

# Evolution of Reinforcement Learning in Uncertain Environments: A Simple Explanation for Complex Foraging Behaviors

Yael Niv<sup>1</sup>, Daphna Joel<sup>1</sup>, Isaac Meilijson<sup>2</sup>, Eytan Ruppin<sup>3</sup>

<sup>1</sup>*Department of Psychology, Tel Aviv University*

<sup>2</sup>*School of Mathematical Sciences, Tel Aviv University*

<sup>3</sup>*School of Computer Sciences & Sackler School of Medicine, Tel Aviv University*

Reinforcement learning is a fundamental process by which organisms learn to achieve goals from their interactions with the environment. Using evolutionary computation techniques we evolve (near-)optimal neuronal learning rules in a simple neural network model of reinforcement learning in bumblebees foraging for nectar. The resulting neural networks exhibit efficient reinforcement learning, allowing the bees to respond rapidly to changes in reward contingencies. The evolved synaptic plasticity dynamics give rise to varying exploration/exploitation levels and to the well-documented choice strategies of risk aversion and probability matching. Additionally, risk aversion is shown to emerge even when bees are evolved in a completely risk-less environment. In contrast to existing theories in economics and game theory, risk-averse behavior is shown to be a direct consequence of (near-)optimal reinforcement learning, without requiring additional assumptions such as the existence of a nonlinear subjective utility function for rewards. Our results are corroborated by a rigorous mathematical analysis, and their robustness in real-world situations is supported by experiments in a mobile robot. Thus we provide a biologically founded, parsimonious, and novel explanation for risk aversion and probability matching.

**Keywords** reinforcement learning · risk aversion · probability matching · evolutionary computation · dopamine · neuromodulation · heterosynaptic plasticity

## 1 Introduction

Reinforcement learning (RL) is a process by which organisms learn from their interactions with the environment to achieve a goal (Sutton & Barto, 1998). In RL, learning is contingent upon a scalar reinforcement signal that provides evaluative information about how good an action is in a certain situation, without providing an instructive supervising cue as to which would be the preferred behavior in the situation.

Behavioral research indicates that RL is a fundamental means by which experience changes behavior in both vertebrates and invertebrates, as most natural learning processes are conducted in the absence of an explicit supervisory stimulus (Donahoe & Packard-Dorsel, 1997). Several brain regions have been implicated in RL, including the midbrain dopaminergic neurons of the substantia nigra pars compacta (SNc) and the ventral tegmental area (VTA) in rats and primates, and their target areas in the basal ganglia

Correspondence to: Y. Niv, Department of Psychology, Tel Aviv University, Tel Aviv 69978, Israel.

E-mail: yaeln@cns.tau.ac.il

Tel: +972-3-6407864, Fax: +972-3-6409113

<http://www.cns.tau.ac.il/yaeln>

Copyright © 2002 International Society for Adaptive Behavior (2002), Vol 10(1): 5–24.

[1059–7123 (200201) 10:1; 5–24; 031999]

(e.g., Graybiel & Kimura, 1995; Houk & Wise, 1995; Schultz, 1998). A computational understanding of neuronal RL will enhance the understanding of learning processes in the brain and can contribute widely to the design of autonomous artificial learning agents.

RL has attracted ample attention in computational neuroscience, yet a fundamental question regarding the underlying mechanism has not been sufficiently addressed, namely, what are the optimal synaptic learning rules for maximizing reward in RL? In this article, we use evolutionary computation techniques to derive (near-)optimal neuronal learning rules that give rise to efficient RL in uncertain environments. We then investigate the behavioral strategies that emerge as a result of (near-)optimal RL.

RL has been demonstrated and studied extensively in foraging bees, thus we have chosen bee foraging as a model system for studying synaptic learning rules for RL. Real (1991, 1996) showed that when foraging for nectar in a field of blue and yellow artificial flowers, bumblebees exhibit efficient RL, rapidly switching their preference for flower type when reward contingencies were switched between the flowers. The bees also manifested risk-averse behavior: In a situation in which blue flowers contained 2  $\mu$ l sucrose solution, and yellow flowers contained 6  $\mu$ l sucrose in one-third of the flowers and zero in the rest, about 85% of the bees' visits were to the blue constant-rewarding flowers, although the mean return from both flower types was identical. Risk-sensitive and risk-averse choice behavior has also been demonstrated extensively in other animals (see Kacelnik & Bateson, 1996, for a review).

Bees foraging for nectar are faced with highly variable ecological conditions during the course of the year and in different habitats: Parameters such as the weather, the season, and competition all affect the availability of rewards from different kinds of flowers. In such an environment rapid learning is crucial for successful foraging, as the foraging individual cannot be prepared genetically for the ecological conditions of a particular habitat (Menzel & Muller, 1996). An uncertain environment also implies a "multi-armed bandit" type scenario, in which the bee collects food and information simultaneously (Greggers & Menzel, 1993). The foraging bee's choices are guided not only by the search for food but also by the search for information regarding the content of different food sources. This implies a trade-off between exploitation and

exploration (Wilson, 1996), as the bee's choices directly affect the "training examples" that it will encounter through the learning process. Thus in a multi-armed bandit situation a bee must devise a policy for choosing between exploiting available knowledge or exploring for more information, at every trial.

In a previous neural network (NN) model, Montague, Dayan, Person, and Sejnowski (1995) simulated bee foraging in a three-dimensional arena of blue and yellow flowers, based on a neurocontroller modeled after an identified interneuron (VUMmx1) in the honeybee suboesophageal ganglion (Hammer, 1993). This neuron's activity represents the reward value of gustatory stimuli, and similar to midbrain dopaminergic neurons of primates, it is activated by unpredicted rewards and by reward-predicting stimuli, and it is not activated by predicted rewards (Hammer, 1997). In their model, this neuron is modeled as a linear unit  $P$ , which receives visual information regarding changes in the percentages of yellow, blue, and neutral colors in the visual field and computes a prediction error. According to  $P$ 's output, the bee decides whether to continue flight in the same direction, or to change heading direction randomly. Upon landing, the reward value given to the network is not the nectar content of the chosen flower, but a transformation of this value according to the bee's subjective utility function for nectar (Harder & Real, 1987). At this time step, the synaptic weights of the network are updated according to a special anti-Hebbian-like learning rule in which the postsynaptic factor selects the direction of change (Montague, 1997). As a result, the values of the weights come to represent the expected subjective rewards from each flower type. The behavior of the decision unit  $P$  in this model is similar to the firing patterns characteristic of dopaminergic neurons in the SNc and VTA of primates, as have been recorded in a monkey performing a delayed-match-to-sample task (Montague, Dayan, & Sejnowski, 1996; Schultz, Dayan, & Montague, 1997), and is thus consistent with the hypothesis that these dopaminergic neurons compute a prediction error (Schultz, 2000).

Although this model replicates Real's foraging results and provides a basic and simple NN architecture to solve RL tasks, many aspects of the model, first and foremost the handcrafted synaptic learning rule, are arbitrarily specified and their optimality with respect to RL questionable. Toward this end, we use a generalized and parameterized version of this model to evolve optimal synaptic learning rules for RL (with

respect to maximizing nectar intake) using a genetic algorithm (Mitchell, 1997). Evolutionary computation techniques are especially suitable for RL problems, as they involve an artificial decision-making agent acting in an environment to achieve a goal. The solution is “evolved” in much the same way as the biological solutions have been evolved by natural selection—through a “parallel search” for effective solutions in a large population of individuals.

In contrast to common evolutionary computation applications involve NNs with evolvable synaptic weights or architectures (e.g., Ackley & Littman, 1991; Floreano & Mondada, 1996; Nolfi, Elman & Parisi, 1994), we set upon the task of evolving the network’s neuronal learning rules. Previous attempts at evolving neuronal learning rules have used heavily constrained network dynamics and very limited sets of learning rules, or evolved only a subset of the learning rule parameters: Chalmers (1990), in one of the first papers describing evolution of learning rules, evolved supervised learning rules for a fixed feed-forward architecture, using a parameterized version of the standard backpropagation learning rule up to second power, with a bit-encoded genome that severely limited the possible parameter values. Fontanari and Meir (1991) used Chalmers’ approach to evolve a learning algorithm for a single-layer network (perceptron) with binary weights, going up to the third power of the local rule. Baxter (1992) also evolved supervised learning rules, but for a binary fully connected NN with hidden units. Synapses could take the values of 0, 1, or (−1), and only nonzero synapses were modifiable. Baxter showed that only evolutions of networks with at least three hidden units converged, and in all successful evolutions the Hebbian learning rule was evolved. Unemi et al. (1994) evolved RL within a Q-learning framework but only evolved a subset of the learning rule parameters (learning rate, discount rate, and exploration rate), in a simple maze-exploration task. They demonstrated that learning ability emerges only if environmental change is faster than the alternating generations. Floreano and Mondada (1998) evolved learning rules using real robots in a framework that consisted of a fixed architecture in which each synapse could be driving or modulatory, excitatory or inhibitory, and could take one of four learning rules and one of four learning rates. In several recent papers, Floreano and Urzelai (e.g., Floreano & Urzelai, 2000, 2001) compared genetically determined synapses to

adaptive synapses, in a sequential “light switching” problem. Using a fully recurrent discrete-time NN with no hidden neurons that controls a real robot, they compared networks that use node-encoded learning rules (out of four possible Hebbian-based rules) to networks with genetically determined synapses and showed that on-line adaptation is advantageous for the studied task and allows for robustness to environmental changes.

In the present article, we define a general framework for evolving Hebbian learning rules, which essentially encompasses all heterosynaptic and monosynaptic Hebbian learning rules and also allows for complex neuromodulatory interactions of gating of synaptic plasticity. Via the genetic algorithm we select bees based solely on their nectar-gathering ability in a changing environment. The uncertainty of the environment ensures that efficient foraging can only be a result of learning throughout a bee’s lifetime, thus promoting the evolution of efficient learning rules.

