

Data Mining
– Midterm, April 21, 2016 –

1. (30%) Answer the following questions with respect to matrix decomposition.
 - (a) Compare/contrast the SVD and NFM methods.
 - (b) Decompose a matrix $[[5,3,0,1],[4,0,0,1],[1,1,0,5],[1,0,0,4]]$ into $[[2.21, 0.13], [1.76, 0.18], [0.32, 1.94], [0.36, 1.55]] * [[2.24, 1.33, 0.56, 0.31], [0.13, 0.3, 0.12, 2.5]]$ by the NFM. What are the meanings of the rows or columns in these three matrices?
 - (c) Map both vectors $[1,1,0,0]$ and $[1,0,1,0]$ into the transformed space by using the matrices shown in (b)?
2. (30%) Given the following transactions $T_1=\{a,b,c,d,e\}$, $T_2=\{b,c,d,e\}$, $T_3=\{a,b,c,d\}$, $T_4=\{a,b,c\}$ and $T_5=\{a,b\}$, where $Sup_{min}=40\%$ and $Conf_{min}=60\%$,
 - (a) Mine all frequent patterns.
 - (b) Find all closed frequent patterns.
 - (c) Find all association rules whose confidence levels are between 60% and 80%.
3. (20%) Answer the following questions with respect to discretization.
 - (a) What is the entropy-based discretization method? What is the natural partitioning method?
 - (b) Compare/contrast both methods.
4. (20%) Predict if there is a game by the Bayesian classification method by using the following dataset.
 - (a) when temperature is hot, humidity is low and wind is light.
 - (b) when temperature is cold and wind is strong.

Temperature	Humidity	Wind	Rain	Game
hot	low	light	yes	no
mild	low	strong	no	yes
cold	low	light	no	yes
cold	high	strong	yes	no
hot	high	light	no	yes

國立臺灣大學期中學期考試答案卷

National Taiwan University Midterm/Final Examination Answer Sheet

課程編號

Course no.

科目 資料探勘

管理

學院 資管所

學系

考試日期 2017 年 4 月 20 日

學號

姓名

Date 2017 Y M D Student ID no. _____ Name _____

從此處開始寫起。試卷用紙務須節用。非經主試認可不得續用其他紙張作答。

Please write from here.

記分	Score	教師簽名或蓋章 Lecturer's signature
----	-------	---------------------------------

PD+2

組 一 年級 Year

記分欄

$$1. (a)$$

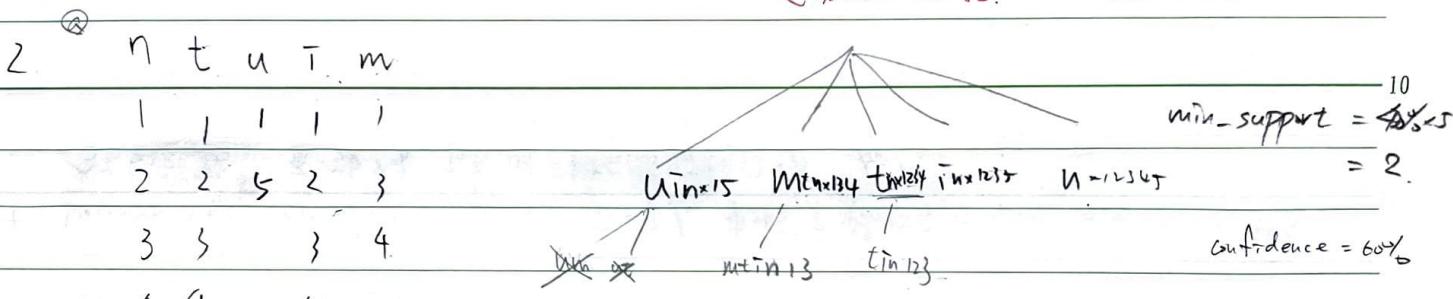
$$\text{② } S = \begin{pmatrix} 1, -1, 2, -2, -3, 3, -4, 4 \end{pmatrix}$$

$$L(x) = \sqrt{S[x] + S[x+1]}$$

$$H(x) = \sqrt{S[G_x] - S[G_{x+1}]}$$

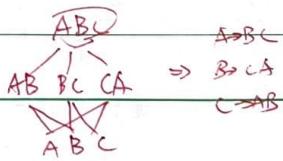
$$\Rightarrow (0, 0, 0, 0, \sqrt{2}, 2\sqrt{2}, -3\sqrt{2}, -4\sqrt{2})$$

