

# Домашняя работа #1 RL

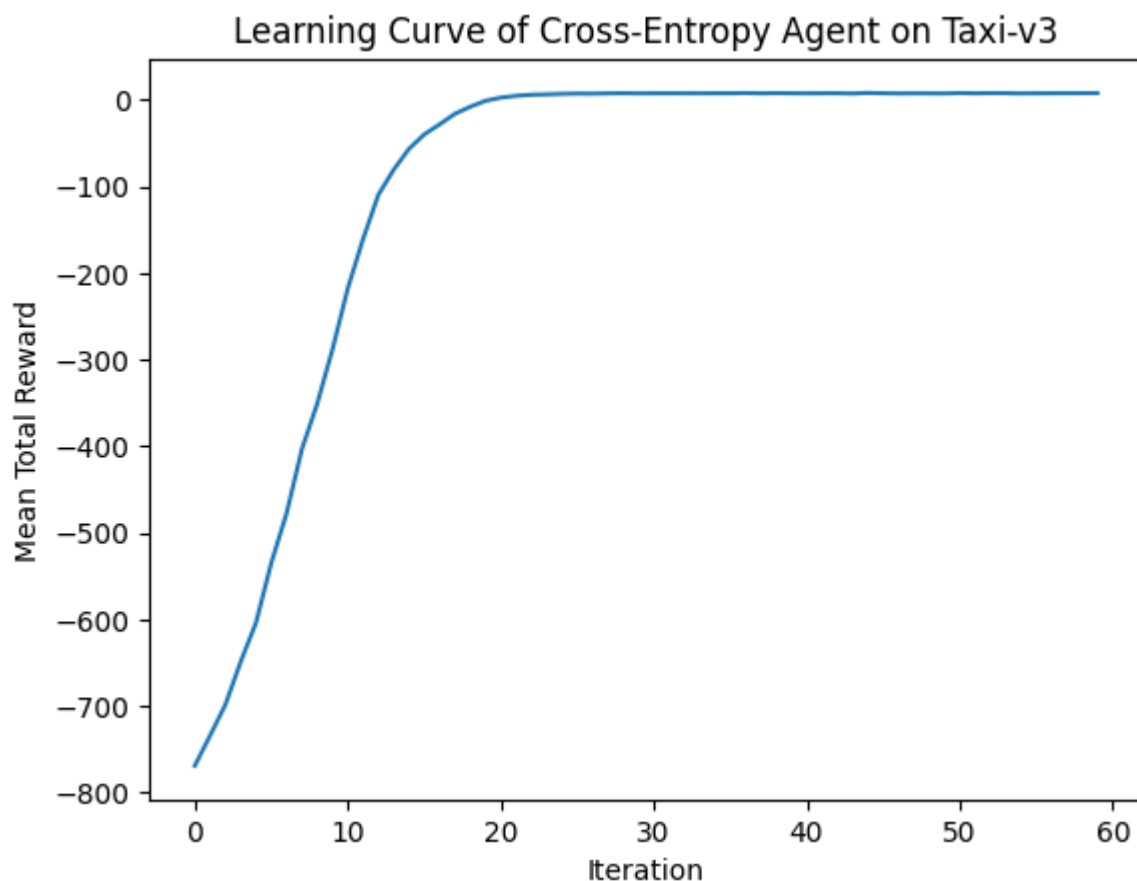
## Введение

## Первое задание

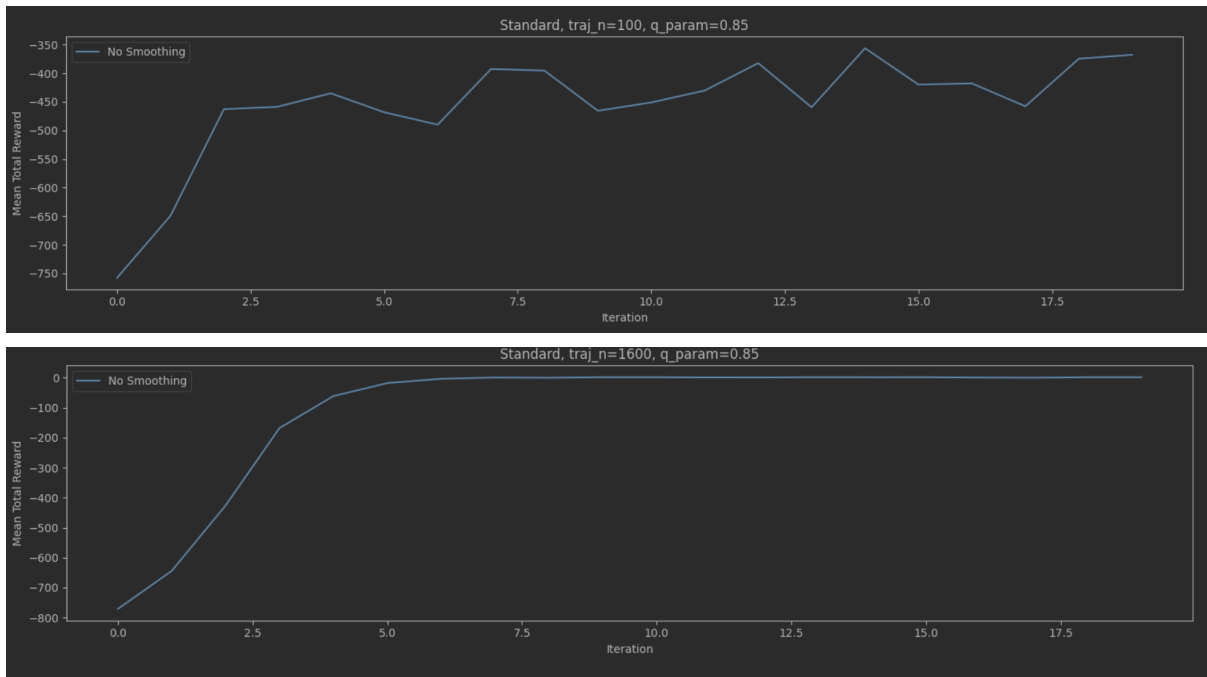
## Обучение

### Эксперимент 1

Для начала я просто запустил модель с квантилем 0.3, и получил очень неплохой результат. Скорее всего он был случайным и я решил поиграться.



Вот некоторые из графиков:



Примечательно то, что чем больше траекторий тем более сглаженной и стабильнее график. высокие параметры  $q$  проигрывали при меньшем количестве траекторий, но с ростом оных все уравнивалось(больше графиков во втором `ipynb` ноутбуке)

