

Matched, mismatched and semiparametric inference in elliptical distributions

Stefano Fortunati

Laboratoire des signaux et systèmes (L2S) & DR2I-IPSA,
Paris, France.

Uncertainty Quantification

Paris-Saclay

May 25, 2023, 2-3 pm

Outline of the presentation

Parametric models

- The CRB and the Maximum Likelihood (ML) estimator
- Parametric estimation in CES distributions

Misspecified models

- The MCRB and the Mismatched ML estimator
- Misspecified estimation in CES distributions

Semiparametric models

- The SCRB and the R -estimators
- Semiparametric inference in CES distributions

Introduction

- ▶ Let $\mathbf{x}_1, \dots, \mathbf{x}_L$ be the set of L observations collected from a random experiment.
- ▶ We indicate with $p_0(\mathbf{x}_1, \dots, \mathbf{x}_L)$ their joint “true” probability density function (pdf).
- ▶ In point estimation, we are interested in evaluating some functional of $p_0(\mathbf{x}_1, \dots, \mathbf{x}_L)$, say $\nu(p_0)$.
- ▶ However, $p_0(\mathbf{x}_1, \dots, \mathbf{x}_L)$ is generally unknown, at least to some extent.
- ▶ The lack of *a priori* knowledge on $p_0(\mathbf{x}_1, \dots, \mathbf{x}_L)$ can be formalized in the concept of *statistical models*.

Parametric models

- ▶ The most widely used statistical models are the *parametric* ones.
- ▶ A parametric model \mathcal{P}_θ is defined as a set of pdfs that are parametrized by a finite-dimensional parameter vector θ :

$$\mathcal{P}_\theta \triangleq \{p_X(\mathbf{x}_1, \dots, \mathbf{x}_L | \theta), \theta \in \Theta \subseteq \mathbb{R}^q\}.$$

- ▶ The underlying parametric assumption is that there exists $\theta_0 \in \Theta$, such that:

$$\mathcal{P}_\theta \ni p_X(\mathbf{x}_1, \dots, \mathbf{x}_L | \theta_0) = p_0(\mathbf{x}_1, \dots, \mathbf{x}_L). \quad (\text{A1})$$

- ▶ The (lack of) knowledge about the random experiment of interest is summarized in θ that needs to be estimated.

The Maximum Likelihood (ML) estimator

- ▶ We suppose that our observations $\mathbf{x}_1, \dots, \mathbf{x}_L$ are *iid* with “true” distribution $p_0(\mathbf{x})$, i.e. $\mathbf{x}_l \sim p_0, \forall l$.
- ▶ The *Maximum Likelihood* (ML) estimator, defined on the parametric model \mathcal{P}_θ , is given by:

$$\hat{\theta}_{L,ML} \triangleq \operatorname{argmax}_{\theta \in \Theta} \prod_{l=1}^L p_X(\mathbf{x}_l | \theta), \quad \mathbf{x}_l \sim p_0.$$

- ▶ The ML estimator is a cornerstone of the parametric estimation due to the following two optimality properties:
 1. *Consistency*,
 2. *Asymptotic Gaussianity and efficiency*.
- ▶ To well understand them, first we need to introduce the *Fisher Information Matrix* (FIM) $\mathbf{I}(\theta)$.

Fisher Information and Cramér-Rao Bound

- ▶ “Under some regularity conditions”¹, and under Assumption (A1), the FIM is defined as:

$$\begin{aligned}\mathbf{I}(\boldsymbol{\theta}) &\triangleq E \left\{ \nabla_{\boldsymbol{\theta}} \ln p_{\mathbf{X}}(\mathbf{x}|\boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}}^T \ln p_{\mathbf{X}}(\mathbf{x}|\boldsymbol{\theta}) \right\} \\ &\triangleq -E \left\{ \nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}}^T \ln p_{\mathbf{X}}(\mathbf{x}|\boldsymbol{\theta}) \right\}, \quad \mathbf{x} \sim p_0.\end{aligned}$$

Cramér-Rao Bound: Any unbiased estimator $\hat{\boldsymbol{\theta}}_L$ of $\boldsymbol{\theta}_0$, derived in $\mathcal{P}_{\boldsymbol{\theta}}$ from $\{\mathbf{x}_l \sim p_0\}_{l=1}^L$ iid observations, satisfies:

$$L \cdot E \left\{ (\hat{\boldsymbol{\theta}}_L - \boldsymbol{\theta}_0)(\hat{\boldsymbol{\theta}}_L - \boldsymbol{\theta}_0)^T \right\} \geq \mathbf{I}(\boldsymbol{\theta}_0)^{-1} \triangleq \text{CRB}(\boldsymbol{\theta}_0),$$

where the unbiasedness condition must hold, i.e. $\forall L \in \mathcal{N}$:

$$E_0\{\hat{\boldsymbol{\theta}}_L\} \triangleq \int \hat{\boldsymbol{\theta}}_L(\mathbf{x}_1, \dots, \mathbf{x}_L) p_0(\mathbf{x}_1, \dots, \mathbf{x}_L) d\mathbf{x}_1, \dots, d\mathbf{x}_L = \boldsymbol{\theta}_0.$$

¹Due to the limited time of the talk, we will not discuss them here. Moreover, we will omit to repeat this “magic” sentence in the following derivations.

- Why is the ML estimator so popular in applications?

Under Assumption (A1), the ML estimator $\hat{\boldsymbol{\theta}}_{L,ML}$ is:

1. \sqrt{L} -consistent:

$$\sqrt{L} \left(\hat{\boldsymbol{\theta}}_{L,ML} - \boldsymbol{\theta}_0 \right) = O_P(1).^2$$

2. *Asymptotically Gaussian and efficient*:

$$\sqrt{L} \left(\hat{\boldsymbol{\theta}}_{L,ML} - \boldsymbol{\theta}_0 \right) \underset{L \rightarrow \infty}{\overset{d}{\rightsquigarrow}} \mathcal{N}(\mathbf{0}, \mathbf{I}(\boldsymbol{\theta}_0)^{-1}) = \mathcal{N}(\mathbf{0}, \text{CRB}(\boldsymbol{\theta}_0)),$$

where $\underset{L \rightarrow \infty}{\overset{d}{\rightsquigarrow}}$ indicates the convergence in distribution.

² Let x_l be a sequence of random variables. Then $x_l = O_P(1)$ if for any $\epsilon > 0$, there exists a finite $N > 0$ and a finite $L > 0$, s.t. $\Pr \{|x_l| > N\} < \epsilon, \forall l > L$ (stochastic boundedness).

