



Depp Reinforcement Learning of Marked Temporal Point Processes

Utkarsh Upadhyay
MPI-SWS
utkarshu@mpi-sws.org

Abir De
MPI-SWS
ade@mpi-sws.org

Manuel Gomez-Rodriguez
MPI-SWS
manuelgr@mpi-sws.org

(Published in NeurIPS-2018)

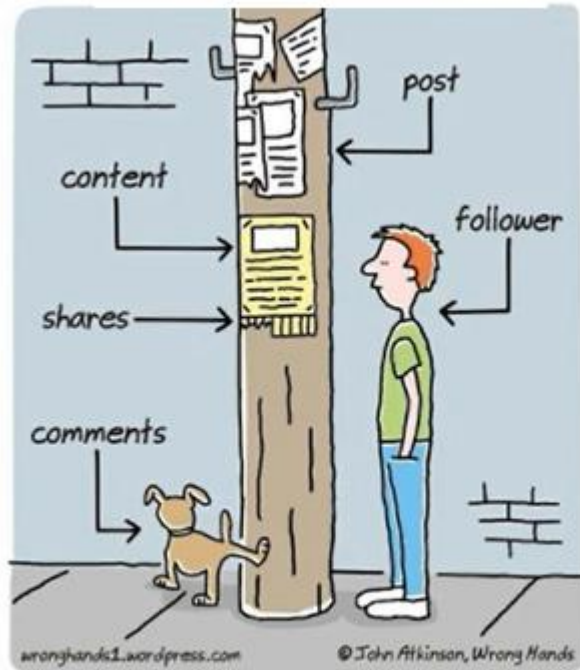
孙佳佳



Outline

- **Background**
- Representation of MTPPs
- Reinforcement learning models
- Policy Optimization
- Evaluation

