



TOPIC 16
TRAVEL SUPPLY-DEMAND
MODELLING

DEVELOPMENT OF ARTIFICIAL NEURAL NETWORK MODELS FOR AUTOMATED DETECTION OF FREEWAY INCIDENTS

HUSSEIN DIA

Department of Civil Engineering
Monash University
Clayton Vic 3168
AUSTRALIA

GEOFF ROSE

Department of Civil Engineering
Monash University
Clayton Vic 3168
AUSTRALIA

Abstract

The high cost of congestion caused by incidents has prompted a growing worldwide interest in developing efficient and effective automated incident detection methods. This paper describes the development of new incident detection techniques based on artificial neural networks (ANNs). These models have the potential to provide faster and more fault-tolerant operation.

INTRODUCTION

The high cost of congestion caused by incidents, mainly in terms of traffic delays, air pollution and deteriorated safety conditions, has prompted a growing worldwide interest in developing efficient and effective automated incident detection methods. Incidents are defined as non-recurring events such as accidents, disabled vehicles, spilled loads, maintenance work and other events that disrupt the normal traffic flow and result in a capacity reduction of a facility. Such incidents are believed to constitute about 50-60% of the total delays on US freeways (Lindley, 1987) and this is also expected to increase as facilities become more congested. Therefore, the benefits to be derived from early incident detection and quick response can drastically reduce traffic delays and improve road safety and real-time traffic control. Motorists should be informed by providing real time traveller information to allow for alternate routing of traffic and timely dispatch of emergency services. Intelligent Transport Systems (ITS) technologies are structured to address these needs through Advanced Traffic Management Systems (ATMS) and Advanced Traveller Information Systems (ATIS). For these systems to be effective, it is necessary to develop procedures for detecting incidents which are both reliable and quick to respond.

This paper presents an overview of Automatic Incident Detection (AID) systems and describes some of the most widely used AID algorithms. The framework for a new AID algorithm, based on Artificial Neural Networks (ANNs) is then discussed and initial incident detection results presented.

AUTOMATIC INCIDENT DETECTION

Automated incident detection systems involve two main components: a traffic detection system and an incident detection algorithm. The traffic detection system provides the traffic information necessary for detection while the incident detection algorithm interprets that information and ascertains the presence or absence of incidents. Inductive loop detectors embedded in the freeway pavement are typically used to obtain traffic data, primarily on occupancy and volume. Dual loop installations also provide speed data. These data form the input to an incident detection algorithm which would raise an alarm to indicate the presence of an incident on the facility.

Algorithm types

A number of automated incident detection algorithms have been developed or proposed over the last two decades. Their structure varies in the degree of sophistication, complexity and data requirements but can generally be grouped into main categories as discussed below.

Comparative or pattern comparison algorithms

The core logic of these algorithms is based on comparing observed upstream and/or downstream traffic data, within and between lanes, with pre-established threshold values to declare the occurrence of an incident. These are perhaps the most widely used incident detection algorithms and include the California-type algorithms (Payne and Tignor, 1978; Levin and Krause, 1979), the UK High Occupancy (HIOCC) and Pattern Recognition (PATREG) algorithms (Collins et al. 1983), the ARRB/VicRoads incident detection algorithm (Luk and Sin, 1992; Sin and Snell, 1992) and the Minnesota algorithm (Stephanedes and Chassiakos, 1993). A form of comparative algorithm has also been implemented using data from a wide-area video detection technology known as AUTOSCOPE (Michalopoulos et al. 1993).

McMaster algorithm

The McMaster algorithm is based on a catastrophe theory of traffic flow describing the relationship between speed, flow and occupancy (Persaud and Hall, 1989) and utilises data observed at a single detector station. The volume-occupancy plot of detector data is separated into four areas corresponding to different states of traffic conditions (Gall and Hall, 1989). Incidents are detected after observing specific changes of the traffic state (movement of data points between the four regions in the plot) in a short time period (two intervals). The algorithm is based on a single station detection logic, using 30 second data from the median (fast) lane. The McMaster algorithm has also been combined with a comparative-type algorithm to produce the AUTOSCOPE Incident Detection Algorithm (AIDA) (Michalopoulos et al. 1993). This combined algorithm is used in conjunction with the AUTOSCOPE video detection technology.

Time series algorithms

This general class of algorithms is based on the logic of using the recent history of traffic variables and employing time-series models to provide short-term traffic forecasts. An incident alarm is raised if significant deviations (typically 2 or more standard deviations) between field and forecast values are observed. Three most widely known algorithms in this class are the Standard Normal Deviation algorithm (Dudek and Messer, 1974), the Double Exponential algorithm (Cook and Cleveland, 1974) and the Auto-Regressive Integrated Moving Average (ARIMA) algorithm (Ahmed and Cook, 1982).

