



## OpenAIR@RGU

### The Open Access Institutional Repository at Robert Gordon University

<http://openair.rgu.ac.uk>

This is an author produced version of a paper published in

Journal of Financial Management of Property and Construction (ISSN  
1366-4387)

This version may not include final proof corrections and does not include  
published layout or pagination.

#### Citation Details

##### Citation for the version of the work held in 'OpenAIR@RGU':

AHIAGA-DAGBUI, D. D. and SMITH, S. D., 2014. Rethinking  
construction cost overruns: cognition, learning and estimation.  
Available from *OpenAIR@RGU*. [online]. Available from:  
<http://openair.rgu.ac.uk>

##### Citation for the publisher's version:

AHIAGA-DAGBUI, D. D. and SMITH, S. D., 2014. Rethinking  
construction cost overruns: cognition, learning and  
estimation. *Journal of Financial Management of Property and  
Construction*, 19 (1), pp. 38-54.

#### Copyright

Items in 'OpenAIR@RGU', Robert Gordon University Open Access Institutional Repository,  
are protected by copyright and intellectual property law. If you believe that any material  
held in 'OpenAIR@RGU' infringes copyright, please contact [openair-help@rgu.ac.uk](mailto:openair-help@rgu.ac.uk) with  
details. The item will be removed from the repository while the claim is investigated.

2014

## Rethinking Construction Cost Overruns: Cognition, Learning and Estimation

Dominic D Ahiaga-Dagbui<sup>1</sup>, Simon D Smith

*School of Engineering, The University of Edinburgh, William Rankine Building, Edinburgh, EH9 3JL, Scotland, United Kingdom*

### ***AUTHOR'S COPY – CITATION***

Ahiaga-Dagbui D.D, Smith S.D(2014) [Rethinking construction cost overruns: cognition, learning and estimation](#), Journal of Financial Management of Property and Construction 2014 **19**:1 , 38-54

---

<sup>1</sup> Corresponding author: [d.d.ahiaga-dagbui@rgu.ac.uk](mailto:d.d.ahiaga-dagbui@rgu.ac.uk)

## Abstract:

- **Purpose:** Drawing on mainstream arguments in the literature, the paper presents a coherent and holistic view on the causes of cost growth, and the dynamics between cognitive dispositions, learning and estimation. A cost prediction model has also been developed using data mining for estimating final cost of projects.
- **Design:** A mixed-method approach was adopted: a qualitative exploration of the causes of cost overrun followed by an empirical development of a final cost model using Artificial Neural Networks (ANN).
- **Findings:** A conceptual model to distinguish between the often conflated causes of underestimation and cost overruns on large publicly funded projects. The empirical model developed in this paper achieved an average absolute percentage error of 3.67% with 87% of the model predictions within a range of  $\pm 5\%$  of the actual final cost.
- **Practical implications:** The model developed can be converted to a desktop package for quick cost predictions and the generation of various alternative solutions for a construction project in a sort of what-if analysis for the purposes of comparison. The use of the model could also greatly reduce the time and resources spent on estimation.
- **Originality:** A thorough discussion on the dynamics between cognitive dispositions, learning and cost estimation has been presented. It also presents a conceptual model for understanding two often conflated issues of cost overrun and under-estimation.
- **Keywords:** Cost Overruns, Optimism Bias, Strategic Misrepresentation, Data Mining, Dunning-Kruger Effects, Prospect Theory, Referenced Class Forecasting.

## INTRODUCTION

Cost performance on a construction project remains one of the main measures of the success of a construction project (Atkinson 1999, Chan and Chan 2004). Reliable cost estimates are important for several reasons – for organisational budgeting purposes, for loan application if a project has to be funded through credit facilities, to estimate likely cost of financing loans (interest payments), for estimating commercial feasibility or viability of the project. The present economic conditions also impose a parsimonious approach to spending on most organisations and governments. However, estimating the final cost of construction projects can be extremely difficult due to the complex web of cost influencing factors that need to be considered. These include type of project, material costs, likely design and scope changes, ground conditions, duration, size of project, type of client, tendering method and so on

(Ahiaga-Dagbui and Smith 2012). Trying to work out the influence of most of these variables at the inception stage of a project when cost targets are set, can be an exhaustive task, if not futile; while ignoring them altogether creates a recipe for cost overruns, disputes, law suits and even project termination in some cases. There is also a high level of uncertainty around most of these factors at the initial stages of the project as noted by Jennings (2012).

Table 1 shows major public projects that have experienced significant cost growth. Flyvbjerg *et al.* (2004) report that nine out of 10 infrastructure projects overrun their budgets and that infrastructure projects have an 86% likelihood of exceeding their budgets. The on-going Edinburgh Trams project, has already far exceeded its initial budget leading to significant scope reduction to curtail the ever-growing cost (Miller 2011, Railnews 2012). The recent 2012 London Olympics bid was awarded at circa £2.4 billion in 2005; was adjusted to about £9.3 billion in 2007 after significant scope changes; and was completed at £8.9 billion in 2010 (Gidson 2012, National Audit Office 2012). These statistics have often led to extensive claims, disputes and lawsuits in some cases within the industry (Love *et al.* 2010).

[Table 1 here]

Cost overrun in the construction industry has been attributed to a number of sources including technical error in design or estimation, managerial incompetency, risk and uncertainty, suspicions of foul play, deception and delusion, and even corruption. A recent debate on the *Construction Network of Building Researchers* (CNBR) on whether or not construction cost overruns could be attributed to error in estimation, or lies by project sponsors and estimators, raised more questions than answers (See the November 2012 CNBR archive online). For instance: How accurate or reliable can cost estimates be? What is the best measure of cost overrun? Might there be need to change how cost performance is presently measured? Should the estimator be absolved of the responsibility of producing reasonably accurate estimates? Should the industry even bother about cost overruns at all, if project goals are met in the long run?

While drawing on the works of some contemporary authorities on the subject, different schools of thought on causes of construction cost overruns have been synthesized in this study, to provide a coherent and holistic view of the problem. Recurring themes have been expanded upon, challenging traditional paradigms of assessing cost performance on construction projects while offering emerging frameworks of reckoning cost growth. It is proposed that there is a conflation of two quite different issues in the understanding of cost

growth: cost under-estimation and cost overrun. The paper presents a conceptual model for understanding these issues and then presents the development of a validated cost model using data mining and artificial neural networks. It is hoped that the continuous and effective application of data mining techniques might be one of the possible avenues for alleviating the problem of project cost overruns within the construction.

## **SOURCES OF COST GROWTH**

Causes of cost growth have been attributed to several sources including improperly managed risk and uncertainty (Okmen and Öztas 2010), scope creep (Love *et al.* 2011, Gil and Lundrigan 2012), optimism bias (Lovallo and Kahneman 2003, Jennings 2012) and suspicions of foul-play and corruption (Wachs 1990, Flyvbjerg 2009). While not attempting to provide a definitive list of all possible sources, the following section of the paper provides a synthesis of mainstream arguments on the causes of cost growth to provide a holistic view of the subject.

