

536 **Appendix**

537 **A General Convergence Considerations for MDPs in Finite State Space**

538 In this section, we will discuss the infinite batch limit and compare the value function obtained with
 539 TD to the ground truth value function. We will, for simplicity, consider in this section a Markov
 540 reward process with transition matrix $p(s_{t+1} = s' | s_t = s) = \Pi(s, s')$. The general theory described
 541 in the main text does not only apply to MDPs, but the convergence analysis for MDPs is much more
 542 straightforward so we describe it here. In this case, the ground truth value function satisfies

$$V(s) = R(s) + \gamma \sum_{s'} \Pi(s, s') V(s') \quad (\text{A.1})$$

543 which gives the vector equation $\mathbf{V} = (\mathbf{I} - \gamma \mathbf{\Pi})^{-1} \mathbf{R}$ for $\mathbf{V}, \mathbf{R} \in \mathbb{R}^{|\mathcal{S}|}$. Suppose the limiting
 544 distribution over states is $\mathbf{p} \in \mathbb{R}^{|\mathcal{S}|}$ which has entries $p(s) = \frac{1}{T} \sum_{t=1}^T p(s_t = s)$. The fixed point of
 545 TD dynamics is

$$\Psi \text{diag}(\mathbf{p}) \Psi^\top \mathbf{w}_{TD} = \Psi \text{diag}(\mathbf{p}) \mathbf{R} + \gamma \Psi \text{diag}(\mathbf{p}) \mathbf{\Pi} \Psi^\top \mathbf{w}_{TD}. \quad (\text{A.2})$$

546 We now consider the two possible cases for this fixed point condition.

547 **Case 1: Underparameterized Regime** First, if the feature dimension N is smaller than the size
 548 of the state space $|\mathcal{S}|$ and the features are maximal rank, then the TD learning fixed point is

$$\mathbf{w}_{TD} = (\Psi \text{diag}(\mathbf{p}) \Psi^\top - \gamma \Psi \text{diag}(\mathbf{p}) \mathbf{\Pi} \Psi^\top)^{-1} \Psi \text{diag}(\mathbf{p}) \mathbf{R} \quad (\text{A.3})$$

549 In this case, the value function is not learned perfectly, as can be seen by computing $\hat{\mathbf{V}} = \Psi^\top \mathbf{w}_{TD}$
 550 and comparing to the ground truth $\mathbf{V} = (\mathbf{I} - \gamma \mathbf{\Pi})^{-1} \mathbf{R}$. In this case, we would say that TD learning
 551 has an *irreducible value error* due to capturing only a N dimensional projection of the value function.

552 **Case 2: Overparameterized Regime** Alternatively, if the feature dimension exceeds the total
 553 number of states, then the fixed point equation for TD is underspecified. However, throughout TD
 554 learning $\mathbf{w}_{TD} \in \text{span}\{\psi(s)\}_{s \in \mathcal{S}}$ so we can instead consider the decomposition $\mathbf{w}_V = \sum_s \alpha(s) \psi(s)$,
 555 where $\alpha \in \mathbb{R}^{|\mathcal{S}|}$ satisfies

$$\text{diag}(\mathbf{p})(\mathbf{I} - \gamma \mathbf{\Pi}) \mathbf{K} \alpha = \text{diag}(\mathbf{p}) \mathbf{R} \quad (\text{A.4})$$

556 where $\mathbf{K} \in \mathbb{R}^{|\mathcal{S}| \times |\mathcal{S}|}$ is the kernel computed with features $K(s, s') = \psi(s) \cdot \psi(s')$. The solution to the
 557 above equation is unique and the learned value function $\hat{\mathbf{V}} = \Psi^\top \mathbf{w}_{TD} = \mathbf{K} \mathbf{K}^{-1} (\mathbf{I} - \gamma \mathbf{\Pi})^{-1} \mathbf{R} =$
 558 $(\mathbf{I} - \gamma \mathbf{\Pi})^{-1} \mathbf{R} = \mathbf{V}$. Therefore, in the over-parameterized limit, the irreducible value error for TD
 559 learning is zero. This limit was considered dynamically in the infinite batch (vanishing SGD noise)
 560 setting by [41].

561 **B Derivation of Learning Curves**

562 In this section, we now consider the dynamics of TD learning when B random episodes are sampled
 563 at a time. In this calculation, the finite batch of episodes leads to non-negligible SGD effects which
 564 can cause undesirable plateaus in TD dynamics.

565 **B.1 Field Theory Derivation**

566 In this section we use a Gaussian field theory formalism to compute the learning curve in the high
 567 dimensional asymptotic limit $N, B \rightarrow \infty$ with $B/N = \alpha$. The episode length T is treated as
 568 $\mathcal{O}(1)$. While this paper focuses on the online setting, where fresh trajectories $\{\tau_n^\mu\}$ are sampled at
 569 each iteration n , this model can be straightforwardly extended to the case where a fixed number of
 570 experience trajectories $\{\tau^\mu\}$ are replayed repeatedly during TD learning. We leave the experience

571 replay dynamic mean field theory calculation for future work. The starting point of our analysis is
 572 tracking the moment generating function for the iterate dynamics

$$Z[\{\mathbf{j}_n\}] = \mathbb{E}_{\{\mathbf{w}_n, \{s_n^\mu(t)\}\}} \exp\left(i \sum_{n=0}^{\infty} \mathbf{j}_n \cdot \mathbf{w}_n\right). \quad (\text{B.1})$$

