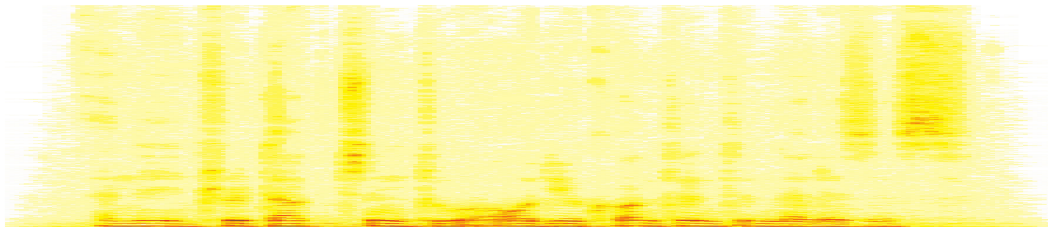


Introduction to Audio Content Analysis

Module 7.1: Audio-to-Audio & Audio-to-Score Alignment

alexander lerch



introduction

overview

corresponding textbook section

Chapter 7: Alignment (pp. 146–150)

- **lecture content**

- Audio-to-Audio alignment
 - use cases
 - features
 - distance measures
 - typical accuracy
- Audio-to-Score alignment

- **learning objectives**

- elaborate on possible use cases for audio-to-audio alignment
- give examples for features and distance measures for alignment
- discuss differences between audio-to-audio and audio-to-score alignment



introduction

overview

corresponding textbook section

Chapter 7: Alignment (pp. 146–150)

- **lecture content**

- Audio-to-Audio alignment
 - use cases
 - features
 - distance measures
 - typical accuracy
- Audio-to-Score alignment

- **learning objectives**

- elaborate on possible use cases for audio-to-audio alignment
- give examples for features and distance measures for alignment
- discuss differences between audio-to-audio and audio-to-score alignment



audio-to-audio alignment

introduction

- **objective**

- align two sequences of audio

- **use cases**

- *quick browsing* for certain parts in recordings
- *timing adjustment* (backing vocals, loops, ...)
- *automated dubbing*
- *musicological analysis* (timing of several performances)

- **processing steps**

- extract suitable features
- compute distance matrix
- compute alignment path

audio-to-audio alignment

introduction

- **objective**

- align two sequences of audio

- **use cases**

- *quick browsing* for certain parts in recordings
- *timing adjustment* (backing vocals, loops, ...)
- *automated dubbing*
- *musicological analysis* (timing of several performances)

- **processing steps**

- extract suitable features
- compute distance matrix
- compute alignment path

audio-to-audio alignment

introduction

- **objective**

- align two sequences of audio

- **use cases**

- *quick browsing* for certain parts in recordings
- *timing adjustment* (backing vocals, loops, ...)
- *automated dubbing*
- *musicological analysis* (timing of several performances)

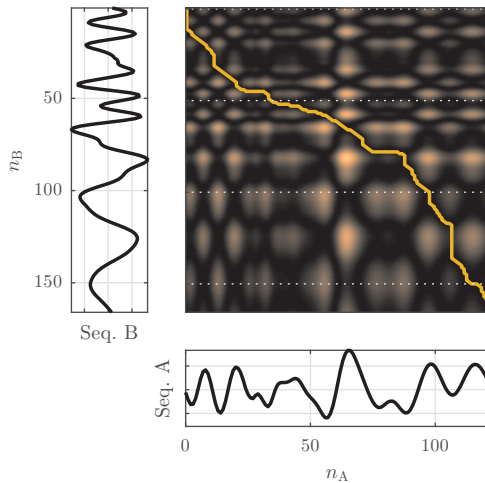
- **processing steps**

- extract suitable features
- compute distance matrix
- compute alignment path

audio-to-audio alignment

alignment path computation

→ prerequisite: Module 7.0—Dynamic Time Warping



audio-to-audio alignment

features

- **use case examples**

- **quick browsing** — find the same part across files
⇒ use *pitch based* features
- **timing adjustment** — backing vocals to lead vocals
⇒ use *intensity based* features
- **automated dubbing** — same speaker several recordings
⇒ use *intensity based* and *timbre based* features

