

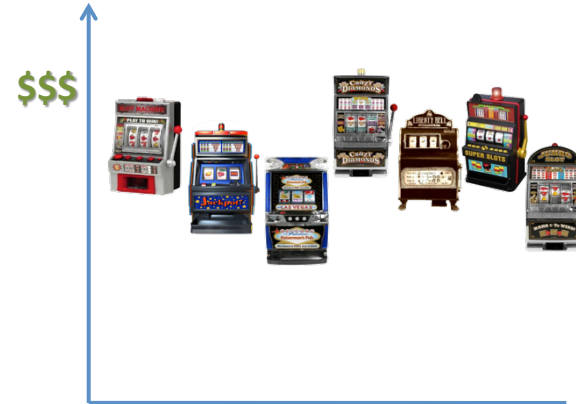
## Multi-Armed Bandit Applications to Marketing Experiments



Eric M. Schwartz  
ericmsch@umich.edu  
[ericmichaelschwartz.com](http://ericmichaelschwartz.com)

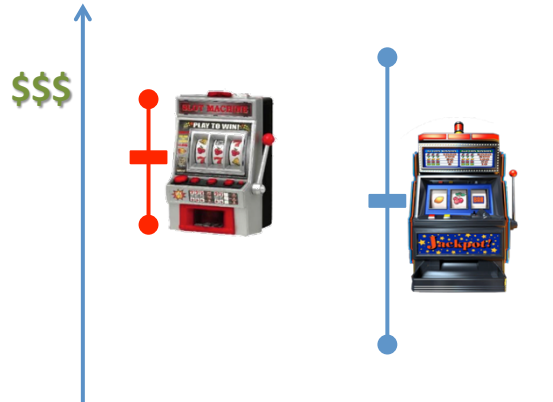


## Introducing the multi-armed bandit



2

## Introducing the multi-armed bandit



3

## Bandit Problem: Key Ideas

### Optimization

- Sequentially select actions to maximize reward (minimize regret)

### Classic Tradeoff

- Learn vs. Earn (Explore vs. Exploit)
- Long-run vs. Immediate payoff



4

## Why is a bandit problem different?

Typical machine learning

- One sample for inference
- “Offline learning”

Bandit problem

- Adaptive sampling
- Active learning
- Partial information / limited feedback
  - Given current data, what data should we collect to optimize an outcome?

## Modern business experimentation

Harvard  
Business  
Review

### How to Design Smart Business Experiments

by Thomas H. Davenport

FROM THE FEBRUARY 2009 ISSUE

INNOVATION

### A Step-by-Step Guide to Smart Business Experiments

by Eric T. Anderson and Duncan Simester

FROM THE MARCH 2011 ISSUE

INNOVATION

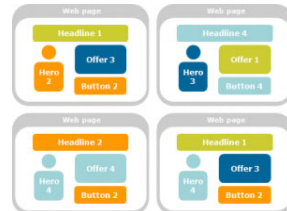
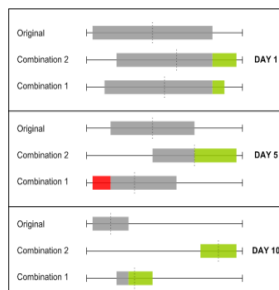
### The Discipline of Business Experimentation

by Stefan Thomke and Jim Manz

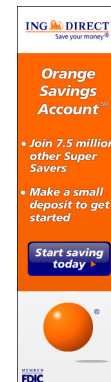
FROM THE DECEMBER 2014 ISSUE

APT  
TUCKER SCHOOL

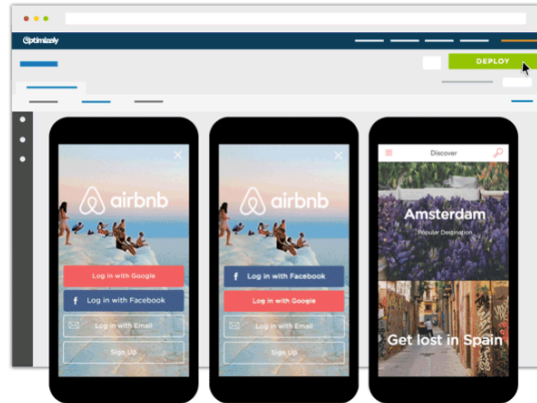
## A/B testing



## Current examples: Online ads

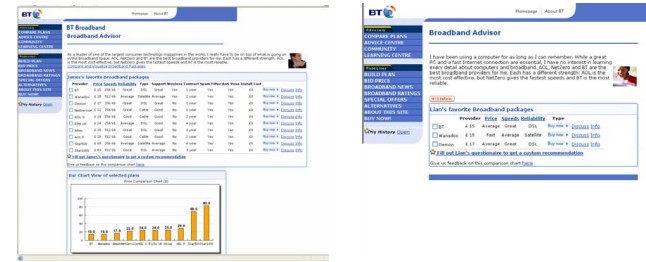


## Current examples: App design



## Current examples: Website style

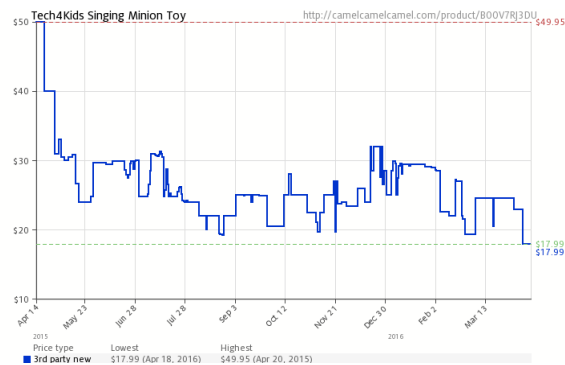
### Comparison of Two Morphs for a Website Advisor