### ***Risk and Uncertainty***

The nature of a construction project makes it particularly prone to the effects of risk and uncertainty – it is complex and dynamic; each project has many parties with differing business and project objectives; projects are exposed to the weather (not in a controlled environment); and total project duration can spread over several years. It is no surprise then that risk, simply defined here as the measure of exposure to financial loss, or gain (Akintoye 2000), has been heavily cited as one of the main causes of failure to meet cost targets on construction projects (Skitmore and Ng 2003, Öztas 2004, Okmen and Öztas 2010). Arguably, the construction industry is perhaps one of the most risk prone industries, with project cost being one of main areas susceptible to its effects. Almost all types of risk (including scope changes, inclement weather, unsuitable ground conditions, disputes, client's cash flow problems, etc.) present financial ramifications.

Ahiaga-Dagbui and Smith (2012) noted that effective cost planning relates the design of facilities to their cost, so that while taking full account of quality, risks, likely scope changes, utility and appearance, the cost of a project is planned to be within the economic limit of expenditure. This stage in a project life-cycle is particularly crucial as decisions made during

the early stages of the development process carry more far-reaching economic consequences than the relatively limited decisions which can be made later in the process. Despite the importance of cost estimation, it is undeniably not simple, nor straightforward, because of the lack of information in the early stages of the project (Hegazy 2002). To achieve accuracy, the estimator has to be able to predict the future – something even the best technologies cannot achieve with certainty. This is because accurate reasoning is only possible in a world where information is complete and certain, and where cause and effect links are accurately known. Risk and uncertainty thus deeply pervade the construction industry and continue to cause unending controversy and debate. As Baccharini (2005) suggests, all too often risks are either ignored or dealt with in a completely arbitrary manner using rules-of-thumb or percentages. Flanagan and Norman (1993) also point out that the task of risk management or response in most cases is thus so poorly performed, that far too much risk is passively retained, ultimately resulting in cost escalation during project delivery.

### ***Strategic Misrepresentation and Optimism Bias***

Some authorities on the subject of cost overrun have proposed more depressing explanations to the phenomenon. Flyvbjerg *et al.*, suggest that overruns are chiefly due to ‘strategic misrepresentations’, i.e. outright lying (Flyvbjerg *et al.* 2002) and ‘optimism bias’ (Flyvbjerg 2007). Flyvbjerg *et al.* compared the cost of projects at the time of the decision to build to the cost at completion and found inaccuracies in cost forecasts for transportation infrastructure projects to be on average 44.7% for rail, 33.8% for bridges and tunnels, 20.4% for roads – concluding that nine out of 10 projects outrun their cost targets. Overruns beyond 100% of original cost are also not uncommon (Trost and Oberlander 2003, Odeck 2004).

In order to get a project approved, sponsors and estimators, especially on public works, tend to intentionally underestimate the true cost of the project in what has been described as the ‘Machiavelli factor’ (Flyvbjerg 2003). “*By routinely overestimating benefits and underestimating costs, promoters make their projects look good on paper, which helps get them approved and built*” (Flyvbjerg *et al.* 2005). It makes little reasoning to stop the project once a considerable amount of money has already been spent to get it started, Flyvbjerg (2004) claims. Wach (1989) was even more forthright in his paper ‘*When planners lie with numbers*’ and later advocated for better ethics in forecasting for public works (Wachs 1990).

If cost overruns cannot be explained by intentional underestimation, optimism bias might be a likely culprit (Flyvbjerg 2007). Optimism bias can be explained as the cognitive disposition to evaluate future events in a fairer light than they might actually be in reality (Lovallo and Kahneman 2003). Unlike strategic misrepresentation, this might not be born out of deceptive intent, but also often leads to underestimating true cost, overestimation of benefits, and overlooking the potential of error and uncertainty. The potential gains of the project thus become overwhelmingly enticing, and almost blinding to likely pitfalls. It also leads to underestimating the full extent of certain risk events, should they occur.

In effect, delusion and deception are complementary explanations of the failure of large infrastructure projects, causing works such as diverting existing utilities, environmental impacts and foreseeable risks to be continually underestimated in construction (Flyvbjerg 2009). This line of diagnosis of the problem of cost overrun might seem appealing, at least on first thought, especially in terms of large capital intensive public projects or those that are likely to make high political statements. Flyvbjerg's far-reaching work on cost overruns led to the endorsement of his 'Reference Class Forecasting' by the American Planning Association in 2005 (cf. APA 2005, Flyvbjerg 2007). This will be discussed in more detail in this paper.

### ***Going beyond Strategic Misrepresentation and Optimism Bias***

Even though deception and delusion might be plausible explanations for cost overruns, particularly in large publicly funded or politically motivated projects, they are not easily generalisable to all types of projects undertaken within the construction industry. Researchers, including Love (2012), rebut Flyvbjerg's conclusions as simplistic, largely misleading and not an accurate reflection of reality. Love *et al.*'s rejoinder suggests a move beyond optimism bias and strategic misrepresentation to focus on intermediary events, actions, the so-called 'pathogens' that occur between project inception and completion. At the core of Love's argument is that many events and actions that are not accounted for in initial estimates, tend to drive up cost. This school of thought is largely supported by Aibinu and Pasco (2008), Odeck (2004) and Odeyinka *et al* (2012). Love's case study of social infrastructure projects suggest that foul-play, as suggested by Flyvbjerg and Wach, might not be best explanations of cost overruns; and that the fingers point at events that occur before and during the project delivery stage (Love *et al.* 2011). Besides, it is almost impossible to

draw valid distinctions along a continuum of motivation when promoting a project from reasonable optimism, through over-enthusiasm, culpable error, to deliberate deceit using statistical analysis, as adopted in the Flyvbjerg's works.

Research on leadership and governance of construction projects by Gil and Lundrigan (2012), perhaps offers a more holistic assessment of cost growth that aligns closely with the views of Love, *et al* above. That projects evolve, is essentially, the core of their defence. Very often, construction projects change considerably in scope and design between conception, to inception and completion, often due to a client's proposed changes or technically imposed changes. This suggests that it might be erroneous to simply compare the cost of a project at inception, A, with the cost at completion, B, and wherever  $B > A$ , then overruns have occurred and estimators of A either lied or were incompetent. A and B are essentially very different. More robust explanations of overruns need to factor-in process and product, as well as sources of changes to scope. For Love and Gil *et al* (*op. cit.*), project overruns are not really a case of projects not going according to plan (budget), but the other way round – plans not going according to project.

Gil and Lundrigan (2012) propose a 'relay race' framework for understanding cost growth, particularly on mega projects such as the London Olympics Project, Scottish Parliament or Terminal 2 project at Heathrow Airport, all of which seemed to have suffered the curse of cost growth, at least on a perfunctory examination. In the relay race of construction delivery, the baton of project leadership is passed on from one person(s) or organisation at the different stages of the project delivery. The aims and scope of the project, as well as skills and competencies of the project sponsors and promoters (project governors) at the conceptual stage, are often very different from their counterparts at the project design or delivery stage. Also, it is not unusual for most public projects to have long gestation periods, stretching over several years before final approval is reached, by which time project budget would also have changed a number of times. The Scottish Parliament Building is a paragon in this respect – the *circa* £40 million submitted by the Scottish Office as likely final cost did not take into consideration project location, or the building of a completely new parliament building. It is no wonder the final cost of the project was 10 times this initial proposed cost (Fraser 2004).

