# Calculating Project Completion in Polynomial Processing Time

Professor Luis Copertari
Computer Engineering Program, Autonomous University of Zacatecas, México

## Abstract

Technology-based organizations and knowledge organizations rely on large activity networks to manage Research & Development (R&D) projects. Avoiding optimistic completion times due to the characteristic Program Evaluation and Review Technique (PERT) assumptions is a problem that can grow exponentially in complexity with the number of activities. A recursive technique that solves the problem in a polynomial number of steps has been developed, assuming that all duration times follow beta distributions. It is important to notice that the only two 100% valid approaches to calculate the project completion time are simulation and the stochastic sum for each and every path in the network. Nevertheless, both require finding the shape parameters, and that is precisely the main contribution of this paper: a system of equations to calculate the shape parameters of each activity and the overall project.

## Keywords

Project management, completion time, scheduling, complex networks

## 1. Introduction

One of the most important theoretical problems in project management is to obtain the distribution of the total completion time in project networks. The main approaches used are the Program Evaluation and Review Technique (PERT) and the Critical Path Method (CPM). CPM was developed in the 1950s by researchers at Du Pont and Sperry Rand (Meredith & Mantel, 2008), and PERT was developed in the late 1950s by consultants working on the development of the Polaris missile (Chachra et al., 1979; Evans & Minieka, 1992; Phillips & Diaz, 1981; and Wiest & Levy, 1997). PERT assumes three-point estimates for probabilistic activity duration times in order to approximate project completion and the relative probability at each milestone, using the normal distribution (Fisher et al., 1985). CPM focuses on the criticality of each activity and the time–cost

tradeoff in deterministic activity networks (Smolin, 1981). For practical and managerial purposes, what matters is the criticality of each activity within a PERT network, which can be assessed using a sound approach to calculate the completion time (Wolf, 1985). Critical activities are activities that, if delayed, would delay the entire project. A sequence of critical activities throughout the network is called a critical path. The critical path is the longest path in the network and it is possible to have more than one critical path at the same time. But unlike CPM, in stochastic activity networks the duration time of individual activities varies, so activities are critical for some combinations of duration times but may not be critical for other combinations. Therefore, activities have a given probability of being critical (i.e., being part of the longest path). The probability of each activity being on the critical path is defined as its criticality. The focus of this paper is to describe an analytical method for calculating the theoretical distribution of the project completion time.

Malcolm et al. (1959) rely on the central limit theorem to postulate that the completion time can be portrayed using a normal distribution as a function of the cumulative mean and variance of all the activities within the longest path. Unfortunately, this results in unreliable (typically less than actual) completion times (Dodin, 1985). Martin (1965) calculates the completion time by approximating task duration density functions using polynomials. Although accurate, Martin's method requires considerable calculation and is not easily adaptable for software implementation. Kleindorfer (1971) and Devroye (1979), among others, obtain lower bounds to the expected duration of the total project, based on node criticality, whereas Dodin & Elmaghraby (1985) approximate such criticality indices. The latter is not entirely correct from a theoretical point of view, but the advantage of bounding the mean completion time from below is that closed-form solutions can be obtained. Also, Dodin (1984) tries to determine the k most critical paths as opposed to calculating completion times for each path. Monte Carlo simulation (Touran & Wiser, 1992; Van Slyke, 1963) is valid from a theoretical point of view, but it requires considerable calculation, which makes it impractical in the case of complex networks.

It will be shown that the PERT assumption of normally distributed project completion time typically leads project managers to optimistic planning, based on less than actual project

completion estimates, due to a failure to consider the absolute bounds to project completion (Donaldson, 1965; Grubbs, 1962; MacCrimmon & Ryavec, 1964; Sasieni, 1986). These bounds arise from the fact that the actual project completion time is the maximum sum of the duration of each and every path, which in turn is the result of adding the actual duration of its activities. It is common practice in PERT to estimate activity durations by using beta distributions (Fisher et al., 1985). Project completion cannot be an unbounded random variable because the sum of bounded (e.g. beta distributed) activity duration times yields bounded path (and project) completion times. The normal distribution cannot give upper and lower bounds on project completion times. PERT uses the same completion time algorithm as CPM, but applied to the mean. The problem is that this algorithm tends to yield inaccurate results.

The PERT textbook formula to calculate expected (mean) activity duration times, which are assumed to follow beta density functions, considers three parameters (minimum, most likely, and maximum), when in fact the beta distribution has four parameters: two range parameters and two shape parameters (MacCrimmon & Ryavec, 1964). The PERT formula used to calculate the mean as a function of the minimum (a), most likely (mode or m), and maximum (b) activity duration time estimates, (a+4m+b)/6, ignores how the biases to the right or left (related to the mode) affect the shape of the beta distribution.

## 2. PERT/CPM networks

Network models can be used to schedule complex projects that consist of many activities. CPM can be used when the duration of each activity is known with certainty to determine the duration of the entire project. It can also be used to determine how long activities in the project can be delayed without delaying the entire project. If the duration of the activities is not known with certainty, PERT can be used to estimate the probability of the project being completed at any given deadline.

A project is a combination of interrelated tasks or activities that must be executed in some pre-specified sequence. Projects are described using probabilistic or deterministic activity networks, which are directed acyclic graphs. Let $\aleph$ denote the adjacency matrix of a

