

# Learning Latent Dynamics for Planning from Pixels

Danijar Hafner<sup>1 2</sup> Timothy Lillicrap<sup>3</sup> Ian Fischer<sup>4</sup> Ruben Villegas<sup>1 5</sup>  
David Ha<sup>1</sup> Honglak Lee<sup>1 5</sup> James Davidson<sup>1</sup>

## 【强化学习 40】PlaNet



张楚琦

清华大学 交叉信息院博士在读

25 人赞同了该文章

Google Brain和DeepMind合作的一个model-based强化学习算法——Deep Planning Network (PlaNet)。

### 原文传送门

Hafner, Danijar, et al. "Learning Latent Dynamics for Planning from Pixels." arXiv preprint arXiv:1811.04551 (2018).

### 特色

一种基于模型（model-based）的强化学习算法，能够直接图像作为输入来控制机器人；规划（planning）的过程在隐空间（latent space）上进行，速度较快；采取了一些措施避免了长距离规划中状态估计不准确的问题。

### 过程

#### 1. 总体流程

既然是基于模型的方法，那么其两大问题主要就是如何学习模型和如何利用模型来做规划。

首先，模型包含些什么？

这里考虑一个基于图像的输入，这是一个高维度的输入，因此我们认为这是一个POMDP。图像输入是观测量  $o_t$ ，状态是一个隐变量  $s_t$ 。基于模型的方法中的“模型”指的是transition model  $p(s_t|s_{t-1}, a_{t-1})$  和 reward model  $p(r_t|s_t)$ ，POMDP还带来了observation model  $p(o_t|s_t)$ 。另外我们还需要一个encoder，快速地利用可观测量  $o_{G1}, a_{G1}$  来估计隐变量  $s_t$ ， $p(s_t|o_{G1}, a_{G1})$ 。

总结一下，要学习以下四个模型的参数 transition model, reward model, observation model, encoder。

规划的流程如下，使用学习到的模型来针对不同的行动预测收益，并选取最好的下一步的行动。

$o_{G1}, a_{G1} \rightarrow \text{encoder} \rightarrow s_t \rightarrow \text{transition model} \rightarrow s_{t+1}, s_{t+2}, \dots \rightarrow \text{reward model} \rightarrow R \rightarrow \text{CEM} \rightarrow a_t$

模型学习的流程如下

$o_{1:T}, o_{1:T} \rightarrow \text{SGD on variational bound} \rightarrow \text{observation and transition model, encoder}$

$o_{1:T}, o_{1:T} \rightarrow \text{encoder} \rightarrow s_{1:T} (+r_{1:T}) \rightarrow \text{reward model}$

整体算法如下

---

### Algorithm 1: Deep Planning Network (PlaNet)

---

**Input :**

$R$	Action repeat	$p(s_t   s_{t-1}, a_{t-1})$	Transition model
$S$	Seed episodes	$p(o_t   s_t)$	Observation model
$C$	Collect interval	$p(r_t   s_t)$	Reward model
$B$	Batch size	$q(s_t   o_{\leq t}, a_{< t})$	Encoder
$L$	Chunk length	$p(\epsilon)$	Exploration noise
$\alpha$	Learning rate		

```

1 Initialize dataset  $\mathcal{D}$  with  $S$  random seed episodes.
2 Initialize model parameters  $\theta$  randomly.
3 while not converged do
    // Model fitting
4   for update step  $s = 1..C$  do
5     Draw sequence chunks  $\{(o_t, a_t, r_t)_{t=k}^{L+k}\}_{i=1}^B \sim \mathcal{D}$ 
        uniformly at random from the dataset.
6     Compute loss  $\mathcal{L}(\theta)$  from Equation 7.
7     Update model parameters  $\theta \leftarrow \theta - \alpha \nabla_{\theta} \mathcal{L}(\theta)$ .

    // Data collection
8    $o_1 \leftarrow \text{env.reset}()$ 
9   for time step  $t = 1..\lceil \frac{T}{R} \rceil$  do
10    Infer belief over current state  $q(s_t | o_{\leq t}, a_{< t})$  from
        the history.
11     $a_t \leftarrow \text{planner}(q(s_t | o_{\leq t}, a_{< t}), p)$ , see
        Algorithm 2 in the appendix for details.
12    Add exploration noise  $\epsilon \sim p(\epsilon)$  to the action.
13    for action repeat  $k = 1..R$  do
14       $r_t^k, o_{t+1}^k \leftarrow \text{env.step}(a_t)$ 
15     $r_t, o_{t+1} \leftarrow \sum_{k=1}^R r_t^k, o_{t+1}^k$ 
16   $\mathcal{D} \leftarrow \mathcal{D} \cup \{(o_t, a_t, r_t)_{t=1}^T\}$ 

