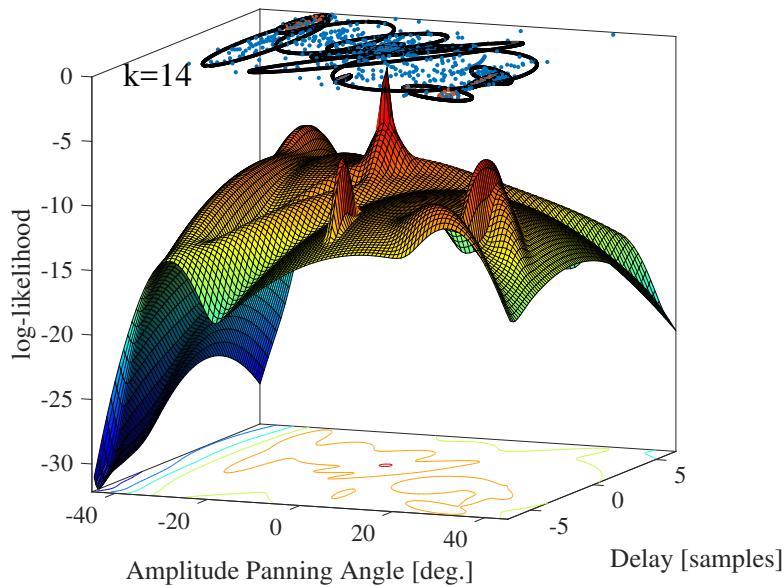

Estimation of Source Parameters and Segmentation of Stereophonic Music Mixtures.

May 2017

Jacob Møller Hjerrild

Project Report



Aalborg University
Sound & Music Computing
Rendsburgsgade 14
DK-9000 Aalborg



Title:

Panning Parameter Estimation
4th Semester M.Sc.in Sound and
Music Computing, Aalborg Uni-
versity.

Theme:

SMC Master Thesis

Project Period:

Spring Semester 2017

Authors:



Jacob Møller Hjerrild

Supervisor:

Mads Græsbøll Christensen

Number of Prints:

none

Number of Pages:

70

Date for Submission:

May, 2017

Abstract:

In this report, we propose a novel source panning estimator for stereophonic mixtures, allowing for panning estimation on multi-channel audio even if the source pitches and harmonic amplitudes are unknown. The presented method does not require prior knowledge of the number of sources present in the mixture. The estimator is formulated using an unsupervised learning framework, allowing for optimal segmentation of the stereophonic signal, based on maximum a posteriori clustering of source parameters.

The content of this report is not public available, and publication may only be pursued due to agreement with the author.

Preface

The semester theme is the master's thesis of the Sound and Music Computing program. The topic must fall within the general area of Sound and Music Computing. The semester is organized as a 30 ECTS-points within Architecture Design & Media Technology at Aalborg University. The target group of this report are future researchers of audio technologies.

Appendices

Appendices are found after the main report in the digital hand-in.

All figures, tables and equations are referred to by the number of the chapter they are used in, followed by a number indicating the number of figure, table or equation in the specific chapter. Hence, each figure has a unique number, which is also printed at the bottom of the figure along with a caption. An example is Figure 2.1, which means the first figure in chapter 2. The same applies to tables and equations, the latter of which have no captions. Appendices are referred to by capital letters instead of chapter numbers.

Bibliography

At the very end of the present documentation, a Bibliography is listed which contains all sources of research used in the study. In the Bibliography books are indicated with author, title, publisher and year. Web pages are indicated with author and title. All information sources are referred to by the number which they feature in the list. This will look like this: [number].

Nomenclature

x Scalar

AIC Akaike's information criterion

BIC Bayes information criterion

BSS Blind Source Separation

EM Expectation-Maximization

GMM Gaussian Mixture Model

MAP Maximum a Posteriori

MDL Minimum Description Length

ML Maximum likelihood

MMDL mixture-MDL

VRC Variance Ratio Criterion

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 1 |
| 1.1 | Related Work | 1 |
| 1.2 | Introduction | 2 |
| 1.2.1 | Motivation | 2 |
| 1.3 | The Panning Parameters | 3 |
| 1.3.1 | Delay Panning | 3 |
| 1.3.2 | Amplitude Panning | 3 |
| 1.3.3 | Stereophonic Amplitude Panning | 4 |
| 1.4 | Signal Model | 7 |
| 1.4.1 | Estimating the Panning Parameters | 7 |
| 1.4.2 | Visualizing the Amplitude Angle as a Panogram | 9 |
| 2 | Clustering | 11 |
| 2.1 | Estimation of Source Parameters | 11 |
| 2.2 | Finite Mixtures | 11 |
| 2.2.1 | Parameterization of the finite mixtures of Gaussians | 12 |
| 2.2.2 | The Gaussian mixture model as a likelihood | 13 |
| 2.2.3 | Fitting the Gaussian mixture as a maximum likelihood solution | 14 |
| 2.3 | Maximizing the likelihood | 16 |
| 2.4 | Initialization using K-Means Clustering | 18 |
| 3 | Model Order Selection | 23 |
| 3.0.1 | The Asymptotic MAP criterion | 23 |
| 3.0.2 | Suitable Prior on the Mixing Probabilities | 26 |
| 3.1 | Mixture MDL | 26 |
| 3.2 | Model Pruning by Component Annihilation | 28 |
| 3.2.1 | Overfitting a Gaussian Mixture Model to Panning Parameter Space | 28 |
| 3.2.2 | Clustering and Component Metrics | 28 |
| 3.2.3 | Annihilation Steps | 29 |
| 3.2.4 | Thresholding on the Cluster Angle and size | 32 |

| | | |
|----------|---|-----------|
| 3.2.5 | Bayesian Interpretation of the Threshold | 33 |
| 4 | Segmentation of the Stereophonic Signal | 37 |
| 4.1 | Signal Segmentation | 37 |
| 5 | Experiments | 39 |
| 5.1 | Experiments | 39 |
| 5.1.1 | Segmentation | 39 |
| 5.1.2 | Source Parameter Estimation | 40 |
| | Bibliography | 43 |
| | Appendix | 47 |
| A | Estimation the Amplitude Panning Angle | 49 |
| A.1 | Initial test of optimal segmentation based on GMM and BIC | 49 |
| B | Initial Project Proposal | 55 |
| B.1 | Abstract | 55 |
| B.2 | Introduction | 55 |
| B.3 | State of the Art | 56 |
| B.4 | Thesis | 56 |
| B.5 | Implementation and Methodology | 57 |
| B.6 | Time plan | 58 |
| B.7 | Notes for article to WASPAA | 58 |
| B.7.1 | Algorithm description | 59 |

Chapter 1

Introduction

1.1 Related Work

***** THIS IS A WORK IN PROGRESS Should be mixed with application*****

Weiss [1] discusses a broad range of pitch estimators in the context of multi pitch estimation, and proposes a novel maximum likelihood multi-pitch estimator, that utilizes the gain parameters; delay and amplitude gain. Weiss assumes that these delay and amplitude parameters are known and the method to estimate these is still an unsolved problem in pitch-estimation. Pitch estimation has applications in problems such as separation [2], enhancement [3], compression [4], transcription [5], classification [6] and source localization [7].

The idea of separating sources from a multi-channel mixture is applied within source separation and array processing [8] research, where audio material is consisting of speakers in different angles and distances to a microphone array setup. The structure of speech can be assumed to be W-disjoint orthogonal [9], however musical mixes builds upon harmonic structures between a variety of sound sources, hence spectral overlap is a common known problem within multi-pitch estimation of musical content. It is interesting that most music is available in stereo only, and this is why the panning parameters are interesting to exploit and implement in the research area of multi-pitch estimation and source separation within music.

The blind source separation of speech signals [10] are based upon the W-disjoint orthogonality and histogramming in the time-frequency domain, based on uniform signal segmentation and manual parameter estimation by visual inspection.

Something about optimal segmentation originally proposed by Prandoni in [11], which is also used for pitch estimation in [12].

In this work we propose an automatic estimator based on clustering by using optimal signal segmentation in time-frequency domain.

1.2 Introduction

The problem of separating sources from a given mixture, is important and often multi-pitch estimation [13], array processing [14, 8], time direction of arrival (TDOA) [7, 2, 15] and blind source separation (BSS) [16] are methods for source separation, depending on the given context. Many methods are well suited to speech mixtures, due to the sparse structure of speech and also fullfills the assumption of approximate W-disjoint orthogonality [9, 17]. However, music mixtures involves interdependent harmonic structures, between sources, and spectral overlap is a common problem within multi-pitch estimation of musical content [18, 19, 20]. The problem of pitch estimation is important in applications such as separation, enhancement, transcription, classification, and source localization. Especially, parametric pitch estimators outperforms non-parametric methods in matters of distinguishing between the fundamental pitch period and multiples of it, and in performance under noisy conditions. The characteristics of the stereophonic mixture created in recording studios, have been shown as beneficial as parameters for multi-pitch estimation [1].

Virtual positioning by panning can be described with amplitude and delay ratios between the stereophonic channels. Amplitude ratios have been estimated in time-frequency domain, by producing an energy histogram for each time frame [16, 21]. Amplitude and delay ratios have been estimated succeeding a pitch estimate, using convex optimization [1, 22]. No solution exists for explicit estimation of amplitude and delay panning parameters without pitch information. This article describes such an algorithm.

1.2.1 Motivation

The problem of pitch estimation is important in applications such as separation, enhancement, transcription, classification, and source localization. Parametric pitch estimators outperforms non-parametric methods in matters of distinguishing between the fundamental pitch period and multiples of it, and in performance under noisy conditions. Recently, researchers within parametric pitch estimation, have made use of the stereophonic mixture created in recording studios and have implied the benefits of knowing the panning parameters. Panning can be described with amplitude and delay ratios between the stereophonic channels. Amplitude ratios have been estimated in time-frequency domain, by producing an energy histogram for each time frame. Amplitude and delay ratios have been estimated succeeding a pitch estimate, using convex optimization. No solution exists for explicit estimation of amplitude and delay panning parameters without pitch information. This project proposes the development of such an algorithm.

The project will involve panning parameter estimation in time-frequency domain, by clustering of amplitude and delay ratios. An optimal segmentation will be applied, which uses a metric such as the Bayesian Information Criteria or similar. The clusters

could be modeled as a mixture of Gaussians or similar. The project is carried out as a scientific research, using MATLAB for simulation and testing.

1.3 The Panning Parameters

The panning parameters discussed in this report is a product of the mixing process applied in sound studios. To enhance the sound quality and to ease the virtual perceived separation of sound sources in a stereo mixture the sound engineer can apply various effects, such as amplitude and delay panning. Other effects such as reverb, equalization and dynamic effects are usually also applied, but are of no interest in the remaining report. Amplitude and delay panning is exactly the two parameters that we estimate in this report, since they carry spatial information that is equivalent to direction and positioning in a real geometric setup.

1.3.1 Delay Panning

The delay panning parameter is directly related to the delay humans experience when a sound signal is received at the ears at separate time instances. Such a delay changes the direction of the source direction for the listener [23]. It has been shown that a constant time delay to one of the speakers is frequency dependent in terms of source positioning [24]. Though amplitude panning is the normal post processing “way to go” for sound engineers, delays are added as part of post processing mixing procedure both to correlate phases of microphones and to change directivity of sources and lastly delays longer than 1 ms can be applied to achieve depth and dimension, by placing the source mostly in the channel where the signal arrives first [25].

1.3.2 Amplitude Panning

Amplitude panning is the general method for changing the perceived direction of a sound source in a sound field between two or more loud speakers. Amplitude panning is an approximation of source localisation and its application ranges from stereophonic amplifiers to multichannel speaker setup and professional multi-channel mixing desks. Most often the user/engineer of a mixing desk can configure the the perceived direction of each individual sound source in the mix by turning one knob, attached to a trim pot that controls the signal voltage level to each speaker output. If the desk is digital, the user has a similar digital knob or slider interface.

Amplitude panning can be applied to multi-speaker setups, while the most common speaker configuration is a stereo setup, consisting of a left and a right speaker, with two audio channels being played back (one for each loudspeaker), whether it is a home audio hi-fi system, PA (Public Adress) system, headphone system etc. The stereophonic configuration is shown in Figure 1.1, where the listener is placed in oregon

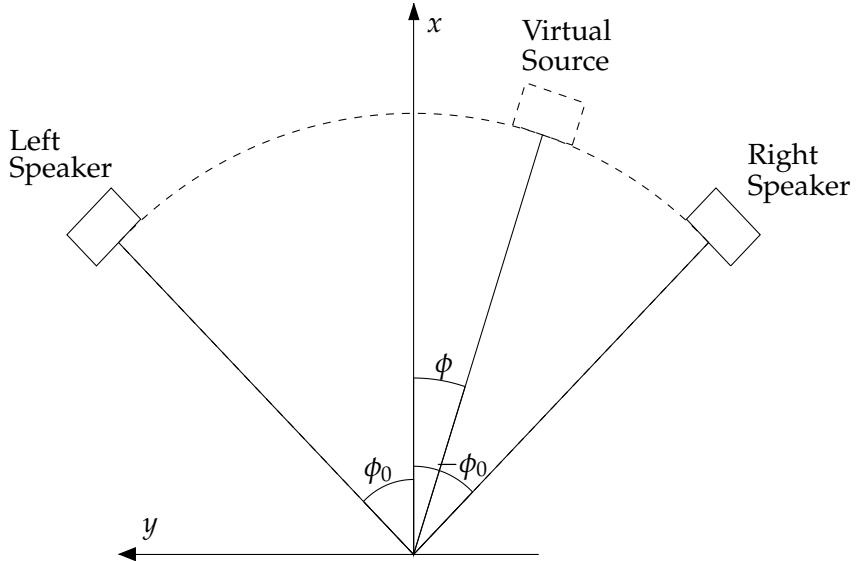


Figure 1.1: Stereophonic Configuration

($x = y = 0$). Amplitude panning in the stereophonic configuration is explained in the following sub section.

1.3.3 Stereophonic Amplitude Panning

Figure 1.1 shows the stereophonic sound configuration patented by Blumlein[26]. The listener is situated equidistant to each speaker in orego. The listener perceives an illusion of an auditory event, that is placed in a specific point on a two dimensional arc between the two speakers. The auditory event is moved by changing the signal amplitudes of the signal in the left and right channel. Amplitude panning is described by Ville Pulkki [27] in a vector based framework that allows two- and three dimensional speaker setups. Amplitude panning can be formulated at time t , by applying a signal $x(t)$ to both loudspeakers with different amplitudes, and gain factors for left and right channel respectively. In general the signal $x_i(t)$ is then

$$x_i(t) = g_i x(t), \quad i = 1, 2, \dots, N \quad (1.1)$$

where $x_i(t)$ is the signal applied to the i^{th} loudspeaker and g_i is the gain factor of the corresponding channel and $N = 2$ is the number of speakers in stereo configuration. While the virtual source is moving along the arc, the distance to the the listener should be constant. For the stereophonic configuration the vectorial distance of the gain factors g_1 and g_2 equals a constant C

$$g_1^2 + g_2^2 = C \quad (1.2)$$

The relation between the gain factors and the perceived virtual source direction has been derived for panning in the stereophonic configuration by Bauer [28] as the “stereophonic law of sines”, where the acoustic shadow of the head is not taken in to account and the sine law is assumed valid at all frequencies. For the sine law, the listener is situated symmetrically between the speakers in orego, facing along the x -axis in Figure 1.1. The stereophonic sine law is described by the ratio of the difference and sum of the gain factors as,

$$\frac{\sin \phi}{\sin \phi_0} = \frac{g_1 - g_2}{g_1 + g_2} \quad (1.3)$$

where ϕ is the perceived angle and ϕ_0 is the speaker base angle. It is required that $0^\circ < \phi_0 < 90^\circ$, $-\phi_0 \leq \phi \leq \phi_0$ and $g_1, g_2 \in [0, 1]$. An extension of the sine law is the tangent law, originally proposed by Bernfeld [29] as

$$\frac{\tan \phi}{\tan \phi_0} = \frac{g_1 - g_2}{g_1 + g_2} \quad (1.4)$$

The tangent law behaves similar to the sine law with very small difference, taking some of the listeners head complexity into account. Ville Pulkki [27] formulates the vector based approach as a reformulation of the tangent law, called the vector based amplitude panning (VBAP). Figure 1.2 visualizes the vector based framework of the stereo configuration, that is used in the remaining of this report to describe the estimated amplitude panning angle, shown in results and in figures.

