



Tutorial
Morning (March, 8)
9:00-12:00, GMT+2

Beyond Probability Ranking Principle: Modeling the Dependencies among Documents

Liang Pang¹, Qingyao Ai², Jun Xu³

1. CAS Key Lab of Network Data Science and Technology

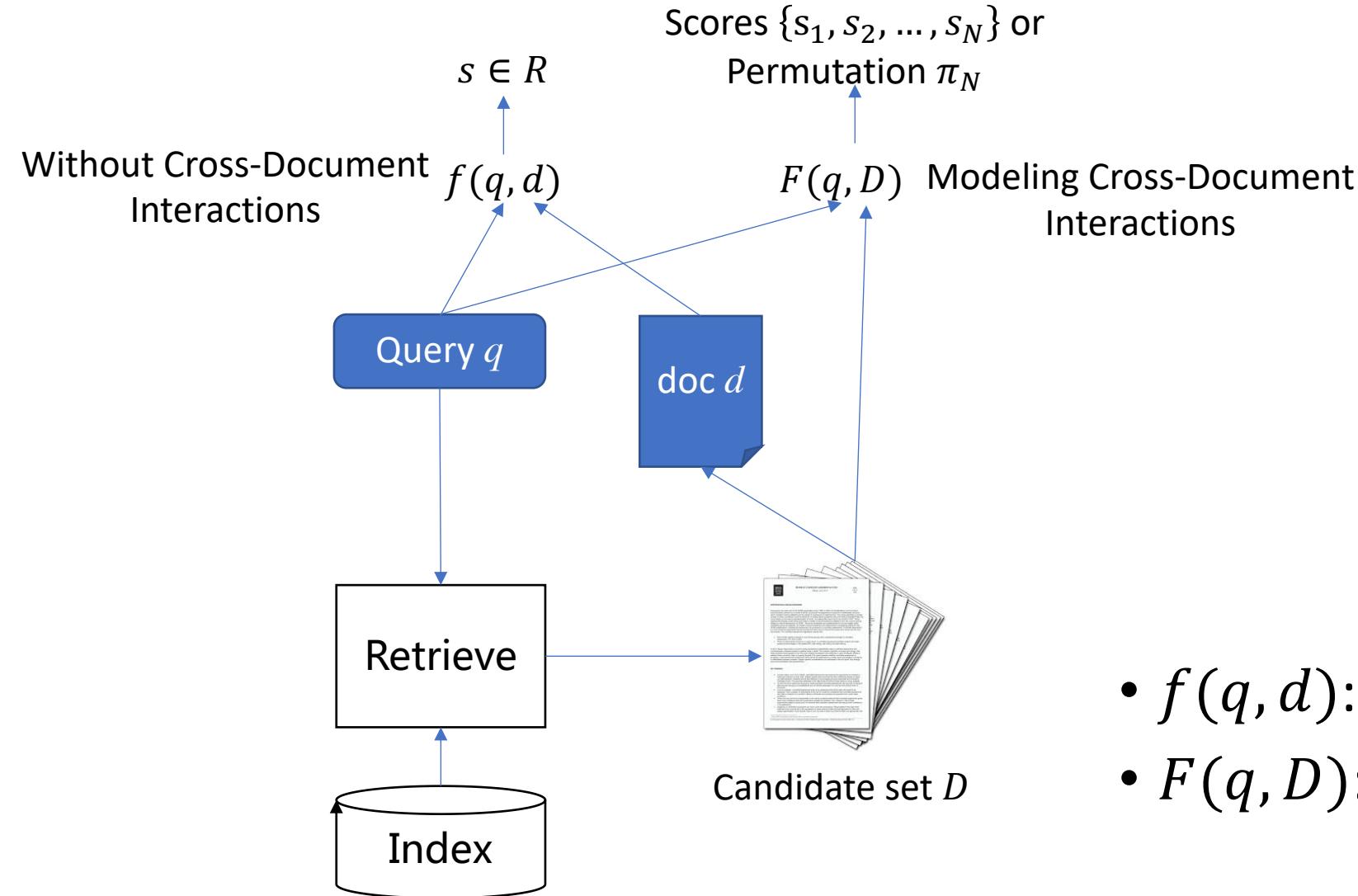
Institute of Computing Technology, Chinese Academy of Sciences, Beijing, China

2. The University of Utah, USA

3. Renmin University of China, Beijing, China



I Modeling the Dependency in IR



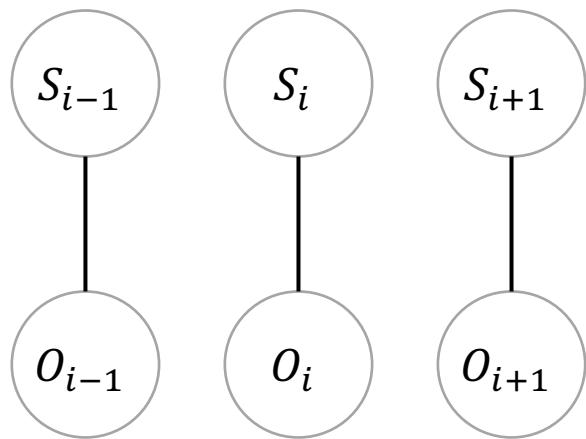
The screenshot shows a search results page for "CCIR2020" with the following details:

- CCIR2020** (Search bar)
- ALL** (Selected tab)
- IMAGES**
- VIDEOS**
- 274,000 Results**
- Any time**
- 全国信息检索学术会议 CCIR2020**
www.cvnis.net/ccir2020/index.html • Translate this page
最新通知: ccir2020微信群, qq群现可加入, 点击查看详情。CCIR Poster模板ppt已上传, 请到征稿界面下载。会议简介-CCIR2020
- CCIR 2020 - Eventegg.com**
https://eventegg.com/ccir •
With a considerable program covering a expansive array of matters such as Climate Change, Scientific Evidence, Human Impacts and Divergent Ecosystems, CCIR 2020 will be totally a must-attend event. ...
- CCIR-2020 | CCIR**
https://www.ccir.it/author/CCIR-2020 •
CCIR-2020. 461 POSTS 0 COMMENTS . Form di iscrizione al Webinar: INDUSTRIA MECCANICA-STRATEGIE DI MARKETING E DIGITALIZZAZIONE... luglio 2, 2020. RICERCA DI UN COORDINATORE DIDATTICO PER LA SCUOLA ITALIANA "ITALO CALVINO"... giugno 22, 2020. 18 GIUGNO 2020, ONLINE. giugno 18, 2020 ...
- CCIR-2020 | CCIR | Page 47**
https://www.ccir.it/author/CCIR-2020/page/47 •
CCIR-2020. 461 POSTS 0 COMMENTS . PRODUZIONE DI PATATE IN RUSSIA. gennaio 6, 2014. 1... 45 46 47 Page 47 of 47. Corso Semiponte, 32/B 20154 Milano. tel.: +39 02 86995240 info@ccir.it. CF 80076750159 – P.IVA 12058220158. SERVIZI IN EVIDENZA. Servizio di ricerca di partner russi;
- Home Page [www.ccir-ccrra.org]**
https://www.ccir-ccrra.org •
The Canadian Council of Insurance Regulators (CCIR) is an inter-jurisdictional association of insurance regulators. The mandate of the CCIR is to facilitate and promote an efficient and effective insurance regulatory system in Canada to serve the public interest
- Cadastro Rural - CCIR**
www.incri.gov.br/pt/cadastro-rural-ccir.html • Translate this page
O Certificado de Cadastro do Imóvel Rural (CCIR) é o documento expedido pelo Incra que comprova a regularidade cadastral do imóvel rural. O certificado contém informações sobre o titular, a área, a localização, a exploração e a classificação fundiária do imóvel rural.

- $f(q, d)$: independent model
- $F(q, D)$: dependent model

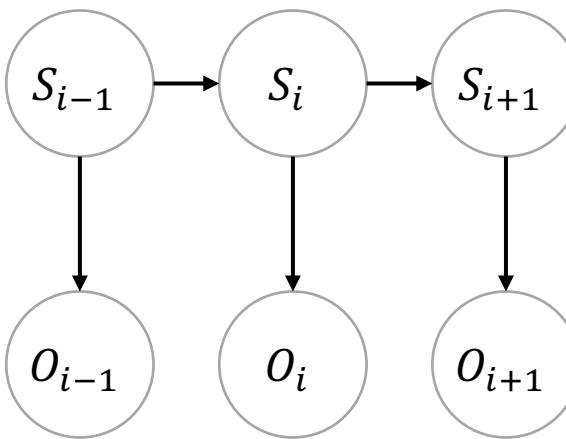
| Modeling the Sequential Dependency: Experiences from Sequence Prediction

No Dependency



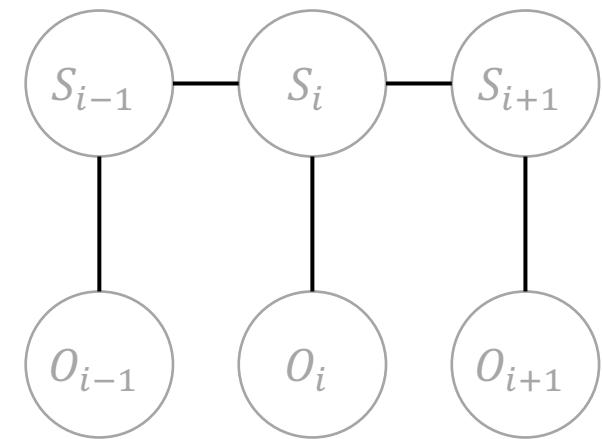
Logistic Regression

Sequential Dependency



HMM

Global Dependency



CRF



THE 14TH ACM INTERNATIONAL CONFERENCE ON
WEB SEARCH AND DATA MINING

Tutorial
Morning (March, 8)
9:00-12:00, GMT+2

4. Ranking with Sequential Dependency

4.1 Heuristic Sequential Ranking Models

4.2 Learning Sequential Ranking Models

4.3 Challenges

Search Result Diversification

Query: jaguar

A screenshot of a Google search results page for the query "jaguar". The results are diverse, including:

- Market Selector | Jaguar | View the site in your preferred language (<https://www.jaguar.com/>)
- Jaguar (@Jaguar) - Twitter (<https://twitter.com/Jaguar>)
- Jaguar UK: Luxury Sports Cars, Executive Saloons and SUVs (<https://www.jaguar.co.uk/>)
- JAGUAR HONG KONG (www.jaguar.com.hk/)
- Images for jaguar
- Jaguar - Home | Facebook (<https://www.facebook.com/jaguar>)
- Jaguar - YouTube (<https://www.youtube.com/user/JaguarCarsLimited>)
- Jaguar MENA: Explore Jaguar the High Performance Luxury Cars (<https://www.jaguar-me.com/>)
- Jaguar Las Vegas | New & Used Car Dealer Las Vegas, NV (www.jaguarlv.com/)

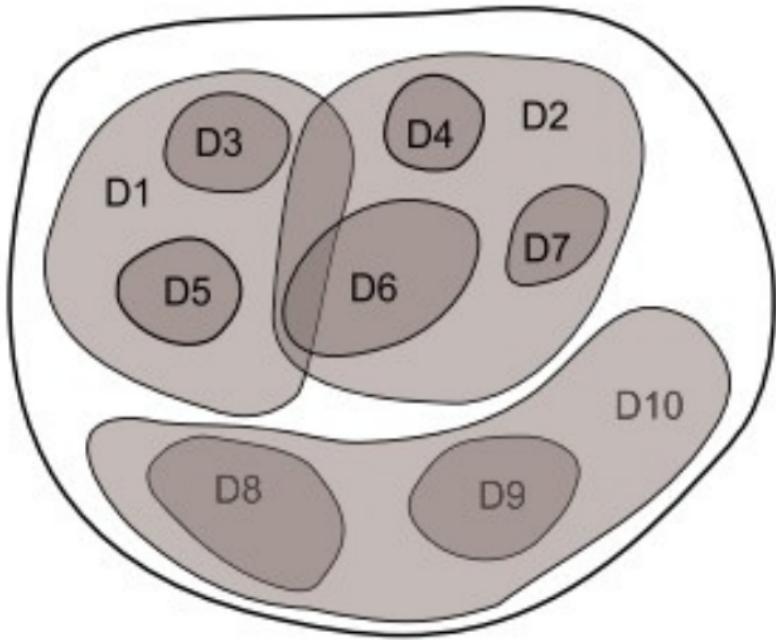
A screenshot of a Google search results page for the query "jaguar". The results are diverse, including:

- Market Selector | Jaguar | View the site in your preferred language (<https://www.jaguar.com/>)
- Jaguar (@Jaguar) - Twitter (<https://twitter.com/Jaguar>)
- Jaguar - Wikipedia (<https://en.wikipedia.org/wiki/Jaguar>)
- Images for jaguar
- Fender American Pro Jaguar®, Rosewood Fingerboard, 3-Color ... (shop.fender.com/en-US/electric-guitars/jaguar.../jaguar.../jaguar0114010700.html)
- Jaguar | jaguarswisswatches.com/
- Brands > Jaguar - MENRAD (<https://www.menrad.de/en/collection/jaguar/>)
- Jaguar Mining Inc.: Home (<https://www.jaguarmining.com/>)
- Jaguar Cars - Wikipedia (https://en.wikipedia.org/wiki/Jaguar_Cars)

- Different user needs
 - Ambiguous queries , e.g., Apple, Jaguar, Band
 - Multi-faceted needs, e.g., Britney spears (news, videos, photos...)
- Information redundancy
 - Many duplicate or similar results

