

Hamed Ajorlou

97101167

February 7, 2021

# ML HW5 - Report

## Imitation Learning

قصد داریم راجع به اجزا و عملکرد کلاس `Replay buffer` و `utils` صحبت کنیم. در ابتدا هدف کلی این کلاس را ذکر می‌نماییم و سپس به جزئیات و به متغیرهای مختلف این کلاس می‌پردازیم. `Agent` در حال یادگیری، مسیرهای مختلفی را بنابر `policy` تعیین شده طی می‌نماید. حال برای ما مطلوب است که این مسیرهای طی شده را در حافظه ذخیره سازی بنماییم تا بعداً این داده‌ها را تحلیل بنماییم و مدل را `train` کنیم. `path` به منظور ذخیره سازی این مسیر طی شده می‌باشد. در واقع کلاس `path` اطلاعات هر `Rollout` را دسته بندی کرده و در یک دیکشنری قرار می‌دهد. با هر قدم مشاهدات جدیدی به `agent` داده می‌شود که توسط دستور `append` در لیست مشاهدات `obs` ذخیره می‌شود و بنابر این مشاهدات، جوایز جدید تحت عنوان `rewards` به لیست جوایز `rewards` افزوده می‌شود. در `RL` به روش‌های تکرار شونده تا رسیدن به نتیجه دلخواه و سرمنزل مقصود اصطلاحاً `Roll out` گفته می‌شود. در این کلاس متغیری تحت عنوان `Rollout_done` در نظر گرفته شده است که بنابر طول مسیری که قصد ذخیره آن را داریم و یا رسیدن به نتیجه، مقداردهی می‌شود و حلقه خاتمه می‌یابد. تابع `convert_listofrollouts` نیز اجزای مختلف دیکشنری‌ها را دریافت می‌کند و در آرایه‌های جدا جدا آنها را دسته بندی می‌کند. از آن در کلاس `Replay Buffer` استفاده شده و اجزای مختلف مسیر در آرایه‌های مربوطه دسته بندی می‌شوند سپس در `add_rollouts` به شی‌بافر اضافه می‌شوند. تابع `sample_random_data` که در `Replay Buffer` تعریف شده است به تعدادی که در ورودی به آن داده می‌شود به طور تصادفی اجزای مختلف مسیر را باز می‌گرداند ولی در تابع دیگر `sample_recent_data` مشخصاً مراحل اخیر که در ورودی تابع ذکر شده است را باز می‌گرداند.

حال به کلاس `logger` می‌پردازیم و درباره‌ی متغیرهای ذکر شده توضیحاتی تقدیم می‌نماییم. آنچه در `init` ذکر شده است یک `Summarywriter` می‌سازد که رخدادها و `summary`ها را در `directory` مشخص شده توسط `log_dir` می‌نویسد بنابر این `log_dir` این محل ذخیره سازی را تعیین می‌نماید. آنچه باید ثبت شود در انتظار می‌ماند و پس از گذشت زمانی که توسط `flush_secs` تعیین می‌شود بر روی دیسک نوشته می‌شود و سائز این صف که در انتظار نوشته شدن می‌ماند توسط `max_queue` تعیین می‌شود. (برای کسب اطلاعات بیشتر از این منبع کمک گرفتیم: [pytorch.org](https://pytorch.org))

حال قصد داریم آنچه در حین کامل کردن کد انجام دادیم را شرح دهیم و توضیحاتی پیرامون آن تقدیم بنماییم . از بخش gym-introduction آغاز می نماییم که در این قسمت observation ها دریافت می شود و وزن هایی رندم مربوط به ۴ پارامتر تاثیرگذار ایجاد می نماییم و با ضرب کردن آنها و ساین گرفتن از نتیجه ی آن ، action حاصل می شود . حال به شکلی تکرار شونده این action های بدست آمده را دوباره به environment می دهیم تا دوباره observation دریافت کنیم این روند ادامه می یابد تا به نقطه ی Done برسیم و حلقه متوقف شود و بازی تمام شود برای مثال میله افقی شود . همچنین بیشترین مرحله ای که از شروع تا خاتمه ی بازی میگذرد و در متغیر cnt ثبت می شود را به عنوان بهترین وزن ذخیره مینماییم و در انتها یک بار با این وزن ها الگوریتم را اجرا می نماییم .