573 To compute this object over random draws of training trajectories, we express the joint average over
 574 $\mathbf{w}_n, \{s_n^\mu(t)\}$ into conditional averages over $\mathbf{w}_n, \{\Delta_n^\mu(t)\} | \{\psi_n^\mu(t)\}$. To simplify the computation, in
 575 this section, we will compute the learning curve for mean zero features $\boldsymbol{\mu}(s) = 0$ and

$$\begin{aligned} Z = & \mathbb{E}_{\{\psi_n^\mu(t)\}} \int \prod_n d\mathbf{w}_n \delta\left(\mathbf{w}_{n+1} - \mathbf{w}_n - \frac{\eta}{\sqrt{BT}} \sum_{\mu t} \Delta_n^\mu(t) \psi_n^\mu(t)\right) \exp\left(i \sum_{n=0}^{\infty} \mathbf{j}_n \cdot \mathbf{w}_n\right) \\ & \times \int \prod_{t\mu n} d\Delta_n^\mu(t) \delta\left(\Delta_n^\mu(t) - \frac{1}{\sqrt{N}}(\mathbf{w}_R - \mathbf{w}_n) \cdot \psi_n^\mu(t) - \frac{\gamma}{\sqrt{N}} \mathbf{w}_n \cdot \psi_n^\mu(t+1)\right) \end{aligned} \quad (\text{B.2})$$

576 Expressing the Dirac-delta function as a Fourier integral $\delta(z) = \int \frac{d\hat{z}}{2\pi} \exp(i\hat{z}z)$ for each of our
 577 constraints. Under the *Gaussian equivalence ansatz*, we can easily average over Gaussian ψ to obtain

$$\begin{aligned} Z = & \int \mathcal{D}\Delta \mathcal{D}\hat{\Delta} \mathcal{D}\mathbf{w} \mathcal{D}\hat{\mathbf{w}} \exp\left(-\frac{\eta^2}{2BT^2} \sum_{n\mu} \sum_{t'} \Delta_n^\mu(t) \Delta_n^\mu(t') \hat{\mathbf{w}}_n^\top \boldsymbol{\Sigma}(t, t') \hat{\mathbf{w}}_n\right) \\ & \exp\left(i \sum_n \hat{\mathbf{w}}_n \cdot (\mathbf{w}_{n+1} - \mathbf{w}_n)\right) \\ & \exp\left(-\frac{1}{2N} \sum_{n\mu t t'} \left[(\mathbf{w}_R - \mathbf{w}_n) \hat{\Delta}_n^\mu(t)\right] \boldsymbol{\Sigma}(t, t') \left[(\mathbf{w}_R - \mathbf{w}_n) \hat{\Delta}_n^\mu(t')\right]\right) \\ & \exp\left(-\frac{\gamma^2}{2N} \sum_{n\mu t t'} \hat{\Delta}_n^\mu(t-1) \hat{\Delta}_n^\mu(t'-1) \mathbf{w}_n^\top \boldsymbol{\Sigma}(t, t') \mathbf{w}_n\right) \\ & \exp\left(-\frac{\gamma}{N} \sum_{n\mu t t'} \hat{\Delta}_n^\mu(t-1) \hat{\Delta}_n^\mu(t') \mathbf{w}_n^\top \boldsymbol{\Sigma}(t, t') (\mathbf{w}_R - \mathbf{w}_n)\right) \\ & \exp\left(-\frac{\eta}{\sqrt{N}BT} \sum_{n\mu t t'} \left[\hat{\Delta}_n^\mu(t) (\mathbf{w}_R - \mathbf{w}_n) + \gamma \hat{\Delta}_n^\mu(t-1) \mathbf{w}_n\right]^\top \boldsymbol{\Sigma}(t, t') \hat{\mathbf{w}}_n \Delta_n^\mu(t')\right) \\ & \exp\left(i \sum_{n\mu t} \hat{\Delta}_n^\mu(t) \Delta_n^\mu(t) + i \sum_n \mathbf{j}_n \cdot \mathbf{w}_n\right) \end{aligned} \quad (\text{B.3})$$

578 where we adopted the shorthand $\mathcal{D}\Delta = \prod_{\mu, n, t} d\Delta_n^\mu(t)$ for the measure for the collection of variables
 579 $\{\Delta_n^\mu(t)\}$. Likewise one should interpret $\mathcal{D}\mathbf{w} = \prod_n d\mathbf{w}_n$. To analyze the high dimensional limit of
 580 the above moment generating function, we introduce order parameters for the theory

$$\begin{aligned} Q_n(t, t') &= \frac{1}{B} \sum_{\mu=1}^B \Delta_n^\mu(t) \Delta_n^\mu(t'), \quad C_n(t, t') = \frac{1}{N} \mathbf{w}_n^\top \boldsymbol{\Sigma}(t, t') \mathbf{w}_n \\ C_n^R(t, t') &= \frac{1}{N} \mathbf{w}_R \boldsymbol{\Sigma}(t, t') \mathbf{w}_n, \quad D_n(t, t') = -\frac{i}{N} \hat{\mathbf{w}}_n^\top \boldsymbol{\Sigma}(t, t') \mathbf{w}_n, \quad D_n^R(t, t') = -\frac{i}{N} \hat{\mathbf{w}}_n^\top \boldsymbol{\Sigma}(t, t') \mathbf{w}_R \end{aligned} \quad (\text{B.4})$$

581 For each of these order parameters, we enforce the definition of the order parameter using the Fourier
 582 representation of a Dirac-delta function

$$\begin{aligned}
 1 &= B \int dQ_n(t, t') \delta \left(BQ_n(t, t') - \sum_{\mu} \Delta_n^{\mu}(t) \Delta_n^{\mu}(t') \right) \\
 &= B \int \frac{dQ_n(t, t') d\hat{Q}_n(t, t')}{4\pi i} \exp \left(\frac{B}{2} \hat{Q}_n(t, t') Q_n(t, t') - \frac{1}{2} \sum_{\mu} \Delta_n^{\mu}(t) \Delta_n^{\mu}(t') \hat{Q}_n(t, t') \right).
 \end{aligned} \tag{B.5}$$

