

Исследование алгоритма обучения функциональных систем в стохастической среде

Сорокин А. Ю.
Руководитель: Бурцев М. С.

Кафедра математических и информационных технологий
Академический университет

2014

Нейроморфные системы

Под нейроморфными системами понимаются модели искусственных нейронных сетей, архитектура или свойства, которых основаны на особенностях реальных нейробиологических систем.



Цель и Задачи

Цель:

Исследовать модель функциональных систем на возможность обучения в стохастической среде

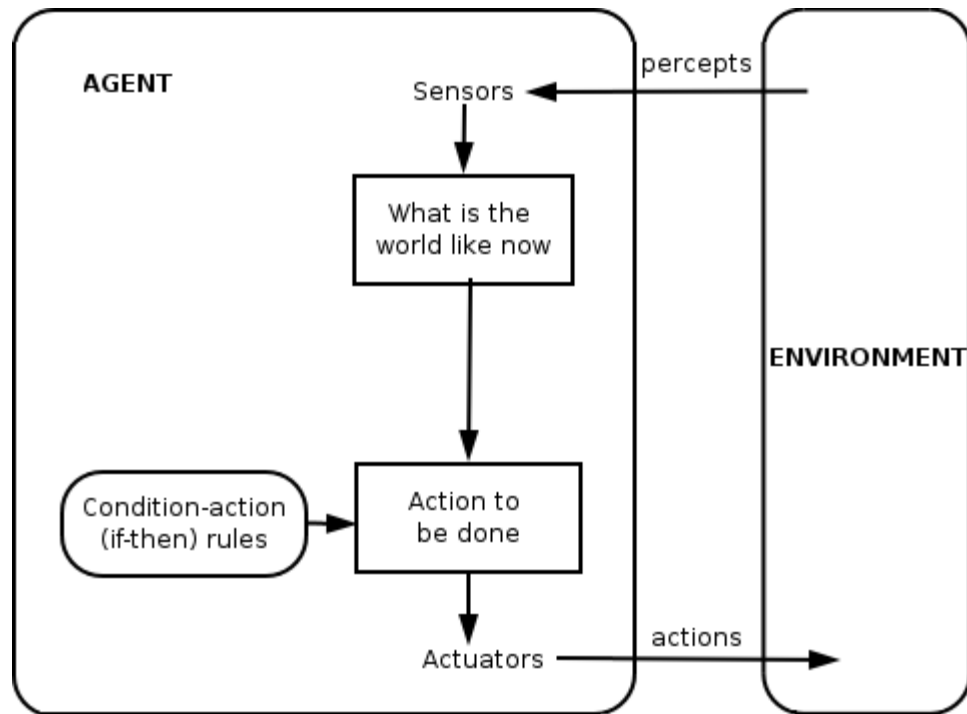
Задачи:

1. Реализовать библиотеку на Python
2. Воспроизвести результаты из статьи[[1](#)]
3. Исследовать обучение модели в стохастической среде

Агентный подход

Цикл взаимодействия:

1. Агент получает информацию о текущем состоянии среды и возможных вознаграждениях
2. Выбирает действие для изменения среды
3. Среда меняет свое в ответ на действие агента.



Функциональные системы

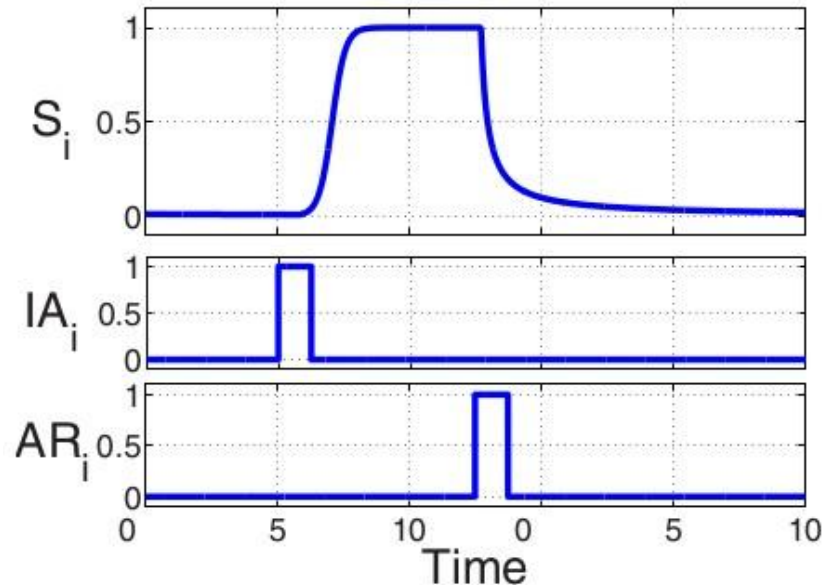
Функциональная система - динамический элемент с двумя стационарными состояниями

