



Contextual Combinatorial Cascading Bandits

Authors: Shuai Li , Baoxiang Wang , Shengyu Zhang , Wei Chen

(ICML'16: Proceedings of the 33rd International Conference on International Conference on Machine Learning - Volume 48, June 2016)

Instructor: Dr. Chicheng Zhang

Presenter: Rupal Jain

Multi-Armed Bandits Problem

- System of base arms, K
- Reward of each arm i are random samples from unknown distributions with unknown means.
 - best arm $\mu^* = \max \mu_i$ (mean = μ_i)
- In each round t , the learning agent selects one arm i_t to play and observe the reward, $R_t(i_t)$
- Regret after playing T rounds:
$$\text{Regret} = T\mu^* - \mathbb{E}\left[\sum_{t=1}^T R_t(i_t)\right]$$

Goal: Minimize the Cumulative regret

Problem of MAB: dealing with trade-off between exploitation and exploration





Feedbacks

At every time step, a learning agent chooses a subset of ground items (super arm) under certain combinatorial constraints.

1. **Bandit feedback:** only reward of chosen super arm is obtained
2. **Semi-bandit feedback:** observe outcomes of the individual base arms in super arm
3. **Cascading feedback:** can obtain the reward of the super arm and the weights of some base arms in the chosen subset of arm, according to some stopping criterion.

Related Work

Contextual Combinatorial Multi-Armed Bandit

- Involves Semi-bandit feedback and nonlinear reward
- May have other combinatorial constraints

Ref: Qin et al., 2014

Cascading Bandit

- Feasible action is Sequential list
- Stopping at first satisfactory item from the list
- Feedback is recorded as Clicks
- Disjunctive objective: reward of an action is 1 if there is at least one “good” item in the list

Ref: Kveton et al. 2015a



Machine Learning

[All](#) [Images](#) [Books](#) [News](#) [Videos](#) [More](#)

About 2,230,000,000 results (0.37 seconds)

Wikipedia
https://en.wikipedia.org/wiki/Machine_learning

Machine learning

Machine learning (ML) is a field of study in artificial intelligence concerned with the development and study of statistical algorithms that can effectively ...

[Outline of machine learning](#) · [Timeline](#) · [Machine Learning \(journal\)](#) · [Data mining](#)

IBM
<https://www.ibm.com/topics/machine-learning>

What is Machine Learning?

Machine learning is a branch of artificial intelligence (AI) and computer science on the use of data and algorithms to imitate the way that ...

[What is machine learning?](#) · [Machine Learning vs. Deep...](#)

MIT Sloan
<https://mitsloan.mit.edu/ideas-made-to-matter/machine-learning>

Machine learning, explained

Apr 21, 2021 — **Machine learning** is a subfield of artificial intelligence that has the ability to learn without explicitly being programmed. "In ...

[Why It Matters](#) · [How Businesses Are Using...](#) · [How Machine Learning Works](#)

SAS Institute
https://www.sas.com/SAS_Insights/Analytics

Machine Learning: What it is and why it matters

Machine learning is a method of data analysis that automates analytical model building. It is a branch of artificial intelligence based on the idea that ...

[Machine Learning In Today's...](#) · [Learn More About Industries...](#) · [How It Works](#)

Related Work - Cont.

Position Discount

- Combes et al. 2015a considered a **cascading model** with a particular case of **contextual information**,
- User is recommended with **K** items.
- The position discounts introduced makes the recommended set in the decreasing order of UCBs, which is more realistic
- Considered in the list order, so that the agent's reward is discounted depending on the position where the stopping criterion is met.

Recommended For You



Wei zhuang zhe
(TV series 2015-)
Drama
★8.4
Kai Wang, Dong Jin, Ge Hu



The Terminator (1984)
R 1 h 47 min - Action | Sci-Fi
★8.1 83 Metascore
Arnold Schwarzenegger, Linda Hamilton, Michael Biehn



Back to the Future (1985)
PG 1 h 56 min - Adventure | Comed...
★8.5 86 Metascore
Michael J. Fox, Christopher Lloyd, Lea Thompson



Breaking Bad
(TV series 2008-2013)
49 min - Crime | Drama | Thriller
★9.5
Bryan Cranston, Aaron Paul, Anna...



Argo (2012)
R - 2 h - Biography | Drama | History 5

Related Work - Cont.

Combinatorial Cascading Multi-Armed Bandit

- Feasible action is subset of items under combinatorial constraints
- Can involve Semi-bandit feedback
- Conjunctive objective: reward of an action is 1 if all the items are “good” in the list
- Challenges:
 - Exponential number of actions
 - Offline optimization may already be hard

Ref: Kveton et al., 2015c

Recommended For You



Wei zhuang zhe
(TV series 2015-)
Drama

★8.4

Kai Wang, Dong Jin, Ge Hu



The Terminator (1984)

R 1 h 47 min - Action | Sci-Fi

★8.1 **83** Metascore

Arnold Schwarzenegger, Linda Hamilton, Michael Biehn



Back to the Future (1985)

PG 1 h 56 min - Adventure | Comed...

