

DQN with PER on MsPacman-v0

Fyodor Karpeyev

April 2022

Содержание

1	Аннотация	3
2	Введение	3
3	Обучение с подкреплением	4
4	Реплеи	5
5	Эксперимент	7
5.1	Реплей с приоритетом	7
5.2	Использование горячего старта	8
5.3	Методы сжатия входных данных	9
6	Выводы и дальнейшая работа	11
7	Заключение	11

Ключевые слова

обучение с подкреплением, реплей, реплей с приоритетом, DQN, Пакман, MsPacman-v0

1 Аннотация

Обучение с подкреплением сейчас является одной из важнейших составляющих машинного обучения. Данная работа представляет из себя исследование обучения агента в среде MsPacman-v0 с помощью DQN.

2 Введение

Обучение с подкреплением [1] сейчас является важной составляющей машинного обучения. Оно широко используется для решения огромного количества различных сред и может комбинироваться с широким набором дополнительных методов для улучшения показателей скорости и/или стабильности обучения. Существует множество модификаций DQN, но среди них существует большая группа моделей использующая как один из методов реплей.

Данная работа предлагает рассмотреть процесс улучшения модели DQN с использованием реплеев в среде Atari MsPacman-v0. Исследуется влияние разных методов и начальных значений на итоговое качество и скорость обучения модели.

3 Обучение с подкреплением

В модели мы используем обучение с подкреплением основывающиеся на Марковском процессе принятия решений [2]. Для автоматического обучения оптимальной стратегии, модель совершает действия и обновляется используя новые полученные данные о среде, в процессе максимизируется функция награды полученная от среды.

Мы формально определяем Марковский процесс принятия решений как набор из 4 элементов (S, A, P_a, R_a) , где S - это набор состояний среды, A - набор действий в данной среде, которые могут быть выполнены отталкиваясь от текущего состояния, $P_a(s, s') = P(s_{t+1} = s' | s_t = s, a_t = a)$ - вероятность того что действия a в состоянии s в момент времени t приведет к состоянию s' в момент времени $t+1$, $R_a(s, s')$ - вознаграждение, получаемое после перехода из состояния s в состояние s' с помощью действия a .

В процессе обучения модель формирует собственную «политику», другими словами стратегию. Выбор агента можно преставить как:

$$\begin{aligned}\pi : S \times A &\rightarrow [0, 1] \\ \pi(a|s) &= P(a_t = a | s_t = s)\end{aligned}$$

где π - вероятностная политика выбора действий, S - состояния, A - действия.

Нейронные сети используются в обучении с подкреплением в качестве аппроксиматоров, например в качестве аппроксиматора Q функции, определяющей качество выполнения определённого действия в определённом состоянии:

$$Q : S \times A \rightarrow \mathbb{R}$$

DQN(Deep Q-network [3]) основывается на Q-learning [4]. Для каждой пары состояние и действие в момент времени насчитывается опеределенное число $Q(s_t, a_t)$, которое изменяется в ходе обучения по формуле:

$$Q^{new}(s_t, a_t) \leftarrow \underbrace{Q(s_t, a_t)}_{old\ value} + \alpha \cdot \delta$$

Где α это гиперпараметр отвечающий за то, насколько сильно изменяются значения Q -функция за 1 шаг. А δ или TD (Temporal Difference) - это

ошибка временной разницы которую и пытается минимизировать модель, и которая вычисляется по формуле:

$$\underbrace{\underbrace{r_t}_{\text{reward}} + \underbrace{\gamma}_{\text{discount factor}} \cdot \underbrace{\max_a Q(s_{t+1}, a)}_{\text{estimate of optimal future value}} - \underbrace{Q(s_t, a_t)}_{\text{old value}}}_{\text{new value(temporal difference target)}}$$

Где γ - это параметр, который отвечает за то, насколько сильно обесценивается долгосрочная награда.

Существует модификация модели DQN, которая называется DDQN-Double Deep Q-Network. Она отличается наличием второй сети, а именно таргет сети, которая является той же сетью, но отдельно скопированной с основной, какое-то время назад. Модификация отличается от основной функцией ошибки, а именно имеет функцию:

$$\delta = r_t + \gamma Q_{\text{target}}(s_{t+1}, \operatorname{argmax}_{a \in \mathcal{A}} Q(s_{t+1}, a)) - Q(s_t, a_t)$$

4 Реплеи

Для улучшения качества обучения используются разные методы и эвристики, одна из таких это реплеи. Они запоминают и используют старые результаты игр для стабилизации и ускорения обучения. Формально реплей - набор записей о состояниях Марковского процесса принятия решений, которые наблюдает агент в процессе обучения, а именно: состояний, действий, полученных состояний, наград и дополнительных данных. Его размер ограничен, поэтому самая бесполезная информация каким-то образом удаляется для записи новых данных.