Gain Vector Relation to Virtual Sound Source Positioning

To ease the understanding of the amplitude panning parameter it is convenient for the human reader to consider the parameter as a perceived angle in a carthesian coordinate system, since a music listener is normally placed in front of two speakers as mentioned in 1.3.3. To present the gain ratios as angles we apply the stereo vector base virtual sound source positioning [27]. A backwards amplitude panning algorithm serves the purpose of estimating the gain parameters. As visualized in Figure 1.1, each loudspeaker has a base angle $\phi_0 = \pm 45^\circ$ to the x -axis direction that the listener is facing towards; the listener is situated equidistant to each speaker in ($x = y$). The angle ϕ describes the virtual source position respective to the x -axis. The trigonometric functions are used for the panning gain since they fit the unit circle, thus they retain unity power along an arc as $1 = \cos^2 + \sin^2$. The gains are then

$$g_x = \cos \theta \quad (1.5)$$

$$g_y = \sin \theta \quad (1.6)$$

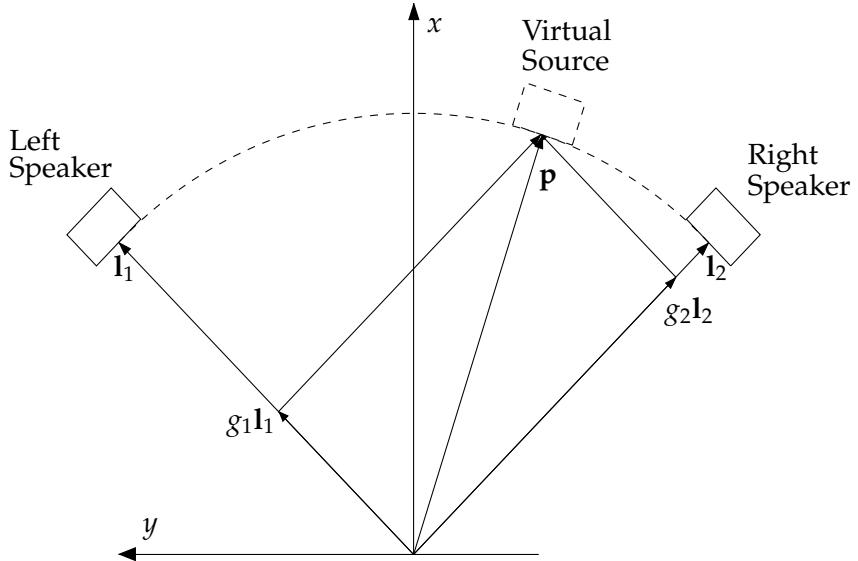


Figure 1.2: Stereophonic configuration with vector formulation

where $\theta = \phi + \phi_0$. If we define a loudspeaker base matrix \mathbf{L}

$$\mathbf{L} = [\mathbf{l}_1 \mathbf{l}_2]^T \quad (1.7)$$

consisting of two unit length loudspeaker vectors $\mathbf{l}_1 = [l_{11} l_{12}]^T$ and $\mathbf{l}_2 = [l_{21} l_{22}]^T$ pointing toward each speaker. In Figure 1.2, the unit vector \mathbf{p} points towards the virtual source as a linear combination of the gained loudspeaker vectors

$$\mathbf{p}^T = \mathbf{g}\mathbf{L} \quad (1.8)$$

This equation can be solved for the gain vector, by applying the inverse loudspeaker base matrix

$$\mathbf{g} = \mathbf{p}^T \mathbf{L}^{-1} \quad (1.9)$$

The loudspeaker base matrix \mathbf{L} is unitary and \mathbf{L}^{-1} exists under the conditions $0^\circ < \phi_0 < 90^\circ$, $-\phi_0 \leq \phi \leq \phi_0$ and $g_1, g_2 \in [0, 1]$. Finally, we can estimate the panning angle $\hat{\theta}$ as

$$\hat{\theta} = \arctan \frac{p(1)}{p(2)} \quad (1.10)$$

The amplitude panning angle applied to sources in a stereo mixture, can be estimated as shown from the obtained gain factors. The trigonometric functions used in this computation, estimates within the domain of the loudspeaker base matrix \mathbf{L} with a span of 90° . In professional studios the aperture of loudspeakers is typically 60° . However, this relation between to the tangent law is linear and is simply solved by normalization to a wider domain by multiplication.

1.4 Signal Model

In the following the signal model and assumptions are introduced. Consider an M -channel music mixture consisting of K unknown sources embedded in noise at time instant n . The data in the m^{th} channel is represented as $\mathbf{x}_m(n) \in \mathbb{R}^N$,

$$\mathbf{x}_m(n) = [x_m(n) \quad x_m(n+1) \quad \dots \quad x_m(n+N-1)]^T \quad (1.11)$$

for $m = 1, \dots, M$. The signals captured by channel m , relating to the k^{th} source are attenuated by gain coefficient $g_{m,k}$ and delayed by $\tau_{m,k}$ depending on their perceptual virtual positioning, given by the panning parameters. The signal mixture is modelled as a linear superposition of K attenuated and delayed sources embedded in noise $\mathbf{e}_{m,k}(n)$,

$$\mathbf{x}_m(n) = \sum_{k=1}^K g_{m,k} \mathbf{s}_k(n - f_s \tau_{m,k}) + \mathbf{e}_{m,k}(n) \quad (1.12)$$

where $g_{m,k}$ and $\tau_{m,k}$ are the attenuation and delay applied to the source $\mathbf{s}_k(n)$, respectively and f_s is the sampling frequency. Considering stereophonic mixtures with $M = 2$ for a stereo loudspeaker setup, amplitude panning is the normal procedure [28, 27] for virtual source positioning. In the post-processing of a every music production, delays can be added to enhance the spatial perception [23]. The trigonometric functions are often used for the panning attenuation because they induce a constant perceived distance between listener and the virtual source, described by $1 = \cos^2 + \sin^2$. The gains for channel m are expressed as [30],

$$g_m = \begin{cases} \cos \theta_k, & \text{for } m = 1 \\ \sin \theta_k, & \text{for } m = 2 \end{cases} \quad (1.13)$$

where $\theta = \phi + \phi_0$ is a sum of the perceived angle ϕ and the speaker base angles $\pm \phi_0 = 45^\circ$. Under the conditions $0^\circ < \phi_0 < 90^\circ$, $-\phi_0 \leq \phi \leq \phi_0$ and $g_1, g_2 \in [0, 1]$ the gains can be expressed as,

$$\mathbf{g}_k = \mathbf{p}_k \mathbf{L}^{-1} \quad (1.14)$$

where the unit-vector \mathbf{p} points towards the virtual source with \mathbf{L} as a unitary loudspeaker base matrix. For the stereophonic mixture ($M = 2$), we simplify notation by modelling attenuation and delay parameters as ratios between the frequency representations of active sources in the two channels.

1.4.1 Estimating the Panning Parameters

When only source k is active, the frequency representation in the two channels is,

$$S_{1,k}(\omega) = \sum_{n=1}^N s_k(n) e^{-j\omega n}, \quad (1.15)$$

$$S_{2,k}(\omega) = \sum_{n=1}^N \gamma_k s_k(n) e^{-j\omega n \delta_k}, \quad (1.16)$$

ω is the frequency grid, $\delta_k = f_s \tau_k$ is the relative delay of source k between the channels and γ_k is the relative attenuation factor corresponding to the ratio of attenuation of source k between the channels. The panning parameters γ_k and δ_k , that are associated with active sources in each frequency point can be computed as,

$$(\gamma_k, \delta_k) = \left(\left| \frac{S_{2,k}(\omega)}{S_{1,k}(\omega)} \right|, \frac{1}{\omega} \angle \frac{S_{1,k}(\omega)}{S_{2,k}(\omega)} \right) \quad (1.17)$$

where we must ensure that,

$$|\omega_{\max} \delta_{\max}| < \pi \quad (1.18)$$

to avoid phase ambiguity. Our aim is to estimate the panning parameters (γ, δ) for all K sources, along with an optimal segment length N , given only the stereophonic mixture in (1.4). The k th panning parameter is associated with only the k th source component, under the assumption that only one source is dominant at each frequency point. This is described by the approximate disjoint orthogonality expressed as [9],

$$S_{1,k}(\omega) S_{1,i}(\omega) \approx 0 \quad \forall \omega, k \neq i \quad (1.19)$$

To ease on this assumption, we apply a segmentation of the signal $\mathbf{x}_m(n)$ into segments of size N , that provides an optimal separation of the clusters in (1.4.1). The optimal segmentation is described in Section 5.1.1. The estimated amplitude and delay ratios, which often is referred to as the measurement vectors, are described from the spectral content of each channel in the stereophonic mixture as,

$$(\hat{\gamma}, \hat{\delta}) = \left(\left| \frac{X_2(\omega)}{X_1(\omega)} \right|, \frac{1}{\omega} \angle \frac{X_1(\omega)}{X_2(\omega)} \right) \quad (1.20)$$

where $X_m(\omega)$ is the discrete Fourier transform of $\mathbf{x}_m(n)$. Music mixtures often have a long duration of several minutes and we assume that such mixtures have stationary panning parameters throughout the full mixture. We will collect measurement vectors and perform segmentation to select parts of the signal that carries relevant information of the measurement vectors. Therefore we use a threshold that removes great part of the noisefloor of the spectrum and lowers computational complexity. We define an indicator function $b(\omega)$ as,

$$b(\omega) \begin{cases} 1, & |X_1(\omega)| |X_2(\omega)| > |\mathbf{X}_1|^T \mathbf{X}_2 / N \\ 0, & \text{otherwise} \end{cases} \quad (1.21)$$

where $X_m(\omega)$ is the pre-whitened DFT of $\mathbf{x}_m(n)$. It is possible to pick a specific amount of samples by increasing the threshold on the indicator function and improve on complexity. This was a description of the signal model and the main assumptions in this context. We will end this chapter by introducing a time-panning domain visualization, which we call the panogram. The next chapter will explain the clustering of measurement vectors.

1.4.2 Visualizing the Amplitude Angle as a Panogram

A visual output of the panning angle can be used to identify the various sources in a stereo mix based on their panning coefficient. This can be accomplished via the amplitude panning ratio in (1.4.1). The computation is very fast and the output is shown in Figure 1.3 for a multi-pitch mixture of two instruments, trumpet and horn, playing the notes C4 (262 Hz) and F#4 (370 Hz), respectively. The mixture is also used in an experiment in [31]. The algorithm for the panogram in Figure 1.3 is based on

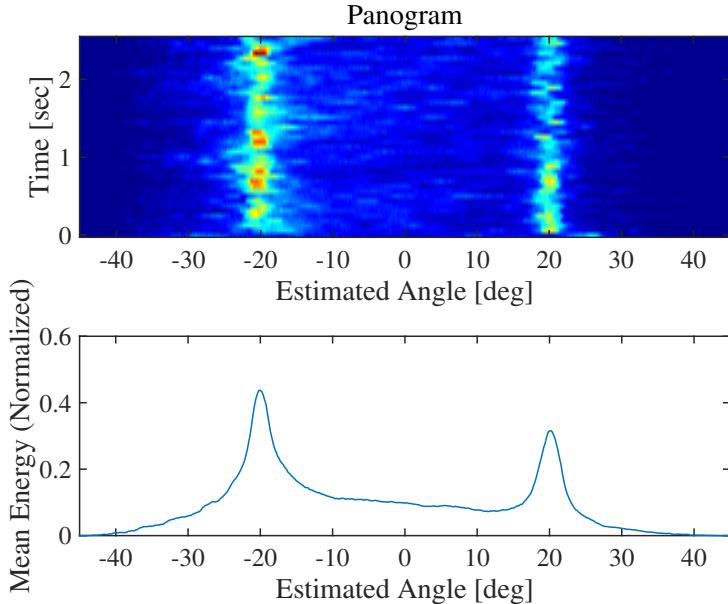


Figure 1.3: Panogram of a multi-pitch mixture of two instruments, trumpet and horn.

searching through the power ratio of the absolute discrete Fourier transform of the two stereo channels. By realizing that,

$$\frac{g_x}{g_y} = \tan \theta \quad (1.22)$$

where g_x and g_y are defined in (1.3.3) and (1.3.3). The ratio expressed as $\tan \theta$ can be considered as a search space given by the aperture of the loudspeakers as a function of θ . Consider a stereo input signal $x(n)$ consisting of both $x_1(n)$ and $x_2(n)$ at time instant n . At each time instant we compute the amplitude ratio $\gamma_n(\omega)$ using 1.4.1 and lastly marginalize by summing over all frequencies and at each time instant n , the panogram $p(\theta)$ is a vector function of θ and can be expressed the power at each panning angle. More Panograms can be found in Appendix A. The visual inspection and manual peak finding in an objective function or histogram created from time-frequency domain ratios, is used within the research area of blind source separation [9, 21, 17] and also used for channel upmix techniques [32], however the aim in the following is to automatically estimate the panning parameters.

Chapter 2

Clustering

Once the measurement space containing the distribution of estimated panning parameters is well defined, it is the aim to estimate the number of sources and the source parameters as an unsupervised learning task, with no prior information given of the source parameters. The problem of finding clusters in a set of measurement vectors can be approached by using probabilistic techniques or non-probabilistic techniques. An example of a non-probabilistic clustering technique is the k -means algorithm [33]. The immediate K -means algorithm requires an input that specifies the number of clusters to estimate. We have modelled the source parameter distribution using the probabilistic clustering technique, as a mixture of Gaussians.

2.1 Estimation of Source Parameters

The source parameters are estimated by maximizing the likelihood. The maximum likelihood estimates are the parameters of the model that describe the observed measurement vectors the best, i.e. the parameters that maximizes the probability of the observed data, \mathbf{x} , given the parameters,

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}|\mathbf{x}) = \arg \max_{\boldsymbol{\theta}} p(\mathbf{x}|\boldsymbol{\theta}) \quad (2.1)$$

where $\hat{\boldsymbol{\theta}}$ is a vector containing the model parameters. In the following the probabilistic model is described along with the K -means clustering algorithm that is used for initialization. We describe the maximum likelihood estimator, using latent variables and finally we consider the model order selection as both a probabilistic and non-probabilistic method.

2.2 Finite Mixtures

The following section is a brief description of the general model of finite mixtures, which the Gaussian mixture model belongs to. We have found the Gaussian mixture

to be well suited for modelling the source parameter distribution. The research issue of order selection is relevant, when aiming to jointly estimate source parameters and number of sources in the stereophonic mixture. We can describe the stereophonic mixture as a finite mixture of K random sources described as probability density functions,

$$p_k(\mathbf{x}), \quad k = 1, \dots, K \quad (2.2)$$

We observe a set of random independent distributed samples, coming from these probability density functions. We define the prior probability of observing data from source s_k as $p(s_k) = \alpha_k$, and the conditional probability of the data given source s_k is $p(\mathbf{x}|s_k) = p_k(\mathbf{x})$, thus the joint probability $p(\mathbf{x}, s_k)$ is expressed as $\alpha_k p_k(\mathbf{x})$. Finally, the unconditional probability density is,

$$p(\mathbf{x}) = \sum_{k=1}^K \alpha_k p_k(\mathbf{x}) \quad (2.3)$$

Which means we that the unconditional density is a finite mixture of component densities $p_k(\mathbf{x})$ weighted by their prior, referred to as the mixing probabilities which we denote α_k for the k th source. The mixing probabilities has the general constraint of summing to one.

2.2.1 Parameterization of the finite mixtures of Gaussians

The unknown parameter vector is denoted by θ . In general for a finite mixture model it will be consisting of the mixing probability and the unknown parameters,

$$\Theta = \{\theta_1, \theta_2, \dots, \theta_k, \alpha_1, \alpha_2, \dots, \alpha_k\}$$

The conditional densities related to the source components are then given by,

$$p(\mathbf{x}|\theta) = \sum_{k=1}^K \alpha_k p(\mathbf{x}|\theta_k)$$

By assuming that sources are Gaussian distributions with arbitrary covariance the conditional density is modelled as,

$$p(\mathbf{x}|\theta_k) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \mathbf{C}_k)$$

the parameter vector contains the mean $\boldsymbol{\mu}_k$ and covariance \mathbf{C}_k for $i = k, \dots, K$,

$$\Theta = \{\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \dots, \boldsymbol{\mu}_k, \mathbf{C}_1, \mathbf{C}_2, \dots, \mathbf{C}_k, \alpha_1, \alpha_2, \dots, \alpha_k\}$$

The aim is now to estimate the parameter set θ from the given observations. An order selection procedure will estimate the given number of sources, while panning parameters are given by the mean of the mixture components. However, the task

of assigning points to mixture components is not trivial to do automatically, since the observed data with unknown classes, can be clustered in to an arbitrary number of classes, dependent on the choice of model and how the model is being fitted to the observations. The aim in such an unsupervised learning task, by model based clustering is that each component models one cluster.

$$\sum_{k=1}^K \alpha_k = 1 \quad \text{and} \quad 0 \leq \alpha_k \leq 1 \quad (2.4)$$

as any probability the mixing probability is required to take a value between 0 and 1. The finite mixture in this general form is possible to parameterize with the unknown parameters and by applying some model to the distribution, we can build a convenient estimator, as we will do in the following section.

2.2.2 The Gaussian mixture model as a likelihood

As discussed in Section 2.4, the K-means assigns every measurement vector uniquely to one cluster as a hard assignment. However, it is not clear that a measurement vector which is placed midway between two cluster centers is assigned appropriately, relative to the cluster center which can affect the precision of the parameter estimates. By using probabilistic models such as the Gaussian mixture model (GMM), the assignments can reflect this level of uncertainty as a soft assignment of measurement vectors to clusters. Furthermore, the mixture model is good at representing class conditional densities in supervised learning, because mixtures can approximate arbitrary densities, i.e. two strongly non-Gaussian classes, can be modelled by mixtures of each class conditional density [34]. On the contrary, in the unsupervised learning task it is a matter of fitting the model sparsely to the data without overfitting to parameter space. Therefore, the Gaussian mixture model will firstly be described as a likelihood, followed by an interpretation as an a posteriori distribution, penalizing higher model orders.

Using the GMM framework, the full parameter space is modelled as a Gaussian mixture distribution i.e. a linear superposition of Gaussians,

$$p(\mathbf{x}) = \sum_{k=1}^K \alpha_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{C}) \quad (2.5)$$

$$p(\mathbf{x}) = \sum_{k=1}^K \alpha_k \frac{(2\pi)^{-\frac{d}{2}}}{|\mathbf{C}_k|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \mathbf{C}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\} \quad (2.6)$$

where $\boldsymbol{\mu}_k$ is the mean and \mathbf{C}_k is the covariance of the k th Gaussian. The mixing probabilities $\{\alpha_1, \dots, \alpha_K\}$ are constrained to

$$\sum_{k=1}^K \alpha_k = 1, \quad 0 \leq \alpha_k \leq 1 \quad (2.7)$$

and can be interpreted as the prior probabilities of having the class k ,

$$p(\mathbf{x}) = \sum_{k=1}^K p(k) p(\mathbf{x}|\boldsymbol{\theta}_k) = \sum_{k=1}^K \alpha_k \frac{(2\pi)^{-\frac{d}{2}}}{|\mathbf{C}_k|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \mathbf{C}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\} \quad (2.8)$$

where each θ_k is the parameter specifying the k th component. The parameter vector is defined as,

$$\boldsymbol{\Theta} \equiv \{\alpha_1, \dots, \alpha_K, \boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_K, \mathbf{C}_1, \dots, \mathbf{C}_K\} \quad (2.9)$$

The parameter vector specifies the full mixture as the complete set of parameters. Observing a set of N independent distributed samples $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, the log-likelihood function corresponding to a K-source mixture is,

$$\mathcal{L}(\boldsymbol{\theta}|\mathbf{X}) = \ln p(\mathbf{X}|\boldsymbol{\alpha}, \boldsymbol{\mu}, \mathbf{C}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \alpha_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \mathbf{C}_k) \right\} \quad (2.10)$$

Maximizing the log-likelihood of (2.10), turns out to be a complex problem mainly due to the summation inside the logarithm. The logarithm function of (2.10) does not act directly on the Gaussian, but also on the summation over k . If we differentiate the log-likelihood and set it to zero it will not have a closed form solution. However, we can maximize the likelihood function with the expectation-maximization (EM) algorithm. In the following section we proceed with a general description of the EM in the context of fitting a mixture of Gaussians to the measurent vectors.

