

Effectiveness of OLAP-based cost data management in construction cost estimate

S.W. Moon ^{a,*}, J.S. Kim ^b, K.N. Kwon ^{a,b}

^a Department of Civil Engineering, Pusan National University, Pusan, 609-735, Republic of Korea

^b Civil Design Team, Hyundai Development Co, Seoul, 135-881, Republic of Korea

Accepted 20 July 2006

Abstract

Historical cost data offer important information on the performance of past construction work, while current information technology makes it easier to develop databases. The database, however, needs to be utilized more, by providing a functional environment of probability analysis. The objective of this paper is to improve the effectiveness of utilizing historical cost data in an analytical OLAP (On-Line Analytical Processing) environment. A probability model has been prepared to support the functionality of OLAP. A prototype of the cost data management environment, Cost Data Management System (CDMS), has been developed to test the advantage of OLAP in cost estimate. A case study was carried out on the cost estimate of a sample project using the prototype. The results show that the OLAP environment can help understand the uncertainties in construction cost estimate, and provide a way for projecting more reliable construction costs.

© 2006 Elsevier B.V. All rights reserved.

Keywords: Historical cost data; Construction cost estimate; Uncertainty; Probability model; Data warehouse; Data cube; OLAP

1. Introduction

1.1. Problems and needs

Evolution of construction information systems has made it easier to gather historical cost data. Historical cost databases, which are being widely developed, offer a means to make an objective, rather than subjective, decision, by providing information on the performance of previous construction works [1].

Despite the increased development of databases in the construction industry, efforts to utilize cost data are still, however, less than satisfactory. Data stored in cost databases are vast so that an analytical approach is required to extract meaningful information. OLAP (On-Line Analytical Processing) supports rapid analysis by storing and refining a massive volume of cost data in a data warehouse. This can serve as a technical structure with which a construction manager can predict construction costs more accurately.

In this study, an OLAP environment is presented to understand the probabilistic nature of construction work and provide more reliable cost estimation.

1.2. Objectives

The objective of this paper is to improve the effectiveness of utilizing historical cost data in an analytical OLAP environment. A probability model has been prepared to support the functionality of OLAP. The uniqueness of the model is that the correlation between construction work items is applied to adjust the estimated construction cost.

A prototype of the cost data management environment, Cost Data Management System (CDMS), has been developed to test the advantage of OLAP in cost estimate. Here, the CDMS works as an enabler for the probability cost evaluation model described in this paper. Using the functions provided in the CDMS, the cost data stored in the data warehouse can be converted into more meaningful information for decision making.

2. Literature review

In the construction sector, research on data warehousing and OLAP has been conducted to effectively manage the data generated during the construction process. The analytical function of OLAP, combined with the data storage capability of a data warehouse, provides an opportunity to improve data

* Corresponding author.

E-mail addresses: sngwmoon@pusan.ac.kr (S.W. Moon),
jskim@hyundai-dvp.com (J.S. Kim), kinamkwon@hanmail.net (K.N. Kwon).

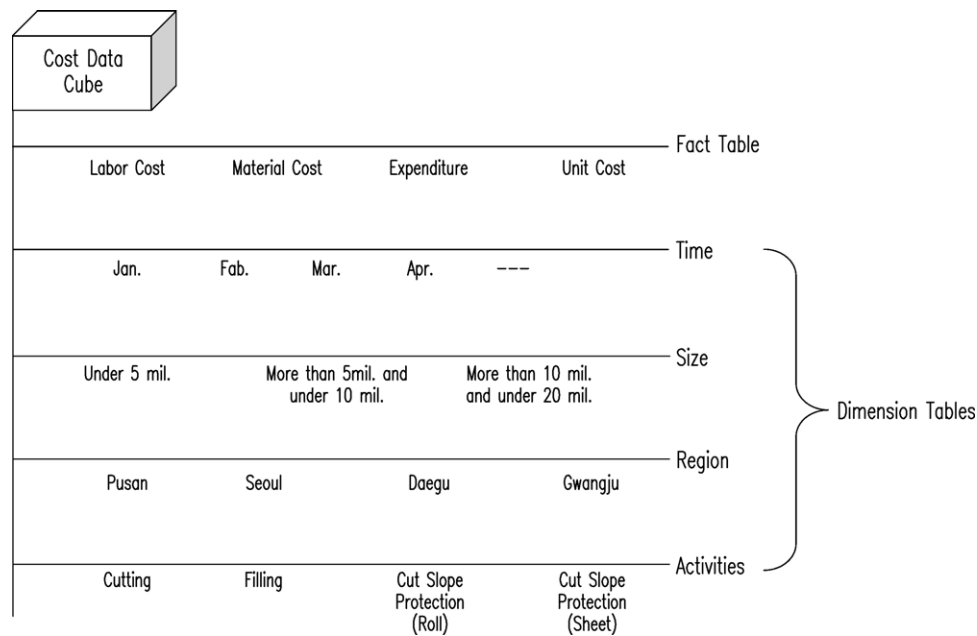


Fig. 1. Composition of cost data cube.

control in the construction sector. Confidence in the merit of utilizing this environment is growing due to the rapid development of construction information technologies [2].

The data warehousing and OLAP technologies have been often applied to assist the decision support systems in construction. Chau et al. [3] studied the applicability of integrating data warehousing with a construction decision support system. Ahmad et al. [4] tried to implement the data warehousing technology in developing a prototype decision support system for selecting an optimal residential housing site. Also, using an OLAP tool, they tried to generate a variety of information out of material inventory data. Zhiliang et al. [5] explored the utilization of electronic documents using the technology.

