

Statistical Parametric Speech Processing

Solving problems with the model-based approach

September 14 2016

Mads Græsbøll Christensen

Audio Analysis Lab, AD:MT
Aalborg University
Denmark



AALBORG UNIVERSITY
DENMARK



Outline

Introduction

- Motivation

- Harmonic Model

Estimating Parameters

- Parameter Estimation Bounds

- Maximum Likelihood Method

- Subspace Method

- Filtering Method

Some Examples

- Multi-Channel Modeling

- Noise Reduction

- Non-Stationary Speech

Discussion and Applications

References

Section 1

Introduction

Motivation



- ▶ Parametric speech processing is processing based on parametric models.
- ▶ Signal models described in terms of physically meaningful parameters.
- ▶ Parametric speech models have been around for many years (e.g., linear prediction in the 70s, sinusoidal model in the 80s).
- ▶ Skeptics argue that the models are (always) wrong and that it is not possible to estimate the model parameters well enough under adverse conditions.
- ▶ Parametric models can, however, be used for many things and in different ways.
- ▶ As an example, we will here take our starting point in the harmonic model.

Motivation



All models are wrong; some models are useful. (G. Box)



Motivation

Methodology:

- ▶ Methods rooted in estimation theory.
- ▶ Based on parametric models of the signal of interest.
- ▶ Analysis of estimation and modeling problems as mathematical problems.

Why parametric methods?

- ▶ They lead to robust, tractable methods whose properties can be analyzed and understood.
- ▶ A full parametrization of the signal of interest is obtained.
- ▶ Back to basics... how can we hope to solve complicated problems if we cannot solve the simple ones?

Motivation



Some interesting questions:

- ▶ Under which conditions can a method be expected to work?
- ▶ How does performance depend on the acoustic environment?
- ▶ Is the method optimal (and what does optimal mean)?
- ▶ How do we improve the method?

Only possible to answer if assumptions are made explicit! Often the assumptions are sufficient conditions but not necessary.

Non-parametric methods are hard to analyze and understand.



Harmonic Model

The harmonic model is given by (for $n = 0, \dots, N - 1$)

$$x(n) = s(n) + e(n) = \sum_{l=1}^L a_l e^{j\omega_0 l n} + e(n). \quad (1)$$

Definitions:

$s(n)$ is voiced speech

$e(n)$ is the noise/stochastic parts

ω_0 is the fundamental frequency

$\psi_l = \omega_0 l$ is the frequency of the l th harmonic

$a_l = A_l e^{j\phi_l}$ is the complex amplitude

$\theta = [\omega_0 \ A_1 \ \phi_1 \ \dots \ A_L \ \phi_L]^T$



Harmonic Model

The model can also be written as (with $\mathbf{x}(n)$ being a snapshot)

$$\mathbf{x}(n) = \mathbf{Z}(n)\mathbf{a} + \mathbf{e}(n) \quad (2)$$

$$= \mathbf{Z}\mathbf{D}^n\mathbf{a} + \mathbf{e}(n) \quad (3)$$

$$= \mathbf{Z}\mathbf{a}(n) + \mathbf{e}(n), \quad (4)$$

with the following definitions:

$$\mathbf{x}(n) = [x(n) \cdots x(n + M - 1)]^T$$

$$\mathbf{z}(\omega) = [1 e^{j\omega} \cdots e^{j\omega(M-1)}]^T$$

$$\mathbf{Z} = [\mathbf{z}(\omega_0) \cdots \mathbf{z}(\omega_0 L)]$$

$$\mathbf{D} = \text{diag}([e^{j\omega_0} e^{j\omega_0^2} \cdots e^{j\omega_0 L}])$$

$$\mathbf{a} = [a_1 \cdots a_L]^T$$



Harmonic Model

The covariance matrix of $\mathbf{x}(n)$ is

$$\mathbf{R} = \text{E} \{ \mathbf{x}(n) \mathbf{x}^H(n) \}. \quad (5)$$

Written in terms of the harmonic model, we get

$$\mathbf{R} = \mathbf{Z} \text{E} \{ \mathbf{a}(n) \mathbf{a}^H(n) \} \mathbf{Z}^H + \text{E} \{ \mathbf{e}(n) \mathbf{e}^H(n) \} \quad (6)$$

$$= \mathbf{Z} \mathbf{P} \mathbf{Z}^H + \mathbf{Q}, \quad (7)$$

which is called the covariance matrix model. Note that often it is assumed that $\mathbf{Q} = \sigma^2 \mathbf{I}$.

\mathbf{P} is the covariance matrix for the amplitudes, which can be shown to be (under certain conditions)

$$\mathbf{P} \approx \text{diag} \left([A_1^2 \ \cdots \ A_L^2] \right). \quad (8)$$



Harmonic Model

What's wrong with this model?

- ▶ It does not take non-stationarity into account
- ▶ Background noise is rarely white (and not always Gaussian)
- ▶ The model order is unknown and time-varying
- ▶ Even if stationary, signals are not perfectly periodic
- ▶ The model does not differentiate between background noise and unvoiced speech
- ▶ It is single-channel

Can this be dealt with? Does it matter?

Section 2

Estimating Parameters



Parameter Estimation Bounds

An estimate $\hat{\theta}_i$ of θ_i (i.e., the i th element of $\boldsymbol{\theta} \in \mathbb{R}^P$) is unbiased if

$$\mathbb{E} \left\{ \hat{\theta}_i \right\} = \theta_i \quad \forall \theta_i, \quad (9)$$

and the difference (if any) is referred to as the bias. The Cramér-Rao lower bound (CRLB) is then given by

$$\text{var}(\hat{\theta}_i) \geq [\mathbf{I}^{-1}(\boldsymbol{\theta})]_{ii}, \quad (10)$$

where the Fisher Information Matrix (FIM) $\mathbf{I}(\boldsymbol{\theta})$ is given by

$$[\mathbf{I}(\boldsymbol{\theta})]_{ii} = -\mathbb{E} \left\{ \frac{\partial^2 \ln p(\mathbf{x}; \boldsymbol{\theta})}{\partial \theta_i \partial \theta_i} \right\}, \quad (11)$$

with $\ln p(\mathbf{x}; \boldsymbol{\theta})$ being the log-likelihood function for $\mathbf{x} \in \mathbb{C}^N$.



