

# 強化学習勉強会

Lihong Li et al., 2010, A Contextual-Bandit Approach  
to Personalized News Article Recommendation

田中一樹

# この論文を選んだ理由

- ・ ちょっと前にスマホプッシュ配信にBanditを使った記事が出て少し盛り上がった記憶
- ・ 知っておきたい（？）
- ・ Webサイトや広告メールにのせる記事や商品を実際に出しながらユーザごとに適した対応がしたい場合多そう
- ・ CTR(Click Through Rate)やConversion率を高めたい
- ・ Bandit使えるってなってるけど実際どういう風に特徴量とかアルゴリズム適用してるの??

# 今日の目標

- ・ Contextual-Banditアルゴリズムの雰囲気を知る
- ・ 具体的なデータへの適用方法と有用性を知る
- ・ どんなフィーチャやデータセット、状況に応用可能かふんわり考える

# 目次

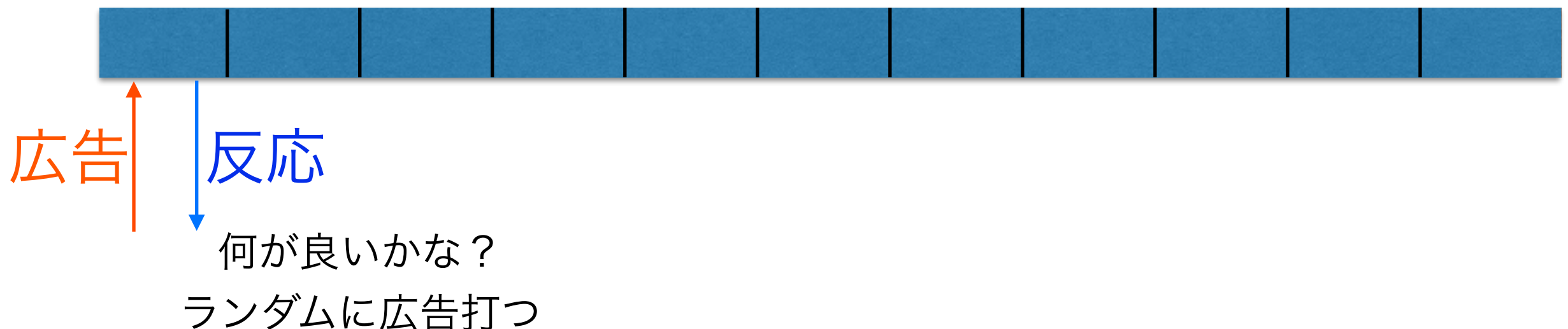
1. WEBにおけるコンテンツ配信の難しさ (Motivation)
2. Bandit Algorithm
  - ・ Contextual-Bandit Algorithmの概要
  - ・ Contextual-Free Bandit Algorithmの概要
3. 提案手法
  - ・ LinUCB with Disjoint Linear Models
  - ・ LinUCB with Hybrid Linear Models
4. 評価方法
5. 実験
6. 結論

## WEBにおけるコンテンツ配信の難しさ (Motivation)

- ・ **コンテンツが大量**にあり、常に追加、削除を繰り返しながら**変化**
  - ・ 古いニュース記事の影響を徐々に消していきながら、速報の人気度を迅速に特定するとか
- ・ 一般的にコンテンツだけの情報で人気度をモデル化するのは困難
- ・ 実際には、ユーザのフィードバックを集めながら新コンテンツの人気度を評価する

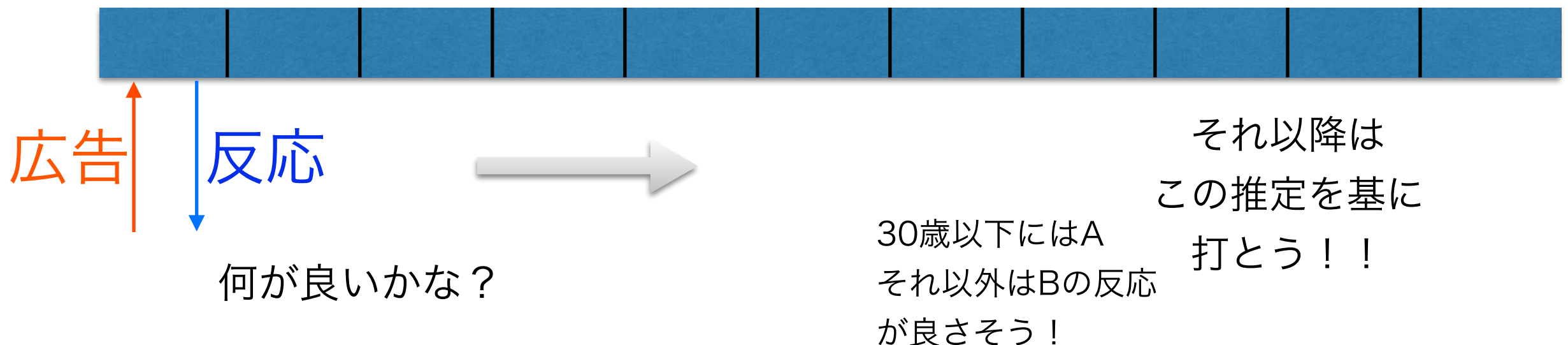
# 例

- ・ 小さなトラフィックで試して、ユーザの反応を基に残りのトラフィックで人気のものを探していく
  - ・  $\epsilon$  - greedy, EXP3, UCB1, ...
  - ・ 新しいコンテンツに対してより多くのトラフィックを掛ける必要がある（当然）



# 例

- ・ 小さなトラフィックで試して、ユーザの反応を基に残りのトラフィックで人気のものを探していく
  - ・  $\epsilon$  - greedy, EXP3, UCB1, ...
  - ・ 新しいコンテンツに対してより多くのトラフィックを掛ける必要はある（当然）



# 例

- ・ 小さなトラフィックで試して、ユーザの反応を基に残りのトラフィックで人気のものを探していく
  - ・  $\epsilon$  - greedy, EXP3, UCB1,...
  - ・ 新しいコンテンツに対してより多くのトラフィックを掛ける必要がある（当然）

全ユーザ





# 協調フィルタリングもあるが

- ・ 多くのWEBサービスでは、コンテンツの多様性は頻繁に変化し、コンテンツの人気は時間とともに変化する
  - ・ cold-startやし無理ゲー。。
- ・ でも新ユーザ、新コンテンツ、もしくは両方成り立つ場合に対して良いマッチングを探すことは必須
  - ・ 現実的に考えると次の課題が生まれる

# 本研究のモチベーション

- ・ 興味とコンテンツ間のマッチング情報を集めながら
- ・ ユーザ満足度を長い目で見て最大化したい

“最終的に”って意味合い



Contextual Bandit Problemや！

新しいアルゴリズムLinUCB提案したよ！

# でも

- ・ インタラクティブだから評価するのが難しいお
- ・ オフラインで前もって性能評価できるような方法も考えたお！（理論的にもOKだお！）
- ・ 色んなアルゴリズムで評価したお！

