

# DRL Course Домашнее задание 1

## Алексей Матушевский

В этой практической работы мы изучаем работу алгоритма Cross Entropy Method (CEM) совмещенного с нейронными сетями для поиска оптимальных параметров в среде с непрерывным пространством с, не обязательными, непрерывными действиями. Непрерывные значения для действий и состояний не позволяют использовать предыдущие подходы и их улучшения, так как предыдущие подходы опирались на матрицы стратегий. Каждая строка такой матрицы представляло собой конкретное состояние системы, а столбец - действия. Непрерывные значения — не позволяют работать с матрицами, или попросту создают огромной матрице.

В таких случаях лучше альтернатива - нейронные сети. Они подходят для работы с непрерывными значениями.

Работа в среде Lunar Lander.

В этой среде мы управляем посадочным модулем 4мя дискретными действиями. Среда описывается 8мя непрерывными состояниями.

Для решения этой задачи я использовал CEM +NN с разными наборами параметров

Episode Len — количество эпизодов тренировки

Q-param — квантиль отбираемые

Trajectory N — количество траекторий в одном эпизоде

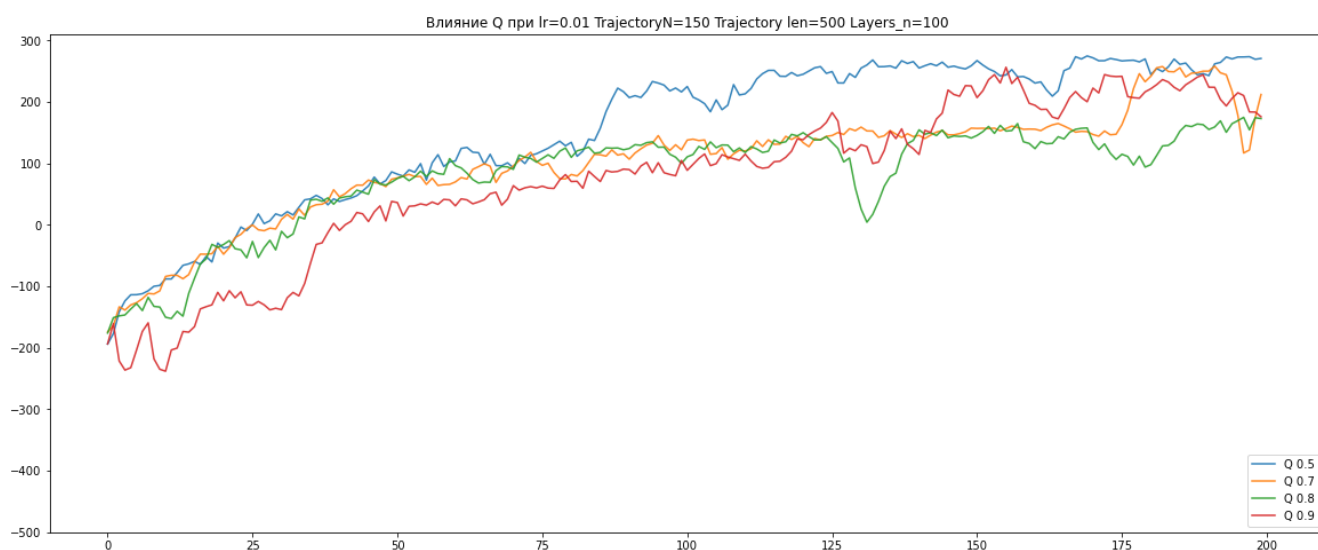
Trajectory Len — длина одной траектории

Layers N — количество нейронов в каждом слое

Learning Rate — темп изменения параметров оптимизатором

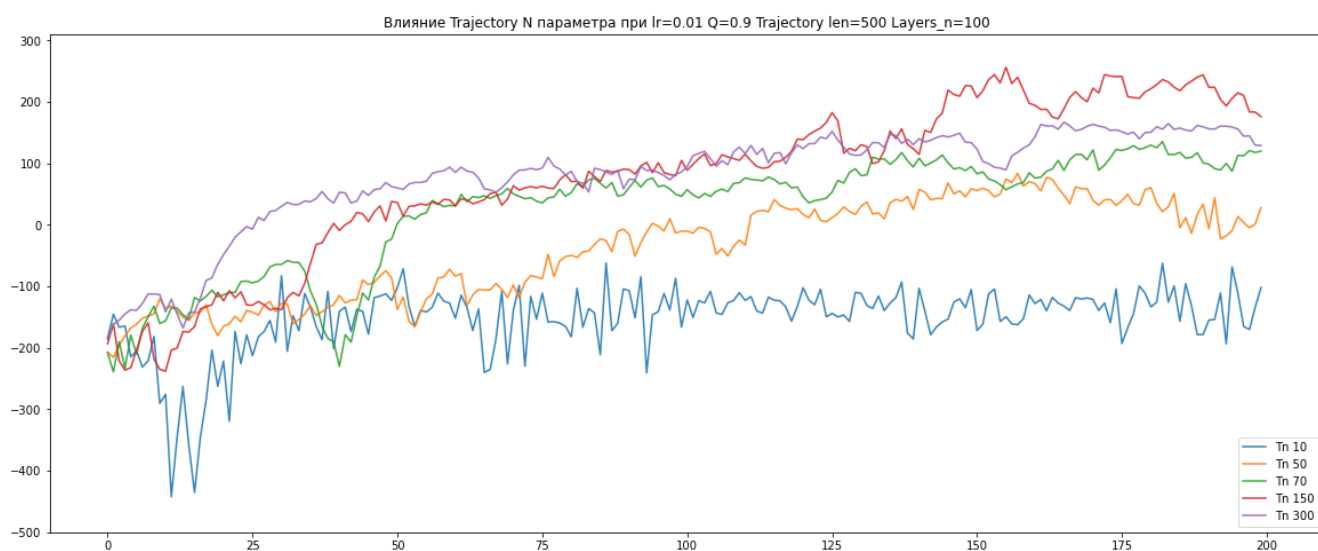
### **Layers N при разных Q**

Проанализируем влияния Q на достижения большой награды. Как видно с увеличением Q достижение высокой награды занимает большее количество эпизодов тренировки.

