

## 使用資料探勘演算法預測非肺小細胞肺癌患者存活情形及其效能比較

翁紹宏<sup>a\*</sup> 陳麗帆<sup>b</sup> 朱基銘<sup>b</sup> 白璐<sup>c</sup> 楊燦<sup>d</sup> 劉立<sup>a</sup> 孫建安<sup>e</sup>

<sup>a</sup> 台北醫學大學醫學資訊研究所

<sup>b</sup> 國防醫學院公共衛生學系暨研究所

<sup>c</sup> 台北醫學大學傷害防治學研究所及受聘為萬芳醫院急診醫學科

<sup>d</sup> 美和技術學院健康事業管理系

<sup>e</sup> 輔仁大學公共衛生學系所

\*通訊作者：翁紹宏，wengsh@ndmctsgh.edu.tw

### 摘要

本研究使用決策樹、類神經網路和邏輯斯迴歸模型三種資料探勘演算法來探討影響非小細胞肺癌的預後因子及影響模型預測能力的因素（不同資料庫、不同死因、單或複合模型和樣本大小）。本研究的研究對象為美國癌症登記資料檔（the Surveillance, Epidemiology, and End Results, SEER），選取自西元1988年至2004年間診斷為非小細胞肺癌患者，並依據死因的不同將其分為死於肺癌與死於轉移癌；資料經過修整後，共有16個自變項納入分析，並根據其存活情形分為一年、三年和五年存活情形。模型的評估指標為準確率（accuracy, ACC）、ROC曲線下的面積（area under the ROC curve, AUC）和外推力（external generalization）。為避免隨機抽樣造成資料的誤差，本研究將對資料庫進行十折交叉驗證（10-fold cross-validation）。

研究結果顯示，影響美國非小細胞肺癌患者死於肺癌的一年、三年和五年存活情形預後因子為手術種類、臨床分組和腫瘤擴散程度；影響非小細胞肺癌患者死於轉移癌的一年、三年和五年存活情形預後因子為手術種類、臨床分期和檢驗淋巴結個數。

三個模型的預測力表現以類神經網路模型的表現較好，外推能力則以邏輯斯迴歸模型表現較好。樣本人數建議為至少為3500人，其中以邏輯氏回歸模型最容易受到小樣本的影響；決策樹則是會因為資料庫提供的訊息不足而無法成樹。複合模型的部份，其結果顯示，當決策樹的測試組的ACC值較另外兩個模型好時，則複合模型的測試組AUC值就會提高。

**關鍵字：**lung cancer, decision tree, artificial neural network, logistic regression, SEER

### 1、前言

根據世界衛生組織(World Health Organization; WHO)的統計，近年來全世界的癌症死亡率正快速的上升中，其中以肺癌的死亡率增加最為快速。而肺癌發生的主要原因和機轉，目前尚未完全明瞭，除高死亡率外，其預後狀況亦不佳，故肺癌的研究和防治在公共衛生研究上的重要性不可言喻。

因腫瘤的預後情形受到許多因子影響，且不同的因子間存在著許多交互作用，所以預測因子研究的第一個工作就是釐清各個變項間的關係並將其獨立出來。而一般常見的測試模型為邏輯斯迴歸(Logistic

regression)及Cox proportional hazards regression。近十年左右資料探勘(Data mining)技術廣為醫學領域研究者所應用，其中以類神經網路模型(Artificial neural networks)運用於癌症患者存活預測分析最為常見，因類神經網路除能準確的預測肺癌病患存活情形，也能應用在個人存活預測上。此外，決策樹(Decision tree)也是不錯的資料探勘模型之一，其優點在於它可以處理複雜的資料，且不受自變項間的影響，並將各自變項之於依變項的結果做一個清晰路徑分布呈現。故本研究將利用這些演算法(類神經網路模型、邏輯斯迴歸模型和決策樹模型)進行非小細胞肺癌存活預後因子的分析，期能藉由這些方法將可能影響非小細胞肺癌病患存活情形的因子分析出來。