★8.5 **86** Metascore

Michael J. Fox, Christopher Lloyd, Lea Thompson



Breaking Bad
(TV series 2008-2013)

49 min - Crime | Drama | Thriller

★9.5

Bryan Cranston, Aaron Paul, Anna...



Argo (2012)

RR - 2 h - Biography | Drama | History

Author's Contribution

- Formulate the **Contextual Combinatorial Cascading Bandits** problem which handles
 - contextual information
 - cascading feedback
 - position discount
 - general reward function (non-linear)
- Proposed C^3 - Algorithm
- Satisfy Monotonicity and Lipschitz continuity conditions

	context	cascading	Position Discount	General Reward
Combinatorial UCB ¹	No	Yes	No	Yes
Contextual Combinatorial UCB ²	Yes	No	No	Yes
Comb-Cascade ³	No	Yes	No	No
C^3 -UCB	Yes	Yes	Yes	Yes

Ref:

1: Chen et al. 2016

2: Qin et al., 2014

3: Kveton et al., 2015c

Learning Protocol

- $E=\{1,...,L\}$: set of base arms
- **Action** $A=(a_1,...,a_k): a_1,...,a_k \in E$; a sequence of base arms (each with length of k)
 - There is a feasible action set \mathcal{S} with length at most K
- At each time $t \geq 1$
 - set of contexts $\{x_{t,a}\}_{a \in E}$ are given (e.g. user/keyword features) and is ≤ 1
 - learning agent recommends a feasible action to user, $A_t = (a_1^t, \dots, a_{|A_t|}^t) \in \mathcal{S}$
 - Cascading Feedback Model: The user checks from the first item and stops at O_t -th item.
 - Feedback: observe weights (“quality”) of first O_t items in A_t at time t , $w_t(a_k^t)$, $k \leq O_t$.

$$\mathbb{E}[w_t(a)|H_t] = \theta_*^T x_{t,a} = w_{t,a}$$

Fixed but unknown

Learning Protocol - Cont.

- Assume the expected reward of an action A is a function $f(A, \mathbf{w})$ of expected weight, $\mathbf{w}_t = \{\mathbf{w}_{t,a}\}_{a \in E} = (\theta_*^T \mathbf{x}_{t,a})_{a \in E}$ of each base arm.
- Satisfy the following assumptions:
 - **Monotonicity:** expected reward function $f(A, \mathbf{w})$ is a non-decreasing with respect to \mathbf{w} : for any $\mathbf{w}, \mathbf{w}' \in [0, 1]^E$, if $w(a) \leq w'(a)$, we have $f(A, \mathbf{w}) \leq f(A, \mathbf{w}')$.
 - **Lipschitz continuity:** The expected reward function $f(A, \mathbf{w})$ is B -Lipschitz continuous with respect to \mathbf{w} together with position discount parameters $\gamma_k \in [0, 1]$, $k \leq K$. For any $\mathbf{w}, \mathbf{w}' \in [0, 1]^E$, we have

$$|f(A, \mathbf{w}) - f(A, \mathbf{w}')| \leq B \sum_{k=1}^{|A|} \gamma_k |w(a_k) - w'(a_k)|$$

where $A = (a_1, \dots, a_{|A|})$.



Learning Protocol - Cont.

- An Oracle $\mathcal{O}_{\mathcal{S}}(\mathbf{w})$ is called an α -approximation oracle for some $\alpha \leq 1$, if on given input \mathbf{w} , the oracle returns an action $\mathbf{A} = \mathcal{O}_{\mathcal{S}}(\mathbf{w}) \in \mathcal{S}$ satisfying $f(\mathbf{A}, \mathbf{w}) \geq \alpha f(\mathbf{A}^*, \mathbf{w})$ where $\mathbf{A} = \mathbf{argmax}_{\mathbf{A} \in \mathcal{S}} f(\mathbf{A}, \mathbf{w})$
- Regret in T rounds

$$\text{Regret} = \sum_{t=1}^T f_t^* - \mathbb{E} \left[\sum_{t=1}^T f(\mathbf{A}_t, \mathbf{w}_t) \right]$$

Best cumulative reward

f_t^* : max expected reward in round t

Example - Movie Recommendation

- Each Movie i has a feature vector \mathbf{m}_i .
- At time t ,
 - A random user comes with feature vector \mathbf{u}_t
 - Use $\mathbf{x}_{i,t} = \mathbf{g}(\mathbf{m}_i, \mathbf{u}_t)$, a function of \mathbf{m}_i and \mathbf{u}_t , as context
 - The learning agent recommends a list of movies \mathbf{A}_t .
 - The user checks from the first movie and stops at the attractive one.
 - The learning agent received reward γ_k if the user stops at position k .