Applying data warehousing and OLAP, Yoo et al. [6] developed a performance evaluation system for construction management. In the system, the individual conditions of construction work items are taken into consideration to verify the degree of performance. This study prepared a balanced performance table based on a multidimensional analysis cube of data schema and analyzed the project performance from a variety of managerial standpoints.

One particular area that requires the application of OLAP is Customer Relationship Management (CRM) [7]. CRM is designed to manage customers' data for marketing purposes. Comprehensive data search and analysis is possible for construction marketing with the application of data warehousing and OLAP technologies.

As shown from its use in the construction sector, data warehousing and OLAP can process and analyze diverse and complicated data, and supply information that can help construction managers make more sensible decisions. In this paper, OLAP technology, centered on a cost data warehouse, is applied to build a user environment for storing and analyzing construction cost data.

3. Technologies for cost data analysis and application

3.1. Composition of a cost data cube in data warehouse

A data warehouse represents a data storage area in which data are collected and stored on a subject-by-subject basis from the individual corporate databases [8,9]. A data warehouse stores data in a multidimensional structure called a data cube. In the data cube, a dimension refers to the storage area for each subject, while a fact signifies the type of data included in the area for the subject. As Russell [10] discussed, the sorting and storage technique reduces the needs to search related databases every time they are required from many different perspectives.

In the case of cost data, a data cube can have four types of dimensions: (1) time, (2) size, (3) region, and (4) work breakdown structure (WBS) (Fig. 1). The time dimension stores construction costs by year or month. The size dimension classifies cost data according to the size of construction works, e.g. 500 million, one billion or 10 billion Korean won. The region dimension stores information about location (e.g. city or country) in which construction work is executed. Finally, the WBS dimension stores the classification of work items based on different levels of division. This dimension is used for calculating unit cost of individual work items.

Time, size and regional dimensions are applicable for developing cost indices. However, they will not be considered in this paper. Rather, this paper will focus on the probability aspect of a cost data set.

3.2. Composition of schema

A star schema is a type of multidimensional data structure [11]. In the star schema, a single fact table is joined by several

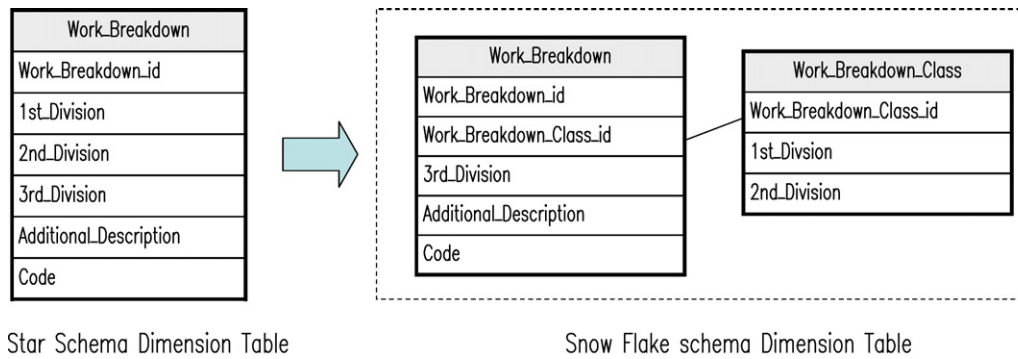


Fig. 2. Star and snowflake schemata for work breakdown structure.

dimension tables to form a one-to-many relationship. Usually, the fact table, the center of the star schema, is abnormalized. According to the various attributes of each dimension, the numerical probability analysis data can be stored in the fact table. For example, cost-related attributes, such as time, location or WBS, are stored in the dimension table, while the individual unit costs, as well as their probability data, can be stored in the fact table. In this way, associated data can be retrieved from several data tables through one single inquiry.

When the contents of the dimension table are oversized, they can be divided into smaller tables of specific purposes, i.e. a star schema can be broken down into a snowflake schema [12]. Fig. 2 shows a star schema table for ‘work_breakdown’, divided into two tables of snowflake schema for ‘work_breakdown’ and ‘work_breakdown_class’. Being composed of multiple connections of normalized small dimension tables in the snowflake schema, the storage space can be minimized, leading to improved search and reference performance.

More importantly, the structure of data cubes is not in a typical format. Rather, schema should be designed to meet the specific requirements when the OLAP environment is implemented. Most OLAP tools support these functions [13].

3.3. OLAP

OLAP accesses a data warehouse, analyzes data on a specific subject and accumulates them in a storage area of the data warehouse [14]. OLAP can describe information from data analysis in a variety of formats to provide useful information so that construction managers can make rational decisions. Therefore, the user can minimize the effort in finding and analyzing the

meaning of data. The role of OLAP is increasing because it is useful for retrieving significant information by analyzing the vast volume and variety of data generated in information management systems.

In this study, OLAP is designed to implement an application environment for the efficient exploitation of cost data. Meanwhile, a data warehouse is designed to improve data storage in terms of accuracy, consistency and availability. Thus, OLAP and data warehousing perform mutually complementary functions to ensure rapid data processing in its storage and utilization of cost data.

4. Probability analysis model of cost data

4.1. Cost data for construction work items

To explain the probability analysis method of cost data, a sample of four work types, i.e. (1) cutting, (2) embankment, (3) slope protection (roll) and (4) slope protection (sheet), are presented in Table 1. Five sets of historical cost data for the work items were collected from five different highway construction projects of a private company.

The means and deviations of the data sets were calculated and used as unit costs. For a sample project, each work item is given arbitrary quantities. The amounts, which are the mean costs for each item, were obtained by multiplying the quantity by the mean unit cost.

