

Terminating Information Search

Stefan Appelhoff

January 31, 2017

Contents

List of Figures

1 Introduction

As agents in different environments, humans often find themselves in situations that require decisions based on limited knowledge about the available choice options. The topic of the present project is the process, how humans selectively sample choice options to arrive at a more refined knowledge about the properties of the options. Specifically, the interest lies in self-directed information search in sequential decision making problems. A short example illustrates the topic further:

Imagine that your laptop is broken and you want to acquire a new one. In the process, you may sequentially sample several available models and their respective properties such as prices, specifications, and reviews. To avoid indefinite browsing, you will need to terminate your search at some point and select a new laptop - most likely without having obtained full certainty about the best option.

For the project described here, the aforementioned example will be transformed into an experimental paradigm that allows for the examination of behavioral and neural aspects of choice in a controlled manner. Similar to the real life example, the experimental paradigm exhibits the following characteristics:

- There is an environment in which several options exist that will yield outcomes upon selection.
- There exists limited knowledge about the options and the frequencies of their associated outcomes.

- There will be the opportunity to obtain information about the options through a sampling of the options.
- The overall goal is to maximize the positive outcomes obtained from the options in the long run.

In the paradigm, participants will make sequential decisions about which option to sample. In these trial to trial decisions, the question of interest to this project lies in how participants set their termination criteria. I.e., when do participants decide to terminate the information search in one option and switch to search another option? When do participants decide to stop sampling altogether and decide for one final option to be optimal based on current knowledge? The remainder of this document will be a concise formulation of the research questions followed by an outline of the experimental paradigm. Finally, the implementation of the experimental paradigm on a computer program will be described.

2 Research Questions

1. How are sampling efforts sequentially allocated to the different options?
2. How do we arrive at a decision to terminate the information search?
3. What is the quantitative link between neural signals obtained by EEG with the behavioral data of participants performing the experimental paradigm.

3 The Experimental Paradigm

The proposed experiment will be based on a combination of two different paradigms that are well known in the literature.

1. The n-armed Bandit Paradigm
2. The Sampling Paradigm

In the following, these two paradigms will be described in detail. Figure ?? provides an overview of the paradigms. Afterwards, the combination of both paradigms for the present experiment will be outlined.

3.1 The n-armed Bandit Paradigm

The n-armed bandit paradigm was first formalized by Robbins in 1952 [robbins1952]. It has since been the subject of interest in a variety of fields, see [berry1985] for a review. Generally, participants are presented with n action that represent initially unknown outcome distributions with a set number of outcomes. Upon selection of an

action, an outcome is yielded from the distribution and presented to the participant as a feedback. With this gained knowledge, the participant may continue with selecting the same or a different action.

The bandit paradigm offers a concise framework to study the exploration exploitation tradeoff common to the field of reinforcement learning. Here, exploration represents the goal to gather information about actions and their associated outcomes, while exploitation represents the goal to make use of the acquired information. A player who only explores will never exploit upon the accumulated knowledge, while a player who only exploits will most likely select suboptimal actions based on the inadequate knowledge about all actions. Hertwig & Erev label the bandit paradigm as Partial-Feedback-Paradigm [hertwig2009].

3.2 The Sampling Paradigm

The sampling paradigm is similar to the bandit paradigm. The main difference is that exploration and exploitation are separated. This is achieved through permitting a participant to sample all actions without consequences for as long as the participant wishes. This stage represents purely the exploration, because the outcomes do not count towards a payoff. Once the participant decides to terminate the sampling process, there is one last opportunity to select any action. For this last choice, the outcome will affect the payoff and thus, this stage represents purely the exploitation. For more information, see [hertwig2011, erev2014, hertwig2009]

3.3 Combination of the Two Paradigms

In the present project, we will apply a paradigm that combines the aforementioned bandit and sampling paradigms into a single experimental flow. We will further introduce another factor, namely the mode of interaction of the participant with the paradigm. Each of the factors will have two levels. In brief, we are thus talking about a 2x2 factorial design of an experiment as visualized in figure ???. The two factors are:

- Paradigm condition
- Interaction condition

3.3.1 Paradigm condition

The paradigm condition is governed by the bandit or sampling paradigm as the two levels respectively. Thus in the bandit condition, participants will perform a bandit task, whereas in the sampling condition, participants will perform a sampling paradigm task. This will allow us to uncover potential differences in behavior when humans are faced with a task that either combines or separates exploration and exploitation into stages.