در قسمت torch utils یک مدل نورال نت طراحی می نماییم به مانند تمرین قبلی و نکته ی مهم و قابل ذکر در مفهوم جدیدی به نام xavier initialization می باشد . دلیل اینکه از xavier initialization استفاده می نماییم آن است که سطح نورال نت ساخته شده لوکال مینیمم های فراوانی دارد و برای آنکه به نتیجه ی مطلوب برسد به وزن های اولیه توزیع احتمالاتی خاصی نسبت می دهیم . من با جستجو در این باره از توزیع احتمالاتی یونیفورم استفاده کردم . طبق آنچه صحبت کردیم policy مدل از یک نورال نت نشئت می گیرد و در قسمت MLP policy قصد داریم ازین شبکه action ها را استخراج کنیم . برای این کار observation ها را به عنوان ورودی به مدل می دهیم و خروجی آن را دریافت می نماییم . با توجه به خروجی ، در صورتی که environment گسسته باشد ، خروجی شبکه را به عنوان یک مدل احتمالاتی چند جمله ای در نظر میگیریم و در صورتیکه environment پیوسته باشد ، اکشن از یک مدل گاوسی گرفته می شود که مطابق زیر است .

$$a_i \sim \text{mean}(s_i) + e^{\log \text{std}} \mathcal{N}(0, 1)$$

سپس برای دو حالت گسسته و پیوسته با Loss function های تعیین شده در کامنت ، optimizer طراحی می نماییم و در بخش update این Loss function را مینیمایز می نماییم . در بخش loaded\_gaussian\_policy مدل های ارائه شده به عنوان expert در مدل load میشوند و در قسمت rl\_trainer در iteration=0 تعیین مینماییم که از policy تعیین شده توسط expert تبعیت کند و در iteration های بعدی از collect policy و در ادامه observation ها را به expert policy می دهیم و action های بدست آمده را جایگزین میکنیم با action هایی که خود expert پیش بینی کرده بود که ممکن است اشتباه باشد و به نتیجه ی دلخواه converge نکنند . با تکمیل همه ی این قسمت ها و ران کردن کد های ذکر شده در فایل Read me نتایج بدست آمده ، که در ادامه ذکر می شود و همچنین تصاویر تنسور بورد که از local host دریافت می شود در انتهای گزارش درج شده است .

```
The best length is: 79.55
The best length is: 130.86
The best length is: 139.48
The best length is: 200.0
The best length is: 200.0
The best length is: 200.0
The best length is: 200.0
The best length is: 200.0
The best length is: 200.0
The best length is: 200.0
game lasted 10 moves
```

سپس کد های ذکر شده مربوط به یادگیری تقلیدی را به ترتیب ران میکنیم :

```
Collecting data for eval...
Eval_AverageReturn : 500.0
Eval_StdReturn : 0.0
Eval_MaxReturn : 500.0
Eval_MinReturn : 500.0
Eval_AverageEpLen : 500.0
Train_AverageReturn : 496.3999938964844
Train_StdReturn : 5.083305835723877
Train_MaxReturn : 500.0
Train_MinReturn : 487.0
Train_AverageEpLen : 496.4
Train_EnvstepsSoFar : 0
TimeSinceStart : 2.9023168087005615
Initial_DataCollection_AverageReturn : 496.3999938964844
Done logging...
```

```
Collecting data for eval...
Eval_AverageReturn : 114.39260864257812
Eval_StdReturn : 131.46856689453125
Eval_MaxReturn : 245.86117553710938
Eval_MinReturn : -17.075965881347656
Eval_AverageEpLen : 158.5
Train_AverageReturn : 218.96922302246094
Train_StdReturn : 87.94571685791016
Train_MaxReturn : 280.5846252441406
Train_MinReturn : 55.61897277832031
Train_AverageEpLen : 230.2222222222223
Train_EnvstepsSoFar : 0
TimeSinceStart : 3.4859776496887207
Initial_DataCollection_AverageReturn : 218.96922302246094
Done logging...
```

```
Collecting data for eval...
Eval_AverageReturn : 258.10260009765625
Eval_StdReturn : 0.0
Eval_MaxReturn : 258.10260009765625
Eval_MinReturn : 258.10260009765625
Eval_AverageEpLen : 331.0
Train_AverageReturn : 172.6307373046875
Train_StdReturn : 15.753379821777344
Train_MaxReturn : 188.38412475585938
Train_MinReturn : 156.8773651123047
Train_AverageEpLen : 1000.0
Train_EnvstepsSoFar : 0
TimeSinceStart : 5.518812417984009
Initial_DataCollection_AverageReturn : 172.6307373046875
Done logging...
```

