

行政院國家科學委員會補助專題研究計畫 成果報告
 期中進度報告

基因及螞蟻規則探勘模式-以事故分析及事故鑑定為例 (II/III)
Developing Genetic and Ant-based Rule Mining Models- Case
Studies on Accident Analysis and Appraisal (II/III)

計畫類別： 個別型計畫 整合型計畫

計畫編號：NSC 97-2628-E-009-035-MY3

執行期間：97年8月1日至100年7月31日

計畫主持人：邱裕鈞 交通大學交通運輸研究所 副教授

計畫參與人員：傅強 交通大學交研所博士班研究生

謝志偉、邱煜勝、鍾佩儒、高筑韻、蕭任谷、魏嘉儀 交通大
學交研所碩士班研究生

成果報告類型(依經費核定清單規定繳交)： 精簡報告 完整報告

本成果報告包括以下應繳交之附件：

- 赴國外出差或研習心得報告一份
- 赴大陸地區出差或研習心得報告一份
- 出席國際學術會議心得報告及發表之論文各一份
- 國際合作研究計畫國外研究報告書一份

處理方式：除產學合作研究計畫、提升產業技術及人才培育研究計畫、列
管計畫及下列情形者外，得立即公開查詢

涉及專利或其他智慧財產權， 一年 二年後可公開查詢

執行單位：國立交通大學

中華民國 99 年 5 月 25 日

一、摘要

1.1 中文摘要

傳統以個體角度進行事故分析之方法，例如，判別分析 (discrimination analysis)、羅吉斯迴歸 (logistic regression)、次序普羅比 (ordered probit)、羅吉特 (logit) 及混合羅吉特 (mixed probit) 等模式，大多僅能探討單一危險因素之影響程度。事實上，事故嚴重與否大多係由多項因素同時發生所導致。此一綜合多項因素之危險情況，在統計分析上，甚難加以窮舉分析。基此，本計畫乃於第一年期提出基因規則探勘模式 (Genetic rule mining, GRM)，可由探勘所得之規則的前半部，判斷何謂危險情況，進而加以避免。惟本研究所提出之 GRM 必須先固定規則數量，再同時進行最佳規則組合之尋優。因此，具有染色體長度過長，尋優效果不佳，以及探勘過多衝突或重覆規則的傾向，進而導致規則難以詮釋，無法提出具體之安全改善策略。

有鑑於此，本計畫第二年期乃提出改良式的基因規則探勘模式 (Genetic rule mining, GRM)，稱為逐步基因規則探勘模式 (Stepwise GRM, SGRM)。SGRM 一次僅挑選使事故嚴重度預測率精確率最高的一條規則，再以此規則為基礎，進行下一條規則之選取，直到精確率無法再改善為止。如此，即可避免選擇規則過多，且相互重覆或矛盾的問題。此外，由於不同類型事故之影響因素與危險情況不一定相同，因此，有必要加以區隔分析。本年度以先以總計 5563 件單車事故 (single vehicle accident) 為分析基礎。結果顯示，本模式共選擇了 38 條規則，其訓練準確度達 75.1%，而驗證準確度則達 73.8% 均遠高於決策樹之預測結果。而影響事故嚴重度之危險情況也加以確認，並研提改善策略。

關鍵字：事故分析、逐步基因規則探勘、事故嚴重度、決策樹。

1.2 Abstract

Conventional individual approach to conduct accident analysis is to associate the crash severity with driver, vehicle and roadway factors by using discrimination analysis, logistic regression, ordered probit, logit and mixed logit models. Although statistic models are the commonly used methods in the context of crash data analysis, most of them have their own assumptions and complexity in the model estimation and interpretation. Once the assumptions were violated, the model could lead to erroneous estimation results, especially for the individual approach wherein most variables explaining the individual crashes are categorical. It is difficult to develop parametric statistical models based upon the categorical data. In addition, most of statistical methods only provide calibrated parameters with significance tests, which are then used to examine the effects of the corresponding variables on crash counts or crash severity. The interrelationship among explanatory factors cannot be examined in details. According to "error chain theory" a crash is often caused by a series of errors, not solely by a single factor. As such, mining the explanatory rules is deemed necessary for crash data analysis. To this end, the first research year of this project has proposed genetic rule mining models to discover the key rules (i.e. risky conditions). However, since the proposed GRM models simultaneously select the rule combinations under a given upper limit of rule number and tend to mine too many conflict or redundant rules, making the rule interpretation difficult.

Based on this, the second year of this project further propose a stepwise GRM (SGRM) model, which select the optimal one rule at a time and iteratively proceed to select the next best rule based on the selected rules until model performance (accuracy) can't not improved. Since the risky conditions and contributory factors of various types of crashes will significantly vary, the analysis is conducted on each type of accidents separately. Taking single-vehicle accident for instance, a total of 5,563 crashes on Taiwan's freeway network from 2003 to 2007 are collected, where numbers of A1 (fatal crash), A2 (injury crash), and A3 (property damage only crash) are 226, 1,593, and 3,744,

respectively - an uneven distribution commonly seen in the context of crash analysis. A total of 38 rules have been mined which can achieve overall correct rates of 75.1% in training and of 73.8% in validation, respectively, much higher than those yield by the decision tree model. Risky conditions along with their corresponding improvement strategies have been identified.

Key Words: *Crash analysis, stepwise genetic rule mining, crash severity, decision tree analysis.*

二、主要研究成果

2.1 Introduction

Crash data analysis can be carried out by two main approaches: collective approach and individual approach (Abdel-Aty and Pande, 2007). The collective approach is characterized by crash frequency modeling. Frequency of crashes is aggregated over specific time periods (months or years) and locations (segments or intersections). Most of these studies attempt to explore the relationship between crash counts and explanatory variables, such as roadway geometry, traffic control facilities, traffic conditions, and so on by using Poisson or Negative Binomial regression models (e.g. Poch and Mannering, 1996; Milton and Mannering, 1998; Ivan *et al.*, 1999; Abdel-Aty and Radwan, 2000; Greibe, 2003; Abdel-Aty and Pande, 2007; Wong *et al.*, 2007). For the collective approach, however, individual contributing factors to the crash (e.g., driver demographics, driver behaviors, vehicle types) are not considered and factors affecting the crash severity cannot be identified either. Therefore, some studies employed individual approach to crash data analysis. The individual approach is characterized by each individual crash case. The main focus of these studies was to associate the crash severity with driver, vehicle and roadway factors by using ordered probit/logit model or logistic regression (e.g., Shanker and Mannering, 1996; Dissanayake *et al.*, 2002; Al-Ghamdi, 2002; Delen, *et al.*, 2002; Tay and Rifaat, 2007; Sze and Wong, 2007). More advanced logit-based approaches, such as nested logit model or mixed logit model, were also employed to analyze the same issue (e.g. Shanker, *et al.*, 1996; Chang and Mannering, 1999; Milton, *et al.*, 2008).

