

Robot Learning

Winter Semester 2017/2018, Homework 4

Prof. Dr. J. Peters, D. Tanneberg, M. Ewerton



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Total points: 60 + 10 bonus

Due date: Wednesday, 31 January 2018 (before the lecture)

Name, Surname, ID Number

Markus Lamprecht, 2424163

Moritz Knaust, 2430801

Problem 4.1 Expectation-Maximization [25 Points + 5 Bonus]

In this exercise your task is to control a 2-DoF planar robot to throw a ball at a specific target. You will use an episodic setup, where you first specify the parameters of the policy, evaluate them on the simulated system, and obtain a reward. The robot will be controlled with the Dynamic Motor Primitives (DMPs). The goal state of the DMPs is pre-specified and the weights of the DMP $\theta_i, i = 1 \dots 10$ are the open parameters of the control policy. Each DoF of the robot is controlled with a DMP with five basis functions. The ball is mounted at the end-effector of the robot and gets automatically released at time step t_{rel} . We define a stochastic distribution $\pi(\theta|\omega) = \mathcal{N}(\theta|\mu, \Sigma)$, with $\omega = \{\mu, \Sigma\}$.

Your task is to update the parameters ω of the policy using EM to maximize the expected return. In this exercise we will not modify the low-level control policy (the DMPs) of the robot.

A template for the simulation of the 2-DoF planar robot and plotting functions can be found at the course website. For the programming exercises, attach snippets of your code.

a) Analytical Derivation [5 Points]

Using the weighted ML estimate,

$$\omega_{k+1} = \arg \max_{\omega} \left\{ \sum_i w^{[i]} \log \pi(\theta^{[i]}; \omega) \right\}, \quad (1)$$

derive analytically the update rule for our policy $\pi(\theta|\omega)$ for the mean μ . Show your derivations.

$$\frac{\partial \omega_{k+1}}{\partial \mu} = 0 \quad (7)$$

$$\pi(\theta^{[i]}; w) = N(\theta|\mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^k |\Sigma|}} \exp\left(-\frac{1}{2}(\mu - \theta)^T \Sigma^{-1}(\mu - \theta)\right) \quad (8)$$

$$\log(\pi(\theta^{[i]}; w)) = -\log\left(\sqrt{(2\pi)^k |\Sigma|}\right) - \frac{1}{2}(\mu - \theta)^T \Sigma^{-1}(\mu - \theta) \quad (9)$$

$$\frac{\partial}{\partial \mu} \log(\pi(\theta^{[i]}; w)) = -\frac{1}{2}(2\Sigma^{-1}\mu - 2\Sigma^{-1}\theta) \quad (10)$$

$$\frac{\partial \omega_{k+1}}{\partial \mu} = 0 \rightarrow \sum_i w^{[i]} \Sigma^{-1} \mu_i = \sum_i w^{[i]} \Sigma^{-1} \theta_i \rightarrow \mu_i = \frac{\sum_i w^{[i]} \theta_i}{\sum_i w^{[i]}} \quad (11)$$

b) Programming Exercise [15 Points]

Implement the EM algorithm using the provided framework. Your goal is to throw the ball at $\mathbf{x} = [2, 1]m$ at time $t = 100$. Use the weighted Maximum Likelihood (ML) solutions for updating the parameters of the policy. For the mean, use the equation derived at the previous exercise. For the covariance, the update rule is

$$\Sigma_{\text{new}} = \frac{\sum_i w^{[i]} (\boldsymbol{\theta}^{[i]} - \boldsymbol{\mu}_{\text{new}})(\boldsymbol{\theta}^{[i]} - \boldsymbol{\mu}_{\text{new}})^T}{\sum_i w^{[i]}}. \quad (12)$$