Covariance/scatter matrix estimation

- ▶ Estimating the correlation structure, i.e. the covariance matrix, of a dataset is a central problem in many applications:
 1. Dimensionality reduction and Principal Component Analysis,
 2. Signal/Image Denoising,
 3. Adaptive detection in radar/sonar systems,
 4. Graph signal processing,
 5. ...
- ▶ A general working assumption (motivated by the CLT) consists of assuming the data as Gaussian-distributed.
- ▶ However, this assumption is generally violated in practical applications where the data may be better characterized by heavy-tailed distributions.

A set of heavy-tailed distributions

- ▶ A family of non-Gaussian/heavy-tailed distribution is the class of **Complex Elliptically Symmetric (CES)** distributions.
- ▶ Thanks to their flexibility, CES distributions represent a reliable data model in many applications. ³
- ▶ The complex Gaussian, Generalized Gaussian, K -distribution, complex t -distribution and all the compound-Gaussian distributions belong to the CES class.
- ▶ The CES model is particularly useful in applications with *impulsive noise* and/or *spiky data*.

³E. Ollila, D. E. Tyler, V. Koivunen and H. V. Poor, "Complex Elliptically Symmetric Distributions: Survey, New Results and Applications", *IEEE Trans. on Signal Processing*, vol. 60, no. 11, pp. 5597-5625, Nov. 2012.

CES distributions (1/2)

- ▶ A CES distributed random vector $\mathbf{x} \in \mathbb{C}^N$ admits a pdf:

$$p_{\mathbf{x}}(\mathbf{x}) = |\Sigma|^{-1} h((\mathbf{x} - \boldsymbol{\mu})^H \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu})) \triangleq CES_N(\boldsymbol{\mu}, \Sigma, h).$$

- ▶ $h \in \mathcal{H}$, $h: \mathbb{R}_0^+ \rightarrow \mathbb{R}^+$ is the *density generator*,
 - ▶ $\boldsymbol{\mu} \in \mathbb{C}^N$ is the location vector,
 - ▶ $\Sigma \in \mathcal{M}_N$ is the (full rank) scatter matrix.
- ▶ Note that Σ and h are not jointly identifiable:

$$CES_N(\boldsymbol{\mu}, \Sigma, h(t)) \equiv CES_N(\boldsymbol{\mu}, c\Sigma, h(ct)), \quad \forall c > 0.$$

- ▶ To avoid this identifiability problem, we introduce the *shape matrix* as a normalized version of Σ :

$$\mathbf{V} \triangleq \Sigma / s(\Sigma).$$

CES distributions (2/2)

- ▶ Typical examples of *scale function* $s(\cdot)$ are:

$$s(\Sigma) = [\Sigma]_{11}, \quad s(\Sigma) = \text{tr}(\Sigma)/N \quad s(\Sigma) = |\Sigma|^{1/N}.$$

- ▶ Not that, *under finite second order moments*, if the scale $s(\Sigma) = \text{tr}(\Sigma)/N$ is adopted, we have that:

$$\mathbf{C} \triangleq E\{(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^H\} \equiv \Sigma = \sigma^2 \mathbf{V},$$

where:

1. \mathbf{C} is the covariance matrix of the CES-distributed vector $\mathbf{x} \sim \text{CES}_N(\boldsymbol{\mu}, \Sigma, h)$,
 2. $\sigma^2 = \text{tr}(\Sigma)/N = E\{\mathbf{x}^H \mathbf{x}\}/N$ is the statistical power of \mathbf{x} .
- ▶ Unless otherwise stated, in the following we always implicitly adopt the scale $s(\Sigma) = \text{tr}(\Sigma)/N$.

The complex t -distribution

- ▶ We suppose to collect $\{\mathbf{x}_l\}_{l=1}^L$, zero mean, iid observations that we assume to be t -distributed.
- ▶ The pdf of t -distributed data can be obtained from the CES family by specifying the density generator:

$$h_0(t) = \frac{1}{\pi^N} \frac{\Gamma(N + \lambda)}{\Gamma(\lambda)} \left(\frac{\lambda}{\eta}\right)^\lambda \left(\frac{\lambda}{\eta} + t\right)^{-(N+\lambda)},$$

1. λ : *shape parameter* controlling the data non-Gaussianity,
 2. η : *scale parameter* controlling the data power $\sigma^2 = \frac{\lambda}{\eta(\lambda-1)}$.
- ▶ To guarantee the finiteness of the second order moments (i.e. the existence of the covariance matrix), we need $\lambda > 1$.
 - ▶ Note that for values of $\lambda \rightarrow 1$ the data are heavy-tailed, while for $\lambda \rightarrow \infty$ the data tends to be Gaussian.

The parametric t -model

- ▶ Under this (zero-mean) t -assumption, the parametric model characterizing the random experiment is:

$$\mathcal{P}_{\theta} \triangleq \left\{ p_{\mathbf{X}}(\mathbf{x}|\theta) = |\Sigma|^{-1} h_0 \left(\mathbf{x}^H \Sigma^{-1} \mathbf{x} \right), \theta \in \Theta \right\}.$$

- ▶ The parameter space is defined as:

$$\theta \in \Theta \triangleq \{ \theta = \text{vec}(\mathbf{V}) | \mathbf{V} = N\Sigma / \text{tr}(\Sigma) \}.$$

- ▶ Optimal inference in this t -model will require the derivation of the ML estimator for the three parameter in $\theta \in \Theta$.
- ▶ Note that $\theta \in \Theta \subseteq \mathbb{C}^{N^2}$ is a complex-valued vector. In the following, we will implicitly use the Wirtinger calculus.

The parametric t -model: ML estimator for \mathbf{V}

- ▶ The ML estimator of the scatter matrix Σ of CES-distributed data can be expressed in term of a fixed-point equation. ⁴
- ▶ Then, an estimator of the constrained shape matrix \mathbf{V} is given by the convergence point of the iterative procedure:

$$\begin{cases} \hat{\Sigma}^{(k+1)} = \frac{N+\lambda}{L} \sum_{l=1}^L \frac{\mathbf{x}_l^H \mathbf{x}_l}{\mathbf{x}_l^H [\hat{\Sigma}^{(k)}]^{-1} \mathbf{x}_l + \lambda/\eta} , \\ \hat{\mathbf{V}}_{SML}^{(k+1)} \triangleq N \hat{\Sigma}^{(k+1)} / \text{tr}(\hat{\Sigma}^{(k+1)}) \end{cases}$$

where, as starting point, we use $\Sigma^{(0)} = \mathbf{I}_N$.

