

## **Learning Complex Neural Network Policies with Trajectory Optimization**

#### Sergey Levine

SVLEVINE@CS.STANFORD.EDU

Computer Science Department, Stanford University, Stanford, CA 94305 USA

Vladlen Koltun

VLADLEN@ADOBE.COM

Adobe Research, San Francisco, CA 94103 USA

## 【强化学习算法 14】GPS



张楚珩 🗸

清华大学 交叉信息院博士在读

10 人赞同了该文章

这里的GPS不是用来导航的GPS,其全称是Guided Policy Search,是一种局域的model-based RL 算法,各种版本基本上都是Levine大神搞出来的。

### 原文传送门:

[Importance-sampled GPS] Levine, Sergey, and Vladlen Koltun. "Guided policy search." International Conference on Machine Learning. 2013.

【Variational GPS】 Levine, Sergey, and Vladlen Koltun. "Variational policy search via trajectory optimization." Advances in Neural Information Processing Systems. 2013.

【Constraint GPS】 Levine, Sergey, and Vladlen Koltun. "Learning complex neural network policies with trajectory optimization."International Conference on Machine Learning. 2014.

【Unknown Dynamics GPS】 Levine, Sergey, and Pieter Abbeel. "Learning neural network policies with guided policy search under unknown dynamics." Advances in Neural Information Processing Systems. 2014.

这里有好多版本的, 暂时这里只讲第三篇工作。

## 特色:

这里讲了一种model-based的算法,主要思路是先找一条state-action的路径,然后交替地1)优化路径,使得该路径既离现有的控制策略不远,又能够最小化cost; 2)优化策略,使得策略离被优化过的路径更近。通过这样trajectory optimization + policy search的方式在bipedal push recovery和walking on uneven terrain问题上取得了较好的实验结果。

#### 分类:

continuous state space continuous action space model-based with known model dynamics

#### 过程:

#### 1. 定义优化问题

首先从一个类似SQL算法的优化目标开始

$$egin{aligned} \min_{ heta,q( au)} D_{\mathrm{KL}}(q( au) \| 
ho( au)) \ & ext{s.t.} \ q(\mathbf{x}_1) = p(\mathbf{x}_1), \ q(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t) = p(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{u}_t), \ D_{\mathrm{KL}}(q(\mathbf{x}_t) \pi_{\theta}(\mathbf{u}_t | \mathbf{x}_t) \| q(\mathbf{x}_t, \mathbf{u}_t)) = 0. \end{aligned}$$

首先先明确一下轨迹,是状态。和控制。交替的一个序列。其中  $\rho(r) \propto \exp(-l(r))$  是最优轨迹分布,注意到当cost幅值比较大的时候这个最优轨迹更趋近于确定性的轨迹。通过约束,轨迹分布 q(r) 就是参数控制的策略  $\pi_0$  下形成的轨迹。之所以不直接对最优轨迹  $\rho(r)$  优化策略,而是要经过中间的 q(r),是由于在提供了环境模型的情况下,轨迹更容易被优化。

## 2. dual gradient descent求解

接下来使用dual gradient descent方法来求解上述优化问题,得到算法的大框架

$$\mathcal{L}(\theta, q, \lambda) = D_{\text{KL}}(q(\tau) \| \rho(\tau)) +$$
 
$$\sum_{t=1}^{T} \lambda_t D_{\text{KL}}(q(\mathbf{x}_t) \pi_{\theta}(\mathbf{u}_t | \mathbf{x}_t) \| q(\mathbf{x}_t, \mathbf{u}_t)).$$
 知乎 @张建珩

$$\lambda_t \leftarrow \lambda_t + \eta D_{\text{KL}}(q(\mathbf{x}_t)\pi_{\theta}(\mathbf{u}_t|\mathbf{x}_t)||q(\mathbf{x}_t,\mathbf{u}_t)).$$
 (2)

### Algorithm 1 Constrained guided policy search

- 1: Initialize the trajectories  $\{q_1(\tau), \ldots, q_M(\tau)\}$
- 2: for iteration k = 1 to K do
- 3: Optimize each  $q_i(\tau)$  with respect to  $\mathcal{L}(\theta, q_i(\tau), \lambda_i)$
- 4: Optimize  $\theta$  with respect to  $\sum_{i=1}^{M} \mathcal{L}(\theta, q_i(\tau), \lambda_i)$
- 5: Update dual variables  $\lambda$  using Equation 2
- 6: end for
- 7: **return** optimized policy parameters  $\theta$  知乎 @张楚珩

