
Learning the Reward Function for a Misspecified Model: Supplemental Material

Erik Talvitie¹

A. Hallucinated DAGger-MC Details

Hallucinated DAGger-MC, like earlier variations on DAGger, requires the ability to reset to the initial state distribution μ and also the ability to reset to an “exploration distribution” ν . The exploration distribution ideally ensures that the agent will encounter states that would be visited by a good policy. The performance bound for H-DAGger-MC depends in part on the quality of the selected ν .

In addition to assuming a particular form for the planner (one-ply MC with a blind rollout policy), H-DAGger-MC requires the dynamics model to be “unrolled”. Rather than learning a single \hat{P} , H-DAGger-MC learns a set $\{\hat{P}^1, \dots, \hat{P}^{T-1}\} \subseteq \mathcal{P}$, where model \hat{P}^i is responsible for predicting the outcome of step i of a rollout, given the state sampled from \hat{P}^{i-1} . While this impractical condition is important theoretically, Talvitie (2017) showed that in practice a single \mathcal{P} can be used for all steps; the experiments in Section 5 make use of this practical alteration.

Algorithm 1 augments H-DAGger-MC to learn a reward model as well as a dynamics model. In particular, H-DAGger-MC proceeds in iterations, each iteration producing a new plan, which is then used to collect data to train a new model. In each iteration state-action pairs are sampled using the current plan and the exploration distribution (lines 7-13), and then the world and model are rolled out in parallel to generate hallucinated training examples (lines 14-21). The resulting data is used to update the model. We simply add a reward model learning process, and collect training examples along with the state transition examples during the rollout. After both parts of the model have been updated, a new plan is generated for the subsequent iteration. Note that while the dynamics model is “unrolled”, there is only a single reward model that is responsible for predicting the reward at every step of the rollout. We assume that the reward learning algorithm is performing a weighted regression

¹Department of Computer Science, Franklin & Marshall College, Lancaster, Pennsylvania, USA. Correspondence to: Erik Talvitie <erik.talvitie@fandm.edu>.

Algorithm 1 Hallucinated DAGger-MC (+ reward learning)

Require: LEARN-DYNAMICS, LEARN-REWARD, exploration distr. ν , MC-PLANNER(blind rollout policy ρ , depth T), # iterations N , # rollouts per iteration K .

- 1: Get initial datasets $\mathcal{D}_1^{1:T-1}$ and \mathcal{E}_1 (maybe using ν)
- 2: Initialize $\hat{P}_1^{1:T-1} \leftarrow \text{LEARN-DYNAMICS}(\mathcal{D}_1^{1:T-1})$.
- 3: Initialize $\hat{R}_1 \leftarrow \text{LEARN-REWARD}(\mathcal{E}_1)$.
- 4: Initialize $\hat{\pi}_1 \leftarrow \text{MC-PLANNER}(\hat{P}_1^{1:T-1}, \hat{R}_1)$.
- 5: **for** $n \leftarrow 2 \dots N$ **do**
- 6: **for** $k \leftarrow 1 \dots K$ **do**
- 7: With probability... \triangleright First sample from ξ
- 8: $1/2$: Sample $(x, b) \sim D_{\mu}^{\hat{\pi}_n}$
- 9: $1/4$: Reset to $(x, b) \sim \nu$.
- 10: $(1-\gamma)/4$: Sample $x \sim \mu, b \sim \hat{\pi}_n(\cdot | x)$.
- 11: $\gamma/4$: Reset to $(y, c) \sim \nu$
- 12: Sample $x \sim P(\cdot | y, c), b \sim \hat{\pi}_n(\cdot | x)$
- 13: Let $s \leftarrow x, z \leftarrow x, a \leftarrow b$.
- 14: **for** $t \leftarrow 1 \dots T - 1$ **do** \triangleright Parallel rollouts...
- 15: Sample $s' \sim P(\cdot | s, a)$.
- 16: Add $\langle z, a, s' \rangle$ to \mathcal{D}_n^t .
 \triangleright (DAGger-MC adds $\langle s, a, s' \rangle$)
- 17: Add $\langle z, a, R_s^a, \gamma^{t-1} \rangle$ to \mathcal{E}_n .
 \triangleright (Standard approach adds $\langle s, a, R_s^a, \gamma^{t-1} \rangle$)
- 18: Sample $z' \sim \hat{P}_{n-1}^t(\cdot | z, a)$.
- 19: Let $s \leftarrow s', z \leftarrow z'$, and sample $a \sim \rho$.
- 20: **end for**
- 21: Add $\langle z, a, R_s^a, \gamma^{T-1} \rangle$ to \mathcal{E}_n .
 \triangleright (Standard approach adds $\langle s, a, R_s^a, \gamma^{T-1} \rangle$)
- 22: **end for**
- 23: $\hat{P}_n^{1:T-1} \leftarrow \text{LEARN-DYNAMICS}(\hat{P}_{n-1}^{1:T-1}, \mathcal{D}_n^{1:T-1})$
- 24: $\hat{R}_n \leftarrow \text{LEARN-REWARD}(\hat{R}_{n-1}, \mathcal{E}_n)$
- 25: $\hat{\pi}_n \leftarrow \text{MC-PLANNER}(\hat{P}_n^{1:T-1}, \hat{R}_n)$.
- 26: **end for**
- 27: **return** the sequence $\hat{\pi}_{1:N}$