- ▶ Constraining the ML estimator of the scatter matrix lead to a sub-optimal estimator for the shape matrix, i.e. $\hat{\mathbf{V}}_{SML}$.

⁴E. Ollila, D. E. Tyler, V. Koivunen and H. V. Poor, "Complex Elliptically Symmetric Distributions: Survey, New Results and Applications", *IEEE Trans. on Signal Processing*, vol. 60, no. 11, pp. 5597-5625, Nov. 2012.

The parametric t -model: performance

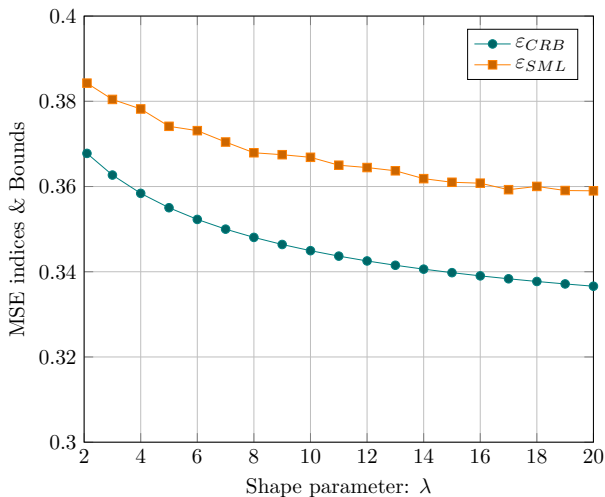
- ▶ Let $\boldsymbol{\theta}_0 = \text{vec}(\mathbf{V}_0)$ be the true parameter vector.
- ▶ The constrained joint Cramér-Rao Bound $\text{CRB}(\boldsymbol{\theta}_0)$ for $\boldsymbol{\theta} \in \Theta$ has been derived in ⁵ and it is not reported here.
- ▶ We compare the performance of the joint ML algorithm for the estimation of \mathbf{V}_0 in terms of *Mean Squared Error (MSE)*

$$\varepsilon_{JML} = \|E\{\text{vec}(\widehat{\mathbf{V}}_{SML} - \mathbf{V}_0)\text{vec}(\widehat{\mathbf{V}}_{SML} - \mathbf{V}_0)^H\}\|_F.$$

- ▶ As performance bound, we plot: $\varepsilon_{CRB} = \|\text{CRB}(\boldsymbol{\theta}_0)\|_F$.
- ▶ Number of observations: **finite sample regime** $L = 5N$.

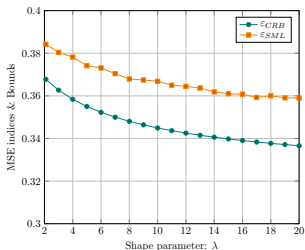
⁵S. Fortunati, F. Gini and M. Greco, "Matched, mismatched, and robust scatter matrix estimation and hypothesis testing in complex t -distributed data," *EURASIP J. Adv. Signal Process.* 2016, 123 (2016).

The parametric t -model: performance



- ▶ Under the “matched model” Assumption (A1), the MSE of $\hat{\mathbf{V}}_{SML}$ estimator is close to the CRB.

The parametric t -model: performance



- ▶ However, the SML is not fully efficient, i.e. its MSE is not exactly equal to the CRB, for two main reasons:
 1. The constraint on Σ has been imposed in a suboptimal way, without relying on a constrained optimization procedure.
 2. We are not in the asymptotic regime. In fact, we assumed $L = 5N$, i.e. L does not tend to infinity for a given N .
- ▶ **What if the model assumption is wrong?**

Outline of the presentation

Parametric models

- The CRB and the Maximum Likelihood (ML) estimator
- Parametric estimation in CES distributions

Misspecified models

- The MCRB and the Mismatched ML estimator
- Misspecified estimation in CES distributions

Semiparametric models

- The SCRB and the R -estimators
- Semiparametric inference in CES distributions

Model misspecification

- ▶ **Classical “matched” assumption:** the true data model and the model assumed to derive the estimation algorithm are the same, i.e. the model is correctly specified.
- ▶ All the results on the ML estimator and the CRB rely on this implicit assumption.
- ▶ However, much evidence from everyday practice shows that this assumption is often violated.
- ▶ **Model misspecification:** the assumed data model (i.e. the data pdf) differs from the true model.

Model misspecification

- ▶ There are two main reasons for model misspecification:
 1. An **imperfect knowledge** of the true data model that leads to a wrong specification of the data pdf.
 2. The true data model is known but it is **too involved** to pursue the optimal “matched” estimator.
- ▶ One may be forced (1) or may prefer (2) to derive an estimator by assuming a *simpler* but *misspecified* data model.
- ▶ This suboptimal procedure may lead to some degradation in the overall system performance.

Formal description of the misspecification

- ▶ Our observations $\{\mathbf{x}_l\}_{l=1}^L$ are *iid* with “true” distribution $p_0(\mathbf{x})$ belonging to a possibly non-parametric model \mathcal{P} .
- ▶ To characterize the statistical behavior of $\mathbf{x}_l, \forall l$, we adopt a different parametric pdf, say $\mathbf{x}_l \sim f_X(\mathbf{x}|\gamma)$, with $\gamma \in \Gamma \subseteq \mathbb{R}^p$.
- ▶ The adopted pdf $f_X(\mathbf{x}|\gamma)$ is assumed to belong to a *possibly misspecified* parametric model :

$$\mathcal{F}_\gamma \triangleq \{f_X(\mathbf{x}|\gamma), \gamma \in \Gamma\}.$$

- ▶ The classical “matched” assumption (A1) requires:

$$\exists \bar{\gamma} \in \Gamma, f_X(\mathbf{x}|\bar{\gamma}) = p_0(\mathbf{x}),$$

or, equivalently, that $p_0(\mathbf{x}) \in \mathcal{F}_\gamma$.

- ▶ If the previous assumption is violated, the model \mathcal{F}_γ is *misspecified*. Formally: ⁶

$$\forall \gamma \in \Gamma, f_X(\mathbf{x}|\gamma) \neq p_0(\mathbf{x}),$$

or, equivalently, that $p_0(\mathbf{x}) \in \mathcal{P} \not\subseteq \mathcal{F}_\gamma$.

