



EXPECTED SHORTEST PATHS IN DYNAMIC AND STOCHASTIC TRAFFIC NETWORKS

LIPING FU*

Department of Civil Engineering, University of Waterloo, Waterloo, ON, Canada, N2L 3G1

and

L. R. RILETT

Department of Civil Engineering, Room 303D CE/TTI Tower, Texas A&M University, College Station, TX, 77843-3136, U.S.A.

(Received 25 July 1997; in revised form 25 March 1998)

Abstract—The dynamic and stochastic shortest path problem (DSSPP) is defined as finding the expected shortest path in a traffic network where the link travel times are modeled as a continuous-time stochastic process. The objective of this paper is to examine the properties of the problem and to identify a technique that can be used to solve the DSSPP given information that will be available in networks with Intelligent Transportation System (ITS) capabilities. The paper first identifies a set of relationships between the mean and variance of the travel time of a given path and the mean and variance of the dynamic and stochastic link travel times on these networks. Based on these relationships it is shown that the DSSPP is computationally intractable and traditional shortest path algorithms cannot guarantee an optimal solution. A heuristic algorithm based on the k -shortest path algorithm is subsequently proposed to solve the problem. Lastly, the trade-off between solution quality and computational efficiency of the proposed algorithm is demonstrated on a realistic network from Edmonton, Alberta. © 1998 Elsevier Science Ltd. All rights reserved

Keywords: shortest path problem, dynamic and stochastic network, k -shortest path problem, traffic network, Intelligent Transportation Systems, Route Guidance Systems

1. INTRODUCTION

In recent years there has been a resurgence of interest in the minimum path problem in transportation engineering. This is directly attributed to the recent developments in Intelligent Transportation Systems (ITS), and especially in the field of in-vehicle Route Guidance Systems (RGS). Central to any RGS is the algorithm that is used to find the optimal route from the origin to the destination. For most RGS currently under development the optimal route between an origin and destination is defined as the one with minimum expected travel time. This optimal route is commonly calculated by applying a Dijkstra type shortest path algorithm where the link travel times are treated deterministically (Dijkstra, 1959). The disadvantage of this type of deterministic treatment is that while it makes the shortest path problem tractable, it may, in fact, generate sub-optimal solutions.

Conversely, when both the dynamic and stochastic nature of link travel times are modeled explicitly, the optimal shortest path algorithms can become computationally inefficient and/or impractical for use within an actual application. The objective of this paper is to investigate the problem of finding the expected shortest path in a traffic network where the dynamic and stochastic nature of link travel times is modeled explicitly and to develop an algorithm which can provide improved solutions without significantly adding to the overall computation time.

The shortest path problem has been studied extensively in the fields of computer science, operations research, and transportation engineering. Most of the literature has focused on the problem in which the link travel cost (or weight) is assumed to be static and deterministic. Many

*Author for correspondence. Fax: 001 519 888 6197; e-mail: lfu@uwaterloo.ca

efficient algorithms have been developed (Bellman, 1958; Dijkstra, 1959; Dreyfus, 1969) and in this paper, these algorithms are referred to as the standard shortest path algorithms. It should be noted that the standard shortest path algorithms also have been found to be applicable to compute shortest paths in time-dependent (but *not* stochastic) networks (Dreyfus, 1969; Orda and Rom, 1990; Kaufman *et al.*, 1993; Ziliaskopoulos and Mahmassani, 1993; Chabini, 1997).

Frank (1969) and Mirchandani (1976) studied the problem of determining the probability distribution of the shortest path length in a stochastic network where link travel times are random variables but not time dependent. Loui (1983), Mirchandani and Soroush (1986), and Murthy and Sarkar (1996) studied the variations of the shortest path problems in stochastic networks by considering different types of cost functions. It was found that if the objective is to identify the expected shortest path (or linear cost function), then the problem simply reduces to a deterministic shortest path problem in a network where the random link travel times are replaced by their expected values. Therefore, the efficient standard shortest path algorithms still can be used to find the expected shortest paths in a static and stochastic network.

Hall (1986) first investigated the shortest path problem in a transportation network where the link travel times are random and time-dependent and demonstrated that the standard shortest path algorithm may fail to find the expected shortest path in these networks. An optimal dynamic programming based algorithm was proposed to find the shortest paths and this algorithm was demonstrated on a small transit network example. Although not explicitly stated, the paper only considers the case where link travel times are modeled as discrete-time stochastic processes and the proposed algorithm is viable only for solving problems of small networks because of computational constraints.

The objective of this paper is to extend the shortest path problem in dynamic and stochastic networks (DSSPP) where link travel times are defined as continuous time stochastic processes. The paper first defines the DSSPP. Next, a general probability-based formula for calculating the mean and variance of the travel time for a given path is developed. The emphasis is on estimating these path parameters using the mean and variance of link travel time as a function of time of day—information typically available in transportation networks with ITS. An extensive analysis of the properties associated with the DSSPP is provided and a heuristic algorithm based on the k -shortest path algorithm is proposed. Finally, the trade-off between solution quality and computational efficiency of the proposed algorithm is demonstrated on a realistic network from Edmonton, Alberta.

2. PROBLEM DEFINITION

Consider a traffic network represented by a directed graph consisting of a finite set of nodes and links. Each link in the network has an associated generalized cost which could be a combination of travel time, direct cost and travel distance. Without loss of generality, this paper will use travel time to represent this generalized cost. It is assumed that the link travel times on some or all of the links in the network are random variables and the probability distributions of link travel times are dependent on the time of day (i.e. the time a link is entered). Consequently, the travel time on these type of links can be modeled as a *continuous-time stochastic process*. Denote $\{X_\alpha(t), t \in T\}$ as a stochastic process of travel time on link α , where $X_\alpha(t)$ is the travel time for vehicles entering link α at time t , and T is a continuous parameter set. In this paper T is the time range examined although for generality it will be assumed that $T = \mathcal{R}_+ = [0, \infty)$. For each time instance t , $X_\alpha(t)$ is considered as a continuous random variable with its *first-order* probability density function (PDF) denoted by $f_{X_\alpha}(x_\alpha, t)$. A network where the link travel time is modeled as a stochastic process is referred to as a *dynamic and stochastic network* in this paper.

Furthermore, it is assumed that the travel times on individual links at a particular point in time are statistically independent. It should be noted that this assumption does not mean that the correlation between link travel times will be ignored. Because the probability distributions of the link travel times are modeled as functions of the time of day, the time of day correlation between individual links is explicitly taken into account. For example, on a typical network it would be expected that the travel times on individual links would all be higher than average during peak periods and lower than average during off-peak periods because of temporal changes in traffic volume. This time of day correlation in link travel times is modeled directly within the individual

links' stochastic processes. It should be noted that in order to consider the correlation between individual links travel times *at a particular point in time* a more disaggregate analysis and more comprehensive data than that typically collected and stored in ITS facilities would be required.

The mean of the stochastic process $\{X_a(t), t \in T\}$ corresponding to its first-order PDF is represented by $\mu_{X_a}(t)$ and is defined as follows:

$$\mu_{X_a}(t) = E[X_a(t)] = \int_0^{+\infty} x_a f_{X_a}(x_a, t) dx_a \quad (1)$$

One measure of dispersion of the random variable $X_a(t)$ about its mean $\mu_{X_a}(t)$ is the variance denoted as $\nu_{X_a}(t)$ which is defined in eqn (2)

$$\nu_{X_a}(t) = E[(X_a(t) - \mu_{X_a}(t))^2] = \int_0^{+\infty} (x_a - \mu_{X_a}(t))^2 f_{X_a}(x_a, t) dx_a \quad (2)$$

It is important to note that the PDF of the link travel time, $f_{X_a}(x_a, t)$, will not be available in most practical situations. However, for transportation networks with ITS capabilities where link travel time data are automatically collected, the mean and standard deviation of link travel times over discrete periods will be readily available. That is, estimates of the mean and variance shown in eqns (1) and (2) will be available over pre-defined discrete, rather than continuous, time intervals.

Let p denote a simple path from an origin node 1 to a destination node N and P is the set of all paths p . The problem is to find path p^* which has the lowest expected travel time from node 1 to node N corresponding to a given departure time at node 1. This problem will be referred to as *dynamic and stochastic shortest path problem* (DSSPP) in this paper. If the random variable W_p denotes the travel time on path p then the expected travel time along path p is defined in eqn (3).

$$E[W_p] = \int_0^{+\infty} w_p f_{W_p}(w_p) dw_p \quad (3)$$

where $f_{w_p}(w_p)$ is the PDF of W_p .