```

---

知乎 @张楚珩

其中第4~7行表示的是模型训练的部分，第11行表示的是模型的规划部分，第8~16行表示的模型的规划连同行动和采样的部分。

## 2. 如何利用模型来做规划？

先假设我们前面提到的各个模型的参数已经学习到了，那么如何来利用这些模型来做规划呢？这里使用了一个最为常用的框架，叫做model-predictive control（MPC）。简单说来就是每一步都往后面模拟很多步，然后选取一个收益最高的方案，然后只采取一步的行动；等到下一步的时候再重新规划并且选取新的一步；即并不是规划出一串行动之后开环地做着一串，而是规划一串只走一步。

回顾刚刚提到的流程框架

$$o_{s1}, a_{c1} \rightarrow \text{encoder} \rightarrow s_t \rightarrow \text{transition model} \rightarrow s_{t+1}, s_{t+2}, \dots \rightarrow \text{reward model} \rightarrow R \rightarrow \text{CEM} \rightarrow a_t$$

我们先基于先前的观察使用encoder得到当前状态的估计，然后对于未来的行动序列进行采样，并结合transition model来预测未来的状态序列，同时使用reward model来估计获得奖励的期望值。这样对于一个给定的未来行动序列就能够得到一个期望收益了，把这一套当做一个黑盒子我们就能利用CEM来规划出一个好的未来行动序列了。取这个行动序列的第一个行动作为下一步我们选择的行动即可。具体算法如下。

---

**Algorithm 2:** Latent planning with CEM

---

**Input :**  $H$  Planning horizon distance       $q(s_t | o_{\leq t}, a_{< t})$  Current state belief  
          $I$  Optimization iterations       $p(s_t | s_{t-1}, a_{t-1})$  Transition model  
          $J$  Candidates per iteration       $p(r_t | s_t)$  Reward model  
          $K$  Number of top candidates to fit

```
1 Initialize factorized belief over action sequences  $q(a_{t:t+H}) \leftarrow \text{Normal}(0, \mathbb{I})$ .
2 for optimization iteration  $i = 1..I$  do
3     // Evaluate  $J$  action sequences from the current belief.
4     for candidate action sequence  $j = 1..J$  do
5          $a_{t:t+H}^{(j)} \sim q(a_{t:t+H})$ 
6          $s_{t:t+H+1}^{(j)} \sim q(s_t | o_{1:t}, a_{1:t-1}) \prod_{\tau=t+1}^{t+H+1} p(s_\tau | s_{\tau-1}, a_{\tau-1}^{(j)})$ 
7          $R^{(j)} = \sum_{\tau=t+1}^{t+H+1} E[p(r_\tau | s_\tau^{(j)})]$ 
8         // Re-fit belief to the  $K$  best action sequences.
9          $\mathcal{K} \leftarrow \text{argsort}(\{R^{(j)}\}_{j=1}^J)_{1:K}$ 
10         $\mu_{t:t+H} = \frac{1}{K} \sum_{k \in \mathcal{K}} a_{t:t+H}^{(k)}$ ,  $\sigma_{t:t+H} = \frac{1}{K-1} \sum_{k \in \mathcal{K}} |a_{t:t+H}^{(k)} - \mu_{t:t+H}|$ .
11         $q(a_{t:t+H}) \leftarrow \text{Normal}(\mu_{t:t+H}, \sigma_{t:t+H}^2 \mathbb{I})$ 
12 return first action mean  $\mu_t$ .
```

---

知乎 @张楚珩

### 3. 如何学习模型？

通过采样得到的数据是观察到的图像  $o_{1:T}$  和  $a_{1:T}$ ，因此我们需要同时训练多个模型使得出现该数据的概率最大（maximum log-likelihood）。

