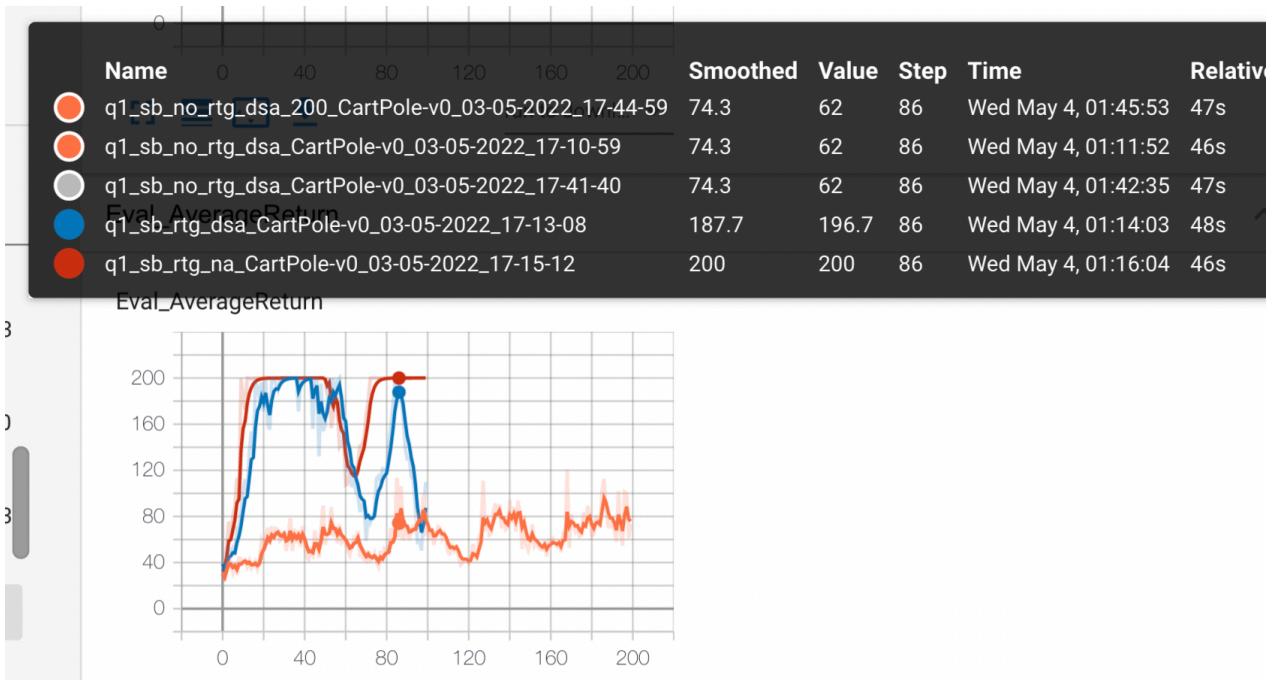


# hw2

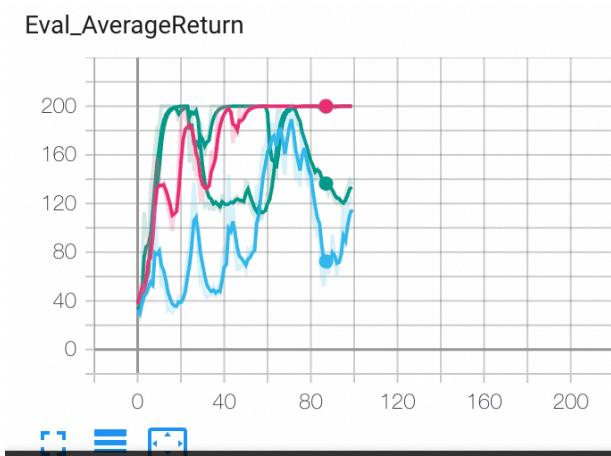
## sec1

### 1. Deliverables for report:

- Create two graphs:
  - – In the first graph, compare the learning curves (average return at each iteration) for the experiments prefixed with q1\_sb\_. (The small batch experiments.)



