

Selection of a Probability Distribution for Modeling Labor Productivity during Overtime**

Woo, Sungkwon*

Abstract

Construction labor productivity, which is the greatest source of variation in overall construction productivity, is the critical factor for determining the project performance in terms of time and cost, especially during scheduled overtime when extra time and cost are invested. The objective of this research is to select an appropriate type of probability distribution function representing the variability of daily labor productivity during overtime. Based on the results of statistical data analysis of labor performance during different weekly work hours, lognormal distribution is selected in order to take advantage of easiness of generating correlated random numbers. The selected lognormal distribution can be used for development of a simulation model in construction scheduling, cost analysis, and other applications areas where representation of the correlations between variables are essential.

키워드 : Labor Productivity, Overtime, Probability Distribution Function, Goodness-of-fit

1. Introduction

Construction labor productivity is not deterministic in nature since it is a fundamental human characteristic that the performance of a person engaged in a certain task will change according to the situation. Also, it is a critical factor for determining the project performance in terms of time and cost, especially during scheduled overtime when extra time and cost are invested. Numerous studies about scheduled overtime and productivity have been conducted in the construction industry because labor productivity is usually the greatest source of variation in overall construction productivity (Halligan et al, 1994).

Generally labor productivity is simplified by

assuming and using a single-fixed value. However, this deterministic approach has the problem of ignoring and not representing the variable's real characteristics. Labor productivity, the variable involving uncertainty, should be modeled as an uncertain variable to reflect the real system. The variability or randomness of actual labor performance can be represented by a probability distribution function. Therefore, stochastic modeling of labor productivity which incorporates the effects of uncertainty is a more reasonable approach to the problem rather than assuming and using single-fixed value estimates.

As a result of the effort of including the uncertain natures involved in a construction project, the stochastic scheduling approach using Monte Carlo simulation has been adopted and studied by many researchers in construction scheduling and cost analysis. However, as Touran and Wiser (Touran and Wiser 1992) pointed out, the problem in most Monte Carlo simulation applications which have been developed in

* 중신회원, 인하대학교 토목공학과, 전임강사, 공학박사

** This work was supported by INHA UNIVERSITY Research Grant. (INHA-30377). Their support is gratefully acknowledged. The author also express thanks to Prof. Calin Popescu in the Univ. of Texas at Austin for his encouragement through this research.

construction scheduling and cost analysis have neglected the correlations between variables. According to them, disregarding correlation between variables in a simulation model when the correlation is significant results in underestimation of variance.

Since a model is the representation of an actual system, not only an uncertain component consisting of the actual system should be modeled as a variable to represent accurately and realistically, but also the correlation between variables, if exists, should be correctly modeled in it. Otherwise, it can be a serious source of error in simulation modeling and analysis. The objective of this research is to select an appropriate type of probability distribution function representing the variability of daily labor productivity during overtime. In addition, because the selected probability distribution function is planned to be used in future simulation modeling and analysis, the aspect of representing the correlation between variables will be considered in the selection process and criteria.

2. Description of Productivity Data

2.1 Productivity Data Overview

Since the selection of the correct probability distribution affects the accuracy of the simulation model, the availability of quality data collected from the actual system is critical. The productivity data used in this research is obtained from the previous construction overtime study reported to the Construction Industry Institute (CII) in 1994. H. Randolph Thomas and Karl A. Raynar from Pennsylvania State University conducted the study under the guidance of the CII Overtime Task Force. The objective of the CII construction overtime study conducted by Thomas and Raynar was the comparison of labor performance during different weekly work hours.

Regarding the source of data, this CII overtime

study is considered reliable since the sources of data, including the types of project and crews, used in this 1994 CII report are clearly identified. Also, other pertinent information such as project environmental conditions and labor situations are known. According to the report (CII 1994), each project had a tranquil labor environment, no inordinate number of changes, no experimental, unique, or poorly managed projects. The early phases and the start-up phase were not included in the study since the stage of construction can affect labor productivity. The overtime schedule was used to maintain schedule, not to attract labor. Therefore, the other potential influences on labor productivity other than the overtime effect were minimized in this study that provided the data to this research.

