
The RL Perceptron: Generalisation Dynamics of Policy Learning in High Dimensions

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Reinforcement learning (RL) algorithms have proven transformative in a range of
2 domains. To tackle real-world domains, these systems often use neural networks
3 to learn policies directly from pixels or other high-dimensional sensory input. By
4 contrast, much theory of RL has focused on discrete state spaces or worst-case
5 analysis, and fundamental questions remain about the dynamics of policy learning
6 in high-dimensional settings. Here, we propose a solvable high-dimensional model
7 of RL that can capture a variety of learning protocols, and derive its typical
8 dynamics as a set of closed-form ordinary differential equations (ODEs). We derive
9 optimal schedules for the learning rates and task difficulty—analogueous to annealing
10 schemes and curricula during training in RL—and show that the model exhibits rich
11 behaviour, including delayed learning under sparse rewards; a variety of learning
12 regimes depending on reward baselines; and a speed-accuracy trade-off driven by
13 reward stringency. Experiments on a variant of the Procgen game “Bossfight” also
14 show such a speed-accuracy trade-off in practice. Together, these results take a
15 step towards closing the gap between theory and practice in high-dimensional RL.

16 Recent years have seen rapid progress in Reinforcement Learning (RL): algorithmic and engineering
17 breakthroughs led to super-human performance in a variety of domains, for example complex games
18 like Go [Silver et al., 2016, Mnih et al., 2015]. Despite these practical successes, our theoretical
19 understanding of RL for high-dimensional problems requiring non-linear function approximation
20 is still limited. While comprehensive theoretical results exist for tabular RL, where the state and
21 action spaces are discrete and small enough for value functions to be represented directly, the curse
22 of dimensionality limits these methods to low-dimensional problems. The lack of a clear notion of
23 similarity between discrete states further means that tabular methods do not address the core question
24 of generalisation: how are values and policies extended to unseen states and across seen states [Kirk
25 et al., 2023]? As a consequence, much of this theoretical work is far from the current practice of RL,
26 which increasingly relies on deep neural networks to approximate and generalise value functions,
27 policies and other building blocks of RL. Moreover, while RL theory has often addressed “worst-case”
28 performance and convergence behaviour, the *typical* behaviour has received comparatively little
29 attention (cf. further related work below). Meanwhile, a growing sub-field of deep learning theory
30 has employed tools from statistical mechanics to analyse various supervised learning paradigms
31 in the average-case, see Seung et al. [1992], Engel and Van den Broeck [2001], Carleo et al. [2019],
32 Bahri et al. [2020], Gabrié et al. [2023] for classical and recent reviews. While this approach has
33 recently been extended to curriculum learning [Saglietti et al., 2022], continual learning [Asanuma
34 et al., 2021, Lee et al., 2021, 2022], few-shot learning [Sorscher et al., 2022] and transfer learning
35 [Lampinen and Ganguli, 2018, Dhifallah and Lu, 2021, Gerace et al., 2022], RL has not been
36 analysed yet using statistical mechanics—a gap we address here by studying the high-dimensional
37 generalisation dynamics of a simple neural network trained on a reinforcement learning task.

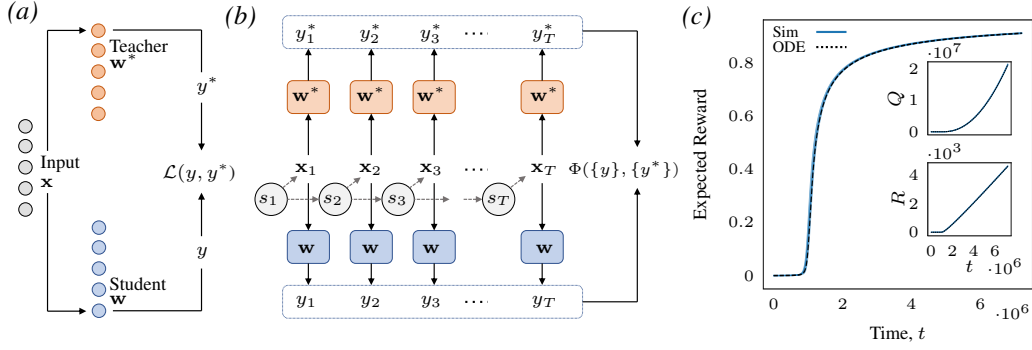


Figure 1: **The RL-Perceptron is a model for policy learning in high dimensions.** (a) In the classic teacher-student model for supervised learning, a neural network called the student is trained on inputs x whose label y^* is given by another neural network, called the teacher. (b) In the RL setting the student moves through states s_t making a series of T choices given in response to inputs x_t . The RL-perceptron is an extension of the teacher-student model as we assume there is a ‘right’ choice y_t on each timestep given by a teacher network. The student receives a reward after T decisions according to a criterion Φ that depends on the choices made and the corresponding correct choices. (c) Example learning dynamics in the RL-perceptron for a problem with $T = 12$ choices where the reward is given only if all the decisions are correct. The plot shows the expected reward of a student trained in the RL perceptron setting in simulations (solid) and for our theoretical results (dashed) obtained from solving the dynamical equations eqs. (5) and (6). Finite size simulations and theory show good agreement. We reduce the stochastic evolution of the high dimensional student to the study of deterministic evolution of two scalar quantities R and Q (more details in Sec. 2.1), their evolution are shown in the inset. *Parameters:* $D = 900$, $\eta_1 = 1$, $\eta_2 = 0$, $T = 12$.

38 **The RL perceptron:** In the classic teacher-student model of supervised learning [Gardner and
39 Derrida, 1989, Seung et al., 1992], a neural network called the student is trained on inputs x whose
40 labels y^* are given by another neural network called the teacher (see fig. 1a). The goal of the
41 student is to learn the function represented by the teacher from samples (x, y^*) . In RL, agents face
42 a sequential decision-making task in which a sequence of correct intermediate choices is required
43 to successfully complete an episode. We translate this process into the RL perceptron, a solvable
44 model for a high-dimensional, sequential policy learning task shown in fig. 1b. The student with
45 weights w takes a sequence of T choices over an episode. The correct choices are governed by
46 the same teacher network w^* , i.e. the same underlying rule throughout every time-step of every
47 episode. Crucially, unlike in the supervised learning setting, the student does not observe the correct
48 choice for each input; instead, it receives a reward which depends on whether earlier decisions are
49 correct. For instance, the student could receive a reward only if all T choices are correct, and no
50 reward otherwise—a learning signal that is considerably less informative than in supervised learning.
51 In addition to introducing the RL perceptron, our **main contributions** are as follows:

- 52 • We derive an asymptotically exact set of Ordinary Differential Equations (ODEs) that
53 describe the typical learning dynamics of policy gradient RL agents by building on classic
54 work by Saad and Solla [1995], Biehl and Schwarze [1995], see section 2.1.
- 55 • We use these ODEs to characterize learning behaviour in a diverse range of scenarios:
 - 56 – We explore several sparse delayed reward schemes and investigate the impact of
57 negative rewards (section 2.2)
 - 58 – We derive optimal learning rate schedules and episode length curricula, and recover
59 annealing strategies typically used in practice (section 2.3)
 - 60 – At fixed learning rates, we identify ranges of learning rates for which learning is
61 ‘easy,’ and ‘hybrid-hard’—possibly causing a critical slowing down in the dynamics
62 (section 2.4)
 - 63 – We identify a speed-accuracy trade-off driven by reward stringency (section 2.5)
- 64 • Finally we demonstrate that a similar speed-accuracy trade-off exists in simulations of high-
65 dimensional policy learning from pixels using the procgen environment “Bossfight” [Cobbe
66 et al., 2019], see section 3.