Одним из улучшений базового реплея, основанном на том, что удаляются самые старые записи, является реплей с приоритетом [5]. В нем объекты на удаление выбираются согласно приоритетам $P(i) = \frac{p_i^\alpha}{\sum_k p_k^\alpha}$, где α - гиперпараметр приоритезации. Таким образом он сохраняет только самые полезные записи.

Одним из вариантов реплея с приоритетом является реплей в котором $p_i = |\delta_i| + \epsilon$, где ϵ это маленькая константа, гарантирующая, то что все элементы реплея имеют не нулевой шансы быть выбранными.

Для стабилизации реплея с приоритетом используется выборка с важностью. Для каждого объекта в реплее мы добавим вес:

$$w_i = \left(\frac{1}{N} \cdot \frac{1}{P(i)} \right)^\beta$$

Где N - это количество записей в реплее в данный момент, а β - параметр, отвечающий за то, насколько сильно влияет шанс на взятие объекта из реплея на его вес.

5 Эксперимент

Все эксперименты поставлены с использованием DDQN в среде MsPacman-v0.

5.1 Реплей с приоритетом

Рассмотрим влияние реплея с приоритетом(и выборки с важностью) на качество обучения.

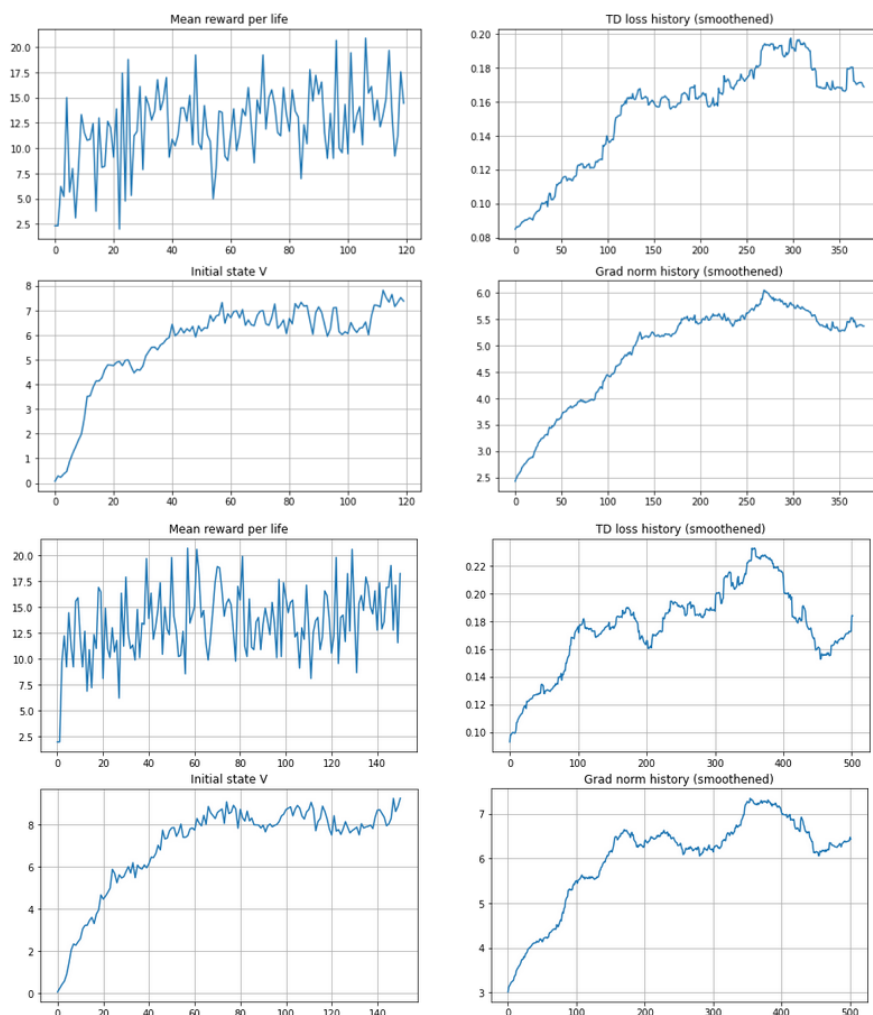


Рис. 1: Сравнение графиков обучения. Сверху с обычным реплеем. Снизу с реплеем с приоритетом.

