

Audio Source Separation via the Grouping of Partials in the Sum-of-Sinusoids Model

Nicholas Esterer



Department of Music Research
Schulich School of Music
McGill University
Montreal, Canada

December 2016

A thesis submitted to McGill University in partial fulfilment of the requirements of the
degree of Master of Arts.

© 2016 Nicholas Esterer

Abstract

This thesis explores a strategy of audio source separation that relies on a classification of partials via their measured modulations. Perceptual studies have shown that signals with common amplitude- and frequency-modulation are heard as coming from the same source [35], [34]. To include these modulations in a sum-of-sinusoids model a nonlinear polynomial phase function is adopted whose parameters are estimated using the Distribution Derivative Method (DDM) [2]. For better estimation accuracy, a window is designed that has lower side-lobes than the canonical Hann window but that is also once-differentiable — a requirement of the DDM. These estimated parameters are used in a new partial tracking algorithm based on linear programming. The resulting partials are classified using the clustering technique of Gaussian mixture models [14] on frequency- and amplitude-modulation data. Principal components analysis is used to emphasize the parameter on which it would be best to perform classification. Once the partials have been classified into sources, the sources are synthesized from the measured sinusoidal parameters.

The additional information provided by the DDM (namely the frequency and log-amplitude slope) is incorporated into interpolating polynomials for the phase and amplitude of sinusoids. The quality of different model-orders for these polynomials is assessed on synthetic signals. The source separation system is evaluated on both simulated data and on a mixture of real recordings of percussive and plucked string instruments. In this latter case, it is shown that using amplitude-modulation is a good criterion for separation when there is little frequency-modulation.

Résumé

Dans ce mémoire nous explorons une stratégie de séparation de sources sonores s'appuyant sur une classification de partiels selon leurs modulations observées. Des études perceptives ont montré que des signaux dont les modulations d'amplitude et de fréquence sont communes sont perçus comme provenant d'une même source [35], [34]. Afin d'inclure ces modulations dans le modèle de synthèse additif, la phase est représentée par une fonction polynomiale non-linéaire dont les paramètres sont estimés par la méthode de distribution dérivée (Distribution Derivative Method - DDM) [2]. Afin d'améliorer la qualité d'estimation, nous avons conçu une fenêtre dont la résolution dynamique est meilleure que celle de la fenêtre canonique de Hann, tout en étant dérivable sur tout son domaine, propriété requise par la technique DDM. Les paramètres ainsi estimés sont utilisés par un nouvel algorithme de suivi de partiels fondé sur le principe d'optimisation par programmation linéaire. Les partiels trouvés sont alors classifiés en différentes sources par une technique de mélange gaussiens appliquée aux paramètres de modulation de fréquence et d'amplitude. Au préalable une analyse en composantes principales (PCA) est utilisée afin de faire ressortir les paramètres les mieux appropriés pour la classification. Une fois les partiels regroupés et classifiés en sources, celles-ci sont synthétisées en fonctions des paramètres associés aux trajets de partiels.

Plus précisément, les informations additionnelles fournies par la DDM (dérivée de la fréquence et de la log-amplitude) sont prises en compte selon plusieurs stratégies impliquant des polynômes de reconstruction de phase et d'amplitude d'ordres différents. La qualité des signaux re-synthétisés est alors évaluée pour chacune de ces stratégies.

Enfin, ce système de séparation de sources est testé sur un mélange de signaux synthétiques puis sur un mélange de signaux instrumentaux réels de percussion et de corde pincée. Dans ce dernier cas, nous montrons que la prise en compte de la modulation d'amplitude aide à la classification en l'absence de modulation en fréquence.

Acknowledgments

My time at McGill University has been a long one, so I have many people to thank.

First and foremost, I would like to thank Professor Philippe Depalle for his mentorship, support, and patience during the writing of this thesis, especially during the final months. He has a tremendous talent for understanding and organizing thought. I would not have been able to produce a document of this quality without his help.

I am grateful to Professor Gary Scavone both for the teaching assistance opportunities he granted me and for inspiring me in those first courses I took with him to further pursue Music Technology. I thank Professor Ichiro Fujinaga and Professor Stephen McAdams for the opportunities they provided to work on Music Technology projects, furthering my training in information technology and signal processing.

Much of my interest in music technology and signal processing was awoken by my colleague and friend, Scott Monk, to whom I express much gratitude.

Finally I would like to thank my parents, Denise and Rob Esterer, for their encouragement in all my endeavours.

Thank you all.

Contents

1	Introduction	1
1.1	Motivation	1
1.2	Applications	2
1.3	Organization	2
1.4	Notation	3
1.4.1	Vectors and matrices	3
1.4.2	Operators	3
1.4.3	Random variables	5
1.4.4	Complex numbers	5
1.4.5	Logarithms	6
1.4.6	Ultimate values	6
2	Methodology	7
2.1	Additive sinusoidal model	7
2.2	Multiple fundamental frequency estimation	7
2.3	Principal latent component analysis (PLCA) and non-negative matrix factorizations (NMF)	8
2.3.1	Motivation	8
2.3.2	Approaches	9
2.3.3	PLCA	9
2.3.4	NMF	10
2.3.5	Synthesis of factorized spectrograms	12
2.3.6	Extensions and shortcomings	12
2.4	An approach using amplitude- and frequency-modulation	13

3	Signal Modeling	15
3.1	Introduction	15
3.2	Time-frequency representations	15
3.3	Polynomial phase models	18
3.3.1	Sinusoidal Representations	18
3.4	Polynomial phase parameter estimation	19
3.5	Choosing atom ψ	21
3.5.1	The Hann window	22
3.5.2	Continuous Blackman-Harris windows	22
3.6	Conclusion	26
4	Partial Tracking	29
4.1	A greedy method	29
4.2	An optimal method	31
4.2.1	L shortest paths via linear programming	33
4.2.2	Complexity	36
4.3	Partial paths on an example signal	37
4.4	Conclusion	41
5	The extended phase and amplitude model	43
5.1	Partial synthesis	43
5.2	The interpolating analysis-synthesis system	45
5.3	$\mathcal{S}_{1,3}$: the McAulay-Quatieri method	45
5.3.1	Analysis: linear phase and constant amplitude	45
5.3.2	Synthesis: cubic phase and linear log-amplitude	46
5.4	$\mathcal{S}_{2,3}$ and $\mathcal{S}_{2,5}$: the DDM-based methods	47
5.4.1	Analysis: quadratic phase and log-amplitude	47
5.4.2	Synthesis: cubic order ($\mathcal{S}_{2,3}$)	48
5.4.3	Synthesis: quintic order ($\mathcal{S}_{2,5}$)	51
5.5	Evaluation	53
5.5.1	Evaluation on a sinusoid of cubic phase and quartic log-amplitude	53
5.5.2	Evaluation on sinusoid of exponential phase	57
5.6	Conclusion	60

5.6.1	Polynomial phase and log-amplitude function	60
5.6.2	Exponential phase function	61
6	Experiment: Partial grouping by amplitude- and frequency-modulation	63
6.1	Introduction	63
6.2	Methodology	63
6.3	Evaluation	64
6.4	Synthesis	64
6.5	Analysis	68
6.5.1	The amplitude signal	70
6.6	Computation of principal components	73
6.7	Preparing data for clustering	75
6.8	Clustering	76
6.9	Results	76
6.10	Conclusion	81
7	Experiment: Separation of two sources using partial decay rate	85
7.1	Introduction	85
7.2	Description of problem	85
7.3	Motivation	89
7.4	Classification	91
7.5	Synthesis	94
7.6	Conclusion	94
8	Conclusion	97
8.1	Results	97
8.1.1	Quality of polynomial models for analysis and synthesis	97
8.1.2	The use of amplitude- and frequency-modulation in audio source separation	98
8.2	Contributions	98
8.2.1	Design of continuous windows with lower side-lobes	98
8.2.2	Partial tracking using linear programming	99
8.3	Future extensions	99
8.3.1	Continuous analysis windows	99

8.3.2	Partial tracking in an optimization framework	99
8.3.3	Signal modeling with nonlinear amplitude and phase polynomials .	100
Appendices		101
A Principal components analysis (PCA)		103
A.1	Motivation	103
A.2	Computation of principal components	103
B Gaussian mixture models (GMM)		105
C The normal distribution		107
References		109

List of Figures

3.1	Minimum 4-Term Blackman-Harris: Time-domain.	23
3.2	Minimum 4-Term Blackman-Harris: Frequency-domain.	24
3.3	Comparison of endpoints of window in time-domain.	25
3.4	\mathcal{C}^1 4-Term Blackman-Harris: Time-domain.	26
3.5	\mathcal{C}^1 4-Term Blackman-Harris: Frequency-domain.	27
4.1	Possible graph connections.	32
4.2	Two shortest paths using the greedy method.	33
4.3	Two shortest paths using the LP method.	35
4.4	Compare greedy and LP partial tracking on chirps in noise, SNR 20 dB. . .	38
4.5	Compare greedy and LP partial tracking on chirps in noise, SNR 15 dB. . .	39
4.6	Compare greedy and LP partial tracking on chirps in noise, SNR 10 dB. . .	40
5.1	Spectrograms of $\mathcal{T}_{3,4}$, $\mathcal{S}_{1,3}$, $\mathcal{S}_{2,3}$ and $\mathcal{S}_{2,5}$	53
5.2	$\mathcal{T}_{3,4}$ vs. $\mathcal{S}_{1,3}$, $\mathcal{S}_{2,3}$ and $\mathcal{S}_{2,5}$: Upper error bound.	54
5.3	$\mathcal{T}_{3,4}$, $\mathcal{S}_{1,3}$, $\mathcal{S}_{2,3}$ and $\mathcal{S}_{2,5}$: Log-amplitude functions.	55
5.4	$\mathcal{T}_{3,4}$ vs. $\mathcal{S}_{1,3}$, $\mathcal{S}_{2,3}$ and $\mathcal{S}_{2,5}$: Log-amplitude function error.	56
5.5	$\mathcal{T}_{3,4}$ vs. $\mathcal{S}_{1,3}$, $\mathcal{S}_{2,3}$ and $\mathcal{S}_{2,5}$: Phase function error.	57
5.6	Spectrograms of $\mathcal{T}_{\text{exp.}}$, $\mathcal{S}_{1,3}$, $\mathcal{S}_{2,3}$ and $\mathcal{S}_{2,5}$	58
5.7	$\mathcal{T}_{\text{exp.}}$ vs. $\mathcal{S}_{1,3}$, $\mathcal{S}_{2,3}$ and $\mathcal{S}_{2,5}$: Upper error bound.	59
5.8	$\mathcal{S}_{2,5}$ evaluation error bound.	61
6.1	Original data-points.	66
6.2	Amplitude function for each source (true).	67
6.3	Original and spurious data-points.	68

6.4	Principal components and their classification.	72
6.5	Source 1 (estimated).	74
6.6	Source 2 (estimated).	75
6.7	Source 1 (estimated) after smooth amplitude path search.	77
6.8	Source 2 (estimated) after smooth amplitude path search.	78
6.9	Source 1 (estimated) after smooth frequency path search.	79
6.10	Source 2 (estimated) after smooth frequency path search.	80
6.11	Smoothed paths.	81
6.12	Partitioning example.	82
7.1	Spectrogram of acoustic guitar.	86
7.2	Spectrogram of xylophone.	86
7.3	Spectrogram of mixture.	87
7.4	Spectrogram of mixture and partial trajectories.	88
7.5	Partial trajectories.	88
7.6	Path cost vs. length and thresholding boundary.	89
7.7	Log-partial-length vs. frequency: principal components.	90
7.8	Estimated memberships.	92
7.9	Spectrogram of source separated acoustic guitar.	93
7.10	Spectrogram of source separated xylophone.	94

Chapter 1

Introduction

1.1 Motivation

This thesis was written in the information age where digital signals can be produced easily and are produced in great volumes. As stated in the introduction of [21], signals are the means by which information is transmitted. In the past, producing a digital signal was costly and required specialized equipment, motivating the user of the equipment to carefully plan the process that was to be documented by encoding the measurements into the signal. Now these tools are widely available and accessible to everyone, increasing both the variance in quality, but also the potency of information. For this reason, new signal processing techniques require ways of removing extraneous information, that is, data that have meaning and structure, but that are not pertinent to the information of interest. In this thesis, a producer of information is a *source*, and we have many sources transmitting information in a single signal.

Source separation is a difficult problem because it involves simultaneously estimating characteristics of the sources while separating them: for improved estimation, interference from other sources should be minimized; in order to remove interfering sources, their characteristics must be known. For the estimation problem, we must resort to using prior information: we assume we know the structure of the sources and can quantify in some way their characteristics, penalizing characterizations estimated by our system that do not match presumptions. For the separation problem, we often must resort to suboptimal solutions. These solutions may be adequate to aid a human in manual refinement of the sources, or serve as input to another technique.

1.2 Applications

Applications of source separation exist in a variety of disciplines and entire conferences are dedicated to the subject (see, for example, [64]). Here we will only consider applications pertinent to audio, acoustics and music. One of the most popular applications of audio source separation is for automated music transcription (see e.g., [1]). Having access to both a representation of the musical score and its constituent sounds would be very convenient for composers and sound engineers. Those interested in isolating individual voices in a recording of multiple speakers (for honourable or dishonourable purposes) would benefit from audio source separation — [42] discusses a strategy using characterizations of language to aid in the separation. There are no doubt many more applications of source separation in the field of audio signal processing.

1.3 Organization

The general strategy explored in this thesis is an iterative process with four distinct steps. (1) a model is chosen of the signals of interest. (2) realisations of this model are identified in the measurement signal. (3) once these have been identified, the parameters of the realisations are estimated. (4) the estimations are used to classify these realisations as one of a smaller set of higher-level objects. The structure of these objects is used to inform the selection of the new model, whose parameters are then estimated, etc.; the process can be repeated to build up successively higher-level models.

This thesis is structured according to these steps and is as follows. Chapter 2 discusses previous approaches to audio source separation and introduces the method adopted here. The following, Chapter 3, describes the signal model chosen to describe musical signals — the additive sinusoidal model with polynomial log-amplitude and phase — and describes a technique for estimating the parameters of models of arbitrary order. We introduce a new analysis window to have more control over the estimation process. Techniques to identify these models in signals are discussed in Chapter 4, where a classical technique of partial tracking is compared to a new linear programming formulation. Chapter 5 shows how the separated sources can be synthesized from the estimated model parameters using the additional information provided by the higher-order polynomial model. Finally, two experiments are carried out that demonstrate the use of amplitude- and frequency-

modulation to classify sources. The first, in Chapter 6 is on synthetic data and the second in Chapter 7 on a recording of percussion and plucked string instruments. It is in these latter chapters where classification is performed and its adaptation to the particular audio source separation problem is discussed. The Appendices A through C explain elements of these classification techniques.

1.4 Notation

1.4.1 Vectors and matrices

While scalars are typeset normally — x is an example of a scalar — vectors and matrices are typeset in a boldface font, with matrices written with a capital letter, e.g., \mathbf{x} is a vector and \mathbf{X} a matrix. If a number is written instead of a symbol, we mean a vector whose entries are all that number, e.g., $\mathbf{1}$ is the vector of all 1s, $\mathbf{0}$ the vector of all 0s. The i th entry of a vector \mathbf{x} is written x_i and the entry in the i th row and j th column of a matrix \mathbf{X} is written $X_{i,j}$. Both are scalars and therefore typeset normally. Sometimes we might find it convenient to extract a column vector or row vector from the matrix \mathbf{X} . We write $\mathbf{x}_{i,:}$ to extract all columns from the i th row and $\mathbf{x}_{:,j}$ to extract all rows from the j th column. These are the i th row vector and j th column vector respectively. The orientation of a vector will be clear from the context, but in general \mathbf{x} is a column vector while $\mathbf{y}_{i,:}$ and \mathbf{x}^T are row vectors.

1.4.2 Operators

Inner product

We will be dealing with objects in vector spaces. The operator $\langle x, y \rangle$ takes two objects in a vector space V , $x, y \in V$ and maps them to an element $k \in \mathbb{K}$ of a field \mathbb{K} . For this thesis, the field will always be the field of real numbers \mathbb{R} or complex numbers \mathbb{C} . The vector space can be the set of vectors of N elements in \mathbb{R}^N or \mathbb{C}^N , in which case the inner product is defined, for $\mathbf{x}, \mathbf{y} \in \mathbb{K}^N$, $k \in \mathbb{K}$

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y} = k$$

The inner product is also defined on the vector space of functions Φ mapping from a set S to a field \mathbb{K} , $\Phi : \forall f \text{ s.t. } f(s) = k, s \in S, k \in \mathbb{K}$ in which case the inner product on $g, f \in \Phi$ is defined as

$$\langle g, f \rangle = \int_{-\infty}^{\infty} g(x) \overline{f(x)} dx$$

and \bar{a} gives the complex conjugate of a .

General outer operators

The outer operator $\cdot \otimes_{\mathcal{O}} \cdot$ will only be defined for vectors in this thesis. It operates on the two vectors $\mathbf{x}, \mathbf{y} \in \mathbb{K}^N$ and is defined as

$$\mathbf{x} \otimes_{\mathcal{O}} \mathbf{y} = \mathbf{W}$$

where the i th row and j th column of \mathbf{W} are

$$w_{i,j} = \mathcal{O}(x_i, y_j)$$

Canonically, the operator \mathcal{O} is multiplication in which case

$$\mathbf{x} \otimes_{\times} \mathbf{y} = \mathbf{W}$$

where the i th row and j th column of \mathbf{W} are

$$W_{i,j} = x_i y_j$$

This outer product is also known as the *Kronecker product*, and we will omit the operator subscript when that is the case, i.e., we will simply write \otimes . As stated above, however, \mathcal{O} can be defined arbitrarily as any function taking two inputs and returning a single output.

Point-wise operators

If an operator on matrices \circ is written with a period preceding it, i.e., $\cdot \circ$ it means perform that operation on each element individually. Some examples follow.

For matrix $\mathbf{X} \in \mathbb{K}^{M,N}$ and $p \in \mathbb{K}$

$$\mathbf{X}^{\cdot p} = \mathbf{W}$$

where

$$W_{i,j} = X_{i,j}^p$$

For matrices $\mathbf{X}, \mathbf{Y} \in \mathbb{K}^{M,N}$

$$\mathbf{X}.\mathbf{Y} = \mathbf{W}$$

where

$$W_{i,j} = X_{i,j}Y_{i,j}$$

(contrast these with canonical matrix multiplication).

1.4.3 Random variables

Many authors denote random variables with a normally typeset uppercase letter. We will use this convention only when convenient, but will always state explicitly that a certain variable is random. We distinguish between discrete and continuous random variables in our notation.

Discrete random variables

If a random variable X can only take on values in a discrete set, we say that this random variable is discrete-valued. Formally a discrete set Γ of size N is one for which there exists an isomorphism \mathcal{J} that maps Γ on to the subset of the integers $[1, \dots, N]$. The probability that X takes on the value x is written $p(X = x)$ for discrete random variables.

Continuous random variables

If a random variable X can take on values in a set Γ isomorphic to \mathbb{R} we say this random variable is continuous-valued. The probability that X takes on the value x is written $p(x)$ for continuous random variables.

1.4.4 Complex numbers

A complex number $z \in \mathbb{C}$ can be described in Cartesian notation as

$$z = a + jb, a, b \in \mathbb{R}$$

or in polar notation as

$$z = \alpha \exp(j\omega), \alpha, \omega \in \mathbb{R}$$

where $j = \sqrt{-1}$. j is also often used to denote an index variable. It will be clear from the context when the imaginary number is meant and when the index is meant.

1.4.5 Logarithms

The logarithm base- e ¹ of x is written $\log(x)$. The logarithm in any other base b will be denoted as such: $\log_b(x)$.

1.4.6 Ultimate values

This thesis presents a number of iterative algorithms whose steps will be counted using an index l and whose solutions will take on the values x_l . We use the $*$ notation to refer to these values at the final iterates: the index of the final iteration is written² l^* and its value $x_{l^*} = x^*$.

¹ e is Euler's constant.

²Arbitrary letters can be used.

Chapter 2

Methodology

Most audio source separation strategies use some combination of two general methodologies: at one end of this continuum are those that use physical or structural models of the sources and at the other, those that use purely statistical models. Here we give a brief overview of some previously proposed techniques.

2.1 Additive sinusoidal model

The additive sinusoidal model [52], [36], is a convenient model with wide-spread use in the computer music community. Various authors have applied this model to the source separation problem. For example, in [61] a prior estimation of the fundamental frequencies is used in tandem with temporal and spectral smoothness constraints to separate sources estimated via an additive model. Similar to this thesis, [6] uses an additive sinusoidal model to provide frequency-modulation cues to a latent component technique, and in [31], common amplitude-modulation and fundamental frequency estimation are used to separate sources, the additive sinusoidal model being used to reconstruct the phases of overlapping harmonics using the fundamental frequencies. The additive sinusoidal model is also the model adopted in this thesis and its use is justified in Chapter 3.

2.2 Multiple fundamental frequency estimation

This technique assumes the signal considered can be described in a format akin to the musical score — a collection of notes each with times indexing their beginning and end

and a frequency, the fundamental, describing the perceived pitch of the note. Multiple fundamental frequency estimation for the purposes of music transcription dates back to the 1970s [20, ch. 20] [41]. This is related to audio source separation because the resulting high-level representation — the estimated score — can be used to synthesize signals corresponding to subsets of notes in the score, e.g., if a particular instrument is desired, its notes are extracted and then a signal is synthesized using stored recordings of the instrument or instrument-modeling synthesis techniques. The technique has become quite developed, see [20, ch. 20] for a review of modern techniques.

There are some drawbacks to the technique. Many musical signals of interest such as unpitched percussion, do not always have a perceivable fundamental frequency. A musical score describes notes as having distinct boundaries in time and frequency, which is not always true when one considers musical gestures such as portamento or *dal niente*¹. Nevertheless, the production of even a crude score from a musical signals is useful in applications such as automated music transcription (e.g., [51]) or music catalogue query [37].

2.3 Principal latent component analysis (PLCA) and non-negative matrix factorizations (NMF)

2.3.1 Motivation

Finding it difficult to specify a detailed model for all types of acoustic sources, we can attempt source separation with a more general criterion: that different sources only combine constructively in a mixture. This is a plausible criterion as the power spectrum of a signal and consequently its spectrogram are always non-negative valued. If a signal is considered as consisting of a sum of original source signals, these source signals will have non-negative spectrograms as well. The following two techniques attempt to determine these spectrograms solely from a spectrogram of their mixture using techniques for determining latent variables. We will refer to techniques of this sort as *latent variable models*. The technique is discussed in a bit more detail in the following to demonstrate how it differs from the additive technique explored in this thesis.

¹“Out of nothing”: usually accompanying a crescendo and indicating that the player start from silence and gradually increase their playing dynamic.

2.3.2 Approaches

The PLCA and NMF approaches to the audio source separation problem are very popular. An early and highly cited work that applies NMF to polyphonic music transcription is [54]. Since then many variations on this approach have been proposed. A technique using smoothness based on spectral difference and sparseness as regularization terms is presented in [62]. In [60] vectors in the resulting matrices are forced to be a linear combination of predefined “basis spectra”, chosen for their harmonic and perceptual properties. [1] explores the uses of different divergences for the purposes of up-mixing and noise removal.

2.3.3 PLCA

In this formulation, the spectrogram (defined in Section 3.2), being non-negative like a probability distribution, is considered as such

$$c|X(t, f)| = p(t, f)$$

where c is a constant so that the distribution integrate to 1. Explicitly, we consider the probability that energy in the spectrogram lie in the vicinity of time t and frequency f . With the hope of retrieving the spectrograms of the P underlying sources, it is proposed that the distribution is actually the distribution of K random variables each being chosen with probability $p(Z = k)$. The pair of random variables from component k , T_k and F_k are assumed independent, i.e., $p(t, f|Z = k) = p(t|Z = k)p(f|Z = k)$. Each random variable, it is hoped, describes a source ($K = P$) or a part of a source ($K > P$). In addition, each of these underlying distributions has marginal distributions $p(t|Z = k)$ and $p(f|Z = k)$. The marginal distributions and $p(Z = k)$ can be estimated using the expectation maximization algorithm with the following update rules for the l th iteration of the algorithm [55]

$$p_{l+1}(Z = k|t, f) = \frac{p_l(Z = k)p_l(t|Z = k)p_l(f|Z = k)}{\sum_{j=0}^{K-1} p_l(Z = j)p_l(t|Z = j)p_l(f|Z = j)}$$

$$p_{l+1}(t|Z = k) = \frac{\sum_f p(t, f)p_{l+1}(Z = k|t, f)}{\sum_s \sum_f p(s, f)p_{l+1}(Z = k|s, f)}$$

$$p_{l+1}(f|Z = k) = \frac{\sum_t p(t, f)p_{l+1}(Z = k|t, f)}{\sum_t \sum_g p(t, g)p_{l+1}(Z = k|t, g)}$$

$$p_{l+1}(Z = k) \frac{\sum_t \sum_f p(t, f) p_{l+1}(Z = k | t, f)}{\sum_{j=0}^{K-1} \sum_t \sum_f p(t, f) p_{l+1}(Z = j | t, f)}$$

After convergence, the marginal distribution $p_{l^*}(t|Z = k)$ gives the distribution of energy of the k th component over time. Similarly, the marginal distribution $p_{l^*}(f|Z = k)$ gives the distribution of energy of the k th component over frequency. Once the set of components $\{\tilde{k}\}$ belonging to the p th source has been determined, we can synthesize the spectrogram of this source as

$$|X_p(t, f)| = \frac{1}{c} \sum_{j \in \{\tilde{k}\}} p_{l^*}(t, f | Z = j) p_{l^*}(Z = j)$$

PLCA can be extended by the use of “kernel distributions” that allow the specification of marginal distributions with both time and frequency extent, and “entropic priors” that encourage sparsity in the resulting marginal distributions [53].

2.3.4 NMF

Instead of considering $|X(t, f)|$ as a probability distribution, we consider it at discrete frequencies mc_f and times nc_t with $m, n \in \mathbb{N}$, the entry at the m th row and n th column of matrix $V_{m,n} = |X(nc_t, mc_f)|$ with non-negative entries. We seek an approximate factorization of $\mathbf{V} \in \mathbb{R}_+^{M \times N}$ into matrices $\mathbf{W} \in \mathbb{R}_+^{M \times K}$ and $\mathbf{H} \in \mathbb{R}_+^{K \times N}$ such that

$$\mathbf{V} \approx \mathbf{W}\mathbf{H} \tag{2.1}$$

This can be done by solving the program

$$\min \mathcal{D}(\mathbf{V}, \mathbf{W}\mathbf{H})$$

subject to

$$\mathbf{V} \geq \mathbf{0}$$

$$\mathbf{W} \geq \mathbf{0}$$

$$\mathbf{H} \geq \mathbf{0}$$

The particular choice of function (\mathcal{D}) that measures divergence leads to different update equations in the iterative procedure for finding \mathbf{W} and \mathbf{H} .

The Kullback-Leibler divergence [30]

The *Kullback-Leibler* divergence function for measuring the divergence between two matrices \mathbf{X} and \mathbf{Y} is

$$\mathcal{D}_{\text{KL}}(\mathbf{X}, \mathbf{Y}) = \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} X_{m,n} \log(Y_{m,n}) - X_{m,n} + Y_{m,n}$$

and can be minimized using the following update equations for the l th iteration²

$$H_{a,b}^{l+1} = H_{a,b}^l \frac{\sum_{j=0}^{M-1} W_{j,a}^l V_{j,b}^l / (W^l H^l)_{j,m}}{\sum_{j=0}^{M-1} W_{j,a}^l}$$

$$W_{a,b}^{l+1} = W_{a,b}^l \frac{\sum_{j=0}^{N-1} H_{b,j}^{l+1} V_{a,j}^l / (W^l H^{l+1})_{a,j}}{\sum_{j=0}^{N-1} H_{b,j}^{l+1}}$$

The Itakura-Saito divergence [11]

Another divergence popular for factorizing spectrograms is the *Itakura-Saito* divergence

$$\mathcal{D}_{\text{IS}}(\mathbf{X}, \mathbf{Y}) = \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} \frac{X_{m,n}}{Y_{m,n}} - \log \left(\frac{X_{m,n}}{Y_{m,n}} \right) - 1$$

This divergence is scale-invariant, meaning that $\mathcal{D}_{\text{IS}}(c\mathbf{X}, c\mathbf{Y}) = \mathcal{D}_{\text{IS}}(\mathbf{X}, \mathbf{Y})$, which makes it well suited for audio signals which have a large dynamic range. Put another way, divergences involving large values in \mathbf{V} and \mathbf{WH} will be weighted similarly to divergences involving small values, which is not the case with the Kullback-Leibler divergence. The Itakura-Saito divergence is minimized through the following update equations

$$\mathbf{H}^{l+1} = \mathbf{H}^l \cdot \frac{\mathbf{W}^{lT} ((\mathbf{W}^l \mathbf{H}^l)^{-2} \cdot \mathbf{V}^l)}{\mathbf{W}^{lT} (\mathbf{W}^l \mathbf{H}^l)^{-1}}$$

$$\mathbf{W}^{l+1} = \mathbf{W}^l \cdot \frac{((\mathbf{W}^l \mathbf{H}^{l+1})^{-2} \cdot \mathbf{V}^l) \mathbf{H}^{l+1T}}{(\mathbf{W}^l \mathbf{H}^{l+1})^{-1} \mathbf{H}^{l+1T}}$$

Once convergence has been obtained the k th column of matrix \mathbf{W} will contain values representing the level of activation of the k th component at the frequency corresponding

²It can be shown that these update equations are equivalent to those for PLCA [53].

to the row index and the k th row of \mathbf{H} the level of activation of the k th component at the time corresponding to the column index. If the set of components $\{\tilde{k}\}$ belonging to the p th source has been determined, we can synthesize the spectrogram of this source as

$$|X_p(nc_t, mc_f)| = \sum_{j \in \{\tilde{k}\}} W_{:,j} H_{j,:}$$

2.3.5 Synthesis of factorized spectrograms

Synthesizing the original signal is less straightforward as the phase information contained in the STFT was discarded to obtain a non-negative spectrogram. We can simply use the original phases of the STFT used to compute the spectrum with the new magnitude information from $|X_p(t, f)|$ but the resulting signal may have artifacts due to the phase information of unwanted sources that remains in the STFT. A technique to lessen these artifacts using constrained Wiener filtering has been proposed in [29]. One may also choose to invert the spectrogram without any phase information by using an algorithm that iteratively reconstructs the phase part of the STFT while minimizing the error between the spectrogram of the reconstructed signal and its original power spectrum, starting from an initial guess [17]. Each iteration requires transforming the signal to its spectrogram and then back to a time-domain signal, requiring considerable computational effort.

