

Emergence of Addictive Behaviors in Reinforcement Learning Agents

Vahid Behzadan,¹ Roman V. Yampolskiy,² Arslan Munir¹

¹Kansas State University

²University of Louisville

behzadan@ksu.edu, roman.yampolskiy@louisville.edu, amunir@ksu.edu

Abstract

This paper presents a novel approach to the technical analysis of wireheading in intelligent agents. Inspired by the natural analogues of wireheading and their prevalent manifestations, we propose the modeling of such phenomenon in Reinforcement Learning (RL) agents as psychological disorders. In a preliminary step towards evaluating this proposal, we study the feasibility and dynamics of emergent addictive policies in Q-learning agents in the tractable environment of the game of Snake. We consider a slightly modified version of this game, in which the environment provides a “drug” seed alongside the original “healthy” seed for the consumption of the snake. We adopt and extend an RL-based model of natural addiction to Q-learning agents in these settings, and derive sufficient parametric conditions for the emergence of addictive behaviors in such agents. Furthermore, we evaluate our theoretical analysis with three sets of simulation-based experiments. The results demonstrate the feasibility of addictive wireheading in RL agents, and provide promising venues of further research on the psychopathological modeling of complex AI safety problems.

A necessary requirement for both current and emerging forms of Artificial Intelligence (AI) is the need for robust specification of objectives for AI agents. Currently, a prominent framework for goal-based control of intelligent agents is Reinforcement Learning (RL) (Sutton and Barto 2018). At its core, the objective of an RL agent is to optimize its actions such that an externally-generated reward signal is maximized. However, RL agents are prone to various types of AI safety problems, among which wireheading is subject to growing interest (Yampolskiy 2014). This problem is generally defined as the manifestation of behavioral traits that pursue the maximization of rewards in ways that do not align with the long-term objectives of the system (Yampolskiy 2016). Considering the roots of this paradigm in neuroscientific literature, (Yampolskiy 2014) presents the argument that wireheading is a commonly observed behavior among humans, instances of which are manifested in traits such as substance addiction (Montague, Hyman, and Cohen 2004). This argument is further supplemented with an investigation of wireheading in AI, leading to the conclusion that wireheading in rational self-improving optimizers is a real and open problem. In recent years, various studies have emphasized on the vitality of this problem in the domain of AI safety (e.g., (Amodei et al. 2016)), and some

have proposed solutions for limited instances of wirehead in RL agents (e.g., (Everitt and Hutter 2016)). Yet, the growing complexity of the current and emerging application settings for RL gives rise to the need for tractable approaches to the analysis and mitigation of wireheading in such agents.

In response to this growing complexity, a recent paper by the authors (Behzadan, Munir, and Yampolskiy 2018) presents an analogy between AI safety problems and psychological disorders, and proposes the adoption of a psychopathological abstraction to capture the problems arising from the deleterious behaviors of AI agents in a tractable framework based on the available tools and models of psychopathology. In particular, (Behzadan, Munir, and Yampolskiy 2018) mentions that the RL framework, which itself is inspired by the neuroscientific models of the dopamine system (Sutton and Barto 2018), has been adopted by neuroscientists to develop models of psychological disorders such as schizophrenia and substance addiction (Montague, Hyman, and Cohen 2004). Accordingly, the authors propose to exploit this bidirectional relationship to investigate the complex problems of AI safety.

To study the feasibility of the proposals in (Behzadan, Munir, and Yampolskiy 2018), this paper adopts the RL-based model of substance addiction in natural agents (Redish 2004) to analyze the problem of wireheading in RL agents. To this end, we investigate the emergence of addictive behaviors in a case study of an RL agent training to play the well-known game of Snake (Martti 2002) in an environment that provides a “drug” seed in addition to the typical, “healthy” seed for the snake. By extending the formulation of (Redish 2004) to Q-learning, we analyze the sufficient conditions for the emergence of addictive behavior, and verify this theoretical analysis via simulation-based experiments. The remainder of this paper provides the required background on RL and RL-based modeling of addiction, details our theoretical analysis, and presents the experimental results. The paper concludes with remarks on the significance and potentials of the results.

