

# EE/ACM 150 - Applications of Convex Optimization in Signal Processing and Communications

## Lecture 15

Andre Tkacenko

Signal Processing Research Group  
Jet Propulsion Laboratory

May 22, 2012



Caltech

# Outline

- 1 Parametric Distribution Estimation
  - Maximum Likelihood Estimation
  - Maximum A Posteriori Probability Estimation
- 2 Nonparametric Distribution Estimation
- 3 Hypothesis Testing and Optimal Detector Design
  - Deterministic and Randomized Detectors
  - Optima Detector Design
  - Multicriterion Formulation and Scalarization
  - Binary Hypothesis Testing
  - Robust Detectors

# Maximum Likelihood and the Log-Likelihood Function

Let  $\mathbf{y} \in \mathbb{R}^m$  be a random variable whose distribution depends on a parameter  $\mathbf{x} \in \mathbb{R}^n$ . We will denote this distribution by  $p_{\mathbf{y}|\mathbf{x}}(\mathbf{y})$  and refer to it as the *likelihood function*.

**Example:** If  $\mathbf{v} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , the distribution of  $\mathbf{v}$  depends on  $\boldsymbol{\mu} \in \mathbb{R}^m$  and  $\boldsymbol{\Sigma} \in \mathbb{S}_+^m$ .

As many common distributions are log-concave, and often independent observations  $y_i$  will be made to form the vector observation  $\mathbf{y}$ , it is convenient to work with the logarithm of the distribution  $p_{\mathbf{y}|\mathbf{x}}(\mathbf{y})$  called the *log-likelihood function*, given by

$$\ell(\mathbf{x}) \triangleq \log p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}) .$$

In *maximum likelihood (ML) estimation*, the parameter  $\mathbf{x}$  is determined as the argument which maximizes the likelihood (or log-likelihood) function.

## Maximum Likelihood (ML) Estimate:

$$\hat{\mathbf{x}}_{\text{ml}} = \underset{\mathbf{x}}{\operatorname{argmax}} p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}) = \underset{\mathbf{x}}{\operatorname{argmax}} \ell(\mathbf{x}) .$$

- Here,  $\mathbf{y}$  is the observed value.
- We can add constraints that  $\mathbf{x} \in \mathcal{C}$  explicitly, or define  $p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}) = 0$  for  $\mathbf{x} \notin \mathcal{C}$ .
- Determining the ML estimate is a convex optimization problem if  $\log p_{\mathbf{y}|\mathbf{x}}(\mathbf{y})$  is concave in  $\mathbf{x}$  for fixed  $\mathbf{y}$ .

# Linear Measurements with i.i.d. Noise

In many cases, observations are made following a *linear measurement model* given by

$$y_i = \mathbf{a}_i^T \mathbf{x} + v_i, \quad i = 1, \dots, m.$$

- $\mathbf{x} \in \mathbb{R}^n$  is the vector of unknown parameters.
- $v_i$  are *independent and identically distributed (i.i.d.)* measurement noise with density  $p_v(v)$ .
- $y_i$  is the measurement:  $\mathbf{y} \in \mathbb{R}^m$  has density  $p_{\mathbf{y}|\mathbf{x}}(\mathbf{y}) = \prod_{i=1}^m p_v(y_i - \mathbf{a}_i^T \mathbf{x})$ .

**ML estimation:** An ML estimate  $\hat{\mathbf{x}}_{\text{ml}}$  is any solution  $\mathbf{x}$  of

$$\text{maximize } \ell(\mathbf{x}) = \sum_{i=1}^m \log p_v(y_i - \mathbf{a}_i^T \mathbf{x}) \quad ,$$

with observations  $y_i$  for  $i = 1, \dots, m$ .

**Examples:**

- Gaussian noise: If  $v_i \sim \mathcal{N}(0, \sigma^2)$ , then  $p_v(v) = (2\pi\sigma^2)^{-1/2} e^{-v^2/(2\sigma^2)}$  and so,

$$\ell(\mathbf{x}) = -(m/2) \log(2\pi\sigma^2) - (1/(2\sigma^2)) \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 .$$

Thus, the ML estimate is the least-squares (LS) solution.

- Laplacian noise: If  $p_v(v) = (1/(2a)) e^{-|v|/a}$  for some  $a > 0$ , then we have,

$$\ell(\mathbf{x}) = -m \log(2a) - (1/a) \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_1 .$$

Hence, the ML estimate is the  $\ell_1$ -norm solution.

- Uniform noise: If  $p_v(v) = 1/(2a)$  for  $v \in [-a, a]$  and  $p_v(v) = 0$  otherwise, then we have

$$\ell(\mathbf{x}) = \begin{cases} -m \log(2a) , & \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_\infty \leq a \\ -\infty , & \text{otherwise} \end{cases} .$$

So, an ML estimate is any  $\mathbf{x}$  with  $\|\mathbf{A}\mathbf{x} - \mathbf{y}\|_\infty \leq a$ .

# ML Interpretation of Penalty Function Approximation

We can interpret any penalty function approximation problem of the form

$$\text{minimize } \sum_{i=1}^m \phi(b_i - \mathbf{a}_i^T \mathbf{x}) ,$$

as an ML estimation problem with a linear measurement model with i.i.d. noise. Specifically, these two problems are identical when the noise density  $p_v(v)$  satisfies

$$p_v(v) = \frac{e^{-\phi(v)}}{\int e^{-\phi(u)} du} \iff \phi(v) = -C \log p_v(v) ,$$

and the observation is  $\mathbf{y} = \mathbf{b}$ .

