

DATA DRIVEN DYNAMIC MODELING AND REINFORCEMENT LEARNING BASED CONTROL OF BATCH PROCESSES

TANUJA JOSHI



DEPARTMENT OF CHEMICAL ENGINEERING
INDIAN INSTITUTE OF TECHNOLOGY DELHI
DECEMBER 2022

© Indian Institute of Technology Delhi (IITD), New Delhi, 2022

DATA DRIVEN DYNAMIC MODELING AND
REINFORCEMENT LEARNING BASED CONTROL
OF BATCH PROCESSES

by

TANUJA JOSHI

Department of Chemical Engineering

Submitted

in fulfilment of the requirements of the degree of Doctor of Philosophy

to the



INDIAN INSTITUTE OF TECHNOLOGY DELHI

DECEMBER 2022

THESIS CERTIFICATE

This is to certify that the thesis titled **Data Driven Dynamic Modeling and Reinforcement Learning based Control of Batch Processes**, submitted by **Tanuja Joshi (2018CHZ8100)**, to the Indian Institute of Technology Delhi, for the award of the degree of **Doctor of Philosophy**, is a bonafide record of the research work done by her under my supervision. The contents of this thesis, in full or in parts, have not been submitted to any other Institute or University for the award of any degree or diploma.

Prof. Hariprasad Kodamana
Research Supervisor
Associate Professor
Dept. of Chemical Engineering &
Yardi School of Artificial Intelli-
gence
IIT-Delhi, 110016

Date: 7 July 2022

ACKNOWLEDGEMENTS

Foremost, I would like to express my sincere gratitude to my advisor Prof. Hariprasad Kodamana for the continuous support during my Ph.D study, for his patience, motivation, enthusiasm, and immense knowledge. His guidance, timely advise and scientific approach is a source of inspiration and helped me in all the time of research and writing of this thesis. His dedication and keen interest and above all his overwhelming attitude to always help his students had been mainly responsible for completing this thesis.

Besides my advisor, I would like to thank my research committee members Prof. Manojkumar C. Ramteke, Prof. Anupam Shukla, and Prof. Subashish Datta for their encouragement and insightful comments. I especially want to thank Prof. Harikumar Kandath, IIIT Hyderabad, for his valuable inputs, vast knowledge and motivating guidance in this project. Many thanks to Prof. Niket Kaisare, IIT Madras to participate in my research by providing invaluable feedback and make this work possible.

In addition, I owe my thanks to my current fellow labmates, Avan, Reena, Nikita, Jyoti, Umang, Ahtesham, Nasre, Deepak and all of my past labmates, Devansh, Abhyansh, Shikhar, Vishesh, Vikas, Mayuna in the Control & Analytics of Process System (CAPS) Lab, IITD for the stimulating discussions and the positive environment they have created in the lab. Furthermore, I would like to express my sincere thanks to all the office staff members of the the Chemical Engineering Department, IITD.

Last but not the least I would like to thank the almighty god, my parents, family members, and my friends for their unconditional love, prayers and constant support throughout my life.

ABSTRACT

Development of control relevant models with excellent predictive capability is a challenging task for a batch process due to the non-linear dynamics and varying operating conditions within the batch and batch-to-batch. First principles-based models do not provide a remedy to this problem in most cases as they are tedious to develop and strenuous to solve for complex processes. Incidentally, industrial processes are warehouses of a large amount of multivariate data. Hence, data-driven modelling approaches can be conveniently used in this juncture. Traditional approaches of the data-based modeling methods focus on global approaches, however, developing a single global model for the entire batch process may not be a reasonable option because operating conditions change dynamically in industrial processes and may warrant updating the model online. The first part of the thesis focuses on the development of a novel data-driven dynamic model based on the local modeling approach i.e., Just-in-Time Learning (JITL) framework for batch process modeling. The accuracy of the underlying model under the JITL framework is based on the ‘similarity measure’ used to extract the relevant data from the massive historical database similar to the query point, and the ‘weighting strategies’ adopted. The proposed formulation incorporates a ‘query profile’ instead of a ‘query point’, with a view to bring in the dynamic modeling capabilities. A new ‘searching strategy’ based on ‘profile similarity’ is proposed for taking into account the time-varying dynamics of the batch processes. Further, a new ‘weighting strategy’ is introduced to ensure that the complete dynamics of the batch processes are captured and also to accommodate the outliers in the historical database.