## Вывод

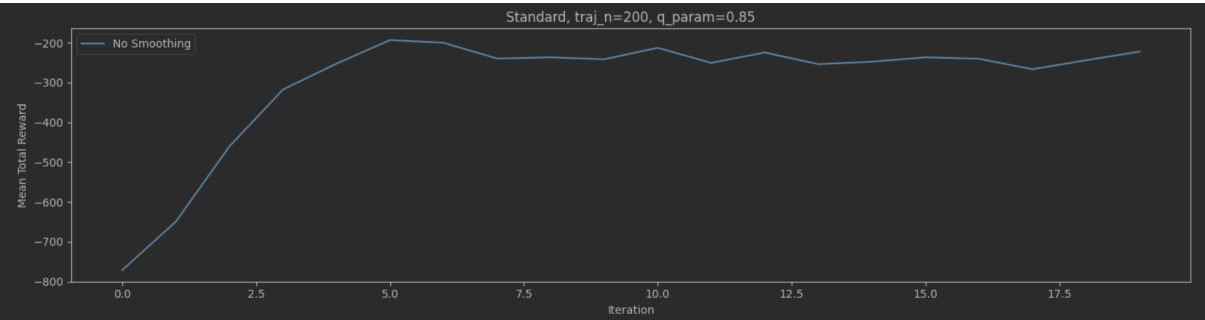
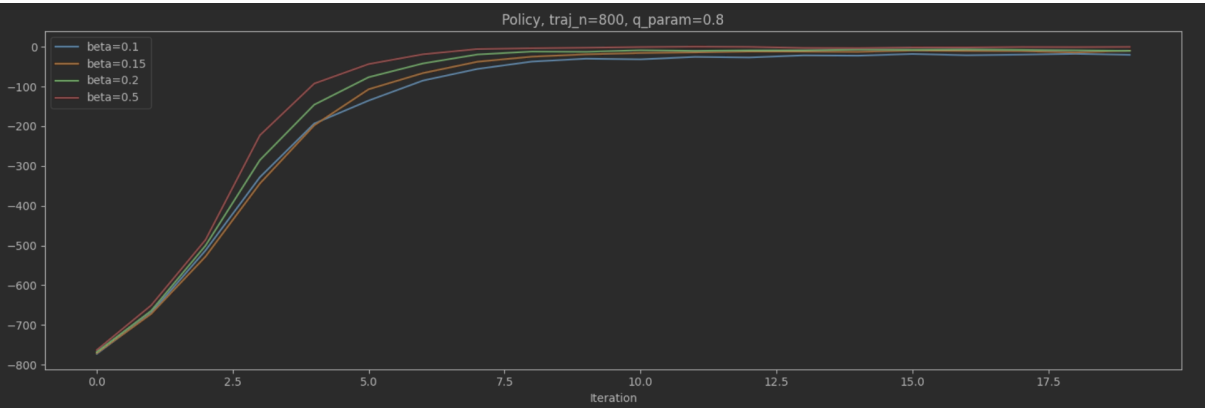
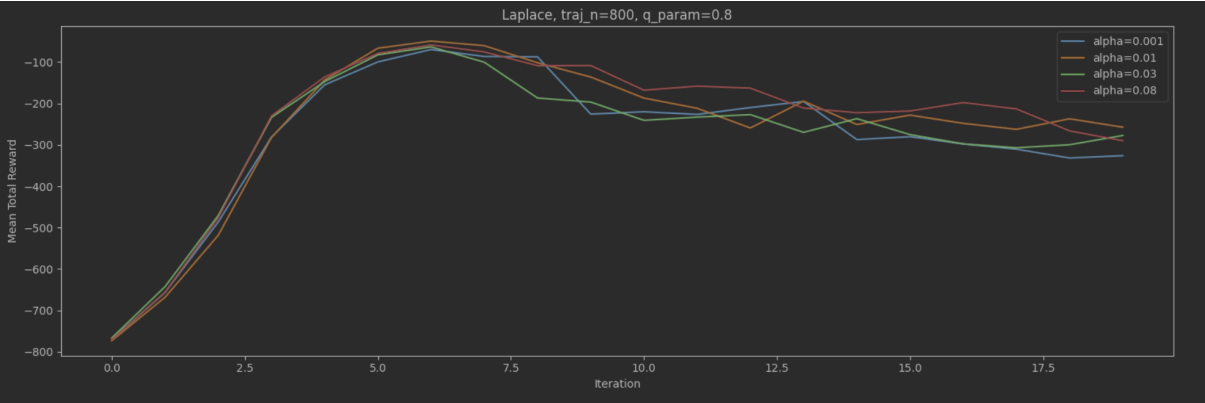
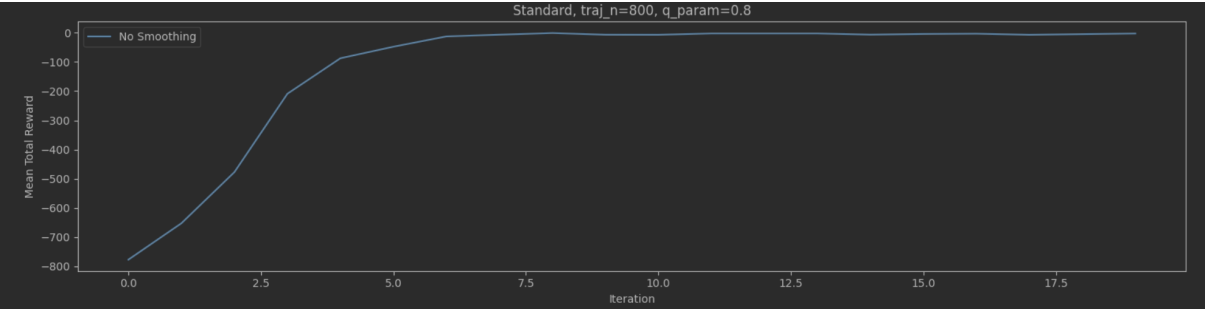
Так как среда, алгоритм, и проделанные действия достаточно простые, то все обучилось довольно быстро, и каких-то особых выводов за исключением уже названных нет.

