

Raja Grewal

raja_grewal1@pm.me

<https://github.com/r-grewal>

Compiled on 02:02 (AEST) Saturday 14th August, 2021.

This document is a work in progress.

Abstract

Deep reinforcement learning is rapidly advancing the training of immortal agents that will inevitably continue to surpass all human experts in an ever-increasing number of intellectually intensive activities. This accelerating growth has led to the wide-spread adoption of numerous entrenched practices throughout the literature irrespective of algorithm, domain, and environment. In this work we directly examine four deeply held convictions for which we can construct alternatives using model-free, off-policy, learning in continuous action spaces.

The first is the ubiquitous use of MSE loss functions for critic evaluation as there exists a myriad of alternatives. We find that its use is acceptable as a starting point but its choice becomes another tuneable hyperparameter with potential to accelerate performance. Research into asymptotic convergence has also highlighted the inadequacy of the Monte-Carlo approach in aggregating mini-batch losses where the obtained empirical mean may not be appropriate if the underlying distribution contains rare, but very large outliers, known as fat tails. Their presence is gauged if the tail exponent is less than unity which might imply the true mean is formally undefined, in which case, a heuristically estimated shadow mean could be a more accurate representation. We find in all cases that the tail exponents are strictly less than unity, however, we are unable to present an alternative except for stating the Monte-Carlo mean must underestimate the true mean. Furthermore, we extend the use of multi-step returns to continuous domains where it behaves in line with all expectations and becomes yet another hyperparameter to tune. This also sets the stage to next examine whether its puzzling coupling to the size of the experience replay buffer is also present. The code is publicly available at <https://github.com/r-grewal/nonergodic-rl> with preliminary results in Sections 9-10 of the report at <https://github.com/r-grewal/data5709/blob/main/documents>.

Since its inception, reinforcement learning has been formulated to maximise the sum of future discounted rewards. Pioneering work on decision theory has shown that this approach is not suitable for the overwhelming majority of real-world domains where rewards are multiplicative and scale proportionally with the amount existing cumulative previously obtained rewards. For this dynamic, losses have an asymmetrically larger effect on performance compared to equal magnitude gains. Attention then shifts to identifying the singular path that maximises growth with probabilities of outcomes becoming irrelevant. To accommodate such environments, we show Q-learning, both stochastic and deterministic policy gradient theorems, soft Q-learning, and soft-policy iteration remain relatively unchanged, though very likely will lead to completely different optimisation priorities. Many of these results are also implicitly valid in a subset of non-stationary regimes incorporating all MDPs. The agent now aims to maximise the future discounted compounding return through the avoidance of steep losses. This theoretical development requires rigorous experimental verification starting with the creation of environments that accommodate its principles.

Keywords Reinforcement Learning · Model-Free · Off-Policy · Critic Loss Evaluation · Extreme Value Theory · Multi-Step Returns · Ergodicity · Multiplicative Dynamics · Non-Markovian Decision Processes

Acknowledgments The Sydney Informatics Hub and the University of Sydney's high performance computing cluster, Artemis, for providing the computing resources that have contributed to the results reported herein.

Notation The agents sequential interaction with the environment is characterised by a tuple representing the history $h_t \equiv s_1 a_1 r_1 \dots s_{t-1} a_{t-1} r_{t-1} s_t$ where the trajectory develops as $s_1 \rightarrow a_1 \rightarrow r_1 s_2 \rightarrow a_2 \rightarrow r_2 s_3 \rightarrow \dots \rightarrow r_{t-1} s_t$.

Contents

Abstract	ii
1 Can You Access Yourself in Infinite Parallel Universes?	1
1.1 Path Preferences with Identical Expectations	1
1.2 Optimal Leverage for a Simple Gamble	3
1.3 Conflating Probabilities with Payoffs	12
1.4 Implications and Anecdotal Outcomes	13
1.5 A Permanent Solution	18
2 Introduction	20
3 Background	24
3.1 Preliminaries	24
3.2 Policy Gradient Theorems	24
3.3 Actor-Critic Methods	25
3.4 Soft Actor-Critic	26
3.5 Robust Critic Evaluation	28
3.6 Preasymptotics and the Tail Exponent	29
3.7 Multi-Step Returns	37
3.8 History Decision Processes	37
3.9 Agent Performance Evaluation	39
4 Non-Ergodicity in Reward Accumulation	41
4.1 Ergodic Special Case	43
5 Q-Learning with Multiplicative Dynamics	44
5.1 Model-Free Return Maximisation	44
5.2 Proof of Convergence and Uniqueness	47
5.3 Clipped Double Q-Learning	51
5.4 Multi-Step Targets	52
6 Policy Gradients with Multiplicative Dynamics	55
6.1 Stochastic Actors	55
6.2 Deterministic Actors	57
7 Maximum Causal Entropy with Multiplicative Dynamics	60
7.1 Soft Learning	60
8 Energy Efficient Agent Inference	64
8.1 Multi-Stage Policy Control	66
9 Related Work	68

10 Experiments	70
References	71
A Agent Algorithms: SAC and TD3	81
B Assorted Applications	84
B.1 Overview	84
B.2 Robotic Control for Medical Surgery	85
B.3 Supply Chain Management	86
B.4 Guidance Systems	86
B.5 Education	87
B.6 Portfolio Management	87
B.7 Systems Control	89
C Taming Hubris	90

1 Can You Access Yourself in Infinite Parallel Universes?

Changes in wealth, health, and the life of any random individual or institution are more often than not compounding, and are represented by relative changes between values using rates of returns. The overwhelming majority of contemporary decision theory is however formulated under the assumption of fixed (additive) absolute changes in values [1–6]. This approach is invalid when the true nature of the environment is multiplicative where percentage decreases in value are not reverted by equal percentage increases and vice versa. Designing of optimal risk-taking strategies in these domains requires the abandonment of deeply entrenched methods involving the maximisation of ‘expectation’ values calculated using a probabilistic approach. Instead, the focus shifts towards finding the optimal path that maximises return under the constraint of avoiding steep losses, quantified by the time-average growth rate.

In this section we outline two motivating examples. The first provides a very simple introduction to why the path taken to identical expectation values even under certainty matters when examined using a multiplicative lens. The second illustrates a simple gamble that an investor behaving perfectly in tune with existing decision theory is virtually certain to go bankrupt if they take a probabilistic approach calculating expected returns at each time step when determining how to allocate risky capital. We discuss these in the context of how a random singular individual would behave if they could replay events multiple times or have only chance to pick a path. Then we present a more formal explanation of why payoffs are more important than probabilities. The section concludes with a discussion on implications and motivations behind this work.

1.1 Path Preferences with Identical Expectations

Assume we start with \$100 at time t_0 , if we were to first gain \$50 at t_1 and then lose \$50 at t_2 , and then vice versa with equal probability, under additive dynamics, the expectation value would remain \$100 at all times. The existing field of economic decision-making dictates that a ‘rational’ person should be indifferent between both paths. In reality if you were to go outside and randomly ask people without formal training in STEM or economics which path they would prefer, would you get a 50-50 split? Chances are you would not, as the overwhelming majority of the general public would prefer to first gain \$50 rather than lose \$50 [7–9]. Not satisfied with this empirical fact and similarly puzzling inconsistencies, the fields of behavioural economics and behavioural finance were concocted to explain many of these experimental realities by introducing a large variety of cognitive ‘biases’ that ultimately resulted in labelling the person on the street as “irrational”.

Let us now revisit both paths of this scenario using multiplicative dynamics that aim to maximise compounding growth while avoiding steep losses. The first situation with gaining \$50 then losing \$50 can be represented as $100\% \cdot (1 + 50\%) \cdot (1 - 0.33\%) = 100\%$ where the left hand side indicates change in value from the previous value. The second path is $100\% \cdot (1 - 50\%) \cdot (1 + 100\%) = 100\%$. Both paths yield the same final (nill) change in value. Visually this is shown in Fig. 1. Which one avoids large crashes? The first path increases value by 50% and then suffers a -33% loss. The second path starts with a -50% loss and then requires a 100% gain (doubling) to get back where we started. Assuming equal time intervals for both paths, as $\delta t = t_2 - t_0 \rightarrow 0$ the preference is likely to be identical as the expectation approach suggests. However, as δt becomes large, empirical evidence conclusively suggest preference for the path with the smaller drawdown [7–9]. The explanation for this that the -50% takes us far closer to complete irrecoverable ruin (of having \$0) and correspondingly we would need a doubling (100%) of value to get back to where we started. This appears to be a concrete art of human nature honed over millenniums of evolution

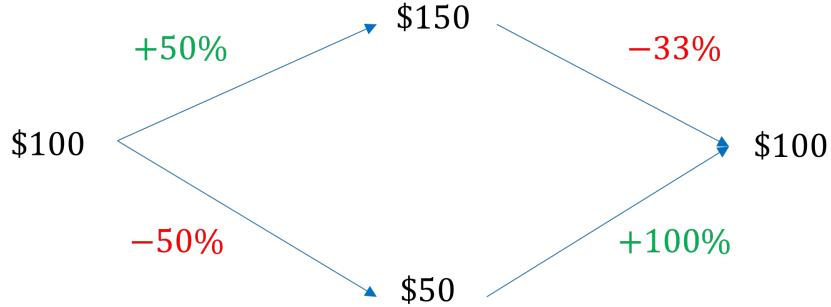


Figure 1: Visual depiction of both paths of the gamble with unchanged expectation values and noticing that to recover from the initial steep loss, a doubling of value is required.

where we rightfully concern ourselves with payoffs, or prefer “cash in hand”, even though the final outcome remains unchanged.

If we were to modify the absolute change $\$50 \rightarrow \0 identical preferences would result as the difference between the compounding percentages for both paths would reduce to be negligible. Likewise, as $\$50 \rightarrow \100 there would be even stronger preference for the first path. Therefore, while the two paths are exactly equivalent using additive dynamics, under multiplicative dynamics there can exist distinctly superior paths. Preference is then for the path that maximises reward under the constraint of avoiding ruin for an individual that only has once chance to decide on which direction to take. This is opposed to having the luxury of being able to repeat the event multiple times and taking the averaged outcome, obtained by uniformly averaging the outcome of each repetition, that is, the Monte Carlo approach which is assumed to converge in the infinite repetition limit to the ‘expectation’ value as the repeated trials should generate the underlying probability distribution.

Notice one curious fact with multiplicative dynamics, at no point did we use probabilities when deciding which path is optimal. This is because probabilities do not matter when a singular individual or institution evaluates which path to take on a gamble if the payoff is non-linear which we will show in the next section. In these situations, probabilities are completely artificial constructions based on what would result if we could replay an event from the start an infinite number of times, with expectation value not necessarily having any connection to the final result a random singular entity achieves. For these situations, the only thing that matters is the payoff structure and all strategies should be constructed based on where in a random world, a random singular individual can conceivably end up in the infinite time limit.

The label “expectation value” $\mathbb{E}[x]$ is hence a misnomer, it should be better named “ensemble average” $\langle x \rangle$ more in line with its actual origin in statistical mechanics [1, 4, 9–12]. Therefore, in gambles where leverage can be controlled to either amplify and reduce the payoff, calculating probabilities and hence expectation values will be shown to not only be entirely meaningless, but an absolutely absurd activity for a random individual. This is because we are solely concerned with how a random individuals valuation changes over time, not the average final valuations across all individuals. The former is what any sensible person would consider their performance to resemble, the latter is very prone to severe overestimation bias.

1.2 Optimal Leverage for a Simple Gamble

A standard example used to highlight failure of conventional wisdom when attempting to maximise reward signals using expectations involves a betting biased coin [4, 5]. This also highlights that the Gamblers Paradox is not really a paradox as it incorrectly assumes we have access and can pool returns from infinite versions of ourselves. A corollary is that observable world class gambling performance obtained via maximising expectation values is really nothing more than the living embodiment of survivorship bias, that is, the overwhelming majority of attempts to achieve such success through this method ends in failure, though there exists a very small minority who succeed due to pure ‘luck’. Therefore, for any random investor, utilising this approach is a fool’s errand.

Consider a gamble where a ‘rational’ investor initially has a portfolio of worth $V_0 = \$100$ and has the option to place a leveraged bet on a coin flip. The possible changes in wealth for this game are held constant across time where the payoffs at each flip are

$$V_{t+1} = \begin{cases} (1 + 50\%)l_t V_t + (1 - l_t)V_t = (1 + 0.5 \cdot l_t)V_t, & p_u = \frac{1}{2} \\ (1 - 40\%)l_t V_t + (1 - l_t)V_t = (1 - 0.4 \cdot l_t)V_t, & p_d = \frac{1}{2} \end{cases} \quad (1)$$

and $l_t \in \mathbb{R}$ is the amount of leverage the investor can apply and either outcome is equally probable. This represents an option since the investor can select $l_t = 0$ and so takes no risk at time t . Note this clearly indicates multiplicative dynamics as at each time step the worth of the portfolio increases by ratio rather than an additive fixed amount. The returns on investment (payoffs) are then

$$r_{t+1} = \begin{cases} +0.5 \cdot l_t, & p_u = \frac{1}{2} \\ -0.4 \cdot l_t, & p_d = \frac{1}{2} \end{cases} \quad (2)$$

Given the payoffs remain unchanged and the game can be played indefinitely, the sequence of investment management decisions is then:

1. Should investors consider playing this game at $t = 0$?
 2. If so, what direction should they place their bet, long ($l_0 > 0$) or short ($l_0 < 0$)?
 3. How much leverage l_0 should they apply?
 4. Should investors consider playing this game at $t = 1$?
- ⋮

This either continues ad infinitum ($t \rightarrow \infty$) or when the game ends where the bounds are $V_t \in [0, \infty)$. The lower bound is strict as once bankrupt the game forcibly concludes and there is no possibility of recovery. Let us also define a ‘stop-loss’ V_{\min} where if $V_\tau \leq V_{\min}$ the investor stops playing at time $t = \tau$ to live to fight another day. The exact value of $V_{\min} = \lambda V_0$ will be dependent on an investors preference for $\lambda \in [0, 1)$ and so we can say $V_t \in (V_{\min} - \varepsilon_\tau, \infty)$ where the lower bound is $V_\tau = V_{\min} - \varepsilon_\tau$ is time dependent. Notice that for some $l_t \in [l_{\min}, l_{\max}]$ where l_{\min} and l_{\max} are such that it impossible to go bankrupt as at worst $V_t \rightarrow V_{\min}$, the investor is technically always in the game and can ‘bounce back’ if the favourable state occurs sufficiently often as $t \rightarrow \infty$. Similarly, for leverages outside this bound, losses can not only exceed the predefined stop-loss, but also the initial deposit leaving the investor in debt.

How would a ‘rational’ investor then answer these questions? They would calculate the expectation value of course! This yields $\mathbb{E}[r_{t+1}] = rl_t = 5\% \cdot l_t \forall t$. If they are feeling especially fancy they would also calculate the ‘standard’ deviation of returns $\sigma_{t+1} = 45\% \cdot l_t \forall t$. In this special case $\sigma_{t+1} = \text{MAD}_{t+1}$ as there are only two components. MAD is discussed in detail in Section 3.9. Therefore, since $\mathbb{E}[r_t] \neq 0$, the investor should consider playing the game at all times. Next, as $\mathbb{E}[r_t] > 0$ the investor would clearly consider going long with $l_t > 0$ for all time. The next question is by far the most important, how should they calculate $l_t \forall t$?

To answer this we consider four types of rational investors increasing in their level of sophistication.

- (i) Investor 1: This first class of investors that only undertake in favourable bets and select leverage $|l_t| \leq 1$ based on maximising return while avoiding bankruptcy with certainty. Therefore, they have $|l_t| \in (0, 1]$, and because $V_t > 0 \forall t$ they will have $V_{\min} = 0$ since they can always ‘bounce back’. As these investors believe this is a favourable bet there is no reason it should not be played indefinitely, however for computational purposes we truncate time to a fixed maximum investment horizon T . A single investors compounding return is then

$$V_T = V_0(1 + \bar{g})^T = V_0 \prod_{t=0}^T (1 + rl_t) \quad (3)$$

where V_T is their final portfolio value and the time-averaged growth rate is \bar{g} . For a fixed leverage $l_t = l$, they ‘expect’ $\bar{g} = rl$ in the long-term as the horizon $T \rightarrow \infty$. Therefore, the choice of leverage appears to be trick question as this appears to be a situation where more leverage is superior and having a stop-loss is not necessary since the game can always be played. The supposed optimal leverage is then simply $l^* = 1$.

- (ii) Investor 2: These investors improve on the first by applying the optimal time-dependent leverage l_t^* while always maintaining a fixed $V_{\min} = \lambda V_0$ for $\lambda \in [0, 1)$ to keep their maximum loss capped at all times. By re-balancing at each step, these investors ensure optimal risk-reward while not being stopped-out. To prevent losses from exceeding the stop-loss we institute a constraint on leverage, where for $|l_t| > 1$ the change in portfolio value V_t from a single time step to V_{t+1} can be very sizeable. If the investor is long and wrong, we have danger threshold $l_t > \frac{5}{2} \left(1 - \frac{V}{V_t}\right) > 0$, if short and wrong, the threshold is $l_t < 2 \left(1 - \frac{V}{V_t}\right) < 0$, where $V = V_{\min}$ to exceed stop-loss and $V = 0$ to exceed deposit. Therefore for a rational investor not to be stopped out in a single step, they should have time-dependent leverage bounds

$$l_t \in \left(-2 \left(1 - \frac{V_{\min}}{V_t}\right), \frac{5}{2} \left(1 - \frac{V_{\min}}{V_t}\right)\right) \quad (4)$$

This shows that as $V_t \rightarrow \infty$ they are able to take substantially more risk at a maximum of $l_t \rightarrow \frac{5}{2}$ than if $V_t \rightarrow V_{\min}$. These investors have therefore successfully constructed open bounds on the leverage they can take at any point. Notice very importantly that the bounds are independent of the probability of either outcome, the only thing that matters are the final payoffs. The question still remains, on what exactly the leverage should be when taking the bet. Since the payoffs are fixed, we formally arrive at the optimisation problem

$$l_t^* = \arg \max_{l_t} V_{t-1}(1 + rl_t) = \arg \max_{l_t} l_t \quad (5)$$

under the constraint of Eq. (4) that must be solved at each step. The solution to this is clearly to approach $l_t^* \rightarrow \frac{5}{2} \left(1 - \frac{V_{\min}}{V_t}\right)$ since $r > 0$, that is, take the maximum amount of leverage possible while ensuring they are

not stopped out. We also assume leverage is available at no additional cost. These investors therefore simply calculate l_t^* at each step and bet $l_t^* V_{t-1}$ portion of their existing portfolio at each step.

- (iii) Investor 3: Improving on the second, these investors apply the same leverage principles but institute a ‘rolling’ stop-loss wherein

$$V_{\min,t} = \begin{cases} \lambda V_0, & \text{if } V_t \leq V_0 \\ V_0 + \phi(V_t - V_0), & \text{otherwise} \end{cases} \quad (6)$$

attempting to retain a fixed portion $\phi \in [0, 1)$ of their winning at each time step if they are profitable. The general leverage bounds for any fixed binary payout are then

$$l_t \in \left(-\frac{1}{|r^u|} \left(1 - \frac{V_{\min,t}}{V_t} \right), \frac{1}{|r^d|} \left(1 - \frac{V_{\min,t}}{V_t} \right) \right) \quad (7)$$

where the r^u and r^d are the static up and down returns, and the result is again independent of probabilities. For optimal leverage, if $|r_u| > |r_d|$ the upper bound is approached and vice versa. Through this type of judicious risk management, they expect that after they cross the $V_t > V_0$ threshold, they will constantly keep build their portfolio wealth as time goes on since this is favourable bet. For the proposed gamble, half of investors should start at this capital accumulation stage while a portion of the others will also eventually reach this stage.

- (iv) Investor 4: The final category of risk-takers involves those for which everything is a variable that they optimise at each time step. Through experience and instinct, they will select optimal the $\lambda_t^*, \phi_t^*, l_t^*$ to maximise returns, expressed very generally at each step as

$$\lambda_t^*, \phi_t^*, l_t^* \in \arg \max_{\lambda_t, \phi_t, l_t} V_{t-1} (1 + g_t(\lambda_t, \phi_t, l_t)) \quad (8)$$

with it not being clear whether a unique solution exists. The growth rate g_t is the key variable that incorporates changes in wealth. Explicitly, from initiation $t = 0$ the problem is to obtain the set variables

$$\begin{aligned} \lambda_1^*, \phi_1^*, l_1^*, \dots, \lambda_T^*, \phi_T^*, l_T^* &\in \arg \max_{\lambda_1, \phi_1, l_1, \dots, \lambda_T, \phi_T, l_T} \prod_{t=1}^T (1 + g_t(\lambda_t, \phi_t, l_t)) \\ &\in \arg \max_{\lambda_1, \phi_1, l_1, \dots, \lambda_T, \phi_T, l_T} \sum_{t=1}^T \ln |1 + g_t(\lambda_t, \phi_t, l_t)| \\ &\in \arg \max_{\lambda_1, \phi_1, l_1, \dots, \lambda_T, \phi_T, l_T} \sum_{t=1}^T (1 + g_t(\lambda_t, \phi_t, l_t)) \end{aligned} \quad (9)$$

since the logarithm is monotonic. The final performance is then measured by the time-averaged (exponential) growth rate \bar{g} defined by

$$\ln |1 + \bar{g}| = \frac{1}{T} \ln \left| \frac{V_T}{V_0} \right| \rightarrow 1 + \bar{g} = \exp \left[\frac{1}{T} \ln \left| \frac{V_T}{V_0} \right| \right] \quad (10)$$

which has no explicit dependence on any variable, incorporates all effects, and is a clear measure of average change in wealth achieved per time step. To select the optimal set from Eq. (9) across all strategies S , the superior set of variables are determined by strategies $S_i \in S$ where $\bar{g}_i \in \max_{S_j \in S} \bar{g}_j$ which need not be unique.

Identically we can identify the optimal set by looking at the final value $\bar{g}_i \in \max_{S_j \in S} V_T^j$. We can also write the time-average as $\bar{g} = \sqrt[T]{V_T/V_0} - 1$ which can be numerically difficult to evaluate for large T .

Another way to represent this is in terms of maximising the future exponential growth rate at each time step

$$g_{t+1} = \frac{\Delta \ln V_t}{\Delta t} \quad (11)$$

which has reduced computational feasibility but can be used for simpler gambles. Next, given that in general $\bar{g} \neq \mathbb{E}[r]$, the difference between them is expressed as

$$\bar{g} = \mathbb{E}[r] + \nu(\lambda, \phi, l, \sigma) \quad (12)$$

where $\nu(\cdot)$ is referred to as the ‘volatility tax’ [13–19]. This can be interpreted as the reward or punishment associated with certain variables. Clearly for any gamble where $0 \leq \sigma_1 < \sigma_2$, $0 \leq l_1 < l_2$, $0 \leq \lambda_1 < \lambda_2$, and $0 \leq \phi_1 < \phi_2$, we have the following general relations

$$\nu(\sigma_2) < \nu(\sigma_1) \leq 0, \quad \nu(l_2) < \nu(l_1) \leq 0 \quad (13)$$

$$\nu(\lambda_2) > \nu(\lambda_1) \geq 0, \quad \nu(\phi_2) > \nu(\phi_1) \geq 0 \quad (14)$$

The first two are expected as large variability in returns due to either the inherent gamble or increasing leverage negatively effects \bar{g} . The latter two indicate that less tolerance to loss will positively effect \bar{g} . More aggressive risk-taking will lead to a volatility tax ‘cost’, while having less tolerance to loss will yield a volatility tax ‘benefit’. The exact nature and impact of volatility tax on the longer-term return will be entirely dependent on the gamble and the parameters that can be controlled. With this formulation the primary concern is simply the final outcome \bar{g} , the success of strategy is then measured by its simulated performance across many trials. This is largely an abstract intractable problem but serves to highlight the decision-making process.

For the gamble at hand, we can optimise Eq. (11) directly assuming no stop-loss or retention ratio, that is, taking the same fixed blind betting strategy across all time as with the first investor. This yields

$$\frac{dg_{t+1}}{dt} = \frac{d}{dt} \left(p_u \ln |V_t(1 + r_u l_t)| + p_d \ln |V_t(1 + r_d l_t)| \right) \quad (15)$$

$$l_t^* = - \frac{(p_u r_u + p_d r_d)}{r_u r_d} = \frac{p_u r_u - p_d |r_d|}{r_u |r_d|} \quad (16)$$

$$p_u \geq \frac{|r_d| (l_t^* r_u + 1)}{r_u + |r_d|} \quad (17)$$

where the last line indicates the minimum up-probability required for any given fixed $l_t^* \forall t$ to be profitable. Hence for our given parameters they find a fixed $l_t^* = 25\% \forall t$ to match $p_u = \frac{1}{2}$ in order to maximise growth \bar{g} in the infinite time limit while seemingly having a heavily reduced $\mathbb{E}[r] = 1.25\% \forall t$.

Let us now simulate this gamble with $N = 150,000$ random investors for $T = 3,000$ steps from each of the four categories. For each category, all N investors will utilise exactly the same strategy and so the difference in final portfolio values will be purely due to the random sequence of payoffs they experience, in other words, the cards they were dealt. The goal is then to accurately assess how robust each of the strategies are for the average investor in N . To accomplish this, we split the sample into two sub-samples, the first contains the bottom 99.99% of adjusted average values V_T^A while the other consists of the top 0.01% of performers with maximum average V_T^M for each time step. This conservatism avoids upwards bias present when working with random variables since it removes the

impact of unlikely cases, a total extreme being $V_T = V_0(1 + r_u)^T$. The choice of large N therefore ensures we capture the behaviour any random rational investor, using T steps will be shown to sufficient in determining the long-term outcome, and removing the top 15 investors will give us a more complete picture for performance for 9,999 out of total 10,000 investors.

For the first type of investor the results for fixed leverage ranging from as little as 10% all the way to betting entirely 100% on favourable bet are shown in Fig. 2. Due to the differences in magnitude, we present results in terms of the log averages $\log_{10} V_t$. Recall that changes in tick increments on log axes represent changes by magnitudes of 10x for each increment.

These are a puzzling results as Fig. 2(a) shows that for $l \gtrsim 75\%$ the $V_T^A \rightarrow 0$ with certainty at rates increasing with leverage. The optimal constant leverage is found to be in the range $l^* \in [35\%, 45\%]$ so we take $l^* \approx 40\%$ which is incredibly difficult to predict a priori. The expected optimal leverage of unity leads to complete loss of capital when they ‘expected’ $V_T^A = V_0 \cdot 1.05^T \sim \10^{63} or $\log_{10} V_T^A \sim 63$. In Fig. 2(b) we see the box plot distribution for log values at maturity across leverages that further reveals even for $l \gtrsim 50\%$ at least half the investors are below the starting value and that only for $l \lesssim 15\%$ are they all guaranteed to make a profit. We therefore refer to $l^* \approx 40\%$ as the optimal leverage and $l_s^* \approx 15\%$ as the safe leverage. We also plot the trajectory from the initial value (2, 2) of the log average maximum value against the log average adjusted value in Fig. 2(c) that reveals how astronomically larger these top values are in magnitude. It is these abnormally large values this that drive the largest contributions to the mean forcing it to not approach zero instantly, and so removing them is crucial to accurately assess performance. Finally, in Fig. 2(d) we show all three MADs and standard deviations (STD) added to the mean, highlighting why STD is inferior to MAD as large outliers hugely inflate volatility.

Overall, we arrive at a completely non-trivial result where the optimal decision for the first investor generates $\bar{g} \approx 0.56\%$ return per step while always holding on to 60% of their capital at all times, while ‘expecting’ $\mathbb{E}[r] \approx r \cdot 0.40 = 2.00\%$. The volatility tax is $\nu = \bar{g} - \mathbb{E}[r] \approx -1.44\%$. This is an absolutely astounding result, the asymmetric effect of losses offset the majority of the ‘expected’ upside, they receive only 25% of what they predicted. Similarly for the safe leverage $\bar{g}_s \approx 0.32\%$ with $\nu = -0.43\%$ while ‘expecting’ $r \cdot l_s^* = 0.75\%$, where steep losses erase 66% of the ‘expected’ gain with everyone still benefiting from the gamble.

This is not only true for $|r_d| < |r_u|$ but is even more significant for the opposite case. For fixed leverage across time, this is still the optimal play as it maximises time-averaged return in the long-run through the avoidance of steep losses. Effectively, we can say in reality $\bar{g} \ll \mathbb{E}[r]$ and conclude that maximisation of the ‘expected’ return is not the quantity of interest, rather we only care about the performance of random investor over time. Notice to arrive at this conclusion, at no point were probabilities calculated, only the paths were simulated.

Moving on to the second category of investors that automatically calculate optimal leverage for a range of fixed stop-loss values and are the kinds of people that calculate standard deviations. The results for any $\lambda \in [0, 1)$ for $V_{\min} = \lambda V_0$ are shown in Fig. 3. We see in Fig. 3(a) that $V_t \rightarrow V_{\min} \forall \lambda$ with complete certainty (zero MAD) within roughly 25 time steps and the approach rate of V_t^A is significantly quicker than V_t^M . Notice that this seemingly optimal strategy is able to intermediately generate enormous profits for the top 0.01%, up to $V_t^M \sim \$10,000,000$, while inevitably these also eventually crash to their stop-loss. While they may ‘cash out’ at any time, it is ‘irrational’ to do so since this is a favourable bet, and the game should be played till maturity. This also reveals how heavily such a small minority can raise total average. The changes in optimal leverage in Fig. 3(b) show this again show this behaviour where the maintenance of the maximum leverage by the 0.01% occurs only for roughly the 15 intermediate

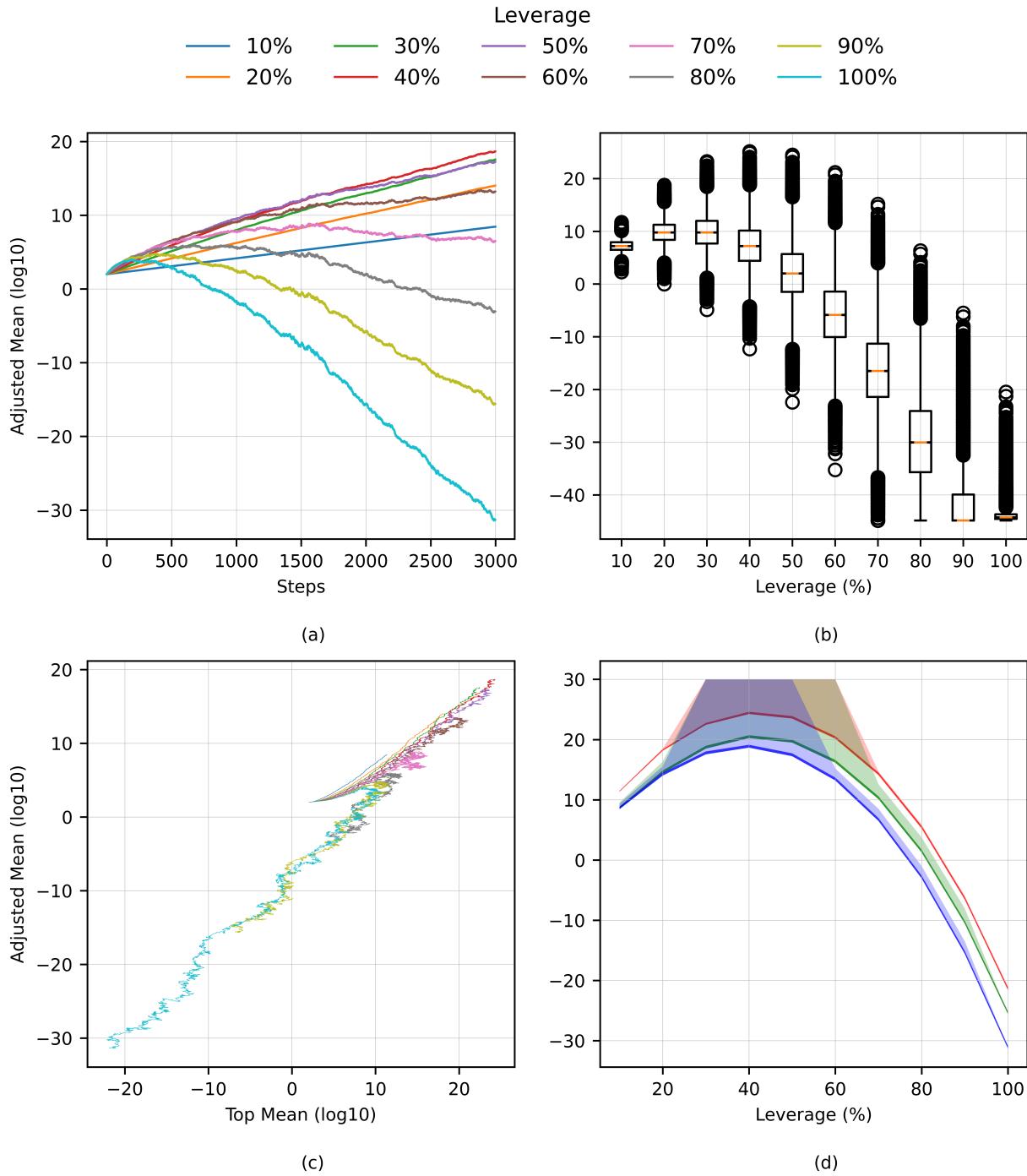


Figure 2: Summary of investor one results: (a) The trajectories of the adjusted (bottom 99.99%) log means for various leverages showing non-ergodicity at display, (b) Box plots of the distribution of adjusted values at maturity of $T = 3,000$ step for various leverages. Note 10^{-40} is not a lower bound, rather it is a numerical accuracy limit, (c) The trajectory of the adjusted mean along with the top (0.01%) log mean for various leverages all initiated at $(2, 2)$. Notice in all cases how astronomically larger the top values are compared to the adjusted values, and (d) Plots the MAD (dark shading) and STD (light shading) added to the mean for various leverages across all three subgroup adjusted (blue), complete (green), and top (red). Observe how STD grossly inflates the true volatility of the gamble in the profitable regions due to the small number of high performers in every group.

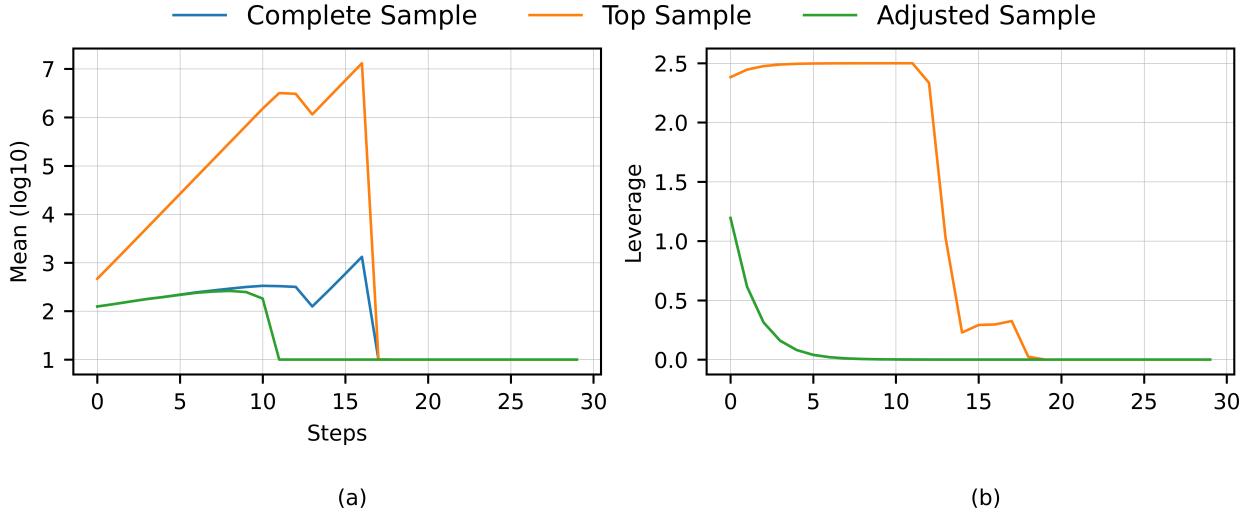


Figure 3: Summary of investor two results: (a) Reveals the mean values for each of the three subgroups highlighting how the top 0.01% are able to misrepresent the complete samples' performance by temporarily greatly skewing it upwards, and (b) The mean leverages for the same groups noting the inevitable crash to zero.

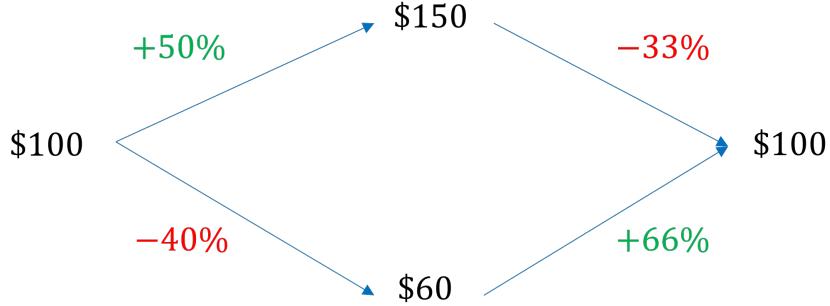


Figure 4: Visual depiction of the simple gamble for two paths to leading to no change in initial value, highlighting the pitfalls of using expectation values for multiplicative processes. Based initially on the example in [4, 5].

steps coinciding with the huge valuations. Therefore, these investors' strategy also fails to consistently increase valuations over time.

Why does this occur? Recall the example in Section 1.1 for why paths to the same valuation matter. In this case we have again equally probable up and down paths, both expressed as $100\% \cdot (1 + 50\%) \cdot (1 - 40\%) = 90\%$ and $100\% \cdot (1 - 40\%) \cdot (1 + 50\%) = 90\%$. Notice the final change in value is a 10% decrease from the initial value. For no change in value we must have $100\% \cdot (1 + 50\%) \cdot (1 - 33\%) = 100\%$ and $100\% \cdot (1 - 40\%) \cdot (1 + 66\%) = 100\%$ if the up and down state occur first respectively as shown in Fig. 4. In a compounding (multiplicative) world, this asymmetry that lies at the heart of everything. We need a 66% gain to recover from a loss, and a -33% loss to revert a gain. Due then to the power of compounding over time, the losses have an asymmetrically larger effect on portfolio worth.

The source of this error lies in the fundamental assumptions of Eq. (5) where they express optimal leverage to be linear in 'expected' return. As such, they lose all sense of the multiplicative nature of this process. This can be considered analog to case of incorrectly using a linear utility function for investor wealth when a logarithmic function

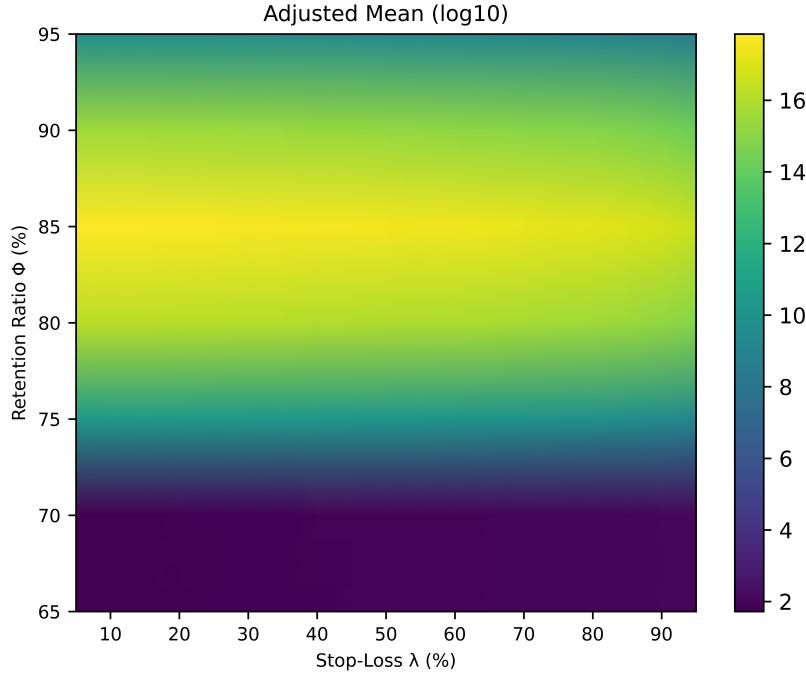


Figure 5: Investor three adjusted mean as a function of retention ratio and stop-loss at maturity of 3,000 steps. Observing that a maintaining a fixed retention ratio of 85% of all winnings at each time step and betting a portion of the remainder appears to generate the highest valuations across all stop-loss levels.

in Eq. (11) would be far more appropriate as it correctly penalises large steep losses [1–6].