关于dual gradient descent方法可以参考这篇帖子: Dual Gradient Descent

这个主算法里面的第3行是在固定  $_{n}$  和  $_{\lambda}$  的情况下优化轨迹  $_{q(r)}$  ,第4行是在固定  $_{q(r)}$  和  $_{\lambda}$  的情况下优化策略  $_{n}$  ,下面分别介绍如何进行这两步操作。

假定轨迹 q(r) 是根据线性策略  $q(u|z_1) \sim N(K_1z_1+k_1,k_1)$  和线性化的动力学  $q(z_1,z_1) \sim N(f_{1t}z_1+f_{1t}u_1,R_1)$  得到的,其整个轨迹上都是高斯分布,令该轨迹分布的均值为  $f = (g_{1t},g_{1t})$  。说明:用  $g_{1t} = \begin{bmatrix} z_1 \\ z_1 \end{bmatrix}$  表示拼起来的向量;下标都表示求导;为了和文中一致下标 表示对  $g_{1t} = g_{1t}$  和 上标的表示原来的轨迹,没有 hat 上标的表示相对于原来轨迹的扰动(即小量)。

轨迹是由哪些量描述的呢?我们这里说的轨迹是一个轨迹的概率分布,它是一个高斯型的分布,因此决定它的除了均值,,之外,还有,的协方差矩阵

$$\Sigma_t = \begin{bmatrix} \mathbf{S}_t & \mathbf{S}_t \mathbf{K}_t^{\mathrm{T}} \\ \mathbf{K}_t \mathbf{S}_t & \mathbf{K}_t \mathbf{S}_t \mathbf{K}_t^{\mathrm{T}} + \mathbf{A}_t \end{bmatrix}.$$

其中  $s_t = Cov[x_t]$  为定义, 其他元素容易求出。

后面的过程比较复杂,我们先要明确一下要我们要干嘛。我们的目标是把问题化为LQG问题(动力学线性、损失函数二次型),然后利用类似iLQG的最优控制公式求解当前轨迹 $_{(r)}$  附近的更好的新轨迹 $_{(r)}$ 。我们这里说的轨迹是由如下参数构成的 $_{(r)}=\{r,\Sigma(S,A,K)\}$ 。

首先对轨迹进行优化的时候认为 和 x 是固定的,我们可以把损失函数写成只与轨迹有关的形式,通过一些计算之后能够得到(注意,这里省略了对于时间 t 的求和)

$$\mathcal{L}(q) \approx \sum_{t=1}^{T} \frac{1}{2} \begin{bmatrix} \hat{\mathbf{x}}_{t} \\ \hat{\mathbf{u}}_{t} \end{bmatrix}^{\mathrm{T}} \ell_{\mathbf{x}\mathbf{u},\mathbf{x}\mathbf{u}t} \begin{bmatrix} \hat{\mathbf{x}}_{t} \\ \hat{\mathbf{u}}_{t} \end{bmatrix} + \begin{bmatrix} \hat{\mathbf{x}}_{t} \\ \hat{\mathbf{u}}_{t} \end{bmatrix}^{\mathrm{T}} \ell_{\mathbf{x}\mathbf{u}t} + \frac{1}{2} \operatorname{tr} \left( \sum_{t} \ell_{\mathbf{x}\mathbf{u},\mathbf{x}\mathbf{u}t} \right) - \frac{1}{2} \log |\mathbf{A}_{t}| + \frac{\lambda_{t}}{2} \log |\mathbf{A}_{t}| + \frac{\lambda_{t}}{2} \log |\mathbf{A}_{t}| + \frac{\lambda_{t}}{2} (\hat{\mathbf{u}}_{t} - \mu_{t}^{\pi}(\hat{\mathbf{x}}_{t}))^{\mathrm{T}} \mathbf{A}_{t}^{-1} (\hat{\mathbf{u}}_{t} - \mu_{t}^{\pi}(\hat{\mathbf{x}}_{t})) + \frac{\lambda_{t}}{2} \operatorname{tr} \left( \mathbf{A}_{t}^{-1} \sum_{t}^{\pi} \right) + \frac{\lambda_{t}}{2} \operatorname{tr} \left( \mathbf{S}_{t} \left( \mathbf{K}_{t} - \mu_{\mathbf{x}t}^{\pi}(\hat{\mathbf{x}}_{t}) \right)^{\mathrm{T}} \mathbf{A}_{t}^{-1} \left( \mathbf{K}_{t} - \mu_{\mathbf{x}t}^{\pi}(\hat{\mathbf{x}}_{t}) \right) \right),$$