583 Repeating this procedure for all order parameters $q = \{Q, \hat{Q}, C, \hat{C}, C^R, \hat{C}^R, D, \hat{D}, D^R, \hat{D}^R\}$ and
 584 disregarding irrelevant prefactors, we have the following formula for the moment generating function

$$Z \propto \int \mathcal{D}q \exp \left(\frac{N}{2} S[q] \right) \tag{B.6}$$

585 where the action S has the form

$$\begin{aligned}
 S &= \sum_n \sum_{tt'} \left[\alpha Q_n(t, t') \hat{Q}_n(t, t') + C_n(t, t') \hat{C}_n(t, t') + C_n^R(t, t') \hat{C}_n^R(t, t') \right] \\
 &\quad - 2 \sum_n \sum_{tt'} \left[D_n(t, t') \hat{D}_n(t, t') + D_n^R(t, t') \hat{D}_n^R(t, t') \right] + \frac{2}{N} \ln \mathcal{Z}_w + 2\alpha \ln \mathcal{Z}_{\Delta} \\
 \mathcal{Z}_w &= \int \mathcal{D}\mathbf{w} \mathcal{D}\hat{\mathbf{w}} \exp \left(-\frac{\eta^2}{2T^2} \sum_{ntt'} Q_n(t, t') \hat{\mathbf{w}}_n^{\top} \boldsymbol{\Sigma}(t, t') \hat{\mathbf{w}}_n + i \sum_n \hat{\mathbf{w}}_n \cdot (\mathbf{w}_{n+1} - \mathbf{w}_n) \right) \\
 &\quad \exp \left(-\frac{1}{2} \sum_{ntt'} \hat{C}_n(t, t') \mathbf{w}_n^{\top} \boldsymbol{\Sigma}(t, t') \mathbf{w}_n - \frac{1}{2} \hat{C}_n^R(t, t') \mathbf{w}_R^{\top} \boldsymbol{\Sigma}(t, t') \mathbf{w}_n \right) \\
 &\quad \exp \left(-i \sum_{ntt'} \hat{D}_n(t, t') \hat{\mathbf{w}}_n^{\top} \boldsymbol{\Sigma}(t, t') \mathbf{w}_n - i \sum_{ntt'} \hat{D}_n^R(t, t') \hat{\mathbf{w}}_n^{\top} \boldsymbol{\Sigma}(t, t') \mathbf{w}_R \right) \\
 \mathcal{Z}_{\Delta} &= \int \mathcal{D}\Delta \mathcal{D}\hat{\Delta} \exp \left(-\frac{1}{2} \sum_{ntt'} \hat{Q}_n(t, t') \Delta_n(t) \Delta_n(t') + i \sum_{nt} \hat{\Delta}_n(t) \Delta_n(t) \right) \\
 &\quad \exp \left(-\frac{1}{2} \sum_{ntt'} \hat{\Delta}_n(t) \hat{\Delta}_n(t') \left[\frac{1}{N} \mathbf{w}_R^{\top} \boldsymbol{\Sigma}(t, t') \mathbf{w}_R + C(t, t') \right] \right) \\
 &\quad \exp \left(\frac{1}{2} \sum_{ntt'} \hat{\Delta}_n(t) \hat{\Delta}_n(t') [C^R(t, t') + C^R(t', t)] \right) \\
 &\quad \exp \left(-\gamma \sum_{t, t'} \hat{\Delta}_n(t) \hat{\Delta}_n(t' - 1) C_n^R(t, t') \right) \\
 &\quad \exp \left(-\frac{\gamma^2}{2} \sum_{t, t'} \hat{\Delta}_n(t - 1) \hat{\Delta}_n(t' - 1) C_n(t, t') \right) \\
 &\quad \exp \left(-\frac{\eta i}{\sqrt{\alpha} T} \sum_{nt, t'} \hat{\Delta}_n(t) [D_n^R(t', t) - D_n(t', t) + \gamma D_n(t', t + 1)] \Delta_n(t') \right)
 \end{aligned} \tag{B.7}$$

586 The function \mathcal{Z} has the interpretation of an effective partition function conditional on order parameters
 587 q . To study the $N \rightarrow \infty$ limit, we use the steepest descent method and analyze the saddle point

588 $\frac{\partial S}{\partial q} = 0$. These saddle point equations give

$$\begin{aligned}
\frac{\partial S}{\partial \hat{Q}_n(t, t')} &= \alpha Q_n(t, t') - \alpha \langle \Delta_n(t) \Delta_n(t') \rangle = 0 \\
\frac{\partial S}{\partial Q_n(t, t')} &= \alpha \hat{Q}_n(t, t') - \frac{\eta^2}{T^2 N} \langle \hat{\mathbf{w}}_n^\top \boldsymbol{\Sigma}(t, t') \hat{\mathbf{w}}_n \rangle = 0 \\
\frac{\partial S}{\partial \hat{C}_n(t, t')} &= C_n(t, t') - \frac{1}{N} \langle \mathbf{w}_n^\top \boldsymbol{\Sigma}(t, t') \mathbf{w}_n \rangle = 0 \\
\frac{\partial S}{\partial C_n(t, t')} &= \hat{C}_n(t, t') - \alpha \langle \hat{\Delta}_n(t) \hat{\Delta}_n(t') + \gamma^2 \hat{\Delta}_n(t-1) \hat{\Delta}_n(t'-1) \rangle = 0 \\
\frac{\partial S}{\partial \hat{C}_n^R(t, t')} &= C_n^R(t, t') - \frac{1}{N} \langle \mathbf{w}_R^\top \boldsymbol{\Sigma}(t, t') \mathbf{w}_n \rangle = 0 \\
\frac{\partial S}{\partial C_n(t, t')} &= \hat{C}_n(t, t') - \alpha \langle \hat{\Delta}_n(t) \hat{\Delta}_n(t') + \gamma \hat{\Delta}_n(t) \hat{\Delta}_n(t'-1) \rangle = 0 \\
\frac{\partial S}{\partial \hat{D}_n(t, t')} &= -2D_n(t, t') - \frac{2i}{N} \langle \hat{\mathbf{w}}_n^\top \boldsymbol{\Sigma}(t, t') \mathbf{w}_n \rangle = 0 \\
\frac{\partial S}{\partial \hat{D}_n^R(t, t')} &= -2D_n^R(t, t') - \frac{2i}{N} \langle \hat{\mathbf{w}}_n^\top \boldsymbol{\Sigma}(t, t') \mathbf{w}_n \rangle = 0 \\
\frac{\partial S}{\partial D_n(t, t')} &= -2\hat{D}_n(t, t') - \frac{2\alpha\eta i}{\sqrt{\alpha}T} \langle \gamma \hat{\Delta}_n(t-1) \Delta_n(t') - \hat{\Delta}_n(t) \Delta_n(t') \rangle = 0 \\
\frac{\partial S}{\partial D_n^R(t, t')} &= -2\hat{D}_n^R(t, t') - \frac{2\alpha\eta i}{\sqrt{\alpha}T} \langle \hat{\Delta}_n(t) \Delta_n(t') \rangle = 0
\end{aligned} \tag{B.8}$$