Although statistic models are the commonly used methods in the context of crash data analysis either collectively or individually, most of them have their own assumptions and complexity in the model estimation and interpretation. Once the assumptions were violated, the model could lead to erroneous estimation results, especially for the individual approach wherein most variables explaining the individual crashes are categorical (e.g., driver gender, road type, lighting condition, violation, weather condition, and severity degree, among others). It is difficult to develop parametric statistical models based upon the categorical data. Therefore, a number of distribution-free methods, such as decision tree (Chang and Chen, 2005; Chang and Wang, 2006) and artificial neural network (Chiou, 2006; Delen *et al.*, 2006), were adopted to deal with the classification and prediction problems. However, two gaps still remain. First, the interpretations of classification results with such methods are weak. The knowledge lying in the crash data cannot be fully discovered, because artificial neural network is in essence a black box and the prediction error of decision tree is usually high. Second, most of statistical methods only provide calibrated parameters with significance tests, which are then used to examine the effects of the corresponding variables on crash counts or crash severity. The interrelationship among explanatory factors cannot be examined in details. According to “error chain theory,” a crash is often caused by a series of errors, not solely by a single factor. As such, mining the explanatory rules is deemed necessary for crash data analysis. It is shown in Figure 1 that limited information could be mined from the influence of single variable on crash severity. In contrast, combination of multiple variables would reveal explicit tendency in crash severity as shown in Figure 2 (The four rules in it is selected from the final rule set in this study).

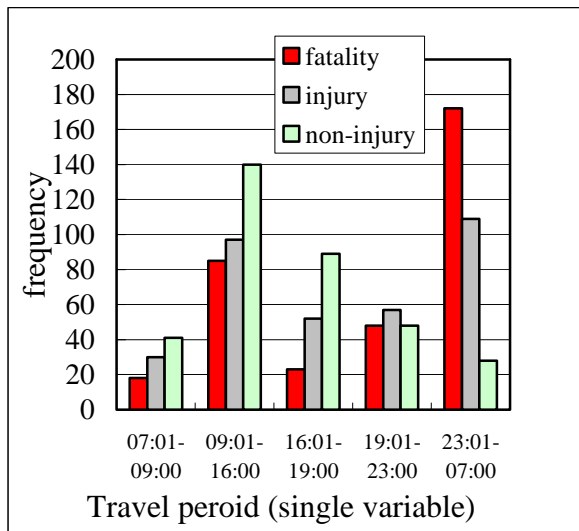


Figure 1 analysis of single variable

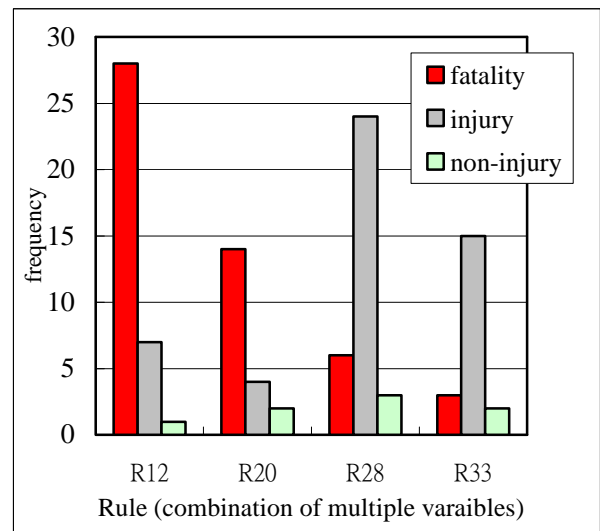


Figure 2 applying rules to analyze

Rule mining, also known as rule generation, rule recovery, or classification/association rule mining, is one of data mining techniques intended to mine for knowledge from available databases and toward decision support. Rule mining is naturally modeled as multi-objective problems with three criteria: (1) predictive accuracy, (2) comprehensibility, and (3) interestingness (Freitas, 1999; Ghosh and Nath, 2004). To automatically search for the optimal combination of rules from a considerable number of potential rules, genetic algorithms (GAs) are perhaps the most commonly used method. By employing GAs to learn of rules is named as genetic mining rule (GMR) (e.g. Freitas, 1999; Shin and Lee, 2002; Ghosh and Nath, 2004; Dehuri and Mall, 2006; Chen and Hsu, 2006). The performances of rule mining algorithms have been proven and applied in many fields. Thus, this paper aims to develop GMR model that can determine the optimal combination of decision rules to achieve the following goals: (1) to discover the key rules that determine the combination of contributing factors' level to crash severity; (2) to provide the possibility of post-adjustment (fine-tune) of the rules mined; (3) to accurately predict the crash severity. Previous relevant studies have seldom considered the problem of conflict and redundancy among the rules mined, our proposed GMR model will account for the conflict and redundancy in addition to conventional objectives: coverage ratio and predictive accuracy.

2.2 Data

The crash data were collected from 2003-2007 National Traffic Accident Investigation Reports compiled by National Police Agency, Taiwan. Each accident investigation report has been digitized and maintained in the database from which detailed individual crash data of freeway accidents are obtained. The individual crash data include detailed information regarding injury severity of each involved individual, time of accident, driver demographics (age, gender, driver sobriety), involved vehicle types, roadway geometry, traffic control condition, weather condition (clear, rain, fog), pavement conditions (wet, dry), lighting condition, and vehicle actions (moving straight, right-turn, left-turn, lane-change).

Considering the characteristics of crash occurrence may differ in collision type, the single-vehicle accident data are chosen to diminish the heterogeneity of crash data. Single-vehicle accidents are those in which only a single vehicle is involved. There are 5,563 single-vehicle crash cases occurring on Taiwan's freeways from 2003 to 2007. The injury severity of crashes is determined according to the injury degree of the worst-injured victims in the accident. Table 1 presents the definition and description of potential explanatory variables to crash severity.

Table 1 Crash data summarized from police accident investigation reports

Information	Variable	Type	Description
Surface condition	x_1	Categorical	1, dry; 2, wet or slippery
Signal control	x_2	Categorical	1, none; 2, yes
Driver gender	x_3	Categorical	1, male; 2, female
Weather	x_4	Categorical	1, sunny; 2, cloudy; 3, rain, storm, fog, etc.
Obstacle	x_5	Categorical	1, none; 2, work zone; 3, others
Lighting condition	x_6	Categorical	1, daytime; 2, dawn or dusk; 3, nighttime with illumination; 4, nighttime without illumination
Speed limit	x_7	Categorical (discretized)	1, 110 KPH; 2, 100KPH; 3, 90-70KPH; 4, 60-40KPH
Road status	x_8	Categorical	1, straight road; 2, grade and curved road; 3, tunnel, bridge, culvert, overpass; 4, others
Marking	x_9	Categorical	1, lane line with marker; 2, lane line without marker; 3, no lane-changing line; 4, no lane line
Use of safety belt	x_{10}	Categorical	1, safety belt fastened; 2, safety belt not fastened; 3, others or unknown
Use of cell phone	x_{11}	Categorical	1, use; 2, not in use; 3, others or unknown
License	x_{12}	Categorical	1, with license; 2, without license; 3, unknown
Driver occupation	x_{13}	Categorical	1, in job; 2, student; 3, jobless; 4, unknown
Driver age	x_{14}	Categorical (discretized)	1, under 30 years old; 2, 30-40 years old; 3, 40-50 years old; 4, 50-65 years old; 5, above 65 years old
Travel period	x_{15}	Categorical (discretized)	1, 07:01-09:00 morning peak; 2, 09:01-16:00 day off-peak; 3, 16:01-19:00 afternoon peak; 4, 19:01-23:00 night-peak; 5, 23:01-07:00 midnight to morning
Location	x_{16}	Categorical	1, fast lane, general lane; 2, shoulder, edge; 3, median; 4, accelerating or decelerating lane, ramp; 5, toll plaza and others
Vehicle type	x_{17}	Categorical	1, passenger car; 2, truck; 3, bus; 4, heavy truck, trailer truck, tractor; 5, others
Action	x_{18}	Categorical	1, forward; 2, left lane-change; 3, right lane-change; 4, urgent deceleration or stop; 5, others
Alcoholic use	x_{19}	Categorical	1, no; 2, under 0.25 mg/l (or 0.05%); 3, over 0.25 mg/l (or 0.05%); 4, cannot be tested; 5, unknown
Journey purpose	x_{20}	Categorical	1, work trip or school trip; 2, business trip; 3, transportation activity; 4, visiting, shopping; 5, others or unknown
Major cause	x_{21}	Categorical	1, improper lane-change; 2, speeding; 3, fail to keep a safe distance; 4, alcoholic use; 5, fail to pay attention to the front; 6, other driver's liability; 7, factors not attributed to drivers
Severity	y	Categorical	1, fatality; 2, injury; 3, no-injury

