

# MARLTS (Multi agent Reinforcement Learning trading system)

I)

- Initially:

- In `train-multi.py` → "policy\_fn": (`env`, `obs_space`, `act_space`,  $\{ \cdot \}$ ),
  - Instantiates two neural networks to become the "brains" of the agent.
  - Two because Using PPO → Actor-Critic method
  - These are

Soldier → • Actor → learns the policy  $\pi_\theta(a|s)$

- takes a state (observation),  $s$ , and decides on action,  $a$ , to take. In this project it takes observation vector and outputs desired portfolio weights,  $\theta$ . Initially,  $\theta$  (weights) are initialized randomly → New, no skills, etc.

Strategist → • Critic, learns Value function  $V_\phi(s)$

- Takes a state,  $s$ , and estimates total expected future reward from that point onwards. Provides a strategic assessment of how "good" current situation is, but initially,  $\theta$  are randomized, signifying no strategic sense.

- The other policies do not get trained, and do not learn, they're basically NPG's. The RL agent is the hero.

Then:

- Deploying our "soldiers" and "strategists" to "corporations".
- Config = config.env\_scanners (num\_env\_scanners = 2, ... )
  - sets up multiple, identical "battlefields"
  - Instead of learning from one "battle" at a time, the agent can fight in multiple simultaneously and pool together the intelligence gained from them

II

- Lets go through a single episode. Here its a sequence of interactions lasting for an episode-length = 252 timesteps
- Timestep,  $t = 0$ 
  - Environment picks random start date from historical data,  $\frac{1}{4}$  RL agents portfolio is initialized with cash = 100\_000, holdings = [0, 0, 0, 0]  $\rightarrow$  env.reset()

2. Environment generates first observation vector  $s_0$ ,

This vector contains historical log returns, all-zero holdings and normalized cash value.

3.  $s_0$  is sent to Policy-f1's "Actor" (soldier)

Network. Right now, this network's weights,  $\Theta$ , are random, so it outputs basically a random action,  $a_0$ , which is a vector of portfolio weights

4. Environment also gets actions from the heuristic agents based on the same market data (or rolling mean, maximizing, etc. agents)

5. Market reaction. The environment aggregates, or combines all these actions into net-notional, calculates the market impact, determines the exec-prices and updates everyone's portfolios.

6. Environment calculates reward  $r_0$ , and the next observation  $s_1$ . Stored as tuple  $(s_0, a_0, r_0, s_1) \Rightarrow$  "Experience"



- Timestep  $t$  = "Some generic step in the episode", agent has some cash and holdings from previous random actions.

1. perception: agent receives current observation vector  $s_t$ .

2.  $s_t$  is now fed to both Actor & Critic networks

↳ • Actor  $\rightarrow$  outputs prob. distr. over actions.

The function "compute-single-action" then "samples" an action,  $a_t$ , from this distribution. Agent still exploring, but decisions now slightly biased by small updates received from previous learning steps

↳ • Critic  $\rightarrow$  outputs a single number  $V(s_t)$ ,

which is the current, yet still poor, estimate of the total future reward it expects to get from the current state

3. Environment collects all agent's actions,  $a_t$ , as well as deterministic actions from mom, mom, etc. agents

↳ impact = np.clip(1.0 + (self.impact\_coeff \* net\_national / (liquidity + 1)),

$\text{exec-prices} = \text{prev-prices} \downarrow \text{instruct}$   
 $\downarrow$   
execution

4. Step function makes the trade at the imputed prices and calculates reward  $r_t$ . Reward is MQL-aware, taking feedback both from profit & performance relative to the graph.

- This goes on for all steps  $t=251$ . Results in a "Complete Report":  $\left[ (s_0, a_0, r_0, V(s_0)), (s_1, a_1, r_1, V(s_1)), \dots \right]$

I II I

PPO Learning Phase:

- Learning begins once parallel environments have collected a batch of trajectories.

