



Statistical Signal Processing

Lecture 1

- course arrangements
- course overview
- chapter 1: parameter estimation
 - stochastic parameters
 - the multivariate Gaussian distribution; Gauss-Markov theorem
 - the parameter estimation problem
 - Bayes estimation: the MMSE, absolute value and uniform cost functions
 - examples: Gaussian mean in Gaussian noise, Poisson process



Course Arrangements

- course assistants:

Francesco Negro, Siouar Bensaid

- theoretical part: exam (counts for 14/20), open notes and handouts

- practical part:

- problem sessions (TDs): 2 not graded
- computer sessions (TPs): 1 (Matlab) graded (counts for 3/20)
 \Rightarrow attendance compulsory
- homework (HW): 1 graded (counts for 3/20)

- emphasis (in exam) on problem solving, knowing how to apply the theory
 \Rightarrow TDs important

- course notes and copies of viewgraphs will appear throughout the term, also on:
 web: <http://teaching.eurecom.fr/courses/SSP> or less updated:

http://intranet.eurecom.fr/academicAffairs/cursus_d/courses_content/techcourses/ssp_d.htm



Course Overview

- Chapter 0: Background Material
- Chapter 1: Parameter Estimation
- Chapter 2: Spectrum Estimation
- Chapter 3: Optimal Filtering
- Chapter 4: Adaptive Filtering
- Chapter 5: Estimation of Sinusoids in Noise



Chapter 0: Background Material

See also the MathEng (Mathematical Methods for Engineers) and Optim (Fundamentals of Optimization) courses and a

Linear Algebra review at <http://teaching.eurecom.fr/courses/SSP/>

- probability
- linear algebra
- multivariate Gaussian distribution, Gauss-Markov theorem
- complex Gaussian variables, circularity
- vector spaces, inner products, norms
- optimization

Chapter 1: Parameter Estimation

- **random parameters**, Bayesian estimation, minimum mean squared error (MMSE) estimation, maximum a posteriori (MAP) estimation, equivalent estimators
- MMSE and MAP estimation for vector parameters, Cramer-Rao bound, linear MMSE estimators, orthogonality principle, the linear model
- **deterministic parameters**, uniformly minimum variance unbiased estimators (UMVUE), maximum likelihood (ML) estimation, properties of estimators (bias, efficiency, consistency), Cramer-Rao bound, Estimation-Maximization (EM) algorithm
- best linear unbiased estimator (BLUE), least-squares techniques, method of moments, the linear model
- Advanced: model order selection, case of complex measurements and/or parameters, ML optimization techniques (steepest-descent, Gauss-Newton, alternating opt., scoring), Estimation-Maximization (EM), Variational Bayesian techniques (VB), Compressed Sensing

Chapter 2: Spectrum Estimation

- spectrum estimation = parameter estimation when parameters = spectrum
- non-parametric techniques, periodogram, windowing, spectral leakage, averaged periodogram, smoothed periodogram
- parametric random process models, autoregressive (AR) processes, moving average (MA) processes, autoregressive moving average (ARMA) processes
- parametric techniques, linear prediction, autoregressive modeling, maximum entropy, Levinson & Schur algorithms (relation to triangular covariance matrix factorization), lattice filters
- time and frequency domain localization, short-time Fourier transform, QMF, subbands, perfect reconstruction filter banks, wavelet transform, hierarchical signal representation/approximation

Chapter 3: Optimal Filtering

- optimal filtering = Bayesian parameter estimation when parameters = signal
- **Wiener filtering**: unrealizable (non-causal), causal and FIR; lattice/ladder filters
- application to linear and decision-feedback equalization
- some elements from optimization theory, steepest descent algorithm
- linear state-space models
- **Kalman filtering**
developed for space applications, relation to Levinson/Schur algorithms, Chandrasekhar equations
- applications (channel estimation)

Chapter 4: Adaptive Filtering

- adaptive filtering = optimal filtering in absence of statistical knowledge
- adaptive FIR filtering
- least-mean-square (LMS) algorithm
- recursive least-squares (RLS) algorithm
- tracking of time-varying parameters, performance analysis
- applications

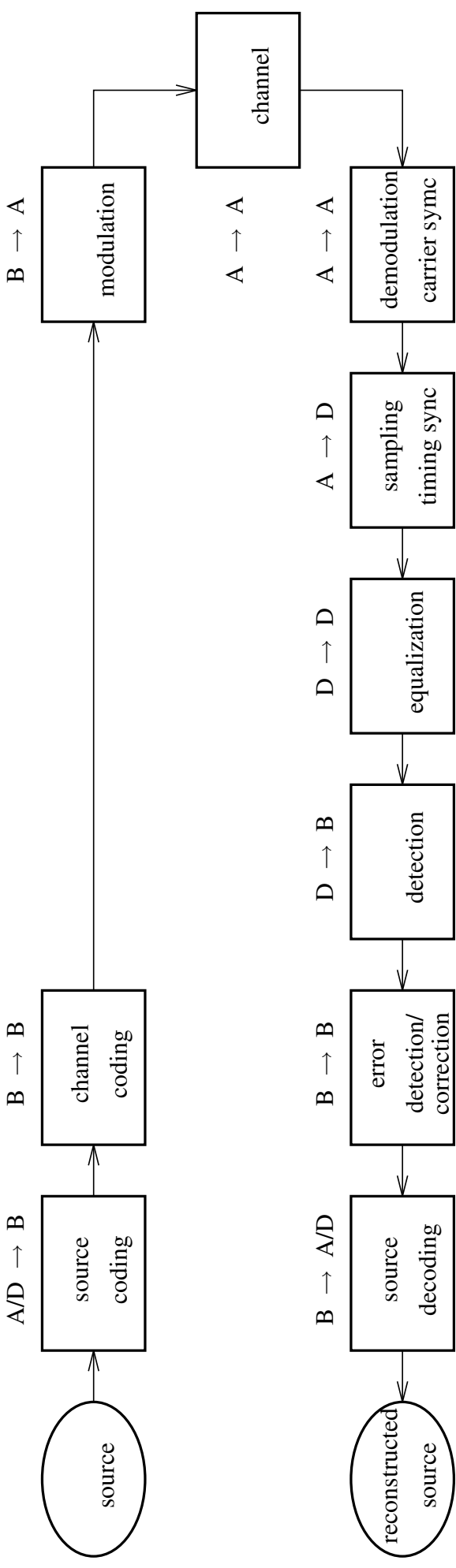
Chapter 5: Estimation of Sinusoids in Noise

- special focus on the estimation of the sinusoid frequencies
- application of deterministic parameter estimation techniques:
 - maximum likelihood: IQML, DIQML, PQML
 - subspace fitting (method of moments)
 - linear prediction: Prony's method, Pisarenko's method
 - linear constrained minimum variance (LCMV): Capon's method
- application of adaptive filtering:
 - adaptive notch filtering

DSP Relevance to Eurecom

- leading thread in DSP course: *statistical signal processing*. Signals as realizations of stochastic processes.
- key building block in telecommunications: *modem*
 - key issue: bit detection: Digital Communications course
 - carrier synchronization: recovery of the carrier frequency and phase
 - timing synchronization: recovery of baud rate and phase
 - channel estimation (or its inverse) for equalization, echo or interference cancellation

MODEM Building Blocks



- modem: modulation-demodulation: digital transmission over analog medium
- signals appear in analog (A), digital (D) or binary/coded (B) form
- DigCom: channel (de)coding, (de)modulation, sampling, detection, equalization
- SSP: source (de)coding, equalization/channel estimation, parameter estimation for timing and carrier synchronization

Applications in Corporate Communications

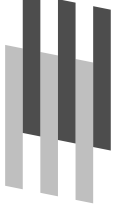
- source compression: (e.g. Unix commands compress and gzip), appetite for information steadily increasing (multimedia)
- hierarchical image compression: nested structure: information at a higher definition level consists of the information at a lower level plus complementary information.
 - result: info split up into streams corresponding to the details at various levels of definition
 - interaction with the network: send different streams separately, with different levels of transmission quality (success of arrival, delay,...)
 - consequences for security: not all streams need to be encrypted (to the same extent)
- communication network performance analysis: queuing theory. parameterized statistical models (e.g. Poisson processes for arrival of packets) need an estimation of their parameter values.

Applications in Multimedia Communications

- lossy source coding for audio and video as opposed to lossless compression of data files
- videoconferencing systems
 - adaptive filters for acoustic echo cancellation (marketing studies: audio more important than video)
 - parametric face models (low bitrate TX via facial cloning), face recognition (biometry), gesture recognition, etc.
- speech recognition for voice-controlled operation: hampered by background noise (other speakers etc.). solution: speech enhancement and directive acoustic microphone arrays
- synchronizaton of multiple information streams (audio & video & ...)
- transmission of continuous sources (speech, video) over packet-mode networks (internet)
- xDSL (Digital Subscriber Loop): high-speed connectivity over the last mile

Applications in Mobile Communications

- digital mobile communications more reliable and flexible than analog systems: requires voice coding; very limited bandwidth \Rightarrow sophisticated low bit rate coding techniques required; changing environment \Rightarrow adaptive multi-rate coding.
- National regulations prohibiting simultaneous driving and calling \Rightarrow handsfree telephone systems \Rightarrow audio conferencing systems and acoustic echos problem. Handsfree talking and dialing: speech recognition techniques in a noisy environment
- multipath propagation in a mobile environment: fast channel equalization techniques needed, possibly of limited complexity, carrier and symbol rate synchronization, Doppler and fading tracking
- user localization (911): position estimation on basis of delay (ToA) and/or direction (DoA) estimates (Time/Direction of Arrival)
- cellular communications and frequency reuse: communications limited, not by noise, but by other users (interferers). Sophisticated antenna array processing permits Spatial Division Multiple Access (SDMA): simultaneous interference reduction and multipath reduction, DoA estimation; spatio-temporal processing; interference reduction in Code DMA (CDMA) systems.