589 The brackets $\langle \rangle$ denote averaging over the stochastic processes defined by moment generating
590 functions $\mathcal{Z}_\Delta, \mathcal{Z}_w$. After these saddle point equations are solved the order parameters q are treated as
591 non-random and a Hubbard-Stratonovich transformation is employed. For example,

$$\exp \left(-\frac{1}{2} \hat{\mathbf{w}}_n \left[\frac{\eta^2}{T^2} \sum_{tt'} Q_n(t, t') \boldsymbol{\Sigma}(t, t') \right] \hat{\mathbf{w}}_n \right) = \mathbb{E}_{\mathbf{u}_n^w} \exp \left(i \sum_n \mathbf{u}_n^w \cdot \hat{\mathbf{w}}_n \right) \tag{B.9}$$

592 where the average is over $\mathbf{u}_n^w \sim \mathcal{N}(0, \eta^2 T^{-2} \sum_{tt'} Q_n(t, t') \boldsymbol{\Sigma}(t, t'))$. After introducing these
593 Hubbard fields \mathbf{u}_n^w and $u_n^\Delta(t)$, we can perform the integrals over $\hat{\mathbf{w}}_n$ and $\hat{\Delta}_n(t)$ which collapse to
594 Dirac-Delta functions. The resulting identities of the delta functions define the following stochastic
595 processes on \mathbf{w}_n and u_n^Δ

$$\begin{aligned}
\mathbf{w}_{n+1} &= \mathbf{w}_n + \mathbf{u}_n^w + \sum_{tt'} \hat{D}_n^R(t, t') \boldsymbol{\Sigma}(t, t') \mathbf{w}_R + \sum_{t, t'} \hat{D}_n(t, t') \boldsymbol{\Sigma}(t, t') \mathbf{w}_n \\
\Delta_n(t) &= u_n^\Delta(t) + \frac{\eta}{\sqrt{\alpha}T} \sum_{tt'} [D_n^R(t, t') - D_n(t, t') - \gamma D_n(t', t+1)] \Delta_n(t').
\end{aligned} \tag{B.10}$$

596 Using a similar trick, we can show that for any observable depending on \mathbf{w}_n or $\{\Delta_n(t)\}$ that

$$\begin{aligned}
-i \langle \hat{\mathbf{w}}_n O(\mathbf{w}_n) \rangle &= \left\langle \frac{\partial}{\partial \mathbf{u}_n^w} O(\mathbf{w}_n) \right\rangle \\
-i \langle \hat{\Delta}_n(t) O(\{\Delta_n(t')\}) \rangle &= \left\langle \frac{\partial}{\partial u_n^\Delta(t)} O(\{\Delta_n(t')\}) \right\rangle
\end{aligned} \tag{B.11}$$

597 Since \mathbf{w}_n is independent. This can be used to conclude

$$D_n(t, t') = 0, \quad D_n^R(t, t') = 0 \tag{B.12}$$

598 which implies that $\Delta_n(t) = u_n^\Delta(t)$. Consequently the response functions have trivial structure

$$\hat{D}_n(t) = -\frac{\eta\sqrt{\alpha}}{T} [\delta(t-t') - \gamma\delta(t-1-t')] , \quad \hat{D}_n^R(t, t') = \frac{\sqrt{\alpha}\eta}{T} \delta(t-t'). \tag{B.13}$$

599 We therefore obtain a stochastic process of the form

$$\begin{aligned}
\mathbf{w}_{n+1} &= \mathbf{w}_n + \mathbf{u}_n^w + \frac{\eta\sqrt{\alpha}}{T} \sum_t \boldsymbol{\Sigma}(t, t) \mathbf{w}_R - \frac{\eta\sqrt{\alpha}}{T} \sum_t [\boldsymbol{\Sigma}(t, t) - \gamma \boldsymbol{\Sigma}(t, t+1)] \mathbf{w}_n \\
\mathbf{u}_n &\sim \mathcal{N}\left(0, \frac{\eta^2}{T^2} \sum_{tt'} Q_n(t, t') \boldsymbol{\Sigma}(t, t')\right), \quad \{\Delta_n(t)\} \sim \mathcal{N}(0, \mathbf{Q}_n) \\
Q_n(t, t') &= \langle \Delta_n(t) \Delta_n(t') \rangle = \frac{1}{N} \mathbf{w}_R^\top \boldsymbol{\Sigma}(t, t') \mathbf{w}_R - C^R(t, t') - C^R(t', t) + C(t, t') \\
C_n(t, t') &= \frac{1}{N} \langle \mathbf{w}_n^\top \boldsymbol{\Sigma}(t, t') \mathbf{w}_n \rangle, \quad C_n^R(t, t') = \frac{1}{N} \langle \mathbf{w}_R^\top \boldsymbol{\Sigma}(t, t') \mathbf{w}_n \rangle
\end{aligned}$$

