

Thompson Sampling Tutorial

“The rediscovery of a swiss army knife”



Hanan Shteingart, PhD
Playtika



Alternative Name for the Lecture ☺

“Squeezing a 96-page review into a 25 min talk”

Based on Russo, D. J., Van Roy, B., Kazerouni, A., Osband, I., & Wen, Z. (2018). **A tutorial on Thompson sampling.** *Foundations and Trends® in Machine Learning*, 11(1), 1-96.

A Tutorial on Thompson Sampling

Daniel J. Russo¹, Benjamin Van Roy², Abbas Kazerouni², Ian Osband³ and Zheng Wen⁴

¹*Columbia University*

²*Stanford University*

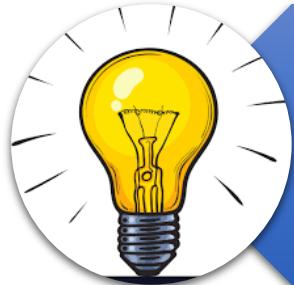
³*Google DeepMind*

⁴*Adobe Research*

https://web.stanford.edu/~bvr/pubs/TS_Tutorial.pdf

https://github.com/iosband/ts_tutorial

The Rediscovery of Thompson Sampling



TS was first proposed in 1933 for allocating experimental effort arising in clinical trials



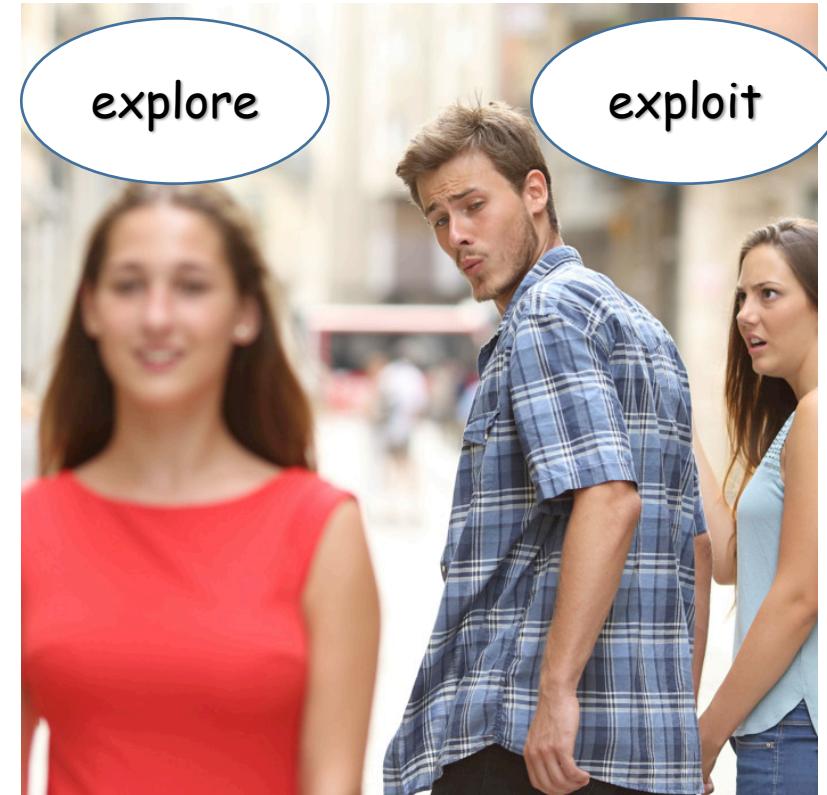
The algorithm was largely ignored in the academic literature until recently (although it was independently rediscovered several times in the interim)



After 8 decades... surge of interest. Especially because of (Chapelle and Li, 2011; Scott, 2010) which displayed the algorithm's strong empirical performance

Why Thompson Sampling?

- Algorithm for **online decision** problems
- **Actions** are taken **sequentially** in a manner that must **balance between exploiting and exploring**
- Addresses a **broad range** of problems in a computationally efficient manner and is therefore enjoying **wide** use
 - Bernoulli bandit problems, shortest path problems, product recommendation, assortment, active learning with neural networks, and reinforcement learning in Markov decision processes



Sorry if this image offends someone, it's just a joke...

Applications & Companies

That's why I was interested in it...



Revenue management
(Ferreira et al., 2015)

Marketing
(Schwartz et al., 2017)

web site optimization
(Hill et al., 2017)

Monte Carlo tree search
(Bai et al., 2013)

A/B testing
(Graepel et al., 2010)

Internet advertising
(Graepel et al., 2010;
Agarwal, 2013;
Agarwal et al., 2014)

Recommendation systems
(Kawale et al., 2015)

hyperparameter tuning
(Kandasamy et al., 2018)

arcade games
(Osband et al., 2016a)

Adobe, Amazon
(Hill et al., 2017)

Facebook, Google
(Scott, 2010; Scott,
2015)

LinkedIn
(Agarwal, 2013;
Agarwal et al., 2014)