## Interpretations:

- Penalty functions which strongly (weakly) penalize large residuals will lead to noise densities with small (large) tails.
- Laplacian noise density has much larger tails than Gaussian noise density. Hence, the Laplacian noise based ML estimation method will not penalize large residuals as severely as that corresponding to the Gaussian noise density.
- Residuals corresponding to ML estimation with a Gaussian density will tend to be small, whereas those for a Laplacian density will sparse (many residuals near zero and a few large outliers).

# Counting Problems with Poisson Distribution

In many cases, the random variable  $y$  has *Poisson distribution* with mean  $\mu > 0$ :

$$\Pr\{y = k\} = \frac{e^{-\mu} \mu^k}{k!}.$$

Often,  $y$  will represent the count or number of events of a Poisson process over some period of time (i.e., photon arrival counts, traffic accidents, etc.). In a simple statistical model,  $\mu$  is modeled as an affine function of a vector  $\mathbf{u} \in \mathbb{R}^n$ :

$$\mu = \mathbf{a}^T \mathbf{u} + b.$$

- $\mathbf{a} \in \mathbb{R}^n$  and  $b \in \mathbb{R}$  are *model parameters*, whereas  $\mathbf{u} \in \mathbb{R}^n$  is an observed vector of *explanatory variables*.
- The goal is to estimate  $\mathbf{a}$  and  $b$  from  $m$  observations  $(u_i, y_i)$  for  $i = 1, \dots, m$ .

From this, the log-likelihood function is given by the following.

$$\ell(\mathbf{a}, b) = \sum_{i=1}^m \left( y_i \log(\mathbf{a}^T u_i + b) - (\mathbf{a}^T u_i + b) - \log(y_i!) \right).$$

We can find an ML estimate of  $\mathbf{a}$  and  $b$  by solving the convex optimization problem

$$\text{maximize} \quad \sum_{i=1}^m \left( y_i \log(\mathbf{a}^T u_i + b) - (\mathbf{a}^T u_i + b) \right),$$

with variables  $\mathbf{a} \in \mathbb{R}^n$  and  $b \in \mathbb{R}$ .

# Logistic Regression

Consider a random variable  $y \in \{0, 1\}$  with

$$\Pr\{y = 1\} = p = \frac{\exp(\mathbf{a}^T \mathbf{u} + b)}{1 + \exp(\mathbf{a}^T \mathbf{u} + b)}, \quad \Pr\{y = 0\} = 1 - p = \frac{1}{1 + \exp(\mathbf{a}^T \mathbf{u} + b)}.$$

- $\mathbf{a} \in \mathbb{R}^n, b \in \mathbb{R}$  are model parameters;  $\mathbf{u} \in \mathbb{R}^n$  is an observed explanatory variable vector.
- The goal is to estimate  $\mathbf{a}$  and  $b$  from  $m$  observations  $(u_i, y_i)$  for  $i = 1, \dots, m$ .

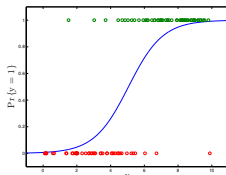
**Log-likelihood function:** Assuming the data has been ordered such that  $y_1 = \dots = y_k = 1$  and  $y_{k+1} = \dots = y_m = 0$ , we get the following.

$$\begin{aligned} \ell(\mathbf{a}, b) &= \log \left( \prod_{i=1}^k \frac{\exp(\mathbf{a}^T u_i + b)}{1 + \exp(\mathbf{a}^T u_i + b)} \prod_{i=k+1}^m \frac{1}{1 + \exp(\mathbf{a}^T u_i + b)} \right), \\ &= \sum_{i=1}^k (\mathbf{a}^T u_i + b) - \sum_{i=1}^m \log(1 + \exp(\mathbf{a}^T u_i + b)), \end{aligned}$$

which is concave in  $\mathbf{a}$  and  $b$ .

**Example:** ( $n = 1, m = 50$  measurements)

- Circles show 50 points  $(u_i, y_i)$ .
- Solid curve is the ML estimate of  $p = \exp(au + b) / (1 + \exp(au + b))$ .



# Covariance Estimate for Gaussian Random Variables

Suppose  $\mathbf{y} \in \mathbb{R}^n$  is a Gaussian random vector with zero mean and covariance matrix  $\mathbf{R} = E[\mathbf{y}\mathbf{y}^T]$ . In many cases, it may be imperative to estimate  $\mathbf{R}$  based on independent observations of  $\mathbf{y}$ . Note that here, the density of  $\mathbf{y}$  is given by

$$p_{\mathbf{y}|\mathbf{R}}(\mathbf{y}) = (2\pi)^{-n/2} (\det(\mathbf{R}))^{-1/2} \exp\left(-\mathbf{y}^T \mathbf{R}^{-1} \mathbf{y}/2\right).$$

Based on  $N$  independent observations  $\mathbf{y}_1, \dots, \mathbf{y}_N \in \mathbb{R}^n$ , we would like to find an ML estimate of  $\mathbf{R}$ . The log-likelihood function is given by the following.