600 These are the final equations defining the stochastic evolution of \mathbf{w}_n and $\Delta_n(t)$.

601 B.2 Simplifying the Saddle Point Equations

602 Using the above saddle point equations, we see that the variables $\{\Delta_n(t)\}$ and $\{\mathbf{w}_n\}$ will be Gaussian
603 random variables. It thus suffices to track their mean and covariance. The $\{\Delta_n(t)\}$ variables have
604 zero mean and covariance given by the $Q_n(t, t')$ function. The $\{\mathbf{w}_n\}$ variables have the following
605 mean evolution

$$\begin{aligned}
\langle \mathbf{w}_{n+1} \rangle &= \langle \mathbf{w}_n \rangle + \eta\sqrt{\alpha} [\bar{\boldsymbol{\Sigma}} \mathbf{w}_R - [\bar{\boldsymbol{\Sigma}} - \gamma \bar{\boldsymbol{\Sigma}}_+] \langle \mathbf{w}_n \rangle] \\
&= \langle \mathbf{w}_n \rangle + \eta\sqrt{\alpha} [\bar{\boldsymbol{\Sigma}} - \gamma \bar{\boldsymbol{\Sigma}}_+] [\mathbf{w}_{TD} - \langle \mathbf{w}_n \rangle]
\end{aligned} \tag{B.14}$$

606 where $\mathbf{w}_{TD} = [\bar{\boldsymbol{\Sigma}} - \gamma \bar{\boldsymbol{\Sigma}}_+]^{-1} \bar{\boldsymbol{\Sigma}} \mathbf{w}_R$ is the fixed point of the TD dynamics. We next compute
607 $\mathbf{M}_n = \langle (\mathbf{w}_n - \mathbf{w}_{TD})(\mathbf{w}_n - \mathbf{w}_{TD})^\top \rangle$ which admits the recursion

$$\mathbf{M}_{n+1} = (\mathbf{I} - \eta\sqrt{\alpha} [\bar{\boldsymbol{\Sigma}} - \gamma \bar{\boldsymbol{\Sigma}}_+]) \mathbf{M}_n (\mathbf{I} - \eta\sqrt{\alpha} [\bar{\boldsymbol{\Sigma}} - \gamma \bar{\boldsymbol{\Sigma}}_+]) + \frac{\eta^2}{T^2} \sum_{tt'} Q_n(t, t') \boldsymbol{\Sigma}(t, t') \tag{B.15}$$

608 To obtain our formulas which hold for finite batch size, we rescale the learning rate by $\eta \rightarrow \eta/\sqrt{\alpha}$
609 giving the following evolution

$$\begin{aligned}
\langle \mathbf{w}_{n+1} \rangle &= \langle \mathbf{w}_n \rangle + \eta [\bar{\boldsymbol{\Sigma}} - \gamma \bar{\boldsymbol{\Sigma}}_+] [\mathbf{w}_{TD} - \langle \mathbf{w}_n \rangle] \\
\mathbf{M}_{n+1} &= (\mathbf{I} - \eta [\bar{\boldsymbol{\Sigma}} - \gamma \bar{\boldsymbol{\Sigma}}_+]) \mathbf{M}_n (\mathbf{I} - \eta [\bar{\boldsymbol{\Sigma}} - \gamma \bar{\boldsymbol{\Sigma}}_+])^\top + \frac{\eta^2}{T^2 \alpha^2} \sum_{tt'} Q_n(t, t') \boldsymbol{\Sigma}(t, t') \tag{B.16}
\end{aligned}$$

610 After this rescaling, we see that the mean evolution for \mathbf{w}_n is independent of α but that the variance
611 picks up an additive term on each step on the order of $\mathcal{O}(\eta^2 \alpha^{-2})$ which vanishes in the infinite batch
612 limit $B/N \rightarrow \infty$. The error for value learning can be obtained from \mathbf{M}_n with $\mathcal{L}_n = \frac{1}{N} \text{Tr} \mathbf{M}_n \bar{\boldsymbol{\Sigma}}$.
613 Lastly, we note that we can express the formula for $Q_n(t, t')$ entirely in terms of \mathbf{M}_n and $\langle \mathbf{w}_n \rangle$. This