(b) 我們可以觀察到 x' 為 x , x'' 為 x' 的近似。但是當做一次轉換時，因為是將每面添加加減去，有資訊丟失。
 故會使得局部的變化縮小。像回小題原本可能有 $(1, -1), (2, -2)$ 的振幅，透過 Haar wavelet transform 後，只剩下 0.少了震盪的幅度。而因為他是每個一組奇偶各一組，所以做 HWT 後的近似值 LPF 仍可以保有 sequence 的變化。(振盪平衡是振盪的情況)，轉換後長度、距離、內積均相同。(LPF 長度)
 趨勢的線性會保留。(這個成長度)



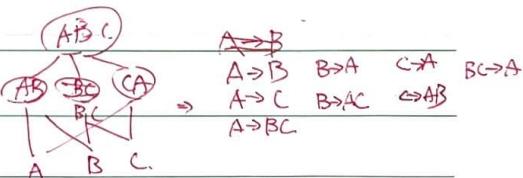
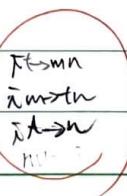
closed patterns: n, in, tn, tin, mtn, uin, utin.

$n \rightarrow \bar{n}$	$\frac{2}{5} = 40\%$	$\bar{n} \rightarrow n$	$\frac{2}{4} = 50\%$	$t \rightarrow n$	$\frac{2}{4} = 50\%$	$m \rightarrow tn$	$\frac{1}{5} = 20\%$	$u \rightarrow \bar{in}$	$\frac{2}{5} = 40\%$
$n \rightarrow t$	$\frac{2}{5} = 40\%$	$\bar{n} \rightarrow tn$	$\frac{2}{4} = 50\%$	$t \rightarrow \bar{in}$	$\frac{2}{4} = 50\%$	$m \rightarrow tin$	$\frac{2}{5} = 40\%$	$u \rightarrow \bar{tn}$	$\frac{2}{5} = 40\%$
$n \rightarrow mt$	$\frac{3}{5} = 60\%$	$\bar{n} \rightarrow \bar{mt}$	$\frac{2}{4} = 50\%$	$t \rightarrow mn$	$\frac{2}{4} = 50\%$	$m \rightarrow \bar{tn}$	$\frac{2}{5} = 40\%$	$u \rightarrow \bar{mt}$	$\frac{2}{5} = 40\%$
$n \rightarrow \bar{u}\bar{t}$	$\frac{2}{5} = 40\%$	$\bar{n} \rightarrow u\bar{t}$	$\frac{2}{4} = 50\%$	$t \rightarrow \bar{mn}$	$\frac{2}{4} = 50\%$	$m \rightarrow \bar{tn}$	$\frac{2}{5} = 40\%$	$u \rightarrow \bar{u}\bar{t}$	$\frac{2}{5} = 40\%$
$n \rightarrow \bar{t}\bar{t}$	$\frac{2}{5} = 40\%$	$\bar{n} \rightarrow \bar{t}\bar{t}$	$\frac{2}{4} = 50\%$	$t \rightarrow \bar{mn}$	$\frac{2}{4} = 50\%$	$m \rightarrow \bar{tn}$	$\frac{2}{5} = 40\%$	$u \rightarrow \bar{u}\bar{t}$	$\frac{2}{5} = 40\%$



Association rules.

$n \rightarrow \bar{n}$ $\bar{n} \rightarrow n$ $t \rightarrow n$ $m \rightarrow tn$ $u \rightarrow \bar{in}$.
 $n \rightarrow t$ $\rightarrow \bar{n} \rightarrow tn$ $t \rightarrow \bar{in}$ $m \rightarrow tin$
 $n \rightarrow mt$ $\rightarrow \bar{n} \rightarrow \bar{mt}$ $t \rightarrow mn$
 $n \rightarrow \bar{t}\bar{t}$



① 第一層只做 closed.

② 往上做 $X \rightarrow Y$, X 只找非 closed pattern.