To avoid the possible confusion of terms, we make a distinction between the notions of heterosynaptic learning (Dittman & Regehr, 1997; Schacher, Wu & Sun, 1997; Vogt & Nicoll, 1999) and neuromodulation of plasticity (Bailey, Giustetto, Huang, Hawkins, & Kandel, 2000; Fellous & Linster, 1998). The classic monosynaptic Hebbian learning rule is an activity-dependent learning rule in which a synapse is updated only when there is activity both in the presynaptic and in the postsynaptic neurons. In contrast, heterosynaptic Hebbian learning allows for a less-restricted modification of synapses, such that a synapse can also be updated when only the presynaptic or the postsynaptic component has been active, and more generally, even when neither have been active. We term this rule “heterosynaptic” modification as it allows for the firing of a neuron to affect all its synapses, regardless of the activity of the other neurons connected to it. Neuromodulation of synaptic plasticity further enhances the learning rule by allowing a three-factor interaction in the learning process: Through neuromodulation the activity of a neuron can gate the plasticity of a synapse between two other neurons. Providing the neuromodulatory neuron has fired, the synapse can be updated (according to the heterosynaptic learning rule pertaining to the pre- and postsynaptic neurons). Otherwise, plasticity in the synapse is shut off. Both heterosynaptic plasticity and neuromodulatory gating of synaptic plasticity have been

demonstrated in neural tissues (Bailey et al., 2000; Dittman & Regehr, 1997; Fellous & Linster, 1998; Schacher et al., 1997; Vogt & Nicoll, 1999) and have been recognized to increase the computational complexity of synaptic learning (Bailey et al., 2000; Fellous & Linster, 1998; Wickens & Kötter, 1995). By allowing for both heterosynaptic learning and neuro-modulation of plasticity, we define a very large search space in which the genetic algorithm can search for optimal synaptic learning rules.

In the following section we describe the model and the evolutionary dynamics. Section 3 reports the results of our simulations: In Section 3.1 we describe the successful evolution of RL. Section 3.2 describes the evolved synaptic update rules, and their influence on the exploration/exploitation trade-off of the evolved bees. In Section 3.3 we analyze the foraging behavior resulting from the learning dynamics and find that when tested in a new environment, our evolved bees manifest risk-averse behavior. Although this choice strategy was not selected for, we rigorously prove that risk aversion emerges directly from RL, in contrast to the conventional explanations of risk-averse behavior that rely on the existence of subjective utility functions. The strength of the evolved learning mechanism is further demonstrated in Section 3.4, which describes the emergence of probability matching behavior. This behavior, which was previously thought to result from competition for food resources, is shown to emerge in a noncompetitive scenario as a result of the learning rule dynamics alone. Section 3.5 describes a minirobot implementation of the model, aimed at assessing its robustness. We conclude with a discussion of our results in Section 4.

## 2 The Model

A simulated bee-agent flies in a three-dimensional arena, over a flower patch composed of  $60 \times 60$  randomly scattered yellow and blue squares representing two types of flowers. A bee's life consists of 100 trials. Each trial begins with the bee placed in a random location above the flower patch and with a random heading direction. The bee starts its descent from a random height of 8–9 units above the flower patch and advances in steps of 1 unit that can be taken in any downward direction ( $360^\circ$  horizontal,  $90^\circ$  vertical). The bee views the world through a cyclopean eye

( $10^\circ$  cone view), and in each time step it decides whether to maintain the current heading direction or to reorient randomly, based on the visual input (Figure 1a). Upon landing (the field has no boundaries, and the bee can land on a flower or outside the flower patch, on neutral-colored ground), the bee consumes any available nectar in one time step and another trial begins. The evolutionary goal (the fitness criterion) is to maximize nectar intake.

In the neural network controlling the bee's flight (Figure 1b), which is an extension of Montague et al.'s (1995) network, three modules ("regular," "differential" and "reward") contribute their input via synaptic weights, to a linear neuron  $P$ . The regular input module reports the percentage of the bee's field of view filled with yellow  $[X_y(t)]$ , blue  $[X_b(t)]$ , and neutral  $[X_n(t)]$ . The differential input module reports temporal differences of these percentages  $[X_i(t) - X_i(t-1)]$ . The reward module reports the actual amount of nectar received from a flower  $[R(t)]$  in the nectar-consuming time step (in this time step it is also assumed that there is no new input  $[X_i(t) = 0]$ ), and zero during flight. Note that, in contrast to Montague et al. (1995), we do not incorporate any form of nonlinear utility function with respect to the reward. Thus,  $P$ 's continuous-valued output is

$$P(t) = R(t) + \sum_{i \in \text{regular}} W_i X_i(t) + \sum_{j \in \text{differential}} W_j [X_j(t) - X_j(t-1)] \quad (1)$$

The bee's action is determined according to the output  $P(t)$  using Montague et al.'s (1995) probabilistic action function (Figure 1c):

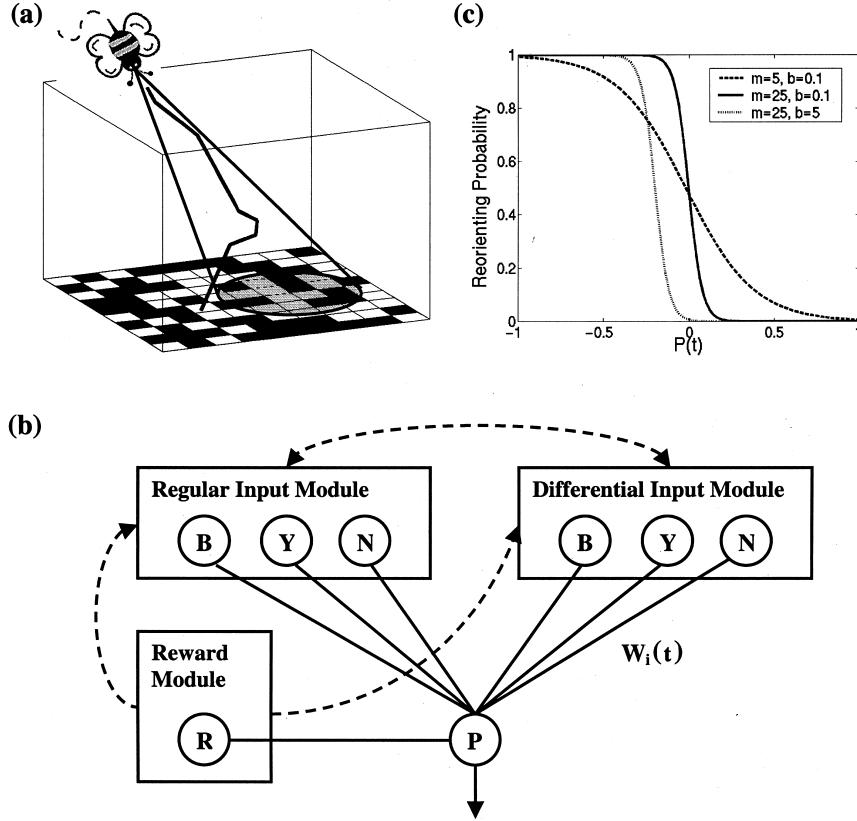
$$p(\text{change direction}) = \frac{1}{1 + \exp[m \cdot P(t) + b]} \quad (2)$$

where  $m$  and  $b$  are real-valued evolvable parameters.

During the bee's "lifetime" the synaptic weights of the regular and differential modules are modified via a heterosynaptic Hebb learning rule of the form

$$\Delta W_i(t) = \begin{cases} \eta [A \cdot V_i(t) P(t) + B \cdot V_i(t) + C \cdot P(t) + D] & \text{dependencies met} \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

where  $\eta$  is a global learning rate parameter,  $V_i(t)$  and  $P(t)$  are the presynaptic and the postsynaptic



**Figure 1** (a) The simulated setting. Depicted is a portion of the arena with black and white squares representing the yellow and blue flowers, the bee's field of view when starting flight, and a typical trajectory. (b) The bee's neural network controller. The weights  $W_i(t)$  of the regular and differential modules are modifiable. (c) The bee's action function. Probability of reorienting direction of flight as a function of  $P(t)$  for different values of parameters  $m, b$ .

values respectively,  $W_i$  is their connection weight, and  $A - D$  are real-valued evolvable parameters. In addition, learning in one module can be dependent on another module (dashed arrows in Figure 1b), such that if module  $M$  depends on module  $N$ ,  $M$ 's synaptic weights will be updated according to Equation 3 only if module  $N$ 's neurons have fired, and if it is not dependent, the weights will be updated on every time step. In this respect, a dependency on the reward module is satisfied when the reward neuron fires, and dependencies on the regular and differential modules are satisfied neuron-wise; that is, when a neuron fires, the synapses connected to the respective (same color) neurons in the dependent module can be updated. Synapses of a module dependent on two other modules can only be updated when satisfying both dependency conditions. Thus the bee's "brain" is capable of a nontrivial neuromodulatory gating of synaptic plasticity.

The simulated bee's genome consists of a string of 28 genes, each representing a parameter governing the network architecture or learning dynamics. Fifteen genes specify the bee's brain at time of "birth" (before the first trial): seven Boolean genes determine whether each synapse exists or not; six real-valued genes (range  $[-1, 1]$ ) specify the initial weights of the regular and differential module synapses (the synaptic weight of the reward module is clamped to 1, effectively scaling the other network weights); and two real-valued genes specify the action-function parameters  $m$  and  $b$  (initialized in ranges  $[5, 45]$  and  $[0, 5]$  respectively). Thirteen remaining genes specify the learning dynamics of the network: The regular and differential modules each have a different learning rule specified by four real-valued genes (parameters  $A - D$  of Equation 3, initialized in range  $[-0.2, 0.2]$ ); The global learning rate of the network  $\eta$  is specified by a real-valued gene (initialized in range  $[0, 1]$ ); and four

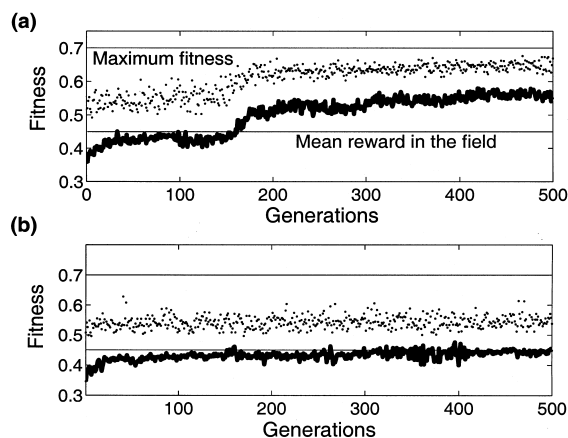
Boolean genes specify dependencies of the visual input modules on each of the other two modules. Apart from the synaptic weights, which are bound, parameter values are unrestricted.