Background

This section presents an overview of RL and the relevant terminology, as well as a summary of the work by Redish (Redish 2004) on modeling addiction using the RL framework. Readers interested in further details of either topics may refer

fer to (Sutton and Barto 2018) and (Montague, Hyman, and Cohen 2004).

Reinforcement Learning

Reinforcement learning is concerned with agents that interact with an environment and exploit their experiences to optimize a decision-making policy. The generic RL problem can be formally modeled as a Markov Decision Process (MDP), described by the tuple $MDP = (S, A, R, P)$, where S is the set of reachable states in the process, A is the set of available actions, R is the mapping of transitions to the immediate reward, and P represents the transition probabilities (i.e., dynamics), which are initially unknown to RL agents. At any given time-step t , the MDP is at a state $s_t \in S$. The RL agent’s choice of action at time t , $a_t \in A$ causes a transition from s_t to a state s_{t+1} according to the transition probability $P_{s_t, s_{t+1}}^{a_t}$. The agent receives a reward r_{t+1} for choosing the action a_t at state s_t . Interactions of the agent with MDP are determined by the policy π . When such interactions are deterministic, the policy $\pi : S \rightarrow A$ is a mapping between the states and their corresponding actions. A stochastic policy $\pi(s)$ represents the probability distribution of implementing any action $a \in A$ at state s . The goal of RL is to learn a policy that maximizes the expected discounted return $E[R_t]$, where $R_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k}$; with r_t denoting the instantaneous reward received at time t , and γ is a discount factor $\gamma \in [0, 1]$. The value of a state s_t is defined as the expected discounted return from s_t following a policy π , that is, $V^\pi(s_t) = E[R_t | s_t, \pi]$. The action-value (Q-value) $Q^\pi(s_t, a_t) = E[R_t | s_t, a_t, \pi]$ is the value of state s_t after using action a_t and following a policy π thereafter.

As a value function-based solution to the RL problem, the Q-learning method estimates the optimal action policies by using the Bellman formulation $Q_{i+1}(s, a) = \mathbf{E}[R + \gamma \max_a Q_i]$ as the iterative update of a value iteration technique. Practical implementation of Q-learning is commonly based on function approximation of the parametrized Q-function $Q(s, a; \theta) \approx Q^*(s, a)$. A common technique for approximating the parametrized non-linear Q-function is via neural network models whose weights correspond to the parameter vector θ . Such neural networks, commonly referred to as Q-networks, are trained such that at every iteration i , the following loss function is minimized:

$$L_i(\theta_i) = \mathbf{E}_{s, a \sim \rho(\cdot)} [(y_i - Q(s, a; \theta_i))^2] \quad (1)$$

where $y_i = \mathbf{E}[R + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$, and $\rho(s, a)$ is a probability distribution over states s and actions a .

RL Model of Addiction

One of the earliest computational models of addiction is the seminal work of Redish in (Redish 2004). In this paper, Redish assumes the hypothesis that addictive drugs access the same neurophysiological mechanisms as natural learning systems, which can be modeled through the Temporal-Difference RL (TDRL) algorithm (Sutton and Barto 2018). TDRL learns to predict rewards by minimizing a prediction error (i.e., reward-error signal), which, in the natural

brain, is believed to be carried by dopamine. Many addictive substances, such as cocaine, increases the dopamine levels. Redish hypothesizes that this noncompensable drug-induced increase of dopamine may lead to incorrect optimizations in TDRL. Considering that the goal of TDRL is to correctly learn the value of each state ($V(s_t)$), TDRL learns the value function by calculating two equations per each action taken by the agent. If the agent leaves state s_t and enters state s_{t+1} and received the reward r_{t+1} , then the corresponding reward-error signal, denoted by δ , is given by:

$$\delta(t+1) = \gamma[R(s_{t+1}) + V(s_{t+1})] - V(s_t) \quad (2)$$

Then, $V(s_t)$ is updated as:

$$V(s_t) \leftarrow V(s_t) + \nu \delta, \quad (3)$$

where ν is a learning rate parameter. The TDRL algorithm stops when the value function correctly predicts the rewards. The value function can be seen as a compensation for the reward, as the change in the perceived value of taking action a_t leading to the state transition $s_t \rightarrow s_{t+1}$ counterbalances the reward achieved on entering state s_{t+1} . This happens when $\delta = 0$. However, cocaine and similar addictive drugs produce a transient surge in dopamine, which can be explained by the hypothesis that the drug-induced surge in δ cannot be compensated by changes in the value. In other words, the effect of addictive drugs is to induce a positive reward-error signal regardless of the change in value function, thus making it impossible for the agent to learn a value function that cancels out this positive error. As a result, the agent learns to assign more value to the states leading to the dopamine surge, thus giving rise to the drug-seeking behavior of addicted agents.

Case Study : RL Addiction in Snake

To investigate the feasibility of addictive wireheading in RL agents, we consider the game of Snake (Martti 2002) for formal and experimental analysis. The most basic form of Snake is played by one player who controls the direction of a constantly-moving snake in a grid, with the goal of consuming as many seeds as possible by running the snake into them. The seeds appear in random positions on the grid, and the consumption of each seed increases the length of the snake’s tail. The game is terminated if the snake runs into the grid walls or its own tail, thus maneuvering becomes progressively more difficult as the snake consumes more seeds.

In this study, the game is modified to include two types of edible items: one is the classical seed that increases the length of snake L_s by 1 unit, and a “drug” seed that increases L_s by u units. The instantaneous reward values in this setting is defined by:

$$r_t = \begin{cases} r_c & \text{if agent consumes a seed,} \\ k \cdot r_c & \text{if agent consumes a drug,} \\ 0 & \text{otherwise} \end{cases} \quad (4)$$

The objective of the agent is to maximize the return, defined as $R = \sum_{t=0}^T r_t$, where T is the terminal time of an episode. We adopt the formalism of Q-learning as an instance of the TD-learning approach.

The questions that we target in this study are two-fold: first is to analyze whether addictive behaviors may emerge in a Q-learning agent training in this environment, and second is to establish the parametric boundaries of the reward function for such behavior to emerge. The following section presents a formal analysis of these two problems.

Analysis

First, we define addictive behavior as the compulsive pursuit of trajectories that may maximize short-term rewards, but defy the core objective of maximizing the long-term cumulative reward of the agent. At a state s_d where the agent can take action a_m to consume a drug (i.e., move into a cell that contains a drug seed), the Q -value is given by:

$$Q(s_d) = k.r_c + \gamma V(s_{d+1}^m), \quad (5)$$

where $\gamma \in [0, 1]$ is the discount factor and $V(s_{d+1}^m)$ is the value of the resulting state s_{d+1}^m . Alternatively, if the agent takes any action a_g other than a_m , the Q -value is given by:

$$Q(s_d, a_g) = r_c + \gamma V(s_{d+1}^g) \quad (6)$$

The manifestation of addiction can be formulated as:

$$\gamma V(s_{d+1}^m) < \gamma V(s_{d+1}^g) \quad (7)$$

$$Q(s_d, a_m) > Q(s_d, a_g) \quad (8)$$

Eq. (7) can be reformulated as

$$V(s_{d+1}^m) = V(s_{d+1}^g)/l_{d+1}, \quad (9)$$

where $l_{d+1} > 1$. From Eq. (8) we have:

$$k.r_c + \gamma V(s_{d+1}^g)/l_{d+1} > r_c + \gamma V(s_{d+1}^g) \quad (10)$$

which can be rearranged as:

$$\frac{(k-1).r_c}{\gamma(1-1/l_{d+1})} > V(s_{d+1}^g) \quad (11)$$

To obtain a sufficient upper bound for emergence of addiction, we find the maximum possible value of $V(s_{d+1}^g)$ as follows: in an $n \times n$ grid, the maximum possible score is achieved when all elements of the grid are filled with the length of the agent. Considering the assumption in Eq. (7), an upper bound for the game score (and hence for state value) is $V_{max} = r_c(n^2 - L_0)$, where L_0 is the initial length of the snake. Therefore, a sufficient condition on k , r_c , and γ for manifestation of addiction is:

$$\frac{(k-1)}{\gamma} > n^2 - L_0. \quad (12)$$

Also, for the condition of Eq. (7) to hold, it is necessary for k to be set such that:

$$k.r_c(n^2 - L_0)/u < r_c(N^2 - L_0)/1 \implies k/u < 1 \quad (13)$$

Experimental Verification

To evaluate the validity of our analysis, we developed the environment of Snake according to the previously discussed specifications. The environment is comprised of an $n = 8 \times 8$ grid, and the initial length of the snake is set to $L_0 = 4$