614 gives the lengthy expression

$$\begin{aligned}
Q_n(t, t') &= \frac{1}{N} \langle (\mathbf{w}_R - \mathbf{w}_n)^\top \boldsymbol{\Sigma}(t, t') (\mathbf{w}_R - \mathbf{w}_n) \rangle + \frac{\gamma}{N} \langle (\mathbf{w}_R - \mathbf{w}_n)^\top \boldsymbol{\Sigma}(t, t' + 1) \mathbf{w}_n \rangle \\
&+ \frac{\gamma}{N} \langle \mathbf{w}_n^\top \boldsymbol{\Sigma}(t + 1, t') (\mathbf{w}_R - \mathbf{w}_n) \rangle + \frac{\gamma^2}{N} \langle \mathbf{w}_n^\top \boldsymbol{\Sigma}(t + 1, t' + 1) \mathbf{w}_n \rangle \\
&= \frac{1}{N} \text{Tr} \mathbf{M}_n \boldsymbol{\Sigma}(t, t') + \frac{1}{N} (\mathbf{w}_{TD} - \langle \mathbf{w}_n \rangle) [\boldsymbol{\Sigma}(t, t') + \boldsymbol{\Sigma}(t', t)] (\mathbf{w}_R - \mathbf{w}_{TD}) \\
&+ \frac{1}{N} (\mathbf{w}_R - \mathbf{w}_{TD})^\top \boldsymbol{\Sigma}(t, t') (\mathbf{w}_R - \mathbf{w}_{TD}) \\
&- \frac{\gamma}{N} \text{Tr} \mathbf{M}_n [\boldsymbol{\Sigma}(t, t' + 1) + \boldsymbol{\Sigma}(t + 1, t')] \\
&+ \frac{\gamma}{N} (\mathbf{w}_{TD} - \langle \mathbf{w}_n \rangle) [\boldsymbol{\Sigma}(t, t' + 1) + \boldsymbol{\Sigma}(t + 1, t')] \mathbf{w}_{TD} \\
&+ \frac{\gamma}{N} (\mathbf{w}_R - \mathbf{w}_{TD})^\top [\boldsymbol{\Sigma}(t, t' + 1) + \boldsymbol{\Sigma}(t + 1, t')] \langle \mathbf{w}_n \rangle \\
&+ \frac{\gamma^2}{N} \text{Tr} \mathbf{M}_n \boldsymbol{\Sigma}(t + 1, t' + 1) + \frac{2\gamma^2}{N} (\langle \mathbf{w}_n \rangle - \mathbf{w}_{TD}) \boldsymbol{\Sigma}(t + 1, t' + 1) \mathbf{w}_{TD} \\
&+ \frac{\gamma^2}{N} \mathbf{w}_{TD}^\top \boldsymbol{\Sigma}(t + 1, t' + 1) \mathbf{w}_{TD}
\end{aligned} \tag{B.17}$$

615 B.3 Final Result

616 Below we state in compact form the full final result for our TD learning curves. The below equations
617 give the evolution of the first and second moments of \mathbf{w}_n obtained from the mean-field density of the
618 previous section. Concretely, these moments obey dynamics

$$\begin{aligned}
\langle \mathbf{w}_{n+1} \rangle &= \langle \mathbf{w}_n \rangle + \eta [\bar{\boldsymbol{\Sigma}} - \gamma \bar{\boldsymbol{\Sigma}}_+] [\mathbf{w}_V - \langle \mathbf{w}_n \rangle] \\
\mathbf{M}_{n+1} &= [\mathbf{I} - \eta \bar{\boldsymbol{\Sigma}} + \eta \gamma \bar{\boldsymbol{\Sigma}}_+] \mathbf{M}_n [\mathbf{I} - \eta \bar{\boldsymbol{\Sigma}} + \eta \gamma \bar{\boldsymbol{\Sigma}}_+]^\top + \frac{\eta^2}{\alpha^2 T^2} \sum_{tt'} Q_n(t, t') \boldsymbol{\Sigma}(t, t') \\
Q_n(t, t') &= \frac{1}{N} \langle (\mathbf{w}_R - \mathbf{w}_n)^\top \boldsymbol{\Sigma}(t, t') (\mathbf{w}_R - \mathbf{w}_n) \rangle + \frac{\gamma}{N} \langle (\mathbf{w}_R - \mathbf{w}_n)^\top \boldsymbol{\Sigma}(t, t' + 1) \mathbf{w}_n \rangle \\
&+ \frac{\gamma}{N} \langle \mathbf{w}_n^\top \boldsymbol{\Sigma}(t + 1, t') (\mathbf{w}_R - \mathbf{w}_n) \rangle + \frac{\gamma^2}{N} \langle \mathbf{w}_n^\top \boldsymbol{\Sigma}(t + 1, t' + 1) \mathbf{w}_n \rangle.
\end{aligned} \tag{B.18}$$

619 These equations can be solved iteratively for $\bar{\mathbf{w}}_n$, \mathbf{M}_n , Q_n . Finite dimensional versions of this result
620 can be obtained by replacing α with B/N as written in the main text. The value estimation error is

$$\mathcal{L}_n = \frac{1}{N} \text{Tr} \mathbf{M}_n \bar{\boldsymbol{\Sigma}}. \tag{B.19}$$

621 B.4 Non-Zero Mean Feature

622 We can also simply modify the DMFT equations if the mean feature is nonvanishing $\boldsymbol{\mu}(s) \neq 0$. In this
623 case, when averaging over all possible trajectories through state space, there is a mean feature vector
624 at each episodic time $\boldsymbol{\mu}(t)$. The above equations are exact for non-zero mean features if $\boldsymbol{\Sigma}(t, t')$ is
625 regarded as the (non-centered) correlation matrix $\langle \boldsymbol{\psi}(t) \boldsymbol{\psi}(t') \rangle$.

626 B.5 Tracking Iterate Moments with Direct Recurrence Relation

627 In this section we give a direct calculation of the first two moments of \mathbf{w} over the collection of
628 randomly sampled features $\{\boldsymbol{\psi}_n^\mu(t)\}$ and show which terms can be disregarded.

