

ECE595 / STAT598: Machine Learning I

Lecture 10 Minimum Probability of Error Rule

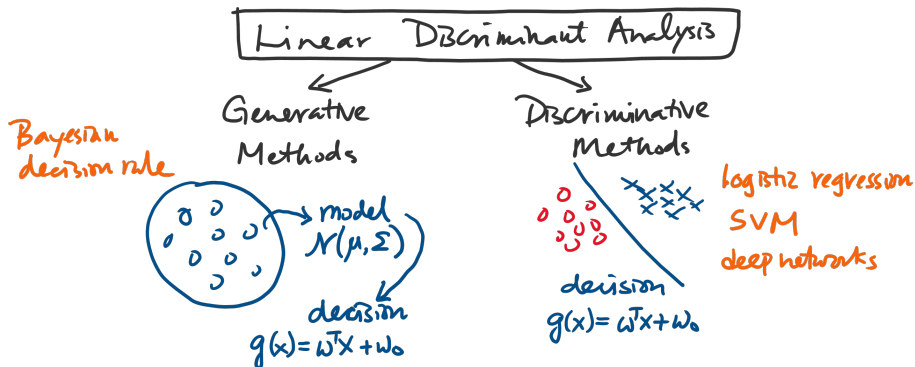
Spring 2020

Stanley Chan

School of Electrical and Computer Engineering
Purdue University



Overview



- In linear discriminant analysis (LDA), there are generally two types of approaches
- **Generative approach:** Estimate model, then define the classifier
- **Discriminative approach:** Directly define the classifier

Outline

Generative Approaches

- Lecture 9 Bayesian Decision Rules
- **Lecture 10 Evaluating Performance**
- Lecture 11 Parameter Estimation
- Lecture 12 Bayesian Prior
- Lecture 13 Connecting Bayesian and Linear Regression

Today's Lecture

- **The Three Cases**
 - $\Sigma_j = \sigma^2 I$ (Last Lecture)
 - **$\Sigma_j = \Sigma$**
 - **General Σ_j**
- Evaluating Performance
 - Probability of Error
 - Bayesian Decision Rule is also Minimum Error Rule
 - ROC Curve

Three Cases of Gaussians

Discriminant function of Gaussian:

$$\begin{aligned}g_i(\mathbf{x}) &= \log p_{\mathbf{X}|Y}(\mathbf{x}|i) + \log \pi_i \\ &= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1}(\mathbf{x} - \boldsymbol{\mu}_i) - \frac{1}{2} \log |\boldsymbol{\Sigma}_i| + \log \pi_i.\end{aligned}$$

- $\boldsymbol{\Sigma}_i = \sigma^2 \mathbf{I}$
 - All Gaussians have the same covariance matrix
 - The covariance matrix is diagonal and same variance
- $\boldsymbol{\Sigma}_i = \boldsymbol{\Sigma}$
 - All Gaussians have the same covariance matrix
 - The covariance matrix can be anything
- arbitrary $\boldsymbol{\Sigma}_i$
 - Any positive semi-definite covariance matrix

Last Lecture

Theorem

If $\Sigma_i = \sigma^2 \mathbf{I}$, then the separating hyperplane is given by

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = 0,$$

where

$$\mathbf{w} = \frac{\boldsymbol{\mu}_i - \boldsymbol{\mu}_j}{\sigma^2}, \quad \text{and} \quad w_0 = -\frac{\|\boldsymbol{\mu}_i\|^2 - \|\boldsymbol{\mu}_j\|^2}{2\sigma^2} + \log \frac{\pi_i}{\pi_j}.$$

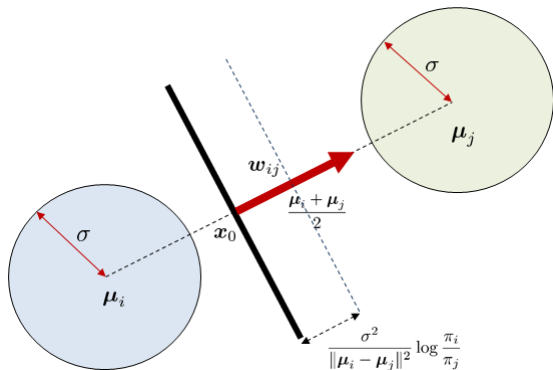
- You tell me the two Gaussians: $\boldsymbol{\mu}_i, \boldsymbol{\mu}_j, \pi_i, \pi_j, \sigma$
- I return you a separating hyperplane