- Состояние активности
- Состояние покоя

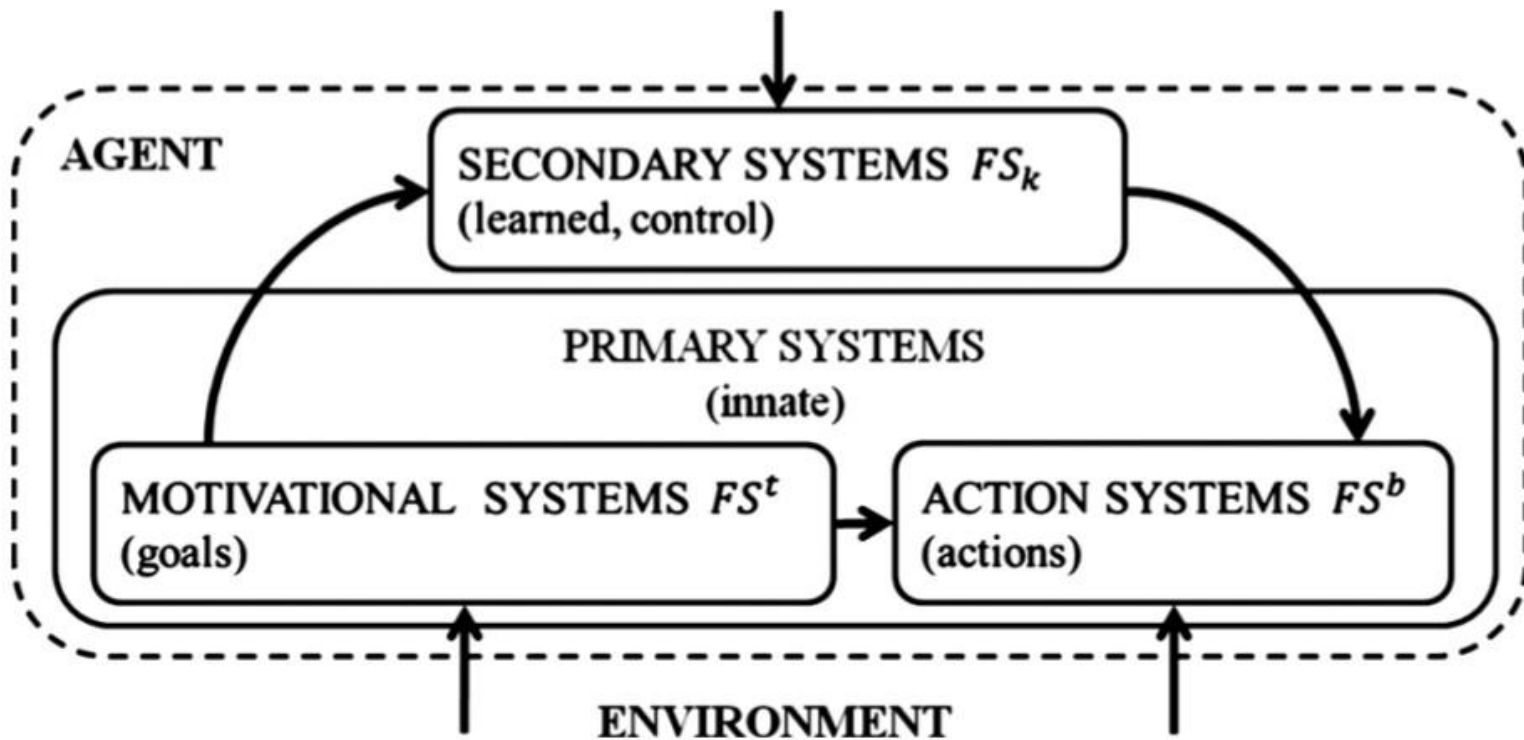
S - активность функциональной системы

IA - входная афферентация

AR - акцептор результата действия



Сеть функциональных систем



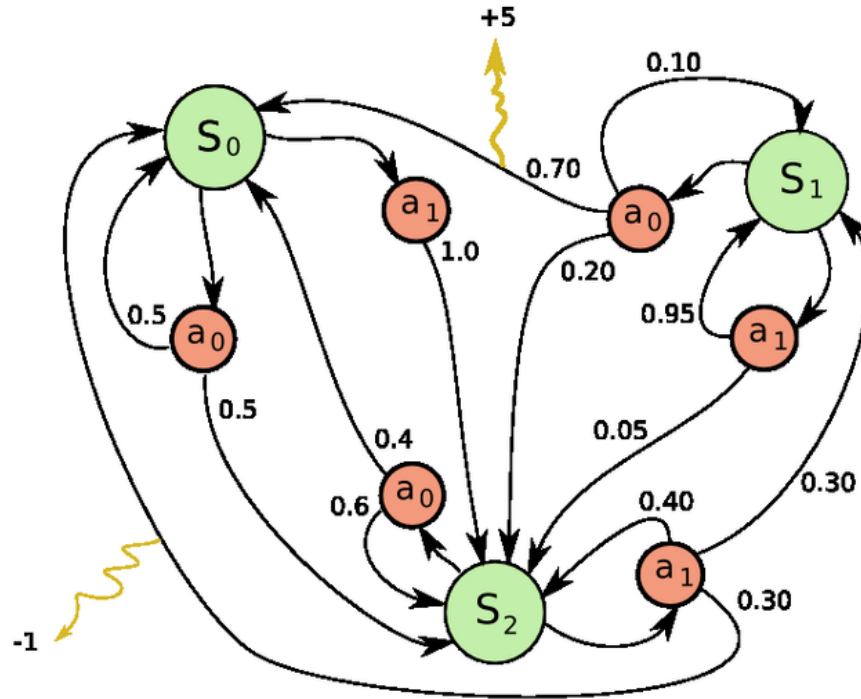
Среда

Среда - граф, где вершины соответствуют различным возможным состояниям, а исходящие ребра являются возможными действиями в этом состоянии.

Стохастическая среда:

- Некоторым переходы в среде присутствуют лишь с определенной вероятностью.
