

Towards Deep Interpretability (MUS-ROVER II): Learning Hierarchical Representations of Tonal

Haizi Yu¹ Lav R. Varshney¹

¹University of Illinois at Urbana-Champaign

ICLR, 2017

Presenter: Xueying Bai

- 1 Introduction
 - Task
 - Related Work
- 2 Mus-Rover Overview
 - Representation of Data
 - Self-Learning Loop
- 3 Mus-Rover II
 - Hierarchical Rule Learning
 - Memory Selection
- 4 Experiments
 - Experiment Results

- 1 Introduction
 - Task
 - Related Work
- 2 Mus-Rover Overview
 - Representation of Data
 - Self-Learning Loop
- 3 Mus-Rover II
 - Hierarchical Rule Learning
 - Memory Selection
- 4 Experiments
 - Experiment Results

Task: Study underlying patterns beneath the music surface.

- **Discovering compositional rules from raw music data:** music theorists develop concepts and rules to describe the regularity in music compositions;
Computer scientists translated these rules into programs that automatically generate music.
- **Forming hierarchical concepts:** music theorists have devised multi-level analytical methods to emphasize the hierarchical structure of the music.

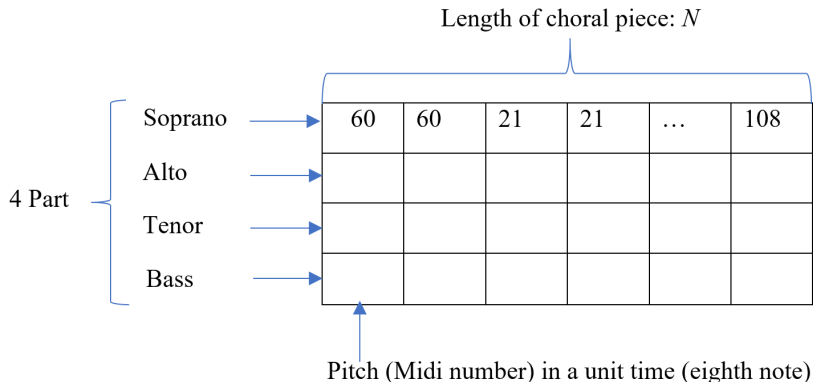
- 1 Introduction
 - Task
 - Related Work
- 2 Mus-Rover Overview
 - Representation of Data
 - Self-Learning Loop
- 3 Mus-Rover II
 - Hierarchical Rule Learning
 - Memory Selection
- 4 Experiments
 - Experiment Results

- **Adversarial or Collaborative:** self-learning loop is similar to GAN.
- **Interpretable Feature Learning:** first recover disentangled representations, then interpret the semantics of the learned features.
- **Automatic Musicians:** models to automate the interaction of music theory and composition.

Outline

- 1 Introduction
 - Task
 - Related Work
- 2 Mus-Rover Overview
 - Representation of Data
 - Self-Learning Loop
- 3 Mus-Rover II
 - Hierarchical Rule Learning
 - Memory Selection
- 4 Experiments
 - Experiment Results

Musical Raw Representation



- Midi number set Ω : $\{21, 22, \dots, 108\}$. Correspond to 88 piano keys.
- N : Length of choral piece.
- Midi matrix $X \in \Omega^{4 \times N}$: Musical raw representation.

Interpretable Feature Representation

- Each sonority p : $p \in \Omega^4$. Each column in Midi matrix.
- Part selection window $w_I : \Omega^4 \mapsto \Omega^{|I|}$. $w_{1,4}(p) = (p_1, p_4)$.
- Basic descriptor B : $\{order, diff, sort, mod_{12}\}$. Atomic arithmetic operations.
- Descriptor with length k : $d_{(k)} = b_k \circ \dots \circ b_1, b_i \in B$
- All windows: $W = \{w_i | I \in 2^{\{1,2,3,4\}} \setminus \{\emptyset\}\}$
- All descriptors with length k : $D^{[k]} = \{d_{(k')} | 0 \leq k' \leq k\}$
- Feature universe: $\Phi = \{d \circ w | w \in W, d \in D^{[k]}\}$
- Any feature $\phi \in \Phi, \Omega^4 \mapsto \phi(\Omega^4)$