To calculate the weights w_i for each sample, transform the returned rewards by

$$w^{[i]} = \exp((R^{[i]}(\boldsymbol{\theta}) - \max(\mathbf{R}))\beta) \quad (13)$$

and normalize them, where the vector $\mathbf{R} \in \mathbb{R}^{N \times 1}$ is constructed from the rewards of the N samples. The parameter β is set to

$$\beta = \frac{\lambda}{\max(\mathbf{R}) - \min(\mathbf{R})}, \quad (14)$$

where $\lambda = 7$. Start with the initial policy

$$\pi(\boldsymbol{\theta}|\boldsymbol{\omega}) = \mathcal{N}(\boldsymbol{\theta}|\mathbf{0}, 10^6 \mathbf{I}) \quad (15)$$

and calculate the average reward at each iteration. Iterate until convergence or for a maximum number of 100 iterations. We assume that convergence is achieved when the average return does not change much at every iteration, i.e.,

$$| \langle \mathbf{R}_{i-1} \rangle - \langle \mathbf{R}_i \rangle | < 10^{-3},$$

where $\langle \cdot \rangle$ denotes the average operator. At each iteration use $N = 25$ samples. In order to avoid getting stuck to a local optimum, we force the algorithm to explore a bit more by adding a regularization factor to the covariance of the policy,

$$\Sigma'_{\text{new}} = \Sigma_{\text{new}} + \mathbf{I}. \quad (16)$$

What is the average return of the final policy? Use `animate_fig` from `Pend2dBallThrowDMP.py` to show how the final policy of the robot looks like (attach all five screenshots).

Average Return of the final policy: -38.89

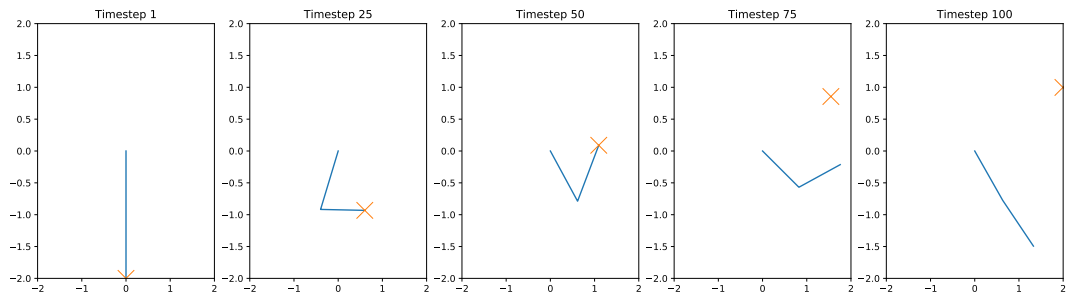


Figure 2: Policy of the Robot's arm (blue) and the Ball's position (yellow cross)

```

1 numDim = 10
2 numSamples = 25
3 maxIter = 100
4 numTrials = 10
5
6 def EM(lam = 7):

```

```

7   # For example, let initialize the distribution...
   Mu_w = np.zeros(numDim)
9   Sigma_w = np.eye(numDim) * 1e6

11  R_mean_old = 0
   # Do your learning
13  R_avg = np.zeros(maxIter)
   for i in xrange(maxIter):
15     # reward
       R = np.zeros((numSamples, 1))
       theta = np.zeros((numSamples, numDim))
17     for j in xrange(numSamples):
         theta[j, :] = np.random.multivariate_normal(Mu_w, Sigma_w)
         R[j] = env.getReward(theta[j, :])
21     R_avg[i] = np.mean(R)

23     # weights
       beta = lam / (np.max(R) - np.min(R))
       w = np.exp((R - np.max(R)) * beta)

27     Mu_w = np.sum(w * theta, axis=0) / np.sum(w)

29     temp = np.zeros((numDim, numDim))

31     # for k in xrange(numDim):
       #     diff = theta - Mu_w[k]

33     for j in xrange(numSamples):
         diff = (theta[j, :] - Mu_w).reshape((numDim, 1))
         temp += w[j] * np.dot(diff, diff.T)
35     Sigma_w = temp / np.sum(w) + np.eye(numDim)

37     if i > 1:
         if np.abs(R_avg[i] - R_avg[i-1]) < 1e-3:
39             break
41     print 'end'
43     return Mu_w, Sigma_w, R_avg

```

Listing 2: EM algorithm

c) Tuning The Temperature [5 Points]

Repeat the learning 10 times and plot the mean of the average return of all runs with 95% confidence for temperatures $\lambda = 25$, $\lambda = 7$ and $\lambda = 3$. How does the value of λ affect the convergence of the algorithm in theory and what do you observe in the results? Use the logarithmic scale for your plot.

