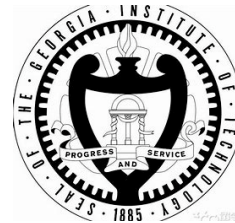


# Hierarchical Reinforcement Learning for Course Recommendation in MOOCs

Jing Zhang, Bowen Hao, Bo Chen, Cuiping Li, Hong Chen  
(Renmin University)

Jimeng Sun (Georgia Institute of Technology)



# Course Recommendation



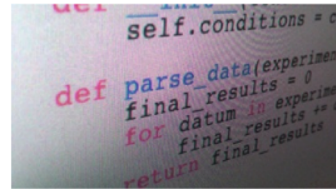
人人都能学计算机：计算机科学入门与Python编程

来自于：哈维穆德学院 | 分类：计算机(502)



相关课程

Related courses



计算机科学和Python编程导论  
(自主模式)



Python程序设计

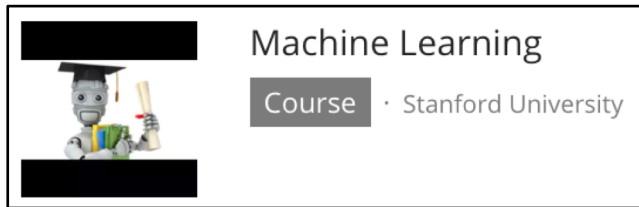


MyCS: 计算机科学入门

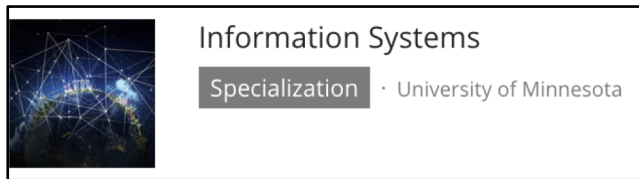


# Problem Definition

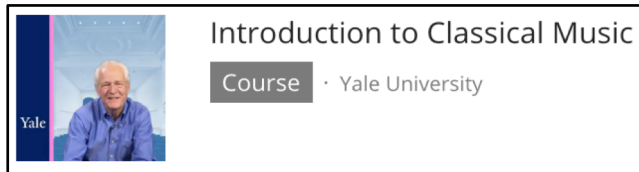
**Input:** Historical enrolled courses of a user before  $t$



Machine Learning  
Course · Stanford University



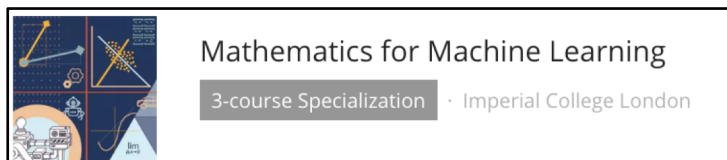
Information Systems  
Specialization · University of Minnesota



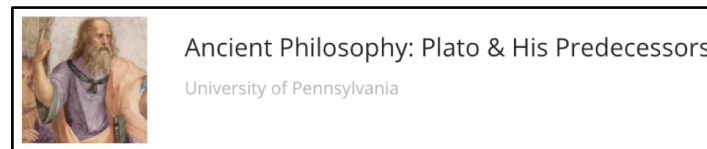
Introduction to Classical Music  
Course · Yale University

**Output:**

Most possible courses to be enrolled at  $t+1$

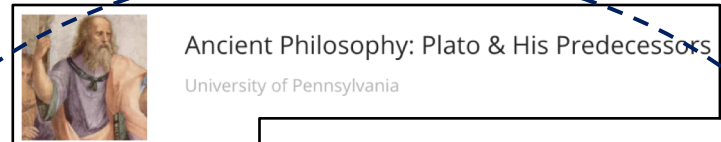


Mathematics for Machine Learning  
3-course Specialization · Imperial College London