$$1 = \gamma_k \geq \dots \geq \gamma_1 \geq 0$$

Recommended for you

23



All Hallows' Eve 2013

★★ =

22



Scream 1996

★★★★ =

21



An American Werewolf in Paris
1997

★★ =

20



Arthur: Malediction 2022

★½ =

20



Godzilla vs. Kong 2021

★★★★½ =

19



White Men Can't Jump 2023

★½ =

18



**Ant-Man and the Wasp:
Quantumania** 2023

★★★ =

17



F9 2021

★½ =

17



Super 2010

★★★★½ =

15



Scary Movie 4 2006

★½ =

Algorithm

- To estimate the expected reward, we get an estimate θ_* using ridge regression problem on context vector \mathbf{x} and observed rewards \mathbf{w} .
- We get an ℓ^2 -regularized least-squares estimation of θ_* with regularization parameter, $\lambda > 0$

$$\hat{\theta}_t = (\mathbf{X}_t^T \mathbf{X}_t + \lambda \mathbf{I})^{-1} \mathbf{X}_t^T \mathbf{Y}_t$$

Where $\mathbf{X}_t \in \mathbb{R}^{(\sum_{s=1}^t O_s) \times d}$: matrix with rows as $\gamma_k \mathbf{x}_{s,a_k^s}^T$ and \mathbf{Y}_t : column vector with elements as $\gamma_k \mathbf{w}_s(a_k^s), k \in [O_s], s \in [t]$

- Let \mathbf{V}_t is a symmetric positive definite matrix, $\in \mathbb{R}^{d \times d}$

$$\mathbf{V}_t = \mathbf{X}_t^T \mathbf{X}_t + \lambda \mathbf{I} = \lambda \mathbf{I} + \sum_{s=1}^t \sum_{k=1}^{O_s} \gamma_k^2 \mathbf{x}_{s,a_k^s} \mathbf{x}_{s,a_k^s}^T$$

Algorithm

- **Lemma 1:** $\beta_t(\delta) = R\sqrt{\ln\left(\frac{\det(V_t)}{\lambda^d \delta^2}\right)} + \sqrt{\lambda}$. For any $\delta > 0$, with probability at least $1 - \delta$, for all $t > 0$, we have $\|\hat{\theta}_t - \theta_*\|_{V_t} \leq \beta_t(\delta)$

Proof: By theorem 2 in (Abbasi-Yadkori et al., 2011), we can obtain a good estimate of difference between θ_* and $\hat{\theta}_t$.

Interpretation: With high probability, the estimate $\hat{\theta}$ lies in the ellipsoid centered at θ_* with confidence radius $\beta_t(\delta)$ under V_t norm.

- The upper confidence bound (UCB) of the expected weight for each base arms:

$$U_t(a) = \min \left\{ \hat{\theta}_{t-1}^T x_{t,a} + \beta_{t-1}(\delta) \|x_{t,a}\|_{V_{t-1}^{-1}}, 1 \right\}$$

- **Lemma 2:** When $\|\hat{\theta}_t - \theta_*\|_{V_t} \leq B_t(\delta)$ holds for time $t-1$, we have

$$0 \leq U_t(a) - w_{t,a} \leq 2\beta_{t-1}(\delta) \|x_{t,a}\|_{V_{t-1}^{-1}}$$

Proof: Note $w_{t,a} = (\theta_*^T x_{t,a})$

By Holder's inequality

$$\begin{aligned} |\hat{\theta}_{t-1}^T x_{t,a} - \theta_*^T x_{t,a}| &= |[V_{t-1}^{1/2}(\hat{\theta}_{t-1} - \theta_*)]^T (V_{t-1}^{-1/2} x_{t,a})| \\ &\leq \|V_{t-1}^{1/2}(\hat{\theta}_{t-1} - \theta_*)\|_2 \|V_{t-1}^{-1/2} x_{t,a}\|_2 \\ &= \|\hat{\theta}_{t-1} - \theta_*\|_{V_{t-1}} \|x_{t,a}\|_{V_{t-1}^{-1}} \\ &\leq \beta_{t-1}(\delta) \|x_{t,a}\|_{V_{t-1}^{-1}} \end{aligned}$$

Because $1 - \theta_*^T x_{t,a} \geq 0$ and

$$\begin{aligned} 0 &\leq (\hat{\theta}_{t-1}^T x_{t,a} + \beta_{t-1}(\delta) \|x_{t,a}\|_{V_{t-1}^{-1}}) - \theta_*^T x_{t,a} \\ &\leq 2\beta_{t-1}(\delta) \|x_{t,a}\|_{V_{t-1}^{-1}} \end{aligned}$$

Claimed result is obtained...