# Bandit Algorithm

# Contextual-Bandit Algorithm

- ・ Context情報を用いたMulti-armed banditとしてパーソナライズされた記事推薦の問題をモデル化
- ・ Multi-armed banditは割愛
- ・ 既存研究[1]に従ってContextual banditと呼ぶ
- ・ 便宜上Contextual-Bandit Algorithmを”A”と示す

[1]J.Langford et al., 2008, The epoch-greedy algorithm for contextual multi-armed bandits

# アルゴリズム

- ・ 離散時間の試行を $t=1,2,3,\dots$ とすると、試行 $t$ で
  1. “A”は現在のユーザ $u_t$ 、アームまたは行動の集合 $A_t$ を観測。同時にすべての行動 $a$ に対し、その特徴ベクトル $x_{t,a}$ も観測。 $x_{t,a}$ はユーザとアーム両方の要約情報であり、“context”と呼ばれる
  2. 以前の試行のpayoffに基づき、“A”は行動 $a_t$ を選択しpayoff  $r_{t,a_t}$ を受け取る。 $r_{t,a_t}$ の期待値はユーザ $u_t$ とアーム $a_t$ の両方に依存する。

“A”: Contextual-bandit algorithm

# アルゴリズム

3. そして新しい観測  $(x_{t,a}, a_t, r_{t,a_t})$  を使ってアーム選択の戦略を向上させる。ここでは選択されていないアームのフィードバックは観測されないことに注意。

このプロセスを経た、“A”のtotal T-trial payoffは

$$\sum_{t=1}^T r_{t,a_t}$$

“A”: Contextual-bandit algorithm

# アルゴリズム

(再掲) “A”のtotal T-trial payoff

$$\sum_{t=1}^T r_{t,a_t}$$

同様に、optimal expected T-trial payoffも定義

$$\mathbf{E} \left[ \sum_{t=1}^T r_{t,a_t^*} \right]$$

$a_t^*$ は試行tでの最大期待payoffをもつアーム

(目標) 上の期待値を最大にするように“A”を作る



# regretを最小化する

(目標) **payoffの期待値を最大にするよう“A”を作る**

- “A”のregretを $R_A(T)$ と定義すると

$$R_A(T) \stackrel{\text{def}}{=} \underset{\text{optimal}}{\mathbf{E}} \left[ \sum_{t=1}^T r_{t,a_t^*} \right] - \mathbf{E} \left[ \sum_{t=1}^T r_{t,a_t} \right] \quad (1)$$

実際のpayoff

- (目標)はregretを最小化することと等価になる！

# 重要な特殊ケースを考えておく

- ・ 任意の $t$ で、アーム集合 $A_t$ が変化せず、 $K$ 個のアームをもち、
- ・ ユーザ $u_t(\text{context}(x_{t,1}, \dots, x_{t,K}))$ が同じ $K$ -armed banditは、
- ・ 普通のbandit algorithmと変わらない
- ・ このContextual banditの特殊ケースは**context-free bandit**とされる

# 記事推薦のcontext

- ・ 記事をアームのpoolとみなす
- ・ 表示された記事がクリック→ $\text{payoff}=1$ 、他は0
- ・ この定義により  $\text{payoff}$  の期待値はCTRになる
- ・ この研究でのbandit定式化において

CTRを最大化する記事を選択すること



total expected payoffの最大化すること

# WEBサービスでは ユーザ情報獲得が簡単

- ・ 例) 若い男どもは退職後の計画よりもiPodの記事が好きだ
- ・ コンパクトに表される有益な情報からユーザと記事を”summarize”
- ・ そうすることでbandit algorithmはユーザ・記事間のCTR情報を汎用的にでき、
- ・ 特に新規ユーザ・記事に対する良い記事の選択を学習できる

# 探索と搾取のバランス

- ・ bandit問題の基本的な課題は”探索”と”搾取”のバランス
- ・ regretを最小化するために、“A”は最適と思われるアームを過去の経験から**”搾取”**する
- ・ 反対に、もっと情報を得るために準最適なアームも**”探索”**する
- ・ “探索”は短期的regretを増加させるが、長期的regretを減少させる

# Contextual-Free Bandit Algorithm

- $\varepsilon$  - greedy
  - $1 - \varepsilon$  で“搾取”、 $\varepsilon$  で“搾取”
  - “unguided”な探索戦略
- UCB
  - “guided”な探索戦略

$$\underbrace{|\hat{\mu}_{t,a} - \mu_a|}_{\text{payoffの平均} \quad \text{payoff}} < \underbrace{c_{t,a}}_{\text{信頼区間}}$$

$$\underbrace{a_t}_{\text{アーム}} = \arg \max_a (\hat{\mu}_{t,a} + c_{t,a})$$

# こんなものもあるお

- ・ EXP4( $\tilde{O}(\sqrt{T})$ )
- ・ epoch-greedy( $\tilde{O}(T^{2/3})$ )
- ・ LinRel( $\tilde{O}(\sqrt{T})$ )
- ・ Bayesian approach
- ・ 計算量のオーダーが良いお

# 提案手法



# LinUCB

- ・ UCBアルゴリズムと似たようなものを考案
- ・ LinUCB
  - ・ パラメトリックなpayoff関数が線形
  - ・ 信頼区間が閉形式で効率的に計算可能
  - ・ 一般的なものは”expensive”だから

# LinUCB with Disjoint Linear Models

- ・ アーム  $a$  の expected payoff

$$\mathbf{E}[r_{t,a} | \mathbf{x}_{t,a}] = \mathbf{x}_{t,a}^\top \boldsymbol{\theta}_a^*$$

- ・  $\mathbf{x}_{t,a}$ :  $d$ 次元のユーザと記事間の context ベクトル
- ・  $\boldsymbol{\theta}_a^*$ : 未知の係数ベクトル  $d$ : 記事  $a$  の context の次元
- ・ Disjoint  $\rightarrow$  パラメタがアーム間で共有されていない

広告

$$\text{CTR} = \overset{\text{context}}{20\text{代}} * \underset{\text{係数}}{0.01} + 30\text{代} * 0.02 + 40\text{代以上} * 0.1$$

# LinUCB with Hybrid Linear Models

- ・  $m$  を context の数、 $d$  を context の次元とした時

計画行列  $\mathbf{D}_a: \mathbb{R}^{m \times d}$       ターゲット  $\mathbf{b}_a: \in \mathbb{R}^m$

- ・ リッジ回帰により係数を推定

c と b 間違えてる(?)

$$\hat{\boldsymbol{\theta}}_a = (\mathbf{D}_a^\top \mathbf{D}_a + \mathbf{I}_d)^{-1} \mathbf{D}_a^\top \mathbf{c}_a$$

# 信頼区間による行動選択

- ・ [2]により、

$$\left| \mathbf{x}_{t,a}^\top \hat{\boldsymbol{\theta}}_a - \mathbf{E}[r_{t,a} | \mathbf{x}_{t,a}] \right| \leq \alpha \sqrt{\mathbf{x}_{t,a}^\top (\mathbf{D}_a^\top \mathbf{D}_a + \mathbf{I}_d)^{-1} \mathbf{x}_{t,a}}$$
$$\alpha = 1 + \sqrt{\ln(2/\delta)/2}$$

