

Introduction to Detection Theory

Reading:

- Ch. 3 in Kay-II.
- Notes by Prof. Don Johnson on detection theory,
see <http://www.ece.rice.edu/~dhj/courses/elec531/notes5.pdf>.
- Ch. 10 in Wasserman.

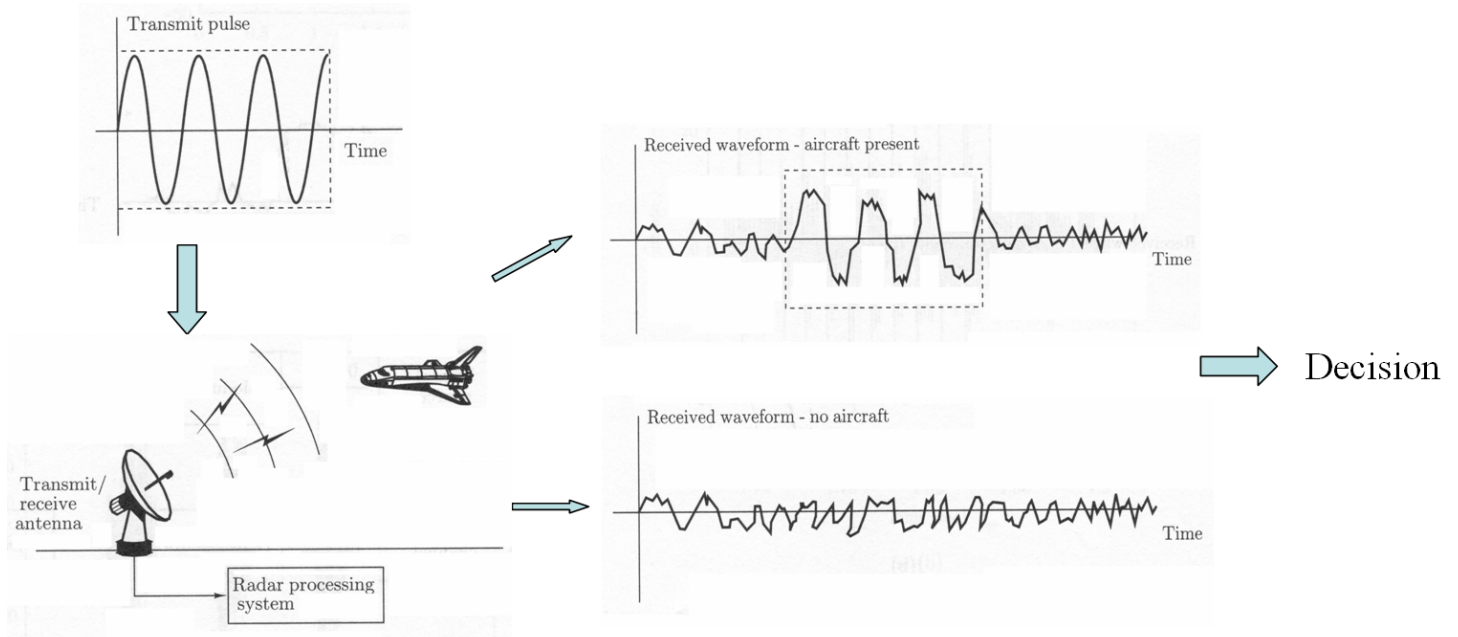
Introduction to Detection Theory (cont.)

We wish to make a decision on a signal of interest using noisy measurements. Statistical tools enable systematic solutions and optimal design.

Application areas include:

- Communications,
- Radar and sonar,
- Nondestructive evaluation (NDE) of materials,
- Biomedicine, etc.

Example: Radar Detection. We wish to decide on the presence or absence of a target.



Introduction to Detection Theory

We assume a parametric measurement model $p(x | \theta)$ [or $p(x; \theta)$, which is the notation that we sometimes use in the classical setting].

In point estimation theory, we estimated the parameter $\theta \in \Theta$ given the data x .

Suppose now that we choose Θ_0 and Θ_1 that form a *partition* of the parameter space Θ :

$$\Theta_0 \cup \Theta_1 = \Theta, \quad \Theta_0 \cap \Theta_1 = \emptyset.$$

In detection theory, we wish to identify *which* hypothesis is true (i.e. make the *appropriate decision*):

$$\mathcal{H}_0 : \theta \in \Theta_0, \quad \text{null hypothesis}$$

$$\mathcal{H}_1 : \theta \in \Theta_1, \quad \text{alternative hypothesis.}$$

Terminology: If θ can only take two values,

$$\Theta = \{\theta_0, \theta_1\}, \quad \Theta_0 = \{\theta_0\}, \quad \Theta_1 = \{\theta_1\}$$

we say that the hypotheses are *simple*. Otherwise, we say that they are *composite*.

Composite Hypothesis Example: $\mathcal{H}_0 : \theta = 0$ versus $\mathcal{H}_1 : \theta \in (0, \infty)$.

The Decision Rule

We wish to design a decision rule (function) $\phi(\mathbf{x}) : \mathcal{X} \rightarrow (0, 1)$:

$$\phi(\mathbf{x}) = \begin{cases} 1, & \text{decide } \mathcal{H}_1, \\ 0, & \text{decide } \mathcal{H}_0. \end{cases}$$

which partitions the data space \mathcal{X} [i.e. the support of $p(\mathbf{x} | \theta)$] into two regions:

$$\text{Rule } \phi(\mathbf{x}): \quad \mathcal{X}_0 = \{\mathbf{x} : \phi(\mathbf{x}) = 0\}, \quad \mathcal{X}_1 = \{\mathbf{x} : \phi(\mathbf{x}) = 1\}.$$

Let us define probabilities of false alarm and miss:

$$P_{\text{FA}} = \mathbb{E}_{\mathbf{x} | \theta}[\phi(\mathbf{X}) | \theta] = \int_{\mathcal{X}_1} p(\mathbf{x} | \theta) d\mathbf{x} \quad \text{for } \theta \text{ in } \Theta_0$$

$$\begin{aligned} P_{\text{M}} &= \mathbb{E}_{\mathbf{x} | \theta}[1 - \phi(\mathbf{X}) | \theta] = 1 - \int_{\mathcal{X}_1} p(\mathbf{x} | \theta) d\mathbf{x} \\ &= \int_{\mathcal{X}_0} p(\mathbf{x} | \theta) d\mathbf{x} \quad \text{for } \theta \text{ in } \Theta_1. \end{aligned}$$

Then, the probability of detection (correctly deciding \mathcal{H}_1) is

$$P_{\text{D}} = 1 - P_{\text{M}} = \mathbb{E}_{\mathbf{x} | \theta}[\phi(\mathbf{X}) | \theta] = \int_{\mathcal{X}_1} p(\mathbf{x} | \theta) d\mathbf{x} \quad \text{for } \theta \text{ in } \Theta_1.$$

Note: P_{FA} and $P_{\text{D}}/P_{\text{M}}$ are generally functions of the parameter θ (where $\theta \in \Theta_0$ when computing P_{FA} and $\theta \in \Theta_1$ when computing $P_{\text{D}}/P_{\text{M}}$).

More Terminology. Statisticians use the following terminology:

- False alarm \equiv “Type I error”
- Miss \equiv “Type II error”
- Probability of detection \equiv “Power”
- Probability of false alarm \equiv “Significance level.”

Bayesian Decision-theoretic Detection Theory

Recall (a slightly generalized version of) the *posterior expected loss*:

$$\rho(\text{action} | \mathbf{x}) = \int_{\Theta} L(\theta, \text{action}) p(\theta | \mathbf{x}) d\theta$$

that we introduced in handout # 4 when we discussed Bayesian decision theory. Let us now apply this theory to our easy example discussed here: *hypothesis testing*, where our action space consists of only two choices. We first assign a loss table:

decision rule $\phi \downarrow$	true state \rightarrow	Θ_1	Θ_0
$\mathbf{x} \in \mathcal{X}_1$		$L(1 1) = 0$	$L(1 0)$
$\mathbf{x} \in \mathcal{X}_0$		$L(0 1)$	$L(0 0) = 0$

with the loss function described by the quantities $L(\text{declared} | \text{true})$:

- $L(1 | 0)$ quantifies loss due to a false alarm,
- $L(0 | 1)$ quantifies loss due to a miss,
- $L(1 | 1)$ and $L(0 | 0)$ (losses due to correct decisions) — typically set to zero in real life. Here, we adopt zero losses for correct decisions.

Now, our posterior expected loss takes two values:

$$\begin{aligned}
 \rho_0(\mathbf{x}) &= \int_{\Theta_1} L(0 | 1) p(\theta | \mathbf{x}) d\theta \\
 &\quad + \int_{\Theta_0} \underbrace{L(0 | 0)}_0 p(\theta | \mathbf{x}) d\theta \\
 &= \int_{\Theta_1} L(0 | 1) p(\theta | \mathbf{x}) d\theta \\
 &\stackrel{L(0 | 1) \text{ is constant}}{=} L(0 | 1) \underbrace{\int_{\Theta_1} p(\theta | \mathbf{x}) d\theta}_{P[\theta \in \Theta_1 | \mathbf{x}]}
 \end{aligned}$$

and, similarly,

$$\begin{aligned}
 \rho_1(\mathbf{x}) &= \int_{\Theta_0} L(1 | 0) p(\theta | \mathbf{x}) d\theta \\
 &\stackrel{L(1 | 0) \text{ is constant}}{=} L(1 | 0) \underbrace{\int_{\Theta_0} p(\theta | \mathbf{x}) d\theta}_{P[\theta \in \Theta_0 | \mathbf{x}]} .
 \end{aligned}$$

We define the *Bayes' decision rule* as the rule that *minimizes the posterior expected loss*; this rule corresponds to choosing the data-space partitioning as follows:

$$\mathcal{X}_1 = \{\mathbf{x} : \rho_1(\mathbf{x}) \leq \rho_0(\mathbf{x})\}$$

or

$$\mathcal{X}_1 = \left\{ \mathbf{x} : \frac{\overbrace{\int_{\Theta_1} p(\theta | \mathbf{x}) d\theta}^{P[\theta \in \Theta_1 | \mathbf{x}]}}{\underbrace{\int_{\Theta_0} p(\theta | \mathbf{x}) d\theta}_{P[\theta \in \Theta_0 | \mathbf{x}]}} \geq \frac{L(1|0)}{L(0|1)} \right\} \quad (1)$$

or, equivalently, upon applying the Bayes' rule:

$$\mathcal{X}_1 = \left\{ \mathbf{x} : \frac{\int_{\Theta_1} p(\mathbf{x} | \theta) \pi(\theta) d\theta}{\int_{\Theta_0} p(\mathbf{x} | \theta) \pi(\theta) d\theta} \geq \frac{L(1|0)}{L(0|1)} \right\}. \quad (2)$$

0-1 loss: For $L(1|0) = L(0|1) = 1$, we have

decision rule $\phi \downarrow$	true state \rightarrow	Θ_1	Θ_0
$\mathbf{x} \in \mathcal{X}_1$		$L(1 1) = 0$	$L(1 0) = 1$
$\mathbf{x} \in \mathcal{X}_0$		$L(0 1) = 1$	$L(0 0) = 0$

yielding the following Bayes' decision rule, called the maximum *a posteriori* (MAP) rule:

$$\mathcal{X}_1 = \left\{ \mathbf{x} : \frac{P[\theta \in \Theta_1 | \mathbf{x}]}{P[\theta \in \Theta_0 | \mathbf{x}]} \geq 1 \right\} \quad (3)$$

or, equivalently, upon applying the Bayes' rule:

$$\mathcal{X}_1 = \left\{ \mathbf{x} : \frac{\int_{\Theta_1} p(\mathbf{x} | \theta) \pi(\theta) d\theta}{\int_{\Theta_0} p(\mathbf{x} | \theta) \pi(\theta) d\theta} \geq 1 \right\}. \quad (4)$$

Simple hypotheses. Let us specialize (1) to the case of simple hypotheses ($\Theta_0 = \{\theta_0\}, \Theta_1 = \{\theta_1\}$):

$$\mathcal{X}_1 = \left\{ \mathbf{x} : \underbrace{\frac{p(\theta_1 | \mathbf{x})}{p(\theta_0 | \mathbf{x})}}_{\text{posterior-odds ratio}} \geq \frac{L(1 | 0)}{L(0 | 1)} \right\}. \quad (5)$$

We can rewrite (5) using the Bayes' rule:

$$\mathcal{X}_1 = \left\{ \mathbf{x} : \underbrace{\frac{p(\mathbf{x} | \theta_1)}{p(\mathbf{x} | \theta_0)}}_{\text{likelihood ratio}} \geq \frac{\pi_0 L(1 | 0)}{\pi_1 L(0 | 1)} \right\} \quad (6)$$

where

$$\pi_0 = \pi(\theta_0), \quad \pi_1 = \pi(\theta_1) = 1 - \pi_0$$

describe the prior probability mass function (pmf) of the binary random variable θ (recall that $\theta \in \{\theta_0, \theta_1\}$). Hence, for binary simple hypotheses, the prior pmf of θ is the Bernoulli pmf.

