

SAS Enterprise Miner (SAS EM) 開始最重要的四步驟
by TKU IM EMBA Students (2012)
Version 2012.04.13

SAS_EM 開始最重要四步驟：

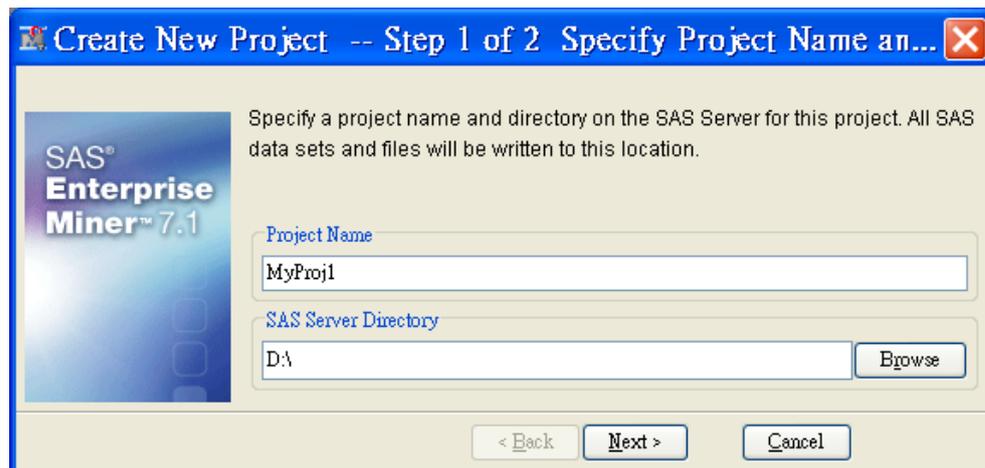
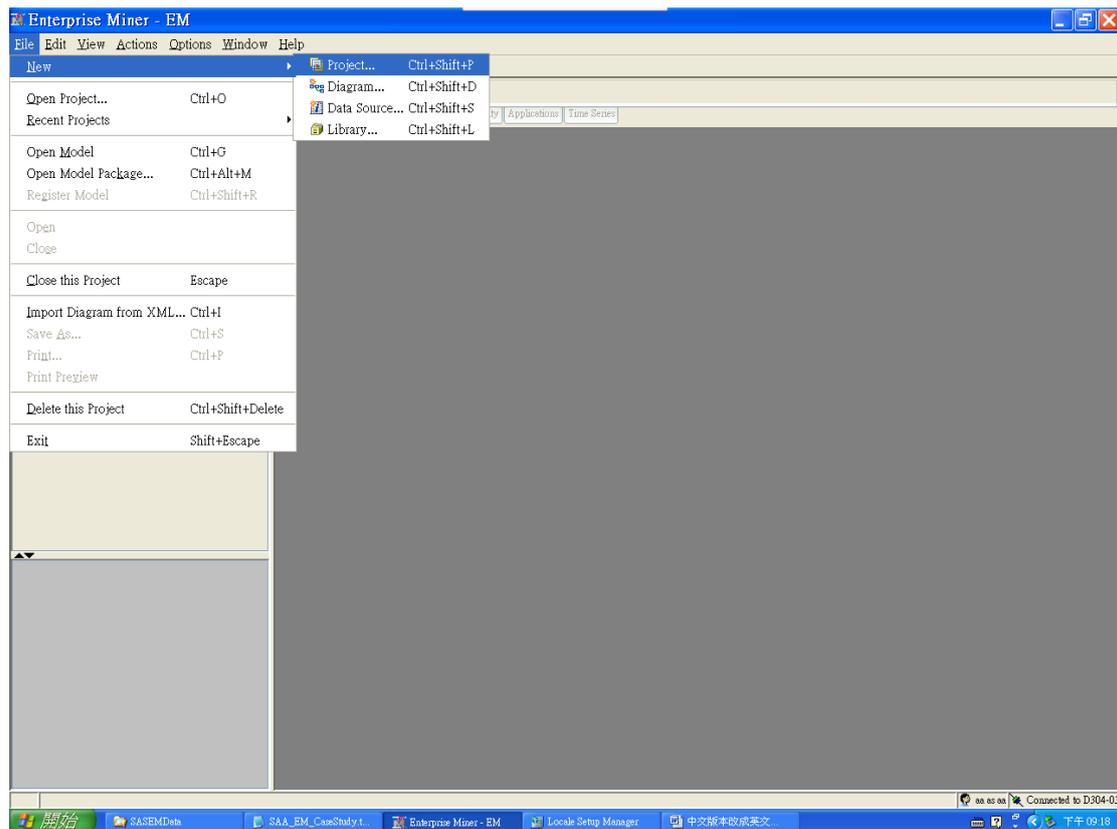
Step 1. Create Project (建 project)

Step 2. Define Library (建 Library)

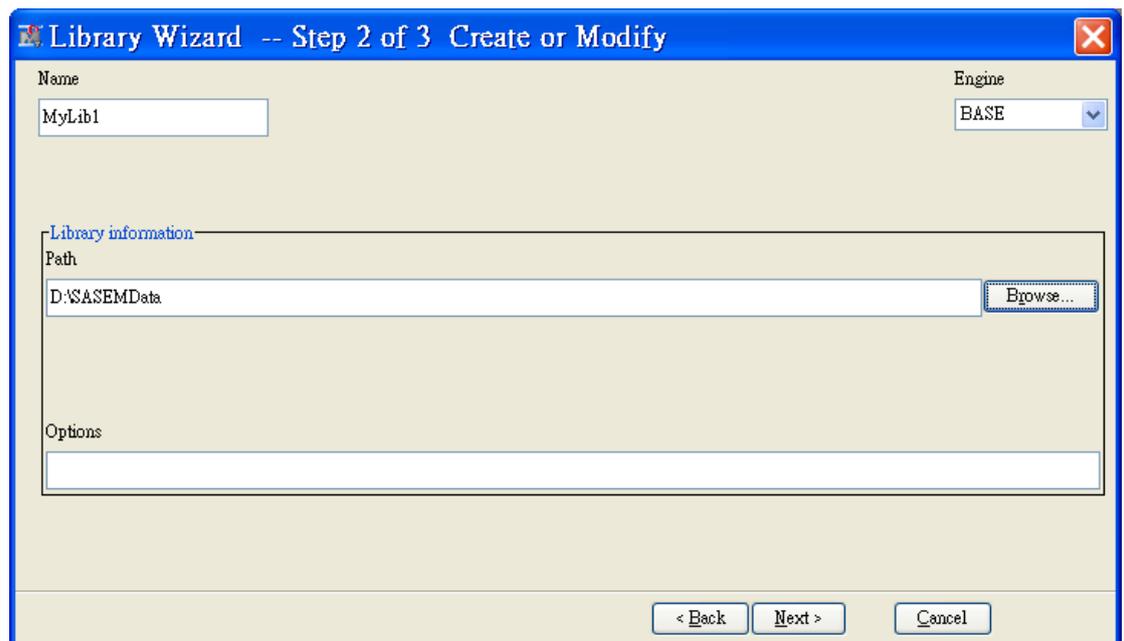
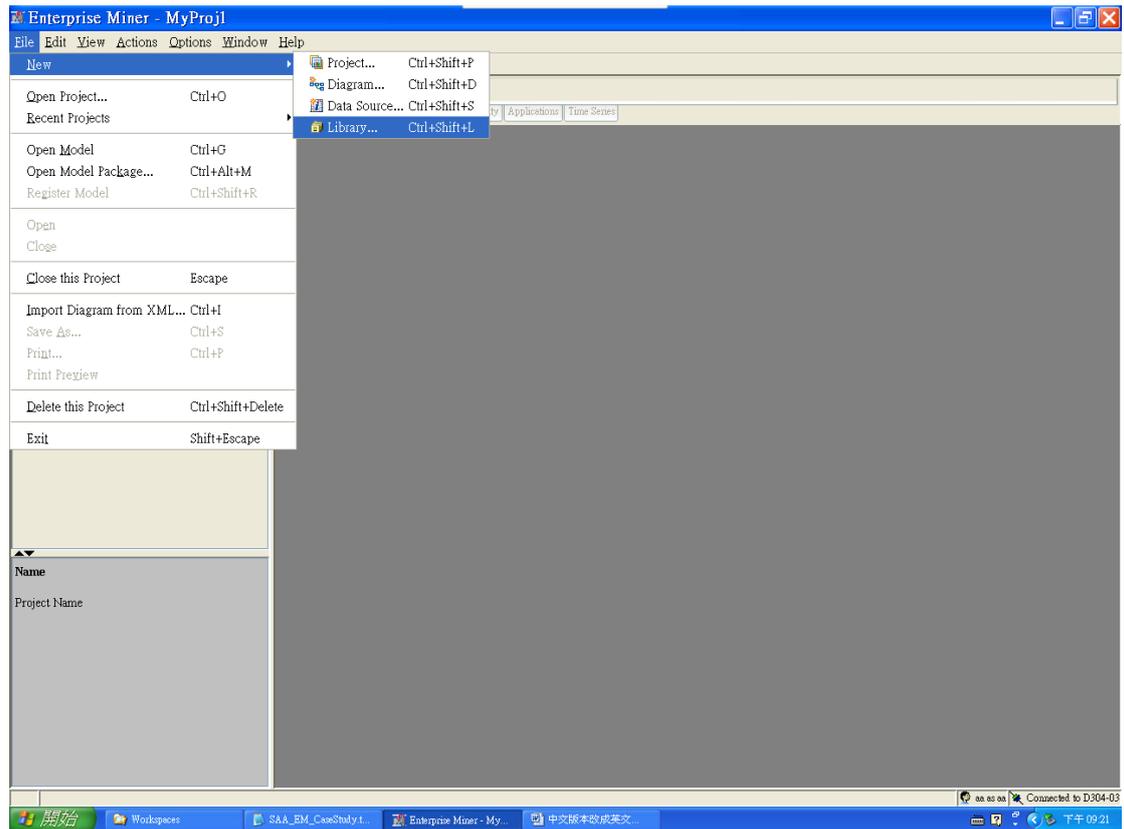
Step 3. New Data Source (Input Data Set)(指定 db source)

Step 4. Create Diagram (Diagram 流程圖)

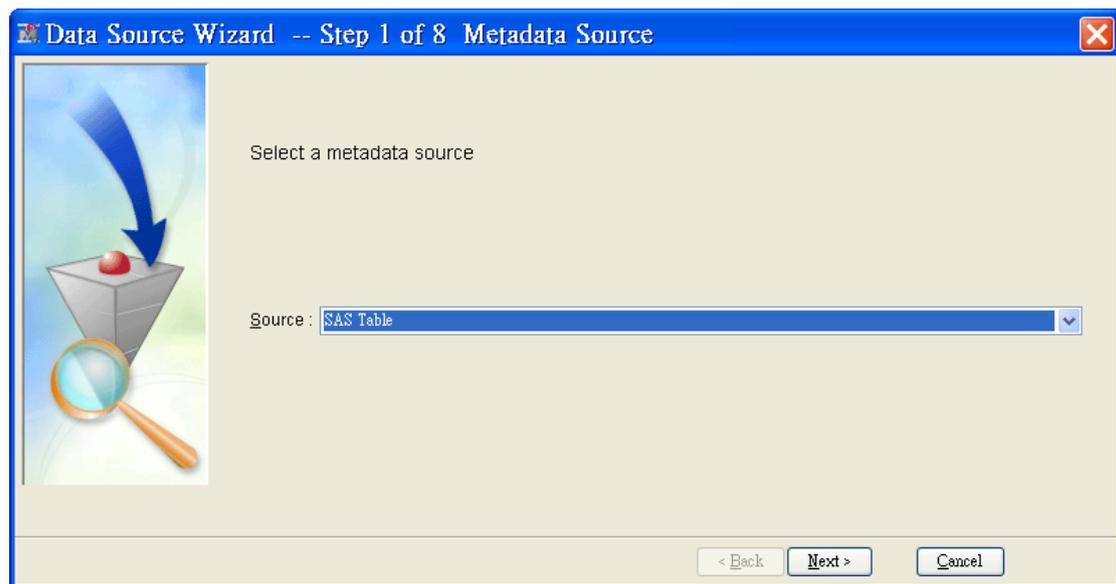
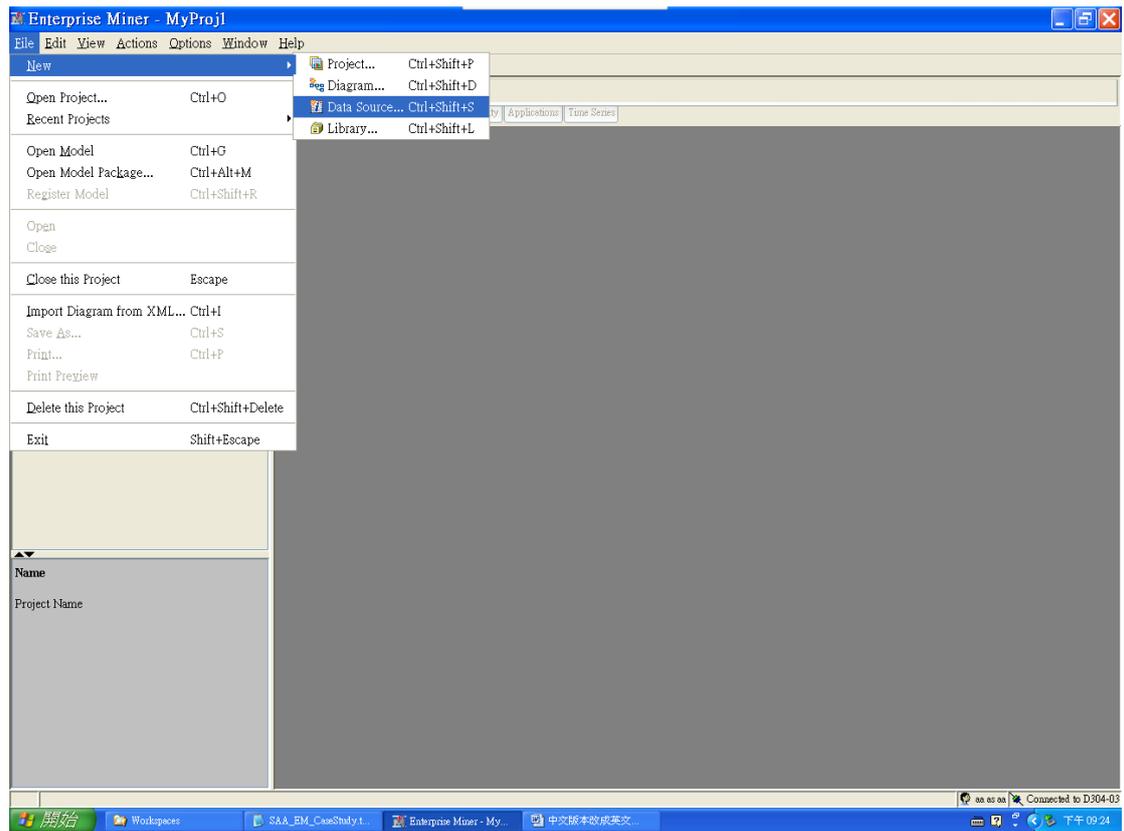
1. 建 project

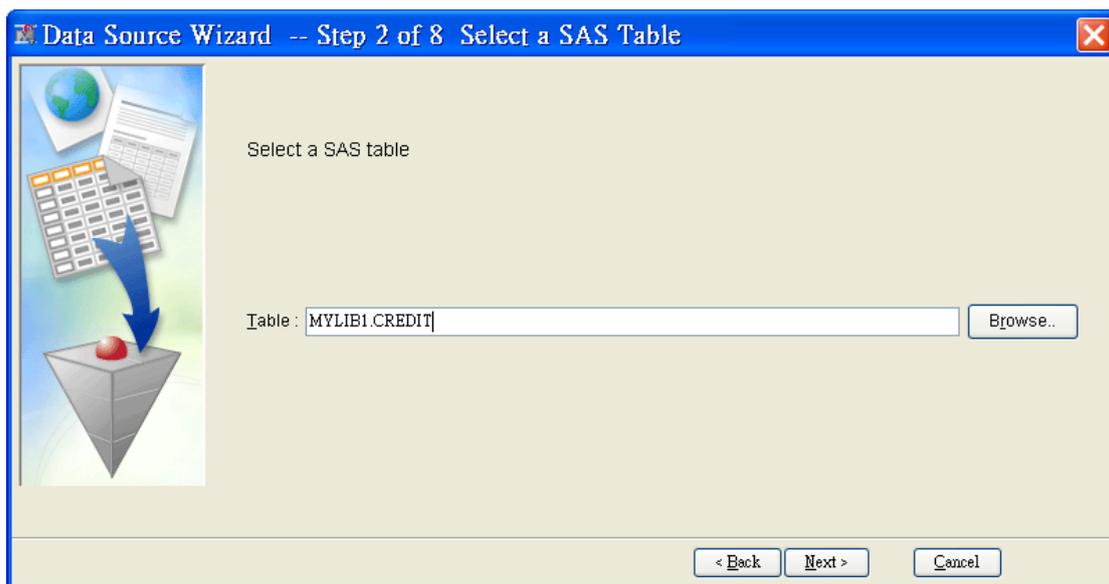
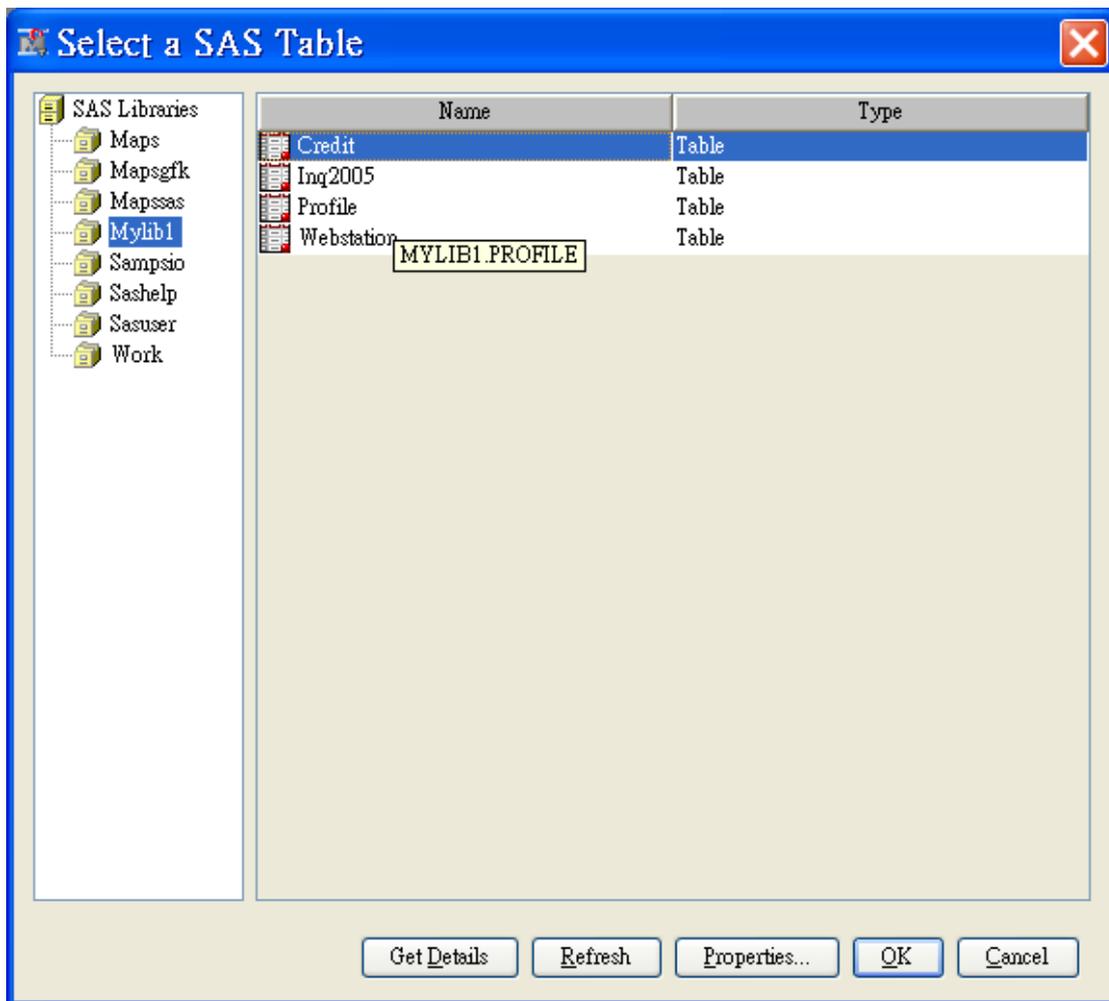