- ・ UCBと同じようにアーム選択をすることができる

$$a_t \stackrel{\text{def}}{=} \arg \max_{a \in \mathcal{A}_t} \left( \mathbf{x}_{t,a}^\top \hat{\boldsymbol{\theta}}_a + \alpha \sqrt{\mathbf{x}_{t,a}^\top \mathbf{A}_a^{-1} \mathbf{x}_{t,a}} \right)$$

expected payoffの標準偏差

$$\mathbf{A}_a \stackrel{\text{def}}{=} \mathbf{D}_a^\top \mathbf{D}_a + \mathbf{I}_d$$

- ・ リッジ回帰の性質などから $\mathbf{x}_{t,a}$ の精度向上の寄与度が評価できる
- ・ 上式はpayoff推定とモデル不確実性の削減の間のtrade-offとみなせる

# LinUCB with Disjoint Linear Models

---

**Algorithm 1** LinUCB with disjoint linear models.

---

```
0: Inputs:  $\alpha \in \mathbb{R}_+$ 
1: for  $t = 1, 2, 3, \dots, T$  do
2:   Observe features of all arms  $a \in \mathcal{A}_t$ :  $\mathbf{x}_{t,a} \in \mathbb{R}^d$ 
3:   for all  $a \in \mathcal{A}_t$  do
4:     if  $a$  is new then
5:        $\mathbf{A}_a \leftarrow \mathbf{I}_d$  ( $d$ -dimensional identity matrix)
6:        $\mathbf{b}_a \leftarrow \mathbf{0}_{d \times 1}$  ( $d$ -dimensional zero vector)
7:     end if
8:      $\hat{\boldsymbol{\theta}}_a \leftarrow \mathbf{A}_a^{-1} \mathbf{b}_a$ 
9:      $p_{t,a} \leftarrow \hat{\boldsymbol{\theta}}_a^\top \mathbf{x}_{t,a} + \alpha \sqrt{\mathbf{x}_{t,a}^\top \mathbf{A}_a^{-1} \mathbf{x}_{t,a}}$ 
10:   end for
11:   Choose arm  $a_t = \arg \max_{a \in \mathcal{A}_t} p_{t,a}$  with ties broken arbitrarily, and observe a real-valued payoff  $r_t$ 
12:    $\mathbf{A}_{a_t} \leftarrow \mathbf{A}_{a_t} + \mathbf{x}_{t,a_t} \mathbf{x}_{t,a_t}^\top$ 
13:    $\mathbf{b}_{a_t} \leftarrow \mathbf{b}_{a_t} + r_t \mathbf{x}_{t,a_t}$ 
14: end for
```

理論的に定常じゃない入力にも有効らしいお…!

# LinUCB with hybrid linear models

- Disjointはアーム固有のfeaturesだけを用いた
- 多くの応用では、アーム間で共有されるfeaturesを用いることは有用→じゃあHybrid modelや

$$\mathbf{E}[r_{t,a} | \mathbf{x}_{t,a}] = \mathbf{z}_{t,a}^\top \boldsymbol{\beta}^* + \mathbf{x}_{t,a}^\top \boldsymbol{\theta}_a^*;$$

- $\mathbf{z}_{t,a}$ : 現在のユーザ・記事の組み合わせのfeatures ( $\mathbb{R}^k$ )  
ちょっとトリッキー(?)
- $\boldsymbol{\beta}$ : すべての記事で共有される係数

# LinUCB with hybrid linear models

- ・ Algorithm 1 はもはや使えない
  - ・ 共有featuresが独立でないから
- ・ でも幸運な事にブロック行列の逆行列のテクニック使えばいける！
- ・ 紙面の関係上擬似コードだけ載せるで（詳細はこれからのFull-paperで）

# LinUCB with hybrid linear models

---

**Algorithm 2** LinUCB with hybrid linear models.
 

---

```

0: Inputs:  $\alpha \in \mathbb{R}_+$ 
1:  $\mathbf{A}_0 \leftarrow \mathbf{I}_k$  ( $k$ -dimensional identity matrix)
2:  $\mathbf{b}_0 \leftarrow \mathbf{0}_k$  ( $k$ -dimensional zero vector)
3: for  $t = 1, 2, 3, \dots, T$  do
4:   Observe features of all arms  $a \in \mathcal{A}_t$ :  $(\mathbf{z}_{t,a}, \mathbf{x}_{t,a}) \in \mathbb{R}^{k+d}$ 
5:    $\hat{\boldsymbol{\beta}} \leftarrow \mathbf{A}_0^{-1} \mathbf{b}_0$ 
6:   for all  $a \in \mathcal{A}_t$  do
7:     if  $a$  is new then
8:        $\mathbf{A}_a \leftarrow \mathbf{I}_d$  ( $d$ -dimensional identity matrix)
9:        $\mathbf{B}_a \leftarrow \mathbf{0}_{d \times k}$  ( $d$ -by- $k$  zero matrix)
10:       $\mathbf{b}_a \leftarrow \mathbf{0}_{d \times 1}$  ( $d$ -dimensional zero vector)
11:    end if
12:     $\hat{\boldsymbol{\theta}}_a \leftarrow \mathbf{A}_a^{-1} (\mathbf{b}_a - \mathbf{B}_a \hat{\boldsymbol{\beta}})$ 
13:     $s_{t,a} \leftarrow \mathbf{z}_{t,a}^\top \mathbf{A}_0^{-1} \mathbf{z}_{t,a} - 2\mathbf{z}_{t,a}^\top \mathbf{A}_0^{-1} \mathbf{B}_a^\top \mathbf{A}_a^{-1} \mathbf{x}_{t,a} +$   

 $\mathbf{x}_{t,a}^\top \mathbf{A}_a^{-1} \mathbf{x}_{t,a} + \mathbf{x}_{t,a}^\top \mathbf{A}_a^{-1} \mathbf{B}_a \mathbf{A}_0^{-1} \mathbf{B}_a^\top \mathbf{A}_a^{-1} \mathbf{x}_{t,a}$ 
14:     $p_{t,a} \leftarrow \mathbf{z}_{t,a}^\top \hat{\boldsymbol{\beta}} + \mathbf{x}_{t,a}^\top \hat{\boldsymbol{\theta}}_a + \alpha \sqrt{s_{t,a}}$ 
15:  end for
16:  Choose arm  $a_t = \arg \max_{a \in \mathcal{A}_t} p_{t,a}$  with ties broken arbitrarily, and observe a real-valued payoff  $r_t$ 
17:   $\mathbf{A}_0 \leftarrow \mathbf{A}_0 + \mathbf{B}_{a_t}^\top \mathbf{A}_{a_t}^{-1} \mathbf{B}_{a_t}$ 
18:   $\mathbf{b}_0 \leftarrow \mathbf{b}_0 + \mathbf{B}_{a_t}^\top \mathbf{A}_{a_t}^{-1} \mathbf{b}_{a_t}$ 
19:   $\mathbf{A}_{a_t} \leftarrow \mathbf{A}_{a_t} + \mathbf{x}_{t,a_t} \mathbf{x}_{t,a_t}^\top$ 
20:   $\mathbf{B}_{a_t} \leftarrow \mathbf{B}_{a_t} + \mathbf{x}_{t,a_t} \mathbf{z}_{t,a_t}^\top$ 
21:   $\mathbf{b}_{a_t} \leftarrow \mathbf{b}_{a_t} + r_t \mathbf{x}_{t,a_t}$ 
22:   $\mathbf{A}_0 \leftarrow \mathbf{A}_0 + \mathbf{z}_{t,a_t} \mathbf{z}_{t,a_t}^\top - \mathbf{B}_{a_t}^\top \mathbf{A}_{a_t}^{-1} \mathbf{B}_{a_t}$ 
23:   $\mathbf{b}_0 \leftarrow \mathbf{b}_0 + r_t \mathbf{z}_{t,a_t} - \mathbf{B}_{a_t}^\top \mathbf{A}_{a_t}^{-1} \mathbf{b}_{a_t}$ 
24: end for
    
```