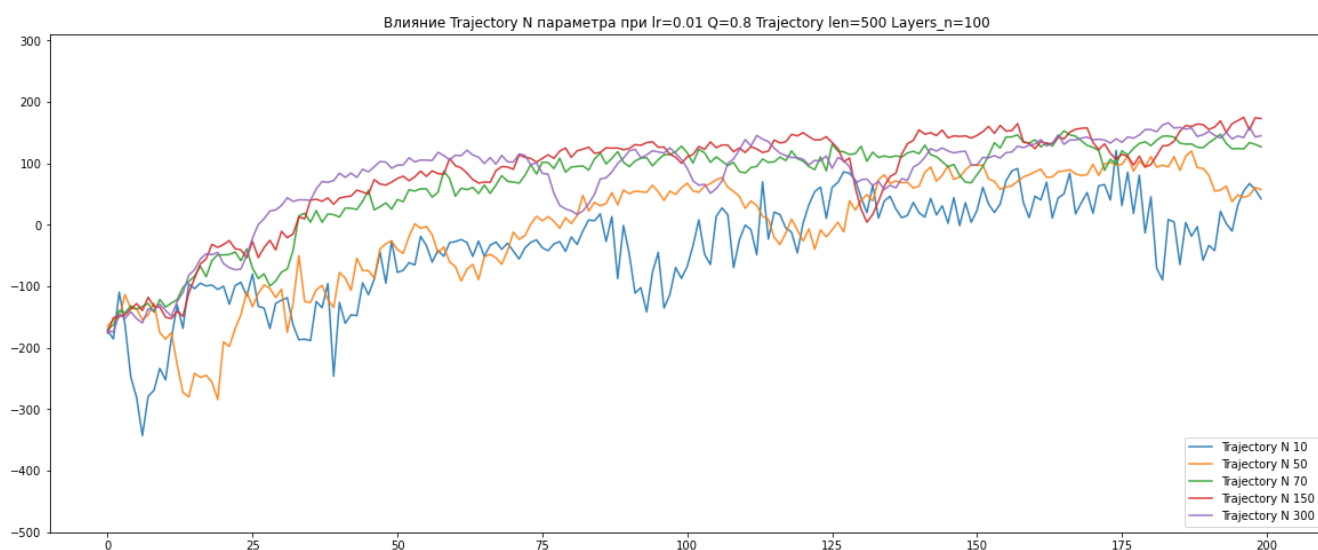


## Trajectory N

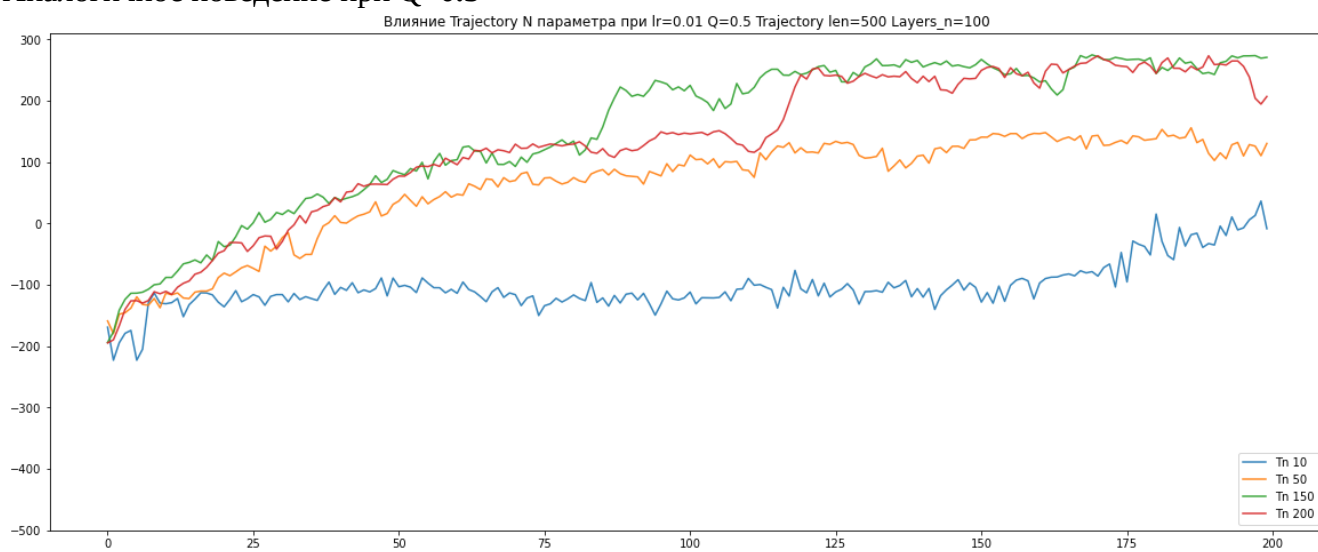
Увеличение Trajectory N при  $q=0.9$  Trajectory len 500. с увеличением - уменьшатся дисперсия в росте оценки.



