# USING QUANTITY ESTIMATE STATISTICAL MODELS FOR INFRASTRUCTURE LONG RANGE COST MANAGEMENT

**Jui-Sheng Chou [1], Min Peng [2], James T. O'Connor [3] and Khali R. Persad [4]**

[1] Assistant Professor, Department of Business Administration, National Chung Cheng University, Chia-Yi, TAIWAN
[2] Ph.D. Candidate, Department of Civil, Architectural and Environmental Engineering, University of Texas, Austin, TX, USA
[3] Professor, Department of Civil, Architectural and Environmental Engineering, University of Texas, Austin, TX, USA
[4] Associate Researcher, Center for Transportation Research, University of Texas, Austin, TX, USA

Correspond to rayson.chou@mail.utexas.edu

**ABSTRACT :** Effective cost management requires reliable cost estimates at every stage of project development. The primary purpose of this research is to develop systematic modeling procedures and an automatic computing program for infrastructure estimating in the Texas Department of Transportation (TxDOT). The computing system toggles between project input information and segregated district unit prices for highway work item quantity estimates associated with earthwork and landscape, subgrade treatments and base, surface courses and pavement, structures, miscellaneous construction, and lighting, signing, markings and signals. This quantity-based approach was chosen because of the conventional approach lacking of quantity information until primary design is complete.

*Key words : Quantity Estimates, Data Analysis, Highway Projects, Cost Management, Database*

## 1. INTRODUCTION

Preliminary cost estimates are crucial to the viability of a project progressing beyond the planning stage, yet little data are available to develop an accurate budget [1]. Many State Departments of Transportation (DOTs) have experienced highly visible projects that have suffered from excessive cost overruns. Inaccurate preliminary cost estimates for highway projects deeply affect financial operations of these organizations in the United States due to marginal budgets [2]. Many studies of project cost estimates have found the final total cost incurred in designing and constructing projects of all types *almost always* exceeds the amounts estimated .

A research study of 258 transportation infrastructure projects among modernized countries led to the following observations [3]: Costs are underestimated in 9 out of 10 transportation infrastructure projects; For road projects, actual costs are on average 20% higher than estimated costs with a standard deviation of 30%; Cost underestimation appears to be a global phenomenon.

Estimating accuracy is closely related to the extent of information available at the time the estimate is developed. The conceptual estimate is often misleading because of the paucity of available information. In particular, a stronger case should be made for predicting and early quantity tracking for eliminating cost error by better exploiting readily available recent district work item unit prices.

## 2. RESEARCH OBJECTIVES

Many studies have investigated unit cost estimating relationships (CERs) between cost per unit quantity ($/lane mile) and plan quantity (quantity take-off) using either neural network or statistical methods [1][2][4][5][6][7][8]. These kinds of approaches avoid highly complex time-consuming "wild guesses" and generate acceptable preliminary cost estimates. But they do not exploit the accuracy of recent unit prices nor do they establish a quantity formulation upon which a more accurate, robust system could be built.

The motivation for use of parametric quantities for preliminary estimates is to exploit the accuracy of and access to historical unit prices and in addition, to promote effective continuous cost tracking and control by initiating quantity estimates during preliminary phases. This research study has pursued parametric quantity estimating for this purpose.

## 3. LITERATURE REVIEW

### 3.1 Timeframes of Estimating for Infrastructure Projects

Infrastructure projects often take longer to develop than other projects. Figure 1 is a simplified illustration of the major stages of development of larger projects. Estimates are required at several points in the process .
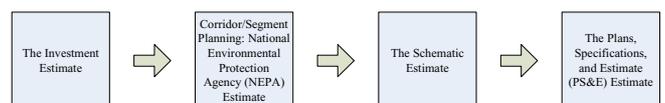


**Figure 1**. Stages of Development of a Major Transportation Project

A project typically first appears in the DOT's long-range plan which documents planned work for the next twenty years. Sources and adequacy of funding must be identified at the investment estimate stage. This initial estimate is usually based on historical costs per lane mile, and it is used to gauge the project's economic feasibility. High-feasibility projects are usually given priority and advanced faster. In effect, an unfavorable estimate can delay or terminate a project. Yet often such decisions are based on erroneous cost estimates.