629 Letting $\mathbf{A} = \bar{\boldsymbol{\Sigma}} - \gamma \bar{\boldsymbol{\Sigma}}_+$, we note that the average evolution of \mathbf{w} has the form

$$\langle \mathbf{w}_{n+1} \rangle = (\boldsymbol{\Sigma} - \gamma \boldsymbol{\Sigma}_+) (\mathbf{w}_{TD} - \langle \mathbf{w}_n \rangle) \tag{B.20}$$

630 Thus, if we disregarded fluctuations in \mathbf{w}_n due to SGD, the model will converge to the correct fixed
631 point. Next, we look at $\mathbf{M}_n = \langle (\mathbf{w}_n - \mathbf{w}_{TD}) (\mathbf{w}_n - \mathbf{w}_{TD}) \rangle$. Under the Gaussian equivalence

632 ansatz, we have

$$\begin{aligned}
\mathbf{M}_{n+1} &= \mathbf{M}_n - \eta \mathbf{A} \mathbf{M}_n - \eta \mathbf{M}_n \mathbf{A}^\top + \frac{\eta^2}{T^2 B^2} \sum_{\mu\nu tt'} \langle \Delta_n^\mu(t) \Delta_n^\nu(t') \boldsymbol{\psi}_n^\mu(t) \boldsymbol{\psi}_n^\nu(t') \rangle \\
&= (\mathbf{I} - \eta \mathbf{A}) \mathbf{M}_n (\mathbf{I} - \eta \mathbf{A})^\top - \frac{\eta^2}{B} \mathbf{A} \mathbf{M}_n \mathbf{A}^\top + \frac{\eta^2}{T^2 B} \sum_{tt'} \langle \Delta_n(t) \Delta_n(t') \boldsymbol{\psi}(t) \boldsymbol{\psi}(t')^\top \rangle \\
&= (\mathbf{I} - \eta \mathbf{A}) \mathbf{M}_n (\mathbf{I} - \eta \mathbf{A})^\top + \frac{\eta^2}{T^2 B} \sum_{tt'} Q_n(t, t') \boldsymbol{\Sigma}(t, t') \\
&\quad + \frac{\eta^2}{T^2 B} \sum_{tt'} \langle \Delta_n(t') \boldsymbol{\psi}(t) \rangle \langle \Delta_n(t) \boldsymbol{\psi}(t')^\top \rangle
\end{aligned} \tag{B.21}$$

633 The mean field theory derived from saddle point integration consists of the first two terms in the
634 final expression. Therefore mean field theory disregards the last term which computes cross time
635 correlations of RPEs with features, effectively making the approximation

$$\frac{\eta^2}{T^2 B} \sum_{tt'} \langle \Delta_n(t') \boldsymbol{\psi}(t) \rangle \langle \Delta_n(t) \boldsymbol{\psi}(t')^\top \rangle \approx 0. \tag{B.22}$$

636 After making this approximation, we recover the learning curve obtained in the previous Section B.3.
637 We show in our experiments that dropping this term does not significantly alter the learning curves.

638 B.6 Scaling of Asymptotic Fixed Points

639 To identify fixed points in the value error dynamics, we can seek non-vanishing fixed points for the
640 weight error covariance $\mathbf{M} = \langle (\mathbf{w} - \mathbf{w}_{TD})(\mathbf{w} - \mathbf{w}_{TD}) \rangle$. We note that $\langle \mathbf{w} \rangle \sim \mathbf{w}_{TD}$ asymptotically.
641 Again, letting $\mathbf{A} = \bar{\boldsymbol{\Sigma}} - \gamma \bar{\boldsymbol{\Sigma}}_+$, we obtain the following fixed point condition for \mathbf{M} under these
642 assumptions

$$\begin{aligned}
\mathbf{A} \mathbf{M} + \mathbf{M} \mathbf{A}^\top - \eta \mathbf{A} \mathbf{M} \mathbf{A}^\top &= \frac{\eta}{B T^2} \sum_{tt'} Q(t, t') \boldsymbol{\Sigma}(t, t') \\
Q(t, t') &= \text{Tr} \mathbf{M} \boldsymbol{\Sigma}(t, t') - \gamma \text{Tr} \mathbf{M} [\boldsymbol{\Sigma}(t, t' + 1) + \boldsymbol{\Sigma}(t + 1, t')] + \gamma^2 \text{Tr} \mathbf{M} \boldsymbol{\Sigma}(t + 1, t' + 1) \\
&\quad + \gamma^2 \mathbf{w}_{TD}^\top \bar{\boldsymbol{\Sigma}}^{-1} \bar{\boldsymbol{\Sigma}}_+ \boldsymbol{\Sigma}(t, t') \bar{\boldsymbol{\Sigma}}_+ \bar{\boldsymbol{\Sigma}}^{-1} \mathbf{w}_{TD} + \gamma^2 \mathbf{w}_{TD}^\top \boldsymbol{\Sigma}(t + 1, t' + 1) \mathbf{w}_{TD} \\
&\quad + \gamma^2 \mathbf{w}_{TD}^\top \bar{\boldsymbol{\Sigma}}^{-1} \bar{\boldsymbol{\Sigma}}_+ [\boldsymbol{\Sigma}(t, t' + 1) + \boldsymbol{\Sigma}(t + 1, t')] \mathbf{w}_{TD}.
\end{aligned} \tag{B.23}$$

643 Where we used the formula for $Q_n(t, t')$ from Appendix B.5, evaluated at $\langle \mathbf{w} \rangle = \mathbf{w}_{TD}$ and used the
644 fact that $\mathbf{w}_R = \mathbf{w}_{TD} - \gamma \bar{\boldsymbol{\Sigma}}^{-1} \bar{\boldsymbol{\Sigma}}_+ \mathbf{w}_{TD}$. The solution $\mathbf{M} = 0$ is a valid fixed point for \mathbf{M} in the
645 $\eta \rightarrow 0$ and $B \rightarrow \infty$ limits because the constant terms on the right-hand side vanish. Similarly, if
646 $\gamma = 0$ (which corresponds to the standard supervised learning case), the right hand side is linear in
647 \mathbf{M} , allowing $\mathbf{M} = 0$ to be a valid fixed point.

648 However, for finite B and non-zero η and γ , there exists a solution to the above fixed point equation.
649 For small $\frac{\eta \gamma^2}{B}$, we can easily deduce that \mathbf{M} must satisfy a self-consistent asymptotic scaling of the
650 form

$$\mathbf{M} = \mathcal{O} \left(\frac{\eta \gamma^2}{B} \right) \tag{B.24}$$

651 implying an asymptotic value error scaling of $\mathcal{L} \sim \text{Tr} \mathbf{M} \bar{\boldsymbol{\Sigma}} \sim \mathcal{O} \left(\frac{\eta \gamma^2}{B} \right)$. These scalings are examined
652 in Figure 3 where experiments obey the expected behavior.