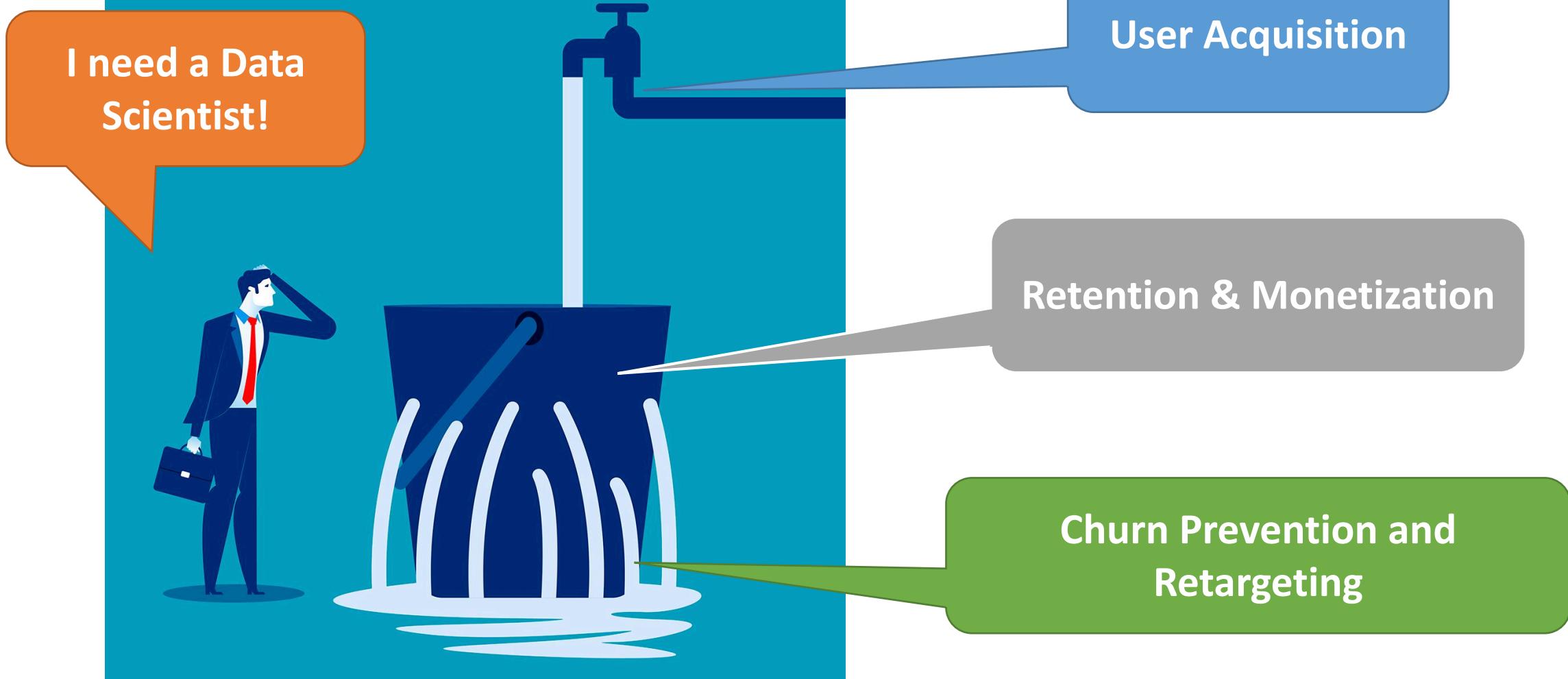
Microsoft
(Graepel et al., 2010)

Netflix

Twitter

Why Do We Need Marketing in Playtika?

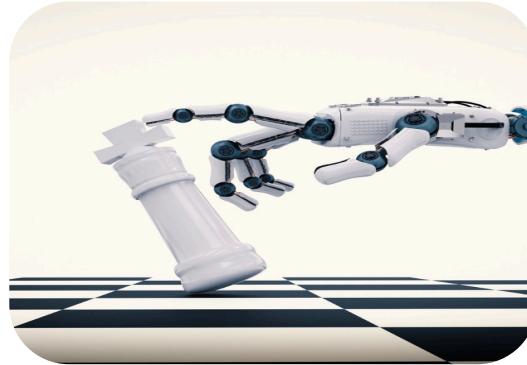
- The leaky bucket model for B2C companies



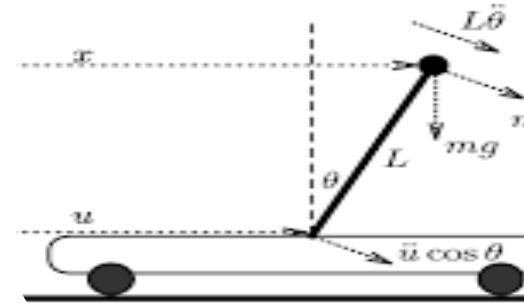
Fields of DS in Playtika (for example)



Player prediction



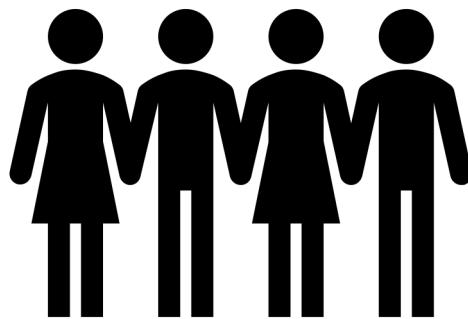
Artificial Player



Optimization &
Control, e.g. User
Acquisition



Recommendatio
n &
Personalization



Several Teams
 $\times 10$ of data scientists
 $\times 10$ of data engineers

6 TB/day
30M unique users /month

Multi Arm Bandit (MAB) Problem

- Imagine a gambler (octopus?!?!) going into a casino with different slot machine.
- Tradeoff: Pull the lever with the best historical payoffs or explore other levers to gain knowledge?
- A naïve approach (sample then decide) is wasteful.