Итоговые результаты по этим графикам это соответственно 450 и 700. И хотя такое большое различие в результатах скорее всего является следствием случайности, на нижнем графике средней награды видно, что реплей с приоритетом позволяет быстрее достигать хороших значений и стабильнее их сохранять.

5.2 Использование горячего старта

С помощью горячего старта от одного из лучших агентов обученных за все время экспериментов с результатом 750.

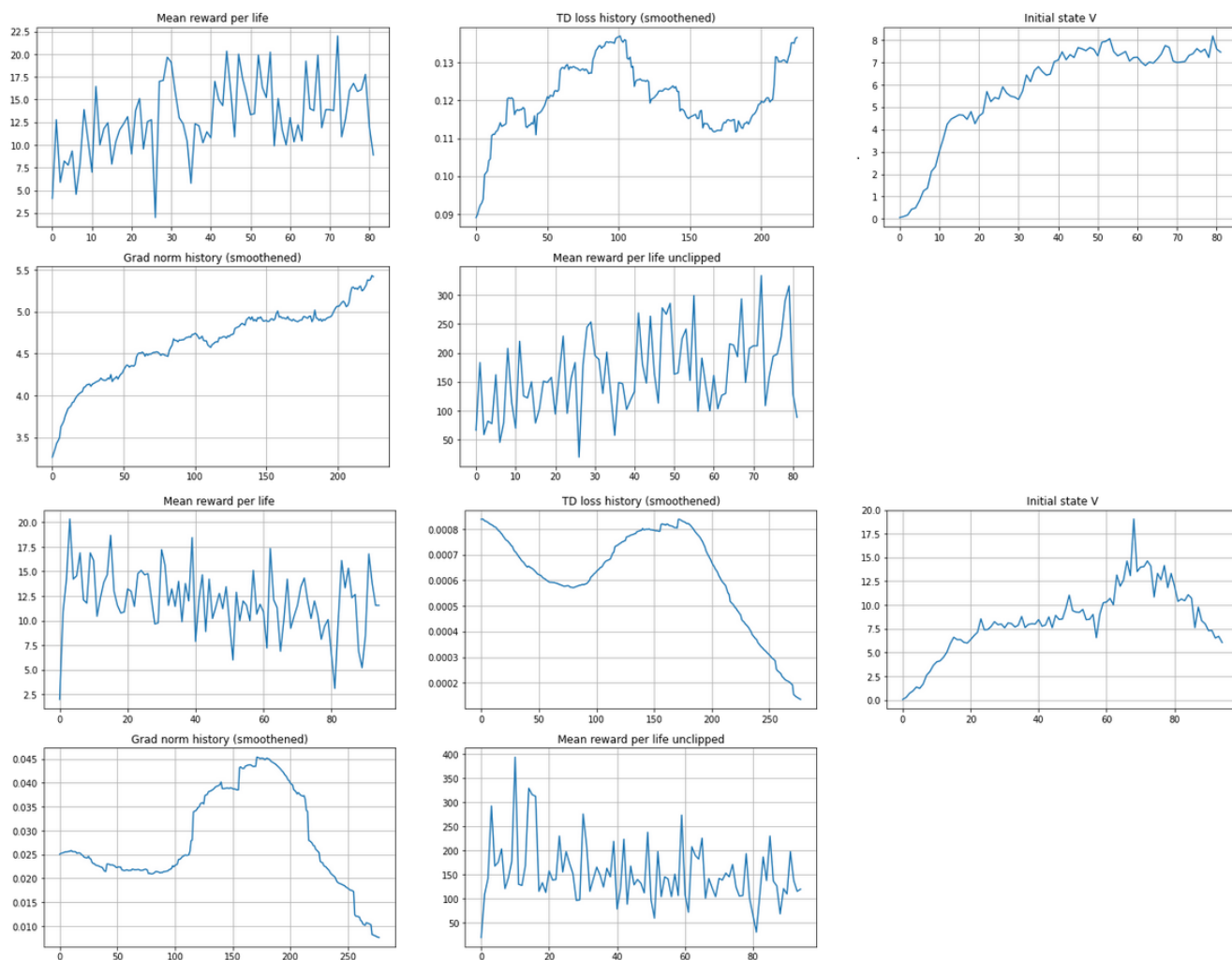


Рис. 2: Сравнение графиков обучения. Сверху без. Снизу с горячим стартом.