$$\ell(\mathbf{R}) = -\frac{Nn}{2} \log(2\pi) - \frac{N}{2} \log \det \mathbf{R} - \frac{N}{2} \text{tr}(\mathbf{R}^{-1} \mathbf{Y}), \text{ where } \mathbf{Y} \triangleq \frac{1}{N} \sum_{k=1}^N \mathbf{y}_k \mathbf{y}_k^T.$$

Here,  $\mathbf{Y}$  is the *sample covariance matrix* for the observations  $\mathbf{y}_1, \dots, \mathbf{y}_N \in \mathbb{R}^n$ . The log-likelihood function is *not* a concave function of  $\mathbf{R}$ , however it *is* a concave function of  $\mathbf{S} \triangleq \mathbf{R}^{-1}$ , which is called the *information matrix*. In this case, the new log-likelihood function  $\ell(\mathbf{S})$  has the form

$$\ell(\mathbf{S}) = -(Nn/2) \log(2\pi) + (N/2) \log \det \mathbf{S} - (N/2) \text{tr}(\mathbf{S}\mathbf{Y}),$$

and so, the ML estimate of  $\mathbf{S}$  (and hence  $\mathbf{R}$ ) can be found by solving the problem

$$\begin{aligned} & \text{maximize} && \log \det \mathbf{S} - \text{tr}(\mathbf{S}\mathbf{Y}) \\ & \text{subject to} && \mathbf{S} \in \mathcal{S} \end{aligned},$$

where  $\mathcal{S}$  represents our prior knowledge of  $\mathbf{S} = \mathbf{R}^{-1}$ .



# Bayesian Version of ML Estimation

In *maximum a posteriori probability (MAP) estimation*, which can be considered as a Bayesian version of ML estimation, we assume that the underlying parameter  $\mathbf{x}$  to be estimated is random with some *a priori* or prior probability density  $p_{\mathbf{x}}(\mathbf{x})$ .

If  $\mathbf{x}$  (the vector to be estimated) and  $\mathbf{y}$  (the observation) are random variables with joint probability density function (pdf)  $p_{\mathbf{x},\mathbf{y}}(\mathbf{x}, \mathbf{y})$ , then we have the following.

$$p_{\mathbf{x},\mathbf{y}}(\mathbf{x}, \mathbf{y}) = p_{\mathbf{x}|\mathbf{y}}(\mathbf{x}, \mathbf{y}) p_{\mathbf{y}}(\mathbf{y}) = p_{\mathbf{y}|\mathbf{x}}(\mathbf{x}, \mathbf{y}) p_{\mathbf{x}}(\mathbf{x}) ,$$

where, for instance,  $p_{\mathbf{x}|\mathbf{y}}(\mathbf{x}, \mathbf{y})$  is the conditional density of  $\mathbf{x}$  given  $\mathbf{y}$ . Specifically, for an observed value  $\mathbf{y} = \mathbf{y}$ ,  $p_{\mathbf{x}|\mathbf{y}}$  is the *a posteriori* or posterior density of  $\mathbf{x}$ . For MAP estimation, we seek the vector  $\mathbf{x}$  which maximizes this quantity.

## MAP Estimate:

$$\hat{\mathbf{x}}_{\text{map}} = \underset{\mathbf{x}}{\operatorname{argmax}} p_{\mathbf{x}|\mathbf{y}}(\mathbf{x}, \mathbf{y}) = \underset{\mathbf{x}}{\operatorname{argmax}} p_{\mathbf{y}|\mathbf{x}}(\mathbf{x}, \mathbf{y}) p_{\mathbf{x}}(\mathbf{x}) = \underset{\mathbf{x}}{\operatorname{argmax}} p_{\mathbf{x},\mathbf{y}}(\mathbf{x}, \mathbf{y}) .$$

Taking logarithms, we can express the MAP estimate as

$$\hat{\mathbf{x}}_{\text{map}} = \underset{\mathbf{x}}{\operatorname{argmax}} (\log p_{\mathbf{y}|\mathbf{x}}(\mathbf{x}, \mathbf{y}) + \log p_{\mathbf{x}}(\mathbf{x})) .$$

- From the last expression for the MAP estimate, we see that the second term penalizes unlikely choices of  $\mathbf{x}$ , according to the prior density  $p_{\mathbf{x}}(\mathbf{x})$ .
- If the log-likelihood function is concave and the prior density for  $\mathbf{x}$  is log-concave, the resulting MAP estimation problem will be convex.

# Examples of MAP Estimation Problems

**Linear measurements with i.i.d. noise:** For the linear measurement model

$$y_i = \mathbf{a}_i^T \mathbf{x} + v_i, \quad i = 1, \dots, m,$$

where the  $v_i$  are i.i.d. with density  $p_v$ , the MAP estimation problem becomes

$$\text{maximize} \quad \sum_{i=1}^m \log p_v(y_i - \mathbf{a}_i^T \mathbf{x}) + \log p_{\mathbf{x}}(\mathbf{x}) \quad .$$

If  $p_v$  and  $p_{\mathbf{x}}$  are log-concave, this problem is convex.

**MAP estimation with perfect linear measurements:** Suppose we have  $m$  deterministic linear measurements  $\mathbf{y} = \mathbf{A}\mathbf{x}$ . In this case, the MAP estimation problem becomes the following.

$$\begin{aligned} &\text{maximize} && \log p_{\mathbf{x}}(\mathbf{x}) \\ &\text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{y} \quad . \end{aligned}$$

If  $p_{\mathbf{x}}$  is log-concave, this is a convex problem.