(where each training example is weighted by γ^{t-1} for the rollout step t in which it occurred).

A.1. Analysis of H-DAGger-MC

We now derive theoretical guarantees for this new version of H-DAGger-MC. The analysis is similar to that of existing

DAgger variants (Ross & Bagnell, 2012; Talvitie, 2015; 2017), but the proof is included for completeness. Let H_n^t be the distribution from which H-DAgger-MC samples a training example at depth t (lines 7-13 to pick an initial state-action pair, lines 14-21 to roll out). Define the average error of the dynamics model at depth t to be

$$\bar{\epsilon}_{prd}^t = \frac{1}{N} \sum_{n=1}^N \mathbf{E}_{(s,z,a) \sim H_n^t} [1 - \hat{P}_n^t(\sigma_s^a | z, a)].$$

Let $\epsilon_{\hat{R}_n}(s, z, a) = |R(s, a) - \hat{R}_n(z, a)|$ and let

$$\bar{\epsilon}_{hrwd} = \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^T \gamma^{t-1} \mathbf{E}_{(s,z,a) \sim H_n^t} [\epsilon_{\hat{R}_n}(s, z, a)]$$

be the average reward model error. Finally, let D_n^t be the distribution from which H-DAgger-MC samples s and a during the rollout in lines 14-21. The error of the reward model with respect to these environment states is

$$\bar{\epsilon}_{erwd} = \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^T \gamma^{t-1} \mathbf{E}_{(s,a) \sim D_n^t} [|R(s, a) - \hat{R}(s, a)|].$$

For a policy π , let $c_\nu^\pi = \sup_{s,a} \frac{D_{\mu,\pi}(s,a)}{\nu(s,a)}$ represent the mismatch between the discounted state-action distribution under π and the exploration distribution ν . Now, consider the sequence of policies $\hat{\pi}_{1:N}$ generated by H-DAgger-MC. Let $\bar{\pi}$ be the uniform mixture over all policies in the sequence. Let $\bar{\epsilon}_{mc} = \frac{1}{N} \frac{4}{1-\gamma} \sum_{n=1}^N \|\bar{Q}_n - \hat{Q}_{T,n}^\rho\|_\infty + \frac{2}{1-\gamma} \|BV_T^\rho - V_T^\rho\|_\infty$ be the error induced by the choice of planning algorithm, averaged over all iterations.