- ▶ This misspecified scenario raises two main questions:
 1. How will the classical statistical properties of an estimator, e.g. *unbiasedness*, *consistency* and *efficiency*, change in this misspecified model framework?
 2. Is it still possible to derive lower bounds on the error covariance of any mismatched estimator?

⁶S. Fortunati, F. Gini, M. S. Greco and C. D. Richmond, "Performance Bounds for Parameter Estimation under Misspecified Models: Fundamental Findings and Applications", *IEEE Signal Processing Magazine*, vol. 34, no. 6, pp. 142-157, Nov. 2017.

The pseudo-true parameter vector

- ▶ “Under some regularity conditions”, there exist a unique interior point γ_0 of Γ , such that:

$$\gamma_0 \triangleq \operatorname{argmin}_{\gamma \in \Gamma} \{-E_0\{\ln f_X(\mathbf{x}|\gamma)\}\} = \operatorname{argmin}_{\gamma \in \Gamma} \{D_\gamma(p_0 \parallel f_X)\},$$

where $E_0\{g(\mathbf{x})\} \triangleq \int g(\mathbf{x})p_0(\mathbf{x})d\mathbf{x}$ and

$$D_\gamma(p_0 \parallel f_X) \triangleq \int \ln \left(\frac{p_0(\mathbf{x})}{f_X(\mathbf{x}|\gamma)} \right) p_0(\mathbf{x})d\mathbf{x}$$

is the **Kullback-Leibler divergence** (KLD) between the true pdf and the assumed pdf.

- ▶ The *pseudo-true parameter vector* γ_0 is the point the minimizes the KLD between the true and the assumed pdfs.

Information matrices under misspecification

- ▶ Let \mathbf{A}_{γ_0} be the matrix defined as:

$$\mathbf{A}_{\gamma_0} \triangleq E_0 \left\{ \nabla_{\gamma} \nabla_{\gamma}^T \ln f_X(\mathbf{x}|\gamma_0) \right\}, \quad \mathbf{x} \sim p_0.$$

- ▶ Let \mathbf{B}_{γ_0} be the matrix defined as:

$$\mathbf{B}_{\gamma_0} \triangleq E_0 \left\{ \nabla_{\gamma} \ln f_X(\mathbf{x}|\gamma_0) \nabla_{\gamma}^T \ln f_X(\mathbf{x}|\gamma_0) \right\}, \quad \mathbf{x} \sim p_0.$$

- ▶ If the model is correctly specified, i.e. if $\exists \bar{\gamma} \in \Gamma$ such that $f_X(\mathbf{x}|\bar{\gamma}) = p_0(\mathbf{x})$, then:
 1. $\gamma_0 = \bar{\gamma}$, i.e. the pseudo-true parameter is equal to the true one (in the classical “matched” sense),
 2. $\mathbf{B}_{\gamma_0} = -\mathbf{A}_{\gamma_0} = \mathbf{I}(\bar{\gamma})$, where $\mathbf{I}(\bar{\gamma})$ is the Fisher Information Matrix for the (“matched” in this case) model $\mathcal{F}_{\bar{\gamma}}$.

The Misspecified Cramér-Rao Bound

- ▶ Given our $\{\mathbf{x}_l \sim p_0\}_{l=1}^L$ iid observations, let's build an estimator $\hat{\gamma}_L$ assuming the possibly misspecified model \mathcal{F}_γ .
- ▶ **Misspecified (MS)-unbiasedness property:** the estimator $\hat{\gamma}_L$ is said to be MS-unbiased iff:

$$E_0\{\hat{\gamma}_L\} \triangleq \int \hat{\gamma}_L(\mathbf{x}_1, \dots, \mathbf{x}_L) p_0(\mathbf{x}_1, \dots, \mathbf{x}_L) d\mathbf{x}_1, \dots, d\mathbf{x}_L = \gamma_0.$$

Misspecified CRB: Any MS-unbiased estimator $\hat{\gamma}_L$ of γ_0 , derived in \mathcal{F}_γ from $\{\mathbf{x}_l \sim p_0\}_{l=1}^L$ iid observations, satisfies:^{7,8,9}

$$L \cdot E_0 \left\{ (\hat{\gamma}_L - \gamma_0)(\hat{\gamma}_L - \gamma_0)^T \right\} \geq \mathbf{A}_{\gamma_0}^{-1} \mathbf{B}_{\gamma_0} \mathbf{A}_{\gamma_0}^{-1} \triangleq \text{MCRB}(\gamma_0).$$

⁷Q. H. Vuong, "Cramér-Rao bounds for misspecified models", *Working paper 652, Division of the Humanities and Social Sciences, Caltech*, October 1986.

⁸S. Fortunati, F. Gini, M. S. Greco, "The Constrained Misspecified Cramér-Rao Bound", *IEEE Signal Process. Letters*, vol. 23, No. 5, pp. 718-721, May 2016.

⁹S. Fortunati, "Misspecified Cramér-Rao Bounds for Complex Unconstrained and Constrained Parameters," EUSIPCO 2017, Kos, Greece, 28 Aug. 2017-2 Sept. 2017

The Mismatched ML estimator (MML)

- ▶ The MML estimator, defined on the possibly misspecified parametric model \mathcal{F}_γ , is given by:

$$\hat{\gamma}_{L,MML} \triangleq \operatorname{argmax}_{\gamma \in \Gamma} \prod_{l=1}^L f_X(\mathbf{x}_l | \gamma), \quad \mathbf{x}_l \sim p_0.$$

Properties: the MML estimator $\hat{\gamma}_{L,MML}$ is: ^{10,11}

1. \sqrt{L} -MS-consistent:

$$\sqrt{L}(\hat{\gamma}_{L,MML} - \gamma_0) = O_P(1).$$

2. Asymptotically Gaussian and MS-efficient:

$$\sqrt{L}(\hat{\gamma}_{L,MML} - \gamma_0) \underset{L \rightarrow \infty}{\overset{d}{\rightsquigarrow}} \mathcal{N}(\mathbf{0}, \mathbf{A}_{\gamma_0}^{-1} \mathbf{B}_{\gamma_0} \mathbf{A}_{\gamma_0}^{-1}) = \mathcal{N}(\mathbf{0}, \text{MCRB}(\gamma_0)),$$

¹⁰ P. J. Huber, "The behavior of Maximum Likelihood Estimates under Nonstandard Conditions," *Proc. of the Fifth Berkeley Symposium in Mathematical Statistics and Probability*. Berkeley: University of California Press, 1967

¹¹ H. White, "Maximum likelihood estimation of misspecified models", *Econometrica* vol.50, pp.1-25, Jan. 1982.