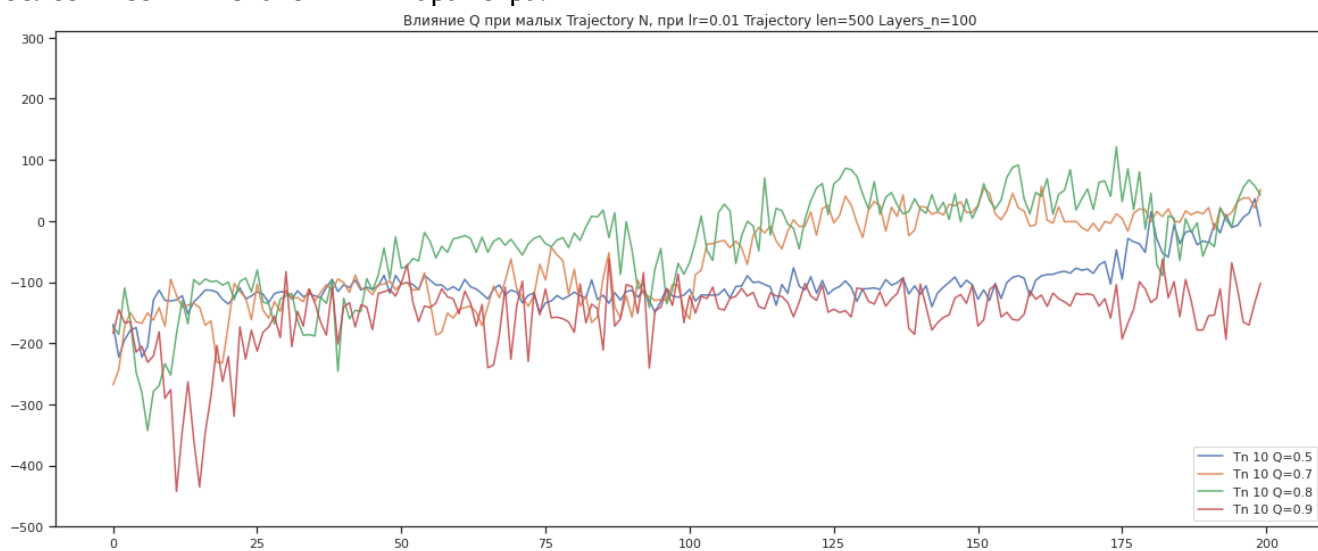
Увеличение Trajectory N при  $q=0.8$  Trajectory len 500. Так-же с увеличением - уменьшатся дисперсия в росте награды.



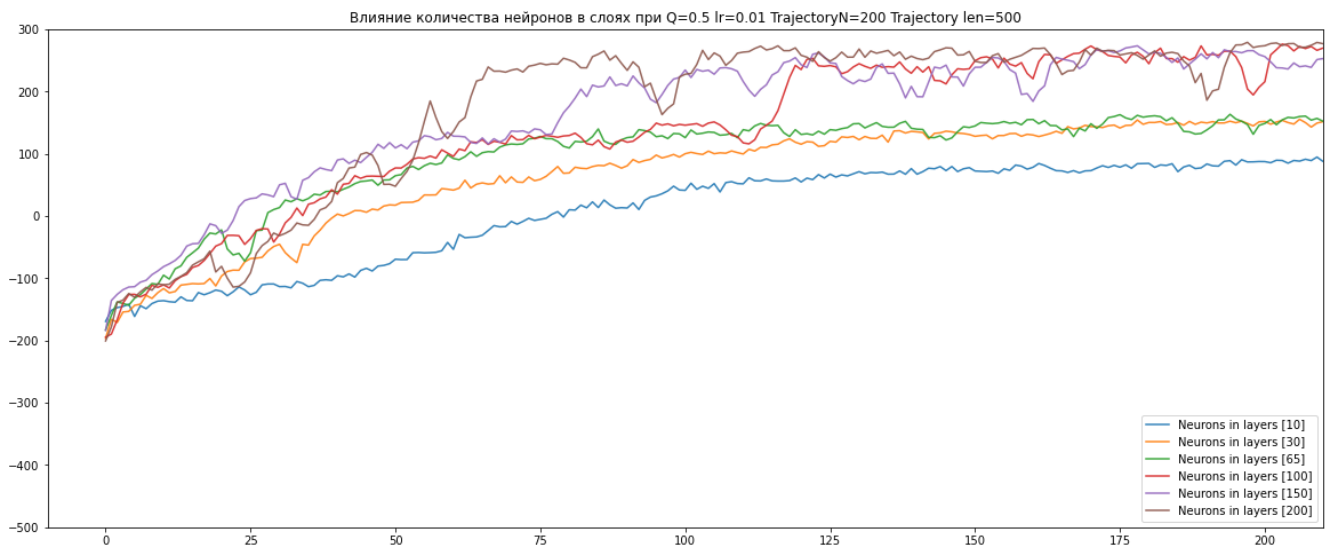
## Аналогичное поведение при $Q=0.5$



При очень низких Trajectory N алгоритм не приходит к стабильному состоянию в сравнении с более высокими значениями параметра.



