



Technische Universität München

TUM School of Computation, Information and Technology

**Asymptotically Optimal Channel Estimation and
Learning Suitable Compressive Sensing Matrices**

Michael Koller

Vollständiger Abdruck der von der TUM School of Computation, Information and Technology der Technischen Universität München zur Erlangung des akademischen Grades eines

Doktors der Ingenieurwissenschaften

genehmigten Dissertation.

Vorsitz: Prof. Dr. Reinhard Heckel
Prüfende der Dissertation: 1. Prof. Dr.-Ing. Wolfgang Utschick
2. Prof. Dr.-Ing. Robert Schober

Die Dissertation wurde am 09.06.2022 bei der Technischen Universität München eingereicht und durch die TUM School of Computation, Information and Technology am 22.11.2022 angenommen.

Abstract

We study a sequence of conditional mean estimators (CMEs) and a method to learn an observation matrix for compressive sensing purposes. In both cases, we aim to recover a channel vector from noisy linear observations. The channel is assumed to be distributed according to an unknown probability density function (PDF) and we assume we have access to a training data set of channel samples.

First, we approximate the unknown channel PDF by means of a Gaussian mixture model (GMM). This is motivated by universal approximation properties of GMMs according to which for any given continuous PDF, there exists a sequence of GMMs which converges uniformly to it. For each GMM sequence element, a CME for channels distributed according to the GMM sequence element can be analytically calculated. Hence, the sequence of GMMs implies a sequence of CMEs. Assuming an invertible observation matrix, we show the pointwise convergence of the sequence of CMEs to the optimal CME which corresponds to the unknown channel PDF. However, since the channel PDF is unknown, this optimal CME cannot be computed, which makes one of the GMM-based CMEs an attractive approximation thereof. The pointwise convergence result holds more generally: Whenever a sequence of PDFs converges uniformly to the channel PDF, the corresponding sequence of CMEs converges pointwise to the optimal CME. We study the case of noninvertible observation matrices in numerical experiments and observe a convergent behavior.

Second, the goal is to learn an observation matrix with the restricted isometry property (RIP). To this end, we interpret a matrix with the RIP as a matrix which maps vectors of interest from a high-dimensional unit hypersphere to a low-dimensional unit hypersphere. In addition to this property, we argue that the vectors on the low-dimensional hypersphere should be uniformly distributed to combat noise. We formulate the goal to obtain such a matrix as the learning problem of matching the distribution of the channel in the range space of the matrix to a uniform distribution on the low-dimensional hypersphere. The distance between the distributions is measured by means of a maximum mean discrepancy (MMD) metric. Lastly, we observe that the GMM-based CME can be used in conjunction with the learned observation matrices.

Contents

Acronyms	v
1 Introduction	1
1.1 Asymptotically Optimal Estimation	1
1.2 Learning a Compressive Sensing Matrix	2
1.3 Motivation and Application	3
2 Preliminaries	5
2.1 Notation	5
2.2 Signal Model	5
2.3 Compressive Sensing	7
2.4 Gaussian Mixture Models	8
2.5 Maximum Mean Discrepancy	9
2.6 Distribution Matching	11
2.7 Random Search	12
3 Asymptotically Optimal Channel Estimation	13
3.1 A Sequence of Conditional Mean Estimators	13
3.2 Pointwise Convergence	15
3.3 Example: Gaussian Mixture Models	17
3.4 Channel Models for Numerical Evaluation	22
3.5 Other Channel Estimators	24
3.6 Numerical Evaluation	26
4 Learning an Observation Matrix	39
4.1 Proposed Method	39
4.2 Related Literature	43
4.3 Learning-Based Compressive Subsampling	46
4.4 Projected Gradient Descent Method	48
4.5 Comparison	49

Contents

4.6	Numerical Evaluation	50
5	Learned Observation Matrix and the Gaussian Mixture Estimator	61
5.1	Simulation Setup	61
5.2	Random Matrices	62
5.3	Learned Matrices	63
5.4	Different Evaluation Functions for Algorithm 6	66
6	Outlook	71
6.1	Asymptotically Optimal Channel Estimation	71
6.2	Learning a Compressive Sensing Matrix	72
A	Proof of Theorem 2	75
B	Discussion for Noninvertible Matrices	81
B.1	Integral in Lemma 1	81
B.2	Uniform Convergence	82
	Bibliography	83

Acronyms

AMP	approximate message passing
CAE	concrete autoencoder
CME	conditional mean estimator
CNN	convolutional neural network
DFT	discrete Fourier transform
EM	expectation-maximization
GAN	generative adversarial network
GMM	Gaussian mixture model
LBCS	learning-based compressive subsampling
LMMSE	linear minimum mean square error
LOS	line of sight
LS	least squares
MIMO	multiple-input multiple-output
MMD	maximum mean discrepancy
MSE	mean square error
NLOS	non-line of sight
O2I	outdoor-to-indoor
OMP	orthogonal matching pursuit
PDF	probability density function
PGDM	projected gradient descent method

Acronyms

RIP	restricted isometry property
SIMO	single-input multiple-output
SISO	single-input single-output
SNR	signal-to-noise ratio
ULA	uniform linear array

Introduction 1

In this dissertation, we are interested in the equation

$$\mathbf{y} = \mathbf{A}\mathbf{h} + \mathbf{n}$$

where \mathbf{A} is a matrix and the remaining quantities are vectors. In the studied setting, the vector \mathbf{y} is given and \mathbf{n} represents an unknown noise realization. The noise is modeled by means of a known Gaussian probability density function (PDF). The goal is to estimate the vector \mathbf{h} and we measure the estimation error in terms of the mean square error (MSE). The vector \mathbf{h} is a realization of a random variable with a PDF $f_{\mathbf{h}}$ which is continuous. The PDF $f_{\mathbf{h}}$ is not assumed to be known but we assume a (training) data set $\{\mathbf{h}_t\}_{t=1}^{T_{\text{tr}}}$ of realizations \mathbf{h}_t to be given.

In this context, we discuss the following two main topics:

1. We propose an algorithm to estimate \mathbf{h} and study its optimality (Chapter 3).
2. We propose an algorithm to obtain a matrix \mathbf{A} suitable for the problem of estimating \mathbf{h} (Chapter 4).

Since we work with a given data set, both topics make use of machine learning methods. The main tools used throughout the dissertation are introduced in Chapter 2. Chapter 5 presents a numerical simulation where the two proposed algorithms are combined. Lastly, Chapter 6 concludes the dissertation with an outlook.

1.1 Asymptotically Optimal Estimation

The first main topic (Chapter 3) is an algorithm to estimate \mathbf{h} . We use Gaussian mixture models (GMMs) to estimate the PDF $f_{\mathbf{h}}$ corresponding to \mathbf{h} . Every GMM is a PDF $f^{(K)}$ which consists of $K \in \mathbb{N}$ components. GMMs are universal approximators and can therefore approximate any continuous PDF arbitrarily well [1]. More precisely, in theory, there always exists a sequence $(f^{(K)})_{K=1}^{\infty}$ of GMMs which converges uniformly to the PDF $f_{\mathbf{h}}$. In practice, we take this as a motivation to use an expectation-maximization (EM) algorithm together with the given data set to obtain a K -component

GMM $f^{(K)}$ that approximates f_h . If the approximation is perfect, i.e., if $f^{(K)} = f_h$ holds, then the conditional mean estimator (CME), which achieves the minimum MSE, for \mathbf{h} can be computed in closed form. If the approximation is not perfect, we still compute a closed-form estimator based on $f^{(K)}$ and we can observe that this GMM-based estimator converges to the CME for $K \rightarrow \infty$. We prove this result for the case where the matrix \mathbf{A} is invertible. The noninvertible case is evaluated in numerical experiments and also shows a convergent behavior.

1.2 Learning a Compressive Sensing Matrix

The second main topic (Chapter 4) is an algorithm to design a matrix \mathbf{A} for the problem of estimating \mathbf{h} . Here, we are mainly interested in underdetermined systems where \mathbf{A} is a wide matrix that maps from an N -dimensional space to an m -dimensional space with $m < N$. Compressive sensing theory studies the problem of estimating \mathbf{h} in this setting when \mathbf{h} lies in a suitable form of a low-dimensional subspace. In particular, compressive sensing theory provides properties which, if a matrix \mathbf{A} fulfills them, guarantee good estimates. One of these properties is the restricted isometry property (RIP). The RIP comes along with a restricted isometry constant which is a nonnegative number where a smaller number indicates a better matrix. We interpret a matrix \mathbf{A} with the RIP as a function which maps normalized vectors $\mathbf{h}/\|\mathbf{h}\|$ into a thin layer around the m -dimensional unit hypersphere. The layer's thickness determines the restricted isometry constant so that in the best case, all $\mathbf{A}\mathbf{h}/\|\mathbf{h}\|$ lie exactly on the m -dimensional unit hypersphere (restricted isometry constant equal to zero). In addition to requiring this hypersphere mapping property, we argue that it is desirable to distribute the points $\mathbf{A}\mathbf{h}/\|\mathbf{h}\|$ isotropically around the origin to combat the influence of the noise. Assuming that all \mathbf{h} are realizations of a random variable \mathbf{h} , in the best case, the random variable $\mathbf{A}\mathbf{h}/\|\mathbf{h}\|$ is thus uniformly distributed on the m -dimensional unit hypersphere. This leads us to propose an optimization problem which aims to determine the matrix \mathbf{A} such that the distance between the distribution of $\mathbf{A}\mathbf{h}/\|\mathbf{h}\|$ and the uniform distribution on the m -dimensional unit hypersphere is minimized. We use a kernel-based maximum mean discrepancy (MMD) [2] to measure the distance between the two distributions and to formulate the distribution matching optimization problem as a learning problem. Finally, we incorporate the constraint that all entries of the matrix \mathbf{A} should have the same modulus into the MMD minimization problem. This constant modulus constraint on \mathbf{A} is motivated by one of the applications considered in the dissertation.

1.3 Motivation and Application

The main motivation for studying the described problems and the application presented in the numerical experiments is channel estimation [3]. In this context, the vector \mathbf{y} is called observation and the goal is to estimate the channel \mathbf{h} . The observation corresponds, e.g., to the receive signal at a base station with a certain number of antennas. The matrix \mathbf{A} is called observation matrix and it can, for example, represent (the positions of) pilot symbols. It can also represent analog phase shifters which can, for example, be used in order to connect a large number of antennas to a smaller number of receiver chains. In this case, the matrix \mathbf{A} is a wide matrix where the number of rows corresponds to the number of receiver chains and the number of columns corresponds to the number of antennas. All elements of a phase shifters matrix have the same constant modulus. This case is the main motivation for designing such a matrix in Chapter 4.

Preliminaries 2

The notation and concepts introduced here are used in the main chapters.

2.1 Notation

We write $\mathcal{N}(\boldsymbol{\mu}, \mathbf{C})$ for the probability density function (PDF) of a real-valued Gaussian random variable with mean vector $\boldsymbol{\mu}$ and covariance matrix \mathbf{C} . We write $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \mathbf{C})$ to evaluate this density at a suitable vector \mathbf{x} :

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \mathbf{C}) = \frac{1}{\sqrt{\det(2\pi\mathbf{C})}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \mathbf{C}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right). \quad (2.1)$$

Analogously, we use $\mathcal{N}_{\mathbb{C}}$ to denote the PDF of a complex-valued circularly-symmetric Gaussian random variable. Generally, we write $z \sim p$ to indicate that z is a random variable with PDF p .

The supremum norm of a continuous function $f : \mathbb{K}^N \rightarrow \mathbb{K}$ is $\|f\|_{\infty} = \sup_{\mathbf{x} \in \mathbb{K}^N} |f(\mathbf{x})|$ with $\mathbb{K} \in \{\mathbb{R}, \mathbb{C}\}$. A sequence $(f^{(K)})_{K=1}^{\infty}$ of continuous functions $f^{(K)} : \mathbb{K}^N \rightarrow \mathbb{K}$ is said to converge uniformly to $f : \mathbb{K}^N \rightarrow \mathbb{K}$ if

$$\lim_{K \rightarrow \infty} \|f - f^{(K)}\|_{\infty} = 0 \quad (2.2)$$

holds. For a vector $\mathbf{d} \in \mathbb{K}^N$, $\text{diag}(\mathbf{d}) \in \mathbb{K}^{N \times N}$ is the $N \times N$ matrix with \mathbf{d} on its main diagonal and zeros elsewhere. For a matrix $\mathbf{D} \in \mathbb{K}^{N \times N}$, $\text{diag}(\mathbf{D}) \in \mathbb{K}^N$ is the vector on its main diagonal. Further, $\|\cdot\|$ denotes the Euclidean norm, $\mathbf{0} \in \mathbb{K}^N$ or $\mathbf{0}_N$ denotes the zero vector in \mathbb{K}^N , $\mathbf{I} \in \mathbb{K}^{N \times N}$ or \mathbf{I}_N is the identity matrix in $\mathbb{K}^{N \times N}$, \mathbf{A}^H denotes the adjoint (conjugate transpose) of a matrix \mathbf{A} , and $\mathbf{A} \otimes \mathbf{B}$ denotes the Kronecker product of two matrices \mathbf{A} and \mathbf{B} . Lastly, $\text{vec}(\mathbf{A}) \in \mathbb{K}^{mN}$ denotes stacking the columns of a matrix $\mathbf{A} \in \mathbb{K}^{m \times N}$ into a long vector, and j is the imaginary unit.

2.2 Signal Model

The signal model considered throughout the dissertation is

$$\mathbf{y} = \mathbf{A}\mathbf{h} + \mathbf{n} \quad (2.3)$$

where $\mathbf{y} \in \mathbb{C}^m$ is called the *observation*, $\mathbf{A} \in \mathbb{C}^{m \times N}$ is called the *observation matrix*, $\mathbf{h} \in \mathbb{C}^N$ is called the *channel*, and $\mathbf{n} \in \mathbb{C}^m$ is called the *noise*. The noise is assumed to be a realization of a Gaussian random variable $\mathbf{n} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{0}, \mathbf{\Sigma})$ with mean vector $\mathbf{0} \in \mathbb{C}^m$ and covariance matrix $\mathbf{\Sigma} \in \mathbb{C}^{m \times m}$.

One topic of interest in the context of the signal model (2.3) is to design a channel estimation algorithm. There, we assume that the observation \mathbf{y} , the observation matrix \mathbf{A} , and the noise PDF are given, and the goal is to obtain an estimate of the channel \mathbf{h} . We propose such an algorithm in Chapter 3. Another topic of interest is to design the observation matrix \mathbf{A} . There, one goal could be to design \mathbf{A} such that it harmonizes well with a certain channel estimation algorithm. A particular challenge when designing \mathbf{A} can be to incorporate structural matrix constraints. In Chapter 4, we propose an algorithm to find an observation matrix which can harmonize with compressive sensing channel estimation algorithms and which is constrained to have constant modulus entries. In contrast to this, the standard procedure in compressive sensing theory is to draw \mathbf{A} randomly, see Section 2.3 for details. Advantages of designed matrices can be a reduced computational complexity of the corresponding channel estimation algorithms, an improved estimation performance as compared to random matrices, or better compatibility for analog implementations.

2.2.1 Special Cases

The receive signal

$$\mathbf{y} = \mathbf{h} + \mathbf{n} \quad (2.4)$$

corresponding to a single-antenna mobile terminal which transmits a pilot to an N -antenna base station is a special case of the signal model in (2.3), see, e.g., [3, 4]. This case is interesting for us because potential channel estimation algorithms can be studied without having to consider the difficulty of finding a suitable observation matrix. Another reason is that the observation matrix (the identity matrix) is invertible which addresses a special case considered in Chapter 3.

The multiple-input multiple-output (MIMO) signal model

$$\mathbf{Y} = \mathbf{H}\mathbf{P} + \mathbf{N} \quad (2.5)$$

with the channel matrix $\mathbf{H} \in \mathbb{C}^{N_{\text{rx}} \times N_{\text{tx}}}$, the pilot matrix $\mathbf{P} \in \mathbb{C}^{N_{\text{tx}} \times N_p}$, and the noise matrix $\mathbf{N} \in \mathbb{C}^{N_{\text{rx}} \times N_p}$, where an N_{tx} -antenna mobile terminal transmits N_p pilots to an N_{rx} -antenna base station (see, e.g., [4]), is another special case of (2.3). This can be seen using the definitions $\mathbf{h} = \text{vec}(\mathbf{H})$, $\mathbf{y} = \text{vec}(\mathbf{Y})$, $\mathbf{n} = \text{vec}(\mathbf{N})$, and $\mathbf{A} = \mathbf{P}^T \otimes \mathbf{I}_{N_{\text{rx}}}$ such that we have $m = N_{\text{rx}}N_p$ and $N = N_{\text{tx}}N_{\text{rx}}$.

Another special case is the single-input single-output (SISO) transmission over a frequency-selective fading channel $\mathbf{H} \in \mathbb{C}^{N_c \times N_t}$ where N_c is the number of subcarriers and N_t is the number of time slots. If only N_p of the $N_c N_t$ resource elements are occupied with pilots, then these can be selected via a matrix $\mathbf{A} \in \{0, 1\}^{N_p \times N_c N_t}$ which represents the pilot positions. Defining $\mathbf{h} = \text{vec}(\mathbf{H}) \in \mathbb{C}^{N_c N_t}$, the observed signal is given by (2.3) with $m = N_p$ and $N = N_c N_t$.

The cases mentioned so far are studied in Chapter 3. The next case is what we are concerned with in Chapter 4. There, we consider the setting where a single-antenna mobile terminal transmits a pilot to an N -antenna base station which does not have a dedicated receiver chain for every antenna. Instead, the N antennas are connected to $m < N$ receiver chains via an analog mixing network which is represented by a matrix $\mathbf{A} \in \mathbb{C}^{m \times N}$ such that $\mathbf{y} = \mathbf{A}\mathbf{h} + \mathbf{n}$ is the base station's receive signal, see, e.g., [5–9]. The analog mixing network is assumed to consist of phase shifters so that every entry of \mathbf{A} has the form

$$[\mathbf{A}]_{k,l} = \frac{1}{\sqrt{m}} \exp(j \phi_{k,l}) \quad (2.6)$$

with a phase $\phi_{k,l} \in \mathbb{R}$. We refer to such a matrix as a *constant modulus matrix* because all entries have a constant modulus constraint: $|[\mathbf{A}]_{k,l}| = \frac{1}{\sqrt{m}}$. Such a setting is often approached from a compressive sensing perspective.

2.3 Compressive Sensing

One part of compressive sensing theory studies the problem of recovering the channel \mathbf{h} from an observation of the form (2.3). There exist recovery guarantees if the matrix has the restricted isometry property (RIP) [10]. Assuming the channel lies in an (infinite) union of subspaces or some abstract subset $\mathcal{H} \subset \mathbb{C}^N$ (see, e.g., [11–13]), the observation matrix \mathbf{A} is said to have the RIP if the restricted isometry constant $\delta \geq 0$ is small in

$$(1 - \delta)\|\mathbf{h}\|^2 \leq \|\mathbf{A}\mathbf{h}\|^2 \leq (1 + \delta)\|\mathbf{h}\|^2 \quad (2.7)$$

for all $\mathbf{h} \in \mathcal{H}$.

In the most prominent case, \mathcal{H} is the set of all p -sparse vectors $\mathbf{s} \in \mathbb{C}^N$ (p nonzero entries). If \mathbf{h} is not sparse itself, it may be sparsely represented in a basis $\Psi \in \mathbb{C}^{N \times N}$ such that $\mathbf{h} = \Psi\mathbf{s}$ holds. Writing

$$\mathbf{y} = \mathbf{A}\Psi\mathbf{s} + \mathbf{n} \quad (2.8)$$

Algorithm 1 Orthogonal Matching Pursuit (OMP) [10]**Require:** matrix $\mathbf{C} \in \mathbb{C}^{m \times L}$, observation \mathbf{y} , sparsity p

- 1: $S^{(0)} \leftarrow \emptyset, \mathbf{s}^{(0)} \leftarrow \mathbf{0}$
- 2: **for** $i = 1$ to p **do**
- 3: $j^* \leftarrow \arg \max_{j \in \{1, \dots, L\}} \{[\mathbf{C}^H(\mathbf{y} - \mathbf{C}\mathbf{s}^{(i-1)})]_j\}$
- 4: $S^{(i)} \leftarrow S^{(i-1)} \cup \{j^*\}$
- 5: $\mathbf{s}^{(i)} \leftarrow \arg \min_{\tilde{\mathbf{s}}} \{\|\mathbf{y} - \mathbf{C}\tilde{\mathbf{s}}\|, \text{support}(\tilde{\mathbf{s}}) \subset S^{(i)}\}$
- 6: **end for**
- 7: **return** $\mathbf{s}^{(p)}$ // p -sparse vector

the greedy recovery algorithm orthogonal matching pursuit (OMP) [14–16] can be used to obtain an estimate $\hat{\mathbf{s}}$ of \mathbf{s} and the channel is then estimated as $\hat{\mathbf{h}} = \mathbf{\Psi}\hat{\mathbf{s}}$. OMP is summarized in Algorithm 1 and would be used with $\mathbf{C} = \mathbf{A}\mathbf{\Psi}$ here.

Random matrices provide a well-known source for matrices with the RIP. For example, a Gaussian random matrix where the entries are drawn independently with $[\mathbf{A}]_{k,l} \sim \mathcal{N}_{\mathbb{C}}(0, \frac{1}{\sqrt{m}})$ is known to have the RIP with high probability [10]. Furthermore, bounded random variables are sub-Gaussian random variables [10] and because sub-Gaussian random matrices have the RIP with high probability [10], there also exist constant modulus matrices with the RIP. To this end, the entries of \mathbf{A} are drawn independently and have the form

$$\frac{1}{\sqrt{m}} \exp(j\phi) \quad \text{with} \quad \phi \sim \mathcal{U}([0, 2\pi]) \quad (2.9)$$

where ϕ is uniformly distributed in the interval $[0, 2\pi]$.

2.4 Gaussian Mixture Models

A Gaussian mixture model (GMM) is a PDF of the form [17]

$$f^{(K)} : \mathbb{R}^N \rightarrow \mathbb{R}, \mathbf{x} \mapsto f^{(K)}(\mathbf{x}) = \sum_{k=1}^K p(k) \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \mathbf{C}_k). \quad (2.10)$$

The summands represent the K GMM *components*. Each component consists of a *mixing coefficient* $p(k)$ with $\sum_{k=1}^K p(k) = 1$ and of a Gaussian PDF with mean vector $\boldsymbol{\mu}_k \in \mathbb{R}^N$ and covariance matrix $\mathbf{C}_k \in \mathbb{R}^{N \times N}$. Given a realization $\mathbf{x} \in \mathbb{R}^N$, the probability $p(k | \mathbf{x})$ of it stemming from component k can be computed [17]:

$$p(k | \mathbf{x}) = \frac{p(k) \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \mathbf{C}_k)}{\sum_{i=1}^K p(i) \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_i, \mathbf{C}_i)}. \quad (2.11)$$

Being able to calculate these *responsibilities* (2.11) plays an important role in Chapter 3.

In Chapter 3, we are also interested in approximating an unknown continuous PDF f via a K -component GMM. To this end, all mixing coefficients $\{p(k)\}_{k=1}^K$, all mean vectors $\{\boldsymbol{\mu}_k\}_{k=1}^K$, and all covariance matrices $\{\mathbf{C}_k\}_{k=1}^K$ of the K -component GMM need to be determined. Maximum likelihood estimates of all these parameters can be obtained using an expectation-maximization (EM) algorithm and a data set $\{\mathbf{h}_t\}_{t=1}^{T_{\text{tr}}}$ of realizations \mathbf{h}_t of a random variable $\mathbf{h} \sim f$. Such a data set is assumed to be given. Details about the GMM fitting process can be found, e.g., in [17]. A list of convergence properties of the EM algorithm can be found, e.g., in [18].

The motivation for approximating the unknown PDF f via a GMM comes from [1, Theorem 5] according to which GMMs are able to approximate any continuous PDF arbitrarily well. We write

$$\mathcal{C} = \left\{ f : \mathbb{R}^N \rightarrow \mathbb{R} : f \geq 0, \int f(\mathbf{x}) d\mathbf{x} = 1, f \text{ is continuous} \right\} \quad (2.12)$$

for the set of all continuous PDFs. Then, one part of [1, Theorem 5] says that for any $f \in \mathcal{C}$, there exists a sequence of GMMs which converges uniformly to f :

Theorem 1 (Universal Approximation Property). *Denoting the standard Gaussian density by $\mathcal{N}(\mathbf{0}, \mathbf{I})$, let*

$$\mathcal{M}_K = \left\{ h : h(\mathbf{x}) = \sum_{k=1}^K c_k \frac{1}{\sigma_k^N} \mathcal{N}\left(\frac{\mathbf{x} - \boldsymbol{\mu}_k}{\sigma_k}; \mathbf{0}, \mathbf{I}\right) \right\}$$

be the class of K -component location-scale finite Gaussian mixtures with

$$\boldsymbol{\mu}_k \in \mathbb{R}^N, \sigma_k > 0, c_k \geq 0 \quad \text{for all } k \in \{1, \dots, K\}$$

and $\sum_{k=1}^K c_k = 1$. For any $f \in \mathcal{C}$, there exists a sequence $(f^{(K)})_{K=1}^{\infty}$ with $f^{(K)} \in \mathcal{M}_K$ which converges uniformly to f .

2.5 Maximum Mean Discrepancy

Maximum mean discrepancy (MMD) can be used to compute a distance between probability distributions and is used in Chapter 4. If \mathcal{X} is a metric space and \mathbf{x} and \mathbf{y} are two random variables with respective probability distributions p and q defined on \mathcal{X} , MMD quite literally computes a maximization of the discrepancy between two means [2, Definition 2]:

$$\text{MMD}_{\mathcal{F}}(p, q) = \sup_{f \in \mathcal{F}} (\mathbb{E}_{\mathbf{x} \sim p}[f(\mathbf{x})] - \mathbb{E}_{\mathbf{y} \sim q}[f(\mathbf{y})]). \quad (2.13)$$

Here, \mathcal{F} is a suitably chosen set of functions $f : \mathcal{X} \rightarrow \mathbb{R}$. Depending on \mathcal{F} , (2.13) can be a metric and then we have $\text{MMD}_{\mathcal{F}}(p, q) = 0$ if and only if $p = q$, and $\text{MMD}_{\mathcal{F}}(p, q) > 0$ otherwise. An overview of suitable function sets \mathcal{F} can be found, e.g., in [19].

Solving the optimization problem (2.13) can be challenging, depending on \mathcal{F} . One particularly interesting set \mathcal{F} in this regard is the unit ball in a reproducing kernel Hilbert space. A brief introduction to reproducing kernel Hilbert spaces can be found in [2]. Every such Hilbert space is associated with a kernel $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$. Conversely, every kernel defines an associated Hilbert space. Intuitively speaking, a kernel k can represent all functions which are an element of the reproducing kernel Hilbert space. For this reason, the optimization with respect to the whole set \mathcal{F} in (2.13) can be expressed using the kernel only and the computation of the supremum is avoided.

Using a suitable kernel k , (2.13) can be expressed as [2, Lemma 6]

$$\text{MMD}_k^2(p, q) = \mathbb{E}_{\mathbf{x}, \mathbf{x}'}[k(\mathbf{x}, \mathbf{x}')] - 2 \mathbb{E}_{\mathbf{x}, \mathbf{y}}[k(\mathbf{x}, \mathbf{y})] + \mathbb{E}_{\mathbf{y}, \mathbf{y}'}[k(\mathbf{y}, \mathbf{y}')]. \quad (2.14)$$

Here, \mathbf{x} and \mathbf{x}' are independent random variables with distribution p , and \mathbf{y} and \mathbf{y}' are independent random variables with distribution q . Among the best known kernels is the Gaussian kernel:

$$k_{\sigma} : \mathbb{R}^l \times \mathbb{R}^l \rightarrow \mathbb{R}, (\mathbf{x}, \mathbf{y}) \mapsto k_{\sigma}(\mathbf{x}, \mathbf{y}) = e^{-\frac{1}{2\sigma^2} \|\mathbf{x} - \mathbf{y}\|^2} \quad (2.15)$$

with a parameter $\sigma > 0$. The Gaussian kernel is a typical choice when MMD_k is used, as, e.g., [20–22] demonstrate. One reason is that the associated reproducing kernel Hilbert space enjoys the property of being universal which makes (2.14) a metric [2, Theorem 5].

Compared to (2.13), computing (2.14) no longer involves solving an optimization problem. It is however still necessary to know both p and q . In some applications (in particular in Chapter 4), only realizations $\{\mathbf{x}_i\}_{i=1}^M$ and $\{\mathbf{y}_j\}_{j=1}^N$ of $\mathbf{x} \sim p$ and $\mathbf{y} \sim q$, respectively, are given and we are still interested in determining whether $p = q$ or whether p and q are at least close to each other. In this situation, a biased estimate of (2.14) can be computed [2, Equation (6)]:

$$\begin{aligned} \text{MMD}_k^2(\{\mathbf{x}_i\}_{i=1}^M, \{\mathbf{y}_j\}_{j=1}^N) = \\ \sum_{i=1}^M \sum_{j=1}^M \frac{k(\mathbf{x}_i, \mathbf{x}_j)}{M^2} - 2 \sum_{i=1}^M \sum_{j=1}^N \frac{k(\mathbf{x}_i, \mathbf{y}_j)}{MN} + \sum_{i=1}^N \sum_{j=1}^N \frac{k(\mathbf{y}_i, \mathbf{y}_j)}{N^2}. \end{aligned} \quad (2.16)$$

The smaller the value of $\text{MMD}_k^2(\{\mathbf{x}_i\}_{i=1}^M, \{\mathbf{y}_j\}_{j=1}^N)$, the closer p and q are to each other. One main result of [2] is that (2.16) converges in probability to (2.14) at a rate of $\mathcal{O}((M + N)^{-\frac{1}{2}})$.

2.6 Distribution Matching

The biased MMD estimate (2.16) is differentiable because the Gaussian kernel (2.15) is differentiable. In the literature, this fact is utilized in the context of machine learning applications, e.g., [20–23]. One possible application, which is related to the considerations in Chapter 4, can be described as follows. A set $\{\mathbf{x}_i\}_{i=1}^M$ of natural images \mathbf{x}_i is given and it is assumed that they are all samples of a random variable \mathbf{x} with distribution p . The goal is to draw new samples of \mathbf{x} , or, in other words, to generate new natural images. To this end, a set $\{g_\phi\}_{\phi \in \Phi}$ of functions g_ϕ is defined where every g_ϕ is differentiable with respect to the parameter ϕ and where Φ is a constraint set for the parameter. For instance, every g_ϕ could be a neural network where ϕ collects all learnable parameters (all weights and biases). Then, an auxiliary latent random variable \mathbf{y}' with the distribution q' is chosen (e.g., a standard Gaussian random variable) and corresponding random variables $g_\phi(\mathbf{y}')$ with respective distributions q_ϕ are defined. Typically, \mathbf{y}' is chosen such that sampling this random variable is simple. The goal is to find $\phi^* \in \Phi$ such that the distribution q_{ϕ^*} is close (ideally equal) to p . This amounts to solving an optimization problem:

$$\text{MMD}_k^2(p, q_{\phi^*}) = \min_{\phi \in \Phi} \text{MMD}_k^2(p, q_\phi). \quad (2.17)$$

New images are now generated by first drawing samples $\mathbf{y}'_1, \mathbf{y}'_2, \dots$ of \mathbf{y}' , which is simple, and computing the images $g_{\phi^*}(\mathbf{y}'_1), g_{\phi^*}(\mathbf{y}'_2), \dots$, which are then samples of the distribution q_{ϕ^*} .

Since p is not known and q_ϕ is generally difficult to compute, (2.17) cannot be solved directly. However, if samples (training data) $\{\mathbf{x}_i\}_{i=1}^M$ and $\{\mathbf{y}'_j\}_{j=1}^N$ are given, (2.16) can be used to define an alternative optimization problem:

$$\min_{\phi \in \Phi} \text{MMD}_k^2(\{\mathbf{x}_i\}_{i=1}^M, \{g_\phi(\mathbf{y}'_j)\}_{j=1}^N). \quad (2.18)$$

A problem of this form is considered in Chapter 4. Such optimization problems are solved via stochastic gradient descent, e.g., [20–23], because the objective function is differentiable, and are considered easy to optimize [23]. Convergence guarantees are part of ongoing research. Solving an optimization problem via stochastic gradient descent often involves some form of a grid or a random search. This is discussed next.

2.7 Random Search

For the following explanation of the grid and the random search procedures [24], we assume we want to train a classical feed-forward neural network using a stochastic gradient descent algorithm. We distinguish between the optimization parameters of the neural network and the hyperparameters of the neural network and of the stochastic gradient descent algorithm. The values of the optimization parameters are derived during the training (or learning) process whereas the hyperparameters are set before the training process begins. Weights and biases constitute the neural network's optimization parameters. Examples of hyperparameters are the number of hidden layers of the neural network, the stochastic gradient descent algorithm's learning rate (i.e., the gradient step size), or the batch size (i.e., the number of samples used to compute the stochastic gradient). The optimization parameters are typically randomly initialized (e.g., from a uniform or a normal distribution) and the hyperparameters are typically determined via a grid or a random search procedure [24].