The second part of this thesis focuses on the data-driven control of batch processes.

Control of batch processes is difficult due to their complex nonlinear dynamics and unsteady-state operating conditions within batch and batch-to-batch. Advanced control strategies like Model Predictive Control (MPC) used in the process industries for batch process control employs mathematical programming to solve a constrained, possibly non-convex, optimisation problem. The key issue here is that the online computational burden at each time step to obtain the optimal control input profile is very high. This limits the implementation of these approaches for complex nonlinear, high-dimensional dynamical systems, despite the advancement in computational hardware and numerical methods. Further, the performance of a model-based controller is highly dependent on the availability of an accurate process model. For a complex and nonlinear process, the availability of such a model is a limitation as it requires significant prior knowledge and expertise. Plant-model mismatch occurs even if there is a slight variation in the real process and its approximate model, leading to inaccurate predictions of the performance variable, such as product yield. Even if the process model is available, the controller performance deteriorates in the presence of uncertainties and process drifts. Batch-to-batch variations that occur due to raw material fluctuations, cleaning, etc., are a source of uncertainty that further deteriorate the closed-loop performance. In this juncture, the development of a control strategy that does not entirely rely on the knowledge of the accurate dynamics of the system and can handle the stochastic dynamics and plant-model mismatches is extremely useful. Model-free Reinforcement Learning (RL), where the agent (analogous to the controller) learns the optimal control action (analogous to control input) by directly interacting with the operational environment (analogous to the process), offers a potential alternative to traditional model-based approaches for process control. RL frameworks with actor-critic architecture have recently become popular for the control process systems where both the state and action spaces are continuous. Subsequent works focus on the development of Actor-Critic RL based controller by developing two novel Actor-Critic RL algorithms, namely, (i) Twin Actor Twin Delayed Deterministic Policy Gradient (TATD3), a deterministic RL algorithm, and (ii) Twin Actor Soft Actor-Critic (TASAC), a stochastic RL algorithm, for batch process control by incorporating an ‘en-

semble of actors' in the Actor-Critic algorithms for training the policies to achieve an overall optimal policy for batch process control.

The efficacies of the developed approaches in this work are evaluated by simulation studies on batch process case studies. In a nutshell, the overall objective is to develop a data-driven dynamic model and a model-free RL-based controller for complex, nonlinear batch processes to ensure optimal operation of the batch processes.

KEYWORDS: data-driven model; just-in-time learning; reinforcement learning; deep-Q-learning; deep deterministic policy gradient; actor-critic algorithms; batch process control