Parameter Estimation Bounds

The CRLBs can be derived for the harmonic model (for WGN):

$$\text{var}(\hat{\omega}_0) \geq \frac{6\sigma^2}{N(N^2 - 1) \sum_{l=1}^L A_l^2 l^2} \quad (12)$$

$$\text{var}(\hat{A}_l) \geq \frac{\sigma^2}{2N} \quad (13)$$

$$\text{var}(\hat{\phi}_l) \geq \frac{\sigma^2}{2N} \left(\frac{1}{A_l^2} + \frac{3l^2(N-1)^2}{\sum_{m=1}^L A_m m^2 (N^2 - 1)} \right). \quad (14)$$

These depend on the following quantity:

$$\text{PSNR} = 10 \log_{10} \frac{\sum_{l=1}^L A_l^2 l^2}{\sigma^2} \text{ [dB]}. \quad (15)$$

For colored noise, pre-whitening should be employed.



Parameter Estimation Bounds

Such bounds are useful for a number of reasons:

- ▶ An estimator attaining the bound is optimal.
- ▶ The bounds tell us how performance can be expected to depend on various quantities.
- ▶ The bounds can be used as benchmarks in simulations.
- ▶ Provide us with “rules of thumb”.

Caveat emptor: The CRLB does not accurately predict the performance of non-linear estimators under adverse conditions.

It is possible to compute *exact* CRLBs, where no asymptotic approximations are used. These predict more complicated phenomena.



Parameter Estimation Bounds

It is possible to relate estimation errors to reconstruction errors. Let the observed signal be given by

$$\mathbf{x} = \mathbf{s}(\boldsymbol{\theta}) + \mathbf{e} \quad (16)$$

Suppose an estimate $\hat{\boldsymbol{\theta}}$ of $\boldsymbol{\theta}$ is used to reconstruct the i th sample as $\hat{s}_i = s_i(\hat{\boldsymbol{\theta}})$, which can be approximated as

$$s_i(\boldsymbol{\theta} + \boldsymbol{\epsilon}) \approx s_i(\boldsymbol{\theta}) + \left(\frac{\partial s_i(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right)^H \boldsymbol{\epsilon}. \quad (17)$$

The mean squared error (MSE) is then

$$E \left\{ (s_i(\boldsymbol{\theta}) - s_i(\boldsymbol{\theta} + \boldsymbol{\epsilon}))^2 \right\} = \left(\frac{\partial s_i(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right)^H E \{ \boldsymbol{\epsilon} \boldsymbol{\epsilon}^H \} \left(\frac{\partial s_i(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right). \quad (18)$$



Parameter Estimation Bounds

If a MLE is used (for sufficiently high N), then

$$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}^{-1}(\boldsymbol{\theta})), \quad (19)$$

where $\mathbf{I}(\boldsymbol{\theta})$ is the FIM! For Gaussian signals with $\mathbf{x} \sim \mathcal{N}(\mathbf{s}(\boldsymbol{\theta}), \mathbf{Q})$ where \mathbf{Q} is the noise covariance matrix, the FIM is given by

$$[\mathbf{I}(\boldsymbol{\theta})]_{nm} = \frac{\partial \mathbf{s}^H(\boldsymbol{\theta})}{\partial \theta_n} \mathbf{Q}^{-1} \frac{\partial \mathbf{s}(\boldsymbol{\theta})}{\partial \theta_m}. \quad (20)$$

The MSE can then be seen to be

$$E \left\{ (s_i(\boldsymbol{\theta}) - s_i(\boldsymbol{\theta} + \boldsymbol{\epsilon}))^2 \right\} = \left(\frac{\partial s_i(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right)^H \mathbf{I}^{-1}(\boldsymbol{\theta}) \left(\frac{\partial s_i(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right). \quad (21)$$



Maximum Likelihood Method

For Gaussian signals, the likelihood function is

$$p(\mathbf{x}(n); \theta) = \frac{1}{\pi^M \det(\mathbf{Q})} e^{-(\mathbf{x}(n) - \mathbf{Za}(n))^H \mathbf{Q}^{-1} (\mathbf{x}(n) - \mathbf{Za}(n))}. \quad (22)$$

If the noise is i.i.d., the likelihood of $\{\mathbf{x}(n)\}_{n=0}^{G-1}$ can be written as

$$p(\{\mathbf{x}(n)\}; \theta) = \prod_{n=0}^{G-1} p(\mathbf{x}(n); \theta). \quad (23)$$

The log-likelihood function is $\mathcal{L}(\theta) = \ln p(\{\mathbf{x}(n)\}; \theta)$ and the maximum likelihood estimator (MLE) is

$$\hat{\theta} = \arg \max \mathcal{L}(\theta). \quad (24)$$



Maximum Likelihood Method

For white Gaussian noise ($\mathbf{Q} = \sigma^2 \mathbf{I}$) with $M = N$ the log-likelihood function is

$$\mathcal{L}(\theta) = -N \ln \pi - N \ln \sigma^2 - \frac{1}{\sigma^2} \|\mathbf{x} - \mathbf{Z}\mathbf{a}\|_2^2. \quad (25)$$

The concentrated MLE is given by

$$\hat{\omega}_0 = \arg \max_{\omega_0} \mathcal{L}(\omega_0) = \arg \max_{\omega_0} \mathbf{x}^H \mathbf{Z} (\mathbf{Z}^H \mathbf{Z})^{-1} \mathbf{Z}^H \mathbf{x} \quad (26)$$

$$\approx \arg \max_{\omega_0} \sum_{l=1}^L \left| \sum_{n=0}^{N-1} x(n) e^{-j\omega_0 l n} \right|^2. \quad (27)$$

This can be computed using an FFT (i.e., using *harmonic summation*)!