2.3.6 Extensions and shortcomings

As with PLCA, certain extensions can be integrated into NMF to encourage particular solutions. For example, to promote sparseness in the resulting matrices, i.e., to encourage that fewer entries be non-zero, one can add the regularization term

$$\alpha \sum_{m=0}^{M-1} \sum_{k=0}^{K-1} W_{m,k}^2 + \beta \sum_{k=0}^{K-1} \sum_{n=0}^{N-1} H_{k,n}^2$$

to the divergence to penalize matrices with large entries. α and β are terms that control the influence of this regularization. A variety of regularization terms are possible [5, ch. 3].

Equation 2.1 can be seen as the sum

$$V \approx \sum_{k=-0}^{K-1} \mathbf{w}_{:,k} \otimes \mathbf{h}_{k,:}$$

i.e., the sum of all the outer products between columns of \mathbf{W} and rows of \mathbf{H} . From this perspective, we can see the columns of \mathbf{W} as a collection of K spectral “templates” and the rows of \mathbf{H} as their time-varying gains. What this means is, sounds exhibiting frequency-modulation require a large number of columns in \mathbf{W} . This large number of columns makes the organizing of templates into sources a difficult task.

2.4 An approach using amplitude- and frequency-modulation

Perceptual studies have shown that sounds modulated synchronously in amplitude or frequency are heard as one sound, whereas asynchronously modulated sounds are heard as distinct [35] [34]. Here we define the modulation of parameters θ_i and θ_j as being synchronous if they are given as functions of time, $\theta_i = f_i(t)$ and $\theta_j = f_j(t)$ and there is an affine transform \mathcal{A} such that $\mathcal{A}\{f_i\}(t) \approx Af_j(t) + B$ where A and B are constants that do not vary with time (at least for the time that we observe the signal). If we can accurately measure these parameters and they are typical of the sounds we are trying to separate, then we can design techniques to reliably separate these sounds from acoustic mixtures. This involves picking those parameters classified as belonging to the same sound, discarding the rest, and resynthesizing from these parameters. The task of audio source separation therefore comprises the following tasks:

- Decide on a signal model for the sound of interest, with parameters that can be estimated and that are similar for similar sounds. We have chosen the additive sinusoidal model with a higher-order polynomial description of amplitude and phase.
- Locate regions in the signal where these signals are thought to be present using a peak-picking technique. Estimate the signal parameters at these locations. Here we use the Distribution Derivative Method [2] to estimate these parameters.

- Use these measurements as input to a partial tracking algorithm. We compare the effectiveness of the original peak matching procedure of McAulay and Quatieri [36] and a new linear programming formulation of the problem.
- Classify the partials and group them as sets of parameters coming from the same source. One of these sources will be the sound of interest. The classification is carried out on the frequency- and amplitude-modulation parameters. Principal components analysis is used to preprocess the data before classification.
- Choose a group of parameters (partials) representing the sound of interest and synthesize the separated signal from them. We compare the quality of synthesis for three models. The first assumes constant amplitude and linear phase at the analysis step resulting in linear amplitude and cubic phase at the synthesis step [36]. The second assumes quadratic amplitude and phase at the analysis step but constrains the amplitude to be cubic at the synthesis step. The final model assumes quadratic amplitude and phase at the analysis step resulting in a quartic model for phase and amplitude at the synthesis step.

It should be noted that the strategy for source separation relies on the same perceptual principle as in [6]. The work here differs in many respects. We have chosen to use solely the additive sinusoidal model as a model of the signals considered. In addition to frequency-modulation, we incorporate amplitude-modulation to aid in the source separation. The source separation itself differs technically from their approach as well. Chapter 7 uses only the amplitude-modulation to classify partials into sources. This is similar to the strategy explored in [31], but does not use a prior estimation of the fundamental frequencies of the sources.

Chapter 3

Signal Modeling

3.1 Introduction

To build tools to analyse and synthesize signals some structure must be imposed on the signals. The structures chosen can reflect something about the behaviour of these signals as observed in the field, as we will see with sinusoidal models. Other structures are chosen phenomenologically — we do not really know the underlying mechanism behind the production of these signals, but a particular structure is chosen for its mathematical or conceptual convenience, such as is the case when we consider higher-order models for sinusoidal phase and amplitude.

We begin the chapter with what could be seen as a mathematical analog of the musical score: time-frequency representations. Through this we will justify the adoption of a sinusoidal model for musical signals. Finding this inadequate to describe the signals of interest with sufficient quality, the sinusoidal model is generalized to incorporate modulations in frequency and amplitude. A technique is described to estimate the parameters of these more complex models which requires windows that are everywhere differentiable — we design a new window having desirable properties close to those of well-known optimal windows, but that is everywhere differentiable.

3.2 Time-frequency representations

As most musical instruments are resonating media, and excited resonating media are well described as linear time-invariant (LTI) auto-regressive (AR) structures, many popular

models of musical audio are some variation of this description [12]. Strictly speaking, incorporating moving-average (MA) structures into a model of musical signals could improve its quality, but such a model would preclude the sum-of-sinusoids model adopted later in this thesis.

An LTI auto-regressive structure is a signal that can be described using the following *difference equation*:

$$x(n) = \sum_{k=1}^K a_k x(n-k) + b_0 v(n)$$

Here x is the output of the system (what is heard or measured) and v is the input. K is the order of the model. Both are general functions of time which, in the case of properly sampled digital audio, can be considered at discrete times $n \in \mathbb{Z}$ without any loss of information [7, ch. 2]. $a_k, b_0 \in \mathbb{C}$ and are constants. Casually you can think of the output of the system at time n as being a linear combination of past outputs, plus some of the scaled input.

AR structures are excited in various ways: some are bowed, others struck, etc. To characterize the above structure we excite it with a simple signal, the *Kronecker delta*

$$\delta(n) = \begin{cases} 1 & n = 0 \\ 0 & \text{otherwise} \end{cases}$$

This Kronecker delta input will yield its *impulse response* from which we can derive many properties of the AR structure.

As an example take the case where $K = 1$ and $a_1 = r \exp(j\omega)$, $r, \omega \in \mathbb{R}$, $|r| < 1$. Then the difference equation is

$$x(n) = r \exp(j\omega) x(n-1) + v(n)$$

Exciting this with the Kronecker delta we get

$$\begin{aligned} x(0) &= 1 \\ x(1) &= r \exp(j\omega) \\ x(n) &= r^n \exp(j\omega n) \end{aligned}$$

which is a complex exponential starting at $n = 0$ and periodic in $n_T = \frac{2\pi}{\omega}$ multiplied by the real-valued exponential r^n . In other words, the output is a damped sinusoid. From this it is not hard to see that if we can estimate the coefficients a_k , we can then know the frequencies, amplitudes and damping factors of the sinusoids that are output when this structure is excited by an impulse (the Kronecker delta). This principle is presented as a motivation for the following techniques and is not pursued here. The interested reader is referred to [33] for more information.

An alternative method for determining the frequencies and amplitudes of sinusoids in mixture is to take the inner product of the signal with a complex exponential of known frequency

$$X(\omega) = \sum_{n=-\infty}^{\infty} x(n) \exp(-j\omega n)$$

The function $X(\omega)$ will be large if $x(n)$ contains a complex exponential of frequency ω and small if it does not, effectively indicating which sinusoidal functions are present in the signal. This transformation of a signal as a function of time n into one as a function of frequency ω is known as the *Discrete-time Fourier Transform* (DTFT).

To create a variety of pitches and timbres, typically the media of musical instruments are not static, but vary in time. That means the sets of sinusoids describing the state of the media and its excitation also change in time. To account for this we consider many small intervals of signal where we assume its characteristics are roughly static. We can then piece these time-intervals together afterwards to get a description of the signal in both time and frequency. To do this, we multiply the signal by a window w which makes the signal 0 outside of the interval of interest. We then test what sinusoids with frequencies ω are present at different times τ , giving a function of two variables

$$X(\tau, \omega) = \sum_{n=-\infty}^{\infty} x(n)w(n - \tau) \exp(-j\omega n) \quad (3.1)$$

This transformation of a signal of time n into one of time τ and frequency ω is known as the *Discrete-time Short-time Fourier Transform* (DTSTFT).

One further point about the window should be discussed. The Fourier transform of the product of two functions, like we have in Equation 3.1, is equal to the convolution of the Fourier transform of each function separately. If we denote the Fourier transform operator

as \mathcal{F} , for functions g and f we have

$$\mathcal{F}(gf) = \mathcal{F}(g) * \mathcal{F}(f)$$

where $*$ is the convolution operator. The value $X(\tau, \omega)$, which will be a complex number, can be seen as describing the amplitude and phase of a sinusoid at that frequency and time. If the Fourier transform of the window function is not purely real its imaginary part will offset the phase of this sinusoid. It is usually simpler to avoid this complication. The Fourier transform of a real even function

$$f(n) = f(-n), f(n) \in \mathbb{R}$$

is real, so we choose windows with this property. See [19, p. 52] for a concrete illustration of this.

3.3 Polynomial phase models

The DTFT and DTSTFT are very useful because they are invertible [45]¹ and fast algorithms exist for their computation by digital computer [59]. If the presence of a sinusoid is determined, e.g., by finding τ^* and ω^* such that X is maximized, that signal can be removed or altered easily.

One drawback of these transforms is they only project onto sinusoidal functions of linear phase, i.e., functions of constant frequency. In general, musical signals are not linear combinations of sinusoids of constant frequency (consider, for example, vibrato). We could decide to project onto a different family of functions and considerable effort has been devoted to finding alternatives (see [28] for a review). In the case of musical signals, however, we have some prior information about the mechanics of their production and can make certain assumptions about the underlying functions.

3.3.1 Sinusoidal Representations

Many musical acoustic signals are quasi-harmonic [12], meaning that they consist of a sum of sinusoids whose frequencies are roughly integer multiples of a fundamental frequency.

¹Provided proper sampling in time and frequency.

For these kinds of signals, most of the energy can be attributed to sinusoids and so the signal can be described by a small number P of sinusoids with slowly varying amplitude and phase, plus some noise. The model is

$$x(n) = \sum_{p=1}^P A_p(n) \exp(j\phi_p(n)) + \varepsilon \quad (3.2)$$

where $\varepsilon \sim \mathcal{N}(0, \sigma^2)$, σ^2 quantifies the power of the noise, and $A_p(n), \phi_p(n) \in \mathbb{R}$ are the functions of amplitude and phase respectively for the p th sinusoid. In the following, we consider equivalent sinusoidal mixtures of complex-valued polynomial phase and log-amplitude exponentials

$$x(n) = \sum_{p=1}^P \exp(\mathcal{P}_p(n)) + \varepsilon$$

where \mathcal{P}_p is the complex valued polynomial of order Q describing the argument of the p th sinusoid, i.e.,

$$\mathcal{P}_p(n) = \sum_{q=0}^Q c_q n^q$$

and $c_q \in \mathbb{C}$. Note the form in Equation 3.2 can be retrieved as

$$\mathcal{P}_p(n) = \log(A_p(n)) + j\phi_p(n) = \Re\{\mathcal{P}_p(n)\} + j\Im\{\mathcal{P}_p(n)\}$$

In Chapter 5 we will see how estimations of these parameters at different times can be interpolated to create higher-order phase and log-amplitude functions with improved synthesis quality.

3.4 Polynomial phase parameter estimation

Assuming this sinusoidal model, how can we then estimate the parameters describing a signal of interest? Recently a set of techniques has been developed that use some combination of derivatives of the analysis window w or the signal x to estimate the polynomial coefficients directly [18]. For this thesis we will only consider a technique that does not estimate derivatives of the signal and only requires a once-differentiable analysis window as it is relatively easy to implement and suits our purposes.

The following is adapted from [2]. Consider the inner product of the signal $x(n) = \exp(\mathcal{P}_p(n))$ and a known analysis atom $\psi(n)$

$$\langle x, \psi \rangle = \int_{-\infty}^{\infty} x(n) \bar{\psi}(n) dn$$

Differentiating with respect to n , we obtain by the product rule

$$\frac{dx}{dn}(n) \bar{\psi}(n) + x(n) \frac{d\bar{\psi}}{dn}(n) = \left(\sum_{q=1}^Q qc_q n^{q-1} \right) x(n) \bar{\psi}(n) + x(n) \frac{d\bar{\psi}}{dn}(n)$$

If $\psi(t)$ is 0 outside of some interval $n \in [-T, T]$ then

$$\sum_{q=1}^Q qc_q \int_{-T}^T n^{q-1} x(n) \bar{\psi}(n) dn + \left\langle x, \frac{d\bar{\psi}}{dn} \right\rangle = 0$$

If we define the operator $\mathcal{T}^\alpha : (\mathcal{T}^\alpha x)(n) = n^\alpha x(n)$ we can write

$$\sum_{q=1}^Q qc_q \langle \mathcal{T}^{q-1} x, \bar{\psi} \rangle = - \left\langle x, \frac{d\bar{\psi}}{dn} \right\rangle$$

From this we can see that to estimate the coefficients c_q , $1 \leq q \leq Q$ we simply need R atoms with $R \geq Q$ to solve the linear system of equations

$$\sum_{q=1}^Q qc_q \langle \mathcal{T}^{q-1} x, \bar{\psi}_r \rangle = - \left\langle x, \frac{d\bar{\psi}_r}{dn} \right\rangle \quad (3.3)$$

for $1 \leq r \leq R$. To estimate c_0 we write the signal we are analysing as

$$s(n) = \exp(c_0) \exp \left(\sum_{q=1}^Q c_q n^q \right) + \eta(n)$$

$\eta(n)$ is the error signal, or the part of the signal that is not explained by our model. We also define a function $\gamma(n)$, the part of the signal whose coefficients have already been estimated

$$\gamma(n) = \exp\left(\sum_{q=1}^Q c_q n^q\right)$$

Computing the inner product $\langle s, \gamma \rangle$, we have

$$\langle s, \gamma \rangle = \langle \exp(c_0) \gamma, \gamma \rangle + \langle \eta, \gamma \rangle$$

The inner-product between η and γ is 0, by the orthogonality principle [26, ch. 12]. Furthermore, because $\exp(c_0)$ does not depend on n , we have

$$\langle s, \gamma \rangle = \exp(c_0) \langle \gamma, \gamma \rangle$$

so we can estimate c_0 as

$$c_0 = \log(\langle s, \gamma \rangle) - \log(\langle \gamma, \gamma \rangle) \quad (3.4)$$

The estimation of the coefficients of a phase and log-amplitude polynomial using this method is known as the *Distribution Derivative Method (DDM)*.

3.5 Choosing atom ψ

As we are dealing with mixtures of sinusoids of small bandwidth, in addition to the finite-time support constraint, we desire atoms whose inner-product is only significant within a finite bandwidth of interest. To construct these atoms, we multiply the Fourier atom by the window w

$$\psi_{\tau, \omega}^{\mathcal{F}w}(n) = w(n - \tau) \exp(-j\omega(n - \tau))$$

A good overview of different windows and their properties is given in [19]. We require that the window be at least once-differentiable and zero outside of a certain interval, therefore, somewhat informally, we require

$$\lim_{n \rightarrow T} \psi(n) = \psi(T) = 0$$

3.5.1 The Hann window

The *Hann* window possesses this property

$$w_h(n) = \begin{cases} 0.5 + 0.5 \cos\left(\frac{n}{T}\pi\right) & -T \leq n \leq T \\ 0 & \text{otherwise} \end{cases}$$

The Hann window is a member of a class of windows constructed by summing scaled harmonically related cosine functions, subject to the constraint that the scaling coefficients sum to 1 so that the window have a value of 1 at $n = 0$. Letting $T = N/2$, where N is the length of the window

$$w(n) = \begin{cases} \sum_{m=0}^{M-1} a_m \cos\left(\frac{2\pi}{N}mn\right) & -\frac{N}{2} \leq n \leq \frac{N}{2} \\ 0 & \text{otherwise} \end{cases}$$

With $M = 2$ and $a_0 = a_1 = 0.5$, we have the Hann window.

The simple expression for its calculation and good trade-off between main-lobe width and side-lobe height make the Hann window a popular choice in many signal processing applications. The expression for its Fourier transform is such that fast digital implementations of windowing a signal by a Hann window involve no multiplies [19, p. 183]. A recursive implementation of the DTSTFT is possible when windowing with the Hann window, which is important for applications where little storage is available [58, p. 102]. In spite of all its merits, other windows have been proposed that have superior qualities, such as lower side-lobes.

3.5.2 Continuous Blackman-Harris windows

A family of windows with certain properties superior to the Hann window is the *Blackman-Harris* family of windows. These are also sum-of-cosine windows and so are easily tabulated. To design these windows, optimization techniques were used to search for coefficients giving minimum height of the highest side-lobe (maximum out-of-band rejection) [47]. The 4-term window whose coefficients a are listed in Table 3.1 has a maximum side-lobe level of -92 dB, just shy of the quantization noise of a 16-bit linear pulse code modulated signal (-96 dB). As can be seen in Figure 3.2, this window has a very large main-lobe which means

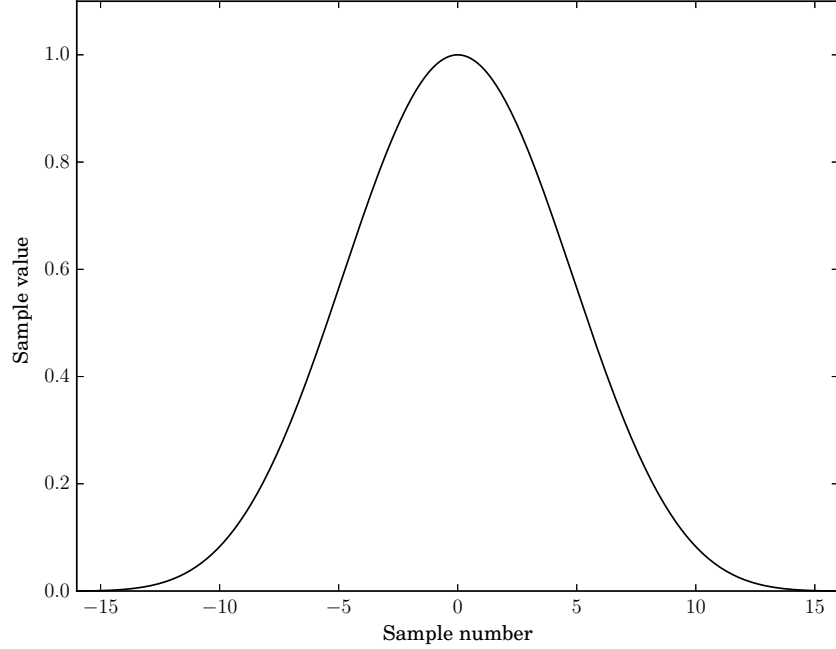


Fig. 3.1: Minimum 4-Term Blackman-Harris: Time-domain

two sinusoids of similar frequency will be difficult to resolve. Furthermore, as shown in Figure 3.3, the window has a discontinuity at its boundaries, e.g., $w\left(\frac{N}{2}\right) \neq 0$, and is not once-differentiable. In any case the window is valuable in that it effectively nulls any influence of signals outside of a bandwidth of interest. The shape of the Blackman-Harris window in the time- and frequency-domains can be seen in Figures 3.1 and 3.2 respectively.

It should be clarified that when we compare the widths of the main-lobes of two windows, we compare two windows of the same length. Of course, the bandwidth of a window can also be decreased by increasing its length, at the expense of time-resolution. When searching for windows superior to the Hann window, we are motivated by our ability to describe the

Table 3.1

Window	a_0	a_1	a_2	a_3
Minimum	0.35857	0.48829	0.14128	0.01168
\mathcal{C}^1	0.35874	0.48831	0.14127	0.01170

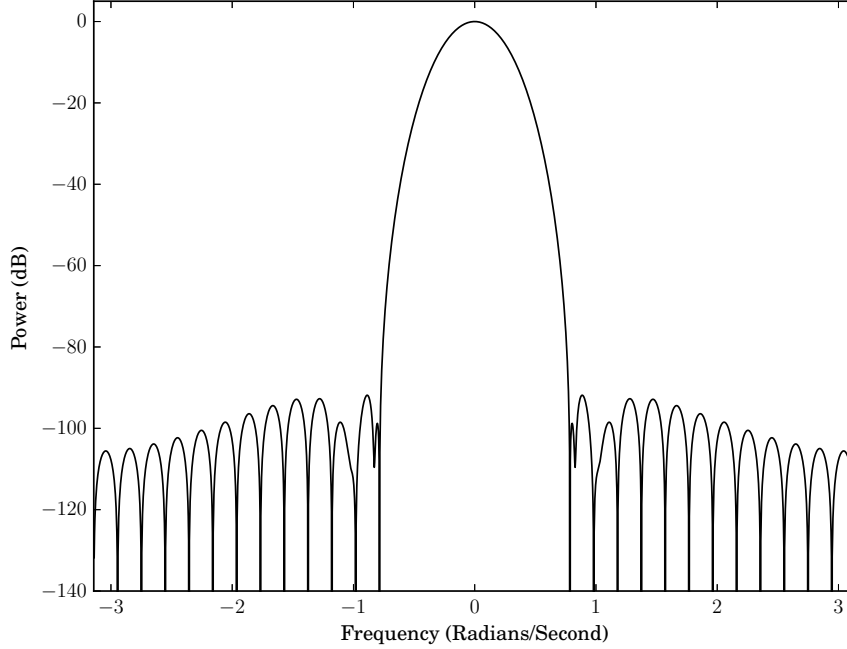


Fig. 3.2: Minimum 4-Term Blackman-Harris: Frequency-domain

signal with more detail between two analysis frames than would be possible with a simple linear-phase sinusoid model. Using a longer window to decrease the main-lobe width is not problematic in this case, but we would still like a high level of signal rejection outside of the bandwidth of interest for improved estimation accuracy of the signal parameters. For this reason, we search for windows that have very low side-lobe height and are also once-differentiable, without caring so much about the width of the main-lobe. To find a window with properties similar to the 4-term Blackman-Harris window but without a discontinuity, we solve the optimization problem

$$\min \|a - \tilde{a}\|_2$$

subject to

$$w_{\tilde{a}}\left(\frac{N}{2}\right) = w_{\tilde{a}}\left(\frac{-N}{2}\right) = 0$$

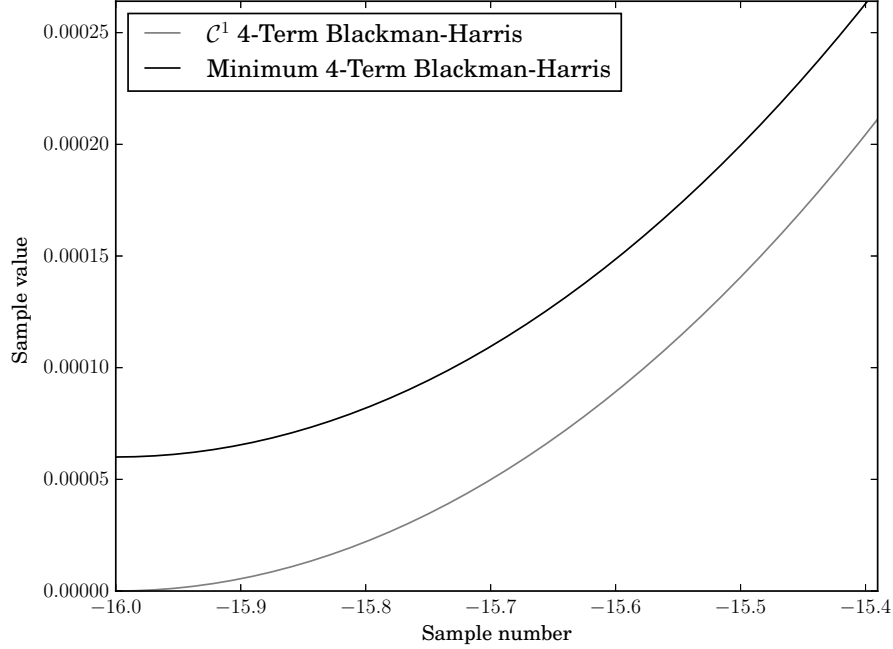


Fig. 3.3: Comparison of endpoints of window in time-domain. Here we observe that the \mathcal{C}^1 4-term Blackman-Harris window approaches 0 at the boundary of its time support, whereas the Minimum 4-Term Blackman-Harris window does not.

$$\sum_{m=0}^{M-1} a_m = 1$$

where

$$w_{\tilde{a}}(n) = \begin{cases} \sum_{m=0}^{M-1} \tilde{a}_m \cos\left(\frac{2\pi}{N} mn\right) & -\frac{N}{2} \leq n \leq \frac{N}{2} \\ 0 & \text{otherwise} \end{cases}$$

The solution \tilde{a}^* is given in Table 3.1 and time- and frequency-domain plots are given in Figures 3.4 and 3.5 respectively. This window will be referred to as the \mathcal{C}^1 4-Term Blackman-Harris window. Some figures of merit for the two windows are compared in Table 3.2 in a similar fashion to [19]. We see from comparing Figures 3.1 and 3.2 with 3.4 and 3.5 that the \mathcal{C}^1 4-Term Blackman-Harris window is not too different from the Minimum 4-term Blackman-Harris window, but has the additional desirable property of differentiability everywhere in its domain. A comparison of the windows's endpoints is presented in Figure 3.3.

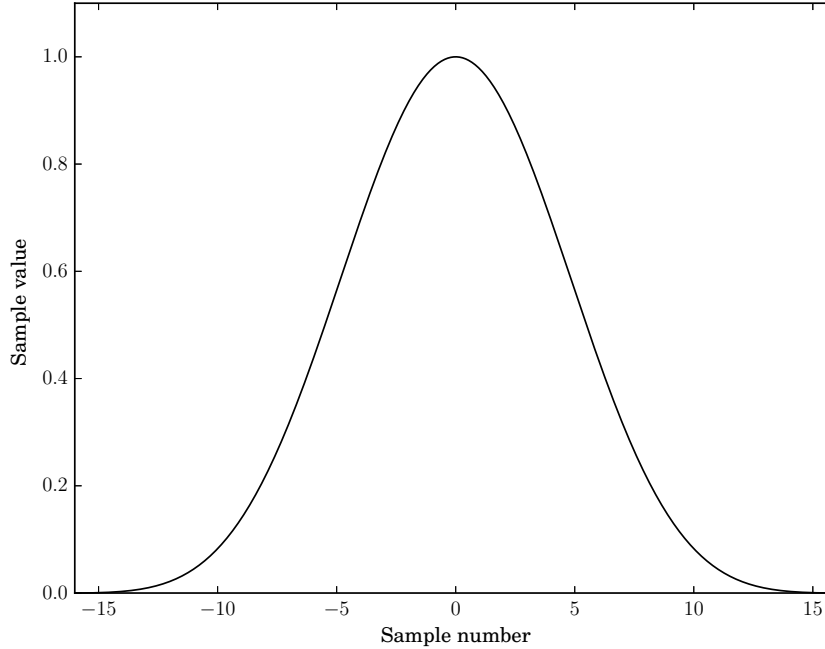


Fig. 3.4: \mathcal{C}^1 4-Term Blackman-Harris: Time-domain

3.6 Conclusion

In this chapter we have developed the rationale for adopting the sinusoidal model when analysing musical signals. In turn we have presented some techniques for estimating the parameters of these signals. Obviously these techniques only work as well as their assumptions are true — for the best results we should use these techniques only on signals that indeed contain sinusoids. For the signals considered in this thesis, we assume this to be true.

Table 3.2

Window	Highest side-lobe level (dB)	6-dB bandwidth in bins	Side-lobe fall-off (dB/octave)
Minimum	-92	2.72	6
\mathcal{C}^1	-90	2.66	12

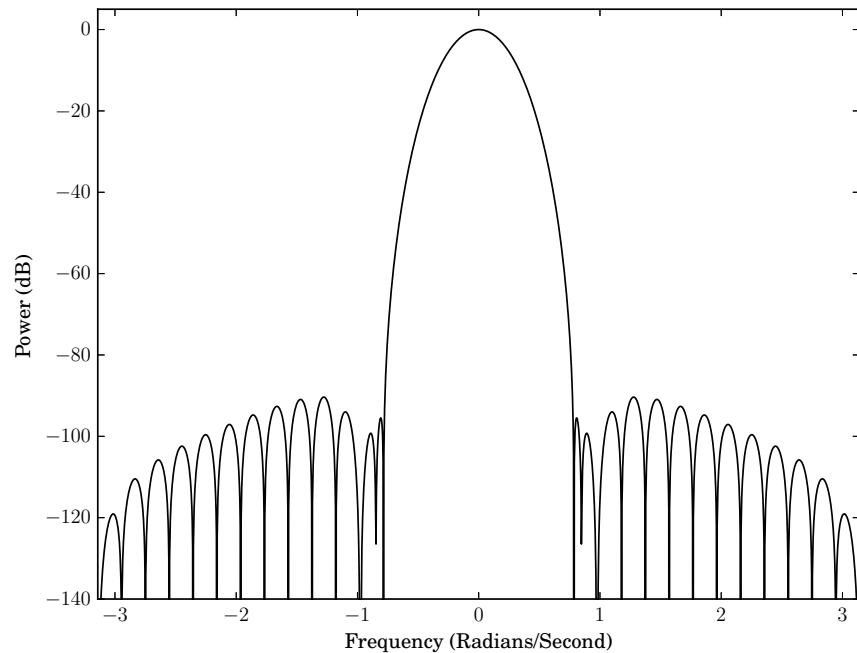


Fig. 3.5: C^1 4-Term Blackman-Harris: Frequency-domain

The techniques presented in this chapter are typically used to estimate the parameters of short signals as we see in the use of window functions to limit the time-frequency extent of our analysis. For the source separation problem we are interested in larger sinusoidal objects — partials — whose global properties are more readily classified. We could use a very large analysis window and a high-order polynomial for phase when solving for the coefficients, but the size of the linear system to be solved in Equation 3.3 will increase quadratically in the order of the model. Furthermore, it is difficult to account for situations where the signal is corrupted or briefly absent. In these situations we may prefer to use interpolation to reconstruct the signal in the corrupted region. For these reasons, we prefer to make multiple estimations of the parameters of low-order models and connect those estimations thought of as belonging to a single partial. We will see in Chapter 5 that this will allow postulating a higher-order phase model. Before that is possible, however, we must determine how to connect multiple estimations to form partials.