The projects were conducted in the 1989–1992 time frame. A total of 151 weeks of productivity data were initially collected from four construction projects that were considered average industrial projects with respect to owner involvement, design, and construction management. Table 1 shows the type of the projects and crews studied. The manufacturing and paper mill projects were the existing facilities where old systems and equipment were removed and new ones were installed. The process plant was a spacious, outdoor, grass-root facility, and the refinery involved the rebuilding of an existing facility.

Table 1. Project and crew description (CII 1994)

Project	Crew No.	Craft	No. of Weeks
Process Plant	9181	Mechanical	13
Process Plant	9182	Electrical	14
Manufacturing	9183	Electrical	11
Manufacturing	9184	Mechanical	12
Paper Mill	9185	Electrical	8
Refinery	9186	Mechanical	16
Refinery	9187	Mechanical	15
Refinery	9188	Mechanical	16
Refinery	9189	Mechanical	16
Refinery	9190	Mechanical	15
Refinery	9191	Mechanical	15
Total			151

2.2 Productivity Data Format

The format of the productivity data obtained is the performance factor. The performance factors show the relative comparison between the baseline productivity during normal 40-hour weeks and actual weekly productivity measured. The reason for converting from productivity data to performance factor was for the purpose of comparison of the labor performances between the standard 40-hour week schedule and the 50- and 60-hour overtime work schedules. The productivities actually measured from the crews were converted to equivalent quantities using conversion factors. A more detailed description of the conversion process is explained in the CII report (CII 1994). A value of performance factor greater than one means that actual productivity measured is better than the expected productivity.

$$\text{Performance factor} = \frac{\text{baseline productivity}}{\text{actual productivity}}$$

2.3 Summary of Previous Data Analysis

As mentioned earlier, the objective of the CII construction overtime study conducted by Thomas and Raynar was the comparison of labor performance during different weekly work hours. Table 2 shows the average performance factor of each crew during 40-, 50- and 60-hour

Table 2. Average performance factor for each crew

Crew No.	40 Hours		50 Hours		60 Hours	
	Avg. PF	N	Avg. PF	N	Avg. PF	N
9181	0.896	5	0.831	5	0.750	2
9182	2.245	7	1.795	4	2.094	2
9183	0.999	1	0.792	3	0.792	2
9184	1.000	5	0.783	4	0.579	3
9185	1.031	5	0.595	1	-	-
9186	1.029	5	1.418	3	0.642	4
9187	0.963	5	1.117	5	2.216	1
9188	1.142	4	1.361	4	0.775	2
9189	1.050	4	0.458	4	0.869	3
9190	1.015	6	0.802	5	0.973	1
9191	1.138	3	0.831	3	1.102	1
Avg.	1.192	50	1.001	41	1.102	21

workweeks. Also, Figure 1 shows the box plot comparison of labor performances during difference work schedules. The box plot shows the average performance differences during different work schedules. The effects of frequent work schedule changes between standard and overtime schedules on the labor performance are ignored in the data analysis.

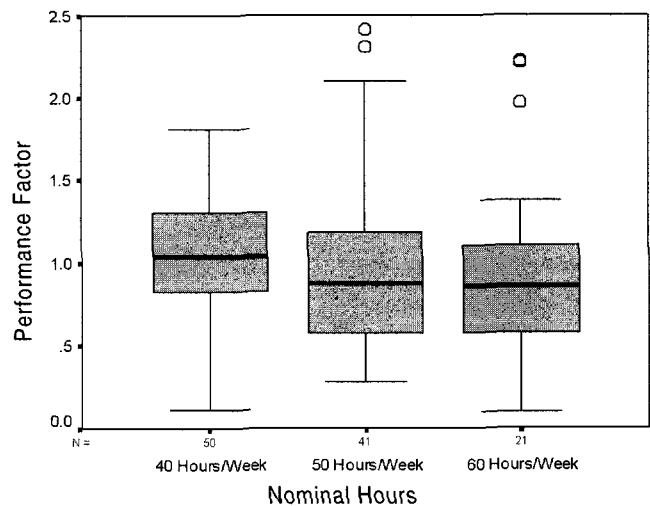


Figure 1. Performance factor comparison using box plot

3. Candidate Probability Distribution Selection

In order to determine a correct probability distribution, the techniques of statistical inference are used to fit a theoretical distribution to the collected data. The hypothesis tests are performed to determine if a particular distribution with certain values for its parameters is a good model for the random variable in a model.