Gene values were optimized using a genetic algorithm. A first generation of bees was produced by randomly generating 100 genome strings. Each bee performed 100 trials independently (no competition) and received a fitness score according to the average amount of nectar gathered per trial. To form the next generation, 50 pairs of parents were chosen (with returns) with a bee's fitness specifying the probability of it being chosen as a parent. Each two parents gave birth to two offspring, which inherited their parents' genome after recombination was performed and random mutations added. Mutations were performed by adding a uniformly distributed value in the range of  $[-0.1, 0.1]$  to real-valued genes, and reversing of Boolean genes. Mutation rate was high for the first generations (16% for real-valued genes and 3.2% for Boolean genes), gradually decaying to a lower value (four-fold decay in real-valued genes to 4%, and eight-fold decay in Boolean genes to 0.4%, the mutation rate decayed linearly in four steps, every 100 generations). The mutation rate for Boolean genes was chosen to be considerably smaller than that of real-valued genes, as a mutation on a real-valued gene only resulted in a small perturbation of the gene value, whereas a mutation on a Boolean gene completely reversed it. Recombination was performed via a uniform crossover of the genes ( $p = 0.25$ ). It is important to note that there was no Lamarckian inheritance—learned weights were not passed on to offspring. One hundred offspring were thus created and once again tested in the flower field. This process continued for 500 generations.

### 3 Results

#### 3.1 Evolution of Reinforcement Learning

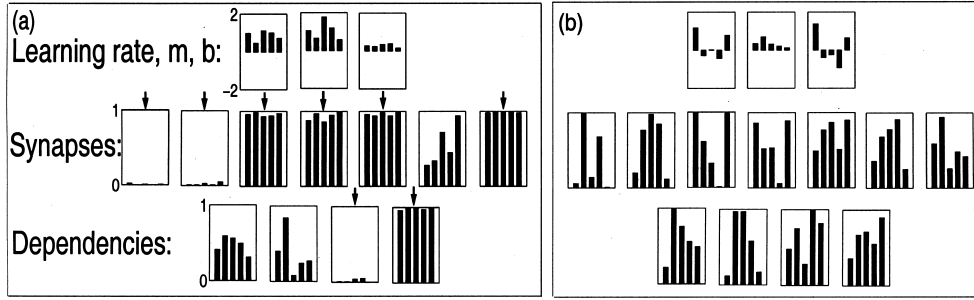
To promote the evolution of learning, bees were evolved in an “uncertain” world: In each generation one of the two flower types was randomly assigned as a constant-yielding high-mean flower (containing  $0.7 \mu\text{l}$  nectar), and the other a variable-yielding low-mean flower ( $1 \mu\text{l}$  nectar in one-fifth of the flowers and zero otherwise). The reward contingencies were switched between the



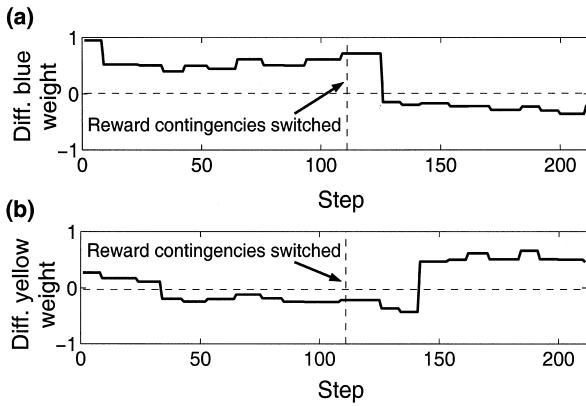
**Figure 2** Typical fitness scores of a successful run (a) and an unsuccessful run (b) of 500 generations. Solid line: mean fitness; dots line: maximum fitness in each generation.

two flower types in a randomly chosen trial during the second or third quarter of each bee's life.

Ten separate evolutionary runs were performed, of which half were successful. The different runs show one of two types of fitness curves: Successful runs, defined as runs in which reward-dependent choice behavior is successfully evolved, are characterized by two distinct evolutionary jumps (Figure 2a). Unsuccessful runs, which produce behavior that is not dependent on rewards, show only the first jump (Figure 2b). This first jump is due to the acquisition of a negative neutral synapse responsible for a “do not land on neutral-colored ground” rule. The second evolutionary jump characteristic of successful runs is due to the almost simultaneous evolution of eight genes governing the network structure and learning dependencies, which are essential for producing efficient learning in the bees' uncertain environment: All successful networks have a specific architecture that includes only four synapses (the regular neutral, differential blue and differential yellow, and the reward synapse), as well as a dependency of the differential module on the reward module, conditioning modification of these synapses on the presence of zreward (Figure 3a). Agents that have almost all the crucial alleles, but are missing one or two, are nevertheless unsuccessful (Figure 3b). Thus we find that in our framework, only a network architecture similar to that used by Montague et al. (1995) can produce



**Figure 3** Mean value of several genes in the last generation of (a) five successful and (b) five unsuccessful runs. Each subfigure shows the mean value of one gene in the last generation of each of five runs. Genes shown (from left to right): Top row—the learning rate gene and the two action function parameters  $m$  and  $b$ . Middle row—the Boolean genes governing the existence of the different synapses: regular blue, regular yellow, and regular neutral input synapses, differential blue, differential yellow, and differential neutral input synapses, and the reward input synapse. Bottom row—Boolean genes determining the dependencies of the regular module on the differential module and on the reward module, and the dependencies of the differential module on the regular module and the reward module. Genes crucial for successful reinforcement learning are marked with an arrow.



**Figure 4** Synaptic weight values of an evolved bee performing 20 test trials. Synaptic weight values for the differential blue (a) and yellow (b) synapses were recorded in a bee from the last generation of a successful evolutionary run, during each of the approximately 210 steps needed to complete 20 test trials. In each trial, the bee tended to choose the color associated with the larger weight. Upon landing and receiving reward, the weights changed according to the evolved learning rules. Blue was the initial constant-rewarding flower ( $0.7 \mu\text{l}$ ) and yellow the variably rewarding flower ( $1 \mu\text{l}$  in one-fifth of the flowers). Reward contingencies were switched after trial 10. The weight values can be seen to follow approximately the rewards expected from each flower type.

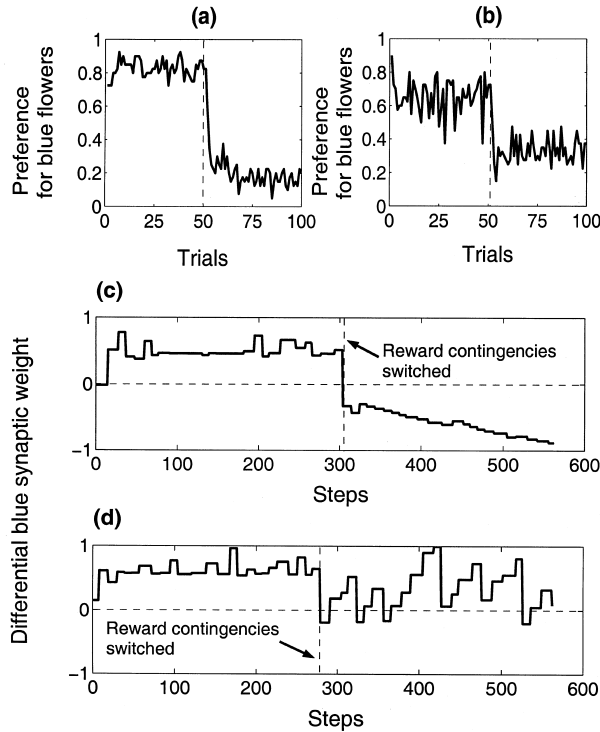
above-random foraging behavior, supporting their choice as an optimal one. However, our evolved networks utilize a family of heterosynaptic learning

rules different from the monosynaptic rule used by Montague et al., giving rise to several important behavioral strategies.

### 3.2 Exploration/Exploitation Trade-off and Heterosynaptic Modification

To understand the evolved learning rule, we examined the foraging behavior of individual bees from the last generation of successful runs. In general, the bees manifest efficient reinforcement learning, showing a marked preference for the high-mean-rewarding flower, with a rapid transition of preferences after the reward contingencies are switched between the flower types. The values of the synaptic weights are also indicative of learning based on reward contingencies, as they follow the rewards expected from each flower (Figure 4).

A more detailed inspection of the behavior of bees from different evolutionary runs reveals that the bees differ in their degree of exploitation of the high-rewarding flowers versus exploration of the other flowers (Figure 5a, b). These individual differences in the foraging strategies employed by the bees, result from an interesting relationship between the micro-level Hebb rule coefficients and the exploration/exploitation trade-off characteristic of the macro-level behavior. According to the dependencies evolved in all the successful evolutionary runs, learning (synaptic



**Figure 5** (a, b) Preference for blue flowers for two different bees from the last generation of two different successful runs, averaged over 40 test bouts, each consisting of 100 trials. Blue was the initial constant-rewarding high-mean flower. Reward contingencies were switched between flower types after trial 50. The bees shown represent extreme cases of evolved “exploitation-inclined” (a) and “exploration-inclined” (b) behavior. (c, d) Fluctuations of the differential blue synaptic weight during a typical flight of the “exploiting” bee (c) and the “exploring” bee (d). Synaptic weight values for the differential blue synapse were recorded while the bees performed 50 test trials (about 600 flight steps). Blue was the initial constant-rewarding flower; reward contingencies were switched after trial 25. The different fluctuation patterns in the last 25 trials result from the different heterosynaptic Hebb rule coefficients. Hebb rule coefficients for the “exploiting” bee (c) were  $A = -0.82$ ,  $B = 0.15$ ,  $C = 0.24$ ,  $D = -0.04$  and for the “exploring” bee (d) were  $A = -0.92$ ,  $B = 0.39$ ,  $C = 0.16$ ,  $D = 0.25$ .

updating) occurs only upon landing, and we can analyze the evolved heterosynaptic learning rule of the differential module as follows: In the common case, upon landing the bee sees only one color, thus all inputs are zero except the differential input corresponding to the color of the chosen flower, which is, in

the absence of new visual input in the landing step,  $X_i(t) - X_i(t-1) = 0 - 1 = -1$ . Thus the output of  $P$  in this step is

$$P(t) = R(t) + (-1) \cdot W_{\text{chosen}}(t) = R(t) - W_{\text{chosen}}(t) \quad (4)$$

Therefore, the synaptic update rule for the differential synapse corresponding to the chosen flower color is

$$\Delta W_{\text{chosen}}(t+1) = \eta[(A - C) \cdot (-1) (R(t) - W_{\text{chosen}}(t)) + (D - B)] \quad (5)$$

leading to an effective monosynaptic coefficient of  $(A - C)$ , and a general weight decay coefficient  $(D - B)$ . For the other differential synapses the synaptic update rule is

$$\Delta W_j(t+1) = \eta[C \cdot (R(t) - W_{\text{chosen}}(t)) + D] \quad (6)$$

Table 1 summarizes the values of the coefficients of the heterosynaptic learning rule of the differential module, as they were evolved in the five successful evolutionary runs. The synaptic weight of the chosen color is affected by the values of  $(A - C)$  and  $(D - B)$ . The large negative values of  $(A - C)$  found specify an anti-Hebbian learning rule (i.e., an anti-correlation learning rule) for the chosen flower. This is a result of the negative presynaptic value while landing. The values of  $D$  and  $B$  further modulate the synaptic strength: In trials in which no prediction error is encountered (i.e., the postsynaptic value is near zero in the learning step), a negative value of  $(D - B)$  results in weakening of the synapse, and a positive value results in strengthening of the synapse. The synaptic weights corresponding to colors that were not chosen are only influenced heterosynaptically, by the values of  $C$  and  $D$ . A positive value of  $C$  results in a global strengthening of these synapses when the postsynaptic value is positive and a global weakening when the postsynaptic value is negative, that is, a generalization of the “good/bad surprise” that was encountered, to the other colors. A negative value of  $C$  will produce the opposite effect. A positive or negative value of  $D$  results in a global strengthening or decay of all the synapses at every learning step.



