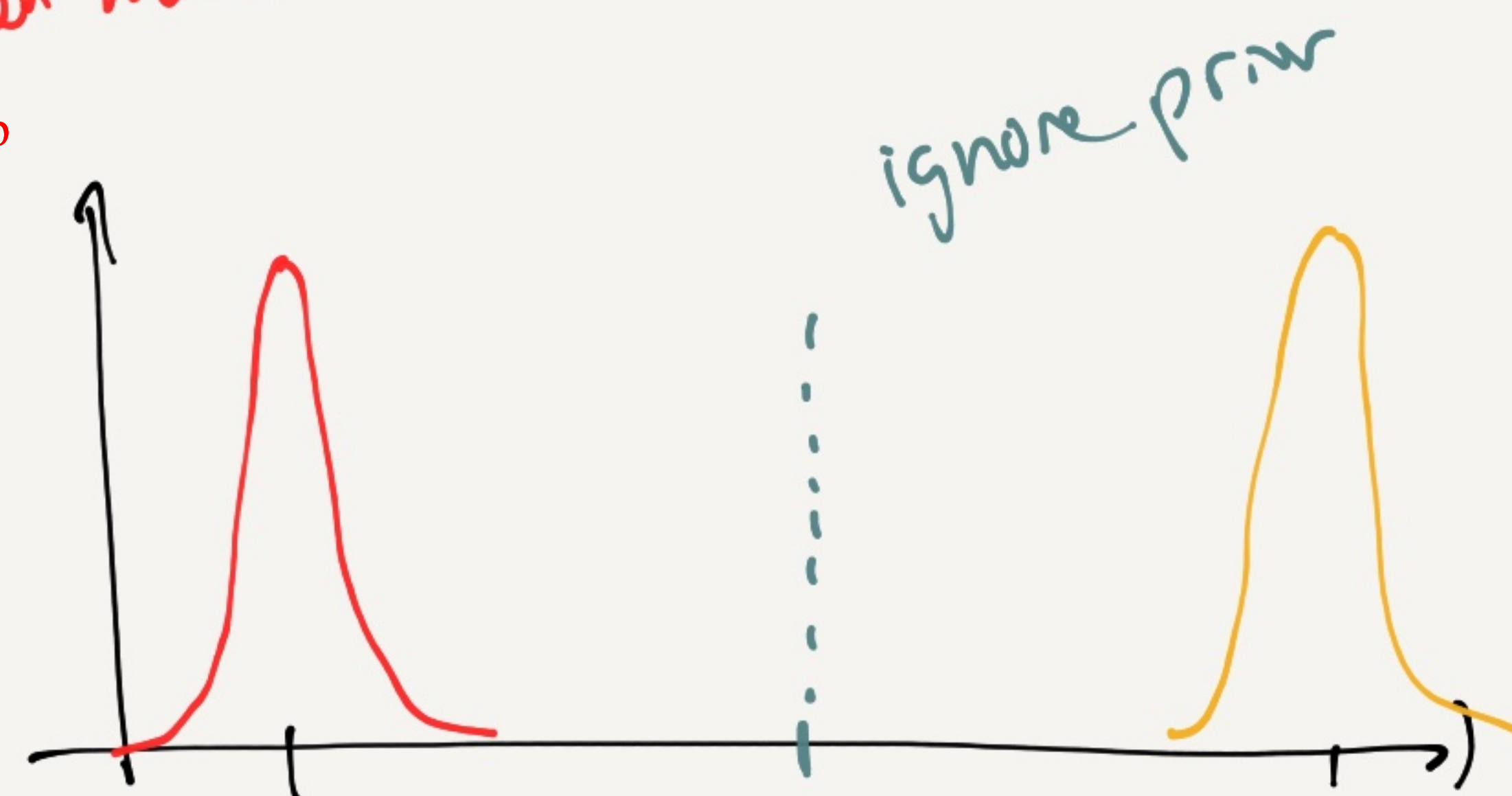$\Rightarrow$ if means are close, then move T a lot (use priors)



original midpoint T

$\mu_0$     $\mu_1$    x

$\frac{p(0)}{=}$
$p(1)$

$p(0) > p(1) \Rightarrow$ shift threshold "capture" more 0's.

ignore prior



use priors

(Need to be very certain it is 1 to choose it)

# Gaussian Classifier

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

prior $p(Y=j) = \pi_j$ ⟵ can be estimated from data.

$X \in \mathbb{R}^d$

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

BDR: $g(x) = \underset{j}{\arg\max} \ \log p(x|j) + \log p(j)$

$\quad = \underset{j}{\arg\max} \ -\frac{1}{2}\|x - \mu_j\|_{\Sigma_j}^2 - \frac{1}{2}\log|\Sigma_j| - \frac{1}{2}\log(2\pi)^d + \log \pi_j$

## Special cases

i) assume $\Sigma_j = \sigma^2 I$   (shared isotropic covariances)

Define $\quad g_j(x) = \omega_j^T x + b_j$

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

$\Rightarrow g^*(x) = \underset{j}{\arg\max} \ g_j(x)$

# Geometric Meaning

classes $i \, \& \, j$ share a boundary if

$\quad g_i(x) = g_j(x)$

$\Rightarrow \omega_i^T x + b_i = \omega_j^T x + b_j$

$\quad \vdots$

$\quad \omega^T (x - x_0) = 0$ ⟵ hyperplane passing through $x_0$ & normal to $\omega$.

$\begin{cases} \omega = \frac{1}{\sigma^2}(\mu_i - \mu_j) \ ⟵ \text{vector btwn } \mu_j \& \mu_i \\ \\ x_0 = \frac{\mu_i + \mu_j}{2} + (\mu_j - \mu_i)\left[ \frac{\sigma^2}{\|\mu_i - \mu_j\|^2} \cdot \log \frac{\pi_i}{\pi_j} \right] \end{cases}$

$\underbrace{\frac{\mu_i+\mu_j}{2}}_{\text{midpoint btwn means}}$ $\underbrace{(\mu_j - \mu_i)}_{\substack{\text{vector from} \\ \mu_i \text{ to } \mu_j \\ \text{class } j}}$ $\underbrace{\left[\frac{\sigma^2}{\|\mu_i-\mu_j\|^2}\right.}_{\substack{1/\text{normalized} \\ \text{distance}}} \underbrace{\left.\log\frac{\pi_i}{\pi_j}\right]}_{\text{priors}}$

if $\pi_i > \pi_j \Rightarrow > 0$



$\pi_i > \pi_j \Rightarrow$ move the decision boundary away from $\mu_i$ to capture more space.