Preposterior (Bayes) Risk

The preposterior (Bayes) risk for rule $\phi(\mathbf{x})$ is

$$\begin{aligned} E_{x,\theta}[\text{loss}] &= \int_{\mathcal{X}_1} \int_{\Theta_0} L(1|0) p(\mathbf{x}|\theta) \pi(\theta) d\theta d\mathbf{x} \\ &\quad + \int_{\mathcal{X}_0} \int_{\Theta_1} L(0|1) p(\mathbf{x}|\theta) \pi(\theta) d\theta d\mathbf{x}. \end{aligned}$$

How do we choose the rule $\phi(\mathbf{x})$ that minimizes the preposterior

risk?

$$\begin{aligned}
 & \int_{\mathcal{X}_1} \int_{\Theta_0} L(1|0) p(\mathbf{x}|\theta) \pi(\theta) d\theta d\mathbf{x} \\
 & \quad + \int_{\mathcal{X}_0} \int_{\Theta_1} L(0|1) p(\mathbf{x}|\theta) \pi(\theta) d\theta d\mathbf{x} \\
 = & \int_{\mathcal{X}_1} \int_{\Theta_0} L(1|0) p(\mathbf{x}|\theta) \pi(\theta) d\theta d\mathbf{x} \\
 & \quad - \int_{\mathcal{X}_1} \int_{\Theta_1} L(0|1) p(\mathbf{x}|\theta) \pi(\theta) d\theta d\mathbf{x} \\
 & \quad + \int_{\mathcal{X}_0} \int_{\Theta_1} L(0|1) p(\mathbf{x}|\theta) \pi(\theta) d\theta d\mathbf{x} \\
 & \quad + \int_{\mathcal{X}_1} \int_{\Theta_1} L(0|1) p(\mathbf{x}|\theta) \pi(\theta) d\theta d\mathbf{x} \\
 = & \quad \underbrace{\text{const}} \\
 & \quad \text{not dependent on } \phi(\mathbf{x}) \\
 & \quad + \int_{\mathcal{X}_1} \left\{ L(1|0) \cdot \int_{\Theta_0} p(\mathbf{x}|\theta) \pi(\theta) d\theta \right. \\
 & \quad \left. - L(0|1) \cdot \int_{\Theta_1} p(\mathbf{x}|\theta) \pi(\theta) d\theta \right\} d\mathbf{x}
 \end{aligned}$$

implying that \mathcal{X}_1 should be chosen as

$$\left\{ \mathcal{X}_1 : L(1|0) \cdot \int_{\Theta_0} p(\mathbf{x}|\theta) \pi(\theta) d\theta - L(0|1) \cdot \int_{\Theta_1} p(\mathbf{x}|\theta) \pi(\theta) d\theta < 0 \right\}$$

which, as expected, is the same as (2), since

minimizing the posterior expected loss

\iff minimizing the preposterior risk for every \mathbf{x}

as showed earlier in handout # 4.

0-1 loss: For the 0-1 loss, i.e. $L(1|0) = L(0|1) = 1$, the preposterior (Bayes) risk for rule $\phi(\mathbf{x})$ is

$$\begin{aligned} E_{\mathbf{x},\theta}[\text{loss}] &= \int_{\mathcal{X}_1} \int_{\Theta_0} p(\mathbf{x} | \theta) \pi(\theta) d\theta d\mathbf{x} \\ &\quad + \int_{\mathcal{X}_0} \int_{\Theta_1} p(\mathbf{x} | \theta) \pi(\theta) d\theta d\mathbf{x} \end{aligned} \quad (7)$$

which is simply the *average error probability*, with averaging performed over the joint probability density or mass function (pdf/pmf) or the data \mathbf{x} and parameters θ .

Bayesian Decision-theoretic Detection for Simple Hypotheses

The Bayes' decision rule for simple hypotheses is (6):

$$\underbrace{\Lambda(\mathbf{x})}_{\text{likelihood ratio}} = \frac{p(\mathbf{x} | \theta_1)}{p(\mathbf{x} | \theta_0)} \stackrel{\mathcal{H}_1}{\gtrless} \frac{\pi_0 L(1|0)}{\pi_1 L(0|1)} \equiv \tau \quad (8)$$

see also Ch. 3.7 in Kay-II. (Recall that $\Lambda(x)$ is the sufficient statistic for the detection problem, see p. 37 in handout # 1.) Equivalently,

$$\log \Lambda(\mathbf{x}) = \log[p(\mathbf{x} | \theta_1)] - \log[p(\mathbf{x} | \theta_0)] \stackrel{\mathcal{H}_1}{\gtrless} \log \tau \equiv \tau'.$$

Minimum Average Error Probability Detection: In the familiar 0-1 loss case where $L(1|0) = L(0|1) = 1$, we know that the preposterior (Bayes) risk is equal to the *average error probability*, see (7). This average error probability greatly

simplifies in the simple hypothesis testing case:

$$\begin{aligned}
 \text{av. error probability} &= \int_{\mathcal{X}_1} \underbrace{L(1|0)}_1 p(\mathbf{x} | \theta_0) \pi_0 d\mathbf{x} \\
 &+ \int_{\mathcal{X}_0} \underbrace{L(0|1)}_1 p(\mathbf{x} | \theta_1) \pi_1 d\mathbf{x} \\
 &= \pi_0 \cdot \underbrace{\int_{\mathcal{X}_1} p(\mathbf{x} | \theta_0) d\mathbf{x}}_{P_{FA}} + \pi_1 \cdot \underbrace{\int_{\mathcal{X}_0} p(\mathbf{x} | \theta_1) d\mathbf{x}}_{P_M}
 \end{aligned}$$

where, as before, the averaging is performed over the joint pdf/pmf of the data \mathbf{x} and parameters θ , and

$$\pi_0 = \pi(\theta_0), \quad \pi_1 = \pi(\theta_1) = 1 - \pi_0 \quad (\text{the Bernoulli pmf}).$$

In this case, our Bayes' decision rule simplifies to the MAP rule (as expected, see (5) and Ch. 3.6 in Kay-II):

$$\frac{p(\theta_1 | \mathbf{x})}{\underbrace{p(\theta_0 | \mathbf{x})}_1} \stackrel{\mathcal{H}_1}{\geq} 1 \quad (9)$$

posterior-odds ratio

or, equivalently, upon applying the Bayes' rule:

$$\frac{p(\mathbf{x} | \theta_1)}{\underbrace{p(\mathbf{x} | \theta_0)}_1} \stackrel{\mathcal{H}_1}{\geq} \frac{\pi_0}{\pi_1}. \quad (10)$$

likelihood ratio

which is the same as

- (4), upon substituting the Bernoulli pmf as the prior pmf for θ and
- (8), upon substituting $L(1|0) = L(0|1) = 1$.

Bayesian Decision-theoretic Detection Theory: Handling Nuisance Parameters

We apply the same approach as before — integrate the nuisance parameters (φ , say) out!

Therefore, (1) still holds for testing

$$\begin{aligned}\mathcal{H}_0 &: \theta \in \Theta_0 \quad \text{versus} \\ \mathcal{H}_1 &: \theta \in \Theta_1\end{aligned}$$

but $p_{\theta|\mathbf{x}}(\theta|\mathbf{x})$ is computed as follows:

$$p_{\theta|\mathbf{x}}(\theta|\mathbf{x}) = \int p_{\theta,\varphi|\mathbf{x}}(\theta, \varphi|\mathbf{x}) d\varphi$$

and, therefore,

$$\frac{\overbrace{\int_{\Theta_1} \int p_{\theta,\varphi|\mathbf{x}}(\theta, \varphi|\mathbf{x}) d\varphi d\theta}^{p(\theta|\mathbf{x})}}{\underbrace{\int_{\Theta_0} \int p_{\theta,\varphi|\mathbf{x}}(\theta, \varphi|\mathbf{x}) d\varphi d\theta}_{p(\theta|\mathbf{x})}} \stackrel{\mathcal{H}_1}{\geq} \frac{L(1|0)}{L(0|1)} \quad (11)$$

or, equivalently, upon applying the Bayes' rule:

$$\frac{\int_{\Theta_1} \int p_{\mathbf{x} | \theta, \varphi}(\mathbf{x} | \theta, \varphi) \pi_{\theta, \varphi}(\theta, \varphi) d\varphi d\theta}{\int_{\Theta_0} \int p_{\mathbf{x} | \theta, \varphi}(\mathbf{x} | \theta, \varphi) \pi_{\theta, \varphi}(\theta, \varphi) d\varphi d\theta} \stackrel{\mathcal{H}_1}{\geq} \frac{L(1 | 0)}{L(0 | 1)}. \quad (12)$$

Now, if θ and φ are independent *a priori*, i.e.

$$\pi_{\theta, \varphi}(\theta, \varphi) = \pi_{\theta}(\theta) \cdot \pi_{\varphi}(\varphi) \quad (13)$$

then (12) can be rewritten as

$$\frac{\int_{\Theta_1} \pi_{\theta}(\theta) \overbrace{\int p_{\mathbf{x} | \theta, \varphi}(\mathbf{x} | \theta, \varphi) \pi_{\varphi}(\varphi) d\varphi}^{p(\mathbf{x} | \theta)} d\theta}{\int_{\Theta_0} \pi_{\theta}(\theta) \underbrace{\int p_{\mathbf{x} | \theta, \varphi}(\mathbf{x} | \theta, \varphi) \pi_{\varphi}(\varphi) d\varphi}_{p(\mathbf{x} | \theta)} d\theta} \stackrel{\mathcal{H}_1}{\geq} \frac{L(1 | 0)}{L(0 | 1)}. \quad (14)$$

Simple hypotheses and independent priors for θ and φ : Let us specialize (11) to the simple hypotheses ($\Theta_0 = \{\theta_0\}$, $\Theta_1 = \{\theta_1\}$):

$$\frac{\int p_{\theta, \varphi | \mathbf{x}}(\theta_1, \varphi | \mathbf{x}) d\varphi}{\int p_{\theta, \varphi | \mathbf{x}}(\theta_0, \varphi | \mathbf{x}) d\varphi} \stackrel{\mathcal{H}_1}{\geq} \frac{L(1 | 0)}{L(0 | 1)}. \quad (15)$$