4.2. Correlation between construction work items

Historical cost data, which are a set of past data, have probability distributions. Therefore, when construction costs are

Table 1
Example of construction work items

| Work items | Unit | Unit cost | Quantities | Cost amounts (expected mean) | Deviations |
|--------------------------|------|-----------|------------|------------------------------|------------|
| Cutting | M3 | 576 | 81,289 | 46,822,464 | 4,672,527 |
| Embankment | M3 | 1100 | 284,203 | 312,623,300 | 25,419,889 |
| Slope protection (roll) | M2 | 2500 | 5123 | 12,807,500 | 857,242 |
| Slope protection (sheet) | M2 | 4040 | 16 | 64,640 | 1280 |
| Sum | | | | 372,317,904 | 25,859,970 |

Table 2
Correlation coefficients between work items

| Work items | Cutting | Embankment | Slope protection (roll) | Slope protection (sheet) |
|--------------------------|---------|------------|-------------------------|--------------------------|
| Cutting | 1.00 | 0.47 | −0.10 | 0.64 |
| Embankment | 0.47 | 1.00 | 0.13 | 0.56 |
| Slope protection (roll) | −0.10 | 0.13 | 1.00 | −0.30 |
| Slope protection (sheet) | 0.64 | 0.56 | −0.30 | 1.00 |

Table 3
Adjustment of construction cost under equivalent uncertainties

| Description | $\rho=+1$ | $\rho=0$ | $\rho=-1$ | $\rho=\text{actual}$ |
|-----------------|-------------|-------------|-------------|----------------------|
| Adjusted amount | 448,546,985 | 372,317,904 | 282,363,217 | 405,969,144 |
| Increment | +76,229,081 | 0 | -89,954,687 | +33,651,240 |

estimated by simply applying unit costs in terms of median or mean, the uncertainty inherent in the construction process cannot be understood. To consider uncertainties, the probability distribution of deviation or variance should be taken into consideration, together with the mean value.

Probability values for specific construction work can be calculated using Eqs. (1)–(3) [15]. From the viewpoint of probability, the expected total construction cost can be obtained from the sum of the mean costs for each item. On the other hand, the overall deviation can be obtained by summing and square rooting the variances of work items. Therefore, in Table 1, the sum of cost amounts is the mere aggregation of individual cost amounts, while the sum of deviations obtained using Eq. (3).

Expected total construction cost:

$$X = x_1 + x_2 + x_3 + \cdots + x_n \quad (1)$$

Overall variance:

$$\sigma^2 = \sigma_1^2 + \sigma_2^2 + \sigma_3^2 + \cdots + \sigma_n^2 = \sum_{i=1}^n \sigma_i^2 \quad (2)$$

Overall deviation:

$$\sigma = \sqrt{\sigma_1^2 + \sigma_2^2 + \sigma_3^2 + \cdots + \sigma_n^2} = \sqrt{\sum_{i=1}^n \sigma_i^2} \quad (3)$$

Considering the fact that a construction project represents a process in which the activities of predecessors and successors are interconnected, construction work items are mutually correlated. The degree of correlation varies depending on the types of construction project. In other words, there are correlations between construction work items and, therefore, the interface between work items becomes a major management factor. In the case of concrete pavement, for example, when the earth work

