

V-trace, PopArt Normalization, Partially Observable MDPs

Milan Straka

🖬 January 7, 2019





EUROPEAN UNION European Structural and Investment Fund Operational Programme Research, Development and Education Charles University in Prague Faculty of Mathematics and Physics Institute of Formal and Applied Linguistics



unless otherwise stated

IMPALA

Ú FÂL

Impala (Importance Weighted Actor-Learner Architecture) was suggested in Feb 2018 paper and allows massively distributed implementation of an actor-critic-like learning algorithm.

Compared to A3C-based agents, which communicates gradients with respect to the parameters of the policy, IMPALA actors communicates trajectories to the centralized learner.



If many actors are used, the policy used to generate a trajectory can lag behind the latest policy. Therefore, a new **V-trace** off-policy actor-critic algorithm is proposed.

IMPALA – V-trace



Consider a trajectory $(S_t, A_t, R_{t+1})_{t=s}^{t=s+n}$ generated by a behaviour policy b. The *n*-step V-trace target for S_s is defined as

$$v_s \stackrel{ ext{def}}{=} V(S_s) + \sum_{t=s}^{s+n-1} \gamma^{t-s} \left(\prod_{i=s}^{t-1} c_i
ight) \delta_t V,$$

where $\delta_t V$ is the temporal difference for V

$$\delta_t V \stackrel{ ext{def}}{=}
ho_t ig(R_{t+1} + \gamma V(s_{t+1}) - V(s_t) ig),$$

and ho_t and c_i are truncated importance sampling ratios with $ar{
ho} \geq ar{c}$:

$$ho_t \stackrel{ ext{\tiny def}}{=} \min\left(ar{
ho}, rac{\pi(A_t|S_t)}{b(A_t|S_t)}
ight), \quad c_i \stackrel{ ext{\tiny def}}{=} \min\left(ar{c}, rac{\pi(A_i|S_i)}{b(A_i|S_i)}
ight)$$

Note that if $b=\pi$ and assuming $ar{c}\geq 1$, v_s reduces to n-step Bellman target.

NPFL122, Lecture 11 IMPALA PopArt Normalization POMDPs MERLIN CTF-FTW

IMPALA – V-trace



Note that the truncated IS weights ρ_t and c_i play different roles:

- The ρ_t appears in the definition of $\delta_t V$ and defines the fixed point of the update rule. For $\bar{\rho} = \infty$, the target is the value function v_{π} , if $\bar{\rho} < \infty$, the fixed point is somewhere between v_{π} and v_b . Notice that we do not compute a product of these ρ_t coefficients.
- The c_i impacts the speed of convergence (the contraction rate of the Bellman operator), not the sought policy. Because a product of the c_i ratios is computed, it plays an important role in variance reduction.

The paper utilizes $ar{c}=1$ and out of $ar{
ho}\in\{1,10,100\}$, $ar{
ho}=1$ works empirically the best.

IMPALA – V-trace



Consider a parametrized functions computing $v(s; \theta)$ and $\pi(a|s; \omega)$. Assuming the defined *n*-step V-trace target

$$v_s \stackrel{ ext{def}}{=} V(S_s) + \sum_{t=s}^{s+n-1} \gamma^{t-s} \left(\prod_{i=s}^{t-1} c_i
ight) \delta_t V,$$

we update the critic in the direction of

$$ig(v_s - v(S_s; oldsymbol{ heta})ig)
abla_{oldsymbol{ heta}} v(S_s; oldsymbol{ heta})$$

and the actor in the direction of the policy gradient

$$ho_s
abla_{oldsymbol{\omega}}\log \pi(A_s|S_s;oldsymbol{\omega})ig(R_{s+1}+\gamma v_{s+1}-v(S_s;oldsymbol{ heta})ig).$$

Finally, we again add the entropy regularization term $H(\pi(\cdot|S_s; \theta))$ to the loss function.