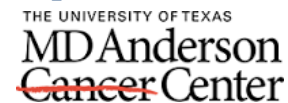
(a) General content, large-load, graphical morph

(b) Focused, small-load, verbal morph

## Current examples: Pricing



## Current examples: Medial clinical trials



- Patients have life threatening disease
  - Patients receive one of two different treatments in a clinical trial
- Ethical dilemma: Collective v Every Individual
  - Learn the best treatment (Clinical Research)
  - Treat the patient better (Clinical Practice)

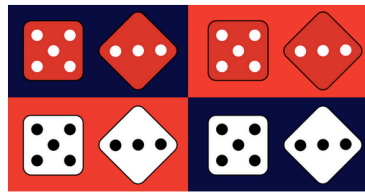
## Current examples



### Netflix Shares a Decade of A/B Test Learning 2015 SXSW Interactive Session Spotlight: A/B Testing at Scale – Minimizing UI Complexity

**Big experiments: Big data's friend for making decisions**

By Eitan Baron on Thursday, April 3, 2014 at 1:31pm [↗](#)



## Current examples



"We're in a world now where you can't wait a whole day to get your data. This has been a huge boon for the growth of our business."



"Testing is a way of life for us. We're increasing our testing resources, and building testing into our long-term, strategic plans."



## Big Data?

### *The Atlantic* The Big Data Boom Is the Innovation Story of Our Time

AN EXPERIMENT EVERY SECOND ERIK BRYNJOLFSSON AND ANDREW MCAFEE | NOV 21 2011, 9:50 AM ET

Science has been dominated by the experimental approach for nearly 400 years. Running controlled experiments is the gold standard for sorting out cause and effect. But experimentation has been difficult for businesses throughout history because of cost, speed and convenience. It is only recently that businesses have learned to run real-time experiments on their customers. The key enabler was the Web.

Consider two versions of Amazon's new Kindle Fire. Using this new service or design, statistically

**Hal Varian: Google "is running on the order of 100-200 experiments on any given day, as they test new products and services, new algorithms and alternative designs."**

This ability and approach model hypothesis an answer

**Greg Linden (Amazon): "Constant, ubiquitous experimentation is the most important thing."**

According to Google economist Eric Varian, his company is running on the order of 100-200 experiments on any given day, as they test new products and services, new algorithms and alternative designs. An iterative review process aggregates findings and frequently leads to further rounds of more targeted experimentation.

## Real-time experiment methods in practice

- A/B testing (A/B/C/.../n)
- Multivariate testing (MVT)
- Stopping rules
- Multi-armed bandit experiments
  - Adaptive allocation of observations

## Example

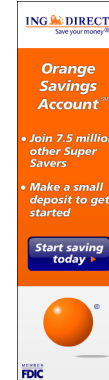
### Customer Acquisition via Display Advertising Using Multi-Armed Bandit Experiments

Schwartz, Bradlow, and Fader (2016)