(a) Entity Relation Diagram

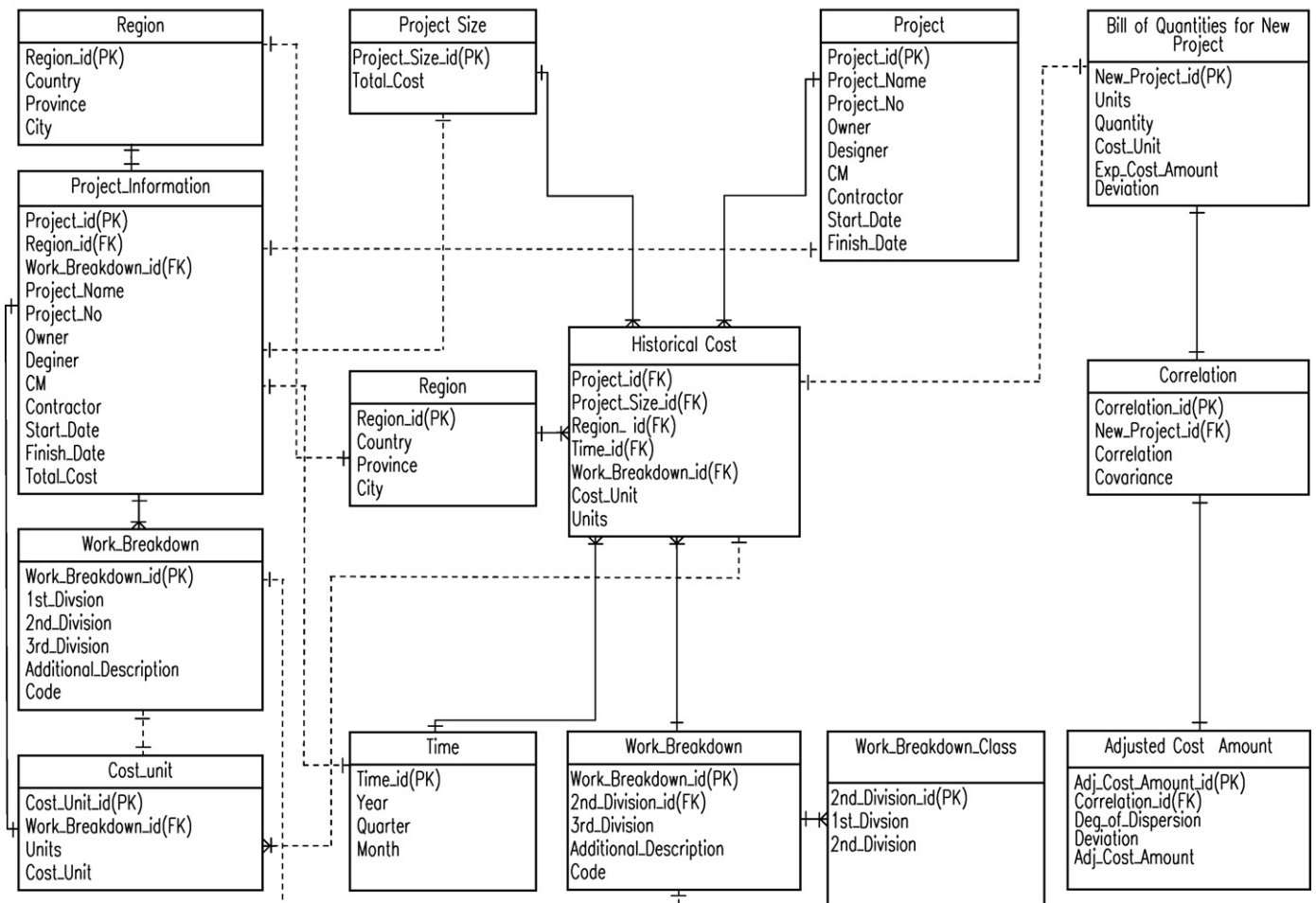


Fig. 3. Structure of cost data. (a) Entity relation diagram. (b) Cost data cube.

(b) Cost Data Cube

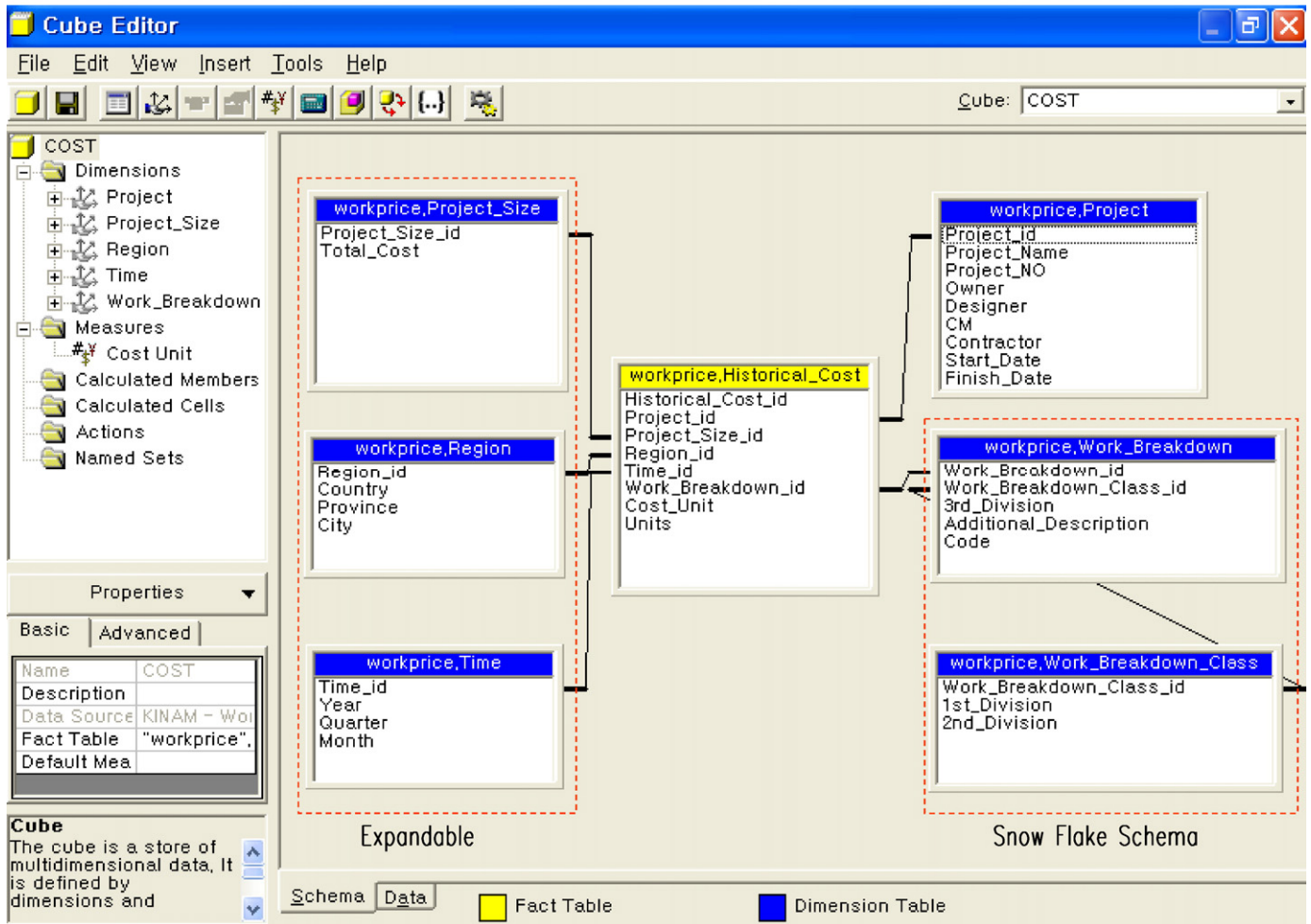


Fig. 3 (continued).

proceeds rapidly, the concrete pavement work also proceeds quickly. Accordingly, the correlation should be taken into account for more accurate cost estimation.

