

**Change Orders Predictability in Construction Projects and Ways to
Improve it – A Data-Driven Study**

Nariman Nabipour

A Thesis in
The
Department of
Building, Civil, and Environmental Engineering

Presented in Partial Fulfillment of the
Requirements for the Degree of Master of
Applied Science in Building Engineering at
Concordia University
Montreal, Québec, Canada

November 2024

© Nariman Nabipour, 2024

CONCORDIA UNIVERSITY
School of Graduate Studies

This is to certify that the thesis prepared

By: Nariman Nabipour

Entitled: **Change Orders Predictability in Construction Projects and Ways to Improve it – A Data-Driven Study**

and submitted in partial fulfillment of the requirements for the degree of
Master of Applied Science (Building Engineering)

complies with the regulations of the University and meets the accepted standards with respect to originality and quality.

Signed by the final Examining Committee:

_____ Chair
Dr. Osama Moselhi

_____ Examiner
Dr. Rebecca Dziedzic

_____ Examiner
Dr. Osama Moselhi

_____ Supervisor
Dr. Mazdak Nik-Bakht

Approved by _____
Dr. Mohamed Ouf, Graduate Program Director

Dr Mourad Debbabi, Dean of Faculty

Abstract

Change Orders Predictability in Construction and Ways to Improve it – A Data-driven Study

Nariman Nabipour

Change orders are formal modifications to the scope of construction projects and play a critical role in shaping project outcomes. However, predicting change orders early in a project's lifecycle remains challenging due to their complex nature and the diverse factors influencing their occurrence. These include project specifications, spatial properties, and the performance of involved actors. Despite extensive research in construction change order management, predictive models for change orders have received limited attention.

This study aims to address this gap by exploring the correlation between project attributes and change order occurrences. It focuses on abstract and easily accessible project attributes such as project type and spatial features, which are more readily shared across the industry. A new set of attributes, derived from domain knowledge, is introduced to quantify project-specific change performance and enhance the prediction of change severity.

The study also tackles the challenge of change timing prediction by modeling the temporal dependence of change orders. Using a Markov Chain approach, it simplifies the relationship between change orders issued in different project phases, assuming each phase's outcome depends solely on the preceding phase. The validity of this assumption is tested through the Chapman-Kolmogorov equation across projects of varying durations.

Results demonstrate a 15% improvement in change severity prediction performance, highlighting the effectiveness of the introduced attributes and feature selection techniques. The findings also confirm the suitability of the Markov Chain model for capturing temporal dependencies in change orders' severity, offering valuable insights for early change prediction and better project planning.

Acknowledgment

Foremost, I wish to express my deepest gratitude to my research supervisor, Dr. Mazdak Nik-Bakht, for his continuous support, patience, enthusiasm, and motivation. His guidance and trust have been invaluable throughout my research career at Concordia University and in my personal growth. His inspiration has been a constant source of motivation during my studies and in my life.

I would also like to extend my heartfelt thanks to the members of my committee, Dr. Rebecca Dziedzic and Dr. Osama Moselhi, for their valuable time and insightful comments.

My sincere appreciation goes to Dr. Shahin Karimi for his kind support and guidance during my studies at Concordia University. I am equally grateful to my friends and colleagues at COMPLECCiTY lab for their camaraderie and support, making the lab feel like a second family.

I also acknowledge the support from the Natural Sciences and Engineering Research Council of Canada (NSERC) and our industry partner, AEDO AI.

Finally, special thanks to my partner, Sahar, for her unwavering support and understanding throughout my studies. This path couldn't have been completed without her, and my other family members' support.

Table of Contents

Table of Figures and Tables.....	vii
Chapter 1 – Introduction.....	1
1.1 Motivation and background.....	1
1.2 Problem statement.....	2
1.3 Objectives.....	3
1.4 Organization of the thesis.....	4
Chapter 2 – Literature review	5
2.1 Change reasons	5
2.2 Change propagation	6
2.3 Change impacts.....	6
2.4 Other factors influencing change impact.....	7
2.5 Change management	8
2.6 Change prediction	9
2.7 Literature gaps	12
Chapter 3 – Methods	15
3.1 Scope of study.....	15
3.2 High-level methodology.....	16
3.3 Change severity prediction.....	17
3.3.1 Initial feature selection and modeling.....	17
3.3.2 Multi-step feature selection and model tuning.....	19
3.4 Analysis of COs’ sequential behavior.....	21
Chapter 4 – Data understanding and preparation.....	24
4.1 Data understanding.....	24
4.1.1 Data acquisition.....	24
4.1.2 Data description.....	25
4.1.3 Data quality	26
4.2 Data preparation.....	27
4.2.1 Data cleaning.....	27
4.2.2 Record selection	28
4.2.3 Data aggregation and integration.....	30
4.2.4 Handling missing data.....	32
4.2.5 Data labeling, encoding, and standardization.....	33
Chapter 5 – Change Orders Cost Impact Prediction	35

5.1 Modelling.....	35
5.2 Multi-step feature selection	37
5.2.1 Attribute type evaluation.....	37
5.2.2 Collinearity removal	39
5.2.3 Permutation-based feature selection	41
Chapter 6 – Analysis of change orders sequential behavior	43
6.1 PCIC model	43
6.2 Testing for Markovian property.....	44
Chapter 7 – Discussion	47
7.1 Feature importance.....	47
7.2 Data mining perspective	50
7.3 Change orders sequential behavior	51
Chapter 8 – Concluding remarks.....	54
8.1 Summary of the study and its key findings	54
8.3 Research contributions	55
8.4 Limitations of this thesis	56
8.5 Future studies.....	57
References.....	59
Appendices.....	71
Appendix 1) 1-step transition matrices for SCI and CCI tests	71
Appendix 2: Observed and expected 4-step transition matrices.....	73
Appendix 3: Data sample	75

Table of Figures and Tables

List of Figures

Figure 3-1: High-level methodology	17
Figure 4-1: Procure change management process (“Change Workflow Diagrams” 2017)..	25
Figure 4-2: Distribution of missing segment of project duration in change orders data.....	30
Figure 5-1: PCA cumulative variation per PC.....	35
Figure 5-2: Level of correlation between predictors.....	40
Figure 5-3: Performance improvement vs clipping threshold	40
Figure 5-4: Permutation and Entropy Results.....	41
Figure 5-5: Performance variation vs permutation cycles.....	42
Figure 7-1: Entropy results before and after feature selection steps.....	49
Figure 7-2: Periodic prediction with/without Markov assumption.....	52

List of Tables

Table 2-1: Review of predictive change order impact studies.....	12
Table 4-1: The raw attributes exist in the Datasets	26
Table 4-2: Sample data points with noise sources	27
Table 4-3: Rules for identification of couples and duplicates	28
Table 4-4: Strategies to deal with the missing values	32
Table 5-1: Dominant attributes in the first five PCs	36
Table 5-2: Chi-squared test results.....	36
Table 5-3: Models' performance comparison	37
Table 5-4: Attribute groups and contained features.....	38
Table 5-5: Attribute-type level feature importance analysis	39
Table 6-1: Project categories by project duration	45
Table 6-4: CHI-Squared results	45
Table 7-1: Dominant attributes of PCA after feature selection	50

Chapter 1 – Introduction

1.1 Motivation and background

Change Orders (COs) have been a challenge to construction projects from the early formation of the industry (Ibbs et al. 2001). Several changes occur throughout the project lifecycle because of fragmented design and construction processes (Shabani Ardakani and Nik-Bakht 2021). They are known as being major causes of labor productivity loss (Al-Kofahi et al. 2022; Golnaraghi et al. 2020a; b; Ibbs 2005, 2012, 2021; Love et al. 2017; Rathnayake and Middleton 2023; Zhao 2023), schedule growth (Chua and Hossain 2012; Fathi et al. 2020; Shrestha et al. 2022, 2017; Shrestha and Fathi 2019; Shrestha and Maharjan 2018, 2019), cost overrun (Anastasopoulos et al. 2010; Fathi et al. 2020; Ibbs et al. 2001; Senouci et al. 2017; Serag et al. 2010; Shrestha and Maharjan 2019), conflicts between stakeholders (Bordat et al. 2004; Kim et al. 2020), and construction disputes (Chou et al. 2013; Fathi et al. 2020; Ibbs 2021; Semple et al. 1994; Shrestha and Fathi 2019; Shrestha and Maharjan 2019; Wang et al. 2023). Zooming on the cost impact, e.g., COs on average increase projects' cost by 10% to 15% of the contract value (Sun and Meng 2009). This was equal to 50-58 B CAD in 2022 only for the Canadian construction industry. Comparing this amount with the 6.7 B \$ property damage of the 2001 World Trade Center attacks (Grossi 2009), it would be translated to the cost of almost five times demolishing and reconstructing the Twin Towers in the Canadian construction industry, and this is only the direct cost impact of change orders in Canada in 2022.

Projects impacted by change orders do not suffer from the diverse effects at the same level. The impact level varies significantly depending on the time and severity of change orders. (Hanna et al. 1999a; b; Ibbs 2005; Ibbs et al. 2007; Leonard et al. 1988; Moselhi et al. 1991, 2005; Thomas and Napolitan 1995; Zhao 2023) have shown that the impacts of change orders on labor efficiency are highly dependent on the magnitude and timing of COs. (Ibbs 2005) found that projects with zero change orders can have labor productivity up to 6% higher than planned values and can mitigate the impact of up to 3% change orders. However, when change orders go beyond 3% the productivity degradation can rise to 40% of planned values depending on the magnitude of change orders. On the other hand, when change orders are issued late toward the end, the productivity never reaches the planned values. The productivity loss effects are twice more in projects with late change orders compared to projects with early change orders. (Alnuaimi et al. 2010; Assaf and Al-Hejji 2006; Bordat et al. 2004; Shrestha et al. 2017; Shrestha and Maharjan 2018) found that, as the COs values grow, the cost and schedule increase significantly. (Serag et al. 2010) identified the timing of COs as one of the most influential features of the cost increase due to change order.

The importance of having reliable expectations on the magnitude and timing of change orders and their impact on labor productivity more clearly emerges knowing that labor costs on average take 30% to 50% of the construction project budget (Toomey et al. 2015), and COs can decrease labor productivity by up to 50% (Ibbs 2005; Thomas and Napolitan 1995). As a result, it is important to be able to predict the overall project's expected COs severity and timing early before a project starts to more accurately measure the labor productivity in the planning phase and monitor it during the project's progress. Additionally, cost performance projection on the project timeline has been an important factor for project managers for years, and the critical value of early cost prediction in project success is confirmed by the industry (Attalla and Hegazy 2003; Du et al. 2016b; Kim et al. 2004; Trost and Oberlender 2003). Change predictive analysis can enhance the cost performance analysis methods by providing a more reliable evaluation of upcoming COs.

Predicting the occurrence and/or magnitude of future change orders is a difficult task, due to the complex nature of construction projects leading to the uniqueness of each project. The National Council of Canada names construction projects as “one-off” projects, to emphasize the fact that each construction project is unique with its distinct characteristics (Hao et al. 2008). The uniqueness of construction projects and, as a result, varied change performance of projects, which have roots in the complexity of the construction projects, are associated with several factors, related to the project's physical, spatial, and managerial specifications (Dao et al. 2017). Looking into common routes of construction COs, project characteristics and related features such as location and project type are among those attributes (Heravi and Charkhakan 2015). (Cantarelli et al. 2012; Flyvbjerg et al. 2003b; a) investigated the influence of project location and project types on the cost performance of projects, it was concluded that different project types have distinct cost performances based on the location of projects. The main goal of predictive analysis of COs is to identify the patterns and features that are important in identifying and grouping projects with similar change performance and accordingly draw a line between project groups and their change performance. Shifting the point of view to the drivers of change, COs occur for various reasons. More than 80 change reasons have been identified in past studies (Padala et al. 2022), each related to one or more of the project properties, categorizable in location, actors, and project-related features. For example, “Unforeseen Site conditions”, as a common change reason (Kim et al. 2020) is logically describable by attributes related to the location of projects, such as “availability of as-built underground infrastructure in that location” or “the number of unforeseen conditions identified in previous projects in the area”. One may argue that the same term is also explainable by the attributes related to the specific project attributes regardless of the location of the project, such as whether extensive inspections are performed before starting the project, or at an even more abstract level, by the project inspection budget (Kim et al. 2020). Either looking into complexity or the change reason attributes of construction projects, a great number of features, related to the location, actors, and the project itself need to be collected to capture differences in projects resulting in varied change behavior. These specifications need to be identified, collected, and analyzed for suitability to predictive analysis.

Moreover, addressing the need for change timing prediction uncovers an additional layer of uncertainty associated with COs. The barrier to performance analysis of processes associated with uncertainties, like change orders, lies in the significant fluctuations in performance over time, which exhibit stochastic behavior. Stochastic performance analysis is the method of choice for evaluating the performance of processes evolving over time (Du et al. 2016b). The analysis requires periodic measurement of the performance factor over time and analysis to capture the correlation of the future performance factors to the earlier measured values (Barraza et al. 2004). One variation of stochastic processes which has received great attention among scholars is the Markov chain. A Markov chain is a sequence of events in which knowing the state of the current event the future state is only dependent on the current state and not the previous ones (Li and Zhang 2009).

1.2 Problem statement

A few studies in the past have tackled the problem of change prediction to predict the cost impact of change orders early before a project starts. However, a systematic approach toward assessing the correlation of potentially correlated attributes to COs prediction is missing. On one hand, the construction industry historically suffers from a lack of enriched data, and there is no

short-term relief to this problem. On the other hand, a predictive approach to such a complex problem like COs is highly dependent on the number and quality of data. Though, there is a need to investigate ways to get the most from the available data by investigating the importance of attributes. This helps to revise the data collection strategies to include more appropriate attributes. Accordingly, more relevant data would be collected, and COs prediction would be more achievable. Due to the high level of information available in construction projects, a well-established methodology is required to collect critical data for tracing changes and accurately predicting COs throughout the project process (Motawa et al. 2006). Each project has unique characteristics and conditions, and as a result, has distinct change performance. Therefore, the prediction models need to consider proper project features for prediction models (Motawa et al. 2006) to describe their unique change behavior.

Moreover, although the literature on construction change orders in general, and change prediction studies in specific, have emphasized the critical role of change timing, the predictive studies of this subject are not properly investigated. Change timing prediction models, in addition to project characteristics, need to consider information regarding the prior issued change orders at each stage of the project. The complex nature of change orders and their temporal co-dependence makes this task more challenging. Therefore, as suggested in the literature, it is beneficial to group change orders and analyze the interdependencies between these groups to enhance predictive models for change timing. By Considering the integrated impacts of groups of change orders, the complexity of their co-dependencies can be simplified, making the change timing prediction feasible.(Chua and Hossain 2011; Motawa et al. 2006; Padala and Maheswari 2023).

1.3 Objectives

This study aims to take a data-driven predictive approach toward COs in construction projects to enhance predictive analysis by improving the change severity¹ prediction performance, under the current limitations of available construction data. The study goes beyond project-specific predictors and also covers assessing the “sequentiality” of COs to enhance the change timing prediction. It investigates the applicability of Markov chain theory for modeling the co-dependence of periodic cost impact of COs².

The following four Research Objectives (ROs) are followed to achieve the research goal:

RO 1: Investigating the correlation between change frequency and change severity in similar projects – This will help to understand how the change performances of similar past projects in terms of frequency of change can be used as an indicator for predicting the change severity of future projects. As a result, the probability of change occurrence in different groups of projects can be described to the model and accordingly improving change prediction,

RO 2: Examining the influence of project location features on change prediction – This is important, because the location of projects is associated with several change reasons, and linking projects to the important features of their location can help to identify projects that are more

¹ ‘Change severity in this study refers to the class of cost impact of COs issued in a project such as high, medium, or low.

² Periodic cost impact refers to the cost impact of COs that are issued in a certain phase/period of project.

vulnerable to those change reasons. These kinds of features are usually accessible from open data sources and as a result, can improve the data context with the least level of investment.

RO 3: Improving CO prediction performance by identifying and employing highly influential predictors. This will be done through feature selection techniques and enriching the available attributes with more context to better describe their influence on the occurrence of COs.

RO 4: Investigating the sequential behavior of change orders to understand the dependence of the cost impact of COs on the COs that occurred earlier in the project to enable change timing prediction.

1.4 Organization of the thesis

This thesis is organized into eight chapters, each detailing specific aspects of the research. After introducing the problem statement, goal, and objectives of the study in the current chapter (Chapter 1), a comprehensive review of concepts related to change orders is presented in Chapter 2. It explores existing literature, definitions, and theoretical frameworks, providing a solid foundation for the subsequent analysis. This review helps contextualize the current study within the broader field of change order management, articulating existing gaps in the literature and identifying the scope of the study.

Chapter 3 describes the methodology employed in this study to analyze and improve the predictability of change orders. This chapter describes the high-level research process, and analytical approaches used to identify patterns and factors influencing change orders.

Chapter 4 contains a detailed description of the data used in this study, including its sources and characteristics. It also addresses the challenges encountered in data preparation and outlines the methods used to preprocess the data, ensuring its suitability for analysis. This chapter emphasizes the importance of data quality and preparation in conducting reliable research.

Chapter 5 presents the detailed methods and results of the change degree prediction. It describes the statistical and computational techniques used to forecast the extent of changes and the attributes correlated to the prediction.

Chapter 6 describes the methods followed to investigate the sequential behavior of change orders in the periodic cost impact of change orders. This chapter focuses on methods followed to investigate the co-dependence between the cost impact of change orders.

Chapter 7 provides a discussion of the findings in Chapters 5 and 6. This chapter offers a critical interpretation of the findings, exploring their implications for change order management practices and decision-making in construction projects.

Finally, Chapter 8 contains the conclusions and recommendations derived from the research. It summarizes the key findings, discusses their significance, and suggests practical applications for improving change order management. Additionally, this chapter provides recommendations for future research work, identifying areas where further investigation could enhance the understanding and predictability of change orders in construction projects.

Chapter 2 – Literature review

COs strategy has been a part of construction contracts from the early formation of the modern construction industry. Before the 1930s, construction was mostly executed by local contractors. After this period, projects started growing both in size and duration, which needed local contractors to work as joint ventures. The rise in the complexity of projects and the involvement of several stakeholders gradually evolved the industry resulting in more COs yielding to claims and disputes to deal with (Carty 1995; Stamatiou et al. 2019). To manage the growing number of COs in the industry, in the mid-1970s, a US committee introduced a change condition clause in which contractors were asked to notify the owner of upcoming potential changes as soon as they encountered a different condition from those specified in the contract (Carty 1995). From then, several definitions are introduced for change orders. The Construction Industry Institute (CII) defines a change order as “any event which results in a modification of the original scope, execution time, or cost of work being inevitable on most construction projects due to the uniqueness of each project and the limited resources of time and money available for planning” (Hanna 2001). COs can be initiated by various stakeholders involved in a project, however, they are mostly initiated by owners and contractors (Senouci et al. 2017). Regardless of the parties that initiated the change, various causes can introduce a change order.

2.1 Change reasons

COs may occur for many reasons, change causes are widely studied in the literature. Major causes include but are not limited to poor estimation, design errors, political issues, and client-base variations (Padala et al. 2020; Shrestha and Maharjan 2018). COs are an inevitable part of construction projects because of the unique characteristics of the industry. The National Research Council of Canada (NRCC) identifies construction projects as “one-off” projects, which means each project has its distinct characteristics, processes, and goals. This specific feature plus multiple stakeholders involved in each project, short-term and project-based supply chains, different delivery methods and payment methods, legal issues, and the extension of construction project scope from design and construction to life cycle have made construction project management more complex and Change orders have become an inevitable part of construction projects to manage uncertainties (Hao et al. 2008).

Several studies have investigated the change causes (Alaryan et al. 2014; Alnuaimi et al. 2010; Hsieh et al. 2004; Khalafallah and Shalaby 2019; Khanzadi et al. 2018; Padala et al. 2022; Sun and Meng 2009). (Padala et al. 2020), gathered 85 change reasons and categorized them into five groups; i.e., client-related, design, construction, performance-related, and external COs. The first four categories cover COs related to the poor performances of actors in projects, or specific conditions of projects. Actors-related COs include variations due to poor financial strength, documentation, constructability analysis, communication, and performance of parties involved, and any incomplete or uninformed decision resulting in a variation in the initial scope of projects. The external change orders, however, cover areas mostly related to the location of projects such as unforeseen site conditions, weather conditions, government policies, residential requirements, labor shortage, price inflation, government policies, or introduction of new technologies.

In addition to the explained change causes, it is known that change order sources are not limited to the change initiators. COs themselves can initiate more COs by affecting future project activities. The co-dependence of COs is discussed in the literature as term change propagation.

2.2 Change propagation

According to (Burati et al. 1992; Lee et al. 2005; Love and Li 2000; Zhao 2021, 2023), change orders can trigger more change orders in downstream activities. For example, poor investigation of spatial conditions of projects, such as soil stability can initiate design COs which can be later propagated to the next activities such as foundation design, and cause more COs (Wu et al. 2005). (Rahmani Mirshekarlou 2012) divides the change impact into change single impact and change ripple impact, while single impact covers the direct effects of each individual change order, the ripple impact refers to the effects that a change order has on subsequent COs. (Moayeri 2017) describes change ripple effects as the propagation of a change order to other parts of the project as a series of change orders. Generally, in any integrated system, in which the components are connected, a change in one component can transform the change into other correlated components (Clarkson et al. 2004).

As stated by (Mokhtar et al. 1998), one of the main challenges in change management is how to monitor change order propagation. According to (Chua and Hossain 2012), tracing change propagation is an important factor in reducing the diverse effects of COs. However, there is a gap in identifying how a change order's impact propagates to other change orders within design and management tools (Habhouba et al. 2009; Isaac and Navon 2008). (Clarkson et al. 2004) developed a mathematical model to assess the risk of change propagation in terms of likelihood and degree of change on activity level during a helicopter design process. However, they also acknowledged the dependency of their work on having a clear understanding of the underlying dependency of sub-systems which is not easy to achieve in a highly complex system such as construction projects. (Chua and Hossain 2011) proposed a model for change propagation to show the likelihood of change propagation to the other stages of a project at different degrees of COs. They mentioned that the probability and the level of propagation of change order to downstream activities depend on the degree of change order at the upstream activity. They categorized change degrees in Low, Medium, and High, and developed transition matrices to capture the dependency of adjacent activities and suggested generating the transition matrices from historical data or expert judgments. However, due to the complexity of capturing the interaction between each pair of change orders, they considered the combined change effect as the representative of a group of change orders. (Motawa et al. 2006) also confirmed the complexity of capturing the correlation between project changes and the problems that occur in later phases of projects. To overcome this problem, (Padala and Maheswari 2023) suggested an integrated approach to the quantification of COS co-dependency rather than COs entity-level analysis.

2.3 Change impacts

Despite the essential role of COs in construction projects, their significant adverse effects are proven facts that can be categorized into tangible and intangible effects from one perspective (Padala et al. 2020; Sun and Meng 2009). Cost and time overrun are among the tangible effects of COs, while quality issues, safety issues, and conflicts in relationships between stakeholders are examples of intangible ones (Padala et al. 2020). The magnitude of the two tangible diverse effects of COs; cost, and time overrun, are widely investigated by different researchers. One study on Canadian construction projects showed more than 30% cost overrun in 50% of projects while having up to 80% time extension as a result of COs (Semple et al. 1994). Another study has noted a more than 23% increase in project cost due to change orders (Love et al. 2017). A Survey showed

almost 10% cost and schedule growth in building construction contracts (Shrestha and Maharjan 2018). Several other researchers have confirmed the severe effect of COs on the cost and schedule overrun (Alnuaimi et al. 2010; Choi et al. 2016; Shrestha et al. 2017; Shrestha and Maharjan 2018).

Moreover, in addition to the direct impacts, COs increase the disruption of the cost and time performances of projects indirectly, by lowering the laborers' productivity. Projects impacted by COs have lower labor efficiency than those not impacted (Hanna et al. 1999a; b; Moselhi et al. 1991, 2005). This reduction depending on the change characteristics, can be up to 50% (Thomas and Napolitan 1995). COs reduce the labor productivity from multiple sources such as; stop and go operations, out-of-sequence works, loss in rhythm of productivity, demotivation of the workforce, loss in learning curves, unbalanced crews, excessive labor fluctuations, and unbalancing of successive operations (Leonard et al. 1988).

Although COs are considered a general challenge for all construction projects, the level of impact varies between projects as a result of their unique characteristics. The uniqueness of construction projects and as a result, varied change performance of projects which has roots in the complexity of the construction projects is associated with several factors, related to the project's physical, spatial, and managerial specifications (Dao et al. 2017). (Cantarelli et al. 2012; Flyvbjerg et al. 2003b; a) investigated the influence of project location and project types on the cost performance of projects and found that based on the type of project and its location the level of cost overrun varies.

2.4 Other factors influencing change impact

In addition to the project characteristics that influence the change impacts, various studies have identified additional attributes that can influence the magnitude of the COs' diverse effects. These attributes include change-timing (Hanna et al. 1999a; b; Leonard et al. 1988; Moselhi et al. 2005; Thomas and Napolitan 1995), change severity (Hanna et al. 1999a; b; Ibbs 2012; Moselhi et al. 1991, 2005), change frequency (Hanna et al. 1999a), work complexity (Leonard et al. 1988), dependency of work sequences (Leonard et al. 1988), work intensity (Leonard et al. 1988), change impaction (Hanna et al. 1999a), management team experience (Hanna and Iskandar 2017; Moselhi et al. 2005), and type of work (Moselhi et al. 2005). As a result, the influential factors on change impacts can be categorized as project type, location, actors related, and COs performance attributes. It is noticeable that in addition to the specific project spatial, type, and actor specifications, the level of impact also correlates with change performance attributes such as the timing, magnitude, and frequency of COs.

Depending on the magnitude (Hanna et al. 2002), frequency (Ibbs 2005), and timing (Hanna et al. 2002) of COs, the severity of the diverse effects on productivity, cost, and schedule performance of projects vary between projects. Regarding productivity loss, a study indicated a significant positive correlation between productivity loss and the amount of change (Hanna et al. 2002); i.e., the higher amount of COs results in more productivity losses (Hanna et al. 1999b; Ibbs 2005). The timing of COs is also an important factor; projects with early change implementation can absorb productivity reduction effects for up to 3% change percentage³, however, projects with late COs lack this ability. Moreover, late changes can increase the productivity loss by up to 100%

³ Cost impact of change orders relative to the total contract value

(Ibbs 2005). COs occurring in between 50% and 75% of project progress, have much more severe negative effects than others (Ibbs et al. 2007).

(Serag et al. 2010) investigated the relation between change-timing and project cost performance and found a significant relation between them. (Shrestha and Maharjan 2018) found the same trend for the amount and number of COs in the highway construction projects. They stated if the number and percentage of COs exceed 20, and 5% of contract value respectively, the effect of COs on cost and schedule growth would be significant and severe. To mitigate the diverse effects of change orders it is required to have a well-established change management system to identify change orders as early as possible and deal with their consequences.