IMPALA

|--|

Architecture	CPUs	GPUs ¹	FPS ²		
Single-Machine			Task 1	Task 2	
A3C 32 workers	64	0	6.5K	9K	
Batched A2C (sync step)	48	0	9K	5K	
Batched A2C (sync step)	48	1	13K	5.5K	
Batched A2C (sync traj.)	48	0	16K	17.5K	
Batched A2C (dyn. batch)	48	1	16K	13K	
IMPALA 48 actors	48	0	17K	20.5K	
IMPALA (dyn. batch) 48 $actors^3$	48	1	21K	24K	
Distributed					
A3C	200	0	46K	50K	
IMPALA	150	1	80K		
IMPALA (optimised)	375	1	200K		
IMPALA (optimised) batch 128	500	1	2	50K	

¹ Nvidia P100 ² In frames/sec (4 times the agent steps due to action repeat). ³ Limited by amount of rendering possible on a single machine.

Table 1 of the paper "IMPALA: Scalable Distributed Deep-RL with Importance Weighted Actor-Learner Architectures" by Lasse Espeholt et al.

IMPALA – Population Based Training

^ÚF_AL

For Atari experiments, population based training with a population of 24 agents is used to adapt entropy regularization, learning rate, RMSProp ε and the global gradient norm clipping threshold.



IMPALA – Population Based Training

^ÚF_AL

For Atari experiments, population based training with a population of 24 agents is used to adapt entropy regularization, learning rate, RMSProp ε and the global gradient norm clipping threshold.

In population based training, several agents are trained in parallel. When an agent is *ready* (after 5000 episodes), then:

- it may be overwritten by parameters and hyperparameters of another agent, if it is sufficiently better (5000 episode mean capped human normalized score returns are 5% better);
- and independently, the hyperparameters may undergo a change (multiplied by either 1.2 or 1/1.2 with 33% chance).

IMPALA





Figure 4 of the paper "IMPALA: Scalable Distributed Deep-RL with Importance Weighted Actor-Learner Architectures" by Lasse Espeholt et al.

NPFL122, Lecture 11 IMPALA PopA

PopArt Normalization

POMDPs MERLIN CTF-FTW

IMPALA – Learning Curves

Figures 5, 6 of the paper "IMPALA: Scalable Distributed Deep-RL with Importance Weighted Actor-Learner Architectures" by Lasse Espeholt et al.

NPFL122, Lecture 11 IMPALA PopArt Normalization POMDPs MERLIN CTF-FTW

Human Normalised Return	Median	Mean
A3C, shallow, experts A3C, deep, experts	54.9% 117.9%	285.9% 503.6%
Reactor, experts	187%	N/A
IMPALA, shallow, experts IMPALA, deep, experts	93.2% 191.8%	466.4% 957.6%
IMPALA, deep, multi-task	59.7%	176.9%

Table 4 of the paper "IMPALA: Scalable Distributed Deep-RL with Importance Weighted Actor-Learner Architectures" by Lasse Espeholt et al.

NPFL122, Lecture 11

IMPALA PopArt Normalization

on POMDPs ME

MERLIN CTF-FTW

 \mathcal{E}

 ε

N

	Task 1	Task 2	Task 3	Task 4	Task 5	
Vithout Replay						
-trace	46.8	32.9	31.3	229.2	43.8	
-Step	51.8	35.9	25.4	215.8	43.7	
-correction	44.2	27.3	4.3	107.7	41.5	
lo-correction	40.3	29.1	5.0	94.9	16.1	
Vith Replay						
-trace	47.1	35.8	34.5	250.8	46.9	
-Step	54.7	34.4	26.4	204.8	41.6	
-correction	30.4	30.2	3.9	101.5	37.6	
lo-correction	35.0	21.1	2.8	85.0	11.2	

Tasks: rooms_watermaze, rooms_keys_doors_puzzle,

lasertag_three_opponents_small,

explore_goal_locations_small, seekavoid_arena_01

Table 2 of the paper "IMPALA: Scalable Distributed Deep-RL with Importance Weighted Actor-Learner Architectures" by Lasse Espeholt et al.

NPFL122, Lecture 11 IMPALA PopArt Normalization POMDPs MERLIN CTF-FTW

IMPALA – Ablations

Figure E.1 of the paper "IMPALA: Scalable Distributed Deep-RL with Importance Weighted Actor-Learner Architectures" by Lasse Espeholt et al.