In Taiwan, crash severity in police investigation report is classified into three degrees: A1 (fatal crash), A2 (injury crash), and A3 (non-injury crash). The cases for these three degrees of crash severity are 226, 1,593, and 3,744, respectively—an uneven distribution commonly seen in the context of crash analysis. Furthermore, 70% of these 5,563 crash cases are randomly chosen for training (i.e., 3,895 cases) and the remaining 1,668 cases are used for model validation. χ^2 -test is performed and the result shows that severity distributions between training and validation datasets do not significantly differ.

2.3 Genetic rule mining model

Genetic mining rule (GMR), which can automatically learn of comprehensive rules from available dataset and toward decision support, is useful in accident analysis (Clarke *et al.*, 1998). The encoding method, fitness function, genetic operators, and rule selection of the proposed GMR model are narrated below.

2.3.1 Encoding method

To represent the relationship between explanatory variables and crash severity, each chromosome is used to represent a potential if-then rule. The conditions associated in the “if part” are termed as antecedence part and those in the “then part” are named as consequent part. Besides, the antecedent part consists of at least one variable, but at most 21 variables, selected from Table 1. And the consequent part is composed by, of course, only one variable: severity degree. In general, a rule is a

knowledge representation of the form “If A Then C ,” where A is a set of cases satisfying the conjunction of predicting attribute values and C is a set of cases with the same predicted degree. Thus, a typical rule i can be of the form: Rule i : If $x_1=a_{i1}$ and $x_2=a_{i2}$...and $x_j=a_{ij}$... and $x_{21}=a_{i21}$ Then $y=g_i$. Or, in a shorter form: Rule i : If A_i Then C_i , where a_{ij} is the categorical value of j^{th} attribute variable in rule i . g_i is the value of classification variable in rule i , which ranges from 1 to 3 representing three degrees of crash severity. A_i and C_i are the sets of parties satisfying the antecedent part and consequent part of rule i , respectively.

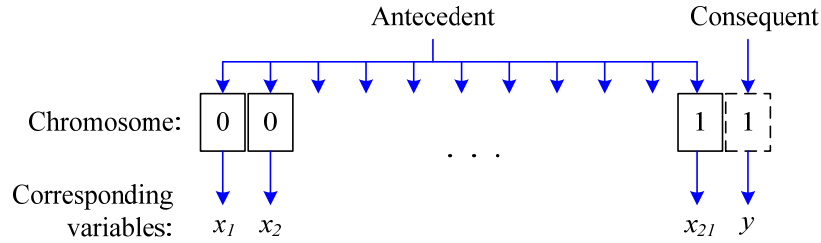


Figure 3 Encoding method of the proposed GMR model

By encoding a rule as a chromosome, each gene is used to represent a corresponding variable. Since the number of potential variables of antecedent and consequent is respectively 21 and one, the length of a chromosome is 22. Each gene will then take one of the categorical values of the corresponding variable. Because the ranges of all variables are different, the ranges of gene values also vary. Moreover, if a gene in a rule antecedent takes a value of 0, it represents the corresponding variable not considered by the rule. If all genes representing the rule antecedent simultaneously take 0 or if the gene representing the rule consequent is 0, then the rule is not included.

Based on this, a rule of “If surface condition=dry and occupation=in job and actions=left lane-change and Then degree of severity=injury” can be encoded as 1000000000001000020002. This rule also contains a family of 4.838×10^{10} offspring rules in total, which can be represented by “If $x_1=1$ and $x_2=\{0, 1, 2\}$ and $x_3=\{0, 1, 2\}$ and $x_4=\{0, 1, \dots, 3\}$ and $x_5=\{0, 1, \dots, 3\}$ and $x_6=\{0, 1, \dots, 4\}$ and $x_7=\{0, 1, \dots, 4\}$ and $x_8=\{0, 1, \dots, 4\}$ and $x_9=\{0, 1, \dots, 4\}$ and $x_{10}=\{0, 1, \dots, 4\}$ and $x_{11}=\{0, 1, \dots, 4\}$ and $x_{12}=\{0, 1, \dots, 3\}$ and $x_{13}=1$ and $x_{14}=\{0, 1, \dots, 5\}$ and $x_{15}=\{0, 1, \dots, 5\}$ and $x_{16}=\{0, 1, \dots, 5\}$ and $x_{17}=\{0, 1, \dots, 5\}$ and $x_{18}=2$ and $x_{19}=\{0, 1, \dots, 5\}$ and $x_{20}=\{0, 1, \dots, 5\}$ and $x_{21}=\{0, 1, \dots, 7\}$ and Then $y=2$.” That is, any case satisfying any one of the offspring rules will certainly also satisfy their parent rule. Generally, the more variable present in the antecedent part (taking non-zero values), the more specific of a rule is (less number of parties will satisfy the rule).

The proposed algorithm aims to select a set of rules which can best predict the severity degree based upon these twenty one explanatory variables. The total number of potential rules equals $3 \times 3 \times 3 \times 4 \times 4 \times 5 \times 5 \times 5 \times 5 \times 4 \times 4 \times 4 \times 5 \times 6 \times 6 \times 6 \times 6 \times 6 \times 6 \times 8 = 1.935 \times 10^{14}$. Obviously, it is barely possible to compare all rule combinations through a total enumeration approach.

2.3.2 Fitness function

An individual chromosome, a rule, with a higher fitness function value has a higher probability to be selected for reproducing offspring. The role of fitness function is to evaluate the quality of the rule numerically. To determine the fitness function, there are three common factors frequently taken into consideration: coverage, completeness and confidence of the rule. The coverage ratio of rule i (*i.e.*, the cases satisfied by the rule antecedent) is denoted by $|A|$: the cardinality of set A (the number of elements in set A). The completeness of the rule (*i.e.*, the proportion of cases of the target class covered by the rule) is given by $|A \cap C|/|C|$. The confidence of rule i (*i.e.*, the predictive accuracy) is given by $|A \cap C|/|A|$ (Freitas, 1999). Shin and Lee(2002) adopted hit ratio(confidence)

as the fitness function which is also defined as predictive accuracy plus coverage in another study (Kim and Han, 2003). However, it is the performance of the entire rule set that should be emphasized instead of those ones of individual rules themselves. In other words, the good performances of individual rules do not guarantee that the combination of these rules also performs well. It results from the redundancy and conflict between rules. In order to overcome this problem, the fitness function is set in this paper as the increase of correctly classified cases by the rule set combining the previous mined rules and the new rule, which can be expressed as follows:

$$f_i = N_{nrs} - N_{prs} \quad (1)$$

where, N_{nrs} is the number of cases that are correctly classified by the rule set combining the previous mined rules with the rule i , and N_{prs} is the number of cases that are correctly classified by the previous mined rules.