1. Calculating advantage

- For a given state  $s_t$ , was the action  $a_t$  that was taken better or worse than the average action we could've taken?  
→
- Advantage function:  $A(s, a)$

- Estimating advantage using "Temporal Difference" Error



(why? - reduces variance  
"one off's")

$$\hat{A}_t = r_t + \gamma V(s_{t+1}) - V(s_t)$$

•  $V(s_t)$  was critics prediction before action

- $r_t + \gamma V(s_{t+1})$  is the improved estimate after seeing the real immediate reward and critics assessment of the next state  $s_{t+1}$ .  $\gamma$  is a factor that "discards" predicted rewards, and values immediate ones higher

- $\hat{A}_t > 0 \rightarrow$  Outcome better than expected,  $a_t$  was good
- $\hat{A}_t < 0 \rightarrow$  Outcome worse than expected,  $a_t$  was bad

## 2. Training Actor

- Goal:  $\uparrow$  probability of actions that had a positive advantage and  $\downarrow$  probability of those with negative advantage

- Conventionally, a single "bad batch" of skewed experiences could spur a policy update.

- PPO Solution  $\Rightarrow$  Clipping :

$$L^{CLIP}(\theta) = \hat{E}_t \left[ \min \left( c_t(\theta) \hat{A}_t, \text{clip} \left( c_t(\theta), 1-\epsilon, 1+\epsilon \right) \hat{A}_t \right) \right]$$

? ? ? ? ? ? ? ? ? ? ? ? ?

- $\hat{A}_t$ : for all the experiences in the batch
- $r_t(\theta) = \frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{old}}(a_t | s_t)}$ : Probability Ratio, measures how much more likely an new, updated policy ( $\pi_\theta$ ) is to take the same action,  $a_t$ , compared to old policy, ( $\pi_{\theta_{old}}$ ) that gathered the data.
- $r_t(\theta) \hat{A}_t$ : Standard objective. If advantage  $\hat{A}_t$  is positive, we want to push  $\pi(\theta)$  up ↑
- $\text{Clip}(r_t(\theta), 1-\epsilon, 1+\epsilon) \hat{A}_t$ : Secondary objective, where probability ratio is "clipped" or confined to a narrow window around 1, (like 0.8 or 1.2, where  $\epsilon=0.2$ ).
- $\min()$ : take the minimum of the two objectives

Why Clip?: Assuming we get a really high  $\hat{A}_t$ , the normal objective would tell us to increase probability of action by a lot. But the clipped objective puts a limit to that, preventing the agent from ...

... getting too "cocky".  $\Rightarrow$  Ensures stable, reliable learning.

- This results in a "optimistic" policy update
  - To maximize  $\rightarrow$  use Gradient Ascent to nudge Actor's weights  $\theta$  in the right direction, Policy clips the update so each "new" model doesn't differ too much from the last

- Through gradient ascent, actor hopes to adjust weights in  $\delta_t(\theta)$ , and more fully:  $L^{critic}(\theta)$

- New weights  $\downarrow$  learning rate  $\downarrow$  gradient  $\nabla$  of objective function  
$$\theta_{k+1} = \theta_k + \alpha \nabla_{\theta} L^{CLIP}(\theta_k)$$

$\nwarrow$  Current weights

### 3. What about the critic?

- Needs to get better at predicting value function
  - prediction was  $V_\phi(s_t)$
  - "true" value (target) is  $V_t^{\text{target}} = r_t + \gamma V(s_{t+1})$
  - Loss is simply mean squared error:  $(V_\phi(s_t) - V_t^{\text{target}})^2$
  - Use Gradient descent to minimize loss,  $(L^{\text{VF}}(\phi))$

↓

Critics predict. become more accurate

↓

Better critic leads to better advantage estimate

↓  
Actor learns more effectively

- For Critic Network Update:

- Need:

1. States ( $s_t$ )

2. Value Targets ( $V_{\text{target},t}$  or  $\hat{R}_t$ )

$$\rightarrow V_{\text{target},t} = \hat{A}_t + V_{\phi_{\text{old}}}(s_t)$$

$$\bullet L^{\text{VF}}(\phi) = \frac{1}{N} \sum_{t=1}^N (V_\phi(s_t) - V_{\text{target},t})^2$$

$\leftarrow$  loss gradient

$$\therefore \phi_{k+1} = \phi_k - \beta \nabla_\phi L^{\text{VF}}(\phi_k)$$

$\uparrow$   
Critic learning rate

\* After num-sgd-int passes over batch, Actor & Critic networks now slightly smarter  $\rightarrow$  Updated policy sent to the env  $\rightarrow$  Cycle repeats, and the initially random agent slowly gets "smarter" compared to others.

