

Reinforcement Learning - Exercises Lectures 1-5

Maurice Frank

11650656

maurice.frank@posteo.de

Code: [github](#)

September 17, 2019

Lecture 0:

0.1 Linear algebra and multivariable derivatives

1.

$$AB = \begin{bmatrix} a_{11} & 0 \\ 0 & a_{22} \end{bmatrix} \cdot \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix} \quad (1)$$

$$= \begin{bmatrix} a_{11}b_{11} & a_{11}b_{12} \\ a_{22}b_{21} & a_{22}b_{22} \end{bmatrix} \quad (2)$$

$$AB^T = \begin{bmatrix} a_{11} & 0 \\ 0 & a_{22} \end{bmatrix} \cdot \begin{bmatrix} b_{11} & b_{21} \\ b_{12} & b_{22} \end{bmatrix} \quad (3)$$

$$= \begin{bmatrix} a_{11}b_{11} & a_{11}b_{21} \\ a_{22}b_{12} & a_{22}b_{22} \end{bmatrix} \quad (4)$$

$$d^T B d = \begin{bmatrix} d_1 & d_2 \end{bmatrix} \begin{bmatrix} b_{11} & b_{12} \\ b_{21} & b_{22} \end{bmatrix} \begin{bmatrix} d_1 \\ d_2 \end{bmatrix} \quad (5)$$

$$= d_1^2 b_{11} + d_1 d_2 b_{12} + d_1 d_2 b_{21} + d_2^2 b_{22} \quad (6)$$

2.

$$A^{-1} = \begin{bmatrix} a_{11}^{-1} & 0 \\ 0 & a_{22}^{-1} \end{bmatrix} \quad (7)$$

$$B^{-1} = \frac{1}{b_{11}b_{22} - b_{21}b_{12}} \begin{bmatrix} b_{22} & -b_{21} \\ -b_{12} & b_{11} \end{bmatrix} \quad (8)$$

3.

$$\frac{\partial c}{\partial x} = \begin{bmatrix} -2x \\ \frac{1}{yx} \end{bmatrix} \quad (9)$$

$$\frac{\partial c}{\partial e} = \begin{bmatrix} -2x & 1 \\ \frac{1}{yx} & -\ln(x)y^{-2} \end{bmatrix} \quad (10)$$

4.

$$f(\mathbf{x}) = \sum_i^N ix_i \quad (11)$$

$$\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} = [1, \dots, N] \quad (12)$$

0.2 Probability theory

1.

$$\mathbb{E}[X + \alpha Y] = \mathbb{E}[X] + \alpha \mathbb{E}[Y] \quad (13)$$

$$= \mu + \alpha \nu \quad (14)$$

2.

$$\text{Var}[X + \alpha Y] = \text{Var}[X] + \alpha^2 \text{Var}[Y] + 2\alpha \text{cov}[X, Y] \quad (15)$$

3.

At last we have σ^2 which is just the variance of the noise in the measurements. This is model independent (we can not train it away). The bias term tells us how good our estimator estimates the sample data points. The estimator variance tells us how jumpy our estimator is. If the model has little parameters ('smooth') it will have high bias but low variance (if regularizes over the evident points but doesn't estimate the data that good anymore). If the model is a complex one with high capacity it will have low bias but high variance.

4.

This is called the bias-variance trade-off as in machine learning we are mostly interested in building a model which has high bias and high variance. The trade-off shows us that these two are opposite in their objective and that optimizing both is not easy.

0.3 OLS, linear projection and gradient descent

We have training set $\mathbf{X} \in \mathbb{R}^{n \times m}$ with targets $\mathbf{y} \in \mathbb{R}^n$. We have a linear model $f_{\beta}(\mathbf{X}) = \mathbf{X} \cdot \beta$.