POMDPs

NPFL122, Lecture 11

PopArt Normalization

IMPALA

MERLIN CTF-FTW

PopArt Normalization

Ú F_AL

An improvement of IMPALA from Sep 2018, which performs normalization of task rewards instead of just reward clipping. PopArt stands for *Preserving Outputs Precisely, while Adaptively Rescaling Targets*.

Assume the value estimate $v(s; \theta, \sigma, \mu)$ is computed using a normalized value predictor $n(s; \theta)$

$$v(s;oldsymbol{ heta},\sigma,\mu) \stackrel{ ext{\tiny def}}{=} \sigma n(s;oldsymbol{ heta}) + \mu$$

and further assume that $n(s; \boldsymbol{\theta})$ is an output of a linear function

$$n(s;oldsymbol{ heta}) \stackrel{ ext{def}}{=} oldsymbol{\omega}^T f(s;oldsymbol{ heta} - \{oldsymbol{\omega},b\}) + b.$$

We can update the σ and μ using exponentially moving average with decay rate β (in the paper, first moment μ and second moment v is tracked, and standard deviation is computed as $\sigma = \sqrt{v - \mu^2}$; decay rate $\beta = 3 \cdot 10^{-4}$ is employed).

PopArt Normalization

Ú_F≩L

Utilizing the parameters μ and σ , we can normalize the observed (unnormalized) returns as $(G - \mu)/\sigma$ and use an actor-critic algorithm with advantage $(G - \mu)/\sigma - n(S; \theta)$.

However, in order to make sure the value function estimate does not change when the normalization parameters change, the parameters ω, b computing the unnormalized value estimate are updated under any change $\mu \to \mu'$ and $\sigma \to \sigma'$ as:

$$oldsymbol{\omega}' \stackrel{ ext{def}}{=} rac{\sigma}{\sigma'} oldsymbol{\omega}, \quad b' \stackrel{ ext{def}}{=} rac{\sigma b + \mu - \mu'}{\sigma'}.$$

In multi-task settings, we train a task-agnostic policy and task-specific value functions (therefore, μ , σ and $n(s; \theta)$ are vectors).

PopArt Results

NPFL122, Lecture 11

IMPALA PopArt Normalization

POMDPs MERLIN CTF-FTW

16/32

PopArt Results

NPFL122, Lecture 11

IMPALA

PopArt Normalization POMDPs

Ps MERLIN CTF-FTW

PopArt Results

Figures 4, 5 of paper "Multi-task Deep Reinforcement Learning with PopArt" by Matteo Hessel et al.

NPFL122, Lecture 11

IMPALA PopArt Normalization

POMDPs MERLIN

CTF-FTW

Partially Observable MDPs

Ú F_ÅL

Recall that a *Markov decision process* (MDP) is a quadruple (S, A, p, γ) , where:

- ${\mathcal S}$ is a set of states,
- \mathcal{A} is a set of actions,
- $p(S_{t+1} = s', R_{t+1} = r | S_t = s, A_t = a)$ is a probability that action $a \in \mathcal{A}$ will lead from state $s \in \mathcal{S}$ to $s' \in \mathcal{S}$, producing a *reward* $r \in \mathbb{R}$,
- $\gamma \in [0,1]$ is a discount factor.

Partially observable Markov decision process extends the Markov decision process to a sextuple (S, A, p, γ, O, o) , where in addition to an MDP

- ${\mathcal O}$ is a set of observations,
- $o(O_t|S_t, A_{t-1})$ is an observation model.

In robotics (out of the domain of this course), several approaches are used to handle POMDPs, to model uncertainty, imprecise mechanisms and inaccurate sensors.

Partially Observable MDPs

ÚFÁ

In Deep RL, partially observable MDPs are usually handled using recurrent networks. After suitable encoding of input observation O_t and previous action A_{t-1} , a RNN (usually LSTM) unit is used to model the current S_t (or its suitable latent representation), which is in turn utilized to produce A_t .

a. RL-LSTM

Figure 1a of paper "Unsupervised Predictive Memory in a Goal-Directed Agent" by Greg Wayne et al.

However, keeping all information in the RNN state is substantially limiting. Therefore, *memory*augmented networks can be used to store suitable information in external memory (in the lines of NTM, DNC or MANN models).