Chapter 4

Partial Tracking

In the previous chapter, we saw how to estimate parameters of sinusoids with polynomial phase. While theoretically applicable to signals of arbitrary length, for reasons of flexibility and efficiency, we usually estimate the local parameters of the signal under a low-order model and connect multiple estimations to form a partial. We will call these local estimations “analysis points” or “parameter sets”.

This chapter presents an interpretation of the *peak matching* procedure of McAulay and Quatieri [36], a classical approach to discovering partials. Our interpretation allows for the specification of an arbitrary cost function measuring the plausibility that a set of analysis points forms the path of a partial. With this path interpretation, we were able to design a technique that finds the optimal set of paths under a constraint on the number of paths. The chapter concludes with an example of partial tracking on a synthetic signal.

Typically the DTSTFT is computed for a block of contiguous samples, called a *frame* and these frames are computed every H samples, H being the *hop-size*. We will denote the M sets of parameters at local maxima in frame h as $\theta_0^h, \dots, \theta_{M-1}^h$ and the N in frame $h + 1$ as $\theta_0^{h+1}, \dots, \theta_{N-1}^{h+1}$ where h and $h + 1$ refer to adjacent frames. We are interested in paths that extend across K frames where each path touches only one parameter set and each parameter set is either exclusive to a single path or is not on a path.

4.1 A greedy method

In this section, we present the McAulay-Quatieri method of peak matching. It is conceptually simple and a set of short paths can be computed quickly, but it can be sensitive to

spurious peaks and is optimal only in the sense that the set of paths computed contains the best path possible — the quality of the other paths may be compromised under this criterion.

In [36, p. 748] the peak matching algorithm is described in a number of steps; we summarize them here in a way comparable with the linear programming formulation to be presented in the sequel. In that paper, the parameters of adjacent frames h and $h + 1$ are the instantaneous amplitude, phase, and frequency and are indexed by frequency as $\omega_0^h, \dots, \omega_{M-1}^h$ and $\omega_0^{h+1}, \dots, \omega_{N-1}^{h+1}$ but we will allow for arbitrary parameter sets. Define a distance function $\mathcal{D}(\theta_i, \theta_j)$ that computes the similarity between $K = 2$ sets of parameters. We will now consider a method that finds L pairs of parameters that are closest.

We compute the cost matrix \mathbf{C}

$$\mathbf{C} = \theta^h \otimes_{\mathcal{D}} \theta^{h+1}$$

so that the i th row and j th column contain $C_{i,j} = \mathcal{D}(\theta_i^h, \theta_j^{h+1})$. For each $l \in [0 \dots L - 1]$, find the indices i_l and j_l corresponding to the shortest distance, then remove the i_l th row and j_l th column from consideration and continue until L pairs have been determined or the distances exceed some threshold Δ . This is summarized in Algorithms 1

Algorithm 1: A generalized McAulay-Quatieri peak-matching algorithm.

Input: the cost matrix \mathbf{C}

Output: L pairs of indices Γ_i and Γ_j

$\Gamma_i \leftarrow \emptyset;$

$\Gamma_j \leftarrow \emptyset;$

for $l \leftarrow 0$ **to** $L - 1$ **do**

$i_l, j_l = \arg \min_{i \in [0, \dots, M-1] \setminus \Gamma_i, j \in [0, \dots, M-1] \setminus \Gamma_j} C_{i,j};$

if $C_{i_l, j_l} > \Delta$ **then**

return Γ_i, Γ_j

end

$\Gamma_i \leftarrow \Gamma_i \cup i_l;$

$\Gamma_j \leftarrow \Gamma_j \cup j_l;$

end

return Γ_i, Γ_j

This is a greedy algorithm because on every iteration the smallest cost is identified and its indices are removed from consideration. Perhaps choosing a slightly higher cost in one iteration would allow smaller costs to be chosen in successive iterations. This algorithm does not allow for that. In other terms, the algorithm does not find a set of pairs that represent a globally minimal sum of costs. Another drawback of the algorithm is that it only works between two successive frames. The cost function could be extended to consider K frames (K arbitrary) of parameter sets, constructing a K -dimensional tensor instead of a matrix, but assuming equal numbers of parameter sets in all frames, the search space would grow exponentially with K . Nevertheless, the method is simple to implement, computationally negligible when K is small, and works well with a variety of signals encountered in audio [36] [57].

4.2 An optimal method

There is a way to find a set of paths over multiple frames ($K > 2$) having the lowest total cost if we restrict the search to exactly L paths. Instead of indexing parameters by their frame number h , we make h part of the parameter set so that it can be used by the distance function \mathcal{D} . Assume that over K frames there are M total parameter sets. In this context we will consider them as nodes in a graph. We define the vector $\mathbf{c} \in \mathbb{R}^{M^2}$ where the entry $\mathbf{c}_{i+Mj} = \mathcal{D}(\theta_i, \theta_j)$. If we have a set of connections $\Gamma_{i,j}$ we can calculate the total cost of these connections by defining the vector

$$\mathbf{x}_{i+Mj} = \begin{cases} 1 & \text{there is a connection between } i \text{ and } j \\ 0 & \text{otherwise} \end{cases}$$

and then forming the inner product

$$c_{\text{total}} = \langle \mathbf{c}, \mathbf{x} \rangle$$

Note that a node cannot be connected to itself. The question is how to find \mathbf{x}^* so that c_{total} is minimized. If no constraints are placed on \mathbf{x} , the solution is trivial, but not useful. How do we constrain \mathbf{x} to give us a solution to the partial tracking problem? Let us consider an example.

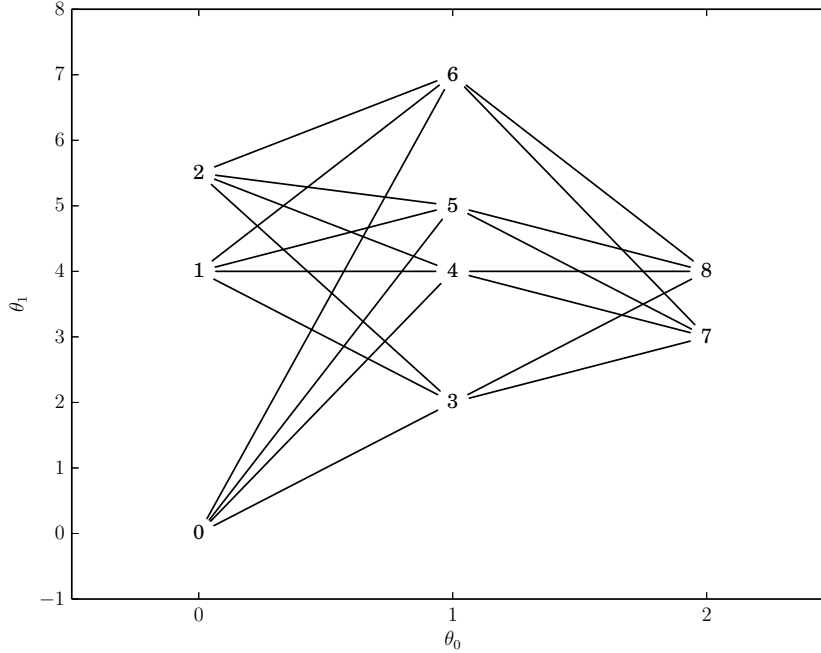


Fig. 4.1: Possible graph connections

In Figure 4.1 we have an example of a simple graph or lattice. Such a graph represents a plausible partial tracking situation: vertically aligned nodes are parameter sets estimated from the same analysis frame and we would like to connect these parameter sets between frames. The numbers are indices of nodes in the graph and the possible connections between them are indicated by lines, or *edges*. Imagine that we would like to find the two shortest paths. We will now examine the resulting paths from two algorithms using different criteria for shortness.

In Figure 4.2 we find the paths using an algorithm similar to Algorithms 1 but search instead over a tensor of distances $C \in \mathbb{R}^{3 \times 4 \times 2}$ whose entry $C_{i,j,h}$ represents the cost of travelling on the path connecting the i th node in layer 0, the j th node in layer 1 and the h th node in layer 2. This cost is the sum of the Euclidean distances giving the lengths of the connections. This is the greedy method of searching for the best paths whose optimality criterion is to find the set of best paths containing the absolute best path. We see in Figure 4.2 that the absolute shortest path, $1 \rightarrow 4 \rightarrow 8$, is discovered, followed by the second shortest path not using the nodes of the first path, $2 \rightarrow 5 \rightarrow 7$.

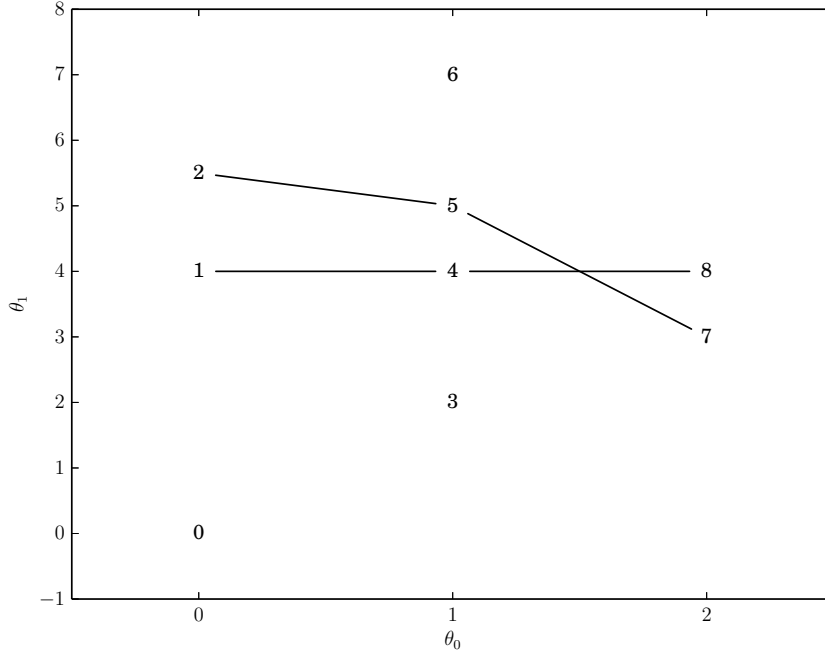


Fig. 4.2: Two shortest paths using the greedy method

4.2.1 L shortest paths via linear programming

To find a set of paths minimizing the total cost, we instead search for total solutions \mathbf{x} that describe all paths in the graph. Assume for now that we can guarantee that the entries of \mathbf{x} will be either 0 or 1. To find a set of constraints for our search, we consider the structure of a valid solution \mathbf{x}^* . To maintain that paths not overlap, a valid solution's nodes are only allowed to have one edge entering — coming from a node in a previous frame — and one edge leaving — going to a node in a successive frame. To translate this into a constraint, consider the node i and its possible R_i successive connecting nodes $j_0 \dots j_{R_i-1}$. Define the vector¹

$$a_{i+Mj_r}^{s,i} = \begin{cases} 1 & \forall j_r \in [j_0 \dots j_{R_i-1}] \\ 0 & \text{otherwise} \end{cases}$$

¹The superscript s stands for “successive”.

As all the entries of \mathbf{x} are either 0 or 1, we have

$$0 \leq \langle \mathbf{a}^{s,i}, \mathbf{x} \rangle \leq 1$$

so we can make this a constraint to ensure that a node has at most one path leaving. Similarly, if we consider the node j and its possible R_j previous connecting nodes $i_0 \dots i_{R_j-1}$, the vector²

$$\mathbf{a}_{i_r+Mj}^{p,j} \begin{cases} 1 & \forall i_r \in [i_0 \dots i_{R_j-1}] \\ 0 & \text{otherwise} \end{cases}$$

constrains that node j have only one path entering through the constraint

$$0 \leq \langle \mathbf{a}^{p,j}, \mathbf{x} \rangle \leq 1$$

A node on a path will also have an edge entering and an edge leaving. To translate this into a constraint, we define a vector that counts the number of edges entering a node and subtracts then the number of edges leaving a node. The result should always be 0 for an equal number of edges entering and exiting a node. If r is the index of the node considered, the vector is simply³

$$\mathbf{a}^{b,r} = \mathbf{a}^{p,r} - \mathbf{a}^{s,r}$$

and the constraint

$$\langle \mathbf{a}^{b,r}, \mathbf{x} \rangle = 0$$

Finally we want to constrain that there be only L paths. We do this by noticing that if this is true, there will be L edges between frames h and $h+1$. We constrain the number of paths going from edges Γ_h in frame h to Γ_{h+1} by forming the vector⁴

$$\mathbf{a}^{c,h} = \sum_{j \in \Gamma_h} \mathbf{a}^{s,j}$$

and asserting the constraint

$$\langle \mathbf{a}^{c,h}, \mathbf{x} \rangle = L$$

²The superscript p stands for “previous”.

³The superscript b stands for “balanced”.

⁴The superscript c stands for “connections”.

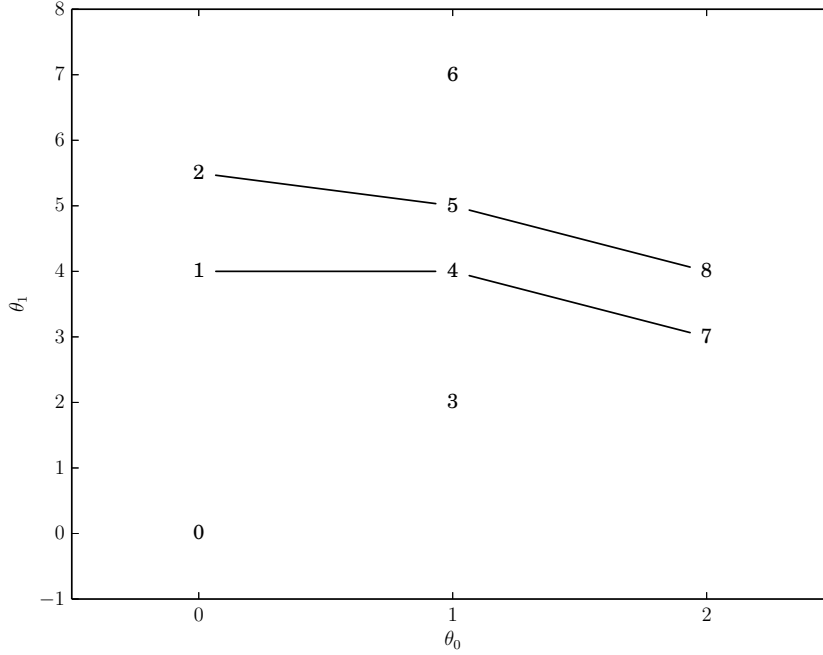


Fig. 4.3: Two shortest paths using the LP method

The length of \mathbf{x} is M^2 so the total size of all the constraints is not insignificant, but most entries in the constraint vectors will be 0 and therefore the resulting constraint matrices very sparse, so sparse linear algebra routines can be used in computations. Furthermore, the \mathbf{a}^b and \mathbf{a}^c constraints are derived from \mathbf{a}^p and \mathbf{a}^s , so only the latter need to be stored.

The complete *linear program (LP)* solving the L shortest paths problem is then

$$\min_{\mathbf{x}} \langle \mathbf{c}, \mathbf{x} \rangle$$

subject to

$$\begin{aligned} \mathbf{0} &\leq \begin{bmatrix} \mathbf{A}_s \\ \mathbf{A}_p \end{bmatrix} \mathbf{x} \leq \mathbf{1} \\ \begin{bmatrix} \mathbf{A}_b \\ \mathbf{A}_c \end{bmatrix} \mathbf{x} &= \begin{bmatrix} \mathbf{0} \\ L\mathbf{1} \end{bmatrix} \\ \mathbf{0} &\leq \mathbf{x} \leq \mathbf{1} \end{aligned}$$

where \mathbf{A}_s is the matrix with $\mathbf{a}^{s,m}$ as its rows for $m \in [0 \dots M - 1]$ and \mathbf{A}_p is the matrix with $\mathbf{a}^{p,m}$ as its rows, etc.

The solution of the two best paths using the LP formulation is shown in Figure 4.3 and a comparison of the total costs is shown in Table 4.1

The LP formulation is inspired by a multiple object tracking algorithm for video [23]. A proof that the solution \mathbf{x}^* will have entries equal to either 0 or 1 can be found in [44, p. 167]. The theoretical computational complexity of the linear program is polynomial in the number of variables, see [25] for a proof and the demonstration of a fast algorithm for finding its solution. In practice, to extract paths from the solution, we do not test equality with 0 or 1 but rather test if the solution vector's values are greater than some threshold. This may mean that suboptimal solutions may still be close enough. The tolerance of the solutions to suboptimality should be investigated, as if they are tolerant, fewer iterations of a barrier-based algorithm would be required to solve the problem. More information on linear programming and optimization in general can be found in [3].

4.2.2 Complexity

The LP formulation of the L -best paths problem gives results equivalent to the solution to the L -best paths problem proposed in [63]. The complexity of our algorithm is different. Assuming we use the algorithm in [25] to solve the LP, our program has a complexity of $O(M^7 B^2)$ where M is the number of nodes (parameter sets) and B is the number of bits used to represent each number in the input. The complexity of the algorithm by Wolf in [63] is equivalent to the Viterbi algorithm for finding the single best path through a trellis whose h th frame has $\binom{N_h}{L} \binom{N_{h+1}}{L} L!$ connections where N_h and N_{h+1} are the number of nodes in two consecutive frames of the original lattice. Therefore, assuming a constant number N of nodes in each frame, its complexity is $O(((\binom{N}{L})^2 L!)^2 T)$. If there are few nodes in each frame and a small number of paths are searched, Wolf's formulation is superior as its complexity increases linearly with the number of frames in the lattice. On the other hand,

Table 4.1. Comparison of total costs in Figure 4.3

Greedy	LP
5.354102	4.946461

if each frame has a large number of nodes or many paths are searched, the LP formulation is superior. Informally we have found this to agree with reality — both algorithms were tried when producing the figures in Section 4.3. Indeed the Wolf formulation took prohibitively long to compute when many paths were desired, as did the LP when many frames were considered.

It should be noted that in the special case that only 1 shortest path is searched an algorithm exists that requires on the order of N^2T calculations [46] where N is the number of nodes in each frame and T is the number of frames (assuming the same number of nodes in each frame): this algorithm is known as the Viterbi algorithm [13].

4.3 Partial paths on an example signal

We compare the greedy and LP based methods for peak matching on a synthetic signal. The signal is composed of $K = 6$ chirps of constant amplitude, the k th chirp s at sample n described by the equation

$$s_k(n) = \exp(j(\phi_k + \omega_k n + \frac{1}{2}\psi_k n^2))$$

The parameters for the 6 chirps are presented in Table 4.2.

Two 1 second long signals are synthesized at a sampling rate of 16000 Hz, the first with chirps 0–2, the second with chirps 3–5. We add Gaussian distributed white noise at several SNRs to evaluate the technique in the presence of noise.

Table 4.2. Parameters of k th chirp. f_0 and f_1 are the initial and final frequency of the chirp in Hz.

k	ϕ_k	ω_k	ψ_k	f_0	f_1
0	0	0.20	2.45×10^{-6}	500	600
1	0	0.39	4.91×10^{-6}	1000	1200
2	0	0.59	7.36×10^{-6}	1500	1800
3	0	0.27	-7.36×10^{-6}	700	400
4	0	0.55	-1.47×10^{-5}	1400	800
5	0	0.82	-2.21×10^{-5}	2100	1200

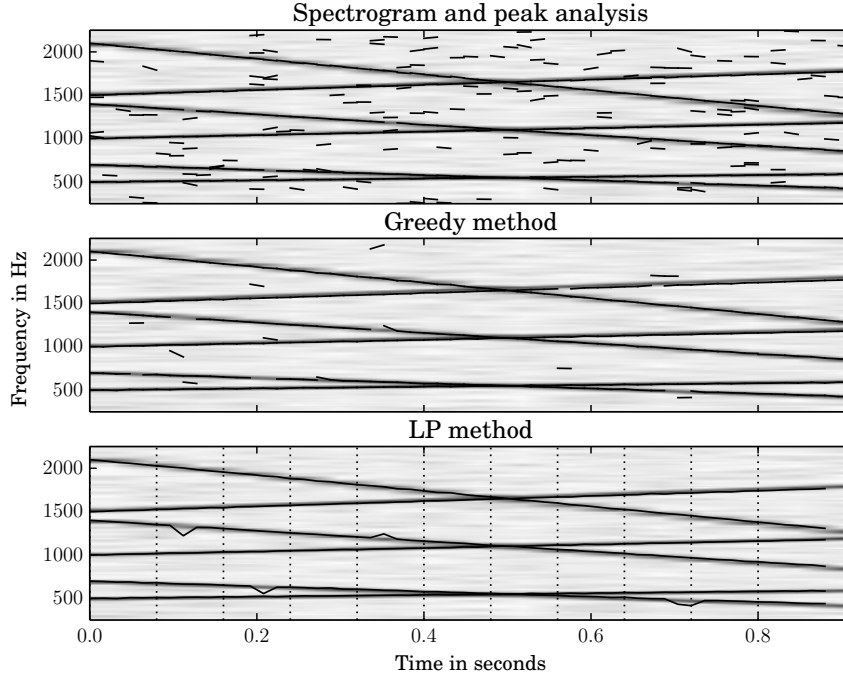


Fig. 4.4: Compare greedy and LP partial tracking on chirps in noise, SNR 20 dB. Line-segments representing the frequency and frequency-slope at local spectrogram maxima. In the bottom two plots the line segments not deemed by the respective algorithms as belonging to a partial path are discarded, revealing the estimated partial trajectories. See Table 4.2 for the chirp parameters.

A spectrogram of each signal is computed with an analysis window length of 1024 samples and a hop-size H of 256 samples. Local maxima are searched in 150 Hz wide bands spaced 75 Hz apart. A local maximum is only accepted if its amplitude is greater than -20 dB. At each local maximum the DDM is used to estimate the local chirp parameters, the i th set of parameters in frame h denoted $\theta_i^h = \{\phi_i^h, \omega_i^h, \psi_i^h\}$. The results of the analyses of both signals are lumped together and it is on this lumped data that we perform partial tracking.

We search for partial tracks using both the greedy and LP strategies. Both algorithms use the distance metric $\mathcal{D}_{\text{pr.}}$ between two parameters sets:

$$\mathcal{D}_{\text{pr.}}(\theta_i^h, \theta_j^{h+1}) = (\omega_i^h + \psi_i^h H - \omega_j^{h+1})$$

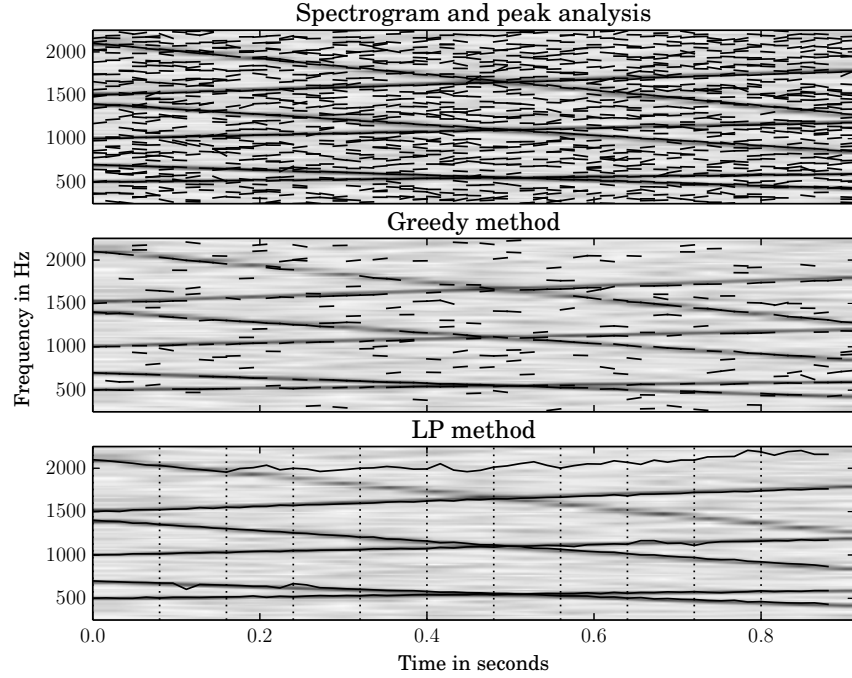


Fig. 4.5: Compare greedy and LP partial tracking on chirps in noise, SNR 15 dB. Line-segments representing the frequency and frequency-slope at local spectrogram maxima. In the bottom two plots the line segments not deemed by the respective algorithms as belonging to a partial path are discarded, revealing the estimated partial trajectories. See Table 4.2 for the chirp parameters.

which is the error in predicting j th frequency in frame $h + 1$ from the i th parameters in frame h . For the greedy method, the search for partial paths is restricted to one frame ahead like in [36]. For the LP method, to keep the computation time reasonable, we search over 6 frames for 6 best paths⁵. To maintain connected paths, the search on the next frames uses the end nodes of the last search as starting points. For both methods, the search is restricted to nodes between frequencies 250 to 2250 Hz.

Figures 4.4, 4.5, and 4.6 show discovered partial trajectories for signals with a SNR of 20, 15, and 10 dB, respectively. It is seen that while the greedy method begins to perform poorly at a SNR of 15dB, the LP method still gives plausible partial trajectories for SNRs of 10 and 15 dB. At lower SNRs, the LP formulation gives some paths that do not correspond to an underlying partial. These could be filtered out by examining the cost of these paths

⁵The number of paths does not affect the computation time.

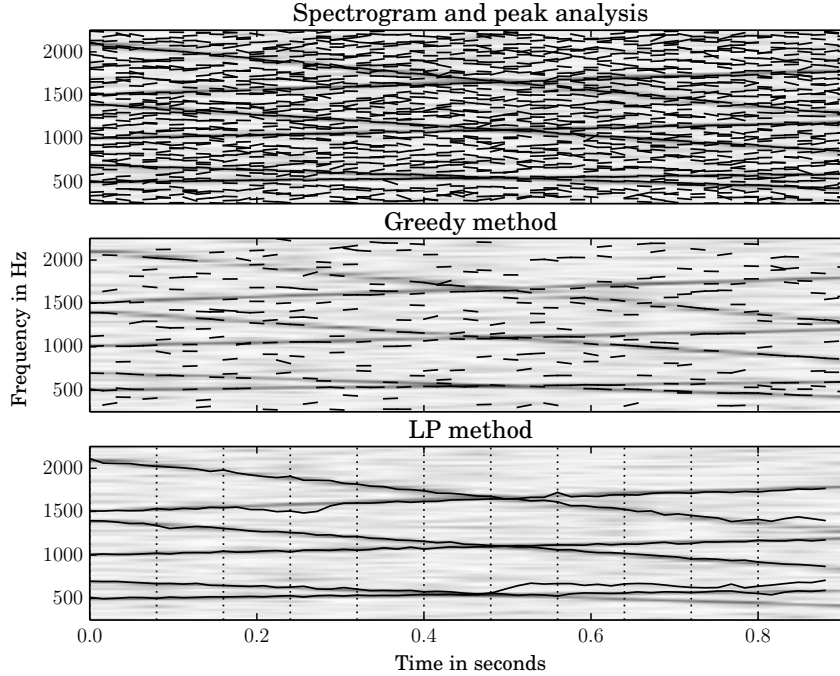


Fig. 4.6: Compare greedy and LP partial tracking on chirps in noise, SNR 10 dB. Line-segments representing the frequency and frequency-slope at local spectrogram maxima. In the bottom two plots the line segments not deemed by the respective algorithms as belonging to a partial path are discarded, revealing the estimated partial trajectories. See Table 4.2 for the chirp parameters.

and comparing them to the costs of the others. Those that deviate from a mean cost more than a certain amount should be rejected. This is the strategy used in Chapter 7 and illustrated in Figure 7.6. In any case, the lower SNRs are relatively challenging for any partial tracking technique.

But why did the LP discover a path not present in the underlying signal? This is due to the cost function, which finds a path with minimum prediction error in using the frequency and frequency slope coefficients of one node to predict another node's frequency coefficient. When there are many nodes in the original analysis it is not surprising that some unexpected path exists. An attribute of these erroneous paths is that they are not smooth. To deter the algorithm from finding such paths, regularization could be used like in Section 3.3.1 that minimizes the integral of the squared estimate of the path's second derivative. More on regularization in optimization can be found in [3, ch. 6.3].

4.4 Conclusion

In this chapter we reformulated the classical greedy algorithm of McAulay and Quatieri and showed that it can be seen as a greedy algorithm for finding the L shortest paths in a lattice. An algorithm was then proposed minimizing the sum of the L paths, using a linear programming approach. It was shown on synthetic signals that the new approach finds plausible paths in lattices with a large number of spurious nodes.

There are problems with the proposed approach. As discussed in 4.2, “jagged” paths should be removed using regularization. There are also situations where it is undesirable to have paths extend throughout the entire lattice. Acoustic signals produced by striking media, such as strings or bars, exhibit a spectrum where the upper partials decay more quickly than the lower ones (e.g., see Figure 7.1) — it would be desirable in these situations to have shorter paths for the upper partials, those decaying more quickly. This could be addressed as in [9] where the signal is divided into overlapping sequences of frames and partial paths are connected between sequences.