- – In the second graph, compare the learning curves for the experiments prefixed with q1\_lb\_. (The large batch experiments.)

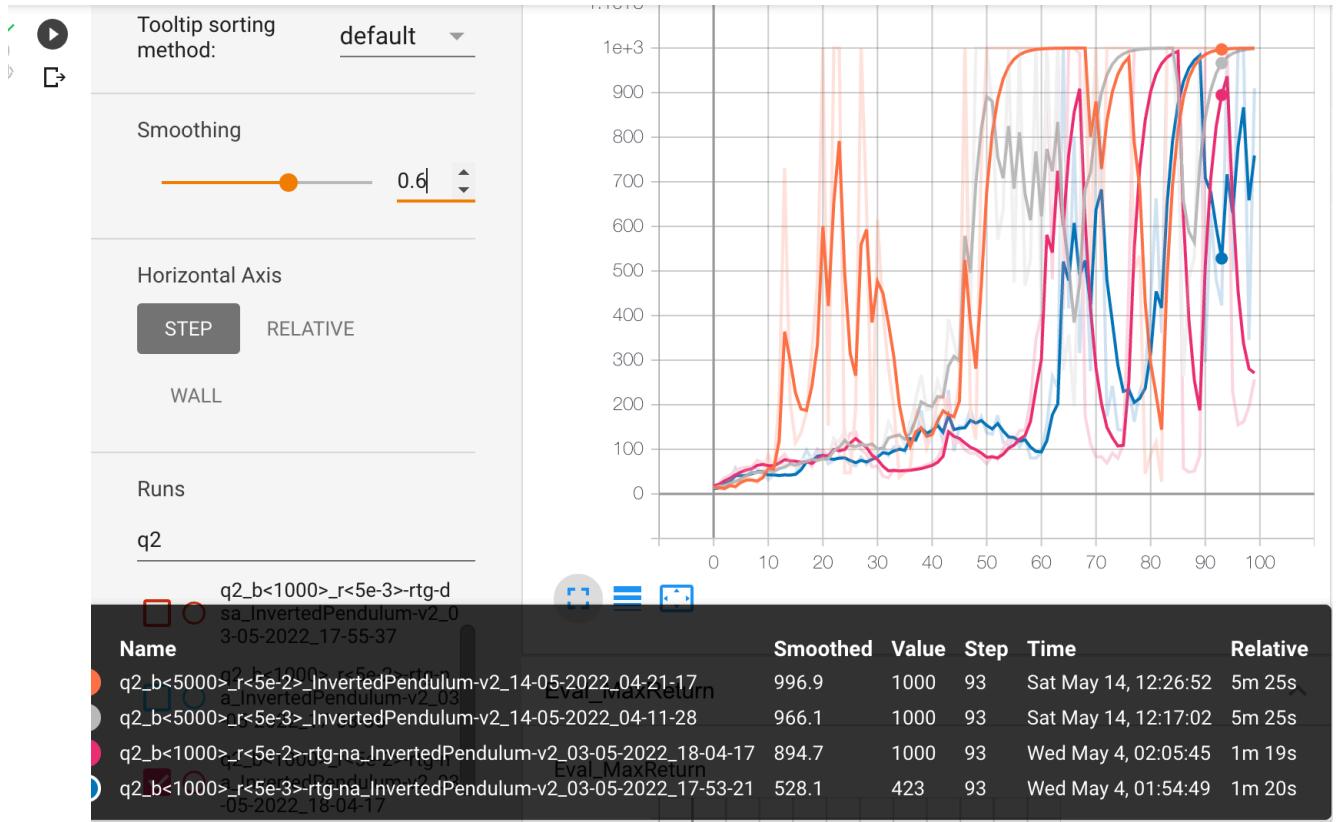


- Answer the following questions briefly:

- Which value estimator has better performance without advantage-standardization: the trajectory-centric one, or the one using reward-to-go?
- Did advantage standardization help? Yes, but indistinctive.
- Did the batch size make an impact?