Ancient Philosophy: Plato & His Predecessors  
University of Pennsylvania

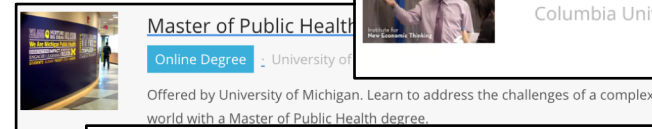
**Candidate Courses:**



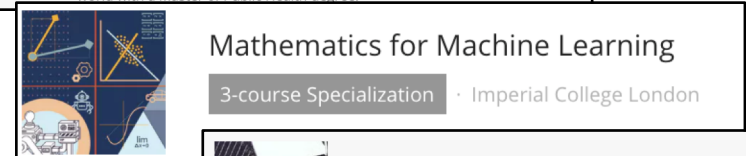
Ancient Philosophy: Plato & His Predecessors  
University of Pennsylvania



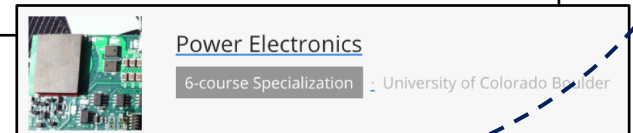
Economics of Money and Banking  
Columbia University



Master of Public Health  
Online Degree · University of Michigan  
Offered by University of Michigan. Learn to address the challenges of a complex world with a Master of Public Health degree.

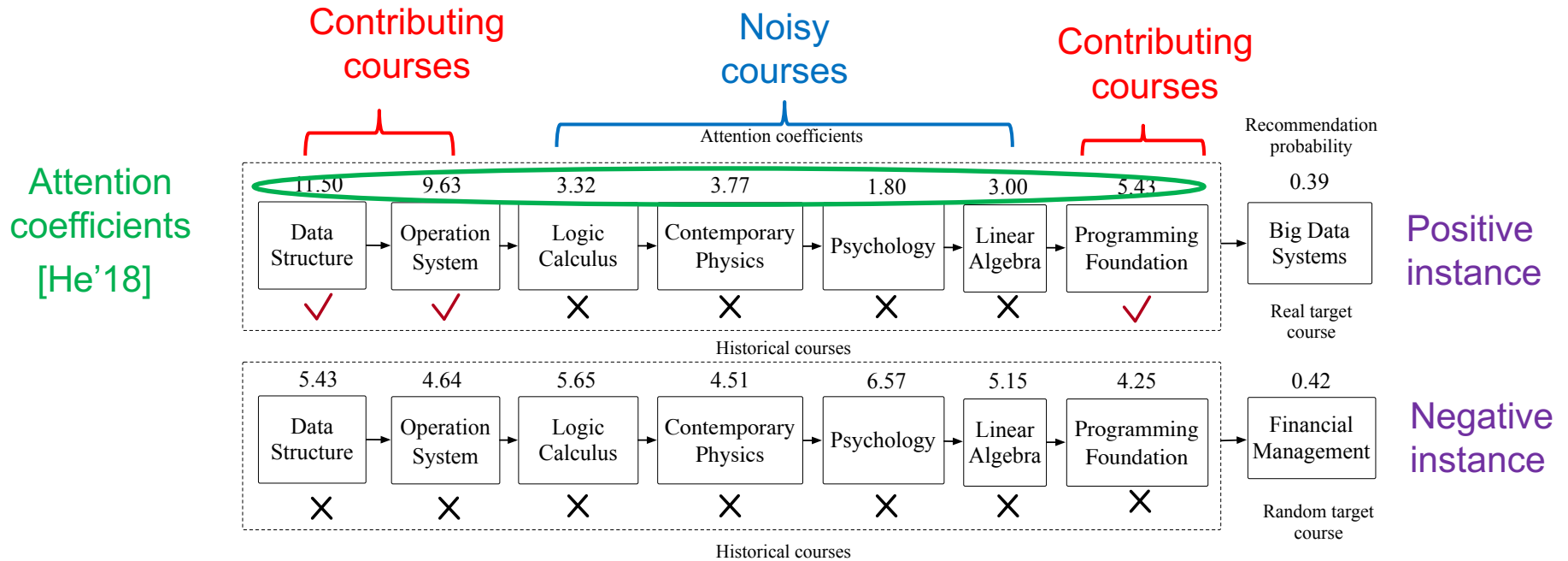


Mathematics for Machine Learning  
3-course Specialization · Imperial College London