Discrete events in continuous time



Posts and shares



Test and diagnoses



Purchases

Marked temporal point processes



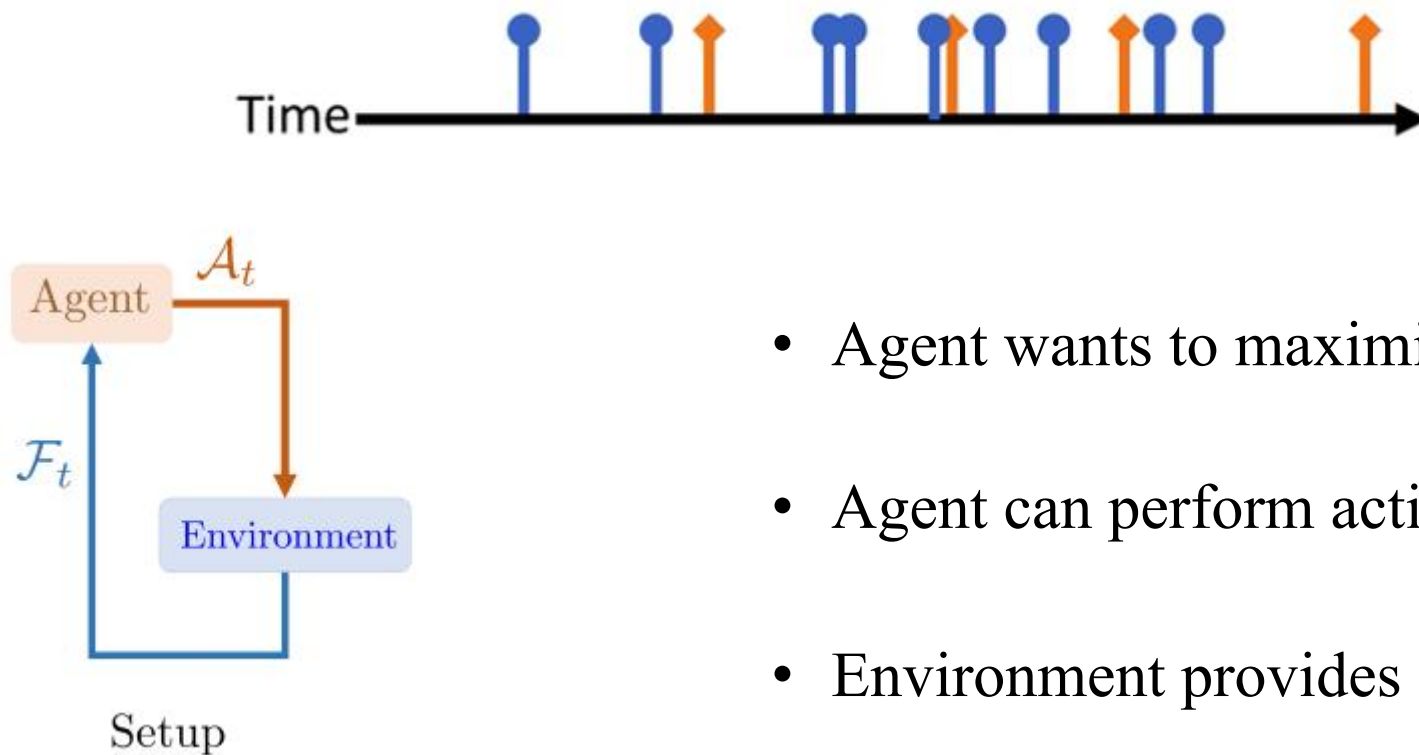
$t_i \in \mathbb{R}$: Times

$z_i \in \mathbb{Z}$: Marks

$$\mathcal{H}_t = \{e_0 = (t_0, z_0), \dots, e_n = (t_n, z_n)\}$$

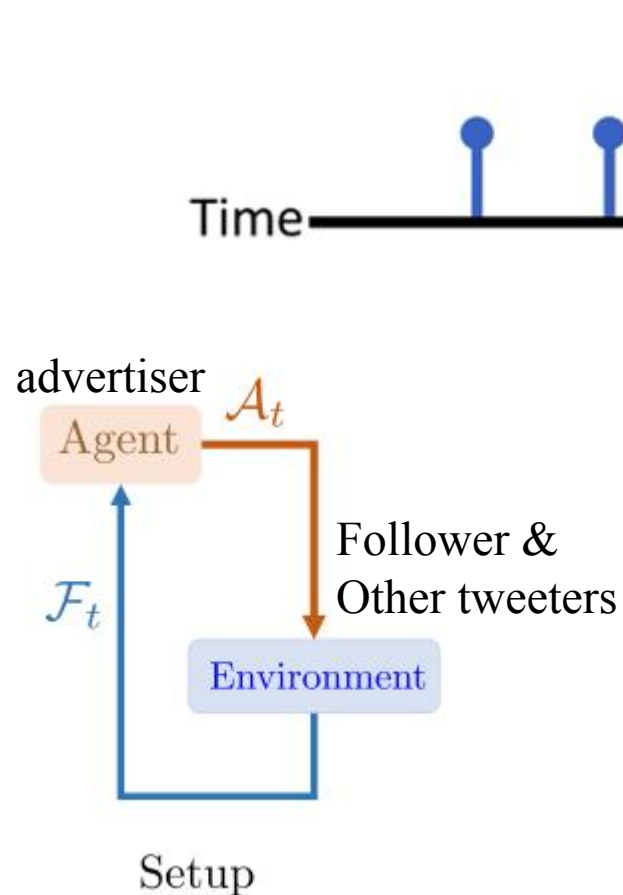
- Sequence of events of type z_i with their times t_i
 - Continuous time
 - Discrete/continuous marks

What can MTPPs model?



- Agent wants to maximize som reward
- Agent can perform action y_j at time t_j
- Environment provides feedback at t_i
- Reward is calculated after each episode