Now, the DSSPP can be stated formally as

$$(\text{DSSPP}) \quad p^* = \arg \min_{p \in P} E[W_p] \quad (4)$$

In order to find a solution to the DSSPP for a given O-D pair, it is a necessary condition to calculate the expectation of travel time of a given path (e.g. $E[W_p]$) based on the travel times of individual links within the path. This could be done easily if the travel time on links are only stochastic but not time-dependent or only time-dependent but not stochastic. However, it is not a trivial problem when the link travel times are both time-dependent (or dynamic) and stochastic. As a result, the problem of estimating the expected travel time of a given path based on the link travel time information will be analyzed first in the following section.

3. TRAVEL TIME ON A GIVEN PATH

Consider a particular path p from an origin node 1 to a destination node N in a dynamic and stochastic network as shown in Fig. 1 and assume that a series of travel experiments are conducted along this path. Each experiment represents a travel departing from the origin node at the exact same time and traveling along the path to the destination node N . Due to the dynamic and stochastic attributes of the network as defined above, the outcomes of the experiment, which include both the arrival time at each node and the travel time on each link, are random variables. The travel time on the given path also will be a random variable and its distribution will depend on the link travel time distribution of each link on the route and the departure time at the origin node. As

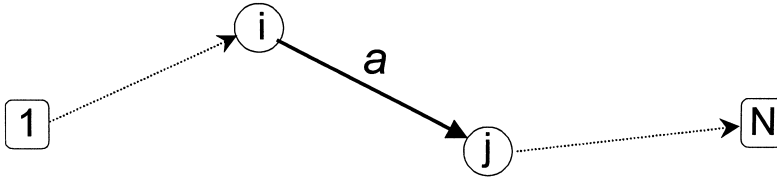


Fig. 1. A path p from origin node 1 to destination node N .

the estimation of the route travel time is equivalent to the estimation of the arrival time at the destination node, the arrival time will be used in the following discussion.

Let the random variable Y_i denote the arrival time at node i . Assuming that there is no waiting time at the node, then Y_i is equal to the departure time at node i or the time link a is entered. The departure at the origin node 1, Y_1 , is assumed to be deterministic and known *a priori*. It should be noted that this deterministic assumption might be relaxed without loss of generality. The PDF of Y_i is represented by $f_{Y_i}(y_i)$. As a random variable, the travel time is completely specified by its distribution ($f_{Y_i}(y_i)$). The problem is therefore to estimate $f_{Y_i}(y_i)$ ($i = 2 \dots N$) based on the departure time at node 1 and the travel times on the individual links. However, it is not a trivial problem to derive the PDF of Y_i when the link travel times are both dynamic and stochastic (or stochastic process), as it is illustrated in the following example.

A two-link network is shown in Fig. 2. The travel time on link a (t_a) is normally distributed with a mean travel time of 5 min and a standard deviation of 1 min. The travel time on link b , t_b , is a function of the time entering the link (T_2): $t_b = 10 + 0.5(T_2 - 5)^2$. If it is assumed that a trip departs at Node 1 at time zero, T_2 will be equal to t_a . Based on the given information on the travel time, the arrival time at Node 3 (T_3) is equal to $T_2 + 10 + 0.5(T_2 - 5)^2$. In spite of the fact that T_3 is a simple function of the normally distributed random variable t_a , its PDF is not easily obtainable. It is clear that identifying the probability density function of the arrival time at downstream links will become quickly intractable and this problem will occur even if simple functions are used.

Given the above problems in estimating the path travel time PDF, the focus of this paper is on estimating the mean and variance of the arrival time at node i and these will be denoted as $E[Y_i]$ and $Var[Y_i]$, respectively. These parameters are defined in eqns (5) and (6).

$$E[Y_i] = \int_0^{+\infty} y_i f_{Y_i}(y_i) dy_i \quad (5)$$

$$Var[Y_i] = E[(Y_i - E[Y_i])^2] = \int_0^{+\infty} (y_i - E[Y_i])^2 f_{Y_i}(y_i) dy_i \quad (6)$$

Denote random variable Z_a as the travel time on link a under a given experiment. It should be noted the distribution of Z_a is conditional to the specific experiment and therefore cannot be known directly until the parameters of the experiment, which include both the path and the departure time, are defined. This is in contrast to the link travel time which is represented as a

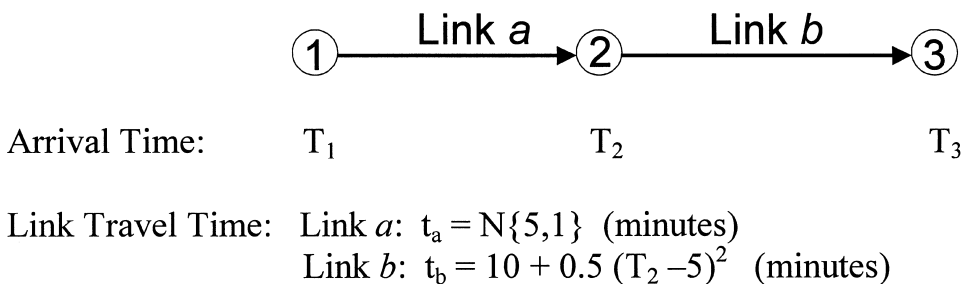


Fig. 2. A two-link dynamic and stochastic network.

continuous-time stochastic process $\{X_a(t), t \in T\}$, and is solely a function of the time entering the link. Because the link travel time on link a is only dependent on the time the link is entered, the probability distribution of Z_a under a given time (i.e. $Y_i = y_i$) will be the same as the probability distribution of $X_a(y_i)$ as shown in eqn (7)

$$P\{Z_a < x | Y_i = y_i\} = P\{X_a(y_i) < x\}, \quad x \in \mathcal{R} \quad (7)$$

Equation (7) also implies that the random variable $Z_a | Y_i = y_i$ has the same mean and variance as the random variable $X_a(y_i)$, or

$$E[Z_a | Y_i = y_i] = \mu_{X_a}(y_i) \quad (8)$$

$$\text{Var}[Z_a | Y_i = y_i] = \nu_{X_a}(y_i) \quad (9)$$

For each experiment described above, the path travel time is equal to the sum of the travel time of all the links along the path. This path travel time can be obtained by calculating the arrival time at each node along the path using a recursive formula until the destination node is reached (i.e. $j = N$) as shown in eqn (10).

$$Y_j = Y_i + Z_a \quad (10)$$

The following sections present several approximation models for estimating the mean and variance of the route travel time based on eqn (10) when the link travel time is dynamic and stochastic.

3.1. Mean of route travel time

Considering eqn (10), the relationship between the mean arrival times at node i and the mean arrival time at the downstream node j is:

$$E[Y_j] = E[Y_i] + E[Z_a] \quad (11)$$

This relationship can be transformed further (Ross, 1989):

$$E[Y_j] = E[Y_i] + E[E[Z_a | Y_i]] \quad (12)$$

Based on eqn (8), the recursive formula for calculating the expected travel time of the route is therefore:

$$E[Y_j] = E[Y_i] + E[\mu_{X_a}(Y_i)] \quad (13)$$

The second term in eqn (13) is defined by the integral shown in eqn (14)

$$E[\mu_{X_a}(Y_i)] = \int_0^{+\infty} \mu_{X_a}(y_i) f_{Y_i}(y_i) dy_i \quad (14)$$

Obviously eqn (14) may be applied only if the PDF of the arrival time, $f_{Y_i}(y_i)$, is available, and this means a recursive formula for estimating the PDF of the arrival time also must be derived. As discussed above, this is impossible in realistic transportation networks because only information available about the link travel times is the estimate of the first two moments of the link travel time PDF. These estimates, which are in fact the historical sample mean and sample variance or forecast mean and variance, are not continuous but rather are calculated for discrete periods throughout the day. Furthermore, even if the PDF of the link travel time as a function of the time of day could be derived from existing data, the derivation of the PDF of the arrival time shown in eqn (13) is mathematically impractical, as illustrated (Fig. 2). Consequently, a recursive formula for calculating the arrival time at a particular node based on approximation techniques is required.