The third category of investors proceeds unaware of the performance of the previous two and are the kind of people that use words like “alpha”. Their modified approach is based on calculating optimal leverage as function of the tuple (λ_t, ϕ_t) at each step. With this prudent structure they seek to steadily accumulate wealth by betting only a fixed portion of their winnings at each step. In Fig. 5 we show the results of the grid search using $\lambda \in [0, 1]$ and $\phi \in [0, 1]$ represented with density plot at maturity T . The optimal value are found to be in the vicinity $(\lambda_T^*, \phi_T^*) \approx (5\%, 85\%)$. The key variable being then retention level ϕ as there appears to be a steady rise in performance as it increases reaching a maximum at holding on to 85% of your achieved profits at each time step when calculating the leverage for the next step. The stop-loss at 5% of V_0 is also clearly understood as it gives the investors the greatest chances of escaping the $V_t > V_0$ threshold while not being stopped out. Figures for final leverage, MAD and STD look identical in density. The optimal leverage appears to converge around $l_T^* \approx 0.18$ and the adjusted MAD for this grid search is once again comparable if not and order of magnitude larger than adjusted averages as with first investor, highlighting that there is a very real possibility of still going bankrupt.

So how did the first three investors perform? The first two would be catastrophically crushed by an entirely risk averse individual that keeps the \$100 in their pocket. The third will achieve outstanding performance of $V_T^A \sim \$10^{17}$ with $\bar{g} \approx 0.53\%$ and $\nu_T = \bar{g} - rl_T^* \approx -0.37\%$. Despite this, they still 10x underperform the very simple strategy with fixed optimal leverage $l^* = 0.40$ with no stop-loss or retention that achieves $V_T^A \sim \$10^{18}$ with $\bar{g} \approx 0.56\%$ and $\nu \approx -1.44\%$. They do however beat the safe leverage $l_s^* = 15\%$ holding $V_T^A \sim \$10^{11}$ with $\bar{g}_s \approx 0.32\%$ and $\nu = -0.43\%$. We emphasise how this seemingly small difference in \bar{g} compounds over T to result in such a large difference. For

$T = 5,000$, investor three underperforms by 100x compared to the optimal leverage which is then $l^* \approx 35\%$. What we see is that as $T \rightarrow \infty$, leverage l^* decreases which is a simple resolution of the Gambler Paradox.

What we see here is that the highest valuation is achieved by maximising $\mathbb{E}[r] + \nu$ across all time. This under-performance is puzzling, but the explanation has to do with short-cut learning, that is, the sophisticated investor gets the correct result, but achieve it using the wrong logic. Their (rightful) success is due entirely to their ironclad (fully-automated) risk management, where even on a favourable bet, they have minimal tolerance for steep losses and so limit their leverage so that they accumulate wealth slowly. This prevents them from losing value even though the strategy they have implemented is identical to the completely failing second investor.

We now turn our attention to the final fourth investor that solves the abstract Eqs. (9, 11). The simple betting strategy may form a non-unique solution to these equation. It is this type of decision-making we are interested in obtaining as the expectation value-based methods have been shown to wholly inappropriate. For the optimal fixed $l_t^* = 25\%$, they obtain $V_T^A \sim \$10^{16}$ with $\bar{g} \approx 1.09\%$ and $\nu = \bar{g} - \mathbb{E}[r] \approx -0.16\%$. This is strategy has the lowest volatility tax and so would be the optimal strategy for any random investor in the infinite time limit. This result effectively generalises and solves the problem outlined in [4, 5].

Why does the simple blindly betting 25% of wealth at each step with no stop loss or profit retention strategy outperform all others? The reason is that the first three investors do not have any understanding on how to take risks in a compounding world. By conflating probabilities with payoffs, the root of their problem lies again in using expectation values which are only valid in additive environments. Probabilities have absolutely no bearing on the real-world performance when the change in value is multiplicative. In these cases, when a random investor calculates an expectation value, they are implicitly announcing via megaphone that they have the ability to pool results across the entire sample N and receive the average which we have shown can be grossly inflated by the performance of the ‘lucky’ top 0.01%. This is presumably accomplished by the singular random investor travelling to N parallel universes and forging a contractual agreement with all themselves to split the winnings evenly amongst each selves. However, for investors confined to a singular reality, the only feature of any concern is the payoff structure with strategies evaluated using time-average growth rate $\bar{g} = \mathbb{E}[r] + \nu$ of the bottom 99.99% as it accurately reveals the outcome for any random singular investor.

Note we do not form explicit risk-reward ratios such as Sharpe, as these are entirely worthless and misleading metrics for past performance. This is because they are independent of leverage, do not account for higher order statistical moments skew and kurtosis, and assume a equivalent importance for both variables where the whole purpose of this example to show the far greater importance of volatility. Furthermore, STD is not an accurate measure of volatility for the majority of underlying probability distributions. Accounting for these factors [9] reveals that high Sharpe ratios, if anything, are a much better predictor for who will go bust in the future rather than achieve superior performance. Ideally, the only metric that should be used to compare performance across any and all strategies that is manipulation proof, correctly penalises for leverage and volatility, and is effortless to communicate is the time-averaged growth rate \bar{g} in Eq. (10) [13–20]. Reporting of past performance via any other measure should merit scepticism.

While we could further explore this simple toy model, such as the interesting short-term profitability for roughly the first 200-300 steps, the actual concern is real-world problems and whether this approach would generalise to other compounding domains. By fully acknowledging the multiplicative nature of the problem and shedding all pre-existing notions of value maximisation, whether a robust methodology exists is an open question.

1.3 Conflating Probabilities with Payoffs

Formally the error in Eq. (5) can be expressed by the inseparability of probability [9]. Given the probability mass function $\pi(\cdot) : \mathcal{A} \rightarrow [0, 1]$ and measurable payoff function $y : \mathbb{R} \rightarrow \mathbb{R}$, for a subset $\mathcal{A}' \subset \mathcal{A}$ we have both the continuous and discrete cases

$$\int_{\mathcal{A}'} dx \pi(x)y(x) \neq \int_{\mathcal{A}'} dx \pi(x)y\left(\int_{\mathcal{A}'} dx\right) \equiv \mathbb{E}_{x \sim \pi(x)}[y(x)] \quad (18)$$

$$\sum_{x \in \mathcal{A}'} \pi(x)y(x) \neq \sum_{x \in \mathcal{A}'} \pi(x)y\left(\frac{1}{N} \sum_{x \in \mathcal{A}'} x\right) \equiv \mathbb{E}_{x \sim \pi(x)}[y(x)] \quad (19)$$

if the payoff function $y(x)$ is non-linear by Jensen's inequality [21] where the right-hand side is the payoff of an average event weighted by the probability of the event. By non-linear we mean that the response to a change of any magnitude $x_2 - x_1 \neq y(x_2) - y(x_1)$, with the discrepancy becoming much worse as the difference expands. The payoff function, or exposure, is in general a non-linear function, while probabilities on their own are nothing more than reductionist beliefs. Since any investor only cares about their end performance which is non-linear in payoff, risk management and leverage allocation should never operate in probability space [9, 13–16].

The only situation where the expectation value can be effectively used is strictly when the payoff function is fixed where $y(x) = \theta(x - K)$ is the standard general Heaviside step function. This situation clearly represents the realm of additive dynamics which forms the entire basis of contemporary decision theory.

The question is then whether this simplification can potentially be used in multiplicative environments. Suppose we are interested in calculating the expected non-linear payoff exceeding a point $K \in \mathbb{R}$, we define the two integrals

$$I_1 \equiv \int_K^\infty dx \pi(x)y(x) \quad (20)$$

$$I_2 \equiv y(K) \int_K^\infty dx \pi(x) = p(x \geq K)y(K) \quad (21)$$

where I_1 is the true outcome and I_2 is the pseudo-expected value. The substitution of integral I_2 for I_1 can be considered valid only for thin-tailed distributions defined by criterion $\eta \equiv \lim_{K \rightarrow \infty} \frac{I_1}{I_2} = 1$ [9]. This essentially implies that $p(x \geq K) = p(x = K) + p(x > K) \approx p(x = K)$ as the speed of the probability of extreme values crashing to zero is not offset by their absolute magnitude. The normalised and centered Gaussian distribution does satisfy this limit condition if at most the payoff is linear $y(x) = x$.

Distributions with $\eta > 1$ are defined as having fat tails where large deviations from the mean are far less likely but exist on astronomically larger scales. These distributions operate with an inverted logic, the typical behaviour in the vicinity of the mean is treated as the noise, while the signal is the exceedingly rare extreme event that entirely defines the long-term behaviour of the environment. An example of this is selling uncapped insurance, most of the time the insurer will generate steady premium with minimal volatility, however when the flood comes, all past performance will be entirely erased by one event. The key fact of fat tailed distributions is that after obtaining more empirical data, you can only ever revise them to be more ‘fatter’, never the opposite.

Unfortunately, not only is virtually every single conceivable payoff function in the real-world is non-linear, the overwhelming majority of underlying distributions in reality, at the bare minimum, have kurtosis in excess of the Gaussian [9]. While for very small changes in the $\delta x = x_2 - x_1$ we may be able at best to take a perturbative approach using linear Taylor approximation of the payoff so that $\delta x \approx y^{(1)}(x_2) - y^{(1)}(x_1)$. This however is not useful in the

long run as the linear payoff hides the true nature of the process. Furthermore, since $\pi(x)$ can only become less thin as more data is collected, $\frac{d\eta}{dt} \geq 0$ and so $I_1 \geq I_2$ with difference only ever permitted to increase over time for all environments.

This can be seen through simple explanation in Section 1.1 where the final return is

$$y_{t_2}(x) = \begin{cases} 1 + \frac{x}{V_{t_1}}, & \text{if } V_{t_1} = \$100 - x \\ 1 - \frac{x}{V_{t_1}}, & \text{if } V_{t_1} = \$100 + x \end{cases} \quad (22)$$

If the absolute change in value $x = \$50 \rightarrow \0 the difference between both responses clearly becomes negligible and so using the expectation value to conclude both paths are equally preferable can be considered accurate. Whereas if $x = \$50 \rightarrow \100 , we have $|\delta y_{t_2}(x)| \gg 0$ so we cannot simplify this payoff to be non-linear and therefore we have a distinct preference for the path avoiding the steep loss.

Regarding the simple gamble in Section 1.2 we have the non-linear (absolute) payoff at each step

$$y_{t+1}(r_u, r_d, \lambda_t, \phi_t, l_t) = \begin{cases} V_t - V_{\min,t} + V_t \cdot r_u l_t(\lambda_t, \phi_t), & p_u = \frac{1}{2} \\ V_t - V_{\min,t} + V_t \cdot r_d l_t(\lambda_t, \phi_t), & p_d = \frac{1}{2} \end{cases} \quad (23)$$

and the expectation value can similarly be utilised only if both returns $r \rightarrow 0$ which ensures $y_{t+1}(\cdot) - y_t(\cdot) \rightarrow 0 \forall t$ regardless of leverage if there is a stop-loss. The first three investors do not understand this requirement, despite $r_u + r_d > 0$ and $p_u \geq p_d$, the gamble is not additive, rather it scales with V_t . If the outcomes were fixed $V_t \cdot r_u l_t = R_u$ and $V_t \cdot r_d l_t = R_d \forall t$ then maximising expectations would be perfectly adequate.

The correct approach is in recognising this and taking the approach of the fourth investor that finds the parameters that simply maximise the time-average growth rate \bar{g} using Eq. (9) as it incorporates all effects. No attempt should be made to disentangle probabilities with non-linear payoffs. Therefore, we pose that the default stance should be to never use expectation values to model the outcome of any process, contrasting sharply with the contemporary prescription from decision theory. Nature and human civilisation are governed by multiplicative processes that by definition are non-additive since they scale dependent on the existing value. The actual response to changing of underlying variables can then easily be $I_1 \gg I_2$ if the variability in returns is significant.

1.4 Implications and Anecdotal Outcomes

We unscientifically proposed and aggregated the responses to the following gamble:

“If I were to give you \$100 and then offer you game where I flip a fair coin, if heads you make 50%, if tails you lose 40%, and lets say we play this game for 3,000 rounds. Would you play the game? If no, you get to keep the \$100 and continue with your day. If yes, as percentage, how much of your money would you bet at the start of each round and why?”

Furthermore, if the respondents were familiar with leverage, they were informed they could also take unlimited leverage at no additional cost. In total, depending on you measure it, we received 155 very informal responses.

In our experience the overwhelming majority of respondents would agree to play the game with fixed leverage $l = 1$ in line with the first investor as most were (thankfully) unfamiliar with the concept of leverage. This group consisted of people with very diverse backgrounds, with and without tertiary education. Most held undergraduate degrees

in finance, economics and a whole host of STEM fields from leading universities. Individuals with more advanced training, often in the form of postgraduate level education in finance or STEM recognised this as a ‘favourable’ gamble and universally proposed to take constant leverage $l \rightarrow \frac{5}{2}$ with no stop-loss $\lambda = 0$ across all time. Respondents with active financial markets experience also came to the exact same conclusion. Both of these groups then form a variant of the second investor. We received no responses resembling the strategy of the third investor.

Finally, a sizeable portion of people also choose not to play the game opting to keep the free \$100 as they in one form or another disliked the odds of the gamble. Very often they were unable to even express any reasoning for their risk aversion, only that they did not want to participate. A major difference between the composition of this group and the former groups was the overall lack of any formal training or education in economics, finance, or any of the STEM fields. The results for this game are of course known and discussed thoroughly in Section 1.2.

It should also never be forgotten that the first known causality of this process was the author in April 2020. Understanding this error was the principle motivation for the entire Master’s degree. This work should be interpreted as nothing other than a monument that catastrophic mistake.

A very frequent rebuttal is of the form:

“While this is a very interesting result, it is really nothing more than a trick question about compounding. It does not significantly affect much, however it would make a very good interview question.”

This game offers a laser-guided surgical exposé of the most deep-rooted underlying decision-making abilities of any individual. It is not about whether they just understand compounding, rather it’s about whether the person realises that the consequences of decisions they make now can have an everlasting (non-linear) influence on the final outcome. While everyone would say they understand this, actions speak louder than words, and failing to operate in line with this supposed understanding for such a simple gamble reveals huge inconsistency.

As virtually every single process in the real-world is governed by fat tails with multiplicative non-linear payoffs, not with fixed additive payoffs even remotely resembling bell curves, there is no good reason to believe that these people are operating correctly in practice. Furthermore, in the majority of cases, the confidence exhibited by respondents mimicking the first two investors was so immense that if it could be converted to electricity, a dozen of these people could cleanly power the planet.

Empirically we have also determined that the speed at which someone comes to the conclusion to dismiss these results is very positively correlated with their age and professional experience. We find immediate dismissal to be more common for participants over the ages of 29-33. A likely reason for their ‘knee-jerk’ reaction is that they are simultaneously finding out that not only is their belief system about formalised risk-taking incomplete, but that the implication of this is that any prior success they have achieved using this incomplete method is quite likely due to luck rather than skill. They utilised the same strategies as everyone else was taught, where others might have failed, they potentially achieved success, but likely only due to the cards they were dealt. C'est La Vie.

The source of this error is not the individual, rather the formal education they have received is lacking crucial content [1–6, 13–19]. They were taught probability theory and (thin tail) statistics but were not told that neither is particularly useful when making personal decisions that cannot be reversed. Probability lets them indulge in fantasises of escaping to more favourable parallel universes, thin tails raise them in a safety net where actions do not have long-term consequences. Most are also completely unaware of the concepts of non-ergodicity and time-averages.

This is clearly seen since the group with the least education had best performance while taking zero risk. Recall the final outcomes of the first two investors in this case would be \\$0, making the risk averse individuals' with \\$100 performance literally ' $\infty\%$ ' superior. While a portion of these people would never play the game, it is unlikely this includes all of them. We refer to the portion that would choose to play if the odds appeared (subjectively) better to them as 'risk averse'.

We hypothesise that since the risk averse person has no defective education, they are forced to totally surrender their decision-making to their subconscious brain, or 'gut' instinct. The human brain is not foolish enough conflate payoffs with probabilities, a skill acquired over millenniums of evolution. Its underlying goal is to survive and is very acutely aware of the fact that it has only one life, hence its ability to assess risk may involve intuitively processing Eqs. (15-17). While it is doubtful they are thinking in terms of optimal leverage, perhaps they abstractly recognise in Eq. (17) that for $l = 1$, they instinctively require better odds ($p_u \geq \frac{2}{3}$). These people are not "irrational", far from it, in fact the quality of their inexplicable human instincts is quite exceptional.

Another way to see this is by building on the analogy in [9]. When walking through the forest alone, the risk-averse investor has a tendency to mistake stones for bears but absolutely never the opposite. The person calculating expectations correctly realises the probability of encountering a bear is minuscule and so walks with a concern weighted by this infinitesimal chance. If they do encounter a bear, the only defence they have are a collection of A3 pages stamped with logos of old buildings. This of course is not an issue if they can travel to the overwhelming majority of parallel universes where the bear is not present. We repeat again for the final time, for multiplicative dynamics probabilities do not matter!

A major implication rigorously discussed in [9] is that this largely invalidates psychology research. This field, not known for its mathematical prowess, almost always assumes in Eqs. (20-21) that $I_1 = I_2$ and therefore mainly uses I_2 when researching human decision-making. One key claim they have is the cognitive bias where empirically, people assume higher probabilities for rare events than what model parameters would predict, implying they are excessively risk averse [7–9, 19]. We now know that characterising this as a mistake can be catastrophic outside of textbooks since payoffs are non-linear with fat tails. All conclusions they have arrived at using this simplification and countless others must be re-examined from scratch. Therefore, given that the whole profession just might turn out to be pseudoscience, its tendency to frequently label people as "irrational" may be very incorrect.

A point worth emphasising is that generations of Sveriges Riksbank Prize recipients in highly authoritative positions, openly announced, and proudly taught countless students that the overwhelming majority of the public they would encounter in the lives were "irrational". When in fact it appears their models were never seriously questioned as the unadjusted empirical data seems correct [1–9]. This may add considerable weight to the common belief that the average person has a far firmer grip on reality than academics.

Keep in mind that modern psychologists have the legal authority to involuntarily administer medications and provide advice that is taken seriously in court rooms which is truly terrifying as one day they may be seen as no different to the witch doctors of old. There is however another field of concern that is far more intertwined with every person's life. Most Western governments mandate all citizens to have an externally managed superannuation or pension scheme where portion of their salary is given to 'professional' money managers, with the alternative of managing your own wealth penalised with increasing fees. As many of these managers often underperform compared to extremely low-fee passive market products, the average random citizen is in real trouble since a sizeable portion of their retirement savings is entirely dependent on these managers.

Let us first reflect on the fact that these managers, their subordinates, their applicants, their superiors, and even their regulators' entire education is based on maximising expectation values. What could possibly go wrong?

Next lets recognise that for the very simple gamble in Section 1.2 based on just maximising the value of portfolio that has binary returns, if they did not implement a very carefully optimised rolling stop-loss and took the approach of investor two by applying enough leverage to not be wiped out by a single move, after 3,000 steps they would have ‘expected’ a valuation of $V_T \sim 10^{153}$. However, within about 30 steps they would all, without exception, be at their stop-loss. At this current time, we are unaware of any profession where a tolerance for error using the default textbook approach of 153 orders of magnitude would be considered acceptable. For reference, there are estimated to be 10^{80} atoms in the observable universe [22], and the size of the search tree required to recursively identify the winning strategy for the board games of Chess and Go are approximately 10^{124} and 10^{360} respectively.

Despite this failure, recall that for roughly the first 20 time steps the top 0.01% would achieve performance of $V_T^M \sim 10^7$ entirely due to the cards they were dealt, before eventually joining their peers at their stop loss. Note that at no point have we specified the unit of time, it could be seconds, minutes, days, months, years, or anything. If we use larger time increments these individuals will increasingly appear to be ‘legendary’ investors, while in practice the fundamental basis of their strategy is exactly identical to the other 99.99%.

Financial markets in the real world are of course far more complicated, there are huge amounts of assets, returns are variable, and simulation is difficult. However, the principle of focusing solely on maximising \bar{g} as opposed to $\mathbb{E}[r]$ remains valid throughout. The choice of leverage to take in any bet that does not have a fixed constant payoff is obviously non-trivial, but step one is recognising the multiplicative nature of the process.

Therefore, we present a heuristic that could be used to evaluate anyone that professionally manages money based on their response to the simple gamble. This should not be considered financial advice, but neither should anything they recommend be either. We classify their responses in terms of the previous investor categories. Responses mimicking the first two investors are obviously inadequate and so we focus on the third:

- (i) Investor 3: Question their “alpha generating strategies” until they have clearly demonstrated an understanding of multiplicative dynamics. However, never bet against them. The simulations reveal the terrifying levels of success these investors can achieve by using prudent and completely automated risk management principles. They can accomplish this while not having any better understanding of the gamble than the previous two investors, that is, minimal clue of what they are doing.

As they attempt to predict ex-ante $\mathbb{E}[\alpha]$, unless they are solely assessed on forecasting ability, they then have to actually allocate capital to capture it in a timely manner. Typically, the weights they select for each security in a portfolio are proportional to each securities ‘expectation’ which obviously does not correctly model the non-linearity of the payoff. Their performance is then measured by an ex-post $\hat{\alpha}$ defined relative to obtained return \hat{r} , this can be taken as the excess to a benchmark r_B , so that $\hat{\alpha} = \hat{r} - r_B$. However, the quantities $\mathbb{E}[\alpha]$ and $\hat{\alpha}$ are not even remotely equivalent.

This is the prescription detailed by all contemporary financial education programs offered globally. The same theories and models are taught, students are awarded high distinction averages and charters for accurately reciting these ideas, and then they proceed to act in the real-world in line with these beliefs validated now by A3 pages that simply act as testaments to their incomplete education.

One should therefore be very sceptical of investors dedicated to seeking “alpha” as most of them have likely been

taught to use expectation values in a compounding world. Therefore, they essentially operate on the assumption that the payoff they receive from taking actions in world around them is the (probability-weighted) average of what they predict, when in fact, they are confined to a singular future [13–19]. The “expected alpha” they chase, using their definition, can not be replicated for a portfolio constructed with variable leverage.

This would be akin to extremely competent brain surgeons not knowing the function of the brain. Does it matter if they always get the surgery correct? Not really. But if they make a mistake, devoid is their knowledge of the consequences.

A potential entirely self-admitted candidate for this type is [23] who outlines many principles with “Life Principle 5.6” being “make your decisions as expected value calculations” along with many other probability-based claims. Surely he would get this simple gamble correct? If not, his prior unmatched performance would therefore highlight his beyond outstanding world-leading risk management, while being somewhat inept at risk-taking. His selections of stop-losses and retention ratios are phenomenal, though his long/short predictions should be not be blindly replicated or taken at face value. An alternative explanation is that these “principles” are brilliant disinformation campaign to prevent the formation of competitors.

Regardless, suppose someone at this pinnacle level of achievement is not only publicly conflating probabilities with payoffs, but has proceeded to write a 600 page book receiving countless awards and limitless amounts of praise from highly authoritative sources, that ultimately serves to concretely display their misunderstanding. What are then the chances exponentially less prominent individuals, world-wide, that have a fiduciary responsibility to manage large portions of people’s life-long retirement savings also get it wrong?

The correct behaviour, instead, should be to make investment management decisions today that allows each of them embrace the one irreversible, unpredictable, volatile, and unknown singular fate that awaits them all. Amor fati [17, 19].

- (ii) The risk averse person that does not want to gamble with provided odds: Commend them for having better risk-taking instincts than the first three investors.
- (iii) All other responses: Since the gamble is simple, have an open mind and simulate the outcome.

Furthermore, despite all the perils of using expectation values, there exists an even more dangerous quantity, standard deviation σ . While an investor using expectations without prudent risk management will go bust, individuals using variance to quantify volatility will not only go bust, but will take counterparties down with them.

This occurs because variance is incredibly unstable under fat tails and all subsequent computations that utilise variance as an input are also then contaminated. Utilisation of variance minimising machine learning tools are not robust to shocks. This includes all forms of regression even with ridge and lasso regularisation. GARCH predictions are spurious as they cannot adapt to environments with large natural kurtosis without forcibly truncating the data. Value at risk measures combining both standard deviation and probability are ineffective. Similarly, one can misinterpret the natural existence of fat tails to be conclusive evidence of heteroscedasticity. Pearson correlations are at best uninformative. Non-parametric methods lead to even worse out-of-sample robustness. See [9] for further details.

1.5 A Permanent Solution

Two valid criticisms of the preceding text are whether the toy problem generalises to more real-world scenarios and that we have not provided a systematic method to solve Eq. (9) for any environment. Obviously, the simple optimisation in Eqs. (15-17) to find $l^* = 25\%$ is not possible in systems with huge numbers of continuous parameters.

This type of analysis was first done in the more useful case of finding optimal leverage for a financial security following geometric Brownian motion (GBM). GBM is often considered as the fundamental basis for most theoretical derivative pricing models in mathematical finance [24]. Using stochastic calculus and Itô's Lemma they found the explicit inverted parabolic relation $\bar{g} = \mathbb{E}[r] + \nu = lr - l^2\sigma^2/2$ revealing the importance of selecting the correct leverage [1, 10, 11]. Clearly then we have $l^* = r/\sigma^2$ implying that $\bar{g}_{\max} = r^2/2\sigma^2$ which is reminiscent to the exponent of the log-normal probability density function. Note that while $\mathbb{E}[r] \propto l$, we have $\nu \propto -l^2$, which is probably the shortest explanation in history for the dangers of leverage.

GBM is also a key model for representing self-reproducing entities, which may be considered as the definition of life as the dynamics it induces are of interest to those concerned with living systems from biology to economics [11]. They then determine very interesting long and medium-term behaviour of the simulated paths in terms of the speeds at which they will ultimately converge to \bar{g} .

In the real-world where all payoffs are non-linear and there are infinite possible multiplicative problems, the optimisation problem in Eq. (9) can be very generally expressed as

$$x_1^*, \dots, x_T^* \in \arg \max_{x_1, \dots, x_T} \sum_{t=1}^T (1 + g_t(x_t)) \quad (24)$$

where $x_t \in \mathbb{R}^n$ incorporates the range of all possible n actions that can be taken to modify the existing state of the system at each time step t in order to maximise the future rewards. Maximisation of the time-average growth rate is clearly a formidable problem and appears that it needs to be solved on a case-by-case basis. The key point is that at time $t = 0$ we must solve across all time steps from $t = 1 \rightarrow T$ to properly maximise $\bar{g} = \mathbb{E}[r] + \nu$.

Therefore while [1–6, 9–11, 13–19] have presented a concrete general case for the need for a reformulation of existing decision theory, the majority of the practical applications of their suggestions need to be explicitly constructed using Eq. (24) and then effectively solved. This is no easy task and requires bespoke domain knowledge for each problem.

To encourage this adoption, for over a decade they have presented their findings in hope to reform and re-educate existing practitioners on the validity of multiplicative dynamics, and more importantly, on the perils of using additive dynamics in a compounding world. Mainstream acceptance however has been at best mixed, despite an ever-growing community forming under the brand Ergodicity Economics. Regardless, we hope momentum continues to build incorporating multiplicative dynamics in environments where it is necessary to correctly model the scenario. The relatively recent endorsement [4] by the world-renowned late theoretical physicist Murray Gell-Mann will ideally serve as an accelerator for its adoption given that his existing contributions (to fields such as Quantum Chromodynamics) will be taught almost certainly, at the very least, till the end of time. Out of everyone, someone of this calibre has the potential to hold a candle to combined intellect [25, 26] of every economist that has ever lived regarding the foundations of economics.

The lack of widespread admission of the failings of the contemporary methods is probably well summarised by the age-old adage, “it is difficult to get a man to understand something, when his salary depends on his not understanding

it” (Sinclair 1935). This would be especially valid for senior practitioners that dominate and control their fields as this would involve acknowledging their entire careers have been built on an error, which is in line with our crude empirical findings in Section 1.4. This inability to decouple probabilities from payoffs not only effects them, and those around them, but can have devastating consequence for human civilisation since no matter how “data-driven” they are, expectation values do not capture the non-linearity of the outcome for any decision.

For example, returning again to pension funds management. OECD data [27, 28] reveals that despite turbulent markets, total global assets under management annually rose 9% to USD\$34.2 trillion at end-2020. For the countless managers of these funds, regardless of how well-meaning their intentions are and how seriously they take their fiduciary responsibility, the moment their decision-making process incorporates probability estimates, they essentially become delusional. No longer are they concerned with each client’s financial wealth, rather they aspire to maximise the wealth of the average of each individual client. The difference is subtle, but the impact astronomical. Unwittingly, the bulk of these retirement savings are effectively being mismanaged by people acting as schizophrenics. Much the same can be said for the countless other trillions being managed in more discretionary funds. While it is perverse that this has continued for so long, this situation breeds unimaginably profitable opportunities. The existence of these hordes of oblivious managers will permit the construction of strategies that will appear non-sensical to them, however, are capable generating inconceivably high returns for a given level of risk. One non-unique method involves a careful combination of largely a passive market portfolio and a small extremely convex insurance portfolio [13–19].

At the same time, we are also left with the intractable Eq. (24) and are unable to offer any systematic and tangible alternatives to the contemporary approach. Therefore, we require a method to generally enforce the principles of multiplicative dynamics that is compatible with maximising compounding growth rate. To be applicable to any environment, it also needs to be self-learning and fully autonomous to analyse the cause and effects of each x_t . This is necessary as directly solving stochastic partial differential equations as with GBM in [1, 10, 11] is not possible, let alone tractable in general. Successful construction of this approach also demands performance that eventually exceeds of all existing practitioners in each real-world domain. Only then will we be able to facilitate beginning the elimination and replacement of all practitioners that conflate probabilities with payoffs and occupy positions of serious power. While extreme, they have already had a decade to respond and so the only logical option left is to begin charting this course, since evidently, it seems “you can’t teach an old dog new tricks” (Heywood 1546).

The remainder of this work is dedicated to accomplishing this formidable task. To achieve this goal, a general problem-solving approach is required. Reading extremely carefully the words under Eq. (24) again, ‘take correct actions to modify the existing state to maximise future rewards’ — this is the language of reinforcement learning.

2 Introduction

Reinforcement learning is generally formulated under the assumptions of ergodic Markov decision processes (MDPs) that are stationary through time. One interpretation of ergodicity is that it implies the time average reward of an agent’s trajectory through state-action space is exactly equal to the expectation value of that reward. This simplification is useful to formulate the theory and prove numerous convergence criterions [29]. The difference between classical dynamics programming methods and reinforcement learning is that in the latter, large scale approximation methods are necessary as exact modelling of the MDPs becomes intractable [29–32]. Performance in environments across various algorithms is generally compared by additive dynamics using the averaging cumulative summations of the rewards the agent receives per time step in each evaluation episode. This approach combined with the use of deep neural networks acting as universal function approximators has led to highly promising advancements over the last decade with gains across the board.

The most well-known of these include classic Atari video games [33–42], the board games of Chess, Go and Shogi [43–48], StarCraft II [49], and the protein folding problem [50–52]. For continuous action spaces pertaining largely to continuous locomotive control tasks [53–61], actor-critic methods combining advancements in Q-learning [62–64] with stochastic or deterministic [65–71] policies have steadily matured. An alternative approach using the principle of maximum causal entropy [72] has led to the soft actor-critic algorithms [73–79] incredibly robust results requiring minimal hyperparameter tuning. In parallel, for these continuous action spaces, on-policy model-free methods [80–82] have also achieved decent performance, and augmented random search [83] has once again proven to be remarkably effective in simpler environments.

During the Q-learning phase, the critic minimises the difference between the Q-value and a target Q-value obtained using the Bellman equation. The vast majority of literature constructs this off-policy critic loss through aggregating the (mean-square) Bellman error. This use of the MSE loss function is considered by [84] to be a “reasonable choice” while acknowledging there is a ‘lack of a good understanding for this measure’, but regardless, believes that “most conclusions would hold for different measures”. There does not appear to exist any such analysis experimentally validating these claims. The use of the MSE functional form is likely due to several reasons, firstly the overwhelming majority of machine learning is built upon this variant of L_2 -norm which traces back to its ability in finding optimal coefficient for linear regression. Secondly, the compatible function theorem and several historical proofs have used MSE from which further developments solidified its prevalence. There likely also exists many other reasons unaware to us at this time. Based on recent advances in non-negative matrix factorisation (NMF) proving the benefit of other loss functions [85], there exists a possibility that similar gains may be achieved.

An additional feature that is globally utilised throughout all statistics and machine learning is when constructing of the empirical mean, it is assumed to converge by the strong law of large numbers to the true unknown population mean in the infinite sample limit. All samples in the real-world are finite and so this convergence never formally occurs, regardless the (equal-weighted) arithmetic mean is still utilised. This approach is appropriate if underlying distribution is well-behaved, meaning that it has thin-tails, and a true mean exists. If it is fat-tailed, not only may the variance be undefined, but so might the mean. To combat this, we need to first test whether a distribution is fat-tailed and if so, a ‘shadow’ mean should be estimated as a more faithful representation of the true mean [9]. This is a very recent development and its applications to statistics and machine learning remain a very open question that merits our investigation.

Off-policy learning is characterised through the use of an experience replay buffer \mathcal{D} [86]. Mini-batch sampling is often either uniform $U(\mathcal{D})$ or with prioritisation proportional to absolute TD-error [87]. The vast majority of algorithms generally use a fixed buffer containing the 10^6 most recent transitions, mainly due to its historical use in [34]. For discrete action spaces size of the buffer is found to be highly impact on overall learning especially when combined with other features [88, 89]. The size of the buffer is interpreted by [89] as the degree of ‘off-policy-ness’ as more data allows the agent to search further back in history and therefore reduces the chance of overfitting.

Using a modified Rainbow agent called Dopamine [41] on classic Atari games, [89] finds performance increases as both the size of buffer increases and the age of the oldest policy reduces which is intuitive as more recent transitions are higher performing. Therefore, there exists a trade-off between buffer size and how relatively ‘on-policy’ the transitions are. They also show that prioritised experience replay gives no additional benefits with increasing buffer size which is an incredibly unintuitive result as one would expect preferential sampling to be very useful. Interestingly, [89] determines multi-step returns [62] to be absolutely crucial using ablative trials while acknowledging there exists no theoretical or well-grounded explanation its importance. Any bootstrapping is superior to the vanilla one-step case when simultaneously increasing buffer size, with the optimal level found via grid search. One explanation for this puzzling result is that gains achieved from increasing the buffer size are positively correlated with the variance of target returns and therefore increasing steps achieves this randomness. There does not appear to exist an equivalent analysis on whether their results hold in continuous action domains.

Under additive dynamics, an algorithms overall performance is evaluated using the summation of the rewards received at each time step. This approach is not suitable for a large class of environments where losses have an asymmetrically larger effect on performance than identical magnitude gains. This is particularly crucial in situations where the agent is operating in mission critical scenarios where time order matters and the agent has only ‘one chance’ as opposed to assembly line production tasks. To model such environments, multiplicative dynamics [1, 4] must be used where evaluation episode performance is expressed as cumulative compounding returns of the rewards received at each time step. Much of the applications for such environments are in non-ergodic domains where the time average is not equal to expectation value or more accurately, the ensemble value.

Pioneering work on such tasks has been done over the last decade [1, 4, 9–11] with direct applications to finance and economics [2, 3, 5, 6, 13–19, 90–95] with experimental psychology validation into optimal human decision-making [7–9, 19]. This type of modelling is crucial in environments such as medicine, supply chains, guidance systems, economic policy, financial portfolio management, and systems control where we may want to encode the asymmetric effect of large negative rewards. The crux of this approach is that the correct way to represent non-ergodic domains is with multiplicative dynamics using compounding products of returns over time as opposed summations.

These ideas were initially correctly posited in 1965 by the well-known Kelly Criterion [96] obtained using information theory and used for making optimal bets based on prior performance where the probability of bankruptcy approaches unity as games of maximum (expected) rate of return are repeated. This certainty of bankruptcy is the well-known Gambler’s Paradox [97], which was also a prime motivator for development of the causal conditioned entropy methods [72] used to derive the soft actor-critic algorithm where time ordering is crucial [73–79]. The current paradigm however is entirely built on a parallel development in 1952 that bases decisions on incorrectly maximising expectations using Markowitz (or Modern) Portfolio Theory [98–100]. Kelly emphasised maximising the time average reward by avoiding steep losses, while Markowitz based risk preferences on completely subjective utility functions that are dependent on personal circumstances, see [1–6] for the complete history.

A direct application of this is seen by the fact that the average economic growth of a country is not equal to the economic growth of a random citizen over time. Economics is entirely concerned with the former and implicitly states that it is exactly equal at all times to the latter [5, 6, 95], and then proceeds to casually build an entire discipline around this assumption. The responsibility for this mistake traces to using a 300-year-old result of Bernoulli that contained an error, and then subsequently a failure to pick up on its 200-year-old (minimally advertised) correction by Laplace [4]. The failure of contemporary academics managing global economies and world trade to recognise this glaring weakness, amongst several other crippling issues [101, 102], is quite alarming.

As an aside, the Markowitz approach also throughout utilises probabilities bounded within the closed interval $[0, 1]$ to construct expectations. Results from the path formulation of quantum mechanics, reveal that probabilities meaningfully exist outside both bounds of this domain in nature [103, 104]. The interpretation of these values in terms of Bayesian inference where decisions on whether information gathering and utilising systems should take a bet is determined by whether it can be first settled, and then to isolate all possible realities [103]. Therefore, it is debatable whether expectation values should be used at all if one does not normalise across the true complete probability space. Accounting for this may offer a path to reconciliation where the additive case approaches multiplicative dynamics.

Existing work on extending the applicability of Q-learning to non-ergodic state process beyond MDPs can be written in the language of Feature Reinforcement Learning and state aggregation [105–109]. Much of this is inspired by earlier work on incorporating partially observed MDPs (POMDPs) and other non-Markovian decision processes (non-MDPs) into reinforcement learning [110–113]. Other more recent approaches in this area is discussed in [114–117]. Using extreme state aggregation [107], any non-MDP can be modelled as a finite-state MDP if there exists a feature map aggregating different histories to states. Construction of this feature mapping is non-trivial in general but correctly reduces to all the well-known results if the underlying process is an MDP. One way to incorporate non-ergodicity while retaining much of the existing machinery [109] involves introducing a new class of Q-Value Uniform Decision Processes (QDPs) specifying several constraints. Under these constraints they show that Q-learning under the additive dynamics case converges to the optimal action-value using slight modifications of the usual methods and so this approach can be used in a subset of non-stationary domains called QDPs. One unknown is that the ability of the QDP approach to perform while using function approximators such as deep neural networks which is essential to many practical domains.

Then by modifying the existing formulation of model-free reinforcement learning to be compatible with maximising compounding growth rate, we will be able to construct full autonomous, self-learning, and risk-reward maximising algorithms that can operate in multiplicative domains. This is necessary as the path to artificial general intelligence (AGI) will be of little value if these AIs also get simple gambles such as those in Section 1.2 incorrect since reinforcement learning as a field is also entirely based on maximising additive sums of all predicted future rewards. Therefore, through making fairly straight-forward modification to well-known algorithms, it might be possible to reformulate the way in which these agents learn. Perhaps modifying their risk-taking to be more consistent with reality [7–9] may lead to them engaging in interesting activities.

As the field of reinforcement learning matures and eventually begins to enter the real-world environments where the cardinality of state and action spaces expands immensely, concerns regarding computational efficiency of agent training and operating will begin to emerge. Generally, the efficiency of training is not of pivotal concern as it can be parallelised over dedicated supercomputing clusters. Agent operation and inference one the other hand is what will occur in practice. For applications where the agent operates using a battery, by minimising power consumption,

operating time will be extended which will inevitably reduce costs. There are two broad methods to reduce power: 1. Increasing hardware computational efficiency, and 2. More efficient agent learning algorithms. The first is outside the scope of this work. The second is our focus and is expected to naturally occur as time progresses. For the special case of agents operating in environments separated by distinct phases or stages, we propose ‘linking’ several agents sharing the same state space but different action spaces. Over millions of agent decisions per agent, the reduction in computational resources is likely to steadily accumulate resulting in overall lower energy consumption permitting lengthier operation that would be especially important in mission-critical situations.

This project is structured with Section 3 presenting a comprehensive review of background material. Section 4 outlines differences between additive and multiplicative dynamics in the context of non-ergodicity. Sections 5 and 6 modify critic Q-learning and actor policies to incorporate multiplicative dynamics respectively. Section 7 incorporates these changes into the soft actor-critic algorithm. In Section 8 we provide an introduction to a new framework for efficient reinforcement learning for segregated action spaces. Section 9 provides a brief recap of motivations and related work justifying our originality. Experiments are conducted in Section 10 covering several of the research areas. An overall discussion summarising all findings is presented in Section ??.