The previous mined rules are also called the temporary rule set in this study. By means of the fitness function above, the effect caused by redundancy or conflict between rules would be effectively reduced in rule mining process. When a new rule is extracted from the final population, it would certainly increase the performance of entire rule set as the new rule set combines the new rule with the temporary rule set.

2.3.3 Genetic operators

Because the genes in our GMR model are not encoded binary, simple genetic algorithms proposed by Goldberg (1989) cannot be used. Instead, we employ the max-min-arithmetical crossover proposed by Herrera *et al.* (1998) and the non-uniform mutation proposed by Michalewicz (1992). A brief description is given below.

(1) Max-min-arithmetical crossover

Let $G_w^t = \{ g_{w1}^t, \dots, g_{wk}^t, \dots, g_{wK}^t \}$ and $G_v^t = \{ g_{v1}^t, \dots, g_{vk}^t, \dots, g_{vK}^t \}$ be two chromosomes selected for crossover, the following four offsprings can be generated:

$$G_1^{t+1} = aG_w^t + (1-a)G_v^t \quad (2)$$

$$G_2^{t+1} = aG_v^t + (1-a)G_w^t \quad (3)$$

$$G_3^{t+1} \text{ with } g_{3k}^{t+1} = \min\{g_{wk}^t, g_{vk}^t\} \quad (4)$$

$$G_4^{t+1} \text{ with } g_{4k}^{t+1} = \max\{g_{wk}^t, g_{vk}^t\} \quad (5)$$

where a is a parameter ($0 < a < 1$) and t is the number of generations.

(2) Non-uniform mutation

Let $G_t = \{ g_1^t, \dots, g_k^t, \dots, g_K^t \}$ be a chromosome and the gene g_k^t be selected for mutation (the domain of g_k^t is $[g_k^l, g_k^u]$), the value of g_k^{t+1} after mutation can be computed as follows:

$$g_k^{t+1} = \begin{cases} g_k^t + \Delta(t, g_k^u - g_k^t) & \text{if } b = 0 \\ g_k^t - \Delta(t, g_k^t - g_k^l) & \text{if } b = 1 \end{cases} \quad (6)$$

where b randomly takes the binary value of 0 or 1. The function $\Delta(t, z)$ returns to a value in the range of $[0, z]$ such that the probability of $\Delta(t, z)$ approaches to 0 as t increases:

$$\Delta(t, z) = z(1 - r^{(1-t/T)^h}) \quad (7)$$

where r is a random number in the interval $[0, 1]$, T is the maximum number of generations and h is a given constant. In eq. (7), the value returned by $\Delta(t, z)$ will gradually decrease as the evolution

progresses.

2.3.4 Rule selection

The method of extracting rules has profound effects on their accompanied performance. Conventionally, a group of different rules is obtained simultaneously from the final results as the stopping criterion is met. Generally speaking, it is an important issue to avoid selecting redundant or conflicting rules during the rule selection process. The redundancy or conflict between the selected rules would lead to reduce the performance of the prediction model, as well as increasing the difficulty in interpreting the causal relationship between explanatory variables and crash severity. However, it is probably difficult to avoid this condition and little information could be found in the literature on dealing with this issue (Shin and Lee, 2002; Kim and Han, 2003; Chen and Hsu, 2006). On the other hand, the mined rules are often too complicated to be understood instead of being interpretable, shorter, and simpler. In order to improve these problems, a learn-one-rule function combining with a neighborhood search was introduced over the rule mining process in this study. Instead of searching a good rule set at a time, a stepwise rule set building procedure with a greedy strategy is proposed. Applying the learn-one-rule function combining with a neighborhood search, the rule set is constructed according to the following steps (as shown in Figure 4):

- Step 1: Rank rules in the final population according to their fitness values in a descending order.
- Step 2: Select the rule with the highest fitness value and perform a neighborhood search with improvement and parsimony principle for rule modification.
- Step 3: Update the temporary rule set by the modified rule.
- Step 4: Terminate until the number of rules in the temporary rule set hit the preset number. Otherwise, implement the GAs for another run and go to Step 1.

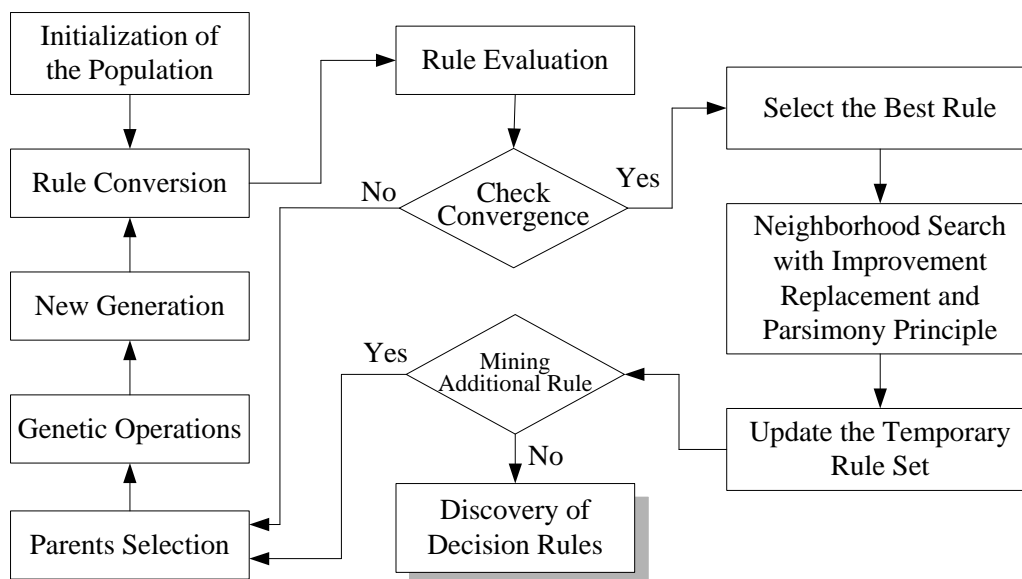


Figure 4 The GA based mining approach

After a rule is selected, a rule modification scheme is introduced. There are two mechanisms in the rule modification process, including improvement replacement and parsimony principle. Due to the characteristic of stochastic operation in evolutionary process, it is understandable that there might be some better points existing near the current solution point in the search space. Based on this, Comparative rules are created by enumerating all other attribute values of one variable controlling all other variables. In the mechanism of improvement replacement, when the predictive accuracy of a comparative rule combining with the previous rule set is better than the raw rule in the same condition, the value of the checked variable would be substituted by the value of the same variable in that comparative rule, as shown in the left part in Figure 5. If there is no better point found, the

mechanism of parsimony principle will hold. When the original value of the checked variable is not zero, but the value of the checked variable is zero in comparative rule with the same predictive accuracy in the same condition, the value of the checked variable would be substituted by zero, as shown in the right part in Figure 5. In this study, the order of checking all explanatory variables is from x_1 to x_{2j} . After all explanatory variables are checked, the last adjusted rule will be put into the temporary rule set for next rule mining if needed.