The Bernoulli Case



Model

- K actions, reward is success or a failure.
- The success probabilities $P(1, \dots, K)$ are unknown to the agent, but are fixed over time.
- The objective, is to maximize the cumulative number of successes over T periods, where T is relatively large compared to the number of arms K .

Real World Application: Publisher side

- The “arms” in this problem might represent different banner ads
- A success is associated either with a click on the ad, or with a conversion
- The parameters are click through or conversion rate

Small-scale experiments are now a core tool at most leading Internet companies.

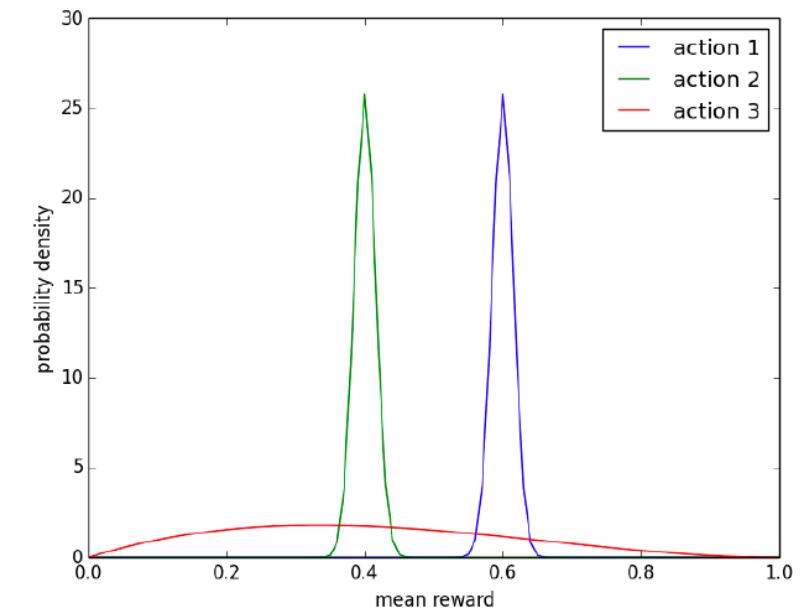
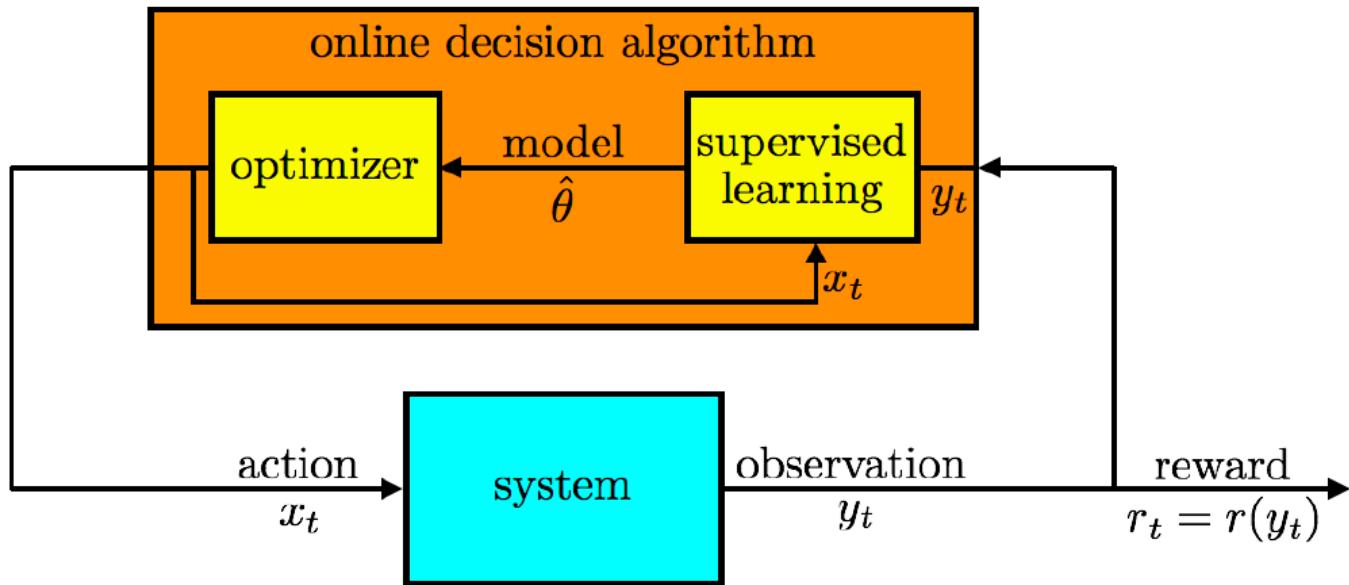
<http://ad-tech.com/>

How Is this different from Supervised Learning?