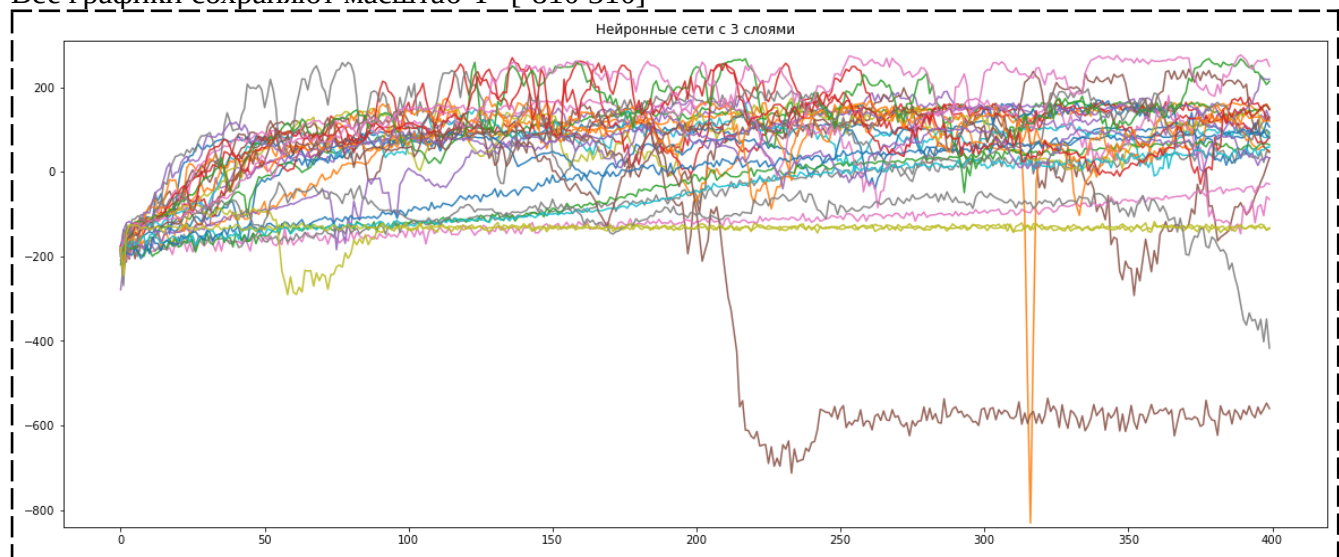
## Влияние количества нейронов Layers N



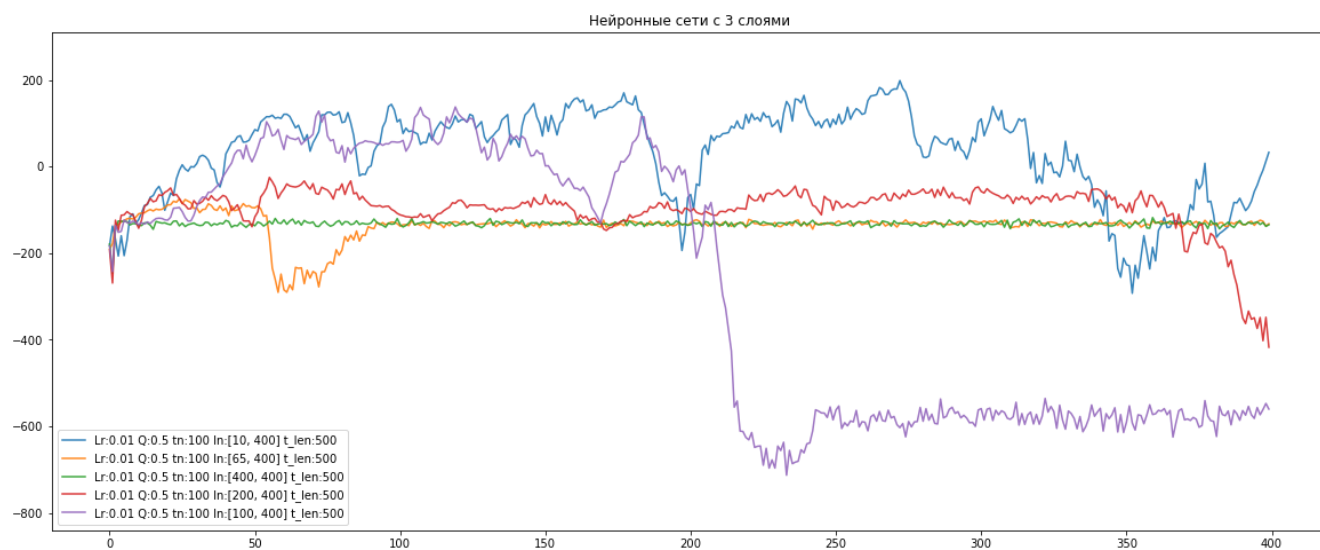
При одних и тех же параметрах  $Q=0.5$ , Trajectory N=200, и Trajectory Len=500 — создали сети с разным количеством нейронов. Как видно из графика с увеличением количества нейронов алгоритм быстрее достигает высокой награды, но при этом увеличивается дисперсия.