Assuming the parameters  $x_i$  are i.i.d. with density  $p_x(x)$ , the MAP estimation problem becomes

$$\begin{aligned} &\text{maximize} && \sum_{i=1}^n \log p_x(x_i) && \iff && \text{maximize} && \sum_{i=1}^n \phi(x_i) \\ &\text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{y} && && \text{subject to} && \mathbf{A}\mathbf{x} = \mathbf{b} \quad , \end{aligned}$$

where  $\phi(u) = -\log p_x(u)$  and  $\mathbf{b} = \mathbf{y}$ . In other words, the MAP estimation problem can be expressed as a least-penalty problem. Conversely, any least-penalty problem can be expressed as a MAP estimation problem with  $m$  perfect measurements  $\mathbf{b} = \mathbf{A}\mathbf{x}$  and  $x_i$  i.i.d. with density

$$p_x(x) = \frac{e^{-\phi(x)}}{\int e^{-\phi(u)} du} .$$

# Introduction to Nonparametric Distribution Estimation

In nonparametric distribution estimation, we are interested in estimating the pdf of a continuous random variable or the probability mass function (pmf) of a discrete random variable.

Here, we will consider a discrete random variable  $X$  with values in the finite set  $\{\alpha_1, \dots, \alpha_n\} \subseteq \mathbb{R}$ .

- The distribution of  $X$  is characterized by  $\mathbf{p} \in \mathbb{R}^n$ , where  $p_k = \Pr\{X = \alpha_k\}$ , and satisfies

$$\mathbf{p} \succeq \mathbf{0}, \mathbf{1}^T \mathbf{p} = 1.$$

- Conversely, if  $\mathbf{p} \in \mathbb{R}^n$  satisfies  $\mathbf{p} \succeq \mathbf{0}$  and  $\mathbf{1}^T \mathbf{p} = 1$ , then it defines a probability distribution for a random variable  $X$ , for which  $\Pr\{X = \alpha_k\} = p_k$ .
- Therefore, the probability simplex

$$\left\{ \mathbf{p} : \mathbf{p} \succeq \mathbf{0}, \mathbf{1}^T \mathbf{p} = 1 \right\},$$

is in one-to-one correspondence with all possible probability distributions for a random variable  $X$  taking values in  $\{\alpha_1, \dots, \alpha_n\}$ .

Given a combination of prior information and, possibly, observations and measurements, we can often estimate the distribution  $\mathbf{p}$  using convex optimization.

# Examples of Nonparametric Distribution Estimation

## Prior information:

- *Expected values:* As  $E[f(X)] = \sum_{i=1}^n p_i f(\alpha_i) = \mathbf{c}^T \mathbf{p}$ , where  $c_i = f(\alpha_i)$  for  $i = 1, \dots, n$ , we can incorporate expected values of functions as either objectives or constraints.
- *Nonlinear functions of expected values:* In some cases, nonlinear functions of expected values can be used to create convex inequality constraints in  $\mathbf{p}$ . For example,

$$\text{var}(X) = E[X^2] - (E[X])^2 = -\mathbf{p}^T \mathbf{a} \mathbf{a}^T \mathbf{p} + \mathbf{b}^T \mathbf{p}, \text{ where } a_i = \alpha_i^2, b_i = \alpha_i,$$

is concave in  $\mathbf{p}$ , and so constraints of the form  $\text{var}(X) \geq \beta$  are convex. Similarly,

$$\Pr\{X \in \mathcal{A} | X \in \mathcal{B}\} = \Pr\{X \in \mathcal{A} \cap \mathcal{B}\} / \Pr\{X \in \mathcal{B}\} = \mathbf{c}^T \mathbf{p} / \mathbf{d}^T \mathbf{p},$$

where  $c_i = 1$  if  $\alpha_i \in \mathcal{A} \cap \mathcal{B}$  and 0 otherwise and  $d_i = 1$  if  $\alpha_i \in \mathcal{B}$  and 0 otherwise, is linear-fractional in  $\mathbf{p}$ , and so constraints of the form,  $\ell \leq \Pr\{X \in \mathcal{A} | X \in \mathcal{B}\} \leq u$ , are linear inequality constraints,  $\ell \mathbf{d}^T \mathbf{p} \leq \mathbf{c}^T \mathbf{p} \leq u \mathbf{d}^T \mathbf{p}$ , and hence convex.