Outline

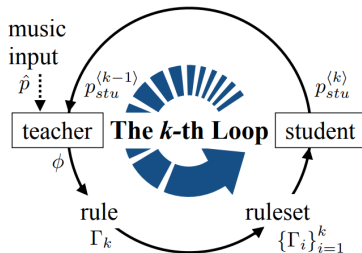
- 1 Introduction
 - Task
 - Related Work
- 2 **Mus-Rover Overview**
 - Representation of Data
 - **Self-Learning Loop**
- 3 Mus-Rover II
 - Hierarchical Rule Learning
 - Memory Selection
- 4 Experiments
 - Experiment Results

Self-Learning Loop

The teacher solves:

$$\begin{aligned} &\text{maximize} && D(p_{\phi,stu}^{(k-1)} \parallel \hat{p}_{\phi}) \\ &\text{subject to} && \phi \in \Phi \setminus \Phi^{(k-1)} \end{aligned}$$

(discrete optimization)



The student solves:

$$\begin{aligned} &\text{maximize} && S_q(p_{stu}^{(k)}) \\ &\text{subject to} && p_{stu}^{(k)} \in \Gamma_1 \\ &&& \dots \\ &&& p_{stu}^{(k)} \in \Gamma_k \end{aligned}$$

(linear least-squares)

- S_q : Tsallis entropy, which achieves a maximum when p is uniform.
- Γ_k : k th rule, containing $(\phi_i, \hat{p}_{\phi_i})$.
- p_{stu}^k : Sonority distribution. (n -gram)
- $p_{\phi,stu}^{<k-1>}$: Feature distribution.
- To get the feature distribution p_{ϕ} from sonority distribution $p(x)$:

$$p_{\phi}(y) = \sum_{x \in \phi^{-1}(\{y\})} p(x)$$

Outline

- 1 Introduction
 - Task
 - Related Work
- 2 Mus-Rover Overview
 - Representation of Data
 - Self-Learning Loop
- 3 Mus-Rover II
 - Hierarchical Rule Learning
 - Memory Selection
- 4 Experiments
 - Experiment Results

- **Feature induced partition of the input domain Ω^4 :**

$$P_\phi = \{\phi^{-1}(\{y\}) \mid y \in \phi(\Omega^4)\}$$

- For two partitions $P, Q \in P_\phi$, P is finer than Q (Q is coarser), written as $P \succeq Q$, if for all $p, p' \in \Omega_4$, p, p' are in the same cluster under $P \Rightarrow p, p'$ are in the same cluster under Q .
 P is strictly finer is written as $P \succ Q$.

Conceptual Hierarchy

Based on the relation \succ , construct the conceptual hierarchy for the partition family P_ϕ as a directed acyclic graph:

- Node: partitions.
- Edge: there's an edge between any pair of nodes v, v' if and only if the partition referred by v is (strictly) finer than v' .
- A higher level feature induces a coarser partition.

Informational Hierarchy

- Given a trace of extracted rules by the k th iteration of loop, a feature ϕ is informationally implied from the extracted rules with tolerance $\gamma > 0$ if:

$$\text{gap}(p_{\phi, dtu}^{<k>} || \hat{p}_{\phi}) := D(p_{\phi, stu}^{<k>} || \hat{p}_{\phi}) < \gamma$$

and

$$\text{gap}(p_{\phi, stu}^{<k'>} || \hat{p}_{\phi}) \geq \gamma, \forall k' < k$$

- This rule is beyond user's satisfaction.

Hierarchical Filters

- Add hierarchical filters when selecting rules to prune hierarchically entangled features and speed up feature selection.
- Teacher's optimization can be represented as:

$$\underset{\phi \in \Phi}{\text{maximize}} \text{gap}(p_{\phi,stu}^{<k-1>} || \hat{p}_{\phi})$$

$$H(\hat{p}_{\phi}) \leq \delta$$

$$\phi \notin C^{(k-1)} := \{\phi | P_{\phi} \prec P_{\phi'}, \phi' \in \Phi^{<k-1>}\}$$

$$\phi \notin I^{(k-1)} := \{\phi | \text{gap}(p_{\phi,stu}^{<k-1>} || \hat{p}_{\phi}) < \gamma\}$$

