

# User Models and Interactive IR

**ESSIR 2022** 

**Romain Deffayet**<sup>1,2</sup> and **Maarten de Rijke**<sup>2</sup> <sup>1</sup>Naver Labs Europe, <sup>2</sup>University of Amsterdam July 19, 2022, 09.00–10.30

r.e.deffayet@uva.nl, m.derijke@uva.nl



Romain Deffayet PhD student at Naver Labs Europe and the University of Amsterdam



Maarten de Rijke University professor at the University of Amsterdam Based on joint work and conversations with Ali Vardasbi, Harrie Oosterhuis, Jean-Michel Renders, Maria Heuss, Shashank Gupta

Materials based in part on (Chuklin et al., 2015; Oosterhuis and de Rijke, 2018; Oosterhuis et al., 2020; Saito and Joachims, 2021)

#### **Objectives**

• We will cover basic concepts and fundamental methods of learning from interactions for search and recommendation

#### Outcomes

• As a result of participating in this tutorial, students will be able to implement and work with basic (counterfactual) learning to rank, bandits and reinforcement learning for search and recommendation

Go to https://irlab.science.uva.nl/2022/07/17/ essir-2022-tutorial-on-user-models-and-interactive-ir/

- Slides (PDF)
- Bibliography (BIB)



## 09.00 Start **09.00–09.05** Domestic matters – Maarten and Romain **09.05–09.20** Setting the scene – Maarten 09.20–09.40 Counterfactual learning-to-rank – Maarten 09.40-10.15 Bandits & Reinforcement learning in IR - Romain **10.15–10.25** Conclusion – Maarten and Romain 10.25–10.30 Final Q&A 10.30 End

## Setting the scene

- Interactions with users our perspective on information retrieval: technology to connect people to information
- Core concepts and examples
- Core problems: learning and evaluating
- Core distinctions: on-policy vs off-policy









How can a search engine or recommender system or conversational assistant get better by interacting with its users?



How can a search engine or recommender system or conversational assistant get better by interacting with its users?

- Context x user history, user profile, query, time of day, ....
- Policy  $\pi$  that selects action a answer, item, result list, ...
- Rewards r that are returned clicks, downloads, purchases, ...

Policy can be ...

- deterministic function from contexts to actions:  $\pi(x) = a$
- stochastic conditional probability of action given context:  $\pi(a \mid x)$

Policy can be ...

- deterministic function from contexts to actions:  $\pi(x) = a$
- stochastic conditional probability of action given context:  $\pi(a \mid x)$

Policy interacts with environment and produces log data  $\mathcal{D} = \{(x_i, a_i, r_i)\}_{i=1}^n$ 

- Observe context  $x x \sim P(x)$
- Select action  $a a \sim \pi(a \mid x)$
- Observe reward  $r p(r \mid x, a)$

 Ad hoc search (not personalized): context – query; action – ranked list of documents; reward – clicks, dwell time

- Ad hoc search (not personalized): context query; action ranked list of documents; reward – clicks, dwell time
- Product search (personalized): context query, user profile, past interactions; action – grid of items; reward – clicks, conversion

- Ad hoc search (not personalized): context query; action ranked list of documents; reward – clicks, dwell time
- Product search (personalized): context query, user profile, past interactions; action – grid of items; reward – clicks, conversion
- Ad placement (personalized): context user profile; action a slate of ads; reward – clicks, conversion

- Ad hoc search (not personalized): context query; action ranked list of documents; reward – clicks, dwell time
- Product search (personalized): context query, user profile, past interactions; action – grid of items; reward – clicks, conversion
- Ad placement (personalized): context user profile; action a slate of ads; reward – clicks, conversion
- Conversational recommendation: context user history, conversation history; action item; reward task completion, conversion

#### **Policy learning**

• Find a new policy that improves upon the current policy

#### **Policy evaluation**

• Determine the quality – often expressed as the (online) performance – of a given policy

In CS, algorithms that receive input sequentially operate in online modality

• Typically includes tasks that involve sequences of decisions, like when you choose how to serve incoming queries in a stream

In CS, algorithms that receive input sequentially operate in online modality

• Typically includes tasks that involve sequences of decisions, like when you choose how to serve incoming queries in a stream

Batch or offline processing does not need human interaction

- E.g., batch learning proceeds as follows:
  - Initialize the weights
  - Repeat the following steps: (Process all the training data; Update the weights)

In CS, algorithms that receive input sequentially operate in online modality

• Typically includes tasks that involve sequences of decisions, like when you choose how to serve incoming queries in a stream

Batch or offline processing does not need human interaction

- E.g., batch learning proceeds as follows:
  - Initialize the weights
  - Repeat the following steps: (Process all the training data; Update the weights)

Typical offline computations in information retrieval:

• Any processing that is not query dependent (crawling, indexing, ...)

In CS, algorithms that receive input sequentially operate in online modality

• Typically includes tasks that involve sequences of decisions, like when you choose how to serve incoming queries in a stream

Batch or offline processing does not need human interaction

- E.g., batch learning proceeds as follows:
  - Initialize the weights
  - Repeat the following steps: (Process all the training data; Update the weights)

Typical offline computations in information retrieval:

• Any processing that is not query dependent (crawling, indexing, ...)

