# Supplementary Material:
# Autoregressive DRL with Learned Intrinsic Rewards for Portfolio Optimisation

**Magdalene Hui Qi Lim**
School of Physical and Mathematical
Sciences, Nanyang Technological
University, Singapore
magdalenehq.lim@gmail.com

**Nixie S Lesmana**
School of Physical and Mathematical
Sciences, Nanyang Technological
University, Singapore
Department of Physics, Faculty of Science,
National University of Singapore,
Singapore
nixiesap001@e.ntu.edu.sg

**Chi Seng Pun**
School of Physical and Mathematical
Sciences, Nanyang Technological
University, Singapore
cspun@ntu.edu.sg

## 1 PSEUDOCODE FOR AUGMENTED AUTOA2C

To supplement our proposed algorithm as outlined in section 4 of the main paper, we provide pseudocodes for Augmented AutoA2C. We address the MMDP decomposition in Algorithm 1, and training and trading specifications in Algorithm 2.

---

**Algorithm 1** Rollout for Autoregressive Model

---

1: **Init:** initialize $\mathcal{D} = \{\}$ and action vector $\mathbf{a_t} = (0, 0, \dots, 0) \in \mathbb{R}^{30}$ corresponding to day $t$ of MDP
2: **for** stock $d = 1$ to $30$ **do**
3:     Input $u_{d-1}^{s_t} = \text{concat}(s_t, \mathbf{a_t})$
4:     Sample action $a_t^d \in \mathbb{R}$ using $\pi_\theta$
5:     Update action vector $\mathbf{a_t}[\text{position } d] \leftarrow a_t^d$
6:     Observe intrinsic reward $r^{\text{in}}(u_{d-1}^{s_t}, a_t^d, u_d^{s_t}) \sim \pi_\eta$
7:     **if** $d < 30$ **then**
8:         Get reward $r^{\text{ex}}(u_{d-1}^{s_t}, a_t^d) = 0$
9:     **else**
10:         Execute action $\mathbf{a_t}$ in environment
11:         Observe next state $s_{t+1}$
12:         Observe inter-day reward $r^{\text{ex}}(u_{29}^{s_t}, a_t^{30}) = (b_{t+1} + \mathbf{p_{t+1}}^T\mathbf{h_{t+1}}) - (b_t + \mathbf{p_t}^T\mathbf{h_t})$
13:         Update $s_t \leftarrow s_{t+1}$
14:     **end if**
15:     Append $(u_{d-1}^{s_t}, a_t^d, u_d^{s_t}, r^{\text{ex}}, r^{\text{in}})$ to $\mathcal{D}$
16: **end for**
17: Return trajectory $\mathcal{D}$

---

For training, we maintain a reward buffer $\mathcal{R}$ for the purpose of transforming the inter-day rewards into some external reward function $U(\cdot)$ (see section 5.2) dependent on some history of inter-day rewards.

---

**Algorithm 2** Augmented AutoA2C

---

1: **Input:** step size parameters $\alpha, \beta$
2: **for** each moving window $\mathcal{W}$ **do**
3:     Initialize network parameters $\theta, \eta$
4:     Initialize day $t = 0$
5:     Initialize reward buffer $\mathcal{R}$
6:     # Train
7:     **repeat**
8:         Sample trajectory $\mathcal{D} \sim \pi_\theta, \pi_\eta$ by Algorithm 1
9:         Store inter-day reward $r^{\text{ex}}(s_t, \mathbf{a_t})$ in $\mathcal{R}$
10:        Replace $r^{\text{ex}}(s_t, \mathbf{a_t})$ in $\mathcal{D}$ with $U(s_t, \mathbf{a_t}, s_{t+1}; \mathcal{R})$ or $U(s_{0:t+1}, \mathbf{a_{0:t}}; \mathcal{R})$
11:        Approximate $\nabla_\theta J^{\text{ex+in}}(\theta; \mathcal{D})$ by Equation 6
12:        Update policy actor $\theta' \leftarrow \theta + \alpha \nabla_\theta J^{\text{ex+in}}(\theta; \mathcal{D})$
13:        Update policy critic $\theta_v$ by Equation 7
14:        Approximate $\nabla_{\theta'} J^{\text{ex}}(\theta'; \mathcal{D})$ by Equation 11
15:        Approximate $\nabla_\eta \theta'$ by Equation 12
16:        Compute $\nabla_\eta J^{\text{ex}} = \nabla_{\theta'} J^{\text{ex}}(\theta'; \mathcal{D}) \nabla_\eta \theta'$
17:        Update intrinsic reward actor $\eta' \leftarrow \eta + \beta \nabla_\eta J^{\text{ex}}$
18:        Update intrinsic reward critic $\eta_v$ by Equation 8
19:        $t \leftarrow t + 1$
20:     **until** done
21:     # Trade
22:     **for** trade day $t = 1$ to 63 **do**
23:         Sample $d$ subactions $a^1, \ldots, a^d$ corresponding to $s_t$ according to $\pi_\theta$
24:         Execute action $\mathbf{a_t} = [a^1, \ldots, a^d]$ in environment
25:         Record portfolio value $b_t + \mathbf{p_t}^T \mathbf{h_t}$)
26:     **end for**
27: **end for**

---