

MUS-ROVER

An Automatic Music Theorist and Pedagogue

Haizi Yu, Lav R. Varshney

University of Illinois at Urbana-Champaign



Supported in part by C3SR

A Music Story





A Music Story







A Music Story





BWV 124





What makes music music?





"The path towards Mount Parnassus"



"The path towards Mount Parnassus"



the home of the muses

"The path towards Mount Parnassus"



the home of the muses



"The path towards Mount Parnassus"



the home of the muses



HOME LOGIN MUS-ROVER Automatic Pathfinder to Mount Parnassus

About
Demo
Team

Image: Demo

Image: Demo

A New Learning Problem

Automatic Concept Learning

Input

Output



Concepts Rules Laws

• • •

A New Learning Paradigm

New Concept Representation New Learning Algorithm

Representation: Data Space

Data space: (X, p_X) or (X, p) for short

Representation: Data Space

Data space: (X, p_X) or (X, p) for short

Assume a data point $x \in X$ is an i.i.d. sample drawn from a probability distribution p.

Representation: Data Space

Data space: (X, p_X) or (X, p) for short

Assume a data point $x \in X$ is an i.i.d. sample drawn from a probability distribution p.

However, the data distribution p, or an estimation of it, is **known**.

The goal here is not to estimate p but to explain it.

Representation: Chord Space

Considering **chords** from Bach's four-part chorales recorded in sheet music.

Here, a **chord** is any collection of four simultaneously sounding pitches.

Chord space: $X = \mathbb{Z}^4$ chord: $x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} \in X$ pitch: $x_i \in \mathbb{Z}$ (C4 \rightarrow 60) voice: $i \in \{1, 2, 3, 4\}$ S A T B Soprano Hato Bass Bass Bass C3 \rightarrow

An **abstraction** \mathcal{A} is a partition of the data space X.

An **abstraction** \mathcal{A} is a partition of the data space X.

 $X = \{x_1, x_2, x_3, x_4, x_5, x_6\}$ $\mathcal{A} = \{\{x_1, x_6\}, \{x_3\}, \{x_2, x_4, x_5\}\}$

An **abstraction** \mathcal{A} is a partition of the data space X.

 $X = \{x_1, x_2, x_3, x_4, x_5, x_6\}$ $\mathcal{A} = \{\{x_1, x_6\}, \{x_3\}, \{x_2, x_4, x_5\}\}$ cells (or less formally, clusters)

An **abstraction** \mathcal{A} is a partition of the data space X.

 $X = \{x_1, x_2, x_3, x_4, x_5, x_6\}$ $\mathcal{A} = \{\{x_1, x_6\}, \{x_3\}, \{x_2, x_4, x_5\}\}$ cells (or less formally, clusters)

An **concept** is a partition cell.

An **abstraction** \mathcal{A} is a partition of the data space X.

 $X = \{x_1, x_2, x_3, x_4, x_5, x_6\}$ $\mathcal{A} = \{\{x_1, x_6\}, \{x_3\}, \{x_2, x_4, x_5\}\}$ cells (or less formally, clusters)

An **concept** is a partition cell.

A **partition matrix** A is a concise way of representing an abstraction A.

An **abstraction** \mathcal{A} is a partition of the data space X.

 $X = \{x_1, x_2, x_3, x_4, x_5, x_6\}$ $\mathcal{A} = \{\{x_1, x_6\}, \{x_3\}, \{x_2, x_4, x_5\}\}$ cells (or less formally, clusters)

An **concept** is a partition cell.

A partition matrix A is a concise way of representing an abstraction A.

$$A = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 1 & 0 \end{bmatrix}$$

An **abstraction** \mathcal{A} is a partition of the data space X.

 $X = \{x_1, x_2, x_3, x_4, x_5, x_6\}$ $\mathcal{A} = \{\{x_1, x_6\}, \{x_3\}, \{x_2, x_4, x_5\}\}$ cells (or less formally, clusters)

An **concept** is a partition cell.

A partition matrix A is a concise way of representing an abstraction A.

$$A = \begin{bmatrix} 1 & 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 \text{ st cell} \\ 2 \text{ nd cell} \\ 3 \text{ rd cell} \end{bmatrix}$$

Representation: Probabilistic Rules

A **probabilistic rule** is a pair:

 $(\mathcal{A}, p_{\mathcal{A}})$

where \mathcal{A} is an abstraction (partition); $p_{\mathcal{A}}$ is a probability distribution over the abstracted concepts (cells).



Representation Summary

Abstraction:a partitionConcept:a partition cell (cluster)Probabilistic rule:abstraction & distribution

Representation Summary

Abstraction:a partitionConcept:a partition cell (cluster)Probabilistic rule:abstraction & distribution

Automatic concept learning: ?

Automatic Concept Learning is

the process of learning probabilistic rules

But how?











Teacher: a Discriminative Model

Want to find the abstraction which exhibits the largest statistical difference between the student and the input data.