**Bounding probabilities and expected values:** Given prior information  $\mathbf{p} \in \mathcal{P}$ , where  $\mathcal{P}$  is convex, upper and lower bounds on  $E[f(X)]$  can be found by solving the convex problems,

$$\begin{array}{ll} \text{minimize/maximize} & \sum_{i=1}^n f(\alpha_i) p_i \\ \text{subject to} & \mathbf{p} \in \mathcal{P} \end{array} .$$

**ML estimation:** Suppose we observe  $N$  independent samples  $x_1, \dots, x_N$  from the distribution for  $X$ . Let  $k_i$  denote the number of samples with value  $\alpha_i$ , so that  $k_1 + \dots + k_n = N$ . The log-likelihood function is then  $\ell(\mathbf{p}) = \sum_{i=1}^n k_i \log p_i$ , which is concave in  $\mathbf{p}$ . Thus, the ML estimate of  $\mathbf{p}$  can be found by solving the convex problem

$$\begin{array}{ll} \text{maximize} & \ell(\mathbf{p}) = \sum_{i=1}^n k_i \log p_i \\ \text{subject to} & \mathbf{p} \in \mathcal{P} \end{array} .$$

# Minimum Kullback-Leibler Divergence Problem

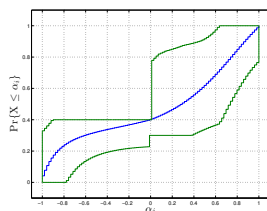
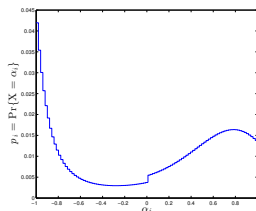
The problem of finding a distribution  $\mathbf{p} \in \mathbb{R}^n$  that has the minimum *Kullback-Leibler divergence* from a given prior distribution  $\mathbf{q} \in \mathbb{R}^n$ , with the prior information that  $\mathbf{p} \in \mathcal{P}$ , leads to the problem

$$\begin{aligned} & \text{minimize} && \sum_{i=1}^n p_i \log(p_i/q_i) \\ & \text{subject to} && \mathbf{p} \in \mathcal{P} \end{aligned}$$

When the set  $\mathcal{P}$  is convex, the problem is convex. In the special case where  $\mathbf{q}$  is the uniform distribution, i.e.,  $\mathbf{q} = (1/n) \mathbf{1}$ , the resulting solution  $\mathbf{p}^*$  is called the *maximum entropy distribution*.

**Example:** Consider 100 equidistant points  $\alpha_i$  in  $[-1, 1]$  with prior information:

$$E[X] \in [-0.1, 0.1], \quad E[X^2] \in [0.5, 0.6], \quad E[3X^3 - 2X] \in [-0.3, -0.2], \quad \Pr\{X < 0\} \in [0.3, 0.4].$$



- Left: maximum entropy distribution that satisfies the prior information constraints.
- Right: bottom and top curves show minimum and maximum possible values of the cumulative distribution function (cdf)  $\Pr\{X \leq \alpha_i\}$ , while the middle curve is the cdf of the maximum entropy distribution.

# Hypothesis Testing and Types of Detectors

## Hypothesis Testing:

Suppose  $\mathbf{x}$  is a random variable with  $n$  values  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  with a distribution that depends on a parameter  $\mathbf{v}$  with  $m$  values  $\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$ . This is characterized by a matrix  $\mathbf{P} \in \mathbb{R}^{n \times m}$  with

$$P_{k,\ell} = \Pr\{\mathbf{x} = \mathbf{x}_k | \mathbf{v} = \mathbf{v}_\ell\}, \quad k = 1, \dots, n, \quad \ell = 1, \dots, m.$$

- The values  $\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$  represent  $m$  *hypotheses*.
- Based on an observation  $x$  of  $\mathbf{x}$ , we wish to guess the correct value  $\mathbf{v}_\ell$ , a process called *hypothesis testing*.

## Types of Detectors:

There are two main types of detectors, namely *deterministic* and *random* detectors.

- For deterministic detectors, if we observe  $\mathbf{x}_k$ , we have  $\hat{\mathbf{v}} = \psi(k)$ . One obvious choice for  $\psi(k)$  is the ML detector, given by

$$\hat{\mathbf{v}} = \psi_{\text{ml}}(k) = \mathbf{v}_{\ell^*(k)} \quad \text{where} \quad \ell^*(k) = \operatorname{argmax}_\ell P_{k,\ell}.$$

- For randomized detectors, the estimate  $\hat{\mathbf{v}}$  is random with a distribution that depends on the observed value of  $\mathbf{x}$ . It is characterized by a matrix  $\mathbf{T} \in \mathbb{R}^{m \times n}$ :

$$T_{i,k} = \Pr\{\hat{\mathbf{v}} = \mathbf{v}_i | \mathbf{x} = \mathbf{x}_k\}, \quad i = 1, \dots, m, \quad k = 1, \dots, n.$$

The columns  $\mathbf{t}_k$  of  $\mathbf{T}$  must satisfy the probability constraints  $\mathbf{t}_k \succeq \mathbf{0}$  and  $\mathbf{1}^T \mathbf{t}_k = 1$ .

# Detection Probability Matrix

For the randomized detector characterized by the matrix  $\mathbf{T}$ , we define the *detection probability matrix* as  $\mathbf{D} \triangleq \mathbf{TP} \in \mathbb{R}^{m \times m}$ . We have

$$D_{i,\ell} = [\mathbf{TP}]_{i,\ell} = \Pr \{ \hat{\mathbf{v}} = \mathbf{v}_i | \mathbf{v} = \mathbf{v}_\ell \}, \quad i = 1, \dots, m, \quad \ell = 1, \dots, m.$$

Hence,  $D_{i,\ell}$  is the probability of guessing  $\hat{\mathbf{v}} = \mathbf{v}_i$ , when in fact  $\mathbf{v} = \mathbf{v}_\ell$  was the true hypothesis.