Marketing Science, forthcoming



## Which ad will bring in the most new customers?



*Always be testing and learning*

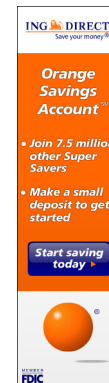
*Always be earning while learning*



## How should we allocate impressions across many ads (served on many websites) to acquire more customers?



*Always be testing and learning*

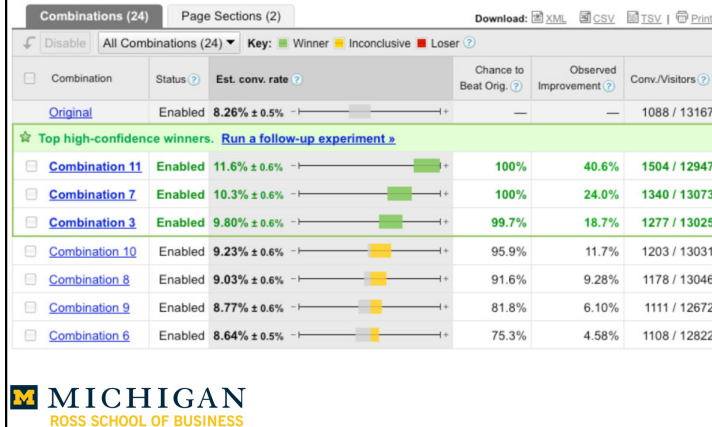


*Always be earning while learning*

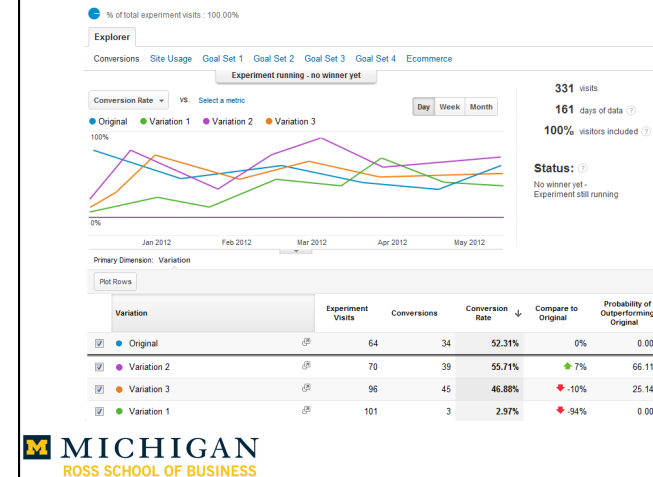


20

## Typical analytics for A/B tests or MVT

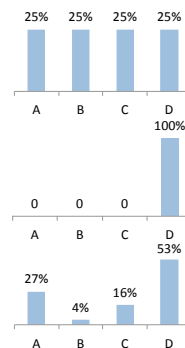


## Google Content Experiments



## Managing the multi-armed bandit

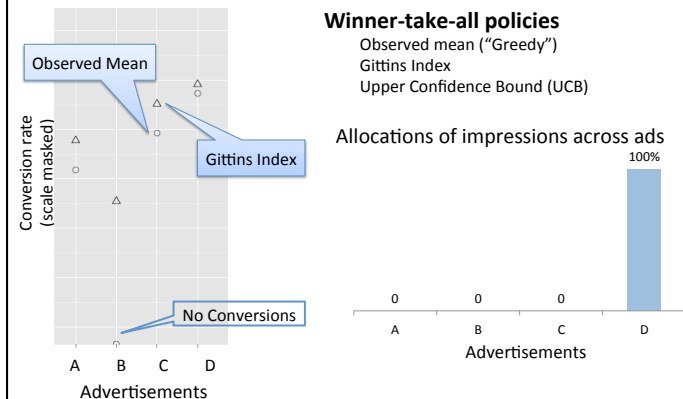
1. Static, balanced design (equal allocation)
2. Adaptive, "greedy" methods (winner take all)
3. Adaptive, randomized (smooth allocation)



Agarwal et al. 2008; Auer et al. 2002; Bertsimas and Mersereau 2007; Lai 1987; Scott 2010; Rumsmevichientong and Tsitsiklis 2010

Gonul and Shi 1998; Gonul and ter Hofstede 2006; Hauser et al. 2009; Montoya et al. 2011; Simester et al. 2006; Sun et al. 2006; Urban et al. 2014

## Managing the multi-armed bandit

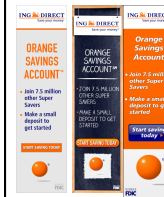


## List of bandit algorithms

- Learn-then-earn
- epsilon-Greedy
- Logit/Boltzman transform (Softmax)
- Exponential Weight (EXP3)
- \*Gittins index\*
- \*Upper Confidence Bound (UCB)\*
- \*Thompson Sampling\*

## The problem we can solve with Gittins index

### Stylized bandit problem



$x_k$  are  $K$ -dimensional indicator vectors

$$\theta = (\mu_1, \dots, \mu_K)$$

$$J = 1$$

$$M_{it} = 1 \text{ for all } t$$

$$\mu_k = E_f(y|x_k)$$

$$Y_{kt} \sim f(y|x_k)$$

$$Y_{kt} = R(A_t, \mu)$$

Objective function

$$\max_{A_t \in 1, \dots, K} \int_{\mu_1} \dots \int_{\mu_K} E_f \left\{ \sum_{t=1}^{\infty} \gamma^t R(A_t, \mu) \right\} p(\mu_1) \dots dp(\mu_K) d\mu_1 \dots d\mu_K$$

**Solved! Exactly optimal policy = Gittins index policy**

Gittins 1979, 1989; Gittins and Jones 1974; Tsitsiklis 1986

Independent actions.  
Not attribute-based.

No hierarchical structure.  
One context.

No batched decisions.  
One-at-a-time.

## The problem we can solve with Gittins index

Stylized bandit problem (ex. Bernoulli bandit with beta priors)



$$a_{kt} = a_{k0} + \sum_{\tau=1}^t y_{k\tau} \text{ and } b_{kt} = b_{k0} + \sum_{\tau=1}^t (m_{k\tau} - y_{k\tau})$$

$$\Pr(Y_{kt} = 1 | a_{kt}, b_{kt}) = E_{p(\mu_k)}(\mu | a_{kt}, b_{kt}) = \frac{a_{kt}}{a_{kt} + b_{kt}}$$

$$V(a_{kt}, b_{kt}, \gamma) = \max \left\{ \frac{G_{kt}}{1 - \gamma}, \frac{[1 + \gamma V(a_{kt} + 1, b_{kt}, \gamma)] \frac{a_{kt}}{a_{kt} + b_{kt}} + [0 + \gamma V(a_{kt}, b_{kt} + 1, \gamma)] \frac{b_{kt}}{a_{kt} + b_{kt}}}{1 - \gamma} \right\}$$

$K$  separable "one-and-a-half-armed bandit"

- $G$  is the Gittins index, which is exactly optimal value.
- Represents "certainty equivalent," "option value," etc.
- Policy: At any state, play the action with the highest Gittins index.