Subspace Method

Recall that the model is

$$\mathbf{x}(n) = \mathbf{Z}\mathbf{a}(n) + \mathbf{e}(n), \quad (28)$$

and that the covariance matrix then is

$$\mathbf{R} = \mathbb{E} \{ \mathbf{x}(n)\mathbf{x}^H(n) \} = \mathbf{Z}\mathbf{P}\mathbf{Z}^H + \sigma^2\mathbf{I}, \quad (29)$$

where $\mathbf{Z}\mathbf{P}\mathbf{Z}^H$ has rank L and

$$\mathbf{P} = \text{diag} ([A_1^2 \ \cdots \ A_L^2]).$$



Subspace Method

Let $\mathbf{R} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^H$ be the EVD of the \mathbf{R} , and let \mathbf{G} be formed as

$$\mathbf{G} = [\mathbf{u}_{L+1} \quad \cdots \quad \mathbf{u}_M] \quad (30)$$

i.e., from the eigenvectors \mathbf{u}_k corresponding to the $M - L$ smallest eigenvalues. Then we have that $\mathbf{Z}^H \mathbf{G} = \mathbf{0}$.

By measuring the angles between subspaces, we can obtain an estimate as

$$\hat{\omega}_0 = \arg \min_{\omega_0} \|\mathbf{Z}^H \mathbf{G}\|_F^2 = \arg \min_{\omega_0} \sum_{l=1}^L \|\mathbf{z}^H(\omega_0 l) \mathbf{G}\|_2^2. \quad (31)$$

This maximizes the angles between the subspaces $\mathcal{R}(\mathbf{Z})$ and $\mathcal{R}(\mathbf{G})$.



Filtering Method

Let the output signal $y(n)$ of a filter having coefficients $h(n)$ be defined as

$$y(n) = \sum_{m=0}^{M-1} h(m)x(n-m) = \mathbf{h}^H \mathbf{x}(n), \quad (32)$$

with $M \leq N$ and where \mathbf{h} is a vector formed from $\{h(n)\}$. The output power is then $E\{|y(n)|^2\} = \mathbf{h}^H \mathbf{R} \mathbf{h}$.

The filtered output can be seen to be

$$\mathbf{h}^H \mathbf{x}(n) = \mathbf{h}^H \mathbf{Z} \mathbf{D}^n \mathbf{a} + \mathbf{h}^H \mathbf{e}. \quad (33)$$

If $\mathbf{h}^H \mathbf{Z} = \mathbf{1}^T$ with $\mathbf{1} = [1 \ \dots \ 1]^T$ the voiced speech would pass undistorted and the noise term $\mathbf{h}^H \mathbf{e}$ could be minimized!



Filtering Method

We would thus like to design a filter as

$$\min_{\mathbf{h}} \mathbf{h}^H \mathbf{R} \mathbf{h} \quad \text{s.t.} \quad \mathbf{h}^H \mathbf{Z} = \mathbf{1}^T. \quad (34)$$

This has the solution

$$\mathbf{h} = \mathbf{R}^{-1} \mathbf{Z} (\mathbf{Z}^H \mathbf{R}^{-1} \mathbf{Z})^{-1} \mathbf{1}. \quad (35)$$

We can use this filter to estimate the pitch as

$$\hat{\omega}_0 = \arg \max_{\omega_0} \mathbf{1}^H (\mathbf{Z}^H \mathbf{R}^{-1} \mathbf{Z})^{-1} \mathbf{1}. \quad (36)$$

Comments



- ▶ These methods are more robust to noise than non-parametric methods (YIN stops working below 10 dB, these work for -5 dB).
- ▶ They are better for low fundamental frequencies too and get better for higher SNR and N .
- ▶ The model order varies and has to be found on a per segment basis.
- ▶ Fast implementations that make the exact NLS as fast as harmonic summation exist.
- ▶ Colored noise can be dealt with.
- ▶ They can be extended to multiple pitches, although not always trivially.

Section 3

Some Examples

Multi-Channel Modeling

Introduction



- ▶ A myriad of different pitch estimators exist, but very few have been proposed for multiple channels except a few heuristic ones.
- ▶ We will now take a look at a method for multi-channel pitch estimation based on a parametric model.
- ▶ The signals in the various channels share the same fundamental frequency but can have different amplitudes, phases, and noise characteristics.
- ▶ This means that the model allows for different conditions in the various channels!



Multi-Channel Modeling

Signal Model

The method operates on snapshots $\mathbf{x}_k(n) \in \mathbb{C}^M$ for the k th channel.

These are modeled as sums of sinusoids in Gaussian noise \mathbf{e}_k with covariance \mathbf{Q}_k , i.e.,

$$\mathbf{x}_k(n) = \mathbf{Z}(n)\mathbf{a}_k + \mathbf{e}_k(n), \quad (37)$$

with $\mathbf{a}_k = [A_{k,1}e^{j\phi_{k,1}} \cdots A_{k,L}e^{j\phi_{k,L}}]^T$. Let θ_k be the parameter vector for the k th channel. The likelihood function is then

$$p(\mathbf{x}_k(n); \theta_k) = \frac{1}{\pi^M \det(\mathbf{Q}_k)} e^{-\mathbf{e}_k^H(n) \mathbf{Q}_k^{-1} \mathbf{e}_k(n)}. \quad (38)$$



Multi-Channel Modeling

Signal Model

If the deterministic part is stationary and $\mathbf{e}_k(n)$ is i.i.d. over n and independent over k , the combined likelihood is

$$p(\{\mathbf{x}_k(n)\}; \{\boldsymbol{\theta}_k\}) = \prod_{k=1}^K \frac{1}{\pi^{MG} \det(\mathbf{Q}_k)^G} e^{-\sum_{n=0}^{G-1} \mathbf{e}_k^H(n) \mathbf{Q}_k^{-1} \mathbf{e}_k(n)}. \quad (39)$$

For simplicity, we assume that the noise is white in each channel but has different σ_k^2 , i.e., $\mathbf{Q}_k = \sigma_k^2 \mathbf{I}$.