To perform a grid search, we determine sets of possible hyperparameter values, for example, $P_{hl} = \{5, 10, 25, 50\}$ for the number of hidden layers, $P_{lr} = \{10^{-5}, 10^{-3}, 10^{-1}\}$ for the learning rate, and $P_{bs} = \{10, 50, 100, 250, 500\}$ for the batch size. Then, we train the neural network once for every combination of the hyperparameter values, i.e., for every element of the Cartesian product $P_{hl} \times P_{lr} \times P_{bs}$. In this example, we would therefore train the neural network $4 \cdot 3 \cdot 5 = 60$ times. In the end, a held out validation data set is used to choose the best of the 60 trained neural networks.

The authors of [24] argue that the hyperparameters space can be explored more efficiently if a random search procedure is employed. To this end, we would also train the neural network, e.g., 60 times but this time the hyperparameter values would be randomly drawn for every neural network. If $\text{uniform}([a, b])$ denotes drawing a sample from a continuous uniform distribution on the interval $[a, b]$, a hyperparameter value can be determined via

$$e^{\text{uniform}([\ln(a), \ln(b)])} \tag{2.19}$$

followed by rounding to the nearest integer if necessary [24]. Using the numerical example for the learning rate from before, we would, e.g., set $a = 10^{-5}$ and $b = 10^{-1}$ to determine the learning rate as $e^{\text{uniform}([\ln(10^{-5}), \ln(10^{-1})])}$ for a given neural network. Here, too, a held out validation data set is used to choose the best of the 60 trained neural networks.

Asymptotically Optimal Channel Estimation 3

This chapter is mainly based on our work in [25, 26]. In this chapter, we distinguish between a random variable \mathbf{x} and a realization \mathbf{x} thereof.

3.1 A Sequence of Conditional Mean Estimators

We write $f_{\mathbf{h}}$ for the probability density function (PDF) of the channel random variable \mathbf{h} . Further, we write $f_{\mathbf{n}} = \mathcal{N}_{\mathbb{C}}(\mathbf{0}, \mathbf{\Sigma})$ for the noise PDF and $f_{\mathbf{y}}$ for the observation PDF in the model

$$\mathbf{y} = \mathbf{A}\mathbf{h} + \mathbf{n} \quad (3.1)$$

with $\mathbf{A} \in \mathbb{C}^{m \times N}$.

A good channel estimator is a function $\tilde{\mathbf{h}} : \mathbb{C}^m \rightarrow \mathbb{C}^N$ which leads to a small mean square error (MSE)

$$\mathbb{E} \left[\left\| \mathbf{h} - \tilde{\mathbf{h}}(\mathbf{y}) \right\|^2 \right]. \quad (3.2)$$

It is known that a minimizer is given by the conditional mean estimator (CME) [27]

$$\hat{\mathbf{h}} : \mathbb{C}^m \rightarrow \mathbb{C}^N, \mathbf{y} \mapsto \hat{\mathbf{h}}(\mathbf{y}) = \mathbb{E}[\mathbf{h} | \mathbf{y} = \mathbf{y}] = \int \mathbf{h} f_{\mathbf{h}|\mathbf{y}}(\mathbf{h}|\mathbf{y}) d\mathbf{h} \quad (3.3)$$

where $f_{\mathbf{h}|\mathbf{y}}(\mathbf{h} | \mathbf{y})$ denotes the conditional PDF of \mathbf{h} given $\mathbf{y} = \mathbf{y}$ evaluated at \mathbf{h} . For our purposes, we first rewrite the conditional PDF via

$$f_{\mathbf{h}|\mathbf{y}}(\mathbf{h} | \mathbf{y}) = \frac{f_{\mathbf{y}|\mathbf{h}}(\mathbf{y} | \mathbf{h}) f_{\mathbf{h}}(\mathbf{h})}{f_{\mathbf{y}}(\mathbf{y})} = \frac{f_{\mathbf{n}}(\mathbf{y} - \mathbf{A}\mathbf{h}) f_{\mathbf{h}}(\mathbf{h})}{f_{\mathbf{y}}(\mathbf{y})} \quad (3.4)$$

and the CME is then given by

$$\hat{\mathbf{h}} : \mathbf{y} \mapsto \hat{\mathbf{h}}(\mathbf{y}) = \mathbb{E}[\mathbf{h} | \mathbf{y} = \mathbf{y}] = \int \mathbf{h} \frac{f_{\mathbf{n}}(\mathbf{y} - \mathbf{A}\mathbf{h}) f_{\mathbf{h}}(\mathbf{h})}{f_{\mathbf{y}}(\mathbf{y})} d\mathbf{h}. \quad (3.5)$$

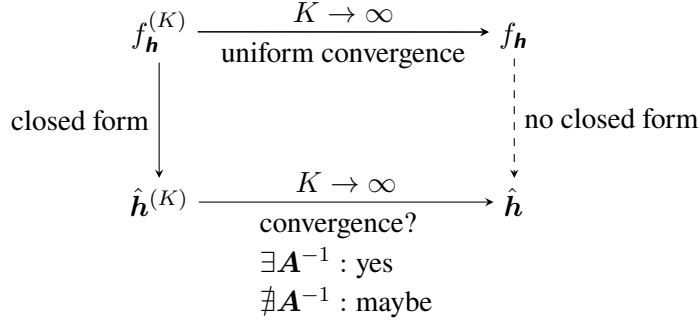


Figure 3.1: In Chapter 3, we discuss the following question: If a sequence $(f_{\mathbf{h}}^{(K)})_{K=1}^{\infty}$ of probability density functions converges uniformly to the channel probability density function $f_{\mathbf{h}}$, does the corresponding sequence $(\hat{\mathbf{h}}^{(K)})_{K=1}^{\infty}$ of conditional mean estimators converge to the channel conditional mean estimator $\hat{\mathbf{h}}$?

Computing (3.5) analytically can be difficult. On top of this, even if all involved PDFs are given analytically, there may not exist an analytic expression of the integral and therefore of the CME $\hat{\mathbf{h}}$. Another problem might be that only the noise PDF $f_{\mathbf{n}}$ —but not $f_{\mathbf{h}}$ —is known, which is the assumption in this dissertation, and then, computing (3.5) is not possible. However, even if the channel PDF $f_{\mathbf{h}}$ is unknown, we can reasonably assume to have access to a sequence $(f_{\mathbf{h}}^{(K)})_{K=1}^{\infty}$ of PDFs which converges uniformly to $f_{\mathbf{h}}$:

$$\lim_{K \rightarrow \infty} \|f_{\mathbf{h}} - f_{\mathbf{h}}^{(K)}\|_{\infty} = 0. \quad (3.6)$$

A prime example is the case where $(f_{\mathbf{h}}^{(K)})_{K=1}^{\infty}$ is a sequence of Gaussian mixture models (GMMs). This case is studied in Section 3.3.

For every $K \in \mathbb{N}$, we now introduce an auxiliary channel random variable $\mathbf{h}^{(K)}$ with PDF $f_{\mathbf{h}}^{(K)}$, we introduce an auxiliary signal model

$$\mathbf{y}^{(K)} = \mathbf{A}\mathbf{h}^{(K)} + \mathbf{n} \quad (3.7)$$

and we write $f_{\mathbf{y}}^{(K)}$ for the PDF of $\mathbf{y}^{(K)}$. For every $K \in \mathbb{N}$, the CME

$$\hat{\mathbf{h}}^{(K)} : \mathbf{y} \mapsto \hat{\mathbf{h}}^{(K)}(\mathbf{y}) = \mathbb{E}[\mathbf{h}^{(K)} \mid \mathbf{y}^{(K)} = \mathbf{y}] = \int \mathbf{h} \frac{f_{\mathbf{n}}(\mathbf{y} - \mathbf{A}\mathbf{h})f_{\mathbf{h}}^{(K)}(\mathbf{h})}{f_{\mathbf{y}}^{(K)}(\mathbf{y})} d\mathbf{h} \quad (3.8)$$

is then an optimal channel estimator for the model (3.7). In the case where $f_{\mathbf{h}}^{(K)}$ is a GMM, $\hat{\mathbf{h}}^{(K)}$ can be calculated analytically (see Section 3.3). Generally, if the sequence

$(f_{\mathbf{h}}^{(K)})_{K=1}^{\infty}$ is chosen such that every CME $\hat{\mathbf{h}}^{(K)}$ can be computed analytically, it is interesting to investigate whether $\hat{\mathbf{h}}^{(K)}$ can be used as an approximation of the optimal $\hat{\mathbf{h}}$, which is not available analytically. That is, we ask whether the uniform convergence of the PDFs in (3.6) implies that the sequence $(\hat{\mathbf{h}}^{(K)})_{K=1}^{\infty}$ converges to $\hat{\mathbf{h}}$ in some suitable sense. A main result in this chapter is that we have at least pointwise convergence if the observation matrix \mathbf{A} is invertible, see Theorem 2. Figure 3.1 schematically summarizes this result.

3.2 Pointwise Convergence

The PDF of a complex-valued random variable can be written in terms of its real and imaginary parts. For this reason, this section (and only this section) considers real-valued quantities only and the results apply to the setting discussed in Section 3.1 by stacking real and imaginary parts appropriately. The proof of the following theorem can be found in Appendix A. The theorem uses the notation from Section 3.1.

Theorem 2. *Let $\mathcal{C} = \{f : \mathbb{R}^N \rightarrow \mathbb{R} : f \geq 0, \int f(\mathbf{x})d\mathbf{x} = 1, f \text{ is continuous}\}$ be the set of all continuous PDFs. Let $f_{\mathbf{h}} \in \mathcal{C}$ and $f_{\mathbf{n}} = \mathcal{N}(\mathbf{0}, \mathbf{\Sigma}) \in \mathcal{C}$. Further, let $\mathbf{A} \in \mathbb{R}^{N \times N}$ be invertible and let $(f_{\mathbf{h}}^{(K)})_{K=1}^{\infty}$ be a sequence of PDFs in \mathcal{C} which converges uniformly to $f_{\mathbf{h}}$. For every $K \in \mathbb{N}$, define a function $\hat{\mathbf{h}}^{(K)} : \mathbb{R}^N \rightarrow \mathbb{R}^N$ via*

$$\hat{\mathbf{h}}^{(K)} : \mathbf{y} \mapsto \hat{\mathbf{h}}^{(K)}(\mathbf{y}) = \int \mathbf{h} \frac{f_{\mathbf{n}}(\mathbf{y} - \mathbf{A}\mathbf{h})f_{\mathbf{h}}^{(K)}(\mathbf{h})}{f_{\mathbf{y}}^{(K)}(\mathbf{y})} d\mathbf{h}. \quad (3.9)$$

Then, the sequence $(\hat{\mathbf{h}}^{(K)})_{K=1}^{\infty}$ converges pointwise to $\hat{\mathbf{h}} : \mathbb{R}^N \rightarrow \mathbb{R}^N$ with

$$\hat{\mathbf{h}} : \mathbf{y} \mapsto \hat{\mathbf{h}}(\mathbf{y}) = \int \mathbf{h} \frac{f_{\mathbf{n}}(\mathbf{y} - \mathbf{A}\mathbf{h})f_{\mathbf{h}}(\mathbf{h})}{f_{\mathbf{y}}(\mathbf{y})} d\mathbf{h}. \quad (3.10)$$

That is, we have

$$\lim_{K \rightarrow \infty} \|\hat{\mathbf{h}}(\mathbf{y}) - \hat{\mathbf{h}}^{(K)}(\mathbf{y})\| = 0 \quad (3.11)$$

for every $\mathbf{y} \in \mathbb{R}^N$. Further, (3.11) holds uniformly for all $\mathbf{y} \in \mathcal{B}_r$ with $\mathcal{B}_r = \{\mathbf{y} \in \mathbb{R}^N : \|\mathbf{y}\| \leq r\}$ for any radius $r > 0$.

3.2.1 Discussion of Theorem 2

As a result of Theorem 2, the channel estimator $\hat{\mathbf{h}}^{(K)}$ can be viewed as an approximation of the optimal channel estimator $\hat{\mathbf{h}}$ for a $K \in \mathbb{N}$. Let $r > 0$ be given and let $\mathcal{B}_r = \{\mathbf{y} \in \mathbb{R}^N : \|\mathbf{y}\| \leq r\}$. According to Theorem 2, if we want the approximation

error to be bounded by a given threshold $\varepsilon_{\text{thr}} > 0$, there exists a $K_r \in \mathbb{N}$ which guarantees

$$\|\hat{\mathbf{h}}(\mathbf{y}) - \hat{\mathbf{h}}^{(K)}(\mathbf{y})\| \leq \varepsilon_{\text{thr}} \quad \text{for all } \mathbf{y} \in \mathcal{B}_r \quad (3.12)$$

if we choose $K \geq K_r$. This does not mean that the approximation error exceeds ε_{thr} for $\mathbf{y} \notin \mathcal{B}_r$. Due to (3.11), for any $\mathbf{y} \in \mathbb{R}^N$, there exists a $K_{\mathbf{y}} \in \mathbb{N}$ such that the approximation error falls below ε_{thr} . However, for some $\mathbf{y} \notin \mathcal{B}_r$, the sequence index $K_{\mathbf{y}}$ needs to be strictly larger than K_r to achieve the goal. Generally, we can expect the channel estimation performance of $\hat{\mathbf{h}}^{(K)}$ to be better for larger K .

The proof of Theorem 2 relies on \mathbf{A} being invertible. If this is not the case, two problems immediately arise. These are discussed in detail in Appendix B. One problem is that the uniform convergence of $(f_{\mathbf{h}}^{(K)})_{K=1}^{\infty}$ to $f_{\mathbf{h}}$ does generally not imply the uniform convergence of $(f_{\mathbf{y}}^{(K)})_{K=1}^{\infty}$ to $f_{\mathbf{y}}$ if \mathbf{A} is not invertible (see Appendix B.2). As a consequence, a crucial step in the proof of Theorem 2 cannot be shown. Nonetheless, numerical studies in Section 3.6 suggest that even if the matrix \mathbf{A} is not invertible, the performance of the estimator $\hat{\mathbf{h}}^{(K)}$ improves with increasing K .

For a noninvertible \mathbf{A} , it seems difficult to infer a statement about the convergence of the PDFs $f_{\mathbf{y}}^{(K)}$, which correspond to $\mathbf{y}^{(K)}$, if only the uniform convergence of the PDFs $f_{\mathbf{h}}^{(K)}$, which correspond to $\mathbf{h}^{(K)}$, is given. However, we can say something about the convergence of the respective distribution functions. To this end, first observe that the uniform convergence of $(f_{\mathbf{h}}^{(K)})_{K=1}^{\infty}$ to $f_{\mathbf{h}}$ implies the pointwise convergence as well. By Scheffe's lemma (e.g., [28]), the sequence of random variables $(\mathbf{h}^{(K)})_{K=1}^{\infty}$ then converges to the random variable \mathbf{h} in distribution, i.e., we have the pointwise convergence of the sequence of distribution functions corresponding to $\mathbf{h}^{(K)}$ to the distribution function corresponding to \mathbf{h} . We denote this by $\mathbf{h}^{(K)} \xrightarrow{d} \mathbf{h}$. Since \mathbf{n} is independent of $\mathbf{h}^{(K)}$ and of \mathbf{h} , also $(\mathbf{h}^{(K)}, \mathbf{n}) \xrightarrow{d} (\mathbf{h}, \mathbf{n})$ holds. Applying the continuous mapping theorem (e.g., [29]) using the continuous mapping

$$s : \mathbb{R}^{N+m} \rightarrow \mathbb{R}^{N+m}, (\mathbf{h}, \mathbf{n}) \mapsto (\mathbf{h}, \mathbf{A}\mathbf{h} + \mathbf{n}) \quad (3.13)$$

implies $s(\mathbf{h}^{(K)}, \mathbf{n}) = (\mathbf{h}^{(K)}, \mathbf{y}^{(K)}) \xrightarrow{d} s(\mathbf{h}, \mathbf{n}) = (\mathbf{h}, \mathbf{y})$. In other words, with $(\mathbf{h}^{(K)}, \mathbf{y}^{(K)})$ we have an approximation of both the channel \mathbf{h} and the observation \mathbf{y} . The author of [30] investigates whether in this case $\mathbb{E}[\mathbf{h}^{(K)} | \mathbf{y}^{(K)}]$ is an approximation of $\mathbb{E}[\mathbf{h} | \mathbf{y}]$. This is answered affirmatively for suitable distributions. Since we assume $f_{\mathbf{h}}$ to be unknown, a discussion of [30] is omitted. Note, $(\mathbb{E}[\mathbf{h}^{(K)} | \mathbf{y}^{(K)}])_{K=1}^{\infty}$ is a sequence of random variables and [30] provides conditions under which this sequence converges in distribution to the random variable $\mathbb{E}[\mathbf{h} | \mathbf{y}]$. In contrast, we investigate the pointwise convergence of the functions $(\mathbf{y} \mapsto \mathbb{E}[\mathbf{h}^{(K)} | \mathbf{y}^{(K)} = \mathbf{y}])_{K=1}^{\infty}$ to the

function $\mathbf{y} \mapsto \mathbb{E}[\mathbf{h} \mid \mathbf{y} = \mathbf{y}]$. Thus, Theorem 2 and [30] are related but study different mathematical objects and different types of convergence. In particular, Theorem 2 is not a consequence of [30].

We still need to address the question of how to obtain a sequence $(f_{\mathbf{h}}^{(K)})_{K=1}^{\infty}$ of PDFs which converges uniformly to $f_{\mathbf{h}}$ and how to compute the estimator $\hat{\mathbf{h}}^{(K)}$ so that it can be used for channel estimation in practice. In the remainder of this chapter, we study the example of a sequence of GMMs which shows the desired convergence behavior and leads to a closed-form expression for $\hat{\mathbf{h}}^{(K)}$ (left column in Figure 3.1).

3.3 Example: Gaussian Mixture Models

We come back to the notation introduced in Section 3.1 where all quantities are complex-valued. Theorem 2 assumes a channel PDF $f_{\mathbf{h}}$ to be given and then requires a sequence $(f_{\mathbf{h}}^{(K)})_{K=1}^{\infty}$ of PDFs which converges uniformly to $f_{\mathbf{h}}$. In light of Theorem 1 (in Section 2.4), there always exists a sequence of GMMs which fulfills that requirement. Henceforth, $(f_{\mathbf{h}}^{(K)})_{K=1}^{\infty}$ denotes a suitable sequence of GMMs. As discussed next, GMMs have the benefit that $\hat{\mathbf{h}}^{(K)}$ from (3.8) can be computed in closed form for every $K \in \mathbb{N}$ which results in a practically useful channel estimator.

3.3.1 The GMM Estimator

For a given $K \in \mathbb{N}$, we are interested in obtaining an analytic expression for $\hat{\mathbf{h}}^{(K)}$ from (3.8). Since this estimator is based on GMMs for the remainder of Chapter 3, we call it the *GMM estimator* from now on. Let

$$f_{\mathbf{h}}^{(K)} : \mathbb{C}^N \rightarrow \mathbb{R}, \mathbf{h} \mapsto f_{\mathbf{h}}^{(K)}(\mathbf{h}) = \sum_{k=1}^K p(k) \mathcal{N}_{\mathbb{C}}(\mathbf{h}; \boldsymbol{\mu}_k, \mathbf{C}_k) \quad (3.14)$$

be the K th GMM. For notational convenience, we omit the sequence index K when we write the GMM parameters so that we have $p(k), \boldsymbol{\mu}_k, \mathbf{C}_k$ instead of, e.g., $p^{(K)}(k), \boldsymbol{\mu}_k^{(K)}, \mathbf{C}_k^{(K)}$. Let \mathbf{k} be the discrete mixing variable such that $(\mathbf{h}^{(K)} \mid \mathbf{k} = k) \sim \mathcal{N}_{\mathbb{C}}(\boldsymbol{\mu}_k, \mathbf{C}_k)$ is the k th Gaussian in the GMM (3.14). Since also the noise $\mathbf{n} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{0}, \boldsymbol{\Sigma})$ is Gaussian, the observation $\mathbf{y}^{(K)} = \mathbf{A}\mathbf{h}^{(K)} + \mathbf{n}$ conditioned on $\mathbf{k} = k$ is Gaussian as well:

$$(\mathbf{y}^{(K)} \mid \mathbf{k} = k) \sim \mathcal{N}_{\mathbb{C}}(\tilde{\boldsymbol{\mu}}_k, \tilde{\mathbf{C}}_k) \quad (3.15)$$

$$\text{with } \tilde{\boldsymbol{\mu}}_k = \mathbf{A}\boldsymbol{\mu}_k \text{ and } \tilde{\mathbf{C}}_k = \mathbf{A}\mathbf{C}_k\mathbf{A}^H + \boldsymbol{\Sigma}. \quad (3.16)$$

Consequently, the PDF of $\mathbf{y}^{(K)}$ is the following GMM:

$$f_{\mathbf{y}}^{(K)} : \mathbb{C}^m \rightarrow \mathbb{R}, \mathbf{y} \mapsto f_{\mathbf{y}}^{(K)}(\mathbf{y}) = \sum_{k=1}^K p(k) \mathcal{N}_{\mathbb{C}}(\mathbf{y}; \tilde{\boldsymbol{\mu}}_k, \tilde{\mathbf{C}}_k). \quad (3.17)$$

The quantities obtained so far can now be used to compute the desired channel estimator $\hat{\mathbf{h}}^{(K)}$ from (3.8). To this end, let $\mathbf{y} \in \mathbb{C}^m$ and employ the law of total expectation:

$$\hat{\mathbf{h}}^{(K)}(\mathbf{y}) = \mathbb{E}[\mathbf{h}^{(K)} \mid \mathbf{y}^{(K)} = \mathbf{y}] \quad (3.18)$$

$$= \sum_{k=1}^K p(k \mid \mathbf{y}^{(K)} = \mathbf{y}) \mathbb{E}[\mathbf{h}^{(K)} \mid \mathbf{y}^{(K)} = \mathbf{y}, \mathbf{k} = k]. \quad (3.19)$$

The first factor in (3.19) can be identified as a responsibility (cf. (2.11)) of the GMM $f_{\mathbf{y}}^{(K)}$ and can be calculated in closed form:

$$p(k \mid \mathbf{y}^{(K)} = \mathbf{y}) = \frac{p(k) \mathcal{N}_{\mathbb{C}}(\mathbf{y}; \tilde{\boldsymbol{\mu}}_k, \tilde{\mathbf{C}}_k)}{\sum_{i=1}^K p(i) \mathcal{N}_{\mathbb{C}}(\mathbf{y}; \tilde{\boldsymbol{\mu}}_i, \tilde{\mathbf{C}}_i)}. \quad (3.20)$$

Since both $\mathbf{h}^{(K)}$ as well as $\mathbf{y}^{(K)}$ conditioned on $\mathbf{k} = k$ are Gaussian random variables, the second factor in (3.19) is the CME known as linear minimum mean square error (LMMSE) estimator and has a well-known closed-form expression. For later reference, we give this estimator the name $\hat{\mathbf{h}}_{\text{LMMSE},k}$:

$$\hat{\mathbf{h}}_{\text{LMMSE},k} : \mathbb{C}^m \rightarrow \mathbb{C}^N, \mathbf{y} \mapsto \hat{\mathbf{h}}_{\text{LMMSE},k}(\mathbf{y}) = \mathbb{E}[\mathbf{h}^{(K)} \mid \mathbf{y}^{(K)} = \mathbf{y}, \mathbf{k} = k] \quad (3.21a)$$

$$= \mathbf{C}_k \mathbf{A}^H \tilde{\mathbf{C}}_k^{-1} (\mathbf{y} - \tilde{\boldsymbol{\mu}}_k) + \boldsymbol{\mu}_k. \quad (3.21b)$$

In summary, the closed-form expression of the GMM estimator $\hat{\mathbf{h}}^{(K)}$ reads:

$$\hat{\mathbf{h}}^{(K)} : \mathbb{C}^m \rightarrow \mathbb{C}^N \quad (3.22a)$$

$$\mathbf{y} \mapsto \hat{\mathbf{h}}^{(K)}(\mathbf{y}) = \sum_{k=1}^K \frac{p(k) \mathcal{N}_{\mathbb{C}}(\mathbf{y}; \tilde{\boldsymbol{\mu}}_k, \tilde{\mathbf{C}}_k)}{\sum_{i=1}^K p(i) \mathcal{N}_{\mathbb{C}}(\mathbf{y}; \tilde{\boldsymbol{\mu}}_i, \tilde{\mathbf{C}}_i)} \hat{\mathbf{h}}_{\text{LMMSE},k}(\mathbf{y}) \quad (3.22b)$$

with $\hat{\mathbf{h}}_{\text{LMMSE},k}$ defined in (3.21).

The formula in (3.22) is known and can also be found, e.g., in [31]. Therein, the authors assume that $f_{\mathbf{h}}^{(K)}$ in (3.14) is equal to the PDF $f_{\mathbf{h}}$ of the channel \mathbf{h} (for some fixed K) so that (3.22) is the optimal channel estimator. They then study this estimator in the high signal-to-noise ratio (SNR) regime to derive suitable pilot symbols. In contrast to this, we view $f_{\mathbf{h}}^{(K)}$ as an approximation of the unknown channel PDF $f_{\mathbf{h}}$.

Algorithm 2 GMM Estimator

Offline GMM Training Phase**Require:** training data $\{\mathbf{h}_t\}_{t=1}^{T_{\text{tr}}}$, number K of GMM components

- 1: $(\{p(k)\}_{k=1}^K, \{\boldsymbol{\mu}_k\}_{k=1}^K, \{\mathbf{C}_k\}_{k=1}^K) \leftarrow \text{EM}(\{\mathbf{h}_t\}_{t=1}^{T_{\text{tr}}}, K)$ // an EM-algorithm computes all parameters of $f_{\mathbf{h}}^{(K)}$
-

Online Channel Estimation Phase**Require:** observation \mathbf{y} , observation matrix \mathbf{A} , noise matrix $\boldsymbol{\Sigma}$

- 2: $\hat{\mathbf{h}}^{(K)} \leftarrow \mathbf{0}$
 - 3: **for** $k = 1$ to K **do**
 - 4: $p(k | \mathbf{y}^{(K)} = \mathbf{y}) \leftarrow \frac{p(k)\mathcal{N}_{\mathbb{C}}(\mathbf{y}; \mathbf{A}\boldsymbol{\mu}_k, \mathbf{A}\mathbf{C}_k\mathbf{A}^H + \boldsymbol{\Sigma})}{\sum_{i=1}^K p(i)\mathcal{N}_{\mathbb{C}}(\mathbf{y}; \mathbf{A}\boldsymbol{\mu}_i, \mathbf{A}\mathbf{C}_i\mathbf{A}^H + \boldsymbol{\Sigma})}$
 - 5: $\tilde{\mathbf{h}} \leftarrow \mathbf{C}_k\mathbf{A}^H(\mathbf{A}\mathbf{C}_k\mathbf{A}^H + \boldsymbol{\Sigma})^{-1}(\mathbf{y} - \mathbf{A}\boldsymbol{\mu}_k) + \boldsymbol{\mu}_k$
 - 6: $\hat{\mathbf{h}}^{(K)} \leftarrow \hat{\mathbf{h}}^{(K)} + p(k | \mathbf{y}^{(K)} = \mathbf{y})\tilde{\mathbf{h}}$
 - 7: **end** // many quantities in the loop can be precomputed
 - 8: **return** $\hat{\mathbf{h}}^{(K)}$ // estimated channel, see (3.22)
-

The approximation becomes better as K increases and due to Theorem 2, we can view $\hat{\mathbf{h}}^{(K)}$ as an approximation of the optimal channel estimator $\hat{\mathbf{h}}$. In that sense, we study $\hat{\mathbf{h}}^{(K)}$ in the “high number of components regime”. The contribution of this chapter lies in providing a strong motivation to use $\hat{\mathbf{h}}^{(K)}$ for channel estimation even if the channel is not distributed according to $f_{\mathbf{h}}^{(K)}$.

3.3.2 Computing the GMM Estimator in Practice

At first glance, the GMM estimator $\hat{\mathbf{h}}^{(K)}$ now seems to be given analytically so that it can be used for channel estimation. However, we still need to discuss how to obtain the required parameters which appear in the formula in (3.22) in practice. To this end, we assume a set $\{\mathbf{h}_t\}_{t=1}^{T_{\text{tr}}}$ of realizations of the channel $\mathbf{h} \sim f_{\mathbf{h}}$ to be given. This training data set can then be used to fit a K -component GMM as explained in Section 2.4. The result of the fitting process are the mixing coefficients $\{p(k)\}_{k=1}^K$, the mean vectors $\{\boldsymbol{\mu}_k\}_{k=1}^K$, and the covariance matrices $\{\mathbf{C}_k\}_{k=1}^K$ of the GMM $f_{\mathbf{h}}^{(K)}$. Thereafter, the mean vectors $\{\tilde{\boldsymbol{\mu}}_k\}_{k=1}^K$ and covariance matrices $\{\tilde{\mathbf{C}}_k\}_{k=1}^K$ of the GMM $f_{\mathbf{y}}^{(K)}$ can be determined via (3.16). After this offline initialization process, all quantities required to compute $\hat{\mathbf{h}}^{(K)}$ are given and $\hat{\mathbf{h}}^{(K)}$ can be used for online channel estimation. The whole process is summarized in Algorithm 2.

Note that since the covariance matrices $\{\tilde{\mathbf{C}}_k\}_{k=1}^K$ of $f_{\mathbf{y}}^{(K)}$ depend on the noise covariance matrix $\boldsymbol{\Sigma}$ (cf. (3.16)), the estimator $\hat{\mathbf{h}}^{(K)}$ also depends on it. In particular,

if the SNR changes, the estimator $\hat{\mathbf{h}}^{(K)}$ needs to be updated. However, this only requires updating $\{\tilde{\mathbf{C}}_k\}_{k=1}^K$ —it is not necessary to fit a new GMM. The fitting process only determines $f_{\mathbf{h}}^{(K)}$, which does not depend on the SNR, and it only needs to be run once for a given $K \in \mathbb{N}$. Generally, it makes sense to precompute the quantities in (3.22) for a number of different SNRs so that computing $\hat{\mathbf{h}}^{(K)}$ online only amounts to matrix-vector multiplications. The complexity of computing the GMM estimator $\hat{\mathbf{h}}^{(K)}$ is discussed next.

3.3.3 Computational Complexity and Number of Parameters

As mentioned in Section 3.3.2, once the GMM fitting process, which yields $\{p(k)\}_{k=1}^K$, $\{\boldsymbol{\mu}_k\}_{k=1}^K$, and $\{\mathbf{C}_k\}_{k=1}^K$, is done, the remaining quantities in (3.22) can be precomputed for different SNRs. This is, in particular, also true for the costly inverses $\{\tilde{\mathbf{C}}_k^{-1}\}_{k=1}^K$ and determinants $\{\det(\tilde{\mathbf{C}}_k)\}_{k=1}^K$. Both of these quantities are required to evaluate the responsibilities in (3.22) because a Gaussian PDF with mean vector $\tilde{\boldsymbol{\mu}} \in \mathbb{C}^m$ and covariance matrix $\tilde{\mathbf{C}} \in \mathbb{C}^{m \times m}$ is given by

$$\mathbf{x} \mapsto \mathcal{N}_{\mathbb{C}}(\mathbf{x}; \tilde{\boldsymbol{\mu}}, \tilde{\mathbf{C}}) = \frac{\exp(-(\mathbf{x} - \tilde{\boldsymbol{\mu}})^{\text{H}} \tilde{\mathbf{C}}^{-1} (\mathbf{x} - \tilde{\boldsymbol{\mu}}))}{\pi^n \det(\tilde{\mathbf{C}})}. \quad (3.23)$$

It can now be seen that for $m \leq N$, an evaluation of $\mathcal{N}_{\mathbb{C}}(\tilde{\boldsymbol{\mu}}, \tilde{\mathbf{C}})$ has a complexity of $\mathcal{O}(N^2)$. The last ingredient to computing $\hat{\mathbf{h}}^{(K)}$ is the LMMSE formula (3.21b) which costs $\mathcal{O}(N^2)$. Since this needs to be done K times (cf. (3.22)), the overall complexity of computing the GMM estimator $\hat{\mathbf{h}}^{(K)}$ is $\mathcal{O}(KN^2)$.

In some cases, we are interested in reducing the number of parameters of the GMM $f_{\mathbf{h}}^{(K)}$. The parameters are the mixing coefficients $\{p(k)\}_{k=1}^K$, the mean vectors $\{\boldsymbol{\mu}_k\}_{k=1}^K$, and the covariance matrices $\{\mathbf{C}_k\}_{k=1}^K$ so that we have $K + KN + K \frac{N(N+1)}{2}$ scalar parameters if we take symmetries in $\{\mathbf{C}_k\}_{k=1}^K$ into account. One reason to reduce the number of parameters is to reduce the computational complexity. As demonstrated in Section 3.3.3.1, cases exist where the complexity of computing $\hat{\mathbf{h}}^{(K)}$ can be reduced to $\mathcal{O}(KN \log(N))$. Another reason to reduce the number of parameters is to improve the GMM fitting process. We can expect to need a larger number T_{tr} of training data $\{\mathbf{h}_t\}_{t=1}^{T_{\text{tr}}}$ in order to successfully fit a GMM with a larger number of parameters. For a given amount of training data, we can either choose a smaller K to control the number of GMM parameters, or, as discussed next, we can impose structure on the covariance matrices. The following two examples are motivated by commonly encountered matrix structures in mobile communications.

