

---

# Homework

Stochastic Models and Adaptive Algorithms

---

Réda Vince

Z697LX

January 6, 2021

# Contents

<b>1</b>	<b>Linear regression</b>	<b>2</b>
1.1	Function approximation with least squares . . . . .	2
1.1.1	Ordinary least-squares . . . . .	2
1.1.2	Recursive least-squares . . . . .	2
1.1.3	Least-norm problem . . . . .	3
1.2	Approximating auto-regressive series . . . . .	4
<b>2</b>	<b>Kernel methods</b>	<b>5</b>
<b>3</b>	<b>Reinforcement learning</b>	<b>6</b>
3.1	The environment . . . . .	6
3.2	Model based methods . . . . .	6
3.2.1	Model generating . . . . .	6
3.2.2	Linear programming . . . . .	6
3.2.3	Value iteration . . . . .	6
3.2.4	Policy iteration . . . . .	7
3.3	Model-free methods . . . . .	7

# 1 Linear regression

## 1.1 Function approximation with least squares

### 1.1.1 Ordinary least-squares

The best linear approximation of a function can be calculated using least-squares.

At first, a noisy sample of  $(x, y)$  pairs is generated such that  $y_i = c x_i \sin(c x_i) + \epsilon_i$ , where  $\epsilon_i \sim \mathcal{N}(0, 1)$ .

Let  $[\Phi]_{ij} = f_j(x_i)$  be the transformed input vector and  $\mathbf{y}$  the output, where  $f_j$  is a basis function. From this the  $\Phi(x)$  matrix can be generated after selecting a suitable  $f$ . A number of these were tried, and the best one for the problem seemed to be the polynomial one, that is  $f_i(x) = x^{i-1}$ . As a parameter,  $d = 10$  was used.

Now we have to find the optimal  $\hat{\theta}$  parameter vector, for which  $\Phi \theta = \mathbf{y} = [y_1 \dots y_n]^T$ . This is done like so:  $\theta^* \approx \hat{\theta} = \Phi^+ \mathbf{y}$ . The  $\Phi$  matrix of the sampled inputs and of the LS estimate are the same, so the function is evaluated at the same  $x$  values.

A computationally cheaper method for the pseudoinverse is QR decomposition. For this, the "economic" mode of scipy's `qr` function is used. Then the pseudoinverse is  $\Phi^+ = \mathbf{R}^{-1} \mathbf{Q}^T$ . Because the pseudoinverse of a matrix is unique, this method gives the same result as the previous one.

The results are shown in figure [1a](#).

### 1.1.2 Recursive least-squares

Next, more of the above described  $(x, y)$  pairs is sampled and  $\hat{\theta}_n$  calculated periodically using recursive least-squares. In the code, all of the samples are measured beforehand for simplicity.

The equation for  $\hat{\theta}$  can be reformulated in the following way:

$$\hat{\theta} = \underbrace{\left[ \sum_{i=1}^n \varphi \varphi^T \right]^{-1}}_{\Psi_n} \underbrace{\sum_{i=1}^n y_i \varphi_i}_{z_n}. \quad (1)$$

Now we have an update rule for both  $\Psi_{n+1} = (\Psi + \varphi_{n+1} \varphi_{n+1}^T)^{-1}$ , and  $z_{n+1} = z_n + \varphi_{n+1} y_{n+1}$ . Both  $\Psi_0$  and  $z_0$  are set to zero. Also, in the code  $\hat{\theta}_n$  is calculated only when plotted for speed.

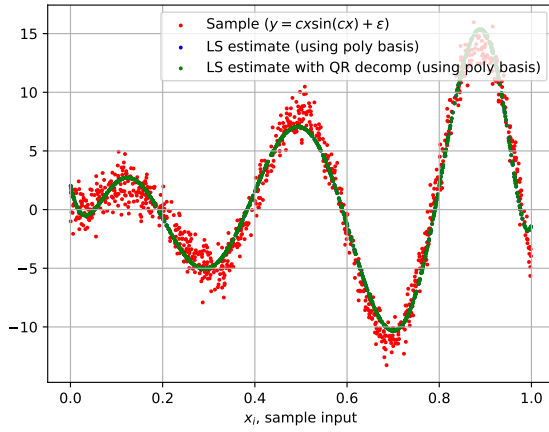
The resulting plots are shown in figure [1c](#)., taken at  $n \in [25, 50, 75, 100]$ .