类比 iLQG 算法里面的损失函数  $c_{\sim (g+Ge)u+\frac{1}{2}u^THu}$  ,有最优控制  $u_{i}=k_{i}+K_{i}z_{i}$  。我们的思路就是针对上面损失函数把相应的项找出来,其方法也很简单就是把上面的损失函数对控制  $u_{i}$  求导,可以得到

$$\mathcal{L}_{\mathbf{u}t} = Q_{\mathbf{u},\mathbf{u}t}\hat{\mathbf{u}}_t + Q_{\mathbf{u},\mathbf{x}t}\hat{\mathbf{x}}_t + Q_{\mathbf{u}t} + \lambda_t \mathbf{A}_t^{-1}(\hat{\mathbf{u}}_t - \mu_t^{\pi}(\hat{\mathbf{x}}_t))$$

$$\mathcal{L}_{\mathbf{u},\mathbf{u}t} = Q_{\mathbf{u},\mathbf{u}t} + \lambda_t \mathbf{A}_t^{-1},$$

但是当你求的时候会发现,求出来的形式跟文章一样,但是你式子里面出现单状态损失函数  $_{\bf q}$  (这个说法不是规范的说法,为了方便叙述我自己编的)的时候,文章里面都写的是  $_{\bf q}$  。为什么这样呢?因为当你改变了某一处控制函数的时候,它后面相应的轨迹也都发生了改变,因此当某一处控制函数的改变导致了单状态损失函数的时候,应该连带上它对于后续轨迹上损失函数的改变,即用下面的  $_{\bf q}$  来替代即可

$$Q_{\mathbf{x}\mathbf{u}t} = \ell_{\mathbf{x}\mathbf{u}t} + f_{\mathbf{x}\mathbf{u}}^{\mathrm{T}} \mathcal{L}_{\mathbf{x}t+1}$$
$$Q_{\mathbf{x}\mathbf{u},\mathbf{x}\mathbf{u}t} = \ell_{\mathbf{x}\mathbf{u},\mathbf{x}\mathbf{u}t} + f_{\mathbf{x}\mathbf{u}}^{\mathrm{T}} \mathcal{L}_{\mathbf{x},\mathbf{x}t+1} f_{\mathbf{x}\mathbf{u}},$$

接下来,比对 iLQG 的算法公式能够得到新的最优控制扰动

$$\mathbf{k}_{t} = -\left(Q_{\mathbf{u},\mathbf{u}t} + \lambda_{t}\mathbf{A}_{t}^{-1}\right)^{-1}\left(Q_{\mathbf{u}t} - \lambda_{t}\mathbf{A}_{t}^{-1}\mu_{t}^{\pi}(\hat{\mathbf{x}}_{t})\right). \tag{3}$$

$$\mathbf{K}_{t} = -\left(Q_{\mathbf{u},\mathbf{u}t} + \lambda_{t} \mathbf{A}_{t}^{-1}\right)^{-1} \left(Q_{\mathbf{u},\mathbf{x}t} - \lambda_{t} \mathbf{A}_{t}^{-1} \mu_{\mathbf{x}t}^{\pi}(\hat{\mathbf{x}}_{t})\right). \tag{4}$$

注意到,其中Q的计算是包含损失函数在后一步上的导数的,因此该过程需要通过逆序的动态规划来求解。计算Q会用到的损失函数计算公式如下。

$$\mathcal{L}_{\mathbf{x}t} = Q_{\mathbf{x},\mathbf{x}t}\hat{\mathbf{x}}_{t} + Q_{\mathbf{x},\mathbf{u}t}(\mathbf{k}_{t} + \mathbf{K}_{t}\hat{\mathbf{x}}_{t}) + \mathbf{K}_{t}^{\mathrm{T}}Q_{\mathbf{u},\mathbf{x}t}\hat{\mathbf{x}}_{t} + \mathbf{K}_{t}^{\mathrm{T}}Q_{\mathbf{u},\mathbf{x}t}(\mathbf{K}_{t}\hat{\mathbf{x}}_{t} + \mathbf{k}_{t}) + Q_{\mathbf{x}t} + \mathbf{K}_{t}^{\mathrm{T}}Q_{\mathbf{u}t} + \lambda_{t}(\mathbf{K}_{t} - \mu_{\mathbf{x}t}^{\pi}(\hat{\mathbf{x}}_{t}))^{\mathrm{T}}\mathbf{A}_{t}^{-1}(\mathbf{K}_{t}\hat{\mathbf{x}}_{t} + \mathbf{k}_{t} - \mu_{t}^{\pi}(\hat{\mathbf{x}}_{t}))$$