3.1 Limitations and Assumptions

The data from the previous study performed by Thomas were analyzed to select the probability distribution representing random labor performance under overtime schedule. However, several drawbacks in the data made the determination of a probability distribution for modeling difficult. First, although the number of weekly work hours can vary, the overtime schedules included in the data are only 50 and 60

hours per week. Therefore, only those two overtime schedules are analyzed and modeled in this study. Second, frequent variations in work schedules resulted in a smaller sample size of performance data for distribution selection. Also, the frequent work schedule variations made the overtime durations short. Because the overtime schedule did not last long enough, there was not enough data for analyzing the labor performances for certain periods of time after the initiation of the overtime schedule. As a result, only the labor performance at the first week of each overtime schedule was analyzed.

Since the probability distribution function representing the labor performance at each time point could not be determined from the data, the same probability distribution determined from the first week performance data is assumed to be applicable to the following weeks of overtime schedule.

3.2 Data Screening & Selection Processes

Identifying the first week performance data for the each overtime schedule was not a clean-cut process because of the frequent work schedule variations in the projects. Any possible effects of transition between the overtime and standard schedules needed to be minimized by selecting the performance data that meet the following criteria.

For the first week of the 50-hour schedule, only the performance data from the 50-hour schedule following right after 40-hour standard week were selected for data analysis. The performance data from the weeks following after 60-hour overtime week were discarded to eliminate the unnecessary transition effects. For the first weeks of the 60-hour schedule, the performance data from the 60-hour week following right after the 40-hour standard week were selected. The data from the 60-hour weeks, which follow after the 50-hour overtime schedule, were also selected to increase the sample size for the 60-hour schedule data analysis.

After the selection process is completed, the sample size of $n_{50\text{-hour}} = 19$ and $n_{60\text{-hour}} = 17$ for 50- and 60-hour overtime schedule, respectively, was found. The selected performance data from the first weeks of 50- and 60-hour overtime schedules were analyzed to determine the appropriate probability distribution representing the distribution of random labor performance. The histograms of the selected performance data from the first weeks of each schedule are shown in Figure 2.

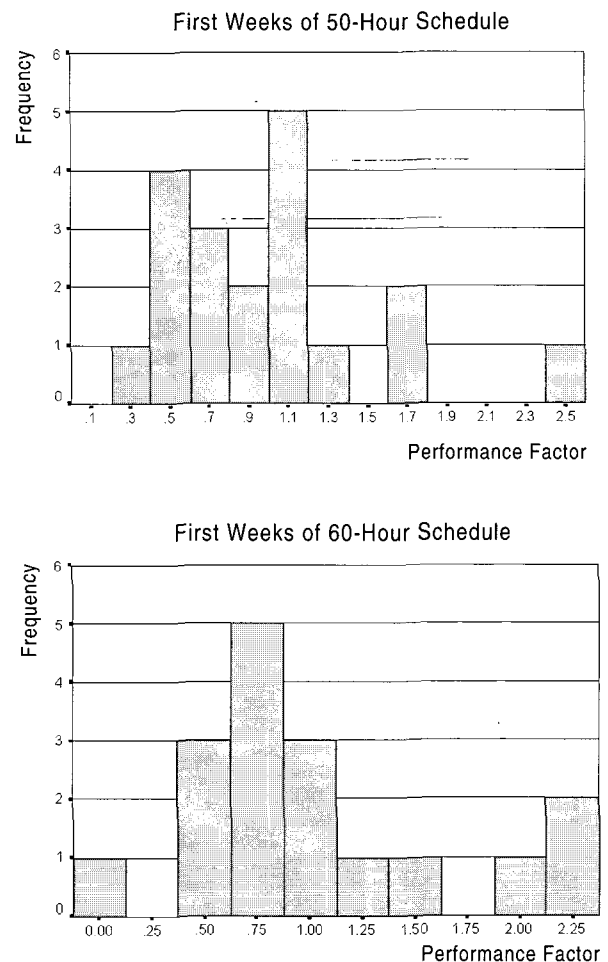


Figure 2. Histogram of performance data

3.3 Hypothesizing Distributions

Before making any inference from the estimates of any statistics from the selected data, prior knowledge about the variable in the model also can provide a helpful idea of how the selected probability distribution should look like without

requiring any data at all. The candidate distributions are pre-selected based on the following properties of the required distribution.