## Gittins index and UCB

Stylized bandit problem

**Gittins index (exactly optimal solution)**

Expected Immediate Reward

Expected Option Value for Learning

$$\max \left\{ \frac{G_{kt}}{1 - \gamma}, \frac{a_{kt}}{a_{kt} + b_{kt}} + \gamma \left[ V(a_{kt} + 1, b_{kt}, \gamma) \frac{a_{kt}}{a_{kt} + b_{kt}} + V(a_{kt}, b_{kt} + 1, \gamma) \frac{b_{kt}}{a_{kt} + b_{kt}} \right] \right\}$$

**Upper confidence bound (asymptotically optimal)**

$$UCB1 = \frac{1}{n_k(t)} \sum_{i=1}^{n_k(t)} y_{ki} + \sqrt{\frac{2 \log(t)}{n_k(t)}} \quad \text{Regret}_t = \sum_{\tau=1}^t (\mu_* - \mu_{A_\tau})$$

Agarwal et al. 2008; Auer et al. 2002; Brezzi and Lai 2002; Lai 1987

## The problems we can solve

### Attribute-based bandit problem

(still no unobserved context heterogeneity, no batching)



29

## The problems we can solve

### Attribute-based bandit problem

(still no unobserved context heterogeneity  $J=1$ , no batching)

Use spillover learning to leverage similarity of actions (regression framework) and predict performance of not yet chosen actions (a la conjoint).

$$\text{UCB-GLM}_{kt} = \text{link}^{-1}(x_k^t \beta_t) + \|x_k\|_{Q_{t-1}^{-1} \rho(t)}$$

$$\|\beta - \hat{\beta}_t\|_{Q_t} \leq \rho^*(t)$$

$$\|x_k\|_{Q_t^{-1}} = \sqrt{x_k^t Q_t^{-1} x_k}$$

$$Q_t = \sum_{\tau=1}^{t-1} x_{A_\tau} x_{A_\tau}^t$$

Linear regression structure

" $(X'X)^{-1}$ "

Nests the UCB1 algorithm and any myopic GLM (regression model) without learning.

Minimizes "regret" with high probability (Dani et al. 2008; Filippi et al. 2010; Rusmevich et al. 2010; Tsitsiklis 2010)

30

## Life is more complicated...

- Natural extensions and complications in testing
  - Large set of candidates to compare (e.g., creative content, display ads)
  - Very rare events (e.g., transactions or customer acquisitions via display ads)
  - Batches of decisions (e.g., "chunky" allocations)
  - Different contexts (e.g., websites, customer segments differ)
  - Long-run customer value (e.g., impact of actions beyond one transaction)

## The problem we want to solve

How should we **allocate** impressions across many ads served on many websites to acquire more customers?



Hierarchical attribute-based bandit

32

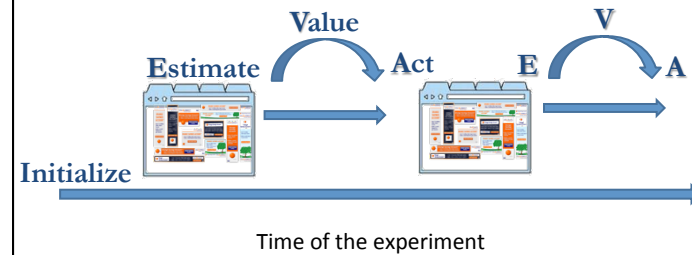


## Overview of solution

- Improve efficiency of online advertising
- Broaden the class of earn-and-learn problems through managing multi-armed bandit with
  - heterogeneity (hierarchical structure)
  - attribute structure (non-independent actions)
  - batched decisions
  - rare events
 by extending the use of Thompson Sampling.
- Document performance of methods
  - Field experiment implementing policy
  - Counterfactual simulations

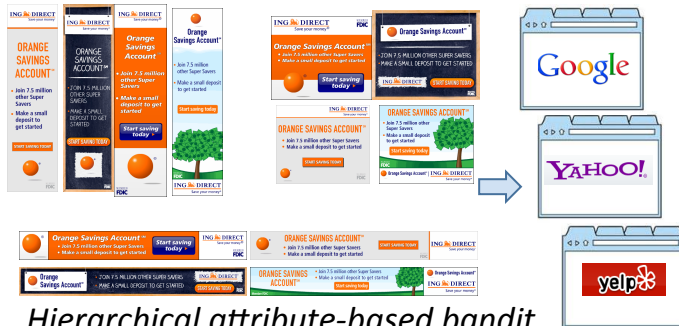
## Adaptive experiment timeline

Data are conversions and impressions for each website-ad combo.



## The problem we want to solve

How should we **allocate** impressions across many ads served on many websites to acquire more customers?



*Hierarchical attribute-based bandit*

## Preliminaries for optimization problem

*Hierarchical attribute-based K-armed bandit with J contexts and batching*

$$\{y_{jkt}, \dots, y_{jkt}\} \sim f(y|x_k, \beta_j) \forall j, k$$

Attribute-based

$$x_k = (x_{k1}, \dots, x_{kd})$$

Hierarchical structure

$$\beta_j \sim g(\bar{\beta}, \Sigma) \forall j \quad \theta = (\{\beta_j\}_1^J, \bar{\beta}, \Sigma)$$

$$\mu_{jk} = E_f(y|x_k, \beta_j) = \text{link}^{-1}(x'_k \beta_j) \forall j, k$$

Batched decisions

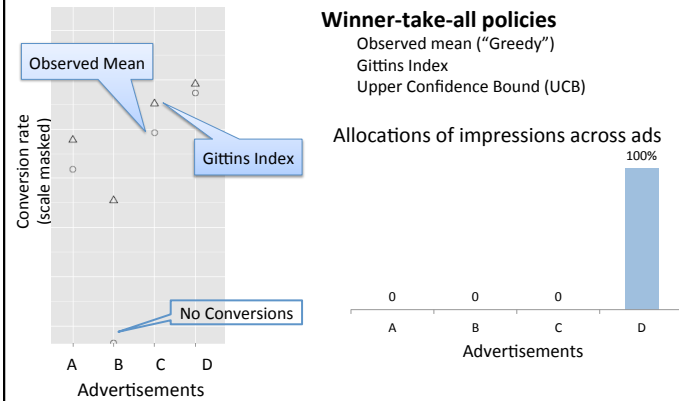
$$1 = \sum_{k=1}^K w_{jkt} \forall j, t \quad \text{batch of } M_{jt} \text{ observations}$$

