



Forward Model Learning for Motion Control Tasks - Challenges for Learning and Planning -

Rudolf Kruse Faculty of Computer Science Otto von Guericke University Magdeburg Alexander Dockhorn School of Electrical and Computer Engineering Queen Mary University London





Forward Model Learning for Motion Control Tasks

Challenges for Learning and Planning

Alexander Dockhorn and Rudolf Kruse

Motivation

- Computational Intelligence refers to the ability of a computer to learn a specific task from data or experimental observation.
 - steering of self-driving cars and parking aids
 - automating a production pipeline





OTTO VON GUERICKE UNIVERSITÄT MAGDEBURG



Motivation

Problems with real environments

- Real tasks are hard to setup
- Failure of the algorithm has considerable cost
- It becomes hard to study the algorithms underlying characteristics

Games can be simulations of real world tasks

- Quantifiable goals, controllable difficulty, and large data sets
- Digital games are fully accessible to computers





A Short History of Computational Intelligence in Games







Modern Reinforcement Learning Examples







Agent–Environment Interface



A general framework for studying games which consists of the elements

- Agent: the learner and decision-maker
- Environment: anything the agent interacts with, e.g. a game
- Actions: Agent and environment interact continuously interact with each other.
- **Reward:** numerical values provided that the agent tries to maximize over time.





Components of a Game



State $S_t \in S$ can be perceived through multiple sensors $(S_t^{(1)}, S_t^{(2)}, \dots, S_t^{(n)})$.

• a state may not be fully observable (partial information game)

The whole environment can be modelled as a probability distribution of possible outcomes:

$$P(R_{t+1}, S_{t+1} | S_0, A_0, S_1, A_1, \ldots, S_t, A_t)$$

• But we can also model both components separately





Problem Classification

Two popular learning approaches:

- 1) Learn which actions are good
- 2) Learn to anticipate the future

Both allow the definition of learning algorithms as well as approximate solutions and hybrid algorithms.







Most Popular Algorithm Classes

Reinforcement Learning

 Performance depends on the available training time

Simulation-based Search

 Performance depends on the available computation time during evaluation







Current and Proposed Solution

Deep Reinforcement Learning

 Learn an approximate function that maps the state to the expected reward

Prediction-based Search

 Learn an approximate function that predicts the next state







Forward Model Learning

Goal: Learn to predict the upcoming states of the environment.

A forward model fm maps the environment's state S_t and the agent's action A_t at time t to the upcoming state S_{t+1} of the environment

$$\mathit{fm}: (\mathcal{S} imes \mathcal{A}) o \mathcal{S} \qquad (S_t, \mathcal{A}_t) \longmapsto S_{t+1}$$

Above definition applies to environment models that fulfill the Markov Property: $P(S_{t+1} \mid S_0, A_0, S_1, A_1, \dots, S_t, A_t) \Rightarrow P(S_{t+1} \mid S_t, A_t)$





End-To-End Solution

Learning a model by consecutive interaction with

the environment.

Each observed state transition represents a single training example.

Classifiers and regressors can be used for the prediction







End-To-End Solution

Deep Learning Forward Models

- Current systems often make use of deep neural networks to predict upcoming states
- Many training examples are required to fit the networks parameters

The training time is dependent on the size of the state and action space. Both can be enormous!

How to reduce the model's training time?







Decomposed Forward Model

Assumption:

• sensor values can be modelled independently

$$\forall i, j \in 1..n : i \neq j \Rightarrow S_{t+1}^{(i)} \perp L S_{t+1}^{(j)} \mid S_t, A_t$$

Learn on sub-model for each observable sensor value

$$fm_i: (S_t, A_t) \longmapsto S_{t+1}^{(i)}$$

Aggregate the result of each sensor value prediction $fm(S_t, A_t) = (fm_1(S_t, A_t), fm_2(S_t, A_t), \dots, fm_n(S_t, A_t))$ $= (S_{t+1}^{(1)}, S_{t+1}^{(2)}, \dots, S_{t+1}^{(n)}) = S_{t+1}$







Optimal Decomposition

In case enough observation data is available:

- infer the dependencies among the in- and outputs sensor values
- optimal decomposition can be achieved

Alternatively, we can ask experts to model the variables dependencies or use heuristic solutions.

With each independency we:

- reduce the model space considerably
- require less training data
- create a more efficient model





Environment's State					
S ⁽¹⁾	S ⁽²⁾	S ⁽³⁾	S ⁽⁴⁾	•••	$S_t^{(n)}$
???					
FM_1	FM ₂	FM ₃	FM_4		FM _n
???					
$S_{t+1}^{(1)}$	S ⁽²⁾ _{t+1}	$S_{t+1}^{(3)}$	$S_{t+1}^{(4)}$	•••	$S_{t+1}^{\left(n\right)}$
Predicted Environment's State					

Continuous State-Space Forward Models

The decomposed forward model can be used to model continuous state-spaces

• Instead of predicting the resulting state, predict the changes in between states

i-th Component Model: $p(s_t^i | s_{t-1}, a_{t-1}, \dots, s_1, a_1)$ i-th Differential Model: $p(s_{t-1}^i - s_t^i | s_{t-1}, a_{t-1}, \dots, s_1, a_1)$

Example: Predicting the movement of a robot

- *End-to-End*: Predict the robot's final position at time t+1
- *Differential Decomposed Model*: separately predict the change in position of its arms and legs





Model Requirements

Model accuracy:

• accurate predictions are required to simulate future time-steps

Model speed:

- the trained model needs to be fast to facilitate more simulations
 Model size:
- the number of parameters should be low to

