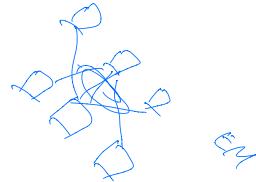
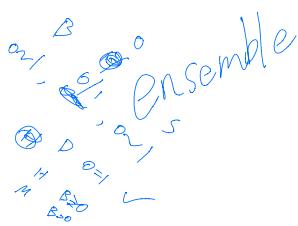


$$\pi_{\text{P}} \leftarrow \pi_m \left( \pi_{\text{S}} \left( \pi_{\text{S}} \left( M, H \right), D \right) \right)$$

$\pi_{\text{P}} \leftarrow \pi_m$   
capable net layer.



与 Hippo 不同：① 输入不同频率数据。

D

H

M

② 频率 subpolicy, subsubpolicy etc.

③ disentangle

④ Option

termination  $\beta$ .  $\beta$  是 sub-policy 学习。  
还是 option

⑤ Each level has own layers of period

D: Null

M: 240/day

H: 4/day

⑥ constraints on action space

D/W instead of LLEP

MRFLESSVM.

$$\mathcal{L}^H = D/W \quad \mathcal{L}^U \mathcal{L}^P \mathcal{L}^H$$

$$\Downarrow$$

$$\oplus \pi^M \pi^H \pi^P$$

Enrich volatility (duration)

Steven Energy Periodic.

$$(5) \text{ visiting prob?} \quad \pi(a|s) = \frac{\pi_\theta}{\pi_{old}}$$

(4) What happens when  $\boxed{done}$ ?

How bound converted in distribution?

L & V Bound.

① action space continuous: could action space be discrete?

② actor.model: roughly

critic.model: discrete

3 so that one model?

Yes. They are re-used variables.

③ some steps: sample step?

which is later steps?

④ action & observation space:

obspace: determine input encode

discrete: one-hot + float

box: float

box could be normalized

policies/policy-build-policy.c

line 129-139

policy.py

⑤ PolicyWithValue.

-init-

self.state = constant

self.initial\_state = None.

step

-evaluate [self.action, self.rv, self.state, self.done]

⑥ run episode:

add a breakpoint to train's forward.

see when it's being called.

build\_policy: policy-learn.

common/policies.py

PolicyWithValue.\_\_init\_\_( ) - - :

line 39: self.\_\_dict\_\_.update(tensors)

extra\_tensors

common/distributions:

make\_pdtype() → distribution()

Categorical: categorical

box: Ding GaussianType

policy\_fn

extra\_tensors } rms 135 normalized things

recurrent\_tensors [5]

for (t in: [ initial\_state

state

(6) epoch? minibatch?

(7) value.net 'copy' 'share'

build\_policy → policy\_fn ← ,

??? 165: todo if not support recurrent net.

sample Advantage

$$L^{VP} = V_\theta - V^{\text{target}}$$

||?

496.  $V_\theta$  - sampled-return.

$$495. \text{ sampled\_return} = V_{\theta, \text{old}} + A_{\theta, \text{old}}$$

$$695, 698 = r_t + \gamma V_{\theta, \text{old}} + \gamma \lambda A_{\theta, \text{old}}$$

return.

(8) How  $V^{\text{target}}$  (return)

ppo2/runner.py line 65

mb\_returns = mb\_advantages + mb\_values.

$$\text{return}_t = V_{\theta, \text{old}} + \gamma V(C_{S_{t+1}}) + \gamma \delta_{t+1} + \dots + (\gamma^T) \delta_T$$

$\delta_t$

Settings:

```

if value_network:
    - None = 'share' (default)
    - 'copy'
    - 'callbk'

```

by default share policy-net's value-network.  
`build_policy()`: policy\_fn.