In Appendix A the two utilised agent algorithms are presented. Appendix B provides a succinct outline of potential applications of multiplicative dynamics utilising reinforcement learning. Finally, in Appendix C a brief summary of the limits of our models in the context of delivering realistic capabilities.

3 Background

3.1 Preliminaries

Standard reinforcement learning is formulated by considering an infinite-horizon Markov decision process $(\mathcal{S}, \mathcal{A}, P, r, \gamma)$ forming a sequence of states $(s_t)_{t \geq 0}$, actions $(a_t)_{t \geq 0}$, and rewards $(r_t)_{t \geq 0}$ experienced by an agent at each time step $t \in \mathbb{Z}^+$. The agents' behaviour is characterised by observing a state $s \in \mathcal{S}$, selecting the next action $a \in \mathcal{A}$ to take from state s based on its current policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$, arriving at the new environment state $s' \in \mathcal{S}$, and also receiving a bounded reward from this transition $r : \mathcal{S} \times \mathcal{A} \rightarrow [r_{\min}, r_{\max}]$ where $r \in \mathcal{R} \subset \mathbb{R}$. The Markov property implies there exists an initial state distribution $p_1(s_1)$ and a stationary transition conditional distribution satisfying $p(s_{t+1}|s_1 a_1 \dots s_t a_t) = p(s_{t+1}|s_t a_t)$. The discount factor $\gamma \in [0, 1]$ determines the priority given to shorter-term rewards. The environment also provides at each time step a Boolean done flag signifying whether the game or episode has concluded. Model-free learning is characterised situations where the agent does not explicitly utilise the transition probability distribution P , rather it takes a trial-and-error approach such as Q-learning.

The policy used to select actions in the MDP could be either stochastic $\pi_\phi(a_t|s_t) : \mathcal{S} \rightarrow P(\mathcal{A})$ or deterministic $\mu_\phi(s_t) : \mathcal{S} \rightarrow \mathcal{A}$ where the policy is generally expressed with a vector of n parameters $\phi \in \mathbb{R}^n$. The reward for all future time steps from a time t is written as an additive discounted cumulative future reward $R_t \equiv \sum_{k=t}^{\infty} \gamma^k r(s_k, a_k)$. For stochastic policies, the state value and action-value functions are defined to be $V^{\pi_\phi}(s) \equiv \mathbb{E}[R_t|s; \pi_\phi]$ and $Q_{\pi_\phi}(s, a) \equiv \mathbb{E}[R_t|s, a; \pi_\phi]$ respectively. The objective to then maximise becomes $J(\phi) = \mathbb{E}_{s \sim \rho^{\pi_\phi}, a \sim \pi_\phi}[R_t] = Q_{\pi_\phi}(s, a)$. In discrete action spaces, a ϵ -greedy approach is generally taken with regards to action selection. For continuous action spaces we are able to explicitly optimise the policy gradient $\nabla_\phi J(\phi)$.

The density at state s' , after transitioning for t time steps from a state s is represented by $p(s \rightarrow s', t, \pi_\phi)$. The (improper) discounted state distribution representing the marginals for the trajectory distribution induced by the policy is $\rho^{\pi_\phi}(s, a) \equiv \int_{\mathcal{S}} \sum_{t=1}^{\infty} \gamma^{t-1} p_1(s) p(s \rightarrow s', t, \pi_\phi) ds$. Therefore, by definition, agent learning is a non-stationary process as the policy improves $\pi'_\phi \leftarrow \pi_\phi + \alpha(\pi'_\phi - \pi_\phi)$ where $\alpha > 0$ is the actor learning rate. All proofs of optimal convergence meanwhile are derived under the MDP assumption that transition probabilities do not vary across time. It should be noted that while this inconsistency is present throughout the field, we are still able to train functioning agents, but the true impact of such an approximation is not well understood.

3.2 Policy Gradient Theorems

For continuous action spaces, there exists two well-known policy gradient methods: stochastic and deterministic. The first is very well-known policy $\pi(a|s)$ and the gradient of the objective is

$$\nabla_\phi J(\phi) = \mathbb{E}_{s \sim \rho^{\pi_\phi}, a \sim \pi_\phi} [\nabla_\phi \ln \pi_\phi(a|s) Q_{\pi_\phi}(s, a)] \quad (25)$$

where the monotonic log likelihood is called the score function or eligibility vector [65]. Deterministic policies are characterised by $\pi_\phi(a|s) \rightarrow \mu_\phi(s)$ and modify the density $\rho^{\pi_\phi} \rightarrow \rho^{\mu_\phi}$. For continuous action spaces, instead of greedy action selection $\mu_\phi(s') = \arg \max_a Q_{\mu_\phi}(s', a)$ at each step which leads to global maximisation, a more efficient approach is to move policy $\mu_\phi(s)$ in the direction of $\nabla Q_{\mu_\phi}(s, \mu_\phi(s))$ for each transition. This lets us approximate the

gradient of the actor learning objective as

$$\nabla_\phi J(\phi) \approx \mathbb{E}_{s \sim \rho^{\mu_\phi}} [\nabla_\phi Q_{\mu_\phi}(s, \mu_\phi(s))] \quad (26)$$

$$= \mathbb{E}_{s \sim \rho^{\mu_\phi}} [\nabla_\phi \mu_\phi(s) \nabla_a Q_{\mu_\phi}(s, \mu_\phi(s))|_{a=\mu_\phi(s)}] \quad (27)$$

where the chain rule is used in the second line [66]. A challenge with both these approaches is determining how to accurately estimate the Q-value at each optimisation step as it directly coupled the policy.

3.3 Actor-Critic Methods

In actor-critic architectures the above policy gradient algorithms are split into two coupled components. The actor updates the policy parameters ϕ by performing the gradient ascent in Eqs. (25-26). To address the concern regarding coupling with Q-value estimation, a critic is introduced using different parameters θ so that $Q_\theta(s, a) \leftarrow Q_{\pi_\phi}(s, a)$. By the compatible function approximation theorem, the reparameteristaion is only exact if i) θ is linear in π_ϕ , and ii) θ is obtained by minimising the MSE error $\sim (Q_{\pi_\phi} - Q_\theta)^2$ [65, 66, 118]. These assumptions are very often relaxed, such as when using deep neural networks as universal function approximators. In practice this leads to brittle weights but works well in practice after several additional stability improvements such as introducing target networks [67, 71, 75, 77, 78].

To learn critic values for both stochastic and deterministic actors, Q-learning is used as a form of temporal difference control [63]. Using the deterministic policy as an example, we know from the Bellman equation [119, 120] that the valuation for the current state-action pair $Q_\theta(s, a)$ is related to the value of the next subsequent state-action pair $Q_\theta(s', a')$. This allows us to construct a target state-action value

$$Q_{\bar{\theta}}(s, a) \leftarrow \mathbb{E}_{s' \sim \rho^{\mu_{\bar{\theta}}}} [r(s, \mu_{\bar{\theta}}(s)) + \gamma Q_{\bar{\theta}}(s', \mu_{\bar{\theta}}(s'))] \quad (28)$$

where the target parameters weights \bar{w} are obtained through delayed Polyak averaging $\bar{w} \leftarrow \tau w + (1-\tau)\bar{w}$ with $\tau \ll 1$. The critic then minimises the difference between these two theoretically equivalent values. The overwhelming majority of literature represents this with standard MSE loss (Bellman error) objective

$$J(\theta) = \mathbb{E}_{s \sim \rho^{\mu_\phi}} [(Q_{\bar{\theta}}(s, a) - Q_\theta(s, \mu_{\bar{\theta}}(s)))^2] \quad (29)$$

This approach led to famous the Deep Deterministic Policy Gradient (DDPG) algorithm [67] along with several other modifications over the years [68–70]. The contemporary stand-out successor is the Twin Delayed DDPG (TD3) algorithm [71]. TD3 improves on DDPG in three key areas: 1. Introduces clipped double-Q learning and uses the minimum of two Q-values for target values to address critic overestimation bias. 2. Delays policy, target policy, and target critic updates to occur every second step. 3. Adds noise to target policy as a form of regularisation to prevent the formation of brittle policies.

Off-policy learning for TD3 utilises a experience replay buffer \mathcal{D} containing tuples (s, a, r, s') of all past steps across all training episodes up to a maximum buffer size. Actor and critic optimisation per respective gradient step involves uniformly sampling a mini-batch of N transitions $U(\mathcal{D})$ from the buffer to construct

$$J(\theta) = \mathbb{E}_{(s, a, r, s') \sim U(\mathcal{D})} [(Q_{\bar{\theta}}(s, a) - Q_\theta(s, a))^2] \quad (30)$$

$$J(\phi) = -\mathbb{E}_{(s,a,r,s') \sim U(\mathcal{D})} [Q_\theta(s, a)] \quad (31)$$

where the negative in the second line is to force gradient ascent. The vectors of parameters ϕ and θ are then adjusted to minimise the empirical losses. Note two explicit assumptions, firstly that using MSE is preferred over a myriad of other loss functions that potentially might offer greater resistance to outliers. The Monte-Carlo or law of large numbers approach is used where as $N \rightarrow \infty$ the sampling of $U(\mathcal{D})$ approaches the true unknown distribution provided the underlying data is i.i.d. and stationary. In most situations this does not hold since in practice a mini-batch of size of a few hundred transitions is always used and policy optimisation by definition is non-stationary process and hence the underlying data is not i.i.d. over training time.

The Kolmogorov theorem on the strong law of large numbers for non-i.i.d. data does however still guarantee this approaches validity provided the mini-batch variable has finite variance [121]. It is worth pointing out that this finite variance assumption applies directly to the true unknown distribution, not the mini-batch as finite sampling will always produce finite moments [9]. For example, if the underlying distribution is Cauchy then it has both undefined variance and mean, if it is Pareto then it also can have undefined variance (if shape $\alpha < 2$) and mean (if $\alpha < 1$). This will be discussed further in Section 3.6.

3.4 Soft Actor-Critic

A alternative formulation based on conditional energy-based models involves the combined use of soft values $V_\theta^{\text{soft}}(s)$ and $Q_\theta^{\text{soft}}(s, a)$, and stochastic policies $\pi_\phi(a|s)$ to maximise a causal entropy objective [72]. Soft generally refers to the value functions being defined by application of softmax functions over all actions whether in discrete or continuous action spaces [72]. As deterministic polices such as TD3 heuristically explore via the injection of noise when selecting actions, stochastic polices have the ability to navigate multi-modal objectives that are common in robotic environments [74]. For example, in situations where two actions appear equally attractive, the policy will commit equal probability mass rather than deterministic selection.

The additive discounted cumulative future reward from a time t is then $R_t = \sum_{k=t}^{\infty} \gamma^k (r(s_k, a_k) + \alpha H(\pi_\phi(\cdot|s_k)))$ using the objective $J(\phi) = \mathbb{E}_{s \sim \rho, a \sim \pi_\phi} [R_t]$. The temperature α is a automatically tuned entropy regularising hyperparameter that controls the relative weighting or importance given to stochastic behaviour and $H(x) = \mathbb{E}_{x \sim P(x)}[-\ln P(x)]$ is the entropy at each state.

The policy is very generally represented as $\pi_\phi(a|s) \propto \exp(\frac{1}{\alpha} Q_\theta^{\text{soft}}(s, a))$ with the definitions

$$Q_\theta^{\text{soft}}(s_t, a) \equiv r_t + \mathbb{E} \left[\sum_{k=t}^{\infty} \gamma^{k-t} (r(s_k, a_k) + \alpha H(\pi_\phi(\cdot|s_k))) \right] \quad (32)$$

$$V_\theta^{\text{soft}}(s_t) \equiv \alpha \ln \int_{\mathcal{A}} da' e^{\frac{1}{\alpha} Q_\theta^{\text{soft}}(s, a')} = \alpha \ln \mathbb{E}_{Z_\theta} \left[\frac{\exp(\frac{1}{\alpha} Q_\theta^{\text{soft}}(s_t, a'))}{Z_\theta(a')} \right] \quad (33)$$

where $Z_\theta(a)$ is the partition function normalising the distribution being independent of policy and does not contribute to the gradients [72, 75, 76]. Notice that $V_\theta^{\text{soft}}(s_t)$ is by definition a softmax function. The optimal policy is very generally expressed by

$$\pi_\phi(a|s) \equiv \exp \left(\frac{1}{\alpha} (Q_\theta^{\text{soft}}(s, a) - V_\theta^{\text{soft}}(s)) \right) \quad (34)$$

with convergence guaranteed using fixed-point iteration [72]. The soft Q-value also satisfies the soft Bellman equation

$$Q_\theta^{\text{soft}}(s_t, a) = r_t + \gamma \mathbb{E}_{s' \sim \rho^{\pi_\phi}} [V_\theta^{\text{soft}}(s')] \quad (35)$$

with the (hard) Bellman equation recovered in the zero entropy policy limit $\alpha \rightarrow 0$ [72]. Note we can also rewrite $V_\theta^{\text{soft}}(s_t) = \mathbb{E}_{a \sim \pi_\phi} [Q_\theta^{\text{soft}}(s_t, a_t) - \alpha \ln \pi_\phi(a_t | s_t)]$.

Stochastic action sampling from the policy $\pi_\phi(\cdot | s_t)$ is defined over some fixed distribution $a_t \sim f_\phi(\epsilon_t; s_t)$. In the case of a Gaussian, the number of output nodes of the ϕ network is $2 \times \dim(\mathcal{A})$ where each actions probability density function is completely defined by a unique mean and standard deviation. One additional modification we require for the sampled actions to be numerically bounded involves a change of variable. Consider action sampling from an unbounded distribution $v_t \sim f_\phi(\epsilon_t; s_t)$ such as a Gaussian, to enforce strict symmetric bounds we can use $a = \tanh(v)$ and the new density can be expressed as $\pi(a|s) = \nu(v|s)|\det(da/dv)|^{-1}$ where the Jacobian is $da/dv = \text{diag}(1 - \tanh^2(v))$ [75, 77]. The transformed log-likelihood is then given by

$$\log \pi(a_t | s_t) = \log \nu(v_t | s_t) - \sum_{j=1}^{\dim(\mathcal{A})} \log (1 - \tanh^2 v_{t,j}). \quad (36)$$

One method of improving the policy $\pi_\phi \rightarrow \pi'_\phi$ using an information-theoretic approach is by directly minimising the Kullback-Leibler divergence $D_{\text{KL}}(\pi_\phi(\cdot | s_t) || \exp(\frac{1}{\alpha} Q_\theta^{\text{soft}}(s_t, \cdot)))$ to give

$$\pi'_\phi(\cdot | s_t) = \arg \max_{\pi_\phi} \mathbb{E} \left[Q_\theta^{\text{soft}}(s_t, a_t) - \alpha \ln \pi_\phi(\cdot | s_t) \mid \pi_\phi \right] \quad (37)$$

which is reminiscent of an advantage function with the baseline being the average across actions soft Q-value.

For automatic entropy adjustment for the α coefficient, we solve the constraint problem of $\max \mathbb{E}_{a_t \sim \pi_\phi} [R_t]$ under $\mathbb{E}_{a_t \sim \pi_\phi} [-\ln \pi_\phi] \geq \bar{H}$. The minimum desired expected entropy of any environment is heuristically set as the cardinality of the action space $\bar{H} = -|\mathcal{A}|$ [77]. Using gradient descent from convex optimisation as an approximation, after updating to π'_ϕ , the optimal solution is $\alpha_t^* = \arg \min_{\alpha_t} \mathbb{E}_{a_t \sim \pi'_\phi} [-\alpha_t (\ln \pi'_\phi + \bar{H})]$ [77, 122].

The Soft Actor-Critic (SAC) algorithm can be expressed at every step the agent takes as sequentially updating three parameters with following objectives to be minimised

$$J(\theta) = \mathbb{E}_{(s, a, r, s') \sim U(\mathcal{D})} \left[\frac{1}{2} \left(Q_\theta^{\text{soft}} - Q_{\bar{\theta}}^{\text{soft}} \right)^2 \right] \quad (38)$$

$$J(\phi) = \mathbb{E}_{(s, a, r, s') \sim U(\mathcal{D})} \left[\alpha \ln \pi_\phi - Q_\theta^{\text{soft}} \right] \quad (39)$$

$$J(\alpha) = \mathbb{E}_{(s, a, r, s') \sim U(\mathcal{D})} \left[-\alpha (\ln \pi_\phi + \bar{H}) \right] \quad (40)$$

where the target action-value function in Eq. (38) is constructed using Eq. (35). SAC also utilises the clipped double-Q learning feature of TD3 for updating both ϕ and θ networks. SAC has the strong advantage of minimal hyperparameters pertaining largely to the neural network architecture. SAC-Discrete [79] is also available that requires minimal changes to policy optimisation where $\pi_\phi(a|s)$ outputs a probability rather than a density.

3.5 Robust Critic Evaluation

The use of the MSE loss function in Eqs. (30) and (38) is deeply entrenched in all algorithms. We propose to investigate the effect of 10 different loss functions. Firstly we examine the effect of MSE and higher even powers where $\mathbb{E}_{U(\mathcal{D})} [(Q_{\bar{\theta}} - Q_{\theta})^{2+n}]$ for $n = 0, 2, 4, 6$. The effects of this loss amplification is to see whether giving substantially larger weights to outliers is beneficial to learning. Outliers in this case represent situations in the mini-batch where the difference between a particular samples current value of a state and its target value is excessively large. Amplifying these effects will force the optimiser to heavily modify the responsible parameters in the neural network. The other six loss functions presented in [85, 123–128] represent different degrees of outlier detection and reduction.

There has been considerable research into the merit these functions in NMF factorisation with a summary and literature review available in [85]. Briefly, NMF is formulated with $V - WH$ where V is the actual (fixed) matrix and WH is the learned representation, the difference is then minimised. While these dictionary learning methods have largely been shelved in preference for deep learning, [85] reveals how truncated Cauchy NMF is vastly superior to existing methods while being computational intensive. We propose to loosely connect this with the difference $Q_{\bar{\theta}} - Q_{\theta}$. An issue is that $Q_{\bar{\theta}}$ is also a learned quantity that varies across time unlike the matrix V . Therefore, TD3 performance should be more stable due to delaying target network updates to occur every second step unlike SAC.

Another argument for shifting away from MSE also based on its exponent is that it resembles standard deviation. An extremely insightful discussion in [9, 129] traces the origin of standard deviation (STD) to its historical roots. They find that preference of STD to mean absolute deviation (MAD) in the statistical sciences is due to a dispute in the 1920s where STD is shown to be 12.5% more asymptotically efficient than MAD if and only if the underlying data set is normally distributed. Under the presence of even minuscule ‘fat tails’, MAD is proven to be overwhelmingly more efficient. Relevant to our discussion, MAD functionally resembles the mean absolute error (MAE) loss function. Since the underlying agent data is also very unlikely to be Gaussian there exists a strong basis to investigate MAE.

Explicitly, we investigate the following functions arranged in increasing levels of outlier suppression

$$\mathbb{E}_{U(\mathcal{D})} [\text{MSE}_n(Q_{\bar{\theta}}, Q_{\theta})] = \frac{1}{N} \sum_{i=1}^N (Q_{\bar{\theta}}^i - Q_{\theta}^i)^{2+n}, \quad \text{for } n = 0, 2, 4, 6 \quad (41)$$

$$\mathbb{E}_{U(\mathcal{D})} [\text{Huber}(Q_{\bar{\theta}}, Q_{\theta})] = \frac{1}{N} \sum_{i=1}^N \begin{cases} \frac{1}{2}(Q_{\bar{\theta}}^i - Q_{\theta}^i)^2, & \text{if } |Q_{\bar{\theta}}^i - Q_{\theta}^i| < 1 \\ |Q_{\bar{\theta}}^i - Q_{\theta}^i| - \frac{1}{2}, & \text{otherwise} \end{cases} \quad (42)$$

$$\mathbb{E}_{U(\mathcal{D})} [\text{MAE}(Q_{\bar{\theta}}, Q_{\theta})] = \frac{1}{N} \sum_{i=1}^N |Q_{\bar{\theta}}^i - Q_{\theta}^i| \quad (43)$$

$$\mathbb{E}_{U(\mathcal{D})} [\text{HSC}(Q_{\bar{\theta}}, Q_{\theta})] = \frac{1}{N} \sum_{i=1}^N \left(\sqrt{1 + (Q_{\bar{\theta}}^i - Q_{\theta}^i)^2} - 1 \right) \quad (44)$$

$$\mathbb{E}_{U(\mathcal{D})} [\text{Cauchy}(Q_{\bar{\theta}}, Q_{\theta}, \omega)] = \frac{1}{N} \sum_{i=1}^N \ln \left(1 + \left(\frac{Q_{\bar{\theta}}^i - Q_{\theta}^i}{\omega} \right)^2 \right) \quad (45)$$

$$\mathbb{E}_{U(\mathcal{D})} [\text{TCauchy}(Q_{\bar{\theta}}, Q_{\theta}, \omega, \xi)] = \frac{1}{N} \sum_{i=1}^N \begin{cases} \ln \left(1 + \omega^{-2} (Q_{\bar{\theta}}^i - Q_{\theta}^i)^2 \right), & \text{if } 0 \leq (Q_{\bar{\theta}}^i - Q_{\theta}^i)^2 \leq \omega^2 \xi \\ \ln(1 + \xi), & \text{otherwise} \end{cases} \quad (46)$$

$$\mathbb{E}_{U(\mathcal{D})} [\text{CIM}(Q_{\bar{\theta}}, Q_{\theta}, \sigma)] = \frac{1}{N} \sum_{i=1}^N \left(1 - \frac{1}{\sqrt{2\pi}\sigma} \exp \left(\frac{-(Q_{\bar{\theta}}^i - Q_{\theta}^i)^2}{2\sigma^2} \right) \right) \quad (47)$$

where HSC, TCauchy, and CIM refer to hypersurface cost-based, truncated Cauchy, and correntropy induced metric

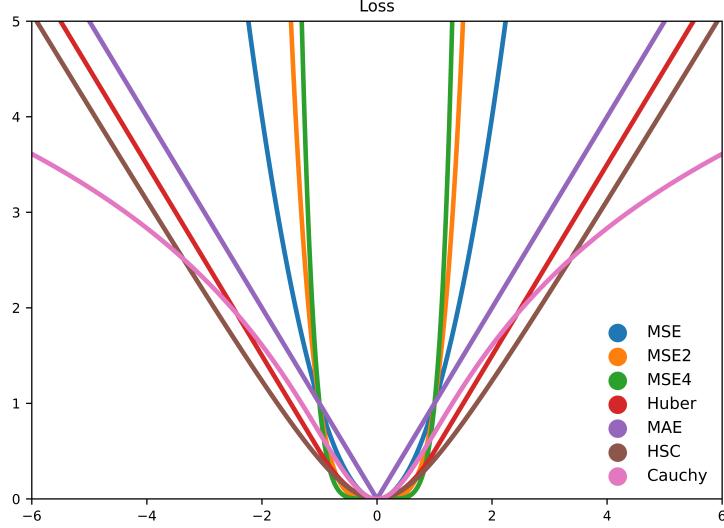


Figure 6: Several of the critic loss function in Eqs. (41-47).

loss functions, respectively.

The scale parameters are empirically evaluated at each training step with the Cauchy scale parameter ω obtained iteratively via the Nagy algorithm [126] with

$$\omega_{t+1} = \omega_t \cdot \left[\left(\frac{1}{N} \sum_{i=1}^N \left(1 + \left(\frac{Q_\theta^i - Q_\theta^i}{\omega_t} \right)^2 \right)^{-1} \right)^{-1} - 1 \right]^{1/2} \quad (48)$$

where $\omega_0 > 0$ and is functionally similar to the previous scale multiplied by a harmonic mean. The CIM kernel size σ is calculated as the empirical standard deviation of the mini-batch [128] with

$$\sigma_t^2 = \frac{1}{N} \sum_{i=1}^N \left((Q_\theta^i - Q_\theta^i) - \mu_t \right)^2 \quad (49)$$

where $\mu_t = \frac{1}{N} \sum_{i=1}^N (Q_\theta^i - Q_\theta^i)$ is the empirical mean. Behaviour of these parameters can be independently evaluated using any loss function. The Cauchy truncation level ξ in Eq. (46) is empirically determined [85] using the 3σ rule where mini-batch samples $|Q_\theta^i - Q_\theta^i| - \mu_t | > 3\sigma_t$ exceeding this value are set to zero $|Q_\theta^i - Q_\theta^i| = 0$ and ignored. In Fig. 6 we provide visualisations for some of the functions.

3.6 Preasymptotics and the Tail Exponent

A key assumption present throughout machine learning and statistics discussed in Section 3.3 is that the Monte-Carlo approach is valid if the data is i.i.d. by the law of large numbers. Reinforcement learning by definition is non-i.i.d. but the Kolmogorov theorem on the strong law of large numbers reveals that convergence is still guaranteed provided the true unknown underlying distribution Ω has finite variance σ^2 [121]. If this is true for a random variable $X \sim \Omega$, we have the condition

$$\frac{1}{M} \sum_{i=1}^M X_i - \frac{1}{M} \sum_{i=1}^M \mathbb{E}_{X \sim \Omega}[X_i] \xrightarrow{M \rightarrow \infty} 0 \quad (50)$$

and so in practice the universal standard is to simply utilise the Monte-Carlo approach

$$\mathbb{E}_{X \sim \Omega}[X] = \frac{1}{M} \sum_{i=1}^M X_i \quad (51)$$

This is especially important in machine learning where parameter optimisation is based entirely on either minimising aggregate losses or maximising aggregate gains. Furthermore, this is usually performed using mini-batch learning $N \subset M$ as it is far more computationally efficient. The mini-batch is often uniformly sampled $N \sim U(M)$ at each learning iteration to ensure overall it is reflective of the actual sample M as the number of iterations increases.

A gigantic and monumental question that appears to have never been seriously asked regarding this method is whether this remains valid for finite M [9]. As the real-world contains only finite sample sizes, the convergence may never formally occur. This problem is even more critical in machine learning as we utilise $N < M$ and so the equality is even more debatable. Essentially, we obtain Ω_N , treat it as an accurate and efficient reflection of Ω_M , which is then assumed to be a perfect representation of the true Ω . In Eq. (51) the left is referred to as the true population mean μ and the right is called the empirical mean \bar{x} . Clearly empirical means are not empirical for $N < M \ll \infty$.

The question is then in what regimes does utilising \bar{x} as substitute for μ remain appropriate and to what degree, that is, what is the preasymptotic behaviour of \bar{x} ? This problem is one of the primary focuses of [9] which provides explicit answers to all these and many more questions. Their focus is mainly on the statistical consequences of fat tail, namely distributions Ω with far larger kurtosis than the Gaussian. As discussed in Section 1.3, these probability density functions are solely defined by their infrequent extreme outliers and are more accurate representations of almost all real-world systems. Hence using \bar{x} is uninformative as ‘empirically’ it will be heavily biased by the noise and fail accurately account for shocks.

Detailed discussion on power laws (fat-tailed distributions) and the Generalised Pareto distribution (GPD) is outside the scope of this work and so we present the final results without proof [9, 130–133]. What is relevant is that the tail exponent α governs the fatness of the tails, lower implies less thin, and is the characteristic feature of a power law. Furthermore, the moment of order p for $X \sim \text{GPD}$ only exists if $\alpha > p$. Hence if $\alpha < 1$ all moments: mean, variance, skew, kurtosis, and higher are formally undefined and cannot be estimated even in the limit $M \rightarrow \infty$. One extreme is the Gaussian with $\alpha \rightarrow \infty$, while Cauchy has $\alpha < 1$. Note MAD only requires the mean to exist.

A very important point is that since sample is finite, even if $\alpha < 1$, we can still ‘empirically’ calculate all moments. These point estimates are therefore entirely misleading and lead to a very false sense of confidence about the nature of distribution Ω_M let alone Ω . A more conservative approach to evaluating them involves first determining the tail exponent α for the sample, only then can we crudely gauge the true nature of Ω .

Assuming a strictly right-tailed fat distribution capped with $X_i \geq 0 \forall i$, we are very likely to obtain $\bar{x} \ll \mu$ for $M < \infty$ as seen in Fig. 7. Using this estimate will eventually lead to a rude awakening when a outlier naturally appears. To counter this gross underestimation [9] propose first constructing a new distribution Ω'_M using the M known samples that exhibits the correct right-tail behaviour. The mean of this distribution $\mathbb{E}_{X \sim \Omega'_M}[X] = \mu_s$ called the ‘shadow’ mean is then considered a far more accurate representation of the true mean μ .

To derive Ω'_M for the finite set of observations $\omega = \{X_1, \dots, X_M\}$ first define the observed maximum $\zeta = \max(\omega)$. Next we require the finite support $X \in [L, H]$ where $L \geq 0$ and H is set to be exceedingly large but finite so that the probability of observing it is minuscule so that $\zeta \ll H < \infty$ shown in Fig. 8. Another feature is defining a threshold $L^* \geq L$ where we are uninterested in the region $X < L^*$ as these occurrences are treated as noise.

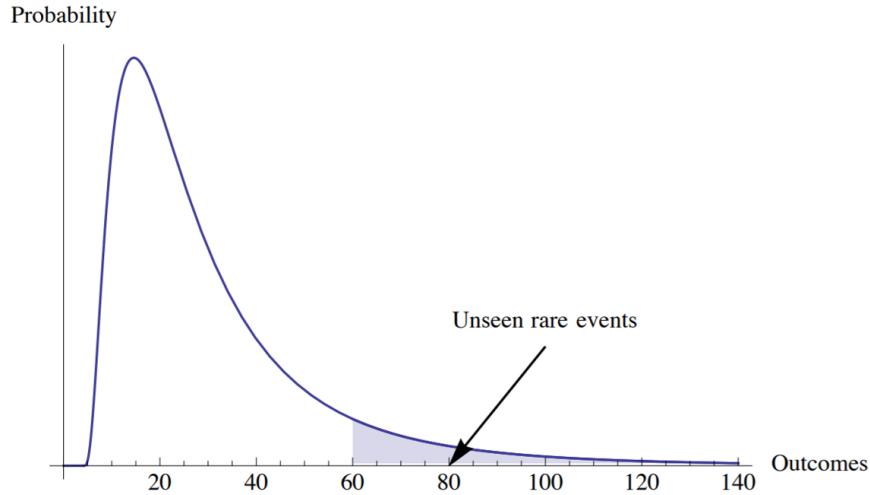


Figure 7: Typical structure of a fat-tailed distribution where the absence of rare large outliers in the obtained finite sample leads to severely underestimated empirical means. Adapted from [9].

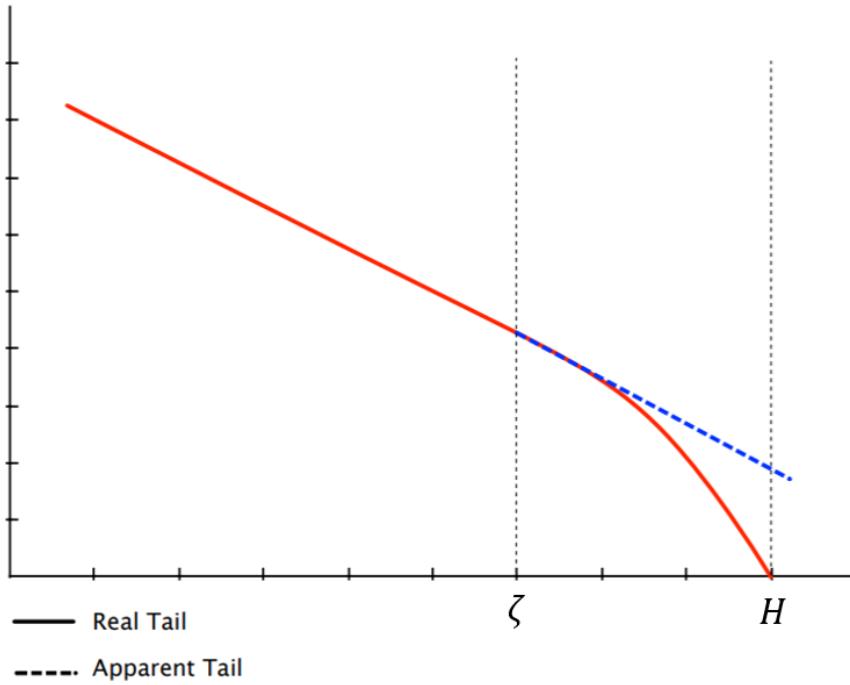


Figure 8: Graphical representation of why a very large finite upper support $H \gg \zeta$ is required to ensure both numerical stability and realistic power law decay when estimating a shadow mean from a finite sample. For all practical purposes the existence of H is ignored as divergence is only evident when we approach the limit, hence using the apparent tail is considered a reasonable approximation. Adapted from [9].

For random variable with a known support $X \in [L, H]$ define a smooth function $\varphi(X)$ requiring $\varphi \in C^\infty$ (differentiable to all degrees), $\varphi^{-1}(\infty) = H$, and $\varphi(L) = \varphi^{-1}(L) = L$, parameterised non-uniquely as

$$\varphi(X) = L - \ln \left| \frac{H-X}{H-L} \right| = L + \ln \left| \frac{H-L}{H-X} \right| \quad (52)$$

Then define a new random variable $Z \equiv \varphi(X)$ with bounds $Z \in [L, \infty)$ and for very large H we can approximate $Z \approx X$ given the first-order Taylor series expansion $\ln|y| \approx y - 1$ about $y = 1$. Therefore, the tail structures in Fig. 8 of the bounded X (real tail) and unbounded Z (apparent tail) are identical up till the vicinity of huge H . Hence the goal is to model the tail of Z and then convert back to X with the inverse

$$X = \varphi^{-1}(Z) = H + (L-H)e^{\frac{L-Z}{H}} \quad (53)$$

as a non-unique method to estimate the shadow moments of the distribution. Note any of the moments of the Z distribution do not formally exist as it is formally unbounded.

Next as we are only interested in the right tail, we focus strictly on values exceeding a interest threshold u where $u = L^* \geq L$. Then define a random variable w as the excess $w \equiv X - u$. In the limit $u \rightarrow \infty$, the cumulative $D(w; \alpha, \varsigma)$ and probability $w \sim d(w; \alpha, \varsigma)$ density functions for $w \geq 0$ can be approximated by a GDP where

$$D(w; \alpha, \varsigma) = \begin{cases} 1 - \left(1 + \frac{w}{\alpha\varsigma}\right)^{-\alpha}, & \text{if } \alpha < \infty \\ 1 - e^{-\frac{w}{\varsigma}}, & \text{otherwise} \end{cases} \quad (54)$$

$$d(w; \alpha, \varsigma) = \frac{1}{\varsigma} \left(1 + \frac{w}{\alpha\varsigma}\right)^{-\alpha-1}, \quad w \in [L^*, \infty) \quad (55)$$

with tail exponent $\alpha \in (-\infty, \infty)$ and the scale parameter $\varsigma \in (0, \infty)$ estimated from the data using a several possible approaches [130–133]. In this case the we have the equivalence of the of both distributions $X \sim f(X\alpha, \varsigma)$ and $Z \sim d(Z; \alpha, \varsigma)$ where

$$\int_{L^*}^{\varphi^{-1}(\infty)} dX f(X - \varphi^{-1}(L^*); \alpha, \varsigma) = \int_{L^*}^{\infty} dZ d(Z - L^*; \alpha, \varsigma) = 1 \quad (56)$$

by definition. Note in particular for both cases the excess difference from the threshold L^* is taken to be compatible with the formulation of the GDP. To find $f(X - \varphi^{-1}(L^*); \alpha, \varsigma)$ we solve the integral equation by substituting Eq. (52) into Eq. (55) and after recalling $\varphi^{-1}(L) = L$ we use the ansatz

$$f(X - \varphi^{-1}(L^*); \alpha, \varsigma) = \frac{H}{\varsigma(H-X)} \left(1 + \frac{H}{\alpha\varsigma} \ln \left| \frac{H-L}{H-X} \right| \right)^{-\alpha-1}, \quad X \in [L^*, H] \quad (57)$$

with integration yielding

$$\begin{aligned} \lim_{B \rightarrow H} \int_{L^*}^B dX f(X - \varphi^{-1}(L^*); \alpha, \varsigma) &= - \lim_{B \rightarrow H} \left(1 + \frac{H}{\alpha\varsigma} \ln \left| \frac{H-L^*}{H-B} \right| \right)^{-\alpha} \Big|_{X=L^*}^{X=B} \\ &= 1 - \lim_{B \rightarrow H} \left(1 + \frac{H}{\alpha\varsigma} \ln \left| \frac{H-L^*}{H-B} \right| \right)^{-\alpha} \\ &= 1 - \left(1 + \frac{H}{\alpha\varsigma} \left(\ln |H-L^*| - \lim_{B \rightarrow H} \ln |H-B| \right) \right)^{-\alpha} \end{aligned} \quad (58)$$

which converges to unity for $\alpha > 0$. For our purposes as we are not overly interested in the behaviour in the vicinity

of H , we can also assume the second term approaches zero before the limit is reached where $\lim_{B \rightarrow E} \frac{H-L^*}{H-B} \rightarrow \infty$ for $\zeta \ll E \ll H < \infty$. Therefore, as $f(X\alpha, \zeta)$ and $d(Z; \alpha, \zeta)$ are one-to-one transformations on each other, they must share the same tail exponent and scale parameter.

The shadow moments of order p of a fat-tailed random variable X conditional on $X \geq L^*$ are then generally

$$\mathbb{E}_{X \sim \Omega'_M} [X^p | X > L^*] \equiv \int_{L^*}^H dX X^p f(X - L^*; \alpha, \zeta) \quad (59)$$

with our desired shadow mean being

$$\begin{aligned} \mathbb{E}_{X \sim \Omega'_M} [X | X > L^*] &= \int_{L^*}^H dX X f(X - L^*; \alpha, \zeta) \\ &\stackrel{(a)}{=} -\frac{H}{\zeta} \lim_{B \rightarrow H} \int_{H-L^*}^{H-B} dt \frac{(H-t)}{t} \left(1 + \frac{H}{\alpha\zeta} \ln \left| \frac{H-L}{t} \right| \right)^{-\alpha-1} \\ &\stackrel{(b)}{=} \lim_{B \rightarrow H} \left(1 + \frac{H}{\alpha\zeta} \ln \left| \frac{H-L^*}{t} \right| \right)^{-\alpha} \\ &\quad \times \left[\alpha(H-L^*) e^{\frac{\alpha\zeta}{H}} \left(\frac{\alpha\zeta}{H} + \ln \left| \frac{H-L^*}{t} \right| \right)^\alpha \Gamma \left(-\alpha, \frac{\alpha\zeta}{H} + \ln \left| \frac{H-L^*}{t} \right| \right) - H \right] \Big|_{t=H-L^*}^{t=H-B} \\ &= \lim_{B \rightarrow H} \left[\alpha(H-L^*) e^{\frac{\alpha\zeta}{H}} \left(\frac{\alpha\zeta}{H} \right)^\alpha \Gamma \left(-\alpha, \frac{\alpha\zeta}{H} + \ln \left| \frac{H-L^*}{t} \right| \right) - H \left(1 + \frac{H}{\alpha\zeta} \ln \left| \frac{H-L^*}{t} \right| \right)^{-\alpha} \right] \Big|_{t=H-L^*}^{t=H-B} \\ &\stackrel{(c)}{=} H - \alpha(H-L^*) e^{\frac{\alpha\zeta}{H}} \left(\frac{\alpha\zeta}{H} \right)^\alpha \Gamma \left(-\alpha, \frac{\alpha\zeta}{H} \right) \\ &\stackrel{(d)}{=} H - (H-L^*) \left[1 - e^{\frac{\alpha\zeta}{H}} \left(\frac{\alpha\zeta}{H} \right)^\alpha \Gamma \left(1-\alpha, \frac{\alpha\zeta}{H} \right) \right] \\ &= L^* + (H-L^*) e^{\frac{\alpha\zeta}{H}} \left(\frac{\alpha\zeta}{H} \right)^\alpha \Gamma \left(1-\alpha, \frac{\alpha\zeta}{H} \right) \end{aligned} \quad (60)$$

where $\Gamma(s, x)$ for strictly $s > 0$ and $x \geq 0$ is the upper incomplete gamma function defined as

$$\Gamma(s, x) \equiv \int_x^\infty dy y^{s-1} e^{-y} \quad (61)$$

and in (a) we perform the change of variable $t = H - X$ and so $dt = -dX$, (b) involves trivial integration by inspection or alternatively, plebeians can perform lengthy and cumbersome manipulation of known results obtained from tables [134] wherein

$$\begin{aligned} \int dx \frac{1-x}{x(1-\ln|x|)^d} &= (1-\ln|x|)^{1-d} \left(\frac{1}{d-1} - e \int_1^\infty dt \frac{e^{-(1-\ln|x|)t}}{t^d} \right) + C \\ &= (1-\ln|x|)^{1-d} \left(\frac{1}{d-1} - e(1-\ln|x|)^{d-1} \Gamma(1-d, 1-\ln|x|) \right) + C \end{aligned} \quad (62)$$

$$= \frac{(1-\ln|x|)^{1-d}}{d-1} - e\Gamma(1-d, 1-\ln|x|) + C \quad (63)$$

with e as Euler's number, in (c) for the upper bound we use the same logic as Eq. (58) and also the limit

$$\lim_{B \rightarrow E} \Gamma \left(-\alpha, \frac{\alpha\zeta}{H} + \ln \left| \frac{H-L^*}{H-B} \right| \right) \rightarrow 0 \quad (64)$$

for $L^* \ll \zeta \ll E \ll H$ as the upper support need never be reached, while the lower bound results in all logarithms yielding zero, and finally (d) uses the known recurrence relation $\Gamma(1+a, x) = a\Gamma(a, x) + x^a e^{-x}$ taught throughout kindergarten [135]. In what follows for simplicity we assume $\zeta = 1$ as its computationally intense estimation [130] is

not practical for mini-batch reinforcement learning over millions of training iterations.