During the National Environmental Protection Act (NEPA) process, estimates are needed to compare alternative layouts and environmental impacts. The cost estimate at NEPA stage is still very rudimentary, as NEPA regulations prohibit DOTs from developing one specific alternative further than others, in order not to prejudice the selection.

Upon environmental approval, a schematic estimate is prepared, sometimes based on project elements. The cost estimate at the schematic estimate stage determines the timing of funding for the project, as well as engineering design fees.

The last major stage of estimating is PS&E estimate. At this stage, findings from site investigations could add to the cost of the project. Design costs could add 5-15% to the total cost. Availability of property may affect the final configuration and significantly alter the project cost. For example, the $300 million M-59 project in Michigan had right of way (R.O.W.) costs of $150 million to relocate residences and businesses. Changes in the project estimate during this stage often affect its chance of receiving early funding.

## 3.2 State DOTs Estimating Approaches

Several approaches of cost estimating are currently adopted by state DOTs in the United States. The *lane-mile historic cost averages* is utilized by thirty-one state DOTs [9]. Because the estimates are based solely on historic lane-mile cost averages for similar projects and some unique characteristics of the project are ignored, the estimates are very rough.

The approach *conventional quantity-take-off and adjusted historical unit price* is used by several state DOTs [9]. Likewise, *component-level parametric unit price range with qualitative adjustment factors* is the fundamental cost component approach to combine unit price with conventional quantity-take-off. The component-level parametric unit price range is related to qualitative adjustment factors such as project location and type. The Connecticut DOT adopts the approach *work item unit price range according to quantity range estimates* by employing work item unit price ranges according to quantity instead of adjustment with qualitative factors. This estimating approach requires unit price ranges to be updated regularly and relatively accurate work quantities despite sketchy conceptual design information. Furthermore, quantitative adjustments for cost escalation, location, and other factors must be made. However, most of these methods described above cannot be applied for preliminary cost estimates when only conceptual design information is available.

## 3.3 Cost Growth Factors in Infrastructure Projects

The most common causes of cost growth for highway projects can be divided into three groups, namely project factors, organization factors, and estimate factors.

Project cost performance is directly related to project conditions. These project factors include changing economic or market conditions, project type, project complexity, project location, project size, duration of construction periods, scope changes, unforeseen engineering complications and constructability challenges, construction accessibility, restricted working hours, use of new technology, method of construction or construction techniques, and experimental or research items and special specifications or provisions [2][3][10][11][12][13].

Projects can be influenced as well by organization factors, including organizational capacity of the owner, designer, and/or contractor, contract type and context of contract, changes in regulatory requirements, disruption or discontinuity within the management team or local political leadership, lack of site familiarity by the design team, expertise of the consultants involved in the project, and poor communication between districts and head office [9][12][13][14][15][16].

Quality and timing of the estimate also influence cost performance. These estimate factors include timing of estimate cost data versus timing of expenditure, estimator-related factors (e.g. cognitive biases), estimating team experience, quality of cost information, time allowed to prepare the estimate, wide variability in contractor's (subcontractors') prices, lack of review of cost estimate by management, lack of adequate guidelines for estimating, and estimators' lack of data processing techniques [3][12][13][14][16][17].

## 3.4 Challenges of Preliminary Cost Estimates

Conventional approach lacks of quantity information and does not initiate the first quantity estimates until 90% to 95% design complete. Figure 2 depicts the problems with quantity "blackout" from project inception to 90% or 95% design complete. Quantity-based estimating system enables periodic quantity adjustment as projects evolve to later phases.
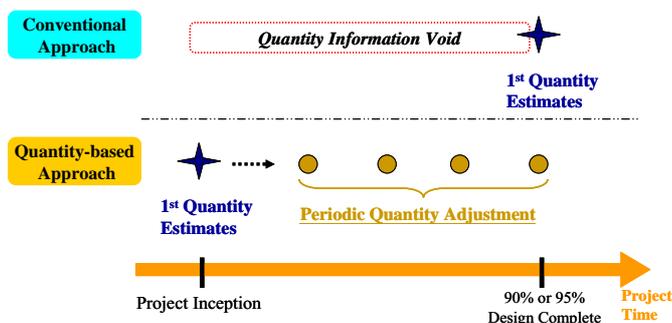