Artificial neural networks

Few of the previously developed algorithms have been implemented in practice due to various limitations and varying operational levels in terms of performance criteria such as detection rate, false alarm rate and time-to-detect. Therefore, the need is pressing for more effective real-time incident detection algorithms that maximise detection rate while only generating an acceptable level of false alarms. Furthermore, desired new-generation algorithms should also lend themselves to implementation on new platforms such as parallel computers and must have the required flexibility for the smooth integration with emerging ITS technologies. One promising approach to address these objectives involves the application of Artificial Neural Networks (ANNs). These are also referred to as parallel distributed processing systems or connectionist systems and have been implemented within recent years as a paradigm of computation and knowledge representation.

Neural Networks, as the name implies, are loosely modelled after the biological structure of the brain. A neural network is constructed from a set of inter-connected simple processing elements (PEs). Each PE performs only a few simple computations such as receiving inputs from other PEs and computing an output value which it sends to other PEs. A neural network is inherently parallel in that many PEs can carry out their computations at the same time. The processing ability of the network, stored in the connection strengths or weights, is obtained by a process of adaptation to, or learning from, a set of training patterns. Neuro-computing differs from other branches of computing in that the algorithms are "data-driven". Rather than the computer working through lists of instructions written by a programmer, it deduces the strengths of different relationships by being exposed to a set of examples of the behaviour concerned. By absorbing patterns in the data, the network learns to generalise.

The neural network approach has a number of strengths which lead to it being explored as a likely solution to incident detection. These strengths include: (1) They are well-suited for parallel implementation because they are structured such that only a few steps are performed per PE. This makes them attractive for real-time pattern recognition and classification applications that need to process large amounts of data very fast (Maren et al. 1990). (2) The network itself develops the relationships by recognising and classifying the spatial and temporal patterns in traffic data and this provides greater flexibility compared to more "rigid" modelling frameworks. (3) The network has the capacity to recognise random fluctuations in traffic flow that cause many false alarms. (4) ANNs are highly fault-tolerant in the sense that given an input pattern with noise or disturbance, they would still be capable of recognising that input and providing an acceptable output.

Ritchie and Cheu (1993), demonstrated the feasibility of using ANNs for incident detection. They tested a multi-layer feed-forward (MLF) ANN on a freeway section using simulated traffic detector data. The results confirmed their hypothesis that spatial and temporal traffic patterns could be recognised and classified by ANNs. However, their results were limited in the sense that they trained the ANN models on simulated traffic detector data, used only volume and occupancy data and did not address operational issues such as the impact of detector malfunction and quality of input data on model performance. The work reported here is part of a research program that aims to address many of these unresolved issues.

Performance measures for incident detection algorithms

The performance of an incident detection algorithm is measured by three criteria: detection rate (DR), false alarm rate (FAR) and time-to-detect (TTD). The DR is defined as the number of incidents detected by the algorithm divided by the total number of incidents known to have occurred during the recorded time. The FAR can be defined in different ways depending on whether it is an off-line or on-line FAR. The on-line FAR (FAR_{ON}) is defined as the number of time intervals (typically provided in 20 or 30-second cycles) which gave false alarms divided by the number of time intervals in the entire data set. The off-line FAR (FAR_{OFF}) is defined as the number of incident-free intervals which gave false alarms divided by the total number of incident free intervals. Finally, the TTD is the difference between the time of occurrence and the time at which the incident was declared or an alarm was raised by the algorithm. When an algorithm is being evaluated, however, it is customary to seek the mean time-to-detect (MTTD) a set of (n) incidents. The occurrence time of an incident is usually not known precisely and an estimate has to be deduced from loop detector data or records kept by police, traffic control centres or towing companies.

The above definitions clearly show that both the DR and FAR measure the effectiveness of the algorithm while the MTTD reflects its efficiency. The detection rate and false alarm rates are, unfortunately, positively correlated. In order to detect more incidents, the algorithm thresholds are relaxed which causes some incident-free intervals to be interpreted as alarms. Since many false alarms are caused by random fluctuations in traffic flow, a persistence test is usually performed by testing warnings in a few consecutive intervals before declaring an alarm. This method, in conjunction with increased duration of the persistence test, has been shown to reduce the FAR. However, this was also found to reduce the efficiency of the algorithm since it increased the MTTD considerably. Clearly the three performance measures are all inter-related. The relative importance of the measures, however, is typically DR, FAR and MTTD.

A FRAMEWORK FOR AUTOMATED INCIDENT DETECTION USING ANNS

As implied by their name, ANN models can be visualised as a network. Consider the section of freeway shown in Figure 1(a) which is defined by upstream and downstream detector locations. A corresponding ANN model structure is shown in Figure 1(b). The detector station data form the input to the ANN. The output is a {0,1} variable indicating the absence or presence of an incident in the freeway section, respectively. The parameters of the ANN model are established through a process known as training. In order to train a neural network to perform incident detection, the network must be presented with input detector data and output states for both incident and incident-free conditions. Therefore, the input to the ANN model comprises real-time speed, flow and occupancy measurements in 20-second intervals from each of the upstream and downstream stations. The output of the ANN model is the traffic state within the section. Output State 1 {0} represents incident-free conditions and output State 2 {1} represents incident-conditions. One of the well-known and widely used neural network models is the back-propagation or multi-layer feed-forward (MLF) network. These models are an outgrowth of earlier work on Perceptrons, with the addition of a hidden layer and use of a more robust and capable learning rule (Rumelhart et al. 1986).