Power Electronics  
6-course Specialization · University of Colorado Boulder

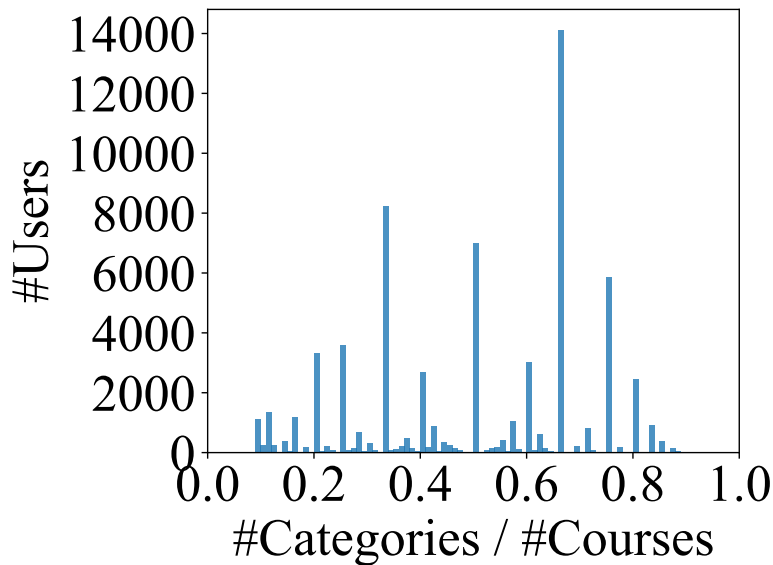
# Challenges



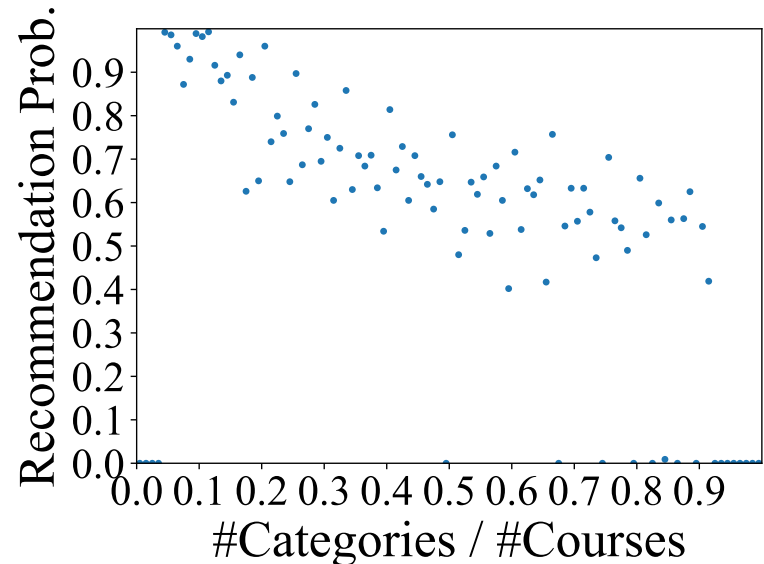
- Users' enrolled courses are usually diverse. The contributing courses may be diluted by noisy courses
- Even if no courses can contribute in predicting a random target course, each historical course will still be assigned an attention coefficient.

# Data Analysis

Users	Courses	Categories	User-course pairs	Time
82,535	1,302	23	458,454	2016.10.1-2018.3.31



A large number of users enrolled diverse courses



The recommendation performance based on the diverse profiles is impacted

(A bigger #categories/#courses indicates the user is more distractive)