Lemma 7. *In H-DAgger-MC, the policies $\hat{\pi}_{1:N}$ are such that for any policy π ,*

$$\begin{aligned} \mathbf{E}_{s \sim \mu} [V^\pi(s) - V^{\bar{\pi}}(s)] &\leq \frac{4}{1-\gamma} c_\nu^\pi \bar{\epsilon}_{hrwd} + \bar{\epsilon}_{mc} \\ &\leq \frac{4}{1-\gamma} c_\nu^\pi \left(\bar{\epsilon}_{erwd} + 2M \sum_{t=1}^{T-1} \gamma^{t-1} \bar{\epsilon}_{prd}^t \right) + \bar{\epsilon}_{mc}. \end{aligned}$$

Proof. Recall that

$$\mathbf{E}_{s \sim \mu} [V^\pi(s) - V^{\bar{\pi}}(s)] = \frac{1}{N} \sum_{n=1}^N \mathbf{E}_{s \sim \mu} [V^\pi(s) - V^{\hat{\pi}_n}(s)].$$

and by Lemma 1 for any $n \geq 1$,

$$\begin{aligned} \mathbf{E}_{s \sim \mu} [V^\pi(s) - V^{\hat{\pi}_n}(s)] &\leq \\ &\frac{4}{1-\gamma} \mathbf{E}_{(s,a) \sim \xi_\mu^{\pi, \hat{\pi}_n}} [|\hat{Q}_{T,n}^\rho(s, a) - Q_T^\rho(s, a)|] + \bar{\epsilon}_{mc}, \end{aligned}$$

where

$$\begin{aligned} \xi_\mu^{\pi, \hat{\pi}_n}(s, a) &= \frac{1}{2} D_{\mu, \hat{\pi}_n}(s, a) + \frac{1}{4} D_{\mu, \pi}(s, a) \\ &\quad + \frac{1}{4} \left((1-\gamma) \mu(s) \hat{\pi}_n(a | s) \right. \\ &\quad \left. + \gamma \sum_{z,b} D_{\mu, \pi}(z, b) P_z^b(s) \hat{\pi}_n(a | s) \right). \end{aligned}$$

Then, combining the above with Theorem 5,

$$\begin{aligned} &\frac{1}{N} \sum_{n=1}^N \frac{4}{1-\gamma} \mathbf{E}_{(s,a) \sim \xi_\mu^{\pi, \hat{\pi}_n}} [|\hat{Q}_{T,n}^\rho(s, a) - Q_T^\rho(s, a)|] + \bar{\epsilon}_{mc} \\ &\leq \frac{1}{N} \sum_{n=1}^N \frac{4}{1-\gamma} \sum_{t=1}^T \gamma^{t-1} \mathbf{E}_{(s,z,a) \sim H_{\xi_\mu^{\pi, \hat{\pi}_n, \rho}}^{t,n}} [\epsilon_{\hat{R}_n}(s, z, a)] + \bar{\epsilon}_{mc} \end{aligned}$$