सार

उत्कृष्ट भविष्यवाणी क्षमता वाले नियंत्रण प्रासंगिक मॉडलों का विकास, गैर-रैखिक गतिशीलता और बैच और बैच-टू-बैच के भीतर बदलती परिचालन स्थितियों के कारण बैच प्रक्रिया के लिए एक चुनौतीपूर्ण कार्य है। पहले सिद्धांत-आधारित मॉडल ज्यादातर मामलों में इस समस्या का समाधान नहीं प्रदान करते हैं क्योंकि वे जटिल प्रक्रियाओं के लिए विकसित करने और हल करने के लिए कठिन हैं। संयोग से, औद्योगिक प्रक्रियाएं बड़ी मात्रा में बहुभिन्नरूपी डेटा के गोदाम हैं। इसलिए, इस मोड़ पर डेटा-संचालित मॉडलिंग दृष्टिकोणों का आसानी से उपयोग किया जा सकता है। डेटा-आधारित मॉडलिंग विधियों के पारंपरिक दृष्टिकोण वैश्विक दृष्टिकोण पर ध्यान केंद्रित करते हैं, हालांकि, संपूर्ण बैच प्रक्रिया के लिए एकल वैश्विक मॉडल विकसित करना एक उचित विकल्प नहीं हो सकता है क्योंकि परिचालन की स्थिति औद्योगिक प्रक्रियाओं में गतिशील रूप से बदलती है और मॉडल को ऑनलाइन अपडेट करने की आवश्यकता हो सकती है। थीसिस का पहला भाग बैच प्रोसेस मॉडलिंग के लिए स्थानीय मॉडलिंग दृष्टिकोण यानी जस्ट-इन-टाइम लर्निंग (JITL) फ्रेमवर्क के आधार पर एक उपन्यास डेटा-संचालित गतिशील मॉडल के विकास पर केंद्रित है। जेआईटीएल ढांचे के तहत अंतर्निहित मॉडल की सटीकता 'समानता माप' पर आधारित है जिसका उपयोग क्वेरी बिंदु के समान विशाल ऐतिहासिक डेटाबेस से प्रासंगिक डेटा निकालने के लिए किया जाता है, और अपनाई गई 'भार रणनीति'। गतिशील मॉडलिंग क्षमताओं को लाने की दृष्टि से प्रस्तावित सूत्रीकरण में 'क्वेरी पॉइंट' के बजाय 'क्वेरी प्रोफाइल' शामिल है। बैच प्रक्रियाओं की समय-भिन्न गतिशीलता को ध्यान में रखते हुए 'प्रोफाइल समानता' पर आधारित एक नई 'खोज रणनीति' प्रस्तावित है। इसके अलावा, यह सुनिश्चित करने के लिए एक नई 'वेटिंग स्ट्रैटेजी' पेश की गई है कि बैच प्रक्रियाओं की पूरी गतिशीलता पर कब्जा कर लिया गया है और ऐतिहासिक डेटाबेस में आउटलेयर को भी समायोजित किया गया है।

इस थीसिस का दूसरा भाग बैच प्रक्रियाओं के डेटा-संचालित नियंत्रण पर केंद्रित है। बैच और बैच-टू-बैच के भीतर उनकी जटिल अरैखिक गतिशीलता और अस्थिर-स्थिति परिचालन

स्थितियों के कारण बैच प्रक्रियाओं का नियंत्रण मुश्किल है। बैच प्रक्रिया नियंत्रण के लिए प्रक्रिया उद्योगों में उपयोग किए जाने वाले मॉडल प्रिडिक्टिव कंट्रोल (एमपीसी) जैसी उन्नत नियंत्रण रणनीतियाँ गणितीय प्रोग्रामिंग को एक विवश, संभवतः गैर-उत्तल, अनुकूलन समस्या को हल करने के लिए नियोजित करती हैं। यहां मुख्य मुद्दा यह है कि इष्टतम नियंत्रण इनपुट प्रोफ़ाइल प्राप्त करने के लिए प्रत्येक चरण में ऑनलाइन कम्प्यूटेशनल बोझ बहुत अधिक है। कम्प्यूटेशनल हार्डवेयर और संख्यात्मक तरीकों में प्रगति के बावजूद, यह जटिल गैर-रैखिक, उच्च-आयामी गतिशील प्रणालियों के लिए इन दृष्टिकोणों के कार्यान्वयन को सीमित करता है। इसके अलावा, मॉडल-आधारित नियंत्रक का प्रदर्शन सटीक प्रक्रिया मॉडल की उपलब्धता पर अत्यधिक निर्भर है। एक जटिल और अरेखीय प्रक्रिया के लिए, ऐसे मॉडल की उपलब्धता एक सीमा है क्योंकि इसके लिए महत्वपूर्ण पूर्व ज्ञान और विशेषज्ञता की आवश्यकता होती है। प्लांट-मॉडल बेमेल तब भी होता है जब वास्तविक प्रक्रिया और उसके अनुमानित मॉडल में थोड़ी भिन्नता होती है, जिससे उत्पाद की उपज जैसे प्रदर्शन चर की गलत भविष्यवाणी होती है। भले ही प्रक्रिया मॉडल उपलब्ध हो, अनिश्चितताओं और प्रक्रिया के बहाव की उपस्थिति में नियंत्रक का प्रदर्शन बिगड़ जाता है। कच्चे माल में उतार-चढ़ाव, सफाई आदि के कारण होने वाली बैच-टू-बैच विविधताएं अनिश्चितता का एक स्रोत हैं जो बंद-लूप प्रदर्शन को और खराब करती हैं। इस समय, एक नियंत्रण रणनीति का विकास जो पूरी तरह से सिस्टम की सटीक गतिशीलता के ज्ञान पर निर्भर नहीं करता है और स्टोकेस्टिक गतिशीलता और प्लांट-मॉडल बेमेल को संभाल सकता है, अत्यंत उपयोगी है। मॉडल-मुक्त सुदृढीकरण सीखना (आरएल), जहां एजेंट (नियंत्रक के अनुरूप) परिचालन वातावरण (प्रक्रिया के अनुरूप) के साथ सीधे बातचीत करके इष्टतम नियंत्रण क्रिया (इनपुट को नियंत्रित करने के अनुरूप) सीखता है, पारंपरिक मॉडल के लिए एक संभावित विकल्प प्रदान करता है। प्रक्रिया नियंत्रण के लिए आधारित दृष्टिकोण। अभिनेता-आलोचक वास्तुकला के साथ आरएल ढांचे हाल ही में नियंत्रण प्रक्रिया प्रणालियों के लिए लोकप्रिय हो गए हैं जहां राज्य और क्रिया स्थान दोनों निरंतर हैं। इसके बाद के कार्यों में दो उपन्यास अभिनेता-आलोचक आरएल एल्गोरिदम विकसित करके अभिनेता-आलोचक आरएल आधारित नियंत्रक के विकास पर ध्यान केंद्रित किया गया है, अर्थात्, (i) द्विन एक्टर द्विन डिलेड डिटर्मिनिस्टिक पॉलिसी ग्रेडिएंट (टीएटीडी3), एक नियतात्मक आरएल एल्गोरिथ्म और (ii) द्विन एक्टर सॉफ्ट द्विन एक्टर सॉफ्ट एक्टर-क्रिटिक (TASAC), एक स्टोकेस्टिक आरएल एल्गोरिथ्म, बैच प्रक्रिया नियंत्रण के लिए एक समग्र इष्टतम नीति प्राप्त करने के लिए नीतियों को प्रशिक्षित