Now, if θ and φ are independent *a priori*, i.e. (13) holds, then

we can rewrite (14) [or (15)] using the Bayes' rule]:

$$\underbrace{\frac{\int p_{\mathbf{x}|\theta,\varphi}(\mathbf{x}|\theta_1,\varphi)\pi_\varphi(\varphi)d\varphi}{\int p_{\mathbf{x}|\theta,\varphi}(\mathbf{x}|\theta_0,\varphi)\pi_\varphi(\varphi)d\varphi}}_{\text{integrated likelihood ratio}} = \overbrace{\frac{p(\mathbf{x}|\theta_1)}{p(\mathbf{x}|\theta_0)} \stackrel{\mathcal{H}_1}{\gtrless} \frac{\pi_0 L(1|0)}{\pi_1 L(0|1)}}^{\text{same as (6)}} \quad (16)$$

where

$$\pi_0 = \pi_\theta(\theta_0), \quad \pi_1 = \pi_\theta(\theta_1) = 1 - \pi_0.$$

Chernoff Bound on Average Error Probability for Simple Hypotheses

Recall that minimizing the average error probability

$$\begin{aligned} \text{av. error probability} &= \int_{\mathcal{X}_1} \int_{\Theta_0} p(\mathbf{x} | \theta) \pi(\theta) d\theta d\mathbf{x} \\ &+ \int_{\mathcal{X}_0} \int_{\Theta_1} p(\mathbf{x} | \theta) \pi(\theta) d\theta d\mathbf{x} \end{aligned}$$

leads to the MAP decision rule:

$$\mathcal{X}_1^* = \left\{ \mathbf{x} : \int_{\Theta_0} p(\mathbf{x} | \theta) \pi(\theta) d\theta - \int_{\Theta_1} p(\mathbf{x} | \theta) \pi(\theta) d\theta < 0 \right\}.$$

In many applications, we *may not be able to obtain* a simple closed-form expression for the minimum average error

probability, but we *can bound it* as follows:

$$\begin{aligned} \text{min av. error probability} &= \int_{\mathcal{X}_1^*} \int_{\Theta_0} p(\mathbf{x} | \theta) \pi(\theta) d\theta d\mathbf{x} \\ &+ \int_{\mathcal{X}_0^*} \int_{\Theta_1} p(\mathbf{x} | \theta) \pi(\theta) d\theta d\mathbf{x} \end{aligned}$$

using the def. of \mathcal{X}_1^*

$$= \int_{\underbrace{\mathcal{X}}_{\text{data space}}} \min \left\{ \int_{\Theta_0} p(\mathbf{x} | \theta) \pi(\theta) d\theta, \int_{\Theta_1} p(\mathbf{x} | \theta) \pi(\theta) d\theta \right\} d\mathbf{x}$$

$$\begin{aligned} &\leq \int_{\mathcal{X}} \left[\overbrace{\int_{\Theta_0} p(\mathbf{x} | \theta) \pi(\theta) d\theta}^{\triangleq q_0(\mathbf{x})} \right]^\lambda \left[\overbrace{\int_{\Theta_1} p(\mathbf{x} | \theta) \pi(\theta) d\theta}^{\triangleq q_1(\mathbf{x})} \right]^{1-\lambda} d\mathbf{x} \\ &= \int_{\mathcal{X}} [q_0(\mathbf{x})]^\lambda [q_1(\mathbf{x})]^{1-\lambda} d\mathbf{x} \end{aligned}$$

which is the *Chernoff bound on the minimum average error probability*. Here, we have used the fact that

$$\min\{a, b\} \leq a^\lambda b^{1-\lambda}, \quad \text{for } 0 \leq \lambda \leq 1, \quad a, b \geq 0.$$

When

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix}$$

with

- x_1, x_2, \dots, x_N **conditionally** independent, **identically** distributed (i.i.d.) **given θ** and
- simple hypotheses (i.e. $\Theta_0 = \{\theta_0\}, \Theta_1 = \{\theta_1\}$)

then

$$q_0(\mathbf{x}) = p(\mathbf{x} | \theta_0) \cdot \pi_0 = \pi_0 \prod_{n=1}^N p(x_n | \theta_0)$$

$$q_1(\mathbf{x}) = p(\mathbf{x} | \theta_1) \cdot \pi_1 = \pi_1 \prod_{n=1}^N p(x_n | \theta_1)$$

yielding

Chernoff bound for N conditionally i.i.d. measurements (given θ) and simple hyp.

$$\begin{aligned} &= \int \left[\pi_0 \prod_{n=1}^N p(x_n | \theta_0) \right]^\lambda \left[\pi_1 \prod_{n=1}^N p(x_n | \theta_1) \right]^{1-\lambda} d\mathbf{x} \\ &= \pi_0^\lambda \pi_1^{1-\lambda} \cdot \prod_{n=1}^N \left\{ \int [p(x_n | \theta_0)]^\lambda [p(x_n | \theta_1)]^{1-\lambda} dx_n \right\} \\ &= \pi_0^\lambda \pi_1^{1-\lambda} \cdot \left\{ \int [p(x_1 | \theta_0)]^\lambda [p(x_1 | \theta_1)]^{1-\lambda} dx_1 \right\}^N \end{aligned}$$

or, in other words,

$$\frac{1}{N} \cdot \log(\text{min av. error probability}) \leq \log(\pi_0^\lambda \pi_1^{1-\lambda}) + \log \int [p(x_1 | \theta_0)]^\lambda [p(x_1 | \theta_1)]^{1-\lambda} dx_1, \quad \forall \lambda \in [0, 1].$$

If $\pi_0 = \pi_1 = 1/2$ (which is almost always the case of interest when evaluating average error probabilities), we can say that, as $N \rightarrow \infty$,

min av. error probability

→

for N cond. i.i.d. measurements (given θ) and simple hypotheses

$$f(N) \cdot \exp \left(- N \cdot \underbrace{\left\{ - \min_{\lambda \in [0,1]} \log \int [p(x_1 | \theta_0)]^\lambda [p(x_1 | \theta_1)]^{1-\lambda} \right\}}_{\text{Chernoff information for a single observation}} \right)$$

where $f(N)$ is a slowly-varying function compared with the exponential term:

$$\lim_{N \rightarrow \infty} \frac{\log f(N)}{N} = 0.$$

Note that the Chernoff information in the exponent term of the above expression *quantifies* the asymptotic behavior of the minimum average error probability.

We now give a useful result, taken from

K. Fukunaga, *Introduction to Statistical Pattern Recognition*, 2nd ed., San Diego, CA: Academic Press, 1990

for evaluating a class of Chernoff bounds.

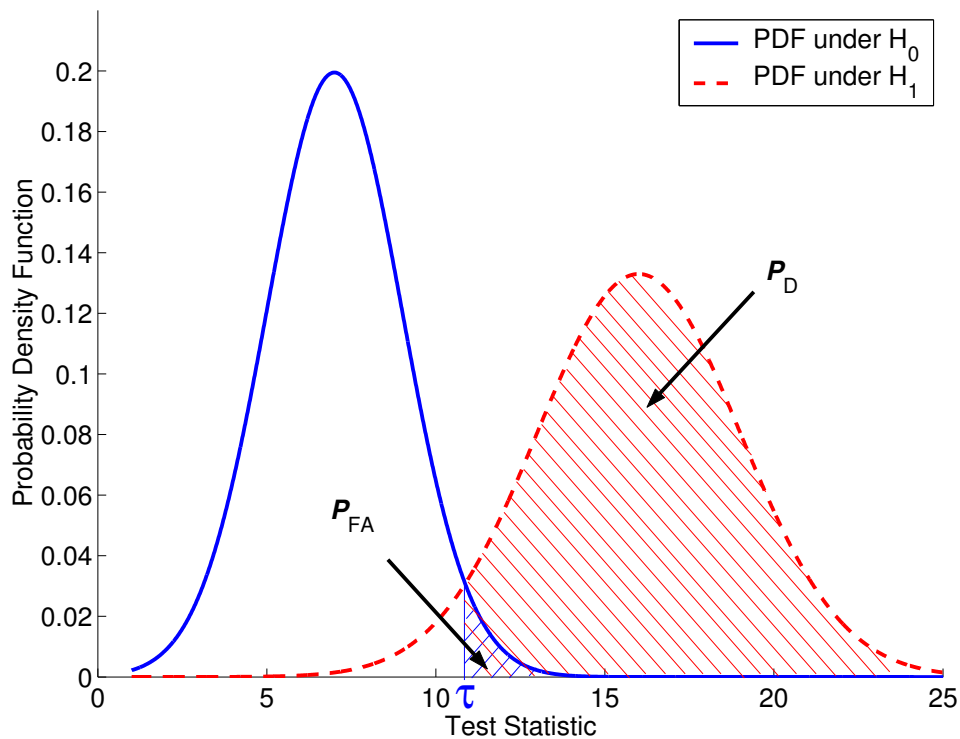
Lemma 1. Consider $p_1(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ and $p_2(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$. Then

$$\int [p_1(\mathbf{x})]^\lambda \cdot [p_2(\mathbf{x})]^{1-\lambda} d\mathbf{x} = \exp[-g(\lambda)]$$

where

$$g(\lambda) = \frac{\lambda(1-\lambda)}{2} \cdot (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1)^T [\lambda \boldsymbol{\Sigma}_1 + (1-\lambda) \boldsymbol{\Sigma}_2]^{-1} (\boldsymbol{\mu}_2 - \boldsymbol{\mu}_1) \\ + \frac{1}{2} \log \left[\frac{|\lambda \boldsymbol{\Sigma}_1 + (1-\lambda) \boldsymbol{\Sigma}_2|}{|\boldsymbol{\Sigma}_1|^\lambda \cdot |\boldsymbol{\Sigma}_2|^{1-\lambda}} \right].$$

Probabilities of False Alarm (P_{FA}) and Detection (P_D) for Simple Hypotheses



$$P_{FA} = P[\underbrace{\text{test statistic}}_{\mathbf{X} \in \mathcal{X}_1} > \tau \mid \theta = \theta_0]$$
$$P_D = P[\underbrace{\text{test statistic}}_{\mathbf{X} \in \mathcal{X}_0} > \tau \mid \theta = \theta_1].$$

Comments:

- (i) As the region \mathcal{X}_1 shrinks (i.e. $\tau \nearrow \infty$), both of the above probabilities shrink towards zero.

- (ii) As the region \mathcal{X}_1 grows (i.e. $\tau \searrow 0$), both of these probabilities grow towards unity.
- (iii) Observations (i) and (ii) *do not imply equality* between P_{FA} and P_{D} ; in most cases, as \mathcal{R}_1 grows, P_{D} grows more rapidly than P_{FA} (i.e. we better be right more often than we are wrong).
- (iv) However, the *perfect case* where our rule is always right and never wrong ($P_{\text{D}} = 1$ and $P_{\text{FA}} = 0$) cannot occur when the conditional pdfs/pmfs $p(\mathbf{x} | \theta_0)$ and $p(\mathbf{x} | \theta_1)$ overlap.
- (v) Thus, to increase the detection probability P_{D} , we must also allow for the false-alarm probability P_{FA} to increase. This behavior
- represents the fundamental tradeoff in hypothesis testing and detection theory and
 - motivates us to introduce a (classical) approach to testing simple hypotheses, pioneered by Neyman and Pearson (to be discussed next).