1.

$$\beta \in \mathbb{R}^m$$

2.

$$\begin{aligned}\hat{\beta} &= \arg \min_{\beta} (\mathbf{y} - f_{\beta}(\mathbf{X}))^2 \\ \frac{\partial}{\partial \hat{\beta}} (\mathbf{y} - f_{\hat{\beta}}(\mathbf{X}))^2 &= \frac{\partial}{\partial \hat{\beta}} (\mathbf{y} - \mathbf{X}\hat{\beta})^2 \\ &= \frac{\partial}{\partial \hat{\beta}} \mathbf{y}^2 - 2\mathbf{y}\mathbf{X}\hat{\beta} + (\mathbf{X}\hat{\beta})^2 \\ &= 2\mathbf{X}^T\mathbf{X}\hat{\beta} - 2\mathbf{y}\mathbf{X} \\ &\stackrel{\text{def}}{=} 0 \\ &\iff \\ \mathbf{X}^T\mathbf{X}\hat{\beta} &= \mathbf{X}^T\mathbf{y} \\ &\iff \\ \hat{\beta} &= (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}\end{aligned}$$

3.

$$\epsilon_{\beta} = \mathbf{y} - \mathbf{X}\hat{\beta}$$

4.

5.

6.

7.

Lecture 1: Introduction

1.1 Introduction

1.

The curse of dimensionality are multiple sad observation called one will make when working with high-dimensional data. In general the problems arise from

the fact that the number of value combinations rises exponentially, with the dimension in the exponential. E.g. in hyper-parameter optimization using grid-search the number of needed models to be tested rises exponentially with the number of hyper-parameters. Another example we often see in machine learning with high-dimensional data. With a limited number of trainings samples the distribution of those might be highly sparse in its space akin like a set of dirac functions. Trying to approximate that might be difficult.

2.

(a)

$$\begin{aligned} N_{\text{states}} &= N_{\text{predator states}} \cdot N_{\text{prey states}} \\ &= 5^2 \cdot 5^2 \\ &= 625 \end{aligned}$$

(b) As it is a toroid we just have to remember the differences of the two entities. So the state is just the offset in toroidal coordinates.

(c)

$$\begin{aligned} N'_{\text{states}} &= 5 \cdot 5 \\ &= 25 \end{aligned}$$

(d) The advantage of this approach is that we have fewer states and now multiple states that have to learn the same response to it. Thus we can assume faster training of our predator.

(e) For Tic-Tac-Toe we could reduce the state space by using the point symmetry of the game board. Of all starting states that are symmetric through the center only keep one.

3.

(a) The greedy agent will perform better. Tic-tac-toe is a solved game as such a trained agent can know the perfect move to any situation and no exploration is necessary.

4.

(a) We decrease the exploration probability ϵ each step with a discount factor η . We can write the exploration probability at step ϵ_t with:

$$\epsilon_t = \epsilon \cdot \eta^t$$

- (b) No, that method would not work if the opponent changes strategy. It continuously decreases exploration over time independent of the game dynamics. We can adapt our strategy by introducing the time step of the last strategy change of the opponent t_{change} . Then we restart from the beginning if the strategy changes:

$$\epsilon_t = \epsilon \cdot \eta^{t-t_{\text{change}}}$$

1.2 Exploration

1.

$$(1 - \epsilon) + \frac{\epsilon}{n}$$

2.

A_3 and A_4 . The first one could be greedy as all states have the same average. Same for the second as 2 and 3 have the greedy average. The third action could be greedy as state 2 has the top average then. The next two are suboptimal thus have to be exploration.

3.

$R_0 = -1$ and $R_1 = +1$. By random we choose A_0 in the first step. See the development of the Q-values (bold is the chooses action with greedy policy):

step	Q_0^{pessi}	Q_1^{pessi}	Q_0^{opti}	Q_1^{opti}
0	-5	-5	5	5
1	-1	-5	-1	5
2	-1	-5	-1	1
3	-1	-5	-1	1

4.

The optimistic initialization leads to the higher return (== 1) than with the pessimistic initialization (== -3). If broken the tie the other way:

step	Q_0^{pessi}	Q_1^{pessi}	Q_0^{opti}	Q_1^{opti}
0	-5	-5	5	5
1	-5	1	5	1
2	-5	1	-1	1
3	-5	1	-1	1

In this case the pessimistic initialization lead to the higher return (== 3) than the optimistic initialization (== 1).

5.

The optimistic initialization leads to the better estimate of the Q-values.

6.

The optimistic initialization works better for exploration as its basic assumption is that any action could be the best until proven otherwise. As such it will lead to everything being tried using the high initialization values.

Lecture 2: MDPs and dynamic programming

2.1 Markov Decision Processes

1.

(a)

2.

(a)

2.2 Homework: Dynamic Programming

1.

2.

3.

4.