The Multivariate Gaussian Distribution

- $\mathbf{X} = [\mathbf{x}_1 \cdots \mathbf{x}_m]^T$: m jointly Gaussian random variables: $\mathbf{X} \sim \mathcal{N}(m_X, C_{XX})$ with mean

$$\begin{aligned} m_X &= E\mathbf{X} = \int dx_1 \cdots \int dx_m f_{\mathbf{X}}(X) X \\ &= \begin{bmatrix} \int dx_1 \cdots \int dx_m f_{\mathbf{X}}(X) x_1 \\ \vdots \\ \int dx_1 \cdots \int dx_m f_{\mathbf{X}}(X) x_m \end{bmatrix} = \begin{bmatrix} \underbrace{\int dx_1 x_1 \int dx_2 \cdots \int dx_m f_{\mathbf{X}}(X)}_{f_{\mathbf{x}_1}(x_1)} \\ \vdots \\ \underbrace{\int dx_m x_m \int dx_1 \cdots \int dx_{m-1} f_{\mathbf{X}}(X)}_{f_{\mathbf{x}_m}(x_m)} \end{bmatrix} = \begin{bmatrix} m_{x_1} \\ \vdots \\ m_{x_m} \end{bmatrix} \end{aligned}$$

and covariance matrix

$$C_{XX} = E(\mathbf{X} - m_X)(\mathbf{X} - m_X)^T = \int dx_1 \cdots \int dx_m f_{\mathbf{X}}(X) \underbrace{(X - m_X)(X - m_X)^T}_{\geq 0, \text{ rank } 1}$$

where $C_{x_i x_j} = E(\mathbf{x}_i - m_{x_i})(\mathbf{x}_j - m_{x_j})$. As weighted average of positive semidefinite matrices: $C_{XX} = C_{XX}^T \geq 0$ symmetric and positive semidefinite:
 $\forall U \in \mathcal{R}^m$: $U^T C_{XX} U = U^T E(\mathbf{X} - m_X)(\mathbf{X} - m_X)^T U = E(U^T (\mathbf{X} - m_X))^2 \geq 0$

- joint Gaussian probability density function (pdf):

$$f_{\mathbf{X}}(X) = (2\pi)^{-\frac{m}{2}} (\det C_{XX})^{-\frac{1}{2}} \exp\left(-\frac{1}{2} [X - m_X]^T C_{XX}^{-1} [X - m_X]\right)$$

Multivariate Gaussian Derivation

- goal: derive multivariate Gaussian distribution from univariate Gaussian distribution and two postulates
- $\mathbf{X} = [\mathbf{x}_1 \cdots \mathbf{x}_m]^T$ real random vector with specified mean and covariance matrix: with mean $m_X = E\mathbf{X} = [m_{x_1} \cdots m_{x_m}]^T$ and covariance matrix

$$\begin{aligned} C_{XX} &= E(\mathbf{X} - m_X)(\mathbf{X} - m_X)^T = \left[E(\mathbf{x}_i - m_{x_i})(\mathbf{x}_j - m_{x_j})^T \right]_{i,j=1}^m \\ &= E(\mathbf{X}\mathbf{X}^T - \mathbf{X}m_X^T - m_X\mathbf{X}^T + m_Xm_X^T) \\ &= (E\mathbf{X}\mathbf{X}^T) - (E\mathbf{X})m_X^T - m_X(E\mathbf{X})^T + m_Xm_X^T \\ &= R_{XX} - m_Xm_X^T - m_Xm_X^T + m_Xm_X^T = R_{XX} - m_Xm_X^T \end{aligned}$$

$R_{XX} = E\mathbf{X}\mathbf{X}^T$ = correlation matrix,

linearity of expectation E exploited (linear operations commute)

- by definition of expectation

$$\begin{aligned} C_{XX} &= E(\mathbf{X} - m_X)(\mathbf{X} - m_X)^T = \int d\mathbf{x}_1 \cdots \int d\mathbf{x}_m f_{\mathbf{X}}(\mathbf{X})(\mathbf{X} - m_X)(\mathbf{X} - m_X)^T \\ \Rightarrow C_{XX} &= C_{XX}^T \text{ symmetric} \end{aligned}$$

Eigendecomposition Covariance Matrix

- *eigenvalues* λ_i and corresponding *eigenvectors* V_i of $C_{XX} : C_{XX} V_i = \lambda_i V_i$
fix *norm* $\|V_i\| = 1, \quad \|V_i\|^2 = V_i^T V_i$
- $(C_{XX} - \lambda_i I_m) V_i = 0 \Rightarrow (C_{XX} - \lambda_i I_m)$ *singular*
 λ_i solution of $\det(C_{XX} - \lambda I_m) = 0 : \text{characteristic equation}$
- $C_{XX} = C_{XX}^T \Rightarrow \lambda_i \in \mathcal{R}, \quad V_i^T V_j = \delta_{ij}, \quad i, j = 1, \dots, n$

$$\text{Kronecker delta} : \delta_{ij} = \begin{cases} 1, & i = j \\ 0, & i \neq j \end{cases}$$

- matrix $V = [V_1 \cdots V_m] \quad (m \times m) :$
 $[V^T V]_{ij} = V_i^T V_j = \delta_{ij} \Rightarrow V^T V = I_m \quad V = \text{orthogonal matrix}$
- C_{XX} real, λ_i real $\Rightarrow V_i$ can be chosen to be real
- $I_m = V^T V \Rightarrow 1 = \det(I_m) = \det(V^T V) = \det(V^T) \det(V) = (\det V)^2$
 $\Rightarrow \det V = \pm 1$. We can choose the signs of the V_i such that $\det V = 1$.

Eigendecomposition Covariance Matrix (2)

- $C_{XX} \geq 0$ positive semidefinite:

$$\forall U \in \mathcal{R}^m : U^T C_{XX} U = E \left(U^T (\mathbf{X} - m_X) \right)^2 \geq 0$$

(positive definite would be: $\forall U \in \mathcal{R}^m \setminus \{0\} : U^T C_{XX} U > 0$)

- $U = V_i : U^T C_{XX} U = V_i^T C_{XX} V_i = \lambda_i V_i^T V_i = \lambda_i \geq 0$
- order the $\lambda_i : \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m \geq 0$
- If $\lambda_m = 0$, C_{XX} is singular $\Rightarrow V_m^T C_{XX} V_m = E \left(V_m^T (\mathbf{X} - m_X) \right)^2 = 0$ mean and variance of $V_m^T (\mathbf{X} - m_X)$ are zero $\Rightarrow V_m^T (\mathbf{X} - m_X) = 0$ in mean square. This means that at least one variable \mathbf{x}_i is a linear combination of the other variables and 1. We shall in general exclude this possibility $\Rightarrow C_{XX} > 0$, $\lambda_i > 0$, $i = 1, \dots, m$

- $C_{XX} V_i - V_i \lambda_i = 0$ are the columns of the matrix $C_{XX} V - V \Lambda = 0$ where $\Lambda = \text{diag} \{ \lambda_1, \dots, \lambda_m \}$. Using $V^{-1} = V^T$, we find

$$C_{XX} = V \Lambda V^T = [V_1 \dots V_m] \text{diag} \{ \lambda_1, \dots, \lambda_m \} [V_1 \dots V_m]^T = \sum_{i=1}^m \lambda_i V_i V_i^T$$

Can show: $C_{XX}^\alpha = V \Lambda^\alpha V^T$ for any α (e.g. $\alpha = 2, -1$)

Multivariate Gaussian Derivation (2)

- Consider now a linear transformation of variables: $\mathbf{Z} = V^T(\mathbf{X} - m_X)$.

Then we find for the first two moments

$$\begin{aligned} m_Z &= E V^T(\mathbf{X} - m_X) = V^T(E\mathbf{X} - m_X) = V^T(m_X - m_X) = 0 \\ C_{ZZ} &= E(\mathbf{Z} - m_Z)(\mathbf{Z} - m_Z)^T = E\mathbf{Z}\mathbf{Z}^T = EV^T(\mathbf{X} - m_X)(\mathbf{X} - m_X)^T V \\ &= V^T(E(\mathbf{X} - m_X)(\mathbf{X} - m_X)^T)V = V^T C_{XX} V = V^T V \Lambda V^T V = \Lambda \end{aligned}$$

Or hence $E\mathbf{z}_i\mathbf{z}_j = \lambda_i\delta_{ij}$: the \mathbf{z}_i are zero mean and uncorrelated.

- At this point: only specified the first two moments of the \mathbf{z}_i . Now specify the rest of their distribution by stating that the \mathbf{z}_i are jointly Gaussian. We furthermore *postulate* that zero mean uncorrelated Gaussian random variables are independent.

$$\left. \begin{array}{l} \text{independent} \\ \text{zero mean} \end{array} \right\} \begin{array}{l} \Rightarrow \\ \Leftarrow \end{array} \left\{ \begin{array}{l} \text{uncorrelated} \\ \text{zero mean} \end{array} \right.$$

Note that in general

Multivariate Gaussian Derivation (3)

- joint distribution of the independent Gaussian r.v.'s \mathbf{z}_i :

$$f_{\mathbf{Z}}(\mathbf{Z}) = \prod_{i=1}^m f_{\mathbf{z}_i}(z_i) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi} \lambda_i} \exp\left(-\frac{z_i^2}{2\lambda_i}\right) = (2\pi)^{-\frac{m}{2}} (\det \Lambda)^{-\frac{1}{2}} \exp\left(-\frac{1}{2} \mathbf{Z}^T \Lambda^{-1} \mathbf{Z}\right)$$

- furthermore *postulate* that a linear transformation of jointly Gaussian random variables produces again jointly Gaussian random variables.
- Since the Jacobian $(\det V^T)$ of the linear transformation between \mathbf{X} and \mathbf{Z} equals one, we get for the joint distribution of the \mathbf{x}_i $\mathbf{Z} = V^T (\mathbf{X} - m_X)$

$$\begin{aligned} f_{\mathbf{X}}(\mathbf{X}) &= f_{\mathbf{Z}}(V^T (\mathbf{X} - m_X)) \\ &= (2\pi)^{-\frac{m}{2}} (\det(V^T C_{XX} V))^{-\frac{1}{2}} \exp\left(-\frac{1}{2} [V^T (\mathbf{X} - m_X)]^T \Lambda^{-1} [V^T (\mathbf{X} - m_X)]\right) \\ &= (2\pi)^{-\frac{m}{2}} ((\det V^T)(\det C_{XX})(\det V))^{-\frac{1}{2}} \exp\left(-\frac{1}{2} [\mathbf{X} - m_X]^T V \Lambda^{-1} V^T [\mathbf{X} - m_X]\right) \\ &= (2\pi)^{-\frac{m}{2}} (\det C_{XX})^{-\frac{1}{2}} \exp\left(-\frac{1}{2} [\mathbf{X} - m_X]^T C_{XX}^{-1} [\mathbf{X} - m_X]\right) \end{aligned}$$

This is the general expression for a multivariate Gaussian probability density function (pdf). Notation: $\mathbf{X} \sim \mathcal{N}(m_X, C_{XX})$: completely specified in terms of the first and second-order moments.