2.5 Change management

According to (Hao et al. 2008), standard change management consists of five steps, identify change, evaluate and propose change, approve change, implement change, and analyze change. Change management aims to mitigate or control the negative effects of COs due to their significant influence on project success factors. (Naji et al. 2022) has directly studied the impact of CO management factors on project success. (Ko et al. 2024) has tackled this area by prioritizing project requirements by analyzing historical data on COs. Another approach involves analyzing common CO management practices through social network analysis followed by (Shabani Ardakani and Nik-Bakht 2021). Better documentation and developing change management systems is another area studied by several researchers (Caldas and Soibelman 2003; Du et al. 2016a; Karimidorabati et al. 2016; Motawa et al. 2007; Wang et al. 2024). These studies focus on finding more efficient tools and methods for change documentation to facilitate the coordination of documents between parties involved and reduce the time of change approval. For example, (Karimidorabati et al. 2016) compared three generations of change management and concluded that the time of change, and traceability are improved as a result of using more advanced platform-based systems. Comparative analysis of change performances of projects such as projects with varied delivery methods has also been explored well (Choi et al. 2016; Hanna 2016; Riley et al. 2005; Rojas and Kell 2008; Shrestha et al. 2012; Shrestha and Fernane 2017). For instance, using 77 building projects, (Shrestha and Fernane 2017) compared the performances of design-build with design-bid-build projects and found that the number of COs in design-build projects is significantly less than in design-bid-build projects. (Hanna 2016) evaluated the performances of integrated project delivery (IPD) and concluded that the IPD projects perform better compared to other delivery methods in terms of COs processing time. Another well-explored area is the effects of COs, particularly their impact on the efficiency of construction planning terms, such as schedule and cost analysis (Serag et al. 2010). Labor productivity estimation is another example in construction planning where accuracy is highly dependent on the timing and magnitude of change orders (Al-Kofahi et al. 2022; Hanna et al. 1999b; Hanna and Iskandar 2017; Ibbs 2005; Ibbs et al. 2007; Moselhi et al. 2005; Rathnayake and Middleton 2023). For example, (Hanna et al. 2002) investigated the projects impacted by COs and found percent change, type of trade, change timing, and percent change related to design issues are the main factors contributing to the project impact. (Moselhi et al. 2005) used data from 33 work packages to develop a model for quantifying the impact of change orders on labor productivity. They found intensity (i.e., number or volume of change), timing, work type, type of impact, onsite management, and project phase among the important factors influencing labor productivity loss.

The scope of most of the mentioned studies is limited to the evaluation of the impact or dealing with the consequences after being recognized. COs prediction has comparatively received less attention, likely due to the lack of rich, structured data in the industry. However, with the growing desire for digitalization in the industry and accordingly availability of more data, this subject can be investigated much more confidently. Predictive analysis of COs before project commencement aims to empower decision-makers by anticipating upcoming COs during the planning phase. This proactive approach enhances the ability to accurately assess their potential impact on projects.

(Arefazar et al. 2022; Hegazy et al. 2001; Motawa et al. 2007; Sun et al. 2006) have highlighted the identification and prediction of change orders as a primary challenge in change management. Proactive change management should be capable of predicting the occurrence and severity of change orders and coordinating changes during the project (Motawa et al. 2006). According to (Love et al. 2015), notifying decision-makers of COs significantly enhances the effectiveness of infrastructure project management, improving both performance and contingency planning. Moreover, having access to this information allows decision-makers to utilize advanced technologies and methods to more effectively mitigate disputes arising from change orders. (Hasanzadeh et al. 2018).

2.6 Change prediction

The predictive approach to evaluating COs impact before happening has been studied in two directions. A group of studies provided a quantification of the probability of change order occurrence at the organizational or project-type level (i.e., grouping projects based on their type of projects) without differentiating between individual projects (Anastasopoulos et al. 2010; Ibbs 2012; Shrestha et al. 2019). (Anastasopoulos et al. 2010) employed statistical methods to analyze the frequency of COs in Indiana highway projects and how it varies between projects with different durations, contract types, and project types. (Shrestha et al. 2019) used data from two studies to develop a baseline for comparing project change performance against other projects. The result of their work was curves showing the probability of having a certain amount of change orders for all projects. (Shrestha et al. 2019) used data from 95 public school building projects and employed a regression model to show the probability of having different levels of COs. The main shortage of their work is the incapability of their work to evaluate the projects based on their unique characteristics. In their approach, all projects receive the same probability values for experiencing a level of change which cannot be true in construction projects. Another approach to this problem is followed based on the machine-learning techniques to utilize AI (Artificial Intelligence) capabilities to analyze CO data from past projects along with considering the specific characteristics of each project to draw a path from features of projects to their change behavior. On one hand, the advantage of this approach is the ability of these methods to perform the analysis on the project level, i.e., the model can predict a separate value for the change performance of each project. On the other hand, the limitation of these models is their great dependence on the size and quality of data (Budach et al. 2022; Foidl and Felderer 2019; Hagedorff 2021). Moreover, due to the complex and dynamic nature of construction projects (Nik Bakht and El-Diraby 2015), the data should encompass various aspects and inherently be complex. The construction industry has historically fallen behind in adopting modern technologies (Bilal et al. 2016; Mitropoulos and Tatum 1999; Sardroud 2015; Yap et al. 2022), and as a result, there may be limited access to high-quality historical data that can be used to train predictive models (Yan et al. 2020).

Few articles exist in the literature on change prediction. (Williams and Gong 2014) used data from 1221 highway projects gathered from the California Department of Transportation to predict the cost growth of projects. They applied text mining methods on the five highest dollar value cost lines of bid documents and generated word tokens to be used as predictors. These tokens along with some numerical attributes such as the number of bidders, the low bid value, and the project's actual cost were used to predict the class of cost overrun. This study mentioned a lack of features related to the location of projects, projects' complexity, and contractors' performance as the limitation of this study (Williams and Gong 2014). It is worth mentioning that "class" is the term used in data mining classification problems for each of the distinct categories the prediction can be assigned to. For example, if the prediction goal is to differentiate between cats and dogs given some photos of both groups, Cat, and Dog are considered as the two classes of the prediction. In change prediction studies, change class refers to change severity, i.e., having a certain range of COs percentage (i.e., the cost impact of change order divided by project value). As an example, considering more than 10% as a high level of change, 'high' can be a class of CO degree. Each study may define change classes based on different criteria, defining certain thresholds, such as 10%, 4%, or minus is one of the methods of defining different classes. This approach usually results in having an unbalanced dataset. Meaning that the number of projects belonging to each class varies a lot compared to other classes. Having an unbalanced dataset is considered as one of the data mining classification challenges which can reduce the performance of the prediction if not being dealt with properly. Another method to define the classes is to define the thresholds that separate the classes based on reaching a balanced dataset, i.e., the thresholds are defined based on how the data points (i.e., projects) are distributed in respect to their prediction attribute (i.e., change percentage). (Williams and Gong 2014) defined three classes with the second approach to obtain an almost balanced dataset. They considered projects with more than 6% change percentage as class high overrun, between 6% and -3% as near original bid class, and less than -3% as underrun.

Another study, by (Ibbs and Chen 2015), used 2,000 construction projects' data from the owner's perspective to predict both the numerical percentage and the class of COs magnitude. For both types of tests, they developed three separate models to have separate predictions for the total change, design change, and construction change. In this study the classes were defined using the clustering method, however, they didn't mention what attributes were considered in clustering. They defined three classes as class high, medium, and low which were projects with more than 10%, between 10% and 1%, and less than 1% COs respectively. One of the main issues with the numerical prediction model of this study, is the great difference between the reported R-squared and the Cross-Validation values, showing up to 50% with an average of 24% difference for the three models. In addition, the reported cross-validation values were the maximum values among all the cross-validation tests, however, the average cross-validation values can show the model performance on an unseen dataset not the maximum. Another issue with this study was having multiple linearly correlated attributes as predictors. As an example, the numerical prediction used ten attributes related to the project budget, contingency, and duration. Four out of the ten are the total, construction, design, and procurement budgets. Obviously, the total project budget is the summation of the other three attributes. The same issue exists with four other features that represent different aspects of the contingency budget. Another problem is the size of the test dataset. Out of 2000 projects, only around 100 projects are used in the training and 21 as the test set. The choice of the prediction performance criterion is another issue with this study. The correlation coefficient comparing predicted change values and the observed values were employed to examine performance and this explains the great R-squared differences between the test and

training. However, this is not a common prediction performance measurement factor and cannot show the level of prediction errors. For the classification model, in addition to these attributes, project complexity, project nature, fast tracking, and the industry group were used as categorical attributes. Accuracy is used as the criteria for prediction performance evaluation. As earlier mentioned, this study uses an unbalanced dataset, however, accuracy is not a proper choice in the case of having an unbalanced dataset. In these cases, other criteria such as F1-score should be used along with other performance metrics to reliably show the performance of the model for all the target classes. The same issue with the number of data points also exists in the classification section.

Another study by (Shafaat et al. 2022), used data from 5,628 transportation projects to classify the projects into two classes. They conducted two sets of binary classification, one targeting the prediction of projects with/without change orders, and another one, projects having less or more than a 5% change percentage. This study considers the location of projects as a categorical attribute with seven classes each representing a district and another attribute to describe the route type of projects along with project type, project awarded price, the awarded price difference with the second bid, the awarded price difference with the initial estimation, project length, contractor name, and project year. In this study, the model for having more than 5% change percentage wasn't able to accurately classify the projects showing 58% accuracy, 35% precision, 76% recall, and 48% F1-score.

One of the most recent studies⁴ by (Nabipour et al. 2023) used 2002 construction and service projects to investigate the predictability of COs and their sensitivity to the number of classes. Two scenarios were defined regarding the number of classes, one a three-class scenario and another one a binary classification. For the multiclass classifier, two thresholds were picked after the literature review and identifying the problematic change percentage levels. The other scenario investigates the prediction performance when being subjected to identifying projects that go beyond a certain level of change percentage. Three different thresholds were picked for the binary classification and the prediction performance was followed separately for each threshold. The results of this study showed that the models were able to predict the change degree with higher confidence in binary classifications.

Table 2-1 shows a brief review of the past studies on the prediction of COs. Exploring the predictors used in these models, it is identical that, due to lack of data, all the mentioned studies have picked their predictors based on the available features in the data, however, the usefulness of these attributes is not investigated and is questionable. One problem reported by all three studies is the class imbalance which affects the prediction performance significantly. It is important to deal with this issue effectively and use proper performance metrics to have a reliable performance measurement. F1-Score is one of the measurement metrics that can reliably report the performance in these cases, whereas accuracy alone cannot effectively describe the prediction performance. Among the three studies, only one has reported F1-Score which varies between 0.47 to 0.51. It is identical that, with the low level of context available in construction project data, it is not an easy task to receive high prediction performance, and other methods and attributes should be investigated to overcome this problem.

⁴ By the time of drafting this thesis

On one hand, the number of change orders is an important factor having a great influence on the cost impact of COs which has not been used in previous studies due to not being known before starting projects. On the other hand, the importance of the attributes related to the location of projects has not been investigated well and needs to be explored. As earlier discussed, the importance of spatial attributes in COs occurrence is not negligible being associated with several change reasons and varied risk levels (El-adaway et al. 2018; Lee et al. 2015).

Table 2-1: Review of predictive change order impact studies

Author	Spatial features	F1-Score	Accuracy
Nariman Nabipour (2023)	Location Name, Demographic features	NA	0.68-0.88
Ali Shafaat (2022)	Location Name, Route Type	0.48-0.51	0.47-0.88
William Ibbs (2015)	NA	NA	0.66-0.81
P. Williams (2014)	NA	NA	0.45

Furthermore, as discussed in section 2.4, change-timing is another critical aspect of COs that significantly affects the impact of COs (e.g., in quantification of labor productivity loss due to COs). Several studies have quantified the impact of change timing on productivity (Ibbs 2005; Moselhi et al. 2005; Zhao 2023). To account for change timing, they divided projects' durations into five equal periods and calculated the change percentage for each period as the cumulative impact of COs issued during each period. This term is then used to consider the time of COs or change timing in labor productivity estimations. In this study, the term Periodic Cost Impact of Change orders (PCIC), refers to the same terminology. It is known that COs can influence the outcomes of upcoming change orders or initiate new ones, which is known as the ripple effect of change orders (Zhao 2021, 2023). However, the correlation between the aggregated cost impact of COs (PCIC) issued in one phase and the next phase is the area that needs further investigation. Understanding the relationship between change impacts of different periods helps to more accurately model the change impact process. This can further be used for change-timing prediction and change monitoring by reducing the complexity of models as discussed in section 2.2.

2.7 Literature gaps

The scope of this study is the prediction of change orders and addressing the limitations on this subject. As mentioned in section 2.1, so far more than 80 change reasons have been identified in the literature and each is associated with different features of projects such as parties involved, project location, and other project-specific properties. To accurately predict a complex phenomenon such as COs, data availability, and quality are among critical factors. However, the construction industry does not do well in terms of data collection and data availability. In addition, the previous studies were not able to achieve acceptable prediction performance and accordingly predict the projects' COs severity.

As a result, it is important to investigate the importance of attributes getting used for such an analysis and make sure that the full potential of available data is unlocked to improve the prediction performance. In addition, as mentioned in section 2.4 the timing of COs is an important

factor that has not been considered in previous COs prediction studies. As explained in section 2.2 the analysis of co-dependence of periodic cost impact of COs can help to reduce the complexity of change timing prediction, which has not been studied earlier. To summarize, the gaps in change prediction can be stated as the following.

- The few studies that have tackled the change prediction problem use different types of attributes as predictors. The suitability of these attributes in change prediction needs to be studied;
- Several studies including the ones on change prediction have mentioned the importance of change timing prediction. The change-timing in these studies refers to the periodic cost impact of COs. The sequential behavior of the cost impact of COs needs to be studied to understand how COs of one phase affect the upcoming phases;
- The past change prediction studies suffer from a lack of reliable results due to low prediction performance and the use of unrealistic performance measurement criteria. It is important to develop prediction models that can reliably predict the change degree using proper performance indicators to make COs prediction operational;
- Change timing is an important aspect of COs that needs to be studied for the possibility of being predicted;
- One of the main aspects of COs that can help project planners to better consider the effect of COs on their projects is the change reason. However, none of the previous studies have studied the possibility of predicting COs with respect to their causes.

Accordingly, this thesis addresses the first three gaps in change order (CO) prediction to enhance the predictability of COs through various feature selection techniques aimed at identifying key attributes, thereby improving change prediction performance. Three groups of attributes are selected: project type, frequency, and spatial features. However, this study does not encompass other attribute types, such as project drawings, bid documents, and performance-related attributes of the actors involved in the projects.

Additionally, by analyzing the sequential behavior of change orders, this research aims to understand the temporal dependency of the periodic cost impact of COs. Understanding the co-dependency of these periodic cost impacts opens avenues for predicting the timing of changes by providing insights into whether change orders issued during a specific period should be considered independent of those issued earlier in the project.

As this research approaches these subjects as a classification problem, it focuses on predicting the class of the cost impact of COs—specifically, high or low—based on whether the cost impact exceeds a defined threshold. Furthermore, the study does not aim to predict individual COs but rather the aggregated cost impact of all COs issued within a project, along with the temporal co-dependency of this aggregated cost impact across different project phases.

Chapter 3 – Methods

In this chapter, based on the gaps identified in the literature (section 2.7), and the research objectives (section 1.3), the methodology to achieve the objectives is elaborated and justified. Firstly, in section 3.1 the scope of this study is structured, followed by section 3.2 which describes the high-level methodology of the study. Secondly, in section 3.3, a comprehensive explanation of the process followed to develop change severity prediction models is described. Finally, section 3.4 explains the fundamentals and methods for assessing the sequential behavior of change orders.

3.1 Scope of study

This study uses data mining methods to answer the research gaps articulated in Chapter 2 to improve change severity prediction early in the planning phase of projects. In other words, it aims to assess the predictability of the aggregated cost impact of all construction change orders issued in a project as a percentage of the contract value as a classification problem, i.e., not each instance of CO. This is followed by evaluating the usefulness of different groups of attributes in the prediction and accordingly using more correlated attributes with the prediction to improve prediction performance. In addition, to account for the importance of change timing prediction this study aims to assess the co-dependence of periodic cost impact of COs. As earlier discussed, the occurrence and severity of COs issues in later phases of projects are under the influence of those issued in upper stream phases. As a result, it is important to investigate the sequential behavior of COs issued in different phases of projects to be able to accurately predict the timing of COs. In this study, the sequentiality of COs is evaluated by assessing the Markovian relation between the Periodic Cost Impact of COs (PCIC)⁵.

Accordingly, the scope of the study is identified as exploring the importance of more contextual attributes regarding the location, and type of projects. Additional context is embedded into the data by adding more features regarding the project location and using the change performances⁶ of past projects to quantify the change performances of similar projects in terms of location and project type. The spatial attributes are considered due to their important role in several change causes such as unforeseen site conditions, and labor shortages. On one hand, past studies have shown that the location of projects is an important factor in the level of cost growth of a project. On the other hand, considering the lack of context in construction project data, spatial features can be integrated with available data from external sources to add the required missing context to the data.

Past studies have mentioned the importance of the frequency of COs and their correlation to the severity of COs. Moreover, several studies have mentioned that projects in different locations, or with distinct project types experience different levels of COs. This study uses these findings to test the suitability of using the average frequency of COs in projects with the same location or type of project as an indicator to quantify the varied change performance in projects with different characteristics. These attributes are analyzed through various feature selection techniques with a look at prediction performance to account for improving the prediction performance. As a result of the feature selection, spatial and frequency attributes along with other features are analyzed to understand the key attributes in change prediction.

⁵ The aggregated cost impact of COs issued in a certain period/phase of construction duration.

⁶ The change performance refers to the number of occurrences or degree of COs projects experience.

Finally, the sequential behavior of COs is investigated in PCIC to understand the co-dependence of PCIC of a phase of a project to its previous phases. The goal of change-timing prediction is to predict the cost impacts of COs issued in a certain phase of a project. So, it is important to know whether the prediction has to be done solely for each phase or whether the COs issued in earlier phases of projects affect the prediction. This study aims to answer this question by assessing the suitability of the Markovian Chain to model the co-dependency of COs issued in different phases

3.2 High-level methodology

The research methodology of this study, as shown in Figure 3-1, consists of two main components, i.e., (i) change degree prediction; and (ii) assessment of sequential behavior in change orders. Similar to other data-driven studies, data needs to be investigated for issues and be pre-processed before analysis. The tasks followed to bring the data to the proper shapes of this study are described in detail in Chapter 4. In addition to the common data preparation tasks, to answer the first two research objectives, (i.e., analyzing the correlation of frequency of change orders and spatial attributes with change degree prediction) two sets of attributes are introduced and integrated with original attributes, (i) Historical Change Performance Indicator (HCPI), and (ii) demographic spatial features. A detailed explanation of these attributes, the concept behind them, and the development process is provided in section 4.2.3. The HCPI attributes are constructed using the domain knowledge of COs to quantify the average frequency of COs based on different project types, and locations. The demographic attributes are integrated using an external dataset to add more context regarding the demographic specifications of projects. To this aim, first, a preliminary feature selection is followed through Principal Component Analysis (PCA), and Chi-Squared methods to select an initial list of attributes to be used as predictors for testing different prediction algorithms. Then, several prediction models are developed, tuned, and compared based on their performance in classifying projects into two classes of change severity, i.e., high, and low. Then the model with the best classification performance is used to conduct different feature selection techniques such as permutation, entropy, and collinearity removal. As a result of the feature selection techniques the model is provided with more contextual features and accordingly the prediction performance improves. Section 3.3, and Chapter 5 provide a detailed explanation of the methods and findings of the change degree prediction, addressing the gap in the literature concerning the need for feature selection.

To answer the fourth research objective (i.e., investigating the sequential behavior of COs) the suitability of Markovian processes in capturing the co-dependence of COs at the construction phase level is analyzed. To this aim, the timeline of construction duration is divided into five equal periods. Then the change percentage is calculated for each period/phase individually to explore how the change severity of a certain period is dependent on the change severity experienced earlier in a project. Then the transition matrices are calculated once for each period and its next period (i.e., 1-step transition matrix) and another time for period one and the last period (4-step transition matrix). The Markovian behavior is evaluated by comparing the observed change severities in the last period and the expected change severities calculated under the Markovian assumption using the Chapman Kolmogorov formula. The comparison is conducted using the Chi-squared method to check the similarity of the observed and expected values and reject or validate based on their level of dissimilarity. The detailed explanation of the methods and the findings are explained in section 3.4 and Chapter 6.

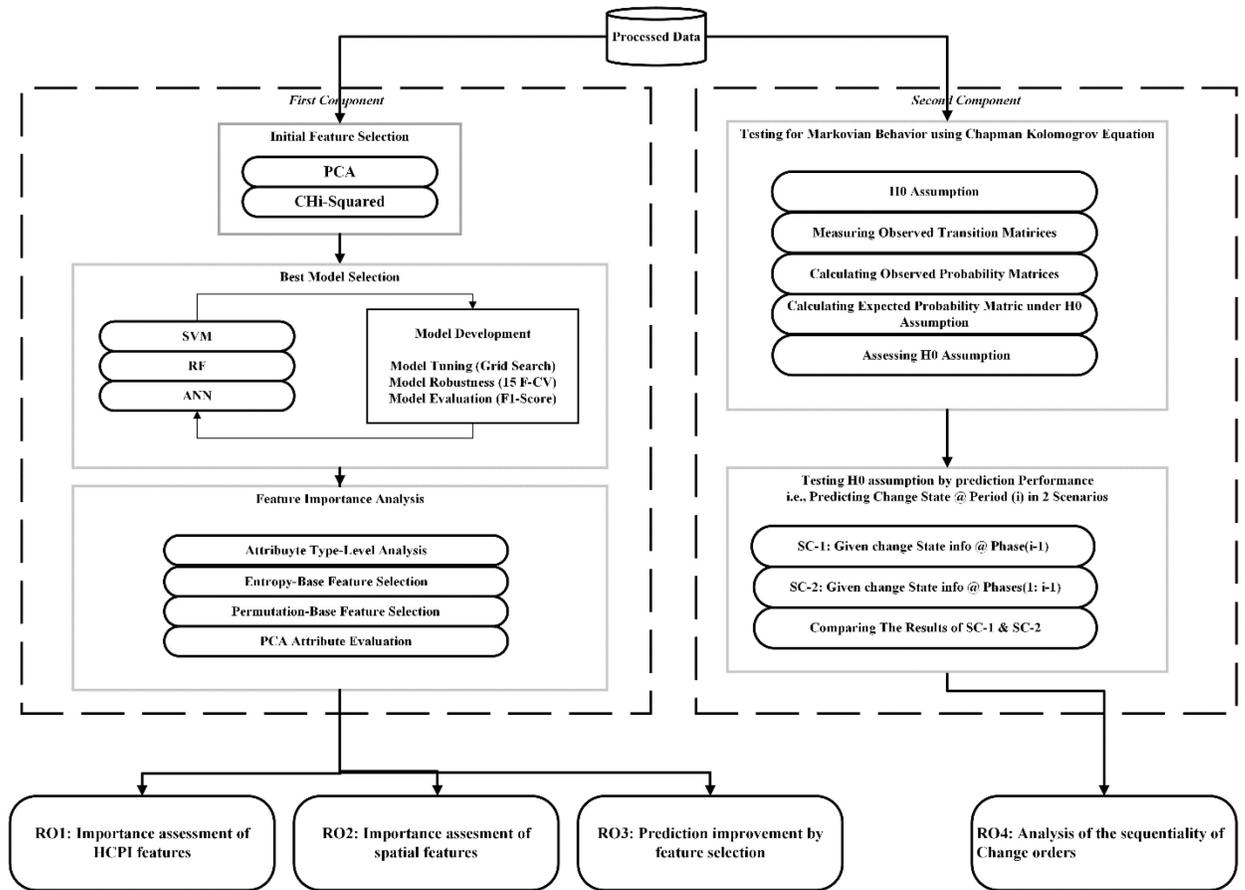


Figure 3-1: High-level methodology

3.3 Change severity prediction

This section provides a brief explanation of the approaches followed to improve the predictability of CIC (Cost Impact of Change orders) and investigates the importance of attributes. First, a preliminary feature selection step described in 3.3.1 is conducted to feed the models with proper attributes and select the best model. Three variants of classifiers, i.e., Support Vector Machines (SVM), Random Forest (RF), and Artificial Neural Networks (ANN) are tested to choose the best model. Second, a multi-step feature selection method is employed to evaluate the importance of attributes and improve the list of predictors. The feature importance analysis is described in section 3.3.2.

3.3.1 Initial feature selection and modeling

To initially evaluate the usefulness of attributes before using them as predictors, two feature selection techniques are employed based on the type of predictors. Principal Component Analysis (PCA) is used for numerical, and the Chi-Squared test is employed for categorical attributes.

PCA is a feature engineering method that projects attributes into another space where variations between data points can be represented using a smaller number of features. This method is frequently used for dimensionality reduction; however, it is also useful in feature selection by

tracing attributes contributing more to the variation of dominant PCs (Principal Components). After applying PCA, PCs are ranked based on their contribution to the variation in the data. By selecting the top-ranked PCs and consequently, dominant attributes correlated with them, the top features are identified. The important attributes of each PC can be selected by checking a coefficient named loading. This term shows the weight of each attribute in forming that PC. The Loading value of 0.3 is considered the threshold for selecting important attributes.

Chi-squared (χ^2) is a statistical test to analyze the dependence between two categorical variables. It works by comparing the expected and observed frequency of co-occurrence for two nominal attributes. A comparatively high χ^2 value shows a high level of discrepancy between categorical target and predictor variables. The test also sets a significance level through a p-value parameter (Fisher 1922). At the initial phase, the test is applied to categorical attributes and those with low χ^2 values and a p-value of more than 0.05 are omitted from the dataset.

As stated earlier, three classification algorithms are employed for this study, SVM, ANN, and RF. It must be emphasized here again that while being committed to high-accuracy prediction; the main goal of the paper is not merely to develop some high-accuracy prediction models. Rather, the main goal is to understand the contribution of various factors. Hence, these algorithms are selected as samples to cover different common data-mining-based prediction approaches. One can employ more advanced techniques and perhaps come up with marginally higher levels of accuracy.

SVM classifiers find the optimal hyperplanes that separate the data into different classes of the target attribute. The standard SVM uses a linear hyperplane to separate data points belonging to different classes, however, the differentiation cannot always happen with a linear hyperplane. In these scenarios, other kernel types, i.e., Polynomial, and RBF (Radial Basis Function) can be a better choice. In this study three distinct SVM models, each using a separate kernel are trained to evaluate each individually, i.e., Linear, RBF, and Polynomial. Several parameters are optimized to tune SVM models and avoid overfitting. They include *regularization parameter* (C) to set a trade-off between a low training error and a low testing error; *Gamma* to adjust the shape complexity of the classification boundary; *Degree*, used for the polynomial kernel to test different degrees of the function; and *class weights*, to address the class imbalance. Class imbalance refers to the situation in which the data is not equally distributed between the classes of the target attribute (Smola and Schölkopf 2004).

RF is an ensemble machine learning algorithm that classifies data based on the combination of decisions made by various decision trees. Each tree splits the data based on predictor values, resulting in high-purity (low entropy) segments of the data. The process continues until reaching pure portions or a stop rule, set as a hyperparameter, is activated (Biau and Scornet 2016). Parameters optimized in this study to achieve high-performance RF include *Number of Trees*, to control the trade-off between prediction confidence and computational cost; *Max depth*, to avoid overfitting; *Min_samples_split*, to control data splitting to avoid imbalanced branching; *Min_samples_leaf*, which identifies the minimum number of samples required in each leaf node; and *criterion*, to measure the quality and purity of a split.