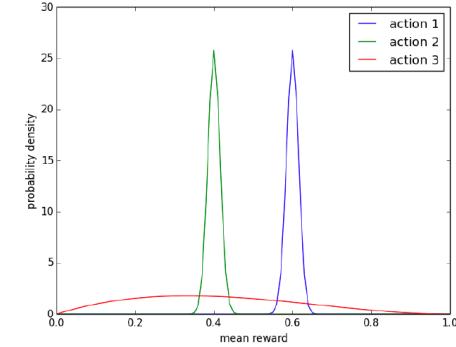
- Supervised machine learning learn **passively** from **historical** data
- Yet in MAB agents **drive the generation of their own training data** (e.g. through interacting with users).
- As a result, there is enormous potential benefit in the design of algorithms that **not only learn from past data**, but also **explore systematically** to generate **useful data that improves future performance**

Why does Greedy Approach Fail?

- “A shortcoming of the greedy approach, which can severely curtail performance, is that it does not actively explore”



Why does Dithering Fail?



- randomly perturbing actions that would be selected by a greedy algorithm.
- ϵ -greedy exploration - greedy action with probability $1 - \epsilon$ and otherwise selects an action uniformly at random
- wastes resources by failing to “write off” actions regardless of how unlikely they are to be optimal
- TS provides an alternative to dithering that more intelligently allocates exploration effort...

Thompson Sampling Main Idea

- Probability Matching: "Optimism in face of uncertainty"
- Probability to select an action is equal to the probability this action is optimal
- By sampling actions according to the posterior probability that they are optimal, the algorithm continues to sample all actions that could plausibly be optimal, while shifting sampling away from those that are unlikely to be optimal.
- **Roughly speaking, the algorithm tries all promising actions while gradually discarding those that are believed to underperform.**

$$\Pr(a|x) = \Pr(\operatorname{argmax}_{a'} r(x, a') = a)$$

Example: Revisit the Bernoulli MAB

- K actions \sim Bernoulli with fixed unknown probability $\theta = (\theta_1, \dots, \theta_K)$

- Assume prior is beta $p(\theta_k) = \frac{\Gamma(\alpha_k + \beta_k)}{\Gamma(\alpha_k)\Gamma(\beta_k)} \theta_k^{\alpha_k-1} (1 - \theta_k)^{\beta_k-1}$

- posterior distribution is also beta (conjugate prior) with parameters that can be updated according to a simple rule

$$(\alpha_k, \beta_k) \leftarrow \begin{cases} (\alpha_k, \beta_k) & \text{if } x_t \neq k \\ (\alpha_k, \beta_k) + (r_t, 1 - r_t) & \text{if } x_t = k. \end{cases} \quad \hat{\theta}_k = \alpha_k / (\alpha_k + \beta_k)$$

Compare with Bernoulli Greedy

Algorithm 1 BernGreedy(K, α, β)

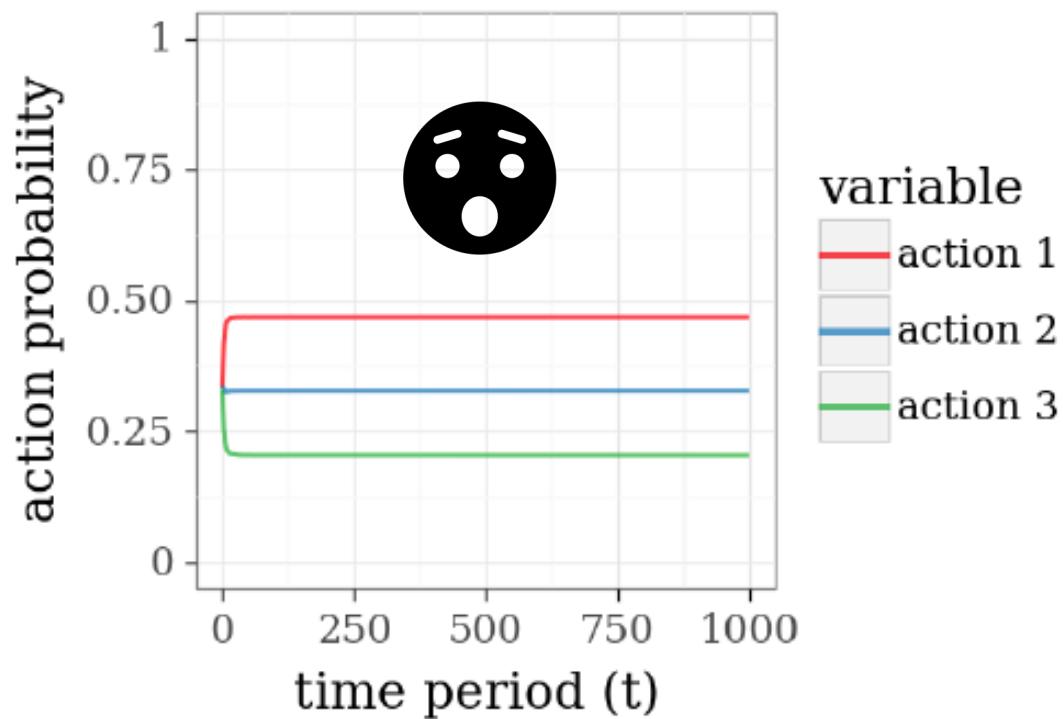
```
1: for  $t = 1, 2, \dots$  do
2:   #estimate model: expectation
3:   for  $k = 1, \dots, K$  do
4:      $\hat{\theta}_k \leftarrow \alpha_k / (\alpha_k + \beta_k)$ 
5:   end for
6:
7:   #select and apply action:
8:    $x_t \leftarrow \operatorname{argmax}_k \hat{\theta}_k$ 
9:   Apply  $x_t$  and observe  $r_t$ 
10:
11:  #update distribution:
12:   $(\alpha_{x_t}, \beta_{x_t}) \leftarrow (\alpha_{x_t} + r_t, \beta_{x_t} + 1 - r_t)$ 
13: end for
```