Как видно по графикам горячий старт помогает быстро достичь хорошей средней награды, но после какого-то времени начинает падать вниз, а графики не награды показывают просто ужастные результаты в отличие от отсутствия горячего старта.

Вот еще 1 график с использованием горячего старта уже для другой конфигурации модели и другого агента для заполнения:

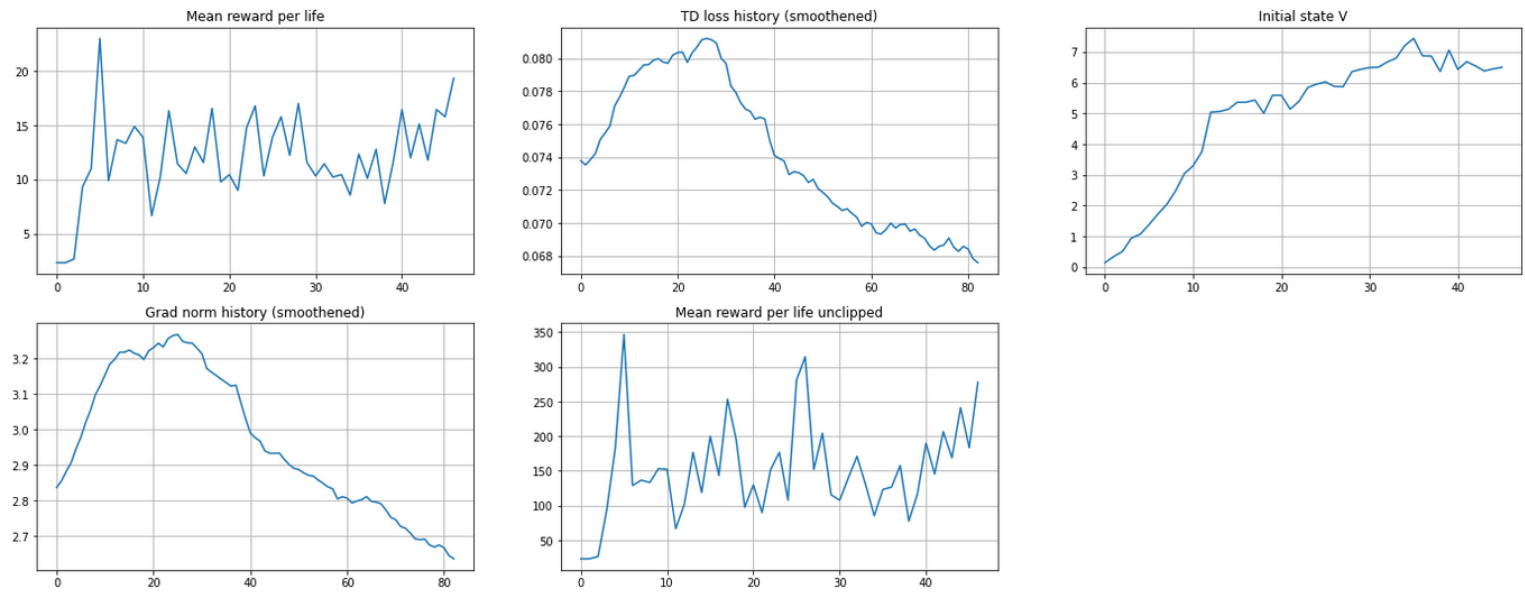


Рис. 3: График обучения с горячим стартом.

Видны те же самые проблемы, и хотя тут график Initial State V еще не упал, виден его отрицательный тренд.