67 **Further related work**

68 **Sample complexity in RL.** An important line of work in the theory of RL focuses on the sample
 69 complexity and other learnability measures for specific classes of models such as tabular RL [Azar
 70 et al., 2017, Zhang et al., 2020b], state aggregation [Dong et al., 2019], various forms of MDPs [Jin
 71 et al., 2020, Yang and Wang, 2019, Modi et al., 2020, Ayoub et al., 2020, Du et al., 2019a, Zhang
 72 et al., 2022], reactive POMDPs [Krishnamurthy et al., 2016], and FLAMBE [Agarwal et al., 2020].
 73 Here, we are instead concerned with the learning dynamics: how do reward rates, episode length, etc.
 74 influence the speed of learning and the final performance of the model.

75 **Statistical learning theory for RL** aims at finding complexity measures analogous to the Rademacher
 76 complexity or VC dimension from statistical learning theory for supervised learning Bartlett and
 77 Mendelson [2002], Vapnik and Chervonenkis [2015]. Proposals include the Bellman Rank Jiang et al.
 78 [2017], or the Eluder dimension [Russo and Van Roy, 2013] and its generalisations [Jin et al., 2021].
 79 This approach focuses on worst-case analysis, which typically differs significantly from practice (at
 80 least in supervised learning [Zhang et al., 2021]). Furthermore, complexity measures for RL are
 81 generally more suitable for value-based methods; policy gradient methods have received less attention
 82 despite their prevalence in practice Bhandari and Russo [2019], Agarwal et al. [2021]. We focus
 83 instead on average-case dynamics of policy-gradient methods.

84 **Dynamics of learning.** A series of recent papers considered the dynamics of temporal-difference
 85 learning and policy gradient in the limit of wide two-layer neural networks Cai et al. [2019], Zhang
 86 et al. [2020a], Agazzi and Lu [2021, 2022]. These works focus on one of two “wide” limits: either
 87 the neural tangent kernel [Jacot et al., 2018, Du et al., 2019b] or “lazy” regime [Chizat et al., 2019],
 88 where the network behaves like an effective kernel machine and does not learn data-dependent
 89 features, which is key for efficient generalisation in high-dimensions. In our setting, the success
 90 of the student crucially relies on learning the weight vector of the teacher, which is hard for lazy
 91 methods [Ghorbani et al., 2019, 2020, Chizat and Bach, 2020, Refinetti et al., 2021]. The other
 92 “wide” regime is the mean-field limit of interacting particles, akin to Mei et al. [2018], Chizat and
 93 Bach [2018], Rotskoff and Vanden-Eijnden [2018], where learning dynamics are captured by a
 94 non-linear partial differential equation. While this elegant description allows them to establish global
 95 convergence properties, it is hard to solve in practice. The ODE description we derive here instead
 96 will allow us to describe a series of effects in the following sections.

97 **1 The RL Perceptron: setup and learning algorithm**

98 We study the simplest possible student network, a perceptron with weight vector \mathbf{w} that takes in
 99 high-dimensional inputs $\mathbf{x} \in \mathbb{R}^D$ and outputs $y(\mathbf{x}) = \text{sgn}(\mathbf{w}^\top \mathbf{x})$. We interpret the outputs $y(\mathbf{x})$ as
 100 decisions, for example whether to go left or right in an environment. Because the student makes
 101 choices in response to high-dimensional inputs, it is analogous to a policy network. To train the
 102 network, we therefore consider a policy gradient learning update analogous to the REINFORCE
 103 algorithm [Sutton et al., 2000] that is adapted to the perceptron. At every timestep t during the μ th
 104 episode of length T , the agent occupies some state s_t in the environment, receives an observation \mathbf{x}_t^μ
 105 conditioned on s_t , and takes an action $y_t^\mu = \text{sgn}(\mathbf{w}^\top \mathbf{x}_t^\mu)$, with $t = 1, \dots, T$. The correct choice for
 106 each input is given by a fixed perceptron teacher with weights \mathbf{w}^* . The crucial point is that the student
 107 does not have access to all the correct choices; it only receives a reward at the end of the episode *if*
 108 it completes the episode successfully, for example by making the correct decision at all times. If it does
 109 not succeed, it *may* receive a penalty; we will see in section 2.4 that receiving penalties is not always
 110 beneficial. In our setup, this translates into a weight update at the end of the μ^{th} episode that is given
 111 by

$$\mathbf{w}^{\mu+1} = \mathbf{w}^\mu + \frac{\eta_1}{\sqrt{D}} \left(\frac{1}{T} \sum_{t=1}^T y_t \mathbf{x}_t \mathbb{I}(\Phi) \right)^\mu - \frac{\eta_2}{\sqrt{D}} \left(\frac{1}{T} \sum_{t=1}^T y_t \mathbf{x}_t (1 - \mathbb{I}(\Phi)) \right)^\mu, \quad (1)$$

112 where \mathbb{I} is an indicator function and Φ is the criterion that determines whether the episode was
 113 completed successfully—for instance, $\mathbb{I}(\Phi) = \prod_t \theta(y_t y_t^*)$ (where θ is the step function) if the student
 114 has to get every decision right in order to receive a reward. The update is general in the sense that the
 115 term proportional to the learning rate $\eta_1 > 0$ prescribes the reward update for the fulfillment of the
 116 condition, while the term proportional to $\eta_2 \geq 0$ gives us the possibility to add a penalty or negative

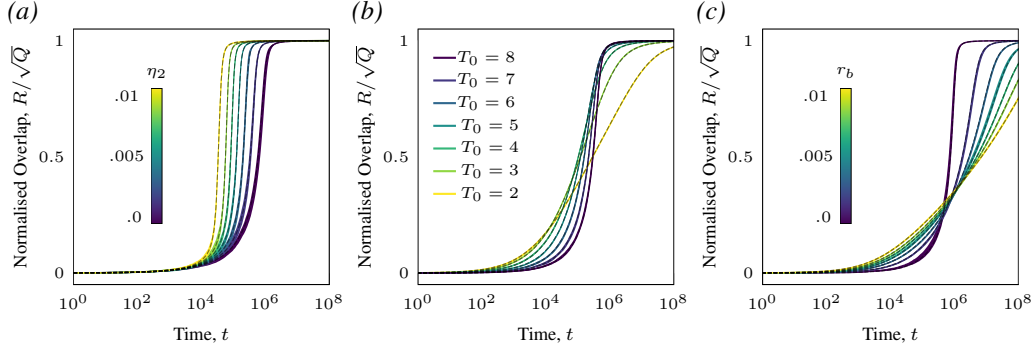


Figure 2: **ODEs accurately describe diverse learning protocols.** Evolution of the normalised student-teacher overlap ρ for the numerical solution of the ODEs (dashed) and simulation (coloured) in three reward protocols. All students receive a reward of η_1 for getting all decisions in an episode correct, and additionally: (a) A penalty η_2 (i.e. negative reward) is received if the agent does not survive until the end of an episode. (b) An additional reward of 0.2 is received if the agent survives beyond T_0 timesteps. (c) An additional reward r_b is received for every correct decision made in an episode. *Parameters:* $D = 900$, $T = 12$, $\eta_1 = 1$.

177 reward should the student not succeed. Note that in the case of $T = 1$, $\eta_2 = 0$, and $\mathbb{I}(\Phi) = \theta(yy^*)$,
 178 the learning rule updates the weight only if the student is correct on a given sample. It can thus be seen
 179 as the “opposite” of the famous perceptron learning rule of supervised learning [Rosenblatt, 1962],
 180 where weights are only updated if the student is wrong. For a more in-detail discussion of the relation
 181 between the weight update in eq. (1) and the REINFORCE algorithm, see appendix A.