Example 1: When to post



- Advertiser wants to gain attention
- Advertiser can **post tweets**
- **Other tweeter** post tweets
- **Avg.rank** is calculated at campaign end

$$\mathcal{A}_t = \{9:10, 9:30, \dots\}$$

$$\mathcal{F}_t = \{(1, 9:15), (2, 9:31), \dots\}$$

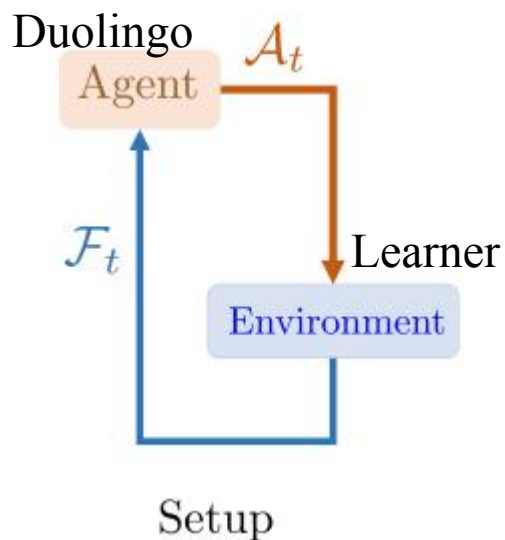
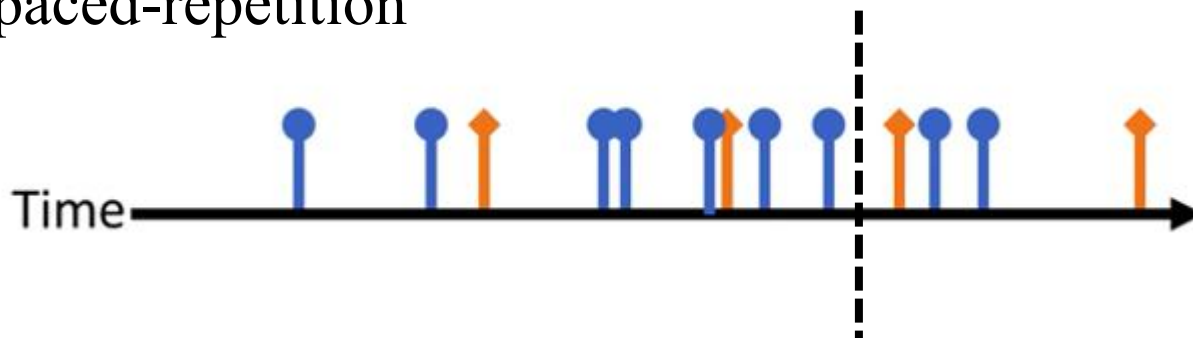
$$R(T) = \frac{1}{T} \int_0^T \text{rank}(t) dt$$



Example 2: spaced-repetition



Online learning platform



- Duolingo wants to teach
- Duolingo can **ask questions**
- Learner can try to **recall**
- **Test score** is calculated at the end of course

$$\mathcal{A}_t = \{(abandon, 9:05), (reinforce, 9:30), \dots\}$$

$$\mathcal{F}_t = \{(\checkmark, 9:05), (\times, 9:30), \dots\}$$

$$R(T) = (exam\ time)$$



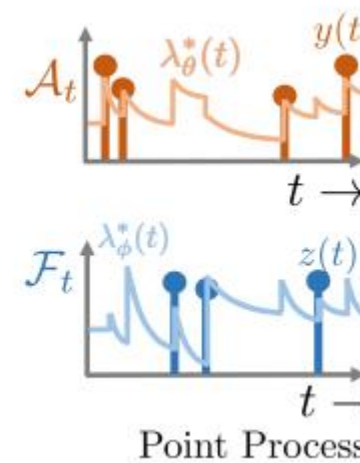
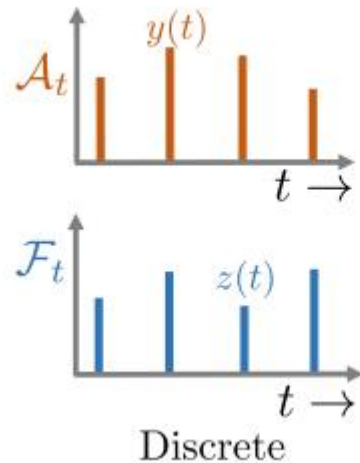
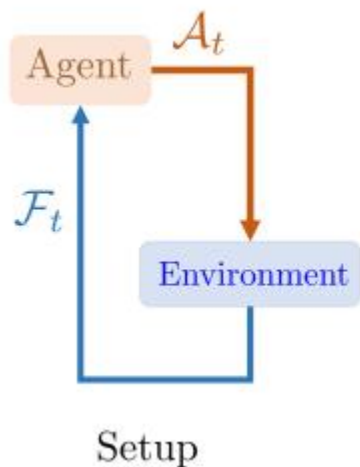
Current methods