Algorithm 2 BernTS(K, α, β)

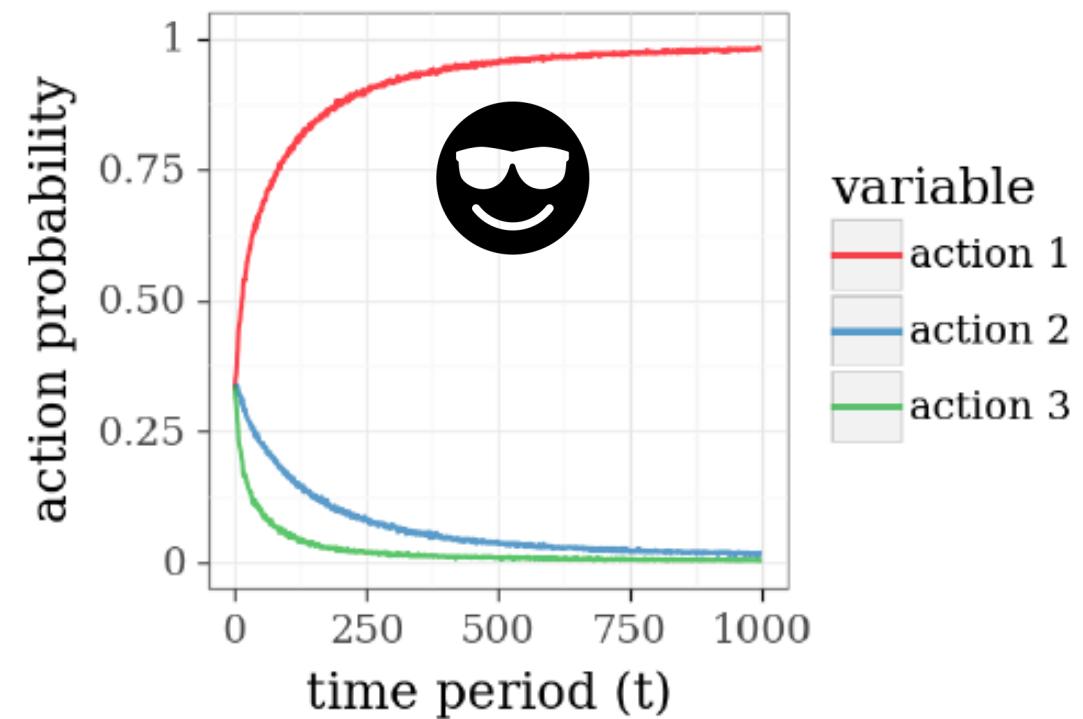
```
1: for  $t = 1, 2, \dots$  do
2:   #sample model: P(best)
3:   for  $k = 1, \dots, K$  do
4:     Sample  $\hat{\theta}_k \sim \text{beta}(\alpha_k, \beta_k)$ 
5:   end for
6:
7:   #select and apply action:
8:    $x_t \leftarrow \operatorname{argmax}_k \hat{\theta}_k$ 
9:   Apply  $x_t$  and observe  $r_t$ 
10:
11:  #update distribution:
12:   $(\alpha_{x_t}, \beta_{x_t}) \leftarrow (\alpha_{x_t} + r_t, \beta_{x_t} + 1 - r_t)$ 
13: end for
```

Results of Simulation

$\theta_1 = 0.9$, $\theta_2 = 0.8$, and $\theta_3 = 0.7$



(a) greedy algorithm



(b) Thompson sampling

Results based on 10,000 independent simulations of each algorithm each over 1,000 time periods

Beyond Bernoulli - Generalized TS

Algorithm 4 Thompson(\mathcal{X}, p, q, r)

