# **6** Supplementary Material

#### 6.1 Gradient Derivation

For ease of reference in this section we give a brief outline of the derivation for the derivative of (1). As in the rest of the paper we focus on the case of an infinite planning horizon with discounted rewards, where other frameworks follow similarly. The first point to note is that for any  $t \in \mathbb{N}$  we have the following identity, often referred to as the 'log-trick',

$$\nabla_{\boldsymbol{w}} p(\boldsymbol{z}_{1:t}; \boldsymbol{w}) = p(\boldsymbol{z}_{1:t}; \boldsymbol{w}) \nabla_{\boldsymbol{w}} \log p(\boldsymbol{z}_{1:t}; \boldsymbol{w}),$$

where the usual limit arguments are made in the case  $p(z_{1:t}; w) = 0$ . Upon interchanging the order of integration and differentiation the gradient takes the form

$$\nabla_{\boldsymbol{w}} U(\boldsymbol{w}) = \sum_{t=1}^{\infty} \mathbb{E}_{p(\boldsymbol{z}_{1:t};\boldsymbol{w})} \bigg[ \gamma^{t-1} R(\boldsymbol{z}_t) \nabla_{\boldsymbol{w}} \log p(\boldsymbol{z}_{1:t};\boldsymbol{w}) \bigg].$$

Due to the Markovian structure of the trajectory distribution (2) this derivative can be written in the equivalent form

$$\nabla_{\boldsymbol{w}} U(\boldsymbol{w}) = \sum_{t=1}^{\infty} \sum_{\tau=1}^{t} \mathbb{E}_{p(\boldsymbol{z}_{\tau}, \boldsymbol{z}_{t}; \boldsymbol{w})} \bigg[ \gamma^{t-1} R(\boldsymbol{z}_{t}) \nabla_{\boldsymbol{w}} \log p(\boldsymbol{a}_{\tau} | \boldsymbol{s}_{\tau}; \boldsymbol{w}) \bigg],$$
$$= \sum_{\tau=1}^{\infty} \sum_{t=\tau}^{\infty} \mathbb{E}_{p(\boldsymbol{z}_{\tau}, \boldsymbol{z}_{t}; \boldsymbol{w})} \bigg[ \gamma^{t-1} R(\boldsymbol{z}_{t}) \nabla_{\boldsymbol{w}} \log p(\boldsymbol{a}_{\tau} | \boldsymbol{s}_{\tau}; \boldsymbol{w}) \bigg],$$

where the second line follows from the first through an interchange of the summations. The chain structure of the trajectory distribution allows the expectation over the marginals of the two time-points of the trajectory distribution to be written as follows

$$\nabla_{\boldsymbol{w}} U(\boldsymbol{w}) = \sum_{\tau=1}^{\infty} \mathbb{E}_{p_{\tau}(\boldsymbol{z};\boldsymbol{w})} \bigg[ \sum_{t=\tau}^{\infty} \mathbb{E}_{p_{t}(\boldsymbol{z}';\boldsymbol{w})} \bigg[ \gamma^{t-1} R(\boldsymbol{z}') \big| \boldsymbol{z}_{\tau} = \boldsymbol{z} \bigg] \nabla_{\boldsymbol{w}} \log p(\boldsymbol{a}|\boldsymbol{s};\boldsymbol{w}) \bigg], \quad (15)$$

where we have used the notation  $p_{\tau}(z; w) \equiv p(z_{\tau} = z; w)$ , for  $\tau \in \mathbb{N}$ . The summation over the inner expectation in (15) can be seen to be equal to the state-action value function scaled by  $\gamma^{\tau-1}$ , *i.e.* 

$$\gamma^{\tau-1}Q(\boldsymbol{z};\boldsymbol{w}) = \sum_{t=\tau}^{\infty} \mathbb{E}_{p_t(\boldsymbol{z}';\boldsymbol{w})} \bigg[ \gamma^{t-1}R(\boldsymbol{z}') \big| \boldsymbol{z}_{\tau} = \boldsymbol{z} \bigg].$$