probabilistic PERT/CPM network composed of nodes (vertices) $\mathbf{N} = \{1,2,\ldots,n\}$ and directed arcs $\mathbf{A} = \{(i,j) \mid i=1\ldots n-1, j=2\ldots n\}$ where n is the total number of nodes. Let m be the total number of activities so that the set of directed arcs, $\mathbf{A}$, can also be denoted as $\mathbf{A} = \{k \mid k=1\ldots m\}$. The duration of arc $(i,j)$ is a random variable $t_{ij}$ with a known probability density function $f_{ij}(t)$ over the closed interval $[a_{ij}, b_{ij}]$ where $\mu_{ij}$ denotes the mean (expected) duration of activity k in arc $(i,j)$ and $\sigma_{ij}^2$ its variance. (Activity on Arc notation or AOA is implicit, where i indicates node of origin and j node of destination.) The completion time at sink node j, $T_j$, is the time at which all activities coming into j have been completed. The completion time at source node i, $T_i$, is the earliest time at which any activity k in arc $(i,j)$ located between nodes i and j is allowed to start. (Notice that $T_i=T_j$ when i and j refer to the same node; i.e., i=j.) $T_i$ (or $T_j$) is a random variable with unknown probability density function $f_i(T)$ (or $f_j(T)$). The purpose of this discussion is to describe how to accurately calculate the relevant probability density functions.
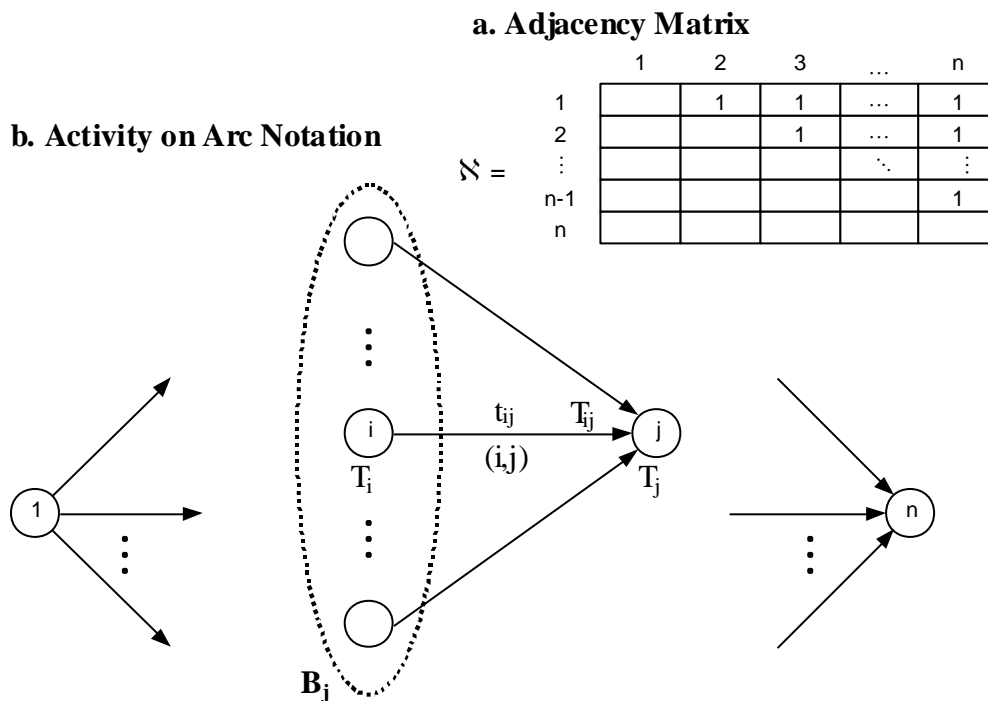
**a. Adjacency Matrix**



**b. Activity on Arc Notation**

**Figure 1.** PERT/CPM network

The adjacency matrix contains all precedence relationships. Figure 1a illustrates the adjacency matrix of a fully connected activity network. Row i indicates the node of origin while column j is the destination node for activity k at coordinates (i,j), explicitly specifying the position within the network for each activity. The nodes in directed acyclic networks are numbered in such way that an arc always leads from a smaller numbered node to a larger one. Let $\mathbf{B_j}$ (j=2…n) denote the set of predecessor nodes connecting to node j. Let i be one such node (i∈$\mathbf{B_j}$). Figure 1b illustrates the notation. The completion time at node i is given by the random variable $T_i$, where $T_i$' is one random occurrence of $T_i$. If i is the only node in $\mathbf{B_j}$ (i.e., |$\mathbf{B_j}$|=1), then $T_j$' is given by the sum $T_i$'+$t_{ij}$', where $t_{ij}$' is one random occurrence of $t_{ij}$. In general, when the number of nodes coming into j is more than one (i.e., |$\mathbf{B_j}$|>1), the resulting completion time at j is the maximum completion time of all incoming arcs as indicated in equation (1). (Notice that T is used to indicate completion time, whereas t indicates duration time; T includes the duration time of all preceding activities.)

$$T_j'= \underset{i\in\mathbf{B_j}}{\text{Max}}\{T_i'+t_{ij}'\} , \text{ j=2,…,n} \tag{1}$$

It is sometimes useful to denote activities using a single number k because it facilitates notation involving sets in which activity k is said to belong to path p for all p=1…w, where w is the total number of paths (a path is a specific sequence of activities beginning at node 1 and ending at node n). Conversely, denoting activities using their nodes of origin and destination facilitates writing equations for forward pass computations such as equation (1). Equation (1) is a stochastic sum across the network. The plus sign is used to denote the addition of two stochastic variables. $T_{ij}=T_i+t_{ij}$ indicates the completion time that would occur at node j if activity (i,j) happens to be critical across the network (longest duration time in a particular combination of random duration times), whereas $t_{ij}$ is the duration time of activity k in arc (i,j). The completion time at node j is by definition the set of all maximum duration time combinations of the set $\mathbf{B_j}$ of all nodes i preceding node j. Random duration times are described using probability density functions. In particular, PERT assumes that each activity duration time is given by a beta density function (Meredith & Mantel, 2008). Range and shape parameters are required to specify beta

density functions. The range parameters are a and b (minimum and maximum), and the shape parameters are $\alpha$ and $\beta$. Let f(x) be a beta density function as defined in equation (2). We chose a beta density function because it is commonly used in project management models to denote activity duration times (Meredith & Mantel, 2008), given the fact that it allows the portrayal of biases to the right or to the left.

$$f(x) = \frac{1}{\int_0^x t^{\alpha-1}(1-t)^{\beta-1} dt} \frac{(x-a)^{\alpha-1}(b-x)^{\beta-1}}{(b-a)^{\alpha+\beta-1}} , \, a < x < b \tag{2}$$

The standardized beta density function varies between 0 and 1 (range parameters given) so that only the shape parameters are required. Range parameters are intuitively easy to understand and it is reasonable to expect decision-makers to use them and to provide their estimates. But shape parameters are difficult to grasp. So instead of specifying the shape parameters, decision-makers are asked to give the range and the most likely duration time (mode). From these the mean and variance are usually approximated in practice by (a+4m+b)/6 and $(((b-a)/6)^2)$, respectively. The problem is how to add the random variables of beta density functions across the activity network accurately in order to obtain a probability density function describing project completion time, when the shape parameters are also required.

## 3. PERT completion time

The coefficients of the function $f_{ij}(T)$ are essentially obtained by adding mean duration times and minimum and maximum times of the preceding activities and adjusting for the variance. PERT involves the addition of mean duration times. But activity networks are a combination of entangled paths and not a single path. The concept of stochastic sum applies only to specific paths. The problem is how to calculate the completion time at nodes with several incoming activities. It is tempting to extend equation (1) and apply it to mean completion times. In fact, that is exactly what PERT is all about. PERT assumes the mean completion time at node j is the maximum of the mean completion time of all the arcs preceding node j (Donaldson, 1965). Let $\lambda_j$ denote the mean completion time at node j for

all j=1…n, where $\lambda_1=0$. Then, the mean completion time at node j in PERT is given according to equation (3)[1].

$$\lambda_j = \underset{i \in \mathbf{B_j}}{\text{Max}}\{\lambda_i + \mu_{ij}\} \quad \forall \; j=2…n \tag{3}$$

Adding the mean completion time of each node i and the corresponding activity in arc (i,j) is statistically acceptable because both means are in sequence and the result would be the mean of activity (i,j) if the activity is critical. But assuming that the mean at node j is the maximum of these is not accurate. This is because we do not know a priori which activity is critical. It may very well be that several activities are critical in different degrees (with different probabilities) for different duration time combinations. Besides, equation (1) applies to random variables and not to expected values. To illustrate, consider two discretely distributed activities arranged in parallel. Assume that the first activity can have a duration time of 5 or 7 with equal probability ($t_1=\{5,7\}$), whereas the duration time of the second activity can be 6 or 8 with equal probability ($t_2=\{6,8\}$). PERT would calculate the mean duration time of the first activity, (5+7)/2=6, and the mean duration time of the second activity, (6+8)/2=7, and assume the mean completion time of both activities to be $\mu_{PERT}=\text{Max}(\mu_1,\mu_2)=\text{Max}(6,7)=7$. But in fact, duration times are random variables, which means that there are four possible combinations for $\text{Max}(t_1,t_2)$ indicating project completion: Max(5,6)=6, Max(5,8)=8, Max(7,6)=7, and Max(7,8)=8. The mean completion time is in fact the average of these: $\mu_{THEORETICAL}=(6+8+7+8)/4=7.25$. In this case PERT underestimates the completion time because it does not consider the probability distributions, which describe the random behavior of activity duration times.