**Figure 2.** Problems with Quantity "Blackout"

# 4. QUANTITY-BASED STATISTICAL MODELS

## 4.1 Item-Level Quantity-Based Approach

Item-level quantity models known at project outset were developed by exploiting cost data stored in the TxDOT Design and Construction Information System (DCIS) from FY2001 through FY2003. The basic estimating parameters, which were identified based upon the statistical analysis, are documented for future input in quantity calculation. This approach provides an opportunity to embed the models within a relational database management system (RDBMS) for computation and data storage purposes.

## 4.2 Data Collection and Preprocessing

The TxDOT statewide computer network, Design and Construction Information System (DCIS) allows all districts within Texas to maintain project data in a standardized format. While a few projects had incomplete data, the DCIS was found to contain extremely useful and reliable data. Thus, the database system was used as the data resource for this research. Unit bid prices associated with corresponding work items were segregated in the data analysis. Hence, price inflation would not be a factor in the later analysis. The project data collected for this study consist of 545,920 records consolidated from 2,222 projects. Eighty-seven fields were extracted from this legacy database. The data analysis in this study was performed with the aid of the software Statistical Package for Social Sciences (SPSS 2003).

## 4.3 Statistical Model Formulation

A general parametric function used in model development was established as shown below through observation of a series of pilot studies.

$$Y = \beta_0 \left( X_1^{\beta_1} X_2^{\beta_2} \cdots X_n^{\beta_n} \right) e^{(\beta_{n+1}D_{N+1} + \beta_{N+2}D_{N+2} \cdots + \beta_{n+m}D_{n+m})} \varepsilon \tag{1}$$

Where,

| | |
|---|---|
| Y: | work item quantity |
| $\beta$'s: | estimated parameters |
| $X_1 \sim X_n$: | predictors representing numerical data |
| e: | exponential constant |
| $D_{n+1} \sim D_{n+m}$: | predictors representing categorical data |
| $\varepsilon$: | error term |

This parametric model is a particular type of nonlinear relationship that has firm grounding in economic theory. It is called the constant elasticity relationship or multiplicative relationship and can be transformed into a linear model by a logarithmic transformation [8][18][19]. In this study, a natural logarithm with base e (i.e. ln) was used. The equation was transformed therefore into the following linear model:

$$\ln(Y) = \ln(\beta_0) + \beta_1\ln(X_1) + \beta_2\ln(X_2) + \ldots + \beta_n\ln(X_n) + \beta_{n+1}D_{n+1} + \beta_{n+2}D_{n+2} + \ldots + \beta_{n+m}D_{n+m} + \ln(\varepsilon) \tag{2}$$

Transformation can clearly reduce the impact of outliers and help allay concerns about violating assumptions of normality and homoscedasticity in regression analysis when raw data exhibit a skewed pattern [8][19][20][21]. On the other hand, these two assumptions were not nearly as crucial as the need for independence. If a transformation is performed in a least squares regression, the resulting statistical properties (e.g., best, linear, unbiased estimates) are true only on the transformed values. Once the results are "back-transformed" to the original units, these statistical niceties are lost [22]. The aforementioned equation can be expressed as a general linear regression model in terms of a linear combination of predictors as below:

$$Y` = X`\beta` + \varepsilon` \tag{3}$$

Where,
$Y` = \ln(Y)$
$X` = [1 \ln(X_1) \ln(X_2)...\ln(X_n) \ln(D_{n+1})..\ln(D_m)]$ ; row vector
$\beta` = [\ln(\beta_0) \ \beta_1 \ \beta_2...\beta_n \ \beta_{n+1}...\beta_{n+m}]^T$ ; column vector
$\varepsilon` = \ln(\varepsilon)$

Cross-product interaction regression models were implemented in the data analyses to improve R-square and to explain interaction effects on the variability of the response variable. When adding interaction terms to the regression models, caution should be exercised on the existence of multicollinearities between some of the predictor variables and some of the interactions terms [19]. A prior knowledge concerning practical interpretation of the interaction terms that are most likely to influence the response variable should be utilized whenever possible.