Typical online computations in information retrieval:

• Any processing that depends on users and their input

#### Evaluation and learning can be **on-policy** or **off-policy**

Evaluation and learning can be **on-policy** or **off-policy** 

**On-policy learning algorithms** evaluate and improve the same policy that is being used to select actions

Evaluation and learning can be on-policy or off-policy

**On-policy learning algorithms** evaluate and improve the same policy that is being used to select actions

**Off-policy learning algorithms** evaluate and improve a policy that is different from the policy that is used for action selection

• Behavior or logging policy: policy that tells the agent what action to take; used to collect actions taken and outcomes; not the target policy in off-policy learning

Deploy two policies  $\pi_A$  and  $\pi_B$  to get an online estimate of performance

- Collect log data  $\mathcal{D}_A = \{(x_i, a_i, r_i)\}_{i=1}^n$  and  $\mathcal{D}_B = \{(x_j, a_j, r_j)\}_{j=1}^m$
- Compute quality as average reward:  $\frac{1}{n}\sum_{i=1}^{n} \{r_i : (x_i, a_i, r_i) \in \mathcal{D}_A\}$  for  $\pi_A$  and  $\frac{1}{m}\sum_{j=1}^{m} \{r_j : (x_j, a_j, r_j) \in \mathcal{D}_B\}$  for  $\pi_B$
- Compare the two average rewards

Again, take two policies  $\pi_A$  and  $\pi_B$  but now

- Given context x, determine most probably actions  $a_A$  and  $a_B$
- Combine actions  $a_A$  and  $a_B$  into action  $a_{A\oplus B}$  and determine credit assignment
- Execute combined action a<sub>A⊕B</sub>, observe reward, following credit assignment function assign credit to π<sub>A</sub> or π<sub>B</sub> (or both or neither)
- Repeat, take average, and compare

Again, take two policies  $\pi_A$  and  $\pi_B$  but now

- Given context x, determine most probably actions  $a_A$  and  $a_B$
- Combine actions  $a_A$  and  $a_B$  into action  $a_{A\oplus B}$  and determine credit assignment
- Execute combined action a<sub>A⊕B</sub>, observe reward, following credit assignment function assign credit to π<sub>A</sub> or π<sub>B</sub> (or both or neither)
- Repeat, take average, and compare

Several design choices

• How to combine, how to satisfy constraints on actions, how to assign credit, ...

Why evaluate online?

- User behavior is indicative of their preferences
  - Avoids issues of labeled data: expensive, unethical to create in privacy-sensitive settings, impossible for small-scale problems, stationary, not necessarily aligned with actual user preferences, ...

Why evaluate online?

- User behavior is indicative of their preferences
  - Avoids issues of labeled data: expensive, unethical to create in privacy-sensitive settings, impossible for small-scale problems, stationary, not necessarily aligned with actual user preferences, ...

Why not to evaluate online?

- Online agents take risks to gain knowledge quickly
- Online evaluations are complex

Why evaluate online?

- User behavior is indicative of their preferences
  - Avoids issues of labeled data: expensive, unethical to create in privacy-sensitive settings, impossible for small-scale problems, stationary, not necessarily aligned with actual user preferences, ...

#### Why not to evaluate online?

- Online agents take risks to gain knowledge quickly
- Online evaluations are complex

#### What if we evaluate off-policy?

- Estimate performance of a policy using only log data collected by a behavior policy
- Compare performance of candidate policies safely and help us decide which policy should be deployed

- Counterfactual evaluation = offline A/B testing = off-policy evaluation
- Counterfactual learning = unbiased learning to rank = off-policy learning

Questions, comments, ...

09.00 Start **09.00–09.05** Domestic matters – Maarten and Romain **09.05–09.20** Setting the scene – Maarten 09.20–09.40 Counterfactual learning-to-rank – Maarten 09.40-10.15 Bandits & Reinforcement learning in IR - Romain **10.15–10.25** Conclusion – Maarten and Romain 10.25–10.30 Final Q&A

10.30 End
# **Counterfactual learning-to-rank**

- Asking counterfactual questions: "what would have been ...."
- Importance sampling
- Learning from logs
- Bias, and correcting for bias

## Movie recommendation from ratings



What to recommend next?

## Extrinsic biases in Missing-Not-At-Random (MNAR) feedback

**Popularity bias** Didn't watch **Positivity bias** didn't rate

Users are more likely to rate movies they have liked.

More popular items are rated more often.

u / s	a user / a user state
d / m / y / a	document / movie / slate / action
r / R	an observed reward $/$ a reward estimate
$\mathcal{D}$	a dataset
$\pi$	a policy
$E_d$	the event "the user examines document $d$ "
$\mathbb{E}$	an expected value
$\mathcal{U} \ / \ \mathcal{N} \ / \ \mathcal{B}$	a uniform / normal / Bernoulli distribution

- Estimate a certain quantity  $\theta$  (e.g., mean, variance) of one distribution by sampling from another
- Here we want 𝔼<sub>(u,m)∼U</sub> [r(u, m)] but we only observe biased instances (u<sup>D</sup>, m<sup>D</sup>) ~ D ...
- **Counterfactual question**: what would have been the average ratings under a uniform distribution?

**IS correction**: In the data, what was the probability of observing  $(u^{\mathcal{D}}, m^{\mathcal{D}})$ ? (Wasserman, 2004)



**IS correction**: In the data, what was the probability of observing  $(u^{\mathcal{D}}, m^{\mathcal{D}})$ ? (Wasserman, 2004)



$$\mathbb{E}_{(u^{\mathcal{D}},m^{\mathcal{D}})\sim\mathcal{D}}\left[\frac{r(u^{\mathcal{D}},m^{\mathcal{D}})}{p^{\mathcal{D}}(u^{\mathcal{D}},m^{\mathcal{D}})}\right] = \sum_{(u,m)} \frac{p^{\mathcal{D}}(u,m)}{p^{\mathcal{D}}(u,m)} r(u,m) \propto \mathbb{E}_{(u,m)\sim\mathcal{U}}\left[r(u,m)\right]$$