The back-propagation algorithm's popularity is due mainly to the solid theoretical foundations on which it rests. The algorithm has been successfully implemented in many pattern recognition applications across many disciplines (Maren et al. 1990). Cheu (1994) tested three ANN architectures suitable for incident detection and real-time classification problems. These included the multi-layer feed-forward (MLF) neural network, the self-organising feature map (SOFM) and the adaptive resonance theory (ART). The MLF, implemented with the back-propagation (BP) training algorithm, proved to be superior to the other architectures tested. The MLF was chosen for implementation in this study based on its earlier success, especially in real-time pattern recognition problems, and based on its demonstrated superior incident detection performance over the other ANN architectures (Cheu, 1994). In particular, the standard three-layer feed-forward neural network has been chosen for this study. It consists of a set of processing elements (PEs) arranged into three layers as shown in Figure 1(b): a layer of (n) input PEs is connected to a layer of (p) "hidden" PEs, which is connected to a layer of (m) output PEs. Each layer comprises at least one processing element.

The detailed structure of the MLF is shown in Figure 1(b). The activity of the input PEs represent the raw information that is fed into the network (the input vector $X = [x_1, x_2, \dots, x_n]$). The activity of each hidden PE (h_j) is determined by the activities of the input vector X and the weights on the connections between the input and hidden PEs (w_{ij}). Similarly, the activity of each output PE (y_k) depends on the activity of the hidden units (h_j) and weights between the hidden and output units (v_{jk}). A typical unit (k) in the output layer determines its activity by following a two step procedure. First, it computes the total weighted input, N_k , using the formula in equation [1]:

$$N_k = \sum_{j=1}^p v_{jk} h_j \quad \forall k = 1, \dots, m \quad (1)$$

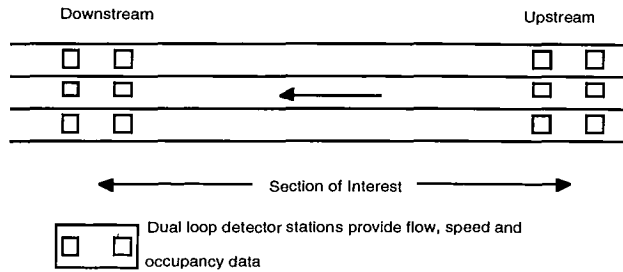
Where h_j is the activity level of the j^{th} unit in the hidden layer and v_{jk} is the weight of the connection between the j^{th} and k^{th} PEs.

Second, the unit calculates its activity y_k using some function of the total weighted input. To obtain the activity level of unit k , a threshold value θ_k is subtracted from the weighted input and the net input is then fed into a transfer function. Typically, the sigmoid function is used:

$$y_k = \frac{1}{1 + e^{-(N_k - \theta_k)}} \quad (2)$$

In order for a neural network to perform some actual task, it must undergo a training process during which the weights on inter-connections (w_{ij}, v_{jk}) and the thresholds associated with the PEs (θ_j, θ_k) are determined. This process begins by assigning random initial values to all the connection weights. Then, each example from the training set is presented to the network and the output vector produced by the network is compared with the desired results. The error between the actual and desired outputs is computed. By applying a learning rule, usually some form of the Generalised Delta Rule (Rumelhart et al. 1986), the inter-connection weights and other network parameters are adjusted in such a way that the error between the desired and actual outputs is reduced. This is achieved by implementing a gradient descent on the error curve of the network's output.

(a) Physical System



(b) ANN Model

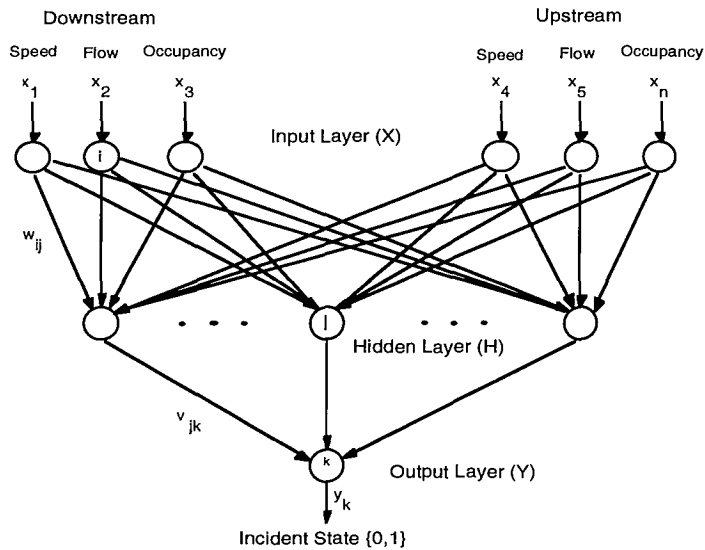


Figure 1 ANN modelling framework

DATA FOR ANN DEVELOPMENT

In order to train a neural network to perform incident detection, the network must be presented with examples of input detector data (speed, flow and occupancy) and output states for both incident and incident-free conditions. Therefore, the data required should at least have a description of the state of traffic along the freeway in addition to detector data comprising traffic flow measurements at regular time intervals for each detector station. In contrast to previous research which has relied on simulated data for model development (Ritchie and Cheu, 1993), this study relies on real data.