The Multivariate Gaussian Distribution (2)

- $\mathbf{X} = [\mathbf{x}_1 \cdots \mathbf{x}_m]^T$ and $\mathbf{Y} = [y_1 \cdots y_n]^T$: $m+n$ jointly Gaussian random variables:

$$\begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} m_X \\ m_Y \end{bmatrix}, \begin{bmatrix} C_{XX} & C_{XY} \\ C_{YX} & C_{YY} \end{bmatrix} \right)$$

with mean and covariance matrix

$$\begin{bmatrix} m_X \\ m_Y \end{bmatrix} = E \begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix}, \quad C = \begin{bmatrix} C_{XX} & C_{XY} \\ C_{YX} & C_{YY} \end{bmatrix} = E \begin{bmatrix} \mathbf{X} - m_X \\ \mathbf{Y} - m_Y \end{bmatrix} \begin{bmatrix} \mathbf{X} - m_X \\ \mathbf{Y} - m_Y \end{bmatrix}^T$$

- joint Gaussian probability density function (pdf): $f_{\mathbf{X}, \mathbf{Y}}(X, Y)$

$$= (2\pi)^{-\frac{m+n}{2}} (\det C)^{-\frac{1}{2}} \exp \left(-\frac{1}{2} \begin{bmatrix} X - m_X \\ Y - m_Y \end{bmatrix}^T \begin{bmatrix} C_{XX} & C_{XY} \\ C_{YX} & C_{YY} \end{bmatrix}^{-1} \begin{bmatrix} X - m_X \\ Y - m_Y \end{bmatrix} \right)$$
- meaning of the joint pdf:

$$\Pr[\mathbf{X} \in (X, X + dX), \mathbf{Y} \in (Y, Y + dY)] = f_{\mathbf{X}, \mathbf{Y}}(X, Y) dX dY$$

where first $dX = [dx_1 \cdots dx_m]^T$ and then $dX = dx_1 dx_2 \cdots dx_m$

The Conditional Gaussian Distribution

- the conditional pdf $f_{\mathbf{X}|\mathbf{Y}}(X|Y)$ means

$$\Pr[\mathbf{X} \in (X, X + dX) | \mathbf{Y} = Y] = f_{\mathbf{X}|\mathbf{Y}}(X|Y) dX$$

- the conditional pdf is defined by Bayes' rule:

$$f_{\mathbf{X}|\mathbf{Y}}(X|Y) = \frac{f_{\mathbf{X},\mathbf{Y}}(X, Y)}{f_{\mathbf{Y}}(Y)} \Leftrightarrow f_{\mathbf{X},\mathbf{Y}}(X, Y) = f_{\mathbf{Y}}(Y) f_{\mathbf{X}|\mathbf{Y}}(X|Y)$$

get the marginal pdf $f_{\mathbf{Y}}(Y)$ from the joint pdf by integrating out X :

$$\begin{aligned} f_{\mathbf{Y}}(Y) &= \int f_{\mathbf{X},\mathbf{Y}}(X, Y) dX \\ &= \int \cdots \int f_{\mathbf{x}_1, \dots, \mathbf{x}_m, y_1, \dots, y_n}(\mathbf{x}_1, \dots, \mathbf{x}_m, y_1, \dots, y_n) d\mathbf{x}_1 \cdots d\mathbf{x}_m \\ &= (2\pi)^{-\frac{n}{2}} (\det C_{YY})^{-\frac{1}{2}} \exp\left(-\frac{1}{2} [Y - m_Y]^T C_{YY}^{-1} [Y - m_Y]\right) \end{aligned}$$

- consider the block Upper Diagonal Lower (UDL) triangular factorization of C :

$$C = \begin{bmatrix} C_{XX} & C_{XY} \\ C_{YX} & C_{YY} \end{bmatrix} = \begin{bmatrix} I & K \\ 0 & I \end{bmatrix} \begin{bmatrix} P & 0 \\ 0 & C_{YY} \end{bmatrix} \begin{bmatrix} I & 0 \\ K^T & I \end{bmatrix}$$

with $K = C_{XY}C_{YY}^{-1}$, $P = C_{XX} - C_{XY}C_{YY}^{-1}C_{YX}$ (*Schur complement*)

The Conditional Gaussian Distribution (2)

- from the UDL factorization of C , we can obtain the LDU factorization of C^{-1}

$$\begin{aligned} C^{-1} &= \begin{bmatrix} I & 0 \\ K^T & I \end{bmatrix}^{-1} \begin{bmatrix} P & 0 \\ 0 & C_{YY} \end{bmatrix}^{-1} \begin{bmatrix} I & K \\ 0 & I \end{bmatrix}^{-1} \\ &= \begin{bmatrix} I & 0 \\ -K^T & I \end{bmatrix} \begin{bmatrix} P^{-1} & 0 \\ 0 & C_{YY}^{-1} \end{bmatrix} \begin{bmatrix} I & -K \\ 0 & I \end{bmatrix} \end{aligned}$$

- we can rewrite the exponent of the joint distribution as

$$\begin{aligned} &\begin{bmatrix} X - m_X \\ Y - m_Y \end{bmatrix}^T \begin{bmatrix} C_{XX} & C_{XY} \\ C_{YX} & C_{YY} \end{bmatrix}^{-1} \begin{bmatrix} X - m_X \\ Y - m_Y \end{bmatrix} \\ &= \begin{bmatrix} X - m_X \\ Y - m_Y \end{bmatrix}^T \begin{bmatrix} I & 0 \\ -K^T & I \end{bmatrix} \begin{bmatrix} P^{-1} & 0 \\ 0 & C_{YY}^{-1} \end{bmatrix} \begin{bmatrix} I & -K \\ 0 & I \end{bmatrix} \begin{bmatrix} X - m_X \\ Y - m_Y \end{bmatrix} \\ &= \begin{bmatrix} X - KY - (m_X - Km_Y) \\ Y - m_Y \end{bmatrix}^T \begin{bmatrix} P^{-1} & 0 \\ 0 & C_{YY}^{-1} \end{bmatrix} \begin{bmatrix} X - KY - (m_X - Km_Y) \\ Y - m_Y \end{bmatrix} \\ &= [X - KY - (m_X - Km_Y)]^T P^{-1} [X - KY - (m_X - Km_Y)] \\ &\quad + [Y - m_Y]^T C_{YY}^{-1} [Y - m_Y] \end{aligned}$$

The Gauss-Markov Theorem

- By taking determinants, we also obtain

$$\begin{aligned}\det C &= \det \begin{bmatrix} I & K \\ 0 & I \end{bmatrix} \det \begin{bmatrix} P & 0 \\ 0 & C_{YY} \end{bmatrix} \det \begin{bmatrix} I & 0 \\ K^T & I \end{bmatrix} \\ &= 1 \cdot \det \begin{bmatrix} P & 0 \\ 0 & C_{YY} \end{bmatrix} \cdot 1 = \det P \det C_{YY}\end{aligned}$$

- **Theorem 0.1 (Gauss-Markov)** *If \mathbf{X} and \mathbf{Y} have the joint Gaussian distribution indicated before, then the conditional distribution is*

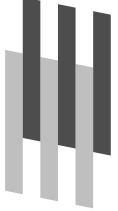
$$\begin{aligned}f_{\mathbf{X}|\mathbf{Y}}(X|Y) &= (2\pi)^{-\frac{m}{2}} (\det P)^{-\frac{1}{2}} \\ &\exp \left(-\frac{1}{2} [X - KY - (m_X - Km_Y)]^T P^{-1} [X - KY - (m_X - Km_Y)] \right)\end{aligned}$$

The conditional distribution is again Gaussian with conditional mean

$$E_{\mathbf{X}|\mathbf{Y}} \mathbf{X} = E[\mathbf{X}|\mathbf{Y} = Y] = m_X + \underbrace{C_{XY}C_{YY}^{-1}}_{=K} (Y - m_Y)$$

and conditional covariance matrix

$$E_{\mathbf{X}|\mathbf{Y}} [\mathbf{X} - E_{\mathbf{X}|\mathbf{Y}} \mathbf{X}] [\mathbf{X} - E_{\mathbf{X}|\mathbf{Y}} \mathbf{X}]^T = P = C_{XX} - C_{XY}C_{YY}^{-1}C_{YX}$$



Example 1.1 Two correlated Gaussian r.v.'s

- $m = n = 1$, zero means $m_X = m_Y = 0$, rename $C_{XX} = \sigma_x^2$, $C_{YY} = \sigma_y^2$.
- Cauchy-Schwarz inequality: $|\langle \mathbf{x}, \mathbf{y} \rangle|^2 \leq \langle \mathbf{x}, \mathbf{x} \rangle \langle \mathbf{y}, \mathbf{y} \rangle$
- $\langle \mathbf{x}, \mathbf{y} \rangle = E\mathbf{xy} \Rightarrow C_{XY}^2 \leq C_{XX}C_{YY}$
- isomorphism: $\mathbf{x} = \underline{x}^T \mathbf{e}$, $\mathbf{y} = \underline{y}^T \mathbf{e}$, $R_{\mathbf{ee}} = I$, $\Rightarrow E\mathbf{xy} = \langle \mathbf{x}, \mathbf{y} \rangle = \langle \underline{x}, \underline{y} \rangle = \underline{x}^T \underline{y}$
- $E(\mathbf{x} + \lambda \mathbf{y})^2 = \lambda^2 r_{yy} + 2\lambda r_{xy} + r_{xx} \stackrel{\lambda = -r_{xy}/r_{yy}}{\geq} r_{xx} - \frac{r_{xy}^2}{r_{yy}} \geq 0 \Rightarrow r_{xy}^2 \leq r_{xx}r_{yy}$
- introduce normalized correlation coefficient $\rho = \frac{C_{XY}}{\sqrt{C_{XX}C_{YY}}} \in [-1, 1]$
can rewrite C_{XY} as $C_{XY} = \rho\sigma_x\sigma_y$.
- The joint Gaussian distribution of \mathbf{x} and \mathbf{y} can be written as

$$f_{\mathbf{x}, \mathbf{y}}(x, y) \leftrightarrow \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} \sigma_x^2 & \rho\sigma_x\sigma_y \\ \rho\sigma_x\sigma_y & \sigma_y^2 \end{bmatrix}\right)$$

- The conditional pdf of \mathbf{x} given \mathbf{y} can be written as

$$f_{\mathbf{x}|\mathbf{y}}(x|\mathbf{y}) \leftrightarrow \mathcal{N}\left(\rho \frac{\sigma_x}{\sigma_y} \mathbf{y}, \sigma_x^2(1-\rho^2)\right)$$



Example 1.1 Two correlated Gaussian r.v.'s (2)

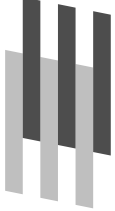
- if \mathbf{x}, y uncorrelated ($\rho = 0$), then $f_{\mathbf{x}|y}(x|y) = f_{\mathbf{x}}(x)$: \mathbf{x}, y independent
- as $|\rho| \rightarrow 1, |E_{\mathbf{x}|y}\mathbf{x}| \nearrow$, conditional variance $\rightarrow 0$:
when $y = y$ is known, it gives us some information about \mathbf{x} and the residual randomness in \mathbf{x} decreases as $|\rho| \rightarrow 1$.

- Extreme cases:

$$\begin{aligned}\rho = 1 : \quad \mathbf{x} &= \frac{\sigma_x}{\sigma_y} \mathbf{y} & \mathbf{x} \text{ and } y \text{ are perfectly correlated,} \\ \rho = -1 : \quad \mathbf{x} &= -\frac{\sigma_x}{\sigma_y} \mathbf{y} & \mathbf{x} \text{ and } y \text{ are perfectly anticorrelated.}\end{aligned}$$

- *concentration ellipse*:

$$\begin{aligned}f_{\mathbf{x},y}(x, y) = c' &\iff \begin{bmatrix} x \\ y \end{bmatrix}^T \begin{bmatrix} \sigma_x^2 & \rho\sigma_x\sigma_y \\ \rho\sigma_x\sigma_y & \sigma_y^2 \end{bmatrix}^{-1} \begin{bmatrix} x \\ y \end{bmatrix} = c \\ &\iff \frac{\left(x - \rho\frac{\sigma_x}{\sigma_y}y\right)^2}{\sigma_x^2(1-\rho^2)} + \frac{y^2}{\sigma_y^2} = c\end{aligned}$$



Example 1.1 Two correlated Gaussian r.v.'s (3)

- the ellipse for $c = 1$

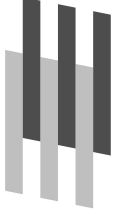
$$\left(\frac{x - \rho \frac{\sigma_x}{\sigma_y} y}{\sigma_x^2 (1 - \rho^2)} \right)^2 + \frac{y^2}{\sigma_y^2} \leq 1$$

contains a significant portion of the probability mass.