So how can the expected project completion time be estimated accurately? One way is to consider all path combinations, calculate the duration time of each path by adding the duration time of its activities, obtaining the joint probability density function of these and calculating its mean. Unfortunately, the number of paths grows exponentially as the number of nodes increases. In other words, the computational effort increases as the complexity of

---

[1] In PERT, $\mu = (a+4m+b)/6$ and $\sigma^2 = ((b-a)/6)^2$.

the network increases. Simulation can and often is used to approximate the theoretical completion time by calculating a large enough number of given duration times for each activity in the network.

## 4. Polynomial completion time

Let the set p denote a path consisting of a sequence of $n_p$ activities and let activity k be one of the activities in path p ($k \in p$). Also, let $f_k(a_k, b_k, \alpha_k, \beta_k)$ or simply $f_k(t)$ be a beta density function with range parameters $a_k$ and $b_k$ ($a_k < t < b_k$) and shape parameters $\alpha_k$ and $\beta_k$ describing the duration of activity k, where $F_k(t)$ is the corresponding cumulative distribution. By definition (Hastings & Peacock, 1975), the mean of the beta distributed duration time for activity k is given according to equation (4).

$$\mu_k = \frac{a_k(\beta_k + 1) + b_k(\alpha_k + 1)}{\alpha_k + \beta_k + 2} \qquad (4)$$

The beta density function can be simplified to the standard beta distribution by assuming that the range parameters are 0 and 1. Let t' = $(t-a_k)/(b_k-a_k)$ be the standardized duration time (0<t'<1), where $a_k$'=0 and $b_k$'=1 denote the range parameters of the standard beta density function. Let $\mu_k$' be the mean of the standardized beta distribution. Clearly, $\mu_k$' is the relative distance between the original mean and the range parameters as indicated in equation (5).

$$\mu_k' = \frac{\mu_k - a_k}{b_k - a_k} \qquad (5)$$

If we assume that the relationship between the shape parameters of the beta density function ($\alpha_k$ and $\beta_k$) and the shape parameters of the standardized beta density function ($\alpha_k$' and $\beta_k$') are given by equations (6) and (7), we can obtain the mean of the standardized beta density function by substituting into equation (4) as indicated in equation (8), since $a_k$' = 0 and $b_k$' = 1 as mentioned above.

$$\alpha_k' = \alpha_k + 1 \tag{6}$$

$$\beta_k' = \beta_k + 1 \tag{7}$$

$$\mu_k' = \frac{a_k'(\beta_k') + b_k'(\alpha_k')}{\alpha_k' + \beta_k'} = \frac{\alpha_k'}{\alpha_k' + \beta_k'} \tag{8}$$

Rearranging the terms of equation (8) yields equation (9).

$$\left(\frac{\beta_k'}{1}\right)\left(\frac{1}{\alpha_k'}\right) = \left(\frac{1}{\mu_k'}\right)\left(\frac{1-\mu_k'}{1}\right) \tag{9}$$

Figure 2 shows all three types of standard beta distributions for different combinations of shape parameters. U-shaped beta distributions occur when the sum of the shape parameters is less than 2. J-shaped beta distributions occur when the sum of the shape parameters is greater than or equal to 2 and less than 3. Bell-shaped beta distributions occur when the sum of the shape parameters is greater than or equal to 4.

It is common practice to portray activity duration times using bell-shaped beta distributions (Malcolm et al., 1959). Only one interpretation of equation (9) provides the simplest system of two equations portraying $\alpha'$ and $\beta'$ as a function of $\mu'$ that guarantees a bell-shaped beta density function (as opposed to U-shaped or J-shaped) for any given value of $\mu'$. The four alternative systems of equations (cases) for $\beta'$ and $1/\alpha'$ consistent with equation (9) are:

1.  $\beta' = \left(\frac{1-\mu'}{\mu'}\right)^{\frac{x}{y}}$ and $\frac{1}{\alpha'} = \left(\frac{1-\mu'}{\mu'}\right)^{\frac{y-x}{y}}$ so that $x/y + (y-x)/y = 1$.

2.  $\beta' = 1 - \mu'$ and $\frac{1}{\alpha'} = \frac{1}{\mu'}$.

3.  $\beta' = 1$ and $\frac{1}{\alpha'} = \frac{1-\mu'}{\mu'}$ if $\mu' > 0.5$ (right bias) or $\alpha' = 1$ and $\beta' = \frac{1-\mu'}{\mu'}$ if $\mu' < 0.5$ (left bias).

4.  $\beta' = \frac{1}{\mu'}$ and $\frac{1}{\alpha'} = \frac{1-\mu'}{1}$.

There are an infinite number of combinations of x and y for the first case, leading to a system of equations for α' and β' consistent with equation (9). According to the scientific precept known as Ockham's razor, attributed to the English philosopher William of Ockham (1990), all things being equal, the simplest explanation tends to be the truth. Clearly, a case in which there are infinite possibilities is not the simplest case, so the first alternative should be discarded. The second alternative is also discarded. It leads to U-shaped beta distributions because $0 < \mu' < 1$ so that $\alpha' = \mu'$ and $\beta' = 1 - \mu'$ must be between 0 and 1 as well (see Figure 2a). The third alternative is also rejected because it corresponds to J-shaped beta distributions since, in that case, either α' or β' equals 1 (see Figure 2b). The last alternative is the only one that ensures a bell-shaped beta distribution in which the sum of α' and β' is greater than or equal to 4 (see Figure 2c). The minimum value for α'+β' occurs when $\mu' = 0.5$ so that $\alpha' = \beta' = 1/0.5 = 2$ and $\alpha' + \beta' = 4$. All other values for μ' lead to values of α'+β' greater than 4. Consequently, the shape parameters of the standardized beta density function describing the standardized duration time of activity k are a function of the standardized mean as indicated in case four. Equations (10) and (11) portray the results from case four.

$$\alpha_k' = \frac{1}{1 - \mu_k'} \tag{10}$$

$$\beta_k' = \frac{1}{\mu_k'} \tag{11}$$

Substituting $\mu_k'$ from equation (5) into equations (10) and (11) and equations (6) and (7) yields equations (12) and (13).

