

Архитектура сети

Вадим Кузьмин

17 июля 2018 г.

1 Описание модели

На вход модели подается изображение, описывающее состояние среды в текущий момент, и one-hot вектор длины *количество опций плюс один*, показывающий то, действие или какая опция выполнялась на предыдущем шаге.

Изображение используется сверточными слоями для получения признакового описания состояния среды, которое далее используется как вход для других частей модели. В модели для реализации иерархии вводятся подзадачи. Каждая подзадача, кроме нулевой, определяет опцию. Нулевая задача пытается решить исходную проблему не зная об опциях. Для выбора задач используется менеджер, который, принимая на вход признаковое описание состояния среды, играет роль Q-функции над задачами, формируя вектор, определяющий подзадачу.

По вектору предыдущего действия решается, передать управление менеджеру, так как до этого не выполнялась опция, или по тому, какая была опция, проверить, закончилась ли она. Если она закончилась, то управление передается менеджеру и формируется вектор, определяющий подзадачу. Если не закончилась, то опция продолжается и вектор, определяющий подзадачу, не меняется. Чтобы проверить, закончилась ли опция, используется отдельный классификатор, который представляет собой стандартный перцептрон.

По вектору, определяющему подзадачу, выбирается выход нужной сети, соответствующей текущей подзадаче. Так формируется выход модели - вектор длины *кол-во базовых действий*, который представляет собой набор значений Q-функции для данного состояния и подзадачи.

Рассмотрим более подробно граф, визуализированный в TensorBoard. Предположим, что используется две опции.

Входные данные:

- input_image - RGB изображение состояния среды;
- prev_action - one-hot вектор длины *число опций плюс один*, показывает, какая задача выполнялась на предыдущем шаге (базовое действие или опция).

Обучаемые модули:

- convolution - сверточные слои, на входе - input_image, на выходе - вектор признаков;
- opt1_checker, opt2_checker - классификаторы для соответствующих опций, на входе классификатора - результат convolution, на выходе - 0/1 (данная опция закончилась или нет);
- task0, task1, task2 - Q-функции соответствующих задач, для каждой на входе - результат convolution, на выходе - вектор длины *число базовых действий*;
- manager - Q-функция над задачами, на входе - результат convolution, на выходе - вектор длины *число опций плюс один*, определяющий задачу.

Программные модули:

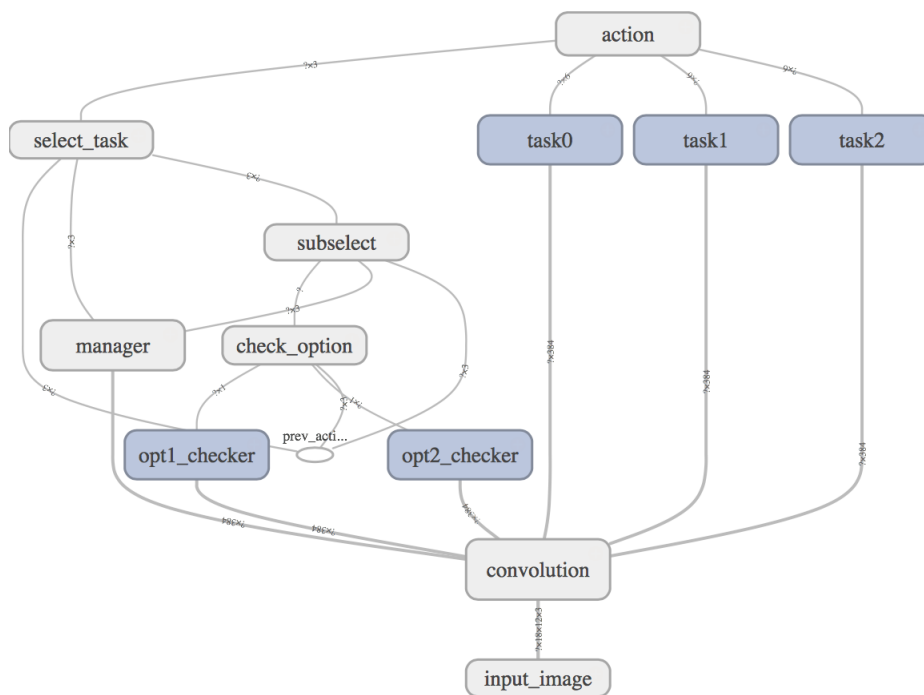


Рис. 1: Модель с двумя опциями

- action - по one-hot вектору, определяющему текущую задачу, выбирает нужный выход из task0, task1, task2;
- select_task - если по prev_action на предыдущем шаге было выполнено базовое действие, то выбирается выход manager, если выполнялась опция - результат работы subselect;
- subselect - выполняется, если на предыдущем шаге выполнялась опция, осуществляется проверка с помощью opt1_checker, opt2_checker, закончилась ли опция, если да, то выбирается выход manager, если не закончилась - выбирается prev_action, так как задача не меняется.

Выходные данные:

- результат выполнения action - вектор длины *число базовых действий*, из которого выбирается действие для агента;
- результат выполнения select_task - будет использован как prev_action.