$$R(w_{jt}, \beta_j) = \sum_{k=1}^K Y_{jkt}$$

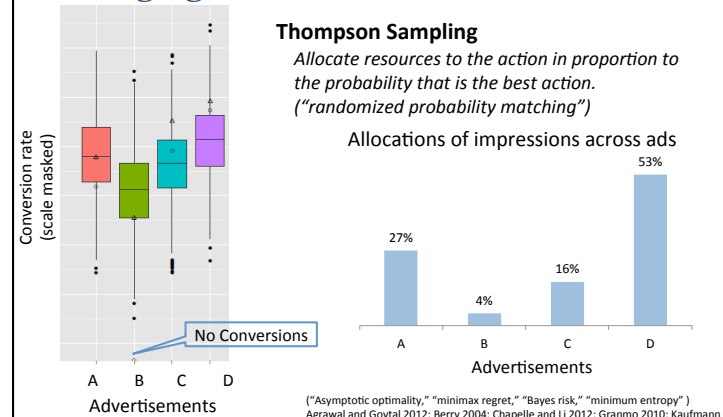
$$E_f R(w_{jt}, \beta_j) = \sum_{k=1}^K w_{jkt} M_{jt} \mu_{jk} = \sum_{k=1}^K w_{jkt} M_{jt} \text{link}^{-1}(x'_k \beta_j) \forall j.$$

$$\text{Regret}_T(w) = \sum_{t=1}^T \sum_{j=1}^J M_{jt} \left( \mu_{j*} - \sum_k w_{jkt} \mu_{jk} \right)$$

## Managing the multi-armed bandit



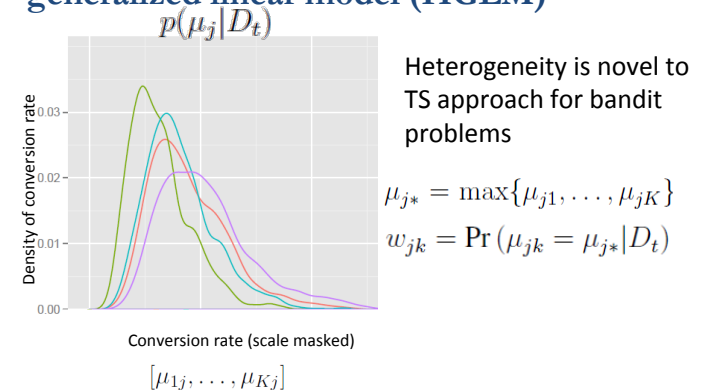
## Managing the multi-armed bandit



## Managing the multi-armed bandit

Online Simulator for Bayesian Bandits with Thompson Sampling  
<https://googledrive.com/host/0B2GQktu-wcTiWDB2R2t2a2tMUG8/>

## Thompson Sampling (TS) with hierarchical generalized linear model (HGLM)



## Thompson Sampling (TS) with hierarchical generalized linear model (HGLM)

$$\mu_{j*} = \max\{\mu_{j1}, \dots, \mu_{jK}\} \quad \text{website-specific "winner" among ads}$$

$$w_{jk} = \Pr(\mu_{jk} = \mu_{j*} | D_t) \quad \text{website-specific probabilities of each ad being the winner}$$










$$w_{jk} = \int_{\mu_j} \mathbf{1}\{\mu_{jk} = \mu_{j*} | \mu_j\} p(\mu_j | D_t) d\mu_j \quad \text{average over uncertainty in each ad's conversion rate}$$

$$w_{jk} \approx \hat{w}_{jk} = \frac{1}{G} \sum_{g=1}^G \mathbf{1}\{\mu_{jk}^{(g)} = \mu_{j*}^{(g)} | \mu_j^{(g)}\}$$

$$(m_{1,j,t+1}, \dots, m_{K,j,t+1}) = (\hat{w}_{j1t}, \dots, \hat{w}_{jKt}) M_{j,t+1}$$

"match" proportional allocations to probabilities

## Field experiment: summary and scope

- **700** million impressions
  - **80** media placements (websites)
  - **15** publishers
- 




  




- Conversion rates were within industry standards (between **1 and 10 out of 1 million** impressions)
  - For each website, **12** ads are described by attributes (**4** ad concepts x **3** ad sizes)

## Field experiment: ad attribute structure

- For each website:  
**4** ad concepts and **3** ad sizes

Conversion rate indexed as percent of average

	Ad A	Ad B	Ad C	Ad D
<b>Tall</b> 160x600	117	88	99	176
<b>Square</b> 300x250	107	72	151	114
<b>Wide</b> 728x90	115	92	66	72
<b>All Sizes</b>	112	80	100	105

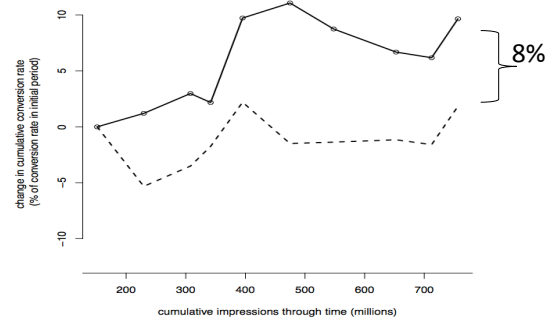
Best concept can differ by size

## How was the field experiment implemented?

- Timing
  - Update every 6 days for 61 days in 2012
- Allocation probabilities are "rotation weights"
  - Receive data and upload weights directly to Google DoubleClick DART (Dynamic Advertising Reporting and Targeting)