$$\alpha_k = \frac{\mu_k - a_k}{b_k - \mu_k} \tag{12}$$

$$\beta_k = \frac{b_k - \mu_k}{\mu_k - a_k} \tag{13}$$

**a. U-Shaped**
$0 < \alpha+\beta < 2$

Left    Symmetrical    Right
Bias                          Bias

$\alpha<\beta$         $\alpha=\beta$         $\alpha>\beta$

**c. Bell-Shaped**
$\alpha+\beta \geq 4$

Left  Symmetrical  Right
Bias                      Bias

$\alpha>1$      $\alpha=\beta$      $\alpha>2$
$\beta>2$                              $\beta>1$

**b. J-Shaped**
$2 \leq \alpha+\beta < 3$

Left    Symmetrical    Right        Left    Symmetrical    Right
Bias                          Bias        Bias                          Bias

$\alpha>1$      $\alpha>1$      $\alpha>1$      $\alpha=1$  $\alpha<1$  $\overline{\alpha=1}$  $\beta<1$  $\alpha=2$
$\beta>1$      $\beta>1$      $\beta>1$      $\beta=2$        $\beta=1$              $\beta=1$
$\alpha<\beta$      $\alpha=\beta$      $\alpha>\beta$

$1<\beta<2$              $1<\alpha<2$
$\alpha=1$                  $\beta=1$

**Figure 2.** Shapes of the beta distribution

The standardized variance is by definition (Hastings & Peacock, 1975) given according to equation (14).

$$\sigma_k'^2 = \frac{\alpha_k'\beta_k'}{\left(\alpha_k'+\beta_k'\right)^2\left(\alpha_k'+\beta_k'+1\right)} \tag{14}$$

Also, the variance is $(b-a)^2$ times the standardized variance as indicated in equation (15).

$$\sigma_k^2 = \left(b_k - a_k\right)^2 \sigma_k'^2 = \left(b - a\right)^2 \frac{\alpha_k'\beta_k}{\left(\alpha_k'+\beta_k'\right)^2\left(\alpha_k'+\beta_k'+1\right)} =$$
$$= \left(b - a\right)^2 \frac{\left(\alpha_k + 1\right)\left(\beta_k + 1\right)}{\left(\alpha_k + \beta_k + 2\right)^2\left(\alpha_k + \beta_k + 3\right)} \tag{15}$$

Clearly, the shape parameters ($\alpha$ and $\beta$) depend on the relationship between the mode (m) and the shape parameters lower and upper bounds (a and b). According to Hastings and Peacock (1975), the mean ($\mu$) in terms of a, b, m, $\alpha$ and $\beta$ are given according to equation (16).

$$\mu_k = \frac{a_k + (\alpha_k + \beta_k)m_k + b_k}{\alpha_k + \beta_k + 2} \tag{16}$$

Also, from equations (12) and (13), we can calculate $\alpha_k' = \alpha_k + 1$ and $\beta_k' = \beta_k + 1$ as shown in equations (17) and (18).

$$\alpha_k' = \alpha_k + 1 = \frac{b_k - a_k}{b_k - \mu_k} \tag{17}$$

$$\beta_k' = \beta_k + 1 = \frac{b_k - a_k}{\mu_k - a_k} \tag{18}$$

For simplicity, for any given activity k, we will not use the subscript k. Thus, substituting equations (12) and (13) and equations (17) and (18) into equation (16) yields equation (19).

$$\mu = \frac{a + \left(\dfrac{\mu - a}{b - \mu} + \dfrac{b - \mu}{\mu - a}\right)m + b}{\dfrac{b - a}{b - \mu} + \dfrac{b - a}{\mu - a}} \tag{19}$$

Simplifying equation (19) yields equation (20).

$$\mu = \frac{(b - \mu)(\mu - a)(a + b) + m\left((\mu - a)^2 + (b - \mu)^2\right)}{(b - a)^2} \tag{20}$$

Solving for $\mu$ from equation (20) yields equation (21).

$$\mu = \frac{4ab - 2mb - 2am}{2(a + b - 2m)} + $$
$$+ \frac{2\sqrt{2b^2a^2 - ab^2m - a^2bm - m^2b^2 + 2abm^2 - a^2m^2 - a^3b + a^3m - ab^3 + mb^3}}{2(a + b - 2m)} \tag{21}$$

Notice that if m is the average of a and b, then equation (21) is not valid because of division by zero. However, in this special case, the mean (μ) equals the mode (m).

In order to avoid calculating optimistic mean duration times, instead of dividing by 2(a+b-2m), the division is made by 1.85(a+b-2m), which results in equation (22). In order to find the 1.85, a considerable number of examples were tried and simulation was used. It was found that dividing by 2(a+b-2m) led to less than actual mean completion times, and so a smaller number was used. This number was found to be 1.85 instead of 2.

$$\mu = \frac{4ab - 2mb - 2am}{1.85(a + b - 2m)} +$$
$$+ \frac{2\sqrt{2b^2a^2 - ab^2m - a^2bm - m^2b^2 + 2abm^2 - a^2m^2 - a^3b + a^3m - ab^3 + mb^3}}{1.85(a + b - 2m)} \quad (22)$$

So now it is possible to calculate the mean duration time of each activity k using equation (22) based on the minimum (a), most likely (m) and maximum (b) time estimates for each activity. Once such calculation has been done, it is possible to apply the PERT procedure to calculate the mean duration time of the entire project.

## 5. Minimum and maximum completion times

Consider a set of paths arranged in parallel with beta distributed duration times, each with a different minimum and maximum. Equation (1) indicates that the resulting completion time is the maximum of all these randomly distributed path duration times. What is the minimum completion time possible? The minimum completion time must be the maximum of the minimum completion time of each and every path. The same reasoning applies to the analysis of maximum path duration times: the maximum completion time must be the maximum of the maximum completion time of each and every path. Therefore, equation (1) can be applied to the range (minimum and maximum) in which the randomly distributed completion time is allowed to vary. Let $A_j$ and $B_j$ be the minimum and maximum completion time at node j for all j=2…n where $A_1$=0 and $B_1$=0 (by definition, the first node

does not indicate completion time). Then, the minimum and maximum completion times at each node j are given according to equations (23) and (24).

$$A_j = \underset{i \in \mathbf{B_j}}{Max}(A_i + a_{ij}) \quad \forall \; j=2\ldots n \tag{23}$$

$$B_j = \underset{i \in \mathbf{B_j}}{Max}(B_i + b_{ij}) \quad \forall \; j=2\ldots n \tag{24}$$

These are absolute bounds to the completion time at node j because no $T_j$ can be less than $A_j$ nor greater than $B_j$ at node j, as shown in equation (25).

$$A_j \le T_j \le B_j \tag{25}$$

## 6. Variance for the completion time

Now we have the mean, the minimum and the maximum completion times for the project. It is time to calculate the variance. Substituting equations (17) and (18) for α' = α+1 and β' = β+1 into equation (15) yields equation (26).

$$\sigma^2 = \frac{(b-a)^2(\alpha+1)(\beta+1)}{(\alpha+\beta+2)^2((\alpha+1)+(\beta+1)+1)} =$$
$$= \frac{(b-a)^2\left(\dfrac{b-a}{b-\mu}\right)\left(\dfrac{b-a}{\mu-a}\right)}{\left(\dfrac{b-a}{b-\mu}+\dfrac{b-a}{\mu-a}\right)^2\left(\dfrac{b-a}{b-\mu}+\dfrac{b-a}{\mu-a}+1\right)} \tag{26}$$

Solving equation (26) for $\sigma^2$ yields equation (27).

$$\sigma^2 = \frac{(b-\mu)^2(\mu-a)^2}{(b-a)^2+(b-\mu)(\mu-a)} \tag{27}$$

Equation (27) can be used to calculate the variance based on the mean, the minimum and the maximum completion times. This, intuitively makes sense, since the further apart the absolute bounds are, the larger the value of the variance will be.