The log-likelihood function then reduces to

$$\ln p(\{\mathbf{x}_k(n)\}; \{\boldsymbol{\theta}_k\}) = -GM \sum_{k=1}^K \ln(\pi \sigma_k^2) - \sum_{k=1}^K \sum_{n=0}^{G-1} \frac{\|\mathbf{e}_k(n)\|^2}{\sigma_k^2}. \quad (40)$$



Multi-Channel Modeling

Estimator

The MLE of the amplitudes for channel k are

$$\hat{\mathbf{a}}_k = \left(\sum_{n=0}^{G-1} \mathbf{z}^H(n) \mathbf{z}(n) \right)^{-1} \sum_{n=0}^{G-1} \mathbf{z}^H(n) \mathbf{x}_k(n). \quad (41)$$

This can be used to form a noise variance estimate as

$$\hat{\sigma}_k^2 = \frac{1}{GM} \sum_{n=0}^{G-1} \|\hat{\mathbf{e}}_k(n)\|^2 = \frac{1}{GM} \sum_{n=0}^{G-1} \|\mathbf{x}_k(n) - \mathbf{z}(n) \hat{\mathbf{a}}_k\|^2. \quad (42)$$

This yields the following log-likelihood for channel k at time n

$$\ln p(\mathbf{x}_k(n); \omega_0) = -M \ln \pi - M \ln \hat{\sigma}_k^2.$$



Multi-Channel Modeling

Estimator

For all n and k , this yields

$$\ln p(\{\mathbf{x}_k(n)\}; \omega_0) = -GMK \ln \pi - GM \sum_{k=1}^K \ln \hat{\sigma}_k^2. \quad (43)$$

The maximum likelihood estimator (MLE) can finally be stated as

$$\hat{\omega}_0 = \arg \min_{\omega_0} \sum_{k=1}^K \ln \hat{\sigma}_k^2. \quad (44)$$

This estimator can then be approximated as

$$\hat{\omega}_0 = \arg \min_{\omega_0} \sum_{k=1}^K \ln \left(\|\mathbf{x}_k\|^2 - \frac{1}{N} \|\mathbf{Z}^H \mathbf{x}_k\|^2 \right), \quad (45)$$

where $\mathbf{x}_k = \mathbf{x}_k(0)$ for $M = N$. This can be computed using FFTs.

Multi-Channel Modeling

Experiments

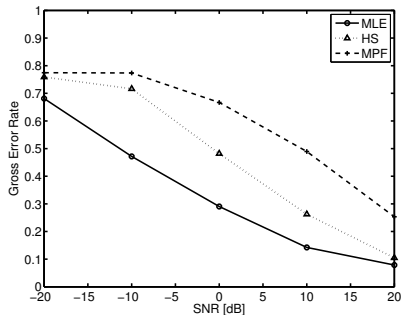
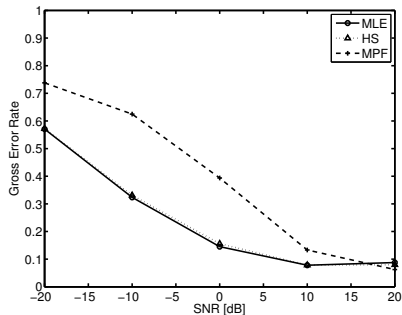


Figure: Gross error rate for (left) symmetrical noise level and (right) asymmetrical noise level (i.e., different noise levels).

Comments



- ▶ As we have seen, it was fairly straightforward to extend the MLE to multiple channels.
- ▶ It works well and under very general conditions.
- ▶ It is fast too.
- ▶ Easy to build in more specific knowledge, like array structure, nearfield, TDOAs, binaural setups.
- ▶ The multi-channel model contains the usual broadband model as a special case with $\omega_0 = 2\pi/N$.
- ▶ Can be used for pitch/DOA estimation and model-based beamforming.

Noise Reduction

Introduction



- ▶ The harmonic signal model has been used for noise reduction in various ways, like the traditional comb filters.
- ▶ We have seen how adaptive and optimal filters can be used for pitch estimation.
- ▶ The same principle can be used for finding noise reduction filters.
- ▶ Some interesting and well-known special cases can be obtained from these filters.



Noise Reduction

Filter Design

As we saw earlier, we get the following model when a filter \mathbf{h} is applied to the observed signal $\mathbf{x}(n)$:

$$\hat{\mathbf{s}}(n) = \mathbf{h}^H \mathbf{x}(n) = \mathbf{h}^H \mathbf{Z} \mathbf{D}^n \mathbf{a} + \mathbf{h}^H \mathbf{e}. \quad (46)$$

This comprises two terms:

- ▶ The filtered voiced speech $\mathbf{h}^H \mathbf{Z} \mathbf{D}^n \mathbf{a}$
- ▶ The filtered noise $\mathbf{h}^H \mathbf{e}$

If $\mathbf{h}^H \mathbf{Z} = \mathbf{1}^T$ then $\mathbf{h}^H \mathbf{Z} \mathbf{D}^n \mathbf{a} = \sum_{l=1}^L a_l e^{j\omega_0 l n}$ while $E\{|\mathbf{h}^H \mathbf{e}|^2\} = \mathbf{h}^H \mathbf{Q} \mathbf{h}$ is minimized, we have distortionless optimal noise reduction!



Noise Reduction

Filter Design

A distortionless filter should have $\mathbf{h}^H \mathbf{z} = \mathbf{1}^T$ and should minimize the residual noise, i.e.,

$$\min_{\mathbf{h}} \mathbf{h}^H \hat{\mathbf{Q}} \mathbf{h} \quad \text{s.t.} \quad \mathbf{z}^H \mathbf{h} = \mathbf{1} \quad (47)$$

The solution can be shown to be

$$\hat{\mathbf{h}} = \hat{\mathbf{Q}}^{-1} \mathbf{z} \left(\mathbf{z}^H \hat{\mathbf{Q}}^{-1} \mathbf{z} \right)^{-1} \mathbf{1}. \quad (48)$$

with $\hat{\mathbf{Q}}$ being a particular *noise* covariance matrix estimate.

These filters are adaptive, optimal comb filters! Unlike the normally used Wiener filter, these do not distort the desired signal.



Noise Reduction

Noise Covariance Matrix

We seek to find a filter such that the MSE is minimized:

$$MSE = \frac{1}{G} \sum_{n=M-1}^{N-1} \left| y(n) - \sum_{l=1}^L a_l e^{j\omega_0 l n} \right|^2 = \frac{1}{G} \sum_{n=M-1}^{N-1} |\mathbf{h}^H \mathbf{x}(n) - \mathbf{a}^H \mathbf{w}(n)|^2,$$

with $\mathbf{w}(n) = [e^{j\omega_0 1n} \dots e^{j\omega_0 Ln}]^T$. Solving for the amplitudes, we get

$$MSE = \mathbf{h}^H \left(\hat{\mathbf{R}} - \mathbf{G}^H \mathbf{W}^{-1} \mathbf{G} \right) \mathbf{h} \triangleq \mathbf{h}^H \hat{\mathbf{Q}} \mathbf{h}, \quad (49)$$

where $\mathbf{G} = \frac{1}{G} \sum_{n=M-1}^{N-1} \mathbf{w}(n) \mathbf{x}^H(n)$ and $\mathbf{W} = \frac{1}{G} \sum_{n=M-1}^{N-1} \mathbf{w}(n) \mathbf{w}^H(n)$.