- its volume is

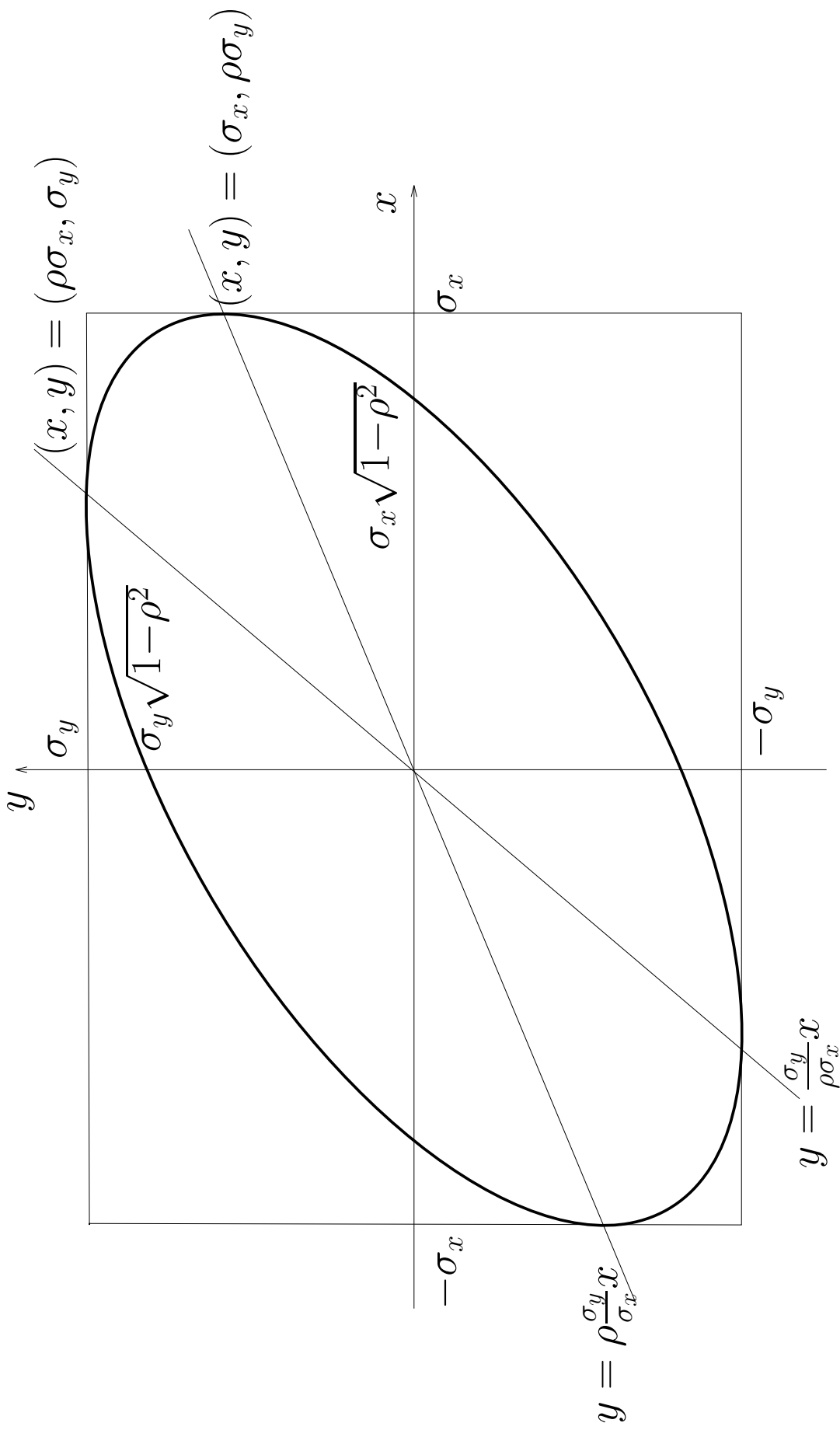
$$V = \pi (\det C)^{\frac{1}{2}} = \pi \sigma_x \sigma_y \sqrt{1 - \rho^2} \in [0, \pi \sigma_x \sigma_y]$$

- for strongly correlated variables, the pair (x, y) takes with high probability values in a fairly small area. Hence, in some sense the randomness (or entropy) of the pair (x, y) decreases as $|\rho| \rightarrow 1$. (later: Gaussian r.v.'s: entropy $\sim \det C$)

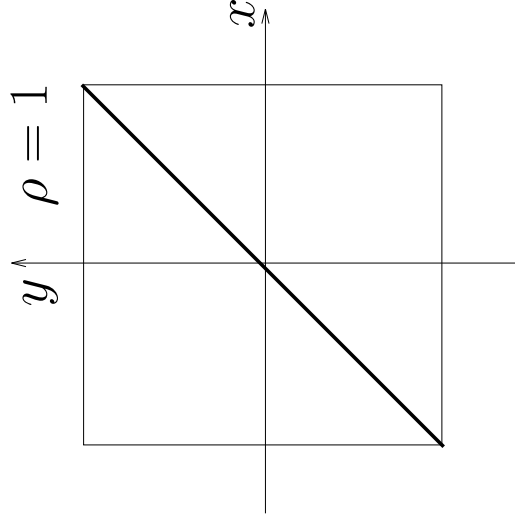
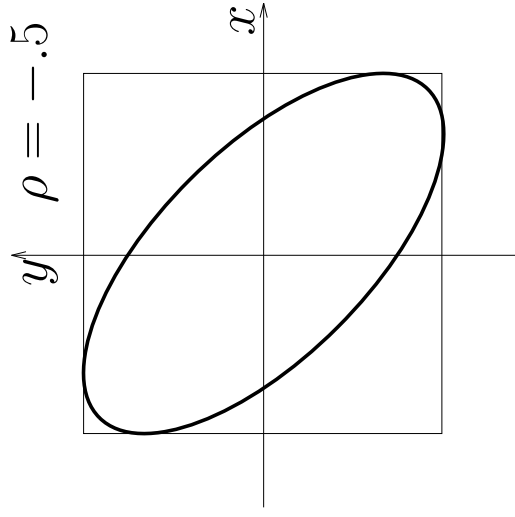
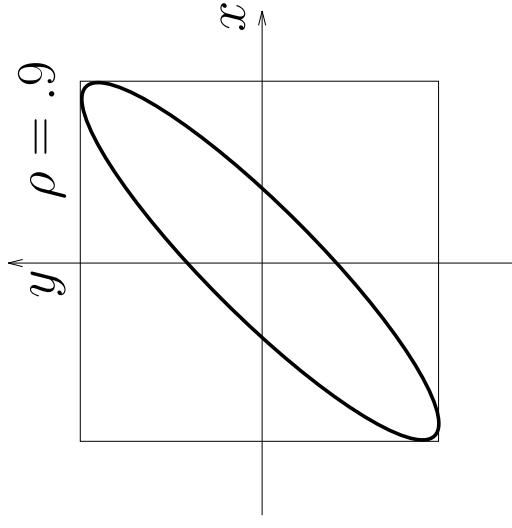
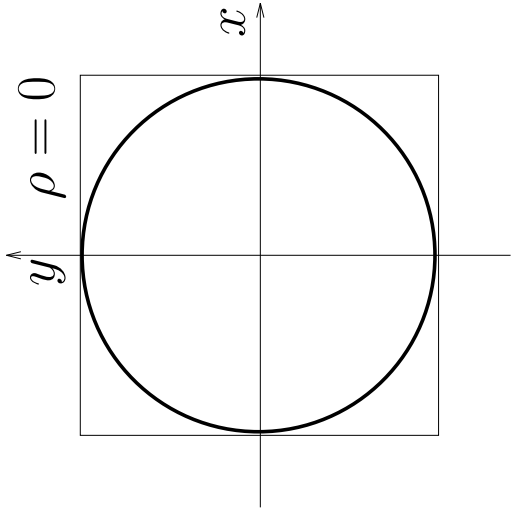


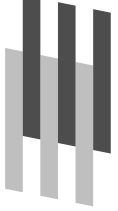
Example 1.1 Two correlated Gaussian r.v.'s (4)

concentration ellipses



Example 1.1 Two correlated Gaussian r.v.'s (5)





Gaussian r.v.'s and Linear Models

- by interchanging X and Y , we find the conditional pdf of Y given X

$$f_{Y|X}(Y|X) \leftrightarrow \mathcal{N}(m_Y + C_{YX}C_{XX}^{-1}(X - m_X), C_{YY} - C_{YX}C_{XX}^{-1}C_{XY})$$

- introduce the Gaussian r. vector V of dimension n also with distribution

$$f_V(V) \leftrightarrow \mathcal{N}(\underbrace{m_Y - C_{YX}C_{XX}^{-1}m_X}_{=m_V}, \underbrace{C_{YY} - C_{YX}C_{XX}^{-1}C_{XY}}_{=C_{VV}})$$

and V independent of X ,

- then X and Y generated as follows

$$\begin{aligned} \begin{bmatrix} X \\ Y \end{bmatrix} &= \begin{bmatrix} I & 0 \\ C_{YX}C_{XX}^{-1} & I \end{bmatrix} \begin{bmatrix} X \\ V \end{bmatrix} \Rightarrow \begin{bmatrix} I & 0 \\ C_{YX}C_{XX}^{-1} & I \end{bmatrix} \begin{bmatrix} m_X \\ m_V \end{bmatrix} = \begin{bmatrix} m_X \\ m_Y \end{bmatrix} \\ \begin{bmatrix} I & 0 \\ C_{YX}C_{XX}^{-1} & I \end{bmatrix} \begin{bmatrix} C_{XX} & \underbrace{C_{XV}}_{=0} \\ \underbrace{C_{VX}}_{=0} & \underbrace{C_{VV}}_{=C_{YY}-C_{YX}C_{XX}^{-1}C_{XY}} \end{bmatrix} \begin{bmatrix} I & C_{XX}^{-1}C_{XY} \\ 0 & I \end{bmatrix} = \begin{bmatrix} C_{XX} & C_{XY} \\ C_{YX} & C_{YY} \end{bmatrix} \end{aligned}$$

(block Lower-Diagonal-Upper (LDU) triangular factorization of C)

are jointly Gaussian and have the correct mean and variance, hence have the correct pdf $f_{X,Y}(X, Y)$.



Gaussian r.v.'s and Linear Models (2)

- we can write explicitly

$$\mathbf{Y} = C_{YX}C_{XX}^{-1}\mathbf{X} + \mathbf{V} .$$

This means that for two sets of correlated Gaussian random variables, one (\mathbf{Y}) can be thought of as being generated from the other (\mathbf{X}) through a linear model ($C_{YX}C_{XX}^{-1}$) and being corrupted by independent Gaussian “measurement” noise (\mathbf{V}).

- Going back to the scalar example ($m = n = 1$), the variance of the linear model part ($C_{YX}C_{XX}^{-1}\mathbf{X}$) is $\rho^2\sigma_y^2$ while the variance of the measurement noise (\mathbf{V}) is $(1-\rho^2)\sigma_y^2$. The two are complementary parts of the total variance σ_y^2 of \mathbf{Y} . We could define a signal to noise ratio (SNR) as the ratio of the two parts which is $\frac{\rho^2}{1-\rho^2}$. The SNR increases from 0 to ∞ as $|\rho|$ increases from 0 to 1, as the correlation between \mathbf{X} and \mathbf{Y} increases.