Data collection

The data were collected from 15 inductive dual-loop stations located within a site 8.5 km in length on the Tullamarine Freeway in Melbourne, Australia. Dual inductive-loop detectors are placed in every lane of the Freeway in both the inbound and outbound directions spaced at distances between 450 and 1070 meters. The data for the study were assembled from two data sets held at the VicRoads Traffic Control and Communications Centre (TCCC). The first data set contained information logged by the operators at the TCCC regarding the incidents that occurred on the Tullamarine Freeway. However, when operators are busy managing incidents, it is not uncommon for important details to be left out from the records. This presented some difficulties when examining the incidents since in many cases the location of the incident or its direction remained unknown from the record. All the incidents that were logged for the selected segments of the Tullamarine Freeway were extracted from this database. These included some 385 incidents for the period between January 1992 and March 1994. The second data set contained detector station data comprising speed, flow and occupancy measurements in 20-second cycles. Each data file obtained comprises the detector measurements for the whole freeway for that specific day. Data files comprising incident-free days were also obtained.

Each of the 385 incidents was then examined individually. Out of the 385 incidents recorded by the operators in the log, only 120 incidents could be confirmed. The rest either occurred outside the 8.5 km segment of the freeway or during low-volume conditions and therefore had no effect on traffic conditions. Others could not be confirmed due to missing information or data or due to faulty detectors. One section in the outbound direction of the Tullamarine Freeway was selected for initial modelling. Figure 2 shows the selected section, between stations S5 and S6, along with related upstream and downstream stations. A total of 12 detectable incidents occurred inside this section, which also had the maximum detector-station separation of 1070 meters. Furthermore, in order to minimise false alarms within the test section, the model should also be trained on incidents that occurred outside the test section. A total of 11 incidents that occurred within the immediate two upstream sections and two incidents that occurred within the immediate downstream section were also included in the database. This resulted in a set of 25 detectable incidents between stations S3 and S7 in Figure 2.

Assignment of desired output states

The input to the ANN model comprises speed, flow and occupancy measurements in 20-second intervals from the upstream station (S6) and the downstream station (S5) as shown in Figure 2. The output of the ANN model is the traffic state within the section (State 1 for incident-free conditions and State 2 for incident-conditions). These desired output states must be provided for each input vector in the data, ie. every 20 seconds. It is therefore essential that the times defining the start and end of incidents are determined such that the correct output states are assigned to the data. As noted earlier, previous studies have used simulated data for ANN model development (Ritchie and Cheu, 1993). In that case, the incident start/end times are known precisely. However, when "real-world" data are being used, the incident start and end times are rarely, if ever, known precisely. In this study, estimates of these times were compiled from the operator's log. These times, however, were reported as the times when the operator detected (or confirmed) the occurrence/clearance of incidents and not the times when the incidents actually occurred or ended. Therefore, it is still necessary to determine the specific 20-second interval that represents the start/end of an incident. The procedure used to arrive at these times is described below. All the input data between the start and end of incidents were assigned desired outputs {1} while the rest of the data were assigned desired outputs {0}.

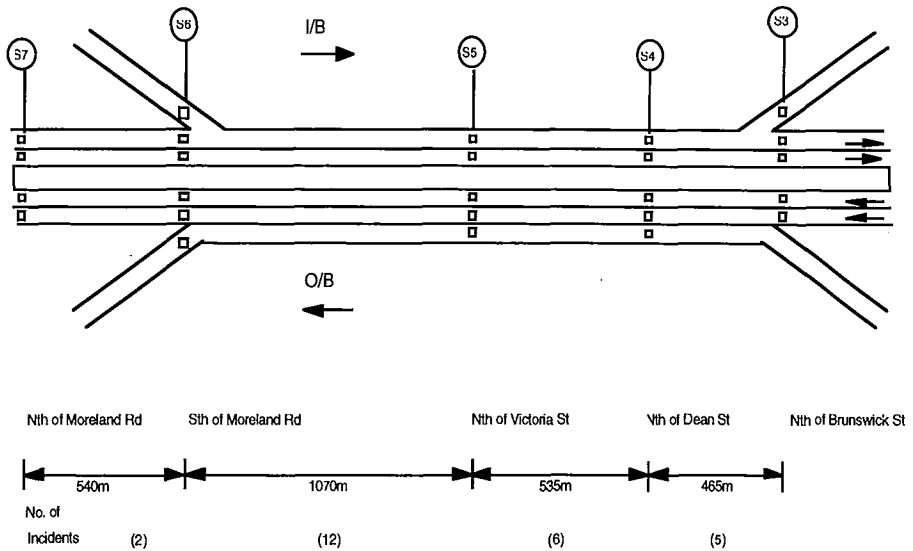


Figure 2 Section S3-S7 of Tullamarine Freeway

Incident start times

For the determination of the first 20-second interval representing the start of an incident, the first three 20-second consecutive intervals (one minute) of non-zero, non-increasing speed and flow and non-decreasing occupancy were identified. If the difference for each parameter between the first and last interval values is greater than a certain threshold (described below), the start of the first 20-second interval is assumed to represent the start of the incident (this will be referred to as the identification criteria). Since there are three parameters to investigate, this method provides three time values for the start of the incident. In many cases, these times do not coincide with one another which highlights that the criteria has not been met for the three parameters simultaneously. Therefore, if the identification criteria is met for the three parameters within the same 20-second interval, the start of that interval is assumed to represent the start of the incident. Otherwise, the earliest of the three times is used.