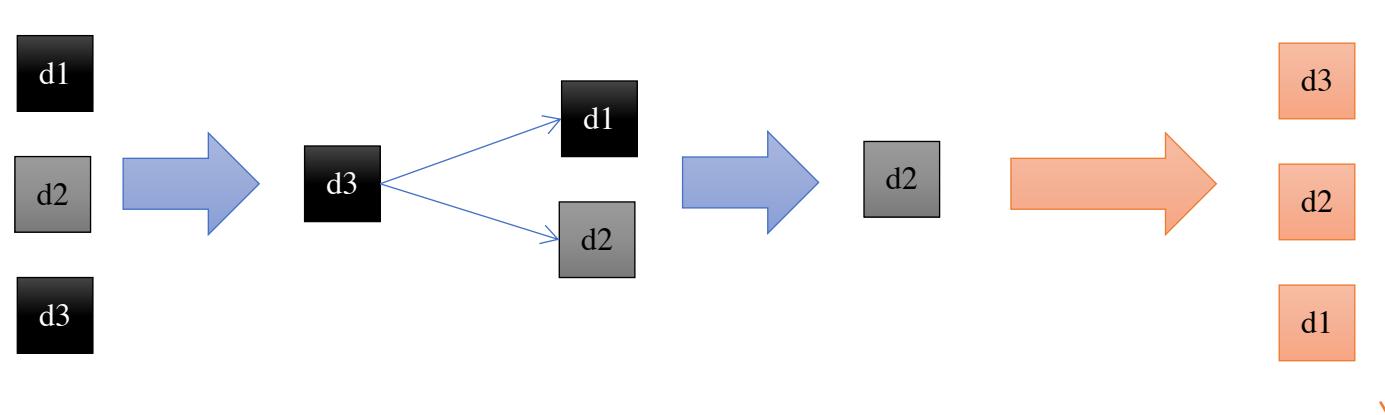
Luxury car
Animal
Electric
Swiss
Eyewear
Mining Inc.

| Greedy approach to diverse ranking



Select k documents to cover more subtopics

subset selection

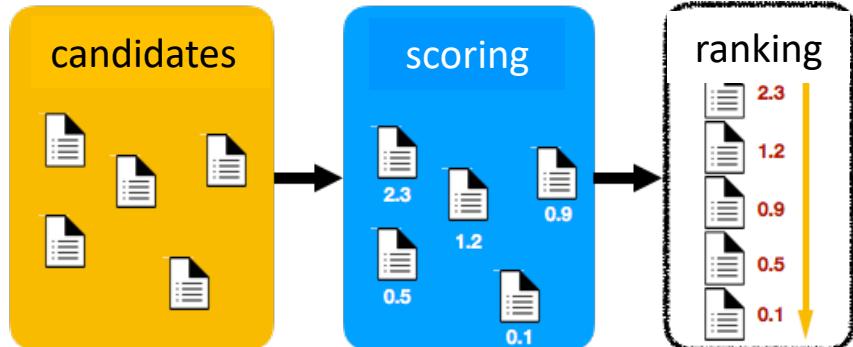


Diversify(q, \mathcal{R}_q, τ)

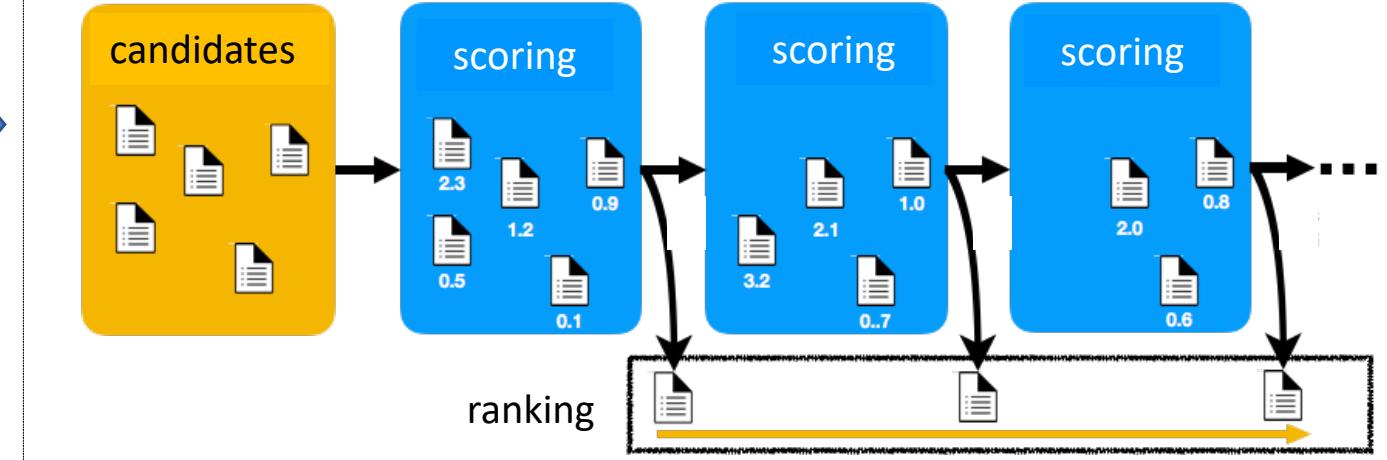
```
1  $\mathcal{D}_q \leftarrow \emptyset$ 
2 while  $|\mathcal{D}_q| < \tau$  do
3    $d^* \leftarrow \arg \max_{d \in \mathcal{R}_q \setminus \mathcal{D}_q} f(q, d, \mathcal{D}_q)$ 
4    $\mathcal{R}_q \leftarrow \mathcal{R}_q \setminus \{d^*\}$ 
5    $\mathcal{D}_q \leftarrow \mathcal{D}_q \cup \{d^*\}$ 
6 end while
7 return  $\mathcal{D}_q$ 
```

| Greedy Approach Models the Dependency Sequentially

Independent Scoring Function



Ranking with Sequential Dependency



Obeys PRP

Beyond PRP

| Pre-defined $f(q, d, D_q)$: MMR (Carbonell and Goldstein, 1998)

- Maximal marginal relevance
 - An implicit diversification model
- Key: defining $f(q, d, D_q)$

Marginal relevance of d given the query q and selected documents D_q

Prefer dissimilar documents

$$f_{\text{MMR}}(q, d, D_q) = \lambda f_1(q, d) - (1 - \lambda) \max_{d_j \in D_q} f_2(d, d_j)$$

Relevance between q and d

aggregation

Similarity between d and a selected document d_j
(dependency between target d and d_j)

Diversify(q, \mathcal{R}_q, τ)

```
1  $\mathcal{D}_q \leftarrow \emptyset$ 
2 while  $|\mathcal{D}_q| < \tau$  do
3    $d^* \leftarrow \arg \max_{d \in \mathcal{R}_q \setminus \mathcal{D}_q} f(q, d, \mathcal{D}_q)$ 
4    $\mathcal{R}_q \leftarrow \mathcal{R}_q \setminus \{d^*\}$ 
5    $\mathcal{D}_q \leftarrow \mathcal{D}_q \cup \{d^*\}$ 
6 end while
7 return  $\mathcal{D}_q$ 
```

| Pre-defined $f(q, d, D_q)$

- Implicit Diversification

- Risk minimization RM (Zhai et al., 2008)

$$f_{\text{RM}}(q, d, \mathcal{D}_q) = f_1(\theta_q, \theta_d)(1 - \lambda - f_2(\theta_d, \theta_{\mathcal{D}_q}))$$

- Explicit Diversification

- Intent-aware selection (IA-Select) (Agrawal et al., 2009)

$$f_{\text{IA-Select}}(q, d, \mathcal{D}_q) = \sum_{c \in \mathcal{T}} f(c|q, \mathcal{D}_q) f(d|q, c)$$

- Proportionality model (PM-2) (Dang and Croft, 2012)

$$f_{\text{PM-2}}(q, d, \mathcal{D}_q) = \sum_{s \in \mathcal{S}_q} b_s \zeta(s|q) p(d|s)$$

- Explicit query aspect diversification (xQuAD) (Santos et al., 2010)

$$f_{\text{xQuAD}}(q, d, \mathcal{D}_q) = (1 - \lambda) P(d|q) + \lambda \sum_{q_i \in Q} \left[P(q_i|q) P(d|q_i) \prod_{d_j \in S} (1 - P(d_j|q_i)) \right]$$

Diversify(q, \mathcal{R}_q, τ)

```
1  $\mathcal{D}_q \leftarrow \emptyset$ 
2 while  $|\mathcal{D}_q| < \tau$  do
3    $d^* \leftarrow \arg \max_{d \in \mathcal{R}_q \setminus \mathcal{D}_q} f(q, d, \mathcal{D}_q)$ 
4    $\mathcal{R}_q \leftarrow \mathcal{R}_q \setminus \{d^*\}$ 
5    $\mathcal{D}_q \leftarrow \mathcal{D}_q \cup \{d^*\}$ 
6 end while
7 return  $\mathcal{D}_q$ 
```



THE 14TH ACM INTERNATIONAL CONFERENCE ON
WEB SEARCH AND DATA MINING

Tutorial
Morning (March, 8)
9:00-12:00, GMT+2

4. Ranking with Sequential Dependency

4.1 Heuristic Sequential Ranking Models

4.2 Learning Sequential Ranking Models

4.3 Challenges

I Learning $f(q, d, D_q)$: R-LTR (Zhu et al., 2014)

- Learning implicit diversification model $f(q, d, D_q)$ from labeled data
- A learnable extension of MMR

$$f_{\text{MMR}}(q, d, \mathcal{D}_q) = \lambda f_1(q, d) - (1 - \lambda) \max_{d_j \in \mathcal{D}_q} f_2(d, d_j)$$

$$f_{\text{R-LTR}}(q, d, D_q) = \langle \mathbf{w}_r, \phi(q, d) \rangle + \langle \mathbf{w}_d, \mathbf{h}_{D_q}(R_d) \rangle$$

Features for describing the relevance
(traditional learning to rank features)

Aggregated features for describing
novelty (dependency) of d given
the selected documents in D_q

Diversify(q, \mathcal{R}_q, τ)

```

1    $\mathcal{D}_q \leftarrow \emptyset$ 
2   while  $|\mathcal{D}_q| < \tau$  do
3        $d^* \leftarrow \arg \max_{d \in \mathcal{R}_q \setminus \mathcal{D}_q} f(q, d, \mathcal{D}_q)$ 
4        $\mathcal{R}_q \leftarrow \mathcal{R}_q \setminus \{d^*\}$ 
5        $\mathcal{D}_q \leftarrow \mathcal{D}_q \cup \{d^*\}$ 
6   end while
7   return  $\mathcal{D}_q$ 

```

$$\mathbf{h}_{D_q}(R_d) = \left\{ \max_{d_j \in D_q} R_{dj1}, \dots, \max_{x_j \in D_q} R_{djK} \right\}$$

Dependency features of d and d_j [$R_{dj1}, R_{dj2}, \dots, R_{djK}$]

- Estimating the parameters \mathbf{w}_r and \mathbf{w}_d with list-wise MLE (*Plackett-Luce based Probability*) or Perceptron (PAMM) (Xia et al., 2015)

• PLSA topic difference: $R_{dj1} = \sqrt{\sum_{k=1}^m (p(z_k|d) - p(z_k|d_j))^2}$

• Text diversity: $R_{dj2} = 1 - \frac{\mathbf{d} \cdot \mathbf{d}_j}{\|\mathbf{d}\| \|\mathbf{d}_j\|}$

I Learning $f(q, d, D_q)$

- SVM-DIV (Yue et al., 2008)
 - Focus on diversity only (subtopic coverage)
 - Defining f as a linear function on features
$$f_{\text{SVM DIV}}(q, d, D_q) = \langle \mathbf{w}^T, \Psi(q, D_q \cup \{d\}) \rangle$$
 - Estimating \mathbf{w} based on structural SVMs
 - Determinantal Point Process (DPP) for diverse ranking (Chen et al., 2018)
 - Selecting item via greedy submodular maximization
 - f is defined based on the item kernel matrix
$$j = \arg \max_{i \in Z \setminus Y_g} f(Y_g \cup \{i\}) - f(Y_g)$$
- See also (Wilhelm et al., 2018)

Algorithm 1 Greedy subset selection by maximizing weighted word coverage

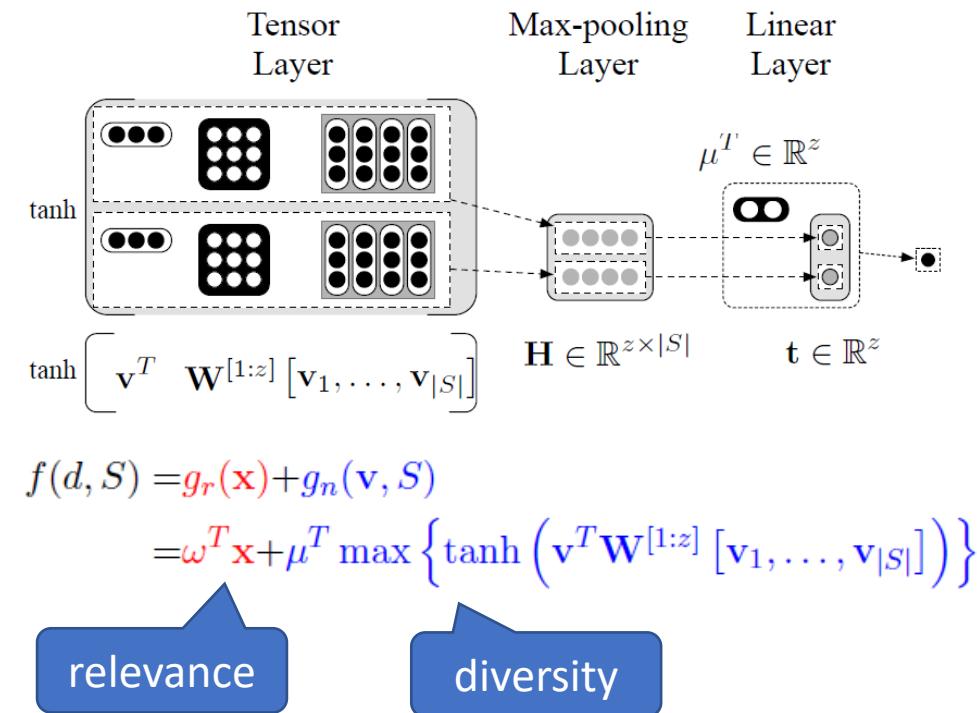
```
1: Input:  $\mathbf{w}, \mathbf{x}$ 
2: Initialize solution  $\hat{\mathbf{y}} \leftarrow \emptyset$ 
3: for  $k = 1, \dots, K$  do
4:    $\hat{x} \leftarrow \operatorname{argmax}_{x: x \notin \hat{\mathbf{y}}} \mathbf{w}^T \Psi(\mathbf{x}, \hat{\mathbf{y}} \cup \{d\})$ 
5:    $\hat{\mathbf{y}} \leftarrow \hat{\mathbf{y}} \cup \{\hat{x}\}$ 
6: end for
7: return  $\hat{\mathbf{y}}$ 
```