故本研究分析美國癌症登記資料檔(The Surveillance, Epidemiology, and End Results, SEER)，期能找出影響非小細胞肺癌患者一年、三年和五年存活情形的預後因子，而本研究之所以限定非小細胞肺癌，主要是因為非小細胞肺癌占肺癌約有80%，而小細胞患者的存活情形、治療方式和分期方式都不同於非小細胞，故本研究限定非小細胞肺癌做為分析對象。此外，本研究亦期望能找出影響模型間表現的因子為何，故會進一步做模型間的比較，即使用準確度(Accuracy)、ROC曲線(Receiver Operating Characteristic curve)下的面積(AUC)和外推力來比較三個模型的效能和外推力(external generalization)，以找出最具有預測能力的模型。

### 材料與方法

#### 資料來源

本研究的資料來源為美國癌症登記資料檔(The Surveillance, Epidemiology, and End Results, SEER)，由SEER資料庫獲得1973年至2004年期間診斷之肺癌病例，該SEER資料庫共包含了481,432例呼吸道相關癌症個案，其中依國際疾病腫瘤分類第三版(ICD-O-3)定義腫瘤部位(Primary site)為C34.0-C34.9的肺癌個案共計有434,274例。依研究目的將研究對象以「死於肺癌」與「死於轉移癌」進行篩選，排除屍體解剖或死亡證明時才發現肺癌者[1-3]和病理組織學為小細胞及其他不常見的肺癌個案，此處的組織病理學挑選是參照Wang所做的研究[1]。

此外，因1988年前的資料有諸多缺漏，故將其剔除，研究對象篩選流程。該資料庫提供的變項大致可

以分為四個部份分別為組織病理學變項、人口統計學變項、臨床診斷治療變項和其他變項，經過整理後，進入分析的自變項共有 16 個，依變項為存活或死亡（一年、三年和五年存活情形）。

#### 預測模型

本研究是使用決策樹、類神經網路和邏輯斯迴歸模型進行分析，使用 SPSS 套裝軟體 Clementine 10.0 和 SPSS 14.0 版進行塑模，各模型的設定介紹如下：

#### 決策樹 (Decision, DT)

決策樹是一個利用分類和歸納的方式進行產生樹狀結構模式的方法之一，可以將輸入的資料做有效率分類，而樹上的每一個節點是為一判斷式，可以將輸入的資料根據特定的自變項去作分類，一直分類到資料無法再繼續分類下去即停止。而常見的決策樹種類以 C5.0 最為常見，主要是因為該模型可以同時應用在類別和連續性的資料上，故本研究採用 C5.0。其模型設定為簡單、偏愛「準確度」，預期雜訊(%)「0」。

#### 類神經網路(Artificial neural network, ANN)

本研究採用最為普遍應用的倒傳遞網路(Back-Propagation Network, BPN)，其架構包含輸入層，用以表現網路的輸入變項;隱藏層，用以表現輸入變項間的交互作用和影響，網路的隱藏層可以不只一層，也可以沒有隱藏層;輸出層，用以表示網路的輸出變項，及研究者所感興趣的依變項。其模型設定為快速，預防過度訓練為 80%，隨機種子為 1，停止於預設。

#### 邏輯斯迴歸 (Logistic regression model, LR)

在定量分析研究中，線性迴歸模型(Linear regression model)是最常用的統計方式，而臨床研究上問題的觀察常為類別資料，而邏輯斯迴歸模型就是針對依變項為二分類變項，如存活或死亡、轉移或未轉移等。本研究邏輯斯迴歸使用 SPSS14.0 進行塑模，方法選擇進入(enter)，且不將常數項納入方程式，因期能得到標準化係數。

#### 複合模型(hybrid models)

將決策樹挑選出來的變項放入類神經網路和邏輯斯迴歸模型即為複合模型，分別為類神經複合模型和邏輯斯迴歸複合模型。

#### 變項重要性比較

本研究變項的相對重要性排序，決策樹以各分支變項的頻率分數來獲得相對重要性順序;類神經網路則根據模型給予各變項的分數(相對重要性)做排序，分數越高者給予越前面的排序;邏輯迴歸模型遇類別變項時，則將各變項的標準化係數取絕對值之後平均，之後根據故變相的標準化係數大小作排序。而經過三個模型排序後，將各個模型所給予的序位加總起來取平均，所得的分數越低者，表示該變項的重要性越高。而上述模型中若遇到不被挑選的變項(決策樹)或是不顯著的變項(邏輯斯迴歸)則給予該變項最高序