### ***Perception and Measuring Overruns***



Perhaps our perception of cost overruns needs to change altogether. What is described as cost overruns at the moment might not be overruns after all if reckoned through the eyes of different procurement routes, for example. It is possibly one of the reasons why cost overrun is not often reported in projects procured through joint ventures or alliancing. Typically, in traditional contracting, design and estimates are first prepared by the Client's Estimator (CE) and then bids are invited from contractors. The lowest bidder often wins the job with the lowest tender value becoming the cost estimate at the beginning of the project (A). The contractor undertakes then to deliver the project at cost, A, and all add-ons are dealt with through change orders or claims until project completion at cost, B. Whenever  $B > A$ , overruns are reported. It is easy to identify how competition, market conditions, optimism bias and the selection by lowest bidder combine to drive down the initial estimate, A, creating a somewhat unrealistic target as likely final cost. For the contractor therefore, winning work at the right price (realistic cost) becomes a very difficult task. To be thorough in estimation would mean including likely cost of most/all risk events in the tender, consequently pricing the contractor out of competition. Most contractors may therefore not include potential risk events in their tenders, so as to increase their likelihood of winning the contract. This was evident in related studies in modelling final cost of construction projects (Ahiaga-Dagbui and Smith 2012).

Some have suggested that the industry move beyond its fixation on measuring project success largely in terms of cost (Bassioni *et al.* 2004, Yeung *et al.* 2008). The CNBR debate was frequently punctuated by the question, 'Why care about cost overruns anyway? If projects run over budget but deliver what the client wants, shouldn't everyone be happy?' After all, cost overruns only represent our human inability to predict future events accurately, or identify risks and quantify their likely impact and cost. Others think perhaps there is a need for a paradigm shift in how projects are evaluated to cover a combination of social, economic, social, usability or value for money (Toor and Ogunlana 2010). The Sydney Opera House experienced large overruns at the time of construction but it is now generally considered a 21<sup>st</sup> century icon of buildings and a popular destination for tourists and opera concerts. Similarly, in spite of the controversies about cost overruns, the Scottish Parliament Building has won several awards, including the coveted Stirling Award in 2005 by the Royal Institute of British Architects for its audacious, highly conceptual and iconic design. Even if cost should be a major factor for assessment, it certainly should not be a simplistic or statistical comparison between awarded contract sum and cost at final accounts.

## *Cognition, Bias and Learning*

Can a science that combines intuition and analysis ever be precise or unbiased? A qualified 'no' is probably the answer to that question, according to Kahneman and Tversky (1979), formulators of Prospect Theory – decision making under risk and uncertainty. The theory suggests people make decisions based on the likely gains, or loss, of a venture, and not necessarily based on the real outcome of the decision. It further proposes that decision making is often flawed by systematic biases and that error in judgement is often systematic and predictable, rather than random. Kahneman, a Noble Prize winner for his works on decision making and behavioural economics, delineates decision making and the illusion of understanding, stating that we often exhibit an excessive confidence in what we believe we know about any situation, and that our inability to acknowledge the full extent of our ignorance and the uncertainty of the world we live in makes us prone to overestimate how much we really understand (Kahneman 2011). Kahneman's work with Lovallo (2003) provides further defence of the Prospect Theory from different business areas. Kahneman's theory holds profound extensions for decision making in the construction industry, especially for large public projects where the effects and cost of risk and uncertainty are particularly heightened. It would also provide large support of Flyvbjerg's arguments on strategic misrepresentation and optimism bias already discussed in this paper. Conceivably, this is one reason why it is easy to err on the side of optimism when promoting a project, or when estimating the outcome of a risk event.

Perhaps even more controversial are the conclusions reached by Kruger and Dunning (2009), that incompetence does not only cause poor performance but also has the dual effect of robbing people of the ability to recognise poor performance. They posit that the metacognitive skills required to judge the accuracy of a decision is the same required to evaluate the error in the same decision – to lack the former, is to fall short in the latter as well (Kruger and Dunning 1999). The result thereof is that the *"incompetent will tend to grossly overestimate their skills and abilities"* (Kruger and Dunning 2009). They tied their conclusion to Darwin's pronouncement: *"ignorance more frequently begets confidence than does knowledge"* (Darwin 1871), a theory largely supported by Ehrlinger *et al.* (2008) and Maki *et al.* (1994).

Herein lies the estimation complex – a combination of optimism bias and prospect theory predisposes us to underestimate true cost, discounting the real effect of uncertainty and error while doing so. At the same time, Dunning-Kruger tendencies blind forecasters to the error in reaching unrealistic estimates for project cost. Juxtapose these with the effect of risk and uncertainty, competition embedded within the culture of lowest-bidder tendering, as well as strategic misrepresentation, and the overruns reported in Table 1 become less surprising. It is easier to understand how most cost estimates can be prepared, or at least reported, with an unjustifiable confidence in their accuracy. If this is the case, then perhaps we might not have to move beyond optimism bias just yet, as suggested by Love (2011). If we are indeed systematically prone to err towards optimism bias in our reasoning, then it might be wise to rethink how that affects our estimates and what needs to be done about it.

Flyvbjerg (2002) also noted that ‘no learning’ seemed to be taking place in the construction industry over the 70 years prior to his study, and that estimation accuracy has not seen much improvement even with the advancements in technology and the proliferation of cost models and project management approaches. Kruger and Dunning (2009), as well as Ehrlinger *et al* (2008) attribute lack of performance improvement to the lack of accurate and constructive feedback. They however observed, that an awareness of limitations of skills and decision making within an environment of uncertainty, helped to improve performance and self-calibration. A lack of learning in the construction industry could be explained in a number of ways: that the mitigating factors causing overruns are ones that the industry absolutely cannot overcome and therefore, has to accept cost overruns as normal part of practice; or, that there is simply very little incentive to reach realistic target inception; or further still that the industry seems largely to miss the opportunities offered by effective knowledge transfer and feedback from previously completed projects (see Hartmann and Dorée, 2013). How is explicit and tacit knowledge captured and utilised within the industry presently? How do project closure reports feed back into the development of new projects for continuous improvement?

## **RETHINKING OVERRUNS**

For the purposes of cost modelling or estimation, it is important to clarify an important point. Existing literature, and recent CNBR debate, on ‘cost overruns’ seems to conflate two related, but different issues – overruns and underestimation. Unfortunately, a lot of cost models do

not make this distinction either and thus become limited in their application in practice. As already pointed out, most large publicly funded projects tend to go through a long gestation period after project conception during which many changes to scope and accompanying costs occur – sometimes the initial scheme bears little likeness to the defined project. The estimated cost at project inception often fails to take into consideration a lot of details and information, largely because much of these are not yet available or uncertain; the case of the initial circa £40 million estimate for the Scottish Parliament. For many large publicly funded projects, this is normally when project sponsors garner for project approval and funding. It is perhaps at this stage the effects of Prospect Theory, Dunning-Kruger effect, optimism bias and strategic misrepresentation are particularly heightened, to keep cost at an attractive low and benefits of undertaking the project high. This might be what accounts for what the authors refer to as underestimation of likely cost – the difference between estimated cost at project inception and cost at the end of project definition phase in Figure 1.

