



The application of temporal difference learning in optimal diet models



Jan Teichmann^{a,*}, Mark Broom^a, Eduardo Alonso^b

^a Department of Mathematical Science, City University London, Northampton Square, London EC1V0HB, United Kingdom

^b Department of Computer Science, City University London, Northampton Square, London EC1V0HB, United Kingdom

AUTHOR - HIGHLIGHTS

- We apply model-free reinforcement learning to optimal diet models.
- The presented model incorporates uncertainty of changing environments.
- The model predicts effects of Batesian mimics and aposematism on predators diet choice and energy intake.
- The model uses a precondition of exploration of the action space for successful aversion formation.
- Conflicting rewards lead to foraging behaviour which is conditionally suboptimal in fixed environments but allows better adaptation in changing environments.

ARTICLE INFO

Article history:

Received 15 April 2013

Received in revised form

28 August 2013

Accepted 30 August 2013

Available online 11 September 2013

Keywords:

Optimal diet

Batesian mimicry

Predator–prey

Taste sampling

Temporal difference learning

ABSTRACT

An experience-based aversive learning model of foraging behaviour in uncertain environments is presented. We use Q-learning as a model-free implementation of Temporal difference learning motivated by growing evidence for neural correlates in natural reinforcement settings. The predator has the choice of including an aposematic prey in its diet or to forage on alternative food sources. We show how the predator's foraging behaviour and energy intake depend on toxicity of the defended prey and the presence of Batesian mimics. We introduce the precondition of exploration of the action space for successful aversion formation and show how it predicts foraging behaviour in the presence of conflicting rewards which is conditionally suboptimal in a fixed environment but allows better adaptation in changing environments.

© 2013 Elsevier Ltd. All rights reserved.

1. Introduction

Predators have to secure a high energy intake in the face of changing and uncertain environments. Through the evolution of predator–prey interactions manifold mechanisms have emerged to avoid predation. The so-called secondary defences commonly involve the possession of toxins or deterrent substances which are not directly observable by predators. However, many defended species use conspicuous signals as warning flags in combination with their secondary defences (aposematism).

There is a wide body of theory which addresses the emergence and evolution of aposematism (Ruxton et al., 2004; Yachi and Higashi, 1998; Broom et al., 2006; Leimar et al., 1986; Lee et al., 2011; Marples et al., 2005). However, the field of aposematism has a renewed interest in the role of the predator and details of the predator's aversive

learning process. In particular, the role of aposematism in memory formation has been widely studied (Speed, 2000; Svádová et al., 2009; Skelhorn and Rowe, 2006; Johnston and Burne, 2008; Speed and Ruxton, 2005). As the selective agent, aversive learning is an important aspect of predator avoidance. It has been shown that predation of defended prey is rather a state dependent decision and predators can increase their attack rates on defended prey e.g. when particularly hungry (Barnett et al., 2007; Sherratt, 2003). There have been suggestions of an interaction of appetitive learning with aversive learning to explain predator behaviour of ingesting toxins in these situations (Hagen et al., 2009).

An interesting perspective is to look at the predator and the consequences of aposematism in combination with aversive learning on the predator's diet and energy intake. In particular, the role of mimics in the evolution of aposematism and their effect on foraging is not very well understood (Gamberale-Stille and Tullberg, 2001; Lev-Yadun and Gould, 2007; Svádová et al., 2009; Holen, 2013). A predator may utilise sampling to distinguish between the toxic model and the mimic (Gamberale-Stille and Tullberg, 2001; Darst, 2006; Holen, 2013).

* Corresponding author. Tel.: +44 7716324160.

E-mail addresses: Jan.Teichmann.1@city.ac.uk (J. Teichmann), Mark.Broom.1@city.ac.uk (M. Broom), e.alonso@city.ac.uk (E. Alonso).

The traditional way of analysing and predicting foraging behaviour is the application of optimal foraging theory (OFT) which maximises the predator's net fitness per unit time (MacArthur and Pianka, 1966; Stephens and Krebs, 1987; Sih and Christensen, 2001). However, OFT has well known limitations: OFT usually fails to correctly predict foraging behaviour on mobile prey in complex environments (Sih and Christensen, 2001; Pyke, 1984; Perry and Pianka, 1997). It can be argued that OFT was never intended for predictions in the case of mobile prey and that the optimisation per unit time omits the uncertainty of more complex environments. There are models which address optimal foraging under the constraints of risk and uncertainty and previously extended OFT with learning (McNamara and Houston, 1985). The two main approaches to optimal behaviour in dynamic decision making are dynamic programming (DP) and stochastic optimal control methods (e.g. Bayesian decision theory) (Houston and McNamara, 1982; Stephens and Charnov, 1982; McNamara and Houston, 1985; Mangel and Clark, 1986; McNamara et al., 2006). Especially dynamic programming found wider application in behavioural ecology and has been used in models of dynamic decision making to identify optimal behaviour numerically (Clark and Mangel, 2000). These models have all in common that they are *model based*: they depend on a representation of the environment in the form of a model developed from expert knowledge and the learning objective is to find the parameters which optimise the representational model.