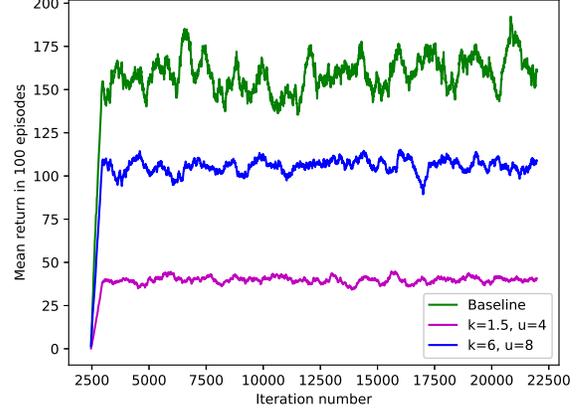


Figure 1: Averaged results of the training process up to 22000 iterations for three experiments

grid cells. At any step, the grid contains two randomly positioned objects, one is the healthy seed (depicted in red), and the other is a drug (colored in blue). Furthermore, we implemented a tabular Q-learning algorithm with iterative update to train in this environment according to the reward function of Eq. (4). The exploration mechanism used in our Q-learning implementation is ϵ -greedy, with the initial value of $\epsilon = 0.99$. We consider a constant discount factor $\gamma = 0.9$, and initialize the table of Q-values to 0. We also consider the instantaneous reward of consuming healthy seeds to be $r_c = 20$. Based on the parametric boundaries derived in Eq. (12) and Eq. (13), we performed three experiments. First, we considered the baseline case in which the consumption of drugs does not produce any rewards or length growth (i.e., $k = u = 0$). For the second experiment, we consider a small value of $k = 1.5$, which does not necessarily abide by the sufficient condition of Eq. (12). Simultaneously, we set $u = 4$, which does satisfy the condition of Eq. (13). In the third experiment, we chose $k = 6$ and $u = 8$ to satisfy both of the derived conditions. To verify the statistical significance of results, the training process of each experiment was repeated 20 times up to 22000 iterations, and the test-time experiments were repeated 100 times each.

Figure 1 demonstrates the training results obtained from the three experiments. It is observed that the baseline case has achieved significantly higher average scores in the same amount of time as the other two cases. Furthermore, the results indicate that the agents training in an environment that includes drug-induced rewards fail to converge towards optimal performance in the observed periods of training. It is also noteworthy that both of the drug-consuming agents reach relatively stable sub-optimal performances in roughly the same time that the healthy agent takes to reach its peak cumulative performance. Moreover, the better performance of the third experiment compared to the second can be explained by the significantly higher instantaneous reward values produced from consuming the drug seeds, which noticeably enhance the average performance in comparison to the

second experiment with lower values of drug-induced rewards.