## 7. Shape parameters for the completion time

The most difficult part of this approach to calculating project completion is the estimation of the shape parameters based on a, b, $\mu$ and $\sigma^2$. This estimation requires working with the fundamental equations of the beta distribution and making some assumptions. First of all, from equation (5) and solving for (b-a) gives equation (28).

$$b - a = \frac{\mu - a}{\mu'} \tag{28}$$

From equation (11) $\beta'$ is substituted into equation (28), yielding equation (29).

$$b - a = \beta'(\mu - a) \tag{29}$$

From equation (4), since $\alpha' = \alpha+1$ as shown in equation (6) and $\beta' = \beta+1$ as shown in equation (7), equation (30) results.

$$\mu = \frac{a\beta' + b\alpha'}{\alpha' + \beta'} \tag{30}$$

Substituting equation (30) into equation (29) results in equation (31).

$$b - a = \beta'\left(\frac{a\beta' + b\alpha'}{\alpha' + \beta'} - a\right) \tag{31}$$

Simplifying equation (31) results in equation (32a). This can also be expressed as shown in equation (32b).

$$\frac{\alpha'\beta'}{\alpha'+\beta'} = 1 \tag{32a}$$

$$\alpha'\beta' = \alpha'+\beta' \tag{32b}$$

Assume the existence of a shape parameter s. This shape parameter reflects how large $\alpha$ and $\beta$ are. As $\alpha$ and $\beta$ become larger the beta distribution becomes taller and thinner. In such cases, the variance will be smaller. A similar but opposite reasoning applies when the beta distribution is shorter and wider. Thus, s is defined as shown in equation (33).

$$s = \frac{1}{\sigma^2} \tag{33}$$

The shape parameter is multiplied by $\alpha$ and $\beta$ wherever they appear, so instead of just having $\alpha$ or $\beta$, there would be $s\alpha$ and $s\beta$. The above is shown in equations (34) and (35).

$$\alpha \leftarrow s\alpha \tag{34}$$

$$\beta \leftarrow s\beta \tag{35}$$

Substituting equations (34) and (35) into equation (15) yields equation (36).

$$\sigma_k{}^2 = (b-a)^2 \frac{s\alpha's\beta'}{(s\alpha'+s\beta')^2(s\alpha'+s\beta'+1)} \tag{36}$$

Simplifying equation (36) yields equation (37a).

$$\sigma_k{}^2 = (b-a)^2 \frac{\alpha'\beta'}{(\alpha'+\beta')^2(s\alpha'+s\beta'+1)} \tag{37}$$

Substituting equation (32b) into equation (37) yields equations (38a). Rearranging equation (38a) yields equation (38b).

$$\sigma_k{}^2 = \frac{(b-a)^2}{\alpha'\beta'(s\alpha'+s\beta'+1)} \tag{38a}$$

$$\alpha'\beta'(\alpha'+\beta')s\sigma_k{}^2 + \alpha'\beta'\sigma^2 = (b-a)^2 \tag{38b}$$

Substituting equation (32b) into equation (38b) yields equation (39).

$$(\alpha'\beta')^2 s\sigma_k{}^2 = (b-a)^2 - \alpha'\beta'\sigma^2 \tag{39}$$

Substituting equation (33) into equation (39) yields equation (40).

$$(\alpha'\beta')^2 = (b-a)^2 - \alpha'\beta'\sigma^2 \tag{40}$$

By conveniently substituting equations (10) and (11) as well as equation (5) into equation (40), it is possible to obtain equations (41) and (42) for the shape parameters, $\alpha'$ and $\beta'$.

$$\alpha' = \left(\frac{\mu-a}{b-a}\right)\left((\mu-a)(b-\mu)-\sigma^2\right) \tag{41}$$

$$\beta' = \left(\frac{b-\mu}{b-a}\right)\left((b-\mu)(\mu-a)-\sigma^2\right) \tag{42}$$

And the original shape parameters, $\alpha$ and $\beta$, can be obtained from equations (6) and (7) as shown in equations (43) and (44).

$$\alpha = \alpha' - 1 \tag{43}$$

$$\beta = \beta' - 1 \tag{44}$$

Notice that the addition of the shape parameter does not changed the standardized mean. The standardized mean is given according to equation (4) and equations (6) and (7) as shown in equation (45), since $a = 0$ and $b = 1$ for the standardized mean.

$$\mu' = \frac{a\beta' + b\alpha'}{\alpha' + \beta'} = \frac{\alpha'}{\alpha' + \beta'} \tag{45}$$

If the shape parameter, s, is added, equation (46) results, which is not different than equation (45); and the shape parameter is simplified.

$$\mu' = \frac{s\alpha'}{s\alpha' + s\beta'} = \frac{s\alpha'}{s(\alpha' + \beta')'} = \frac{\alpha'}{\alpha' + \beta'} \tag{46}$$

## 8. Project complexity

The minimum number of activities, $\underline{m}$, for a serial network of n nodes is n-1. The maximum number of activities, $\overline{m}$, for a fully connected network of n nodes is n(n-1)/2 (Lawler, 1976). Clearly, the minimum number of paths $\underline{w}$ for an all-serial network, is 1. What is the maximum number of paths? The maximum number of paths occurs in fully connected networks of n nodes and n(n-1)/2 activities. All paths must include the first and last nodes. The total number of combinations (subsets) of all intermediate nodes gives the total number of paths in fully connected networks. In a fully connected network of n nodes, there are n-2 intermediate nodes (n nodes minus the first and last nodes), so that the maximum number of paths, $\overline{w}$, is the total number of subsets of a set of size n-2. According to the binomial theorem (Parzen, 1960), this is given by equation (45).

$$\overline{w} = 1 + \binom{n-2}{1} + \binom{n-2}{2} + ... + \binom{n-2}{n-2} = 2^{n-2} \tag{45}$$

Let $\eta$ be defined as the density coefficient for an AOA network with n nodes and m activities, indicating how nearly all-serial or all-parallel (fully connected) the network is, as shown in equation (46). The density coefficient is the proportional distance between the actual and the minimum compared to the maximum number of activities. (Notice, for m=$\underline{m}$, $\eta$=0, and for m=$\overline{m}$, $\eta$=1.)