In order to determine the second part of eqn (13), the function $\mu_{X_a}(t)$ may be expanded as a Taylor's series around the point $t = E[Y_i]$ as shown in eqn (15). Note that this step requires that the function $\mu_{X_a}(t)$ is differentiable at point $t = E[Y_i]$.

$$\mu_{X_a}(t) = \mu_{X_a}(E[Y_i]) + \mu'_{X_a}(E[Y_i]) \cdot (t - E[Y_i]) + \frac{1}{2} \mu''_{X_a}(E[Y_i]) \cdot (t - E[Y_i])^2 + \dots \quad (15)$$

If the series is truncated at the linear terms (or assuming that the second and higher order derivatives are equal to zero) and then applied in eqn (14), the first order approximation of $E[\mu_{X_a}(Y_i)]$ is obtained:

$$\begin{aligned} E[\mu_{X_a}(Y_i)] &\cong \int_0^{+\infty} \left\{ \mu_{X_a}(E[Y_i]) + \mu'_{X_a}(E[Y_i]) \cdot (y_i - E[Y_i]) \right\} f_{Y_i}(y_i) dy_i \\ &= \mu_{X_a}(E[Y_i]) \cdot \int_0^{+\infty} f_{Y_i}(y_i) dy_i + 0 \\ &= \mu_{X_a}(E[Y_i]) \end{aligned} \quad (16)$$

Therefore, the *first order approximation model* of the recursive formula is:

$$E[Y_j] \cong E[Y_i] + \mu_{X_a}(E[Y_i]) \quad (17)$$

The first order approximation model shown in eqn (17) can be improved by including higher order terms of the Taylor series. For example, if the second order term in eqn (15) is included, the second order approximation of $E[\mu_{X_a}(Y_i)]$ is accordingly:

$$\begin{aligned} E[\mu_{X_a}(Y)] &\cong \int_0^{+\infty} \left\{ \mu_{X_a}(E[Y_i]) + \mu'_{X_a}(E[Y_i]) \cdot (y_i - E[Y_i]) + \frac{1}{2} \mu''_{X_a}(E[Y_i]) \cdot (y_i - E[Y_i])^2 \right\} f_{Y_i}(y_i) dy_i \\ &= \mu_{X_a}(E[Y_i]) \cdot \int_0^{+\infty} f_{Y_i}(y_i) dy_i + 0 + \frac{1}{2} \mu''_{X_a}(E[Y_i]) \cdot \int_0^{+\infty} (y_i - E[Y_i])^2 f_{Y_i}(y_i) dy_i \\ &= \mu_{X_a}(E[Y_i]) + \frac{1}{2} \mu''_{X_a}(E[Y_i]) \cdot \text{Var}[Y_i] \end{aligned} \quad (18)$$

Using eqn (18) the *second order approximation model* of the mean arrival time can be obtained as shown below:

$$E[Y_j] \cong E[Y_i] + \mu_{X_a}(E[Y_i]) + \frac{1}{2} \mu''_{X_a}(E[Y_i]) \cdot \text{Var}[Y_i] \quad (19)$$

The relative quality of the first and second approximation models can be illustrated using the same example shown in Fig. 2. The expected arrival time at Node 3 can be directly obtained as shown below:

$$\begin{aligned} E[T_3] &= E[t_a + 10 + 0.5 \cdot (t_a - 5)^2] \\ &= 0.5 \cdot E[t_a^2] - 4 \cdot E[t_a] + 22.5 \\ &= 15.5 \end{aligned}$$

Alternatively, $E[T_3]$ can be estimated by using the first order approximation model [eqn (17)] and it may be seen that the expected arrival time is underestimated in this example.

$$\begin{aligned}
E[T_3] &\cong E[T_2] + \mu_b(E[T_2]) \\
&= E[T_1] + \mu_a(E[T_1]) + \mu_b(E[T_2]) \\
&= 0 + 5 + 10 + 0.5 \times (5 - 5)^2 = 15
\end{aligned}$$

Conversely, if the second order approximation model [eqn (19)] is used the expected arrival time is equal to the actual arrival time as shown below.

$$\begin{aligned}
E[T_3] &\cong E[T_2] + \mu_b(E[T_2]) + \frac{1}{2}\mu_b''(E[T_2])\text{Var}[T_2] \\
&= E[T_1] + \mu_a(E[T_1]) + \mu_b(E[T_2]) + \frac{1}{2}\mu_b''(E[T_2])\sigma_a^2 \\
&= 0 + 5 + 10 + 0.5 \times (5 - 5)^2 + \frac{1}{2} \times 1.0 \times 1^2 = 15.5
\end{aligned}$$

Based on the two approximation models that were developed above the following observations may be obtained.

Remark 1. The first order approximation model for the expected arrival time at node i is equivalent to that obtained assuming a dynamic and deterministic treatment. That is, the expected route travel time is found by substituting the average link travel time given a particular arrival time in place of the link travel time random variable and then summing over all links to calculate the route travel time. From eqn (19), it can be expected that this model may be acceptable when the variance of the arrival time is small relative to the mean, or the mean link travel time is approximately a linear function of time of day (i.e. $\mu_{X_a}(E[Y_i]) \approx 0$).

Remark 2. The reasonableness of the second order approximation model shown in eqn (19) can be illustrated using the following simple example. Consider a network with one link connecting two nodes i and j . Assume that the time entering the link, Y_i , may be modeled as a normally distributed random variable with mean represented by $E[Y_i]$. The link travel times are assumed to be deterministic and dynamic and the second derivative of the link travel time when the vehicle enters the link at time y_i is zero (i.e. $\mu_{X_a}''(y_i) = 0$). In this example the three potential dynamic link travel time patterns are linear, convex, and concave, as shown in Fig. 3. It can be found that when the link travel time is constant, the link travel time is the same for any arrival time. Consequently, the expected link travel time is $\mu(E[Y_i])$, which is the same as that obtained from eqn (19) when

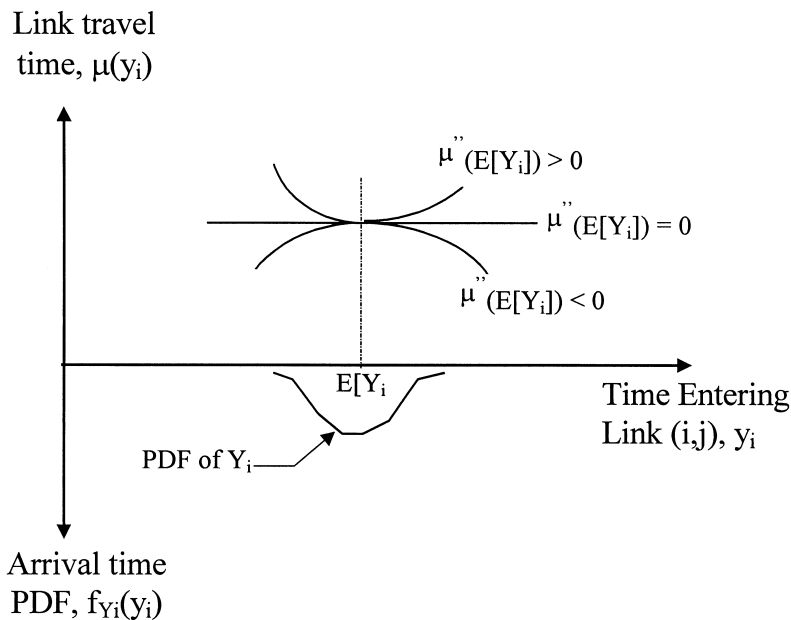


Fig. 3. The effect of link travel time pattern on the estimation of the expected link travel time.

$\mu''(E[Y_i]) = 0$. When the link travel time is a convex function of arrival time, i.e. $\mu''(E[Y_i]) > 0$, then the travel time under any realization of the arrival time always will be greater than $\mu(E[Y_i])$ and therefore the expected travel time should be greater than $\mu(E[Y_i])$. This result is compatible with the result from eqn (19) because the last term in eqn (19), $\frac{1}{2}\mu''(E[Y_i])\text{Var}[Y_i]$, is always greater than zero. A similar explanation for the case when the link travel time is concave also may be demonstrated.