Formal definition of the counterfactual estimator:

$$\tilde{r} = \frac{1}{N_u \cdot N_m} \sum_{(u^{\mathcal{D}}, m^{\mathcal{D}}) \in \mathcal{D}} \frac{r(u^{\mathcal{D}}, m^{\mathcal{D}})}{p^{\mathcal{D}}(u^{\mathcal{D}}, m^{\mathcal{D}})}$$

• 
$$\mathbb{E}[\tilde{r}] = \mathbb{E}_{(u,m)\sim \mathcal{U}}[r(u,m)]$$
:  $\tilde{r}$  is unbiased

- $\tilde{r}$  can have high variance if  $p^{\mathcal{D}}$  are small
- Tricks to decrease the variance: clipping/normalization (Saito and Joachims, 2021)

#### Learning from search logs: a screenshot

lisbon

Q AI 
Q Maps 
☐ Images 
H News 
Videos 
¡ More

Tools

× J. Q

About 322,000,000 results (0.54 seconds)

https://en.wikipedia.org > wiki > Lisbon

#### Lisbon - Wikipedia



Lisbon is the capital and the largest city of Portugal, with an estimated population of 544,851 within its administrative limits in an area of 100.05 km<sup>2</sup>.

Siege of Lisbon: 1147 CE Country: Portugal Area code(s): (+351) 21 XXX XXXX Historic province: Estremadura

Lisbon Metro - Lisbon Airport - Tourism in Lisbon - Lisbon District

#### People also ask

What is Lisbon famous for?	~
Is Lisbon unsafe?	~
Is Lisbon worth visiting?	~
Is Lisbon a poor city?	Foodback

https://www.visitlisboa.com > ...

#### Visit Lisboa: Lisboa OFFICIAL Site

It's your turn to conquer this monumental castle in the Lisbon region. Take a trip to Palmela to get to know the area and the Arrábida hills which surround it.

https://www.britannica.com > ... > Cities & Towns H-L

#### Lisbon | History, Culture, Economy, & Facts | Britannica

Lisbon, Portuguese Lisboa, city, port, capital of Portugal, and the centre of the Lisbon metropolitan area. Located in western Portugal on the estuary of ...



Lisbon Capital of Portugal

Area: 100 km²

Elevation: 2 m

Weather: 30°C, Wind E at 10 km/h, 35% Humidity weather.com

Local time: Wednesday 10:36

Population: 504,718 (2016) United Nations

Metro population: 2,871,133

#### Plan a trip

Things to do

3-star hotel averaging €116, 5-star averaging €240

h 1 h 15 min flight, from €104

### Learning from search logs: a snapshot

	AOL-user-ct-collection — less — 116×53	H2
710766	www.peoplesearch.comww.reviewplace.seardh 2006-05-30 22:10:13	
710766	www.peoplesearch.comww.reviewplace.seardh 2006-05-30 22:10:33	
711391	can not sleep with snoring husband 2006-03-01 01:24:00	
711391	cannot sleep with snoring husband 2006-03-01 01:24:07 9 http://www.wjla.com	
/11391	cannot sleep with snoring husband 2006-03-01 01:24:07 9 http://www.wjla.com	
/11391	cannot sleep with snoring husband 2006-03-01 01:33:06 1 http://www.epinions.com	
/11391	jackie zeanan nude 2006-03-01 15:26:27	
/11391	jackie zeman nude 2006-03-01 15:26:38	
11391	strange cosmos 2006-03-01 16:07:15 1 http://www.strangecosmos.com	
11391	mansfield first assembly 2006-03-01 16:09:20 1 http://www.mansfieldfirstassembly.org	
/11391	mansfield first assembly 2006-03-01 16:09:20 3 http://netministries.org	
/11391	reverend harry nyers 2006-03-01 16:10:07	
11391	reverend harry nyers 2006-03-01 16:10:30	
11391	National enguirer 2006-03-01 1/:13114 1 http://www.nationalenguirer.com	
11391	NOW TO KILL MOCKINGDINGS 2000-03-01 17:18:11	
11391	NOW TO RILL BOCKINGDINGS 2000-03-01 17:10:33	
11391	Now to kill annoying birds in your yards 2006-03-01 1/18:58	
11391	now to kitt annoying birds in your yards 2000-03-01 17:19:53 2 http://www.sortprice.com	
11391	Now to rid your yard of noisy annoying birds 2006-03-01 1/23108 3 http://shopping.msh.com	
11391	now to rid your yard of noisy annoying birds 2000-03-01 17:23:08 10 http://www.bergen.org	
11391	Now to rid your yard of noisy annoying Dirds 2006-03-01 1/124:35 15 http://www.saterbrand.co	
11391	now do i get mocking birds out of my yard 2000-03-01 1/12/11/	
11391	Now do 1 get mockingbirds out of my yard 2006-03-01 1/12/136 9 http://www.asrl.org	
11391	now do i get mockingDiros out of my yard 2000-03-01 17:30114	
11391	Now to get rid of noisy loud birds 2006-03-01 1/130152 3 http://www.bird-x.com	
11391	now to get rid of holsy loud birds 2000-03-01 1/130152 1 http://forumsz.gardenweb.com	
11391	Now to get rid of noisy loud birds 2006-03-01 1/:30152 10 http://www.birding.com	
11391	mansfield first assembly 2000-03-01 18:31:30 3 http://hetministries.org	
11391	Detrimotre 2000-03-01 19:42:41 1 http://www.tprot.org	
11391	judy daker ministries 2006-03-01 19:49:03 2 http://www.empracinggrace.com	
11391	god witt futfilt your nearts desires 2000-03-01 19:39:00 10 http://www.pureintimacy.org	
11391	untare relevantaja can be very special 2000-03-01 23:09:37	
11391	online friendsnips can be very special 2000-03-0123:09:57	
11391	United Filenuships 2000-03-01 2010024	
11391	Cypress fairbanks iso 2000-03-02 0/:00:03 1 http://www.cfisd.net	
11391	people are not always now they seen over the internet 2000-03-02 00:31:51	
11391	friends untile can be different in person 2000-03-02 00:32:42	
11201	haston butte can be darierent an person 2000-03-02 00:33:04 is nttp://www.satun.com	
11391	computer christian church houston ty 2006-02-02 16:07:52	
11201	community christenn christian two 2006-03-02 10:00/133	
11391	yay churches in houseon (x 2000-03-02 ionoiza)	
11201	to manage the one has a loss of the control to the	
11301	Howston (X is one hot place $2000-05-02$ 10:04:44 howston (X is one hot place to live $2000-05-02$ 10:04:55 0 http://traval.vahoo.com	
11301	houston tx is one hot place to the $2000-03-02$ $18:16:05$ 1 http://net.houston-texe_online.	
11391	texas hill country and sinhts around san antonio ty $2896-83-21$ R:19:08 5 bttp://www.nouscome.exas-ontcher-	
	texas nice councily and signed around san anconico ex- 2000-03-02 10:19:00 5 nice;//www.answe	
11201	can liver problems cause you to loose your bais 2006-02-02 10:27:04	
11391	can liver problems cause you to losse your har 2006-03-02 18:27:64	-
11201	ctrange compared 2002 you to 201 1 http://www.dskuuctiish.u	em)
11391	strange cosmos 2000-03-02 is 29.34 a necessita necessita a necessita necessita a necessita	
11001	white hard dry skin on face $2005 - 03 - 02 - 03 - 02 - 03$	