[Figure 1 here]

Overruns however, are aptly described as the difference in cost at project completion and project definition stage (refer Figure 1). This is usually as a result of further scope changes, normally not as significant as those at project definition stage, ground conditions, technical and managerial difficulties, material or labour price changes or estimation error. These are the factors that Love *et al* (2011) describe as ‘pathogens’. So, whereas, Flyvbjerg’s work mainly deals with underestimation, Love’s explanations for cost growth largely covers the latter phases of the construction project. It is important to note however that Figure 1 is not necessarily wholly applicable for small, non-political and routine projects where the effects of the political and cognitive causes of cost growth are less heightened. Much of the media hype on cost overruns however is often based on a comparison between cost at inception and cost at completion, almost ignoring the mediating phases of project gestation and definition.

## REFERENCE CLASS FORECASTING

Flyvbjerg developed a practical method for forecasting cost of large projects based on Reference Class Forecasting (RCF) formulated by Kahneman and Tversky (1979, 1994). RCF attempts to use ‘distributional information’ (knowledge) from previous projects similar to the new project being undertaken, the so-called taking of an ‘outside view’ of planned actions,

based on actual past performance. Kahneman and Flyvbjerg reckon this approach might somehow help to bypass optimism bias and strategic misrepresentation in decision making (Flyvbjerg 2007). The methodology involves three steps, summarised simply here as:

- a. Identify a reference class of past, similar projects.
- b. Estimate a probability distribution for the selected reference class, and
- c. Establish likely cost of the new project using the reference class distribution.

The first instance of its application was on Edinburgh Tram project by the UK Government – the original forecast by the Transport Initiatives Edinburgh (tie), the project promoter was about £255 million but the RCF indicated this could rise up to £400 million and warned that the final cost could even be exceedingly higher (Flyvbjerg 2007). Recent estimates now indicate that the final construction cost of the Trams could be around £776 million (Miller 2011, Railnews 2012). The RCF has reportedly been applied to the £15 billion London Crossrail and £7.5 million Taunton Third Way projects in the UK (Flyvbjerg 2007).

Even though RCF remains to be widely tested or adopted, it might be a step in the right direction especially in dealing with the root causes of *underestimation*, (as opposed to cost overrun) as shown in Figure 1, i.e. optimism bias, Prospect Theory, Dunning-Kruger effect and strategic misrepresentation. However, as pointed out by Flyvbjerg, RCF is largely applicable to large, non-routine or one-off projects such as stadiums, museums, dams, etc. On smaller, less political, or frequent projects however, a fairly similar but more established method of forecasting that employs previous experience and incremental learning is data mining. This has been extensively used in other industries including finance (Kovalerchuk and Vityaev 2000), medicine (Bellazzi and Zupan 2008, Koh and Tan 2011) and business (Apte *et al.* 2002), but is yet to see widespread adoption in the construction industry. Notwithstanding, it has been applied to construction knowledge management (Yu and Lin 2006), for estimating the productivity of construction equipment (Yang *et al.* 2003), study of occupational injuries (Cheng, Leu, *et al.* 2012), alternative dispute resolution (Fan and Li in press) and prediction of the compressive strength high performance concrete (Cheng, Chou, *et al.* 2012). Data mining is used to develop final cost models in the next section of this paper, in a manner that addresses the *overruns* part of Figure 1.

## FINAL COST MODEL DEVELOPMENT USING DATA MINING

Data mining is the analytic process for exploring large amounts of data in search of consistent patterns, correlations and/or systematic relationships between variables; and to then validate the findings by applying the detected patterns to new subsets of data (StatSoft Inc 2008). Data mining attempts to scour databases to discover hidden patterns and relationships in order to find predictive information for business improvement. Similar to reference class forecasting, data mining starts with the selection of relevant data from a data warehouse that contains information on organisation and business transactions of the firm (Ngai *et al.* 2009). The selected data set is then pre-processed before actual data mining commences. Data pre-processing typically involves steps such as sub-sampling, clustering, transformation, de-noising, normalisation or feature extraction (StatSoft Inc. 2011), to ensure that the data are structured and presented to the model in the most suitable way for effective modelling.

[Figure 2 here]

The next stage, as shown in Figure 2, involves the actual modelling, where one or a combination of data mining techniques is applied to scour down the dataset to extract useful knowledge. The results obtained are then evaluated and presented into some meaningful form to aid business decision making. This final step might involve graphical representation or visualisation of the model for easy communication. Artificial neural networks (ANN) is used for the modelling aspect of this study mainly because of its learning and generalisation capabilities (Anderson 1995).

### **Data**

The data used for the models in this paper were supplied by an industry partner with its primary operation in the delivery of water infrastructure and utility in the UK. Approximately 1600 projects completed between 2004 and 2012, with cost range of between £4000 to £15 million, comprising newly built, upgrade, repair or refurbishment projects. Fifteen project cases were selected using stratified random sampling to be used for independent testing of the final models. The remaining data were then split in an 80:20% ratio for training and testing of the models, respectively.

Cost values were normalised to a 2012 baseline with base year 2000 using the infrastructure resources cost indices by the Building Cost Information Services (Building Cost Information Services 2012). Numerical predictors were further standardized to *zScores* using

$$zScore = \frac{x_i - \mu}{\sigma} \quad (\text{Equation 1})$$

where: *zScore* is the standardized value of a numerical input,  $x_i$ ;  $\mu$  is the mean of the numerical predictor; and  $\sigma$  is the standard deviation of the numerical predictor.

This allowed numerical inputs to be squashed into a smaller range of variability, potentially improving the numerical condition of the optimization process of the model (StatSoft Inc 2008). If one input has a range of 0 to 1, while another has a range of 0 to 30 million, as was the case in the data that were used in this analysis, the neural net will expend most of its effort learning the second input to the possible exclusion of the first. All categorical variables were coded using a binary (0,1) coding system. Data screening using scree test, factor analysis and optimal binning allowed for the selection of six initial predictors (primary purpose of project, project scope, project delivery partners, operating region, project duration, and initial estimated cost) to be used for the actual modelling using ANN. Invariant variables, such as payment method, fluctuation measure and type of client, were removed from the variable set as they would only increase the model complexity and yet offer no useful information for model performance.

### ***Model Development***

The final model was developed after an iterative process of fine-tuning the network parameters and/or inputs until acceptable error levels were achieved or when the model showed no further improvement. First, the automatic network search function of Statistica 10<sup>®</sup> software was used to optimise the search for the best network parameters, after which customized networks were developed using the optimal parameters identified. Five activation functions<sup>i</sup> were iterated in both hidden and output layers, using gradient descent, conjugate descent and Quasi-Newton (BFGS) training algorithms. About two thousand multi-layer perceptron networks were trained at each iteration stage, retaining the 5 best before further tweaking to investigate possible model improvement.

Early stopping, the process of halting training when the test error stops decreasing, was used to prevent memorising or over-fitting the dataset in order to improve generalization. Over-

fitted models perform very well on training and testing data, but fail to generalise satisfactorily when new ‘unseen’ cases are used to validate their performance. The best networks at each stage were selected based on their overall performance, measured using the correlation coefficient between predicted and output values as well as the Sum of Squares (SOS) of errors. SOS is defined here as:

$$SOS = \sum(T_i - O_i)^2 \quad (\text{Equation 2})$$

Where:  $O_i$  is the predicted final cost of the  $i$ th data case (Output); and  $T_i$  is the actual final cost of the  $i$ th data case (Target).