3.2. Variance of route travel time

In order to use the second order approximation method the variance of the arrival time is required. Additionally, when deciding among competing routes it also would be desirable to have an estimate of the variability of the arrival time at node N (which is equivalent to the route travel time variability). The variance of the arrival time at a downstream node may be derived using eqn (20) which is based on the recursive formula shown in eqn (10).

$$\text{Var}[Y_j] = \text{Var}[Y_i] + \text{Var}[Z_a] + 2\text{COV}(Y_i, Z_a) \quad (20)$$

The last two parts of eqn (20) can be transformed further based on eqns (8) and (9) (Ross, 1989):

$$\begin{aligned} \text{Var}[Z_a] &= E[\text{Var}[Z_a|Y_i]] + \text{Var}[E[Z_a|Y_i]] \\ &= E[\text{Var}[X_a(Y_i)]] + \text{Var}[E[X_a(Y_i)]] \\ &= E[v_{X_a}(Y_i)] + \text{Var}[\mu_{X_a}(Y_i)] \end{aligned} \quad (21)$$

and

$$\begin{aligned} \text{COV}(Y_i, Z_a) &= E[Y_i \cdot Z_a] - E[Y_i]E[Z_a] \\ &= E[E[Y_i \cdot Z_a|Y_i]] - E[Y_i][E[Z_a|Y_i]] \\ &= E[Y_i \cdot E[Z_a|Y_i]] - E[Y_i]E[\mu_{X_a}(Y_i)] \\ &= E[Y_i \mu_{X_a}(Y_i)] - E[Y_i]E[\mu_{X_a}(Y_i)] \end{aligned} \quad (22)$$

The variance of the arrival time at node j is shown in eqn (23).

$$\text{Var}[Y_j] = \text{Var}[Y_i] + E[v_{X_a}(Y_i)] + \text{Var}[\mu_{X_a}(Y_i)] + 2E[Y_i \mu_{X_a}(Y_i)] - 2E[Y_i]E[\mu_{X_a}(Y_i)] \quad (23)$$

Employing similar reasoning as that in the mean route travel time analysis, it is easy to show that it will be a significant challenge to identify the functions $\mu_{X_a}(Y_i)$ and $v_{X_a}(Y_i)$ under realistic assumptions. If the functions $\mu_{X_a}(Y_i)$ and $v_{X_a}(Y_i)$ are replaced with truncated Taylor series expansions about point $E[Y_i]$, then approximation models of the recursive model shown in eqn (23) can be obtained. The *first order approximation model* is obtained by assuming that the second and higher derivatives of $\mu_{X_a}(t)$ and $v_{X_a}(Y_i)$ are equal to zero, as shown in eqn (24).

$$\text{Var}[Y_j] \cong A \cdot \text{Var}[Y_i] + v_{X_a}(E[Y_i]) \quad (24)$$

where

$$A \cong \left\{ 1 + \mu'_{X_a}(E[Y_i]) \right\}^2$$

By assuming the third and higher derivatives of $\mu_{X_a}(t)$ and $v_{X_a}(Y_i)$ are equal to zero, the *second order approximation* of eqn (23) can be obtained as shown in eqn (25).

$$\begin{aligned} \text{Var}[Y_j] &\cong \left\{ \left(1 + \mu'_{X_a}(E[Y_i]) \right)^2 + \frac{1}{2} v''_{X_a}(E[Y_i]) \cdot \frac{1}{4} \mu''_{X_a}(E[Y_i]) \cdot \text{Var}[Y_i] \right\} \cdot \text{Var}[Y_i] \\ &\quad + v_{X_a}(E[Y_i]) \\ &\quad + \left(1 + \mu'_{X_a}(E[Y_i]) \right) \cdot \mu''_{X_a}(E[Y_i]) \cdot E[(Y_i - E[Y_i])^3] \\ &\quad + \frac{1}{4} \mu''_{X_a}(E[Y_i])^2 \cdot E[(Y_i - E[Y_i])^4] \end{aligned} \quad (25)$$

As shown in eqn (25), the second order approximation, although potentially leading to a better solution, would involve identifying the third and fourth central moments of the arrival time Y_i . This means that in order to use eqn (24) a recursive formula for estimating the third and fourth order moments of the arrival times, and by extension the third and fourth order model of the link travel times, are required. However, if we assume that an estimate of the coefficient of skewness (β_1) and the coefficient of kurtosis (β_2) is available, then the second order model is shown below.

$$\begin{aligned} \text{Var}[Y_j] \cong & \left\{ \left(1 + \mu'_{X_a}(E[Y_i]) \right)^2 + \frac{1}{2} v''_{X_a}(E[Y_i]) - \frac{1}{4} \mu''_{X_a^2}(E[Y_i]) \cdot \text{Var}[Y_i] \right\} \cdot \text{Var}[Y_i] \\ & + v_{X_a}(E[Y_i]) \\ & + \left(1 + \mu'_{X_a}(E[Y_i]) \right) \cdot \mu''_{X_a}(E[Y_i]) \cdot (\beta_1 \text{Var}^3[Y_i])^{\frac{1}{3}} \\ & + \frac{1}{4} \mu''_{X_a^2}(E[Y_i]) \cdot \beta_2 \text{Var}^2[Y_i] \end{aligned}$$

Additionally, if it is further assumed that the arrival time is symmetric ($\beta_1 = 0$) and that Y_i is neither platykurtic or leptokurtic (i.e. $\beta_2 = 3$) then the second order approximation model can be further simplified as shown below. Note that these latter assumptions would be good approximations if the arrival time was approximately normally distributed.

$$\text{Var}[Y_j] \cong (A + B) \cdot \text{Var}[Y_i] + v_{X_a}(E[Y_i]) \quad (26)$$

where A was defined in eqn (24) and B is defined below.

$$B = \frac{1}{2} \left\{ v''_{X_a}(E[Y_i]) + \mu''_{X_a^2}(E[Y_i]) \cdot \text{Var}[Y_i] \right\}$$

The application of the approximation models can be illustrated using the example shown in Fig. 2. The real value of the variance of the arrival time at Node 3 can be directly obtained as shown below.

$$\begin{aligned} \text{Var}[T_3] &= \text{Var}[t_a + 10 + 0.5 \times (t_a - 5)^2] \\ &= E[(t_a + 10 + 0.5 \times (t_a - 5)^2)^2] - \{E[t_a + 10 + 0.5 \times (t_a - 5)^2]\}^2 \\ &= 0.25 \times E[(t_a - 5)^4] + E[(t_a - 5)^3] + 16 \times E[(t_a - 5)^2] + 30 \times E[(t_a - 5)] + 225 - 15.5^2 \\ &= 0.25 \times 3 \times 1^4 + 0 + 16 + 0 - 15.25 \\ &= 1.5 \end{aligned}$$

If the first order approximation model [eqn (24)] is used, the estimate of $\text{Var}[T_3]$ is 50% smaller than the true value as shown below.

$$\begin{aligned} A &\cong \{1 + \mu'_b(E[T_2])\}^2 \\ &= \{1 + \mu'_b(E[t_b])\}^2 \\ &= \{1 + 0\}^2 \\ &= 1 \end{aligned}$$

$$\begin{aligned} \text{Var}[Y_j] &\cong A \cdot \text{Var}[T_2] + v_b(E[T_2]) \\ &= 1 \times 1 + 0 \\ &= 1 \end{aligned}$$

If the second order approximation model [eqn (26)] is used, then estimate of $\text{Var}[T_3]$ is equal to the correct value.

$$\begin{aligned}
B &\cong \frac{1}{2} \{v_b''(E[T_2]) + \mu_b''^2(E[T_2]) \cdot \text{Var}[T_2]\} \\
&= \frac{1}{2} \{0 + 1 \times 1\} \\
&= 0.5
\end{aligned}$$

$$\begin{aligned}
\text{Var}[Y_j] &\cong (A + B) \cdot \text{Var}[T_2] + v_b([T_2]) \\
&= (1 + 0.5) \times 1 + 0 \\
&= 1.5
\end{aligned}$$

Based on the above approximation models, the following observations may be obtained.

Remark 3. The first order approximation model of eqn (24) shows that the variance of the route travel time is dependent not only on the variance associated with the link travel time but also on the variation associated with the time of day. This model can be partially verified by using an example similar to that described in Remark 2. Consider a one-link situation. Assume that the time entering the link, Y , may be modeled as a uniformly distributed random variable over the range a, b (i.e. $U\{a, b\}$). In addition, the link travel time (X) is assumed to be deterministic and can be represented by a linear function of the time the link is entered as shown in Fig. 4. If the link is entered at time y_i , the link travel time is equal to ky_i . It is relatively straightforward to show that the arrival time at the exit node of the link also is uniformly distributed but with different parameters, (i.e. $U\{a + ka, b + kb\}$). Therefore, the variance of the arrival time at the exiting node of the link is essentially $(1 + k)^2$ times greater than the variance of the arrival time at the entering node of the link. The same conclusion can be obtained directly from eqn (24) as shown in eqn (27).

$$\text{Var}[Y_j] = \{1 + \mu'_{x_a}(E[Y_i])\}^2 \cdot \text{Var}[Y_i] = \{1 + k\}^2 \cdot \text{Var}[Y_i] \quad (27)$$

Remark 4. It can be anticipated that the difference between the second order approximation model [eqn (26)] and the first order approximation model [eqn (24)] may be trivial in applications on realistic traffic networks. This is because the variance in link travel time for a particular departure time is relatively small compared to the variance in link travel time associated with the time of day. Therefore, the second order derivatives could be expected to be negligible. However, in the