First of all, the probability distribution should be a continuous distribution. The continuous distribution should have a range of values equal to or greater than zero since the lower bound value of labor performance factor in the model should be equal to 0. Furthermore, the long-term average value of the distribution would be close to but smaller than 1, due to the expected performance drop under overtime work. Therefore, the probability distribution to be selected can not be a symmetric distribution because the upper limit of the distribution can be any certain positive value reasonably greater than 2. The probability distribution to be selected should be a distribution skewed to the right. As a result, the proposed input probability distribution should be continuous, non-symmetric, unimodal, and positively skewed distribution.

The approximate shape of the proposed density function is shown in Figure 3. The shape of distribution shown in Figure 3 conforms the shape the histograms in Figure 2. When the proposed probability distribution function is skewed to the right and has a density with a shape similar to that in the Figure 3, the possible candidates among the various continuous distributions are Gamma, Exponential, Weibull, and Lognormal (Law and Kelton 1991). Therefore, those four probability distributions are pre-selected as candidate distributions.

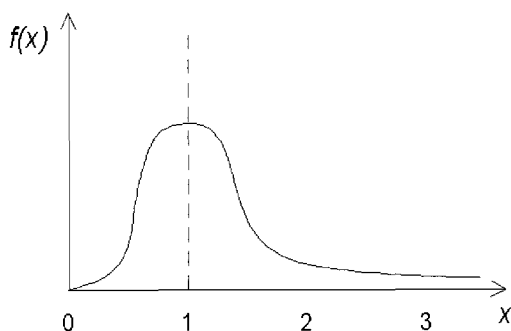


Figure 3. Expected shape of a distribution

For the performance data sets selected for each overtime schedule, the summary statistics were estimated from the samples (Table 3). The summary or descriptive statistics sometimes provide very useful information about the true distribution because they show the important properties of the data.

Table 3. Summary statistics of the first week data

	50 Hours	60 Hours
Sample size, n	19	17
Min.	0.353	0.100
Max.	2.414	2.222
Median	0.984	0.861
Sample Mean	1.007	1.04824
Sample Variance	0.263	0.3735
Coefficient of Variation	0.509	0.5830
Skewness	1.242	0.7977

According to Law & Kelton (1991), the coefficient of variation (cv), a measure of variability like the variance, is particularly useful in selecting a distribution among the exponential, gamma, Weibull, and lognormal distributions. For the exponential distribution, $cv = 1$ regardless of the scale parameter β . Since the estimated coefficient of variation for 50- and 60-hour data are not close to 1, the exponential distribution is ruled out from the candidate distributions.

For the gamma and Weibull distributions, cv is greater than, equal to, or greater than 1 when the shape parameter is less than, equal to, or greater than 1, respectively (Law and Kelton 1991). These two distributions will have a shape similar to the density function in Figure 3 when the shape parameter $\alpha > 1$, which implies that $cv < 1$ (Law and Kelton 1991). Therefore, the gamma and Weibull distributions are the candidate for the input probability distribution. Also, the lognormal distribution always has a density with a shape similar to that in Figure 3, but its cv can be any positive real number. If there were more observations of large performance factors for the 50- and 60-hour schedules, it is very likely to

have the coefficient of variation greater than 1. Therefore, considering the fact that it is very likely to have few, if any, observations from the right tail of the true underlying distribution when the sample size n is not very large, the lognormal distribution is also a candidate distribution. Therefore, the gamma, Weibull, and lognormal distributions are pre-selected as the candidate distribution for modeling the random labor performance under 50- and 60-hour overtime schedules.

4. Determination of Best Fitted Distribution

There are various ways of determining if a hypothesized distribution fits the data. In this study, visual inspection or comparison of the productivity data and three candidate distributions and goodness-of-fit tests are performed to determine the best fitted distribution. In order to make a comparison the specific values of the parameters should be calculated from the data set.

4.1 Estimation of Parameters

Because three distributions, lognormal, gamma and Weibull, are hypothesized, the specific values of their parameters were estimated from the data.

The maximum-likelihood estimators(MLE) were used to calculate the parameters. The Table 4 shows the estimated parameters of the lognormal, gamma, and Weibull distributions.