### 1.1.3 Least-norm problem

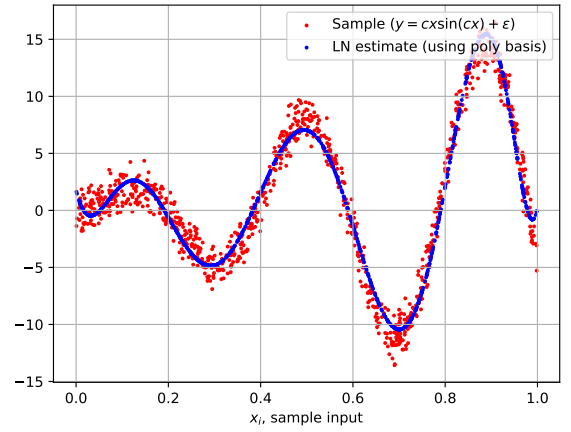
Now let's make  $d > n$ , specifically  $n = 100$  and  $d = 200$ , which makes  $\Phi$  fat. We make the assumption that  $\Phi$  is still full-rank. The solution is the same, except we are going to use singular value decomposition (SVD) for the pseudoinverse of  $\Phi$ .

The SVD is calculated like this:  $\Phi_{d \times n} = \mathbf{U}_{d \times d} \mathbf{\Sigma}_{d \times n} \mathbf{V}_{n \times n}^T$ . Let's denote the matrix of column vectors of the normalized eigen-vectors of matrix  $\Phi$  by  $\text{eig}(\Phi)$ . Let's denote the eigenvalues by  $\text{eigval}(\Phi)$ . Then,  $\mathbf{U} = \text{eig}(\Phi \Phi^T)$ ,  $\mathbf{V} = \text{eig}(\Phi^T \Phi)$  and  $\mathbf{\Sigma} = \text{diag}(\text{eigval}(\Phi \Phi^T))_{d \times n}$ . Then, the pseudoinverse is  $\Phi^+ = \mathbf{V} \mathbf{\Sigma}^+ \mathbf{U}^T$ . For  $\mathbf{\Sigma}^+$ , we take the inverse of the non-zero elements of  $\mathbf{\Sigma}$ , and add zeros such that it has the shape of  $d \times n$ .

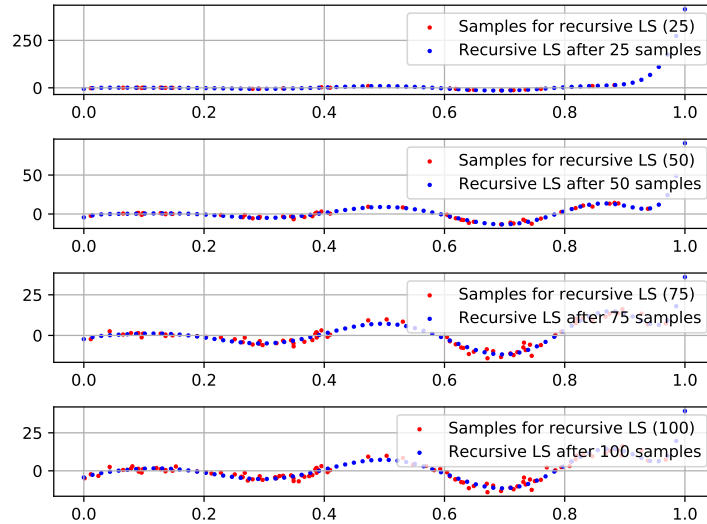
The resulting plots are shown in figure 1b.