# PPO Learning → By Hand

## • Setup:

- Env: MultiAsset Market Env: Continuous action in  $[0, 1]^n$

↓

normalized to weights → trades execute simultaneously

↓

- Agents: "mom-0", "mean-0", "mm-0", "rl-trader-0" [Reward is shaped]

- PPO trained: "policy-61" only in "train-multi.py"

- Action: RLlib Library PPO uses a diagonal Gaussian policy,

skashes box bounds → env renormalizes to sum to 1. Going to ignore tanh Jacobian term for hand calc simplicity.

## • Reward shaping in code:

$$\text{base} = (\text{cash} + \text{holdings.newPrices}) - (\text{cash} + \text{holdings.prevPrices})$$

$$\text{mean_val} = \text{Mean}(\text{all_vals}), \text{std_val} = \text{Std}(\text{all_vals}) + 1e-9$$

$$\# \text{cooperation} \rightarrow \text{coop_bouns} = \text{coop_weight} * (1.0 / (1.0 + \text{std_val}))$$

$$\# \text{competition} \rightarrow \text{comp_bouns} = \text{comp_weight} * \left( (\text{mean_val} - \text{all_vals}[a]) / (\text{std_val} + 1e-9) \right)$$

$$\text{shaped_reward}[a] = \text{base} + \text{coop_bouns} - \text{comp_bouns}$$

↓

(coop\_weight = 0.05) (comp\_weight = 0.1)

Our pseudo-market: 2 assets, one step, no commission, no impact

- prices at  $t$ :  $p_t = [100, 50]$
- prices at  $t+1$ :  $p_{t+1} = [10, 40]$
- each agent starts w cash = 100,000; holdings =  $[0, 0]$
- (example) agents weights at  $t$ :
  - mom-0:  $[0.5, 0.5]$
  - mean-0:  $[0.3, 0.7]$
  - mm-0:  $[0.05, 0.95]$
  - rL-trader-0:  $[0.6, 0.4]$

## A. Converting weights to trades

- target-dollar (in each asset) = weights \* total-portfolio-value  
at time  $t$
- For PPO (rL-trader-0):
  - total  $V_t$  (value at  $t$ ) = 100,000
  - target dollars =  $[60,000, 40,000]$
  - Current \$ in assets (prior to trading) =  $[0, 0]$
  - dollar-diff =  $[60k, 40k]$
  - qty to buy =  $(\text{dollar-diff}) / (\text{price-}t) = \left[ \frac{60k}{100}, \frac{40k}{50} \right] = [600, 800]$

- Would then do this for all other agents

Agent	Weights	target \$	Qty at t ( $P = [100, 80]$ )
mom-0	[0.5, 0.5]	[50k, 50k]	[500, 1000]
mean-0	[0.3, 0.7]	[20k, 70k]	[300, 1400]
mm-0	[0.05, 0.95]	[5k, 95k]	[50, 1900]
rL-trader-0	[0.6, 0.4]	[60k, 40k]	[600, 800]

- In code/practice, would include a small price impact based on `impact_coeff` and `Liquidity-Scale` and then  $\text{exec\_price} = \text{prev\_price} * \text{impact}$ . For hand calc  $\rightarrow \text{impact} = 1.0$

- Post trade, holdings  $\downarrow$ : Cash  $\rightarrow$  New holdings = [600, 800]  
 $\text{Cash} = [100,000 - 60,000 - 40,000] = 0$

B

- Environment advances to  $t+1$  and computes base profit and loss ( $P; L$ ) forward

- Code:

$$\text{Prev-Val} = \text{cash} + \text{holdings} * P_t$$

$$\text{New-Val} = \text{cash} + \text{holdings} * P_{t+1}$$

$$\text{base} = \text{New-Val} - \text{Prev-Val}$$

PPD:

$$\cdot \text{PPO-Val} = 0 + [600, 800] \cdot [100, 50] = 60,000 + 40,000 = 100,000$$

$$\cdot \text{New-Val} = 0 + [600, 800] \cdot [101, 49] = 60,600 + 39,200 = 99,800$$

$$\cdot \text{base} = 99,800 - 100,000 = -200 \rightarrow \text{PPO agent lost \$200}$$

• Had 40% in asset that fell

Other agents:

$$\text{mom\_0: } \cdot \text{PPO} = 100k, \text{ New} = (500 \cdot 101) + (1000 \cdot 49) = 99,500, \text{ base} = -500$$

$$\text{mean\_0: } \text{New} = 98,900 \mid \text{base} = -1100$$

$$\text{mm\_0: } \text{New} = 98,150 \mid \text{base} = -1850$$

Raw P:L rankings: PPO > mom > mean > mm

C

• Reward Shaping

• Stats using New Values  $V_{t+1}$ :