Stepwise regression using probability of F-test as entry (0.05) and removal (0.10) criteria was employed throughout the model development in predictor selection. Multicollinearity exists as a fairly strong linear relationship between two or more independent variables [18]. The degree to which the predictors are correlated among themselves can affect regression results and make estimation unstable. The strength of collinearity among the independent variables in the models was measured by a statistic called tolerance. SPSS [23] suggests if any of the tolerances are less than 0.1, multicollinearity may be a problem.

Upon obtaining regression model results, several assumptions ought to be revisited alongside hypothesis testing for the regression model (F-test) and regression coefficients (t-test) with confidence interval. Another important part of linear regression modeling is checking whether the required assumptions of linearity and i.i.d. (independence and identical distribution) of observations are met. Although the validity of the assumptions can never be entirely certain, there are ways to check for gross violations by analyzing residuals. The prevailing techniques employed in this study to diagnose residuals include box plots, Q-Q plots, scatter plots, partial regression plots, and residuals versus predicted values.

## 4.4 Results of Selected Quantity Models

The scope of analyzed construction activities in this study include: (a) earthwork and landscape, (b) subgrade treatments and base, (c) surface courses and pavement, (d) structures, (e) miscellaneous construction, and (f) lighting, signing, markings and signals.   Work items, as selected by their significant contribution to project cost were analyzed using multivariate regression within several project types up to eight.      Derived item-level quantity models were tested for goodness of fit and statistical validity.

Six quantity-based models are selected as demonstration in the following (the remaining sixty-two models are not shown here).   The resulting predictive models should be utilized appropriately within specified data range.

*Earthwork and Landscape - Item 132 Embankment:*

$$Q_{132} = PL^{(1.480-0.954RER-1.031BWR-1.091INC-0.898UGN-0.788WNF-0.801NNF-0.828WF)}PW^{0.382}(PercentTrucks*AdtPresent)^{0.103}e^{(8.125 - 2.036RER+0.545INC+0.367TrunkSysFlag- 0.610UrbanRural-1.031WNF-1.069UGN-1.180BWR)}$$

(4)

*Subgrade Treatments and Base - Item 247 Flexible Base:*

$$Q_{247} = PL^{(1.061-0.851NNF-0.793BWR-0.610INC-0.234RER+1.667WFPW0.273WNF+ 0.289NNF)}e^{(8.129 +1.858INC)}$$

(5)

*Surface Courses and Pavement - Item 305 Salvaging, and Stockpiling Reclaimable Asphalt Pavement:*

$$Q_{305} = PL^{0.697}(PercentTrucks*AdtPresent)^{0.088}e^{9.981}$$

(6)

*Structures - Item 450 Railing:*

$$Q_{450} = PL^{(0432-0.213RER-0.359NNF)}e^{[7.737+0.297DIV+0.107NOB+ (7.2E-7)EBDA-2.458RER-2.051WNF-1.421BR-1.073NNF-0.546BWR]}$$

(7)

*Miscellaneous Construction - Item 512 Portable Concrete Traffic Barrier:*

$$Q_{512} = PL^{(0.512PW1.297)}e^{(3.856+0.771UrbanRural-1.861WNF-1.359RER)}$$

(8)

*Lighting, Signing, Markings, and Signals - Item 662 Work Zone Pavement Markings:*

$$Q_{662} = PL^{(0.818RER+0.897UGN+0.825WNFPW0.635-0.49RER}e^{7.254+0.952UGN+0.672WNF)}$$

(9)

Where,

Q_{132}:   Quantity of Embankment, (cubic yards)
Q_{247}:   Quantity of Flexible Base, (cubic yards)
Q_{305}:   Quantity of Salvaging and Stockpiling Reclaimable Asphalt Pavement, (square yards)
Q_{450}:   Quantity of Railing, (linear feet)
Q_{512}:   Quantity of Portable Concrete Traffic Barrier, (linear feet)
Q_{662}:   Quantity of Work Zone Pavement Markings, (linear feet)
PL:   Project Length, (miles)
PW:   Project Width, (feet)
BR:   Bridge Replacement
BWR:   Bridge Widening or Rehabilitation
INC:   Interchange
NNF:   New Location Non-Freeway
RER:   Rehabilitation of Existing Road
UGN:   Upgrade to Non-Freeway Standards
WF:          Widen Freeway