122 2 Theoretical Results

123 2.1 A set of dynamical equations captures the learning dynamics of an RL perceptron exactly

124 The goal of the student during training is to emulate the teacher as closely as possible; or in other
 125 words, have a small number of disagreements with the teacher $y(\mathbf{x}) \neq y^*(\mathbf{x})$. The generalisation
 126 error is given by the average number of disagreements

$$\epsilon_g \equiv \langle y(\mathbf{x})y^*(\mathbf{x}) \rangle = \left\langle \text{sgn} \left(\mathbf{w}^* \cdot \mathbf{x} / \sqrt{D} \right) \text{sgn} \left(\mathbf{w} \cdot \mathbf{x} / \sqrt{D} \right) \right\rangle = \langle \text{sgn}(\nu) \text{sgn}(\lambda) \rangle \quad (2)$$

127 where the average $\langle \cdot \rangle$ is taken over the inputs \mathbf{x} , and we have introduced the scalar pre-activations
 128 for the student and the teacher, $\lambda \equiv \mathbf{w} \cdot \mathbf{x} / \sqrt{D}$ and $\nu \equiv \mathbf{w}^* \cdot \mathbf{x} / \sqrt{D}$, respectively. We can therefore
 129 transform the high-dimensional average over the inputs \mathbf{x} into a low-dimensional average over the
 130 pre-activations (λ, ν) . The average in eq. (2) can be carried out by noting that the tuple (λ, ν) follow
 131 a jointly Gaussian distribution with means $\langle \lambda \rangle = \langle \nu \rangle = 0$ and covariances

$$Q \equiv \langle \lambda^2 \rangle = \frac{\mathbf{w} \cdot \mathbf{w}}{D}, \quad R \equiv \langle \lambda \nu \rangle = \frac{\mathbf{w} \cdot \mathbf{w}^*}{D} \quad \text{and} \quad S \equiv \langle \nu^2 \rangle = \frac{\mathbf{w}^* \cdot \mathbf{w}^*}{D}. \quad (3)$$

132 These covariances, or overlaps as they are sometimes called in the literature, have a simple interpreta-
 133 tion. The overlap S is simply the length of the weight vector of the teacher; in the high-dimensional
 134 limit $D \rightarrow \infty$, $S \rightarrow 1$. Likewise, the overlap Q gives the length of the student weight vector; however,
 135 this is a quantity that will vary during training. For example, when starting from small initial weights,
 136 Q will be small, and grow throughout training. Lastly, the “alignment” R quantifies the correlation
 137 between the student and the teacher weight vector. At the beginning of training, $R \approx 0$, as both the
 138 teacher and the initial condition of the student are drawn at random. As the student starts learning,
 139 the overlap R increases. Evaluating the Gaussian average in eq. (2) shows that the generalisation
 140 error is then a function of the normalised overlap $\rho = R/\sqrt{Q}$, and given by

$$\epsilon_g = \frac{1}{\pi} \arccos \left(\frac{R}{\sqrt{Q}} \right) \quad (4)$$

141 The crucial point here is that we have reduced the description of the high-dimensional learning
 142 problem from the D parameters of the student weight \mathbf{w} to two time-evolving quantities, Q and R .
 143 We now discuss how to analyse their dynamics.

144 **The dynamics of order parameters.** At any given point during training, the value of the order
 145 parameters determines the test error via eq. (4). But how do the order parameters evolve during
 146 training with the update rule eq. (1)? We followed the approach of Kinzel and Ruján [1990], Saad
 147 and Solla [1995], Biehl and Schwarze [1995] to derive a set of dynamical equations that describe
 148 the dynamics of the student in the high-dimensional limit where the input dimension goes to infinity.
 149 We give explicit dynamics for different reward conditions Φ , namely requiring all decisions correct
 150 in an episode of length T ; requiring n or more decisions correct in an episode of length T ; and
 151 receiving reward for each correct response. Due to the length of these expressions, we report the
 152 generic expression of the updates in the supplementary material in appendix B. Below, we state a
 153 version of the equations for the specific reward condition where the agent must survive until the end
 154 of an episode to receive a reward, $\mathbb{I}(\Phi) = \prod_t^T \theta(y_t y_t^*)$. The ODEs for the order parameters then read

$$\frac{dR}{d\alpha} = \frac{\eta_1 + \eta_2}{\sqrt{2\pi}} \left(1 + \frac{R}{\sqrt{Q}} \right) P^{T-1} - \eta_2 R \sqrt{\frac{2}{\pi Q}} \quad (5)$$

$$\frac{dQ}{d\alpha} = (\eta_1 + \eta_2) \sqrt{\frac{2Q}{\pi}} \left(1 + \frac{R}{\sqrt{Q}} \right) P^{T-1} - 2\eta_2 \sqrt{\frac{2Q}{\pi}} + \frac{(\eta_1^2 - \eta_2^2)}{T} P^T + \frac{\eta_2^2}{T}, \quad (6)$$

155 where $\alpha \equiv \mu/D$ serves as a continuous time variable in the limit $D \rightarrow \infty$ (not to be confused
 156 with t which counts episode steps), and $P = (1 - \cos^{-1}(R/\sqrt{Q})/\pi)$ is the probability of a single
 157 correct decision. While our derivation of the equations follow heuristics from statistical physics, we
 158 anticipate that their asymptotic correctness in the limit $D \rightarrow \infty$ can be established rigorously using
 159 the techniques of Goldt et al. [2019], Veiga et al. [2022], Arnaboldi et al. [2023]. We illustrate the
 160 accuracy of these equations already in finite dimensions ($D = 900$) in fig. 1c, where we show the
 161 expected reward, as well as the overlaps R and Q , of a student as measured during a simulation and
 162 from integration of the dynamical equations (solid and dotted lines, respectively).

163 The derivation of the dynamical equations that govern the learning dynamics of the RL perceptron
 164 are our first main result. Equipped with this tool, we now analyse several phenomena exhibited by
 165 the RL perceptron through a detailed study of these equations.

166 2.2 Learning protocols

167 The RL perceptron allows for the characterization of different RL protocols by adapting the reward
 168 condition Φ . We considered the following three settings:

169 **Vanilla:** The dynamics in the ‘standard’ case without penalty, $\eta_2 = 0$, is shown in fig. 5a and fig. 5b.
 170 Rewards are sparsest in this protocol, and as a result we observe a characteristic initial plateau in
 171 expected reward followed by a rapid jump. The length of this plateau increases with T , consistent
 172 with the notion that sparse rewards make exploration hard and slow learning [Bellemare et al., 2016].
 173 Plateaus during learning, which arise from saddle points in the loss landscape, have also been studied
 174 for (deep) neural networks in the supervised setting [Saad and Solla, 1995, Dauphin et al., 2014],
 175 but do not arise in the supervised perceptron. Hence the RL setting can qualitatively change the
 176 learning trajectory. The benefit of withholding penalties is that while slower, the perceptron reaches
 177 the highest level of expected reward in this case. This is a first example of a speed-accuracy trade-off
 178 that we will explore in more detail in section 2.5 and that we also found in our experiments with
 179 Bossfight in section 3.