位。

#### 模型評估

模型評估的指標有三，準確率(Accuracy, ACC)、ROC 曲線下的面積(area under the ROC, AUC)和外推力(external generalization)。準確率為正確預測實際存活或死亡狀態的機率值;AUC 即為 ROC 曲線下的面積，而該面積值越大即表示該模型的預測情形越好，一般而言，當 ROC 曲線下的面積越接近 1.0 時，表示該診斷的真實度越高，若接近 0.5，即表示該模型的真實性越差;而外推力是為訓練組的 ACC 和 AUC 值減去測試組的 ACC 和 AUC，該差值越小，表示該模型的外推能力越好，可以廣泛的應用於其他非訓練樣本的資料。此三個指標的判斷順序，先以 ACC，當 ACC 相同時再比較 AUC。

此外，模型中的訓練組和測試組人數設定各半，由電腦隨機分組;為避免隨機誤差影響結果，各在進行 SEER 資料庫分析時，本研究訓練組和測試組會隨機分組十次，此後將此十次的結果做平均。

#### 結果

##### 基本人口學描述分析

該資料庫經過整理和篩檢後，死於肺癌的肺癌患者共有 123,972，其年齡分佈大多在 65-74 歲 (34.5%)，其次為 45-64 歲 (32.8%) 大於(含)75 歲以上(30.1%);種族以白人居多(82.6%);婚姻狀況大多為已婚狀況(55.4%);性別的部分則以男性(59.9%)居多，女性(40.1%)，腫瘤部位大多為上肺葉 (46.9%)，組織病理學主要為腺細胞(37.8%)，其次為鱗狀細胞(24.8%);腫瘤分化程度除不詳外(48.3%)，大多為分化不佳(31.3%);腫瘤擴散情形主要為 distant(45.1%)，其次為 regional(34.8%);而腫瘤病變程度除少數 54 位個案為原位癌外，其餘皆為惡性腫瘤;側性則以右邊開始占多數(53.6%);手術方面，大多數未接受過手術(52.6%);放射治療部分，接受過放射治療有 48.6%，其次為沒有接受過放射治療 46.7%;而手術放療順序部分，大多為沒有接受過手術或放療(91.6%);淋巴結檢驗個數部分則有 79.8%的病患為無，而臨床分期大部分的個案為第四期(42.9%)，其次為第三期(29.4%);腫瘤大小的部份其平均值為 47.4mm(最小值為 4mm，最大值為 979mm)。

死於轉移癌的肺癌患者共有 7,285 位，該資料庫的患者年齡該資料庫的患者其年齡分布大多在 65-74 歲(33.9%)，其次為 45-64 歲(33.4%)及大於(含)75 歲以上(29.3%);種族以白人居多(81.6%);婚姻狀況大多為已婚狀況(53.2%);性別的部分則以男性(58.3%)居多，女性(41.7%)，腫瘤部位大多為上肺葉 (39.1%)，組織病理學主要為腺細胞(42.8%);腫瘤分化程度大多為不詳(52.6%)，其次為分化不佳(27.5%);腫瘤擴散情形主要為 distant(55%);側性則以右邊開始占多數(48.6%);手術方面，大多數未接受過手術(51.9%)，放射治療部分，沒有接受過放射治療有 59.4%;而手術放療順序部分，大多為沒有接受過手術或放療(91.1%);淋巴結檢驗個數部分則有 75.6%的病患為無，而臨床分期

