

Разработка алгоритма на основе методов динамического программирования.

1. Введение в динамическое программирование.

Динамическое программирование – метод решения сложных задач, в котором вместо одной сложной задачи решаются несколько более простых подзадач, на которые она разбита, после чего для каждой подзадачи сохраняются решения [3].

В обучении с подкреплением под динамическим программированием понимают семейство алгоритмов, которые используются для вычисления оптимальных стратегий, при условии, что имеется модель окружающей среды в виде марковского процесса принятия решений [1].

Ключевая идея ДП и ОП вообще заключается в использовании функций ценности для организации и структурирования поиска хороших стратегий. Стратегию можно представить функцией $\pi(s): S \rightarrow A$, указывающей какое действие должно выполняться в каждом состоянии. Конечной целью является нахождение оптимальной стратегии, задающей для каждого состояния действие, максимизирующее конечный доход. Функция ценности состояния указывает насколько хорошо для агента пребывание в конкретном состоянии при стратегии π , то есть оценивает ценность состояния при следовании стратегии и определяется как:

$$\begin{aligned} v_{\pi}(s) &= E_{\pi}[G_t \mid S_t = s] \\ &= E_{\pi}[R_{t+1} + \gamma \sum_a \pi(a|s) \sum_{\acute{s}, r} p(\acute{s}, r|s, a) [r + \gamma E_{\pi}[G_{t+1} \mid S_{t+1} = \acute{s}]]] \\ &= \sum_a \pi(a|s) \sum_{\acute{s}, r} p(\acute{s}, r|s, a) [r + \gamma v_{\pi}(\acute{s})] \text{ для всех } s \in S \end{aligned} \quad (1)$$

где действия a берутся из множества $A(s)$, следующие состояния \acute{s} берутся из множества S , в которые переходит агент после совершенного действия, и вознаграждения r берутся из множества R . Равенство (1) называется уравнением Беллмана для функции ценности состояния s при стратегии π , обозначаемая v_{π} [1].

При этом для детерминированной стратегии и с учетом ограниченного количества шагов получим

$$v_{\pi}(s) = \sum_{\acute{s}, r} p(\acute{s}, r|s, a) [r + \gamma v_{\pi}(\acute{s})] \quad (2)$$

Данное равенство представляет собой систему линейных уравнений, которую можно представить следующим образом

$$\begin{aligned} v_{\pi}(s_n) &= r_n + \gamma * v_{\pi}(s_{n-1}) \\ v_{\pi}(s_{n-1}) &= r_{n-1} + \gamma * v_{\pi}(s_{n-2}) \\ &\dots \\ v_{\pi}(s_1) &= r_1 \end{aligned} \tag{3}$$

Однако для решения системы уравнений (3) требуется прямолинейное длительное вычисление с использованием большого объема памяти, экспоненциально растущего в зависимости от n . В этой связи для решения уравнения Беллмана [1] разработаны итеративные методы:

- итерация по стратегиям;
- итерация по ценности.

1.2 Решение задачи поиска оптимальной стратегии итерацией по стратегиям.

При итерации по стратегиям алгоритм начинается с инициализации ценности всех состояний S случайной величиной (например 0), после чего инициализируется случайная стратегия и проводится расчет ценности по формуле 2. Далее необходимо проверить функцию $V(s)$ на оптимальность, и если функция не оптимальна, провести улучшение стратегии с последующим итерационным вычислением функции ценности. Типовой алгоритм итерации по стратегиям представлен на рисунке 2.

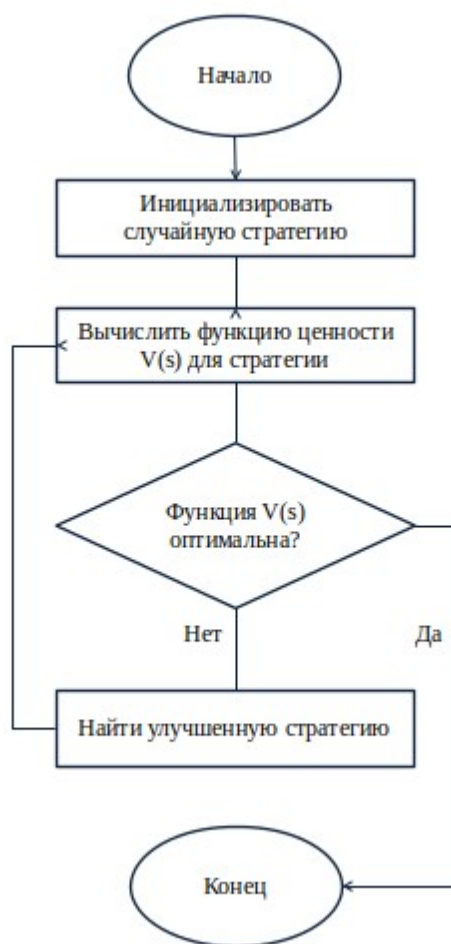


Рисунок 2.