リッジ回帰

リッジ回帰

信頼区間計算



# 評価方法

# オフラインで評価したい

- ・ でも過去ログは選択された記事だけのログだったり、試したいアルゴリズムと違ったポリシーで貯められてる場合が多いから何も考えないでは使えない
- ・ 過去に配信していない記事は評価できない
- ・ シミュレータを作ってもいいがbiasが入りやすい
- ・ 著者らはログデータを用いてunbiasな方法を提案

# 評価アルゴリズム

0行目(入力)

T: 得たい履歴数

$\pi$ : 使用Algorithm

ログ

1,2行目

得た履歴を格納する配列

payoffの初期化

---

**Algorithm 3** Policy\_Evaluator.

---

```
0: Inputs:  $T > 0$ ; policy  $\pi$ ; stream of events
1:  $h_0 \leftarrow \emptyset$  {An initially empty history}
2:  $R_0 \leftarrow 0$  {An initially zero total payoff}
3: for  $t = 1, 2, 3, \dots, T$  do
4:   repeat
5:     Get next event  $(\mathbf{x}_1, \dots, \mathbf{x}_K, a, r_a)$ 
6:   until  $\pi(h_{t-1}, (\mathbf{x}_1, \dots, \mathbf{x}_K)) = a$ 
7:    $h_t \leftarrow \text{CONCATENATE}(h_{t-1}, (\mathbf{x}_1, \dots, \mathbf{x}_K, a, r_a))$ 
8:    $R_t \leftarrow R_{t-1} + r_a$ 
9: end for
10: Output:  $R_T/T$ 
```

---

# 評価アルゴリズム

ログ(例)

1/1, ユーザA, 広告A, Click  
1/1, ユーザA, 広告B, Non-click  
1/1, ユーザB, 広告A, Click  
1/1, ユーザC, 広告A, Non-click  
1/1, ユーザD, 広告B, Non-click  
1/1, ユーザC, 広告B, Non-click  
1/1, ユーザA, 広告A, Click  
1/1, ユーザB, 広告B, Non-click

---

**Algorithm 3** Policy\_Evaluator.

---

0: Inputs:  $T > 0$ ; policy  $\pi$ ; stream of events  
1:  $h_0 \leftarrow \emptyset$  {An initially empty history}  
2:  $R_0 \leftarrow 0$  {An initially zero total payoff}  
3: **for**  $t = 1, 2, 3, \dots, T$  **do**  
4:     **repeat**  
5:         Get next event  $(\mathbf{x}_1, \dots, \mathbf{x}_K, a, r_a)$   
6:     **until**  $\pi(h_{t-1}, (\mathbf{x}_1, \dots, \mathbf{x}_K)) = a$   
7:      $h_t \leftarrow \text{CONCATENATE}(h_{t-1}, (\mathbf{x}_1, \dots, \mathbf{x}_K, a, r_a))$   
8:      $R_t \leftarrow R_{t-1} + r_a$   
9: **end for**  
10: Output:  $R_T/T$

---

# 評価アルゴリズム

t=1

過去の履歴とユーザAの  
contextを基に

$\pi$

履歴

ログを順番に見ていく

ログ(例)



1/1, ユーザA, 広告A, Click

1/1, ユーザA, 広告B, Non-click

1/1, ユーザB, 広告A, Click

1/1, ユーザC, 広告A, Non-click

1/1, ユーザD, 広告B, Non-click

1/1, ユーザC, 広告B, Non-click

1/1, ユーザA, 広告A, Click

1/1, ユーザB, 広告B, Non-click

# 評価アルゴリズム

t=1

過去の履歴とユーザAの  
contextを基に

$\pi$  → 広告B出す！

実際は広告Aを出してるから  
広告Bの反応が見れないので追加しない

履歴

ログ(例)

1/1, ユーザA, 広告A, Click  
1/1, ユーザA, 広告B, Non-click  
1/1, ユーザB, 広告A, Click  
1/1, ユーザC, 広告A, Non-click  
1/1, ユーザD, 広告B, Non-click  
1/1, ユーザC, 広告B, Non-click  
1/1, ユーザA, 広告A, Click  
1/1, ユーザB, 広告B, Non-click

# 評価アルゴリズム

$t=1$

過去の履歴とユーザAの  
contextを基に

$\pi$

履歴

$t$ はそのまま次ログに行く

ログ(例)

1/1, ユーザA, 広告A, Click

→ 1/1, ユーザA, 広告B, Non-click

1/1, ユーザB, 広告A, Click

1/1, ユーザC, 広告A, Non-click

1/1, ユーザD, 広告B, Non-click

1/1, ユーザC, 広告B, Non-click

1/1, ユーザA, 広告A, Click

1/1, ユーザB, 広告B, Non-click

# 評価アルゴリズム

t=1

過去の履歴とユーザAの  
contextを基に

$\pi$  → 広告B出す！

実際に広告Bを出してるから  
履歴に追加する

履歴

ユーザA, 広告B, Non-click

ログ(例)

1/1, ユーザA, 広告A, Click

1/1, ユーザA, 広告B, Non-click

1/1, ユーザB, 広告A, Click

1/1, ユーザC, 広告A, Non-click

1/1, ユーザD, 広告B, Non-click

1/1, ユーザC, 広告B, Non-click

1/1, ユーザA, 広告A, Click

1/1, ユーザB, 広告B, Non-click



# 評価アルゴリズム

t=2

過去の履歴とユーザBの  
contextを基に

$\pi$  → 広告A出す！

実際に広告Aを出してるから  
履歴に追加する

履歴

ユーザA, 広告B, Non-click

ユーザB, 広告A, Click

追加したらt+1

ログ(例)