大部分的個案為第四期(54%);腫瘤大小的平均值為40.13mm(最小值為4mm, 最大值為330mm)。

#### 模型變項挑選

死於肺癌一年存活情形決策樹挑出的前五名變項依序為手術種類、臨床分期、放射治療、檢驗淋巴結個數和性別;類神經網路篩選出的前五名變項依序為手術種類和臨床分期並列第一, 其次依序為腫瘤擴散程度、腫瘤大小、年齡分組和放射治療;邏輯斯迴歸篩選出的前五名變項為手術種類、腫瘤病變程度、臨床分期、年齡分組和腫瘤擴散程度。綜合三個模型變項排序的結果, 其前五名排序為手術種類、臨床分期、放射治療、腫瘤擴散程度和年齡分組;決策樹三年存活情形篩選出的變項前五名依序為手術種類、腫瘤擴散程度、臨床分期、腫瘤大小和手術放療順序;類神經網路篩選出的前五名變項依序為手術種類、臨床分期、腫瘤擴散程度、檢驗淋巴結個數和分化程度;邏輯斯迴歸篩選出的前五名變項為手術種類、腫瘤病變程度、手術放療順序、臨床分期和分化程度。綜合三個模型變項排序的結果, 其前五名排序為手術種類、臨床分期、腫瘤擴散程度、分化程度和手術放療順序;類神經網路五年存活情形篩選出的前五名變項依序為手術種類、臨床分期、年齡分組、腫瘤擴散程度和分化程度;邏輯斯迴歸篩選出的前五名變項為腫瘤病變程度、手術種類、手術放療順序、臨床分期和放射治療。綜合三個模型變項排序的結果, 其前五名排序為手術種類、臨床分期、年齡分組、分化程度和腫瘤擴散程度, 而決策樹模型於此階段無法成樹。

死於轉移癌的部份, 決策樹一年存活情形篩選出的變項前五名依序為臨床分期、腫瘤部位、手術種類、放射治療和檢驗淋巴結個數;類神經網路篩選出的前五名變項依序為手術種類、腫瘤擴散程度、腫瘤部位和檢驗淋巴結個數;邏輯斯迴歸篩選出的前五名變項為手術種類、臨床分期、腫瘤擴散情形、放射治療和性別。綜合三個模型變項排序的結果, 其前五名排序為, 臨床分期和手術種類並列第一, 其次依序為放射治療、腫瘤部位、腫瘤擴散程度和檢驗淋巴結個數;決策樹三年存活情形篩選出的變項前五名依序為臨床分期、手術種類、手術放療順序、腫瘤大小和分化程度;類神經網路篩選出的前五名變項依序為手術種類、臨床分期、腫瘤擴散程度、檢驗淋巴結個數和年齡分組;邏輯斯迴歸篩選出的前五名變項為手術種類和手術放療順序並列第一, 其次依序為臨床分期、放射治療、年齡分組和腫瘤大小。綜合三個模型變項排序的結果, 其前五名排序為手術種類、臨床分期、手術放療順序、腫瘤擴散情形和年齡分組(並列第四)及檢驗淋巴結個數;決策樹五年存活情形篩選出的變項前五名依序為手術種類、臨床分期、年齡分組、種族和檢驗淋巴結個數;類神經網路篩選出的前五名變項依序為手術種類、臨床分期、檢驗淋巴結個數、腫瘤擴散程度和年齡分組;邏輯斯迴歸篩選出的前五名變項為手術種類、年齡分組、臨床分期、放射治療和種族。綜合三個模型變項排序的結果, 其前五名排序為手術種類、臨床分期、年齡分組、種族和腫瘤部位(並列第四)及

檢驗淋巴結個數。

#### 模型評估

死於肺癌的部份, 由表 8 可知, 一年存活情形時其 ACC 表現類神經網路和決策樹的表現相當, 進而比較 AUC, 結果以類神經網路的表現最好。外推能力則是以邏輯斯迴歸模型表現最好。所以一年存活情形的模型以類神經網路表現較好, 外推能力以邏輯斯迴歸模型的表現最差好;三年存活情形, 就 ACC 來看, 決策樹和類神經網路準確率相同, 且優於邏輯斯迴歸, 因此比較 AUC, 結果類神經網路優於決策樹, 因此就三年存活情形的模型以類神經網路最好, 而外推能力則以邏輯斯迴歸模型的表現最好;五年存活情形, 就 ACC 而言, 以類神經網路的表現較邏輯好, 而就 AUC 而言, 則是以邏輯斯迴歸較好;外推力則是以類神經網路的表現較好。