2. 建 Library (db) 以 table 放在資料夾的(.sas7bdat)
請指定資料夾，非資料夾內的檔案 (D:\SASEMData)



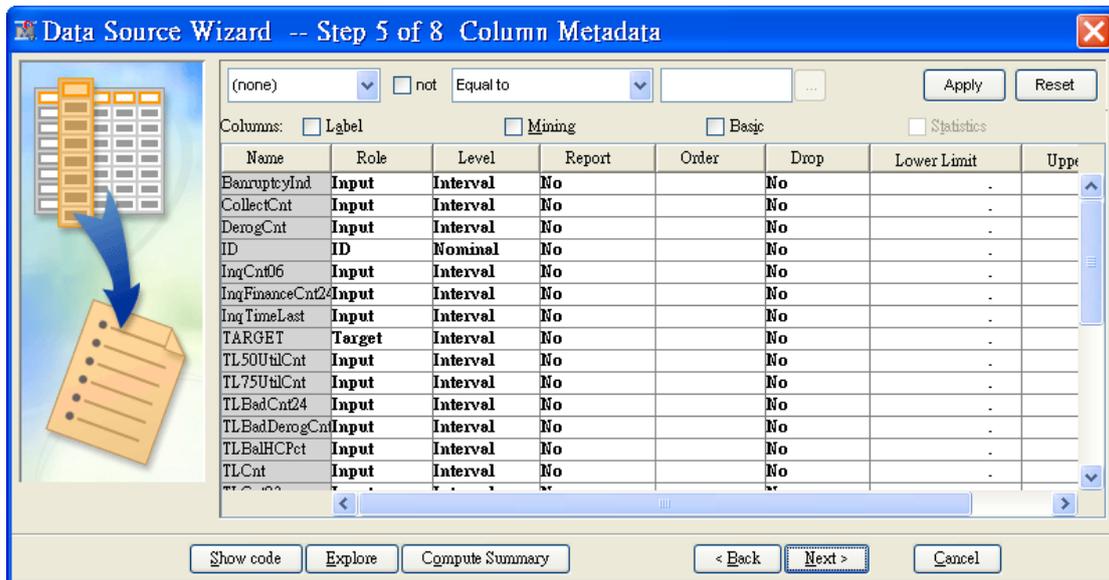
3. 指定 db source





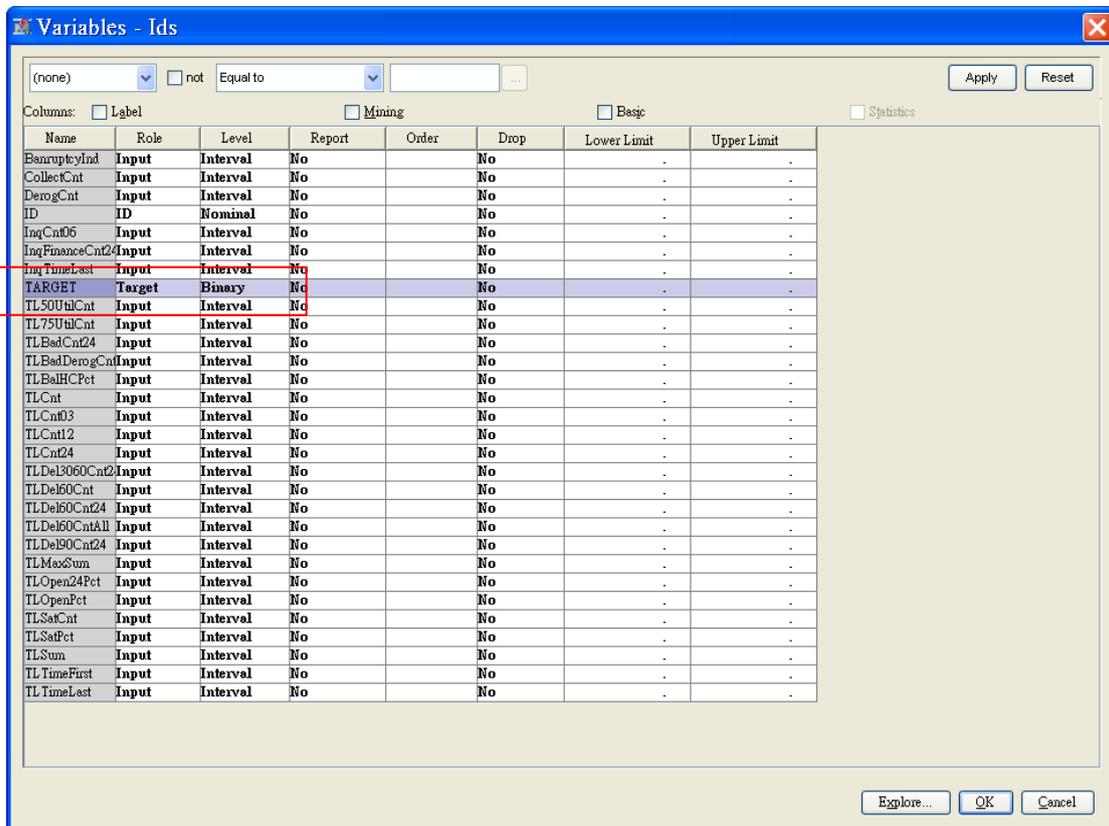
欄位名稱及角色

1. 一定要有 Target 當 Y 變數
2. input 為 X 變數

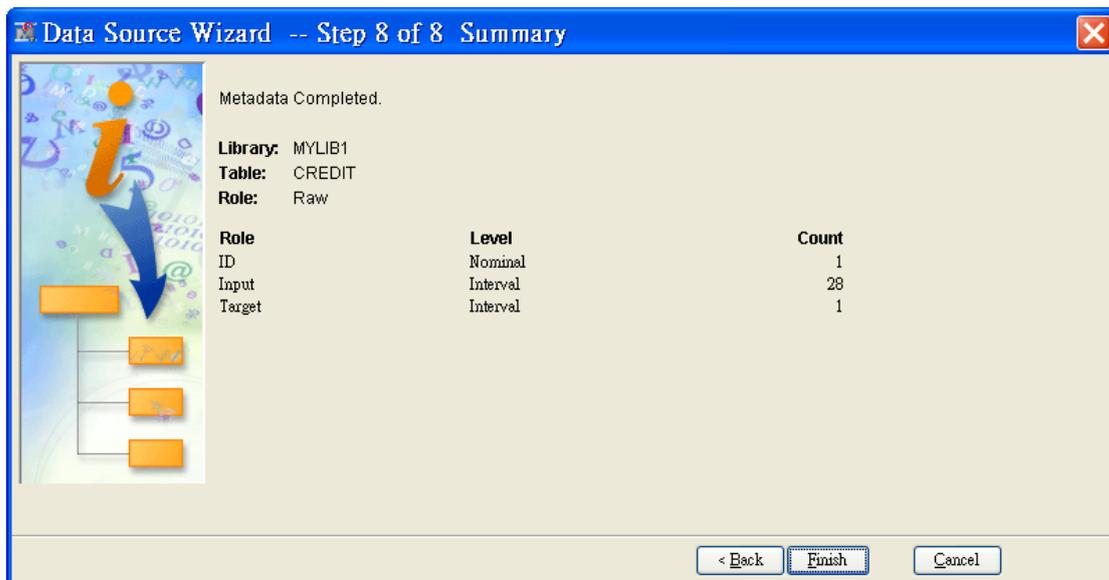
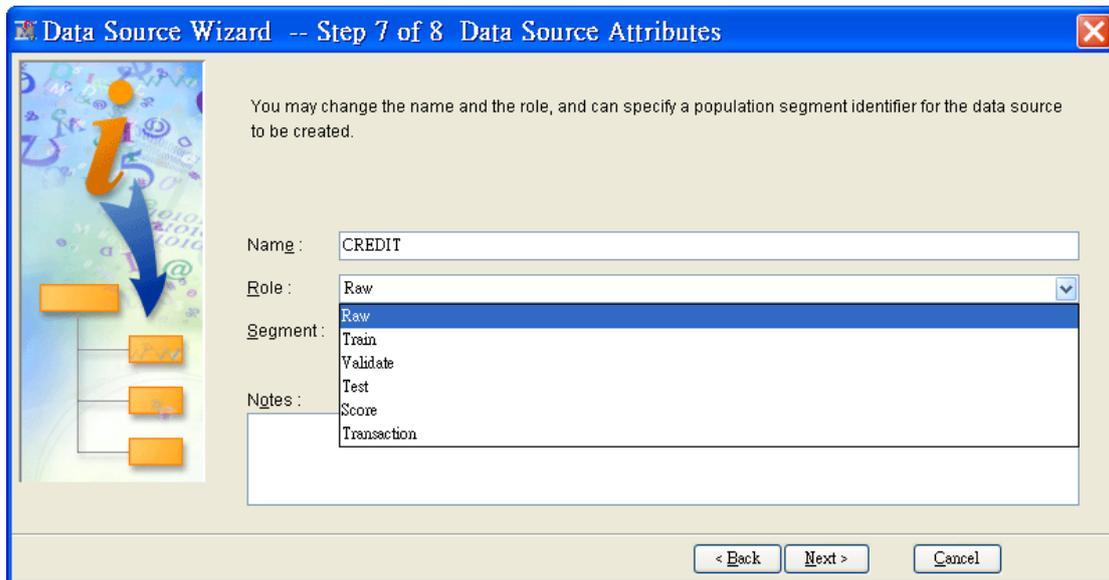


Interval：在 decision tree 會產生平均值，如：購買幾次

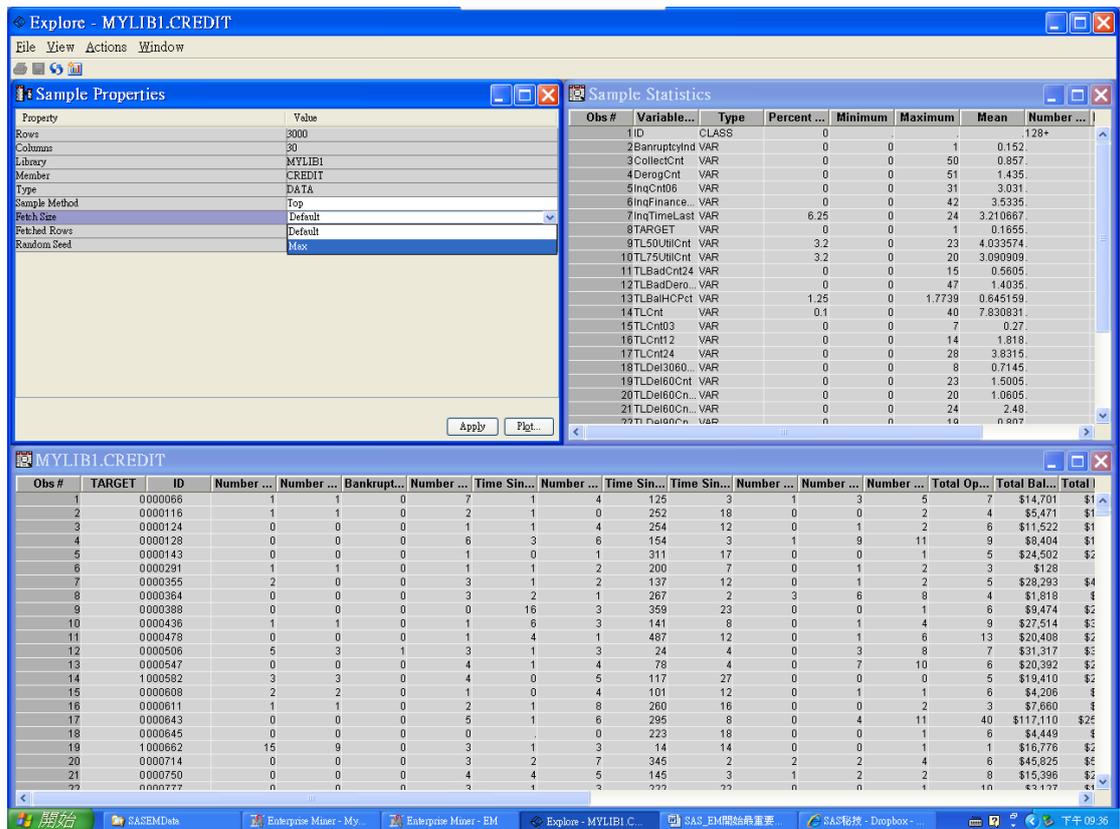
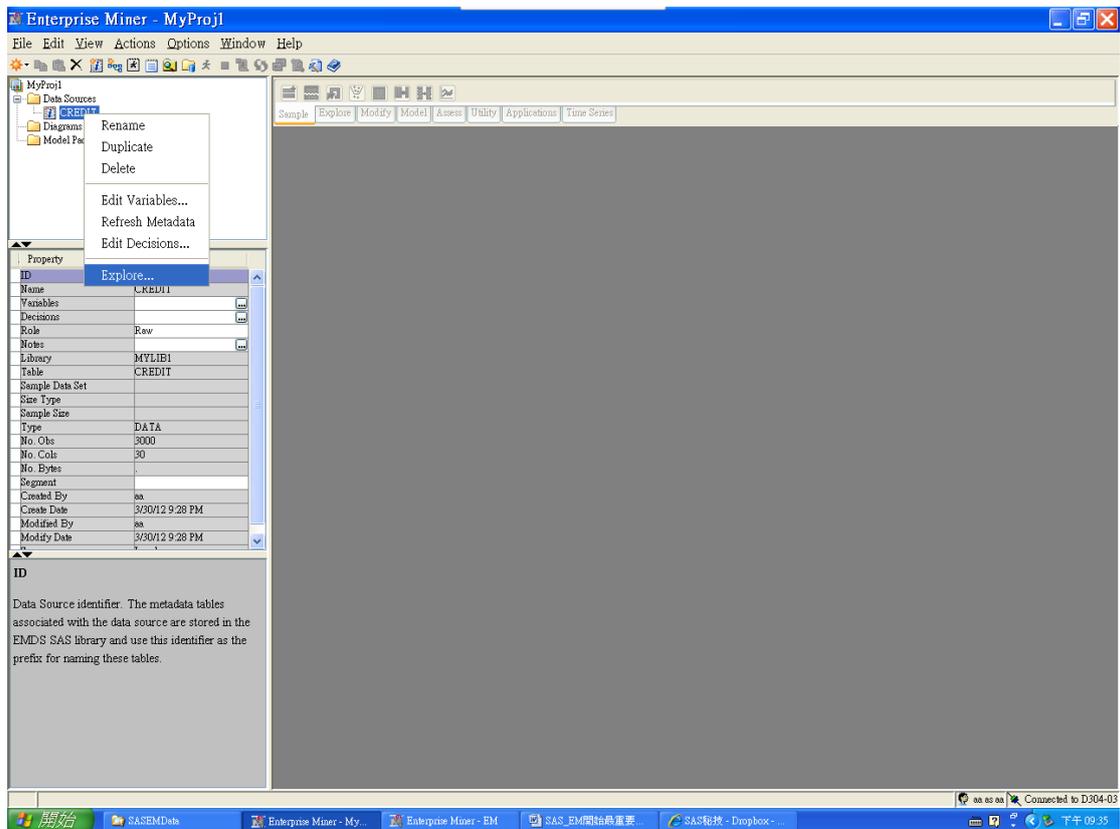
若改為 Binary 為 0/1 or yes/no（如：1=好 0=壞客戶、yes=會 no=不會買）