• New\_vals:

• mom: 99,800

• mean: 98,900

• mm: 98,150

• PPO: 99,800

• mean\_vals =  $(99500 + 98900 + 98150 + 99800)/4$   
 $= 99,087.5$

• deviations:

$$[412.5, -109.5, -137.5, 712.5] \rightarrow (x - \mu)$$

• Squares:

$$[170, 156, 35, 156, 878, 906, 509, 656]$$

• mean square =

$$(1,591,875)/4 = 397,968.75$$

• Std =  $\sqrt{397,968.75} = 631.0 \quad (\sigma)$

• Bonus given by PPO:

$$\text{• coop-bonus} = 0.05 \cdot \left( \frac{1}{1+std} \right) = 0.05 \cdot \frac{1}{632} = 0.000079$$

$$\begin{aligned} \text{• Comp-bonus} &= 0.1 \left( \frac{\text{mean-val} - \text{PPO-val}}{\text{std}} \right) \\ &= 0.1 \left( \frac{-712.5}{631} \right) \approx -0.113 \end{aligned}$$

• Final Reward:

$$\begin{aligned} f_t(\text{ppo}) &= \text{base} + \text{coop-bonus} - \text{comp-bonus} \\ &\approx -200 + 0.000079 - (-0.113) \\ &\approx -199.89 \end{aligned}$$

• (Subtracting comp bonus Comp Computed as negative, so to cancel out)

$$\sigma = \sqrt{\frac{\sum (x - \mu)^2}{n}}$$

D

- Returns  $\hat{A}_t$ : Advantages
  - discount  $\gamma = 0.99$
  - GAE parameter  $\lambda = 0.95$
  - $\hat{A}_t = r_t + \gamma V(s_{t+1}) - V(s_t)$

Model Outputs for Critic:

- $V(s_t) = \hat{V}_t = +100.0$  (Remember, Critic outputs forward expectation, so this indicates small positive)
- $V(s_{t+1}) = \hat{V}_{t+1} = +97.02$

Advantage Function

TD-error  $\delta_t \rightarrow$  Temporal Difference (discrepancy between current reward & acc. and next step)

$$\delta_t = r_t + \gamma \hat{V}_{t+1} - \hat{V}_t = -199.887 + 97.02 - 100 = -202.867$$

GAE advantage  $A_t$ : For one step  $A_t \approx \delta_t$ , (longer horizons:  $\gamma \delta_{t+h} + \dots$ )

So in:  $A_t \approx -202.867$

Actual  
Return

Monte-Carlo Return  $R_t$ :  $R_t = r_t + \gamma \hat{V}_{t+1} = -199.887 + 97.02 = -102.867$

- Critic loss will later try to fit  $\hat{V}_t \rightarrow R_t$   
(expected returns  $\rightarrow$  Actual returns)  
(Through NN Back Prop)

## E: Log Probs, ratios and PPO Objective

### Calculating Updates

- Assuming 2-dimension action (pre-squashing), independent Gaussians.

- At  $s_t$ , the old policy - which generated the trajectory - had:

- Means:  $\mu_{\text{old}} = [0.55, 0.45]$
- Std's  $\sigma_{\text{old}} = [0.10, 0.10]$

- Sample Action to produce weights initially:  $a_t = [0.62, 0.41]$

- Current policy (after one gradient step):

means  $\mu = [0.58, 0.43]$

stds  $\sigma = [0.10, 0.10]$  (same for ex)

- Using log Probability of a diagonal gaussian

$$? \rightarrow r_t(\theta) = \frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \rightarrow \log r_t(\theta) = \log \pi_\theta(a_t | s_t) - \log \pi_{\theta_{\text{old}}}(a_t | s_t)$$

Probability Ratio

$$\log \pi(a_t | s_t) = \sum_i \left[ -\frac{1}{2} \left( \frac{a_i - \mu_i}{\sigma_i} \right)^2 - \log(\sqrt{2\pi} \cdot \sigma_i) \right]$$