$$= Q_{\mathbf{x},\mathbf{u}t}\mathbf{k}_{t} + \mathbf{K}_{t}^{\mathrm{T}}Q_{\mathbf{u},\mathbf{u}t}\mathbf{k}_{t} + Q_{\mathbf{x}t} + \mathbf{K}_{t}^{\mathrm{T}}Q_{\mathbf{u}t} + \lambda_{t}(\mathbf{K}_{t} - \mu_{\mathbf{x}t}^{\pi}(\hat{\mathbf{x}}_{t}))^{\mathrm{T}}\mathbf{A}_{t}^{-1}(\mathbf{k}_{t} - \mu_{t}^{\pi}(\hat{\mathbf{x}}_{t}))$$

$$\mathcal{L}_{\mathbf{x},\mathbf{x}t} = Q_{\mathbf{x},\mathbf{x}t} + Q_{\mathbf{x},\mathbf{u}t}\mathbf{K}_{t} + \mathbf{K}_{t}^{\mathrm{T}}Q_{\mathbf{u},\mathbf{x}t} + \mathbf{K}_{t}^{\mathrm{T}}Q_{\mathbf{u},\mathbf{u}t}\mathbf{K}_{t} + \lambda_{t}(\mathbf{K}_{t} - \mu_{\mathbf{x}t}^{\pi}(\hat{\mathbf{x}}_{t}))^{\mathrm{T}}\mathbf{A}_{t}^{-1}(\mathbf{K}_{t} - \mu_{\mathbf{x}t}^{\pi}(\hat{\mathbf{x}}_{t})), \quad \text{with } \mathbf{x} \in \mathcal{L}_{t}^{\mathrm{T}}(\mathbf{x}_{t})$$

由此我们能够确定新的轨迹的均值  $f_{new} = f + \tau$  ,而该新轨迹的协方差矩阵  $\Sigma_{new}$  也需要被确定,观察上面关于协方差矩阵的定义,如果我们求出了  $s_i$  、  $A_i$  和  $K_i$  ,就可以确定轨迹的协方差矩阵了。

知道新轨迹的均值之后怎么求出相应的  $s_a$ 、  $a_a$  和  $a_b$  呢?自然想到的就是对损失函数相对于要求的这些量求导,并且令导数为零。同样,为了避免轨迹协方差变化对于后续轨迹产生的影响,我们需要把单状态损失函数替换为

$$Q_{\mathbf{x}\mathbf{u},\mathbf{x}\mathbf{u}t} = \ell_{\mathbf{x}\mathbf{u},\mathbf{x}\mathbf{u}t} + 2f_{\mathbf{x}\mathbf{u}t}^{\mathrm{T}} \mathcal{L}_{\mathbf{S}t+1} f_{\mathbf{x}\mathbf{u}t}.$$

求导可得

$$\mathcal{L}_{\mathbf{A}t} = \frac{1}{2} Q_{\mathbf{u},\mathbf{u}t} + \frac{\lambda_t - 1}{2} \mathbf{A}_t^{-1} - \frac{\lambda_t}{2} \mathbf{A}_t^{-1} \mathbf{M} \mathbf{A}_t^{-1}$$

$$\mathcal{L}_{\mathbf{K}t} = Q_{\mathbf{u},\mathbf{u}t} \mathbf{K}_t \mathbf{S}_t + Q_{\mathbf{u},\mathbf{x}t} \mathbf{S}_t + \lambda_t \mathbf{A}_t^{-1} \mathbf{K}_t \mathbf{S}_t - \lambda_t \mathbf{A}_t^{-1} \mu_{\mathbf{x}t}^{\pi}(\hat{\mathbf{x}}_t) \mathbf{S}_t$$