$$\eta = \frac{m - \underline{m}}{\overline{m} - \underline{m}} = \frac{2(m - n + 1)}{(n-1)(n-2)} \,, n \geq 3 \tag{46}$$
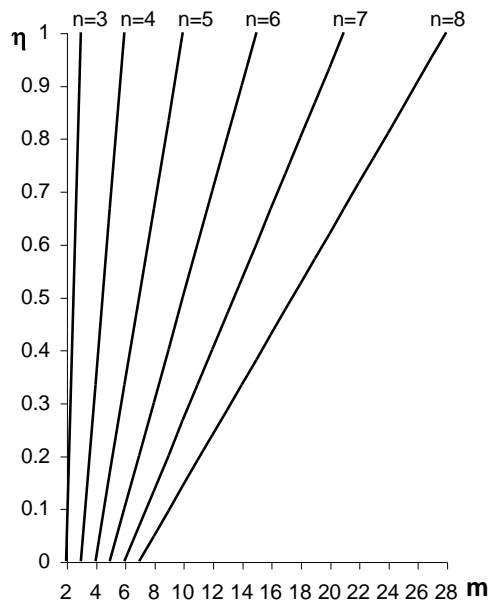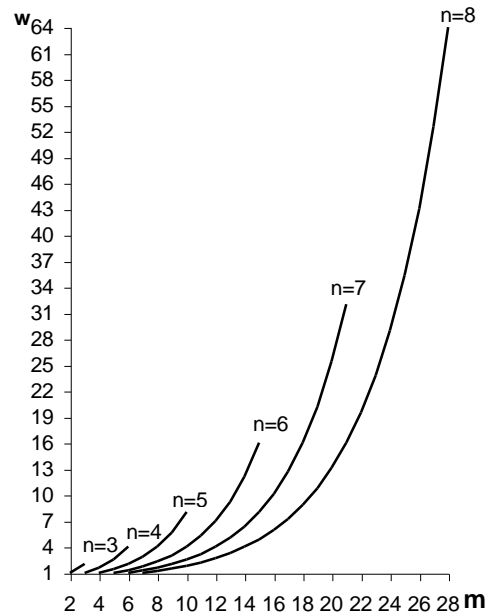
The total number of paths, w, is then defined as a function of both n and $\eta$ according to equation (47), where $\lfloor x \rfloor$ is the floor function (truncation of the fractional component) of x.

$$w \approx \left\lfloor 2^{\eta(n-2)} \right\rfloor \tag{47}$$

When $\eta=0$ for minimally connected (all-serial) networks, $w = \underline{w} = 2^0 = 1$ and when $\eta=1$ for fully connected (all-parallel) networks, $w = \overline{w} = 2^{n-2}$. Equation (47) portrays exponential growth mediated by $\eta$. The network density coefficient from equation (46) measures network complexity and equation (47) provides the number of paths for a given combination of activities and nodes.

Figure 3a plots network density ($\eta$ on the vertical axis) as a function of the number of activities (m on the horizontal axis) and network size (n as different lines in the graph) according to equation (46). The relationship between the number of activities (m) and network density ($\eta$) is linear. It is clear from Figure 3a and equation (46) that the rate (given by the slope) at which network density ($\eta$) grows will decrease as network size (n) increases.

Figure 3b plots the number of network paths (w on the vertical axis) as a function of the number of activities and network density (m which determines $\eta$ along the horizontal axis) for different network sizes (n indicated as different lines) according to equation (47). Although the number of paths (w) increases exponentially as network size (n) and network density ($\eta$) increase, the rate at which such exponential growth occurs ($\eta$) decreases as network size (n) increases.

**a. Network Density versus Activities**        **b. Paths versus Activities**



**Figure 3.** Network density and complexity

The total number of paths is important because more paths increase the computational effort required to estimate the expected project completion time and approximate the probability distribution of such completion time. Figure 3 shows that, although the computational effort tends to increase exponentially as network size increases (exponential time problem), the range for the number of activities ($\overline{m} - \underline{m}$) increases for larger networks. Consequently, the probability of having complex networks and significant exponential growth in the number of paths decreases as the number of nodes increases. The total number of paths is further constrained by the variability within the network of activity duration times and the resulting range and variability in pathway duration times (absolute bounds to completion time). When the mean and variance of the network is calculated according to equations (22) to (42), we can use CPM and assume that the duration time of each activity is the mean duration time ($\mu$). In this way, the computational effort is significantly reduced and becomes polynomial because there is no need to calculate the completion time according to a truly stochastic approach to project completion[2].

---

[2] Keeping in mind that the number of paths increases exponentially with the number of activities in the network.

By using the system of equations proposed to calculate the project completion time in a manner that is similar to PERT, except that now the variance and the shape parameters are considered, the number of calculations required for the forward pass computations is the same as the number of activities in the network. This is because only arcs (activities) require computational effort in order to incorporate such duration time into the node's completion time. Solving for m from equation (46) yields equation (48).

$$w = m = \frac{(n-1)(\eta(n-2)+2)}{2} = \frac{\eta n^2}{2} + n\left(1 - \frac{3\eta}{2}\right) + \eta - 1 \qquad (48)$$

Equation (48) is a quadratic polynomial in n. Hence, calculating project completion is a problem with polynomial complexity (requiring a polynomial number of steps as a function of problem size to solve the problem), not exponential complexity (requiring an exponential number of steps as a function of problem size to solve the problem).

## 9. Example

Consider a relatively small (although illustrative) example. A larger example would be more desirable, but the computational effort required to do the simulation would be much more time-consuming. The only way to approximate the true distribution of the project completion time is by simulating a very large number of randomly beta distributed duration times and calculate in each case the duration time of the project. In this case, a total of 20,000 different sets were used and, after calculating the mean and variance of the duration time for each set, 10 runs (each of 20,000 different sets) were obtained.

The activity names, precedence, and estimates for the minimum, most likely, and maximum duration times for each activity is shown in Table 1.

**Table 1.**  Data for the illustrative example

| Activity | Precedence | a | M | b |
|----------|------------|---|---|---|
| A | - | 2 | 6 | 7 |
| B | - | 4 | 8 | 10 |
| C | B | 3 | 7 | 8 |
| D | A, C | 5 | 10 | 11 |
| E | A, C | 6 | 10 | 12 |
| F | B | 6 | 10 | 11 |
| G | D | 3 | 7 | 8 |
| H | E, F | 3 | 7 | 8 |

The information from Table 1 leads for the Activity on Arc (AOA) representation shown in Figure 4.



**Figure 4.** AOA representation for the illustrative example

However, there is a fundamental problem. In order to do the simulation, the shape parameters of each activity are required in order to calculate random duration times for each activity in the AOA network so that a random occurrence of the project completion time can be obtained. Which shape parameters should be used? Should the shape parameters in equations (10) and (11) or the shape parameters in equations (41) and (42) be used?