Using then the provided sample ω , the shadow mean of the known distribution Ω'_M is expressed as

$$\mathbb{E}_{X \sim \Omega'_M} [X | X > L^*] = L^* + (H - L^*) e^{\frac{\alpha}{H}} \left(\frac{\alpha}{H} \right)^\alpha \Gamma \left(1 - \alpha, \frac{\alpha}{H} \right) \quad (65)$$

where $\mathbb{E}_{X \sim \Omega'_M} [X | X > L] = \mathbb{E}_{X \sim \Omega'_M} [X]$. Furthermore, given that by construction we should have $\frac{\alpha}{H} \ll 1$, the fact that numerical integration of the usual (complete) gamma function $\Gamma(s)$ is far more common and hence easier to perform, and the known relation with the lower incomplete gamma function $\gamma(s, x)$ where $\Gamma(s) = \Gamma(s, x) + \gamma(s, x)$, it may be easier to numerically evaluate the difference

$$\Gamma \left(1 - \alpha, \frac{\alpha}{H} \right) = \Gamma(1 - \alpha) - \int_0^{\frac{\alpha}{H}} dy y^{-\alpha} e^{-y} \quad (66)$$

$$\approx \int_0^{\infty} dy y^{-\alpha} e^{-y} \quad (67)$$

where great care must be taken when using the approximation as it should only be considered if $\frac{\alpha}{H} \rightarrow 0$ with certainty. Therefore, a smaller α (larger kurtosis of sample ω) correctly increases the value obtained from integration by both reducing the impact of $\gamma(s, x)$, and resulting in greater divergence of $y^{-\alpha} e^{-y}$ near the origin.

This entire approaches rests on α , commonly represented in literature with the scale parameter as $\gamma = \alpha^{-1}$, and whose estimation is no easy task as there are bountiful methods for obtaining an empirical value. Extreme value theory offers general approaches involve using MLE or the method of moments [130]. In ω let $X_{1,M} \leq \dots \leq X_{M,M}$ be the order statistics, then define

$$H_{k,M}^{(j)} \equiv \frac{1}{k} \sum_{i=0}^{k-1} (\ln |X_{M-i,M}| - \ln |X_{M-k,M}|)^j \quad (68)$$

where $X_{M-k,M}$ is an empirically determined intermediate order statistics used as threshold for when the power law dominates the distribution. Importantly, for asymptotic reasons we take $1 \leq k < M$ so that for $k \rightarrow \infty$, $k/n \rightarrow 0$ as $n \rightarrow \infty$ by the strong law of large numbers [130]. However, since we are constrained with requiring $0 < \gamma < 1$, the simplest method is the Hill estimator [136] where

$$\begin{aligned} \hat{\gamma} &= H_{k,M}^{(1)} = \frac{1}{k} \left[\ln \left| \prod_{i=0}^{k-1} X_{M-i,M} \right| - k \ln |X_{M-k,M}| \right] \\ &= \ln \left| \left(\prod_{i=0}^{k-1} X_{M-i,M} \right)^{-k} \right| - \ln |X_{M-k,M}| \end{aligned} \quad (69)$$

observing that this represents the difference between the natural logarithm of geometric means. The method of moments yields the following estimator

$$\hat{\gamma} = H_{k,M}^{(1)} + 1 - \frac{1}{2} \left(1 - \frac{(H_{k,M}^{(1)})^2}{H_{k,M}^{(2)}} \right)^{-1} \quad (70)$$

to then obtain $\hat{\alpha} = \hat{\gamma}^{-1}$ [137, 138]. This entire process rests on the correct selection of the k extreme values. Selecting small k causes large variance in $\sigma_\gamma^2 \sim \gamma^2$ while for large k the entire sample is unlikely to adhere to a power law. The choice of k is then dependent on the properties of the underlying process X and appears a very subjective decision.

There are however a plethora of methods of doing so that are not exact but offer a more objective approach, however these are also subject to other bias. Many of these approaches require expensive computational optimisation and so are not feasible for mini-batch learning.

There is though far more simpler approach that is consistent with the derivation of Eq. (65) as the GPD assumption assumes a very strong power law in a region where it totally dominates the behaviour. At this point the distribution can be simplified by approximating it as Pareto with $\hat{\gamma}$ estimated as the slope of a Zipf or Pareto quantile (Q-Q) plot that is expected to be nearly linear [139–142]. This is done by plotting

$$(r_k, s_k) = \left(\ln \left| \frac{M+1}{k} \right|, \ln |X_{M-k+1,M}| \right), \quad k = 1, \dots, M \quad (71)$$

with its gradient approximately equal to $\hat{\gamma}$. A more refined version involves allowing the power law to only dominate at the infinite limit where the generalised quantile plot becomes

$$(r_k, s_k) = \left(\ln \left| \frac{M+1}{k} \right|, \ln |X_{M-k+1,M} \cdot H_{k,M}^{(1)}| \right), \quad k = 1, \dots, M \quad (72)$$

which is ultimately linear for small k . The slope is estimated with the usual linear regression coefficient

$$\hat{\gamma} = \frac{\sum_{i=1}^k (r_i - \bar{r})(s_i - \bar{s})}{\sum_{i=1}^k (r_i - \bar{r})^2} \quad (73)$$

noticing that by construction we will always obtain $\hat{\gamma} \geq 0$ as we are using ordered statistics.

Furthermore, using Zipf plots suppose we construct an estimate using $k = b \ll c$ points we will generally find that $\hat{\gamma}_b \geq \hat{\gamma}_c$ or $\hat{\alpha}_b \leq \hat{\alpha}_c$ if the sample actually contains some extreme values. This is because the values are sorted in terms of largest to smallest and so including fewer values will lead to steeper gradients as there are fewer points r_1, \dots, r_k with tighter spacing. In other words, we will see comparable (albeit a little less) rise, for less run, causing a greater incline in the line of best fit. Instead of then tuning for appropriate k , if we take the complete sample $k = M$ we are able to construct the upper bound for the true value $\hat{\alpha} \leq \hat{\alpha}_M$. Additionally, if $\hat{\alpha}_M < 1$ using Eq. (71), then it likely must also be the case using Eq. (72). Hence if $\hat{\alpha}_M < 1$, it will be satisfactory to state that the sample contains extreme values and so justification for using the shadow mean is strong.

Suppose we have either determined the optimal intermediate $X_{M-k,M}$ or have used either of the Zipf plots to obtain $\hat{\gamma} = \hat{\alpha}^{-1}$. For $\hat{\alpha} \in (0, \infty)$ use of Eq. (65) is considered acceptable. If $0 < \alpha < 1$ use of Eq. (51) is entirely inappropriate as the sample appears to exhibit fat tails with an undefined traditional mean and so only the shadow mean should be utilised. If $\hat{\alpha} \leq 0$, then we have no choice but to estimate with the empirical mean in Eq. (51). Overall, we then have the deterministic shadow mean estimate

$$\mu_s(L^*, H, \hat{\alpha}) \equiv L^* + (H - L^*) e^{\frac{\hat{\alpha}}{H}} \left(\frac{\hat{\alpha}}{H} \right)^{\hat{\alpha}} \int_{\frac{\hat{\alpha}}{H}}^{\infty} dy y^{-\hat{\alpha}} e^{-y} \quad (74)$$

In machine learning it is an open question as to whether this approach could potentially be utilised. Most common are situations in supervised learning where $X \sim (Y - \bar{Y})^2$ are ‘empirical’ losses from a mini-batch and then aggregated using Eq. (51). In this case clearly $X_i \geq 0 \forall i$ and so we have a situation with unbounded right-tailed skew. Perhaps this preasymptotic approach may lead to superior results.

In reinforcement learning we have critic and soft-critic losses of the form $X \sim (Q_\theta - Q_{\bar{\theta}})^2$ which certainly might

be amenable to this approach. The fact that this represents the difference between quantities that are only separated by one or two learning iterations, it is likely that it might have somewhat stable properties. In this case, the extreme values will be isolated solely to samples in the mini-batch where valuations are most divergent. Using the shadow mean would then amplify the impact of these outliers and force the optimisation process to place greater emphasis on adjusting the network parameters to account for them. This can therefore be considered somewhat analogous to the case in Eq. (41) for $n > 0$ amplifying the effect of these outliers. This approach however is far more sophisticated as it does not inflate the smaller values.

In this case the backpropagation of the graph will center on adjusting parameters so the tail exponent $\hat{\alpha}$ increases leading to a reduced shadow mean μ_s . Furthermore, for each mini-batch, a systematic scheme for estimating the supports $[L^*, H]$ is required. Setting $L^* = L = 0$ intuitively appears to the strictest requirement as therefore we are interested in minimising all critic errors regardless of how minor they are in a given mini-batch. Selection of H is a more tricky endeavour, as we require $\zeta \ll H < \infty$, we must utilise both the known sample maximum ζ and the our domain knowledge to estimate its value.

Theoretically, uniformly sampling experiences (s_j, a_j, r_j, s'_j) for $j = 1, \dots, N$ from the replay buffer \mathcal{D} during training iteration t , the critic loss $X_{j,t} = (Q_{\theta,t}(s_j, a_j) - Q_{\bar{\theta},t}(s'_j, a'))^2$ are naturally uncapped $X \in [0, \infty)$ and form the mini-batch $\omega_t = \{X_{1,t}, \dots, X_{N,t}\}$. However we also expect as the number of training iterations increases that $X_{j,t+\delta t} \leq X_{j,t}$ for $\delta t \gg 0$ as indication of successful agent learning. Hence we demand that the upper support is dependant on both time and the known largest loss $\zeta_t = \max(\omega_t)$ contained in the mini-batch $H_t = H(t, \zeta_t)$. We can further simplify this approach by assuming that ζ_t completely encapsulates the time-dependence so that $H_t = H(\zeta_t)$. Recall again that the exact value of the finite maximum loss H is not required as we do not ever expect it to be approached. Therefore a very simple approach would be to utilise a constant linear multiplier $H_t = \lambda_t \zeta_t$ for $\lambda_t \in (1, \infty)$ and its exact value is empirically determined on a case-by-case basis.

One crude method of potentially gauging a lower bound λ_t^l would be to estimate the time-dependent value required for equivalence between both empirical and shadow means

$$\frac{1}{N} \sum_{j=1}^N X_{j,t} = L^* + (\lambda_t^l \zeta_t - L^*) e^{\frac{\hat{\alpha}_t}{\lambda_t^l \zeta_t}} \left(\frac{\hat{\alpha}_t}{\lambda_t^l \zeta_t} \right)^{\hat{\alpha}_t} \int_{\frac{\hat{\alpha}_t}{\lambda_t^l \zeta_t}}^{\infty} dy y^{-\hat{\alpha}_t} e^{-y} \quad (75)$$

Numerically solving this transcendental equation would likely find $\lambda_t^l \leq 1$ if ω_t is fat-tailed with $0 < \hat{\alpha}_t < 1$. Then examining the time-dependent structure of λ_t^l for successful agents would allow us to gain insight on the overall process. Generally a larger λ_t^l would necessarily demand a higher estimate of the finite support $H_t = \lambda_t \zeta_t$ where $\lambda_t^l \ll \lambda_t \forall t$. Overall, this is a ad hoc procedure and in practice setting a large constant for all time $\lambda = \lambda_t$ is likely the most efficient approach.

Regarding actor losses, deterministic policy gradient are of the form $X \sim -Q_\theta$ and soft-policies are $X \sim \alpha \ln \pi_\phi - Q_\theta$. The former seems applicable if one applies the negative factor at the final step, while latter is unlikely to be bounded by zero. Finally, entropy temperature $X \sim -\alpha (\ln \pi_\phi + \overline{H})$ might also be feasible with this approach as we can again factor out the negative multiplier.

Therefore, to utilise this approach there are several hyperparameters to tune in contrast with global standard in Eq. (51). For actors and each of the twin critics (and also temperature if using SAC) we need to determine suitable L^* and H , and also have a systematic procedure to identify the intermediate order statistics $X_{M-k,M}$ for each network

or utilise the Zipf plot. None of these selections are remotely trivial. Use of Zipf plots therefore is likely the superior choice as the computations are routine. An additional complication when mini-batch learning over many iterations is that these selections must be done at each iteration and so the procedure needs to be computational efficient.

3.7 Multi-Step Returns

Multi-step targets or returns [62] up to the m -step ahead of current step t are expressed as

$$Q_{\bar{\theta},t}^{(m)}(s,a) \equiv \sum_{k=0}^{m-1} \gamma^k r_{t+k} + \gamma^m Q_{\bar{\theta}}(s_{t+m}, a') \quad (76)$$

where the action a' is sampled from either the SAC or TD3 policies. One advantage to this approach is that tuning for appropriate m can potentially rapidly accelerate learning [29, 39, 62, 89].

3.8 History Decision Processes

Consider a general environment with finite observations $o \in \mathcal{O}$ where $\mathcal{S} \in \mathcal{O}$ which is typical for partially observed non-Markovian decision processes. Often due to the complexity or the sheer size of \mathcal{O} , the agent is unable to directly learn from this space and so has to utilise \mathcal{S} when interacting with the environment. To perform this task the agent must have access to a mapping from the complete observation space to the operating access space [105–109].

As the agent continues through time in each episode, it constructs tuples of histories $h_t \in \mathcal{H}_t : (\mathcal{O} \times \mathcal{R} \times \mathcal{A})^{t-1} \times \mathcal{O}$ where $h_t \equiv o_1 a_1 r_1 \dots o_{t-1} a_{t-1} r_{t-1} o_t$ or $h_t \equiv h_{t-1} a_{t-1} r_{t-1} o_t$. Explicitly, the trajectory sequentially develops as $o_1 \rightarrow a_1 \rightarrow r_1 o_2 \rightarrow a_2 \rightarrow r_2 o_3 \rightarrow \dots \rightarrow a_{t-1} \rightarrow r_{t-1} o_t$. Note this definition is different to the one presented in [107–109] but the results remain valid. The set of all finite histories is denoted by $\mathcal{H}^* = \bigcup_t \mathcal{H}_t$ and the empty set is ϵ . Histories also need not be unique, but as the agent becomes more successful, we typically expect the length of these histories to increase in environments where there is no time limit. The transition probability at any time t to the next state o_{t+1} is a function of the history-action pair (h_t, a_t) not the state-action pair (s_t, a_t) . This leads to the notion of history-based decision process (HDPs) with transition probabilities $P : \mathcal{H}^* \times \mathcal{A} \rightsquigarrow \mathcal{O} \times \mathcal{R}$ where \rightsquigarrow denotes a stochastic mapping [109].

This allows us to define $R_t \equiv \sum_{k=0}^{\infty} \gamma^k r(s_k, a_k)$ for stochastic policies where the state value and action-value functions are $V^\pi(h_t) \equiv \mathbb{E}[R_t | h_t; \pi]$ and $Q_\pi(h_t, a) \equiv \mathbb{E}[R_t | h_t, a; \pi]$. The optimal values are $V^*(h_t) = \max_\pi V_\pi(h_t)$ and $Q^*(h_t, a_t) = \max_\pi Q_\pi(h_t, a_t)$ where $\pi^* \in \arg \max_\pi V_\pi(\epsilon)$. Note that the optimal policy may not be unique and so is denoted as an element. New Bellman equations can be written

$$Q_\pi(h_t, a_t) = \sum_{o_{t+1}, r_t} P(o_{t+1}, r_t | h_t, a_t) [r_t + \gamma V_\pi(h_{t+1})] \quad (77)$$

$$V_\pi(h_t) = Q_\pi(h_t, \pi(h_t)) \quad (78)$$

$$Q^*(h_t, a_t) = \sum_{o_{t+1}, r_t} P(o_{t+1}, r_t | h_t, a_t) [r_t + \gamma V^*(h_{t+1})] \quad (79)$$

$$V^*(h_t) = \max_{a_t \in \mathcal{A}} Q^*(h_t, a_t) \quad (80)$$

$$\pi^*(h_t) \in \arg \max_{a_t \in \mathcal{A}} Q^*(h_t, a_t) \quad (81)$$

where they are considered pseudo-recursive and not self-consistent as length of h_{t+1} is larger than h_t and so an

algorithm based on frequency of visits cannot be used [107]. Utilising this approach is impractical as the cardinality $|\mathcal{H}^*|$ is extremely large since the histories are unlikely to ever repeat. To handle this situation [107] introduce a surjective aggregation feature mapping $\phi : \mathcal{H} \rightarrow \mathcal{S}$ so that $s = \phi(h) \in \mathcal{S}$ where \mathcal{S} is by definition finite and small enough to work within. The history is then reduced $h_t \equiv h_{t-1}a_{t-1}r_{t-1}o_t \rightarrow h_t \equiv h_{t-1}a_{t-1}r_{t-1}s_t$ where the (s_t, a_t) pair is sufficient. This obviously describes a MDP however [107] find its applicability to be far more general.

For the feature map ϕ , the transition probabilities are constructed via marginalisation

$$P_\phi(s_{t+1}, r_t | h_t, a_t) = \sum_{o_{t+1}: \phi(h_t a_t r_t o_{t+1}) = s_{t+1}} P(o_{t+1}, r_t | h_t, a_t) \quad (82)$$

where $P_\phi \in \text{MDP}$ and therefore is stationary if $\exists p : P_\phi(s_{t+1}, r_t | h_t, a_t) = p(s_{t+1}, r_t | s_t, a_t) \forall \phi(h_t) = s_t$ where $p : \mathcal{S} \times \mathcal{A} \rightsquigarrow \mathcal{O} \times \mathcal{R}$. The exact aggregation for $P_\phi \in \text{MDP}$ can then be shown to yield $V_\pi(h_t) = V_\pi(s_t)$, $Q_\pi(h_t, a_t) = Q_\pi(s_t, a_t)$ where $\pi(h_t) = \pi(s_t)$, and $V^*(h_t) = V^*(s_t)$, $Q^*(h_t, a_t) = Q^*(s_t, a_t)$ where $\pi^*(h_t) = \pi^*(s_t)$ [107]. This allows us to recover the well-known Bellman equations

$$Q_\pi(s_t, a_t) = \sum_{s_{t+1}, r_t} p(s_{t+1}, r_t | s_t, a_t) [r_t + \gamma V_\pi(s_{t+1})] = \sum_{s_{t+1}} p(s_{t+1} | s_t, a_t) \mathbb{E}[r_t | s_t, a_t] \quad (83)$$

$$V_\pi(s_t) = Q_\pi(s_t, \pi(s_t)) \quad (84)$$

$$Q^*(s_t, a_t) = \sum_{s_{t+1}, r_t} p(s_{t+1}, r_t | s_t, a_t) [r_t + \gamma V^*(s_{t+1})] = \sum_{s_{t+1}} p(s_{t+1} | s_t, a_t) \mathbb{E}[r_t | s_t, a_t] \quad (85)$$

$$V^*(s_t) = \max_{a_t \in \mathcal{A}} Q^*(s_t, a_t) \quad (86)$$

$$\pi^*(s_t) \in \arg \max_{a_t \in \mathcal{A}} Q^*(s_t, a_t) \quad (87)$$

that are identical to standard reinforcement learning and are widely amenable to optimal solutions using the usual iterative frequency-based convergence [29, 30, 143–145].

For the general case with transition probabilities P not necessarily forming MDPs, approximate aggregation results reveal that they can also be modelled by MDPs [107]. They find several insightful results with the most relevant to our discussion being that assuming

$$|Q^*(h, a) - Q^*(\tilde{h}, a)| \leq \epsilon \quad \forall \phi(h) = \phi(\tilde{h}) \quad \forall a, \quad (88)$$

the bounds for the optimal value functions for some $\epsilon, \gamma > 0$ can be proven to be

$$|Q^*(h, a) - Q^*(s, a)| \leq \frac{\epsilon}{1 - \gamma} \quad (89)$$

$$|V^*(h) - V^*(s)| \leq \frac{\epsilon}{1 - \gamma} \quad (90)$$

with the condition $\epsilon = 0$ if $\pi^*(h) = \pi^*(s)$. Determining the exact structure of this feature map ϕ however is in general a highly non-trivial process requiring significant calibration for each environment. They go further discussing extreme state aggregation wherein any process P can theoretically be represented using small finite-state MDPs. In particular they also present an argument explaining the surprisingly robust performance of MDP-based Q-learning methods when applied to non-MDP domains is very likely due to the generality of the law of large numbers [121] in the infinite sampling limit.

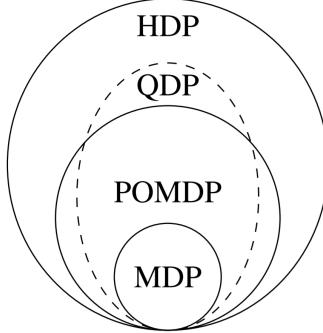


Figure 9: Q-Value Uniform Decision Processes (QDPs) in terms intersections with a broader class of more well-known processes. Adapted from [109].

To convert HDPs into practical algorithms we must introduce additional definitions and constraints [109]. The state process p_h for an $\phi(h) = s$ is similarly

$$p_h(s_{t+1}, r_t | s_t, a_t) = \sum_{o_{t+1}: \phi(h_t a_t r_t o_{t+1}) = s_{t+1}} P(o_{t+1}, r_t | h_t, a_t) \quad (91)$$

where for any $h \neq \tilde{h}$ clearly $p_h(s', r' | s, a) \neq p_{\tilde{h}}(s', r' | s, a)$ in general. This encodes the nature of non-MDPs where path-dependent histories are essential to future decisions and can incorporate non-stationary domains. If the state process is a MDP then p_h is independent of history and standard convergence result in Eqs. (83-87) follow.

To keep the state process history-dependent while making Q-learning independent of history, [109] introduce a subset of HDPs called Q-Value Uniform Decision Processes (QDPs) shown in Fig. 9 that incorporate all MDPs, have a non-empty intersection with all POMDPs, and are defined by a state-uniformity condition. The condition specifies that for any action, if any two histories h, \tilde{h} map to the same state s i.e. $\phi(h) = \phi(\tilde{h}) = s$, then the optimal Q-value of the underlying HDP of these histories is identical $Q^*(h, a) = Q^*(\tilde{h}, a)$ [109]. This class is defined by Eq. (88) where $\epsilon = 0$ as $\pi^*(h) = \pi^*(s)$, and therefore by Eq. (89) the Q-values $Q^*(h, a) = Q^*(s, a)$ are state-uniform. QDPs then allow modelling of non-stationary domains, which accurately represents agent learning, and can be interpreted as history-independent Q-values of POMDPs.

The proof of QDP Q-learning convergence assumes that the state-process is ergodic where in all states are reachable under any policy from the current state after sufficiently many steps, meaning different unique histories can all reach the same state while the underlying HDP process P is still clearly non-ergodic [109]. This point is worth repeating, only the theoretical optimal values need to be equivalent at convergence that may never be reached, all intermediate values need not be equal at any fixed point or time. Also note for practical purposes the history of an agent is already stored within the experience replay buffer.

3.9 Agent Performance Evaluation

The current reward aggregation paradigm in all of reinforcement learning utilises additive dynamics when assessing the performance of the agent over an episode [29]. This means that the final cumulative reward is composed of the independent summation of rewards received at every time step $\Gamma^+(s_T, a_{T-1} | h_T) = \sum_{t=1}^T r_t$ where the “+” exponent indicates additive scheme. Improvements from tweaking existing algorithms, and comparison between models are

generally measured using this method under two approaches. The first is generally a tabular presentation of averaged scores and standard deviations

$$\bar{\Gamma}^+ \equiv \frac{1}{n} \sum_{i=1}^n \Gamma^+(s_T^i, a_{T-1}^i | h_T^i) \quad (92)$$

$$\sigma^+ \equiv \left(\frac{1}{n} \sum_{i=1}^n (\Gamma^+(s_T^i, a_{T-1}^i | h_T^i) - \bar{\Gamma}^+)^2 \right)^{1/2} \quad (93)$$

across a range of environments using the final trained agent across n runs with each generating a unique history h_t^i . The second approach is similar to the first and usually graphically displays the average scores and standard deviations across evaluation episodes occurring at fixed intervals during agent training to better highlight performance over training time.

Before proceeding we clearly define what is meant by comparable performance. In practice this is done by conducting N evaluation episodes every fixed training interval using the identical parameters weights ϕ_r, θ_r for policies π_{ϕ_r} or μ_{ϕ_r} and action-values Q_{θ_r} . This process is then repeated for M trials where the agent is trained from scratch again likely generating different weights ϕ, θ each time due to random initialisation where the theoretical optimal weights are only reached in the infinite time limit guaranteeing identical parameters. Ultimately this process results in $n = NM$ unique histories at each evaluation interval whose performances are summarised using Eqs. (92-93). The reason for this design is that agent performance over multiple runs is notorious for being brittle and sensitive to hyperparameter selection [146, 147].

This universally accepted procedure is perfectly reasonable provided that the nature of the environment can be characterised by summing independent uncorrelated reward signals at each time step, and if the volatility in agent learning can be represented by a “standard” deviation. In general, it is unlikely that accurate performance measurement in all conceivable environments can be reduced to this way. Furthermore, quantifying the scale of volatility of a metric using standard deviation (STD) is only acceptable to alternative measures under strict requirements. As discussed in Section 3.5, a origin story for near-universal use of STD in all areas of science is presented in [9, 129] where they reveal it is only the preferred measure if the underlying data is also normally distributed. For all other situations, mean absolute deviation (MAD) is superior as it is far more asymptotically efficient, meaning that it is more robust to ‘fat tails’, and it does not require the variance of the true unknown underlying distribution to be finite. Therefore, it would be unwise to blindly assume the Gaussian relationship $\Gamma^+(s_T^i, a_{T-1}^i | h_T^i) \sim N(\bar{\Gamma}^+, (\sigma^+)^2)$. Mean deviation MAD defined as

$$\text{MAD}^+ \equiv \frac{1}{n} \sum_{i=1}^n |\Gamma^+(s_T^i, a_{T-1}^i | h_T^i) - \bar{\Gamma}^+| \quad (94)$$

should therefore be used as the default volatility measure instead. Noting in particular the bound $\sigma^+ \geq \text{MAD}^+$. Other common names for MAD are mean absolute error (MAE), and average deviation (AVEDEV) used in Microsoft Excel. It is also functionally equivalent to the L_1 -norm used throughout machine learning as in Eq. (43).

4 Non-Ergodicity in Reward Accumulation

The additive reward scheme is not in general appropriate for accurately modelling rewards in all environments. There exists a vast array of environments where the amount of risk-taking (leverage) can either be amplified or reduced in order to magnify or shrink potential rewards, where time order matters, and losses have an asymmetrically larger effect on performance compared to equally sized gains. These are domains where simply maximising the total reward is not the objective, rather the goal is to maximise the total reward while avoiding steep losses. Much of the applications for such environments are in non-ergodic domains where the time average is not equal to expectation or ensemble value as seen in Section 1.

To construct reward signals for agents in multiplicative dynamic environments there are three possibilities dependent on what information the agent receives. Firstly, as these situations involve change in valuations across time steps, an initial portfolio value $V_0 \neq 0$ must be specified as a baseline from which all future returns are measured against. The signal received by the agent at a time t can then be in terms of: absolute reward received r_t so the new valuation is $V_t = (V_{t-1} + r_t)/V_{t-1} = 1 + r_t/V_{t-1}$, directly as a growth percentage g_t so that $V_t = V_{t-1}(1 + g_t)$, or as a multiplier λ_t that gives $V_t = V_{t-1} \times \lambda_t$. In what follows, this work will focus solely on the first type of signal.

This formulation allows the valuation to crash towards zero if a substantial negative rewards are encountered. This allows us to create a ‘game over’ or bankruptcy criterion where for example, if $V_\tau \leq V_{\min}$, the episode ends at $t = \tau$ with history $h_\tau = h_{\tau-1}a_{\tau-1}r_{\tau-1}s_\tau$. This is essentially the stop-loss from Section 1 but we refrain from using this terminology to keep the discussion general. To evaluate agent performance, we then measure the compounding changes in valuation over all time steps in a evaluation episode. For both dynamics we can write cumulative rewards up to a final step t as

$$\Gamma_t^+(s_t, a_{t-1}|h_t) = \Gamma_t^+ \equiv \sum_{k=1}^t r(\Gamma_{k-1}^+, s_k, a_k) = \Gamma_{t-1}^+ + r_t \quad (95)$$

$$1 + \Gamma_t^\times(s_t, a_{t-1}|h_t) = (1 + \Gamma_t^\times) \equiv \prod_{k=1}^t \frac{\Gamma_{k-1}^+ + r(\Gamma_{k-1}^+, s_k, a_k)}{\Gamma_{k-1}^+} = (1 + \Gamma_{t-1}^\times) \cdot \frac{\Gamma_{t-1}^+ + r_t}{\Gamma_{t-1}^+} \quad (96)$$

for $t > 0$ with the initial value $\Gamma_0^+ = V_0 > V_{\min}$. Note to ease computations we can also decompose the product $\Gamma_t^\times = \exp[\sum_{k=1}^t \ln|1 + r_k/\Gamma_{k-1}^+|] - 1$. This allows explicit encoding of the ‘game over’ condition at $t = \tau$, if $r_\tau \leq (V_{\min} - \Gamma_{\tau-1}^+)$ then the episode is over and we set compounding return to be $\Gamma_{t \geq \tau}^\times = 0$.

We strongly highlight that in general, the rewards $r(\Gamma_{k-1}^+, s_k, a_k)$ do not necessarily equal the standard rewards discussed in Section 3.1 so that $r(\Gamma_{k-1}^+, s_k, a_k) \neq r(s_k, a_k)$. This how we explicitly encode past performance and the concept of leverage into future rewards. Functionally, the tuple (Γ_{t-1}^+, s_t) act as a new effective ‘state’ of the system combining the previous cumulative reward with the current environment.

For Eq. (96) to be a valid representation we must have the crude relation where, in general, $|r(\Gamma_{k-1}^+, s_k, a_k)| \ll \Gamma_{k-1}^+$, meaning that the reward at an time $t = k$ should be somewhat comparable the cumulative reward. This is achieved when the agent has access to Γ_{k-1}^+ as another state of the environment and so the ‘aggressiveness’ of the next action to maximise the reward is dependent on the existing value. We define a successful policy in the limit $t \rightarrow \infty$ as one that exhibits $\Gamma_{t-1}^+ \rightarrow \infty$ where the agent is more comfortable in taking more risky actions that have greater potential in maximising the reward but are also correspondingly capable of causing larger absolute losses.

For a fixed successful policy we then have strict requirement

$$\lim_{t \rightarrow \infty} |r(\Gamma_{t-1}^+, s_k, a_t)| \gg |r(\Gamma_{t-1}^+, s_k, a_t)| \quad (97)$$

to ensure the reward appropriately scales with the cumulative value due partly to more aggressive actions taken for the same underlying state. We refer to this as the asymptotic reward scaling condition. Environments where this condition cannot be envisioned to occur are not multiplicative and should be treated as additive.

We can also rewrite the product in the form of compounding returns as

$$(1 + \Gamma_t^\times) \equiv \prod_{k=1}^t \left(1 + \frac{r(\Gamma_{k-1}^+, s_k, a_k)}{\Gamma_{k-1}^+} \right) = (1 + \Gamma_{t-1}^\times) \cdot \left(1 + \frac{r(\Gamma_{k-1}^+, s_k, a_k)}{\Gamma_{t-1}^+} \right) \quad (98)$$

$$\ln |1 + \Gamma_t^\times| \equiv \sum_{k=1}^t \ln \left| 1 + \frac{r(\Gamma_{k-1}^+, s_k, a_k)}{\Gamma_{k-1}^+} \right| = \ln |1 + \Gamma_{t-1}^\times| + \ln \left| 1 + \frac{r(\Gamma_{t-1}^+, s_t, a_t)}{\Gamma_{t-1}^+} \right| \quad (99)$$

Note without any loss of generality we can use $\Gamma_t^\times \leftarrow V_0 (1 + \Gamma_t^\times)$ when comparing overall performance. Regarding the two other types of reward signals we can also equivalently write $(1 + \Gamma_t^\times) = \prod_{k=1}^t (1 + g_t) = \prod_{k=1}^t \lambda_t$.

For any general environment the agents performance for episodes is ergodic if it satisfies the Birkhoff theorem

$$V_0(1 + \bar{g}) \equiv \underbrace{\lim_{t \rightarrow \infty} \frac{1}{t} \int_0^t dt \Gamma(s_t^j, a_{t-1}^j | h_t^j)}_{\text{Time average of } \Gamma} \stackrel{?}{=} \underbrace{\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \Gamma(s_i^j, a_{t-1}^i | h_i^j)}_{\text{Expectation value of } \Gamma} \equiv \mathbb{E}[\Gamma] \quad (100)$$

where the equality implies the systems is ergodic [5, 148, 149]. The left-hand side states that for any random history h_t^j we calculate the cumulative reward in the infinite time limit or till the episode ends and take the average over time. The right-hand side represents the ensemble average across infinite histories at any points in time. To test whether this holds in general, take $h_t^j = h_{t-1} a_{t-1} r_{t-1} s_t$ with $h_T^j = h_t^j a_t r_t s_T$ meaning the game ends at the next step. Simplifying notion $\Gamma(s_t^j, a_{t-1}^j | h_t^j) = \Gamma(h_t^j)$, for both the additive and multiplicative dynamics we have

$$\frac{1}{T} \left(\Gamma^+(h_t^j) + r_T \right) \stackrel{?}{=} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \Gamma^+(h_i^j) \quad (101)$$

$$\frac{1}{T} \left((1 + \Gamma^\times(h_t^j)) \times \frac{\Gamma^+(h_t^j) + r_T}{\Gamma^+(h_t^j)} \right) \stackrel{?}{=} \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n (1 + \Gamma^\times(h_i^j)) \quad (102)$$

As we are free to select any h_t^j arbitrarily, by contradiction, it is unlikely that either equality holds in general for any non-trivial environment. Hence under both dynamics, agent learning is very likely a non-ergodic process where the time-average performance of any single path should not be assumed to exactly equal the average performance of all paths at any time. The difference is again the volatility tax $\nu = \bar{g} - \left(\frac{\mathbb{E}[\Gamma]}{V_0} - 1 \right)$.

Returning to the discussion in Section 3.9 regarding how in practice $n = NM$ where M is the amount total training trials and N is the amount of evaluations episodes per trial. We can see that as $n \rightarrow \infty$, if M is held constant the ergodic case is far more likely to be approached compared to if N is held constant. This is again due to the fact that the agent policy ϕ and action value θ parameters are shared for each M . For all theoretical results in this work, whenever such limits are considered we assume $N = 1$ and $M \rightarrow \infty$ as this is more in line with reality.

Regarding evaluation on exactly the same environments via $\bar{\Gamma}$ and $\text{MAD}(\Gamma)$ we can also hypothesise several results

comparing agent training on both the dynamics. We expect the following three occurrences

$$\text{MAD}(\Gamma^+) \geq \text{MAD}(\Gamma^\times) \quad (103)$$

$$\text{MAD}_u(\Gamma^+) \leq \text{MAD}_u(\Gamma^\times) \quad (104)$$

$$\text{MAD}_d(\Gamma^+) \geq \text{MAD}_d(\Gamma^\times) \quad (105)$$

where the up-side and down-side MADs are calculated considering only the $\Gamma_j \geq \bar{\Gamma}$ and $\Gamma_j \leq \bar{\Gamma}$ respectively. This is because multiplicative dynamics is more suitable for steady growth expect there to be less overall volatility while also anticipating there exists more positive, and fewer negative deviations from the mean. Overall, this is consistent with stable increasing learning while avoiding of steep losses.

4.1 Ergodic Special Case

There is one scenario where the system can become ergodic if the agent has learned the optimal deterministic policy $\mu(s) \rightarrow \mu^*(s)$ at some time $t \rightarrow t^*$. During on- or off-policy model-free learning for all times $t \geq t^*$, Then the agent will take identical actions when presented identical states. If the game has no natural end and the agent can continue for $t \rightarrow \infty$, the agent states and actions will convergence to a optimal sequence. Under this special case, for some critical time $t_c \gg t^*$ we can say that

$$\lim_{T \rightarrow \infty} \frac{1}{T - t_c} \int_{t_c}^T dt \Gamma(h_t^j) \approx \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \Gamma(h_{t \geq t_c}^i) \quad (106)$$

where the approximation is very loose and assumes that majority of the cumulative rewards are obtained from times $t > t_c$. Another way this can be seen is through the situation where longer-term rewards are more significant for both dynamics

$$\Gamma^+(h_T^j) = \Gamma^+(h_{t < t_c}^j) + \Gamma^+(h_{t \geq t_c}^j) \approx \Gamma^+(h_{t \geq t_c}^j) \quad (107)$$

$$1 + \Gamma^\times(h_T^j) = (1 + \Gamma^\times(h_{t < t_c}^j)) \cdot (1 + \Gamma^\times(h_{t \geq t_c}^j)) \approx 1 + \Gamma^\times(h_{t \geq t_c}^j) \quad (108)$$

for all histories. Given the purpose of multiplicative dynamics is to penalise losses significantly we can also likely assume that at least $t_c^\times \geq t_c^+$ in general. This then allows us to assume ergodicity amongst all histories wherein we can crudely state that $\lim_{t \rightarrow \infty} \frac{1}{t} \Gamma(h_t^j) \approx \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \Gamma(h_i^j)$ as both limits are functionally equivalent. Therefore we obtain in the limit $\nu \rightarrow 0$ for any randomly selected history h_j .

This situation could be considered an implicit assumption made in all existing literature results since it reasonable under stationary MDP convergence proofs. The current formulation of reinforcement learning using additive reward signals in MDP environments can therefore be described as ergodic if we assume the agent has learned and is operating with an optimal deterministic policy at all times. For non-MDP processes, stochastic policies, and sub-optimal policies, all bets are off, and the system should be considered non-ergodic.

5 Q-Learning with Multiplicative Dynamics

Recall standard Q-learning the agent performs an action-value update to the (s, a) -estimate using the rule

$$Q_{t+1}(s, a) \leftarrow (1 - \alpha_t(s, a)) Q_t(s, a) + \alpha_t(s, a) \left(r_t(s, a) + \gamma \max_{a'} Q_{t+1}(s', a') \right) \quad (109)$$

where $\alpha_t(s, a) > 0$ is the time- and (s, a) -dependent learning rate under a ϵ -greedy policy using Bellman's principle of optimality [119, 120]. Essentially, the second term performs a small adjustment to the Q-value in the direction of optimal value. The convergence of this approach to unique optimal fixed point is very well-known result [29, 30, 63, 84, 143–145].

For additive dynamics we can adjust $r_t(s, a) \leftarrow r_t(s, a) + \Gamma_{t-1}^+$ where Γ_{t-1}^+ is a known constant and all the results in [107, 109] remain unchanged. This is valid as the next reward to be maximised is treated independent from the previous cumulative rewards.

5.1 Model-Free Return Maximisation

To incorporate multiplicative dynamics into this formulation we must resort to using the QDPs discussed in Section 3.8 in order to retain the histories of rewards while keeping the evaluation of Q-values dependent only on the (s, a) -pair in Markovian fashion.

