

# Отчет

По домашнему заданию №1

Pavel Nakaznenko

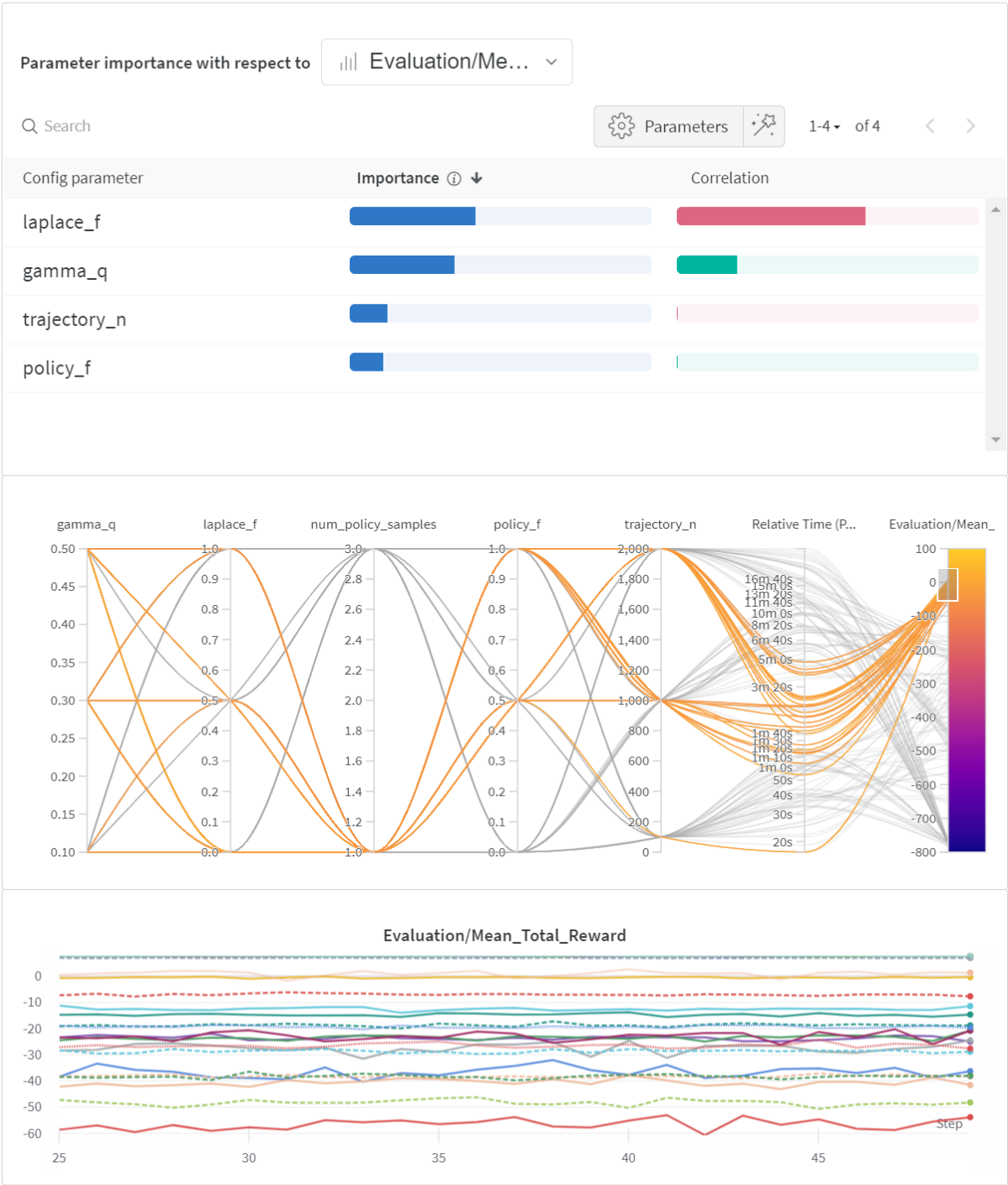
Гиперпараметры:

- `gamma_q` - гамма параметр квантиля
- `laplace_f` - фактор интерполяции по Лапласу
- `policy_f` - фактор линейной интерполяции от предыдущей политики к новой
- `trajectory_n` - количество траекторий в выборке
- `iteration_n` - (фиксированный в данном отчете, равен 25) количество итераций обучения
- `max_trajectory_len` - (фиксированный в данном отчете, равен 200) максимально допустимая длина траектории в выборке. Значение выбрано в соответствии с критерием останова симуляции.

Эксперименты:

1. Имплементация классического Cross Entropy Method
2. Имплементация сглаживания по Лапласу (`laplace_f`) и сглаживания политики (`policy_f`)
3. Имплементация множественного сэмпинга (`num_policy_samples`) стохастической среды (`num_policy_samples > 1`)
4. Определение теоретически возможной максимальной средней награды (эмпирически измеряем `mean total reward` для всех возможных изначальных значений `state`)
5. Интерполяция гиперпараметров `trajectory_n` и `gamma_q` от эпизода к эпизоду
6. `grid search` для определения влияния гиперпараметров

Результаты:



Наблюдения и выводы:

1. Для оптимальной производительности важна "насмотренность": такое сочетание гиперпараметров `trajectory_n` и `gamma_q`, при котором агент "наблюдает" достаточно различных последствий своих действий.
2. Эмпирически посчитанный максимально возможный `mean total reward` 7.93
3. Лучший результат обучения методом кросс энтропии находится в окрестности теоретически возможного на `eval`-сете с большим количеством эпизодов
4. Оптимальные гиперпараметры находятся в окрестности следующих значений: `trajectory_n`  $\geq 1000$ , `gamma_q`  $\leq 0.3$ , `iteration_n`  $\geq 100$
5. Техники по улучшению результата обучения на стохастических средах в среднем дают ожидаемое ухудшение производительности на детерминированной среде, а так же увеличивают время сходимости
6. Интерполяция гиперпараметров `trajectory_n` и `gamma_q` от более обширного обхвата до более элитарного в среднем увеличивает скорость сходимости обучения, сохраняя баланс между результатом и скоростью обучения.