## Увеличение количества слоев в сети.

Оставляем те же параметры  $Q=0.5$ , Trajectory N = 100, Trajectory Len=500 и **400** эпизодов  
Все графики сохраняют масштаб  $Y=[-810\ 310]$

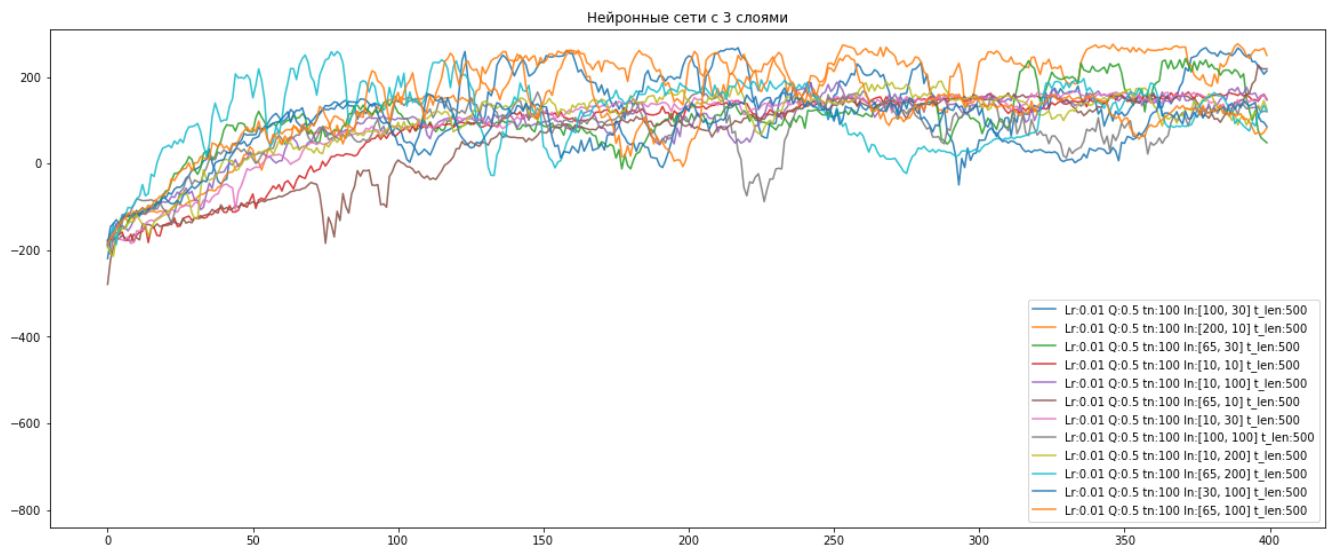


При очень большом количестве нейронов, сеть не может удержаться на высокой награде или вообще не достигает положительной награды.

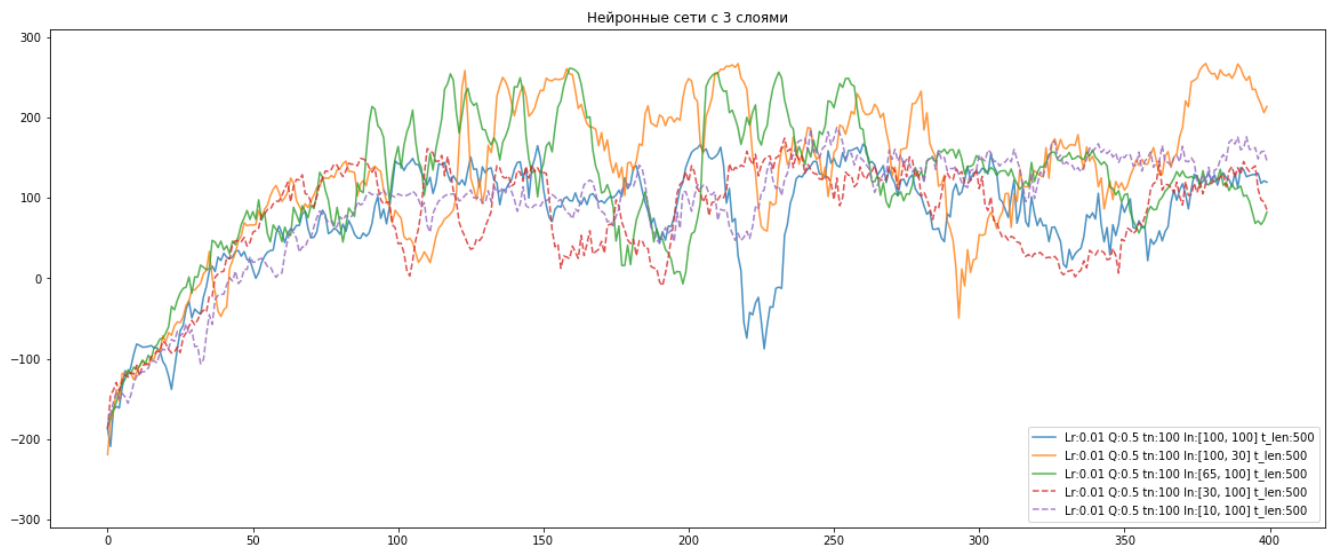
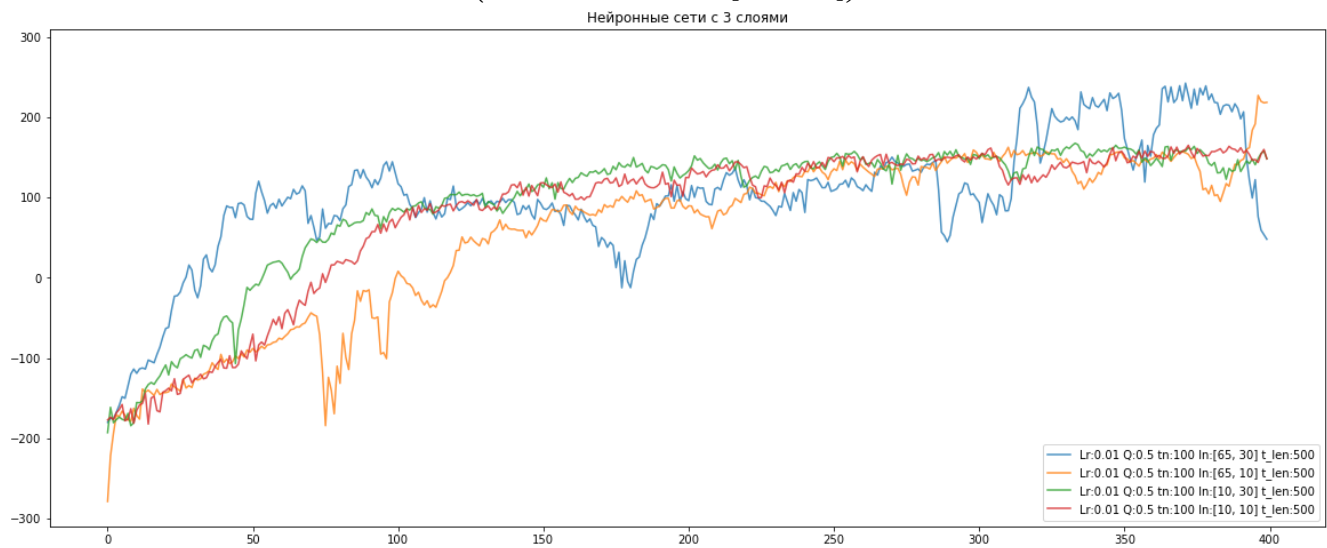