Thus we can estimate \mathbf{Q} as $\hat{\mathbf{Q}} = \hat{\mathbf{R}} - \mathbf{G}^H \mathbf{W}^{-1} \mathbf{G}$



Noise Reduction

Special Cases

Special cases:

- ▶ Setting $\mathbf{W} = \mathbf{I}$ yields the usual noise covariance matrix estimate.
- ▶ Capon-like filters can be obtained from $\hat{\mathbf{Q}} = \hat{\mathbf{R}}$, i.e.,

$$\hat{\mathbf{h}} = \hat{\mathbf{R}}^{-1} \mathbf{Z} \left(\mathbf{Z}^H \hat{\mathbf{R}}^{-1} \mathbf{Z} \right)^{-1} \mathbf{1}.$$
- ▶ Setting $\hat{\mathbf{R}} = \sigma^2 \mathbf{I}$ yields $\hat{\mathbf{h}} = \mathbf{Z} \left(\mathbf{Z}^H \mathbf{Z} \right)^{-1} \mathbf{1}$.
- ▶ Noting that $\lim_{M \rightarrow \infty} M \mathbf{Z} \left(\mathbf{Z}^H \mathbf{Z} \right)^{-1} = \mathbf{Z}$, we get $\hat{\mathbf{h}} = \frac{1}{M} \mathbf{Z} \mathbf{1}$.
- ▶ Binary masking can also be obtained using these principles.

Noise Reduction

Examples

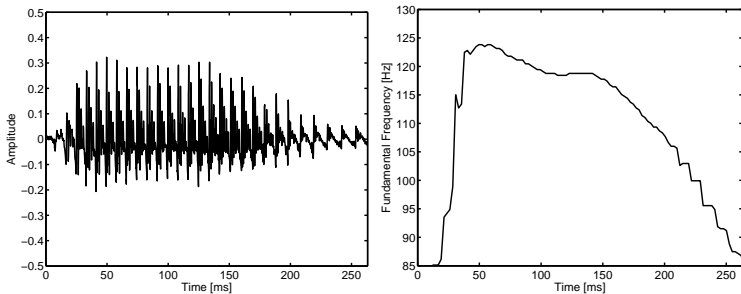


Figure: The original voiced speech signal and the estimated pitch.

Noise Reduction

Examples

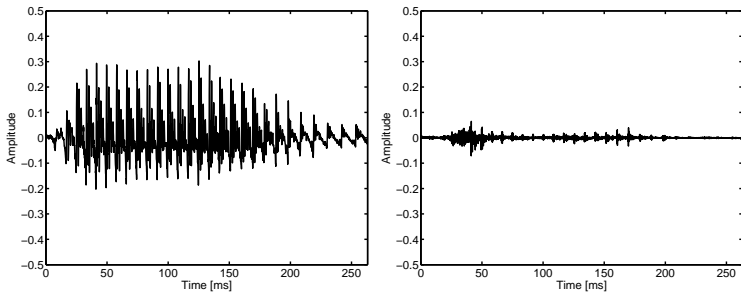


Figure: The extracted signal and the difference between the two signals, i.e., the part of the signal that was not extracted.

Noise Reduction

Examples

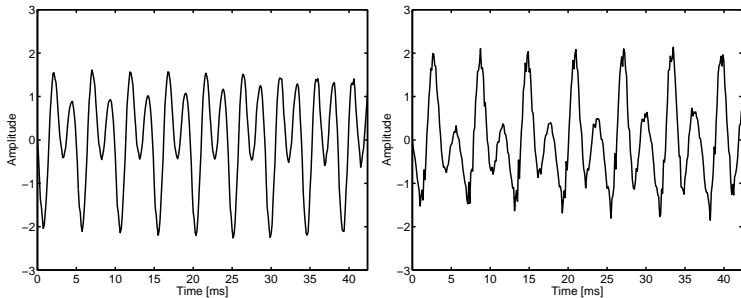


Figure: The voiced speech signal of sources 1 and 2.

Noise Reduction

Examples

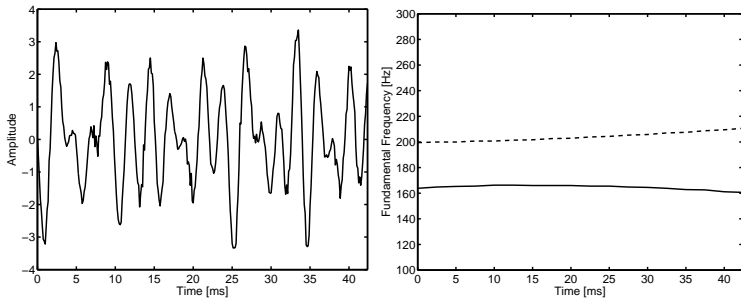


Figure: The mixture of the two signals and the estimated pitch tracks for source 1 (dashed) and 2 (solid).

Noise Reduction

Examples

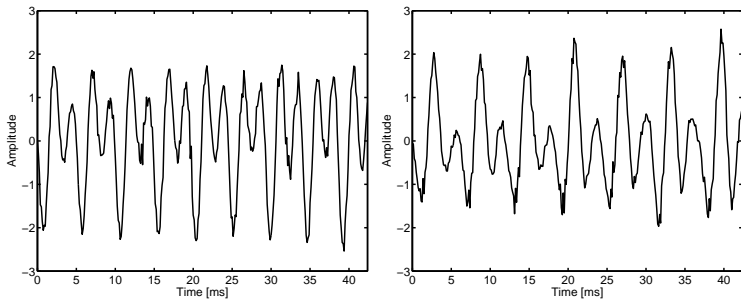


Figure: The estimate of sources 1 and 2 obtained from the mixture.

Comments



- ▶ We have seen how the harmonic model can be used for designing filters for noise reduction.
- ▶ The filters are distortionless, i.e., they let the signal of interest pass undistorted.
- ▶ Meanwhile, the noise is attenuated as much as possible.
- ▶ The resulting filters are thus optimal in terms of output SNR and minimum distortion!
- ▶ They do not require a priori knowledge of noise statistics.
- ▶ They can be generalized to multiple channels.