- В каждом конкретном испытании, такой переход либо возможен, либо нет.
- Существуют переходы среди, которых в каждый конкретный момент присутствует только один.
- Путь

Марковский процесс принятия решений



Существующие подходы

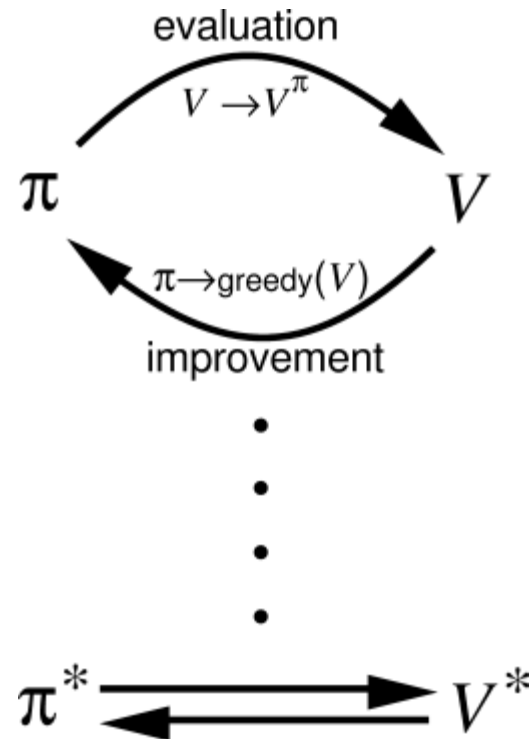
- Алгоритмы обучения нейронных сетей
 - Алгоритм обратного распространения ошибки
 - Недостатки:
 - Необходимо значительно изменить структуру сети
 - Сложно разделить случаи вознаграждения и наказания в испытании
- Динамическое программирование
 - Недостатки:
 - Нужно знать точную модель среды