Inserting this form for this inner expectation into (15) gives

$$\nabla_{\boldsymbol{w}} U(\boldsymbol{w}) = \sum_{\tau=1}^{\infty} \mathbb{E}_{p_{\tau}(\boldsymbol{z};\boldsymbol{w})} \left[ \gamma^{\tau-1} Q(\boldsymbol{z};\boldsymbol{w}) \nabla_{\boldsymbol{w}} \log p(\boldsymbol{a}|\boldsymbol{s};\boldsymbol{w}) \right]$$
$$= \mathbb{E}_{p_{\gamma}(\boldsymbol{z};\boldsymbol{w})Q(\boldsymbol{z};\boldsymbol{w})} \left[ \nabla_{\boldsymbol{w}} \log p(\boldsymbol{a}|\boldsymbol{s};\boldsymbol{w}) \right],$$

where the second line follows from the definition of  $p_{\gamma}(z; w)$ . This completes the derivation of (4). This derivation is different, but equivalent, to the standard derivation in [27].

## 6.2 Hessian Derivation

Following similar calculations to those used section(6.1) it is simple to see that the Hessian takes the form

$$\mathcal{H}(\boldsymbol{w}) = \sum_{t=1}^{\infty} \mathbb{E}_{p(\boldsymbol{z}_{1:t};\boldsymbol{w})} \bigg[ \gamma^{t-1} R(\boldsymbol{z}_t) \bigg( \nabla_{\boldsymbol{w}} \log p(\boldsymbol{z}_{1:t};\boldsymbol{w}) \nabla_{\boldsymbol{w}}^T \log p(\boldsymbol{z}_{1:t};\boldsymbol{w}) + \nabla_{\boldsymbol{w}} \nabla_{\boldsymbol{w}}^T \log p(\boldsymbol{z}_{1:t};\boldsymbol{w}) \bigg) \bigg],$$

where the first term is obtained by a second application of the 'log-trick'. This is the form of the Hessian given in (5). The form of  $\mathcal{H}_2(w)$  given in (9) follows from analogous calculations to those performed in section(6.1) of the supplementary material.

#### 6.3 Expectation Maximisation Derivation

For ease of reference in this section we give a brief derivation of the application of Expectation Maximisation to Markov Decision Processes. As in the rest of the paper we focus on the case of an infinite planning horizon with discounted rewards, where other frameworks follow similarly. A quantity of central importance in this derivation is the following non-negative function

$$\tilde{p}(\boldsymbol{z}_{1:t}, t; \boldsymbol{w}) = \gamma^{t-1} R(\boldsymbol{z}_t) p(\boldsymbol{z}_{1:t}; \boldsymbol{w}).$$

This function is non-negative because the reward function is considered as non-negative. Note that this function is trans-dimensional, *i.e.* it is a function of both  $t \in \mathbb{N}$  and the state-action variables up until the  $t^{\text{th}}$  time-point. It can be seen from (1) and (2) that the 'normalisation constant' of  $\tilde{p}(\boldsymbol{z}_{1:t},t;\boldsymbol{w})$  is equal to  $U(\boldsymbol{w})$ , *i.e.*  $U(\boldsymbol{w}) = \sum_{t=1}^{\infty} \sum_{\boldsymbol{z}_{1:t}} \tilde{p}(\boldsymbol{z}_{1:t},t;\boldsymbol{w})$ . We denote the normalised version of this function by  $\hat{p}(\boldsymbol{z}_{1:t},t;\boldsymbol{w})$ .

To obtain the lower-bound on the log-objective we note that the Kullback-Leibler divergence between any two distributions is always non-negative and so calculating the Kullback-Leibler divergence between the variational distribution,  $q(z_{1:t}, t)$ , and  $\hat{p}(z_{1:t}, t; w)$  gives

$$KL(q||\hat{p}) = \log U(\boldsymbol{w}) - \mathcal{H}_{entropy}(q(\boldsymbol{z}_{1:t}, t)) - \mathbb{E}_{q(\boldsymbol{z}_{1:t}, t)} \left[\log \tilde{p}(\boldsymbol{z}_{1:t}, t; \boldsymbol{w})\right] \ge 0.$$

The lower-bound (10) is now immediately obtained from the definition of  $\tilde{p}(\boldsymbol{z}_{1:t}, t; \boldsymbol{w})$ .