: Dagger

۳ نتیجه از ۱۰ بار تکرار برای هر یک از دگر ها در گزارش ذکر شده است :

## LunarLander

```
Training agent using sampled data from replay buffer...
train step # 0   loss: tensor(0.5199, grad_fn=<NllLossBackward>)
train step # 200 loss: tensor(0.5926, grad_fn=<NllLossBackward>)
train step # 400 loss: tensor(0.6996, grad_fn=<NllLossBackward>)
train step # 600 loss: tensor(0.5138, grad_fn=<NllLossBackward>)
train step # 800 loss: tensor(0.4728, grad_fn=<NllLossBackward>)
itr # 7 :   loss: tensor(0.5944, grad_fn=<NllLossBackward>)
```

Beginning logging procedure...

Collecting data for eval...

```
Eval_AverageReturn : 306.7672119140625
Eval_StdReturn : 0.0
Eval_MaxReturn : 306.7672119140625
Eval_MinReturn : 306.7672119140625
Eval_AverageEplen : 256.0
Train_AverageReturn : 277.0559997558594
Train_StdReturn : 10.193094253540039
Train_MaxReturn : 293.86517333984375
Train_MinReturn : 263.8361511230469
Train_AverageEplen : 216.4
Train_EnvstepsSoFar : 7730
TimeSinceStart : 24.794593334197998
Initial_DataCollection_AverageReturn : 218.96922302246094
Done logging...
```

```
Training agent using sampled data from replay buffer...
train step # 0   loss: tensor(0.7175, grad_fn=<NllLossBackward>)
train step # 200 loss: tensor(0.5548, grad_fn=<NllLossBackward>)
train step # 400 loss: tensor(0.6474, grad_fn=<NllLossBackward>)
train step # 600 loss: tensor(0.4478, grad_fn=<NllLossBackward>)
train step # 800 loss: tensor(0.5638, grad_fn=<NllLossBackward>)
itr # 8 :   loss: tensor(0.5479, grad_fn=<NllLossBackward>)
```

Beginning logging procedure...

Collecting data for eval...

```
Eval_AverageReturn : 268.0810852050781
Eval_StdReturn : 0.0
Eval_MaxReturn : 268.0810852050781
Eval_MinReturn : 268.0810852050781
Eval_AverageEplen : 251.0
Train_AverageReturn : 172.53292846679688
Train_StdReturn : 107.80461883544922
Train_MaxReturn : 270.86602783203125
Train_MinReturn : 28.444793701171875
Train_AverageEplen : 219.8
Train_EnvstepsSoFar : 8829
TimeSinceStart : 28.025717735290527
Initial_DataCollection_AverageReturn : 218.96922302246094
Done logging...
```

```

Training agent using sampled data from replay buffer...
train step # 0   loss: tensor(0.5691, grad_fn=<NllLossBackward>)
train step # 200 loss: tensor(0.6448, grad_fn=<NllLossBackward>)
train step # 400 loss: tensor(0.5738, grad_fn=<NllLossBackward>)
train step # 600 loss: tensor(0.5859, grad_fn=<NllLossBackward>)
train step # 800 loss: tensor(0.5347, grad_fn=<NllLossBackward>)
itr # 9 :   loss: tensor(0.5369, grad_fn=<NllLossBackward>)

Beginning logging procedure...

Collecting data for eval...
Eval_AverageReturn : 250.9696502685547
Eval_StdReturn : 0.0
Eval_MaxReturn : 250.9696502685547
Eval_MinReturn : 250.9696502685547
Eval_AverageEpLen : 225.0
Train_AverageReturn : 194.44985961914062
Train_StdReturn : 127.45550537109375
Train_MaxReturn : 308.0882263183594
Train_MinReturn : 37.22319412231445
Train_AverageEpLen : 204.6
Train_EnvstepsSoFar : 9852
TimeSinceStart : 30.928033351898193
Initial_DataCollection_AverageReturn : 218.96922302246094
Done logging...