Teacher: a Discriminative Model

The teacher solves an optimization problem:

$$\begin{array}{ll} \underset{\mathcal{A} \in \mathfrak{P}_{X}}{\text{maximize}} & D_{KL} \left(p_{\mathcal{A},stu}^{\langle k-1 \rangle} \parallel p_{\mathcal{A}} \right) \\ \text{subject to} & \text{the abstraction } \mathcal{A} \text{ satisfying} \\ & \text{the memorability condition} \\ & \text{the hierarchy condition} \end{array}$$

. . .
The teacher solves an optimization problem:

$$\begin{array}{ll} \underset{\mathcal{A} \in \mathfrak{P}_{X}}{\text{maximize}} & D_{KL} \left(p_{\mathcal{A},stu}^{\langle k-1 \rangle} \parallel p_{\mathcal{A}} \right) \\ \text{subject to} & \text{the abstraction } \mathcal{A} \text{ satisfying} \\ & \text{the memorability condition} \\ & \text{the hierarchy condition} \end{array}$$

What does \mathfrak{P}_X come from?

 \mathfrak{P}_X : abstraction universe

\mathfrak{P}_X : abstraction universe

Mathematically, a partition lattice, which is a special type of partially ordered set

\mathfrak{P}_X : abstraction universe

Mathematically, a partition lattice, which is a special type of partially ordered set

Pictorially, a directed acyclic graph (vertex: partition; edge: coarser than)

\mathfrak{P}_X : abstraction universe

Mathematically, a partition lattice, which is a special type of partially ordered set

Pictorially, a directed acyclic graph (vertex: partition; edge: coarser than)



How to construct $\mathcal{A} \in \mathfrak{P}_X$ (abstraction universe) ?

How to construct $\mathcal{A} \in \mathfrak{P}_X$ (abstraction universe) ?

- Feature-Induced Partition
- Symmetry-Generated Partition

How to construct $\mathcal{A} \in \mathfrak{P}_X$ (abstraction universe) ?

- Feature-Induced Partition
- Symmetry-Generated Partition

In both cases, the underlying mechanism that generates the partition is human-interpretable.

We do not consider arbitrary partitions.

How to construct $\mathcal{A} \in \mathfrak{P}_X$ (abstraction universe) ?

- **Solution** Feature-Induced Partition
 - Symmetry-Generated Partition

In both cases, the underlying mechanism that generates the partition is human-interpretable.

We do not consider arbitrary partitions.

How to construct $\phi \in \Phi$ (feature universe) ? $\phi = d \circ w$

How to construct $d \in D$ (discriptors) and $w \in W$ (windows) ?

$$d = b_k \circ \dots \circ b_1, \quad b_i \in B$$
$$w = w_I, \quad I \subseteq \{1, 2, 3, 4\}$$

How to construct $\phi \in \Phi$ (feature universe) ? $\phi = d \circ w$

How to construct $d \in D$ (discriptors) and $w \in W$ (windows) ?

$$d = b_k \circ \dots \circ b_1, \quad b_i \in B$$
$$w = w_I, \quad I \subseteq \{1, 2, 3, 4\}$$

For example: $d = \text{mod}_{12} \circ \text{diff}, \quad w = w_{\{1,4\}}$

The teacher solves an optimization problem:

$$\begin{array}{ll} \underset{\mathcal{A} \in \mathfrak{P}_{X}}{\text{maximize}} & D_{KL} \left(p_{\mathcal{A},stu}^{\langle k-1 \rangle} \parallel p_{\mathcal{A}} \right) \\ \text{subject to} & \text{the abstraction } \mathcal{A} \text{ satisfying} \\ & \text{the memorability condition} \\ & \text{the hierarchy condition} \end{array}$$

. . .

The teacher solves an optimization problem:

$$\begin{array}{ll} \underset{\mathcal{A} \in \mathfrak{P}_{X}}{\text{maximize}} & D_{KL} \left(p_{\mathcal{A},stu}^{\langle k-1 \rangle} \parallel p_{\mathcal{A}} \right) \\ \text{subject to} & \text{the abstraction } \mathcal{A} \text{ satisfying} \\ & \text{the memorability condition} \\ & \text{the hierarchy condition} \end{array}$$



The teacher solves an optimization problem:

$$\begin{array}{ll} \underset{\mathcal{A} \in \mathfrak{P}_{X}}{\text{maximize}} & D_{KL} \left(p_{\mathcal{A},stu}^{\langle k-1 \rangle} \parallel p_{\mathcal{A}} \right) \\ \text{subject to} & \text{the abstraction } \mathcal{A} \text{ satisfying} \\ & \text{the memorability condition} \\ & \text{the hierarchy condition} \end{array}$$



The teacher solves an optimization problem:

$$\begin{array}{ll} \underset{\mathcal{A} \in \mathfrak{P}_{X}}{\text{maximize}} & D_{KL} \left(p_{\mathcal{A},stu}^{\langle k-1 \rangle} \parallel p_{\mathcal{A}} \right) \\ \text{subject to} & \text{the abstraction } \mathcal{A} \text{ satisfying} \\ & \text{the memorability condition} \\ & \text{the hierarchy condition} \end{array}$$



Apply probabilistic rules, which is known as the rule realization problem.