3.3.3.1 Circulant covariance matrices

In the first example, we force every covariance matrix of the GMM $f_{\mathbf{h}}^{(K)}$ to have the circulant structure

$$\mathbf{C}_k = \mathbf{F}^H \text{diag}(\mathbf{c}_k) \mathbf{F} \quad (3.24)$$

where $\mathbf{F} \in \mathbb{C}^{N \times N}$ is the discrete Fourier transform (DFT) matrix and $\mathbf{c}_k \in \mathbb{C}^N$. This significantly reduces the number of GMM parameters to $K + KN + KN$. The motivation comes from the fact that channel covariance matrices are Toeplitz structured if the base station employs a uniform linear array (ULA) and that Toeplitz matrices are well approximated by circulant matrices [32]. The charm of circulant matrices is that thanks to fast Fourier transforms, matrix-vector multiplications can be performed in $\mathcal{O}(N \log(N))$ time. This is particularly interesting if we have $\mathbf{A} = \mathbf{I}$ and $\mathbf{\Sigma} = \sigma^2 \mathbf{I} = \sigma^2 \mathbf{F}^H \mathbf{F}$ because the LMMSE formula (3.21b) simplifies to

$$\hat{\mathbf{h}}_{\text{LMMSE},k}(\mathbf{y}) = \mathbf{F}^H \text{diag}(\mathbf{d}_k) \mathbf{F}(\mathbf{y} - \boldsymbol{\mu}_k) + \boldsymbol{\mu}_k \quad (3.25)$$

where the i th entry of the vector \mathbf{d}_k is given by $[\mathbf{d}_k]_i = \frac{[\mathbf{c}_k]_i}{[\mathbf{c}_k]_i + \sigma^2}$, such that (3.25) can be calculated in $\mathcal{O}(N \log(N))$ time. On top of this, with a circulant covariance matrix, the Gaussian density can be written as

$$\mathbf{x} \mapsto \mathcal{N}_{\mathbb{C}}(\mathbf{x}; \boldsymbol{\mu}, \mathbf{F}^H \text{diag}(\mathbf{c}) \mathbf{F}) = \frac{\exp(-(\mathbf{F}(\mathbf{x} - \boldsymbol{\mu}))^H \text{diag}(\mathbf{c})^{-1} \mathbf{F}(\mathbf{x} - \boldsymbol{\mu}))}{\pi^n \prod_{i=1}^N [\mathbf{c}]_i}. \quad (3.26)$$

In contrast to (3.23), this can be computed in $\mathcal{O}(N \log(N))$ time. Consequently, computing $\hat{\mathbf{h}}^{(K)}$ has a complexity of $\mathcal{O}(KN \log(N))$. Even if $\mathbf{A} \neq \mathbf{I}$, such that we may not arrive at a complexity of $\mathcal{O}(KN \log(N))$, circulant covariance matrices can still be interesting to control the number of GMM parameters and thereby also the GMM fitting process. In Section 3.6, we investigate how the amount of training data and the number of GMM parameters influence the performance of $\hat{\mathbf{h}}^{(K)}$.

To implement the circulant structure (3.24), e.g., in Python, we can make use of the fact that diagonal covariance matrices are a commonly used constraint for GMM components so that there exist specialized versions of the expectation-maximization (EM) fitting algorithm for this case. We can fit a GMM with diagonal covariance matrices on Fourier transformed training data $\{\mathbf{F} \mathbf{h}_t\}_{t=1}^{T_{\text{tr}}}$, which yields the mean vectors $\{\mathbf{F} \boldsymbol{\mu}_k\}_{k=1}^K$ and covariance matrices $\{\text{diag}(\mathbf{c}_k)\}_{k=1}^K$, so that we obtain the corresponding quantities for the original data $\{\mathbf{h}_t\}_{t=1}^{T_{\text{tr}}}$ as $\{\boldsymbol{\mu}_k\}_{k=1}^K$ and $\{\mathbf{F}^H \text{diag}(\mathbf{c}_k) \mathbf{F}\}_{k=1}^K$ by making use of the DFT matrix. A numerical evaluation of the circulant GMM estimator is presented in Section 3.6.2.

3.3.3.2 Kronecker covariance matrices

Another common channel covariance matrix structure results from the assumption that the scattering in the vicinity of the transmitter and of the receiver are independent of each other, cf. [33]. In that case, we can force the covariance matrices to have the form

$$\mathbf{C}_k = \mathbf{C}_{\text{tx},k} \otimes \mathbf{C}_{\text{rx},k} \quad (3.27)$$

where $\mathbf{C}_{\text{tx},k}$ and $\mathbf{C}_{\text{rx},k}$ represent a transmit and a receive side spatial covariance matrix, respectively. Using the signal model (2.5), if mobile terminals have N_{tx} antennas and the base station has N_{rx} antennas, the training data $\{\mathbf{H}_t\}_{t=1}^{T_{\text{tr}}}$ consist of channel matrices $\mathbf{H}_t \in \mathbb{C}^{N_{\text{rx}} \times N_{\text{tx}}}$. If no constraints on the covariance matrices of the GMM are to be taken into account, we fit a K -component GMM on the vectorized training data $\{\text{vec}(\mathbf{H}_t)\}_{t=1}^{T_{\text{tr}}}$, see also Section 2.2.1.

In contrast, in order to obtain a GMM with matrices of the form (3.27), we first fit two GMMs: We fit a K_{tx} -component GMM on all rows of $\{\mathbf{H}_t\}_{t=1}^{T_{\text{tr}}}$, and we fit a K_{rx} -component GMM on all columns of $\{\mathbf{H}_t\}_{t=1}^{T_{\text{tr}}}$. This results in two sets of GMM parameters: $\{(\boldsymbol{\mu}_{\text{tx},i}, \mathbf{C}_{\text{tx},i})\}_{i=1}^{K_{\text{tx}}}$ and $\{(\boldsymbol{\mu}_{\text{rx},j}, \mathbf{C}_{\text{rx},j})\}_{j=1}^{K_{\text{rx}}}$. These are now combined to obtain the GMM $f_{\mathbf{h}}^{(K)}$ with $K = K_{\text{tx}}K_{\text{rx}}$ components: $\{(\boldsymbol{\mu}_{\text{tx},i} \otimes \boldsymbol{\mu}_{\text{rx},j}, \mathbf{C}_{\text{tx},i} \otimes \mathbf{C}_{\text{rx},j}) : 1 \leq i \leq K_{\text{tx}}, 1 \leq j \leq K_{\text{rx}}\}$. Lastly, given these K mean vectors and covariance matrices, we compute the mixing coefficients $\{p(k)\}_{k=1}^K$ corresponding to $f_{\mathbf{h}}^{(K)}$ by performing one E-step (cf. [17]) of the EM algorithm. Overall, the GMM obtained as described has a significantly reduced number of parameters: $K + K_{\text{tx}}N_{\text{tx}} + K_{\text{rx}}N_{\text{rx}} + K_{\text{tx}}\frac{N_{\text{tx}}(N_{\text{tx}}+1)}{2} + K_{\text{rx}}\frac{N_{\text{rx}}(N_{\text{rx}}+1)}{2}$ in contrast to $K + KN_{\text{tx}}N_{\text{rx}} + K\frac{N_{\text{tx}}N_{\text{rx}}(N_{\text{tx}}N_{\text{rx}}+1)}{2}$. Plugging the Kronecker product matrices into the LMMSE formula (3.21b) does not necessarily lead to a computational advantage due to the inverse being a matrix which is generally not a Kronecker product of two smaller matrices. This inverse, however, can be approximated by means of a Kronecker product, cf., e.g., [34], if the computational complexity is an issue. A numerical evaluation of the Kronecker GMM estimator is presented in Section 3.6.2.

3.4 Channel Models for Numerical Evaluation

Section 3.6 presents a numerical evaluation of the GMM estimator. We consider the single-input multiple-output (SIMO), the multiple-input multiple-output (MIMO), and the wideband single-input single-output (SISO) signal models described in Section 2.2.1. The channels are generated according to one of the following two models.

3.4.1 3GPP

To generate the SIMO and the MIMO channels, we use a spatial channel model [35,36] where channels are modeled as $(\mathbf{h} \mid \mathbf{d} = \mathbf{d}) \sim \mathcal{N}_{\mathbb{C}}(\mathbf{0}, \mathbf{C}_{\mathbf{d}})$. That is, given a realization \mathbf{d} of a random variable \mathbf{d} which models the angles of arrival/departure and the path gains of the main propagation clusters between a mobile terminal and the base station, the channel is assumed to be Gaussian with a covariance matrix $\mathbf{C}_{\mathbf{d}}$ which depends on the realized angles and gains. To generate T channel samples $\{\mathbf{h}_t\}_{t=1}^T$, we realize T different \mathbf{d}_t and corresponding $\mathbf{C}_{\mathbf{d}_t}$. For every $\mathbf{C}_{\mathbf{d}_t}$, we then realize one Gaussian vector \mathbf{h}_t from $\mathcal{N}_{\mathbb{C}}(\mathbf{0}, \mathbf{C}_{\mathbf{d}_t})$.

If both the mobile terminal and the base station employ ULAs, the transmit and receive side spatial channel covariance matrices are given by

$$\mathbf{C}_{\{\text{tx,rx}\},\mathbf{d}} = \int_{-\pi}^{\pi} g_{\{\text{tx,rx}\}}(\theta; \mathbf{d}) \mathbf{a}_{\{\text{tx,rx}\}}(\theta) \mathbf{a}_{\{\text{tx,rx}\}}(\theta)^H d\theta. \quad (3.28)$$

Here,

$$\mathbf{a}_{\{\text{tx,rx}\}}(\theta) = [1, e^{j\pi \sin(\theta)}, \dots, e^{j\pi(N_{\{\text{tx,rx}\}}-1)\sin(\theta)}]^T \quad (3.29)$$

is the array steering vector for an angle of arrival/departure θ and $g_{\{\text{tx,rx}\}}$ is a power density consisting of a sum of weighted Laplace densities whose standard deviations describe the angle spread of the propagation clusters [35]. Assuming independent scattering in the vicinity of the transmitter and the receiver, cf., e.g., [33], the full covariance matrices are then given as $\mathbf{C}_{\mathbf{d}} = \mathbf{C}_{\text{tx},\mathbf{d}} \otimes \mathbf{C}_{\text{rx},\mathbf{d}}$. In the SIMO case, we have $\mathbf{C}_{\mathbf{d}} = \mathbf{C}_{\text{rx},\mathbf{d}}$.

3.4.2 QuaDRiGa

We use version 2.4 of the QuaDRiGa channel simulator [37, 38]. The base station has a height of 25 meters and covers a 120° sector in an urban macrocell scenario. It is equipped with an N_{rx} -antenna ULA with “3GPP-3D” antennas. Mobile terminals are equipped with an N_{tx} -antenna ULA with “omni-directional” antennas. The mobile terminals have a minimum and a maximum distance to the base station of 35 meters and 500 meters, respectively. In 80% of the cases, the mobile terminals are located indoors at different floor levels. The mobile terminals’ height is 1.5 meters in the case of outdoor locations. For the SIMO and the MIMO simulations, we generate single-carrier channels at a center frequency of 2.53 GHz. For the wideband SISO simulations, we have a typical 5G frame structure (cf. [39]) with $N_c = 24$ carriers over a bandwidth of 360 kHz and with $N_t = 14$ time symbols over slots of 1 ms duration.

The center frequency is again 2.53 GHz. Each mobile terminal moves at a certain velocity v in a random direction.

The QuaDRiGa simulator models a channel corresponding to a carrier c (with frequency f_c) and a time symbol t as $\mathbf{H}_{c,t} = \sum_{l=1}^L \mathbf{G}_l e^{-2\pi j f_c \tau_{l,t}}$, where l is the path number. The l th path delay corresponding to the time symbol t is $\tau_{l,t}$. The number of multi-path components L depends on whether there is line of sight (LOS), non-line of sight (NLOS), or outdoor-to-indoor (O2I) propagation: $L_{\text{LOS}} = 37$, $L_{\text{NLOS}} = 61$, or $L_{\text{O2I}} = 37$, cf. [40]. The attenuation of the l th path, the antenna radiation pattern weighting, and the polarization are collected in the coefficients matrix \mathbf{G}_l , cf. [40]. The generated channels are post-processed to remove the path gain as described in the QuaDRiGa manual [38].

3.5 Other Channel Estimators

We now discuss the algorithms which we compare the GMM estimator $\hat{\mathbf{h}}^{(K)}$ to. The numerical evaluation is presented in Section 3.6.

Using the Moore-Penrose pseudoinverse \mathbf{A}^\dagger , a simple baseline algorithm is given by the least squares (LS) channel estimation:

$$\mathbf{y} \mapsto \hat{\mathbf{h}}_{\text{LS}}(\mathbf{y}) = \mathbf{A}^\dagger \mathbf{y}. \quad (3.30)$$

Further, using a set $\{\mathbf{h}_t\}_{t=1}^{T_{\text{tr}}}$ of training data, we can compute a sample covariance matrix $\hat{\mathbf{C}} = \frac{1}{T_{\text{tr}}} \sum_{t=1}^{T_{\text{tr}}} \mathbf{h}_t \mathbf{h}_t^H$ and employ this matrix together with the LMMSE formula such that we have the following channel estimator:

$$\mathbf{y} \mapsto \hat{\mathbf{h}}_{\text{sample cov.}}(\mathbf{y}) = \hat{\mathbf{C}} \mathbf{A}^H (\mathbf{A} \hat{\mathbf{C}} \mathbf{A}^H + \mathbf{\Sigma})^{-1} \mathbf{y}. \quad (3.31)$$

When using the 3GPP channel model from Section 3.4.1, every channel is generated from a Gaussian PDF $\mathcal{N}_{\mathbb{C}}(\mathbf{0}, \mathbf{C}_d)$ so that we have a channel covariance matrix \mathbf{C}_d for every channel sample. Given this matrix (or, equivalently, given the realization \mathbf{d}), the LMMSE estimator

$$\mathbf{y} \mapsto \hat{\mathbf{h}}_{\text{gen. LMMSE}}(\mathbf{y}) = \mathbb{E}[\mathbf{h} \mid \mathbf{y} = \mathbf{y}, \mathbf{d} = \mathbf{d}] = \mathbf{C}_d \mathbf{A}^H (\mathbf{A} \mathbf{C}_d \mathbf{A}^H + \mathbf{\Sigma})^{-1} \mathbf{y} \quad (3.32)$$

is optimal. Note, however, that this estimator is not equal to the CME $\hat{\mathbf{h}}$ from (3.5) because of the additional knowledge of the realization $\mathbf{d} = \mathbf{d}$. Since, in practice, the channel covariance matrix \mathbf{C}_d is not given, we speak of a *genie-aided LMMSE estimator* (gen. LMMSE). For this reason, $\hat{\mathbf{h}}_{\text{gen. LMMSE}}$ yields a performance bound for all estimators in this chapter and we can generally not expect to reach this bound with any of the other algorithms.

Another channel estimator is based on a convolutional neural network (CNN). The estimator is introduced in [36] for the SIMO case (cf. (2.4)) and [41] presents our generalization of it to the MIMO case (cf. (2.5)). For example, assuming the noise covariance matrix is $\Sigma = \sigma^2 \mathbf{I} \in \mathbb{C}^{N \times N}$, the SIMO CNN estimator is given by [36]

$$\mathbf{y} \mapsto \hat{\mathbf{h}}_{\text{CNN}}(\mathbf{y}) = \mathbf{F}_2^{\text{H}} \text{diag} \left(\mathbf{a}^{(2)} \star \text{ReLU} \left(\mathbf{a}^{(1)} \star \frac{1}{\sigma^2} |\mathbf{F}_2 \mathbf{y}|^2 + \mathbf{b}^{(1)} \right) + \mathbf{b}^{(2)} \right) \mathbf{F}_2 \mathbf{y} \quad (3.33)$$

where the columns of the matrix $\mathbf{F}_2 \in \mathbb{C}^{2N \times N}$ consist of the first N columns of the $2N \times 2N$ DFT matrix, ReLU denotes the rectified linear unit, the absolute value in $|\mathbf{F}_2 \mathbf{y}|^2$ is applied elementwise, \star denotes circular convolution, and $\mathbf{a}^{(1)}, \mathbf{a}^{(2)}, \mathbf{b}^{(1)}, \mathbf{b}^{(2)} \in \mathbb{R}^{2N}$ are learnable parameters which are optimized based on an MSE loss function. We use this estimator for SIMO channel estimation. For the MIMO case, we use the estimator as described in [41] where the activation function is also the rectified linear unit. In both cases, the computational complexity is $\mathcal{O}(N \log(N))$. During the training phase, samples corresponding to the channel model on which the estimator is later tested are used. Note that the CNN estimator depends on σ^2 and thereby on the SNR such that a new estimator needs to be trained for every SNR value.

Many modern channel estimation algorithms focus on compressive sensing approaches, cf., e.g., the surveys [42,43]. For channel estimation by means of compressive sensing algorithms, every channel realization is modeled as (approximately) sparse: $\mathbf{h} \approx \Psi \mathbf{s}$. The *dictionary* $\Psi \in \mathbb{C}^{N \times L}$ is typically an oversampled DFT matrix (e.g., [44]), and $\mathbf{s} \in \mathbb{C}^L$ is a sparse vector. Given an observation \mathbf{y} , compressive sensing algorithms then introduce the auxiliary model $\mathbf{y} = \mathbf{A} \Psi \mathbf{s} + \mathbf{n}$ and recover an estimate $\hat{\mathbf{s}}$ of \mathbf{s} so that the channel can be estimated as $\hat{\mathbf{h}}_{\text{CS}} = \Psi \hat{\mathbf{s}}$. One example algorithm is orthogonal matching pursuit (OMP) [14–16] which is detailed in Algorithm 1 (in Section 2.3) and would be used with $\mathbf{C} = \mathbf{A} \Psi$. Since OMP requires the sparsity order (number of nonzero elements in \mathbf{s}) and since this is not known for the two channel models described in Section 3.4, we provide the true channel realization to the OMP algorithm so that it can find the optimal sparsity order, see Algorithm 3. This genie-aided approach yields a performance bound for OMP. Another compressive sensing algorithm we compare our results with is approximate message passing (AMP) [45,46], which does not need to know the sparsity order.

For wideband channel estimation (cf. Section 2.2.1), we compare our results with the concrete autoencoder (CAE) algorithm from [47]. The vectorized wideband channels have a dimension of $N_c N_t$ and the pilot selection matrix $\mathbf{A} \in \{0, 1\}^{N_p \times N_c N_t}$ allows to observe $N_p < N_c N_t$ of the channel elements. The idea of the CAE is to replace the encoder of an autoencoder by a *concrete selector layer* which extracts

Algorithm 3 Genie-Aided Orthogonal Matching Pursuit (gen. OMP)

Require: observation matrix \mathbf{A} , dictionary $\Psi \in \mathbb{C}^{N \times L}$

Require: observation \mathbf{y} and corresponding true channel \mathbf{h}

```

1:  $S^{(0)} \leftarrow \emptyset, \mathbf{s}^{(0)} \leftarrow \mathbf{0}, e \leftarrow \infty$ 
2: for  $i = 1$  to  $L$  do
3:    $j^* \leftarrow \arg \max_{j \in \{1, \dots, L\}} \{[(\mathbf{A}\Psi)^H(\mathbf{y} - \mathbf{A}\Psi\mathbf{s}^{(i-1)})]_j\}$ 
4:    $S^{(i)} \leftarrow S^{(i-1)} \cup \{j^*\}$ 
5:    $\mathbf{s}^{(i)} \leftarrow \arg \min_{\tilde{\mathbf{s}}} \{\|\mathbf{y} - \mathbf{A}\Psi\tilde{\mathbf{s}}\|, \text{support}(\tilde{\mathbf{s}}) \subset S^{(i)}\}$ 
6:   if  $\|\mathbf{h} - \Psi\mathbf{s}^{(i)}\| < e$  then
7:      $e \leftarrow \|\mathbf{h} - \Psi\mathbf{s}^{(i)}\|$  //  $\Psi\mathbf{s}^{(i)}$  is closer to the true  $\mathbf{h}$  than  $\Psi\mathbf{s}^{(i-1)}$ 
8:   else
9:     return  $\Psi\mathbf{s}^{(i-1)}$  // the new estimate  $\Psi\mathbf{s}^{(i)}$  is worse than the previous one
10:  end if
11: end for
12: return  $\Psi\mathbf{s}^{(L)}$ 

```

the N_p most important features of the $N_c N_t$ -dimensional input. This corresponds to designing the matrix \mathbf{A} . The decoder of the CAE then serves as the channel estimator. The training phase uses noisy channel data so that the CAE performs both a denoising of the input as well as a reconstruction of the $N_c N_t$ -dimensional channels. Since the training data is noisy, a new CAE needs to be trained for every SNR. We do not employ further denoising networks after the decoder of the CAE.

Another wideband channel estimator, called *ChannelNet*, can be found in [48]. It consists of 2D CNNs and combines both an image super-resolution network and an image restoration network. The super-resolution network has three layers, the image restoration network has twenty layers. In this context, the receive signal is interpreted as a low-resolution image and the channel is the image that needs to be restored. Again, a new ChannelNet needs to be trained for every SNR.

3.6 Numerical Evaluation

We look at SIMO, MIMO, and wideband channel estimation scenarios (cf. Section 2.2.1) involving the GMM estimator $\hat{\mathbf{h}}^{(K)}$ from (3.22). To obtain the GMM estimator, we fit one GMM on the given training data via an EM algorithm, see Section 3.3.2. Thereafter, the inverses and the determinants in (3.22) are precomputed for every SNR. In comparison, the neural network-based estimators from Section 3.5 all require new estimators to be trained for every SNR. In particular, for every SNR, a

hyperparameter search (cf. Section 2.7) is necessary.

Channel estimation performance is measured in terms of the normalized MSE (nMSE). To this end, $T_{\text{tst}} = 10^4$ N -dimensional test channel samples $\{\mathbf{h}_t\}_{t=1}^{T_{\text{tst}}}$ are generated and used to get a corresponding number of noisy observations of the form (2.3) from which channel estimates $\{\hat{\mathbf{h}}_t\}_{t=1}^{T_{\text{tst}}}$ are computed with every algorithm of interest. Then, we calculate $\text{nMSE} = \frac{1}{NT_{\text{tst}}} \sum_{t=1}^{T_{\text{tst}}} \|\mathbf{h}_t - \hat{\mathbf{h}}_t\|^2$. The channels are normalized such that $\mathbb{E}[\|\mathbf{h}\|^2] = N$ holds. The noise covariance matrix is $\mathbf{\Sigma} = \sigma^2 \mathbf{I} \in \mathbb{C}^{N \times N}$ which leads us to define the SNR as $\frac{1}{\sigma^2}$. For algorithms which require training data (including the sample covariance estimator (3.31)), we generate $T_{\text{tr}} = 10^5$ samples unless another T_{tr} is stated. Generally, we choose T_{tr} large enough such that increasing it does not lead to better results in the testing phase.

3.6.1 SIMO Channel Estimation

In the SIMO case (Figures 3.2 to 3.6), the observation matrix is the identity matrix $\mathbf{A} = \mathbf{I} \in \mathbb{C}^{N \times N}$ and the dictionaries of OMP and AMP are $L = 4N$ and $L = 2N$ times oversampled DFT dictionaries, respectively. These numbers lead to the best performance in our simulations. We use a training data size of $T_{\text{tr}} = 19 \cdot 10^4$ to fit the GMM unless another T_{tr} is stated.

The first experiment, displayed in Figure 3.2, uses the 3GPP channel model from Section 3.4.1 with one propagation cluster and the number N of antennas at the base station is varied. The genie-aided LMMSE estimator from (3.32) provides a utopian performance bound for all estimators. The GMM estimator with $K = 128$ components comes close to this performance bound. At around $N = 96$ antennas, the CNN estimator (cf. (3.33)) starts to become better than the genie-aided OMP estimator. The reason is that the assumptions which are made in [36] to derive the CNN estimator are better fulfilled for larger N . Generally, the qualitative behavior of all estimators does not change much as N varies. For this reason, the number of antennas is fixed at $N = 128$ in the other experiments.

In Figure 3.3, we again use the 3GPP channel model from Section 3.4.1 with one propagation cluster. From the mid- to high-SNR range, the two compressive sensing algorithms show a similar performance and they are among the better estimators. For most SNRs, the GMM estimator with $K = 128$ components comes close to the genie-aided LMMSE performance bound. In Figure 3.4, where we have three propagation clusters, the CNN estimator's strong performance in the mid-SNR range is prominent. In the high-SNR range, all estimators tend to show a similar performance to the LS estimator. The only exception is the GMM estimator which still improves

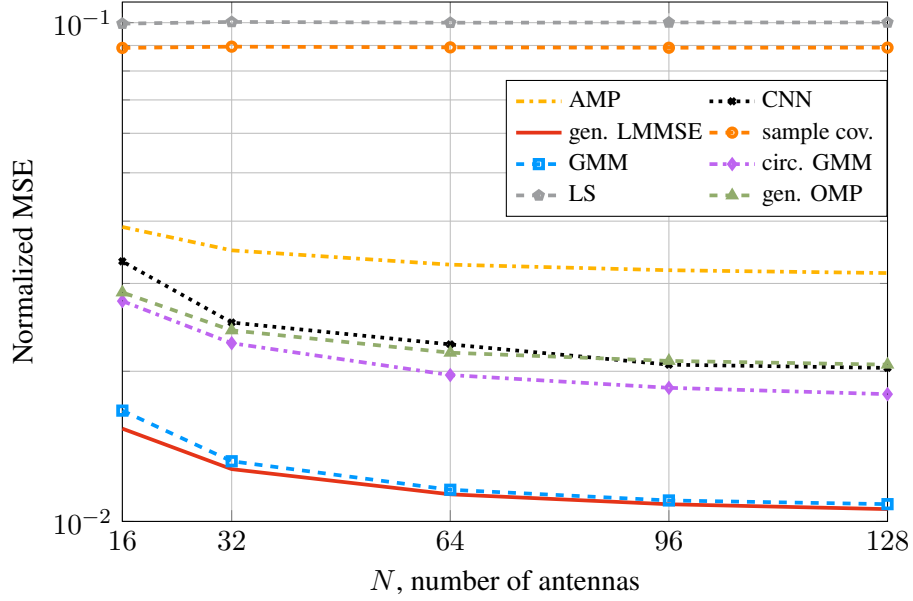


Figure 3.2: SIMO signal model and 3GPP channel model (Section 3.4.1) with **one** propagation cluster at 10 dB SNR. The performance of the circulant GMM estimator (“circ. GMM”, $-\diamond-$, Section 3.3.3.1) is shown, too. In both cases, $K = 128$ components are used.

upon the LS estimator. Recall that we generally cannot expect any estimator to reach the genie LMMSE curve because it has more knowledge (the true covariance matrix for every sample, cf. (3.32)) than is usually available.

In Figure 3.5, we consider the QuaDRiGa channel model (cf. Section 3.4.2) where the channel covariance matrices and therefore the genie LMMSE curve are no longer available. The two compressive sensing algorithms’ performances now differ considerably which may be attributed to a lack of sparsity of the channels. Both the CNN and the GMM estimators generally show a strong performance mainly differing in the low-SNR range where the GMM estimator is better.

In Figure 3.6, we investigate the interplay between the number K of GMM components and the number T_{tr} of training data used to fit the K -component GMM. We concentrate on the 3GPP channel model with three propagation clusters at an SNR of 10 dB (see also Figure 3.4). As expected, since increasing K leads to a larger number of GMM parameters, more training data is necessary for a good fit. As long as T_{tr} is large enough, the behavior indicated by Theorem 2 can be observed: increasing K improves the estimator’s performance. Theorem 2 can be applied because we have

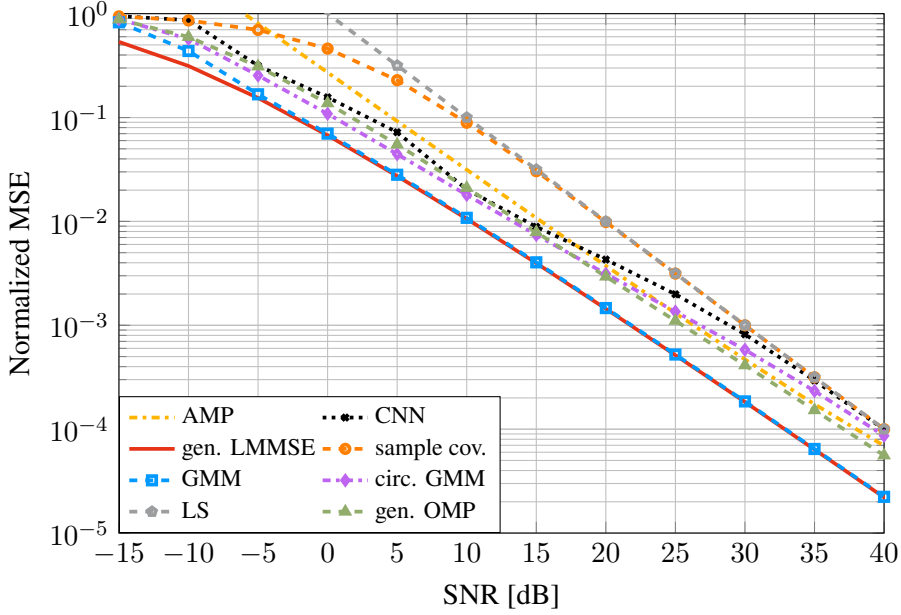


Figure 3.3: SIMO signal model and 3GPP channel model (Section 3.4.1) with **one** propagation cluster and $N = 128$ antennas. The performance of the circulant GMM estimator (“circ. GMM”, $- \blacklozenge -$, Section 3.3.3.1) is shown, too. In both cases, $K = 128$ components are used.

$\mathbf{A} = \mathbf{I}$ which is invertible. At this point, note that we generally cannot expect the GMM estimator to converge to the genie LMMSE estimator (3.32) (which is displayed in Figure 3.4). The genie LMMSE estimator has more knowledge (namely the true channel covariance matrix \mathbf{C}_d for every channel realization) and is therefore not the CME $\mathbf{y} \mapsto \hat{\mathbf{h}}(\mathbf{y}) = \mathbb{E}[\mathbf{h} \mid \mathbf{y} = \mathbf{y}]$ which we want to approximate in Theorem 2. The CME $\hat{\mathbf{h}}$ cannot be computed (because $f_{\mathbf{h}}$ is not known) which is the motivation to study the GMM estimator in the first place.

3.6.2 GMM Estimator with Structured Covariance Matrices

We briefly study the circulant and Kronecker GMM estimators from Section 3.3.3.1 and Section 3.3.3.2, respectively. The circulant GMM estimator uses circulant covariance matrices in the GMM. Figures 3.3 to 3.5 contain the corresponding curves. The circulant GMM estimator has a significantly reduced number of parameters and a computational complexity of $\mathcal{O}(N \log(N))$. Its performance is not as good as that of

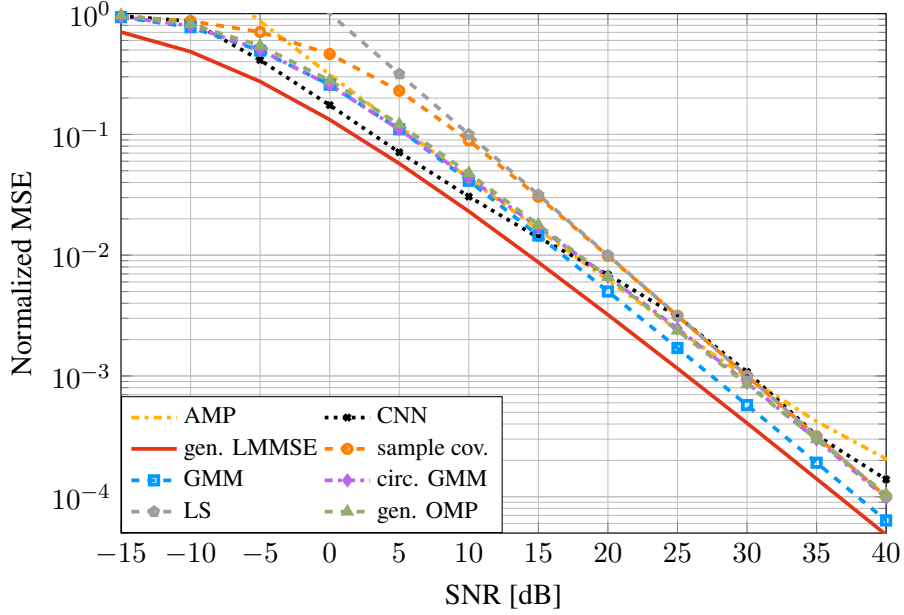


Figure 3.4: SIMO signal model and 3GPP channel model (Section 3.4.1) with **three** propagation clusters and $N = 128$ antennas. The performance of the circulant GMM estimator (“circ. GMM”, $- \cdot \diamond \cdot -$, Section 3.3.3.1) is shown, too. In both cases, $K = 128$ components are used.

the GMM estimator with unconstrained covariance matrices but it can compete with the other depicted estimators. In particular, in Figure 3.5 the circulant GMM estimator performs well and is even comparable to the unconstrained GMM estimator.