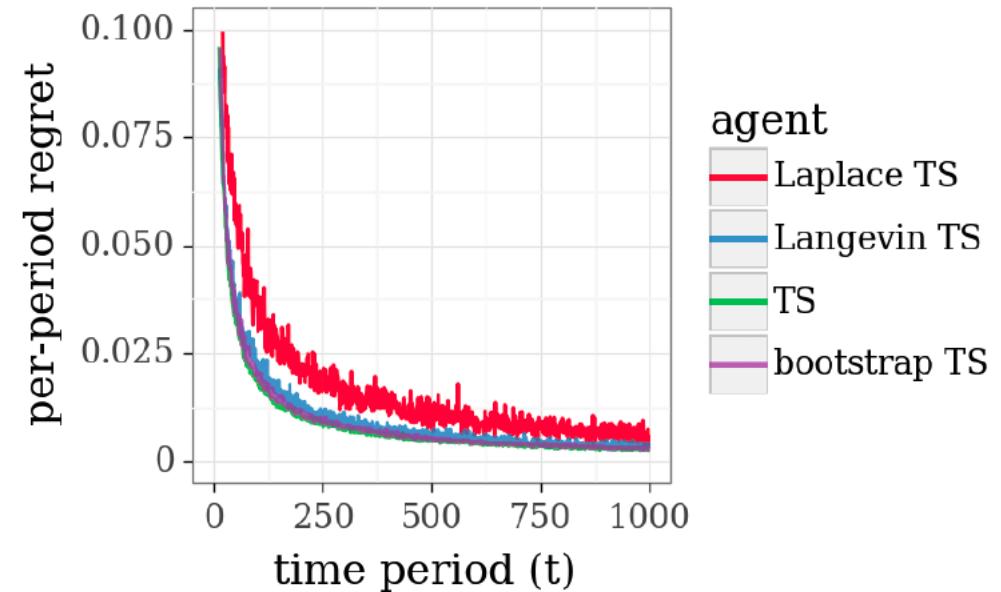
```
1: for  $t = 1, 2, \dots$  do
2:   #sample model:
3:   Sample  $\hat{\theta} \sim p$ 
4:
5:   #select and apply action:
6:    $x_t \leftarrow \operatorname{argmax}_{x \in \mathcal{X}} \mathbb{E}_{q_{\hat{\theta}}} [r(y_t) | x_t = x]$ 
7:   Apply  $x_t$  and observe  $y_t$ 
8:
9:   #update distribution:
10:   $p \leftarrow \mathbb{P}_{p, q}(\theta \in \cdot | x_t, y_t)$ 
11: end for
```

Approximation

- Many practical contexts call for more complex models for which exact Bayesian inference is computationally intractable.
- Fortunately, there are reasonably **efficient** and **accurate** methods that can be used to **approximately sample from posterior distributions**:
 - Gibbs sampling
 - Langevin Monte Carlo
 - sampling from a Laplace approximation
 - **Bootstrap**

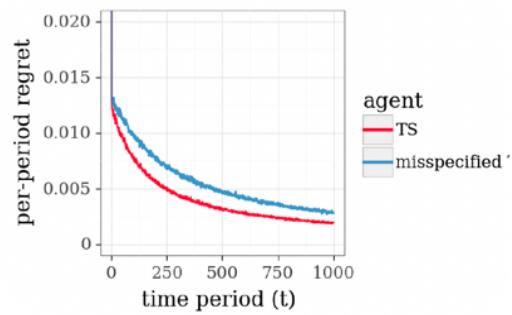
Bootstrap TS

- The method generates a **hypothetical history** which is **made up** of $t - 1$ action-observation pairs, each sampled uniformly with replacement from $H_{t-1} \rightarrow \hat{H}_{t-1} = \{(\hat{x}_1, \hat{y}_1), \dots, (\hat{x}_{t-1}, \hat{y}_{t-1})\}$
- We then maximize the likelihood of estimator under the hypothetical history
- there is a lack of theoretical justification for bootstrap approaches or even understanding of whether there are nontrivial problem classes for which they are guaranteed to perform well. Yet see [1]



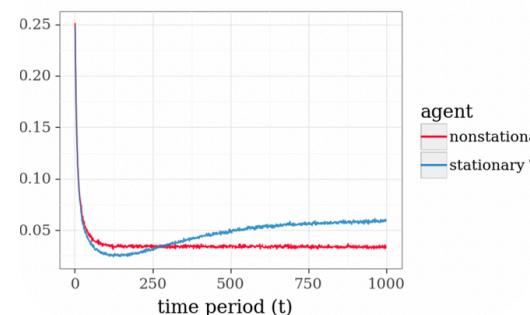
(a) Bernoulli bandit

Practical Considerations



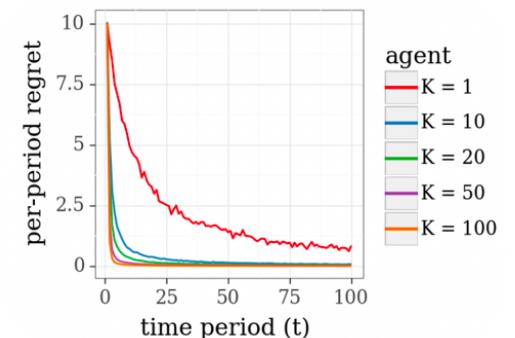
Priors – coherent vs.
mis-specified prior

$$p_{\theta}(\cdot | x_t, z_t)$$



Context – additional
information

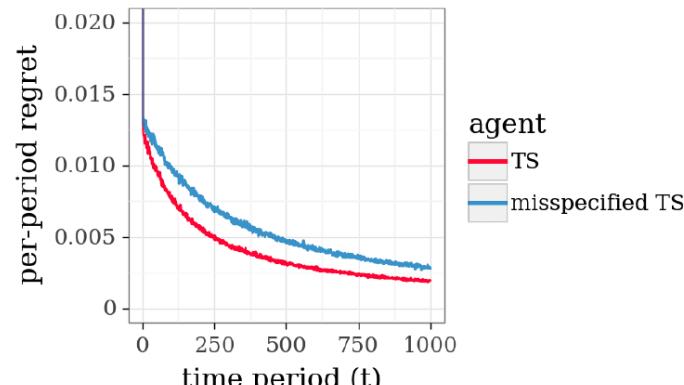
Non stationarity –
environment changes
over time