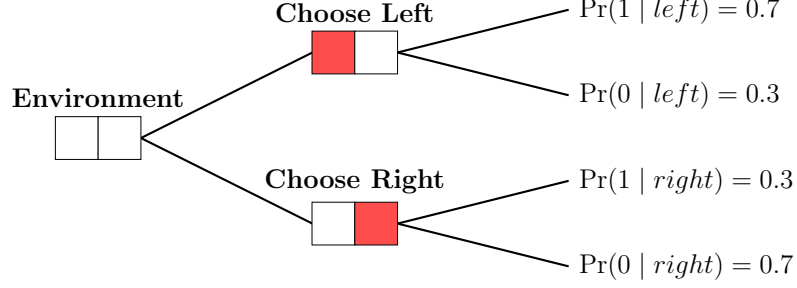
3.3.2 Interaction condition

The interaction condition is governed by active or passive interaction respectively. During the active condition, participants will actively select the actions and steer the flow of the paradigm. During the passive condition, participants will merely watch a replay of previous games. This allows us to make a distinction between (a) cognitive and motor involvement in active tasks and (b) neither cognitive nor motor involvement in the passive tasks. During the passive tasks there is no cognitive involvement, because participants are not solving or learning the underlying decision paradigm.

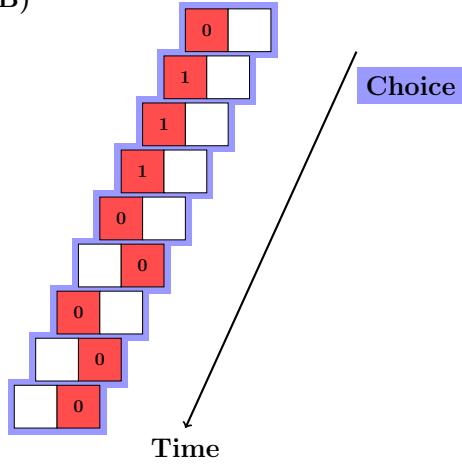
3.3.3 Orthogonal task

To make sure that participants pay attention even during the passive conditions, we introduce an orthogonal task to the 2x2 experimental design. Here, participants will have to react to a stimulus of a different color by a quick button press. If the reaction times are too slow, we can assume that the participants did not pay sufficient attention.

A)



B)



C)

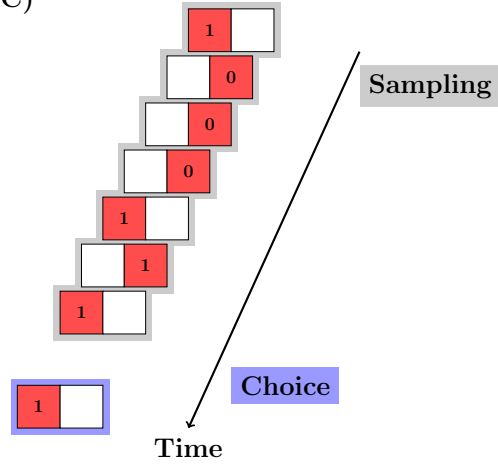


Figure 1: **A)** Depiction of a two-armed bandit. Note that the number of options in the environment can be extended arbitrarily to form an n -armed bandit. Each option contains its own probability mass function (pmf) of outcomes. These pmfs are unknown to a subject performing the bandit task. **B)** A typical run of trials in a bandit task. A subject sequentially chooses among options and is provided with feedback. Each trial represents a exploration-exploitation problem. **C)** A typical run of trials in the sampling paradigm. A subject is allowed to sample the available options without an impact of the outcomes on the final payoff. At some point, the subject can decide to terminate sampling and choose one of the options, which will impact the final payoff. In the sampling paradigm, exploration and exploitation are thus distinct.

Paradigm Condition			
Interaction Condition		Sampling Paradigm	Bandit Paradigm
	Active	Active Sampling	Active Bandit
	Passive	Passive Sampling	Passive Bandit

Figure 2: The 2x2 factorial design of the experiment.