In the results we observe, that for $\lambda = 25$ the average return converges faster than for smaller lambdas. Moreover it should be noted that for $\lambda = 7$ we have the smallest average return at episode=100.

In Theory: With the temperature λ one can select how the maximum of the weights is selected:

$$\lim_{\lambda \rightarrow \infty} w^{[i]} : \text{weights are selected deterministically} \quad (19)$$

$$\lim_{\lambda \rightarrow 0} w^{[i]} : \text{weights are selected stochastically} \quad (20)$$

In the plots we can prove that for small lambdas the convergence is slower because the weights are selected stochastically. However the mean of the average reward for the last episode should be the smallest if the lambda is very small. This theoretical fact differs from the plot above.

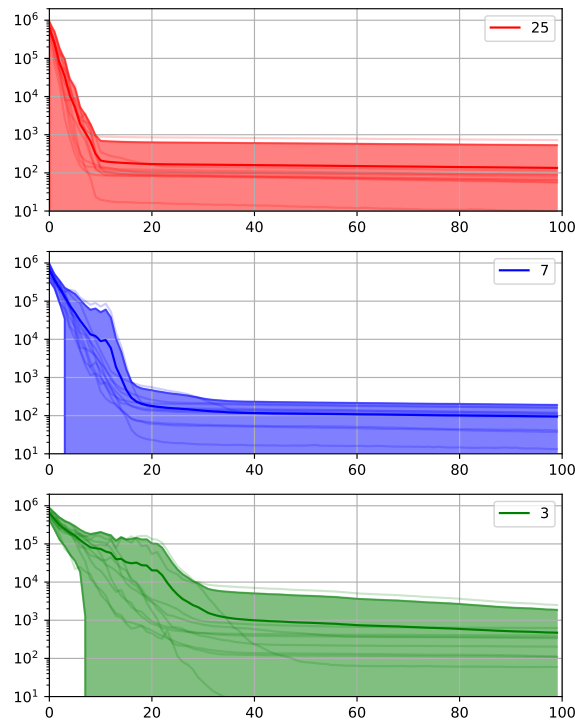


Figure 3: The average return of different λ over the episodes of the EM-Algorithm.

d) Optimal Temperature [5 Bonus Points]

The problem of choosing λ falls under a more general RL problem. Which one? Which algorithm would allow you to automatically select the optimal λ ? Explain it with your words. Does it still have hyperparameters to be set?

RL problem: Policy Search Methods using Policy Gradients finding the right weights scaling factor.

Algorithm that selects automatically the right lambda: Reward Weighted Regression.

The Reward Weighted Regression

is an algorithm that generalizes the EM algorithm. This algorithm reduces the problem of learning with immediate rewards to a reward-weighted regression problem with an adaptive, integrated reward transformation for faster convergence. The resulting algorithm is efficient, learns smoothly without dangerous jumps in solution space, and works well in applications of complex high degree-of-freedom robots.

Hyperparameters:

A transformed Hyperparameter still has to be set in RWR. However this parameter is automatically updated. With this parameter the convergence speed can be selected.

Problem 4.2 Policy Gradient [25 Points + 5 Bonus]

In this exercise, you are going to solve the same task of the previous exercise but using policy gradient. For the programming exercises, you are allowed to reuse part of the EM code. Attach snippets of your code.

a) Analytical Derivation [5 Points]

You have a Gaussian policy with diagonal covariance, i.e., $\pi(\boldsymbol{\theta}|\boldsymbol{\omega}) = \mathcal{N}(\boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2))$, where $\boldsymbol{\omega} = [\boldsymbol{\mu}, \boldsymbol{\sigma}]$. Compute analytically the gradient of the logarithm of the policy with respect to the parameters $\boldsymbol{\omega}$, i.e., $\nabla_{\boldsymbol{\omega}} \log \pi(\boldsymbol{\theta}|\boldsymbol{\omega})$. (Hint: consider the properties of the diagonal covariance matrix.)