Table 4. Estimated parameters of distributions

	Parameter	50 hour	60 hour
Lognormal	Mean	-0.107123	-0.152324
	Std. deviation	0.47879	0.715593
Gamma	shape $\hat{\alpha}$	4.526194	2.662018
	scale $\hat{\beta}$	0.222576	0.393775
Weibull	shape $\hat{\alpha}$	2.148714	1.826482
	scale $\hat{\beta}$	1.142017	1.178666

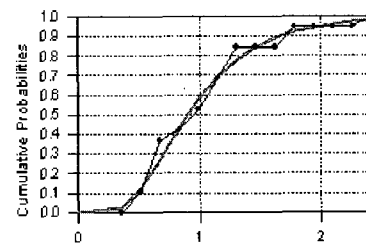
4.2 Visual Comparison of Productivity Data with Distributions

For the visual inspection of how closely the data agree with the candidate distributions, the

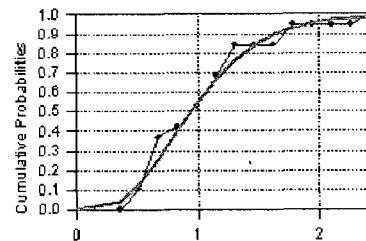
cumulative distributions of the productivity data and the three hypothesized distributions with estimated parameters are drawn and compared. Figure 4 and 5 show the comparison of 50- and 60-hour workweek data with the three hypothesized distributions.

After the visual comparison of the data with the distributions, little disparity is detected between the data and those three distributions. Also, the amount of deviation does not differ greatly among the distributions. Therefore, it was concluded that the three pre-selected distributions agree fairly well with the labor performance data. And, there is no big enough difference to draw a distinction between those distributions from the visual inspection.

Comparison of 50-Hour-Input Distribution and Lognormal(-0.107, 0.479)



Comparison of 50-Hour-Input Distribution and Gamma(4.526,0.223)



Comparison of 50-Hour-Input Distribution and Weibull(2.149, 1.142)

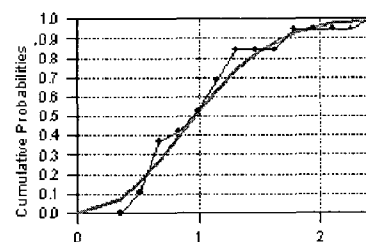
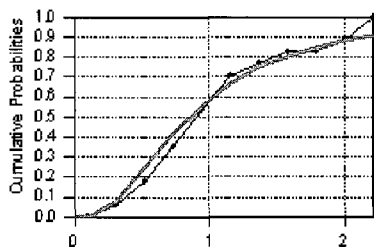
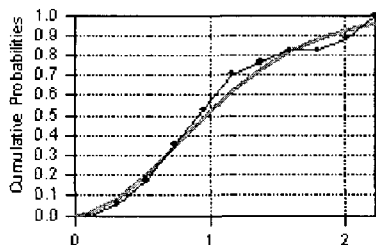


Figure 4. Comparison of 50-hour productivity data with hypothesized distributions

Comparison of 60-Hour-Input Distribution and Lognormal(-0.152, 0.716)



Comparison of 60-Hour-Input Distribution and Weibull(2.149, 1.142)



Comparison of 60-Hour-Input Distribution and Gammal(2.662, 0.394)

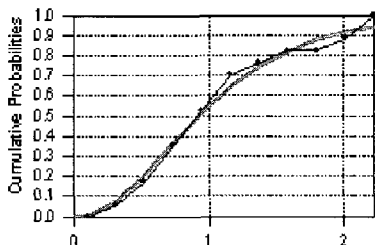


Figure 5. Comparison of 60-hour productivity data with hypothesized distributions

4.3 Goodness-of-fit Tests

A goodness-of-fit test is a statistical hypothesis test for assessing formally if the observations x_1, x_2, \dots, x_n are an independent sample from a particular distribution with distribution function F . The null hypothesis tested in goodness-of fit tests is

$$H_0: \text{The } x_i\text{'s are IID random variables with distribution function } F$$

The common goodness-of-fit tests are chi-square test, Kolmogorov-Smirnov (K-S) test, and Anderson-Darling (A-D) test. While the chi-square may be the most widely used technique, it requires data grouped into categories. On the other hand, the K-S test and A-D test work with individual values rather than groups of values. The chi-square test is applicable to any hypothesized

distribution, while the other tests have limited applicability.