A common misspecified scenario in CES data

- ▶ Our *iid* observations $\{\mathbf{x}_l \sim p_0\}_{l=1}^L$ are CES-distributed, that is $p_0 \sim CES_N(\boldsymbol{\mu}, \boldsymbol{\Sigma}_0, h_0)$.
- ▶ The “true” density generator is supposed to be the one of the *t*-distribution: $h_0(t) = \frac{1}{\pi^N} \frac{\Gamma(N+\lambda)}{\Gamma(\lambda)} \left(\frac{\lambda}{\eta}\right)^\lambda \left(\frac{\lambda}{\eta} + t\right)$.
- ▶ The practitioner decides to build a ML estimator for $\boldsymbol{\Sigma}_0$ on the misspecified Gaussian model \mathcal{F}_γ , such that:

$$\mathcal{F}_\gamma = \left\{ f_X(\mathbf{x}|\gamma) = |\boldsymbol{\Sigma}|^{-1} g\left(\mathbf{x}^H \boldsymbol{\Sigma}^{-1} \mathbf{x}\right), \gamma \in \Gamma \right\}.$$

where $g(t) = (\pi\sigma_X^2)^{-N} \exp(-t/\sigma_X^2)$ and

$$\gamma \in \Gamma \triangleq \left\{ \gamma = (\text{vec}(\mathbf{V})^T, \sigma_X^2)^T \mid \mathbf{V} = N\boldsymbol{\Sigma}/\text{tr}(\boldsymbol{\Sigma}) \right\},$$

- ▶ Clearly, $\forall \gamma \in \Gamma, f_X(\mathbf{x}|\gamma) \neq p_0(\mathbf{x})$: model mismatch!

Misspecified scatter matrix estimation

- ▶ Let's recall the Sample Covariance Matrix, i.e. the ML estimator under Gaussian assumption as:

$$\text{SCM} \triangleq \frac{1}{L} \sum_{l=1}^L \mathbf{x}_l \mathbf{x}_l^H.$$

- ▶ The Mismatched ML (MML) estimator can be derived as:

$$\begin{cases} \hat{\mathbf{V}}_{MML} = \frac{N}{\text{tr}(\text{SCM})} \text{SCM} \\ \hat{\sigma}_X^2 = \frac{1}{N \cdot L} \sum_{l=1}^L \mathbf{x}_l^H \hat{\mathbf{V}}_{MML}^{-1} \mathbf{x}_l \end{cases}.$$

- ▶ The practitioner should now answer the following questions:
 1. Is $\hat{\mathbf{V}}_{MML}$ a MS-consistent estimator for $\mathbf{V}_0 = N\boldsymbol{\Sigma}_0/\text{tr}(\boldsymbol{\Sigma}_0)$?
 2. Is it efficient wrt the MCRB?
 3. Is its performance loss wrt the matched case acceptable?

- ▶ To answer the first question, we need to evaluate the pseudo-true parameter vector

$$\gamma_0 \triangleq \operatorname{argmin}_{\gamma \in \Gamma} \{D_\gamma(p_0 \parallel f_X)\},$$

where p_0 is the t -distribution and f_X is the Gaussian one.

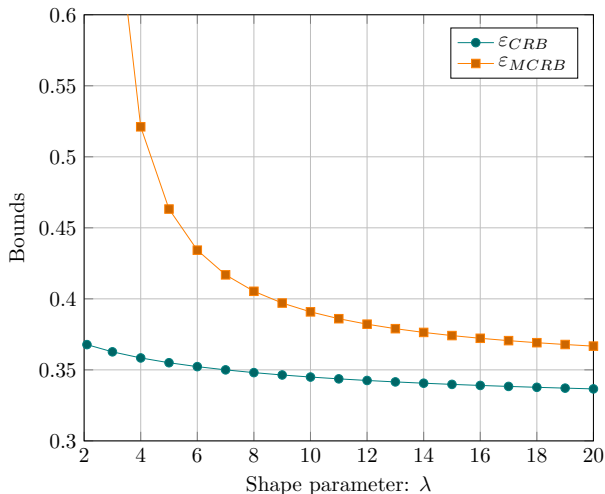
- ▶ It can be shown ¹² that $\gamma_0 = (\operatorname{vec}(\mathbf{V}_0)^T, \sigma_0^2)^T$ then:

1. $\sqrt{L}(\widehat{\mathbf{V}}_{MML} - \mathbf{V}_0) = O_P(1),$
2. $\sqrt{L}(\widehat{\sigma}_X^2 - \sigma_0^2) = O_P(1),$ where $\sigma_0^2 = \frac{\lambda_0}{\eta_0(\lambda_0 - 1)}.$

- ▶ The practitioner can use $\widehat{\mathbf{V}}_{MML}$ since it converge to the true shape matrix!

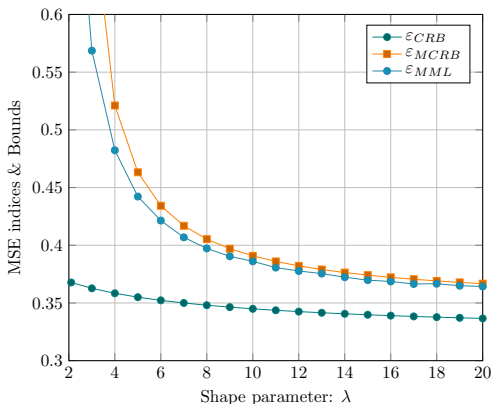
¹²S. Fortunati, F. Gini, M. S. Greco, "The Misspecified Cramér-Rao Bound and its Application to the Scatter Matrix estimation in Complex Elliptically Symmetric distributions," IEEE Trans. Signal Processing, vol. 64, no. 9, pp. 2387-2399, 2016.

Mispecified estimation performance: Bounds



- For small values of λ (highly non-Gaussian data), the estimation losses due to model mismatching rapidly increase!