Algorithm

1. Learning agent computer (UCBs), $\mathbf{U}_t \in [0,1]^E$ for expected weights of all base arms in E .
2. Uses \mathbf{U}_t to select an action $\mathbf{A}=(\mathbf{a}_1^t, \dots, \mathbf{a}_{|\mathbf{A}_t|}^t)$
3. User selects from the first base arm in \mathbf{A}_t and stops at O_t -th base arm,
Agent observes $\mathbf{w}_t(\mathbf{a}_k^t)$, $k \leq O_t$
4. Learning agent updates $\mathbf{V}_t, \mathbf{X}_t, \mathbf{Y}_t$ to get newer estimates $\hat{\boldsymbol{\theta}}_t$ and $\boldsymbol{\theta}_*$ and new confidence radius, $\beta_t(\delta)$

Algorithm 1 C³-UCB

- 1: Parameters:
- 2: $\{\gamma_k \in [0, 1]\}_{k \leq K}; \delta = \frac{1}{\sqrt{n}}; \lambda \geq C_\gamma = \sum_{k=1}^K \gamma_k^2$
- 3: Initialization:
- 4: $\hat{\boldsymbol{\theta}}_0 = 0, \beta_0(\delta) = 1, \mathbf{V}_0 = \lambda \mathbf{I}, \mathbf{X}_0 = \emptyset, \mathbf{Y}_0 = \emptyset$
- 5: **for all** $t = 1, 2, \dots, n$ **do**
- 6: Obtain context $x_{t,a}$ for all $a \in E$
- 7: $\forall a \in E$, compute
- 8: $\mathbf{U}_t(a) = \min\{\hat{\boldsymbol{\theta}}_{t-1}^\top x_{t,a} + \beta_{t-1}(\delta) \|x_{t,a}\|_{\mathbf{V}_{t-1}^{-1}}, 1\}$
- 9: //Choose action \mathbf{A}_t using UCBs \mathbf{U}_t
- 10: $\mathbf{A}_t = (\mathbf{a}_1^t, \dots, \mathbf{a}_{|\mathbf{A}_t|}^t) \leftarrow \mathcal{O}_S(\mathbf{U}_t)$
- 11: Play \mathbf{A}_t and observe $\mathbf{O}_t, \mathbf{w}_t(\mathbf{a}_k^t), k \in [\mathbf{O}_t]$
- 12: //Update statistics
- 13: $\mathbf{V}_t \leftarrow \mathbf{V}_{t-1} + \sum_{k=1}^{\mathbf{O}_t} \gamma_k^2 x_{t,\mathbf{a}_k^t} x_{t,\mathbf{a}_k^t}^\top$
- 14: $\mathbf{X}_t \leftarrow [\mathbf{X}_{t-1}; \gamma_1 x_{t,\mathbf{a}_1^t}^\top; \dots; \gamma_{\mathbf{O}_t} x_{t,\mathbf{a}_{\mathbf{O}_t}^t}^\top]$
- 15: $\mathbf{Y}_t \leftarrow [\mathbf{Y}_{t-1}; \gamma_1 \mathbf{w}_t(\mathbf{a}_1^t); \dots; \gamma_{\mathbf{O}_t} \mathbf{w}_t(\mathbf{a}_{\mathbf{O}_t}^t)]$
- 16: $\hat{\boldsymbol{\theta}}_t \leftarrow (\mathbf{X}_t^\top \mathbf{X}_t + \lambda \mathbf{I})^{-1} \mathbf{X}_t^\top \mathbf{Y}_t$
- 17: $\beta_t(\delta) \leftarrow R \sqrt{\ln(\det(\mathbf{V}_t)/(\lambda^d \delta^2))} + \sqrt{\lambda}$
- 18: **end for** t

Result

Theorem: Suppose the expected reward function $f(\mathbf{A}, \mathbf{w})$ is a function of expected weights and satisfies the requirements of monotonicity and B-Lipschitz continuity. Then the α -regret of **C³-UCB** algorithm, satisfies

$$\text{Regret}^\alpha(n) = O\left(\frac{dB}{p^*} \sqrt{TK \ln(T)}\right)$$

Proof:

$$\begin{aligned} R^\alpha(n) &\leq \frac{2\sqrt{2}B}{p^*} \sqrt{nKd \ln(1 + C_\gamma n / (\lambda d))} \\ &\quad \cdot (R \sqrt{\ln[(1 + C_\gamma n / (\lambda d))^d n]} + \sqrt{\lambda}) + \alpha \sqrt{n} \\ &= O\left(\frac{dBR}{p^*} \sqrt{nK \ln(C_\gamma n)}\right) \end{aligned}$$

$p_{t,A}$ be the probability of full observation of A .
 $p^* = \min_{1 \leq t \leq T} \min_{A \in S} p_{t,A}$: Minimal probability that action has all base arms observed all time.

d : dimension of latent and feature vectors;

K : largest length of the sequence

$n = T$ (number of rounds)

R is the sub-Gaussian constant

$C_\gamma \leq K$ (sum of discounts \leq number of positions)

$$C_\gamma = \sum_{k=1}^K \gamma_k^2 \leq K$$

Experiment 1

Synthetic Data: compare C^3 -UCB to CombCascade (Combinatorial Cascading) on synthetic problems.