Non-Stationary Speech

Introduction



- ▶ Parametric methods based on the harmonic model have proven to overcome the problems of correlation-based methods.
- ▶ However, as mentioned earlier, there might be concerns about the stationarity within segments.
- ▶ To investigate whether this is a problem, we will take a closer look at the harmonic chirp model and derive an estimator for determining its parameters.



Non-Stationary Speech

Signal Model

For a segment of a speech signal with $n = n_0, \dots, n_0 + N - 1$ the new harmonic chirp model is given by

$$x(n) = \sum_{l=1}^L A_l e^{j\theta_l(n)} + e(n) \quad (50)$$

where

- ▶ L is the number of harmonics (assumed known).
- ▶ A_l the l th is the amplitude.
- ▶ $\theta_l(n)$ is the instantaneous phase of the l th harmonic.
- ▶ $e(n)$ are the stochastic parts of the observed signal.
- ▶ n_0 is the start index.



Non-Stationary Speech

Signal Model

The instantaneous phase $\theta_l(\cdot)$ is given by

$$\theta_l(t) = \int_0^t l\omega_0(\tau) d\tau + \phi_l, \quad (51)$$

where $\omega_0(t)$ is the time-varying pitch and ϕ_l is the phase of the l th harmonic. In the harmonic model (HM) we have that $\omega_l(t) = l\omega_0$.

If the pitch is slowly varying, i.e., $\omega_0(t) = \alpha_0 t + \omega_0$, we get

$$\theta_l(t) = \frac{1}{2}\alpha_0 l t^2 + \omega_0 l t + \phi_l, \quad (52)$$

where α_0 is the fundamental chirp rate.

The resulting model is called the harmonic chirp model (HCM).



Non-Stationary Speech

NLS Estimator

Define a vector with $n_0 = -(N - 1)/2$ as

$$\mathbf{x} = [x(n_0) \quad x(n_0 + 1) \quad \dots \quad x(n_0 + N - 1)]. \quad (53)$$

and a matrix as

$$\mathbf{Z} = [\mathbf{z}(\omega_0, \alpha_0) \quad \mathbf{z}(2\omega_0, 2\alpha_0) \quad \dots \quad \mathbf{z}(L\omega_0, L\alpha_0)], \quad (54)$$

with columns

$$\mathbf{z}(l\omega_0, l\alpha_0) = \left[e^{j(\frac{1}{2}\alpha_0 l n_0^2 + \omega_0 l n_0)} \quad \dots \quad e^{j(\frac{1}{2}\alpha_0 l (n_0 + N - 1)^2 + \omega_0 l (n_0 + N - 1))} \right]^T.$$

For convenience, we introduce $\mathbf{\Pi}_{\omega_0, \alpha_0} = \mathbf{Z} (\mathbf{Z}^H \mathbf{Z})^{-1} \mathbf{Z}^H$.



Non-Stationary Speech

NLS Estimator

As before, the nonlinear least squares (NLS) estimator can be used:

$$\{\hat{\alpha}_0, \hat{\omega}_0\} = \arg \min_{\alpha_0, \omega_0} \|\mathbf{x} - \mathbf{Z} (\mathbf{Z}^H \mathbf{Z})^{-1} \mathbf{Z}^H \mathbf{x}\|^2. \quad (55)$$

We solve this iteratively as follows (with i being the iteration index). First obtain an estimate $\hat{\alpha}_0^{(i)}$ from $\hat{\omega}_0^{(i-1)}$ for $i = 1, 2, \dots$ as

$$\hat{\alpha}_0^{(i)} = \arg \max_{\alpha_0} \left\{ \mathbf{x}^H \mathbf{\Pi}_{\hat{\omega}_0^{(i-1)}, \alpha_0} \mathbf{x} \right\}, \quad (56)$$

and then update the estimate of the fundamental frequency, ω_0 , as

$$\hat{\omega}_0^{(i)} = \arg \max_{\omega_0} \left\{ \mathbf{x}^H \mathbf{\Pi}_{\omega_0, \hat{\alpha}_0^{(i)}} \mathbf{x} \right\}. \quad (57)$$

This is then repeated for $i = 1, 2, \dots$ until convergence. We initialize with $\alpha_0^{(0)} = 0$.

Non-Stationary Speech

Experiments

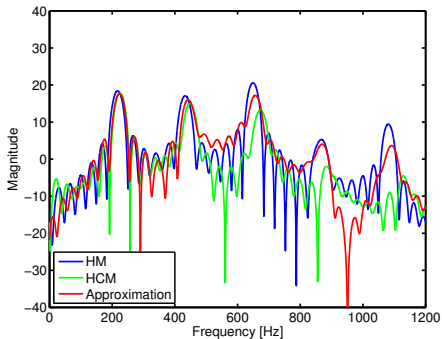


Figure: Spectrum of harmonic model, harmonic chirp model, and an approximation.

Non-Stationary Speech

Experiments

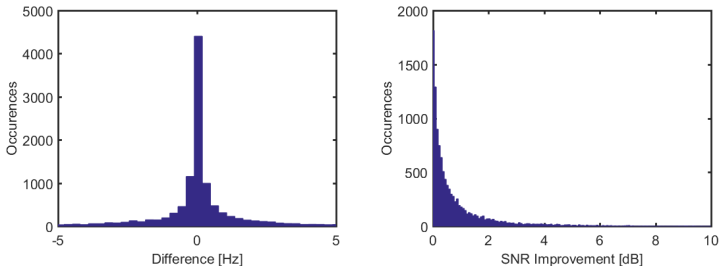


Figure: Histogram of differences in pitch estimates (left) and reconstruction SNRs (right) between HM and HCM for 30 sentences.

Comments



- ▶ As we have seen, it is quite easy to account for non-stationarity.
- ▶ Although the differences in pitch estimates are small, they may matter.
- ▶ There exists fast implementations for the exact NLS for the harmonic chirp model too!
- ▶ It is also possible to use HCM with the distortionless filters, meaning we can design filters that account for the non-stationarity of speech.