3. Apriori 是一個来找 frequent pattern 的 algorithm, 這個演算法的最大原則是 "如果一個 itemset 不是 frequent, 則它的 super set 也不是 frequent", 所以從最小的長度起, 慢慢往上推進, 可以避免掉必須 traverse >1 種組合的情況, e.g. AB 不是 frequent, 則 ABC, ABD... 都不需要考慮。

Partition 的做法是首先將 data set 分塊, 分別找出在各個 partition 裡的 local frequent pattern, 然後這些 local frequent pattern 檢驗是否為 global frequent. 其原理是 "當一個 set 為 global frequent, 則它必為某個 partition 的 local pattern", 反之亦然。假設全部有 N 筆, $Sup_{min} = \alpha\%$, 假設平均 10 筆每個 N 筆, 假設存在一個 set 皆為 local, 但為 global 則 $(\frac{N}{10} \times \alpha\% - 1) \times 10 > N \times \alpha\%$, 需成立。但明顯不行。而檢查 local 是否為 global 的效率相對 partition 之後好, 因為某種 set 都分在一區並非 local, 但其實 global 時, 不為 frequent.

至於這兩者的最大差別 在於計算上的速度。Apriori 從最小長度一直往上算, 每次計算就要 traverse 整個 data set, 計算成本很高。然而, partition 只要一开始 partition 就 traverse 一次, 之後再做一次 traverse 檢查 local 是否為 global, 總共只需兩次。計算 global 時 candidate 量會暴增。

4. Decision tree: 可以利用 Entropy 或是 Gini Index 來找到最代表性的 attribute 作為 attribute 做為樹的分支的節點, 一直往下分, 直到將資料分離的夠仔細。divide and conquer. top-down.

Naive Bayesian: 利用機率模型計算各個結果的可能性, 最後根據最大的機率來指定分類。

$$P(H|X) = \frac{P(X|H)P(H)}{P(X)}, + \text{比較每個 class } H \rightarrow P(X|H)P(H)$$

Naive Bayesian 利用機率模型計算, 有理論基礎, 計算也容易。但是他的前提是假說是各 attribute 之間是獨立的。e.g. 年紀, 薪資, 這個假說其實是有問題的, 因為 attribute 之間可能會有 dependency。而 decision tree 跟然做法有點像 greedy 找最好, 但是沒有額外的假說前提。

Decision tree 需要看過資料。Naive 不需要資料, 處理 sparse。

Decision tree 可以產生 tree rule.

Decision tree 1st entropy. Time complexity ↗

Naive. Incremental

5. 公司可以收集兩大類型的數據加以分析。第一大類是每個 user 對於各 video 的評分。第二大類是收集 user 的 information 及 user 所瀏覽過的 video list，以下分別對這兩類資料提供 mining method 及可應用的方式。

① 評分：我們可以在收集到評分(0~5分)時，利用 NMF 的方式，拆成兩個矩陣逼近希望得到兩種結果。⁵使用者對未看過的 video 的預期評價。使用者對於新 feature (video 類型。e.g. 搞笑、動作...) 的偏好程度。第二個結果可以應用在推薦系統上面，藉由推薦 user 可能已喜歡的 video 來留住 user。第二個可以幫助，當公司想要生產劇時能夠觀察 user 的取向來針對取向設計劇情內容。

② user information 及瀏覽紀錄。可以利用 naive Bayesian 朴素模型的分析，得到每個背景的使用者 (e.g. 年紀、性別) 會偏好哪種類型的劇。這個分析結果可以運用在將廣告引入系統時。雖然在使用者方面，公司可能需要提供免費或是降低費用的權益。但是因為可以知道觀看此 video 的使用者族群背景，所以可能向外找尋廣告投放。因為有精準的 TA，更可以吸引廣告主在此平台上投放廣告。