180 **Penalty:** The initial plateau can be reduced by providing a penalty or negative reward ($\eta_2 > 0$) when
 181 the student fails in the task. This change provides weight updates much earlier in training and thus
 182 accelerates the escape from the plateau. The dynamics under this protocol are shown in fig. 2a. It is
 183 clear the penalty provides an initial speed-up in learning, as expected if the agent were to be unaligned
 184 and more likely to commit an error. However, a high penalty can create additional sub-optimal fixed
 185 points in the dynamics leading to a low asymptotic performance (more on this in section 2.4). In the
 186 simulations, finite size effects occasionally permit escape from the sub-optimal fixed point and jumps
 187 to the optimal one, leading to a high variance in the results.

188 **Subtask and breadcrumbs:** The model is also able to capture the dynamics of more complicated
 189 protocols: fig. 3b shows learning under the protocol where a smaller sub-reward is received if the
 190 agent survives beyond a shorter duration $T_0 < T$, i.e. some reward is still received even if the
 191 agent does not survive for the entire episode. Another learning protocol we can capture is that of
 192 ‘graded-breadcrumbs’, where the agent receives a small reward r_b for every correct decision made

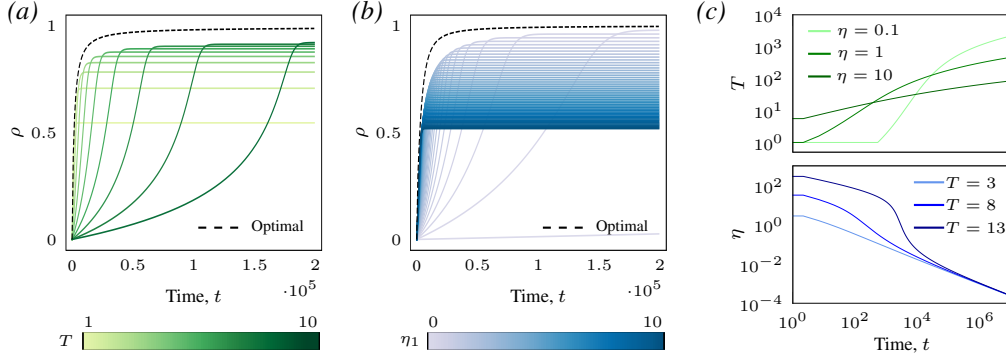


Figure 3: **Optimal schedules for episode length T and learning rate η .** (a) Evolution of the normalised overlap under optimal episode length scheduling (dashed) and various constant episode lengths (green). (b) Evolution of the normalised overlap under optimal learning rate scheduling (dashed) and various constant learning rates (blue). (c) Evolution of optimal T (green) and η (blue) over learning. *Parameters:* $D = 900$, $Q = 1$, $\eta_2 = 0$, (a) $\eta = 1$, (b) $T = 8$.

193 in an episode, i.e. like the previous method some reward is still received even if the agent does not
 194 survive for the entire episode, these dynamics are captured in fig. 3c.

195 2.3 Optimal hyper-parameter schedules: make episodes longer and anneal your learning rate

196 Hyper-parameter schedules are crucial for successful training of RL agents. In our setup, the two
 197 most important hyper-parameters are the learning rates and the episode length. In the RL perceptron,
 198 we can derive optimal schedules for both hyper-parameters. For simplicity, here we report the
 199 results in the spherical case, where the length of the student vector is fixed at \sqrt{D} (we discuss the
 200 unconstrained case in the appendix C), then $Q(\alpha) = 1$ at all times and we only need to track the
 201 teacher-student overlap $\rho = R/\sqrt{Q}$, which quantifies the generalisation performance of the agent.
 202 Keeping the choice $\mathbb{I}(\Phi) = \prod_{t=1}^T \theta(y_t y_t^*)$ and turning off the penalty term ($\eta_2 = 0$), we find that the
 203 teacher-student overlap is governed by the equation

$$\frac{d\rho}{d\alpha} = \frac{\eta}{\sqrt{2\pi Q}}(1 - \rho^2) \left(1 - \frac{1}{\pi} \cos^{-1}(\rho)\right)^{T-1} - \frac{\eta^2}{2TQ} \rho \left(1 - \frac{1}{\pi} \cos^{-1}(\rho)\right)^T \quad (7)$$

204 The optimal schedules over episodes for T and η can then be found by maximising the change
 205 in overlap at each update, i.e. setting $\frac{\partial}{\partial T} \left(\frac{d\rho}{d\alpha}\right)$ and $\frac{\partial}{\partial \eta} \left(\frac{d\rho}{d\alpha}\right)$ to zero respectively. After some
 206 calculations, we find the optimal schedules to be

$$T_{\text{opt}} = \left\lfloor \frac{\sqrt{\pi}}{2} \frac{\eta \rho P}{(1 - \rho^2) \sqrt{2Q}} \left[1 + \sqrt{1 - \frac{\sqrt{2Q}}{\eta \rho} \frac{4(1 - \rho^2)}{\sqrt{\pi} P \ln(P)}} \right] \right\rfloor \quad \text{and} \quad \eta_{\text{opt}} = \sqrt{\frac{Q}{2\pi}} \frac{T(1 - \rho^2)}{\rho P} \quad (8)$$

207 where $\lfloor \cdot \rfloor$ indicates the floor function.

208 Figure 3a shows the evolution of ρ under the optimal episode length schedule (dashed) compared to
 209 other constant episode lengths (green). Similarly, fig. 3b shows the evolution of ρ under the optimal
 210 learning rate schedule (dashed) compared to other constant learning rates (blue). The functional
 211 forms of T_{opt} and η_{opt} over time are shown in fig. 3c.

212 During learning the student seeks increasingly refined information to improve its expected reward.
 213 This simple observation explains the monotonic increase of the optimal episode length and the
 214 decrease in learning rates. Starting from the episode duration, we can observe that given the discrete
 215 nature of the decisions, information obtained from the rewards simply pushes the decision boundary
 216 towards a partition of the input space. This partition is determined by the episode length T and
 217 correspond to a fraction $1/2^T$ of the entire input space. Therefore a positive reward conveys T bits of
 218 information. At a fixed learning rate, when the student becomes proficient in the task it will not be
 219 able to improve further the decision boundary, and will fluctuate around the optimal solution unless
 220 longer episodes are provided.

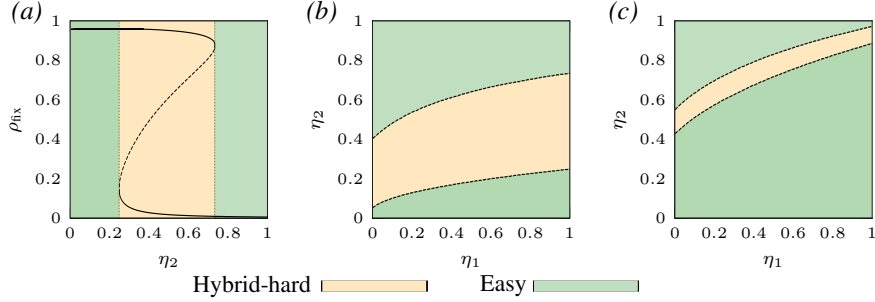


Figure 4: **Phase plots characterising learnability.** In the case where all decisions in an episode of length T must be correct in order to receive a reward. (a) the fixed points of ρ for $T = 13$ and $\eta_1 = 1$, the dashed portion of the line denotes where the fixed points are unstable. (b) Phase plot showing regions of hardness for $T = 13$. (c) Phase plot showing regions or hardness for $T = 8$. The green regions represent the *Easy* phase where with probability 1 the algorithm naturally converges to the optimal ρ_{fix} from a random initialisation. The orange region indicates the *Hybrid-hard* phase, where with high probability the algorithm converges to the sub-optimal ρ_{fix} from random initialisation. *Parameters:* $D = 900$, $Q = 1$.