$$\nabla_{\boldsymbol{\omega}} \log \pi(\boldsymbol{\theta}|\boldsymbol{\omega}) = \nabla_{\boldsymbol{\omega}} \left(-\log \left(\sqrt{(2\pi)^k} |\text{diag}(\boldsymbol{\sigma}^2)| \right) - \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\mu})^T \text{diag}(\boldsymbol{\sigma}^2)^{-1} (\boldsymbol{\theta} - \boldsymbol{\mu}) \right) \quad (28)$$

$$\nabla_{\boldsymbol{\omega}} \log \pi(\boldsymbol{\theta}|\boldsymbol{\omega}) = \nabla_{\boldsymbol{\omega}} \left(-\log \left(\sqrt{(2\pi)^k} (\boldsymbol{\sigma}^{2k}) \right) - \frac{1}{2\boldsymbol{\sigma}^2} \sum_{i=1}^k (\theta_i - \mu_i)^2 \right) \quad (29)$$

$$\nabla_{\boldsymbol{\omega}} \log \pi(\boldsymbol{\theta}|\boldsymbol{\omega}) = \begin{bmatrix} \frac{\partial}{\partial \boldsymbol{\mu}} \log \pi(\boldsymbol{\theta}|\boldsymbol{\omega}) \\ \frac{\partial}{\partial \boldsymbol{\sigma}} \log \pi(\boldsymbol{\theta}|\boldsymbol{\omega}) \end{bmatrix} \quad (30)$$

$$\frac{\partial}{\partial \boldsymbol{\mu}} \log \pi(\boldsymbol{\theta}|\boldsymbol{\omega}) = \frac{1}{\boldsymbol{\sigma}^2} \sum_{i=1}^k (\theta_i - \mu_i) \quad (31)$$

$$\frac{\partial}{\partial \boldsymbol{\sigma}} \log \pi(\boldsymbol{\theta}|\boldsymbol{\omega}) = -\frac{2k}{\boldsymbol{\sigma}} + \frac{1}{\boldsymbol{\sigma}^3} \sum_{i=1}^k (\theta_i - \mu_i)^2 \quad (32)$$

Note that:

$$\mathcal{N}(\boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2)) = \frac{1}{\sqrt{(2\pi)^k} |\text{diag}(\boldsymbol{\sigma}^2)|} \exp \left(-0.5 (\boldsymbol{\theta} - \boldsymbol{\mu})^T \text{diag}(\boldsymbol{\sigma}^2)^{-1} (\boldsymbol{\theta} - \boldsymbol{\mu}) \right) \quad (33)$$

$$(\boldsymbol{\theta} - \boldsymbol{\mu})^T \text{diag}(\boldsymbol{\sigma}^2)^{-1} (\boldsymbol{\theta} - \boldsymbol{\mu}) = \frac{1}{\boldsymbol{\sigma}^2} \sum_{i=1}^k (\theta_i - \mu_i)^2 \quad (34)$$

b) Programming Exercise [5 Points]

Consider the same robotic task as in the previous exercise and this time solve it with policy gradient. Use an initial mean of $\mu_0 = [0 \dots 0]$ and a fixed $\sigma = \text{diag}([10 \dots 10])$ (i.e., do not update σ). Set the learning rate to $\alpha = 0.1$ and use the same hyperparameters as before (25 episodes sampled for each iteration and max 100 iterations).

Repeat the learning 10 times and plot the mean of the average return of all runs with 95% confidence. Use the logarithmic scale for your plot. Comment your results.

$$\omega_{k+1} = \omega_k + \alpha \nabla_{\omega} J_{\omega} \quad (38)$$

Update rule for θ :

$$\sigma_{k+1} = \sigma_k + \alpha \frac{\partial}{\partial \sigma} J(\sigma_k) \quad (39)$$

Update rule for μ :