Notes:

policy has step() & value()  
 $\uparrow$   
 $\uparrow$  self if  
 $\uparrow$  policy\_fn has step & value  
 $\uparrow$  X, value\_network: callable  
 $\uparrow$  nenv is calculated by `batch_policy(X)`  
 $\uparrow$  action, `action = policy_fn(...).sample()`, `policy_fn`  
 $\uparrow$  value is calculated by `value_fn(X)`  
 $\uparrow$  `value_fn = value_network(X)`  
 $\uparrow$  `value_fn = 'copy'` or `'callbk'` (S3-16)  
 $\uparrow$  `value_fn = 'share'`  
 $\uparrow$  `value_fn = policy_fn(X)`  
 $\uparrow$  `value_fn = 'share'`  
 $\uparrow$  `value_fn = 'share'`  
 $\uparrow$  because `of.value = vnet(encoded_X)`.  
 $\uparrow$  `v.net` g copy, item needs 'non-' paren  
 $\uparrow$  callable, custom fn need to handle  
 $\uparrow$  nenv or nsteps  
 $\uparrow$  `nbatch = X.shape[0]`

How does build\_policy handles nbatch?  
 $nbatch = nenv * nsteps$  Work?  
 $nbatch = nenv * nsteps = 6 \times 20 = 120$

models.py  
 $M_t$  = placeholder  
 $m_t, x_t$   
 $initial\_state$   
 $nbatch\_train > nsteps$        $nbatch > nsteps$

$x_s, m_s$  [list [ $x_{t_1}, \dots, x_T$ ]  
 $[M_{t_1}, \dots M_T]$ ]

④. mask: True/False  
 $\downarrow$   
 Done

(1) In LSTM `a2c/utils.py/lstm()`:

cell state  $t-1 \times (1 - \text{mask}_t)$   
 hidden state  $t-1$

(2) In `advs/returns/ppo2/runner.py/run()`

next nonterminal =  $1 - \text{self}.dones$ . `np.newaxis`  
 $[1, nenv]$

Sample: Model, act-model differs from train-model  
 [env, ntrn] self.dones [1, num] in batch size & nsteps  
 model.step(X, States, M) [b1, -] parameters are reused.  
 [envs, obs] For sampling, nsteps always be 1.

① sample  $a \sim \pi(a|x)$

i. build  $p(a|x, \mu, \sigma)$   
 with policy\_fn  
 $\mu = fc_x(\text{lstm}(\text{preproc}(X), s, M))$   
 $\sigma = tf.given$  : 有向图 X : 等待输入？不确定？ to be confirmed.

② estimate value  $V_{\theta}$

$v_f(S_t)$   
 $= fc_v(\text{lstm}(\text{preproc}(X), s, M))$  if copy  
 $= fc_v(\text{custom\_value\_net}(X, S, M))$

③ recurrent state: LSTM's state, not env's state

PolicyWithValue.\_\_init\_\_( ) :  
 line 39: self.\_\_dict\_\_.update(tensors)  
 extra\_tensors  
 policy\_fn  
 extra\_tensors } rms 135 normalized things  
 recurrent\_tensors | 5 |  
 for lstm: { initial\_state  
 state : concat [ h\_t, C\_t ]

Train:

$$\text{advs} = \frac{\hat{a} - \bar{a}}{\text{std}(a) + 1e-8}$$

## Stateful LSTM

Stateful: states from previous batch will be used  
 batch\_size batch\_size in batch working index of samples



Qs: 1. How big frame batch (batch, timestep, feat\_dim)  
 batch\_size size work in here

2. LSTM size does not match input time\_steps  
 its cell size, not timesteps

## KDD 20 Map

Input { stable baseline  
Hippo  
dev by Sunday

Hippo Paper  
GAE Book & Blogs.  
↓  
Option PPO + advantage

Stable baselines.

Init.

Model

Policy

Base RL Model

get policy

check env. { num  
action, input space  
VecEnv If.