करने के लिए अभिनेता-आलोचक एल्गोरिदम में अभिनेताओं के समूह को शामिल करके बैच प्रक्रिया नियंत्रण के लिए।

इस कार्य में विकसित दृष्टिकोणों की प्रभावशीलता का मूल्यांकन बैच प्रक्रिया मामले के अध्ययन पर सिमुलेशन अध्ययन द्वारा किया जाता है। संक्षेप में, समग्र उद्देश्य बैच प्रक्रियाओं के इष्टतम संचालन को सुनिश्चित करने के लिए जटिल, गैर-रैखिक बैच प्रक्रियाओं के लिए एक डेटा-संचालित गतिशील मॉडल और एक मॉडल-मुक्त आरएल-आधारित नियंत्रक विकसित करना है।

कीवर्ड: डेटा-संचालित मॉडल; समय-समय पर सीखना; सुदृढीकरण सीखना; डीप-क्यू-लर्निंग; गहरी नियतात्मक नीति प्रवणता; अभिनेता-आलोचक एल्गोरिदम; बैच प्रक्रिया नियंत्रण

Contents

ACKNOWLEDGEMENTS	i
ABSTRACT	ii
LIST OF TABLES	xii
LIST OF FIGURES	xiv
ABBREVIATIONS	xv
NOTATION	xvi
1 Introduction	1
1.1 Motivation	1
1.2 Research Contributions	5
1.3 Outline of the Thesis	8
2 Literature Review	9
2.1 Just-in-Time Learning	9
2.2 Reinforcement Learning	17
2.2.1 Markov Decision Process	19
2.2.2 Monte Carlo and Temporal Difference learning	20
2.2.3 Value-based Approach	22
2.2.4 Policy-based RL Approach	29
2.2.5 Actor-Critic Methods	32
2.2.6 Deep Deterministic Policy Gradient	35