2.2.3 Fitting the Gaussian mixture as a maximum likelihood solution

A powerful method for finding the maximum likelihood solutions to models with latent variables is the EM algorithm [35, 36]. Due to the inner sum of (2.10), it is necessary to view the problem by defining a K -dimensional binary latent variable \mathbf{z} that for a given n has k latent variables where only one of these is equal to 1, while the rest are equal to 0. This means that the vector \mathbf{z} has K possible states and $z_k \in \{0, 1\}$ and $\sum_{k=1}^K z_k = 1$. We can then view α_k as the prior probability $p(z_k = 1) = \alpha_k$, i.e. the probability of z_k equals 1. In these terms the marginal distribution over \mathbf{z} can be written in the form,

$$p(\mathbf{z}) = \prod_{k=1}^K \alpha_k^{z_k} \quad (2.11)$$

The conditional distribution of \mathbf{x} given \mathbf{z} is also a Gaussian and can be described by,

$$p(\mathbf{x}|\mathbf{z}) = \prod_{k=1}^K \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_k, \mathbf{C}_k)^{z_k} \quad (2.12)$$

We are now able to work with the joint distribution $p(\mathbf{x}, \mathbf{z}) = p(\mathbf{z})p(\mathbf{x}|\mathbf{z})$. By summing the joint distribution over all possible states of \mathbf{z} , we can obtain the marginal distribution of \mathbf{x} as,

$$p(\mathbf{x}) = \sum_{\mathbf{z}} p(\mathbf{z})p(\mathbf{x}|\mathbf{z}) = \sum_{k=1}^K \alpha_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \mathbf{C}_k)$$

Which is equivalent to the form of the Gaussian mixture expressed as a linear superposition of Gaussian distributions as given by (2.2.2), only now there is a corresponding latent variable for each measurement vector \mathbf{x}_n . Observing a set of N independent distributed samples $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$, the log-likelihood function corresponding to a K -source mixture can now be expressed for the complete measurement vectors $\{\mathbf{X}, \mathbf{Z}\}$ containing both the observed data \mathbf{X} and the latent variable \mathbf{Z} [37]. The log-likelihood is then expressed as,

$$\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\alpha}, \boldsymbol{\mu}, \mathbf{C}) = \sum_{n=1}^N \sum_{k=1}^K z_{n,k} \{\ln\{\alpha_k\} + \ln \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \mathbf{C}_k)\} \quad (2.13)$$

Since the logarithm now acts directly on the Gaussian distribution, it leads to much simpler solution for the maximum likelihood. In practice, the values of the latent variables are unknown, thus we consider the expectation with respect to the posterior distribution of the latent variables, which takes the form,

$$p(\mathbf{Z}|\mathbf{X}, \boldsymbol{\alpha}_k, \boldsymbol{\mu}_k, \mathbf{C}_k) \propto \prod_{n=1}^N \prod_{k=1}^K (\alpha_k \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \mathbf{C}_k))^{z_{n,k}} \quad (2.14)$$

Where $\alpha_k = \frac{1}{N} \sum_{n=1}^N z_{n,k}$. The expected value of the complete data log-likelihood function is now,

$$\mathbb{E}_{\mathbf{Z}}[\ln p(\mathbf{X}, \mathbf{Z}|\boldsymbol{\alpha}, \boldsymbol{\mu}, \mathbf{C})] = \sum_{n=1}^N \sum_{k=1}^K \beta(z_k) \{\ln\{\alpha_k\} + \ln \mathcal{N}(\mathbf{x}_n|\boldsymbol{\mu}_k, \mathbf{C}_k)\} \quad (2.15)$$

Where $\beta(z_k)$ is a quantity that plays an important role as the conditional probability of \mathbf{z} given \mathbf{x} . By viewing α_k as the prior probability of $z_k = 1$ and $\beta(z_k)$ as the corresponding posterior once we have observed \mathbf{x} . The quantity $\beta(z_k)$ is also referred to as the responsibility that component k takes for explaining the observation of \mathbf{x} .

$$\beta(z_k) \equiv p(z_k = 1|\mathbf{x}) = \mathbb{E}[z_k] = \frac{\alpha_k \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_k, \mathbf{C}_k)}{\sum_{j=1}^K \alpha_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \mathbf{C}_j)} \quad (2.16)$$

The responsibility $\beta(z_k)$ applies different weight for each parameter estimate, which turns out to be crucial for the model selection procedure for the mixture model, as described in Section ??.

2.3 Maximizing the likelihood

Now that we have defined the log-likelihood by using latent variables to describe the complete data, we are ready to apply the EM-algorithm for the Gaussian mixture models. The condition that must be satisfied at the maximum of a likelihood function is found by setting the derivatives of $\ln p(\mathbf{X}|\boldsymbol{\alpha}, \boldsymbol{\mu}, \mathbf{C})$ in (2.10) to zero. First the mean parameter:

$$\frac{d}{d\boldsymbol{\mu}_k} \ln p(\mathbf{x}|\boldsymbol{\alpha}, \boldsymbol{\mu}, \mathbf{C}) = 0 \quad (2.17)$$

$$\sum_{n=1}^N \frac{\alpha_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \mathbf{C}_k)}{\sum_j \alpha_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \mathbf{C}_j)} \mathbf{C}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) = 0 \quad (2.18)$$

Where it is interesting that the responsibility of (2.2.3) appears naturally, and the expression is equivalent to,

$$\sum_{n=1}^N \beta(z_{n,k}) \mathbf{C}_k^{-1} (\mathbf{x}_n - \boldsymbol{\mu}_k) = 0 \quad (2.19)$$

When we multiply by \mathbf{C}_k we can rearrange the expression to,

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \beta(z_{n,k}) \mathbf{x}_n \quad (2.20)$$

where $N_k = \sum_{n=1}^N \beta(z_{n,k})$, can be interpreted as the effective number of points assigned to cluster k . Therefore, $\boldsymbol{\mu}_k$ for the k th Gaussian component is obtained by taking a weighted mean of all the points in the measurement vectors. The weight is given by the posterior probability $\beta(z_{n,k})$ that component k was for generating \mathbf{x}_n .

The maximum likelihood solution for the covariance \mathbf{C}_k is found by,

$$\frac{d}{d\mathbf{C}_k} \ln p(\mathbf{x}|\boldsymbol{\alpha}, \boldsymbol{\mu}, \mathbf{C}) = 0 \quad (2.21)$$

$$\mathbf{C}_k = \frac{1}{N_k} \sum_{n=1}^N \beta(z_{n,k}) (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T \quad (2.22)$$

where each measurement vector also is weighted by the responsibility $\beta(z_{n,k})$. The maximum likelihood solution for the mixing probability is derived in [37] and it is,

$$\alpha_k = \frac{N_k}{N} \quad (2.23)$$

where $N_k = \sum_{n=1}^N \beta(z_{n,k})$. This means that the mixing coefficient for the k th component is given by the average responsibility which the component takes for explaining the measurement vectors. It is now possible to proceed with the EM-algorithm to obtain the maximum likelihood estimate for the particular case of the Gaussian mixture model.

EM-algorithm for the complete measurement vectors

1. Choose an initial value for the parameter vector θ^{old} .
2. (**E-step**). Evaluate $p(\mathbf{Z}|\mathbf{X}, \theta^{\text{old}})$, by evaluation of the responsibilities of the current parameter values.

$$\beta(z_{n,k}) = \frac{\alpha_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \mathbf{C}_k)}{\sum_j \alpha_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \mathbf{C}_j)} \quad (2.24)$$

3. (**M-step**). Evaluate θ^{new} , re-estimating the parameters using the current probabilities,

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n=1}^N \beta(z_{n,k}) \mathbf{x}_n$$

$$\mathbf{C}_k = \frac{1}{N_k} \sum_{n=1}^N \beta(z_{n,k}) (\mathbf{x}_n - \boldsymbol{\mu}_k) (\mathbf{x}_n - \boldsymbol{\mu}_k)^T$$

$$\alpha_k = \frac{N_k}{N}$$

where $N_k = \sum_{n=1}^N \beta(z_{n,k})$.

4. Evaluate the log-likelihood

$$\mathcal{L}(\theta|\mathbf{x}) = \ln p(\mathbf{X}|\boldsymbol{\alpha}, \boldsymbol{\mu}, \mathbf{C}) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \alpha_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \mathbf{C}_k) \right\}$$

Check for convergence. If no convergence, then update, $\theta^{\text{old}} \leftarrow \theta^{\text{new}}$ and go to step 2.

We have considered how to use the EM-algorithm to maximize the likelihood, when there are discrete latent variables. This has been derived for the Gaussian mixture model. In the following we will also use the EM-algorithm for finding the maximum a posterior solutions (MAP). Since our aim is to estimate a given number of clusters and their respective source parameters from the measurement vectors alone, we will use the MAP model. The MAP model adds a prior $p(\theta)$ to the log-likelihood expression in (2.10). The prior is defined over the parameters and a suitable choice of the prior will improve the model selection.

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} \{ \ln p(\mathcal{X}|\theta) + \ln p(\theta) \} \quad (2.25)$$

Model selection is explained in Section 3. In the following section we will describe the non-probabilistic clustering method of K -means. We use this algorithm for initialization and furthermore, the model selection criteria of Calinski-Harabasz fits well to the K -means algorithm and we will explain this connection also.

2.4 Initialization using K-Means Clustering

Given the unlabelled measurement vectors, the aim is to estimate the corresponding unknown parameter vector θ , which can be done using the non-probabilistic method of K-means. Once we have estimated which points go to which cluster, we can estimate a Gaussian mean and covariance for that cluster. It is unlikely that the guess is right the first time, but based on the initial estimates of parameters, it is possible to make a better guess at pairing points with components, in an iterative procedure using the EM-algorithm. We consider the problem of identifying clusters of measurement vectors in a multidimensional space. We observe N observations of the measurement vectors $\mathbf{x}_1, \dots, \mathbf{x}_N$. In general the variable \mathbf{x} is D -dimensional. However, we have defined two parameters in this study ($D = 2$). Each cluster center is represented by μ_k after we have assigned each point in the measurement vectors to a given cluster. The assignment of a measurement vector \mathbf{x}_n to cluster k is described by the binary indicator variable $b_{n,k} \in \{0, 1\}$. The aim is to minimize the sum of squares distance from each measurement vector to its closest center vector μ_k . We can now describe a cost function J as,

$$J = \sum_{n=1}^N \sum_{k=1}^K b_{n,k} \|\mathbf{x}_n - \mu_k\|^2 \quad (2.26)$$

Finding the values of $b_{n,k}$ and μ_k that will minimize J is done by using an iterative optimization procedure, involving two steps for each iteration. To begin the iterations, some initial values are assigned to μ_k . The two iterative steps are,

- Minimize J with respect to $b_{n,k}$, with μ_k fixed.
- Minimize J with respect to μ_k , with $b_{n,k}$ fixed. .

The assignment of the n th measurement vector to the closest cluster center can be expressed as,

$$b_{n,k} = \begin{cases} 1, & \text{if } k = \arg \min_j \|\mathbf{x}_n - \mu_j\|^2 \\ 0, & \text{otherwise} \end{cases} \quad (2.27)$$

Since the cost function J is a quadratic function of μ_k we differentiate with respect to μ_k and set it to zero,

$$\frac{d}{d\mu_k} J = 2 \sum_{n=1}^N b_{n,k} (\mathbf{x}_n - \mu_k) = 0 \quad (2.28)$$

and solve for μ_k ,

$$\mu_k = \frac{\sum_{n=1}^N b_{n,k} \mathbf{x}_n}{\sum_{n=1}^N b_{n,k}} \quad (2.29)$$

which expresses that μ_k is the mean of all measurement vectors \mathbf{x}_n assigned to cluster k . The iteration over these two steps are guaranteed to reach convergence. The K-means

assigns every measurement vector uniquely to one cluster, and it is not clear that a measurement vector which is placed midway between two cluster centers is assigned appropriately, but by using probabilistic models such as the Gaussian mixture model (GMM), the assignments can reflect this level of uncertainty. For initialization of the EM-algorithm by deliberately overfitting, i.e. choosing a K much larger than the expected value, the K-means algorithm assures that the true parameter are among the estimates, making it convenient to use it for initialization of the GMM-EM algorithm before applying the MAP model selection to the GMM-model.

Model selection using K-means

It is possible to evaluate the k -means for different number of clusters and then choose the optimal number of clusters based on the variance ratio criterion [38]. The variance ratio criterion (VRC) is based on the ratio between the overall between-cluster variance and the overall within-cluster variance. We run a short experiment with 7 sources. From the scatter plot in Figure 2.1, we can see that the correct number of clusters have

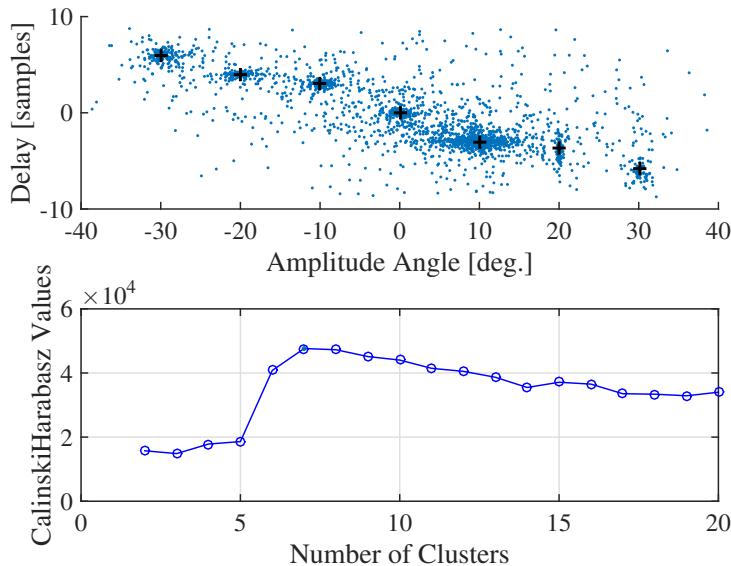


Figure 2.1: Model selection using the Calinski Harabasz criterion on mixture of seven sources.

been found in this specific case. However, the K -means clustering algorithm can be stuck in a local minimum rather than the global and it is therefore dependent on the initialization to be well considered. An initialization of the K -means have been proposed as the K -means++ algorithm by [39], a variant that chooses centers at random from the measurement vectors, but weighs the measurement vectors according to their squared distance, squared from the closest center that has already been chosen. This gives a faster convergence and overcomes some of the local minimum problems. Although the K -means clustering algorithm offers no accuracy guarantee, its simplicity

is very appealing in practice, thus it is widely used for clustering.

Calinski Harabasz Evaluation

The selected clusters in the measurement vectors shown with black plus signs in Figure 2.1, were subjected to a cluster validation algorithm called the Calinski-Harabasz [38] or the Variance Ratio Criterion (VRC), which is similar to the Inter-Intra class distance [40]. The validation algorithm selects the subset of clusters that maximizes the cluster separability. It is based on the Euclidean distance measure between measurement vectors in the measurement vectors. The assumption of mutually exclusive clusters leads to the assumption that the expectation vectors of the different cluster centroids are discriminating [40]. The optimal measure is a monotonically increasing function of the distance between expectation vectors and an increasing function of the scattering around the expectations. The conditional expectation of the measurement vectors given the cluster is the sample mean $\hat{\mu}_k$:

$$\hat{\mu}_k = \frac{1}{N_k} \sum_{n=1}^{N_k} \mathbf{x}_{k,n} \quad (2.30)$$

where $\mathbf{x}_{k,n}$ are measurement vectors from cluster C_k . The unconditional expectation of the measurement vector \mathbf{x} is the sample mean of the full measurement vectors $\hat{\mu}$:

$$\hat{\mu} = \frac{1}{N_s} \sum_{n=1}^{N_s} \mathbf{x}_n \quad (2.31)$$

where $N_s = \sum N_k$ is all samples in the set. The scattering of vectors from a given class C_k is:

$$\mathbf{S}_k = \frac{1}{N_k} \sum_{n=1}^{N_k} (\mathbf{x}_{k,n} - \hat{\mu}_k)(\mathbf{x}_{k,n} - \hat{\mu}_k)^T \quad (2.32)$$

It is analogous to a covariance matrix. The scatter matrix describing the noise is called the within-scatter matrix \mathbf{S}_w . Averaged over all classes it describes the average scatter within classes.

$$\mathbf{S}_w = \frac{1}{N_s} \sum_{k=1}^K N_k \mathbf{S}_k \quad (2.33)$$

Complementary to the within-scatter \mathbf{S}_w is the between scatter matrix \mathbf{S}_b that describes the scattering of class dependent sample means around the overall average:

$$\mathbf{S}_b = \frac{1}{N_s} \sum_{k=1}^K N_k (\hat{\mu}_k - \hat{\mu})(\hat{\mu}_k - \hat{\mu})^T \quad (2.34)$$

With these definitions we can express the Calinski-Harabasz criterion as,

$$CH_k = \frac{\mathbf{S}_b}{\mathbf{S}_w} \frac{N_s - k}{k - 1} \quad (2.35)$$

To determine the optimal number of clusters, we maximize CH_k with respect to k . The optimal number of clusters is the solution with the highest Calinski-Harabasz index value. The rightmost fraction of 2.4 is different from the inter-intra class distance ratio of [40]. Basically, this fraction expresses that we maximize the criterion by explaining large amount of observations by few clusters. The scatter within and scatter between ratio can be regarded as a signal to noise ratio, but it does not alter the underlying tendency of the K-means clustering algorithm, to overfit the measurement space and get stuck in a local miximum. Therefore, it seems reasonable to use the probabilistic method for cluster evaluation and only the K-means for initialization. Figure 2.2 shows the correct estimated model order, however the K-Means clustering algorithm is stuck in a local minimum and has therfore missed one of the true clusters. Since it is a

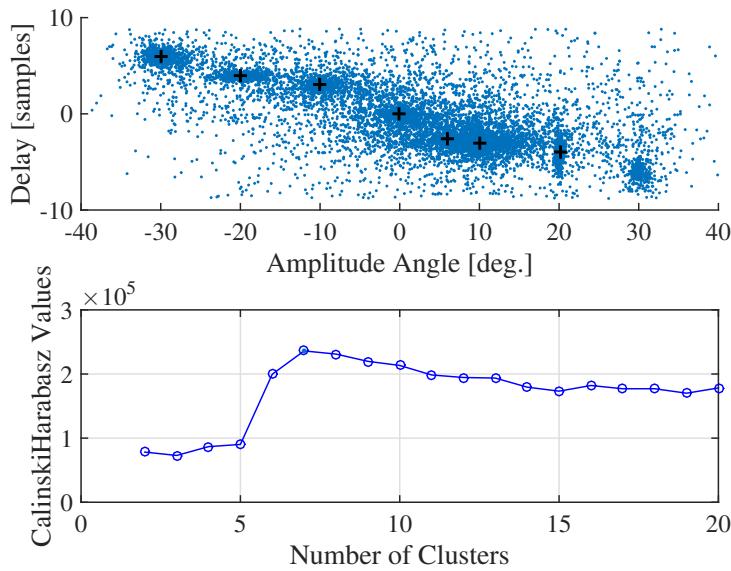


Figure 2.2: Model selection using the Calinski Harabasz criterion on a mixture of seven sources. In this specific evaluation the K-means is overfitting to the measurements compared to Figure 2.1.

strong criterion for non-probabilistic model order selection which fits well to the K-means clustering, which is a minimizer of the squared error, we will test its ability to be used for segmentation, only then we normalize the Calinski-Harabsz criterion to the measurement space. The normalized objective functions of the two given examples of Figure 2.1 and 2.2 are shown in Figure 2.3.