將決策樹所挑選出來的變項分別放入類神經網路和邏輯斯迴歸, 其一年存活情形的結果可知, 就 ACC 而言, 複合模型並未增加預測力, 但 AUC 的部份, ANN 和 LR 都有增加;其三年存活情形結果和單模型比較起來, ANN 複合模型的表現較單模型時差, 而邏輯斯複合模型的 AUC 則略較單模型時高;外推力的部份, 複合模型的外推力都較單模型時好。

死於轉移癌的部份, 一年存活情形, 由表 11 可知, 就 ACC 而言, 類神經網路和決策樹的表現一樣好, 但以類神經網路的 AUC 較決策樹好, 所以一年存活情形以類神經網路表現最好, 而外推能力部份, 就 $\Delta$ ACC 而言, 以邏輯斯迴歸表現最好, 決策樹最差;就 $\Delta$ AUC 而言, 則以邏輯斯迴歸表現最不好, 決策樹和類神經網路一樣好, 故整體來說, 以類神經網路的表現較好;三年存活情形, 就 ACC 而言, 決策樹的表現較好, 而外推能力的表現則是以類神經網路的表現最好;五年存活情形, 決策樹和類神經網路的 ACC 相同, 而 AUC 則以類神經網路表現較佳, 且外推力也以類神經網路較優, 所以五年存活情形以類神經網路的表現最好。

複合模型的部份, 由表 12 可知, 一年存活情形, 就 ACC 而言, ANN 和 LR 複合模型都較單模型預測力好, 就 AUC 而言, ANN 複合模型的預測力並未提升, LR 複合模型卻提升;而外推力的部份, 整體而言, 複合模型的外推能力都較單模型時好;三年存活情形, 由表 13 可知和單模型比較起來, ANN 和 LR 複合模型的表現都較單模型時好, 外推力複合模型的表現亦較單模型時好。

上述的複合模型分析皆不進行五年存活情形分析, 主要是因為五年存活情形時死亡率較高, 且死於肺癌的五年存活情形決策樹無法成樹, 故不做五年存活情形的複合模型。

#### 討論

從上述的結果中不難發現, 三個模型以類神經網路所挑選出來的變項及其排序和綜合變項排序較相似, 而三個模型對於變項的挑選和排序卻都有所不同, 而此結果和 Jefferson 的研究<sup>4</sup>與 Hanai 的研究[5]結果相近, 結果都顯示類神經網路和邏輯斯迴歸挑選

變項的順序都大不相同。而兩個模型的變項順序不同可能是因為受到變項間的交互作用所導致的，因為相較於邏輯斯迴歸，類神經網路完全不會受到變項的交互作用的影響，但是邏輯斯迴歸模型則會受到影響，因此才會出現不一致的排序。因此對於多樣化的醫學資料庫，則可以選擇類神經網路模型來處理這類多變化的資料[6]。除了變項間的相關性外，模型間不同的演算方式或許也是導致這樣的結果的原因之一。此外，在 Jefferson 和 Hanai 的研究中，類神經網路和邏輯斯迴歸的結果變項排序前三名幾乎都有臨床分期，即便該兩位學者放入模型的變項都與本研究不相同，但是結果都顯示臨床分期是非常重要的變項。