For a brief discussion of the Kronecker GMM estimator, we consider a MIMO simulation with a scaled DFT matrix as pilot matrix \mathbf{P} , cf. (2.5). Our main interest is to compare the Kronecker GMM estimator and the GMM estimator with unconstrained covariance matrices. First, we generate $T_{\text{tr}} = 10^5$ training data $\{\mathbf{H}_t\}_{t=1}^{T_{\text{tr}}}$. Then, in order to obtain the GMM estimator with unconstrained covariance matrices, we fit a GMM with $K = 32$ components on the vectorized data $\{\text{vec}(\mathbf{H}_t)\}_{t=1}^{T_{\text{tr}}}$. In contrast, in order to get the Kronecker GMM estimator, we fit two GMMs: a GMM with $K_{\text{tx}} = 4$ components using all of the rows of $\{\mathbf{H}_t\}_{t=1}^{T_{\text{tr}}}$ and a GMM with $K_{\text{rx}} = 8$ components using all of the columns of $\{\mathbf{H}_t\}_{t=1}^{T_{\text{tr}}}$. These two GMMs are then combined to obtain the Kronecker GMM estimator with $K_{\text{tx}}K_{\text{rx}} = 32$ components. The details are explained in Section 3.3.3.2.

A MIMO simulation with $(N_{\text{rx}}, N_{\text{tx}}) = (32, 4)$ can be found in Figure 3.7. It can

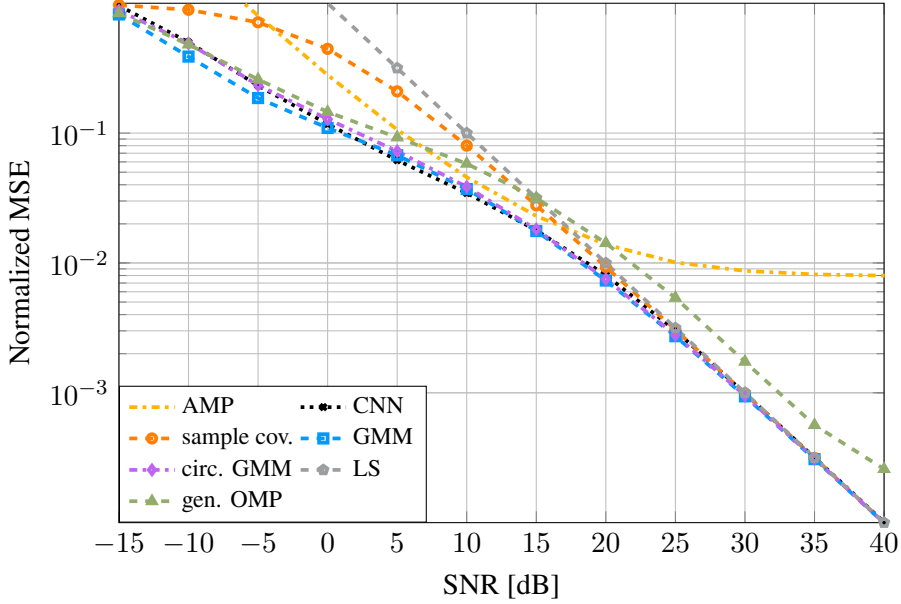


Figure 3.5: SIMO signal model and QuaDRiGa channel model (Section 3.4.2) with $N = 128$ antennas. The performance of the circulant GMM estimator (“circ. GMM”, $-\diamond-$, Section 3.3.3.1) is shown, too. In both cases, $K = 128$ components are used.

be seen that the two GMM estimators show a similar performance. This is interesting because these two estimators have a different number of parameters. In particular, the GMM estimator with unconstrained covariance matrices has $K \frac{N(N+1)}{2} = 264,192$ covariance parameters with $N = N_{\text{rx}}N_{\text{tx}} = 128$. By comparison, the Kronecker GMM estimator has only $K_{\text{rx}} \frac{N_{\text{rx}}(N_{\text{rx}}+1)}{2} + K_{\text{tx}} \frac{N_{\text{tx}}(N_{\text{tx}}+1)}{2} = 4,224 + 40 = 4,264$ covariance parameters. As the number of parameters differs, we can expect the Kronecker GMM estimator to show a better performance if the amount of training data is relatively small. This behavior can be seen in Figure 3.8: The performance of the two different GMM estimators is very different for a small amount of training data and the unconstrained GMM estimator starts to outperform the Kronecker GMM estimator when enough training data are available. Another observation in Figure 3.7 is that also the CNN estimator performs similarly to the GMM estimators with a minor performance loss in the high-SNR regime. All three mentioned estimators outperform the channel covariance matrix-based estimator as well as the genie-aided OMP algorithm. The OMP algorithm uses a Kronecker product of two two-times oversampled DFT matrices as a dictionary in this simulation.

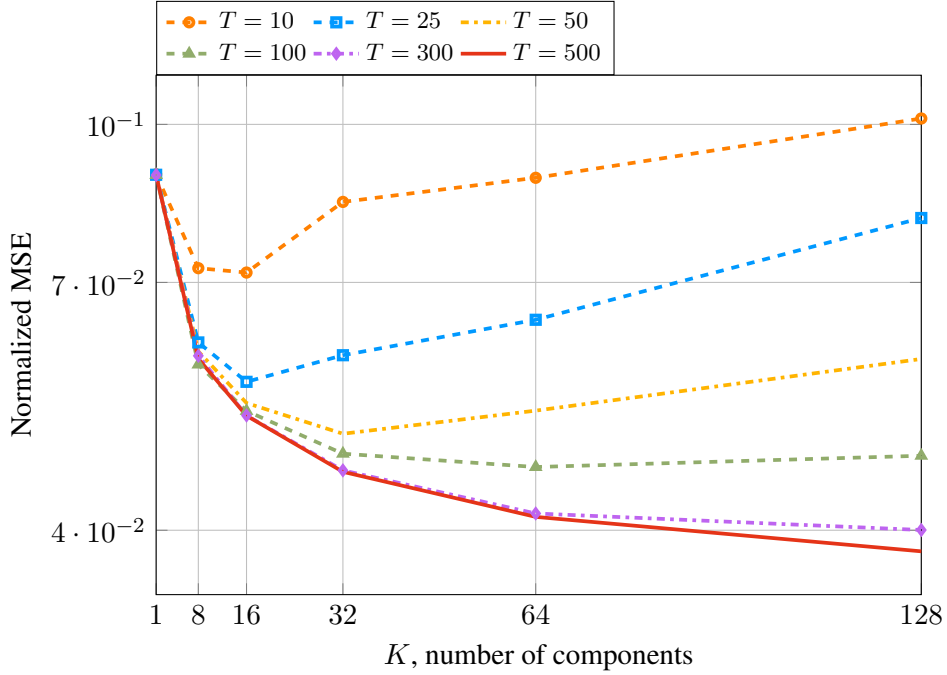


Figure 3.6: SIMO signal model and 3GPP channel model (Section 3.4.1) with **three** propagation clusters and $N = 128$ antennas. The SNR is 10 dB. The GMM estimator is fit using $T_{\text{tr}} = T \cdot 10^3$ training data.

3.6.3 Wideband Channel Estimation

For wideband simulations, we focus on the 5G frame structure described in Section 3.4.2. The number of pilots is $N_p = 50$ so that 50 out of the $24 \cdot 14 = 336$ resource elements are occupied with pilots. Thus, here, the observation matrix $\mathbf{A} \in \{0, 1\}^{50 \times 336}$ is not invertible so that Theorem 2 does not apply. Nonetheless, the following Figures 3.9 to 3.11 show the desired behavior of an improving GMM estimator if the number K of components is increased. In this subsection, the considered algorithms are the GMM estimator, the sample covariance matrix-based estimator in (3.31), ChannelNet, and the CAE estimator (see Section 3.5). All estimators are trained on $T_{\text{tr}} = 10^5$ channel samples. There are three commonly used types of pilot arrangements: block, comb, lattice (see, e.g., [49, 50]). We indicate in the figures which of these is used.

Figure 3.9 shows results of a simulation where the mobile terminals move at a constant velocity of $v = 3$ km/h. The block-type pilot arrangement is used with the CAE estimator as the only exception because it uses its own optimized pilot

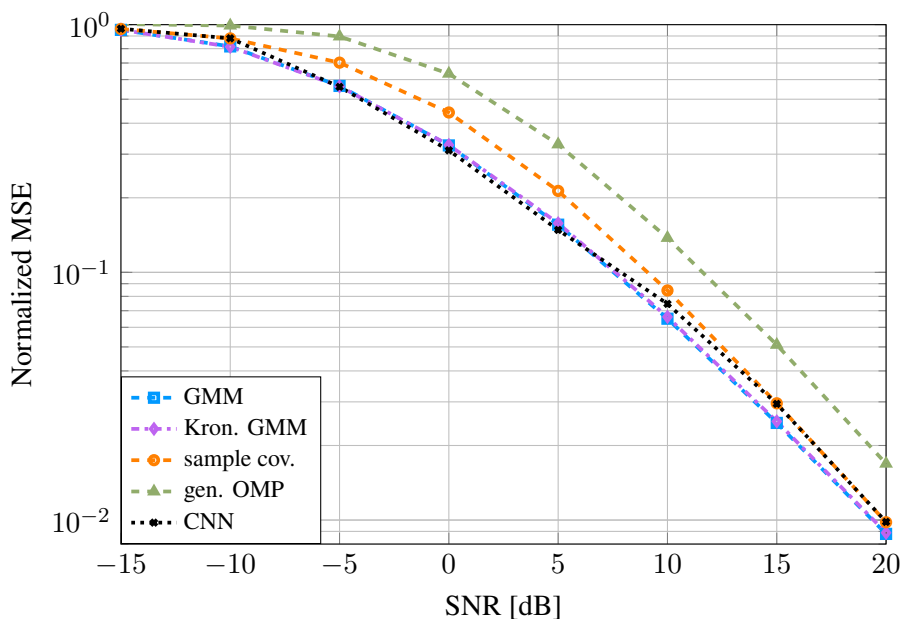


Figure 3.7: MIMO signal model and QuaDRiGa channel model (Section 3.4.2) with $(N_{\text{rx}}, N_{\text{tx}}) = (32, 4)$. The GMM estimator (“GMM”, -■-) uses $K = 32$ components and the Kronecker GMM estimator (“Kron. GMM”, -◆·-, Section 3.3.3.2) uses $K = 4 \cdot 8$ components.

pattern. Both ChannelNet and the CAE estimator behave similarly and outperform the sample covariance matrix-based estimator. The GMM estimator can compete with or outperforms both of those estimators. Noticably, an increase of the number of components from $K = 8$ to $K = 128$ leads to a performance improvement even though the observation matrix \mathbf{A} is not invertible.

In Figure 3.10, the mobile terminals’ velocities are chosen uniformly at random from the interval $[0, 300]$ km/h. In this setting, the lattice-type pilot arrangement is used. Overall, channel estimation is more difficult now and the gap between the sample covariance matrix-based estimator and the other estimators increases. Further, choosing a larger number K of components for the GMM estimator is necessary for a competitive performance. Again, the GMM estimator shows the desired behavior of an improved performance for increased K . The last figure, Figure 3.11, investigates the GMM estimator’s performance for different numbers of components and highlights the importance of a large enough amount of training data. Overall, the considerations in this subsection provide some numerical evidence for the GMM estimator’s convergence

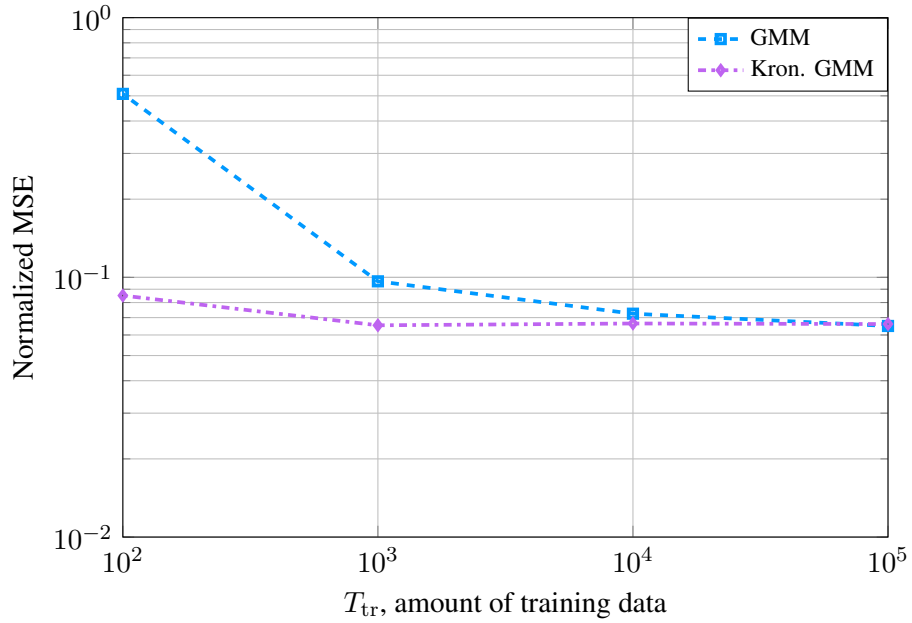


Figure 3.8: MIMO signal model and QuaDRiGa channel model (Section 3.4.2) with $(N_{rx}, N_{tx}) = (32, 4)$. The SNR is 10 dB. The Kronecker product GMM estimator (“Kron. GMM”, $- \blacklozenge \cdot -$, Section 3.3.3.2) has fewer parameters than the GMM estimator with unconstrained covariance matrices. Both estimators have $K = 32$ components.

in case of noninvertible observation matrices.

3.6.4 Conclusion

The main contribution of this chapter is a strong motivation to employ the GMM estimator even if the channel PDF $f_{\mathbf{h}}$ is not a GMM. On the one hand, we have a theoretical motivation in terms of Theorem 2. On the other hand, we have a motivation in terms of the numerical simulations. Overall, the numerical simulations seem to show the performance promised by Theorem 2. In the SIMO and in the MIMO settings, Theorem 2 can be applied directly and the desired behavior is evident. In the wideband setting, the observation matrix is not invertible so that we do not know whether Theorem 2 holds. Nonetheless, the GMM estimator’s performance improves when the number of components is increased.

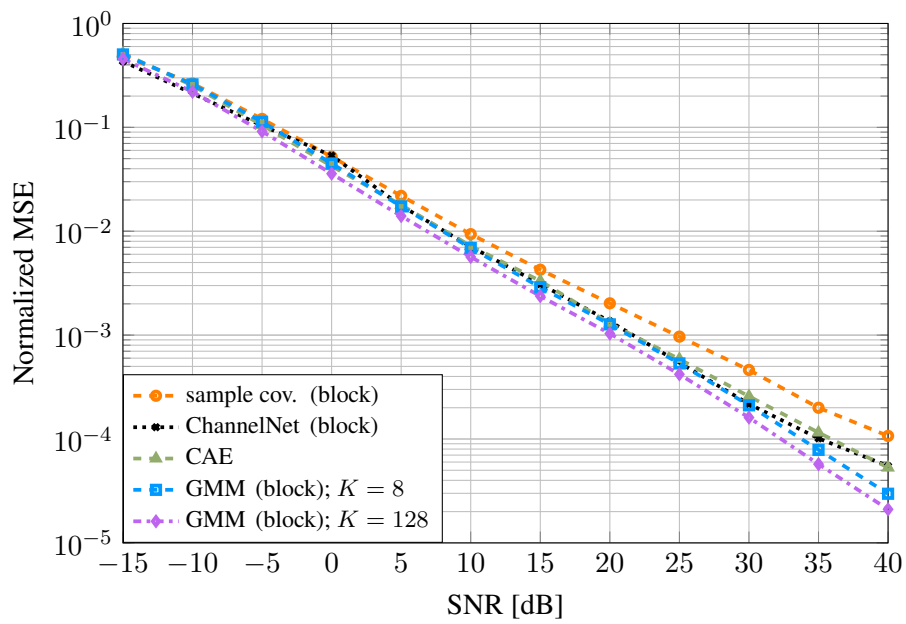


Figure 3.9: Wideband signal model and QuaDRiGa channel model (Section 3.4.2) with $N_p = 50$ pilots for $N_c = 24$ carriers and $N_t = 14$ time symbols for $v = 3$ km/h.

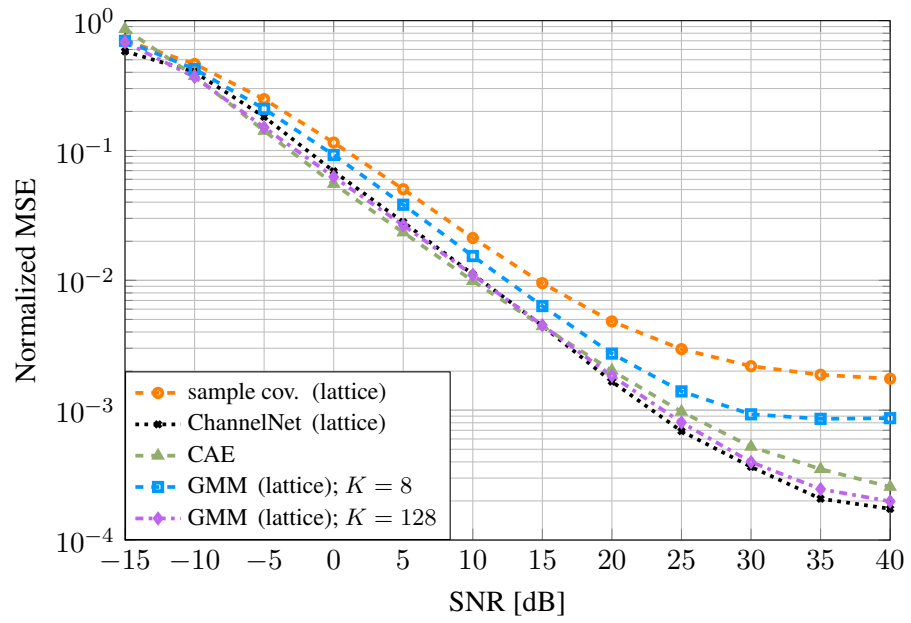


Figure 3.10: Wideband signal model and QuaDRiGa channel model (Section 3.4.2) with $N_p = 50$ pilots for $N_c = 24$ carriers and $N_t = 14$ time symbols for $v \in [0, 300]$ km/h.

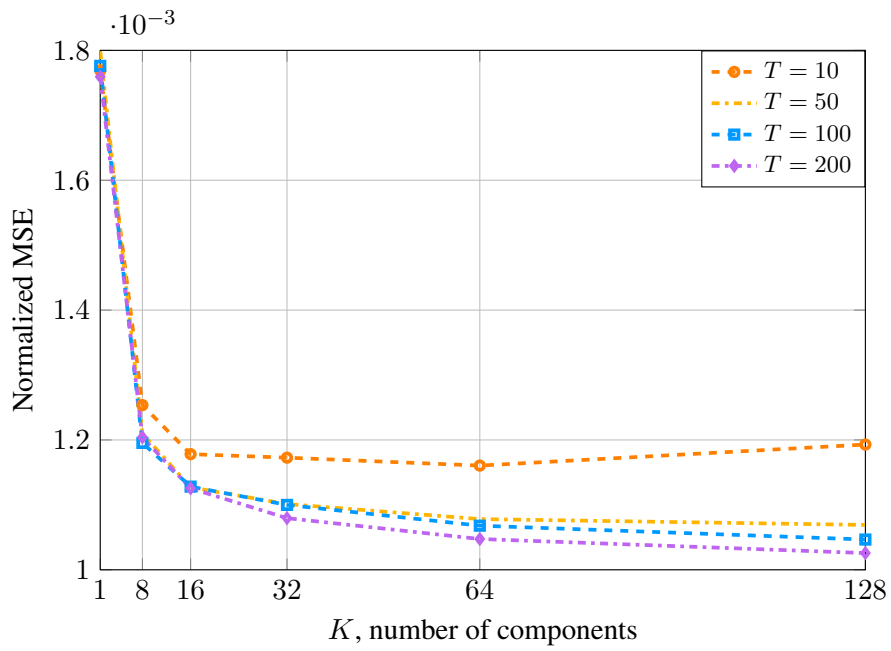


Figure 3.11: Wideband signal model and QuaDRiGa channel model (Section 3.4.2) with $N_p = 50$ pilots for $N_c = 24$ carriers and $N_t = 14$ time symbols for $v = 3$ km/h. The block-type pilot arrangement is used. The SNR is 20 dB. The GMM estimator is fit using $T_{\text{tr}} = T \cdot 10^3$ training data.

Learning an Observation Matrix

This chapter is mainly based on our work in [51]. In Chapter 3, we study a channel estimator for the signal model

$$\mathbf{y} = \mathbf{A}\mathbf{h} + \mathbf{n}. \quad (4.1)$$

In this chapter, we propose a method to find a suitable observation matrix $\mathbf{A} \in \mathbb{C}^{m \times N}$ to recover \mathbf{h} given this model. The focus lies on constant modulus matrices, see Section 2.2.1. For ease of exposition, we define the set of constant modulus matrices in $\mathbb{C}^{m \times N}$ as $\mathcal{M}_{\text{const}}^{m \times N} = \{\mathbf{A} \in \mathbb{C}^{m \times N} : |[\mathbf{A}]_{k,l}| = \frac{1}{\sqrt{m}}\}$. We assume that all channel realizations are an element of a set $\mathcal{H} \subset \mathbb{C}^N$. For example, \mathcal{H} could be the set of all p -sparse vectors for some $p < N$; or, \mathcal{H} could, e.g., contain vectors which have a sparse representation in some unknown basis, see also Section 2.3. The goal in this chapter is to obtain a fixed constant modulus observation matrix suitable for recovering vectors in \mathcal{H} from observations of the form (4.1).

4.1 Proposed Method

Compressive sensing theory provides answers to the question what properties of an observation matrix $\mathbf{A} \in \mathbb{C}^{m \times N}$ with $m < N$ can facilitate the recovery of \mathbf{h} from (4.1). One widely used property is the restricted isometry property (RIP) in (2.7). Since $\mathbf{0} \in \mathbb{C}^N$ fulfills (2.7) for any $\delta \geq 0$ and any $\mathbf{A} \in \mathbb{C}^{m \times N}$, we exclude it from the set \mathcal{H} and then write the RIP condition as follows:

$$(1 - \delta) \leq \left\| \mathbf{A} \frac{\mathbf{h}}{\|\mathbf{h}\|} \right\|^2 \leq (1 + \delta) \quad \text{for all } \mathbf{h} \in \mathcal{H}. \quad (4.2)$$

Matrices which fulfill (4.2) with δ as small as possible are desired.

The vector $\bar{\mathbf{h}} = \frac{\mathbf{h}}{\|\mathbf{h}\|} \in \mathbb{C}^N$ in (4.2) lies on the (unit) hypersphere in \mathbb{C}^N and according to (4.2), \mathbf{A} maps $\bar{\mathbf{h}}$ to a vector $\mathbf{A}\bar{\mathbf{h}} \in \mathbb{C}^m$ which lies (approximately) on the (unit) hypersphere in \mathbb{C}^m —the Euclidean distance between $\mathbf{A}\bar{\mathbf{h}}$ and the closest point to it on the hypersphere in \mathbb{C}^m is at most δ . In this analogy, finding an \mathbf{A} with a small δ corresponds to determining a matrix which brings all $\mathbf{A}\bar{\mathbf{h}}$ close to the hypersphere in \mathbb{C}^m . In this chapter, we aim to find a fixed matrix which has this

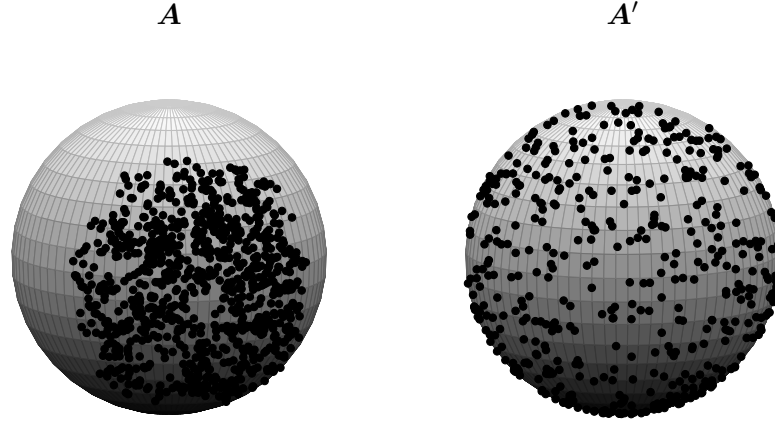


Figure 4.1: The image of 500 points $\mathbf{h}_i/\|\mathbf{h}_i\|$ with $\mathbf{h}_i \in \mathcal{H}$ under two different matrices \mathbf{A} (left) and \mathbf{A}' (right).

property. Construction methods to obtain a matrix with the RIP deterministically are generally unknown and one typically resorts to drawing random matrices which have the RIP with high probability, cf. Section 2.3. In view of this, it seems advisable to design a probabilistic method to achieve our goal.

To motivate the design of our proposed method, consider the sketch in Figure 4.1. With the hyperspheres interpretation of (4.2) in mind, imagine two matrices are given, \mathbf{A} and \mathbf{A}' with the RIP, and assume $\delta = 0$ holds for both matrices. That is, both \mathbf{A} and \mathbf{A}' map all $\bar{\mathbf{h}}$ perfectly onto the hypersphere in \mathbb{C}^m . Further, imagine \mathbf{A} maps all $\bar{\mathbf{h}}$ in such a way onto the hypersphere in \mathbb{C}^m that the points $\mathbf{A}\bar{\mathbf{h}}$ are all clustered around the point $\mathbf{e}_1 = (1, 0, 0, \dots)^\top \in \mathbb{C}^m$ (i.e., all $\mathbf{A}\bar{\mathbf{h}}$ are clustered around a “pole”). In contrast, assume \mathbf{A}' maps all $\bar{\mathbf{h}}$ in such a way onto the hypersphere in \mathbb{C}^m that the points $\mathbf{A}'\bar{\mathbf{h}}$ are isotropically distributed around the origin.

We would now like to argue that the matrix \mathbf{A}' offers more robustness against measurement noise than the matrix \mathbf{A} . To this end, let $\mathbf{h}_1, \mathbf{h}_2 \in \mathcal{H}$ be two vectors with $\|\mathbf{h}_1\| \neq \|\mathbf{h}_2\|$ so that we also have $\|\mathbf{A}\mathbf{h}_1\| \neq \|\mathbf{A}\mathbf{h}_2\|$ as well as $\|\mathbf{A}'\mathbf{h}_1\| \neq \|\mathbf{A}'\mathbf{h}_2\|$ (because $\delta = 0$ is assumed). This means that both \mathbf{A} and \mathbf{A}' map \mathbf{h}_1 and \mathbf{h}_2 onto two different noiseless observations — $\mathbf{A}\mathbf{h}_1 \neq \mathbf{A}\mathbf{h}_2$ and $\mathbf{A}'\mathbf{h}_1 \neq \mathbf{A}'\mathbf{h}_2$ — so that there is a chance of uniquely recovering \mathbf{h}_1 and \mathbf{h}_2 in both cases. At the same time, note that we have $\|\mathbf{A}\mathbf{h}_1\| = \|\mathbf{A}'\mathbf{h}_1\|$ and $\|\mathbf{A}\mathbf{h}_2\| = \|\mathbf{A}'\mathbf{h}_2\|$ (because $\delta = 0$ is assumed). Therefore, the main difference between the pairs of noiseless observations $(\mathbf{A}\mathbf{h}_1, \mathbf{A}\mathbf{h}_2)$ and $(\mathbf{A}'\mathbf{h}_1, \mathbf{A}'\mathbf{h}_2)$ are their relative directions. That is, given the “pole clustering property” of \mathbf{A} , the vectors $\mathbf{A}\mathbf{h}_1$ and $\mathbf{A}\mathbf{h}_2$ tend to point in the same direction (in the

direction of e_1) whereas the vectors $A'h_1$ and $A'h_2$ can point in more dissimilar directions and are thus farther apart from one another. In the presence of noise, i.e., if we have observations $\mathbf{y}_i = \mathbf{A}\mathbf{h}_i + \mathbf{n}_i$ and $\mathbf{y}'_i = \mathbf{A}'\mathbf{h}_i + \mathbf{n}_i$, it can occur that the noise realizations map \mathbf{y}_1 and \mathbf{y}_2 close to each other, making them numerically indistinguishable, which, in turn, makes recovering distinct estimates of \mathbf{h}_1 and \mathbf{h}_2 difficult. At the same time, \mathbf{y}'_1 and \mathbf{y}'_2 might still be distinguishable because $A'h_1$ and $A'h_2$ are likely to be farther apart from one another. Generally, we believe that observations under A' are more robust with respect to noise than under A , making A' preferable.

In view of this, we aim to find a matrix A which maps all $\bar{\mathbf{h}}$ to the hypersphere in \mathbb{C}^m and which at the same time distributes the points isotropically around the origin. The first requirement corresponds to satisfying the RIP condition with a small δ and the latter requirement aims to combat measurement noise. Therefore, in summary, we want all $A\bar{\mathbf{h}}$ to be uniformly distributed on the hypersphere in \mathbb{C}^m to achieve both requirements. To express this goal in terms of an optimization problem, we let all vectors in \mathcal{H} be realizations of a random variable \mathbf{h} so that also $\bar{\mathbf{h}}$ and $A\bar{\mathbf{h}}$ are random variables. Note that $A\bar{\mathbf{h}}$ is a random variable because $\bar{\mathbf{h}}$ is a random variable whereas A is not random. Further, let \mathbf{u} be a random variable with a uniform distribution on the (unit) hypersphere in \mathbb{C}^m . That is, we have $\|\mathbf{u}\| = 1$ and \mathbf{u} is isotropically distributed around the origin. If we find a matrix which achieves the stated goal, then $A\bar{\mathbf{h}}$ has the same distribution as \mathbf{u} . Generally, we aim to find an A which brings the distribution of $A\bar{\mathbf{h}}$ close to the distribution of \mathbf{u} . Using the maximum mean discrepancy (MMD) as a metric between probability distributions (cf. Section 2.5), we can encapsulate our goal in the following distribution matching (cf. Section 2.6) optimization problem:

$$\min_{A \in \mathbb{C}^{m \times N}} \text{MMD}_k^2(p_{\mathbf{u}}, q_{A\bar{\mathbf{h}}}) \quad (4.3)$$

where we write $p_{\mathbf{u}}$ and $q_{A\bar{\mathbf{h}}}$ for the distributions of \mathbf{u} and $A\bar{\mathbf{h}}$, respectively. Further, k in (4.3) denotes a kernel (cf. Section 2.5).

A main challenge with (4.3) is that in order to evaluate the objective function, the distribution $q_{A\bar{\mathbf{h}}}$ needs to be available in closed form. Even if that is the case, the evaluation can be too difficult. This is a challenge which MMD shares with alternative distances like, e.g., the Kullback-Leibler divergence. An advantage of MMD however is the possibility to estimate $\text{MMD}_k^2(p_{\mathbf{u}}, q_{A\bar{\mathbf{h}}})$ based on samples of \mathbf{u} and $\bar{\mathbf{h}}$ (and thereby $A\bar{\mathbf{h}}$), which can easily be generated (in case of \mathbf{u}) or are assumed to be available (in case of $\bar{\mathbf{h}}$). This, together with differentiability, makes MMD charming and enables us to employ machine learning techniques to solve the optimization problem (4.3).

4.1.1 Learning a Constant Modulus Matrix

Coming back to the goal of finding a suitable constant modulus matrix, we want to solve

$$\min_{\substack{\mathbf{A} \in \mathbb{C}^{m \times N} \\ \|\mathbf{A}\|_{k,l} = \frac{1}{\sqrt{m}}}} \text{MMD}_k^2(p_{\mathbf{u}}, q_{\mathbf{A}\bar{\mathbf{h}}}). \quad (4.4)$$

This optimization problem aims to find a constant modulus matrix $\mathbf{A} \in \mathcal{M}_{\text{const}}^{m \times N}$ that brings $q_{\mathbf{A}\bar{\mathbf{h}}}$ as close to a uniform distribution on the (unit) hypersphere in \mathbb{C}^m as possible. Obtaining an analytic solution seems difficult in general, for example, because an analytic expression of $q_{\mathbf{A}\bar{\mathbf{h}}}$ is not available.