Neyman-Pearson Test for Simple Hypotheses

Setup:

- Parametric data models $p(\mathbf{x}; \theta_0)$, $p(\mathbf{x}; \theta_1)$,
- Simple hypothesis testing:

$$\mathcal{H}_0 : \theta = \theta_0 \quad \text{versus}$$

$$\mathcal{H}_1 : \theta = \theta_1.$$

- No prior pdf/pmf on θ is available.

Goal: Design a test that maximizes the probability of detection

$$P_D = P[\mathbf{X} \in \mathcal{X}_1; \theta = \theta_0]$$

(equivalently, minimizes the miss probability P_M) under the constraint

$$P_{\text{FA}} = P[\mathbf{X} \in \mathcal{X}_1; \theta = \theta_0] = \alpha' \leq \alpha.$$

Here, we consider simple hypotheses; classical version of testing composite hypotheses is much more complicated. The Bayesian version of testing composite hypotheses is trivial (as

is everything else Bayesian, at least conceptually) and we have already seen it.

Solution. We apply the Lagrange-multiplier approach: maximize

$$\begin{aligned} L &= P_D + \lambda \cdot (P_{FA} - \alpha') \\ &= \int_{\mathcal{X}_1} p(\mathbf{x}; \theta_1) d\mathbf{x} + \lambda \cdot \left[\int_{\mathcal{X}_1} p(\mathbf{x}; \theta_0) d\mathbf{x} - \alpha' \right] \\ &= \int_{\mathcal{X}_1} [p(\mathbf{x}; \theta_1) - \lambda p(\mathbf{x}; \theta_0)] d\mathbf{x} - \lambda \cdot \alpha'. \end{aligned}$$

To maximize L , set

$$\mathcal{X}_1 = \{ \mathbf{x} : p(\mathbf{x}; \theta_1) - \lambda \cdot p(\mathbf{x}; \theta_0) > 0 \} = \left\{ \mathbf{x} : \frac{p(\mathbf{x}; \theta_1)}{p(\mathbf{x}; \theta_0)} > \lambda \right\}.$$

Again, we find the likelihood ratio:

$$\Lambda(\mathbf{x}) = \frac{p(\mathbf{x}; \theta_1)}{p(\mathbf{x}; \theta_0)}.$$

Recall our constraint:

$$\int_{\mathcal{X}_1} p(\mathbf{x}; \theta_0) d\mathbf{x} = P_{FA} = \alpha' \leq \alpha.$$

If we increase λ , P_{FA} and P_{D} go down. Similarly, if we decrease λ , P_{FA} and P_{D} go up. Hence, to maximize P_{D} , choose λ so that P_{FA} is as big as possible under the constraint.

Two useful ways for determining the threshold that achieves a specified false-alarm rate:

- Find λ that satisfies

$$\int_{\mathbf{x} : \Lambda(\mathbf{x}) > \lambda} p(\mathbf{x}; \theta_0) d\mathbf{x} = P_{\text{FA}} = \alpha$$

or,

- expressing in terms of the pdf/pmf of $\Lambda(\mathbf{x})$ under \mathcal{H}_0 :

$$\int_{\lambda}^{\infty} p_{\Lambda; \theta_0}(l; \theta_0) dl = \alpha.$$

or, perhaps, in terms of a monotonic function of $\Lambda(\mathbf{x})$, say $T(\mathbf{x}) = \text{monotonic function}(\Lambda(\mathbf{x}))$.

Warning: We have been implicitly assuming that P_{FA} is a continuous function of λ . Some (not insightful) technical adjustments are needed if this is not the case.

A way of handling nuisance parameters: We can utilize the integrated (marginal) likelihood ratio (16) under the Neyman-Pearson setup as well.

Chernoff-Stein Lemma for Bounding the Miss Probability in Neyman-Pearson Tests of Simple Hypotheses

Recall the definition of the *Kullback-Leibler (K-L) distance* $D(\mathbf{p} \parallel \mathbf{q})$ from one pmf (\mathbf{p}) to another (\mathbf{q}):

$$D(\mathbf{p} \parallel \mathbf{q}) = \sum_k p_k \log \frac{p_k}{q_k}.$$

The complete proof of this lemma for the discrete (pmf) case is given in

Additional Reading: T.M. Cover and J.A. Thomas, *Elements of Information Theory*. Second ed., New York: Wiley, 2006.

Setup for the Chernoff-Stein Lemma

- Assume that x_1, x_2, \dots, x_N are conditionally i.i.d. given θ .
- We adopt the Neyman-Pearson framework, i.e. obtain a decision threshold to achieve a fixed P_{FA} . Let us study the asymptotic $P_{\text{M}} = 1 - P_{\text{D}}$ as the number of observations N gets large.
- To keep P_{FA} constant as N increases, we need to make our decision threshold (γ , say) vary with N , i.e.

$$\gamma = \gamma_N(P_{\text{FA}})$$

Now, the miss probability is

$$P_{\text{M}} = P_{\text{M}}(\gamma) = P_{\text{M}}\left(\gamma_N(P_{\text{FA}})\right).$$

Chernoff-Stein Lemma

The Chernoff-Stein lemma says:

$$\lim_{P_{\text{FA}} \rightarrow 0} \lim_{N \rightarrow \infty} \frac{1}{N} \log P_{\text{M}} = - \underbrace{D\left(p(X_n | \theta_0) \parallel p(X_n | \theta_1)\right)}_{\text{K-L distance for a single observation}}$$

where the K-L distance between $p(x_n | \theta_0)$ and $p(x_n | \theta_1)$

$$D\left(p(X_n | \theta_0) \parallel p(X_n | \theta_1)\right) = -\mathbb{E}_{p(x_n | \theta_0)} \left[\log \frac{p(X_n | \theta_1)}{p(X_n | \theta_0)} \right]$$

discrete (pmf) case

$$= - \sum_{x_n} p(x_n | \theta_0) \log \left[\frac{p(x_n | \theta_1)}{p(x_n | \theta_0)} \right]$$

does not depend on the observation index n , since x_n are conditionally i.i.d. given θ .

Equivalently, we can state that

$$P_{\text{M}} \xrightarrow{N \rightarrow +\infty} f(N) \cdot \exp \left[-N \cdot D\left(p(X_n | \theta_0) \parallel p(X_n | \theta_1)\right) \right]$$

as $P_{\text{FA}} \rightarrow 0$ and $N \rightarrow \infty$, where $f(N)$ is a slowly-varying function compared with the exponential term (when $P_{\text{FA}} \rightarrow 0$ and $N \rightarrow \infty$).

Detection for Simple Hypotheses: Example

Known positive DC level in additive white Gaussian noise (AWGN), Example 3.2 in Kay-II.

Consider

$$\mathcal{H}_0 : \quad x[n] = w[n], \quad n = 1, 2, \dots, N \quad \text{versus}$$

$$\mathcal{H}_1 : \quad x[n] = A + w[n], \quad n = 1, 2, \dots, N$$

where

- $A > 0$ is a known constant,
- $w[n]$ is zero-mean white Gaussian noise with known variance σ^2 , i.e.

$$w[n] \sim \mathcal{N}(0, \sigma^2).$$

The above hypothesis-testing formulation is EE-like: noise versus signal plus noise. **It is similar to the on-off keying scheme in communications**, which gives us an idea to rephrase it so that it fits our formulation on p. 4 (for which we have developed all the theory so far). Here is such an

alternative formulation: consider a family of pdfs

$$p(\mathbf{x}; a) = \frac{1}{\sqrt{(2\pi\sigma^2)^N}} \cdot \exp \left[-\frac{1}{2\sigma^2} \sum_{n=1}^N (x[n] - a)^2 \right] \quad (17)$$

and the following (equivalent) hypotheses:

$$\begin{aligned} \mathcal{H}_0 : & \quad a = 0 \quad (\text{off}) \quad \text{versus} \\ \mathcal{H}_1 : & \quad a = A \quad (\text{on}). \end{aligned}$$

Then, the likelihood ratio is

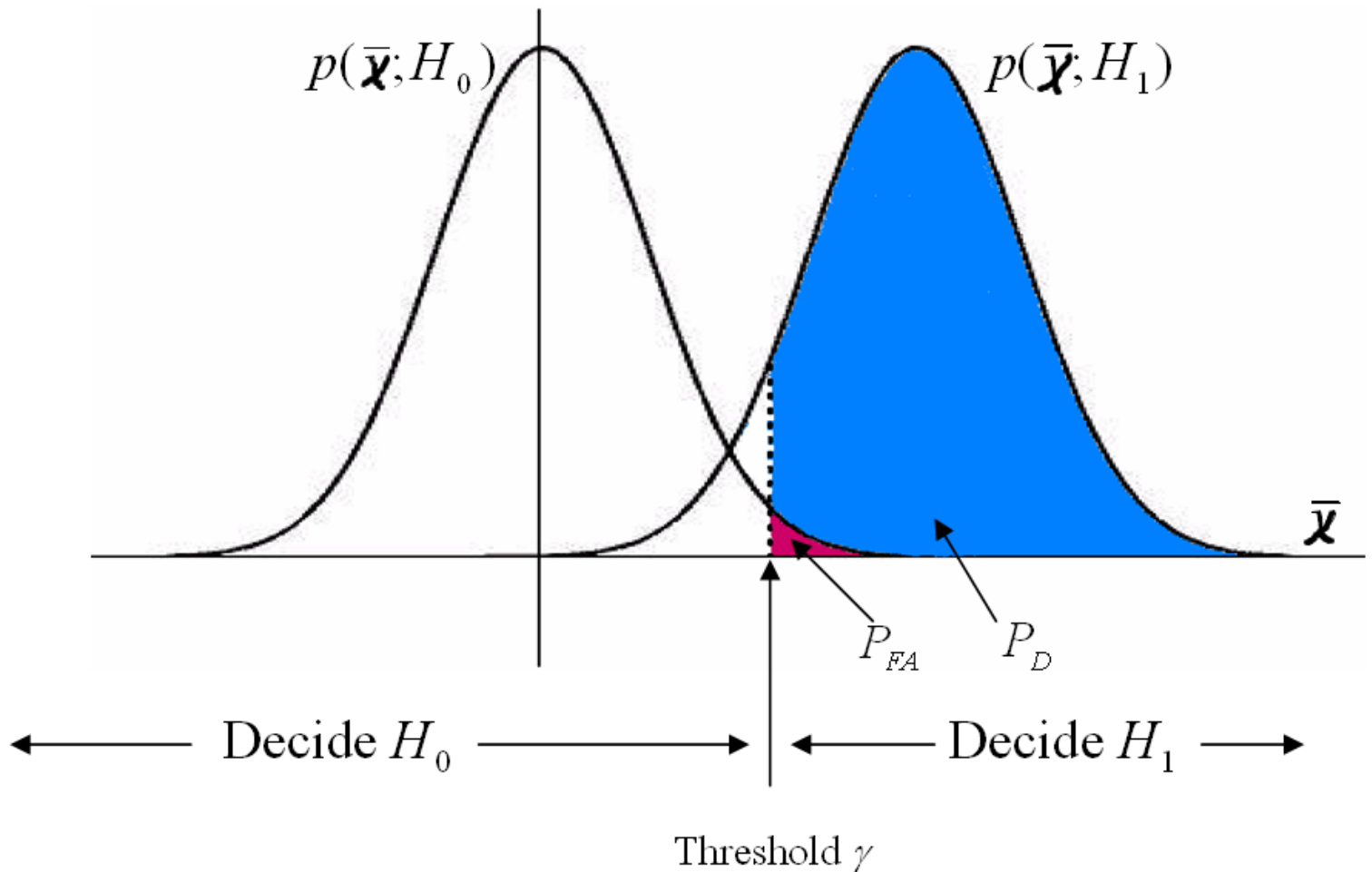
$$\begin{aligned} \Lambda(\mathbf{x}) &= \frac{p(\mathbf{x}; a = A)}{p(\mathbf{x}; a = 0)} \\ &= \frac{1/(2\pi\sigma^2)^{N/2} \cdot \exp[-\frac{1}{2\sigma^2} \sum_{n=1}^N (x[n] - A)^2]}{1/(2\pi\sigma^2)^{N/2} \cdot \exp(-\frac{1}{2\sigma^2} \sum_{n=1}^N x[n]^2)}. \end{aligned}$$

Now, take the logarithm and, after simple manipulations, reduce our likelihood-ratio test to comparing

$$T(\mathbf{x}) = \bar{x} \triangleq \frac{1}{N} \sum_{n=1}^N x[n]$$

with a threshold γ . [Here, $T(\mathbf{x})$ is a monotonic function of $\Lambda(x)$.] If $T(\mathbf{x}) > \gamma$ accept \mathcal{H}_1 (i.e. reject \mathcal{H}_0), otherwise accept

\mathcal{H}_0 (well, not exactly, we will talk more about this decision on p. 59).