An EM-algorithm is obtained through coordinate-wise optimisation of the bound (10) w.r.t. the variational distribution (the E-step) and the policy parameters (the M-step). In the E-step the lower-bound (10) is optimised when the Kullback-Leibler divergence is zero, namely when  $q(\mathbf{z}_{1:t}, t) \equiv \hat{p}(\mathbf{z}_{1:t}, t; \mathbf{w}_k)$ , where  $\mathbf{w}_k$  are the current policy parameters. In the E-step the lower-bound is optimised w.r.t.  $\mathbf{w}$ , which is equivalent to optimising the quantity

$$\mathbb{E}_{\tilde{p}(\boldsymbol{z}_{1:t},t;\boldsymbol{w}_{k})}\left[\log \tilde{p}(\boldsymbol{z}_{1:t},t;\boldsymbol{w})\right] = \mathbb{E}_{\tilde{p}(\boldsymbol{z}_{1:t},t;\boldsymbol{w}_{k})}\left[\sum_{\tau=1}^{t}\log \pi(\boldsymbol{a}_{\tau}|\boldsymbol{s}_{\tau};\boldsymbol{w})\right] + \text{terms independent of } \boldsymbol{w}.$$

Using the definition of  $\tilde{p}(\boldsymbol{z}_{1:t}, t; \boldsymbol{w}_k)$  we have

$$\mathbb{E}_{\tilde{p}(\boldsymbol{z}_{1:t},t;\boldsymbol{w}_{k})}\left[\sum_{\tau=1}^{t}\log\pi(\boldsymbol{a}_{\tau}|\boldsymbol{s}_{\tau};\boldsymbol{w})\right] = \sum_{t=1}^{\infty}\sum_{\tau=1}^{t}\mathbb{E}_{p(\boldsymbol{z}_{\tau},\boldsymbol{z}_{t};\boldsymbol{w}_{k})}\left[\gamma^{t-1}R(\boldsymbol{z}_{t})\log\pi(\boldsymbol{a}_{\tau}|\boldsymbol{s}_{\tau};\boldsymbol{w})\right],$$

so that using the same manipulations as those used in section(6.1) of the supplementary material we have that the E-step is equivalent to maximising the function

$$\mathcal{Q}(\boldsymbol{w}, \boldsymbol{w}_k) = \mathbb{E}_{p_{\gamma}(\boldsymbol{z}; \boldsymbol{w}_k) Q(\boldsymbol{z}; \boldsymbol{w}_k)} \bigg[ \log \pi(\boldsymbol{a} | \boldsymbol{s}; \boldsymbol{w}) \bigg],$$

*w.r.t.* to the first parameter, w. This completes the derivation of the EM-algorithm for Markov Decision Processes with an infinite planning horizon and discounted rewards.

## 6.4 Affine Invariance of Approximate Newton Method

To show that the use of  $\mathcal{H}_2(w)$  in place of the true Hessian maintains the linear (affine) invariance of the Newton method we use the following formulae

$$\nabla_{\boldsymbol{w}}\tilde{f}(\boldsymbol{w}) = \mathcal{T}^T \nabla_{\boldsymbol{v}} f(\boldsymbol{v}), \qquad \nabla_{\boldsymbol{w}} \nabla_{\boldsymbol{w}}^T \tilde{f}(\boldsymbol{w}) = \mathcal{T}^T \nabla_{\boldsymbol{v}} \nabla_{\boldsymbol{v}}^T f(\boldsymbol{v}) \mathcal{T}_{\boldsymbol{v}}$$

where f is some twice differentiable function of w,  $\tilde{f}(w) = f(\mathcal{T}w)$  and  $v = \mathcal{T}w$ . Using these formulae we have the following two identities