Algorithm 1 Fast Greedy MAP Inference

```
1: Input: Kernel  $\mathbf{L}$ , stopping criteria
2: Initialize:  $\mathbf{c}_i = [], d_i^2 = \mathbf{L}_{ii}, j = \arg \max_{i \in Z} \log(d_i^2), Y_g = \{j\}$ 
3: while stopping criteria not satisfied do
4:   for  $i \in Z \setminus Y_g$  do
5:      $e_i = (\mathbf{L}_{ji} - \langle \mathbf{c}_j, \mathbf{c}_i \rangle) / d_j$ 
6:      $\mathbf{c}_i = [\mathbf{c}_i \quad e_i], d_i^2 = d_i^2 - e_i^2$ 
7:   end for
8:    $j = \arg \max_{i \in Z \setminus Y_g} \log(d_i^2), Y_g = Y_g \cup \{j\}$ 
9: end while
10: Return:  $Y_g$ 
```

I Learning deep $f(q, d, D_q)$

- Neural Tensor Network for representing diversity (Xia et al., 2016)
 - Implicit diversification based on R-LTR and PAMM
 - Modeling the novelty part of f with a neural tensor network (without linear and bias terms)
- Document Sequence with Subtopic Attention, DSSA (Jiang et al., 2017)
 - Explicit Diversification
 - Use RNNs with attention to model novelty part of f
 - Training by minimizing list pairwise loss



$$\begin{aligned} S_{\text{DSSA}}(q, d_t, C_{t-1}, \mathcal{I}_q) &= s_{d_t} = \\ (1 - \lambda) S^{\text{rel}}(\mathbf{v}_{d_t}, \mathbf{v}_q) + &\quad \Rightarrow \text{relevance} \\ \lambda S^{\text{div}}\left(\mathbf{v}_{d_t}, \mathbf{v}_{i_{(\cdot)}}, \underbrace{\mathcal{A}\left(\mathcal{H}([\mathbf{v}_{d_1}, \dots, \mathbf{v}_{d_{t-1}}]), \mathbf{v}_{i_{(\cdot)}}\right)}_{\text{subtopic attention}}\right), &\quad \Rightarrow \text{diversity} \end{aligned}$$

I Learning deep $f(q, d, D_q)$ (cont')

- DVGAN (Liu et al., 2020)
 - Explicit & implicit diversification
 - Use RNNs with attention to model novelty part of f (same as in DSSA)
 - Training the parameters with GAN
 - Generator: DSSA
 - Discriminator: R-LTR

$$f_{\theta}(d_t|q, S) = (1 - \lambda)S^{\text{rel}}(d_t, q) + \lambda \sum_{i \in I_q} A(i|S) * S^{\text{sub}}(d_t, i),$$

$$S^{\text{rel}}(d_t, q) = \mathcal{S}(e_{d_t}, e_q) + w_r^T(g) * x_{d_t, q},$$

$$S^{\text{sub}}(d_t, i_k) = \mathcal{S}(e_{d_t}, e_{i_k}) + w_r^T(g) * x_{d_t, i_k},$$

Diversify(q, \mathcal{R}_q, τ)

```
1   $\mathcal{D}_q \leftarrow \emptyset$ 
2  while  $|\mathcal{D}_q| < \tau$  do
3       $d^* \leftarrow \arg \max_{d \in \mathcal{R}_q \setminus \mathcal{D}_q} f(q, d, \mathcal{D}_q)$ 
4       $\mathcal{R}_q \leftarrow \mathcal{R}_q \setminus \{d^*\}$ 
5       $\mathcal{D}_q \leftarrow \mathcal{D}_q \cup \{d^*\}$ 
6  end while
7  return  $\mathcal{D}_q$ 
```

Inference by the learned generator

Multi-Slot Bandits for Ranking

Multi-armed bandit (MAB)

Arm 1	Arm 2	Arm 3
<ul style="list-style-type: none">• Current average reward 70%	<ul style="list-style-type: none">• Current average reward 55%	<ul style="list-style-type: none">• Current average reward 65%

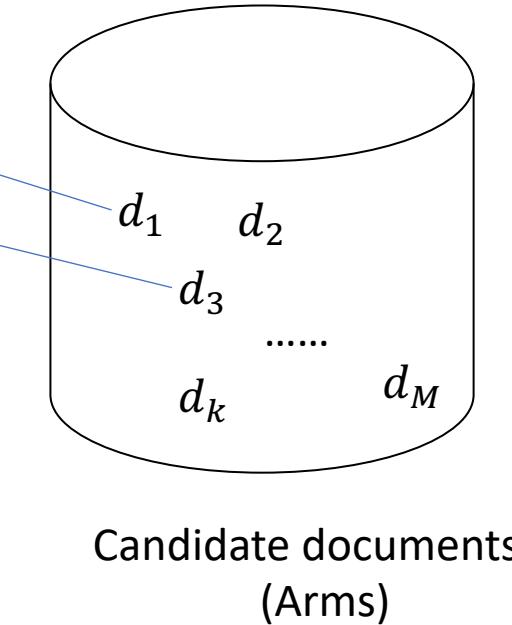


Next choice?

- K arms, unknown sequence of *stochastic* rewards
- For each round $t = 1, 2, \dots,$
 - Choose arm A_t
 - Obtain reward

rank	doc
1	MAB ₁
2	MAB ₂
...	...
...	...
k	MAB _k

Ranked document list



- Multi-slot MAB for ranking
 - MAB_i: model for the i -th rank
 - Arm j : the j -th document
 - Reward: online collecting user clicks
 - Goal: maximize the expected clicks

Ranked Bandit Algorithm [Radlinski et al., ICML '08]

- Runs an MAB instance at each rank
 - MAB₁ is responsible for choosing document shown at rank 1
 - MAB₂ is responsible for choosing document shown at rank 2
 - ... until top K documents are selected
- Show top K to users and receive response
 - If click slot i
 - MAB_i gets reward 1
 - MAB_j ($j < i$) get rewards 0
- Update MABs according to the received feedback

Document selection
for k positions

Update bandits

Algorithm 2 Ranked Bandits Algorithm

```
1: initialize MAB1(n), ..., MABk(n)           Initialize MABs
2: for t = 1 ... T do
3:   for i = 1 ... k do                         Sequentially select documents
4:      $\hat{b}_i(t) \leftarrow$  select-arm (MABi)
5:     if  $\hat{b}_i(t) \in \{b_1(t), \dots, b_{i-1}(t)\}$  then      Replace repeats
6:        $b_i(t) \leftarrow$  arbitrary unselected document
7:     else
8:        $b_i(t) \leftarrow \hat{b}_i(t)$ 
9:     end if
10:   end for
11:   display  $\{b_1(t), \dots, b_k(t)\}$  to user; record clicks
12:   for i = 1 ... k do                          Determine feedback for MABi
13:     if user clicked  $b_i(t)$  and  $\hat{b}_i(t) = b_i(t)$  then
14:        $f_{it} = 1$ 
15:     else
16:        $f_{it} = 0$ 
17:     end if
18:     update (MABi, arm =  $\hat{b}_i(t)$ , reward =  $f_{it}$ )
19:   end for
20: end for
```

Problem Abstraction: Learning of Assignments

[Streeter et al., NIPS '09]

- Ranking: selecting an assignment of documents to positions
 - K positions: K sets of candidate documents
 - Dependences described by *submodular functions*
 - Repeatedly executing the greedy selection step

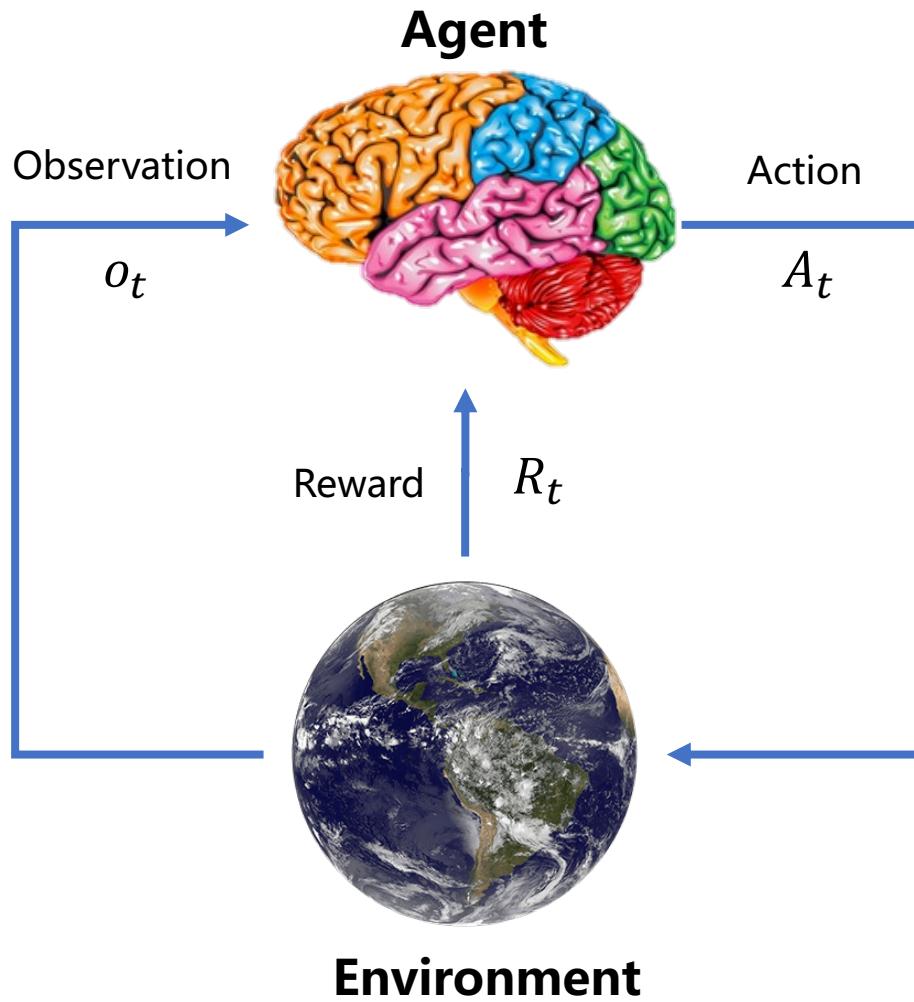
Algorithm: TABULARGREEDY

Input: integer C , sets P_1, P_2, \dots, P_K , function $f : 2^{\mathcal{V}} \rightarrow \mathbb{R}_{\geq 0}$ (where $\mathcal{V} = \bigcup_{k=1}^K P_k$)

```
set  $G := \emptyset$ .
for  $c$  from 1 to  $C$  do
    for  $k$  from 1 to  $K$  do
        set  $g_{k,c} = \arg \max_{x \in P_k \times \{c\}} \{F(G + x)\}$       /* For each color */
        set  $G := G + g_{k,c}$ ;                                /* For each partition */
    /* Greedily pick  $g_{k,c}$  */
for each  $k \in [K]$ , choose  $c_k$  uniformly at random from  $[C]$ .
return sample $_{\vec{c}}(G)$ , where  $\vec{c} := (c_1, \dots, c_K)$ .
```

I Markov Decision Process (MDP)

Markov Decision Process (MDP)