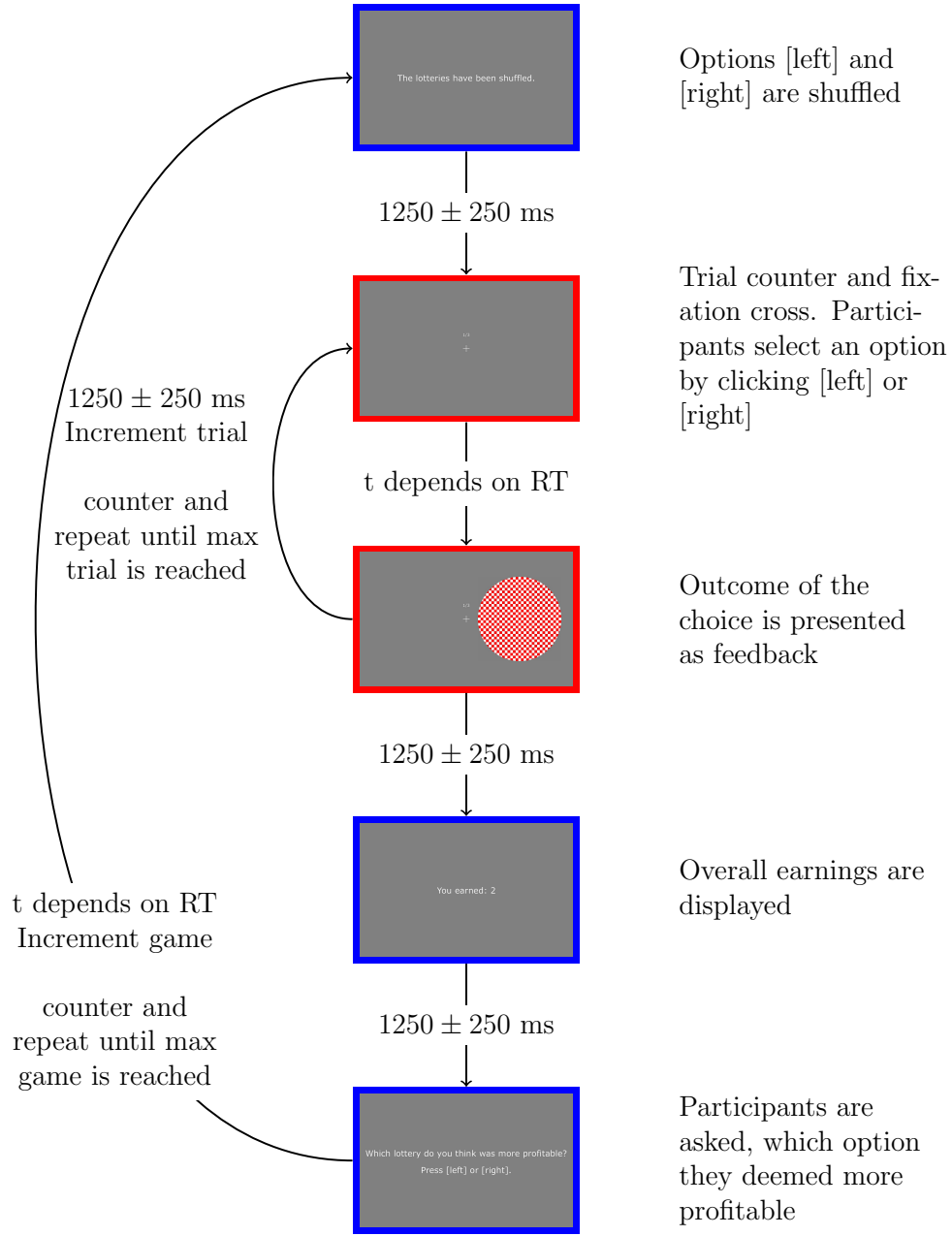


Figure 3: Experimental flow of the bandit paradigm. Red colors indicate the trial loop where participants explore and exploit the options. Once all trials have been spent within the trial loop, a transition to the blue game loop occurs, resetting the environment for a new game and leading to the trial loop again, until all games are exhausted. Note: RT=Reaction time, ms=milliseconds, t=time