**One-step predictive distribution** The variational bound for latent dynamics models  $p(o_{1:T}, s_{1:T} | a_{1:T}) = \prod_t p(s_t | s_{t-1}, a_{t-1}) p(o_t | s_t)$  and a variational posterior  $q(s_{1:T} | o_{1:T}, a_{1:T}) = \prod_t q(s_t | o_{\leq t}, a_{< t})$  follows from importance weighting and Jensen's inequality as shown,

$$\begin{aligned} \ln p(o_{1:T} | a_{1:T}) &\triangleq \ln E_{p(s_{1:T} | a_{1:T})} \left[ \prod_{t=1}^T p(o_t | s_t) \right] \\ &= \ln E_{q(s_{1:T} | o_{1:T}, a_{1:T})} \left[ \prod_{t=1}^T p(o_t | s_t) p(s_t | s_{t-1}, a_{t-1}) / q(s_t | o_{\leq t}, a_{< t}) \right] \\ &\geq E_{q(s_{1:T} | o_{1:T}, a_{1:T})} \left[ \sum_{t=1}^T \ln p(o_t | s_t) + \ln p(s_t | s_{t-1}, a_{t-1}) - \ln q(s_t | o_{\leq t}, a_{< t}) \right] \\ &= \sum_{t=1}^T \left( \underbrace{E[\ln p(o_t | s_t)]}_{\text{reconstruction}} - E \left[ \underbrace{\text{KL}[q(s_t | o_{\leq t}, a_{< t}) \parallel p(s_t | s_{t-1}, a_{t-1})]}_{\text{complexity}} \right] \right). \end{aligned} \quad (8)$$

知乎 @张楚珩

这里完全是按照ELBO（evidence lower bound）的推导。注意，不等式的右边包含encoder

$q(s_t | o_{\leq t}, a_{< t})$ ，observation model  $p(o_t | s_t)$  和 transition model  $p(s_t | s_{t-1}, a_{t-1})$ 。通过对于右式做随机梯度上升即可学到这三个模型。

Reward model的学习就比较简单了，类似于普通的有监督学习。

至此，大致的框架已经讲完了，但是本文的创新点是提出了以下两种技术，使得planning过程中状态的估计更为准确。

#### 4. Recurrent State Space Model (RSSM)

前面我们提到状态的转移模型 transition model 是表示为  $a_t \sim p(a_t | s_{t-1}, a_{t-1})$ ，文中都是使用的diagonal Gaussian distribution，即由神经网络生成下一步状态的均值和方差，并从高斯模型中进行采样。如果每一个单步的transition model都是完美的，那么逐步估计也会比较准确，但是单步的模型并不完美，它受限于不充足的训练以及较弱的模型表示能力。每一步都完全是随机采样，在经过很多步之后，前面信息就会丢失掉。因此，这里作者还引入了确定性（非随机）的部分。（当然了，如果神经网络自己学出来某一部分的方差就为0，那么就等效于加上确定性部分了，不过不能期望神经网络能够自己学出来）

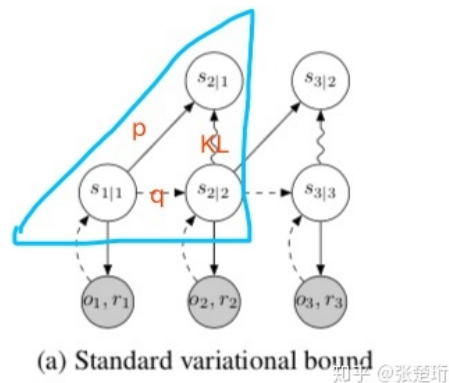
通过在状态隐变量中加上确定性的部分，就得到了RSSM。transition model变为了

$h_t = f(h_{t-1}, a_{t-1}, a_{t-1})$ ， $a_t \sim p(a_t | h_t)$ ，即原来的  $a_t$  拆成了  $(h_t, a_t)$  两个部分。相应地，encoder也进行了一些改变， $q(a_t | o_{1:t}, a_{1:t}) \rightarrow q(a_t | o_{1:t}, a_{1:t}) = \prod_i q(a_i | h_i, a_i) = \prod_i q(a_i | f(h_{i-1}, a_{i-1}, a_{i-1}), a_i)$ 。

这样的改变只是改变了网络的结构，训练的目标和方式没有大的变化。

#### 5. Latent Overshooting

前面推导到的下界（lower bound）是优化的目标，transition model的学习主要是和encoder  $q(a_t | o_{1:t}, a_{1:t})$  的比较得来的，它们之间的比较是一步一步地比较。如下图所示，蓝色框中的虚线代表encoder一步步的推断，实线代表transition model的生成过程，波浪线代表需要对两者生成的  $a_t$  分布进行比较，并且最小化它们之间的KL散度。