On the contrary, a normative framework of rational decision making in a changing and complex environment is reinforcement learning (RL). RL combines the computational task of maximising rewards and the algorithmic implementation of natural learning without an explicit supervisory control signal (Mitchell, 1997; Sutton and Barto, 1998).

Neural correlates of behaving animals show that reinforcement signals in the brain represent the reward prediction error rather than a direct reward-reinforcement relation. Temporal difference (TD) learning reflects these insights by representing states and actions in terms of predictions about future rewards (Niv, 2009; Berns et al., 2001). Additionally, TD learning is *model-free*: the environment is represented by moving targets rather than by a model and the learning objective is to iteratively update the targets towards its true values based on experience from interactions with the environment. TD learning has been widely used in artificial systems to choose appropriate actions in complex non-stationary environments. Furthermore, the computational theories are increasingly supported by experimental data describing the activity of dopaminergic neurons, mediate reward-processing and reward-dependent learning (Schultz et al., 1997; Montague et al., 2004; Daw et al., 2006; Dayan and Niv, 2008). In the greater picture of learning algorithms, TD learning resides between dynamic programming and Monte Carlo methods (Sutton and Barto, 1998).

The rest of the paper is structured as follows. In the next two sections we apply a TD learning algorithm in a model of predator's interaction with conspicuous prey to gain insights on how aversive learning influences foraging in uncertain environments, and present the results. Next we discuss the main findings and discuss similarities and differences to the optimisation approach of traditional OFT. In particular, we will compare TD learning with McNamara and Houston (1985) and Sherratt (2003). We will conclude that TD learning is a new approach to OFT which is better suited for modelling foraging in dynamic environments with learning.

2. Methodology

In our model the predator interacts with its environment to find an optimal foraging strategy to optimise its rewards. The predator's environment offers a stable background of alternative food sources. Additionally, the predator has the choice to include a

conspicuous looking type of prey into its diet. However, the conspicuous prey population may consist of an aposematic model species and a Batesian mimic species. We assume the environment to be uncertain with non-stationary parameters over a predator's lifespan.

2.1. Temporal difference learning

The predator is not able to distinguish models and mimics based on their appearance and utilises experience to learn the optimal foraging behaviour. Based on the growing understanding of learning at the computational and neural level we use Temporal difference (TD) learning to implement the predator's aversive learning: in particular, we use Q-learning (Watkins, 1989). The learning process consists of a reward prediction termed the *action-value function* (1) of taking action a in state s at iteration k :

$$Q(s, a) = E\{R_k | s_k = s, a_k = a\}. \quad (1)$$

The condition for the action-value function and Q-learning is for the Markov property to hold

$$P\{s_{k+1} = s', r_{k+1} = r | s_k, a_k\}. \quad (2)$$

The reinforcement signal consists of the TD error of the reward prediction based on experienced rewards following an undertaken action a . Finally, the Q-learning update rule is utilised in order to minimise the prediction error (Barto et al., 1983; Sutton and Barto, 1998).

Each action taken has a state dependent subsequent reward signal termed r_{k+1} . The predator not only takes immediate rewards into account but also the sum of discounted future rewards (3) with K being the end of an episode and γ being the discount factor. This combines an ubiquitous interest into rewards with the uncertainty of future events as follows:

$$\begin{aligned} R_k &= \sum_{i=0}^K \gamma^i r_{k+i+1} = r_{k+1} + \sum_{i=1}^K \gamma^i r_{k+i+1} \\ &= r_{k+1} + \gamma \sum_{i=0}^T \gamma^i r_{k+i+2} \\ &= r_{k+1} + \gamma R_{k+1}. \end{aligned} \quad (3)$$

The predator uses the experienced immediate reward r_{k+1} to minimise the prediction error by updating its state dependent action-value function using the Q-learning method. The algorithmic representation of the Q-learning update process is presented in (4) with α being the learning rate following the derivation in (3) as follows:

$$Q'(s_k, a_k) \leftarrow Q(s_k, a_k) + \alpha \left(\underbrace{r_{k+1} + \gamma \max_{a_{k+1}} Q(s_{k+1}, a_{k+1})}_{\text{target}} - Q(s_k, a_k) \right). \quad (4)$$

Q-learning is an iterative algorithm which uses the immediate experienced reward to form a target with Q' being the new estimate for Q . Thereby, Q-learning bases its update partially on a prevailing estimate $Q(s_{k+1}, a_{k+1})$ which is known as bootstrapping. Q-learning is widely used to model Markov decision problems and under certain conditions, Q-learning has been proved to converge to optimality (Watkins and Dayan, 1992). For a more detailed introduction of the Q-learning algorithm we refer to the supplementary material in Appendix A.