- Position bias: users are more likely to observe item on top of the page
- **Item-selection bias**: users cannot observe items which are not returned by the earch engine
- **Trust bias**: users may trust the search engine to return relevant results and are therefore more likely to click on top documents

• . . .

Slate: combination of multiple items offered to user at same time

• Search engine result page, playlist, grid with products, carousel, banner ads, ....

Slate: combination of multiple items offered to user at same time

• Search engine result page, playlist, grid with products, carousel, banner ads, ...

**Logging policy**  $\pi_L$ :  $\pi_L(\mathbf{y} \mid u)$  probability that the system chooses slate  $\mathbf{y}$  for user u. Under the logging policy, what was the probability of ...

... that slate being returned? (Precup et al., 2000)

$$R^{\rm IS} = \frac{1}{|\mathcal{D}|} \sum_{(u,\mathbf{y},\mathbf{r})\in\mathcal{D}} \frac{1}{\pi_{\mathcal{L}}(\mathbf{y} \mid u)} \cdot \sum_{j=1}^{k} r^{j}$$

**Logging policy**  $\pi_L$ :  $\pi_L(\mathbf{y} \mid u)$  probability that the system chooses slate  $\mathbf{y}$  for user uUnder the logging policy, what was the probability of ...

... that document being placed at that position in the slate? (McInerney et al., 2020)

$$R^{\text{pos-IS}} = \frac{1}{|\mathcal{D}|} \sum_{(u, \mathbf{y}, \mathbf{r}) \in \mathcal{D}} \sum_{j=1}^{k} \frac{1}{\frac{\pi_{L}^{j}(y_{j} \mid u)}{\frac{\pi_{L}^{j}(y_{j} \mid u)}{\frac{\pi$$

**Examination hypothesis:** A clicked document is both **examined** and **relevant** Under the logging policy, what was the probability of . . .

... the user examining that document? (Joachims et al., 2017)

$$R^{\text{IPS}} = \frac{1}{|\mathcal{D}|} \sum_{(u, \mathbf{y}, \mathbf{r}) \in \mathcal{D}} \sum_{d \in \mathbf{y}} \underbrace{\frac{1}{\mathcal{P}(E_d = 1 \mid u)}}_{\text{examination prob.}} \cdot r^d$$

For example, position-based model:  $P(E_d = 1 \mid u, \operatorname{rank}(d) = j) = \gamma_j$ 

Improvements to handle item-selection and trust bias: policy-aware and affine estimators (Oosterhuis and de Rijke, 2021).

- Direct (biased) estimate of the user response:  $\tilde{r}(d|u)$
- Simple idea: get a rough, biased estimate of user response and apply the correction only on the difference with this estimate (less variance-related risk)

$$R^{\text{DRE}} = \frac{1}{|\mathcal{D}|} \sum_{(u,\mathbf{y},\mathbf{r})\in\mathcal{D}} \sum_{d\in\mathbf{y}} \underbrace{\left[ \underbrace{\tilde{r}(d|u)}_{\text{high-bias, low-variance}} + \underbrace{\frac{1}{P(E_d = 1|u)} \cdot (r^d}_{\text{unbiased, high-variance}} - \tilde{r}(d|u)) \right]}_{\text{unbiased, low-variance}}$$



# Eye tracking experiments

(Joachims et al., 2007)

### Fitting importance weights: 3 solutions



```
doc1 CTR=0.8

doc2 CTR=0.3

doc2 CTR=0.6

doc1 CTR=0.4
```

Eye tracking experiments (Joachims et al., 2007)

Swap interventions (Joachims et al., 2017)

### Fitting importance weights: 3 solutions



CTR=0.8 doc doc2 CTR=0.3 CTR=0.6 doca CTR=0.4 doc



Eye tracking experiments (Joachims et al., 2007)

Swap interventions (Joachims et al., 2017)

Learning from logs: Click models

# Click models (Chuklin et al., 2015)

- Interpretable structure with latent variables and parameterized causal relations
- Learn parameters by Expectation-Maximization or Gradient Descent

#### **Example: Dynamic Bayesian Networks**



#### **Evaluation of click models**

#### doc1 click p=0.8 doc2 p=0.4 skip doc3 click p=0.5 doc4 p=0.1 skip Perplexity: $PPL@i = 2^{-\sum_{i=1}^{n} c_{j}^{i} \log_{2}(\tilde{p}_{j}^{i}) + (1 - c_{j}^{i}) \log_{2}(1 - \tilde{p}_{j}^{i})}$ $PPL = \frac{1}{k} \sum_{i=1}^{k} PPL@j \approx 1.51$

**Click prediction** 

What to do in Lisbon?