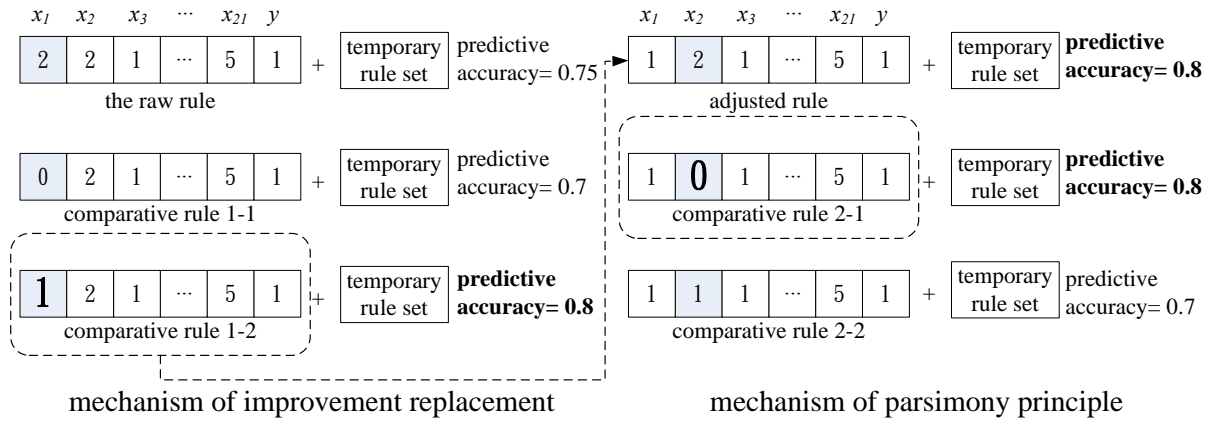


Figure 5 Rule modification process

It is almost inevitable that two or more rules with different predicted classes may be simultaneously fired by a crash case. In this situation, the case is would be predicted as the class of the rule with the highest accuracy if two or more rules are applied to the case at the same time.

2.4 Results

The parameters of the proposed GMR model are set as follows: population size=50, crossover rate=0.85, mutation rate=0.08, and maximum number of generations=1,000 (the stopping criterion). The number of rules to mine is set as 55. The learning process of the GMR model is shown in Figure 6.

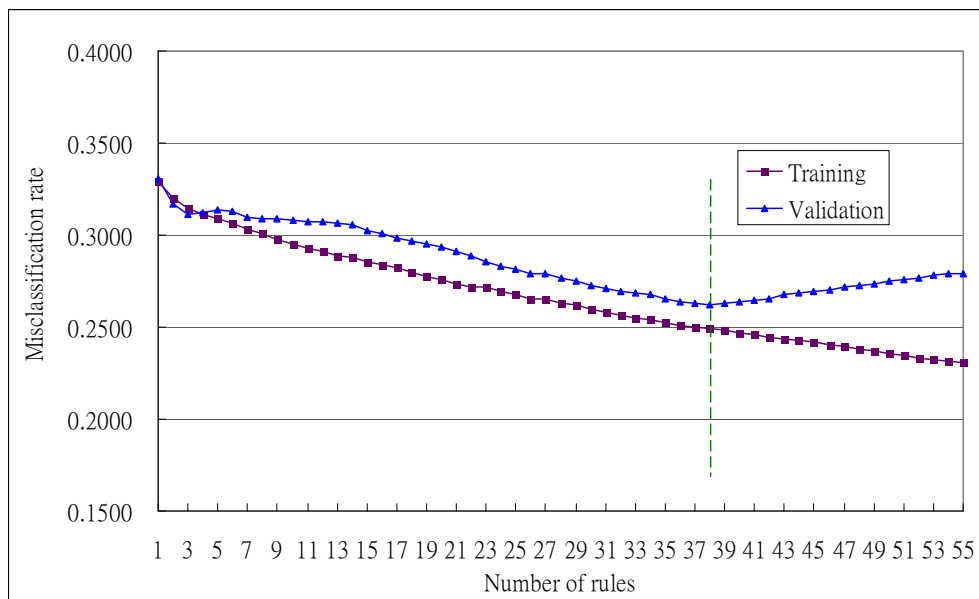


Figure 6 Learning process of the GMR model

Theoretically, the misclassification rate can be lowered to zero monotonically by increasing the number of rules in the GMR model. However, a good classification model should not only fit the

training data well, it must also accurately classify records it has never seen before. To avoid model overfitting, 38 rules are selected in the GMR model as the misclassification rate of validation data hit the lowest value. Table 2 shows the final selected rules along with its corresponding performance indices. Note that a total of 38 rules are selected with a descending order according to PA_i . In terms of predictive accuracy (PA_i), the top twenty five rules have remarkably higher values than the rest of thirteen rules. In terms of coverage ratio (CR_i), R23 can explain 3,800 cases, followed by R30 (1,460 cases) and R31 (529 cases). In contrast, some rules cover only very few cases, such as R1 (6 cases), R6 (6 cases) or R7 (6 cases).

The importance of variable can be identified by the number of its presence in all rules. The number of variables with values other than 0 (*i.e.* the variable is not considered by the rule) in all rules is then calculated. In this regard, x_{13} (driver occupation) is the most important variable which appears in 16 rules, followed by x_{16} (location), x_{15} (travel period), and x_{17} (vehicle type). Two variables are shown in less than three rules, which are x_2 (signal control) and x_8 (road status), indicating their least significance to crash severity. There are six rules associated with A1 crash, twenty-eight rules with A2 crash, and four rules with A3 crash.

Most of the rules could be readily inspected and explained by the if-then relationship of the rules themselves. Taking R1 for instance, the rule indicates that when speed limit is 40~60 KPH and driver's age is over 65 years old, it tends to lead A2 crash. R2 shows when drivers are male, in job and under 30 years old, speed limit is 100 KPH, travel period is midnight to morning, and major cause is alcoholic, it tends to lead A2 crash. As to R19, when safety belt is not fastened with driver's speeding, it tends to cause A1 crash. In contrast to R19, R23 reveals when safety belt is fastened, it tends to be less severe (A3 crash). The rest may be deduced by analogy. More exploration of the potential implications of the rules is depicted as the following. In regard to driver characteristics, it is interesting that jobless driver combining with specific conditions would tend to cause A2 crash. The conditions include cloud (R3), nighttime with illumination, under 30 years old, and midnight to morning (R20), and no obstacle (R26). Regarding Behavior and environment factors, when safety belt is not fastened with driver's speeding, it tends to cause A1 crash (R19). Use of cell phone combining with the antecedents of R14 and R35 tends to lead A2 crash. The alcoholic use has positive correlation in crash severity. On the other hand, wet or slippery surface condition and obstacle do not have significant effects on crash severity. About vehicle type, truck combining with the antecedents of R6, R13, R18, and R27 is likely to lead A2 crash. As to trip characteristic, midnight to morning combining with the antecedents of R2, R5, R20, R21, and R29 also tends to lead A2 crash. The above-rule interpretations might be useful references for law enforcement or management by the related authorities.