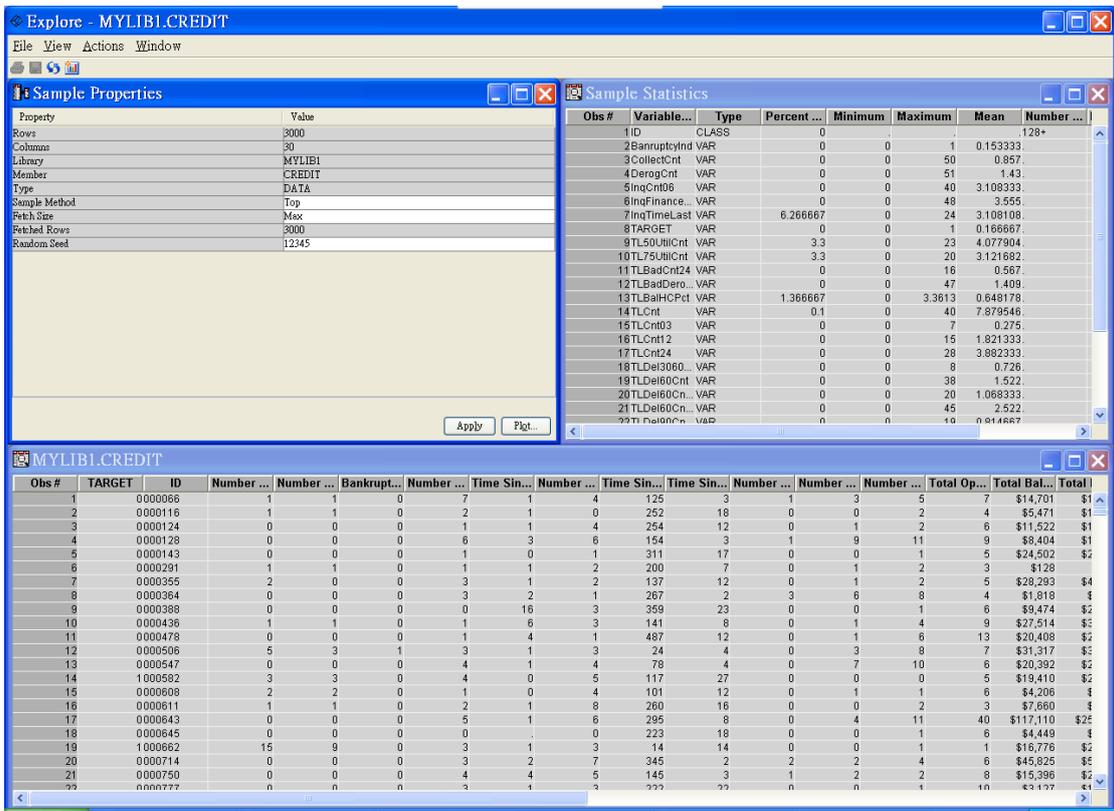


選 Raw data

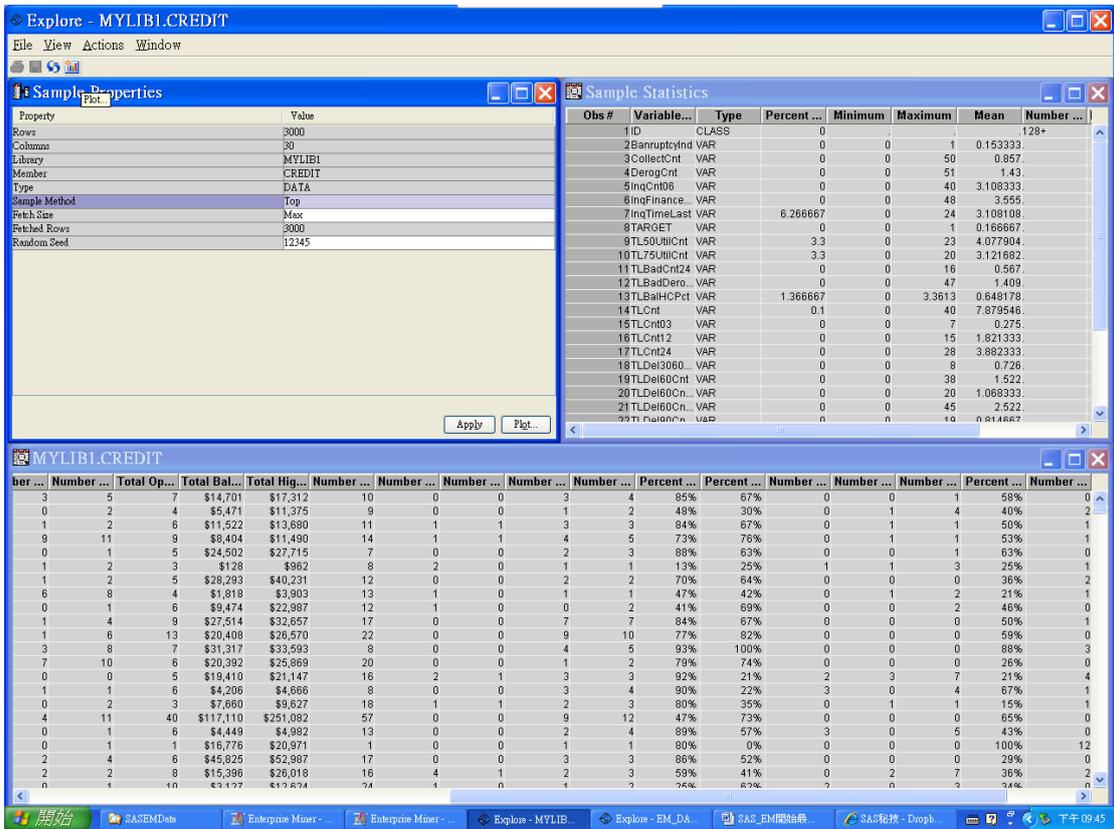


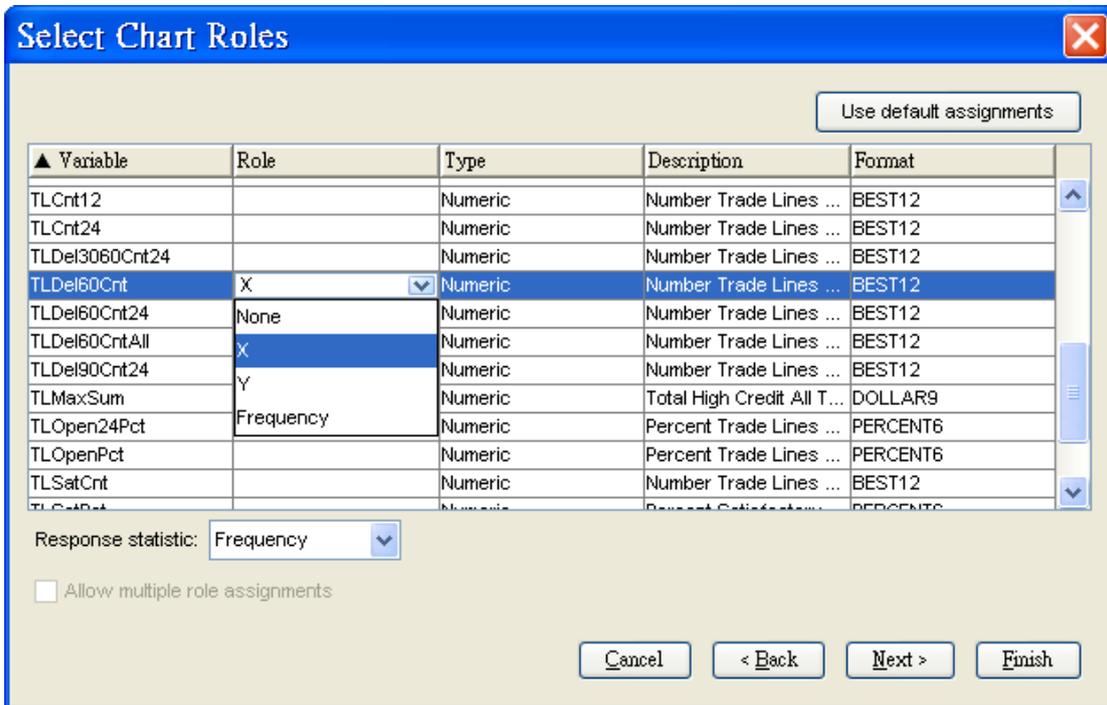
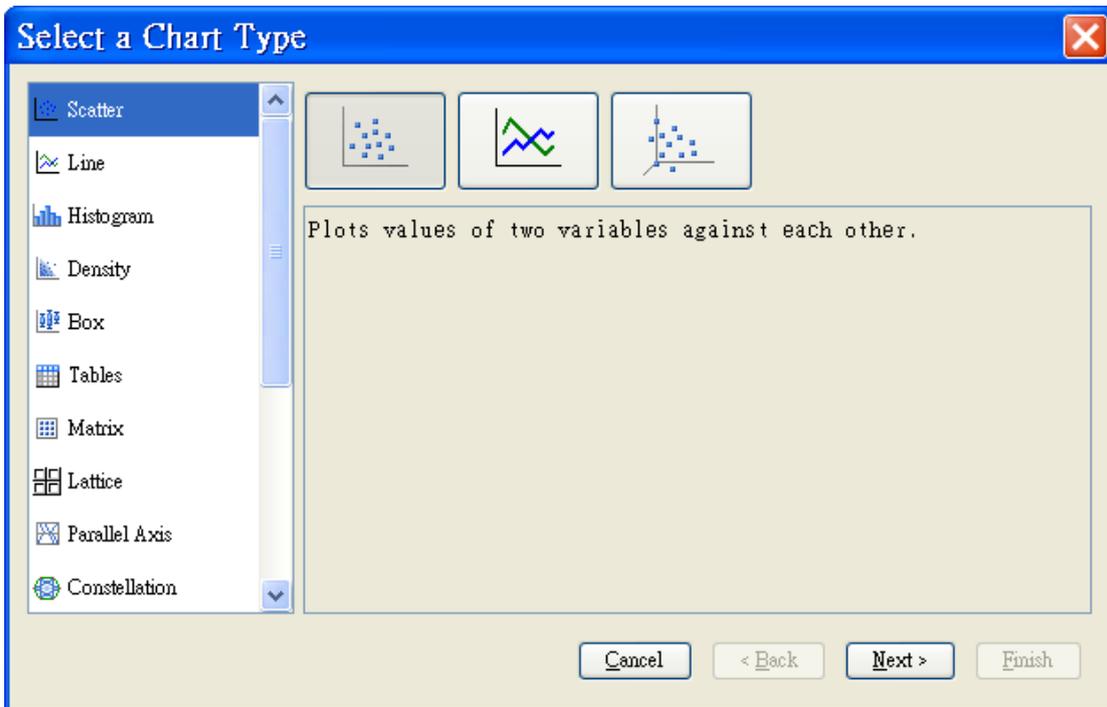
觀察資料：改最大值(Max)，講義後面章節 P.72 有資料庫欄位說明（預設為 2000 筆，改完為 3000 筆）

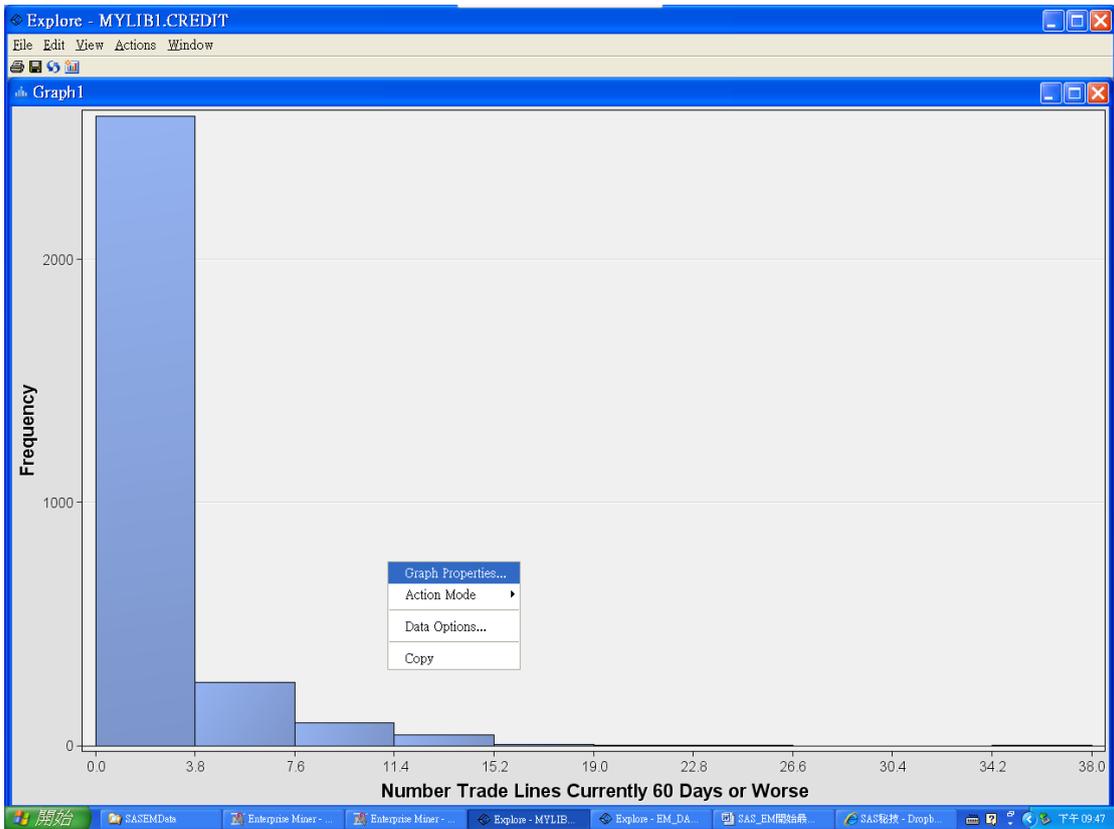
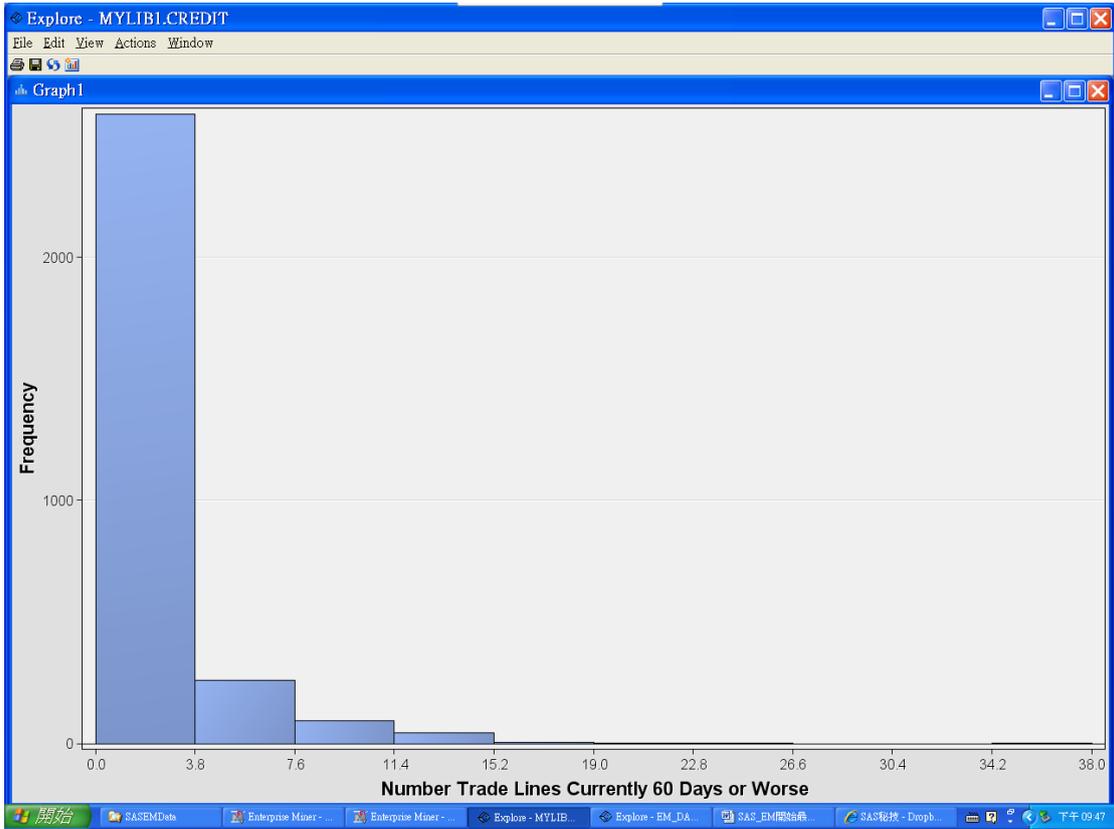