221 Our analysis shows that a polynomial increase in the episode length gives the optimal performance
 222 in the RL perceptron, see fig. 3c (top); increasing T in the RL perceptron is akin to increasing task
 223 difficulty, and the polynomial scheduling of T_{opt} specifies a curriculum. Curricula of increasing task
 224 difficulty are commonly used in RL to give convergence speed-ups and learn problems that otherwise
 225 would be too difficult to learn *ab initio* Narvekar et al. [2020]. Analogously, the fluctuations can be
 226 reduced by annealing the learning rate and averaging over a larger number of samples. Akin to work in
 227 RL literature studying adaptive step-sizes [Dabney, 2014, Pirota et al., 2013], we find that annealing
 228 the learning rate during training is beneficial for greater speed and generalisation performance. For the
 229 RL perceptron, a polynomial decay in the learning rate gives optimal performance as shown in fig. 3c
 230 (bottom), consistent with work in the parallel area of high-dimensional non-convex optimization
 231 problems [d’Ascoli et al., 2022], and stochastic approximation algorithms in RL [Dalal et al., 2017].

232 2.4 Phase Space

233 With a non-zero penalty (η_2), the generalisation performance of the agent can enter different regimes
 234 of learning. This is most clearly exemplified in the spherical case, where the number of fixed points of
 235 the ODE governing the dynamics of the overlap exist in distinct phases determined by the combination
 236 of reward and penalty. For the simplest case ($\mathbb{I}(\Phi) = \prod_t^T (y_t y_t^*)$) these phases are shown in fig. 4.
 237 Figure 4a shows the fixed points achievable over a range of penalties for a fixed $\eta_1 = 1$ (obtained from
 238 a numerical solution of the ODE in ρ). There are two distinct regions: 1) *Easy*, where there is a unique
 239 fixed point and the algorithm naturally converges to this optimal ρ_{fix} from a random initialisation,
 240 2) a *Hybrid-hard* region (given the analogy with results from inference problems Ricci-Tersenghi
 241 et al. [2019]), where there are two stable (1 good and 1 bad) fixed points, and 1 unstable fixed point,
 242 and either stable point is achievable depending on the initialisation of the student (orange). The
 243 ‘hybrid-hard’ region separates two easy regions with very distinct performance levels. In this region
 244 the algorithm with high probability converges to ρ_{fix} with the worse performance level. These two
 245 regions are visualised in (η_1, η_2) space in fig. 4b for an episode length of $T = 13$. The topology
 246 of these regions are also governed by episode length, with a sufficiently small T reducing the the
 247 area of the ‘hybrid-hard’ phase to zero, meaning there is always 1 stable fixed point which may not
 248 necessarily give ‘good’ generalisation. Figure 4c shows the phase plot for $T = 8$, where the orange
 249 (hybrid-hard) has shrunk, this corresponds to the s-shaped curve in fig. 4a becoming flatter (closer
 250 to monotonic). Learning with η_2 This is not a peculiarity specific to the spherical case, indeed, we
 251 observe different regimes in the learning dynamics in the setting with unrestricted Q which we report
 252 in appendix C.

253 These phases show that at a fixed η_1 increasing η_2 will eventually lead to a first order phase transition,
 254 and the speed benefits gained from a non-zero η_2 will be nullified due to the transition into the
 255 hybrid-hard phase. In fact, when taking η_2 close to the transition point, instead of speeding up
 256 learning there is the presence of a critical slowing down, which we report in appendix C.

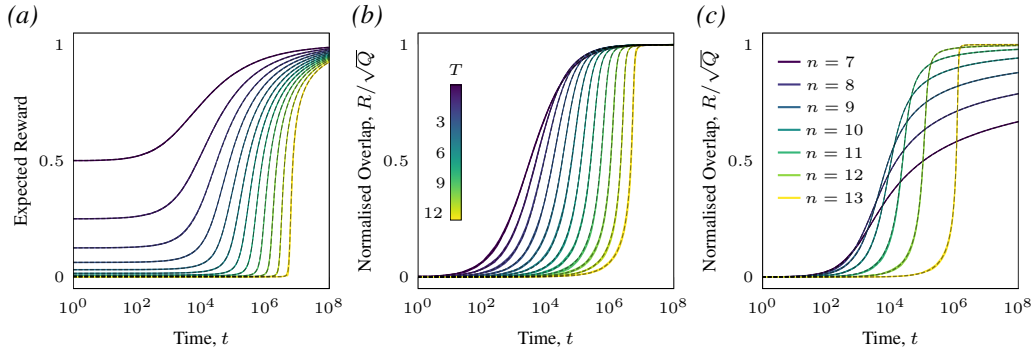


Figure 5: **Speed-accuracy tradeoff.** Evolution of (a) the expected reward and (b) corresponding normalised overlap for simulation (solid) and ODE solution (dashed) over a range of T when all decisions in an episode of length T are required correct, and $\eta_2 = 0$. (c) Evolution of the normalised overlap between student and teacher weights for simulation (solid) and ODE solution (dashed) for the case where n or more decisions in an episode of length 13 are required correct for an update with $\eta_2 = 0$. More stringent reward conditions slow learning but can improve performance. *Parameters:* $D = 900$, $\eta_1 = 1$, $\eta_2 = 0$.

257 A common problem with REINFORCE is high variance gradient estimates leading to bad
 258 performance [Marbach and Tsitsiklis, 2003, Schulman et al., 2015]. The reward (η_1) and punishment
 259 (η_2) magnitude alters the variance of the updates, and we show that the interplay between reward,
 260 penalty and reward-condition and their effect on performance can be probed within our model. This
 261 framework opens the possibility for studying phase transitions between learning regimes [Gamarnik
 262 et al., 2022].

263 2.5 Speed-accuracy trade-off

264 Figure 5c shows the evolution of normalised overlap $\rho = R/\sqrt{Q}$ between the student and teacher
 265 obtained from simulations and from solving the ODEs in the case where n or more decisions must
 266 be correctly made in an episode of length $T = 13$ in order to receive a reward (with $\eta_2 = 0$). We
 267 observe a speed-accuracy trade-off, where decreasing n increases the initial speed of learning but
 268 leads to worse asymptotic performance; this alleviates the initial plateau in learning seen previously
 269 in fig. 5b at the cost of good generalisation. In essence, a lax reward function is probabilistically
 270 more achievable early in learning; but it rewards some fraction of incorrect decisions, leading to
 271 lower asymptotic accuracy. By contrast a stringent reward function slows learning but eventually
 272 produces a highly aligned student. For a given MDP, it is known that arbitrary shaping applied
 273 to the reward function will change the optimal policy (reduce asymptotic performance) [Ng et al.,
 274 1999]. Empirically, reward shaping has been shown to speed up learning and help overcome difficult
 275 exploration problems [Gullapalli and Barto, 1992]. Reconciling these results with the phenomena
 276 observed in our setting is an interesting avenue for future work.