Since the standardized beta density function is implicitly used, this involves $\alpha'$ and $\beta'$, and the standardized beta density function is given in the simulation by equation (49).

$$f(t) = f(\alpha', \beta', a, b) = a + (b-a)f(\alpha', \beta', 0, 1) \qquad (49)$$

Using the set of shape parameters from equations (10) and (11) for the simulation results in the probabilistic density function shown in Figure 5 and the cumulative probability density function shown in Figure 6.
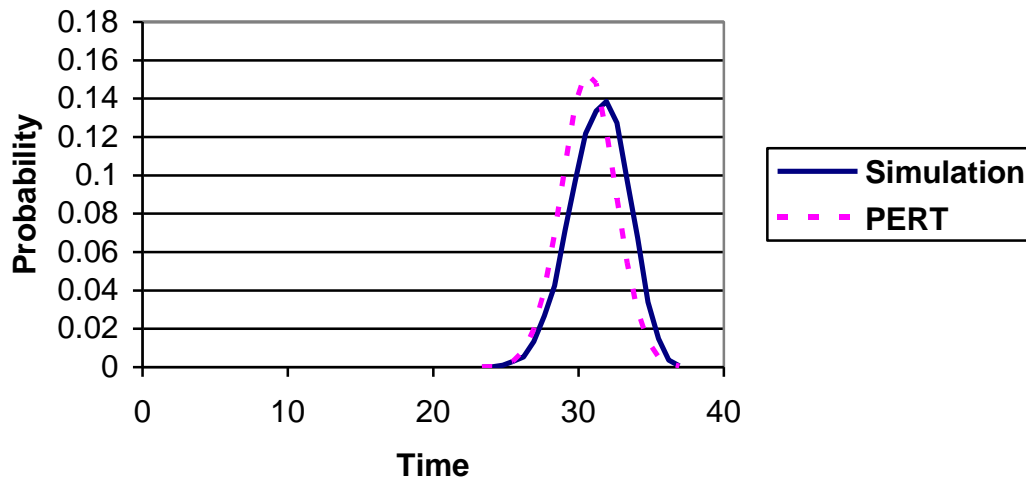


**Figure 5.** Probability density functions of the illustrative example for the shape parameters in equations (10) and (11)
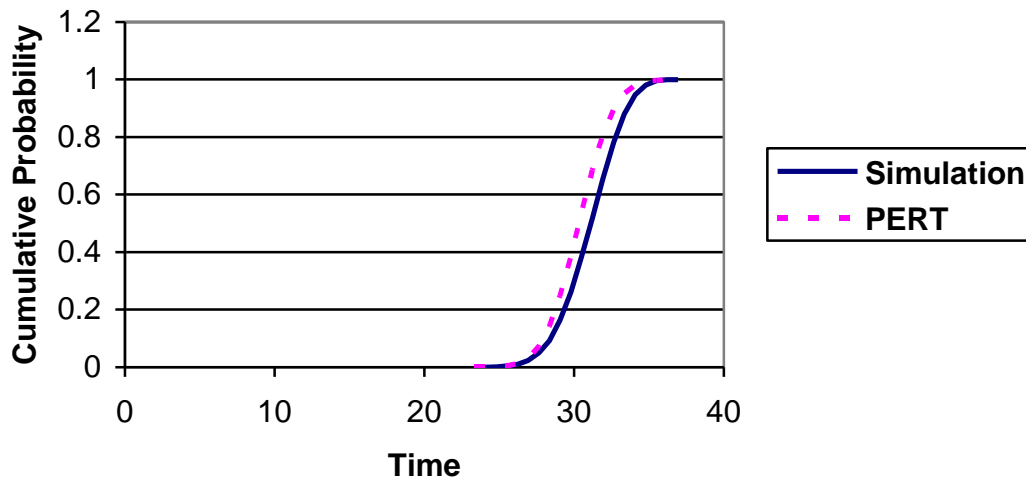
**Figure 6.** Cumulative probability density functions of the illustrative example for the shape parameters in equations (10) and (11)

Also, the mean in the simulation for the 10 sets (each of 20,000 runs) is 31.04 and the variance is 3.93. For PERT, the mean is 30.33 and the variance is 3.39.

The other (more complete) approach of using the shape parameters from equations (41) and (42), the mean completion time from equation (22) and the variance from equation (27) for each activity results in Figure 7. The cumulative density function is shown in Figure 8.
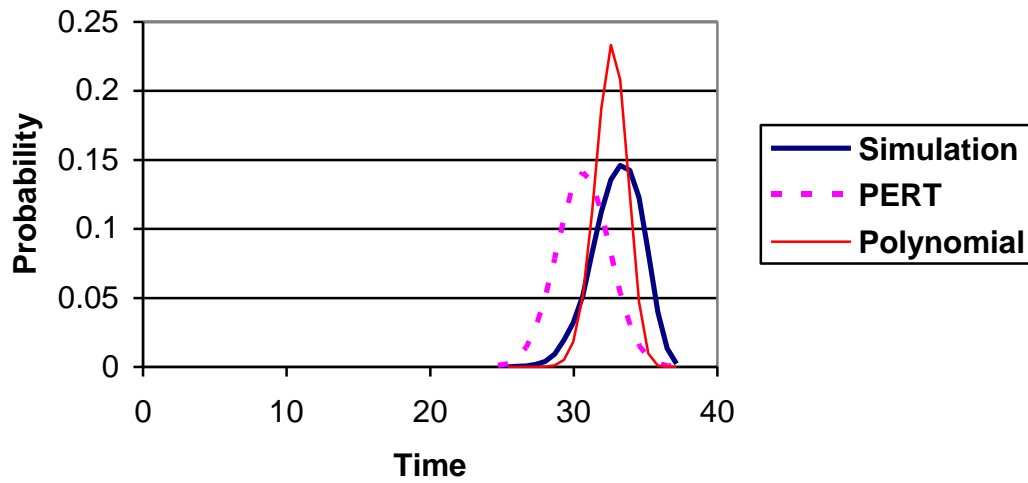
**Figure 7.** Probability density functions of the illustrative example for the shape parameters in equations (41) and (42)

The mean completion time for the simulation from the 10 sets of 20,000 runs each is 32.70 and the variance is 2.97. The mean completion time for PERT again is 30.33 and the variance is 3.39. The mean completion time for the polynomial approach is 32.10 and the variance is 3.95.
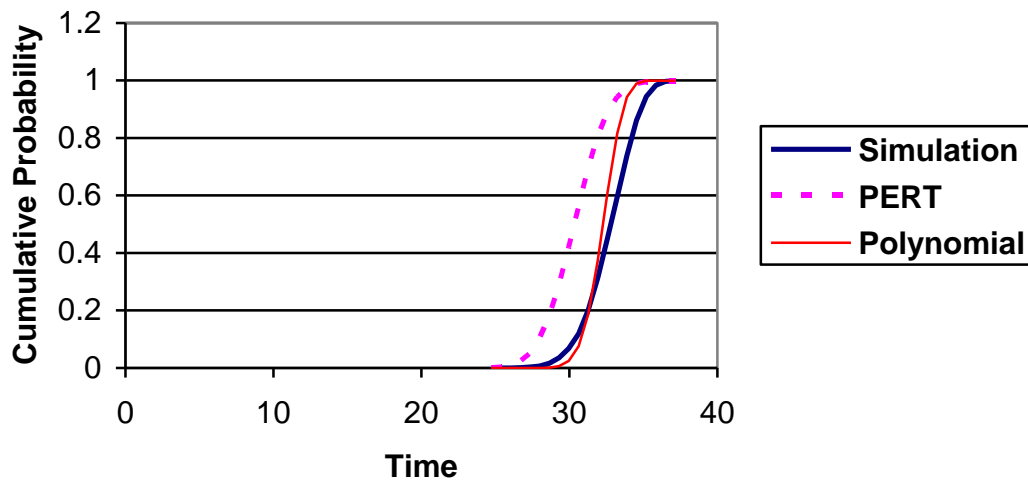