Существующие подходы

- Методы Монте-Карло
 - Подходят, но в среднем работают хуже, чем TD-обучение
- Обучение на основе временных различий (TD-обучение)
 - Q-обучение и алгоритмы с разделенной оценкой стратегий
 - Недостатки: не подходят по идейным соображениям
 - SARSA и алгоритмы с интегрированной оценкой
 - Алгоритмы Исполнитель-Критик

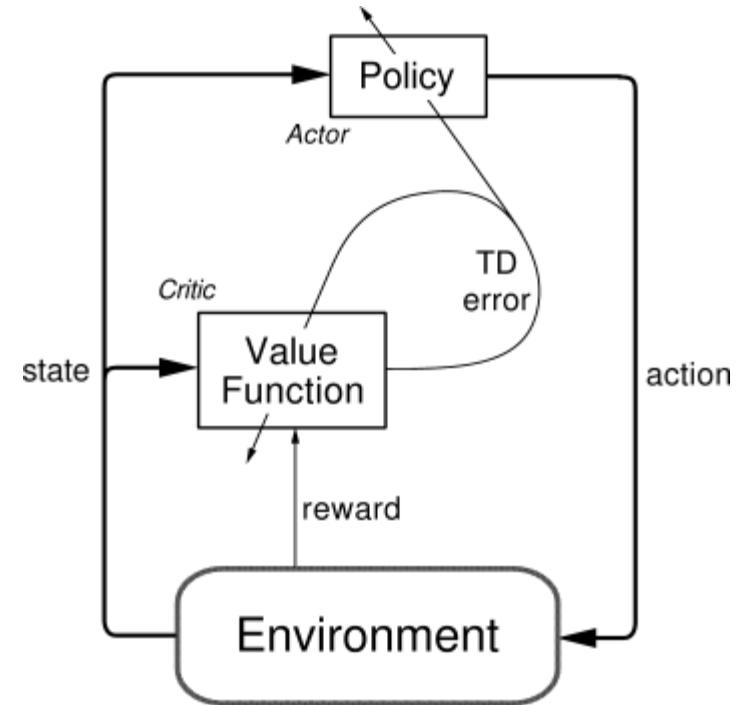
Обобщенная итерация по стратегиям

- Большинство описанных методов является частными случаями Обобщенной Итерации по Стратегиям
- Основные фазы:
 - Оценивание стратегии
 - Улучшение стратегии



Исполнитель - Критик

- Методы Исполнитель-Критик имеют различную структуру памяти для явного изменения стратегии не зависимо от функции ценности.
- Исполнитель и Критик взаимодействуют только, через передачу значения ошибки от критика исполнителю.