## sec2

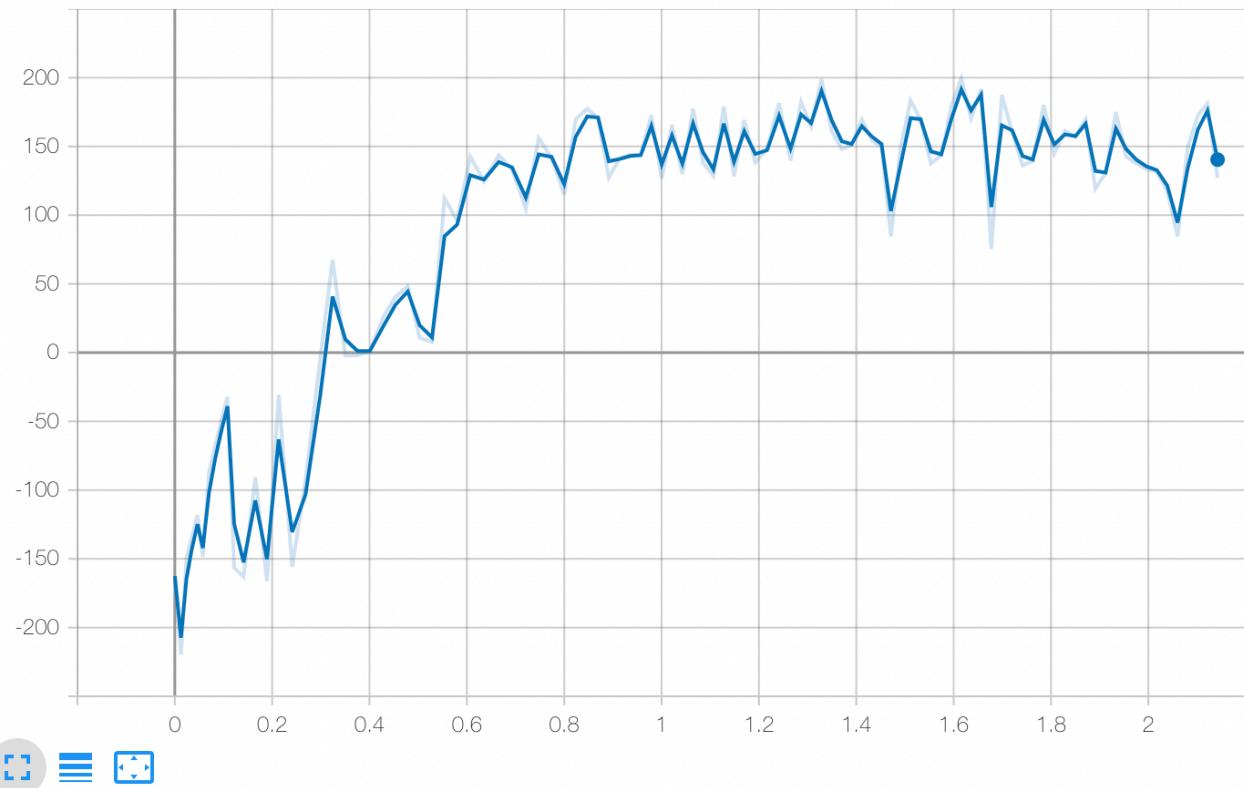
- Large batch size: stable, low variance.
- Small learning rate: converges quicker but unstable