However, as detailed in Section 2.6 where we go from problem (2.17) to problem (2.18), we can use samples $\{\bar{\mathbf{h}}_t\}_{t=1}^{T_{\text{tr}}}$ and $\{\mathbf{u}_t\}_{t=1}^{T_{\text{tr}}}$ of both $\bar{\mathbf{h}}$ and \mathbf{u} to solve an approximation of (4.4). A random variable \mathbf{u} with uniform distribution on the (unit) hypersphere in \mathbb{C}^m is given by normalization of an isotropic Gaussian random variable, e.g., $\mathbf{u} = \frac{\mathbf{v}}{\|\mathbf{v}\|}$ with $\mathbf{v} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{0}, \mathbf{I})$, which means that sampling is simple. If

$$\text{stk} : \mathbb{C}^m \rightarrow \mathbb{R}^{2m}, \mathbf{z} \mapsto \text{stk}(\mathbf{z}) = [\Re(\mathbf{z})^T, \Im(\mathbf{z})^T]^T \quad (4.5)$$

denotes stacking real and imaginary parts of a complex vector into a real vector, we can formulate an approximation of (4.4):

$$\min_{\substack{\mathbf{A} \in \mathbb{C}^{m \times N} \\ \|\mathbf{A}\|_{k,l} = \frac{1}{\sqrt{m}}}} \text{MMD}_k^2 \left(\{\text{stk}(\mathbf{u}_t)\}_{t=1}^{T_{\text{tr}}}, \{\text{stk}(\mathbf{A}\bar{\mathbf{h}}_t)\}_{t=1}^{T_{\text{tr}}} \right). \quad (4.6)$$

We would further like to express (4.6) in terms of real-valued optimization variables. Every entry of $\mathbf{A} \in \mathcal{M}_{\text{const}}^{m \times N}$ has the form $[\mathbf{A}]_{k,l} = \frac{1}{\sqrt{m}} \exp(j \phi_{k,l})$ with a phase $\phi_{k,l} \in \mathbb{R}$ so that $\mathbf{A} \in \mathcal{M}_{\text{const}}^{m \times N} \subset \mathbb{C}^{m \times N}$ is parameterized by mN real-valued parameters which we collect in the matrix $\Phi \in \mathbb{R}^{m \times N}$. Applying sin and cos elementwise, we can write

$$\mathbf{A} = \Re(\mathbf{A}) + j \Im(\mathbf{A}) = \frac{1}{\sqrt{m}} \cos(\Phi) + j \frac{1}{\sqrt{m}} \sin(\Phi) \quad (4.7)$$

and $\mathbf{A}\bar{\mathbf{h}} = (\Re(\mathbf{A}) + j \Im(\mathbf{A}))(\Re(\bar{\mathbf{h}}) + j \Im(\bar{\mathbf{h}}))$ yields

$$\text{stk}(\mathbf{A}\bar{\mathbf{h}}) = \begin{bmatrix} \Re(\mathbf{A}\bar{\mathbf{h}}) \\ \Im(\mathbf{A}\bar{\mathbf{h}}) \end{bmatrix} = \underbrace{\begin{bmatrix} \frac{\cos(\Phi)}{\sqrt{m}} & -\frac{\sin(\Phi)}{\sqrt{m}} \\ \frac{\sin(\Phi)}{\sqrt{m}} & \frac{\cos(\Phi)}{\sqrt{m}} \end{bmatrix}}_{=\mathbf{A}(\Phi)} \begin{bmatrix} \Re(\bar{\mathbf{h}}) \\ \Im(\bar{\mathbf{h}}) \end{bmatrix} = \mathbf{A}(\Phi) \text{stk}(\bar{\mathbf{h}}) \quad (4.8)$$

Algorithm 4 Learning $\mathbf{A} \in \mathcal{M}_{\text{const}}^{m \times N}$

Require: training data $\{(\bar{\mathbf{h}}_t, \mathbf{u}_t)\}_{t=1}^{T_{\text{tr}}}$

- 1: randomly initialize $\Phi \in \mathbb{R}^{m \times N}$
- 2: **while** termination criterion not met **do**
- 3: draw a batch of T samples uniformly from the training data: $\{(\bar{\mathbf{h}}_t, \mathbf{u}_t)\}_{t=1}^T \subset \{(\bar{\mathbf{h}}_t, \mathbf{u}_t)\}_{t=1}^{T_{\text{tr}}}$
- 4: compute the stochastic gradient

$$\mathbf{g}_{\Phi} = \frac{\partial}{\partial \Phi} \text{MMD}_k^2 \left(\{\text{stk}(\mathbf{u}_t)\}_{t=1}^T, \{\mathbf{A}(\Phi) \text{stk}(\bar{\mathbf{h}}_t)\}_{t=1}^T \right)$$

- 5: update Φ with a gradient algorithm using \mathbf{g}_{Φ}
 - 6: **end while**
-

where the notation $\mathbf{A}(\Phi) \in \mathbb{R}^{2m \times 2N}$ emphasizes that \mathbf{A} only depends on the mN real parameters in Φ . The optimization problem with real-valued optimization parameters is thus

$$\min_{\Phi \in \mathbb{R}^{m \times N}} \text{MMD}_k^2 \left(\{\text{stk}(\mathbf{u}_t)\}_{t=1}^{T_{\text{tr}}}, \{\mathbf{A}(\Phi) \text{stk}(\bar{\mathbf{h}}_t)\}_{t=1}^{T_{\text{tr}}} \right). \quad (4.9)$$

As explained in Section 2.6, such a problem is typically solved via stochastic gradient descent, see Algorithm 4. An example of how the kernel k can be chosen and of how good hyperparameters of the gradient algorithm in Algorithm 4 can be found is presented in Section 4.6.3.

Calculating the gradient in Algorithm 4 has the same order of computational complexity as evaluating MMD_k^2 in the forward pass if backpropagation is used [52]. One evaluation of MMD_k^2 for a batch of T samples requires $\mathcal{O}(T^2)$ evaluations of the kernel, see (2.16), where the vectors have dimension $2m$ due to stacked real and imaginary parts. Evaluating the kernel can be done in $\mathcal{O}(m)$, see (2.15). Computing $\mathbf{A}(\Phi) \text{stk}(\bar{\mathbf{h}})$ can be done in $\mathcal{O}(mN)$. Therefore, the computational complexity of evaluating MMD_k^2 is $\mathcal{O}(mNT + mT^2)$. To save computation time, it can be interesting to use a linear (in T) version of MMD which can be found in [2, Lemma 14]. This would result in a complexity of $\mathcal{O}(mNT)$.

4.2 Related Literature

The following reviews different approaches to design an observation matrix. It is worth pointing out that the approaches focus on real-valued observation matrices without

constraints, which differs from the approach presented in Section 4.1 where we are looking for a complex-valued constant modulus matrix in $\mathcal{M}_{\text{const}}^{m \times N}$. The only exception is the very recently introduced algorithm in [53] which also studies constant modulus matrices, see Section 4.4. The comparison of Algorithm 4 to the method in [53] is new in this dissertation and not present in [51]. Further, we modify an algorithm from [54] such that it also yields constant modulus matrices, see Section 4.3.1.

The proposed Algorithm 4 and all reviewed approaches have a data dependency in common. This is the case because either data samples are used explicitly or because the obtained matrix is based on a dictionary which, in turn, is typically suitable for a particular data set. In contrast, random matrices, which are the default choice in compressive sensing applications (cf. Section 2.3), enjoy a universality property which makes them applicable with almost any dictionary and thus data [10]. The data dependency, therefore, introduces a tradeoff where universality is traded for a potentially better performance on a given data set.

4.2.1 Literature Review

The real-valued signal model $\mathbf{y} = \mathbf{A}\mathbf{h} + \mathbf{n}$ with $\mathbf{A} \in \mathbb{R}^{m \times N}$ is studied in [55]. Given a set $\{\mathbf{h}_t\}_{t=1}^{T_{\text{tr}}}$ of samples, the proposed *Uncertain Component Analysis* determines

$$\mathbf{A} = \arg \max_{\tilde{\mathbf{A}} \in \mathbb{R}^{m \times N}, \tilde{\mathbf{A}}\tilde{\mathbf{A}}^T = \mathbf{I}} \prod_{t=1}^{T_{\text{tr}}} \Pr(\mathbf{h}_t | \mathbf{y}_t; \tilde{\mathbf{A}}) \quad (4.10)$$

with $\mathbf{y}_t = \tilde{\mathbf{A}}\mathbf{h}_t$ and where $\Pr(\mathbf{h}_t | \mathbf{y}_t; \tilde{\mathbf{A}})$ is the posterior probability of the data. Two fixed-point equations are used to solve the optimization problem algorithmically.

The author of [56] works with the signal model $\mathbf{y} = \mathbf{A}\Psi\mathbf{s}$ with $\mathbf{A} \in \mathbb{R}^{m \times N}$ and where $\Psi \in \mathbb{R}^{N \times L}$ ($N < L$) is a dictionary such that \mathbf{s} is sparse. The author proposes minimizing the t -averaged mutual coherence $\mu_t(\mathbf{A}\Psi)$ with respect to \mathbf{A} . To this end, the columns of $\mathbf{A}\Psi$ are normalized and then the Gram matrix $\mathbf{G} = (\mathbf{A}\Psi)^T(\mathbf{A}\Psi)$ is computed. The t -averaged mutual coherence is now given by the average of all entries with $|\mathbf{G}_{k,l}| \geq t$. An iterative algorithm which alternates between two stages is developed to achieve the goal. The algorithm chooses a *down-scaling factor* $\gamma \in (0, 1)$ and starts with an initial $\mathbf{A}^{(1)}$. In iteration i , the entries of $\mathbf{G}^{(i)} = (\mathbf{A}^{(i)}\Psi)^T(\mathbf{A}^{(i)}\Psi)$, which have an absolute value larger than t , are multiplied by γ to *shrink* them. This yields a matrix $\hat{\mathbf{G}}^{(i)}$ which is decomposed in the algorithm's second stage. The decomposition aims to find $\mathbf{A}^{(i+1)}$ such that $(\mathbf{A}^{(i+1)}\Psi)^T(\mathbf{A}^{(i+1)}\Psi)$ is close to $\hat{\mathbf{G}}^{(i)}$. The two stages are alternated for a predetermined number of iterations.

In the same setting as [56], the Gram matrix $\mathbf{G} = (\mathbf{A}\Psi)^T(\mathbf{A}\Psi)$ is utilized to find $\mathbf{A} \in \mathbb{R}^{m \times N}$ in [57]. There, the goal is to bring \mathbf{G} close to the identity matrix. After a random initialization, a KSVD-like algorithm determines \mathbf{A} . Further, [57] proposes a modification of the approach where also the dictionary Ψ can be optimized.

Also the authors of [58] work with a Gram matrix. Their goal is to determine $\mathbf{A} \in \mathbb{R}^{m \times N}$ such that a bi-Lipschitz condition is fulfilled. Specifically, given a set $\{\mathbf{h}_t\}_{t=1}^{T_{\text{tr}}}$ of samples, the goal is to establish the RIP for all differences $\mathbf{h}_t - \mathbf{h}_s$. To this end, the secant set

$$\mathcal{S} = \left\{ \frac{\mathbf{h}_t - \mathbf{h}_s}{\|\mathbf{h}_t - \mathbf{h}_s\|} : 1 \leq t < s \leq T_{\text{tr}} \right\} \quad (4.11)$$

is defined and the bi-Lipschitz criterion is expressed as

$$\left| \|\mathbf{A}\mathbf{v}\|^2 - 1 \right| = |\mathbf{v}^T \mathbf{A}^T \mathbf{A} \mathbf{v} - 1| = |\mathbf{v}^T \mathbf{G} \mathbf{v} - 1| \leq \delta \quad (4.12)$$

for all $\mathbf{v} \in \mathcal{S}$ and for $\mathbf{G} = \mathbf{A}^T \mathbf{A}$. The motivating optimization problem then is

$$\min_{\mathbf{G} \in \mathbb{R}^{N \times N}} \sup_{\mathbf{v} \in \mathcal{S}} |\mathbf{v}^T \mathbf{G} \mathbf{v} - 1| \quad \text{s.t.} \quad \mathbf{G} \succeq 0, \quad \text{rank}(\mathbf{G}) = r, \quad \text{trace}(\mathbf{G}) = b \quad (4.13)$$

and a relaxation thereof is solved. The related problem

$$\min_{\mathbf{G} \in \mathbb{R}^{N \times N}, \mathbf{G} = \mathbf{G}^T, \mathbf{G} \succeq 0} \text{trace}(\mathbf{G}) \quad \text{s.t.} \quad \sup_{\mathbf{v} \in \mathcal{S}} |\mathbf{v}^T \mathbf{G} \mathbf{v} - 1| \leq \delta \quad (4.14)$$

is studied in [59]. In both cases, an eigendecomposition of the optimal \mathbf{G} yields the desired \mathbf{A} .

Neural networks are used in [60] in the context of the signal model $\mathbf{y} = \mathbf{A}\mathbf{s}$ with nonnegative sparse vectors $\mathbf{s} \in \mathbb{R}_+^N$ and $\mathbf{A} \in \mathbb{R}^{m \times N}$. The matrix \mathbf{A} is learned by *unrolling* the update rule of a projected subgradient method. The authors of [61] then employ the same unrolling method for channel state information feedback. There, the matrix $\mathbf{A} \in \mathbb{R}^{m \times 2N}$ is used to compress a channel $\mathbf{h} \in \mathbb{C}^N$ as $\mathbf{y} = \mathbf{A}[\Re(\mathbf{h})^T, \Im(\mathbf{h})^T]^T \in \mathbb{R}^m$, allowing the base station to recover \mathbf{h} from the fed back signal \mathbf{y} using compressive sensing methods. The approach is modified in [62] for channel estimation. To this end, real and imaginary parts of channels $\mathbf{h} \in \mathbb{C}^N$ are treated as independent data points such that a real-valued matrix $\mathbf{A} \in \mathbb{R}^{m \times N}$ can be learned.

Algorithm 5 Learning-Based Compressive Subsampling (LBCS) [54]

Require: training data $\{\mathbf{h}_t\}_{t=1}^{T_{\text{tr}}}$

Require: matrix $\mathbf{V} \in \mathbb{C}^{N \times N}$, number m of rows to select

1: normalize the training data: $\{\bar{\mathbf{h}}_t\}_{t=1}^{T_{\text{tr}}} \leftarrow \left\{ \frac{\mathbf{h}_t}{\|\mathbf{h}_t\|} \right\}_{t=1}^{T_{\text{tr}}}$

2: **for** $r = 1$ to N **do**

3: $\alpha_r \leftarrow \sum_{t=1}^{T_{\text{tr}}} |\mathbf{v}_r^* \bar{\mathbf{h}}_t|^2$ // \mathbf{v}_r is row r of \mathbf{V}

4: **end for**

5: $\Omega \leftarrow$ indices r which correspond to m largest α_r

6: **return** $\mathbf{P}_\Omega \mathbf{V}$ // learned observation matrix

4.3 Learning-Based Compressive Subsampling

The methods summarized in Section 4.2 lead to an unconstrained real-valued matrix $\mathbf{A} \in \mathbb{R}^{m \times N}$ whereas we are looking for a complex-valued constant modulus matrix $\mathbf{A} \in \mathcal{M}_{\text{const}}^{m \times N}$. The methods which are based on a decomposition $\mathbf{G} = \mathbf{A}^T \mathbf{A}$ of a Gram matrix seem ill-suited for our task because such a decomposition will generally not exist if the entries are constrained to have a constant modulus. However, it seems simple to adapt the algorithm in [54] to the constant modulus constraint. After introducing the algorithm, we present a modified version in Section 4.3.1.

The signal model in [54] is

$$\mathbf{y} = \mathbf{P}_\Omega \mathbf{V} \mathbf{h} \quad (4.15)$$

for a unitary matrix $\mathbf{V} \in \mathbb{C}^{N \times N}$. Here, $\Omega \subset \{1, 2, \dots, N\}$ is an index set and $\mathbf{P}_\Omega \in \mathbb{C}^{m \times N}$ selects the rows of \mathbf{V} corresponding to the indices in Ω . The observation matrix is then given by $\mathbf{A} = \mathbf{P}_\Omega \mathbf{V}$. Selecting the indices randomly is a known strategy to obtain a suitable \mathbf{A} . However, the authors of [54] propose to select the row indices deterministically based on a set $\{\bar{\mathbf{h}}_t\}_{t=1}^{T_{\text{tr}}}$ of normalized samples since this strategy can improve upon random indices. The authors of [54] select the m rows $\mathbf{v}_r \in \mathbb{C}^{1 \times N}$ of \mathbf{V} which maximize the average captured energy

$$\frac{1}{T_{\text{tr}}} \sum_{t=1}^{T_{\text{tr}}} \sum_{r \in \Omega} |\mathbf{v}_r^* \bar{\mathbf{h}}_t|^2 \quad (4.16)$$

of the training set. We call the algorithm ‘‘learning-based compressive subsampling (LBCS)’’ after the paper’s title. The algorithm is summarized in Algorithm 5.

Algorithm 6 Monte Carlo LBCS

Require: training data $\{\mathbf{h}_t\}_{t=1}^{T_{\text{tr}}}$, validation data $\{\mathbf{h}_t\}_{t=1}^{T_{\text{val}}}$ **Require:** evaluation parameters par , number m of rows to select, iterations I

- 1: $\mathbf{A}_{\text{MC}} \leftarrow \mathbf{0}$, $\text{error}_{\text{MC}} \leftarrow \infty$
 - 2: **for** $i = 1$ to I **do**
 - 3: draw a random unitary matrix $\mathbf{V} \in \mathbb{C}^{N \times N}$
 - 4: divide all entries of \mathbf{V} by their absolute values
 - 5: normalize the rows of \mathbf{V}
 - 6: $\mathbf{A} \leftarrow \text{LBCS}(\{\mathbf{h}_t\}_{t=1}^{T_{\text{tr}}}, \mathbf{V}, m)$
 - 7: $\text{error} \leftarrow \text{Evaluation}(\mathbf{A}, \{\mathbf{h}_t\}_{t=1}^{T_{\text{val}}}, \text{par})$
 - 8: **if** $\text{error} < \text{error}_{\text{MC}}$ **then**
 - 9: $\text{error}_{\text{MC}} \leftarrow \text{error}$, $\mathbf{A}_{\text{MC}} \leftarrow \mathbf{A}$
 - 10: **end if**
 - 11: **end for**
 - 12: **return** \mathbf{A}_{MC}
-

4.3.1 Modified LBCS Algorithm

Since Algorithm 5 selects rows of a given matrix $\mathbf{V} \in \mathbb{C}^{N \times N}$, we get a constant modulus matrix whenever \mathbf{V} is a constant modulus matrix. Thus, the first modification is to initialize Algorithm 5 with a constant modulus matrix. The second modification is due to the fact that the algorithm was derived for noiseless observations whereas we work with the signal model (4.1). The main idea is to run Algorithm 5 multiple times (I times in Algorithm 6) for multiple random initializations, to then evaluate all the obtained matrices on a held out set $\{\mathbf{h}_t\}_{t=1}^{T_{\text{val}}}$ of samples, and to return the best performing matrix. The modified algorithm is outlined in Algorithm 6. An algorithm for generating random unitary matrices (cf. Line 3) can be found, e.g., in [63]. The evaluation (cf. Line 7) depends on the desired goal. Since we are ultimately interested in channel estimation, in Section 4.6, we evaluate the current observation matrix by running a channel estimation algorithm on the data set and by returning the corresponding mean square error (MSE). In this case, we have $\text{error} = \text{MSE}$. Generally, the evaluation needs to return an error and the matrix which achieves the smallest error is returned by Algorithm 6. An example of this part of the algorithm is explained in Section 4.6.

Algorithm 7 PGDM [53]

Require: initial $\mathbf{A}^{(0)} \in \mathbb{C}^{m \times N}$ with $|\mathbf{A}^{(0)}_{k,l}| = 1$, dictionary $\Psi \in \mathbb{C}^{N \times L}$

Require: $\alpha \geq 1$, step size ζ , iterations I

- 1: $\beta \leftarrow \sqrt{\frac{L-m}{m(L-1)}}$
- 2: **for** $i = 1$ to I **do**
- 3: $\mathbf{Q} \leftarrow \mathbf{A}^{(i-1)} \Psi$
- 4: $\mathbf{D} \leftarrow \text{diag}\left(\frac{1}{\|\mathbf{q}_1\|}, \dots, \frac{1}{\|\mathbf{q}_L\|}\right)$ // \mathbf{q}_j is column j of \mathbf{Q}
- 5: $\mathbf{E} \leftarrow \mathbf{D} \mathbf{Q}^H \mathbf{Q} \mathbf{D} - \mathbf{I}$
- 6: compute $\tilde{\mathbf{E}}$ with elements

$$[\tilde{\mathbf{E}}]_{k,l} \leftarrow \begin{cases} 0, & |[\mathbf{E}]_{k,l}| < \alpha\beta \\ \frac{[\mathbf{E}]_{k,l}}{||[\mathbf{E}]_{k,l}|} (|[\mathbf{E}]_{k,l}| - \alpha\beta), & \text{otherwise} \end{cases}$$

- 7: $\mathbf{C} \leftarrow 2\tilde{\mathbf{E}} \mathbf{D} \mathbf{Q}^H \mathbf{Q} \mathbf{D}^3$
- 8: compute the gradient

$$\mathbf{g} \leftarrow 4\mathbf{Q} \mathbf{D} \tilde{\mathbf{E}} \mathbf{D} \Psi^H - 2\mathbf{A}^{(i-1)} \Psi \text{diag}(\mathbf{C}) \Psi^H$$

- 9: $\tilde{\mathbf{A}} \leftarrow \mathbf{A}^{(i-1)} - \zeta \mathbf{g}$ // gradient descent
 - 10: compute $\mathbf{A}^{(i)}$ with elements $[\mathbf{A}^{(i)}]_{k,l} \leftarrow \frac{[\tilde{\mathbf{A}}]_{k,l}}{||[\tilde{\mathbf{A}}]_{k,l}|}$ // projection
 - 11: **end for**
 - 12: **return** $\frac{1}{\sqrt{m}} \mathbf{A}^{(I)}$
-

4.4 Projected Gradient Descent Method

In [53], the setting of this chapter is studied and the goal is to obtain a constant modulus matrix. A dictionary $\Psi \in \mathbb{C}^{N \times L}$ for the data of interest is assumed to be given. Using the abbreviation $\Phi = \mathbf{A} \Psi$, the authors want to minimize the mutual coherence

$$\mu(\Phi) = \max_{k \neq l} \frac{|\phi_k^H \phi_l|}{\|\phi_k\| \|\phi_l\|} \quad (4.17)$$

where ϕ_k is column k of Φ . Since the corresponding optimization problem is difficult, the authors instead solve

$$\min_{\mathbf{A} \in \mathbb{C}^{m \times N}} \|\Phi^H \Phi - \mathbf{I}\|_F^2 \quad \text{s.t.} \quad \|\phi_k\| = 1 \quad \text{and} \quad |\mathbf{A}_{k,l}| = 1 \quad (4.18)$$

Algorithm	Goal	Method	Data Dependency
Alg. 4 (MMD)	spherical uniform distribution	distribution matching with maximum mean discrepancy	via training data
Alg. 6 (LBCS)	maximize average captured energy	selecting rows	via training data
Alg. 7 (PGDM)	minimize coherence	projected gradient descent on surrogate objective	via dictionary

Table 4.1: Three algorithms to obtain a constant modulus matrix. Algorithm 4 is proposed in this chapter, Algorithm 6 is our modification of [54], and Algorithm 7 is from [53].

via a projected gradient descent method (PGDM), where $\|\cdot\|_F$ is the Frobenius norm.

The details of PGDM are summarized in Algorithm 7. The algorithm’s hyperparameters α, ζ, I need to be determined via a grid or a random search (cf. Section 2.7). The hyperparameters which lead to the smallest coherence are considered to be the best. The initialization of $\mathbf{A}^{(0)}$ is another hyperparameter. In [53], the m leading eigenvalues and corresponding eigenvectors of the eigenvalue decomposition $\mathbf{U}\mathbf{\Lambda}\mathbf{U}^H = \mathbf{\Psi}\mathbf{\Psi}^H$ are used to compute $\tilde{\mathbf{A}} = \mathbf{\Lambda}_m^{-1/2}\mathbf{U}_m^H$ which is then projected as in Line 10 of Algorithm 7. However, since $\tilde{\mathbf{A}}$ can contain zeros, the projection and therefore this initialization is not always possible. In that case, we initialize $\mathbf{A}^{(0)}$ with independent random elements of the form $e^{j\phi}$ with ϕ uniformly drawn in the interval $[0, 2\pi]$ in the numerical experiments.

4.5 Comparison

With Algorithm 4, Algorithm 6, and Algorithm 7, we have three different methods to obtain a constant modulus matrix. Table 4.1 briefly summarizes the different approaches.

All algorithms share a grid search or a random search aspect: Algorithms 4 and 7 depend on a good choice of their hyperparameters (e.g., learning rate) which are found by running the algorithms multiple times with different hyperparameter values, and Algorithm 6 searches over multiple random matrix initializations. Another similarity of the algorithms is data dependency: Algorithms 4 and 6 require training data and Algorithm 7 requires a dictionary which is typically tailored for a particular data set.

The necessity of a dictionary can be viewed as an advantage and as a disadvantage. On the one hand, if the obtained matrix is used in conjunction with a recovery algorithm which also requires a dictionary, then the matrix and the recovery algorithm might harmonize well. On the other hand, if the recovery algorithm does not need a dictionary or if no dictionary for a given data set is known, it might not be beneficial or it might be difficult to employ Algorithm 7.

A difference between Algorithms 4 and 7 and Algorithm 6 is that the latter requires an evaluation function to decide for the best matrix (Line 7 in Algorithm 6). In the following numerical simulations, we evaluate the current matrix via a channel estimation algorithm. As such, the obtained matrix depends on this algorithm (and on the signal-to-noise ratio (SNR)) which might not be desirable. We discuss the dependency of Algorithm 6 on the evaluation function in more detail in Chapter 5 where we combine the matrix design algorithms with the channel estimation algorithm from Chapter 3. In the remainder of this chapter, we compare the three algorithms in Table 4.1 in compressive sensing recovery tasks.

4.6 Numerical Evaluation

We investigate three data models which have an (approximate) sparsity in common and are thus typical for compressive sensing applications. We continue to use the nomenclature established in the dissertation so far. Only the third data model can be considered a channel model but we speak of channel models and channel estimation nonetheless. We are interested in assessing how well-suited the matrices determined via Algorithm 4, Algorithm 6, or Algorithm 7 are for the problem of channel estimation from observations of the form (4.1). To this end, we employ the well-known compressive sensing algorithm orthogonal matching pursuit (OMP) (see Algorithm 1 in Section 2.3) to perform channel estimation using different matrices and evaluate the corresponding MSE as the figure of merit. Ultimately, the choice to use OMP as a channel estimation algorithm is arbitrary. It is, however, a widely used and fast algorithm. Our main interest lies in observing the effect of different matrices on the channel estimation performance. In Chapter 5, we discuss the performance of the Gaussian mixture model (GMM) estimator from Chapter 3 when it is used with matrices learned via the algorithms in this chapter.

4.6.1 Channel Models and Estimation

The first channel model is

$$\mathbf{h} = \mathbf{s} \tag{4.19}$$

where $\mathbf{s} \in \mathbb{C}^N$ is a sparse vector with p nonzero entries. Second, we consider channels

$$\mathbf{h} = \mathbf{F}\mathbf{s} \quad (4.20)$$

which are p -sparse in the discrete Fourier transform (DFT) basis $\mathbf{F} \in \mathbb{C}^{N \times N}$. In both cases, the nonzero entries are drawn independently from $\mathcal{N}_{\mathbb{C}}(0, \frac{1}{p})$. Third, channels

$$\mathbf{h} = \sum_{k=1}^p s_k \mathbf{a}(\theta_k) \quad (4.21)$$

are modeled as a weighted sum of steering vectors

$$\mathbf{a}(\theta) = [1, e^{j\pi \sin(\theta)}, \dots, e^{j\pi(N-1)\sin(\theta)}]^T. \quad (4.22)$$

The path angles θ_k are drawn uniformly in $[0, 2\pi]$ and the path gains s_k are drawn from $\mathcal{N}_{\mathbb{C}}(0, \frac{1}{p})$.

As explained in Section 2.3, OMP makes use of a dictionary for channel estimation. For the channel models (4.19) and (4.20), we can use the dictionaries $\Psi = \mathbf{I}_N$ and $\Psi = \mathbf{F}$, respectively, to express the channel as $\mathbf{h} = \Psi\mathbf{s}$. For the channel model (4.21), we define a dictionary

$$\Psi_L = [\mathbf{a}(\hat{\theta}_1) \quad \dots \quad \mathbf{a}(\hat{\theta}_L)] \in \mathbb{C}^{N \times L} \quad (4.23)$$

of L steering vectors which correspond to L equidistantly sampled angles $\hat{\theta}_l$ in $[0, 2\pi]$. We can then approximate a channel by means of a linear combination of the steering vectors in Ψ_L : $\mathbf{h} \approx \Psi_L \mathbf{s}$ where $\mathbf{s} \in \mathbb{C}^L$ is a p -sparse vector, cf., e.g., [44]. In all three cases, we write

$$\mathbf{y} = \mathbf{A}\Psi\mathbf{s} + \mathbf{n} \quad (4.24)$$

with either $\Psi = \mathbf{I}_N$, $\Psi = \mathbf{F}$, or $\Psi = \Psi_L$, and we use OMP from Algorithm 1 with $\mathbf{C} = \mathbf{A}\Psi$ to recover a p -sparse estimate $\hat{\mathbf{s}}$ of \mathbf{s} so that the channel estimate is given by $\hat{\mathbf{h}} = \Psi\hat{\mathbf{s}}$, see also Section 2.3. In summary, the three channel models represent: (i) sparse vectors, (ii) vectors which are sparse in a nontrivial basis, and (iii) vectors which are approximately sparse.

4.6.2 Simulation Setup

The noise in (4.1) is Gaussian: $\mathbf{n} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{0}_m, \sigma^2 \mathbf{I}_m)$. We define the SNR as

$$\text{SNR} = \frac{\|\mathbf{A}\mathbf{h}\|^2}{\mathbb{E}[\|\mathbf{n}\|^2]} = \frac{\|\mathbf{A}\mathbf{h}\|^2}{m\sigma^2} \quad (4.25)$$

where σ^2 is adjusted for every pair (\mathbf{A}, \mathbf{h}) to achieve the desired SNR. To measure the channel estimation performance, we generate $T_{\text{tst}} = 10^4$ test channel samples $\{\mathbf{h}_t\}_{t=1}^{T_{\text{tst}}}$, generate corresponding observations according to (4.1), and compute channel estimates $\{\hat{\mathbf{h}}_t\}_{t=1}^{T_{\text{tst}}}$ using OMP. While the channel estimation algorithm is always OMP, there is a difference in the matrix \mathbf{A} that is used to obtain the observations. We distinguish five cases for a given channel model:

- *MMD*: A single observation matrix is learned via Algorithm 4 and all observations at all SNRs are computed using this one matrix.
- *PGDM*: A single observation matrix is obtained via Algorithm 7 and all observations at all SNRs are computed using this one matrix.
- *MMD with PGDM init.*: Algorithm 4 is initialized with the matrix obtained via Algorithm 7 and all observations at all SNRs are computed using the resulting matrix.
- *LBCS*: For every SNR, a matrix is obtained via Algorithm 6 and all observations at a given SNR are computed using the corresponding matrix.
- *random*: For every SNR and every channel \mathbf{h}_t , a new random constant modulus matrix with entries according to (2.9) is drawn and used to compute a corresponding observation.

The channel estimation performance is always measured in terms of relative MSE: $\sum_{t=1}^{T_{\text{tst}}} \|\mathbf{h}_t - \hat{\mathbf{h}}_t\|^2 / \sum_{t=1}^{T_{\text{tst}}} \|\mathbf{h}_t\|^2$. The whole procedure is summarized in Algorithm 8 for a given SNR. For the *random* case, a new random constant modulus matrix is drawn between Line 2 and Line 3 in Algorithm 8. It is worth pointing out that the random case has a higher computational complexity because the product $\mathbf{A}\Psi$ has to be computed for every channel, see Line 5 in Algorithm 8, whereas it can be precomputed in all other cases.

4.6.3 Simulation Parameters

When OMP is used in conjunction with the channel model in (4.21), we use the dictionary $\Psi = \Psi_L$ from (4.23) with $L = 16N$.

Algorithm 6 is used with $T_{\text{tr}} = 10^5$ and $T_{\text{val}} = 10^3$. The OMP algorithm is used for the evaluation in Line 7. The evaluation parameters are the SNR, the sparsity p , and the OMP dictionary Ψ . We always run $I = 100$ iterations.