$$g(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

- This is the best possible hyperplane according to posterior distribution

Case 1: $\Sigma_i = \sigma^2 I$: Geometry

$$\mathbf{w} = \frac{\boldsymbol{\mu}_i - \boldsymbol{\mu}_j}{\sigma^2}, \quad \text{and} \quad \mathbf{x}_0 = \frac{\boldsymbol{\mu}_i + \boldsymbol{\mu}_j}{2} - \frac{\sigma^2}{\|\boldsymbol{\mu}_i - \boldsymbol{\mu}_j\|^2} \left(\log \frac{\pi_i}{\pi_j} \right) (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j),$$



Case 2: $\Sigma_i = \Sigma$

Make all $\Sigma_i = \Sigma$.

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_i) - \frac{1}{2} \log |\boldsymbol{\Sigma}| + \log \pi_i.$$

In this case, we can show that

$$\begin{aligned} g_i(\mathbf{x}) &= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_i) - \frac{1}{2} \log |\boldsymbol{\Sigma}| + \log \pi_i \\ &= -\frac{1}{2} \left(\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x} - 2\boldsymbol{\mu}_i^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}_i^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_i \right) + \log \pi_i \\ &= -\frac{1}{2} \left(\cancel{\mathbf{x}^T \boldsymbol{\Sigma}^{-1} \mathbf{x}} - 2\boldsymbol{\mu}_i^T \boldsymbol{\Sigma}^{-1} \mathbf{x} + \boldsymbol{\mu}_i^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_i \right) + \log \pi_i \\ &= \underbrace{\boldsymbol{\mu}_i^T \boldsymbol{\Sigma}^{-1} \mathbf{x}}_{\mathbf{w}_i^T \mathbf{x}} \quad \underbrace{-\frac{1}{2} \boldsymbol{\mu}_i^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_i + \log \pi_i}_{w_{i0}} \end{aligned}$$

Case 2: $\Sigma_i = \Sigma$

The discriminant function is therefore:

$$\begin{aligned}g(\mathbf{x}) &= g_i(\mathbf{x}) - g_j(\mathbf{x}) \\&= (\mathbf{w}_i - \mathbf{w}_j)^T \mathbf{x} + (w_{i0} - w_{j0}) \\&= \underbrace{[\Sigma^{-1}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)]^T \mathbf{x}}_{\mathbf{w}^T \mathbf{x}} - \underbrace{\frac{1}{2} (\boldsymbol{\mu}_i^T \Sigma^{-1} \boldsymbol{\mu}_i - \boldsymbol{\mu}_j^T \Sigma^{-1} \boldsymbol{\mu}_j)}_{w_0} + \log \frac{\pi_i}{\pi_j}.\end{aligned}$$

If we want to write $g(\mathbf{x}) = \mathbf{w}^T (\mathbf{x} - \mathbf{x}_0)$, we can show that

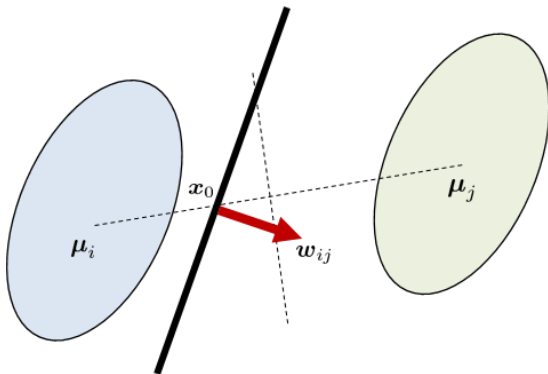
$$\begin{aligned}\mathbf{w} &= \Sigma^{-1}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j), \\ \mathbf{x}_0 &= \frac{\boldsymbol{\mu}_i + \boldsymbol{\mu}_j}{2} - \frac{\log \frac{\pi_i}{\pi_j}}{(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)^T \Sigma^{-1} (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)} (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j).\end{aligned}$$