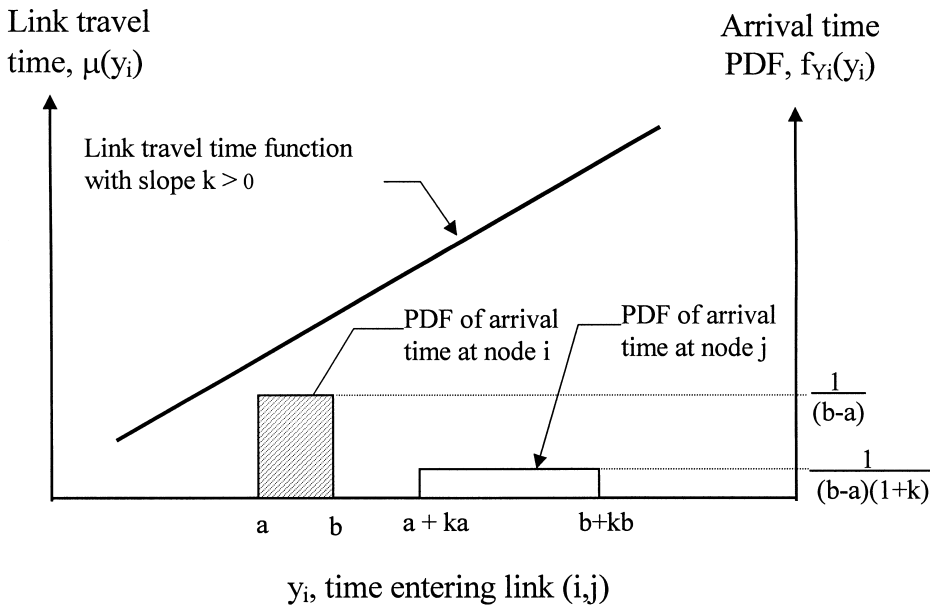


Fig. 4. Link travel time pattern and link travel time variance.

situation where the link travel times on the network changes relatively rapidly, such as when the peak period is starting or ending, the second order approximation will provide a better estimate.

3.3. Link travel time approximation

The application of the approximation models developed above requires that the derivatives of the mean and variance of the link travel time must be available. While the mean and standard deviation of link travel time data are available in most realistic transportation networks they are typically stored in a discrete form. Consequently, the discrete link travel time mean and variance have to be approximated by smooth functions before they can be used in the procedure developed in this paper. The following section discusses how the link travel time may be approximated using a differential function representing the recurring traffic congestion situation.

Under normal traffic situations, the link travel times may be assumed to be stable from day to day and therefore the historical link travel data can be used. These historical data may be obtained from various data sources such as roadside detectors, probe vehicles, or even traffic models. Due to inherent fluctuations in traffic demand and errors in measurement related to each data source, the link travel time obtained is not a fixed value even for the same time moment of two similar days. Figure 5(a) shows a hypothetical example of the travel time data for a link. To use these data to estimate the mean and variance of the link travel time, the time horizon is usually divided into time periods so that the number of data points for each period is high enough to provide statistically confident estimates and to minimize data management problems. For example, the mean and variance of link travel times from the Houston TransGuide ATMS project and the San Antonio TranStar ATMS project is available in 5 min intervals (Turner *et al.*, 1997). As shown in Fig. 5(b), the time dependent link travel time is modeled as a discrete or step function. When the arrival time falls into a specific time period, the average link travel time of that time period can be used. An improvement on this method is to use a pairwise linear function, as shown in Fig. 5(c) (Rilett, 1992).

In this paper the application of the second order approximation models requires the mean and variance of the link travel time to have a second order derivative and therefore a second order polynomial is used to approximate the link travel time. If the mean link travel time is $\mu(Y)$ where Y is the time entering the link, then the general form of the function is:

$$\mu = b_0 + b_1 Y + b_2 Y^2 \quad (28)$$

For the purpose of the route travel time estimation, the major point of interest is the link travel time pattern in the vicinity of the mean arrival time on the link. It will be neither necessary nor efficient to fit all the data with one continuous function and therefore a three-point approximation method will be used in this paper. If the time in which the link is entered falls in time period k , the link travel times from time period $k - 1$ to time period $k + 1$ are approximated by eqn (28) which goes through these three points as shown in Fig. 5(d). Because the system consists of three linear equations with three variables, the parameters can be readily identified. If μ_k represents the mean link travel time for interval k with the middle time of the interval noted as y_k , the solution is shown in eqn (29).

$$\{\mathbf{b}\} = [\mathbf{T}]^{-1}\{\mu\} \quad (29)$$

where:

$$\begin{aligned} \{\mathbf{b}\} &= \{b_0, b_1, b_2\}' \\ [\mathbf{T}] &= \{(1, y_{k-1}, y_{k-1}^2), (1, y_k, y_k^2), (1, y_{k+1}, y_{k+1}^2)\}' \\ \{\mu\} &= \{\mu_{k-1}, \mu_k, \mu_{k+1}\}' \end{aligned}$$

Therefore, once the arrival time at node i is known, the mean and variance of the link travel time and their derivatives can be quickly calculated without any significant additional computation burden.

It should be noted that the approximation method proposed above is a relatively simple approach and a more comprehensive method could be used to approximate the link travel time

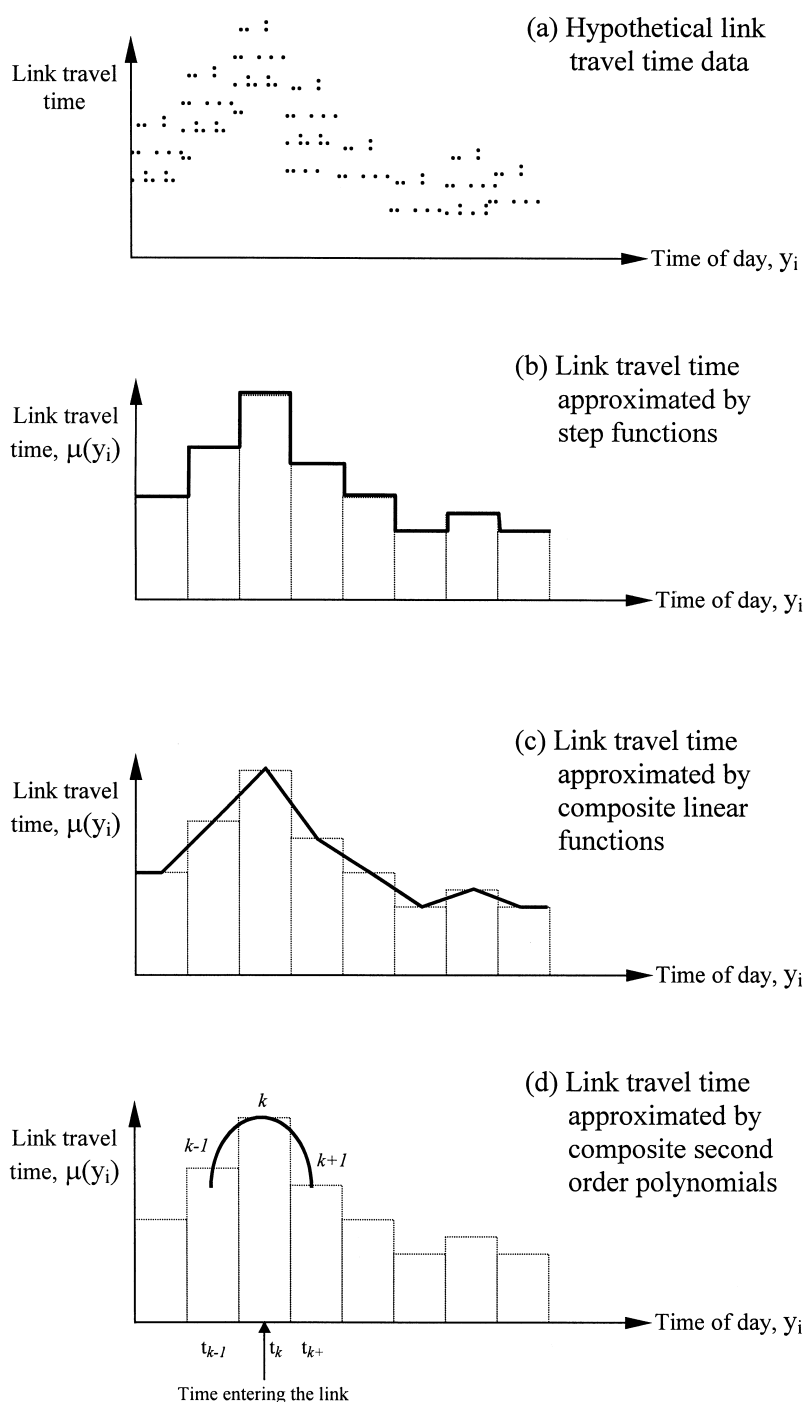


Fig. 5. Schematic illustration of link travel time approximation methods.

under different situations. For example, the variance of the arrival time and the size of the link travel time interval also can be taken into account during the approximation procedure. For situations where the travel time interval is smaller or the arrival time variance is larger, more than three intervals may need to be considered in the approximation method. The underlying relation is schematically illustrated in Fig. 6. In this situation it can be seen that a function which fits the intervals from $k - 2$ to $k + 2$ would result in a better approximation. It should be noted that the models developed in this paper can be easily adapted to handle these latter approximation methods. The only modification will be to change the function format shown in eqn (28) and the parameter estimation methods shown in eqn (31).