Recalling the definitions in Section 3.7 that the rewards are now also dependent on the cumulative previous additive reward where $r(\Gamma_{t-1}^+, s_t, a_t) \leftarrow r(s_t, a_t)$ in order to internally control leverage by factoring the existing value along with the usual environment. This way the agent is able to keep track of final outcome of all past actions summarised with a single scalar. For this to be occur the asymptotic reward scaling condition in Eq. (97) must universally hold. As Γ_{t-1}^+ is always known at any time t we are free to define a new environment state $\xi_t = s_t \cup \Gamma_{t-1}^+$ with the fixed cardinality relation $|\xi_t| = |s_t| + 1$ and formally $\xi_t \in \Xi : \mathcal{S} \times \mathcal{R}$.

The agent history is then $h_t \equiv \xi_1 a_1 r_1 \dots \xi_{t-1} a_{t-1} r_{t-1} \xi_t$ where $h_t \in \mathcal{H}_t : (\Xi \times \mathcal{R} \times \mathcal{A})^{t-1} \times \Xi$ remains unchanged and the trajectory develops as $\xi_1 \rightarrow a_1 \rightarrow r_1 \xi_2 \rightarrow a_2 \rightarrow r_2 \xi_3 \rightarrow \dots \rightarrow r_{t-1} \xi_t$. Practically, the only change required is the optimisation procedure also uses the current additive reward as an input which treated as an additional state. In what follows we will always utilise the notion $r_t = r(h_t, a_t)$ for convenience, it should never be confused for $r_t = r(s_t, a_t)$ and it will be made clear where the latter applies.

To reformulate Q-learning to be compatible with multiplicative dynamics and prove convergence we must derive results from the most basic levels again. The value functions are then defined as $V^\pi(h_t) \equiv \mathbb{E}[1 + R_t^\times | h_t; \pi]$ and $Q_\pi(h_t, a) \equiv \mathbb{E}[1 + R_t^\times | h_t, a; \pi]$, the optimal values are $V^*(h_t) = \max_\pi V_\pi(h_t)$ and $Q^*(h_t, a_t) = \max_\pi Q_\pi(h_t, a_t)$ where $\pi^* \in \arg \max_\pi V_\pi(\epsilon)$. In this case the discounted future compounding rewards are

$$\begin{aligned} 1 + R_t^\times &\equiv (1 + \Gamma_{t-1}^\times) \cdot \frac{\Gamma_{t-1}^+ + r_t}{\Gamma_{t-1}^+} \cdot \frac{(\Gamma_{t-1}^+ + r_t) + \gamma r_{t+1}}{\Gamma_{t-1}^+ + r_t} \cdot \frac{(\Gamma_{t-1}^+ + r_t + \gamma r_{t+1}) + \gamma^2 r_{t+2}}{\Gamma_{t-1}^+ + r_t + \gamma r_{t+1}} \cdot \dots \\ &= (1 + \Gamma_t^\times) \cdot \frac{\Gamma_t^+ + \gamma r_{t+1}}{\Gamma_t^+} \cdot \frac{(\Gamma_t^+ + \gamma r_{t+1}) + \gamma^2 r_{t+2}}{\Gamma_t^+ + \gamma r_{t+1}} \cdot \dots \\ &= (1 + \Gamma_t^\times) \cdot \prod_{k=1}^{\infty} \frac{\Gamma_{t-1+k}^+ + \gamma^k r_{t+k}}{\Gamma_{t-1+k}^+} \end{aligned} \quad (110)$$

where $\Gamma_{t-1+k}^+ = \Gamma_{t-1}^+ + \sum_{\lambda=0}^{k-1} \gamma^\lambda r_{t+\lambda}$ is the discounted additive reward up to time step $t = k - 1$, the current return

Γ_t^\times is completely known, and in Section 5.2 we show that the initial value $\Gamma_0^+ = V_0 \geq \gamma > V_{\min}$ is required. We can then rewrite this more clearly as

$$\begin{aligned} 1 + R_t^\times &\equiv (1 + \Gamma_t^\times) \cdot \left(1 + \frac{\gamma r_{t+1}}{\Gamma_t^+}\right) \cdot \left(1 + \frac{\gamma^2 r_{t+2}}{\Gamma_{t+1}^+}\right) \cdot \left(1 + \frac{\gamma^3 r_{t+3}}{\Gamma_{t+2}^+}\right) \cdots \\ &= (1 + \Gamma_t^\times) \cdot \prod_{k=1}^{\infty} \left(1 + \frac{\gamma^k r_{t+k}}{\Gamma_{t-1+k}^+}\right) \end{aligned} \quad (111)$$

For the other two reward signals, we can express this equivalently as $(1 + R_t^\times) = (1 + g_t) \prod_{k=1}^t (1 + \gamma^k g_{t+k})$ or $(1 + R_t^\times) = \lambda_t \prod_{k=1}^t \gamma^k \lambda_{t+k}$. This represents the general compounding return that the agent seeks to maximise over time through learning an optimal policy π^* . A more practical way to characterise this is

$$1 + R_t^\times \equiv (1 + \Gamma_t^\times) \cdot \exp \left[\sum_{k=1}^{\infty} \ln \left| 1 + \frac{\gamma^k r_{t+k}}{\Gamma_{t-1+k}^+} \right| \right] \quad (112)$$

$$= (1 + \Gamma_t^\times) \cdot e^{\sum_{k=1}^{\infty} \psi_{t+k}} \quad (113)$$

which represents the objective as (exponentially) continuously compounding growth with the logarithmic return in each time period being ψ_t . The total future return is then the summation of all these compounding rates $\Psi_t = \sum_{k=1}^{\infty} \psi_{t+k}$. Furthermore, given the logarithm and exponential are monotonically increasing functions, the optimisation process can be greatly simplified as we can treat the agent decision at every time step individually. This way maximisation of Eq. (112) involves simply maximising each ψ_t at each time t separately for all time. In terms of optimal policy we explicitly have $\pi^* \in \arg \max_{\pi} \Psi_t$.

We can then directly write the proportionality

$$\begin{aligned} 1 + R_t^\times &\propto \left(1 + \frac{\gamma r_{t+1}}{\Gamma_t^+}\right) + \left(1 + \frac{\gamma^2 r_{t+2}}{\Gamma_{t+1}^+}\right) + \left(1 + \frac{\gamma^3 r_{t+3}}{\Gamma_{t+2}^+}\right) + \dots \\ &\propto \sum_{k=1}^{\infty} \left(1 + \frac{\gamma^k r(h_{t+k}, a_{t+k})}{\Gamma_{t-1+k}^+}\right) \end{aligned} \quad (114)$$

where we retain the addition of unity at each step to highlight the reward is a return relative to an existing value and therefore is bounded $\gamma^k r_{t+k} \geq (V_{\min} - \Gamma_{t-1+k}^+)$ with equality implying $(1 + R_t^\times) = 0$ ‘game over’ at time $t + k$. We have therefore reframed the desired Eq. (24) in the language of reinforcement learning.

Observe that by doing this we have left multiplicative dynamics and returned back to the realm of additive dynamics given we now maximise a summation [29]. However, we continue to refer to this situation as multiplicative due to three reasons: 1. The rewards at each time step are coupled to the cumulative prior rewards forming a relative return and so performance is not measured on absolute terms independently, rather it is designed to maximise rate of growth, 2. The state $\xi_t = s_t \cup \Gamma_{t-1}^+$ explicitly incorporates the existing cumulative value as a factor determining the reward, and 3. The objective is still ultimately a compounding rate of return that is proportional to a summation.

Then by using an iterative process at each time step we can numerically approximate using the Bellman equation in Eq. (28) all future returns at any time step to be

$$1 + \frac{Q_\pi(h_{t+k}, a_{t+k})}{\Gamma_{t-1+k}^+} \leftarrow 1 + \frac{\mathbb{E}_{a_{t+k+1} \sim \pi} [r(h_{t+k}, a_{t+k}) + \gamma Q_\pi(h_{t+k+1}, a_{t+k+1})]}{\Gamma_{t-1+k}^+} \quad (115)$$

where the value functions are not equivalent to the additive case, however since they are an artificial construction, we are free to define them as we choose. Notice that we can also write the right-hand side as

$$\frac{\Gamma_{t+k}^+ + \gamma V_\pi(h_{t+1+k})}{\Gamma_{t-1+k}^+} = \left(1 + \frac{r_{t+k}}{\Gamma_{t-1+k}^+}\right) + \frac{\gamma V_\pi(h_{t+1+k})}{\Gamma_{t-1+k}^+} = 1 + \frac{r_{t+k} + \gamma V_\pi(h_{t+1+k})}{\Gamma_{t-1+k}^+} \quad (116)$$

where $\Gamma_{t+k}^+ = \Gamma_{t-1+k}^+ + r_{t+k}$. This can be interpreted as the sum of both the one-holding period return and the forecasted discounted perpetuity ratio for all later periods. The ‘game over’ criterion is also built-in where if $\Gamma_t^+ \leq V_{\min}$ the episode is terminated. The pseudo-Bellman relations in Eqs. (77-81) can then be written as

$$1 + \frac{Q_\pi(h_t, a_t)}{\Gamma_{t-1}^+} = \sum_{o_{t+1}, r_t} P(o_{t+1}, r_t | h_t, a_t) \left[1 + \frac{r_t + \gamma V_\pi(h_{t+1})}{\Gamma_{t-1}^+} \right] \quad (117)$$

$$V_\pi(h_t) = Q_\pi(h_t, \pi(h_t)) \quad (118)$$

$$1 + \frac{Q^*(h_t, a_t)}{\Gamma_{t-1}^+} = \sum_{o_{t+1}, r_t} P(o_{t+1}, r_t | h_t, a_t) \left[1 + \frac{r_t + \gamma V^*(h_{t+1})}{\Gamma_{t-1}^+} \right] \quad (119)$$

$$V^*(h_t) = \max_{a_t \in \mathcal{A}} Q^*(h_t, a_t) \quad (120)$$

$$\pi^*(h_t) \in \arg \max_{a_t \in \mathcal{A}} Q^*(h_t, a_t) \quad (121)$$

To make this amenable to Q-learning we change notation not to confuse iterative updates of a Q-value estimate from $Q_{t+1}(h, a) \leftarrow Q_t(h, a)$ for any history-action pair (h, a) defined as as $(h_i, a_i) \leftarrow (h_t, a_t)$ with known $\Gamma_{i-1}^+ \leftarrow \Gamma_{t-1}^+$. Standard Q-learning in Eq. (109) is then modified to

$$\left(1 + \frac{Q_{t+1}(h, a)}{\Gamma_{i-1}^+}\right) \leftarrow (1 - \alpha_t(h, a)) \left(1 + \frac{Q_t(h, a)}{\Gamma_{i-1}^+}\right) + \alpha_t(h, a) \left(1 + \frac{1}{\Gamma_{i-1}^+} \left(r_t(h, a) + \gamma \max_{a'} Q_{t+1}(h', a)\right)\right) \quad (122)$$

which forms a general non-stationary HDP using a ϵ -greedy policy. This can also be crudely expressed as

$$\psi_{t+1} \leftarrow \psi_t + \frac{\alpha_t(h, a)}{\Gamma_{i-1}^+} \left(r_t + \gamma \max_{a'} Q_{t+1}(h', a) - Q_t(h, a)\right) \quad (123)$$

where the second term on the right is the improvement in the exponential compounding rate in one time step from a single learning update. To make this tractable we convert the standard results in [109] to multiplicative rewards and use their following assumptions:

1. The state-process is ergodic meaning that all states are reachable under any policy from the current state after sufficiently many steps. Strictly speaking this is false due to the ‘game over’ criterion, however if $\Gamma_{i-1}^+ + r_i \gg V_{\min}$ for any conceivable $r_i < 0$, the agent can afford to go any state for some $\delta t > 0$ where $\delta t \rightarrow 0$ as $\Gamma_{i-1+\delta t}^+ \rightarrow V_{\min}$. Therefore if we initialise the value $\Gamma_0^+ = V_0 \gg V_{\min}$ this assumption may be considered reasonable. Another way to state this is $|r_i| \ll V_0 \forall r_i < 0$, that is, the change in maximum absolute downwards change in valuation from a single step is small relative to the initial value. No such constraint is required $\forall r_i > 0$. Practically this can be enforced if each step represents a very small change in time and so there is a limit on the maximum change in cumulative additive value that can reasonably occur.
2. The rewards are bounded $r \in [r_{\min}, r_{\max}]$ which is standard to ensure stable convergence. In our case the lower bound also varies with time where $r_{i,\min} > V_{\min} - \Gamma_i^+$ otherwise the episode ends.

3. The state-process is a QDP where $Q^*(h, a) = Q^*(\tilde{h}, a)$ for some feature map $\phi(h) = \phi(\tilde{h}) = \xi$, and therefore $Q^*(h, a) = Q^*(\xi, a)$. We do not explicitly specify the feature map ϕ but assume one exists. This is a reasonable assumption since the scope of this work encapsulates only keeping track of the reward history while making no use of prior states and actions for future decisions. Importantly, this condition allows $Q(h_i, a_i) \neq Q(\xi_i, a_i)$ for all intermediate action-values and so is very flexible.

To prove the convergence of Eq. (122) by repeated updates we re-purpose standard methods [30, 84, 109, 143–145]. Without loss of generality, we can reparameterise the $Q_t(\xi, a) \leftarrow \left(1 + \frac{Q_t(\xi, a)}{\Gamma_{i-1}^+}\right)$ as it is a artificially constructed value that we seek to maximise. The update rule is then rewritten as

$$Q_{t+1}(\xi, a) = (1 - \alpha_t(\xi, a)) Q_t(\xi, a) + \alpha_t(\xi, a) (T_{h_t}^\pi Q)(\xi, a) \quad (124)$$

where we define $T_{h_t}^\pi$ to be the Bellman history-based operator for a non-stationary policy π that generally incorporates both the decision process history h_i and the history of all learning step sizes $\alpha_t(\xi, a)$ with the condition $\alpha_t(\xi_i, a_i) = 0 \forall (\xi, a) \neq (\xi_i, a_i)$. Application of this operator to a (ξ, a) -pair yields

$$(T_{h_t}^\pi Q)(\xi, a) = 1 + \frac{Q_{\pi, t+1}(\xi_i, a_i)}{\Gamma_{i-1}^+} = \sum_{\xi_{i+1}, r_i} p_{h_i}(\xi_{i+1}, r_i | s_i, a_i) \left[1 + \frac{r_i + \gamma V_\pi(\xi_{i+1})}{\Gamma_{i-1}^+} \right] \quad (125)$$

$$= \sum_{\xi_{i+1}} p_{h_i}(\xi_{i+1} | s_i, a_i) \mathbb{E} \left[1 + \frac{r_i + \gamma V_\pi(\xi_{i+1})}{\Gamma_{i-1}^+} \mid \xi_i, a_i \right] \quad (126)$$

$$= \mathbb{E}_{p_{h_i}} \left[1 + \frac{1}{\Gamma_{i-1}^+} \left(r_i + \gamma \max_{a'} Q_t(\xi_{i+1}, a') \right) \mid T_t^\pi \right] \quad (127)$$

where in the final line we assume a ϵ -greedy policy p_i and T_t^π is a complete history of the algorithm including h_i and all the steps $(\alpha_k)_{k \leq t}$. In Section 5.2 we prove under a strict criteria the convergence of Eq. (124) where $(T_{h_t}^\pi Q)(\xi, a) \propto Q_{t+1}(\xi, a) \rightarrow Q^*(\xi, a)$ with probability 1 (w.p.1.) as infinite updates $t \rightarrow \infty$ are applied.

5.2 Proof of Convergence and Uniqueness

The convergence and uniqueness of multiplicative Q-learning is demonstrated by first converting Eq. (124) to a standard form [145]. Consider a stochastic process $(\alpha_t(\xi, a), \Delta_t, F_t), t > 0$ for $t \in \mathbb{Z}^+$ where $\alpha_t(\xi, a), \Delta_t, F_t : s, a \rightarrow \mathbb{R}$. Let $T_{h_t}^\pi$ be a sequence of increasing σ -fields such that $\alpha_0(\xi, a)$ and Δ_0 are $T_{h_0}^\pi$ -measurable and $\alpha_t(\xi, a), \Delta_t$ and F_{t-1} are $T_{h_t}^\pi$ -measurable $\forall t$. Begin by defining

$$\Delta_t \equiv Q_{\pi, t+1}(\xi, a) - Q_t^*(\xi, a) \quad (128)$$

$$F_t \equiv (T_{h_t}^\pi Q)(\xi, a) - Q_t^*(\xi, a) \quad (129)$$

$$\Delta_{t+1} \equiv (1 - \alpha_t(\xi, a)) \Delta_t(\xi, a) + \alpha_t(\xi, a) F_t \quad (130)$$

A sequence $(Q_t(\xi, a))_{t \in \mathbb{N}}$ then generated by the iteration of Eq. (124) represented using Eqs. (128-130) converges to the optimal action-value under multiplicative dynamics $Q^*(\xi, a) = Q^*(h, a)$ of a QDP state-process w.p.1. if the following conditions are satisfied:

1. The agent state space $\xi \in \Xi : \mathcal{S} \times \mathcal{R}$ is finite.
2. The discount factor is bounded $\gamma \in [0, 1)$ for all steps.

3. Infinite learning updates are possible where $\Delta_{t+1} \rightarrow \Delta_\infty$ in the limit $t \rightarrow \infty$.

4. The Robbins-Monro (RM) conditions

$$\sum_{t=0}^{\infty} \alpha_t(\xi, a) = \infty, \text{ and } \sum_{t=0}^{\infty} \alpha_t^2(\xi, a) < \infty \quad (131)$$

for learning rates are satisfied which requires $\alpha_t(\xi, a) \in (0, 1]$ and $\alpha_t(\xi, a) = 0 \forall (\xi, a) \neq (\xi_t, a_t)$ [150, 151]. This also requires the state-process to be ergodic and the step size asymptotically decreases to converge to a fixed point though never ceases $\alpha_t(\xi, a) \neq 0$ learning in order to avoid local maxima.

5. There exists a $Q^*(\xi, a)$ such that

$$\| (T_{h_t}^\pi Q)(\xi, a) - Q^*(\xi, a) \|_\infty \leq \gamma \| Q_{\pi,t}(\xi, a) - Q^*(\xi, a) \|_\infty \quad \forall t \quad (132)$$

to prove that $Q^*(\xi, a)$ is a unique fixed point of the contraction T_h^π and converges to the optimal solution in the limit $t \rightarrow \infty$. This condition follows from the usual Banach's fixed-point theorem applied to MDPs.

6. The F_t term satisfies in expectation

$$\| \mathbb{E}_{p_{h_t}}[F_t | T_t^\pi] \|_\infty \leq \kappa \| \Delta_t \|_\infty + c_t \quad (133)$$

where $\kappa \in [0, 1]$ and $c_t \rightarrow 0$ w.p.1. as $t \rightarrow \infty$.

7. The noise is bounded if the conditional variance of F_t satisfies

$$\text{Var}(F_t | T_t^\pi) \leq \kappa (1 + \| \Delta_t \|_\infty)^2 \quad (134)$$

where κ is a constant.

The first condition is satisfied as it is the primary reason for converting a general HDP to a tractable QDP where \mathcal{S} is a small finite subset of $\mathcal{S} \in \mathcal{O}$. The third condition condition will be assumed valid in line with all literature derivations despite there existing clear episode termination criteria.

In practice RM conditions are violated as constant iterative step sizes $\alpha_t(\xi, a) = \alpha \forall t$ are usually used. This simplification works well as π_ϕ is non-stationary and that when using a mini-batch for learning, the policy parameters converge $\phi \rightarrow \phi^*$ as $t \rightarrow \infty$ with the variance of convergence proportional to α^2 [84].

To prove the fifth condition we must first show that $T_{h_t}^\pi$ is a max-norm contraction and that the fixed point equation $(T_{h_t}^\pi Q)(\xi, a) \propto Q_{t+1}^\pi(\xi, a)$ has a unique solution for L -Lipschitzian

$$\| (T_{h_t}^\pi Q)(\xi, a) - (T_{h_t}^\pi Q')(\xi, a) \| \leq L \| Q_{\pi,t}(\xi, a) - Q'_{\pi,t}(\xi, a) \| \quad (135)$$

where L is called the contraction error and $T_{h_t}^\pi$ is called a non-expansion if $L \leq 1$ or a contraction if $L < 1$. This is done by proving Banach's fixed point theorem for multiplicative dynamics where $(T_{h_t}^\pi Q)(\xi, a) \propto Q_{\pi,t+1}(\xi, a) \rightarrow Q^*(\xi, a)$ at a geometric rate as $t \rightarrow \infty$ expressed as

$$\| Q_{\pi,n}(\xi, a) - Q_\pi(\xi, a) \| \leq \gamma^n \| Q_{\pi,0}(\xi, a) - Q_\pi(\xi, a) \| \quad (136)$$

which then poses two additional questions: 1. Whether the $Q(\xi, a)$ pair is a fixed point of T_h^π in action-value space, and 2. Whether the Q-values $Q(\xi, a) = Q'(\xi, a)$ implying the (ξ, a) -pair can be uniquely represented.

For a fixed history h_t , the operator $T_{h_t}^\pi$ for a ϵ -greedy policy is shown to be a contraction mapping by

$$\begin{aligned}
& \| (T_{h_t}^\pi Q)(\xi, a) - (T_{h_t}^\pi Q')(\xi, a) \|_\infty \\
&= \max_{\xi, a} \left| \mathbb{E}_{p_{h_i}} \left[1 + \frac{1}{\Gamma_{i-1}^+} \left(r_i + \gamma \max_{a'} Q_t(\xi', a) \right) \middle| T_t^\pi \right] - \mathbb{E}_{p_{h_i}} \left[1 + \frac{1}{\Gamma_{i-1}^+} \left(r_i + \gamma \max_{a'} Q'_t(\xi', a) \right) \middle| T_t^\pi \right] \right| \\
&\stackrel{(a)}{=} \frac{\gamma}{\Gamma_{i-1}^+} \max_{\xi, a} \left| \mathbb{E}_{p_{h_i}} \left[\max_{a'} Q_t(\xi', a) \middle| s, a \right] - \mathbb{E}_{p_{h_i}} \left[\max_{a'} Q'_t(\xi', a) \middle| s, a \right] \right| \\
&\stackrel{(b)}{\leq} \frac{\gamma}{\Gamma_{i-1}^+} \max_{\xi, a} \max_{\xi'} \left| \max_{a'} Q_t(\xi', a) - \max_{a'} Q'_t(\xi', a) \right| \\
&\stackrel{(c)}{\leq} \frac{\gamma}{\Gamma_{i-1}^+} \max_{\xi, a} |Q_t(\xi, a) - Q'_t(\xi, a)| \\
&\stackrel{(d)}{=} \frac{\gamma}{\Gamma_{i-1}^+} \|Q_t(\xi, a) - Q'_t(\xi, a)\|_\infty
\end{aligned} \tag{137}$$

in which (a) for a fixed history h_i has the same $\Gamma_{i-1}^+ \geq \Gamma_0^+ > V_{\min}$ and the QDP assumption assures the expectation depends only on the (ξ, a) -pair, (b) establishes a upper bound by removing the expectation, (c) increases the upper bound by no longer demanding the smallest optimal maximal difference, and (d) is the usual max-norm contraction $\|x\|_\infty = \sup_{x \in \chi} (f(x)) = \max_{x \in \chi} |x|$ by definition.

For Eq. (137) to be valid, clearly we must have $\Gamma_{i-1}^+ > V_{\min} \geq \gamma \forall i$ to ensure $\gamma < \Gamma_{i-1}^+$ which enforces the bound for multiplicative dynamics. This is clearly an artificial construction on the dimensionless reward scheme we must use for convergence. The impact of this is subtlety will be highly dependent on the domain. One interpretation for this requirement is that the reward scheme could be defined in units of the discount factor where the zero point is γ . Regardless, this allows us to say

$$\| (T_{h_t}^\pi Q)(\xi, a) - (T_{h_t}^\pi Q')(\xi, a) \|_\infty \leq \gamma \|Q_{\pi, t}(\xi, a) - Q'_{\pi, t}(\xi, a)\|_\infty \tag{138}$$

which is now in the desired form in Eq. (135). To prove that this represented unique fixed points we follow the discussion in [84] where we can show for subsequent update steps $k > 0$ that

$$\begin{aligned}
\|Q_{t+k}^\pi(\xi, a) - Q_t^\pi(\xi, a)\| &= \| (T_{h_{t-1+k}}^\pi Q)(\xi, a) - (T_{h_{t-1+k}}^\pi Q)(\xi, a) \| \\
&\leq \gamma \|Q_{\pi, t-1+k}(\xi, a) - Q_{\pi, t-1}(\xi, a)\|_\infty = \gamma \| (T_{h_{t-2+k}}^\pi Q)(\xi, a) - (T_{h_{t-2+k}}^\pi Q)(\xi, a) \| \\
&\quad \vdots \\
&\leq \gamma^k \|Q_{\pi, 0}(\xi, a) - Q_{\pi, 0}(\xi, a)\|_\infty
\end{aligned} \tag{139}$$

as shown in Eq. (136). Next use the triangle inequality and recall $\sum_{\lambda=0}^{\infty} \gamma^\lambda = (1 - \gamma)^{-1}$ to get

$$\begin{aligned}
\|Q_{\pi, k}(\xi, a) - Q_{\pi, 0}(\xi, a)\| &\leq \sum_{j=1}^k \|Q_{\pi, j}(\xi, a) - Q_{\pi, j-1}(\xi, a)\|_\infty \\
&\leq \sum_{j=1}^k \gamma^{j-1} \|Q_{\pi, 1}(\xi, a) - Q_{\pi, 0}(\xi, a)\|_\infty
\end{aligned}$$

$$\leq \frac{1}{1-\gamma} \|Q_{\pi,1}(\xi, a) - Q_{\pi,0}(\xi, a)\|_\infty \quad (140)$$

Therefore

$$\|Q_{\pi,t+k}(\xi, a) - Q_{\pi,t}(\xi, a)\| \leq \frac{\gamma^t}{1-\gamma} \|Q_{\pi,1}(\xi, a) - Q_{\pi,0}(\xi, a)\|_\infty \quad (141)$$

is of the form of a standard Cauchy sequence in the limit of infinite learning updates $t \rightarrow \infty$. Due to exponential suppression by $\gamma \in [0, 1)$ this leads to for $k > 0$ the well-known convergence

$$\lim_{t \rightarrow \infty} \|Q_{\pi,t+k}(\xi, a) - Q_{\pi,t}(\xi, a)\| = 0 \quad (142)$$

allowing us to state $Q_{\pi,t+k}(\xi, a) \rightarrow Q_\pi(\xi, a) \forall k > 0$ that is the optimal value. Next recall the definition $(T_{h_t}^\pi Q)(\xi_i, a_i) \propto Q_{\pi,t+1}(\xi_i, a_i)$ and take the limit on both sides

$$\lim_{t \rightarrow \infty} \|(T_{h_t}^\pi Q)(\xi_i, a_i)\| = \lim_{t \rightarrow \infty} \left\| 1 + \frac{Q_{\pi,t+1}(\xi_i, a_i)}{\Gamma_{i-1}} \right\| = \lim_{t \rightarrow \infty} \left\| 1 + \frac{Q_{\pi,t}(\xi_i, a_i)}{\Gamma_{i-1}} \right\| = 1 + \frac{Q_\pi(\xi_i, a_i)}{\Gamma_{i-1}} \quad (143)$$

to see that $Q_\pi(\xi, a)$ must be a fixed point of the L -Lipschitzian continuous contraction $T_{h_t}^\pi$. Finally, regarding the uniqueness of Q-values for any (ξ_i, a_i) -pair we can solve the equation that

$$\begin{aligned} \|(T_h^\pi Q)(\xi_i, a_i) - (T_h^\pi Q')(\xi_i, a_i)\| &= \frac{1}{\Gamma_{i-1}} \|Q_\pi(\xi_i, a_i) - Q'_\pi(\xi_i, a_i)\| \\ &\leq \frac{\gamma}{\Gamma_{i-1}} \|Q_\pi(\xi_i, a_i) - Q'_\pi(\xi_i, a_i)\| \\ (1-\gamma) \|Q_\pi(\xi_i, a_i) - Q'_\pi(\xi_i, a_i)\| &\leq 0 \end{aligned} \quad (144)$$

so either $\|Q_\pi(\xi, a) - Q'_\pi(\xi, a)\| = 0$ or $Q_\pi(\xi, a) = Q'_\pi(\xi, a)$ (unique). To show that it is the latter case we once again construct the bound

$$\begin{aligned} \|Q_{\pi,t}(\xi, a) - Q'_\pi(\xi, a)\| &= \|(T_{h_{t-1}}^\pi Q)(\xi, a) - (T_{h_{t-1}}^\pi Q')(\xi, a)\| \\ &\leq \gamma \|Q_{\pi,t-1}(\xi, a) - Q'_\pi(\xi, a)\|_\infty = \gamma \|(T_{h_{t-2}}^\pi Q)(\xi, a) - (T_{h_{t-2}}^\pi Q')(\xi, a)\| \\ &\leq \gamma^t \|Q_{\pi,0}(\xi, a) - Q'_\pi(\xi, a)\|_\infty \\ &\leq \frac{\gamma^t}{1-\gamma} \|Q_{\pi,0}(\xi, a) - Q'_\pi(\xi, a)\|_\infty \end{aligned} \quad (145)$$

which forms a Cauchy sequence and by the same logic as earlier we can prove that as $t \rightarrow \infty$ we must have $Q_{\pi,t}(\xi, a) \rightarrow Q_\pi(\xi, a) \forall t$ and therefore $Q_\pi(\xi, a) = Q'_\pi(\xi, a)$ is a unique fixed point.

Condition six requires that as the agent samples from the underlying HDP, the expected difference between the application of the Bellman history-based operator $(T_{h_t}^\pi Q)(\xi, a)$ and the optimal Q-value $Q^*(\xi, a)$ is finite and bounded as we apply infinite updates. Application of Eq. (138) directly into Eq. (133) gives

$$\begin{aligned} \|\mathbb{E}_{p_{h_t}}[F_t | T_{h_t}^\pi]\|_\infty &= \|\mathbb{E}_{p_{h_t}}[(T_{h_t}^\pi Q)(\xi, a) - Q_t^*(\xi, a) | T_t^\pi]\|_\infty \\ &\stackrel{(a)}{\leq} \|\mathbb{E}_{p_{h_t}}[Q_{\pi,t+1}(\xi, a) - Q_t^*(\xi, a) | \xi, a]\|_\infty \\ &\stackrel{(b)}{\leq} \gamma \|Q_{\pi,t}(\xi, a) - Q_t^*(\xi, a)\|_\infty \\ &= \gamma \|\Delta_t\|_\infty \end{aligned} \quad (146)$$

as for (a) recall the reparameterisation $Q_t(\xi_i, a_i) \leftarrow \left(1 + \frac{Q_t(\xi_i, a_i)}{\Gamma_{t-1}^+}\right)$ and for (b) $\mathbb{E}_{p_{h_i}}[Q_{\pi, t+1}(\xi, a)|\xi_i, a_i] = Q_{\pi, t+1}(\xi, a)$ since the Q-value is already an expectation as seen in Eq. (126). Therefore condition six is satisfied since $\gamma \in [0, 1]$ as one would expect by our definition of F_t .

The final condition ensures the conditional variance between difference between $(T_{h_t}^\pi Q)(\xi, a)$ and $Q^*(\xi, a)$ is also finite and bounded as we apply infinite updates. This is seen with

$$\begin{aligned} \text{Var}(F_t|T_t^\pi) &= \mathbb{E}_{p_{h_i}}[F_t^2|T_t^\pi] - \mathbb{E}_{p_{h_i}}[F_t|T_t^\pi]^2 \\ &\stackrel{(a)}{\leq} \frac{1}{4} (\max(F_t) - \min(F_t))^2 \\ &\stackrel{(b)}{=} \frac{1}{4} \max_{s,a} |(T_{h_t}^\pi Q)(\xi, a) - Q_t^*(\xi, a)|^2 \\ &\stackrel{(c)}{\leq} \frac{\gamma}{4} \|Q_{\pi, t}(\xi, a) - Q_t^*(\xi, a)\|_\infty^2 \\ &\leq \frac{\gamma}{4} \|\Delta_t\|_\infty^2 \end{aligned} \tag{147}$$

where (a) utilises Popoviciu's inequality [152, 153] on variances $\text{Var}(X) \leq \frac{1}{4} (\max(X) - \min(X))^2$, (b) uses the fact the $\min(F_t) = 0$ when $(T_{h_t}^\pi Q)(\xi, a) \propto Q^*(\xi, a)$ and (c) is again the reparameterisation. An alternative to construct the strict upper bound in (a) may involve using the far more general Bhatia-Davis inequality [154] $\text{Var}(X) \leq (\max(X) - \text{Mean}(X))(\text{Mean}(X) - \min(X))$. Hence condition seven is also satisfied as $\gamma \in [0, 1]$.

Therefore we have shown that for multiplicative dynamics if $\Gamma_{t-1}^+ > V_{\min} \geq \gamma \forall t$ then $Q_t(\xi, a) \rightarrow Q^*(\xi, a)$ w.p.1. as the number of updates increases indefinitely. This convergence is valid if we assume the state-process is ergodic and forms at the very least a MDP. However, unlike the standard stationary MDP proof, our formulation uses the contraction operator $T_{h_t}^\pi$ which is dependent on history and allows to proof to scale to state-processes that are QDPs which can be non-stationary. Formally, if $Q^*(h, a) = Q^*(\tilde{h}, a)$ for some feature map $\phi(h) = \phi(\tilde{h}) = \xi$, then $Q^*(h, a) = Q^*(\xi, a)$. This concludes the proof for the convergence of Q-learning using multiplicative dynamics that valid for a large class of decision problems including MDPs.

5.3 Clipped Double Q-Learning

The utility of using two Q-values was first in the context of Double-Q learning to enhance learning stability in the tabular regime [64]. This approach combined with deep neural networks to more expressively represent action-values was successfully used to obtain seminal results in classic Atari video games [33–39], the board games of Chess, Go and Shogi [43–46]. In actor-critic methodologies discussed in Section 3.3 this approach was combined with actors to generate Q-value estimates [62, 66, 67]. A constant trend seen is the overestimation bias in the Q-values due to accumulating propagation of errors as learning continues. To combat this [71] introduced clipped double Q-learning that uses the minimum of two Q-values for target values which works extremely well compared to all prior formulations as seen across locomotive continuous control tasks [71, 75, 77, 78, 155].

With some modification of the proofs in Section 5.2 and using the methods of [71] we can show the convergence of clipped double Q-learning under multiplicative dynamics assuming the state-process is once again a QDP. For the twin Q-values $Q_{\pi, t}^A$ and $Q_{\pi, t}^B$ where the contraction $T_{h_t}^\pi$ operates uniquely on the minimum of two. The optimal action a^* for both Q-values at the same state s_i is defined by using only the first Q-value where for a ϵ -greedy policy

$a^* = \arg \max_a Q_{\pi,t}^A$. Therefore, using

$$(T_{h_t}^\pi Q)(\xi, a) = 1 + \frac{1}{\Gamma_{i-1}^+} \min \left(Q_{\pi,t+1}^A(\xi, a^*), Q_{\pi,t+1}^B(\xi, a^*) \right) \quad (148)$$

and the same reparameterisation $Q_t(\xi_i, a_i) \leftarrow \left(1 + \frac{Q_t(\xi_i, a_i)}{\Gamma_{i-1}^+} \right)$. We then modify Eqs. (128-130) to yield

$$\Delta_t \equiv Q_{\pi,t+1}^A(\xi, a) - Q_t^*(\xi, a) \quad (149)$$

$$\begin{aligned} F_t &\equiv (T_{h_t}^\pi Q)(\xi, a) - Q_t^*(\xi, a) \\ &= (T_{h_t}^\pi Q)(\xi, a) - Q_t^*(\xi, a) + (T_{h_t}^\pi Q^A)(\xi, a) - (T_{h_t}^\pi Q^A)(\xi, a) \\ &= F_t^Q + c_t \end{aligned} \quad (150)$$

$$\Delta_{t+1} \equiv (1 - \alpha_t(\xi, a)) \Delta_t(\xi, a) + \alpha_t(\xi, a) F_t \quad (151)$$

where $F_t^Q = (T_{h_t}^\pi Q^A)(\xi, a) - Q_t^*(\xi, a)$ is the standard Q-learning in Eq. (129) which we have proven to converge in Eqs. (146-147). The second term $c_t = (T_{h_t}^\pi Q)(\xi, a) - (T_{h_t}^\pi Q^A)(\xi, a)$ is purposely designed to resemble Eq. (133) and so $c_t \rightarrow 0$ w.p.1. which is expected given we expect $Q_{\pi,t+1}^A(\xi, a^*)$ and $Q_{\pi,t+1}^B(\xi, a^*)$ to eventually converge to be equal as $t \rightarrow \infty$. Formally we show this by defining $\Delta_t^{AB} \equiv Q_{\pi,t}^A(\xi, a) - Q_{\pi,t}^B(\xi, a)$ so that

$$\begin{aligned} \Delta_{t+1}^{AB} &= \Delta_t^{AB} + \alpha_t(\xi, a) \left((T_{h_t}^\pi Q)(\xi, a) - Q_{\pi,t}^A(\xi, a) - (T_{h_t}^\pi Q)(\xi, a) - Q_{\pi,t}^B(\xi, a) \right) \\ &= \Delta_t^{AB} + \alpha_t(\xi, a) \left(Q_{\pi,t}^B(\xi, a) - Q_{\pi,t}^A(\xi, a) \right) \\ &= (1 - \alpha_t(\xi, a)) \Delta_t^{AB} \end{aligned} \quad (152)$$

which clearly proves $c_t \rightarrow 0$ convergence as the learning rate decreases. All other aspects of Q-learning convergence remain unchanged and so we can state that $Q_t^A(\xi, a) \rightarrow Q^*(\xi, a)$ w.p.1. under the same criteria as before. Noting that in our derivation at no point did we explicitly assume additive or multiplicative dynamics and so can be used for either. Hence we have extended the applicability of clipped double Q-learning for both dynamics and confirmed that it is functional for history-dependent non-stationary QDP domains.

5.4 Multi-Step Targets

Multi-step targets discussed in Section 3.7 can be easily extended to both dynamics. Firstly the additive dynamics case adjusts Eq. (76) to include the cumulative additive episodic sum so that the target Q-value becomes

$$Q_{\bar{\theta}}^{(m)}(s_t, a_t) \equiv \Gamma_{t-1}^+ + R_t^{(m)} + \gamma^m Q_{\bar{\theta}}(s_{t+m}, a') \quad (153)$$

where $R_t^{(m)} = \sum_{k=0}^{m-1} \gamma^k r_{t+k}$. For multiplicative dynamics there are two ways to achieve this. The first involves calculating the compounding return for $m > 1$ as

$$\Gamma_{t+m-2}^x \equiv \prod_{k=0}^{m-2} \frac{\Gamma_{t-1+k}^+ + \gamma^k r_{t+k}}{\Gamma_{t-1+k}^+} = \exp \left[\sum_{k=0}^{m-2} \ln \left| 1 + \frac{\gamma^k r_{t+k}}{\Gamma_{t-1+k}^+} \right| \right] \quad (154)$$

where at each m -step bootstrapping, $\Gamma_{t-1+k}^+ + \gamma^k r_{t+k} > V_{\min} \geq \gamma$ must be explicitly confirmed. This then allows us to write the target value as

$$1 + \frac{Q_{\bar{\theta}}^{(m)}(\xi_t, a_t)}{\Gamma_{t+m-2}^+} \equiv \Gamma_{t+m-2}^\times \cdot \frac{\Gamma_{t+m-2}^+ + \gamma^{m-1} (r_{t+m-1} + \gamma Q_{\bar{\theta}}(\xi_{t+m}, a'))}{\Gamma_{t+m-2}^+} \quad (155)$$

as the product of compounding growth rates at each step with the estimated Q-value being the heavily discounted future value. Calculating this at each learning step when randomly sampling from an experience replay buffer \mathcal{D} for a relatively large mini-batch size will be computationally expensive as there are two products to calculate while confirming each episode is not terminated at each bootstrapping step due to the ‘game over’ condition.