Artificial Neural Networks are combinations of interconnected layers, including an input, an output, and one or more hidden layers of micro-computation units, known as neurons. Weights and biases of neurons are set, usually through iterations (epochs) of backpropagation, to minimize a loss function measuring the accuracy of classification (Goodfellow 2016). To reduce the computational cost of training, a number of samples (i.e., batches) are processed together. In this

study, several ANN models are trained with different numbers of layers and complexity of hidden layers. The models' performances are monitored with the growing number of layers as well as the complexity of layers to figure out the best architecture for this study. It was observed that two hidden layers each containing four neurons provide the most efficient architecture. *Relu* function is considered as the activation function for the input layer, while *Sigmoid* is considered as the activation function of the remaining layers. The learning rate coefficient after several tries and errors is considered 0.01 constant, and the Adam optimization function is used to optimize the bias values. Binary-cross entropy is considered as the loss function and the F1-Score is chosen as the performance evaluation criteria. The batch size is assumed to be 32 and the number of epochs is limited to 200. The optimal number of epochs for each scenario is defined by plotting the performance of the model on both training and test sets and selecting the best number of epochs resulting in the highest performance while not being overfitted.

Model Robustness and Hyper Parameter Tuning – Model robustness refers to the ability of the model to perform well on new/unseen data (i.e., test set) to the same degree it works well on the training data (i.e., training set). In this study, 25% of data points are set aside as the test set, and the models are trained on the other 75% using the Grid Search algorithm (for SVM and RF) with 15-fold cross-validation. Since the ANN parameters are defined manually as described in the previous section, the cross-validation technique is applied without the Grid Search algorithm. F1-Score is chosen here as the performance metric, to address the issues with imbalanced classes. Additionally, classes (of the target attribute) are weighted reciprocal to their frequency, to reward the model for correctly predicting instances of minority classes. The weights of classes are defined using Equation 3-1.

$$\text{Class_Weight} = \frac{\text{Total Number of Records}}{\text{Number of Classes} \times \text{Number of records in the Class}} \quad \text{Equation 3-1}$$

3.3.2 Multi-step feature selection and model tuning

After Evaluating the models and selecting the model with the best performance, a multi-step feature selection and model improvement process is designed to investigate the importance of attributes while enhancing the performance of the model. This process contains three levels of feature importance analysis and model tuning. First, attributes are investigated for effectiveness at the attribute type level, and then, a concurrent correlation and permutation analysis is performed to remove collinearity and better highlight the effectiveness of attributes. The importance of attributes is followed based on their Entropy generated from the RF model. Entropy is a term that measures the impurity of data, a high entropy shows that an attribute can perfectly divide the data into different classes. Whereas a zero entropy means that the feature is not able to separate the classes. Entropy is a concept that originated from thermodynamics with various applications in different fields. Entropy quantifies the level of uncertainty or disorder in possible outcomes of a variable. One of its applications in the field of artificial intelligence is to quantify the level of information a variable carries (Ali et al. 2023).

Third, permutation analysis is employed to further investigate the importance of attributes and improve the model efficiency. Permutation is a powerful method for analyzing the randomness of results, which can be used for feature engineering. It works by removing the dependency of an attribute with the prediction by shuffling the values of that attribute and then following the

prediction performance before and after permutation. Permuting a highly correlated attribute should result in a higher level of reduction in performance.

Step 1-Attribute Type Evaluation- In this step, after selecting the best model applicable to the problem, the relevance of attributes in predicting the target variable is assessed at the attribute type level. To enhance the interpretability of analysis and account for the diversity of attributes in our dataset, the attributes are categorized into five distinct groups based on their characteristics and origin. The categories are defined to investigate the influence of specific types of attributes and their potential impact on the prediction task. The model was then run and tuned under different scenarios, each focusing on certain categories of attributes. These scenarios aimed to systematically investigate whether a particular type of attribute held significant importance in predicting the target variable.

Step 2- Collinearity removal- To highlight the important attributes in the prediction of COs severity, it was needed to remove redundant features in the model to avoid confusing the model which can lead to unrealistic feature importance outcomes, and reduced performance. In the short term, the objectives are to underscore the significance of attributes, address collinearity issues, reduce dimensionality, and improve prediction performance.

Given the presence of numerous correlated attributes within the model, each of which exhibits correlations with multiple other attributes, it is essential to identify groups of attributes correlated to each other and systematically remove those with the least level of contribution to the prediction. The correlation matrix of attributes is created using the Spearman correlation method and subsequently, the correlation matrix is converted to the distance matrix. Third, a dendrogram is employed to visualize attribute correlations and identify closely related and potentially redundant attributes. A dendrogram is a hierarchical inverted-tree-like diagram in which attributes are distributed as leaves or branches at the bottom edges of the tree, on the horizontal axis. Each attribute is connected to the central axis of the tree by a line, and the vertical axis shows the distance between attributes by linking the lines of attributes at a height equal to the distance between them. Then, correlated attributes are removed, by manually picking a distance as a threshold and ‘clipping’ the correlations less than the threshold. From each clipped cluster only one attribute with the highest Entropy is retained and other correlated attributes are removed.

Finally, the model is trained again with the remaining attributes, and the performance of the model in terms of F1-Score and Accuracy is measured and compared to the previous stage before removing those attributes. The threshold selection and clipping process are repeated until the performance of the model drops by omitting a new attribute.

Step 3-Permutation Base Feature Selection- In this step, permutation analysis is performed in conjunction with our Random Forest model to identify and remove attributes with the least contribution to the prediction performance. In simple words, permutation analysis evaluated the contribution of each attribute by randomly shuffling the values of an attribute multiple times and comparing the results to the original values. The level of degradation in the results implies the contribution of that attribute in the results. At the end of this process, a set of performance differences, each for a feature, is generated. Then one attribute with the least level of contribution to the prediction task is omitted, and the performance of the model is compared to the previous stage before removing that attribute. This iterative procedure continues until no further improvement in model performance is observed.

3.4 Analysis of COs' sequential behavior

As mentioned earlier, this study aims to lay the foundation for simplified PCIC modeling for change timing prediction. When dividing a construction project's duration into equal segments, the change orders' performances of projects vary across different periods. The goal of change-timing prediction is to forecast the PCIC in different phases of the project's duration. However, each period's performance may depend on one or more previous period(s). Therefore, it is necessary to understand how COs from different phases are linked to be able to model PCIC and accordingly address the gap in the literature regarding change timing prediction. In this study, a stochastic process known as a Markov Chain is examined for its suitability in modeling the PCIC.

Markov Chain is a simplified branch of stochastic processes that has been extensively used in studying the stochastic behavior of systems under uncertainties. It assumes that the future outcomes of a sequence of events are dependent only on their exact prior step and not the others. This method has been applied to numerous problems to analyze the trends and forecast future outcomes such as credit risk (Yang et al. 2019; Yu et al. 2019), stock market trends (Czech and Wielechowski 2021; Lupu 2015; Mahmoudi and Ghaneei 2022; Wang et al. 2021), pandemic effects (Adekoya et al. 2021; Athari et al. 2023), human geography (Li and Zhang 2009), and medical services (Alanis et al. 2013; Kechagias et al. 2024; Lee and Lee 2018).

Markov Chain models have been also widely employed in the construction industry's literature to address problems in various subjects. In the context of project management and controls Markov chain models are used for invoice processing (Younes et al. 2015), cost performance (Du et al. 2016b), supply chain (Tian et al. 2012), construction simulation (Du et al. 2024), maintenance management (Kim et al. 2015; Scherer and Glagola 1994; Tao et al. 2021a), and tunneling progress and productivity (Touran 1997). In the transportation field, this method is used in driving cycle analysis (Shi et al. 2016; Zhang et al. 2019a; b), the impact of transportation infrastructure on CO₂ emissions (Lu and Du 2024), fault diagnosis of railway pantograph (Ma et al. 2024), and vehicular traffic prediction (Williams n.d.). In structural health monitoring and prediction, this method is employed for bridge deck systems performance prediction (Morcous 2006), and structure displacement evaluation (Qiao and Liu 2012). In Material production, the Markov theory is applied in aggregate production operation (Hajjar and AbouRizk 1998), reinforced bars fatigue prediction (Mantawy and Ravuri 2024), and GFRP fatigue prediction (Thomas and Sobanjo 2013).

An application of the Markov models which has received great attention in the Architecture, Engineering & Construction (AEC) industry for various subjects is performance monitoring and prediction. Markov chain models can capture the temporal dependencies of variables and due to the stochastic nature of many construction processes, this method has been widely used in performance monitoring and prediction papers. Examples of Markovian-based performance prediction and monitoring subjects are bridge structural health (Fang and Sun 2019; Jiang et al. 1988; Morcous 2006; Morcous et al. 2003; Tao et al. 2021b; Yang et al. 2024), pavement health (Abaza 2016, 2017; Butt et al. 1987; Thomas and Sobanjo 2013; Wasiq and Golroo 2024; Wei et al. 2022), project cost (Du et al. 2016b; Leu et al. 2023), pipe health (Micevski et al. 2002; Ossai et al. 2016; Sempewo and Kyokaali 2019), Energy (Hamzehei et al. 2022; Todorov 2024), and infrastructural systems (Kleiner et al. 2006).

For example, (Du et al. 2016b) employed Markov theory to propose an alternative method to traditional earn value analysis (EVA) for cost performance monitoring to overcome the gaps in EVA methods. They mentioned the assumption of static cost performance during project construction as the limitation of EVA which results in a lower prediction of the cost performance. To overcome this gap, a Markovian Cost Performance Monitoring (MCPM) method was proposed to consider the stochastic nature of the problem. The proposed model improved the prediction performance four times more than the model not considering the stochastic behavior.

An Important aspect of Markovian analysis is verifying the validity of the Markov property (Chen and Hong 2012). Testing the Markovian property ensures that a system's past information doesn't provide more knowledge of the process. In other words, given information on the current step of a process, future events can be predicted with the same level of confidence as having access to information on previous steps. Although numerous studies have utilized Markov Chain models to tackle different problems, few of them have investigated the validity of Markov property in their field (Chen and Hong 2012). However, the effectiveness of the Markovian models is under question as long as the Markovian property is not validated. Investigating the Markovian Property helps to understand the underlying temporal dependencies of variables, and accordingly more efficient systems analysis and prediction.

For Example, (Aït-Sahalia 1996; Yang et al. 2019) tested the Markovian property for credit risk and interest rate assessment respectively by checking the validity of the Chapman-Kolmogorov equation. (Morcou 2006) validated the main assumption of bridge performance prediction models which is the dependency of the future deck condition only on the current condition and being independent from past conditions. In this study, the Chapman Kolmogorov equation is used to evaluate the suitability of Markovian processes in capturing the sequential behavior of COs. This method is capable of testing the Markovian assumption by comparing the expected outcomes of a sequence of events under the Markovian assumption with observed outcomes.

Although there are different variations of Markovian processes, in this study, the Markovian property is considered as a dependency of future step (t+1) outcomes to only the state of the current step (t) and not the previous ones. This is also known in the literature as *memoryless* property which indicates the same concept.

Assuming a process with n time steps $t = 1, 2, \dots, n$ and $S(t)$ being the state space at time t, the Markovian property can be formulated in Equation 3-2:

$$P(S_{t+1} = s | S_1, S_2, S_3, \dots, S_t) = P(S_{t+1} = s | S_t) \quad \text{Equation 3-2}$$

considering $p_{ij}(t, t + 1)$ to be the probability of moving from state i in time (t) to state j at time (t+1). Then for each i, j, and t, the 1-step probability of observing stage j at time (t+1), given being at state i in time (t) can be calculated by Equation 3-3:

$$p_{ij}(t, t + 1) = P(S_{t+1} = j | S_t = i) \quad \text{Equation 3-3}$$

In this equation, p represents the transition probability, and i and j belong to the state space. Forming $\mathbf{P}(t,t+1)$ (the transition matrix between time t and its immediate next step) with p_{ij} being the corresponding values of the elements, the n -step transition matrix can be formulated in EQ

$$\mathbf{P}(t, t + n) = \mathbf{P}(t, t + 1) \times \mathbf{P}(t + 1, t + 2) \times \dots \times \mathbf{P}(t + n - 1, t + n) \quad \text{Equation 3-4}$$

Equation 3-4, referred to the Chapman-Kolmogorov equation, is not valid for non-Markov processes. As a result, if the Markovian property exists in the process, the expected n -step transition matrix (\mathbf{P}_e), can be calculated using this equation by multiplying the 1-step transition matrices. The 1-step transition matrices can be formed by looking into the historical data and counting the number of migrations between states in successive steps (t) and ($t+1$). However, another way to calculate the n -step transition matrix is to observe the migrations directly between the step (t) and ($t+n$) and form the n -step observed probability matrix, (\mathbf{P}_o). Then, the Markovian assumption can be validated by comparing the two matrices and ensuring that the difference between the two is not significant. The Chi-Squared test is one of the methods to investigate the similarity of the expected and observed values. By assuming a threshold for the level of significance of the results, (i.e., 0.05) the values higher than the significance level denote the significant similarity of the expected and observed values and as a result, confirm the null hypothesis of validity of Markov property assumption. The 0.05 significance level is a common threshold used by most researchers. It means that there is a 5% chance to reject the null hypothesis while it is actually true.

Chapter 4 – Data understanding and preparation

This chapter is constructed in two sections, (i.e., data understanding, and data preparation) . Section 4.1 presents an overview of the data in terms of the source of data, datasets, contained features, and the challenges to bringing data into the proper shape of this study. This section provides a foundation to better understand the data and the context and identify the required modifications to prepare the data for the purpose of this study. The methods followed to prepare the data and bring it to the proper shape are presented in section 4.2.

4.1 Data understanding

This section is presented in two sub-sections, each presenting a generic task in CRISP-DM (CRoss Industry Standard Process for Data Mining) (Wirth and Hipp 2000) under the data understanding phase; Data description, and Data quality. The CRISP-DM is the widely accepted standard for data analysis to ensure a systematic approach toward the analysis of data. It is divided into six phases each focusing on a specific objective of data mining. This section starts with section 4.1.1 (i.e., Data acquisition) and introduces the data collection process. Then in section 4.1.2; Data Description, a comprehensive overview of the datasets, attributes, terms, and definitions is provided. It outlines the general characteristics and context of the information. The challenges related to information quality, which guide the roadmap for the data preparation phase, are discussed in Section 4.1.3, Data Quality.

4.1.1 Data acquisition

The data used for this research was a group of datasets generated from an industry partner's project management system (i.e., Procure project Management). This firm assists construction companies by providing them with business insights. Nine datasets are collected from this system, each providing information on a separate aspect of construction projects. These projects are collected from the projects executed in Canada. The names of the companies and the original datasets are confidential and cannot be published. Two main datasets are employed for this study, namely *Projects* and *Change Orders*. In addition, an external dataset is integrated to add more context regarding the location of attributes, named Canada Cities. A sample of data from each dataset is presented in appendix 3.

The Procure change management process contains several steps to manage a change during its life cycle. Although the exact process can vary for different types of COs and project types, the overall process is as described in Figure 4-1. It starts with creating a change event from an observed condition affecting the sub-contractor's work scope. Then the subcontractor either requests a quote to collect the required cost information or prepares it. At this stage, a change request would be submitted to the GC (General Contractor) as a commitment change order (i.e., a change that affects the sub-contractor contract). The submitted commitment change order would be analyzed by the responsible team, if the change is in the scope of the original contract (i.e., the contract between the GC and the owner) but not in the scope of the contract between GC and the sub-contractor, the commitment change would be accepted. In the alternative case, when the change is not in the scope of the original contract (i.e., Prime contract) the GC submits a separate change request and creates a potential prime contract change order, and upon approval of the prime contract change order, approves the commitment change order as well.

Two distinct datasets, *change events*, and *change orders* have logged this process. *Change events* dataset logs the lifecycle of a CO from when being recognized till getting final approval.

The *change orders* dataset keeps the information on accepted or rejected COs. Since the scope of this study is the prediction of the approved change orders, only the *change orders* data frame is employed for this study (“Request for Change Orders” 2024).

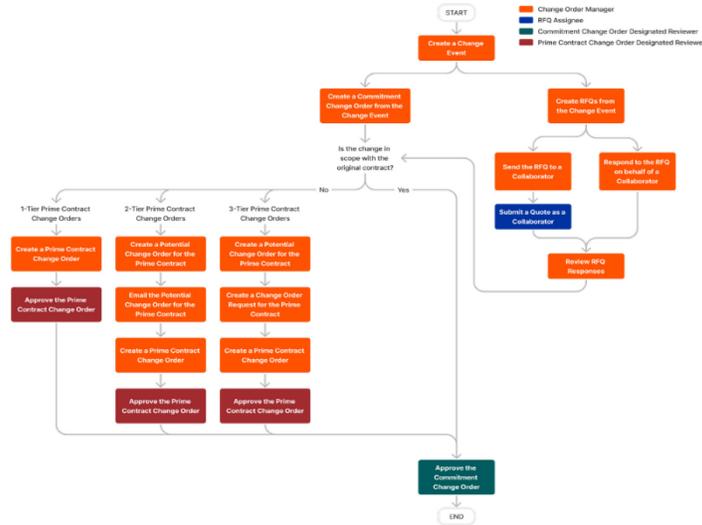


Figure 4-1: Procure change management process (“Change Events - Workflow Diagrams” 2017)

4.1.2 Data description

This study uses data from real construction projects in Canada, logged between January 2017 and August 2022 provided by a prime contractor. Data from different aspects of projects are provided, however, two main datasets, i.e., *Projects* and *Change Orders*, are being used in this research. *Change Orders* dataset provides details on each instance of COs, mainly the cost impact, change description, and implementation time. The *Projects* dataset describes projects by providing information such as project value, duration, location, and type. Moreover, to enrich the data with more contextual specifications of projects’ locations, an external dataset, i.e., *Canada Cities* (“canadacities.csv” 2023) is integrated which provides information on the population and population density of projects. Population and population density of the location of projects are considered, as they are known as being associated with various change reasons such as labor shortage, site access issues, residential requirements, and weather conditions. Table 4-1 summarizes the attributes, definitions, classes, and data distribution of these datasets. In the following sections, these data sets are described briefly.

Projects- The *Projects* keeps the information on the characteristics of the projects. The original dataset contains 3,065 data points. Each data point represents an individual contract, and its details are logged in attributes including Expected Start Date, Expected End Date, Base Contract Value, City, Province, Type, Operating Unit, Billing Type, Project Classification, Project Parent ID, and Project ID. The Type identifies if a project is a construction or service project. The Operating Unit is an inter-organizational project categorization attribute that takes six values and indicates the importance of the project. The project classification has two levels of information in the values assigned to the projects (i.e., Industry group and project type). This attribute is explained more in detail in section 4.2.3. The project ID and parent ID identify the unique IDs assigned to projects and their higher-level parent projects.

Change_Order- *Change orders* contain approximately 9M data points, divided into nine data frames (tables). Each data point is a change order at a certain stage. *Change Orders* not only provides information on each accepted or denied CO but also logs them at different stages in their lifecycle. A CO might be logged several times at different stages. Attributes such as Date Created, Amount, Status, and Project ID are the attributes that identify COs at different stages of their lifecycle and are related to this study.

Canada Cities- The *Projects* data frame which keeps the information on project properties, has two attributes *Project City*, and *Project Province*. To consider the effect of project location on the occurrence of COs from a demographic point of view, an external dataset, i.e., *Canada Cities* is linked with the available datasets that is merged with the final dataset before modeling. Two attributes are added to the dataset as a result of integrating this dataset, i.e., *Population*, and *Population Density*. These two attributes are integrated to add demographic context to the location of projects. According to (Padala et al. 2020), environmental and social causes of COs include noise pollution, the opposition of neighboring communities, labor shortage, and residential requirements which can be influenced by the demographic properties of projects' location.

Table 4-1: The raw attributes exist in the Datasets

Dataset	Attribute Name	Type	Values Description
Projects	Base Contract Value	Numerical	Between 1,095 - 36,000,000 CAD
	Duration	Numerical	Between 1 – 465 Days
	Province	Categorical	4 Provinces (QC, ON, AB, BC)
	City	Categorical	Include 134 Cities
	Type	Categorical	Construction / Services
	Classification	Categorical	34 Classes including Residential, Plant, Military, Infrastructure, etc.
	Operating Unit	Categorical	6 Classes including Controls, Special, Major, Corporate, etc.
Change Orders	Amount	Numerical	Change Order Amount
	Date	Date	Change Order Date
	Description	String	Change Order Explanation
Canada Cities	City	Categorical	Name of City
	Population	Numerical	Population of City
	Density	Numerical	0 - 5492 person / Sq-ft

4.1.3 Data quality

Working with real data is always challenging, and this research is not an exception. Following the CRISP-DM, this section presents the challenges in adapting the data with the required format including Descriptive Data, Duplicates, and Couples that form the foundations for the data preparation section.

Descriptive Data – *Change orders* contain an attribute that provides descriptive context (i.e., Change Description) of each individual change order. Textual data must be converted to a machine-readable format before it can be utilized. The worst-case scenario occurs when no standard format is followed for logging the data, particularly when multiple loggers are involved. Looking into Table 4-2, it is identical that the values of the ‘Change Description’ attribute needs

cleansing to get rid of noises before being used. Given the importance of this attribute to this research, a significant portion of time is dedicated to preparing and standardizing it.

Table 4-2: Sample data points with noise sources

	Change Description	Change Amount	Change Date
1	<p>Construction Contract for Fire Protection at Project Name.</p> <p> </p>	3,330,000 \$	2022/05/16
2	<p>Construction Contract for Fire at Project Name.</p> <p> </p>	-3,330,000 \$	2022/05/16
3	<p>Construction Contract for Fire Protection at Project Name.</p> <p> </p>	3,330,000 \$	2022/05/16
4	<p>Construction Contract for Fire at Project Name.</p> <p> </p>	3,330,000 \$	2022/05/16

Duplicates – One of the first and most significant issues faced was the vast number of duplicated COs in the dataset. Upon exploring the dataset, it became apparent that the number of COs per project did not align with what has been mentioned in the literature. According to (Anastasopoulos et al. 2010), although the number of COs varies from one project to another, construction projects experience five COs on average, with a maximum of 55. However, the provided dataset contains more than 2,900 change orders per project. A deeper investigation revealed that not all the change orders were unique and relevant. Numerous data points with the same or similar values in the *Amount*, *Date Created*, and *Description* were logged to capture the events within the change orders life cycle. As an example, the change orders shown in rows 1, 3, and 4 all represent one CO, and are Duplicates that need to be dealt with.

Couples – In this study, ‘Couples’ are defined as two change orders logged in the system to neutralize the effect of a canceled change order. These groups of COs share the same, or similar values in the Description, and Date Created, but the values of the *Amount* have an opposite sign (i.e., Row 2 of Table 4-2 is the couple for the other three). Couples are created by the logger mostly to correct or remove a pre-logged change order during or after its creation. The clue to this finding is found in the description of these entities which was mostly describing a correction in the change orders, omitting a change order, or transferring a change order from one contract to another one.

4.2 Data preparation

In this section, the methods followed to overcome the challenges and pre-process the data are shown including data cleaning techniques, filtering projects and COs, attribute construction, aggregation of datasets, and dealing with missing values.

4.2.1 Data cleaning

The initial raw datasets contained various sources of noise that required thorough cleaning before analysis. The noise in the *Change Orders* fell into two main categories, (i) noise in

unstructured textual information in the description of COs; and (ii) Duplicates and Couples. While duplicates are the result of multiple logging (due to mistakes, or revisions), couples are usually the results of canceling a CO (for various reasons). When a CO is canceled, sometimes the operator introduces a dummy CO of the same value and opposite sign to neutralize the cost effect of canceled COs that had been already logged in the database.

Unstructured change order description – The attribute 'Change Description' in the *Change Orders* provides an explanation of the source of COs in a textual format and is found to be greatly noisy. Several steps are followed to bring the values of this attribute into a structured format for further analysis. Tokenization, lower casing, stop word removal, stemming, and lemmatization, are the five common steps in most text data preparation practices and are employed to clean this attribute. Tokenization can be simply described as the process of breaking sentences into words; lower casing brings data to a standard format; stop word removal aims to eliminate words not adding a semantic value; and stemming and lemmatization are the processes of bringing words to their root format.

Couples and Duplicates shall be removed from the dataset as they fall outside the focus of this study. A two-step algorithm is implemented for this purpose. Initially, for each instance of COs, the algorithm searches for co-occurrences of tokens (words) in descriptions of the other COs issued for the same project. The Jaccard similarity is used for evaluating the similarity among descriptions. Jaccard similarity is a statistic used to quantify the similarity between two groups. It calculates the size of the intersection divided by the size of the union of the two groups. In this study, this term is used to measure the similarity between descriptions of COs and accordingly identify and group the same COs. To this aim, first, the descriptions are converted to word tokens, and each CO is compared to other COs in the same project. Then the number of common word tokens is divided by the number of tokens to measure the similarity of the two COs. Subsequently, by applying the rules shown in Table 4-3, duplicates and couples are identified and removed from the dataset.

Table 4-3: Rules for identification of couples and duplicates

Category	Date of Issue	Amount	Measure of similarity
<i>Duplicates</i>	Same	Same or less than 10% difference	Jaccard Similarity > 0.8
<i>Couples</i>	Same	Same with a negative sign	Jaccard Similarity > 0.8

4.2.2 Record selection

Having the noise removed, the two provided data frames are investigated for the alignment of data points with the scope of this study. Various types of problematic records are identified from each dataset, and those instances and their corresponding records in the other dataset are removed. In this section, each type is explained along with the reason and methods used. The exclusion criteria encompassed the following projects and COs.

Change orders with no project data – Since the main goal of this research is to predict the magnitude of COs by using project information available in the planning phase, the availability of this information is critical. These pieces of information are recorded in the *Projects*. However, for some COs, the corresponding project information is not logged. Tracking each CO to its project

can be done through a key (an attribute in common among both datasets) named Project ID, and accordingly, those records for which the corresponding Project ID is not recorded in the *Projects*, are removed from the *Change Orders*.

'Child-Parent relationship' for projects and their corresponding COs – From the viewpoint of cross-project dependencies, each project can be classified into a 'Parent' or 'Child' class. The two types are logged and specified in the *Projects* through an attribute named Parent Project ID. An analysis of the COs of Parent projects and their Child Projects was performed to find relations between the COs of sub-jobs and parent projects. The results revealed that the same COs are logged for all the sub-jobs of a parent project as well as the parent project itself. Having several projects with same amount and number of COs affects the classifier's performance. In such cases, if the classifier is fed by one of those projects (either the parent or its sub-jobs) during the training phase, it will simply memorize that data point instead of learning the pattern, and as a result, it will correctly predict all other instances from the same family. Although such a classifier will result in synthetically higher performance metrics, those data points are removed since no additional value exists in them.

Change orders of service projects – Two categories of projects exist in the data, i.e., 'Construction' and 'Service' jobs. As these two types of projects have a completely different nature, and a great portion of service projects have a short duration without any COs, the *Change Orders* was cleansed to be limited to COs of construction projects.

Zero contract values – As discussed previously, project properties logged in the *Projects* are the key features to predict COs. To understand the importance of available predictors, different feature selection techniques are performed which are explained in detail in the modeling section. The result of feature selection showed that the Contract Value is one of the most correlated attributes to the magnitude of COs. Hence, missing values of this attribute are not substituted, and accordingly, projects with missing contract values are removed. Moreover, some projects with zero contract value exist in the dataset which are related to performed jobs under warranty responsibilities. Those are also removed from the dataset because of their different nature (i.e., Remaining responsibilities of past projects rather than representing a new contract).

Open change orders – The *Change Orders* contains information on each instance of COs after receiving the initial approvals. However, some COs were not approved at the time of data collection by this study. Those 'open' COs are also excluded from the dataset.