277 3 Experiments

278 To verify that our theoretical framework captures qualitative features of more general settings, we
 279 train agents from pixels on the Procgen [Cobbe et al., 2019] game ‘Bossfight’ (example frame, fig. 6a
 280 (top)). To remain close to our theoretical setting, we consider a modified version of the game where
 281 the agent cannot defeat the enemy and wins only if it survives for a given duration T . On each
 282 timestep the agent has the binary choice of moving left/right and aims to dodge incoming projectiles.
 283 We give the agent h lives, where the agent loses a life if struck by a projectile and continues an
 284 episode if it has lives remaining. This reward structure reflects the sparse reward setup from our
 285 theory and is analogous to requiring n out of T decisions to be correct within an episode. We further
 286 add asteroids at the left and right boundaries of the playing field which destroy the agent on contact,
 287 such that the agent cannot hide in the corners. Observations, shown in fig. 6a (bottom), are centred on
 288 the agent and downsampled to size 35×64 with three colour channels, yielding a 6720 dimensional
 289 input. The pixels corresponding to the agent are set to zero since these otherwise act as near-constant
 290 bias inputs not present in our model. The agent is endowed with a shallow policy network with

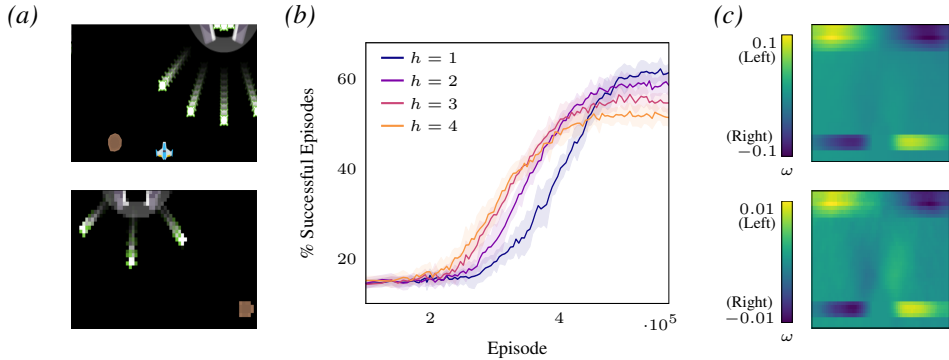


Figure 6: **Empirical speed-accuracy tradeoff in Bossfight.** (a) Top: Screenshot from a frame of ‘Bossfight.’ Bottom: Example observation provided to the agent’s policy network. In our variant, the agent can move left or right and aims to survive for a given duration T . Collision with projectiles or asteroids costs one life, and the agent has h lives before an episode terminates. (b) Performance during training, measured on evaluation episodes with $h = 3$ lives. Agents trained in stringent conditions ($h = 1$) learn slowly but eventually outperform agents trained in lax conditions ($h = 4$), an instance of the speed-accuracy tradeoff. Shaded regions indicate SEM over 10 repetitions. (c) Policy network weights for an agent with (top) $h = 4$ lives and (bottom) $h = 1$ life. For simplicity, one colour channel (red) is shown. Training with fewer lives increases the weight placed on dodging projectiles (see text). *Parameters:* $T = 100, \eta_1 = 8.2e - 5, \eta_2 = 0$.

291 logistic output unit that indicates the probability of left or right action. The weights of the policy
 292 network are trained using the policy gradient update of eq. (1) under a pure random policy.

293 To study the speed-accuracy trade-off, we train agents with different numbers of lives. As seen
 294 in fig. 6b, we observe a clear speed-accuracy trade-off mediated by agent health consistent with
 295 our theoretical findings (c.f. fig. 3c). Figure 6c shows the final policy weights for agents trained
 296 with $h = 1$ and $h = 4$. These show interpretable structure, roughly split into thirds vertically: the
 297 weights in the top third detect the position of the boss and centre the agent beneath it; this causes
 298 projectiles to arrive vertically rather than obliquely, making them easier to dodge. The weights
 299 in the middle third dodge projectiles. Finally, the weights in the bottom third avoid asteroids near
 300 the agent. Notably, the agent trained in the more stringent reward condition ($h = 1$) places greater
 301 weight on dodging projectiles, showing the qualitative impact of reward on learned policy. Hence
 302 similar qualitative phenomena as in our theoretical model can arise in more general settings.

303 4 Concluding perspectives

304 The RL perceptron provides a framework to investigate high-dimensional policy gradient learning in
 305 RL for a range of plausible sparse reward structures. We derive closed ODEs that capture the *average-*
 306 *case* learning dynamics in high-dimensional settings. The reduction of the high-dimensional learning
 307 dynamics to a low-dimensional set of differential equations permits a precise, quantitative analysis
 308 of learning behaviours: computing optimal hyper-parameter schedules, or tracing out phase diagrams
 309 of learnability. Our framework offers a starting point to explore additional settings that are closer
 310 to many real-world RL scenarios, such as those with conditional next states. Furthermore, the RL
 311 perceptron offers a means to study common training practices, including curricula; and more advanced
 312 algorithms, like actor-critic methods. We hope to extract more analytical insights from the ODEs,
 313 particularly on how initialization and learning rate influence an agent’s learning regime. Our findings
 314 emphasize the intricate interplay of task, reward, architecture, and algorithm in modern RL systems.

315 References

316 Alekh Agarwal, Sham Kakade, Akshay Krishnamurthy, and Wen Sun. Flambe: Structural complexity
 317 and representation learning of low rank mdps. *Advances in neural information processing systems*,
 318 33:20095–20107, 2020.

- 319 Alekh Agarwal, Sham M Kakade, Jason D Lee, and Gaurav Mahajan. On the theory of policy
320 gradient methods: Optimality, approximation, and distribution shift. *The Journal of Machine*
321 *Learning Research*, 22(1):4431–4506, 2021.
- 322 Andrea Agazzi and Jianfeng Lu. Global optimality of softmax policy gradient with single hidden layer
323 neural networks in the mean-field regime. In *International Conference on Learning Representations*,
324 2021. URL <https://openreview.net/forum?id=bB2drc7DPuB>.
- 325 Andrea Agazzi and Jianfeng Lu. Temporal-difference learning with nonlinear function approximation:
326 lazy training and mean field regimes. In Joan Bruna, Jan Hesthaven, and Lenka Zdeborova, editors,
327 *Proceedings of the 2nd Mathematical and Scientific Machine Learning Conference*, volume 145
328 of *Proceedings of Machine Learning Research*, pages 37–74. PMLR, 16–19 Aug 2022. URL
329 <https://proceedings.mlr.press/v145/agazzi22a.html>.
- 330 Luca Arnaboldi, Ludovic Stephan, Florent Krzakala, and Bruno Loureiro. From high-dimensional &
331 mean-field dynamics to dimensionless odes: A unifying approach to sgd in two-layers networks.
332 *arXiv preprint arXiv:2302.05882*, 2023.
- 333 Haruka Asanuma, Shiro Takagi, Yoshihiro Nagano, Yuki Yoshida, Yasuhiko Igarashi, and Masato
334 Okada. Statistical mechanical analysis of catastrophic forgetting in continual learning with teacher
335 and student networks. *Journal of the Physical Society of Japan*, 90(10):104001, 2021.
- 336 Alex Ayoub, Zeyu Jia, Csaba Szepesvari, Mengdi Wang, and Lin Yang. Model-based reinforcement
337 learning with value-targeted regression. In *International Conference on Machine Learning*, pages
338 463–474. PMLR, 2020.
- 339 Mohammad Gheshlaghi Azar, Ian Osband, and Rémi Munos. Minimax regret bounds for rein-
340 forcement learning. In *International Conference on Machine Learning*, pages 263–272. PMLR,
341 2017.
- 342 Yasaman Bahri, Jonathan Kadmon, Jeffrey Pennington, Sam S. Schoenholz, Jascha Sohl-Dickstein,
343 and Surya Ganguli. Statistical mechanics of deep learning. *Annual Review of Condensed Matter*
344 *Physics*, 11(1):501–528, 2020. doi: 10.1146/annurev-conmatphys-031119-050745. URL <https://doi.org/10.1146/annurev-conmatphys-031119-050745>.
- 346 Peter L Bartlett and Shahar Mendelson. Rademacher and gaussian complexities: Risk bounds and
347 structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- 348 Marc Bellemare, Sriram Srinivasan, Georg Ostrovski, Tom Schaul, David Saxton, and Remi Munos.
349 Unifying count-based exploration and intrinsic motivation. *Advances in neural information*
350 *processing systems*, 29, 2016.
- 351 Jalaj Bhandari and Daniel Russo. Global optimality guarantees for policy gradient methods. *arXiv*
352 *preprint arXiv:1906.01786*, 2019.
- 353 Michael Biehl and Holm Schwarze. Learning by on-line gradient descent. *Journal of Physics A:*
354 *Mathematical and general*, 28(3):643, 1995.
- 355 Qi Cai, Zhuoran Yang, Jason D Lee, and Zhaoran Wang. Neural temporal-difference learning
356 converges to global optima. *Advances in Neural Information Processing Systems*, 32, 2019.
- 357 Giuseppe Carleo, Ignacio Cirac, Kyle Cranmer, Laurent Daudet, Maria Schuld, Naftali Tishby, Leslie
358 Vogt-Maranto, and Lenka Zdeborová. Machine learning and the physical sciences. *Reviews of*
359 *Modern Physics*, 91(4):045002, 2019.
- 360 L. Chizat and F. Bach. On the global convergence of gradient descent for over-parameterized
361 models using optimal transport. In *Advances in Neural Information Processing Systems 31*, pages
362 3040–3050, 2018.
- 363 Lenaic Chizat and Francis Bach. Implicit bias of gradient descent for wide two-layer neural networks
364 trained with the logistic loss. In *Conference on Learning Theory*, pages 1305–1338. PMLR, 2020.