4.3.1 Chi-Square Tests

In order to perform chi-square tests for both data sets, the input distribution of performance data from the 1st week of 50- and 60-hour schedules are divided into 13 and 10 intervals, respectively, and compared with the fitted lognormal, gamma and Weibull distributions. As a result, the test statistics are obtained and compared with the critical values. The test statistics obtained and critical values are provided in Table 5.

The value of test statistic χ^2 will be small if the distribution is a good fit. Therefore, the null hypothesis H_0 is rejected if χ^2 is too large. The null hypothesis H_0 is rejected only if the test statistic χ^2 is greater than the critical value $\chi^2_{k-1, 1-\alpha}$. Since the values of test statistics for those three distributions are all smaller than the critical values, none of the distributions are rejected.

Table 5. Chi-square test statistics & critical values

		χ^2	$1-\alpha$			
			0.900	0.950	0.975	0.990
50 hour	lognormal	15.706154	18.549348	21.02607	23.336664	26.216967
	gamma	16.862615	18.549348	21.02607	23.336664	26.216967
	Weibull	17.539217	18.549348	21.02607	23.336664	26.216967
60 hour	lognormal	8.200708	18.549348	21.02607	23.336664	26.216967
	gamma	7.091397	18.549348	21.02607	23.336664	26.216967
	Weibull	6.913416	18.549348	21.02607	23.336664	26.216967

4.3.2 Kolmogorov-Smirnov Tests

Kolmogorov-Smirnov (K-S) test for goodness-of fit compares an empirical distribution function from data with the hypothesized distribution function. The detailed description of the K-S test procedures and its characteristics can be found in Law & Kelton (1991).

Because the K-S test measures how closely the empirical distribution and the fitted distribution match, larger the distance, it is more likely to reject

the hypothesis. The adjusted statistics and critical values $C'_{1-\alpha}$ for normal, Weibull, and gamma distributions are given in the Table 6. The values of adjusted K-S statistics of data from the 50- and 60-hour schedules with three hypothesized distributions were computed using the formula, and the results obtained are shown on Table 7.

Since the values of the test statistics for 50- and 60-hour data with three hypothesized distributions are smaller than the corresponding critical values, the null hypothesis H_0 for all three distributions are not rejected at level 0.15. Therefore, the hypotheses, lognormal, Weibull and gamma distributions represents the labor performance under the overtime schedules, are not rejected.

Table 6. Modified critical values for adjusted statistic.

Distribution	Adjusted test statistic	1- α				
		0.850	0.900	0.950	0.975	0.990
$N(\bar{x}(n), S^2(n))$	$(\sqrt{n}-0.01+\frac{0.85}{\sqrt{n}})Dn$	0.775	0.819	0.895	0.955	1.035
Weibull	$\sqrt{n}Dn$	-	0.779	0.843	0.907	0.973
gamma	$(\sqrt{n}+0.12+\frac{0.11}{\sqrt{n}})Dn$	1.138	1.224	1.358	1.480	1.628

(Law and Kelton 1991)

Table 7. K-S Test statistics calculated from data

Distribution	K-S Statistics Calculated	
	50-hour	60-hour
lognormal	0.143201	0.143124
gamma	0.149274	0.098313
Weibull	0.148414	0.118843

4.3.3 Anderson-Darling Test

The Anderson-Darling (A-D) test is designed to detect discrepancies in the tails by giving larger weights for $\hat{F}(x)$ close to 1 and $\hat{F}(x)$ close to 0, whereas the K-S test gives the same weight to the difference between $F_n(x)$ and $\hat{F}(x)$ for all values of x . In A-D test, when normal (or lognormal) or Weibull distribution is assumed and the parameters are estimated from the data, the adjusted statistic and the modified critical values are used for testing hypothesis. The A-D test procedure is described in detail by Law & Kelton

(1991).

The adjusted A-D statistics calculated for 50- and 60-hour schedules, shown in Table 9, were compared with the modified critical values in Table 8. Since the adjusted A-D statistics calculated are smaller than the corresponding critical values, the H_0 is not rejected at level 0.10. Therefore, the lognormal, Weibull and gamma distributions reasonably represent the labor performance under the overtime schedules.

