# 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, ???

---

### 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 $\theta_i, i = 1 \ldots 10$ 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 $t_{\mathrm{rel}}$. We define a stochastic distribution $\pi(\boldsymbol{\theta}|\boldsymbol{\omega}) = \mathcal{N}(\boldsymbol{\theta}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$, with $\boldsymbol{\omega} = \{\boldsymbol{\mu}, \boldsymbol{\Sigma}\}$.

Your task it to update the parameters $\boldsymbol{\omega}$ of the policy using EM to maximize the expected return. In this exercises 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,

$$\boldsymbol{\omega}_{k+1} = \arg\max_{\boldsymbol{\omega}}\{\sum_i w^{[i]} \log \pi(\boldsymbol{\theta}^{[i]}; \boldsymbol{\omega})\}, \tag{1}$$

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

$$\frac{\partial \boldsymbol{w}_{k+1}}{\partial \boldsymbol{\mu}} = 0 \tag{7}$$

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

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

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

$$\frac{\partial \boldsymbol{w}_{k+1}}{\partial \boldsymbol{\mu}} = 0 \rightarrow \sum_i w^{[i]}\boldsymbol{\Sigma}^{-1}\mu_i = \sum_i w^{[i]}\boldsymbol{\Sigma}^{-1}\theta_i \rightarrow \mu_i = \frac{\sum_i w^{[i]}\theta_i}{\sum_i w^{[i]}} \tag{11}$$

b) Programming Exercise [15 Points]

Implement the EM algorithm using the provided framework. Your goal is to throw the ball at $\boldsymbol{x} = [2, 1]m$ 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

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

To calculate the weights $\boldsymbol{w}_i$ for each sample, transform the returned rewards by

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

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

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

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

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

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.,

$$| < R_{i-1} > - < R_i > | < 10^{-3},$$

where $< \cdot >$ 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,

$$\Sigma'_{\text{new}} = \Sigma_{\text{new}} + \boldsymbol{I}. \tag{16}$$

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).

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 $\lambda$ affect the convergence of the algorithm in theory and what do you observe in the results? Use the logarithmic scale for your plot.

d) Optimal Temperature [5 Bonus Points]

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

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.)

hallo

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 $\boldsymbol{\mu}_0 = [0 \dots 0]$ and a fixed $\boldsymbol{\sigma} = \text{diag}([10 \dots 10])$ (i.e., do not update $\boldsymbol{\sigma}$). 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.

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.

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 $\alpha$ with 95% confidence. How does the value of $\alpha$ affect the convergence of the algorithm? Use the logarithmic scale for your plot.

e) Variable Variance [5 Points]

Try to improve the optimization process by learning also the variance $\boldsymbol{\sigma}$. 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.

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?

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.