The choice of γ depends on the approach that we take. For the Bayes' decision rule, γ is a function of π_0 and π_1 . For the Neyman-Pearson test, γ is chosen to achieve (control) a desired P_{FA} .

Bayesian decision-theoretic detection for 0-1 loss (corresponding to minimizing the average error

probability):

$$\log \Lambda(\mathbf{x}) = -\frac{1}{2\sigma^2} \sum_{n=1}^N (x[n] - A)^2 + \frac{1}{2\sigma^2} \sum_{n=1}^N (x[n])^2 \stackrel{\mathcal{H}_1}{\geq} \log \left(\frac{\pi_0}{\pi_1} \right)$$

$$\Leftrightarrow \frac{1}{2\sigma^2} \sum_{n=1}^N \underbrace{(x[n] - x[n] + A)}_0 (x[n] + x[n] - A) \stackrel{\mathcal{H}_1}{\geq} \log \left(\frac{\pi_0}{\pi_1} \right)$$

$$\Leftrightarrow 2A \cdot \left(\sum_{n=1}^N x[n] \right) - A^2 N \stackrel{\mathcal{H}_1}{\geq} 2\sigma^2 \log \left(\frac{\pi_0}{\pi_1} \right)$$

$$\Leftrightarrow \left(\sum_{n=1}^N x[n] \right) - \frac{AN}{2} \stackrel{\mathcal{H}_1}{\geq} \frac{\sigma^2}{A} \log \left(\frac{\pi_0}{\pi_1} \right) \quad \text{since } A > 0$$

and, finally,

$$\bar{x} \stackrel{\mathcal{H}_1}{\geq} \underbrace{\frac{\sigma^2}{NA} \cdot \log(\pi_0/\pi_1)}_{\gamma} + \frac{A}{2}$$

which, for the practically most interesting case of equiprobable hypotheses

$$\pi_0 = \pi_1 = \frac{1}{2} \tag{18}$$

simplifies to

$$\bar{x} \stackrel{\mathcal{H}_1}{\geq} \underbrace{\frac{A}{2}}_{\gamma}$$

known as the *maximum-likelihood test* (i.e. the Bayes' decision rule for 0-1 loss and *a priori* equiprobable hypotheses is defined as the maximum-likelihood test). This maximum-likelihood detector *does not* require the knowledge of the noise variance

σ^2 to declare its decision. However, the knowledge of σ^2 is *key* to assessing the detection performance. Interestingly, these observations will carry over to a few maximum-likelihood tests that we will derive in the future.

Assuming (18), we now derive the minimum average error probability. First, note that $\bar{X} | a = 0 \sim \mathcal{N}(0, \sigma^2/N)$ and $\bar{X} | a = A \sim \mathcal{N}(A, \sigma^2/N)$. Then

$$\begin{aligned}
 \text{min av. error prob.} &= \frac{1}{2} \underbrace{P[\bar{X} > \frac{A}{2} | a = 0]}_{P_{\text{FA}}} + \frac{1}{2} \underbrace{P[\bar{X} < \frac{A}{2} | a = A]}_{P_{\text{M}}} \\
 &= \frac{1}{2} P \left[\underbrace{\frac{\bar{X}}{\sqrt{\sigma^2/N}}}_{\substack{\text{standard} \\ \text{normal} \\ \text{random var.}}} > \frac{A/2}{\sqrt{\sigma^2/N}}; a = 0 \right] \\
 &\quad + \frac{1}{2} P \left[\underbrace{\frac{\bar{X} - A}{\sqrt{\sigma^2/N}}}_{\substack{\text{standard} \\ \text{normal} \\ \text{random var.}}} < \frac{A/2 - A}{\sqrt{\sigma^2/N}}; a = A \right] \\
 &= \frac{1}{2} \underbrace{Q\left(\sqrt{\frac{N A^2}{4 \sigma^2}}\right)}_{\substack{\text{standard} \\ \text{normal complementary cdf}}} + \frac{1}{2} \underbrace{\Phi\left(-\sqrt{\frac{N A^2}{4 \sigma^2}}\right)}_{\substack{\text{standard} \\ \text{normal cdf}}} = Q\left(\frac{1}{2} \sqrt{\frac{N A^2}{\sigma^2}}\right)
 \end{aligned}$$

since

$$\Phi(-x) = Q(x).$$

The minimum probability of error decreases monotonically with $N A^2/\sigma^2$, which is known as the *deflection coefficient*. In this case, the Chernoff information is equal to $\frac{A^2}{8\sigma^2}$ (see Lemma 1):

$$\begin{aligned} & \min_{\lambda \in [0,1]} \log \left\{ \int \underbrace{[\mathcal{N}(0, \sigma^2/N)]^\lambda}_{p_0(x)} \underbrace{[\mathcal{N}(A, \sigma^2/N)]^{1-\lambda}}_{p_1(x)} dx \right\} \\ &= \min_{\lambda \in [0,1]} \left[\frac{\lambda(1-\lambda)}{2} \cdot \frac{A^2}{\sigma^2} \right] = \frac{A^2}{8\sigma^2} \end{aligned}$$

and, therefore, we expect the following asymptotic behavior:

$$\text{min av. error probability} \xrightarrow{N \rightarrow +\infty} f(N) \cdot \exp\left(-N \cdot \frac{A^2}{8\sigma^2}\right) \quad (19)$$

where $f(N)$ varies slowly with N when N is large, compared with the exponential term:

$$\lim_{N \rightarrow \infty} \frac{\log f(N)}{N} = 0.$$

Indeed, in our example,

$$Q\left(\sqrt{\frac{N A^2}{4\sigma^2}}\right) \xrightarrow{N \rightarrow +\infty} \frac{1}{\underbrace{\sqrt{\frac{N A^2}{4\sigma^2}} \cdot \sqrt{2\pi}}_{f(N)}} \cdot \exp\left(-\frac{N A^2}{8\sigma^2}\right)$$

where we have used the asymptotic formula

$$Q(x) \xrightarrow{x \rightarrow +\infty} \frac{1}{x\sqrt{2\pi}} \cdot \exp(-x^2/2) \quad (20)$$

given e.g. in equation (2.4) of Kay-II.

Known DC Level in AWGN: Neyman-Pearson Approach

We now derive the Neyman-Pearson test for detecting a known DC level in AWGN. Recall that our likelihood-ratio test compares

$$T(\mathbf{x}) = \bar{x} \triangleq \frac{1}{N} \sum_{n=0}^{N-1} x[n]$$

with a threshold γ .

- If $T(\mathbf{x}) > \gamma$, decide \mathcal{H}_1 (i.e. reject \mathcal{H}_0),
- otherwise decide \mathcal{H}_0 (see also the discussion on p. 59).

Performance evaluation: Assuming (17), we have

$$T(\mathbf{x}) | a \sim \mathcal{N}(a, \sigma^2/N).$$

Therefore, $T(\mathbf{X}) | a = 0 \sim \mathcal{N}(0, \sigma^2/N)$, implying

$$P_{\text{FA}} = P[T(\mathbf{X}) > \gamma; a = 0] = Q\left(\frac{\gamma}{\sqrt{\sigma^2/N}}\right)$$

and we obtain the decision threshold as follows:

$$\gamma = \sqrt{\frac{\sigma^2}{N}} \cdot Q^{-1}(P_{\text{FA}}).$$

Now, $T(\mathbf{X}) | a = A \sim \mathcal{N}(A, \sigma^2/N)$, implying

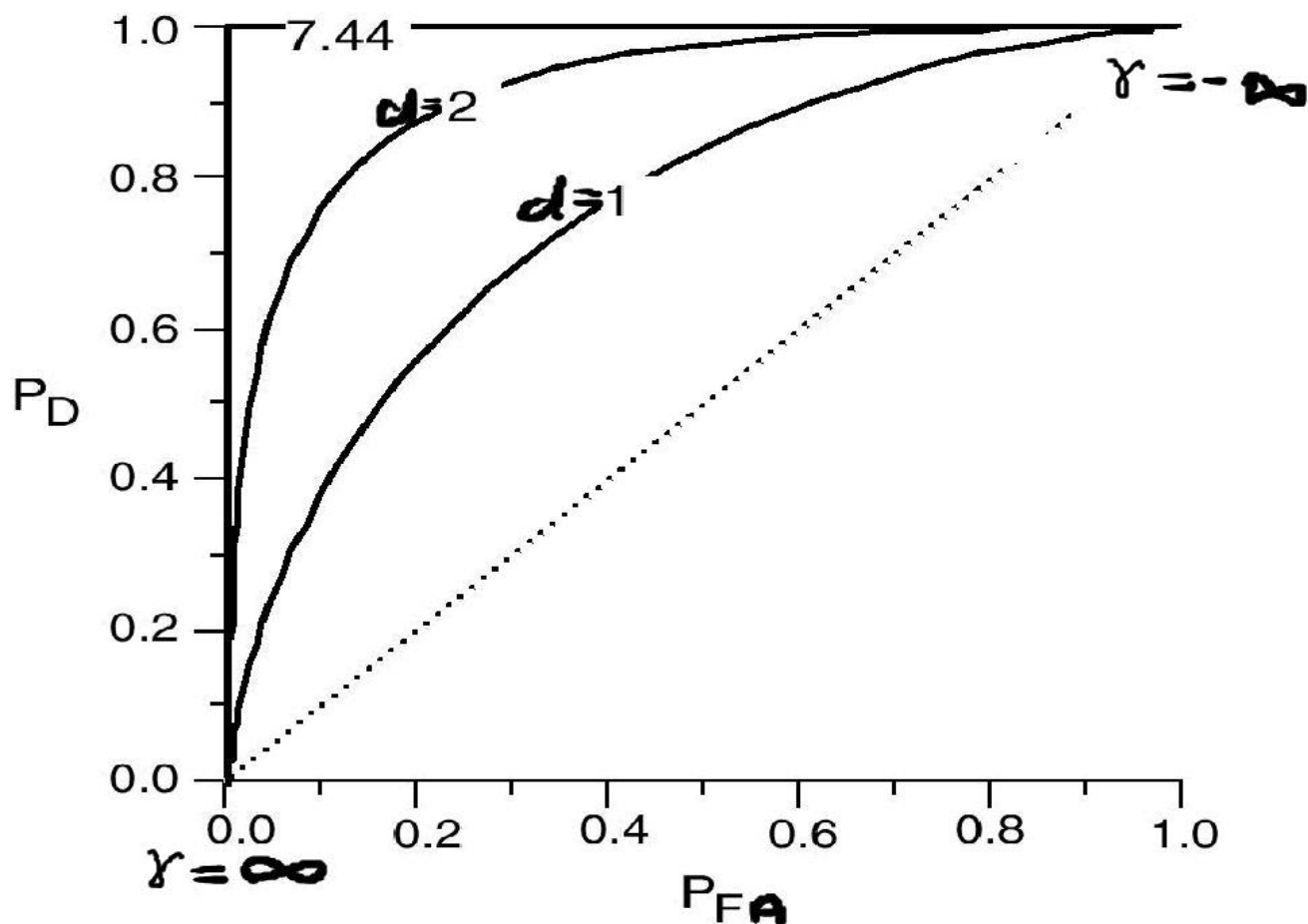
$$\begin{aligned} P_{\text{D}} &= P(T(\mathbf{X}) > \gamma | a = A) = Q\left(\frac{\gamma - A}{\sqrt{\sigma^2/N}}\right) \\ &= Q\left(Q^{-1}(P_{\text{FA}}) - \sqrt{\frac{A^2}{\sigma^2/N}}\right) \\ &= Q\left(Q^{-1}(P_{\text{FA}}) - \sqrt{\underbrace{NA^2/\sigma^2}_{\triangleq \text{SNR} = d^2}}\right). \end{aligned}$$