模型預測力的表現部分，整體來說以類神經網路的表現較好較穩定，而以邏輯斯迴歸模型的表現較差較不穩定，但外推力的部份則是以邏輯斯迴歸模型表現最好。Jefferson 對 620 位接受手術的非小細胞肺癌患者使用類神經網路和邏輯斯迴歸來預測存活情形並比較模型，其結果發現不論在哪個時間點上(6 個月、12 個月、18 個月和 24 個月)，類神經網路的表現都較邏輯斯迴歸的表現好[4]; Hanai 學者對 125 位接受過治癒性手術的非小細胞肺癌患者使用邏輯斯迴歸和類神經網路去預測存活情形，並比較兩個模型，結果顯示一年、三年和五年存活情形，除了一年存活情形是邏輯斯迴歸的表現較好外，其餘都是類神經網路的準確率都較邏輯斯迴歸高<sup>5</sup>。由上述文獻和本文研究可知，類神經網路在許多的情況下表現都較邏輯氏好。Long 使用邏輯斯迴歸和決策樹來預測是否會有急性心臟局部缺血，結果顯示雖然兩個模型的表現都不錯，但邏輯斯迴歸模型還是表現較決策樹好[7]; Rudolfer 使用邏輯斯迴歸模型和決策樹模型來診斷腕道症候群並比較兩個模型，結果顯示兩個模型表現一樣好[8]; Abu-Hanna 和 de Keizer 使用傳統工具邏輯斯迴歸模型和合成模型決策樹來預測加護病房病患的預後情形，其結果顯示決策樹的表現較邏輯氏模型好[9]; Schwarzer 使用模糊推論(Fuzzy inference)、決策樹和邏輯斯迴歸模型去預測舌頭癌患者的頸淋巴結轉移情形，其結果顯示，模糊理論和決策樹的精確率較高[10]。Delen 使用決策樹、類神經網路和邏輯斯迴歸模型三種資料探來預測乳癌患者的存活情形，其結果顯示決策樹模型的表現最好(精確率為 93.6%)，其次為類神經網路(精確率為 91.2%)，邏輯斯迴歸模型表現最差(精確率為 89.2%)[11]。從上述的文獻和本研究的結果發現，以類神經網路的表現好的次數多於另外兩個模型。

雖然結果顯示三個模型以類神經網路的表現較好，但決策樹、類神經網路和邏輯斯迴歸模型是各具特色和優點，如決策樹的分類過程很清晰，不像類神經網路有所謂的黑盒子(隱藏層)，可以很清楚的知道輸入變項是怎麼被使用和分類;而類神經網路雖然無法得知黑盒子過程，但是這段黑盒子過程卻學習力佳，可以接受很複雜的資料庫;而邏輯斯迴歸模型則可以讓研究者清楚的知道每個變項在控制其他的自變項後對依變項的影響。因此每個模型都有其專屬特色，應該在特定的狀況下選擇適合的模型進行分析。

複合模型方面，本研究的結果發現，並非所有的

複合模型其預測能力表現都較單模型時來得好，而外推力整體來說是有提升的，此結果和許意絃的研究結果大致相同，該研究顯示並非所有的複合模型的表現皆優於單模型的表現，但外推力的部分則有增加的現象。此外該作者提出若決策樹模型和類神經網路的準確率近似時，其複合模型預測能力之準確度不會提高反而會降低[12]。而本研究的結果卻顯示，當決策樹和類神經網路的 ACC 相近時，其 AUC 反而會較單模型時低，反之，當決策樹的準確率較邏輯斯迴歸模型高時，其 AUC 就會較單模型時高。雖然本研究的結果不同於許意絃所提出的結論，但是結果都表示，當決策樹的表現較另外兩個模型表現好時，若要提升預測力是可以採用複合模型的方式，此外，Dwinnell 也認為決策樹 CART 可以當作是執行類神經網路前的變項挑選工具，或是 pre-processor，因為如此一來可以增加類神經網路的準確度，並縮短類神經網路訓練的時間<sup>13</sup>。

## 結論

影響肺癌患者死於肺癌的一年、三年和五年存活情形的預後因子主要為手術種類、臨床分期、年齡分組和腫瘤擴散程度等因子。而影響肺癌患者死於轉移癌的一年、三年和五年存活情形的預後因子主要為手術種類、臨床分期和檢驗淋巴結個數。

整體來說，三個模型的預測能力以類神經網路模型的表現最好，而外推能力以邏輯斯迴歸模型的表現較好;而複合模型的部份，當該模型的 ACC 與決策樹的 ACC 值相近時，該複合模型 AUC 值不會提升，可能下降，反之若決策樹 ACC 值表現較好時，則會提升該複合模型的 AUC 值。死因部分，整體而言以死於轉移癌的模型的 AUC 表現較好，原因可能是因為輸入的變項比較能解釋肺癌患者死於轉移癌的關係。