The higher the SOS value, the poorer the network at generalisation, whereas the higher the correlation coefficient, the better the network. The  $p$ -values of the correlation coefficients were also computed to measure their statistical significance. The higher the  $p$ -value, the less reliable the observed correlations. Overall, about 30 networks were retained, which were then validated using the 15 separate projects that were selected using stratified sampling at the beginning of the modelling exercise. Figure 3 shows the performance of the best 7 out the 30 validated models.

[Figure 3 here] and [Table 2 here]

Table 2 shows the performance of overall best model (model 33). It compares the final cost forecasts reached by the model with the actual final cost recorded at the end of the project. This model was an MLP 8-11-1, i.e. a multilayer perceptron with eight nodes in the input layer, 11 hidden units and one output (final cost). It was trained with a Quasi-Newton (BFGS) training algorithms and had a hyperbolic tangent (tanH) activation function in both hidden and output layers. The tanh activation function, defined in equation 3, squashes continuous variables into a range of (-1,+1) for more effective training of the neural network models:

$$f(x) = \frac{\sinh x}{\cosh x} = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{e^{2x} - 1}{e^{2x} + 1} \quad (\text{Equation 3})$$

The final predictors in this model were the purpose of the project, the construction delivery partner used by the client, the estimated duration, an early scheme estimate of final cost and



scope of the project. The average APE achieved by this model was 3.67% across the 15 validation cases. Its APEs ranged between 0.04% and 15.85%. It was observed that the worst performances of the model were achieved on projects with the smallest values in the validation set (cases 13 and 15). This might be because a majority of the projects used for the model training had values in excess of £5 million. However, the actual monetary errors on these predictions were deemed satisfactory as they were relatively small (about £3500 and £2500 for models 13 and 15 respectively). Eighty-seven per cent of the validation predictions of the best model were within  $\pm 5\%$  of the actual cost of the project. The authors are now exploring avenues of transforming the models into standalone desktop applications for deployment within the operations of the industry partner that collaborated in this research.

## CONCLUSION

Cost estimate reliability and accuracy on construction projects continues to receive a lot of attention from both industry and academia. The industry is faced with a complex web of causes, which we propose fall into two distinct yet often conflated categories – cost underestimation and cost overrun summarised as follows.

### *Underestimation*

- Optimism Bias- a propensity to believe and act on a notion that all will go well leading to the underestimation the role of uncertainty in outcomes;
- Prospect Theory- making decisions based on likely gains and loss rather than the actual outcome of the decision;
- Strategic misrepresentation- outright lying and corruption;
- Dunning-Kruger effect- the bend to overestimate competency or accuracy in judgement and the inability to see past our own errors; competition to win projects.

### *Overrun*

- Scope changes, whether mandated by circumstances or requested by client;
- Managerial and technical difficulties;
- Risk and uncertainty; and
- Ground conditions, price changes (etc.).

Most of these, especially the cognitive and psychological factors, tend to work together to drive down the true cost of the project during the initial stages, creating a false and unreliable estimate as target to reach. We have attempted to provide a holistic view of the problem of cost growth, while presenting a conceptual model to distinguish between these often conflated ideas of underestimation and overruns on construction projects. Reference Class Forecasting was discussed as a possible means of addressing underestimation, particularly on large publicly funded projects. The development of a final cost prediction model using data mining and artificial neural networks was then presented as a possible avenue of addressing cost overruns in the construction industry. The best model achieved an average absolute percentage error of 3.67% with 87% of the validation predictions falling within an error range of  $\pm 5\%$ . These methods can be used to develop decision support systems especially at early stages of the construction project as well as complement traditional methods of estimation in order to reach more accurate and reliable cost estimates.

Clients can play a crucial role in ensuring the quality and reliability of cost estimates in the construction industry. As indicated by the Commercial Manager of one of the biggest construction companies in the UK, "*winning a tender is easy. But winning at the right price is difficult*". Unless clients start demanding *realistic* estimates, rather than the lowest estimates at the early stages of a project, the problem of cost overrun might remain with the industry for a long time to come. Cultural changes within the industry towards the search for realistic targets might incentivise contractors to flag up potential estimating pitfalls early-on. Questions about who has the responsibility on behalf of the client to govern the project always has profound implications on cost growth from inception to completion and needs to be addressed very early on a project. This is particularly important on mega projects.

Project knowledge capture and its utilisation would also be crucial in tackling cost overruns. Some data mining techniques like neural networks are particularly useful in modelling both explicit and tacit knowledge within extensive databases. This can be used to complement traditional cost estimation methods or RFC to reach more realistic and reliable estimates. Finally, and perhaps more importantly, is the creation of a culture of critical questioning, measures of accountability, with checks and balances to make sure that cost is managed to be within reasonable budget limits.

## REFERENCES

- Ahiaga-Dagbui, D.D. and Smith, S.D. (2012) Neural networks for modelling the final target cost of water projects. *in* Smith, S. D. (ed.) *Procs 28th Annual ARCOM Conference*. Edinburgh, UK 3-5th September 2012 Association of Researchers in Construction Management, pp.307-316.
- Ahiaga-Dagbui, D D and Smith, S D (2013) "My cost runneth over": Data mining to reduce construction cost overruns. In: *Procs 29th Annual ARCOM Conference*, Smith, S D and Ahiaga-Dagbui, D D, Eds.), Reading, UK: Association of Researchers in Construction Management, 559-68.
- Aibinu, A.A. and Pasco, T. (2008) The accuracy of pre-tender building cost estimates in Australia. *Construction Management and Economics*, **26**(12), pp.1257 - 1269.
- Akintoye, A. (2000) Analysis of factors influencing project cost estimating practice. *Construction Management & Economics*, **18**(1), pp.77-89.
- Anderson, J.A. (1995) *An Introduction to Neural Networks*. Cambridge, Massachusetts: MIT Press.
- APA (2005) *JAPA article calls on planners to help end inaccuracies in public project revenue forecasting*. [online]. American Planning Association (APA). Available at: <<http://goo.gl/I7DLA>>.
- Apte, C., Liu, B., Pednault, E.P.D. and Smyth, P. (2002) Business applications of data mining. *Communications of the ACM*, **45**(8), pp.49-53.
- Atkinson, R. (1999) Project management: cost, time and quality, two best guesses and a phenomenon, its time to accept other success criteria. *International Journal of Project Management*, **17**(6), pp.337-342.
- Baccarini, D. (2005) Estimating Project Cost Contingency – Beyond the 10% syndrome. *2005 Australian Institute of Project Management Conference* [online] AIPM. [Accessed 12/08/2011]. Available at: <<http://www.aipm.com.au/resource/Baccarini-AIPMconf05.pdf>>.
- Bassioni, H., Price, A. and Hassan, T. (2004) Performance measurement in construction. *Journal of management in engineering*, **20**(2), pp.42-50.
- BCIS (2012) *BIS Construction Price and Cost Indices*. <http://www.bcis.co.uk>: Building Cost Information Services, UK.
- Bellazzi, R. and Zupan, B. (2008) Predictive data mining in clinical medicine: current issues and guidelines. *International Journal of Medical Informatics*, **77**(2), pp.81-97.

