

Model based reinforcement learning at FERMI FEL

Simon Hirlaender

University of Salzburg, Salzburg, Austria

Niky Bruchon

University of Trieste, Trieste, IT

(Dated: November 23, 2020)

In this paper we discuss a model based reinforcement learning approach in comparison to a model free reinforcement learning approach applied at the FERMI FEL system. Both algorithms and approaches are new in this context and the main purpose of this paper is to show how reinforcement learning can be applied on an operational level in a feasible training time on real accelerator physics problems. In terms of sample-complexity the model-based approach is faster, while the final performance of the model free method is superior - the so called asymptotic performance. The model-based algorithm is done in a Dyna-style using an uncertainty aware model and the model-free algorithm is based on tailored deep Q-learning using some tricks to increase the sample efficiency

- Noise benchmarks with pendulum
- Benchmark NAF2 with pendulum
- The first time model based on accelerator problem compared to a novel state of the art MFRL: NAF2
- Tailored to accelerators - short horizons are standard?
- Improvements of the algorithms:
 - Adding the noise feature
 - Jumping feature using SAC
 - Noise in the policy
 - discuss different approaches to e.g. dynamic waiting time

I. INTRODUCTION AND MOTIVATION

In particle accelerators one main goal is to provide stable and reproducible performance. In order to achieve this, a number of control problems have to be considered simultaneously. Especially if there is no way to model the physics, one might use optimisation techniques as e.g. derivative free optimisers (DFOs) or model-based optimisations as Gaussian processes to restore or maintain the performance. Another promising way is to apply reinforcement learning which shows several advantages over optimisation methods:

- It covers a larger class of problems, RL optimises a sequence of decisions.
- It memorizes the problem and does not start always from scratch as an DFO.
- Existing data can be used.
- The underlying structure of the problem might be deduced.

On critical aspect using RL is the number of iteration to train a controller. In this paper, we present the study carried out to solve the maximisation-problem of the radiation intensity generated by a seeded Free-Electron Laser (FEL) on the Free Electron laser Radiation for Multidisciplinary Investigations (FERMI) at Elettra Sincrotrone Trieste. Three algorithms were applied successfully showing different advantages, which will be discussed in the following.

In a seeded free-electron laser one of the most critical parameters is the temporal and transverse overlap of the electron and laser beam in the magnetic section called modulator undulator.

II. THE PHYSICAL SET-UP

At FERMI several beam-based feedback systems are deployed to control the beams trajectories shot to shot with the main purpose of guaranteeing a steady intensity of the light radiation. Nevertheless, in the last years various approaches have been studied to investigate their applicability in tuning operations. The paper is organised as follows.

- Description of the problem set-up at FERMI
- Overview of RL
- Details of the implementations used in these studies and theoretical concerns
- Results
- Summary and outlook

III. DEEP REINFORCEMENT LEARNING

Assume states \mathbf{s} the set of all states \mathcal{S} and actions \mathbf{a} all actions \mathcal{A} and rewards r and the set of all rewards

\mathcal{R} an initial state distribution d_0 . Goal of reinforcement learning is to find a $\pi : \mathbf{s} \mapsto \mathbf{a}$, which is the solution of:

$$\max J(\pi) = \mathbb{E}_{\tau \sim p_\pi} \left[\sum_{t=0}^H \gamma^t r(\mathbf{s}_t, \mathbf{a}_t) \right], \quad (1)$$

where $p_\pi = d_0(\mathbf{s}_0) \prod_{t=0}^H \pi(\mathbf{a}_t, \mathbf{s}_t) T(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t)$ is the distribution of all trajectories $\tau := (\mathbf{s}_0, \mathbf{a}_0, \mathbf{s}_1, \mathbf{a}_1, \dots, \mathbf{s}_H, \mathbf{a}_H)$ drawn by π with a horizon H . In the modern field of RL one distinguishes if the policy $\pi(\mathbf{a}) \approx \pi_\phi(\mathbf{s})$ or the state-action value function $Q(\mathbf{s}, \mathbf{a})$ is approximated using a high capacity function approximator, as e.g. a deep neural network. In the first case one speaks about policy gradient methods in the latter about approximate dynamic programming, which we now discuss.