## 研究限制

SEER 資料庫雖然是個 SEER 資料庫雖然是個很龐大且完整的資料庫，但可惜該資料庫所包含的變項中未納入化學治療相關變項，此外影響肺癌病患預後的因素除本研究所討論的病理和臨床變項外，尚有一些生物標誌(Biological marker)，如 P27 和 P53 標誌[5, 14]等，是肺癌很重要的預後因子之一，若可以將本研究發現的重要預測因子和上述的生物標誌都納入研究分析，相信定可以增加模型的預測能力。

## 未來展望

期本研究的結果可供台灣國健局在建構台灣癌症資料庫時一個參考範本，除納入 SEER 所納入的變項外，亦可加入化學治療和基因相關的變項，供相關研究學者分析討論，以促進台灣肺癌國人的存活品質。此外，因 COX 迴歸分析方法亦是許多學者在探討存活分析時常用到的分析方法之一，故建議未來研究可將 COX 迴歸模型納入分析比較。

## 參考文獻

- [1] Wang SJ, Fuller CD, Thomas CR, Jr. Ethnic disparities in conditional survival of patients with

- non-small cell lung cancer. *J Thorac Oncol.* 2007; **2**: 180-90.
- [2] Chen KY, Chang CH, Yu CJ, Kuo SH, Yang PC. Distribution according to histologic type and outcome by gender and age group in Taiwanese patients with lung carcinoma. *Cancer.* 2005; **103**: 2566-74.
- [3] Skuladottir H, Olsen JH. Conditional survival of patients with the four major histologic subgroups of lung cancer in Denmark. *J Clin Oncol.* 2003; **21**: 3035-40.
- [4] Jefferson MF, Pendleton N, Lucas SB, Horan MA. Comparison of a genetic algorithm neural network with logistic regression for predicting outcome after surgery for patients with nonsmall cell lung carcinoma. *Cancer.* 1997; **79**: 1338-42.
- [5] Hanai T, Yatabe Y, Nakayama Y, Takahashi T, Honda H, Mitsudomi T, *et al.* Prognostic models in patients with non-small-cell lung cancer using artificial neural networks in comparison with logistic regression. *Cancer Sci.* 2003; **94**: 473-7.
- [6] Rodvold DM, McLeod DG, Brandt JM, Snow PB, Murphy GP. Introduction to artificial neural networks for physicians: taking the lid off the black box. *Prostate.* 2001; **46**: 39-44.
- [7] Long WJ, Griffith JL, Selker HP, D'Agostino RB. A comparison of logistic regression to decision-tree induction in a medical domain. *Comput Biomed Res.* 1993; **26**: 74-97.
- [8] Rudolfer SM, Paliouras G, Peers IS. A comparison of logistic regression to decision tree induction in the diagnosis of carpal tunnel syndrome. *Comput Biomed Res.* 1999; **32**: 391-414.
- [9] Abu-Hanna A, de Keizer N. Integrating classification trees with local logistic regression in Intensive Care prognosis. *Artif Intell Med.* 2003; **29**: 5-23.
- [10] Schwarzer G, Nagata T, Mattern D, Schmelzeisen R, Schumacher M. Comparison of fuzzy inference, logistic regression, and classification trees (CART). Prediction of cervical lymph node metastasis in carcinoma of the tongue. *Methods Inf Med.* 2003; **42**: 572-7.
- [11] Delen D, Walker G, Kadam A. Predicting breast cancer survivability: a comparison of three data mining methods. *Artif Intell Med.* 2005; **34**: 113-27.
- [12] 許意絃. 輔助冠狀動脈心臟病診斷之邏輯斯迴歸、決策樹、類神經網路及複合分析模型效能比較. 公共衛生研究所 流行病學組 國防醫學院. 2007.
- [13] Dwinnell W. Data enhancement — filtering for neural network success. *PC AI.* 2000; **14**: 20-23.
- [14] Hsia TC, Chiang HC, Chiang D, Hang LW, Tsai FJ, Chen WC. Prediction of survival in surgical unresectable lung cancer by artificial neural networks including genetic polymorphisms and clinical parameters. *J Clin Lab Anal.* 2003; **17**: 229-34.