1/1, ユーザA, 広告A, Click

1/1, ユーザA, 広告B, Non-click

→ 1/1, ユーザB, 広告A, Click

1/1, ユーザC, 広告A, Non-click

1/1, ユーザD, 広告B, Non-click

1/1, ユーザC, 広告B, Non-click

1/1, ユーザA, 広告A, Click

1/1, ユーザB, 広告B, Non-click

# 評価アルゴリズム

t=3

過去の履歴とユーザCの  
contextを基に

$\pi$  → 広告A出す！

実際に広告Aを出してるから  
履歴に追加する

履歴

ユーザA, 広告B, Non-click

ユーザB, 広告A, Click

ユーザC, 広告A, Non-click

ログ(例)

1/1, ユーザA, 広告A, Click

1/1, ユーザA, 広告B, Non-click

1/1, ユーザB, 広告A, Click

→ 1/1, ユーザC, 広告A, Non-click

1/1, ユーザD, 広告B, Non-click

1/1, ユーザC, 広告B, Non-click

1/1, ユーザA, 広告A, Click

1/1, ユーザB, 広告B, Non-click

# 評価アルゴリズム

t=100

最終的な履歴を見て  
 $CTR = 1/3 = 33.3\%$

## 履歴

ユーザA, 広告B, Non-click

ユーザB, 広告A, Click

ユーザC, 広告A, Non-click

ログ(例)

1/1, ユーザA, 広告A, Click

1/1, ユーザA, 広告B, Non-click

1/1, ユーザB, 広告A, Click

1/1, ユーザC, 広告A, Non-click

1/1, ユーザD, 広告B, Non-click

1/1, ユーザC, 広告B, Non-click

1/1, ユーザA, 広告A, Click

1/1, ユーザB, 広告B, Non-click

...



# 数学的な証明

- ・ (定理)長さ $T$ の履歴を得るために必要なログデータの期待値は $KT$
- ・  $K$ はアーム、記事の数
- ・ 証明されてる。

# 評価方法の感想

賢いちゃん

実験

setup



# Yahoo! Today-Module

ここに何の広告を  
出すか決める

**Featured** | Entertainment | Sports | Life

**McNair's final hours revealed**  
**STORY**  
Police release 50 text messages that depict the late NFL player's alleged killer as losing control. » [Details](#)  
• [UConn murder victim mourned](#)  
🔍 [Find Steve McNair murder case](#)

**F1** Steve McNair's final hours revealed

**F2** Cindy Crawford stays fierce in a black mini

**F3** Watch for dozens of 'shooting stars' tonight

**F4** At team's big moment, star player isn't around

» More: [Featured](#) | [Buzz](#)

F1に出される広告は  
STORYにでかく表示される

# Data Collection

- Tuning data: 1/May/2009
  - 最適なパラメタ決定に使用
  - 4.7 million events
- Evaluation data: 03-09/May/2009
  - 36 million events

# Data Collection

- ・ ユーザのインタラクションイベントは3つ
- ・ ランダムに選択された記事
- ・ ユーザ・記事の情報
- ・ ユーザがSTORYにあるその広告をクリックしたか

# Feature Construction

- ・ 0.1以上の“支持度”をもつ特徴量を抜き出す
- ・ 各**ユーザ**は1000以上のraw featuresで表される
  - ・ demographic: 性別、年齢（10個にセグメント）
  - ・ geographic: 世界の200都市、U.S.州
  - ・ behavioral: Yahoo!でのユーザの消費をまとめた1000バイナリカテゴリ

# Feature Construction

- ・ 各記事は約100個のfeatureをもつ
  - ・ URL カテゴリ
  - ・ 編集カテゴリ: エディタによるタグ付けトピック
- ・ [3]によりfeatureをbinaryにencode
- ・ 最後[0, 1]にnormalize
- ・ ☆ 値が1のconstant featureも追加  
重要そう (感想)

# Feature Reduction

- ・ 現実的には小さな特徴空間でやらないと辛いだろう
- ・ ということで[4]使って、ユーザ特徴量を記事カテゴリに射影して、K-means使ってクラスタ化した
  - ・ コンジョイント分析
- ・ ユーザと記事をそれぞれ5個にセグメント

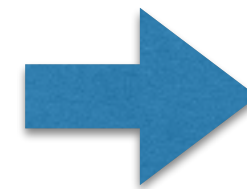
5つのユーザセグメント

$$u \begin{bmatrix} 0.1 & 0.9 & 0.0 & 0.0 & 0.0 \\ 0.1 & 0.3 & 0.3 & 0.2 & 0.1 \\ \dots & & & & \end{bmatrix}$$

5つの記事セグメント

$$a \begin{bmatrix} 0.0 & 0.4 & 0.0 & 0.4 & 0.1 \\ 0.9 & 0.0 & 0.0 & 0.1 & 0.0 \\ \dots & & & & \end{bmatrix}$$

これに1ベクトル  
足して外積



$Z_{t,a}$   
共有特徴量

比較するアルゴリズム

# 比較するアルゴリズム

## 1. Context-freeなAlgorithm

- ・ random
- ・  $\epsilon$  - greedy
- ・ UCB
- ・ omniscient(後知恵から得る過去一番のCTR記事)

## 2. Algorithm with “warm-start”

- ・ Context-freeなアルゴにユーザ固有のCTR情報を足し合わせる

## 3. オンラインなユーザ固有のAlgorithm

- ・  $\epsilon$  - greedy, UCB (セグメント分けされたユーザ使用)
- ・  $\epsilon$  - greedy, UCB (disjoint & hybrid)