Given the false-alarm probability P_{FA} , the detection probability P_{D} depends only on the *deflection coefficient*:

$$d^2 = \frac{NA^2}{\sigma^2} = \frac{\{E[T(\mathbf{X}) | a = A] - E[T(\mathbf{X}) | a = 0]\}^2}{\text{var}[T(\mathbf{X} | a = 0)]}$$

which is also (a reasonable definition for) the signal-to-noise ratio (SNR).

Receiver Operating Characteristics (ROC)



$$P_D = Q(Q^{-1}(P_{FA}) - d).$$

Comments:

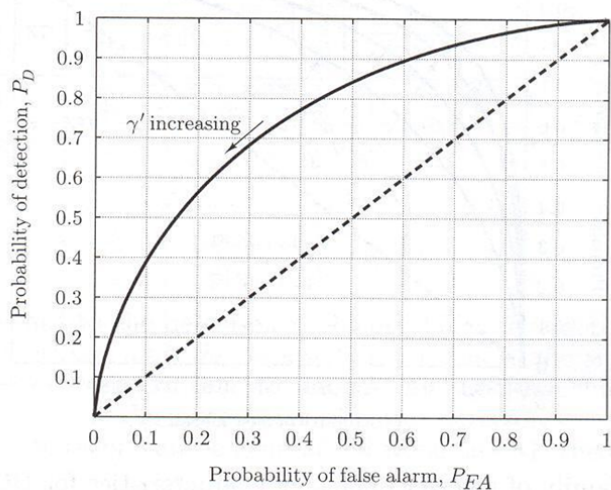
- As we raise the threshold γ , P_{FA} goes down but so does P_D .
- ROC should be above the 45° line — otherwise we can do better by flipping a coin.
- Performance improves with d^2 .

Typical Ways of Depicting the Detection Performance Under the Neyman-Pearson Setup

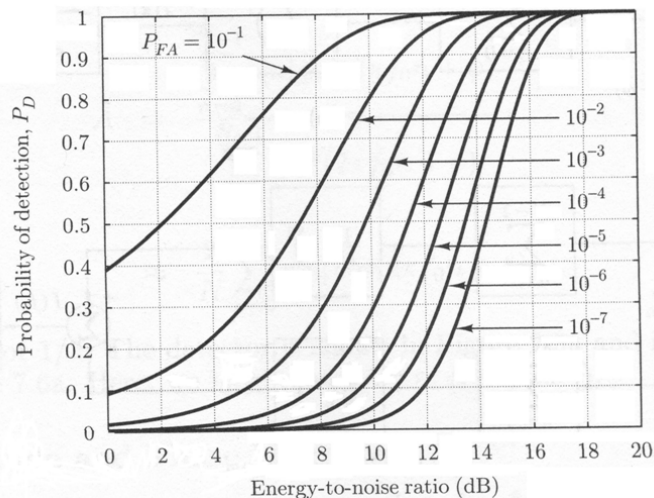
To analyze the performance of a Neyman-Pearson detector, we examine two relationships:

- Between P_D and P_{FA} , for a given SNR, called the *receiver operating characteristics* (ROC).
- Between P_D and SNR, for a given P_{FA} .

Here are examples of the two:



(a) ROC: P_D vs P_{FA}



(b) P_D vs SNR

see Figs. 3.8 and 3.5 in Kay-II, respectively.

Asymptotic (as $N \rightarrow \infty$ and $P_{\text{FA}} \searrow 0$) P_{D} and P_{M} for a Known DC Level in AWGN

We apply the Chernoff-Stein lemma, for which we need to compute the following K-L distance:

$$\begin{aligned} D\left(p(X_n | a = 0) \parallel p(X_n | a = A)\right) \\ = -\mathbb{E}_{p(x_n | a=0)} \left\{ \log \left[\frac{p(X_n | a = A)}{p(X_n | a = 0)} \right] \right\} \end{aligned}$$

where

$$\begin{aligned} p(x_n | a = 0) &= \mathcal{N}(0, \sigma^2) \\ \log \left[\frac{p(x_n | a = A)}{p(x_n | a = 0)} \right] &= -\frac{(x_n - A)^2}{2\sigma^2} + \frac{x_n^2}{2\sigma^2} = \frac{1}{2\sigma^2} \cdot (2Ax_n - A^2). \end{aligned}$$

Therefore,

$$D\left(p(X_n | \theta_0) \parallel p(X_n | \theta_1)\right) = \frac{A^2}{2\sigma^2}$$

and the Chernoff-Stein lemma predicts the following behavior of the detection probability as $N \rightarrow \infty$ and $P_{\text{FA}} \rightarrow 0$:

$$P_{\text{D}} \approx 1 - \underbrace{f(N) \cdot \exp\left(-N \cdot \frac{A^2}{2\sigma^2}\right)}_{\approx P_{\text{M}}}$$

where $f(N)$ is a slowly-varying function of N compared with the exponential term. In this case, the exact expression for P_M (P_D) is available and consistent with the Chernoff-Stein lemma:

$$\begin{aligned}
P_M &= 1 - Q\left(Q^{-1}(P_{\text{FA}}) - \sqrt{NA^2/\sigma^2}\right) \\
&= Q\left(\sqrt{NA^2/\sigma^2} - Q^{-1}(P_{\text{FA}})\right) \\
&\xrightarrow{N \rightarrow +\infty} \frac{1}{[\sqrt{NA^2/\sigma^2} - Q^{-1}(P_{\text{FA}})] \cdot \sqrt{2\pi}} \\
&\quad \cdot \exp\left\{-[\sqrt{NA^2/\sigma^2} - Q^{-1}(P_{\text{FA}})]^2/2\right\} \\
&= \frac{1}{[\sqrt{NA^2/\sigma^2} - Q^{-1}(P_{\text{FA}})] \cdot \sqrt{2\pi}} \\
&\quad \cdot \exp\left\{-NA^2/(2\sigma^2) - [Q^{-1}(P_{\text{FA}})]^2/2\right. \\
&\quad \left.+ Q^{-1}(P_{\text{FA}})\sqrt{NA^2/\sigma^2}\right\} \\
&= \frac{1}{[\sqrt{NA^2/\sigma^2} - Q^{-1}(P_{\text{FA}})] \cdot \sqrt{2\pi}} \\
&\quad \cdot \exp\left\{-[Q^{-1}(P_{\text{FA}})]^2/2 + Q^{-1}(P_{\text{FA}})\sqrt{NA^2/\sigma^2}\right\} \\
&\quad \cdot \underbrace{\exp[-NA^2/(2\sigma^2)]}_{\text{as predicted by the Chernoff-Stein lemma}} \tag{21}
\end{aligned}$$

where we have used the asymptotic formula (20). When P_{FA}

is small and N is large, the first two (green) terms in the above expression make a slowly-varying function $f(N)$ of N . Note that the exponential term in (21) *does not* depend on the false-alarm probability P_{FA} (or, equivalently, on the choice of the decision threshold). The Chernoff-Stein lemma asserts that *this is not a coincidence*.

Comment. For detecting a known DC level in AWGN:

- The slope of the exponential-decay term of P_{M} is

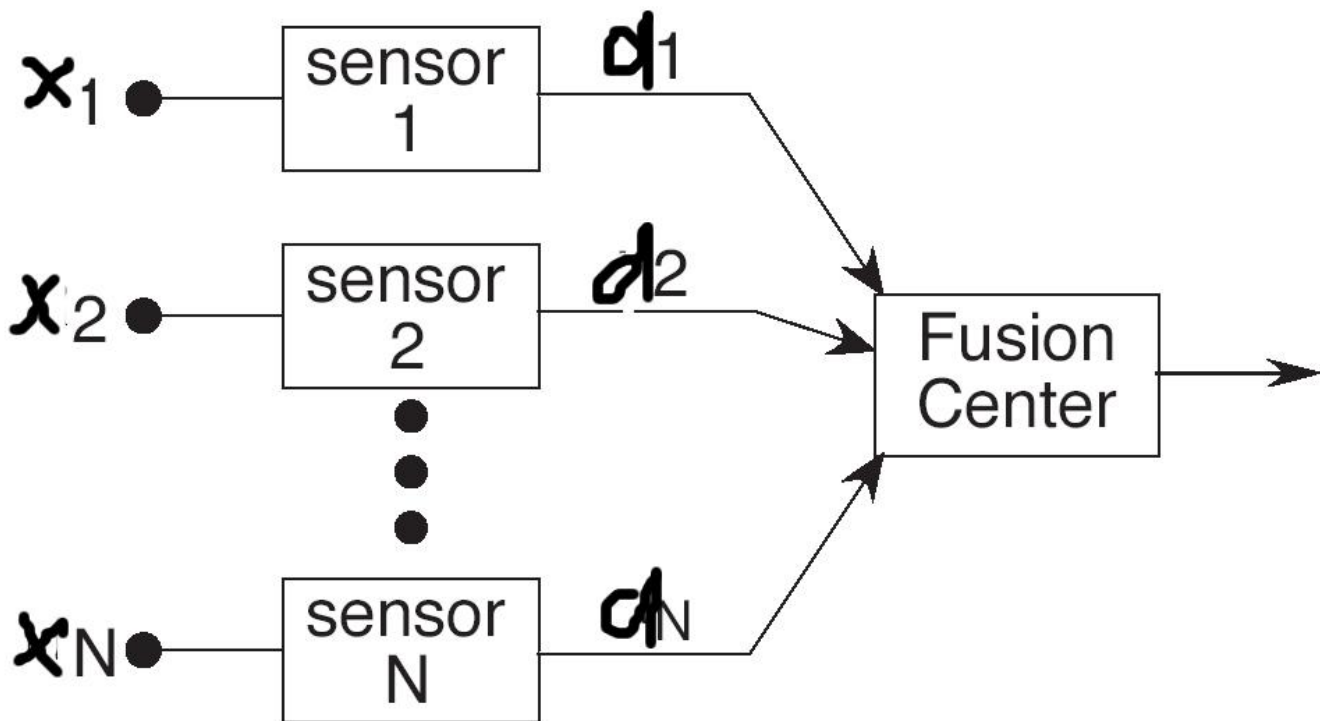
$$A^2/(2\sigma^2)$$

which is *different from* (larger than, in this case) the slope of the exponential decay of the minimum average error probability:

$$A^2/(8\sigma^2).$$

Decentralized Neyman-Pearson Detection for Simple Hypotheses

Consider a sensor-network scenario depicted by



Assumption: Observations made at N spatially distributed sensors (nodes) follow the same marginal probabilistic model:

$$\mathcal{H}_i : x_n \sim p(x_n | \theta_i)$$

where $n = 1, 2, \dots, N$ and $i \in \{0, 1\}$ for binary hypotheses.