**Relevance estimation** 

- We want to learn new policy from user feedback but **historical data contains biases**
- Importance Sampling applies a **correction to observed feedback** to recover unbiased estimates
- IS is unbiased but suffers from **high variance**: we need to leverage user **models** (inverse propensity scoring, doubly robust estimator)
- Propensity weights must be **adequately computed and evaluated** (eye tracking, swap interventions, click models)

Questions, comments, ...

# 09.00 Start **09.00–09.05** Domestic matters – Maarten and Romain **09.05–09.20** Setting the scene – Maarten 09.20–09.40 Counterfactual learning-to-rank – Maarten 09.40-10.15 Bandits & Reinforcement learning in IR - Romain **10.15–10.25** Conclusion – Maarten and Romain 10.25-10.30 Final Q&A 10.30 End

# Bandits & Reinforcement learning in IR

- Learning to interact with new users
- More complex feedback loops
- Long-term satisfaction
- Learning from logs

Cold start



How to quickly find Maarten's and Romain's preferences?

For every action a in the set A, we define:

- the reward  $r(a) \in \{0, 1\}$ : like or dislike,
- the **regret**  $\overline{r}(a) = r(a^*) r(a)$  with  $a^* = \arg \max_{a \in \mathcal{A}} \mathbb{E}[r(a)]$ .

The goal can now be formulated as finding a strategy minimizing the **expected cumulative regret**  $\overline{R}$ :

$$\overline{R}(T) = \mathbb{E}\left[\sum_{t=1}^{T} \overline{r}(a_t)\right]$$

Past rewards :

Past actions ;



#### Greedy

#### **Exploration-Exploitation dilemma:**



#### Greedy

#### **Exploration-Exploitation dilemma:**



#### Greedy

#### **Exploration-Exploitation dilemma:**



#### Greedy

#### **Exploration-Exploitation dilemma:**



#### Greedy

#### **Exploration-Exploitation dilemma:**



#### **Exploration-Exploitation dilemma**:



#### **Exploration-Exploitation dilemma**:
# Multi-Armed Bandits (MAB): greedy and *e*-greedy



### **Exploration-Exploitation dilemma**:

*exploiting* the knowledge acquired from interactions lowers the regret ... ... but *exploring* different actions multiple times is required to find the optimal action

# Multi-Armed Bandits (MAB): greedy and *e*-greedy



### **Exploration-Exploitation dilemma**:

*exploiting* the knowledge acquired from interactions lowers the regret ... ... but *exploring* different actions multiple times is required to find the optimal action

# Multi-Armed Bandits (MAB): greedy and *e*-greedy



### **Exploration-Exploitation dilemma**:

*exploiting* the knowledge acquired from interactions lowers the regret ... ... but *exploring* different actions multiple times is required to find the optimal action

### The expected cumulative regret admits a logarithmic lower bound:



It is impossible to find a MAB algorithm with bounded expected cumulative regret!











## Short digression on Bayesian inference

How to maintain an estimate of the full distribution of rewards, and update it after having interacted ?

Use parameterized families of distributions for the rewards, like  $\mathcal{N}(\mu, \sigma^2)$  or  $\mathcal{B}(\mu)$ .



- **Probability matching**: we want to select an action according to its probability of being optimal.
- Thompson Sampling does this by
  - 1. sampling a parameter value for each action
  - 2. selecting the best action under the chosen parameters
  - 3. observing the reward and updating the corresponding action

- Rewards are clicks/skips:  $r(a) \sim \mathcal{B}(\mu_a)$
- Prior and posterior distributions on  $\mu_a$ ?  $\rightarrow$  **Conjugate prior** is  $\mu_a \sim B(\alpha_a, \beta_a)$
- Simple update:  $\alpha_a^{t+1} = \alpha_a^t + r$  and  $\beta_a^{t+1} = \beta_a^t + (1-r) \ (\alpha_a^0 = \beta_a^0 = 0)$



- Rewards are clicks/skips:  $r(a) \sim \mathcal{B}(\mu_a)$
- Prior and posterior distributions on  $\mu_a$ ?  $\rightarrow$  **Conjugate prior** is  $\mu_a \sim B(\alpha_a, \beta_a)$
- Simple update:  $\alpha_a^{t+1} = \alpha_a^t + r$  and  $\beta_a^{t+1} = \beta_a^t + (1-r) \ (\alpha_a^0 = \beta_a^0 = 0)$



- Rewards are clicks/skips:  $r(a) \sim \mathcal{B}(\mu_a)$
- Prior and posterior distributions on  $\mu_a$ ?  $\rightarrow$  **Conjugate prior** is  $\mu_a \sim B(\alpha_a, \beta_a)$
- Simple update:  $\alpha_a^{t+1} = \alpha_a^t + r$  and  $\beta_a^{t+1} = \beta_a^t + (1-r) \ (\alpha_a^0 = \beta_a^0 = 0)$



- Rewards are clicks/skips:  $r(a) \sim \mathcal{B}(\mu_a)$
- Prior and posterior distributions on  $\mu_a$ ?  $\rightarrow$  **Conjugate prior** is  $\mu_a \sim B(\alpha_a, \beta_a)$
- Simple update:  $\alpha_a^{t+1} = \alpha_a^t + r$  and  $\beta_a^{t+1} = \beta_a^t + (1-r) \ (\alpha_a^0 = \beta_a^0 = 0)$



- Rewards are clicks/skips:  $r(a) \sim \mathcal{B}(\mu_a)$
- Prior and posterior distributions on  $\mu_a$ ?  $\rightarrow$  **Conjugate prior** is  $\mu_a \sim B(\alpha_a, \beta_a)$
- Simple update:  $\alpha_a^{t+1} = \alpha_a^t + r$  and  $\beta_a^{t+1} = \beta_a^t + (1-r) \ (\alpha_a^0 = \beta_a^0 = 0)$



## What about dwell-time ?



How to quickly find Maarten's and Romain's preferences?