Projects with partial change data availability – The *Change Orders* dataset contains change orders issued between January 2017 and August 2022. This means that projects with start or end dates out of this period may have COs that are not logged in the dataset. 202 projects fall in this category; however, the missing period issue is not critical for all of them at the same level. Some projects may miss only one day of the duration, while others might miss more than half of their duration. After plotting the distribution of projects with missing change data versus the number of missing days of change data as a percentage of project duration, a decision was made to keep the projects with missing change data equal to or less than 25% of their duration. This approach resulted in keeping in almost 80 projects out of 202 while maintaining the reliability of data (by not missing a great period of data). Figure 4-2 shows histograms of the missing portion of project duration for those 202 projects, from the lower bound (i.e., Projects started before January 2017) in Figure 4-2(a), and from the upper bound (i.e., Projects ended after August 2022) in Figure 4-2(b). The horizontal axis shows, the portion of projects' durations falling outside the data

availability period while having the number of projects with the same portion of missing data on the vertical axis.

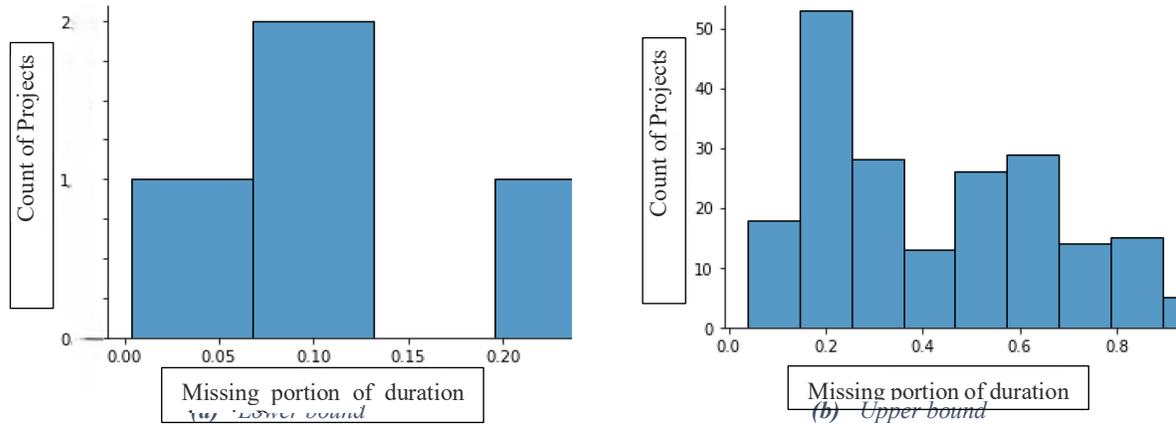


Figure 4-2: Distribution of missing segment of project duration in change orders data

4.2.3 Data aggregation and integration

As mentioned, each record in the *Change Orders* is representative of one single change order; however, prediction of individual COs is not the aim of this study, as it depends on several factors and is extremely complex, if ever possible, to predict it from the project attributes alone. Alternatively, this study tends to predict the severity of Cost Impact of Change orders (CIC) as a classification problem. As a result, after aggregating COs on the project level, a new data frame replaced the original *Change Orders*, containing four attributes showing the Project ID, CIC as percentage of contract value, Number of Additive COs, and Number of Deductive COs. For the 1st component of this study, the CIC is calculated by directly summing up the values of all COs issued in a project and normalizing it by the contract value. For the second component, a similar term is calculated for each period of the project individually, named Periodic Cost Impact of COs (PCIC). To calculate PCIC, first, the duration of a project is divided into 5 equal phases. Then, the PCIC of each phase is calculated in two ways regarding the change impact aggregation. Once Cumulatively by aggregating COs of a phase and all previous COs as the percentage of contract value, this term is named PCIC-C. Second, singularly by calculating the PCIC for each period separately and named PCIC-S (both terms are calculated as the percentage of contract value).

The Number of ‘Deductive’, and ‘Additive’ COs show the number of change orders with a negative or positive value sign, respectively. This term is gathered to transform categorical attributes into semantic numeric values that show the expected frequency of COs in a certain group of projects. This step is further described in the HCPI attributes section. Afterward, the three fragmented datasets are merged and flattened to form one single data frame containing all the information needed in one place (*Projects* and *Change Orders* through the Project ID, and the *Canada Cities* through the name of cities).

Auxiliary Attributes – Two new groups of attributes are constructed and added to the dataset. First, using the change domain knowledge a set of additional attributes is introduced to capture a valuable piece of information getting lost during data preprocessing and developing a term to quantify varied change performances of distinct groups of projects. It is known that the magnitude of COs varies by project type and geographical location of projects (Cantarelli et al. 2012; Love et al. 2017). With this in mind, a set of attributes is proposed to quantify varied change performances of different groups of projects. Second, an existing attribute in the *Projects*, named

Project Classification, is divided into two distinct attributes to capture a higher level of information compared to the original raw feature.

Historical Change Performance Indicator (HCPI) – The main idea behind this set of auxiliary attributes is inspired from a data preparation method known as ‘Ordinal Coding’. This method converts categorical attributes to numeric format to quantify the differences between different categories of nominal attributes. It works based on assigning each category to semantic numerical values to represent the existing ranking/order in their nature, if exist. For instance, categories such as 'high,' 'medium,' and 'low' possess a ranking by default and can be encoded as integers 3, 2, and 1, respectively. The three classes of high, medium, and low are treated as distinct categories by the machine; however, when represented numerically, the machine can more effectively interpret the ordinal relationships between these classes. This numerical representation allows the machine to recognize the relative magnitude of each class, enhancing its ability to process and analyze the data. On the contrary, the categories of an attribute showing the name of cities, are of no ordinal nature and hence cannot be meaningfully assigned to numerical values.

The categorical features used in this study derive from the second scenario, where there is no inherent semantic relationship among the categories. The overarching goal of this set of proposed attributes (HCPI attributes) is to impart an ordinal nature to these categories. Consequently, for each categorical attribute, four new features are created. These new attributes replace the original categories with the average frequency of change orders (COs) in projects within that specific category. Two of the new features substitute the original categorical values with the average number of positive and negative COs per unit of project duration for projects belongs to that group. The other two features replace the original values with the number of positive and negative change orders per unit of contract value. For example, the Operating Unit column includes a category named "Electrical," which represents electrical projects. For this attribute, four new features are constructed, replacing the word "Electrical" with the average frequency of COs in projects belonging to Electrical group.

The number of change orders, which was also neglected in the previous prediction studies, is chosen due to its significant correlation with the total dollar value of change orders of projects. The Pearson correlation coefficient for these two attributes is 0.78 showing a significant direct correlation between the two terms, and accordingly, the change frequency can be a good representative of COs performance of projects which is also along with the past studies (Shrestha et al. 2023) . The number of COs cannot be feasibly used as a predictor of the target attribute due to not being known in the planning phase before the project starts. However, the number of COs in similar past projects is known through historical data and can be used as an indicator to quantify the ‘expected’ number of COs in a group of projects. It can be hence identified and used as a predictor for projects in the same category. The values calculated for each category, when compared to other categories, show the difference in the expected number of change orders or change performance between projects in those groups.

There are various ways to define the similarity criterion for grouping similar projects, (e.g., a combination of project type, and a range of contract values). In this study, the similarity criterion is defined at the simplest level, which is the categories available in the categorical attributes. This means projects in the same group of categorical attributes such as ‘Casino’ in the ‘Building Type’ attribute are considered as one group. To this aim, two temporary columns are created; one takes the number of COs with positive cost impact, and the other takes the number of COs with negative

cost impact as values. Since the frequency of COs is reported as highly dependent on the size and duration of projects (Anastasopoulos et al. 2010), it is needed to neutralize the effect of the project size (in terms of duration and contract value). As a result, from the two temporary columns (the number of COs with negative and positive cost impact) two new columns are constructed; one normalizes the values of that column based on the project duration, and another one by the project's contract values. In short, these four columns at this stage describe the number of COs with positive/negative cost impact per unit contract value/duration of each project. Then for each of the four column, the average of values for projects in same group (groups are defined from the values of the initial nominal column) is considered as average frequency of COs in that group of projects. This value is now the representation of the average negative/positive value change frequency per unit of project value/duration. The four values calculated for each group of project are then assigned to all projects of same group in four separate columns.

The naming convention for these auxiliary attributes includes a first letter (P/N, referring to the Positive or Negative cost impact), followed by DU/CV (standing for normalized by DURATION or Contract Value). For example, HCPI-P-CV shows the frequency of COs with positive cost impact, and the values are normalized by contract value.

Breaking down the 'Project Classification' attribute – An attribute of the *Change Orders* dataset, named *Project Classification* contains two levels of information: (i) the construction type, such as residential, commercial, and industrial; and (ii) the building use, e.g. hospital, office, or casino. To capture interrelationships within project types and industry groups, this attribute is separated into two distinct columns, named Industry Group, and Project Type.

4.2.4 Handling missing data

Missing values mainly coming from the Project dataset are treated separately for each attribute, based on their nature. Several common methods of dealing with missing values were tested, including replacing missing values with in-class (or the whole dataset's) statistics such as mean, median, or mode; random sampling; imputation through regression; multiple imputations; or eliminating the observations with missing values. It is essential to choose the right method (based on the type and nature of data and missing values, the number of missing values, and the analysis's main objective) to ensure accurate results. As a result, missing values of each attribute are dealt with separately to minimize the risk of harming the data. The adopted strategy for each feature is explained in Table 4-4.

Table 4-4: Strategies to deal with the missing values

Dataset	Attribute	Method
Projects	Commercial Group (Classification 1)	Replaced by the mode of the same attribute for records in the same province
	Project Type (Classification 2)	Replaced by mode of the same attribute for records in the same Commercial Group
	Billing Type	Replaced by the mode of all records
	Project Importance (Operating Unit)	Mode of the same attribute for records with the same value in the Project Classification attribute
	Project Province	The city's corresponding Province matched automatically
	Project City	Replaced with the mode of records in the same province
	Duration	Imputing missing values using a regression model
	Contract Value	Records with missing values were omitted.

Canadian Cities	City	Missing Cities are added from data recorded on the Canada Census portal
-----------------	------	---

It must be noted that although both numerical attributes of the Project dataset (i.e., Duration and Contract Value) contain missing values, it was decided not to impute both. As the target attribute is labeled based on the cumulative amount of COs, divided by the project’s contract value, imputing the Contract Value would affect the target attribute. Since the main scope of this study is the planning phase, and as in such phases the estimated contract value is usually more accurate than the estimated duration, imputing the duration of projects is expected to have less effect on the results. As a result, the projects with missing values, or values equal to zero in the contract value are omitted and missing values of duration are imputed.

To impute the missing values of the Duration attribute, several algorithms including Artificial Neural Network (ANN), Support Vector Machine Regressor (SVR), Random Forest (RF) Regressor, KNN Regressor, and Multi Linear Regressor (MLR) are utilized to perform the imputation task by applying 15 folds of cross-validation. Among the models tested, the RF Regressor was best able to impute the missing values, achieving an R-squared of 80%, and as a result, missing values are imputed using this model. To impute the missing values, first, the dataset is divided into projects with missing values in the Duration attribute (test set) and those without missing values (training set). Then, the other attributes in the *Project* dataset are used as predictors to train the model on the training set and accordingly use the model to predict the missing values in the test set.

4.2.5 Data labeling, encoding, and standardization

Since this study approaches the COs prediction as a classification problem, the target attribute (i.e., change percentage) must be categorical and the data needs to be labeled. For the first component, i.e., change severity prediction, projects are labeled in two classes, i.e., *high* and *low*. Projects with CIC exceeding 4% of the contract value are considered class high, and the others are labeled as class low, i.e., 68% of projects belong to class low and 32% to class high. For the second component, i.e., analysis of the sequential behavior of COs, projects are labeled into three classes as *high*, *medium*, and *low*. The 10% and 4% are considered thresholds for defining the three classes, and as a result, 24% of data belongs to class high, 8% to class medium, and 68% to class low. For both components, the thresholds can be adjusted based on the expectations of each organization and how they define high, medium, and low.

Since the algorithms used in this study only work with numerical features, the categorical attributes are encoded to numerical values using integer coding. All attributes are then standardized for Artificial Neural Network (ANN), Random Forest (RF), and Principal Component Analysis (PCA) algorithms. Standardization brings the numerical attributes into the same scale with a mean of zero and a unit variance to prevent those with higher scales from dominating the others. The standardization transformation is conducted using Equation 4-1. In this equation, X is the raw value of each data point, μ is the mean of values in that attribute, and σ is the standard deviation of the attribute.

$$z = \frac{X - \mu}{\sigma} \quad \text{Equation 4-1}$$

After performing all the data preparation steps, the size of data after performing the data preparation steps is reduced from 8.8M data points to 52K for the *change orders* dataset, and from 3,066 to 1029 for the *Projects* data frame.

Chapter 5 – Change Orders Cost Impact Prediction

This chapter outlines details of applying the methods, and the results of the first component of this study, which aims to investigate the predictability of the Cost Impact of Change Orders (CIC) and explore ways to enhance predictability. Additionally, by evaluating the significance of the attributes used as predictors, the features correlated with the CIC are identified.

First, as described in section 5.1, an initial feature selection and modeling are conducted to choose the best model. Several predictive algorithms are trained and tested to compare their performance and select one for the next steps. Second, as explained in section 5.2, the importance of attributes is examined in greater detail by a multi-step feature selection technique.

5.1 Modelling

Before training models and selecting the most effective one, it is necessary to provide the models with a proper list of attributes. The dataset contains two types of attributes based on their value types, i.e., categorical and numerical. The Principal Component Analysis (PCA) algorithm is used to investigate the importance of numerical attributes by examining the amount of variation they explain. For the categorical attributes, the Chi-squared test is employed to assess their significance. Figure 5-1 shows the results of the PCA analysis, illustrating the cumulative level of variation captured by different principal components (PCs) for the numerical attributes.

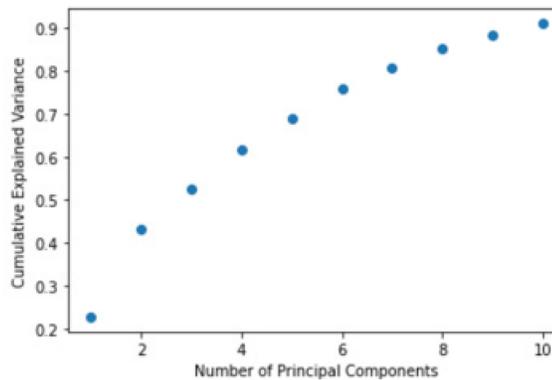


Figure 5-1: PCA cumulative variation per PC

It is assumed that, for the initial model training, a 70% variation is enough to select the most influential PCs and accordingly identify the most correlated attributes. It is identical to Figure 5-1 that the first five PCs explain 70% variation in the data. As a result, the most dominant attributes in these PCs are selected as the most correlated numerical attributes to train the models.

Looking into Figure 5-1, it is identical that, the first Principal Component (PC1) accounts for the highest level of variation in the dataset, explaining approximately 25% of the total variance. The subsequent PCs gradually capture decreasing levels of variance, with PC2 explaining approximately 20% of the variation, PC3 explaining around 10%, and so on. Collectively, the first five PCs account for over 70% of the total variance, suggesting that the majority of the information can be retained by considering the first five components. The explained variance ratio also shows that the first three PCs are the most influential, contributing to more than 55% of the total variance.

Table 5-1 shows the list of dominant attributes of the first five PCs excluding the ones already selected in previous PCs.

Table 5-1: Dominant attributes in the first five PCs

PC-1	PC-2	PC-3	PC-4	PC-5
HCPI-Industry-P-DU	HCPI-OperatingUnit-P-DU	Population	HCPI-BuildingType-N-DU	HCPI-Classification-P-DU
HCPI-Province-P-DU	HCPI-BuildingType-N-DU	Density	HCPI-Industry-N-DU	
HCPI-City-N-CV	HCPI-Classification-P-DU	HCPI-City-P-DU		
HCPI-Province-N-DU		Contract Value		
		Duration		

Table 5-2 shows the results of the Chi-Squared test coupled with their corresponding p-values for the categorical attributes. Attributes with both p-values less than 0.05 and a high Chi-squared value are considered to be correlated attributes to the target attribute. As identical to the P-values, in Table 5-2, City, Industry Group, Province, and Billing Type do not show a significant correlation to the target attribute. Looking into the Chi-Squared values, Province, Industry Group, and Billing Type do not show an efficient contribution to the prediction. As a result, the Province, Industry Group, and billing Type are omitted from the list of attributes for the initial model evaluation. However, the importance of these attributes is analyzed again after choosing the best model.

Table 5-2: Chi-squared test results

Feature	CHI-VALUE	P-VALUE
City	101.975	0.482
Classification	63.069	0.001
Building Type	55.415	0.001
Importance	17.563	0.004
Industry Group	9.343	0.053
Province	5.831	0.212
Billing Type	4.142	0.042

Then several predictive algorithms (i.e., ANN, RF, SVM-Linear, SVM-Poly, and SVM-RBF) are developed and tested to identify the most effective model. Table 5-3 shows the performance of the algorithms developed using the selected list of attributes. To address the class imbalance problem and avoid overfitting, the models are tuned to achieve the highest possible F1-Score on the test set. Among the others, the RF model showed better fitness to the data by well-addressing the class in balance achieving 61% F1-Score and 70% accuracy on the test set. As a result, this model is chosen for the next steps.

Table 5-3: Models' performance comparison

	Accuracy Train	Accuracy Test	F1-Score Train	F1-Score Test
SVM-Linear	0.70	0.62	0.59	0.54
SVM-Poly	0.75	0.69	0.65	0.59
SVM-RBF	0.70	0.64	0.68	0.62
RF	0.79	0.70	0.70	0.61
ANN	0.79	0.75	0.64	0.50

5.2 Multi-step feature selection

To evaluate the importance of attributes in more detail, a multi-step feature selection method is proposed to check the features from different perspectives. First, as explained in section 5.2.1, attributes are evaluated on the attributes type level to check if introducing certain groups of attributes influences the prediction performance. Second, as described in section 5.2.2, to better understand the significance of attributes and reduce the effect of collinearity between predictors, a dendrogram is utilized to remove highly correlated attributes by considering their prediction contribution. Third, in section 5.2.3 another step of feature selection is followed to investigate the importance of attributes based on permutation analysis.

5.2.1 Attribute type evaluation

To evaluate the usefulness of attributes on the attribute type level, the features are categorized into five categories as described in Table 5-4. Then the RF model is trained and tuned several times each with a subset of attributes to check if adding a certain group significantly affects the model performance. Table 5-5 shows the performance of the RF model in different scenarios each with a subset of attribute groups. *Project Size* consists of the Duration and Contract Value of projects to capture the importance of the size of projects. According to (Choi et al. 2016), occurrences of COs are highly dependent on these two attributes (i.e., Project Duration, Contract Value). *Project Type* encompasses *building type*, *industry group*, *importance*, and *classification* to describe variations in project characteristics. *Location* includes *City* and *Province* to capture specific properties of project locations. The *Project Type* and *Location* attributes are included as projects of different types and locations experience varying levels of cost variation (Cantarelli et al. 2012; Love et al. 2017).

The last two groups of features (i.e., Demographic Properties, and HCPI features) were not among the existing attributes of the dataset. The Demographic Properties (DP) group contains *Population* and *Population Density*. These attributes are included to provide demographic context to the project locations which can influence the probability of change occurrences. According to (Padala et al. 2020), environmental and social causes of COs include noise pollution, opposition from neighboring communities, labor shortages, and residential requirements. These attributes can be linked to demographic specifications of the location of projects. The *HCPI* attributes, which are explained in detail in section 4.2.3 tend to add another layer of context to the data using the change domain knowledge. According to (Shrestha et al. 2023), the frequency of change orders is directly correlated with the severity of COs, meaning that more change orders result in a higher change percentage. In addition, as mentioned earlier, the magnitude and number of COs may vary based on the type and location of projects. This set of attributes adds a layer of information regarding the average frequency of COs based on project type and location to restore the linkage between the change severity and frequency of COs. For each existing categorical attribute (i.e., categorical features define the type of projects and their spatial difference) four new attributes are constructed to quantify the average change frequency of projects belonging to the same category. These attributes are constructed to capture both the frequency of positive and negative value COs. Since the size of a project (i.e., duration and project value) is in correlation with the number of change orders in a project, the number of negative, and positive value COs are normalized once per duration and another time with the value of projects. As a result, the four attributes show the expected number of positive/negative value COs in units of project duration, or value. The four types of attributes are shown in Table 5-4, “ATR” in their naming structure is the placeholder for the name of the main categorical attribute it is created from.

Table 5-4: Attribute groups and contained features

Project Size	Project Type	Location	DP	HCPI
Duration	Building Type	City	Population	HCPI-ATR-P-DU
Contract Value	Classification	Province	Population Density	HCPI-ATR-N-DU
	Industry Group			HCPI-ATR-P-SZ
	Importance			HCPI-ATR-N-SZ

It is evident from Table 5-5 that the added context due to the integration of HCPI attributes is effectively correlated with the target attribute. This is equivalent to almost 8% prediction performance improvement in both terms of F1-score and Accuracy on the test sets. Moreover, the DP attributes better describe the specific characteristics of project locations, rather than the name of the city, or province by showing a 2% prediction performance improvement.

Table 5-5: Attribute-type level feature importance analysis

	Attribute Types					Model parameters					Performance Indicators			
	Project Size	Project Type	Location	DP	HCPI	Max depth	Min sample leaf	Min sample split	Number of trees	criteria	Acc Train	Acc Test	F1 Train	F1 Test
1	1	1	0	0	0	7	7	15	400	Gini	0.65	0.58	0.56	0.52
2	1	1	1	0	0	9	7	15	400	Gini	0.71	0.64	0.61	0.54
3	1	1	1	1	0	8	7	15	400	Gini	0.77	0.63	0.68	0.56
4	1	1	1	1	1	8	7	15	400	Gini	0.79	0.70	0.70	0.61
5	1	1	0	1	1	8	7	15	1000	Gini	0.79	0.73	0.71	0.63
6	1	1	1	0	1	8	7	15	400	Gini	0.78	0.73	0.70	0.63
7	1	1	0	0	1	8	7	15	400	Gini	0.77	0.70	0.67	0.61

5.2.2 Collinearity removal

The entropy results generated from the RF model are employed to investigate the importance of attributes. When being used for feature importance analysis, entropy is vulnerable to multicollinearity (Soofi 1990). Since the level of contribution of each attribute in prediction of the target feature is distributed among several correlated attributes, the dominance of each individual feature decreases. To avoid this issue, it is important to check the predictors for collinearity and only keep the features that are most correlated with the target attribute. To check the multicollinearity of attributes a hierarchical clustering method is utilized on the Spearman rank-order correlations. Then, with the aid of a dendrogram, the proper threshold for clipping is identified to form the most correlated cluster of attributes. From attributes collected in each cluster, only the one with the highest permutation level remains.

Figure 5-2 presents the attributes correlated to each other along with their level of similarity based on the distance between them. The link between attributes defines the level of difference between a group of correlated attributes. As a result, the lower levels of linkage show higher similarity or correlation of attributes. For example, *Duration*, and *Contract Value* form a correlation cluster at a distance equal to 0.4. However, the same attributes form another cluster at a distance equal to 1.1 with *Bill Type*. A lower distance connection shows higher levels of correlation. In this step of feature selection, at each iteration, thresholds are picked based on the lowest level of observable dissimilarity and from each cluster with a level of difference less than the threshold only one feature with the highest entropy is kept and the others are removed. Looking at Figure 5-3 which shows the performance of the model after each iteration, it is observable that the prediction performance improves continuously in terms of F1-Score and Accuracy in the first four iterations from 0.726 to 0.74 and 0.636 to 0.657, relatively and drops sharply after the fifth iteration. As a result, the first step of the feature selection process is cut at the fourth round.