- 365 Lénaïc Chizat, Edouard Oyallon, and Francis Bach. On lazy training in differentiable programming.
366 In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché Buc, E. Fox, and R. Garnett, editors,
367 *Advances in Neural Information Processing Systems 32*, pages 2933–2943. Curran Associates, Inc.,
368 2019.
- 369 Karl Cobbe, Christopher Hesse, Jacob Hilton, and John Schulman. Leveraging procedural generation
370 to benchmark reinforcement learning. *arXiv preprint arXiv:1912.01588*, 2019.
- 371 William M. Dabney. Adaptive step-sizes for reinforcement learning. 2014.
- 372 Gal Dalal, Balázs Szörényi, Gagan Thoppe, and Shie Mannor. Concentration bounds for two timescale
373 stochastic approximation with applications to reinforcement learning. *CoRR*, abs/1703.05376,
374 2017. URL <http://arxiv.org/abs/1703.05376>.
- 375 Stéphane d’Ascoli, Maria Refinetti, and Giulio Biroli. Optimal learning rate schedules in high-
376 dimensional non-convex optimization problems, 2022.
- 377 Yann N Dauphin, Razvan Pascanu, Caglar Gulcehre, Kyunghyun Cho, Surya Ganguli, and Yoshua
378 Bengio. Identifying and attacking the saddle point problem in high-dimensional non-convex
379 optimization. *Advances in neural information processing systems*, 27, 2014.
- 380 Oussama Dhifallah and Yue M Lu. Phase transitions in transfer learning for high-dimensional
381 perceptrons. *Entropy*, 23(4):400, 2021.
- 382 Shi Dong, Benjamin Van Roy, and Zhengyuan Zhou. Provably efficient reinforcement learning with
383 aggregated states. *arXiv preprint arXiv:1912.06366*, 2019.
- 384 Simon Du, Akshay Krishnamurthy, Nan Jiang, Alekh Agarwal, Miroslav Dudik, and John Langford.
385 Provably efficient rl with rich observations via latent state decoding. In *International Conference*
386 *on Machine Learning*, pages 1665–1674. PMLR, 2019a.
- 387 S.S. Du, X. Zhai, B. Póczos, and A. Singh. Gradient descent provably optimizes over-parameterized
388 neural networks. In *International Conference on Learning Representations*, 2019b.
- 389 Andreas Engel and Christian Van den Broeck. *Statistical mechanics of learning*. Cambridge
390 University Press, 2001.
- 391 Marylou Gabrié, Surya Ganguli, Carlo Lucibello, and Riccardo Zecchina. Neural networks: from the
392 perceptron to deep nets. *arXiv preprint arXiv:2304.06636*, 2023.
- 393 David Gamarnik, Cristopher Moore, and Lenka Zdeborová . Disordered systems insights on compu-
394 tational hardness. *Journal of Statistical Mechanics: Theory and Experiment*, 2022(11):114015,
395 nov 2022. doi: 10.1088/1742-5468/ac9cc8. URL [https://doi.org/10.1088/1742-5468/](https://doi.org/10.1088/1742-5468/ac9cc8)
396 [2Fac9cc8](https://doi.org/10.1088/1742-5468/ac9cc8).
- 397 Elizabeth Gardner and Bernard Derrida. Three unfinished works on the optimal storage capacity of
398 networks. *Journal of Physics A: Mathematical and General*, 22(12):1983, 1989.
- 399 Federica Gerace, Luca Saglietti, Stefano Sarao Mannelli, Andrew Saxe, and Lenka Zdeborová.
400 Probing transfer learning with a model of synthetic correlated datasets. *Machine Learning: Science*
401 *and Technology*, 3(1):015030, 2022.
- 402 Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Limitations of lazy
403 training of two-layers neural network. In *Advances in Neural Information Processing Systems*,
404 volume 32, pages 9111–9121, 2019.
- 405 Behrooz Ghorbani, Song Mei, Theodor Misiakiewicz, and Andrea Montanari. When do neural
406 networks outperform kernel methods? In *Advances in Neural Information Processing Systems*,
407 volume 33, 2020.
- 408 Sebastian Goldt, Madhu Advani, Andrew M Saxe, Florent Krzakala, and Lenka Zdeborová. Dynamics
409 of stochastic gradient descent for two-layer neural networks in the teacher-student setup. *Advances*
410 *in neural information processing systems*, 32, 2019.