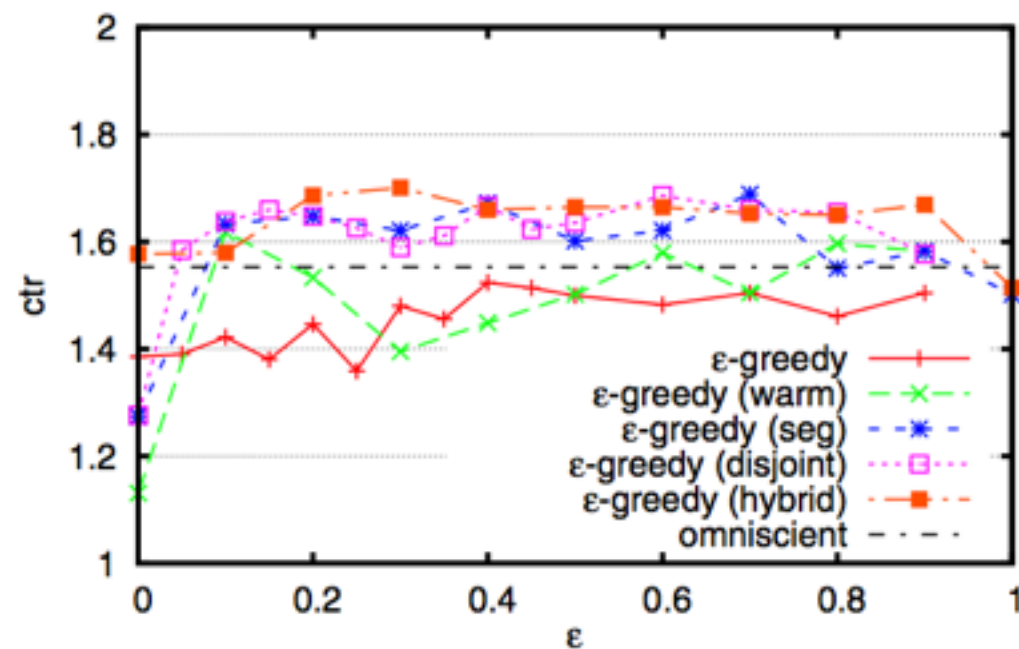


Performance Metric

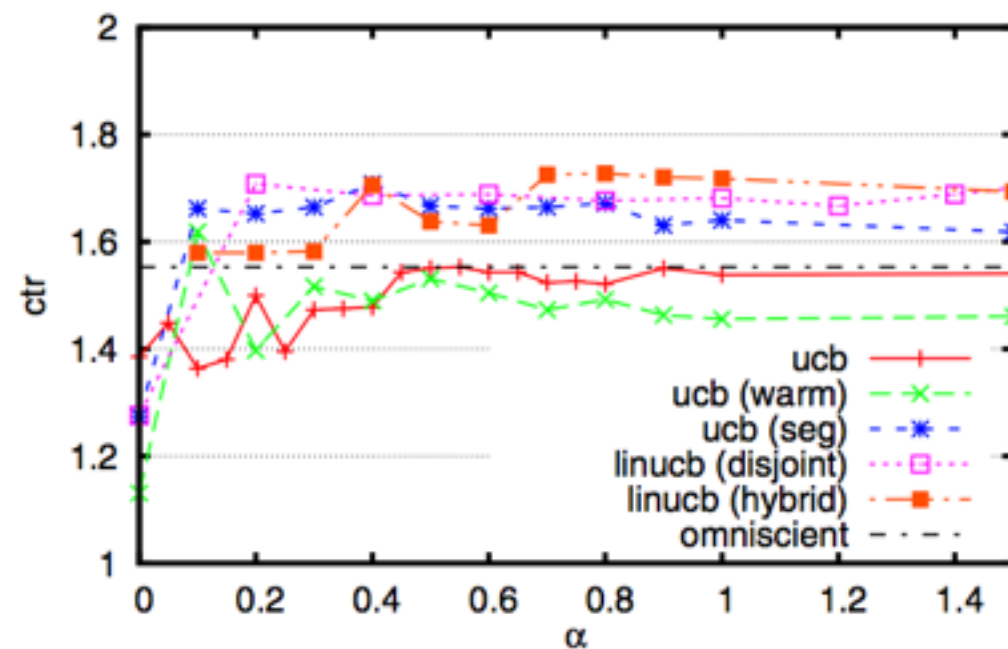
# Performance Metric

- ・ ページへの全トラフィックをランダムに2つに分割
  1. learning bucket : アルゴリズム学習用
  2. deployment bucket : アルゴリズム実験用
- ・ deployの方が重要だから大きなbucketになりがち
- ・ だけど、learningもより速い学習(より小さな regret)を示唆するので両方の結果を示す

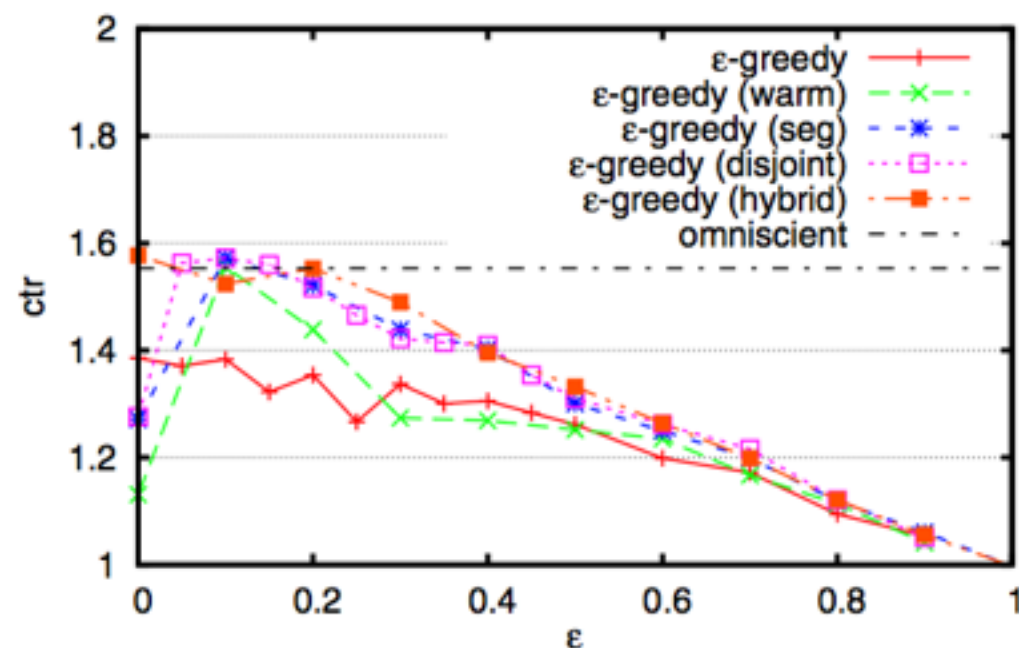
# チューニングデータでの結果



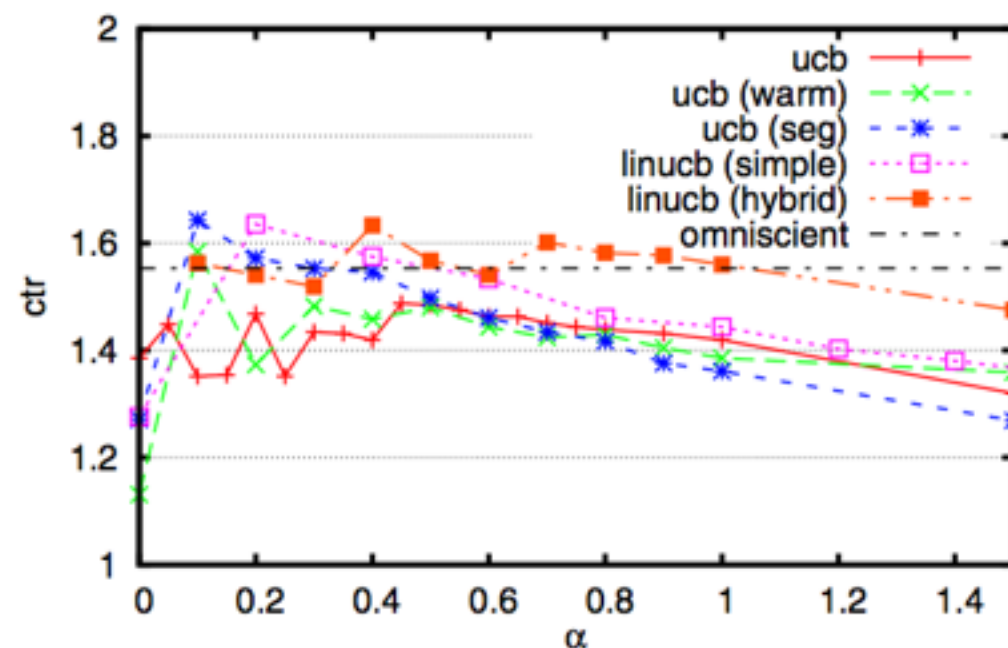
(a) Deployment bucket.