- Per dim ( $a_t = [0.62, 0.41]$ )  $M_{old} = [0.55, 0.45]$

• Old:

$$dim1 = (0.62 - 0.55) / 0.1 = 0.7 \rightarrow -\frac{1}{2}(0.7)^2 = -0.245$$

$$dim2 = (0.41 - 0.45) / 0.1 = -0.4 \rightarrow -\frac{1}{2}(-0.4)^2 = -0.08$$

- Same  
(constants)
- $-\log(\sqrt{2\pi} \cdot 0.1) \approx -\log(0.25006) \approx +1.383$
  - Sum:  $(-0.245 + 1.383) + (-0.08 + 1.383)$   
 $\hookrightarrow = 2.441$

• Current:

$$-\frac{1}{2}(0.4)^2$$

$$dim1: (0.62 - 0.58) / 0.1 = 0.4 \rightarrow -0.08$$

$$dim2: (0.41 - 0.43) / 0.1 = -0.2 \rightarrow -0.02$$

• Same constants

$$Sum: (0.08 + 1.383) + (0.02 + 1.383) = 2.666$$

$$\log \pi_{old}(a_t | s_t) = 2.441 \quad ; \quad \log \pi(a_t | s_t) = 2.666$$

$$e^{\lambda} (\log f_t(\theta)) \rightarrow f_t(\theta) = e^{\lambda} (\log \pi - \log \pi_{old}) = \exp(2.666 - 2.441)$$

$$= \exp(0.225) \approx 1.252$$

$$\hookrightarrow f_t(\theta) = 1.252$$

## Clipping PPO Objective:

$$L_t^{\text{CLIP}} = \min(\hat{r}_t A_t, \text{clip}(\hat{r}_t, 1-\epsilon, 1+\epsilon) \cdot A_t), \epsilon = 0.2$$

$$\bullet A_t = -202.867$$

$$\bullet \hat{r}_t A_t = 1.252(-202.867) = -254.39$$

$$\bullet \text{Clip } \hat{r}_t \text{ to } [0.8, 1.2] \rightarrow \hat{r}_t^{\text{clipped}} = 1.2 \\ -0.2, +0.2$$

$$\bullet \text{Clipped term} = 1.2 * (-202.867) = -243.44$$

Pessimistic  $\rightarrow$  Taking min:

$$L_t^{\text{CLIP}} = \min(\hat{r}_t A_t, \text{clip}(\hat{r}_t, 1, 1))$$
$$L_t^{\text{CLIP}} = -254.39$$

- Because  $A_t < 0$ , increasing  $\hat{r}_t$  beyond the clip hurts the surrogate more  $\rightarrow$  Clipping protects against overly large policy changes.
- In a batch, w/ multiple results, PPO sums  $L_t$  over timesteps and takes gradient ascent. We're only doing 1 step.

F

- Critic Loss

$$\hookrightarrow L_t^{\text{VF}} = (R_t - \hat{V}_t)^2 = (-102.867 - 100)^2 = (-202.867)^2 \approx 41,156$$

- Entropy  $\rightarrow$  Encourages Exploration  $\rightarrow$  more risk bonus

$\hookrightarrow$  more reward?

$$\hookrightarrow L_t^{\text{ENT}} = \beta \cdot H[\pi(\cdot | s_t)] \text{ with small weight } \beta (0.0-0.1)$$

- For gaussian (like our distribution):

$$\bullet H = \sum_i \frac{1}{2} \log(2\pi\sigma_i^2)$$

$\bullet \sigma_i = 0.1 \rightarrow H$  is a constant per dimension

- Total PPO objective:

$$J(\theta, \phi) = \mathbb{E} \left[ L_{(\theta), v}^{\text{CLIP}} L_{(\phi), c}^{\text{VF}} \right]$$

$\hookrightarrow$  ( $C_v$  and  $C_c$  are coefficients)

$\hookrightarrow$  Puts everything together.

Optimized updates:

- actor params  $\Theta$  to increase  $L^{CLIP} + C_c H$
- critic params  $\phi$  to decrease loss  $L^{VF}$

→ • Advantage is  $-A_{t+1}$ , PPO decreases (log) prob. of action taken

• Critic lowers  $V(s_t)$  because  $-102.9 L + 100$ , (realizing  $L$  expectation), over time  $V(s_t)$  becomes more accurate

• PPO beat mean, so "comp" mode reward less negative

• GAE will blend many  $s_{t+1}$  so policy learns from patterns gathered than days