$$\mathcal{L}_{\mathbf{S}t} = \frac{1}{2} \left[ Q_{\mathbf{x},\mathbf{x}t} + \mathbf{K}_t^{\mathrm{T}} Q_{\mathbf{u},\mathbf{x}t} + Q_{\mathbf{x},\mathbf{u}t} \mathbf{K}_t + \mathbf{K}_t^{\mathrm{T}} Q_{\mathbf{u},\mathbf{u}t} \mathbf{K}_t + \lambda_t (\mathbf{K}_t - \mu_{\mathbf{x}t}^{\pi}(\hat{\mathbf{x}}_t))^{\mathrm{T}} \mathbf{A}_t^{-1} (\mathbf{K}_t - \mu_{\mathbf{x}t}^{\pi}(\hat{\mathbf{x}}_t)) \right], \quad (6)$$
where  $\mathbf{M} = \Sigma_t^{\pi} + (\hat{\mathbf{u}}_t - \mu_t^{\pi}(\hat{\mathbf{x}}_t))(\hat{\mathbf{u}}_t - \mu_t^{\pi}(\hat{\mathbf{x}}_t))^{\mathrm{T}} + (\mathbf{K}_t - \mu_{\mathbf{x}t}^{\pi}(\hat{\mathbf{x}}_t))^{\mathrm{T}}$ 

零导数为零,可得

$$\mathbf{K}_{t} = -\left(Q_{\mathbf{u},\mathbf{u}t} + \lambda_{t}\mathbf{A}_{t}^{-1}\right)^{-1}\left(Q_{\mathbf{u},\mathbf{x}t} - \lambda_{t}\mathbf{A}_{t}^{-1}\mu_{\mathbf{x}t}^{\pi}(\hat{\mathbf{x}}_{t})\right). \tag{7}$$

$$\mathbf{A}_t Q_{\mathbf{u}, \mathbf{u}t} \mathbf{A}_t + (\lambda_t - 1) \mathbf{A}_t - \lambda_t \mathbf{M} = 0.$$
 (8)

由于 s, 、 A, 和 K, 相互依赖, 因此我们需要迭代多次直到它们数值收敛。

#### 4. 策略优化

有了更优的轨迹之后,我们接下来把轨迹固定,来优化相应的策略,这一步的优化相对来讲就比较直接了:写出关于损失函数关于策略参数。的表示,然后使用梯度方法做优化即可,比如SGD或者LBFGS。

$$\mathcal{L}(\theta) = \sum_{t=1}^{T} \lambda_t \sum_{i=1}^{N} D_{KL}(\pi_{\theta}(\mathbf{u}_t | \mathbf{x}_{ti}) || q(\mathbf{u}_t | \mathbf{x}_{ti}))$$

$$= \sum_{t=1}^{T} \lambda_t \sum_{i=1}^{N} \frac{1}{2} \left\{ tr(\Sigma_t^{\pi}(\mathbf{x}_{ti}) \mathbf{A}_t^{-1}) - \log |\Sigma^{\pi}(\mathbf{x}_{ti})| + (\mathbf{K}_t \mathbf{x}_{ti} + \mathbf{k}_t - \mu^{\pi}(\mathbf{x}_{ti}))^{\mathrm{T}} \mathbf{A}_t^{-1} (\mathbf{K}_t \mathbf{x}_{ti} + \mathbf{k}_t - \mu^{\pi}(\mathbf{x}_t))^{\mathrm{T}} \mathbf{A}_t^{-1} (\mathbf{K}_t \mathbf{x}_{ti} + \mathbf{k}_t - \mu^{\pi}(\mathbf{x}_t))^{\mathrm{T}} \mathbf{A}_t^{-1} (\mathbf{K}_t \mathbf{x}_{ti} + \mathbf{k}_t - \mu^{\pi}(\mathbf{x}_t))^{\mathrm{T}} \mathbf{A}_t^{-1} (\mathbf{K}_t \mathbf{x}_t)^{\mathrm{T}} \mathbf{A}_t^{-1} (\mathbf{K}_t \mathbf{x}_t)^{\mathrm{T}} \mathbf{A}_t^{-1} (\mathbf{K}_t \mathbf{x}_t)^{\mathrm{T}} \mathbf{A}_t^{-1} (\mathbf{K}_t \mathbf{x}_t)^{\mathrm{T}} \mathbf{A}_t^{-$$

算法:

# Algorithm 1 Constrained guided policy search

- 1: Initialize the trajectories  $\{q_1(\tau), \ldots, q_M(\tau)\}$
- 2: **for** iteration k = 1 to K **do**
- 3: Optimize each  $q_i(\tau)$  with respect to  $\mathcal{L}(\theta, q_i(\tau), \lambda_i)$
- 4: Optimize  $\theta$  with respect to  $\sum_{i=1}^{M} \mathcal{L}(\theta, q_i(\tau), \lambda_i)$
- 5: Update dual variables  $\lambda$  using Equation 2
- 6: end for
- 7: **return** optimized policy parameters  $\theta$  知乎 @张楚珩

# Algorithm 2 Trajectory optimization iteration

```
1: Compute f_{\mathbf{x}\mathbf{u}t}, \mu_{\mathbf{x}t}^{\pi}, \ell_{\mathbf{x}\mathbf{u}t}, \ell_{\mathbf{x}\mathbf{u},\mathbf{x}\mathbf{u}t} around \hat{\tau}
  2: for t = T to 1 do
          Compute \mathbf{K}_t and \mathbf{k}_t using Equations 3 and 4
  3:
          Compute \mathcal{L}_{\mathbf{x}t} and \mathcal{L}_{\mathbf{x},\mathbf{x}t} using Equation 5
  4:
  5: end for
  6: Initialize \alpha \leftarrow 1
  7: repeat
          Obtain new trajectory \hat{\tau}' using \mathbf{u}_t = \alpha \mathbf{k}_t + \mathbf{K}_t \mathbf{x}_t
  8:
 9:
          Decrease step size \alpha
10: until \mathcal{L}(\hat{\tau}', \mathbf{K}, \mathbf{A}) > \mathcal{L}(\hat{\tau}, \mathbf{K}, \mathbf{A})
11: Compute f_{\mathbf{xu}t}, \mu_{\mathbf{x}t}^{\pi}, \ell_{\mathbf{xu}t}, \ell_{\mathbf{xu},\mathbf{xu}t} around \hat{\tau}'
12: repeat
          Compute S_t using current A_t and K_t
13:
          for t = T to 1 do
14:
15:
              repeat
16:
                  Compute \mathbf{K}_t using Equation 7
                  Compute A_t by solving CARE in Equation 8
17:
18:
              until A_t and K_t converge (about 5 iterations)
              Compute \mathcal{L}_{St} using Equation 6
19:
20:
          end for
21: until all S_t and A_t converge (about 2 iterations)
22: return new mean \hat{\tau}' and covariance terms \Delta = K
```

其中算法1是主要的算法,算法2是其中第3行的具体算法。其中值得说明的是,第7-10行在最优开环控制上面缩小,以找到确定能改善损失函数的开环控制,其原因是由于环境是非线性的,而我们在局域假设了线性,因此如果开环控制变化的太多会导致假设不成立,形成的损失函数甚至不降反升。另外,第10行的大于符号应该是写错了,应该是小于符号,这里希望做的是最小化损失函数。

从绿框公式到红框公式的推导?

$$\begin{split} D_{KL}(q(\tau)||\rho(\tau)) &= \sum_{t=1}^{T} D_{KL}(q(x_{t}, u_{t})||\rho(x_{t}, u_{t})) \\ &= \sum_{t=1}^{T} \int q(x_{t}, u_{t})[\log q(x_{t}, u_{t}) - \log(\rho(x_{t}, u_{t}))] dx_{t} du_{t} \\ &= \sum_{t=1}^{T} \mathcal{H}(q(x_{t}, u_{t})) + \mathbb{E}_{q(x_{t}, u_{t})}[l(x_{t}, u_{t})] \end{split}$$

注意每一步的熵函数其实并不独立,后面的会依赖前面的,所以上面这样写其实不规范,不过我们 在运算的过程中记住这一点就好了,同时,高斯分布的熵函数为

```
\begin{split} \mathcal{H}(q(x_i, u_i)) &= \int q(x_i, u_i) \log q(x_i, u_i) dx_i du_i \\ &= \int q(x_i, u_i) \log q(x_i | x_{i-1}, u_{i-1}) dx_i du_i + \int q(x_i, u_i) \log q(u_i | x_i) dx_i du_i \\ &= \int q(x_i) \log q(x_i | x_{i-1}, u_{i-1}) dx_i + \int q(u_i) \log q(u_i | x_i) du_i \\ &= -1 - \log 2\pi - \frac{1}{2} |A_i| - \frac{1}{2} |F_i| \\ &\sim -\frac{1}{2} |A_i| \end{split}
```