The proposed algorithm, while faster than algorithms based on the Viterbi algorithm, is still not fast. Assuming the same cost function $\mathcal{D}_{\text{pr.}}$ as in Section 4.3 it would be more efficient to consider narrow bands over which to search for paths when analysing signals with little frequency modulation. However, as we will see in Chapter 6, with different cost functions, the algorithm is useful for solving general L shortest paths problems outside of partial tracking.

Chapter 5

The extended phase and amplitude model

In Chapter 4 techniques were presented for discovering partials in a signal. Each partial is a set of analysis points indexed by time. The information at each analysis point can be used to synthesize a portion of the partial and these are combined to give a signal representing a partial. The various techniques to synthesize these pieces of the signal discussed here differ in the orders used for analysis and those for synthesis. The first technique presented simply synthesizes signals of short duration using the estimated parameters and blends these segments together. We recognize that having multiple estimations of parameters at discrete times within the partial's duration allow us to postulate, via interpolation, functions describing the partial of higher-order than those whose parameters were estimated during analysis. It is shown that this strategy is not always to our advantage — interpolants of higher-order than the underlying function can suffer from errors due to over-fitting. In the case of functions that are always better approximated by higher-order polynomials, we will see that there is an advantage to using high-order interpolation. These cases are illustrated through the analysis and synthesis of synthetic signals.

5.1 Partial synthesis

A popular technique for synthesizing partials from a set of analysis points is the *overlap-and-add* procedure [45], [40]. We assume that in the neighbourhood of τ_r the partial's signal is approximately described by the function $x(n) \approx f_{\tau_r}(n - \tau_r)$. To synthesize an

approximation of x we sum windowed f_{τ_r} at multiple locations, windowed by a function w with finite support so the resulting signal has finite energy and the piecewise assumption is maintained. For simplicity we assume the τ_r are equally spaced by H samples, and $\tau_0 = 0$, so we have $\tau_r = rH$. The length of the window function w is $M = VH + 1$ samples, with $V, H \in \mathbb{N}^1$. The approximate signal at sample n is then

$$\tilde{x}(n) = \sum_{l=L_-}^{L_+} w(n - lH) f_{\tau_l}(n - \tau_l)$$

where

$$L_- = \left\lceil \frac{n}{H} \right\rceil - V$$

and

$$L_+ = \left\lfloor \frac{n}{H} \right\rfloor + V$$

This method has some drawbacks. Usually the function f is an approximation \tilde{f} of the true underlying function. In the case of partial tracking, often partials that are too short are discarded or missed. At amplitude transients, these short partials are important for reproducing sharp attacks that are shorter than the window length. If these partials are missing, the resulting signal takes on a transient similar to the window shape. This could be overcome by choosing a window with a shape similar to the overall amplitude envelope in the attack region when resynthesizing an attack transient.

Another drawback is that no attempt is made to interpolate between the functions estimated at τ_r and τ_{r+1} using the model that the underlying sinusoids are non-stationary². From Equation 3.3 we know we can estimate a polynomial of arbitrary order for the argument of the complex exponential. We will see that using this additional information can give us an interpolating function closer to the underlying model.

¹ M is always odd so one may wonder how the Fast Fourier Transform can be used to invert F_{τ_r} , the frequency-domain representation of f_{τ_r} . Recall that the DTFT can be interpreted as the coefficients of a Fourier series that give the periodic version of the analysed signal. Also recall that we use window functions that are real and even. In practice the edges of the window are often equal to zero so that the length of the non-zero part is equal to the length of the DTFT N . In the case they are not, simply ensuring that values in the window indexed by integer multiples of N are 0, and that the value at the centre of the window is 1 will ensure proper synthesis [45, p. 244]. In that case, the values outside of the part of the window presented to the DTFT are folded into this region using the window indices modulo- N . See [45] and [40] for more details on this procedure.

²This is one of the causes of “pre-echo” when time-stretching using the STFT [49].

5.2 The interpolating analysis-synthesis system

In the following, we investigate the synthesis quality of three interpolating analysis-synthesis systems. The qualifier “interpolating” is used because each system takes the multiple sets of estimated parameters of a smaller-order model and interpolates them with a higher-order model, which is then used for synthesis. This is necessary because we only have values every H samples from the analysis step but require a value for every sample value in the output. The systems will be denoted $\mathcal{S}_{p,q}$ where p is the order of the analysis system and q the order of the synthesis system, e.g., a linear analysis system has $p = 1$, etc.

5.3 $\mathcal{S}_{1,3}$: the McAulay-Quatieri method

5.3.1 Analysis: linear phase and constant amplitude

For the McAulay-Quatieri method, the analysis model is a sinusoid of constant frequency and amplitude (linear phase and constant log-amplitude) in each analysis frame. To estimate the frequency of this sinusoid we find the bin with the most energy and find a refined estimate of the frequency as the maximum of a quadratic interpolating polynomial fit to this bin and its two neighbouring bins. This is a procedure documented in [52, p. 45]. The interpolation is best performed in the log-spectrum and on a spectrum produced using a window, such as a zero-padded Hann window, giving a wide enough main-lobe so that the three points lie on this lobe and not on side-lobes. A refined estimate of the amplitude of the sinusoid is obtained with this procedure as well.

In the original paper by McAulay and Quatieri, they do not use this technique but, as they show an example analysis of a speech signal, instead adjust the analysis window to be a multiple of the period of the glottal pulse. The bins of the DTFT used in the analysis will be integer multiples of the frequency given as the reciprocal of this period. Under the model of the speech signal as harmonically related sinusoids, the best estimate for the frequency is the bin of a local maximum, its amplitude the modulus of the spectrum at this maximum, and the phase the argument.

In our system, we use a fixed window size. To estimate the initial phase then we use Equation 3.4 with

$$\gamma_{\text{MQ}}(n) = \exp(2\pi \frac{k^*}{M}n)$$

where k^* is the bin we have determined to correspond to the frequency of the sinusoid. The initial phase is then $\Im\{c_0\}$.

5.3.2 Synthesis: cubic phase and linear log-amplitude

The phase part

Given two local maxima of the DTSTFT $X(\tau_0, \omega_0)$ and $X(\tau_1, \omega_1)$, where $H = \tau_1 - \tau_0$ we can conjecture a cubic polynomial phase function for the imaginary part of the phase argument

$$\tilde{\phi}(n) = \Im\{c_3\}(n - \tau_0)^3 + \Im\{c_2\}(n - \tau_0)^2 + \Im\{c_1\}(n - \tau_0) + \Im\{c_0\}$$

By noting that we have 2 measurements of the phase and frequency, $\angle\{X(\tau_0, \omega_0)\}$ and $\angle\{X(\tau_1, \omega_1)\}$, and the frequency is the derivative of the phase, we can solve for the coefficients of the polynomial phase function using the following linear system of equations, assuming the DTSTFT was computed using a real and even window

$$\begin{pmatrix} 0 & 0 & 0 & 1 \\ H^3 & H^2 & H & 1 \\ 0 & 0 & 1 & 0 \\ 3H^2 & 2H & 1 & 0 \end{pmatrix} \begin{pmatrix} \Im\{c_3\} \\ \Im\{c_2\} \\ \Im\{c_1\} \\ \Im\{c_0\} \end{pmatrix} = \begin{pmatrix} \angle\{X(\tau_0, \omega_0)\} \\ \angle\{X(\tau_1, \omega_1)\} + 2\pi M \\ \omega_0 \\ \omega_1 \end{pmatrix}$$

We choose M so that

$$\int_0^H \left(\frac{d^2 \tilde{\phi}}{dt^2}(t) \right)^2 dt \quad (5.1)$$

is minimized in order to have a smooth evolution of frequency in the interpolated region. M is necessary because some integer number of periods of a sinusoid will have passed from times τ_0 to τ_1 . Informally we choose M so that a polynomial describing the phase evolution between these two times takes a direct route, which is a plausible criterion because a signal with more radical phase variation would unlikely exhibit a spectrum that could be well described by two points in the time-frequency plane, i.e., the signal would exhibit a large bandwidth. See [36, p. 751] for further clarification.

The amplitude part

As only two measurements of the amplitude of the sinusoid are available, $|X(\tau_0, \omega_0)|$ and $|X(\tau_1, \omega_1)|$, the coefficients c_3 and c_2 are purely imaginary and the real parts of c_1 and c_0

are determined as

$$\begin{pmatrix} 0 & 1 \\ H & 1 \end{pmatrix} \begin{pmatrix} \Re\{c_1\} \\ \Re\{c_0\} \end{pmatrix} = \begin{pmatrix} \log(|X(\tau_0, \omega_0)|) \\ \log(|X(\tau_1, \omega_1)|) \end{pmatrix}$$

5.4 $\mathcal{S}_{2,3}$ and $\mathcal{S}_{2,5}$: the DDM-based methods

Here we extend the $\mathcal{S}_{1,3}$ model of McAulay-Quatieri to account for the additional parameters estimated via the DDM. For the $\mathcal{S}_{2,3}$ model, we must introduce additional constraints into the system as we have more estimated parameters than are available in the synthesis model. It would be possible to solve this system via least-squares, but the proposed constraints simplify analytically the expression maximizing the smoothness of the phase and log-amplitude functions, and give satisfactory results. For the $\mathcal{S}_{2,5}$ model the derivation is straightforward as in the $\mathcal{S}_{1,3}$ case — there are the same number of estimated parameters as there are parameters in the model.

5.4.1 Analysis: quadratic phase and log-amplitude

The DDM is used on segments of the signal to estimate the parameters of a sinusoid with a complex quadratic polynomial argument. This sinusoid has the form

$$x_a(n) = \exp(a_2 n^2 + a_1 n + a_0) \quad (5.2)$$

with $a_i \in \mathbb{C}$.

We can estimate the coefficients of Equation 5.2 using the DDM. We write Equation 3.3 in matrix form with $Q = 2$

$$\begin{pmatrix} \langle \mathcal{T}^0 x, \bar{\psi}_1 \rangle & \langle \mathcal{T}^1 x, \bar{\psi}_1 \rangle \\ \vdots & \vdots \\ \langle \mathcal{T}^0 x, \bar{\psi}_R \rangle & \langle \mathcal{T}^1 x, \bar{\psi}_R \rangle \end{pmatrix} \begin{pmatrix} a_1 \\ 2a_2 \end{pmatrix} = \begin{pmatrix} -\langle x, \frac{d\bar{\psi}_1}{dn} \rangle \\ \vdots \\ -\langle x, \frac{d\bar{\psi}_R}{dn} \rangle \end{pmatrix}$$

From this we recognize we need to define three functions

$$\langle \mathcal{T}^0 x, \bar{\psi}_k \rangle = \sum_{m=0}^{M-1} w(m) x(m + \tau) \exp(-j2\pi \frac{km}{M})$$

$$\langle \mathcal{T}^1 x, \bar{\psi}_k \rangle = \sum_{m=0}^{M-1} mw(m)x(m+\tau) \exp(-j2\pi \frac{km}{M})$$

$$\left\langle x, \frac{d\bar{\psi}_k}{dn} \right\rangle = -j2\pi \frac{k}{M} X_{p_1}(\tau, k) + \sum_{m=0}^{M-1} \frac{dw}{dn}(m)x(m+\tau) \exp(-j2\pi \frac{km}{M})$$

Where M is the length of the window and k is the frequency “bin”³. We also only consider x at the samples $m = [0, \dots, M-1]$ because we can always shift the time reference to view an arbitrary contiguous segment of signal with these indices.

We then find k^* such that X is maximum. If multiple components are present in the signal and are sufficiently separated in frequency, we can split the signal up into frequency bands and find local maxima. A technique for doing so is described in [52, p. 42]. To have a system of equations with a unique solution, we take the two adjacent bins $k-1$ and $k+1$ to have enough unique atoms for Equation 3.3. These bins should only contain energy from the component whose parameters we are interested in measuring — this is true if the components are adequately separated in time and frequency. We could choose only two bins to have a non-singular system, and there are many possibilities for choosing different atoms [2, p. 4639]. We choose three from the same frame to have improved estimation accuracy in situations where components are adequately separated in frequency, while avoiding a more sophisticated local peak selection procedure. Then a_2 and a_1 can be determined by solving the linear system

$$\begin{pmatrix} \langle \mathcal{T}^0 x, \bar{\psi}_{k-1} \rangle & \langle \mathcal{T}^1 x, \bar{\psi}_{k-1} \rangle \\ \langle \mathcal{T}^0 x, \bar{\psi}_k \rangle & \langle \mathcal{T}^1 x, \bar{\psi}_k \rangle \\ \langle \mathcal{T}^0 x, \bar{\psi}_{k+1} \rangle & \langle \mathcal{T}^1 x, \bar{\psi}_{k+1} \rangle \end{pmatrix} \begin{pmatrix} a_1 \\ 2a_2 \end{pmatrix} = \begin{pmatrix} -\left\langle x, \frac{d\bar{\psi}_{k-1}}{dn} \right\rangle \\ -\left\langle x, \frac{d\bar{\psi}_k}{dn} \right\rangle \\ -\left\langle x, \frac{d\bar{\psi}_{k+1}}{dn} \right\rangle \end{pmatrix}$$

With a_1 and a_2 determined, we can use Equation 3.4 to estimate a_0 . We will write a_i^τ to refer to coefficient i determined at time τ .

5.4.2 Synthesis: cubic order ($\mathcal{S}_{2,3}$)

In this section we describe how to obtain a cubic phase and log-amplitude polynomial from local estimations of the coefficients of a quadratic phase and log-amplitude polynomial.

³For tractability, the functions $X(\tau, k)$ are only evaluated at a finite number of frequencies, which are often called “bins” in the signal processing literature.

The phase part

A complex sinusoid with cubic phase has the following form:

$$\beta(n) = \exp(j(b_3 n^3 + b_2 n^2 + b_1 n + b_0)) \quad (5.3)$$

with $b_i \in \mathbb{R}$. This sinusoid has magnitude 1 everywhere, only its phase is changing.

Once the \mathbf{a}^τ have been determined at two times τ_0 and τ_1 , with $H = \tau_1 - \tau_0$, and these times have been determined as connected (see Chapter 4), we can write a system of equations to determine an interpolating cubic phase polynomial. To avoid numerical instabilities and for simplicity, we shift the time origin so that $\tau_0 = 0$. This means $b_0 = \Im\{a_0^{\tau_0}\}$. To reduce the size of the system, we require that

$$\frac{d\phi}{dn} \left(\frac{H}{2} \right) = \frac{1}{2} (\Im\{a_1^{\tau_0}\} + \Im\{a_1^{\tau_1}\})$$

and

$$\frac{d^2\phi}{dn^2} \left(\frac{H}{2} \right) = \frac{1}{2} (\Im\{a_2^{\tau_0}\} + \Im\{a_2^{\tau_1}\})$$

i.e., the frequency and first-order frequency modulation in the middle of the segment are the average of the two measured coefficients. Finally we require that the change in phase from time 0 to H correspond to that what was observed, but account for the cycles that were not observed by adding an integer number of 2π radians. If $\phi(n) = j(b_3 n^3 + b_2 n^2 + b_1 n + b_0)$, then

$$\phi(H) = \Im\{a_0^{\tau_1}\} - \Im\{a_0^{\tau_0}\} + 2\pi U^*$$

where $U^* \in \mathbb{Z}$ is determined to minimize Equation 5.1, in this case:

$$\tilde{U} = \arg \min_U \int_0^H (6b_3 t + 2b_2)^2 dt \quad (5.4)$$

which is then rounded to the nearest integer to give U^* . To summarize we have

$$\begin{pmatrix} H^3 & H^2 & H \\ \frac{3}{4}H^2 & H & 1 \\ 3H & 2 & 0 \end{pmatrix} \begin{pmatrix} b_3 \\ b_2 \\ b_1 \end{pmatrix} = \begin{pmatrix} \Im\{a_0^{\tau_1}\} - \Im\{a_0^{\tau_0}\} + 2\pi U^* \\ \frac{1}{2} (\Im\{a_1^{\tau_0}\} + \Im\{a_1^{\tau_1}\}) \\ \frac{1}{2} (\Im\{a_2^{\tau_0}\} + \Im\{a_2^{\tau_1}\}) \end{pmatrix}$$

Solving for b_1, \dots, b_3 , we have

$$\begin{pmatrix} b_3 \\ b_2 \\ b_1 \end{pmatrix} = \begin{pmatrix} \frac{4}{H^3} (\Im \{a_0^{\tau_1}\} - \Im \{a_0^{\tau_0}\} + 2\pi U^*) - \frac{2}{H^2} (\Im \{a_1^{\tau_1}\} + \Im \{a_1^{\tau_0}\}) \\ \frac{-6}{H^2} (\Im \{a_0^{\tau_1}\} - \Im \{a_0^{\tau_0}\} + 2\pi U^*) - \frac{3}{H} (\Im \{a_1^{\tau_1}\} + \Im \{a_1^{\tau_0}\}) + \frac{1}{4} (\Im \{a_2^{\tau_0}\} + \Im \{a_2^{\tau_1}\}) \\ \frac{-H}{4} (\Im \{a_2^{\tau_0}\} + \Im \{a_2^{\tau_1}\}) + \frac{3}{H} (\Im \{a_0^{\tau_1}\} - \Im \{a_0^{\tau_0}\} + 2\pi U^*) - \Im \{a_1^{\tau_1}\} - \Im \{a_1^{\tau_0}\} \end{pmatrix}$$

and then \tilde{U} is determined using Equation 5.4 to be

$$\tilde{U} = \frac{1}{4\pi} [H (\Im \{a_1^{\tau_1}\} + \Im \{a_1^{\tau_0}\}) - 2 (\Im \{a_0^{\tau_1}\} - \Im \{a_0^{\tau_0}\})]$$

and then rounded to obtain U^* .

The amplitude part

Solving for the cubic polynomial describing the local log-amplitude function

$$\mu(n) = \exp(d_3 n^3 + d_2 n^2 + d_1 n + d_0)$$

with $d_i \in \mathbb{R}$, is more straightforward analytically as it does not require solving to maximize the smoothness of resulting polynomial. To require continuity at the end-points of our polynomial, we require

$$\mu(0) = \Re \{a_0^{\tau_0}\}$$

and

$$\mu(H) = \Re \{a_0^{\tau_1}\}$$

The first constraint is satisfied simply by setting $d_0 = \Re \{a_0^{\tau_0}\}$. The second will be accounted for in a constrained least-squares solution for the other coefficients. The other observations are

$$\frac{d\mu}{dn}(0) = \Re \{a_1^{\tau_0}\}$$

$$\frac{d\mu}{dn}(H) = \Re \{a_1^{\tau_1}\}$$

$$\frac{d^2\mu}{dn^2}(0) = \Re \{a_2^{\tau_0}\}$$

$$\frac{d^2\mu}{dn^2}(H) = \Re \{a_2^{\tau_1}\}$$

The constrained least-squares problem to be solved is then

$$\begin{pmatrix} 0 & 0 & 1 \\ 3H^2 & 2H & 1 \\ 0 & 2 & 0 \\ 6H & 2 & 0 \end{pmatrix} \begin{pmatrix} d_3 \\ d_2 \\ d_1 \end{pmatrix} = \begin{pmatrix} \Re \{a_1^{\tau_0}\} \\ \Re \{a_1^{\tau_1}\} \\ \Re \{a_2^{\tau_0}\} \\ \Re \{a_2^{\tau_1}\} \end{pmatrix}$$

subject to

$$\begin{pmatrix} H^3 & H^2 & H \end{pmatrix} \begin{pmatrix} d_3 \\ d_2 \\ d_1 \end{pmatrix} = \left(\Re \{a_0^{\tau_1}\} \right)$$

This can be solved using numerical methods, in particular, using a specific interpretation of weighted least-squares [16, p. 266].

5.4.3 Synthesis: quintic order ($\mathcal{S}_{2,5}$)

The phase part

Solving for the coefficients of a quintic phase polynomial is done very similarly to Section 5.4.2. As we have the same number of analysis and synthesis parameters, no constraints have to be introduced to solve the system apart from the value U that maximizes smoothness of the phase function. The quintic phase polynomial is⁴

$$\lambda(n) = \exp \left(j \left(u_5 n^5 + u_4 n^4 + u_3 n^3 + u_2 n^2 + u_1 n + u_0 \right) \right) \quad (5.5)$$

with $u_i \in \mathbb{R}$. We have

$$\begin{aligned} \phi(0) &= \Im \{a_0^{\tau_0}\} \\ \frac{d\phi}{dn}(0) &= \Im \{a_1^{\tau_0}\} \\ \frac{d^2\phi}{dn^2}(0) &= \Im \{a_2^{\tau_0}\} \end{aligned}$$

⁴Remember, this sinusoid has constant amplitude of 1 and this function only describes its change of phase.

and solving for the remaining coefficients is done using the linear system of equations:

$$\begin{pmatrix} H^5 & H^4 & H^3 \\ 5H^4 & 4H^3 & 3H^2 \\ 20H^3 & 12H^2 & 6H \end{pmatrix} \begin{pmatrix} u_5 \\ u_4 \\ u_3 \end{pmatrix} = \begin{pmatrix} -\frac{H^2}{2} \Im \{a_2^{\tau_0}\} - H \Im \{a_1^{\tau_0}\} + \Im \{a_0^{\tau_1}\} - \Im \{a_1^{\tau_1}\} + 2\pi U^* \\ -H \Im \{a_2^{\tau_0}\} + \Im \{a_1^{\tau_1}\} - \Im \{a_1^{\tau_0}\} \\ \Im \{a_2^{\tau_1}\} - \Im \{a_2^{\tau_0}\} \end{pmatrix}$$

The smoothness maximizing \tilde{U} is found as

$$\tilde{U} = \frac{1}{80\pi} [20H(\Im \{a_1^{\tau_0}\} + \Im \{a_1^{\tau_1}\}) + H^2(\Im \{a_2^{\tau_0}\} - \Im \{a_2^{\tau_1}\}) + 40(\Im \{a_0^{\tau_0}\} - \Im \{a_0^{\tau_1}\})]$$

and then rounded to produce U^* as above.

The quintic interpolating phase polynomial has been proposed in a previous paper [15] although they do not directly estimate the frequency slope, choosing instead to derive it using the difference in frequency between two analysis frames.

The amplitude part

Solving for the quintic log-amplitude polynomial

$$\rho(n) = \exp(v_5 n^5 + v_4 n^4 + v_3 n^3 + v_2 n^2 + v_1 n + v_0) \quad (5.6)$$

with $v_i \in \mathbb{R}$, is as follows:

$$\begin{aligned} \mu(0) &= \Re \{a_0^{\tau_0}\} \\ \frac{d\mu}{dn}(0) &= \Re \{a_1^{\tau_0}\} \\ \frac{d^2\mu}{dn^2}(0) &= \Re \{a_2^{\tau_0}\} \end{aligned}$$

and solving for the remaining coefficients is done using the linear system of equations:

$$\begin{pmatrix} H^5 & H^4 & H^3 \\ 5H^4 & 4H^3 & 3H^2 \\ 20H^3 & 12H^2 & 6H \end{pmatrix} \begin{pmatrix} v_5 \\ v_4 \\ v_3 \end{pmatrix} = \begin{pmatrix} -\frac{H^2}{2} \Re \{a_2^{\tau_0}\} - H \Re \{a_1^{\tau_0}\} + \Re \{a_0^{\tau_1}\} - \Re \{a_1^{\tau_1}\} \\ -H \Re \{a_2^{\tau_0}\} + \Re \{a_1^{\tau_1}\} - \Re \{a_1^{\tau_0}\} \\ \Re \{a_2^{\tau_1}\} - \Re \{a_2^{\tau_0}\} \end{pmatrix}$$

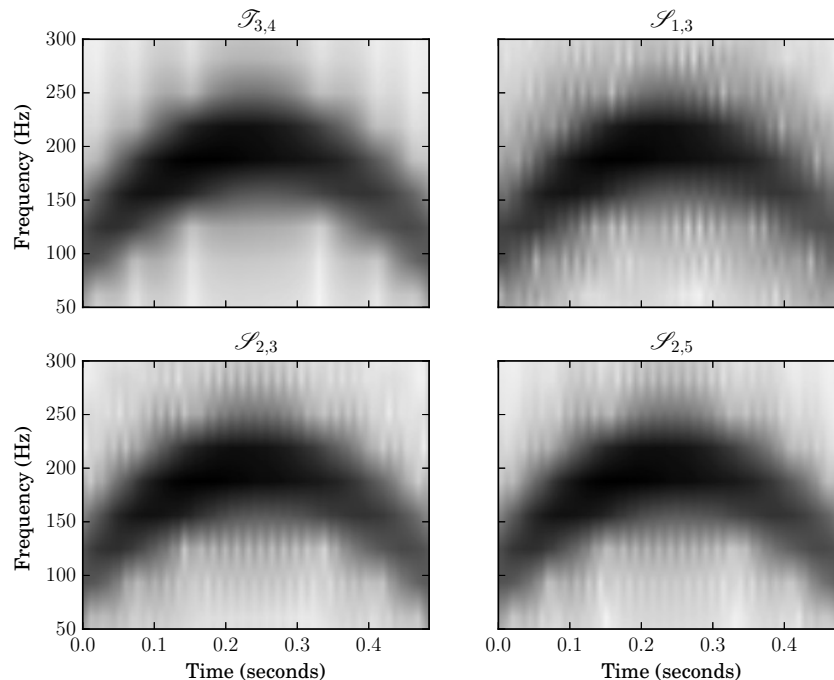


Fig. 5.1: Spectrograms of $\mathcal{T}_{3,4}$, $\mathcal{S}_{1,3}$, $\mathcal{S}_{2,3}$ and $\mathcal{S}_{2,5}$. Spectrograms of the true signal and estimated signals for the $\mathcal{T}_{3,4}$ signal.

5.5 Evaluation

We compared the quality of an analysis-synthesis system using the original $\mathcal{S}_{1,3}$ method, $\mathcal{S}_{2,3}$ method, and the $\mathcal{S}_{2,5}$. Frequency- and amplitude-modulated sinusoids were synthesized and then analysed frame-by-frame using the DDM to estimate their initial phase (amplitude), frequency (amplitude slope), and frequency-modulation (amplitude-modulation). Afterwards, the signals were resynthesized using the estimated parameters and compared to the original. We are interested in seeing in what cases higher-order phase and log-amplitude polynomials will improve the accuracy of synthesis.

5.5.1 Evaluation on a sinusoid of cubic phase and quartic log-amplitude

The initial evaluation illustrates how higher-order phase and log-amplitude polynomials will not necessarily improve the quality of synthesis if the underlying phase and log-amplitude

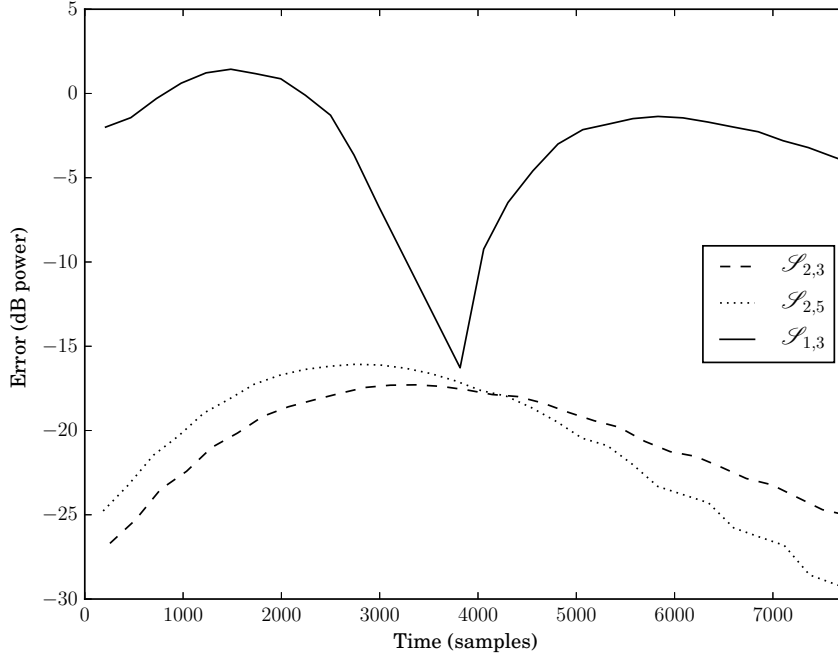


Fig. 5.2: $\mathcal{I}_{3,4}$ vs. $\mathcal{S}_{1,3}$, $\mathcal{S}_{2,3}$ and $\mathcal{S}_{2,5}$: Upper error bound. The power of the error when subtracting the original signal from the estimated signal. The local upper bound on the error was produced by connecting the local maxima in the error data.

functions are a polynomial of lower order than the polynomials used for synthesis. As we will see, the estimated phase and log-amplitude functions suffer from “overfitting”.

The synthesized signal has 3 frequency break-points and an initial phase, therefore its phase function can be interpolated by a cubic polynomial

$$x_\phi(n) = \exp(g_3 n^3 + g_2 n^2 + g_1 n + g_0)$$

the frequency break-points are summarized in Table 5.1. The initial phase is 0 radians.

Table 5.1

Time (seconds)	0	0.25	0.5
Frequency (Hz)	100	200	100

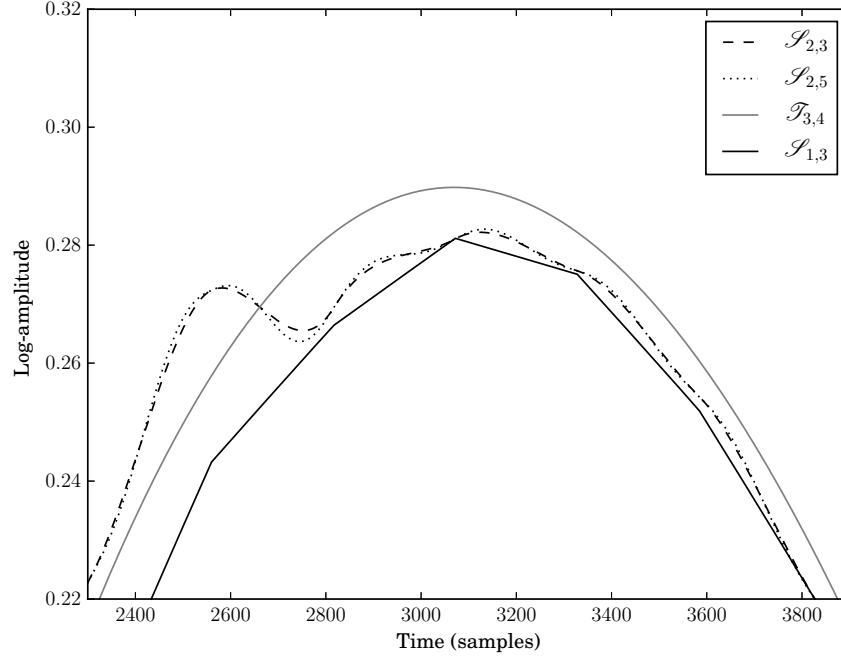


Fig. 5.3: $\mathcal{T}_{3,4}$, $\mathcal{S}_{1,3}$, $\mathcal{S}_{2,3}$ and $\mathcal{S}_{2,5}$: Log-amplitude functions. This compares the original log-amplitude function with the interpolated log-amplitude functions. The log-amplitude functions are considered because these are the real part of the polynomial exponents in the complex sinusoid model.