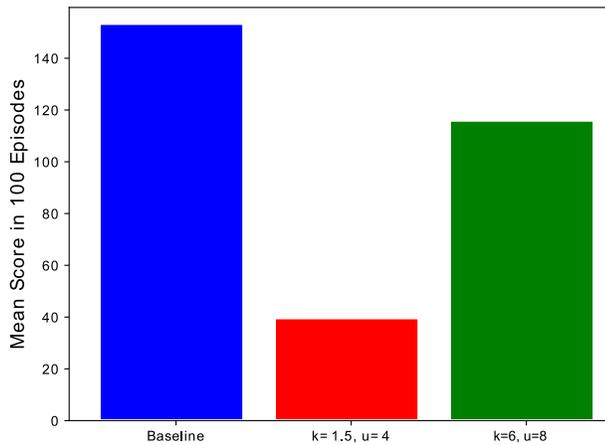


Figure 2: Test-time performance of agents in each experiment

The test-time performance of the agents trained in aforementioned environments is illustrated in Figure 2. These results are in agreement with those of Figure 1, as the baseline agents demonstrate superior performance in gaining cumulative rewards, as opposed to the agents trained under drug-induced rewards. Furthermore, Figure 3 presents a comparison between the number of healthy seeds and drugs consumed by each agent at test-time. As expected, the baseline results demonstrate a significantly higher consumption of healthy seeds, and the minor levels of drug consumption are due to unintended collisions with the drug seeds during game play. It is interesting to note the similarity in the consumption levels of agents trained with drug-induced rewards. In both cases, the agents consume slightly more drugs than healthy seeds, which indicates bias towards the short-term drug-induced surges of rewards over the pursuit of healthy seeds. Although, the difference between the averaged levels of healthy and drug seed consumption is not significant, which may indicate that the agents learned a balanced sub-optimal policy, resulting in confinement within local optima. While this problem can be resolved via enhanced randomization and exploration strategies, one shall consider the effect of this deficiency on sample-efficiency and the consequent limitations of real-world applications.

Conclusion

We studied the feasibility of adopting the RL-based model of substance addiction in natural agents to analyze the dynamics of wireheading in RL-based artificial agents. We presented an analytical extension to a TD-learning based model of addiction, and established sufficient parametric conditions on reward functions for the emergence of addictive behavior in AI agents. To verify this extension, we presented experimental results obtained from Q-learning agents learning to play the game of Snake, which is modified to include

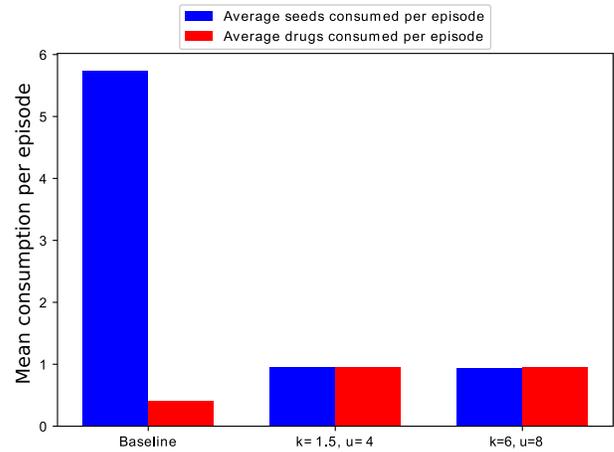


Figure 3: Test-time consumption of healthy seeds and drugs in three experiments

drug-induced surges in instantaneous rewards. The results demonstrate the promising potential of adopting the psychopathological models of mental disorders in the analysis of complex AI safety problems.

References

Amodei, D.; Olah, C.; Steinhardt, J.; Christiano, P.; Schulman, J.; and Mané, D. 2016. Concrete problems in ai safety. *arXiv preprint arXiv:1606.06565*.

Behzadan, V.; Munir, A.; and Yampolskiy, R. V. 2018. A psychopathological approach to safety engineering in AI and AGI. In *Computer Safety, Reliability, and Security - SAFE-COMP 2018 Workshops, Västerås, Sweden, September 18, 2018, Proceedings*, 513–520.

Everitt, T., and Hutter, M. 2016. Avoiding wireheading with value reinforcement learning. In *Artificial General Intelligence*. Springer. 12–22.

Martti, H. 2002. Nokia: the inside story.

Montague, P. R.; Hyman, S. E.; and Cohen, J. D. 2004. Computational roles for dopamine in behavioural control. *Nature* 431(7010):760.

Redish, A. D. 2004. Addiction as a computational process gone awry. *Science* 306(5703):1944–1947.

Sutton, R. S., and Barto, A. G. 2018. *Reinforcement learning: An introduction*. MIT press.

Yampolskiy, R. V. 2014. Utility function security in artificially intelligent agents. *Journal of Experimental & Theoretical Artificial Intelligence* 26(3):373–389.

Yampolskiy, R. V. 2016. Taxonomy of pathways to dangerous artificial intelligence. In *AAAI Workshop: AI, Ethics, and Society*.