



Отчет о проверке

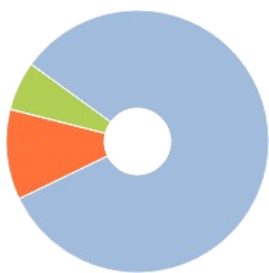
Автор: dinhngoctuan6789@gmail.com / ID: 11491219

Проверяющий:

Название документа: ВКР Динь Нгок Туан R34372

РЕЗУЛЬТАТЫ ПРОВЕРКИ

Тариф: FULL



Совпадения:
11,43%



Оригинальность:
82,1%



Цитирования:
6,47%



Самоцитирования:
0%



1

«Совпадения», «Цитирования», «Самоцитирования», «Оригинальность» являются отдельными показателями, отображаются в процентах и в сумме дают 100%, что соответствует проверенному тексту документа.

- **Совпадения** — фрагменты проверяемого текста, полностью или частично сходные с найденными источниками, за исключением фрагментов, которые система отнесла к цитированию или самоцитированию. Показатель «Совпадения» — это доля фрагментов проверяемого текста, отнесенных к совпадениям, в общем объеме текста.
- **Самоцитирования** — фрагменты проверяемого текста, совпадающие или почти совпадающие с фрагментом текста источника, автором или соавтором которого является автор проверяемого документа. Показатель «Самоцитирования» — это доля фрагментов текста, отнесенных к самоцитированию, в общем объеме текста.
- **Цитирования** — фрагменты проверяемого текста, которые не являются авторскими, но которые система отнесла к корректно оформленным. К цитированиям относятся также шаблонные фразы; библиография; фрагменты текста, найденные модулем поиска «СПС Гарант: нормативно-правовая документация». Показатель «Цитирования» — это доля фрагментов проверяемого текста, отнесенных к цитированию, в общем объеме текста.
- **Текстовое пересечение** — фрагмент текста проверяемого документа, совпадающий или почти совпадающий с фрагментом текста источника.
- **Источник** — документ, проиндексированный в системе и содержащийся в модуле поиска, по которому проводится проверка.
- **Оригинальный текст** — фрагменты проверяемого текста, не обнаруженные ни в одном источнике и не отмеченные ни одним из модулей поиска. Показатель «Оригинальность» — это доля фрагментов проверяемого текста, отнесенных к оригинальному тексту, в общем объеме текста.

Обращаем Ваше внимание, что система находит текстовые совпадения проверяемого документа с проиндексированными в системе источниками. При этом система является вспомогательным инструментом, определение корректности и правомерности совпадений или цитирований, а также авторства текстовых фрагментов проверяемого документа остается в компетенции проверяющего.

ИНФОРМАЦИЯ О ДОКУМЕНТЕ

Номер документа: 2

Тип документа: Не указано

Дата проверки: 12.05.2024 12:52:49

Дата корректировки: Нет

Количество страниц: 41

Символов в тексте: 53516

Слов в тексте: 6988

Число предложений: 3134

Комментарий: не указано

ПАРАМЕТРЫ ПРОВЕРКИ

Выполнена проверка с учетом редактирования: Да

Выполнено распознавание текста (OCR): Нет

Выполнена проверка с учетом структуры: Нет

Модули поиска: Переводные заимствования (KyEn), Цитирование, Коллекция НБУ, СПС ГАРАНТ: нормативно-правовая документация, СПС ГАРАНТ: аналитика, Переводные заимствования (KkEn), Патенты СССР, РФ, СНГ, Переводные заимствования*, Переводные заимствования по Интернету (KyRu), Переводные заимствования (RuEn), Переводные заимствования по коллекции Интернет в английском сегменте, ИПС Адилет, Шаблонные фразы, Публикации РГБ, Диссертации НББ, Библиография, Публикации eLIBRARY, Перефразирования по СПС ГАРАНТ: аналитика, Перефразирования по Интернету (EN), Переводные заимствования издательства Wiley, Перефразированные заимствования по коллекции Интернет в английском сегменте, Кольцо вузов, Медицина, Перефразирования по коллекции издательства Wiley, Переводные заимствования IEEE, Переводные заимствования по Интернету (EnRu), Переводные заимствования по Интернету (KkRu), Переводные заимствования по коллекции Гарант: аналитика, Перефразирования по коллекции IEEE, Публикации eLIBRARY (переводы и перефразирования), IEEE, Перефразированные заимствования по коллекции Интернет в русском сегменте, Издательство Wiley, Сводная коллекция ЭБС, СМИ России и СНГ, Перефразирования по Интернету, Публикации РГБ (переводы и перефразирования), Переводные заимствования по коллекции Интернет в русском сегменте, Интернет Плюс*

❌ Модули, недоступные в рамках тарифа: Интернет Free

ИСТОЧНИКИ

№	Доля в тексте	Доля в отчете	Источник	Актуален на	Модуль поиска	Комментарий
[01]	5,58%	5,58%	не указано	13 Янв 2022	Библиография	
[02]	3,61%	3,56%	Reinforcement Learning and Feed... https://ieeexplore.ieee.org	22 Июн 2023	Переводные заимствования IEEE	
[03]	2,43%	0,46%	Optimal and Autonomous Control ... https://ieeexplore.ieee.org	07 Дек 2017	IEEE	
[04]	2%	0,98%	Р. С. Саттон, Э. Г. Барто ; пер. с а... http://dlib.rsl.ru	01 Дек 2014	Публикации РГБ (переводы и перефразирования)	
[05]	1,9%	1,9%	Предвзятость алгоритмов искусс... http://ivo.garant.ru	30 Окт 2021	Перефразирования по СПС ГАРАНТ: аналитика	
[06]	1,68%	0%	Books and Papers http://uta.edu	31 Янв 2017	Интернет Плюс*	
[07]	1,65%	0%	Reinforcement Learning and Feed... https://ieeexplore.ieee.org	22 Июн 2023	IEEE	
[08]	1,65%	0%	https://iris.poliba.it/retrieve/dd89f... https://iris.poliba.it	12 Мая 2024	Интернет Плюс*	
[09]	1,63%	0%	http://www.ijmlc.org/vol6/594-L00... http://ijmlc.org	12 Мая 2024	Интернет Плюс*	
[10]	1,62%	0,71%	courses:ml:d0%be%d0%b1%d1... https://se.moevm.info	29 Мар 2023	Перефразированные заимствования по коллекции Интернет в русском сегменте	
[11]	1,58%	1,07%	120479 http://biblioclub.ru	15 Апр 2016	Сводная коллекция ЭБС	
[12]	1,58%	0%	Обучение с подкреплением http://biblioclub.ru	21 Янв 2020	Сводная коллекция ЭБС	
[13]	1,55%	0%	Neural Network Based Online Sim... https://ieeexplore.ieee.org	12 Окт 2012	IEEE	
[14]	1,48%	0%	Actor-critic neural network reinfor... https://ieeexplore.ieee.org	18 Дек 2014	IEEE	
[15]	1,43%	0%	Off-Policy Reinforcement Learning... https://ieeexplore.ieee.org	09 Мая 2014	IEEE	
[16]	1,38%	1,38%	Нейронные сети: полный курс - ... https://djvu.online	20 Апр 2024	Перефразированные заимствования по коллекции Интернет в русском сегменте	
[17]	1,37%	0,14%	Reinforcement learning of LQR co... https://ieeexplore.ieee.org	09 Июн 2023	IEEE	
[18]	1,32%	0%	Off-policy reinforcement learning ... https://arxiv.org	10 Июл 2017	Интернет Плюс*	
[19]	1,32%	0%	Neural Network Based Online Sim... https://doi.org	06 Сен 2019	Интернет Плюс*	
[20]	1,31%	0%	https://arxiv.org/pdf/1311.6107.pdf https://arxiv.org	12 Мая 2024	Интернет Плюс*	

[21]	1,3%	0%	基于数据的自学习优化控制:研究进展... http://aas.net.cn	12 Мая 2024	Интернет Плюс*	
[22]	1,29%	0%	https://www.princeton.edu/~nda... https://princeton.edu	06 Окт 2022	Интернет Плюс*	
[23]	1,29%	0%	66325 http://e.lanbook.com	09 Мар 2016	Сводная коллекция ЭБС	
[24]	1,29%	0%	Обучение с подкреплением https://e.lanbook.com	22 Янв 2020	Сводная коллекция ЭБС	
[25]	1,28%	0%	Optimal Adaptive Control of Const... https://ieeexplore.ieee.org	05 Апр 2024	Перефразирования по коллекции IEEE	
[26]	1,28%	0%	[IEEE 2012 American Control Conf... https://doi.org	18 Дек 2018	Интернет Плюс*	
[27]	1,27%	0%	Model-based vs data-driven adap... https://doi.org	31 Мая 2018	Издательство Wiley	
[28]	1,14%	0%	Model-Free H ∞ Control Design for...	31 Мар 2016	Издательство Wiley	
[29]	1,13%	0%	Adaptive Optimal Control of Highl... https://doi.org	19 Окт 2019	Интернет Плюс*	
[30]	1,13%	0%	Online adaptive algorithm for opti... https://doi.org	18 Апр 2021	Перефразирования по коллекции издательства Wiley	
[31]	1%	0%	Model-based reinforcement learn... https://doi.org	01 Окт 2020	Перефразирования по коллекции издательства Wiley	
[32]	0,97%	0%	Data-Driven Adaptive Critic Appro... https://doi.org	31 Янв 2018	Издательство Wiley	Источник исключен. Причина: Маленький процент пересечения.
[33]	0,94%	0%	Reinforcement learning control fo... https://core.ac.uk	30 Окт 2020	Интернет Плюс*	Источник исключен. Причина: Маленький процент пересечения.
[34]	0,9%	0%	https://ieeexplore.ieee.org/ielaam... https://ieeexplore.ieee.org	12 Мая 2024	Интернет Плюс*	Источник исключен. Причина: Маленький процент пересечения.
[35]	0,89%	0,89%	не указано	13 Янв 2022	Шаблонные фразы	
[36]	0,86%	0%	Adaptive optimal tracking control ... https://ieeexplore.ieee.org	29 Окт 2015	IEEE	Источник исключен. Причина: Маленький процент пересечения.
[37]	0,86%	0%	Adaptive neural network-based o... https://doi.org	04 Сен 2013	Издательство Wiley	Источник исключен. Причина: Маленький процент пересечения.
[38]	0,77%	0,77%	CITech_abstracts.pdf http://conf.nsc.ru	04 Мая 2023	Перефразированные заимствования по коллекции Интернет в английском сегменте	
[39]	0,72%	0%	Adaptive Optimal Control of Unkn... https://ieeexplore.ieee.org	03 Мая 2024	Перефразирования по коллекции IEEE	
[40]	0,67%	0%	Online solution of nonlinear two-... https://doi.org	21 Июл 2011	Издательство Wiley	Источник исключен. Причина: Маленький процент пересечения.
[41]	0,62%	0%	Off-policy reinforcement learning ... http://arxiv.org	08 Янв 2018	Перефразирования по Интернету (EN)	
[42]	0,6%	0%	melcer_s_p_issledovanie-sredy-s-...	09 Янв 2024	Кольцо вузов	
[43]	0,59%	0%	Autonomy and Machine Intelligen... http://ece.ucsb.edu	09 Янв 2018	Перефразирования по Интернету (EN)	
[44]	0,56%	0%	http://www.mit.edu/~dimitrib/RL_... http://mit.edu	28 Мая 2022	Интернет Плюс*	Источник исключен. Причина: Маленький процент пересечения.
[45]	0,52%	0%	zisman_i_a_proekt-kooperaciya-eg...	09 Янв 2024	Кольцо вузов	
[46]	0,47%	0%	Review on Reinforcement Learnin... https://ieeexplore.ieee.org	30 Дек 2022	IEEE	Источник исключен. Причина: Маленький процент пересечения.
[47]	0,46%	0,46%	Выборнов, Андрей Олегович дис... http://dlib.rsl.ru	20 Янв 2010	Публикации РГБ (переводы и перефразирования)	
[48]	0,45%	0%	https://hal.science/hal-00756747/f... https://hal.science	24 Фев 2023	Перефразированные заимствования по коллекции Интернет в английском сегменте	
[49]	0,45%	0%	https://hal.archives-ouvertes.fr/ha... https://hal.archives-ouvertes.fr	28 Ноя 2022	Перефразированные заимствования по коллекции Интернет в английском сегменте	
[50]	0,45%	0%	A Survey of Actor-Critic Reinforce... https://ieeexplore.ieee.org	22 Июн 2023	Перефразирования по коллекции IEEE	
[51]	0,44%	0%	Multi Objective Combinatorial Opt...	09 Янв 2024	Кольцо вузов	Источник исключен. Причина: Маленький процент пересечения.
[52]	0,43%	0%	ISBN9785996325009.txt	26 Окт 2017	Кольцо вузов	Источник исключен. Причина: Маленький процент пересечения.
[53]	0,41%	0%	https://spb.hse.ru/mirror/pubs/sh... https://spb.hse.ru	12 Мая 2024	Интернет Плюс*	Источник исключен. Причина: Маленький процент пересечения.

[54]	0,41%	0%	https://spb.hse.ru/mirror/pubs/sh...https://spb.hse.ru	12 Мая 2024	Интернет Плюс*	Источник исключен. Причина: Маленький процент пересечения.
[55]	0,4%	0%	Q-Learning Algorithms: A Compre... https://ieeexplore.ieee.org	13 Сен 2019	IEEE	Источник исключен. Причина: Маленький процент пересечения.
[56]	0,4%	0%	Adaptive Optimal Control of Unkn... https://ieeexplore.ieee.org	03 Мая 2024	Переводные заимствования IEEE	Источник исключен. Причина: Маленький процент пересечения.
[57]	0,4%	0%	Novel Discounted Adaptive Critic ... https://ieeexplore.ieee.org	23 Апр 2024	Переводные заимствования IEEE	Источник исключен. Причина: Маленький процент пересечения.
[58]	0,39%	0%	Robust Adaptive Dynamic Progra... http://arxiv.org	31 Янв 2017	Перефразирования по Интернету (EN)	Источник исключен. Причина: Маленький процент пересечения.
[59]	0,38%	0%	Model-based and Model-free Opti... https://ieeexplore.ieee.org	12 Июл 2023	Переводные заимствования IEEE	Источник исключен. Причина: Маленький процент пересечения.
[60]	0,37%	0%	Предотвращение столкновений ... http://elibrary.ru	10 Янв 2019	Публикации eLIBRARY (переводы и перефразирования)	Источник исключен. Причина: Маленький процент пересечения.
[61]	0,36%	0%	Ascertaining properties of weighti...	10 Ноя 2020	Издательство Wiley	Источник исключен. Причина: Маленький процент пересечения.
[62]	0,36%	0%	Data-Driven Actuator Allocation fo... https://ieeexplore.ieee.org	05 Апр 2024	Перефразирования по коллекции IEEE	Источник исключен. Причина: Маленький процент пересечения.
[63]	0,36%	0%	Adaptive Optimal Stabilization Re... https://ieeexplore.ieee.org	28 Мар 2024	Перефразирования по коллекции IEEE	Источник исключен. Причина: Маленький процент пересечения.
[64]	0,35%	0%	Robust Policy Iteration of Uncertai... https://ieeexplore.ieee.org	19 Апр 2024	Перефразирования по коллекции IEEE	Источник исключен. Причина: Маленький процент пересечения.
[65]	0,34%	0%	Design and Comparison Base Ana... https://doi.org	25 Мая 2021	Перефразирования по коллекции издательства Wiley	Источник исключен. Причина: Маленький процент пересечения.
[66]	0,33%	0%	Data-driven Control and Optimiza... http://uta.edu	06 Янв 2018	Перефразирования по Интернету (EN)	Источник исключен. Причина: Маленький процент пересечения.
[67]	0,33%	0%	Integral Reinforcement Learning f... http://uta.edu	07 Янв 2018	Перефразирования по Интернету (EN)	Источник исключен. Причина: Маленький процент пересечения.
[68]	0,32%	0%	ОБУЧЕНИЕ С ПОДКРЕПЛЕНИЕМ.	19 Окт 2023	Публикации eLIBRARY (переводы и перефразирования)	Источник исключен. Причина: Маленький процент пересечения.
[69]	0,32%	0%	Стратегии исследования окруже... http://machinelearning.ru	03 Июн 2020	Переводные заимствования по коллекции Интернет в русском сегменте	Источник исключен. Причина: Маленький процент пересечения.
[70]	0,31%	0%	https://books.ifmo.ru/file/pdf/258...https://books.ifmo.ru	07 Июн 2022	Интернет Плюс*	Источник исключен. Причина: Маленький процент пересечения.
[71]	0,31%	0%	Data-Driven Adaptive Critic Appro... https://doi.org	31 Янв 2018	Перефразирования по коллекции издательства Wiley	Источник исключен. Причина: Маленький процент пересечения.
[72]	0,31%	0%	Adaptive optimal tracking control ... https://sciencedirect.com	12 Мая 2024	Интернет Плюс*	Источник исключен. Причина: Маленький процент пересечения.
[73]	0,29%	0%	https://arxiv.org/pdf/1908.02077.p...https://arxiv.org	12 Мая 2024	Интернет Плюс*	Источник исключен. Причина: Маленький процент пересечения.
[74]	0,29%	0%	Using Social Cues to Recognize Ta... https://arxiv.org	25 Мар 2023	Перефразированные заимствования по коллекции Интернет в английском сегменте	Источник исключен. Причина: Маленький процент пересечения.
[75]	0,29%	0%	«Стратегии исследования окруж... http://machinelearning.ru	09 Июн 2020	Переводные заимствования по коллекции Интернет в русском сегменте	Источник исключен. Причина: Маленький процент пересечения.
[76]	0,28%	0%	6342-9949-1-SM.pdf	15 Янв 2024	Кольцо вузов	Источник исключен. Причина: Маленький процент пересечения.
[77]	0,28%	0%	https://ctlab.itmo.ru/~vaksenov/di...https://ctlab.itmo.ru	20 Апр 2024	Интернет Плюс*	Источник исключен. Причина: Маленький процент пересечения.
[78]	0,27%	0%	https://www.cs.colostate.edu/~an...https://cs.colostate.edu	20 Мая 2022	Интернет Плюс*	Источник исключен. Причина: Маленький процент пересечения.
[79]	0,27%	0%	Machine Learning-Based Researc... https://link.springer.com	11 Дек 2022	Интернет Плюс*	Источник исключен. Причина: Маленький процент пересечения.
[80]	0,26%	0%	Обучение с подкреплением http://studentlibrary.ru	20 Дек 2016	Медицина	Источник исключен. Причина: Маленький процент пересечения.
[81]	0,25%	0%	chepenko_d_d_offlayn-algoritm-ob...	24 Мая 2022	Кольцо вузов	Источник исключен. Причина: Маленький процент пересечения.
[82]	0,24%	0%	tarelov_v_o_proekt-razrabotka-inte...	01 Мар 2024	Кольцо вузов	Источник исключен. Причина: Маленький процент пересечения.
[83]	0,24%	0%	Распознавание нечётко определ... http://studentlibrary.ru	20 Дек 2016	Медицина	Источник исключен. Причина: Маленький процент пересечения.
[84]	0,23%	0%	Root Locus Analysis for Swinging ... https://ieeexplore.ieee.org	07 Авг 2023	IEEE	Источник исключен. Причина: Маленький процент пересечения.
[85]	0,21%	0%	Solutions for the Linear-Quadratic...	раньше 2011	Публикации eLIBRARY	Источник исключен. Причина: Маленький процент пересечения.
[86]	0,21%	0%	Data-Driven Optimal Assistance C... https://frontiersin.org	06 Июл 2020	СМИ России и СНГ	Источник исключен. Причина: Маленький процент пересечения.