The correlation coefficient (ρ_{xy}) indicates the degree of correlation between two particular work items and can be obtained by Eq. (4). The correlation coefficient can have a value between -1 and 1 . When $\rho_{xy}=0$, no correlation exists between the two work items. When $\rho_{xy}=+1$, the two variables are fully positively correlated and, when $\rho_{xy}=-1$, the two variables are fully negatively correlated [16]. As the correlation coefficient nears zero, the correlation also nears zero. In general, correlations between work items are assumed to be zero for cost estimation.

$$\rho_{xy} = \frac{\sigma_{xy}}{\sigma_x \sigma_y} \quad (4)$$

Here, ρ_{xy} : correlation coefficient of variables x and y

σ_{xy} covariance of variables x and y
 σ_x standard deviation of variable x
 σ_y standard deviation of variable y

The correlation matrix shown in Table 2 represents the correlation coefficients for construction work items. These coefficients were calculated using Eq. (4). The matrix indicates that cutting and slope protection (roll) are negatively correlated (-0.10), and cutting and slope protection (sheet) are positively correlated (0.64).

When correlation coefficients are given, as in Table 2, Eq. (2) should be modified to Eq. (5) to consider the interrelationship between two work items. Therefore, when projecting construction costs, the degree of the uncertainties in investment varies depending on the interactions between construction work activities.

$$\begin{aligned} \sigma_p^2 &= \sigma_1^2 + \sigma_2^2 + \cdots + \sigma_n^2 + 2\sigma_{1,2} + 2\sigma_{2,3} + \cdots + 2\sigma_{n-1,n} \\ &= \sigma_1\sigma_1 + \sigma_2\sigma_2 + \cdots + \sigma_n\sigma_n + 2\rho_{1,2}\sigma_1\sigma_2 \\ &\quad + 2\rho_{2,3}\sigma_2\sigma_3 + \cdots + 2\rho_{n-1,n}\sigma_{n-1}\sigma_n \\ &= \sum_{i=1}^n \sum_{j=1}^n \sigma_{ij} \end{aligned} \quad (5)$$

Here

σ_i standard deviation of work item i
 $\sigma_{i,j}$ covariance of work item i and work item j ($\sigma_{n-1,n} = \sigma_{n,n-1}$)

4.3. Construction costs prediction under equivalent risks

The variance of total construction costs describes the degree of uncertainties, while uncertainties in turn represent the risks inherent in construction works. Table 3 shows that the degree of uncertainties can vary depending on correlation coefficients. To ensure equivalent uncertainties in estimating construction costs, the dispersion in the cost data set should be understood and the costs adjusted accordingly.

The degree of dispersion (DoD) represents the degree to which the data are scattered with the expected mean value as a center point. The DoD can be simply calculated by dividing the deviation by the expected total construction cost. Eq. (6) shows the calculation of DoD when correlation coefficient is zero.

$$\text{When } \sigma = 0, \text{ DoD} = \frac{\text{Deviation}}{\text{Expected Total Construction Cost}} = \frac{25,859,970}{372,317,904} = 0.069 \quad (6)$$

As the DoD increases, the dispersion becomes wider, and vice versa. A larger DoD means that there is a higher possibility that the actual construction cost becomes larger or smaller than cost estimates. Since the decision theory is based on the approach from a conservative standpoint, the larger DoD signifies a higher risk. To reduce the DoD, expected construction costs should be increased.

The figure 0.069, obtained in Eq. (6), represents a DoD when the correlation coefficient is zero. When the actual correlation coefficients, shown in Table 2, are applied, however, the DoD will differ. In this case, the expected total construction cost should be increased or decreased to achieve the same DoD as that attained when the correlation coefficient is zero.

Table 3 shows the adjusted construction costs when correlation coefficients are $\rho = +1$, '0', -1 , and 'actual'. The construction costs are adjusted to make a DoD equivalent to that when the

correlation coefficient is zero. Here, when the actual correlation coefficients in Table 2 are applied, the construction cost is increased by 33,651,240 Korean won to 405,969,144 Korean won to get the same level of DoD, i.e. the overall deviation rose from 25,859,970 to 28,011,871. To cope with the increased uncertainties, the construction cost increased about 9%.

5. Structure of cost data warehouse

Data warehouse modeling is a process to generate multidimensional information. It is a process in which a set of dimensions is established and information organized to meet the needs of data management. Once the data warehouse structure has been established, data should be periodically retrieved and saved from corporate databases.

The data warehouse searches relevant databases and extracts appropriate data, including project information, construction periods, construction amounts and work locations, not to mention cost data. Fig. 3(a) shows the entity relation diagram (ERD) for detailed data structure. The ERD shows the SQL tables associated with 1) historical cost data in a corporate database, 2) star schema in a data warehouse, and 3) cost estimate in OLAP.

Fig. 3(b) shows the data cube for the storage of cost data. The data cube indicates the basic composition of saving cost data and other related information in the data warehouse. In this study, the subjects of information, such as time, size and region, are stored in the dimension tables. The WBS is represented in two normalized snowflake schema tables. The fact table stores units and unit costs.