- Chan, A.P. and Chan, A.P. (2004) Key performance indicators for measuring construction success. *Benchmarking: An International Journal*, **11**(2), pp.203-221.
- Cheng, C.-W., Leu, S.-S., Cheng, Y.-M., Wu, T.-C. and Lin, C.-C. (2012a) Applying data mining techniques to explore factors contributing to occupational injuries in Taiwan's construction industry. *Accident Analysis & Prevention*, **48**(0), pp.214-222.
- Cheng, M.-Y., Chou, J.-S., Roy, A.F.V. and Wu, Y.-W. (2012b) High-performance Concrete Compressive Strength Prediction using Time-Weighted Evolutionary Fuzzy Support Vector Machines Inference Model. *Automation in Construction*, **28**(0), pp.106-115.
- Darwin, C. (1871) *The decent of man*. London: John Murray.
- Ehrlinger, J., Johnson, K., Banner, M., Dunning, D. and Kruger, J. (2008) Why the unskilled are unaware: Further explorations of (absent) self-insight among the incompetent. *Organizational Behavior and Human Decision Processes*, **105**(1), pp.98-121.
- Fan, H. and Li, H. (in press) Retrieving similar cases for alternative dispute resolution in construction accidents using text mining techniques. *Automation in Construction*, (0).
- Flanagan, R and Norman, G (1993) *Risk Management and Construction*. Oxford: Blackwell Science Ltd.
- Flyvbjerg, B. (2003) Machiavellian Tunnelling. *World Tunnelling*, p.43.
- Flyvbjerg, B. (2005) Design by Deception: The Politics of Megaproject Approval. *Harvard Design Magazine*, **22**, pp.50-59.
- Flyvbjerg, B. (2007) Curbing Optimism Bias and Strategic Misrepresentation in Planning: Reference Class Forecasting in Practice. *European Planning Studies*, **16**(1), pp.3-21.
- Flyvbjerg, B. (2009) Survival of the unfittest: why the worst infrastructure gets built—and what we can do about it. *Oxford Review of Economic Policy*, **25**(3), pp.344-367.
- Flyvbjerg, B., Holm, M.K.S. and Buhl, S.L. (2002) Understanding costs in public works projects: Error or lie? *Journal of the American Planning Association*, **68**(279-295).
- Flyvbjerg, B., Holm, M.K.S. and Buhl, S.L. (2004) What Causes Cost Overrun in Transport Infrastructure Projects? *Transport Reviews*, **24**(1), pp.3-18.
- Flyvbjerg, B., Skamris Holm, M.K. and Buhl, S.L. (2005) How (In)accurate Are Demand Forecasts in Public Works Projects?: The Case of Transportation. *Journal of the American Planning Association*, **71**(2), pp.131-146.
- Fraser (2004) *Holyrood Enquiry- A Report by the Rt Hon Lord Fraser of Carmyllie QC on the construction of the Holyrood Building Project presented to the First Minister and Presiding Officer*. (SP Paper No. 205). <http://www.holyrood inquiry.org/>: Scottish Parliamentary Corporate Body 2004.

- Gidson, O. (2012) London 2012 Olympics will cost a total of £8.921bn. *The Guardian* [online]. 23 October 2012. [Accessed 22 April, 2013]. Available at: <<http://goo.gl/sxatK>>.
- Gil, N. and Lundrigan, C. (2012) *The Leadership and Governance of Megaprojects* [online] Centre for Infrastructure Development (CID), Manchester Business School, The University of Manchester. Available at: <<http://goo.gl/cf2ST>>.
- Hartmann, A. and Dorée, A. (2013) Messages in bottles: The fallacy of transferring knowledge between construction projects. in Smith, S. D. and Ahiaga-Dagbui, D. D. (eds.) *Procs 29th Annual ARCOM Conference*. Reading, UK 2-4th September 2013 Association of Researchers in Construction Management, p.In Press.
- Hegazy, T. (2002) *Computer-based construction project management*. Upper Saddle River, NJ: Prentice Hall Inc.
- Jennings, W. (2012) Why costs overrun: risk, optimism and uncertainty in budgeting for the London 2012 Olympic Games. *Construction Management and Economics*, **30**(6), pp.455-462.
- Kahneman, D. (1994) New challenges to the rationality assumption. *Journal of Institutional and Theoretical Economics*, **150**, pp.18 – 36.
- Kahneman, D. (2011) *Thinking, Fast and Slow*. New York: Farrar, Straus and Giroux.
- Kahneman, D. and Tversky, A. (1979) Prospect Theory: An Analysis of Decision under Risk. *Econometrica*, **47**(2), pp.263-291.
- Koh, H.C. and Tan, G. (2011) Data mining applications in healthcare. *Journal of Healthcare Information Management—Vol*, **19**(2), p.65.
- Kovalerchuk, B. and Vityaev, E. (2000) *Data mining in finance*. USA: Kluwer Academic Publisher, Hingham MA.
- Kruger, J. and Dunning, D. (1999) Unskilled and unaware of it: How difficulties in recognizing one's own incompetence lead to inflated self-assessments. *Journal of personality and social psychology*, **77**(6), pp.1121-1134.
- Kruger, J. and Dunning, D. (2009) Unskilled and unaware of it: how difficulties in recognizing one's own incompetence lead to inflated self-assessments. *Psychology*, **1**, pp.30-46.
- Lovullo, D. and Kahneman, D. (2003) Delusions of Success: How Optimism Undermines Executives' Decisions. *Harvard Business Review*, (July 2003).

- Love, P.E.D., Davis, P., Ellis, J. and Cheung, S.O. (2010) Dispute causation: identification of pathogenic influences in construction. *Engineering, Construction and Architectural Management*, **17**(4), pp.404-423.
- Love, P.E.D., Edwards, D.J. and Irani, Z. (2011) Moving beyond optimism bias and strategic misrepresentation: An explanation for social infrastructure project cost overruns.
- Love, P.E.D., Sing, C.-P., Wang, X., Irani, Z. and Thwala, D.W. (2012) Overruns in transportation infrastructure projects. *Structure and Infrastructure Engineering*, pp.1-19.
- Maki, R.H., Jonas, D. and Kallod, M. (1994) The relationship between comprehension and metacomprehension ability. *Psychonomic Bulletin & Review*, **1**(1), pp.126-129.
- Miller, D. (2011) Edinburgh Trams: Half a line at double the cost. *BBC* [online]. [Accessed 18th February, 2013]. Available at: <<http://goo.gl/mfr96>>.
- NAO (2012) *The London 2012 Olympic Games and Paralympic Games: post-Games review*. HC 794- Session 2012-13 National Audit Office, UK.
- Ngai, E.W.T., Xiu, L. and Chau, D. (2009) Application of data mining techniques in customer relationship management: A literature review and classification. *Expert Systems with Applications*, **36**(2), pp.2592-2602.
- Odeck, J. (2004) Cost overruns in road construction—what are their sizes and determinants? *Transport Policy*, **11**(1), pp.43-53.
- Odeyinka, H., Larkin, K., Weatherup, R., Cunningham, G., McKane, M. and Bogle, G. (2012) *Modelling risk impacts on the variability between contract sum and final account (A research report submitted to RICS)*. London, United Kingdom: Royal Institution of Chartered Surveyors.
- Okmen, O. and Öztas, A. (2010) Construction cost analysis under uncertainty with correlated cost risk analysis model. *Construction Management and Economics*, **28**(2), pp.203-212.
- Öztas, A. (2004) Risk analysis in fixed-price design–build construction projects. *Building and Environment*, **39**(2), pp.229-237.
- Railnews (2012) Edinburgh tram costs soar again. *Railnews* [online]. [Accessed 18th February, 2013]. Available at: <<http://goo.gl/M5uZ7>>.
- Skitmore, M.R. and Ng, T.S. (2003) Forecast models for actual construction time and cost. *Building and Environment*, **38**(8), pp.1075-1083.