Now note that for any t and any n ,

$$\begin{aligned} &\mathbf{E}_{(s,z,a) \sim H_{\xi_\mu^{\pi, \hat{\pi}_n, \rho}}^{t,n}} [\epsilon_{\hat{R}_n}(s, z, a)] \\ &= \frac{1}{2} \sum_{s', a'} D_{\mu, \hat{\pi}_n}(s', a') \mathbf{E}_{(s,z,a) \sim H_{s', a', \rho}^{t,n}} [\epsilon_{\hat{R}_n}(s, z, a)] \\ &\quad + \frac{1}{4} \sum_{s', a'} D_{\mu, \pi}(s', a') \mathbf{E}_{(s,z,a) \sim H_{s', a', \rho}^{t,n}} [\epsilon_{\hat{R}_n}(s, z, a)] \\ &\quad + \frac{\gamma}{4} \sum_{s', a'} \sum_{s'', a''} D_{\mu, \pi}(s'', a'') P_{s''}^{a''}(s') \hat{\pi}_n(a' | s') \\ &\quad \mathbf{E}_{(s,z,a) \sim H_{s', a', \rho}^{t,n}} [\epsilon_{\hat{R}_n}(s, z, a)] \\ &\quad + \frac{1-\gamma}{4} \sum_{s', a'} \mu(s') \hat{\pi}_n(a' | s') \\ &\leq \frac{1}{2} \sum_{s', a'} D_{\mu, \hat{\pi}_n}(s', a') \mathbf{E}_{(s,z,a) \sim H_{s', a', \rho}^{t,n}} [\epsilon_{\hat{R}_n}(s, z, a)] \\ &\quad + \frac{1}{4} c_\nu^\pi \sum_{s', a'} \nu(s', a') \mathbf{E}_{(s,z,a) \sim H_{s', a', \rho}^{t,n}} [\epsilon_{\hat{R}_n}(s, z, a)] \\ &\quad + \frac{\gamma}{4} c_\nu^\pi \sum_{s', a'} \sum_{s'', a''} \nu(s'', a'') P_{s''}^{a''}(s') \hat{\pi}_n(a' | s') \\ &\quad \mathbf{E}_{(s,z,a) \sim H_{s', a', \rho}^{t,n}} [\epsilon_{\hat{R}_n}(s, z, a)] \\ &\quad + \frac{1-\gamma}{4} \sum_{s', a'} \mu(s') \hat{\pi}_n(a' | s') \end{aligned}$$

$$\mathbf{E}_{(s,z,a) \sim H_{s', a', \rho}^{t,n}} [\epsilon_{\hat{R}_n}(s, z, a)]$$

$$\begin{aligned}
 &\leq c_\nu^\pi \left(\frac{1}{2} \sum_{s', a'} D_{\mu, \hat{\pi}_n}(s', a') \right. \\
 &\quad \mathbf{E}_{(s, z, a) \sim H_{s', a', \rho}^{t, n}} [\epsilon_{\hat{R}_n}(s, z, a)] \\
 &\quad + \frac{1}{4} \sum_{s', a'} \nu(s', a') \mathbf{E}_{(s, z, a) \sim H_{s', a', \rho}^{t, n}} [\epsilon_{\hat{R}_n}(s, z, a)] \\
 &\quad + \frac{\gamma}{4} \sum_{s', a'} \sum_{s'', a''} \nu(s'', a'') P_{s', a'}^{a''}(s') \hat{\pi}_n(a' | s') \\
 &\quad \mathbf{E}_{(s, z, a) \sim H_{s', a', \rho}^{t, n}} [\epsilon_{\hat{R}_n}(s, z, a)] \\
 &\quad + \frac{1-\gamma}{4} \sum_{s', a'} \mu(s') \hat{\pi}_n(a' | s') \\
 &\quad \left. \mathbf{E}_{(s, z, a) \sim H_{s', a', \rho}^{t, n}} [\epsilon_{\hat{R}_n}(s, z, a)] \right) \\
 &= c_\nu^\pi \mathbf{E}_{(s, z, a) \sim H_{\xi_n, \rho}^{t, n}} [\epsilon_{\hat{R}_n}(s, z, a)].
 \end{aligned}$$

When $t = 1$,

$$\mathbf{E}_{(s, z, a) \sim H_{\xi_n, \rho}^{t, n}} [\epsilon_{\hat{R}_n}(s, z, a)] = \mathbf{E}_{(s, a) \sim \xi_n(s, a)} [\epsilon_{\hat{R}_n}(s, z, a)].$$