Finally, the predator uses the Gibbs soft-max policy which is the probability of taking action a in state s under stochastic policy π to translate its action-value predictions into foraging behaviour:

$$\pi(s, a) = P\{a_k = a | s_k = s\} = \frac{\exp(Q(s, a))}{\sum_a \exp(Q(s, a))}. \quad (5)$$

2.2. The predator's interaction with conspicuous prey

We term the action of falling back on the alternative background food sources as $a=0$ and the action of attacking conspicuous prey as $a=1$.

We assume the population of conspicuous prey consists of a fraction p of Batesian mimics and a fraction $1-p$ of defended models. The reward signal for the alternative stable background food source is $r_{k+1} = \{1 \mid a=0\}$. The reward signal for ingesting a mimic individual is $r_{k+1} = \{2 \mid a=1, i = \text{mimic}\}$ and $r_{k+1} = \{1-t^2 \mid a=1, i = \text{model}\}$ for ingesting a model individual with toxicity t . These reward signals do not have to represent necessarily fitness related entities (Pyke, 1984). In our model we simply assume mimics to be rewarding and that toxicity has a non-linear effect on the reward, which seems like a reasonable assumption.

We consider two different cases (Fig. 1):

1. The predator has the ability to use taste-sampling to distinguish models from mimics assuming that the model's toxicity t operates as a clue to the predator. This foraging strategy is also called *go-slow behaviour* (Guilford, 1994). The probability of rejecting a model based on taste-sampling is given as follows:

$$d(t) = 1 - \frac{1}{1 + d_0 t}. \quad (6)$$

2. The predator has no ability to distinguish mimics and models and the encounter is solely frequency dependent i.e. $d_0=0$ in Eq. (6).

3. Results

In the case of the predator being unable to distinguish models from mimics ($d_0=0$) the average reward signal is solely frequency dependent and given as

$$R = \begin{cases} 1 & \text{if } a=0 \\ 2p + (1-t^2)(1-p) & \text{if } a=1. \end{cases} \quad (7)$$

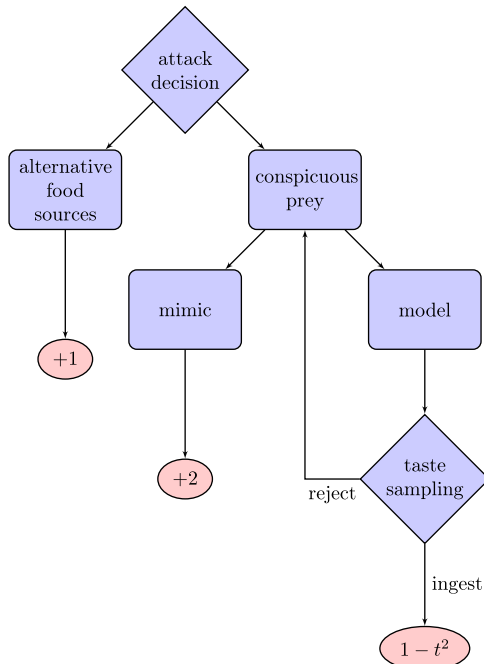


Fig. 1. The predator's interaction with its environment and possible reward signals. The predator has the ability to recognise toxic models by taste-sampling. t stands for the toxicity of defended models.