Each node n makes a *hard local decision* d_n based on its local observation x_n and sends it to the *headquarters* (fusion center), which collects all the local decisions and makes the final *global*

decision about \mathcal{H}_0 or \mathcal{H}_1 . This structure is clearly suboptimal — it is easy to construct a better decision strategy in which each node sends its (quantized, in practice) likelihood ratio to the fusion center, rather than the decision only. However, such a strategy would have a higher communication (energy) cost.

We now go back to the decentralized detection problem. Suppose that each node n makes a local decision $d_n \in \{0, 1\}$, $n = 1, 2, \dots, N$ and transmits it to the fusion center. Then, the fusion center makes the global decision based on the likelihood ratio formed from the d_n s. The simplest fusion scheme is based on the assumption that the d_n s are conditionally independent given θ (which may not always be reasonable, but leads to an easy solution). We can now write

$$p(d_n | \theta_1) = \underbrace{P_{D,n}^{d_n} (1 - P_{D,n})^{1-d_n}}_{\text{Bernoulli pmf}}$$

where $P_{D,n}$ is the n th sensor's local detection probability. Similarly,

$$p(d_n | \theta_0) = \underbrace{P_{FA,n}^{d_n} (1 - P_{FA,n})^{1-d_n}}_{\text{Bernoulli pmf}}$$

where $P_{FA,n}$ is the n th sensor's local detection false-alarm

probability. Now,

$$\begin{aligned} \log \Lambda(\mathbf{d}) &= \sum_{n=1}^N \log \left[\frac{p(d_n | \theta_1)}{p(d_n | \theta_0)} \right] \\ &= \sum_{n=1}^N \log \left[\frac{P_{D,n}^{d_n} (1 - P_{D,n})^{1-d_n}}{P_{FA,n}^{d_n} (1 - P_{FA,n})^{1-d_n}} \right] \stackrel{\mathcal{H}_1}{\geq} \log \tau. \end{aligned}$$

To further simplify the exposition, we assume that all sensors have identical performance:

$$P_{D,n} = P_D, \quad P_{FA,n} = P_{FA}.$$

Define the number of sensors having $d_n = 1$:

$$u_1 = \sum_{n=1}^N d_n$$

Then, the log-likelihood ratio becomes

$$\log \Lambda(\mathbf{d}) = u_1 \log \left(\frac{P_D}{P_{FA}} \right) + (N - u_1) \log \left(\frac{1 - P_D}{1 - P_{FA}} \right) \stackrel{\mathcal{H}_1}{\geq} \log \tau$$

or

$$u_1 \log \left[\frac{P_D \cdot (1 - P_{FA})}{P_{FA} \cdot (1 - P_D)} \right] \stackrel{\mathcal{H}_1}{\geq} \log \tau + N \log \left(\frac{1 - P_{FA}}{1 - P_D} \right). \quad (22)$$

Clearly, each node's local decision d_n is meaningful only if $P_D > P_{FA}$, which implies

$$\frac{P_D \cdot (1 - P_{FA})}{P_{FA} \cdot (1 - P_D)} > 1$$

the logarithm of which is therefore positive, and the decision rule (22) further simplifies to

$$u_1 \stackrel{\mathcal{H}_1}{\gtrless} \tau'.$$

The Neyman-Person performance analysis of this detector is easy: the random variable U_1 is binomial given θ (i.e. conditional on the hypothesis) and, therefore,

$$P[U_1 = u_1] = \binom{N}{u_1} p^{u_1} (1 - p)^{N - u_1}$$

where $p = P_{FA}$ under \mathcal{H}_0 and $p = P_D$ under \mathcal{H}_1 . Hence, the “global” false-alarm probability is

$$P_{FA,global} = P[U_1 > \tau' | \theta_0] = \sum_{u_1 = \lceil \tau' \rceil}^N \binom{N}{u_1} \cdot P_{FA}^{u_1} \cdot (1 - P_{FA})^{N - u_1}.$$

Clearly, $P_{FA,global}$ will be a discontinuous function of τ' and therefore, we should choose our $P_{FA,global}$ specification from

the available discrete choices. But, if none of the candidate choices is acceptable, this means that our current system does not satisfy the requirements and a remedial action is needed, e.g. increasing the quantity (N) or improving the quality of sensors (through changing P_D and P_{FA}), or both.

Testing Multiple Hypotheses

Suppose now that we choose $\Theta_0, \Theta_1, \dots, \Theta_{M-1}$ that form a *partition* of the parameter space Θ :

$$\Theta_0 \cup \Theta_1 \cup \dots \cup \Theta_{M-1} = \Theta, \quad \Theta_i \cap \Theta_j = \emptyset \quad \forall i \neq j.$$

We wish to distinguish among $M > 2$ hypotheses, i.e. identify *which* hypothesis is true:

$$\begin{aligned} \mathcal{H}_0 & : \theta \in \Theta_0 \quad \text{versus} \\ \mathcal{H}_1 & : \theta \in \Theta_1 \quad \text{versus} \\ & \quad \quad \quad \vdots \quad \quad \quad \text{versus} \\ \mathcal{H}_{M-1} & : \theta \in \Theta_{M-1} \end{aligned}$$

and, consequently, our action space consists of M choices. We design a decision rule $\phi : \mathcal{X} \rightarrow (0, 1, \dots, M-1)$:

$$\phi(\mathbf{x}) = \begin{cases} 0, & \text{decide } \mathcal{H}_0, \\ 1, & \text{decide } \mathcal{H}_1, \\ \vdots & \\ M-1, & \text{decide } \mathcal{H}_{M-1} \end{cases}$$

where ϕ partitions the data space \mathcal{X} [i.e. the support of $p(\mathbf{x} | \theta)$] into M regions:

$$\text{Rule } \phi: \quad \mathcal{X}_0 = \{\mathbf{x} : \phi(\mathbf{x}) = 0\}, \dots, \mathcal{X}_{M-1} = \{\mathbf{x} : \phi(\mathbf{x}) = M-1\}.$$

We specify the loss function using $L(i | m)$, where, typically, the losses due to correct decisions are set to zero:

$$L(i | i) = 0, \quad i = 0, 1, \dots, M - 1.$$

Here, we adopt zero losses for correct decisions. Now, our posterior expected loss takes M values:

$$\begin{aligned} \rho_m(\mathbf{x}) &= \sum_{i=0}^{M-1} \int_{\Theta_i} L(m | i) p(\theta | \mathbf{x}) d\theta \\ &= \sum_{i=0}^{M-1} L(m | i) \int_{\Theta_i} p(\theta | \mathbf{x}) d\theta, \quad m = 0, 1, \dots, M - 1. \end{aligned}$$

Then, the Bayes' decision rule ϕ^* is defined via the following data-space partitioning:

$$\mathcal{X}_m^* = \{ \mathbf{x} : \rho_m(\mathbf{x}) = \min_{0 \leq l \leq M-1} \rho_l(\mathbf{x}) \}, \quad m = 0, 1, \dots, M - 1$$

or, equivalently, upon applying the Bayes' rule

$$\begin{aligned} \mathcal{X}_m^* &= \left\{ \mathbf{x} : m = \arg \min_{0 \leq l \leq M-1} \sum_{i=0}^{M-1} \int_{\Theta_i} L(l | i) p(\mathbf{x} | \theta) \pi(\theta) d\theta \right\} \\ &= \left\{ \mathbf{x} : m = \arg \min_{0 \leq l \leq M-1} \sum_{i=0}^{M-1} L(l | i) \int_{\Theta_i} p(\mathbf{x} | \theta) \pi(\theta) d\theta \right\}. \end{aligned}$$

The preposterior (Bayes) risk for rule $\phi(\mathbf{x})$ is

$$\begin{aligned}
 \mathbb{E}_{\mathbf{x},\theta}[\text{loss}] &= \sum_{m=0}^{M-1} \sum_{i=0}^{M-1} \int_{\mathcal{X}_m} \int_{\Theta_i} \mathbf{L}(m|i) p(\mathbf{x}|\theta) \pi(\theta) d\theta d\mathbf{x} \\
 &= \sum_{m=0}^{M-1} \int_{\mathcal{X}_m} \underbrace{\sum_{i=0}^{M-1} \mathbf{L}(m|i) \int_{\Theta_i} p(\mathbf{x}|\theta) \pi(\theta) d\theta}_{h_m(\theta)} d\mathbf{x}.
 \end{aligned}$$

Then, for an arbitrary $h_m(\mathbf{x})$,

$$\left[\sum_{m=0}^{M-1} \int_{\mathcal{X}_m} h_m(\mathbf{x}) d\mathbf{x} \right] - \left[\sum_{m=0}^{M-1} \int_{\mathcal{X}_{m^*}} h_m(\mathbf{x}) d\mathbf{x} \right] \geq 0$$

which verifies that the Bayes' decision rule ϕ^* minimizes the preposterior (Bayes) risk.

Special Case: $L(i|i) = 0$ and $L(m|i) = 1$ for $i \neq m$ (0-1 loss), implying that $\rho_m(\mathbf{x})$ can be written as

$$\begin{aligned}\rho_m(\mathbf{x}) &= \sum_{i=0, i \neq m}^{M-1} \int_{\Theta_i} p(\theta | \mathbf{x}) d\theta \\ &= \underbrace{\text{const}}_{\text{not a function of } m} - \int_{\Theta_m} p(\theta | \mathbf{x}) d\theta\end{aligned}$$

and

$$\begin{aligned}\mathcal{X}_m^* &= \left\{ \mathbf{x} : m = \arg \max_{0 \leq l \leq M-1} \int_{\Theta_l} p(\theta | \mathbf{x}) d\theta \right\} \\ &= \left\{ \mathbf{x} : m = \arg \max_{0 \leq l \leq M-1} P[\theta \in \Theta_l | \mathbf{x}] \right\} \quad (23)\end{aligned}$$

which is the MAP rule, as expected.