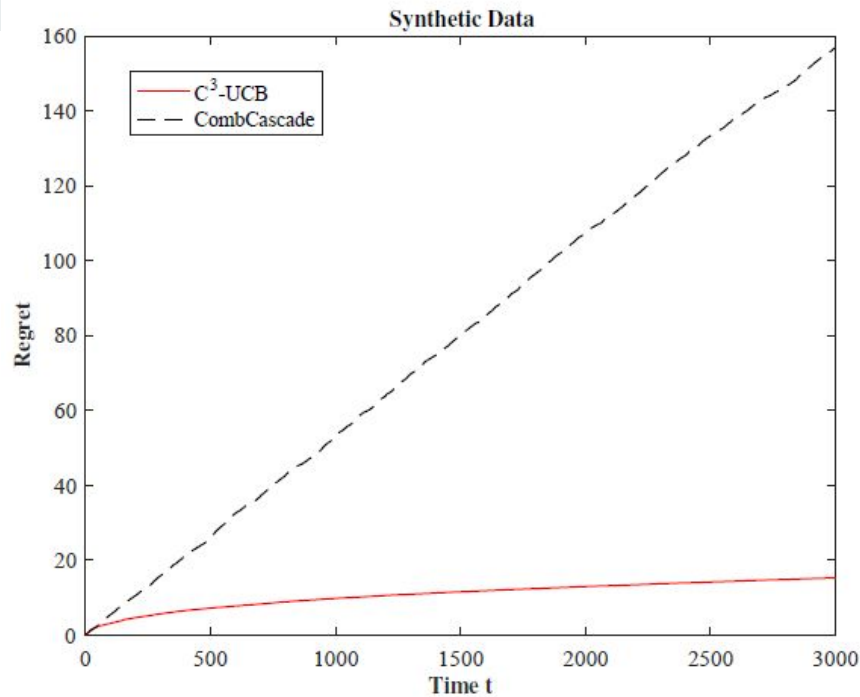
Setup:

1. The problem is a contextual cascading bandit with $L = 200$ items and $K = 4$, $d = 20$ where at each time t .
2. The agent recommends K items to the user.
3. At first, we randomly choose a θ in \mathbb{R}^{d-1} with $\|\theta\|_2 = 1$ and let $\theta_* = (\theta/2, 1/2)$.
4. Then at each time t , we randomly assign $\mathbf{x}'_{t,a}$ in \mathbb{R}^{d-1} with $\|\mathbf{x}'_{t,a}\|_2 = 1$ to arm \mathbf{a} and use $\mathbf{x}_{t,a} = (\mathbf{x}'_{t,a}, 1)$ to be the contextual information for arm \mathbf{a} .

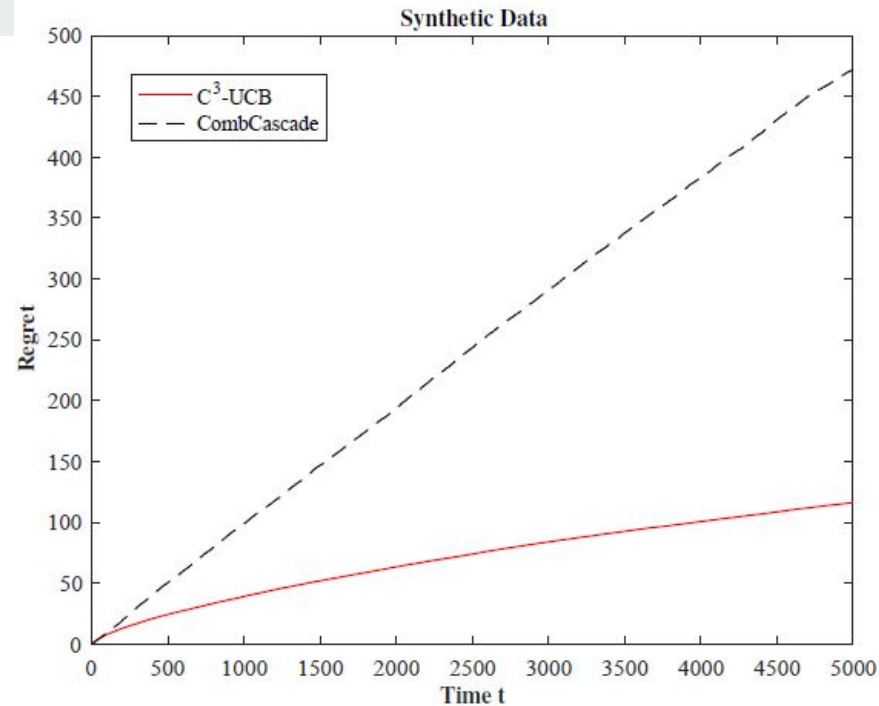
$$\theta_*^\top x_{t,a} = \frac{1}{2}(\theta^\top x'_{t,a} + 1) \in [0, 1]$$

This processing will guarantee the inner product

Next we generate the weight for arm \mathbf{a} at time t by a random sample from the Bernoulli distribution with mean $\theta_*^\top \mathbf{x}_{t,a}$