- *Correct detection or detection probabilities:*

$$P_i^d = \Pr \{ \hat{\mathbf{v}} = \mathbf{v}_i | \mathbf{v} = \mathbf{v}_i \} = D_{i,i}, \quad i = 1, \dots, m.$$

- *Error probabilities:*

$$P_i^e = \Pr \{ \hat{\mathbf{v}} \neq \mathbf{v}_i | \mathbf{v} = \mathbf{v}_i \} = 1 - D_{i,i} = \sum_{\ell \neq i} D_{\ell,i}, \quad i = 1, \dots, m.$$

- If  $\mathbf{D} = \mathbf{I}_m$ , then the detector is perfect: no matter what the parameter  $\mathbf{v}$  is, we correctly guess  $\hat{\mathbf{v}} = \mathbf{v}$ .

# Minimax and Bayes Detector Designs

## Minimax Detector Design:

$$\begin{aligned} & \text{minimize} && \max_{i=1, \dots, m} P_i^e \\ & \text{subject to} && \mathbf{t}_k \succeq \mathbf{0}, \mathbf{1}^T \mathbf{t}_k = 1, k = 1, \dots, n \end{aligned}$$

- The minimax detector minimizes the worst-case probability of error over all  $m$  hypotheses.
- It can be reformulated as an LP.

## Bayes Detector Design:

In Bayes detector design, the hypotheses have a prior distribution given by  $\mathbf{q} \in \mathbb{R}^m$ :

$$q_i = \Pr\{\mathbf{v} = \mathbf{v}_i\}, i = 1, \dots, m.$$

With this, the probability of error of the detector is given by  $\mathbf{q}^T \mathbf{p}^e$ , where  $[\mathbf{p}^e]_i = P_i^e$  for  $i = 1, \dots, m$  is the vector of error probabilities for the hypotheses. The Bayes detector design problem is then as follows.

$$\begin{aligned} & \text{minimize} && \mathbf{q}^T \mathbf{p}^e \\ & \text{subject to} && \mathbf{t}_k \succeq \mathbf{0}, \mathbf{1}^T \mathbf{t}_k = 1, k = 1, \dots, n \end{aligned}$$

- This problem is an LP that has a simple analytical solution.
- For the special case when  $\mathbf{q} = (1/m) \mathbf{1}$ , the Bayes optimal detector minimizes the average probability of error.



# Probability Constraints and Alternate Objectives

## Probability Constraints:

- *Detection probability lower limits:*

$$P_i^d = D_{i,i} \geq L_i, \quad i = 1, \dots, m.$$

- *Error probability upper limits:*

$$D_{k,\ell} \leq U_{k,\ell}, \quad k \neq \ell.$$

**Alternate Objectives or Constraints:** (valid when ordering of hypothesis values  $\mathbf{v} \in \{\mathbf{v}_1, \dots, \mathbf{v}_m\}$  has some significance)

- *Bias:*

$$E_i[(\hat{\mathbf{v}} - \mathbf{v})] = \sum_{k=1}^m (\mathbf{v}_k - \mathbf{v}_i) D_{k,i}, \quad \text{when } \mathbf{v} = \mathbf{v}_i.$$

- *Mean-square error:*

$$E_i[(\hat{\mathbf{v}} - \mathbf{v})^2] = \sum_{k=1}^m (\mathbf{v}_k - \mathbf{v}_i)^2 D_{k,i}, \quad \text{when } \mathbf{v} = \mathbf{v}_i.$$

- *Average absolute error:*

$$E_i[|\hat{\mathbf{v}} - \mathbf{v}|] = \sum_{k=1}^m |\mathbf{v}_k - \mathbf{v}_i| D_{k,i}, \quad \text{when } \mathbf{v} = \mathbf{v}_i.$$

# Formulations of the Optimal Detector Design Problem

## Multicriterion Formulation:

$$\begin{aligned} & \text{minimize (with respect to } \mathbb{R}_+^{m(m-1)} \text{)} && D_{i,\ell}, i, \ell = 1, \dots, m, i \neq \ell \\ & \text{subject to} && \mathbf{t}_k \succeq \mathbf{0}, \mathbf{1}^T \mathbf{t}_k = 1, k = 1, \dots, n \end{aligned}$$

## Scalarization Formulation:

To scalarize the multicriterion problem, we introduce a *loss matrix*  $\mathbf{W} \in \mathbb{R}^{m \times m}$  which satisfies

$$W_{i,i} = 0, i = 1, \dots, m, W_{i,\ell} > 0, i, \ell = 1, \dots, m, i \neq \ell.$$

We then form the weighted sum objective  $\text{tr}(\mathbf{W}^T \mathbf{D})$ . This leads to the scalarization formulation:

$$\begin{aligned} & \text{minimize} && \text{tr}(\mathbf{W}^T \mathbf{D}) \\ & \text{subject to} && \mathbf{t}_k \succeq \mathbf{0}, \mathbf{1}^T \mathbf{t}_k = 1, k = 1, \dots, n \end{aligned}$$

