Abstract
Intro
Related work
{ method }
{ EXP }
Citations

{ How is our work
better than others }

Journal # of building
engineering
Building Simulation

5 days to open the app

⇒ Train on March
Test on Jan and April

⇒ Train on Dec
Test on Nov and Dec. (slight
overlap)

⇒ Exp 13, 14, 15 ⇒ All uses Batyt
and spase reward.
↓
Vanilla policy gradiet. on spase reward

Exp 14 ⇒ Transfer learning ~~using~~ at agent level.

Exp 15 ⇒ Trasf lea. at env-model level.

Exp 4 ⇒ ┐ → Not transfer learning at all. here.
Exp 7 ⇒ ┘  Model is for constant & pricing
         scenario on bop fest - hydronic
         - heat fest.

         ↳ Model is learnt using Dynamic
           pricing
→ Gets from expriment 2 (Policy and
                          model based
                          RL with
                          Dynamic pricing
                          scenario)

( $\Rightarrow$ We may want to add an experiment where pretrained agent (for constat pricing scenario is used)

$\Rightarrow$ Methodology

  $\Rightarrow$ Same a Model based RL with policy gradient / Reinforce

  $\Rightarrow$ Transfer learning

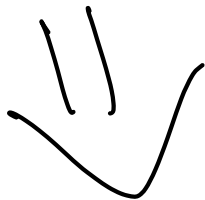  actions are added $\Rightarrow$ N N    $\Rightarrow$ state space remains same.

  $\Rightarrow$ How sparse reward is created and used in policy gradient algorithm.

  $\Rightarrow$

ALGORITHM STARTS BELOW

(1)    $\Theta = \Theta_I$   $\Leftarrow$ Initialized

$$\text{Policy-Loss} = - \sum_{t=0}^{T=1344} \log\left( \pi_{\Theta_I}(s_t) \right) R_t$$

where,

$$a_t = \pi_\Theta(s_t), \quad R_t = 0 \text{ if, } t\%50 \mathrel{!}= 0$$

where,

$\Rightarrow$ The negetive sign is to indicate we need to carry out gradient ascent instead of descent.

After training,   $\Theta = \Theta_T$

(2) Approach 2: Initialize $\Theta^P = \Theta^P_I$.

Pretraining:

$$\text{policy-loss} = -\sum_{t=1}^{T=672} \log\left(P_{\Theta^P}(s_t)\right) \cdot R_t$$

After pretraining: $\Theta^P = \Theta^P_{P.T}$

Fine tuning:

$$\text{policy-loss} = -\sum_{t=1}^{T=1344} \log\left(P_{\Theta^P}(s_t)\right) R_t$$

where,
$$R_t = 0 \quad \text{if} \quad t \% 50 \,! = 0$$

After Fine tuning:
$$\Theta^P = \Theta^P_F$$

③ Approach 3:

Here we pretrain using an environment model and finetune the environment model on bestest-hydronic.
We then use the finetuned environment model to train the RL agent

Initialize: $\Theta^E = \Theta^E_I$

Pretraining:
$$\left(E_{\Theta^E}(S_t, a_t)[0] - S_{t+1}\right)^2$$
$$+ \left(E_{\Theta^E}(S_t, a_t)[1] - R_t\right)^2$$

After pretraining: $\Theta^E = \Theta^E_{P.T}$

Finetuning $_{T=1344}$
$$env\_loss = \sum_{t=0}^{} \left(E_{\Theta_{P.T}}(S_t, a_t)[0] - S_{t+1}\right)^2$$
$$+ \left(E_{\Theta_{P.T}}(S_t, a_t)[1] - R_t\right)^2$$

$\Rightarrow$ only calculate loss when $t\%50 \,!= 0$

After Finetuning: $\Theta^E = \Theta^E_F$

Now training the policy gradient

$$\text{policy\_loss} = - \sum_{t=1}^{T=1344} \log\left(P_{\theta_I^P}(s_t)\right) . R_t$$

where,

If $t \% 50 == 0 \Rightarrow$ Use the actual environment

$t \% 50 != 0 \Rightarrow$ Use the surrogate environment.