[87]	0,21%	0%	Survey of Image Processing Techn... https://frontiersin.org	22 Мая 2020	СМИ России и СНГ	Источник исключен. Причина: Маленький процент пересечения.
[88]	0,19%	0%	Шизофрения и вероятностное п... http://emil.ru	21 Дек 2016	Медицина	Источник исключен. Причина: Маленький процент пересечения.
[89]	0,18%	0%	MDP (марковский процесс приня... https://mlcentre.ru	12 Мая 2024	Интернет Плюс*	Источник исключен. Причина: Маленький процент пересечения.
[90]	0,17%	0%	Качков_магистерская	11 Июн 2018	Кольцо вузов	Источник исключен. Причина: Маленький процент пересечения.
[91]	0,17%	0%	Leveraging explainable machine l... https://doi.org	25 Сен 2023	Издательство Wiley	Источник исключен. Причина: Маленький процент пересечения.
[92]	0,17%	0%	https://nauchkor.ru/uploads/docu... https://nauchkor.ru	12 Мая 2024	Интернет Плюс*	Источник исключен. Причина: Маленький процент пересечения.
[93]	0,17%	0%	Интеллектуальное управление и и... https://nauchkor.ru	10 Мая 2024	Интернет Плюс*	Источник исключен. Причина: Маленький процент пересечения.
[94]	0,16%	0%	Обучение с подкреплением: вве... https://russianblogs.com	19 Мая 2021	Интернет Плюс*	Источник исключен. Причина: Маленький процент пересечения.
[95]	0,16%	0%	Frontiers From classical to quan... https://translated.turbopages.org	12 Мая 2024	Интернет Плюс*	Источник исключен. Причина: Маленький процент пересечения.
[96]	0,16%	0%	smirnov_p_d_obuchenie-s-podkre...	09 Янв 2024	Кольцо вузов	Источник исключен. Причина: Маленький процент пересечения.
[97]	0,15%	0%	Применение глубокого обучения... http://elibrary.ru	01 Янв 2023	Публикации eLIBRARY	Источник исключен. Причина: Маленький процент пересечения.
[98]	0,15%	0%	Многомерные системы с неопре... http://elibrary.ru	28 Авг 2014	Публикации eLIBRARY	Источник исключен. Причина: Маленький процент пересечения.
[99]	0,14%	0%	Диссертация на тему «Метод и а... https://dissercat.com	12 Мая 2024	Интернет Плюс*	Источник исключен. Причина: Маленький процент пересечения.
[100]	0,13%	0%	Что такое обучение усилению? https://machinelearningmastery.ru	09 Мая 2024	Интернет Плюс*	Источник исключен. Причина: Маленький процент пересечения.
[101]	0,11%	0%	https://e-learning.bmstu.ru/iu6/pl... https://e-learning.bmstu.ru	11 Мая 2024	Интернет Плюс*	Источник исключен. Причина: Маленький процент пересечения.
[102]	0,11%	0%	Convergence of the approximate ...	01 Янв 2015	Публикации eLIBRARY	Источник исключен. Причина: Маленький процент пересечения.

ИССЛЕДОВАНИЕ АЛГОРИТМОВ ОПТИМАЛЬНОГО УПРАВЛЕНИЯ, ОСНОВАННЫХ НА ОБУЧЕНИИ С ПОДКРЕПЛЕНИЕМ

Титульный лист

Техническое задание (утвержденное из ИСУ)

Основные вопросы, подлежащие разработке / Key issues to be analyzed

Провести исследование метода обучения с подкреплением и способов его применения в задачах оптимизации, а также основных алгоритмов, основанных на этом методе.

Провести исследование оптимальных регуляторов на основе обучения с подкреплением.

Реализовать оптимальный адаптивный регулятор для динамических систем на основе алгоритмов обучения с подкреплением.

Провести моделирование оптимального алгоритма управления на основе обучения с подкреплением с объектами управления в MATLAB, выполнить анализ полученных результатов.

Аннотация

Цель исследования / Research goal

Исследование алгоритмов оптимального управления, основанных на обучении с подкреплением

Задачи, решаемые в ВКР / Research tasks

Для достижения цели требуется выполнение следующих подзадач: 1. Провести исследование метода обучения с подкреплением и способов его применения в задачах оптимизации, а также основных алгоритмов, основанных на этом методе. 2. Провести исследование оптимальных регуляторов на основе обучения с подкреплением 3. Реализовать оптимальный адаптивный регулятор для динамических систем на основе алгоритмов обучения с интегральным подкреплением 4. Провести моделирование оптимального алгоритма управления на основе обучения с подкреплением с объектами управления в симуляторе Matlab 5. Провести анализу работоспособности регулятора, основанного на обучении с подкреплением, на основе полученных результатов моделирования и сравните их с обычным оптимальным регулятором.

Краткая характеристика полученных результатов / Short summary of results/findings

Изучены типичные идеи обучения с подкреплением. Рассмотрены и исследованы математические модели построения оптимальных адаптивных регуляторов. Моделирование алгоритма выполнено в Matlab. В данной работе приближенное динамическое программирование в сочетании с алгоритмами обучения с подкреплением успешно нашло приближенное решение уравнения Гамильтона Якоби-Беллмана в реальном времени, тем самым было успешно реализовано построение адаптивного оптимального регулятора на основе алгоритмов обучения с интегральным подкреплением для непрерывных систем. Результаты сравнения и анализа показывают, что оптимальный адаптивный регулятор, основанный на обучении с подкреплением ⁵ имеет примерно такую же производительность, как и оптимальный регулятор.

Оглавление

Введение.....	4
1 Обучение с подкреплением.....	5
1.1 Обзор обучения с подкреплением (RL).....	5
1.2 Элемент обучения с подкреплением.....	6
1.3 Марковский процесс принятия решений (МППР) - Markov Decision Processe.....	7
1.3.1 Задачи оптимального последовательного решения.....	8
1.3.2 Обратная рекурсия для значения.....	9
1.3.3 Обзор адаптивного динамического программирования.....	9
1.3.4 Уравнение Беллмана и уравнение оптимальности Беллмана.....	10
1.4 Итерация по стратегии и итерация по критерию.....	12
1.4.1 Итерация по стратегии алгоритм.....	12
1.4.2 Итерации по критерию.....	13
1.4.3 Сравнение итерации по стратегии и итерации по критерию.....	13
2. Регулятора для линейных систем.....	13
2.1 Оптимальный регулятор.....	14
2.2 Адаптивный регулятор.....	14
2.3 Оптимальный адаптивный регулятор.....	15
3. Задача оптимального управления для нелинейных систем.....	17
3.1 Оптимальное управление и уравнение Гамильтона Якоби – Беллмана.....	18
3.2 Адаптивный алгоритм оптимального управления на основе итераций по стратегии (PI). 20	
3.3 Онлайн реализация алгоритма итерации по стратегии.....	20
4. Результаты моделирование.....	22
4.1 Регулятор для линейных систем.....	22
4.2 Исследование влияния времени выборки T на оптимальный адаптивный регулятор.....	26
4.3 Регулятор для нелинейных систем.....	29
4.4 Регулятор линейной системы для перевернутого маятника и тележки.....	32
4.5 Моделирование для маятника.....	35
5. Заключение.....	39
Список использованных источников.....	40

Оптимальные регуляторы обычно разрабатываются в автономном режиме путем решения уравнений Гамильтона-Якоби-Беллмана (HJB), например, уравнения Риккати, с использованием полных знаний о динамике системы. Такой регулятор имеет множество применений в науке, технике и исследованиях. Для расширения применимости оптимального управления необходима разработка алгоритмов оптимального управления, способных адаптироваться к изменению динамических свойств системы. Адаптивные регуляторы учатся в режиме реального времени управлять системами с неизвестными параметрами, используя данные, измеряемые в моменте вдоль траекторий системы. Однако адаптивные регуляторы обычно не сходятся к оптимальным решениям [1]. Метод обучения с подкреплением (RL) и адаптивное динамическое программирование стали прочной основой для разработки адаптивных оптимальных регуляторов. Такой метод предполагает причинно-следственную связь между действиями и вознаграждением или наказанием. Агент взаимодействует со своей средой посредством действия, и за этим действием следует вознаграждение (положительный сигнал RL), что дает снижение затрат на управление, а наказание (отрицательный сигнал RL) - увеличение затрат на управление. Действуя таким образом, алгоритм RL со временем обучается оптимальной стратегии. Алгоритм применяется для решения множества задач оптимального управления, включая обеспечение устойчивости, подавление шумов, оптимальное слежение за траекторией и т. д., без необходимости решения уравнений Гамильтона - Якоби - Беллмана, что позволяет решать задачу управления в условиях неполной информации о динамике объекта.

Список использованных источников:

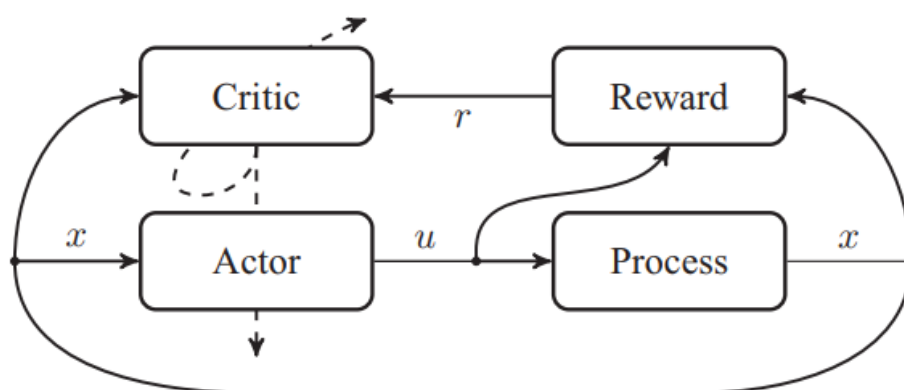
1. Optimal control, Frank L. Lewis, Draguna L. Vrabie, Vassilis L. Syrmos.-3rd ed, 2012.
2. Kyriakos G. Vamvoudakis, Draguna Vrabie, Frank L. Lewis, Online adaptive learning of optimal control solutions using integral reinforcement learning, 2011.
3. Bahare Kiumarsi, Member, IEEE, Kyriakos G. Vamvoudakis, Senior Member, IEEE, Hamidreza Modares, Member, IEEE, and Frank L. Lewis, Fellow IEEE, Optimal and Autonomous Control Using Reinforcement Learning: A Survey, 2018.

1 Обучение с подкреплением

1.1 Обзор обучения с подкреплением (RL)

Обучение с подкреплением — это одно из машинного обучения. Речь идет о принятии правильных действий для получения максимальной выгоды в конкретной ситуации. Обучение с подкреплением (RL) — это наука о принятии решений. Речь идет об обучении оптимальному поведению в окружающей среде для достижения максимального вознаграждения. В области теории управления обучение с подкреплением относится к методу, который позволяет разрабатывать адаптивные регуляторы в реальном времени для решения определяемых пользователем задач оптимального управления.

Обучение с подкреплением подразумевает наличие причинно-следственной связи между действиями и вознаграждением или наказанием. Это структура, в которой агент (или регулятор) оптимизирует свое поведение, взаимодействуя с окружающей средой. Совершив действие в каком-то состоянии, агент получает скалярное вознаграждение или наказание от среды, которое дает агенту представление о качестве этого действия. Алгоритмы RL исходят из идеи, что успешное поведение (приносящее высокие награды) будет запоминаться в том смысле, что оно имеет тенденцию использоваться повторно в последующих случаях [1]. Идея RL возникла из экспериментов обучения в биологии. RL теоретически тесно связан прямо и косвенно с адаптивными методами оптимального управления.



Рисунка 1. Схематический обзор алгоритма актер-критик [2]

Одной из популярных структур RL является структура «Актер-критик» [Barto, Sutton, Anderson 1983], Актерско-критические методы сочетают в себе преимущества методов "только актер" и "только критик" [3], в которой компонент «Актер»⁴ выполняет действия (стратегия управления), влияющие на окружающую среду, а секция «Критик» оценивает это действие. На основе этой оценки используется множество методов для калибровки или улучшения действия, чтобы новое действие создавало большую ценность (награду), чем

предыдущее. Таким образом, структура «Актер-Критик» состоит из двух этапов: оценка поведения и улучшение поведения. Поведенческая оценка проводится путем наблюдения за результатами, полученными из окружающей среды после выполнения определенного поведения.

Метод «Актер-критик» обычно обладают хорошей сходимостью, в отличие от методов, основанных только на критике [4].

1.2 Элемент обучения с подкреплением

В этом разделе представлены некоторые элементы и термины обучения с подкреплением. Можно также рассмотреть некоторые элементы обучения с подкреплением, представленные в [5].

В обучении с подкреплением есть термин «агент», который означает субъект, который взаимодействует с окружающей средой посредством действий.

Среда - это окружающее агента пространство, в котором он существует и взаимодействует.

Действия - это методы агента, которые позволяют ему взаимодействовать с окружающей средой и изменять ее. Основываясь на состоянии $S(t)$ текущей среды, агент будет выполнять действия $a(t)$

Стратегия - это то, что определяет, как агент действует в данный момент времени. Другими словами, стратегия - это отображение состояний среды на действия, которые будут выполняться в этих состояниях. Политика является ядром агента при определении поведения. В некоторых случаях политика может быть простой функцией или таблицей поиска. В других случаях политика может включать в себя обширные вычисления, например процесс поиска.

Вознаграждение - За каждое действие среда посылает агенту определенное вознаграждение. Цель агента - максимизировать общее вознаграждение, которое он получает в течение длительного периода времени. Сигнал о вознаграждении помогает определить, какие события являются хорошими, а какие плохими для агента, а также является основной базой для изменения стратегии. Если действие, выбранное в соответствии с стратегией, приносит низкое вознаграждение, стратегия может быть изменена. В будущем агент будет выбирать другие действия в аналогичных ситуациях.

Исследование и использование - Одна из проблем, возникающих в обучении с подкреплением - это компромисс между использованием и исследованием. Чтобы получить большее вознаграждение, агент должен отдавать предпочтение тем действиям, которые уже пробовал в прошлом и которые помогли ему достичь вознаграждения. Агент рассматривает все возможные действия для данного состояния, а затем выбирает действия, основываясь на максимальном значении этих действий. Это называется использованием, поскольку для принятия решения мы используем доступную информацию.

Кроме того, вместо выбора действий, основанных на максимальном будущем вознаграждении, агент может выбирать действия случайным образом. Случайные действия важны, поскольку они позволяют агенту исследовать и обнаруживать новые состояния, которые не были выбраны во время операции.

Агент должен использовать то, что он испытал, для получения вознаграждения, но он также должен исследовать, чтобы сделать лучший выбор в отношении будущих действий.

1.3 Марковский процесс принятия решений (МППР) - Markov Decision Processe

RL можно выразить с помощью марковского процесса принятия решений (МППР), как показано на рисунке 2. Каждое окружение представлено состоянием, которое отражает происходящее в окружении. Агент RL совершает действия в среде, которые вызывают изменение текущего состояния среды, порождая новое состояние, и получает вознаграждение в зависимости от результатов. Агент получает положительное вознаграждение за хорошие действия и отрицательное - за плохие, что помогает ему оценить выполненное действие в данном состоянии и учиться на опыте. Идея МППР представлена в статье [7].



Рисунок 2.1 - Взаимодействие агента и среды в марковском процессе принятия решений.

Рассмотрим МППР (X, U, P, R) , где X — набор состояний, а U — набор действий или элементов управления. Вероятности перехода $P: X \times U \times X \rightarrow [0,1]$ описывают что для каждого состояния $x \in X$, и действие $u \in U$, условная вероятность $P_{x,x'|u} = P_{r x' | x, u}$ перехода в состояние $x' \in X$ учитывая, что МППР находится в состоянии x и выполняет действие u . Функция ценности $R: X \times U \times X \rightarrow R$ это ожидаемая непосредственная стоимость $R_{x,x'|u}$ уплаченная после перехода в состояние $x' \in X$ учитывая, что МППР запускается в состоянии $x \in X$ и выполняет действие u . Свойство Маркова относится к тому факту, что вероятности перехода $P_{x,x'|u}$ зависят только от текущего состояния x , а не от истории того, как МППР достигла этого состояния.