Interpretability and Reliability:

• risk aware applications require interpretable models





Prediction-based Search







Comparison of Learning Approaches

(Deep) Reinforcement Learning:

• immediately output the action that is best according to previous experiences

Prediction-based Search:

- search for the best action-sequence according to the simulated environment
- Solution can be interpreted since the whole search tree is available
- Evaluating the model's prediction confidence and model the agent's risk
- > Performance scales with the forward model's accuracy and the available search time





Motion Control Test Environments

Five Motion Control Environments of the OpenAl Gym Framework gym.openai.com

- Testing discrete in- and outputs as well as rewards
- All tested environments are fairly low in complexity
 - By studying them we want to achieve a clear picture on prediction-based

search methods and how they perform



Experimental Setup

For each environment we train the agent by:

- Continuously observe the environment after each action
- Update the model at the end of each episode
- The process is repeated 10 times, recording the mean reward per episode

Results are compared to Reinforcement Learning Model

- Discrete action space: SARSA [1], DQN [2], CEM [3]
- Continuous action space: NAF [4]

Training Results I/II

Both, **Cart Pole** and **Acrobot**, can be effectively learned and show stable results

First successful run in **Cart Pole** has been observed in the fifth episode.

In **Acrobot** the prediction accuracy slowly increases for longer search depths.

Training Results II/II

Lunar Lander: the agent quickly learn to land in case the target is directly below it

 But fails to land at distant positions due to missing positive examples

Pendulum: the agent quickly learns to swing up the pole and keeps it balanced

• With limited search depth, the agent does not swing to the opposite direction for speed

Video for Training and Evaluation Runs

Prediction Accuracy Over Time

Comparing the true and the predicted state

allows us to measure the models accuracy:

- Initial states are predicted well since many of them have been sampled during training
- Errors propagate over consecutive predictions
- Rare events or insufficiently sampled states

can yield drastic prediction errors

Adapted from: Freeman, C. D., Metz, L., & Ha, D. (2019). Learning to Predict Without Looking Ahead: World Models Without Forward Prediction. 9, 1–12. http://arxiv.org/abs/1910.13038

Planning Horizon Dilemma I/II

Fixing a planning horizon is a non-trivial task:

- Long planning horizons yield a more accurate estimation of the expected return
- but it also increases the number of states to be analyzed

In addition to this trade-off the model's accuracy limits the effective search depth

- As shown before: errors will accumulate over time reducing our confidence with every layer of the search tree
- In comparison: the true forward model provides an accurate prediction for each level of the search tree

Planning Horizon Dilemma II/II

Probabilistic classifiers may allow us to measure the model's confidence in made predictions

- in case we can measure the confidence of our prediction, we may bound the search to a confidence threshold
- therefore, dynamically limiting the search depth

Even better: learn models that are reliable in the long-term prediction task

Exploration vs. Exploitation during Training

The classic **exploration vs. exploitation** tradeoff affects the training procedure

- exploration: gathering samples of the whole state-space to assure an accurate prediction for all value ranges
- **exploitation**: focus on analyzing promising areas of the state-space to efficiently learn how to act in the unknown environment

The latter can be seen in the training process of Lunar Lander

- the agent quickly learns to land in case the drone is directly above the platform
- but fails in case the platform is in another position

Conclusion

Forward Model Learning can be feasible for motion control applications but lacks

- model centric training approaches
- long-term reliable forward models

Similarly, search-based methods can yield great performance but are missing

• methods for efficiently selecting actions in continuous action spaces

State- and action-abstraction methods may resolve these problems.

Future Work

- How does the presented approach compare with model-based reinforcement learning?
- Extend the efficiency of search-based algorithms in case of continuous state and action-spaces.
- Estimate the reliability of a model's prediction given the current observation.
- Implement application-centric state and action abstractions to speed-up the model building process.

Thank you for your attention!

Our Related Papers

Apeldoorn, D., & Dockhorn, A. (2020). Exception-Tolerant Hierarchical Knowledge Bases for Forward Model Learning. *IEEE Transactions on Games* (accepted). <u>https://doi.org/10.1109/TG.2020.3008002</u>

Dockhorn, A., Lucas, S. M., Volz, V., Bravi, I., Gaina, R. D., & Perez-Liebana, D. (2019). Learning Local Forward Models on Unforgiving Games. 2019 IEEE Conference on Games (CoG), 1–4. <u>https://doi.org/10.1109/CIG.2019.8848044</u>

Dockhorn, A., Tippelt, T., & Kruse, R. (2018). Model Decomposition for Forward Model Approximation. 2018 IEEE Symposium Series on Computational Intelligence (SSCI), 1751–1757. <u>https://doi.org/10.1109/SSCI.2018.8628624</u>

Rudolf Kruse Faculty of Computer Science Otto von Guericke University Magdeburg Alexander Dockhorn School of Electrical and Computer Engineering Queen Mary University London

References

- [1] R. S. Sutton and A. G. Barto, Reinforcement Learning, 2nd ed. Cambridge: The MIT Press, 2018.
- [2] V. Mnih et al., "Human-level control through deep reinforcement learning," Nature, vol. 518, no. 7540, pp. 529-533, 2015.
- [3] I. Szita and A. Lorincz, "Learning tetris using the noisy cross-entropy method," Neural Computation, vol. 18, no. 12, pp. 2936-2941, 2006.
- [4] S. Gu, T. Lillicrap, U. Sutskever, and S. Levine, "Continuous deep q-learning with model-based acceleration,"33rd International Conference on Machine Learning, ICML 2016, vol. 6, pp. 4135-4148, 2016

