# An empirical study of the naïve REINFORCE algorithm for predictive maintenance of industrial machines

Rajesh Siraskar

June 20, 2023

## Abstract

The 14th-century friar, William of Ockham, has been attributed for giving us the "Occam's razor" principle – when there are two competing "theories" predicting the same phenomenon, one should prefer the *simpler* of two.

In this empirical study, we document the performance of a simple, early reinforcement learning algorithm, REINFORCE, implemented for a predictive maintenance problem. We compare a very naive implementation of REINFORCE against the predictions of industry-grade Stable-Baselines3 (SB-3) implementations of three advanced algorithms, namely, Deep Q-Network (DQN), Advantage Actor-Critic (A2C) and Proximal Policy Optimization (PPO). Our broad goal was to understand the performance under various scenarios such as a simulation-based environment, three sets of real tool-wear data, added noise levels, and a random chance of break-down. Model performance was measured by how accurately the predictive maintenance agent suggested tool replacement compared to a deterministic preventive maintenance rule based on the tool-wear threshold.

Our findings indicate that the REINFORCE performs significantly well for this particular problem. Across variants of the environment, the REINFORCE algorithm demonstrated an average F1 performance of 0.836 against 0.383 for A2C, 0.471 for DQN, and 0.402 for PPO. As a measure of stability, the overall standard deviation for REINFORCE was 0.041, while A2C, DQN, and PPO standard deviations were 0.059, 0.029, and 0.070, respectively. Across precision on tool replacement, REINFORCE was better

by 0.354 basis points than the best of advanced algorithms and demonstrated a variance lower by 0.004. While the REINFORCE demonstrated better performance for each variant, it was observed that the training was unstable, occasionally producing poor performance models. On the other hand, the SB-3 implementations training was more stable, almost always producing models with an F1 in the range 0.47-0.50.

# 1 Introduction

Introduced in 1992, the REINFORCE algorithm is considered as a basic reinforcement learning algorithm. It is a policy-based, on-policy as well as off-policy algorithm, capable of handling both discrete and continuous observation and action domains.

In practice the REINFORCE algorithm is considered as "weak" learner and superseded by several algorithms developed since. Most notably the Q-Learning and its deep-neural network version, the DQN, followed by Actor-Critic and one of the most robust modern day algorithms, the PPO.

Reinforcement Learning in Robotics: A Survey - Jens Kober J. Andrew Bagnell Jan Peters - Initial gradient-based approaches such as finite differences gradients or REINFORCE (Williams, 1992) have been rather slow. The weight perturbation algorithm is related to REINFORCE but can deal with non-Gaussian distributions which significantly improves the signal to noise ratio of the gradient (Roberts et al., 2010). Recent natural policy gradient approaches (Peters and Schaal, 2008c,b) have allowed for faster convergence which may be advantageous for robotics as it reduces the learning time and required real-world interactions.

| Model | REINFORCE | | | SB-3 A2C | | | SB-3 DQN | | | SB-3 PPO | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Precision | Recall | F1 | Precision | Recall | F1 | Precision | Recall | F1 | Precision | Recall | F1 |
| Simulated - No noise | 1 | 0.7 | 0.823 | 0.487 | 0.44 | 0.46 | 0.067 | 0.01 | 0.017 | 0.457 | 0.28 | 0.343 |
| Simulated - Low NBD | 0.938 | 0.9 | 0.918 | 0 | 0 | 0 | 0.45 | 0.03 | 0.055 | 0.367 | 0.06 | 0.103 |
| Simulated - High NBD | 0.895 | 1 | 0.944 | 0.498 | 0.53 | 0.511 | 0.497 | 0.96 | 0.655 | 0.584 | 0.15 | 0.237 |
| PHM C01 SS - No noise | 0.907 | 0.96 | 0.932 | 0.523 | 0.56 | 0.538 | 0.367 | 0.03 | 0.055 | 0.498 | 0.25 | 0.332 |
| PHM C01 SS - Low NBD | 0.886 | 0.8 | 0.836 | 0.38 | 0.13 | 0.192 | 0.4 | 0.02 | 0.038 | 0.463 | 0.14 | 0.207 |
| PHM C01 SS - High NBD | 0.783 | 0.93 | 0.849 | 0 | 0 | 0 | 0.508 | 0.96 | 0.664 | 0.519 | 0.21 | 0.287 |
| PHM C04 SS - No noise | 0.821 | 0.96 | 0.885 | 0.513 | 0.09 | 0.149 | 0.497 | 0.97 | 0.657 | 0.489 | 0.47 | 0.478 |
| PHM C04 SS - Low NBD | 0.739 | 0.99 | 0.846 | 0.474 | 0.51 | 0.487 | 0.706 | 0.72 | 0.712 | 0.532 | 0.28 | 0.366 |
| PHM C04 SS - High NBD | 0.671 | 0.78 | 0.72 | 0.522 | 0.55 | 0.534 | 0.5 | 0.98 | 0.662 | 0.472 | 0.25 | 0.325 |
| PHM C06 SS - No noise | 1 | 0.65 | 0.785 | 0.465 | 0.46 | 0.461 | 0.511 | 0.98 | 0.671 | 0.379 | 0.14 | 0.198 |
| PHM C06 SS - Low NBD | 0.978 | 0.84 | 0.902 | 0.493 | 0.53 | 0.508 | 0.951 | 0.58 | 0.72 | 0.5 | 0.38 | 0.432 |
| PHM C06 SS - High NBD | 0.715 | 0.88 | 0.789 | 0.505 | 0.47 | 0.482 | 0.505 | 0.98 | 0.667 | 0.369 | 0.11 | 0.169 |
| PHM C01 MS - No noise | 0.798 | 0.92 | 0.852 | 0.514 | 0.59 | 0.549 | 0.584 | 0.97 | 0.729 | 0.486 | 0.24 | 0.313 |
| PHM C04 MS - No noise | 0.774 | 0.69 | 0.728 | 0.476 | 0.53 | 0.5 | 0.506 | 0.97 | 0.664 | 0.495 | 0.34 | 0.401 |
| PHM C06 MS - No noise | 1 | 0.57 | 0.725 | 0.491 | 0.6 | 0.536 | 0.492 | 0.96 | 0.651 | 0.446 | 0.48 | 0.46 |