## Live field experiment: TS *versus* Balanced



- Test: Adaptive experiment (TS)
- Control: Static **balanced** experiment

ING DIRECT

M MICHIGAN  
ROSS SCHOOL OF BUSINESS

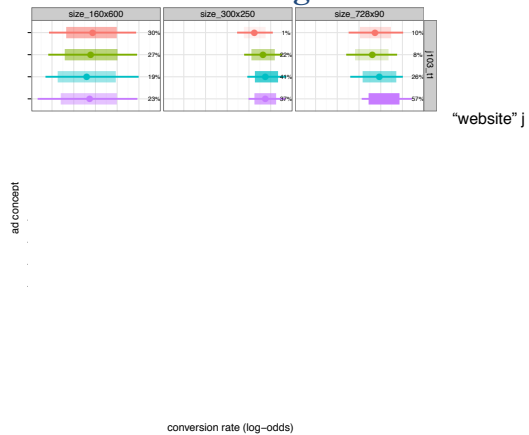
45

## How does TS learn for heterogeneous allocation?

\*4 concepts within 3 sizes, 2 "websites"

\*Boxplots show posterior of conversion rates. (right is better)

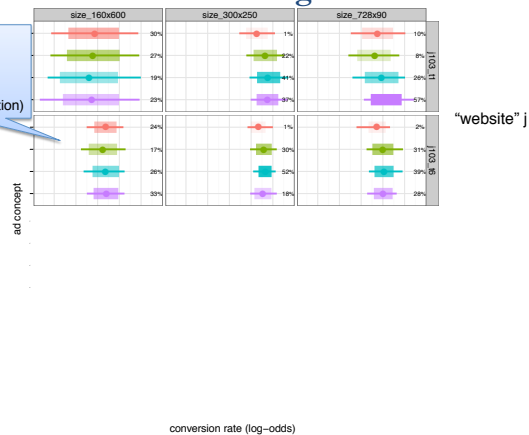
\*Darker coloring shows more allocation of impressions.



M MICHIGAN  
ROSS SCHOOL OF BUSINESS

## How does TS learn for heterogeneous allocation?

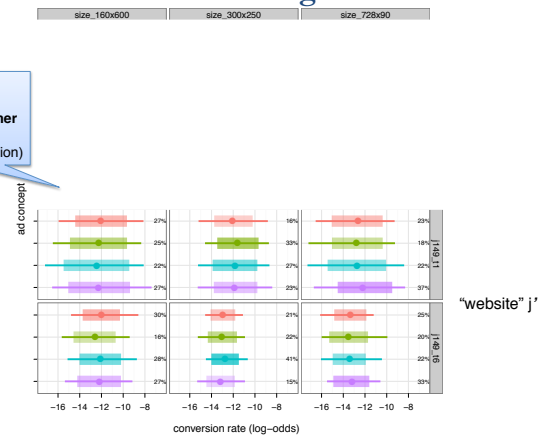
Learning each ad conversion rate over time for a website (uncertainty reduction)



M MICHIGAN  
ROSS SCHOOL OF BUSINESS

## How does TS learn for heterogeneous allocation?

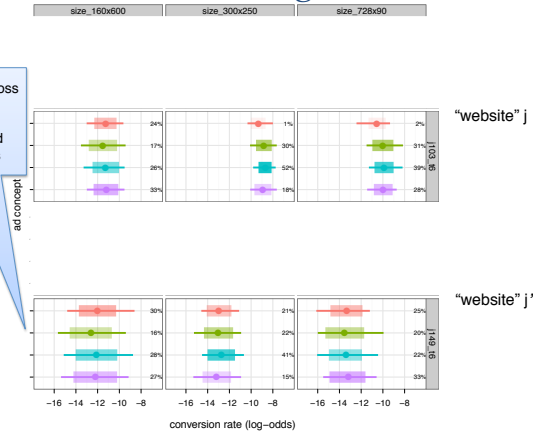
Learning each ad conversion rate over time for **another** website (uncertainty reduction)



M MICHIGAN  
ROSS SCHOOL OF BUSINESS

## How does TS learn for heterogeneous allocation?

Heterogeneity across websites: the same ads have different conversion rate and different allocations



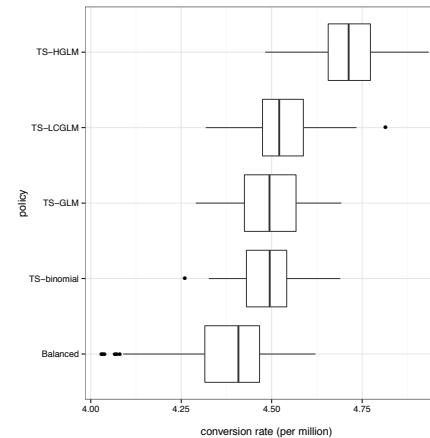
## Counter-factual simulations based on field experiment

## Counter-factual simulations based on field experiment

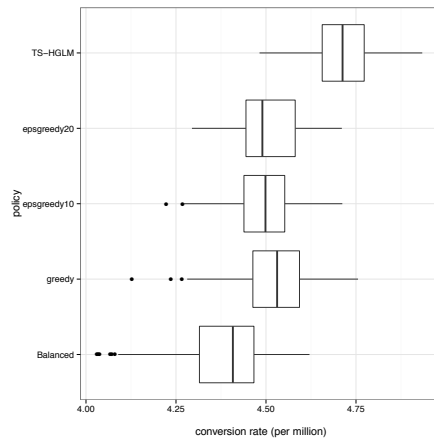
Run any policy to sequentially reallocate impressions given budget and timing constraints.

