

경영데이터분석기초

- SPSS, Excel을 활용한 통계분석 -

유 진 호

jhyoo@smu.ac.kr

$$X_1 \dots X_n \rightarrow Y \in \{0, 1\}$$

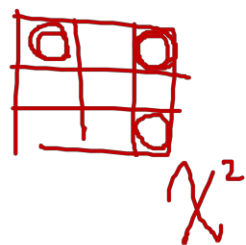
$$\frac{1}{0}$$

의사결정나무(Decision Tree)는 의사결정 규칙을 나무 구조로 나타내어 전체 자료를 몇 개의 소그룹으로 분류하거나 예측할 때 수행하는 분석방법이다. Decision Tree는 시각적이고 명시적인 방법으로 의사결정 과정과 결정된 의사를 보여주는 데 사용된다. Decision Tree 분석방법은

지도 분류 학습에서 가장 유용하게 사용되고 있는 기법 중 하나로서, 노드(집단)내에서는 동질성이 커지고, 노드(집단)간에는 이질성이 가장 커지도록 개체를 분류한다. $\rightarrow F = \frac{\text{집단간 제곱합}}{\text{집단내 제곱합}}$ \uparrow

Supervised
① ②
↓ supervised
unsupervised.

목표변수가 이산형인 경우 상위노드에서 가지분할을 수행할 때, 분류 기준변수와 분류 기준값의 선택방법으로 카이제곱 통계량의 p-값, 지니 지수(Gini index), 엔트로피 지수(Entropy index) 등이 사용된다.

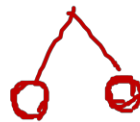


$$F = \frac{\text{between group variance}}{\text{within group variance}}$$





$$p \rightarrow 0$$



CHAID (Chi-square Automatic Interaction Detection) 알고리즘은 최적의 분할, 즉 최적의 예측변수를 선택하는데 있어 카이제곱 검정 또는 F-검정을 이용하여 데이터를 분리하는 방법을 사용한다. 범주형 변수에 대해서는 카이제곱 검정을 사용하고, 연속형 변수에 대해서는 F-검정을 사용한다.

$$F = \frac{\text{Between}}{\text{within}} \uparrow$$

CART (Classification and Regression Trees) 알고리즘은 지니 지수 (Gini Index) 또는 분산의 감소량을 사용하여 나무의 가지를 이진(Binary) 분리한다. 범주형 변수에 대해서는 지니 지수를 사용하고, 연속형 변수에 대해서는 분산의 감소량을 사용한다.

bank.sav [DataSet1] - IBM SPSS Statistics Data Editor

File Edit View Data Transform **Analyze** Direct Marketing Graphs Utilities Add-ons Window Help

Reports
Descriptive Statistics
Tables
Compare Means
General Linear Model
Generalized Linear Models
Mixed Models
Correlate
Regression
Loglinear
Neural Networks
Classify ✓
Dimension Reduction
Scale
Nonparametric Tests
Forecasting
Survival
Multiple Response
Missing Value Analysis...
Multiple Imputation
Complex Samples
Simulation...
Quality Control
ROC Curve...
IBM SPSS Amos...

TwoStep Cluster...
K-Means Cluster...
Hierarchical Cluster...
Tree... ✓
Discriminant...
Nearest Neighbor...

	age	job	default	balance	housing	loan	contact	day	month
1	30	unemployed	no	1787	no	no	cellular	19	oct
2	33	services	no	4789	yes	yes	cellular	11	may
3	35	management	no	1350	yes	no	cellular	16	apr
4	30	management	no	1476	yes	yes	unknown	3	jun
5	59	blue-collar	no	0	yes	no	unknown	5	may
6	35	management	no	747	no	no	cellular	23	feb
7	36	self-employed	no	307	yes	no	cellular	14	may
8	39	technician	no			no	cellular	6	may
9	41	entrepreneur	no			no	unknown	14	may
10	43	services	yes		yes	yes	cellular	17	apr
11	39	services	no			no	unknown	20	may
12	43	admin.	no			no	cellular	17	apr
13	36	technician	no			no	cellular	13	aug
14	20	student	no			no	cellular	30	apr
15	31	blue-collar	no	300	yes	yes	cellular	29	jan
16	40	management	no	194	no	yes	cellular	29	aug
17	56	technician	no	4073	no	no	cellular	27	aug
18	37	admin.	no	2317	yes	no	cellular	20	apr
19	25	blue-collar	no	-221	yes	no	unknown	23	may
20	31	services	no	132	no	no	cellular	7	jul
21	38	management	no	0	yes	no	cellular	18	nov
22	42	management	no	16	no	no	cellular	19	nov
23	44	services	no	106	no	no	unknown	12	jun
24	44	entrepreneur	no	93	no	no	cellular	7	jul