The Parameter Estimation Problem

description of stochastic processes known up to some parameters

- tone detector: a push button telephone emits sinusoids of which the frequencies are distinct for different buttons. The unknown parameter of interest is the sinusoid frequency. This frequency determination problem may in fact be better approached as a detection problem as we shall see shortly
- the carrier phase and timing instants in linear digital modulation schemes
- the impulse response of the transmission channel. We may take a parameterized model (e.g. FIR model) for this impulse response and then the channel identification problem becomes a problem of estimating the parameters in its model
- similar impulse response identification problems occur in the problems of the cancellation of electrical echos caused by hybrids connecting 2-wire and 4-wire sections of telephone line, and acoustic echos in teleconferencing systems

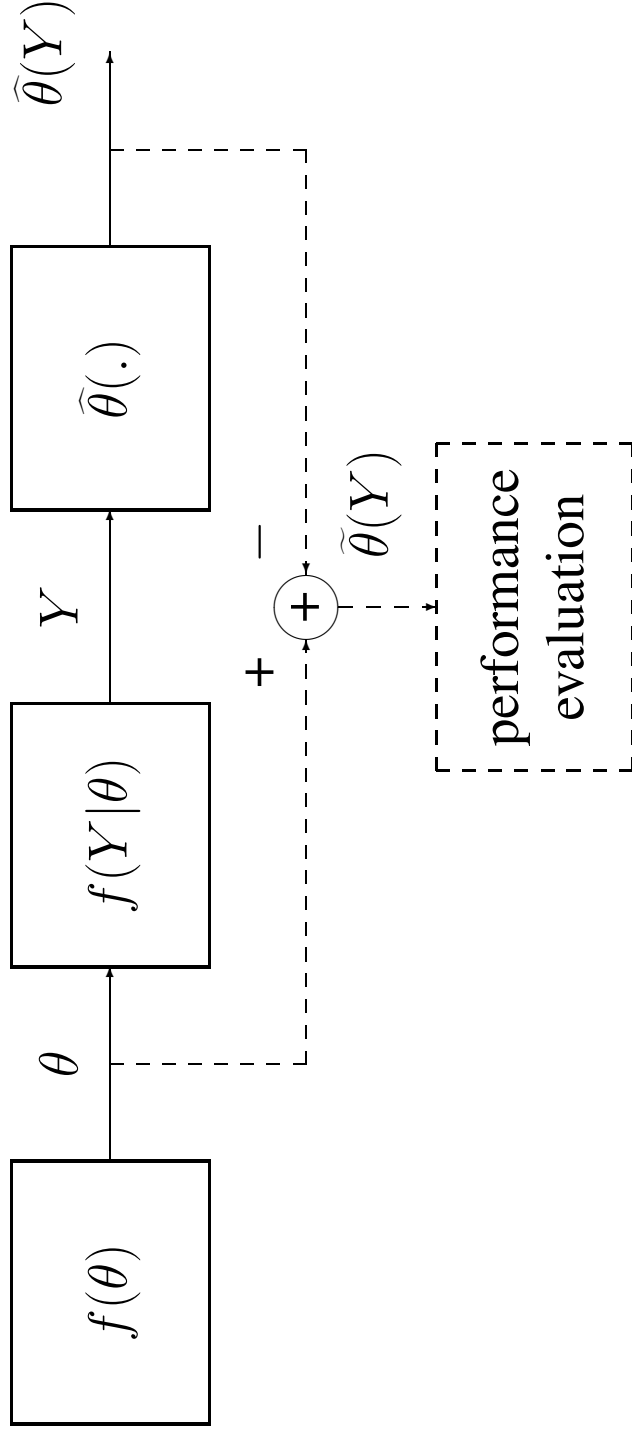
The Parameter Estimation Problem (2)

description of stochastic processes known up to some parameters

- the average rate in a Poisson process occurring in queuing theory and network performance analysis
- the pixel value in the noisy capture of an image by a satellite
- the position and orientation of an object in an image
- security in mobile communications: the (parameterized) distribution of features characterizing the behavior of users. The deviation from typical habits may be a reason for alarm.

Bayesian Framework

- Bayesian approach: consider parameters θ to be random variables
- the parameters have some *a priori* distribution $f_{\theta}(\theta)$. This distribution is called a priori because it summarizes the knowledge we have about θ before making any measurement.
- Next we make a *measurement* Y . The measurement is not a deterministic function of the unknown parameters. The random aspect of this relation is captured by the conditional distribution $f_{Y|\theta}(Y|\theta)$.



Ex 1.2: Additive indep. measurement noise

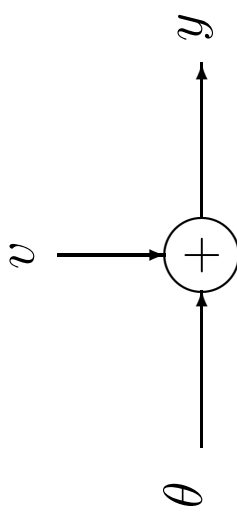
$$\bullet f_{y|\boldsymbol{\theta}}(y|\boldsymbol{\theta}) = \frac{f_{y,\boldsymbol{\theta}}(y,\boldsymbol{\theta})}{f_{\boldsymbol{\theta}}(\boldsymbol{\theta})} = \frac{f_{v,\boldsymbol{\theta}}(y-\boldsymbol{\theta},\boldsymbol{\theta})}{f_{\boldsymbol{\theta}}(\boldsymbol{\theta})} = \frac{f_v(y-\boldsymbol{\theta}) f_{\boldsymbol{\theta}}(\boldsymbol{\theta})}{f_{\boldsymbol{\theta}}(\boldsymbol{\theta})} = f_v(y-\boldsymbol{\theta})$$

- first identity: Bayes' rule
- second identity: transformation of variables

$$\begin{bmatrix} v \\ \theta \end{bmatrix} = \begin{bmatrix} 1 & -1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} y \\ \theta \end{bmatrix}$$

linear transformation \Rightarrow Jacobian = determinant of the transformation matrix, which equals 1.

- third identity: independence of v and θ



Estimators

- *estimator* : a function $\widehat{\theta}(\cdot)$ to be applied to the measurement \mathbf{Y}
- *estimate* : $\widehat{\theta}(Y)$ is the estimator $\widehat{\theta}(\mathbf{Y})$ evaluated at $\mathbf{Y} = Y$.
- we consider here a *point estimator*, i.e. it delivers one value $\widehat{\theta} = \widehat{\theta}(Y) = t(y_1, \dots, y_n)$ that we hope to be close to θ .
- The function $\widehat{\theta}(\mathbf{Y}) = t(\mathbf{y}_1, \dots, \mathbf{y}_n)$ of the measurement data is called a *statistic*.
- Another type of estimator would be an *interval estimator*: two statistics $t_1(y_1, \dots, y_n)$ and $t_2(y_1, \dots, y_n)$ are used to define an interval such that

$$\Pr \{ \theta \in (t_1(y_1, \dots, y_n), t_2(y_1, \dots, y_n)) \} = \gamma$$

can be determined. γ is called the *confidence coefficient*.

Ex: Estimators for the mean of Gaussian r.v.'s

- Let the $y_i \sim \mathcal{N}(\theta, 5)$ be i.i.d. (independent and identically distributed) normal random variables with unknown mean θ and variance equal to 5.
- The arithmetic mean is a *point estimator* of θ :

$$\hat{\theta} = \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i \sim \mathcal{N}\left(\theta, \frac{5}{n}\right).$$

- Now, since $\bar{y} \sim \mathcal{N}(\theta, \frac{5}{n})$ or $\bar{y} = \theta + v$ with $v \sim \mathcal{N}(0, \frac{5}{n})$, we can also state

$$\Pr\left\{\theta \in \left(\bar{y} - 2\sqrt{\frac{5}{n}}, \bar{y} + 2\sqrt{\frac{5}{n}}\right)\right\} = \Pr\left\{\bar{y} \in \left(\theta - 2\sqrt{\frac{5}{n}}, \theta + 2\sqrt{\frac{5}{n}}\right)\right\} = 0.95$$

since

$$\bar{y} - 2\sqrt{\frac{5}{n}} < \theta < \bar{y} + 2\sqrt{\frac{5}{n}} \Leftrightarrow \theta - 2\sqrt{\frac{5}{n}} < \bar{y} < \theta + 2\sqrt{\frac{5}{n}}$$

so that $(\bar{y} - 2\sqrt{\frac{5}{n}}, \bar{y} + 2\sqrt{\frac{5}{n}})$ is an *interval estimator* for θ with confidence coefficient 0.95.

Estimation Considerations

- The design of interval estimators is the result of the following compromise. On the one hand, one wants a small interval so that θ can be known fairly accurately. On the other hand, a small interval necessarily leads to a small confidence coefficient since in general we can only state with a low probability level that θ is contained in a small interval. But we would like the confidence coefficient to be high so that we would be able to say with high probability where θ is located. This again would lead to a large interval etc.
- We will restrict the further discussion to point estimators. In general, the estimators we shall consider will asymptotically (for many measurements, large n) have a Gaussian distribution so that we can easily construct an interval estimator from the point estimator as in the example above.
- As a final remark, we shall assume that the parameters (and their estimators) can take on a continuous range of values so that their distribution is described by a probability density function. In case the parameters can take on only a *discrete set of values* and the task is to decide which one of the values the parameters actually take, then the estimation problem is called a *detection* or *decision* problem (see DIGICOM course).