**Table 1** Heterosynaptic learning rule coefficients in five successful evolutionary runs. Columns 2–5: Mean (standard deviation) of the evolved learning rule coefficients of the differential module heterosynaptic learning rule, in the last generation of five successful runs. Columns 6–9: Number of bees in the last generation with a positive value of  $B$ ,  $C$ ,  $D$  and  $(D - B)$ , respectively.

Run	$A$	$B$	$C$	$D$	$(B > 0)$	$(C > 0)$	$(D > 0)$	$(D - B > 0)$
1	-1.2 (0.1)	0.1 (0.1)	-0.1 (0.1)	-0.1 (0.1)	73	24	1	15
2	-1.2 (0.2)	-0.2 (0.1)	0.2 (0.1)	-0.1 (0.2)	0	100	47	61
3	-1.7 (0.2)	0.3 (0.1)	-0.3 (0.1)	0.0 (0.1)	94	0	82	7
4	-1.3 (0.1)	0.2 (0.1)	0.3 (0.1)	-0.2 (0.1)	98	100	1	0
5	-1.0 (0.1)	0.0 (0.3)	0.3 (0.1)	-0.5 (0.1)	66	100	0	2
Mean	-1.3 (0.3)	0.1 (0.2)	0.1 (0.3)	-0.2 (0.2)	66	65	26	17

Thus the family of heterosynaptic learning rules evolved cover a large range of synaptic dynamics that affect all the synapses of the differential module at every learning step and influence the bee's exploration/exploitation trade-off (Figure 5c, d). A positive value of  $D$  results in "spontaneous" strengthening of competing synapses, leading to an exploration-inclined bee (Figure 5d). A positive  $C$  value further enhances this behavior. Conversely, negative values of  $C$  and  $D$  will result in a declining tendency to visit competing flower types, leading to exploitation-inclined behavior (Figure 5c).

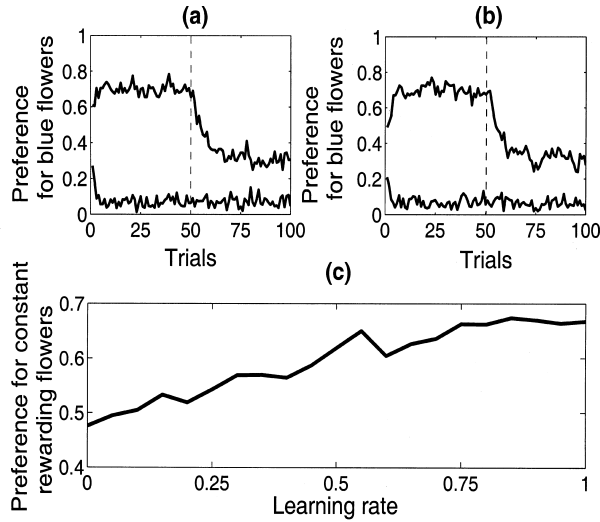
### 3.3 Emergence of Risk Aversion

A prominent strategy exhibited by the evolved bees is risk aversion. Figure 6a shows the choice behavior of 200 previously evolved bees (40 bees from the last generation of each of the five successful runs), tested in a new environment in which the mean rewards of the two kinds of flowers are identical, but their variance is different. Although the situation does not call for any flower preference, as on average the bee would receive the same nectar amount regardless of its choice strategy, the bees consistently prefer the constantly rewarding flower over the higher-content but variably rewarding flower.

Theoretically, this preference could have been a result of the fact that during evolution, the constantly rewarding flower was also the high-mean-rewarding flower, so that the bees would come to prefer constantly over variably rewarding flowers. However, this does not seem to be the case, as we have found that risk-averse behavior can also be evolved without

introducing any risk to the environment during the evolutionary process: Bees evolved in an environment containing two constantly rewarding flowers yielding different amounts of nectar also exhibit risk-averse behavior when tested in a variably rewarding flower scenario. Figure 6b shows the choice preferences of five new populations of bees evolved in an environment that contained two constantly rewarding flower types, one yielding 0.8  $\mu$ l nectar, and the other 0.3  $\mu$ l. Although risk-less, this environment was also an uncertain one, in that the flower color that yielded the higher reward was randomly chosen in each generation, and the reward contingencies were switched between the two flower types in a random time step within the second or third quarter of the bees' lifetime, similar to the previously described evolutionary runs. After completing 500 generations of evolution, 40 bees from the last generation of each of five successful evolutionary runs were tested in the risk-aversion scenario in which both flowers yielded the same mean reward, but one was variably rewarding. Confronted with variably rewarding flowers for the first time, the bees showed a pronounced preference for the constantly rewarding flower type (Figure 6b). Thus we can observe emergent risk-averse behavior in bees evolved in a risk-less environment. Below we prove the emergence of risk-aversion from RL analytically, relying only on the assumption of Hebbian learning. This verifies that the emergence of risk aversion is indeed not dependent on the amount of risk or the quality of the risky versus the nonrisky flower during the evolutionary process.<sup>1</sup>

Risk aversion has been studied extensively in the fields of economics and game theory. The risk-averse



**Figure 6** (a) Risk aversion. Preference for blue flowers in 100 test trials averaged over 200 previously evolved bees (40 bees from the last generation of each of the five successful runs), now tested in conditions different from those in which they were evolved. Although both flower types yield the same mean reward (blue: 0.5  $\mu$ l nectar, yellow: 1  $\mu$ l in half the flowers, contingencies switched after trial 50), the mean (top) and standard deviation (bottom) of the proportion of bees in each evolutionary run who preferred blue flowers in each trial, reveal a marked preference for the constant-yielding flower. (b) Risk aversion in bees evolved in a risk-less environment. Preference for blue flowers in 100 test trials averaged over 200 bees (40 bees from each of five successful evolutionary runs) evolved in a risk-less environment that contained two constantly-rewarding flower types. As in (a), the mean (top) and standard deviation (bottom) of the proportion of bees in each evolutionary run who preferred blue flowers in each trial, show that risk aversion is prominent, although the two flower types yield the same mean reward. (c) Risk aversion is ordered by learning rate. An illustration of the general principle of the dependency of risk aversion on learning rate, proven mathematically (see Appendix). Each point represents the percentage of visits to constant-rewarding flowers in 50 test trials averaged over 40 previously evolved bees (all from one successful evolutionary run), tested with a clamped learning rate.

behavior that has been observed in many choice scenarios in animals as well as in humans has traditionally been accounted for by hypothesizing the existence of a nonlinear concave “utility function” for reward. In bees, for instance, such a subjective utility function for nectar can result from a concave relationship between nectar volume and net energy intake, between net

energy intake and fitness, or between the actual and perceived nectar volume (Harder & Real, 1987; Smallwood, 1996). Montague et al. (1995) incorporate this explanation into their model to reproduce Real’s (1991) risk-aversion results, by directly applying a nonlinear utility function to the nectar content and feeding its result to the reward module. In contradistinction to this conventional explanation of risk aversion, our model does not include any form of nonlinear utility for reward. What then brings about risk-averse behavior in our model? Corroborating previous numerical results (March, 1996), we prove analytically that this foraging strategy is a direct consequence of Hebbian learning dynamics in a two-armed banditlike RL situation.

During a trial, a bee makes a series of choices regarding its flight direction, to choose which flower to land on. As the bee does not learn (i.e., there is no synaptic plasticity) during flight, all the choices throughout one trial are influenced by the same weight values. Thus under a certain rewarding regime, that is, in between changes in reward contingencies, the bee’s stochastic foraging decisions can be modeled as choices between a variably rewarding ( $v$ ) and a constantly rewarding ( $c$ ) flower, based on synaptic weights  $W_v$  and  $W_c$ .

For simplicity, we examine the case of simple monosynaptic anti-Hebbian learning: In this case, the synaptic update rule is the well-known temporal difference (TD) learning rule (Sutton & Barto, 1998)

$$\Delta W(t) = \eta(R(t) - W(t-1)). \quad (7)$$

The synaptic weights are in effect a “memory” mechanism, as they are a function of the rewards previously obtained from each of the two flower types.  $W_v$ , representing the reward expected from the variable flower, is thus an exponentially weighted average of  $[v_1, v_2, \dots]$ , the previous rewards obtained from ( $v$ ):

$$W_v = W_v(\eta) = \eta(v_t + (1 - \eta)v_{t-1} + (1 - \eta)^2 v_{t-2} + \dots) \quad (8)$$

$W_c$ , as an exponentially weighted average of rewards obtained from the constantly rewarding flower type, is constant.

Bearing this notion in mind, the Appendix provides a mathematical proof of the emergence of risk aversion. We compute  $f_v$ , the frequency of visits to

variably rewarding flowers, and show that it is a function of  $p_v(W_v)$ , defined as the probability of choosing ( $v$ ) in a trial in which the synaptic weight corresponding to the variably rewarding flower is  $W_v$ . We prove that  $W_v$ , as a function of the learning rate, is risk ordered, such that for higher learning rates  $W_v(\eta)$  is riskier than for lower learning rates. We then use the mathematical definition of riskiness to show that under relatively mild assumptions regarding the bee's choice function, the frequency of visits to variably rewarding flowers is lower for higher learning rates than for lower learning rates and is ordered by learning rate. Finally we show that the risk order property of  $W_v(\eta)$  always implies risk-averse behavior.<sup>2</sup> That is, for every learning rate, the frequency of visits to the variable flower ( $f_v$ ) is less than 50%, further decreasing under higher learning rates. Our simulations corroborate these analytical results (Figure 6c).