Alternatively, we can use a simplification that is far more computationally efficient but strictly speaking is not a multiplicative process. Instead of products, using Eq. (116) as a guide we can construct

$$1 + \frac{Q_{\bar{\theta}}^{(m)}(\xi_t, a_t)}{\Gamma_{t+m-2}^+} \approx \left(1 + \frac{R_t^{(m)}}{\Gamma_{t-1}^+}\right) + \frac{\gamma^m Q_{\bar{\theta}}(\xi_{t+m}, a')}{\Gamma_{t-1}^+} = 1 + \frac{R_t^{(m)} + \gamma^m Q_{\bar{\theta}}(\xi_{t+m}, a')}{\Gamma_{t-1}^+} \quad (156)$$

as the sum of the m -holding period return and the forecasted discounted perpetuity ratio for all later periods. It is unlikely that these two methods would produce identical learning outcomes since intermediate rewards in $R_t^{(m)}$ could trigger the ‘game over’ criterion. An exception to this is obviously the usual $m = 1$ case

$$1 + \frac{Q_{\bar{\theta}}^{(1)}(\xi_t, a_t)}{\Gamma_{t-1}^+} \equiv 1 + \frac{Q_{\bar{\theta}}(\xi_t, a_t)}{\Gamma_{t-1}^+} \leftarrow 1 + \frac{r_t + \gamma Q_{\bar{\theta}}(\xi_{t+1}, a')}{\Gamma_{t-1}^+} \quad (157)$$

with no additional bootstrapping that can be considered the multiplicative analogue of Eq. (28). For actor-critic methods, the critic optimisation in Eq. (30) can then be written

$$\begin{aligned} \frac{Q_{\bar{\theta}}^{(m)}(\xi_t, a_t)}{\Gamma_{t+m-2}^+} - \frac{Q_{\theta}(\xi_t, a_t)}{\Gamma_{t-1}^+} &\approx \left(1 + \frac{R_t^{(m)} + \gamma^m Q_{\bar{\theta}}(\xi_{t+m}, a')}{\Gamma_{t-1}^+}\right) - \left(1 + \frac{Q_{\theta}(\xi_t, a_t)}{\Gamma_{t-1}^+}\right) \\ &= \frac{1}{\Gamma_{t-1}^+} \left(R_t^{(m)} + \gamma^m Q_{\bar{\theta}}(\xi_{t+m}, a') - Q_{\theta}(\xi_t, a_t)\right) \end{aligned} \quad (158)$$

as the difference to be minimised. This is the key point, for off-policy learning with a mini-batch using an MSE loss function we have the objectives

$$J(\theta^+) = \mathbb{E}_{U(\mathcal{D}^+)} \left[\left(R_t^{(m)} + \gamma^m Q_{\bar{\theta}}(s_{t+m}, a') - Q_{\theta}(s_t, a_t) \right)^2 \right] \quad (159)$$

$$J(\theta^\times) \approx \mathbb{E}_{U(\mathcal{D}^\times)} \left[\left(\frac{1}{\Gamma_{t-1}^+} \left(R_t^{(m)} + \gamma^m Q_{\bar{\theta}}(\xi_{t+m}, a') - Q_{\theta}(\xi_t, a_t) \right) \right)^2 \right] \quad (160)$$

with uniform sampling from the experience replay buffers returning the tuples $(s_t, a_t, R_t^{(m)}, s_{t+m}) \sim U(\mathcal{D}^+)$ and $(\Gamma_{t-1}^+, s_t, a_t, R_t^{(m)}, s_{t+m}) \sim U(\mathcal{D}^\times)$ respectively, with multiplicative case restricting sampling to include only cases where $\Gamma_{t-1}^+ > V_{\min} \geq \gamma$ and so only learns from ‘living’ states $\mathcal{D}^\times \subseteq \mathcal{D}^+$, or in terms of cardinality $|\mathcal{D}^\times| \leq |\mathcal{D}^+|$. This is crucial for many environments where the agent is not permitted to go ‘temporarily bankrupt’ in order to maximise returns, it must strictly find the optimal policy under the constraint of staying ‘alive’ at all times. This type of behaviour is required for all situations where there may exist alternative paths to the same destination. It should be noted that this approach will then by definition be more sample inefficient as it does not use the complete buffer, this will be especially important during the early stages when the agent is repeatedly failing.

Each sample within the mini-batch will have a unique Γ^+ and so the difference between the two aggregation schemes is non-trivial. The largest contributions for the additive case will be when the absolute difference is large, while for multiplicative, it will be for those samples that have the largest difference in returns. There is no reason to assume that the largest contributors will be shared between the dynamics.

6 Policy Gradients with Multiplicative Dynamics

Policy gradients introduced in Section 3.2 come in two varieties, stochastic $\pi_\phi(a|s)$ and deterministic $\mu_\phi(s)$ action sampling when the agent is provided a state. These are essential if we desire to use actors capable of performing more complex manoeuvres than the ϵ -greedy approach. In both cases, the action-values are reparameterised $Q_{\pi_\phi}(s, a) \rightarrow Q_\theta(s, a)$ to enhance stability by reducing coupling.

For additive dynamics, as with Q-learning, all existing literature results holds as we simply scale the rewards $r_t \leftarrow r_t + \Gamma_{t-1}^+$. In this section we prove that for multiplicative dynamics policy gradient theorems also remain relatively unchanged as there are no temporal differences between Q-values at any point. The results will again resemble Eqs. (159-160) with the functional forms of both dynamics looking similar but experimental results will likely be very different.

6.1 Stochastic Actors

For stochastic policies the objective to maximise $J(\pi_\phi) = \mathbb{E} [1 + R_t^\times | \pi_\phi]$ through gradient ascent using $\nabla_\phi J(\pi_\phi) = \mathbb{E} [\nabla_\phi \pi_\phi(a|s) (1 + R_t^\times)] = \mathbb{E} [\nabla_\phi \ln \pi_\phi(a|s) (1 + R_t^\times) | \pi_\phi]$ or in other words $\phi^* = \arg \max_\phi \mathbb{E} [(1 + R_t^\times) | \pi_\phi]$ where R_t^\times is defined in Eqs. (110-112). Using the usual (improper) discounted state visitation distribution for policy π_ϕ as ρ^{π_ϕ} which can be interpreted as representing the marginals for the trajectory distribution, the policy gradient becomes

$$\begin{aligned} \nabla_\phi J(\pi_\phi^\times) &= \mathbb{E}_{\xi_t \sim \rho^{\pi_\phi}, a_t \sim \pi_\phi} \left[\nabla_\phi \pi_\phi(a_t | \xi_t) \left(1 + \frac{Q_{\pi_\phi}(\xi_t, a_t)}{\Gamma_{t-1}^+} \right) \right] \\ &= \mathbb{E}_{\xi_t \sim \rho^{\pi_\phi}} \left[\nabla_\phi \ln \pi_\phi(a_t | \xi_t) \left(1 + \frac{Q_{\pi_\phi}(\xi_t, a_t)}{\Gamma_{t-1}^+} \right) \middle| \pi_\phi \right] \end{aligned} \quad (161)$$

where the future return is due to the Bellman equation in Eq. (115) being of the correct monotonic structure to maximise R_t^\times , but not an exact representation of the compounding growth rate.

To prove that this theorem works in multiplicative dynamics we must explicitly show that gradient ascent works $\nabla_\phi J(\pi_\phi^+) \propto \nabla_\phi \pi_\phi(a_t | s_t) Q_{\pi_\phi}(s_t, a_t)$ as seen in Eq. (25). This is done by using the Bellman equation and unrolling the Q-values following the original derivation [29, 65]. We first remove the expectation value and operate in a discrete tabular state-action space. We also denote the additive cumulative reward $\Gamma_{t-1}^+(s_{t-1}, a_{t-2} | h_{t-1}) = \Gamma_{t-1}^+$ as a constant at each step for stochastic sampling $a \in \mathcal{A}$ since it is solely based on a known history. Observe then

$$\begin{aligned} \nabla_\phi J(\pi_\phi^\times) &= \mathbb{E}_{a_t \sim \pi_\phi} \left[\nabla_\phi \pi_\phi(a | \xi_t) \left(1 + \frac{Q_{\pi_\phi}(\xi_t, a_t)}{\Gamma_{t-1}^+} \right) \right] \\ &= \sum_a \left[\nabla_\phi \pi_\phi(a | \xi_t) \left(1 + \frac{Q_{\pi_\phi}(\xi_t, a_t)}{\Gamma_{t-1}^+} \right) + \pi_\phi(a | \xi_t) \nabla_\phi \left(1 + \frac{Q_{\pi_\phi}(\xi_t, a_t)}{\Gamma_{t-1}^+} \right) \right] \\ &\stackrel{(a)}{=} \sum_a \left[\nabla_\phi \pi_\phi(a | \xi_t) \left(1 + \frac{Q_{\pi_\phi}(\xi_t, a_t)}{\Gamma_{t-1}^+} \right) \right. \\ &\quad \left. + \pi_\phi(a | \xi_t) \nabla_\phi \sum_{\xi_{t+1}} \left[p_{h_t}(\xi_{t+1} | \xi_t, a_t) \left(1 + \frac{r(\xi_t, a') + \gamma Q_{\pi_\phi}(\xi_{t+1}, a')}{\Gamma_{t-1}^+} \right) \right] \right] \\ &\stackrel{(b)}{=} \sum_a \left[\nabla_\phi \pi_\phi(a | \xi_t) \left(1 + \frac{Q_{\pi_\phi}(\xi_t, a_t)}{\Gamma_{t-1}^+} \right) \right] \end{aligned}$$

$$\begin{aligned}
& + \pi_\phi(a|\xi_t) \nabla_\phi \sum_{\xi_{t+1}} \left[p_{h_t}(\xi_{t+1}|\xi_t, a_t) \left(1 + \frac{\gamma Q_{\pi_\phi}(\xi_{t+1}, a')}{\Gamma_{t-1}^+} \right) \right] \\
& \stackrel{(c)}{=} \sum_a \left[\nabla_\phi \pi_\phi(a|\xi_t) \left(1 + \frac{Q_{\pi_\phi}(\xi_t, a_t)}{\Gamma_{t-1}^+} \right) + \pi_\phi(a|\xi_t) \sum_{\xi_{t+1}} [p_{h_t}(\xi_{t+1}|\xi_t, a_t) \right. \\
& \quad \times \sum_{a'} \left[\nabla_\phi \pi_\phi(a'|\xi_{t+1}) \left(1 + \frac{\gamma Q_{\pi_\phi}(\xi_{t+1}, a')}{\Gamma_{t-1}^+} \right) \right. \\
& \quad \left. \left. + \pi_\phi(a'|\xi_{t+1}) \nabla_\phi \sum_{\xi_{t+2}} \left[p_{h_{t+1}}(\xi_{t+2}|\xi_{t+1}, a_{t+1}) \left(1 + \frac{r(\xi_{t+1}, a'') + \gamma^2 Q_{\pi_\phi}(\xi_{t+2}, a')}{\Gamma_{t-1}^+} \right) \right] \right] \right] \\
& = \sum_a \left[\nabla_\phi \pi_\phi(a|\xi_t) \left(1 + \frac{Q_{\pi_\phi}(\xi_t, a_t)}{\Gamma_{t-1}^+} \right) + \pi_\phi(a|\xi_t) \sum_{\xi_{t+1}} [p_{h_t}(\xi_{t+1}|\xi_t, a_t) \right. \\
& \quad \times \sum_{a'} \left[\nabla_\phi \pi_\phi(a'|\xi_{t+1}) \left(1 + \frac{\gamma Q_{\pi_\phi}(\xi_{t+1}, a')}{\Gamma_{t-1}^+} \right) \right. \\
& \quad \left. \left. + \pi_\phi(a'|\xi_{t+1}) \nabla_\phi \sum_{\xi_{t+2}} \left[p_{h_{t+1}}(\xi_{t+2}|\xi_{t+1}, a_{t+1}) \left(1 + \frac{\gamma^2 Q_{\pi_\phi}(\xi_{t+2}, a')}{\Gamma_{t-1}^+} \right) \right] \right] \right] \tag{162}
\end{aligned}$$

which clearly leads to infinite recursion. In (a) we use the Bellman equation seen in Eq. (126), (b) we assume $\nabla_\phi(r(\xi_t, a')/\Gamma_{t-1}^+) = 0$ as action sampling for constant rewards is random which is not exactly true since policy learning is a non-stationary endeavour, and (c) we apply the Bellman equation again further revealing the recursion.

As usual we define $p(\xi_t \rightarrow \xi', k, \pi_\phi)$ to be the probability of $\xi_t \rightarrow \xi'$ in k time steps under policy π_ϕ . We can summarise this recursion by repeated application of the Bellman equation with

$$\nabla_\phi J(\pi_\phi^\times) = \sum_{\xi_{t+1}} \left(\sum_{t=0}^{\infty} p(\xi_t \rightarrow \xi_{t+1}, t, \pi_\phi) \right) \sum_{a_t} \nabla_\phi \pi_\phi(a|\xi_{t+1}) \left(1 + \frac{\gamma^t Q_{\pi_\phi}(\xi_{t+1}, a)}{\Gamma_{t-1}^+} \right) \tag{163}$$

where the policy quantifies Eq. (115). When performing gradient ascent, the impact of constant values does not effect the optimisation process. Hence if modify the unrolled returns such that

$$\gamma^t \left(1 + \frac{Q_{\pi_\phi}(\xi_{t+1}, a)}{\Gamma_{t-1}^+} \right) \leftarrow \left(1 + \frac{\gamma^t Q_{\pi_\phi}(\xi_{t+1}, a)}{\Gamma_{t-1}^+} \right) \tag{164}$$

there will be no change on the final results, only the speed of convergence especially if there are no local maximums. The γ^t term geometrically decreases for subsequent bootstrapping and so poses no issues. It would be incorrect to exactly factor out γ^t as inclusion of a γ^{-t} would lead to divergence given $\gamma \in [0, 1)$ when $t \rightarrow \infty$. For practical purposes, policy learning occurs at a fixed step and so this will still ensure the weights ϕ are updated in the correct direction regardless.

This allows us to utilise existing machinery and construct the usual density

$$\rho^{\pi_\phi}(\xi') \equiv \int_S d\xi \sum_{t=1}^{\infty} \gamma^{t-1} p_1(\xi) p(\xi \rightarrow \xi', t, \pi_\phi) \tag{165}$$

so that

$$\nabla_\phi J(\pi_\phi^\times) \approx \nabla_\phi \int_S \rho^{\pi_\phi}(\xi) \int_A da d\xi \pi_\phi(a|\xi_t) \left(1 + \frac{Q_{\pi_\phi}(\xi_t, a_t)}{\Gamma_{t-1}^+} \right)$$

$$\begin{aligned}
&= \int_{\mathcal{S}} \rho^{\pi_\phi}(\xi) \int_{\mathcal{A}} da d\xi \nabla_\phi \pi_\phi(a_t | \xi_t) \left(1 + \frac{Q_{\pi_\phi}(\xi_t, a_t)}{\Gamma_{t-1}^+} \right) \\
&= \mathbb{E}_{\xi_t \sim \rho^{\pi_\phi}, a_t \sim \pi_\phi} \left[\nabla_\phi \ln \pi_\phi(a_t | \xi_t) \left(1 + \frac{Q_{\pi_\phi}(\xi_t, a_t)}{\Gamma_{t-1}^+} \right) \right]
\end{aligned} \tag{166}$$

where the approximation is due to Eq. (164). This is of the required form [29, 65, 66] and note we have purposely retained the $(1 + R^\times)$ structure to highlight that the policy is maximising the future discounted (exponentially) continuous compounding return in Eq. (112).

Effectively this is saying that as long as the underlying parameters policy ϕ^\times are updated in the direction of $\nabla_\phi \ln \pi_\phi(a_t | \xi_t) \left(1 + \frac{Q_{\pi_\phi}(\xi_t, a_t)}{\Gamma_{t-1}^+} \right)$, the optimal policy is approached. The difference between this and the standard $\nabla_\phi J(\pi_\phi^+) \propto \nabla_\phi \pi_\phi(a_t | s_t) Q_{\pi_\phi}(s_t, a_t)$ result is simply here we have represented the value as return over the existing cumulative episodic return. Therefore objectives to be maximised for off-policy stochastic actors are then

$$J(\phi^+) = \mathbb{E}_{U(\mathcal{D}^+)} [\ln \pi_\phi(a_t | s_t) Q_{\pi_\phi}(s_t, a_t)] \tag{167}$$

$$J(\phi^\times) \approx \mathbb{E}_{U(\mathcal{D}^\times)} \left[\ln \pi_\phi(a_t | \xi_t) \left(1 + \frac{Q_{\pi_\phi}(\xi_t, a_t)}{\Gamma_{t-1}^+} \right) \right] \tag{168}$$

where the discussion in Section 5.4 on the differences between both dynamics is applicable. Overall, we see that multiplicative dynamics maximises the future return relative to the existing value while additive dynamics simply maximises the future value. The former then can be interpreted as specifically prioritising policies that avoid steep losses at all time steps even though they may not yield the largest valuations.

6.2 Deterministic Actors

Deterministic polices have the advantage of being more computationally efficient as they remove random policy distribution sampling at each gradient step. The objective to maximise is $J(\mu_\phi) = \mathbb{E} [1 + R_t^\times | \mu_\phi]$ through gradient ascent using $\nabla_\phi J(\mu_\phi) = \mathbb{E} [\nabla_\phi \mu_\phi (1 + R_t^\times)]$ or equivalently $\phi^* = \arg \max_\phi \mathbb{E} [(1 + R_t^\times) | \mu_\phi]$. Similarly the (improper) discounted state visitation distribution for policy μ_ϕ is ρ^{μ_ϕ} , the deterministic policy gradient is then

$$\nabla_\phi J(\mu_\phi^\times) \approx \mathbb{E}_{\xi_t \sim \rho^{\mu_\phi}} \left[\nabla_\phi \left(1 + \frac{Q_{\mu_\phi}(\xi_t, \mu_\phi(\xi_t))}{\Gamma_{t-1}^+} \right) \right] \tag{169}$$

where the approximation is again due to local maximisation and we must then show $\nabla_\phi J(\mu_\phi^+) \propto \nabla_\phi Q_{\pi_\phi}(s_t, a_t)$ as seen in Eq. (26). This is done similarly to the stochastic case but does not require action sampling $a \in \mathcal{A}$. To avoid hassles with changing order of integrals, we use a discrete tabular version of the original derivation [66]. The additive cumulative reward $\Gamma_{t-1}^+(s_{t-1}, \mu_\phi(s_{t-2}) | h_{t-1}) = \Gamma_{t-1}^+$ is also a constant at each step given the history is completely known. We then express the policy gradient as

$$\begin{aligned}
\nabla_\phi J(\mu_\phi^\times) &\approx \nabla_\phi \left(1 + \frac{Q_{\mu_\phi}(\xi_t, \mu_\phi(\xi_t))}{\Gamma_{t-1}^+} \right) \\
&\stackrel{(a)}{=} \nabla_\phi \sum_{\xi_{t+1}} \left[p_{h_t}(\xi_{t+1} | \xi_t, \mu_\phi(\xi_t)) \left(1 + \frac{r(\xi_t, \mu_\phi(\xi_t)) + \gamma Q_{\mu_\phi}(\xi_{t+1}, \mu_\phi(\xi_{t+1}))}{\Gamma_{t-1}^+} \right) \right] \\
&\stackrel{(b)}{=} \nabla_\phi \frac{r(\xi_t, \mu_\phi(\xi_t))}{\Gamma_{t-1}^+} + \nabla_\phi \sum_{\xi_{t+1}} \left[p_{h_t}(\xi_{t+1} | \xi_t, \mu_\phi(\xi_t)) \left(1 + \frac{\gamma Q_{\mu_\phi}(\xi_{t+1}, \mu_\phi(\xi_{t+1}))}{\Gamma_{t-1}^+} \right) \right]
\end{aligned}$$

$$\begin{aligned}
&\stackrel{(c)}{=} \frac{\nabla_\phi \mu_\phi(\xi_t) \nabla_a r(\xi_t, a_t) |_{a_t=\mu_\phi(\xi_t)}}{\Gamma_{t-1}^+} \\
&+ \sum_{\xi_{t+1}} \left[\nabla_\phi \mu_\phi(\xi_t) \nabla_a p_{h_t}(\xi_{t+1} | \xi_t, \mu_\phi(\xi_t)) |_{a_t=\mu_\phi(\xi_t)} \left(1 + \frac{\gamma Q_{\mu_\phi}(\xi_{t+1}, \mu_\phi(\xi_{t+1}))}{\Gamma_{t-1}^+} \right) \right. \\
&\quad \left. + p_{h_t}(\xi_{t+1} | \xi_t, \mu_\phi(\xi_t)) \nabla_\phi \left(1 + \frac{\gamma Q_{\mu_\phi}(\xi_{t+1}, \mu_\phi(\xi_{t+1}))}{\Gamma_{t-1}^+} \right) \right] \\
&\stackrel{(d)}{=} \nabla_\phi \mu_\phi(\xi_t) \nabla_a \left(\frac{r(\xi_t, a_t)}{\Gamma_{t-1}^+} + \sum_{\xi_{t+1}} \left[p_{h_t}(\xi_{t+1} | \xi_t, a_t) \left(1 + \frac{\gamma Q_{\mu_\phi}(\xi_{t+1}, a_{t+1})}{\Gamma_{t-1}^+} \right) \right] \right) \Big|_{a_i=\mu_\phi(\xi_i)} \\
&+ \sum_{\xi_{t+1}} p_{h_t}(\xi_{t+1} | \xi_t, \mu_\phi(\xi_t)) \nabla_\phi \left(1 + \frac{\gamma Q_{\mu_\phi}(\xi_{t+1}, \mu_\phi(\xi_{t+1}))}{\Gamma_{t-1}^+} \right) \\
&= \nabla_\phi \mu_\phi(\xi_t) \nabla_a \sum_{\xi_{t+1}} \left(p_{h_t}(\xi_{t+1} | \xi_t, a_t) \left(1 + \frac{r(\xi_t, a_t) + \gamma Q_{\mu_\phi}(\xi_{t+1}, a_{t+1})}{\Gamma_{t-1}^+} \right) \right) \Big|_{a_i=\mu_\phi(\xi_i)} \\
&+ \sum_{\xi_{t+1}} p_{h_t}(\xi_{t+1} | \xi_t, \mu_\phi(\xi_t)) \nabla_\phi \left(1 + \frac{\gamma Q_{\mu_\phi}(\xi_{t+1}, \mu_\phi(\xi_{t+1}))}{\Gamma_{t-1}^+} \right) \\
&\stackrel{(e)}{=} \nabla_\phi \mu_\phi(\xi_t) \nabla_a \left(1 + \frac{Q_{\mu_\phi}(\xi_t, a_t)}{\Gamma_{t-1}^+} \right) \Big|_{a_i=\mu_\phi(\xi_i)} + \Lambda(\xi_t, \mu_\phi(\xi_t)) \tag{170}
\end{aligned}$$

where we define the infinitely recursive summation as $\Lambda(\xi_t, \mu_\phi(\xi_t))$. In (a) the Bellman equation in Eq. (126) is used, (b) the rewards are independent of the next state probability distributions, (c) we use the chain rule on the first term and both the product rule on the second term, (d) terms are grouped according to the respective gradients, and (e) the first term is aggregated since it is a constant at a fixed action and the second term includes all additional unrolling via the Bellman equation. This term can be rewritten as

$$\Lambda(\xi_t, \mu_\phi(\xi_t)) = \sum_{\xi_{t+1}} p(\xi_t \rightarrow \xi_{t+1}, 1, \mu_\phi) \nabla_\phi \left(1 + \frac{\gamma Q_{\mu_\phi}(\xi_{t+1}, \mu_\phi(\xi_{t+1}))}{\Gamma_{t-1}^+} \right) \tag{171}$$

with $p(\xi_t \rightarrow \xi', k, \mu_\phi)$ again to be the probability of $\xi_t \rightarrow \xi'$ in k time steps. The recursion can be seen through

$$\begin{aligned}
\Lambda(\xi_t, \mu_\phi(\xi_t)) &= \nabla_\phi \mu_\phi(\xi_t) \nabla_a \sum_{\xi_{t+1}} \left(p(\xi_t \rightarrow \xi_{t+1}, 1, \mu_\phi) \left(1 + \frac{r(\xi_t, a_t) + \gamma Q_{\mu_\phi}(\xi_{t+1}, a_{t+1})}{\Gamma_{t-1}^+} \right) \right) \Big|_{a_i=\mu_\phi(\xi_i)} \\
&+ \sum_{\xi_{t+1}} p(\xi_t \rightarrow \xi_{t+1}, 1, \mu_\phi) \sum_{\xi_{t+2}} p(\xi_{t+1} \rightarrow \xi_{t+2}, 1, \mu_\phi) \nabla_\phi \left(1 + \frac{\gamma^2 Q_{\mu_\phi}(\xi_{t+2}, \mu_\phi(\xi_{t+2}))}{\Gamma_{t-1}^+} \right) \\
&= \nabla_\phi \mu_\phi(\xi_t) \nabla_a \sum_{\xi_{t+1}} \left(p(\xi_t \rightarrow \xi_{t+1}, 1, \mu_\phi) \left(1 + \frac{r(\xi_t, a_t) + \gamma Q_{\mu_\phi}(\xi_{t+1}, a_{t+1})}{\Gamma_{t-1}^+} \right) \right) \Big|_{a_i=\mu_\phi(\xi_i)} \\
&+ \sum_{\xi_{t+1}} p(\xi_t \rightarrow \xi_{t+1}, 2, \mu_\phi) \nabla_\phi \left(1 + \frac{\gamma^2 Q_{\mu_\phi}(\xi_{t+1}, \mu_\phi(\xi_{t+1}))}{\Gamma_{t-1}^+} \right) \\
&\vdots \\
&= \sum_{\xi_{t+1}} \sum_{t=1}^{\infty} p(\xi_t \rightarrow \xi_{t+1}, t, \mu_\phi) \nabla_\phi \mu_\phi(\xi_t) \nabla_a \left(1 + \frac{\gamma^t Q_{\mu_\phi}(\xi_{t+1}, a_{t+1})}{\Gamma_{t-1}^+} \right) \Big|_{a_{t+1}=\mu_\phi(\xi_{t+1})} \tag{172}
\end{aligned}$$

Therefore we can write

$$\nabla_\phi J(\mu_\phi^\times) = \sum_{\xi_{t+1}} \sum_{t=0}^{\infty} p(\xi_t \rightarrow \xi_{t+1}, t, \mu_\phi) \nabla_\phi \mu_\phi(\xi_t) \nabla_a \left(1 + \frac{\gamma^t Q_{\mu_\phi}(\xi_{t+1}, a_{t+1})}{\Gamma_{t-1}^+} \right) \Big|_{a_{t+1}=\mu_\phi(\xi_{t+1})} \quad (173)$$

As before, the impact of constant values does not affect the optimisation process and so we use the approximations in Eqs. (164) which allows us to then construct the usual density

$$\rho^{\mu_\phi}(\xi') \equiv \int_S d\xi \sum_{t=1}^{\infty} \gamma^{t-1} p_1(\xi) p(\xi \rightarrow \xi', t, \mu_\phi) \quad (174)$$

and so finally

$$\begin{aligned} \nabla_\phi J(\mu_\phi^\times) &\approx \nabla_\phi \int_S ds \rho^{\mu_\phi}(s) \left(1 + \frac{Q_{\pi_\phi}(s, \mu_\phi(s))}{\Gamma_{t-1}^+} \right) \\ &= \mathbb{E}_{\xi_t \sim \rho^{\mu_\phi}} \left[\nabla_\phi \left(1 + \frac{Q_{\pi_\phi}(s, \mu_\phi(s))}{\Gamma_{t-1}^+} \right) \right] \end{aligned} \quad (175)$$

$$= \mathbb{E}_{\xi_t \sim \rho^{\mu_\phi}} \left[\frac{\nabla_\phi \mu_\phi(\xi_t) \nabla_a Q_{\pi_\phi}(\xi_t, a_t)}{\Gamma_{t-1}^+} \Big|_{a_t=\mu_\phi(\xi_t)} \right] \quad (176)$$

which is of the desired form [66]. We also assume without proof the limiting case theorem that the stochastic policy actor approaches the deterministic policy actor as variance is reduced

$$\lim_{\sigma \rightarrow 0} \nabla_\phi J(\pi_{\phi,\sigma}^\times) = \nabla_\phi J(\mu_\phi^\times) \quad (177)$$

which is reasonable as the functional form for the multiplicative gradients is identical to the stochastic gradients up to a constant and scaled by cumulative reward. Therefore we do not expect these predefined known values to have any impact on limiting convergence.

The parameters ϕ^\times are therefore updated in the direction of $\nabla_\phi \left(1 + \frac{Q_{\pi_\phi}(\xi_t, a_t)}{\Gamma_{t-1}^+} \right)$ with difference between this and the standard $\nabla_\phi J(\pi_\phi) \propto \nabla_\phi Q_{\pi_\phi}(s_t, a_t)$ being now using return over the existing cumulative episodic return. Therefore, for popular off-policy algorithms such as DDPG [67] and TD3 [71] the deterministic actors maximise the following objectives

$$J(\phi^+) \approx \mathbb{E}_{U(D^+)} [Q_{\pi_\phi}(s_t, \mu_\phi(s_t))] \quad (178)$$

$$J(\phi^\times) \approx \mathbb{E}_{U(D^\times)} \left[\left(1 + \frac{Q_{\pi_\phi}(\xi_t, \mu_\phi(\xi_t))}{\Gamma_{t-1}^+} \right) \right] \quad (179)$$

with the same implications as discussed earlier. Both are approximations due to formulation of deterministic policy gradients not performing global maximisation of Q-values [66], while the second requires an additional approximation due to Eq. (164). Notice that for both the stochastic and deterministic policy gradients, if we performed the simple reparameterisation $Q(s, a) \leftarrow \left(1 + \frac{Q(\xi, a)}{\Gamma_{t-1}^+} \right)$, the existing derivations [29, 65, 66] would have been completely valid and left unchanged, allowing us to effortlessly arrive at the final results. This shortcut however would not have elucidated the strict requirements in Eq. (164) that are needed to represent the gradients in a tractable manner. Furthermore, with the explicit proof we are unable to find any interesting hidden structure.

7 Maximum Causal Entropy with Multiplicative Dynamics

The energy based polices discussed in Section 3.4 are built around soft updates using an maximum entropy objective [72–79]. Modification of this approach to incorporate multiplicative dynamics uses a very similar method to that in Section 5.1. We also explicitly assume the underlying processes is a QDP and therefore we formally extend the soft-actor critic algorithm to this much larger class of process already encompassing all MDPs.

The discounted future compounding rewards with state-dependent entropies are

$$\begin{aligned} 1 + R_t^\times &\equiv (1 + \Gamma_{t-1}^\times) \cdot \frac{\Gamma_{t-1}^+ + (r_t + \alpha H[\pi(\cdot|\xi_t)])}{\Gamma_{t-1}^+} \cdot \frac{(\Gamma_{t-1}^+ + r_t) + \gamma[r_{t+1} + \alpha H(\pi(\cdot|\xi_{t+1}))]}{\Gamma_{t-1}^+ + r_t} \cdot \dots \\ &= (1 + \Gamma_t^\times) \cdot \prod_{k=1}^{\infty} \left(1 + \frac{\gamma^k [r_{t+k} + \alpha H(\pi(\cdot|\xi_{t+k}))]}{\Gamma_{t-1+k}^+} \right) \\ &= (1 + \Gamma_t^\times) \cdot \exp \left[\sum_{k=1}^{\infty} \ln \left| 1 + \frac{\gamma^k [r_{t+k} + \alpha H(\pi(\cdot|\xi_{t+k}))]}{\Gamma_{t-1+k}^+} \right| \right] \end{aligned} \quad (180)$$

$$\propto \sum_{k=1}^{\infty} \left(1 + \frac{\gamma^k [r_{t+k} + \alpha H(\pi(\cdot|\xi_{t+k}))]}{\Gamma_{t-1+k}^+} \right) \quad (181)$$

where $\Gamma_{t-1+k}^+ = \Gamma_{t-1}^+ + \sum_{\lambda=0}^{k-1} \gamma^\lambda r_{t+\lambda}$ do not include entropies, the current return Γ_t^\times is completely known, and we will show again that the initial value $\Gamma_0^+ = V_0 > V_{\min} \geq \gamma$ is required. The final proportionality is based on same logic as in Q-learning in Eq. (114). One important difference is that now the bound is $\gamma^k r_{t+k} \geq (V_{\min} - \Gamma_{t-1+k}^+ - \gamma^k \alpha H(\pi(\cdot|\xi_{t+k})))$ where the equality again ends the episode. This is a particularly robust mechanism to train agents as it adds an additional state-dependent noise to the obtained reward signal, leading to overall much less brittle learning. In other words, this artificially simulates volatility in returns when learning a policy to encourages exploration. The automatically tuned temperature parameter α then governs the magnitude of randomness, where ideally we desire it to increase as the agents simulated performance increase and vice versa.

7.1 Soft Learning

The maximum entropy objective can then be written by in terms of $\pi^* \in \arg \max_{\pi^* \in \Pi} (1 + R_t^\times)$. Formally, first we define the entropy augmented reward $r_{t+k}^\pi \equiv r_{t+k} + \alpha H(\pi(\cdot|\xi_{t+k}))$ that so that

$$1 + R_t^\times \propto \sum_{k=1}^{\infty} \left(1 + \frac{\gamma^k r_{t+k}^\pi}{\Gamma_{t-1+k}^+} \right) \quad (182)$$

which is of the exact form in Eq. (114) and so much of the earlier results apply. Recall the definition for soft values functions in Eqs. (32-33) which we modify to

$$1 + \frac{Q_\pi^{\text{soft}}(\xi_t, a_t)}{\Gamma_{t-1}^+} \equiv \mathbb{E}_{\xi_{t+1} \sim \rho^\pi, a_{t+1} \sim \pi} \left[\sum_{k=1}^{\infty} \left(1 + \frac{\gamma^k r_{t+k}^\pi}{\Gamma_{t-1}^+} \right) \right] \quad (183)$$

$$1 + \frac{V_\theta^{\text{soft}}(\xi_t)}{\Gamma_{t-1}^+} \equiv \alpha \ln \int_{\mathcal{A}} da' \exp \left[\left(\frac{1}{\alpha} \left(1 + \frac{Q_\pi^{\text{soft}}(\xi_t, a')}{\Gamma_{t-1}^+} \right) \right) \right] \quad (184)$$

$$= \alpha \ln \mathbb{E}_{Z_\pi \sim a'} \left[\frac{1}{Z_\pi(a')} \exp \left(\frac{1}{\alpha} \left(1 + \frac{Q_\pi^{\text{soft}}(\xi_t, a')}{\Gamma_{t-1}^+} \right) \right) \right] \quad (185)$$

Next, we use the state distribution ρ^π to define again the approximate (soft) Bellman equation

$$1 + \frac{Q_\pi^{\text{soft}}(\xi_{t+k}, a_{t+k})}{\Gamma_{t-1+k}^+} \leftarrow 1 + \frac{1}{\Gamma_{t-1+k}^+} \mathbb{E}_{\xi_{t+k+1} \sim \rho^\pi, a_{t+k+1} \sim \pi} [r^\pi(\xi_{t+k}, a_{t+k}) + \gamma Q_\pi^{\text{soft}}(\xi_{t+k+1}, a_{t+k+1})] \quad (186)$$

where r_{t+k}^π is again a known constant and so in the limit $\alpha \rightarrow 0$ this result is exactly equivalent to Eq. (115). Based on the MDP additive dynamics case set out in [73–75, 77], we can introduce a soft Bellman history operator $\mathcal{T}_{h_t}^\pi$ that is valid under all previous QDP assumptions to perform soft policy evaluation. This is defined as

$$(\mathcal{T}_{h_t}^\pi Q^{\text{soft}})(\xi, a) = 1 + \frac{Q_\pi^{\text{soft}}(\xi_t, a_t)}{\Gamma_{t-1}^+} = 1 + \frac{\mathbb{E}_{\xi_{t+1} \sim \rho^\pi, a_{t+1} \sim \pi} [r^\pi(\xi_t, a_t) + \gamma Q_\pi^{\text{soft}}(\xi_{t+1}, a_{t+1})]}{\Gamma_{t-1}^+} \quad (187)$$

To prove under a strict criteria the convergence $(\mathcal{T}_{h_t}^\pi Q^{\text{soft}})(\xi, a) = Q_{t+1}^{\text{soft}}(\xi, a) \rightarrow Q^{\text{soft*}}(\xi, a)$ w.p.1. as infinite updates $t \rightarrow \infty$ we must introduce an additional condition that the cardinality of the action space $|\mathcal{A}| < \infty$ is finite to ensure the entropy term is bounded. Now combining this condition, the QDP assumptions, and the conditions in Section 5.2, we see that if condition five can be proven then we can utilise all but the prior results without any additional changes to prove that soft Q-learning converges to the optimal soft Q-value in QDP domains.

To show this we adapt the ideas in [73, 74]. Suppose for a fixed history h_t and fixed point (ξ_t, a_t) the max-norm exists $\epsilon = \frac{1}{\Gamma_{t-1}^+} \|Q_{\pi,t}(\xi, a) - Q'_{\pi,t}(\xi, a)\|_\infty$, if then $Q_{\pi,t} \leq \epsilon + Q'_{\pi,t}$ and similarly $Q_{\pi,t} \geq -\epsilon + Q'_{\pi,t}$ then by the definition of the max-norm

$$\|(\mathcal{T}_{h_t}^\pi Q)(\xi, a) - (\mathcal{T}_{h_t}^\pi Q')(\xi, a)\|_\infty \leq \gamma \epsilon = \frac{\gamma}{\Gamma_{t-1}^+} \|Q_{\pi,t}(\xi, a) - Q'_{\pi,t}(\xi, a)\|_\infty \quad (188)$$

and therefore by the same constraints on Γ_{t-1}^+ we have

$$\|(\mathcal{T}_{h_t}^\pi Q)(\xi, a) - (\mathcal{T}_{h_t}^\pi Q')(\xi, a)\|_\infty \leq \gamma \|Q_{\pi,t}(\xi, a) - Q'_{\pi,t}(\xi, a)\|_\infty \quad (189)$$

hence all other convergence results follow, proving that there exists a unique optimal solution to soft Q-learning that is greedily approached via soft value iteration.