A. Approximate dynamic programming

The state-value function $Q(\mathbf{s}, \mathbf{a})$ is expressed as $Q_\theta(\mathbf{s}, \mathbf{a})$, where θ denotes the parameters of the function approximator. By satisfying the Bellmann-optimality equation Q_θ can be trained towards the optimal $Q^*(\mathbf{s}, \mathbf{a})$:

$$\min_{\theta} \left(\tilde{Q}_\theta - \mathcal{B}^* \tilde{Q}_\theta \right)^2. \quad (2)$$

And π can be calculated via:

$$\pi_\theta(\mathbf{s}) = (\delta(\mathbf{a}) - \arg\max_{\mathbf{a}} Q_\theta(\mathbf{s}, \mathbf{a})). \quad (3)$$

The Bellman operator \mathcal{B}^* has a unique fixed point but is non-linear, due to the max - operator:

$$\mathcal{B}^* \tilde{Q}_\theta(\mathbf{s}_t, \mathbf{a}_t) := r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \max_{\mathbf{a}} (Q_\theta(\mathbf{s}_{t+1}, \mathbf{a}) - Q_\theta(\mathbf{s}_t, \mathbf{a}_t)). \quad (4)$$

The form of this equation can cause overestimation and other complications, when using a function approximator. Several methods exist which try to mitigate the problems [1]. One way to reduce the effect is to take a simple analytical form of the Q -function.

B. NAF

If a specific quadratic form of the Q function is assumed:

$$Q_\theta(\mathbf{s}, \mathbf{a}) = -\frac{1}{2}(\mathbf{a} - \mu_\theta(\mathbf{s}))P_\theta(\mathbf{a} - \mu_\theta(\mathbf{s}))^T + V_\theta(s). \quad (5)$$

One modification, which is used in these tests is a twin network (weights denoted by $\tilde{\theta}$), where only one is used to obtain the policy, while the other is used for the update rule to avoid over-estimation. It is motivated by double Q -learning [2]. The maximum is given analytically as $\max_{\mathbf{a}} Q(\mathbf{s}, \mathbf{a}) = V(\mathbf{s})$, hence from eq. (2):

$$\min_{\theta} \left(\left(r(\mathbf{s}_t, \mathbf{a}_t) + \gamma \min_{1,2} V_{\theta_{i_{\text{targ}}}}(\mathbf{s}_{t+1}) - (1 + \gamma)Q_\theta(\mathbf{s}_t, \mathbf{a}_t) \right)^2 \right) \quad (6)$$

θ_{targ} are the weights of a target network, which is softly updated. To stabilize the network training a small artificial noise is added to the actions in the update. This algorithm has an extremely good sample-efficiency for suitable problems as can be found in accelerators and is used as the baseline for the considered control problem.

Algorithm 1: On-policy policy gradient with Monte Carlo estimator

```

1: initialise  $\theta_0$  for iteration  $k \in [0, \dots, K]$  do
  2:
    sample trajectories  $\{\tau_i\}$  by running  $\pi_{\theta_k}(\mathbf{a}_t | \mathbf{s}_t)$ 
     $\triangleright$  each  $\tau_i$  consists of  $\mathbf{s}_{i,0}, \mathbf{a}_{i,0}, \dots, \mathbf{s}_{i,H}, \mathbf{a}_{i,H}$ 