1. Compare TS-based policies with different models  
-attribute-based: logit (GLM), latent class (LCGLM)
2. Compare greedy and epsilon-greedy policies  
-Play the winner, always exploit (greedy)  
-Randomize: Explore ( $\epsilon\%$ ) and exploit ( $1-\epsilon\%$ ) (epsilon-greedy)
3. Compare intuitive test-rollout policies  
-Balanced experiment, then play the winner (explore, then exploit)

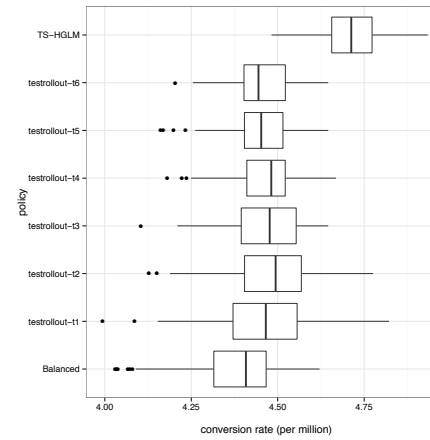
## TS-based policies with different models



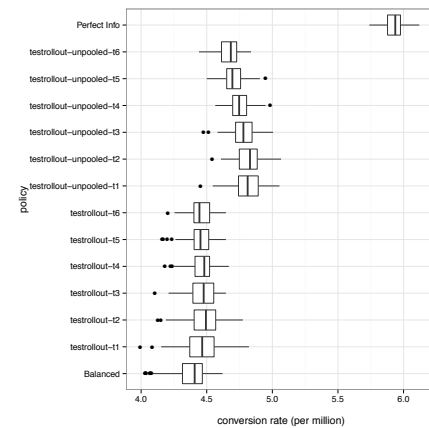
## Greedy and epsilon-greedy policies



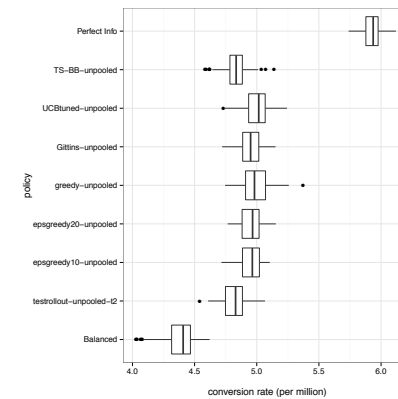
## Test-rollout policies



## Model heterogeneity vs. Allocation rule



## Heterogeneity matters most unpooled and partially pooled beat pooled



### Summary of counterfactuals

- Accounting for heterogeneity across context improves performance
- Hierarchical model (partial pooling) with bandit outperforms standard (pooled) models with bandit
- Even more flexible forms of heterogeneity can beat hierarchical bandit
- The modeling story matters more than particular bandit allocation rule...  
... so get the story right!



57

### Summary

- Large scale real-time adaptive field experiment and simulations show benefits of Thompson Sampling with appropriate model.
- Firms that experiment adaptively and systematically for continuous improvement should be “earning while learning.”



58

### Summary

- Large scale real-time adaptive field experiment and simulations show benefits of Thompson Sampling with appropriate model.
- The proposed hierarchical bandit beats standard models used with bandit algorithms.
- Firms that experiment adaptively and systematically for continuous improvement should be “earning while learning.”



59

### Future directions

- Dynamic (robust) firm pricing over many periods can improve when combining bandit methods with econ theory
- Market research survey techniques that adapt to respondents answers can more efficiently utilize its sample with bandit learning
- Consider lifetime value when exploring and exploiting new sources of customer acquisition



60

## Takeaways

- Active learning and adaptive sampling
- Learn vs. Earn tradeoff
- Many algorithms (Gittins, UCB, Thompson, egreedy) from in many disciplines
- Real-world problems motivate complicated versions with “bells and whistles”
- Always be earning while learning



**THANK YOU!**

[ericmsch@umich.edu](mailto:ericmsch@umich.edu)  
[ericmichaelschwartz.com](http://ericmichaelschwartz.com)

## Additional Details on Multi-Armed Bandit Optimization Methods



## Dynamic Robust Pricing with Multi-Armed Bandits

- How should we set price with limited information about demand to maximize profit over time?

## Dynamic Robust Pricing Overview

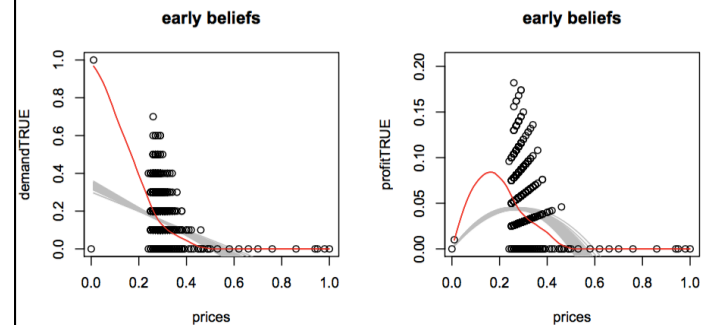
Research objective: provide a scalable method for price experimentation with computer science algorithms and economic theory