Сети которые удерживаются выше нуля, имеем слои не больше 200 нейронов в одном из слоев



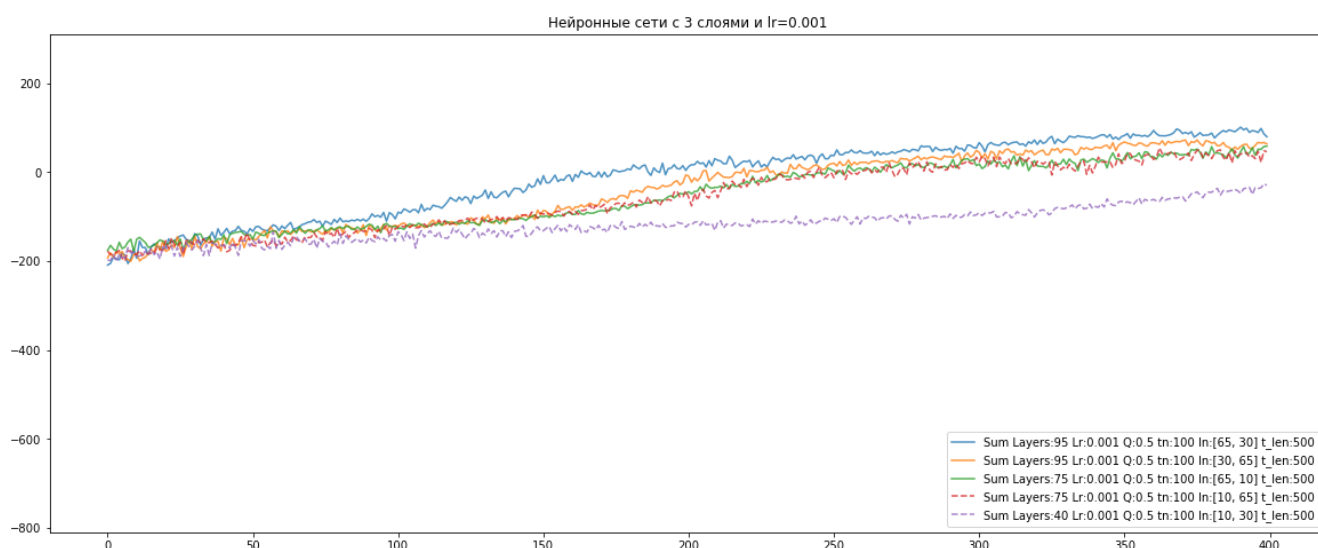


Сети с меньшим в сумме количеством нейронов дают меньшую дисперсию в росте награды  
(масштаб изменен на [-310 300])





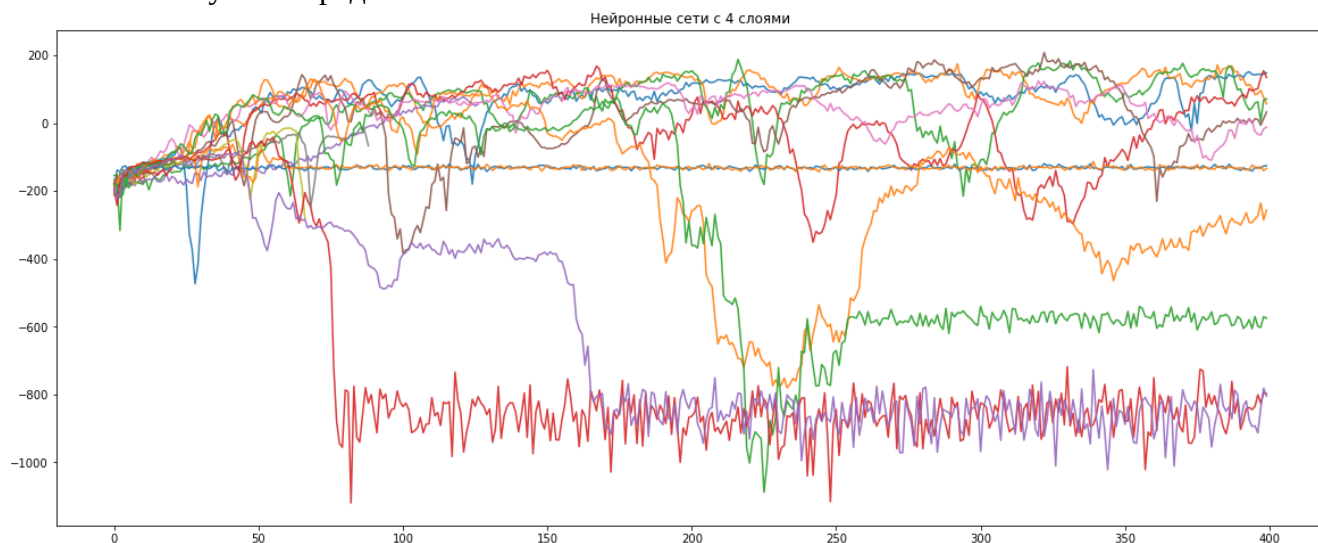
при изменении Learning Rate — скорость схождения сети уменьшается пропорционально изменению  $lr$ , сеть с меньшим количеством нейронов дает наименьшую скорость схождения  
(масштаб изменен на [-810 300])