Rewards are positive real numbers sampled from a normal distribution with fixed variance: r(a) ~ N(μ<sub>a</sub>, σ<sup>2</sup>)

 $\rightarrow \sigma$  can be interpreted as an exploration parameter.

• Prior and posterior distributions on  $\mu_a? \rightarrow$  **Conjugate prior** is  $\mu_a \sim \mathcal{N}(\nu_a, \frac{\sigma^2}{\lambda^2})$ 

• Update: 
$$(\lambda_a^{t+1})^2 = (\lambda_a^t)^2 + \frac{1}{\sigma^2}$$
 and  $\mu_a^{t+1} \left[ 1 + \frac{1}{\sigma^2 (\lambda_a^t)^2} \right] = \mu_a^t + \frac{r}{\sigma^2 (\lambda_a^t)^2}$ 

Normal-Normal Thompson Sampling achieves a logarithmic lower bound! (Agrawal and Goyal, 2017)

## Contextual click/dwell-time maximization



What to recommend next to Maarten and Romain?



We still have no prior knowledge about the current user ...

... but we use knowledge from previous users!

- Rewards are sampled from a normal distribution with fixed variance and where the mean is a linear combination of context features: r<sup>t</sup>(a) ~ N(X<sub>t</sub><sup>T</sup> μ<sub>a</sub>, σ<sup>2</sup>).
- Prior and posterior distributions on  $\mu_a$ ?

 $\rightarrow$  **Conjugate prior** is a multivariate normal distribution  $\mu_a \sim \mathcal{N}(\nu_a, \sigma^2 \cdot \Lambda_a^{-1})$ .

• Update:  $\Lambda_a^{t+1} = \Lambda_a^t + X_t^T X_t$  and  $\Lambda_a^{t+1} \nu_a^{t+1} = \Lambda_a^t \nu_a^t + r_t \cdot X_t$ .

## Contextual Bandits: LinTS extensions (Riquelme et al., 2018)

$$\sigma_a^2 \sim \Gamma^{-1}(\alpha_a, \beta_a)$$
$$\mu_a \mid \sigma_a^2 \sim \mathcal{N}(\nu_a, \sigma_a^2 \cdot \Lambda_a^{-1})$$
$$\alpha_a^{t+1} = \alpha_a^t + 1/2$$
$$\beta_a^{t+1} = \beta_a^t + 1/2 \left( r_t^2 - \mu_a^t {}^T \Lambda_a^t \mu_a^t \right)$$

### Unknown variance



### **Non-linear rewards**

## Long-term user engagement



How to satisfy Maarten and Romain on the long run?

The objectives match (maximize sum of rewards / minimize cumulative regret) ...

... but the methods we used assume that actions are independent of future rewards conditioned on the current context, i.e., **the user is static** 

 $\rightarrow$  need explicitly capturing causal effect of recommendations on future user states: Reinforcement Learning

## **MDPs and POMDPs**

### Partially-observable Markov Decision Process:

- States  $s \in S$ : user's mind
- Observations  $o \in \mathcal{O}$ : history, extrinsic context, ...
- Actions  $a \in \mathcal{A}$ : recommendations
- Reward function  $r : S \times A \mapsto \mathbb{R}$ : click, dwell-time, etc.
- Transition probabilities T(s'|s, a): how recommendations influence the user
- Initial state distribution  $S(s_1)$ : user state when they arrive on the platform
- Observation probabilities  $\Omega(o|s, a)$ : how the user's mind is revealed

### Goal: Maximize the expected cumulative rewards

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=1}^{T} r(s_t, a_t) \right] \quad , \tau = (s_1, a_1, \dots, s_T, a_T)$$

# RL101: Policy gradients (Sutton and Barto, 2018)

• Objective function 
$$J(\pi_{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}} \left[ \sum_{t=1}^{T} r(s_t, a_t) \right]$$

• Policy gradient theorem:

$$abla_{ heta} J(\pi_{ heta}) = \mathbb{E}_{t \sim \pi_{ heta}} \left[ \sum_{t=1}^{T} 
abla_{ heta} \log \pi_{ heta}(a_t | s_t) \sum_{t'=t}^{T} r(s_{t'}, a_{t'}) 
ight]$$

• **REINFORCE** algorithm alternate between two steps:

Collect a trajectory  $\tau^k$  from  $\pi_{\theta^k} \leftrightarrow$  improve the policy  $\theta^{k+1} = \theta^k + \alpha \nabla_{\theta} J(\pi_{\theta}) \mid_{\theta = \theta^k}$ 

# RL101: Dynamic programming (Sutton and Barto, 2018)

- Q-function:  $Q^{\pi}(a|s) = \mathbb{E}_{\tau \sim \pi, s_1 = s, a_1 = a} \left[ \sum_{t=1}^{T} r(s_t, a_t) \right]$  $\rightarrow$  how good action *a* is in state *s*
- Bellman Equation:  $Q^{\pi}(a|s) = r(s,a) + \mathbb{E}_{a' \sim \pi(\cdot|s')} \left[Q(a'|s')\right]$
- Q-Learning alternates between two steps:

Collect experience (s, a, r, s') from an  $\epsilon$ -greedy policy w.r.t  $Q \leftrightarrow$ improve the Q-function  $Q(s, a) \leftarrow (1 - \alpha)Q(s, a) + \alpha [r + \max_{a' \in \mathcal{A}} Q(s', a')]$ 

• Actor Critic combines policy gradient and dynamic programming (actually slightly more complex than that ...)