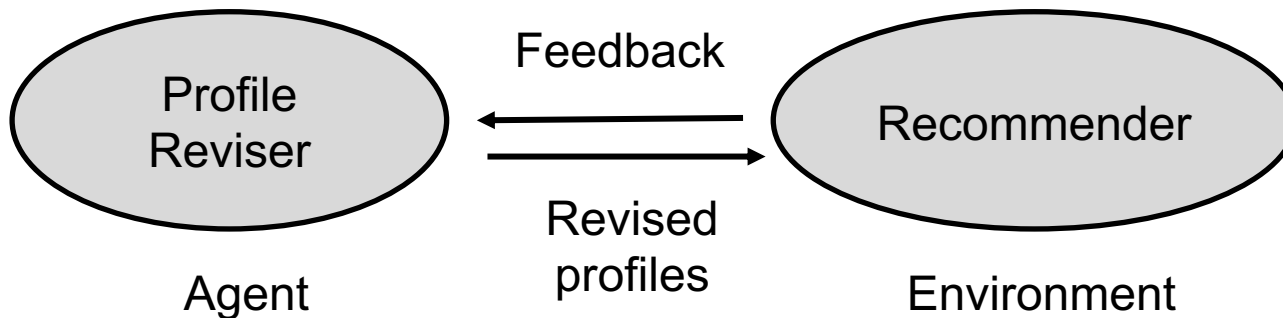
# Idea to Deal with the Challenges

---

- Revise the user profiles by removing the noisy courses instead of assigning an attention coefficient to each of them.
  - But how to determine which courses should be removed?
  - Without the supervised information, can we automatically learn the pattern?

# Solution

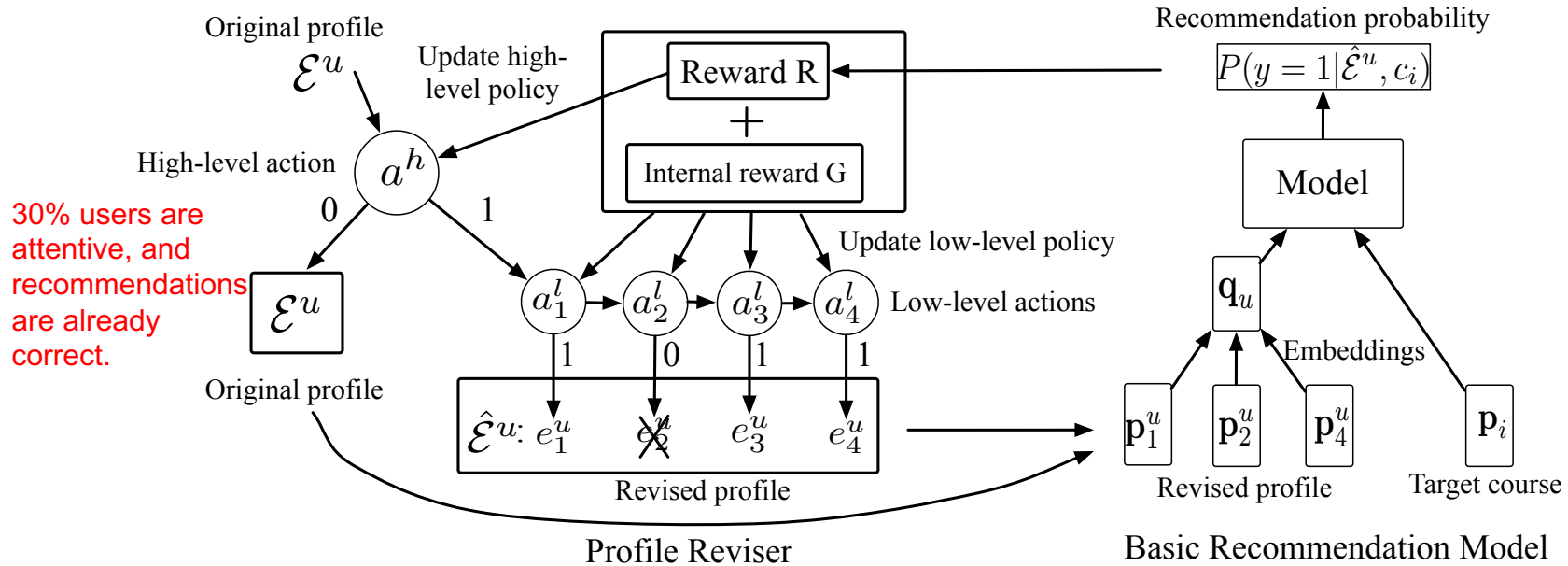
Revise the user profiles based on the feedbacks from the recommender



The revising process of a user profile: a sequential decision process

- The profile reviser (agent)
  - Revise the profile
  - Gets a delayed reward from the recommender
  - Update its policy
- The recommender (environment)
  - Update its parameters based on the profiles revised by the profile reviser

# Framework



- A hierarchical Markov Decision Process
  - The agent firstly performs a high-level task to determine whether to revise the whole profile or not
  - If it decides to revise, the agent performs a low-level task of multiples actions to determine whether to remove each historical course or not
  - The overall task is finished after the low-level task is finished or the high-level task decides to make no revision.