Section 4

Discussion and Applications

Summary



We have seen how

- ▶ the problem of finding the parameters of the harmonic model can be analyzed.
- ▶ the parameters of the harmonic model can be found in various ways.
- ▶ the harmonic model and its estimators can be extended to multiple channels under quite general conditions.
- ▶ the harmonic model can be used for designing optimal and distortionless filters that do not require knowledge of noise statistics.
- ▶ it is fairly straightforward to take the non-stationary nature of speech into account.

Applications



These ideas are/can be used in many applications, including:

- ▶ Hearing aids
- ▶ Voice over IP
- ▶ Telecommunication
- ▶ Reproduction systems
- ▶ Voice analysis
- ▶ Intelligence, law enforcement, defense
- ▶ Music equipment/software

Some Other Results



- ▶ Parametric models can be used for speech/audio compression.
- ▶ Model-based interpolation/extrapolation can be used for packet losses/corrupt data.
- ▶ Feedback cancellation can be improved using a model of the near-end signal.
- ▶ Beamforming can be improved with the model-based approach.
- ▶ Jointly optimal segmentation and parameter estimates can be found with dynamic programming.
- ▶ Optimal filters can be designed for the chirp model too.
- ▶ We have recently shown that fast implementations can be found!

Conclusion



- ▶ Parametric models have shown promise for several problems, but they are not (yet) widespread.
- ▶ An argument against the usage of such models is that they do not take various phenomena into account.
- ▶ However, we can only have this discussion because the assumptions are explicit.
- ▶ And it is often fairly easy to improve the model and methods, if needed.
- ▶ There are many more speech processing problems that could probably benefit from this approach!
- ▶ These include applications with multiple channels, adverse conditions or where the fine details matter.



References I

M. G. Christensen and A. Jakobsson, *Multi-Pitch Estimation*, ser. Synthesis Lectures on Speech & Audio Processing. Morgan & Claypool Publishers, 2009, vol. 5, 160 pages.

S. M. Nørholm, J. R. Jensen, and M. G. Christensen, “Instantaneous pitch estimation with optimal segmentation for non-stationary voiced speech,” *IEEE Trans. Audio, Speech, Language Process.*, 2016, accepted.

S. Karimian-Azari, J. R. Jensen, and M. G. Christensen, “Computationally efficient and noise robust DOA and pitch estimation,” *IEEE Trans. Audio, Speech, Language Process.*, vol. 24(9), pp. 1613–1625, 2016.

S. M. Nørholm, J. R. Jensen, and M. G. Christensen, “Enhancement and noise statistics estimation for non-stationary voiced speech,”



References II

IEEE Trans. Audio, Speech, Language Process., vol. 24(4), pp. 645–658, 2016.

J. R. Jensen, M. G. Christensen, J. Benesty and S. H. Jensen, “Joint spatio-temporal filtering methods for DOA and fundamental frequency estimation,” *IEEE Trans. Audio, Speech, Language Process.*, vol. 23(1), pp. 174–185, 2015.

S. I. Adalbjörnsson, A. Jakobsson, and M. G. Christensen, “Multi-pitch estimation exploiting block sparsity,” *Signal Processing*, vol. 109, pp. 236–247, Apr. 2015.

J. K. Nielsen, M. G. Christensen, A. T. Cemgil, and S. H. Jensen, “Bayesian model comparison with the g-prior,” *IEEE Trans. Signal Process.*, vol. 62(1), pp. 225–238, 2014.



References III

M. G. Christensen, “Accurate estimation of low fundamental frequencies,” *IEEE Trans. Audio, Speech, Language Process.*, vol. 21(10), pp. 2042–2056, 2013.

Z. Zhou, H. C. So, and M. G. Christensen, “Parametric modeling for damped sinusoids from multiple channels,” *IEEE Trans. Signal Process.*, vol. 61(15), pp. 3895–3907, 2013.

J. R. Jensen, G.-O. Glentis, M. G. Christensen, A. Jakobsson, and S. H. Jensen, “Fast LCMV-based methods for fundamental frequency estimation,” *IEEE Trans. Signal Process.*, vol. 61(12), pp. 3159–3172, 2013.

M. G. Christensen, “Metrics for vector quantization-based parametric speech enhancement and separation,” *J. Acoust. Soc. Am.*, vol. 133(5), pp. 3062–3071, 2013.



References IV

K. Ngo, T. van Waterschoot, M. G. Christensen, M. Moonen, and S. H. Jensen, “Improved prediction error filters for adaptive feedback cancellation in hearing aids,” *Signal Processing*, vol. 91(11), pp. 3062–3075, 2013.

J. R. Jensen, M. G. Christensen, and S. H. Jensen, “Nonlinear least squares methods for joint DOA and pitch estimation,” *IEEE Trans. Audio, Speech, Language Process.*, vol. 21(5), pp. 923–933, 2013.

J. K. Nielsen, M. G. Christensen, and S. H. Jensen, “Default Bayesian estimation of the fundamental frequency,” *IEEE Trans. Audio, Speech, Language Process.*, vol. 21(3), pp. 598–610, 2013.

P. Mowlaee, R. Saeidi, M. G. Christensen, Z.-H. Tan, T. Kinnunen, P. Fränti, and S. H. Jensen, “A joint approach for single-channel speaker identification and speech separation,” *IEEE Trans. Audio, Speech, Language Process.*, vol. 20(9), pp. 2586–2601, 2012.



References V

J. R. Jensen, J. Benesty, M. G. Christensen, and S. H. Jensen, “Enhancement of single-channel periodic signals in the time-domain,” *IEEE Trans. Audio, Speech, Language Process.*, vol. 20(7), pp. 1948–1963, 2012.

D. Giacobello, M. G. Christensen, M. N. Murthi, S. H. Jensen, and M. Moonen, “Sparse linear prediction and its applications to speech processing,” *IEEE Trans. Audio, Speech, Language Process.*, vol. 20(5), pp. 1644–1657, 2012.

J. X. Zhang, M. G. Christensen, S. H. Jensen, and M. Moonen, “Joint DOA and multi-pitch estimation based on subspace techniques,” *EURASIP J. on Advances in Signal Process.*, vol. 2012(1), pp. 1–11, 2012.

M. G. Christensen, J. L. Højvang, A. Jakobsson, and S. H. Jensen, “Joint fundamental frequency and order estimation using optimal



References VI

filtering,” *EURASIP J. on Advances in Signal Process.*, vol. 2011(1), pp. 1–13, 2011.

