

# Multi Armed Bandits and Applications

Dr. M.C.Kaptein  
Assistant Professor AI

4 December 2015

Background

The multi armed bandit problem

Thompson sampling & optimal design

The contextual Multi-armed Bandit problem

Bootstrap Thompson Sampling

Streaming Bandit: software

Applications of the software:

Future

# Section 1

## Background

# Background

- ▶ MSc. Economic Psychology, Tilburg University
- ▶ PdEng. User System Interaction, University of Eindhoven
- ▶ Ph.D. Industrial Design, University of Eindhoven & Stanford University
- ▶ Post Doc. Marketing, Aalto School of Economics, Helsinki
- ▶ Assistant Professor Statistics, Tilburg University, Tilburg
- ▶ Founder PersuasionAPI.

## Section 2

### The multi armed bandit problem

# Slot machines



# Formal presentation

- ▶ For  $t = 1, \dots, t = T$
- ▶ Select an action  $a_t$  out of  $\mathcal{A}_t$ . Often actions  $k = 1, \dots, k = K$ .
- ▶ Observe reward  $r_t$  (generated by some unknown distribution  $F_k(r|\theta_k)$ )
- ▶ Play according to some policy  $\Pi : \{a_1, \dots, a_{t-1}, r_1, \dots, r_{t-1}\} \mapsto a_t$

## Aim of a “good” policy

- ▶ Well, get as much reward as possible!
- ▶ Thus, maximize  $\sum_{t=1}^T r_t$
- ▶ Or, minimize regret:  $\sum_{t=1}^T (\Pi^*(t) - \Pi(t))$



# Exploration-Exploitation tradeoff

- ▶ Suppose observations  $X_k \sim \text{Bern}(p_k)$
- ▶ Explore:  $p_1 > p_2$ ? Play alternating arms to learn.
- ▶ Exploit: Play arm 1.

Very general trade-off: Exploring the outcomes of uncertain actions, versus choosing actions that one believes to be good.

# Omnipresence of the tradeoff

Exploration vs. exploitation found in many places:

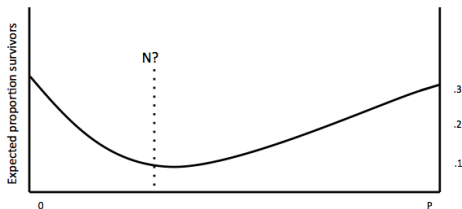
- ▶ Clinical trial: which medicine to subscribe?
- ▶ Online content selection: Which ad, news article, or product to show?
- ▶ Job choices: Try something new, vs. stick to what you have?
- ▶ Food choices: Try a new dish, stick to one you like
- ▶ Etc. etc.

Also known as *Earning vs. Learning*.

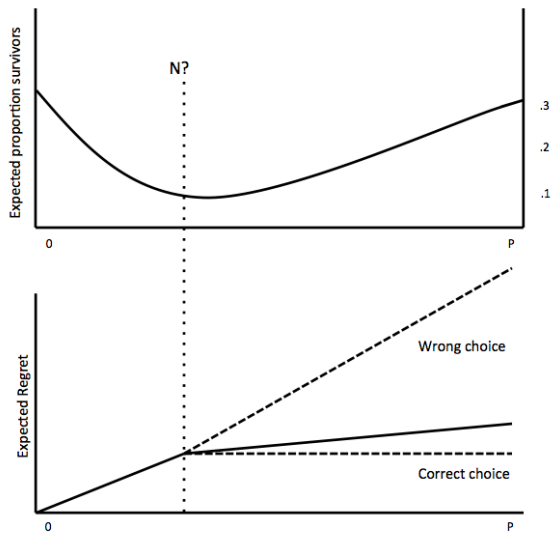
# The randomized clinical trial

What is the regret of a simple RCT choosing between two medications?

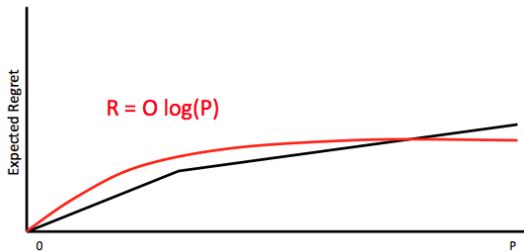
- ▶ Simple choices:  
 $X_1 \sim \text{Bern}(p_1 = .1)$ ,  
 $X_2 \sim \text{Bern}(p_2 = .5)$
- ▶ Obviously:  $p_1, p_2$  unknown at start
- ▶  $N$  patients in trial,  $P$  total patients in population
- ▶ Always  
 $Pr(\text{Wrong}) > 0$



# The randomized clinical trial: Regret



# Optimal policies



## Section 3

### Thompson sampling & optimal design

# Thompson sampling

Compute or sample from  $\Pr(\theta|\mathcal{D})$ . We can then select an action according to its probability of being optimal:

$$\int \mathbf{1} \left[ \mathbb{E}(r|a, \theta) = \max_{a'} \mathbb{E}(r|a', \theta) \right] \Pr(\theta|\mathcal{D}) d\theta \quad (1)$$

where  $\mathbf{1}$  is the indicator function.