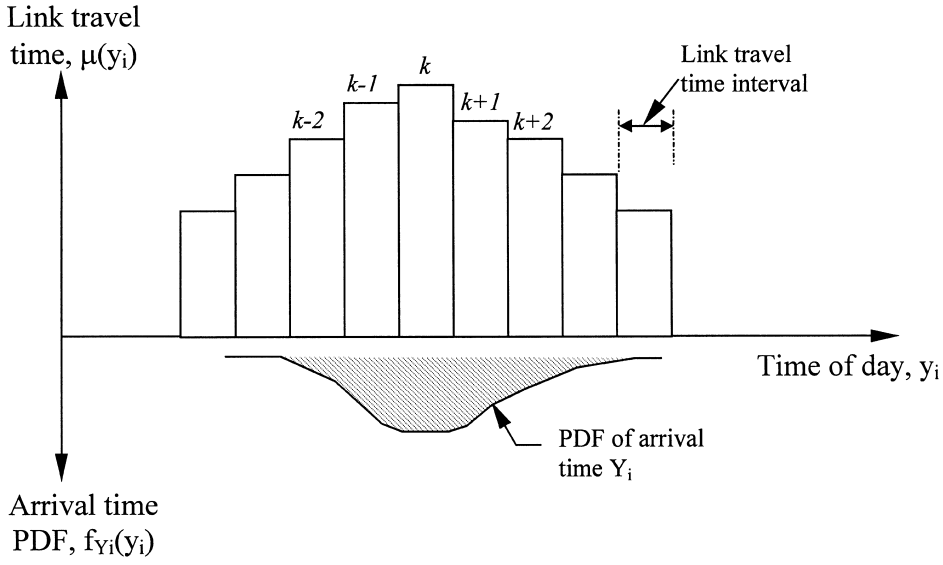


Fig. 6. Link travel time approximation and arrival time pattern.

4. DSSPP: PROPERTIES

From eqn (19), the following properties of the DSSPP may be observed.

Property 1. If the mean of the link travel time of at least one link in a network is non-linear, the standard shortest path algorithms may fail to find the expected shortest path between two nodes in the network.

This property can be illustrated using an example network shown in Fig. 7. The network is composed of two sub paths ($p1$ and $p2$) from the origin node s to an intermediate node i , and one link (i, j) from node i to the destination node j . Assume that the travel time on $p1$ is deterministic and that the travel time on $p2$ is stochastic. The travel time on link (i, j), $\mu_{(i,j)}$, is deterministic but

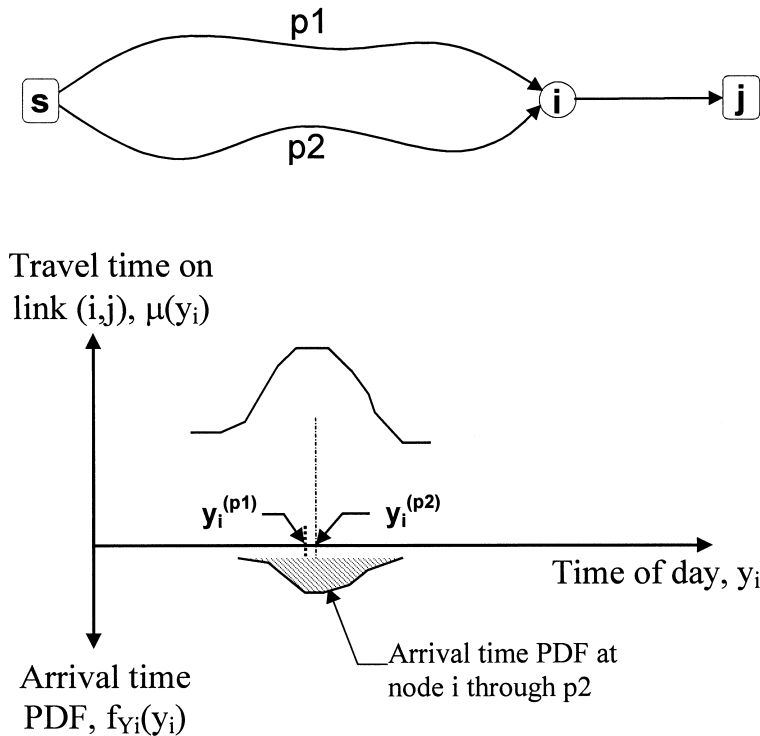


Fig. 7. A simple dynamic and stochastic network.

changes with time in a non-linear fashion as shown in Fig. 7. If the expected arrival time at node i through $p1$, $y_i^{(p1)}$, is marginally less than through $p2$ ($y_i^{(p2)}$) then subpath $p1$ is the minimum expected route from node s to node i . On the other hand, eqn (19) shows that the expected minimum arrival time at node j not only depends on the expected arrival time at node i , but also on the variance of the arrival time at node i and the second derivative of the mean travel time on link (ij) . Given that the travel time on link (ij) is concave and hence its second derivative is negative, it is possible that subpath $p2$ is on the expected minimum path from s to j . In short, Bellman's "principle of optimality" which states that any subpath of a shortest path must be a shortest path (Denardo, 1982), does not hold in a DSSPP.

The above property also implies that the standard shortest path algorithms could be applicable if the link travel time function $\mu(t)$ is linear, or in a more relaxed sense, close to linear within the local range.

Property 2. The DSSPP is computationally intractable. Consider a simple acyclic type of network shown in Fig. 8. The network has N nodes sequentially labelled from 1 to N , with 2 links between each successive pair of nodes. The problem is to find the expected shortest path from node 1 to node N . If the network is deterministic, the problem can be solved by finding the shortest path to node 2 first, then node 3, and onward until node N . The computation time is $O(n)$. This type of procedure will not work in a dynamic and stochastic network because the optimal path to node i does not have to be part of the optimal path to node $i + 1$. This means that all the 2^{N-1} paths from 1 to N must be examined before the optimal path can be definitely identified. The computation time therefore grows exponentially with the number of nodes N , which indicates that the DSSPP is computationally intractable (Gary and Johnson, 1979).

5. A HEURISTIC ALGORITHM TO CALCULATE THE EXPECTED SHORTEST PATH

In Section 4 it was shown that standard shortest path algorithm may not identify the expected shortest path on a dynamic and stochastic network. In addition, the DSSPP is intractable in the sense that there is no polynomial time algorithm, like the standard shortest path algorithm, to solve this problem. Therefore, a heuristic algorithm was developed in this paper to identify 'optimal' routes.

The heuristic algorithm proposed in this paper is based on the idea of examining the potential better paths instead of enumerating all the possible paths. The algorithm is based on the k -shortest path algorithm and has a parameter K indicating that K shortest paths will be examined.

The algorithm proceeds as follows:

1. Find the first through K shortest paths from the origin node to the destination node, based on the mean link travel times in the network and store in list P . These are stored in ascending order with respect to travel time in list P .
2. Set $k = 1$; take the k th shortest path from list P and call it p^* . Calculate the expected travel time over p^* using eqns (19) and (26) and denote this value as w^* .
3. If $k > K$: p^* is the 'optimal' path, w^* is the 'minimum' expected travel time. Stop. Otherwise, go to step 4.
4. Set $k = k + 1$, take the k th shortest path from A , and call it p_k . Calculate the expected travel time over p_k by using (19) and (26), and denote this value as w_k . If $w_k < w^*$ then $p^* = p_k$ and $w^* = w_k$. Go to step 3.

There are three issues that need to be addressed before this algorithm can be implemented. The first issue is to identify the technique for finding the K shortest paths. In this paper Shier's k -shortest path algorithm (a label setting algorithm) was adopted based on efficiency considerations (Shier, 1979).



Fig. 8. An acyclic network.

The second issue was to identify the value of K . From a practical point of view the appropriate K value can be based on an empirical sensitivity study. The use of a larger value for K will increase the chances of finding the optimum expected shortest path, but at same time will require greater computational effort. A sensitivity analysis of the solution quality and computation cost vs the K value will be performed in the example problem.