- **feature categories**

- **intensity**: energy, onset probability, ...
- **tonal**: pitch chroma, ...
- **timbral**: MFCCs, spectral shape, ...

plot from¹

¹H. Kirchhoff and A. Lerch, "Evaluation of Features for Audio-to-Audio Alignment," *Journal of new music research*, vol. 40, no. 1, pp. 27–41, 2011. DOI: [10.1080/09298215.2010.529917](https://doi.org/10.1080/09298215.2010.529917).

audio-to-audio alignment

features

- **use case examples**

- **quick browsing** — find the same part across files
⇒ use *pitch based* features
- **timing adjustment** — backing vocals to lead vocals
⇒ use *intensity based* features
- **automated dubbing** — same speaker several recordings
⇒ use *intensity based* and *timbre based* features

- **feature categories**

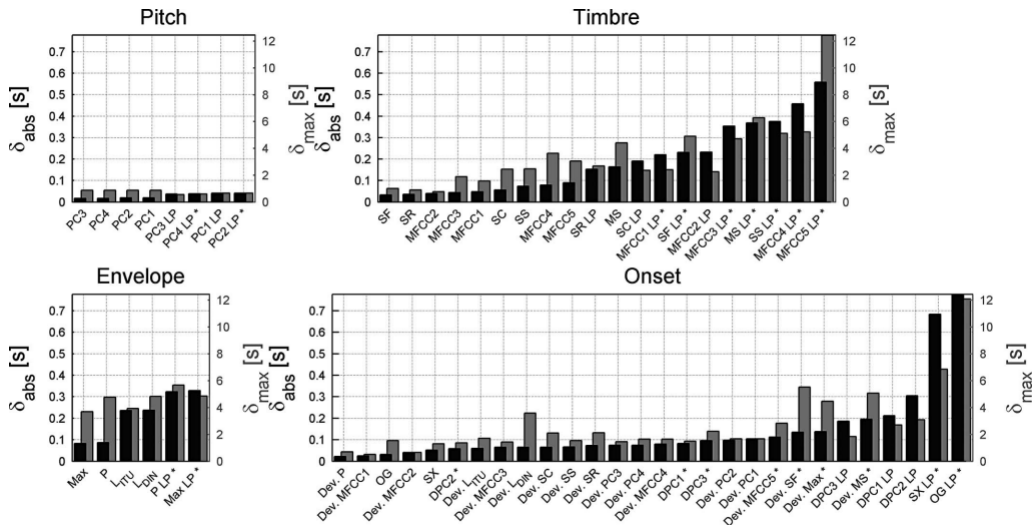
- **intensity**: energy, onset probability, ...
- **tonal**: pitch chroma, ...
- **timbral**: MFCCs, spectral shape, ...

plot from¹

¹H. Kirchhoff and A. Lerch, "Evaluation of Features for Audio-to-Audio Alignment," *Journal of new music research*, vol. 40, no. 1, pp. 27–41, 2011. DOI: [10.1080/09298215.2010.529917](https://doi.org/10.1080/09298215.2010.529917).

audio-to-audio alignment

features



audio-to-audio alignment

compute distance matrix — distance measures

- typical distance measures

- Euclidean distance:* $d_E(s) = \sqrt{\sum_{j=0}^{11} (\nu_e(j) - \nu_{t,s}(j))^2}$

- Manhattan distance:* $d_M(s) = \sum_{j=0}^{11} |\nu_e(j) - \nu_{t,s}(j)|$

- Cosine distance:* $d_C(s) = 1 - \left(\frac{\sum_{j=0}^{11} \nu_e(j) \cdot \nu_{t,s}(j)}{\sqrt{\sum_{j=0}^{11} \nu_e(j)^2} \sqrt{\sum_{j=0}^{11} \nu_{t,s}(j)^2}} \right)$

- Kullback-Leibler divergence:* $d_{KL}(s) = \sum_{j=0}^{11} \nu_e(j) \cdot \log \left(\frac{\nu_e(j)}{\nu_{t,s}(j)} \right)$