(b) Deployment bucket.



(c) Learning bucket.



(d) Learning bucket.

全オンライン学習アルゴリズムはOmniscientに勝ってる

# 評価データでの結果

全データ



よりスパースなデータ

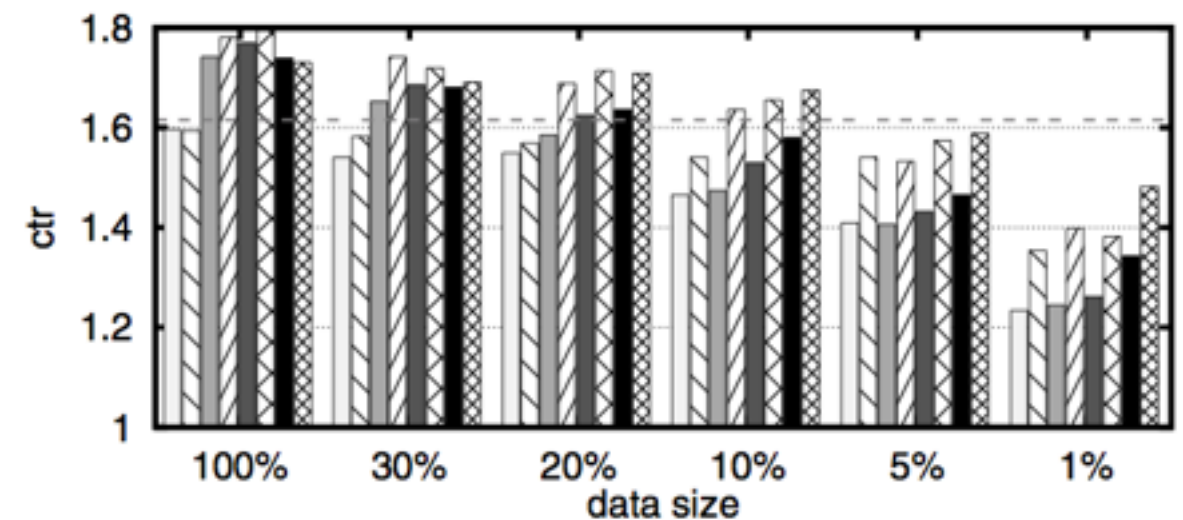


algorithm	size = 100%		size = 30%		size = 20%		size = 10%		size = 5%		size = 1%	
	deploy	learn	deploy	learn	deploy	learn	deploy	learn	deploy	learn	deploy	learn
$\epsilon$ -greedy	1.596 0%	1.326 0%	1.541 0%	1.326 0%	1.549 0%	1.273 0%	1.465 0%	1.326 0%	1.409 0%	1.292 0%	1.234 0%	1.139 0%
ucb	1.594 0%	1.569 18.3%	1.582 2.7%	1.535 15.8%	1.569 1.3%	1.488 16.9%	1.541 5.2%	1.446 9%	1.541 9.4%	1.465 13.4%	1.354 9.7%	1.22 7.1%
$\epsilon$ -greedy (seg)	1.742 9.1%	1.446 9%	1.652 7.2%	1.46 10.1%	1.585 2.3%	1.119 -12%	1.474 0.6%	1.284 -3.1%	1.407 0%	1.281 -0.8%	1.245 0.9%	1.072 -5.8%
ucb (seg)	1.781 11.6%	1.677 26.5%	1.742 13%	1.555 17.3%	1.689 9%	1.446 13.6%	1.636 11.7%	1.529 15.3%	1.532 8.7%	1.32 2.2%	1.398 13.3%	1.25 9.7%
$\epsilon$ -greedy (disjoint)	1.769 10.8%	1.309 -1.2%	1.686 9.4%	1.337 0.8%	1.624 4.8%	1.529 20.1%	1.529 4.4%	1.451 9.4%	1.432 1.6%	1.345 4.1%	1.262 2.3%	1.183 3.9%
linucb (disjoint)	1.795 12.5%	1.647 24.2%	1.719 11.6%	1.507 13.7%	1.714 10.7%	1.384 8.7%	1.655 13%	1.387 4.6%	1.574 11.7%	1.245 -3.5%	1.382 12%	1.197 5.1%
$\epsilon$ -greedy (hybrid)	1.739 9%	1.521 14.7%	1.68 9%	1.345 1.4%	1.636 5.6%	1.449 13.8%	1.58 7.8%	1.348 1.7%	1.465 4%	1.415 9.5%	1.342 8.8%	1.2 5.4%
linucb (hybrid)	1.73 8.4%	1.663 25.4%	1.691 9.7%	1.591 20%	1.708 10.3%	1.619 27.2%	1.675 14.3%	1.535 15.8%	1.588 12.7%	1.507 16.6%	1.482 20.1%	1.446 27%