Bayes Estimation

- nonnegative cost function $\mathcal{C}(\theta, \hat{\theta})$
- since $\theta, \hat{\theta}(Y)$ are random variables, we can never hope to make this cost function small for every outcome of θ and Y . All we can hope for is to minimize the expected value of the cost function, also called the risk

$$\mathcal{R}(\hat{\theta}(.)) = E \mathcal{C}(\theta, \hat{\theta}(Y)) = E_{\theta, Y} \mathcal{C}(\theta, \hat{\theta}(Y)) .$$

So the estimator is that function $\hat{\theta}(.)$ that minimizes the Bayes risk

$$\hat{\theta}(.) = \arg \min_{\hat{\theta}(.)} \mathcal{R}(\hat{\theta}(.))$$

It is common practice to limit the choice of the cost function to a nonnegative cost function of the parameter estimation error

$$\tilde{\theta} = \theta - \hat{\theta}(Y)$$

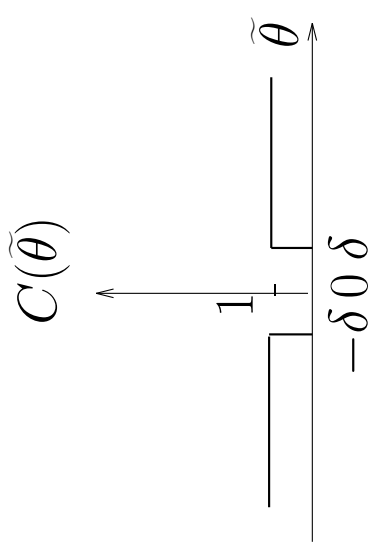
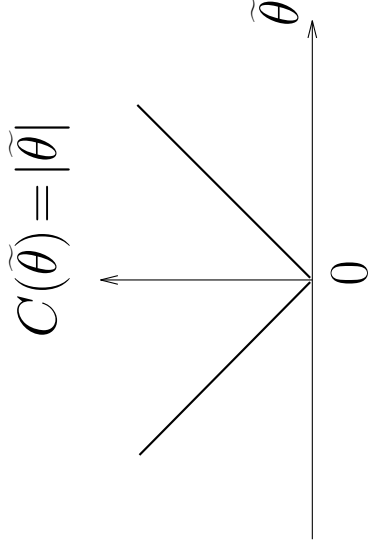
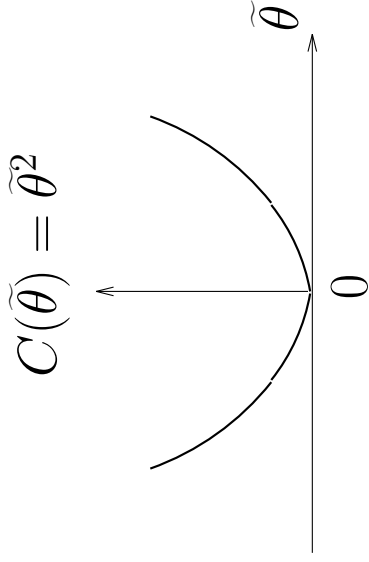
only, in which case the estimation problem becomes

$$\hat{\theta}(.) = \arg \min_{\hat{\theta}(.)} E \mathcal{C}(\tilde{\theta})$$



Bayes Cost Functions

- Three popular Bayes cost functions: squared parameter deviation, absolute parameter deviation, and the uniform cost function. All three choices assign no cost when there is no error.



Bayes Risk Minimization

- manipulate the Bayes optimization criterion

$$\begin{aligned}\min_{\hat{\theta}(\cdot)} \mathcal{R}(\hat{\theta}(\cdot)) &= \min_{\hat{\theta}(\cdot)} E \mathcal{C}(\hat{\theta}) = \min_{\hat{\theta}(\cdot)} \int \int f(Y, \theta) \mathcal{C}(\theta - \hat{\theta}(Y)) dY d\theta \\ &= \min_{\hat{\theta}(\cdot)} \underbrace{\int f(Y) dY}_{\geq 0} \underbrace{\int f(\theta|Y) \mathcal{C}(\theta - \hat{\theta}(Y)) d\theta}_{\mathcal{R}(\hat{\theta}(\cdot)|Y)} = \min_{\hat{\theta}(\cdot)} E_Y \mathcal{R}(\hat{\theta}(\cdot)|Y) = E_Y \min_{\hat{\theta}(Y)} \mathcal{R}(\hat{\theta}(Y)|Y)\end{aligned}$$

$\mathcal{R}(\hat{\theta}(\cdot)|Y) = \mathcal{R}(\hat{\theta}(Y)|Y)$ is a function of Y .

- Since the contributions of the conditional risk $\mathcal{R}(\hat{\theta}(Y)|Y)$ for every Y are combined via the nonnegative weighting factor $f(Y)$ to obtain the global risk $\mathcal{R}(\hat{\theta}(\cdot))$, we can minimize $\mathcal{R}(\hat{\theta}(\cdot))$ by minimizing $\mathcal{R}(\hat{\theta}(Y)|Y)$ for every particular Y .
- While the minimization of the global risk $\mathcal{R}(\hat{\theta}(\cdot))$ is with respect to (w.r.t.) the estimator function $\hat{\theta}(\cdot)$, the minimization of the conditional risk $\mathcal{R}(\hat{\theta}(Y)|Y)$ is w.r.t. the particular estimate $\hat{\theta}(Y)$, which is simply a number (the estimator function $\hat{\theta}(\cdot)$ evaluated at $Y = Y$).

Bayes Risk Minimization (2)

- This minimization of the conditional risk $\mathcal{R}(\hat{\theta}(Y)|Y)$ requires the *a posteriori* distribution $f(\theta|Y)$ of θ given the measurement Y .
- The a posteriori distribution describes the randomness left in θ after we have made a measurement of (the related quantity) Y .
- The a posteriori distribution $f(\theta|Y)$ can be determined from the conditional distribution $f(Y|\theta)$ and the a priori distribution $f(\theta)$, which are normally given in the problem description, as follows (using Bayes' rule):

$$f(\theta|Y) = \frac{f(Y, \theta)}{f(Y)} = \frac{f(Y|\theta)f(\theta)}{\int_{\Theta} f(Y, \theta)d\theta} = \frac{f(Y|\theta)f(\theta)}{\int_{\Theta} f(Y|\theta)f(\theta)d\theta}$$

where Θ is the region of support for θ .

The MMSE Criterion

- using the quadratic cost function $\mathcal{C}_{MMSE}(\hat{\theta}) = |\hat{\theta}|^2$, minimizing the conditional Bayes risk yields

$$\min_{\hat{\theta}(Y)} \mathcal{R}_{MMSE}(\hat{\theta}(Y)|Y) = \min_{\hat{\theta}(Y)} \int_{-\infty}^{\infty} f(\theta|Y) (\theta - \hat{\theta}(Y))^2 d\theta$$

- to take the derivative w.r.t. $\hat{\theta}$, recall Leibnitz's rule:

$$\frac{\partial}{\partial u} \int_{g_1(u)}^{g_2(u)} h(u, v) dv = \int_{g_1(u)}^{g_2(u)} \frac{\partial h(u, v)}{\partial u} dv + \frac{dg_2(u)}{du} h(u, g_2(u)) - \frac{dg_1(u)}{du} h(u, g_1(u))$$

- Using Leibnitz's rule, we obtain

$$\frac{\partial}{\partial \hat{\theta}} \mathcal{R}_{MMSE}(\hat{\theta}|Y) = 2 \int_{-\infty}^{\infty} f(\theta|Y) (\hat{\theta} - \theta) d\theta = 0$$

where we put the derivative equal to zero in order to find an extremum.

The MMSE Criterion (2)

- We can rewrite as

$$\widehat{\theta}(Y) \underbrace{\int_{-\infty}^{\infty} f(\theta|Y) d\theta}_{=1} = \int_{-\infty}^{\infty} \theta f(\theta|Y) d\theta \Rightarrow \widehat{\theta}_{MMSE}(Y) = E(\theta|Y)$$

which is the *mean* of the a posteriori distribution of θ given Y .

- To know whether this extremum is a local minimum, we can apply Leibnitz's rule a second time to obtain

$$\frac{\partial^2}{\partial \widehat{\theta}^2} \mathcal{R}_{MMSE}(\widehat{\theta}|Y) = 2 \int_{-\infty}^{\infty} f(\theta|Y) d\theta = 2 > 0$$

Hence the extremum at $\widehat{\theta}(Y) = E(\theta|Y)$ is a local minimum and it is furthermore the global minimum since it is the unique local minimum. Note that $E(\theta|Y)$ is indeed a function of Y .

The absolute value cost function

- For the absolute value cost function $\mathcal{C}_{ABS}(\hat{\theta}) = |\hat{\theta}|$, the minimization problem of the conditional Bayes risk becomes

$$\begin{aligned} \min_{\hat{\theta}(Y)} \mathcal{R}_{ABS}(\hat{\theta}(Y)|Y) &= \min_{\hat{\theta}(Y)} \int_{-\infty}^{\infty} f(\theta|Y) |\theta - \hat{\theta}(Y)| d\theta \\ &= \min_{\hat{\theta}(Y)} \left[\int_{-\infty}^{\hat{\theta}} f(\theta|Y) (\hat{\theta} - \theta) d\theta + \int_{\hat{\theta}}^{\infty} f(\theta|Y) (\theta - \hat{\theta}) d\theta \right]. \end{aligned}$$

- We shall again use Leibnitz's rule to take the derivative. For the first integral, we let $h(\hat{\theta}, \theta) = f(\theta|Y) (\hat{\theta} - \theta)$ and we get

$$h(\hat{\theta}, g_2(\hat{\theta})) = h(\hat{\theta}, \hat{\theta}) = f(\hat{\theta}|Y) (\hat{\theta} - \hat{\theta}) = 0$$

and $\frac{d g_1(\hat{\theta})}{d \hat{\theta}} = 0$ since the lower limit does not depend on $\hat{\theta}$. The corresponding terms for the second integral are similarly equal to zero.

- So the only terms in the derivative remaining are those obtained by differentiating the integrands directly:

$$\frac{\partial}{\partial \hat{\theta}} \mathcal{R}_{ABS}(\hat{\theta}|Y) = \int_{-\infty}^{\hat{\theta}} f(\theta|Y) d\theta - \int_{\hat{\theta}}^{\infty} f(\theta|Y) d\theta = 0$$

where again we put the derivative equal to zero in order to find an extremum.

The absolute value cost function (2)

- So the condition for an extremum becomes

$$F(\widehat{\theta}|Y) = \int_{-\infty}^{\widehat{\theta}} f(\theta|Y) d\theta = \int_{\widehat{\theta}}^{\infty} f(\theta|Y) d\theta = 1 - F(\widehat{\theta}|Y)$$

where $F(\theta|Y)$ is the a posteriori cumulative distribution function (cdf) of θ given Y .

- We can rewrite this also as

$$F(\widehat{\theta}_{ABS}(Y)|Y) = \frac{1}{2}$$

which means that $\widehat{\theta}_{ABS}(Y)$ is the *median* of the a posteriori distribution of θ given Y .

The absolute value cost function (3)

- Again, to know whether this extremum is a local minimum, we can apply Leibnitz's rule a second time to obtain

$$\frac{\partial^2}{\partial \hat{\theta}^2} \mathcal{R}_{ABS}(\hat{\theta}|Y) \Big|_{\hat{\theta}=\hat{\theta}_{ABS}} = 2 f(\hat{\theta}_{ABS}|Y) \geq 0 .$$

If $f(\hat{\theta}_{ABS}|Y) = 0$ then, since $f(\theta|Y) \geq 0$, the first nonzero derivative w.r.t. θ of $f(\theta|Y)$ will be of even order and positive (this reasoning can be extended to the case where $f(\theta|Y)$ would not be sufficiently differentiable). Hence the extremum at $\hat{\theta}_{ABS}(Y)$ is a local minimum and it is furthermore the global minimum since it is the unique local minimum. Note again that $\hat{\theta}_{ABS}(Y)$ is indeed a function of Y .

The uniform cost function

- For the uniform cost function we get

$$\begin{aligned}\min_{\hat{\theta}(Y)} \mathcal{R}_{UNIF}(\hat{\theta}(Y)|Y) &= \min_{\hat{\theta}(Y)} \left[\left(\int_{-\infty}^{\hat{\theta}-\delta} + \int_{\hat{\theta}+\delta}^{\infty} \right) f(\theta|Y) d\theta \right] \\ &= \min_{\hat{\theta}(Y)} \underbrace{\left[\int_{-\infty}^{\infty} f(\theta|Y) d\theta \right]}_{=1} - \int_{\hat{\theta}-\delta}^{\hat{\theta}+\delta} f(\theta|Y) d\theta\end{aligned}$$

- Since we take δ to be arbitrarily small, the optimization problem becomes

$$\max_{\hat{\theta}(Y)} \int_{\hat{\theta}-\delta}^{\hat{\theta}+\delta} f(\theta|Y) d\theta \approx \max_{\hat{\theta}(Y)} 2\delta f(\hat{\theta}|Y) = 2\delta \max_{\hat{\theta}(Y)} f(\hat{\theta}|Y)$$

Hence, for δ arbitrarily small, $\mathcal{R}_{UNIF}(\hat{\theta}(Y)|Y)$ is minimized by choosing $\hat{\theta}$ to be the *mode* (location of the maximum) of the a posteriori distribution of θ given Y .