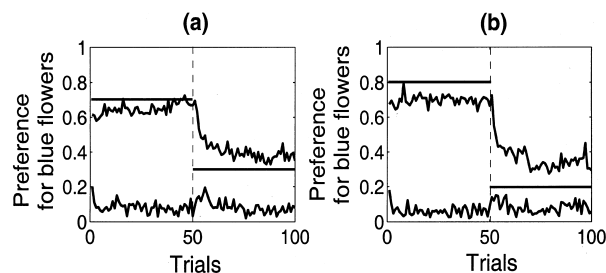
Intuitive insight into these results can be gained by observing that in the learning process, the bee makes its decisions based on finite time windows and does not compute the long-term mean reward obtained from each flower. This is even more pronounced with high learning rates such as those evolved ( $\sim 0.8$ ). With such a learning rate, after landing on an empty flower of the variably rewarding type, the bee updates the reward expectation from this flower type (i.e., updates the corresponding synaptic weight according to the evolved heterosynaptic Hebb update rule) to near zero, and as a result, prefers the constantly rewarding flower, from which it constantly expects (and receives) a reward of  $0.5 \mu\text{l}$ . As long as the bee chooses the constantly rewarding flower, it will not update the expectation from the variably rewarding flower, which will remain near zero. Even after an occasional "exploration" trial in which a visit to the variable flower yields a high reward, the preference for this flower will be short lived, lasting only until the next unrewarded visit. Note that rapid learning such as has been evolved here is essential for obtaining high fitness in a highly variable environment (Menzel & Muller, 1996), and such abnormally high learning rates have been hypothesized by Real (1991), and were also used in Montague et al.'s (1995) model. Nevertheless, the above mathematical analysis shows that even with low learning rates, as long as the bee is a reinforcement-learning bee (i.e., its learning rate greater than zero), it will manifest risk-averse behavior.

### 3.4 Emergence of Probability-Matching Behavior

Another notable strategy by which bumblebees (and other animals) optimize choice in multiarmed bandit situations is probability matching. "Probability matching" refers to the phenomenon observed in situations in which the different alternatives offer similar rewards, but with different probabilities (Bitterman, 1965).<sup>3</sup> In such cases, probability matching predicts that the different alternatives will be chosen according to the ratio of their reward probabilities. Probability matching has been shown to describe the behavior of some animals (e.g., Bitterman, 1965), but not others (e.g., Herrnstein & Loveland, 1975). Keasar, Rashkovich, Cohen, and Shimida (in press) have shown that when faced with variably rewarding flowers offering the same rewards with different probabilities, bees match the choices of the different flower types to their reward ratio, in accordance with probability matching (see also Greggers & Menzel, 1993).

This seemingly "irrational" behavior with respect to optimization of reward intake (Herrnstein & Loveland, 1975; Herrnstein, 1997) was explained as an evolutionarily stable strategy (ESS) for the individual forager, when faced with competitors (Thuijsman, Peleg, Amitai, & Shmida, 1995). In a multi-animal competitive setting, probability matching produces an ideal free distribution (IFD) in which the average intake of food is the same at all food sources, and no animal can improve its payoff by feeding at another source. Using evolutionary computation techniques, Seth (1999) evolved battery-driven agents that competed for two different battery refill sources and showed that indeed matching behavior emerges in a multi-agent scenario, whereas when evolved in isolation, agents choose only the high-probability refill source.

Surprisingly, our evolved bees also demonstrate probability matching behavior. Figure 7 shows the performance of 200 previously evolved bees (40 from the last generation of each of five successful evolutions), when tested in probability-matching experiments in which all flowers yield  $1 \mu\text{l}$  nectar, but with different reward probabilities. In both conditions, the bees show near-matching behavior, preferring the high-probability flower to the low-probability one, by a ratio that closely matches the reward probability ratios. In contrast to the "overmatching" described by Kaesar et al. (in press), our simulated bees exhibit



**Figure 7** Probability matching. Preference for blue flowers in 100 test trials averaged over 200 previously evolved bees (40 bees from the last generation of each of the five successful runs), now tested in probability-matching conditions. All flowers yielded 1  $\mu$ l nectar but with different reward probabilities. Reward probabilities for blue and yellow flowers, respectively, were (a) 0.8, 0.4 and (b) 0.8, 0.2 (contingencies switched after trial 50). Mean (top) and standard deviation (bottom) of the proportion of bees in each evolutionary run who preferred blue flowers in each trial. Horizontal lines: behavior predicted by perfect probability matching [i.e., a 2:1 choice ratio in (a) and a 4:1 choice ratio in (b)].

the more commonly found (Domjan, 1998) “undermatching,” as they prefer the higher-probability flower slightly less than predicted by perfect probability matching. In our simulations, this emergent behavior is not a result of competitive conditions (as the bees were evolved and tested in isolation), but rather a direct result of the evolved reinforcement learning dynamics: Due to the high learning rate, the fluctuating weights representing the expected yield from each flower will essentially move back and forth from zero to one. When both are zero, the two flowers are chosen randomly, but the high-yielding flower has a greater chance of yielding reward, after which its weight will be updated to 1, and this flower is preferred to the other. When both weights are 1, the less-yielding flower has a greater chance of having its weight updated to zero, again resulting in preference for the high-yielding flower. Thus, as an alternative to previous accounts, probability-matching behavior can be evolved in a noncompetitive setting, again as a direct consequence of (near-)optimal RL.

### 3.5 Robot Implementation

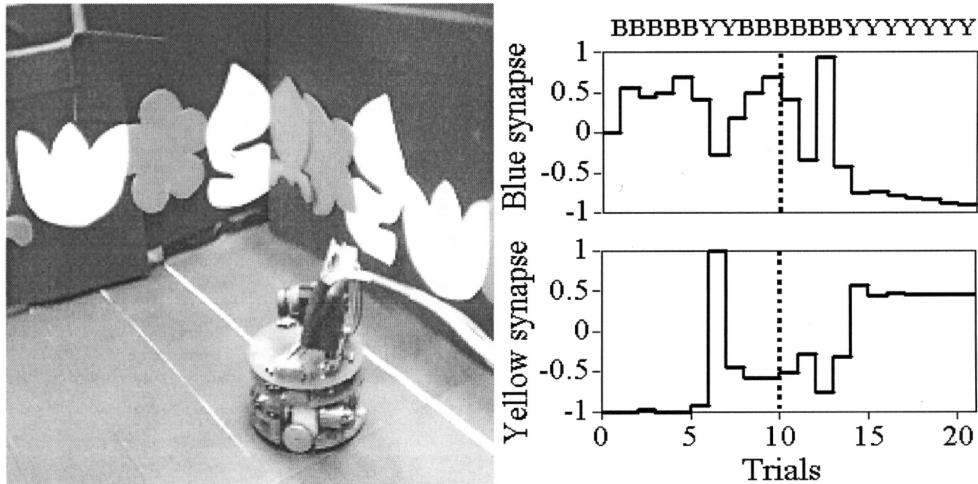
To assess the robustness of the evolved RL algorithm, we implemented it in a mobile minirobot by letting the robot’s actions be governed by a NN controller

similar to that evolved in successful bees, and by having its synaptic learning dynamics follow the previously evolved RL rules. A Khepera minirobot foraged in a  $70 \times 35$  cm arena whose walls were lined with flowers, viewing the arena via a low-resolution CCD camera ( $200 \times 200$  pixels), moving at a constant velocity and performing turns according to the action function (Equation 2) to choose flowers, in a manner completely analogous to that of the simulated bees. The NN controller was identical to that evolved for the simulated bees, except that it received no “neutral” inputs. All calculations were performed in real time on a Pentium-III 800 Mhz computer (256 MB RAM) in tether mode. Moving with continuous speed and performing all calculations in real time, the foraging robot exhibited rapid reinforcement learning and risk-averse behavior, analogous to that of the simulated bees (Figure 8). Thus the algorithms and behaviors evolved in the virtual bees’ simulated environment using discrete time steps hold also in the different and noisy environment of real foraging minirobots operating in continuous time.

## 4 Discussion

The interplay between learning and evolution has been previously investigated in the field of evolutionary computation. Much of this research has been directed toward elucidating the relationship between evolving traits (such as synaptic weights) versus learning them (e.g. Ackley & Littman, 1991; Hinton & Nowlan, 1987). A relatively small amount of research has been devoted to the evolution of the learning process itself, most of which was constrained to choosing the appropriate learning rule from a limited set of predefined rules (Baxter, 1992; Chalmers, 1990; Floreano & Mondada, 1996). In this work we show for the first time that (near-)optimal learning rules for RL in a general class of multi-armed bandit situations can be evolved in a general Hebbian learning framework. The evolved learning rules are by no means trivial, as they are heterosynaptic and employ synaptic plasticity modulation.

As a result of the evolved learning rules, several complex foraging behaviors emerge, demonstrating the strength of evolutionary computation as a neuroscience research methodology that links together phenomena on the neuronal and behavioral levels. We



**Figure 8** Synaptic weights of a mobile robot incorporating a neural network controller of one of the previously evolved bees, performing 20 foraging trials (blue flowers: 0.5  $\mu$ l nectar, yellow: 1  $\mu$ l in half the flowers, contingencies switched after trial 10). (Left) The foraging robot. (Right) Blue and yellow weights in the differential module represent the rewards expected from the two flower colors along the trials. Top: Flower color chosen in each trial.

show that in our model the fundamental macro-level strategies of risk aversion and probability matching are a direct result of the micro-level synaptic learning dynamics, which also control the trade-off between exploration and exploitation. These behavioral strategies have not been explicitly evolved but emerge in the model as “side-effects” of RL, making additional assumptions conventionally used to explain these behaviors unnecessary. Furthermore, in this work we not only show that risk-averse behavior can exist in the absence of a nonlinear utility function for reward, but also that this behavior can emerge spontaneously in a risk-free environment. Risk aversion is shown to be a direct consequence of adaptation to a changing, uncertain world that induces a form of learning in which the estimation of expected rewards is biased, as it takes the more recent sampling into account more strongly. When acting in an environment in which intraflower risk is added to the interflower uncertainty, this bias is expressed as risk-averse behavior. This result is important not only to the fields of evolutionary computation and animal learning theories, but also to the fields of economics and game theory, as it provides a new perspective on the well-studied paradigm of risk-sensitive behavior.