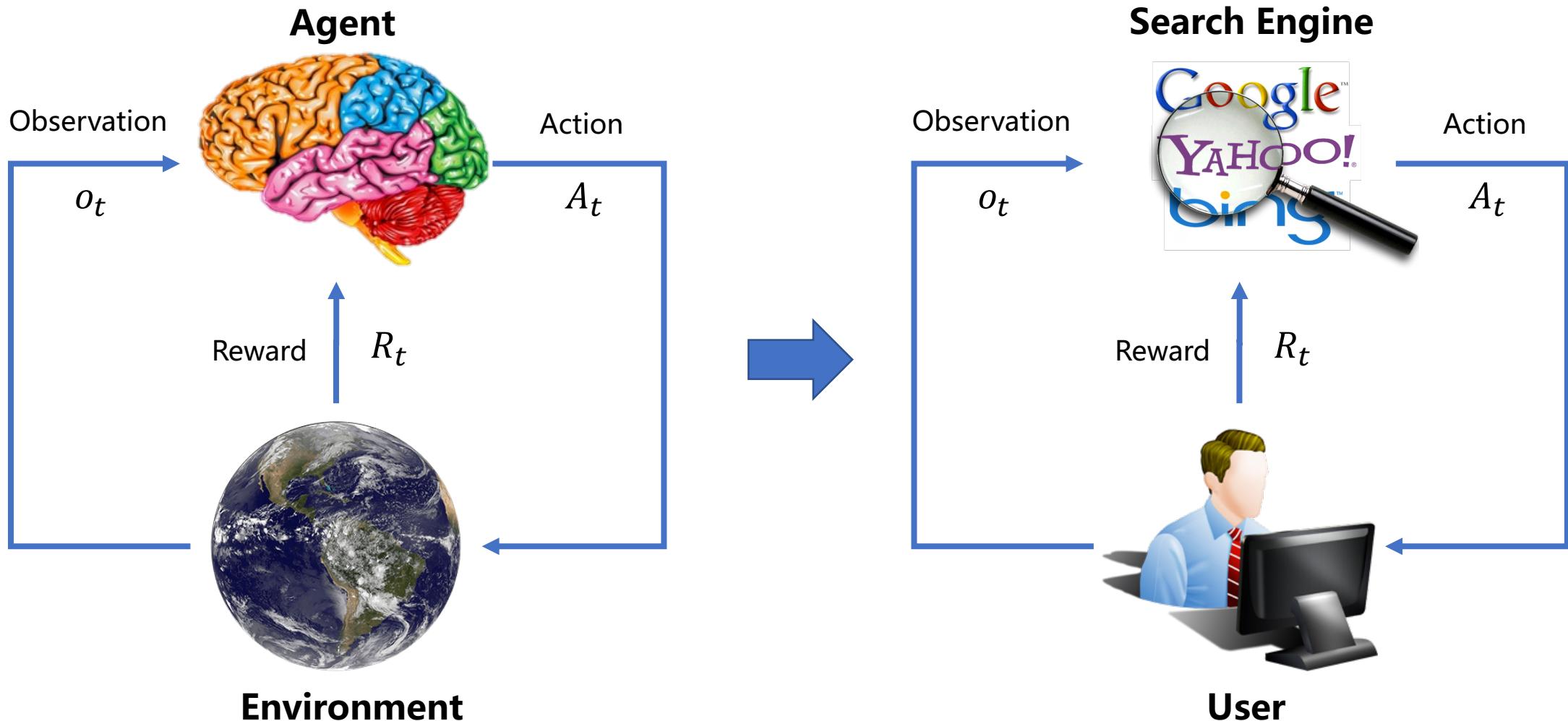
- An **agent** with the capacity to **act**
- Each **action** influences the **environment state**
- Success is measured by a scalar **reward** signal

Goal: select actions to maximize future reward

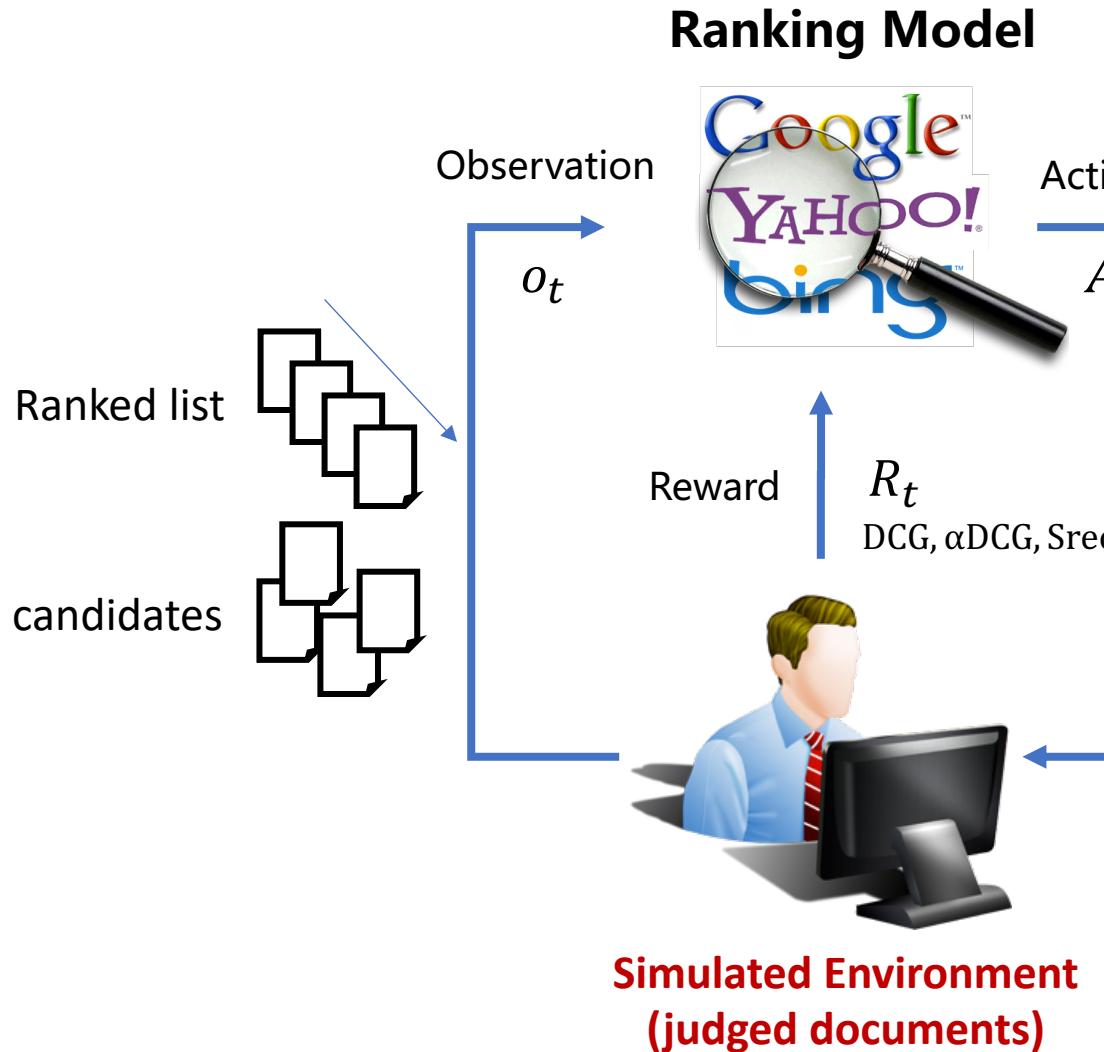
- At each step t the agent:
 - Executes action A_t
 - Receives observation o_t
 - Receives scalar reward r_t
- The environment:
 - Receives action A_t
 - Emits observation o_{t+1}
 - Emits scalar reward r_{t+1}

I Modeling Sequential Dependency with MDP

Markov Decision Process (MDP)



I Modeling Sequential Dependency with MDP



- **Action A_t :**
 - Selects a document and puts ranking list
- **Observation o_t :**
 - query, top t ranked list, candidate set
- **Reward R_t :**
 - designed based on rank evaluation measures

MDP configurations for diverse ranking [Xia et al., 2017]

$\mathbf{x}_{m(a_t)}$: document embedding

MDP factors	Corresponding diverse ranking factors
Time steps	The ranking positions
State	$s_t = [Z_t, X_t, \mathbf{h}_t]$
Policy	$\pi(a_t s_t = [Z_t, X_t, \mathbf{h}_t]) = \frac{\exp\{\mathbf{x}_{m(a_t)}^T \mathbf{U} \mathbf{h}_t\}}{Z}$
Action	Selecting a doc and placing it to rank $t + 1$
Reward	Based on evaluation measure αDCG, SRecall etc. For example: $R = \alpha \text{DCG}[t + 1] - \alpha \text{DCG}[t];$
State Transition	$\begin{aligned} s_{t+1} &= T(s_t = [Z_t, X_t, \mathbf{h}_t], a_t) \\ &= [Z_t \oplus \{\mathbf{x}_{m(a_t)}\}, X_t \setminus \{\mathbf{x}_{m(a_t)}\}, \sigma(\mathbf{V} \mathbf{x}_{m(a_t)} + \mathbf{W} \mathbf{h}_t)] \end{aligned}$

- States $s_t = [Z_t, X_t, \mathbf{h}_t]$
 - Z_t : sequence of t preceding documents, $Z_0 = \phi$
 - X_t : set of candidate documents,
 - $X_0 = X$
 - $\mathbf{h}_t \in R^K$: latent vector
 - initialized with query: $\mathbf{h}_0 = \sigma(\mathbf{V}_q \mathbf{q})$

I Learning with Policy Gradient

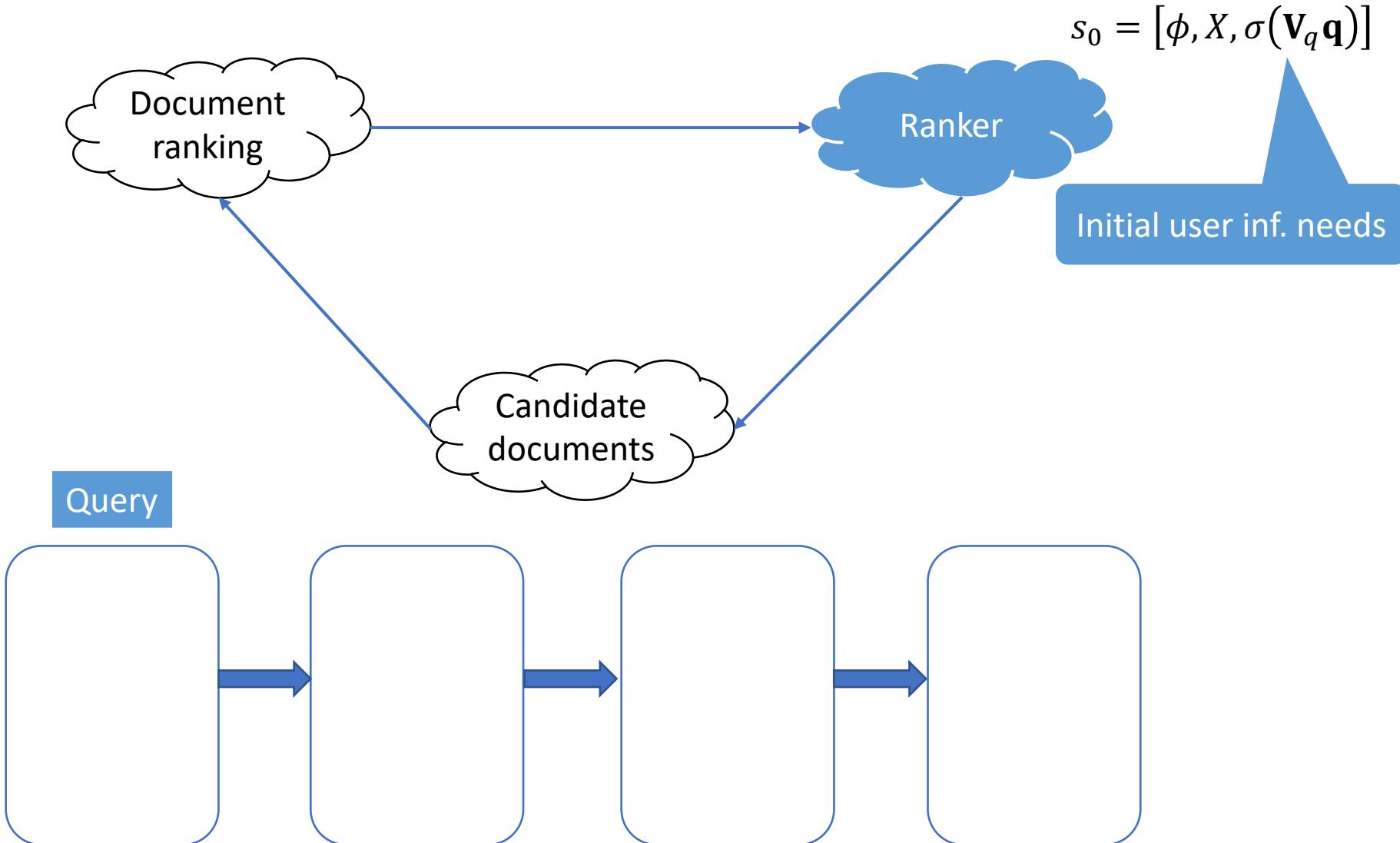
- Model parameters $\Theta = \{\mathbf{V}_q, \mathbf{U}, \mathbf{V}, \mathbf{W}\}$
- Learning objective: maximizing expected return (discounted sum of rewards) of each training query

$$\max_{\Theta} v(\mathbf{q}) = E_{\pi} G_0 = E_{\pi} \left[\sum_{k=0}^{M-1} \gamma^k r_{k+1} \right]$$