2.3	Maximum entropy Reinforcement Learning	38
2.4	Conclusion from the literature	38
3	Data-Driven Model of batch processes under the Just-in-Time Learning framework	40
3.1	Formulation of the proposed dynamic JIT framework	42
3.1.1	Identification of dynamic sequence similar to query sequence	42
3.1.2	Attribution of weights	45
3.1.3	Local model development in the JIT framework	51
3.2	Illustrative Examples	53
3.2.1	Example 1: Numerical example	53
3.2.2	Example 2: Batch polymerisation of methyl methacrylate (PMMA)	55
3.2.3	Example 3: Batch transesterification	59
3.3	Discussions	61
3.3.1	Performance of the proposed approach in the presence of outliers	61
3.3.2	Effect of window length on the prediction Performance	62
3.3.3	Effect of batch-to-batch Variation	63
3.3.4	Benefits of space, time and batch weights	65
3.4	Conclusion	66
4	Twin Actor Twin Delayed Deep Deterministic Policy Gradient (TATD3) learning for batch process control	67
4.1	Preliminaries	68
4.2	Main Contribution: TATD3 - the methodology	70
4.2.1	Reward function	72
4.2.2	TATD3 - steps involved	75
4.3	Results and Discussion	80
4.3.1	Case study 1: Batch transesterification process	81
4.3.2	Case study 2: Exothermic batch process	87
4.3.3	Effect of batch-to-batch variation	89

4.4 Conclusion	89
5 TASAC: a twin-actor reinforcement learning framework with stochastic policy for batch process control	96
5.1 Preliminaries	98
5.2 The proposed framework : TASAC	99
5.2.1 Strategies for action selection and TVs	102
5.3 Results and Discussion	104
5.3.1 System Description	106
5.3.2 TASAC training details	108
5.3.3 Reward function	108
5.3.4 Performance comparison under nominal condition	110
5.3.5 Performance in the presence of measurement noise	110
5.3.6 Effect of batch-to-batch variations	112
5.4 Conclusion	115
6 Conclusion and future directions	117
6.1 Summary and Conclusions	117
6.2 Future Directions	119
BIBLIOGRAPHY	120
LIST OF PUBLICATIONS	131

List of Tables

3.1	MSE comparison for three methods for Numerical Example	54
3.2	MSE comparison of three different methods for batch PMMA process	59
3.3	MSE comparison of three different methods for batch transesterification process	59
3.4	MSE comparison of the proposed approach in presence of outliers for batch transesterification process	62
3.5	MSE comparison of the proposed approach in presence of outliers for Numerical Example	62
3.6	MSE comparison of three different methods for different window length in Numerical Example	63
3.7	MSE comparison of three different methods for batch-to-batch variation(k) for batch PMMA process	64
3.8	MSE comparison of three different methods for batch-to-batch variation(I) for batch PMMA process	64
3.9	MSE comparison of three different methods for batch-to-batch variation (k) in batch transesterification process	64
3.10	MSE comparison for three methods for batch-to-batch variation(k) in Numerical Example	65
3.11	MSE comparison of three different methods for batch transesterification process (only space weights)	66
3.12	MSE comparison of three different methods for batch transesterification process (space and time weights)	66
4.1	Hyperparameters for TD3 algorithm	83
4.2	Tracking error (in terms of RMSE values) comparison of five different RL algorithms for batch transesterification process	84
4.3	Variability in control action for four different RL algorithm for the PID reward for batch transesterification process	87

4.4 Comparison of control effort for continuous action space algorithms for batch transesterification process	87
4.5 Tracking error (in terms of RMSE values) of three different RL algorithms for exothermic batch process	88
4.6 Comparison of control effort for continuous action space algorithmsbatch exothermic process	88
5.1 Hyperparameters used in DNNs : TASAC	108
5.2 Comparison of controller performance for different studies for TASAC algorithm (for 10 different random seeds)	109
5.3 Comparison of controller performance for TASAC & SAC algorithm (under nominal conditions)	110
5.4 Comparison of controller performance for TASAC and SAC algorithm in the presence of measurement noise	112
5.5 Comparison of controller performance for TASAC and DDPG algorithm (under batch to batch variation)	113