- Rely on SDE based modeling of the **environment**
- Rely on carefully chosen **reward** functions

limitations

- **Feedback** process model is complex/unknown
- **Reward function** can be arbitrarily complex
- Cannot leverage advances in deep learning

Classical RL is too restrictive



Classical RL

- Actions and feedbacks at **discrete time steps**
- Policy determines only the mark of actions

RL of MTPPs

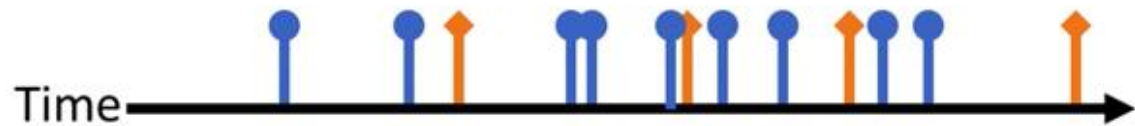
- **Asynchronous** actions and feedbacks in continuous time
- Policy determines the time as well as mark of the next action
- **Arbitrary reward function**



Outline

- Background
- **Representation of MTPPs**
- Reinforcement learning models
- Policy Optimization
- Evaluation

How to represent MTPPs



$t_i \in \mathbb{R}$: Times

$z_i \in \mathbb{Z}$: Marks

$$\mathcal{H}_t = \{e_0 = (t_0, z_0), \dots, e_n = (t_n, z_n)\}$$

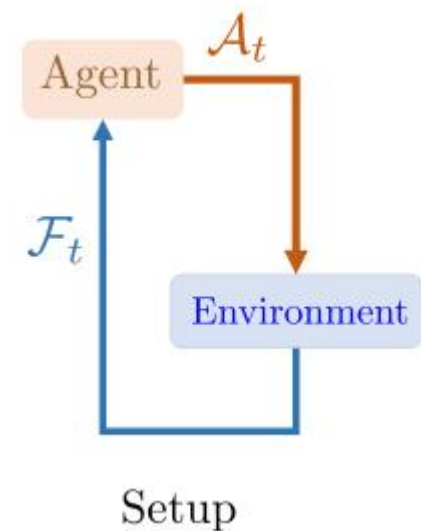
- Intensity function: $\lambda^*(t) := \mathbb{P}\{\text{event in } [t, t + dt) \mid \mathcal{H}_t\},$
- Markd distribution: $m(z \mid \mathcal{H}_t) = m^*(z)$
- Likelihood of $\mathcal{A}_T \subseteq \mathcal{H}_T$

$$\mathbb{P}(\mathcal{A}_T) := \left(\prod_{e_i \in \mathcal{A}_T} \underbrace{\lambda^*(t_i)}_{\text{Prob. of an action at } t_i} \underbrace{m^*(z_i)}_{\text{Prob. of mark } z_i} \right) \overbrace{\exp \left(- \int_0^T \lambda^*(s) ds \right)}^{\text{Prob. of no actions at } t \in [0, T] \setminus \{t_i\}}$$

How to represent our problem as MTPPs



- Actions: $\mathcal{A} = \{e_i = (t_i, y_i)\}$, where $(t_i, y_i) \sim p_{\mathcal{A};\theta}^* = (\lambda_{\theta}^*, m_{\theta}^*)$
- Feedbacks: $\mathcal{F} = \{f_i = (t_i, z_i)\}$, where $(t_i, z_i) \sim p_{\mathcal{F};\phi}^* = (\lambda_{\phi}^*, m_{\phi}^*)$



$$\underset{p_{\mathcal{A};\theta}^*(\cdot)}{\text{maximize}} \quad \mathbb{E}_{\mathcal{A}_T \sim p_{\mathcal{A};\theta}^*(\cdot), \mathcal{F}_T \sim p_{\mathcal{F};\phi}^*(\cdot)} [R^*(T)]$$