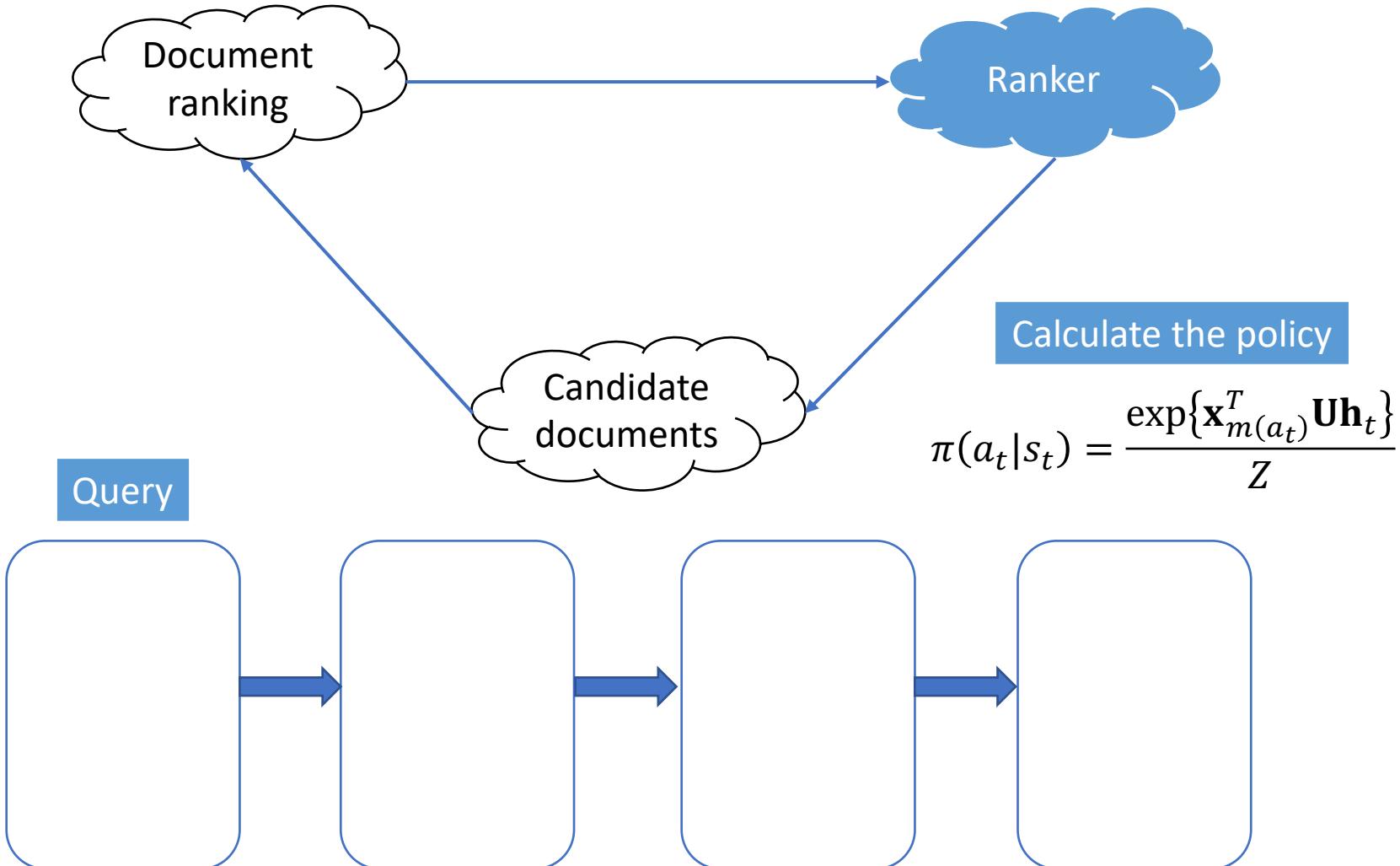
- Directly optimizes evaluation measure as $G_0 = \alpha \text{DCG}@M$
- Monte-Carlo stochastic gradient ascent is used to conduct the optimization (REINFORCE algorithm)

$$\widehat{\nabla_{\Theta} v(\mathbf{q})} = \gamma^t G_t \nabla_{\Theta} \log \pi(a_t | s_t; \Theta)$$

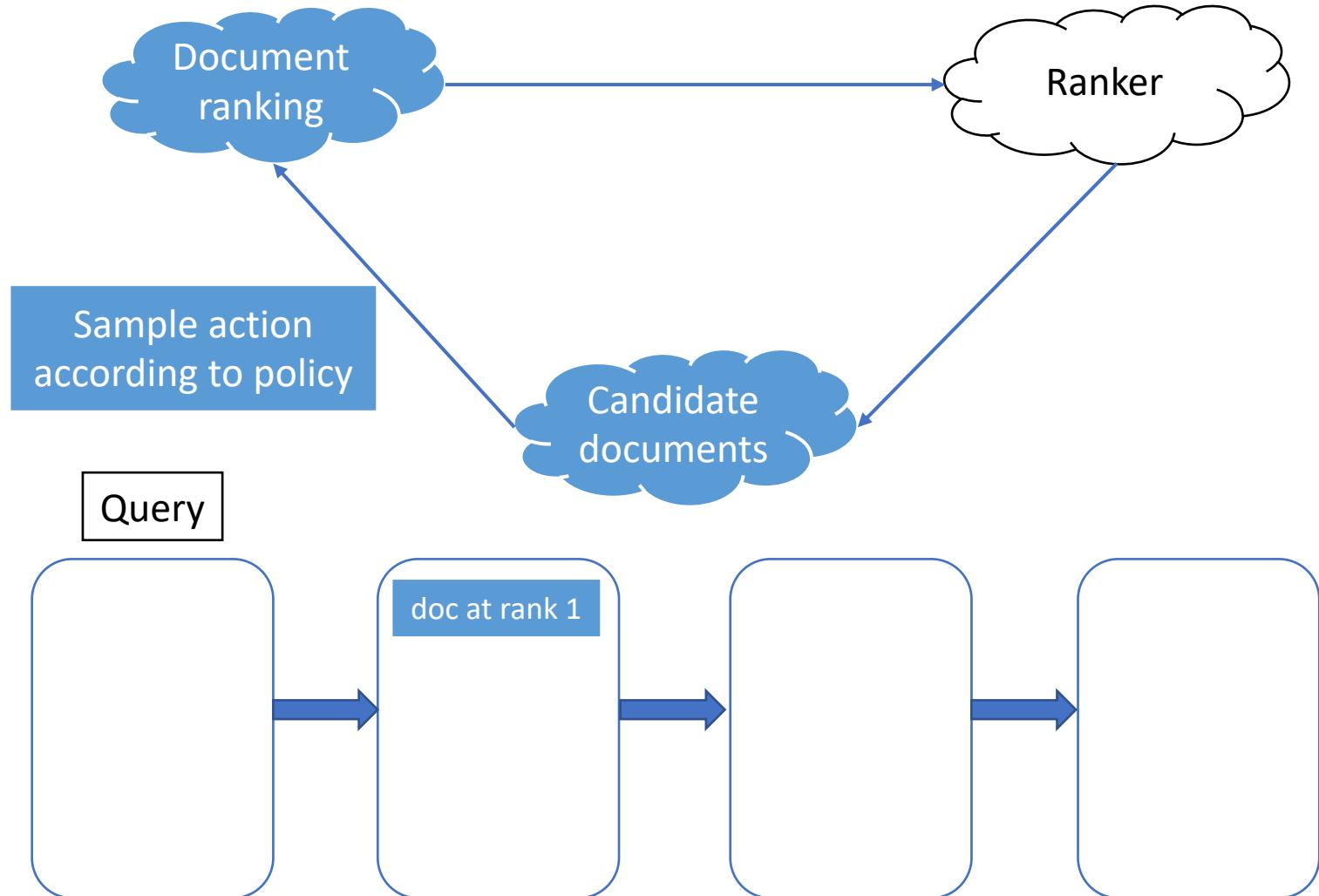
| Greedy Sequential Decision: Initialize State