(a) Least-squares estimate for thin  $\Phi$



(b) Least-norm estimate for fat  $\Phi$



(c) Recursive least-squares

Figure 1: Experiments with ordinary least-squares

## 1.2 Approximating auto-regressive series

A recursive time-series is generated from the give equation:  $y_t = a y_{t-1} + b y_{t-2} + \epsilon_t$ . Let's calculate the least-squares estimate using  $\varphi_t = [y_{t-1}, y_{t-2}]$ ,  $\Phi = [\varphi_1 \dots \varphi_n]$ . Then  $\hat{\theta} = \Phi^+ y$ .

The time plots can be seen on figure 2a.

Let's calculate the inverse of the covariance matrix:s  $\Gamma_n = \frac{1}{n} \Phi^T \Phi$ . Now define  $\Delta\theta := (\theta - \hat{\theta}_n)$ . The confidence ellipsoid is given by

$$\Delta\theta^T \Gamma_n \Delta\theta \leq \frac{q \hat{\sigma}_n^2}{n}, \quad (2)$$

where  $q$  is calculated from the inverse of the cumulative  $\chi^2$  distribution function given the  $p$  probabilities ( $q = F(p)_{\chi^2(d)}^{-1}$ ). In this problem,  $d = 2$ . Eq. 2 means that with probability  $p$ , the optimal  $\theta^*$  is at most  $\Delta\theta$  distance from  $\hat{\theta}$ .

Now we assume that  $\hat{\sigma}_n = 1$ , and let  $\Gamma_{ij}/n = [\Gamma_n]_{ij}$ . Then we have an equation that outputs an ellipse for a  $p$  probability. Written out:

$$\Delta\theta_1^2 \Gamma_{11} + 2 \Delta\theta_1 \Delta\theta_2 \Gamma_{12} + \Delta\theta_2^2 \Gamma_{22} = q. \quad (3)$$

For the plotting, this equation is transformed so that the axis distances and the rotation angle is known with the function `ellipse_transform[1]`.

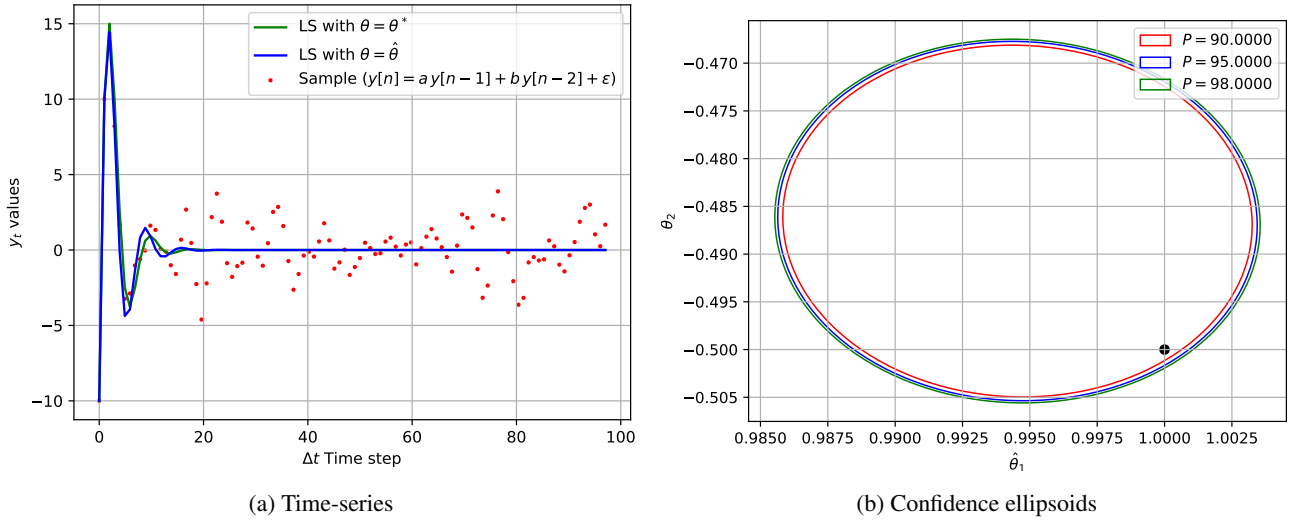


Figure 2: Experiments with auto-regressive series

## 2 Kernel methods

## 3 Reinforcement learning

### 3.1 The environment

For the environment, the cliff walking with a size of  $12 \times 4$  was chosen. Specifically, a modified version of caburu's `gym-cliffwalking`[2] is used. Two main modifications have been made to the original code:

- Originally, the state-space had a size of 48, though 10 of these are not real states. The cliff states were excluded.
- Rendering has been added.