**Figure 8.** Cumulative probability density functions of the illustrative example using the shape parameters in equations (41) and (42)

## 10. Discussion and conclusions

Figure 5 shows a relatively good fit of the PERT approach when compared with the simulation, although the simulation shows a larger variance than the PERT approach (the variance of the simulation is 3.93 and the variance of PERT is 3.39). Nevertheless, PERT still slightly underestimates the completion time (seen in Figure 5), since the PERT probability density function is slightly to the left of the simulated density function. The situation can be further illustrated by plotting the cumulative probability functions shown in Figure 6. Figure 6 clearly shows that the PERT approach overestimates the completion time at the beginning but as times goes on it tends to fit the simulation. However, even in this case, in order to be able to run the simulation of 20,000 sets of duration time occurrences for each activity, the shape parameters of the beta density function are required. These shape parameters can be calculated according to equations (10) and (11) or according to equations (41) and (42), but they are always required. Which approach is better?

If equations (10) and (11) are used instead of equations (41) and (42), the variance is not taken into account. This is incorrect, because clearly a larger variance should result in smaller shape parameters and vice versa. This means that the situation illustrated in Figure 7 and Figure 8 should prevail. In this case, the so-called polynomial approach leads to a better solution because it tends not to underestimate the project completion time, whereas PERT clearly underestimates the completion time.

But why are different results obtained when a different set of equations for the standardized shape parameters are used? This is because equations (10) and (11) implicitly assume a value for $\alpha'+\beta'$ closer to 4, which is exactly the PERT assumption. Nevertheless, the use of equations (41) and (42) for the shape parameters is better from a theoretical point of view, because they take into account the variance (even though the variance of the polynomial approach was higher than the one obtained in the simulation).

The recommendation to the practitioner is to avoid calculating the mean and the variance according to the PERT textbook formula, $\mu=(a+4m+b)/6$ where m is the mode, a is the minimum activity duration and b is the maximum activity duration and $\sigma^2=((b-a)/6)^2$. These

are then simply used to calculate the maximum mean and variance at each node. In fact, the textbook formula assumes a fixed value for the sum of the shape parameters ($\alpha+\beta=4$) to calculate the mean, and it calculates the variance as an approximation to that assumption. Furthermore, PERT does not consider the variance when determining which path is the longest, since the variance of the project completion time is assumed to be the same as the variance of the path with the longest sum of mean duration times. All these assumptions typically lead to optimistic planning due to less than actual project completion times. Instead, the practitioner can use equation (22) to calculate the mean duration time of each activity, equation (27) to calculate the variance of each activity, and equations (41) and (42) to calculate the shape parameters of each activity. Then, the PERT principles are used to calculate the mean and variance of the project completion time, except for the use of the values previously mentioned for each activity.

**Bibliography**

Chachra, V., et al. 1979, *Applications of Graph Theory Algorithms*, New York: North Holland.

Devroye, L. 1979, Inequalities for the completion times of stochastic PERT networks, *Mathematics of Operations Research*, vol. 4, no. 4, 441-447.

Dodin, B. & Elmaghraby S. 1985, Approximating the criticality indices of the activities in PERT networks, *Management Science*, vol. 31, no. 2, 207-223.

Dodin, B. 1985, Bounding the project completion time distribution in PERT networks, *Operations Research*, vol. 33, no. 4, 862-881.

Dodin, B. 1984, Determining the k most critical paths in PERT networks, *Operations Research*, vol. 32, no. 4, 859-877.

Donaldson, W. 1965, The estimation of the mean and variance of a 'PERT' activity time, *Operations Research*, vol. 13, 382-385.

Evans, T. & Minieka E. 1992, *Optimization Algorithms for Network and Graphs*, New York: Dekker.

Fisher, D., et al. 1985, Stochastic PERT networks: OP diagrams, critical paths and the project completion time, *Computers & Operations Research*, vol. 12, no. 5, 471-482.

Grubbs, F. 1962, Attempts to validate certain PERT statistics or 'picking on PERT', *Operations Research*, vol. 10, 912-915.

Hastings, N. & Peacock, J. 1975, *Statistical Distributions*, New York: John Wiley & Sons.

Kleindorfer, G. 1971, Bounding Distributions for a Stochastic Acyclic Network, *Operations Research*, vol. 19, 1589-1601.

Lawler, E. 1976, *Combinatorial Optimization: Networks and Matroids*, Chicago: Holt, Rinehart & Winston.

MacCrimmon, K. & Ryavec, C. 1964, An analytical study of the PERT assumptions, *Operations Research*, vol. 12, 16-37.

Malcolm, D., et al. 1959, Application of a Technique for Research and Development Program Evaluation, *Operations Research*, vol. 7, 646-669.

Martin, J. 1965, Distribution of the Time Through a Directed, Acyclic Network, *Operations Research*, vol. 13, 46-66.

Meredith, J. & Mantel, S. 2008, *Project Management: A Managerial Approach*, New York: John Wiley & Sons.

Ockham, W. 1990, *Philosophical Writings: A Selection*, Indianapolis: Hackett Pub. Co.

Parzen, E. 1960, *Modern Probability Theory and its Applications*, New York: John Wiley and Sons.

Phillips, D. & Diaz A. 1981, *Fundamentals of Network Analysis*, Englewood Cliffs, N.J.: Prentice Hall.

Sasieni, M. 1986, A note on PERT times, *Management Science*, vol. 32, vo. 12, 1652-1653.

Smolin, R. 1981, A critical path method for project control, *Interface Age, Computing for Business*, vol. 6, no. 8, 106-110.

Touran, A. & Wiser, E. 1992, Monte Carlo technique with correlated variables, *Journal of Construction Engineering and Management*, vol. 118, no. 2, 258-272.

Van Slyke, R. 1963, Monte Carlo methods and the PERT problem, *Operations Research*, vol. 11, no. 5, 839-860.

Wiest, J., & Levy, F. 1997, *A Management Guide to PERT/CPM*, Englewood Cliffs, N.J.: Prentice Hall.

Winston, W. 1994, *Operations Research: Applications and Algorithms*, Belmont, California: Duxbury Press.

Wolf, G. 1985, PERT/CPM: A tool for managing projects, *Industrial Management*, vol. 27, no. 1, 22-25.