Table 2 Combination of rules mined by GMR model

Rules	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8	x_9	x_{10}	x_{11}	x_{12}	x_{13}	x_{14}	x_{15}	x_{16}	x_{17}	x_{18}	x_{19}	x_{20}	x_{21}	y	CR_i	PA_i	
R1	1					1	4		1					5					1			2	6	1.000	
R2			1	1			2			1			1	1	5			1			4	2	12	0.917	
R3				2									3									2	12	0.917	
R4			2				1				1							1	2			2	10	0.900	
R5	1		1												3	1				3		2	12	0.833	
R6				1									1		5	2	2					2	6	0.833	
R7	1				1			1		3	3				2							2	6	0.833	
R8																		5		3		2	12	0.833	
R9			1			3			3				1			1		1				2	11	0.818	
R10				1			2		1					4				1		5	7	2	16	0.813	
R11				1			1					2					1					2	16	0.813	
R12						1	3						1								5	3	64	0.813	
R13	1											2			2		2					2	10	0.800	
R14					1		3				1									2		2	15	0.800	
R15			1	1												1	4		1			3	239	0.799	
R16	1				1	4							1	3			2					2	22	0.773	
R17												1				2			2	5		2	12	0.750	
R18			2										1		2		2				5	2	12	0.750	
R19										2												2	1	11	0.727
R20						3							3	1	5							2	11	0.727	
R21													4		5	1	2					2	25	0.720	
R22	1						3		1					2								2	14	0.714	
R23										1												3	3800	0.687	
R24						4								1		1			1			3	201	0.687	
R25											3											1	106	0.613	
R26					1								3									2	154	0.435	
R27				1													2				1	2	77	0.429	
R28																			4			1	47	0.426	
R29	1												1	3	5				1			2	91	0.374	
R30	1				1			1	1			1				1						2	1460	0.325	
R31						4			1				1									2	529	0.319	
R32							3					1	1				1					2	305	0.302	
R33															2				2			2	64	0.297	
R34	2														4		1					2	149	0.262	
R35						1	1			1	1	1	1			2		1	1			2	121	0.215	
R36	1												1								2	1	97	0.196	
R37			1	1		3								2		1						1	75	0.080	
R38													1			2		1	1			1	267	0.064	
<i>m</i>	10	0	7	8	5	9	10	2	6	5	4	7	16	9	11	13	11	8	10	5	4	-	-	-	

Note: *m* is the number of variable presence in the selected 38 rules.

Table 3 gives the distribution of cases with degree of severity predicted by GMR model and with real degree of severity. As shown in Table 3, in the training dataset, the proposed GMR model can actually predict the A3 crash (correct rate 80.77%), followed by A2 crash (64.90%) and A1 (53.13%). The overall correct rate of the proposed GMR model in training has achieved 75.10%. In the validation dataset, the overall correct rate has achieved 73.80%.

Table 3 Number of cases with degree of severity predicted by GMR

Datasets	Real severity	Predicted severity			Total
		A1	A2	A3	
Training	A1	<u>85 (53.13%)</u>	46 (28.75%)	29 (18.13%)	160 (100.00)
	A2	32 (2.87%)	<u>723 (64.90%)</u>	359 (32.23%)	1114 (100.00)
	A3	22 (0.84%)	482 (18.39%)	<u>2117 (80.77%)</u>	2621 (100.00)
	Total		139	1251	2505
Validation	A1	<u>37 (56.06%)</u>	15 (22.73%)	14 (21.21%)	66 (100.00)
	A2	3 (0.63%)	<u>307 (64.09%)</u>	169 (35.28%)	479 (100.00)
	A3	11 (0.98%)	225 (20.04%)	<u>887 (78.98%)</u>	1123 (100.00)
	Total		51	547	1070

Note: The percentages are given in the parentheses.

2.5 Comparisons

For comparison purpose, a decision tree (DT) model is also used to mine the rules explaining the same crash dataset. The DT model is performed by SAS Enterprise Miner Release 4.3. Several settings of the DT model are tried and the best performed settings are as follows. Splitting criterion is Gini reduction. Minimum number of observations in a leaf is 1. Observations required for a split search is 8. Maximum number of branches from a node is 2. Maximum depth of tree is 6. Splitting rules saved in each node is 5. The learning process of the DT model is depicted in Figure 7. Note that the misclassification rate decreases as the number of leaves gets larger.

Table 4 presents the number of cases with various degrees of severity predicted by the DT model. Note that the DT model performs better in predicting the A3 crash (correct rates in training and validation are 97.71% and 97.15%, respectively) than the proposed GMR model. However, the DT model performs much worse than the proposed GMR model while predicting both A1 and A2 crashes. Averagely, the overall correct rates of the DT model in training and validation are 70.24% and 69.54%, respectively, which are inferior to the proposed GMR model.

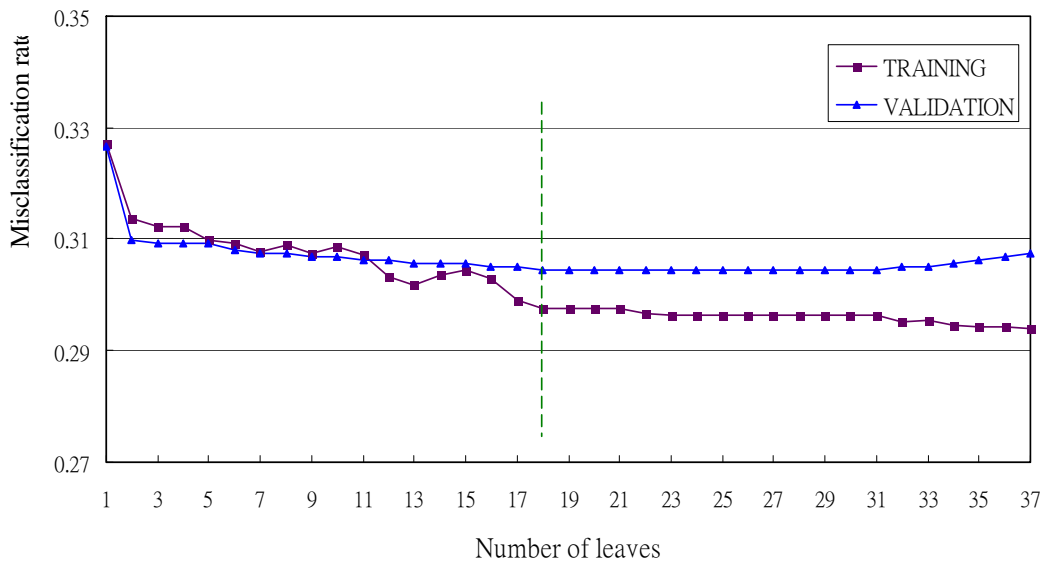


Figure 7 Learning process of the DT model

Table 4 Number of cases with degree of severity predicted by DT based on balanced dataset

Datasets	Real severity	Predicted severity			Total
		A1	A2	A3	
Training	A1	71 (44.38%)	10 (6.25%)	79 (49.38%)	160 (100.00)
	A2	34 (3.05%)	104 (9.34%)	976 (87.61%)	1114 (100.00)
	A3	10 (0.38%)	50 (1.91%)	2561 (97.71%)	2621 (100.00)
	Total	115	164	3616	3895
Validation	A1	36 (54.55%)	1 (1.52%)	29 (43.94%)	66 (100.00)
	A2	7 (1.46%)	33 (6.89%)	439 (91.65%)	479 (100.00)
	A3	7 (0.62%)	25 (2.23%)	1091 (97.15%)	1123 (100.00)
	Total	50	59	1559	1668

Note: The percentages are given in the parentheses.