- data-driven approach: train classifier with 2-class problem¹

¹H. Kirchhoff and A. Lerch, "Evaluation of Features for Audio-to-Audio Alignment," *Journal of new music research*, vol. 40, no. 1, pp. 27–41,

audio-to-audio alignment

compute distance matrix — distance measures

- typical distance measures

- Euclidean distance:* $d_E(s) = \sqrt{\sum_{j=0}^{11} (\nu_e(j) - \nu_{t,s}(j))^2}$

- Manhattan distance:* $d_M(s) = \sum_{j=0}^{11} |\nu_e(j) - \nu_{t,s}(j)|$

- Cosine distance:* $d_C(s) = 1 - \left(\frac{\sum_{j=0}^{11} \nu_e(j) \cdot \nu_{t,s}(j)}{\sqrt{\sum_{j=0}^{11} \nu_e(j)^2} \sqrt{\sum_{j=0}^{11} \nu_{t,s}(j)^2}} \right)$

- Kullback-Leibler divergence:* $d_{KL}(s) = \sum_{j=0}^{11} \nu_e(j) \cdot \log \left(\frac{\nu_e(j)}{\nu_{t,s}(j)} \right)$

- data-driven approach: train classifier with 2-class problem¹

¹H. Kirchhoff and A. Lerch, "Evaluation of Features for Audio-to-Audio Alignment," *Journal of new music research*, vol. 40, no. 1, pp. 27–41,

audio-to-audio alignment

compute distance matrix — distance measures

- typical distance measures

- Euclidean distance*: $d_E(s) = \sqrt{\sum_{j=0}^{11} (\nu_e(j) - \nu_{t,s}(j))^2}$

- Manhattan distance*: $d_M(s) = \sum_{j=0}^{11} |\nu_e(j) - \nu_{t,s}(j)|$

- Cosine distance*: $d_C(s) = 1 - \left(\frac{\sum_{j=0}^{11} \nu_e(j) \cdot \nu_{t,s}(j)}{\sqrt{\sum_{j=0}^{11} \nu_e(j)^2} \sqrt{\sum_{j=0}^{11} \nu_{t,s}(j)^2}} \right)$

- Kullback-Leibler divergence*: $d_{KL}(s) = \sum_{j=0}^{11} \nu_e(j) \cdot \log \left(\frac{\nu_e(j)}{\nu_{t,s}(j)} \right)$

- data-driven approach: train classifier with 2-class problem¹

¹H. Kirchhoff and A. Lerch, "Evaluation of Features for Audio-to-Audio Alignment," *Journal of new music research*, vol. 40, no. 1, pp. 27–41,

audio-to-audio alignment

compute distance matrix — distance measures

- typical distance measures

- Euclidean distance:* $d_E(s) = \sqrt{\sum_{j=0}^{11} (\nu_e(j) - \nu_{t,s}(j))^2}$

- Manhattan distance:* $d_M(s) = \sum_{j=0}^{11} |\nu_e(j) - \nu_{t,s}(j)|$

- Cosine distance:* $d_C(s) = 1 - \left(\frac{\sum_{j=0}^{11} \nu_e(j) \cdot \nu_{t,s}(j)}{\sqrt{\sum_{j=0}^{11} \nu_e(j)^2} \sqrt{\sum_{j=0}^{11} \nu_{t,s}(j)^2}} \right)$

- Kullback-Leibler divergence:* $d_{KL}(s) = \sum_{j=0}^{11} \nu_e(j) \cdot \log \left(\frac{\nu_e(j)}{\nu_{t,s}(j)} \right)$

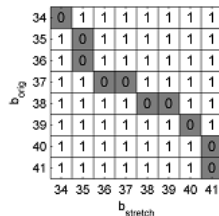
- data-driven approach: train classifier with 2-class problem¹

¹H. Kirchhoff and A. Lerch, "Evaluation of Features for Audio-to-Audio Alignment," *Journal of new music research*, vol. 40, no. 1, pp. 27–41,

audio-to-audio alignment

compute distance matrix — distance measures

- typical distance measures
- data-driven approach: train classifier with 2-class problem¹