Misspecified estimation performance: MSE



- ▶ The MSE of $\hat{\mathbf{V}}_{MML}$ is slightly below the MCRB because of the residual bias.
- ▶ **How can we overcome the misspecification?**

Outline of the presentation

Parametric models

- The CRB and the Maximum Likelihood (ML) estimator
- Parametric estimation in CES distributions

Misspecified models

- The MCRB and the Mismatched ML estimator
- Misspecified estimation in CES distributions

Semiparametric models

- The SCRB and the R -estimators
- Semiparametric inference in CES distributions

Semiparametric models

- ▶ A semiparametric model $\mathcal{P}_{\theta,h}$ is a set of pdfs characterized by a finite-dimensional parameter $\theta \in \Theta$ along with a *function*, i.e. an infinite-dimensional parameter, $h \in \mathcal{H}$:

$$\mathcal{P}_{\theta,h} \triangleq \{p_X(\mathbf{x}_1, \dots, \mathbf{x}_L | \theta, h), \theta \in \Theta \subseteq \mathbb{C}^q, h \in \mathcal{H}\}.$$

- ▶ Usually, θ is the (finite-dimensional) parameter of interest while h can be considered as a nuisance parameter.
- ▶ Most of the SP inference problems can be cast in the semiparametric framework:
 1. Inference in CES distributions (as we will see here),
 2. Estimation with missing data,
 3. Non-linear regression and inverse problems,
 4. Time series analysis, ...

CES distributions as semiparametric model

- ▶ The CES distributions family is a perfect example of *semiparametric model*.
- ▶ The (zero-mean and *iid*) CES semiparametric model can be obtained from the parametric one by relaxing the *unrealistic assumption* on the a-priori knowledge of the density generator:

$$\mathcal{P}_{\theta, h} \triangleq \left\{ p_{\mathbf{X}}(\mathbf{x}|\theta) = |\Sigma|^{-1} h\left(\mathbf{x}^H \Sigma^{-1} \mathbf{x}\right), \theta \in \Theta, h \in \mathcal{H} \right\},$$

where the parameter of interest is

$$\theta \in \Theta \triangleq \{\theta = \text{vec}(\mathbf{V}) | \mathbf{V} = N\Sigma / \text{tr}(\Sigma)\},$$

while $h \in \mathcal{H}$ is a nuisance function.

Inference in semiparametric model

- ▶ Semiparametric inference requires sophisticated tools in functional analysis and asymptotic statistics (e.g. the *Hájek–Le Cam convolution theorem*).
- ▶ Here, only a very short and non-exhaustive introduction will be provided aiming at highlighting some **crucial results**.
- ▶ As rigorously discussed in ¹³, any semiparametric inference scheme is based on the following key ingredients:
 1. The score vector of the parameter of interest $\mathbf{s}(\mathbf{x}; \boldsymbol{\theta}, h)$,
 2. The *nuisance tangent space* \mathcal{T}_h ,
 3. The *efficient score vector* $\bar{\mathbf{s}}(\mathbf{x}; \boldsymbol{\theta}, h)$.

¹³P.J. Bickel, C.A.J. Klaassen, Y. Ritov and J.A. Wellner, *Efficient and Adaptive Estimation for Semiparametric Models*, Johns Hopkins University Press, 1993.

The basic ingredients

- ▶ The score vector of the parameter of interest is defined as in the parametric case as:

$$\mathbf{s}(\mathbf{x}; \boldsymbol{\theta}, h) \triangleq \nabla_{\boldsymbol{\theta}} \ln p_{\mathcal{X}}(\mathbf{x}|\boldsymbol{\theta}, h).$$

- ▶ To define the *nuisance tangent space* \mathcal{T}_h and the associated projection operator $\Pi(\cdot|\mathcal{T}_h)$ we need the notion of *regular parametric sub-models*.
- ▶ The **efficient score vector** is defined as the residual of $\mathbf{s}(\mathbf{x}; \boldsymbol{\theta}, h)$ after projecting it on the nuisance tangent space \mathcal{T}_h :

$$\bar{\mathbf{s}}(\mathbf{x}; \boldsymbol{\theta}, h) \triangleq \mathbf{s}(\mathbf{x}; \boldsymbol{\theta}, h) - \Pi(\mathbf{s}(\mathbf{x}; \boldsymbol{\theta}, h)|\mathcal{T}_h),$$

- ▶ Let us finally introduce the **efficient information matrix** as:

$$\bar{\mathbf{I}}(\boldsymbol{\theta}|h) \triangleq E_0\{\bar{\mathbf{s}}(\mathbf{x}; \boldsymbol{\theta}, h)\bar{\mathbf{s}}(\mathbf{x}; \boldsymbol{\theta}, h)^T\}.$$

A lower bound in semiparametric estimation

- ▶ Let $\{\mathbf{x}_l\}_{l=1}^L$ be a set of iid observations, such that $\mathbf{x}_l \sim p_0(\mathbf{x}; \boldsymbol{\theta}_0, h_0) \in \mathcal{P}_{\boldsymbol{\theta}, h} \forall l$.
- ▶ The class of *Regular and Asymptotically Linear (RAL) estimators* is defined as:
 1. \sqrt{L} -consistent: $\sqrt{L}(\hat{\boldsymbol{\theta}}_L - \boldsymbol{\theta}_0) = O_P(1)$,
 2. Asymptotically normal: $\sqrt{L}(\hat{\boldsymbol{\theta}}_L - \boldsymbol{\theta}_0) \underset{L \rightarrow \infty}{\overset{d}{\rightsquigarrow}} \mathcal{N}(\mathbf{0}, \boldsymbol{\Xi}(\boldsymbol{\theta}_0, h_0))$.
- ▶ The ML and all the robust estimators belong to this class.

Semiparametric CRB (SCRIB): Any RAL estimator $\hat{\boldsymbol{\theta}}_L$ of $\boldsymbol{\theta}_0$, derived in $\mathcal{P}_{\boldsymbol{\theta}, h}$ from $\{\mathbf{x}_l \sim\}_{l=1}^L$ iid observations, satisfies: ¹⁴

$$\boldsymbol{\Xi}(\boldsymbol{\theta}_0, h_0) \geq \bar{\mathbf{I}}(\boldsymbol{\theta}_0 | h_0)^{-1} \triangleq \text{SCRIB}(\boldsymbol{\theta}_0 | h_0).$$

¹⁴P.J. Bickel, C.A.J. Klaassen, Y. Ritov and J.A. Wellner, *Efficient and Adaptive Estimation for Semiparametric Models*, Johns Hopkins University Press, 1993.

- ▶ Let us focus on the efficient estimation of the parameter of interest $\theta \in \Theta$ in the presence of the unknown function $h \in \mathcal{H}$.
- ▶ Clearly, Maximum Likelihood estimation is not an option.
- ▶ Is there any other “optimal” procedure for deriving asymptotically efficient estimates other than the ML one?
- ▶ The answer is positive and it is given by the semiparametric *rank-based (R-) Le Cam's “one step” estimators*.^{15,16,17}

¹⁵ L. Le Cam, “Locally asymptotically normal families of distributions,” *University of California Publications Statist.*, vol. 3, 1960, pp. 37-98.