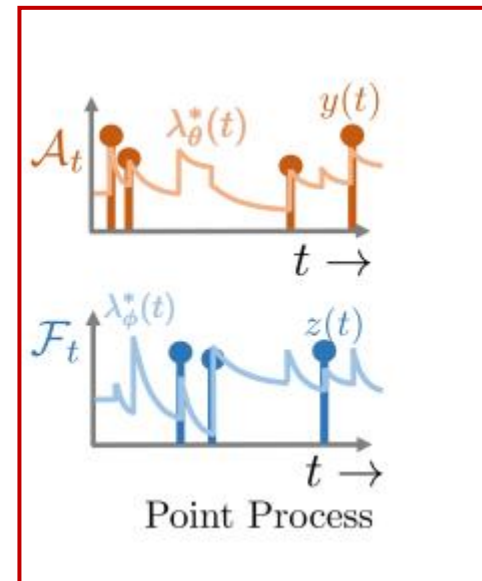
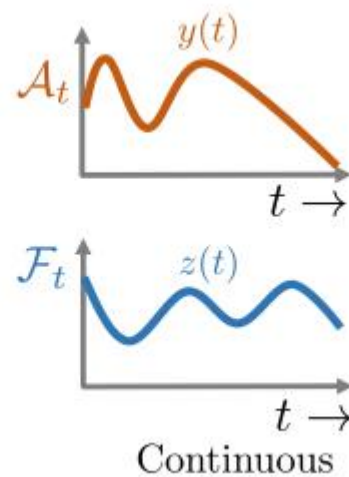
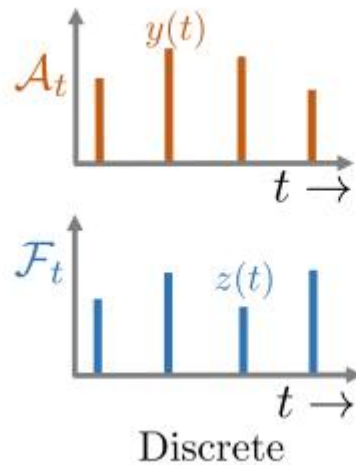
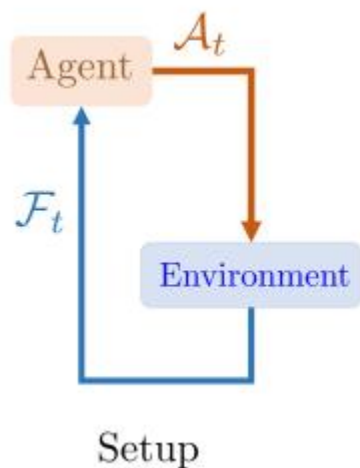
Model free reinforcement
learning



Outline

- Background
- Representation of MTPPs
- **Reinforcement learning models**
- Policy Optimization
- Evaluation

Reinforcement Learning model



- Design the policy
- Realize the policy

- Actions and feedback are **asynchronous**
- Characterized by **intensity function**

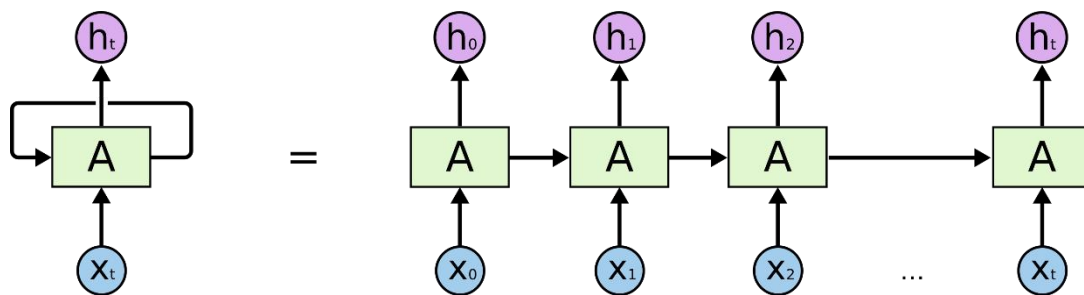
Policy parametrization

$$\lambda^*(t) := \mathbb{P}\{\text{event in } [t, t + dt) \mid \mathcal{H}_t\},$$

$$\mathcal{H} = \{e_0 = (t_0, z_0), e_1 = (t_1, z_1), \dots, e_n = (t_n, z_n)\}$$