We now describe an approach used by Merlin architecture (Unsupervised Predictive Memory in a Goal-Directed Agent DeepMind Mar 2018 paper).

NPFL122, Lecture 11

IMPALA MERLIN PopArt Normalization POMDPs

CTF-FTW

MERLIN – Memory Module

Let $oldsymbol{M}$ be a memory matrix of size $N_{mem} imes 2|z|$.

Assume we have already encoded observations as e_t and previous action a_{t-1} . We concatenate them with K previously read vectors and process by a deep LSTM (two layers are used in the paper) to compute h_t .

Then, we apply a linear layer to \boldsymbol{h}_t , computing K key vectors $\boldsymbol{k}_1, \ldots, \boldsymbol{k}_K$ of length 2|z| and K positive scalars β_1, \ldots, β_K .

b. RL-MEM

Figure 1b of paper "Unsupervised Predictive Memory in a Goal-Directed Agent" by Greg Wayne et al.

Reading: For each *i*, we compute cosine similarity of k_i and all memory rows M_j , multiply the similarities by β_i and pass them through a softmax to obtain weights ω_i . The read vector is then computed as Mw_i .

Writing: We find one-hot write index \boldsymbol{v}_{wr} to be the least used memory row (we keep usage indicators and add read weights to them). We then compute $\boldsymbol{v}_{ret} \leftarrow \gamma \boldsymbol{v}_{ret} + (1 - \gamma) \boldsymbol{v}_{wr}$, and update the memory matrix using $\boldsymbol{M} \leftarrow \boldsymbol{M} + \boldsymbol{v}_{wr}[\boldsymbol{e}_t, 0] + \boldsymbol{v}_{ret}[0, \boldsymbol{e}_t]$.

MERLIN — Prior and Posterior

However, updating the encoder and memory content purely using RL is inefficient. Therfore, MERLIN includes a *memory-based predictor* (*MBP*) in addition to policy. The goal of MBP is to compress observations into low-dimensional state representations z and storing them in memory.

According to the paper, the idea of unsupervised and predictive modeling has been entertained for decades, and recent discussions have proposed such modeling to be connected to hippocampal memory.

We want the state variables not only to faithfully represent the data, but also emphasise rewarding elements of the environment above irrelevant ones. To accomplish this, the authors follow the hippocampal representation theory of Gluck and Myers, who proposed that hippocampal representations pass through a compressive bottleneck and then reconstruct input stimuli together with task reward.

In MERLIN, a *prior* distribution over z_t predicts next state variable conditioned on history of state variables and actions $p(z_t|z_{t-1}, a_{t-1}, \ldots, z_1, a_1)$, and *posterior* corrects the prior using the new observation o_t , forming a better estimate $q(z_t|o_t, z_{t-1}, a_{t-1}, \ldots, z_1, a_1)$.

MERLIN — Prior and Posterior

To achieve the mentioned goals, we add two terms to the loss.

- We try reconstructing input stimuli, action, reward and return using a sample from the state variable posterior, and add the difference of the reconstruction and ground truth to the loss.
- We also add KL divergence of the prior and posterior to the loss, to ensure consistency between the prior and posterior.