As  $\text{tr}(\mathbf{W}^T \mathbf{D}) = \text{tr}(\mathbf{W}^T \mathbf{T} \mathbf{P}) = \text{tr}(\mathbf{P} \mathbf{W}^T \mathbf{T}) = \sum_{k=1}^n \mathbf{c}_k^T \mathbf{t}_k$ , where  $\mathbf{c}_k$  is the  $k$ -th column of  $\mathbf{W} \mathbf{P}^T$ , this objective is *separable* in  $\mathbf{t}_k$ . Thus, we can solve this problem by separately solving

$$\begin{aligned} & \text{minimize} && \mathbf{c}_k^T \mathbf{t}_k \\ & \text{subject to} && \mathbf{t}_k \succeq \mathbf{0}, \mathbf{1}^T \mathbf{t}_k = 1 \end{aligned}$$

for  $k = 1, \dots, n$ . But this LP has a simple analytic solution: find index  $i$  such that  $c_{k,i} = \min_{\ell} c_{k,\ell}$  and then take  $\mathbf{t}_k^* = \mathbf{e}_i$ . This optimal solution corresponds to a deterministic detector: when  $\mathbf{x} = \mathbf{x}_k$  is observed, our estimate is

$$\hat{\mathbf{v}} = \mathbf{v}_{\ell^*(k)} \text{ where } \ell^*(k) = \underset{\ell}{\text{argmin}} \left[ \mathbf{W} \mathbf{P}^T \right]_{\ell,k}.$$

# MAP and ML Detector Design Problems

## MAP Detector:

For a Bayes detector design with prior distribution  $\mathbf{q}$ , the mean probability of error is

$$\mathbf{q}^T \mathbf{p}^e = \sum_{\ell=1}^m q_{\ell} \sum_{k \neq \ell} D_{k,\ell} = \sum_{k,\ell=1}^m W_{k,\ell} D_{k,\ell},$$

if we define the weight matrix  $\mathbf{W}$  as

$$W_{k,\ell} = q_{\ell}, \quad k, \ell = 1, \dots, m, \quad k \neq \ell, \quad W_{k,k} = 0, \quad k = 1, \dots, m.$$

Thus, we have,

$$\left[ \mathbf{W} \mathbf{P}^T \right]_{\ell,k} = \sum_{i \neq \ell} q_i P_{k,i} = \sum_{i=1}^m q_i P_{k,i} - q_{\ell} P_{k,\ell}.$$

Note that the first term is independent of  $\ell$ . So, when  $\mathbf{x} = \mathbf{x}_k$  is observed, the optimal detector is

$$\hat{\mathbf{v}} = \mathbf{v}_{\ell^*(k)} \quad \text{where } \ell^*(k) = \underset{\ell}{\operatorname{argmax}} (P_{k,\ell} q_{\ell}).$$

Since  $P_{k,\ell} q_{\ell}$  given the probability that  $\mathbf{v} = \mathbf{v}_{\ell}$  and  $\mathbf{x} = \mathbf{x}_k$ , this detector is a MAP detector.

## ML Detector:

For the special case  $\mathbf{q} = (1/m) \mathbf{1}$ , which corresponds to a uniform prior distribution on  $\mathbf{v}$ , this MAP detector reduces to an ML detector:

$$\hat{\mathbf{v}} = \mathbf{v}_{\ell^*(k)} \quad \text{where } \ell^*(k) = \underset{\ell}{\operatorname{argmax}} P_{k,\ell}.$$

Hence, an ML detector minimizes the average or mean probability of error.

# Likelihood-Ratio Test and Neyman-Pearson Lemma

For the special case in which  $m = 2$ , we have *binary hypothesis testing*. In this case,  $\mathbf{x}$  is generated from one of two distributions,  $\mathbf{p} \in \mathbb{R}^n$  or  $\mathbf{q} \in \mathbb{R}^n$ . Also, we have

$$\mathbf{D} = \begin{bmatrix} 1 - P_{\text{fp}} & P_{\text{fn}} \\ P_{\text{fp}} & 1 - P_{\text{fn}} \end{bmatrix}.$$

- $P_{\text{fp}}$  is the probability of a *false positive*, or *false alarm probability*.
- $P_{\text{fn}}$  is the probability of a *false negative*, or *missed detection probability*.
- The optimal trade-off curve between  $P_{\text{fp}}$  and  $P_{\text{fn}}$  is called the *receiver operating characteristic (ROC)*.

For a weight matrix  $\mathbf{W}$ , an optimal detector, assuming  $\mathbf{x} = \mathbf{x}_k$  is observed, is given by

$$\hat{\mathbf{v}} = \begin{cases} \mathbf{v}_1, & p_k/q_k > W_{1,2}/W_{2,1} \\ \mathbf{v}_2, & p_k/q_k \leq W_{1,2}/W_{2,1} \end{cases},$$

which is a *likelihood ratio threshold test* with likelihood ratio  $p_k/q_k$  and threshold  $W_{1,2}/W_{2,1}$ . Choosing different thresholds leads to different Pareto optimal detectors that give different levels of false positive and false negative error probabilities. This result is known as the *Neyman-Pearson lemma*.