5.3 Методы сжатия входных данных

Изначально картинка в Пакамане имеет вид: особа сжатия цветов картинки в оттенки серого, и в либо черный либо белый цвет:



И помимо обрезания и сжатия, которые банально ускоряют процесс, за-
счет меньшего числа данных, и использования бафера из 4 последовательных
кадров, для понимания направления и скорости движения. Мы рассмотрим
2 способа сжатия цветов картинки в оттенки серого, и в либо черный либо
белый цвет:

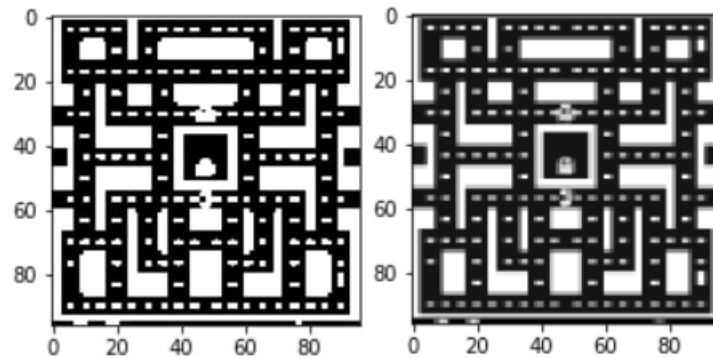


Рис. 4: Версия просле обработки слева только черный или белый, справа оттенки серого.

На левой картинке не хватает некоторых деталей, и это подтверждается
обучением. Не понятным мне образом приведение картинки в бинарный вид
пикселей вызывает взрыв коэффициентов даже при рабочем нормализаторе:

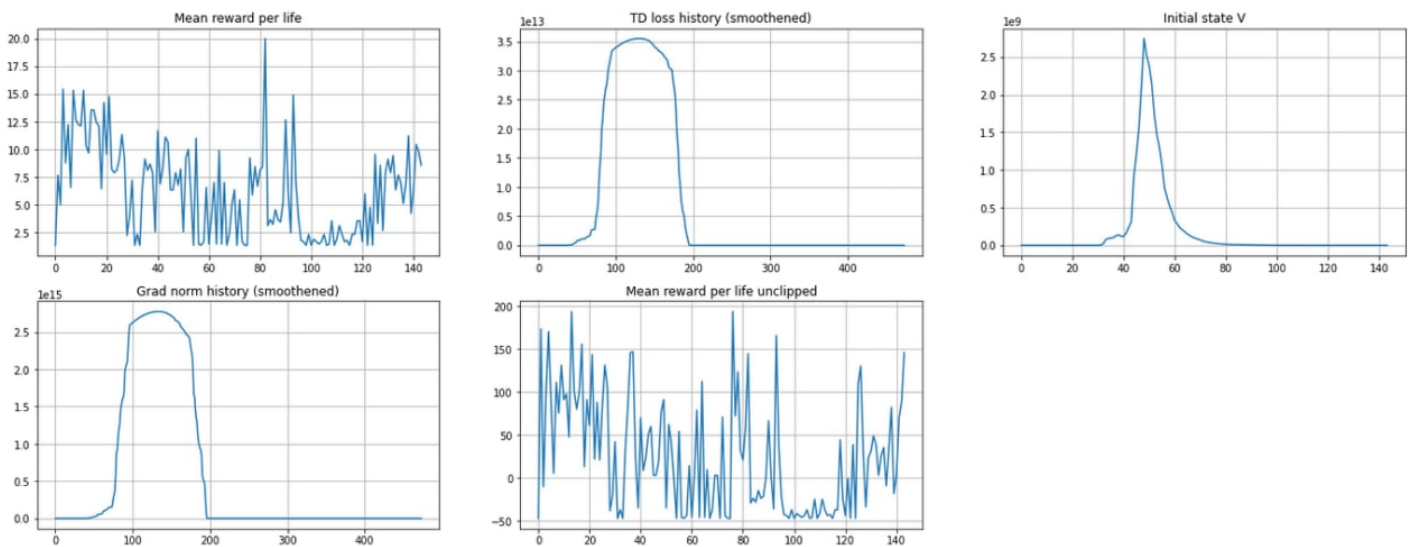


Рис. 5: График обучения со взрывом коэффициентов.

Несмотря на взрыв коэффициентов, агент обученный таким образом пока-
зывает не лучший, но неплохой результат.

Можно посмотреть на еще один подобный пример, с немного другой внут-
ренней архитектурой и другими параметрами сети.