Case 2: $\Sigma_i = \Sigma$: Geometry

The discriminant function is $g(\mathbf{x}) = \mathbf{w}^T(\mathbf{x} - \mathbf{x}_0)$, with

$$\mathbf{w} = \Sigma^{-1}(\mu_i - \mu_j),$$

$$\mathbf{x}_0 = \frac{\mu_i + \mu_j}{2} - \frac{\log \frac{\pi_i}{\pi_j}}{(\mu_i - \mu_j)^T \Sigma^{-1}(\mu_i - \mu_j)} (\mu_i - \mu_j).$$



Case 3: Arbitrary Σ_i

This is the most general setting.

$$g_i(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1}(\mathbf{x} - \boldsymbol{\mu}_i) - \frac{1}{2} \log |\boldsymbol{\Sigma}_i| + \log \pi_i.$$

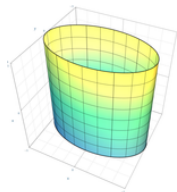
We can show that

$$\begin{aligned} g_i(\mathbf{x}) &= \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1}(\mathbf{x} - \boldsymbol{\mu}_i) + \frac{1}{2} \log |\boldsymbol{\Sigma}_i| - \log \pi_i \\ &= \underbrace{\frac{1}{2} \mathbf{x}^T \boldsymbol{\Sigma}_i^{-1} \mathbf{x}}_{\frac{1}{2} \mathbf{x}^T \mathbf{W}_i \mathbf{x}} - \underbrace{\boldsymbol{\mu}_i^T \boldsymbol{\Sigma}_i^{-1} \mathbf{x}}_{\mathbf{w}_i^T \mathbf{x}} + \underbrace{\frac{1}{2} \boldsymbol{\mu}_i^T \boldsymbol{\Sigma}_i^{-1} \boldsymbol{\mu}_i + \frac{1}{2} \log |\boldsymbol{\Sigma}_i| - \log \pi_i}_{w_{i0}}. \end{aligned}$$

Therefore, the discriminant function is

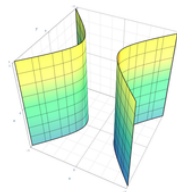
$$\begin{aligned} g(\mathbf{x}) &= g_i(\mathbf{x}) - g_j(\mathbf{x}) \\ &= \frac{1}{2} \mathbf{x}^T (\mathbf{W}_i - \mathbf{W}_j) \mathbf{x} + (\mathbf{w}_i - \mathbf{w}_j)^T \mathbf{x} + (w_{i0} - w_{j0}). \end{aligned}$$

Case 3: Arbitrary Σ_j : Geometry



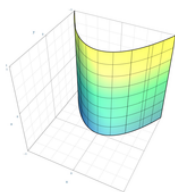
$$W = \begin{bmatrix} a & 0 & 0 \\ 0 & b & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$w = 0$$



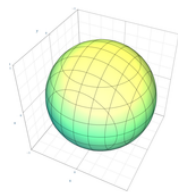
$$W = \begin{bmatrix} a & 0 & 0 \\ 0 & -b & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$w = 0$$



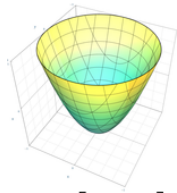
$$W = \begin{bmatrix} a & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$w \neq 0$$



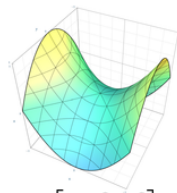
$$W = \begin{bmatrix} a & 0 & 0 \\ 0 & a & 0 \\ 0 & 0 & a \end{bmatrix}$$

$$w = 0$$



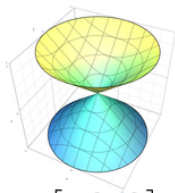
$$W = \begin{bmatrix} a & 0 & 0 \\ 0 & a & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$w \neq 0$$



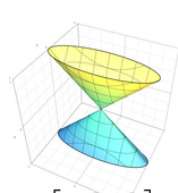
$$W = \begin{bmatrix} a & 0 & 0 \\ 0 & -b & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$w \neq 0$$



$$W = \begin{bmatrix} a & 0 & 0 \\ 0 & a & 0 \\ 0 & 0 & -b \end{bmatrix}$$

$$w = 0$$



$$W = \begin{bmatrix} a & 0 & 0 \\ 0 & b & 0 \\ 0 & 0 & -c \end{bmatrix}$$

$$w = 0$$

Case 3: Arbitrary Σ_j : Computing

Quadratic classifier is not a big deal, on computer!