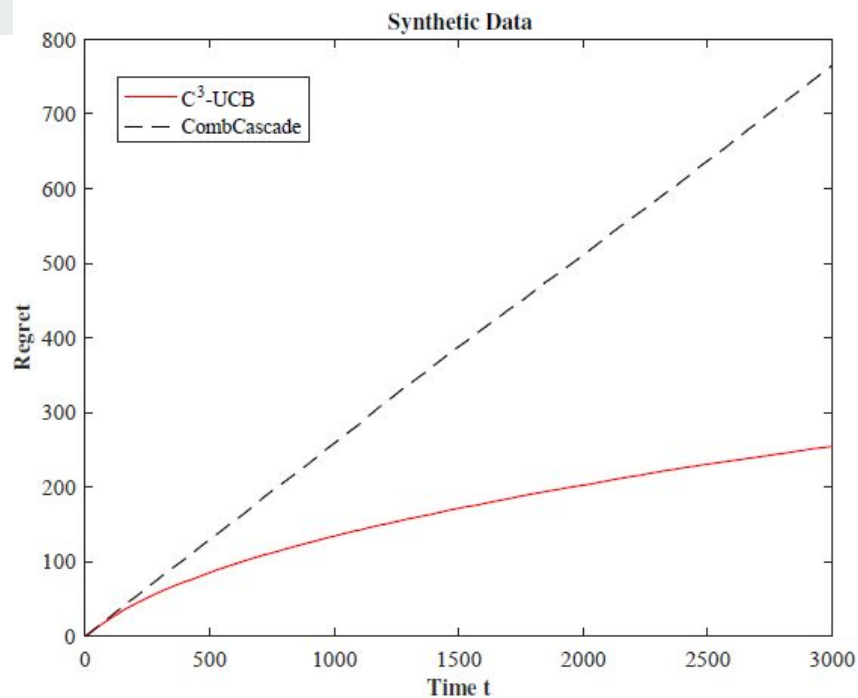


Disjunctive, $\gamma = 1$, 9.77%

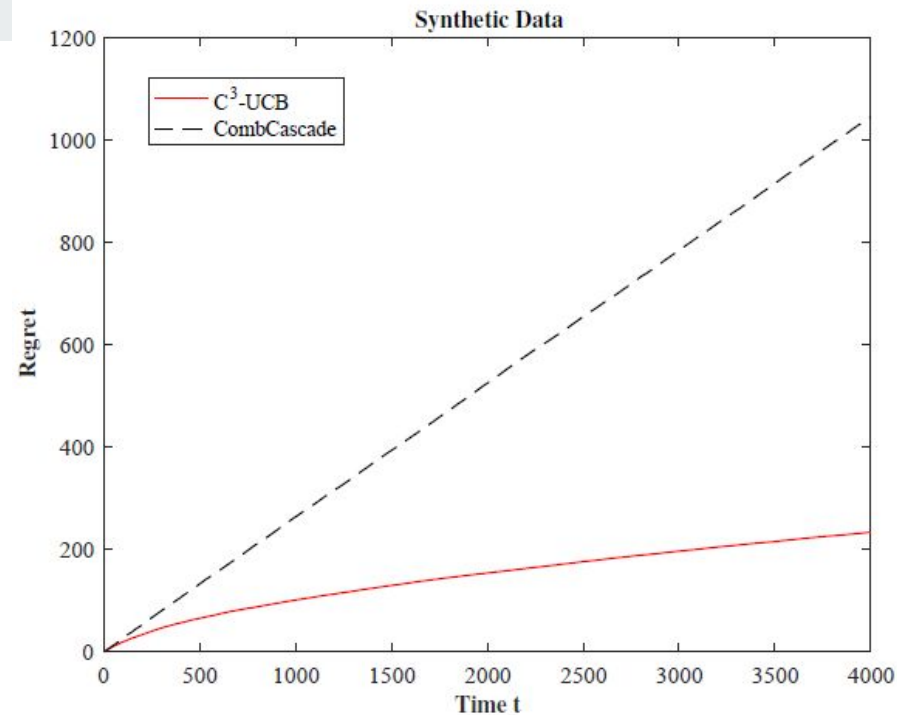


Disjunctive, $\gamma = 0.9^{k-1}$, 24.6%

The above two settings are of disjunctive objective where the learning agent chooses a set of K items out of L ground items and observes a prefix of the chosen K items until the first one with weight 1;



Conjunctive, $\gamma = 1$, 3.33%



Conjunctive, $\gamma = 0.9^{k-1}$, 22.4%

The above two settings are of conjunctive objective where the learning agent chooses a set of K items out of L ground items and observes from the first item until the first one with weight 0.