# Thompson sampling Bernoulli Bandit

Thompson sampling in practice for the k-armed Bernoulli Bandit

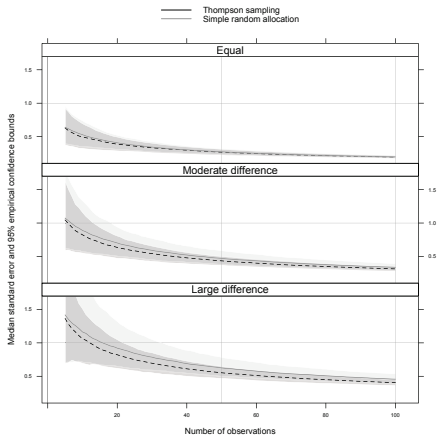
- ▶ Suppose again  $X_k \sim \text{Bern}(p_k)$
- ▶ Use (independent)  $\text{Beta}(\alpha_k = 1, \beta_k = 1)$  priors
- ▶ Generate random draw  $d_k$  from each  $k = 1, \dots, k = K$   $\text{Beta}()$  distributions
- ▶ Select arm  $k' = \max_k d_k$
- ▶ Update posterior  $\text{Beta}(\alpha_{k'} + r_t, \beta_{k'} + 1 - r_t)$

Thompson sampling is an asymptotically optimal strategy.



# Experimental design as exploration vs. exploitation <sup>1</sup>

- ▶ Thompson sampling for optimal design
- ▶ Sample for most “informative” datapoints
- ▶ Assume heterogeneity of variances
- ▶ Select treatments based on posterior variance estimates



<sup>1</sup>Kaptein, M.C. (2014).

## Section 4

### The contextual Multi-armed Bandit problem

## Extension: contexts

- ▶ For  $t = 1, \dots, t = T$
- ▶ **Observe the world**,  $x_t \in \mathcal{X}_t$
- ▶ Select and action  $a_t$  out of  $\mathcal{A}_t$ . Often actions  $k = 1, \dots, k = K$ .
- ▶ Observe reward  $r_t$  (generated by some unknown distribution  $F_k(r|\theta_k)$ )
- ▶ Play according to some policy  $\Pi : \{x_1, \dots, x_{t-1}, a_1, \dots, a_{t-1}, r_1, \dots, r_{t-1}\} \mapsto a_t$

# Examples of contextual bandits

- ▶ Clinical trial: which medicine to subscribe to a **specific patient**?
- ▶ Online content selection: Which ad, news article, or product to show to a **user**?
- ▶ Job choices: You have information regarding the jobs
- ▶ Food choices: You know the ingredients of the dish
- ▶ Etc. etc.

Interesting model for (e.g.,) **treatment personalization**.

## Section 5

# Bootstrap Thompson Sampling

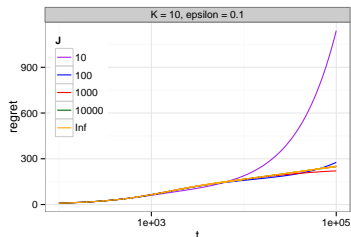
# Thompson sampling for contextual bandits

- ▶ Setup some model  $r = f(a, x; \theta)$
- ▶ Choose prior  $Pr(\theta)$
- ▶ Observe  $\mathcal{D} = (x_t, a_t, r_t)$
- ▶ Use Bayes rule and sample  $\theta_{t'}$  from  $Pr(\theta|\mathcal{D})$
- ▶ Select action  $a$  that maximizes  $f(a, x; \theta)$  given  $x_{t'}$  and  $\theta_{t'}$

This can be hard if  $Pr(\theta|\mathcal{D})$  is hard to sample from.

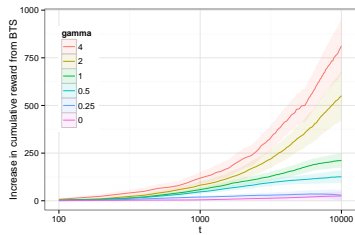
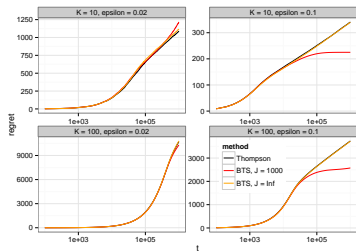
# Bootstrapped Bandit<sup>2</sup>

- ▶ Thompson sampling works well if posterior is known
- ▶ Not the case for complex models
- ▶ What about the (double or nothing) bootstrap distribution?
- ▶ For  $1, \dots, J$  online bootstrapped replicates



<sup>2</sup>Kaptein, & Eckles (2014).