WNF:          Widen Non-Freeway
PercentTrucks:   Percent Trucks, %
AdtPresent:   Present Average Daily Traffic, vehicle/day
TrunkSysFlag:   Trunk System Flag (Yes=1; No=0)
NOB:          Number of Bridges
UrbanRural:   Project Location (Urban=1; Rural=0)
DIV:          Divided Roadway (Divided Roadway: DIV=1; Undivided Roadway: DIV=0)
EBDA:          Existing Bridge Deck Area, (square feet)

Validity of the general linear models were inspected through scatter plots of the predicted versus observed values as shown in Figure 3 as an example.   The derived predictive models were calibrated in terms of goodness of fit. From the figures, most of the predicted and historical values cluster approximately around the diagonal line.
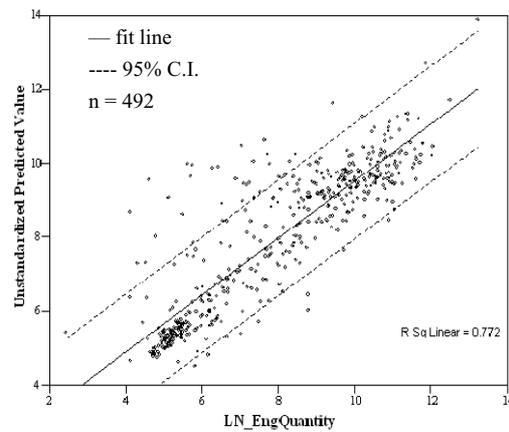


**Figure 3.** Predicted Values of Logarithmic Engineering Quantity vs. Historical Values for Flexible Base

## 4.5 Predictive Model Validation

The construction Industry Institute (CII) has published recommended practices regarding various estimate classes with correspondent accuracy [24].   The Association for the Advancement of Cost Engineering International (AACE) also proposed five classes of cost estimate definitions [25][26].   For estimates at conceptual stages, accuracy of +/-30% to +/-50 was suggested based on the project definition.   The validation results in this study indicate that the selected quantity models performed well for the new highway projects examined with individual prediction error ranging from -26.0% to +27.3% and average prediction error, -10.8% to +10.1%.   These percentage error ranges meet the accuracy suggested by CII and AACE considered preliminary estimates.

## 5. WINDOWS-BASED COMPUTING PROGRAM

An automated estimating system was developed for aforementioned quantity models to be used at the earliest stages of projects.   This system was developed with Microsoft ACCESS, Visual Basic programming and SQL.

### 5.1 Development of Computing Database System

The design concept driving the development of this automated system was to exploit work item historical unit

811

prices rather to continue to guess based on subjective experience and initiate preliminary estimates with quantity-based models. The final result is a semi-detailed preliminary cost estimate with a listing of project information, item quantity, unit price, item cost, and total project estimate.

## 5.2 Computer Program Implementation

As the user proceeds, the system displays a graphical interface as shown in Figure 4. The user can select to create a new estimate, track a previous estimate or exit the system.

Once the user enters the system and selects "Create New Estimate", the program prompts the users to select a project type as shown in Figure 5. The graphical interface offers the users a selection among eight project types. Within the system, a list of major work items is generated through Structured Query Language (SQL) based on project type.

Upon selection of a project, the program requests project basis information as input parameters for preliminary cost estimation. A sample input is illustrated in Figure 6. A message box pops up once the quantity and unit price calculation is completed successfully with the provided information upon clicking the calculating button.

As the user proceeds to the next page, the computer program generates a list of major work items related to the selected project. In this window, the user can preview preliminary cost estimating results by selecting or inputting control section job (CSJ) number as shown in Figure 7.