If the predator utilises taste-sampling it can distinguish models from mimics based on the model's toxicity and will not ingest the toxic model with probability $d(t)$ given in (6). After the predator rejects a conspicuous prey individual it will stay in the locality and forage for another conspicuous prey individual. The average reward signal incorporating taste sampling derives from the geometric series and is given as follows:

$$R = \begin{cases} 1 & \text{if } a=0 \\ 2p \frac{1}{1-(1-p)d(t)} + (1-t^2)(1-p) \frac{(1-d(t))}{1-(1-p)d(t)} & \text{if } a=1. \end{cases} \quad (8)$$

To obtain the optimal diet we find the correct, discounted action-value function by solving the TD learning problem:

$$0 = R + \gamma \max_{a_{k+1}} Q(s_{k+1}, a_{k+1}) - Q(s_k, a_k). \quad (9)$$

Figs. 2 and 3 show the probability of an experienced predator attacking conspicuous prey based on the frequency of mimics (p) and the model's toxicity (t). We define aversiveness as $\pi(a=1) < 0.5$ with the threshold toxicity (t^*) given in (10) for which conspicuous prey becomes aversive and $R(a=0, t^*) = R(a=1, t^*)$ holds as follows:

$$t^* = \begin{cases} \sqrt{\frac{p}{p-1}} & \text{if } d_0 = 0 \\ -\frac{\sqrt{p^2 d_0^2 - 4p^2 + 4p + p d_0}}{2p-2} & \text{otherwise.} \end{cases} \quad (10)$$

We see that taste-sampling lowers the aversiveness of defended conspicuous prey when mimics are present.

Figs. 4 and 5 show the average reward (R) of an experienced predator. Mimics increase the average reward of the predator through increased foraging on non-averse conspicuous prey. Conversely, increasing toxicity of the models reduces the average reward for the predator until the increasing toxicity intake from mistakenly ingested models becomes aversive.

4. Discussion

We apply Q-learning to the problem of optimal foraging behaviour of an experienced predator in an uncertain environment. Our motivation lays in the recognised importance of aversive learning in aposematism and the difficulties of the classical OFT approach to predict foraging behaviour on mobile prey (Sih and Christensen, 2001). In the case of mobile prey additional factors of prey handling and uncertainty need to be considered, making the OFT model increasingly complex (Holen, 2013). Instead, reinforcement learning offers a normative framework of rational decision making in a changing and complex environment with growing evidence of neural correlates.

The TD learning based approach puts the emphasis on experience including discounted future rewards and requires exploration of the action space. This is fundamentally different from the OFT models of net fitness maximisation per unit time. It has been long argued that a learning animal cannot be foraging optimally and vice versa (Ollason, 1980).

We hypothesise that a non-stationary environment introduces great uncertainty on the prey-population's parameters t and p which selects for learning in evolving predators to adapt quicker to their changing environment. Evidence for this claim has to come from an evolutionary model and is subject to future work. To coincide widely with the original OFT methodology, we assume that the learning process is sufficiently faster than the frequency of change of the environment to concentrate solely on the experienced predator and to exclude the iterative learning phase. Furthermore, we assume that the conspicuous prey inhabits a distinct locality. These assumptions allow us to solve the TD learning problem directly (9) and we present the policy a predator adopts through Q-learning.

In the context of previous foraging models which incorporated learning, our learning methodology is model-free. Relevant models,

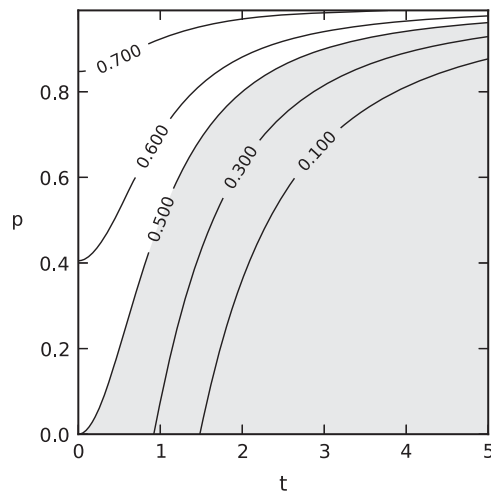


Fig. 2. Predator attack probability (π) of conspicuous prey without taste-sampling ($d_0=0$) and discount rate $\gamma=0.5$ following soft-max policy (5). t stands for the toxicity of models and p for the fraction of mimics. The shaded area indicates aversive toxicity.

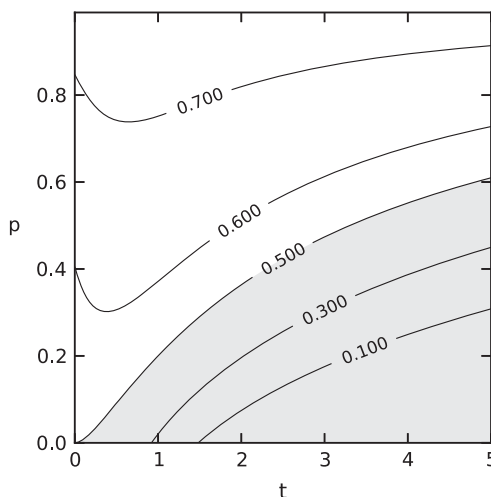


Fig. 3. Predator attack probability (π) of conspicuous prey utilising taste-sampling ($d_0=3$) (6) and discount rate $\gamma=0.5$ following Gibbs soft-max policy (5). t stands for the toxicity of models and p for the fraction of mimics. The shaded area indicates aversive toxicity.