Основная проблема МППР состоит в том, чтобы найти отображение $\pi: X \times U \rightarrow [0,1]$, которое дает для каждого состояния x и действия u условную вероятность $\pi_{x,u} = P_{u|x}$ выполнения действия u при условии, что МППР находится в состоянии x . Такое отображение называется управлением с обратной связью или стратегией действий или стратегией. Стратегия $\pi_{x,u} =$

$P_{ru|x}$ называется стохастической или смешанной, если существует ненулевая вероятность выбора более одного элемента управления в состоянии x . Смешанные стратегии можно рассматривать как векторы распределения вероятностей, имеющие в качестве компонента i вероятность выбора i -го управляющего воздействия в состоянии $x \in X$. Если отображение $\pi_{x,u} = P_{ru|x}$ допускает только одно управление с вероятностью единица, то в каждом состоянии x отображение называется детерминированной стратегией. Тогда $\pi_{x,u} = P_{ru|x}$ соответствует функции, отображающей состояния в управления $\mu_x : X \rightarrow U$.

МППР, которые имеют конечные пространства состояний и действий, называются конечными МППР.

1.3.1 Задачи оптимального последовательного решения

Обозначим значения состояния и действия в момент времени k через x_k, u_k . Зачастую желательно, чтобы системы, спроектированные человеком, были оптимальными с точки зрения экономии ресурсов, таких как стоимость, время, топливо и энергия.

Определение ценности этапа во время k по формуле:

$$r_k = r(x_k, u_k, x_{k+1}) \quad (1)$$

Тогда функция ценности:

$$R_{x,x'u} = E[r_k | x_k = x, u_k = u, x_{k+1} = x'] \quad (2)$$

с E - оператором ожидаемого значения.

Определение индекса производительности как сумму будущих затрат за интервал времени $[k, k+T]$:

$$J_{k,T} = \sum_{i=0}^{T-1} \gamma^i r_{k+i} = \sum_{i=0}^{T-1} \gamma^i [r(x_{k+i}, u_{k+i}, x_{k+i+1})] \quad (3)$$

где γ - коэффициент дисконтирования, уменьшающий вес затрат, понесенных в дальнейшем.

Предположим, что агент выбирает стратегию управления $\pi_k(x_k, u_k)$, которая используется на каждом этапе k МППР. Нас в первую очередь интересуют стационарные стратегии, где условные вероятности $\pi_k(x_k, u_k)$ не зависят от k . Тогда $\pi_{x,u} = \pi_{x,u} = P_{ru|x}$ для всех k . Нестационарная детерминированная стратегии имеет вид $\pi = \{\mu_0, \mu_1, \dots\}$, где каждая запись представляет собой функцию $\mu_k : X \rightarrow U$; $k=0, 1, \dots$. Стационарные детерминированные стратегии не зависят от времени, т. е. имеют вид $\pi = \mu, \mu, \dots$.

Выберите фиксированную стационарную стратегию $\pi_{x,u} = P_{ru|x}$. Тогда "замкнутая" МППР сводится к цепи Маркова с пространством состояний X . То есть вероятности переходов между состояниями фиксированы, и дальнейшая свобода выбора действий отсутствует. Вероятности переходов этой цепи Маркова задаются следующим образом:

$$p_{x,x'} \equiv P_{x,x'} = \sum_u P_{ru|x} P_{x',u} = \sum_u \pi_{x,u} P_{x',u} \quad (4)$$

где используется тождество Чепмена-Колмогорова [3].

Ценность стратегии определяется как условное ожидаемое значение будущих затрат при запуске в состоянии x в момент k и после этого следовании стратегии $\pi(x,u)$:

$$V_{\pi}(x) = E \{ J_k, T(x_k) = x \} = E \{ \sum_{i=k}^{\infty} \gamma^i c_i - k r_i \mid x_k = x \} \quad (5)$$

где $E \{ \dots \}$ - ожидаемое значение при условии, что агент придерживается стратегии $\pi(x,u)$, и $V_{\pi}(x)$ называется функцией ценности для стратегии $\pi(x,u)$, что является ценностью пребывания в состоянии x с учетом того, что является ценностью пребывания в состоянии x , учитывая, что стратегия $\pi(x,u)$.

Основная цель МППР — определить стратегию $\pi(x,u)$, позволяющую минимизировать ожидаемые будущие затраты

$$\pi^*(x,u) = \operatorname{argmin}_s V_{\pi}(x) \\ = \operatorname{argmin}_{\pi} E \{ \sum_{i=k}^{\infty} \gamma^i c_i - k r_i \mid x_k = x \} \quad (6)$$

Такая стратегия называется оптимальной, а соответствующее оптимальное значение задается как

$$V_k^*(x) = \min_{\pi} V_{\pi}(x) = \min_{\pi} E \{ \sum_{i=k}^{\infty} \gamma^i c_i - k r_i \mid x_k = x \} \quad (7)$$

1.3.2 Обратная рекурсия для значения

Используя тождество Чепмена-Колмогорова и свойство Маркова, значение стратегии $\pi(x,u)$ можно записать как

$$V_{\pi}(x) = E \{ J_k, T(x_k) = x \} = E \{ \sum_{i=k}^{\infty} \gamma^i c_i - k r_i \mid x_k = x \}$$

$$V_{\pi}(x) = E \{ r_k + \gamma V_{\pi}(x_{k+1}) - k r_k \mid x_k = x \},$$

$$V_{\pi}(x) = u_{\pi}(x,u) + P_{xx'}(u) R_{xx'}(u) + \gamma V_{\pi}(x_{k+1}) \mid x_{k+1} = x' \quad (8)$$

Поэтому функция ценности для стратегии $\pi(x,u)$ удовлетворяет:

$$V_{\pi}(x) = u_{\pi}(x,u) + P_{xx'}(u) R_{xx'}(u) + \gamma V_{\pi}(x_{k+1}) \quad (9)$$

Это уравнение обеспечивает обратную рекурсию для значения в момент времени k в терминах значения в момент времени $k+1$.

1.3.3 Обзор адаптивного динамического программирования

Оптимальную стоимость можно записать как

$$V_k^*(x) = \min_{\pi} V_{\pi}(x) = \min_{\pi} [u_{\pi}(x,u) + P_{xx'}(u) R_{xx'}(u) + \gamma V_{\pi}(x_{k+1})] \quad (10)$$

Принцип оптимальности Беллмана гласит: "Оптимальная стратегия обладает тем свойством, что независимо от того, какими были предыдущие управляющие воздействия, оставшиеся управления представляют собой оптимальную стратегию по отношению к состоянию, возникшему в результате этих предыдущих управлений". Поэтому мы можем записать:

$$V_k^*(x) = \min_{\pi} [u_{\pi}(x,u) + P_{xx'}(u) R_{xx'}(u) + \gamma V_{\pi}(x_{k+1})] \quad (11)$$

Предположим, что в момент времени k применяется произвольное управление u , а начиная с момента времени $k+1$ применяется оптимальная. Тогда принцип оптимальности Беллмана гласит, что оптимальное управление в момент времени k задается следующим образом:

$$\pi^*(x_k=x, u) = \operatorname{argmin}_{\pi} [u_{\pi}(x,u) + P_{xx'}(u) R_{xx'}(u) + \gamma V_{\pi}(x_{k+1})] \quad (12)$$

В предположении, что цепь Маркова, соответствующая каждой стратегии, с вероятностями перехода, является эргодической, каждая МППР имеет стационарную детерминированную оптимальную стратегию. Тогда мы можем

эквивалентно ² минимизировать условное ожидание по всем действиям u в состоянии x . Поэтому:

$$V_k^* x = \min_u \{ P_{xx'} u R_{xx'} u + \gamma V_{k+1}^* x' \} \quad (13)$$

$$u_k^* = \operatorname{argmin}_u \{ P_{xx'} u R_{xx'} u + \gamma V_{k+1}^* x' \} \quad (14)$$

Обратная рекурсия (11), (13) составляет основу динамического программирования (ДП), которое дает автономные методы для работы в обратном направлении во времени для определения оптимальной стратегии. ДП - это автономная процедура поиска оптимального значения и оптимальных стратегий, которая требует знания полной динамики системы в виде вероятностей перехода $P_{x,x'} u = P_{x'x} x, u$ и ожидаемых затрат $R_{xx'} u = E\{ r_k \mid x_k = x, u_k = u, x_{k+1} = x' \}$.

1.3.4 Уравнение Беллмана и уравнение оптимальности Беллмана

Обучение с подкреплением включает в себя определение решений для управления в реальном времени и с упреждением во времени. Ключом к этой проблеме является уравнение Беллмана, которое будет представлено ниже.

Динамическое программирование — это метод поиска оптимального значения и стратегия в обратном времени. В отличие от этого, обучение с подкреплением направлено на поиск оптимальной стратегии на основе каузального опыта путем выполнения последовательных решений, которые улучшают управляющие действия на основе наблюдаемых результатов использования текущей стратегии. ²

Эта процедура требует выведения методов поиска оптимальных значений и оптимальных стратегий, которые могут быть выполнены в прямом времени.

Чтобы получить методы, направленные в прямом времени, для поиска оптимальных значений и оптимальных стратегий, устанавливается временной горизонт T равным бесконечности и определяется стоимость бесконечного горизонта

$$J_k = \sum_{i=0}^{\infty} \gamma^i r_{k+i} = \sum_{i=k}^{\infty} \gamma^i r_i - k r_i \quad (15)$$

Бесконечная функция ценности связанного горизонта для стратегии $\pi(x, u)$ равна

$$V_{\pi} x = E_{\pi} \{ J_k x_k = x \} = E_{\pi} \{ \sum_{i=k}^{\infty} \gamma^i r_i - k r_i \mid x_k = x \} \quad (16)$$

Используя (4) с $T = \infty$, можно увидеть, что функция ценности стратегии $\pi(x, u)$ удовлетворяет уравнению Беллмана

$$V_{\pi} x = u_{\pi} x, u \{ x' P_{xx'} u R_{xx'} u + \gamma V_{\pi} x' \} \quad (17)$$

Важность этого уравнения заключается в том, что в нем с обеих сторон фигурирует одна и та же функция стоимости, что связано с тем, что используется бесконечный горизонт стоимости. Поэтому уравнение Беллмана (17) можно интерпретировать как уравнение согласованности, которому должна удовлетворять функция стоимости на каждом временном этапе. Оно выражает связь между текущей стоимостью нахождения в состоянии x и стоимостью нахождения в следующем состоянии x' при использовании стратегии $\pi(x, u)$. ²

Если МППР конечен и имеет N состояний, то уравнение Беллмана (17) ³⁸ представляет собой систему N одновременных линейных уравнений для

значения $v_\pi(x)$ пребывания в каждом состоянии x при текущей стратегии π, u .

Оптимальное значение бесконечного горизонта удовлетворяет

$$V^* x = \min_{\pi} v_{\pi}(x) = \min_{\pi} u_{\pi} x, u x' P_{xx'} u R_{xx'} u + \gamma v_{\pi} x' \quad (18)$$

Тогда принцип оптимальности Беллмана приводит к уравнению оптимальности Беллмана

$$V^* x = \min_{\pi} v_{\pi}(x) = \min_{\pi} u_{\pi} x, u x' P_{xx'} u R_{xx'} u + \gamma V^* x' \quad (19)$$

Эквивалентно, при условии эргодичности цепей Маркова, соответствующих каждой стратегии, уравнение оптимальности Беллмана можно записать как

$$V^* x = \min_{u x' P_{xx'} u R_{xx'} u + \gamma V^* x' \quad (20)$$

Это уравнение известно как уравнение Гамильтона – Якоби – Беллмана (HJB) в системах управления. Если МППР конечен и имеет N состояний, то уравнение оптимальности Беллмана представляет собой систему N нелинейных уравнений для оптимального значения $v^*(x)$ пребывания в каждом состоянии. Оптимальное управление определяется выражением

$$u^* = \operatorname{argmin}_{u x' P_{xx'} u R_{xx'} u + \gamma V^* x' \quad (21)$$

1.4 Итерация по стратегии и итерация по критерию

Чтобы найти оптимальную стратегию для данной задачи, в обучении с подкреплением обычно используются два метода: итерация по стратегии и итерация по критерию.

16

Эти два метода являются методами автономной оптимизации в обучении с подкреплением. Чтобы иметь возможность реализовать методы онлайн-оптимизации, необходимо объединить их с методами приближенного динамического программирования. Это то, что нам нужно для построения алгоритма оптимального адаптивного регулятора в реальном времени.

1.4.1 Итерация по стратегии алгоритм

Чтобы понять алгоритм итерации по стратегии, нам сначала нужно обобщить функцию значения.

$$v_{\pi} x = u_{\pi} x, u x' P_{xx'} u R_{xx'} u + \gamma v_{\pi} x' \quad (22)$$

Это равенство называется равенством Беллмана и упрощает вычисление функции ценности за счет использования динамического программирования.

Алгоритм итерации по стратегии представлен следующим образом:

Выберите начальную стратегию $\pi_0(x, u)$. Начиная с $j = 0$, повторите по j до сходимости

Оценка стратегии (обновление ценности)

$$v_j x = u_{\pi_j} x, u x' P_{xx'} u R_{xx'} u + \gamma v_j x', \forall x \in X \quad (23)$$

Улучшение стратегии (обновление стратегии)

$$\pi_{j+1} x, u = \operatorname{argmin}_{\pi} \pi x' P x x' u R x x' u + \gamma V_j \pi x', \forall x \in X \quad (24)$$

На каждой итерации j , алгоритм итерации по стратегии определяет решение уравнения Беллмана для определения значения функции ценности $v_j(x)$, соответствующей текущей стратегии $\pi_j(x, u)$. Тогда улучшение стратегии осуществляется благодаря уравнению (24). Эти шаги повторяются до тех пор, пока и функция ценности, и стратегия придут к своим оптимальным значениям. Чтобы алгоритм работал корректно, нам нужно обратить внимание на начальную инициализацию стратегии $\pi_0(x, u)$ и начальное значение $v_1 \leq v_0$.

1.4.2 Итерации по критерию

Алгоритм итерация по критерию основан на чрезвычайно простой, но очень эффективной идее. Рассматривая алгоритм, итерация по стратегии, мы видим, что он одновременно поддерживает и стратегию, и функции ценности, что делает вычисления громоздкими и трудоемкими. Алгоритм итерация по критерию начинает с того, что сначала пытается оптимизировать функцию ценности, затем стратегия, соответствующая оптимальной функции ценности, конечно же, также будет оптимальной стратегией.

Алгоритм итерации по стратегии представлен следующим образом:

Выберите начальную стратегию $\pi_0(x, u)$. Начиная с $j = 0$, итерации по j до сходимости

Обновление ценности

$$V_{j+1} x = \pi_j x, u x' P x x' u R x x' u + \gamma V_j \pi x', \forall x \in S_j \in X \quad (25)$$

Улучшение стратегии

$$\pi_{j+1} x, u = \operatorname{argmin}_{\pi} \pi x' P x x' u R x x' u + \gamma V_{j+1} \pi x', \forall x \in S_j \in X \quad (26)$$

В каждом цикле j алгоритм итерация по критерию оценивает значение функции $v_{j+1}(x)$, соответствующее стратегии управления $\pi_j(x, u)$, на основе значения функции ценности на предыдущем этапе. Обновление стратегии производится по формуле (26). Повторяйте шаги алгоритма до сходимости.

1.4.3 Сравнение итерации по стратегии и итерации по критерию

В то время как значение $v_j(x')$ в алгоритме итерации по стратегии представляет собой фактическое значение текущей стратегии π_j , в алгоритме

итерации по критерию его можно рассматривать как оценку стоимости перехода из состояния x в будущее состояние x' .

В алгоритме итерации по критерию в каждом цикле мы обновляем значение функции ценности $V(x)$, и когда $V(x)$ достигает оптимального значения, стратегия автоматически становится оптимальной. В отличие от итерации по критерию, итерация по стратегии пытается определить эффективность текущей стратегии с помощью оценки стратегии, а затем обновляет новую стратегию, после чего стратегия и значение функции ценности постоянно обновляются до тех пор, пока не будет достигнуто оптимальное значение.

Итерации по стратегии обычно сходятся быстрее, чем итерации по критерию. Поэтому во многих случаях он предпочтительнее. При определенных условиях гарантируется сходимость алгоритма итерации по стратегии за ограниченное число шагов. Однако количество шагов, необходимых для сходимости алгоритм итерации по критерию, не обязательно ограничено.

2. Регулятора для линейных систем

Рассмотрим непрерывную линейную и инвариантную систему во времени систему с моделью в пространстве состояний.

$$\dot{x} = Ax(t) + Bu(t), \quad (2.1)$$

Где $B \in \mathbb{R}^{n \times m}$, $A \in \mathbb{R}^{n \times n}$, $x(t) \in \mathbb{R}^n$ – состояния системы, $u(t) \in \mathbb{R}^m$ – входной сигнал управления, пара матрицы (A, B) – стабилизируемость.

Задача заключается в том, что нужно разрабатывать стратегию управления с обратной связью $u = -Kx$ так, чтобы

$$\lim_{t \rightarrow \infty} |x(t)| = 0 \quad (2.2)$$

2.1 Оптимальный регулятор

Одним из наиболее известных типов оптимальных регуляторов является линейный квадратичный регулятор (LQR). Для системы (2.1), задача LQR с обратной связью по состоянию $u = -Kx$ состоит в поиске оптимального регулятора u^* , который минимизирует функцию стоимости $J(x, u)$ на бесконечном горизонте, связанную с системой:

$$J_{x,u} = \int_0^{\infty} x^T(t)Qx(t) + u^T(t)Ru(t) dt \quad (2.3)$$

Здесь можно показать, что оптимальное управление является линейной функцией от x

$$u = -Kx = R^{-1} B^T P x \quad (2.4)$$

где P — определенное положительное решение алгебраического уравнения Риккати

$$A^T P + P A - P B R^{-1} B^T P + Q = 0 \quad (2.5)$$

Если $\lambda \geq 0$, $\lambda > 0$, (A, B) стабилизируемость и (Q, A) обнаруживаемость, то алгебраическое уравнение Риккати имеет единственное симметричное положительное решение и соответствующий ему регулятор делает систему (3.1) асимптотически устойчивой.

2.2 Адаптивный регулятор

В случае адаптивного регулятора для системы (2.1) матрица A считается неизвестной. Полагая $A = A_0 + b \theta^T$, уравнение системы (2.1) принимает вид

$$\dot{x} = A_0 x + b \theta^T x + u \quad (2.6)$$

Где θ - неизвестный параметр и A_0 — эталонная матрица Гурвица.