Want to find the probabilistic model which enables novelty while at the same time satisfies all the rules.

The student solves another optimization problem:

$$\begin{array}{ll} \underset{p_{stu}^{\langle k \rangle} \in \Delta_{|X|}}{\text{maximize}} & S_q(p_{stu}^{\langle k \rangle}) := (q-1)^{-1} \left(1 - \| p_{stu}^{\langle k \rangle} \|_q^q \right) \\ \text{subject to} & A^{(i)} p_{stu}^{\langle k \rangle} = p_{\mathcal{A}^{(i)}}, \quad i = 1, \dots, k \end{array}$$

The student solves another optimization problem:

Tsallis entropy: measures randomness maximize $P_{stu}^{\langle k \rangle} \in \Delta_{|X|}$ subject to $A^{(i)}p_{stu}^{\langle k \rangle} = p_{\mathcal{A}^{(i)}}, \quad i = 1, \dots, k$

The student solves another optimization problem:

Tsallis entropy: measures randomness maximize $figs_q^{\langle k \rangle}(p_{stu}^{\langle k \rangle}) := (q-1)^{-1} \left(1 - \|p_{stu}^{\langle k \rangle}\|_q^q\right)$ subject to $A^{(i)}p_{stu}^{\langle k \rangle} = p_{\mathcal{A}^{(i)}}, \quad i = 1, \dots, k$ partition matrix: represents abstraction

The student solves another optimization problem:

Tsallis entropy: measures randomness maximize $S_q(p_{stu}^{\langle k \rangle}) := (q-1)^{-1} \left(1 - \|p_{stu}^{\langle k \rangle}\|_q^q\right)$ subject to $A^{(i)}p_{stu}^{\langle k \rangle} = p_{\mathcal{A}^{(i)}}, \quad i = 1, \dots, k$ partition matrix: represents abstraction linear equality constraint

The student solves another optimization problem:

Tsallis entropy: measures randomness maximize $S_q(p_{stu}^{\langle k \rangle}) := (q-1)^{-1} \left(1 - \|p_{stu}^{\langle k \rangle}\|_q^q\right)$ subject to $A^{(i)}p_{stu}^{\langle k \rangle} = p_{\mathcal{A}^{(i)}}, \quad i = 1, \dots, k$ partition matrix: represents abstraction linear equality constraint

q = 2: gini impurity function

The student solves another optimization problem:

Tsallis entropy: measures randomness maximize $S_q(p_{stu}^{\langle k \rangle}) := (q-1)^{-1} \left(1 - \|p_{stu}^{\langle k \rangle}\|_q^q\right)$ subject to $A^{(i)}p_{stu}^{\langle k \rangle} = p_{\mathcal{A}^{(i)}}, \quad i = 1, \dots, k$ partition matrix: represents abstraction linear equality constraint

q = 2: gini impurity function Linear Least-Squares Problem!

How MUS-ROVER Self-Evolves?

Context-Free Rules (1-gram) Rules and Outcomes

Student 0







(Spacing) Almost always, the soprano pitch is above the alto, alto above tenor, and tenor above bass.

Student 1





Rule 2: $mod_{12} \circ w_1$



Rule 2: $mod_{12} \circ w_1$

(Scale) The soprano voice is drawn from a diatonic scale with high probability.

Student 2


Fundamentals (1-gram)

Student 22



Fundamentals (1-gram)

Student 22



Context-Specific Rules (n-gram)

Bach's Music Brain

Part Writing (n-gram): 14th Chord









Part Writing (n-gram): 14th Chord • unlearned 1-gram Loop 4 • 3-gram 10-gram • 6-gram 7-gram docoqu 000

Part Writing (n-gram): 14th Chord



Part Writing (n-gram): 14th Chord



Part Writing (n-gram): 14th Chord







Part Writing (n-gram): 14th Chord



Part Writing (n-gram): 14th Chord



Part Writing (n-gram): 14th Chord



Generalizing to Other Topic Domains

- Physics
- Cancer biology
- and more

References

- Haizi Yu, Igor Mineyev, and Lav R. Varshney. A Group-Theoretic Approach to Abstraction: Hierarchical, Interpretable, and Task-Free Clustering. *arXiv:1807.11167v1 [cs.LG], 2018*.
- Haizi Yu, Tianxi Li, and Lav R. Varshney. Probabilistic rule realization and selection. In *Proceedings of NIPS 2017*.
- Haizi Yu, and Lav R. Varshney. Towards deep interpretability (MUS-ROVER II): learning hierarchical representations of tonal music. In *Proceedings of ICLR 2017*.
- Haizi Yu, Lav R. Varshney, Guy E. Garnett, and Ranjitha Kumar. MUS-ROVER: A self-learning system for musical compositional rules. In *Proceedings of MUME 2016*.
- Haizi Yu, Lav R. Varshney, Guy E. Garnett, and Ranjitha Kumar. Learning interpretable musical compositional rules and traces. *In ICML* 2016 Workshop on Human Interpretability in Machine Learning.