$$\begin{aligned}g(\mathbf{x}) &= g_i(\mathbf{x}) - g_j(\mathbf{x}) \\ &= \frac{1}{2}\mathbf{x}^T(\mathbf{W}_i - \mathbf{W}_j)\mathbf{x} + (\mathbf{w}_i - \mathbf{w}_j)^T\mathbf{x} + (w_{i0} - w_{j0}).\end{aligned}$$

Recall: A hypothesis function is

$$h(\mathbf{x}) = \begin{cases} 1, & \text{if } g(\mathbf{x}) > 0 \\ 0, & \text{if } g(\mathbf{x}) < 0 \\ \text{either,} & \text{if } g(\mathbf{x}) = 0 \end{cases}$$

- All you need is to check whether $g(\mathbf{x}) > 0$ or $g(\mathbf{x}) < 0$.
- No need to obtain a closed form solution.

Outline

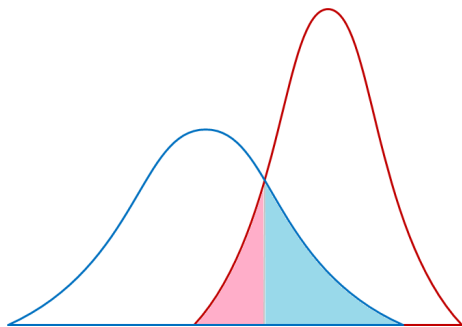
Generative Approaches

- Lecture 9 Bayesian Decision Rules
- **Lecture 10 Evaluating Performance**
- Lecture 11 Parameter Estimation
- Lecture 12 Bayesian Prior
- Lecture 13 Connecting Bayesian and Linear Regression

Today's Lecture

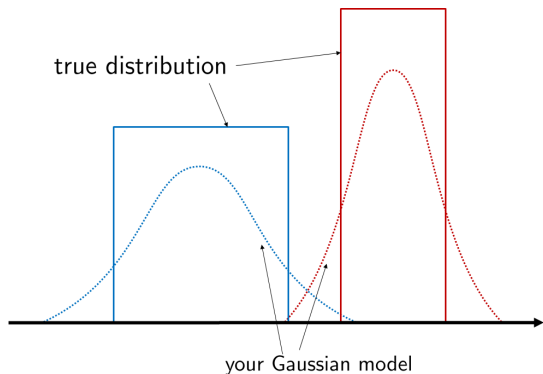
- The Three Cases
 - $\Sigma_j = \sigma^2 I$ (Last Lecture)
 - $\Sigma_j = \Sigma$
 - General Σ_j
- **Evaluating Performance**
 - **Probability of Error**
 - **Bayesian Decision Rule is also Minimum Error Rule**
 - **ROC Curve**

How Well does Bayesian Decision Rule Perform?



- Gaussian is a model.
- If the underlying true distribution is Gaussian, then there is always overlapping.
- No matter how you do, there is always classification error.

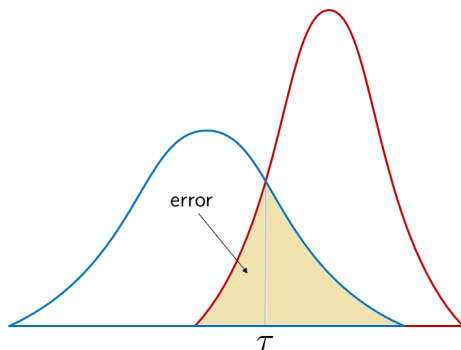
What if the distributions do not match



- If the true distribution is not overlapping, then the overlapping of the Gaussian is not a problem.
- In this lecture we assume the true distribution is indeed Gaussian.

Probability of Error

Suppose we have found a decision boundary, which is a threshold τ in 1D.