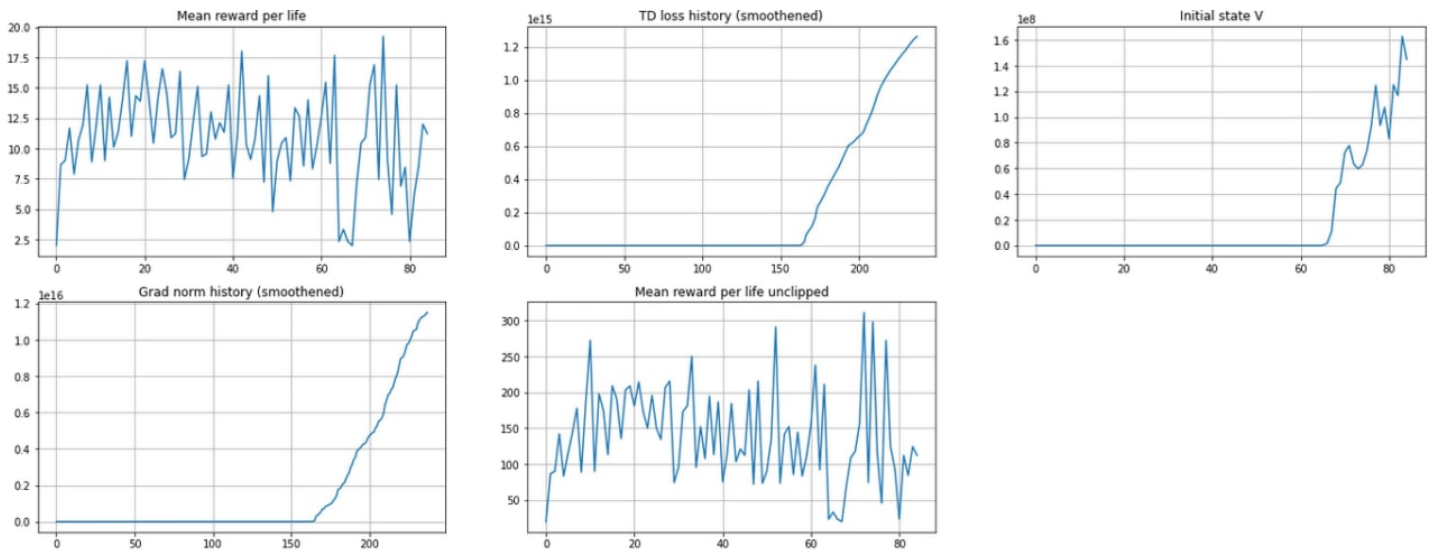


Рис. 6: Еще один график обучения со взрывом коэффициентов.

6 Выводы и дальнейшая работа

Лучший результат у агента, которого мне удалось достичь это 750 очков, что является достаточно неплохим результатом для достаточно простой модели без использования разных трюков. Делая выводы из проделанной работы реплей с приоритетом улучшает качество обучения, но с другими методами надо быть осторожнее. Любой метод, который в теории должен улучшать качество, может наоборот ухудшить его или вызвать непредвиденные последствия. С точки зрения DQN расширения не в глубь, а в ширь чаще всего не дает сильно хороших результатов, так например увеличение реплея, или размера выборки из реплея, или учащение обновления таргет сети, или расширение сети на несколько слоев дает почти незаметное улучшение качества, но при этом модель обучается намного дольше. Исходя из этого, можно сказать что для простой модели типа DDQN достойным способом улучшения является добавление дополнительных методов и эвристик.

Дальнейшая работа состоит в том, чтобы изучить разнообразные методы и эвристики и рассмотреть их влияние на качество обучения этой модели и других, а так же провести подобные эксперименты в других средах.

7 Заключение

В данной работе рассматривается исследование поведения модели DDQN на среде MsPacman-v0. Было достигнуто хорошее относительно возможно-

стей итоговое качество, а полученные знания и методы можно применить для других моделей и сред.

Список литературы

- [1] Vincent François-Lavet и др. «An Introduction to Deep Reinforcement Learning». Англ. В: *CoRR* abs/1811.12560 (2018). arXiv: 1811.12560. URL: <http://arxiv.org/abs/1811.12560>.
- [2] «Markov decision process». В: (). URL: https://en.wikipedia.org/wiki/Markov_decision_process.
- [3] Volodymyr Mnih и др. «Playing Atari with Deep Reinforcement Learning». Англ. В: *CoRR* abs/1312.5602 (2013). arXiv: 1312.5602. URL: <http://arxiv.org/abs/1312.5602>.
- [4] «Q-learning». В: (). URL: <https://en.wikipedia.org/wiki/Q-learning>.
- [5] Tom Schaul и др. «Prioritized Experience Replay». Англ. В: *ICLR*. 2016.