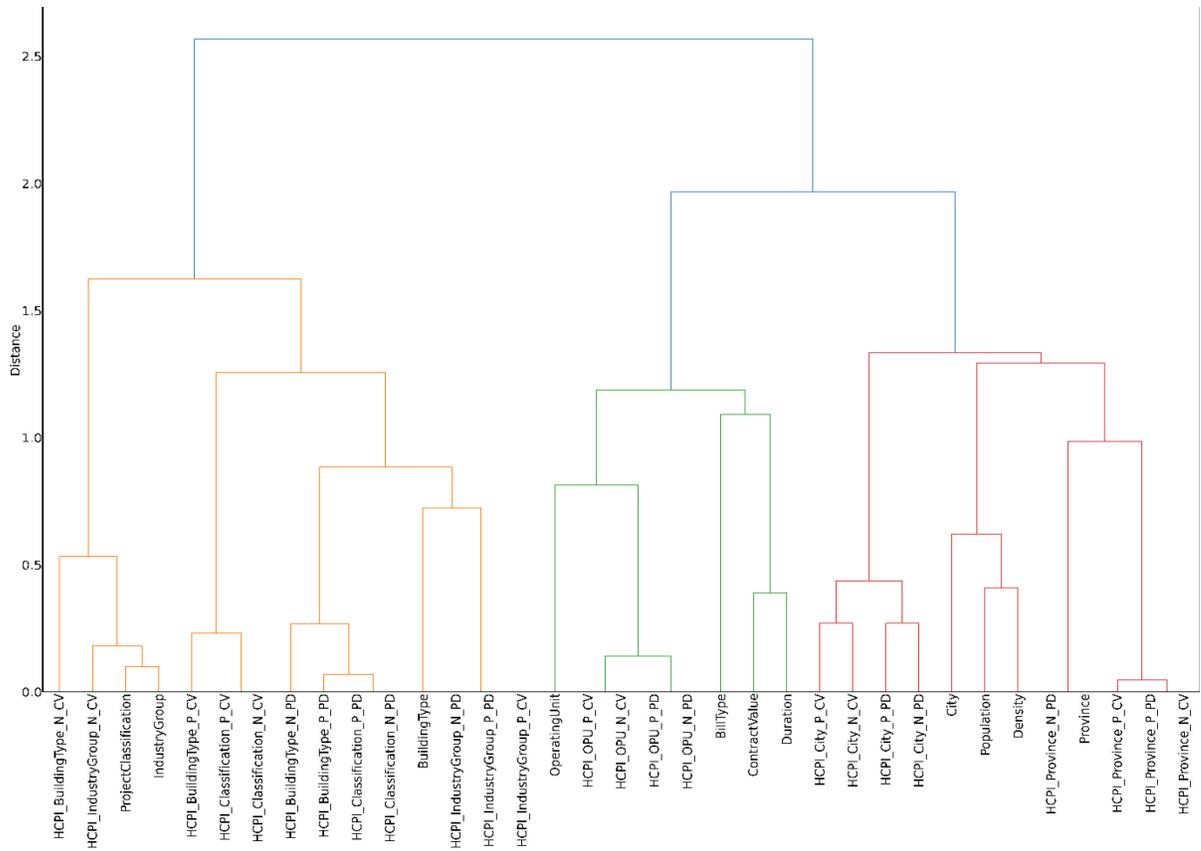


Figure 5-2: Level of correlation between predictors

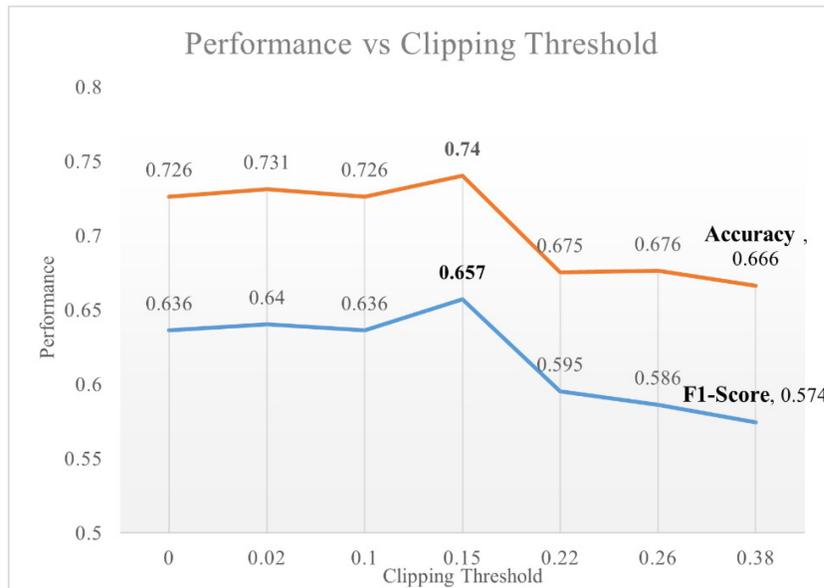
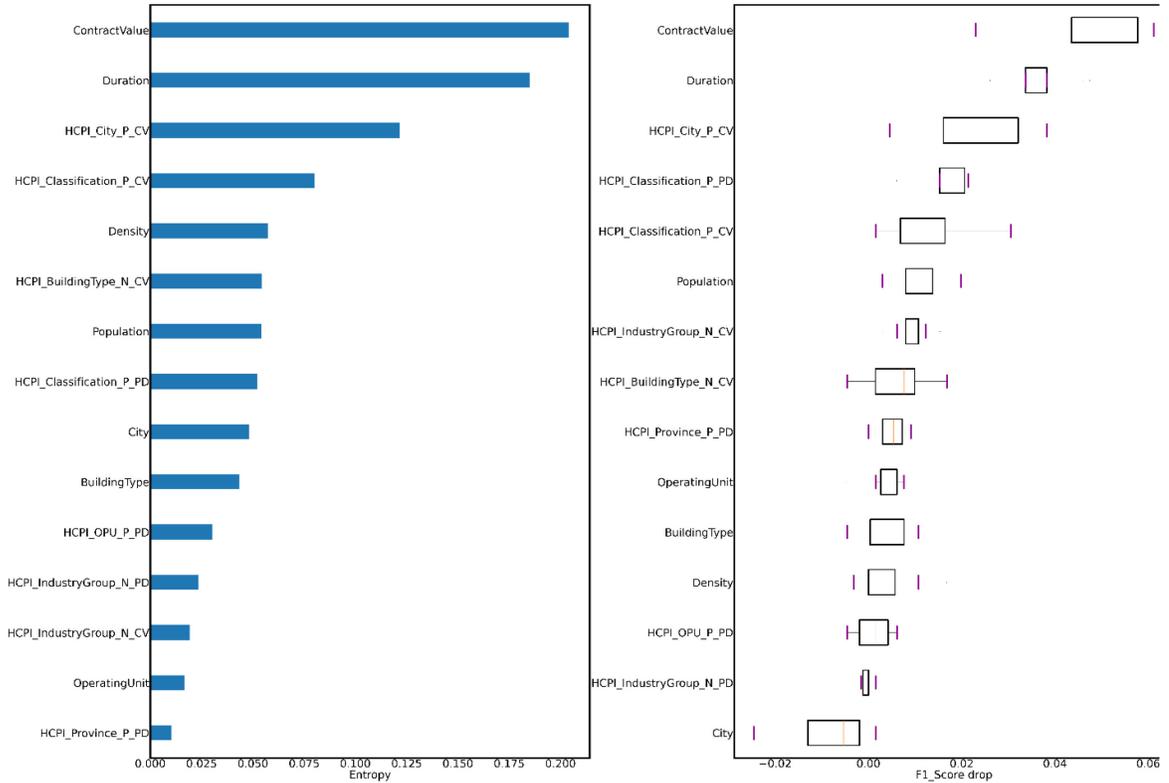


Figure 5-3: Performance improvement vs clipping threshold

5.2.3 Permutation-based feature selection

In the next step, to further check the attribute's importance, the permutation analysis is employed on the remaining attributes using the RF model and F1-Score as the performance criterion. First, attributes are rank order based on their level of contribution generated from permutation analysis, as shown in Figure 5-4. Then at each step, one feature with the lowest importance is removed from the dataset and the model is trained on the remaining subset of attributes. This process continues till further omission of attributes reduces the model performance. The remaining attributes are then identified as the attributes with the highest correlation to the target attribute.



(a)-Entropy result

(b)-Permutation Result

Figure 5-4: Permutation and Entropy Results

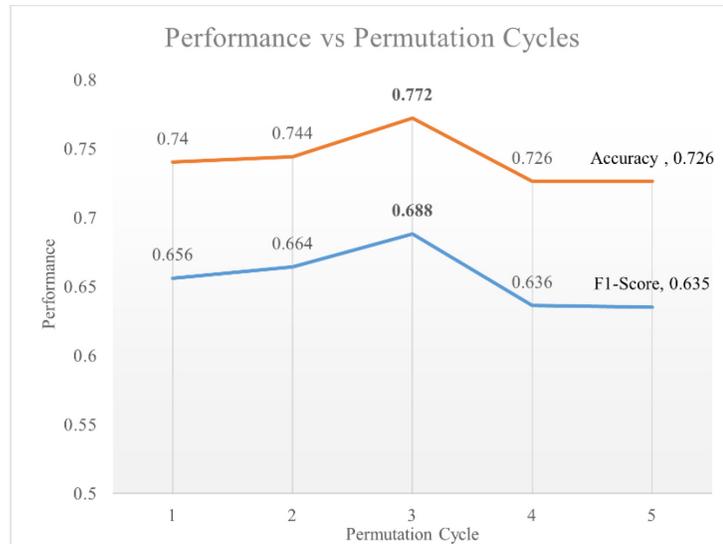


Figure 5-5: Performance variation vs permutation cycles

In addition to the attribute importance now being better understood as a result of the multi-step feature selection method, as observable in Figure 5-5, the first three rounds of the permutation step resulted in an improvement of accuracy and F1-Score from 0.74 to 0.77 and 0.65 to 0.69 percent respectively. This means that the process not only assists in a better understanding of the important attributes in change prediction but also increases the prediction performance. As a result, the process is cut at the third step. The remaining list of attributes describes the most dominant features in change order prediction.

As illustrated in Figure 5-4 the results from both methods align in identifying the most dominant attributes. "Duration" and "Contract Value" emerge as the two highest correlated attributes based on entropy and permutation analysis, with the contract value demonstrating a stronger correlation to the predictions compared to duration. Following these, the HCPI transformation of project location, measured from the number of positive values (COs) normalized by a unit of the contract value, ranks as the next most significant attribute. The subsequent important group includes the HCPI transformation of "Classification," which describes the number of positive value COs. Interestingly, demographic attributes, such as "Population" in permutation analysis and "Population Density" in entropy results, rank higher than "Project Type", "Industry Group", and "Operating Unit", showing their high correlation to the target attribute.

Chapter 6 – Analysis of change orders sequential behavior

This chapter explores the co-dependence between the severity of COs issued in different phases of projects or in other words ‘change sequentiality’. Specifically, it addresses whether the PCICs (Periodic Cost Impact of Change orders) projected on the construction timeline follow a Markovian process. PCIC and CIC (Cost Impact of Change orders) are the terms used in this study to distinguish between the cost impact of COs issued in a certain phase/period of a project and the cost impact of all COs of a project. CIC aggregates the cost impact of all COs whereas PCIC aggregates COs based on the phase they are issued. A similar terminology is defined by (Du et al. 2016b) for cost performance to differentiate between the overall cost performance of a project and the cost performance at a certain time of a project. The term phase refers to a portion of the construction duration and not the project phases such as design, and construction phase. In performance forecasting, specifically in construction projects, it is critical to accurately deal with correlated variables. Common methods treat performance correlation by assuming the same performance for future activities as measured for past activities (Barraza et al. 2004). The alternative method is to employ historical change performances of projects to measure future change performance using the trends captured in past projects. This requires measuring the change performance at different phases of projects and following the variations through time. Moreover, as described earlier (in section 2.2), due to the complexity of analyzing dependencies among individual change orders (COs), it is recommended to investigate the aggregated effects of groups of COs (rather than individual instances). To simplify this process and to account for the need to evaluate the PCIC as outlined in section 2.6, the analysis is conducted at the phase level. This involves grouping change orders issued within a specific construction phase and looking into the group of COs as a single entity (i.e., the cost impact of the group of COs (PCIC of that phase)). The analysis of the co-dependence between PCICs issued at different phases is the aim of the second component of this study. It particularly investigates the suitability of Markov Chains for modeling the relation between PCICs. Section 6.2 describes details of the approach followed to assess the Markovian property and answer the question of whether PCIC in various phases of a project can be modeled as a Markovian process.

6.1 PCIC model

Investigating the dependencies among each instance of change orders (COs) is a complex challenge due to the complex nature of projects. Change orders can arise at any stage of a project for a variety of reasons, such as design alterations, unforeseen site conditions, stakeholder requests, or compliance with new regulations. This variability complicates the analysis, as each change order may be influenced by separate factors specific to its context and timing. Moreover, the reasons behind change orders are not always clearly defined. Often, the rationale for a change may be ambiguous or influenced by a combination of factors, making it difficult to pinpoint a singular cause. This lack of clarity can lead to challenges in establishing accurate dependencies between current change orders and past events.

To simplify the complexity, this study explores the Markovian property in aggregated cost impact of change orders issued during different phases, referred to as the Periodic Cost Impact of Change Orders (PCIC) in this study, rather than analyzing each individual CO.

To model the Cost Impact of COs (CIC) as a stochastic process, the project timeline is divided into five phases, with each phase representing 20% of the total project duration. Each phase is considered a step in the time-space continuum, denoted as $t = 0, 1, 2, 3, 4, 5$. At each time

step t , a project exists in a specific change state that belongs to the change state space. In this study, the state space comprises three distinct states:

1. **Low**: Change percentage less than 4% of the project value.
2. **Medium**: Change percentage between 4% and 10% of the project value.
3. **High**: Change percentage exceeding 10% of the project value.

As outlined in Section 3.4, to evaluate the Markovian behavior of change state evolution between states 0 to 5, it is essential to first assume that the process follows a Markovian process. In a Markovian process, the state at time $t+1$ is solely dependent on the state at time t . Thus, the Markov chain can be described using its transition probability matrix, which captures the probabilities of transitioning between adjacent time steps (i.e., the 1-step transition matrix). The 2-step transition matrix between times t and $t+2$, denoted as $P(t, t+2)$, can be characterized by multiplying the two 1-step transition probability matrices (i.e., $P(t, t+1)$, and $P(t+1, t+2)$). Consequently, the 4-step transition matrix between time 1 and time 5 can be computed by multiplying the four 1-step matrices. This 4-step transition matrix, derived using this method, is known as the Chapman-Kolmogorov equation.

In the case of validity of Markovian assumption, this shows the overall expected transitions between steps 1 and 5, $P_e(0,5)$. In the case of an ideal Markov process, this expected transition matrix should be equal to the observed transitions between steps 1 and 5 when skipping the times in between, $P_o(0,5)$. The 4-step transition matrix can be calculated also from the real observations of state migrations when looking into time 1 to time 5. The goal is to explore the similarity of the expected and observed five-step transition matrices, $P_o(0,5)$, and $P_e(0,5)$. The comparison can be done through a Chi-squared test by considering a 0.05 dissimilarity significance level .

Entries of the one-step transition matrices can be easily calculated using Equation 6-1:

$$p_{ij}(s, t) = \frac{E_{ij}(s, t)}{N_i} \quad \text{Equation 6-1}$$

Where $E_{ij}(s, t)$ is the number of migrations from state i to state j , when moving from time s to time t .

6.2 Testing for Markovian property

Two types of tests are considered to assess the Markovian behavior of PCIC. The first group considers Cumulative Cost Impact (CCI) as the term used to define the state of projects at each period. CCI values are measured by aggregating the values of all former periods. The second group of tests does not cumulate the values of previous periods and as a result, the state of projects is defined using Singular Cost Impact (SCI) values. Moreover, since the duration of projects is not the same, to reduce the effects of differences in project durations, projects are categorized into four groups based on the length of projects, and the two types of tests are conducted for each of the four groups separately. The four categories of projects are shown in Table 6-1.

Table 6-1: Project categories by project duration

Category	Project Duration (PD)
PD0	All projects
PD1	PD < 150 days
PD2	150 days < PD < 365 days
PD3	PD > 365 days

Then, for each of the two types of tests, the one-step transition matrices are shaped by looking into the historical data and counting the number of transitions from one state to another between adjacent periods. The 1-step probability matrices for both CCI and SCI values can be found in Appendix 1. Having the 1-step matrices shaped, the 4-step observed and expected probability matrices are developed from two methods. Observed transitions by counting the number of migrations from one state to another when looking into only phases 1, and 5. The expected values, however, are calculated using the Chapman-Kolmogorov equation. This group of matrices is shown in Appendix 2.

Then to evaluate the Markovian property, it is needed to compare the 4-step observed and expected values. Utilizing the Chi-Squared method, the expected and observed transition matrices for each of the eight tests are compared to check the similarity. By considering a confidence level of 0.95, the p-values less than 0.05 confirm the validity of the null hypothesis (i.e., Markovian behavior).

Table 6-4 show the result of the Chapman-Kolmogorov equation validation, to assess the Markovian property in SCI, and CCI. It is identical to the p-values of this table that in all of the eight tests, the p-values are greater than 0.05 significance level, and as a result, the Markovian assumption is valid. The four groups of project duration were defined to investigate the sensitivity of the assumption to the project duration differences. The results confirm that regardless of project duration the assumption is valid for both CCI and SCI with a 0.05 confidence level assumption.

Table 6-2: CHI-Squared results

Test Type	Duration Group	CHI-Squared	P-Value
SCI	PD-1	0.615	0.99
SCI	PD-2	12.54	0.051
SCI	PD-3	6.34	0.385
SCI	PD-0	7.59	0.27
CCI	PD-1	3.45	0.75
CCI	PD-2	4.39	0.62
CCI	PD-3	4.89	0.56
CCI	PD-0	9.05	0.17

The results indicate that while the Markovian behavior is validated for each of the eight tests, the Markovian property for the SCI values of mid-term projects (i.e., PD-2) is validated with less confidence. It is evident that, excluding the SCI test for PD-2 projects, the level of confidence

diminishes as project duration increases. Among the various project types, the PD-0 projects exhibit the lowest confidence level, characterized by a non-homogeneous sample in terms of project duration (excluding the SCI test for PD-2 projects). In contrast, the PD-1 projects, which are short-term projects lasting less than 150 days, demonstrate the highest level of similarity across both tests. This finding aligns with the P-Values of the PD-1 projects, which are 0.99 for SCI and 0.75 for CCI, respectively. Interestingly, the results show no significant difference in confidence levels between the SCI and CCI tests, as both tests validate two groups of projects with higher confidence. However, it was anticipated that the CCI values would exhibit greater confidence due to their formulation, which inherently incorporates the historical data from the previous phase. Specifically, the CCI for each phase is based on the change percentage from the former phase, combined with the change percentage of new change orders issued in the current phase.

Chapter 7 – Discussion

This chapter discusses the results of this study and how they contribute to the research objectives. One of the key objectives was to examine the importance of features that influence the change performance of projects. Therefore, the chapter begins with a discussion of the findings related to different types of attributes in section 7.1 to understand the project specifications with more influence on change severity prediction while assessing the effectiveness of the methods used to improve change predictability. Since another objective of this study was to improve the change prediction, in section 7.2, the results of the first component of this study are discussed from a data-mining perspective. Then section 7.3 presents a discussion of the results of the second component of this study, focusing on the sequential behavior of change orders and their influence on the prediction of COs.

7.1 Feature importance

Construction projects are unique, and consequently, and so are the changes they experience. The uniqueness of projects can be described by their specific characteristics of projects. The attributes used in this study as predictors are the representers of some aspects of projects that make a project unique from another. Analyzing the importance of these features in change prediction assists in understanding the features with more influence on the level of change that a project experiences. As mentioned in section 2.4 the location of projects, actors, and the types of projects are the three categories of features that can explain the uniqueness of projects. Due to data limitations, only the spatial and project attributes were investigated in this study. According to (Cantarelli et al. 2012; Love et al. 2017), based on location and type, projects experience different levels of cost variation. This means projects with the same spatial and project-type characteristics are more likely to have similar change experiences putting aside other important factors such as the size of projects. In section 4.2.3, a group of attributes was introduced to capture and quantify the change behavior of groups of projects with the same location and project type properties. The term used in this study to capture the change performance of these project groups is the number of COs as had been identified as a directly correlated term with change order cost impact (Shrestha et al. 2022). To neutralize the effect of the size of projects this term was then normalized once by project duration and another time by the contract value to be able to focus only on the spatial, and project-type attributes.

To highlight the importance of attributes and remove the redundant features a multi-step feature selection process was followed. The usefulness of attributes was investigated with several methods (i.e., Prediction performance, Entropy, PCA, Chi-squared, and permutation). Figure 7-1 compares the results of Entropy analysis before and after performing this process. As a result of this process, the number of attributes was reduced from 33 to 13, resulting in better show up of important attributes while reducing the dimensionality. Looking into the entropy level of attributes before and after this process, the significant features are more clearly shown. Table 7-1 also shows the importance of attributes from PCA point of view. In the following, these attributes are discussed more in detail.

Looking into Entropy and Permutation results in Figure 5-4, project size attributes (i.e., ‘Duration’ and ‘Project Value’) showed the highest level of correlation to the change prediction. Even though they are less important in terms of the level of variation they initiate into the data, as identical to the PCA results in Table 7-1, they are the most influential in terms of entropy and permutation as shown in Figure 7-1. This finding is aligned with a study by (Choi et

al. 2016), which showed that the occurrences of change orders are highly dependent on the duration and value of projects. Among the project size features, contract value ranked higher than duration meaning that the change severity is more influenced by the dollar value of projects rather than duration, however, the duration is still the second most important attribute. The critical role of contract value becomes clear when looking into the next correlated attributes from the entropy and permutation point of view. The next important feature, HCPI-City-P-CV is showing the average change frequency per dollar value of these groups of projects. This means in projects in the same location, the amount of dollars spent on the project governs the probability of occurrences of COs rather than days passed. As a result, project planners should pay more attention to taking proactive approaches to reduce changes specifically when dealing with accelerated projects in which higher dollar amounts are spent in a shorter duration.

Following the ‘Duration’, and ‘Contract Value’ was the ‘HCPI-P-City-CV’. This attribute shows the average number of COs with positive values in past projects of the same location. Having this attribute as the next critical feature after the project size attributes shows the significant role of the location of projects, however, this aspect of projects has received less attention in past studies. A high correlation of the frequency of COs per dollar values of projects in the same location with the prediction means that projects executed in the same location are more likely to show the same performance in terms of the number of change orders and accordingly the magnitude of COs. Another point to consider is the number of positive value change orders has more effect on the severity of COs compared to the number of Negative value COs.

The next correlated attribute in change prediction was ‘Classification’. This attribute is an internal categorization of projects based on projects’ importance to the organization. Having importance as one of the key attributes in change prediction can be due to allocating more resources to projects with higher strategic value to the construction firms. As a result, construction planners when analyzing past projects for decision-making regarding the change performances of future projects should consider gaining information from projects with the same level of significance to the organization. In the same direction as the location of projects, the importance feature appeared in the format of HCPI attributes (i.e., the number of positive value change orders per dollar value). This can indicate that the more important projects get more planning attention and as a result, fewer unseen conditions occur.

Following Classification are the demographic attributes of the location of projects (i.e., ‘Population’ and ‘Population Density’). The first point to consider is the critical role of the location of projects and their specifications in predicting change orders of projects. The presence of spatial attributes even in higher rank compared to the type of projects illustrates the significant role of the location of projects and their specifications in CO performances of projects. The second critical point is the importance of demographic features and accordingly, the change reasons associated with them. Change reasons such as labor shortage, weather conditions, and opposition of local communities are among the change reasons under the influence of demographic properties of project location. Having these attributes among the dominant attributes in change prediction confirms the critical role of these attributes and their capacity to define the severity of change.

It is quite common in construction projects to use the outcomes of past projects to plan for upcoming projects. This includes contingency, cost, schedule, and labor efficiency planning. Since all of these terms are highly dependent on the change performance of projects, the construction

planers should pay more attention to making their decisions from similar projects in terms of demographic specifications of their location, project importance, and project type.

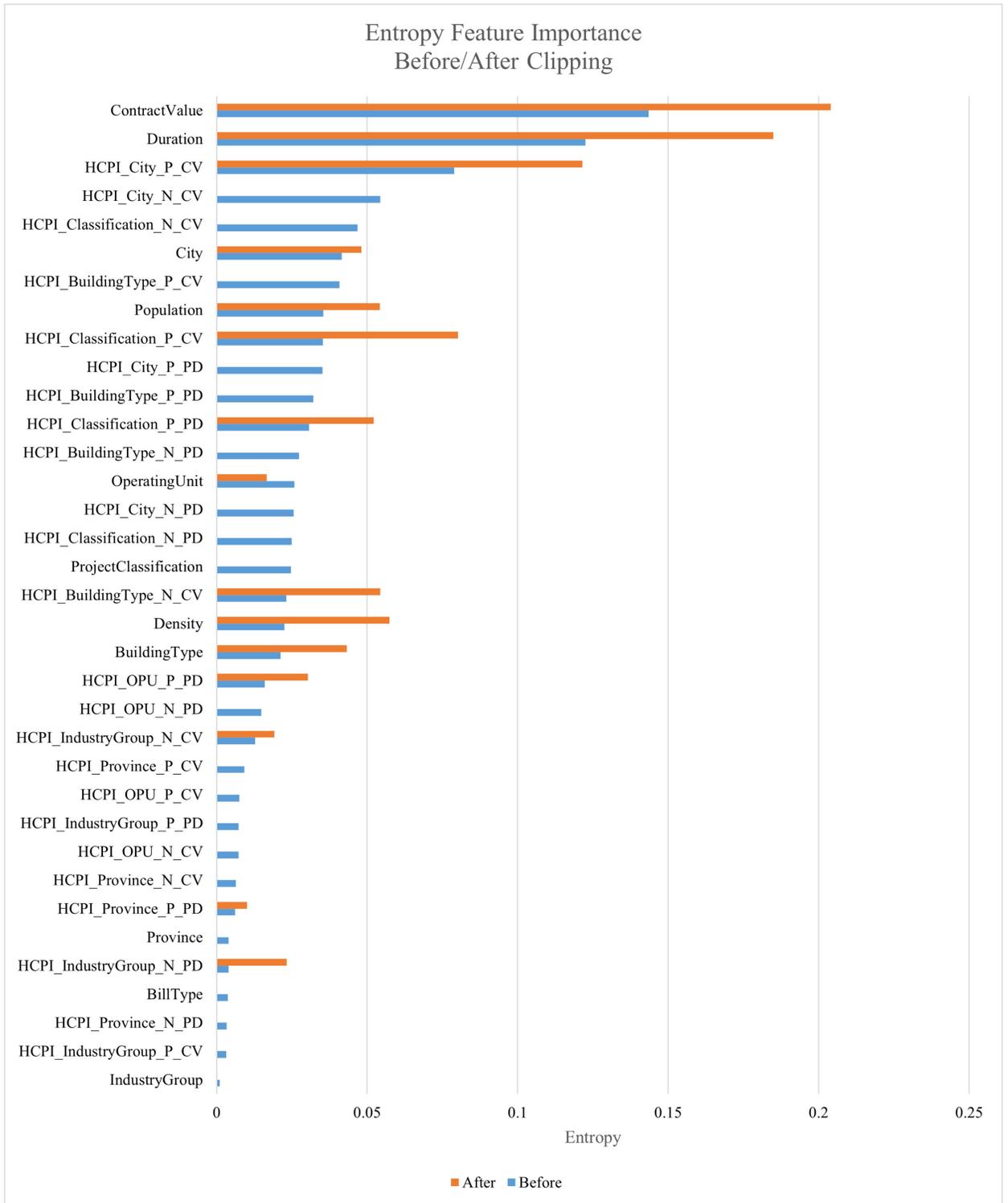


Figure 7-1: Entropy results before and after feature selection steps

Table 7-1: Dominant attributes of PCA after feature selection

PC-1	PC-2	PC-3
HCPI_IndustryGroup_N_DU	Duration	HCPI_OPU_P_DU
HCPI_IndustryGroup_N_CV	ContractValue	
HCPI_Classification_P_CV	HCPI_Province_P_DU	
HCPI_Classification_P_DU	HCPI_City_P_CV	
HCPI_BuildingType_N_CV		

7.2 Data mining perspective

In this study, a multi-step feature selection method was followed to investigate the usefulness of attributes and analyze the importance of feature selection in the prediction of the severity of COs. This is remarkable that although 60% of attributes are removed from data as a result of the feature selection process, the performance of the model has increased by almost 14%. As a result of the transformation made by HCPI features, more layers of information are added to the data, and accordingly, the performance of the prediction increased by 7%, and 8% in terms of F1-Score, and Accuracy relatively (i.e., by comparing SC-2 VS SC-6, SC-1 VS SC-7, and SC-3 VS SC-4 in Table 5-5) solely due to the integration of HCPI attributes into the data. Removing the highly correlated attributes using a Spearman-based dendrogram and permutation analysis added an additional 2% improvement to the prediction performance. Finally, by utilizing permutation analysis, more redundant attributes are identified and omitted from the dataset. This approach has increased the prediction performance by 3% resulting in final prediction performance reaching 77% and 69% in terms of Accuracy, and F1-Score respectively.

Looking into the dominant attributes after the feature selection process in Figure 5-4 from the Entropy and Permutation point of view, the HCPI attributes were the most dominant attributes along with Duration and Project Value. The Principal Component Analysis (PCA) results revealed that the HCPI attributes emerged as the highest contributors to the variance within the data by being the only dominant attributes shaping the first two PCs. Second, remarkably, a unique pattern is observable in the importance of these new attributes. While original Spatial attributes, the Project Classification, and the Industry Group did not exhibit statistical significance when assessed through Chi-Squared analysis, the HCPI attributes generated from these seemingly less important categorical attributes turned out to be among the highest contributors in the PCA results. This underscores the valuable transformation these attributes underwent, allowing them to capture essential information that was not immediately evident when working with the raw categorical data.

Looking into HCPI attributes from another perspective, it is clear that they significantly increased the dimensionality of the model, it is noteworthy that this augmentation did not come at the cost of overfitting. This observation becomes evident when comparing the models' performance on both the test and training subsets. The close values of performance on these subsets indicate that the model has maintained its generalization capabilities, showing a balance between fitting the training data and making accurate predictions on unseen data. Furthermore, an examination of the parameters of the Random Forest (RF) model reveals that the complexity of the model has either remained the same or even decreased in some cases. Remarkably, the model

achieved higher performance with the same depth and number of trees, underscoring the efficiency and effectiveness of the new attributes in enhancing the predictive power of our model.

Looking into the PCA results in Table 7-1, as clearly noticeable, the first Principal Component (PC1), the most dominant PC, is primarily influenced by a subset of the HCPI attributes, generated by project type attributes. The second PC however, is governed by mostly spatial attributes along with Project Value and Duration. Considering the significant importance of the Duration, and Contract Value, it underscores the importance and correlation of project type and location of project attributes with the prediction.