Сигнал управления имеет следующую формулу

$$u = -\hat{\theta}^T x \quad (2.7)$$

Где $\hat{\theta}$ - оценка вектора, динамическая модель системы будет иметь вид

$$\dot{x} = A_0 x + b \hat{\theta}^T x - b \theta^T x = A_0 x + b \tilde{\theta}^T x \quad (2.8)$$

Где $\tilde{\theta} = \hat{\theta} - \theta$ - вектор параметрических ошибок.

Алгоритм адаптивной стабилизации объекта (2.1) имеет вид

$$\dot{\hat{\theta}} = \gamma x^T P b = \gamma x^T \hat{b}^T P x \quad (2.9)$$

Где симметричная положительно определенная матрица P является решением уравнения Ляпунова

$$A_0^T P + P A_0 + Q = 0, \quad Q = Q^T > 0 \quad (2.10)$$

2.3 Оптимальный адаптивный регулятор

Для системы (2.1) задача заключается в том, что, нужно найти оптимальный регулятор u^* , который минимизирует квадратичную функцию стоимости бесконечного горизонта, связанную с системой.

$$V(x(t), t) = \int_t^{\infty} x^T(\tau) Q x(\tau) + u^T(\tau) R u(\tau) d\tau \quad (2.11)$$

Где для системы (2.1) матрица A неизвестна, $Q \geq 0$, $R \geq 0$ и Q, A – обнаруживаемость.

Тогда оптимальный регулятор записывается следующим образом:

$$u^*(t) = \arg \min_{u(t)} V(t, x(t), u(t)), \quad t_0 \leq t \leq \infty \quad (2.12)$$

Решение этой задачи оптимального управления, определенное в соответствии с принципом оптимизации Беллмана, дается с помощью обратной связи по состоянию $u(t) = -Kx(t) = -R^{-1} B^T P x(t)$.

Пусть K – стабилизирующий коэффициент обратной связи по состоянию для (2.1) такой, что $x = A - BK$ представляет собой устойчивую замкнутую систему. Тогда квадратичная бесконечная стоимость или ценность на бесконечном горизонте определяется выражением

$$V(x(t)) = \int_t^{\infty} (x^T(\tau) Q x(\tau) + u^T(\tau) R u(\tau)) d\tau, \quad u = -Kx$$

$$V(x(t)) = \int_t^{\infty} x^T(\tau) Q + K^T R K x(\tau) d\tau = x^T(t) P x(t) \quad (2.13)$$

Где P – симметричная положительно определенная матрица, является решением матричного уравнения Ляпунова

$$A - BK)^T P + P (A - BK) = -(K^T R K + Q) \quad (2.14)$$

Тогда $V(x(t))$ служит функцией Ляпунова для (2.1) с коэффициентом усиления регулятора K . Функция ценности (2.8) можно записать в следующем виде

$$V(x(t)) = \int_t^{\infty} x^T(\tau) Q + K^T R K x(\tau) d\tau + V(x(t+T)) \quad (2.15)$$

Это уравнение Беллмана для задачи LQR. Используя уравнение Беллмана IRL, можно обойти отмеченные в разделе (1.3.2) проблемы применения обучения с подкреплением к системам с непрерывным временем.

Интегральное армирование

$$\rho(x(t), t, T) = \int_t^{t+T} x^T(\tau) Q + K^T R K x(\tau) d\tau \quad (2.16)$$

Формула (2.15) называется уравнением Беллмана обучения с интегральным подкреплением (IRL).

Далее показано, как найти оптимальный регулятор, применяя интегральное обучение с подкреплением к линейным системам (итерации по стратегии) (2.1).

$$W_k^T [\phi(x(t)) - \phi(x(t+T))] = \int_t^{t+T} x^T(\tau) (Q + K_k^T R K_k) d\tau.$$

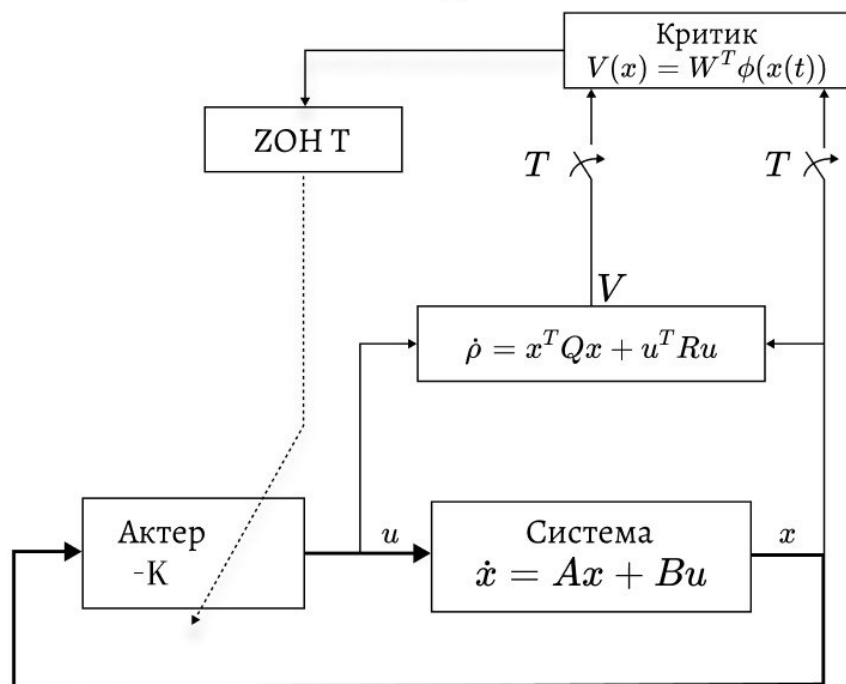


Рисунок 2.1 Структура системы с оптимальным адаптивным регулятором для линейных систем.

На основе интегрального обучения с подкреплением (2.15) можно написать следующий алгоритм RL:

$$V_{x(t)} = \frac{1}{T} \int_t^{t+T} x^T(\tau) (Q + K_i^T R K_i) d\tau + x^T(t+T) P x(t+T) \quad (2.16)$$

Параметры функции ценности $V_i(x(t)) = x^T(t) P x(t)$ на итерации i IRL являются элементами симметричной ядерной матрицы P_i . Они должны быть найдены на каждой итерации i путем измерения данных $(x(t), x(t+T), u(t))$ в

моменты времени t вдоль траекторий системы. Для вычисления этих параметров термин $x^T(t) P_i x(t)$ записывается как

$$x^T(t) P_i x(t) = w^T \varphi(x(t)) \quad (2.17)$$

Где $\varphi(x(t))$ – вектор функции активации. $\varphi(x(t))$ обозначает базисный вектор квадратичного полинома произведения Кронекера

$$\varphi(x(t)) = x_1^2 x_1 x_2 \dots x_1 x_n x_2^2 x_2 x_3 \dots x_n^2 \quad (2.18)$$

Применяя формулу аппроксимации значения функции (2.17), можно записать (2.15) следующим образом

$$w_{j+1}(\varphi(x(t)) - \varphi(x(t+T))) = \int_t^{t+T} x^T \tau Q + K_i^T R K_i x \, d\tau = d(\varphi(x(t)), K_i), \quad (2.17)$$

Применяя к этой проблеме структуру «Актер-Критик», Критик использует уравнение (2.17) для вычисления матрицы P_i . Когда Критик сходится к решению P_i , Актер рассчитает новое значение управления для управления системой на основе P_i по формуле:

$$u_{i+1} = -R^{-1} B^T P_i x = K_{i+1} x \quad (2.16)$$

С помощью использования алгоритма итерации по стратегии мы можем найти оптимальный регулятор без необходимости знания матрицы для линейной системы.

3. Задача оптимального управления для нелинейных систем

К непрерывным системам обучение с подкреплением применить сложнее, чем к дискретным, и оно дает меньше результатов. Один из методов, описанных в [13], [14], называется обучением с интегральным подкреплением. Этот метод может быть использован для разработки методов оптимального управления непрерывными системами, не требуя полного понимания динамики системы. Этот метод называется методом оптимального адаптивного управления. Ниже изучается и рассматривается оптимальный адаптивный регулятор для нелинейных систем.

3.1 Оптимальное управление и уравнение Гамильтона Якоби – Беллмана

Рассмотрим нелинейную динамическую систему с непрерывным временем

$$\dot{x} = f(x) + g(x)u \quad (3.1)$$

38

Где $x(t) \in R^n$ – состояние системы, входное управление $u(t) \in R^m$, точка равновесия $x=0$, $f(0)=0$ и $f(x) + g(x)u$ Липшиц (Lipschitz) на съемочной площадке $\Omega \in R^n$. Предположим, что система стабилизируема на Ω , т. е. существует непрерывная функция управления $u(t)$ такая, что замкнутая система асимптотически устойчива на Ω .

Определите меру эффективности или функцию ценности, которая имеет значение, связанное со стратегией управления с обратной связью $u = \mu(x)$, заданное как

$$V_\mu(x) = \int_0^\infty r(x, u) dt \quad (3.2)$$

Где $r(x, u) = Q(x) + u^T R u$, положительно определенный $Q(x)$, т. е. $Q(x) > 0$ для всех x и $x=0 \Rightarrow Q(x)=0$ и положительно определенная матрица $R = R^T \in R^{m \times m}$.

Для LQR с непрерывным временем приведенные выше выражения имеют вид как показан в пункте 2.

$$\dot{x} = Ax + Bu \quad (3.3)$$

$$V_\mu(x) = \int_0^\infty (x^T Q x + u^T R u) dt. \quad (3.4)$$

Определение 1. Стратегия управления $\mu(x)$ определяется как допустимая относительно (3.2) на Ω , обозначаемая $u \in \Psi(\Omega)$, если $\mu(x)$ непрерывна на Ω , $u(0)=0$, $\mu(x)$ стабилизирует (3.1) на Ω и $V(x_0)$ конечно $\forall x_0 \in \Omega$.

Для любой допустимой стратегии управления $\mu \in \Psi(\Omega)$, если соответствующая функция стоимости $V(x_0)$ равна c_1 , то бесконечно малым эквивалентом (3.2) является уравнение Беллмана

$$0 = r(x, \mu(x)) + \nabla V_\mu^T (f(x) + g(x)\mu(x)), \quad V_\mu(0) = 0 \quad (3.5)$$

Где ∇V_μ , взятый здесь как вектор-столбец, обозначает вектор градиента функции ценности V_μ по x . Учитывая допустимый регулятор $\mu(x) \in \Psi(\Omega)$, уравнение Беллмана можно решить через соответствующую функцию ценности $V_\mu(x)$. В линейном случае оно становится уравнением Ляпунова. Учитывая, что $\mu(x)$ является допустимой стратегией управления, если $V_\mu(x)$ удовлетворяет

(3.5) с $r(x, \mu(x)) \geq 0$, то можно доказать, что $v_\mu(x)$ является функцией Ляпунова для системы (1) со стратегией управления $\mu(x)$.

Теперь можно сформулировать задачу оптимального управления: для непрерывной системы (3.1), множества $u \in \Psi(\Omega)$ допустимых стратегий управления и функции ценности на бесконечном горизонте (3.2) найти допустимую стратегию управления такую, что индекс ценности (3.2), связанная с системой (3.1), минимизируется.

Определите гамильтониан

$$H(x, \mu, \nabla v_\mu) = r(x, \mu) + \nabla v_\mu^T (f(x) + g(x)u) \quad (3.6)$$

Тогда оптимальная функция ценности $v^*(x)$ удовлетворяет уравнению НЖВ

$$0 = \min_\mu H(x, \mu, \nabla v^*) \quad (3.7)$$

и удовлетворяющие оптимальному управлению

$$\mu^* = \operatorname{argmin}_\mu H(x, \mu, \nabla v^*) \quad (3.8)$$

$$\mu^* = - (R^{-1})^T g^T(x) \nabla v^*(x) \quad (3.9)$$

Подставляя эту стратегию оптимального управления в гамильтониан, мы получаем формулировку уравнения НЖВ в терминах $v^*(x)$:

$$0 = Q(x) + \nabla v^{*T}(x) f(x) - \frac{1}{2} \nabla v^{*T}(x) g(x) R^{-1} g^T(x) \nabla v^*(x), \quad (3.10)$$

$$\nabla v^*(0) = 0$$

Это достаточное условие функции оптимального значения [10]. В случае линейной системы с квадратичной функцией стоимости это уравнение НЖВ становится уравнением Риккати.

Чтобы найти оптимальное управляющее решение задачи, нужно всего лишь решить НЖВ (3.10) для функции ценности, а затем подставить решение в (3.9) для получения оптимального управления. Однако решение уравнения НЖВ, как правило, затруднено, поскольку оно представляет собой нелинейное дифференциальное уравнение второго порядка, следующее за градиентом функции стоимости, и его решение также требует полного знания системы системной динамики (т. е. динамики системы, описываемой уравнением

функции $f(x)$ и $g(x)$ должны быть известны). Более того, глобально гладких решений может не существовать.

3.2 Адаптивный алгоритм оптимального управления на основе итераций по стратегии (PI).

Итерация по стратегии (PI) — это итерационный метод обучения с подкреплением [4, 12] для решения (3.10) на основе простых уравнений. PI состоит из двух этапов: улучшение стратегии на основе (3.9) и оценка стратегии на основе Беллмана (3.5). Говорят, что PI имеет структуру обучения с подкреплением «Актер-Критик», где улучшение стратегии выполняется актером, а оценка стратегии выполняется критиком, решая уравнение Беллмана. Алгоритм итерация по стратегии задается следующим образом.

Пусть $u(x)$ — допустимая стратегия для (3.1), такая что замкнутая система асимптотически устойчива на Ω . Тогда стоимость бесконечного горизонта для любого $x \in \Omega$ определяется формулой (3.2), а $V(x(t))$ служит функцией Ляпунова для (3.1). Функцию стоимости (3.2) можно записать в виде

$$V_{\mu}(x) = \int_0^{\infty} \text{Tr}(x^T \mu x + V_{\mu}(x)) dt \quad (3.11)$$

На основании (3.11) и (3.7), учитывая начальную допустимую стратегию управления $u_0(x)$, можно вывести следующую схему итерации по стратегии

1. Решить для $V_{\mu_i}(x)$ с использованием

$$V_{\mu_i}(x) = \int_0^{\infty} \text{Tr}(x^T \mu_i x + V_{\mu_i}(x)) dt \quad (3.12)$$

2. Обновить стратегию управления с помощью

$$\mu_{i+1}(x) = \arg\min_{\mu} \{H(x, \mu, V_{\mu_i})\} \quad (3.13)$$

Где

$$\mu_{i+1}(x) = - (R + \gamma \int_0^{\infty} \text{Tr}(x^T \mu_i x + V_{\mu_i}(x)) dt)^{-1} \nabla V_{\mu_i}(x) \quad (3.14)$$

Уравнения (3.13) и (3.14) формулируют новый алгоритм итерации по стратегии для поиска оптимального управления без использования каких-либо знаний о внутренней динамике системы $f(x)$.

3.3 Онлайн реализация алгоритма итерации по стратегии

Применяя аппроксимационные возможности параллельной нейронной сети для уменьшения объема вычислений по сравнению со структурой «Актер-критик», этот алгоритм использует нейронную сеть для аппроксимации

оптимальной функции стоимости $V(x)$ с помощью $x \in \Omega$ следующим образом:

$$V(x) = w^T \varphi(x) \quad (3.15)$$

Где $w^T \in \mathbb{R}^N$ – неизвестная веса, N - количество нейронов, $\varphi(x)$ – вектор функции активации.

Используя нейронную сеть для аппроксимации оптимальной функции ценности, подставьте формулу (3.15) в формулу (3.12), чтобы получить:

$$w^T \varphi(x(t)) = t + Tr_{xt}, u_{ix} dt + w^T \varphi(x(t+T)) \quad (3.16)$$

Появляется ошибка $e(t)$, которая является ошибкой аппроксимации функции Беллмана.

$$e(t) = w^T \varphi(x(t+T)) - \varphi(x(t)) = -t + Tr_{xt}, u_{ix} dt \quad (3.17)$$

Обозначим что

$$h(t) = \varphi(x(t+T)) - \varphi(x(t)) \quad (3.18)$$

$$y(t) = -t + Tr_{xt}, u_{ix} dt \quad (3.19)$$

Уравнение (3.16) можно написать в следующий вид

$$e(t) = w^T h(t) + y(t) \quad (3.20)$$

$$\text{Где } H = [h(t_1), \dots, h(t_N)]^T, Y = [y(t_1), \dots, y(t_N)]^T \quad (3.21)$$

Уравнение (3.20) является линейной функцией по параметру w . Следовательно, мы можем применить алгоритм наименьших квадратов ошибки, чтобы найти оптимальное значение для w .

Данные системы собирается из N различных выборок за период времени T , поэтому мы вычисляем (3.19) в n точках от $t_1 \rightarrow t_N$, чтобы получить следующие функции:

$$H = [h(t_1), \dots, h(t_N)]^T \quad (3.22)$$

$$Y = [y(t_1), \dots, y(t_N)]^T \quad (3.23)$$

Чтобы определить веса W нейронной сети, аппроксимирующей функцию V , что приводит к минимизации следующей целевой функции

$$S = \int_{\Omega} e_{x,T} dx \quad (3.24)$$

Произведение $\langle f, g \rangle = \int_{\Omega} f g dx$ интеграла Лебега можно записать:

$$\langle e_{x,T}, dW \rangle = 0 \quad (3.25)$$

Использовать уравнение (3.20) для уравнения (3.25), получили

$$H_{NW+Y} = 0 \quad (3.26)$$

$$W = -N^{-1} H^T Y \quad (3.27)$$

Алгоритм онлайн обучения с интегральным подкреплением использует нейронные сети представлен следующим образом

Шаг 1. $\forall x \in \Omega_x$, инициализировать допустимый закон управления $u(x) \in \Psi(\Omega)$. Поместите в систему управляющий сигнал u_0 и соберите необходимые данные системы о состоянии и управляющих сигналах на N различных выборках за период времени T . Назначьте $i \leftarrow 0$, инициализируйте ϵ_w

Шаг 2. Используйте данные, собранную о системе, для расчета H и Y . Определить W из уравнения (3.27)

Шаг 3: обновите закон управления для следующего цикла

$$u_{i+1}(x) = -\frac{1}{2} R^{-1} g^T \nabla_v x W_i \quad (3.28)$$

Если критерий сходимости удовлетворяется так, что $\|W_{i+1} - W_i\| < \epsilon_w$, алгоритм завершается. Если не удовлетворено, присвойте $i \leftarrow i+1$, подайте сигнал u_i в систему и соберите необходимую информацию системы о состоянии и управляющих сигналах в N различных выборках за период T , затем вернитесь к шагу 2.