ActorCriticModel

Setup\_ModelC >: nenv, n\_steps, batch\_size  
act\_model = policy(nenv, 1, nenv)  
train\_model = policy( $\frac{nenv}{minibatch}$ , n\_steps,  $\frac{nenv * nsteps}{minibatch}$ )

policies.py  
 q-value (D) policies 420  
 distributions 242. q-value shape?  
 value\_fn policies 428  
 n\_actions  
 q-value: linear(nn.output, 'f', self.size)  
 self.value\_fn: linear(nn.output, 'f', 1)  
 self.neglogp: pd.neglogp(self.action)  
 self.value\_fn: self.value\_fn[:, 0]  
 action: pd.sample  
 deterministic: pd.mode  
 policy\_proba: Categorical  
 Diag Gaussian: [mean, std]

distributions.py

self.action: pd.sample >  
 self.neglogp: self.neglogp(action)  
 self.

# Hippo

policy param.

$\pi(\theta)$ :  
path:  $\rightarrow s_t, r_t, a_t, \log p_t$   
 $a_t = \text{sample} >$

$s_t, r_t = Env(a_t)$   
 $\pi_{\text{roll}}(\theta)$ :  
[rlab/sampler/utils.py/rollouts](#)

$$\delta_t = r_t + \gamma \delta_{t+1} \quad V_t = 9.876543210$$

$$\delta_1 = r_1 + \gamma$$

A: Base Sampler / process\\_samples  
A manager

$$\delta_2 = r_2 + 0.5 \times 9 = 12.5$$

Hier Batch Sampler / process\\_samples:  
A skill

$$\delta_3 = 7 + 0.5 \times 12.5 = 13.25$$

$$6.25$$

$$\delta_4 = 6 + 0.5 \times 13.25 = 12.625$$

$$6.625$$

hier\_batch\_sampler.py  
skill-dependent baseline  $V_f$ -skill

```

# add noise
obs = [T, 1] M action space
extended_obs [T, N+1]
latents [T, M]

new = obs [T, N+1] * latents [T, 1, M]
new = new.reshape [T, (N+1)*M]
j=1...M
new_j = obs * [latents[:, j]]
new = concat ([new, ..., new_M])

M+N+1 |-----| contact scalar
       j|M
       new_j = obs * [latent[:, j]]



|     |     |     |     |              |     |        |
|-----|-----|-----|-----|--------------|-----|--------|
| t=1 | --- | N+1 | --- | (N+1)(M-1)+1 | ... | (N+1)M |
| t=2 | --- | --- | --- | ---          | --- | ---    |
| t=T | --- | N+1 | --- | ---          | --- | M      |


now = concat [new, extended_obs, latents]
(N+1)XM   N+1   M.

now [T, (N+1)XM + (N+1) + M]

```

$$V_f \text{ skill} = V_f \text{ skill}(new)$$

## H PPO

$\max_{\pi^m}$  over  $(obs_t, z_t, a_t)_{t=1:T}$   
 $\pi_{\text{old}}^{\text{full}}(obs_t) = \pi_{\text{old}}^m \theta_{\text{old}}^m$   
 $\pi_{\text{old}}^{\text{full}}(obs_t, z_t) = \pi_{\text{old}}^z \theta_{\text{old}}^z$   
 T and overlaps  
 H merger has been called H times  $H \in \{t_1, t_2, t_3, \dots, t_T\}$   
 P current skill total period  $P \sim U(\text{period\_min}, \text{period\_max})$   
 $p_t$  remaining time  $p_t = \frac{P - t}{\text{max\_period}}$

$$\begin{aligned}
 \text{Loss} &= \frac{E[\dots H]}{L_{\text{mean\_period}}} + \frac{E[\dots T]}{L_{\text{skill}}} \\
 L_a &= E\left[\sum_{t=1}^T A_t^m(z_t, s_t)\right] + E[H] \\
 &= E\left[\exp(\log \pi_{\text{old}}^m(z_t, s_t) - \log \pi^m(z_t, s_t)) A_t^m(z_t, s_t, a_t)\right] + E[H] \\
 \pi_{\text{old}}^m(z_t | s_t) &\leftarrow \text{tex} \\
 s_t &= \text{concat}(obs_t, p_t) \quad \text{neural\_net.} \\
 z_t &\sim \pi^m(z | s_t)
 \end{aligned}$$