A quartic polynomial is used for the log-amplitude function

$$x_{\mu}(n) = \exp(h_4 n^4 + h_3 n^3 + h_2 n^2 + h_1 n + h_0)$$

and its amplitude break-points are summarized in Table 5.2. This will be referred to as the $\mathcal{T}_{3,4}$ signal.

Table 5.2

Time (seconds)	0	0.1	0.3	0.5
Amplitude (dB)	-10	0	0	-10

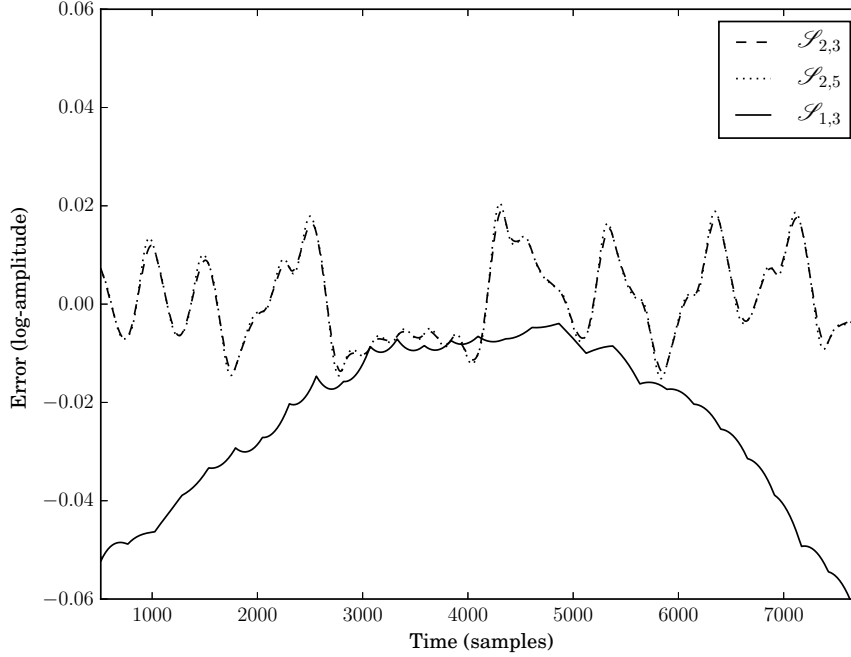


Fig. 5.4: $\mathcal{I}_{3,4}$ vs. $\mathcal{S}_{1,3}$, $\mathcal{S}_{2,3}$ and $\mathcal{S}_{2,5}$: Log-amplitude function error. This shows the error of the interpolated log-amplitude functions when compared with the original log-amplitude function for the three proposed methods.

These polynomials are chosen because their orders are greater than or equal to the order of the synthesis model in the $\mathcal{S}_{2,3}$ system and less than the order of the synthesis model in the $\mathcal{S}_{2,5}$ system.

The signal was sampled with a sampling rate of 16000 Hz and was analysed every 256 samples with an analysis window of length 1024 samples. For the DDM method, the \mathcal{C}^1 4-Term Blackman-Harris was used (see Section 3.5).

The results of the evaluation are presented in Figures 5.1 through 5.5. In Figure 5.1 we see informally that the synthesis quality is good for all model orders even though the $\mathcal{S}_{1,3}$ model assumes linear phase in its analysis. Figure 5.2 shows how accurately the different model orders reconstruct the original signal. We see that the $\mathcal{S}_{2,5}$, although of higher-order, is not systemically a better interpolator of the true underlying signal. Figure 5.3 shows the estimated log-amplitude functions along with the original and Figure 5.3 their respective errors in approximating the true log-amplitude function. Finally, Figure 5.5 shows that

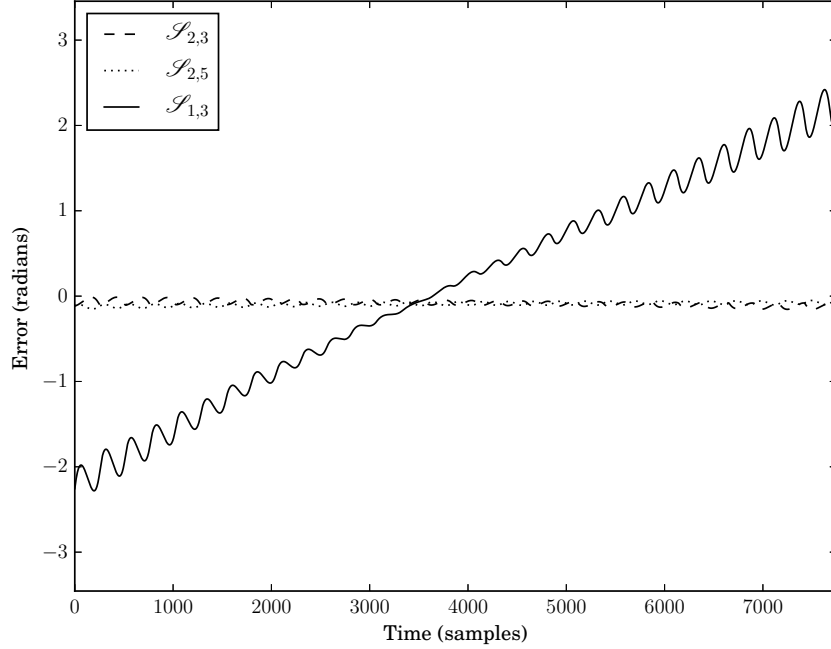


Fig. 5.5: $\mathcal{S}_{3,4}$ vs. $\mathcal{S}_{1,3}$, $\mathcal{S}_{2,3}$ and $\mathcal{S}_{2,5}$: Phase function error. This compares the theoretical phase function with the interpolated phase functions. The phase functions are considered because these are the imaginary part of the polynomial exponents in the complex sinusoidal model. The errors are “wrapped” to lie between $-\pi$ and π . The errors stem from both the estimation of the phase and the interpolation of phase between analysis points. As a frequency modulated sinusoid is considered, it is not surprising that the stationary frequency assumption of the $\mathcal{S}_{1,3}$ model exhibits the most errors.

the models incorporating the non-stationary assumption in their analysis perform better at approximating the true phase function.

5.5.2 Evaluation on sinusoid of exponential phase

The previous evaluation of this analysis-synthesis system was on a sinusoid with small-order polynomial phase and log-amplitude. In the cases where we observed overfitting, the polynomial used for synthesis was of higher-order than the true underlying one — the interpolating polynomials were more times differentiable than the true polynomial. We propose evaluating the system on an infinitely differentiable and analytic phase function. The rationale behind this stems from the definition of an analytic function: one whose

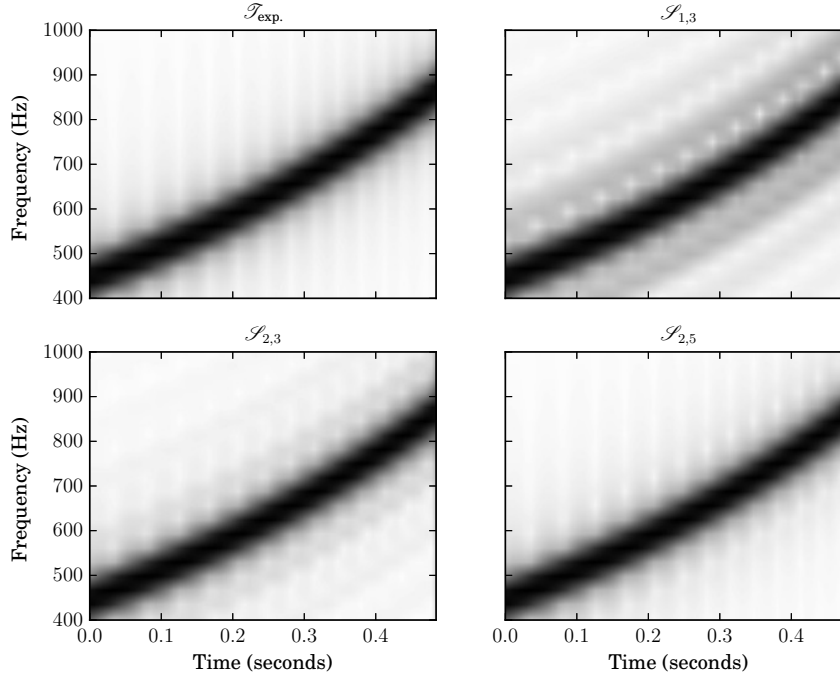


Fig. 5.6: Spectrograms of $\mathcal{T}_{\text{exp.}}$, $\mathcal{S}_{1,3}$, $\mathcal{S}_{2,3}$ and $\mathcal{S}_{2,5}$. Spectrograms of the true signal and estimated signals for the exponential phase signal.

power series representation (a polynomial) converges to the function as the number of terms approaches infinity. What this means is, in the region of convergence, the larger the number of terms in the approximating polynomial, the better the approximation to the true underlying function. The exponential function

$$y = \exp(x), x, y \in \mathbb{R}$$

is one such function whose power series is

$$\exp(x) = \sum_{n=0}^{\infty} \frac{x^n}{n!}$$

To find the radius of convergence, we use the ratio test

$$\lim_{n \rightarrow \infty} \frac{|1/n!|}{|1/(n+1)!|} = \lim_{n \rightarrow \infty} n = \infty$$

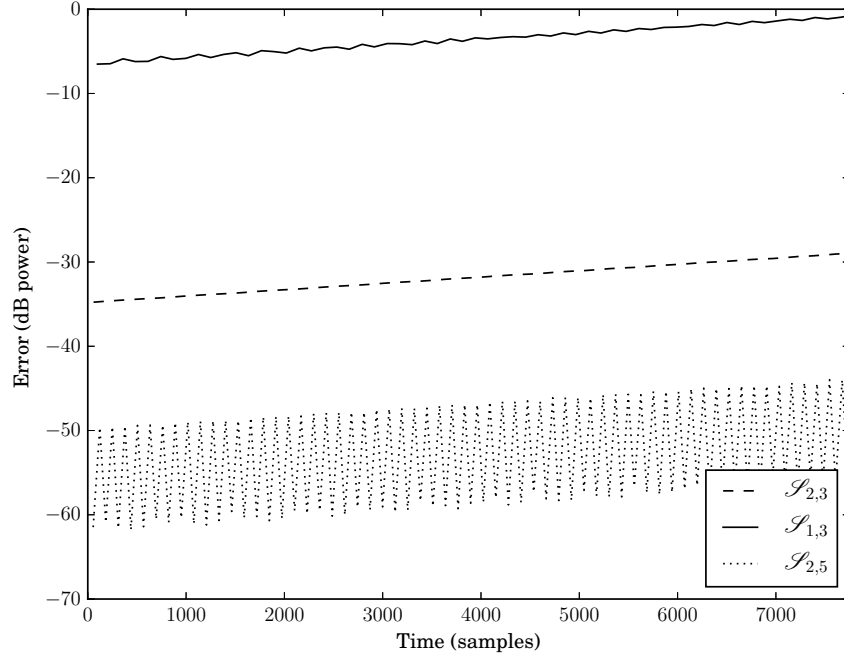


Fig. 5.7: $\mathcal{T}_{\text{exp.}}$ vs. $\mathcal{S}_{1,3}$, $\mathcal{S}_{2,3}$ and $\mathcal{S}_{2,5}$: Upper error bound. The power of the error when subtracting the original signal from the estimated signal for the signals of exponential phase. The local upper bound on the error was produced by connecting the local maxima in the error data.

i.e., the power series of the exponential function converges everywhere and so using more terms of its power series will improve its approximation for any $x \in \mathbb{R}$.

The exponential function arises in music. In 12-tone equal temperament tuning, to find the frequency f_1 of a pitch b -semitones away from the frequency f_0 we compute

$$f_1 = f_0 2^{\frac{b}{12}} = f_0 \exp(\log(2) \frac{b}{12})$$

So a linear transition from pitch b_0 to b_1 (in semitones) is an exponential change in frequency. This could be observed in recordings of performances of the *portamento* gesture.

We synthesize a sinusoid of length N samples with exponential phase and use the same analysis system as in Section 5.5.1 to evaluate the synthesis accuracy for piece-wise

interpolating polynomials of cubic and quartic order for phase. The signal x is defined

$$x(n) = \exp(j \frac{2\pi f_0}{c_1} \exp(c_1 n + c_0))$$

with

$$c_0 = \log(2) \frac{b_0}{12}$$

and

$$c_1 = \log(2) \frac{b_1}{12N}$$

i.e., a signal starting at pitch b_0 with frequency f_0 and arriving at pitch b_1 in N samples. We keep the amplitude of the signal constant in this evaluation as we are interested in the accuracy of the phase reconstruction. The same procedure as in Section 5.5.1 is used to estimate the parameters of piece-wise interpolating phase polynomials. We can see in Figure 5.7 that the resynthesis accuracy is greater for higher-order polynomials. Spectrograms of the original and estimated signals are plotted in Figure 5.6.

5.6 Conclusion

5.6.1 Polynomial phase and log-amplitude function

For the polynomial signal model considered, we observed that increased orders for analysis and synthesis models did not necessarily improve signal reproduction fidelity. For instance, we observe overfitting by the proposed cubic and quintic models in Figure 5.3. This is consistent with the results of [15]. From Figure 5.5 it is clear that the DDM-based methods provide superior estimation of the phase function — this is not the case for the log-amplitude function. Depending on the underlying signal, perhaps better results can be obtained by postulating a lower-order log-amplitude function and higher-order phase function. The possibility of errors arising from numerical accuracy when evaluating the quintic polynomials has been ruled out. We evaluated these polynomials using an implementation of Horner’s method that keeps track of the error bound [22, p. 95]: the errors are negligible, see Figure 5.8 for the results.

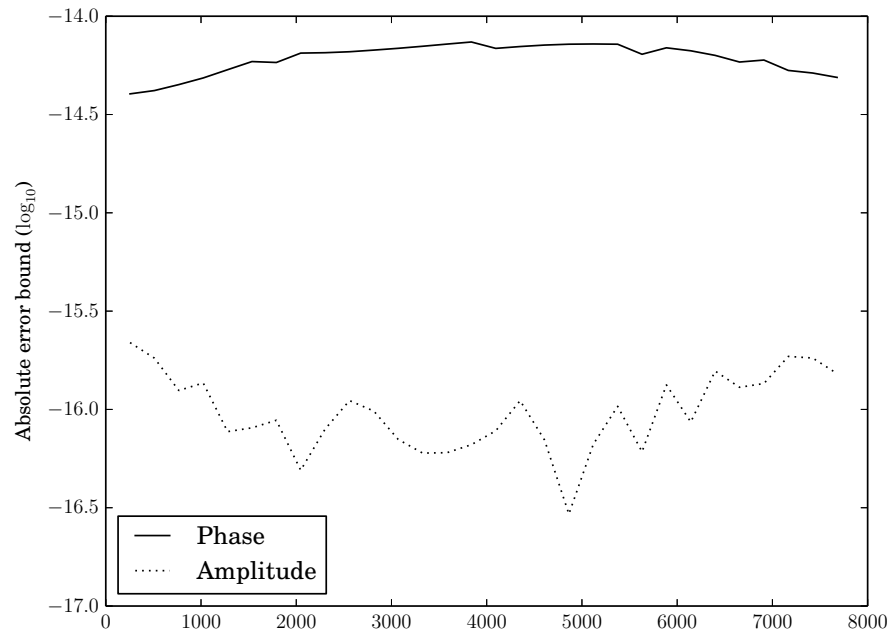


Fig. 5.8: $\mathcal{S}_{2,5}$ evaluation error bound. The error bound in evaluating the quintic log-amplitude and phase polynomials of $\mathcal{S}_{2,5}$ using Horner’s method. This plot was produced by plotting only the local maxima of the error bound data in order to reduce the plot’s range.

5.6.2 Exponential phase function

The quintic interpolation $\mathcal{S}_{2,5}$, the polynomial of highest order, performs the most accurate resynthesis. This is consistent with the analytic property of the exponential function and an encouraging result as it confirms that analytic phase functions can be approximated with arbitrary accuracy simply by increasing the order of the interpolating polynomials. Many models of musical gestures involve such functions, apart from the portamento gesture modeled by an exponential phase function, vibrato can be modeled as a sinusoid with sinusoidal phase [32]. The DDM-based analysis system combined with the higher-order polynomial phase and log-amplitude synthesis system presented here allows for accurate modeling of these gestures.

Chapter 6

Experiment: Partial grouping by amplitude- and frequency-modulation

6.1 Introduction

To evaluate whether the grouping of partials with common amplitude- and frequency-modulation (AM and FM) parameters is plausible, we synthesize a set of parameters corrupted with noise and add to this signal spurious sets of parameters that should not belong to any sources. We then attempt to retrieve the original parameters.

6.2 Methodology

We assume parameters have been estimated already so we start from theoretical values for the amplitude, frequency, frequency-modulation and amplitude-modulation. On each frame of analysis data, i.e., for parameters belonging to the same time instant, we consider each data-point as a multi-dimensional random variable. With these random variables, we compute principal components in order to produce a variable with maximum variance. This variable is classified using a clustering algorithm and we evaluate the results. A summary follows:

- Parameters are synthesized from a theoretical mixture of AM and FM sinusoids. Spurious data are added to these parameters.

- Principal components analysis is carried out on the parameters happening at one time instance.
- A histogram is made of the first principal components. Values sharing a bin with too few other values are discarded to remove spurious data points.
- Initial means and standard deviations for the GMMs (see Appendix B) are made by dividing the histogram into equal parts by area and choosing the centres of these parts.
- The EM algorithm (see Appendix B) for GMMs is carried out to classify the sources.

6.3 Evaluation

The algorithm is run on a typical source separation problem to evaluate its plausibility. Two sources are synthesized, each exhibiting both frequency- and amplitude-modulation. The amplitude and frequency of the frequency-modulation are chosen to be realistic with respect to musical instrument sounds — ± 12.5 cents surrounding the fundamental at around 6 Hz (see Table 6.2), similar to the measurements obtained in [38] for violin vibrato. Because musical sounds exhibit a wide variety of amplitude envelopes, one is chosen that is realistic, but that is not based on any particular instrument or recording. For this process all that is important to carry out source separation is the relative amplitude-modulation of the two sources — a reasonable assumption for recordings of a mixture of two different instruments or two performers of the same instrument.

To have control over the frequency- and amplitude-modulation separately, we compute the parameters of a function describing the amplitude envelope, and one describing oscillatory part. The parameters are combined when carrying out the classification.

6.4 Synthesis

Our model makes available the parameters summarized in Table 6.1. To incorporate inharmonicity often observed in real string instruments where the strings exhibit some stiffness,

H	Duration between data-point calculations in samples (i.e., the hop size).
N	Number of sources.
p	Which source.
$f_{0,p}$	Fundamental frequency.
K_p	Number of harmonics.
$k_{60,p}$	Harmonic number 60 dB lower than the first.
B_p	The inharmonicity coefficient.
$\phi_{0,p}$	Initial phase.
$\phi_{0,f,p}$	Initial FM phase.
$t_{60,p}$	Time until amplitude of partial has dropped 60 dB.
$t_{\text{attack},p}$	Time duration of attack portion.
$A_{f,p}$	Amplitude of FM.
$f_{f,p}$	Frequency of FM.
s_p	The signal representing the p th source.
$a_{60,p}$	The slope of the line in the argument of the exponential describing the amplitude variation.
$a_{k,60,p}$	The coefficient of the harmonic number in the argument of the exponential describing the initial amplitude of a harmonic as a function of its harmonic number.

Table 6.1. Synthesis parameters. Time values are in seconds, frequency values are in Hz and phase values are in radians.

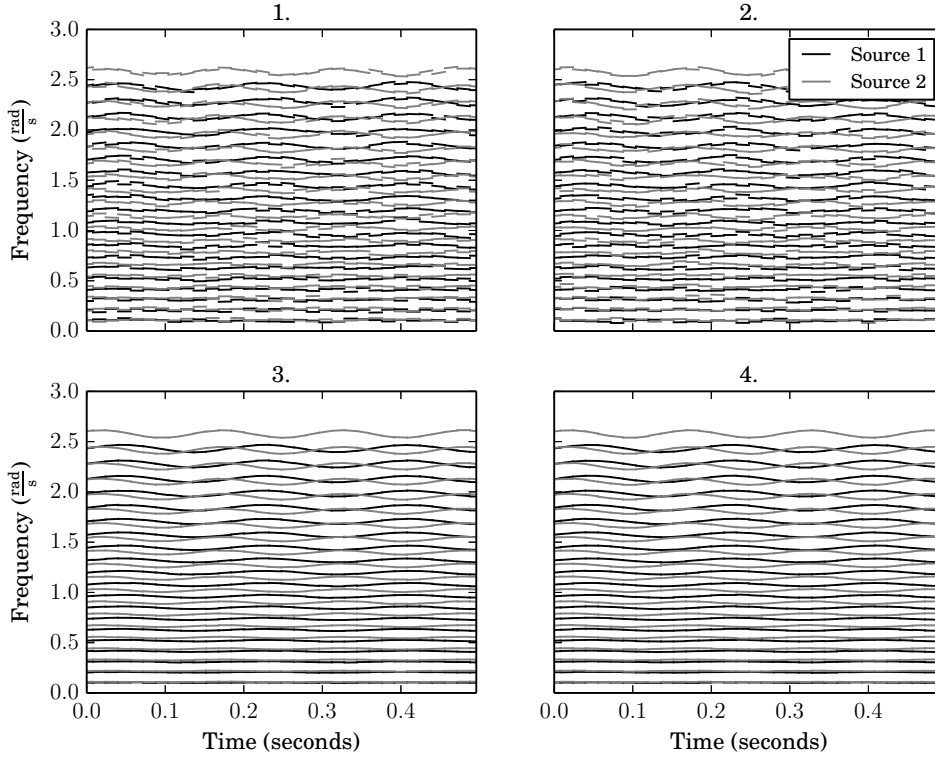


Fig. 6.1: Original data-points. Line-segments describing the frequency and frequency-modulation of both sources with no added spurious parameters. The amount of noise added and the corresponding plot title number are summarized in Table 6.3.

we define the *stretched* harmonic numbers as follows [56]¹

$$\mathcal{K}_B(k) = k(1 + Bk^2)^{\frac{1}{2}}$$

Each source is synthesized using the following equation:

$$s_p(t) = \sum_{k=1}^{K_p} A_p(k, t) \exp(j((2\pi f_{0,p}t - \frac{A_{f,p}}{f_{f,p}} \cos(2\pi f_{f,p}t + \phi_{0,f,p}))\mathcal{K}_{B_p}(k) + \phi_{0,p}))$$

¹http://ccrma.stanford.edu/~jos/pasp/Dispersion_Filter_Design.I.html

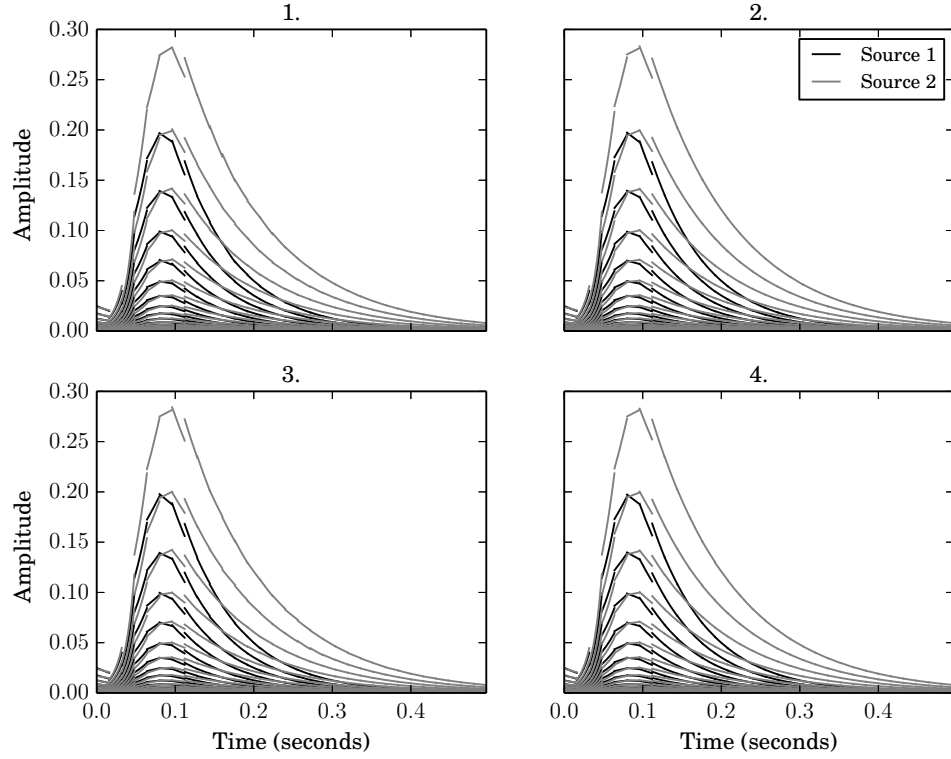


Fig. 6.2: Amplitude function for each source (true). Line-segments representing the instantaneous amplitude and amplitude slope at analysis time points. The amount of noise added and the corresponding plot title number are summarized in Table 6.3.

where

$$A_p(k, t) = \begin{cases} \exp(a_{60,p}t + a_{k,60,p}k) \cos^2\left(\frac{\pi}{2}\left(\frac{t}{t_{\text{attack},p}} - 1\right)\right) & \text{if } t \leq t_{\text{attack},p}, \\ \exp(a_{60,p}t + a_{k,60,p}k) & \text{if } t > t_{\text{attack},p}, \\ 0 & \text{otherwise.} \end{cases}$$

$$a_{60,p} = \frac{\log(10^{-3})}{t_{60,p}}$$

$$a_{k,60,p} = \frac{\log(10^{-3})}{k_{60,p}}$$

The piecewise amplitude function is based on the amplitude function of the *Formant Wave Function (FOF)*² described in [48, p. 19].

²FOF stands for *Forme d'Onde Formantique*.

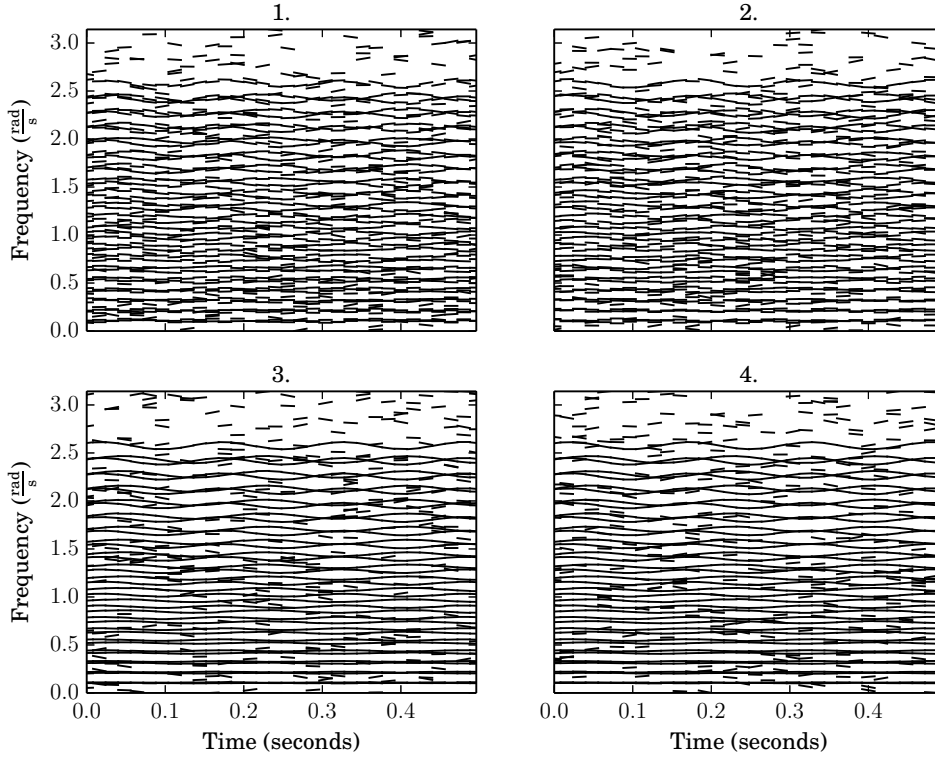


Fig. 6.3: Original and spurious data-points. Line-segments describing the frequency and frequency-modulation of the original data and the spurious data. The amount of noise added and the corresponding plot title number are summarized in Table 6.3.

6.5 Analysis

The estimation of these parameters is a separate problem addressed by the DDM (see Section 3.4). We use theoretical values calculated directly from the model signals. For interpretation, and to make it possible to simply replace the theoretical values with those obtained from an analysis, we compute parameters that correspond to a model whose parameters could be estimated through a technique such as the DDM.