4. Результаты моделирование

В этом разделе для применения и исследования алгоритмов, представленных ранее в разделах 2 и 3, была выбрана система перевернутого маятника и тележки [17][18].

4.1 Регулятор для линейных систем

Рассмотрим непрерывную инвариантную во времени линейную систему

$$\dot{x} = Ax + Bu, \quad (4.1)$$

$$A = \begin{bmatrix} 0 & 1 \\ -10 & 0 \end{bmatrix}, B = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

Где $x \in \mathbb{R}^2$ – состояние системы, $u \in \mathbb{R}^1$ – сигнал управления системы

Задача оптимального управления состоит в оптимизации следующей функции

$$J_{x,u} = \int_0^{\infty} x^T Q x + u^T R u dt \quad (4.2)$$

$$Q = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, R = 1$$

В случае оптимального регулятора (LQR), решая уравнение Риккати

$$A^T P + P A - P B R^{-1} B^T P + Q = 0$$

Мы получали матрицу $P = \begin{bmatrix} 2.0739 & 0.2110 \\ 0.2110 & 0.0672 \end{bmatrix}$, следовательно обратная матрица K вычислена

$$K = R^{-1} B^T P = \begin{bmatrix} -2.2840 & -0.278 \end{bmatrix}$$

Параметры для адаптивного регулятора

$$\gamma = 1, A_0 = A - B K = \begin{bmatrix} -2.2840 & 0.7217 \\ 0.7151 & -10.2783 \end{bmatrix}, P_a = P = \begin{bmatrix} 2.0739 & 0.2110 \\ 0.2110 & 0.0672 \end{bmatrix}$$

Далее представлены результаты моделирования

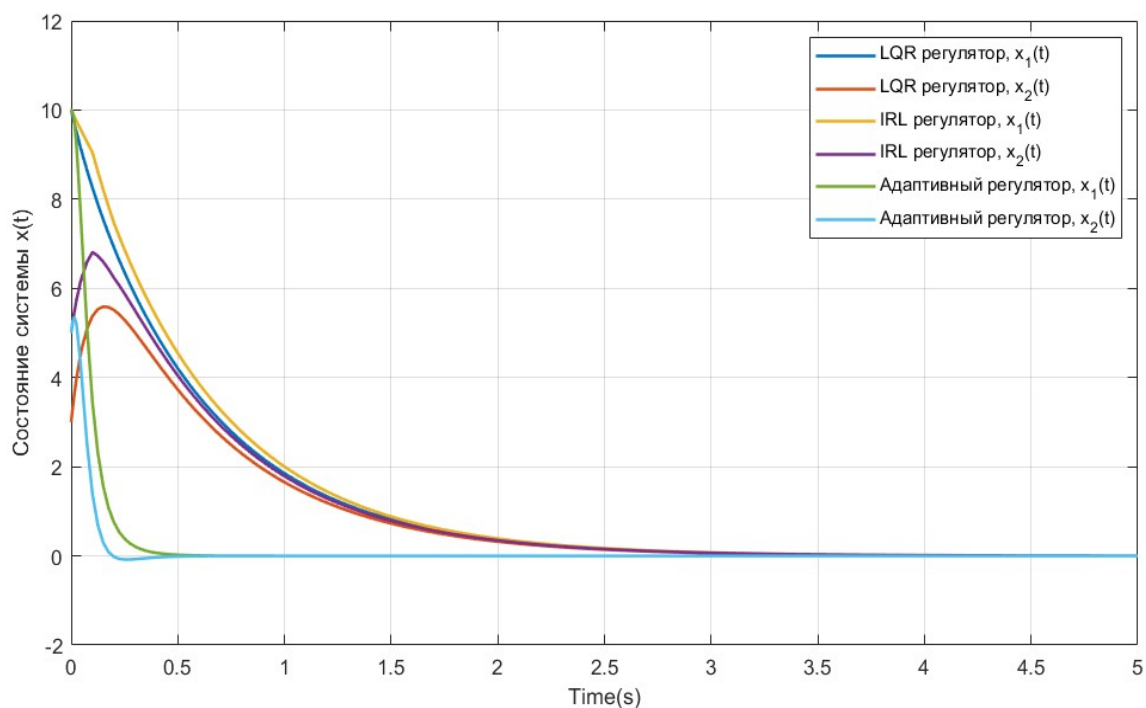


Рисунок 4.1 - Графика вектор состояния $x(t)$ для оптимального регулятора, адаптивный регулятор и оптимального адаптивного регулятора.

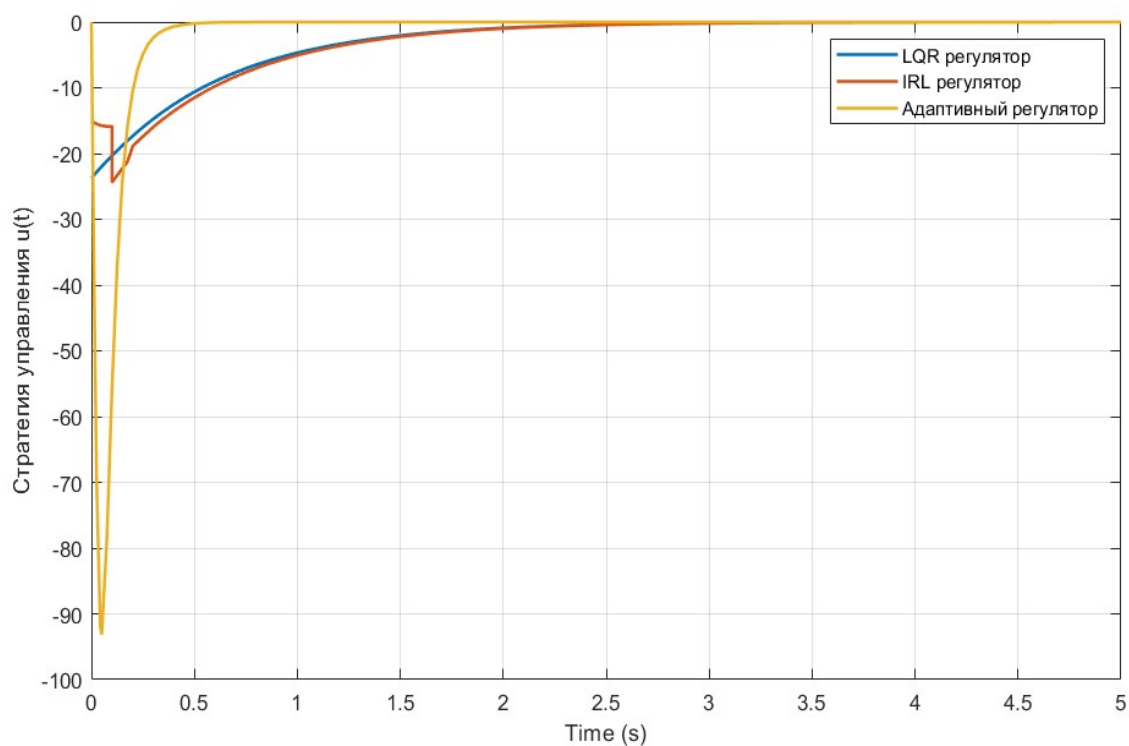


Рисунок 4.2 - Графика сигнал управления для оптимального регулятора, адаптивный регулятор и оптимального адаптивного регулятора (IRL)

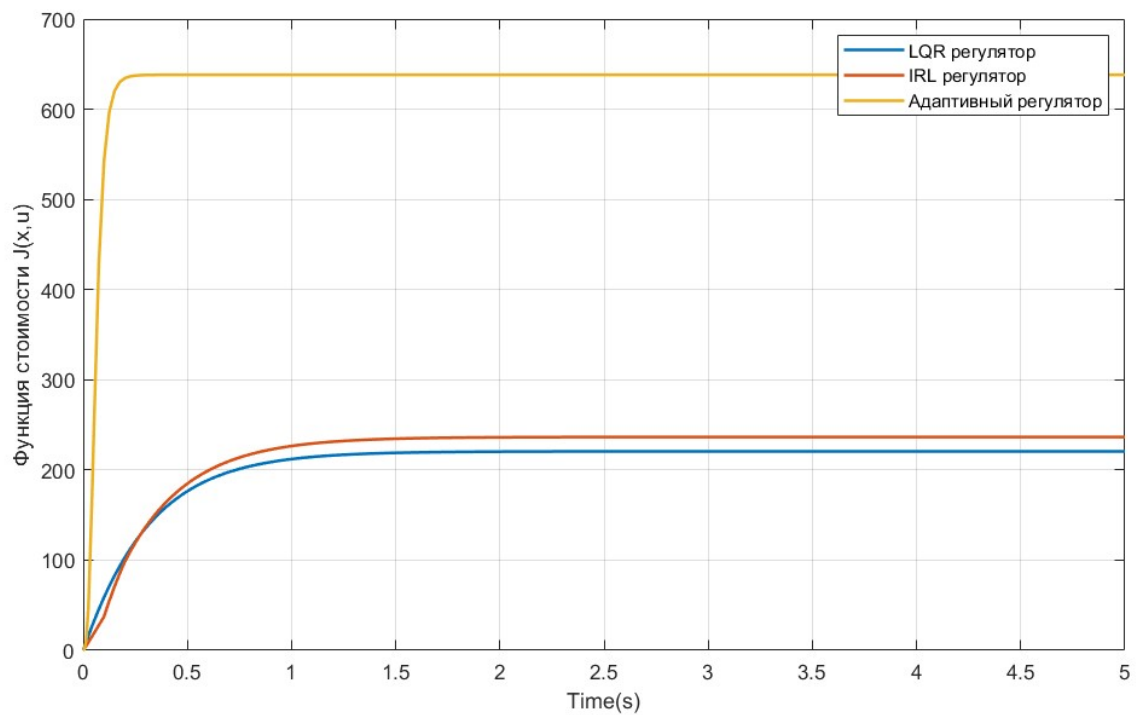


Рисунок 4.3 - Значение функции стоимости для оптимального регулятора, адаптивный регулятор и оптимального адаптивного регулятора (IRL)

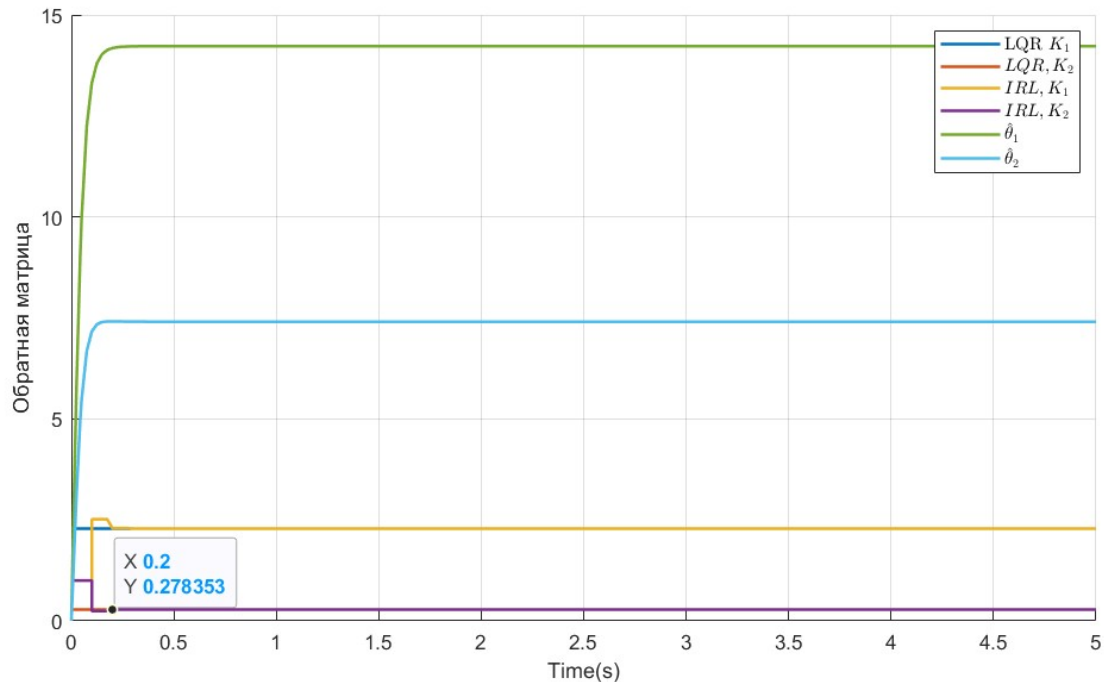


Рисунок 4.4 – Графика матрицы регулятора K (LQR), параметры θ для адаптивного регулятора и сходимости значения обратной матрицы оптимального регулятора (IRL)

Приведенные выше результаты моделирования показывают, что с помощью алгоритма итерации по стратегии обучения с подкреплением оптимальный адаптивный регулятор сошелся к асимптоте оптимального управления через 0.2 секунды. При этом адаптивный регулятор не сходится к оптимальному управлению.

Кроме того, производительность адаптивного регулятора зависит от выбора параметров γ и ρ_a для системы. Поэтому при изменении значений матрицы Q и R нам приходится заново выбирать параметры оптимального регулятора. Для оптимального адаптивного регулятора нам не нужно настраивать параметры при изменении Q и R .

4.2 Исследование влияния времени выборки T на оптимальный адаптивный регулятор.

При применении алгоритма обучения с интегральным подкреплением к оптимальному адаптивному регулятору нам нужно точно аппроксимировать оптимальную матрицу P , используя данные по траектории системы на интервалах времени T . Графики влияния времени выборки T к производительности оптимального адаптивного регулятора представлены ниже.

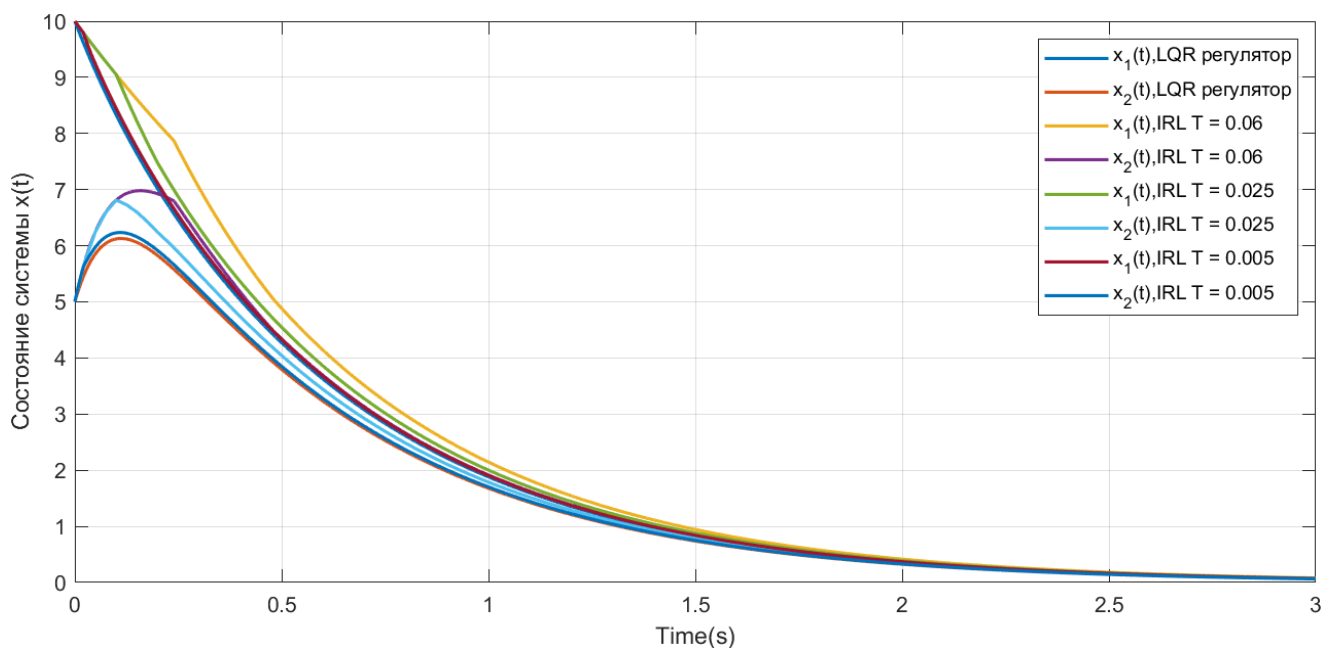


Рисунок 4.5 – Графика состояния системы $x(t)$ оптимального регулятора и оптимального адаптивного регулятора при некоторых значениях T