When $t > 1$,

$$\begin{aligned}
 &\mathbf{E}_{(s, z, a) \sim H_{\xi_n, \rho}^{t, n}} [\epsilon_{\hat{R}_n}(s, z, a)] \\
 &= \sum_{s_t, z_t, a_t} \mathbf{E}_{(s_1, a_1) \sim \xi_n} \left[\sum_{a_{1:t-1}} \rho(a_{2:t} | a_1) \right. \\
 &\quad \left. P_{s_1}^{a_{0:t-1}}(s_t | s_1, a_{0:t-1}) \hat{P}_n^{1:t-1}(z_t | s_1, a_{0:t-1}) \right. \\
 &\quad \left. \epsilon_{\hat{R}_n}(s_t, z_t, a_t) \right] \\
 &= \mathbf{E}_{(s, z, a) \sim H_n^t} [\epsilon_{\hat{R}_n}(s, z, a)].
 \end{aligned}$$

Thus, putting it all together, we have shown that

$$\begin{aligned}
 &\mathbf{E}_{s \sim \mu} [V^\pi(s) - V^{\hat{\pi}}(s)] \\
 &\leq \frac{4}{1-\gamma} c_\nu^\pi \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^T \gamma^{t-1} \mathbf{E}_{(s, z, a) \sim H_n^t} [\epsilon_{\hat{R}_n}(s, z, a)] \\
 &\quad + \bar{\epsilon}_{mc} \\
 &= \frac{4}{1-\gamma} c_\nu^\pi \bar{\epsilon}_{hrwd} + \bar{\epsilon}_{mc}.
 \end{aligned}$$

Thus we have proven the first inequality. Furthermore, by

Theorem 6,

$$\begin{aligned}
 \bar{\epsilon}_{hrwd} &= \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^T \mathbf{E}_{(s, z, a) \sim H_n^t} [\hat{R}_n(s, z, a)] \\
 &\leq \frac{1}{N} \sum_{n=1}^N \left(\sum_{t=1}^T \gamma^{t-1} \mathbf{E}_{(s, a) \sim D_n^t} [|R(s, a) - \hat{R}_n(s, a)|] \right. \\
 &\quad \left. + 2M \sum_{t=1}^{T-1} \gamma^{t-1} \mathbf{E}_{(s, z, a) \sim H_n^t} [1 - \hat{P}_n^t(\sigma_s^a | z, a)] \right) \\
 &= \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^T \gamma^{t-1} \mathbf{E}_{(s, a) \sim D_n^t} [|R(s, a) - \hat{R}_n(s, a)|] \\
 &\quad + 2M \sum_{t=1}^{T-1} \gamma^{t-1} \frac{1}{N} \sum_{n=1}^N \mathbf{E}_{(s, z, a) \sim H_n^t} [1 - \hat{P}_n^t(\sigma_s^a | z, a)] \\
 &\leq \bar{\epsilon}_{erwd} + 2M \sum_{t=1}^{T-1} \gamma^{t-1} \epsilon_{prd}^t.
 \end{aligned}$$

This gives the second inequality. \square

Note that this result holds for *any* comparison policy π . Thus, if $\bar{\epsilon}_{mc}$ is small and the learned models have low error, then if ν is similar to the state-action distribution under *some* good policy, $\bar{\pi}$ will compare favorably to it. That said, Lemma 7 shares the limitations of the comparable results for the other DAgger algorithms. It focuses on the L1 loss, which is not always a practical learning objective. It also assumes that the expected loss at each iteration can be computed exactly (i.e. that there are infinitely many samples per iteration). It also applies to the average policy $\bar{\pi}$, rather than $\hat{\pi}_N$. Ross & Bagnell (2012) discuss extensions that address more practical loss functions, finite sample bounds, and results for $\hat{\pi}_N$.