Algorithm 8 Evaluation**Require:** observation matrix \mathbf{A} , data $\{\mathbf{h}_t\}_{t=1}^{T_{\text{tst}}}$ **Require:** SNR, sparsity p , dictionary Ψ

- 1: $e \leftarrow 0, \Delta \leftarrow 0$
- 2: **for** $t = 1$ to T_{tst} **do**
- 3: using (SNR, \mathbf{A}, \mathbf{h}_t), get σ via (4.25), draw $\mathbf{n} \sim \mathcal{N}_{\mathbb{C}}(\mathbf{0}_m, \sigma^2 \mathbf{I}_m)$
- 4: $\mathbf{y} \leftarrow \mathbf{A}\mathbf{h}_t + \mathbf{n}$ *// compute the observation*
- 5: $\hat{\mathbf{h}} \leftarrow \Psi \text{OMP}(\mathbf{A}\Psi, \mathbf{y}, p)$ *// estimate the channel \mathbf{h}_t using OMP*
- 6: $e \leftarrow e + \|\mathbf{h}_t\|^2, \Delta \leftarrow \Delta + \|\hat{\mathbf{h}} - \mathbf{h}_t\|^2$ *// square error*
- 7: **end for**
- 8: **return** $\frac{\Delta}{e}$ *// relative mean square error*

Algorithm 7 is used with the OMP dictionary Ψ and with an initial $\mathbf{A}^{(0)}$ as described in Section 4.4. The number of iterations is always $I = 500$ and the remaining hyperparameters α and ζ are determined as follows. First, we run Algorithm 7 with all combinations of

$$\alpha \in \{1.0 + i \cdot 0.1 : i \in \{0, 1, \dots, 10\}\} \cup \{5.0, 10.0, 50.0, 100.0\} \quad (4.26)$$

and $\zeta \in \{10^i : i \in \{-7, \dots, 1\}\}$ and determine the best (smallest coherence, see Section 4.4) combination (α^*, ζ^*) . Then, we run Algorithm 7 another 10 times with $\alpha \in [0.9\alpha^*, 1.1\alpha^*]$ and $\zeta \in [0.9\zeta^*, 1.1\zeta^*]$ chosen randomly (random search, cf. Section 2.7). Among all these 145 evaluations of Algorithm 7, the matrix which leads to the smallest coherence is used for the numerical experiments.

For Algorithm 4, we use MMD_k in (2.16) with the kernel

$$k : \mathbb{R}^{2m} \times \mathbb{R}^{2m}, (\mathbf{x}, \mathbf{y}) \mapsto k(\mathbf{x}, \mathbf{y}) = \sum_{\sigma \in \mathcal{S}} k_{\sigma}(\mathbf{x}, \mathbf{y}) \quad (4.27)$$

where k_{σ} is the Gaussian kernel from (2.15). The kernel parameter $\sigma > 0$ is difficult to choose optimally [20]. However, a mixture of kernels [20] can lead to a satisfying performance. Here, we use $\mathcal{S} = \{2, 5, 10, 20, 40, 80\}$.

We make use of Pytorch [64] to implement the stochastic gradient optimization in Algorithm 4. The employed optimizer is Adam [65]. The stochastic gradient algorithm has three hyperparameters: batch size T , learning rate (gradient step size) l_r , and exponential learning rate decay β . We determine these hyperparameters via random search [24]. As explained in Section 2.7, the strategy is to draw $(T, l_r, \beta) \in [150, 1500] \times [10^{-6}, 5 \cdot 10^{-3}] \times [0.94, 1]$ randomly and run Algorithm 4. In our case,

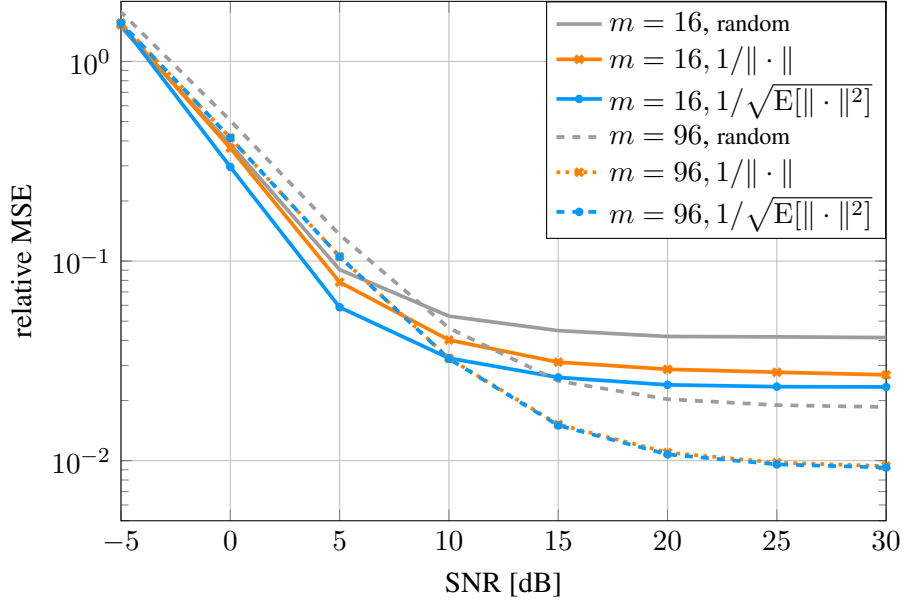


Figure 4.2: We have $\mathbf{A} \in \mathcal{M}_{\text{const}}^{m \times 128}$. Dashed curves show evaluation with model (4.21) for $m = 96$ and $p = 10$, solid curves for $m = 16$ and $p = 1$. Evaluation with random matrices is displayed in dark gray (“random”). Blue curves refer to learning a matrix with training data normalization $\mathbf{h}/\|\mathbf{h}\|$, orange curves refer to normalization $\mathbf{h}/\sqrt{\mathbb{E}[\|\mathbf{h}\|^2]}$.

this is repeated 100 times and validation data is then used to pick the best of the 100 so-obtained matrices. The termination criterion for Algorithm 4 is early stopping [66] with a patience of 150. We always have $T_{\text{tr}} = 5 \cdot 10^4$ training data in all cases and settings. However, instead of normalizing the data such that we work with $\{\frac{\mathbf{h}_t}{\|\mathbf{h}_t\|}\}_{t=1}^{T_{\text{tr}}}$, we normalize the training data as follows:

$$\bar{\mathbf{h}}_t = \frac{\mathbf{h}_t}{\sqrt{\frac{1}{T_{\text{tr}}} \sum_{i=1}^{T_{\text{tr}}} \|\mathbf{h}_i\|^2}}, \quad \text{for } t = 1, 2, \dots, T_{\text{tr}}. \quad (4.28)$$

The reason for this is explained in Section 4.6.4.

4.6.4 Training Data Normalization

The motivation to propose Algorithm 4 is to write the RIP condition as

$$(1 - \delta) \leq \left\| \mathbf{A} \frac{\mathbf{h}}{\|\mathbf{h}\|} \right\|^2 \leq (1 + \delta) \quad (4.29)$$

for all $\mathbf{h} \in \mathcal{H}$. The proposed optimization problem

$$\min_{\mathbf{A} \in \mathcal{M}_{\text{const}}^{m \times N}} \text{MMD}_k^2(p_{\mathbf{u}}, q_{\mathbf{A}\bar{\mathbf{h}}}) \quad (4.30)$$

aims to match the distributions $q_{\mathbf{A}\bar{\mathbf{h}}}$ to $p_{\mathbf{u}}$, where $p_{\mathbf{u}}$ denotes the uniform distribution on the unit hypersphere in \mathbb{C}^m and where $\bar{\mathbf{h}} = \mathbf{h}/\|\mathbf{h}\|$. We can generally not expect to achieve a minimum of zero. This is because both $\|\bar{\mathbf{h}}\| = 1$ and $\|\mathbf{u}\| = 1$ hold so that a minimum of zero implies $\|\mathbf{A}\bar{\mathbf{h}}\| = 1$ which corresponds to a matrix with restricted isometry constant $\delta = 0$. However, in particular, when the number m of rows is small, we can expect to require a larger $\delta > 0$.

It might now be beneficial to introduce some variation in the norm of $\bar{\mathbf{h}}$ to allow $\|\mathbf{A}\bar{\mathbf{h}}\|$ to be closer to one. Since the optimization problem not only aims to match the norm but also aims for isotropy, varying the norm may allow for a better isotropy. These considerations lead us to compare the training of Algorithm 4 (MMD) using two different training data normalizations: (i) all training channels \mathbf{h}_t are normalized as $\mathbf{h}_t/\|\mathbf{h}_t\|$, and (ii) all training channels are normalized as described in (4.28). In the latter case, the training data only has “on average” a norm equal to one.

Figure 4.2 displays the result of running Algorithm 4 (MMD) with the two different normalizations. We consider two edge cases: (i) a small number m of rows and only one path, and (ii) a large number m of rows and many paths. In both cases, the channel model from (4.21) is used. We observe that the training data normalization (4.28) is at least as good as or better than the normalization $\mathbf{h}_t/\|\mathbf{h}_t\|$. The difference is particularly pronounced when the number of rows is small. This suits the intuition that the restricted isometry constant δ is potentially larger for smaller m . In these cases, alleviating the norm one constraint can be beneficial. On the other hand, as m increases, the matrix \mathbf{A} gets closer to a square matrix where an isometry with $\delta = 0$ is possible. Due to the observed effect, in all other experiments involving Algorithm 4 (MMD), the training data normalization (4.28) is employed.

4.6.5 Simulation Results

In all simulations, we consider constant modulus matrices $\mathbf{A} \in \mathcal{M}_{\text{const}}^{m \times N}$ with $m \in \{16, 32, 64, 96\}$ and $N = 128$ and we vary the number of paths (i.e., the sparsity value) $p \in \{1, 5, 10\}$. First, results with the channel models (4.19) and (4.20) are discussed, where the channels have a sparse representation. Thereafter, we discuss the case of approximate sparsity with channel model (4.21).

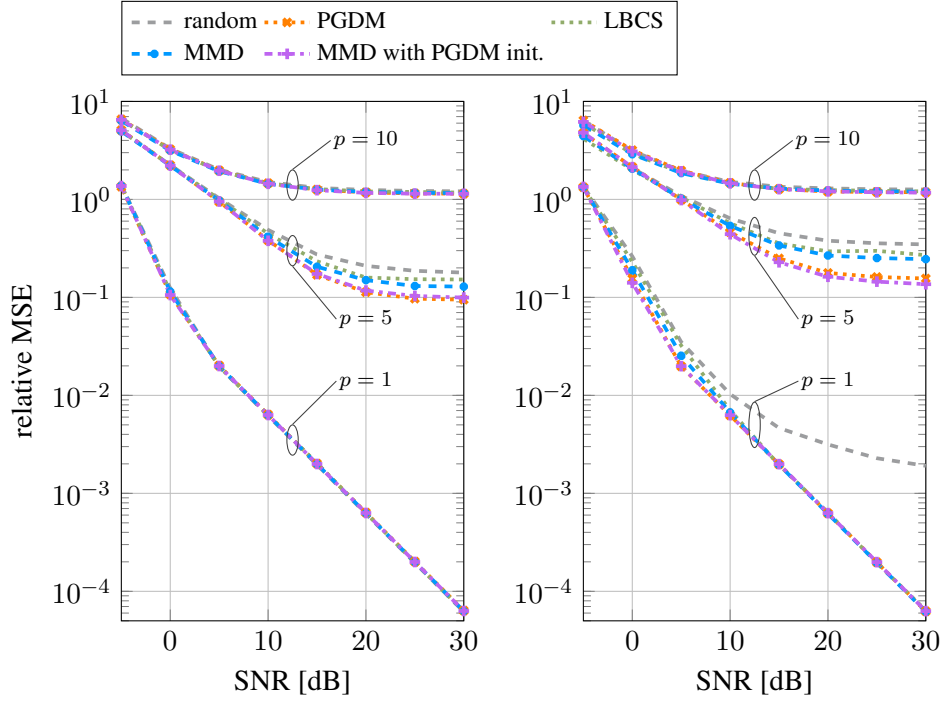


Figure 4.3: Evaluation with channel models (4.19) (left) and (4.20) (right). We have $\mathbf{A} \in \mathcal{M}_{\text{const}}^{16 \times 128}$ and $p \in \{1, 5, 10\}$. The legend is explained in Section 4.6.2.

4.6.5.1 Sparse channels

The left plot in Figure 4.3 shows evaluation with the channel model from (4.19) and the right plot shows evaluation with the channel model from (4.20). In both cases, the number of rows is $m = 16$. Already $p = 5$ seems to be a difficult setting and for $p = 10$, the relative MSE is above 1.0 which means that channel estimation is not possible. The channel models in Figure 4.3 have true sparsity with respect to an orthonormal basis ($\Psi = \mathbf{I}$ and $\Psi = \mathbf{F}$, respectively) in common so that the performance of random matrices should be similar in both cases. This is because Subgaussian random matrices (and constant modulus matrices are Subgaussian) are universal which means that they can be used for sparse recovery in an arbitrary orthonormal basis [10]. However, in particular, the case $p = 1$ behaves very differently in the two plots. One reason might be that both the random observation matrix as well as the basis $\Psi = \mathbf{F}$ have entries of the form $e^{j\phi}$ which can lead to a too large coherence between \mathbf{A} and Ψ . Interestingly, all non-random algorithms are able to find matrices which do not suffer from this effect. This emphasizes that universality can be traded for a performance improvement.

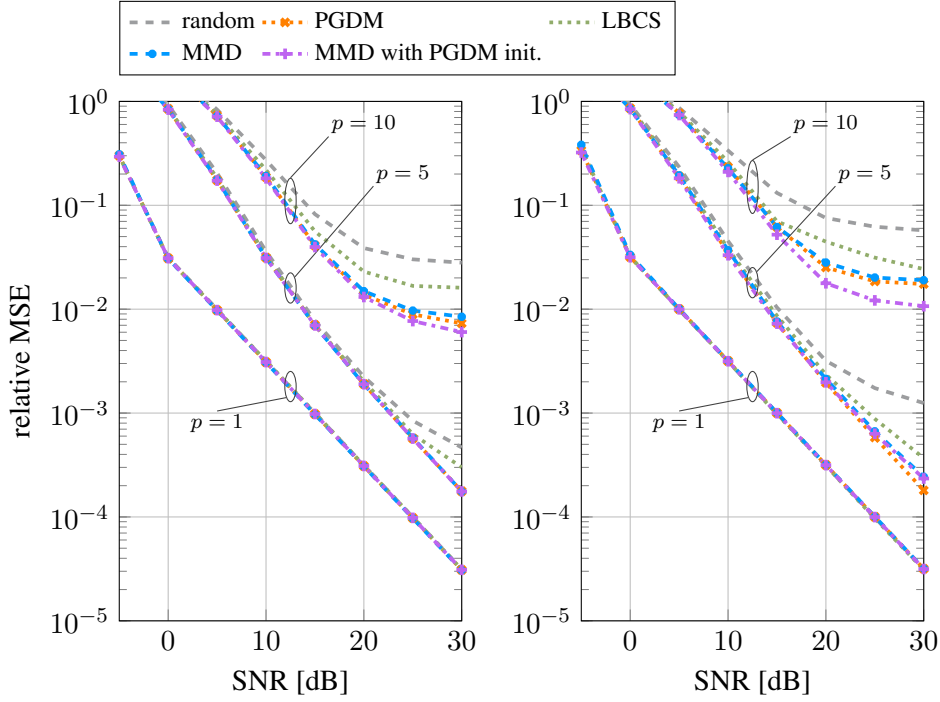


Figure 4.4: Evaluation with channel models (4.19) (left) and (4.20) (right). We have $\mathbf{A} \in \mathcal{M}_{\text{const}}^{32 \times 128}$ and $p \in \{1, 5, 10\}$. The legend is explained in Section 4.6.2.

For $p = 5$, all algorithms in Figure 4.3 outperform random matrices. Algorithm 7 (PGDM) is particularly strong.

Figure 4.4 considers the same setting as Figure 4.3 but with $m = 32$. The number m of rows appears to be large enough to facilitate channel estimation also for $p = 10$. Furthermore, the non-random algorithms yield good matrices for the cases $p = 1$ and $p = 5$ where for higher SNRs the remaining estimation error seems to be due to the presence of noise only. However, random matrices do not quite show this behavior in the case where $p = 5$ and the coherence problem seemingly still exists. In almost all cases, Algorithm 4 (MMD) and Algorithm 7 (PGDM) achieve a similar performance with Algorithm 7 (PGDM) in a slight advantage in some settings. Interestingly, the combination of using Algorithm 7 (PGDM) as an initialization of Algorithm 4 (MMD) can lead to a performance which is superior to the one which is obtained if both algorithms are used independently, see $p = 10$ in the left and right plot of Figure 4.4. Unfortunately, as the case $p = 5$ in the right plot of Figure 4.4 highlights for high SNRs, the combination is not guaranteed to yield the best matrix but it is at least as

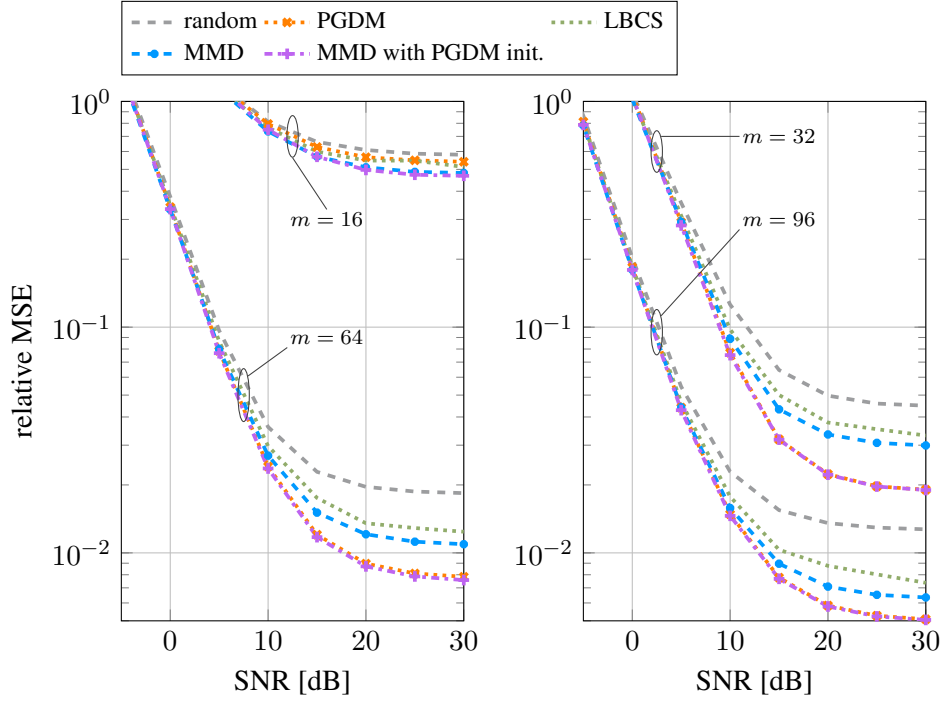


Figure 4.5: Evaluation with channel model (4.21). We have $\mathbf{A} \in \mathcal{M}_{\text{const}}^{m \times 128}$ for $m \in \{16, 32, 64, 96\}$ and $p = 5$. The legend is explained in Section 4.6.2.

good as the result of initializing Algorithm 4 (MMD) randomly.

The cases $m = 64$ and $m = 96$ lead to overlapping curves for all algorithms and for both channel models, and all three settings $p \in \{1, 5, 10\}$ show straight lines. These results are not visualized.

4.6.5.2 Approximately sparse channels

Figures 4.5 and 4.6 focus on the channel model in (4.21). A first observation is that in all of the considered settings, the curves saturate even at high SNRs. This is a result of the approximate sparsity: Whenever one of the steering vectors $\mathbf{a}(\theta_k)$ in the channel (4.21) corresponds to an angle $\theta_k \notin \{\hat{\theta}_1, \dots, \hat{\theta}_L\}$ which is not present in the dictionary Ψ_L from (4.23), then this channel cannot be estimated perfectly. There exist algorithms which focus on *grid-less* channel estimation (see, e.g., [67, Chapter 11] for an overview). However, since we are only interested in the behavior of different observation matrices, working with the on-grid OMP algorithm suffices for our purpose.

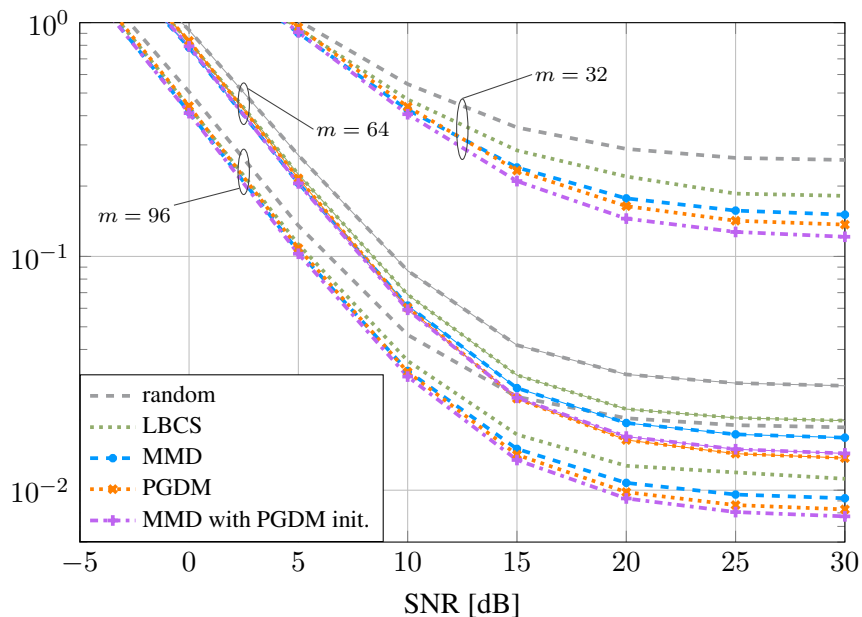


Figure 4.6: Evaluation with channel model (4.21). We have $\mathbf{A} \in \mathcal{M}_{\text{const}}^{m \times 128}$ for $m \in \{32, 64, 96\}$ and $p = 10$. The legend is explained in Section 4.6.2.

It is interesting to see how large the gap between the random matrices approach and all non-random algorithms is. For example, a comparison of the left and right plot in Figure 4.5 reveals that the best matrix with $m = 32$ rows (right plot) performs almost as well as random matrices with $m = 64$ rows (left plot). What is more, the best matrix with $m = 64$ rows (left plot) outperforms random matrices with $m = 96$ rows (right plot). In that sense, computing non-random matrices can save 32 rows (32 receiver chains) in both cases. Figure 4.6 displays a similar effect (random with $m = 96$ rows vs. non-random with $m = 64$ rows). In Figure 4.6, the combination of Algorithm 4 (MMD) and Algorithm 7 (PGDM) can improve upon the individual algorithms.

4.6.5.3 Conclusion

Generally, in all considered experiments, Algorithm 4 (MMD) and, in particular, the recently introduced Algorithm 7 (PGDM) show a strong performance. It seems to be a good strategy to initialize Algorithm 4 (MMD) with Algorithm 7 (PGDM) in the random search training procedure: In one half of the random searches, the initialization

would be random, in the other half, the initialization would be based on Algorithm 7 (PGDM). It is worth pointing out that Algorithm 7 (PGDM) makes use of a dictionary for the considered channels whereas this is not required for the training of Algorithm 4 (MMD). Thus, in contrast to Algorithm 7 (PGDM), in cases where no dictionary is available, Algorithm 4 (MMD) can still be used.

In all experiments, also Algorithm 6 (LBCS) improves upon random matrices. It lags behind the other non-random approaches. However, one advantage of the method is that it does not explicitly require a dictionary for the training procedure as well. Of course, this is only true if the evaluation algorithm used in Algorithm 6 (LBCS) does not require a dictionary. We consider such a setting in Chapter 5. As a consequence, if Algorithm 7 (PGDM) cannot be employed, it is interesting to initialize Algorithm 4 (MMD) with Algorithm 6 (LBCS) in one half of the random searches.

Learned Observation Matrix and the Gaussian Mixture Estimator

5

One motivation to prefer a deterministic observation matrix over random matrices is to save computational complexity during the channel estimation. For example, as already pointed out in Chapter 4, if the orthogonal matching pursuit (OMP) algorithm is used, employing random matrices requires the computation of the product $\mathbf{A}\Psi$ between a random observation matrix \mathbf{A} and the dictionary Ψ for every channel estimate, see Line 5 in Algorithm 8. In contrast, this product can be precomputed if a matrix from one of the three algorithms in Chapter 4 is used to learn an observation matrix.

Being able to precompute quantities involving \mathbf{A} is also interesting if the Gaussian mixture model (GMM) estimator from Chapter 3 is used. The GMM estimator computes a weighted sum of linear minimum mean square error (LMMSE) channel estimates (cf. (3.22)), and to compute an LMMSE channel estimate, the matrices $\tilde{\mathbf{C}}_k = \mathbf{A}\mathbf{C}_k\mathbf{A}^H + \Sigma$ from (3.16) need to be inverted, cf. (3.21b). Further, to compute the weights (or responsibilities), the determinants of all $\tilde{\mathbf{C}}_k$ are necessary, too, cf. (3.20). The inverses and determinants have to be computed for every channel estimate if random matrices are used. Again, in contrast, these quantities can be precomputed if \mathbf{A} is not a random matrix.

In Section 5.3, we combine the GMM estimator from Chapter 3 with the matrices obtained via the algorithms from Chapter 4. Further, we study what influence the evaluation function which is employed during the iterations of Algorithm 6 (LBCS) (see Line 7) has on the resulting observation matrix. To this end, we run Algorithm 6 (LBCS) using three different evaluation functions, cf. Section 5.4.

5.1 Simulation Setup

It seems natural to also study the matrices obtained via one of the algorithms from Chapter 4 in the context of a channel model which does not admit a sparse representation by construction. One such setting is presented in Figure 3.3 in Chapter 3 where we evaluate various channel estimators with a 3GPP channel model with one propagation

cluster and where we have $\mathbf{A} = \mathbf{I}$. In this setting, (i) the displayed compressive sensing algorithms show a considerable gap to the genie lower bound and (ii) the GMM estimator comes close to the genie lower bound. It is now interesting to investigate what happens if $\mathbf{A} = \mathbf{I}$ is replaced with a constant modulus observation matrix. In particular, it is interesting to see whether (i) the gap between the estimators remains when a constant modulus observation matrix is present, (ii) the GMM estimator still comes close to the genie lower bound, and (iii) we can again improve upon random matrices.

In the following simulations, the signal model is $\mathbf{y} = \mathbf{A}\mathbf{h} + \mathbf{n}$ with a constant modulus matrix $\mathbf{A} \in \mathcal{M}_{\text{const}}^{m \times N}$ with $N = 128$ and the channels are generated with one propagation cluster as explained in Section 3.4.1. We work with three channel estimation algorithms: the genie LMMSE estimator from (3.32), the genie OMP estimator (Algorithm 3), and the GMM estimator from (3.22). The GMM estimator always uses $K = 128$ components and $19 \cdot 10^4$ training data. The channel estimators are evaluated as described in Algorithm 8. The only difference is that the OMP algorithm in Line 5 is replaced with one of the three mentioned estimators. The algorithms from Chapter 4 are used as explained in Section 4.6.3. As in Chapter 4, when random observation matrices are used, a new random constant modulus matrix is drawn for every channel.

Algorithm 6 (LBCS) evaluates the current observation matrix on a validation set of channel samples, see Line 7. In Chapter 4, OMP (Algorithm 1) is used for channel estimation during the evaluation of Algorithm 6 (LBCS). In this chapter, we instead employ either the genie OMP estimator, the genie LMMSE estimator, or the GMM estimator during the evaluation in Line 7 in order to gauge the estimators' effect on the resulting observation matrix. We indicate this as follows:

- *LBCS(algo)*: Algorithm 6 (LBCS) uses the estimator

$$\text{algo} \in \{\text{gen. OMP, gen. LMMSE, GMM}\}$$

for the evaluation in Line 7 of its iterations.

5.2 Random Matrices

Figure 5.1 shows the behavior of the three channel estimators genie OMP, genie LMMSE, and GMM estimator if random constant modulus matrices with $m \in \{16, 32, 64, 96\}$ rows are used. All cases display a considerable gap between genie OMP and the other two estimators. Further, with $m = 32$ or more rows, the GMM

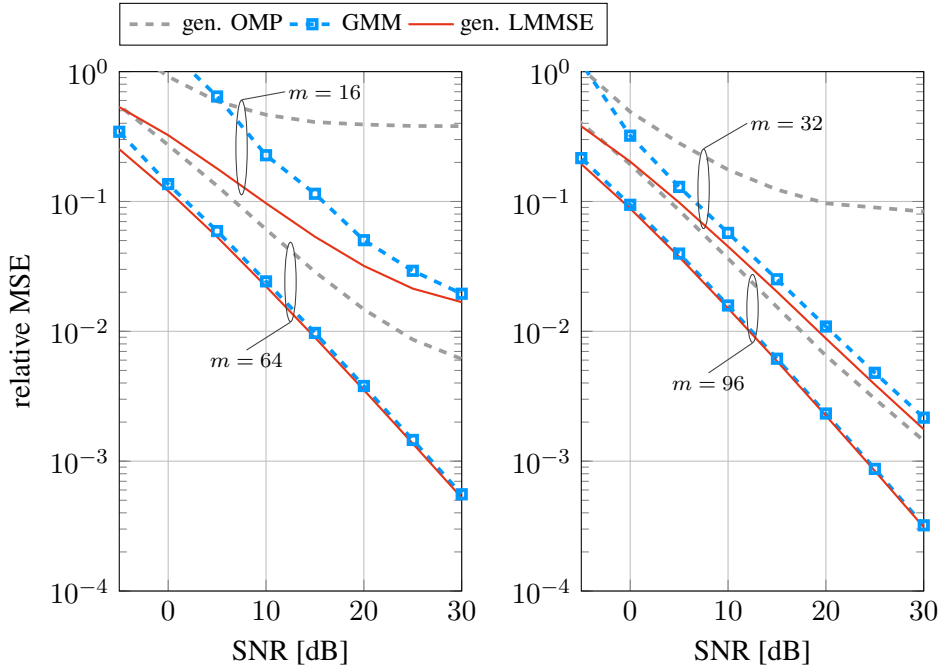


Figure 5.1: Channel estimation with the genie OMP, the genie LMMSE, and the GMM estimator. In all cases, random observation matrices $\mathbf{A} \in \mathcal{M}_{\text{const}}^{m \times 128}$ with $m \in \{16, 64\}$ (left) and $m \in \{32, 96\}$ (right) are used.

estimator starts to come close to the genie LMMSE estimator. This re-emphasizes the strength of the GMM estimator even for cases where the observation matrix is not invertible. Note, of the three algorithms, only the GMM estimator can be implemented in practice as is whereas the other two require genie knowledge.

5.3 Learned Matrices

Figure 5.2 compares channel estimation using the GMM estimator when different observation matrices designed according to one of the algorithms from Chapter 4 are used, cf. Section 4.6.2. As another reference curve, the genie LMMSE estimator in conjunction with random observation matrices is displayed as well. As a first observation, only in the cases $m = 16$ and $m = 32$ there is a noteworthy gap between genie LMMSE estimation and GMM estimation. Further, for $m = 64$ and $m = 96$, the GMM estimator with learned matrices can slightly outperform the genie LMMSE estimator with random matrices. In all cases, (slight) improvements upon random

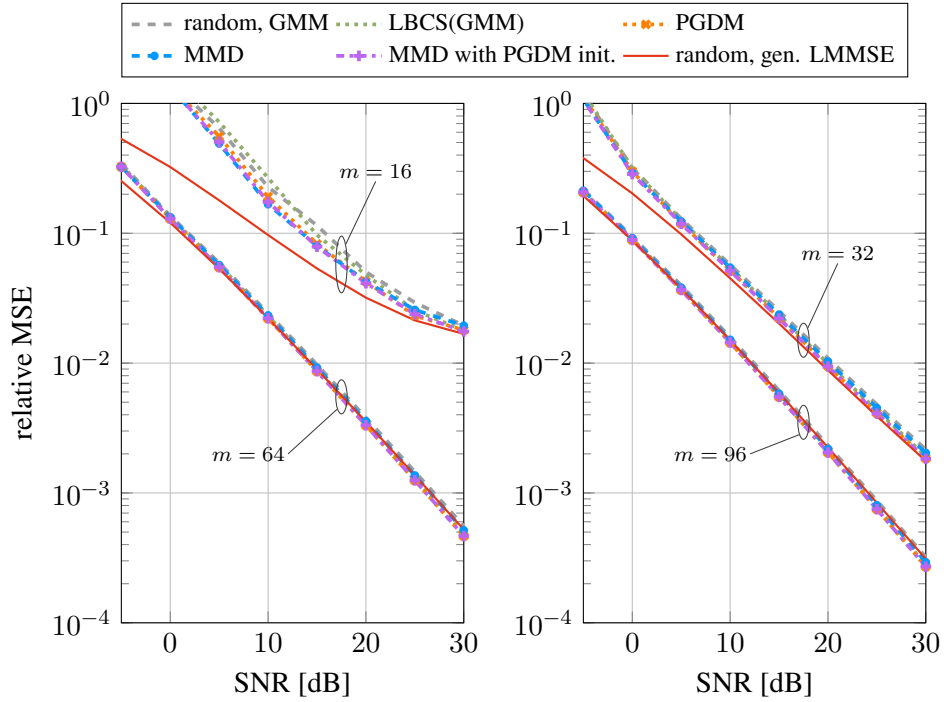


Figure 5.2: Channel estimation using the GMM estimator. In all cases, observation matrices $\mathbf{A} \in \mathcal{M}_{\text{const}}^{m \times 128}$ with $m \in \{16, 64\}$ (left) and $m \in \{32, 96\}$ (right) are used. “random” refers to employing random observation matrices with either the GMM estimator (“random, GMM”) or the genie LMMSE estimator (“random, gen. LMMSE”). The rest of the legend is explained in Section 4.6.2.

matrices are visible with only a little difference between the matrix design algorithms.