Embed the history using RNN

- Each event updates the hidden state \mathbf{h}
- Intensity function is determined by $\lambda_k(t) = f_k(\mathbf{w}_k^\top \mathbf{h})$



Embed the history using RNN

Input
layer

$$\begin{aligned}\tau_i &= \mathbf{W}_t(t_i - t_{i-1}) + \mathbf{b}_t, \\ \mathbf{b}_i &= \mathbf{W}_a(1 - e_i) + \mathbf{W}_f e_i + \mathbf{b}_b,\end{aligned}$$

$$\begin{aligned}\mathbf{y}_i &= \mathbf{W}_y y_i + \mathbf{b}_y \text{ if } e_i = 0 \\ \mathbf{z}_i &= \mathbf{W}_z z_i + \mathbf{b}_z \text{ if } e_i = 1\end{aligned}$$

Hidden
layer

$$\mathbf{h}_i = \tanh(\mathbf{W}_h \mathbf{h}_{i-1} + \mathbf{W}_1 \tau_i + \mathbf{W}_2 \mathbf{y}_i + \mathbf{W}_3 \mathbf{z}_i + \mathbf{W}_4 \mathbf{b}_i + \mathbf{b}_h)$$

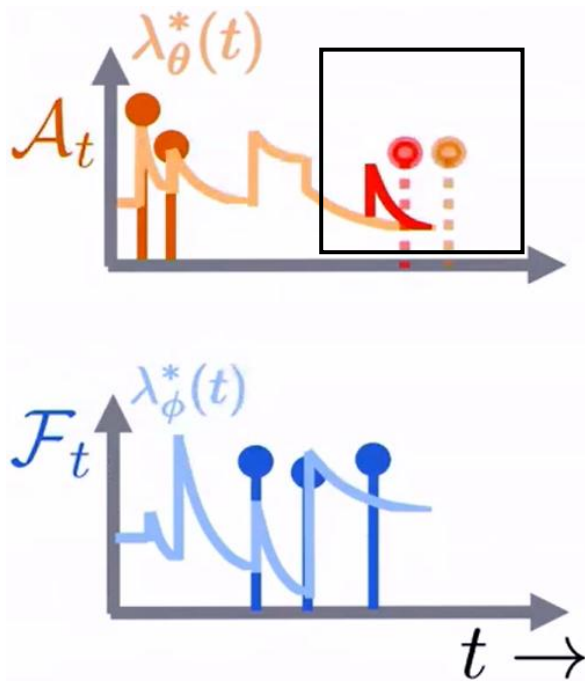
Output
layer

$$\lambda_{\theta}^*(t) = \exp(b_{\lambda} + w_t(t - t_i) + \mathbf{V}_{\lambda} \mathbf{h}_i) \quad \mathbb{P}[y_{i+1} = c] = \frac{\exp(\mathbf{V}_{c,:}^y \mathbf{h}_i)}{\sum_{l \in \mathcal{Y}} \exp(\mathbf{V}_{l,:}^y \mathbf{h}_i)}$$

$$\begin{aligned}\theta = \{ & \mathbf{V}_{\lambda}, \mathbf{V}^y, w_t, b_{\lambda}, \\ & \mathbf{W}_z, b_z, \mathbf{W}_t, b_t, \mathbf{W}_b, b_b, \\ & \mathbf{W}_h, b_h, \mathbf{W}_y, b_y \}\end{aligned}$$