- For this reason, the estimator corresponding to a uniform cost function is nominally called the *Maximum A Posteriori* (MAP) estimator:

$$\hat{\theta}_{MAP}(Y) = \arg \max_{\theta} f(\theta|Y)$$

which is again a function of Y .

MAP Estimator: Remarks

- We may note that the same estimator is obtained by choosing the cost function $\mathcal{C}(\hat{\theta}) = 1 - \delta(\hat{\theta})$. While this cost function is cleaner in that it involves no limiting operation with δ , it might be considered non-intuitive since it is not nonnegative. However, one should keep in mind that adding or subtracting a constant to a cost function does not influence its minimizing argument.
- Instead of maximizing $f(\theta|Y)$, any strictly increasing function of it may be maximized. Since $f(\theta|Y)$ is often given in factored form and often contains exponential distributions, a convenient choice is to maximize

$$\ln f(\theta|Y) = \ln f(Y|\theta) + \ln f(\theta) - \ln f(Y) .$$

- Often $f(\theta|Y)$ satisfies certain regularity conditions so that $\hat{\theta}_{MAP}$ is a solution of

$$\frac{\partial}{\partial \theta} \ln f(\theta|Y) = 0 = \frac{\partial}{\partial \theta} \ln f(Y|\theta) + \frac{\partial}{\partial \theta} \ln f(\theta) .$$

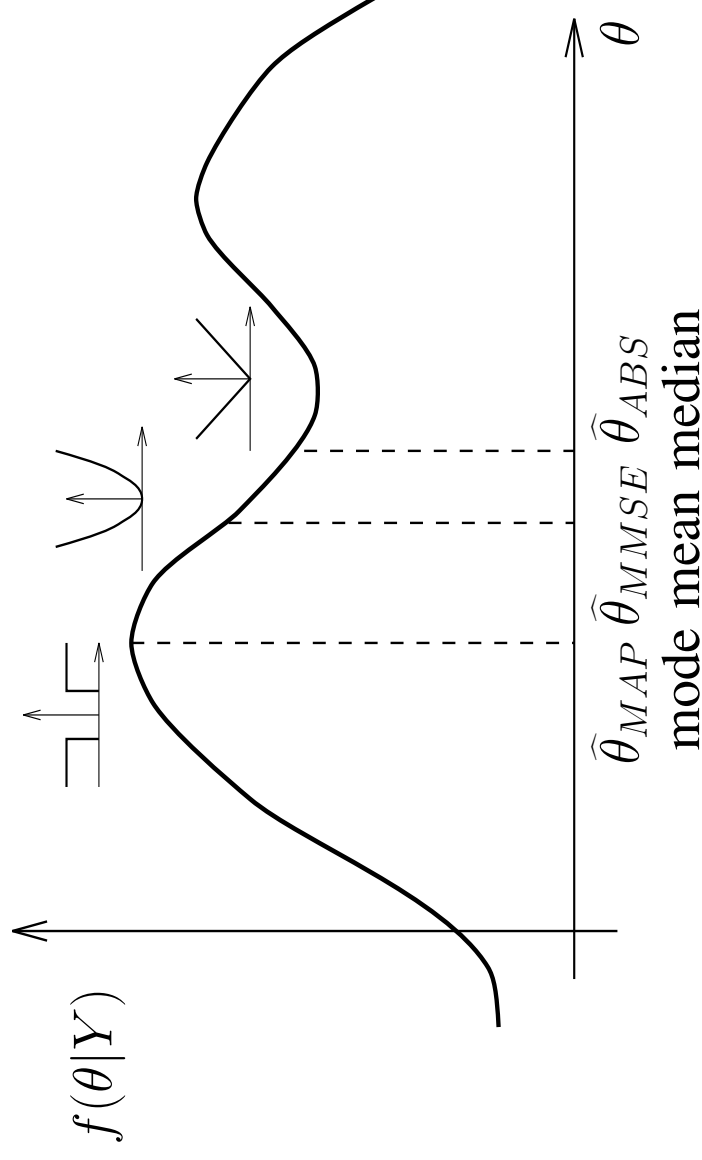
Note that $\frac{\partial \ln f(\theta|Y)}{\partial \theta} = \frac{1}{f(\theta|Y)} \frac{\partial f(\theta|Y)}{\partial \theta}$. The conditions for a maximum (rather than another form of extremum) need to be verified of course.

MAP Estimator: Remarks (2)

- The previous equation indicates that, starting from the description of the problem which includes $f(Y|\theta)$ and $f(\theta)$, the calculation of $\hat{\theta}_{MAP}$ should be relatively straightforward.
- The MAP estimator is given by the *global* maximum of $f(\theta|Y)$. If there are several local maxima, all of them need to be examined and compared to find the global maximum.
- Even if $f(\theta|Y)$ satisfies regularity conditions, the maximum may occur at the boundary of the parameter space Θ (which may not necessarily be $(-\infty, \infty)$). In that case, the maximum is not a local extremum.

3 Bayes estimators

- in general, the three Bayes estimators may differ



Ex: Gaussian mean in Gaussian noise

- estimating an unknown dc level in additive Gaussian noise with zero mean: e.g. satellite transmission of a digital image (accumulate many noisy images)
- for any given pixel, the measurement problem can be modeled as

$$y_i = \theta + v_i, \quad i = 1, \dots, n$$

where θ is the true grey value of the pixel, y_i are the consecutive noisy measured grey values, the $v_i \sim \mathcal{N}(0, \sigma_v^2)$ are i.i.d. (Central Limit Theorem).

- Even though the images vary very slowly, the image of one group of $n = 100$ shots will not differ very much from the image of the previous 100 shots. Therefore, we can consider the image estimate (with value m_θ at the pixel considered) from the previous group to be prior information for the problem of estimating the image of the current group (with value θ at the same pixel).

Ex: Gaussian mean in Gaussian noise (2)

- This prior information is not perfect due to the estimation variance associated with the processing of the previous group and also the variation from one image (group) to the next. Therefore we model the prior information as $\theta \sim \mathcal{N}(m_\theta, \sigma_\theta^2)$ where σ_θ^2 reflects the joint effect of estimation variability and variability due to change in time. Also, θ is assumed to be independent of the v_i .
- We can write the measurements in vector form

$$Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} = \theta \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} + \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} = \theta \mathbf{1} + V$$

With $X = \theta$, we can consider the following invertible transformation

$$\begin{bmatrix} X \\ Y \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ \mathbf{1} & I_n \end{bmatrix} \begin{bmatrix} \theta \\ V \end{bmatrix}.$$

Now, since θ and V are independent and both Gaussian, they are jointly Gaussian. A linear transformation of jointly Gaussian random variables leads again to Gaussian random variables. Hence, $X = \theta$ and Y are jointly Gaussian.

Ex: Gaussian mean in Gaussian noise (3)

- to determine a Bayes estimator, we need the a posteriori pdf $f(\theta|Y) = f(X|Y)$. Since X and Y are jointly Gaussian, the Gauss-Markov theorem tells us that $f(X|Y)$ is also Gaussian and it also tells us how to compute it from the first and second-order moments of X and Y . So we shall compute these moments.

- First-order moments:
$$\begin{aligned} m_X &= E X = m_\theta \\ m_Y &= E Y = E(\theta \mathbf{1} + V) = (E\theta) \mathbf{1} + E V = m_\theta \mathbf{1} \end{aligned}$$

- $C_{VV} = EVV^T = \sigma_v^2 I_n$ since $(EVV^T)_{i,j} = Ev_i v_j = \sigma_v^2 \delta_{i,j} = (\sigma_v^2 I_n)_{i,j}$

- The joint covariance matrix is

$$\begin{aligned} C &= E \begin{bmatrix} X - m_X \\ Y - m_Y \end{bmatrix} \begin{bmatrix} X - m_X \\ Y - m_Y \end{bmatrix}^T = E \begin{bmatrix} \theta - m_\theta \\ (\theta - m_\theta) \mathbf{1} + V \end{bmatrix} \begin{bmatrix} \theta - m_\theta \\ (\theta - m_\theta) \mathbf{1} + V \end{bmatrix}^T \\ &= \begin{bmatrix} \sigma_\theta^2 & \sigma_\theta^2 \mathbf{1}^T \\ \sigma_\theta^2 \mathbf{1} & \sigma_\theta^2 \mathbf{1} \mathbf{1}^T + \sigma_v^2 I \end{bmatrix} = \begin{bmatrix} C_{XX} & C_{XY} \\ C_{YX} & C_{YY} \end{bmatrix} \end{aligned}$$

where we used $E\{(\theta - m_\theta)V\} = E\{(\theta - m_\theta)\}$ $E\{V\} = 0$. $0 = 0$,

$$\begin{aligned} E[(\theta - m_\theta) \mathbf{1} + V][(\theta - m_\theta) \mathbf{1} + V]^T &= E(\theta - m_\theta)^2 \mathbf{1} \mathbf{1}^T + E V V^T + \mathbf{1} E V^T (\theta - m_\theta) \\ &\quad + E V (\theta - m_\theta) \mathbf{1}^T = \sigma_\theta^2 \mathbf{1} \mathbf{1}^T + \sigma_v^2 I + 0 + 0 \end{aligned}$$

Ex: Gaussian mean in Gaussian noise (4)

- A key quantity that appears in the expression for $f(X|Y)$ is $C_{XY}C_{YY}^{-1}$. We get

$$\begin{aligned} C_{XY}C_{YY}^{-1} &= \sigma_\theta^2 \mathbf{1}^T [\sigma_v^2 I + \sigma_\theta^2 \mathbf{1} \mathbf{1}^T]^{-1} \\ &= \frac{\sigma_\theta^2}{\sigma_v^2} \mathbf{1}^T \left[I + \frac{\sigma_\theta^2}{\sigma_v^2} \mathbf{1} \mathbf{1}^T \right]^{-1}. \end{aligned}$$

In order to compute the inverse of the matrix in brackets, we shall use the following identity.

- **Lemma 0.1 (Matrix Inversion Lemma)** *If A and C are respectively $n \times n$ and $m \times m$ invertible matrices and B and D are respectively $n \times m$ and $m \times n$ matrices, then*

$$[A + BCD]^{-1} = A^{-1} - A^{-1}B[DA^{-1}B + C^{-1}]^{-1}DA^{-1}$$

if the inverses exists.

- This lemma is very useful when we have to compute $[A + BCD]^{-1}$, we know A^{-1} , and m is (much) smaller than n . We shall apply this lemma with $A = I$, $B = \mathbf{1}$, $C = \frac{\sigma_\theta^2}{\sigma_v^2}$ and $D = \mathbf{1}^T$. Hence $m = 1$.