653 C Reward Shaping

654 In this section, we consider the role of reward shaping on the dynamics of TD learning. As discussed
655 in the main text, we consider potential based shaping with potential function decomposable in the
656 features $\phi(s) = \mathbf{w}_\phi \cdot \boldsymbol{\psi}(s)$. We first describe the change to the average weight evolution $\langle \mathbf{w}_n \rangle$ and

657 then describe the dynamics of the correlations. In potential based shaping, the TD errors take the
658 form

$$\Delta(t) = R(s(t)) + \phi(s(t)) - \gamma\phi(s(t+1)) + \gamma\hat{V}(s(t+1)) - \hat{V}(s(t)) \quad (\text{C.1})$$

659 Computing from the DMFT equations the evolution of $\langle \mathbf{w}_n \rangle$ we have

$$\begin{aligned} \langle \mathbf{w}_{n+1} \rangle &= \langle \mathbf{w}_n \rangle + \eta \bar{\Sigma} (\mathbf{w}_R + \mathbf{w}_\phi - \langle \mathbf{w}_n \rangle) + \gamma \eta \bar{\Sigma}_+ (\langle \mathbf{w}_n \rangle - \mathbf{w}_\phi) \\ &= \langle \mathbf{w}_n \rangle - \eta \mathbf{A} [\mathbf{w}_{TD} + \mathbf{w}_\phi - \langle \mathbf{w}_n \rangle]. \end{aligned} \quad (\text{C.2})$$

660 We see that including the reward shaping function ϕ offsets the fixed point of the algorithm to be
661 $\mathbf{w}_{TD} + \mathbf{w}_\phi$. This occurs precisely because the potential-based reward shaping generates an additive
662 correction to the target value function by $\phi(s)$ [61]. When we predict value at evaluation, we use the
663 reshifted value $\hat{V}(s) - \phi(s)$. The natural quantity to track at the level of the mean field equations is
664 the adapted version of M_n

$$\mathbf{M}_n = \left\langle (\mathbf{w}_n - \mathbf{w}_{TD} - \mathbf{w}_\phi) (\mathbf{w}_n - \mathbf{w}_{TD} - \mathbf{w}_\phi)^\top \right\rangle. \quad (\text{C.3})$$

665 This correlation matrix has dynamics

$$\mathbf{M}_{n+1} = (\mathbf{I} - \eta \mathbf{A}) \mathbf{M}_n (\mathbf{I} - \eta \mathbf{A})^\top + \frac{\eta^2}{BT^2} \sum_{tt'} Q_n(t, t') \Sigma(t, t') \quad (\text{C.4})$$

666 and the TD-error correlations $Q_n(t, t')$ have the form

$$\begin{aligned} Q_n(t, t') &= \left\langle (\mathbf{w}_R + \mathbf{w}_\phi - \mathbf{w}_n)^\top \Sigma(t, t') (\mathbf{w}_R + \mathbf{w}_\phi - \mathbf{w}_n) \right\rangle \\ &\quad + \gamma \left\langle (\mathbf{w} - \mathbf{w}_\phi)^\top [\Sigma(t, t') + \Sigma(t', t)] (\mathbf{w}_R + \mathbf{w}_\phi - \mathbf{w}_n) \right\rangle \\ &\quad + \gamma^2 \left\langle (\mathbf{w}_n - \mathbf{w}_\phi)^\top \Sigma(t+1, t'+1) (\mathbf{w}_n - \mathbf{w}_\phi) \right\rangle \end{aligned} \quad (\text{C.5})$$

667 The value estimation error is again $\mathcal{L}_n = \text{Tr} \mathbf{M}_n \bar{\Sigma}$. We see that the two primary ways that reward
668 shaping alters the loss dynamics is

- 669 • A change in the initial condition for M_n to be $M_0 = (\mathbf{w}_{TD} + \mathbf{w}_\phi)(\mathbf{w}_{TD} + \mathbf{w}_\phi)^\top$
- 670 • A change in the TD error covariance term $Q_n(t, t')$

671 Both effects can generate significant changes in the dynamics and plateaus of the model.

672 D Numerical methods and additional details

673 The code to generate the Figures is provided in the Supplementary Material as a Jupyter Notebook.
674 Here, we briefly highlight some of the parameter choices.

675 For Figures 3 and 4 we use diagonally decoupled, but temporally correlated power law features
676 with $\Sigma_{k\ell}(t, t') = \delta_{k\ell} k^{-1.2} \exp(-|t - t'|/\tau_k)$ with $\tau_k = \frac{10}{k+1}$ and $w_k^R = k^{-1.1}$ for $k \in [N]$ with
677 $N = 300$. This type of feature structure is especially easy to evaluate the theoretical learning curves
678 for. Unless otherwise stated, these figures used $\gamma = 0.9$ and batch size $B = 10$.

679 For the 2D MDP grid world, we defined a discrete set of states on a 17×17 grid. The agent starts in
680 the middle position and follows a random diffusion policy where each possible movement (up, down,
681 left, right) is taken with equal probability. The features were generated as bell-shaped place cells
682 (shown). We computed $\Sigma(t, t')$ for the theory by sampling 5000 random draws of length $T = 50$.
683 The Gaussian learning curve is obtained with TD learning with $\psi_G \sim \mathcal{N}(0, \Sigma)$.

684 Numerical experiments were performed on a NVIDIA SMX4-A100-80GB GPU. Together numerical
685 experiments (both preliminary experiments and those presented in the paper) took less than 1 hour of
686 compute time.