太遠在執行階段，可以多描述一些在 pre-process 時的工作。

Concept hierarchy 來做資料分群。

資料離散化。

15

20

25

30

Data Mining

– Midterm, April 26, 2018 –

1. (30%) Consider a book review matrix R , where each entry represents the review score for each book, and the review score of each user is denoted as one row. Decompose R into three matrices by SVD, $U=\{\{-0.4, 0.6, -0.2, 0.1, -0.7\}, \{-0.4, -0.2, -0.8, 0.3, 0.3\}, \{-0.5, 0.4, 0.2, -0.5, 0.5\}, \{-0.6, -0.6, 0.2, -0.3, -0.4\}, \{-0.4, 0, 0.5, 0.7, 0.2\}\}$, $\Lambda=\{\{2.5, 0, 0, 0, 0\}, \{0, 0.8, 0, 0, 0\}, \{0, 0, 0.7, 0, 0\}, \{0, 0, 0, 0.4, 0\}, \{0, 0, 0, 0, 0.25\}\}$, $V^T=\{\{-0.4, -0.2, -0.5, -0.7, -0.3\}, \{-0.6, 0, -0.2, 0.2, 0.7\}, \{-0.4, -0.5, 0.8, -0.2, 0\}, \{-0.5, 0, 0, 0.6, -0.6\}, \{-0.2, 0.9, 0.3, -0.3, -0.1\}\}$
 - (a) What are the axes of the new data space if the dataset is transformed into a three-dimensional space?
 - (b) How/what is the original data transformed into a three-dimensional space?
 - (c) Transform $(0, 1, 1, 0, 1)$ into a three-dimensional space.
2. (20%) Given the following closed patterns J:5, JU:4, JUN:3 and JUNE:2, where the number after the colon denotes the support of the closed pattern,
 - (a) restore all frequent patterns and their supports.
 - (b) Find all association rules, where the minimum confidence threshold is 80%.
3. (20%) Given the following sequence database $\{\langle \text{MAY} \rangle, \langle \text{M(MA)(MAY)A} \rangle, \langle \text{M(MAY)A} \rangle\}$,
 - (a) find all sequential patterns, where the minimum support threshold is 60%.
 - (b) Find all association rules, where the minimum confidence threshold is 80%.
4. (30%) Answer the following questions with respect to an online food ordering and delivery platform (such as Grubhub, UberEATS), which provides an online and mobile platform for restaurant pick-up and delivery orders. Describe the content of data that might be collected. Then, design a framework to mine some business insights from the dataset, and discuss the expected results and possible value-added services.

Data Mining

– Midterm, April 25, 2019 –

1. (20%) How does SVD work? How is it used to reduce the number of dimensions? What are the ~~axes~~^轴 of the reduced space? How is a ~~new~~^新 data record transformed into the reduced space?
2. (20%) Design a neural network to classify customers into 3 classes, where each customer profile contains 20 attributes. What will you concern when choosing the number of hidden layers and the number of nodes in each hidden layer?
^{設計 NN}
3. (20%) Why are the sigmoid and hyperbolic tangent functions often used as the activation function in a neural network? Compare/contrast both functions.
 $\Sigma / 2 \tanh$
4. (20%) What is LSTM? What is each function in LSTM used for?
^{設計原理}
5. (20%) What is the philosophy of devising SVM? How does SVM resolve a non-linear classification problem?
^{kernel func}

Data Mining

– Midterm, April 29, 2021 –

1. (30%) Answer the following questions with respect to the non-negative matrix factorization (NMF) method, where an $n \times m$ matrix V is decomposed into WH , each column of V denotes a record of data, W is an $n \times k$ matrix and H is a $k \times m$ matrix.
 - (a) What do the columns (or rows) of W and H stand for?
 - (b) How do you apply the NMF to deal with the data sparsity?
 - (c) How do you obtain the representation of a new vector $X^T = (x_1, x_2, \dots, x_n)$ in the transformed space?
2. (30%) Answer the following questions with respect to the decision trees.
 - (a) Why is the ID3 decision tree method biased towards a multi-level attribute in attribute selection?
 - (b) How can the C4.5 decision tree method resolve the problem described in 2(a)?
 - (c) Why does the C4.5 prefer the attribute with unbalanced splits?
3. (40%) Answer the following questions with respect to the neural networks.
 - (a) (10%) Why is the sigmoid or hyperbolic tangent function frequently used in a neural network model?
 - (b) (10%) When will the ReLU activation function be used? Why?
 - (c) (20%) How are the biases and weights updated in each iteration of the training process of a multi-layer feed forward model if the hyperbolic tangent function is used?

SVD

V

$$\begin{bmatrix} 0.4, 0.6, -0.2, 0.1, -0.7 \\ -0.4, -0.2, -0.8, 0.3, 0.3 \\ 0.5, 0.4, 0.2, -0.5, 0.5 \\ 0.6, -0.6, 0.2, -0.3, -0.4 \\ -0.4, 0, 0.5, 0.7, 0.2 \end{bmatrix}$$

3

A

$$\begin{bmatrix} 2.5, 0, 0, 0, 0 \\ 0, 0.8, 0, 0, 0 \\ 0, 0, 0.7, 0, 0 \\ 0, 0, 0, 0.4, 0 \\ 0, 0, 0, 0, 0.25 \end{bmatrix}$$

3

V

$$\begin{bmatrix} 0.4, -0.2, -0.5, -0.7, -0.3 \\ -0.6, 0, -0.2, 0.2, 0.1 \\ 0.4, -0.5, 0.8, -0.2, 0 \\ -0.5, 0, 0, 0.6, -0.6 \\ -0.2, 0.9, 0.3, -0.3, -0.1 \end{bmatrix}$$

5x5

5x3

3