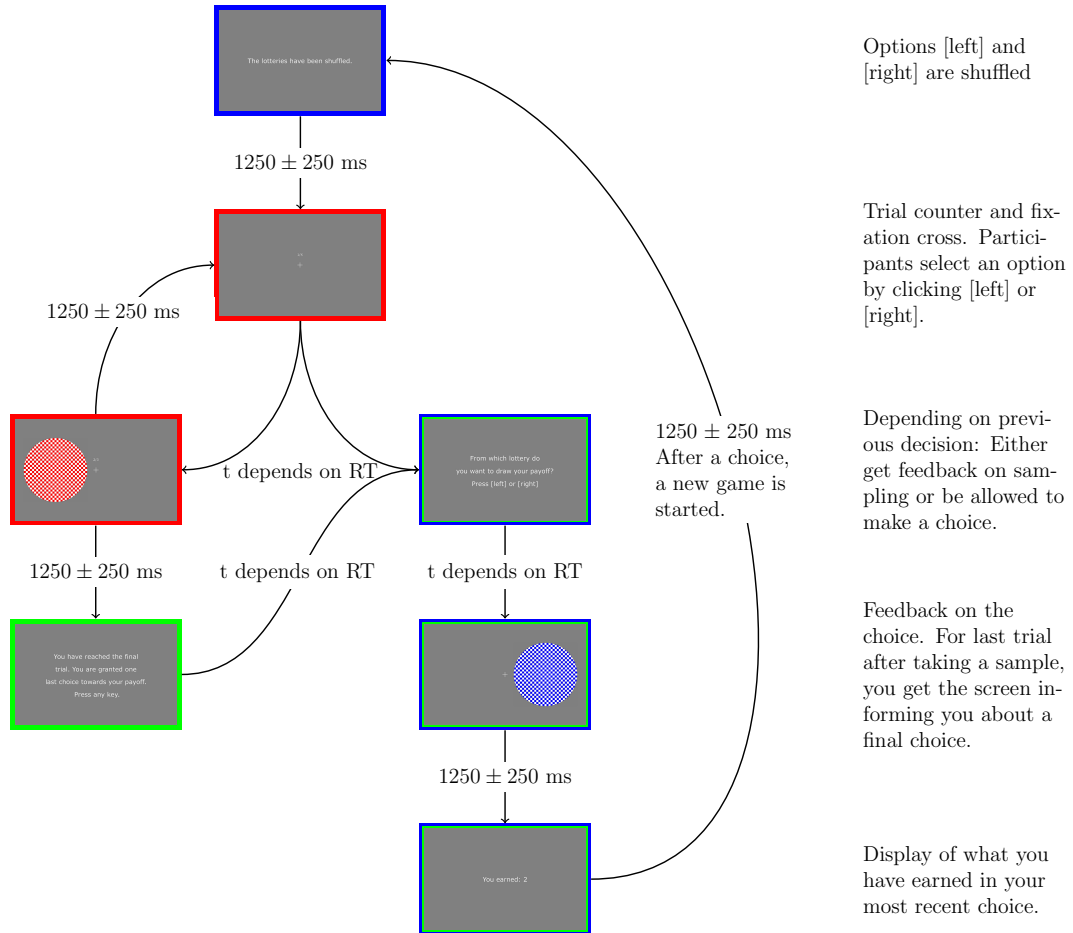


Figure 4: Experimental flow of the sampling paradigm. Red colors indicate the sampling loop, from which one can transition to the blue choice loop, leading back to the sampling loop. Once the maximum of trials has been reached, a transition to the green loop for a last choice occurs. Note: RT=Reaction time, ms=milliseconds, t=time.

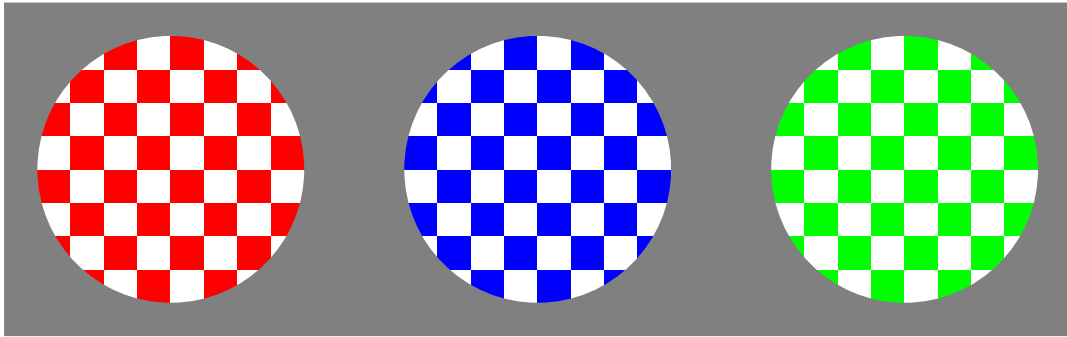


Figure 5: Three stimuli against a background as used in the experiment . The red and blue stimuli represent either a win or lose outcome. The green stimulus represents a distractor stimulus.