Furthermore if we very generally represent the policy distribution as general energy-based policy of the form of a Boltzmann distribution [12] where

$$\pi(a_t | \xi_t) \equiv \exp \left[\frac{1}{\alpha} \left(1 + \frac{Q_\pi^{\text{soft}}(\xi_t, a_t)}{\Gamma_{t-1}^+} \right) - \ln |Z_\pi(\xi_t)| \right] \quad (190)$$

the soft policy improvement can be defined by the objective to be maximised

$$J(\pi) \approx \mathbb{E}_{\xi_t \sim \rho^\pi, a_t \sim \pi} \left[\left(1 + \frac{Q_\pi^{\text{soft}}(\xi_t, a_t)}{\Gamma_{t-1}^+} \right) \right] \quad (191)$$

with the approximation due to Eq. (186). Using the Kullback-Leibler divergence we can improve the policy $\pi(a_t | \xi_t)$ in Eq. (190) to $\pi(a_t | \xi_t) \rightarrow \pi'(a_t | \xi_t)$ with

$$\begin{aligned} \pi'(\cdot | \xi_t) &= \arg \min_{\pi' \in \Pi} D_{\text{KL}}(\pi'(\cdot | \xi_t) || \pi(\cdot | \xi_t)) \\ &= \arg \min_{\pi' \in \Pi} \mathbb{E}_{a_t \sim \pi'} \left[\ln \pi'(a_t | \xi_t) - \left(\frac{1}{\alpha} \left(1 + \frac{Q_\pi^{\text{soft}}(\xi_t, a_t)}{\Gamma_{t-1}^+} \right) - \ln |Z_\pi(\xi_t)| \right) \right] \end{aligned}$$

$$\leq \arg \min_{\pi \in \Pi} \mathbb{E}_{a_t \sim \pi} \left[\ln \pi(a_t | \xi_t) - \left(\frac{1}{\alpha} \left(1 + \frac{Q_\pi^{\text{soft}}(\xi_t, a_t)}{\Gamma_{t-1}^+} \right) - \ln |Z_\pi(\xi_t)| \right) \right] \quad (192)$$

$$= 0 \quad (193)$$

as we are always free to select $\pi'(a_t | \xi_t) = \pi(a_t | \xi_t)$ and since the expectation is state-independent by definition of the policy in Eq. (190). We can then rewrite this more clearly with the normalising definition of the value function in Eq. (185) as

$$\mathbb{E}_{a_t \sim \pi'} \left[\left(1 + \frac{Q_\pi^{\text{soft}}(\xi_t, a_t)}{\Gamma_{t-1}^+} \right) - \alpha \ln \pi'(a_t | \xi_t) \right] \geq \alpha \mathbb{E}_{a_t \sim \pi'} [\ln |Z_\pi(\xi_t)|] = \left(1 + \frac{V_\pi^{\text{soft}}(\xi_t)}{\Gamma_{t-1}^+} \right) \quad (194)$$

which defines policy improvement per time step. The equality holds when $\pi'(a_t | \xi_t) = \pi(a_t | \xi_t)$, hence we have the update rule for the soft value

$$\left(1 + \frac{V_\pi^{\text{soft}}(\xi_t)}{\Gamma_{t-1}^+} \right) = \mathbb{E}_{a_t \sim \pi} \left[\left(1 + \frac{Q_\pi^{\text{soft}}(\xi_t, a_t)}{\Gamma_{t-1}^+} \right) - \alpha \ln \pi(a_t | \xi_t) \right] \quad (195)$$

The soft Bellman equation in Eq. (186) can be converted using the fact $Q_\pi^{\text{soft}}(\xi_t, a_t) = \mathbb{E}_{a_{t+1} \sim \pi} [V_\pi^{\text{soft}}(\xi_{t+1})]$ to yield

$$1 + \frac{Q_\pi^{\text{soft}}(\xi_t, a_t)}{\Gamma_{t-1}^+} = 1 + \frac{r(\xi_t, a_t) + \gamma \mathbb{E}_{\xi_{t+1} \sim \rho^\pi} [V_\pi^{\text{soft}}(\xi_{t+1})]}{\Gamma_{t-1}^+} \quad (196)$$

To show that this can generalise to policy iteration and allows the soft Q-values to monotonically improve we can apply this fact to the Bellman equation in Eq. (196) repeatedly and construct the

$$\begin{aligned} 1 + \frac{Q_\pi^{\text{soft}}(\xi_t, a_t)}{\Gamma_{t-1}^+} &= 1 + \frac{r(\xi_t, a_t)}{\Gamma_{t-1}^+} + \frac{\gamma}{\Gamma_{t-1}^+} \mathbb{E}_{\xi_{t+1} \sim \rho^\pi} [V_\pi^{\text{soft}}(\xi_{t+1})] \\ &\leq 1 + \frac{r(\xi_t, a_t)}{\Gamma_{t-1}^+} + \frac{\gamma}{\Gamma_{t-1}^+} \mathbb{E}_{\xi_{t+1} \sim \rho^\pi} \left[\mathbb{E}_{a_{t+1} \sim \pi'} \left[\left(1 + \frac{Q_\pi^{\text{soft}}(\xi_{t+1}, a_{t+1})}{\Gamma_{t-1}^+} \right) - \alpha \ln \pi'(a_{t+1} | \xi_{t+1}) \right] \right] \\ &\vdots \\ &\leq 1 + \frac{Q_{\pi'}^{\text{soft}}(\xi_t, a_t)}{\Gamma_{t-1}^+} \end{aligned} \quad (197)$$

hence $\pi(a_t | \xi_t) \rightarrow \pi'(a_t | \xi_t)$ monotonically improves the soft Q-value. The proof that $\pi(a_t | \xi_t) \rightarrow \pi^*(a_t | \xi_t)$ to a unique optimal policy is left unchanged from [75, 77]. The arguments discussing the connection of this approach to policy gradients in Section 6 shown in [74] also remain valid.

Therefore, for off-policy learning with a mini-batch using an MSE loss function the critic objectives to be minimised for both dynamics are

$$J(\theta^+) = \mathbb{E}_{U(\mathcal{D}^+)} \left[\left(R_t^{(m)} + \gamma^m Q_\theta^{\text{soft}}(s_{t+m}, a') - Q_\theta^{\text{soft}}(s_t, a_t) \right)^2 \right] \quad (198)$$

$$J(\theta^\times) \approx \mathbb{E}_{U(\mathcal{D}^\times)} \left[\left(\frac{1}{\Gamma_{t-1}^+} \left(R_t^{(m)} + \gamma^m Q_\theta^{\text{soft}}(\xi_{t+m}, a') - Q_\theta^{\text{soft}}(\xi_t, a_t) \right) \right)^2 \right] \quad (199)$$

where the targets $Q_\theta^{\text{soft}}(s_{t+m}, a')$ are calculated by applying Eq. (195) into Eq. (196). The stochastic policy objectives

to be maximised are then also

$$J(\phi^+) = \mathbb{E}_{U(\mathcal{D}^+)} \left[Q_\theta^{\text{soft}}(s_t, a_t) - \alpha \ln \pi_\phi(a_t | s_t) \right] \quad (200)$$

$$J(\phi^\times) = \mathbb{E}_{U(\mathcal{D}^\times)} \left[\left(1 + \frac{Q_\theta^{\text{soft}}(\xi_t, a_t)}{\Gamma_{t-1}^+} \right) - \alpha \ln \pi_\phi(a_t | \xi_t) \right] \quad (201)$$

with same implications for the differences between the dynamic discussed in Section 5.4 for critics and Section 6 for actors. Overall, we have shown that energy-based policies, specifically the soft actor critic algorithm under the assumption of QDPs can be used to maximise the objective in Eq. (180) with minor modifications. The updating of critic values were approximated using Eq. (186) while the policy update procedure was left unchanged, this is unlike the requirement for policy gradients in Eq. (164). One key difference to policy gradients that effectively maximise the return from a given action, is that this structure maximises the advantage, namely the difference between the return from a given action and the average return across all actions.

8 Energy Efficient Agent Inference

Real-world environments have near-infinite, but ultimately finite number possible states and can be accurately represented using HDPs discussed in Section 3.8. These are also accompanied by a finite number of known possible actions the agent has control over at each time interval. As discussed earlier, this problem is often intractable and so we generally simplify the system to be a QDP, crudely considered a history-independent POMDP [109].

Furthermore, for the combined results of Sections 5-7 to be valid, one key requirement is the both the state space $\Xi : \mathcal{S} \times \mathcal{R}$ (where recall $|\Xi| = |\mathcal{S}| + 1$) and action space \mathcal{A} are finite. In terms of the cardinality of these spaces ($|\mathcal{S}|, |\mathcal{A}|$) for locomotive continuous control tasks, they are quite small to the simulations to be practical. The MuCoJo environments [53] tested in [71, 75, 77, 78] reveal how exponentially more difficult training agents become as you increase the both spaces. The more contemporary PyBullet engine [61] interfaced through the OpenAI gym [54] also features a range of $(|\mathcal{S}|, |\mathcal{A}|) = (4, 1) \rightarrow (44, 17)$ with the highly non-linear difficulty scaling similarly seen in [155]. The DeepMimic simulations [57] can range from $(|\mathcal{S}|, |\mathcal{A}|) = (197, 36) \rightarrow (418, 94)$ and result in far more realistic and natural movements of trained agents with an uncanny resemblance to humans.

A more realistic example would be that of the popular esport video game of StarCraft II wherein the trained agent AlphaStar competes in real-time against human players and achieves a rank exceeding 99.8% of the community [49]. It achieves such success by operating with constraints very similar to that of the reaction time of top-tier human players and vision constraints identical to any human player. Due to the vision of each player being limited to their particular position and a crude ‘minimap’ revealing the general positioning only if they have structures or units in a region of the complete map, therefore it is an imperfect game approximated as a MDP. Each state leads to a set of sequential action decisions: action type (several hundred), who to issue action to (units and or structures), where to target, where to observe and act next, and where to move the agents camera view. Ultimately this leads to approximately 10^{26} possible choices at each step. Interestingly, it achieves such high-levels success while having on average fewer actions per minute than humans which is often considered a key performance metric as it is correlated with the amount of multi-tasking a player performs in order to favourably fine-tune each outcome.

All these systems consist of two separate processes. The first is the training of the agent that can occur at a central location with extensive computational throughput and so while increasing learning speed via algorithmic efficiency is highly sought after, it can also be achieved using with superior quality and or larger quantities of hardware. These dedicated supercomputing clusters are notoriously well-known for their high levels power consumption, leading to immense heat generation, and subsequently requiring constant cooling to prevent, at best thermal throttling, or at worst catastrophic hardware failure coupled with extensive damage to components. Next, once these agents have been trained, decision-making involves simply processing the input state and then taking the optimal action according to the trained policy.

Modern reinforcement learning extensively utilised the use of deep neural networks due to their remarkable ability to act as universal function approximator that allow far superior parametrisation compared to the tabular approach used prior [29]. Regardless of network architecture, training involves improving all network parameters using backpropagation against a target value [156–159], evaluation involves forward propagation through up to or less than all the network parameters.

The simplest architecture is a fully-connected feed-forward neural network has known number of inputs I and outputs O , n hidden layers indexed $k = 1, \dots, n$ with each layer having h_k nodes, and features the inclusion of a

singular bias unit in all but the output layer that are not connected to the previous layer. The total number of trainable weights in this setup are

$$N = \underbrace{Ih_1 + \sum_{k=1}^{n-1} h_k h_{k+1} + h_n O}_{\text{Fully Connected}} + \underbrace{\sum_{k=1}^n h_k + O}_{\text{Bias Connections}} \quad (202)$$

Reinforcement learning using the SAC and TD3 agent learning algorithms presented in Appendix A consist of three feed-forward neural networks, two identically structured critic Q-value approximators, and a single actor policy parametrisation. All three networks consist of two hidden layers but the value and policy networks have differing quantities of inputs and outputs. Using multiplicative dynamics we have total parameters

$$N_Q = h_1(|\Xi| + |\mathcal{A}| + h_2 + 1) + h_2 + 1 \quad (203)$$

$$N_\pi^\lambda = h_1(|\Xi| + h_2 + 1) + \lambda|\mathcal{A}|(h_2 + 1) + h_2 \quad (204)$$

where $\lambda = 1, 2$ for TD3 and SAC actor policies respectively. Therefore, for each agent, the total number of parameters needed to be trained are

$$\begin{aligned} N^\lambda &= 2N_Q + N_\pi^\lambda \\ &= h_1(3|\Xi| + 2|\mathcal{A}| + 3h_2 + 3) + \lambda|\mathcal{A}|(h_2 + 1) + 5h_2 + 2 \end{aligned} \quad (205)$$

$$N_\pi^{\text{SAC}} = 768|\Xi| + 1,026|\mathcal{A}| + 198,658 \quad (206)$$

$$N^{\text{TD3}} = 1,200|\Xi| + 1,101|\mathcal{A}| + 362,702 \quad (207)$$

using the hyperparameters in Table 1. Importantly, training all these parameters can occur at dedicated supercomputers. Once trained, only forward propagation through the policy network is required which is a substantially easier task as it essentially breaks down to matrix multiplication from the input to output layer. This agent receives state $\xi \in \Xi$ and takes actions $a \in \mathcal{A}$ to maximise the reward signal indefinitely. In this case the number of parameters to be utilised for inference are

$$N_\pi^{\text{SAC}} = 256|\Xi| + 514|\mathcal{A}| + 66,048 \quad (208)$$

$$N_\pi^{\text{TD3}} = 400|\Xi| + 301|\mathcal{A}| + 120,700 \quad (209)$$

and so linearly scale with the size of both state and action spaces. For example, inference of the DeepMimic simulations will have parameter count ranges of $134,984 \rightarrow 221,372$ and $210,336 \rightarrow 316,194$ for SAC and TD3 respectively. While any shallow learning model with this many parameters almost certainly would lead to overfitting, deep learning has remained triumphant in its ability to achieve world-class accuracy [159]. Furthermore, any modern or even somewhat dated computing hardware can effortlessly process these calculations, however, continuous operation will again consume power, output heat, and naturally demand cooling.

In comparison, AlphaStar has $139 \cdot 10^6$ trainable parameters while only $55 \cdot 10^6$ parameters are utilised for inference [49]. AlphaStar also performed this forward propagation on slightly dated computing hardware that can be purchased from any consumer electronics store. This is a remarkable feat, the ability for a regular computer to defeat professional gamers in arguably the most globally competitive real-time strategy game, heralds far greater implications than the

better known successes of AlphaGo and its boardgame counterparts.

Suppose now we designing agents to operate in environments where repair, part-replacement, or any contact post-deployment is not possible. Furthermore, lets assume the agent is installed on battery power and so the maximum operating time of the device is directly proportional to its built-in power level. There is also interplay with cooling, a device that utilises large amounts of continuous electron flux will require greater heat dissipation than the same device with lower flow. Coolant regulation also may utilise the on-board fixed power such as fans. This is a highly coupled problem, the desire is then to design agents that ultimately consume the least power for the equivalent performance, that is, the best bang for the buck.

The question of measuring energy consumed by a neural network can be simplistically be considered proportional to the number of multiply-accumulate (MAC) operations as a proxy for the number of floating-point operations, combined with the number of weights to model the the number of main memory accesses [160, 161]. As the weights have to be loaded from memory, they have a high relative energy cost relative to MAC operations [162]. Therefore, reducing the number of parameters to forward propagate through can reduce power consumption during inference [163–166]. There are numerous methods of doing so such as pruning [162, 164], compression [167], and compacting [168] the architecture.

Therefore, for reinforcement learning agents operating in these extreme environments, by truncating the number of model parameters for inference we can extended their useful life while operating on a fixed battery. While such applications are likely to only emerge in the somewhat distant future once highly-practical and fully functioning agents are used in the real-world, we believe it is crucial to begin developing the theory given the accelerating progress in reinforcement learning is only going to increase.

Our purpose is also to construct complete end-to-end agents that can operate with zero external inputs. A further key requirement we limit ourselves to however is that the agent must achieve at worst equivalent performance to the ‘inefficient’ agent while using fewer nodes. Obviously there may be cases where we may achieve great power savings for slightly reduced performance, these more interesting situations are outside the scope of our analysis.

8.1 Multi-Stage Policy Control

The open question is then in what environments are we able to reduce the number of parameters for agent inference without affecting performance. Using SAC and TD3 in Eqs. (208–209) as examples, we can either reduce $|\Xi|$, $|\mathcal{A}|$, or both. Suppose then all parameters are still required for operation, but not all of them at all times. In such situations assume we are able to segregate environment into p separate phases or stages indexed $\varrho = 1, \dots, p$.

Consider the original agent requiring inference through N_π parameters, and that the number of parameters required for any stage p are however $N_\pi^p < N_\pi \forall p$. In what environments would this be possible? This would involve situations where certain states, actions, or both are not required for inference.

A prime example of this in the future would probes for deep space exploration. Contemporarily, their simplified operation involves: Earthly launch through atmosphere, reaching escape velocity, autonomously charting course to destination using $G_{\mu\nu} + \Lambda g_{\mu\nu} = \frac{8\pi G}{c^4} T_{\mu\nu}$, deploying probes, and all while communicating back to mission control. This is an extreme environment where the agent has access to a relatively fixed battery as the utility of solar energy will be diminished at large distances.

Extending the operational life of this probes is equivalent to increasing the amount of scientific data they can

uniquely provide, unique because the alternative of sending another probe to their location would take time. Another advantage of end-to-end agents would be that at no time would mission control be required to intervene, hence communication can be limited to be solely one-way. It is much easier for a object to target Earth which travels on a fixed trajectory, then it is for us to target a distant probe, further reducing power consumption. This effect would reduce heat output and require less coolant, enabling either lighter vessels or more equipment. Note that the vacuum of space prohibits any heat transfer other than infrared photon emissions.

Given the feature of distinct sequential stages the probe must adhere to at all times, we can effectively use different agents for each phase. Furthermore, for simplicity we consider only irreversible transitions to the next stage with no possibility of the reverse. One agent can govern take-off, another controls upper atmosphere manoeuvres, the next manages the journey, and the final organises the delivery of probes to their optimal destinations. This procedure could also be adopted for other *more explosive* applications.

Regarding states and actions, each phase is likely to share common features while also containing unique variables. For example, states and actions pertaining to launch tied to events such as wind, weather, and atmospheric composition, are not necessary after this phase. This segregation of responsibilities is a somewhat trivial problem but complexity arises on how to teach the agent to learn to recognise the transitions.

— OLD BELOW

For environments where actions can be separated into distinct phases, a general framework for constructing all manner of guidance systems from simple goal-based systems to far more complex and intelligent tools is needed in order to make the problem tractable by isolating the relevant action spaces. This will be done by building a three-stage agent with distinct phases: standby, active, and engage. We call this prototype method SARAEnⁿ. The acronym is for state-action-reward, then if reward meets a threshold, the agent initiates an ‘active’ phase (distinct second set of actions) which directly leads to an engage phase (distinct third set of actions). The $n = 2$ index denotes that the threshold creates only two different decision loops, the $n = 1$ case reduces to the usual existing formulation. This naming is clearly a modification of the well-known SARSA algorithm [29].

Applications involve rockets for deep-space exploration and atmospheric entry, self-driving cars, and items that have exponentially larger potential to explosively change the landscape. In Figs. we show the blueprints and schematics for this algorithm. Take particular note of how the agent is able to determine what stage it is in and able to take optimal actions for this situation.

9 Related Work

There are numerous areas where existing research has been done and served as inspirations for our results. Most of these related areas were discussed in their corresponding discussions in Sections 2-3. In this brief section we provide a succinct review of these motivations and the locations where our work aims to address.

There appears to be no prior work in the area trying different critic loss functions other than the standard MSE in reinforcement learning. Whether it be model-based or model-free algorithms, everyone appears to use MSE. With the authoritative [84] stating 15 years ago that similar results are to be expected for other loss functions. Our analysis reimagines the use of different losses to be analogous to the penalty applied to an agent when it makes a perceived error in valuing a given state-action pair. An a large outlier in this case is interpreted to be large ‘mistake’ relative to the rest of the mini-batch, this outlier then correspondingly forms a larger portion of the aggregated loss to be minimised, and so the agent places more emphasis on adjusting its future parameters to minimise the impact of the same outlier occurring again. Recall this perceived error is nothing more than an artifact of Eqs (28). The target is simply composed of the current known reward and the discounted current Q-value that utilising delay parameter updates, it is essentially a moving target, unlike the fixed true known values that exist in NMF. The question is then whether then Eqs. (41-47) that offer different degrees of outlier smoothing change the final result. While all these have been thoroughly tested in NMF with outstanding success [85, 123–128], there does not appear to exist and conclusive results in reinforcement learning. The use of stronger functions will therefore be a empirical test of whether relatively large mistakes are important to learning.

On the topic of shadow means there is nothing but emptiness in the literature regarding its utility in not just machine learning, but all of statistics [9]. The utility of this approach for practical purposes has yet to be tested in reality, as it stands, it is a theoretical construct for estimating in a probabilistic manner the true population mean for a fat tailed, right-skewed distribution. Its estimation can be split into two components. First we need to estimate the tail coefficient $\hat{\alpha}$, and if $\hat{\alpha} < 1$, the shadow mean estimate can be used. There is ample literature on tail estimation, expressed mainly in the field of Extreme Value Theory (EVT) [130]. As discussed in Section 3.1, there are plethora of choices, one is to determine the appropriate intermediate order statistic either using the Hill estimate [136], or more advanced [130, 137, 138] method of moments or MLE techniques. The other is to directly observe the gradient of a Zipf plot of certain number of descending order statistic [139–142]. EVT is a mature field but is yet to offer an objective approach to not only estimating $\hat{\alpha}$, but even agreeing on the validity of some of these approaches. In our case, the first test is whether the estimated $\hat{\alpha} < 1$ conclusively which will indicate critic losses are fat-tailed and have no true mean when modeled using GPDs, and then whether the shadow mean [9] can be utilised by the agent to achieve greater success. While [9] have shown it to be useful in redefining the way we examine extreme events, applications to reinforcement learning and most of statistics are non-existent.

Multi-step returns have a very long history of being used to accelerate agent learning, but only for discrete action spaces [29, 39, 62, 89]. Whether they work for continuous action domains is a complete unknown in the history of reinforcement learning. Unlike discrete domains where the Q-value is explicitly used to determine the next action via ϵ -greedy polices, for TD3 and SAC, the Q-values are used to update the policy usually in the direction of maximal value. This is similar to the discrete case, but due to the addition of large amounts of noise added across the mini-batch it is unclear what the final outcome will become. Next the coupling with to experience replay [89] is an extremely puzzling phenomenon. This effect is incredibly difficult to formulate a well-grounded theoretical explanation for as

there is no clear reason why tuning for m -steps and size of the buffer can improve performance across the board.

There has also been zero attempt to incorporate true multiplicative dynamics into reinforcement learning. While the reformulation of contemporary decision theory is discussed in detail in [1–6, 9–11, 90–95], designing of reinforcement learning algorithms that can fully autonomously self-learn these principles is unknown. To test this, we require training environments where the asymptotic reward scaling condition in Eq. (97) must universally hold. To the best of our knowledge there does not exist any publicly available open-source environments that admit this sort of structure. This is both very unfortunate and a great opportunity, it is a detraction as we are unable to experimentally verify majority of our theoretical developments, however, this also means we are the first to come across this deficiency. Therefore, while this work will not be able to test these results, we will pave the way for everyone else by constructing these suitable environments at a later date.

Overall, of the four areas of investigation, none have been thoroughly examined, most have never been considered. Analysis of different critic losses interpreted as whether large mistakes are important for learning has not been considered. Shadow means and tail exponents appear to be vacant in machine learning literature and are usually omitted from the general study of statistics despite most distributions in the real-world being fat tailed. Multi-step returns have been thoroughly analysed for discrete action spaces, yet there does not appear to exist any such analysis for continuous action spaces that are far more relevant to the real-world. Multiplicative dynamics is non-existent in majority of all discussions conducted on Earth even though it is the definition of how wealth, health, and life evolve for any random individual or institution.

10 Experiments

Currently running very lengthy experimental trials. Due Artemis HPC allocation limits, unable to give a clear time-frame for completion, probably late September to mid-October. Preliminary results and discussion are available at <https://github.com/r-grewal/data5709>. Ignore all Sections other than 9-10 as they have all been greatly superseded.

References

- [1] Peters, O. “Optimal leverage from non-ergodicity”. *Quantitative Finance* 11, 11 (2011), pp. 1593–1602.
- [2] Peters, O. “The time resolution of the St Petersburg paradox”. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences* 369, 1956 (Dec. 2011), pp. 4913–4931.
- [3] Peters, O. *Menger 1934 revisited*. 2011. arXiv: [1110.1578 \[q-fin.RM\]](https://arxiv.org/abs/1110.1578).
- [4] Peters, O. and Gell-Mann, M. “Evaluating gambles using dynamics”. *Chaos: An Interdisciplinary Journal of Nonlinear Science* 26, 2 (2016), p. 023103. ISSN: 1089-7682.
- [5] Peters, O. “The ergodicity problem in economics”. *Nature Physics* 15, 12 (Dec. 2019), pp. 1216–1221.
- [6] Peters, O. and Adamou, A. *The time interpretation of expected utility theory*. 2021. arXiv: [1801.03680 \[q-fin.EC\]](https://arxiv.org/abs/1801.03680).
- [7] Meder, D. et al. *Ergodicity-breaking reveals time optimal decision making in humans*. 2020. arXiv: [1906.04652 \[econ.GN\]](https://arxiv.org/abs/1906.04652).
- [8] Peters, O. et al. *What are we weighting for? A mechanistic model for probability weighting*. 2020. arXiv: [2005.00056 \[econ.TH\]](https://arxiv.org/abs/2005.00056).
- [9] Taleb, N. N. *Statistical Consequences of Fat Tails: Real World Preasymptotics, Epistemology, and Applications (Technical Incerto)*. STEM Academic Press, 2020. ISBN: 978-1544508054.
- [10] Peters, O. and Klein, W. “Ergodicity Breaking in Geometric Brownian Motion”. *Phys. Rev. Lett.* 110 (10 2013), p. 100603.
- [11] Peters, O. and Adamou, A. *The sum of log-normal variates in geometric Brownian motion*. 2018. arXiv: [1802.02939 \[cond-mat.stat-mech\]](https://arxiv.org/abs/1802.02939).
- [12] Landau, L. D. and Lifshitz, E. M. *Statistical Physics, Part 1: Volume 5 (Course of Theoretical Physics Series)*. Third Edition. Butterworth-Heinemann, 1980. ISBN: 978-0750633727.
- [13] Spitznagel, M. *Safe Haven Investing - Part One: Not All Risk Mitigation is Created Equal*. Universa Investments L.P., 2017.
- [14] Spitznagel, M. *Safe Haven Investing - Part Two: Not All Risk is Created Equal*. Universa Investments L.P., 2017.
- [15] Spitznagel, M. *Safe Haven Investing - Part Three: Those Wonderful Tenbaggers*. Universa Investments L.P., 2017.
- [16] Spitznagel, M. *Safe Haven Investing - Part Four: The Volatility Tax*. Universa Investments L.P., 2018.
- [17] Spitznagel, M. *Safe Haven Investing: Amor Fati (The Love of One's Fate)*. Universa Investments L.P., 2019.
- [18] Spitznagel, M. *Safe Haven Investing: Why Do People Still Invest in Hedge Funds?* Universa Investments L.P., 2020.
- [19] Spitznagel, M. *Interim Decennial Letter*. Universa Investments L.P., 2020.
- [20] Goetzmann, W. et al. “Portfolio Performance Manipulation and Manipulation-proof Performance Measures”. *Review of Financial Studies* 20, 5 (May 2007), pp. 1503–1546.

- [21] Jensen, J. L. W. V. “Sur les fonctions convexes et les inégalités entre les valeurs moyennes”. *Acta Mathematica* 30, 0 (1906), pp. 175–193.
- [22] Planck Collaboration et al. “Planck 2018 results - VI. Cosmological parameters”. *Astronomy & Astrophysics* 641 (Sept. 2020), A6.
- [23] Dalio, R. *Principles: Life and Work*. Simon & Schuster, Sept. 2017. ISBN: 978-1501124020.
- [24] Hull, J. *Options, futures, and other derivatives*. 10th Edition. Upper Saddle River, NJ: Pearson Prentice Hall, 2018. ISBN: 978-0131977051.
- [25] Doctor, J. N., Wakker, P. P., and Wang, T. V. “Economists’ views on the ergodicity problem”. *Nature Physics* 16, 12 (Dec. 2020), pp. 1168–1168.
- [26] Peters, O. “Reply to: Economists’ views on the ergodicity problem”. *Nature Physics* 16, 12 (Dec. 2020), pp. 1169–1169.
- [27] Despalins, R., Antolin, P., and Payet, S. *Pension Markets in Focus 2020*. OECD, 2020.
- [28] Despalins, R., Antolin, P., and Payet, S. *Pension Funds in Figures 2021*. OECD, 2021.
- [29] Sutton, R. S. and Barto, A. G. *Reinforcement Learning: An Introduction*. Second Edition. The MIT Press, 2018. ISBN: 978-0262039246.
- [30] Bertsekas, D. *Dynamic Programming and Optimal Control, Vol. II: Approximate Dynamic Programming*. Fourth Edition. Belmont, Mass: Athena Scientific, 2012. ISBN: 978-1886529441.
- [31] Bertsekas, D. *Dynamic Programming and Optimal Control, Vol. I*. Fourth Edition. Belmont, Mass: Athena Scientific, 2017. ISBN: 978-1886529397.
- [32] Bertsekas, D. *Reinforcement Learning and Optimal Control*. Belmont, Mass: Athena Scientific, 2019. ISBN: 978-1886529434.
- [33] Mnih, V. et al. *Playing Atari with Deep Reinforcement Learning*. 2013. arXiv: [1312.5602 \[cs.LG\]](https://arxiv.org/abs/1312.5602).
- [34] Mnih, V. et al. “Human-level control through deep reinforcement learning”. *Nature* 518, 7540 (Feb. 2015), pp. 529–533.
- [35] Hasselt, H. v., Guez, A., and Silver, D. “Deep Reinforcement Learning with Double Q-Learning”. *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence*. AAAI’16. Phoenix, Arizona: AAAI Press, 2016, pp. 2094–2100.
- [36] Wang, Z. et al. “Dueling Network Architectures for Deep Reinforcement Learning”. *Proceedings of The 33rd International Conference on Machine Learning*. Ed. by Balcan, M. F. and Weinberger, K. Q. Vol. 48. Proceedings of Machine Learning Research. New York, New York, USA: PMLR, 2016, pp. 1995–2003.
- [37] Bellemare, M. G., Dabney, W., and Munos, R. “A Distributional Perspective on Reinforcement Learning”. *Proceedings of the 34th International Conference on Machine Learning*. Ed. by Precup, D. and Teh, Y. W. Vol. 70. Proceedings of Machine Learning Research. PMLR, 2017, pp. 449–458.
- [38] Fortunato, M. et al. *Noisy Networks for Exploration*. 2019. arXiv: [1706.10295 \[cs.LG\]](https://arxiv.org/abs/1706.10295).
- [39] Hessel, M. et al. *Rainbow: Combining Improvements in Deep Reinforcement Learning*. 2017. arXiv: [1710.02298 \[cs.AI\]](https://arxiv.org/abs/1710.02298).

- [40] Dabney, W. et al. *Distributional Reinforcement Learning with Quantile Regression*. 2017. arXiv: [1710.10044 \[cs.AI\]](#).
- [41] Castro, P. S. et al. *Dopamine: A Research Framework for Deep Reinforcement Learning*. 2018. arXiv: [1812.06110 \[cs.LG\]](#).
- [42] Hafner, D. et al. *Mastering Atari with Discrete World Models*. 2021. arXiv: [2010.02193 \[cs.LG\]](#).
- [43] Silver, D. et al. “Mastering the game of Go with deep neural networks and tree search”. *Nature* 529, 7587 (Jan. 2016), pp. 484–489.
- [44] Silver, D. et al. “Mastering the game of Go without human knowledge”. *Nature* 550, 7676 (Oct. 2017), pp. 354–359.
- [45] Silver, D. et al. *Mastering Chess and Shogi by Self-Play with a General Reinforcement Learning Algorithm*. 2017. arXiv: [1712.01815 \[cs.AI\]](#).
- [46] Silver, D. et al. “A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play”. *Science* 362, 6419 (Dec. 2018), pp. 1140–1144.
- [47] Schrittwieser, J. et al. “Mastering Atari, Go, chess and shogi by planning with a learned model”. *Nature* 588, 7839 (Dec. 2020), pp. 604–609.
- [48] Tomašev, N. et al. *Assessing Game Balance with AlphaZero: Exploring Alternative Rule Sets in Chess*. 2020. arXiv: [2009.04374 \[cs.AI\]](#).
- [49] Vinyals, O. et al. “Grandmaster level in StarCraft II using multi-agent reinforcement learning”. *Nature* 575, 7782 (Oct. 2019), pp. 350–354.
- [50] Senior, A. W. et al. “Improved protein structure prediction using potentials from deep learning”. *Nature* 577, 7792 (Jan. 2020), pp. 706–710.
- [51] Jumper, J. et al. “Highly accurate protein structure prediction with AlphaFold”. *Nature* (July 2021).
- [52] Tunyasuvunakool, K. et al. “Highly accurate protein structure prediction for the human proteome”. *Nature* (July 2021).
- [53] Todorov, E., Erez, T., and Tassa, Y. “MuJoCo: A physics engine for model-based control”. *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, Oct. 2012.
- [54] Brockman, G. et al. *OpenAI Gym*. 2016. arXiv: [1606.01540 \[cs.LG\]](#).
- [55] Kenneally, G., De, A., and Koditschek, D. E. “Design Principles for a Family of Direct-Drive Legged Robots”. *IEEE Robotics and Automation Letters* 1, 2 (July 2016), pp. 900–907.
- [56] Tan, J. et al. *Sim-to-Real: Learning Agile Locomotion For Quadruped Robots*. 2018. arXiv: [1804.10332 \[cs.RO\]](#).
- [57] Peng, X. B. et al. “DeepMimic”. *ACM Transactions on Graphics* 37, 4 (Aug. 2018), pp. 1–14.
- [58] Hafner, R. et al. *Towards General and Autonomous Learning of Core Skills: A Case Study in Locomotion*. 2020. arXiv: [2008.12228 \[cs.RO\]](#).
- [59] Springenberg, J. T. et al. *Local Search for Policy Iteration in Continuous Control*. 2020. arXiv: [2010.05545 \[cs.LG\]](#).

- [60] Team, O. E. L. et al. *Open-Ended Learning Leads to Generally Capable Agents*. 2021. arXiv: [2107.12808 \[cs.LG\]](#).
- [61] Coumans, E. and Bai, Y. *PyBullet, a Python module for physics simulation for games, robotics and machine learning*. 2016–2021.
- [62] Sutton, R. S. “Learning to predict by the methods of temporal differences”. *Machine Learning* 3, 1 (Aug. 1988), pp. 9–44.
- [63] Watkins, C. J. C. H. and Dayan, P. “Q-learning”. *Machine Learning* 8, 3-4 (May 1992), pp. 279–292.
- [64] Hasselt, H. van. “Double Q-learning”. *Advances in Neural Information Processing Systems*. Ed. by Lafferty, J. et al. Vol. 23. Curran Associates, Inc., 2010.
- [65] Sutton, R. S. et al. “Policy Gradient Methods for Reinforcement Learning with Function Approximation”. *Proceedings of the 12th International Conference on Neural Information Processing Systems*. NIPS’99. Denver, CO: MIT Press, 1999, pp. 1057–1063.
- [66] Silver, D. et al. “Deterministic Policy Gradient Algorithms”. *Proceedings of the 31st International Conference on Machine Learning*. Ed. by Xing, E. P. and Jebara, T. Vol. 32. Proceedings of Machine Learning Research 1. Bejing, China: PMLR, 2014, pp. 387–395.
- [67] Lillicrap, T. P. et al. *Continuous control with deep reinforcement learning*. 2019. arXiv: [1509.02971 \[cs.LG\]](#).
- [68] Popov, I. et al. *Data-efficient Deep Reinforcement Learning for Dexterous Manipulation*. 2017. arXiv: [1704.03073 \[cs.LG\]](#).
- [69] Lowe, R. et al. *Multi-Agent Actor-Critic for Mixed Cooperative-Competitive Environments*. 2020. arXiv: [1706.02275 \[cs.LG\]](#).
- [70] Barth-Maron, G. et al. *Distributed Distributional Deterministic Policy Gradients*. 2018. arXiv: [1804.08617 \[cs.LG\]](#).
- [71] Fujimoto, S., Hoof, H. van, and Meger, D. *Addressing Function Approximation Error in Actor-Critic Methods*. 2018. arXiv: [1802.09477 \[cs.AI\]](#).
- [72] Ziebart, B. D. “Modeling purposeful adaptive behavior with the principle of maximum causal entropy”. PhD thesis. Carnegie Mellon University, 2010.
- [73] Fox, R., Pakman, A., and Tishby, N. “Taming the Noise in Reinforcement Learning via Soft Updates”. *Proceedings of the Thirty-Second Conference on Uncertainty in Artificial Intelligence*. UAI’16. Jersey City, New Jersey, USA: AUAI Press, 2016, pp. 202–211. ISBN: 978-0996643115.
- [74] Haarnoja, T. et al. “Reinforcement Learning with Deep Energy-Based Policies”. *Proceedings of the 34th International Conference on Machine Learning - Volume 70*. ICML’17. Sydney, NSW, Australia: JMLR.org, 2017, pp. 1352–1361.
- [75] Haarnoja, T. et al. *Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor*. 2018. arXiv: [1801.01290 \[cs.LG\]](#).
- [76] Haarnoja, T. et al. *Composable Deep Reinforcement Learning for Robotic Manipulation*. 2018. arXiv: [1803.06773 \[cs.LG\]](#).

- [77] Haarnoja, T. et al. *Soft Actor-Critic Algorithms and Applications*. 2019. arXiv: [1812.05905 \[cs.LG\]](#).
- [78] Haarnoja, T. et al. *Learning to Walk via Deep Reinforcement Learning*. 2019. arXiv: [1812.11103 \[cs.LG\]](#).
- [79] Christodoulou, P. *Soft Actor-Critic for Discrete Action Settings*. 2019. arXiv: [1910.07207 \[cs.LG\]](#).
- [80] Schulman, J. et al. *Trust Region Policy Optimization*. 2017. arXiv: [1502.05477 \[cs.LG\]](#).
- [81] Schulman, J. et al. *Proximal Policy Optimization Algorithms*. 2017. arXiv: [1707.06347 \[cs.LG\]](#).
- [82] Li, Z. et al. *Reinforcement Learning for Robust Parameterized Locomotion Control of Bipedal Robots*. 2021. arXiv: [2103.14295 \[cs.RO\]](#).
- [83] Mania, H., Guy, A., and Recht, B. *Simple random search provides a competitive approach to reinforcement learning*. 2018. arXiv: [1803.07055 \[cs.LG\]](#).
- [84] Szepesvári, C. *Algorithms for Reinforcement Learning*. 2009. ISBN: 978-1608454921.
- [85] Guan, N. et al. “Truncated Cauchy Non-Negative Matrix Factorization”. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 41, 1 (Jan. 2019), pp. 246–259.
- [86] Lin, L.-J. “Self-improving reactive agents based on reinforcement learning, planning and teaching”. *Machine Learning* 8, 3-4 (May 1992), pp. 293–321.
- [87] Schaul, T. et al. *Prioritized Experience Replay*. 2016. arXiv: [1511.05952 \[cs.LG\]](#).
- [88] Zhang, S. and Sutton, R. S. *A Deeper Look at Experience Replay*. 2018. arXiv: [1712.01275 \[cs.LG\]](#).
- [89] Fedus, W. et al. “Revisiting Fundamentals of Experience Replay”. *Proceedings of the 37th International Conference on Machine Learning*. Ed. by III, H. D. and Singh, A. Vol. 119. Proceedings of Machine Learning Research. PMLR, 2020, pp. 3061–3071.
- [90] Peters, O. and Adamou, A. *Insurance makes wealth grow faster*. 2017. arXiv: [1507.04655 \[q-fin.RM\]](#).
- [91] Peters, O. and Adamou, A. *An evolutionary advantage of cooperation*. 2018. arXiv: [1506.03414 \[nlin.AO\]](#).
- [92] Adamou, A. T. I. et al. *Microfoundations of Discounting*. 2020. arXiv: [1910.02137 \[econ.TH\]](#).
- [93] Peters, O. and Adamou, A. *Leverage efficiency*. 2020. arXiv: [1101.4548 \[q-fin.GN\]](#).
- [94] Berman, Y., Peters, O., and Adamou, A. “Wealth Inequality and the Ergodic Hypothesis: Evidence from the United States”. *SSRN Electronic Journal* (2020).
- [95] Adamou, A., Berman, Y., and Peters, O. *The Two Growth Rates of the Economy*. 2020. arXiv: [2009.10451 \[econ.GN\]](#).
- [96] Kelly, J. L. “A new interpretation of information rate”. *The Bell System Technical Journal* 35, 4 (1956), pp. 917–926.
- [97] Cover, T. M. and Thomas, J. A. *Elements of information theory*. Hoboken, N.J: Wiley-Interscience, 2006. ISBN: 978-0471241959.
- [98] Markowitz, H. M. “Portfolio selection”. *The Journal of Finance* 7, 1 (Mar. 1952), pp. 77–91.
- [99] Markowitz, H. M. “Investment for the long run: New evidence for an old rule”. *The Journal of Finance* 31, 5 (Dec. 1976), pp. 1273–1286.

- [100] Markowitz, H. M. *Portfolio selection : efficient diversification of investments*. Cambridge, Mass: B. Blackwell, 1991. ISBN: 978-1557861085.
- [101] L'Her, J.-F., Masmoudi, T., and Krishnamoorthy, R. K. "Net Buybacks and the Seven Dwarfs". *Financial Analysts Journal* 74, 4 (Sept. 2018), pp. 57–85.
- [102] Phalippou, L. "An Inconvenient Fact: Private Equity Returns & The Billionaire Factory". *SSRN Electronic Journal* (June 2020).
- [103] Gell-Mann, M. and Hartle, J. B. "Decoherent histories quantum mechanics with one real fine-grained history". *Physical Review A* 85, 6 (2012). ISSN: 1094-1622.
- [104] Gell-Mann, M. and Hartle, J. B. "Adaptive coarse graining, environment, strong decoherence, and quasiclassical realms". *Physical Review A* 89, 5 (2014).
- [105] Hutter, M. "Feature Reinforcement Learning: Part I. Unstructured MDPs". *Journal of Artificial General Intelligence* 1, 1 (Jan. 2009).
- [106] Nguyen, P., Sunehag, P., and Hutter, M. "Feature Reinforcement Learning in Practice". *Recent Advances in Reinforcement Learning*. Ed. by Sanner, S. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012, pp. 66–77. ISBN: 978-3642299469.
- [107] Hutter, M. "Extreme state aggregation beyond Markov decision processes". *Theoretical Computer Science* 650 (Oct. 2016), pp. 73–91.
- [108] Daswani, M., Sunehag, P., and Hutter, M. "Q-learning for history-based reinforcement learning". *Proceedings of the 5th Asian Conference on Machine Learning*. Ed. by Ong, C. S. and Ho, T. B. Vol. 29. Proceedings of Machine Learning Research. Australian National University, Canberra, Australia: PMLR, 2013, pp. 213–228.
- [109] Majeed, S. J. and Hutter, M. "On Q-learning Convergence for Non-Markov Decision Processes". *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI-18*. International Joint Conferences on Artificial Intelligence Organization, July 2018, pp. 2546–2552.
- [110] Pendrith, M. D. and McGarity, M. "An Analysis of non-Markov Automata Games: Implications for Reinforcement Learning". 1997.
- [111] Pendrith, M. D. and McGarity, M. "An Analysis of Direct Reinforcement Learning in Non-Markovian Domains". *Proceedings of the Fifteenth International Conference on Machine Learning*. ICML '98. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 1998, pp. 421–429. ISBN: 1558605568.
- [112] Perkins, T. J. and Pendrith, M. D. "On the Existence of Fixed Points for Q-Learning and Sarsa in Partially Observable Domains". *Proceedings of the Nineteenth International Conference on Machine Learning*. ICML '02. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 2002, pp. 490–497. ISBN: 1558608737.
- [113] Li, L., Walsh, T., and Littman, M. "Towards a Unified Theory of State Abstraction for MDPs". *ISAIM*. 2006.
- [114] Leike, J. "Nonparametric General Reinforcement Learning". PhD thesis. Australian National University, 2016.
- [115] Azizzadenesheli, K. "Reinforcement Learning in Structured and Partially Observable Environments". PhD thesis. University of California, Irvine, 2019.