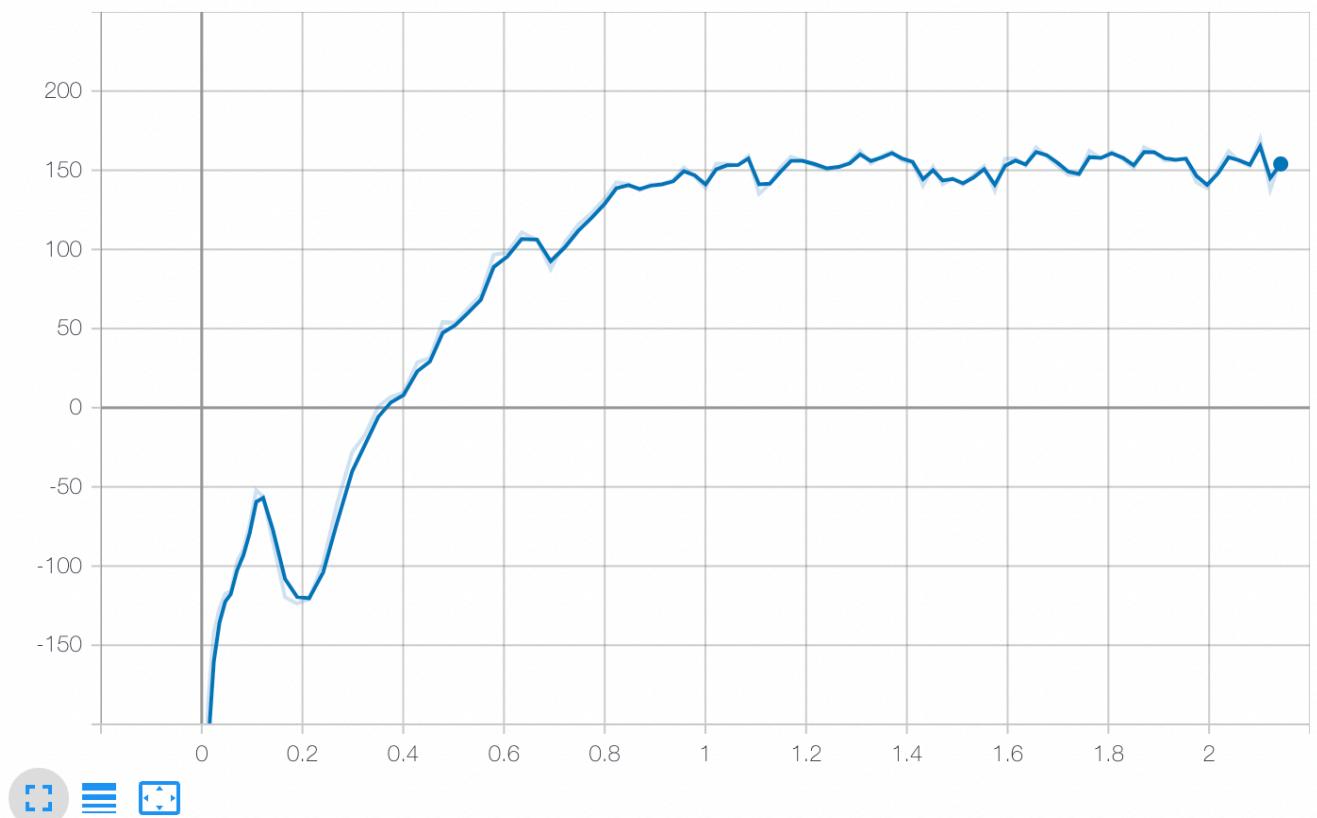
## sec3

- maxreturn is over 200

Eval\_AverageReturn



Train\_AverageReturn

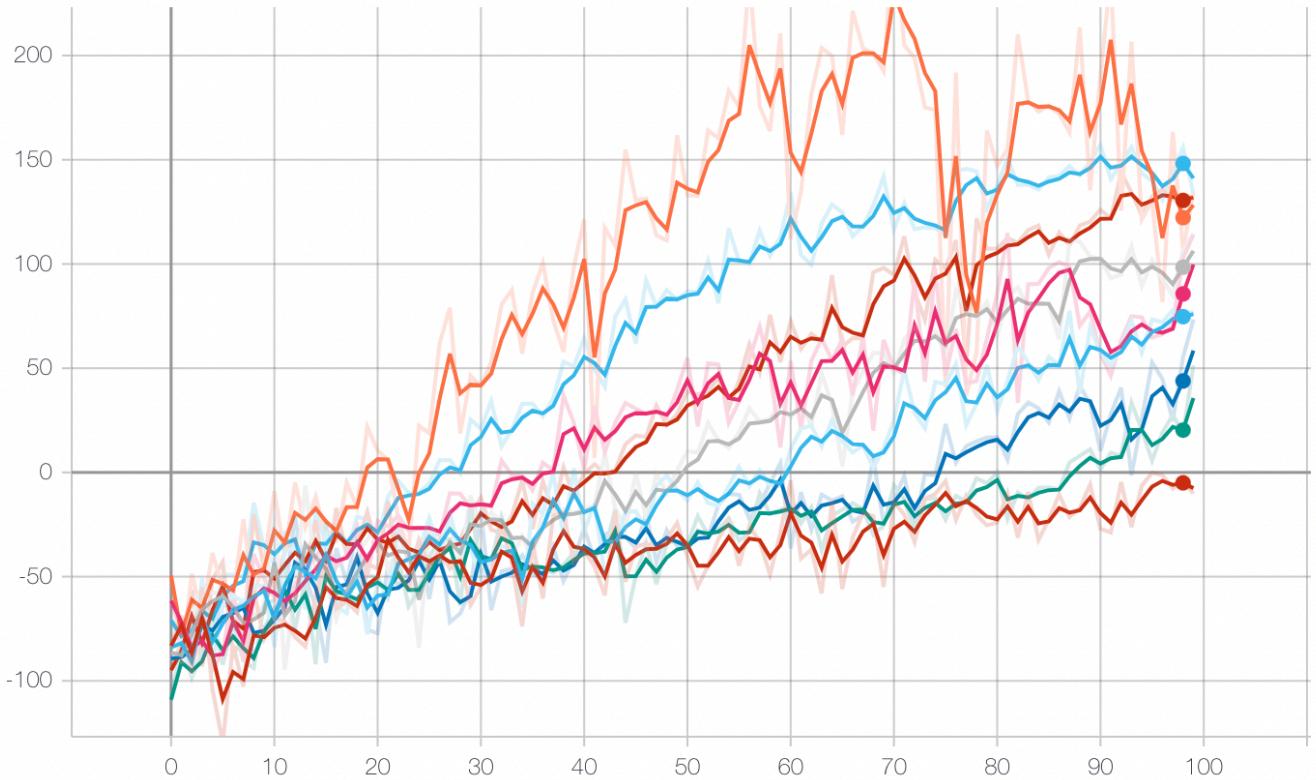


(lec7)要不知道mdp，要不在一个状态用同个策略做好多次（来模拟出状态转移概率）。

而利用q-function，不依赖于策略，只依赖于mdp,在任何策略下采样都可以(模拟mdp)，不知道mdp也无妨

## sec4

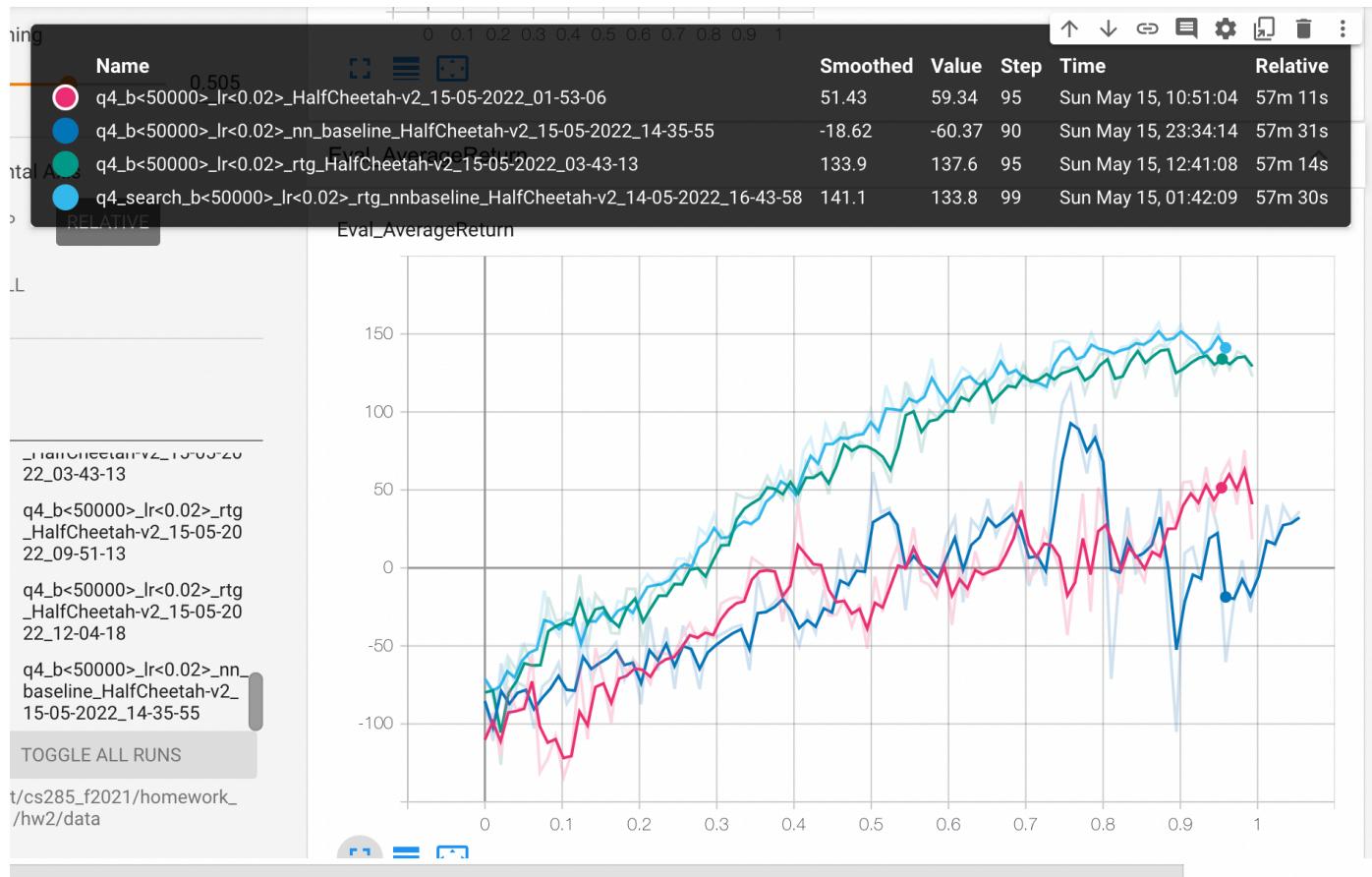
Eval\_AverageReturn



Name	Smoothed	Value	Step	Time	Relative
RELATIVE	-29.28	-24.19	49	Sat May 14, 19:07:37	5m 44s
q4_search_b<10000>_lr<0.005>_rtg_nnbaseline_HalfCheetah-v2_14-05-2022_11-01-38	-8.62	-8.496	49	Sat May 14, 19:26:54	5m 46s