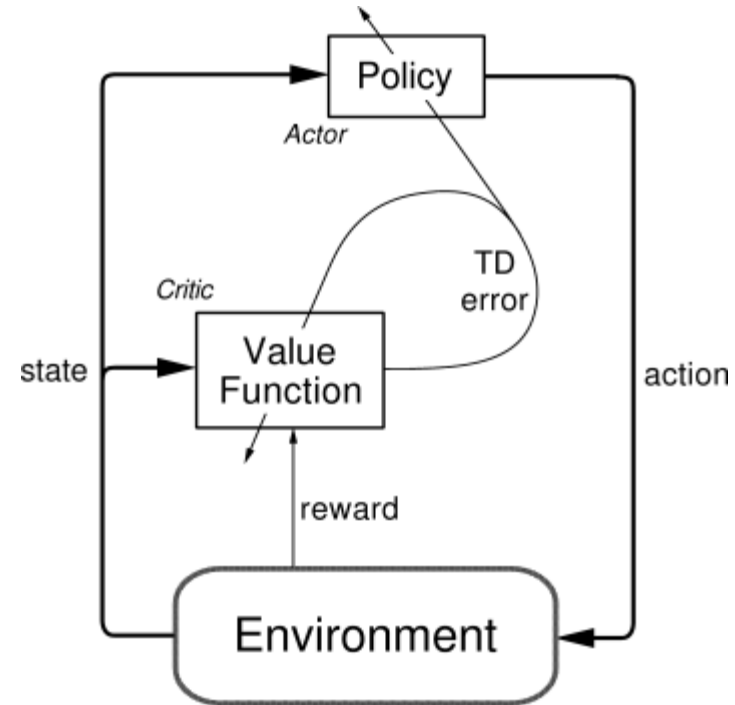


Исполнитель - Критик

В качестве исполнителя выбран алгоритм Direct Actor

- В отличие от других алгоритмов обновляет не только значение совершенного действия, но и всех альтернативных вариантов действия

В качестве критика используется модифицированная версия SARSA



Недостатки Методов TD-обучения (1)

Медленное распространение ошибки:

- Пошаговое обновление
- Угасание ошибки при распространении:

$$\Delta Q(s_t, a_t) = \alpha[R_t - Q(s_t, a_t)] \quad R_t = r_{t+1} + \gamma Q(s_{t+1}, a_{t+1})$$

Решение:

- Обновление всех посещенных состояний по окончании эпизода, как в методах Монте-Карло
- Модификация техники n-шагового дублирования вознаграждений:

$$R_t = r_{t+1} + \dots + r_{t+n} + \gamma Q(s_{t+n}, a_{t+n})$$

Недостатки Методов TD-обучения (2)

Эпизод всегда заканчивается удовлетворением потребности:

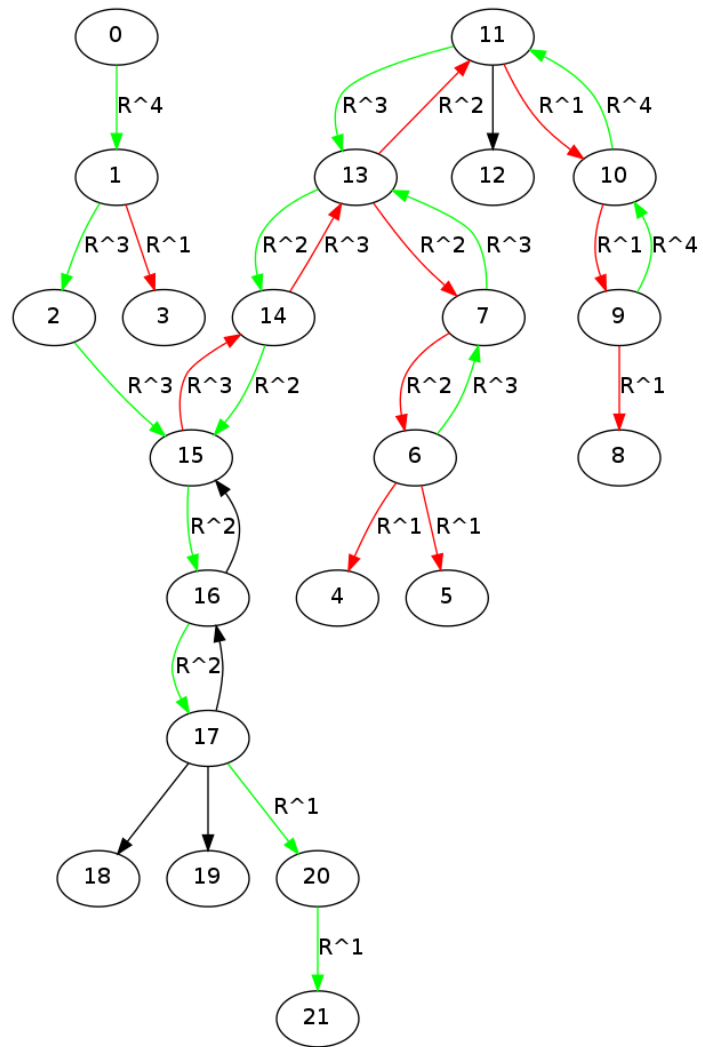
- Сложно определить какие именно действия в эпизоде были полезны, а какие вредны.