So the environment has a state-space of size 38, and an action-space of size 4 (right, down, left, up).

### 3.2 Model based methods

#### 3.2.1 Model generating

Firstly, the model has been generated, which is basically a database of transitions and rewards.  $m : \mathbb{X} \times \mathbb{A} \rightarrow \mathbb{R} \times \mathbb{X}$ , where  $\mathbb{X}$  is the state-space and  $\mathbb{A}$  is the action-space. A random policy was used here. A figure of this can be seen running `3-Reinforcement_learning/show_scene.py` (taking the given action in the given state, on the left, the numbers mean the rewards, on the right, the next states).

#### 3.2.2 Linear programming

The optimal solution is given by linear programming. Let  $V$  be an arbitrary value-function. The optimal  $V^*$  is obtained by minimizing  $\sum V_S$ , given

$$V_S \geq g(S, A) + \beta V_{S+1} \text{ and } V_{\text{goal}} = 0, \quad (4)$$

where  $\beta$  is the discount factor, and  $g(S, A)$  is the immediate reward. The optimization is done with the python library `cvxpy`.

The resulting optimal value-function can be seen on figure 3.

#### 3.2.3 Value iteration

An iterative solution is using value iteration. Start with an arbitrary value-function (all zeros in this case), and repeatedly sweep through the state-space. For all the states

$$V_S^{n+1} = \max_{A \in \mathbb{A}} (R + \beta V_{S+1}^n), \quad (5)$$

where  $(S, A) \Rightarrow R$  and  $(S, A) \Rightarrow S + 1$  can be queried from the model. Then,  $\lim_{n \rightarrow \infty} V^n \rightarrow V^*$ .

The value-functions at some iterations are shown on figure 4.

### 3.2.4 Policy iteration

An other iterative method is policy iteration. Here we start with an arbitrary policy, then we evaluate it by calculating its value-function. This is done similarly to 5, only  $A$  is not the one with the maximal reward, but the one given by the current policy.

After this, we make the policy greedy with respect to the calculated value-function:  $p(S) = \operatorname{argmax}_{A \in \mathbb{A}} (R + \beta V_{S+1})$ .

The value-functions at some iterations are shown on figure 5.

The (euclidean) distances of both  $V_{VI}$  and  $V_{PI}$  from  $V^*$  are shown on figure 9a.

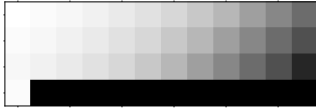


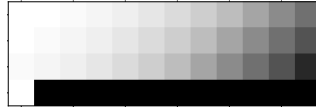
Figure 3: Linear programming



(a) Iteration 2

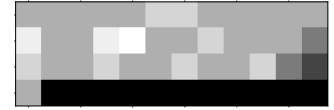


(b) Iteration 5

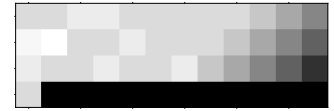


(c) Iteration 12

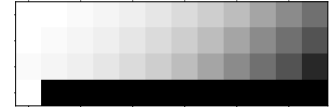
Figure 4: Value iteration



(a) Iteration 2



(b) Iteration 5



(c) Iteration 12

Figure 5: Policy iteration

## 3.3 Model-free methods

In this section, online Q-learning is going to be implemented. The update rule of Watkins' Q-learning is as follows:

$$Q_{n+1}(S, A) = (1 - \gamma_n) Q_n(S, A) + \gamma_n (R + \beta \max_{B \in \mathbb{A}} Q_n(\tilde{S}, B)), \quad (6)$$

where  $\gamma_n = \frac{1}{n+1}$  is the learning rate at step  $n$ , and  $\tilde{S}$  is the next state. At every step,  $A = p(S)$  is given by the policy.