- Online interactions are expensive : we cannot just return random results to users during training!
- In large scale RS or search engines, we must train from logs while accounting for the interactivity of system and user ...

## Learning from logs & session optimization

#### lisbon

Q AI O Maps E Images E News T Videos I More

0 Tools

×

About 322,000,000 results (0.54 seconds)

https://en.wikipedia.org.).wiki > Lishon

#### Lisbon - Wikipedia

Lisbon is the capital and the largest city of Portugal, with an estimated population of 544 851 within its administrative limits in an area of 100 05 km<sup>2</sup>

Siege of Lisbon: 1147 CE Country: Portugal

Area code(s): (+351) 21 XXX XXXX Historic province: Estremadura

Lishon Metro - Lishon Aimort - Tourism in Lishon - Lishon District

#### People also ask

What is Lisbon famous for?	`
Is Lisbon unsafe?	`
Is Lisbon worth visiting?	· · · · · · · · · · · · · · · · · · ·
Is Lisbon a poor city?	

#### https://www.visitlisboa.com > ... 1

#### Visit Lisboa: Lisboa OFFICIAL Site

It's your turn to conquer this monumental castle in the Lisbon region. Take a trip to Palmela to get to know the area and the Arrábida hills which surround it.

#### https://www.britannica.com > .... > Cities & Towns H-L



Lisbon, Portuguese Lisboa, city, port, capital of Portugal, and the centre of the Lisbon metropolitan area. Located in western Portugal on the estuary of ...



#### Lisbon Capital of Portugal

Lisbon is Portugal's hilly, coastal capital city. From imposing São Jorne Castle, the view encompasses the old city's pastel-colored buildings, Tagus Estuary and Ponte 25 de Abril suspension bridge Nearby, the National Azuleio Museum displays 5 centuries of decorative ceramic files. Just outside Lisbon is a string of Atlantic heaches, from Cascais to Estoril - Google

#### Area: 100 km<sup>a</sup>

Elevation: 2 m

Weather: 30°C. Wind E at 10 km/h. 35% Humidity weather com Local time: Wednesday 10:36 Population: 504 718 (2016) United National

Metro population: 2.871.133

#### Plan a trin

Things to do

- 3-star hotel averaging €116, 5-star averaging €240 1±100
- $\mathbf{+}$ 1 h 15 min flight, from €104

We need to cover many aspects of the query and anticipate for future needs

Learning from logged interactions, without interventions, is hard:

- Optimizer's curse: A maximization process is likely to select an overestimated solution. With n items of expected rewards r<sub>1</sub>,..., r<sub>n</sub>, we can have E[r̂<sub>k</sub>] = r<sub>k</sub> and yet E[r̂<sub>k\*</sub> r<sub>k\*</sub>] > 0 with k\* = arg max<sub>k</sub> r̂<sub>k</sub>
  - $\rightarrow$  We will be disappointed. (Jeunen and Goethals, 2021)
- **Deadly Triad**: Optimizer's curse is much worse when 3 conditions are satisfied: (van Hasselt et al., 2018)
  - Off-Policy training
  - Dynamic programming
  - Q-function approximation

Under the logging policy...

... what would have been the probability of observing that sequence of rankings ?

$$R^{IS} = \frac{1}{|\mathcal{D}|} \sum_{\tau \in \mathcal{D}} \sum_{t=1}^{T} \prod_{i=1}^{t} \frac{\pi(a_i|s_i)}{\pi_L(a_i|s_i)} r_t$$
product of past IS weights

Even worse than in traditional CLTR!

 $\rightarrow$  variance grows exponentially with horizon length.

**Offline RL**: Find a way to recover policies staying within the support of the logging policy (Levine, 2021)

**Offline RL** 



Pessimistic Q-Functions

(Kumar et al., 2020)

Filtered Behavior cloning (Chen et al., 2021)

- The environment is highly stochastic: two users with same history may react differently to a recommendation
- There are a lot of actions: from thousands to billions and above
- How to integrate user models?

- We must **balance exploration and exploitation** to find user preferences quickly and reliably (Bandits)
- We can augment bandits algorithms with **context features to leverage knowledge from other users** (Contextual Bandits)
- We must capture the **effect of recommendations on future user behavior** to enable long-term satisfaction (Reinforcement Learning)

Questions, comments, ...
## 09.00 Start **09.00–09.05** Domestic matters – Maarten and Romain **09.05–09.20** Setting the scene – Maarten 09.20–09.40 Counterfactual learning-to-rank – Maarten 09.40-10.15 Bandits & Reinforcement learning in IR - Romain 10.15–10.25 Conclusion – Maarten and Romain 10.25-10.30 Final Q&A 10.30 End



- Taking stock
- Directions not covered
- Challenges

- Using user interactions to evaluate or optimize interactive systems
  - Online vs off-policy
  - Counterfactual evaluation / learning
- Counterfactual learning to rank
- Bandits and reinforcement learning

#### What we have not covered

- Recent advances in bias-variance trade-offs
- Complex (very large) action spaces
- Tuning hyperparameters
- Working with multiple logging policies
- Dealing with distributional shifts
- Combinations of online & offline, with occasional online exploration to collect new data (Oosterhuis and de Rijke, 2021)
- Limitations of de-biasing in counterfactual learning to rank (Oosterhuis, 2022)
- Simulation environments
- Libraries and packages

Guarantees on ...