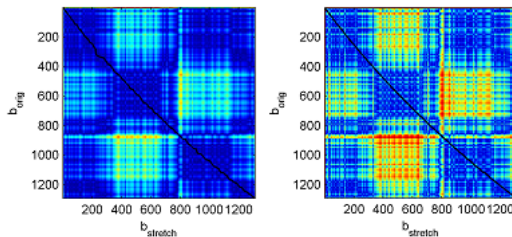


¹H. Kirchhoff and A. Lerch, "Evaluation of Features for Audio-to-Audio Alignment," *Journal of new music research*, vol. 40, no. 1, pp. 27–41, 2011. DOI: [10.1080/09298215.2010.529917](https://doi.org/10.1080/09298215.2010.529917).

audio-to-audio alignment

compute distance matrix — distance measures

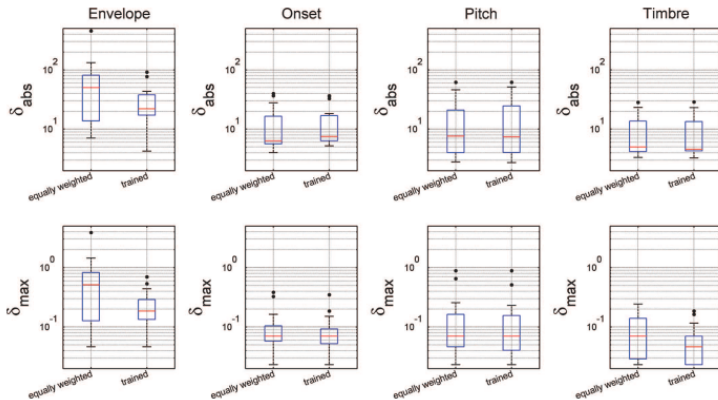
- typical distance measures
- data-driven approach: train classifier with 2-class problem¹



¹H. Kirchhoff and A. Lerch, "Evaluation of Features for Audio-to-Audio Alignment," *Journal of new music research*, vol. 40, no. 1, pp. 27–41, 2011. DOI: [10.1080/09298215.2010.529917](https://doi.org/10.1080/09298215.2010.529917).

audio-to-audio alignment

typical results



originals

synced

left: instrumental

right: a capella



audio-to-score alignment

overview

- **objective**

- align an audio sequence with a score sequence

- **use cases**

- score viewer
- music education
- identify matching score/audio via cost function
- musicological analysis

- **processing steps**

- see audio-to-audio alignment

audio-to-score alignment

overview

- **objective**

- align an audio sequence with a score sequence

- **use cases**

- score viewer
- music education
- identify matching score/audio via cost function
- musicological analysis

- **processing steps**

- see audio-to-audio alignment

audio-to-score alignment

overview

- **objective**

- align an audio sequence with a score sequence

- **use cases**

- score viewer
- music education
- identify matching score/audio via cost function
- musicological analysis

- **processing steps**

- see audio-to-audio alignment

audio-to-score alignment challenges

- features from **different domains** (no timbre and proper loudness information in the score)
 - *approach 1*: convert score into audio-like representation
 - MIDI-to-audio
 - use model for harmonics and ADSR
 - *approach 2*: convert audio into score-like representation
 - audio-to-MIDI
 - pitch chroma
 - event-based segmentation
- pauses and rests
 - DTW algorithm has no graceful way of dealing with pauses

audio-to-score alignment challenges

- features from **different domains** (no timbre and proper loudness information in the score)
 - *approach 1*: convert score into audio-like representation
 - MIDI-to-audio
 - use model for harmonics and ADSR
 - *approach 2*: convert audio into score-like representation
 - audio-to-MIDI
 - pitch chroma
 - event-based segmentation
- pauses and rests
 - DTW algorithm has no graceful way of dealing with pauses

summary

lecture content

- **audio-to-audio alignment**

- ① extract features
- ② create distance matrix with suitable distance measure
- ③ use DTW to find alignment path
- ④ (use time-stretching to actually align the sequences)

- **audio-to-score alignment**

- ① extract usually pitch-based features
- ② distance measure
- ③ use DTW, HMM, etc to extract alignment path