```

## LunarLandercontinuous

```

Training agent using sampled data from replay buffer...
train step # 0   loss: tensor(0.5436, grad_fn=<MseLossBackward>)
train step # 200 loss: tensor(0.2331, grad_fn=<MseLossBackward>)
train step # 400 loss: tensor(0.3271, grad_fn=<MseLossBackward>)
train step # 600 loss: tensor(0.1652, grad_fn=<MseLossBackward>)
train step # 800 loss: tensor(0.2435, grad_fn=<MseLossBackward>)
itr # 7 :   loss: tensor(0.2926, grad_fn=<MseLossBackward>)

Beginning logging procedure...

Collecting data for eval...
Eval_AverageReturn : 180.29721069335938
Eval_StdReturn : 94.48086547851562
Eval_MaxReturn : 274.778076171875
Eval_MinReturn : 85.81634521484375
Eval_AverageEpLen : 195.5
Train_AverageReturn : 255.3571014404297
Train_StdReturn : 14.89013957977295
Train_MaxReturn : 271.17547607421875
Train_MinReturn : 232.22781372070312
Train_AverageEpLen : 258.0
Train_EnvstepsSoFar : 8435
TimeSinceStart : 36.78365659713745
Initial_DataCollection_AverageReturn : 172.6307373046875
Done logging...

```



```
Training agent using sampled data from replay buffer...
train step # 0   loss: tensor(0.2470, grad_fn=<MseLossBackward>)
train step # 200 loss: tensor(0.2110, grad_fn=<MseLossBackward>)
train step # 400 loss: tensor(0.2203, grad_fn=<MseLossBackward>)
train step # 600 loss: tensor(0.3742, grad_fn=<MseLossBackward>)
train step # 800 loss: tensor(0.2181, grad_fn=<MseLossBackward>)
itr # 8 :   loss: tensor(0.3115, grad_fn=<MseLossBackward>)
```

Beginning logging procedure...

Collecting data for eval...

```
Eval_AverageReturn : 308.8662109375
Eval_StdReturn : 0.0
Eval_MaxReturn : 308.8662109375
Eval_MinReturn : 308.8662109375
Eval_AverageEpLen : 280.0
Train_AverageReturn : 265.43487548828125
Train_StdReturn : 17.038076400756836
Train_MaxReturn : 285.86065673828125
Train_MinReturn : 241.5272979736328
Train_AverageEpLen : 270.25
Train_EnvstepsSoFar : 9516
TimeSinceStart : 40.414562702178955
Initial_DataCollection_AverageReturn : 172.6307373046875
Done logging...
```

```
Training agent using sampled data from replay buffer...
train step # 0   loss: tensor(0.4508, grad_fn=<MseLossBackward>)
train step # 200 loss: tensor(0.3913, grad_fn=<MseLossBackward>)
train step # 400 loss: tensor(0.2752, grad_fn=<MseLossBackward>)
train step # 600 loss: tensor(0.2460, grad_fn=<MseLossBackward>)
train step # 800 loss: tensor(0.2213, grad_fn=<MseLossBackward>)
itr # 9 :   loss: tensor(0.1992, grad_fn=<MseLossBackward>)
```