其中随后一步是因为,求导的损失函数是要对于轨迹参数求导的,无关的常数项就可以不写了,有

$$\mathcal{L}(q) = \sum_{t=1}^{T} E_{q(\mathbf{x}_{t}, \mathbf{u}_{t})}[\ell(\mathbf{x}_{t}, \mathbf{u}_{t})] - \frac{1}{2} \log |\mathbf{A}_{t}| + \lambda_{t} E_{q(\mathbf{x}_{t})}[D_{\text{KL}}(\pi_{\theta}(\mathbf{u}_{t}|\mathbf{x}_{t}) || q(\mathbf{u}_{t}|\mathbf{x}_{t}))]$$
 ②张楚珩

注意到 [[ɾɨ]=テュ , 有

$$\begin{split} \mathbb{E}_{q(\mathbf{x}_{t}, \mathbf{u}_{t})}[l(x_{t}, u_{t})] &= \mathbb{E}_{q(\mathbf{x}_{t}, \mathbf{u}_{t})}[l(0, 0) + l_{out}^{T} \tau_{t} + \frac{1}{2} \tau_{t}^{T} l_{sumut} \tau_{t}] \\ &= l(0, 0) + l_{out}^{T} \dot{\tau}_{t} + \frac{1}{2} \mathbb{E}_{q(\mathbf{x}_{t}, \mathbf{u}_{t})}[(\tau_{t} - \dot{\tau}_{t} + \dot{\tau}_{t})^{T} l_{sumut}(\tau_{t} - \dot{\tau}_{t} + \dot{\tau}_{t})] \\ &= l(0, 0) + l_{out}^{T} \dot{\tau}_{t} + \frac{1}{2} \mathbb{E}_{q(\mathbf{x}_{t}, \mathbf{u}_{t})}[(\tau_{t} - \dot{\tau}_{t})^{T} l_{sumut}(\tau_{t} - \dot{\tau}_{t})] + \frac{1}{2} \mathbb{E}_{q(\mathbf{x}_{t}, \mathbf{u}_{t})}[(\tau_{t} - \dot{\tau}_{t})^{T} l_{sumut} \dot{\tau}_{t}] \\ &+ \frac{1}{2} \mathbb{E}_{q(\mathbf{x}_{t}, \mathbf{u}_{t})}[\dot{\tau}_{t}^{T} l_{sumut}(\tau_{t} - \dot{\tau}_{t})] + \frac{1}{2} \mathbb{E}_{q(\mathbf{x}_{t}, \mathbf{u}_{t})}[\dot{\tau}_{t}^{T} l_{sumut} \dot{\tau}_{t}] \\ &= l(0, 0) + l_{out}^{T} \dot{\tau}_{t} + \frac{1}{2} tr(\Sigma_{t} l_{sumut}) + 0 + 0 + \frac{1}{2} \dot{\tau}_{t}^{T} l_{sumut} \dot{\tau}_{t} \\ &\sim l_{out}^{T} \dot{\tau}_{t} + \frac{1}{2} tr(\Sigma_{t} l_{sumut}) + \frac{1}{2} \dot{\tau}_{t}^{T} l_{sumut} \dot{\tau}_{t} \end{split}$$

其中最后一个等号的成立可以参看这个课件的第六页。

同时注意到高斯函数的交叉熵

有

```
\begin{split} & \mathbb{E}_{q(\mathbf{e}_t)}[D_{KL}(\pi_{\theta}(u_t|x_t)||q(u_t|x_t))] \\ = & \mathbb{E}_{q(\mathbf{e}_t)}[\int \pi_{\theta}(u_t|x_t)\log \pi_{\theta}(u_t|x_t)du_t] - \mathbb{E}_{q(\mathbf{e}_t)}[\int \pi_{\theta}(u_t|x_t)\log q(u_t|x_t)du_t] \\ \sim & 0 + \frac{1}{2}tr(A_t^{-1}\Sigma_t^x) + \frac{1}{2}\log|A_t| + \frac{1}{2}\mathbb{E}_{q(\mathbf{e}_t)}[(u_t - \mu_t^x(x_t))^TA_t^{-1}(u_t - \mu_t^x(x_t))] \\ = & \frac{1}{2}tr(A_t^{-1}\Sigma_t^x) + \frac{1}{2}\log|A_t| + tr(S_t(K_t - \mu_{xt}^x)^TA_t^{-1}(K_t - \mu_{xt}^x)) + (\hat{u}_t - \mu_x(\hat{x}_t))^TA_t^{-1}(\hat{u}_t - \mu_x(\hat{x}_t)) \end{split}
```