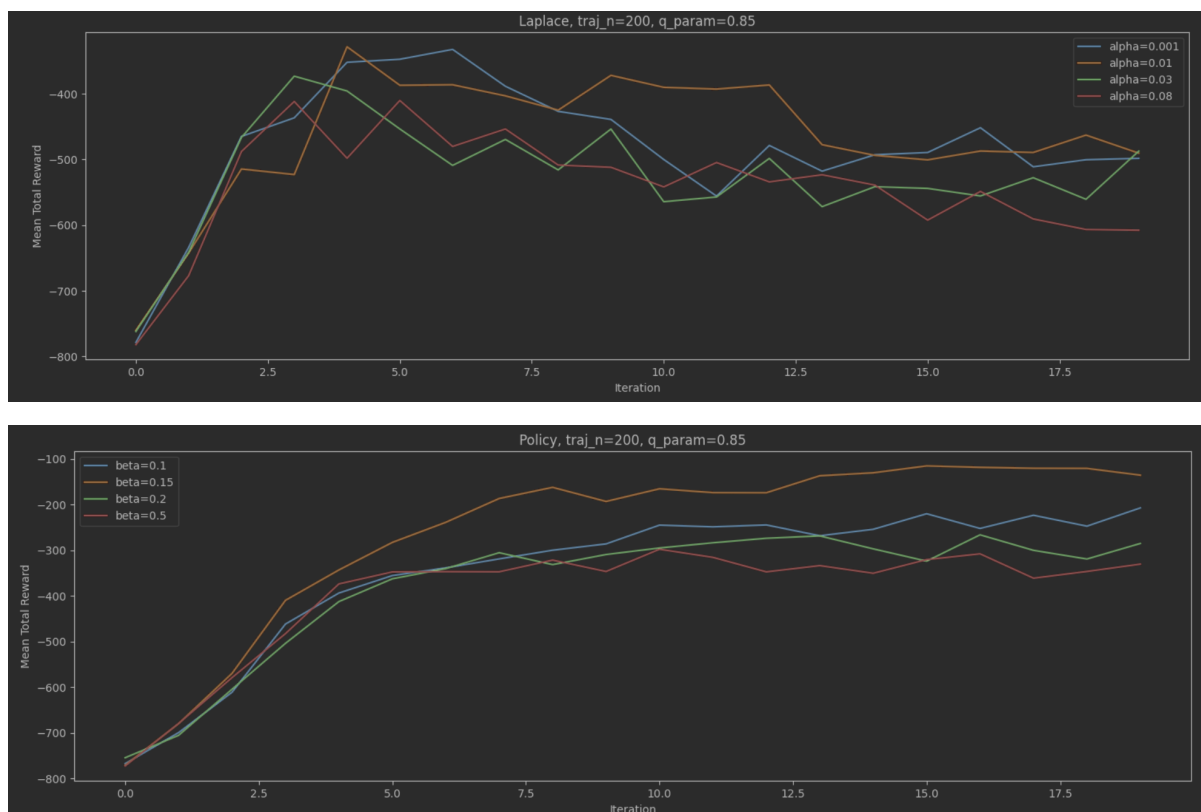
## Второе задание

### Обучение

#### Эксперимент 1

Во втором задании необходимо было применить два способа сглаживания и сравнить с обычным алгоритмом. Ниже графики(в ноутбуке больше).