¹⁶ P.J. Bickel, C.A.J. Klaassen, Y. Ritov and J.A. Wellner, *Efficient and Adaptive Estimation for Semiparametric Models*, Johns Hopkins University Press, 1993.

¹⁷ M. Hallin, B. J. M. Werker, “Semi-parametric efficiency, distribution-freeness and invariance,” *Bernoulli*, vol. 9, no. 1, pp. 137-165, 2003.

- ▶ Let $\{\mathbf{x}_l \sim p_0(\mathbf{x}; \boldsymbol{\theta}_0, h_0)\}_{l=1}^L$ be a set of iid observations.
- ▶ In the seminal paper ¹⁸, a *rank-based* (R -) class of one step estimators has been proposed:

$$\hat{\boldsymbol{\theta}}_{L,R} = \hat{\boldsymbol{\theta}}_L^* + L^{-1/2} \hat{\mathbf{Y}}_{\hat{\boldsymbol{\theta}}_L^*}^{-1} \tilde{\boldsymbol{\Delta}}_{\hat{\boldsymbol{\theta}}_L^*}$$

- ▶ $\boldsymbol{\theta}_L^*$ is a sub-optimal (consistent, not efficient) estimator of $\boldsymbol{\theta}$,
 - ▶ $\hat{\mathbf{Y}}_{\hat{\boldsymbol{\theta}}_L^*}$ is a *rank-based*, \sqrt{L} -consistent estimator of the *efficient information matrix* $\bar{\mathbf{I}}(\boldsymbol{\theta}_0|h_0)$,
 - ▶ $\tilde{\boldsymbol{\Delta}}_{\hat{\boldsymbol{\theta}}_L^*} \triangleq \sum_{l=1}^L \tilde{\boldsymbol{\varphi}}(\mathbf{x}_l, \boldsymbol{\theta}_L^*)$, where $\tilde{\boldsymbol{\varphi}}$ is a distributionally-free, *rank-based* approximation of the *efficient score vector*.
- ▶ No non-parametric estimator \hat{h}_L of $h \in \mathcal{H}$ is required!

¹⁸ M. Hallin, B. J. M. Werker, "Semi-parametric efficiency, distribution-freeness and invariance," *Bernoulli*, vol. 9, no. 1, pp. 137-165, 2003.

- ▶ An R -estimator built on $\mathcal{P}_{\theta,h}$ satisfies the same optimality properties of the ML estimator built on \mathcal{P}_{θ} ! ¹⁹

Under any possible $h \in \mathcal{H}$, the R -estimator $\hat{\theta}_{L,R}$ is:

1. \sqrt{L} -consistent:

$$\sqrt{L} \left(\hat{\theta}_{L,R} - \theta_0 \right) = O_P(1).$$

2. *Asymptotically Gaussian and "efficient"*:

$$\sqrt{L} \left(\hat{\theta}_{L,R} - \theta_0 \right) \underset{L \rightarrow \infty}{\overset{d}{\rightsquigarrow}} \mathcal{N}(\mathbf{0}, \bar{\mathbf{I}}(\theta_0 | h_0)^{-1}) = \mathcal{N}(\mathbf{0}, \text{SCR}(\theta_0 | h_0)).$$

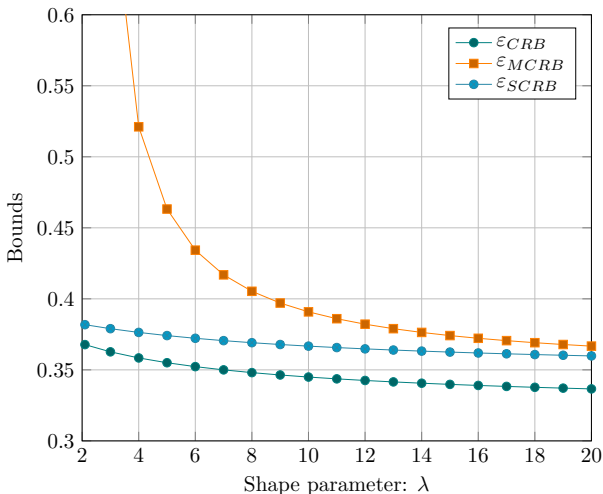
- ▶ **Note:** A classical alternative to R -estimators are the robust M -estimators. However, **the M -estimators are not efficient!**

¹⁹M. Hallin, B. J. M. Werker, "Semi-parametric efficiency, distribution-freeness and invariance," *Bernoulli*, vol. 9, no. 1, pp. 137-165, 2003.

- ▶ The R -estimator in real elliptical data has been proposed by Hallin, Oja and Paindaveine in ²⁰.
- ▶ Our recent works aimed:
 - ▶ Extension to complex-valued CES-distributed data,
 - ▶ Discussion about optimal setting,
 - ▶ Robustness to outliers,
 - ▶ Extensive comparison with other robust estimators,
 - ▶ Application to classical array processing estimation problems.

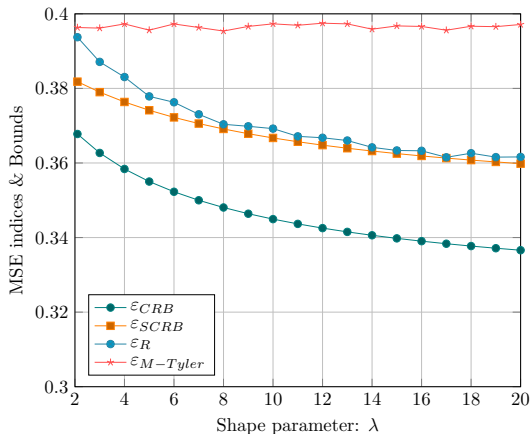
²⁰ M. Hallin, H. Oja, and D. Paindaveine, "Semiparametrically efficient rank-based inference for shape II. Optimal R-estimation of shape," *The Annals of Statistics*, vol. 34, no. 6, pp. 2757–2789, 2006.

Semiparametric estim. performance: Bounds



- The SCRb is in between the “unrealistic” CRB and the MCRB as expected.

Semiparametric estimation performance: MSE



► Regarding the MSE of the estimators:

1. The popular M -Tyler's covariance estimator is not efficient.
2. The R -estimator is (almost) semiparametrically efficient!