Care must be taken in the selection of the threshold value such that the random variations in the data are not mistaken for incident conditions on one-hand and such that the data is not pre-conditioned in a manner that could influence the training of the neural network. It was observed that random fluctuations do not last as long as incidents and generally result in variations not exceeding 20%, while incident conditions usually last longer and result in greater deterioration of the parameters (normally exceeding 25%). Therefore, incident times were generated for the 25 incidents for thresholds between 20% and 90%. It was found that lower threshold values resulted in earlier incident times and that the identification criteria was often not met at all for many incidents when higher thresholds were used. This is due to the fact that larger thresholds are only met when the conditions have deteriorated significantly. The use of high thresholds may therefore impede the detection of less severe incidents that only result in a moderate deterioration of conditions. As a result, the identification criteria was implemented using a threshold of 20%.

As for the number of consecutive intervals over which the parameters were observed, the choice of one-minute (three 20-second consecutive intervals) was arrived at after careful examination of the detector data. It was observed that the data (as provided in 20-second cycles), exhibited random fluctuations for 20-40 seconds. Therefore, the choice of two consecutive intervals is inappropriate since these are inherent in the data. Similarly, and due to the random nature of the

data, the use of four intervals (80 seconds) was also deemed inappropriate. This is attributed to the difficulty in locating such four non-zero consecutive intervals where the criteria for the parameters would be met and even then, the incident start times generated were not earlier than those with three 20-second intervals. Therefore, the identification criteria was implemented using three 20-second intervals.

Incident end times

For the determination of the first 20-second interval representing the end of an incident, the first three 20-second consecutive intervals (one minute) of non-zero, non-decreasing speed and flow and non-increasing occupancy are identified. If the identification criteria is met for the three parameters within the same 20-second interval, the end of that interval is assumed to represent the end of incident. Otherwise, the earliest of the three time values is used to signal the end of incident. The choice of threshold and consecutive intervals was determined in a similar fashion as for the incident start times described previously. Out of the three parameters, speed was found to be the limiting factor in the determination of the incident start times (80% of cases) and incident end times (96% of cases).

Creation of training and training test data sets

The next activity involved compiling the training and training test sets. The training set is used for determining the network parameters. The training test set is used to prevent the network from learning the specific patterns in the training set and thereby enables the ANN model to generalise better. The training test set issue will be discussed in more detail in later sections.

The ANN models should be trained on a set of incidents that are representative of the population to which the network will ultimately be applied. The same applies to the training test set. Training an ANN model with a wide range of incidents that include different patterns (location, severity and duration) under a variety of flow conditions (heavy, moderate and light) helps improve its robustness in detecting incidents under varying conditions. Therefore, the 25 incidents were stratified accordingly and two sets were selected randomly into the training and training test sets. The next sections describe the procedures adopted to randomly select 13 and 12 representative incidents for the training and training test sets respectively.

Incident severity

Three incidents that occurred in Section S5-S6 (Figure 2) did not affect the fast lane. One was placed in the training set and the other two in the training test set. All other incidents resulted in all lanes being affected in and upstream of the respective sections where those incidents occurred.

Incident duration

The incidents were stratified into four categories of 30-minute duration. The distribution of incidents in the training and training test sets according to duration is as shown in Tables 1 and 2, respectively.

Table 1 The distribution of incidents in the training set according to incident duration

	Total Incidents	t ≤ 30	Duration (minutes)	t > 90
S4-S5	3	1	1	0
S5-S6	6	1	3	1
S6-S7	1	0	1	0
Total	13	2	7	1

Table 2 The distribution of incidents in the training test set according to incident duration

Section	Total Incidents	$t \leq 30$	Duration	(minutes)	$t > 90$
S3-S4	2	0	1	1	0
S4-S5	3	0	1	2	0
S5-S6	6	3	2	1	0
S6-S7	1	0	0	0	1
Total	12	3	4	4	1

Flow conditions

The flow conditions (maximum vehicle per hour per lane (vphpl) calculated at the immediate upstream station from the incident during the last 15 minutes prior to the start of incidents) were used to stratify the incidents into one of three categories: incidents occurring during heavy, moderate or light flow conditions. The criteria for levels of service on basic freeway sections (AUSTROADS, 1988) were used as a guide in stratifying the incidents. Flows exceeding 1550 vphpl were considered heavy flows, while those less than 700 vphpl were considered light flows. The distribution of incidents in the training and training test sets according to pre-incident flow conditions is as shown in Table 3 and Table 4, respectively.