- [116] Chandak, Y. et al. “Optimizing for the Future in Non-Stationary MDPs”. *Proceedings of the 37th International Conference on Machine Learning*. Ed. by III, H. D. and Singh, A. Vol. 119. Proceedings of Machine Learning Research. PMLR, 2020, pp. 1414–1425.
- [117] Chandak, Y. et al. “Towards Safe Policy Improvement for Non-Stationary MDPs”. *Advances in Neural Information Processing Systems*. Ed. by Larochelle, H. et al. Vol. 33. Curran Associates, Inc., 2020, pp. 9156–9168.
- [118] Thrun, S. and Schwartz, A. “Issues in Using Function Approximation for Reinforcement Learning”. *Proceedings of the 1993 Connectionist Models Summer School*. Ed. by Mozer, M. et al. Erlbaum Associates, 1993.
- [119] Bellman, R. E. “On the Theory of Dynamic Programming”. *Proceedings of the National Academy of Sciences* 38, 8 (Aug. 1952), pp. 716–719.
- [120] Bellman, R. E. *Dynamic Programming*. Jan. 2003. ISBN: 978-0486428093.
- [121] Fazekas, I. and Klesov, O. “A General Approach to the Strong Law of Large Numbers”. *Theory of Probability & Its Applications* 45, 3 (Jan. 2001), pp. 436–449.
- [122] Boyd, S. and Vandenberghe, L. *Convex Optimization*. Cambridge University Press, 2004. ISBN: 978-0521833783.
- [123] Samson, C. et al. “A variational model for image classification and restoration”. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 22, 5 (May 2000), pp. 460–472.
- [124] Hamza, A. and Brady, D. “Reconstruction of reflectance spectra using robust nonnegative matrix factorization”. *IEEE Transactions on Signal Processing* 54, 9 (Sept. 2006), pp. 3637–3642.
- [125] Liu, W., Pokharel, P. P., and Principe, J. C. “Correntropy: Properties and Applications in Non-Gaussian Signal Processing”. *IEEE Transactions on Signal Processing* 55, 11 (Nov. 2007), pp. 5286–5298.
- [126] Nagy, F. “Parameter Estimation of the Cauchy Distribution in Information Theory Approach”. *Journal of Universal Computer Science* 12, 9 (Sept. 28, 2006), pp. 1332–1344.
- [127] Pokharel, R. and Principe, J. C. “Kernel classifier with Correntropy loss”. *The 2012 International Joint Conference on Neural Networks (IJCNN)*. IEEE, June 2012.
- [128] Du, L., Li, X., and Shen, Y.-D. “Robust Nonnegative Matrix Factorization via Half-Quadratic Minimization”. *2012 IEEE 12th International Conference on Data Mining*. IEEE, Dec. 2012.
- [129] Huber, P. J. *Robust Statistics*. Mar. 1981. ISBN: 978-0471418054.
- [130] Haan, L. de and Ferreira, A. *Extreme Value Theory: An Introduction*. Springer, 2006. ISBN: 978-0387239460.
- [131] Maronna, R. A., Martin, D. R., and Yohai, V. J. *Robust Statistics*. May 2006. ISBN: 978-0470010921.
- [132] Falk, M., Husler, J., and Reiss, R.-D. *Laws of Small Numbers: Extremes and Rare Events*. Third Edition. Oct. 2010. ISBN: 978-3034800082.
- [133] Embrechts, P., Kluppelberg, C., and Mikosch, T. *Modelling Extremal Events*. Jan. 2013. ISBN: 978-3540609315.
- [134] Abramowitz, M. and Stegun, I. A. *Handbook of Mathematical Functions: with Formulas, Graphs, and Mathematical Tables*. Dover Publications, 2012.
- [135] Jameson, G. J. O. “The incomplete gamma functions”. *The Mathematical Gazette* 100, 548 (June 2016), pp. 298–306.

- [136] Hill, B. M. “A Simple General Approach to Inference About the Tail of a Distribution”. *The Annals of Statistics* 3, 5 (Sept. 1975).
- [137] Dekkers, A. L. M., Einmahl, J. H. J., and Haan, L. D. “A Moment Estimator for the Index of an Extreme-Value Distribution”. *The Annals of Statistics* 17, 4 (Dec. 1989).
- [138] Christopeit, N. “Estimating parameters of an extreme value distribution by the method of moments”. *Journal of Statistical Planning and Inference* 41, 2 (Sept. 1994), pp. 173–186.
- [139] Beirlant, J., Vynckier, P., and Teugels, J. L. “Excess Functions and Estimation of the Extreme-Value Index”. *Bernoulli* 2, 4 (Dec. 1996), p. 293.
- [140] Beirlant, J., Vynckier, P., and Teugels, J. L. “Tail Index Estimation, Pareto Quantile Plots, and Regression Diagnostics”. *Journal of the American Statistical Association* 91, 436 (Dec. 1996), p. 1659.
- [141] Beirlant, J. et al. “Tail Index Estimation and an Exponential Regression Model”. *Extremes* 2, 2 (1999), pp. 177–200.
- [142] Beirlant, J., Dierckx, G., and Guillou, A. “Estimation of the extreme-value index and generalized quantile plots”. *Bernoulli* 11, 6 (Dec. 2005).
- [143] Jaakkola, T., Jordan, M. I., and Singh, S. P. “On the Convergence of Stochastic Iterative Dynamic Programming Algorithms”. *Neural Computation* 6, 6 (Nov. 1994), pp. 1185–1201.
- [144] Ribeiro, C. H. C. and Szepesvári, C. “Q-learning Combined with Spreading: Convergence and Results”. *Proceedings of ISRF-IEE International Conference: Intelligent and Cognitive Systems, Neural Networks Symposium*. Tehran, Iran, 1996, pp. 32–36.
- [145] Singh, S. et al. “Convergence Results for Single-Step On-Policy Reinforcement-Learning Algorithms”. *Machine Learning* 38, 3 (2000), pp. 287–308.
- [146] Duan, Y. et al. “Benchmarking Deep Reinforcement Learning for Continuous Control”. *Proceedings of The 33rd International Conference on Machine Learning*. Ed. by Balcan, M. F. and Weinberger, K. Q. Vol. 48. Proceedings of Machine Learning Research. New York, New York, USA: PMLR, 2016, pp. 1329–1338.
- [147] Henderson, P. et al. *Deep Reinforcement Learning that Matters*. 2019. arXiv: [1709.06560 \[cs.LG\]](https://arxiv.org/abs/1709.06560).
- [148] Birkhoff, G. D. “Proof of the Ergodic Theorem”. *Proceedings of the National Academy of Sciences* 17, 12 (Dec. 1931), pp. 656–660.
- [149] Birkhoff, G. D. “What is the Ergodic Theorem?” *The American Mathematical Monthly* 49, 4 (Apr. 1942), pp. 222–226.
- [150] Robbins, H. and Monro, S. “A Stochastic Approximation Method”. *The Annals of Mathematical Statistics* 22, 3 (Sept. 1951), pp. 400–407.
- [151] Blum, J. R. “Approximation Methods which Converge with Probability one”. *The Annals of Mathematical Statistics* 25, 2 (June 1954), pp. 382–386.
- [152] Popoviciu, T. “Sur certaines inégalités qui caractérisent les fonctions convexes”. *Sectia I a Mat* 11 (1965), pp. 155–164.

- [153] Sharma, R., Gupta, M., and Kapoor, G. “Some better bounds on the variance with applications”. *Journal of Mathematical Inequalities* 3 (2010), pp. 355–363.
- [154] Bhatia, R. and Davis, C. “A Better Bound on the Variance”. *The American Mathematical Monthly* 107, 4 (Apr. 2000), pp. 353–357.
- [155] Raffin, A., Kober, J., and Stulp, F. *Smooth Exploration for Robotic Reinforcement Learning*. 2021. arXiv: [2005.05719 \[cs.LG\]](#).
- [156] Bengio, Y., Courville, A., and Vincent, P. “Representation Learning: A Review and New Perspectives”. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 35, 8 (Aug. 2013), pp. 1798–1828.
- [157] Schmidhuber, J. “Deep learning in neural networks: An overview”. *Neural Networks* 61 (Jan. 2015), pp. 85–117.
- [158] LeCun, Y., Bengio, Y., and Hinton, G. “Deep learning”. *Nature* 521, 7553 (May 2015), pp. 436–444.
- [159] Goodfellow, I., Bengio, Y., and Courville, A. *Deep Learning*. Cambridge, Massachusetts: The MIT Press, 2016. ISBN: 978-0262035613.
- [160] García-Martín, E. et al. “How to Measure Energy Consumption in Machine Learning Algorithms”. *ECML PKDD 2018 Workshops*. Springer International Publishing, 2019, pp. 243–255.
- [161] García-Martín, E. et al. “Estimation of energy consumption in machine learning”. *Journal of Parallel and Distributed Computing* 134 (Dec. 2019), pp. 75–88.
- [162] Han, S. et al. “Learning Both Weights and Connections for Efficient Neural Networks”. *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 1*. Ed. by Cortes, C. et al. Vol. 28. NIPS’15. Montreal, Canada: MIT Press, 2015, pp. 1135–1143.
- [163] Rouhani, B. D., Mirhoseini, A., and Koushanfar, F. “DeLight: Adding Energy Dimension To Deep Neural Networks”. *Proceedings of the 2016 International Symposium on Low Power Electronics and Design*. ACM, Aug. 2016.
- [164] Yang, T.-J., Chen, Y.-H., and Sze, V. “Designing Energy-Efficient Convolutional Neural Networks Using Energy-Aware Pruning”. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, July 2017.
- [165] Cai, E. et al. *NeuralPower: Predict and Deploy Energy-Efficient Convolutional Neural Networks*. 2017. arXiv: [1710.05420 \[cs.LG\]](#).
- [166] Rodrigues, C., Riley, G., and Luján, M. “SyNERGY: An energy measurement and prediction framework for Convolutional Neural Networks on Jetson TX1”. English. *PDPTA’18 - The 24th International Conference on Parallel and Distributed Processing Techniques and Applications*. June 2018. ISBN: 1-60132-487-1.
- [167] Han, S., Mao, H., and Dally, W. J. *Deep Compression: Compressing Deep Neural Networks with Pruning, Trained Quantization and Huffman Coding*. 2016. arXiv: [1510.00149 \[cs.CV\]](#).
- [168] Iandola, F. N. et al. *SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and <0.5MB model size*. 2016. arXiv: [1602.07360 \[cs.CV\]](#).
- [169] Kingma, D. P. and Ba, J. *Adam: A Method for Stochastic Optimization*. 2017. arXiv: [1412.6980 \[cs.LG\]](#).
- [170] Loshchilov, I. and Hutter, F. *Decoupled Weight Decay Regularization*. 2019. arXiv: [1711.05101 \[cs.LG\]](#).

- [171] Fukushima, K. “Neocognitron: A self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position”. *Biological Cybernetics* 36, 4 (Apr. 1980), pp. 193–202.
- [172] Rumelhart, D. E., Hinton, G. E., and McClelland, J. L. “A General Framework for Parallel Distributed Processing”. *Parallel Distributed Processing: Explorations in the Microstructure of Cognition, Vol. 1: Foundations*. Cambridge, MA, USA: MIT Press, 1986, pp. 45–76. ISBN: 026268053X.
- [173] Malik, J. and Perona, P. “Preattentive texture discrimination with early vision mechanisms”. *Journal of the Optical Society of America A* 7, 5 (May 1990), p. 923.
- [174] Nair, V. and Hinton, G. E. “Rectified Linear Units Improve Restricted Boltzmann Machines”. *Proceedings of the 27th International Conference on International Conference on Machine Learning*. ICML’10. Haifa, Israel: Omnipress, 2010, pp. 807–814. ISBN: 9781605589077.
- [175] Glorot, X., Bordes, A., and Bengio, Y. “Deep Sparse Rectifier Neural Networks”. *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*. Ed. by Gordon, G., Dunson, D., and Dudík, M. Vol. 15. Proceedings of Machine Learning Research. Fort Lauderdale, FL, USA: PMLR, 2011, pp. 315–323.
- [176] Maas, A. L., Hannun, A. Y., and Ng, A. Y. “Rectifier nonlinearities improve neural network acoustic models”. *ICML Workshop on Deep Learning for Audio, Speech and Language Processing*. 2013.
- [177] Sun, Y., Wang, X., and Tang, X. *Deeply learned face representations are sparse, selective, and robust*. 2014. arXiv: [1412.1265 \[cs.CV\]](https://arxiv.org/abs/1412.1265).
- [178] Mitchell, M. *Why AI is Harder Than We Think*. 2021. arXiv: [2104.12871 \[cs.AI\]](https://arxiv.org/abs/2104.12871).

A Agent Algorithms: SAC and TD3

Algorithm 1 Soft Actor-Critic (SAC) [77]

```

1: procedure SAC( $\theta_1, \theta_2, \phi, \tau, d, \gamma, \mathcal{S}, \mathcal{D}, m, L$ )
2:   Initialise soft critic networks  $Q_{\theta_1}, Q_{\theta_2}$  with random parameters  $\theta_1, \theta_2$      $\triangleright$  Twin value function approximators
3:   Initialise actor network  $\pi_\phi$  with random parameter  $\phi$                                  $\triangleright$  Policy density function approximation
4:   Initialise target networks  $\bar{\theta}_i \leftarrow \theta_i$ 
5:   Initialise empty replay buffer of fixed size  $\mathcal{D} \leftarrow$ 
6:   for each environment step  $t = 1$  to  $T$  do
7:     Select action  $a_t \sim \pi_\phi(\cdot | s_t)$                                       $\triangleright$  Symmetric zero-mean random policy distribution  $\mathcal{S}$ 
8:     Observe reward  $r_t$ , next state  $s_{t+1}$ , and done flag  $e_t$ 
9:     Store  $\mathcal{D} \leftarrow \mathcal{D} \cup \{s_t, a_t, r_t, s_{t+1}, e_t\}$                           $\triangleright$  Add current transition to buffer
10:    for each gradient step do
11:      Uniformly sample mini-batch of  $N$  transitions  $\{s_k, a_k, R_k^{(m)}, s_{k+m}, e_k\}$  from  $\mathcal{D}$ 
12:      Select next actions  $\bar{a}_{k+m} \sim \pi_\phi(\cdot | s_{k+m})$ 
13:      Evaluate target soft Q-values  $Q_{\bar{\theta}_1}(s_{k+m}, \bar{a}_{k+m}), Q_{\bar{\theta}_2}(s_{k+m}, \bar{a}_{k+m})$ 
14:      Take minimum  $Q_{\bar{\theta}}(s_{k+m}, \bar{a}_{k+m}) = \min(Q_{\bar{\theta}_1}(s_{k+m}, \bar{a}_{k+m}), Q_{\bar{\theta}_2}(s_{k+m}, \bar{a}_{k+m}))$      $\triangleright$  Double-Q clipping
15:      Construct target soft values  $V_{\bar{\theta}}(s_{k+m}) = Q_{\bar{\theta}}(s_{k+m}, \bar{a}_{k+m}) - \alpha \log \pi_\phi(\bar{a}_{k+m} | s_{k+m})$ 
16:      Evaluate target critic Q-value  $Q_{\bar{\theta}}(s_{k+m}, \bar{a}_{k+m}) \leftarrow R_k^{(m)} + \gamma^m V_{\bar{\theta}}(s_{k+m})$ 
17:      Evaluate twin critic Q-values  $Q_{\theta_1}(s_k, a_k), Q_{\theta_2}(s_k, a_k)$ 
18:      Take minimum  $Q_{\theta}(s_k, a_k) = \min(Q_{\theta_1}(s_k, a_k), Q_{\theta_2}(s_k, a_k))$                                  $\triangleright$  Double-Q clipping
19:      Construct order-statistic sets  $\omega_i$  for mini-batch  $L(Q_{\bar{\theta}}(s_{k+m}, \bar{a}_{k+m}), Q_{\theta_i}(s_k, a_k))$      $\triangleright$  Loss functions  $L$ 
20:      Using  $\omega_i$  obtain tail exponent estimate  $\hat{\alpha}_i$                                           $\triangleright$  Zipf plot gradient
21:      Construct empirical critic objectives  $J(\theta_i) = \mathbb{E}_{U(\mathcal{D})} [\omega_i]$ 
22:      Construct shadow critic objectives  $J_s(\theta_i) = \mu_s(L_i^*, H_i, \hat{\alpha}_i)$ 
23:      Backpropagate  $\theta_{1,2} = \arg \min_{\theta_{1,2}} (J(\theta_1) + J(\theta_2))$                        $\triangleright$  Update both critics simultaneously
24:      if  $t \bmod d$  then                                                  $\triangleright$  Delayed actor update interval d
25:        Construct policy loss  $J(\phi) = \mathbb{E}_{U(\mathcal{D})} [\alpha \log \pi_\phi(a_k | s_k) - Q_{\theta}(s_k, a_k)]$ 
26:        Backpropagate  $\phi = \arg \min_{\phi} J(\phi)$ 
27:        Construct temperature loss  $J(\alpha) = \mathbb{E}_{U(\mathcal{D})} [-\alpha (\log \pi_\phi(a_k | s_k) + \bar{H})]$ 
28:        Dual gradient descent  $\alpha = \arg \min_{\alpha} J(\alpha)$ 
29:        Update  $\bar{\theta}_i \leftarrow \tau \theta_i + (1 - \tau) \bar{\theta}_i$                                       $\triangleright$  Polyak update of target critics
30:      end if
31:    end for
32:  end for

```

Algorithm 2 Twin Delayed Deep Deterministic Policy Gradient (TD3) [71]

```

1: procedure TD3( $\theta_1, \theta_2, \phi, \tau, \sigma, \bar{\sigma}, c, d, \gamma, \mathcal{S}, \mathcal{D}, m, L$ )
2:   Initialise critic networks  $Q_{\theta_1}, Q_{\theta_2}$  with random parameters  $\theta_1, \theta_2$             $\triangleright$  Twin value function approximators
3:   Initialise actor network  $\pi_\phi$  with random parameter  $\phi$                                  $\triangleright$  Policy density function approximation
4:   Initialise target networks  $\bar{\theta}_i \leftarrow \theta_i, \bar{\phi} \leftarrow \phi$ 
5:   Initialise empty replay buffer of fixed size  $\mathcal{D} \leftarrow$ 
6:   for each environment step  $t = 1$  to  $T$  do
7:     Generate exploration noise  $\epsilon \sim \mathcal{S}(0, \sigma)$                                  $\triangleright$  Symmetric zero-mean random distribution  $\mathcal{S}$ 
8:     Select action  $a_t \sim \pi_\phi(s_t) + \epsilon$                                           $\triangleright$  Manually inject exploration noise
9:     Observe reward  $r_t$ , next state  $s_{t+1}$ , and done flag  $e_t$ 
10:    Store  $\mathcal{D} \leftarrow \mathcal{D} \cup \{s_t, a_t, r_t, s_{t+1}, e_t\}$                                 $\triangleright$  Add current transition to buffer
11:    for each gradient step do
12:      Uniformly sample mini-batch of  $N$  transitions  $\{s_k, a_k, R_k^{(m)}, s_{k+m}, e_k\}$  from  $\mathcal{D}$ 
13:      Generate target exploration noise  $\bar{\epsilon} \sim \mathcal{S}(0, \bar{\sigma})$                        $\triangleright$  Target policy smoothing
14:      Clip the noise  $\bar{\epsilon} \leftarrow \text{clip}(\bar{\epsilon}, -c, c)$ 
15:      Select next action  $\bar{a}_{k+m} \sim \pi_{\bar{\phi}}(s_{k+m}) + \bar{\epsilon}$ 
16:      Obtain target Q-values  $Q_{\bar{\theta}_1}(s_{k+m}, \bar{a}_{k+m}), Q_{\bar{\theta}_2}(s_{k+m}, \bar{a}_{k+m})$ 
17:      Take minimum  $Q_{\bar{\theta}}(s_{k+m}, \bar{a}_{k+m}) = \min(Q_{\bar{\theta}_1}(s_{k+m}, \bar{a}_{k+m}), Q_{\bar{\theta}_2}(s_{k+m}, \bar{a}_{k+m}))$      $\triangleright$  Double-Q clipping
18:      Evaluate target critic Q-value  $Q_{\bar{\theta}}(s_{k+m}, \bar{a}_{k+m}) \leftarrow R_k^{(m)} + \gamma^m Q_{\bar{\theta}}(s_{k+m}, \bar{a}_{k+m})$ 
19:      Evaluate twin critic Q-values  $Q_{\theta_1}(s_k, a_k), Q_{\theta_2}(s_k, a_k)$ 
20:      Construct order-statistic sets  $\omega_i$  for mini-batch  $L(Q_{\bar{\theta}}(s_{k+m}, \bar{a}_{k+m}), Q_{\theta_i}(s_k, a_k))$      $\triangleright$  Loss functions  $L$ 
21:      Using  $\omega_i$  obtain tail exponent estimate  $\hat{\alpha}_i$                                       $\triangleright$  Zipf plot gradient
22:      Construct empirical critic objectives  $J(\theta_i) = \mathbb{E}_{U(\mathcal{D})}[\omega_i]$ 
23:      Construct shadow critic objectives  $J_s(\theta_i) = \mu_s(L_i^*, H_i, \hat{\alpha}_i)$ 
24:      Backpropagate  $\theta_{1,2} = \arg \min_{\theta_{1,2}} (J(\theta_1) + J(\theta_2))$                    $\triangleright$  Update both critics simultaneously
25:      if  $t \bmod d$  then                                                                $\triangleright$  Delayed actor update interval d
26:        Construct policy loss  $J(\phi) = -\mathbb{E}_{U(\mathcal{D})}[Q_{\theta_1}(s, a)]$ 
27:        Backpropagate  $\phi = \arg \min_{\phi} J(\phi)$ 
28:        Update  $\bar{\theta}_i \leftarrow \tau \theta_i + (1 - \tau) \bar{\theta}_i$                                  $\triangleright$  Polyak update of target critics
29:        Update  $\bar{\phi} \leftarrow \tau \phi + (1 - \tau) \bar{\phi}$                                  $\triangleright$  Polyak update of target actor
30:      end if
31:    end for
32:  end for

```

Parameters	TD3	SAC
Optimiser	Adam ^[169, 170]	Adam
Learning Rates (θ, ϕ, α)	$1 \cdot 10^{-3}$	$3 \cdot 10^{-4}$
Layer 1 Hidden Nodes (θ, ϕ)	400	256
Layer 2 Hidden Nodes (θ, ϕ)	300	256
Non-linearity (θ, ϕ)	ReLU ^[171–177]	ReLU
Mini-Batch Size (N)	100	256
Gradient Iterations per Time Step	1	1
Policy Gaussian Noise (σ)	0.1	
Target Policy Gaussian Noise ($\bar{\sigma}$)	0.2	
Target Policy Noise Clip Rate (c)	0.5	
Minimum Expected Entropy Target (\bar{H})		$-\dim(\mathcal{A})$
Target Update Smoothing Rate (τ)	$5 \cdot 10^{-3}$	$5 \cdot 10^{-3}$
Actor Network Update Interval (d)	2	1
Target Network Update Interval	2	1
Discount Factor (γ)	0.99	0.99
Replay Buffer Size (\mathcal{D})	$1 \cdot 10^6$	$1 \cdot 10^6$

Table 1: Common hyperparameters for TD3 [71] and SAC [77] including deep linear neural network architectures.

B Assorted Applications

There are probably only a high double-digit number of people in the world that truly understand the reward schemes of the locomotive environments (Hopper, Walker2D, Humanoid etc.) we have tested throughout this project. High scores in these environments are completely meaningless to the outside world. Furthermore, the implications of the mathematics derivations are not obvious to anyone not very well-versed in reinforcement learning.

Two distinguishing features of our selected agent algorithms is that they are both model-free and off-policy. Model-free algorithms involve minimal assumptions regarding the environments and therefore are a general approach to any problem. The agent initially knows nothing about the world in which it is operating, what it can control, and how the rewards are tied its actions. Off-policy means that when the agent is performing a task, it recalls many previous times it performed the task, and then bases its next current action by determining what the optimal past actions might have been.

Agent learning can be described to be analogous to how a toddler learns to walk. We do not provide them a (complete or partial) instruction set on controlling their tiny musculoskeletal system, rather, they learn through experience consisting of a fine balance between failure and praise. When it comes to quantifying the scale of failure, a common highly penalising (MSE) method is used in existing models. Recent advancements in other fields have highlighted the effectiveness of relatively softer approaches discussed in Section 3.5, applied here, these might lead to superior toddler learning. The overwhelming majority of existing models also would critique the toddlers walking at every single step they attempt, instead, letting them take multiple steps before their evaluating performance might lead to more stable learning as seen in Section 3.7. The multiplicative aspect can be described as ensuring that when the toddler falls, the magnitude of the fall is not large enough to make recovery questionable. Existing reinforcement learning models do not internally provide any such safety feature, see Sections 1, 5-7. Including this feature requires considerable additional resources for constant supervision, but since toddlers are not fungible, it is the only valid approach. Finally, unlike toddlers that operate on a fixed universal time scale, agent training can be simulated and accelerated using supercomputers.

B.1 Overview

The value of our original approach is largely based around tasks where the amount of risk-taking (leverage) can either be amplified or reduced to magnify or shrink potential rewards. The systems are also non-ergodic essentially meaning that the volatility tax $\nu \neq 0$ and so bouncing back from failure is tougher if not impossible. Domains where $\bar{g} \neq \mathbb{E}[r]$ are defined as non-ergodic and are generally the mark of strong multiplicative dynamics. These are situations where time order matters and losses have an asymmetrically larger effect on total return as discussed in Section 1.

The goal is to find the optimal balance of risk-taking using model-free algorithms that can be generally applied to any environment, so we do not have to teach any specific concept or theory to the agent. In contrast, training of assembly line robots with this approach is inappropriate since they are producing identical products and hence the existing additive dynamics approach is perfectly fine.

As a rule of thumb, if we want to optimise the results for a single random individual or operation, multiplicative dynamics should be used as compounding is key. If the goal is maximising the performance a collective group performing the operations, then additive dynamics is fine since averaged results converge to the “expected” value with virtually no one in the group experiencing this precise value. Problem with the latter is that it has sampling bias

since a small number of constituents can dominate the sample average. Therefore, we should be far more interested in determining what the average performance over time for an any random constituent as this is more informative.

Below we provide several examples where the novel components of this project can be easily applied if:

- (i) The input states form a QDP as discussed in Sections 3.8 and 5.
- (ii) The input states can be cleanly sourced and formatted to be compatible with the OpenAI gym [54] framework.
- (iii) The input states can be represented in a format capable of being fed to deep (linear) neural networks.
- (iv) The input states permit random but sensible simulation.
- (v) There exists an accurate reward signal, or once can be designed.

Note all these descriptions are works in progress and only serve to illustrate shortcomings of existing algorithms and how our approach might be superior. We focus solely on optimal decision-making rather than discussing computer vision and object recognition details.

B.2 Robotic Control for Medical Surgery

Performing non-trivial surgery on someone is a highly non-ergodic process as people's composition, medical history, and needs are unique. For robots to perform this task, first it is a very difficult computer vision problem, assuming this has been solved how would the actions a robot takes during surgery be evaluated? Multiplicative dynamics is necessary as this is certainly a system where mistakes are far more determinantal in relative terms compared to successes.

If one serious injury that would never have occurred by a human surgeon occurs, it would likely be sufficient to getting the robot pulled from supermarket shelves, R&D costs never recouped, and huge and possibly irreversible reputational damage. An upside would be watching to see how the CEO tries to explain that their Doctor Robot was designed and trained to (incorrectly?) operate under the assumption that mistakes do not matter as long as the majority of surgeries are a success. This purely probabilistic approach is unacceptable from a sales perspective as the scale of the mistake (tail event) might greatly exceed the benefit of the procedure.

In terms of reinforcement learning, consider each surgery occurring at a time step and the reward from each step being proportional to how successful the surgery was by some measure. Assuming we know all states, this can be modeled as an MDP consisting of one episode. Existing algorithms would be trained to take the probability approach where the total reward would be the sum of each surgery treated independently. This is the wrong to for the agent to perform surgeries.

The reward should be compounding over correct procedures, and highly penalising for mistakes. This way, to work back up to the total reward prior to the mistake, the agent will have to achieve more relative success. The agent selected to perform surgeries should be the one that is trained to perform well while avoiding large mistakes rather than one that blindly maximises performance. Furthermore, the impact of a mistake depends heavily on time.

For example, assuming the agent remains unchanged over time, suppose the first surgery Doctor Robot gets wrong with the same mistake is #13 compared to #99,323. Under existing models, the agent considers these of equal significance. In reality however, there is a world of a difference to long-term profits from sales of Doctor Robot between the two. The former is far more likely to lead to complete failure than the latter.

This is the exact issue we are trying to tackle in this project. We are trying to encode the importance of time when it comes to decision making. The best agent in this approach is one that takes a more conservative approach to patient health. As far as we are aware, very few people (if any) are attempting to design algorithms that explicitly consider this asymmetry (non-ergodicity) which is beyond crucial in everyday life.

B.3 Supply Chain Management

Management of postal delivery services very quickly becomes an unbelievable complicated problem when trying to derive optimal distribution chains, see the NP-hard Traveling Salesman Problem (TSP). Humans have not been able to solve this and have instead resorted to approximate solutions like most things. Agents using deep neural networks might be able to yield even superior results.

Consider the case of single Sydney-like city under the assumption all postal mail and packages are traceable and there is a given finite number of fungible delivery mechanisms (postmen). The system is a POMDP but we assume QDP holds and so MDP modeling is acceptable. As it stands, existing algorithms using additive dynamics would be fine. To make the situation realistic, introduce the notion of express post where there is guarantee of delivery within a known time frame and there are severe financial penalties for late submissions.

The goal is then to design a agent that effectively allocates collections of items to postmen, and then specifies each of them a delivery trajectory for each day (episode). The priority of the agent is to maximise reward under the constraint express post parcels are a priority. This is a highly non-ergodic system as the long-run effect delay can have damaging repercussions to the post offices reputation and future sales since express post is higher priced. As before, there is a huge difference if the first express parcel the agent misses the delivery deadline on is #13, compared to if it was #99,323. See the previous application.

B.4 Guidance Systems

A highly topical area of discussion is self-driving cars. This is a computer vision problem to accurately map the world through sensors, and a decision problem for the next action to take. The decision making can be made using multi-task learning or reinforcement learning. In this environment agent training should be conservative and there should be no continuing from a failed state such as driving off a cliff. The area of concern is then the performance of a particular agent in the infinite time limit, rather than the final average performance across infinite fixed time trials. Naturally this is again a multiplicative problem, generally with multi-modal solutions and hence the stochastic SAC algorithm is applicable while the deterministic TD3 is limited in its utility. There is a world of a difference in the long-term profits of selling self-driving software if the first accident happens on sale #13, compared to if it was on sale #99,323. For the latter, the seller may have already reached critical mass and be able to downplay the negativity of the crash as being a ‘one-off’, whereas the former will be far more difficult to justify.

We can also create state-of-the-art autonomous AI-guidance systems for projectiles that operate at extremely fast speeds. Combined with efficient action segregation using SARA^E², if sufficient computational resources are provided, the trained agents will be totally dominating in their ability to reach destinations with the act of developing effective countermeasures being an exponentially more challenging problem — catching is more difficult than throwing.

B.5 Education

B.6 Portfolio Management

In Section 1.4 we highlighted how financial markets in the real world are much more complicated than the simple coin flipping gamble. Let us now attempt to generalise multiplicative dynamics combined with reinforcement learning to this setting where we will find it is not nearly as complex as it seems. This description is neither exhaustive nor a complete recipe for shekel accumulation. Rather, it serves to connect our work to real-world implications. In reality, implementing this would involve piecemeal modular additions before a complete end-to-end prescription. We therefore theorise what the optimal hypothetical portfolios might be.

Consider an arbitrary financial portfolio with initial value \$100 and a known fixed number of tradable securities (must be fixed for now). The number of action components $\dim(\mathcal{A})$ would equal the number of tradeable securities. Each component has maximum and minimum values equal to the leveraged position limits in those securities. For example, continuous action space limits $[-5, 8]$ imply we can take up to a maximum leveraged 8x long and -5x short position in a product. The constraint of total funds available can also easily be satisfied if the position is thought of as proportions of funds to be allocated.

The input states would consist of at minimum all these tradeable securities prices plus extra features. The extra features would include all other variables we consider important such as: additional security prices, index values, NLP data, memes, and economic data. This quickly leads to a POMDP, but since we assume the QDP assumption holds, we can simplify the problem to an MDP. This assumption is made throughout finance, see for example geometric Brownian motion in derivatives and options pricing [24]. Regardless, this selection of extra features is going to be key since too many superfluous states add noise and slow learning and too few will lead to underfitting. Some form of extravagant dimension reduction will be needed to represent global markets.

The rewards per time step would simply be the profit and loss of the step given the current overall positioning and the next set of prices. We learn off-policy from a batch of past portfolio histories with total values greater than the threshold of say \$20 set as our stop-loss. Actor policy optimisation is unchanged, its sole purpose is to maximise the total rewards discounted future rewards of each sample in the batch.

Convention would use additive dynamics where the critic would attempt to make sure the dollar differences between current and target Q-values are minimised. Our work on multiplicative dynamics (non-ergodicity) encourages the critic to minimise the change in total portfolio returns between current and target ‘Q-returns’. This is also realistic as there is a reason valuation changes are most often quoted in terms percentage changes rather than dollar changes.

Additionally, our use of different critic loss functions can prevent excessive impact on neural networks parameters (being ‘spooked’) from a small number of very large portfolio return differences across the batch, unlike the standard MSE which exacerbates this occurrence. Multi-step returns can be seen as extending the framework to minimise the difference the current Q-return and the bootstrapped target return. This has shown to extraordinarily effective in many environments with no conclusive theoretical explanation currently available as to why this is the case.

Combining this with non-ergodicity we no longer need to allow for the unrealistic existence of ‘temporarily bankrupt’ intermediate states. This facilitates the construction of fully autonomous, self-learning, risk-reward maximising algorithms that have both robust signal-to-noise detection and far more importantly, explicitly account for the asymmetric effect that large losses have on total compounding portfolio return.

Working backwards, is this method risk-reward maximising? Yes, reward is maximised as per usual by policy

optimisation and risk is correctly expressed as changes in portfolio value with differences across the batch aggregated and minimised using robust loss functions. Furthermore, the use of clipped double Q-learning ensures or discounted future rewards are more conservative. Therefore, the agent learns strictly from histories that are considered relatively ‘successful’ and ignores those where its actions lead to failure. This however increases the need for larger training samples.

Is the algorithm self-learning regarding the environment? Yes, at no point did we input or make assumptions regarding the actual environment other than being a QDP. No mentions of efficient frontier, mean-variance optimisation, risk-free rate of return, investor utility, correlation, security details, any type of DCF, etc. were made. Note we also did not specify whether a particular tradeable security is a position in options, futures or spot. We also never even referred to standard deviation, volatility or even alpha. All the agent knows is that it gets input states (prices and extras) that are all context-less streams of numbers, the agent has a bunch of levers it can pull with each controlling positioning in a tradeable security, and then receives an easily calculable reward signal at every time step. The agent then does its best to maximise this reward over time regardless of what it all means. The deep neural networks being universal function approximators automatically construct optimal variables in a ‘black-box’ fashion.

Finally, is the algorithm fully autonomous? To a certain extent. Assume we can accurately simulate states, select the correct extra features (no easy task), and are able to perform extensive hyperparameter tuning for the details of networks. Some of this can be automated using LTSM/GRU networks. Next, assuming the best version or parameters of the agent algorithm have been found, how would risk-management work in a live environment? Firstly, we are able to set a portfolio stop-loss at for example \$20, the agent can then be train using this threshold. Single position limits can also be set through action space limits. Sector/asset concentration limits or enforcing delta (or other Greek) neutral portfolios cannot be done with this very simplistic approach. These shortcomings can however be addressed through the addition of more constraints. Keeping also in mind that theoretically the optimal agent should learn concentration limits on its own, and will construct a delta neutral portfolio if considered the best possible case. Another important feature is that agent assumes it can make decisions forever with no human interference or overwriting. Hence, actions taken by the agent might seem non-sensical to humans but may turn out to be prophetic.

The assumption fixed number of tradable securities (agent actions) can be relaxed. If adding new products, we can use transfer learning to import existing values for current products and train again which should provide a substantial speed boost. Removal of products would also work similar with exclusion of those parameters and then re-training. Large changes input states however would not be so easy and would most likely require significant re-training. On top of this, adversarial training can also be used against humans to speed up learning.

An uninteresting, but possible outcome might simply be holding a market-capitalisation weighted market portfolio with leverage one. This is justified in prior work on optimal leverage efficiency where any deviations from one are unwise over longer time scales [93]. Furthermore, efficient market hypothesis arguments and contemporary portfolio theory also lend weight to this outcome.

A more novel outcome might be a delta neutral, extremely convex, insurance portfolio constructed using complex net long (mainly put) option payoffs relative to the underlying market portfolio. The basic idea is a portfolio that has struck the perfect balance between incinerating capital (through premium or theta decay), and long enough gamma to astronomically profit from drawdowns. The agent would then govern the rolling over of options contracts where the only downside risk would be a known theta, making the portfolio somewhat ‘risk-free’.

Ultimately, we expect the final result will be a very unequal combination of the above two cases [13–19]. Most

likely, roughly 90-97% will be allocated to a standard market portfolio and the remainder placed in the insurance portfolio. If we were to take the “expected” return of the strategy, it would exceed by many multiples the return of any tradeable risk-free security in existence while maintaining a similar level of risk. This point is worth repeating, by focusing on payoffs and not “expectations”, we have synthetically constructed an essentially risk-free portfolio that outperforms all similarly risky assets (on the efficient frontier), is robust to any crash without requiring prior knowledge, minimal overall rebalancing, and at no point requires “investing” know-how.

The beauty of this approach is that at no point do we have to find ‘alpha’, we simply use non-ergodic theory to find the portfolio that maximises compounding growth through avoidance of steep losses. In the meantime while we wait for inevitable drawdowns, we could create a highly successful neutral market-making business for day-to-day entertainment.

B.7 Systems Control

C Taming Hubris

Addressing concerns in [178]: As we are not aiming to create artificial general intelligence, much of the critiques in this work are less strong. While shortcut learning is defiantly a concern, if we are able to accurately simulate states and train for extended periods of time, we hope this potential weak point is mitigated.

- Fallacy 1: Our goal is not to create human-like general intelligence, rather we aim to create agents with very specific goals when operating in non-ergodic situations, well defined action spaces, and environment states that can be sourced.
- Fallacy 2: The tasks we intend for the agent to perform are not so much measured against humans (yet), instead we compare with existing models to see whether we get better results.
- Fallacy 3: Very valid claim as evaluating the true performance of our agents is not possible. We can only use the benchmarking tools we currently have available, hence any results we achieve should not be assumed easily generalisable to reality.
- Fallacy 4: Agents while considered to be represented by deep neural networks tend to give the impression of a functioning brain. Intelligence is encoded in these network weights and since the action space includes all ways an agent can interact with an environment, we can say that our models multi-sensory when it comes to feedback and response.