# The High-level Task

- Determine whether to revise the whole profile  $\varepsilon^u$  or not

- State:

- The average cosine similarity between the embedding vectors of each historical course in  $\varepsilon^u$  and the target course  $c_i$ .
- The average element wise product between the embedding vectors of each historical course in  $\varepsilon^u$  and the target course  $c_i$ .
- The probability  $P(y = 1|\varepsilon^u, c_i)$  of recommending  $c_i$  to user  $u$  by the basic recommender.

(The lower recommendation probability is, more effort should be taken to revise the profile)

- Action:

- {Revise, keep}

- Delayed reward:

Policy function: two-layer NN

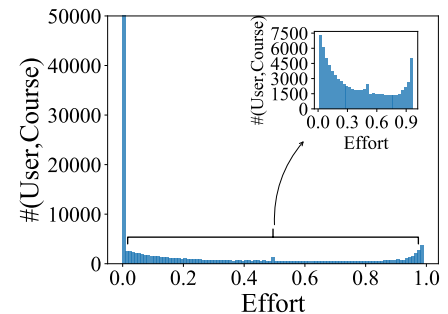
$$\begin{aligned}\mathbf{H}_t^l &= \text{ReLU}(\mathbf{W}_1^l \mathbf{s}_t^l + \mathbf{b}^l), \\ \pi(\mathbf{s}_t^l, a_t^l) &= P(a_t^l | \mathbf{s}_t^l, \Theta^l) \\ &= a_t^l \sigma(\mathbf{W}_2^l \mathbf{H}_t^l) + (1 - a_t^l)(1 - \sigma(\mathbf{W}_2^l \mathbf{H}_t^l)),\end{aligned}\tag{3}$$

The difference between the log-likelihood after and before the profile is revised.

$$R(a_t^l, \mathbf{s}_t^l) = \begin{cases} \log p(\hat{\mathcal{E}}^u, c_i) - \log p(\mathcal{E}^u, c_i). & \text{if } t = t_u; \\ 0 & \text{otherwise,} \end{cases}$$

# The low-level Task

- Determine whether to remove a historical course  $e_t^u \in \varepsilon^u$  or not
  - **State:**
    - The cosine similarity between the embedding vectors of the current historical course  $e_t^u$  and the target course  $c_i$ .
    - The element wise product between the embedding vectors of the current historical course  $e_t^u$  and the target course  $c_i$ .
    - Effort taken in the course.
  - **Action:**
    - {Remove, Keep}
  - **Reward:**
    - Add an internal reward
      - speed up local learning and does not propagate to the high-level.
    - $G(a_t^l, s_t^l)$  : calculate the average cosine similarity between each historical course and the target course after and before the profile is revised.



$$R(a_t^m, s_t^m) + G(a_t^m, s_t^m)$$

Effort: we calculate the ratio between the watch duration and the total duration of a video as the watch ratio, and use the maximal watch ratio of all the videos in a course to represent the effort taken by the user in the course

# Objective Function

- Maximize the expected reward:

$$\Theta^* = \operatorname{argmax}_{\Theta} \sum_{\tau} P_{\Theta}(\tau; \Theta) R(\tau),$$

A sequence of the sampled actions and the transited states

$\{s_1^l, a_1^l, s_2^l, \dots, s_t^l, a_t^l, s_{t+1}^l, \dots\}$

Sampling probability

The reward for the sampled sequence

- Update the policy network with policy gradient:

$$\nabla_{\Theta} = \frac{1}{m} \sum_{m=1}^M \nabla_{\Theta} \log \pi_{\Theta}(\mathbf{s}^m, a^m) R(a_t^m, \mathbf{s}_t^m).$$

# Training Procedure

Pre-train the basic recommendation model;  
Pre-train the profiler reviser by running Algorithm 2  
with the basic recommendation model fixed;  
Jointly train the two models together by running  
Algorithm 2;