among others, are from McNamara and Houston (1985) and Sherratt (2003). McNamara's learning rule describes a Monte Carlo method using past events to learn the maximum possible long-term rate as defined by the marginal value theorem (Charnov, 1976). It uses discounted experience from past interactions with the environment to optimise a current parameter estimation. The corresponding concept in TD learning is termed *eligibility trace* and is bridging TD learning with Monte Carlo methods. Eligibility traces can make TD learning more efficient but as we exclude the iterative learning phase it has no application in our model. Nevertheless, TD learning is conceptually different as its learning objective is based on bootstrapping future rewards rather than optimising the current estimate of a parameter from past events.

Sherratt's (2003) model uses Bayesian learning based on dynamic programming. The learning objective is to infer the Bayesian posterior mean estimate of the fraction of defended prey in an unknown population from past experience. The model uses Beta distributions in the Bayesian inference to represent an assumed underlying binomial distribution of defence in a group of prey. The main assumption for the application of dynamic programming is the existence of a finite time horizon where the

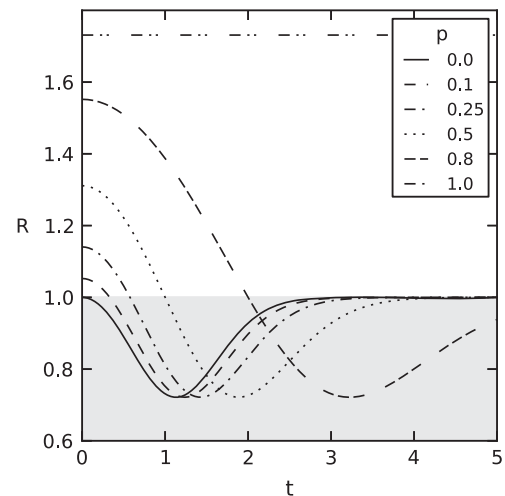


Fig. 4. The predator's average reward (R) from interacting with its environment without taste-sampling ($d_0=0$) and discount rate $\gamma=0.5$. t stands for the toxicity of models and p for representative fractions of mimics. The shaded area indicates suboptimal rewards due to foraging on aversive prey.

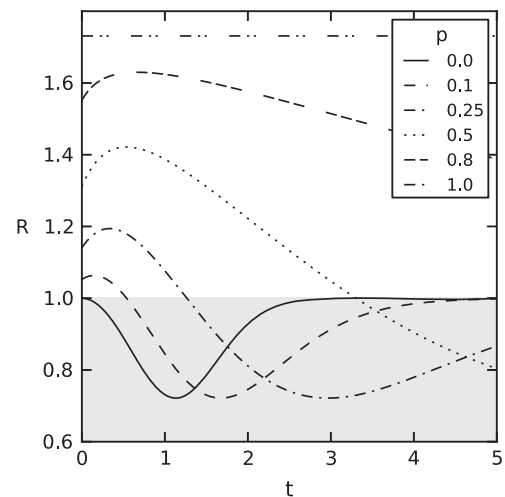


Fig. 5. The predator's average reward from interacting with its environment utilising taste-sampling ($d_0=3$) and discount rate $\gamma=0.5$. t stands for the toxicity of models and p for representative fractions of mimics. The shaded area indicates suboptimal rewards due to foraging on aversive prey.

predator ceases attacking completely. Sherratt's model provides an optimal sampling strategy for novel prey populations with constant values for cost and benefit of an attack. However, the model cannot provide optimal foraging policies in changing populations or when defence is not just binomial distributed.

We conclude that TD learning is a new approach to optimal foraging in dynamic environments where cost-benefit values of attacking prey do not necessarily follow simple distributions. TD learning uses a model free objective which makes it an ideal method for learning in complex and dynamic environments where parameters are subject to constant change.

Our model confirms expected results such as that mimics in general lower the aversiveness of the conspicuous prey population and undermine aposematism. Nevertheless, highly toxic models can sustain aversion even for high frequencies of mimics especially in predators not utilising taste sampling. However, it requires exploration for a predator to gain insights about its environment and to form aversive memory. Therefore, even an aversive prey population experiences some level of predation.

