

Motivation どんな問題に取り組むのか？

音響信号がどのような音響的“パーツ”によってどのように構成されているか(“設計図”)を情報論的アプローチによって推定する計算理論を研究しています。今回は混合音の構成要素となっている音を抽出する複素NMFと呼ぶ手法を紹介します。

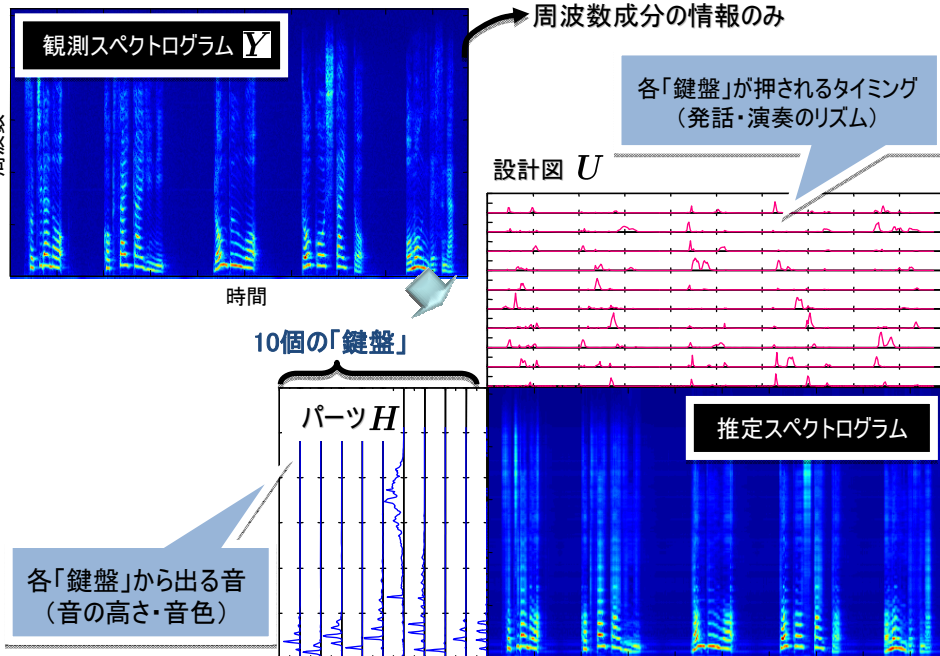
Originality 得られた結果はどう新しいのか？

単純な周波数分析では様々な周波数成分が複雑に重畳するため個々の音の分離抽出は困難でしたが、提案法では繰り返し生起する音を重要な構成要素として学習するアイデアにより、これを可能にします。

Impact この研究が成功した場合のインパクトは？

将来、推定した“パーツ”を選択・置換したり“設計図”を書き換えて音を再構成することで、混合音から目的音を選択的に抽出したり、音の情報の一部、たとえばある音の音色だけを変えたりすることが可能になります。

複素NMF たとえるなら、実世界音響信号を一種の「鍵盤楽器」から生成されたものと捉え、そのもとで、各「鍵盤」が、押された時どんな音を出すか、いつからいつまで押されたか、を観測信号から推定します。



音響信号モデル

(仮定1) 音響信号は I 種類の要素信号だけから構成:

$$F_{k,t} = \sum_{i=1}^I a_{k,t}^{(i)}$$

周波数 \swarrow
時間 \searrow

(仮定2) 各要素信号の周波数成分比は時不変:

$$|a_{k,t}^{(i)}| = H_k^{(i)} U_t^{(i)} \quad (H_k^{(i)} \geq 0, U_t^{(i)} \geq 0),$$

時刻 t でのゲイン \swarrow
周波数成分比 \searrow $\sum_k H_k^{(i)} = 1$

位相スペクトルは時変: $\arg(a_{k,t}^{(i)}) = \phi_{k,t}^{(i)}$

$$Y_{k,t} \simeq F_{k,t} = \sum_i H_k^{(i)} U_t^{(i)} e^{j\phi_{k,t}^{(i)}}$$

観測時間 \swarrow
周波数成分 \searrow

$$\theta := \{H_k^{(i)}, U_t^{(i)}, \phi_{k,t}^{(i)}\}$$

最適化アルゴリズム

$$\begin{aligned} &\text{minimize } f(\theta) := \sum_{k,t} |Y_{k,t} - F_{k,t}|^2 + 2\lambda \sum_{i,t} |U_{i,t}|^p \\ &\text{subject to } \sum_k H_k^{(i)} = 1 \quad (i = 1, \dots, I) \end{aligned}$$

$$\begin{aligned} f^+(\theta, \bar{\theta}) &:= \sum_{i,k,t} \frac{1}{\beta_{k,t}^{(i)}} \left| \bar{Y}_{k,t}^{(i)} - H_k^{(i)} U_t^{(i)} e^{j\phi_{k,t}^{(i)}} \right|^2 \\ &\quad + \lambda \sum_{i,t} \left\{ p |\bar{U}_t^{(i)}|^{p-2} U_t^{(i)2} + (2-p) |\bar{U}_t^{(i)}|^p \right\} \\ \bar{\theta} &:= \{\bar{Y}_{k,t}^{(i)}, \bar{U}_t^{(i)}\} \quad \text{ただし, } \sum_i \bar{Y}_{k,t}^{(i)} = Y_{k,t} \end{aligned}$$

任意の定数 $\beta_{k,t}^{(i)}, p$:
 $\sum_i \beta_{k,t}^{(i)} = 1, \beta_{k,t}^{(i)} \in (0, 1), p \in (0, 2]$

$$\Rightarrow f(\theta) \leq f^+(\theta, \bar{\theta})$$

局所最適解への収束が保証

- Step 1) $\bar{\theta} \leftarrow \operatorname{argmin}_{\bar{\theta}} f^+(\theta, \bar{\theta})$ いずれも閉形式で求まる!
- Step 2) $\theta \leftarrow \operatorname{argmin}_{\theta} f^+(\theta, \bar{\theta})$