## Четырехслойные сети

При увеличении количества слоев до 4 — общее поведение остается тем же. Высокое количество нейронов — дает

- большой разброс в росте награды
- новый минимум в награде



## Learning Rate

Влияние Learning Rate на разные сети. Как и в случай двухслойных чем ниже LR тем дольше сходится алгоритм



### Лучшие результаты

Если выложить все параметры на 3D график,

X — Q-param

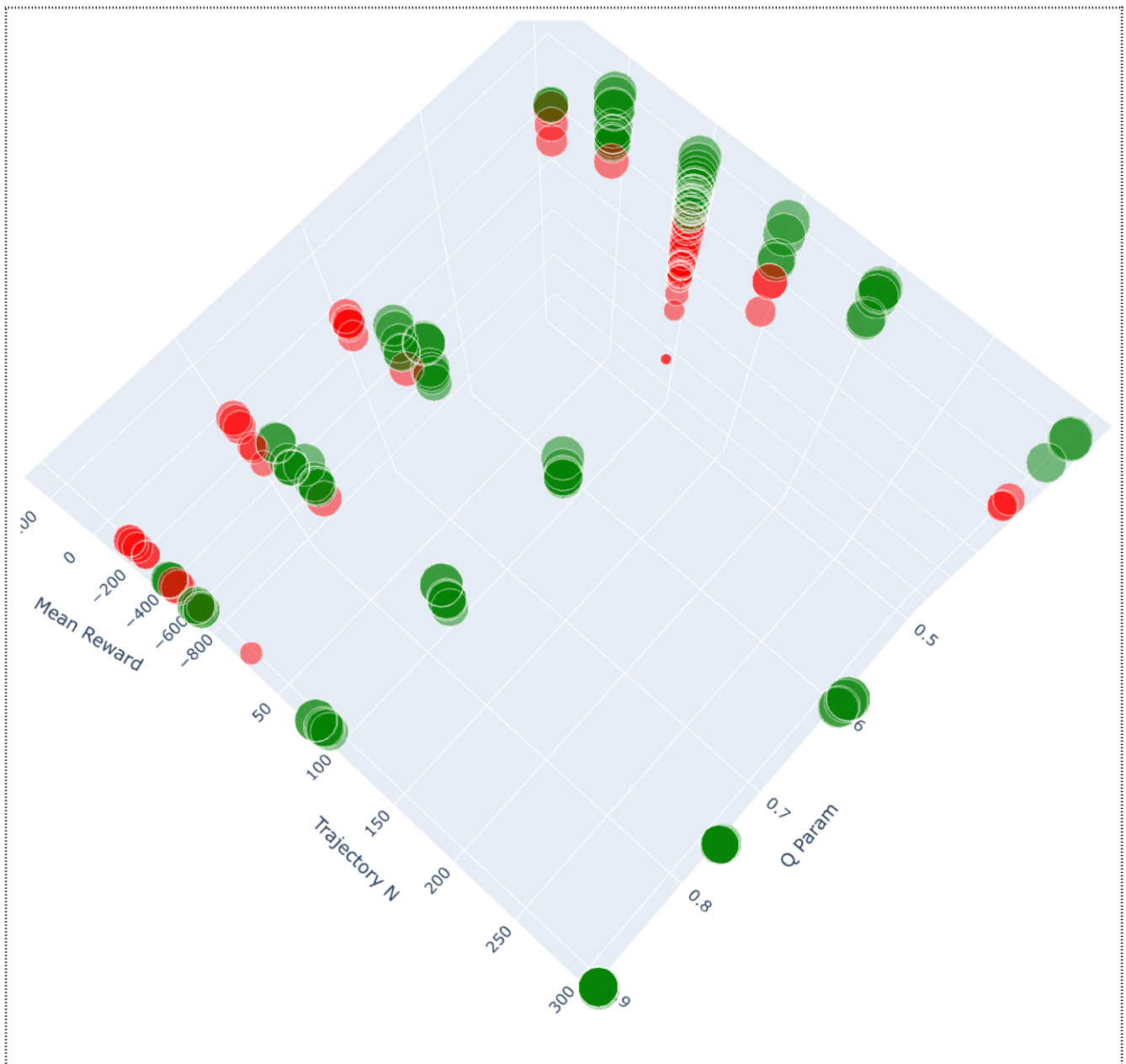
Y — Trajectory N

Z — Mean Reward

Можно увидеть что в обласит Q от 0.6-0.9 и Trajectory N [75-300] нет моделей со средними отрицательными значениями награды.