```

IV. UNCERTAINTY AWARE DYNA-STYLE REINFORCEMENT LEARNING

The original Dyna algorithm [3] is modified here in several aspects. Generally, Dyna style algorithms [4] denote algorithms, where a MFRL algorithm is trained on purely synthetic data from an approximate dynamics model or on a mixture of synthetic and real data. We use only synthetic data to reduce the interaction with the real environment to a minimum. An overview of the used method is shown in fig. 1. At the beginning the data is collected or read in from a pre-collection. An uncertainty aware model is trained, using anchored ensembles [5] on the data, which allows to take the allegorical (measurement errors) as well as the epistemic uncertainty (lack of data) into account. Subsequently a MFRL algorithm is deployed to learn the controller on the learned model by only using the synthetic data, by taking the uncertainty into account. After a defined number of training steps the controller is tested on each model individually. If there is no improvement in a specific ratio < 1 of models, the controller is tested on the real environment for a number of episodes. The training is stopped if the controller solves the problem on the real environment.

A. Critical design decisions

In the following the most important aspects of a successful application of a Dyna-AE algorithm are listed:

- The ANN - including the prior - number of models - noise level - early stopping
- The number of data points and the policy.
- The uncertainty - how is this included?
- The episodic design - avoid long trajectories.
- The MFRL agent as e.g. the TRPO and PPO or the SAC and its training
- Tuning the algorithm on a model.

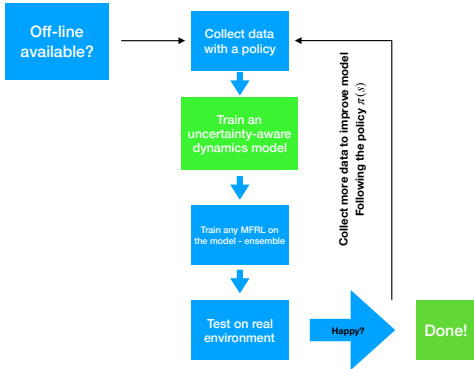


FIG. 1: A schematic overview of the AE-Dyna approach used in this paper.

We try to discuss all items in some detail in the following. The anchoring ensemble methods usually yields good results already with a small number of different networks. The idea behind it, is to mimic the posterior probability of the dynamics model to capture first of all its own uncertainty due to the lack of data. Empirical results show that three to five models were sufficient to see clear advantages over a single network approach and the main goal is not to calculate the exact posterior of the model. A small two-layer network with around 15-25 nodes and the *tanh* activation were used. The last layer was linear and the inputs were normalized to the interval $[-1,1]$. The prior is controlled by the initialization of the weights, which were normalized by the number of nodes in each layer. The noise level is also respected during the initialization. For the moment only homeostatic errors are used. Several methods to set the training method were implemented. Early stopping, a standard techniques showed good results, mainly to the fact, that the number of training steps is increasing with more data. The advantage is that the uncertainty at regions without data shrink in this way leading to a better training performance. A fixed loss value threshold was also tested, as also used in the original anchore-ensemble technique. In the experiments a combination of both was taken.

One of the most crucial points is how this uncertainties are taken into account by the RL algorithms. Several different approaches were tried. It is possible to take the average of the models and add Gaussian noise leveled by the standard deviation of the models. Another straight forward way is to randomly select a model to provide at each training step, which is the original implementation of the ME-TRPO algorithm and was used in the experiments labeled *ME-TRPO*. A pessimistic setting only would select the model resulting in the lowest predicted reward. Good results were obtained by following a randomly selected single model each full episode and were

used in the experiments labeled *AE-Dyna*.

The number of data-points taken every time the model as improved was firstly determined by the number of initial data points, collect using a random policy. The initialization phase has to be chosen not to small to minimize the chance of getting trapped in a local minimum for too long. Afterwards at least on full episode should be taken. We decided to use a short total episode length to reduce the impact of the compound error, the accumulation error following a wrong model, as well known for Dyna-style approaches. The maximal number of steps was ten. During the data collected on the real system, the latest policy was taken with some small noise added to improve the exploration.

To decide which MFRL algorithm to use, two main algorithm classes have to be considered: the on and the off-policy algorithms. Online algorithm show a more stable convergence in general, while off-policy algorithms have the advantage that the data buffer can be filled with real data as well. A very attractive off-policy algorithm is the Soft-actor-critic. This algorithm tries to maximize the entropy of the actions to find a good trade-off between exploration and exploitation. On the other hand the on-policy algorithm TRPO provides some theoretical improvement guarantees. In the latter case the policy was improved over the whole training, while for the SAC the policy was reset, to take advantage of the exploration features.