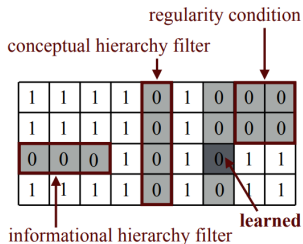
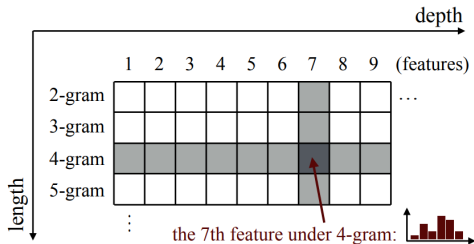
- $H(\hat{p}_{\phi})$ constraint here is based on 'rules should be relatively easy to learn'.
- By selecting hyper-parameters γ and δ , can control the learning pace.

Outline

- 1 Introduction
 - Task
 - Related Work
- 2 Mus-Rover Overview
 - Representation of Data
 - Self-Learning Loop
- 3 Mus-Rover II
 - Hierarchical Rule Learning
 - **Memory Selection**
- 4 Experiments
 - Experiment Results

Two-Dimensional Memory selection

Rule selecting task: pick the right feature under the right n gram.



Outline

- 1 Introduction
 - Task
 - Related Work
- 2 Mus-Rover Overview
 - Representation of Data
 - Self-Learning Loop
- 3 Mus-Rover II
 - Hierarchical Rule Learning
 - Memory Selection
- 4 Experiments
 - Experiment Results

Application

Table 1: Customizing a syllabus (* signifies rules that are skipped in the faster pace)

Rule Trace	Faster Pace ($\gamma = 0.5$)	Slower Pace ($\gamma = 0.1$)
1	<code>order</code> \circ $w_{\{1,2,3,4\}}$	<code>order</code> \circ $w_{\{1,2,3,4\}}$
2	<code>mod</code> ₁₂ \circ $w_{\{1\}}$	<code>order</code> \circ <code>diff</code> \circ <code>sort</code> \circ $w_{\{1,2,4\}}$ *
3	<code>mod</code> ₁₂ \circ <code>diff</code> \circ $w_{\{2,3\}}$	<code>order</code> \circ <code>diff</code> \circ <code>mod</code> ₁₂ \circ $w_{\{1,2,3\}}$ *
4	<code>mod</code> ₁₂ \circ <code>diff</code> \circ $w_{\{3,4\}}$	<code>order</code> \circ <code>diff</code> \circ <code>diff</code> \circ $w_{\{1,2,3,4\}}$ *
5	<code>diff</code> \circ <code>sort</code> \circ $w_{\{2,3\}}$	<code>order</code> \circ <code>sort</code> \circ <code>mod</code> ₁₂ \circ $w_{\{2,3,4\}}$ *
6	<code>mod</code> ₁₂ \circ $w_{\{3\}}$	<code>order</code> \circ <code>sort</code> \circ <code>mod</code> ₁₂ \circ $w_{\{1,3,4\}}$ *
7	<code>mod</code> ₁₂ \circ <code>diff</code> \circ $w_{\{1,2\}}$	<code>order</code> \circ <code>sort</code> \circ <code>mod</code> ₁₂ \circ $w_{\{1,2,3,4\}}$ *
8	<code>mod</code> ₁₂ \circ <code>diff</code> \circ $w_{\{2,4\}}$	<code>mod</code> ₁₂ \circ $w_{\{1\}}$
9	<code>diff</code> \circ $w_{\{1,2\}}$	<code>mod</code> ₁₂ \circ <code>diff</code> \circ $w_{\{2,3\}}$
10	<code>diff</code> \circ <code>sort</code> \circ $w_{\{1,3\}}$	<code>mod</code> ₁₂ \circ <code>diff</code> \circ $w_{\{3,4\}}$

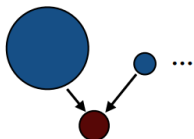
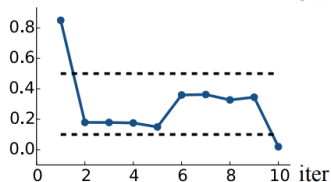
Skip rules that are not so important.