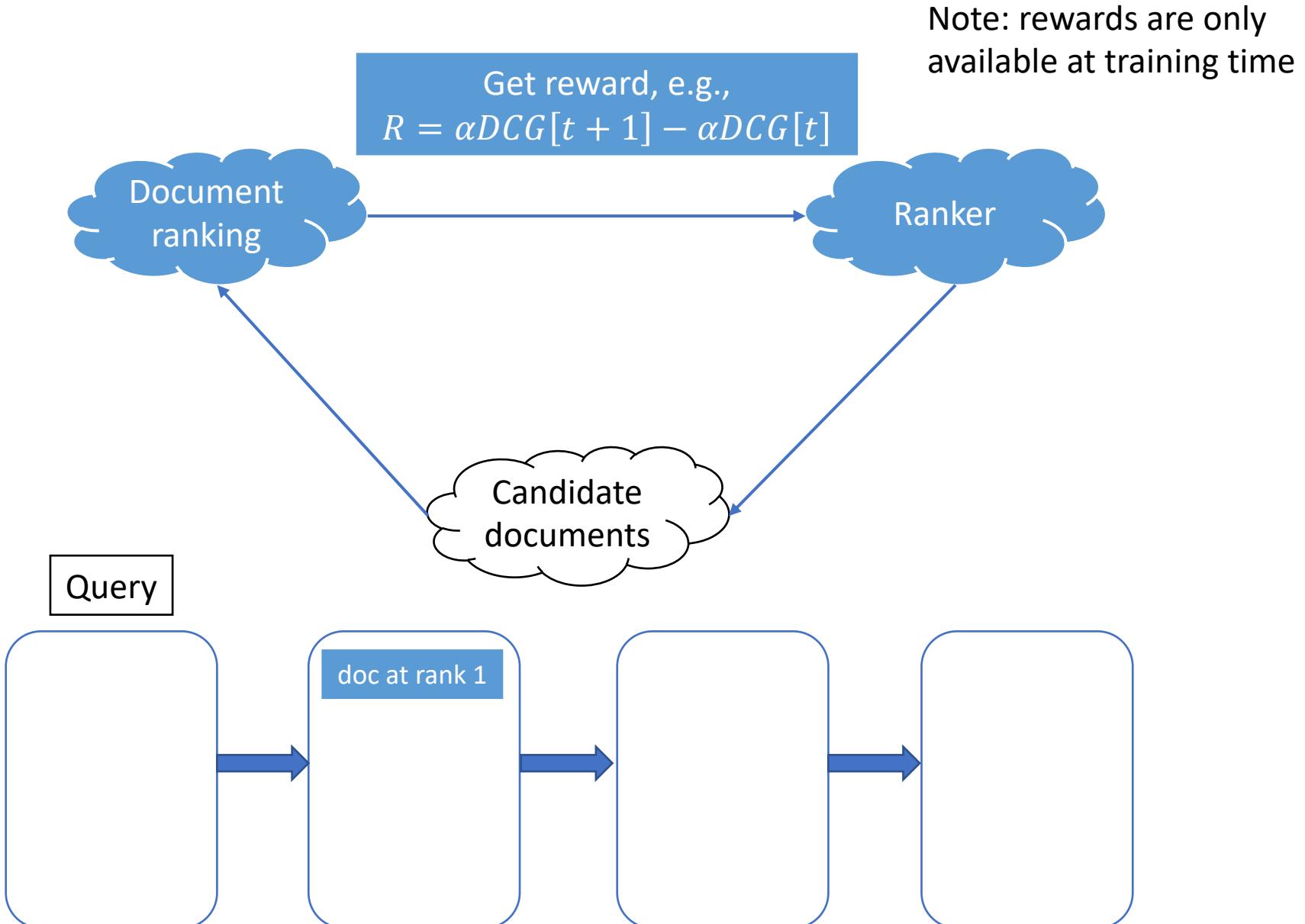
| Greedy Sequential Decision: Policy



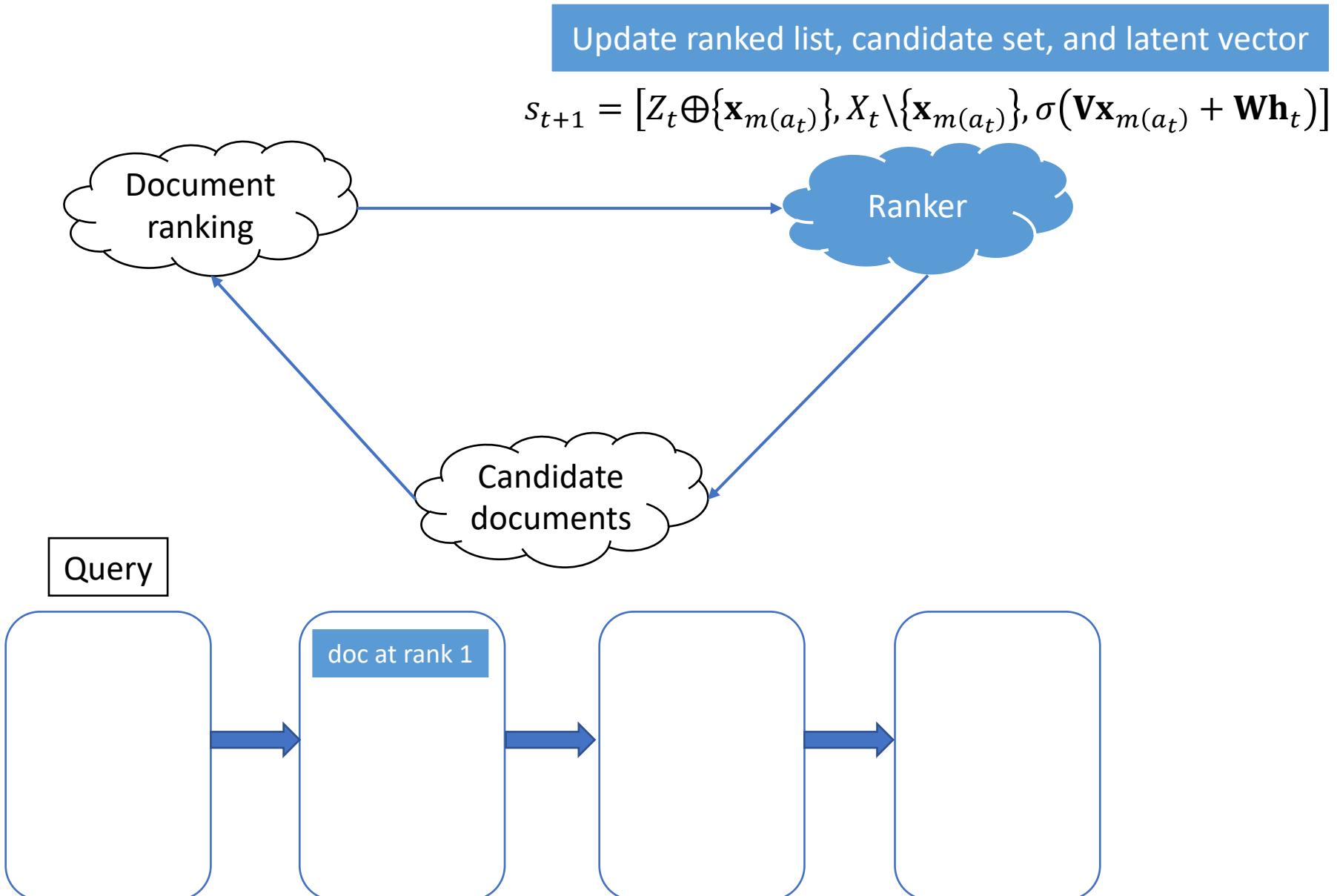
| Greedy Sequential Decision: Action



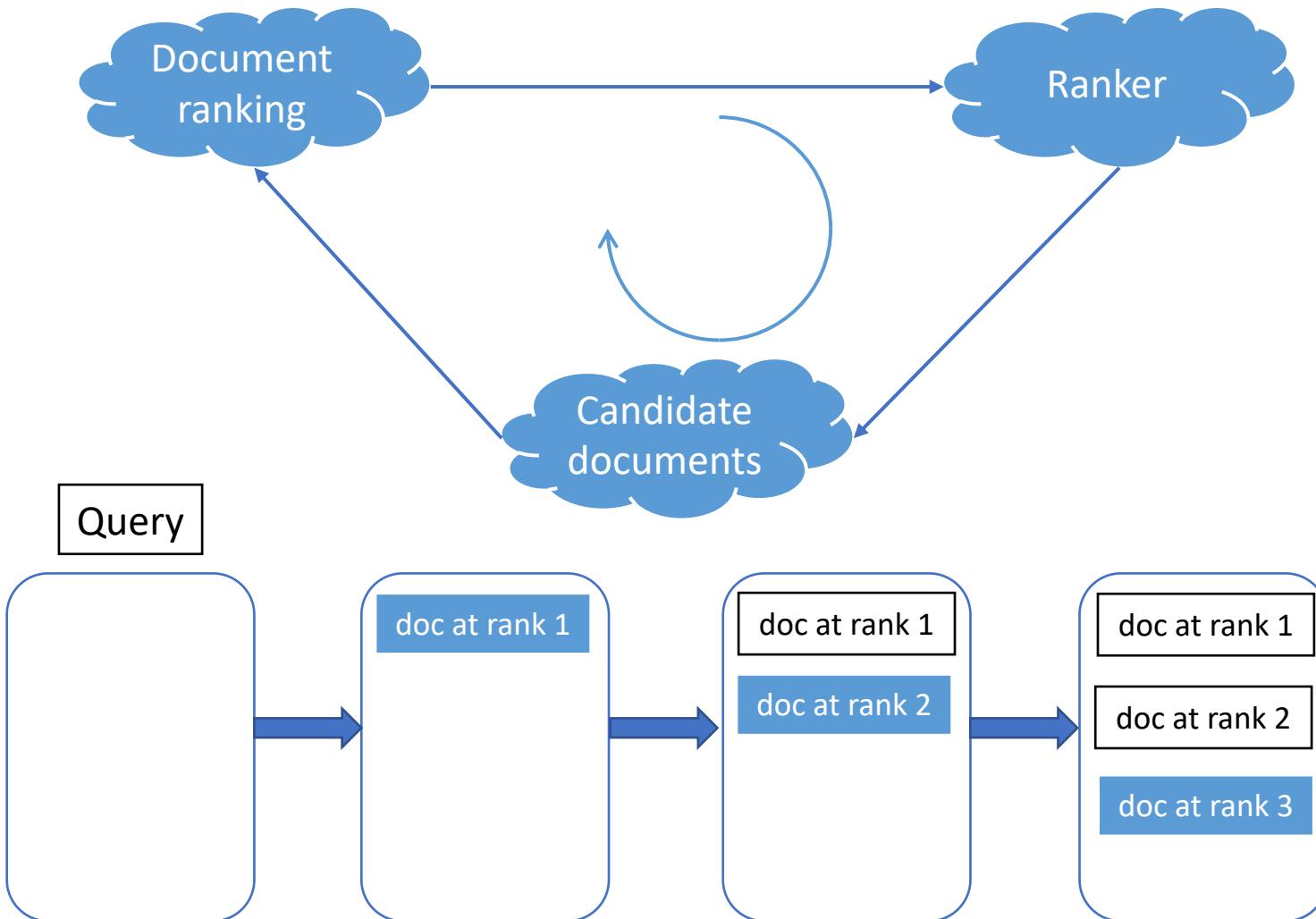
| Greedy Sequential Decision: Reward



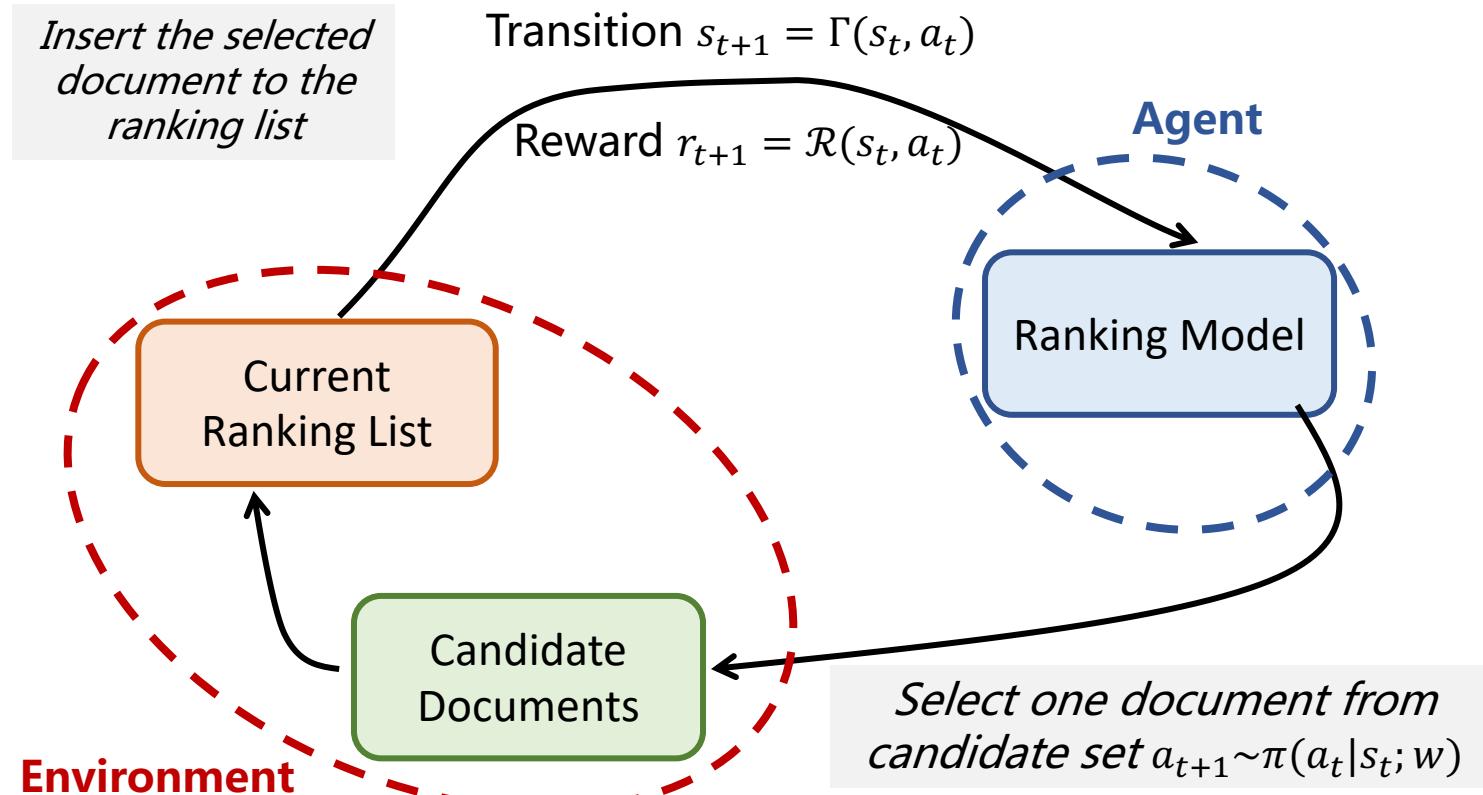
| Greedy Sequential Decision: State Transition



| Greedy Sequential Decision: Iterate



MDPRank – General Learning to Rank (Zeng et al., 2017)



MDP Elements	Definitions in Relevance Ranking
State	$S_t = [t, X_t]$
Transition	$\Gamma([t, X_t], a_t) = [t + 1, X_t / \{d_m(a_t)\}]$
Reward	$\mathcal{R}(s_t, a_t) = [\text{DCG}[t + 1] - \text{DCG}[t]]$
Policy	$\pi(a_t s_t; w) = \frac{\exp\{w^T x_{m(a_t)}\}}{\sum_{a \in A_t} \exp\{w^T x_{m(a)}\}}$

Empirical Evaluations

Method	α -NDCG@5	α -NDCG@10	ERR-IA@5	ERR-IA@10
MMR	0.2753	0.2979	0.2005	0.2309
xQuAD	0.3165	0.3941	0.2314	0.2890
PM-2	0.3047	0.3730	0.2298	0.2814
SVM-DIV	0.3030	0.3699	0.2268	0.2726
R-LTR	0.3498	0.4132	0.2521	0.3011
PAMM(α -NDCG)	0.3712	0.4327	0.2619	0.3029
NTN-DIV(α -NDCG)	0.3962	0.4577	0.2773	0.3285
MDP-DIV(α -DCG)	0.4189	0.4762	0.2988	0.3494

Method	NDCG@1	NDCG@3	NDCG@5	NDCG@10
RankSVM	0.4958	0.4207	0.4164	0.4140
ListNet	0.5326	0.4732	0.4432	0.4410
AdaRank-MAP	0.5388	0.4682	0.4613	0.4429
AdaRank-NDCG	0.5330	0.4790	0.4673	0.4496
SVMMAP	0.5229	0.4663	0.4516	0.4319
MDPRank	0.5925	0.4992	0.4909	0.4587
MDPRank(ReturnOnly)	0.5363	0.4885	0.46949	0.4591

MDP-based ranking for search result diversification and relevance ranking



THE 14TH ACM INTERNATIONAL CONFERENCE ON
WEB SEARCH AND DATA MINING

Tutorial
Morning (March, 8)
9:00-12:00, GMT+2

4. Ranking with Sequential Dependency

4.1 Heuristic Sequential Ranking Models

4.2 Learning Sequential Ranking Models

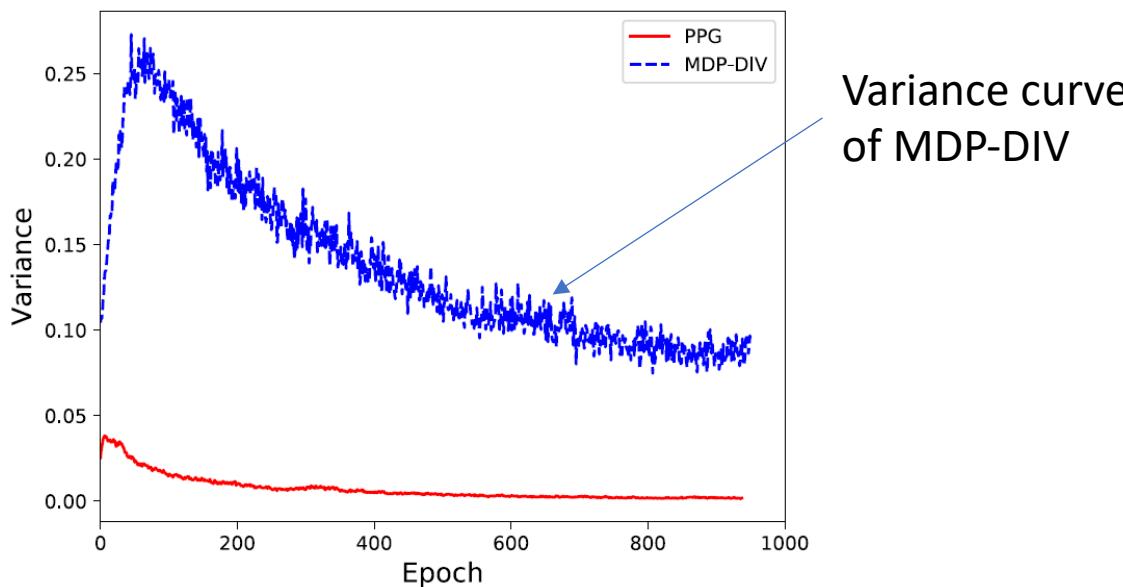
4.3 Challenges

Challenges in Greedy Sequential Decision

1

Low sampling efficiency

REINFORCE (Gradient Policy) is **unbiased**, but is known to have **high variance**. Because of huge difference in episodic rewards.