For this experiment we seek signals $s_k \in \mathbb{C}$ of the following form:

$$s_k(n) = \exp(\log(A_k) + \alpha_k n + j(\phi_k + \omega_k n + \frac{1}{2}\psi_k n^2)) \quad (6.1)$$

Parameter	Source 1 value	Source 2 value
$f_{0,p}$	261.63	277.18
K_p	20	20
$k_{60,p}$	20	20
B_p	0.001	0.001
$\phi_{0,p}$	0	0
$\phi_{0,f,p}$	0	0.8
$t_{60,p}$	0.5	0.75
$t_{\text{attack},p}$	0.1	0.1
$A_{f,p}$	3.805	4.032
$f_{f,p}$	6.5	5.5

Table 6.2. Synthesis parameters for source separation by frequency- and amplitude-modulation. Sources 1 and 2 have the fundamental frequencies of a C_4 and C_4^\sharp respectively. The values for $A_{f,p}$ are found by computing $f_{0,p}2^{1/48} - f_{0,p}$, giving ± 12.5 cents of frequency-modulation centred around the fundamental frequency.

Title number	ψ_{no}	ω_{no}	α_{no}	A_{no}
1	1.0×10^{-2}	1.0×10^{-2}	1.0×10^{-2}	3.2×10^{-5}
2	1.0×10^{-2}	1.0×10^{-2}	1.0×10^{-2}	1.0×10^{-5}
3	1.0×10^{-3}	1.0×10^{-3}	1.0×10^{-3}	3.2×10^{-5}
4	1.0×10^{-3}	1.0×10^{-3}	1.0×10^{-3}	1.0×10^{-5}

Table 6.3. The plot title numbers and the amount of noise added to the synthesis parameters for that realisation.

Here n is the sample number. This is the model of a sinusoid with linear amplitude-modulation and quadratic phase-modulation. We compute from the synthesis model what these parameters would be and add noise to simulate measurement error.

Typically when performing a short-time analysis, the time corresponding to $n = 0$ is made to be the centre of the window, therefore, t is the time at the centre of the window and N_w , in samples, is the length of the middle (usually non-zero) portion of the window.

6.5.1 The amplitude signal

The coefficients describing the amplitude-slope of the k th harmonic of the p th source from our synthetic model are given by

$$\alpha_{k,p}(t) = \frac{a_{60,p}}{f_s}$$

$$A_{k,p}(t) = \exp(a_{60,p}t + a_{k,60,p}k)$$

for the part of the signal after the attack portion. Note that the amplitude-slope is not time-varying.

For the attack portion, we estimate the amplitude parameters of Equation 6.1 using least-squares on a rectangular-windowed signal³. Let

$$\hat{\mathbf{s}}_{k,p}(t_n) = \begin{pmatrix} Q(t_n - \frac{N_w}{2f_s}) \\ Q(t_n - \frac{N_w}{2f_s} + 1) \\ \vdots \\ Q(t_n + \frac{N_w}{2f_s} - 1) \\ Q(t_n + \frac{N_w}{2f_s}) \end{pmatrix}$$

with

$$Q(t) = \exp(a_{60,p}(t) + a_{k,60,p}k) \cos^2\left(\frac{\pi}{2}\left(\frac{t}{t_{\text{attack},p}} - 1\right)\right)$$

the function describing the amplitude during the attack portion of the signal. Then $\log(A_{k,p})$ and $\alpha_{k,p}$ are found as the least-squares solution of

$$\begin{bmatrix} 1 & -\frac{N_w}{2} \\ \vdots & \vdots \\ 1 & \frac{N_w}{2} \end{bmatrix} \begin{pmatrix} \log(A_{k,p}(t_n)) \\ \alpha_{k,p}(t_n) \end{pmatrix} = \log \hat{\mathbf{s}}_{k,p}(t_n)$$

where $t_n = \frac{n}{f_s}$ is the time in seconds at sample n . For the argument parameters (those multiplied by j in Equation (6.1))

$$\phi_k(t) = \left(2\pi f_{0,p}t - \frac{A_{f,p}}{f_{f,p}} \cos(2\pi f_{f,p}t + \phi_{0,f,p})\right) \mathcal{K}_{B_p}(k) + \phi_{0,p}$$

³We cannot simply compute the modulation parameters of Equation 6.1 using the Taylor-series expansion of $Q(t)$ because the attack portion is an exponential function modulated by a raised cosine, which does not match the analysis model.

$$\omega_{k,p}(t) = \frac{2\pi}{f_s} (f_{0,p} + A_{f,p} \sin(2\pi f_{f,p}t + \phi_{0,f,p})) \mathcal{K}_{B_p}(k)$$

$$\psi_{k,p}(t) = \left(\frac{2\pi}{f_s}\right)^2 A_{f,p} f_{f,p} \cos(2\pi f_{f,p}t + \phi_{0,f,p}) \mathcal{K}_{B_p}(k)$$

To simulate the noise that would be present in an estimation of the signal parameters from an arbitrary signal, we create noise corrupted values by substituting the random variables:

- $\tilde{\psi}_{k,p}(t) \sim \psi_{k,p}(t) + \mathcal{N}(0, \psi_{no})$
- $\tilde{\omega}_{k,p}(t) \sim \omega_{k,p}(t) + \mathcal{N}(0, \omega_{no})$
- $\tilde{\alpha}_{k,p}(t) \sim \alpha_{k,p}(t) + \mathcal{N}(0, \alpha_{no})$
- $\tilde{A}_{k,p}(t) \sim A_{k,p}(t) + \mathcal{N}(0, A_{no})$

The θ_{no} (where θ is replaced by ψ etc.) specifies the variance of the particular parameter. Most likely in practice these random variables would be correlated but not knowing the estimation method, we cannot at this point say anything about this correlation. Therefore the noisy parameters are uncorrelated random variables for this experiment.

We also add spurious data-points as a fraction r of the number of true data-points. Their values are drawn from uniform distributions with boundaries θ_{\min} and θ_{\max} , where θ is some parameter above, e.g., ω_{\min} and ω_{\max} for the ω parameter. For this experiment $r = 0.25$, which is quite a large number of spurious points. This value is chosen to show that, given an acceptable accuracy of estimation of the true parameters, good source separation results can be achieved, even with such a high proportion of spurious points. The parameters of the uniformly distributed random variables are given in Table 6.4. Data-points are computed for the times $t_n = 0, \frac{H}{f_s}, \frac{2H}{f_s}, \dots, \frac{\lfloor \frac{N}{H} \rfloor H}{f_s}$.

Table 6.4. Distribution parameters of uniformly distributed random variables

Parameter	θ_{\min}	θ_{\max}
ω	0	π
ψ	-1×10^{-4}	1×10^{-4}
α	-1×10^{-3}	1×10^{-3}

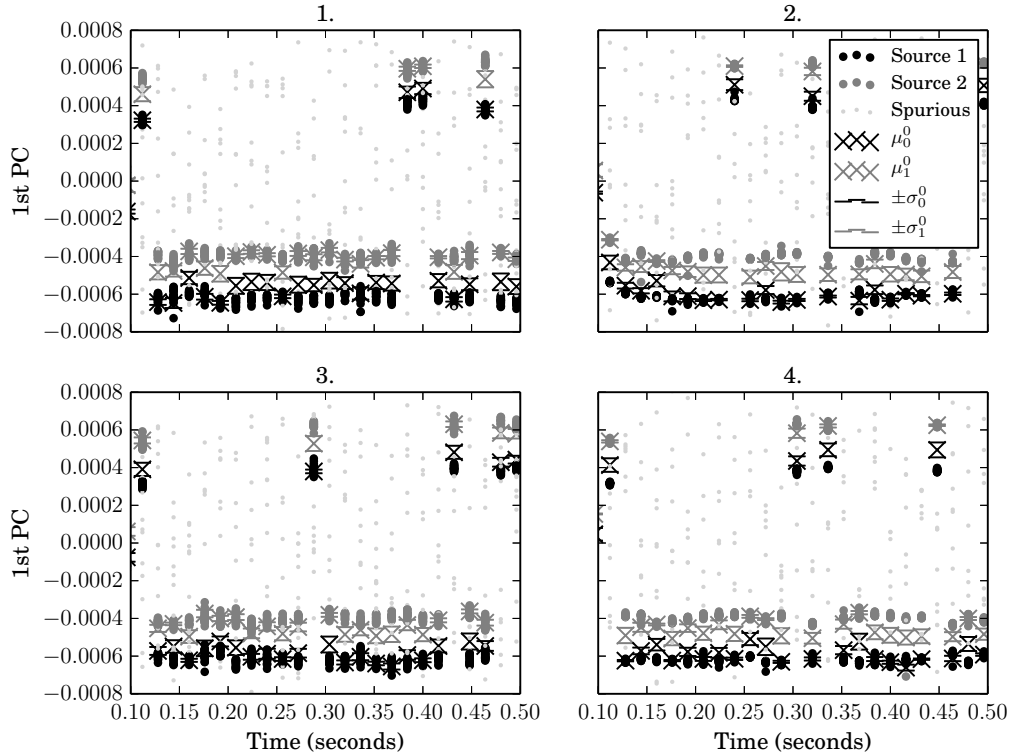


Fig. 6.4: Principal components and their classification. The PCs for each theoretical analysis time-point. μ_i^0 and σ_i^0 are respectively the initial mean and standard deviation guesses for the EM algorithm fitting the GMM parameters to the i th source. These values are also visible on the plot. The spurious points rejected using the process described in Section 6.7 are included for comparison. The amount of noise added and the corresponding plot title number are summarized in Table 6.3.

6.6 Computation of principal components

At each time t we have L data-points. As the source of each data-point is now unknown, we replace the k and p indices with index l . We only consider the amplitude- and frequency-modulation. According to our model, the frequency-modulation is greater for harmonics of greater centre frequency. To take this into consideration, we divide the frequency-modulation estimate $\psi_l(t)$ by the constant frequency estimate $\omega_l(t)$. This is similar to the approach taken in [6]. The amplitude-modulation $\alpha_l(t)$ remains constant for all harmonics of the same source⁴, only its initial value changes according to $k_{60,p}$. We compile the data-points at one time into a set of observations.

$$\mathbf{x}_l(t) = \begin{pmatrix} \frac{\psi_l(t)}{\omega_l(t)} \\ \alpha_l \end{pmatrix}$$

$$\mathbf{X}(t) = [\mathbf{x}_1(t) \dots \mathbf{x}_L(t)]$$

From these L observations the correlation matrix \mathbf{S} is computed⁵. We use the correlation matrix because the values in each row of $\mathbf{x}_l(t)$ do not have the same units, see [24, p. 22] for a discussion about this.

Following the standard technique for producing principal components (see Appendix A and also [24, p. 11]), we obtain a matrix $\mathbf{V}(t)$ of eigenvectors sorted so that the eigenvector corresponding to the largest eigenvalue is in the first column, etc. The principal components $\mathbf{A}(t)$ are then computed as

$$\mathbf{A}(t) = \mathbf{V}^T(t)\mathbf{X}(t)$$

⁴We acknowledge that this might not be realistic for all sounds. If the amplitude-modulation is a function of (normalized) frequency as well as time $\beta_l(t) = \mathcal{A}(\omega)$ we need only perform the transformation $\alpha_l(t) = \mathcal{A}^{-1}(\beta_l(t))$ to obtain the same data points as classified here.

⁵If we have N samples of random variables X_i and X_j , the entry in the i th row and j th column of correlation matrix \mathbf{S} is their estimated correlation, i.e.,

$$S_{i,j} = \frac{\sum_{n=1}^N (x_{i,n} - \bar{x}_i)(x_{j,n} - \bar{x}_j)}{\sqrt{\sum_{n=1}^N (x_{i,n} - \bar{x}_i)^2 \sum_{n=1}^N (x_{j,n} - \bar{x}_j)^2}}$$

where $\bar{x}_i = \frac{1}{N} \sum_{n=1}^N x_{i,n}$, the sample mean [21, p. 66].

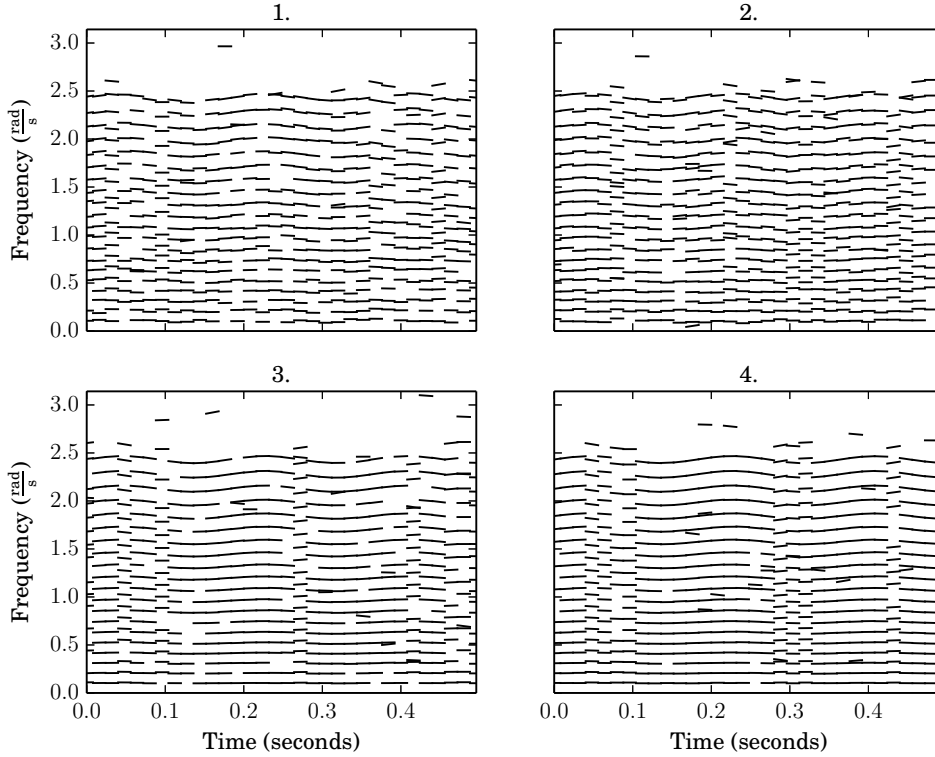


Fig. 6.5: Source 1 (estimated). Line-segments classified as belonging to Source 1. The classification is done on each frame so classifications in consecutive frames may not belong to the same true source. This is because the ordering of the clusters in each frame in Figure 6.4 is not predictable. The amount of noise added and the corresponding plot title number are summarized in Table 6.3.

We have found it sufficient to use only the first principal component and therefore only use the values in the first row of $\mathbf{A}(t)$. The i th principal component of sample l at time t is written $a_{i,l}(t)$.

If we see the $\mathbf{x}_l(t)$ as realizations of a random variable, the above computation of principal components has the effect of projecting realizations of $\mathbf{x}_l(t)$ to points $a_{1,l}(t)$ on a 1-dimensional subspace. It is a fundamental theory of principal components that the transformation above maximizes the expected Euclidean distance between the points $a_{1,l}(t)$. This is desirable for the current problem because it will always produce a variable emphasizing the parameter with the most variance. More specifically, if a scatter plot of the frequency-modulation measurements shows multiple distinct clusters whereas the

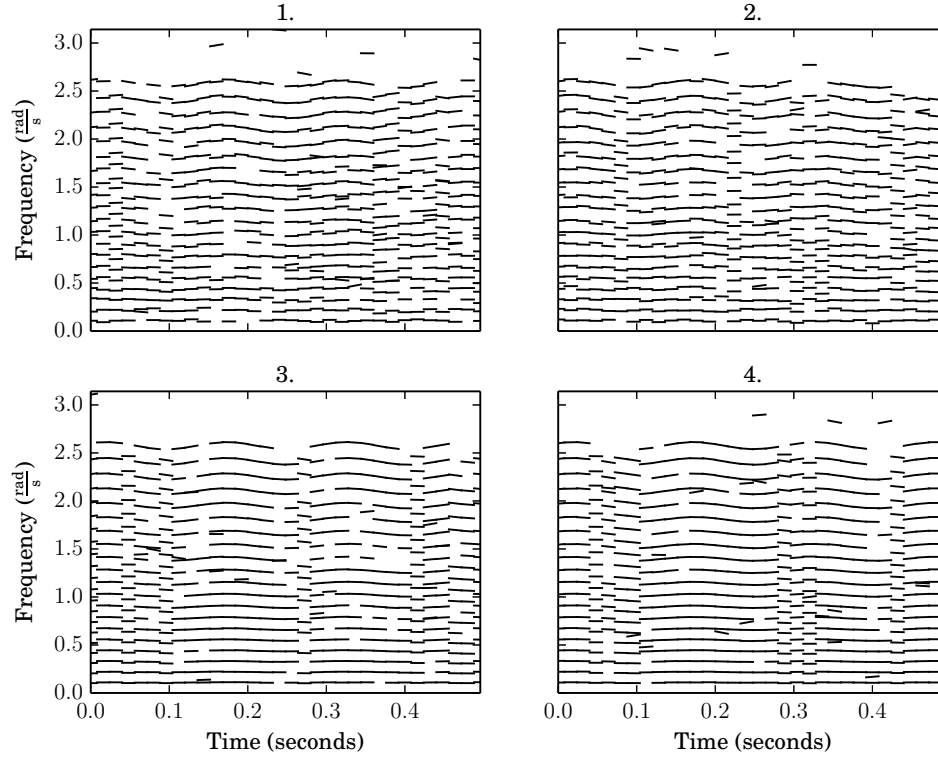


Fig. 6.6: Source 2 (estimated). Line-segments classified as belonging to Source 2. See Figure 6.5 for more information. The amount of noise added and the corresponding plot title number are summarized in Table 6.3.

amplitude-modulation measurements are all close and show only one cluster in a scatter plot, the first PC will emphasize the frequency-modulation measurements, which we desire for ease of clustering. The drawback of this approach is that if one parameter is very noisy and the other is not, the noisy parameter will be emphasized but forming informative clusters will be difficult. In that case it would be better to reject this parameter or use more PCs on which to perform clustering.

6.7 Preparing data for clustering

The EM underlying the Gaussian mixture model parameter estimation can converge to a local maximum [8], therefore, for the best results, we compute a good initial guess and remove obvious outliers before carrying out the clustering algorithm.

The $a_{1,l}(t)$ are compiled into a histogram of N_b bins. The minimum and maximum bin boundaries are computed from the maximum and minimum values of $a_{1,l}(t)$ respectively. Values in a bin with less than λ_h other values are discarded. We find N contiguous sections of equal area in the new histogram omitting the discarded values. We use the centres of these sections as the initial mean guesses and half their width as the distance 3 standard deviations from the mean (roughly 99.7 percent of values drawn from one distribution will lie within this interval if they follow a normal distribution). The initial guesses for the weights are simply $\frac{1}{N}$.

6.8 Clustering

GMM parameter estimation is discussed in Section B. After convergence we have an estimated probability $p(a_{1,l}(t))$ from distribution p . We choose the distribution p for each $a_{1,l}(t)$ that gives the highest probability of it having occurred. The values $\mathbf{x}_t(t)$ corresponding to the $a_{1,l}(t)$ have this same classification. Those sharing the same classification can be interpreted as coming from the same source. The figure shows the results of the above steps carried out on a mixture of two sources synthesized with the parameters summarized in Table 6.2. The length of the signal N is 8000 samples and the analysis hop size H is 256 samples.

6.9 Results

Here we show source separation results for the synthesized signals with noise added to the synthesis parameters ψ , ω , α and A . Each plot shows four realisations with varying amounts of noise. The amount of noise added and the corresponding plot title number are summarized in Table 6.3.

The Figures 6.1, 6.3, 6.4, 6.5, 6.6 and 6.9 summarize the results of the source separation experiment before the smoothing step. Figure 6.1 shows a time-frequency representation of the original data, Figure 6.3 shows the original data with spurious data added, Figure 6.4 shows the principal components at each time frame and the initial guesses for the EM algorithm, Figure 6.5 shows the initial classifications for source 1 and Figure 6.6 the initial classifications for source 2 before the smoothing step.

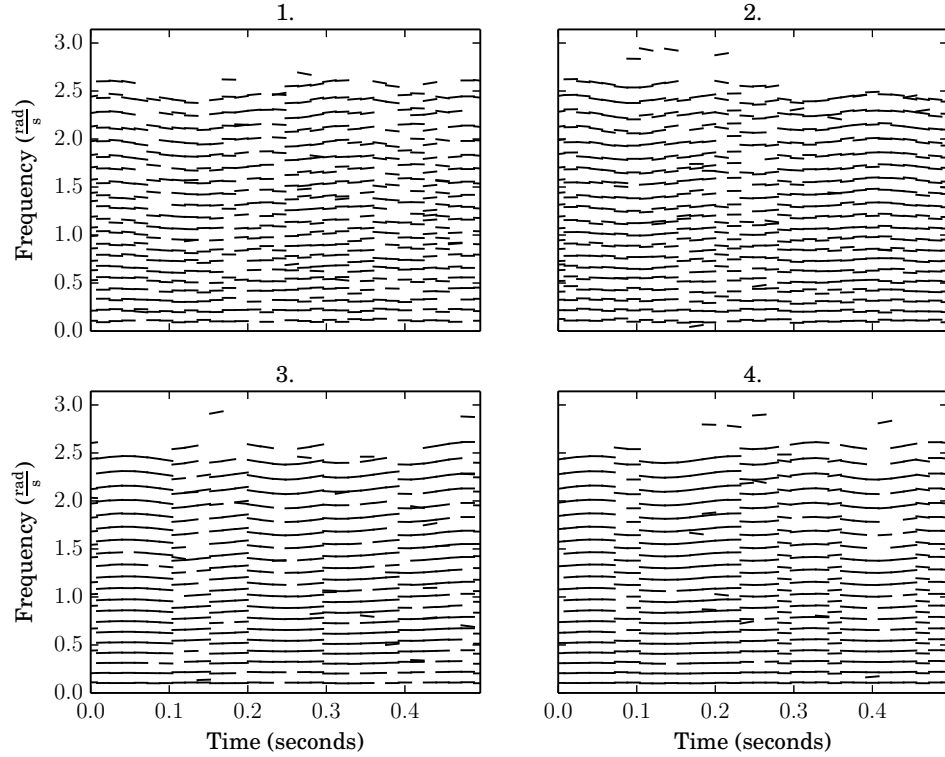


Fig. 6.7: Source 1 (estimated) after smooth amplitude path search. Line-segments classified as belonging to source 1 after smoothing in amplitude using \mathcal{D}_a . The amount of noise added and the corresponding plot title number are summarized in Table 6.3.

As seen in Figure 6.5 and 6.6, while classified well in individual frames, the overall classification does not always correspond to a single source. We must find a collection of frames with high plausibility of belonging to one source. We consider the collections of classified data-points corresponding to each source as a node in a lattice. Each frame of the lattice contains two nodes, one for each source. A best path through the lattice should connect together those nodes belonging to a single source. We use the results of Section 4.2 to find the two best paths through this lattice. We compare two distance metrics for the cost function.

The first prefers smoothness in frequency between two frames. For frame h with initial classification \tilde{p} we have frequency measurements $\omega_{k,\tilde{p}}^h$ and frequency slope measurements $\psi_{k,\tilde{p}}^h$. The set of parameters at time h from initially classified source \tilde{p} we will denote $\theta_{\tilde{p}}^h$. Between frame h and frame $h + 1$ we use Algorithm 1 on the pairs $\{\theta_{\tilde{m}}^h, \theta_{\tilde{n}}^{h+1}\}$ with

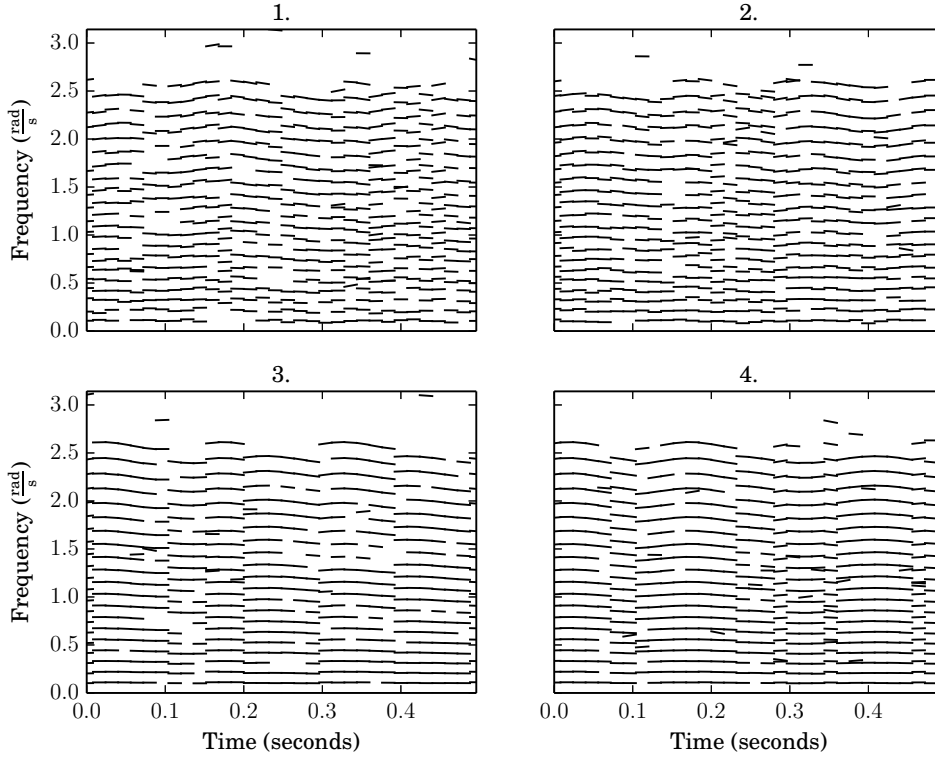


Fig. 6.8: Source 2 (estimated) after smooth amplitude path search. Line-segments classified as belonging to source 2 after smoothing in amplitude. The amount of noise added and the corresponding plot title number are summarized in Table 6.3.

$(m, n) \in \{0, 1\} \times \{0, 1\}$. For each pair, L is set to $\min(\#\theta_{\tilde{m}}^h, \#\theta_{\tilde{n}}^{h+1})$.⁶ The cost function is the absolute error in predicting the frequency in the next frame from parameters in the current frame, i.e.,

$$\mathcal{D}_f(\theta_{i,\tilde{m}}^h, \theta_{j,\tilde{n}}^{h+1}) = |\omega_{i,\tilde{m}}^h + \psi_{i,\tilde{m}}^h H - \omega_{j,\tilde{n}}^{h+1}|$$

where H is the hop-size in samples between the two frames. The second distance metric measures the smoothness in amplitude between two frames by predicting the next frame's amplitude parameters using the amplitude and amplitude-modulation parameters of the current frame. It is given as

$$\mathcal{D}_a(\theta_{i,\tilde{m}}^h, \theta_{j,\tilde{n}}^{h+1}) = |\log(A_{i,\tilde{m}}^h) + \alpha_{i,\tilde{m}}^h H - \log(A_{j,\tilde{n}}^{h+1})|$$

⁶Here, the threshold parameter $\Delta = \infty$, i.e., a connection of any cost is possible.

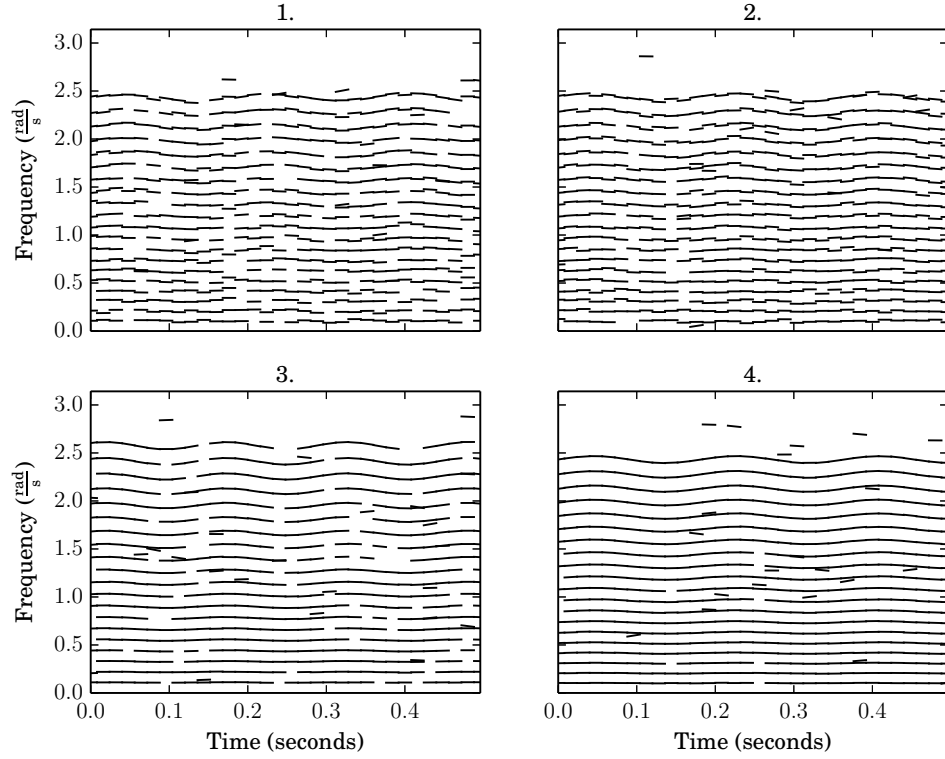


Fig. 6.9: Source 1 (estimated) after smooth frequency path search. Line-segments classified as belonging to source 1 after smoothing in frequency. The amount of noise added and the corresponding plot title number are summarized in Table 6.3.

We have found the absolute error to give better results than the squared error.

The costs of these connections are summed over the index pairs Γ_h to give the entries of the cost vector \mathbf{c} in the LP⁷

$$c_{4h+2m+n} = \sum_{i^*, j^* \in \Gamma_h} \mathcal{D}(\theta_{i^*, \tilde{m}}^h, \theta_{j^*, \tilde{n}}^{h+1})$$

The specification of the constraint matrices is done according to the topology of the lattice and the requirement that we find 2 non-overlapping paths (see Section 4.2). An example of

⁷The indexing of \mathbf{c} is explained as follows. There are 4 possible classification connections between frame h and $h + 1$. Each source m at time h has a cost of being associated with a source n at time $h + 1$, which is stored in the index $4h + 2m + n$. This is merely how the indices are laid out in the array representing the cost vector \mathbf{c} . See Section 4.2 for more about \mathbf{c} .