Looking into the attribute type level results in Table 5-5, it is evident that the location of projects is an important factor in change orders performances of projects when comparing the SC-1 VS (SC-2, and SC-3) observing 2%, and 6% increase in F1-Score and Accuracy relatively due to adding location of projects and an additional 2% to the F1-Score as a result of the DP attributes. Their limited performance improvement shows that more attributes regarding the location of projects should be collected to describe different aspects of project location. However, when looking into, Entropy and permutation results, it is noticeable that the transition made as a result of HCPI attributes has effectively added more context to the location attributes becoming the third important attribute after Duration and Project Value. Moreover, the DP attributes, population, and population density have shown effective contributions to the model performance being among the top six highest influential attributes in both methods.

7.3 Change orders sequential behavior

As described in Chapter 6, this study was able to prove that the Periodic Cost Impact of Change orders (PCIC) can be modeled as a Markovian process by checking the validity of the Chapman-Kolmogorov equation. Two types of tests were conducted to evaluate the Markovian property in PCIC, one for Cumulative Cost Impact (CCI), and another one for Singular Cost Impact (SCI) referring to the method of aggregation of COs cost impact. Each of the methods was tested for four distinct groups of projects based on their duration. All eight tests confirmed the suitability of modeling the PCIC as a Markovian process. However, the PD2, i.e., the mid-term projects from a construction duration perspective, passed the assessment with less level of confidence, by being at a marginal acceptance range. Notably, when excluding the SCI test for PD-2 projects, a decline in confidence levels is observed as the project duration extends. Among the different project categories, PD-0 projects demonstrate the lowest confidence level, characterized by a heterogeneous sample in terms of project duration (excluding the SCI test for PD-2 projects).

Conversely, the PD-1 projects, which are short-term in nature and last less than 150 days, display the highest degree of similarity in both tests. This observation is consistent with the P-Values for PD-1 projects, which stand at 0.99 for SCI and 0.75 for CCI, respectively.

Interestingly, there is no substantial difference in the confidence levels between the SCI and CCI tests, as both methods validate two groups of projects with considerable confidence. It was, however, expected that the CCI values would yield higher confidence levels due to their formulation, which inherently takes into account the historical data from previous phases. Specifically, the CCI for each phase is derived from the change percentage of the preceding phase, combined with the change percentage of new change orders issued in the current phase.

It can be concluded, PCIC is memoryless, meaning that given the cost impact of the current period, the change orders' cost impact of the next period can be forecasted without knowing the state of the previous periods. In other words, the possible outcomes of the next steps are only dependent on the current step and not the previous ones. As a result, a prediction model forecasting the state at the time (t+1) given the state of the current state at the time (t) should have the same output irrespective of being provided with information of the state of previous steps at t=1, 2, ..., t-1.

To test this finding, the same RF model used in Chapter 5 as well as the list of attributes generated from the multi-step feature selection technique are utilized to predict PCIC and in other words change timing. Two scenarios are defined and compared to predict the CCI at period t+1 during the project construction duration. Both scenarios tend to predict the same target attribute but using different predictors. Both scenarios are provided with the same project features used in Chapter 5, however, in addition, the first scenario is provided with information on the current state of the project, while the second scenario is provided with information on all prior and the current states of project.

Figure 7-2 Shows the RF prediction performance in predicting the change state of the next periods in terms of accuracy and F1-Score. As identical, the model is capable of performing the same level of prediction performance for both scenarios, confirming the validity of the assumption. For both performance criteria, the level of variation is negligible showing less than 1% difference.

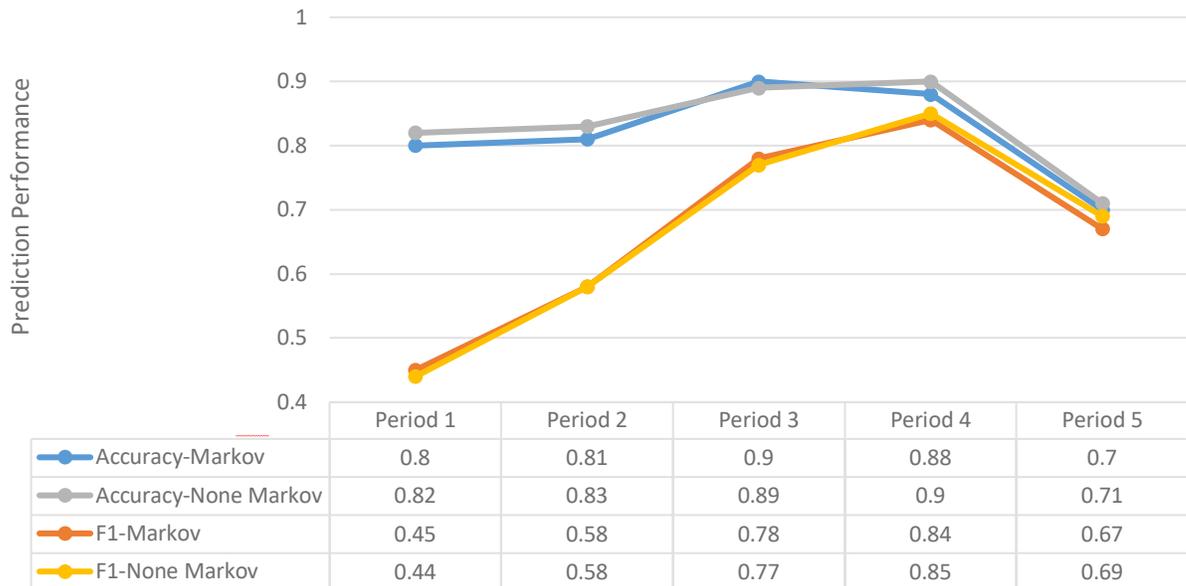


Figure 7-2: Periodic prediction with/without Markov assumption

As discussed earlier in Chapter 5 and sections 7.1, and 7.2, change performances of projects can be predicted using historical change performances of similar projects. It was explained how project type and spatial features are critical in grouping similar projects. The results of the Markovian assessment showed that, in addition to these features, project duration is important to be considered when predicting the periodic change performances of projects. Either using machine learning techniques to analyze past projects or comparing past projects with current projects as project planners to evaluate expected change performances of projects, these decisions should be derived from projects with the same project duration group. This was understood by observing the

Markovian assumption being validated with higher confidence in projects with the same project duration group. Moreover, these results facilitate the change performance monitoring by showing that the change state of the next phase of a project is only dependent on to change state of the current phase and not previous ones. As a result, in change monitoring, when adjusting the change expectations for the next phases of projects based on the current conditions, it is only needed to compare those projects with the ones with the same current state regardless of their different performances in other phases.

Chapter 8 – Concluding remarks

8.1 Summary of the study and its key findings

Change orders (COs) are known as one of the most influential factors in the failure of construction projects. They are disruptive to the efficiency of construction planning decisions from different aspects. However, a great portion of the critical strategic projects' decisions are being made with least level of awareness on change performances of projects. Change management tends to mitigate the diverse effects of change orders by recognizing and implementing them as early as possible. The predictive analysis of COs aims to anticipate change orders' severity early before starting projects. Although the benefits are numerous, less attention has been given to this aspect of COs analysis. This can be as a result of the complexity of COs, and their association with several aspects of projects such as actors, location, and specific characteristics of projects. A change predictive model should be able to anticipate the number of change orders, their magnitude, and timing to help construction planners consider the change impacts in their decisions. Several factors need to be included in construction data collection standards to explain the specific conditions of projects. These features should cover more than 80 change reasons identified so far, and the specific project, actors, and location characteristics. Moreover, to predict the timing of change orders it is important to accurately model the co-dependence of change orders to account for the impact of change orders on initiating more change orders in upcoming activities.

The main goal of this study was to investigate the predictability of the severity of change orders, correlated features, and ways to improve predictability. A review of previous related works was conducted to identify the gaps in the change prediction field. Then, two components were defined for the methodology. The first component, by focusing on change severity prediction, used several feature selection techniques to identify the key correlated attributes with the prediction. Moreover, by integrating two new sets of attributes, more layers of information were added to the data to unlock hidden layers of information and improve prediction. Getting the benefit of the change domain knowledge, a set of features (i.e., HCPI attributes) was introduced to quantify the varied change performances of different groups of projects using the frequency of change orders in past projects. Another set of features integrated specific characteristics of the location of projects from external sources to better explain the project spatial features. The feature selection techniques showed that the HCPI features add a valuable layer of information to the data that can significantly increase change prediction performance. The results confirmed the usefulness of these attributes in several ways. Looking at the performance of models, up to 15% improvement was achieved using these features. The PCA results, indicated the characterization of the most dominant PC (PC-1) by the presence of only the HCPI attributes. Moreover, the only attributes that were presented in the final list of dominant attributes along with these attributes were Duration and Contract Value; the most influential specification of projects in occurrence and severity of change orders (i.e., from a prediction point of view).

Entropy results were also in the same direction as PCA, Showing the significance of Duration, Contract Value, and HCPI features. However, more attention should be given to the spatial attributes while observing the third level of important features belonging to the location-based features and their HCPI representative. It can be assumed that these types of features are

important, however, more contextual information is needed to provide more specification of the location of projects.

By comparing the Entropy and PCA results before and after feature selection, the remaining attributes more accurately reflect the significance of individual attributes, confirming the usefulness of the feature selection process of this study. As a result, feature selection techniques should be given more attention in future studies.

Clarification is needed that this paper is not a game of reaching higher accuracy. The author believes that the performance of models in such complex problems as the prediction of change orders is dependent on a variety of criteria such as quality and size of data, viewpoint of the data collector, and a lot more. However, using the same data, higher accuracies could have been reported in cost of not addressing the imbalanced data reporting accuracy instead of F1-Score. Since there isn't a short-term relief to the problem of lack of data and quality issues, this study is focused on the ways to get the most from currently available data, however, it suggests the parties involved in the industry collect more informative data and make it accessible to the academic community.

The second component, by examining the co-dependence of change orders constructed a foundation to enable change timing prediction. Due to the complex nature of the change orders and specifically the relation between individual COs, it is needed to evaluate the dependence of groups of COs. This study, by dividing projects' durations into five increments assessed the co-dependence of periodic cost impact of change orders (PCIC). In other words, it evaluated the effect of a period's change severity on the same term in the upcoming periods. To simplify the change timing modeling, this study particularly examined the suitability of Markov Chain models for capturing the co-dependence of PCIC. To this aim, the Chapman-Kolmogorov equation was employed to validate the assumption. The test was conducted for different groups of projects based on the duration to account for sensitivity to the size of projects. The results showed that the PCIC can be modeled as a Markovian process for all groups of projects. In the concept of periodic change severity prediction, this means that when predicting a period's change severity only one temporal attribute is needed to provide the immediate previous period's change state.

The results of the Markovian property assessment test revealed that the PCIC can be modeled as a Markovian process. This means given the current period's COs severity, the future period's change severity is only dependent on this period and not the previous ones. This helps to reduce the complexity of change orders timing prediction by keeping the connection of periods to their last former periods and removing the connection to the other periods. Moreover, it also helps a lot in change orders monitoring by knowing that each period only affects its next period and not the others. Although all project categories were able to pass the tests, it was observed that groups of projects with homogenous durations passed the tests with higher confidence. As a result, it is important to account for variation of projects duration.

8.3 Research contributions

On one hand, although few previous studies have tackled the change severity prediction, limited efforts have been made to identify the correlated feature with the change severity prediction. Other studies on the change management domain have shown that based on the type, and location of projects the change performances of projects vary. However, the available construction data features do not include a term to quantify the distinct behavior of different groups

of projects. Moreover, the spatial features collected in construction data do not explain well the different characteristics of their location. On the other hand, several studies have addressed the significant influence of change timing on the change impact level. In the construction literature, the change timing term is defined as the periodic cost impact of change orders. To include change timing prediction in the predictive analysis of change orders, it is required to be able to model the temporal co-dependence of COs. However, the temporal dependence of COs is not completely clear.

To address these gaps in the change prediction field two components were defined for the methodology of this study. (i) In the first component, by developing predictive models and feature selection techniques the importance of features was investigated in detail. The findings of the first components not only helped to improve change prediction models but were also important to change management. Understanding the importance of attributes helps the change management to understand what aspects of projects contribute to their change behavior. Moreover, these findings can be used in developing construction data collection standards and protocols to consider the change prediction requirements. (ii) The second component developed a foundation for change timing prediction by investigating the co-dependence of periodic cost impact of change orders (PCIC). This component particularly examined the suitability of Markov Chain models for modeling the temporal co-dependence of PCIC.

To summarize, the main contributions of this research are:

- Development of an innovative term based on the domain knowledge of change management to quantify the varied change behavior of different groups of projects;
- Highlighting the importance of projects' spatial features in change prediction;
- Assessing the correlation of demographic features in change prediction;
- Highlighting the significant correlation of the average number of change orders in similar past projects in change prediction;
- Examining the importance of feature selection in change prediction;
- Developing a simplified process for change timing modeling and prediction;
- Improving change prediction performance by feature selection techniques.

8.4 Limitations of this thesis

Same as other data-driven studies, the size and quality of data are major factors influencing the outcomes. The initial dataset contained more than 3000 projects, however, due to quality issues, and the different nature of projects, only 1024 projects were used for this study. Moreover, a lot of missing values were dealt with. Although different methods were chosen for each attribute in a way to reduce the impact on data, the uncertainty in dealing with missing values remains a limitation of this study. Another major challenge in the data was linking the change orders to their change reasons which was provided in a distinct dataset. Due to the lack of linkage between the change orders and the other dataset which contained change reasons, this task couldn't be done with high confidence. As a result of that, this study did not consider the change reasons in its prediction. Looking into the data from another perspective, the data was collected from projects executed in Canada by a General contractor, as a result, the findings of this study stay limited to the same project types (electrical and mechanical construction projects), and the Canadian construction industry. Data size in terms of both the number of features and data points is another limitation of this study. More than half of the projects and their corresponding change orders were

removed from the dataset due to missing values and the projects not being aligned with the scope of the research. Another issue with the data was the inaccurate values of project start, end date, and change orders date. Change timing prediction is not addressed directly in this study due to this issue.

One of the major objectives of this study was to analyze the importance of HCPI attributes in change prediction. These features utilize the average change performances of similar projects to quantify the varied change performances of different project types. In this study, the similarity is defined at the most abstract level as belonging to the same category of categorical features. Even at this level, the attributes performed well in change prediction. However, a more detailed analysis of project specifications is needed to identify the required features for grouping similar projects. Another set of features analyzed were those related to the location of projects. In this study, demographic features were considered, but other important aspects such as weather conditions, specific regulations, and underground conditions were not covered. Other types of logically important attributes specifically those linked with the change reasons that were not evaluated in this research include but are not limited to features describing the parties involved in the projects, the complexity of projects, and financial, political, and social conditions. Since the algorithms used to assess the importance of attributes could only deal with numerical values, the categorical attributes are encoded using integer coding, however, the attributes did not have ordinal relationships among their categories. As a result, the integer coding is also another limitation of this study. The predictions were conducted to classify projects into two classes. The definition of these classes and the number of them were other limitations of this study. The same limitation applies to the assessment of the Markovian property regarding the definition of states. Moreover, this study examined Markov behavior by dividing project duration into five phases and did not investigate other ways to define phases, such as monthly divisions.

8.5 Future studies

It needs to be emphasized that in this study, it was decided to keep the HCPI attributes at the simplest level, by defining similarity criteria using a single variable, (i.e., categories of nominal attributes) to better realize the contribution of these attributes. However, future studies should investigate complex rules such as grouping projects based on a combination of attributes, time of execution, size of projects, and more.

Future studies should also consider integrating data from different project types, and contractors to analyze the significant attributes that can describe the differences of project types and contractors' performance. Moreover, a lot more attributes are logically correlated with change reasons, the importance of these attributes specifically those related to the location of projects, involved parties, and project complexity needs to be investigated. The feature importance analysis can also be conducted using other feature importance evaluation methods.

The choice of change order magnitude to define the target classes can be investigated with different looks such as targeting a balanced distribution of data between classes or considering more classes; however, it can be selected to meet the specific needs of organizations. In addition, instead of putting all change orders in one basket, they can be grouped based on change reasons, to also address the change causes in the prediction. This also opens gates to a more detailed analysis of important attributes in change orders prediction, by concurrent feature important analysis using the proposed multi-step feature selection method while predicting change orders magnitude and

their reasons. This will provide more informative insights into the attributes associated with different change reasons.

There are numerous areas related to change timing that can be interesting topics for future research. Change-timing prediction and development of change orders monitoring tools to be added to the change management systems, knowing that the PCIC can be modeled with less complexity as a Markovian process. Other topics include analyzing the key attributes and change causes in different phases of projects shading lights into more probable COs in different stages of projects.

References

- Abaza, K. A. 2016. "Simplified staged-homogenous Markov model for flexible pavement performance prediction." *Road Materials and Pavement Design*, 17 (2): 365–381. <https://doi.org/10.1080/14680629.2015.1083464>.
- Abaza, K. A. 2017. "Empirical Markovian-based models for rehabilitated pavement performance used in a life cycle analysis approach." *Structure and Infrastructure Engineering*, 13 (5): 625–636. <https://doi.org/10.1080/15732479.2016.1187180>.
- Adekoya, O. B., J. A. Oliyide, and G. O. Oduyemi. 2021. "How COVID-19 upturns the hedging potentials of gold against oil and stock markets risks: Nonlinear evidences through threshold regression and markov-regime switching models." *Resources Policy*, 70. <https://doi.org/10.1016/j.resourpol.2020.101926>.
- Aït-Sahalia, Y. 1996. "Testing Continuous-Time Models of the Spot Interest Rate." *The Review of Financial Studies*, 9 (2): 385–426. <https://doi.org/10.1093/rfs/9.2.385>.
- Alanis, R., A. Ingolfsson, and B. Kolfal. 2013. "A markov chain model for an EMS system with repositioning." *Production and Operations Management*, 22 (1): 216–231. <https://doi.org/10.1111/j.1937-5956.2012.01362.x>.
- Alaryan, A., E. Elbeltagi, A. Elshahat, and M. Dawood. 2014. "Causes and Effects of Change Orders on Construction Projects in Kuwait." *Int. Journal of Engineering Research and Applications*.
- Ali, A., S. Naeem, S. Anam, and M. M. Ahmed. 2023. "Shannon Entropy in Artificial Intelligence and Its Applications Based on Information Theory."
- Al-Kofahi, Z. G., A. Mahdavian, and A. Oloufa. 2022. "System dynamics modeling approach to quantify change orders impact on labor productivity 1: principles and model development comparative study." *International Journal of Construction Management*, 22 (7): 1355–1366. <https://doi.org/10.1080/15623599.2020.1711494>.
- Alnuaimi, A. S., R. A. Taha, M. Al Mohsin, and A. S. Al-Harthi. 2010. "Causes, Effects, Benefits, and Remedies of Change Orders on Public Construction Projects in Oman." *Journal of Construction Engineering and Management*, 136 (5): 615–622. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0000154](https://doi.org/10.1061/(ASCE)CO.1943-7862.0000154).
- Anastasopoulos, P. C., S. Labi, A. Bhargava, C. Bordat, and F. L. Mannering. 2010. "Frequency of Change Orders in Highway Construction Using Alternate Count-Data Modeling Methods." *Journal of Construction Engineering and Management*, 136 (8): 886–893. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0000198](https://doi.org/10.1061/(ASCE)CO.1943-7862.0000198).
- Arefazar, Y., A. Nazari, M. R. Hafezi, and S. A. H. Maghool. 2022. "Prioritizing agile project management strategies as a change management tool in construction projects." *International Journal of Construction Management*, 22 (4): 678–689. Taylor & Francis. <https://doi.org/10.1080/15623599.2019.1644757>.
- Assaf, S. A., and S. Al-Hejji. 2006. "Causes of delay in large construction projects." *International Journal of Project Management*, 24 (4): 349–357. <https://doi.org/10.1016/j.ijproman.2005.11.010>.
- Athari, S. A., D. Kirikkaleli, and T. S. Adebayo. 2023. "World pandemic uncertainty and German stock market: evidence from Markov regime-switching and Fourier based approaches." *Quality and Quantity*, 57 (2): 1923–1936. <https://doi.org/10.1007/s11135-022-01435-4>.

- Attalla, M., and T. Hegazy. 2003. "Predicting Cost Deviation in Reconstruction Projects: Artificial Neural Networks versus Regression." *J. Constr. Eng. Manage.*, 129 (4): 405–411. [https://doi.org/10.1061/\(ASCE\)0733-9364\(2003\)129:4\(405\)](https://doi.org/10.1061/(ASCE)0733-9364(2003)129:4(405)).
- Barraza, G. A., W. E. Back, and F. Mata. 2004. "Probabilistic Forecasting of Project Performance Using Stochastic S Curves." *J. Constr. Eng. Manage.*, 130 (1): 25–32. [https://doi.org/10.1061/\(ASCE\)0733-9364\(2004\)130:1\(25\)](https://doi.org/10.1061/(ASCE)0733-9364(2004)130:1(25)).
- Biau, G., and E. Scornet. 2016. "A random forest guided tour." *TEST*, 25 (2): 197–227. <https://doi.org/10.1007/s11749-016-0481-7>.
- Bilal, M., L. O. Oyedele, J. Qadir, K. Munir, S. O. Ajayi, O. O. Akinade, H. A. Owolabi, H. A. Alaka, and M. Pasha. 2016. "Big Data in the construction industry: A review of present status, opportunities, and future trends." *Advanced Engineering Informatics*, 30 (3): 500–521. <https://doi.org/10.1016/j.aei.2016.07.001>.
- Bordat, C., B. McCullough, and K. Sinha. 2004. *An Analysis of Cost Overruns and Time Delays of INDOT Projects*. FHWA/IN/JTRP-2004/07, 2811. West Lafayette, IN: Purdue University.
- Budach, L., M. Feuerpfeil, N. Ihde, A. Nathansen, N. Noack, H. Patzlaff, F. Naumann, and H. Harmouch. 2022. "The Effects of Data Quality on Machine Learning Performance." arXiv.
- Burati, J. L., J. J. Farrington, and W. B. Ledbetter. 1992. "Causes of Quality Deviations in Design and Construction." *J. Constr. Eng. Manage.*, 118 (1): 34–49. [https://doi.org/10.1061/\(ASCE\)0733-9364\(1992\)118:1\(34\)](https://doi.org/10.1061/(ASCE)0733-9364(1992)118:1(34)).
- Butt, A. A., M. Y. Shahin, K. J. Feighan, and S. H. Carpenter. 1987. "PAVEMENT PERFORMANCE PREDICTION MODEL USING THE MARKOV PROCESS." *Transportation Research Record*, (1123).
- Caldas, C. H., and L. Soibelman. 2003. "Automating hierarchical document classification for construction management information systems." *Automation in Construction*, 12 (4): 395–406. [https://doi.org/10.1016/S0926-5805\(03\)00004-9](https://doi.org/10.1016/S0926-5805(03)00004-9).
- "canadacities.csv." 2023. <https://simplemaps.com/data/canada-cities>.
- Cantarelli, C. C., B. Flyvbjerg, and S. L. Buhl. 2012. "Geographical variation in project cost performance: the Netherlands versus worldwide." *Journal of Transport Geography*, Special Section on Theoretical Perspectives on Climate Change Mitigation in Transport, 24: 324–331. <https://doi.org/10.1016/j.jtrangeo.2012.03.014>.
- Carty, G. J. 1995. "Construction." *Journal of Construction Engineering and Management*, 121 (3): 319–328. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)0733-9364\(1995\)121:3\(319\)](https://doi.org/10.1061/(ASCE)0733-9364(1995)121:3(319)).
- "Change Events - Workflow Diagrams." 2017. *Procure*. Accessed September 8, 2024. <https://support.procure.com/products/online/user-guide/project-level/change-events/workflow>.
- Chen, B., and Y. Hong. 2012. "TESTING FOR THE MARKOV PROPERTY IN TIME SERIES." *Econom. Theory*, 28 (1): 130–178. <https://doi.org/10.1017/S0266466611000065>.
- Choi, K., H. W. Lee, J. Bae, and D. Bilbo. 2016. "Time-Cost Performance Effect of Change Orders from Accelerated Contract Provisions." *J. Constr. Eng. Manage.*, 142 (3): 04015085. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0001071](https://doi.org/10.1061/(ASCE)CO.1943-7862.0001071).
- Chou, J.-S., C.-F. Tsai, and Y.-H. Lu. 2013. "PROJECT DISPUTE PREDICTION BY HYBRID MACHINE LEARNING TECHNIQUES." *Journal of Civil Engineering and Management*, 19 (4): 505–517. <https://doi.org/10.3846/13923730.2013.768544>.

- Chua, D. K. H., and Md. A. Hossain. 2012. "Predicting Change Propagation and Impact on Design Schedule Due to External Changes." *IEEE Transactions on Engineering Management*, 59 (3): 483–493. <https://doi.org/10.1109/TEM.2011.2164082>.
- Chua, D. K., and M. A. Hossain. 2011. "Predicting change propagation and impact on design schedule due to external changes." *IEEE Transactions on Engineering Management*, 59 (3): 483–493. IEEE.
- Clarkson, P. J., C. Simons, and C. Eckert. 2004. "Predicting Change Propagation in Complex Design." *Journal of Mechanical Design*, 126 (5): 788–797. <https://doi.org/10.1115/1.1765117>.
- Czech, K., and M. Wielechowski. 2021. "Is the alternative energy sector covid-19 resistant? Comparison with the conventional energy sector: Markov-switching model analysis of stock market indices of energy companies." *Energies*, 14 (4). <https://doi.org/10.3390/en14040988>.
- Dao, B., S. Kermanshachi, J. Shane, S. Anderson, and E. Hare. 2017. "Exploring and Assessing Project Complexity." *Journal of Construction Engineering and Management*, 143 (5): 04016126. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0001275](https://doi.org/10.1061/(ASCE)CO.1943-7862.0001275).
- Du, J., M. El-Gafy, and D. Zhao. 2016a. "Optimization of Change Order Management Process with Object-Oriented Discrete Event Simulation: Case Study." *J. Constr. Eng. Manage.*, 142 (4): 05015018. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0001092](https://doi.org/10.1061/(ASCE)CO.1943-7862.0001092).
- Du, J., B.-C. Kim, and D. Zhao. 2016b. "Cost Performance as a Stochastic Process: EAC Projection by Markov Chain Simulation." *Journal of Construction Engineering and Management*, 142 (6): 04016009. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0001115](https://doi.org/10.1061/(ASCE)CO.1943-7862.0001115).
- Du, Y.-N., Z.-C. Qin, C.-C. Guan, D.-C. Feng, and G. Wu. 2024. "Bayesian model updating of super high-rise building for construction simulation." *Structural Design of Tall and Special Buildings*, 33 (10). <https://doi.org/10.1002/tal.2104>.
- El-adaway, I. H., I. S. Abotaleb, M. S. Eid, S. May, L. Netherton, and J. Vest. 2018. "Contract Administration Guidelines for Public Infrastructure Projects in the United States and Saudi Arabia: Comparative Analysis Approach." *J. Constr. Eng. Manage.*, 144 (6): 04018031. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0001472](https://doi.org/10.1061/(ASCE)CO.1943-7862.0001472).
- Fang, Y., and L. Sun. 2019. "Developing A semi-markov process model for bridge deterioration prediction in Shanghai." *Sustainability (Switzerland)*, 11 (19). <https://doi.org/10.3390/su11195524>.
- Fathi, M., P. P. Shrestha, and B. Shakya. 2020. "Change Orders and Schedule Performance of Design-Build Infrastructure Projects: Comparison between Highway and Water and Wastewater Projects." *J. Leg. Aff. Dispute Resolut. Eng. Constr.*, 12 (1): 04519043. [https://doi.org/10.1061/\(ASCE\)LA.1943-4170.0000353](https://doi.org/10.1061/(ASCE)LA.1943-4170.0000353).
- Fisher, R. A. 1922. "On the Interpretation of χ^2 from Contingency Tables, and the Calculation of P." *Journal of the Royal Statistical Society*, 85 (1): 87. <https://doi.org/10.2307/2340521>.
- Flyvbjerg, B., N. Bruzelius, and W. Rothengatter. 2003a. *Megaprojects and Risk: An Anatomy of Ambition*. Cambridge University Press.
- Flyvbjerg, B., M. K. Skamris holm, and S. L. Buhl. 2003b. "How common and how large are cost overruns in transport infrastructure projects?" *Transport Reviews*, 23 (1): 71–88. Routledge. <https://doi.org/10.1080/01441640309904>.