其中约等于符号是因为前一项是完全关于策略的,与我们这里要优化的轨迹无关,因此不理会这一项,其中比较复杂的最后一项细写下来是

```
\begin{split} &\mathbb{E}_{q(\mathbf{a})}[(K_{t}x_{t} - \mu_{t}^{x}(x_{t}))^{T}A_{t}^{-1}(K_{t}x_{t} - \mu_{t}^{x}(a_{t}))] \\ &= &\mathbb{E}_{q(\mathbf{a})}[(K_{t}x_{t} - K_{t}\hat{x}_{t} + \hat{a}_{t} - \mu_{\tau}(\hat{a}_{t}) + \mu_{\tau}(\hat{a}_{t}) - \mu_{\tau}(x_{t}))^{T}A_{t}^{-1}(\cdots)] \\ &= &\mathbb{E}_{q(\mathbf{a})}[((K_{t} - \mu_{xt}^{x})(x_{t} - \hat{x}_{t}) + (\hat{a}_{t} - \mu_{\tau}(\hat{a}_{t})))^{T}A_{t}^{-1}(\cdots)] \\ &= &\mathbb{E}_{q(\mathbf{a})}[((K_{t} - \mu_{xt}^{x})(x_{t} - \hat{x}_{t})^{T}A_{t}^{-1}(\cdots)] + 0 + 0 + \mathbb{E}_{q(\mathbf{a})}[(\hat{a}_{t} - \mu_{\tau}(\hat{a}_{t}))^{T}A_{t}^{-1}(\cdots)] \\ &= &tr(S_{t}(K_{t} - \mu_{xt}^{x})^{T}A_{t}^{-1}(K_{t} - \mu_{xt}^{x})) + (\hat{a}_{t} - \mu_{x}(\hat{a}_{t}))^{T}A_{t}^{-1}(\hat{a}_{t} - \mu_{x}(\hat{a}_{t})) \end{split}
```

其中轨迹的控制和策略的控制都有直流项的,但是由于后面求协方差并不影响,为了简便,这里就 不写直流项了。

最后拼起来有

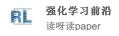
$$\mathcal{L}(q) \approx \sum_{t=1}^{T} \frac{1}{2} \begin{bmatrix} \hat{\mathbf{x}}_{t} \\ \hat{\mathbf{u}}_{t} \end{bmatrix}^{T} \ell_{\mathbf{xu},\mathbf{xu}t} \begin{bmatrix} \hat{\mathbf{x}}_{t} \\ \hat{\mathbf{u}}_{t} \end{bmatrix} + \begin{bmatrix} \hat{\mathbf{x}}_{t} \\ \hat{\mathbf{u}}_{t} \end{bmatrix}^{T} \ell_{\mathbf{xu}t} + \frac{1}{2} \operatorname{tr} \left( \sum_{t} \ell_{\mathbf{xu},\mathbf{xu}t} \right) - \frac{1}{2} \log |\mathbf{A}_{t}| + \frac{\lambda_{t}}{2} \log |\mathbf{A}_{t}| + \frac{\lambda_{t}}{2} \log |\mathbf{A}_{t}| + \frac{\lambda_{t}}{2} (\hat{\mathbf{u}}_{t} - \mu_{t}^{\pi}(\hat{\mathbf{x}}_{t}))^{T} \mathbf{A}_{t}^{-1} (\hat{\mathbf{u}}_{t} - \mu_{t}^{\pi}(\hat{\mathbf{x}}_{t})) + \frac{\lambda_{t}}{2} \operatorname{tr} \left( \mathbf{A}_{t}^{-1} \sum_{t}^{\pi} \right) + \frac{\lambda_{t}}{2} \operatorname{tr} \left( \mathbf{S}_{t} \left( \mathbf{K}_{t} - \mu_{\mathbf{x}t}^{\pi}(\hat{\mathbf{x}}_{t}) \right)^{T} \mathbf{A}_{t}^{-1} \left( \mathbf{K}_{t} - \mu_{\mathbf{x}t}^{\pi}(\hat{\mathbf{x}}_{t}) \right) \right)$$

## 推导过程参考了 【强化学习大讲坛】强化学习进阶 第九讲 引导策略搜索。

发布于 2018-10-01

机器学习 算法 强化学习 (Reinforcement Learning)

## 文章被以下专栏收录



进入专栏