$$\mu_{k+1} = \mu_k + \alpha \frac{\partial}{\partial \mu} J(\mu_k) \quad (40)$$

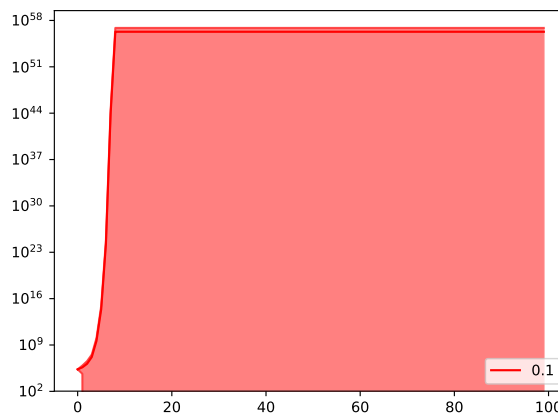


Figure 5: The average return over the episodes for $\lambda = 0.1$

You can see, that the reward explodes and one does not get a stable behavior. This depends on the high variance of the gradient.

```

numDim = 10
numSamples = 25
maxIter = 100
numTrials = 10

# Do your learning
def PolicyGradient(alpha = 0.1):
    # For example, let initialize the distribution...
    Mu = np.zeros(numDim)
    Sigma = np.eye(numDim) * 10

    # Do your learning
    R_avg = np.zeros(maxIter)
    for i in xrange(maxIter):
        # reward
        R = np.zeros((numSamples, 1))
        theta = np.zeros((numSamples, numDim))
        for j in xrange(numSamples):
            theta[j, :] = np.random.multivariate_normal(Mu, Sigma)

```

```

20         R[j] = env.getReward(theta[j, :])
21     R_avg[i] = np.mean(R)
22     #print 'R_mean_{}: {}'.format(i, R_avg)
23
24     temp = np.zeros(theta.shape)
25     for j in xrange(numSamples):
26         temp[j, :] = np.dot(np.linalg.inv(Sigma), theta[j] - Mu) * R[j]
27         # gradient
28     grad_mu = 1./numSamples * np.sum(temp, axis=0)
29
30     Mu = Mu + alpha*grad_mu
31
32     if i>1:
33         if np.abs(R_avg[i] - R_avg[i-1]) < 1e-3:
34             break
35     print 'R_avg_end: {}'.format(R_avg[i-1])
36     return Mu, Sigma, R_avg

```

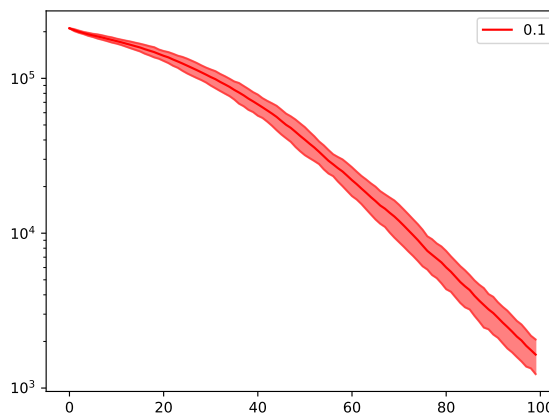
Listing 4: policy gradient algorithm

c) A Little Trick [5 Points]

How would you improve the above implementation? (Beside using a smaller or adaptive learning rate, or using the natural gradient). What is the theory behind this “trick”? Repeat the learning and discuss the results.

$$\nabla_{\omega} J_{\omega} = \sum_{i=1}^N \nabla_{\omega} \log \pi(\theta; \omega) (R_i - b) \quad (42)$$

The baseline has to be subtracted. Thereby the possibly high gradient of the variance can be reduced. It is still unbiased. As baseline we use the average reward.

Figure 7: The average return over the episodes of for $\lambda = 0.1$ with baseline subtracted of the PG algorithm

```

1 numDim = 10
2 numSamples = 25
3 maxIter = 100
4 numTrials = 10
5
6 # Do your learning
7 def PolicyGradient(alpha = 0.1):
8     ...
9     for i in xrange(maxIter):
10         ...
11
12         temp = np.zeros(theta.shape)
13         for j in xrange(numSamples):
14             temp[j, :] = np.dot(np.linalg.inv(Sigma), theta[j] - Mu) * (R[j] - R_avg[i
15         ])
16             # gradient
17         grad_mu = 1./numSamples * np.sum(temp, axis=0)
18
19         Mu = Mu + alpha*grad_mu
20         ...
21     print 'R_avg_end: {}'.format(R_avg[-1])
22     return Mu, Sigma, R_avg

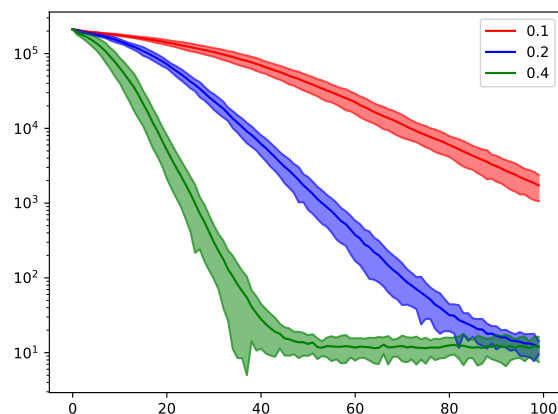
```