The user can export this output into other desired applications such as Microsoft Excel and Word or directly as a report as shown in Figure 8 and 9. For other project types, this program generates similar screens and various outputs.
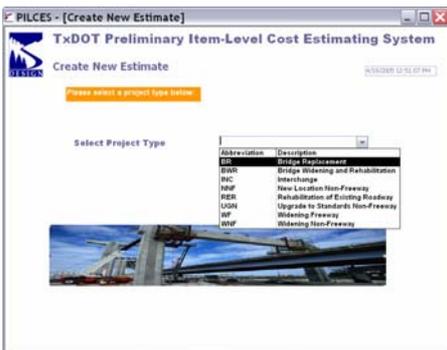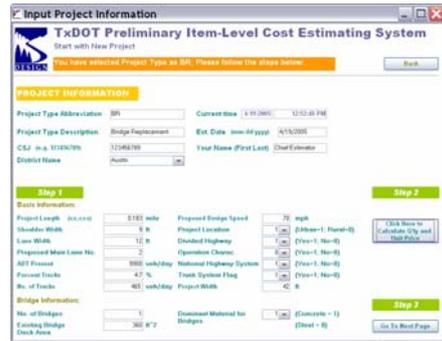


**Figure 6.** Example of Input Data for BR Project



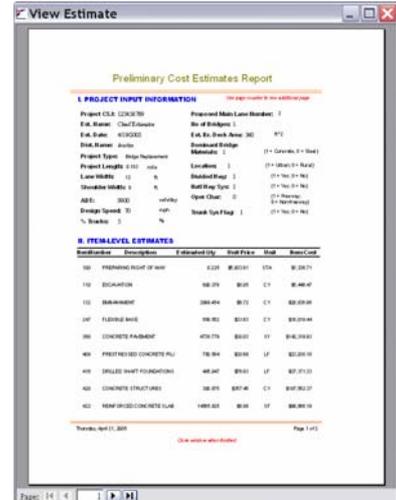**Figure 7.** Windows for Preview, Customize, and Print Preliminary Cost Estimating Results



**Figure 8.** Screen Shot of Estimating Report



**Figure 4.** Program Main Menu



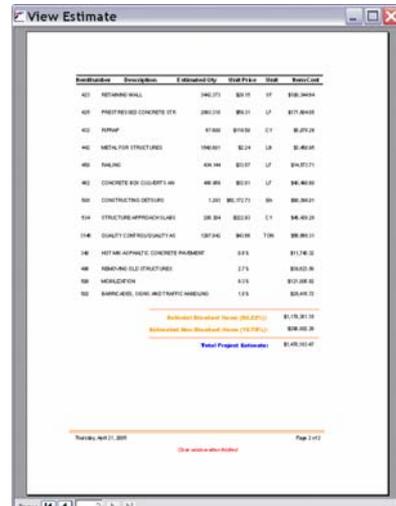**Figure 5.** Project Type Selection Window



**Figure 9.** Screen Shot of Estimating Report (Cont'd)

## 6. CONCLUSIONS

The item-level quantity-base approach employs quantity-based models at project inception to produce preliminary estimates and provides an opportunity for periodic quantity adjustment till design complete. In this study, statistical parametric models were developed to predict quantities of sixty eight identified major work items associated with eight infrastructure project types. A computational database management system was developed for producing item-level cost estimates. This system takes advantage of work item historical unit prices and created quantity-based models to enable subsequent quantity growth tracking. It requires minimal effort and minimal information, and leaves a documentation trail for early estimate input parameters. The stored information can also be used for subsequent cost long range management and control.

## REFERENCES

[1] Bell, L. C., and Bozai, G. A., "Preliminary cost estimating for highway construction projects", *AACE Transactions*, C.6.1-C.6.4, 1987.

[2] Sanders, S. R., Maxwell, R. R., and Glagola, C. R., "Preliminary estimating models for infrastructure projects", *Cost Engineering*, 34(8), pp.7-13, 1992.

[3] Flyvbjerg, B., Holm, M. S., and Buhl, S., "Underestimating costs in public works projects - Error or Lie?", *Journal of the American Planning Association*, 68(3), pp. 279-295, 2002.

[4] Al-Tabtabai, H., Alex, A. P., and Tantash, M., "Preliminary cost estimation of highway construction using neural networks", *Cost Engineering,* 41(3), pp.19-24, 1999.