Simple hypotheses: Let us specialize (23) to simple hypotheses ($\Theta_m = \{\theta_m\}$, $m = 0, 1, \dots, M-1$):

$$\mathcal{X}_m^* = \left\{ \mathbf{x} : m = \arg \max_{0 \leq l \leq M-1} p(\theta_l | \mathbf{x}) \right\}$$

or, equivalently,

$$\mathcal{X}_m^* = \left\{ \mathbf{x} : m = \arg \max_{0 \leq l \leq M-1} [\pi_l p(\mathbf{x} | \theta_l)] \right\}, \quad m = 0, 1, \dots, M-1$$

where

$$\pi_0 = \pi(\theta_0), \quad \dots, \pi_{M-1} = \pi(\theta_{M-1})$$

define the prior pmf of the M -ary discrete random variable θ (recall that $\theta \in \{\theta_0, \theta_1, \dots, \theta_{M-1}\}$). If π_i , $i = 0, 1, \dots, M - 1$ are all equal:

$$\pi_0 = \pi_1 = \dots = \pi_{M-1} = \frac{1}{M}$$

the resulting test

$$\begin{aligned} \mathcal{X}_m^* &= \left\{ \mathbf{x} : m = \arg \max_{0 \leq l \leq M-1} \left[\frac{1}{M} p(\mathbf{x} | \theta_l) \right] \right\} \\ &= \left\{ \mathbf{x} : m = \arg \max_{0 \leq l \leq M-1} \underbrace{p(\mathbf{x} | \theta_l)}_{\text{likelihood}} \right\} \end{aligned} \quad (24)$$

is the *maximum-likelihood test*; this name is easy to justify after inspecting (24) and noting that the computation of the optimal decision region \mathcal{X}_m^* requires *the maximization of the likelihood $p(\mathbf{x} | \theta)$ with respect to the parameter $\theta \in \{\theta_0, \theta_1, \dots, \theta_{M-1}\}$.*

Summary: Bayesian Decision Approach versus Neyman-Pearson Approach

- The Neyman-Pearson approach appears particularly suitable for applications where the null hypothesis can be formulated as absence of signal or perhaps, absence of statistical difference between two data sets (treatment versus placebo, say).
- In the Neyman-Pearson approach, the null hypothesis is treated very differently from the alternative. (If the null hypothesis is true, we wish to control the false-alarm rate, which is different from our desire to maximize the probability of detection when the alternative is true). Consequently, our decisions should also be treated differently. If the likelihood ratio is large enough, we decide to *accept \mathcal{H}_1 (or reject \mathcal{H}_0)*. However, if the likelihood ratio is not large enough, we decide *not to reject \mathcal{H}_0* because, in this case, it may be that either
 - (i) \mathcal{H}_0 is true or
 - (ii) \mathcal{H}_0 is false but the test has low detection probability (power) (e.g. because the signal level is small compared with noise or we collected too small number of observations).

- The Bayesian decision framework is suitable for communications applications as it can easily handle multiple hypotheses (unlike the Neyman-Pearson framework).
 - **0-1 loss:** In communications applications, we typically select a 0-1 loss, implying that all hypotheses are treated equally (i.e. we could change the roles of null and alternative hypotheses without any problems). Therefore, interpretations of our decisions are also straightforward. Furthermore, in this case, the Bayes' decision rule is also optimal in terms of minimizing the average error probability, which is one of the most popular performance criteria in communications.

P Values

Reporting “accept \mathcal{H}_0 ” or “accept \mathcal{H}_1 ” is not very informative. Instead, we could vary P_{FA} and examine how our report would change.

Generally, if \mathcal{H}_1 is accepted for a certain specified P_{FA} , it will be accepted for $P'_{\text{FA}} > P_{\text{FA}}$. Therefore, there exists a smallest P_{FA} at which \mathcal{H}_1 is accepted. This motivates the introduction of the p value.

To be more precise (and be able to handle composite hypotheses), here is a definition of a *size of a hypothesis test*.

Definition 1. The *size of a hypothesis test* described by

$$\text{Rule } \phi: \quad \mathcal{X}_0 = \{x : \phi(x) = 0\}, \quad \mathcal{X}_1 = \{x : \phi(x) = 1\}.$$

is defined as follows:

$$\alpha = \max_{\theta \in \Theta_0} P[\mathbf{x} \in \mathcal{X}_1 \mid \theta] = \text{max possible } P_{\text{FA}}.$$

A hypothesis test is *said to have level α* if its size is less than or equal to α . Therefore, a level- α test is *guaranteed* to have a false-alarm probability less than or equal to α .

Definition 2. Consider a Neyman-Pearson-type setup where our test

$$\text{Rule } \phi_\alpha: \mathcal{X}_{0,\alpha} = \{\mathbf{x} : \phi(\mathbf{x}) = 0\}, \quad \mathcal{X}_{1,\alpha} = \{\mathbf{x} : \phi(\mathbf{x}) = 1\}. \quad (25)$$

achieves a specified size α , meaning that,

$$\begin{aligned} \alpha &= \max \text{ possible } P_{\text{FA}} \\ &= \max_{\theta \in \Theta_0} P[\mathbf{x} \in \mathcal{X}_{1,\alpha} | \theta] \quad (\text{composite hypotheses}) \end{aligned}$$

or, in the simple-hypothesis case ($\Theta_0 = \{\theta_0\}, \Theta_1 = \{\theta_1\}$):

$$\alpha = \overbrace{P[\mathbf{x} \in \mathcal{X}_{1,\alpha} | \theta = \theta_0]}^{P_{\text{FA}}} \quad (\text{simple hypotheses}).$$

We suppose that, for every $\alpha \in (0, 1)$, we have a size- α test with decision regions (25). Then, the p value for this test is the smallest level α at which we can declare \mathcal{H}_1 :

$$p \text{ value} = \inf\{\alpha : \mathbf{x} \in \mathcal{X}_{1,\alpha}\}.$$

Informally, the p value is a measure of evidence for supporting \mathcal{H}_1 . For example, p values less than 0.01 are considered *very strong evidence* supporting \mathcal{H}_1 .

There are a lot of warnings (and misconceptions) regarding p values. Here are the most important ones.

Warning: A large p value is *not* strong evidence in favor of \mathcal{H}_0 ; a large p value can occur for two reasons:

- (i) \mathcal{H}_0 is true or
- (ii) \mathcal{H}_0 is false but the test has low detection probability (power).

Warning: *Do not* confuse the p value with

$$P[\mathcal{H}_0 \mid \text{data}] = P[\theta \in \Theta_0 \mid \mathbf{x}]$$

which is used in Bayesian inference. **The p value is *not* the probability that \mathcal{H}_0 is true.**

Theorem 1. *Suppose that we have a size- α test of the form*

$$\text{declare } \mathcal{H}_1 \text{ if and only if } T(\mathbf{x}) \geq c_\alpha.$$

Then, the p value for this test is

$$p \text{ value} = \max_{\theta \in \Theta_0} P[T(\mathbf{X}) \geq T(\mathbf{x}) \mid \theta]$$

where x is the observed value of \mathbf{X} . For $\Theta_0 = \{\theta_0\}$:

$$p \text{ value} = P[T(\mathbf{X}) \geq T(\mathbf{x}) \mid \theta = \theta_0].$$

In words, Theorem 1 states that

The p value is the probability that, under \mathcal{H}_0 , a random data realization \mathbf{X} is observed yielding a value of the test statistic $T(\mathbf{X})$ that is greater than or equal to what has actually been observed (i.e. $T(\mathbf{x})$).

Note: This interpretation requires that we allow the experiment to be repeated many times. This is what Bayesians criticize by saying that “data that have never been observed are used for inference.”

Theorem 2. *If the test statistics has a continuous distribution, then, under $\mathcal{H}_0 : \theta = \theta_0$, the p value has a $\text{uniform}(0, 1)$ distribution. Therefore, if we declare \mathcal{H}_1 (reject \mathcal{H}_0) when the p value is less than or equal to α , the probability of false alarm is α .*

In other words, if \mathcal{H}_0 is true, the p value is like a random draw from an $\text{uniform}(0, 1)$ distribution. If \mathcal{H}_1 is true and if we repeat the experiment many times, the random p values will concentrate closer to zero.

Multiple Testing

We may conduct many hypothesis tests in some applications, e.g.

- bioinformatics and
- sensor networks.

Here, we perform *many* (typically binary) *tests* (one test per node in a sensor network, say). This is different from *testing multiple hypotheses* that we considered on pp. 54–58, where we performed *a single test of multiple hypotheses*. For a sensor-network related discussion on multiple testing, see

E.B. Ermiş, M. Alanyali, and V. Saligrama, “Search and discovery in an uncertain networked world,” *IEEE Signal Processing Magazine*, vol. 23, pp. 107–118, Jul. 2006.

Suppose that each test is conducted with false-alarm probability $P_{\text{FA}} = \alpha$. For example, in a sensor-network setup, each node conducts a test based on its local data.

Although the chance of false alarm at each node is only α , the chance of at least one falsely alarmed node is much higher, since there are many nodes. Here, we discuss two ways to deal with this problem.

Consider M hypothesis tests:

$$H_{0i} \text{ versus } H_{1i}, i = 1, 2, \dots, M$$

and denote by p_1, p_2, \dots, p_M the p values for these tests. Then, the *Bonferroni method* does the following:

Given the p values p_1, p_2, \dots, p_m , accept \mathcal{H}_{1i} if

$$p_i < \frac{\alpha}{M}.$$

Recall the *union-of-events bound*:

$$P\left[\bigcup_{i=1}^n A_i\right] \leq \sum_{i=1}^n P[A_i].$$

which holds for arbitrary events A_i , $i = 1, 2, \dots, n$, see handout # 2 in EE 420x notes.

Theorem 3. *If we apply the Bonferroni method, the probability of any individual false alarm is less than or equal to α .*

Proof. Denote by A the event that there is at least one false alarm and by A_i the event that the i th node is falsely alarmed.

Then

$$P\{A\} = P\left\{\bigcup_{i=1}^M A_i\right\} \leq \sum_{i=1}^n P\{A_i\} = \sum_{i=1}^M \frac{\alpha}{M} = \alpha.$$

□

Comments: Suppose now that the tests are *statistically independent* and consider the smallest of the M p values. Under \mathcal{H}_{0i} , p_i are uniform(0, 1). Then

$$P\{\min\{P_1, P_2, \dots, P_M\} > x\} \Big| \text{assuming } \mathcal{H}_{0i} = (1 - x)^M \\ i = 1, 2, \dots, M$$

yielding the “proper” p value to be attached to $\min\{p_1, p_2, \dots, p_m\}$ as

$$1 - (1 - \min\{p_1, p_2, \dots, p_M\})^M \quad (26)$$

and if $\min\{p_1, p_2, \dots, p_M\}$ is small and M not too large, (26) will be close to $m \min\{p_1, p_2, \dots, p_M\}$.

False Discovery Rate: Sometimes it is reasonable to control *false discovery rate (FDR)*, which we introduce below.

Suppose that we accept \mathcal{H}_{1i} for all i for which

$$p_i < \text{threshold } \tau$$

and let $m_0 + m_1 = m$ be

number of true \mathcal{H}_{0i} hypotheses + number of true \mathcal{H}_{1i} hypotheses = total number of hypotheses (nodes, say).

# of different outcomes	\mathcal{H}_0 not rejected	\mathcal{H}_1 declared	total
\mathcal{H}_0 true	U	V	m_0
\mathcal{H}_1 true	T	S	m_1
total	$m - R$	R	m

Define the *false discovery proportion* (FDP) as

$$\text{FDP} = \begin{cases} V/R, & R > 0, \\ 0, & R = 0 \end{cases}$$

which is simply the proportion of *incorrect* \mathcal{H}_1 decisions. Now, define

$$\text{FDR} = \text{E}[\text{FDP}].$$

The Benjamini-Hochberg (BH) Method

(i) Denote the ordered p values by $p_{(1)} < p_{(2)} < \cdots < p_{(m)}$.

(ii) Define

$$l_i = \frac{i \alpha}{C_m m} \quad \text{and} \quad R = \max\{i : p_{(i)} < l_i\}$$

where C_m is defined to be 1 if the p values are independent and $C_m = \sum_{i=1}^m (1/i)$ otherwise.

(iii) Define the *BH rejection threshold* $\tau = p_{(R)}$.

(iv) Accept all \mathcal{H}_{1i} for which $p_{(i)} \leq \tau$.

Theorem 4. *(formulated and proved by Benjamini and Hochberg) If the above BH method is applied, then, regardless of how many null hypotheses are true and regardless of the distribution of the p values when the null hypothesis is false,*

$$\text{FDR} = \text{E} [\text{FDP}] \leq \frac{m_0}{m} \alpha \leq \alpha.$$