q4_search_b<10000>_lr<0.01>_rtg_nnbaseline_HalfCheetah-v2_14-05-2022_11-20-55	33.62	39.67	49	Sat May 14, 19:44:31	5m 47s
q4_search_b<30000>_lr<0.005>_rtg_nnbaseline_HalfCheetah-v2_14-05-2022_11-54-53	-36.98	-33.02	49	Sat May 14, 20:12:03	16m 44s
q4_search_b<30000>_lr<0.01>_rtg_nnbaseline_HalfCheetah-v2_14-05-2022_12-29-30	-4.106	0.877	49	Sat May 14, 20:47:20	17m 23s
q4_search_b<30000>_lr<0.02>_rtg_nnbaseline_HalfCheetah-v2_14-05-2022_13-09-41	139	161.8	49	Sat May 14, 21:31:05	20m 49s
q4_search_b<50000>_lr<0.005>_rtg_nnbaseline_HalfCheetah-v2_14-05-2022_14-31-09	-33.34	-35.37	49	Sat May 14, 23:00:21	28m 24s
q4_search_b<50000>_lr<0.01>_rtg_nnbaseline_HalfCheetah-v2_14-05-2022_15-31-34	25.41	27.94	49	Sun May 15, 00:00:26	28m 10s
q4_search_b<50000>_lr<0.02>_rtg_nnbaseline_HalfCheetah-v2_14-05-2022_16-43-58	83.09	82.7	49	Sun May 15, 01:12:50	28m 11s

b<30000> lr<0.02>提升很快，方差极大。最终最好的是b<50000>lr<0.02>

rtg nn-baseline对照



**sec5**

