
Algorithm 1: Training Process for PPO_Train Algorithm

```
1 function ppo_train (env, ppoTrainer, statePrepareList, greedyScores);  
   Input: Environment inheritance from gym.Env, ppoTrainer object, list of all statePrepare  
         for every problem and greedyScores list for problems  
   Data: statePrepares  $\leftarrow$  np.array(statePrepareList)  
   Data: State greedyScores  $\leftarrow$  np.array(greedyScores)  
   Data: best_score  $\leftarrow$  0.0; n_steps  $\leftarrow$  0  
   Data: score_history  $\leftarrow$  []; remain_cap_history  $\leftarrow$  []  
2 for Each i in N_TRAIN_STEP do  
   Data: batchs  $\leftarrow$  ppoTrainer.generate_batch(PROBLEMS_NUM, MAIN_BATCH_SIZE)  
3   for Each batch in batchs do  
4     env.setStatePrepare(statePrepares[batch])  
5     externalObservation, _  $\leftarrow$  env.reset()  
     Data: d  
6     one  $\leftarrow$  False  
7     while not done do  
       Data: internalObservatio, actions, accepted_actions, sumProbs, sumVals,  
             sumRewards, steps  $\leftarrow$  ppoTrainer.make_steps(externalObservation,  
                   env.statePrepares)  
       Data: externalObservation_, externalReward, done, info  $\leftarrow$   
             env.step(accepted_actions)  
8       ppoTrainer.save_step(externalObservation, internalObservatio, actions, sumProbs,  
                             sumVals, sumRewards, steps, done)  
9       n_steps  $\leftarrow$  n_steps + 1  
10      if n_steps % ppoTrainer.config.internal_batch == 0 then  
11        | ppoTrainer.train_minibatch()  
12      end  
13      externalObservation  $\leftarrow$  externalObservation_  
14    end  
    Data: scores, remain_cap_ratios  $\leftarrow$  env.final_score()  
    Data: batch_score_per_greedy  $\leftarrow$  mean of scores/greedyScores[batch]  
15    score_history.append(batch_score_per_greedy)  
16    remain_cap_history.append(np.mean(remain_cap_ratios))  
    Data: avg_score  $\leftarrow$  mean of score_history[-50:]  
17    if avg_score > best_score then  
      | Data: b  
18      | est_score  $\leftarrow$  avg_score ppoTrainer.save_models()  
19    end  
20  end  
21 end
```

Algorithm 2: make_step method in PPOTrainer class

```
1 function make_step (externalObservation, statePrepares);  
  Input: externalObservation, batch of statePrepares  
  Output: internalObservation, actions, accepted_actions, sumProbs, values, sumRewards,  
           steps  
  Data: actions  $\leftarrow$  zero tensor with shape of  $([MAIN\_BATCH\_SIZE, 0, 2])$   
  Data: sumLogProbs  $\leftarrow$  zero tensor with shape of  $([MAIN\_BATCH\_SIZE])$   
  Data: sumRewards  $\leftarrow$  zero tensor with shape of  $([MAIN\_BATCH\_SIZE])$   
  Data: internalObservation  $\leftarrow$  zero tensor with shape of  
         $([MAIN\_BATCH\_SIZE, 0, generat\_link\_number + 1, input\_decode\_dim])$   
  Data: step  $\leftarrow$  tensor of  $([1]MAIN\_BATCH\_SIZE)$   
  Data: steps  $\leftarrow$  zero tensor with shape of  $([MAIN\_BATCH\_SIZE, 0])$   
  Data: prompt  $\leftarrow$  None  
  Data: accepted_actions  $\leftarrow$  numpy array of  
         $([[[-1]2]generat\_link\_number]MAIN\_BATCH\_SIZE)$   
2 for Each  $i$  in generat_link_number do  
3   generatedInstance, generatedKnapsack, prompt  $\leftarrow$  actor_model.generateOneStep(step,  
     externalObservation, prompt)  
   Result: updated prompt if prompt == None and get new distributions as  
     generatedInstance, generatedKnapsack  
4   act, log_prob  $\leftarrow$  _choose_actions(generatedInstance, generatedKnapsack)  