MERLIN — Algorithm

Algorithm 1 MERLIN Worker Pseudocode // Assume global shared parameter vectors θ for the policy network and χ for the memorybased predictor; global shared counter T := 0// Assume thread-specific parameter vectors θ', χ' // Assume discount factor $\gamma \in (0, 1]$ and bootstrapping parameter $\lambda \in [0, 1]$ Initialize thread step counter t := 1repeat Synchronize thread-specific parameters $\theta' := \theta$; $\chi' := \chi$ Zero model's memory & recurrent state if new episode begins $t_{\text{start}} := t$ repeat Prior $\mathcal{N}(\mu_t^p, \log \Sigma_t^p) = p(h_{t-1}, m_{t-1})$ $e_t = \operatorname{enc}(o_t)$ Posterior $\mathcal{N}(\mu_t^q, \log \Sigma_t^q) = q(e_t, h_{t-1}, m_{t-1}, \mu_t^p, \log \Sigma_t^p)$ Sample $z_t \sim \mathcal{N}(\mu_t^q, \log \Sigma_t^q)$ Policy network update $\tilde{h}_t = \text{rec}(\tilde{h}_{t-1}, \tilde{m}_t, \text{StopGradient}(z_t))$ Policy distribution $\pi_t = \pi(\tilde{h}_t, \text{StopGradient}(z_t))$ Sample $a_t \sim \pi_t$ $h_t = \operatorname{rec}(h_{t-1}, m_t, z_t)$ Update memory with z_t by Methods Eq. 2 $R_t, o_t^r = \operatorname{dec}(z_t, \pi_t, a_t)$ Apply a_t to environment and receive reward r_t and observation o_{t+1} t := t + 1; T := T + 1**until** environment termination or $t - t_{\text{start}} = \tau_{\text{window}}$ If not terminated, run additional step to compute $V_{\nu}^{\pi}(z_{t+1}, \log \pi_{t+1})$ and set $R_{t+1} := V^{\pi}(z_{t+1}, \log \pi_{t+1}) //$ (but don't increment counters) Reset performance accumulators $\mathcal{A} := 0$; $\mathcal{L} := 0$; $\mathcal{H} := 0$ for k from t down to t_{start} do 0, if k is environment termination $\gamma_t :=$ γ , otherwise $R_k := r_k + \gamma_t R_{k+1}$ $\delta_k := r_k + \gamma_t V^{\pi}(z_{k+1}, \log \pi_{k+1}) - V^{\pi}(z_k, \log \pi_k)$ $A_k := \delta_k + (\gamma \lambda) A_{k+1}$ $\mathcal{L} := \mathcal{L} + \mathcal{L}_k \text{ (Eq. 7)}$ $\mathcal{A} := \mathcal{A} + A_k \log \pi_k[a_k]$ $\mathcal{H} := \mathcal{H} - \alpha_{\text{entropy}} \sum_{i} \pi_{k}[i] \log \pi_{k}[i]$ (Entropy loss) end for $d\chi' := \nabla_{\chi'} \mathcal{L}$ $d\theta' := \nabla_{\theta'}(\mathcal{A} + \mathcal{H})$ Asynchronously update via gradient ascent θ using $d\theta'$ and χ using $d\chi'$ until $T > T_{max}$

Algorithm 1 of paper "Unsupervised Predictive Memory in a Goal-Directed Agent" by Greg Wayne et al.

NPFL122, Lecture 11 IMPALA PopArt Normalization POMDPs MERLIN CTF-FTW

MERLIN

Figure 2 of paper "Unsupervised Predictive Memory in a Goal-Directed Agent" by Greg Wayne et al.

NPFL122, Lecture 11

IMPALA PopArt Normalization POMDPs

MERLIN CTF-FTW

MERLIN

PopArt Normalization

IMPALA

MERLIN CTF-FTW

POMDPs

MERLIN

Extended Figure 3 of paper "Unsupervised Predictive Memory in a Goal-Directed Agent" by Greg Wayne et al.

NPFL122, Lecture 11

PopArt Normalization

IMPALA

ization POMDPs

MERLIN CTF-FTW

Figure 2 of paper "Human-level performance in first-person multiplayer games with population-based deep reinforcement learning" by Max Jaderber et al.

NPFL122, Lecture 11 IMPALA PopArt N

PopArt Normalization POMDPs

MERLIN CTF-FTW

- Extension of the MERLIN architecture.
- Hierarchical RNN with two timescales.
- Population based training controlling KL divergence penalty weights, slow ticking RNN speed and gradient flow factor from fast to slow RNN.

NPFL122, Lecture 11 IMPALA PopArt Normalization

MERLIN

Figure S10 of paper "Human-level performance in first-person multiplayer games with population-based deep reinforcement learning" by Max Jaderber et al.

NPFL122, Lecture 11 IMPALA PopArt Normalization POMDPs MERLIN CTF-FTW

Figure 4 of paper "Human-level performance in first-person multiplayer games with population-based deep reinforcement learning" by Max Jaderber et al.

NPFL122, Lecture 11

PopArt Normalization

IMPALA

ation POMDPs

MERLIN CTF-FTW