$$\nabla_{\boldsymbol{w}} \log p(\boldsymbol{a}|\boldsymbol{s}; \mathcal{T}\boldsymbol{w}) = \mathcal{T}^{T} \nabla_{\boldsymbol{v}} \log p(\boldsymbol{a}|\boldsymbol{s}; \boldsymbol{v}),$$
$$\nabla_{\boldsymbol{w}} \nabla_{\boldsymbol{w}}^{T} \log p(\boldsymbol{a}|\boldsymbol{s}; \mathcal{T}\boldsymbol{w}) = \mathcal{T}^{T} \nabla_{\boldsymbol{v}} \nabla_{\boldsymbol{v}}^{T} \log p(\boldsymbol{a}|\boldsymbol{s}; \boldsymbol{v}) \mathcal{T}$$

which hold for each  $(s, a) \in S \times A$ . Defining  $\tilde{U}(w) = U(\mathcal{T}w)$ , we have  $\nabla_w \tilde{U}(w) = \mathcal{T}^T U(v)$ . Following calculations almost identical to those in section(6.2) it can be shown that  $\tilde{\mathcal{H}}_2(w)$  takes the form

$$\tilde{\mathcal{H}}_{2}(\boldsymbol{w}) = \mathbb{E}_{p_{\gamma}(\boldsymbol{z}; \mathcal{T}\boldsymbol{w})Q(\boldsymbol{z}; \mathcal{T}\boldsymbol{w})} \bigg| \nabla_{\boldsymbol{w}} \nabla_{\boldsymbol{w}}^{T} \log p(\boldsymbol{a}|\boldsymbol{s}; \mathcal{T}\boldsymbol{w}) \bigg|,$$

which gives the following

$$\begin{split} \tilde{\mathcal{H}}_{2}(\boldsymbol{w}) &= \mathbb{E}_{p_{\gamma}(\boldsymbol{z};\mathcal{T}\boldsymbol{w})Q(\boldsymbol{z};\mathcal{T}\boldsymbol{w})} \bigg[ \nabla_{\boldsymbol{w}} \nabla_{\boldsymbol{w}}^{T} \log p(\boldsymbol{a}|\boldsymbol{s};\mathcal{T}\boldsymbol{w}) \bigg], \\ &= \mathbb{E}_{p_{\gamma}(\boldsymbol{z};\mathcal{T}\boldsymbol{w})Q(\boldsymbol{z};\mathcal{T}\boldsymbol{w})} \bigg[ \mathcal{T}^{T} \nabla_{\boldsymbol{v}} \nabla_{\boldsymbol{v}}^{T} \log p(\boldsymbol{a}|\boldsymbol{s};\boldsymbol{v}) \mathcal{T} \bigg], \\ &= \mathcal{T}^{T} \mathbb{E}_{p_{\gamma}(\boldsymbol{z};\mathcal{T}\boldsymbol{w})Q(\boldsymbol{z};\mathcal{T}\boldsymbol{w})} \bigg[ \nabla_{\boldsymbol{v}} \nabla_{\boldsymbol{v}}^{T} \log p(\boldsymbol{a}|\boldsymbol{s};\boldsymbol{v}) \bigg] \mathcal{T}, \\ &= \mathcal{T}^{T} \mathcal{H}_{2}(\boldsymbol{v}) \mathcal{T}. \end{split}$$

Using these two expressions we have that the parameter updates, under the approximate Newton method, of the objective functions U and  $\tilde{U}$  are related as follows

$$egin{aligned} & m{v}_{\mathsf{new}} = m{v} + lpha \mathcal{H}_2(m{v})^{-1} 
abla_{m{v}} U(m{v}), \ & = \mathcal{T}ig(m{w} + lpha ilde{\mathcal{H}}_2(m{w})^{-1} 
abla_{m{v}} ilde{U}(m{w})ig) \end{aligned}$$

where  $\alpha$  is some step-size parameter. This shows that the approximate Newton method is affine invariant. The invariance of the diagonal approximate Newton method to rescaling of the parameters follows similarly.