   Result: get tensor of act as  $([inst\_act, ks\_act])$  and log_pro as summation of instance log  
     prob with knapsack log prob  
5   actions  $\leftarrow$  concatenate of actions and act  
6   reward  $\leftarrow$  reward(act, accepted_actions, step, prompt, statePrepares)  
   Result: get internalReward and if instance is allocated in the knapsack update  
     accepted_actions, step, prompt  
7   steps  $\leftarrow$  concatenate of steps, step to trac of step in actor model  
8   internalObservation  $\leftarrow$  concatenate of internalObservation, prompt  
9   sumProbs  $\leftarrow$  sumProbs + prob  
10  sumRewards  $\leftarrow$  sumRewards + reward  
11 end
```

Algorithm 3: train_minibatch method in PPOTrainer class

```
1 function train_minibatch ();
2 for Each _ in ppo_epochs do
    Data: batchs ← generate_batch()
3   for Each index in MAIN_BATCH_SIZE do
        Data: obs ← memoryObs[index]
        Data: intObs ← memoryIntObs[index]
        Data: acts ← memoryAct[index]
        Data: probs ← memoryPrb[index]
        Data: rewards ← memoryRwd[index]
        Data: vals ← memoryVal[index]
        Data: stps ← memoryStp[index]
        Data: done ← memoryDon[index]
        Data: advantage ← zero tensor with shape of ([internal_batch])
4     for Each t in internal_batch-1 do
        Data: discount ← 1
        Data: a.t ← 0
5       for Each k between t and internal_batch-1 do
6         a.t += discount * (rewards[k] + config.gamma * vals[k+1] * (1 - int(done[k]))
           - vals[k])
7         discount *= gamma * gae_lambda
8       end
9       advantage[t] ← a.t;
10    end
11    for Each batch in batches do
        Data: batchObs ← obs[batch]
        Data: batchIntObs ← intObs[batch]
        Data: batchActs ← acts[batch]
        Data: batchSteps ← stps[batch]
        Data: batchProbs ← probs[batch].to(device)
        Data: batchVals ← vals[batch].to(device)
12      new_log_probs ← torch.tensor([0] * config.ppo_batch_size, dtype=torch.float64,
        device=device);
13      for Each i in generat_link_number do
14        generatedInstance, generatedKnapsack, _ ←
            actor_model.generateOneStep(batchSteps[:,i], batchObs.to(device),
            batchIntObs[:,i]) Result: get new distributions as generatedInstance,
            generatedKnapsack for external and internal observations
        Data: inst_dist ← Categorical(generatedInstance)
        Data: ks_dist ← Categorical(generatedKnapsack)
        Result: get new torch distributions categorical objects
        Data: inst_log_probs ← inst_dist.log_prob(batchActs[:,i,0])
        Data: ks_log_probs ← ks_dist.log_prob(batchActs[:,i,1])
        Data: newProbs ← inst_log_probs + ks_log_probs
15      new_log_probs += newProbs
16    end
        Data: newVal ← critic_model(batchObs)
        Data: prob_ratio ← exponential of (new_log_probs - batchProbs)
        Data: weighted_probs ← advantage[batch] * prob_ratio
        Data: weighted_clipped_probs ← torch.clamp(prob_ratio, 1 - cliprange, 1 +
            cliprange) * advantage[batch]
        Data: actor_loss ← mean of -min between weighted_probs and
            weighted_clipped_probs
        Data: returns ← advantage[batch] + batchVals
        Data: critic_loss ← mean of euclidean distance between returns and newVal
        Data: total_loss ← actor_loss + 0.5 * critic_loss
17      actor_optimizer.zero_grad()
18      critic_optimizer.zero_grad()
19      total_loss.backward()
20      actor_optimizer.step()
21      critic_optimizer.step()
22    end
23  end
24 end
```