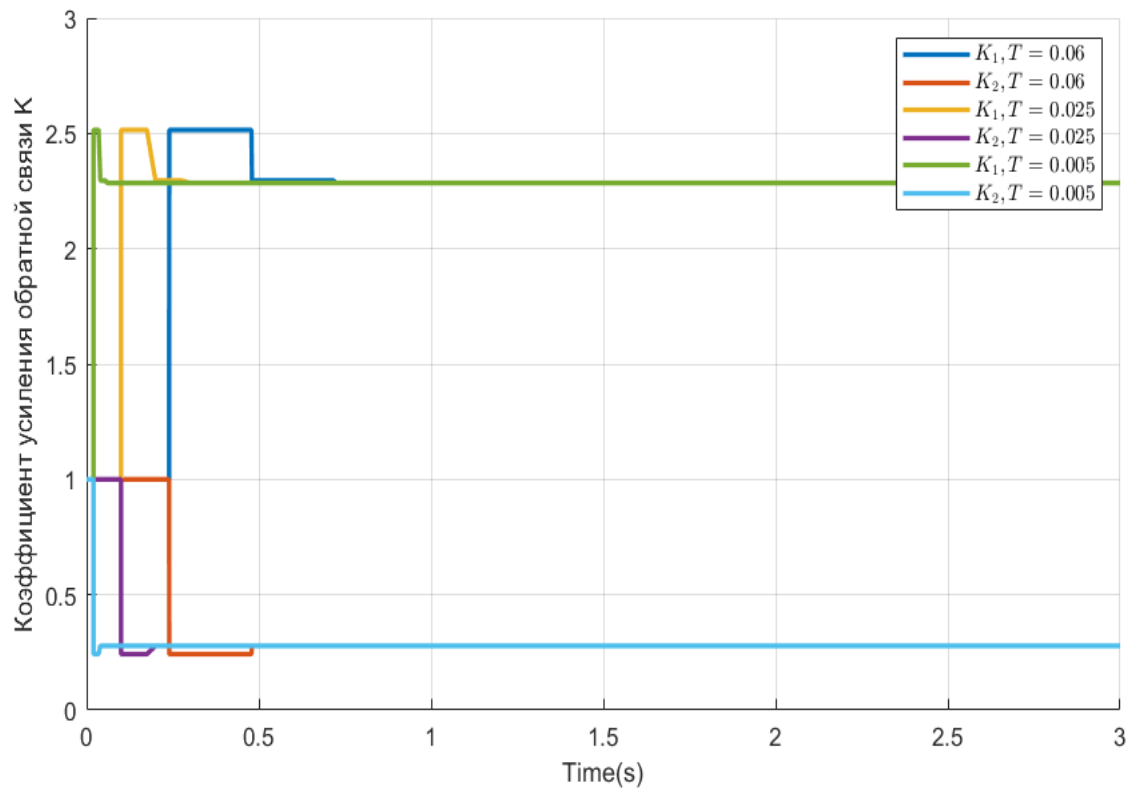


Рисунок 4.6 – Графики коэффициента усиления обратной связи K оптимального адаптивного регулятора для разных значений T

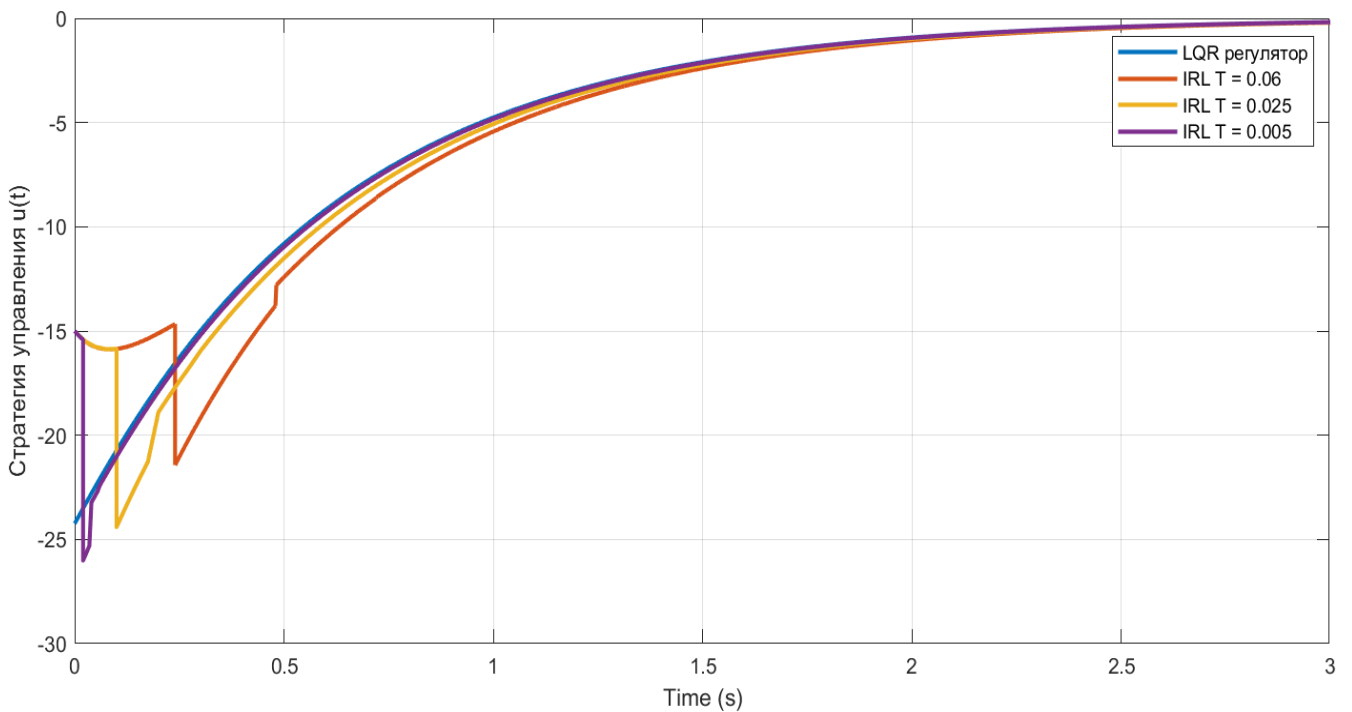


Рисунок 4.7 – График, представляющий управляющий сигнал оптимального регулятора LQR и оптимального адаптивного регулятора с различными значениями T .

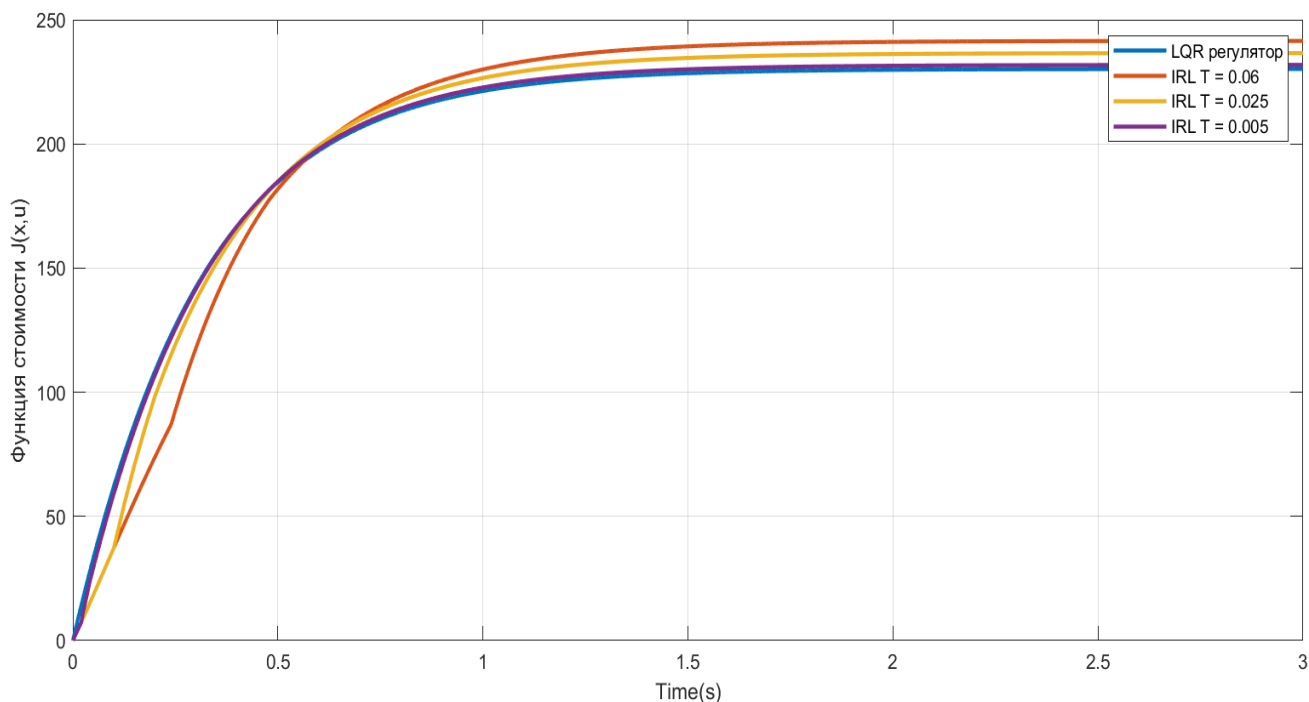


Рисунок 4.8 - График, отражающий значение функции стоимости для оптимального регулятора и оптимального адаптивного регулятора с различными значениями T .

Поскольку алгоритм использует данные, полученные из траектории системы, система не будет устойчивой, когда состояние системы $x(t)$ равно примерно 0, и в этот момент алгоритм будет работать неправильно, что приведет к тому, что управление перестанет быть оптимальным. Поэтому мы должны выключить алгоритм, когда стратегия управления приблизится к оптимальному значению. В этом случае нам необходимо отключить алгоритм, когда коэффициент матрицы обратной связи K приблизится к оптимальному значению. Еще один случай, на который следует обратить внимание, заключается в том, что мы применяем управляющий сигнал u_0 , чтобы сделать систему устойчивой. Поэтому может случиться так, что система сойдётся до того, как будет найдено оптимальное значение управления. Поэтому нам нужно выбрать подходящее время выборки T .

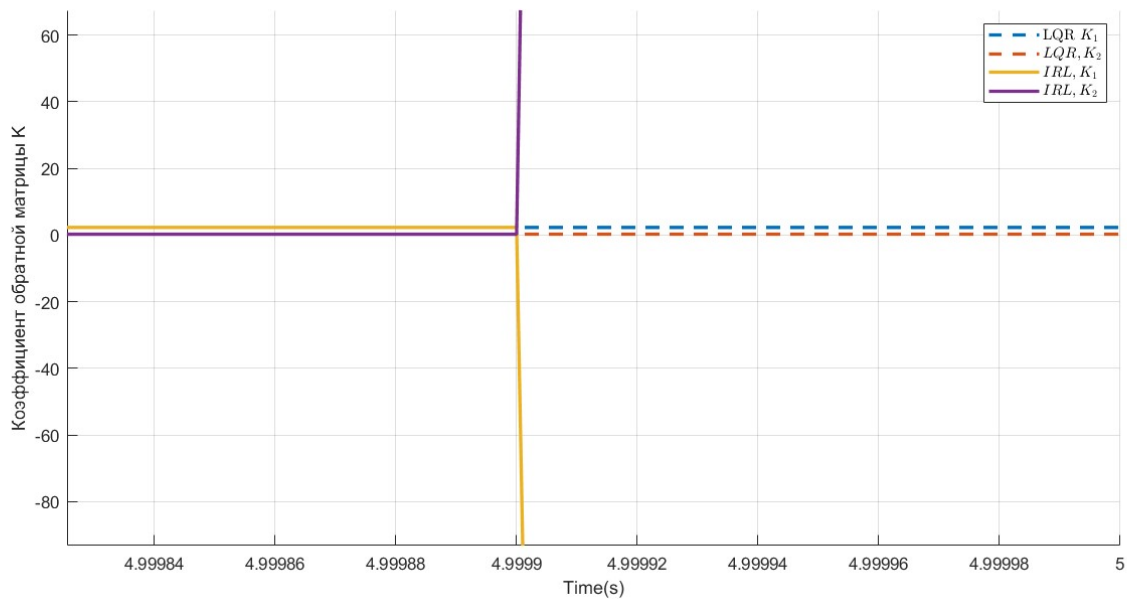


Рисунок 4.9 - Графики коэффициента усиления матрицы обратной связи K оптимального регулятора и оптимального адаптивного регулятора при T=0.25

На графике выше мы видим, что значение коэффициента матрицы обратной связи K приблизилось к оптимальному значению, как и в случае использования LQR. Но когда состояние системы сходится к нулю, алгоритм не отключается, в результате чего значение коэффициента матрицы обратной связи K быстро расходится.

При применении обучения с интегральным подкреплением и приближенного динамического программирования для линейных систем мы часто используем квадратичную функцию стоимости. Пусть n - размер состояния системы, $x \in \mathbb{R}^n$. Тогда объем данных, который необходимо собрать для вычисления новой стратегии управления, равен $nn+12$. Таким образом, для квадратичной системы ($n = 2$) необходимо собрать не менее 3 точек данных, а для квадратичной системы ($n = 4$) необходимо собрать не менее 10 точек данных, чтобы аппроксимировать параметры матрицы P . То есть, собрав необходимое количество точек данных, алгоритм будет аппроксимировать параметры матрицы P . Алгоритм продолжается до тех пор, пока значение P не сойдется с оптимальным значением.

4.3 Регулятор для нелинейных систем

Рассмотрим нелинейных систем

$$\dot{x} = -x_1 + x_2 - 0.5x_1 - 0.5x_2(1 - \cos 2x_1 + 22) + 0 \cos 2x_1 + 2 \text{ и} \quad (4.3)$$

Функция ценности:

$$J = \int_0^{\infty} (x^T Q x + u^T R u) dt \quad (4.4)$$

$$Q = 1001, R = 1$$

Решить уравнение НЖВ мы получили функцию Беллмана и сигнал оптимального управления

$$V^*(x) = w_1 x_1^2 + 2 w_{12} x_1 x_2 + w_2 x_2^2 = 0.5 x_1^2 + x_2^2 \quad (4.5)$$

$$u^*(x) = -\cos 2x_1 + 2x_2 \quad (4.6)$$

Функция вектора активации выбрана

$$\varphi(x) = [x_1^2, x_1 x_2, x_2^2]^T \quad (4.7)$$

Инициализировать начальный вес $W = -205T$ и $x_0 = 12T$

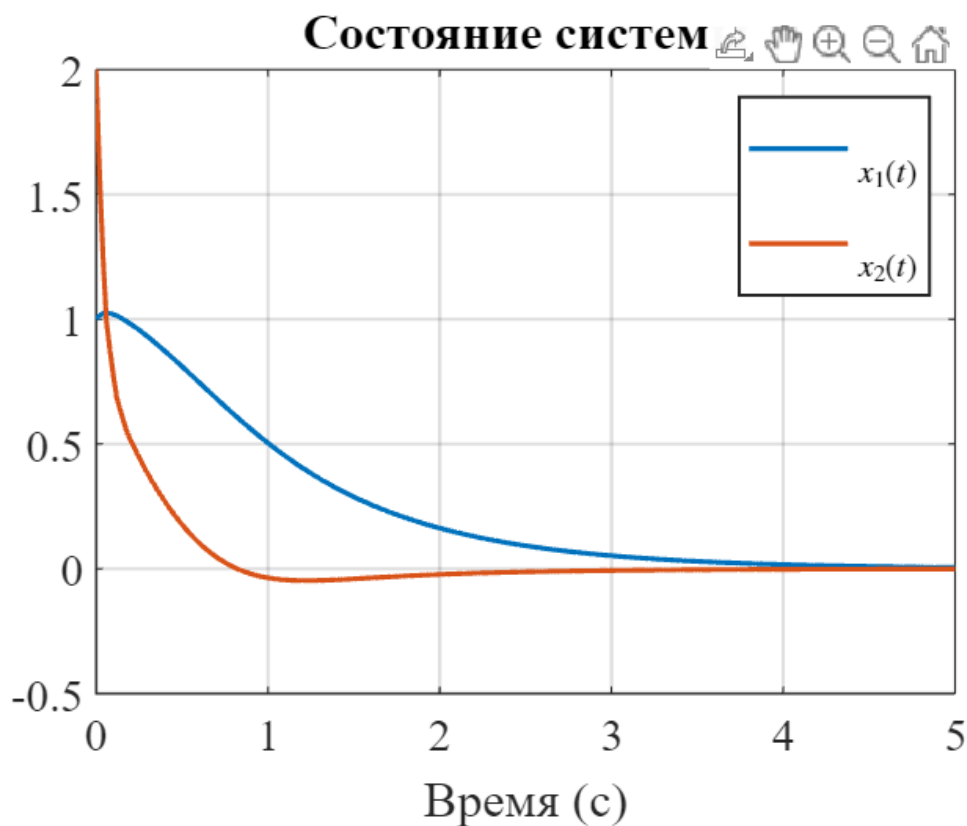


Рисунок 4.10 – График состояния системы с использованием онлайн обучения с интегральным подкреплением