RL with Asynchronous Feedback

- Distribution of actions may **change** because of incoming feedback



Algorithm 1: Returns the next action time

```

1: Input: Parameters  $b_{\lambda}, w_t, \mathbf{V}_{\lambda}, h_i$ , last event time  $t'$ 
2: Output: Next action time  $t$ 
3:  $CDF(\bullet) \leftarrow$  Cumulative distribution of next arrival time
4:  $u \leftarrow \text{UNIF}[0, 1]$ 
5:  $t \leftarrow CDF^{-1}(u)$ 
6: while  $t < T$  do
7:    $(s, z) \leftarrow \text{WAITUNTILNEXTFEEDBACK}(t)$ 
8:   if feedback arrived before  $t$  then
9:      $CDF(\bullet) \leftarrow \text{MODIFY}(CDF(\bullet), s, z)$ 
10:     $t \leftarrow CDF^{-1}(u)$ 
11:   else
12:     return  $t$ 
13:   end if
14: end while
15: return  $t$ 

```



RL problem with PTPPs: summary

- **New:** RL problem in continuous time with discrete events
 - **Embedding state** to real vectors using RNN
 - Efficient and unbiased **sampling** procedure to handle asynchronous feedbacks



Outline

- Background
- Representation of MTPPs
- Reinforcement learning models
- **Policy Optimization**
- Evaluation



Policy Gradient method

- Maximizing the expected reward

$$J(\theta) = \mathbb{E}_{\mathcal{A}_T \sim p_{\mathcal{A};\theta}^*(\cdot), \mathcal{F}_T \sim p_{\mathcal{F};\phi}^*(\cdot)} [R^*(T)]$$

- SGD Updates:

$$\theta_{l+1} = \theta_l + \alpha_l \nabla_{\theta} J(\theta)|_{\theta=\theta_l}$$

- REINFORCE trick:

- No need to model the feedback process
- Reward function can be arbitrary

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\mathcal{A}_T \sim p_{\mathcal{A};\theta}^*(\cdot), \mathcal{F}_T \sim p_{\mathcal{F};\phi}^*(\cdot)} [R^*(T) \nabla_{\theta} \log \mathbb{P}_{\theta}(\mathcal{A}_T)]$$

$$\log \mathbb{P}_{\theta}(\mathcal{A}_T) = \sum_{e_i \in \mathcal{A}_T} (\log \lambda_{\theta}^*(t_i) + \log m_{\theta}^*(z_i)) - \int_0^T \lambda_{\theta}^*(s) ds.$$



Outline

- Background
- Representation of MTPPs
- Reinforcement learning models
- Policy Optimization
- **Evaluation**

Spaced repetition

$$\mathcal{A}_t = \{(abandon, 9:05), (reinforce, 9:30), \dots\}$$

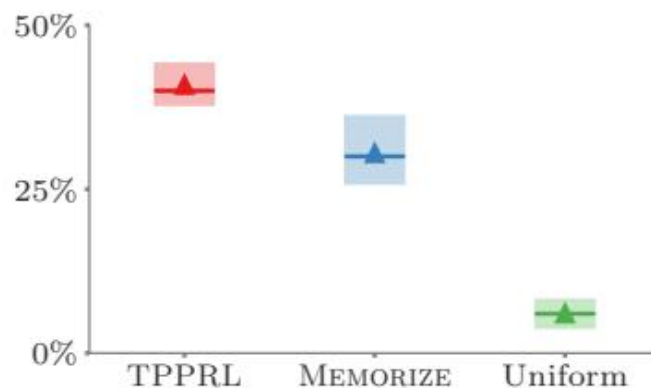
$$\mathcal{F}_t = \{(\checkmark, 9:05), (\times, 9:30), \dots\}$$

$$R(T) = (exam\ time)$$

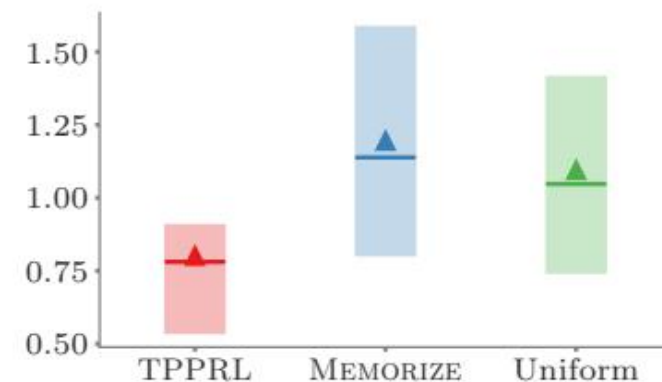


Online learning platform

TPPRL	RL of MTPP; No assumption on feedback; Arbitrarily complex reward function
MEMORIZE	Has full access to the student model; Design specific reward function
Uniform	Choose items uniformly at random