B. Operational deployment

Niky

C. short FERMI introduction

The Free Electron laser Radiation for Multidisciplinary Investigation better known as the FERMI, is the seeded free-electron laser facility next to the third-generation synchrotron facility ELETTRA at the Elettra Sincrotrone Trieste laboratories. A free-electron laser is a fourth-generation light source where the lasing medium consists of very-high-speed electron moving freely through a magnetic structure. The FERMI peculiarity is given by the usage of an external seeding source that provides several advantages, as the increased stability in pulse and photon energy, reduced size of the device, improved longitudinal coherence, more control on the pulse structure, and improved temporal synchronization with external timing systems.

The external optical laser signal provide is contribution to the FEL process in the modulator where it interacts with the relativistic electron beam modulating it in energy. The modulation in energy is the converted into a charge modulation in the dispersive section, and finally the density modulated beams radiation is amplified in the radiators section. The importance of ensuring the

best possible overlapping between the seed laser and the electron beam in the modulator is therefore evident.

For this reason the proposed study focuses on the control of the seed laser trajectory in the modulator, looking at FERMI performance as a reference.

D. our system: the modulator and the seed laser

A more detailed description of the alignment process in the modulator is here provided. The most critical parameters in a seeded free-electron laser are the temporal and transverse overlap of the electron and laser beams in the modulator magnetic section. The problem is simplified by keeping constant the trajectory of the electron beam and the mechanical delay line that controls the temporal alignment. The final problem faced consists in optimising the seed laser trajectory to match the electron beam and consequently increase the intensity of the FEL radiation.

V. PROOF-OF-PRINCIPLE APPLICATION OF RL AGENT AT FERMI TRAJECTORY CORRECTION

Niky

A. Experiment results from FERMI RL tests

Niky and Simon

VI. DISCUSSION AND OUTLOOK

VII. CONCLUSIONS

Appendix A: A non-linear standard control problem

To provide some transparency of these studies for other labs we provide results on a famous classical standard control problem [6], the *inverted pendulum*. It is a non-linear low dimensional unsolved continuous control problem. Unsolved means there is no threshold for the reward to terminate an episode. The episode length is set to 200 steps. In the following several tests were carried out on the *inverted pendulum* to demonstrate the improvements of the selected algorithms, mainly in terms of noise handling. This is of importance when dealing with measurements on real systems.

1. The impact of noise

We add artificial Gaussian noise $\mathcal{N}(0,0.05)$ on the state with and compare the different *NAF* variants: *smoothing-double*, *smoothing single*, *standart*.

Appendix B: NAF2 details

Appendix C: AE training details

Appendix D: Theoretical aspects on the AE-DYNA on FERMI FEL

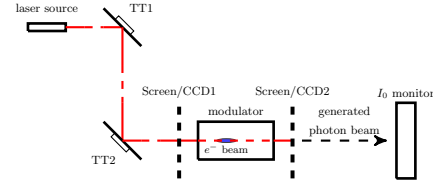


FIG. 2: The training .

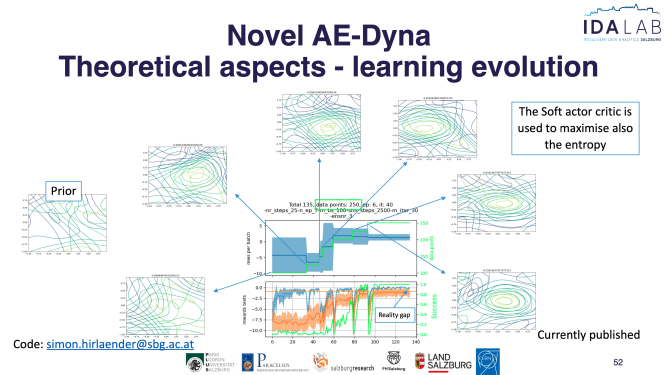


FIG. 3: The training .