Lemma 7 effectively says that *if* the models have low training error, the resulting policy will be good. It does not promise that the models will have low training error. Following Ross & Bagnell (2012) note that $\bar{\epsilon}_{prd}^t$ and $\bar{\epsilon}_{hrwd}$ can each be interpreted as the average loss of an online learner on the problem defined by the aggregated datasets. Then for each horizon depth t let $\bar{\epsilon}_{\mathcal{P}}^t$ be the error of the best dynamics model in \mathcal{P} under the training distribution at that depth, in retrospect. Specifically,

$$\bar{\epsilon}_{\mathcal{P}}^t = \inf_{P' \in \mathcal{P}} \frac{1}{N} \sum_{n=1}^N \mathbf{E}_{(s, z, a) \sim H_n^t} [1 - P'(\sigma_s^a | z, a)].$$

Similarly, let

$$\bar{\epsilon}_{\mathcal{R}} = \inf_{R' \in \mathcal{R}} \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^T \gamma^{t-1} \mathbf{E}_{(s, z, a) \sim H_n^t} [\epsilon_{R'}(s, z, a)]$$

be the error of the best reward model in \mathcal{R} in retrospect.

The average regret for the dynamics model at depth t is $\bar{\epsilon}_{prgt}^t = \bar{\epsilon}_{prd}^t - \bar{\epsilon}_{\mathcal{P}}^t$. For the reward model it is $\bar{\epsilon}_{rrgt} = \bar{\epsilon}_{hrwd} - \bar{\epsilon}_{\mathcal{R}}$. For a no-regret online learning algorithm, average regret approaches 0 as $N \rightarrow \infty$. This gives the following bound on H-Dagger-MC's performance in terms of model regret.

Theorem 8. *In H-Dagger-MC, the policies $\hat{\pi}_{1:N}$ are such that for any policy π ,*

$$\begin{aligned} & \mathbf{E}_{s \sim \mu} [V^\pi(s) - V^{\hat{\pi}}(s)] \\ & \leq \frac{4}{1-\gamma} c_\nu^\pi (\bar{\epsilon}_{\mathcal{R}} + \bar{\epsilon}_{rrgt}) + \bar{\epsilon}_{mc} \\ & \leq \frac{4}{1-\gamma} c_\nu^\pi \left(\bar{\epsilon}_{erwd} + 2M \sum_{t=1}^{T-1} \gamma^{t-1} (\bar{\epsilon}_{\mathcal{P}}^t + \bar{\epsilon}_{prgt}^t) \right) + \bar{\epsilon}_{mc} \end{aligned}$$

and if the learning algorithms are no-regret then as $N \rightarrow \infty$, $\bar{\epsilon}_{rrgt} \rightarrow 0$ and $\bar{\epsilon}_{prgt}^t \rightarrow 0$ for each $1 \leq t \leq T-1$.

Theorem 8 says that if \mathcal{R} contains a low-error reward model relative to the learned dynamics models then, as discussed above, if $\bar{\epsilon}_{mc}$ is small and ν visits important states, the resulting policy will yield good performance. If \mathcal{P} and \mathcal{R} contain perfect models, $\hat{\pi}$ will be comparable to the plan generated by the perfect model.

As noted by Talvitie (2017), this result does *not* promise that H-Dagger-MC will eventually achieve the performance of the best available set of dynamics models. The model at each rollout depth is trained to minimize prediction error given the input distribution provided by the shallower models without regard for the effect on deeper models. It is possible that better overall error could be achieved by *increasing* the prediction error at one depth in exchange for a favorable state distribution for deeper models. Similarly, as discussed in Section 4, H-Dagger-MC will not necessarily achieve the performance of the best available combination of dynamics and reward models. The dynamics model is trained without regard for the impact on the reward model. It could be that a dynamics model with higher prediction error would allow for lower hallucinated reward error. H-Dagger-MC does not take this possibility into account.

References

- Ross, S. and Bagnell, D. Agnostic system identification for model-based reinforcement learning. In *Proceedings of the 29th International Conference on Machine Learning (ICML)*, pp. 1703–1710, 2012.
- Talvitie, E. Agnostic system identification for monte carlo planning. In *Proceedings of the 29th AAAI Conference on Artificial Intelligence (AAAI)*, pp. 2986–2992, 2015.
- Talvitie, E. Self-correcting models for model-based reinforcement learning. In *Proceedings of the Thirty-First*

AAAI Conference on Artificial Intelligence (AAAI), pp. 2597–2603, 2017.