### 6.5 Recurrent State Search Direction Evaluation for Approximate Newton Method

In [31] a sampling algorithm was provided for estimating the gradient of an infinite horizon MDP with average rewards. This algorithm makes use of a recurrent state, which we denote by  $s^*$ . In algorithm(1) we detail a straightforward extension of this algorithm to the estimation the approximate Hessian,  $\mathcal{H}_2(w)$ , in this MDP framework. The analogous algorithm for the estimation of the diagonal matrix,  $\mathcal{D}_2(w)$ , follows similarly. In algorithm(1) we make use of an *eligibility trace* for both the gradient and the approximate Hessian, which we denote by  $\Phi^1$  and  $\Phi^2$  respectively. The estimates (up to a positive scalar) of the gradient and the approximate Hessian are denoted by  $\Delta^1$  and  $\Delta^2$  respectively.

### 6.6 Tetris Experiment Specification

In this section we give a detailed specification of the procedure used for each of the algorithms in the tetris experiment. The same general procedure was used for all the algorithms considered in the experiments. We modelled the environment through with an infinite planning horizon with average rewards. The reward at each time-point is equal to the number of lines deleted. We used a recurrent state formulation [31] of the gradient of the average reward framework to perform the gradient evaluation. We used analogous versions of this recurrent state formulation for natural gradient ascent, the diagonal approximate Newton method and the full approximate Newton method. See section(6.5) of the supplementary material for more details. As in [16] we used the sample trajectories obtained during the gradient evaluation to estimate the Fisher information matrix. Irrespective of the policy a game of tetris is guaranteed to terminate after a finite number of turns, see e.g. [7]. As we are considering the average reward framework a new game starts when the previous one terminates. We used the empty board as a recurrent state where, because a new game starts with an empty board, this state is recurrent.<sup>1</sup> During each training iteration the (approximation of the) search direction was obtained by sampling 1000 games, where these games were sampled using the current policy parameters. Given the current approximate search direction the following basic line search method was used to obtain a step-size: For every step-size in a given finite set of step-sizes sample a set number of games and then return the step-size with the maximal score over these games. In practice, in order to reduce the susceptibility to random noise, we used the same simulator seed for each possible step-size in the set. To avoid over-fitting a different simulator seed was used during each training iteration. In this line search procedure we sampled 1000 games for each of the possible

<sup>&</sup>lt;sup>1</sup>This is actually an approximation because it doesn't take into account that the state is given by the configuration of the board and the current piece, so this particular 'recurrent state' ignores the current piece. Empirically we found that this approximation gave better results, presumably due to reduced variance in the estimands, and there is no reason to believe that it is unfairly biasing the comparison between the various parametric policy search methods.

**Algorithm 1** Recurrent state sampling algorithm to estimate the search direction of the approximate Newton method when applied to a MDP with an infinite planning horizon with average rewards.

Sample a state from the initial state distribution:

$$s_1 \sim p_1(\cdot).$$

for t = 1, ..., N, for some  $N \in \mathbb{N}$ , do

Given the current state, sample an action from the policy:

$$\boldsymbol{a}_t \sim \pi(\cdot | \boldsymbol{s}_t; \boldsymbol{w}).$$

if  $s_t \neq s^*$  then Update the eligibility traces:

opuate the englotinty traces.

$$\mathbf{\Phi}^1 \leftarrow \mathbf{\Phi}^1 + 
abla_{oldsymbol{w}} \log \pi(oldsymbol{a}_t;oldsymbol{w}) \qquad \mathbf{\Phi}^2 \leftarrow \mathbf{\Phi}^2 + 
abla_{oldsymbol{w}} 
abla_{oldsymbol{w}}^T \log \pi(oldsymbol{a}_t|oldsymbol{s}_t;oldsymbol{w})$$

else

reset the eligibility traces:

$$\mathbf{\Phi}^1 = \mathbf{0}, \qquad \mathbf{\Phi}^2 = \mathbf{0}$$

end if

Update the estimates of the gradient and the approximate Hessian:

$$\mathbf{\Delta}^1 \leftarrow \mathbf{\Delta}^1 + R(\boldsymbol{a}_t, \boldsymbol{s}_t) \mathbf{\Phi}^1, \qquad \qquad \mathbf{\Delta}^2 \leftarrow \mathbf{\Delta}^2 + R(\boldsymbol{a}_t, \boldsymbol{s}_t) \mathbf{\Phi}^2$$