Our model predicts that a taste-sampling predator increases its attack rate on mixed conspicuous prey populations in the case of moderately defended models and rewarding mimics. The taste-sampling predator gains increased rewards from moderately defended models as it allows for better discrimination of models and mimics. This is a contrary finding to [Holen \(2013\)](#) in which mimics benefit from moderately defended models. This difference is founded on the representation of toxins as recovery time in the OFT maximisation approach and the missing occasional ingestion of models to maintain aversion for highly toxic models.

An interesting paradox is the foraging behaviour on aversive prey which reduces the reward for the predator further before recovering through increasingly falling back on alternative background food sources. (The adopted attack policy for certain parameters results in an average reward R which lays in the shaded area in [Figs. 4](#) and [5](#), and is suboptimal.) This is a result of the conflicting reward signals of mimics and models and the necessity of exploration of the action space in the face of uncertainty for successful aversion formation. Additionally, an increasing frequency of mimics slows the switching to alternative food sources through further extended uncertainty. Similar results have been observed in counter conditioning and operant conflict situations ([Williams and Barry, 1966](#); [Blaisdell et al., 2000](#); [Mazur and Ratti, 1991](#); [Matsushima et al., 2008](#)). Our model predicts a fixed amount of average long term toxicity intake which a predator tolerates motivated either by the higher reward signal of ingested mimics or as a consequence of uncertainty. (Although the toxicity of immediate rewards which induce switching to alternative food sources depends on the amount of mimics and the specific rewards, see Eq. (10) and [Figs. 2](#) and [3](#), the average reward function described in Eqs. (7) and (8) has a fixed minimum as presented in [Figs. 4](#) and [5](#).) This foraging behaviour on aversive prey for a specific parameter space is conditionally suboptimal in a stationary environment (even if only during an individuals lifetime) but we note that (a) it reflects what real animals do, and (b) it is a good policy precisely because environments are inherently uncertain.

Summarising, our main conclusions are as follows:

- TD learning is a suitable approach to optimal foraging in changing environments.
- Even aversive prey experiences some level of predation as part of the predator's aversive memory formation.
- Taste-sampling lowers the effective aversiveness of conspicuous prey if mimics are present.
- Intermediate toxicity of aposematic models increases the predator's foraging on conspicuous prey through increased discrimination from taste-sampling and higher average rewards when mimics are rewarding.
- The conflicting reward signals from mimics and models cause uncertainty and conditionally suboptimal foraging behaviour on aversive prey.

- The uncertainty is linked to a fixed amount of average toxicity intake which predators tolerate in order to forage on rewarding mimics before switching to mediocre background food sources.
- Taste-sampling extends the range of parameters where sub-optimal foraging occurs.

Appendix A. Q-learning algorithm

Q-learning is a simple algorithmic implementation of reinforcement learning. Particularly, it is a model free method which allows to learn about Markovian environments from experienced rewards without the necessity of building representations of the environment. Instead, the algorithm uses moving target values.

The predator learns from iterative interactions with its environment. We term the current iteration subscript k . At each iteration k the predator finds itself in state s_k of its environment, accordingly, s_k is the encounter with a particular type of prey in our model. The actual learning process targets the predator's reward prediction following action a_k (respectively, attacking conspicuous or alternative prey) in state s_k termed the action-value function $Q(s_k, a_k)$. This action-value function is an approximation of the actual function $Q^*(s, a)$. Consequently, the aim of the learning process is to find $Q(s_k, a_k) \approx Q^*(s, a)$. The predator is basing its decision process on $Q(s_k, a_k)$ following a decision policy $\pi(s_k, Q(s_k, a_k))$, effectively knowing all of the current Q values gives the probability that we choose to attack or not for the next encounter. This involves an iterative update process which is typically formulated in an algorithmic representation because of its origin in computing, as follows:

$$Q'(s_k, a_k) = Q(s_k, a_k) + \alpha \left(\underbrace{r_{k+1} + \gamma \max_{a_{k+1}} Q(s_{k+1}, a_{k+1})}_{\text{target}} - Q(s_k, a_k) \right) \quad (\text{A.1})$$