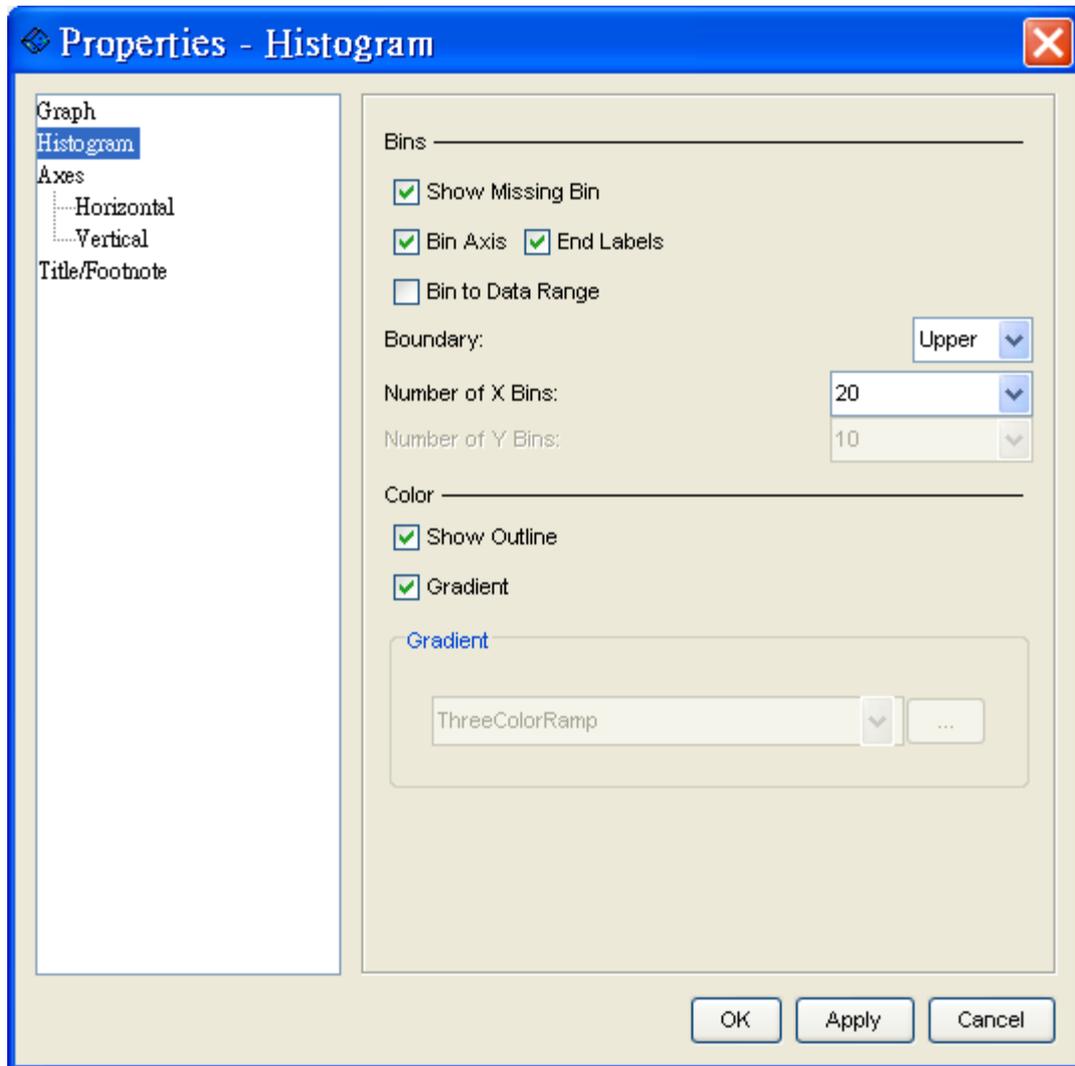
Plot

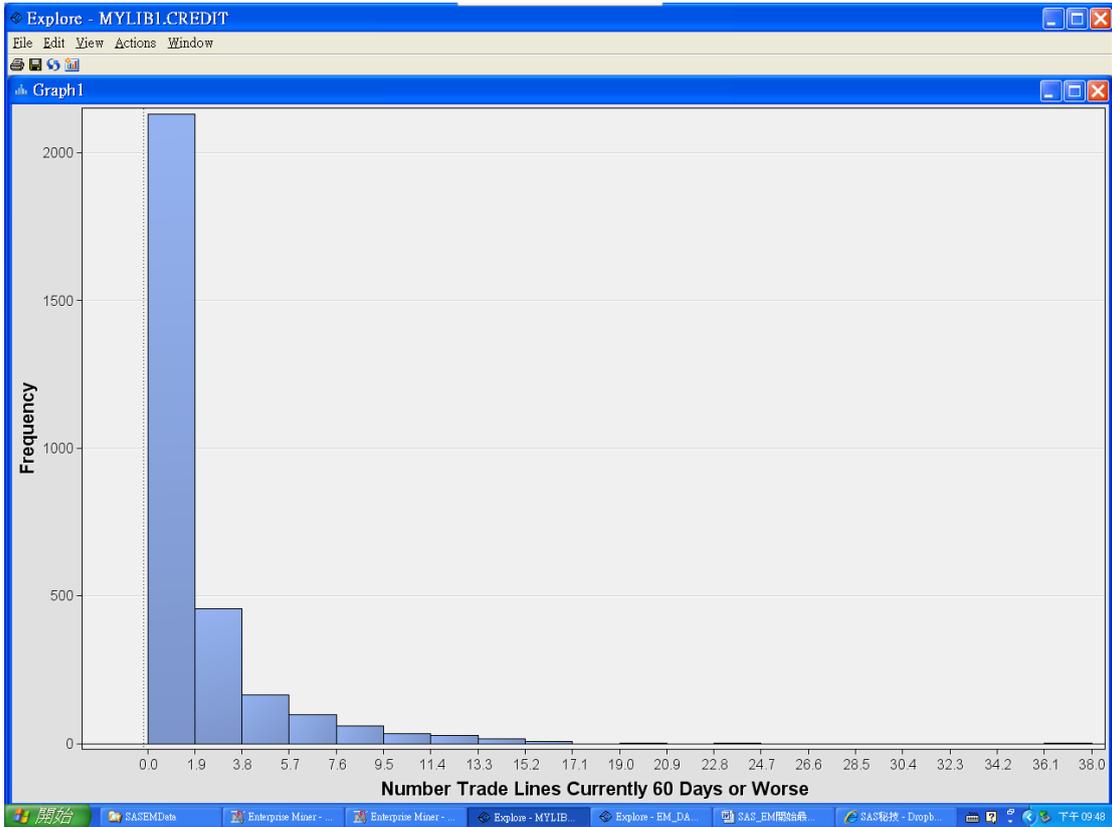




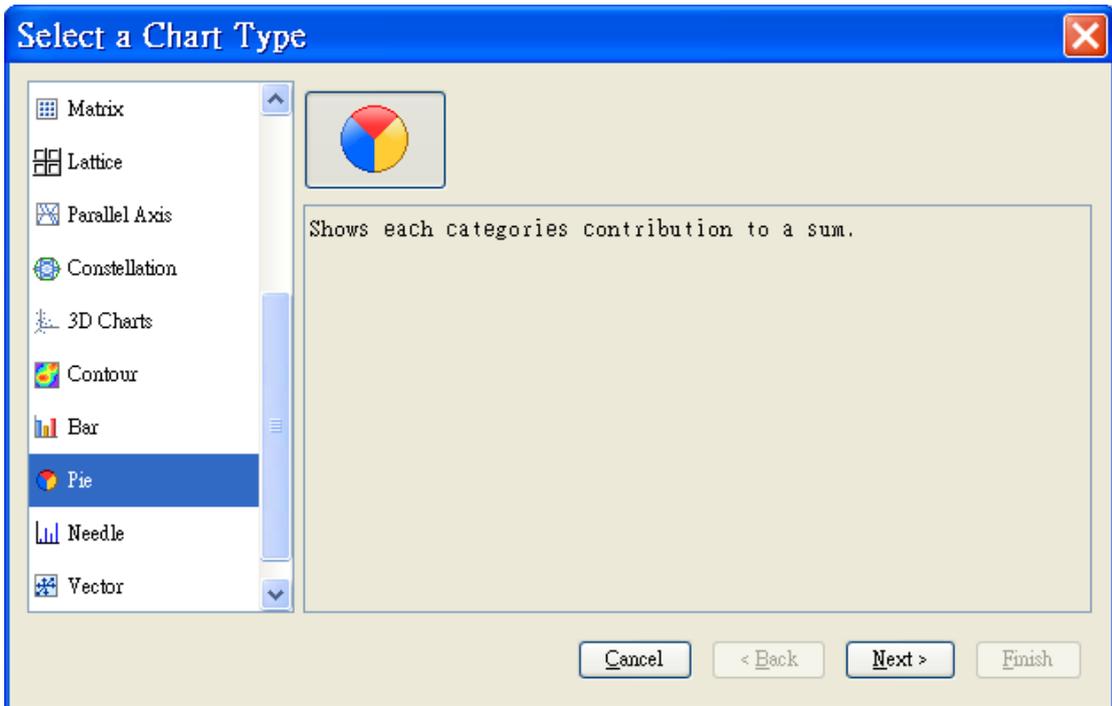


Bin：區塊





Pie



區分好壞客戶(0.1)

Select Chart Roles

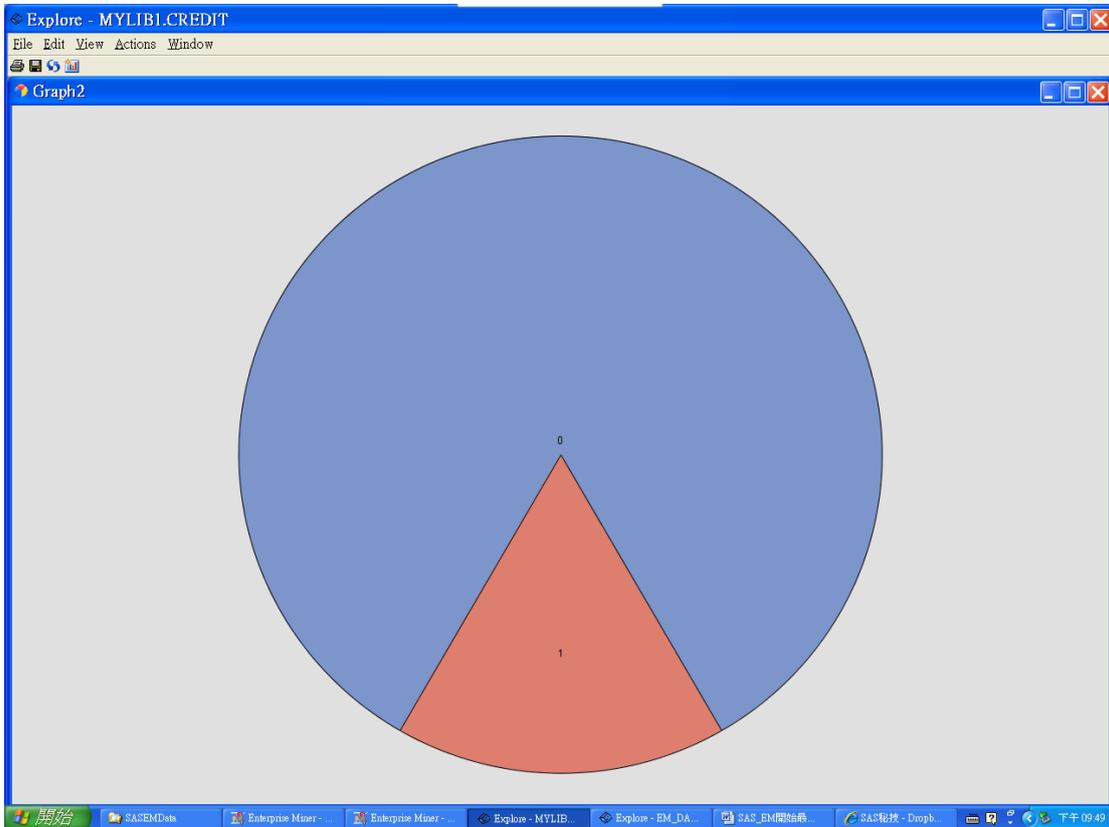
Use default assignments

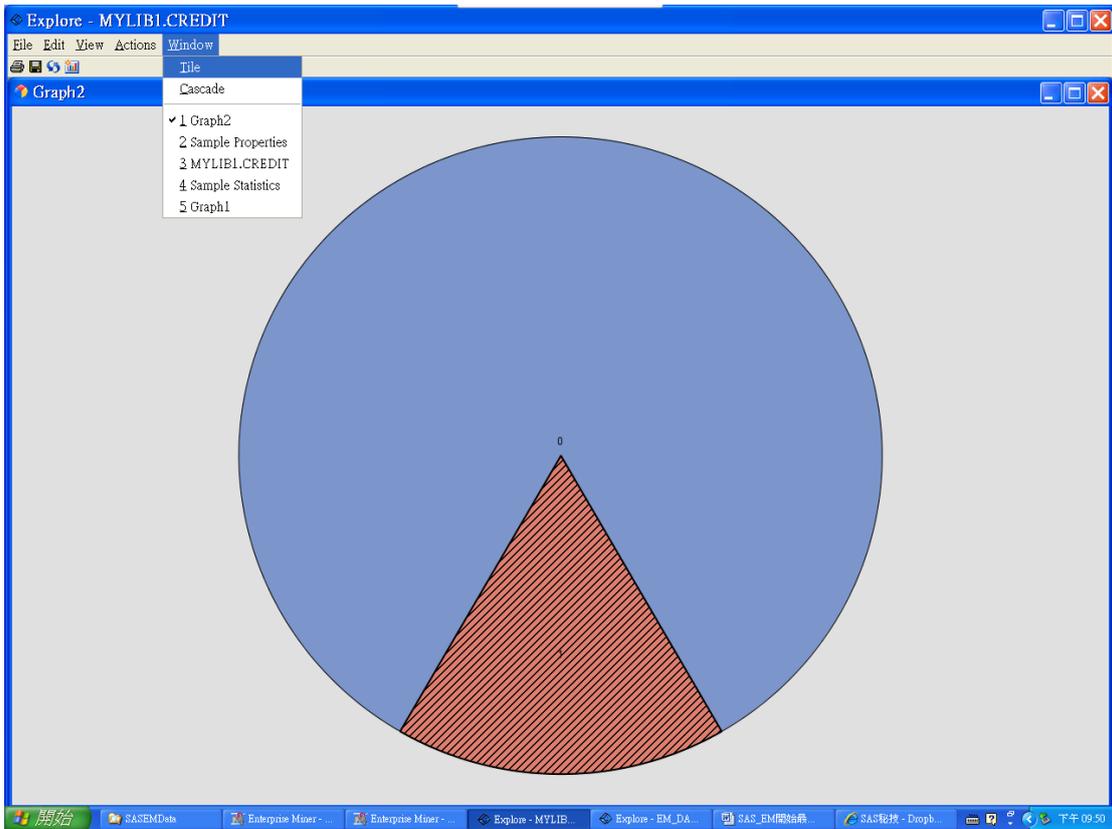
Variable	Role	Type	Description	Format
BanruptcyInd		Numeric	Bankruptcy Indicator	BEST12
CollectCnt		Numeric	Number Collections	BEST12
DerogCnt		Numeric	Number Public Derog...	BEST12
ID		Character	ID	
InqCnt06		Numeric	Number Inquiries 6 M...	BEST12
InqFinanceCnt24		Numeric	Number Finance Inqui...	BEST12
InqTimeLast		Numeric	Time Since Last Inquiry	BEST12
TARGET	Category	Numeric	TARGET	
TL50UtilCnt		Numeric	Number Trade Lines ...	BEST12
TL75UtilCnt		Numeric	Number Trade Lines ...	BEST12
TLBadCnt24		Numeric	Number Trade Lines ...	BEST12
TLBadDerogCnt		Numeric	Number Bad Dept plu...	BEST12
TLBalHCPct		Numeric	Percent Trade Line B...	PERCENT6
TLCnt		Numeric	Total Open Trade Lines	BEST12

Allow multiple role assignments

Cancel < Back Next > Finish

1:壞客戶:500
0:好客戶:2500



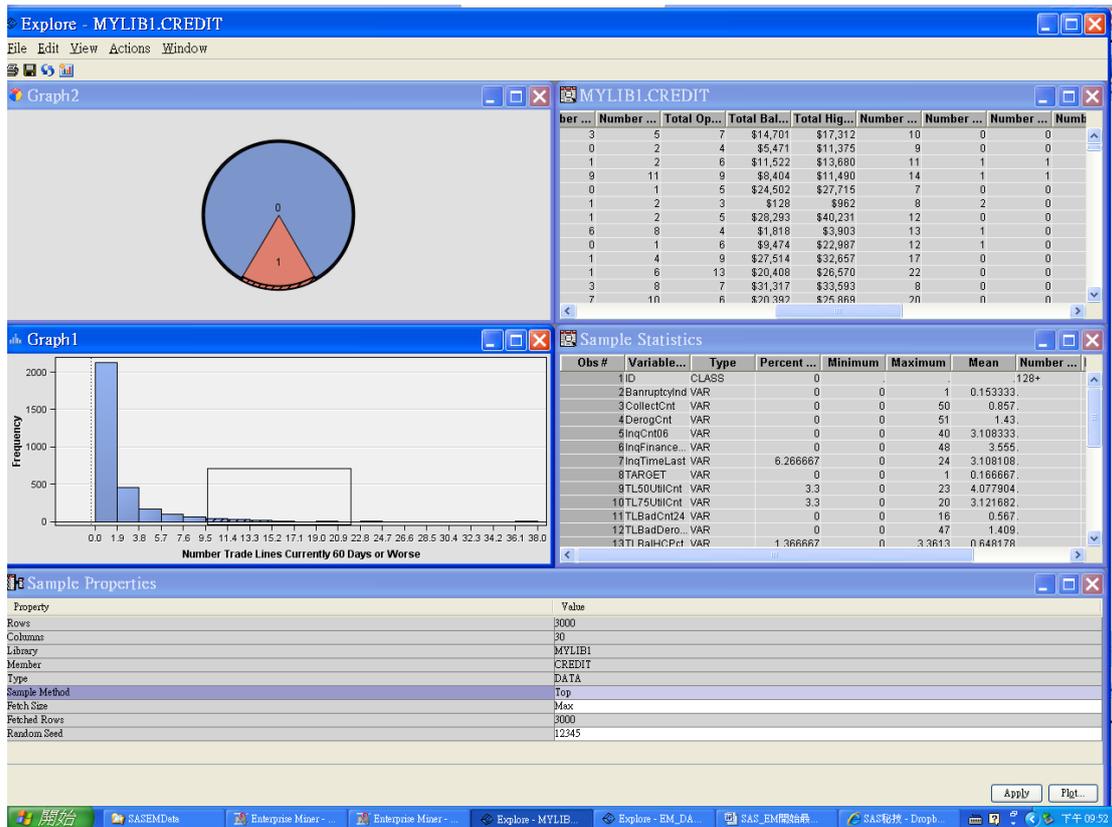


ber ...	Number ...	Total Op...	Total Bal...	Total Hig...	Number ...	Number ...	Number ...	Numb...
3	5	7	\$14,701	\$17,312	10	0	0	0
0	2	4	\$5,471	\$11,375	9	0	0	0
1	2	6	\$11,522	\$13,680	11	1	1	1
9	11	9	\$9,404	\$11,490	14	1	1	1
0	1	5	\$24,502	\$27,715	7	0	0	0
1	2	3	\$128	\$962	8	2	0	0
1	2	5	\$28,293	\$40,231	12	0	0	0
6	8	4	\$1,818	\$3,903	13	1	0	0
0	1	6	\$9,474	\$22,987	12	1	0	0
1	4	9	\$27,514	\$32,657	17	0	0	0
1	6	13	\$20,408	\$26,570	22	0	0	0
3	8	7	\$31,317	\$33,593	8	0	0	0
7	10	6	\$20,382	\$25,668	20	0	0	0