Sample state from the transition dynamics:

$$s_{t+1} \sim p(\cdot | \boldsymbol{a}_t, \boldsymbol{s}_t).$$

end for

Return the estimated gradient and approximate Hessian, which up to a positive scaling are given by  $\Delta^1$  and  $\Delta^2$  respectively.

step-sizes. The same set of step-sizes was used in all of the different training algorithms considered in section(4), where we used the set

$$\{0.1, 0.5, 1.0, 2.0, 4.0, 8.0, 16.0, 32.0, 64.0, 128.0\}.$$

To reduce the amount of noise in the results we used the same set of simulator seeds in the search direction evaluation for each of the algorithms considered in section(4). In particular, we generated a  $n_{\text{experiments}} \times n_{\text{iterations}}$  matrix of simulator seeds, where  $n_{\text{experiments}}$  was the number of repetitions of the experiment and  $n_{\text{iterations}}$  was the number of training iterations in each experiment. We then used this one matrix in all of the different training algorithms, where the element in the  $j^{\text{th}}$  column and  $i^{\text{th}}$  row corresponds to the simulator seed used in the  $j^{\text{th}}$  training iteration of the  $i^{\text{th}}$  experiment. In a similar manner the set of simulator seeds used for the line search procedure was the same for all of the different training algorithms. Finally, to make the line search consistent among all of the different training algorithms the search direction was normlised and the resulting unit vector was the vector used in the line search procedure.

#### 6.7 Linear System Specification

The *N*-link rigid robot arm manipulator is a standard continuous MDP model, consisting of an end effector connected to an *N*-linked rigid body [17]. A typical continuous control problem for such systems is to apply appropriate torque forces to the joints of the manipulator so as to move the end effector into a desired position. The state of the system is given by  $q, \dot{q}, \ddot{q} \in \mathbb{R}^N$ , where  $q, \dot{q}$  and  $\ddot{q}$  denote the angles, velocities and accelerations of the joints respectively, while the control variables

are the torques applied to the joints  $\tau \in \mathbb{R}^N$ . The nonlinear state equations of the system are given by, see *e.g.* [17],

$$M(\boldsymbol{q})\ddot{\boldsymbol{q}} + C(\dot{\boldsymbol{q}}, \boldsymbol{q})\dot{\boldsymbol{q}} + g(\boldsymbol{q}) = \boldsymbol{\tau}$$
(16)

where M(q) is the inertia matrix,  $C(\dot{q}, q)$  denotes the Coriolis and centripetal forces and g(q) is the gravitational force. While this system is highly nonlinear it is possible to define an appropriate control function  $\hat{\tau}(q, \dot{q})$  that results in linear dynamics in a different state-action space. This process is called *feedback linearisation*, see *e.g.* [17], and in the case of an *N*-link rigid manipulator recasts the torque action space into the acceleration action space. This means that the state of the system is now given by q and  $\dot{q}$ , while the control is  $a = \ddot{q}$ . Ordinarily in such problems the reward would be a function of the generalised co-ordinates of the end effector, which results in a non-trivial reward function in terms of q,  $\dot{q}$  and  $\ddot{q}$ . While this reward function can be modelled as a mixture of Gaussians for simplicity we consider the simpler problem where the reward is a function of q,  $\dot{q}$  and  $\ddot{q}$  directly.

In the experiments we considered a 3-link rigid manipulator, which results in a 9-dimensional stateaction space and a 22-dimensional policy. In the experiment we discretised the continuous time dynamics into time-steps of  $\Delta_t = 0.1$  and considered a finite planning horizon of H = 100. The mean of the initial state distribution was set to zero. The elements of the policy parameters m and  $\pi_{\sigma}$  were initialised randomly from the interval [-2, 2] and [1, 2] respectively. The matrix K was initialised to be zero on inter-link entries, while intra-link entries were initialised using rejection sampling. We sampled the parameters for each link independently from the set  $[-400, 40] \times [-50, 10]$  and rejected the sample if the corresponding link was unstable. In the reward function the desired angle of each joint was randomly sampled from the interval  $[\pi/4, 3\pi/4]$ . The covariance matrices of the initial state distribution and state transition dynamics were set to diagonals, where the diagonal elements were initialised randomly from the interval [0, 0.05]. The covariance matrix of the reward function was set to be a diagonal with all entries equal to 0.1.