Probability of error is

$$P_e = \mathbb{P}[X < \tau \text{ and } Y = 1] + \mathbb{P}[X > \tau \text{ and } Y = 0]$$

Minimize the Probability of Error

Let us do some calculation:

$$\begin{aligned}P_e &= \mathbb{P}[X < \tau \text{ and } Y = 1] + \mathbb{P}[X > \tau \text{ and } Y = 0] \\&= \mathbb{P}[X < \tau | Y = 1]\mathbb{P}[Y = 1] + \mathbb{P}[X > \tau | Y = 0]\mathbb{P}[Y = 0] \\&= \pi_1 \mathcal{N}(X < \tau | \mathcal{C}_1) + \pi_0 \mathcal{N}(X > \tau | \mathcal{C}_0) \\&= \pi_1 \int_{-\infty}^{\tau} \frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}} dx + \pi_0 \int_{\tau}^{\infty} \frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{(x-\mu_0)^2}{2\sigma_0^2}} dx\end{aligned}$$

Can we minimize the P_e by finding a good τ ?

$$\begin{aligned}\frac{d}{d\tau} P_e &= \pi_1 \frac{d}{d\tau} \left(\int_{-\infty}^{\tau} \frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}} dx \right) \\&\quad + \pi_0 \frac{d}{d\tau} \left(\int_{\tau}^{\infty} \frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{(x-\mu_0)^2}{2\sigma_0^2}} dx \right)\end{aligned}$$

Minimize the Classification Error

Fundamental Theorem of Calculus:

$$\frac{d}{d\tau} \int_{-\infty}^{\tau} f(x) dx = f(\tau)$$

Therefore, in our problem,

$$\frac{d}{d\tau} \left(\int_{-\infty}^{\tau} \frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}} dx \right) = \frac{1}{\sqrt{2\pi\sigma_1^2}} e^{-\frac{(\tau-\mu_1)^2}{2\sigma_1^2}}$$

$$\frac{d}{d\tau} \left(\int_{\tau}^{\infty} \frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{(x-\mu_0)^2}{2\sigma_0^2}} dx \right) = -\frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{(\tau-\mu_0)^2}{2\sigma_0^2}}$$

Special Case: $\sigma_1 = \sigma_0 = \sigma$

Let us assume that $\sigma_1 = \sigma_0 = \sigma$. Then,

$$\begin{aligned}\frac{d}{d\tau} P_e &= \pi_1 \frac{d}{d\tau} \int_{-\infty}^{\tau} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu_1)^2}{2\sigma^2}} dx + \pi_0 \frac{d}{d\tau} \int_{\tau}^{\infty} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu_0)^2}{2\sigma^2}} dx \\ &= \pi_1 \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\tau-\mu_1)^2}{2\sigma^2}} - \pi_0 \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\tau-\mu_0)^2}{2\sigma^2}} = 0\end{aligned}$$

So finally, we can equate the two sides, and get

$$\begin{aligned}\pi_1 \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\tau-\mu_1)^2}{2\sigma^2}} &= \pi_0 \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(\tau-\mu_0)^2}{2\sigma^2}} \\ \log \left\{ \pi_1 e^{-\frac{(\tau-\mu_1)^2}{2\sigma^2}} \right\} &= \log \left\{ \pi_0 e^{-\frac{(\tau-\mu_0)^2}{2\sigma^2}} \right\} \\ \tau &= \frac{\mu_1 - \mu_0}{2} - \frac{\sigma^2}{\mu_1 - \mu_0} \log \frac{\pi_1}{\pi_0}.\end{aligned}$$

Optimality of Bayesian Decision Rule

Therefore,

$$\tau = \frac{\mu_1 - \mu_0}{2} - \frac{\sigma^2}{\mu_1 - \mu_0} \log \frac{\pi_1}{\pi_0},$$

and this means that

$$h(x) = \begin{cases} 1, & \text{if } x > \frac{\mu_1 - \mu_0}{2} - \frac{\sigma^2}{\mu_1 - \mu_0} \log \frac{\pi_1}{\pi_0}, \\ 0, & \text{if } x < \frac{\mu_1 - \mu_0}{2} - \frac{\sigma^2}{\mu_1 - \mu_0} \log \frac{\pi_1}{\pi_0}, \\ \text{either,} & \text{if } x = \frac{\mu_1 - \mu_0}{2} - \frac{\sigma^2}{\mu_1 - \mu_0} \log \frac{\pi_1}{\pi_0}. \end{cases}$$

This is the **exact same result** as the Bayesian decision rule.