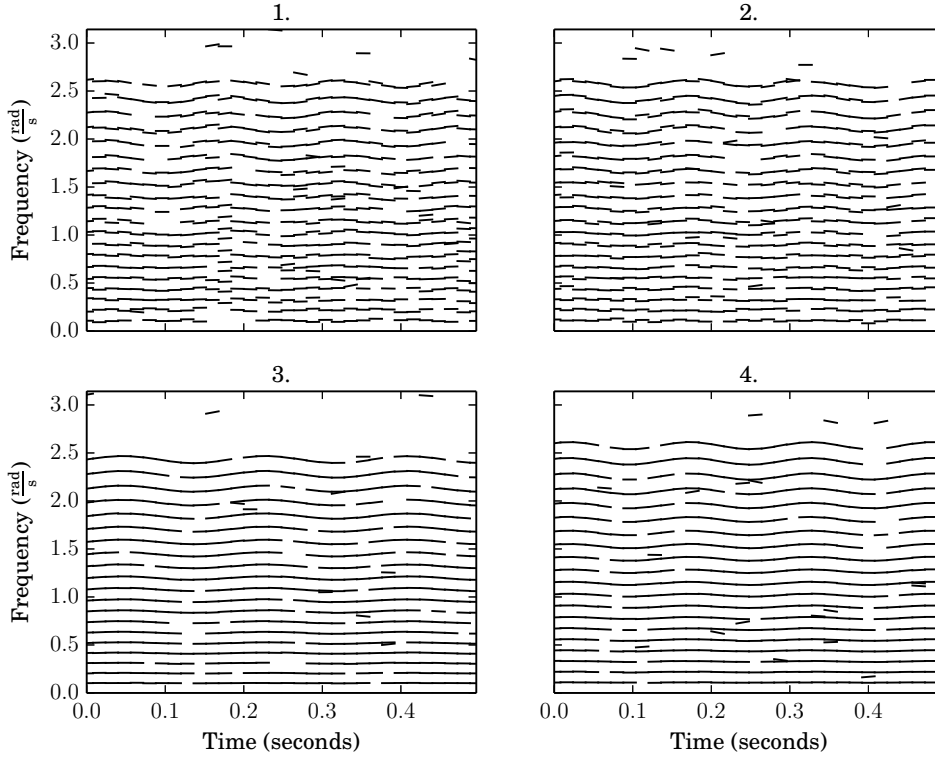


Fig. 6.10: Source 2 (estimated) after smooth frequency path search. Line-segments classified as belonging to source 2 after smoothing in frequency using \mathcal{D}_f . The amount of noise added and the corresponding plot title number are summarized in Table 6.3.

discovered paths is given in Figure 6.11. The estimated sources after smoothing in frequency using \mathcal{D}_f are shown in Figures 6.9 and 6.10. The estimated sources after smoothing in amplitude using \mathcal{D}_a are shown in Figures 6.7 and 6.8. We see that when smoothed in frequency, the results are acceptable. However, when both sets of parameters are close and give close costs, the spurious data-points can influence the cost function causing a false classification. This difficulty is not surprising, looking at Figure 6.1 we see that there are some segments where the frequency slopes are close.

When smoothed in amplitude, the results are less convincing. This is not surprising as smoothness in amplitude is not the best criterion at all time points. In Figure 6.2 we see that the amplitudes of both sources are similar at many points, e.g., at around 0.05 and 0.15 seconds.

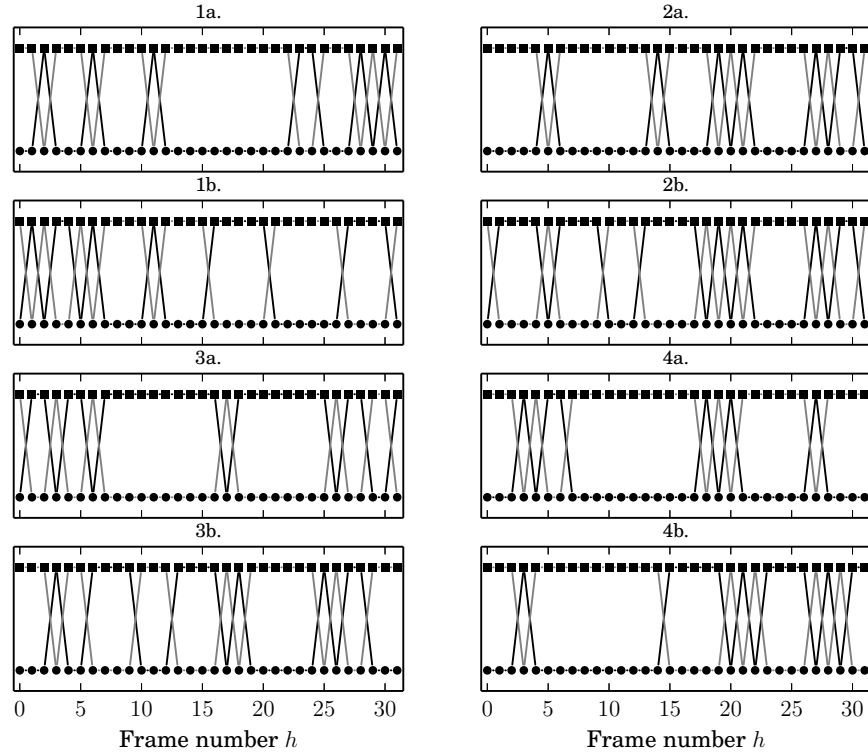


Fig. 6.11: Smoothed paths. The points originally classified as source 1 are marked with circles and those originally classified as source 2 squares. The paths in black or grey connect the points for source 1 or source 2 respectively with optimal smoothness. Plots indicated with an *a* are the paths with frequency smoothness as the criterion and those indicated with *b* are with amplitude smoothness as the criterion. The amount of noise added and the corresponding plot title number are summarized in Table 6.3.

6.10 Conclusion

In this chapter we evaluated the plausibility of separating two mixed sources based on their theoretical frequency- and amplitude-modulation. We obtained acceptable results for signals with small measurement errors. The method is also robust in the presence of spurious data points. A shortcoming of the method is the requirement that the frequency and frequency-modulation of the signals be known. Although for this experiment synthetic data were used, if the signals are sufficiently separated in frequency and have small bandwidth, as shown in Section 3.4, the DDM can be used to estimate these parameters. There are also techniques for estimating amplitude- and frequency-modulation that were

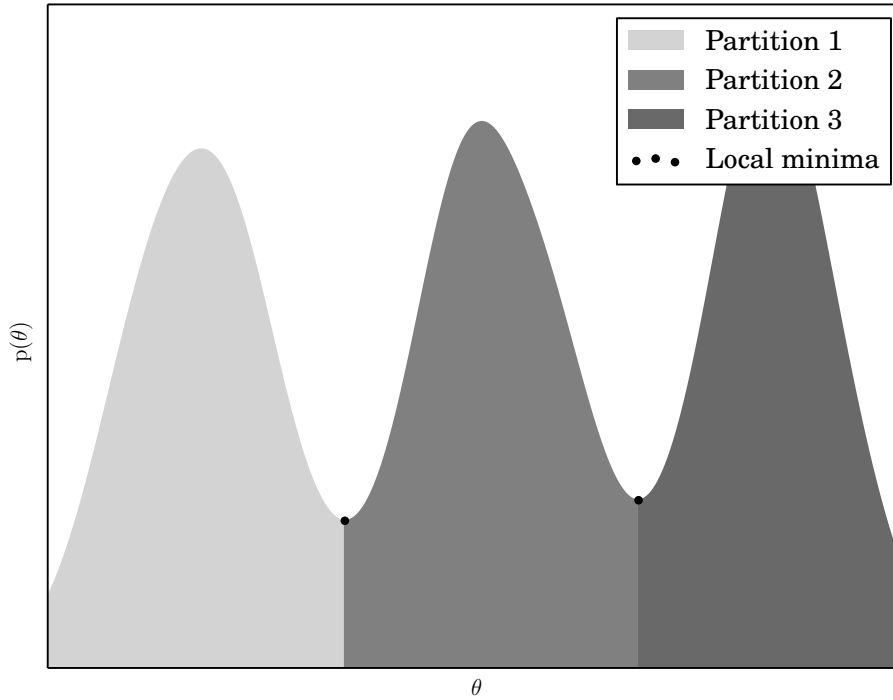


Fig. 6.12: Partitioning example. The data points are convolved with some smoothing kernels giving a function with a small number of extrema. The minima, indicated by circles, are used as boundaries between the partitions which are illustrated with different shades of grey.

not explored in this thesis. If signals are close in frequency, but the number of partials is known, and these exhibit slow modulations, signal subspace methods could be used [50] where the estimations at different time points are connected as in Chapter 4 and the modulation parameters postulated via interpolation similarly to Section 3.3.1. If it is possible to make uncorrupted measurements of the two signals and identify where the signals are not easily identified (e.g., the locations that their partials cross in the time-frequency plane) a strategy might be to use two measurements of one source and extrapolate the parameters of the signal in the part corrupted by the other source. A similar approach is explored in [27] where the Hough transform is used to identify crossing partials. Another shortcoming of the technique presented here is the use of the costly EM algorithm to classify data points using GMM (see Appendix B). A more ad hoc approach could be taken to save on

these computations, perhaps partitioning the data sets using local minima as illustrated in Figure 6.12. In any case, the source separation technique presented here, being iterative, is of a complexity similar to NMF or PLCA but can also resolve the phases of the sinusoids which are discarded in most NMF or PLCA implementations⁸.

⁸See [4] for an approach that does take into consideration the phase information in the spectrogram.

Chapter 7

Experiment: Separation of two sources using partial decay rate

7.1 Introduction

In this section we demonstrate how the techniques described above can be used to perform audio source separation on signals obtained from recordings of acoustic instruments. Specifically, we show that in the absence of frequency-modulation, amplitude-modulation — the decay rate — can be used to classify partials in a mixture of two sources into two groups, each group representing an underlying source.

7.2 Description of problem

We start with a recording of an acoustic guitar playing A_3 and a xylophone playing F_4^\sharp . The recordings are from [43] and have been mixed down to one channel (by simply adding the two signals together) and resampled at 16 kHz, coded simply as a stream of 64-bit floating-point numbers. Spectrograms of the original signals are shown in Figure 7.1 and Figure 7.2. The spectrograms were produced with a Hann window, DFT size of 4096 samples and a hop size of 512 samples. We see that neither source exhibits much frequency-modulation. The spectrogram of the mixture can be seen in Figure 7.3 and the partial paths in Figure 7.4.

The mixture of the two signals was analysed using the DDM for finding the coefficients of a cubic complex phase polynomial. Local maxima in each frame were found using the technique described in [52, p. 42]. For each of these local maxima, the polynomial

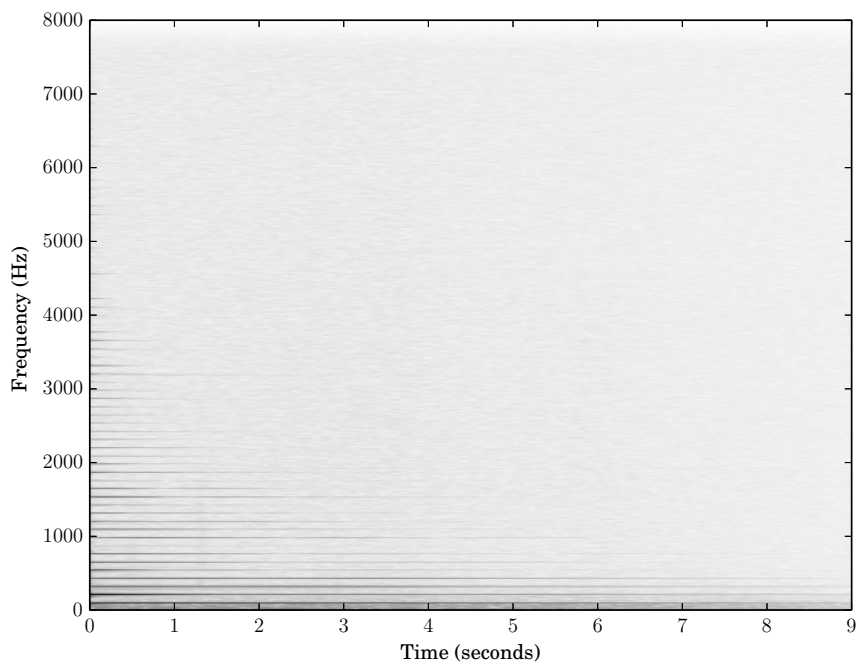


Fig. 7.1: Spectrogram of acoustic guitar. The fundamental is A_3 or 220 Hz.

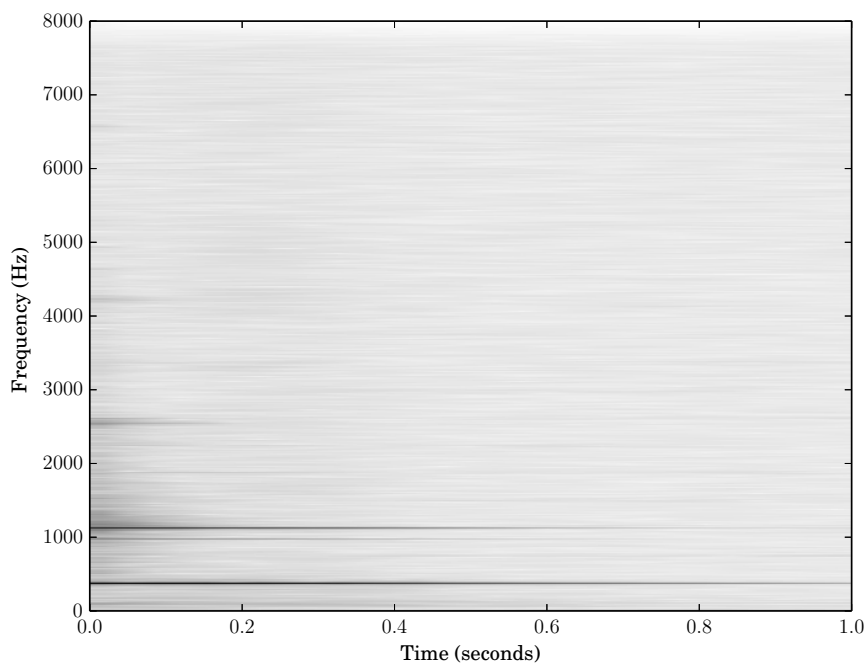


Fig. 7.2: Spectrogram of xylophone. The fundamental is F_4^\sharp or approximately 370 Hz.

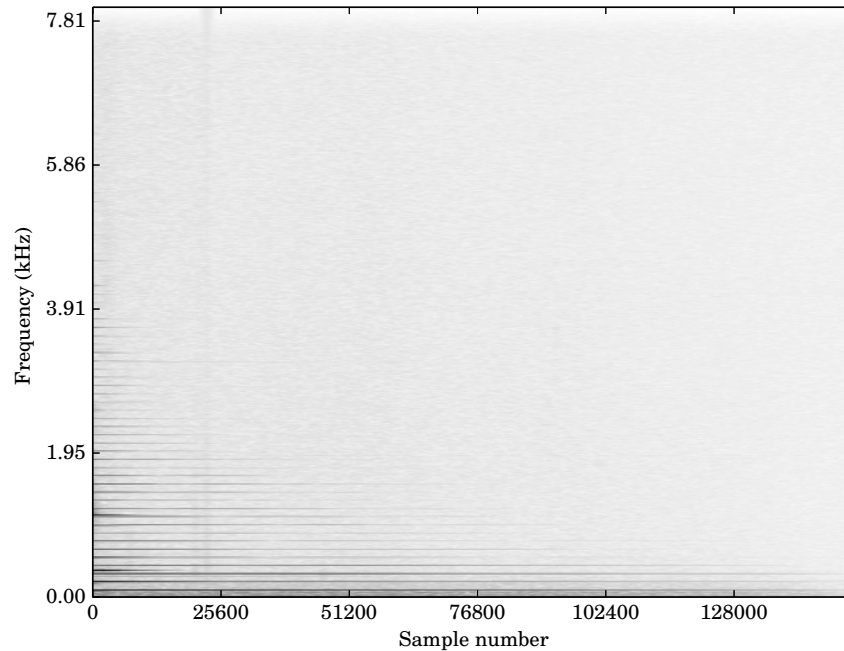


Fig. 7.3: Spectrogram of mixture. This spectrogram is the sum of the spectrograms in Figure 7.1 and 7.2.

coefficients were estimated. The analysis used the \mathcal{C}^1 4-Term Blackman-Harris window that was designed in Chapter 3. To obtain partials it was then necessary to connect the local maxima. As the partials of these two sound sources are quite stable in frequency it sufficed to use the Viterbi algorithm [13] and the cost metric $\mathcal{D}_{\text{pr.}}$ from Section 4.3 to connect local maxima in sub-bands of the spectrum. The cost function is simply the Euclidean distance between the frequencies of two local maxima. Partial starting points are considered in sub-bands of width 15 Hz and these sub-bands overlap by 7.5 Hz. A partial path starts on the first local maximum in the band exceeding -100 dB and ends at the last maximum exceeding -100 dB. The path search algorithm will also look ahead to further frames if no maximum is present in the next frame. Because of this, sometimes unrealistic paths are discovered that jump between spurious maxima. These are filtered out by discarding paths whose cost-length ratio is excessive. See Figure 7.6 to see a plot of these values and the thresholding function.

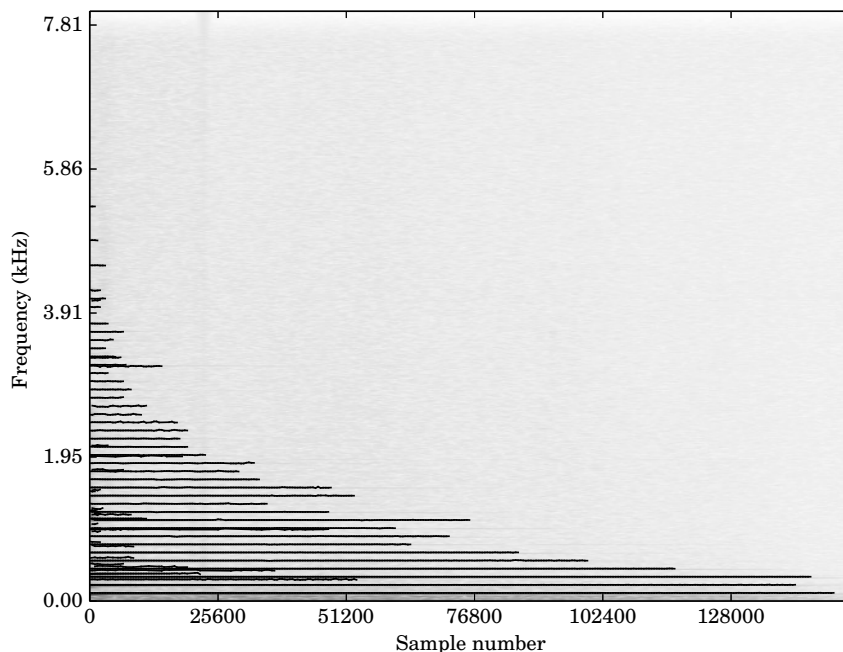


Fig. 7.4: Spectrogram of mixture and partial trajectories

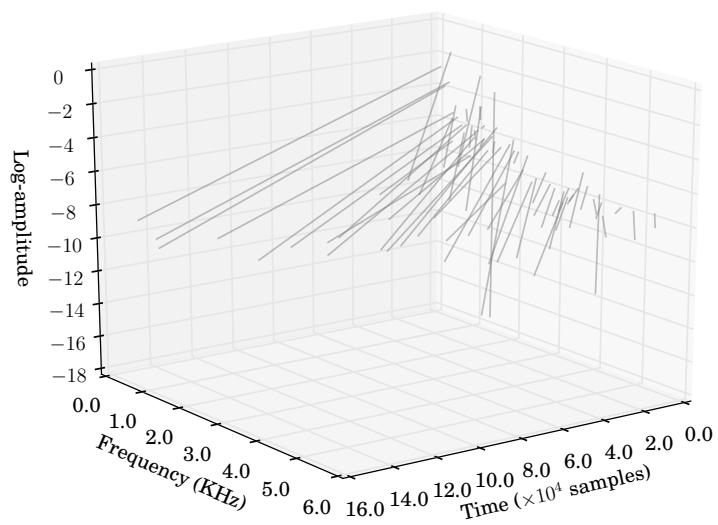


Fig. 7.5: Partial trajectories. Line functions are fit to the partial trajectory data of each partial to examine their general amplitude slopes.

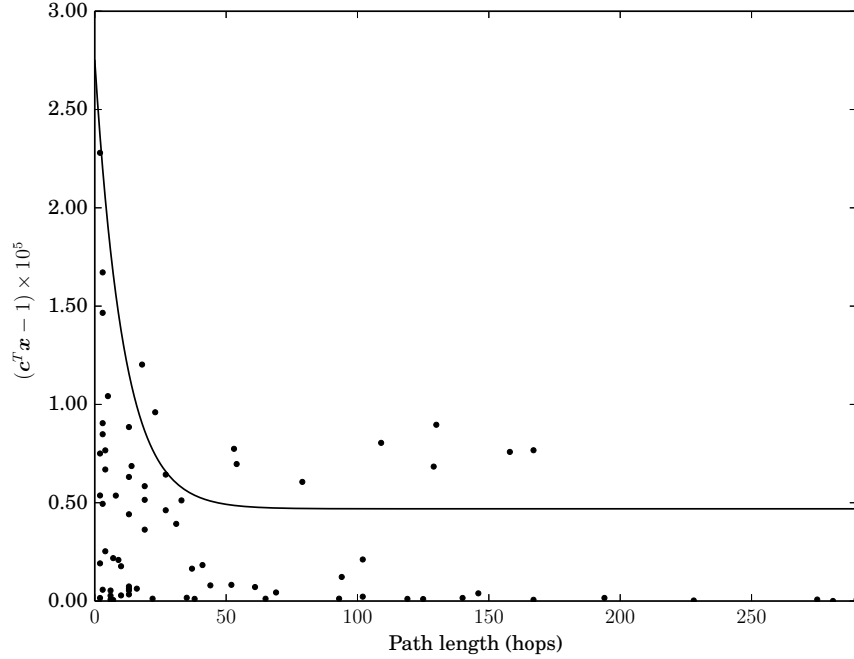


Fig. 7.6: Path cost vs. length and thresholding boundary. Paths, represented as circles, are considered only if their path cost is smaller than the threshold function, the black curve. The threshold is higher for shorter partials as to not reject those that represent the transient region of the sound. The partials during this time typically have rapidly changing frequency- and amplitude-modulation, so their path costs could be disproportionately high.

7.3 Motivation

Line functions (functions of time) are fit to the amplitude and frequency data on each partial trajectory via least-squares, as shown in Figure 7.5. Here roughly two kinds of partial slope with respect to amplitude are observed — those that are steep and brief and others that are longer and more gradual. Our goal is to classify based on the amplitude-modulation of each partial, or to an approximation, the slope of these line functions. We found that examining the log-length of the partials gives better results than examining the slope directly. This is perhaps because the log-length encodes both the starting amplitude and the slope. Recall that the partials start on the first local maximum exceeding an amplitude threshold — those with lower starting amplitude and steeper amplitude slope will be shorter, while those with a higher starting amplitude and shallower slope will be

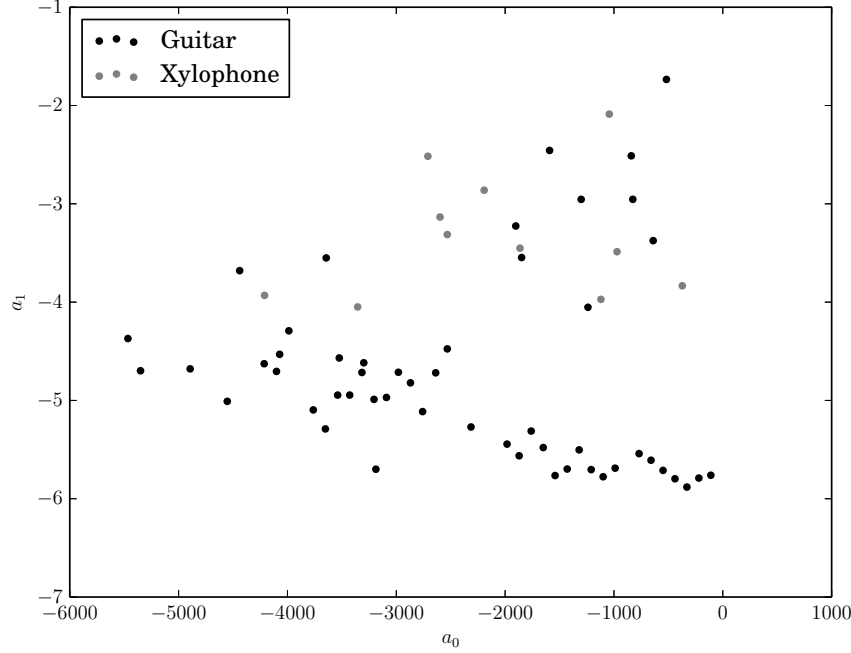


Fig. 7.7: Log-partial-length vs. frequency: principal components. This shows the distribution of partials when plotting their two principal components derived from their mean frequency and log-length. In this case, the source memberships of the partials are known. We see that there is generally a separation of the partials into two clusters corresponding to the two sources.

longer. We see in Figure 7.7 that using both the amplitude slope and the initial amplitude of the partial gives clear separation in a plot of the log-length vs. the average frequency of the partials. These partials are from separate overlaid analyses of the guitar and xylophone signals. The experiment uses an analysis of a signal consisting of a mixture of the sources, of course.

The data-points have the form

$$\mathbf{a}_i = \begin{pmatrix} a_{i,0} \\ a_{i,1} \end{pmatrix}$$

where $a_{i,0}$ is the first principal component and $a_{i,1}$ the second and are computed via a linear transformation of

$$\mathbf{x}_i = \begin{pmatrix} \bar{f}_i \\ \ell_i \end{pmatrix}$$

where \bar{f}_i is the mean frequency of the i th partial and ℓ_i its log-length (see Appendix A for the computation of principal components). The set of principal components will be denoted $\{\mathbf{a}\}$. We see that, for the most part, the partials belonging to the two sources are separated appropriately into two clusters. The partials from the xylophone present in the guitar cluster belong to higher partials, whose omission in the final rendering of the xylophone source would not be detrimental to its perceptual quality. Similarly, partials belonging to the guitar present in the xylophone cluster are short and most likely belong to briefly excited modes of the guitar body.

7.4 Classification

Our intention is now to use GMM (see Appendix B) on a set of unclassified partials to yield a plausible source separation. GMM fitting is sensitive to its initial guess of the parameters as the algorithm can converge to a local maximum of the likelihood function [26, p. 187]. To find an initial guess we convolve the scatter plot with kernel \mathcal{K} , giving a continuous function. \mathcal{K} is defined¹

$$\mathcal{K}(\mathbf{x}, \boldsymbol{\beta}) = \exp\left(-\frac{1}{2}\mathbf{x}^T \boldsymbol{\beta}^{-1} \mathbf{x}\right)$$

Here $\mathbf{x} \in \mathbb{R}^2$ and $\boldsymbol{\beta} \in \mathbb{R}^{2 \times 2}$ controls the extent of the kernel, i.e., how much it smooths in each dimension.

We use the two local maxima of this function as the initial means for the two sought classifying Gaussian distributions. The convolution function evaluated at $\hat{\mathbf{a}}$ is

$$f(\hat{\mathbf{a}}) = \sum_{\mathbf{a}_i \in \{\mathbf{a}\}} \mathcal{K}(\hat{\mathbf{a}} - \mathbf{a}_i, \boldsymbol{\beta}_{\mathbf{a}})$$

To make the variance proportional to the extent of each dimension, $\boldsymbol{\beta}_{\mathbf{a}}$ is defined as

$$\boldsymbol{\beta}_{\mathbf{a}} = \begin{pmatrix} \frac{\Delta_{\mathbf{a}_1}}{\Delta_{\mathbf{a}_0 + \Delta_{\mathbf{a}_1}}} \theta_{\boldsymbol{\beta}_{\mathbf{a}}} & 0 \\ 0 & \frac{\Delta_{\mathbf{a}_0}}{\Delta_{\mathbf{a}_0 + \Delta_{\mathbf{a}_1}}} \theta_{\boldsymbol{\beta}_{\mathbf{a}}} \end{pmatrix}$$

where

$$\Delta_{\mathbf{a}_0} = \max(\mathbf{a}_0) - \min(\mathbf{a}_0)$$

¹Note its similarity to the normal distribution, defined in Appendix C.

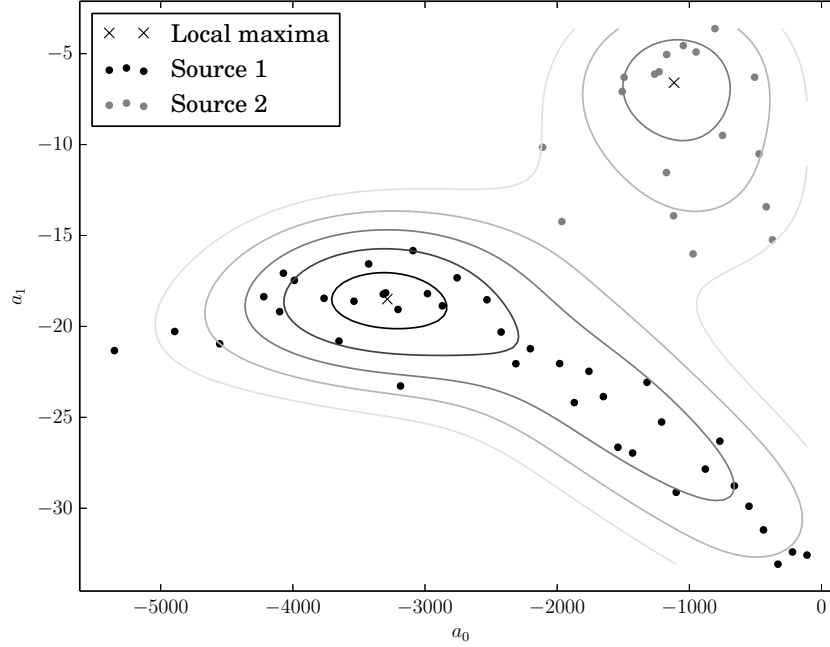


Fig. 7.8: Estimated memberships. The contours of the function resulting from convolving the data-points with kernels are represented by the lines. It is from the local maxima of these functions that the EM algorithm begins its search for the mixture of Gaussians (not shown) that give the classifications. A marker's classification is indicated by its colour.

$$\Delta_{a_1} = \max(a_1) - \min(a_1)$$

and θ_{β_a} is a parameter to control the smoothness of the resulting function, here $\theta_{\beta_a} = 1.2$. A contour plot of the resulting function $f(\hat{\mathbf{a}})$ is shown in Figure 7.8.