Решение:

- Пенальти за каждое попадание в тупик
- Специальное правило для получения оценки ожидаемого вознаграждения:

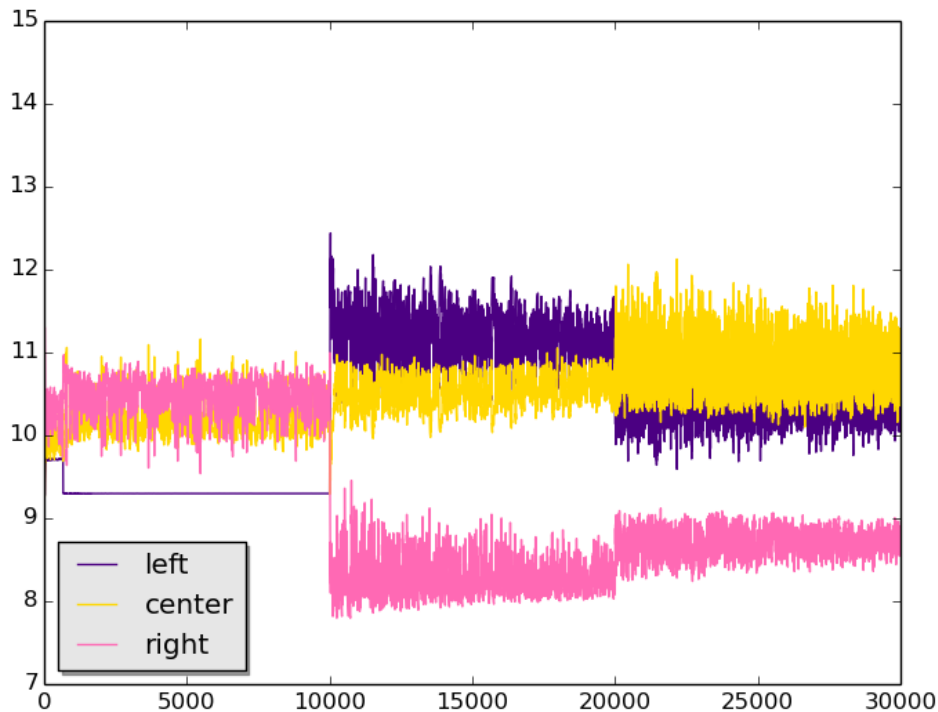
$$a' = \operatorname{argmax}_{a_{t+k}}(R_{t+k}), \forall a_{t+k} \in A(s')$$

$$R_t = r_{t+1} + Q(s', a')$$



Поведение алгоритма в нестационарной среде

- Первый этап:
 - $\text{Prob}(\text{left}) = 0.3$
 - $\text{Prob}(\text{center}) = 0.7$
 - $\text{Prob}(\text{right}) = 0.9$
- Второй этап:
 - $\text{Prob}(\text{left}) = 0.9$
 - $\text{Prob}(\text{center}) = 0.1$
 - $\text{Prob}(\text{right}) = 0.5$
- Третий этап:
 - $\text{Prob}(\text{left}) = 0.9$
 - $\text{Prob}(\text{center}) = 0.1$
 - $\text{Prob}(\text{right}) = 0.5$



Результаты:

- Реализована библиотека на Python для работы с моделью функциональных систем
- Разработан алгоритм способный к адаптивному обучению в условиях стохастической среды и сохраняющий основные свойства модели и способный сравнивать эффективность различных вариантов действий и выбирать более эффективные из них.

Спасибо за внимание!