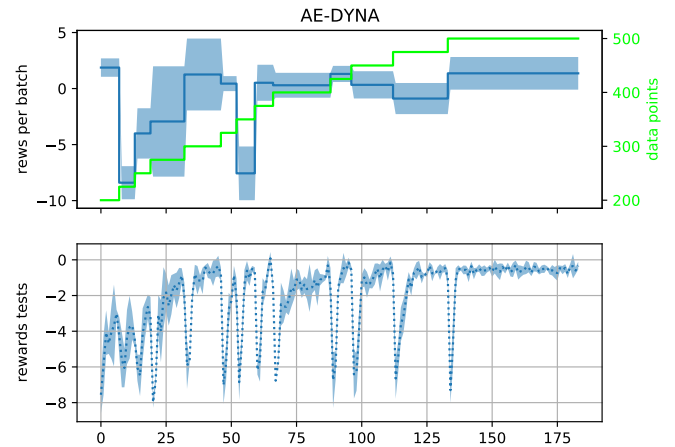


FIG. 4: The training .

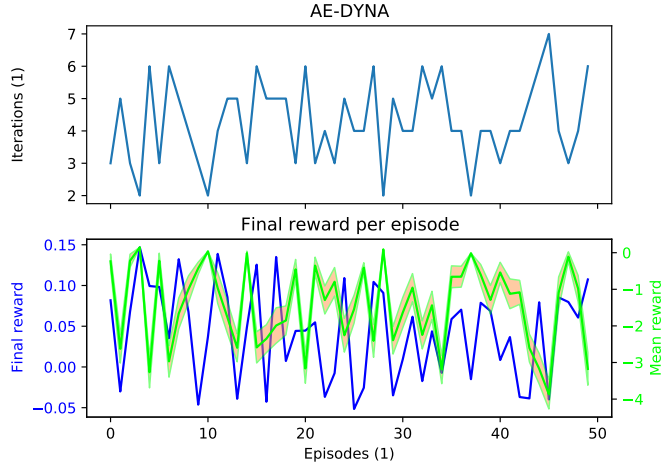


FIG. 5: The training .

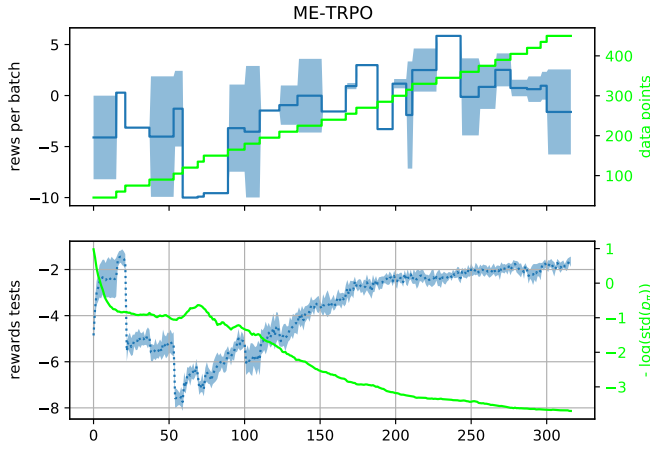


FIG. 6: The training .