To initialize GMM the initial means $\boldsymbol{\mu}^0$ are chosen to be the points corresponding to the local maxima of the smoothed scatter plot². To determine initial weights \mathbf{w}^0 we first determine the value of the function at the two local maxima, $f(\mathbf{a}_0^*)$ and $f(\mathbf{a}_1^*)$. To weight relative to these two values, we compute

$$w_i^0 = \frac{\Theta_w \{f(\mathbf{a}_i^*)\}}{\sum_{p=0}^{R-1} \Theta_w \{f(\mathbf{a}_p^*)\}}$$

²Recall that the superscript here refers to the iteration number of the algorithm.

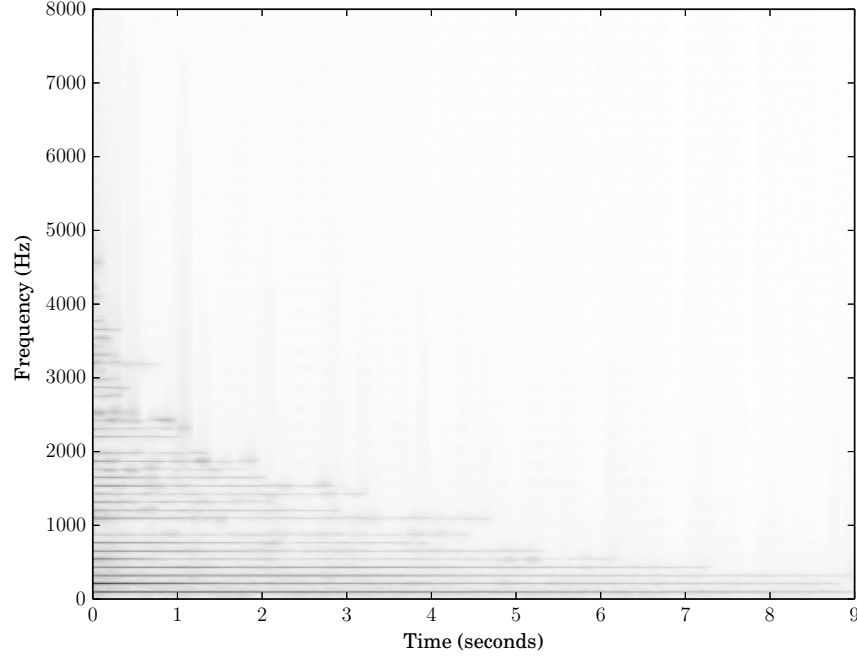


Fig. 7.9: Spectrogram of source separated acoustic guitar. The fundamental is A_3 or 220 Hz.

where Θ_w is some kind of weighting operator to have parametric control over the influence of each function value and R is the number of maxima. Here

$$\Theta_w \{f(\mathbf{a}_i^*)\} = \begin{cases} f(\mathbf{a}_i^*)\theta_w & i = 0 \\ f(\mathbf{a}_i^*) & \text{otherwise} \end{cases}$$

and $R = 2$, i.e., only the first maximum is weighted. For this experiment the parameter set as $\theta_w = 1.1$ gave the best results. The covariance matrix Σ^0 is computed as

$$\Sigma^0 = \mathbf{S}(\{\mathbf{a}\}) + \epsilon \mathbf{I}$$

where \mathbf{S} computes the sample covariance and $\epsilon \mathbf{I}$ is a matrix whose only non-zero entries are on the main diagonal and are equal to a small constant to avoid a singular initial covariance matrix. 100 iterations of the EM algorithm (see Appendix B) are performed to compute classifications. Each point is assigned to its most likely cluster using the final estimated

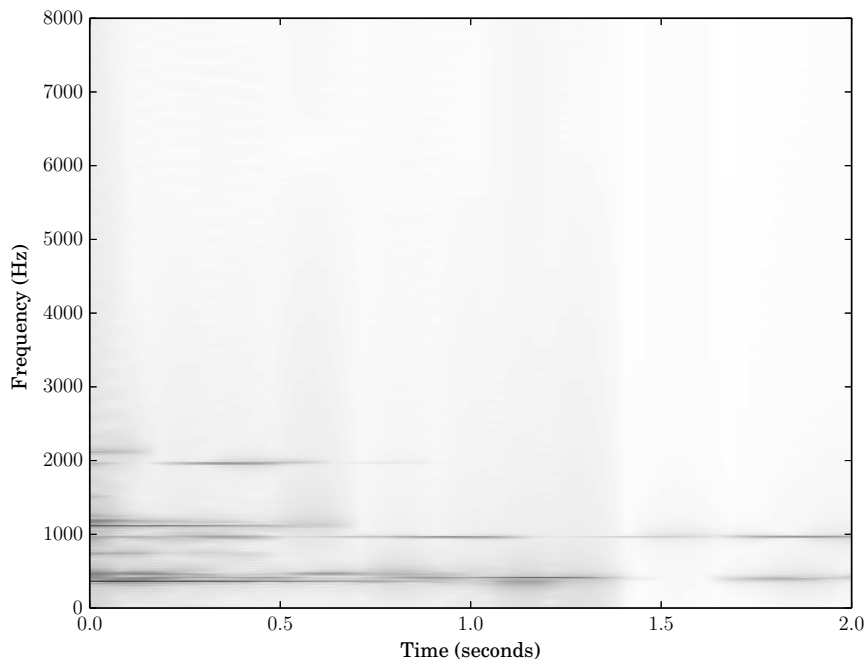


Fig. 7.10: Spectrogram of source separated xylophone. The fundamental is F_4^\sharp or approximately 370 Hz.

Gaussian distributions. The final classifications for this classification task can be seen in Figure 7.8.

7.5 Synthesis

After the classifications have been made, synthesizing the separated sources simply involves only synthesizing the partials classified as belonging to the same source. For the synthesis, we use the technique described in Section 5.4.2. Spectrograms of the source separated signals are shown in Figure 7.9 and Figure 7.10.

7.6 Conclusion

After an informal listening, the source separation is perceptually convincing.³ At least one partial from the guitar can be heard in the xylophone recording, however — it is difficult

³Soundfiles can be downloaded from
<https://drive.google.com/file/d/0B8B4c04j8tBwZDFraEZ1dFZHRFU/view?usp=sharing>

to separate partials that do not have sufficient spatial separation in Figure 7.8. Another drawback of the current technique is that it requires some tuning of the parameters θ_w and θ_{β_a} . From the spectrograms of the resynthesized sources, we see that some of the partials from both sounds were lost in the analysis. Although a shortcoming of the analysis rather than the classification, if partials are not sufficiently separated in time or frequency, they cannot be separated as their analysis will yield simply one partial when there are in fact many. In any case, it is important to see that source separation can be carried out by only considering the amplitude-modulation (in this case, the decay rate) in relation to the partial frequency. Apart from the combination of instruments presented here many plausible situations can be imagined where this technique could be carried out: e.g., with a mixture of sustained instruments such as violin, voice or horns and pitched percussive instruments such as the piano or guitar.

Chapter 8

Conclusion

Here we summarize the results of our work, highlight our contributions to the audio source separation problem, and suggest some possible extensions to the techniques presented in this thesis.

8.1 Results

8.1.1 Quality of polynomial models for analysis and synthesis

In Chapter 5 we explored various orders of polynomial amplitude and phase models for analysis and synthesis and assessed their fidelity in reproducing synthetic signals of with polynomial and exponential phase functions. A quadratic estimation of the phase and log-amplitude in the analysis step improved the synthesis quality over the constant amplitude and linear phase model in all cases considered. The use of a quintic phase polynomial in the synthesis step better approximated a signal with exponential phase. However, the use of the quintic phase and amplitude polynomial did not necessarily improve the synthesis of a signal with cubic phase and quartic amplitude. Nevertheless, the additional information provided by the DDM in the analysis stage was shown useful in postulating interpolating polynomials at the synthesis stage.

8.1.2 The use of amplitude- and frequency-modulation in audio source separation

Chapter 6 explored the contribution of amplitude- and frequency-modulation to the source separation problem on simulated data. Chapter 7 investigated, using real recordings, how to overcome the case where no frequency-modulation was present in either of the sources and so only amplitude-modulation could be used to perform separation.

Using amplitude- and frequency-modulation

In Chapter 6 we considered both the amplitude- and frequency-modulation when classifying into sources. We adaptively selected the best variable on which to classify using principal components analysis and found that amplitude-modulation aided in the case of little frequency-modulation. A post-processing step was required to extract the original sources. This step extracted sources based on the smoothness between analysis frames of the amplitude and frequency measurements. Only the post-processing encouraging smoothness in frequency gave encouraging results — smoothing in amplitude failed due to the rapidly varying amplitude in the transient regions of the sounds.

Real recordings with no frequency-modulation

In Chapter 7 we applied our source separation technique to a recording of a percussive and a plucked string instrument. These instruments exhibit virtually no frequency-modulation. We showed, as in Chapter 6, that it is possible to carry out source separation when the difference in amplitude-modulation of the two sources is sufficient.

8.2 Contributions

8.2.1 Design of continuous windows with lower side-lobes

The DDM requires the use of windows that are once differentiable. The canonical window with this property is the Hann window, however the Blackman-Harris window has better properties in terms of its out-of-band signal rejection, but is not differentiable. In Chapter 3, starting from the description of a Blackman-Harris window, we used optimization techniques to search for a window close to the Blackman-Harris, but whose end-points are 0, thus allowing differentiation everywhere after windowing with a rectangular window (necessary to make the window have finite support). Although no formal quantification of

the improvement was performed, we found this to improve the estimation accuracy of the DDM for signals sufficiently separated in frequency (separated by at least the bandwidth of the main-lobe of the frequency domain representation of the window) when compared with the DDM using the Hann window.

8.2.2 Partial tracking using linear programming

The original peak-matching algorithm of McAulay and Quatieri [36] searches for partials through a lattice of spectral analysis data by considering optimal peak-matches between two adjacent frames. In Chapter 4 we show that this can be generalized to search for L paths between an arbitrary number of frames in a lattice, and that it is a greedy algorithm — one that gives the shortest possible path without considering the lengths of the other paths in the solution. Also, it is shown that the complexity of this algorithm is exponential in the number of frames and therefore impractical to apply to even moderately sized problems. A linear program to search for the L paths through a lattice with shortest total cost is introduced. We show that this algorithm is computationally tractable and that it outperforms the McAulay-Quatieri method on signals with low SNRs.

8.3 Future extensions

8.3.1 Continuous analysis windows

It remains to be seen systematically in what situations the continuous Blackman-Harris window improves the estimation accuracy of the DDM. For this, DDM-based analysis using both the Hann and continuous Blackman-Harris windows should be performed on mixtures of synthetic signals and the accuracy of analysis compared. Further improvements to the window design could be investigated, such as finding a continuous Blackman-Harris window with good trade-off between main-lobe width and side-lobe height.

8.3.2 Partial tracking in an optimization framework

The reinterpretation of partial tracking as an optimization problem allows the application of techniques from mathematical programming, a field with a rich body of knowledge [3]. As discussed in Chapter 4, regularization could be used to encourage less “jagged” partial

trajectories. The merits of partial tracking via linear programming should be assessed on real music and speech signals.

8.3.3 Signal modeling with nonlinear amplitude and phase polynomials

A non-stationary description of signals could make it possible to describe signals using fewer analysis frames as the signal can be more accurately modeled between analysis points. As mentioned in [2], it should be investigated whether or not such a model could contribute to the field of audio coding. Specifically, compression ratios should be compared between analyses using a lower-order model, but with more analysis frames, and those using higher-order models but fewer analysis frames. Possible improvements to applications such as the time-stretching of transients and the extrapolation of missing signals (inpainting) should also be considered.

Appendices

Appendix A

Principal components analysis (PCA)

A.1 Motivation

If data-points consist of more than two dimensions (say p), it becomes burdensome to try and find the single best or two best dimensions on which to examine for grouping. If we consider variables on each of the dimensions that take on the data-point's corresponding values, we are interested in the variables that capture most of the data-points's variance. It turns out we can determine a linear transformation of our original dataset giving p variables and their p variances such that the resulting variable with the highest variance will have the maximum variance achievable, under some constraints that will be explained shortly.

A.2 Computation of principal components

The following development is based on [24]. Say we have a set $\{\mathbf{x}\}$ of data-points and their covariance matrix \mathbf{S} . A linear function of \mathbf{x} , $f_1(\mathbf{x}) = \mathbf{a}_1^T \mathbf{x}$, has variance $\sigma_{\mathbf{a}_1^T \mathbf{x}} = \mathbf{a}_1^T \mathbf{S} \mathbf{a}_1$. Therefore, we desire a vector \mathbf{a} that maximizes $\sigma_{\mathbf{a}_1^T \mathbf{x}}$. We can find this via the program

$$\max \mathbf{a}_1^T \mathbf{S} \mathbf{a}_1$$

subject to

$$\mathbf{a}_1^T \mathbf{S} \mathbf{a}_1 = 1$$

(to obtain a bounded solution).

The resulting function of \mathbf{x} , $f_1(\mathbf{x}) = \mathbf{a}_1^{*T} \mathbf{x}$ is called the first *principal component*. The second principal component $f_2(\mathbf{x}) = \mathbf{a}_2^{*T} \mathbf{x}$ is found similarly to the first, except with the additional constraint that it be uncorrelated (orthogonal) to the first component, i.e., $\mathbf{a}_2^{*T} \mathbf{a}_1 = 0$, and the third is found by requiring orthogonality with the first two principal components, etc.

The principal components now allow us to examine for grouping more easily as the total variance of the dataset has been captured in the first few principal component variables. These transformed data-points can now be classified using a classification algorithm — see Appendix B for the classification algorithm used for the experiments in this thesis.

Appendix B

Gaussian mixture models (GMM)

Consider the data-points \mathbf{X} as realizations of the vector Gaussian distributed random variable X . With a large enough sample and a small enough covariance, we will observe realizations of X as a cluster with some mean (centre point) $\boldsymbol{\mu}$ and a shape described by the covariance matrix $\boldsymbol{\Sigma}$. If we observe multiple clusters this might imply that there are P different distributions each with mean $\boldsymbol{\mu}_p$ and covariance matrix $\boldsymbol{\Sigma}_p$ and on each iteration one is chosen with probability w_p . With N observations \mathbf{x}_n we can estimate, via maximum likelihood, the P sets of parameters using a form of the *expectation maximization (EM)* algorithm [39], [8], which is an algorithm suitable for estimating missing data from known ones. First, define

$$p(\mathbf{x}_n|p) = \frac{\mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_p^k, \boldsymbol{\Sigma}_p^k) w_p^k}{\sum_{l=1}^P \mathcal{N}(\mathbf{x}_n; \boldsymbol{\mu}_l^k, \boldsymbol{\Sigma}_l^k) w_l^k}$$

the probability that \mathbf{x}_n given distribution p (see Appendix C for the definition of $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}^k, \boldsymbol{\Sigma}^k)$). The superscript k indicates the value of this parameter on iteration k . To update w_p^k :

$$w_p^{k+1} = \frac{1}{N} \sum_{n=1}^N p(\mathbf{x}_n|p)$$

which means intuitively that the probability of a data-point having been generated by distribution p is the average probability of observing any \mathbf{x}_n given p . To update $\boldsymbol{\mu}_p^{k+1}$:

$$\boldsymbol{\mu}_p^{k+1} = \frac{\sum_{n=1}^N p(\mathbf{x}_n|p) \mathbf{x}_n}{\sum_{n=1}^N p(\mathbf{x}_n|p)}$$

which is a weighted mean of all the data-points. Those less likely for a given p will weight the mean less and vice versa. A similar computation is made for $\boldsymbol{\Sigma}_p^{k+1}$:

$$\boldsymbol{\Sigma}_p^{k+1} = \frac{\sum_{n=1}^N p(\mathbf{x}_n|p) (\mathbf{x}_n - \boldsymbol{\mu}_p^{k+1}) (\mathbf{x}_n - \boldsymbol{\mu}_p^{k+1})^T}{\sum_{n=1}^N p(\mathbf{x}_n|p)}$$

The algorithm is halted after some number of iterations or when convergence is reached, i.e., the parameters change little each iteration. After convergence, the classification p^* of the data-point \mathbf{x} is simply

$$p^* = \arg \min_p p(\mathbf{x}_n|p)$$

Appendix C

The normal distribution

If a random variable X has mean $\mu = E(X)$ and variance $\sigma^2 = E((X - \mu)^2)$, where E is the expectation operator, and is completely described by these two parameters, we say that draws from X , x , follow a normal distribution. Put another way, given a draw x from X , the probability of obtaining x is given by

$$p_{\mathcal{N}}(x) = \mathcal{N}(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

When a random variable X is distributed by a normal distribution with mean μ and variance σ^2 we write

$$X \sim \mathcal{N}(\mu, \sigma^2)$$

In the case that $\tilde{\mathbf{x}}$ is an N -dimensional random vector¹ characterized by the multidimensional mean $\boldsymbol{\mu} = E(\tilde{\mathbf{x}})$ and covariance matrix $\boldsymbol{\Sigma} = E((\tilde{\mathbf{x}} - \boldsymbol{\mu})(\tilde{\mathbf{x}} - \boldsymbol{\mu})^T)$ then the probability of obtaining a realization \mathbf{x} is given by

$$((2\pi)^N |\boldsymbol{\Sigma}|)^{-1/2} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)$$

$|\boldsymbol{\Sigma}|$ means the determinant of $\boldsymbol{\Sigma}$. Similar to above, when a random vector is distributed by a normal distribution, we write

$$\tilde{\mathbf{x}} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$$

¹We write $\tilde{\mathbf{x}}$ because \mathbf{X} would be a matrix.

The normal distribution, also called the Gaussian distribution, arises often in statistics and signal processing due to the *Central Limit Theorem* [10] and its relationship to least-squares estimation [26].

References

- [1] Nancy Bertin, Roland Badeau, and Emmanuel Vincent. Enforcing harmonicity and smoothness in Bayesian non-negative matrix factorization applied to polyphonic music transcription. *IEEE Transactions on Audio, Speech, and Language Processing*, 18(3):538–549, 2010.
- [2] Michaël Betser. Sinusoidal polynomial parameter estimation using the distribution derivative. *Signal Processing, IEEE Transactions on*, 57(12):4633–4645, 2009.
- [3] Stephen Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge University Press, 2004.
- [4] James Bronson. An exploration of complex matrix factorization as a tool for single-channel musical source separation. Master’s thesis, McGill University, 2014.
- [5] Andrzej Cichocki, Rafal Zdunek, Anh Huy Phan, and Shun-ichi Amari. *Nonnegative matrix and tensor factorizations: Applications to exploratory multi-way data analysis and blind source separation*. John Wiley & Sons, 2009.
- [6] Elliot Creager. Musical source separation by coherent frequency modulation cues. Master’s thesis, McGill University, 2016.
- [7] Ronald E Crochiere and Lawrence R Rabiner. *Multi-Rate signal processing*. Englewood Cliffs, NJ: Prentice Hall, 1983.
- [8] Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society. Series B (methodological)*, pages 1–38, 1977.

-
- [9] Philippe Depalle, Guillermo Garcia, and Xavier Rodet. Tracking of partials for additive sound synthesis using hidden Markov models. In *Acoustics, Speech, and Signal Processing, 1993. ICASSP-93., 1993 IEEE International Conference on*, volume 1, pages 225–228. IEEE, 1993.
 - [10] William Feller. *An introduction to probability theory and its applications*, volume 2. John Wiley & Sons, 2008.
 - [11] Cédric Févotte, Nancy Bertin, and Jean-Louis Durrieu. Nonnegative matrix factorization with the Itakura-Saito divergence: With application to music analysis. *Neural Computation*, 21(3):793–830, 2009.
 - [12] Neville H Fletcher and Thomas Rossing. *The physics of musical instruments*. Springer Science & Business Media, 2012.
 - [13] G David Forney. The Viterbi algorithm. *Proceedings of the IEEE*, 61(3):268–278, 1973.
 - [14] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. *The elements of statistical learning*, volume 1. Springer series in statistics. Berlin: Springer, 2001.
 - [15] Laurent Girin, Sylvain Marchand, Joseph Di Martino, Axel Röbel, and Geoffroy Peeters. Comparing the order of a polynomial phase model for the synthesis of quasi-harmonic audio signals. In *Applications of Signal Processing to Audio and Acoustics, 2003 IEEE Workshop on.*, pages 193–196. IEEE, 2003.
 - [16] Gene H Golub and Charles F Van Loan. *Matrix computations*. The John Hopkins University Press, 3rd edition, 1996.
 - [17] Daniel Griffin and Jae Lim. Signal estimation from modified short-time Fourier transform. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 32(2):236–243, 1984.
 - [18] Brian Hamilton. Non-stationary sinusoidal parameter estimation. Master’s thesis, McGill University, 2011.
 - [19] Fredric J Harris. On the use of windows for harmonic analysis with the discrete Fourier transform. *Proceedings of the IEEE*, 66(1):51–83, 1978.

-
- [20] David Havelock, Sonoko Kuwano, and Michael Vorländer. *Handbook of signal processing in acoustics*. Springer Science & Business Media, 2008.
 - [21] Monson H Hayes. *Statistical digital signal processing and modeling*. John Wiley & Sons, 1999.
 - [22] Nicholas J Higham. *Accuracy and stability of numerical algorithms*. Siam, 2002.
 - [23] Hao Jiang, Sidney Fels, and James J Little. A linear programming approach for multiple object tracking. In *Computer Vision and Pattern Recognition, 2007. CVPR'07. IEEE Conference on*, pages 1–8, 2007.
 - [24] Ian Jolliffe. *Principal component analysis*. Wiley Online Library, 2002.
 - [25] Narendra Karmarkar. A new polynomial-time algorithm for linear programming. In *Proceedings of the Sixteenth Annual ACM Symposium on Theory of Computing*, pages 302–311. ACM, 1984.
 - [26] Steven M Kay. *Fundamentals of statistical signal processing, volume I: Estimation theory*. Prentice Hall, 1993.
 - [27] Corey Kereliuk and Philippe Depalle. Improved hidden Markov model partial tracking through time-frequency analysis. In *Proceedings of the 11th International Conference on Digital Audio Effects (DAFx-08), Espoo, Finland, 2008*.
 - [28] Corey Kereliuk and Philippe Depalle. Sparse atomic modeling of audio: A review. In *Proceedings of the 14th International Conference on Digital Audio Effects (DAFx-11), Paris, France, 2011*.
 - [29] Jonathan Le Roux and Emmanuel Vincent. Consistent Wiener filtering for audio source separation. *IEEE Signal Processing Letters*, 20(3):217–220, 2013.
 - [30] Daniel D Lee and H Sebastian Seung. Algorithms for non-negative matrix factorization. In *Advances in neural information processing systems*, pages 556–562, 2001.
 - [31] Yipeng Li, John Woodruff, and DeLiang Wang. Monaural musical sound separation based on pitch and common amplitude modulation. *IEEE Transactions on Audio, Speech, and Language Processing*, 17(7):1361–1371, 2009.

-
- [32] Robert Maher and James Beauchamp. An investigation of vocal vibrato for synthesis. *Applied Acoustics*, 30(2-3):219–245, 1990.
 - [33] John Makhoul. Linear prediction: A tutorial review. *Proceedings of the IEEE*, 63(4):561–580, 1975.
 - [34] Cécile MH Marin and Stephen McAdams. Segregation of concurrent sounds. II: Effects of spectral envelope tracing, frequency modulation coherence, and frequency modulation width. *The Journal of the Acoustical Society of America*, 89(1):341–351, 1991.
 - [35] Stephen McAdams. Segregation of concurrent sounds. I: Effects of frequency modulation coherence. *The Journal of the Acoustical Society of America*, 86(6):2148–2159, 1989.
 - [36] Robert J McAulay and Thomas F Quatieri. Speech analysis/synthesis based on a sinusoidal representation. *Acoustics, Speech and Signal Processing, IEEE Transactions on*, 34(4):744–754, 1986.
 - [37] Rodger J McNab, Lloyd A Smith, Ian H Witten, Clare L Henderson, and Sally Jo Cunningham. Towards the digital music library: Tune retrieval from acoustic input. In *Proceedings of the first ACM International Conference on Digital Libraries*, pages 11–18. ACM, 1996.
 - [38] Maureen Melody and Gregory H Wakefield. The time-frequency characteristics of violin vibrato: Modal distribution analysis and synthesis. *The Journal of the Acoustical Society of America*, 107(1):598–611, 2000.
 - [39] Todd K Moon. The expectation-maximization algorithm. *IEEE Signal Processing Magazine*, 13(6):47–60, 1996.
 - [40] F Richard Moore. *Elements of computer music*. Prentice-Hall, Inc., 1990.
 - [41] James A Moorer. On the transcription of musical sound by computer. *Computer Music Journal*, pages 32–38, 1977.
 - [42] Gautham J Mysore and Paris Smaragdis. A non-negative approach to language informed speech separation. In *International Conference on Latent Variable Analysis and Signal Separation*, pages 356–363. Springer, 2012.

-
- [43] F Opolko and J Wapnick. McGill University master samples [cd-rom]. *Montreal, Quebec, Canada: jwapnick@music.mcgill.ca*, 1987.
 - [44] R Gary Parker and Ronald L Rardin. *Discrete optimization*. Academic Press, 1988.
 - [45] Michael R Portnoff. Implementation of the digital phase vocoder using the fast Fourier transform. *Acoustics, Speech and Signal Processing, IEEE Transactions on*, 24(3):243–248, 1976.
 - [46] Lawrence R Rabiner. A tutorial on hidden markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2):257–286, 1989.
 - [47] Lawrence R Rabiner, Bernard Gold, and CA McGonegal. An approach to the approximation problem for nonrecursive digital filters. *Audio and Electroacoustics, IEEE Transactions on*, 18(2):83–106, 1970.
 - [48] Xavier Rodet, Yves Potard, and Jean-Baptiste Barrière. The CHANT project: From the synthesis of the singing voice to synthesis in general. *Computer Music Journal*, 8(3):15–31, 1984.
 - [49] Axel Roebel. Transient detection and preservation in the phase vocoder. In *International Computer Music Conference (ICMC)*, pages 247–250, 2003.
 - [50] Robert Roy, Arogyaswami Paulraj, and Thomas Kailath. ESPRIT—a subspace rotation approach to estimation of parameters of cisoids in noise. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 34(5):1340–1342, 1986.
 - [51] Matti P Ryynänen and Anssi P Klapuri. Automatic transcription of melody, bass line, and chords in polyphonic music. *Computer Music Journal*, 32(3):72–86, 2008.
 - [52] Xavier Serra. *A system for sound analysis/transformation/synthesis based on a deterministic plus stochastic decomposition*. PhD thesis, Stanford University, 1989.
 - [53] Madhusudana Shashanka, Bhiksha Raj, and Paris Smaragdis. Probabilistic latent variable models as nonnegative factorizations. *Computational Intelligence and Neuroscience*, 2008.

-
- [54] Paris Smaragdis and Judith C Brown. Non-negative matrix factorization for polyphonic music transcription. In *Applications of Signal Processing to Audio and Acoustics, 2003 IEEE Workshop on.*, pages 177–180. IEEE, 2003.
- [55] Paris Smaragdis, Bhiksha Raj, and Madhusudana Shashanka. A probabilistic latent variable model for acoustic modeling. *Workshop on Advances in Models for Acoustic Processing at NIPS*, 148:8–1, 2006.
- [56] Julius O Smith. *Physical Audio Signal Processing*. <http://ccrma.stanford.edu/~jos/pasp/>, accessed December 15, 2016.
- [57] Julius Orion Smith and Xavier Serra. *PARSHL: An analysis/synthesis program for non-harmonic sounds based on a sinusoidal representation*. CCRMA, Department of Music, Stanford University, 1987.
- [58] Ljubisa Stankovic, Milos Dakovic, and Thayannathan Thayaparan. *Time-frequency signal analysis with applications*. Artech house, 2014.
- [59] Charles Van Loan. *Computational frameworks for the fast Fourier transform*, volume 10. Siam, 1992.
- [60] Emmanuel Vincent, Nancy Bertin, and Roland Badeau. Harmonic and inharmonic nonnegative matrix factorization for polyphonic pitch transcription. In *2008 IEEE International Conference on Acoustics, Speech and Signal Processing*, pages 109–112. IEEE, 2008.
- [61] Tuomas Virtanen. Algorithm for the separation of harmonic sounds with time-frequency smoothness constraint. In *Proceedings of the International Conference on Digital Audio Effects (DAFx)*, pages 35–40, 2003.
- [62] Tuomas Virtanen. Monaural sound source separation by nonnegative matrix factorization with temporal continuity and sparseness criteria. *IEEE Transactions on Audio, Speech, and Language Processing*, 15(3):1066–1074, 2007.
- [63] Jack K Wolf, Audrey M Viterbi, and Glenn S Dixon. Finding the best set of K paths through a trellis with application to multitarget tracking. *IEEE Transactions on Aerospace and Electronic Systems*, 25(2):287–296, 1989.

-
- [64] Vicente Zarzoso, Eric Moreau, Rémi Gribonval, and Emmanuel Vincent. *Latent Variable Analysis and Signal Separation*. Springer, 2010.