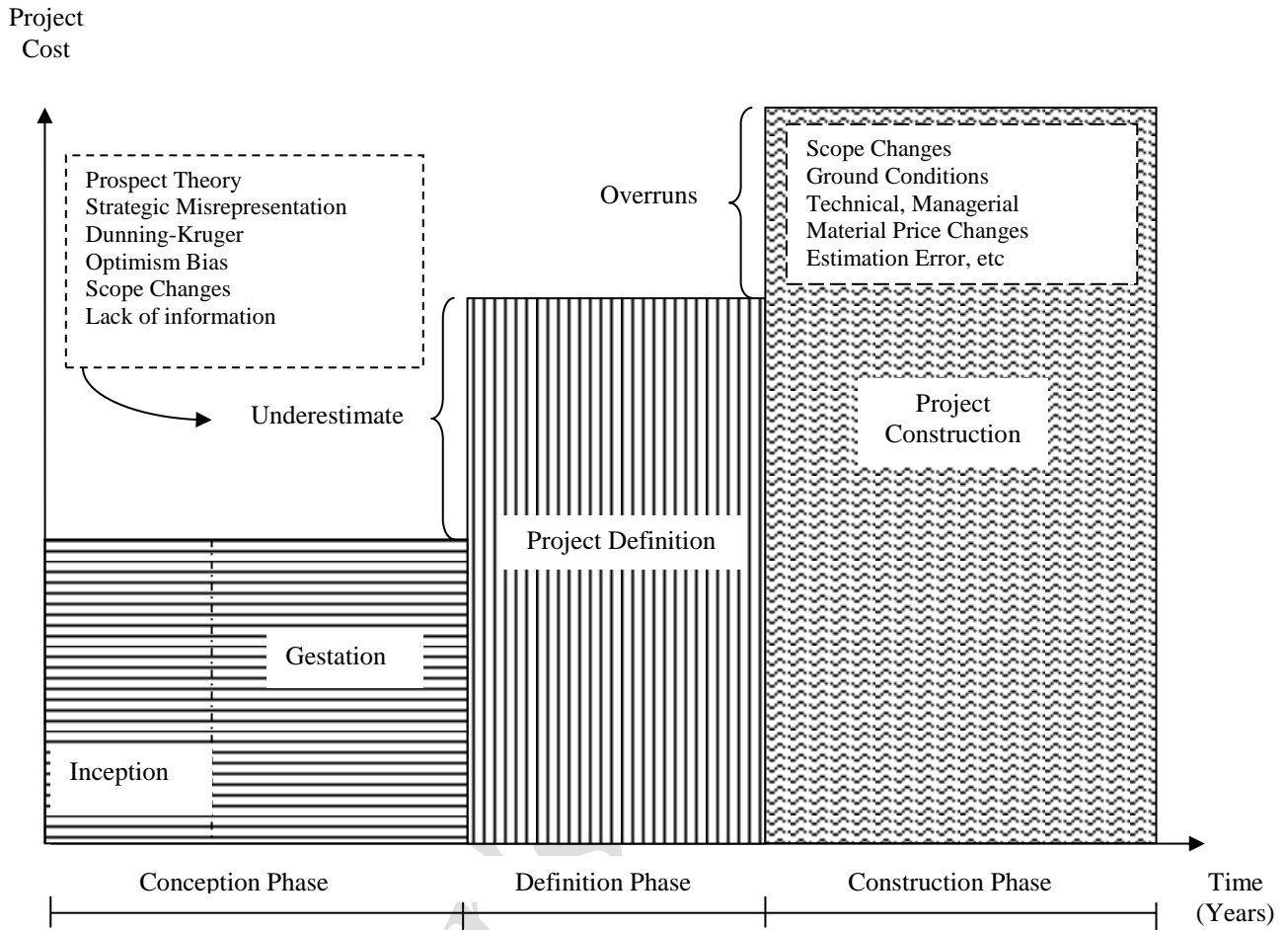
- StatSoft Inc (2008) *A Short Course in Data Mining*. StatSoft, Inc. [online]. [Accessed 25th January, 2012]. Available at: <[http://www.statsoft.com/Portals/0/Products/Data-Mining/data\\_mining\\_tutorial.pdf](http://www.statsoft.com/Portals/0/Products/Data-Mining/data_mining_tutorial.pdf)>.
- StatSoft Inc. (2011) *STATISTICA 10 (data analysis software system)*, [www.statsoft.com](http://www.statsoft.com).
- Toor, S.-u.-R. and Ogunlana, S.O. (2010) Beyond the 'iron triangle': Stakeholder perception of key performance indicators (KPIs) for large-scale public sector development projects. *International Journal of Project Management*, **28**(3), pp.228-236.
- Trost, S.M. and Oberlander, G. (2003) Predicting Accuracy of Early Cost Estimates Using Factor Analysis and Multivariate Regression. *Journal of Construction Engineering and Management*, **129**(2).
- Wachs, M. (1989) When planners lie with numbers. *Journal of the American Planning Association*, **55**(4), pp.476-479.
- Wachs, M. (1990) Ethics and advocacy in forecasting for public policy. *Business and Professional Ethics Journal*, **9**(1-2), pp.141-157.
- Yang, J., Edwards, D.J. and Love, P.E.D. (2003) A computational intelligent fuzzy model approach for excavator cycle time simulation. *Automation in Construction*, **12**(6), pp.725-735.
- Yeung, J.F., Chan, A.P. and Chan, D.W. (2008) Establishing quantitative indicators for measuring the partnering performance of construction projects in Hong Kong. *Construction Management and Economics*, **26**(3), pp.277-301.
- Yu, W.-d. and Lin, H.-w. (2006) A VaFALCON neuro-fuzzy system for mining of incomplete construction databases. *Automation in Construction*, **15**(1), pp.20-32.