- Foidl, H., and M. Felderer. 2019. "Risk-based data validation in machine learning-based software systems." *Proceedings of the 3rd ACM SIGSOFT International Workshop on Machine Learning Techniques for Software Quality Evaluation*, MaLTesQuE 2019, 13–18. New York, NY, USA: Association for Computing Machinery.
- Golnaraghi, S., O. Moselhi, S. Alkass, and Z. Zangenehmadar. 2020a. "Predicting construction labor productivity using lower upper decomposition radial base function neural network." *Engineering Reports*, 2 (2). <https://doi.org/10.1002/eng2.12107>.
- Golnaraghi, S., O. Moselhi, S. Alkass, and Z. Zangenehmadar. 2020b. "Modelling construction labour productivity using evolutionary polynomial regression." *International Journal of Productivity and Quality Management*, 31 (2): 207–226. <https://doi.org/10.1504/IJPQM.2020.110024>.
- Goodfellow, Y. B. 2016. *Deep Learning*. MIT Press.
- Grossi, P. 2009. "Property Damage and Insured Losses from the 2001 World Trade Center Attacks." *Peace Economics, Peace Science and Public Policy*, 15 (2). <https://doi.org/10.2202/1554-8597.1163>.
- Habhoub, D., A. Desrochers, and S. Cherkaoui. 2009. "Agent-based assistance for engineering change management: An implementation prototype." *2009 13th International Conference on Computer Supported Cooperative Work in Design*, 288–293.
- Hagendorff, T. 2021. "Linking Human And Machine Behavior: A New Approach to Evaluate Training Data Quality for Beneficial Machine Learning." *Minds & Machines*, 31 (4): 563–593. <https://doi.org/10.1007/s11023-021-09573-8>.
- Hajjar, D., and S. M. AbouRizk. 1998. "Modeling and Analysis of Aggregate Production Operations." *Journal of Construction Engineering and Management*, 124 (5): 390–401. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)0733-9364\(1998\)124:5\(390\)](https://doi.org/10.1061/(ASCE)0733-9364(1998)124:5(390)).
- Hamzehei, S., O. Akbarzadeh, N. Fasihour, M. Alzgoool, and M. R. Khosravi. 2022. "Simulating a Charging Station for Electric Vehicles (EV) Based on the Concept of the Markov Chain to Analyze the System Performance."
- Hanna, A. S. 2001. *Quantifying the Cumulative Impact of Change Orders for Electrical and Mechanical Contractors: A Report to the Construction Industry Institute, the University of Texas at Austin, from the University of Wisconsin at Madison, Madison, Wisconsin*. Construction Industry Institute.
- Hanna, A. S. 2016. "Benchmark Performance Metrics for Integrated Project Delivery." *J. Constr. Eng. Manage.*, 142 (9): 04016040. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0001151](https://doi.org/10.1061/(ASCE)CO.1943-7862.0001151).
- Hanna, A. S., R. Camlic, P. A. Peterson, and E. V. Nordheim. 2002. "Quantitative Definition of Projects Impacted by Change Orders." *J. Constr. Eng. Manage.*, 128 (1): 57–64. [https://doi.org/10.1061/\(ASCE\)0733-9364\(2002\)128:1\(57\)](https://doi.org/10.1061/(ASCE)0733-9364(2002)128:1(57)).
- Hanna, A. S., and K. A. Iskandar. 2017. "Quantifying and Modeling the Cumulative Impact of Change Orders." *J. Constr. Eng. Manage.*, 143 (10): 04017076. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0001385](https://doi.org/10.1061/(ASCE)CO.1943-7862.0001385).
- Hanna, A. S., J. S. Russell, T. W. Gotzian, and E. V. Nordheim. 1999a. "Impact of Change Orders on Labor Efficiency for Mechanical Construction." *Journal of Construction Engineering and Management*, 125 (3): 176–184. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)0733-9364\(1999\)125:3\(176\)](https://doi.org/10.1061/(ASCE)0733-9364(1999)125:3(176)).

- Hanna, A. S., J. S. Russell, E. V. Nordheim, and M. J. Bruggink. 1999b. "Impact of Change Orders on Labor Efficiency for Electrical Construction." *J. Constr. Eng. Manage.*, 125 (4): 224–232. [https://doi.org/10.1061/\(ASCE\)0733-9364\(1999\)125:4\(224\)](https://doi.org/10.1061/(ASCE)0733-9364(1999)125:4(224)).
- Hao, Q., W. Shen, J. Neelamkavil, and R. Thomas. 2008. "Change management in construction projects." *the Proceedings of International Conference on Information Technology in Construction CIBW78*, 15–17.
- Hasanzadeh, S., B. Esmacili, S. Nasrollahi, G. M. Gad, and D. D. Gransberg. 2018. "Impact of Owners' Early Decisions on Project Performance and Dispute Occurrence in Public Highway Projects." *J. Leg. Aff. Dispute Resolut. Eng. Constr.*, 10 (2): 04518004. [https://doi.org/10.1061/\(ASCE\)LA.1943-4170.0000251](https://doi.org/10.1061/(ASCE)LA.1943-4170.0000251).
- Hegazy, T., E. Zanelidin, and D. Grierson. 2001. "Improving Design Coordination for Building Projects. I: Information Model." *J. Constr. Eng. Manage.*, 127 (4): 322–329. [https://doi.org/10.1061/\(ASCE\)0733-9364\(2001\)127:4\(322\)](https://doi.org/10.1061/(ASCE)0733-9364(2001)127:4(322)).
- Heravi, G., and M. H. Charkhakan. 2015. "Predicting Change by Evaluating the Change Implementation Process in Construction Projects Using Event Tree Analysis." *J. Manage. Eng.*, 31 (5): 04014081. [https://doi.org/10.1061/\(ASCE\)ME.1943-5479.0000325](https://doi.org/10.1061/(ASCE)ME.1943-5479.0000325).
- Hsieh, T., S. Lu, and C. Wu. 2004. "Statistical analysis of causes for change orders in metropolitan public works." *International Journal of Project Management*, 22 (8): 679–686. <https://doi.org/10.1016/j.ijproman.2004.03.005>.
- Ibbs, C. W., C. K. Wong, and Y. H. Kwak. 2001. "Project Change Management System." *J. Manage. Eng.*, 17 (3): 159–165. [https://doi.org/10.1061/\(ASCE\)0742-597X\(2001\)17:3\(159\)](https://doi.org/10.1061/(ASCE)0742-597X(2001)17:3(159)).
- Ibbs, W. 2005. "Impact of Change's Timing on Labor Productivity." *J. Constr. Eng. Manage.*, 131 (11): 1219–1223. [https://doi.org/10.1061/\(ASCE\)0733-9364\(2005\)131:11\(1219\)](https://doi.org/10.1061/(ASCE)0733-9364(2005)131:11(1219)).
- Ibbs, W. 2012. "Construction Change: Likelihood, Severity, and Impact on Productivity." *J. Leg. Aff. Dispute Resolut. Eng. Constr.*, 4 (3): 67–73. [https://doi.org/10.1061/\(ASCE\)LA.1943-4170.0000089](https://doi.org/10.1061/(ASCE)LA.1943-4170.0000089).
- Ibbs, W. 2021. "Update on Quantitative Analysis of Change and Loss of Productivity." *J. Leg. Aff. Dispute Resolut. Eng. Constr.*, 13 (1): 02520002. [https://doi.org/10.1061/\(ASCE\)LA.1943-4170.0000447](https://doi.org/10.1061/(ASCE)LA.1943-4170.0000447).
- Ibbs, W., and C. Chen. 2015. "Proactive Project Change-Prediction Tool." *J. Leg. Aff. Dispute Resolut. Eng. Constr.*, 7 (4): 04515003. [https://doi.org/10.1061/\(ASCE\)LA.1943-4170.0000175](https://doi.org/10.1061/(ASCE)LA.1943-4170.0000175).
- Ibbs, W., L. D. Nguyen, and S. Lee. 2007. "Quantified Impacts of Project Change." *J. Prof. Issues Eng. Educ. Pract.*, 133 (1): 45–52. [https://doi.org/10.1061/\(ASCE\)1052-3928\(2007\)133:1\(45\)](https://doi.org/10.1061/(ASCE)1052-3928(2007)133:1(45)).
- Isaac, S., and R. Navon. 2008. "Feasibility Study of an Automated Tool for Identifying the Implications of Changes in Construction Projects." *Journal of Construction Engineering and Management*, 134 (2): 139–145. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)0733-9364\(2008\)134:2\(139\)](https://doi.org/10.1061/(ASCE)0733-9364(2008)134:2(139)).
- Jiang, Y., M. Saito, and K. C. Sinha. 1988. "BRIDGE PERFORMANCE PREDICTION MODEL USING THE MARKOV CHAIN." *Transportation Research Record*, (1180).
- Karimidorabati, S., C. T. Haas, and J. Gray. 2016. "Evaluation of automation levels for construction change management." *ECAM*, 23 (5): 554–570. <https://doi.org/10.1108/ECAM-01-2015-0013>.

- Kechagias, G. A., A. C. Diamantidis, and T. D. Dimitrakos. 2024. "A semi-Markov decision model for the optimal control of an emergency medical service system." *International Journal of Industrial and Systems Engineering*, 46 (2): 169–194. <https://doi.org/10.1504/IJISE.2024.136414>.
- Khalafallah, A., and Y. Shalaby. 2019. "Change Orders: Automating Comparative Data Analysis and Controlling Impacts in Public Projects." *J. Constr. Eng. Manage.*, 145 (11): 04019064. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0001700](https://doi.org/10.1061/(ASCE)CO.1943-7862.0001700).
- Khanzadi, M., F. Nasirzadeh, and M. S. Dashti. 2018. "Fuzzy Cognitive Map Approach to Analyze Causes of Change Orders in Construction Projects." *J. Constr. Eng. Manage.*, 144 (2): 04017111. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0001430](https://doi.org/10.1061/(ASCE)CO.1943-7862.0001430).
- Kim, G.-H., S.-H. An, and K.-I. Kang. 2004. "Comparison of construction cost estimating models based on regression analysis, neural networks, and case-based reasoning." *Building and Environment*, 39 (10): 1235–1242. <https://doi.org/10.1016/j.buildenv.2004.02.013>.
- Kim, J. J., J. A. Miller, and S. Kim. 2020. "Cost Impacts of Change Orders due to Unforeseen Existing Conditions in Building Renovation Projects." *Journal of Construction Engineering and Management*, 146 (8): 04020094. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0001888](https://doi.org/10.1061/(ASCE)CO.1943-7862.0001888).
- Kim, J. W., G. Choi, J. C. Suh, and J. M. Lee. 2015. "Optimal scheduling of the maintenance and improvement for water main system using Markov decision process." 379–384.
- Kleiner, Y., R. Sadiq, and B. Rajani. 2006. "Modelling the deterioration of buried infrastructure as a fuzzy Markov process." *Journal of Water Supply: Research and Technology - AQUA*, 55 (2): 67–80. <https://doi.org/10.2166/aqua.2006.074>.
- Ko, T., J. Lee, and H. David Jeong. 2024. "Project Requirements Prioritization through NLP-Driven Classification and Adjusted Work Items Analysis." *J. Constr. Eng. Manage.*, 150 (3): 04023171. <https://doi.org/10.1061/JCEMD4.COENG-13655>.
- Lee, H.-R., and T. Lee. 2018. "Markov decision process model for patient admission decision at an emergency department under a surge demand." *Flexible Services and Manufacturing Journal*, 30 (1–2): 98–122. <https://doi.org/10.1007/s10696-017-9276-8>.
- Lee, S., F. Peña-Mora, and M. Park. 2005. "Quality and Change Management Model for Large Scale Concurrent Design and Construction Projects." *Journal of Construction Engineering and Management*, 131 (8): 890–902. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)0733-9364\(2005\)131:8\(890\)](https://doi.org/10.1061/(ASCE)0733-9364(2005)131:8(890)).
- Lee, S., S. Tae, N. Jee, and S. Shin. 2015. "LDA-Based Model for Measuring Impact of Change Orders in Apartment Projects and Its Application for Prerisk Assessment and Postevaluation." *J. Constr. Eng. Manage.*, 141 (7): 04015011. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0000971](https://doi.org/10.1061/(ASCE)CO.1943-7862.0000971).
- Leonard, C. A., P. Fazio, and O. Moselhi. 1988. "Construction Productivity: Major Causes of Impact." *American Association of Cost Engineers. Transactions of the American Association of Cost Engineers*, D.10.1. Morgantown, United States: American Association of Cost Engineers.
- Leu, S.-S., Y. Liu, and P.-L. Wu. 2023. "Project Cost Overrun Risk Prediction Using Hidden Markov Chain Analysis." *Buildings*, 13 (3). <https://doi.org/10.3390/buildings13030667>.
- Li, W., and C. Zhang. 2009. "Markov Chain Analysis." *International Encyclopedia of Human Geography (Second Edition)*, A. Kobayashi, ed., 407–412. Oxford: Elsevier.

- Love, P. E. D., Z. Irani, J. Smith, M. Regan, and J. Liu. 2017. “Cost performance of public infrastructure projects: the nemesis and nirvana of change-orders.” *Production Planning & Control*, 28 (13): 1081–1092. <https://doi.org/10.1080/09537287.2017.1333647>.
- Love, P. E. D., and H. Li. 2000. “Quantifying the causes and costs of rework in construction.” *Construction Management and Economics*, 18 (4): 479–490. <https://doi.org/10.1080/01446190050024897>.
- Love, P. E. D., C.-P. Sing, B. Carey, and J. T. Kim. 2015. “Estimating Construction Contingency: Accommodating the Potential for Cost Overruns in Road Construction Projects.” *Journal of Infrastructure Systems*, 21 (2): 04014035. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)IS.1943-555X.0000221](https://doi.org/10.1061/(ASCE)IS.1943-555X.0000221).
- Lu, C., and Q. Du. 2024. “The heterogeneous effects of transportation infrastructure on trade-embodied CO₂ transfer: A multi-scale perspective.” *Energy*, 302: 131785. <https://doi.org/10.1016/j.energy.2024.131785>.
- Lupu, L. 2015. “European stock markets correlations in a Markov switching framework.” *Romanian Journal of Economic Forecasting*, 18 (3): 103–119.
- Ma, Z., H. Duan, Z. Chen, X. Bin, and L. Jian. 2024. “Intelligent fault diagnosis of railway pantograph using a novel graph construction methodology.” *Meas. Sci. Technol.*, 35 (7): 076117. IOP Publishing. <https://doi.org/10.1088/1361-6501/ad34eb>.
- Mahmoudi, M., and H. Ghaneei. 2022. “Detection of structural regimes and analyzing the impact of crude oil market on Canadian stock market: Markov regime-switching approach.” *Studies in Economics and Finance*, 39 (4): 722–734. <https://doi.org/10.1108/SEF-09-2021-0352>.
- Mantawy, I. M., and N. L. C. Ravuri. 2024. “Predicting low-cycle fatigue-induced fracture in reinforcing bars: A CNN-based approach.” *Structures*, 64. <https://doi.org/10.1016/j.istruc.2024.106509>.
- Micevski, T., G. Kuczera, and P. Coombes. 2002. “Markov Model for Storm Water Pipe Deterioration.” *Journal of Infrastructure Systems*, 8 (2): 49–56. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)1076-0342\(2002\)8:2\(49\)](https://doi.org/10.1061/(ASCE)1076-0342(2002)8:2(49)).
- Mitropoulos, P., and C. B. Tatum. 1999. “Technology Adoption Decisions in Construction Organizations.” *J. Constr. Eng. Manage.*, 125 (5): 330–338. [https://doi.org/10.1061/\(ASCE\)0733-9364\(1999\)125:5\(330\)](https://doi.org/10.1061/(ASCE)0733-9364(1999)125:5(330)).
- Moayeri, V. 2017. “Design Change Management in Construction Projects Using Building Information Modeling (BIM).”
- Mokhtar, A., C. Bédard, and P. Fazio. 1998. “Information Model for Managing Design Changes in a Collaborative Environment.” *J. Comput. Civ. Eng.*, 12 (2): 82–92. [https://doi.org/10.1061/\(ASCE\)0887-3801\(1998\)12:2\(82\)](https://doi.org/10.1061/(ASCE)0887-3801(1998)12:2(82)).
- Morcous, G. 2006. “Performance Prediction of Bridge Deck Systems Using Markov Chains.” *Journal of Performance of Constructed Facilities*, 20 (2): 146–155. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)0887-3828\(2006\)20:2\(146\)](https://doi.org/10.1061/(ASCE)0887-3828(2006)20:2(146)).
- Morcous, G., Z. Lounis, and M. S. Mirza. 2003. “Identification of Environmental Categories for Markovian Deterioration Models of Bridge Decks.” *Journal of Bridge Engineering*, 8 (6): 353–361. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)1084-0702\(2003\)8:6\(353\)](https://doi.org/10.1061/(ASCE)1084-0702(2003)8:6(353)).
- Moselhi, O., I. Assem, and K. El-Rayes. 2005. “Change Orders Impact on Labor Productivity.” *Journal of Construction Engineering and Management*, 131 (3): 354–359. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)0733-9364\(2005\)131:3\(354\)](https://doi.org/10.1061/(ASCE)0733-9364(2005)131:3(354)).

- Moselhi, O., C. Leonard, and P. Fazio. 1991. "Impact of change orders on construction productivity." *Can. J. Civ. Eng.*, 18 (3): 484–492. <https://doi.org/10.1139/191-059>.
- Motawa, I. A., C. J. Anumba, and A. El-Hamalawi. 2006. "A fuzzy system for evaluating the risk of change in construction projects." *Advances in Engineering Software*, 37 (9): 583–591. <https://doi.org/10.1016/j.advengsoft.2006.01.006>.
- Motawa, I. A., C. J. Anumba, S. Lee, and F. Peña-Mora. 2007. "An integrated system for change management in construction." *Automation in Construction*, 16 (3): 368–377. <https://doi.org/10.1016/j.autcon.2006.07.005>.
- Nabipour, N., A. Martinez, and M. Nik-Bakht. 2023. "Predicting Change Order Magnitude in Construction Projects—A Machine Learning Approach." *Proceedings of the Canadian Society for Civil Engineering Annual Conference 2023, Volume 4*, S. Desjardins, G. J. Poitras, and M. Nik-Bakht, eds., 137–150. Cham: Springer Nature Switzerland.
- Naji, K. K., M. Gunduz, and A. F. Naser. 2022. "The Effect of Change-Order Management Factors on Construction Project Success: A Structural Equation Modeling Approach." *J. Constr. Eng. Manage.*, 148 (9): 04022085. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0002350](https://doi.org/10.1061/(ASCE)CO.1943-7862.0002350).
- Nik Bakht, M., and T. E. El-Diraby. 2015. "Synthesis of Decision-Making Research in Construction." *Journal of Construction Engineering and Management*, 141 (9): 04015027. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0000984](https://doi.org/10.1061/(ASCE)CO.1943-7862.0000984).
- Ossai, C. I., B. Boswell, and I. J. Davies. 2016. "Application of Markov modelling and Monte Carlo simulation technique in failure probability estimation - A consideration of corrosion defects of internally corroded pipelines." *Engineering Failure Analysis*, 68: 159–171. <https://doi.org/10.1016/j.engfailanal.2016.06.004>.
- Padala, S. P. S., and J. U. Maheswari. 2023. "Modeling a construction project in a matrix-based framework for managing requirement changes." *International Journal of Construction Management*, 23 (14): 2369–2390. <https://doi.org/10.1080/15623599.2022.2059739>.
- Padala, S. P. S., J. U. Maheswari, and H. Hirani. 2020. "Identification and classification of change causes and effects in construction projects." *International Journal of Construction Management*, 1–20. <https://doi.org/10.1080/15623599.2020.1827186>.
- Padala, S. P. S., J. U. Maheswari, and H. Hirani. 2022. "Identification and classification of change causes and effects in construction projects." *International Journal of Construction Management*, 22 (14): 2788–2807. <https://doi.org/10.1080/15623599.2020.1827186>.
- Qiao, S., and B. Liu. 2012. "Prediction of Ground Displacement and Deformation Induced by Dewatering of Groundwater." 73–79. American Society of Civil Engineers. [https://doi.org/10.1061/40867\(199\)7](https://doi.org/10.1061/40867(199)7).
- Rahmani Mirshekarlou, B. 2012. "A Taxonomy for causes of changes in construction." Master Thesis. Middle East Technical University.
- Rathnayake, A., and C. Middleton. 2023. "Systematic Review of the Literature on Construction Productivity." *Journal of Construction Engineering and Management*, 149 (6). <https://doi.org/10.1061/JCEMD4.COENG-13045>.
- "Request for Change Orders: Adjusting the Contract Scope." 2024. *Procore*. Accessed September 8, 2024. <https://www.procore.com/library/request-for-change-orders>.
- Riley, D. R., B. E. Diller, and D. Kerr. 2005. "Effects of Delivery Systems on Change Order Size and Frequency in Mechanical Construction." *J. Constr. Eng. Manage.*, 131 (9): 953–962. [https://doi.org/10.1061/\(ASCE\)0733-9364\(2005\)131:9\(953\)](https://doi.org/10.1061/(ASCE)0733-9364(2005)131:9(953)).

- Rojas, E. M., and I. Kell. 2008. "Comparative Analysis of Project Delivery Systems Cost Performance in Pacific Northwest Public Schools." *J. Constr. Eng. Manage.*, 134 (6): 387–397. [https://doi.org/10.1061/\(ASCE\)0733-9364\(2008\)134:6\(387\)](https://doi.org/10.1061/(ASCE)0733-9364(2008)134:6(387)).
- Sardroud, J. M. 2015. "Perceptions of automated data collection technology use in the construction industry." *Journal of Civil Engineering and Management*, 21 (1): 54–66. Taylor & Francis. <https://doi.org/10.3846/13923730.2013.802734>.
- Scherer, W. T., and D. M. Glagola. 1994. "Markovian Models for Bridge Maintenance Management." *Journal of Transportation Engineering*, 120 (1): 37–51. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)0733-947X\(1994\)120:1\(37\)](https://doi.org/10.1061/(ASCE)0733-947X(1994)120:1(37)).
- Sempewo, J. I., and L. Kyokaali. 2019. "Comparative performance of regression and the markov based approach in the prediction of the future condition of a water distribution pipe network amidst data scarce situations: A case study of Kampala water, Uganda." *Water Practice and Technology*, 14 (4): 946–958. <https://doi.org/10.2166/wpt.2019.075>.
- Semple, C., F. T. Hartman, and G. Jergeas. 1994. "Construction Claims and Disputes: Causes and Cost/Time Overruns." *Journal of Construction Engineering and Management*, 120 (4): 785–795. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)0733-9364\(1994\)120:4\(785\)](https://doi.org/10.1061/(ASCE)0733-9364(1994)120:4(785)).
- Senouci, A., A. Alsarraj, M. Gunduz, and N. Eldin. 2017. "Analysis of change orders in Qatari construction projects." *International Journal of Construction Management*, 17 (4): 280–292. <https://doi.org/10.1080/15623599.2016.1211973>.
- Serag, E., A. Oloufa, L. Malone, and E. Radwan. 2010. "Model for Quantifying the Impact of Change Orders on Project Cost for U.S. Roadwork Construction." *Journal of Construction Engineering and Management*, 136 (9): 1015–1027. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0000206](https://doi.org/10.1061/(ASCE)CO.1943-7862.0000206).
- Shabani Ardakani, S., and M. Nik-Bakht. 2021. "Functional Evaluation of Change Order and Invoice Management Processes under Different Procurement Strategies: Social Network Analysis Approach." *Journal of Construction Engineering and Management*, 147 (1): 04020155. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0001974](https://doi.org/10.1061/(ASCE)CO.1943-7862.0001974).
- Shafaat, A., F. Marbouti, and T. Mahfouz. 2022. "Early warning system for highway construction projects using GA-SVM." *International Journal of Construction Management*, 0 (0): 1–10. Taylor & Francis. <https://doi.org/10.1080/15623599.2022.2056805>.
- Shi, S., N. Lin, Y. Zhang, J. Cheng, C. Huang, L. Liu, and B. Lu. 2016. "Research on Markov property analysis of driving cycles and its application." *Transportation Research Part D: Transport and Environment*, 47: 171–181. <https://doi.org/10.1016/j.trd.2016.05.013>.
- Shrestha, B., P. P. Shrestha, R. Maharjan, and D. Gransberg. 2022. "Cost, Change Order, and Schedule Performance of Highway Projects." *J. Leg. Aff. Dispute Resolut. Eng. Constr.*, 14 (1): 04521044. [https://doi.org/10.1061/\(ASCE\)LA.1943-4170.0000523](https://doi.org/10.1061/(ASCE)LA.1943-4170.0000523).
- Shrestha, K., P. P. Shrestha, and M. Lohoud. 2023. "Critical Factors Affecting the Bid Cost of Building Construction Projects." *Journal of Legal Affairs and Dispute Resolution in Engineering and Construction*, 15 (3): 04523018. American Society of Civil Engineers. <https://doi.org/10.1061/JLADAH.LADR-960>.
- Shrestha, P. P., and M. Fathi. 2019. "Impacts of Change Orders on Cost and Schedule Performance and the Correlation with Project Size of DB Building Projects." *J. Leg. Aff. Dispute Resolut. Eng. Constr.*, 11 (3): 04519010. [https://doi.org/10.1061/\(ASCE\)LA.1943-4170.0000311](https://doi.org/10.1061/(ASCE)LA.1943-4170.0000311).