Table 3 The distribution of incidents in the training set according to flow conditions

Section	Total Incidents	Flow Conditions		
		Light	Moderate	Heavy
S3-S4	3	0	1	2
S4-S5	3	0	1	2
S5-S6	6	0	2	4
S6-S7	1	0	1	0
Total	13	0	5	8

Table 4 The distribution of incidents in the training test set according to flow conditions

Section	Total Incidents	Flow Conditions		
		Light	Moderate	Heavy
S3-S4	2	0	1	1
S4-S5	3	1	0	2
S5-S6	6	0	2	4
S6-S7	1	0	1	0
Total	12	1	4	7

TRAINING OF THE MLF NEURAL NETWORK

The 13 incidents in the training set had a total duration of about 363 minutes. However, at least 15 minutes prior to incident occurrence and after incident clearance were also included in the training. This is necessary to ensure that enough time was given for traffic conditions to stabilise prior to and after incident occurrence and clearance. This resulted in the training and training test sets comprising 1710 and 1970 minutes, respectively. Table 5 shows the distribution of State 1 (incident conditions) and State 2 (incident-free conditions) vectors in the two sets. Each vector includes the traffic input (speed, flow and occupancy) along with the incident state variable (State 1 or State 2).

Table 5 Incident and incident-free intervals in the training and training test sets

Data Set	Total Number of	Incident-free Intervals	Incident Intervals
	20-Second Intervals	(State 1)	(State 2)
Training Set	5131	4042	1089
Training Test Set	5910	5219	691

Input features

As this study is concerned with the application of ANN models to incident detection, the features to be investigated are therefore related to both the modelling tool (the ANNs) and incident detection parameters. The issues related to incident detection parameters are discussed first.

Incident detection parameters

One of the main issues in incident detection modelling is the selection of input features. The choice of traffic flow variables, detection logic and other related parameters is a function of the desired complexity of the model and the surveillance technology used. Some of the issues related to input features are discussed in the following sections.

Single or dual stations

This issue is related to whether inputs from one or two stations will be required to identify incidents within a section. It is generally accepted that incidents within a section will result in an increase in occupancy upstream and a decrease in flow and occupancy downstream compared to pre-incident conditions. However, some automatic incident detection (AID) systems are known to utilise only inputs from a single station which presumably enhances its generalisation potential. Single and dual-station models will be investigated in this study.

Time intervals

This issue is related to the number of 20-second input values that are needed for each decision regarding the presence or absence of incidents at any time interval (t). Incidents do not result in an instantaneous deterioration of conditions and therefore some time must be allowed before the effects of an incident are detected at a station. Inputs from the current time interval (t) and the previous four intervals (t-1, t-2, t-3, t-4) for each station will be investigated in this study.

Station input

Some of the available AID systems utilise station averages (eg. California algorithms) while others use input from the fast lane (eg. McMaster) or all lanes (ARRB-VicRoads model). The use of station averages or fast lane measurements obviously has the advantage of limiting the number of inputs to the ANN and thus decreasing the model's complexity.

As noted earlier, three of the 25 incidents in the data only affected the slow and middle lanes. The feasibility of using fast lane (or averages across station) data to detect these incidents needs to be investigated. Therefore, station averages, fast lane and all lanes scenarios will be investigated in this study. Table 6 summarises the incident detection parameters. A total of 45 (3x5x3) models will therefore need to be investigated. Only an equal number of time intervals are considered for both the upstream and downstream stations.

Table 6 Incident detection parameters

Stations	Time Intervals	Station Input
Dual Upstream & Downstream	t	Station Average
Only Upstream	t, t-1	Fast Lane
Only Downstream	t, t-1, t-2	All Lanes
	t, t-1, t-2, t-3	
	t, t-1, t-2, t-3, t-4	

ANN features/parameters

The next step after arriving at the structure of the 45 model types was the selection of an appropriate set of ANN features and parameters to use in the appraisal of these models. At this stage, it was necessary to determine the basic features of models (Maren et al. 1990). This included the selection of a training method (supervised/unsupervised), a network model (Back-propagation(BP), Adaptive Resonance Theory (ART), Self Organising Map (SOM) etc.) and a learning rule (delta, cumulative delta, backprop etc.). The following ANN features were used consistently throughout the designed experiments. These features were arrived at after considerable investigation and based on the demonstrated performance of these features for pattern recognition problems in general and for incident detection in particular (Maren et al. 1990; Cheu, 1994).

- Training Method : Supervised
- ANN Model : Logicon Projection Network
- Learning Rule : QuickProp
- Transfer function : Sigmoid
- Output Ranges : 0.2-0.8 (instead of 0-1)
- Objective Function : Classification Rate (the average of the correctly classified incident and non-incident states)

Furthermore, to evaluate a neural network method in a useful way, the conditions that might be relevant to its performance should be varied systematically. These conditions include the number of inputs, training cases and hidden units, as well as the amount of noise, presence of irrelevant inputs and initial weights. It is generally accepted that designing the least complicated network provides good results. This is found to be true especially if the training data are noisy as a complicated network may learn odd relationships between the input and output based on the noise and not the data. The number of nodes in the hidden layer depends primarily on the size and nature of the training data. In general, the network performance can be enhanced by adding more nodes to the hidden layer. However, this only applies up to a certain point beyond which performance would start to deteriorate. More hidden nodes also allow the network to generalise better although this is achieved at the expense of the number of training cycles. It is therefore difficult to determine in advance the number of nodes in the hidden layer. One has to experiment with different nodes and choose the best performing model. Each of the 45 proposed models was therefore trained under a varying number of hidden units. This resulted in the design and training of some 500 models.