- Bayesian decision is optimal in posterior
- Bayesian decision is optimal in probability of error
- Optimal \neq error free
- You still have error if you have infinite training data

Two Levels of Optimality

Error due to **Finite Training Samples**

- You need to estimate μ and Σ
- More samples, better estimate
- Controlled by Hoeffding inequality

Error due to **Linear Non-Separability**

- Error even if you have infinite training data
- Fundamental limit of the model
- Gaussians are not linearly separable
- Minimize error probability \neq zero error

Risk and Probability of Error

Is Your “Optimal” Really Optimal?

$$P_e = \mathbb{P}[X < \tau \mid Y = 1]\pi_1 + \mathbb{P}[X > \tau \mid Y = 0]\pi_0.$$

- $Y = 1$: disease present. $Y = 0$: normal.
- $X < \tau$: report no disease. $X > \tau$: report disease.
- $\mathbb{P}[X < \tau \mid Y = 1]$: There is a disease, but you cannot find it.
- $\mathbb{P}[X > \tau \mid Y = 0]$: There is no disease, but you say there is.
- Which one is more serious?
- We need to define **Risk**

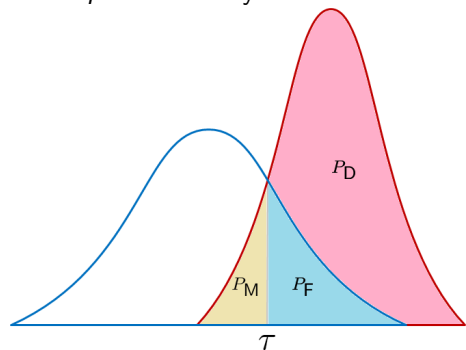
$$\underbrace{R}_{\text{risk}} = C_{01} \underbrace{\mathbb{P}[X < \tau \mid Y = 1]\pi_1}_{\text{miss}} + C_{10} \underbrace{\mathbb{P}[X > \tau \mid Y = 0]\pi_0}_{\text{false alarm}}.$$

- C_{01} : cost of miss. C_{10} : cost of false alarm.

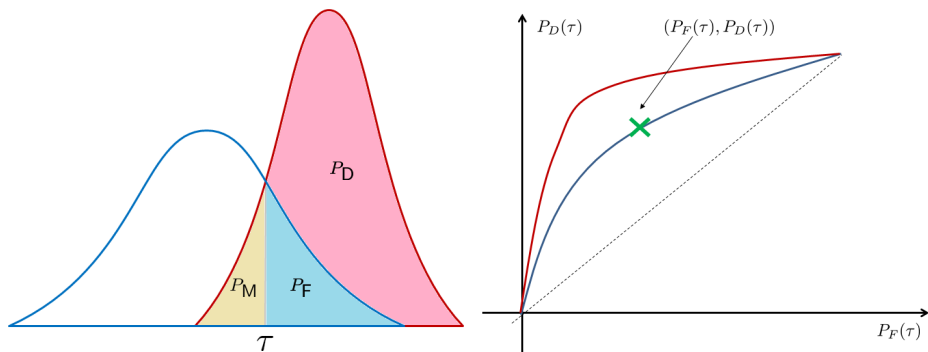
False Alarm and Detection

$$P_e = C_{01} \underbrace{\mathbb{P}[X < \tau \mid Y = 1]}_{P_M} \pi_1 + C_{10} \underbrace{\mathbb{P}[X > \tau \mid Y = 0]}_{P_F} \pi_0.$$

- P_M : Probability of Miss
- $P_D = 1 - P_M$: Probability of Detection
- P_F : Probability of False Alarm

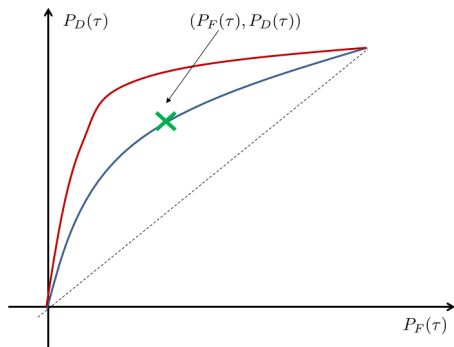


Receive Operating Characteristic (ROC) Curve



- They are all functions of τ .
- If you change τ , they will change correspondingly.
- Can plot $P_F(\tau)$ and $P_D(\tau)$ as τ changes.
- This called a Receive Operating Characteristic (ROC) curve.