We used the minFunc<sup>2</sup> optimisation library in all of the gradient-based algorithms. We found that both the line search algorithm and the step-size initialisation had an effect on the performance of all the algorithms. We therefore tried various combinations of these settings for each algorithm and selected the one that gave the best performance. We tried bracketing line search algorithms with: *stepsize halving*; *quadratic/cubic interpolation from new function values*; *cubic interpolation/extrapolation from new function and gradient values*; *step-size doubling and bisection*; *cubic interpolation/extrapolation with function and gradient values*. We tried the following step-size initialisations: *quadratic initialization using previous function value and new function value/gradient*; *twice the previous step-size*. To handle situations where the initial policy parameterisation was in a 'flat' area of the parameter space far from any optima we set the function and point toleration of minFunc to zero for all algorithms. To handle the large number of training iterations and training iterations to 10,000.

#### 6.8 Non-Linear System Specification

In section we detail the implementation of the experiment on the synthetic two-dimensional nonlinear MDP considered in [29]. The state-space of the problem is two-dimensional,  $s = (s^1, s^2)$ , where  $s^1$  is the agent's position and  $s^2$  is the agent's velocity. The control is one-dimensional and the dynamics of the system is given as follows

$$\begin{split} s^1_{t+1} &= s^1_t + \frac{1}{1+e^{-u_t}} - 0.5 + \kappa, \\ s^2_{t+1} &= s^2_t - 0.1s^1_{t+1} + \kappa, \end{split}$$

where  $\kappa$  is a zero-mean Gaussian random variable with standard deviation  $\sigma_{\kappa} = 0.02$ . The agent starts in the state s = (0, 1), with the addition of Gaussian noise with standard deviation 0.001, and the objective is to transport the agent to the state (0, 0). We use the same policy as in [29], which is given by  $a_t = (w + \epsilon_t)^T s_t$ , where w are the control parameters and  $\epsilon_t \sim \mathcal{N}(\epsilon_t; \mathbf{0}, \sigma_{\epsilon}^2 I)$ . The objective function is non-trivial for  $w \in [0, 60] \times [-8, 0]$ . In the experiment the initial control parameters were sampled from the region  $w_0 \in [0, 60] \times [-8, 0]$ . We considered a finite planning

<sup>&</sup>lt;sup>2</sup>This software library is freely available at http://www.di.ens.fr/~mschmidt/Software/ minFunc.html.

horizon with a planning horizon of H = 80. We used forward sampling to perform the inference, where in all algorithms 50 trajectories were sampled during each training iteration.

We also detail the procedure used for the step-size tuning for the approximate Newton method and natural gradient ascent. The tuning of the step-size sequence in steepest gradient ascent was similar in nature to natural gradient ascent and so is omitted from the discussion. Using the intuition that the approximate Newton method has a natural step-size of one, which corresponds to an EM-step in this problem because the log-policy is quadratic in the control parameters, we considered step-size sequences of the form  $\alpha_k = (1 - \frac{k}{N})\alpha + \frac{k}{N}$ , where N is the total number of training iterations considered and  $\alpha \in \mathbb{R}^+$ . In the experiment we considered the values  $\alpha = 1$ , 6, 12, 18 and 24. The intuition used in this selection was that, provided that the steps were not so large as to cause overshooting in the parameter space, larger steps will increase performance. In steepest and natural gradient ascent we used step-size sequences that satisfied the Robbins-Munro conditions. In both of these methods it was necessary to obtain a gauge of a reasonable scale of a good step-size sequence. For this reason in the experiment we considered step-size sequences of the form  $\alpha_k = \frac{\alpha}{\sqrt{k}}$  with  $\alpha = 0.0001, 0.001, 0.01, 0.1, 1, 2$  and 4. It was found that the sequence  $\alpha = 18$  gave the best results for natural gradient ascent.