6. OLAP-based cost data management

6.1. Structure of OLAP environment

The CDMS prototype presented in this study has been developed based on OLAP and data warehousing technologies (Fig. 4)). Parts of the data warehouse and OLAP were developed by applying Microsoft SQL Server 2000 and Analysis Service, respectively. The subject specific storage area of the cost data warehouse can store cost-related data extracted from various corporate databases, including construction management,

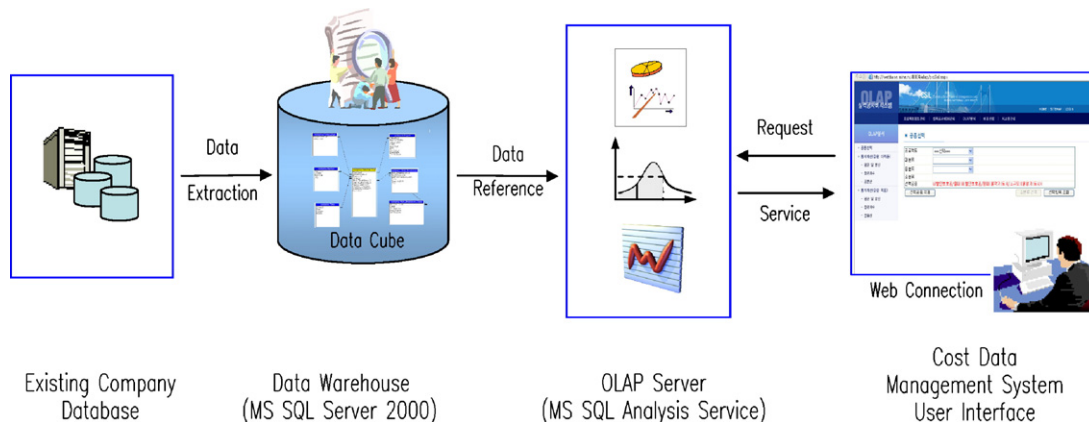


Fig. 4. Architecture of cost data management system.

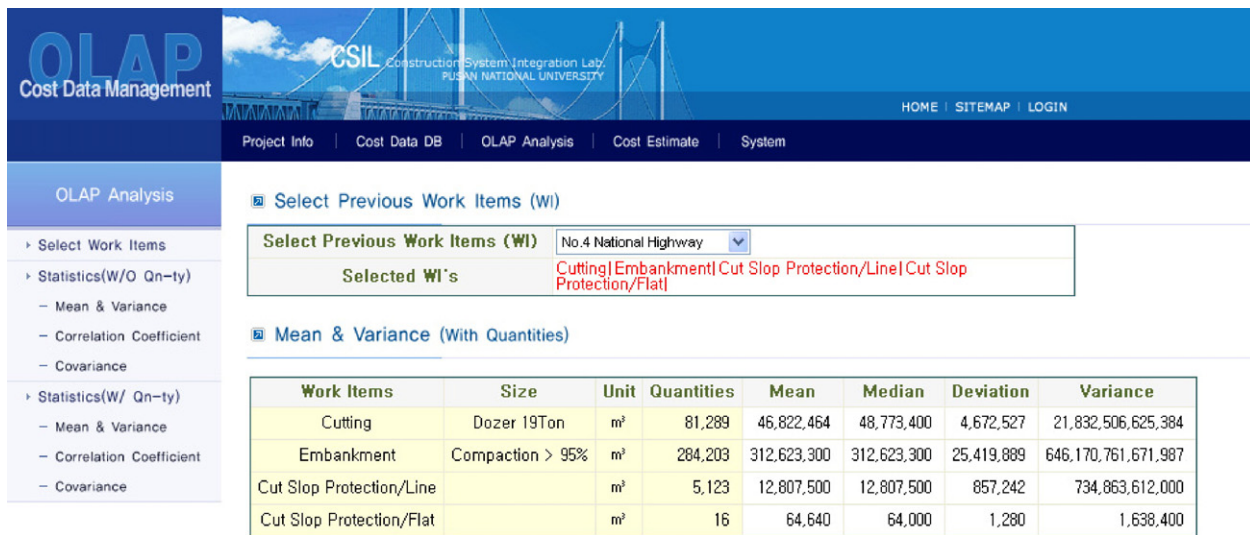


Fig. 5. Estimate of cost in OLAP.

feasibility studies, etc. The analytical capabilities of the CDMS were developed with ease, using the various statistical functions provided by the OLAP tool. These capabilities generate useful information quickly by calculating cost data acquired through SQL queries.

6.2. Probability analysis in OLAP

To implement an OLAP environment, historical cost data should be first refined and stored in the cost data cube. The schema table contains information necessary for construction