Table 1: Model performance comparison. SB-3 algorithms trained for 10K episodes.

|  | Precision | | Recall | | F1-score | |
|---|---|---|---|---|---|---|
|  | Mean | SD | Mean | SD | Mean | SD |
| A2C | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| DQN | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| PPO | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| REINFORCE | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |

Table 2: Model performance comparison.

|  | Precision | | Recall | | F1-score | | F1-beta score | |
|---|---|---|---|---|---|---|---|---|
|  | Mean | SD | Mean | SD | Mean | SD | Mean | SD |
| A2C | 1.000 | 2.000 | 3.000 | 4.000 | 5.000 | 6.000 | 7.000 | 8.000 |
| DQN | 1.000 | 2.000 | 3.000 | 4.000 | 5.000 | 6.000 | 7.000 | 8.000 |
| PPO | 1.000 | 2.000 | 3.000 | 4.000 | 5.000 | 6.000 | 7.000 | 8.000 |
| REINFORCE | 1.000 | 2.000 | 3.000 | 4.000 | 5.000 | 6.000 | 7.000 | 8.000 |

Table 3: Model performance comparison (with F-beta score).

xxx

xxx

xxx

xxx

xxx

## 1.1 Algorithm timelines

- 1947: Monte Carlo Sampling

- 1959: Temporal Difference Learning

- 1989: Q-Learning

- 1992: REINFORCE

- 2013: DQN

- 2016: A3C

- 2017: PPO

# 2 The REINFORCE algorithm

Three key features of any RL algorithm:

1. Policy: $\pi_\theta =$ Probablities of all actions, given a state. Parameterized by $\theta$

2. Objective function:
$$\max_\theta J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta}[R(\tau)] \tag{1}$$

3. Method: Way to udate the parameters = Policy Gradient

## 2.1 Policy gradient numerical computation

1. Plain vanilla:

$$\nabla\theta J(\pi_\theta) = \mathbb{E}_{\tau \sim \pi_\theta}\big[\sum_{t=0}^{T} R_t(\tau)\,\nabla_\theta \ln \pi_\theta(a_t|s_t)\,\big] \tag{2}$$

2. With Monte Carlo sampling and approximation: $\nabla_\theta J(\pi_\theta) \approx \big[\sum_{t=0}^{T} R_t(\tau)\,\nabla_\theta \ln \pi_\theta(a_t|s_t)\,\big]$

3. With baseline: $\nabla_\theta J(\theta) \approx \big[\sum_{t=0}^{T}(R_t(\tau) - b(s_t))\,\nabla_\theta \ln \pi_\theta(a_t|s_t)\,\big]$

4. Where, baseline does not change per time-step, it is for the entire trajectory

5. One baseline option: $V^\pi$ - leads to Actor-Critic algorithm

6. Simpler option: Average returns over trajectory: $b = \frac{1}{T}\sum_{t=0}^{T} R_t(\tau)$

Algorithm

# 3 About Stable-Baselines-3

- SB3- paper (Raffin et al., 2021), Raffin et al. (2021)

- sb-3 main doc – (SB3, b)

- sb-3 ppo doc – (SB3, a)

# 4 Method

We normalize the tool wear and other state features, $x \in [0, \ 1] \subset \mathbb{R}$. This allows for adding white noise of similar magnitudes across experiments of different data-sets

# 5 Results

# 6 Discussion

# 7 Conclusion

# References

Stable-baselines3 docs - reliable reinforcement learning implementations, a. URL https://stable-baselines3.readthedocs.io/en/master/modules/ppo.html#how-to-replicate-the-results.

Ppo, b. URL https://stable-baselines3.readthedocs.io/en/master/index.html. Accessed: 2023-05-14.

Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *J. Mach. Learn. Res.*, 22(1), jan 2021. ISSN 1532-4435.