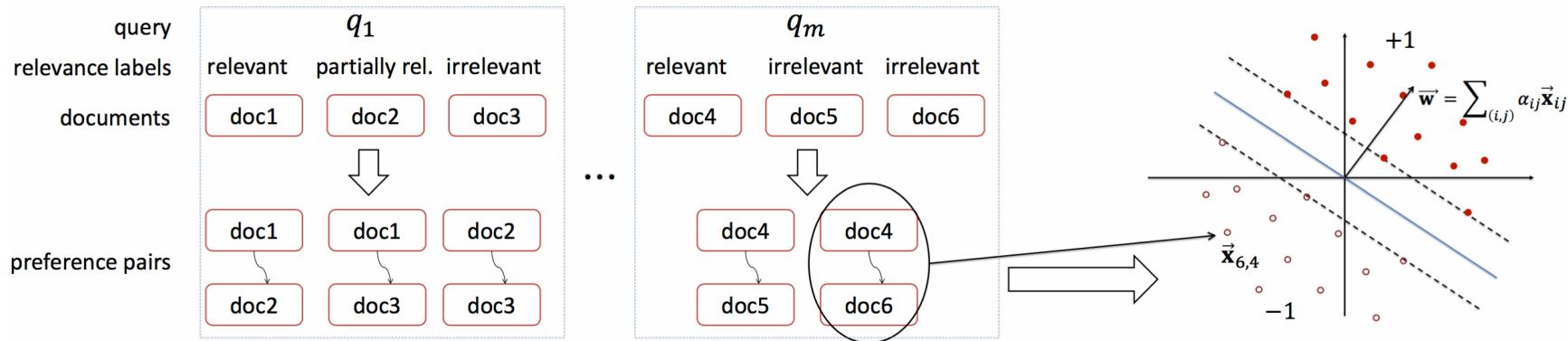
- Possible ways to improve the sampling efficiency
 - Sampling guided by experts
 - E.g., imitation learning
 - Efficient while biased estimation
 - **Control variates method**
 - Reducing variance via comparisons
 - E.g., SVRG, policy gradient with baseline
 - Unbiased estimation

I PPG – Pairwise Policy Gradient (Xu et al., '20)

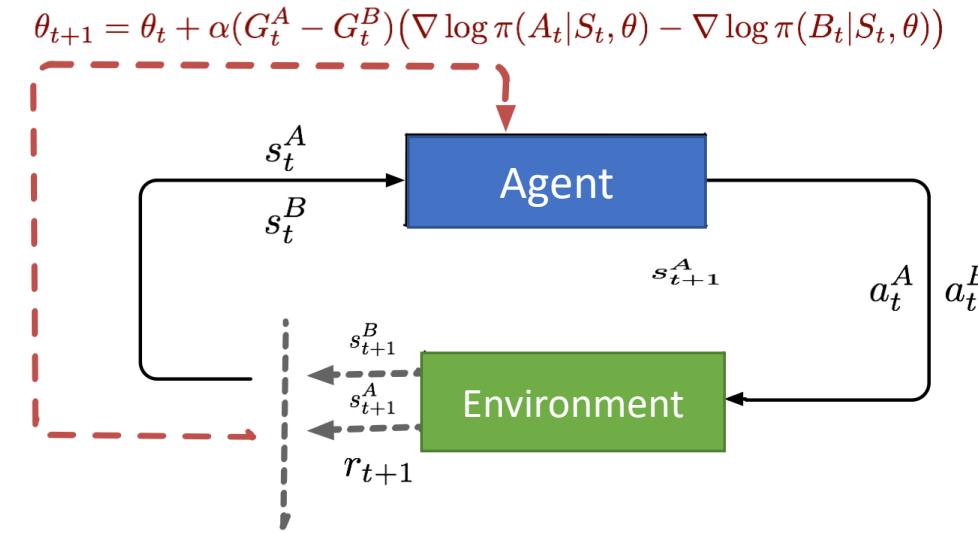
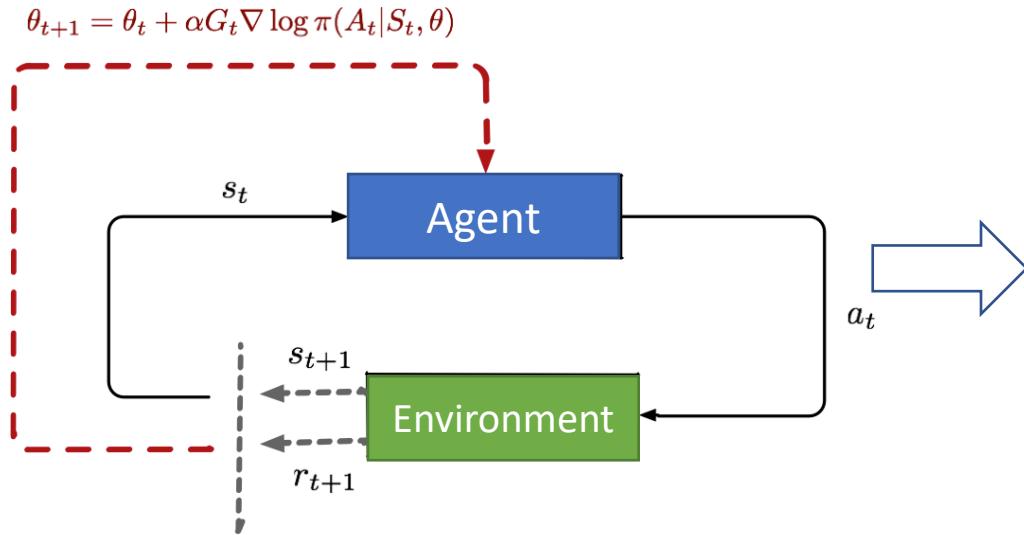
1

Improving sampling efficiency

Pairwise learning
to rank



Pairwise policy
gradient



Theoretical and Empirical Analysis of PPG

- The gradients estimated by PPG are unbiased and have low variance

THEOREM 4.1. The gradient of $J(\theta)$ in Equation (2) can be represented as

$$\nabla J(\theta) \propto \sum_s \mu(s) \sum_a \sum_b (q_\pi(s, a) - q_\pi(s, b)) \cdot (\pi(b|s; \theta) \nabla \pi(a|s; \theta) - \pi(a|s; \theta) \nabla \pi(b|s; \theta)),$$

Unbiased gradient

THEOREM 4.2. Given state $s \sim \mu_\pi$ where μ_π is an on-policy distribution under π , and given two actions $a \sim \pi(\cdot|s; \theta)$ and $b \sim \pi(\cdot|s; \theta)$. Considering the following two representations of the gradient:

$$g_1 = (q_\pi(s, a) - q_\pi(s, b)) \cdot (\nabla \log \pi(a|s; \theta) - \nabla \log \pi(b|s; \theta)),$$
$$g_2 = q_\pi(s, a) \cdot \nabla \log \pi(a|s; \theta) + q_\pi(s, b) \cdot \nabla \log \pi(b|s; \theta).$$

The variances of g_1 and g_2 satisfy

$$\text{Var}(g_1) \leq \text{Var}(g_2),$$

Low variance gradient

- PPG improves the sorting accuracy and convergence speed

Table 4: Performance comparison on LETOR OHSUMED.

Method	NDCG@1	NDCG@3	NDCG@5	NDCG@10
RankSVM	0.4958	0.4207	0.4164	0.4140
ListNet	0.5326	0.4732	0.4432	0.4410
AdaRank	0.4790	0.3730	0.4673	0.4496
MDPRank	0.5743	0.5045	0.4784	0.4558
PPG	0.5771	0.5218	0.4911	0.4664

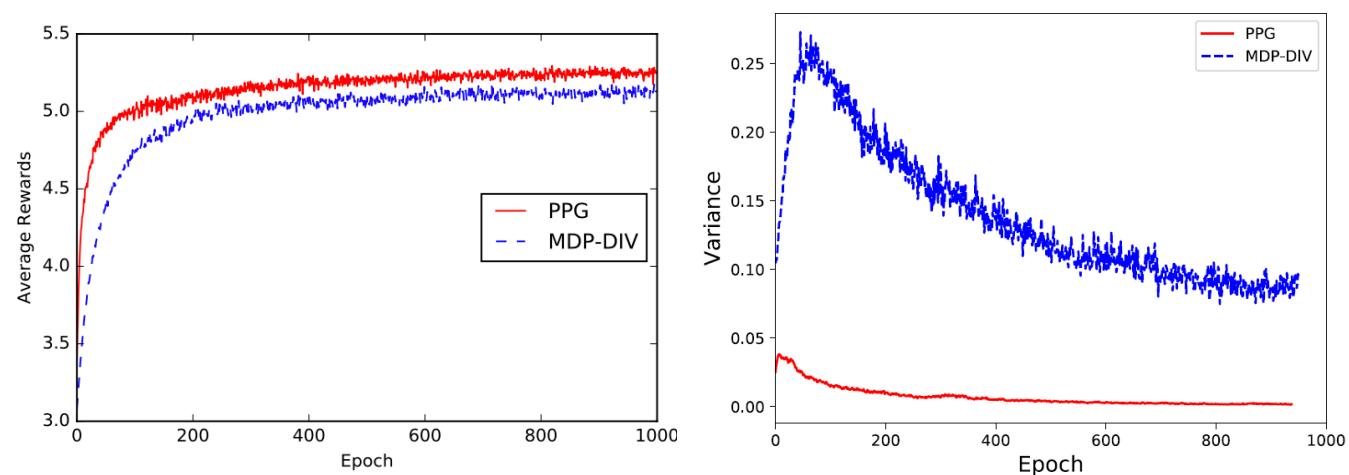


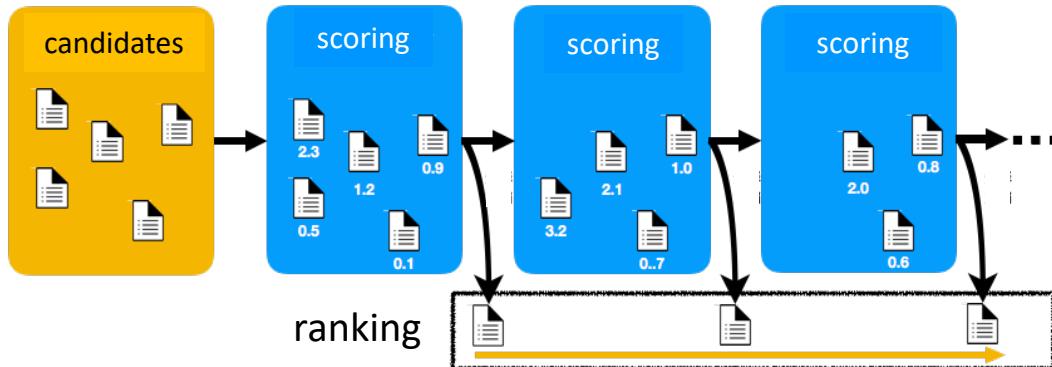
Figure 3: Learning curves of PPG and MDP-DIV.

I Challenges in Greedy Sequential Decision

2 Local Optimal

Sequential decision making is **greedy**.

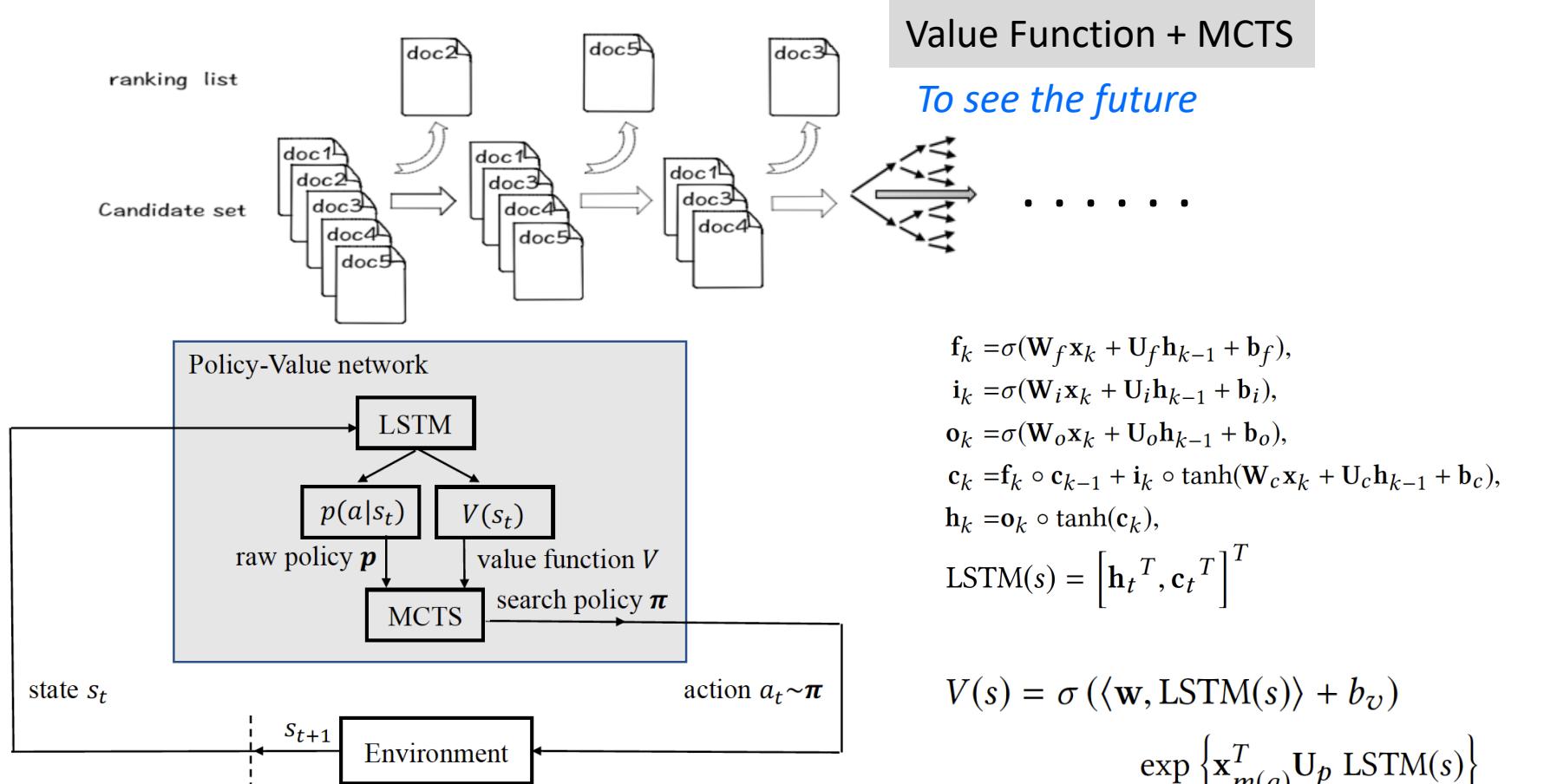
The choice of the previous steps will affect the choice of the next document, but **NOT** vice versa.