TD error

The iterative algorithm expands as follows: at iteration k , the predator interacts with the environment of state s_k which is a realisation from the state space S . Following a certain decision policy π , the predator takes action a_k out of the action space A . As a result of this interaction at iteration k , the predator experiences an immediate reward r_{k+1} . The terminology refers to the experienced reward at the subsequent iteration $k+1$ which emphasise that the reward is in consequence of the predator's action. Next, the predator forms a target value which is a composition of the experienced reward r_{k+1} and discounted future rewards. Thereby, future rewards are a prevailing estimate $Q(s_{k+1}, a_{k+1})$ which is known as *bootstrapping*. The difference between the target value and the estimate at iteration k gives the *temporal-difference (TD) error*. Finally, the Q-learning algorithm updates the estimate $Q(s_k, a_k)$ to $Q'(s_k, a_k)$ towards the formed target value, subsequently reducing the TD error. As the Q-learning algorithm uses bootstrapping, these targets are moving ones. Hence, the update process should progress slowly with α , the learning rate, being a small positive constant. [Fig. A1](#) shows a possible implementation of the Q-learning algorithm as pseudo-code.

```

432  Q ← 0
433  s_k ← s_0
434  WHILE learning DO
435    a_k ← π(s_k, Q)
436    s_(k+1) ← f(s_k, a_k)
437    Q(s_k, a_k) ← Q(s_k, a_k) + α (r_(k+1) +
438      γ max_a Q(s_(k+1), a) - Q(s_k, a_k) )
439    s_k ← s_(k+1)
440
```

Fig. A1. Q-learning algorithm in pseudo-code.

Appendix B. Supplementary data

Supplementary data associated with this article can be found in the online version at <http://dx.doi.org/10.1016/j.jtbi.2013.08.036>.