Listing 6: policy gradient algorithm with subtracted baseline

d) Learning Rate [5 Points]

Repeat the optimization changing the learning rate to $\alpha = 0.4$ and $\alpha = 0.2$ (keep the trick of the previous exercise). Plot in one figure the mean of the average returns for all α with 95% confidence. How does the value of α affect the convergence of the algorithm? Use the logarithmic scale for your plot.

The higher the α the faster the algorithm converges.

Figure 9: The average return over the episodes for different learning rates α

e) Variable Variance [5 Points]

Try to improve the optimization process by learning also the variance σ . Is it easier or harder to learn also the variance? Why?

Without using the natural gradient, tune the learning process to achieve better results. If you think it is necessary, you can impose a lower bound to avoid that the variance collapses to infinitely small values (e.g., if $\sigma(i) < \sigma_{\text{lower}}$ then $\sigma(i) = \sigma_{\text{lower}}$). In one figure, plot the learning trend with confidence interval as done before and compare it to the one achieved with $\alpha = 0.4$ before.

Now use:

$$\sigma_{k+1} = \sigma_k + \alpha \frac{\partial}{\partial \omega} J(\sigma_k) \quad (44)$$

f) Natural Gradient [5 Bonus Points]

Write down the equation of the natural gradient. What is the theory behind it? Is it always easy to use?

The Natural Gradient: (with G the Fischer Information Matrix)

$$\nabla_{\omega} J = G(\omega)^{-1} \nabla_{\omega} J \quad (46)$$

Theory of the Natural Gradient:

KL divergence is a measure of how close a distribution is to another distribution. This brings us to the natural gradient. If we blindly update our network in the direction of its gradients, there are no guarantees the distribution of the new network will be similar to the old one.

To fix this, we first consider all combinations of parameters that result in a new network a constant KL divergence away from the old network. This constant value can be viewed as the step size or learning rate. Out of all these possible combinations, we choose the one that minimizes our loss function.

Basically, we're adding in a constraint to our update, that the new network will behave relatively similar to our old one. Our step size corresponds directly to the actual distribution of the network, not its parameter's values.

This comes with the added benefit of a more stable learning process. Especially when we're updating using a randomly sampled batch, some outliers may make drastic changes to a network's parameters. With natural gradient descent, a single update can't make too much of an impact.

Note that: Replacing a sigmoid activation with a tanh function would change the standard gradient, but not the natural gradient.

Is it always easy to use: No actually NG is much more computationally expensive compared to the simple gradient. Moreover there are two types of KL:

- Moment projection
- Information projection (is not unique selects the mode!)

It is not so easy to select which KL should be used.

Note that KL can be approximated by the FIM which captures information how the single parameters influence the distribution.

Problem 4.3 Reinforcement Learning [10 Points]

a) RL Exploration Strategies [10 Points]

In which spaces can you perform exploration in RL? Discuss the two exploration strategies applicable to RL.

Exploration can be performed in Action Space and Parameter space.

Parameter Space:

In Episode Based Policy Search a distribution is learned over the parameters of the low level control policy. Exploration in parameter space allows more sophisticated strategies and is often very efficient for a small amount of parameters.

A problem is the high variance of the returns, because one add different random variables. If one has a too high variance, the learned policy can become instable. To avoid this, you can use the baseline trick for example.

An advantage is, that you do not have to know the structure of the optimization problem. So it is an "Black-Box Optimizer".

Action Space:

In the Episode based Policy Search you learn an full movement as episode and then get the reward for it. In Step-based Policy Search one decompose the episode in single timesteps. So we get the reward for single state-actions pairs. This reward has a lower variance.

In Step Based Policy Search there is less variance in quality assessment. Moreover as the exploration happens in action space it is less likely to create unstable policies.