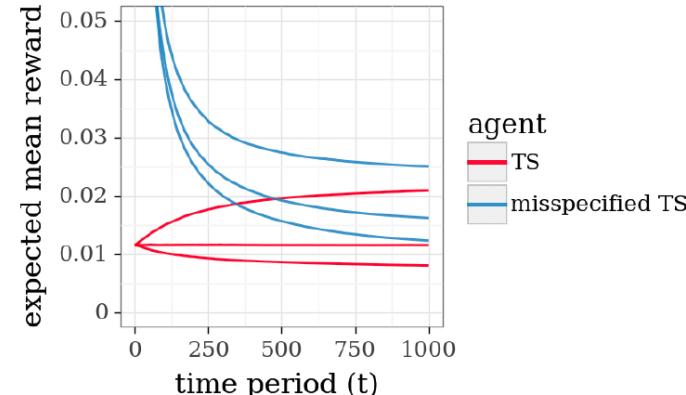
Concurrency – ability
to select multiple
actions per step

Practical Considerations (1. using Priors)

- In particular, this prior represents no understanding of the context, ignoring any useful knowledge from past experience. Taking knowledge into account reduces what must be learned and therefore reduces the time it takes for or TS to identify the most effective ads.



(a) regret



(b) expected mean rewards

Figure 6.2: Comparison of TS for the Bernoulli bandit problem with coherent versus misspecified priors.

2. contextual online decision problems

- In such problems, the response y_t to action x_t
- also depends on an independent random variable z_t that the agent
- observes prior to making her decision. In such a setting, the conditional
- distribution of y_t takes the form $p(\cdot | x_t, z_t)$.

3. Non-Stationarity

- In many practical applications,
- the agent faces a nonstationary system, which is more appropriately
- modeled by time-varying parameters $1, 2, \dots$, such that the response
- y_t to action x_t is generated according to $p_t(\cdot|x_t)$.
- With minor modification, TS remains an effective
- approach so long as model parameters change little over durations that
- are sufficient to identify effective actions

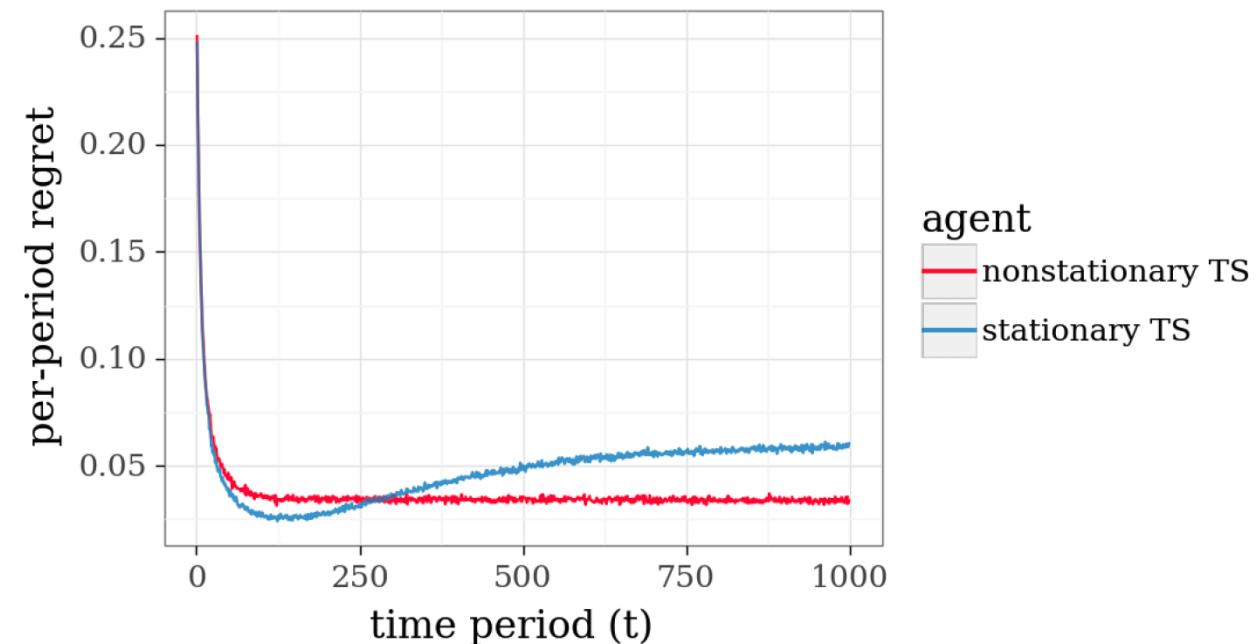
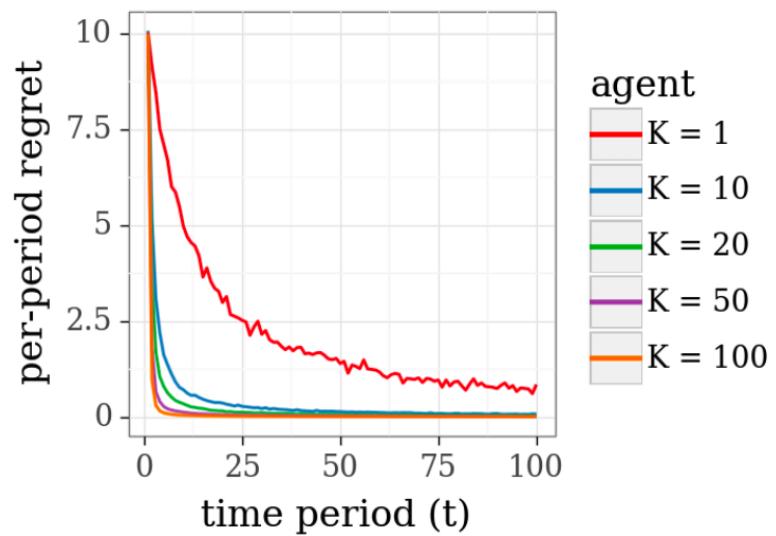