List of Figures

2.1 Comparison between traditional and JIT modeling approaches	13
2.2 RL Framework	18
2.3 Course of RL Algorithms	21
2.4 RL Approaches	23
2.5 Q Table	25
2.6 Deep Q Learning	27
2.7 Actor Critic Architecture	34
3.1 Sample Search Strategy	46
3.2 Structure of the cost matrix (C)	46
3.3 Numerical Example: Generated data for multiple batches	55
3.4 Prediction Performance comparison of the three methods for 50 time step-ahead prediction for numerical example	56
3.5 Batch PMMA process: Generated data for multiple batches	57
3.6 Prediction performance comparison of three methods for 50 time step-ahead prediction for batch PMMA process)	58
3.7 Batch transesterification process :Generated data for multiple batches	60
3.8 Prediction performance comparison of the three methods for 50 time step-ahead prediction for batch transesterification process	61
4.1 Schematic of the TATD3 based controller for batch process	76
4.2 Comparison of tracking performance of a) TATD3, TD3, and DDPG for PID reward and b) TATD3,TD3, and DDPG for PI reward for batch transesterification process	86
4.3 Comparison of tracking performance of a) DQN and GP for PID reward b) DQN and GP for PI reward for batch transesterification process	91

4.4 Comparison of (a) control inputs of all approaches for PID reward (b) penalty for TATD3, TD3 and DDPG for PID reward for batch transesterification process	92
4.5 Comparison of controller performance of TATD3 offline and online learning(for batch transesterification process)	93
4.6 Comparison of tracking performance of a) TATD3, TD3, and DDPG for PID reward and b) TATD3, TD3, and DDPG for PI reward for exothermic batch process	94
4.7 Control input profiles: TATD3, TD3, and DDPG for exothermic batch process	95
4.8 Comparison of tracking performance of TATD3 vs. TD3 vs. DDPG for PID reward(for batch-to-batch variation in batch transesterification process)	95
5.1 Schematic of the TASAC based controller for batch process	100
5.2 Selection of best action	102
5.3 Selection of TV	103
5.4 Average rewards for 10 different seeds a) TASAC based controller b) SAC based controller. TASAC based controller achieves better reward when compared to SAC based controller.	111
5.5 Comparison of tracking performance of SAC and TASAC controller	112
5.6 Comparison of tracking performance of SAC and TASAC controller (measurement noise)	113
5.7 Comparison of tracking performance of TASAC and DDPG controller (batch-to-batch variation)	114
5.8 Percentage improvement in ITAE using TASAC comparing with SAC and DDPG	115

ABBREVIATIONS

JITL	Just-in-Time Learning
RL	Reinforcement Learning
MPC	Model Predictive Control
TD	Temporal Difference
MDP	Markov Decision Process
MC	Monte Carlo
DQN	Deep Q Network
DNN	Deep Neural Network
PG	Policy Gradient
DDPG	Deep Deterministic Policy Gradient
TD3	Twin Delayed Deep Deterministic Policy Gradient
TATD3	Twin Actor Twin Delayed Deep Deterministic Policy Gradient
SAC	Soft Actor Critic
TASAC	Twin Actor Soft Actor Critic
FAME	Fatty Acid Methyl Ester
MSE	Mean Squared Error
TV	Target Value

NOTATION

s, a, r, s'	State, action, reward, next state
γ	Discount factor
G_t	Discounted return
$\phi_{\mathbf{C}_j}$	Q-network weight
$\phi_{\mathbf{C}_j, \mathbf{T}}$	Target Q-network weight
$\phi_{\mathbf{A}_i}$	i^{th} Actor network weight
$\phi_{\mathbf{A}_i, \mathbf{T}}$	Target Actor Network weight
$\mu_{\phi_{\mathbf{A}_i}}$	Deterministic parameterised policy
π	Policy
π_{ϕ_A}	Policy corresponding to parameter ϕ_A
π^*	Optimal policy
$\pi_{\phi_A}(\cdot s_t)$	Stochastic policy parameterised with parameter ϕ_A
$\mu_{\phi_A}(s_t)$	Deterministic policy parameterised with parameter ϕ_A
\tilde{a}_i	Target action
τ	Target update rate
β_C	Critic learning rate
β_A	Actor learning rate
k_i	Rate constant
T_r	Reactor temperature
T_j	Jacket temperature
T_{jin}	Jacket inlet temperature
α	Temperature parameter

$\mathbf{H}(\pi(\cdot|s_t))$ Entropy of policy π