- Accuracy, also for rare phenomena
- Efficiency during both training and inference
- Reliability when assumptions begin to fail (e.g., on user behavior)
- Reproducibility of experimental results
- Resilience against distributional shifts and adversarial attacks
- Safety of user data and proprietary data

## 09.00 Start **09.00–09.05** Domestic matters – Maarten and Romain **09.05–09.20** Setting the scene – Maarten 09.20–09.40 Counterfactual learning-to-rank – Maarten 09.40-10.15 Bandits & Reinforcement learning in IR - Romain 10.15–10.25 Conclusion – Maarten and Romain 10.25-10.30 Final Q&A 10.30 End

# Final Q&A

Questions, comments, ...

#### References i

- S. Agrawal and N. Goyal. Thompson sampling for contextual bandits with linear payoffs. In Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28, ICML'13, pages III–1220–III–1228. JMLR.org, 2013.
- S. Agrawal and N. Goyal. Near-optimal regret bounds for thompson sampling. J. ACM, 64(5), 2017.
- O. Chapelle and L. Li. An empirical evaluation of thompson sampling. In *Proceedings of the 24th International Conference on Neural Information Processing Systems*, pages 2249–2257. Curran Associates Inc., 2011.
- L. Chen, K. Lu, A. Rajeswaran, K. Lee, A. Grover, M. Laskin, P. Abbeel, A. Srinivas, and I. Mordatch. Decision transformer: Reinforcement learning via sequence modeling. In *NeurIPS 2021: Conference* on Neural Information Processing Systems, 2021.
- A. Chuklin, I. Markov, and M. de Rijke. *Click Models for Web Search*. Synthesis Lectures on Information Concepts, Retrieval, and Services. Morgan & Claypool Publishers, August 2015.
- O. Jeunen and B. Goethals. Pessimistic reward models for off-policy learning in recommendation. In *Fifteenth ACM Conference on Recommender Systems*, pages 63–74. ACM, 2021.

#### References ii

- T. Joachims, L. Granka, B. Pan, H. Hembrooke, F. Radlinski, and G. Gay. Evaluating the accuracy of implicit feedback from clicks and query reformulations in web search. ACM Trans. Inf. Syst., 25(2): 7–es, 2007.
- T. Joachims, A. Swaminathan, and T. Schnabel. Unbiased learning-to-rank with biased user feedback. In *Proceedings of the Tenth ACM International Conference on Web Search and Data Mining*, pages 781–789. ACM, 2017.
- A. Kumar, A. Zhou, G. Tucker, and S. Levine. Conservative q-learning for offline reinforcement learning. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*. Curran Associates Inc., 2020.
- S. Levine. A gentle introduction to offline reinforcement learning. https://www.youtube.com/watch?v=tW-BNW1ApN8, 2021.
- J. McInerney, B. Brost, P. Chandar, R. Mehrotra, and B. A. Carterette. Counterfactual evaluation of slate recommendations with sequential reward interactions. In KDD '20: The 26th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, Virtual Event, CA, USA, August 23-27, 2020, pages 1779–1788. ACM, 2020.

#### References iii

- H. Oosterhuis. Reaching the end of unbiasedness: Uncovering implicit limitations of click-based learning to rank. *arXiv preprint arXiv:2206.12204*, 2022.
- H. Oosterhuis and M. de Rijke. Differentiable unbiased online learning to rank. In CIKM 2018: International Conference on Information and Knowledge Management, pages 1293–1302. ACM, October 2018.
- H. Oosterhuis and M. de Rijke. Unifying online and counterfactual learning to rank. In WSDM 2021: 14th International Conference on Web Search and Data Mining. ACM, March 2021.
- H. Oosterhuis, R. Jagerman, and M. de Rijke. Unbiased learning to rank: Counterfactual and online approaches. In *Companion Proceedings of the Web Conference 2020*, pages 299–300. ACM, April 2020.
- D. Precup, R. S. Sutton, and S. P. Singh. Eligibility traces for off-policy policy evaluation. In Proceedings of the Seventeenth International Conference on Machine Learning, ICML '00, pages 759–766, San Francisco, CA, USA, 2000. Morgan Kaufmann Publishers Inc.
- C. Riquelme, G. Tucker, and J. Snoek. Deep bayesian bandits showdown: An empirical comparison of bayesian deep networks for thompson sampling. In *ICLR 2018*, 2018.

#### **References** iv

- Y. Saito and T. Joachims. Counterfactual learning and evaluation for recommender systems: Foundations, implementations, and recent advances. In *Fifteenth ACM Conference on Recommender Systems*, pages 828–830. ACM, 2021.
- F. Sarvi, M. Heuss, M. Aliannejadi, S. Schelter, and M. de Rijke. Understanding and mitigating the effect of outliers in fair ranking. In *WSDM 2022: The Fifteenth International Conference on Web Search and Data Mining.* ACM, February 2022.
- D. Silver. RL Course Lecture 9: Exploitation and exploration. http://youtube.com, 2022.
- R. Sutton and A. Barto. Reinforcement Learning : An Introduction. MIT Press, 2018.
- H. van Hasselt, Y. Doron, F. Strub, M. Hessel, N. Sonnerat, and J. Modayil. Deep reinforcement learning and the deadly triad. *arXiv preprint arXiv:1812.02648*, 2018.
- L. A. Wasserman. All of Statistics. Springer, 2004.



# User Models and Interactive IR

**ESSIR 2022** 

**Romain Deffayet**<sup>1,2</sup> and **Maarten de Rijke**<sup>2</sup> <sup>1</sup>Naver Labs Europe, <sup>2</sup>University of Amsterdam July 19, 2022, 09.00–10.30

r.e.deffayet@uva.nl, m.derijke@uva.nl