References

- Barnett, C., Bateson, M., Rowe, C., 2007. State-dependent decision making: educated predators strategically trade off the costs and benefits of consuming aposematic prey. *Behavioral Ecology* 18 (4), 645–651.
- Barto, A.G., Sutton, R.S., Anderson, C.W., 1983. Neuronlike adaptive elements that can solve difficult learning control problems. *IEEE Transactions on Systems, Man, and Cybernetics* 13 (5), 834–846.
- Berns, G.S., McClure, S.M., Pagnoni, G., Montague, P.R., 2001. Predictability modulates human brain response to reward. *The Journal of Neuroscience* 21 (8), 2793–2798.
- Blaisdell, A.P., Denniston, J.C., Savastano, H.J., Miller, R.R., 2000. Counterconditioning of an overshadowed cue attenuates overshadowing. *Journal of Experimental Psychology: Animal Behavior Processes* 26 (1), 74.
- Broom, M., Speed, M., Ruxton, G., 2006. Evolutionarily stable defence and signalling of that defence. *Journal of Theoretical Biology* 242, 32–34.
- Charnov, E.L., 1976. Optimal foraging, the marginal value theorem. *Theoretical Population Biology* 9 (2), 129–136.
- Clark, C.W., Mangel, M., 2000. *Dynamic State Variable Models in Ecology: Methods and Applications: Methods and Applications*. Oxford University Press, USA.
- Darst, C.R., 2006. Predator learning, experimental psychology and novel predictions for mimicry dynamics. *Animal Behaviour* 71 (4), 743–748.
- Daw, N., Doya, K., et al., 2006. The computational neurobiology of learning and reward. *Current Opinion in Neurobiology* 16 (2), 199–204.
- Dayan, P., Niv, Y., 2008. Reinforcement learning: the good, the bad and the ugly. *Current Opinion in Neurobiology* 18 (2), 185–196.
- Gamberale-Stille, G., Tullberg, B.S., 2001. Fruit or aposematic insect? Context-dependent colour preferences in domestic chicks. *Proceedings of the Royal Society B: Biological Sciences* 268, 2525–2529.
- Guilford, T., 1994. Go-slow signalling and the problem of automimicry. *Journal of Theoretical Biology* 170 (3), 311–316.
- Hagen, E., Sullivan, R., Schmidt, R., Morris, G., Kempter, R., Hammerstein, P., 2009. Ecology and neurobiology of toxin avoidance and the paradox of drug reward. *Neuroscience* 160 (1), 69–84.
- Holen, Ø.H., 2013. Disentangling taste and toxicity in aposematic prey. *Proceedings of the Royal Society B: Biological Sciences* 280, 20122588.
- Houston, A.I., McNamara, J., 1982. A sequential approach to risk-taking. *Animal Behaviour* 30, 1260–1261.
- Johnston, A.N., Burne, T.H., 2008. Aposematic colouration enhances memory formation in domestic chicks trained in a weak passive avoidance learning paradigm. *Brain Research Bulletin* 76 (3), 313–316.
- Lee, T.J., Speed, M.P., Stephens, P.A., 2011. Honest signaling and the uses of prey coloration. *American Society of Naturalists* 178, E1–E9.
- Leimar, O., Enquist, M., Sillen-Tullberg, B., 1986. Evolutionary stability of aposematic coloration and prey unprofitability: A theoretical analysis. *American Society of Naturalists* 128, 469–490.
- Lev-Yadun, S., Gould, K., 2007. What do red and yellow autumn leaves signal? *Botanical Review* 73 (4), 279–289, cited by (since 1996) 30.
- MacArthur, R.H., Pianka, E.R., 1966. On optimal use of a patchy environment. *American Naturalist* 100, 603–609.
- Mangel, M., Clark, C.W., 1986. Towards a unified foraging theory. *Ecology* 67, 1127–1138.
- Marples, N.M., Kelly, D.J., Thomas, R.J., 2005. Perspective: the evolution of warning coloration is not paradoxical. *Evolution* 59 (5), 933–940.
- Matsushima, T., Kawamori, A., Bem-Sojka, T., 2008. Neuro-economics in chicks: foraging choices based on amount delay and cost. *Brain Research Bulletin* 76 (3), 245–252.
- Mazur, J., Ratti, T., 1991. Choice behavior in transition: development of preference in a free-operant procedure. *Animal Learning & Behavior* 19, 241–248.
- McNamara, J.M., Houston, A.I., 1985. Optimal foraging and learning. *Journal of Theoretical Biology* 117 (2), 231–249.
- McNamara, J.M., Green, R.F., Olsson, O., 2006. Bayes theorem and its applications in animal behaviour. *Oikos* 112 (2), 243–251.
- Mitchell, T., 1997. *Machine Learning* (McGraw-Hill International Edit), 1st ed. McGraw-Hill Education, October (ISE editions).
- Montague, P.R., Hyman, S.E., Cohen, J.D., 2004. Computational roles for dopamine in behavioural control. *Nature* 431 (7010), 760–767.
- Niv, Y., 2009. Reinforcement learning in the brain. *Journal of Mathematical Psychology* 53 (3), 139–154.
- Ollason, J., 1980. Learning to forage – optimally? *Theoretical Population Biology* 18 (1), 44–56.
- Perry, G., Pianka, E.R., 1997. Animal foraging: past, present and future. *Trends in Ecology & Evolution* 12 (9), 360–364.
- Pyke, G.H., 1984. Optimal foraging theory: a critical review. *Annual Review of Ecology and Systematics* 15, 523–575.
- Ruxton, G., Sherratt, T., Speed, M., 2004. *Avoiding Attack: The Evolutionary Ecology of Crypsis, Warning Signals and Mimicry*. Oxford University Press.
- Schultz, W., Dayan, P., Montague, P.R., 1997. A neural substrate of prediction and reward. *Science* 275 (5306), 1593–1599.
- Sherratt, T.N., 2003. State-dependent risk-taking by predators in systems with defended prey. *Oikos* 103 (1), 93–100.
- Sih, A., Christensen, B., 2001. Optimal diet theory: when does it work, and when and why does it fail? *Animal Behaviour* 61 (2), 379–390.
- Skelhorn, J., Rowe, C., 2006. Prey palatability influences predator learning and memory. *Animal Behaviour* 71 (5), 1111–1118.
- Speed, M.P., 2000. Warning signals, receiver psychology and predator memory. *Animal Behaviour* 60 (3), 269–278.
- Speed, M., Ruxton, G., 2005. Aposematism: what should our starting point be? *Proceedings of the Royal Society B: Biological Sciences* 272 (1561), 431–438.
- Stephens, D.W., Charnov, E.L., 1982. Optimal foraging: some simple stochastic models. *Behavioral Ecology and Sociobiology* 10 (4), 251–263.
- Stephens, D.W., Krebs, J.R., 1987. *Foraging theory*. Princeton University Press.
- Sutton, R.S., Barto, A.G., 1998. *Reinforcement Learning: An Introduction*. Cambridge University Press.
- Svádová, K., Exnerová, A., Štys, P., Landová, E., Valenta, J., Fucíková, A., Socha, R., 2009. Role of different colours of aposematic insects in learning, memory and generalization of naïve bird predators. *Animal Behaviour* 77 (2), 327–336.
- Watkins, C., 1989. *Learning from Delayed Rewards*. PhD Thesis, King's College, Cambridge.
- Watkins, C.J., Dayan, P., 1992. Q-learning. *Machine Learning* 8 (3), 279–292.
- Williams, D.R., Barry, H., 1966. Counter conditioning in an operant conflict situation. *Journal of Comparative and Physiological Psychology* 61 (1), 154.
- Yachi, S., Higashi, M., 1998. The evolution of warning signals. *Nature* 394 (6696), 882–884.