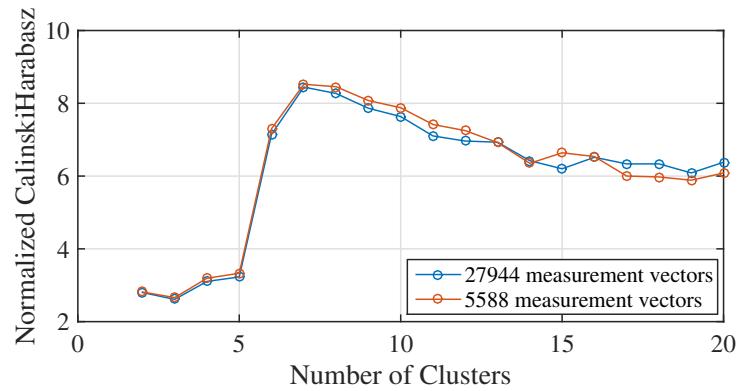


Figure 2.3: The normalized Calinski Harabasz objective functions of Figure 2.1 and 2.2.

Chapter 3

Model Order Selection

One advantage of the mixture model approach to clustering is that it allows the use of approximate Bayes factors to compare models. A thorough comparison of Bayes factors can be read in [41]. The model order selection and the segmentation can be done with a *maximum a posteriori* (MAP) criterion. The MAP estimator is,

$$\hat{\theta}_{\text{MAP}} = \arg \max_{\theta} \{ \ln p(\mathbf{x}|\theta) + \ln p(\theta) \} \quad (3.1)$$

where $p(\theta)$ is the prior on the parameters and \mathbf{x} is the observed data. We will introduce the MAP criterion in the following.

There exists several approaches for finding a solution to the model order estimate. Two of these are very often used [42, 43], namely the AIC and the MDL, which formally coincides with the BIC, why we will describe the MDL as a special case MAP criterion in the following section. The AIC is given as [43],

$$\mathcal{M}_s = \arg \min_{\mathcal{M}_k} \{ -\ln p(\mathbf{x}|\hat{\theta}, \mathcal{M}_k) + N_p \} \quad (3.2)$$

The MDL is,

$$\mathcal{M}_s = \arg \min_{\mathcal{M}_k} \left\{ -\ln p(\mathbf{x}|\hat{\theta}, \mathcal{M}_k) + \frac{N_p}{2} \ln N \right\} \quad (3.3)$$

where \mathcal{M}_s is the selected model, \mathbf{x} is the observed measurement vector, $p(\mathbf{x}|\hat{\theta}, \mathcal{M}_k)$ is the probability density function of the data given the model parameters and the model, θ is the parameter vector and $\hat{\theta}$ is the maximum likelihood of θ and N_p is the dimension of θ .

3.0.1 The Asymptotic MAP criterion

The principle of the MAP is choosing the model \mathcal{M} that maximizes the posterior probability given the observed data \mathbf{x} ,

$$\widehat{\mathcal{M}} = \arg \max_{\mathcal{M}} p(\mathcal{M}|\mathbf{x}) \quad (3.4)$$

expressed by using Bayes method,

$$\widehat{\mathcal{M}} = \arg \max_{\mathcal{M}} \frac{p(\mathbf{x}|\mathcal{M})p(\mathcal{M})}{p(\mathbf{x})} \quad (3.5)$$

Choosing a uniform prior $p(\mathcal{M})$ to not favour any model beforehand and noting that once \mathbf{x} is observed $p(\mathbf{x})$ is constant and the MAP model reduces to the likelihood of the observed data given the model,

$$\widehat{\mathcal{M}} = \arg \max_{\mathcal{M}} p(\mathbf{x}|\mathcal{M}) \quad (3.6)$$

where the likelihood is dependent on the parameters, $\boldsymbol{\theta}$. In the Bayesian framework we obtain the marginal density of the measurements given the model, by integrating the parameters out [42],

$$p(\mathbf{x}|\mathcal{M}) = \int_{\Theta} p(\mathbf{x}|\boldsymbol{\theta}, \mathcal{M})p(\boldsymbol{\theta}|\mathcal{M})d\boldsymbol{\theta} \quad (3.7)$$

The asymptotic approximation to this integral is found by assuming high amounts of data, when the most significant peaks occur in the likelihood function around the maximum likelihood estimates $\hat{\boldsymbol{\theta}}$. (3.7) becomes equal to [43],

$$p(\mathbf{x}|\mathcal{M}) = (2\pi)^{N_p/2} \det(\widehat{H})^{-1/2} p(\mathbf{x}|\hat{\boldsymbol{\theta}}, \mathcal{M})p(\hat{\boldsymbol{\theta}}|\mathcal{M}) \quad (3.8)$$

where \widehat{H} is the Hessian of the log-likelihood function when evaluated at the $\hat{\boldsymbol{\theta}}$,

$$\widehat{\mathcal{H}} = -\left. \frac{\partial^2 \ln(p(\mathbf{x}|\boldsymbol{\theta}, \mathcal{M}))}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \right|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} \quad (3.9)$$

By neglecting terms of order $\mathcal{O}(1)$, the asymptotic MAP expression is found by taking the negative logarithm of (3.8), where the term $2\pi^{N_p/2}$ can be assumed constant for asymptotic signal length N , while a weak prior on $p(\boldsymbol{\theta}|\mathcal{M})$ has been used [43] to obtain the MAP expression [42],

$$\widehat{\mathcal{M}} = \arg \min_{\mathcal{M}} \left\{ -\ln p(\mathbf{x}|\hat{\boldsymbol{\theta}}, \mathcal{M}) + \frac{1}{2} \ln \det(\widehat{\mathcal{H}}) \right\} \quad (3.10)$$

where the first term is the log-likelihood and the last term is the penalty added. The first term of the criterion decreases when the complexity of the model increases, and by contrast, the second term increases and acts as a penalty for using additional parameters to model the data. The penalty term is found by noting that the Hessian in (3.9) can be replaced by the Fisher information matrix since the error it introduces is smaller than the neglected terms of order $\mathcal{O}(1)$ [42, 43]. The Hessian is then,

$$\widehat{\mathcal{H}} \approx -E \left\{ \frac{\partial^2 \ln(p(\mathbf{x}|\boldsymbol{\theta}))}{\partial \boldsymbol{\theta} \partial \boldsymbol{\theta}^T} \right\} \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} \quad (3.11)$$

Under the assumptions of the observed measurement being real, independent and identically distributed, we can write,

$$\det \widehat{\mathcal{H}} = \mathcal{O}(N^{\frac{N_p}{2}}) \quad (3.12)$$

The interested reader can find specific details on this assumption in [43]. The expression in (3.10) then reduces to,

$$\widehat{\mathcal{M}} = \arg \min_{\mathcal{M}} -\ln p(\mathbf{x}|\hat{\theta}, \mathcal{M}) + \frac{N_p}{2} \ln N \quad (3.13)$$

which is the MDL that formally coincides with the BIC. For the case of the multivariate Gaussian distribution with arbitrary covariance $N_p = d + d(d+1)/2$. The expression in (3.13) is not valid for all signal processing families of models. In fact, for the Gaussian mixture model this rule will not be directly appropriate for model order selection without applying priors to the parameters as described in Section 3.0.2. Figure 3.1 and 3.2 shows the AIC and and the asymptotic MAP, referred to as the BIC under the described assumptions. Both figures show the criterion applied to seven sources, us-

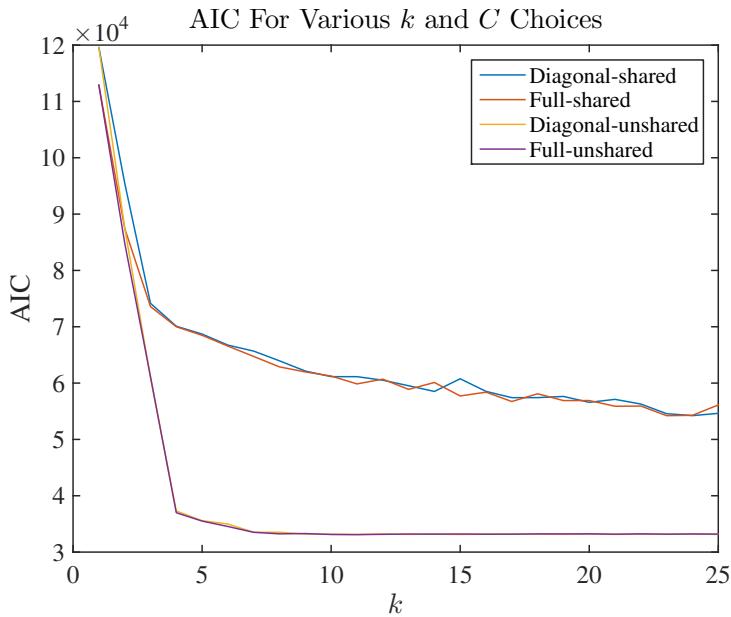


Figure 3.1: AIC Curve for fitting with the Gaussian mixture model.

ing the MAP criterion implemented by the EM-algorithm. From both of the figures it is clear that the criterion results in a monotonically decreasing function of the model order. This tendency to overfitting is due to the fact that the measurement vectors does not have equal weight in each parameter estimate in the Gaussian mixture model. The penalty term dependent on N is not sufficient for a Gaussian mixture model.

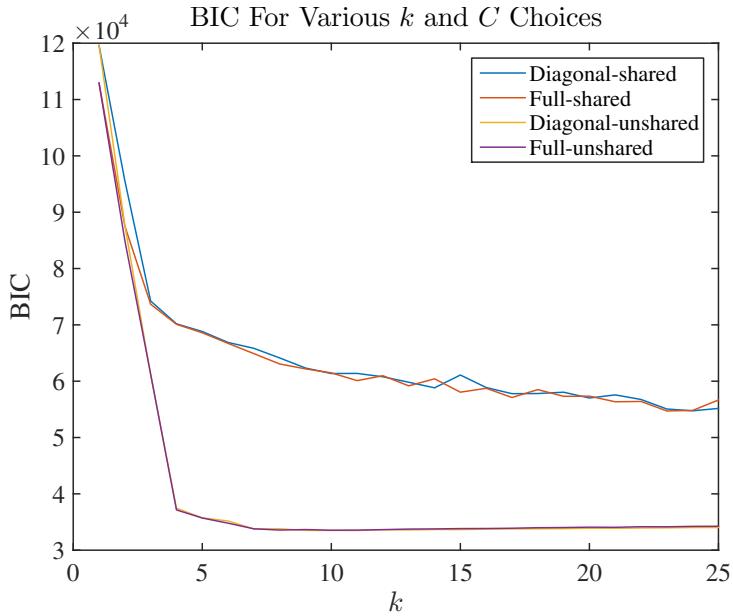


Figure 3.2: BIC Curve for fitting with the Gaussian mixture model.

3.0.2 Suitable Prior on the Mixing Probabilities

With suitable priors on the parameters, the MAP estimator can be used for model selection. In particular, [44] and [45] put the Dirichlet prior on the mixing probabilities, of the components in the Gaussian mixture model, and [46] applied the “entropic prior” on the same parameters to favor models with small entropy. All of these have in common that they used the MAP estimator to drive the mixing probabilities associated with unnecessary components toward extinction. Based on an improper Dirichlet prior, [34] suggested to use minimum message length criterion to determine the number of the components, and further proposed an efficient algorithm for learning a finite mixture from multivariate data which we have adopted for source estimation based on panning parameters. It is the model called mixture-MDL (MMDL), which is described in the following section.

3.1 Mixture MDL

If we recall the parameter vector of the Gaussian mixture model as,

$$\Theta = \{\theta_1, \theta_2, \dots, \theta_k, \alpha_1, \alpha_2, \dots, \alpha_k\}$$

Once we have estimated one source parameter i.e. $\hat{\theta}_k$ the sample size “seen” by this parameter is $N\alpha_k$ due to the mixing probability weighting [34]. The penalty term becomes dependent on not only the number of measurement vectors N , but also on the

mixing probabilities α . The fisher information for θ_k for one observation from component k becomes $N\alpha_k \mathcal{I}(\theta_k)$. Therefore, the prior on the parameters of the asymptotic MAP expression for mixtures is,

$$p(\theta_k) = \frac{k-1}{2} \ln N + \frac{N_p}{2} \sum_k \ln(n\alpha_k) \quad (3.14)$$

which is referred to as the mixture-MDL (MMDL). We will adopt this criterion from [34]. The mixture-MDL is,

$$\hat{\theta}_{k\text{MMDL}} = \arg \min_{\theta_k} \{-\ln p(\mathbf{x}|\theta_k) + p(\theta_k)\} \quad (3.15)$$

The key observation of the MMDL is that the prior $p(\theta_k)$ is not only a function on k and for a fixed k it is not a ML estimate. For fixed k , MMDL has a simple Bayesian interpretation [34]:

$$p(\{\alpha_1, \dots, \alpha_k\}) \propto \exp \left\{ -\frac{N_p}{2} \sum_{k=1}^K (\alpha_k)^{\frac{N_p}{2}} \right\} \quad (3.16)$$

Which is a Dirichlet-type improper prior, which can be used on the mixing probability in the maximum a posteriori (MAP) estimator for model selection.

The procedure is then as follows: We start with a large number of randomly initialized components and search for the MAP solution using the iterative procedure of the EM algorithm. The prior drives the irrelevant components to extinction. In this way, while searching for the MAP solution, the number of components is reduced until convergence is achieved. See [47] for details on Dirichlet type prior relation to the standard MDL. The MMDL minimization criteria is [34],

$$\hat{\theta} = \arg \min_{\theta} \mathcal{L}(\theta|\mathcal{X}) \quad (3.17)$$

with

$$\begin{aligned} \mathcal{L}(\theta|\mathbf{X}) = \arg \min_{\theta} & \left\{ -\log p(\mathbf{X}|\theta) \right. \\ & \left. + \frac{N_p}{2} \sum_{k=1}^K \log \frac{l\alpha_k}{12} + \frac{k}{2} + \log \frac{l}{12} + \frac{k(N_p+1)}{2} \right\} \end{aligned} \quad (3.18)$$

where $\alpha_k > 0$ and N_p is the number of parameters specifying each component. The MMDL mixture model is including the component-wise EM algorithm CEM² [48]. The expected number of measurement vectors generated by the c th component of the mixture is $n\alpha_k$, which is the sample size seen by the θ_k , thus the optimal (in the MDL sense) for each θ_k is $N_p/2\log(n\alpha_k)$ [34]. The MMDL promotes sparseness in the sense that it is initialized with much higher k than expected, and the EM-MMDL will then set some $\alpha_k = 0$ by killing the weakest component and then restart the CEM² algorithm [48].

3.2 Model Pruning by Component Annihilation

The following section is the proposed method for component annihilation for stereopanning estimation. In the following description, this method is described as a post-processing procedure. However, it is desirable to implement the functionality of this method as part of the likelihood-model in the segmentation algorithm which still remains unsolved. In the end of this section the model pruning will be described as a Bayesian interpretation.

3.2.1 Overfitting a Gaussian Mixture Model to Panning Parameter Space

The challenge of overfitting, a Gaussian mixture to the distribution space is ambiguous. It is the case that the MMDL will estimate a model order k that is equal to or larger than the true order. However, the GMM-model is designed to describe every single measurement vector as being part of a Gaussian distribution. The ambiguity is that the overfitting of the Gaussian mixture model can be exploited as a parameter to utilize for initialization of the EM-algorithm, which is also the case for the MMDL [34]. By starting with k , where k is much larger than the true/optimal number of mixture components, the adopted algorithm is robust with respect to initialization of the EM-algorithm. The MMDL algorithm applies component annihilation, by adopting a Dirichlet prior on the mixing probabilities [[araki_stereo](#), 34], and selects the number of components by annihilating the weakest component in the M-step of an iterative component-wise EM (CEM²) [48]. This procedure leads to a smaller model order and still describes every measurement vector as being part of a Gaussian distribution. It is important to notice that every true parameter is then described by at least one or more of the clusters.