- Pricing with incomplete information
- Learning demand and price experimentation
- Prices are bandit arms
- Dynamic pricing
- Robust dynamic pricing
- Combining machine learning and pricing

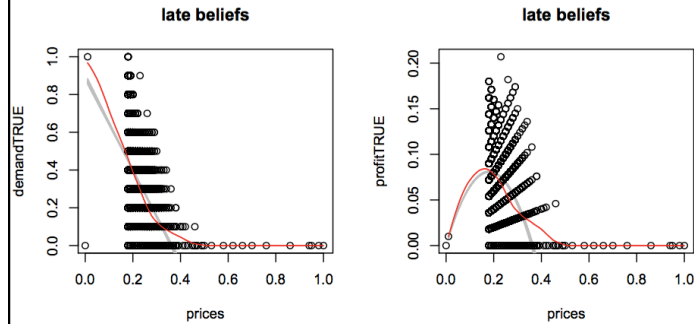
## Pricing with incomplete information

- Consider a firm objective of pricing a new product
  - Firm must set a price
  - Will assume consumers have stable preferences
- Limited data available
- How can a manager set a price?
  - Experiment!

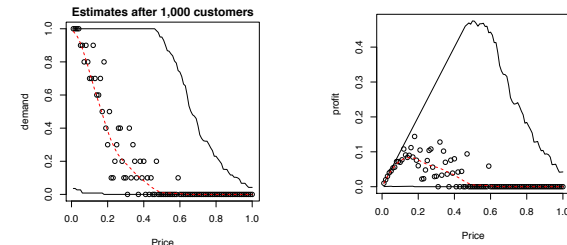
## Learning: assuming linear demand



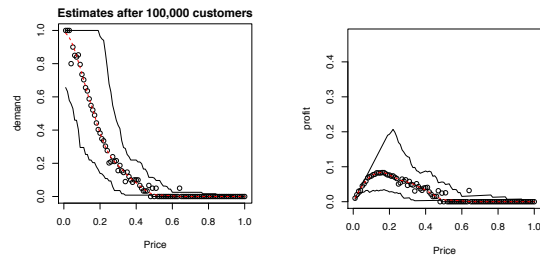
## Learning: assuming linear demand



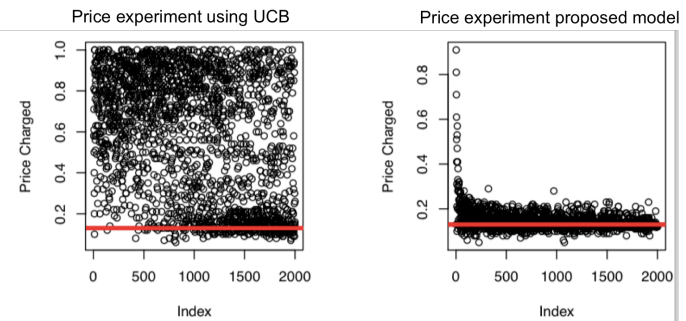
## Learning: no parametric assumption



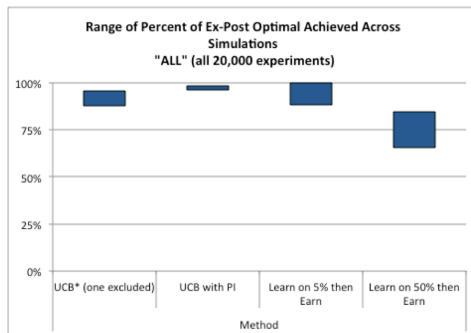
## Learning: no parametric assumption



## Converging to best price faster



## Proposed bandit pricing best controls risk



## Future directions

- Dynamic (robust) firm pricing over many periods can improve when combining bandit methods with econ theory
- Market research survey techniques that adapt to respondents answers can more efficiently utilize its sample with bandit learning
- Consider lifetime value when exploring and exploiting new sources of customer acquisition

## How does TS 'optimally' solve the problem?

For independent actions, no heterogeneity, one-at-a-time decisions ...

- Gittins index is the optimal certainty equivalent of an uncertain arm (1979)

But... that's not our problem!

- Explore/Exploit is balanced by sampling from full distribution of beliefs
- TS is asymptotically optimal in maximizing cumulative reward (i.e., minimax regret shrinks in log time).

Agrawal and Goyal 2012; Berry 2004; Chapelle and Li 2012; Granmo 2010; Kaufmann et al. 2012; May et al 2011; Russo and Van Roy 2014; Scott 2010; Thompson 1933

## Thompson Sampling (TS) with hierarchical generalized linear model (HGLM)

$$\mu_{j*} = \max\{\mu_{j1}, \dots, \mu_{jK}\}$$

website-specific "winner" among ads

$$w_{jk} = \Pr(\mu_{jk} = \mu_{j*} | D_t)$$

website-specific probabilities of each ad being the winner

$$w_{jk} = \int \mathbf{1}\{\mu_{jk} = \mu_{j*} | \mu_j\} p(\mu_j | D_t) d\mu_j$$

average over uncertainty in each ad's conversion rate

$$w_{jk} \approx \hat{w}_{jk} = \frac{1}{G} \sum_{g=1}^G \mathbf{1}\{\mu_{jk}^{(g)} = \mu_{j*}^{(g)} | \mu_j^{(g)}\}$$

$$(m_{1,j,t+1}, \dots, m_{K,j,t+1}) = (\hat{w}_{j1t}, \dots, \hat{w}_{jKt}) M_{j,t+1}$$

"match" proportional allocations to probabilities