The prevailing accounts of risk aversion in the fields of economics and game theory are based on the

notion of nonlinear utility functions, which were introduced into rational choice theory models to explain the apparent divergence from rationality embedded in risk-sensitive behavior. Utility functions were later assumed to result from biological/energetic considerations, an explanation that usually involves postulating and quantifying a large number of parameters (e.g., see Harder and Real, 1987, for the energetic derivation of the utility function for nectar postulated for bumblebees). It should be stressed that, although widely accepted as an explanation for risk aversion, in traditional economics, utility functions are inferred entities (Herrnstein, 1997), and there is no direct evidence for the existence of subjective utility functions other than the behavioral manifestation of risk-sensitive behavior itself. The existence of the assumed underlying subjective utility function and the shape of the function are assessed not by an independent measure of subjective utility, but according to the resulting risk-averse behavior. In contrast, in our model, the well-documented phenomenon of risk aversion, characteristic of choice behavior in many species, emerges directly from the learning dynamics of a simple NN model based on an architecture identified in brains of bumblebees. We have shown empirically and theoretically that a finite memory, which is a direct consequence of Hebbian learning, is sufficient

to account for risk aversion in any model of temporal difference learning that involves a learning rate greater than zero. As nonlinear utility functions are abstract concepts that have not been proven to exist biologically, omitting them from the model has the effect of producing a more generally applicable model that includes less species-specific parameters and functions and can therefore account for risk-averse behavior in many RL situations and in different species.

According to our mathematical analysis and simulation results (see Figure 6c), there is a direct relationship between the model bee's learning rate and the resulting risk-averse behavior. Based on this we can predict that individual differences between bees in the observed risk-averse behavior will be correlated with the learning rate of the bees. This prediction can be empirically tested by assessing the learning rate of individual bees based on a learning task and then measuring the amount of risk aversion displayed by the bees in a foraging task similar to that described in this work.

Keeping in mind that the model we have described is an abstract learning model and therefore should be related to biological learning mechanisms with caution, the evolved learning architecture can be used to shed light on the biological implementation of RL. Thus, the significance of using heterosynaptic learning rules in our model should be noted: Through the learning rules evolved, a synapse can be modified even when its corresponding flower type was not chosen. This allows for nontrivial interactions between flower types when forming predictions as to the reward expected from each flower. For example, the expectation from one flower type can be updated as a function of the disappointment or surprise encountered when choosing the other flower type. Heterosynaptic learning rules have been used in computational modeling of synaptic plasticity and have been recognized as contributing to the complexity of the learning process, but surprisingly few studies have directly explored this phenomenon in the brain. Evidence from cerebellar (Dittman & Regehr, 1997) and hippocampal (Vogt & Nicoll, 1999) synapses shows that heterosynaptic plasticity indeed occurs in the brain and can affect the spiking patterns of neurons through interactions between adjacent synapses. In vitro recordings in *Aplysia* have directly demonstrated heterosynaptic facilitation of synapses from two presynaptic neurons

onto a common postsynaptic target (Schacher et al., 1997). We show here that heterosynaptic plasticity can exert a pronounced effect on the behavior of the organism and specifically illustrate this effect on the trade-off between exploration and exploitation characteristic of multi-armed bandit situations.

The learning mechanism described in this work can be closely related to the "adaptive critic" described by Sutton and Barto in the actor-critic architecture frequently used to implement RL in artificial agents (Barto, 1995; Sutton, 1988). Furthermore, although embedded in a very simple artificial neural network, the learning rules we have studied are biologically plausible and can be related to learning dependent on the dopaminergic system in the basal ganglia (Houk, Adams, & Barto, 1995). In the actor-critic model, an actor subnetwork learns to perform actions so as to maximize the weighted sum of future rewards, which is computed at every time step by a critic subnetwork (Barto, 1995). The critic is adaptive, in that it learns to predict the weighted sum of future rewards based on the current sensory input and the actor's policy, by means of an iterative process in which it compares its own predictions to the actual rewards obtained by the acting agent. The learning rule used by the adaptive critic is the TD learning rule (Sutton, 1988) in which the error between two adjacent predictions (the TD error) is used to update the critic's weights. In our model the synaptic weights come to represent the expected rewards, and as the inputs are differential, the activity of  $P$  represents an ongoing comparison between the expected reward in subsequent time steps. As in the critic model (Barto, 1995), this comparison provides the error measure by which the network updates its weights and learns to better predict future rewards.

Why, then, is learning in this model restricted by the intermodule dependencies only to the landing step, whereas learning in the classic adaptive critic model occurs at every step? The answer to this lies in the simplicity of the model and the learning task. The classic adaptive critic uses the error measure to learn to better predict the consequences of every input scenario, based on similarity of the current sensory input to previously encountered (and learned) inputs (Barto, 1995). This learning scheme is capable of dealing with complex situations in which outcomes depend on a long sequence of actions, but the final outcome of actions cannot be easily predicted by the sensory

inputs in every time step. Our model bee deals with a much simpler problem, as in every time step it can directly assess the final outcome of flying in its current direction. The inputs available to our model bee are not the primary sensory input from the retina, but a preprocessed input consisting of the percentages of reward-predicting features (flower colors) in the primary input. As a result, these inputs are not only sufficient to drive the learning process to generate valid predictions, but once a prediction is generated (according to the reward obtained), this prediction can be completely generalized to all the different input cases encountered by the bee during flight. Given a set of weights (which define the prediction of the network), a consistent prediction can be made for any combination of color percentages in the visual field, so there is no prediction error to learn from during the bee's flight. The prediction itself can be updated only when the bee lands and encounters (or does not encounter) an actual reward, which can be compared to the predicted reward.

Several studies have suggested that dopaminergic neurons in the basal ganglia may constitute a biological implementation of a TD reinforcement signal (Montague, Dayan, Nowlan, Pouget, & Sejnowski, 1993; Montague et al., 1996; Schultz et al., 1997), and together with the striatum, they may implement an actor-critic learning model (Houk et al., 1995; Schultz, 1998). The output of unit *P* in our model, as in that of Montague et al.'s (1995) model, quite accurately captures the essence of the activity patterns of midbrain dopaminergic neurons (Montague et al., 1996; Schultz et al., 1997) in primates and rodents, and the corresponding octopaminergic neurons in bees (Hammer, 1997; Menzel & Muller 1996). These neurons, originating in the SNc and VTA of primates and rodents, presumably project a widespread prediction-error signal to many cortical and subcortical areas, including the striatum (Schultz, 1998; Schultz et al., 1995; Waelti, Dickinson, & Schultz, 2001). In corticostriatal synapses, this signal is implicated in mediating learning by reinforcement, by neuromodulating synaptic plasticity [providing a "now learn" signal that allows synaptic plasticity and long-term memory formation in areas that were active at the same time or immediately before a rewarding input was encountered (Bailey et al., 2000; Wickens & Kötter, 1995)].

The similarity between dopamine-dependent plasticity in corticostriatal synapses (Kimura &

Matsumoto, 1997; Suzuki, Miura, Mishimura, & Aosaki, 2001; Wickens & Kötter, 1995) and the dependencies evolved in our model should be noted. However, this comparison should be done with caution, as ours is a very simplified model. We have found that efficient RL is dependent upon a three-factor Hebbian learning rule, in which the synaptic weights are updated as a function of a neuromodulating reward signal, as well as the presynaptic and post-synaptic factors. Our demonstration of the optimality of this learning rule to RL has bearing on the computational function of dopamine-dependent plasticity in the basal ganglia. Since the current model reflects mainly the critic module of the actor-critic framework and consists only of a very simplistic actor (the probabilistic action function), future work that will focus on elaborating the actor component of the model will undoubtedly increase the relevance of the model to learning in the basal ganglia.

In summary, the significance of this work is three-fold: On the one hand we show the strength of simple evolutionary computation models in evolving fundamental processes such as reinforcement learning, and on the other hand we show that optimal reinforcement learning can directly explain complex behaviors such as risk aversion and probability matching, without need for further assumptions. This is done by means of a very simple (and thus easily analyzed) model, which nonetheless has important bearing on learning in biological neuronal circuits.

## Notes

- 1 We have also tried to evolve bees in a third scenario in which the environment contains two variably rewarding flowers (each rewarding with a certain probability, and empty otherwise). Unfortunately, the bees' fitness remained at random level even after thousands of generations, and no RL was evolved. Apparently, in our framework, such excessive uncertainty of the environment is too difficult for the evolutionary process to solve, and the underlying consistencies cannot be discovered and exploited by the evolving bees. Note, however, that had any learning mechanism that uses Hebbian learning based on an error function been evolved in this environment (e.g., with a lower learning rate that allows the bee to compute the mean rewards of each flower type more accurately), our mathematical analysis below ensures that risk-averse behavior would have also emerged.

- 2 Note that our mathematical analysis applies only to the classical risk-aversion scenario (March, 1996) in which one reward source is constant rewarding (not risky at all), and does not extend to the case in which both flowers types are variably rewarding. The emergence of risk aversion in the general case of two variably rewarding flowers with the same mean reward but with two different probabilities of rewarding, can, however, be shown in simulation in the evolved bees.
- 3 Probability matching should not be confused with the "matching law" (Herrnstein, 1997), which describes asymptotic behavior when choosing between alternatives that show diminishing rewards as they are chosen more frequently (e.g., the depletion of a patch of flowers continuously foraged upon).

## Acknowledgments

We thank the Telluride Workshop on Neuromorphic Engineering (2000) for providing the setting and equipment for the minirobot experiments. We are grateful to Dr. Tamar Keasar for the valuable contribution of ideas and experimental data from probability- matching experiments in bees. We thank our anonymous referees whose thorough comments and constructive suggestions have contributed greatly to this final result.