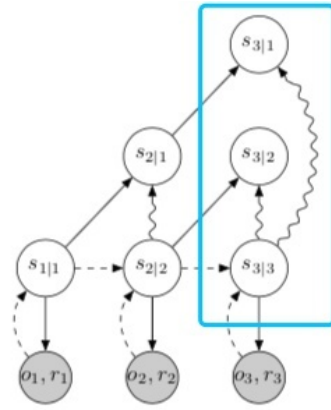


图：原变分下界对应的模型。虚线代表推断，实线代表生成过程，波浪线代表KL散度约束。

可以观察到，在这样的做法中，transition model始终是一步步分离开的，并没有像之后使用的时候那样多步串在一起去预测一个未来较长时间内的状态分布。少了多步串联起来的训练，得到的transition model在多步串联的时候会产生较大的误差，因此文章采用了以下latent overshooting的方法来改善。

Latent overshooting的做法如下图所示。以  $a_2$  为例，不仅仅是一步转移产生的状态  $a_{2|1}$  要和encoder产生的状态  $a_{2|2}$  作比较并减少差距，而且多步转移（这里是两步）产生的状态  $a_{2|3}$  也要和  $a_{2|2}$  作比较并最小化KL散度。这相当于不仅要求transition model在走一步的情况下要和encoder对的上，而且

要求它在串联工作的时候也要和encoder对的上。



(c) Latent overshooting

图：Latent overshooting对应的模型。虚线代表推断，实线代表生成过程，波浪线代表KL散度约束。

考虑了latent overshooting之后的变分下界写作

$$\frac{1}{D} \sum_{d=1}^D \ln p_d(o_{1:T}) \geq \sum_{t=1}^T \left( \underbrace{\frac{\mathbb{E}_{q(s_t | o_{\leq t})} [\ln p(o_t | s_t)]}{\text{reconstruction}}}_{\text{latent overshooting}} - \frac{1}{D} \sum_{d=1}^D \beta_d \mathbb{E} \left[ \text{KL}[q(s_t | o_{\leq t}) \parallel p(s_t | s_{t-1})] \right] \right). \quad (7)$$

其中， $D$  表示transition model最多串联的步数，可以看到，KL散度项被替换为了多对约束的加权平均，权重通过  $\beta_d$  来调节。

## 实验结果

这个工作跑了如下六组实验，这六组实验都有各自的特点，比如有些具有比较复杂的动力学（与环境有接触）、有些只能观察到部分的信息、有些具有稀疏的奖励结构。同时，该工作着重说明了RSSM和Latent overshooting在其中发挥的作用。

此外，该工作还训练了一个模型，使得该模型能够在所有的环境中运行。

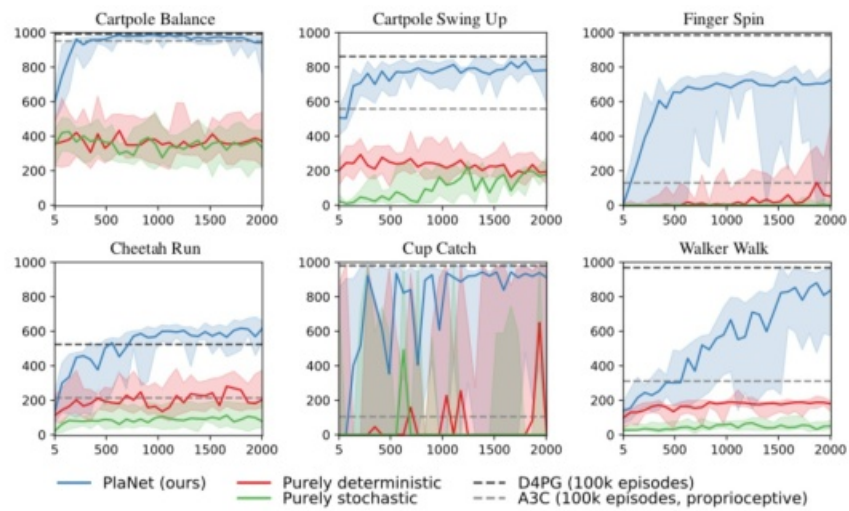


Figure 4: Comparison of PlaNet to model-free algorithms and other model designs. Plots show test performance for the number of collected episodes. We compare PlaNet using our RSSM (Section 3) to purely deterministic (RNN) and purely stochastic models (SSM). The RNN does not use latent overshooting, as it does not have stochastic latents. The plots show medians and the areas show percentiles 5 to 95 over 4 seeds and 10 rollouts.

文章来源：导师推荐

发布于 2019-02-23

强化学习 (Reinforcement Learning)

赞同 25

1 条评论

分享

喜欢

收藏

...

文章被以下专栏收录



强化学习前沿  
读呀读paper

进入专栏