Опять таки в iрunb ноутбуке значительно больше графиков. Что стало понятно. При больших количествах итераций обычный и policy способы сравниваются по точности. Бета почти не влияет. А вот если говорить о маленьких количествах траекторий то тут картина иная, и модели ведут себя по разному. Очень странно поведение Laplace. Он изначально добивается неплохой точности, но потом происходит падение награды. В целом исходя из этой лабораторной не совсем понятно насколько хороший эффект оказывают policy и laplace. Возможно с другой средой они бы повели себя по другому.

## Вывод

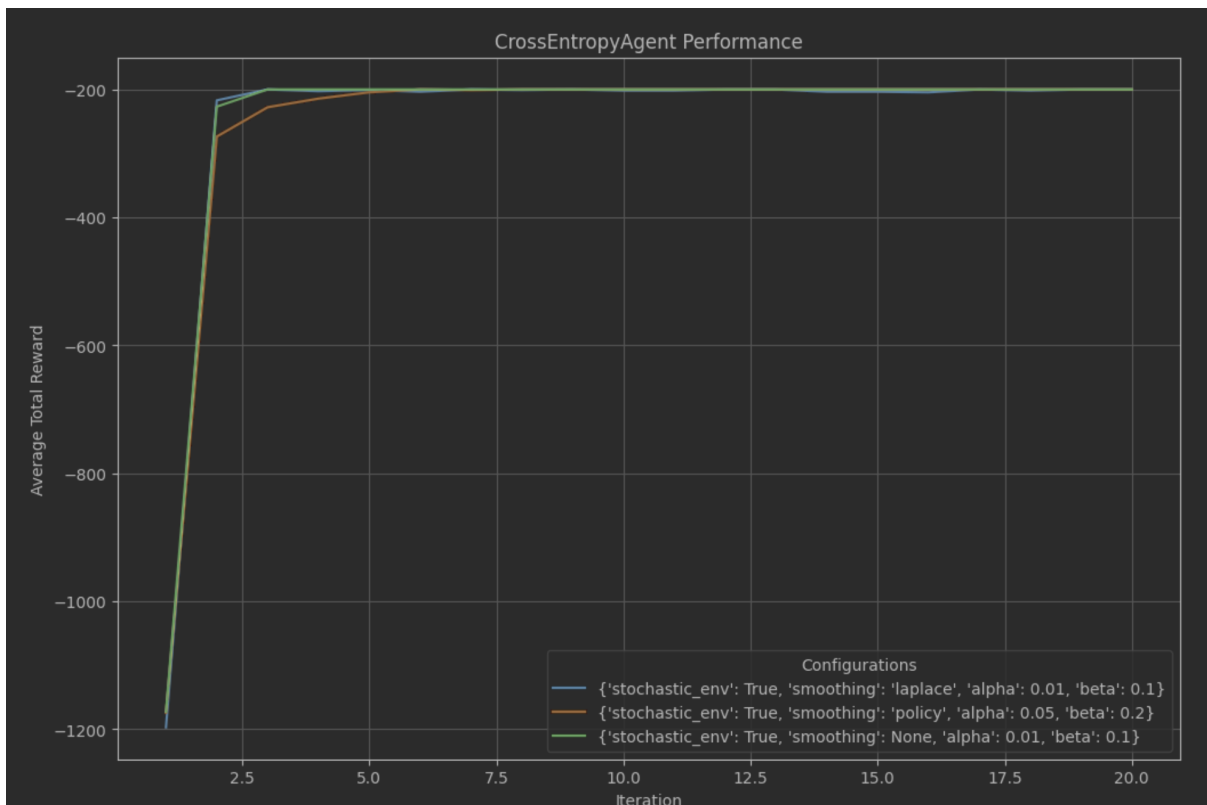
Исходя из этой лабораторной не совсем понятно как policy и laplace влияют. Единственное что видно это то, что laplace почему-то теряет точность. То ли попадает в локальный оптимум то ли еще что-то.