# Bootstrap bandit continued ...



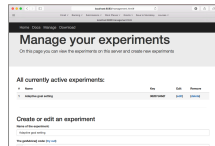
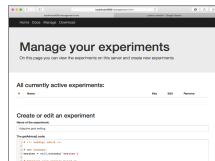
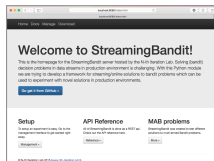


## Section 6

### Streaming Bandit: software

# Streaming Bandit <sup>3</sup>

- ▶ Back end solution for streaming bandits
- ▶ Setup a REST server to handle action selection
- ▶ Recently released first stable version



<sup>3</sup>Kaptein, M.C. & Kruijswijk, J. (2015).

## Design choice: learning vs. choosing

We identify two steps:

1. The *summary* step: In each summary step  $\theta_{t'-1}$  is updated by the new information  $\{x_{t'}, a_{t'}, r_{t'}\}$ . Thus,  $\theta_{t'} = g(\theta_{t'-1}, x_{t'}, a_{t'}, r_{t'})$  where  $g()$  is some update function.
2. The *decision* step: In the decision step, the model  $r = f(a, x_{t'}; \theta_{t'})$  is evaluated for the current context and the possible actions. Then, the recommended action at time  $t'$  is selected.

Implemented in `getAction()` and `setReward()` calls.

# Online learning

Forces an online learning approach.

- ▶ Summation over datapoints:

- ▶ Version 1:  $S_T = \sum_{t=1}^T x_t$

- ▶ Version 2:  $S_T = S_{T-1} + x_t$

Linear vs. Quadratic function of  $T$  to compute.<sup>4</sup>

---

<sup>4</sup>Ippel, L., Vermunt, J., & Kaptein, M.C. (2015)  
Streaming EM approximations. *Under submission*.

## Section 7

Applications of the software:

# Simple experiment using StreamingBandit

## Summarize:

```
import libs.base as base
prop = base.Proportion(self.get_theta(key="version",
    value=self.action["version"]))
prop.update(self.reward["click"])
self.set_theta(prop, key="version", value=self.action["version"])
```

## Decide:

```
import libs.base as base
propl = base.List(self.get_theta(key="version"),
    base.Proportion, ["A", "B"])
if propl.count() > 1000:
    self.action["version"] = propl.max()
else
    self.action["version"] = propl.random()
```

## Simple experiment using StreamingBandit 2

Decide:

```
import libs.thompson as thmp
prop1 = thmp.BBThompsonList(self.get_theta(key="version"),
    Proportion, ["A", "B"])
self.action["version"] = prop1.thompson()
```

# Streaming Bandit in practice: Lock in Feedback

Possible policy for the continuum bandit problem:

- ▶  $a \in \mathbb{R}$
- ▶ Project together with Prof. Dr. Davide Iannuzzi
- ▶ Oscillate  $a$  with a known frequency
- ▶ Amplify  $r$ , and integrate to obtain first derivative.<sup>5</sup>

---

<sup>5</sup>Kaptein, M.C. & Iannuzzi, D. (2015)



# Streaming Bandit in practice: Santander

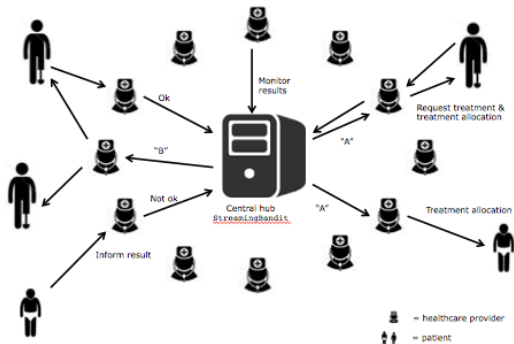
Software currently used for contextual bandit trials

- ▶ Observe features of a customer requesting a loan
- ▶ Select an interest rate
- ▶ Observe acceptance of loan ( $r = f(IR, y)$ )
- ▶ Objective: Choose interest such as to maximize profit

# Section 8

## Future

# Personalized feedback and treatment selection



# Questions?

Maurits Kaptein  
Archipelstraat 13  
6524LK, Nijmegen  
0621262211  
maurits@mauritskaptein.com