# Detector Design Example

Consider a binary hypothesis testing example with  $n = 4$  and

$$\mathbf{P} = \begin{bmatrix} 0.70 & 0.10 \\ 0.20 & 0.10 \\ 0.05 & 0.70 \\ 0.05 & 0.10 \end{bmatrix}$$

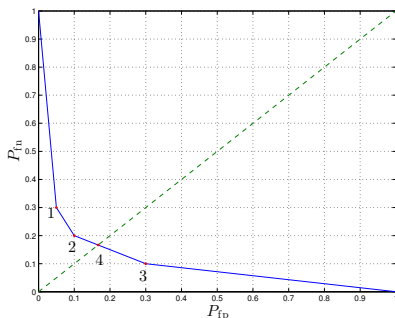
Four Pareto optimal detector matrices are given below.

$$\mathbf{T}^{(1)} = \begin{bmatrix} 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

$$\mathbf{T}^{(2)} = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 \end{bmatrix}$$

$$\mathbf{T}^{(3)} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 \end{bmatrix}$$

$$\mathbf{T}^{(4)} = \begin{bmatrix} 1 & 2/3 & 0 & 0 \\ 0 & 1/3 & 1 & 1 \end{bmatrix}$$



- The vertices labeled 1, 2, and 3 correspond to deterministic detectors.
- The point labeled 4 corresponds to a random detector.
- Here, the dashed line shows  $P_{fn} = P_{fp}$ . The point at which this line intersects with the optimal trade-off curve, namely point 4, corresponds to the minimax detector.

# Worst-Case Detection Probability Matrix

- When  $\mathbf{P}$  is not known exactly, but  $\mathbf{P} \in \mathcal{P}$ , where  $\mathcal{P}$  is a known set of possible distributions, we can attempt to accommodate this uncertainty to develop a *robust* detector.
- Here, error and detection probabilities will be judged by their worst-case values over  $\mathbf{P} \in \mathcal{P}$ .
- This leads us to define a *worst-case detection probability matrix*  $\mathbf{D}^{\text{wc}}$  as follows.

$$D_{k,\ell}^{\text{wc}} = \sup_{\mathbf{P} \in \mathcal{P}} D_{k,\ell}, \quad k, \ell = 1, \dots, m, \quad k \neq \ell, \quad D_{k,k}^{\text{wc}} = \inf_{\mathbf{P} \in \mathcal{P}} D_{k,k}, \quad k = 1, \dots, m.$$

- The off-diagonal entries of  $\mathbf{D}^{\text{wc}}$  give the largest error probabilities, while the diagonal entries give the smallest detection probabilities, over  $\mathbf{P} \in \mathcal{P}$ .
- We define the worst-case probability of error as

$$P_i^{\text{wce}} \triangleq 1 - D_{i,i}^{\text{wc}},$$

which is the largest probability of error, when  $\mathbf{v} = \mathbf{v}_i$ , over all possible distributions in  $\mathcal{P}$ .

- The *robust minimax detector* is defined as the detector that minimizes the worst-case probability of error, over all hypotheses, i.e., minimizes the objective

$$\max_i P_i^{\text{wce}} = \max_{i=1,\dots,m} \sup_{\mathbf{P} \in \mathcal{P}} \left\{ 1 - [\mathbf{TP}]_{i,i} \right\} = 1 - \min_{i=1,\dots,m} \inf_{\mathbf{P} \in \mathcal{P}} \left\{ [\mathbf{TP}]_{i,i} \right\}.$$

# Examples of Robust Minimax Detectors

## Robust minimax detector for finite $\mathcal{P}$ :

When  $\mathcal{P} = \{\mathbf{P}_1, \dots, \mathbf{P}_k\}$ , the robust minimax detector can be found by solving

$$\begin{aligned} & \text{maximize} && \min_{i=1, \dots, m} \inf_{\mathbf{P} \in \mathcal{P}} [\mathbf{TP}]_{i,i} = \min_{i=1, \dots, m} \min_{\ell=1, \dots, k} [\mathbf{TP}_\ell]_{i,i} \\ & \text{subject to} && \mathbf{t}_i \succeq \mathbf{0}, \mathbf{1}^T \mathbf{t}_i = 1, i = 1, \dots, n \end{aligned}$$

- The objective is concave piecewise-linear, so the problem can be posed as an LP.
- The resulting robust minimax detector is also optimal for the polyhedron  $\text{conv}(\mathcal{P})$ .

## Robust minimax detector for polyhedral $\mathcal{P}$

When  $\mathcal{P}$  is a polyhedron of the form

$$\mathcal{P} = \left\{ \mathbf{P} = [ \mathbf{p}_1 \quad \cdots \quad \mathbf{p}_m ] : \mathbf{A}_k \mathbf{p}_k = \mathbf{b}_k, \mathbf{1}^T \mathbf{p}_k = 1, \mathbf{p}_k \succeq \mathbf{0} \right\},$$

the robust minimax detector design problem can be shown to be expressed as the LP

$$\begin{aligned} & \text{maximize} && \gamma \\ & \text{subject to} && \mathbf{b}_i^T \boldsymbol{\nu}_i + \mu_i \geq \gamma, i = 1, \dots, m \\ & && \mathbf{A}_i^T \boldsymbol{\nu}_i + \mu_i \mathbf{1} \preceq \tilde{\mathbf{t}}_i, i = 1, \dots, m \\ & && \mathbf{t}_i \succeq \mathbf{0}, \mathbf{1}^T \mathbf{t}_i = 1, i = 1, \dots, n \end{aligned}$$

with variables  $\boldsymbol{\nu}_1, \dots, \boldsymbol{\nu}_m, \mu_1, \dots, \mu_n$ , and  $\mathbf{T}$  (with columns  $\mathbf{t}_1, \dots, \mathbf{t}_n$  and rows  $\tilde{\mathbf{t}}_1, \dots, \tilde{\mathbf{t}}_m$ ).