Y. N. wishes to dedicate this work to her eternally beloved Jörg Kramer.

## References

- Ackley, D., & Littman, M. (1991). Interactions between learning and evolution. In C. G. Langton, C. Taylor, J. D. Farmer, & S. Rasmussen (Eds.), *Artificial life II*. Redwood City, CA: Addison-Wesley.
- Bailey, C., Giustetto, M., Huang, Y., Hawkins, R., & Kandel, E. (2000). Is heterosynaptic modulation essential for stabilizing hebbian plasticity and memory? *Nature Reviews Neuroscience*, 1, 11–20.
- Barto, A. G. (1995). Adaptive critic and the basal ganglia. In J. C. Houk, J. L. Davis, & D. G. Beiser (Eds.), *Models of information processing in the basal ganglia* (pp. 215–232). Cambridge, MA: MIT Press.
- Baxter, J. (1992). The evolution of learning algorithms for artificial neural networks. In D. Green & T. Bossomaier (Eds.), *Complex systems*. Amsterdam: IOS Press.
- Bitterman, M. E. (1965). Phyletic differences in learning. *American Psychologist*, 20, 396–410.
- Breiman, L. (1968). *Probability*. Reading, MA: Addison-Wesley.
- Chalmers, D. J. (1990). The evolution of learning: An experiment in genetic connectionism. In D. S. Touretzky, J. L. Elman, T. J. Sejnowski, & G. E. Hinton (Eds.), *Proceedings of the 1990 connectionist models summer school* (pp. 81–90). San Mateo, CA: Morgan Kaufmann.
- Dittman J. S., & Regehr, W. G. (1997). Mechanism and kinetics of heterosynaptic depression at a cerebellar synapse. *Journal of Neuroscience*, 17(23), 9048–9059.
- Domjan, M. (1998). *The principles of learning and behavior* (4th ed., pp. 178–185). Pacific Grove, CA: Brooks/Cole.
- Donahoe, J. W., & Packard-Dorsel, V. (Eds.). (1997). *Neural network models of cognition: Biobehavioral foundations*. Amsterdam: Elsevier Science.
- Fellous, J.-M., & Linster, C. (1998). Computational models of neuromodulation: A review. *Neural Computation*, 10, 791–825.
- Floreano, D., & Mondada, F. (1996). Evolution of homing navigation in a real mobile robot. *IEEE Transactions on Systems, Man and Cybernetics*, 26(3), 396–407.
- Floreano, D., & Mondada, F. (1998). Evolutionary neurocontrollers for autonomous mobile robots. *Neural Networks*, 11, 1461–1478.
- Floreano, D., & Urzelai, J. (2000). Evolutionary robots with online self organization and behavioral fitness. *Neural Networks*, 13, 431–443.
- Floreano, D., & Urzelai, J. (2001). Evolution of plastic control networks. *Autonomous Robots*, 11, 311–317.
- Fontanari, J. F., & Meir, R. (1991). Evolving a learning algorithm for the binary perceptron. *Network*, 2(4), 353–359.
- Graybiel, A. M., & Kimura, M. (1995). Adaptive neural networks in the basal ganglia. In J. C. Houk, J. L. Davis, & D. G. Beiser (Eds.), *Models of information processing in the basal ganglia* (pp. 103–116). Cambridge, MA: MIT Press.
- Greggers, U., & Menzel, R. (1993). Memory dynamics and foraging strategies of honeybees. *Behavioral Ecology and Sociobiology*, 32, 17–29.
- Hammer, M. (1993). An identified neuron mediates the unconditioned stimulus in associative learning in honeybees. *Nature*, 366, 59–63.
- Hammer, M. (1997). The neural basis of associative reward learning in honeybees. *Trends in Neuroscience*, 20(6), 245–252.
- Harder, L. D., & Real, L. A. (1987). Why are bumble bees risk averse? *Ecology*, 68(4), 1104–1108.
- Hardy G. H., Littlewood, J. E., & Polya, G. (1934). *Inequalities*. Cambridge, UK: Cambridge University Press.
- Herrnstein, R., & Loveland, D. H. (1975). Maximizing and matching on concurrent ratio schedules. *Journal of the Experimental Analysis of Behavior*, 24, 107–116.
- Herrnstein, R. J. (1997). *The matching law: Papers in psychology and economics*. Cambridge, MA: Harvard University Press.
- Hinton, G. E., & Nowlan, S. J. (1987). How learning guides evolution. *Complex Systems*, 1, 495–502.



- Houk, J. C., Adams, J. L., & Barto, A. G. (1995). A model of how the basal ganglia generate and use neural signals that predict reinforcement. In J. C. Houk, J. L. Davis, & D. G. Beiser (Eds.), *Models of information processing in the basal ganglia* (pp. 249–270). Cambridge, MA: MIT Press.
- Houk, J. C., & Wise, S. P. (1995). Distributed modular architecture linking basal ganglia, cerebellum and cerebral cortex: Their role in planning and controlling action. *Cerebral Cortex*, 5, 95–111.
- Kacelnik, A., & Bateson, M. (1996). Risky theories—the effect of variance on foraging decisions. *American Zoologist*, 36, 402–434.
- Kaesar, T., Rashkovich, E., Cohen, D., & Shmida, A. (in press). Choice behavior of bees in two-armed bandit situations: Experiments and possible decision rules. *Behavioral Ecology*.
- Kimura, M., & Matsumoto, N. (1997). Nigrostriatal dopamine system may contribute to behavioral learning through providing reinforcement signals to the striatum. *European Neurology*, 38(Suppl. 1), 11–17.
- March, J. G. (1996). Learning to be risk averse. *Psychological Review*, 103(2), 309–319.
- Menzel, R., & Muller, U. (1996). Learning and memory in honeybees. From behavior to neural substrates. *Annual Reviews in neuroscience*, 19, 379–404.
- Mitchell, T., (1997). *Machine learning*. New York: McGraw-Hill.
- Montague, P. R. (1997). Biological substrates of predictive mechanisms in learning and action choice. In J. W. Donahoe & V. Packard-Dorsel (Eds.), *Neural network models of cognition: Biobehavioral foundations* (pp. 406–421). Amsterdam: Elsevier Science.
- Montague, P. R., Dayan, P., Nowlan, S. J., Pouget, A., & Sejnowski, T. J. (1993). Using aperiodic reinforcement for directed self-organization. In C. L. Giles, S. J. Hanson, & J. D. Cowan (Eds.), *Advances in neural information processing systems* (Vol. 5, pp. 969–976). San Mateo, CA: Morgan Kaufmann.
- Montague, P. R., Dayan, P., Person, C., & Sejnowski, T. J. (1995). Bee foraging in uncertain environments using predictive Hebbian learning. *Nature*, 377, 725–728.
- Montague, P. R., Dayan, P., & Sejnowski, T. J. (1996). A framework for mesencephalic dopamine systems based on predictive Hebbian learning. *Journal of Neuroscience*, 16, 1936–1947.
- Nolfi, S., Elmann, J. L., & Parisi, D. (1994). Learning and evolution in neural networks. *Adaptive Behavior*, 3, 5–28.
- Real, L. A. (1991). Animal choice behavior and the evolution of cognitive architecture. *Science*, 253, 980–985.
- Real, L. A. (1996). Paradox, performance and the architecture of decision making in animals. *American Zoologist*, 36, 518–529.
- Rothschild, M., & Stiglitz, J. (1970). Increasing risk: I. A definition. *Journal of Economic Theory*, 2, 225–243.
- Schacher, S., Wu, F., & Sun, Z. Y. (1997). Pathway-specific synaptic plasticity: Activity-dependent enhancement and suppression of long-term heterosynaptic facilitation at converging inputs on a single target. *Journal of Neuroscience*, 17(2), 597–606.
- Schultz, W. (1998). Predictive reward signal of dopamine neurons. *Journal of Neurophysiology*, 80, 1–27.
- Schultz, W. (2000). Multiple reward signals in the brain. *Nature Reviews Neuroscience*, 1, 199–207.
- Schultz, W., Dayan, P., & Montague, P. R. (1997). A neural substrate of prediction of reward. *Science*, 275, 1593–1599.
- Schultz, W., Romo, R., Ljungberg, T., Mirenowicz, J., Hollerman, J. R., & Dickinson, A. (1995). Reward-related signals carried by dopamine neurons. In J. C. Houk, J. L. Davis, & D. G. Beiser (Eds.), *Models of information processing in the basal ganglia* (pp. 233–248). Cambridge, MA: MIT Press.
- Seth, A. K. (1999). Evolving behavioral choice: An investigation into Herrnstein's matching law. In D. Floreano, J. Nicoud, & F. Mondada (Eds.), *Advances in artificial life, 5th European conference, ECAL '99* (pp. 225–235). Lausanne, Switzerland: Springer.
- Smallwood, P. D. (1996). An introduction to risk sensitivity: The use of Jensen's inequality to clarify evolutionary arguments of adaptation and constraint. *American Zoologist*, 36, 392–401.
- Sutton, R. S. (1988). Learning to predict by the method of temporal difference. *Machine Learning*, 3, 9–44.
- Sutton, R. S., & Barto, A. G. (1998). *Reinforcement learning: An introduction*. Cambridge, MA: MIT Press.
- Suzuki, T., Miura, M., Mishimura, K., & Aosaki, T. (2001). Dopamine-dependent synaptic plasticity in the striatal cholinergic interneurons. *Journal of Neuroscience*, 21(17), 6492–6501.
- Thuijsman, F., Peleg, B., Amitai, M., & Shmida, A. (1995). Automata, matching and foraging behavior of bees. *Journal of Theoretical Biology*, 175, 305–316.
- Unemi, T., Negayoshi, M., Hirayama, N., Nade, T., Yano, K., & Mausjima, Y. (1994). Evolutionary differentiation of learning abilities—A case study on optimizing parameter values in Q-learning by a genetic algorithm. In R. A. Brooks & P. Maes (Eds.), *Artificial life IV* (pp. 331–336). Cambridge, MA: MIT Press.
- Vogt, K. E., & Nicoll, R. E. (1999). Glutamate and gamma-aminobutyric acid mediate a heterosynaptic depression at mossy fiber synapses in the hippocampus. *Proceedings of the National Academy of Science, USA*, 96, 1118–1122.
- Waelti, P., Dickinson, A., & Schultz, W. (2001). Dopamine responses comply with basic assumptions of formal learning theory. *Nature*, 412, 43–48.
- Wickens, J., & Kötter, R. (1995). Cellular models of reinforcement. In J. C. Houk, J. L. Davis, & D. G. Beiser (Eds.), *Models of information processing in the basal ganglia* (pp. 187–214). Cambridge, MA: MIT Press.

Wilson, S. W. (1996). Explore/exploit strategies in autonomy. In P. Maes, M. Mataric, J. Pollack, J. A. Meyer, & S. Wilson (Eds.), *From animals to animats 4: Proceedings of the Fourth International Conference on the Simulation of Adaptive Behavior* (pp. 325–332). Cambridge, MA: MIT Press/Bradford Books.

