### DRN: A Deep Reinforcement Learning Framework for News Recommendation

Guanjie Zheng, Fuzheng Zhang, Zihan Zheng, Yang Xiang,

Nicholas Jing Yuan, Xing Xie, Zhenhui (Jessie) Li





# Introduction: Why reinforcement recommendation First round

#### Second round



Equal rewarding recommendation for current round





#### Images from:

- 1. https://www.ohio.com/akron/sports/cavs/cavaliers-j-r-smith-expects-lebron-james-to-keep-opening-night-streak-alive-despite-sprained-ankle
- 2. https://www.cnbc.com/2018/03/05/kobe-bryant-has-won-an-oscar-heres-what-he-says-it-takes-to-succeed.html

3. <u>https://twitter.com/NBA</u>

- 4. <u>https://www.stgeorgeutah.com/news/archive/2016/06/11/djg-alert-significant-weather-strong-thunderstorm-rolls-through-washington-county/#.WsqHUdPwb6Y</u>
  - . http://www.clipartpanda.com/clipart\_images/guestion-mark-36565633

## Introduction: News recommendation is dynamic

The life period for news is usually very short.

User's interest may change during time.



### Introduction: Is there more than click/noclick?

User's clicks on news are usually very dense in a short period. Then, user usually leave the app!

User may return everyday!



Introduction: Should we keep recommending similar items?



Lebron James will be the MVP! Tony Parker has come back from injury! Paul Gasol promises to help the Spurs in the playoff.

Will you **get bored** if all the recommended news are from NBA when you are browsing the sports news?

5

### Method: Using reinforcement learning to do recommendation





## Method: Dueling network structure – value and advantage function



# Method: user activeness modeling -- survival analysis

User activeness decay function

User activeness

$$\lambda(t) = \lim_{dt \to 0} \frac{Pr\{t \le T < t + dt | T \ge t\}}{dt}$$

$$S(t) = e^{-\int_0^t \lambda(x) dx}$$

$$\int_{\substack{0.8 \\ \text{substrained} \\ 0.2 \\ 0.0 \\ \ t_1 \\ \ t_2 \\ \ t_3 \\ \ t_4 \\ \ t_5 \\ \ t_6 \\ \ t_7 \\ \ t_8 \\ \ t_9 \\ \ t_{10} \\ \$$

Time

11/9/18

### Method: Effective exploration

- Step1: get recommendation from Q and  $\tilde{Q}$
- Step2: probabilistic interleave these two lists
- Step3: get feedback from user and compare the performance of two network
- Step4: if  $\tilde{Q}$  performs better, update Q towards it



### Dataset

Stage	Duration	# of users	# of news
Offline stage	6 months	541,337	1,355,344
Online stage	1 month	64,610	157,088



### Results: Offline

Method	CTR	nDCG
LR	0.1262	0.3659
FM	0.1489	0.4338
W&D	0.1554	0.4534
LinUCB	0.1447	0.4173
HLinUCB	0.1194	0.3491
DN	0.1587	0.4671
DDQN	0.1662	0.4877
DDQN + U	0.1662	0.4878
DDQN + U + EG	0.1609	0.4723
DDQN + U + DBGD	0.1663	0.4854



### Results: Online

ethod	CTR	Precision@5	nDCG	Method	
	0.0059	0.0082	0.0326	LR	
	0.0072	0.0078	0.0353	FM	0
	0.0052	0.0067	0.0258	W&D	0
В	0.0075	0.0091	0.0383	LinUCB	0
JCB	0.0085	0.0128	0.0449	HLinUCB	0
	0.0100	0.0135	0.0474	DN	0
Ν	0.0111	0.0139	0.0477	DDQN	0.
QN + U	0.0089	0.0110	0.0425	DDQN + U	0.
QN + U + EG	0.0083	0.0100	0.03391	DDQN + U + EC	6 O.
QN + U + DBGD	0.0113	0.0149	0.0492	DDQN + U + DE	BGD <b>0</b> .

11/9/18

Diversity

### Summary of motivation and solution

Mo	otivation	Solution		
Loi	ng term effect in recommendation	Deep reinforcement learning (DRL)		
•	Dynamic nature of news recommendation	<ul> <li>Online learning feature of DRL</li> </ul>		
•	Consider more measures for long term effect	<ul> <li>Reward function design of DRL</li> </ul>		
•	Recommendation diversity	Explore in DRL		

### Conclusion

- We propose a reinforcement learning framework to do online personalized news recommendation, taking care of both immediate and future reward. Our framework can be generalized to many other recommendation problems.
- We consider user activeness to help improve recommendation accuracy, which can provide extra information than simply using user click labels.
- Our system has been deployed online in a commercial news recommendation application. Extensive offline and online experiments have shown the superior performance of our methods.