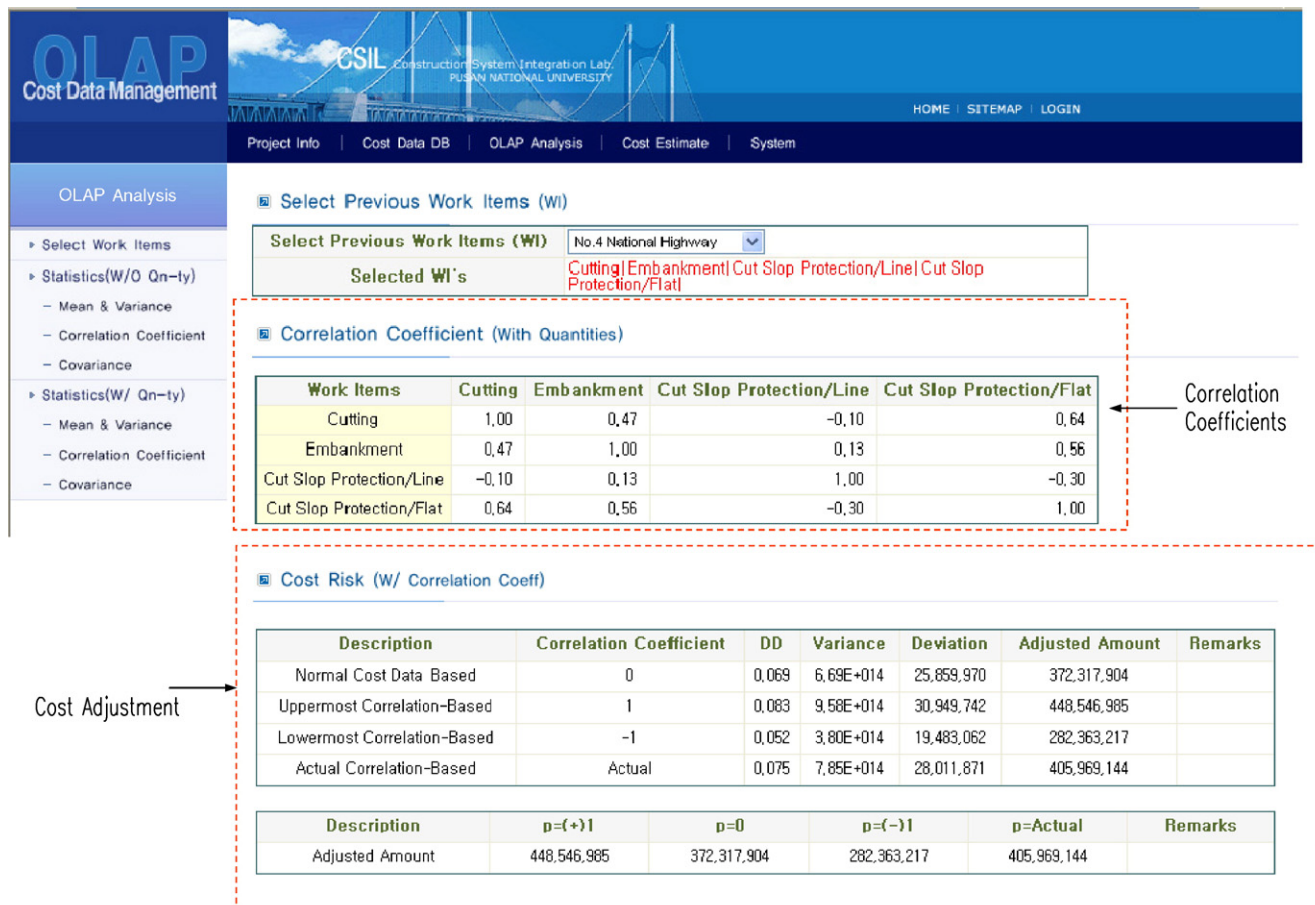


Fig. 6. Adjustment of cost considering correlation coefficients.

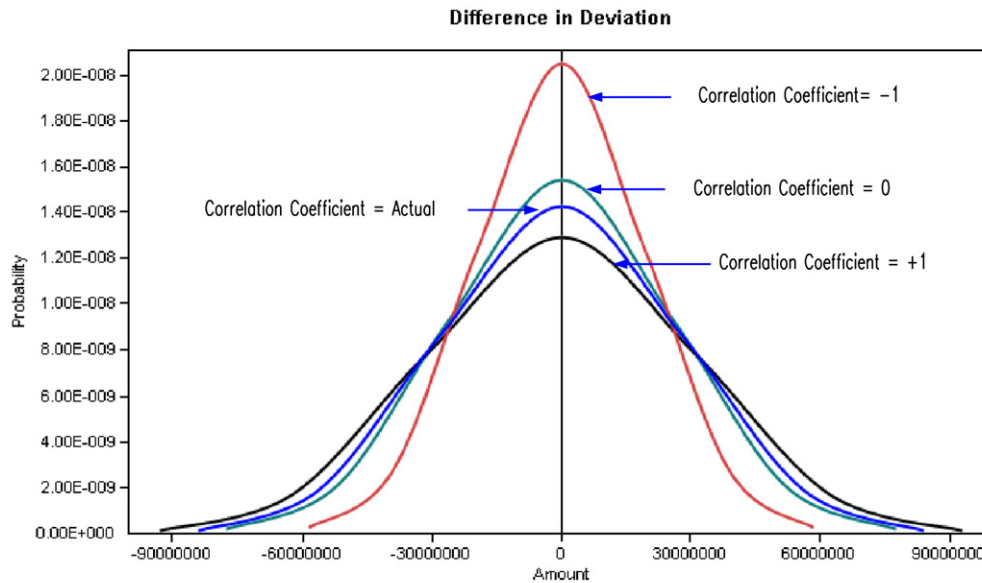


Fig. 7. Cost estimate—differences in variance.

cost estimate, including project information, work breakdown, units and unit costs.

Fig. 5 shows the costs and deviations of four construction work items: (1) cutting, (2) embanking, (3) slope protection (roll) and (4) slope protection (sheet). As discussed previously, when estimating construction costs, the cost amount for each work item can be predicted by multiplying quantities with the probability values of median or mean. The deviation of each work item is obtained multiplying the deviation of unit costs by quantities.

In Fig. 5, the unit costs are presented in terms of mean and median values, which are obtained through SQL queries from the cost data cube in the CDMS. The total construction cost adds up to 372,317,904 Korean won when mean values are applied as unit costs.

6.3. Correlation coefficient and cost adjustment

Various probability analyses, which were not possible in the past, can now be made through cost data accumulation and the OLAP functions can swiftly provide the results of probability analysis for cost data. The correlation coefficients between work items are one of the benefits made available through the OLAP environment.

The correlation part in Fig. 6 indicates the correlation coefficients calculated using Eq. (4), and represent the correlations between work items. The correlation coefficients in Table 2 can be obtained using the OLAP functions, i.e. the values of the coefficients for correlation or covariance were obtained using the statistical functions provided by the OLAP tool, such as 'Correlation()' and 'Covariance()'.