[http://alexeimatusovski.com/odsr/repot\\_2-1-2.html](http://alexeimatusovski.com/odsr/repot_2-1-2.html)

[http://alexeimatusovski.com/odsr/repot\\_2-1-1.html](http://alexeimatusovski.com/odsr/repot_2-1-1.html)



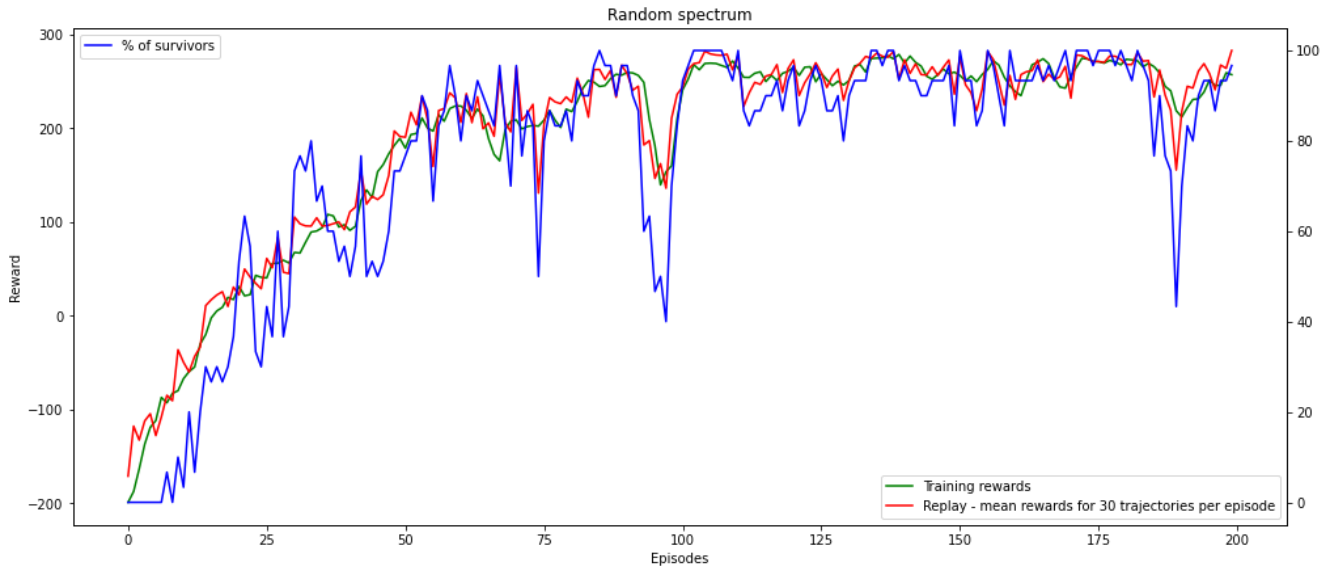
Эксперимент — это запуск алгоритма с определенными параметрами. Критерии выбора лучших результатов основываются на полученных данных из 190 экспериментов.

Для каждого эксперимента мы сохранили:

- параметра эксперимента
- состояние нейросети в каждом его эпизоде
- средние награды в течении эксперимента

Можно выбирать лучшие параметры по средней оценки во время тренировки. Графики тренировок могут содержать «провалы» в награде. Если повторить это эксперимент используя сохраненные состояния весов нейросети, можно получить выборки по средним наградам а также оценить количество «выживших». Выжившие — это те траектории которые не получали награду -100 за конкретный шаг.

Из графика видно что 100% выживаемости (из 30 запусков) достигается только у нейросетей с самыми высокими средними наградами.



Минимальная награда которая дает 100% выживших (из выборки с 30ю запусками) 256. Поэтому можно предположить что эпизоды дающие такие и выше средние награды, будут подходить для лучших результатов.

Reward	Survivors
262.951434	1
269.111621	1
270.265116	1
282.396529	1
279.517435	1
278.370167	1
277.963013	1
277.992915	1
274.999818	1
280.340668	1
275.9605	1
281.941145	1
276.330687	1
282.151603	1
256.508145	1
266.308203	1
278.3295	1
277.618083	1
274.333956	1
271.442862	1
270.33336	1
277.222227	1
273.670604	1
280.847168	1

Итого список топ 7 параметров:

	lr	episode_n	trajectory_n	trajectory_len	q_param	layers_cnt	layers_n1	layers_n2	layers_n3
0	0.01	200	300	500	0.7	1	200	-	-
1	0.01	200	150	500	0.5	1	100	-	-
2	0.01	200	300	500	0.5	1	150	-	-
3	0.01	200	300	500	0.7	1	100	-	-
4	0.01	200	300	1000	0.5	1	100	-	-
5	0.01	400	200	500	0.5	1	100	-	-
6	0.01	400	200	500	0.5	1	150	-	-

```

52.09920158153051
245.97549456956529
28.625970473042997
268.9199202073091
243.1997945499797
276.8018165858673
259.87678350851877
281.9496241211709
310.6052094386491
283.43576523139814
303.2496112156107
262.0132151504142
304.5785939413573
244.05831221649265
273.2977652434016
301.6569009112594
268.8604899021014
273.49488557821303
301.99129078204885
308.8571346939664
292.5819956039746
255.91951396232758
284.80069565927886
266.5035216231822
309.4003962283613
277.81980260532714
284.64525477144525
260.1543781597268
254.71124569247087
273.32057055755723
260.6688140404064
46.92878497920353
299.90095331946975
261.9069606214112
261.3784999829279
273.842968017716
233.75690924217233
285.5291693258896
282.5849450259742

```