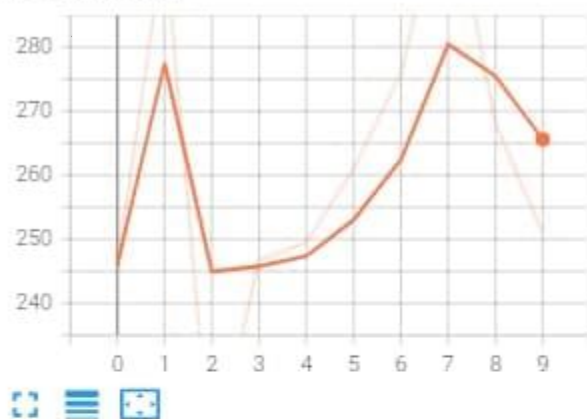
Beginning logging procedure...

Collecting data for eval...

```
Eval_AverageReturn : 43.98357009887695
Eval_StdReturn : 0.0
Eval_MaxReturn : 43.98357009887695
Eval_MinReturn : 43.98357009887695
Eval_AverageEpLen : 255.0
Train_AverageReturn : 226.8162841796875
Train_StdReturn : 120.89540100097656
Train_MaxReturn : 323.73065185546875
Train_MinReturn : 19.594528198242188
Train_AverageEpLen : 283.75
Train_EnvstepsSoFar : 10651
TimeSinceStart : 44.074424266815186
Initial_DataCollection_AverageReturn : 172.6307373046875
Done logging...
```

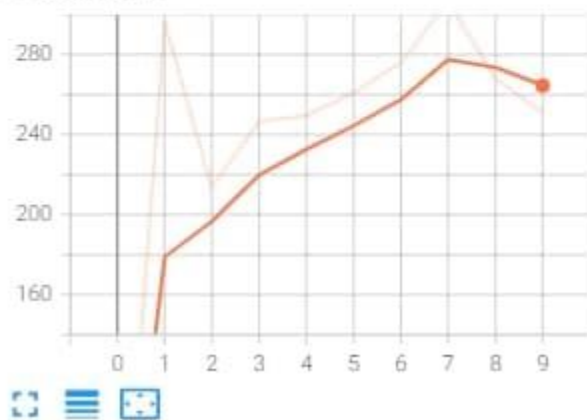
برخی از تنسور بورد های بدست آمده :