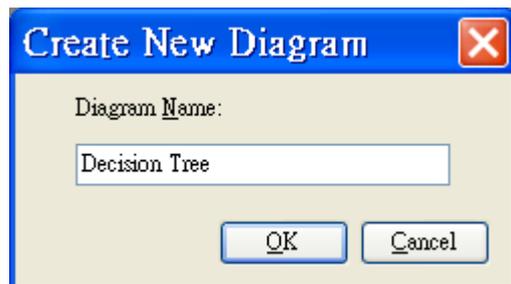
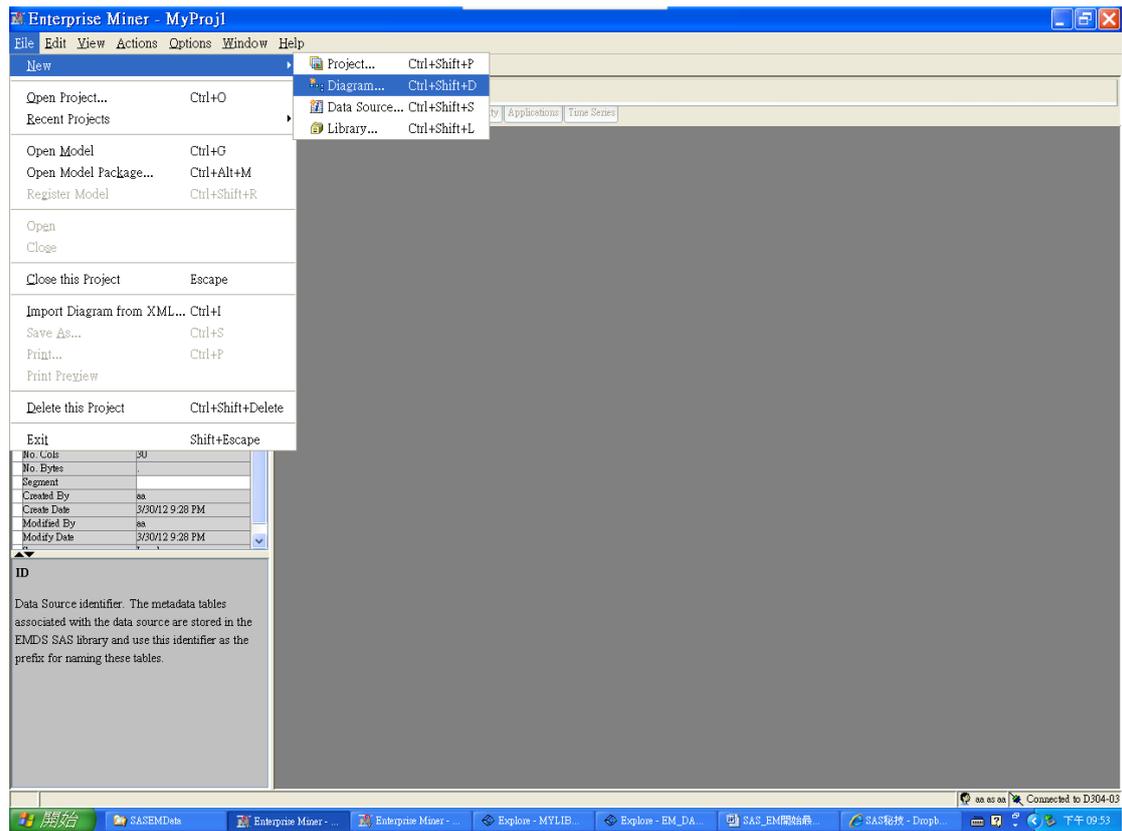
Obs #	Variable...	Type	Percent ...	Minimum	Maximum	Mean	Number ...
1	ID	CLASS	0	.	.	.	128+
2	BanruptyInd	VAR	0	0	1	0.153333	
3	CollectCnt	VAR	0	0	50	0.857	
4	DerogCnt	VAR	0	0	51	1.43	
5	InqCnt06	VAR	0	0	40	3.108333	
6	InqFinance...	VAR	0	0	48	3.555	
7	InqTimeLast	VAR	6.266667	0	24	3.108108	
8	TARQET	VAR	0	0	1	0.166667	
9	TL50UBICnt	VAR	3.3	0	23	4.077904	
10	TL75UBICnt	VAR	3.3	0	20	3.121882	
11	TLBadCnt24	VAR	0	0	16	0.567	
12	TLBadDero...	VAR	0	0	47	1.409	
13	TLRainCOPH	VAR	1.366667	0	3.3613	0.648178	

Property	Value
Rows	3000
Columns	30
Library	MYLIB1
Member	CREDIT
Type	DATA
Sample Method	Top
Fetch Size	Max
Fetch Rows	3000
Random Seed	12345

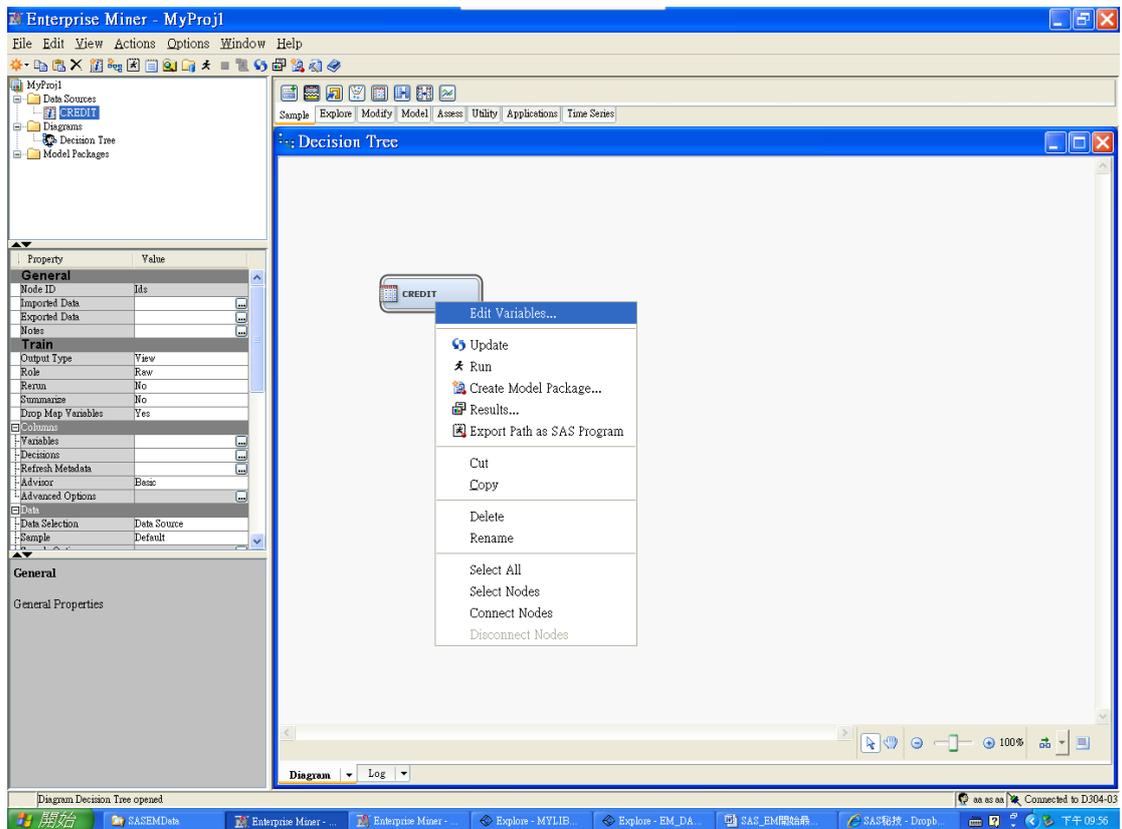
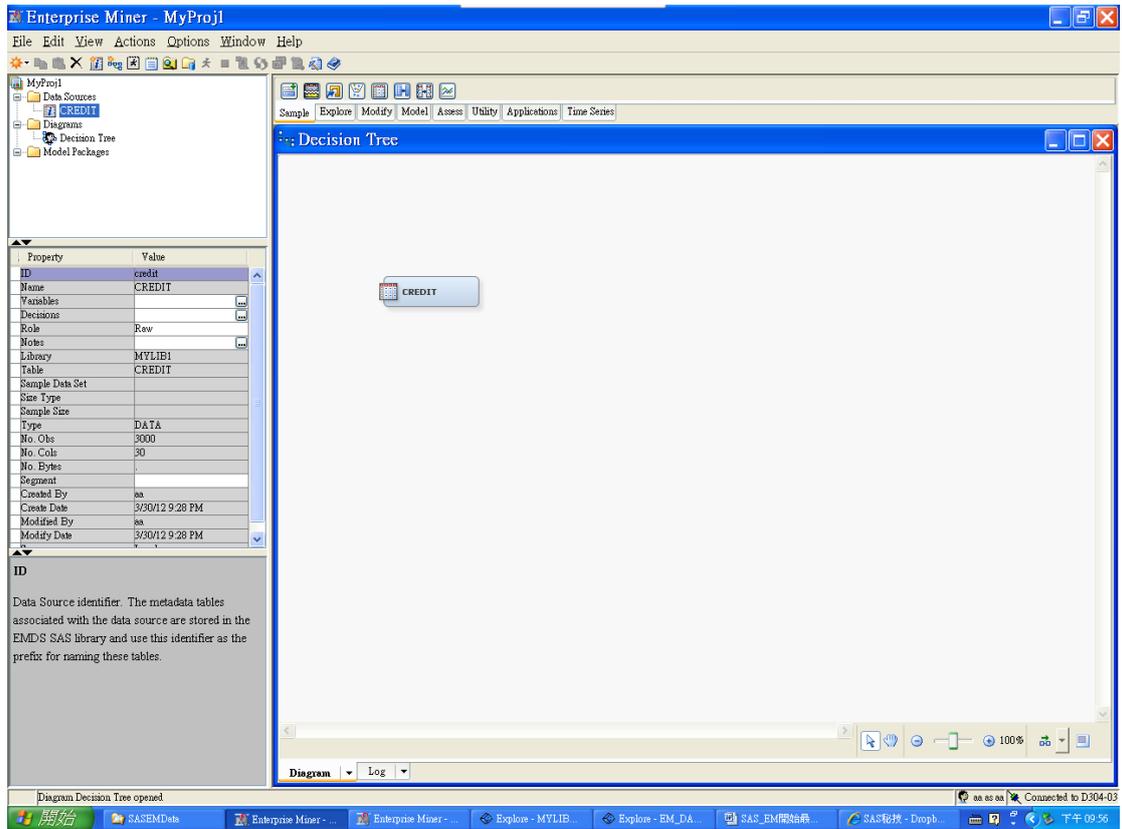
選擇部份區塊，可自動在 Pie 上看出分佈狀況

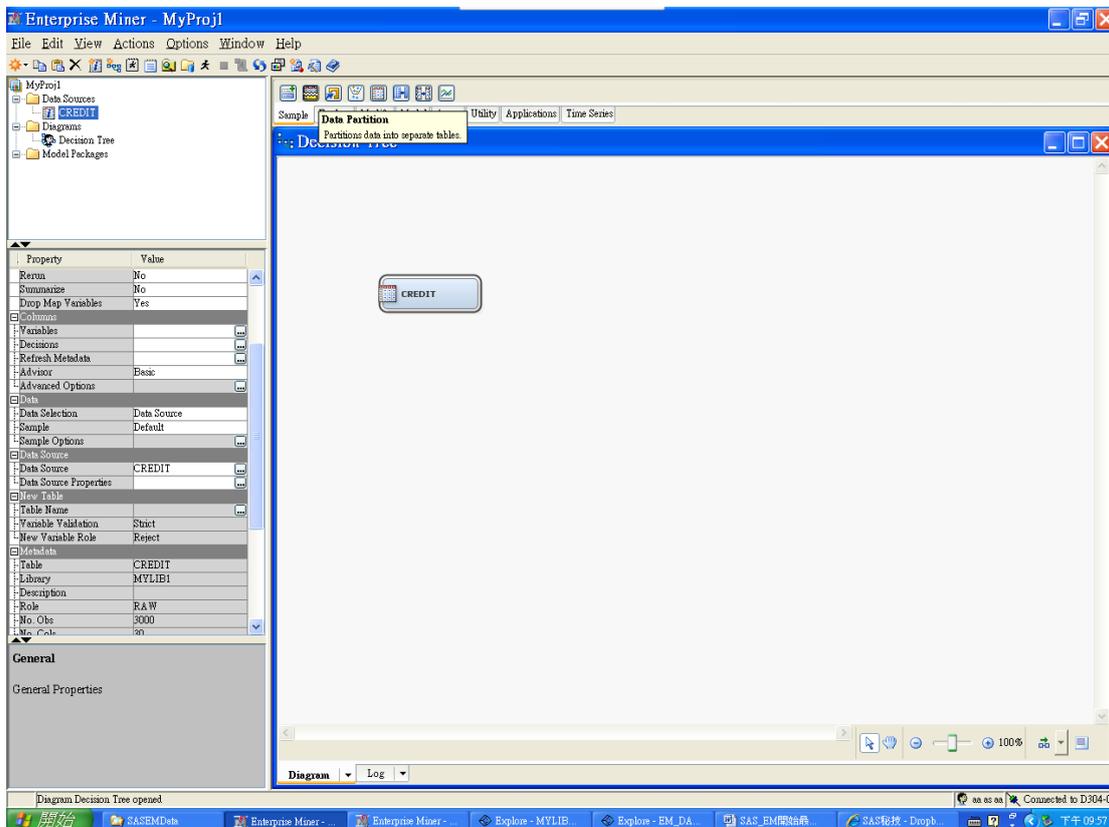
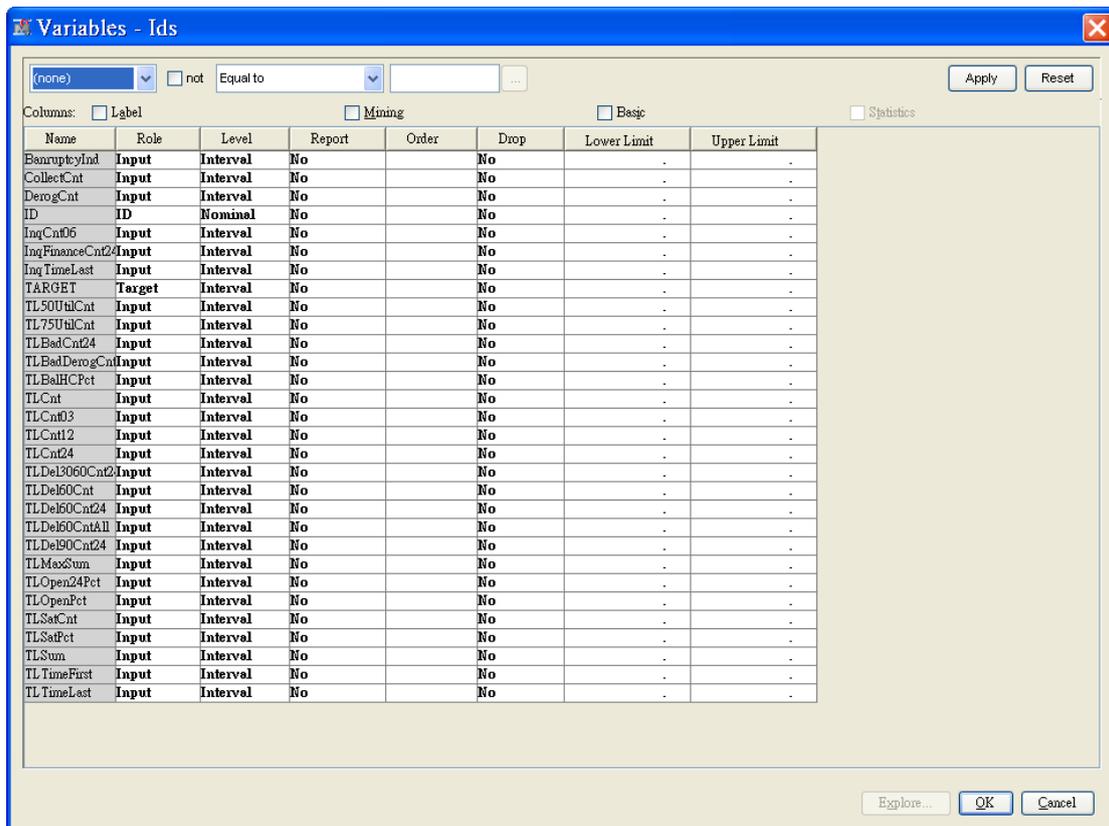