Finally, the proposed heuristic algorithm requires applying eqns (19) and (26) presented in Section 3, which are derived based on the assumptions that the mean and standard deviation of the link travel time are continuous functions of time of day and have at least second order derivatives. A second order polynomial was therefore used to smooth the mean and variance of the link travel times as discussed previously.

6. COMPUTATIONAL ANALYSIS

The objective of this section is to demonstrate the solution quality and computational efficiency of the proposed algorithm with respect to the value of the parameter K used in the algorithm. The heuristic developed in this paper was coded in C++ and executed under the Microsoft Windows operating environment on a 486 compatible with 50 MHz speed and 8 MB RAM.

The experiment was performed on a network from Edmonton, Alberta. This network, composed of 3800 links and 1400 nodes, is used primarily for planning applications. The a.m. peak (6:00a.m.~9:00a.m.) was selected as the study period. Due to a lack of real time data, the dynamic and stochastic travel time patterns in the network were created based on a hypothetical change pattern in travel time during the a.m. peak period at 15-min intervals. The mean link travel times during the analysis period were first generated for each link with maximum travel times around triple of the free flow travel times on the link, as shown in Fig. 9. It was assumed that the coefficient of variation (COV) of the link travel time in the network is a random variable with values uniformly distributed between 0.10 and 0.20. Consequently, the standard deviation of the travel time on a specific link at each time interval could be estimated by multiplying the mean link travel time at that interval with a randomly selected COV for this link. The link travel time data were then represented as a set of discrete means and standard deviations through the a.m. peak period, as shown in Fig. 9.

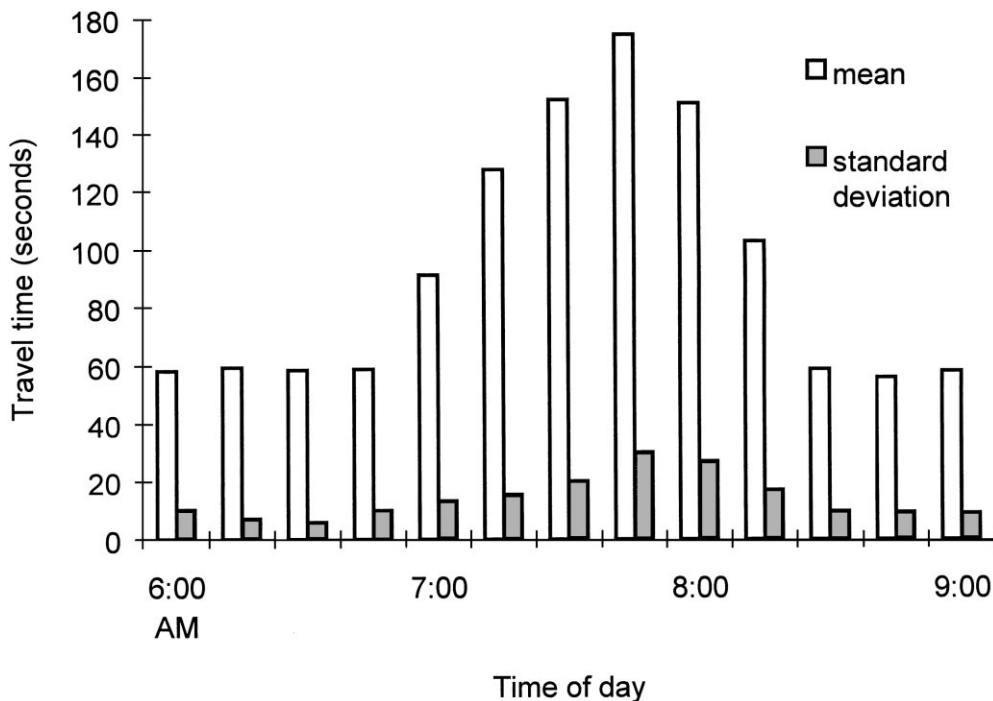


Fig. 9. Link travel time pattern.

To demonstrate the solution quality of the heuristic algorithm a comparison is usually made with the optimal solution. Because it is too computationally intensive to identify the expected optimal path in a large network, a reference path is used for comparison purposes. In this paper the reference path is found by using the proposed algorithm with a pre-specified value of K . A K value of 10 was used indicating that the best path within these ten is the optimal path. Obviously, for a real situation a more detailed computational study should be conducted to identify an appropriate value for K .

Three hundred random O-D pairs with random departure time were generated and their respective expected minimum paths were calculated using the proposed algorithm. Figure 10 shows the relationship between the K value and the percentage of time the optimal path was found. For example, when the K value is equal to one, there is a 30% chance that the minimum

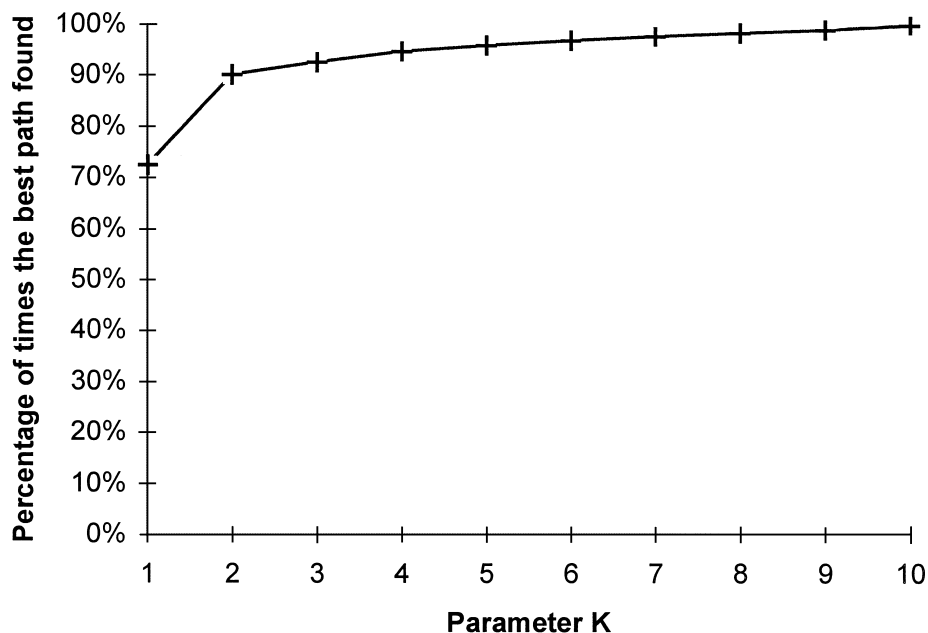


Fig. 10. Percentage correct vs K value.

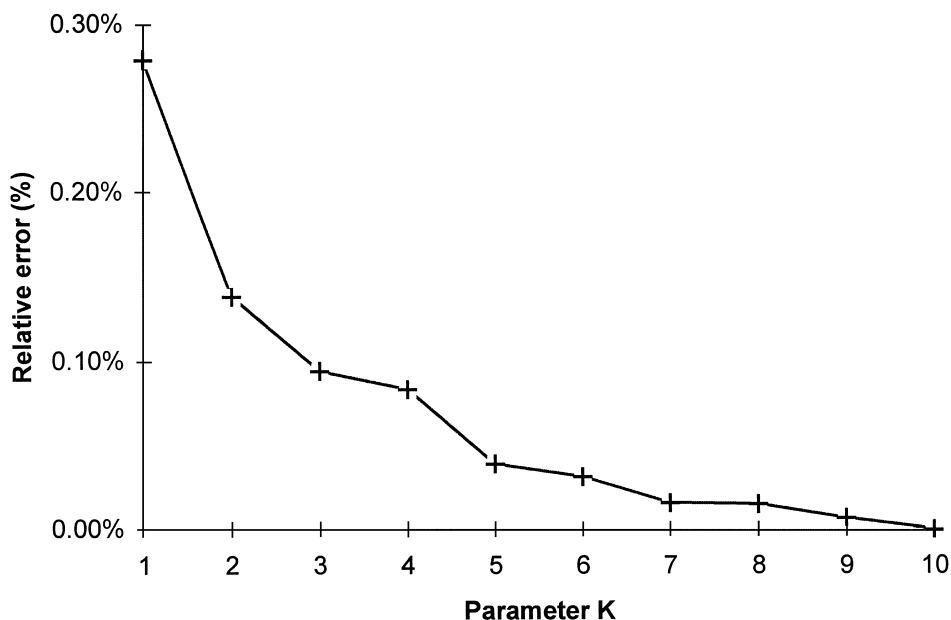


Fig. 11. Relative error vs K value.