$$\begin{aligned}
 \pi_{\text{old}}^s(z_t | s_t, a_t) &\leftarrow \text{tex} \\
 s_t &= \text{concat}(obs_t, p_t) \\
 a_t &\sim \pi_{\text{old}}^s(z_t | s_t, a_t)
 \end{aligned}$$

$$\begin{aligned}
 A_t^m(z) &\leftarrow \text{tex} \\
 \text{for } t=1:T \\
 v_t^m &= V_{\text{old}}^m(obs_t) \\
 \delta_t^m &= r_t^m + \gamma \text{gamma} * A_{t+1}^m \\
 R_t^m &= r_t^m + \text{gamma} * \delta_t^m \\
 A_t^s(z) &\leftarrow \text{tex} \\
 \text{for } t=1:T \\
 s_t &= \text{concat}[obs_t, p_t] \\
 v_t^s &= V_{\text{old}}^s(\text{concat}[\text{concat}[s_t, \text{new\_form}_m(z_t)], s_t, z_t]) \\
 \delta_t^s &= r_t^s + \gamma \text{gamma} * v_{t+1}^s - v_t^s \\
 A_t^s &= \delta_t^s + \text{argmax} A_{t+1}^s \\
 R_t^s &= R_t^m
 \end{aligned}$$

$V_{\text{old}}^m \rightarrow V_{\text{old}}^s$   
 $\text{obj min } (V_{\text{old}}^s V_{\text{old}}^m - R_t)^2$

In ppo with  $(V - V_{target})^2$   
↑  
changed every policy gradient iteration

Both hippo, ppo,

action : unchanged.

A : unchanged.

old\_vpred\_ph: unchanged . only used for clip-vf.

Implement normalize: obs, return, advantages.

② for recurrent,  $\beta$ , time\_remaining is not needed <sup>for</sup> low\_policy hier-base.py, hier-batch-sampler.py.  
③ vf-manager, vf-skill fitted during sample period (process-samples > )

should fit expected return

rather than reward directly.

# Hier Options.

Markov Process,

Semi-MP

Off policy ?? Sample Efficiency

Mult fractal Theory?

value, policy as GANV

GMM Graph HMM.

station 2000  $\frac{\partial Q}{\partial \theta}$  A? forget, forgetvar

[BS] 92 PBR<sub>3</sub>

GAE 可参数化 2000, 92?

O

## TRPO

$$\eta(\pi) = E_{s_0, a_0, \dots} \left[ \sum_{t=0}^{\infty} r^t r(s_t) \right]$$

$$\eta(\pi_{new}) = \eta(\pi_{old}) + \sum_s \underbrace{Q(s)}_{\pi_{new}} \sum_a T_{new}(a|s) A_{\pi}(s, a)$$