4. Diagram 流程圖 (Decision Tree)

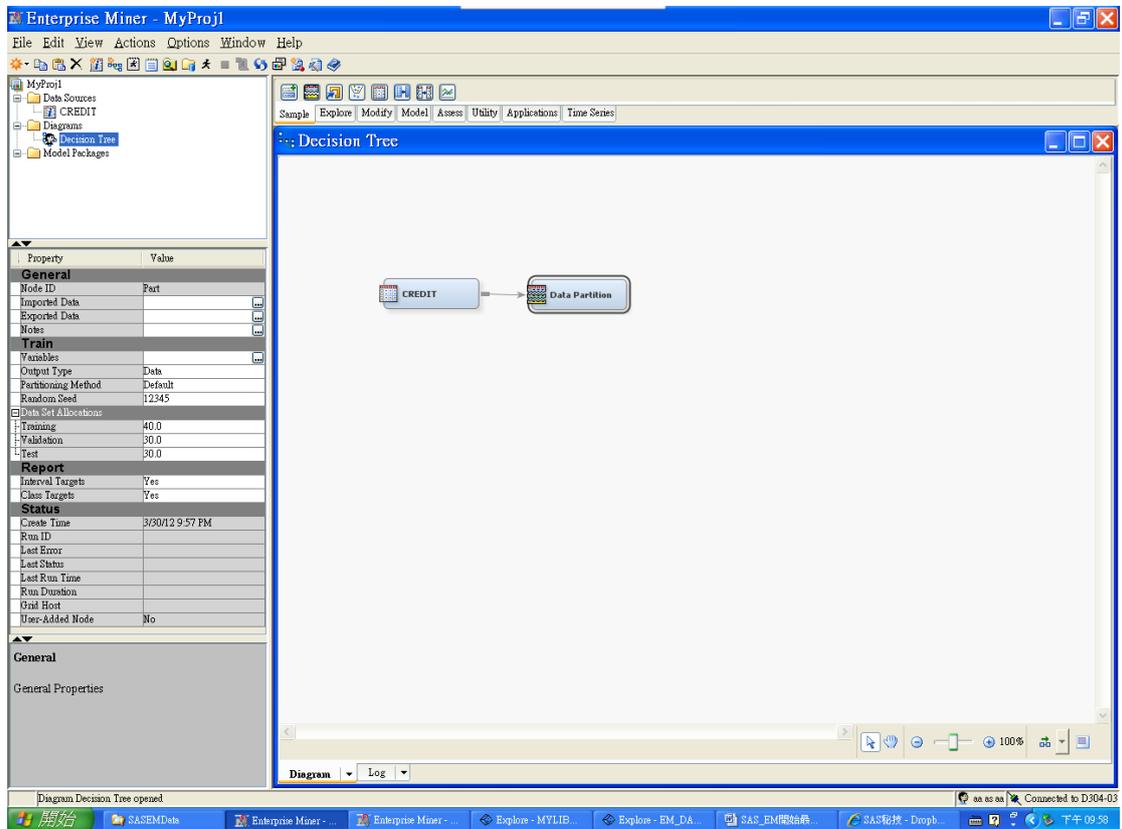
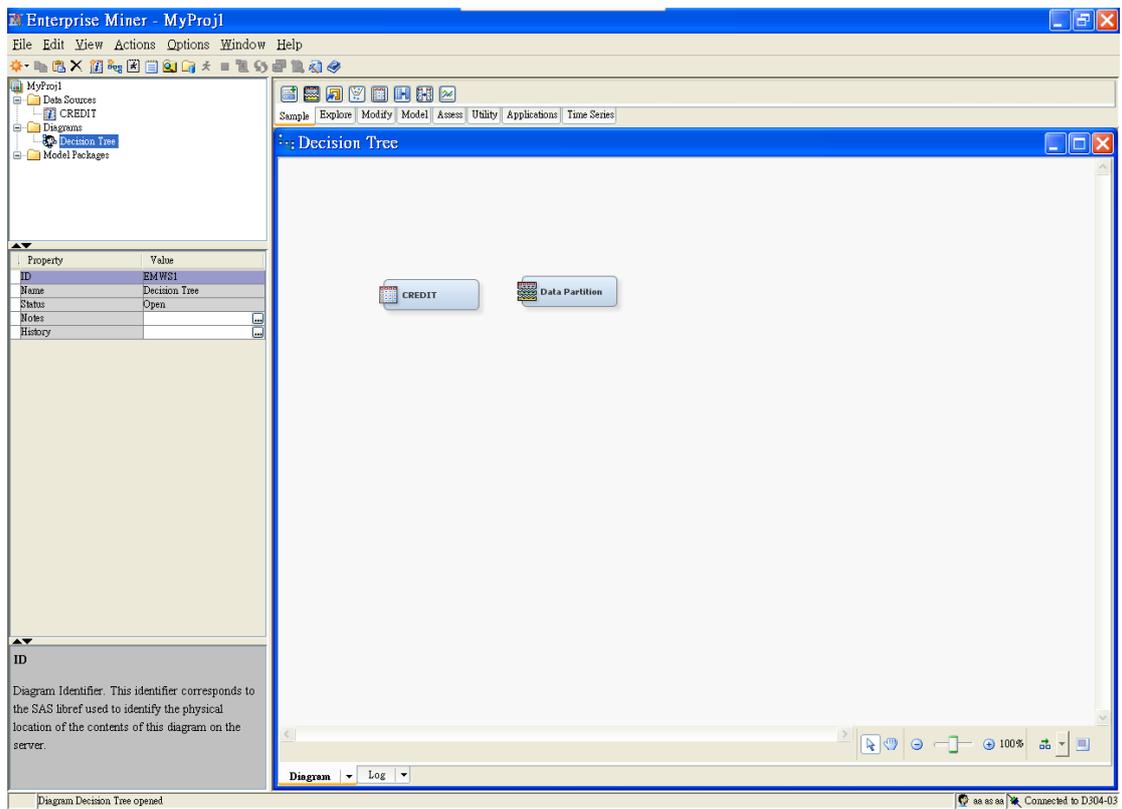


增加節點：(拖拉)

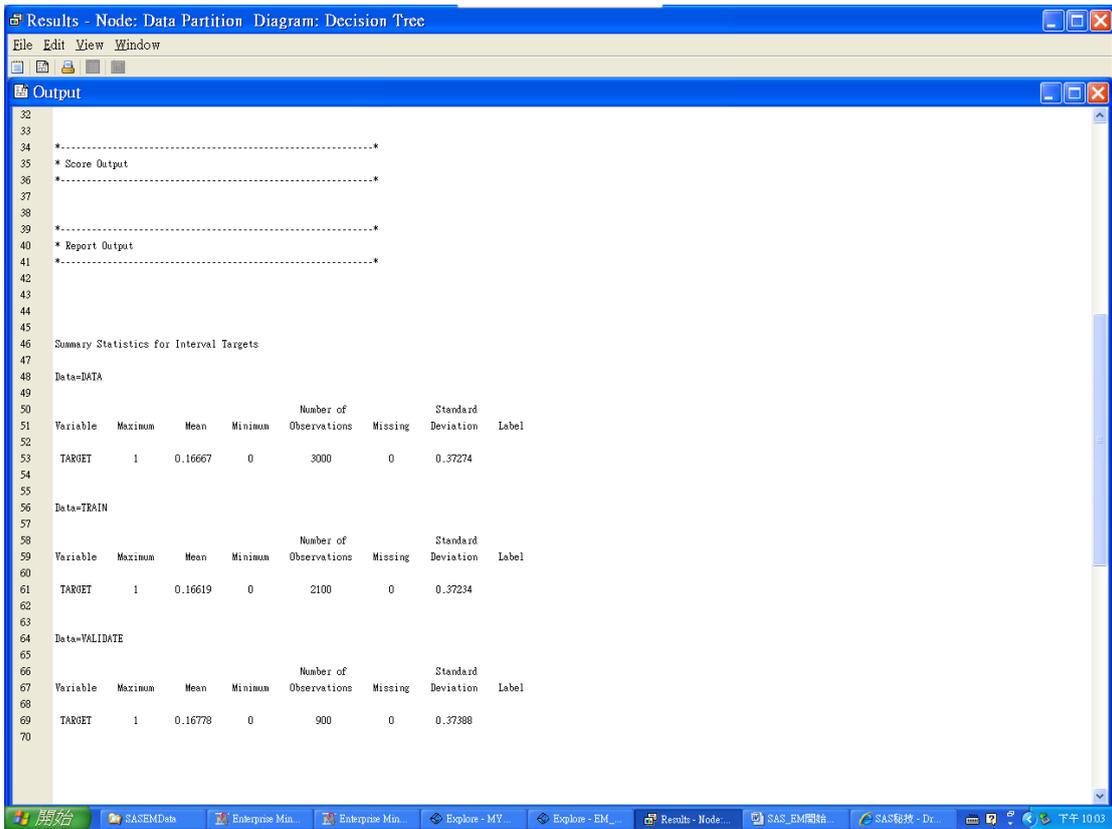




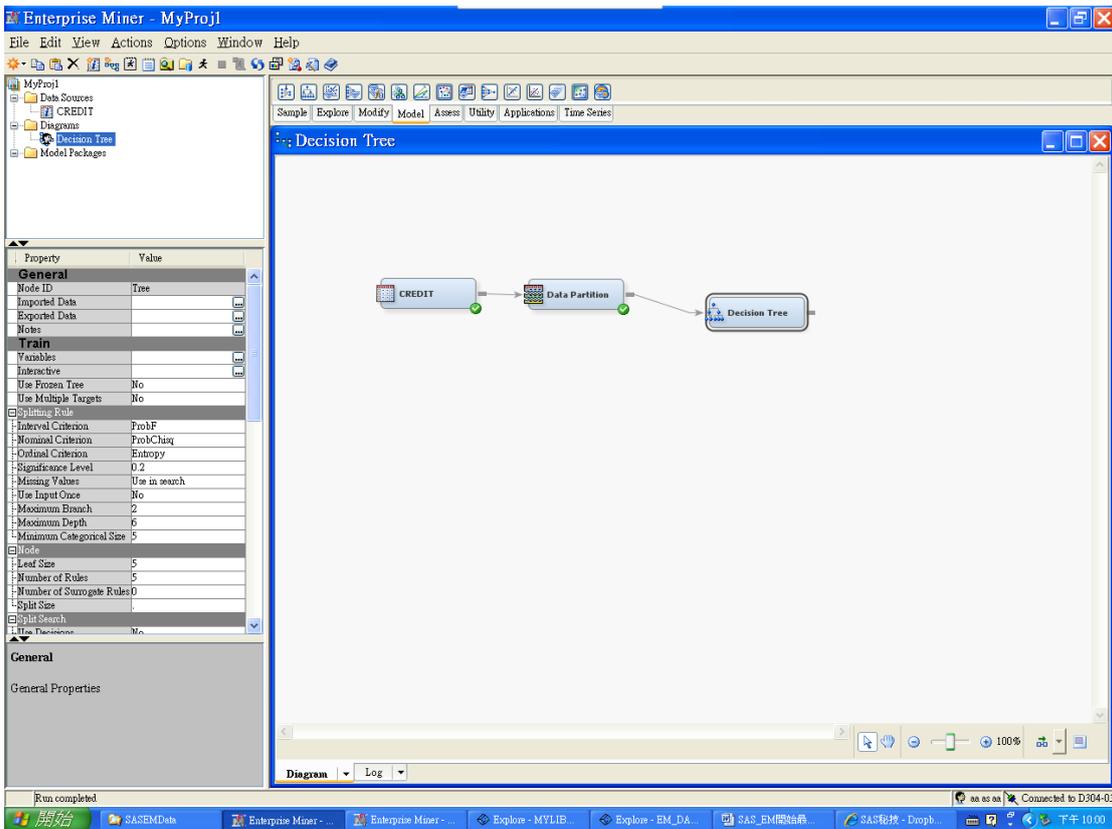
要將 3000 筆分成訓練資料 70%跟驗證資料 30%

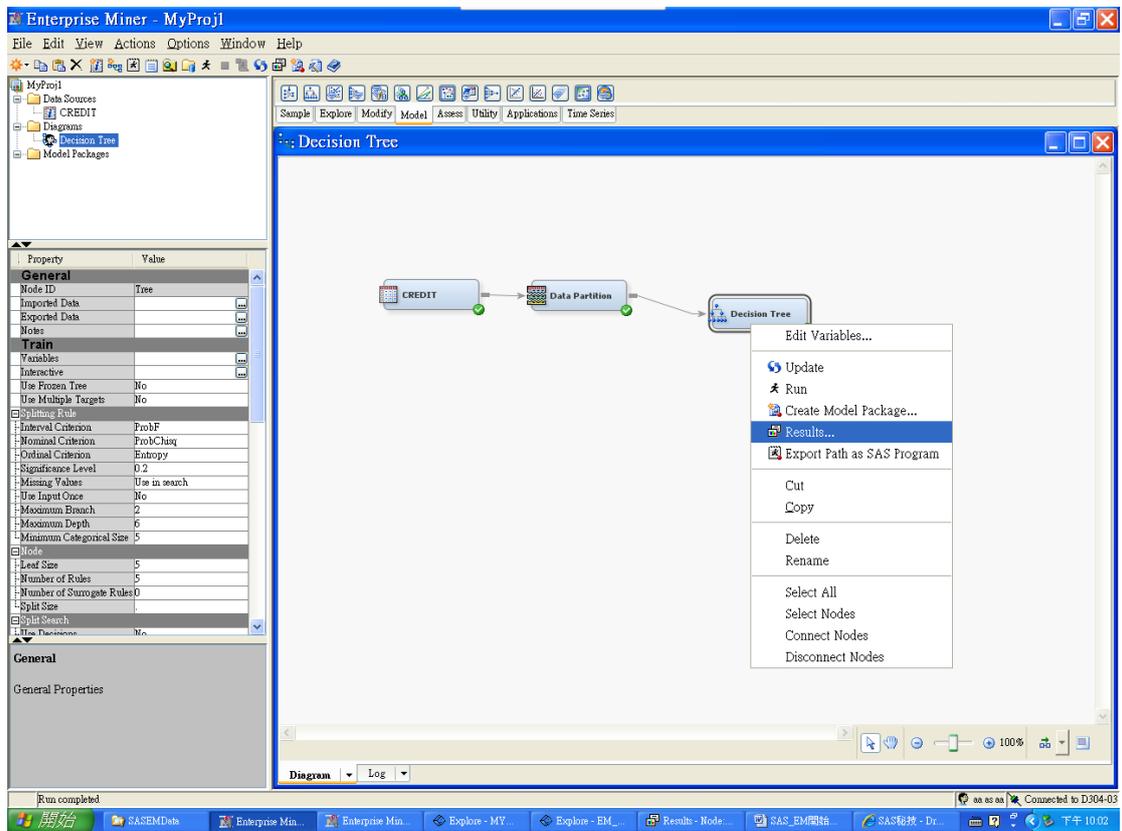
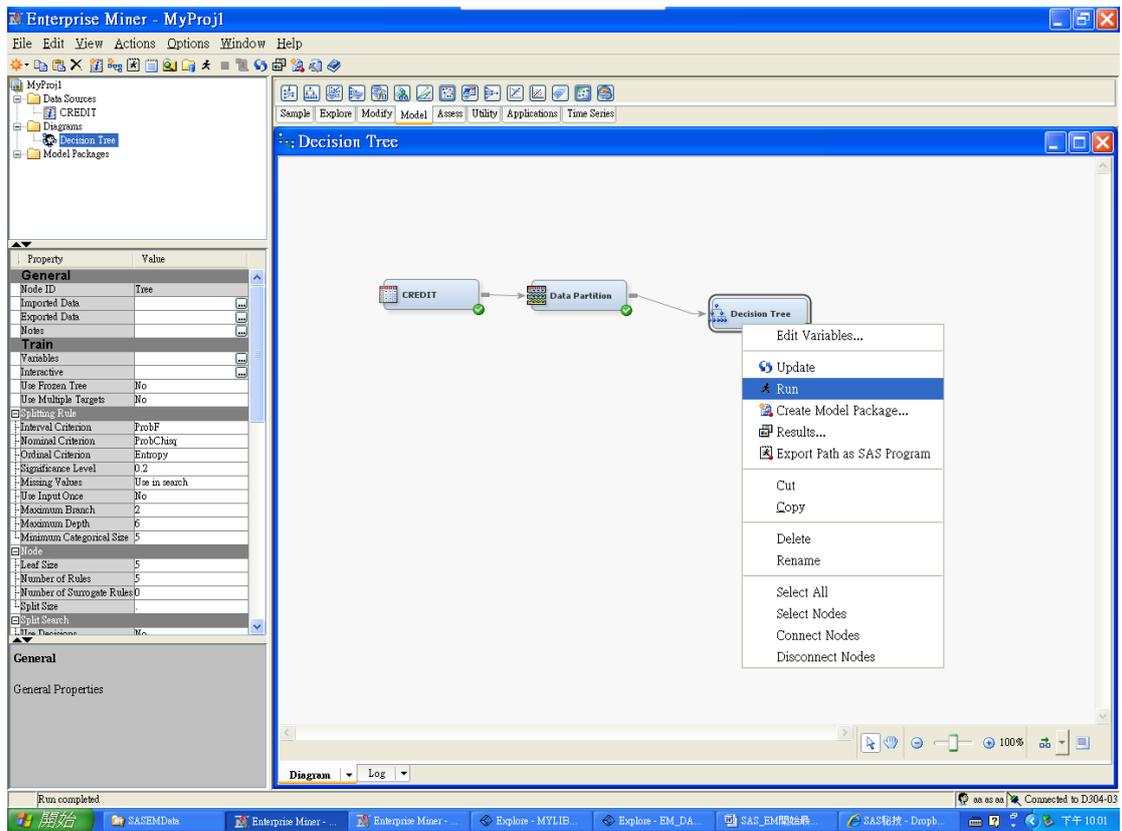