Interpreting the ROC



- Higher ROC is better: Same level of false alarm, higher detection rate
- One classifier has one ROC. One ROC curve for **all** τ .
- Ways to improve ROC:
 - Train with more sample so that you have less training error
 - Find a better model than Gaussian
 - Find a better decision rule than linear classifier

Reading List

Bayesian Decision Rule

- Bishop, Pattern Recognition and Machine Learning, Chapter 4.1
- Duda, Hart and Stork's Pattern Classification, Chapter 2.1, 2.2, 2.6
- Stanford CS 229 Generative Algorithms
<http://cs229.stanford.edu/notes/cs229-notes2.pdf>

Probability of Error

- Duda, Hart and Stork's *Pattern Classification*, Chapter 2.7, 3.1.
- Poor, *Intro to Signal Estimation and Detection*, Chapter 2.

ROC Curve

- ECE645 Note. <https://engineering.purdue.edu/ChanGroup/ECE645Notes/StudentLecture02.pdf>

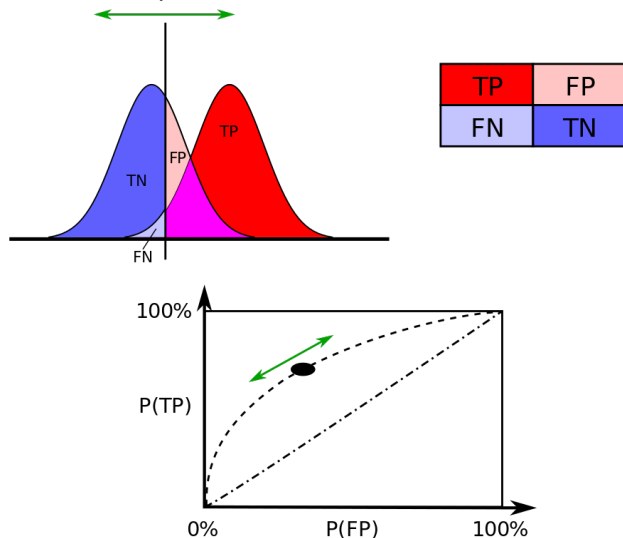
Appendix

Q&A 1 Is ROC curve always concave?

- No.
- The ROC curve is concave when the decision rule comes from Neyman-Pearson. See V. Poor's Intro to Signal Estimation and Detection Chapter 2.
- You can always create a decision rule that performs terribly on one part and very well on the other part.
- Most algorithms will generate ROC curves that shows stair-case effects.
- It is important to report the entire ROC curve.

Q&A 2 Type 1 and Type 2 Error

You probably have heard these two terms before. Below is an image I found from Wikipedia



Q&A 2 Type 1 and Type 2 Error

- A typical naming issue between statistics and engineering
- False Positive = Type 1 Error = False Alarm
- False Negative = Type 2 Error = Miss
- Both are **probabilities**
- In the previous figure, we see that a decision rule is the location of the black line.
- If you move the black line left and right, you will get a point on the ROC curve.

Q&A 3 What is the difference between ROC and Precision-Recall Curve (PRC)?

- ROC: Receiver Operating Curve.
 - Detection: $\text{No. TP} / (\text{No. TP} + \text{No. FN})$
 - False Alarm: $\text{No. FP} / (\text{No. FP} + \text{No. TN})$
 - ROC Curve always goes up.
- PRC: Precision Recall Curve.
 - Precision: $\text{No. TP} / (\text{No. TP} + \text{No. FP})$
 - Recall: $\text{No. TP} / (\text{No. TP} + \text{No. FN})$
 - PRC Curve always goes down.
- There is generally no definitive answer to which of ROC or PRC is better. However, according to the following paper, ROC is better for datasets with balanced classes whereas PRC is better for skewed classes. <https://www.biostat.wisc.edu/~page/rocpr.pdf>