3.2.2 Clustering and Component Metrics

After model order selection by the MMDL algorithm, each true parameter is described by one or more components. Therefore, we have applied a post-processing step to select the true number of clusters from prior spatial knowledge of the conditional distributions. In the following this procedure is described, starting with the practical view and lastly we interpret the model pruning as a Bayesian posterior.

The model pruning post-processing step selects clusters from an analysis on each cluster covariance compared to the number of estimated points embedded in each cluster. We know from [9] that due to the non-disjoint spectral overlap of sources, the variance increases in the amplitude direction. Therefore, we propose to select clusters with largest amount of estimated points, relative to the size of the respective embedding covariance and its rotational angle in the parameter space. We describe this for the k th covariance \mathbf{C}_k in the following. We define \mathbf{C}_k geometrically as an

ellipsoid by applying the singular value decomposition as,

$$\mathbf{C}_k = \mathbf{U}\Sigma\mathbf{V}^T \quad (3.19)$$

where \mathbf{U} and \mathbf{V} are orthonormal rotation matrices and the diagonal of Σ contains the principal axes a^2 and b^2 . We compute the angle θ of the principal axes a to the x-axis,

$$\theta = \tan(\mathbf{u}_2 / \mathbf{u}_1)^{-1} \quad (3.20)$$

we center \mathbf{C}_k by subtracting the mean $\boldsymbol{\mu}_k$ as $\mathbf{d} = \mathbf{C}_k - \boldsymbol{\mu}_k$. The x and y coordinates of each estimate is given as,

$$(x, y) = (\mathbf{d}_1 \cos \theta + \mathbf{d}_2 \sin \theta, -\mathbf{d}_1 \sin \theta + \mathbf{d}_2 \cos \theta) \quad (3.21)$$

We count the number of points inside cluster k . The specific point (x, y) is inside the ellipse k if,

$$\frac{x^2}{a_k^2} + \frac{y^2}{b_k^2} < 1 \quad (3.22)$$

Lastly we compute the size of \mathbf{C}_k as the determinant of \mathbf{C}_k and we compute the "shadow" of the covariance on the x-axis as $s_k = a_k \cos \theta_k + b_k \sin \theta_k$.

3.2.3 Annihilation Steps

The following component annihilation steps uses the cluster parameters from the MMDL and by comparing these to the measurements, the true clusters are selected and the rest of the measurement vectors will be removed. In this initial implementation, implemented as a post processing model pruning algorithm, we have two rules applied which is:

1. If a cluster shares an estimated point with a smaller cluster, all points that is only part of the bigger cluster is removed. The overlapping bigger clusters are referred to as sticky clusters.
2. A geometric threshold is applied based on $\det(\mathbf{C}_k)$, the number of points embedded in \mathbf{C}_k and the shadow s on the x-axis as described by (3.2.4) and (3.2.4).

The two steps described as rule 1 and rule 2, is always carried out with step 1 first. Step 1 removes every "sticky" cluster. Once the sticky clusters have been removed, each remaining clusters is measured with the ratio described by (3.2.4) and (3.2.4). An example of these two steps are shown in Figure 3.3. It is noticeable that only by removing the sticky clusters we have reduced the number of clusters from 14 to 5, when the true numbers of clusters is 4. All clusters shown as ellipsoids which are embedding estimates in blue are the sticky clusters which are now ignored. It is easy to see that one of the orange clusters has a lower density than the remaining 4 clusters.

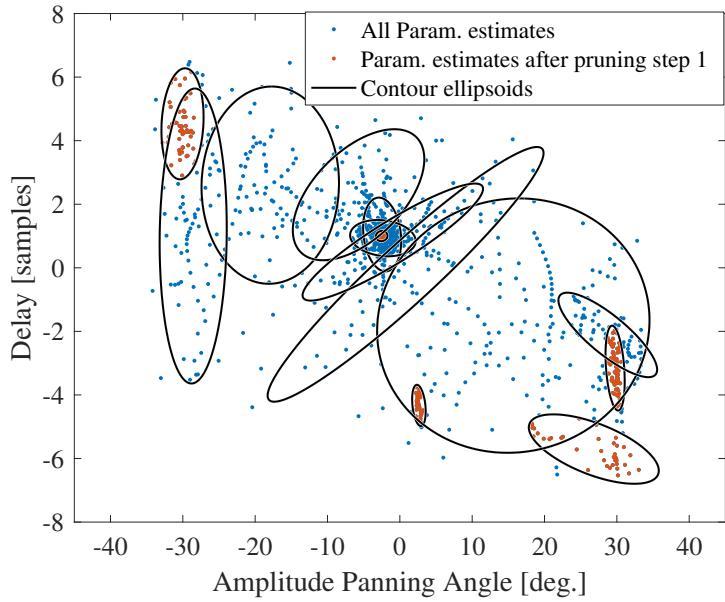


Figure 3.3: Component annihilation step 1 has been applied to the parameter estimates of a mixture of 4 sources. All ellipsoids represent a cluster. The ellipsoids containing orange dots are not sticky and are kept. All blue clusters will be ignored following step 1.

In order to remove clusters with relative low density and high correlation between parameters, we apply step 2. Another point to notice is that the one of the low density clusters also differs from the remaining in the angle of its principal component, which shows a relative higher correlation between the two given features, thus it has a greater variance in both directions since it spans a larger region, but especially the variance in the amplitude panning direction is interpreted a sign of non-disjoint orthogonality in the source mixture [9]. Figure 3.5 shows the ratio function of (3.2.4). We note that between $k = 4$ and $k = 5$ there is a ratio difference on the order of 10^{20} , which often is sufficient for a fixed thresholding. Figure 3.5 shows the parameter space after step 2 has been applied. It can be seen that the 4 correct clusters now has been found after model pruning. Lastly, the full Gaussian mixture is shown in Figure 3.6. From this figure it can be seen that the problem is not a convex one, and there is at least 6 local minimum. It is worth noticing that one of the correct clusters has a quite low peak likelihood value. Thus, it has relative low mixing probability and spans a larger region than the other correct cluster estimates. The procedure of model pruning by component annihilation has kept this specific cluster because it has a low variance in the amplitude direction when comparing to the other clusters, by using (3.2.4).

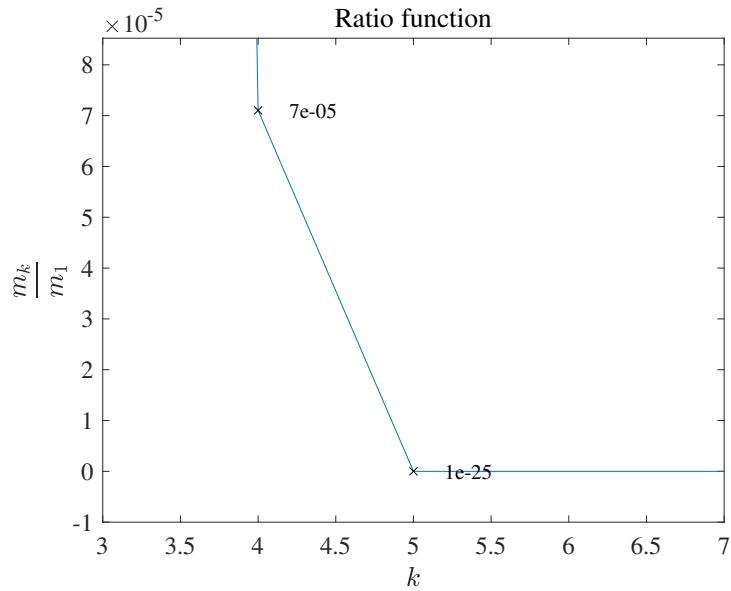


Figure 3.4: The ratio as a function of number of estimated clusters (3.2.4). It has been applied to the mixture of 4 sources.

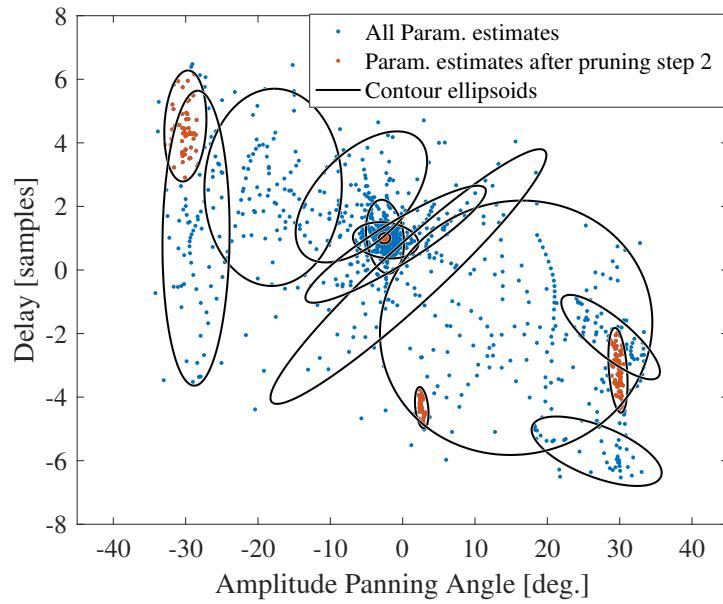


Figure 3.5: Component annihilation step 1 and step 2 has been applied to the parameter estimates of a mixture of 4 sources. All ellipsoids represent a cluster. The ellipsoids containing orange dots are kept after model pruning.

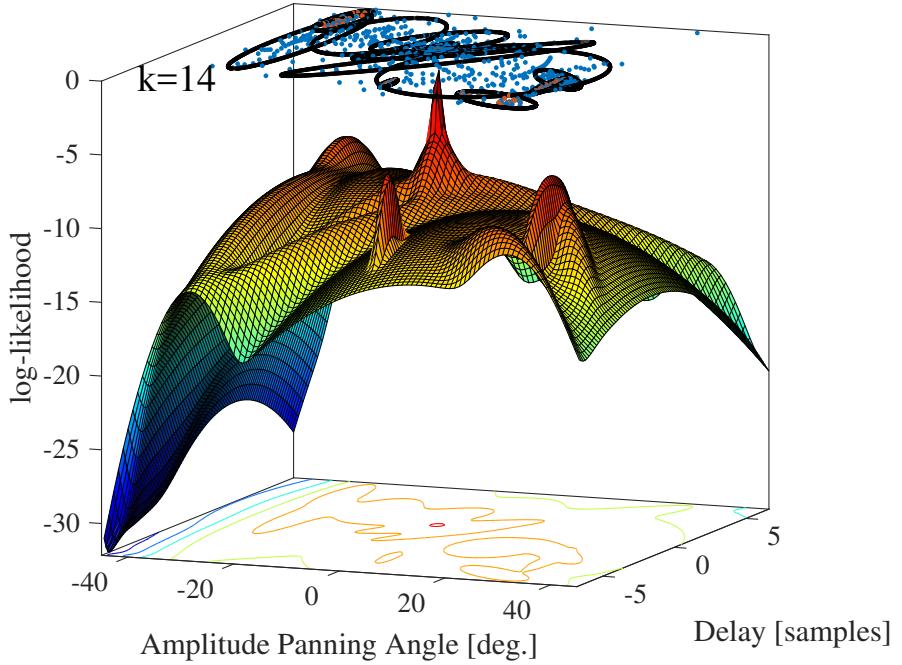


Figure 3.6: The estimated Gaussian mixture-

3.2.4 Thresholding on the Cluster Angle and size

The spectral overlap of mixture sources magnifies the variance in the amplitude panning direction [9]. In the case of disjoint orthogonality, the covariance would be very small and diagonal or have greatest variance in the delay-direction. Therefore we can apply a threshold from the rotation angle θ and the size of the region relative to the given number of estimates in the region. We define a variable $0 \leq p_k \leq 1$ which is the percentage of points that is inside the k th cluster. We notice that the mixing parameter \mathbf{ff} is proportional to $\frac{p}{\det(\mathbf{C})}$. We define a metric of peakiness m for the k th source as,

$$m_k = \frac{p_k}{\det(\mathbf{C}_k)s_k} \quad (3.23)$$

where s_k is the amplitude shadow, $s_k = a_k \cos \theta_k + b_k \sin \theta_k$. The metric m_k carries implicit information of both the size and angle of the k th cluster region, that includes a percentage of all the estimated points (after sticky clusters have been removed). From all metrics \mathbf{m} we define a threshold where m_k is relative to the smallest $m = m_1$. The metric ratio is,

$$ratio_k = \frac{m_k}{m_1} \quad (3.24)$$

Through experiments, we have found that this method of component annihilation has good performance for precisely estimating the number of sources in the mixture and the panning parameters. In the following we will interpret this with Bayesian

terminology as a posterior probability by comparing it to the K -nearest neighbour classifier.

3.2.5 Bayesian Interpretation of the Threshold

In the following we describe the model pruning method as a posterior probability, by comparing it to the K -nearest neighbour technique. To do this, we make use of Bayes theorem and apply the K -nearest neighbour method for classification to each cluster separately. Let us suppose that we have a data set with N samples, where N_k points belongs to class \mathbb{C}_k , so that $\sum_k N_k = N$. If we wish to classify a point \mathbf{x} with the K -nearest neighbour method, we draw a hypersphere that is centered on \mathbf{x} , containing K points irrespective of their class. Suppose this sphere has volume $V(\mathbf{x})$ and contains K_k points belonging to class \mathbb{C}_k . An estimate of the density associated with each class is then [37, 40],

$$\hat{p}(\mathbf{x}|\mathbb{C}_k) = \frac{K_k}{N_k V(\mathbf{x})} \quad (3.25)$$

Similarly, the unconditional density is given by,

$$\hat{p}(\mathbf{x}) = \frac{K}{N V(\mathbf{x})} \quad (3.26)$$

and the class priors are given by,

$$\hat{p}(\mathbb{C}_k) = \frac{N_k}{N} \quad (3.27)$$

We can combine these three equations using Bayes' theorem to obtain the posterior probability of class membership

$$\hat{p}(\mathbb{C}_k|\mathbf{x}) = \frac{\hat{p}(\mathbf{x}|\mathbb{C}_k)\hat{p}(\mathbb{C}_k)}{\hat{p}(\mathbf{x})} = \frac{K_k}{K} \quad (3.28)$$

Which means that we can minimize the risk of misclassification, by assigning the point \mathbf{x} to the class having the largest posterior probability corresponding to K_k/K . Such a classification can be expressed as:

$$\begin{aligned} \hat{\mathbb{C}}_k(\mathbf{x}) &= \mathbb{C}_k \quad \text{with } k = \arg \max_{i=1,\dots,K} \left\{ \hat{p}(\mathbf{x}|\mathbb{C}_k)\hat{p}(\mathbb{C}_k) \right\} \\ &= \arg \max_{i=1,\dots,K} \left\{ \frac{K_i}{N_i V(\mathbf{x})} \frac{N_i}{N_{\text{total}}} \right\} = \arg \max_{i=1,\dots,K} \left\{ \frac{K_i}{N_{\text{total}} V(\mathbf{x})} \right\} \end{aligned} \quad (3.29)$$

We can compare this expression to the threshold in (3.2.4) of the model pruning method, where $m_k = \frac{p_k}{\det(\mathbb{C}_k)s_k}$ with p_k the percentage of points inside cluster k and the volume of the hypersphere determined by $\det(\mathbb{C}_k)$ which is weighted relative to

the shadow s_k , to prefer the clusters that have lowest variance in the amplitude direction. In Figure 3.2.5 showing the parameter space it is clear that the covariance assumption described above in Section 3.2.3 holds and all 7 correct clusters have been estimated correctly. The selected clusters are shown in orange and the unselected clusters are shown in blue. The correct cluster covariances are either very small and close to diagonal or they are larger and have dominant variance mainly in the delay direction (upwards). We notice that the MMDL algorithm in this case chose a $k = 18$. Notice that the 11 "wrong clusters" are large with random covariance structure and rotation angle. Figure 3.9 shows the ratio function in (3.2.4). We note that between $k = 7$ and $k = 8$ there is a ratio difference on the order of 10^{18} , which often is sufficient for a fixed threshold.

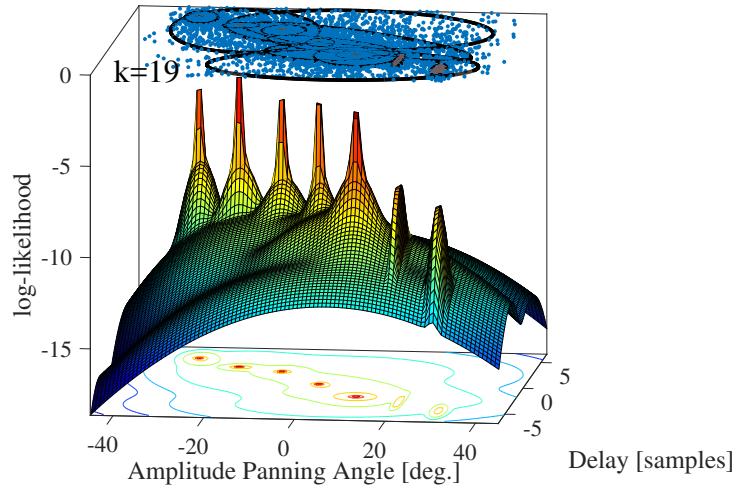


Figure 3.7: Gaussian mixture of 7 sources

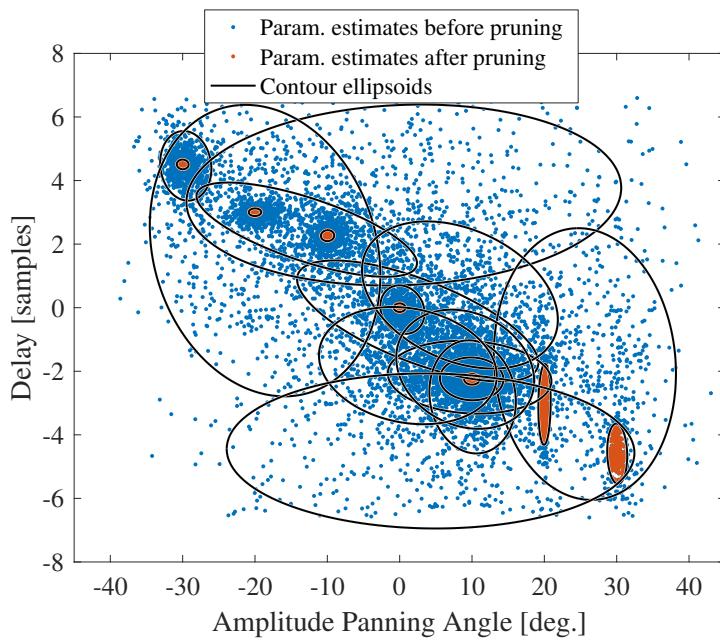


Figure 3.8: Before and after pruning.