Pre-train + joint train

**Algorithm 1:** The Overall Training Process

**Input:** Training data  $\{\mathcal{E}^1, \mathcal{E}^2, \dots, \mathcal{E}^{|\mathcal{U}|}\}$ , a pre-trained  
basic recommendation model and a profile  
reviser parameterized by  $\Phi^0$  and  $\Theta^0$  respectively

Initialize  $\Theta = \Theta^0, \Phi = \Phi^0$ ;

**for** episode  $l=1$  to  $L$  **do**

**foreach**  $\mathcal{E}^u := (e_1^u, \dots, e_{t_u}^u)$  and  $c_i$  **do**

    Sample a high-level action  $a^h$  with  $\Theta^h$ ;

**if**  $a^h = 0$  **then**

$R(s^h, a^h) = 0$

**else**

      Sample a sequence of low-level actions

$\{a_1^l, a_2^l, \dots, a_{t_u}^l\}$  with  $\Theta^l$ ;

      Compute  $R(a_{t_u}^l, s_{t_u}^l)$  and  $G(a_{t_u}^l, s_{t_u}^l)$ ;

      Compute gradients by Eq. (5) and (6);

**end**

**end**

  Update  $\Theta$  by the gradients;

  Update  $\Phi$  in the basic recommendation model;

**end**

Sample actions and  
transited states and  
get rewards

Update profile reviser  
and recommender

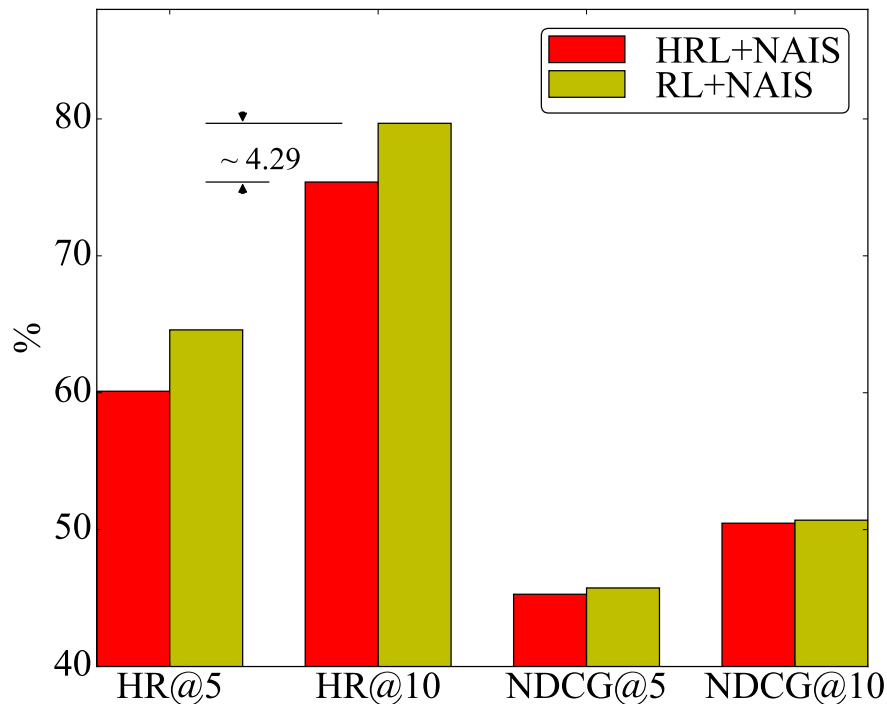
**Algorithm 2:** The Hierarchical Reinforcement Learning

# Experiment Results

Table 1: Recommendation performance (%).

	Methods	HR@5	HR@10	NDCG@5	NDCG@10
User-item model	BPR	46.82	60.73	34.16	38.65
	MLP	52.16	66.29	40.39	44.41
	FM	46.01	61.07	35.28	40.15
Item-item model, average	FISM	52.73	65.64	40.00	44.98
	GRU	52.07	68.63	38.92	46.30
Attention model	NAIS	56.42	69.05	43.73	47.82
	NASR	54.64	69.48	42.39	47.33
RL model	HRL+NAIS	<b>64.59</b>	<b>79.68</b>	<b>45.74</b>	<b>50.69</b>
	HRL+NASR	<b>59.05</b>	<b>74.50</b>	<b>47.51</b>	<b>52.73</b>