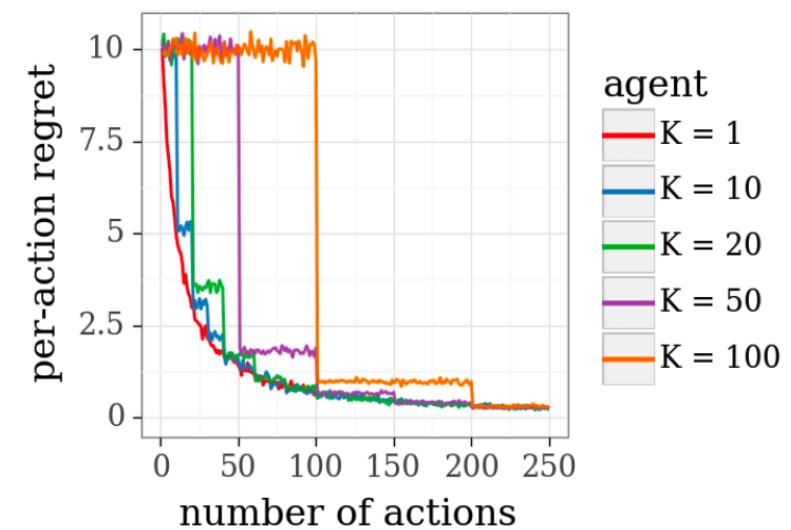
Figure 6.3: Comparison of TS versus nonstationary TS with a nonstationary Bernoulli bandit problem.

4. Concurrent

- concurrency plays an important role in web services, where at any time, a system may experiment by providing different versions of a service to different users



(a) per-action regret over time



(b) per-action regret over actions

Figure 6.4: Performance of concurrent Thompson sampling.

When does it fails?

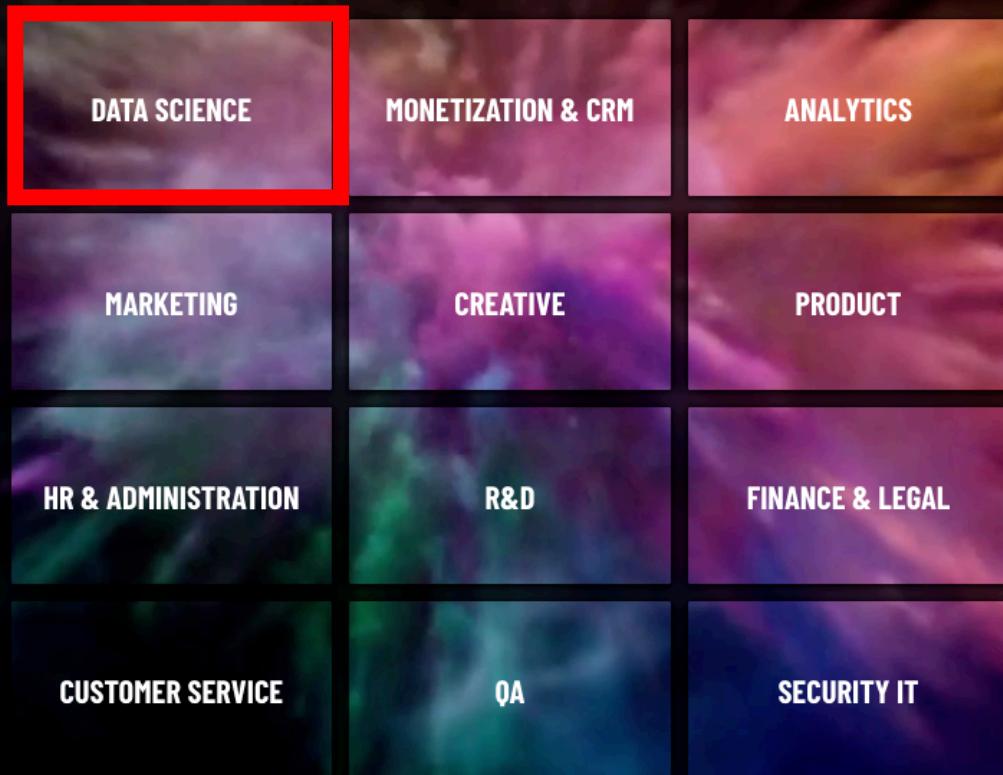
- Problems that do not Require Exploration
- Problems that do not Require Exploitation
- Time Sensitivity
- Problems Requiring Careful Assessment of Information Gain

MAKE YOUR NEXT MOVE

Departments Our offices Life at Playtika

We're looking for product ninjas, marketing wizards, creative rockstars, data masters, business aces, legal gladiators.

If you have passion to build and ambition to grow we want to meet you.



www.playtika.com/careers