改測試筆數比例：

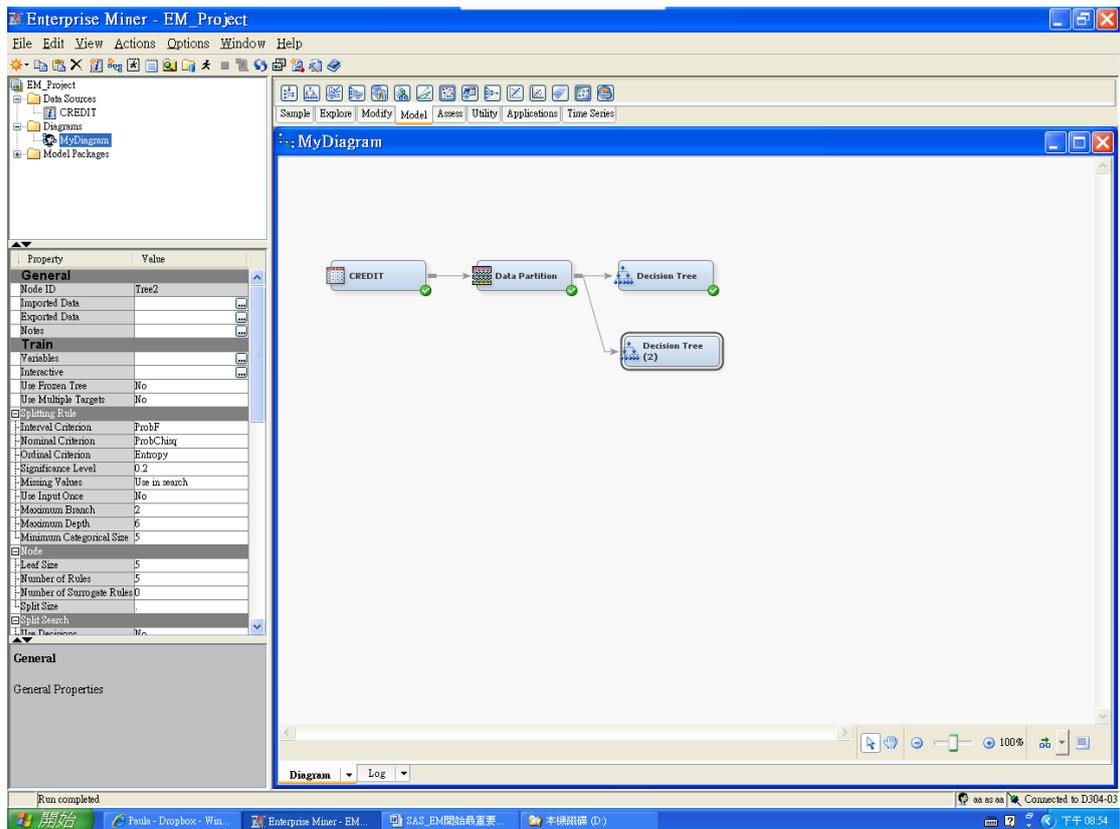


加入決策樹

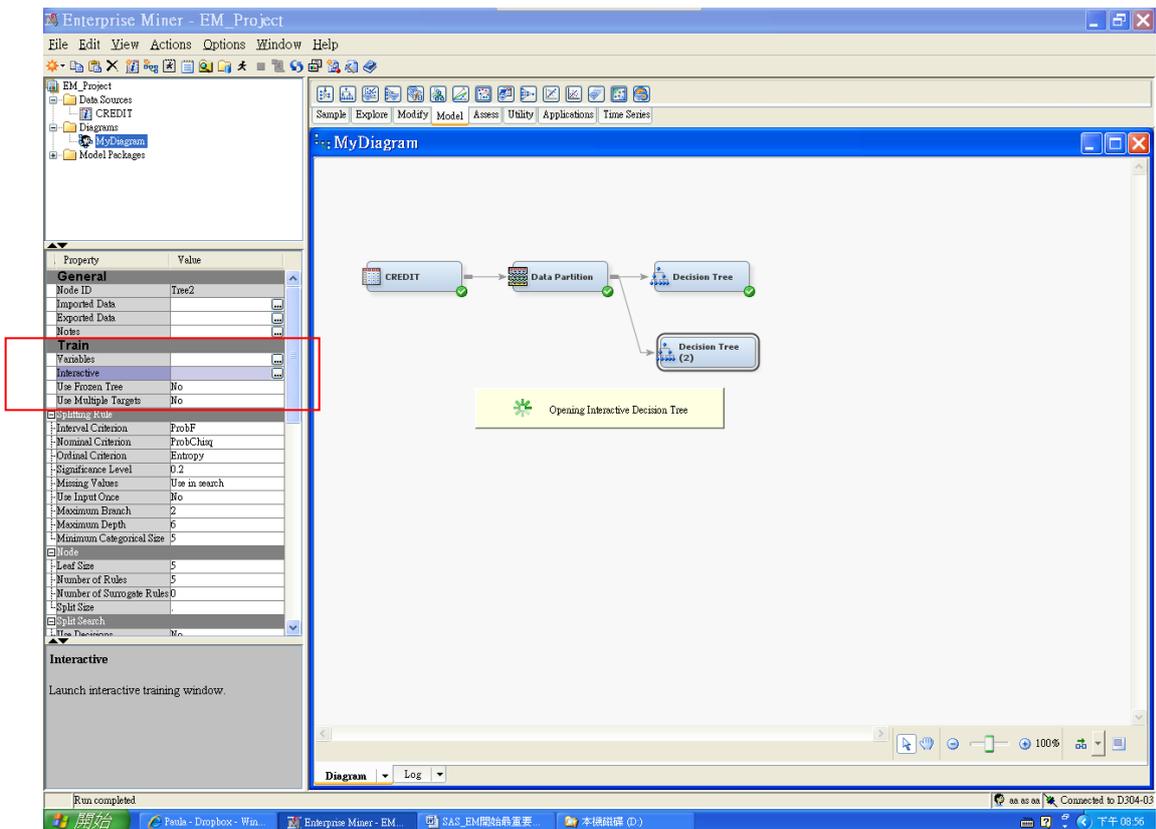




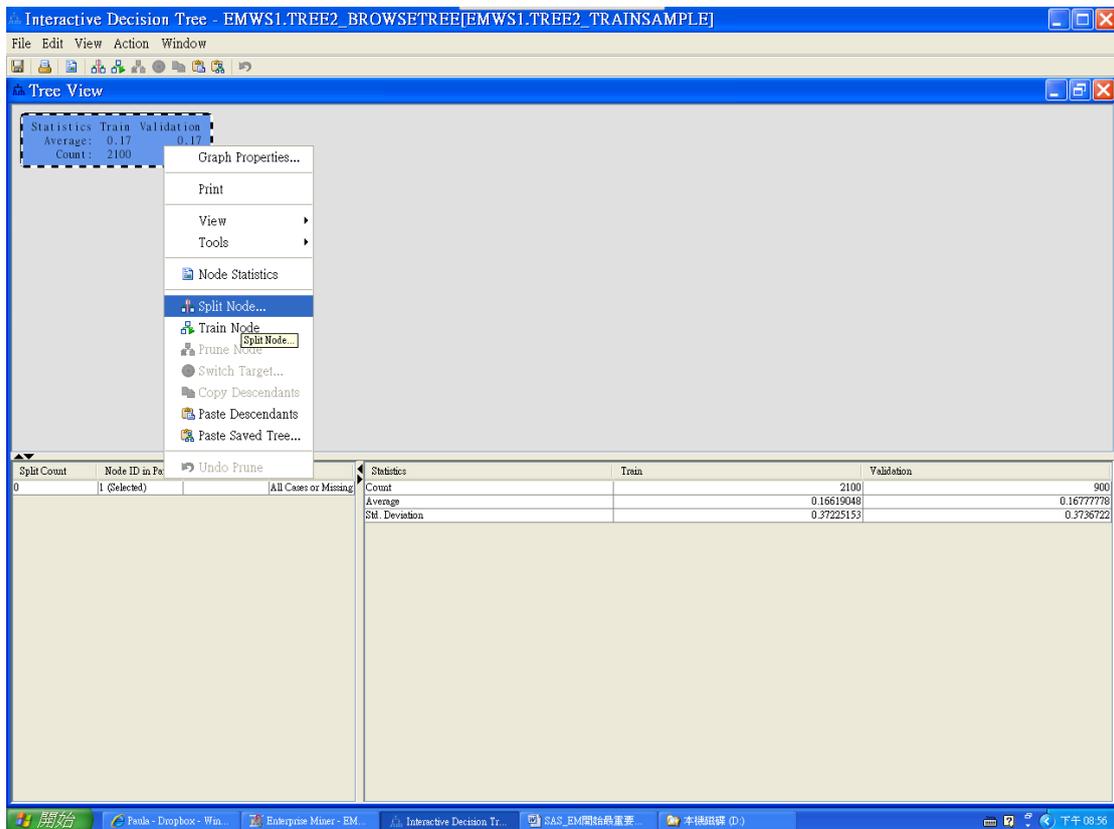
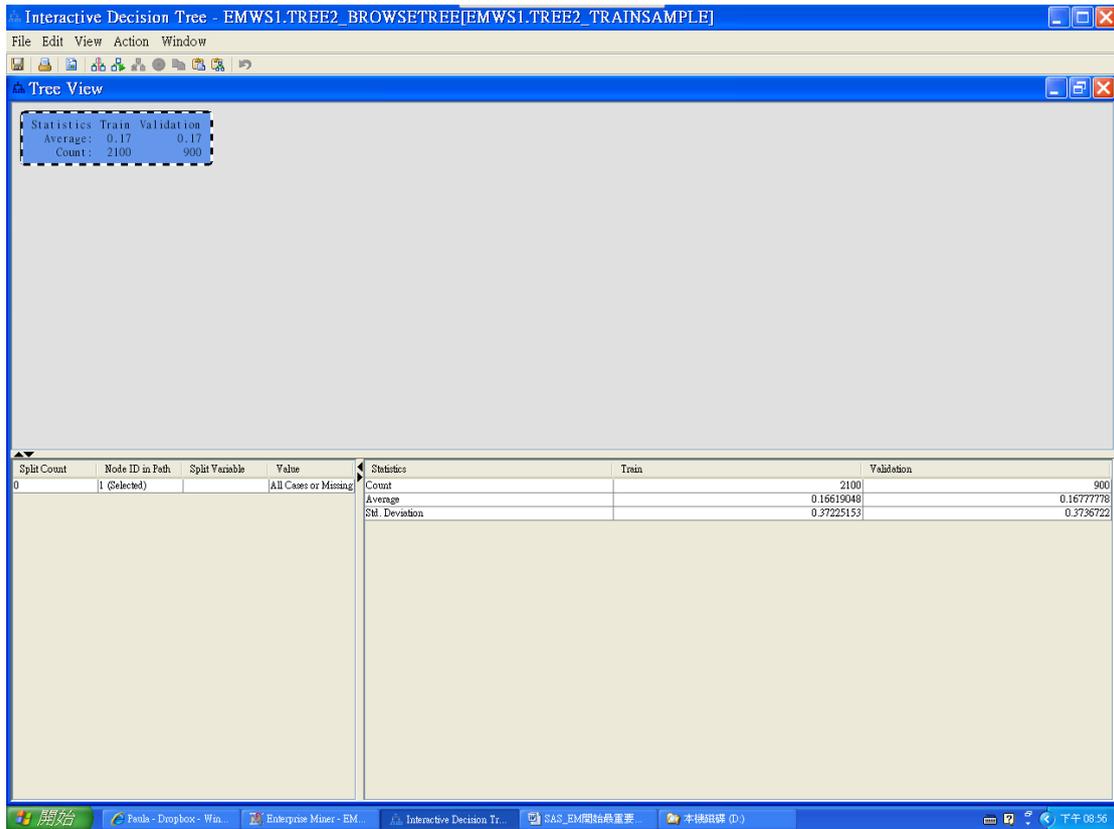
5. 自建決策樹



Train:Interactive



以下為：target 的 level 屬性設為：Interval：在 decision tree 會產生平均值，就會出現以下畫面



-Log(p)愈大，變數愈重要（影響力比較大）

Split Node 1

Target Variable: TARGET

Variable	Variable Description	-Log(p)	Branches
TLDel60Cnt24	Number Trade Lines 60 ...	34.11736551	2
TLDel3060Cnt24	Number Trade Lines 30 ...	22.92869586	2
TLSatPct	Percent Satisfactory to T...	22.84360338	2
TLDel90Cnt24	Number Trade Lines 90+...	22.77306337	2
TLDel60CntAll	Number Trade Lines 60 ...	20.97313146	2
TlBadDerogCnt	Number Bad Dept plus P...	20.32827873	2
TLDel60Cnt	Number Trade Lines Cur...	17.92098999	2
TlBalHCPct	Percent Trade Line Balan...	16.13721611	2
TlBadCnt24	Number Trade Lines Bad...	14.09943862	2
InqFinanceCnt24	Number Finance Inquires...	8.23434114	2
CollectCnt	Number Collections	7.9772365	2
DerogCnt	Number Public Derogato...	5.85954434	2
InqCnt06	Number Inquiries 6 Mon...	4.89069085	2
TL75UtilCnt	Number Trade Lines 75 ...	4.73886237	2
TL50UtilCnt	Number Trade Lines 50 ...	3.3284985	2
TLTimeFirst	Time Since First Trade Li...	2.80421077	2
TLSatCnt	Number Trade Lines Cur...	2.07657783	2

Buttons: Edit Rule..., OK, Cancel, Apply, Refresh

Interactive Decision Tree - EMWS1.TREE2_BROWSETREE[EMWS1.TREE2_TRAINSAMPLE]

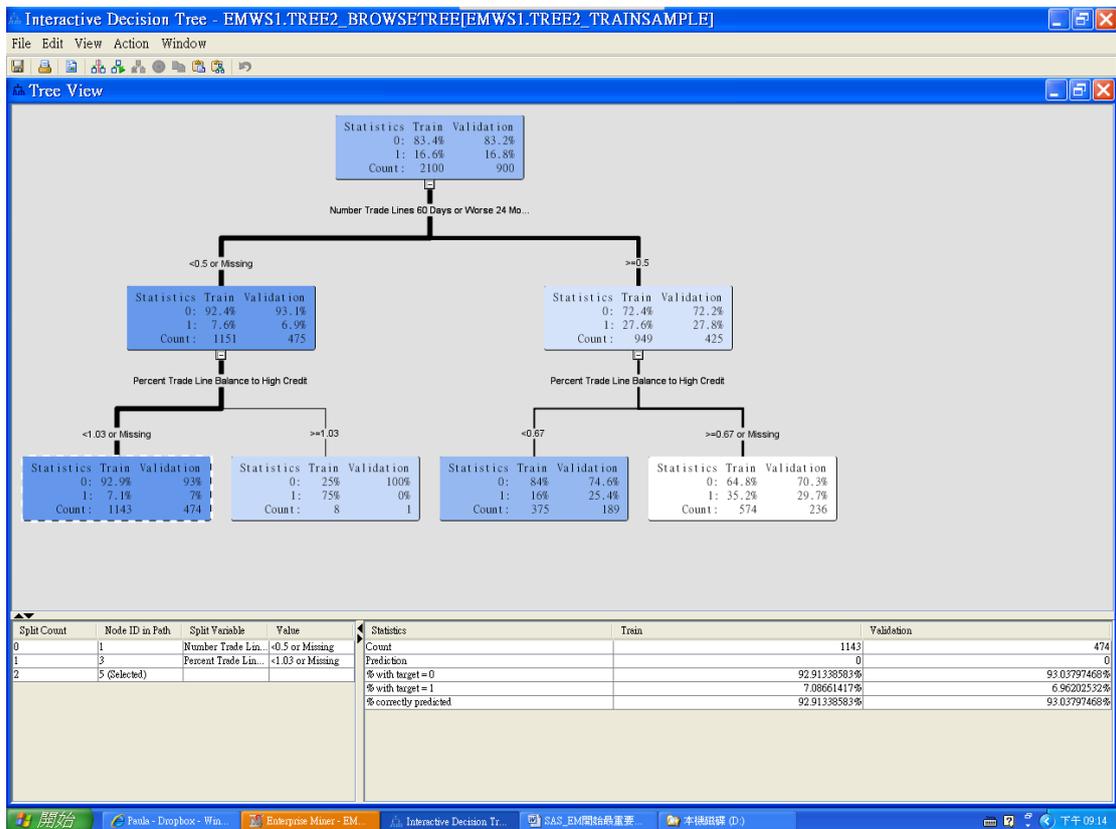
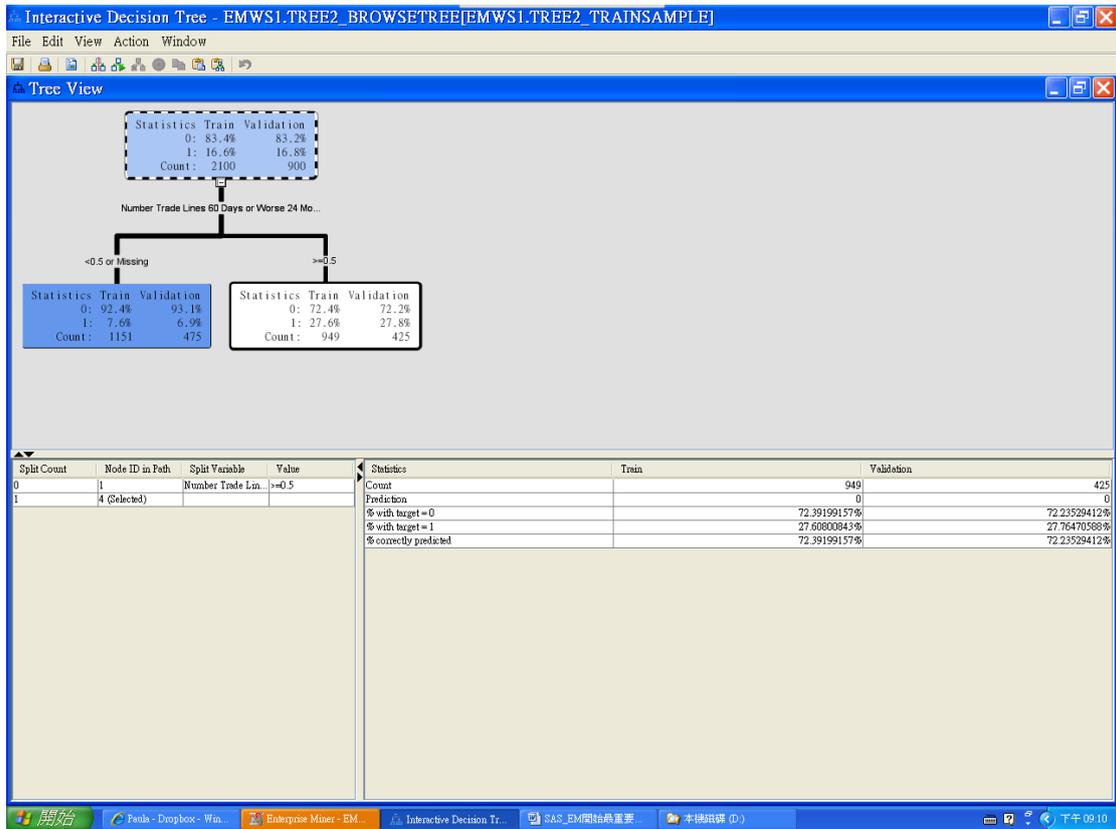
File Edit View Action Window

Tree View

Split Count	Node ID in Path	Split Variable	Value	Statistics	Train	Validation
0	1	Number Trade Lin...	<=0.5	Count	375	189
1	4	Percent Trade Lin...	<=0.67	Average	0.16	0.25396825
2	7 (Selected)			SM. Deviation	0.3666006	0.4453071