[5] Hegazy, T., and Ayed, A., "Neural network model for parametric cost estimation of highway projects", *J. Constr. Eng. Manage.,* 124(3), pp. 210-218, 1998.

[6] Herbsman, Z., "Long-range forecasting highway construction costs", *J. Constr. Engrg. and Mgmt.,* 109(4), 423-435, 1983.

[7] Morcous, G., Bakhoum, M. M., Taha, M. A., and El-Said, M., "Preliminary Quantity Estimate of Highway Bridges using Neural Networks", *Proceedings of the Sixth International Conference on the Application of Artificial Intelligence to Civil & Structural Engineering Computing (Scotland),* pp. 51-52, 2001.

[8] Saito, M., Sinha, K. C., and Anderson, V. L., "Statistical models for the estimation of bridge replacement costs", *Transpn. Res. -A,* 25A(6), pp. 339-350, 1991.

[9] Schexnayder, C.J., Weber, S.L., and Fiori, C., *Project Cost Estimating: A Synthesis of Highway Practice, National Cooperative Research Program,* Transportation Research Board, 2003.

[10] Black, J. H., "Application of parametric estimating to cost engineering", *AACE Transaction,* B.10.1-B.10.5., 1984.

[11] Odeck, J., "Cost Overruns in Road Construction-What are their sizes and determinants?", *Transport Policy, (IN PRESS),* 2003.

[12] Wilmot, C. G., and Cheng, G. "Estimating Future Highway Construction Costs", *J. Constr. Eng. Manage.,* 129(3), pp. 272-279, 2003.

[13] Akintoye, A., and Fitzgerald, E., "A Survey of Current Cost Estimating Practices in the UK", *Construction Management and Economics*, Vol. 18, pp. 161-172, 2000.

[14] Akinici, B., and Fischer, M., "Factors affecting Contractors' Risk of Cost Overburden", *Journal of Management in Engineering,* 14(1), pp. 67-76, 1998.

[15] Akintoye, A., "Analysis of Factors Influencing Project Cost Estimating Practice", *Construction Management and Economics,* Vol.18, pp. 77-89, 1998.

[16] Baloi, D., and Price, A. D. F., "Modeling Global Risk Factors Affecting Construction Cost Performance", *International Journal of Project Management*, Vol. 21, pp. 261-269, 2003.

[17] Trost, S. M., and Oberlender, G. D., "Predicting Accuracy of Early Cost Estimates using Factor Analysis and Multivariate Regression", *J. Constr. Eng. Manage*., 129(2), pp. 198-204, 2003.

[18] Albright, S. C., Winston, W., and Zappe, C. J., *Data Analysis & Decision Making with Microsoft Excel,* Thomson Brooks/Cole, 2003.

[19] Neter, J., Kutner, M. H., Nachtsheim, C. J., and Wasserman, W., *Applied Linear Statistical Models,* IRWIN, 1996.

[20] Chengalur-Smith, I. N., Ballou, D. P., and Pazer, H. L., "Modeling the Costs of Bridge Rehabilitation", *Transpn. Res. -A,* 31(4), pp. 281-293, 1997.

[21] Phaobunjong, K., and Popescu, C. M., "Parametric Cost Estimating Model for Buildings", *AACE International Transaction,* EST.13.1-EST.13.11, 2003.

[22] Bobko, P., *Correlation and Regression Applications for Industrial Organizational Psychology and Management,* Sage Publications, Inc., 2001.

[23] *SPSS for WIndows, release 12.0.0.: Statistical Package for Social Sciences,* SPSS Inc., Chicago.

[24] Trost, S. M., *A Quantitative Model for Predicting the Accuracy of Early Cost Estimates for Construction Projects in the Process Industry,* Ph.D. Dissertation, Oklahoma State University, 1998.

[25] *AACE Recommended Practice No.17R-97: Cost Estimate Classification System,* AACE, Inc., 1997.

[26] Lorance, Randal B., and Wendling, Robert V., "Basic Techniques for Analyzing and Presenting Cost Risk Analysis", *AACE International Transactions,* Paper RISK.01, 1999.