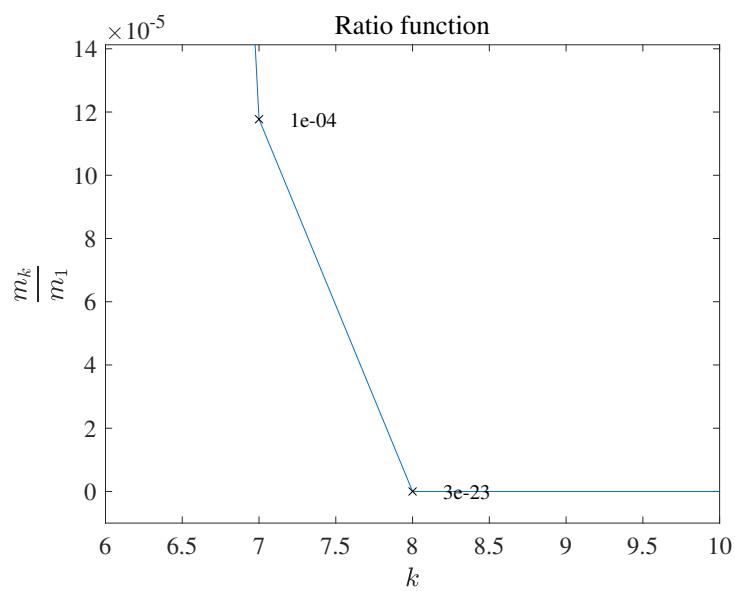


Figure 3.9: Ratio function of 7 sources shown in Figure 3.7

Chapter 4

Segmentation of the Stereophonic Signal

4.1 Signal Segmentation

The characteristics of the observed stereophonic signal are varying over time with different durations, meaning that a fixed segment length is not optimal. Using the MAP criterion. the cost associated with the different outcomes from the set of segment lengths can be compared and the optimal can be chosen as the one that minimizes (3.1). The segmentation is based on the principle in [49, 50, 12] which is outlined in Algorithm 1. A minimal segment length, N_{\min} generating a block of samples and dividing the signal into M blocks. Since this will give 2^{M-1} ways of segmenting the signal into M blocks a maximum number of blocks K_{\max} is defined to ease on computational complexity, since very high segment length is assumed to be generating noise in the distribution. The maximum number of samples in one segment is $N_{\max} = K_{\max}N_{\min}$. A dynamic programming algorithm, computes the optimal segment length k_{opt} for all blocks, $m = 1, \dots, M$, starting at $m = 1$ moving continuously to $m = M$. For every block, the cost of all new block combinations are reused from earlier blocks. When the end of the signal is reached, the optimal segmentation of the signal is found, starting with the last block and continuing through the signal to the beginning. Starting at $m = M$, setting the number of blocks in the last segment to $k_{\text{opt}}(M)$. The next segment ends at block $m = M - k_{\text{opt}}(M)$ and includes $k_{\text{opt}}(M - k_{\text{opt}}(M))$ blocks. This is continued until $m = 0$.

Algorithm 1 Segmentation

```

while  $m \times N_{\text{MIN}} \leq \text{length(signal)}$  do
    Initialize  $K = \min([m, K_{\text{max}}])$ .
    for  $k = 1$  to  $K$  do
        block of signal to use is  $m - k + 1, \dots, m$ 
        estimate  $(\hat{\gamma}, \hat{\delta})$  from (1.4.1) and (1.4.1)
        compute  $\mathcal{L}(\theta|\mathcal{Y})_{(m-k+1)m}$  from (3.1)
        if  $m = 1$  then
             $\mathcal{L}(\theta|\mathcal{Y})_{(k)} = \mathcal{L}(\theta|\mathcal{Y})_{(m-k+1)m} + \mathcal{L}(\theta|\mathcal{Y})_{1(m-k)}$ 
        else
             $\mathcal{L}(\theta|\mathcal{Y})_{(k)} = \mathcal{L}(\theta|\mathcal{Y})_{(m-k+1)m}$ 
        end if
    end for
end while
 $m = M$ 
while  $m > 0$  do
    number of blocks in segment is  $k_{\text{opt}}(m)$ 
     $m = m - k_{\text{opt}}(m)$ 
end while

```

Chapter 5

Experiments

5.1 Experiments

In the following the different proposed methods are tested through simulations on synthetic signals and real audio from the SQAM database [51]. To represent real music the synthetic signal are based on guitar recordings from which the amplitudes and phases have been extracted, by using an inharmonic approximate non-linear least square (ANLS) pitch estimator [18]. By testing the segmentation with synthetic signals, we can create a ground truth to when each source is active. The used signal were generated with 20 harmonic amplitudes and phases. The fundamental frequencies are representing notes that can be played on a guitar in the range $f_0 \in [80, 1700]\text{Hz}$, randomly applied. f_s was set to 44100 Hz. The synthetic signal has a duration of 15 seconds. White Gaussian noise has been applied to the signal with an SNR of 50 dB. The clustering and model

5.1.1 Segmentation

The segmentation is tested on the synthetic signal. The synthetic signal is consisting of two sources with a minimum active signal duration of 300 ms and note duaration as multiples of 300 ms. The signal is segmented according to the MAP MMDL criteria of (3.1), where the minimum segment length $N_{\text{Min}} = 150$ ms and the maximum number of blocks $K_{\text{Max}} = 20$ meaning that the maximum length of a segment is 3 s. A representative example of the chosen segment length as a function of time is shown in Figure 5.1.2 with white vertical lines. In the top the two active sources are shown time domain along with a black horizontal line indicating which source is active at which time (the input segment ground truth). In the background the signal frequency content is shown to give a detailed view of the signal content. Generally the chosen segments are long if the content is not changing. Each segment contains some valuable information and seperates active input segments of the two sources. The four notes played from 0 to 4 sec (with same f_0) will consistently produce two underlying

clusters, and we would expect the segments to be long but random. When the silent period starts a shorter segment length is chosen in all three silent periods starting at [3.6, 5.4, 8] sec. The note at 5 s. is clearly chosen, and the next three notes has an overlap that is segmented in to two parts, where only the second part has two active sources. The following notes after the silence at 8 s is chosen precisely in 300 ms segments each. Lastly, the long note from 12-15 s. is chosen in longer segments of 600 ms, but in order to separate the underlying clusters, the two overlapping notes in the end is chosen in segments in their respective note duration, even though they both overlap with the longer note. This indicates that the panning model, describes the signal in a precise way considering the source panning parameters, independently of the pitch information. The resulting distribution of parameter estimates $p(y|\theta)$ can be seen in Figure 5.1.2. It is clear that $k = 2$, after segmentation and thresholding.

5.1.2 Source Parameter Estimation

The estimation of source parameters are tested on the SQAM-CD signals in 100 iterations. In this part, the estimation of panning parameters are tested for optimal segments and fixed segments. Panning Parameters implicitly has the model order, as dimensionality. For each iteration, a mixture consists of minimum 2 and maximum 5 randomly picked source components, mixed according to (??). Each source signal is normalized to have an absolute maximum amplitude of 1. The duration of each mixture is varying, and is defined from the shortest audio signal in the mixture, which is minimum 16 seconds for files on the SQAM-CD. The files containing pink noise has been removed from the test set. The fixed segment size is set to 600 ms. and all mixtures are passed through to thresholding of (1.4.1). All applied panning parameters are simulation stereo, which means that for every equal amount of sources, they will be panned equally to each side. The results are shown in Table 5.1. We measure the estimated model order, number of correct parameter estimates, and the root mean square error for both source parameters. Clearly, there is an improvement by applying the optimal segmentation scheme, with a correct number of parameter estimates of 96.6% compared to 43% with fixed segments. The model order i.e. the correct estimate of the number of sources is also clearly improved by the segmentation scheme. It seems that by applying the segmentation scheme we offer some precision in the parameter estimates, both for the amplitude and delay estimates.

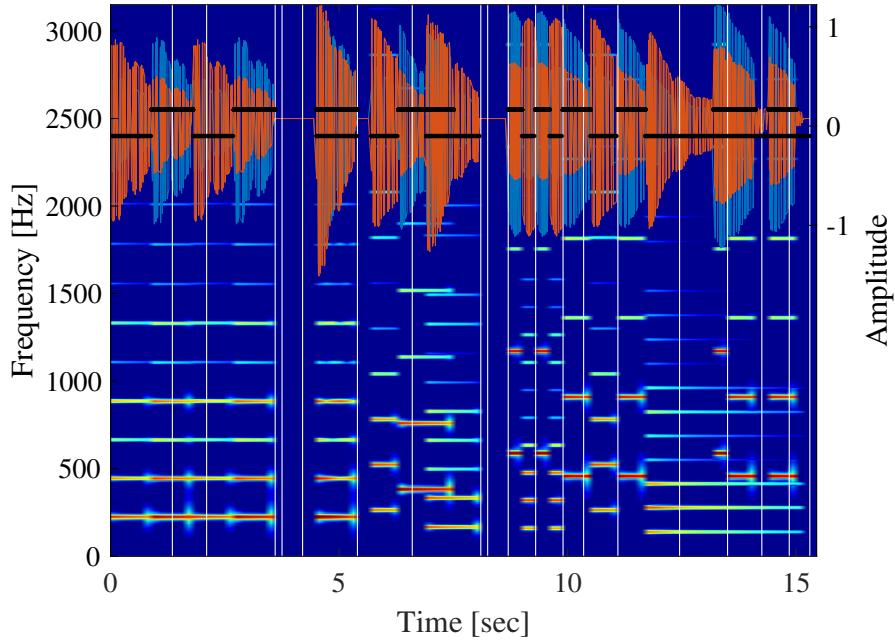


Figure 5.1: Optimal segmentation on two sources with varying pitches.

| Estimates | opt. seg. | fixed |
|---|--------------|-------|
| Correct Parameters (err. $\angle < 0.5^\circ$) | 96.6% | 84.5% |
| Correct Model Order | 94.1% | 58.4% |
| Amplitude Angle (RMSE) | 0.1° | 0.07 |
| Delay (RMSE) | 0.33 samples | 0.03 |

Table 5.1: Source parameter test results.

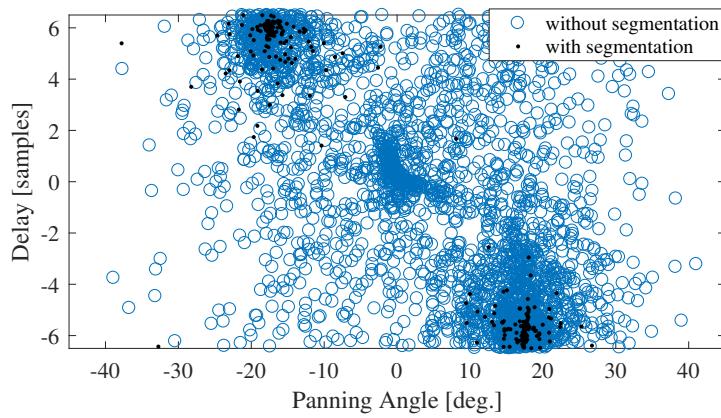


Figure 5.2: Parameter space with and without optimal segmentation and thresholding.

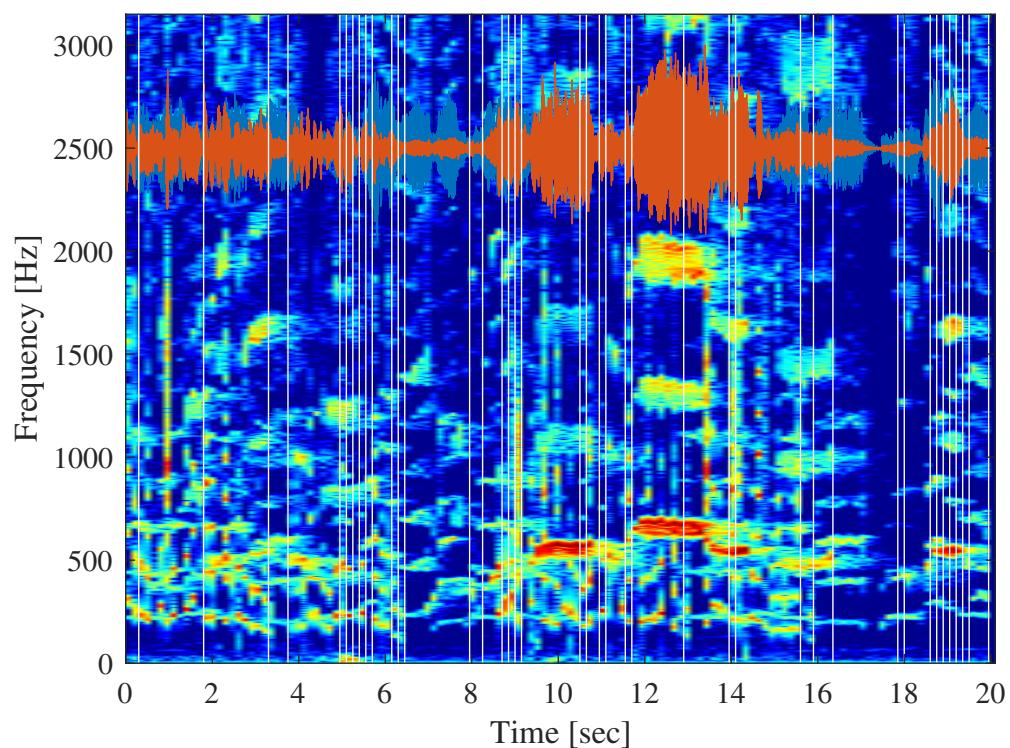


Figure 5.3: Signal segmentation on two sources from the SQAM database.

Bibliography

- [1] Martin Weiss Hansen, Jesper Rindom Jensen, and Mads Græsbøll Christensen. "Estimation of Multiple Pitches in Stereophonic Mixtures using a Codebook-based Approach". In: *I E E E International Conference on Acoustics, Speech and Signal Processing. Proceedings* (Mar. 2017). ISSN: 1520-6149.
- [2] R. J. Weiss M. I. Mandell and D. P. Ellis. "Model-Based Expectation-Maximization Source Separation and Localization". In: *IEEE Trans. Audio, Speech and Language Process.* 18.2 (2010), pp. 384–394.
- [3] J. Benesty J. R. Jensen and M. G. Christensen. "Joint filtering scheme for nonstationary noise reduction". In: *Proc. European Signal Processing Conf.* (2012), pp. 2323–2327.
- [4] E. Bryan George and Mark J. Smith. "Analysis-by-Synthesis/Overlap-Add Sinusoidal Modeling Applied to the Analysis and Synthesis of Musical Tones". In: *J. Audio Eng. Soc* 40.6 (1992), pp. 497–516. URL: <http://www.aes.org/e-lib/browse.cfm?elib=7044>.
- [5] Anssi Klapuri and Manuel Davy, eds. New York: Springer, 2006. ISBN: 0-387-30667-6.
- [6] G. Tzanetakis and P. Cook. "Musical Genre Classification of Audio Signals". In: *IEEE Transactions on Speech and Audio Processing* 10.5 (2002).
- [7] Jesper Rindom Jensen, Mads Græsbøll Christensen, and Søren Holdt Jensen. "Nonlinear Least Squares Methods for Joint DOA and Pitch Estimation". In: *IEEE Trans. Audio, Speech & Language Processing* 21.5 (2013), pp. 923–933.
- [8] Y. Huang J. Benesty J. Chen. *Microphone Array Signal Processing*. Springer, 2008.
- [9] S. Rickard and O. Yilmaz. "On the Appriximate W-Disjoint Orthogonality of Speech". In: *IEEE Acoustics, Speech, and Signal Processing* (2002).
- [10] *IEEE International Conference on Acoustics, Speech, and Signal Processing*. IEEE, 2000. URL: <http://ieeexplore.ieee.org/xpl/mostRecentIssue.jsp?punumber=6939>.