- Possible ways to avoid local optimal
 - Exhaustive search (Brute-force search)
 - Enumerating all possible candidate rankings
 - Checking their performances at each position
 - Exact optimal solution but extremely costly
- **Monte Carlo tree search (MCTS)**
 - Search tree based on random sampling
 - *Near-optimal* solution but much faster

MCTS Enhanced MDP Ranking (Feng et al., 2018)

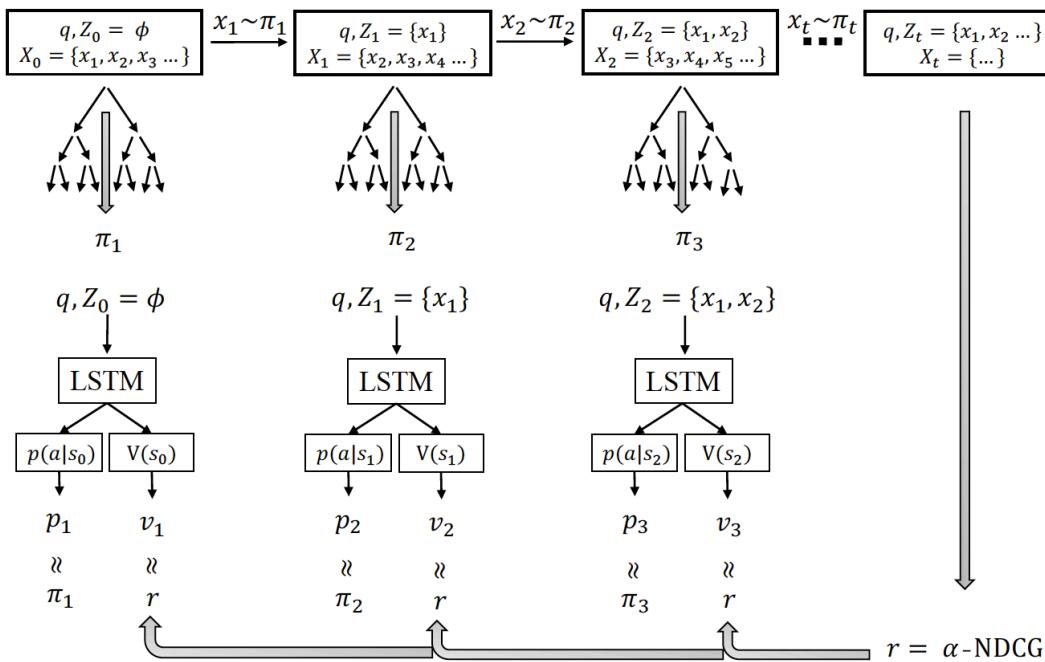
2 Avoiding Local Optimal with MCTS



- Ranking as an MDP
- MCTS guided by the predicted policies and values

Learning the Parameters

$$\ell(E, r) = \sum_{t=1}^{|E|} \left((V(s_t) - r)^2 + \sum_{a \in \mathcal{A}(s_t)} \pi_t(a|s_t) \log \frac{1}{p(a|s_t)} \right)$$



- Predicted value is as close to the real α -NDCG as possible
- Raw policy is as close to the search policy as possible

Evaluation of existing models

Method	α -NDCG@5	α -NDCG@10	ERR-IA@5	ERR-IA@10
MMR	0.2753	0.2979	0.2005	0.2309
xQuAD	0.3165	0.3941	0.2314	0.2890
PM-2	0.3047	0.3730	0.2298	0.2814
SVM-DIV	0.3030	0.3699	0.2268	0.2726
R-LTR	0.3498	0.4132	0.2521	0.3011
PAMM(α -NDCG)	0.3712	0.4327	0.2619	0.3029
NTN-DIV(α -NDCG)	0.3962	0.4577	0.2773	0.3285
MDP-DIV(α -DCG)	0.4189	0.4762	0.2988	0.3494
M ² Div(without MCTS)	0.4386*	0.4835	0.3435*	0.3668*
M ² Div(with MCTS)	0.4424*	0.4852	0.3459*	0.3686*

Heuristic Sequential Models

Supervised Learning of Sequential Ranking Models

MDP-based Sequential Ranking Models

MCTS enhanced MDP

Sequential Ranking Models for Search Result Diversification

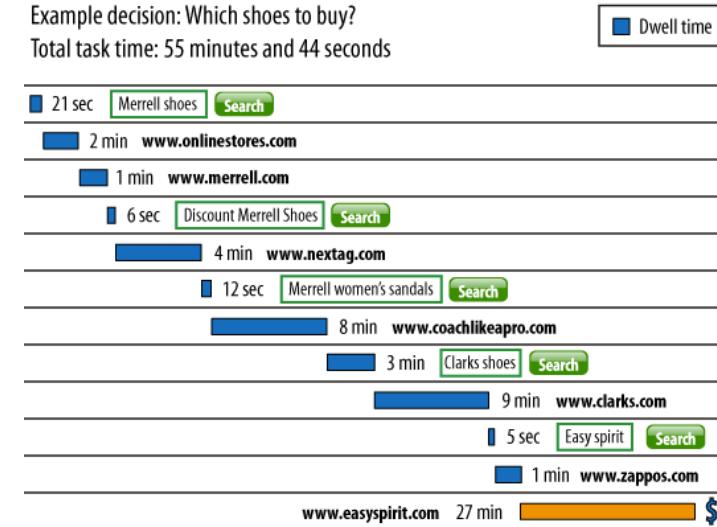
I More Sequential Dependent Ranking Models

- For multi-page search (depends on results in previous result page)
 - MDP-MPS (Zeng et al., ICTIR '18)
 - E-commerce Search as MDP: DPG-FBE (Hu et al., KDD '18)
- For session Search (depends on results in previous search session)
 - Query Change Model (QCM) (Yan et al, '15)
 - Win-Win (Luo et al, '14)
 - DPL (Luo et al, '15)
- Interactively optimizing IR systems
 - Dueling Bandit Gradient Descent (DBGD) (Yue and Joachims, ICML '09)
 - Balancing Exploration and Exploitation (Hofmann et al., IRJ '13)



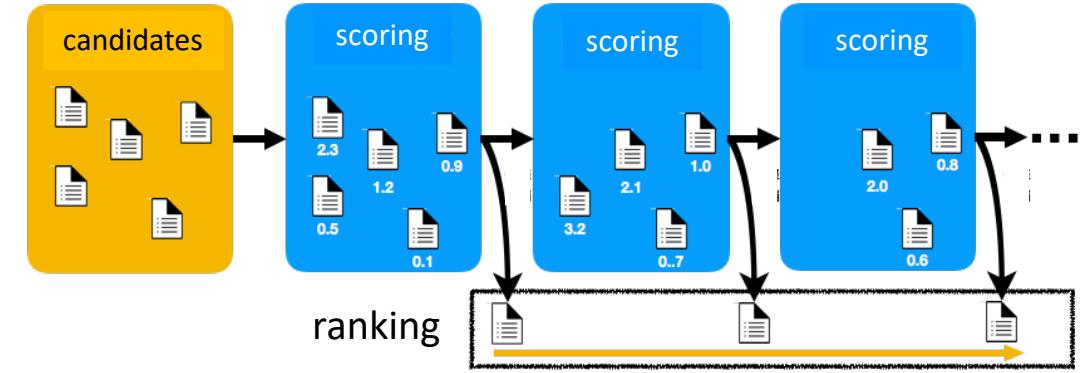
Inside a real query "session"

Example decision: Which shoes to buy?
Total task time: 55 minutes and 44 seconds



Brief Summary: Sequential Dependent Ranking Models

- Ranking as **greedy sequential decision making**
 - Key: the scoring function f
 - Heuristically predefined f
 - Learned f
 - Linear / kernel / deep function
 - Multi-armed bandits
 - Markov decision process
 - Challenges
 - Low efficiency
 - Local optimal
- Next: globally dependent ranking models



```
1  $\mathcal{D}_q \leftarrow \emptyset$ 
2 while  $|\mathcal{D}_q| < \tau$  do
3    $d^* \leftarrow \arg \max_{d \in \mathcal{R}_q \setminus \mathcal{D}_q} f(q, d, \mathcal{D}_q)$ 
4    $\mathcal{R}_q \leftarrow \mathcal{R}_q \setminus \{d^*\}$ 
5    $\mathcal{D}_q \leftarrow \mathcal{D}_q \cup \{d^*\}$ 
6 end while
7 return  $\mathcal{D}_q$ 
```

Reference

- R. L. T. Santos, C. Macdonald, and I. Ounis. Search Result Diversification. FnTIR Vol. 9, No. 1 (2015) 1–90.
- J. Carbonell and J. Goldstein. The use of MMR, diversity-based reranking for reordering documents and producing summaries. SIGIR 1998, 335–336.
- C. Zhai, W. W. Cohen, and J. Lafferty. Beyond independent relevance: Methods and evaluation metrics for subtopic retrieval. SIGIR 2003, 10–17.
- V. Dang and W. B. Croft. Diversity by proportionality: an election-based approach to search result diversification. SIGIR 2012, 65–74.
- R. L. T. Santos, C. Macdonald, and I. Ounis. Exploiting query reformulations for web search result diversification. WWW 2010, 881–890.
- Y. Yue and T. Joachims. Predicting Diverse Subsets Using Structural SVMs, ICML 2008.
- R. Agrawal, S. Gollapudi, A. Halverson, and S. Jeong. Diversifying search results. WSDM 2009, 5–14.
- Y. Zhu , Y. Lan, J.Guo, X. Cheng and S. Niu. Learning for Search Result Diversification. SIGIR 2014, 293-302.
- L. Xia, J. Xu, Y. Lan, J. Guo, and X. Cheng. Learning Maximal Marginal Relevance Model via Directly Optimizing Diversity Evaluation Measures. SIGIR 2015, 113-122.
- L. Xia, J. Xu, Y. Lan, J. Guo, and X. Cheng. Modeling Document Novelty with Neural Tensor Network for Search Result Diversi_cation. SIGIR 2016, 395-404.
- Z. Jiang, J. R. Wen, Z. Dou, W. X. Zhao, J.-Y. Nie, M. Yue. Learning to Diversify Search Results via Subtopic Attention. SIGIR 2017, 545-554
- F. Radlinski, R. Kleinberg, and T. Joachims. Learning diverse rankings with multi-armed bandits. ICML 2008, 784–791.
- L. Xia, J. Xu, Y. Lan, J. Guo, W. Zeng, and X. Cheng. Adapting Markov Decision Process for Search Result Diversification. SIGIR 2017, 535-544.

Reference

- W. Zeng, J. Xu, Y. Lan, J. Guo, and X. Cheng. Reinforcement Learning to Rank with Markov Decision Process. SIGIR 2017, 945-948.
- J. Xu, W. Zeng, L. Xia, Y. Lan, D. Yin, X. Cheng, and J.-R. Wen. Reinforcement Learning to Rank with Pairwise Policy Gradient. SIGIR 2020, 509–518
- Y. Feng, J. Xu, Y. Lan, J. Guo, W. Zeng, X. Cheng. From Greedy Selection to Exploratory Decision-Making: Diverse Ranking with Policy-Value Networks. SIGIR 2018, 125-134.
- M. Streeter, D. Golovin, and A. Krause. Online learning of assignments. In NIPS 2009.
- Y. Hu, Q. Da, A. Zeng , Y. Yu, Y. Xu. Reinforcement Learning to Rank in E-Commerce Search Engine: Formalization, Analysis, and Application. KDD 2018, 368–377.
- W. Zeng, J. Xu, Y. Lan, J. Guo, X. Cheng. Multi Page Search with Reinforcement Learning to Rank. ICTIR 2018, 175-178.
- H. Yang, D. Guan, S. Zhang. The Query Change Model: Modeling Session Search as a Markov Decision Process. TOIS, No. 20, 2015.
- J. Luo, S. Zhang, H. Yang. Win-Win Search: Dual-Agent Stochastic Game in Session Search. SIGIR 2014,
- J. Luo, S. Zhang, X. Dong, H. Yang. Designing states, actions, and rewards for using POMDP in session search. ECIR 2015, 526-537.
- Y. Yue, T. Joachims. Interactively optimizing information retrieval systems as a dueling bandits problem. ICML 2009.
- J. Liu, Z. Dou, X. Wang, S. Lu, and J.-R. Wen. DVGAN: A Minimax Game for Search Result Diversification Combining Explicit and Implicit Features. SIGIR 2020. 479-488.
- M. Streeter, D. Golovin, A. Krause. Online learning of assignments. NIPS 2009. 1794-1802.
- L. Chen, G. Zhang, H. Zhou. Fast Greedy MAP Inference for Determinantal Point Process to Improve Recommendation Diversity. NIPS 2018. 1-12.
- M. Wilhelm, A. Ramanathan, A. Bonomo, S. Jain, E. H. Chi, J. Gillenwater. Practical Diversified Recommendation on Youtube with Determinantal Point Process. CIKM 2018.

Thank you!