## Appendix

### Mathematical Analysis: Risk Aversion Is Ordered by Learning Rate

We consider the bee's long-term choice dynamics as a sequence of  $N$  cycles, each choice of  $(v)$  beginning a cycle. Let  $n_i \geq 0$  be the number of visits to constant flowers in the  $i$ th cycle. The frequency  $f_v$  of visits to  $(v)$  is determined (via Birkhoff's ergodic theorem (Breiman, 1968), an extension to dependent variables of the strong law of large numbers) by the expected number of visits to  $(c)$  in a typical cycle  $[E(n)]$ :

$$\begin{aligned} f_v &= \lim_{N \rightarrow \infty} \frac{N}{N + \sum_{i=1}^N n_i} \\ &= \lim_{N \rightarrow \infty} \frac{1}{1 + \frac{1}{N} \sum_{i=1}^N n_i} = \frac{1}{1 + E(n)} \end{aligned} \quad (A1)$$

As  $W_c$  is constant, the bee's choices are only a function of  $W_v$ , and we can define the bee's choice function as  $p_v(W_v)$ , the probability of choosing  $(v)$  in a trial in which the synaptic weight corresponding to the variably rewarding flower is  $W_v$ . Thus given  $W_v$ ,  $[n_i + 1]$  is geometrically distributed with  $p_v(W_v)$ , giving:

$$\begin{aligned} E(n) &= E[E(n|W_v)] = E\left[\frac{1}{p_v(W_v)} - 1\right] \\ &= E\left[\frac{1}{p_v(W_v)}\right] - 1 \end{aligned} \quad (A2)$$

and so

$$f_v = \frac{1}{E\left[\frac{1}{p_v(W_v)}\right]} \quad (A3)$$

According to the mathematical definition of riskiness that comes from theories of second degree stochastic dominance (Hardy, Littlewood, & Polya, 1934), for

$X$  and  $Y$  with a finite equal mean, Rothschild and Stiglitz (1970) show the following three statements to be equivalent:

1.  $EU(X) \geq EU(Y)$  for every concave function  $U$  for which these expectations exist.
2.  $E[\max(X - x, 0)] \leq E[\max(Y - x, 0)]$  for all  $x \in \mathcal{R}$ .
3. There exists on some probability space two random variables  $X$  and  $Z$  such that  $Y = X + Z$  and  $E(Z | X) = 0$  with probability 1.

Statement (1) provides the mathematical definition of riskiness, that is,  $X$  is less risky than  $Y$  if (1) is true, as for every concave utility function the mean subjective reward obtained from  $X$  is greater than that obtained from  $Y$  so every risk averter would prefer  $X$  to  $Y$ . Statement (2) is an easier condition to check when determining which of two random variables is riskier. Statement (3) is a mathematically equivalent definition of riskiness that we will use later in our analysis to prove that the bee is always risk averse. According to this,

**Lemma:** If  $X, X_1, X_2, X_3, \dots$  are identically distributed (not necessarily independent) random variables with a finite mean,  $Y = \sum_{i=1}^{\infty} \alpha_i X_i$  (where  $\vec{\alpha}_i$  is a probability vector) is less risky than  $X$ .

**Proof:**

$$\begin{aligned} \sum \alpha_i X_i - x &= \sum \alpha_i (X_i - x) \\ &\leq \sum \alpha_i [\max(X_i - x, 0)] \end{aligned} \quad (A4)$$

Since the right-hand side is nonnegative,

$$\begin{aligned} \max\left[\sum \alpha_i X_i - x, 0\right] &\leq \sum \alpha_i [\max(X_i - x, 0)] \end{aligned} \quad (A5)$$

Now taking expectations of both sides

$$\begin{aligned} E[\max(\sum \alpha_i X_i - x, 0)] &\leq \sum \alpha_i E[\max(X_i - x, 0)] \\ &= \sum \alpha_i E[\max(X - x, 0)] \\ &= E[\max(X - x, 0)] \end{aligned} \quad \therefore \quad (A6)$$

As a corollary, we shall prove that exponential smoothers such as  $W_v(\eta)$  are risk ordered such that a

lower learning rate  $\beta$  leads to less risk aversion than a higher learning rate  $\alpha$  ( $0 < \beta < \alpha < 1$ ).

**Lemma:** Let  $W_v(\eta)$  be an exponentially weighted average of identically distributed variables  $V_i$  ( $i = 1, 2, 3, \dots$ ) as in Equation 8, then  $W_v(\alpha)$  is riskier than  $W_v(\beta)$  for every  $0 < \beta < \alpha < 1$ .

**Proof:** Let us define  $W_v^{(k)}(\alpha)$  identically distributed (not independent) variables as the following:

$$\begin{aligned} W_v^{(1)}(\alpha) &= \alpha v_1 + \alpha(1 - \alpha)v_2 \\ &\quad + \alpha(1 - \alpha)^2 v_3 + \dots \\ W_v^{(2)}(\alpha) &= \alpha v_2 + \alpha(1 - \alpha)v_3 \\ &\quad + \alpha(1 - \alpha)^2 v_4 + \dots \\ &\vdots \\ W_v^{(k)}(\alpha) &= \alpha v_k + \alpha(1 - \alpha)v_{k+1} \\ &\quad + \alpha(1 - \alpha)^2 v_{k+2} + \dots \end{aligned} \quad (\text{A7})$$

Let us now choose a specific probability vector  $(\vec{\alpha}_i)$  as follows:

$$\alpha_1 = \frac{\beta}{\alpha}; \quad \alpha_n = \alpha_1(\alpha - \beta)(1 - \beta)^{n-2} \quad (n \geq 2) \quad (\text{A8})$$

We then have

$$\begin{aligned} \sum_{k=1}^{\infty} \alpha_k W_v^{(k)}(\alpha) &= \frac{\beta}{\alpha} \cdot \alpha \cdot [v_1 + (1 - \alpha)v_2 \\ &\quad + (1 - \alpha)^2 v_3 + \dots \\ &\quad + (\alpha - \beta)v_2 + (\alpha - \beta)(1 - \alpha)v_3 \\ &\quad + (\alpha - \beta)(1 - \alpha)^2 v_4 + \dots \\ &\quad + (\alpha - \beta)(1 - \beta)v_3 \\ &\quad + (\alpha - \beta)(1 - \beta)(1 - \alpha)v_4 \\ &\quad + (\alpha - \beta)(1 - \beta)(1 - \alpha)^2 v_5 + \dots] \\ &= \beta[v_1 + (1 - \beta)v_2 \\ &\quad + (1 - \beta)^2 v_3 + \dots] = W_v(\beta) \end{aligned} \quad (\text{A9})$$

Thus  $W_v(\beta)$ , as a weighted sum of  $W_v^{(k)}(\alpha)$ , is less risky than  $W_v(\alpha)$ .  $\therefore$

Now tying this to Equation (A3) and to Statement (1) of the Rothschild and Stiglitz (1970) theorem, if  $\phi(\cdot) = 1/pv(\cdot)$  is convex (and so  $-1/pv(\cdot)$  is concave), then

$$E\left(-\frac{1}{p_v(W_v(\beta))}\right) \geq E\left(-\frac{1}{p_v(W_v(\alpha))}\right) \quad (\text{A10})$$

$$\begin{aligned} \Rightarrow f_v(\alpha) &= \frac{1}{E\left(\frac{1}{p_v(W_v(\alpha))}\right)} \\ &\leq \frac{1}{E\left(\frac{1}{p_v(W_v(\beta))}\right)} = f_v(\beta) \end{aligned} \quad (\text{A11})$$

And the bee will display ordered risk-averse behavior: The higher the learning rate, the lower is the frequency of visits  $f_v$  to the (v) flowers. Convexity of  $1/pv(\cdot)$  is a mild assumption as for every concave increasing  $p_v$ ,  $1/pv(\cdot)$  is strictly convex, so convexity will also be preserved under minor departures from concavity of  $p_v$ .

According to Statement (3) of the Rothschild and Stiglitz (1970) theorem, if  $Y$  is obtained from  $X$  by further fair gambling, then  $X$  is less risky than  $Y$ . Thus when both flower types yield the same mean reward,  $W_v$  is riskier than  $W_c$ . From this follows that even with low learning rates, since  $p_v$  is symmetric with respect to  $W_v$  and  $W_c$  [i.e.,  $p_v(W_c) = 1/2$ ], when both flower types reward with the same mean, the frequency  $f_v$  is always less than  $1/p_v(W_c) = 1/2$ , and the bee is always risk averse.

## About the Authors



**Yael Niv** is currently a Ph.D. student at the Interdisciplinary Center for Neural Computation at the Hebrew University in Jerusalem. She received her master's degree in psychology from Tel Aviv University in 2001, after undergraduate studies in the Interdisciplinary Program for Fostering Excellence. Her research interests focus on understanding the process of reinforcement learning and its relationship to neural mechanisms, specifically to the basal ganglia, using computational modeling techniques. This article summarizes her master's thesis research. *Address:* Department of Psychology, Tel-Aviv University Tel-Aviv 69978, Israel, *E-mail:* yaeln@cns.tau.ac.il



**Daphna Joel** received her Ph.D. in psychology from Tel Aviv University (TAU) in 1998 and joined the faculty of TAU, after receiving the Alon fellowship for young Israeli scientists. Her research interests focus on understanding the involvement of basal ganglia–thalamocortical circuits in normal and abnormal behaviour, using animal models of psychopathology as well as computational approaches. In a recent series of papers she has presented a novel rat model of obsessive compulsive disorder. *Address:* Department of Psychology, Tel Aviv University, Tel-Aviv 69978, Israel, *E-mail:* djoel@post.tau.ac.il



**Isaac Meilijson** was born in Argentina and raised in Chile. He studied at the Hebrew University of Jerusalem (B.Sc. 1965, M. Sc. 1967) and at U.C. Berkeley (Ph.D. 1969). He is Professor of Statistics at Tel Aviv University, where he has been on the faculty since 1971. He defines himself as a stochastic modeler, with research interests ranging over probability, statistics, and operations research, with applications to issues in mathematical economics, biomathematics, and neural computation. *Address:* School of Mathematical Sciences, Tel Aviv University, Tel Aviv 69978, Israel. *E-mail:* isaco@math.tau.ac.il



**Eytan Ruppin** is an associate professor in the School of Computer Sciences and the School of Medicine of Tel Aviv University (TAU). He received his M.D. (1986) M.Sc. (1989) and Ph.D. (1993) degree in computer science from TAU. Receiving the Rothschild fellowship, he was a postdoctorate fellow at the University of Maryland (1993–1995). In 1995 he returned to join the faculty of TAU, receiving the Alon fellowship for young Israeli scientists. Prof. Ruppin's research has focused in the past on computational studies of brain and cognitive disorders, investigating possible causal links between pathological brain alterations and the clinical manifestations of brain disorders such as Alzheimer's, stroke, and schizophrenia. In recent years, his research interests have turned to developing and studying evolutionary models of neural processing in embodied autonomous agents. *Address:* School of Computer Sciences, Tel Aviv University, Tel Aviv 69978, Israel. *E-mail:* ruppin@math.tau.ac.il. <http://www.math.tau.ac.il/~ruppin>