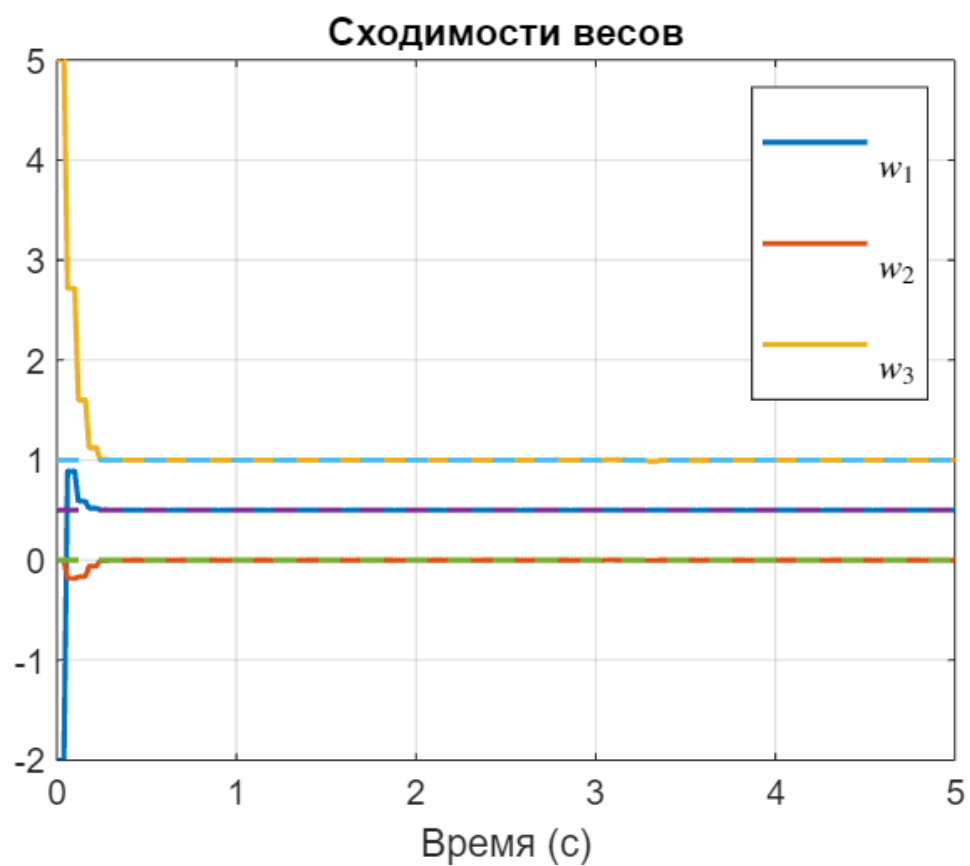


Рисунок 4.11 Сходимость весов W

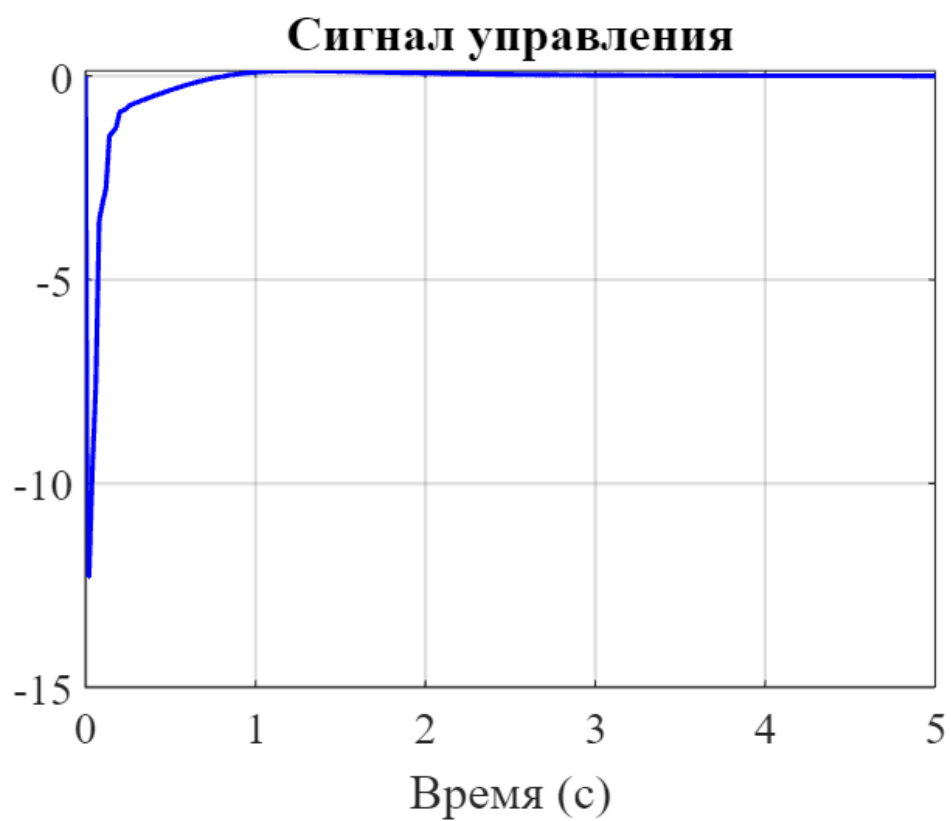


Рисунок 4.12 – Сигнал управления системой с алгоритмом обучения с подкреплением

Как видно из графика, вес W сходится точно к своему оптимальному значению через 0.24 секунды, в то время как сигнал управления по-прежнему помогает системе стабилизироваться на достаточно хорошей скорости. Тем самым показывая корректность алгоритма.

4.4 Регулятор линейной системы для перевернутого маятника и тележки

Данная система

$$\dot{x} = Ax + Bu,$$

$$A = \begin{bmatrix} 0 & 1 & 0 & 0 \\ -\frac{g}{L} & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix}, B = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

Где $x = [\theta \ \dot{\theta} \ x \ \dot{x}]^T$ – состояние системы, $u \in \mathbb{R}^1$ – сигнал управления системы

Параметры для моделирования из [18]

$$M = 2.4 \text{ кг}, m = 0.23 \text{ кг}, g = 9.81 \text{ м/с}^2, L = 0.46 \text{ м}$$

Задача оптимального управления состоит в оптимизации следующей функции

$$J_{x,u} = \int_0^\infty (x^T Q x + u^T R u) dt$$

$$Q = \begin{bmatrix} 10 & 0 & 0 & 0 \\ 0 & 1000 & 0 & 0 \\ 0 & 0 & 10 & 0 \\ 0 & 0 & 0 & 10 \end{bmatrix}, R = 1$$

Матрица оптимального регулятора (LQR)

$$K = R^{-1} B^T P = [-62.7 \ -13.1 \ -1.0 \ -2.92]$$

Результаты моделирования

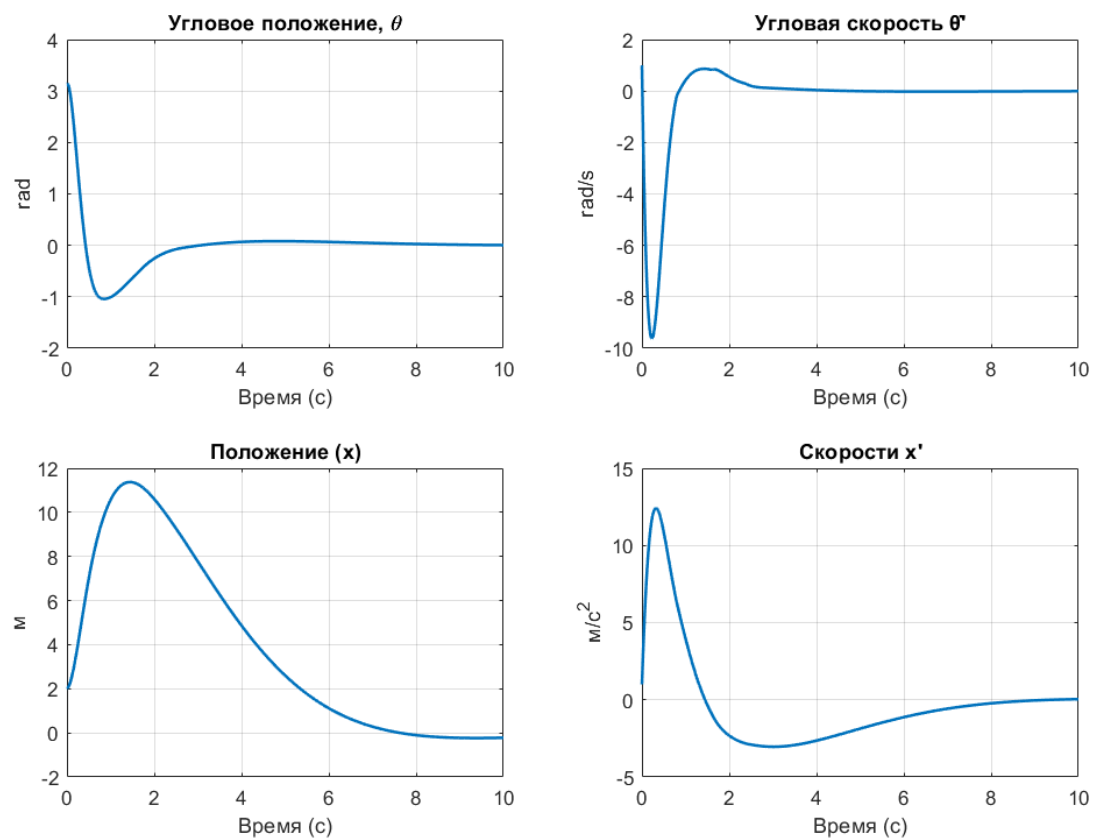


Рисунок 4.13 – Графика состояния системы $x(t)$

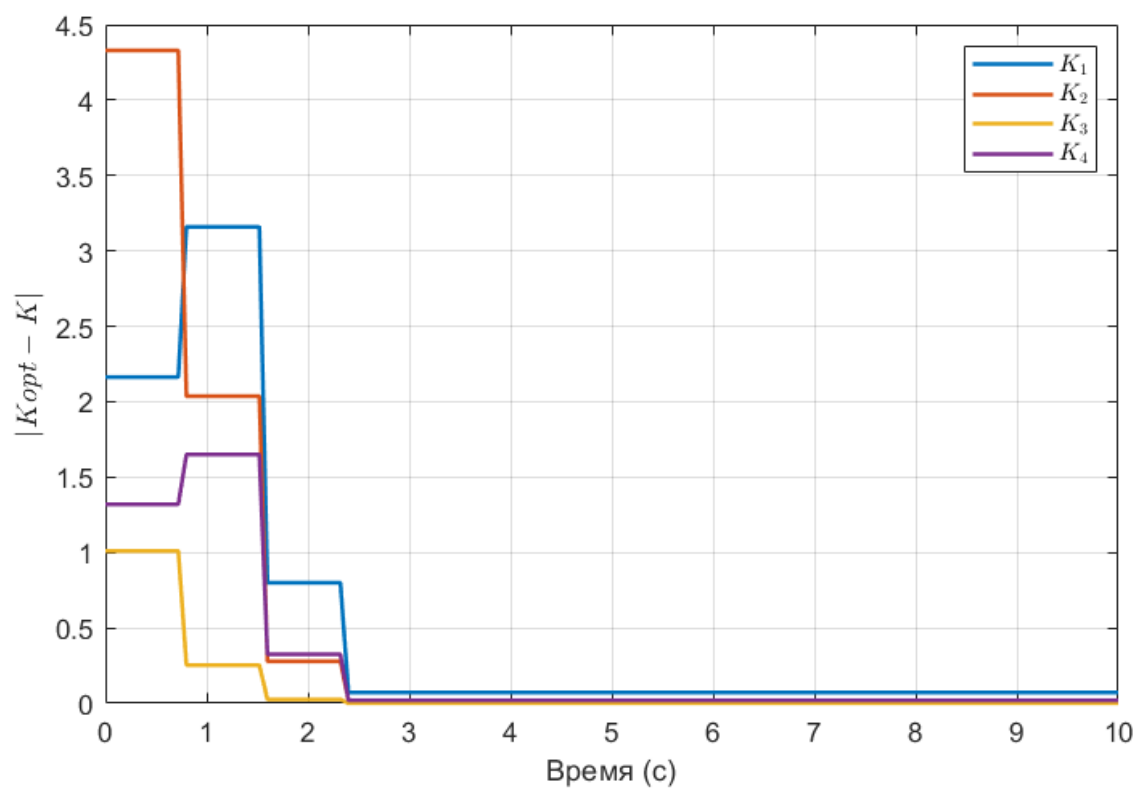


Рисунок 4.14 – График отклонения коэффициентов усиления управления K от их оптимальных значений K_{opt}

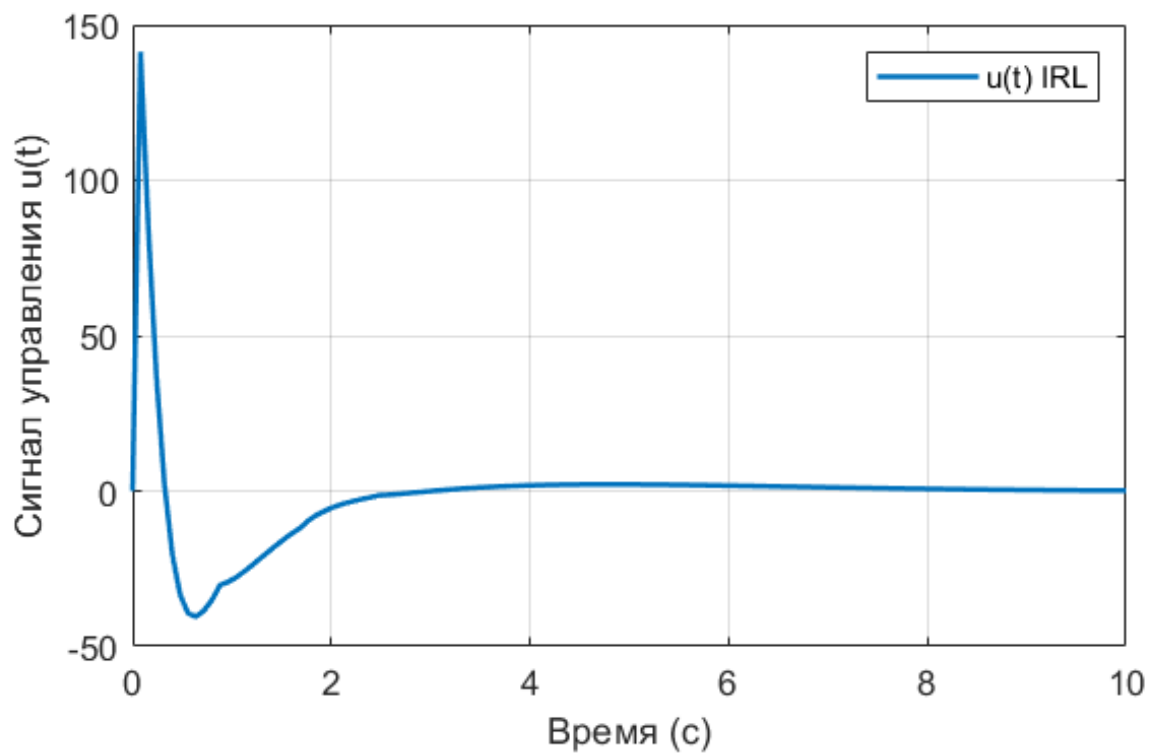


Рисунок 4.15 – График сигнал управления системы $u(t)$

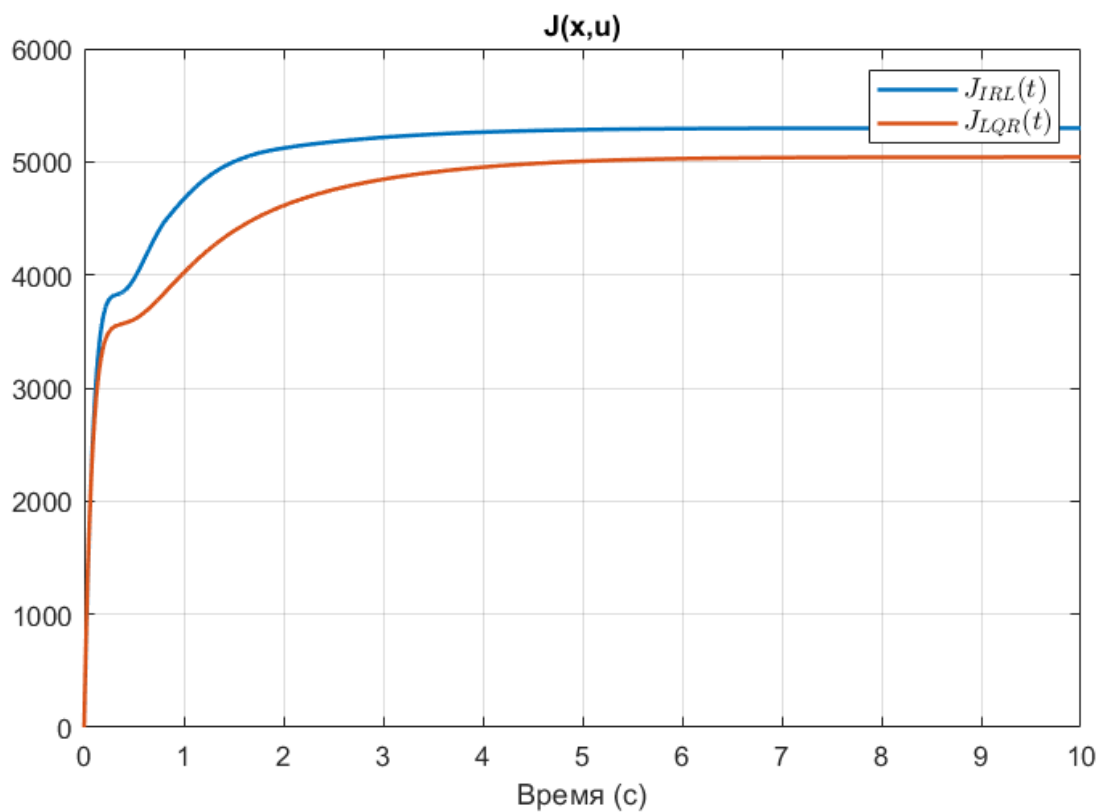


Рисунок 4.16 – График, отражающий значение функции стоимости оптимального регулятора и оптимального адаптивного регулятора