## Третье задание

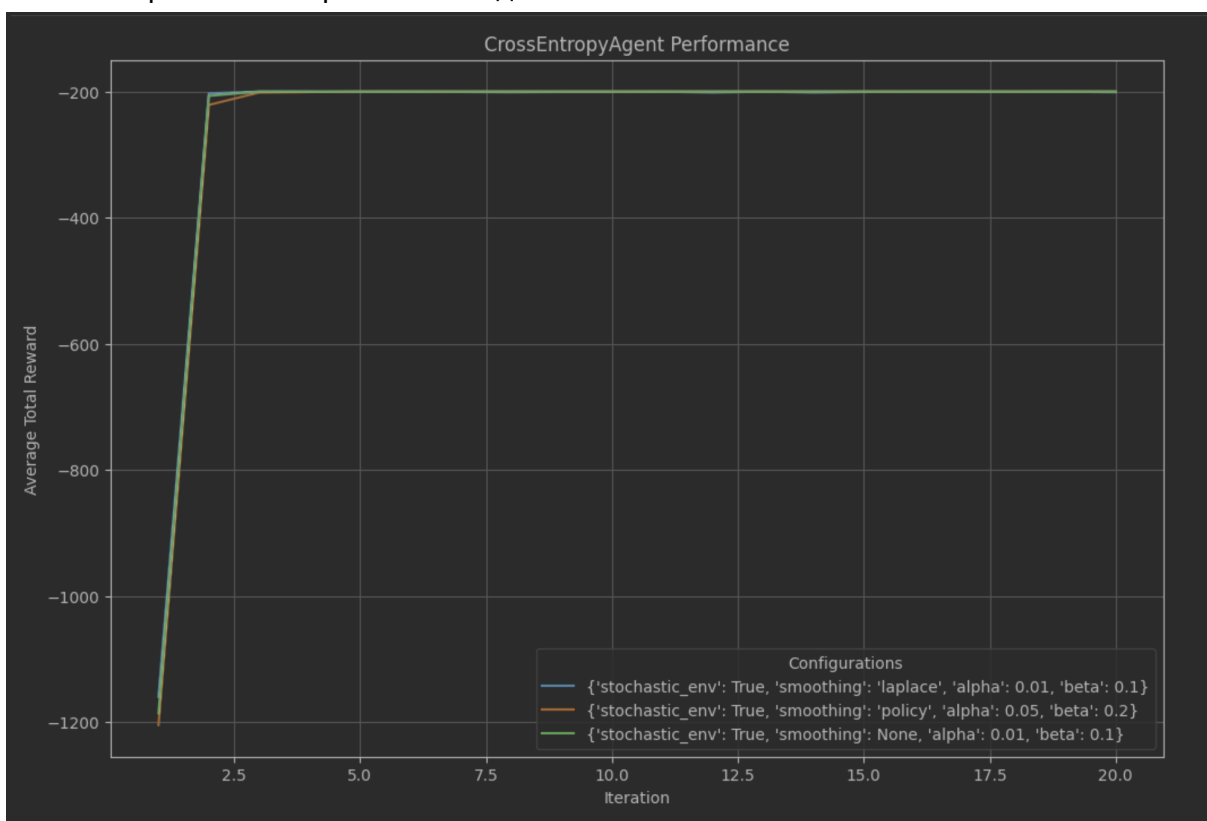
### Обучение

#### Эксперимент 1

В третьем опциональном задании надо было написать комбинировать стохастические и детерминированные политики. Вот что получилось:



Не совсем понял, что здесь произошло, и почему так, но решил увеличить  $m$  в несколько раз и посмотреть что выйдет



Результат тот же самый модел застряла.

## Вывод

Каких-то особых выводов нет, кроме того, что можно по разному прокидывать траектории и политики в модель, и экспериментировать с этим.