Eval\_MaxReturn



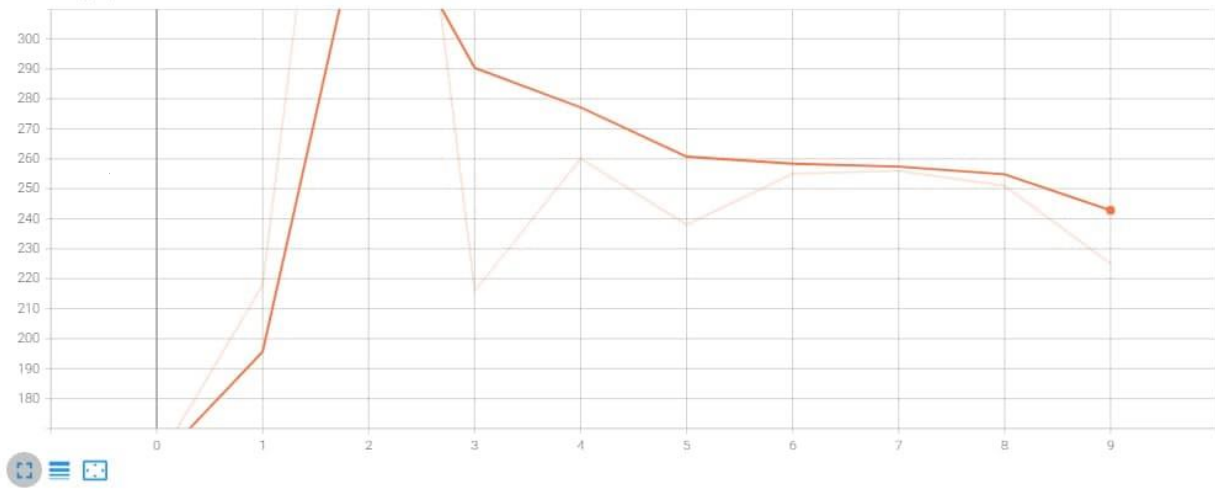
Eval\_MinReturn

Eval\_MinReturn

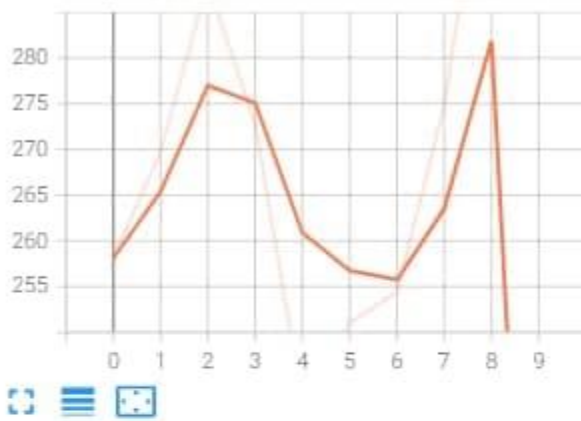




Eval\_AverageEpLen

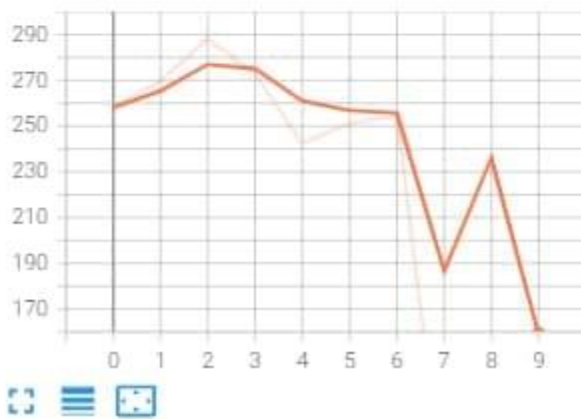


Eval\_MaxReturn

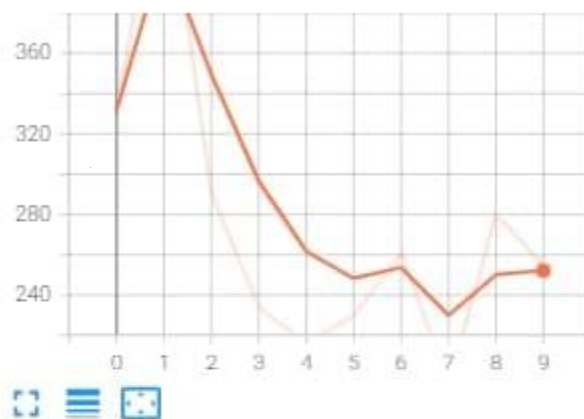


Eval\_MinReturn

Eval\_MinReturn

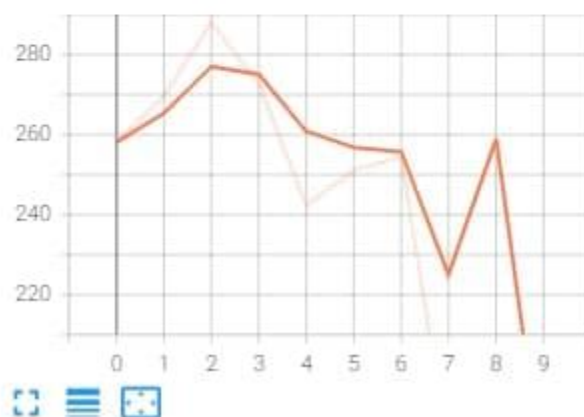


Eval\_AverageEpLen



Eval\_AverageReturn

Eval\_AverageReturn



با تشکر از همراهی شما