# Compared with One-level RL



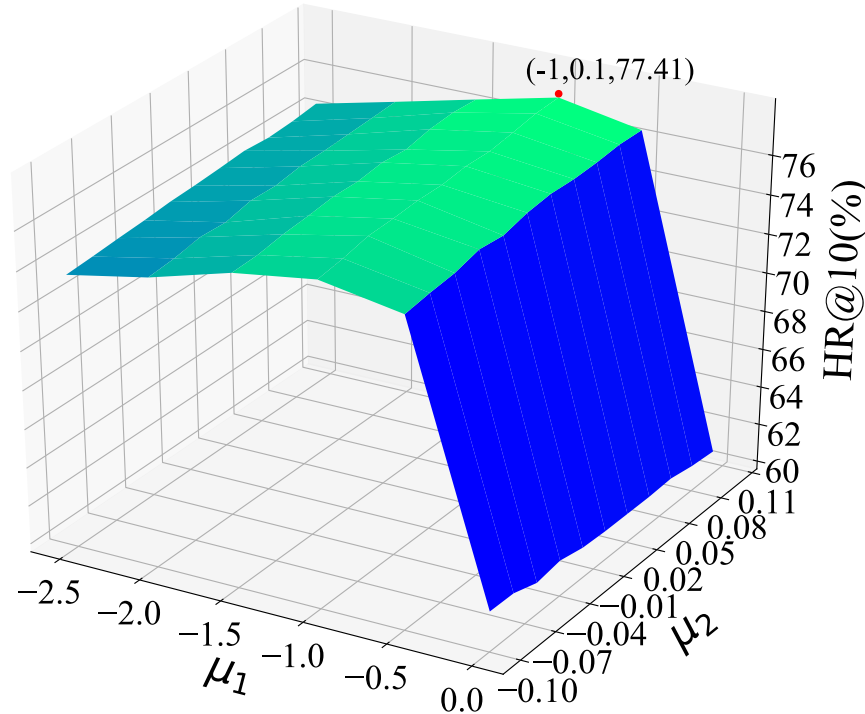
- The average #Categories/#Courses of the revised profiles:
  - Two-level RL: **0.73**
  - One-level RL: **0.75**
- The revised profiles by the two-level HRL are more consistent

The high-level task of the two-level RL:

- The average #Categories/#Courses of
  - the kept profiles: **0.57**
  - the revised profiles: **0.69**
- High-level task tends to keep more consistent profiles

(A bigger #categories/#courses indicate the user is more distractive)

# Compared with Greedy Revision



The greedy reviser

- firstly decides to revise the whole profile if  $P(y = 1 | \varepsilon^u, c_i) < \mu_1$
- and then removes the course  $e_t^u \in \varepsilon^u$  if its cosine similarity with  $c_i$  is less than  $\mu_2$

# Compared with Attentions

Table 2: Case studies of the profiles revised by HRL+NAIS and the attention coefficients learned by NAIS.

Methods	Revised profile or the learned attentions	The target course
HRL+NAIS	<del>Crisis Negotiation, Social Civilization, Web Technology, C++ Program</del>	Web Development
NAIS	Crisis Negotiation(29.61), Social Civilization(29.09), Web Technology(28.32), C++ Program(28.12)	Web Development
HRL+NAIS	<del>Modern Biology, Medical Mystery, Biomedical Imaging, R Program</del>	Biology
NAIS	Modern Biology(37.79), Medical Mystery(37.96), Biomedical Imaging(37.62), R Program(37.84)	Biology
HRL+NAIS	<del>Web Technology, Art Classics, National Unity Theory, Philosophy</del>	Life Aesthetics
NAIS	Web Technology(38.32), Art Classics(35.87), National Unity Theory(40.63), Philosophy(43.69)	Life Aesthetics



# Conclusion

- We present the first attempt to solve the problem of course recommendation in MOOCs platform by a **hierarchical RL model**.
- The model **jointly trains a profile reviser and a basic recommendation model**, which enables the hierarchical RL model effectively to remove the noisy courses to the target course, and enables recommendation model to be improved on revised user profiles by an agent.
- The experimental results on a real dataset collected from XuetangX validate the effectiveness of the proposed model.

---

Thank you !