- [11] M. Vetterli P. Prandoni M. M. Goodwin. "Optimal time segmentation for signal modeling and compression". In: *IEEE Acoustics, Speech, and Signal Processing* (1997), pp. 2029–2032.
- [12] Sidsel Marie Nørholm, Jesper Rindom Jensen, and Mads Græsbøll Christensen. "Instantaneous Fundamental Frequency Estimation With Optimal Segmentation for Nonstationary Voiced Speech". In: *IEEE/ACM Trans. Audio, Speech & Language Processing* 24.12 (2016), pp. 2354–2367.
- [13] Jesper Rindom Jensen et al. "On frequency domain models for TDOA estimation". In: *2015 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2015, South Brisbane, Queensland, Australia, April 19-24, 2015*. IEEE, 2015, pp. 11–15. ISBN: 978-1-4673-6997-8. DOI: 10.1109/ICASSP.2015.7177922. URL: <http://dx.doi.org/10.1109/ICASSP.2015.7177922>.
- [14] Vincent Mohammad Tavakoli et al. "A partitioned approach to signal separation with microphone ad hoc arrays". In: *2016 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2016, Shanghai, China, March 20-25, 2016*. IEEE, 2016, pp. 3221–3225. ISBN: 978-1-4799-9988-0. DOI: 10.1109/ICASSP.2016.7472272. URL: <http://dx.doi.org/10.1109/ICASSP.2016.7472272>.
- [15] Johan Xi Zhang et al. "Joint DOA and multi-pitch estimation based on subspace techniques". In: *EURASIP J. Adv. Sig. Proc.* 2012 (2012), p. 1.
- [16] I. Jafari et al. "Time-frequency clustering with weighted and contextual information for convolutive blind source separation". In: *IEEE Workshop on Statistical Signal Processing (SSP)* (2014).
- [17] S. Rickard A. Jourjine and O. Yilmaz. "Blind separation of disjoint orthogonal signals: Demixing N sources from 2 mixtures". In: *IEEE International Conference on Acoustics, Speech, and Signal Processing*. 5 (2000), pp. 2985–2988.
- [18] Mads Græsbøll Christensen and Andreas Jakobsson. *Multi-Pitch Estimation*. Synthesis Lectures on Speech and Audio Processing. Morgan & Claypool Publishers, 2009. DOI: 10.2200/S00178ED1V01Y200903SAP005. URL: <http://dx.doi.org/10.2200/S00178ED1V01Y200903SAP005>.
- [19] T. Nilsson et al. "Multi-pitch estimation of inharmonic signals". In: *Proc. European Signal Processing Conf.* (2013), pp. 1–5.
- [20] I. Barbancho et al. "Inharmonicity-Based Method for the Automatic Generation of Guitar Tablature". In: *IEEE Trans. Audio, Speech and Language Process.* 20.6 (2012), pp. 1857–1868.
- [21] S. Rickard and O. Yilmaz. "Blind separation of speech mixtures via time-frequency masking". In: *IEEE Transactions on Signal Processing* 52.7 (2004), pp. 1830–1847.
- [22] T. Kronvall et al. "Sparse Multi-Pitch and Panning Estimation of Stereophonic Signals". In: *IEEE Trans. Audio, Speech and Language Process.* (Dec. 2016).

- [23] Jens Blauert. *The Psychophysics of Human Sound Localization*. MIT Press, 2009.
- [24] Duane H. Cooper. "Problems with Shadowless Stereo Theory: Asymptotic Spectral Status". In: *J. Audio Eng. Soc* 35.9 (1987), pp. 629–642. URL: <http://www.aes.org/e-lib/browse.cfm?elib=5188>.
- [25] Robert A. Katz. *Mastering Audio: The Art and the Science*. Butterworth-Heinemann Newton, 2009.
- [26] Alan Dower Blumlein. U.K. Patent 394 (1931); reprinted in *Stereophonic Techniques* (Audio Engineering Society, New York, 1986).
- [27] Ville Pulkki and Matti Karjalainen. "Localization of Amplitude-Panned Virtual Sources". In: *J. Audio Eng. Soc* 49.9 (2001), pp. 739–752.
- [28] Benjamin B. Bauer. "Phasor Analysis of some Stereophonic Phenomena". In: *The Journal of The Acoustical Society of America* 33.11 (1961), pp. 1536–1540.
- [29] Benjamin Bernfeld. "Attempts for Better Understanding of the Directional Stereophonic Listening Mechanism". In: *Audio Engineering Society Convention 44*. 1973. URL: <http://www.aes.org/e-lib/browse.cfm?elib=1743>.
- [30] Ville Pulkki. "Virtual Sound Source Positioning Using Vector Base Amplitude Panning". In: *J. Audio Eng. Soc* 45.6 (1997), pp. 456–466.
- [31] 24th European Signal Processing Conference, EUSIPCO 2016, Budapest, Hungary, August 29 - September 2, 2016. IEEE, 2016. ISBN: 978-0-9928-6265-7. URL: <http://ieeexplore.ieee.org/xpl/mostRecentIssue.jsp?punumber=7740646>.
- [32] Carlos Avendano and Jean-Marc Jot. "Frequency Domain Techniques for Stereo to Multichannel Upmix". In: *AES 22nd international Conference on Virtual, Synthetic and Entertainment Audio* (2002).
- [33] Stuart P. Lloyd. "Least squares quantization in pcm". In: *IEEE Transactions on Information Theory* 28 (1982), pp. 129–137.
- [34] Mário A. T. Figueiredo and Anil K. Jain. "Unsupervised Learning of Finite Mixture Models". In: *IEEE TRANSACTIONS ON PATTERN ANALYSIS AND MACHINE INTELLIGENCE* 24 (2000), pp. 381–396.
- [35] A. P. Dempster, N. M. Laird, and D. B. Rubin. "Maximum likelihood from incomplete data via the EM algorithm". In: *JOURNAL OF THE ROYAL STATISTICAL SOCIETY, SERIES B* 39.1 (1977), pp. 1–38.
- [36] G. McLachlan and T. Krishnan. *The EM Algorithm and Extensions*. Wiley Series in Probability and Statistics. Wiley, 1996. ISBN: 9780471123583. URL: <https://books.google.dk/books?id=iRSWQgAACAAJ>.
- [37] Christopher M. Bishop. *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Secaucus, NJ, USA: Springer-Verlag New York, Inc., 2006. ISBN: 0387310738.

- [38] T. Caliński and J. Harabasz. "A Dendrite Method for Cluster Analysis". In: *Communications in Statistics* 3.1 (1974), pp. 1–27. doi: 10.1080/03610927408827101.
- [39] David Arthur and Sergei Vassilvitskii. "K-means++: The Advantages of Careful Seeding". In: *Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms*. SODA '07. New Orleans, Louisiana: Society for Industrial and Applied Mathematics, 2007, pp. 1027–1035.
- [40] D. de Ridder F. van der Heijden R.P.W. Duin and D.M.J. Tax. *Classification, Parameter Estimation and State Estimation*. 1. ed. John Wiley and Sons, Ltd., 2004, pp. 17–32.
- [41] Robert E. Kass and Adrian E. Raftery. "Bayes Factors". In: *Journal of the American Statistical Association* 90.430 (1995), pp. 773–795.
- [42] Mads Græsbøll and Andreas Jakobsen. *Multi-Pitch Estimation*. 1. ed. Morgan and Claypool, 2009.
- [43] Petar M. Djuric. "Asymptotic MAP criteria for model selection". In: *IEEE Trans. Signal Processing* 46.10 (1998), pp. 2726–2735.
- [44] Dirk Ormoneit and Volker Tresp. "Averaging, maximum penalized likelihood and Bayesian estimation for improving Gaussian mixture probability density estimates". In: *IEEE Trans. Neural Networks* 9.4 (1998), pp. 639–650.
- [45] Zoran Zivkovic and Ferdinand van der Heijden. "Recursive Unsupervised Learning of Finite Mixture Models." In: *IEEE Trans. Pattern Anal. Mach. Intell.* 26.5 (2004), pp. 651–656.
- [46] Matthew Brand. "Structure Learning in Conditional Probability Models via an Entropic Prior and Parameter Extinction". In: *Neural Computation* 11.5 (1999), pp. 1155–1182.
- [47] Andrew Gelman et al. *Bayesian Data Analysis*. Chapman and Hall/CRC, 2003. ISBN: 158488388X.
- [48] Gilles Celeux et al. "A Component-Wise EM Algorithm for Mixtures". In: *Journal of Computational and Graphical Statistics* 10.4 (2001), pp. 697–712.
- [49] Paolo Prandoni, Michael Goodwin, and Martin Vetterli. "Optimal Time Segmentation For Signal Modeling And Compression". In: *In Proc. ICASSP*, pp. 2029–2032.
- [50] Paolo Prandoni and Martin Vetterli. "R/D Optimal Linear Prediction". In: *IEEE Trans. Speech and Audio Proc* 8 (2000), pp. 646–655.
- [51] European Broadcasting Union. *Sound quality assessment material recordings for subjective tests: Users handbook for the EBU SQAM CD*. Tech. Rep. EBU - TECH 3253. 2008.

Appendix

Appendix A

Estimation the Amplitude Panning Angle

It is possible to estimate the panning angle by doing a simple search within the frequency domain. This was an initial estimator of the amplitude panning parameter. What is interesting about the following algorithm is that it requires very little amount of data for a simple estimate and therefore it is interesting to use this simple algorithm to make estimates in time-pan domain as we often see the time-frequency domain referred to as the spectrogram; we refer to time-plot plot as the panogram. The algorithm is introduced in Section 1.4.2. In this section we show the panogram representations of stereo mixture of instruments. All panning parameters have been applied by using the Digital Audio Workstation (DAW) called Logic Pro. Therefore all panning parameters are applied without the use of our model, but only estimated using our model. As it is seen, the mean energy of the amplitude panning parameter estimates are at the correct location. The panning knob interface of a DAW like Logic Pro, is based on a loudspeaker aperture of 60°.

A.1 Initial test of optimal segmentation based on GMM and BIC

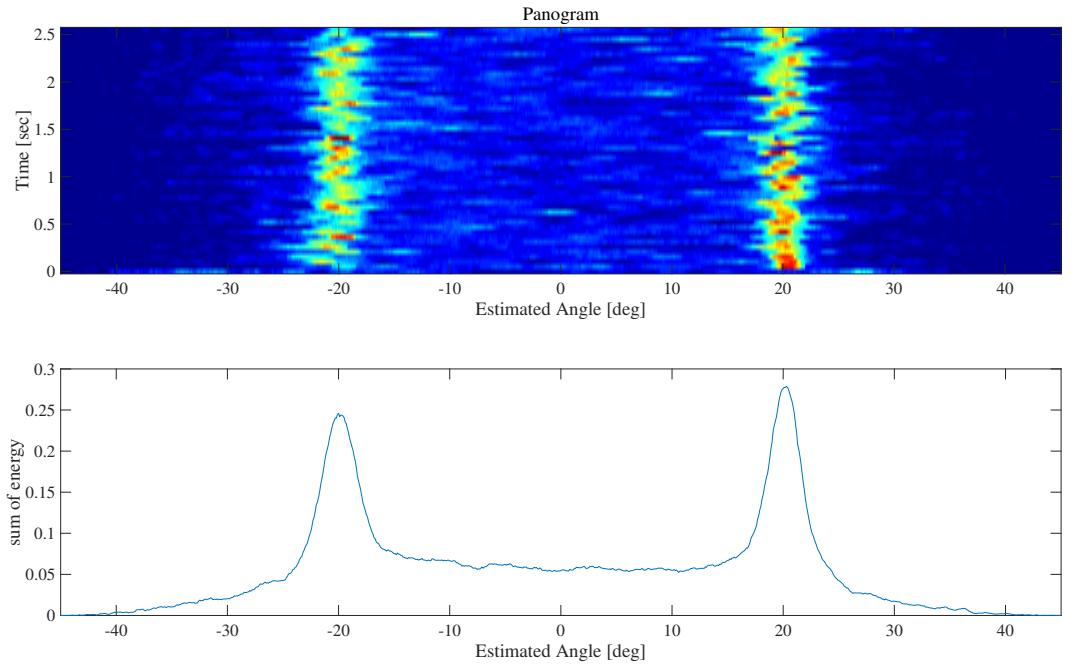


Figure A.1: Panogram of the trumpet mixture refered to in the experiments in [1].

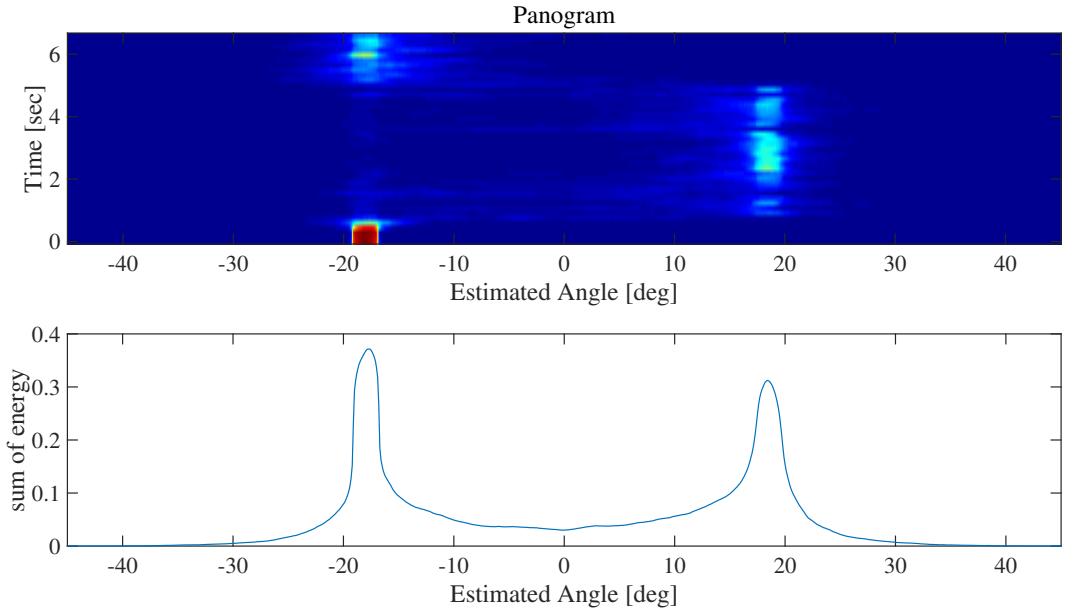


Figure A.2: Panogram of a mix of cellos. In Logic Pro these are panned on the given knob to 28. This fits with a ratio of 60/90 because the estimate is based on a loudspeaker aperture of 90°.

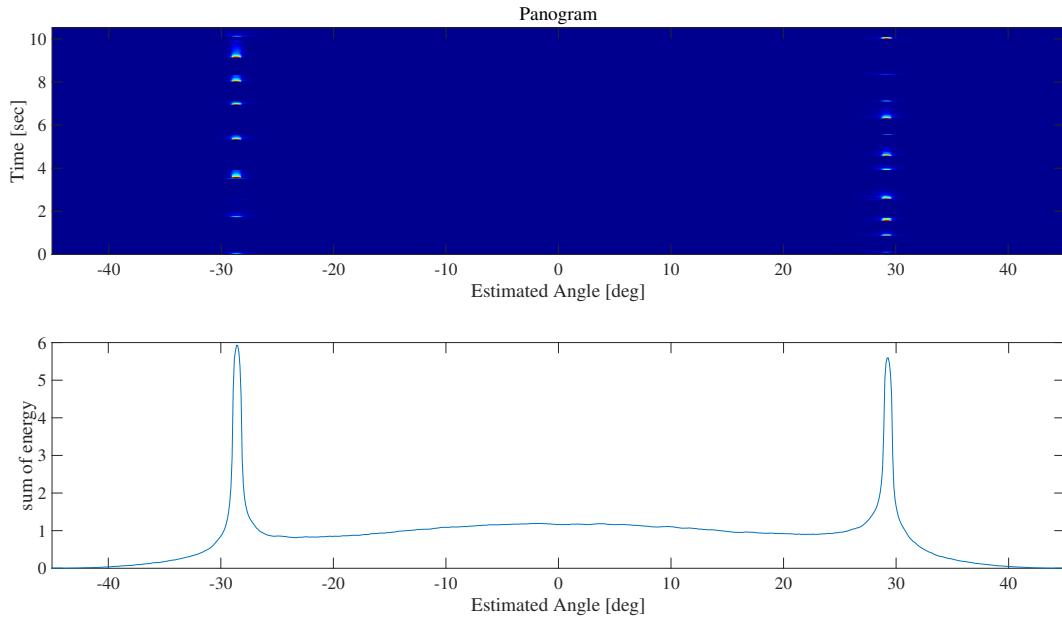


Figure A.3: Panogram of a mix of guitar plucks. In Logic Pro these are panned on the given knob to 43. This fits with a ratio of 60/90 because the estimate is based on a loudspeaker aperture of 90° .

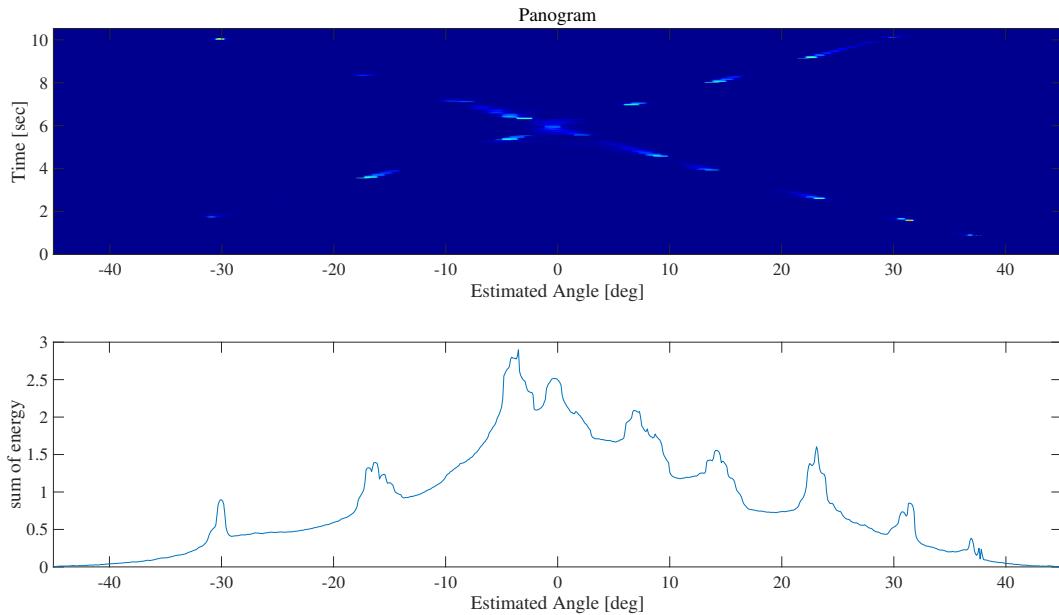


Figure A.4: Panogram of the mix of same guitar plucks from Figure A.3 only now they are changing position over time.