A total of 18 rules are generated by the DT model as follows: two rules associated with A1 crash, six rules with A2 crash, and ten rules with A3 crash.

R1: If $x_{11}=3$ Then $y=1$

R2: If $x_{11}=2$ Then $y=3$

R3: If $x_{21}=2$ and $x_{10} = \{2, 3\}$ and $x_{17} = \{1, 4\}$ and $x_{11}=1$ Then $y=1$

- R4: If $x_3=2$ and $x_4= \{2, 3\}$ and $x_{17}= \{2, 3, 5\}$ and $x_{11}=1$ Then $y=2$
R5: If $x_3=1$ and $x_4= \{2, 3\}$ and $x_{17}= \{2, 3, 5\}$ and $x_{11}=1$ Then $y=3$
R6: If $x_{12}=1$ and $x_{19}=1$ and $x_{10}=1$ and $x_{17}= \{1, 4\}$ and $x_{11}=1$ Then $y=3$
R7: If $x_{21}= \{2, 3, 4, 5, 7\}$ and $x_{19}= \{2, 3, 4, 5\}$ and $x_{10}=1$ and $x_{17}= \{1, 4\}$ and $x_{11}=1$ Then $y=3$
R8: If $x_{15}= \{2, 4, 5\}$ and $x_{21}= \{1, 3, 4, 5, 6, 7\}$ and $x_{10}= \{2, 3\}$ and $x_{17}= \{1, 4\}$ and $x_{11}=1$ Then $y=2$
R9: If $x_{15}= \{1, 3\}$ and $x_{21}= \{1, 3, 4, 5, 6, 7\}$ and $x_{10}= \{2, 3\}$ and $x_{17}= \{1, 4\}$ and $x_{11}=1$ Then $y=3$
R10: If $x_{13}= \{1, 2, 4\}$ and $x_{21}= \{2, 3, 6\}$ and $x_4=1$ and $x_{17}= \{2, 3, 5\}$ and $x_{11}=1$ Then $y=3$
R11: If $x_{13}=3$ and $x_{21}= \{2, 3, 6\}$ and $x_4=1$ and $x_{17}= \{2, 3, 5\}$ and $x_{11}=1$ Then $y=2$
R12: If $x_{20}=3$ and $x_{21}= \{1, 4, 5, 7\}$ and $x_4=1$ and $x_{17}= \{2, 3, 5\}$ and $x_{11}=1$ Then $y=3$
R13: If $x_{21}= \{1, 2, 3, 6, 7\}$ and $x_{12}= \{2, 3\}$ and $x_{19}=1$ and $x_{10}=1$ and $x_{17}= \{1, 4\}$ and $x_{11}=1$ Then $y=3$
R14: If $x_{21}=5$ and $x_{12}= \{2, 3\}$ and $x_{19}=1$ and $x_{10}=1$ and $x_{17}= \{1, 4\}$ and $x_{11}=1$ Then $y=2$
R15: If $x_{14}= \{1, 2\}$ and $x_{21}= \{1, 6\}$ and $x_{19}= \{2, 3, 4, 5\}$ and $x_{10}=1$ and $x_{17}= \{1, 4\}$ and $x_{11}=1$ Then $y=2$
R16: If $x_{14}= \{2, 3, 5\}$ and $x_{21}= \{1, 6\}$ and $x_{19}= \{2, 3, 4, 5\}$ and $x_{10}=1$ and $x_{17}= \{1, 4\}$ and $x_{11}=1$ Then $y=3$
R17: If $x_{15}= \{1, 2, 3, 4\}$ and $x_{20}= \{1, 2, 4, 5\}$ and $x_{21}= \{1, 4, 5, 7\}$ and $x_4=1$ and $x_{17}= \{2, 3, 5\}$ and $x_{11}=1$ Then $y=3$
R18: If $x_{15}=5$ and $x_{20}= \{1, 2, 4, 5\}$ and $x_{21}= \{1, 4, 5, 7\}$ and $x_4=1$ and $x_{17}= \{2, 3, 5\}$ and $x_{11}=1$ Then $y=2$

2.6 Conclusion

This paper identifies risky conditions (joint effects of risk factors) to crash severity by developing a novel genetic mining rule (GMR) model. Three different types of A1, A2 and A3 single-vehicle crash cases are drawn from 2003-2007 Taiwan's freeway accidents dataset. A total of 38 rules have been mined which can achieve an overall correct rate of 75.10% in training and 73.80% in validation, respectively. Our proposed GMR model has demonstrated superior to the conventional decision tree (DT) model, which can only achieve an overall correct rate of 70.24% in training and 69.54% in validation, respectively, with the same database. According to the mined rules, x_{13} (driver occupation), x_{16} (location), x_{15} (travel period), and x_{17} (vehicle type) are the four key factors contributing to crash severity. Consequently, attention must be paid to these four factors to ameliorate the traffic safety.

Some directions for future studies can be identified. First, the neighboring traffic condition of the crash is also an important factor to crash severity; however, the police accident investigation report did not record such information. The crash data may be further matched with the traffic database so as to gain more information regarding the contributing factors to crash severity. Second, in order to lessen the model complexity, various performance indices may be integrated into an overall fitness function; namely, a multi-objective GMR model deserves further elaboration. Last but not least, more comparisons can be made to other commonly used methods (e.g., logistic regression model, ordered Logit model, artificial neural network) to demonstrate the superiority of the proposed model.

三、計畫成果自評

本計畫為三年期計畫。其中，本期中報告已完成第二年期之研究內容，依據本研究計畫書之原訂研究內容，完成之研究成果如下：

1.研究目的與範圍確立

本年度之研究係奠基於第一個研究年期之主要研究成果，原擬進一步建立螞蟻規則探勘模式，以高速公路事故資料為應用實例，進行分析與預測模式之構建，並與基因規則探勘、判別分析、羅吉斯迴歸等模式績效進行比較分析。惟由於螞蟻規則探勘模式之績效表現不佳，

故改提出「改良型逐步基因規則探勘模式 (Stepwise GRM, SGRM)」。後續有關螞蟻規則探勘之相關研究課題，均以 SGRM 模式替代之。

2. 相關文獻彙析

本研究擬利用數位圖書館及網際網路等資源，檢索有關規則探勘、螞蟻多目標數學規劃及事故分析與預測等相關文獻資料，俾供本研究進行模式建構與比較分析之參考。

3. 蒐集高速公路事故資料並篩選重要解釋變數

高速公路事故資料分為 A1 (死亡事故)、A2 (受傷事故)、A3 (財損事故) 三大類。相關變數包括：事故發生時間地點、當時天候狀況資料、當地道路幾何條件資料、事故類型、主要肇事原因、傷亡狀況、交通管制狀況、駕駛人行為與違規狀況等。將利用交叉分析表方式，先作顯著變數之初步篩選。此外，由於不同類型事故之危險因子與危險情況可能差異甚大，因此，本研究乃先加以分類後，再進行應用分析。以單車事故為例，共計蒐集 2003 至 2007 年間，5,563 事故件數，其中，A1、A2 及 A3 各 226 件、1,593 件，以及 3,744 件。