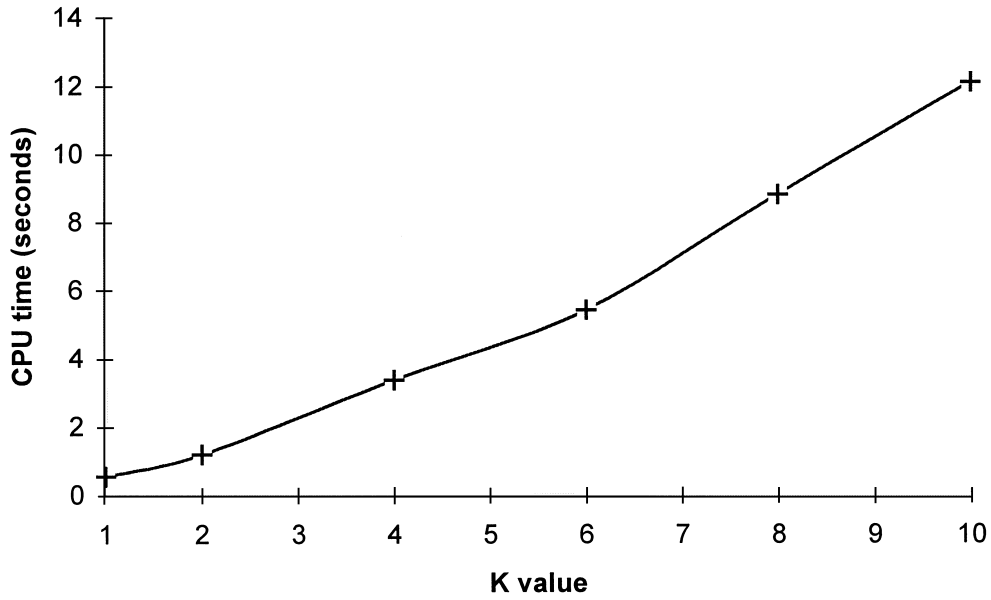


Fig. 12. CPU time vs K value.

path route would not be identified whereas if K is increased to five this percentage decreases to 5%. It should be noted that a K value of one corresponds to the case of simply using a standard minimum path algorithm.

Figure 11 shows the relative error of the solution as a function of the K value. It can be seen that the relative error is very small. For example when K is equal to one, the relative error is very small (less than 0.3%). The average absolute error for a K value of one is approximately 5 s with a maximum error of 120 s with respect to the average travel time 1788 s. It may be seen in Figs 10 and 11 that the greatest jump in accuracy occurs at the lower K values (i.e. more improvement when K increases from one to two than when K increases from nine to ten).

The computation time of the proposed algorithm with respect to the K value is shown in Fig. 12. It can be seen that the increase in CPU time is significant. For example, when K is equal two the increase in CPU time is approximately 90%. However, it should be kept in mind that this algorithm is considerably faster than a complete enumeration.

7. CONCLUDING REMARKS

This paper examined the dynamic and stochastic shortest path problem (DSSPP) of finding the expected shortest paths in a traffic network where the link travel times are modeled as a continuous-time stochastic process. A set of probability-based approximation models was developed to estimate the mean and variance of the travel time of a given path based on the mean and variance of the link travel times. It was shown that the DSSPP is computationally intractable and that it cannot be solved exactly using standard shortest path algorithms. This paper proposed a heuristic algorithm for solving the DSSPP where the dynamic and stochastic attributes of the link travel times are modeled by the mean and variance of the link travel time as a function of time of day. The algorithm is based on k shortest path algorithm and its performance was tested on a realistic network with hypothetical travel times. The following points were illustrated in this paper:

1. The standard shortest path algorithms may fail to find the minimum expected paths in a dynamic and stochastic network. The solution error by the standard shortest path algorithm was shown to be small for the sample problem (approximately 5 s on average). This may result from the fact that the dynamic link travel times used in the sample problem change relatively slowly with respect to time of day.
2. The proposed heuristic algorithm provided improved solutions with an acceptable increase in overall computation time. When the number of paths was increased from one to two the relative error decreased by approximately 18% and the calculation time increased by 90%.

3. As shown in this paper the use of standard shortest path algorithms in dynamic and stochastic traffic networks is incorrect. However, from a practical perspective the standard shortest path algorithms may be acceptable. This will be particularly true if the change of travel time as a function of time of day on the links in the network is moderate.
4. The travel time estimation models developed in this paper provided improved solutions as compared to the traditional methods. The data required by the new models, including the mean and standard deviation of the link travel time as a function of time of day, is information typically available from most ATMS. It can be expected that these models would play an important role in situations where the estimation of the travel time variance is also necessary.

It should be noted that the above conclusions are based on hypothesized link travel time data. It would therefore be necessary to conduct further studies based on real travel time data before any general conclusions may be made. It also would be beneficial to conduct the experiments during incident conditions to identify whether this technique has potential benefits in these situations. Finally, another important direction for further research would be to extend the proposed models and algorithm to take account the correlation of travel times between individual links for a particular point in time.

Acknowledgements—This research was funded through a Natural Science and Engineering Research Council of Canada Operating Grant (OGP 138670) and their assistance is gratefully acknowledged. The authors would like to thank the three anonymous referees for their helpful comments.

REFERENCES

- Bellman, E. (1958) On a routing problem. *Quarterly of Applied Mathematics* **16**, 87–90.
- Chabini, I. (1997) A new algorithm for shortest paths in discrete dynamic networks. *Proceeding of the 8th International Federation of Automatic Control (IFAC) Symposium on Transportation Systems*, Vol. 2, eds M. Papageorgiou and A. Pouliezios, pp. 551–556. Chania, Greece.
- Denardo, E. (1982) *Dynamic Programming: Models and Applications*. Prentice-Hall, Inc, NJ.
- Dijkstra, E. W. (1959) A note on two problems in connection with graphs. *Numerical Mathematics* **1**, 269–271.
- Dreyfus, S. (1969) An appraisal of some shortest path algorithms. *Operations Research* **17**, 395–412.
- Frank, H. (1969) Shortest paths in probability graphs. *Operations Research* **17**, 583–599.
- Fu, L. (1996) Vehicle routing and scheduling in dynamic and stochastic traffic networks. Ph.D. dissertation, University of Alberta, Edmonton, Alberta.
- Gary, R. M. and Johnson, D. S. (1979) *Computers and Intractability: a Guide to the Theory of NP-Completeness*. Freeman, St. Francisco.
- Hall, R. (1986) The fastest path through a network with random time-dependent travel time. *Transportation Science* **20**(3), 182–188.
- Kaufman, E., Lee, J. and Smith, R. L. (1990) Anticipatory Traffic Modeling and Route Guidance in Intelligent Vehicle-Highway Systems. IVHS Technical Report-90-2, University of Michigan.
- Kaufman, E., Lee, J. and Smith, R. L. (1993) Fastest paths in time-dependent networks for intelligent vehicle highway systems application. *IVHS Journal* **1**(1), 1–11.
- Loui, P. (1983) Optimal paths in graphs with stochastic or multidimensional weights. *Comm. ACM* **26**, 670–676.
- Mirchandani, B. P. and Soroush, H. (1986) Routes and flows in stochastic networks. *Advanced Schools on Stochastic in Combinatorial Optimization*, eds G. Angrealtta, F. Mason and P. Serafini, pp. 129–177, CISM. Udine, Italy.
- Mirchandani, B. P. (1976) Shortest distance and reliability of probabilistic networks. *Computer and Operations Research* **3**, 347–676.
- Murthy, I. and Sarkar, S. (1996) A relaxation-based pruning technique for a class of stochastic shortest path problems. *Transportation Science* **30**(3), 220–236.
- Orda, A. and Rom, R. (1990) Shortest-path and minimum-delay algorithms in networks with time-dependent edge-length. *Journal of the ACM* **37**, 607–625.
- Rilett, L. R. (1992) *Modeling of TravTek's Dynamic Route Guidance Logic Using the Integration Model*. Ph.D. dissertation, Queen's University, Kingston, Ontario.
- Ross, M. S. (1989) *Introduction to Probability Models*, 3 edn. Academic Press Inc. San Diego, CA.
- Shier, D. (1979) On algorithms for finding the k shortest paths in a network. *Networks* **9**, 195–214.
- Turner, Shawn M., Brydia Robert E., Liu Jyh C. and Eisele William L. (1997) ITS Data Management System: Year One Activities, Report No. FHWA/TX-98/1752-2. Texas Department of Transportation, Texas Transportation Institute, September, 1997.
- Ziliaskopoulos, A. K. and Mahassani, H. S. (1993) Time-dependent, shortest-path algorithm for real-time intelligent vehicle system applications. *Transportation Research Record* **1408**, TRB, National Research Council, Washington, DC., pp. 94–100.