

# Two-Part Inventions by SARSA

(Insert Logo)

Author's name(s)  
Media Informatics Special Interest Group,  
Affiliations

MIWAI XXXX

E-mail

## Abstract

## Principle

## Implementation

## Results

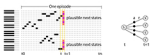
A musical agent learns to generate a two-part invention using SARSA. SARSA is a reinforcement learning technique that learns an optimal policy by sampling the state space to estimate the utility of state-action pairs  $Q(s,a)$  where  $s$  denotes a state,  $a$ , denotes an action,  $r$  denotes a reward,  $\alpha$  denotes the learning rate and  $\gamma$  denotes the discount rate.

$$Q(s,a) \leftarrow Q(s,a) + \alpha(r + \gamma Q(s',a') - Q(s,a))$$

- First, the policy was learned using hand-crafted rules describing the desired characteristics of two-part inventions. These rules could also be discovered using data mining techniques.

- Then, the rules acted as a critic's comments to the generated music. The musical agent would amend its policy based on these comments.

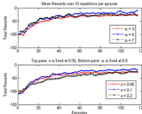
In our approach, each episode was a complete 32-bar two-part counterpoint. Form and other contexts were incorporated into the system via the critic's rules and the usage of context dependent Q-tables.



Policy learning in RL is a powerful concept. An agent explores a partially observable environment until it learns a policy (i.e., how it should react to the environment) that maximises its return, SR $\bar{S}$ . The representation of the state space,  $S$ , and actions,  $A$ , are critical since they are the abstraction of behaviours to be learned. In further work, the following directions could be pursued:

- (i) to improve the handcrafted rules for different composition,
- (ii) to automate rules-acquisition process, and
- (iii) to apply the approach to other genres (e.g., four part writing, Jazz, etc).

In this work, we employed SARSA to generate 32-bar two-part invention pieces. By carefully selecting the representation of states, actions, rules and contexts, a complex problem such as algorithmic composition could be dealt with and reasonable output could be obtained with comparatively less effort.



Criteria	Reward value
Parallel 5th, octave	-0.1
Crossing between parts	-0.1
Spacing between voice more than one octave	-0.1
Repeated notes	-0.1
Repeated consonant major, minor third	-0.1
Repeated consonant major, minor sixth	-0.1
Wide leap interval	-0.1
Dissonant progression, second, seventh	-0.1
Consonant progression major, minor third	0.1
Contrary motion	0.1

SARSA Parameter Settings				
Learning rate ( $\alpha$ )	0.1, 0.2, 0.5, 0.7	0.5		
Discount rate ( $\gamma$ )	0.9	0.9		
( $\lambda$ )-greedy probability	0.1	0.01, 0.1, 0.2		
Max-iteration	200	200		
Max-episode	120	120		

SARSA Parameter Settings				
Learning rate ( $\alpha$ )	0.1, 0.2, 0.5, 0.7	0.5		
Discount rate ( $\gamma$ )	0.9	0.9		
( $\lambda$ )-greedy probability	0.1	0.01, 0.1, 0.2		
Max-iteration	200	200		
Max-episode	120	120		

## References

- Sutton, R.S. and Barto, A.G.: Reinforcement Learning: An Introduction. A Bradford Book, The MIT Press, 1998.
- Watkins, C.J. and Dayan, P.: Q-learning. Machine Learning 8:279-292, 1992.
- Sonmuk Phon-Amnuaisuk: Generating Tonal Counterpoint Using Reinforcement Learning. ICONIP (1) 2009: 580-589