4. 建立高速公路事故分析與預測模式

以第一個研究年期所建構之基因規則探勘模式為基礎，進一步提出改良式逐步基因規則探勘模式，並應用於高速公路事故分析與預測案例。在規則學習與驗證上，本研究採用交叉驗證方式，將所有案例分為兩部份，分別作為訓練資料及驗證資料。

5. 模式績效之比較分析

本研究將除與基因規則探勘模式進行比較外，並將同時與決策樹分析結果進行比較。

6. 推理規則之產生與詮釋

經由比較分析後，可依據表現較佳之規則探勘模式所挑選規則，進行分析及詮釋，並加以整理列表。以深入了解各環境變數群、交通管制變數群、駕駛人行為變數群間對事故嚴重性之聯合效果關係，並據以研提改善策略。

7. 結論與建議

由本年度之研究經驗、求解結果及比較分析，研提具體研究結論與後續研究方向之建議。

上述第二個年期之預期研究成果已順利達成，並為下一年度之研究奠定良好基礎。此外，本計畫之主要成果已分別發表國際期刊 1 篇文章[28]，並已改寫投稿國際研討會及學術期刊中[29, 30]。此外，本計畫亦用以指導一名博士生進行論文寫作[31]。

四、參考文獻

1. Abdel-Aty, M. and Pande, A., 2007. Crash data analysis: Collective vs. individual crash level approach. *Journal of Safety Research* 38, 581-587.
2. Abdel-Aty, M. and Radwan, A.E., 2000. Modeling traffic accident occurrence and involvement. *Accident Analysis and Prevention* 32, 633-542.
3. Al-Ghamdi, A.S., 2002. Pedestrian-vehicle crashes and analytical techniques for stratified contingency tables. *Accident Analysis and Prevention* 34, 205-214.
4. Chang, L.Y. and Mannering, F., 1999. Analysis of injury severity and vehicle occupancy in truck- and non-truck-involved accidents. *Accident Analysis and Prevention* 31, 579-592.
5. Chang, L.Y. and Chen, W.C., 2005. Data mining of tree-based models to analyze freeway accident frequency. *Journal of Safety Research* 36, 365-375.

6. Chang, L.Y. and Wang, H.W., 2006. Analysis of traffic injury severity: An application of non-parametric classification tree techniques. *Accident Analysis and Prevention* 38, 1019-1027.
7. Chen, T.C. and Hsu, T.C., 2006. A GAs-based approach for mining breast cancer pattern. *Expert Systems with Applications* 30, 674-681.
8. Chiou, Y.C., 2006. An artificial neural network-based expert system for the accident appraisal of two-car crash accidents. *Accident Analysis and Prevention* 38, 777-785.
9. Clarke, D.D., Forsyth, R.S. and Wright, R.L., 1998. Behavioural factors in accidents at road junctions: The use of a genetic algorithm to extract descriptive rules from police case files. *Accident Analysis and Prevention* 30, No. 2, 223-234.
10. Dehuri, S. and Mall, R., 2006. Prediction and comprehensible rule discovery using a multi-objective genetic algorithm. *Knowledge-Based System* 19, 413-421.
11. Delen, D., Sharda, R. and Bessonov, M., 2006. Identifying significant predictors of injury severity in traffic accidents using series of artificial neural networks. *Accident Analysis and Prevention* 38, 434-444.
12. Dissanayake, S. and Lu, J.J., 2002. Factors influential in making an injury severity difference to older drivers involved in fixed object-passenger car crashes. *Accident Analysis and Prevention* 34, 609-618.
13. Freitas, A.A., 1999. On rule interestingness measures. *Knowledge-Based Systems* 12, 309-315.
14. Ghosh, A. and Nath, B., 2004. Multi-objective rule mining using genetic algorithms. *Information Sciences* 163, 123-133.
15. Goldberg, D.E., 1989. *Genetic Algorithms in Search, Optimization, and Machine Learning*. Addison-Wesley, New York.
16. Greibe, P., 2003. Accident prediction models for urban roads. *Accident Analysis and Prevention* 35, 273-285.
17. Herrera, F., Lozano, M. and Verdegay, J.L., 1998. A learning process for fuzzy control rules using genetic algorithms. *Fuzzy Sets and Systems* 100, 143-158.
18. Ivan, J.N., Pasupathy, R.K. and Ossenbruggen, P.J., 1999. Differences in causality factors for single and multi-vehicle crashes on two-lane roads. *Accident Analysis and Prevention* 31, 695-704.
19. Kim, M.J. and Han, I., 2003. The discovery of experts' decision rules from qualitative bankruptcy data using genetic algorithms. *Expert Systems with Applications* 25, 637-646.
20. Michalewicz, Z., 1992. *Genetic Algorithms + Data Structures = Evolution Programs*, Springer, Berlin.
21. Milton, J., Shankar, V., and Mannering, F., 2008. Highway accident severities and the mixed logit model: An exploratory empirical analysis, *Accident Analysis and Prevention* 40, 260-266.
22. Poch, M. and Mannering, F., 1996. Negative binomial analysis of intersection. *Journal of Transportation Engineering* 12, 105-113.
23. Shanker V., Mannering, F. and Barfield, W., 1996. Statistical analysis of accident severity on rural freeways, *Accident Analysis and Prevention* 28, 391-401.
24. Shin, K.S. and Lee, Y.J., 2002. A genetic algorithm application in bankruptcy prediction modeling. *Expert Systems with Applications* 23, 321-328.
25. Sze, N.N. and Wong, S.C., 2007. Diagnostic analysis of the logistic model for pedestrian injury severity in traffic crashes. *Accident Analysis and Prevention* 39, 1267-1278.
26. Tay, R. and Rifaat, S.M., 2007. Factors contributing to the severity of intersection crashes. *Journal of Advanced Transportation* 41, 245-265.
27. Wong, S.C., Sze, N.N. and Li, Y.C., 2007. Contributory factors to traffic crashes at signalized intersections in Hong Kong, *Accident Analysis and Prevention* 39, 1107-1113.
28. Chiou, Y.C., Lan, W.L. and Chen, W.B. (2010) "Contributory factors to crash severity in Taiwan freeways: Genetic mining rule approach," *Journal of Eastern Asia Society for Transportation Studies*. (Accepted)
29. Chiou, Y.C., Lan, W.L. and Chen, W.B. (2010) "Identification and estimation the risky

30. Chiou, Y.C., Lan, W.L. and Chen, W.B. (2010) "Identification of risky conditions contributing to crash severity with genetic mining rules," submitted to the *15th Conference of Hong Kong Society for Transportation Studies*.
31. 陳文斌, Identifying contributory factors to crash severity in Taiwan freeways by genetic rough set rule mining, 交通大學交通運輸研究所, 博士論文(進行中), 民國98年。