It is appropriate at this stage to discuss the reasons for the selection of the Logicon Projection Network for the development of the AID models. The basic motivation behind the development of the Logicon Projection Network in the first place was the desire to build a faster and more streamlined network by combining the positive features of closed and open boundary networks. Closed boundary networks, eg. Adaptive Resonance Theory Networks (ART), are fast learning because they properly initialise the network weights and thresholds to prototypes of the training set. On the other hand, open boundary networks, eg. Back-propagation, minimise the output error through gradient descent. Combining these two features results in faster training times. This was favourable since it meant that the number of training cycles needed to train each of the 500 models need not be as large as for the standard back-propagation network. Furthermore, the Logicon

Projection Network has an accurate means of initialisation which reduces the possibility of getting trapped in a local minimum.

The outcome of the investigation of these 500 models will be the determination of the optimal input features. Once these input features are determined, the selected model(s) will undergo further training and modifications to optimise network performance.

Training strategy

After a model is designed, it is trained on the training set for 513100 iterations (this is equivalent to 100 cycles since the training set had 5131 vectors which meant that each vector was presented to the ANN 100 times). As was mentioned earlier, the number of training cycles required for the Logicon Projection Network are generally lower than those for other Back-propagation networks. The choice of 100 cycles was the result of considerable investigation in which it was found that this number of cycles was sufficient for the network to learn the general patterns needed to produce the correct classification. Once the optimal input features are determined, the selected model(s) will undergo further training in which the effect of training cycles, initial weights, learning rules and transfer functions will be tested and evaluated such that the network performance is optimised.

During the training process and every 1000 iterations, the trained model is tested on the training test set. If the classification rate on the training test set improves, the model is saved. Otherwise training continues. If the classification rate does not improve for any consecutive 100 tests, the training is stopped and the last model saved constitutes the best model for the given input features and ANN parameters. This training strategy, in conjunction with limiting the number of hidden units to the point that the network does not have the capacity to learn the specific patterns of the individual samples, was adopted to prevent overtraining or overfitting. This usually occurs when the model adjusts its parameters in such a way to memorise the training set and thus loses generality in classifying the training test set.

RESULTS OF TRAINING

The traffic flow input parameters every 20 seconds are to be classified into one of two classes or states (State 1 or State 2). One measure of network performance is therefore the classification rate which is best illustrated using a classification rate matrix as shown below. When classifying the incoming 20-second data, two types of error may be committed. The first is a Type I Error, where the model concludes that incident conditions are present when in fact they are not. This error basically represents the false alarm rate. The second type of error is a Type II Error, where the network concludes that incident conditions are not present, when in fact they are. The classification rate is typically a function of the average of the correctly classified states. Assuming no persistence checks are applied, an incident is detected when the traffic state changes from State 1 {0} to State 2 {1}.

		Desired Output	
		0	1
ANN Output	1	Type I Error (FAR)	% Correctly Classified Incident Conditions
	0	% Correctly Classified Incident-free Conditions	Type II Error

The classification rate matrix for a sample model comprising average measurements from dual stations is shown below. The correct classification of incident-free conditions is 95.7% while the correct classification of incident conditions is 79%, resulting in an average classification rate of 87.4% on the training set. The corresponding Sum of Squared Error (SSE) between the ANN output and desired output values is 291.8. The 79% classification rate of incident conditions does not correspond to the model incident detection rate, since an incident is successfully detected when the traffic state changes from State 1 to State 2 as described previously. Furthermore, assuming no persistence checks are applied, the false alarm rate is taken directly from the classification rate matrix as 4.26 %.

$$\begin{matrix} & \begin{matrix} 0 & 1 \end{matrix} \\ \begin{matrix} 1 \\ 0 \end{matrix} & \begin{bmatrix} 0.0426 & 0.7906 \\ 0.9574 & 0.2094 \end{bmatrix} \end{matrix}$$

The incident detection measures corresponding to the above classification rate matrix are summarised in Table 7 below. These results clearly indicate the positive correlation between detection rate (DR) and false alarm rate (FAR). Better detection rates are achieved at the expense of higher false alarm rates. Furthermore, these sample results also show the benefits to be derived from the application of persistence tests. The application of a one interval persistence test results in a 67% reduction in false alarm rate at the expense of a 33% reduction in detection rate and a 19% increase in the mean time-to-detect.

The classification rate results and incident detection performance measures were obtained for each of the trained 500 models. These results are currently undergoing analysis using ANOVA techniques to determine whether there are any significant differences between the various factors (Table 6). This is needed to determine the optimal network architecture and the optimal input parameters, in terms of the above mentioned evaluation criteria.