-
- [1] S. Levine, A. Kumar, G. Tucker, and J. Fu, Offline reinforcement learning: Tutorial, review, and perspectives on open problems, 2005.01643.
 - [2] H. Hasselt, Double q-learning, in *Advances in Neural Information Processing Systems*, Vol. 23, edited by J. Lafferty, C. Williams, J. Shawe-Taylor, R. Zemel, and A. Culotta (Curran Associates, Inc., 2010) pp. 2613–2621.
 - [3] T. Kurutach, I. Clavera, Y. Duan, A. Tamar, and P. Abbeel, Model-ensemble trust-region policy optimization, 1802.10592.
 - [4] R. S. Sutton, Dyna, an integrated architecture for learning, planning, and reacting, *ACM SIGART Bulletin* **2**, 160 (1991).
 - [5] T. Pearce, F. Leibfried, A. Brintrup, M. Zaki, and A. Neely, Uncertainty in neural networks: Approximately bayesian ensembling, 1810.05546.
 - [6] K. Furuta, M. Yamakita, and S. Kobayashi, Swing up control of inverted pendulum, in *Proceedings IECON '91: 1991 International Conference on Industrial Electronics, Control and Instrumentation* (IEEE).

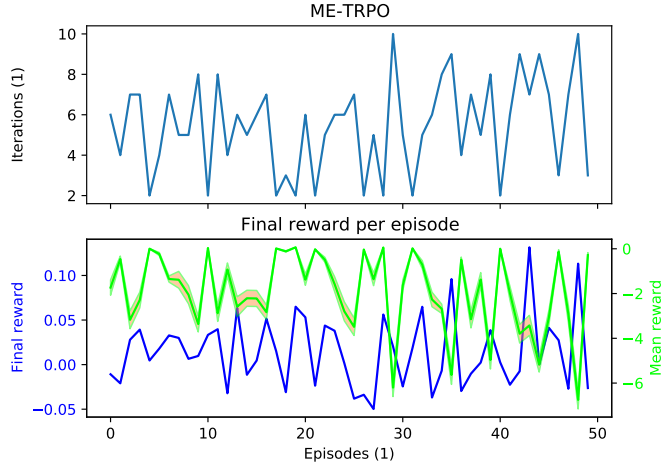


FIG. 7: The training .

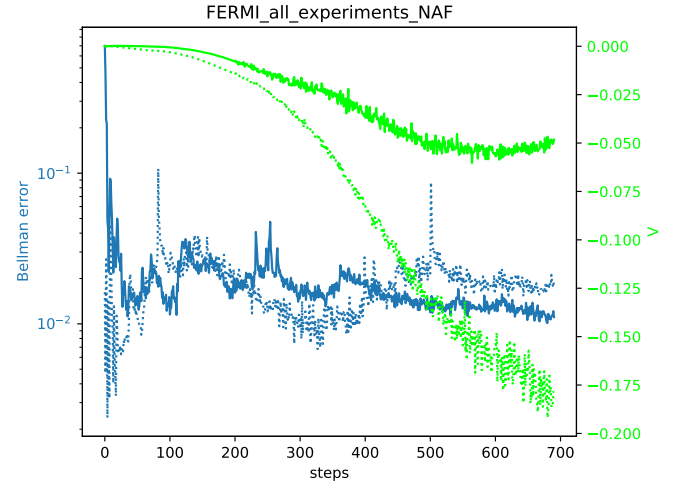


FIG. 9: The training .

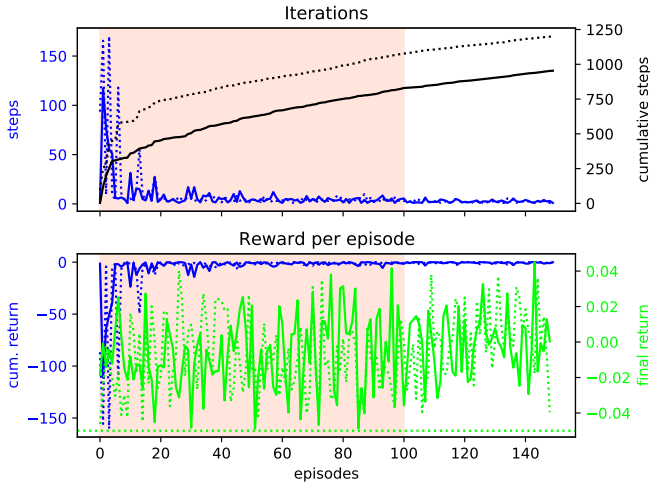


FIG. 8: The training .