Fundamentals: Hierarchical 1-gram

For fundamentals, knowledge independent of the context: 1-gram.

- With respect to a feature, record the gap between student and Bach for each iteration.

gap $\phi_2 = \text{order} \circ \text{diff} \circ \text{sort} \circ w_{\{1,2,4\}}$



gap $\phi_7 = \text{order} \circ \text{sort} \circ \text{mod}_{12} \circ w_{\{1,2,3,4\}}$

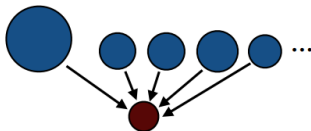
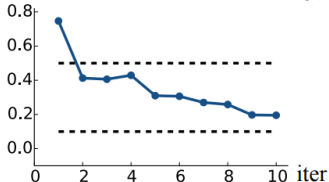


Figure 3: Gap trajectories for two features. The dashed black lines show two different satisfactory gaps ($\gamma = 0.5$ and 0.1). The bottom charts show the informationally implied hierarchies.

Part-writing: n-grams

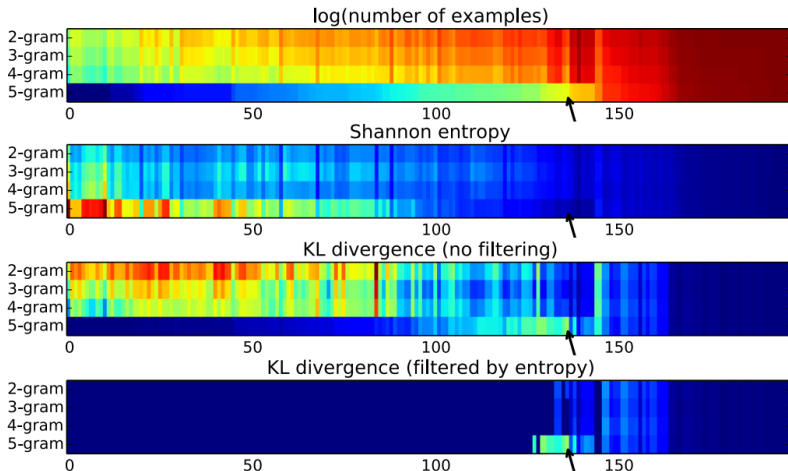


Figure 4: The relative performance of the selected rule (pointed) among the pool of all cells in the 2D memory. A desired rule has: higher confidence (measured by the number of examples, brighter regions in the first row), more regularity (measured by Shannon entropy, darker regions in the second row), and larger style gap (measured by KL divergence, brighter regions in the bottom two rows).

Visualizing Bach's Mind

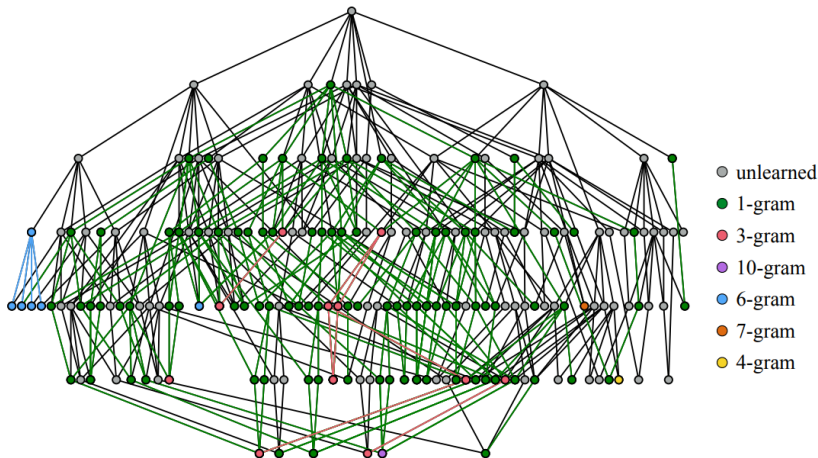


Figure 5: Visualization of Bach's music mind for writing chorales. The underlying DAG represents the conceptual hierarchy (note: edges always point downwards). Colors are used to differentiate rule activations from different n -gram settings. We have enlarged $N = \{1, 2, \dots, 10\}$ to allow even longer-term dependencies.