(a) Recall



(b) Items' difficulty

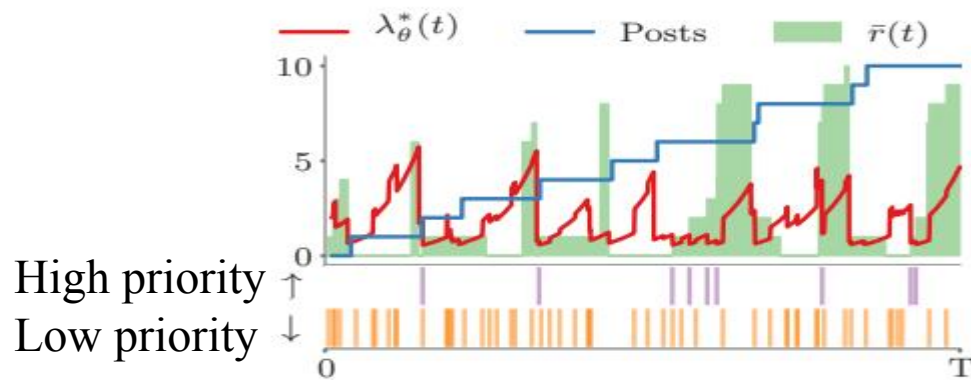
When to post

TPPRL	RL of MTPP; No assumption on feedback; Arbitrarily complex reward function
RQ	feeds sorted in reverse chronological order; minimize the average rank; intensity $\propto \text{rank}_{\text{chrono}}(t)$
RQ^*	feeds sorted in reverse chronological order; minimize the average rank; intensity $\propto \text{rank}_{\text{priority}}(t)$
Karimi	feeds sorted in reverse chronological order; maximize the time at the top

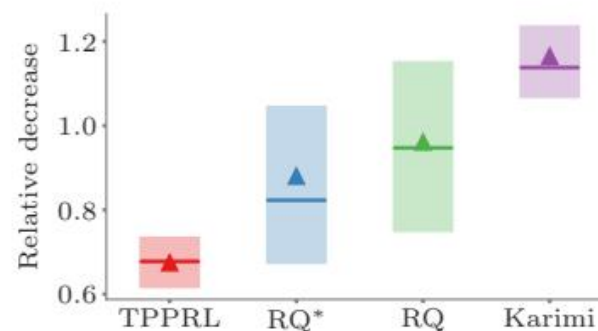
$$\mathcal{A}_t = \{9:10, 9:30, \dots\}$$

$$\mathcal{F}_t = \{(1, 9:05), (2, 9:30), \dots\}$$

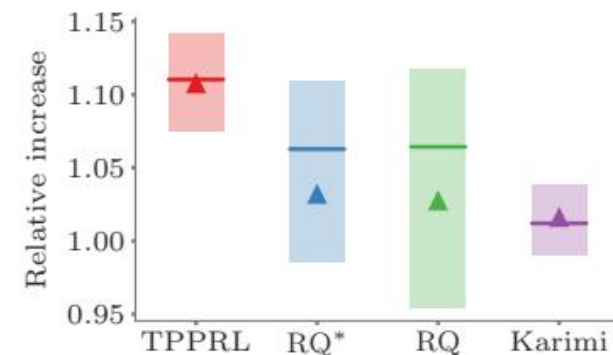
$$R(T) = \frac{1}{T} \int_0^T \text{rank}(t) dt$$



(c) Example



(a) Average rank



(b) Time at top



Conclusion

Reinforcement Learning method for:

- Asynchronous actions/rewards represented as MTPPs
- Arbitrary reward functions

Future work:

- Deriving more sophisticated reinforcement learning algorithms
- Multi-agent reinforcement learning



Thank you!