In Figure 5.3, we see the same experiment but with the genie OMP estimator instead of the GMM estimator. Generally, there is a significant gap between the genie OMP estimator and the genie LMMSE estimator. Further, learning an observation matrix can considerably improve the estimation performance in comparison to employing random matrices. Interestingly, for $m = 64$ and $m = 96$, the genie OMP does not yet show a saturation effect in the high signal-to-noise ratio (SNR) region if learned matrices are used. Figure 5.4 shows the same experiment but with the genie LMMSE estimator. In this case, there is not much difference between all methods but minor improvements can still be seen compared to random matrices.

Lastly, as a summary, Figure 5.5 is similar to Figure 5.1 in that it compares the three different estimators. In Figure 5.1, the observation matrices are always random

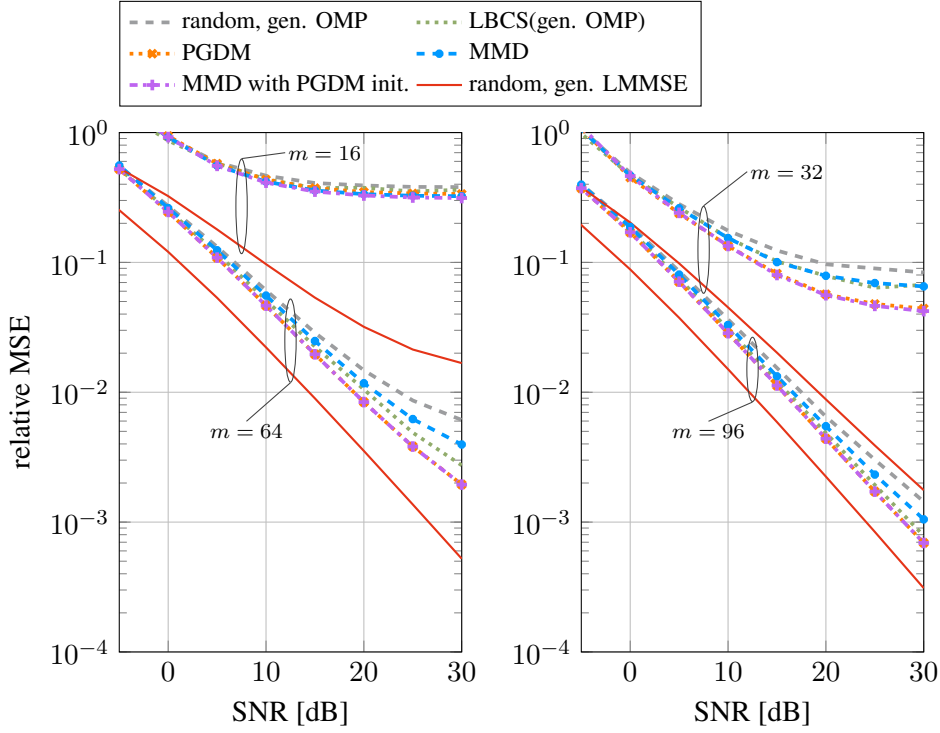


Figure 5.3: Channel estimation using the genie OMP estimator. In all cases, observation matrices $\mathbf{A} \in \mathcal{M}_{\text{const}}^{m \times 128}$ with $m \in \{16, 64\}$ (left) and $m \in \{32, 96\}$ (right) are used. “random” refers to employing random observation matrices with either the genie OMP estimator (“random, gen. OMP”) or the genie LMMSE estimator (“random, gen. LMMSE”). The rest of the legend is explained in Section 4.6.2.

whereas in Figure 5.5, the observation matrices are always determined via Algorithm 4 (MMD) initialized with Algorithm 7 (PGDM). The red curves in both figures are identical: the genie LMMSE estimator is employed with random observation matrices. It is interesting to see that essentially all curves improve from Figure 5.1 to Figure 5.5 and they come closer to the “random, gen. LMMSE” baseline, which highlights the importance of a good observation matrix.

Overall, the difference between the performance of random matrices and the learned matrices seems to be smaller for better estimators. In the shown plots, the genie LMMSE estimator is as good as or better than the GMM estimator which, in turn, is as good as or better than the genie OMP estimator. At the same time, the gap between random matrices and learned matrices increases from genie LMMSE estimation

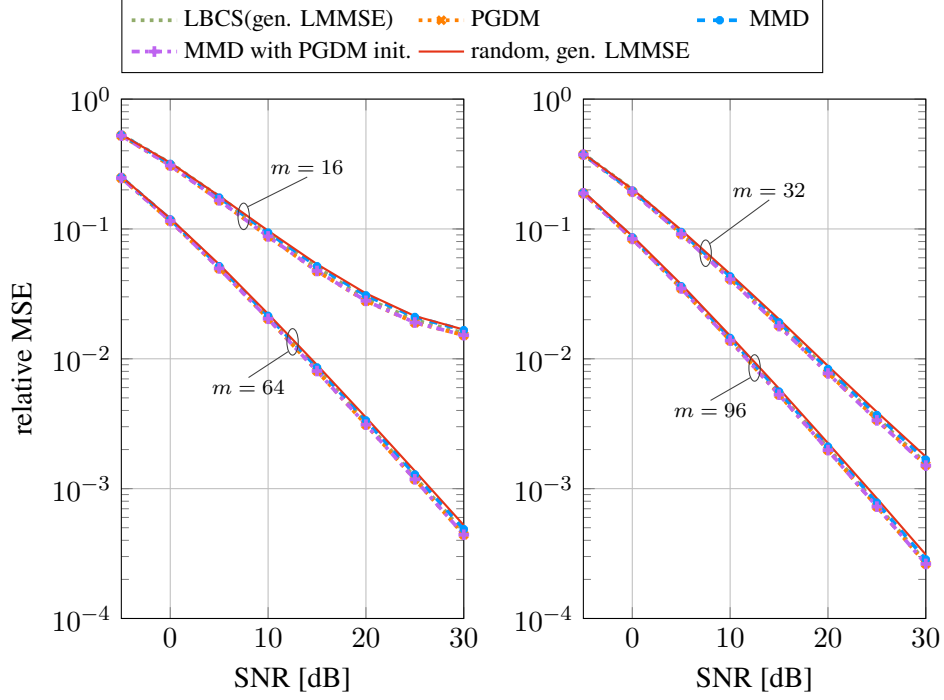


Figure 5.4: Channel estimation using the genie LMMSE estimator. In all cases, observation matrices $\mathbf{A} \in \mathcal{M}_{\text{const}}^{m \times 128}$ with $m \in \{16, 64\}$ (left) and $m \in \{32, 96\}$ (right) are used. “random, gen. LMMSE” refers to employing random observation matrices with the genie LMMSE estimator. The rest of the legend is explained in Section 4.6.2.

through GMM estimation to genie OMP estimation. Nonetheless, all estimators benefit from learned matrices. The recently introduced projected gradient descent method (PGDM) (Algorithm 7) from [53] is strong, indicating that the codebook Ψ_L from (4.23) is well-suited for the considered data set and problem setting. Generally, a sensible strategy seems to be to initialize Algorithm 4 (MMD) with Algorithm 7 (PGDM) in half of the random searches as long as a dictionary for the considered data is available. If this is not the case, Algorithm 5 (LBCS) should serve as an initializer.

5.4 Different Evaluation Functions for Algorithm 6

Lastly, we look at the influence of the evaluation function on the final result of Algorithm 6 (LBCS). We use the notation described in Section 5.1. The left plot of

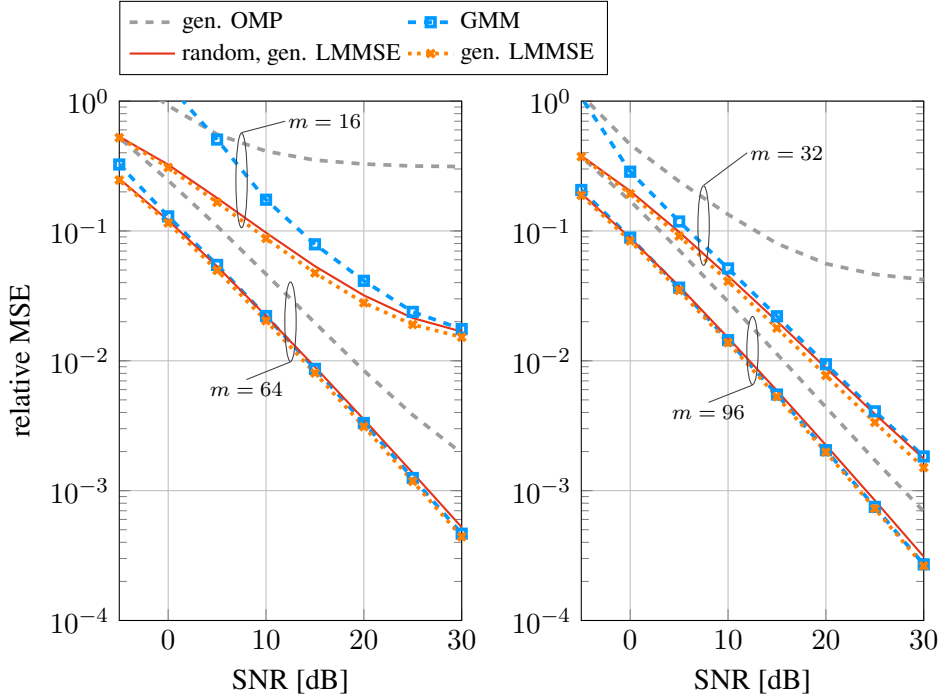


Figure 5.5: Channel estimation with the genie OMP, the genie LMMSE, and the GMM estimator. In all cases, observation matrices $\mathbf{A} \in \mathcal{M}_{\text{const}}^{m \times 128}$ with $m \in \{16, 64\}$ (left) and $m \in \{32, 96\}$ (right) are used. “random, gen. LMMSE” refers to employing random observation matrices with the genie LMMSE estimator. In the remaining cases, the observation matrices are determined via *MMD with PGDM init.*, cf. Section 4.6.2.

Figure 5.6 shows the channel estimation with the genie OMP estimator. While the difference between the three options LBCS(gen. OMP), LBCS(gen. LMMSE), and LBCS(GMM) is small, LBCS(gen. OMP) tends to perform best. In the right plot of Figure 5.6, the channel estimation is performed with the genie LMMSE estimator. In this case, the difference between the three options is even smaller but LBCS(gen. LMMSE) tends to perform best. Finally, Figure 5.7 shows the channel estimation with the GMM estimator. Here, for $m = 16$, no clear preference can be seen, and in the other cases, the three options hardly differ. In summary, we can observe the tendency that it is beneficial to use the same channel estimator during the iterations of Algorithm 6 (LBCS) as is applied afterwards when the resulting matrix is used for channel estimation. Hence, the appropriate version of Algorithm 6 (LBCS) is displayed in Figures 5.1 to 5.4. However, the differences in the experiments are minute.

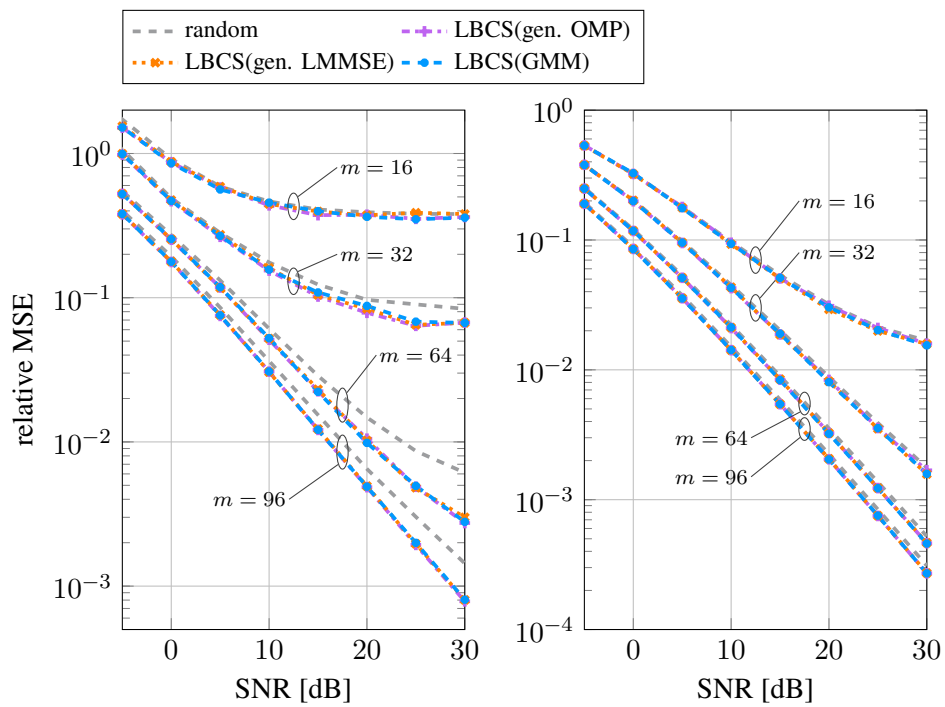


Figure 5.6: Algorithm 6 (LBCS) uses the channel estimator shown in brackets during its iterations. Thereafter, for the evaluation in the plots, channels are estimated with genie OMP (left plot) or genie LMMSE (right plot).

5.4 Different Evaluation Functions for Algorithm 6

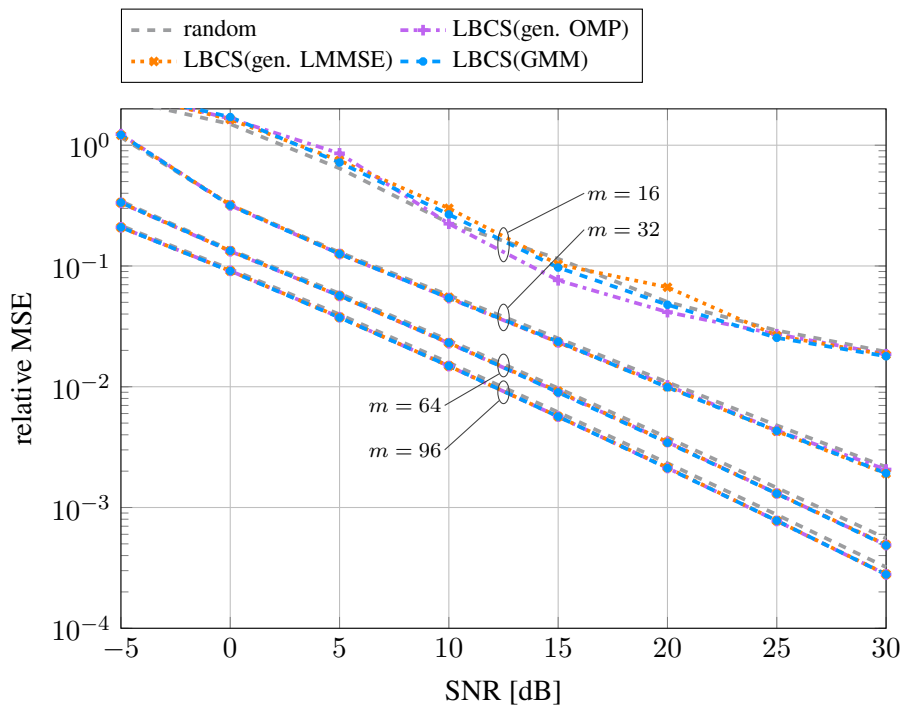


Figure 5.7: Algorithm 6 (LBCS) uses the channel estimator shown in brackets during its iterations. Thereafter, for the evaluation in the plot, channels are estimated with the GMM estimator.

Outlook 6

The following topics related to Chapters 3 and 4 might be worth studying in the future.

6.1 Asymptotically Optimal Channel Estimation

The approach in Chapter 3 is to approximate the channel probability density function (PDF) by means of a uniformly convergent sequence of PDFs. Thereafter, for every sequence element, a conditional mean estimator (CME) is defined, which yields a sequence of CMEs. Theorem 2 then shows under which conditions the sequence of CMEs converges to the optimal CME.

Interestingly, in Theorem 2, there are no further conditions on the sequence of PDFs. One example of such a sequence is given by Gaussian mixture models (GMMs). However, no properties specific to GMMs are used in the proof. Thus, it is possible and interesting to investigate other universal approximators. This could potentially lead to estimators with, e.g., a lower computational complexity or with a better performance or which require a smaller amount of channel training data. Other universal approximators are mentioned in [1]. The challenge lies in fitting the corresponding PDF based on channel training data or in computing the corresponding CME $\hat{\mathbf{h}}^{(K)}$ in closed form.

If the computation of the CME $\hat{\mathbf{h}}^{(K)}$ is based on GMMs, a weighted sum of K linear minimum mean square error (LMMSE) estimators is computed (cf. (3.22)):

$$\hat{\mathbf{h}}^{(K)} : \mathbf{y} \mapsto \hat{\mathbf{h}}^{(K)}(\mathbf{y}) = \sum_{k=1}^K p(k | \mathbf{y}^{(K)} = \mathbf{y}) \hat{\mathbf{h}}_{\text{LMMSE},k}(\mathbf{y}). \quad (6.1)$$

Here, the responsibility $p(k | \mathbf{y}^{(K)} = \mathbf{y})$ (cf. (3.20)) is the probability that the current observation \mathbf{y} was sampled from the k th GMM component. In view of the computational complexity, we could compute only the sum of those $K' < K$ summands in (6.1) which correspond to the K' largest probabilities. It is, in particular, interesting to investigate the estimator's behavior when $K' \in \{1, 2, \dots, K\}$ is varied. Another question in this context is how sensitive the estimator (6.1) is with respect to the responsibilities. For instance, in order to avoid evaluating K Gaussian PDFs to

compute the exact responsibilities (cf. (3.20)), we could try to find a classifier which aims to compute

$$\mathbf{y} \mapsto \begin{pmatrix} p(1 | \mathbf{y}^{(K)} = \mathbf{y}) \\ \vdots \\ p(K | \mathbf{y}^{(K)} = \mathbf{y}) \end{pmatrix} \quad (6.2)$$

with less complexity. For example, using a neural network-based classifier, all responsibilities could be computed simultaneously in a single forward pass. This might save computation time while hardly affecting the sum in (6.1).

Abstractly, in Chapter 3, we study a function $\hat{\mathbf{h}}$ whose computation requires an analytic expression of a PDF. Since this PDF is not given, we approximate it by means of a uniformly convergent sequence and compute the corresponding approximation-based functions $\hat{\mathbf{h}}^{(K)}$. Then we ask whether the sequence $(\hat{\mathbf{h}}^{(K)})_{K=1}^{\infty}$ converges to $\hat{\mathbf{h}}$. This approach could work with other functions (not necessarily channel estimators) as well. It is particularly interesting when GMMs are used as approximating sequence. In the case of computing the CME, we can directly make use of the well-known LMMSE formula because conditioned on a GMM component, we are in a Gaussian setting. More generally, the conditional Gaussianity of GMMs can make it possible to use available closed-form solutions in order to obtain a practically useful approximation of the function $\hat{\mathbf{h}}$ of interest.

6.2 Learning a Compressive Sensing Matrix

In Chapter 4, we study a distribution matching problem which yields a constant modulus observation matrix. For an arbitrary (not necessarily constant modulus) matrix $\mathbf{A} \in \mathbb{C}^{m \times N}$, the optimization problem reads

$$\min_{\mathbf{A} \in \mathbb{C}^{m \times N}} \text{MMD}_k^2 \left(\{\text{stk}(\mathbf{u}_t)\}_{t=1}^{T_{\text{tr}}}, \{\text{stk}(\mathbf{A}\bar{\mathbf{h}}_t)\}_{t=1}^{T_{\text{tr}}} \right). \quad (6.3)$$

Using

$$\text{stk}(\mathbf{A}\bar{\mathbf{h}}) = \begin{bmatrix} \Re(\mathbf{A}\bar{\mathbf{h}}) \\ \Im(\mathbf{A}\bar{\mathbf{h}}) \end{bmatrix} = \underbrace{\begin{bmatrix} \Re(\mathbf{A}) & -\Im(\mathbf{A}) \\ \Im(\mathbf{A}) & \Re(\mathbf{A}) \end{bmatrix}}_{=\tilde{\mathbf{A}}} \begin{bmatrix} \Re(\bar{\mathbf{h}}) \\ \Im(\bar{\mathbf{h}}) \end{bmatrix} \quad (6.4)$$

it can again be expressed by means of real quantities. Interestingly, structural constraints like, e.g., a Toeplitz structure, a circulant structure, or a block diagonal structure can easily be incorporated into the learning algorithm which solves (6.3). To this end, the number of optimization variables in $\tilde{\mathbf{A}}$ is reduced by forcing certain elements to be equal, as dictated by the sought structure. This technique is known as *weight tying*

in machine learning, see, e.g., [66]. A motivation for doing this is, for example, [68] where Toeplitz matrices are used in the context of channel estimation and where it is shown that randomly drawing such matrices can lead to matrices with the restricted isometry property (RIP). Here, we might again be interested in obtaining one fixed matrix instead of drawing random ones.

More generally, the idea in Chapter 4 is to match the distribution of $\mathbf{A}\bar{\mathbf{h}}$ to the distribution of \mathbf{u} where $\bar{\mathbf{h}} = \mathbf{h}/\|\mathbf{h}\|$ are normalized channels and where \mathbf{u} is uniformly distributed on the m -dimensional unit hypersphere. It might be interesting to investigate other distribution matching methods as well. One example is the generative adversarial network (GAN) [69]. In the context of Chapter 4, the normalized channel samples $\bar{\mathbf{h}}_t$ are the GAN's noise samples and the observation matrix \mathbf{A} is the generator. The GAN's discriminator then aims to distinguish between fake samples $\mathbf{A}\bar{\mathbf{h}}_t$ and real samples \mathbf{u}_t .

Proof of Theorem 2 A

The proof of Theorem 2 makes use of and is presented after the following lemma and can also be found in [26].

Lemma 1. For an arbitrary $\mathbf{y} \in \mathbb{R}^N$, it holds

$$\int \|\mathbf{h}\| f_{\mathbf{n}}(\mathbf{y} - \mathbf{A}\mathbf{h}) d\mathbf{h} \leq \sqrt{\det(\mathbf{A}^{-1}\mathbf{A}^{-\text{T}})} \sqrt{\|\mathbf{A}^{-1}\mathbf{y}\|^2 + \text{trace}(\mathbf{A}^{-1}\boldsymbol{\Sigma}\mathbf{A}^{-\text{T}})}.$$

Proof. Recall that $f_{\mathbf{n}} = \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma})$ denotes a Gaussian probability density function (PDF). We have

$$f_{\mathbf{n}}(\mathbf{y} - \mathbf{A}\mathbf{h}) = \frac{\exp(-\frac{1}{2}(\mathbf{y} - \mathbf{A}\mathbf{h})^{\text{T}}\boldsymbol{\Sigma}^{-1}(\mathbf{y} - \mathbf{A}\mathbf{h}))}{\sqrt{(2\pi)^N \det(\boldsymbol{\Sigma})}} \quad (\text{A.1})$$

$$= \frac{\exp(-\frac{1}{2}(\mathbf{A}^{-1}\mathbf{y} - \mathbf{h})^{\text{T}}\mathbf{A}^{\text{T}}\boldsymbol{\Sigma}^{-1}\mathbf{A}(\mathbf{A}^{-1}\mathbf{y} - \mathbf{h}))}{\sqrt{(2\pi)^N \det(\boldsymbol{\Sigma})}} \quad (\text{A.2})$$

$$= \frac{\exp(-\frac{1}{2}(\mathbf{h} - \mathbf{A}^{-1}\mathbf{y})^{\text{T}}(\mathbf{A}^{-1}\boldsymbol{\Sigma}\mathbf{A}^{-\text{T}})^{-1}(\mathbf{h} - \mathbf{A}^{-1}\mathbf{y}))}{\sqrt{(2\pi)^N \det(\boldsymbol{\Sigma})}}. \quad (\text{A.3})$$

Therefore,

$$\frac{f_{\mathbf{n}}(\mathbf{y} - \mathbf{A}\mathbf{h})}{\sqrt{\det(\mathbf{A}^{-1}\mathbf{A}^{-\text{T}})}} = \mathcal{N}(\mathbf{h}; \tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}) \quad (\text{A.4})$$

is a Gaussian PDF with mean vector $\tilde{\boldsymbol{\mu}} = \mathbf{A}^{-1}\mathbf{y}$ and covariance matrix $\tilde{\boldsymbol{\Sigma}} = \mathbf{A}^{-1}\boldsymbol{\Sigma}\mathbf{A}^{-\text{T}}$. Let $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ be a standard Gaussian random variable and let $\tilde{\boldsymbol{\Sigma}}^{\frac{1}{2}}$ be a square root of the covariance matrix $\tilde{\boldsymbol{\Sigma}} = \tilde{\boldsymbol{\Sigma}}^{\frac{\text{T}}{2}}\tilde{\boldsymbol{\Sigma}}^{\frac{1}{2}}$. Then, we are interested in computing

$$\int \|\mathbf{h}\| f_{\mathbf{n}}(\mathbf{y} - \mathbf{A}\mathbf{h}) d\mathbf{h} = \sqrt{\det(\mathbf{A}^{-1}\mathbf{A}^{-\text{T}})} \int \|\mathbf{h}\| \frac{f_{\mathbf{n}}(\mathbf{y} - \mathbf{A}\mathbf{h})}{\sqrt{\det(\mathbf{A}^{-1}\mathbf{A}^{-\text{T}})}} d\mathbf{h} \quad (\text{A.5})$$

$$= \sqrt{\det(\mathbf{A}^{-1}\mathbf{A}^{-\text{T}})} \int \|\mathbf{h}\| \mathcal{N}(\mathbf{h}; \tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}}) d\mathbf{h} \quad (\text{A.6})$$

$$= \sqrt{\det(\mathbf{A}^{-1}\mathbf{A}^{-\text{T}})} \mathbb{E}_{\mathbf{w}}[\|\tilde{\boldsymbol{\mu}} + \tilde{\boldsymbol{\Sigma}}^{\frac{1}{2}}\mathbf{w}\|] \quad (\text{A.7})$$

where we take the expectation with respect to the standard Gaussian random variable \mathbf{w} . To see the last equality, note that $(\tilde{\boldsymbol{\mu}} + \tilde{\boldsymbol{\Sigma}}^{\frac{1}{2}} \mathbf{w}) \sim \mathcal{N}(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}})$ is a Gaussian random variable with mean $\tilde{\boldsymbol{\mu}}$ and covariance matrix $\tilde{\boldsymbol{\Sigma}}$. Jensen's inequality yields:

$$\left(\mathbb{E}_{\mathbf{w}}[\|\tilde{\boldsymbol{\mu}} + \tilde{\boldsymbol{\Sigma}}^{\frac{1}{2}} \mathbf{w}\|] \right)^2 \leq \mathbb{E}_{\mathbf{w}}[\|\tilde{\boldsymbol{\mu}} + \tilde{\boldsymbol{\Sigma}}^{\frac{1}{2}} \mathbf{w}\|^2] \quad (\text{A.8})$$

$$= \|\tilde{\boldsymbol{\mu}}\|^2 + 2\tilde{\boldsymbol{\mu}}^T \tilde{\boldsymbol{\Sigma}}^{\frac{1}{2}} \mathbb{E}_{\mathbf{w}}[\mathbf{w}] + \mathbb{E}_{\mathbf{w}}[\text{trace}(\mathbf{w}^T \tilde{\boldsymbol{\Sigma}}^{\frac{1}{2}} \tilde{\boldsymbol{\Sigma}}^{\frac{1}{2}} \mathbf{w})] \quad (\text{A.9})$$

$$= \|\tilde{\boldsymbol{\mu}}\|^2 + \text{trace}(\tilde{\boldsymbol{\Sigma}} \mathbb{E}_{\mathbf{w}}[\mathbf{w}\mathbf{w}^T]) = \|\tilde{\boldsymbol{\mu}}\|^2 + \text{trace}(\tilde{\boldsymbol{\Sigma}}) \quad (\text{A.10})$$

$$= \|\mathbf{A}^{-1} \mathbf{y}\|^2 + \text{trace}(\mathbf{A}^{-1} \boldsymbol{\Sigma} \mathbf{A}^{-T}). \quad (\text{A.11})$$

Lastly, we use the square root of this bound in (A.7). \square

Proof of Theorem 2. First, we show that the uniform convergence of $(f_{\mathbf{h}}^{(K)})_{K=1}^{\infty}$ to $f_{\mathbf{h}}$ implies the uniform convergence of $(f_{\mathbf{y}}^{(K)})_{K=1}^{\infty}$ to $f_{\mathbf{y}}$. The PDF of $\mathbf{x} = \mathbf{A}\mathbf{h}$ is

$$f_{\mathbf{x}}(\mathbf{x}) = \frac{1}{|\det(\mathbf{A})|} f_{\mathbf{h}}(\mathbf{A}^{-1} \mathbf{x}) \quad (\text{A.12})$$

because \mathbf{A} is invertible. Analogously, we obtain

$$f_{\mathbf{x}}^{(K)}(\mathbf{x}) = \frac{1}{|\det(\mathbf{A})|} f_{\mathbf{h}}^{(K)}(\mathbf{A}^{-1} \mathbf{x}) \quad (\text{A.13})$$

as the PDF of $\mathbf{A}\mathbf{h}^{(K)}$. Since the random variable $\mathbf{y} = \mathbf{x} + \mathbf{n}$ is a sum of two stochastically independent random variables, its PDF can be computed via convolution:

$$f_{\mathbf{y}}(\mathbf{y}) = \int f_{\mathbf{n}}(s) f_{\mathbf{x}}(\mathbf{y} - s) ds. \quad (\text{A.14})$$

Similarly, $f_{\mathbf{y}}^{(K)}$ is obtained by replacing $f_{\mathbf{x}}$ with $f_{\mathbf{x}}^{(K)}$ in (A.14). For later reference, note, because $f_{\mathbf{n}}$ is positive ($f_{\mathbf{n}}(s) > 0$ for all $s \in \mathbb{R}^N$) and $f_{\mathbf{x}}$ as well as $f_{\mathbf{x}}^{(K)}$ are continuous PDFs, the convolution results $f_{\mathbf{y}}^{(K)}$ and $f_{\mathbf{y}}$ are positive, too.

We have

$$|f_{\mathbf{y}}(\mathbf{y}) - f_{\mathbf{y}}^{(K)}(\mathbf{y})| = \left| \int f_{\mathbf{n}}(s) \left(f_{\mathbf{x}}(\mathbf{y} - s) - f_{\mathbf{x}}^{(K)}(\mathbf{y} - s) \right) ds \right| \quad (\text{A.15})$$

$$\leq \int \left| f_{\mathbf{n}}(s) \frac{f_{\mathbf{h}}(\mathbf{A}^{-1}(\mathbf{y} - s)) - f_{\mathbf{h}}^{(K)}(\mathbf{A}^{-1}(\mathbf{y} - s))}{|\det(\mathbf{A})|} \right| ds \quad (\text{A.16})$$

$$\leq \frac{\|f_{\mathbf{h}} - f_{\mathbf{h}}^{(K)}\|_{\infty}}{|\det(\mathbf{A})|} \int |f_{\mathbf{n}}(s)| ds = \frac{\|f_{\mathbf{h}} - f_{\mathbf{h}}^{(K)}\|_{\infty}}{|\det(\mathbf{A})|}. \quad (\text{A.17})$$

The last integral is equal to one because $f_{\mathbf{n}}$ is a PDF. Since $\lim_{K \rightarrow \infty} \|f_{\mathbf{h}} - f_{\mathbf{h}}^{(K)}\|_{\infty} = 0$ holds by the uniform convergence assumption, we have

$$\lim_{K \rightarrow \infty} \|f_{\mathbf{y}} - f_{\mathbf{y}}^{(K)}\|_{\infty} = 0 \quad (\text{A.18})$$

which means that the sequence $(f_{\mathbf{y}}^{(K)})_{K=1}^{\infty}$ converges uniformly to $f_{\mathbf{y}}$.

Let $r > 0$ be arbitrary. In the following, we show that (3.11) holds uniformly for all $\mathbf{y} \in \mathcal{B}_r$. But because r is arbitrary, the pointwise convergence in (3.11) follows immediately.