4.5 Моделирование для маятника

Рассмотрим систему маятника

$$\ddot{x} = -g \sin \theta - b \dot{\theta} + \ddot{\theta} u \quad (4.8)$$

Где $x = \theta$ - состояние системы, u – сигнал управление

Функция ценности:

$$J = \int_0^{\infty} x^T Q x + u^T R u dt \quad (4.4)$$

$$Q = 1001, R = 1$$

Для моделирования, функции вектора активации выбирается следующим образом

$$\varphi(x) = [x_1^2, x_1 x_2, x_2^2]^T$$

Инициализировать начальный вес $W = 0.20T$ и $x_0 = \pi/60T$

Результаты моделирования

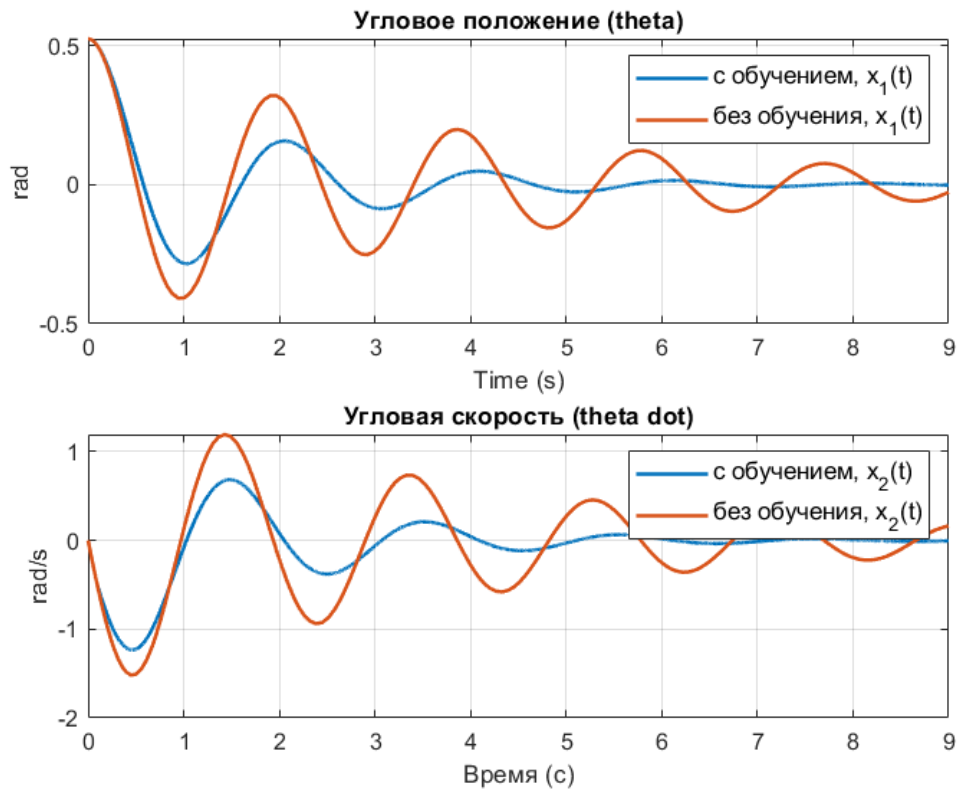


Рисунок 4.17 – Состояние с обучением и без обучения

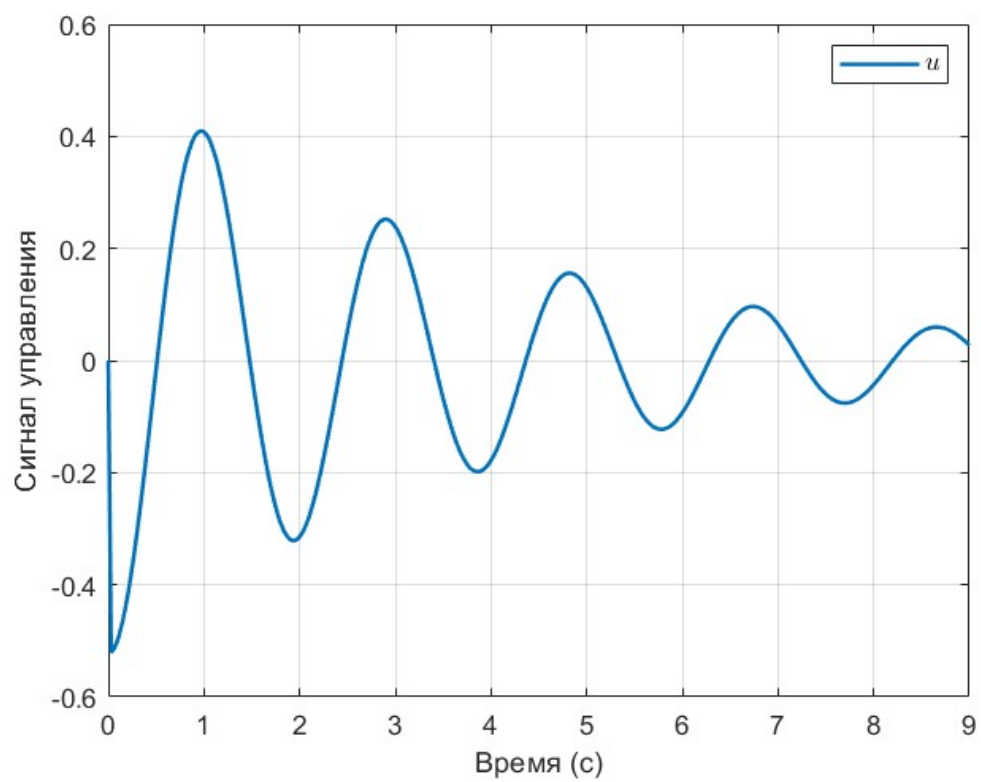


Рисунок 4.18 – Сигнал управления при без обучения

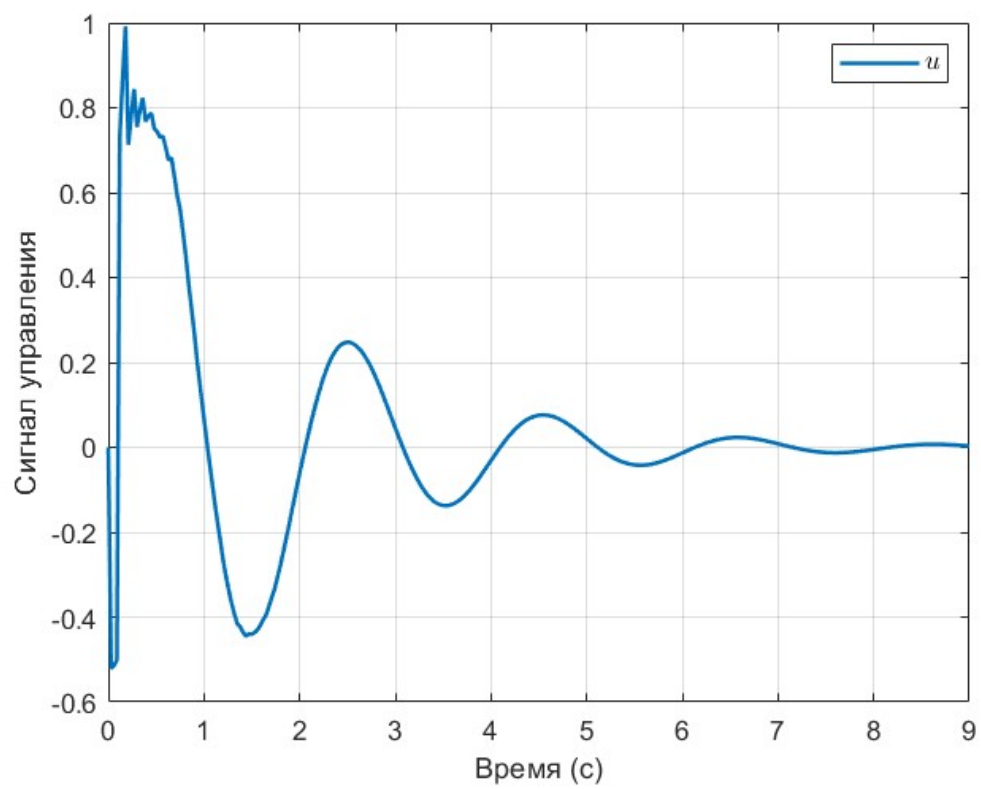


Рисунок 4.19 – Сигнал управления u с обучением

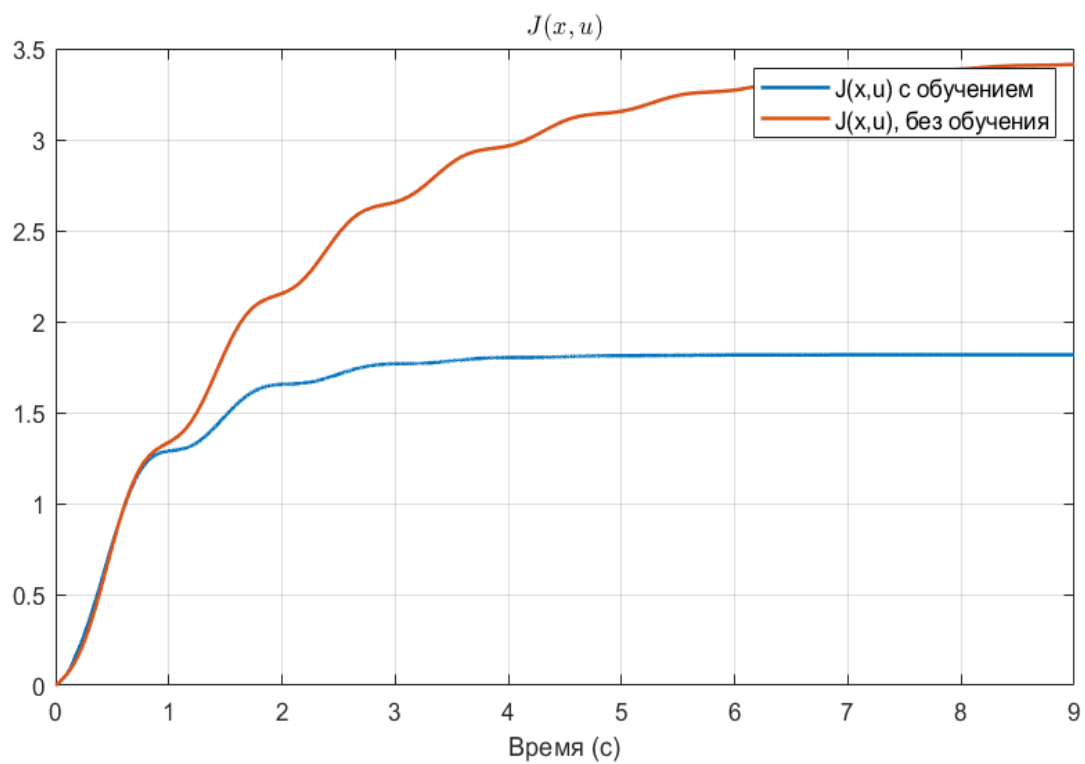


Рисунок 4.20 – Значение функции стоимости $J(x, u)$ с обучением и без обучения

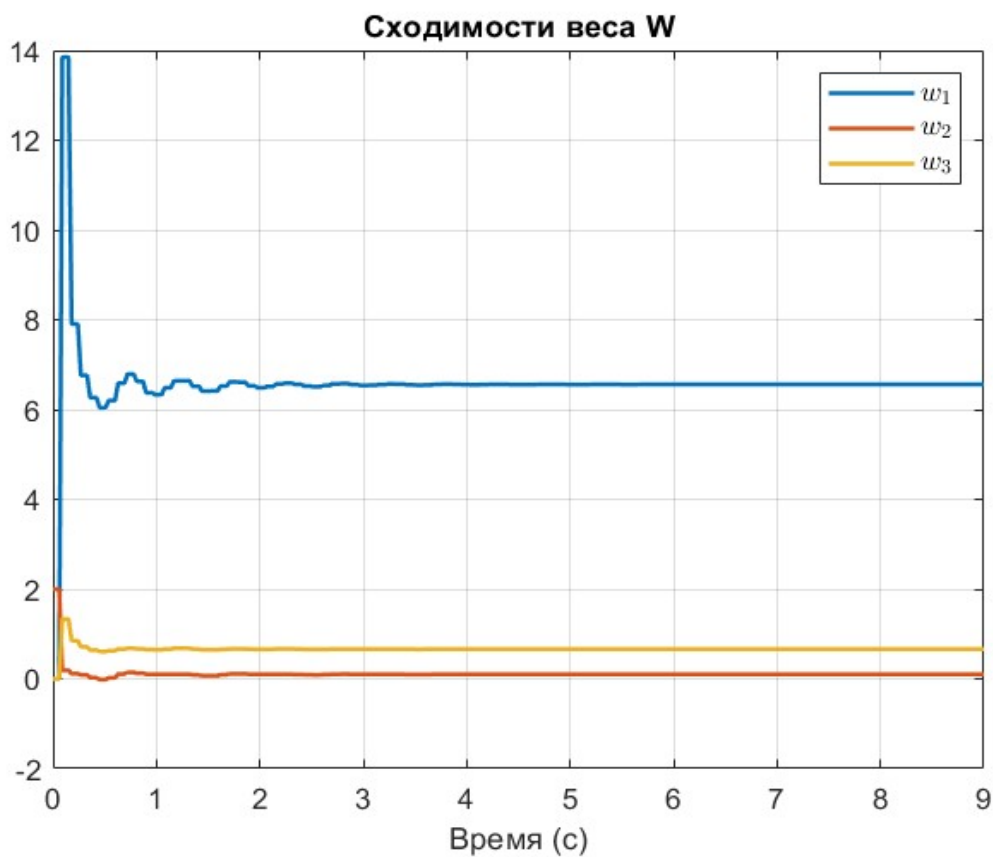


Рисунок 4.21 – Сходимости весов при применении обучения с интегральным подкреплением

Мы видим, что состояние системы сходилось лучше в случае с обучением. Значение функции стоимости также минимизировано, видно, что оно меньше в случае без обучения. Таким образом, при использовании обучения с интегральным подкреплением, регулятор улучшает исходный управляющий сигнал u_0 и минимизирует функцию стоимости $J(x,u)$. Для этой системы сложно решить задачу НЖВ для сравнения с результатами системы. Но в целом алгоритм IRL сошелся к стратегии управления, которая лучше, чем изначально инициализированный сигнал управления. Пример в разделе 4.3 более наглядно показывает сходимость алгоритма к оптимальному значению управления для нелинейной системы.

5. Заключение

В данной работе были изучены основные концепции и идеи, лежащие в основе алгоритма обучения с подкреплением. Были изучены два популярных алгоритма в обучении с подкреплением — это алгоритмы итерации по стратегии и по критерию, предназначенные для решения проблемы поиска оптимальных стратегий.

Представлены и исследованы математические регуляторы, таким как адаптивный регулятор, оптимальный регулятор (LQR) и оптимальный адаптивный регулятор, основанные на обучении с интегральным подкреплением.

Моделирование выполняется в MATLAB. Для линейных систем был реализован анализ и исследование эффективности оптимального адаптивного регулятора на основе обучения с интегральным подкреплением. Проведено сравнение эффективности оптимального адаптивного регулятора с эффективностью оптимального регулятора и адаптивного регулятора. В результате оптимальный адаптивный регулятор имеет эффективность, примерно равную эффективности оптимального регулятора. При этом адаптивный регулятор не сходится к оптимальному управлению. Для нелинейных систем корректность оптимального адаптивного регулятора доказана, когда веса

функции аппроксимируют сходящееся значение к решению уравнения Гамильтона Якоби Беллмана.

Оптимальный адаптивный регулятор, основанный на обучении с интегральным подкреплением, строится на основе данных, собранных по траектории движения системы. Поэтому проанализировано и изучено влияние времени сборки данных T на сходимость оптимального управления оптимального адаптивного регулятора. В результате, чем меньше время T , тем лучше система приближается к оптимальному управлению. Выбор подходящего времени сборки данных был представлен. Если регулятор сходится к оптимальному значению, мы должны прекратить обновление стратегии управления системой, подробности представлены в разделе 4.2.

Таким образом, был успешно исследован оптимальный адаптивный регулятор на основе интегрального обучения с подкреплением. Преимущество этого регулятора в том, что мы получаем оптимальный управляющий сигнал без необходимости полного понимания динамики системы (без знания матрицы A для линейных систем и без знания $f(x)$ для нелинейных систем).

Список использованных источников

1. Frank Lewis and Draguna Vrabie // *Reinforcement learning and adaptive dynamic programming for feedback control. Circuits and Systems Magazine, IEEE*, 9:32 – 50, 01 2009.
2. Ivo Grondman, Lucian Bus, oniu, Gabriel A. D. Lopes, and Robert Babuska // *A Survey of Actor-Critic Reinforcement Learning: Standard and Natural Policy Gradients*
3. Ivo Grondman, Lucian Bus, oniu, Gabriel A.D. Lopes and Robert Babuska // *A Survey of Actor-Critic Reinforcement Learning: Standard and Natural Policy Gradients*
4. V. R. Konda and J. N. Tsitsiklis // “On Actor-Critic Algorithm”, *SIAM Journal on Control and Optimization*, vol. 42, No. 4, pp. 1143–1166, 2003.
5. Richard S. Sutton and Andrew G. Barto // *Reinforcement Learning: An Introduction, Second Edition*, MIT Press, Cambridge, MA, 2018
6. Fadi AlMahamid, and Katarina Grolinger, Department of Electrical and Computer Engineering Western University // *Reinforcement Learning Algorithms: An Overview and Classification*

7. A. Geramifard, T. J. Walsh, S. Tellex, G. Chowdhary, N. Roy, and J. P. How // *A Tutorial on Linear Function Approximators for Dynamic Programming and Reinforcement Learning*, Vol. 6, No. 4 (2013) 375–454, 2013
8. Beakcheol Jang, Myeonghwi Kim, Gaspard Harerimana, and Jong Wook Kim // *Q-Learning Algorithms: A Comprehensive Classification and Applications*, Department of Computer Science, Sangmyung University, Seoul 03016, South Korea
9. Draguna Vrabie, Kyriakos G. Vamvoudakis and Frank L. Lewis // *Optimal Adaptive Control and Differential Games by Reinforcement Learning Principles*
10. Frank L. Lewis, Draguna L. Vrabie, Vassilis L. Syrmos // *Optimal control-3rd ed*, 2012.
11. Bahare Kiumarsi, Member, IEEE, Kyriakos G. Vamvoudakis, Senior Member, IEEE, Hamidreza Modares, Member, IEEE, and Frank L. Lewis, Fellow, IEEE, *Optimal and Autonomous Control Using Reinforcement Learning: A Survey*, 2018.
12. Kyriakos G. Vamvoudakis, Draguna Vrabie, Frank L. Lewis, *Online adaptive learning of optimal control solutions using integral reinforcement learning*, 2011.
13. D. Vrabie, O. Pastravanu, M. Abu-Khalaf, and F. L. Lewis, “Adaptive optimal control for continuous-time linear systems based on policy iteration,” *Automatica*, vol. 45, no. 2, pp. 477–484, 2009.
14. D. Vrabie and F. L. Lewis, “Neural network approach to continuous-time direct adaptive optimal control for partially-unknown nonlinear systems,” *Neural Netw.*, vol. 22, no. 3, pp. 237–246, Apr. 2009.
15. Vamvoudakis KG, Lewis FL. *Online actor-critic algorithm to solve the continuous-time infinite horizon optimal control problem. Automatica* 2010.
16. Kleinman D. *On an iterative technique for Riccati equation computations. IEEE Transactions on Automatic Control* 1968;
17. Uğur Yıldırım, Yildiz Technical University, *Adaptive Control of an Inverted Pendulum by a Reinforcement Learning based LQR Method*.

18. Vinayak Kumar, Ruchi Agarwal, Modeling and Control of Inverted Pendulum cart system using ¹ PID-LQR based Modern Controller, July 01-03, 2022, Prayagraj, India.