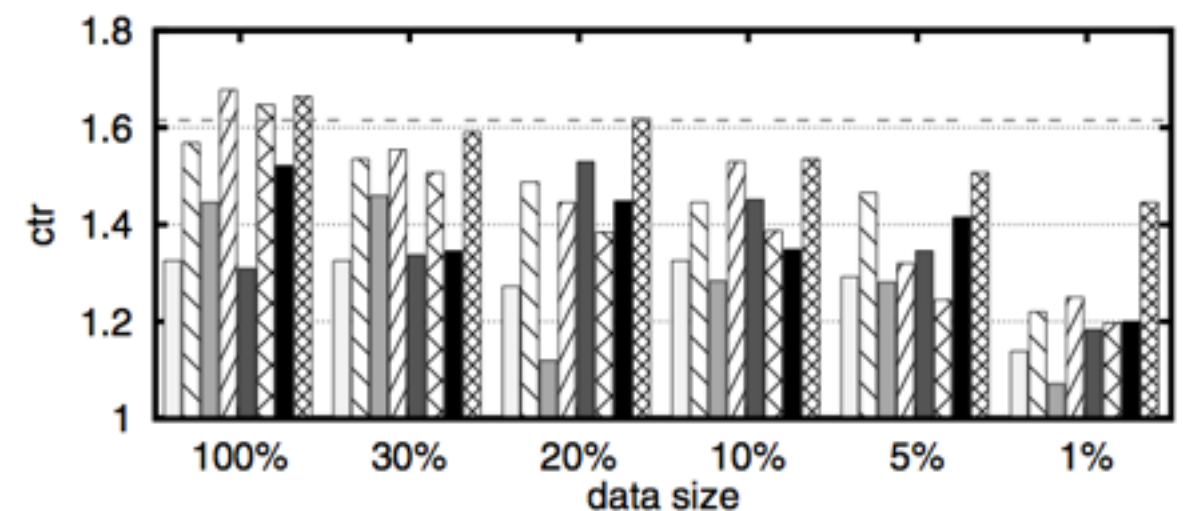
random: 1, omniscient: 1.615

# スパースなデータ

- ・ 下がり気味
- ・ UCB系強い
- ・ hybridは記事間の情報を共有してるから強い



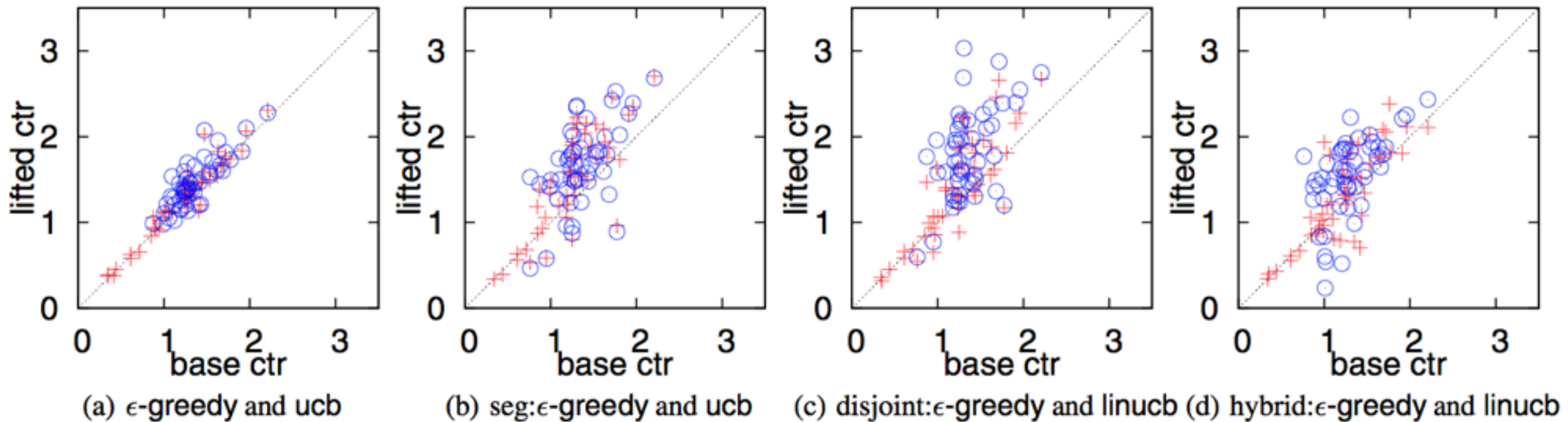
(a) CTRs in the deployment bucket.



(b) CTRs in the learning bucket.



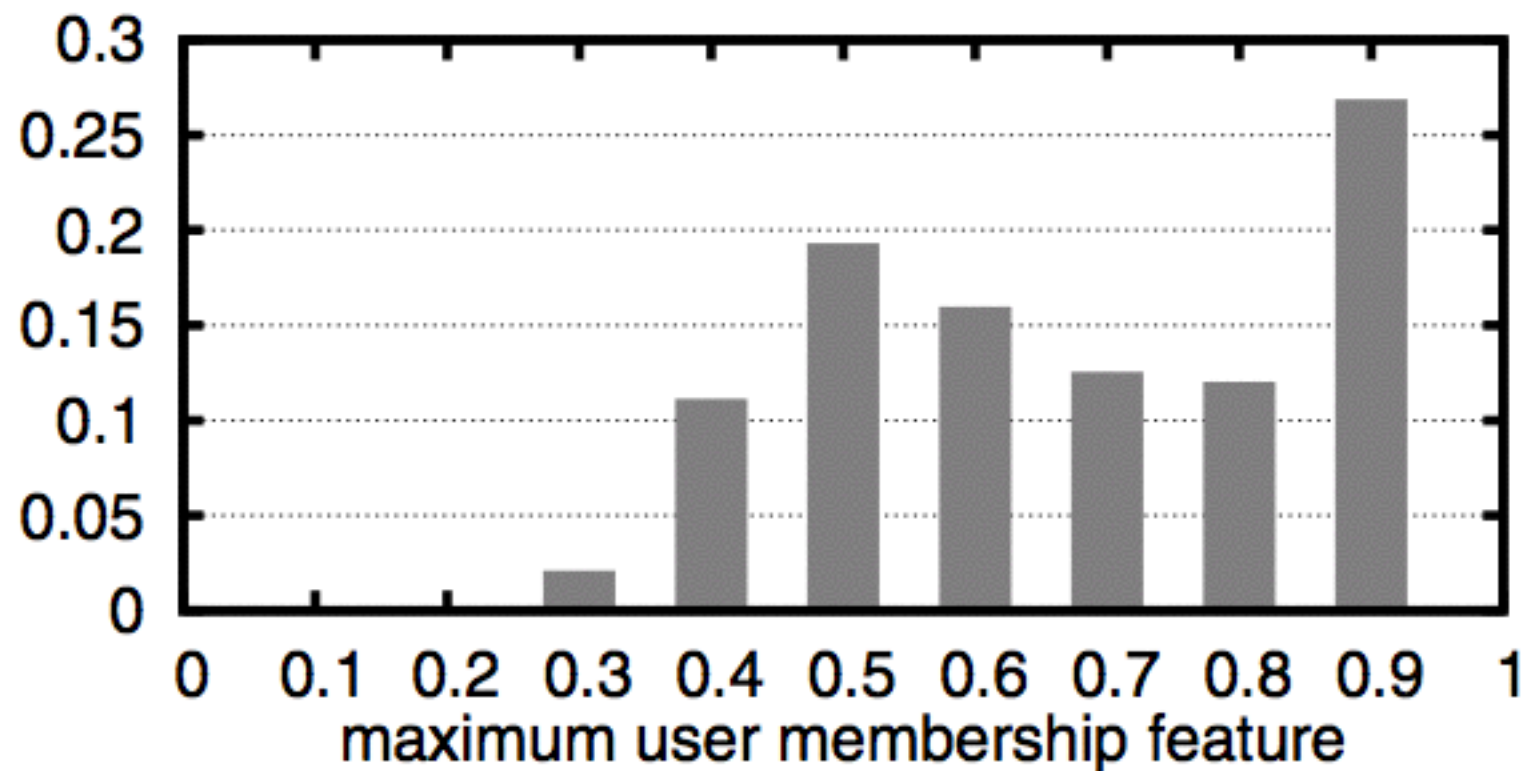
# Context-freeとContextの比較



- ・ 評価データの中で最も選ばれた50個の記事を使用
- ・ Personalized recommendationは有効
- ・ Base ctrはcontext-freeなそのアルゴリズム

# 一応ユーザの分割も良かったよ

- ・ 各ユーザのメンバーシップの最大値のヒストグラム
- ・ 85%が0.5以上、40%が0.8以上



**Figure 5: User maximum membership histogram.**

# 結論



# 結論

- ・ contextual-banditな方法をニュースレコメンドのようなパーソナライズされたWEBサービスに適用した
- ・ LinUCBとシンプルで信頼性のある評価方法提案した

# 参考文献

- Lihong Li et al., 2010, A Contextual-Bandit Approach to Personalized News Article Recommendation
- <http://www.slideshare.net/tsubosaka/contextual-bandit>