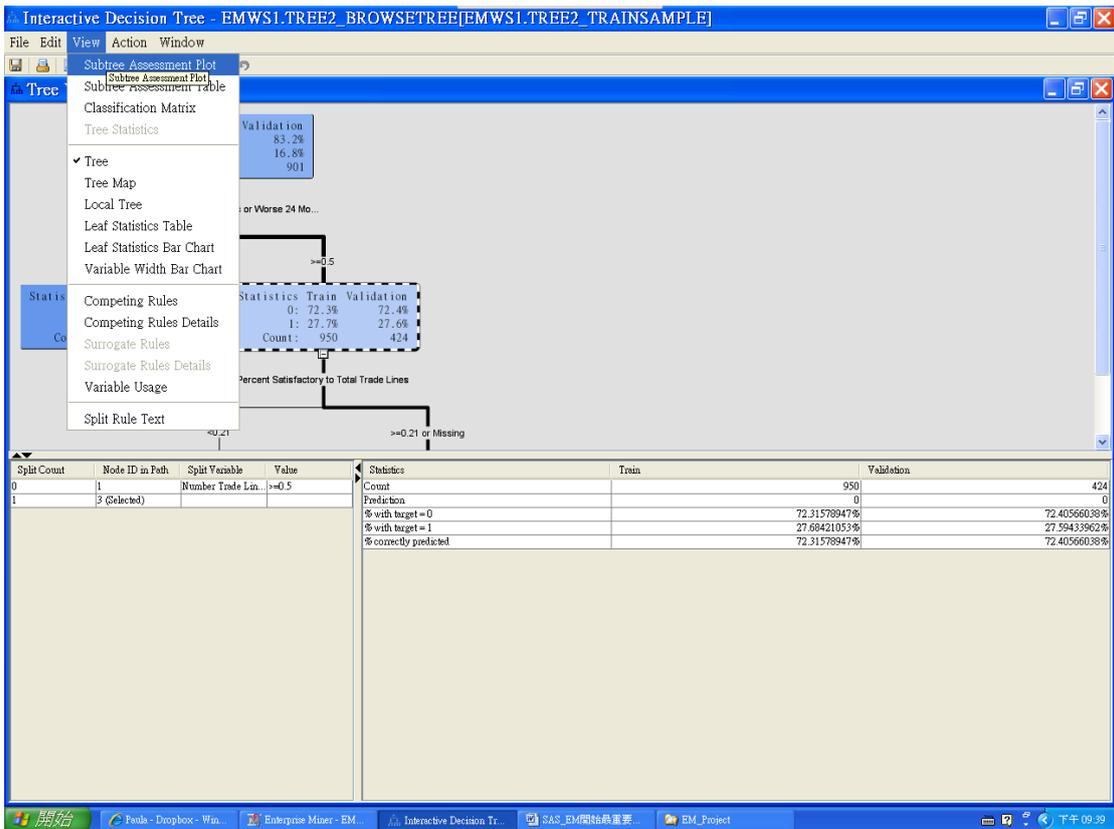
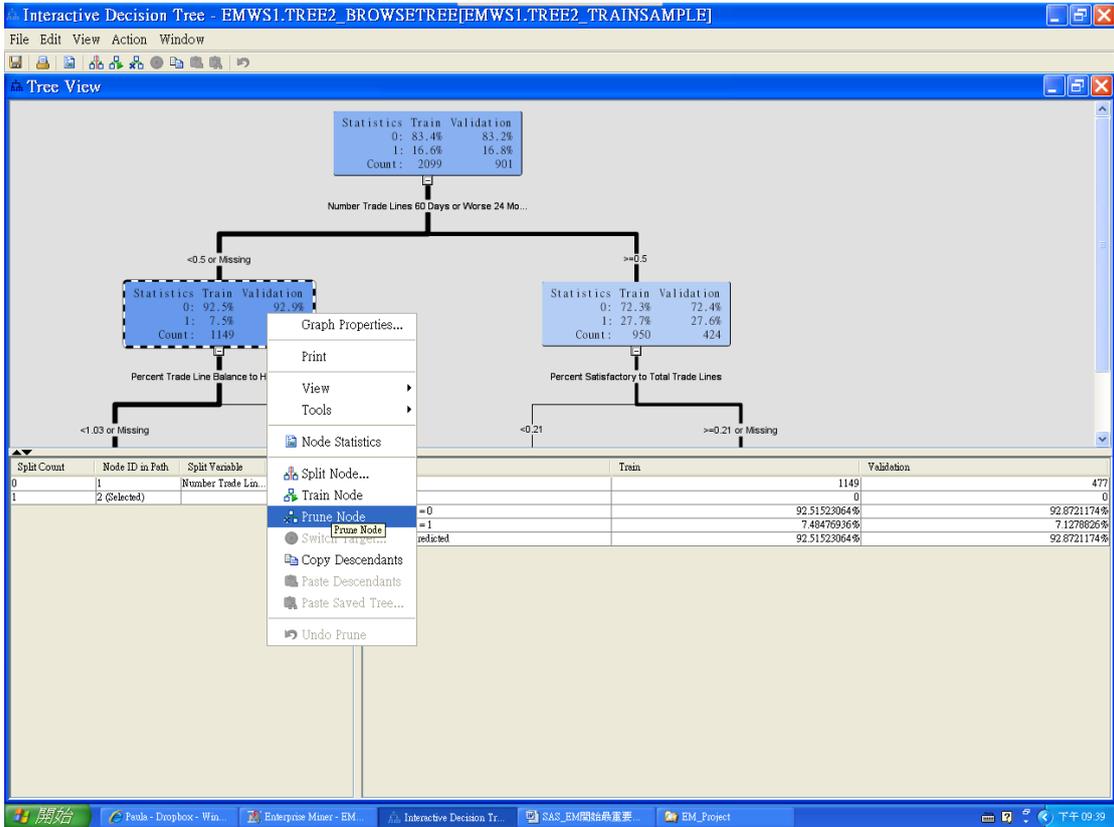
粗線代表：人數最大宗(一般性 rule)，若經費有限就是會先選擇

以下為：target 的 level 屬性設為：Binary 為 0/1 or yes/no (如：0=好 1=壞客戶、yes=會 no=不會買)

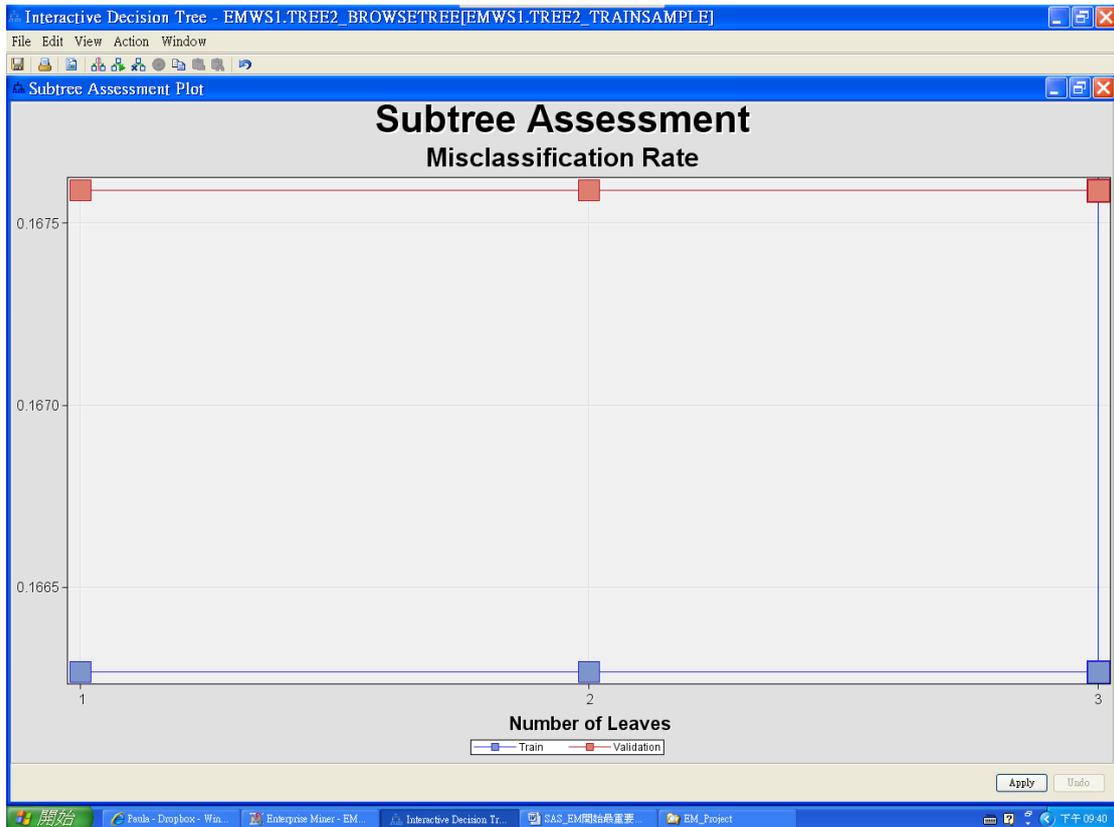


粗線代表：人數最大宗(一般性 rule)，若經費有限就是會先選擇
(粗細代表人數的佔比)

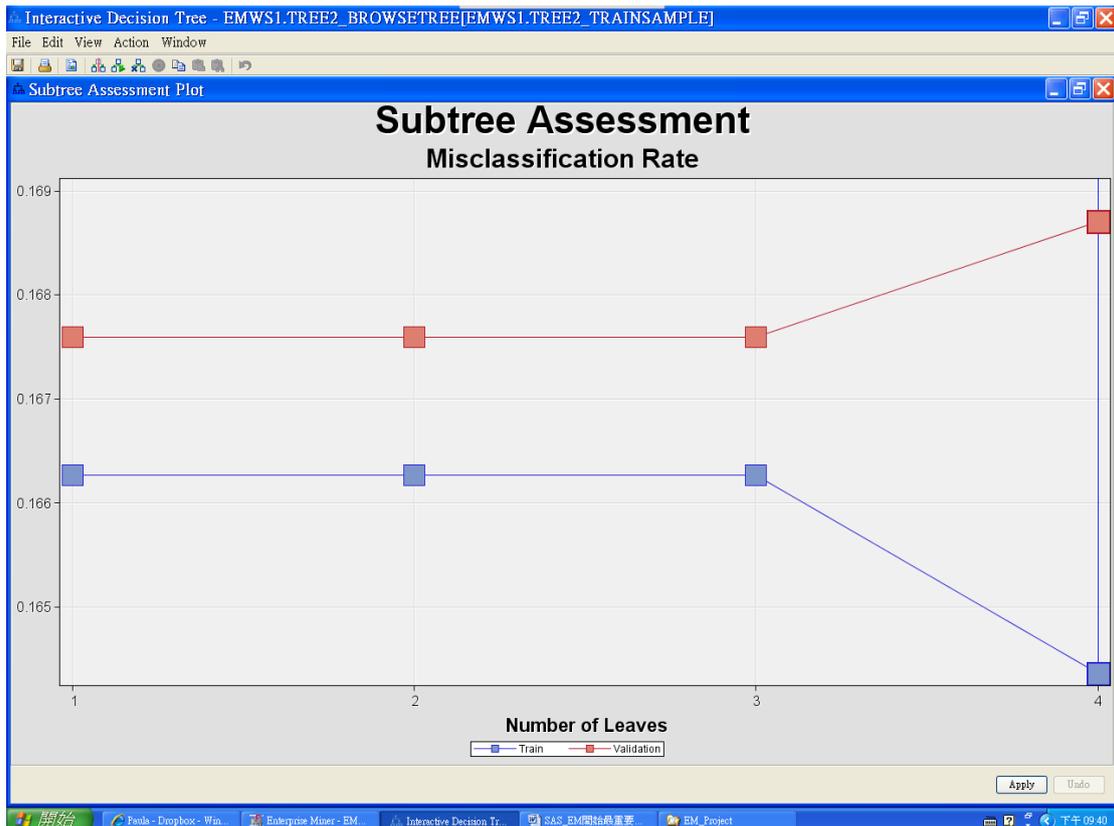
修技：



正常：



不正常或過於極端：以 Validation 為主，二個若發散就為不好，則要從發散後開始修枝（保留 3，從 4 開始刪，Validation 要往 train 方向）



可以從發散的開始修

TLDel60Cnt24 - Interval Split Rule

Target Variable: TARGET

Assign missing values to

A specific branch 1

A separate missing values branch

All branches

Branches

Branch		Split Point
1	<	1
2	<	3
3	>=	3

New split point: Add Branch Remove Branch

OK
Cancel
Apply
Reset

Interactive Decision Tree - EMWS1.TREE2_BROWSETREE[EMWS1.TREE2_TRAINSAMPLE]

File Edit View Action Window

Tree View

Statistics Train Validation

0:	83.4%	83.2%
1:	16.6%	16.8%
Count:	2099	901

Number Trade Lines 60 Days or Worse 24 Mo...

<1 or Missing

0:	92.5%	92.9%
1:	7.5%	7.1%
Count:	1149	477

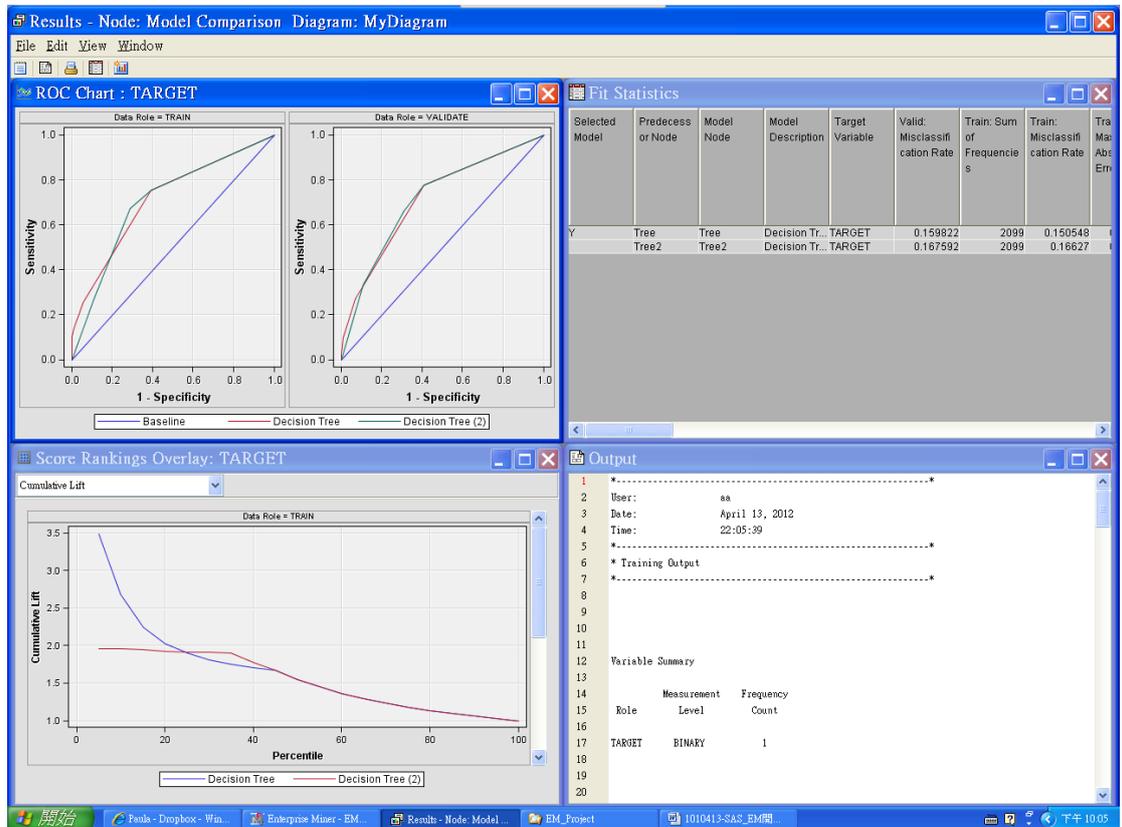
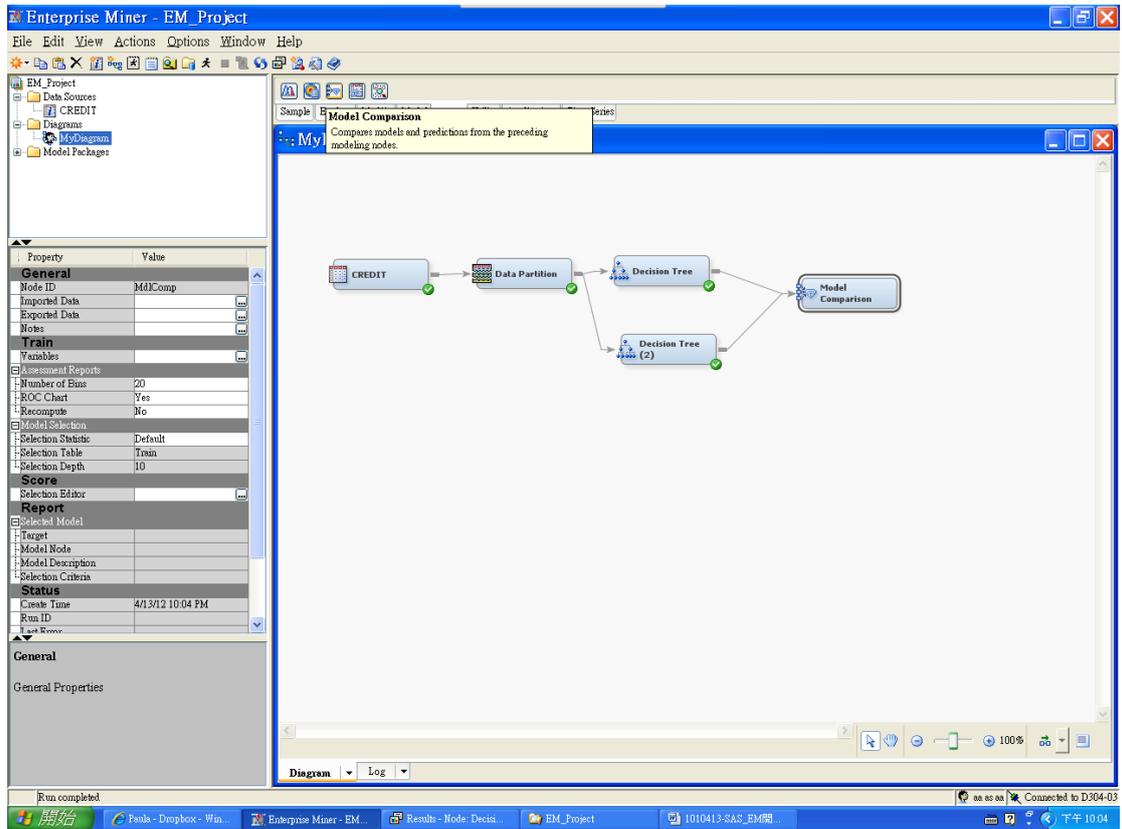
>=1 and <3

0:	74.4%	77.1%
1:	25.6%	22.9%
Count:	664	292

>=3

0:	67.5%	62.1%
1:	32.5%	37.9%
Count:	286	132

Split Count	Node ID in Path	Split Variable	Value	Statistics	Train	Validation
0	1	Number Trade Lin...	>=3	Count	286	132
1	54 (Selected)			Prediction	0	0
				% with target = 0	67.48251748%	62.12121212%
				% with target = 1	32.51748252%	37.87878788%
				% correctly predicted	67.48251748%	62.12121212%



References:

<http://mail.tku.edu.tw/myday/teaching.htm#1002BI>