Table 8. Modified critical values for adjusted statistic.

Distribution	Adjusted test statistic	1- α				
		0.850	0.900	0.950	0.975	0.990
$N(\bar{x}(n), S^2(n))$	$(1+\frac{4}{n}-\frac{25}{n^2})A_n^2$	0.775	0.819	0.895	0.955	1.035
Weibull	$(1+\frac{0.2}{\sqrt{n}})A_n^2$		0.779	0.843	0.907	0.973
gamma	A_n^2 for $n \geq 5$	1.138	1.224	1.358	1.480	1.628

(Law and Kelton 1991)

Table 9. A-D Test statistics calculated from Data

Distribution	A-D Statistics Calculated	
	50-hour	60-hour
lognormal	0.310558	0.517055
gamma	0.341333	0.302958
Weibull	0.443976	0.325521

4.4 Correlation Consideration in Probability Distribution Selection

Even though there is no study or research that has shown the existence of performance correlation between labor trades, it is very likely that certain amount of correlation exists between labor trades in certain situations. Since all types of labor trades are working under the same project condition, it is very likely for the performance or productivity of certain types of labor trades to have correlations, more likely positive. Especially when overtime schedule is initiated, the project will be in accelerated pace and the project conditions will have much more impact on the labor performances. It is a reasonable assumption that activities using same type of resource will show similar pattern in productivity change. Therefore, no matter what distribution is chosen

for representing the random labor performance, the simulation model should be capable of generating correlated random numbers for more realistic representation of the actual system.

In that aspect, the lognormal distribution is the better choice than gamma and Weibull because it has the advantage of generating correlated random numbers. Since multivariate normal distribution can be transformed into multivariate lognormal, one can take advantage of the property that correlated random numbers can be relatively easily generated using multivariate normal distribution. Generation of correlated labor productivity variables from a multivariate lognormal requires only the marginal distribution of each variable, given that covariances between variables are known.

Lognormal distribution has been used to represent random cost variables by Touran (1993) who pointed out the problem of neglecting correlations in Monte Carlo simulation in construction scheduling and cost control. Therefore, in case there is not much difference in goodness-of-fit among lognormal, Weibull and gamma, the lognormal should be selected in order to take advantage of easiness of generating correlated random numbers.

5. Conclusions

The objective of this research is to select an appropriate type of probability distribution function representing the variability of daily labor productivity during overtime. Based on the productivity data collected from the previous construction overtime study by Construction Industry Institute (CII) and the prior knowledge about the properties of variable, four probability distributions, gamma, exponential, Weibull, and lognormal are pre-selected as the candidate distributions. Then, in order to find the best fitted distribution, visual inspection and goodness-of-fit tests are performed to determine the best fitted distribution.

Since the lognormal, gamma and Weibull distributions were not rejected from the goodness-of-fit tests, it can be said that the data are an independent sample from a lognormal, gamma or Weibull distribution. It means that all three distributions provide a reasonably good model for the random labor performance under the overtime schedules, even considering the small sample sizes of the data. However, because the probability distribution function selected is planned to be used in future simulation modeling and analysis, the lognormal distribution is selected due to the advantages in correlated random number generation. The selected lognormal distribution can be used for simulation applications in construction scheduling, cost analysis, and other applications areas where representation of the correlations between variables are essential.

References

1. CII. Effects of Scheduled Overtime on Labor Productivity: A Quantitative Analysis. Source Document 98, Construction Industry Institute, The University of Texas at Austin, 1994
2. Halligan, David W. et al. "Action-response model and loss of productivity in construction." *Journal of Construction Engineering and Management*, 120(1), p.47-63. 1994
3. Law, Averill M. and Kelton, W. David. Simulation Modeling & Analysis. McGraw-Hill, Inc. New York, 1991
4. Touran, Ali. "Probabilistic cost estimating with subjective correlations." *Journal of Construction Engineering and Management*, 119(1), p.58-71. 1993
5. Touran, Ali and Wiser, Edward P. "Monte carlo technique with correlated random variables." *Journal of Construction Engineering and Management*, 118(2), p.258-272. 1992