Concluding remarks

- ▶ In this talk, the estimation of the (normalized) covariance matrix has been addressed under three different frameworks:
 1. *Parametric*: perfect knowledge of the functional form of the density generator,
 2. *Misspecified*: wrong assumption on the density generator.
 3. *Semiparaemtric*: no assumption on the density generator.

- ▶ For the tree cases, we investigate the efficiency:
 1. *Parametric*: SML estimator and CRB,
 2. *Misspecified*: MML estimator and MCRB.
 3. *Semiparaemtric*: R -estimator and SCRB.

**An R -estimator is able to reach almost the SCRB !
Consequently, it is the best estimator when density generator
is unknown!**

▶ Short-term perspectives:

1. Rigorous analysis of the *almost efficiency* of the R -estimator (*Chernoff-Savage result*),
2. Derivation of an R -estimator of the eigenspace of the scatter matrix.

▶ Medium-term perspectives:

1. Semiparametric Mahalanobis distance,
2. Applications to clustering and distance learning.

▶ Long-term perspectives:

1. Semiparametric statistics and missing data,
2. Applications to array processing, image reconstruction, ecc...

R-estimators and SCRB for the CES model

1. S. Fortunati, F. Gini, M. S. Greco, A. M. Zoubir and M. Rangaswamy, "Semiparametric Inference and Lower Bounds for Real Elliptically Symmetric Distributions", *IEEE Transactions on Signal Processing*, vol. 67, no. 1, pp. 164-177, 1 Jan.1, 2019.
2. S. Fortunati, F. Gini, M. S. Greco, A. M. Zoubir and M. Rangaswamy, "Semiparametric CRB and Slepian-Bangs Formulas for Complex Elliptically Symmetric Distributions", *IEEE Transactions on Signal Processing*, vol. 67, no. 20, pp. 5352-5364, 15 Oct.15, 2019.
3. S. Fortunati, A. Renaux, F. Pascal, "Robust semiparametric efficient estimators in complex elliptically symmetric distributions", *IEEE Transactions on Signal Processing*, vol. 68, pp. 5003-5015, 2020.
4. S. Fortunati, A. Renaux, F. Pascal, "Joint Estimation of Location and Scatter in Complex Elliptical Distributions: A robust semiparametric and computationally efficient R -estimator of the shape matrix", *Journal of Signal Processing Systems*, July 2021.
5. S. Fortunati, A. Renaux, F. Pascal, "Properties of a new R -estimator of shape matrices", *EUSIPCO 2020*, Amsterdam, the Netherlands, August 24-28, 2020.
6. S. Fortunati, A. Renaux, F. Pascal, "Robust Semiparametric DOA Estimation in non-Gaussian Environment", *2020 IEEE Radar Conference*, Florence, Italy, September 21-25, 2020
7. S. Fortunati, A. Renaux, F. Pascal, "Robust Semiparametric Joint Estimators of Location and Scatter in Elliptical Distributions", *IEEE International Workshop on Machine Learning for Signal Processing*, Aalto University, Espoo, Finland, September 21-24, 2020.
8. All the code about real and complex R -estimator is provided in my GitHub page.

Thanks for your attention!

Backup slides

- ▶ Let $\mathbf{z} \triangleq (\mathbf{x}^T, \mathbf{y}^T)^T$ be a *complete* dataset, where:
 - ▶ \mathbf{x} is the *observed* (available) dataset.
 - ▶ \mathbf{y} is the *unobservable* (missing) dataset.
- ▶ **Problem:** Estimate $\theta \in \Theta$ from the observed dataset \mathbf{x} when the pdf p_Y of the missing data \mathbf{y} is unknown.
- ▶ The pdf p_X of the observed dataset can be expressed as:

$$p_X(\mathbf{x}|\theta) = \int_{\mathcal{Y}} p_{X,Y}(\mathbf{x}, \mathbf{y}|\theta) d\mathbf{y} = \int_{\mathcal{Y}} p_{X|Y}(\mathbf{x}|\mathbf{y}, \theta) p_Y(\mathbf{y}) d\mathbf{y}.$$

- ▶ The set of all the pdfs of the observed dataset \mathbf{x} is a *semiparametric mixture model* of the form :

$$\mathcal{P}_{\theta, p_Z} \triangleq \{p_X | p_X(\mathbf{x}|\theta, p_Y), \theta \in \Theta, p_Y \in \mathcal{K}\}.$$

Semiparametric models: Non-linear regression

- ▶ Let us consider the general non-linear regression model:

$$\mathbf{x} = f(\mathbf{z}, \boldsymbol{\theta}) + \epsilon,$$

- ▶ $\boldsymbol{\theta} \in \Theta$: parameter vector to be estimated,
 - ▶ $f \in \mathcal{F}$: possibly unknown non-linear function,
 - ▶ \mathbf{z} : random vector with possibly unknown pdf $p_Z \in \mathcal{K}$,
 - ▶ ϵ : random noise with possibly unknown pdf $p_\epsilon \in \mathcal{E}$
- ▶ The set of all pdfs for \mathbf{x} is a semiparametric model of the form:

$$\mathcal{P}_{\boldsymbol{\theta}, f, p_Z, p_\epsilon} \triangleq \{p_X(\mathbf{x}|\boldsymbol{\theta}, f, p_Z, p_\epsilon), \boldsymbol{\theta} \in \Theta, f \in \mathcal{F}, p_Z \in \mathcal{K}, p_\epsilon \in \mathcal{E}\}.$$

- ▶ This model is a general form of a *semiparametric regression model*.

Semiparametric models: Autoregressive processes

- ▶ Consider the $AR(p)$ process:

$$x_n = \sum_{i=1}^p \theta_i x_{n-i} + w_n, \quad n \in (-\infty, \infty)$$

- ▶ $\boldsymbol{\theta} \triangleq [\theta_1, \dots, \theta_p]$: parameter vector to be estimated.
 - ▶ w_n : i.i.d. innovations with unknown pdf $p_w \in \mathcal{W}$,
- ▶ Let $\mathbf{x} \in \mathbb{R}^N$ a vector of N observations from an $AR(p)$.
- ▶ The set of all possible pdfs for $\mathbf{x} \in \mathbb{R}^N$ is a semiparametric model:

$$\mathcal{P}_{\boldsymbol{\theta}, p_w} \triangleq \{p_X | p_X(\mathbf{x} | \boldsymbol{\theta}, p_w), \boldsymbol{\theta} \in \Theta, p_w \in \mathcal{W}\}.$$