- Shrestha, P. P., and J. D. Fernane. 2017. "Performance of Design-Build and Design-Bid-Build Projects for Public Universities." *J. Constr. Eng. Manage.*, 143 (3): 04016101. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0001241](https://doi.org/10.1061/(ASCE)CO.1943-7862.0001241).
- Shrestha, P. P., and R. Maharjan. 2018. "Effects of Change Orders on Cost Growth, Schedule Growth, and Construction Intensity of Large Highway Projects." *J. Leg. Aff. Dispute Resolut. Eng. Constr.*, 10 (3): 04518012. [https://doi.org/10.1061/\(ASCE\)LA.1943-4170.0000264](https://doi.org/10.1061/(ASCE)LA.1943-4170.0000264).
- Shrestha, P. P., and R. Maharjan. 2019. "Effect of Change Orders on Cost and Schedule for Small Low-Bid Highway Contracts." *J. Leg. Aff. Dispute Resolut. Eng. Constr.*, 11 (4): 04519025. [https://doi.org/10.1061/\(ASCE\)LA.1943-4170.0000323](https://doi.org/10.1061/(ASCE)LA.1943-4170.0000323).
- Shrestha, P. P., J. T. O'Connor, and G. E. Gibson. 2012. "Performance Comparison of Large Design-Build and Design-Bid-Build Highway Projects." *J. Constr. Eng. Manage.*, 138 (1): 1–13. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0000390](https://doi.org/10.1061/(ASCE)CO.1943-7862.0000390).
- Shrestha, P. P., K. K. Shrestha, and T. K. Kandie. 2017. "Effects of Change Orders on the Cost and Schedule of Rural Road Maintenance Projects." *J. Leg. Aff. Dispute Resolut. Eng. Constr.*, 9 (3): 04517010. [https://doi.org/10.1061/\(ASCE\)LA.1943-4170.0000227](https://doi.org/10.1061/(ASCE)LA.1943-4170.0000227).
- Shrestha, P. P., K. K. Shrestha, and H. B. Zeleke. 2019. "Probability of change orders and the effect on cost and schedule for new public school buildings." *ECAM*, 26 (6): 1087–1104. <https://doi.org/10.1108/ECAM-01-2018-0017>.
- Smola, A. J., and B. Schölkopf. 2004. "A tutorial on support vector regression." *Statistics and Computing*, 14 (3): 199–222. <https://doi.org/10.1023/B:STCO.0000035301.49549.88>.
- Soofi, E. S. 1990. "Effects of collinearity on information about regression coefficients." *Journal of Econometrics*, 43 (3): 255–274. [https://doi.org/10.1016/0304-4076\(90\)90120-I](https://doi.org/10.1016/0304-4076(90)90120-I).
- Stamatiou, D. R. I., K. A. Kirytopoulos, S. T. Ponis, S. Gayialis, and I. Tatsiopoulou. 2019. "A process reference model for claims management in construction supply chains: the contractors' perspective." *International Journal of Construction Management*, 19 (5): 382–400. Taylor & Francis. <https://doi.org/10.1080/15623599.2018.1452100>.
- Sun, M., A. Fleming, S. Senaratne, I. Motawa, and M. L. Yeoh. 2006. "A Change Management Toolkit for Construction Projects." *Architectural Engineering and Design Management*, 2 (4): 261–271. Taylor & Francis. <https://doi.org/10.1080/17452007.2006.9684621>.
- Sun, M., and X. Meng. 2009. "Taxonomy for change causes and effects in construction projects." *International Journal of Project Management*, 27 (6): 560–572. <https://doi.org/10.1016/j.ijproman.2008.10.005>.
- Tao, W., P. Lin, and N. Wang. 2021a. "Optimum life-cycle maintenance strategies of deteriorating highway bridges subject to seismic hazard by a hybrid Markov decision process model." *Structural Safety*, 89. <https://doi.org/10.1016/j.strusafe.2020.102042>.
- Tao, W., N. Wang, B. Ellingwood, and C. Nicholson. 2021b. "Enhancing bridge performance following earthquakes using Markov decision process." *Structure and Infrastructure Engineering*, 17 (1): 62–73. <https://doi.org/10.1080/15732479.2020.1730410>.
- Thomas, H. R., and C. L. Napolitan. 1995. "Quantitative Effects of Construction Changes on Labor Productivity." *Journal of Construction Engineering and Management*, 121 (3): 290–296. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)0733-9364\(1995\)121:3\(290\)](https://doi.org/10.1061/(ASCE)0733-9364(1995)121:3(290)).
- Thomas, O., and J. Sobanjo. 2013. "Comparison of Markov Chain and Semi-Markov Models for Crack Deterioration on Flexible Pavements." *Journal of Infrastructure Systems*, 19 (2):

- 186–195. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)IS.1943-555X.0000112](https://doi.org/10.1061/(ASCE)IS.1943-555X.0000112).
- Tian, J., M. Yang, and S. Li. 2012. “Research on Coordinated Supply-Production-Distribution Plan Model in Supply Chain under Stochastic Demand.” 1001–1006. American Society of Civil Engineers. [https://doi.org/10.1061/40996\(330\)143](https://doi.org/10.1061/40996(330)143).
- Todorov, M. G. 2024. “A New Approach to the Energy-to-Peak Performance Analysis of Continuous-Time Markov Jump Linear Systems.” *IEEE Control Systems Letters*, 1–1. <https://doi.org/10.1109/LCSYS.2024.3406492>.
- Toomey, D. E., J. S. Marks, and D. T. Zhao. 2015. “Calculating Lost Labor Productivity: Is There a Better Way?” 35 (2).
- Touran, A. 1197. “Probabilistic Model for Tunneling Project Using Markov Chain.” Accessed June 13, 2024. <https://ascelibrary.org/doi/10.1061/%28ASCE%290733-9364%281997%29123%3A4%28444%29>.
- Trost, S. M., and G. D. Oberlender. 2003. “Predicting Accuracy of Early Cost Estimates Using Factor Analysis and Multivariate Regression.” *J. Constr. Eng. Manage.*, 129 (2): 198–204. [https://doi.org/10.1061/\(ASCE\)0733-9364\(2003\)129:2\(198\)](https://doi.org/10.1061/(ASCE)0733-9364(2003)129:2(198)).
- Wang, J., S. Zhang, P. Fenn, X. Luo, Y. Liu, and L. Zhao. 2023. “Adopting BIM to Facilitate Dispute Management in the Construction Industry: A Conceptual Framework Development.” *J. Constr. Eng. Manage.*, 149 (1): 03122010. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0002419](https://doi.org/10.1061/(ASCE)CO.1943-7862.0002419).
- Wang, J., M. Zhou, X. Guo, L. Qi, and X. Wang. 2021. “A Markov regime switching model for asset pricing and ambiguity measurement of stock market.” *Neurocomputing*, 435: 283–294. <https://doi.org/10.1016/j.neucom.2020.12.103>.
- Wang, L., J. Lee, J. Nimawat, K. Han, and A. Gupta. 2024. “Integrated 4D Design Change Management Model for Construction Projects.” *J. Constr. Eng. Manage.*, 150 (5): 04024023. <https://doi.org/10.1061/JCEMD4.COENG-14246>.
- Wasiq, S., and A. Golroo. 2024. “Probabilistic pavement performance modeling using hybrid Markov Chain: A case study in Afghanistan.” *Case Studies in Construction Materials*, 20. <https://doi.org/10.1016/j.cscm.2024.e03023>.
- Wei, B., C. Guo, and M. Deng. 2022. “An Innovation of the Markov Probability Model for Predicting the Remaining Service Life of Civil Airport Rigid Pavements.” *Materials*, 15 (17). <https://doi.org/10.3390/ma15176082>.
- Williams, B. M. n.d. “Modeling and forecasting vehicular traffic flow as a seasonal stochastic time series process.” Ph.D. United States -- Virginia: University of Virginia.
- Williams, T. P., and J. Gong. 2014. “Predicting construction cost overruns using text mining, numerical data and ensemble classifiers.” *Automation in Construction*, 43: 23–29. <https://doi.org/10.1016/j.autcon.2014.02.014>.
- Wirth, R., and J. Hipp. 2000. “CRISP-DM: Towards a Standard Process Model for Data Mining.”
- Wu, C., T. Hsieh, and W. Cheng. 2005. “Statistical analysis of causes for design change in highway construction on Taiwan.” *International Journal of Project Management*, 23 (7): 554–563. <https://doi.org/10.1016/j.ijproman.2004.07.010>.
- Yan, H., N. Yang, Y. Peng, and Y. Ren. 2020. “Data mining in the construction industry: Present status, opportunities, and future trends.” *Automation in Construction*, 119: 103331. <https://doi.org/10.1016/j.autcon.2020.103331>.
- Yang, G., L. Tian, G. Tang, J. Mao, and Y. Du. 2024. “Research on Bridge Performance Degradation Prediction Based on Combination of the D-S Theory and the Markov Chain.”

- Applied Mathematics and Mechanics*, 45 (4): 416–428. <https://doi.org/10.21656/1000-0887.440343>.
- Yang, H., V. N. Nair, J. Chen, and A. Sudjianto. 2019. “Assessing Markov property in multistate transition models with applications to credit risk modeling.” *Appl Stoch Models Bus & Ind*, 35 (3): 552–570. <https://doi.org/10.1002/asmb.2336>.
- Yap, J. B. H., C. G. Y. Lam, M. Skitmore, and N. Talebian. 2022. “Barriers to the adoption of new safety technologies in construction: a developing country context.” *Journal of Civil Engineering and Management*, 28 (2): 120–133. <https://doi.org/10.3846/jcem.2022.16014>.
- Younes, B., A. Bouferguène, M. Al-Hussein, and H. Yu. 2015. “Overdue Invoice Management: Markov Chain Approach.” *Journal of Construction Engineering and Management*, 141 (1): 04014062. American Society of Civil Engineers. [https://doi.org/10.1061/\(ASCE\)CO.1943-7862.0000913](https://doi.org/10.1061/(ASCE)CO.1943-7862.0000913).
- Yu, F.-H., J. Lu, J.-W. Gu, and W.-K. Ching. 2019. “Modeling Credit Risk with Hidden Markov Default Intensity.” *Computational Economics*, 54 (3): 1213–1229. <https://doi.org/10.1007/s10614-018-9869-7>.
- Zhang, J., Z. Wang, P. Liu, Z. Zhang, X. Li, and C. Qu. 2019a. “Driving cycles construction for electric vehicles considering road environment: A case study in Beijing.” *Applied Energy*, 253: 113514. <https://doi.org/10.1016/j.apenergy.2019.113514>.
- Zhang, M., S. Shi, N. Lin, and B. Yue. 2019b. “High-Efficiency Driving Cycle Generation Using a Markov Chain Evolution Algorithm.” *IEEE Trans. Veh. Technol.*, 68 (2): 1288–1301. <https://doi.org/10.1109/TVT.2018.2887063>.
- Zhao, T. 2021. “Modeling with Functions for Cumulative Impact of Changes.” *Journal of Legal Affairs and Dispute Resolution in Engineering and Construction*, 13 (4). [https://doi.org/10.1061/\(ASCE\)LA.1943-4170.0000492](https://doi.org/10.1061/(ASCE)LA.1943-4170.0000492).
- Zhao, T. 2023. “Timing and Severity of Cumulative Impact of Changes on Labor Productivity.” *J. Leg. Aff. Dispute Resolut. Eng. Constr.*, 15 (1): 04522039. [https://doi.org/10.1061/\(ASCE\)LA.1943-4170.0000586](https://doi.org/10.1061/(ASCE)LA.1943-4170.0000586).

Appendices

Appendix 1) 1-step transition matrices for SCI and CCI tests

1-step transition matrices for SCI test

(a) PD-1

P(1,2)			
	1	2	3
1	0.94	0.5	0.01
2	0.77	0.23	0
3	1	0	0

P(2,3)			
	1	2	3
1	0.93	0.05	0.02
2	0.62	0.38	0
3	0.5	0.5	0

P(3,4)			
	1	2	3
1	0.91	0.07	0.02
2	0.68	0.32	0
3	0.57	0.14	0.29

P(4,5)			
	1	2	3
1	0.55	0.27	0.18
2	0.34	0.5	0.16
3	0.3	0.4	0.3

(b) PD-2

P(1,2)			
	1	2	3
1	0.8	0.17	0.03
2	0.79	0.21	0
3	0	1	0

P(2,3)			
	1	2	3
1	0.82	0.16	0.02
2	0.52	0.4	0.07
3	0.6	0.1	0.3

P(3,4)			
	1	2	3
1	0.79	0.18	0.03
2	0.36	0.59	0.05
3	0.46	0.18	0.36

P(4,5)			
	1	2	3
1	0.58	0.26	0.16
2	0.20	0.61	0.19
3	0.53	0.27	0.2

(c) PD-3

P(1,2)			
	1	2	3
1	0.78	0.19	0.03
2	0.39	0.55	0.06
3	0.75	0.25	0

P(2,3)			
	1	2	3
1	0.80	0.17	0.03
2	0.36	0.6	0.04
3	0.33	0.33	0.33

P(3,4)			
	1	2	3
1	0.78	0.16	0.06
2	0.36	0.61	0.03
3	0.15	0.54	0.31

P(4,5)			
	1	2	3
1	0.76	0.15	0.09
2	0.30	0.54	0.16
3	0.39	0.44	0.17

(d) PD-0

P(1,2)			
	1	2	3
1	0.85	0.13	0.02
2	0.55	0.42	0.03
3	0.71	0.29	0

P(2,3)			
	1	2	3
1	0.87	0.11	0.02
2	0.46	0.49	0.05
3	0.48	0.24	0.28

P(3,4)			
	1	2	3
1	0.84	0.12	0.04
2	0.41	0.56	0.03
3	0.36	0.32	0.32

P(4,5)			
	1	2	3
1	0.61	0.24	0.15
2	0.27	0.56	0.17
3	0.42	0.37	0.21

1-step transition matrices for CCI test

(a) PD-1

P(1,2)			
	1	2	3
1	0.94	0.5	0.01
2	0	0.92	0.08
3	0	0	1

P(2,3)			
	1	2	3
1	0.93	0.05	0.02
2	0.03	0.85	0.12
3	0	0.14	0.86

P(3,4)			
	1	2	3
1	0.9	0.08	0.02
2	0	0.98	0.02
3	0	0	1

P(4,5)			
	1	2	3
1	0.58	0.24	0.18
2	0.07	0.7	0.23
3	0.04	0.04	0.92

(b) PD-2

P(1,2)			
	1	2	3
1	0.81	0.16	0.03
2	0	1	0
3	0	0	1

P(2,3)			
	1	2	3
1	0.82	0.16	0.02
2	0	0.92	0.08
3	0	0	1

P(3,4)			
	1	2	3
1	0.81	0.15	0.04
2	0.03	0.93	0.04
3	0.1	0	0.9

P(4,5)			
	1	2	3
1	0.62	0.26	0.12
2	0.02	0.75	0.23
3	0	0	1

(c) PD-3

P(1,2)			
	1	2	3
1	0.8	0.17	0.03
2	0.05	0.87	0.08
3	0	0.13	0.87

P(2,3)			
	1	2	3
1	0.82	0.14	0.03
2	0.05	0.88	0.07
3	0	0	1

P(3,4)			
	1	2	3
1	0.77	0.16	0.07
2	0.01	0.96	0.03
3	0	0	1

P(4,5)			
	1	2	3
1	0.74	0.15	0.11
2	0.03	0.8	0.17
3	0.02	0	0.98

(d) PD-0

P(1,2)			
	1	2	3
1	0.86	0.12	0.02
2	0.03	0.91	0.06
3	0	0.07	0.92

P(2,3)			
	1	2	3
1	0.88	0.1	0.02
2	0.03	0.89	0.08
3	0	0.03	0.97

P(3,4)			
	1	2	3
1	0.84	0.12	0.04
2	0.02	0.95	0.03
3	0.03	0	0.97

P(4,5)			
	1	2	3
1	0.62	0.23	0.15
2	0.03	0.76	0.21
3	0.42	0.37	0.21

Appendix 2: Observed and expected 4-step transition matrices

None-Cumulative test

Observed Transition Matrix	Expected Transition Matrix	CHI-Squared	P-Value																																
(a) PD-1																																			
<table border="1"> <tr><td></td><td>1</td><td>2</td><td>3</td></tr> <tr><td>1</td><td>0.53</td><td>0.29</td><td>0.18</td></tr> <tr><td>2</td><td>0.46</td><td>0.38</td><td>0.15</td></tr> <tr><td>3</td><td>0.5</td><td>0.25</td><td>0.25</td></tr> </table>		1	2	3	1	0.53	0.29	0.18	2	0.46	0.38	0.15	3	0.5	0.25	0.25	<table border="1"> <tr><td></td><td>1</td><td>2</td><td>3</td></tr> <tr><td>1</td><td>0.52</td><td>0.29</td><td>0.18</td></tr> <tr><td>2</td><td>0.52</td><td>0.30</td><td>0.18</td></tr> <tr><td>3</td><td>0.52</td><td>0.3</td><td>0.18</td></tr> </table>		1	2	3	1	0.52	0.29	0.18	2	0.52	0.30	0.18	3	0.52	0.3	0.18	0.615	0.99
	1	2	3																																
1	0.53	0.29	0.18																																
2	0.46	0.38	0.15																																
3	0.5	0.25	0.25																																
	1	2	3																																
1	0.52	0.29	0.18																																
2	0.52	0.30	0.18																																
3	0.52	0.3	0.18																																
(b) PD-2																																			
<table border="1"> <tr><td></td><td>1</td><td>2</td><td>3</td></tr> <tr><td>1</td><td>0.48</td><td>0.35</td><td>0.17</td></tr> <tr><td>2</td><td>0.57</td><td>0.43</td><td>0</td></tr> <tr><td>3</td><td>0</td><td>0</td><td>1</td></tr> </table>		1	2	3	1	0.48	0.35	0.17	2	0.57	0.43	0	3	0	0	1	<table border="1"> <tr><td></td><td>1</td><td>2</td><td>3</td></tr> <tr><td>1</td><td>0.48</td><td>0.35</td><td>0.17</td></tr> <tr><td>2</td><td>0.48</td><td>0.35</td><td>0.17</td></tr> <tr><td>3</td><td>0.45</td><td>0.38</td><td>0.17</td></tr> </table>		1	2	3	1	0.48	0.35	0.17	2	0.48	0.35	0.17	3	0.45	0.38	0.17	12.54	0.051
	1	2	3																																
1	0.48	0.35	0.17																																
2	0.57	0.43	0																																
3	0	0	1																																
	1	2	3																																
1	0.48	0.35	0.17																																
2	0.48	0.35	0.17																																
3	0.45	0.38	0.17																																
(c) PD-3																																			
<table border="1"> <tr><td></td><td>1</td><td>2</td><td>3</td></tr> <tr><td>1</td><td>0.57</td><td>0.30</td><td>0.13</td></tr> <tr><td>2</td><td>0.68</td><td>0.26</td><td>0.05</td></tr> <tr><td>3</td><td>0.87</td><td>0.13</td><td>0.0</td></tr> </table>		1	2	3	1	0.57	0.30	0.13	2	0.68	0.26	0.05	3	0.87	0.13	0.0	<table border="1"> <tr><td></td><td>1</td><td>2</td><td>3</td></tr> <tr><td>1</td><td>0.6</td><td>0.28</td><td>0.12</td></tr> <tr><td>2</td><td>0.56</td><td>0.31</td><td>0.12</td></tr> <tr><td>3</td><td>0.6</td><td>0.28</td><td>0.11</td></tr> </table>		1	2	3	1	0.6	0.28	0.12	2	0.56	0.31	0.12	3	0.6	0.28	0.11	6.34	0.385
	1	2	3																																
1	0.57	0.30	0.13																																
2	0.68	0.26	0.05																																
3	0.87	0.13	0.0																																
	1	2	3																																
1	0.6	0.28	0.12																																
2	0.56	0.31	0.12																																
3	0.6	0.28	0.11																																
(d) PD-0																																			
<table border="1"> <tr><td></td><td>1</td><td>2</td><td>3</td></tr> <tr><td>1</td><td>0.52</td><td>0.31</td><td>0.17</td></tr> <tr><td>2</td><td>0.62</td><td>0.32</td><td>0.06</td></tr> <tr><td>3</td><td>0.64</td><td>0.14</td><td>0.22</td></tr> </table>		1	2	3	1	0.52	0.31	0.17	2	0.62	0.32	0.06	3	0.64	0.14	0.22	<table border="1"> <tr><td></td><td>1</td><td>2</td><td>3</td></tr> <tr><td>1</td><td>0.53</td><td>0.31</td><td>0.16</td></tr> <tr><td>2</td><td>0.51</td><td>0.33</td><td>0.16</td></tr> <tr><td>3</td><td>0.52</td><td>0.32</td><td>0.16</td></tr> </table>		1	2	3	1	0.53	0.31	0.16	2	0.51	0.33	0.16	3	0.52	0.32	0.16	7.59	0.27
	1	2	3																																
1	0.52	0.31	0.17																																
2	0.62	0.32	0.06																																
3	0.64	0.14	0.22																																
	1	2	3																																
1	0.53	0.31	0.16																																
2	0.51	0.33	0.16																																
3	0.52	0.32	0.16																																

Cumulative test

Observed Transition Matrix	Expected Transition Matrix	CHI-Squared	P-Value																																
(a) PD-1																																			
<table border="1"> <tr><td></td><td>1</td><td>2</td><td>3</td></tr> <tr><td>1</td><td>0.47</td><td>0.3</td><td>0.23</td></tr> <tr><td>2</td><td>0</td><td>0.77</td><td>0.23</td></tr> <tr><td>3</td><td>0</td><td>0</td><td>1</td></tr> </table>		1	2	3	1	0.47	0.3	0.23	2	0	0.77	0.23	3	0	0	1	<table border="1"> <tr><td></td><td>1</td><td>2</td><td>3</td></tr> <tr><td>1</td><td>0.47</td><td>0.3</td><td>0.23</td></tr> <tr><td>2</td><td>0.07</td><td>0.57</td><td>0.36</td></tr> <tr><td>3</td><td>0.04</td><td>0.13</td><td>0.83</td></tr> </table>		1	2	3	1	0.47	0.3	0.23	2	0.07	0.57	0.36	3	0.04	0.13	0.83	3.45	0.75
	1	2	3																																
1	0.47	0.3	0.23																																
2	0	0.77	0.23																																
3	0	0	1																																
	1	2	3																																
1	0.47	0.3	0.23																																
2	0.07	0.57	0.36																																
3	0.04	0.13	0.83																																
(b) PD-2																																			
<table border="1"> <tr><td></td><td>1</td><td>2</td><td>3</td></tr> <tr><td>1</td><td>0.35</td><td>0.4</td><td>0.25</td></tr> <tr><td>2</td><td>0.07</td><td>0.86</td><td>0.07</td></tr> <tr><td>3</td><td>0</td><td>0</td><td>1</td></tr> </table>		1	2	3	1	0.35	0.4	0.25	2	0.07	0.86	0.07	3	0	0	1	<table border="1"> <tr><td></td><td>1</td><td>2</td><td>3</td></tr> <tr><td>1</td><td>0.35</td><td>0.41</td><td>0.24</td></tr> <tr><td>2</td><td>0.04</td><td>0.65</td><td>0.31</td></tr> <tr><td>3</td><td>0.06</td><td>0.03</td><td>0.91</td></tr> </table>		1	2	3	1	0.35	0.41	0.24	2	0.04	0.65	0.31	3	0.06	0.03	0.91	4.39	0.62
	1	2	3																																
1	0.35	0.4	0.25																																
2	0.07	0.86	0.07																																
3	0	0	1																																
	1	2	3																																
1	0.35	0.41	0.24																																
2	0.04	0.65	0.31																																
3	0.06	0.03	0.91																																
(c) PD-3																																			
<table border="1"> <tr><td></td><td>1</td><td>2</td><td>3</td></tr> <tr><td>1</td><td>0.41</td><td>0.34</td><td>0.25</td></tr> <tr><td>2</td><td>0</td><td>0.74</td><td>0.26</td></tr> <tr><td>3</td><td>0</td><td>0.13</td><td>0.87</td></tr> </table>		1	2	3	1	0.41	0.34	0.25	2	0	0.74	0.26	3	0	0.13	0.87	<table border="1"> <tr><td></td><td>1</td><td>2</td><td>3</td></tr> <tr><td>1</td><td>0.40</td><td>0.36</td><td>0.24</td></tr> <tr><td>2</td><td>0.08</td><td>0.62</td><td>0.3</td></tr> <tr><td>3</td><td>0.03</td><td>0.09</td><td>0.88</td></tr> </table>		1	2	3	1	0.40	0.36	0.24	2	0.08	0.62	0.3	3	0.03	0.09	0.88	4.89	0.56
	1	2	3																																
1	0.41	0.34	0.25																																
2	0	0.74	0.26																																
3	0	0.13	0.87																																
	1	2	3																																
1	0.40	0.36	0.24																																
2	0.08	0.62	0.3																																
3	0.03	0.09	0.88																																
(d) PD-0																																			
<table border="1"> <tr><td></td><td>1</td><td>2</td><td>3</td></tr> <tr><td>1</td><td>0.42</td><td>0.34</td><td>0.24</td></tr> <tr><td>2</td><td>0.02</td><td>0.77</td><td>0.21</td></tr> <tr><td>3</td><td>0</td><td>0.07</td><td>0.93</td></tr> </table>		1	2	3	1	0.42	0.34	0.24	2	0.02	0.77	0.21	3	0	0.07	0.93	<table border="1"> <tr><td></td><td>1</td><td>2</td><td>3</td></tr> <tr><td>1</td><td>0.41</td><td>0.35</td><td>0.24</td></tr> <tr><td>2</td><td>0.07</td><td>0.6</td><td>0.33</td></tr> <tr><td>3</td><td>0.04</td><td>0.08</td><td>0.88</td></tr> </table>		1	2	3	1	0.41	0.35	0.24	2	0.07	0.6	0.33	3	0.04	0.08	0.88	9.05	0.17
	1	2	3																																
1	0.42	0.34	0.24																																
2	0.02	0.77	0.21																																
3	0	0.07	0.93																																
	1	2	3																																
1	0.41	0.35	0.24																																
2	0.07	0.6	0.33																																
3	0.04	0.08	0.88																																

Appendix 3: Data sample

Change Orders:

ProjectId	DateCreated	Change Order Type	Description	Amount	co_Status
1001	2022-04-27	PrimeContractChangeOrder	mno 006 ot insulation work april 29th	760	3_approved
1001	2022-01-26	PrimeContractChangeOrder	mno #02 demobilization remobilization	1260	3_approved
1002	2022-04-27	PrimeContractChangeOrder	mno 006 ot insulation work april 29th	1440	3_approved
1002	2022-05-25	PrimeContractChangeOrder	mno 007 ot insulation work april 29th credit	-1912	3_approved
1002	2022-01-26	PrimeContractChangeOrder	mno #02 demobilization remobilization	2475	3_approved
1002	2022-04-21	PrimeContractChangeOrder	mno #05 insulation	8640	3_approved

Projects

Project Id	Parent Project Id	City	Classification	Contract Value	Billing Type	Expected StartDate	Expected End Date	Operating Unit	Province	Project Type
1001	NULL	SCARBO ROUGH	COMM. OFFICE.LOW	684755	Fixed Fee	August 5, 2021	Dec 31, 2021	Special Projects	ON	Construction
1001	NULL	Toronto	COMM. OFFICE.HIGH	1600000	Fixed Fee	August 25, 2021	July 28, 2022	Special Projects	ON	Construction
1002	1008	Belleville	INDUST. BUILD	159532	Fixed Fee	May 24, 2021	Aug 31, 2021	Electrical	ON	Construction
1002	NULL	Toronto	COMM. OFFICE.LOW	92500	Fixed Fee	March 1, 2022	June 30, 2022	Special Projects	ON	Construction
1002	NULL	Newtonville	INDUST. WASTE	0	Fixed Fee	October 1, 2016	Sep 30, 2021	Electrical	ON	Construction
1003	1008	Ottawa	COMM. OFFICE.LOW	0	Fixed Fee	April 8, 2020	Fe 19, 2021	Controls	ON	Sub Job

Canadian Cities:

city	city_ascii	province_id	province_name	population	density
Toronto	Toronto	ON	Ontario	5429524	4334.4
Montréal	Montreal	QC	Quebec	3519595	3889
Vancouver	Vancouver	BC	British Columbia	2264823	5492.6
Calgary	Calgary	AB	Alberta	1239220	1501.1
Edmonton	Edmonton	AB	Alberta	1062643	1360.9
Ottawa	Ottawa	ON	Ontario	989567	334