Let $\mathbf{y} \in \mathcal{B}_r$ be arbitrary. With (3.10) and (3.9) in mind, we find the following upper bound:

$$\|\hat{\mathbf{h}}(\mathbf{y}) - \hat{\mathbf{h}}^{(K)}(\mathbf{y})\| \leq \int \|\mathbf{h}\| f_{\mathbf{n}}(\mathbf{y} - \mathbf{A}\mathbf{h}) \left| \frac{f_{\mathbf{h}}(\mathbf{h})}{f_{\mathbf{y}}(\mathbf{y})} - \frac{f_{\mathbf{h}}^{(K)}(\mathbf{h})}{f_{\mathbf{y}}^{(K)}(\mathbf{y})} \right| d\mathbf{h} \quad (\text{A.19})$$

$$\leq \sup_{\mathbf{h} \in \mathbb{R}^N} \left| \frac{f_{\mathbf{h}}(\mathbf{h})}{f_{\mathbf{y}}(\mathbf{y})} - \frac{f_{\mathbf{h}}^{(K)}(\mathbf{h})}{f_{\mathbf{y}}^{(K)}(\mathbf{y})} \right| \int \|\mathbf{h}\| f_{\mathbf{n}}(\mathbf{y} - \mathbf{A}\mathbf{h}) d\mathbf{h}. \quad (\text{A.20})$$

The last integral is independent of K and by Lemma 1, it is finite for any $\mathbf{y} \in \mathbb{R}^N$. It is, in particular, bounded by

$$\sqrt{\det(\mathbf{A}^{-1}\mathbf{A}^{-\text{T}})} \sqrt{\|\mathbf{A}^{-1}\|^2 r^2 + \text{trace}(\mathbf{A}^{-1}\boldsymbol{\Sigma}\mathbf{A}^{-\text{T}})} < \infty$$

for all $\mathbf{y} \in \mathcal{B}_r$, where $\|\mathbf{A}^{-1}\|$ is the spectral norm of the matrix \mathbf{A}^{-1} . Hence, as soon as

$$\lim_{K \rightarrow \infty} \sup_{\mathbf{h} \in \mathbb{R}^N} \left| \frac{f_{\mathbf{h}}(\mathbf{h})}{f_{\mathbf{y}}(\mathbf{y})} - \frac{f_{\mathbf{h}}^{(K)}(\mathbf{h})}{f_{\mathbf{y}}^{(K)}(\mathbf{y})} \right| = 0, \quad \forall \mathbf{y} \in \mathcal{B}_r \quad (\text{A.21})$$

is shown, (3.11) is confirmed for all $\mathbf{y} \in \mathcal{B}_r$.

To prove (A.21), we write

$$\left| \frac{f_{\mathbf{h}}(\mathbf{h})}{f_{\mathbf{y}}(\mathbf{y})} - \frac{f_{\mathbf{h}}^{(K)}(\mathbf{h})}{f_{\mathbf{y}}^{(K)}(\mathbf{y})} \right| = \left| \frac{f_{\mathbf{h}}(\mathbf{h})f_{\mathbf{y}}^{(K)}(\mathbf{y}) - f_{\mathbf{y}}(\mathbf{y})f_{\mathbf{h}}^{(K)}(\mathbf{h})}{f_{\mathbf{y}}(\mathbf{y})f_{\mathbf{y}}^{(K)}(\mathbf{y})} \right| \quad (\text{A.22})$$

for an arbitrary $\mathbf{h} \in \mathbb{R}^N$. Now, we add $0 = f_{\mathbf{y}}^{(K)}(\mathbf{y})f_{\mathbf{h}}^{(K)}(\mathbf{h}) - f_{\mathbf{y}}^{(K)}(\mathbf{y})f_{\mathbf{h}}^{(K)}(\mathbf{h})$ in

the numerator on the right-hand side and apply the triangle inequality to get

$$\left| \frac{f_{\mathbf{h}}(\mathbf{h})f_{\mathbf{y}}^{(K)}(\mathbf{y}) - f_{\mathbf{y}}(\mathbf{y})f_{\mathbf{h}}^{(K)}(\mathbf{h})}{f_{\mathbf{y}}(\mathbf{y})f_{\mathbf{y}}^{(K)}(\mathbf{y})} \right| \quad (\text{A.23})$$

$$\begin{aligned} &\leq \frac{\left| (f_{\mathbf{h}}(\mathbf{h}) - f_{\mathbf{h}}^{(K)}(\mathbf{h})) f_{\mathbf{y}}^{(K)}(\mathbf{y}) \right|}{f_{\mathbf{y}}(\mathbf{y})f_{\mathbf{y}}^{(K)}(\mathbf{y})} + \frac{\left| (f_{\mathbf{y}}^{(K)}(\mathbf{y}) - f_{\mathbf{y}}(\mathbf{y})) f_{\mathbf{h}}^{(K)}(\mathbf{h}) \right|}{f_{\mathbf{y}}(\mathbf{y})f_{\mathbf{y}}^{(K)}(\mathbf{y})} \\ &\leq \frac{\|f_{\mathbf{h}} - f_{\mathbf{h}}^{(K)}\|_{\infty} \|f_{\mathbf{y}}^{(K)}\|_{\infty} + \|f_{\mathbf{y}}^{(K)} - f_{\mathbf{y}}\|_{\infty} \|f_{\mathbf{h}}^{(K)}\|_{\infty}}{f_{\mathbf{y}}(\mathbf{y})f_{\mathbf{y}}^{(K)}(\mathbf{y})}. \end{aligned} \quad (\text{A.24})$$

The goal is now to find upper bounds for all terms in the numerator and lower bounds for all terms in the denominator which in the end allow us to see that (A.22) converges to zero as K approaches infinity.

By the compactness of \mathcal{B}_r and the continuity of $f_{\mathbf{y}}$, there exists a $\mathbf{y}_{\min} \in \mathcal{B}_r$ at which $f_{\mathbf{y}}$ attains a minimum value $f_{\mathbf{y}}(\mathbf{y}_{\min}) > 0$ over \mathcal{B}_r . Due to the uniform convergence (A.18), there exists an index $N_1 \in \mathbb{N}$ such that $|f_{\mathbf{y}}(\mathbf{y}) - f_{\mathbf{y}}^{(K)}(\mathbf{y})| \leq \frac{1}{2}f_{\mathbf{y}}(\mathbf{y}_{\min})$ holds for all $K \geq N_1$ and for all $\mathbf{y} \in \mathcal{B}_r$. The reverse triangle inequality then shows that $f_{\mathbf{y}}^{(K)}(\mathbf{y}) \geq f_{\mathbf{y}}(\mathbf{y}) - |f_{\mathbf{y}}(\mathbf{y}) - f_{\mathbf{y}}^{(K)}(\mathbf{y})| \geq f_{\mathbf{y}}(\mathbf{y}_{\min}) - \frac{1}{2}f_{\mathbf{y}}(\mathbf{y}_{\min}) = \frac{1}{2}f_{\mathbf{y}}(\mathbf{y}_{\min})$ is true. Hence, with $M_1 = \frac{1}{2}f_{\mathbf{y}}(\mathbf{y}_{\min}) > 0$, the inequality

$$f_{\mathbf{y}}^{(K)}(\mathbf{y}) \geq M_1 \quad (\text{A.25})$$

holds for all $\mathbf{y} \in \mathcal{B}_r$ and for all $K \geq N_1$. Further, since $\|f_{\mathbf{h}} - f_{\mathbf{h}}^{(K)}\|_{\infty} \rightarrow 0$ and $\|f_{\mathbf{h}}\|_{\infty} < \infty$, there exist an $M_2 > 0$ and an $N_2 \in \mathbb{N}$ such that

$$\|f_{\mathbf{h}}^{(K)}\|_{\infty} \leq M_2 \quad \text{for all } K \geq N_2. \quad (\text{A.26})$$

Analogously, there exist an $M_3 > 0$ and an $N_3 \in \mathbb{N}$ such that

$$\|f_{\mathbf{y}}^{(K)}\|_{\infty} \leq M_3 \quad \text{for all } K \geq N_3. \quad (\text{A.27})$$

Let $\varepsilon > 0$ be arbitrary. Due to $\|f_{\mathbf{h}} - f_{\mathbf{h}}^{(K)}\|_{\infty} \rightarrow 0$, there exists an index $N_4 \geq \max\{N_1, N_3\}$ such that

$$\|f_{\mathbf{h}} - f_{\mathbf{h}}^{(K)}\|_{\infty} \leq \frac{f_{\mathbf{y}}(\mathbf{y}_{\min})M_1}{2M_3}\varepsilon \quad \text{for all } K \geq N_4. \quad (\text{A.28})$$

Similarly, there exists an index $N_5 \geq \max\{N_1, N_2\}$ with

$$\|f_{\mathbf{y}} - f_{\mathbf{y}}^{(K)}\|_{\infty} \leq \frac{f_{\mathbf{y}}(\mathbf{y}_{\min})M_1}{2M_2}\varepsilon \quad \text{for all } K \geq N_5. \quad (\text{A.29})$$

We can use the last five inequalities to bound (A.24). To this end, $f_{\mathbf{y}}(\mathbf{y}_{\min})$ and (A.25) provide lower bounds on the terms in the denominator, (A.28) and (A.27) bound the first summand in the numerator, and (A.29) and (A.26) bound the second summand in the numerator. In total, this yields an upper bound on (A.22):

$$\left| \frac{f_{\mathbf{h}}(\mathbf{h})}{f_{\mathbf{y}}(\mathbf{y})} - \frac{f_{\mathbf{h}}^{(K)}(\mathbf{h})}{f_{\mathbf{y}}^{(K)}(\mathbf{y})} \right| \leq \frac{f_{\mathbf{y}}(\mathbf{y}_{\min})M_1}{2M_3}\varepsilon \cdot \frac{M_3}{f_{\mathbf{y}}(\mathbf{y}_{\min})M_1} + \frac{f_{\mathbf{y}}(\mathbf{y}_{\min})M_1}{2M_2}\varepsilon \cdot \frac{M_2}{f_{\mathbf{y}}(\mathbf{y}_{\min})M_1} \quad (\text{A.30})$$

for all $K \geq \max\{N_4, N_5\}$ and for all $\mathbf{y} \in \mathcal{B}_r$. We conclude

$$\sup_{\mathbf{h} \in \mathbb{R}^N} \left| \frac{f_{\mathbf{h}}(\mathbf{h})}{f_{\mathbf{y}}(\mathbf{y})} - \frac{f_{\mathbf{h}}^{(K)}(\mathbf{h})}{f_{\mathbf{y}}^{(K)}(\mathbf{y})} \right| \leq \varepsilon \quad (\text{A.31})$$

for all $K \geq \max\{N_4, N_5\}$ and all $\mathbf{y} \in \mathcal{B}_r$, and because ε was arbitrary, (A.21) is confirmed, which concludes the proof. \square

Discussion for Noninvertible Matrices B

This appendix provides more details for the discussion in Section 3.2.1. As mentioned there, two problems arise in the proof of Theorem 2 if \mathbf{A} is not invertible. First, Lemma 1 (see Appendix A) might no longer hold. That is, the integral $\int \|\mathbf{h}\| f_{\mathbf{n}}(\mathbf{y} - \mathbf{A}\mathbf{h}) d\mathbf{h}$ might not be finite for all $\mathbf{y} \in \mathbb{R}^m$. This is discussed in Appendix B.1. Second, a crucial step in the proof of Theorem 2 is showing that $(f_{\mathbf{y}}^{(K)})_{K=1}^{\infty}$ converges uniformly to $f_{\mathbf{y}}$, cf. (A.18). This is possible because the sequence $(f_{\mathbf{x}}^{(K)})_{K=1}^{\infty}$ of probability density functions (PDFs) corresponding to the auxiliary variables $\mathbf{x}^{(K)} = \mathbf{A}\mathbf{h}^{(K)}$ converges uniformly to the PDF $f_{\mathbf{x}}$ of $\mathbf{x} = \mathbf{A}\mathbf{h}$ if \mathbf{A} is invertible. Unfortunately, if \mathbf{A} is not invertible, this will likely not hold for an arbitrary $f_{\mathbf{h}}$. This is discussed in Appendix B.2.

B.1 Integral in Lemma 1

To see why the integral $\int \|\mathbf{h}\| f_{\mathbf{n}}(\mathbf{y} - \mathbf{A}\mathbf{h}) d\mathbf{h}$ in Lemma 1 might not be finite if \mathbf{A} is not invertible, consider the matrix $\mathbf{A} = [\mathbf{I}, \mathbf{0}] \in \mathbb{R}^{m \times N}$ with $m < N$ where $\mathbf{I} \in \mathbb{R}^{m \times m}$ is the identity matrix and the remaining elements of \mathbf{A} are zero. Let us write $\mathbf{h} = [\mathbf{h}_m^{\top}, \mathbf{h}_{N-m}^{\top}]^{\top} \in \mathbb{R}^m \times \mathbb{R}^{N-m}$. Define the set

$$\mathcal{D} = \{\mathbf{h} \in \mathbb{R}^N : \|\mathbf{h}_m\| \geq 1, \|\mathbf{h}_{N-m}\| \geq 1\} \quad (\text{B.1})$$

where the norm of both subvectors \mathbf{h}_m and \mathbf{h}_{N-m} is at least one such that we always have $\|\mathbf{h}\| \geq 1$ on \mathcal{D} . We can now compute:

$$\int_{\mathbb{R}^N} \|\mathbf{h}\| f_{\mathbf{n}}(\mathbf{y} - \mathbf{A}\mathbf{h}) d\mathbf{h} \geq \int_{\mathcal{D}} f_{\mathbf{n}}(\mathbf{y} - \mathbf{A}\mathbf{h}) d\mathbf{h} \quad (\text{B.2})$$

$$= \int_{\|\mathbf{h}_{N-m}\| \geq 1} \int_{\|\mathbf{h}_m\| \geq 1} f_{\mathbf{n}}(\mathbf{y} - \mathbf{h}_m) d\mathbf{h}_m d\mathbf{h}_{N-m}. \quad (\text{B.3})$$

Since $f_{\mathbf{n}}$ is an m -dimensional Gaussian PDF, the inner integral is equal to some constant c with $0 < c < 1$ and it follows that $\int \|\mathbf{h}\| f_{\mathbf{n}}(\mathbf{y} - \mathbf{A}\mathbf{h}) d\mathbf{h}$ is not finite.

B.2 Uniform Convergence

We want to find an expression for the PDF of $\mathbf{x} = \mathbf{A}\mathbf{h}$ when $\mathbf{A} \in \mathbb{R}^{m \times N}$ is a wide matrix with full rank $m < N$. We can assume that the first m columns are linearly independent (otherwise we introduce a permutation matrix for the following argument). This allows us to partition $\mathbf{A} = [\mathbf{A}_i, \mathbf{A}_n]$ into an invertible part $\mathbf{A}_i \in \mathbb{R}^{m \times m}$ and a noninvertible part $\mathbf{A}_n \in \mathbb{R}^{m \times N-m}$. With a corresponding partitioning of $\mathbf{h} = [\mathbf{h}_i^T, \mathbf{h}_n^T]^T \in \mathbb{R}^m \times \mathbb{R}^{N-m}$, we can write

$$\mathbf{x} = \mathbf{A}\mathbf{h} = \mathbf{A}_i\mathbf{h}_i + \mathbf{A}_n\mathbf{h}_n. \quad (\text{B.4})$$

Defining a transformation

$$t : \mathbb{R}^m \times \mathbb{R}^{N-m} \rightarrow \mathbb{R}^m \times \mathbb{R}^{N-m}, (\mathbf{h}_i, \mathbf{h}_n) \mapsto (\mathbf{A}_i\mathbf{h}_i + \mathbf{A}_n\mathbf{h}_n, \mathbf{h}_n) \quad (\text{B.5})$$

with the inverse

$$t^{-1} : (\mathbf{x}, \mathbf{x}') \mapsto (\mathbf{A}_i^{-1}(\mathbf{x} - \mathbf{A}_n\mathbf{x}'), \mathbf{x}') \quad (\text{B.6})$$

we can now compute the joint density $f_{\mathbf{x}, \mathbf{x}'}$ via the usual transformation formula:

$$(\mathbf{x}, \mathbf{x}') \mapsto f_{\mathbf{x}, \mathbf{x}'}(\mathbf{x}, \mathbf{x}') = \frac{f_{\mathbf{h}}(t^{-1}(\mathbf{x}, \mathbf{x}'))}{\left| \det \left(\frac{\partial t}{\partial \mathbf{h}}(t^{-1}(\mathbf{x}, \mathbf{x}')) \right) \right|}. \quad (\text{B.7})$$

Together with $\left| \det \left(\frac{\partial t}{\partial \mathbf{h}}(t^{-1}(\mathbf{x}, \mathbf{x}')) \right) \right| = |\det(\mathbf{A}_i)|$, we can express the PDF of $\mathbf{x} = \mathbf{A}\mathbf{h}$ via marginalization:

$$\mathbf{x} \mapsto f_{\mathbf{x}}(\mathbf{x}) = \int_{\mathbb{R}^{N-m}} f_{\mathbf{x}, \mathbf{x}'}(\mathbf{x}, \mathbf{x}') d\mathbf{x}' = \int_{\mathbb{R}^{N-m}} \frac{f_{\mathbf{h}}(t^{-1}(\mathbf{x}, \mathbf{x}'))}{|\det(\mathbf{A}_i)|} d\mathbf{x}'. \quad (\text{B.8})$$

Analogously, we obtain the PDF $f_{\mathbf{x}}^{(K)}$ of $\mathbf{x}^{(K)} = \mathbf{A}\mathbf{h}^{(K)}$. Given

$$\left| f_{\mathbf{x}}^{(K)}(\mathbf{x}) - f_{\mathbf{x}}(\mathbf{x}) \right| = \frac{\left| \int_{\mathbb{R}^{N-m}} \left(f_{\mathbf{h}}^{(K)}(t^{-1}(\mathbf{x}, \mathbf{x}')) - f_{\mathbf{h}}(t^{-1}(\mathbf{x}, \mathbf{x}')) \right) d\mathbf{x}' \right|}{|\det(\mathbf{A}_i)|} \quad (\text{B.9})$$

for $\mathbf{x} \in \mathbb{R}^m$, we can conjecture that due to the integral over \mathbb{R}^{N-m} the uniform convergence of $(f_{\mathbf{h}}^{(K)})_{K=1}^{\infty}$ to $f_{\mathbf{h}}$ alone is generally not sufficient to infer the uniform convergence of $(f_{\mathbf{x}}^{(K)})_{K=1}^{\infty}$ to $f_{\mathbf{x}}$.

Bibliography

- [1] T. T. Nguyen, H. D. Nguyen, F. Chamroukhi, and G. J. McLachlan, “Approximation by finite mixtures of continuous density functions that vanish at infinity,” *Cogent Math. Statist.*, vol. 7, no. 1, p. 1750861, 2020.
- [2] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola, “A kernel two-sample test,” *J. Mach. Learn. Res.*, vol. 13, no. 1, pp. 723–773, Mar. 2012.
- [3] E. Björnson, J. Hoydis, and L. Sanguinetti, *Massive MIMO Networks: Spectral, Energy, and Hardware Efficiency*. Now Publishers Inc., Nov. 2017, vol. 11, no. 3–4.
- [4] A. Goldsmith, *Wireless Communications*. Cambridge University Press, 2005.
- [5] O. E. Ayach, S. Rajagopal, S. Abu-Surra, Z. Pi, and R. W. Heath, “Spatially sparse precoding in millimeter wave MIMO systems,” *IEEE Trans. Wireless Commun.*, vol. 13, no. 3, pp. 1499–1513, Mar. 2014.
- [6] L. Liang, W. Xu, and X. Dong, “Low-complexity hybrid precoding in massive multiuser MIMO systems,” *IEEE Wireless Commun. Lett.*, vol. 3, no. 6, pp. 653–656, Dec. 2014.
- [7] A. Alkhateeb, O. El Ayach, G. Leus, and R. W. Heath, “Channel estimation and hybrid precoding for millimeter wave cellular systems,” *IEEE J. Sel. Topics Signal Process.*, vol. 8, no. 5, pp. 831–846, Oct. 2014.
- [8] R. Méndez-Rial, C. Rusu, N. González-Prelcic, A. Alkhateeb, and R. W. Heath, “Hybrid MIMO architectures for millimeter wave communications: Phase shifters or switches?” *IEEE Access*, vol. 4, pp. 247–267, 2016.
- [9] K. Ardah, G. Fodor, Y. C. B. Silva, W. C. Freitas, and F. R. P. Cavalcanti, “A unifying design of hybrid beamforming architectures employing phase shifters or switches,” *IEEE Trans. Veh. Technol.*, vol. 67, no. 11, pp. 11 243–11 247, Nov. 2018.

Bibliography

- [10] S. Foucart and H. Rauhut, *A Mathematical Introduction to Compressive Sensing*. Birkhäuser Basel, 2013.
- [11] Y. C. Eldar and M. Mishali, “Robust recovery of signals from a structured union of subspaces,” *IEEE Trans. Inf. Theory*, vol. 55, no. 11, pp. 5302–5316, Nov. 2009.
- [12] T. Blumensath, “Sampling and reconstructing signals from a union of linear subspaces,” *IEEE Trans. Inf. Theory*, vol. 57, no. 7, pp. 4660–4671, Jul. 2011.
- [13] T. Wiese, L. Weiland, and W. Utschick, “Inexact projected gradients on unions of subspaces,” in *Proc. IEEE Int. Symp. Inf. Theory (ISIT)*, 2017, pp. 236–240.
- [14] G. Davis, S. Mallat, and M. Avellaneda, “Adaptive greedy approximations,” *Constr. Approx.*, vol. 13, no. 1, pp. 57–98, Mar. 1997.
- [15] Y. C. Pati, R. Rezaifar, and P. S. Krishnaprasad, “Orthogonal matching pursuit: Recursive function approximation with applications to wavelet decomposition,” in *Proc. Asilomar Conf. on Signals, Syst., and Comput.*, Nov. 1993, pp. 40–44.
- [16] J. A. Tropp, “Greed is good: Algorithmic results for sparse approximation,” *IEEE Trans. Inf. Theory*, vol. 50, no. 10, pp. 2231–2242, Oct. 2004.
- [17] C. M. Bishop, *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Springer-Verlag, 2006.
- [18] C. F. J. Wu, “On the convergence properties of the EM algorithm,” *The Ann. of Statist.*, vol. 11, no. 1, pp. 95–103, 1983.
- [19] B. K. Sriperumbudur, A. Gretton, K. Fukumizu, B. Schölkopf, and G. R. Lanckriet, “Hilbert space embeddings and metrics on probability measures,” *J. Mach. Learn. Res.*, vol. 11, pp. 1517–1561, Aug. 2010.
- [20] Y. Li, K. Swersky, and R. Zemel, “Generative moment matching networks,” in *Proc. Int. Conf. on Mach. Learn. (ICML)*, 2015, pp. 1718–1727.
- [21] G. K. Dziugaite, D. M. Roy, and Z. Ghahramani, “Training generative neural networks via maximum mean discrepancy optimization,” in *Proc. Conf. on Uncertainty in Artif. Intell. (UAI)*, 2015, pp. 258–267.
- [22] Y. Ren, J. Li, Y. Luo, and J. Zhu, “Conditional generative moment-matching networks,” in *Proc. Adv. in Neural Inf. Process. Syst. (NeurIPS)*, 2016, pp. 2928–2936.

- [23] H. Gao and H. Huang, “Joint generative moment-matching network for learning structural latent code,” in *Proc. Int. Joint Conf. on Artif. Intell. (IJCAI)*, Jul. 2018, pp. 2121–2127.
- [24] J. Bergstra and Y. Bengio, “Random search for hyper-parameter optimization,” *J. Mach. Learn. Res.*, vol. 13, pp. 281–305, Feb. 2012.
- [25] M. Koller, B. Fesl, N. Turan, and W. Utschick, “An asymptotically optimal approximation of the conditional mean channel estimator based on Gaussian mixture models,” in *Proc. IEEE Int. Conf. Acoust., Speech, and Signal Process. (ICASSP)*, 2022, pp. 5268–5272.
- [26] —, “An asymptotically MSE-optimal estimator based on Gaussian mixture models,” *IEEE Trans. Signal Process.*, vol. 70, pp. 4109–4123, 2022.
- [27] G. Schay, *Introduction to Probability with Statistical Applications*. Birkhäuser Boston, Jun. 2016.
- [28] S. Resnick, *A Probability Path*. Birkhäuser Boston, 2005.
- [29] A. Klenke, *Probability Theory: A Comprehensive Course*. Springer, 2008.
- [30] E. M. Goggin, “Convergence in distribution of conditional expectations,” *The Ann. of Probability*, vol. 22, no. 2, pp. 1097–1114, 1994.
- [31] Y. Gu and Y. D. Zhang, “Information-theoretic pilot design for downlink channel estimation in FDD massive MIMO systems,” *IEEE Trans. Signal Process.*, vol. 67, no. 9, pp. 2334–2346, May 2019.
- [32] R. M. Gray, “Toeplitz and circulant matrices: A review,” *Found. and Trends in Commun. and Inf. Theory*, no. 3, pp. 155–239, 2006.
- [33] J. Kermoal, L. Schumacher, K. Pedersen, P. Mogensen, and F. Frederiksen, “A stochastic MIMO radio channel model with experimental validation,” *IEEE J. Sel. Areas Commun.*, vol. 20, no. 6, pp. 1211–1226, 2002.
- [34] M. Šimko, C. Mehlführer, M. Wrulich, and M. Rupp, “Doubly dispersive channel estimation with scalable complexity,” in *Proc. Int. ITG Workshop on Smart Antennas (WSA)*, Feb. 2010, pp. 251–256.
- [35] 3GPP, “Spatial channel model for multiple input multiple output (MIMO) simulations,” 3rd Generation Partnership Project (3GPP), Tech. Rep. 25.996 (V16.0.0), Jul. 2020.

Bibliography

- [36] D. Neumann, T. Wiese, and W. Utschick, "Learning the MMSE channel estimator," *IEEE Trans. Signal Process.*, vol. 66, no. 11, pp. 2905–2917, Jun. 2018.
- [37] S. Jaeckel, L. Raschkowski, K. Börner, and L. Thiele, "QuaDRiGa: A 3-D multi-cell channel model with time evolution for enabling virtual field trials," *IEEE Trans. Antennas Propag.*, vol. 62, no. 6, pp. 3242–3256, 2014.
- [38] S. Jaeckel, L. Raschkowski, K. Börner, L. Thiele, F. Burkhardt, and E. Eberlein, "QuaDRiGa: Quasi deterministic radio channel generator, user manual and documentation," Fraunhofer Heinrich Hertz Institute, Tech. Rep., v2.2.0, 2019.
- [39] 3GPP, "NR; physical channels and modulation," 3rd Generation Partnership Project (3GPP), Tech. Spec. 38.211 (V16.7.0), Sep. 2021.
- [40] M. Kurras, S. Dai, S. Jaeckel, and L. Thiele, "Evaluation of the spatial consistency feature in the 3GPP geometry-based stochastic channel model," in *Proc. IEEE Wireless Commun. and Netw. Conf. (WCNC)*, Apr. 2019, pp. 1–6.
- [41] B. Fesl, N. Turan, M. Koller, and W. Utschick, "A low-complexity MIMO channel estimator with implicit structure of a convolutional neural network," in *Proc. IEEE Int. Workshop on Signal Process. Adv. in Wireless Commun. (SPAWC)*, 2021.
- [42] S. A. Busari, K. M. S. Huq, S. Mumtaz, L. Dai, and J. Rodriguez, "Millimeter-wave massive MIMO communication for future wireless systems: A survey," *IEEE Commun. Surveys Tuts.*, vol. 20, no. 2, pp. 836–869, Secondquarter 2018.
- [43] K. Hassan, M. Masarra, M. Zwingelstein, and I. Dayoub, "Channel estimation techniques for millimeter-wave communication systems: Achievements and challenges," *IEEE Open J. Commun. Soc.*, vol. 1, pp. 1336–1363, 2020.
- [44] A. Alkhateeb, G. Leus, and R. W. Heath, "Compressed sensing based multi-user millimeter wave systems: How many measurements are needed?" in *Proc. IEEE Int. Conf. Acoust., Speech, and Signal Process. (ICASSP)*, 2015, pp. 2909–2913.
- [45] D. L. Donoho, A. Maleki, and A. Montanari, "Message passing algorithms for compressed sensing: I. Motivation and construction," in *Proc. IEEE Inf. Theory Workshop on Inf. Theory (ITW)*, Jan. 2010, pp. 1–5.
- [46] A. Maleki, L. Anitori, Z. Yang, and R. G. Baraniuk, "Asymptotic analysis of complex LASSO via complex approximate message passing (CAMP)," *IEEE Trans. Inf. Theory*, vol. 59, no. 7, pp. 4290–4308, Jul. 2013.

- [47] M. Soltani, V. Pourahmadi, and H. Sheikhzadeh, "Pilot pattern design for deep learning-based channel estimation in OFDM systems," *IEEE Wireless Commun. Lett.*, vol. 9, no. 12, pp. 2173–2176, Dec. 2020.
- [48] M. Soltani, V. Pourahmadi, A. Mirzaei, and H. Sheikhzadeh, "Deep learning-based channel estimation," *IEEE Commun. Lett.*, vol. 23, no. 4, pp. 652–655, Apr. 2019.
- [49] S. Coleri, M. Ergen, A. Puri, and A. Bahai, "Channel estimation techniques based on pilot arrangement in OFDM systems," *IEEE Trans. Broadcast.*, vol. 48, no. 3, pp. 223–229, 2002.
- [50] L. Tong, B. Sadler, and M. Dong, "Pilot-assisted wireless transmissions: General model, design criteria, and signal processing," *IEEE Signal Process. Mag.*, vol. 21, no. 6, pp. 12–25, Nov. 2004.
- [51] M. Koller and W. Utschick, "Learning a compressive sensing matrix with structural constraints via maximum mean discrepancy optimization," *Signal Process.*, p. 108553, 2022.
- [52] A. Griewank and A. Walther, *Evaluating Derivatives: Principles and Techniques of Algorithmic Differentiation*. Society for Industrial and Applied Mathematics, 2008.
- [53] K. Ardah and M. Haardt, "Compressed sensing constant modulus constrained projection matrix design and high-resolution DoA estimation methods," in *Proc. Int. ITG Workshop on Smart Antennas (WSA)*, Nov. 2021, pp. 1–5.
- [54] L. Baldassarre, Y. Li, J. Scarlett, B. Gözcü, I. Bogunovic, and V. Cevher, "Learning-based compressive subsampling," *IEEE J. Sel. Topics Signal Process.*, vol. 10, no. 4, pp. 809–822, Jun. 2016.
- [55] Y. Weiss, H. S. Chang, and W. T. Freeman, "Learning compressed sensing," in *Proc. Allerton Conf. Commun., Control, and Comput.*, 2007.
- [56] M. Elad, "Optimized projections for compressed sensing," *IEEE Trans. Signal Process.*, vol. 55, no. 12, pp. 5695–5702, Dec. 2007.
- [57] J. M. Duarte-Carvajalino and G. Sapiro, "Learning to sense sparse signals: Simultaneous sensing matrix and sparsifying dictionary optimization," *IEEE Trans. Image Process.*, vol. 18, no. 7, pp. 1395–1408, Jul. 2009.

Bibliography

- [58] A. Sadeghian, B. Bah, and V. Cevher, “Energy-aware adaptive bi-Lipschitz embeddings,” in *Proc. Int. Conf. Sampling Theory and Appl. (SampTA)*, Jul. 2013.
- [59] C. Hegde, A. C. Sankaranarayanan, W. Yin, and R. G. Baraniuk, “NuMax: A convex approach for learning near-isometric linear embeddings,” *IEEE Trans. Signal Process.*, vol. 63, no. 22, pp. 6109–6121, Nov. 2015.
- [60] S. Wu, A. Dimakis, S. Sanghavi, F. Yu, D. Holtmann-Rice, D. Storchus, A. Rostamizadeh, and S. Kumar, “Learning a compressed sensing measurement matrix via gradient unrolling,” in *Proc. Int. Conf. on Mach. Learn. (ICML)*, vol. 97, Jun. 2019, pp. 6828–6839.
- [61] P. Wu, Z. Liu, and J. Cheng, “Compressed CSI feedback with learned measurement matrix for mmWave massive MIMO,” *arXiv*, 2019. [Online]. Available: <http://arxiv.org/abs/1903.02127>
- [62] P. Wu and J. Cheng, “Acquiring measurement matrices via deep basis pursuit for sparse channel estimation in mmWave massive MIMO systems,” *arXiv*, 2020. [Online]. Available: <http://arxiv.org/abs/2007.05177>
- [63] F. Mezzadri, “How to generate random matrices from the classical compact groups,” *Notices of the Amer. Math. Society*, vol. 54, Oct. 2006.
- [64] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, “Pytorch: An imperative style, high-performance deep learning library,” in *Proc. Adv. in Neural Inf. Process. Syst. (NeurIPS)*, vol. 32, 2019, pp. 8026–8037.
- [65] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” in *Proc. Int. Conf. Learn. Representations (ICLR)*, 2015.
- [66] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016.
- [67] Z. Yang, J. Li, P. Stoica, and L. Xie, “Sparse methods for direction-of-arrival estimation,” in *Academic Press Library in Signal Process.*, 2018, pp. 509–581.
- [68] J. Haupt, W. U. Bajwa, G. Raz, and R. Nowak, “Toeplitz compressed sensing matrices with applications to sparse channel estimation,” *IEEE Trans. Inf. Theory*, vol. 56, no. 11, pp. 5862–5875, Nov. 2010.

Bibliography

- [69] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, “Generative adversarial nets,” in *Proc. Adv. in Neural Inf. Process. Syst. (NeurIPS)*, 2014, pp. 2672–2680.