Experiment 2

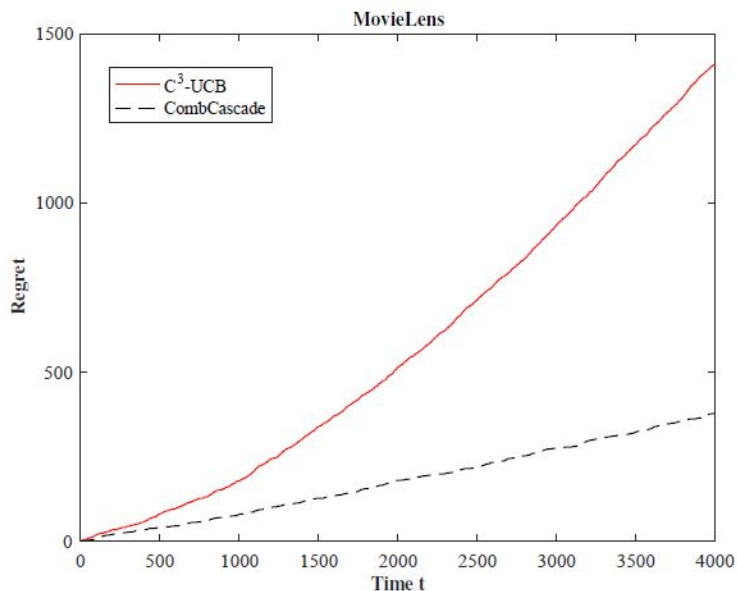
Movie Recommendation: evaluate C3-UCB algorithm with dataset MovieLens (Lam & Herlocker, 2015) of 2015

Problem:

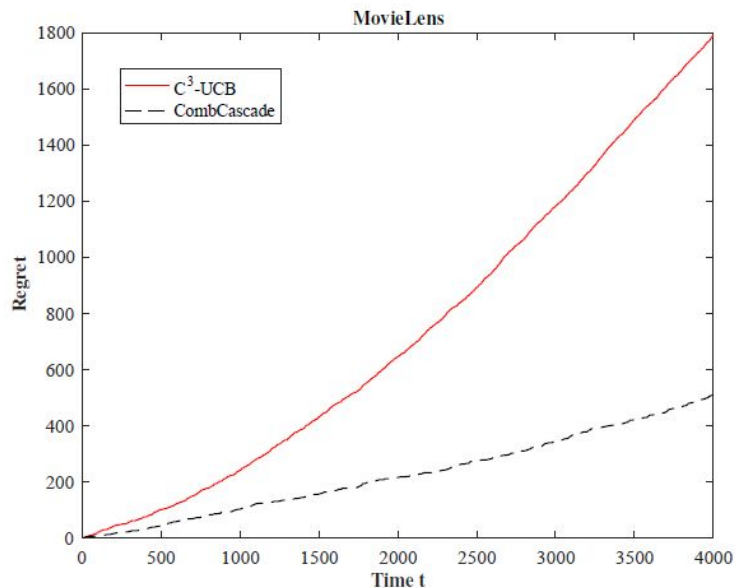
1. There is a big sparse matrix $A \in \{0, 1\}^{N_1 \times N_2}$ where $A(i; j) = 1$ denotes user i has watched movie j .
2. Split A as $H + F$ according to a Bernoulli distribution $\sim \text{Ber}(p)$ for some fixed p . (H : known history “what users have watched” and F as future criterion.)
3. Derive **feature vectors** of both users and movies via SVD, $U = (u_1, \dots, u_{N_1})$ and $M = (m_1, \dots, m_{N_2})$.
4. At every time t , we randomly choose a user $I_t \in [N_1]$. From Li et al., 2010, use $x_{t;j} = u_{I_t} m_j^T$ as the contextual information for each movie j .
5. User provides cascading feedback, stopping at a movie they like and Agent receives discounted reward based on stop position.
6. The real weight of movie j at time t , $w_t(j)$, is $F(I_t, m_j)$.

Setup:

1. The problem is a contextual cascading bandit with $L = 400$ movies and $K = 4$, $d = 400$ where at each time t .
2. We experiment with both $\gamma = 1$ (no position discount) and $\gamma = 0.9$, and compare our algorithm with CombCascade.



$\gamma=1$



$\gamma=0.9$

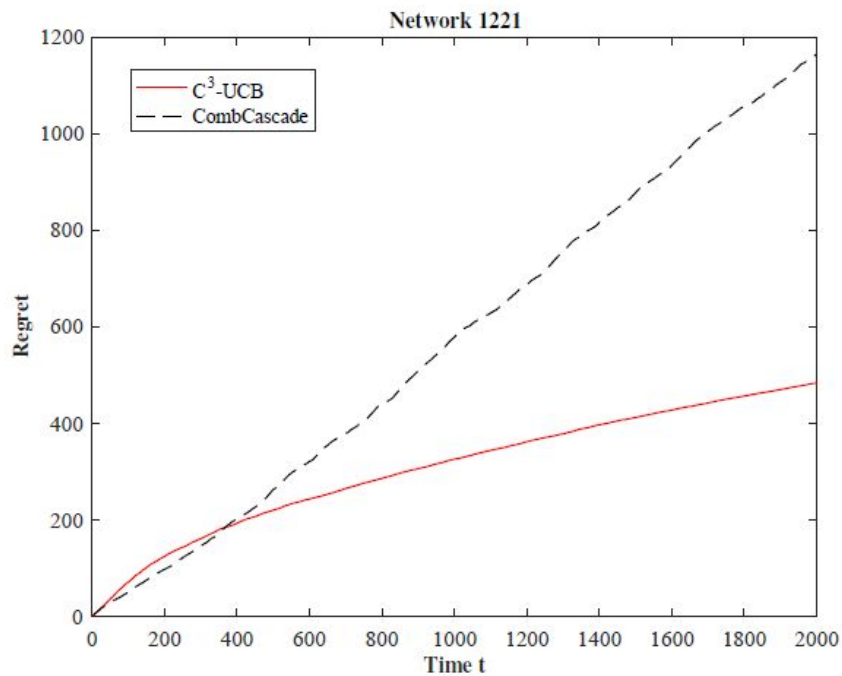
Result: The rewards of our algorithms are **3.52** and **3.736** times of those of CombCascade (for $\gamma = 1$ and 0.9 , respectively), which demonstrate the advantage to involve contextual information in real applications.