- 411 Vijaykumar Gullapalli and Andrew G Barto. Shaping as a method for accelerating reinforcement
412 learning. In *Proceedings of the 1992 IEEE international symposium on intelligent control*, pages
413 554–559. IEEE, 1992.
- 414 A. Jacot, F. Gabriel, and C. Hongler. Neural tangent kernel: Convergence and generalization in neural
415 networks. In *Advances in Neural Information Processing Systems 32*, pages 8571–8580, 2018.
- 416 Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E Schapire. Contextual
417 decision processes with low bellman rank are pac-learnable. In *International Conference on*
418 *Machine Learning*, pages 1704–1713. PMLR, 2017.
- 419 Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement
420 learning with linear function approximation. In *Conference on Learning Theory*, pages 2137–2143.
421 PMLR, 2020.
- 422 Chi Jin, Qinghua Liu, and Sobhan Miryoosefi. Bellman eluder dimension: New rich classes of rl
423 problems, and sample-efficient algorithms. *Advances in neural information processing systems*,
424 34:13406–13418, 2021.
- 425 W. Kinzel and P. Ruján. Improving a Network Generalization Ability by Selecting Examples. *EPL*
426 (*Europhysics Letters*), 13(5):473–477, 1990.
- 427 Robert Kirk, Amy Zhang, Edward Grefenstette, and Tim Rocktäschel. A Survey of Zero-shot
428 Generalisation in Deep Reinforcement Learning. *Journal of Artificial Intelligence Research*, 76:
429 201–264, January 2023. ISSN 1076-9757. doi: 10.1613/jair.1.14174. URL [https://jair.org/
430 index.php/jair/article/view/14174](https://jair.org/index.php/jair/article/view/14174).
- 431 Akshay Krishnamurthy, Alekh Agarwal, and John Langford. Pac reinforcement learning with rich
432 observations. *Advances in Neural Information Processing Systems*, 29, 2016.
- 433 Andrew K Lampinen and Surya Ganguli. An analytic theory of generalization dynamics and transfer
434 learning in deep linear networks. *arXiv preprint arXiv:1809.10374*, 2018.
- 435 Sebastian Lee, Sebastian Goldt, and Andrew Saxe. Continual learning in the teacher-student setup:
436 Impact of task similarity. In *International Conference on Machine Learning*, pages 6109–6119.
437 PMLR, 2021.
- 438 Sebastian Lee, Stefano Sarao Mannelli, Claudia Clopath, Sebastian Goldt, and Andrew Saxe.
439 Maslow’s hammer for catastrophic forgetting: Node re-use vs node activation. *arXiv preprint*
440 *arXiv:2205.09029*, 2022.
- 441 Peter Marbach and John N. Tsitsiklis. Approximate gradient methods in policy-space optimization of
442 markov reward processes. *Discrete Event Dynamic Systems*, 13:111–148, 2003.
- 443 Song Mei, Andrea Montanari, and Phan-Minh Nguyen. A mean field view of the landscape of two-
444 layer neural networks. *Proceedings of the National Academy of Sciences*, 115(33):E7665–E7671,
445 2018.
- 446 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare,
447 Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control
448 through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- 449 Aditya Modi, Nan Jiang, Ambuj Tewari, and Satinder Singh. Sample complexity of reinforcement
450 learning using linearly combined model ensembles. In *International Conference on Artificial*
451 *Intelligence and Statistics*, pages 2010–2020. PMLR, 2020.
- 452 Sanmit Narvekar, Bei Peng, Matteo Leonetti, Jivko Sinapov, Matthew E. Taylor, and Peter Stone.
453 Curriculum learning for reinforcement learning domains: A framework and survey, 2020.
- 454 Andrew Y Ng, Daishi Harada, and Stuart Russell. Policy invariance under reward transformations:
455 Theory and application to reward shaping. In *Icml*, volume 99, pages 278–287. Citeseer, 1999.

- 456 Matteo Pirotta, Marcello Restelli, and Luca Bascetta. Adaptive step-size for policy gradient
457 methods. In C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, ed-
458 itors, *Advances in Neural Information Processing Systems*, volume 26. Curran Associates,
459 Inc., 2013. URL [https://proceedings.neurips.cc/paper_files/paper/2013/file/
460 f64eac11f2cd8f0efa196f8ad173178e-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2013/file/f64eac11f2cd8f0efa196f8ad173178e-Paper.pdf).
- 461 Maria Refinetti, Sebastian Goldt, Florent Krzakala, and Lenka Zdeborova. Classifying high-
462 dimensional gaussian mixtures: Where kernel methods fail and neural networks succeed. In Marina
463 Meila and Tong Zhang, editors, *Proceedings of the 38th International Conference on Machine
464 Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8936–8947. PMLR,
465 18–24 Jul 2021. URL <https://proceedings.mlr.press/v139/refinetti21b.html>.
- 466 Federico Ricci-Tersenghi, Guilhem Semerjian, and Lenka Zdeborová. Typology of phase transitions
467 in bayesian inference problems. *Physical Review E*, 99(4):042109, 2019.
- 468 F. Rosenblatt. *Principles of Neurodynamics*. Spartan, New York, 1962.
- 469 G.M. Rotskoff and E. Vanden-Eijnden. Parameters as interacting particles: long time convergence
470 and asymptotic error scaling of neural networks. In *Advances in Neural Information Processing
471 Systems 31*, pages 7146–7155, 2018. URL <http://arxiv.org/abs/1805.00915>.
- 472 Daniel Russo and Benjamin Van Roy. Eluder dimension and the sample complexity of optimistic
473 exploration. *Advances in Neural Information Processing Systems*, 26, 2013.
- 474 David Saad and Sara A Solla. On-line learning in soft committee machines. *Physical Review E*, 52
475 (4):4225, 1995.
- 476 Luca Saglietti, Stefano Sarao Mannelli, and Andrew Saxe. An analytical theory of curriculum
477 learning in teacher–student networks. *Journal of Statistical Mechanics: Theory and Experiment*,
478 2022(11):114014, 2022.
- 479 John Schulman, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. High-dimensional
480 continuous control using generalized advantage estimation, 2015. URL [https://arxiv.org/
481 abs/1506.02438](https://arxiv.org/abs/1506.02438).
- 482 Hyunjune Sebastian Seung, Haim Sompolinsky, and Naftali Tishby. Statistical mechanics of learning
483 from examples. *Physical review A*, 45(8):6056, 1992.
- 484 David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche,
485 Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering
486 the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- 487 Ben Sorscher, Surya Ganguli, and Haim Sompolinsky. Neural representational geometry under-
488 lies few-shot concept learning. *Proceedings of the National Academy of Sciences*, 119(43):
489 e2200800119, 2022. doi: 10.1073/pnas.2200800119. URL [https://www.pnas.org/doi/abs/
490 10.1073/pnas.2200800119](https://www.pnas.org/doi/abs/10.1073/pnas.2200800119).
- 491 R. S. Sutton, D. Mcallester, S. Singh, and Y. Mansour. Policy gradient methods for reinforcement
492 learning with function approximation. In *Advances in Neural Information Processing Systems 12*,
493 volume 12, pages 1057–1063. MIT Press, 2000.
- 494 Vladimir N Vapnik and A Ya Chervonenkis. On the uniform convergence of relative frequencies of
495 events to their probabilities. *Measures of complexity: festschrift for alexey chervonenkis*, pages
496 11–30, 2015.
- 497 Rodrigo Veiga, Ludovic STEPHAN, Bruno Loureiro, Florent Krzakala, and Lenka Zdeborova.
498 Phase diagram of stochastic gradient descent in high-dimensional two-layer neural networks.
499 In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho, editors, *Advances in
500 Neural Information Processing Systems*, 2022. URL [https://openreview.net/forum?id=
501 GL-3WEdNRM](https://openreview.net/forum?id=GL-3WEdNRM).
- 502 Lin Yang and Mengdi Wang. Sample-optimal parametric q-learning using linearly additive features.
503 In *International Conference on Machine Learning*, pages 6995–7004. PMLR, 2019.

- 504 Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep
505 learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115,
506 2021.
- 507 Xuezhou Zhang, Yuda Song, Masatoshi Uehara, Mengdi Wang, Alekh Agarwal, and Wen Sun.
508 Efficient reinforcement learning in block mdps: A model-free representation learning approach. In
509 *International Conference on Machine Learning*, pages 26517–26547. PMLR, 2022.
- 510 Yufeng Zhang, Qi Cai, Zhuoran Yang, Yongxin Chen, and Zhaoran Wang. Can temporal-difference
511 and q-learning learn representation? a mean-field theory. In H. Larochelle, M. Ranzato, R. Hadsell,
512 M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33,
513 pages 19680–19692. Curran Associates, Inc., 2020a. URL [https://proceedings.neurips.
514 cc/paper_files/paper/2020/file/e3bc4e7f243ebc05d66a0568a3331966-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/e3bc4e7f243ebc05d66a0568a3331966-Paper.pdf).
- 515 Zihan Zhang, Yuan Zhou, and Xiangyang Ji. Almost optimal model-free reinforcement learning via
516 reference-advantage decomposition. *Advances in Neural Information Processing Systems*, 33:
517 15198–15207, 2020b.