Ex: Gaussian mean in Gaussian noise (5)

- We obtain

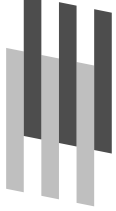
$$C_{XY}C_{YY}^{-1} = \frac{\sigma_\theta^2}{\sigma_v^2} \mathbf{1}^T \left[I - \mathbf{1}(\mathbf{1}^T \mathbf{1} + \frac{\sigma_v^2}{\sigma_\theta^2})^{-1} \mathbf{1}^T \right] = \underbrace{\frac{\sigma_\theta^2}{\sigma_v^2} \left(1 - \frac{n}{n + \frac{\sigma_v^2}{\sigma_\theta^2}} \right)}_{= n + \frac{\sigma_v^2}{\sigma_\theta^2}} \mathbf{1}^T = \frac{1}{n + \frac{\sigma_v^2}{\sigma_\theta^2}} \mathbf{1}^T$$

- From the Gauss-Markov theorem, we can now compute the a posteriori mean

$$\begin{aligned} E[\theta|Y] &= m_\theta + C_{XY}C_{YY}^{-1}(Y - m_Y) = m_\theta + \frac{1}{n + \frac{\sigma_v^2}{\sigma_\theta^2}} \mathbf{1}^T (Y - m_\theta \mathbf{1}) \\ &= \left(\frac{1}{\sigma_\theta^2} m_\theta + \frac{n}{\sigma_v^2} \bar{y} \right) / \left(\frac{1}{\sigma_\theta^2} + \frac{n}{\sigma_v^2} \right) = \frac{\frac{1}{\sigma_\theta^2}}{\left(\frac{1}{\sigma_\theta^2} + \frac{n}{\sigma_v^2} \right)} m_\theta + \frac{\frac{n}{\sigma_v^2}}{\left(\frac{1}{\sigma_\theta^2} + \frac{n}{\sigma_v^2} \right)} \bar{y} \end{aligned}$$

where we introduced the sample mean $\bar{y} = \frac{1}{n} \mathbf{1}^T Y = \frac{1}{n} \sum_{i=1}^n y_i$.

- Note that $E[\theta|Y]$ is of the form $\alpha m_\theta + (1-\alpha)\bar{y}$, which is a convex combination between the a priori mean and the sample mean obtained from the data. The respective weighting factors are proportional to the inverse of the variance associated with each of the two components.



Ex: Gaussian mean in Gaussian noise (6)

- We shall call the inverse of the variance the amount of information available:
 - $\frac{1}{\sigma_\theta^2}$ = information in the prior distribution $f(\theta)$,
 - $\frac{n}{\sigma_v^2}$ = information in the n independent measurements (conditional distribution $f(y_i|\theta) = f_{v_i}(y_i - \theta)$).
- We can consider two extreme cases:
 - $\sigma_\theta^2 \ll \frac{\sigma_v^2}{n} : E[\theta|Y] \approx m_\theta$. In this case, the measurements are so noisy that the reliable a priori information dominates.
 - $\sigma_\theta^2 \gg \frac{\sigma_v^2}{n} : E[\theta|Y] \approx \bar{y}$. In this case, the accurate measurements dominate the unreliable a priori information. Note that this second case will eventually occur when $n \rightarrow \infty$.

Ex: Gaussian mean in Gaussian noise (7)

- From the Gauss-Markov theorem, we can also determine the a posteriori variance, which we shall denote by σ_θ^2 :

$$\sigma_\theta^2 = \text{Var} [\theta|Y] = C_{XX} - C_{XY} C_{YY}^{-1} C_{YX} = \sigma_\theta^2 - \frac{1}{n + \frac{\sigma_v^2}{\sigma_\theta^2}} \mathbf{1}^T \mathbf{1} \sigma_\theta^2 = \dots = \left(\frac{1}{\sigma_\theta^2} + \frac{n}{\sigma_v^2} \right)^{-1}$$

or hence

$$\frac{1}{\sigma_\theta^2} = \frac{1}{\sigma_\theta^2} + \frac{n}{\sigma_v^2}$$

which means that the a priori information and the information in the n independent measurements add up to form the total information in the posterior distribution.

- This allows us to rewrite the a posteriori mean as

$$\frac{1}{\sigma_\theta^2} E [\theta|Y] = \frac{1}{\sigma_\theta^2} m_\theta + \frac{n}{\sigma_v^2} \bar{y}.$$

Ex: Gaussian mean in Gaussian noise (8)

- Using the Gauss-Markov theorem, we are finally ready to write the a posteriori distribution as

$$f(\theta|Y) \leftrightarrow \mathcal{N}(E[\theta|Y], \sigma_\theta^2)$$

where $E[\theta|Y]$ and σ_θ^2 are given above.

- Since the posterior distribution is Gaussian, its mean, mode and median coincide. Hence we get

$$\hat{\theta}_{MMSE} = \hat{\theta}_{MAP} = \hat{\theta}_{ABS} = E[\theta|Y] .$$

- Note that the posterior distribution $f(\theta, Y)$ depends on Y only through \bar{y} . In this case we say that the statistic \bar{y} (a function of the data Y) is a *sufficient statistic*, meaning that, for the purpose of estimating θ , the only thing we need to know about y_1, \dots, y_n is \bar{y} .

Equivalences of Bayes Estimators

One characteristic of Bayes estimation is that a Bayes estimator depends on the cost function it is associated with. In the following, we shall examine three cases in which the Bayes estimator coincides with $\hat{\theta}_{MMSE}$ meaning that the estimator minimizing the risk based on some non-quadratic cost function coincides with $\hat{\theta}_{MMSE}$, the estimator minimizing a quadratic cost function.

1. Consider a cost function $C(\tilde{\theta})$ with the following properties:

- symmetric: $C(\tilde{\theta}) = C(-\tilde{\theta})$
- convex : $C(\alpha\tilde{\theta}_1 + (1-\alpha)\tilde{\theta}_2) \leq \alpha C(\tilde{\theta}_1) + (1-\alpha) C(\tilde{\theta}_2)$, $\alpha \in [0, 1]$

and let $f(\theta|Y)$ be symmetric about $E[\theta|Y]$. Then $\hat{\theta}_C = \hat{\theta}_{MMSE}$.

Equivalences of Bayes Estimators (2)

2. Consider a cost function $C(\hat{\theta})$ with the following properties:

- symmetric: $C(\hat{\theta}) = C'(|\hat{\theta}|)$
- $C'(|\hat{\theta}|)$ is a non-decreasing function of $|\hat{\theta}|$ (this condition is weaker than $C(\hat{\theta})$ being convex)

and let $f(\theta|Y)$ have the following properties

- symmetric about $E[\theta|Y]$
- unimodal (only one local maximum)

and finally let $C(\cdot)$ and $f(\cdot|Y)$ be such that $\lim_{\theta \rightarrow \infty} C(\theta) f(\theta|Y) = 0$, $\forall Y$, then again $\hat{\theta}_C = \hat{\theta}_{MMSE}$.

3. If $f(\theta|Y)$ is Gaussian, then $\hat{\theta}_{MMSE} = \hat{\theta}_{MAP} = \hat{\theta}_{ABS}$.

These equivalences show the relative importance of $\hat{\theta}_{MMSE}$. In general however, $\hat{\theta}_{MAP}$ is the easiest to compute, as is illustrated in the following example.

Ex: Poisson Process

- Consider a communication network node where messages are passing. Let the observation \mathbf{N} be the number of messages that pass during a certain observation period of duration τ . \mathbf{N} has a Poisson distribution

$$\Pr[\mathbf{N} = N|\theta] = (\theta\tau)^N \frac{e^{-\theta\tau}}{N!}, \quad N = 0, 1, 2, \dots$$

where the parameter θ represents the average number of messages passing per second by the node we are considering.

- This average frequency has itself an a priori distribution (over the different nodes in the network) which we take to be exponential with average value $m_\theta = 1/\lambda$:

$$f(\theta) = \begin{cases} \lambda e^{-\lambda\theta} & , \quad \theta > 0 \\ 0 & , \quad \theta \leq 0. \end{cases}$$

Ex: Poisson Process (2)

- In order to find $\hat{\theta}_{MMSE}$, we shall compute the a posteriori distribution of θ given the measurement N . Using Bayes' rule we get

$$f(\theta|N) = \frac{\Pr[\mathbf{N} = N|\theta] f(\theta)}{\Pr[\mathbf{N} = N]} = \frac{1}{\Pr[\mathbf{N} = N]} \lambda \frac{\tau^N}{N!} \theta^N e^{-\theta(\tau+\lambda)} = k(N) \theta^N e^{-\theta(\tau+\lambda)}, \quad \theta > 0$$

where $k(N)$ depends on N but not on θ . Note that $f(\theta|N) = 0$, $\theta \leq 0$.

- In order to determine $k(N)$, we use the constraint

$$\int_0^\infty f(\theta|N) d\theta = 1 = k(N) \underbrace{\int_0^\infty \theta^N e^{-\theta(\tau+\lambda)} d\theta}_{g(N)}.$$

Hence $k(N) = 1/g(N)$. We shall calculate $g(N)$ by partial integration:

$$g(N) = \left[\frac{\theta^N e^{-\theta(\tau+\lambda)}}{-(\tau+\lambda)} \right]_0^\infty + \frac{N}{\tau+\lambda} \int_0^\infty \theta^{N-1} e^{-\theta(\tau+\lambda)} d\theta = \frac{N}{\tau+\lambda} g(N-1) = \dots = \frac{N!}{(\tau+\lambda)^N} g(0)$$

where

$$g(0) = \int_0^\infty e^{-\theta(\tau+\lambda)} d\theta = \left[\frac{e^{-\theta(\tau+\lambda)}}{-(\tau+\lambda)} \right]_0^\infty = \frac{1}{\tau+\lambda}$$

Ex: Poisson Process (3)

- which finally leads to

$$g(N) = \frac{N!}{(\tau + \lambda)^{N+1}}, k(N) = 1/g(N) = \frac{(\tau + \lambda)^{N+1}}{N!}.$$

- Now we can calculate

$$\begin{aligned}\widehat{\theta}_{MMSE} &= E[\theta|N] = \int_0^\infty \theta f(\theta|N) d\theta \\ &= \frac{1}{g(N)} \int_0^\infty \theta^{N+1} e^{-\theta(\tau+\lambda)} d\theta = \frac{g(N+1)}{g(N)} = \frac{N+1}{\tau + \lambda}.\end{aligned}$$

Ex: Poisson Process (4)

- To determine $\hat{\theta}_{MAP}$, consider

$$\ln f(\theta|N) = -\ln g(N) + N \ln \theta - \theta(\tau + \lambda) .$$

Differentiation w.r.t. θ yields

$$\frac{\partial}{\partial \theta} \ln f(\theta|Y) = 0 = \frac{N}{\theta} - (\tau + \lambda) .$$

So we obtain

$$\hat{\theta}_{MAP} = \frac{N}{\tau + \lambda} \neq \frac{N + 1}{\tau + \lambda} = \hat{\theta}_{MMSE} .$$

- Note however that if $\tau \gg \lambda$ (observation duration much longer than the a priori expected time between observations), then with high probability $N \gg 1$ and hence $\hat{\theta}_{MAP} \approx \hat{\theta}_{MMSE} \approx \frac{N}{\tau}$, which is simply the sample average of the number of messages per second.