Appendix A <sup>lamma1.</sup>  
22, 23, 24

$$t=0: r_h - V_0 \quad -A(\pi) = -E_{s_0} [V_{\pi}(s_0)]$$

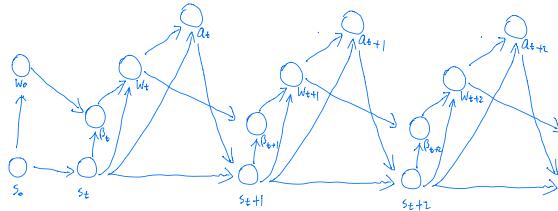
$$t=1: r_h - V_1$$

$$t=2: r_h - V_2$$

$$\vdots$$

$$k-1: V_k$$

$$\begin{aligned}
P(S_{t+1}, W_{t+1} | S_t, W_t) &= \sum_{\alpha} \pi_{w|\alpha}(a_t | s_t) P(S_{t+1} | s_t, a_t) \\
&\quad \underbrace{[ (1 - \beta_{w|\alpha}(s_{t+1})) \prod_{w \in \mathcal{W}} \pi_{w|\alpha} + \beta_{w|\alpha} \pi_{S_{t+1}} \pi_{w|\alpha}(s_{t+1}) ]}_{\Downarrow \beta_w \in \{0, 1\}} \\
&\quad \text{For option } w \\
V1: \quad \sum_{\beta_w} \frac{\sum_{\alpha} P(a_t | s_t, \alpha) P(w_{t+1} | s_t, \alpha)}{P(w_{t+1} | s_t, \alpha) + P(w_{t+1} | s_t, \alpha)^2} \\
&\quad \Downarrow \\
&\quad \sum_{\beta_w} [\mu_{\alpha}(s_{t+1}) P(w_{t+1} | s_{t+1})]^{\beta_w} \\
&\quad \text{For option } w \\
V2: \quad \sum_{\beta_w} \frac{\sum_{\alpha} P(a_t | s_t, \alpha) P(w_{t+1} | s_t, \alpha)}{P(w_{t+1} | s_t, \alpha) + P(w_{t+1} | s_t, \alpha)^2} \\
&\quad \Downarrow \\
&\quad \sum_{\beta_w} \left[ \prod_{w \in \mathcal{W}} \mu_{\alpha}(s_{t+1})^{\beta_w} P(w_{t+1} | s_{t+1}) \right]^{\beta_w}
\end{aligned}$$



$$P(s_0) P(s_1 | s_0) = P(s_{t+1} | s_t, a_t) P(\beta_{t+1} | s_t, w_t) P(w_{t+1} | s_{t+1}, \beta_{t+1}) P(a_{t+1} | w_{t+1}, s_{t+1})$$

$$\begin{aligned}
&P(s_0, w_0, s_1, \beta_1, w_1, a_1, \dots, s_t, \beta_t, w_t, a_t) \\
&= P(w_0 | s_0) \prod_{k=0}^{\infty} P(w_{k+1} | s_{k+1}, \beta_{k+1}) P(a_{k+1} | w_{k+1}, s_{k+1}) P(\beta_{k+1} | s_k, w_k) P(s_0) P(s_1 | s_0) \prod_{t=1}^{\infty} P(s_{t+1} | s_t, a_t)
\end{aligned}$$

$$\begin{aligned}
&P(s_0, s_1, a_1, \dots, s_t, a_t) = \sum_{\beta} \sum_{w} P(s_0, w_0, s_1, \beta_1, w_1, a_1, \dots, s_t, \beta_t, w_t, a_t) \\
&= \sum_{w} \sum_{\beta} P(w_0 | s_0) \prod_{k=0}^{\infty} P(w_{k+1} | s_{k+1}, \beta_{k+1}) P(a_{k+1} | w_{k+1}, s_{k+1}) P(\beta_{k+1} | s_k, w_k) P(s_0) P(s_1 | s_0) \prod_{t=1}^{\infty} P(s_{t+1} | s_t, a_t) \\
&= \sum_{w} P(w_0 | s_0) \prod_{k=0}^{\infty} \left[ P(w_{k+1} | s_{k+1}, \beta_{k+1}) P(a_{k+1} | w_{k+1}, s_{k+1}) \sum_{\beta} P(\beta_{k+1} | s_k, w_k) \right] \left[ P(s_0) P(s_1 | s_0) \prod_{t=1}^{\infty} P(s_{t+1} | s_t, a_t) \right]
\end{aligned}$$

option, critic

Jupyter

temporal abstraction.

temporal in duration time

not in the space time

Survey problem



房间 → 房间: temporal abstraction

show 某种关系, action

但房间本身没有 temporal 关系。

学习问题: 一个由房间上类属的抽象

但是数据间没有 temporal relationship

根据内部也没有。

房间 → 房间: 与探索数据是否?

布天区别? :

数据从人没有。