Experiment 3

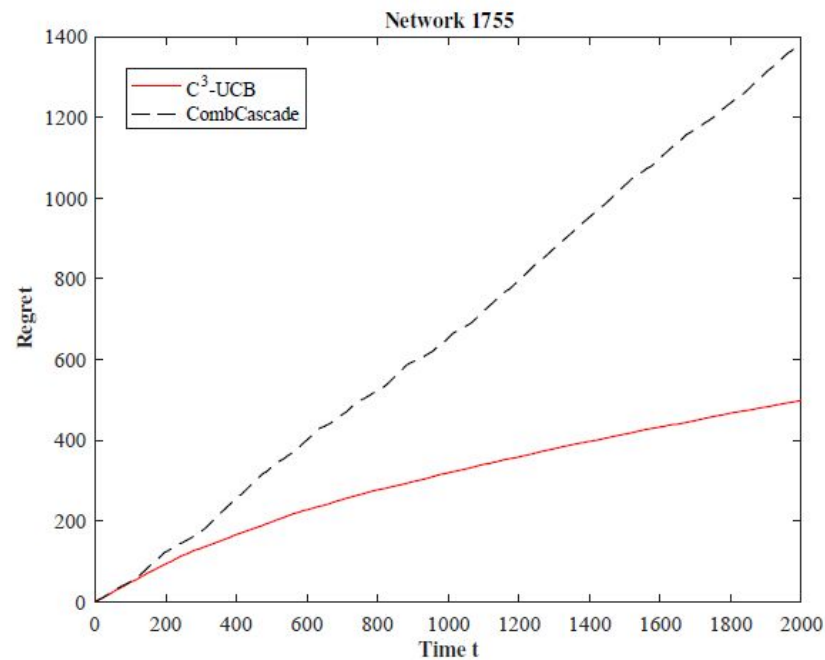
Network Routing Problem: evaluate C^3 -UCB with RocketFuel dataset (Spring et al., 2004).

Setup:

1. The ground set E is the set of links in the network.
2. Before learning, the environment randomly chooses a d -dimensional vector $\theta \in [0, 1]^d$ ($d = 5$)
3. At each time t , a pair of source and destination nodes are randomly chosen
4. The feasible action set S_t at time t contains all simple paths, paths without cycles, between the source and destination.
5. Any edge a in the set S_t is assigned with a random d -dimensional contextual information vector $x_{t,a}$.



Network 1221, $\gamma=1$



Network 1221, $\gamma=0.9$

Result: Both θ and x have been processed same as experiment 1, such that $\theta_*^T x \in [0, 1]$. The weight for each edge a is a sample from Bernoulli distribution with mean $\theta_*^T x_{t,a}$. Then the learning agent recommends a feasible path A to maximize the expected reward in the conjunctive objective. This experiment is on different position discounts.

Conclusion:

- Incorporating contextual information to cascading bandit with position discounts
- Each action is an ordered list and only a prefix of the action is observed each time.
- demonstrate the advantage to involve contextual information and position discounts.
- Application potential
 - Any sequential list recommendation (search, ads, mobile recommendations)
 - Need online (real time) feedback

Future Work:

- investigate on lower bounds of the regret and cascading on general graphs

References:

- <https://icml.cc/2016/reviews/581.txt>
- Kveton, Branislav, Wen, Zheng, Ashkan, Azin, and Szepesvari, Csaba. Combinatorial cascading bandits. Advances in Neural Information Processing Systems, 2015c.
- Kveton, Branislav, Szepesv'ari, Csaba, Wen, Zheng, and Ashkan, Azin. Cascading bandits: learning to rank in the cascade model. In Proceedings of the 32th International Conference on Machine Learning, 2015a.
- Qin, Lijing, Chen, Shouyuan, and Zhu, Xiaoyan. Contextual combinatorial bandit and its application on diversified online recommendation. In Proceedings of the 2014 SIAM International Conference on Data Mining (SDM), 2014.
- Combes, Richard, Magureanu, Stefan, Proutiere, Alexandre, and Laroche, Cyrille. Learning to rank: Regret lower bounds and efficient algorithms. In Proceedings of the 2015 ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Systems, pp. 231–244. ACM, 2015a.