Three policies are going to be put to test.

Firstly, the random policy, which just generates random actions for every state.



Second, the  $\epsilon$ -greedy policy. This acts greedily (see paragraph 3.2.4), most of the time, with acts randomly with an  $\epsilon$  probability so as to encourage exploration.

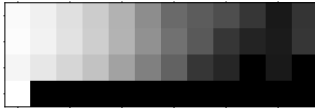
Lastly, the semi-greedy policy (called soft-max in the code) basically acts randomly if it doesn't have a much better choice. The exact probability of choosing action  $A$  is given by

$$\mathbb{P}(\pi_n(S) = A) = \frac{\exp(Q_n(S, A)/\tau)}{\sum_{B \in \mathbb{A}} \exp(Q_n(S, B)/\tau)}, \quad (7)$$

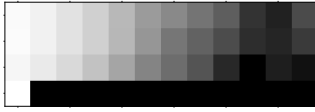
where  $\tau$  is the so-called Boltzmann-temperature, which influences the randomness of the policy.

For all three policies, the value-functions at some iterations are shown on figures 6, 7 and 8.

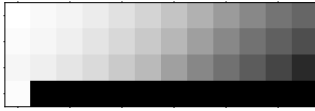
Also the distances of these policies' value-functions from the optimal one can be seen in figure 9b. The sums of the rewards received are also shown in figure 9c.



(a) Iteration 2

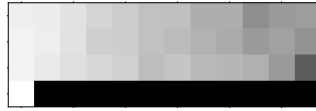


(b) Iteration 5

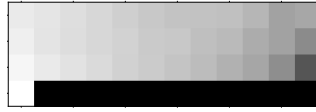


(c) Iteration 30

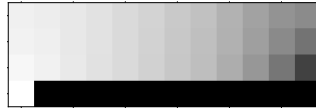
Figure 6: Random policy



(a) Iteration 2



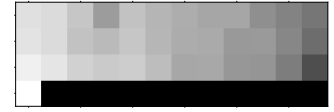
(b) Iteration 5



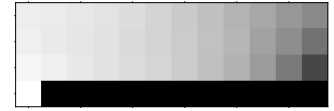
(c) Iteration 30

Figure 7:  $\epsilon$ -greedy policy

(a) Iteration 2

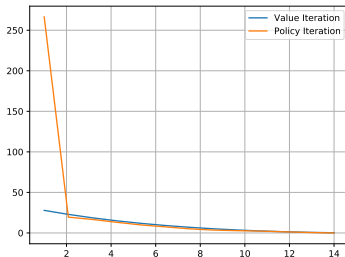


(b) Iteration 5

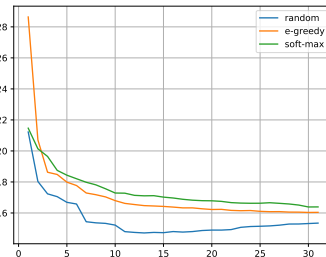


(c) Iteration 30

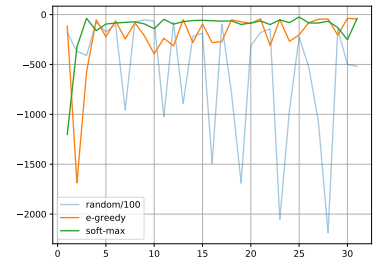
Figure 8: Semi-greedy policy



(a) Model-based



(b) Model-free



(c) Rewards (model-free)

Figure 9: Distance from the optimal value-function, and rewards of policies

## References

[1] Ellipse equation:

[https://www.maa.org/external\\_archive/joma/Volume8/Kalman/General.html](https://www.maa.org/external_archive/joma/Volume8/Kalman/General.html),  
2021. Jan 6., 12:22

[2] Cliffwalking environment:

<https://github.com/caburu/gym-cliffwalking>,  
2021. Jan 6., 12:26