Функция ценности рассматривается как функция состояние/действие и имеет название Q-функция. Она показывает, насколько хорошо для агента выполнять конкретное действие в соответствии со стратегией π и определяется

$$q_{\pi}(s, a) = E_{\pi} [R_{t+1} + \gamma v_{\pi}(S_{t+1}) | S_t = s, A_t = a] = \sum_{s', r} p(s', r | s, a) [r + \gamma v_{\pi}(s')] \quad (4)$$

Q-функция играет важную роль в улучшении стратегии, поскольку если существует такая стратегия, что выполняется следующее неравенство $q_{\pi'}(s, a) > v_{\pi}(s)$ для одного состоя s , при условии, что в остальном стратегия π' соответствует стратегии π , тогда стратегия π' лучше стратегии π . Тогда, при рассмотрении изменения всех возможных действий во всех состояниях, выбирая в каждом состоянии то действие, которое максимизирует Q-функцию, получим новую жадную стратегию π' , определенную следующим образом:

$$\begin{aligned} \pi(s) = \operatorname{argmax}_a q_\pi(s, a) & \quad \operatorname{argmax}_a E[R_{t+1} + \gamma v_\pi(S_{t+1}) | S_t = s, A_t = a] \\ & \quad \operatorname{argmax}_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_\pi(s')] \end{aligned} \quad (5)$$

Жадная стратегия выбирает действие, которое является наилучшим в краткосрочной перспективе согласно функции v_π . Процесс конструирования новой стратегии, улучшающий исходную путем жадного выбора относительно функции ценности исходной стратегии, называется улучшением стратегии.

Итерационное улучшение стратегии необходимо проводить до получения оптимальной стратегии – стратегии максимизирующей вероятный доход. Критерий оптимальности стратегии определяется равенством новой, выработанной на очередном шаге, жадной стратегии π' и стратегии предыдущего шага π , поскольку тогда $v_\pi = v_{\pi'}$ и из формулы 5 следует, что для всех $s \in S$. Таким образом, итерационное улучшение стратегии до оптимальной можно представить следующей схемой:

$$\pi_0 \xrightarrow{E} v_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} v_{\pi_1} \xrightarrow{I} \pi_2 \xrightarrow{E} \dots \xrightarrow{I} \pi_s \xrightarrow{E} v_{\pi_s}$$

где E – это оценивание стратегии, а I – улучшение стратегии.

При этом, в соответствии с вышеизложенным, каждая следующая стратегия лучше предыдущей, только если предыдущая уже не является оптимальной. Для конечного МППР существует конечное число стратегий, поэтому процесс должен сойтись к оптимальной стратегии и оптимальной функции ценности за конечное число итераций. Данный способ нахождения оптимальной стратегии и является способом итерации по стратегиям. Таким образом, полный алгоритм итерации по стратегиям будет выглядеть следующим образом:



Рисунок 3.

Для решения задачи построения траектории облета зон противодействия сил ПВО противника реализовав данный алгоритм в программном обеспечении и рассчитав ценности каждого состояния с учетом текущего расположения ПВО получим:

вставить скриншот ПО

Теперь для построения маршрута облета зон ПВО из любой точки необходимо, всего лишь, построить маршрут по точкам (состояниям) жадно (глубина расчета 1 шаг) выбирая действие максимизирующее значение Q-функции. Выбрав за начальную точку агента $s(2, 2)$, объекты ПВО $s(x, x)$ и $s(x, x)$ и объект назначения $s(x, x)$ получим:

вставить скриншот ПО

Задача построения маршрута облета зон ПВО методом динамического программирования с использованием итерации по стратегиям решена.

1.3 Решение задачи поиска оптимальной стратегии итерацией по ценности.

При итерации по ценности мы сначала инициализируем функцию ценности некоторым случайным значением, после чего перебираем все состояния и находим новую функцию ценности

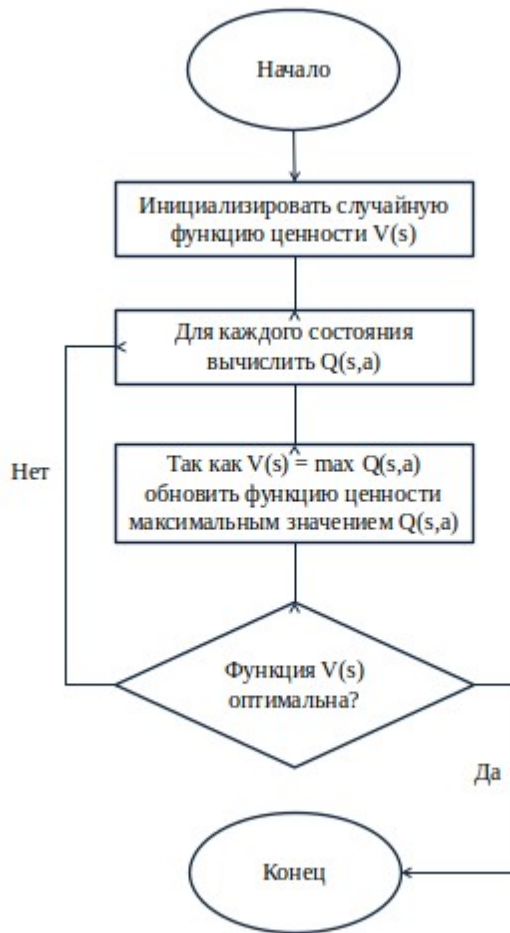


Рисунок 3.

Итерация по ценности – оптимизация происходит без учета стратегии, алгоритм следующий: инициализируем ценность состояний случайными числами, например $V(s)=0$ и перебором смотрим все возможные значения Q-функции (функции ценности пары состояние/действие) с учетом текущего вознаграждения за действие $a \in A$ и дисконтированной ценности следующего состояния $Q(S, A)=R_{ss}^A + \gamma * V(\dot{S})$ при этом ценность текущего состояния определяем как максимальное из возможных результатов Q-функции $V(s)=\operatorname{argmax}_a Q(S, A)$. Поскольку сходимость возможна только в пределе, но стратегии становится оптимальной раньше, вводим ϵ – маленькое положительное число, и при изменении ценности всех состояний на значение меньше ϵ останавливаем итерации. При этом стратегия максимизирующая $V(s)$ будет являться оптимальной.