The cost adjustment part in Fig. 6 shows DoD values, dispersions and deviations obtained by applying the values of the correlation coefficients between two particular work items. To

understand the effects of correlation on construction cost estimate, this paper observed the changes in the dispersion by applying the four correlation coefficients, including '+1', '0', '-1', and 'actual'. The difference in dispersion indicates the changes in DoD values.

When the correlation coefficients were applied, the DoD values for each of the correlation coefficients were found to be 0.069, 0.083, 0.052 and 0.075, respectively. The construction costs were adjusted in consideration of the DoD to obtain the same level of uncertainties. This shows that, when actual correlation coefficients are applied, the construction amount should be increased by 33,651,240 Korean won, from 372,317,904 to 405,969,144 Korean won.

6.4. Graph analysis

Fig. 7 shows four standard normal distribution curves when the four correlation coefficients were applied for cost estimate. In the figure, the correlation coefficient for the sample project in this paper exists between (+)1 and 0. These distribution curves show that, as the correlation coefficient approaches -1, the deviation from the expected total construction cost becomes smaller, increasing the reliability of the cost estimate. This result indicates the effect of correlation in construction cost estimate.

7. Conclusion

Owing to the development of construction information systems, most construction operations are conducted via digital processing, in which vast amounts of cost-related data are generated. Meanwhile, retrieval of valuable and worthwhile information from the accumulated data is a separate problem, where the use of data stored in a database has become an important company asset.

Thus, this study has applied OLAP technology, centered on data warehousing, to present a new approach to construction cost data management. The originality of this paper is in the fact that the study demonstrated a new way of utilizing historical cost data. The originality can be discussed in two aspects. First, the probability model presented in this paper takes advantage of the correlation between two work items. Correlation was used to understand the uncertainties residing in construction work, which in turn used to adjust construction cost estimate. Second, understanding the probability characteristics of the vast amount of cost, however, requires more analytical functions. The prototype of the CDMS showed that the historical cost data can be effectively utilized by implementing the analytical functions of OLAP together with the data storing capability of warehouse.

Implementing the prototype of the CDMS has brought about the following results:

- (a) The uncertainties residing in the cost estimate were measured considering the correlation in the construction work items.
- (b) The variance was used to adjust the cost estimate to have an equivalent level of uncertainties.
- (c) Application of the data analysis function of OLAP allows for various probability analyses of cost data.
- (d) OLAP technology processes the raw cost data to provide useful information for strategic decision making in construction cost estimate.

The limitations of this study are that (1) it does not consider the cost indices for cost adjustment and (2) the prototype was built without linkage to actual corporate databases. Despite these drawbacks, the study presents a guideline to utilize cost data in a situation where systematic exploitation of accumulated data is necessary. The technology of cost data collection and utilization requires further study to support various types of decision making, such as cost estimates, feasibility analysis, etc.

References

- [1] S.W. Moon, S.O. Shin, T.Y. Park, Introduction to the KOLAND land development feasibility study system, International Joint Study and Conference for Comparative Land Use and Urban Developments in the U.S. and Korea, June 1–22, Korea Land Corporation, Daejeon, Korea, 2000, pp. 229–242.
- [2] J.K. Lee, H.S. Lee, Principles and strategies for applying data warehouse technology to construction industry, *Architectural Research* 5 (1) (2003) 61–68.
- [3] K.W. Chau, Y. Cao, M. Anson, J. Zhang, Application of data warehouse and decision support system in construction management, *Automation in Construction* 12 (2002) 213–224.
- [4] I. Ahmad, S. Azhar, P. Lukauskis, Development of a decision support system using data warehousing to assist builders/developers in site selection, *Automation in Construction* 13 (2004) 525–542.
- [5] M. Zhiliang, K.D. Wong, L. Heng, Y. Jun, Utilizing exchanged documents in construction projects for decision support based on data warehousing technique, *Automation in Construction* 14 (2005) 405–412.
- [6] J.H. Yoo, S.H. Song, W.H. Lyoo, H.S. Lee, Development of performance analysis system for construction projects using a data warehousing technology, *Korean Journal of Construction Engineering and Management* 6 (1) (2005) 89–98.
- [7] R. Buck-Emden, P. Zencke, *mySAP CRM: The Official Guidebook to SAP CRM Release 4.0*, SAP Press, Rockville, MD, 2004.
- [8] W.H. Inmon, *Building the Data Warehouse*, 3rd ed Wiley, New York, New York, 2004.
- [9] A. Khan, *Data Warehousing 101: Concepts and Implementation*, The Canton Group, Lincoln, NE, 2003.
- [10] K. Russell, Data cubes, *Computerworld* 38 (2004) 32–32.
- [11] C. Seidman, *Data Mining with Microsoft SQL Server 2000*, Microsoft Press, Buffalo, New York, 2001.
- [12] M. Levene, G. Loizou, Why is the snowflake schema a good data warehouse design? *Information Systems* 28 (2003) 225–240.
- [13] O. Train, R. Jacobson, *Microsoft SQL Server 2000 Analysis Services Step by Step*, Microsoft Press, Buffalo, New York, 2000.
- [14] J. Lechtenborger, G. Vossen, Multidimensional normal forms for data warehouse design, *Information Systems* 28 (2003) 415–434.
- [15] Project Management Institute (PMI), *Project Management Body of Knowledge*, Newtown Square, PA, , 1996.
- [16] J.R. Evans, W.M. Lindsay, *The Management and Control of Quality*, West Publishing Co, St. Paul, MN, 1989.