Table 7 Sample incident detection results

Persistence Test	Detection Rate (DR) (%) (DI/TI)	False Alarm Rate (FAR) (%)	Mean Time-to-detect (MTTD) (Minutes)
0	100 (6/6)	4.26	1.83
1	67 (4/6)	1.39	2.17
2	50 (3/6)	0.49	2.22
3	50 (3/6)	0.37	2.55

Key: DI=Detected Incident, TI=Total Incidents which Occurred in Section S5-S6

CONCLUSIONS AND RESEARCH DIRECTIONS

Incident detection algorithms have an important role to play in freeway incident management. In response to the challenge to develop detection algorithms which are both reliable and quick to respond, recent research interest has focused on the potential of ANN models for incident detection. Previous studies have relied on simulated data for ANN training. In contrast, the results reported here demonstrate that "real world" data can be used to train ANN incident detection models.

Using the results from a number of model runs, rigorous statistical analyses are currently being used to determine the optimal network architecture and input parameters. Once these are determined, the selected network will undergo further investigation, where it will be trained under

a variety of conditions (different initial weights, varying number of training cases, added noise etc.) such that its performance is optimised.

A new validation test set of independent incidents are being compiled for validating the optimal network and investigating the generalisation of the results to other sections of the freeway. The ANN models' performance will also be compared to the algorithm currently implemented by VicRoads on the Tullamarine Freeway.

ACKNOWLEDGMENTS

The authors acknowledge the support of VicRoads for providing the data. In particular, we acknowledge the generous assistance of Mr. Francis Sin and Mr. Anthony Snell.

REFERENCES

- AUSTROADS (1988) *Guide to Traffic Engineering Practice, Roadway Capacity*. AUSTROADS, Sydney.
- Ahmed, S.R. and A.R. Cook (1982) Application of time-series analysis techniques to freeway incident detection. *Transportation Research Record*, 841, 19-21.
- Cheu, R.L. (1994) *Neural Network Models for Automated Detection of Lane-Blocking Incidents on Freeways*. Ph.D. Dissertation, University of California, Irvine.
- Collins, J.F. (1983) Automatic incident detection - experience with TRRL algorithm HIOCC. *TRRL Supplementary Report 775*, Transport and Road Research Laboratory, Crowthorne, Berkshire, U.K.
- Cook, A.R. and D.E. Cleveland (1974) Detection of freeway capacity-reducing incidents by traffic-stream measurements. *Transportation Research Record*, 495, 1-11.
- Dudek, C.L. and C.J. Messer (1974) Incident detection on urban freeways. *Transportation Research Record*, 495, 12-24.
- Gall, A.I. and F.L. Hall (1989) Distinguishing between incident congestion and recurrent congestion: a proposed logic. *Transportation Research Record*, 1232, 1-8.
- Levin, M. and G.M. Krause (1979) Incident detection algorithms. part 1: off-line evaluation; part 2: on-Line evaluation. *Transportation Research Record*, 722, 49-64.
- Lindley, J.A. (1987) Urban freeway congestion: quantification of the problem and effectiveness of potential solutions. *Institute of Transportation Engineers Journal*, 57(1), 27-32.
- Luk, J.Y.K. and F.Y.C. Sin (1992) The calibration of freeway incident detection algorithms. *Working Document No. WD-TE 92/001*, Australian Road Research Board, Nunawading, Victoria, Australia.
- Maren, A.J., C.T. Harston and R.M. Pap (1990) *Handbook of Neural Computing Applications*. Academic Press Inc., San Diego.
- Michalopoulos, P.G., R.D. Jacobson, C.A. Anderson and B. DeBruycker (1993) Automatic incident detection through video image processing. *Traffic Engineering and Control*, 34(2), 66-75.
- Payne, H.J. and S.C. Tignor (1978) Freeway incident detection algorithms based on decision trees with states. *Transportation Research Record*, 682, 30-37.
- Persaud, B.N. and F.L. Hall (1989) Catastrophe theory and patterns in 30-second freeway traffic data - implications for incident detection. *Transportation Research*, 23A(2), 103-113.
- Rumelhart, D.E., G.E. Hinton and R.J. Williams (1986) Learning Internal Representations by Error Propagation. In: *Parallel and Distributed Processing*. (D.E. Rumelhart, J.L. McClelland and the PDP Research Group, eds.), Vol. 1, pp. 318-362, MIT Press, Boston.

TOPIC 16

TRAVEL SUPPLY-DEMAND MODELLING

Ritchie, S.G., and R.L. Cheu (1993) Simulation of freeway incident detection using artificial neural networks. *Transportation Research*, 1C(3), 203-217.

Sin, F. and A. Snell (1992) Implementation of Automatic Incident Detection Systems on the Inner Metropolitan Freeways in Melbourne. In: *Proceedings of the Seventh Road Engineering Association of Asia and Australasia (REAAA) Conference*, Vol. 1, pp. 337-346.

Stephanedes, Y.J. and A.P. Chassiakos (1993) Freeway incident detection through filtering. *Transportation Research*, 1C(3), 219-233.