**Table 1: Some Examples of Cost Growth in Construction Projects**

<i>Project</i>	<i>Estimated Cost</i> <i>(in millions)</i>	<i>Final Cost</i> <i>(in millions)</i>	<i>% Overrun</i>
Sydney Opera House	AUD 7	AUD 102	1357
Nat West Tower	£15	£115	667
Thames Barrier Project	£23	£461	1904
Scottish Parliament	£195*	£414	112
British Library	£142	£511	260

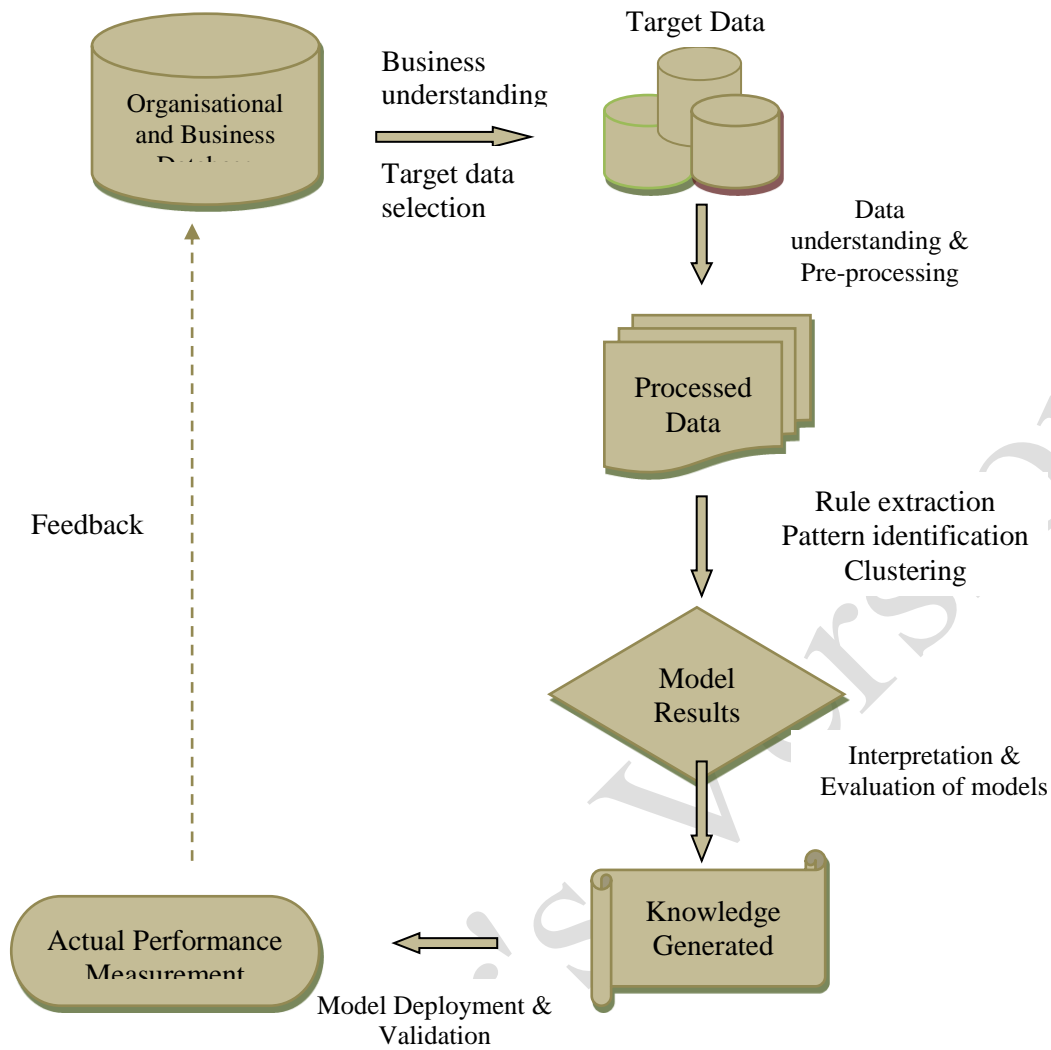
\*September 2000 estimate. Initially stated cost was about £40 million Source: Audit Scotland (2004)





**Figure 1: Conceptual Model for Understanding Cost Growth on Large public Projects**

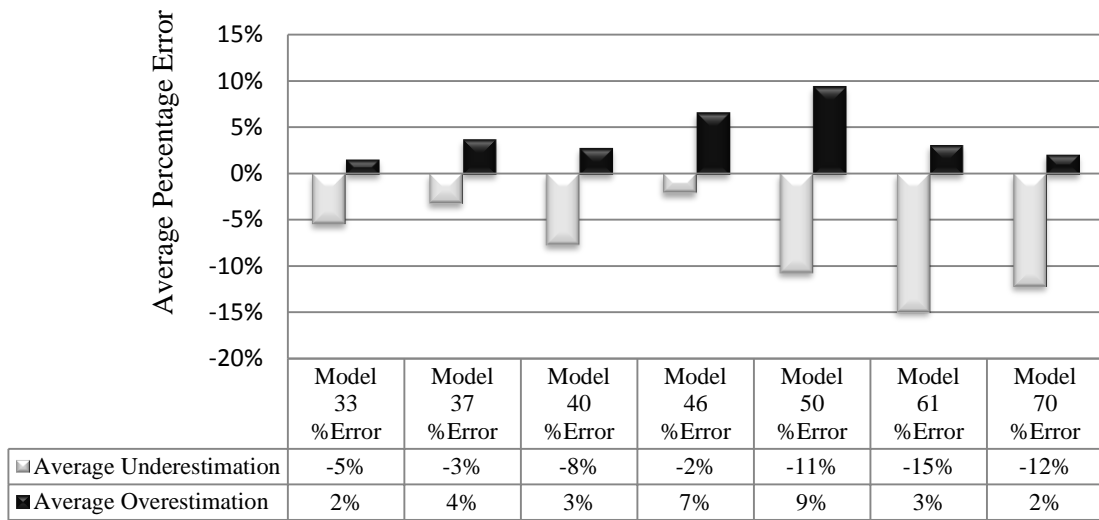
Source: *Ahiaga-Dagbui and Smith (2013)*



**Figure 2: The Generic Data Mining Process**

*Source: Ahiaga-Dagbui and Smith, 2013*

## Performance of Retained Models



**Figure 3: Performance of Selected Models**

**Table 2: Validation Results of the Best Model (Model 33)**

<i>Validation Case</i>	<i>Actual Final Cost</i>	<i>Final Cost predicted</i>	<i>Model Error</i>	<i>Model Absolute % Error</i>
1	£ 4,912,649	£ 5,120,943	-£ 208,294	4.24%
2	£ 1,617,225	£ 1,617,805	-£ 580	0.04%
3	£ 11,277,470	£ 10,743,624	£ 533,846	4.73%
4	£ 2,110,260	£ 2,136,125	-£ 25,865	1.23%
5	£ 5,398,965	£ 5,425,142	-£ 26,177	0.48%
6	£ 180,532	£ 181,214	-£ 681	0.38%
7	£ 2,572,564	£ 2,530,178	£ 42,386	1.65%
8	£ 1,440,593	£ 1,372,864	£ 67,729	4.70%
9	£ 3,842,258	£ 3,793,851	£ 48,407	1.26%
10	£ 4,194,219	£ 4,131,285	£ 62,934	1.50%
11	£ 375,170	£ 387,731	-£ 12,561	3.35%
12	£ 50,637	£ 51,502	-£ 865	1.71%
13	£ 24,479	£ 22,017	£ 2,462	10.06%
14	£ 858,112	£ 824,334	£ 33,779	3.94%
15	£ 21,798	£ 18,344	£ 3,454	15.85%
<b>Average Absolute % Error</b>				<b>3.67%</b>

<sup>i</sup> identity, logistic, tanh, exponential and sine activation functions