J. K. Nielsen, M. G. Christensen, A. T. Cemgil, S. J. Godsill, and S. H. Jensen, “Bayesian interpolation and parameter estimation in a dynamic sinusoidal model,” *IEEE Trans. Audio, Speech, Language Process.*, vol. 19(7), pp. 1986–1998, 2011.

P. Mowlae, M. G. Christensen, and S. H. Jensen, “New results on single-channel speech separation using sinusoidal modeling,” *IEEE Trans. Audio, Speech, Language Process.*, vol. 19(5), pp. 1265–1277, 2011.

M. G. Christensen and A. Jakobsson, “Optimal filter designs for separating and enhancing periodic signals,” *IEEE Trans. Signal Process.*, vol. 58(12), pp. 5969–5983, 2010.



References VII

J. X. Zhang, M. G. Christensen, S. H. Jensen, and M. Moonen, “A robust and computationally efficient subspace-based fundamental frequency estimator,” *IEEE Trans. Audio, Speech, Language Process.*, vol. 18(3), pp. 487–497, 2010.

M. G. Christensen, A. Jakobsson, and S. H. Jensen, “Sinusoidal order estimation using angles between subspaces,” *EURASIP J. on Advances in Signal Process.*, pp. 1–11, 2009, Article ID 948756.

M. G. Christensen, J. H. Jensen, A. Jakobsson and S. H. Jensen, “On optimal filter designs for fundamental frequency estimation,” *IEEE Signal Process. Lett.*, vol. 15, pp. 745–748, 2008.

M. G. Christensen, P. Stoica, A. Jakobsson and S. H. Jensen, “Multi-pitch estimation,” *Signal Processing*, vol. 88(4), pp. 972–983, Apr. 2008.



References VIII

M. G. Christensen, A. Jakobsson and S. H. Jensen, “Joint high-resolution fundamental frequency and order estimation,” *IEEE Trans. Audio, Speech, Language Process.*, vol. 15(5), pp. 1635–1644, 2007.

M. S. Kavalekalam, M. G. Christensen, and J. B. Boldt, “Binaural speech enhancement using a codebook based approach,” in *Proc. Int. Workshop on Acoustic Signal Enhancement*, 2016.

M. W. Hansen, J. R. Jensen, and M. G. Christensen, “Multi-pitch estimation of audio recordings using a codebook-based approach,” in *Proc. European Signal Processing Conf.*, 2016.

J. K. Nielsen, T. L. Jensen, J. R. Jensen, M. G. Christensen, and S. H. Jensen, “Grid size selection for nonlinear least-squares optimization in spectral estimation and array processing,” in *Proc. European Signal Processing Conf.*, 2016.



References IX

J. K. Nielsen, T. L. Jensen, J. R. Jensen, M. G. Christensen, and S. H. Jensen, “Fast and statistically efficient fundamental frequency estimation,” in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2016, pp. 86–90.

J. R. Jensen, J. K. Nielsen, R. Heusdens, and M. G. Christensen, “DOA estimation of audio sources in reverberant environments,” in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2016, pp. 176–80.

J. K. Nielsen, T. L. Jensen, J. R. Jensen, M. G. Christensen, and S. H. Jensen, “A fast algorithm for maximum likelihood-based fundamental frequency estimation,” in *Proc. European Signal Processing Conf.*, 2015.



References X

M. W. Hansen, J. R. Jensen, and M. G. Christensen, “Pitch estimation of stereophonic mixtures of delay and amplitude panned signals,” in *Proc. European Signal Processing Conf.*, 2015.

J. R. Jensen, J. K. Nielsen, M. G. Christensen and S. H. Jensen, “On frequency domain models for TDOA estimation,” in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2015, pp. 11–15.

M. G. Christensen and J. R. Jensen, “Pitch estimation for non-stationary speech,” in *Rec. Asilomar Conf. Signals, Systems, and Computers*, 2014, pp. 1400–1404.

J. R. Jensen and M. G. Christensen, “Near-field localization of audio: A maximum likelihood approach,” in *Proc. European Signal Processing Conf.*, 2014, pp. 895–899.



References XI

- S. M. Nørholm, J. R. Jensen, and M. G. Christensen, “On the influence of inharmonicities in model-based speech enhancement,” in *Proc. European Signal Processing Conf.*, 2013, pp. 1–5.
- M. G. Christensen, J. R. Jensen, J. Benesty, and A. Jakobsson, “Spatio-temporal filtering methods for enhancement and separation of speech signals,” in *Proc. IEEE China Summit & Int. Conf. on Signal and Information Process.*, 2013, pp. 303–307.
- S. Karimian-Azari, J. R. Jensen, and M. G. Christensen, “Fast joint doa and pitch estimation using a broadband MVDR beamformer,” in *Proc. European Signal Processing Conf.*, 2013, pp. 1–5.
- J. K. Nielsen, M. G. Christensen, and S. H. Jensen, “Bayesian model comparison and the BIC for regression models,” in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2013, pp. 6362–6366.



References XII

J. R. Jensen, M. G. Christensen, and S. H. Jensen, “Statistically efficient methods for pitch and DOA estimation,” in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2013, pp. 3900–3904.

J. R. Jensen, M. G. Christensen, J. Benesty, and S. H. Jensen, “Joint filtering scheme for nonstationary noise reduction,” in *Proc. European Signal Processing Conf.*, 2012, pp. 2323–2327.

M. G. Christensen, “A method for low-delay pitch tracking and smoothing,” in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2012, pp. 345–348.

J. K. Nielsen, M. G. Christensen, and S. H. Jensen, “An approximate bayesian fundamental frequency estimator,” in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2012, pp. 4617–4620.



References XIII

M. G. Christensen, “Multi-channel maximum likelihood pitch estimation,” in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, 2012, pp. 409–412.

J. R. Jensen, M. G. Christensen, and S. H. Jensen, “Fundamental frequency estimation using polynomial rooting of a subspace-based method,” in *Proc. European Signal Processing Conf.*, 2010.

J. R. Jensen, J. K. Nielsen, M. G. Christensen, S. H., Jensen and T. Larsen, “On fast implementation of harmonic music for known and unknown model orders,” in *Proc. European Signal Processing Conf.*, 2008.