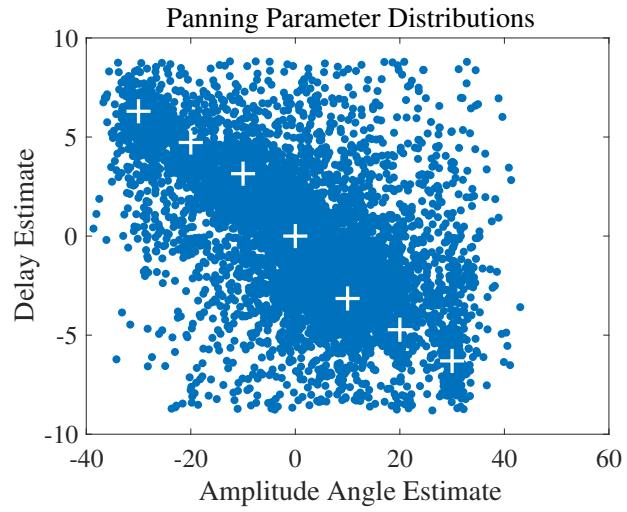


Figure A.5: Distribution of amplitude and delay ratios, based on the STFT, given a threshold for each time slot.

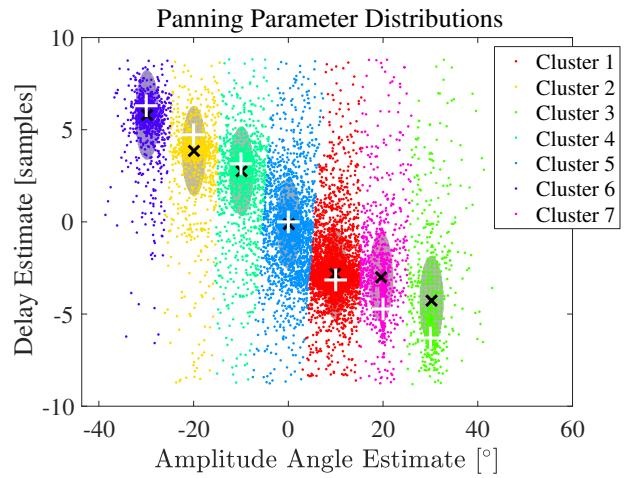


Figure A.6: Data clustered by fitting to a Gaussian Mixture Model, by using a k -means initialization.

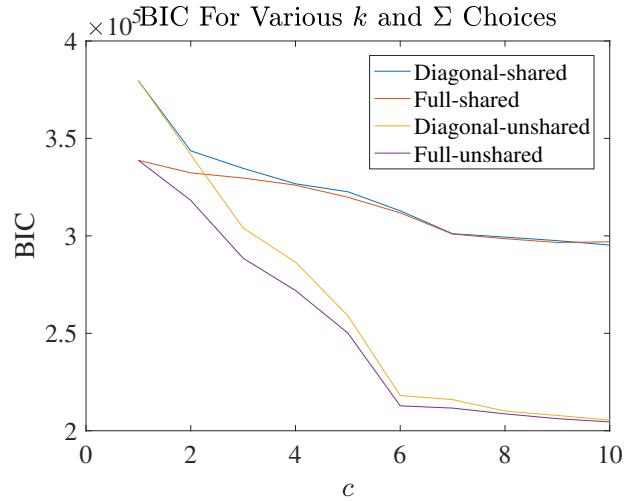


Figure A.7: Bayesian information criteria as a function of number of clusters applied to the GMM. These are shown for four different assumptions.

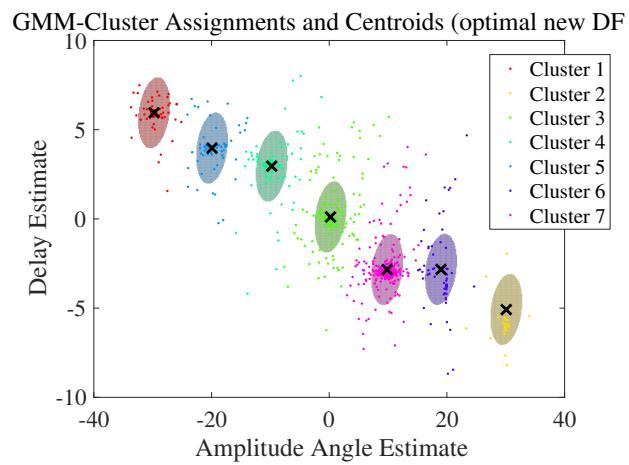


Figure A.8: Subceeding the BIC criteria for optimal segmentation. Data is clustered by fitting to a Gaussian Mixture Model, by using a k -means initialization and new DFTs for each optimal segment.

Appendix B

Initial Project Proposal

B.1 Abstract

An automatic estimation of panning parameters in stereophonic music mixtures would be of benefit to researchers in signal processing who does research within source separation and multi-pitch estimation. This paper proposes the development of such an algorithm. By applying this algorithm, researchers will be able to improve performance when designing parametric estimators within multi-pitch estimation and source separation on musical content.

B.2 Introduction

Music mixture signals is consisting of a variety of audio sources mixed together, which can be expressed as a sum of sinusoidal components. When mixing music, one of the common signal manipulations is panning the content between speaker channels. Panned content gives the listener a directional perception of the given sources in the mixture [28, 30], which separates the sources in the mixture when perceived by the listener. In signal processing, the separation of a given number of sources from a given number of mixtures, is an ongoing research area and has been researched during the recent decades. Many approaches are based on multi-pitch estimation [42] and by array processing [8], time direction of arrival (TDOA) [7, 2, 15] and blind source separation (BSS) [16, 21].

The audio source content in such research is often based on human speakers in a given acoustic environment, often anechoic, where audio material is consisting of speakers in different angles and distances to a microphone array setup. Speech content is very common in communication systems and BSS methods can be particularly well suited to speech mixtures because the time-frequency representation of speech is sparse and leads to the assumption of W-disjoint orthogonality [9, 17], however musical mixtures involves harmonic structures between a variety of sound sources, hence spectral over-

lap is a common known problem within multi-pitch estimation of musical content [42, 19, 20].

Weiss [1] proposes a novel maximum likelihood multi-pitch estimator, that utilizes the panning parameters, for improvement of multi-pitch estimates. Weiss assumes that the panning parameters are known. Kronvall et al. [22] proposes the sparse joint multi-pitch and panning estimator, as a convex optimization solution. By estimating the panning parameters Kronvall et al. argues that the estimation of panning parameters will yield a reduced computational complexity for long segments of audio where the same panning is used.

In this work we propose a panning estimator for stereophonic music that does not require an initial pitch estimate. The estimator uses an optimal signal segmentation scheme [11, 12] in time-frequency domain and is not assuming that sources are W-disjoint. Since most music is commonly available in stereo with a duration of several seconds, the panning parameters can be estimated without estimating initial pitches, with a clustering of amplitude and delay ratio estimates, of the two signal channels. This relatively simple solution utilizes the long duration of music segments, under an assumption of stationary panning parameters for all sources in the stereophonic mixture. Eventually a robust estimator for music will be proposed, that can be applied to a given multi-pitch stereophonic music problem.

B.3 State of the Art

Research in the area of panning parameter estimation is new in the area of pitch estimation and source separation for music mixtures. It has been shown that by applying panning parameters to parametric multi pitch estimation, performance can be improved [weiss2, 1], and it is implied that the panning parameters leads to faster pitch estimators [22]. Whereas panning parameters have been estimated already by kronvall et al. [22], with a search grid using an initial pitch estimate, the estimation of panning parameters preceding a pitch estimate have not been investigated in depth so far, but has been assumed to be known beforehand. Some studies [21, 17, 9] have shown that the demixing of a given number of speakers, can be estimated from 2 microphone signals, requiring some manual inspection of amplitude and delay parameter estimates, but it remains to be shown that the panning parameters can be estimated automatically on stereophonic music mixtures without an initial pitch estimate.

B.4 Thesis

Automatic panning parameter estimation for stereophonic music mixtures, will contribute to better parametric multi-pitch estimators used for source separation etc. A panning estimator for music can be achieved by applying

an optimal segmentation scheme in time-frequency domain. This will result in better automatic source separation in stereophonic music mixtures. Such an automatic method will be of benefit to researchers in signal processing.

B.5 Implementation and Methodology

The project is carried out as a scientific research within signal processing at Audio Analysis Lab at Aalborg University. Simulations and testing will be done using MATLAB. This initial proposal requires the following investigations and measures:

- The clustering of amplitude and delay estimates is found from ratios of the two channels in time-frequency domain. Such an analysis can be transferred from the TDOA community to stereophonic mixtures.
- The clustering of panning parameters is 2-dimensional in a combination of energy ratios (amplitudes) and delay ratios.
- The panning parameter estimator requires a criteria for minimizing the spread within cluster estimates.
- The clusters needs to be modelled such that it is easy to minimize the spread.
- The number of clusters is dependent on the criteria for minimizing the spread.
- The minimization of the cluster spread will be improved by applying an optimal segmentation scheme in time-frequency domain.
- The segmentation scheme is dependent on the criteria for minimizing the spread and the modelling of clusters.

The final estimator:

- Will be tested on simple instrument signals from music recordings.
- Will be measured on its computational complexity by comparing the proposed estimator in combination with a fast parametric pitch estimator to a solution based on convex optimization.
- Will be measured on estimation accuracy on a known number of sources.
- Will be measured on estimation accuracy an estimated number of sources.

B.6 Time plan

The project starts 1. Feb. 2017 (0201) and documentation needs to be handed in 22. May 2017 (0522). Main phases of the project is scheduled as:

- 0201: Start writing the proposal.
- 0203: First draft of proposal to supervisor.
- 0208: Final draft of proposal to supervisor.
- : End of literature review
- : Finish definition of modules to implement,
- : Finish documentation of overall system.
- : ...
- 0401: Prepare the finishing of all tests.
- 0406: Writing conference paper.
- 0420: Final paper hand-in for conference.
- 0522: Hand in final report.

B.7 Notes for article to WASPAA

Up to this point a solution several solutions has been investigated for solving the panning parameter estimates. A solution based on optimal segmentation has been found good. This solution will be described in the following. Firstly an outline of the subjects will be given. The proposed system contains the modules/blocks:

- Model and Selection
 - Panning model based on music mixtures in stereo format
 1. Amplitude panning
 2. Delay Panning, with reference to TDOA systems.
 3. Frequency analysis with DFT or STFT
 4. Compare to Scott Rickard as related work in a context of remixing with spatial information and 2 microphones.
 - * compare to assumptions of disjoint orthogonality versus music tracks with long duration.
 - * note that a manual inspection of data is not desired in this work.

- 5. Threshold selection of data (computational lower cost and precision is higher).
- 6. input stereo track / output clusters in 2D panning space.
- 7. The remaining clusters are quite noisy and will be processed with an optimal segmentation algorithm.
- Segmentation of input signal.
- Clustering
- Experiments

B.7.1 Algorithm description

Since panning parameters are stationary for longer periods of time.

Signal Model

Consider a dual-channel music mixture consisting of k unknown sources embedded in noise at time instant n . The data in the m^{th} channel is represented as

$$\mathbf{x}_m = [x_m(0) \quad x_m(1) \quad \cdots \quad x_m(N-1)] \quad (\text{B.1})$$

The signals captured by channel m , relating to the k^{th} source are attenuated by α_{km} and delayed by δ_{km} depending on their perceptual virtual positioning, given by the panning parameters. This stereo mixture model is generally consisting of a linear superposition of K sources in additive noise $e_m(n)$.

$$x_m(n) = \sum_{k=1}^K \alpha_k s_k(n - \delta_k) + e_m(n) \quad (\text{B.2})$$

where $m = [1, 2]$ for stereo mixtures, δ_k is the relative delay of source k between the channels and α_k is the relative attenuation factor corresponding to the ratio of attenuation of source k between the channels and $e_m(n)$ is independent white Gaussian noise. (DESCRIBE MAXIMAL DELAY and narrow band assumption) (along with panning theory limits). We aim to estimate the original sources and panning parameters without estimating pitch, given only the stereo mixtures.

Method Assumptions

The music source instruments within mixtures are often harmonically related both in pitch and partial amplitudes. Hence, we do not assume that the sources are Windowed-disjoint orthogonal as in [scott rickard and others].

$$X_i(\omega) X_j(\omega) \approx 0 \quad \forall \omega, i \neq j \quad (\text{B.3})$$

where $X_i(\omega)$ is the discrete Fourier Transform (DFT) of the i^{th} source $s_i(n)$ defined as

$$X_i(\omega) = \sum_{n=1}^N W(n)s_i(n)e^{-j\frac{2\pi}{N}k}, \quad k = 0, \dots, N-1 \quad (\text{B.4})$$

Where $W(n)$ is the windowing function. Music instrument mixtures are assumed to have stationary panning parameters and varying pitches, over several minutes of duration. Based on this assumption we exploit the duration by clustering of spectral features which reflect the panning parameters of the mixtures. In order extract useful information, an optimal segmentation scheme is applied to the mixtures based on the MMDL criteria by [Mario].

Signal Model

Consider an M -channel music mixture consisting of K unknown sources embedded in noise at time instant n . The data in the m^{th} channel is represented as $\mathbf{x}_m(n) \in \mathbb{R}^N$,

$$\mathbf{x}_m(n) = [x_m(n) \quad x_m(n+1) \quad \dots \quad x_m(n+N-1)]^T \quad (\text{B.5})$$

for $m = 1, \dots, M$. The signal mixture is modelled as a linear superposition of K sources \mathbf{s}_k in additive noise $\mathbf{e}_{m,k}$. The signals captured by channel m , relating to the k^{th} source are attenuated by gain coefficient $g_{m,k}$ and delayed by $\tau_{m,k}$ depending on their perceptual virtual positioning, given by the panning parameters,

$$\mathbf{x}_m(n) = \sum_{k=1}^K g_{m,k} \mathbf{s}_k(n - f_s \tau_{m,k}) + \mathbf{e}_{m,k}(n) \quad (\text{B.6})$$

where $g_{m,k}$ and $\tau_{m,k}$ are the attenuation and delay applied to the signal, respectively and f_s is the sampling frequency. Considering stereophonic mixtures with $M = 2$ in a stereo loudspeaker setup, the tangent law [Bernfeld,VBAP] describes an amplitude panning angle applied to the K sources, with a linear relation to the gain coefficients g_1 and g_2 . The trigonometric functions are used for the panning attenuation because they induce a constant perceived distance between listener and the virtual source, described by $1 = \cos^2 + \sin^2$. The gains are

$$g_m = \begin{cases} \cos \theta_k, & \text{for } m = 1 \\ \sin \theta_k, & \text{for } m = 2 \end{cases} \quad (\text{B.7})$$

where $\theta = \phi + \phi_0$ is a sum of the perceived angle ϕ and the speaker base angles $\pm\phi_0 = 45^\circ$. Under the conditions $0^\circ < \phi_0 < 90^\circ$, $-\phi_0 \leq \phi \leq \phi_0$ and $g_1, g_2 \in [0, 1]$ the gains can be expressed as,

$$\mathbf{g} = \mathbf{p}_k \mathbf{L}^{-1} \quad (\text{B.8})$$

where the unit-vector \mathbf{p} points towards the virtual source with \mathbf{L} as a unitary loudspeaker base matrix. The amplitude panning angle $\hat{\theta}_k$ of the k^{th} source is found as,

$$\hat{\theta}_k = \arctan \frac{p_k(1)}{p_k(2)} \quad (\text{B.9})$$

It is relevant to note here that panning parameters are applied in a post-processing step of a music production, and delays can be added to enhance the spatial perception [Blauert]. Since all mixing parameters are applied as a post-processing procedure, only the direct path is investigated, and it leads to simpler analysis if the attenuation and delay parameter are modelled as ratios between the two channels. The two channel mixtures are,

$$\mathbf{x}_1(n) = \sum_{k=1}^K \mathbf{s}_k(n) + \mathbf{e}_k(n) \quad (\text{B.10})$$

$$\mathbf{x}_2(n) = \sum_{k=1}^K \alpha_k \mathbf{s}_k(n - \delta_k) + \mathbf{e}_k(n) \quad (\text{B.11})$$

where $\delta_k = f_s \tau_k$ is the relative delay of source k between the channels and α_k is the relative attenuation factor corresponding to the ratio of attenuation of source k between the channels and $\mathbf{e}_k(n)$ is independent white Gaussian noise. We can express the model of eq. B.7.1 and B.7.1 in frequency domain as,

$$X_1(\omega) = \sum_{n=1}^N \mathbf{s}_k(n) e^{-j\omega l}, \quad (\text{B.12})$$

$$X_2(\omega) = \sum_{n=1}^N \alpha_k \mathbf{s}_k(n) e^{-j\omega l \delta_k}, \quad (\text{B.13})$$

with $\omega = \frac{2\pi}{N}$ and $l = 0, \dots, N-1$. The computation of relative panning ratios α_k and δ_k can be expressed as,

$$(\alpha_k, \delta_k) = \left(\left| \frac{X_2(\omega)}{X_1(\omega)} \right|, \frac{1}{\omega} \angle \frac{X_2(\omega)}{X_1(\omega)} \right) \quad (\text{B.14})$$

to avoid phase ambiguity we must ensure that $|\omega_{\max} \delta_{\max}| < \pi$. The assumption that the sources are approximately Windowed-disjoint orthogonal [scott rickard] is satisfied by applying a threshold in frequency domain.

$$X_i(\omega) X_j(\omega) \approx 0 \quad \forall \omega, i \neq j \quad (\text{B.15})$$

where $X_i(\omega)$ is the discrete Fourier Transform (DFT) of the i^{th} source $s_i(n)$ defined as

$$X_i(\omega) = \sum_{n=1}^N W(n) s_i(n) e^{-j \frac{2\pi}{N} k}, \quad k = 0, \dots, N-1 \quad (\text{B.16})$$

Where $W(n)$ is the windowing function. Music instrument mixtures are assumed to have stationary panning parameters and varying pitches, over several minutes of duration. Based on this assumption we exploit the duration by clustering of spectral features which reflect the panning parameters of the mixtures. In order extract useful information, an optimal segmentation scheme is applied to the mixtures based on the MMDL criteria by [Mario].