(의미있는 이산형 변수,
의미있는 연속형 변수들을 가지고
Decision Tree 분석
.분류표 제공(예측력)

6개로 실습

age
pdays
duration
marital -> maritalG
Housing -> HousingG
Loan -> LoanG

bank_logit.Regression.sav [데이터집합1] - IBM SPSS Statistics Data Editor

파일(F) 편집(E) 보기(V) 데이터(D) 변환(T) 분석(A) **다**irect 마케팅(M) 그래프(G) 유틸리티(U) 창(W) 도움말(H)

보고서(P) >
 기술통계량(E) >
 표 >
 평균 비교(M) >
 일반선행모형(G) >
 일반화 선행 모형(Z) >
 혼합 모형(X) >
 상관분석(C) >
 회귀분석(R) >
 로그선행분석(O) >
 신경망(W) >
분류분석(Y) >
 차원 감소(D) >
 척도(A) >
 비모수 검정(N) >
 예측(T) >
 생존확률(S) >
 다중응답(U) >
 결측값 분석(V) >
 다중 대입(T) >
 복합 표본(L) >
 시뮬레이션... >
 품질 관리(Q) >
 ROC 곡선(Y) >

이단계 군집분석(T)...
 K-평균 군집분석(K)...
 계층적 군집분석(H)...
트리(R)...
 판별분석(D)...
 가장 가까운 이웃(N)...

	age	job	marital	education	experience	balance	housing	loan	contact	day	month	
1	30	unemployed	married	primary	no	1787	no	no	cellular	19	oct	
2	33	services	married	primary	no	4789	yes	yes	cellular	11	may	
3	35	management	single	primary	no	1350	yes	no	cellular	16	apr	
4	30	management	married	primary	no	1476	yes	yes	unknown	3	jun	
5	59	blue-collar	married	primary	no	0	yes	no	unknown	5	may	
6	35	management	single	primary	no	747	no	no	cellular	23	feb	
7	36	self-employed	married	primary	no	307	yes	no	cellular	14	may	
8	39	technician	married	primary	no	147	yes	no	cellular	6	may	
9	41	entrepreneur	married	primary	no	221	yes	no	unknown	14	may	
10	43	services	married	primary	no			yes	cellular	17	apr	
11	39	services	married	primary	no			no	unknown	20	may	
12	43	admin.	married	primary	no			no	cellular	17	apr	
13	36	technician	married	primary	no			no	cellular	13	aug	
14	20	student	single	primary	no			no	cellular	30	apr	
15	31	blue-collar	married	primary	no			yes	cellular	29	jan	
16	40	management	married	primary	no			yes	cellular	29	aug	
17	56	technician	married	primary	no			no	cellular	27	aug	
18	37	admin.	single	primary	no	2317	yes	no	cellular	20	apr	
19	25	blue-collar	single	primary	no	-221	yes	no	unknown	23	may	
20	31	services	married	primary	no	132	no	no	cellular	7	jul	
21	38	management	divorced	primary	no	0	yes	no	cellular	18	nov	
22	42	management	divorced	primary	no	16	no	no	cellular	19	nov	
23	44	services	single	primary	no	106	no	no	unknown	12	jun	
24	44	entrepreneur	married	primary	no	93	no	no	cellular	7	jul	
25	26	housemaid	married	tertiary	no	543	no	no	cellular	30	jan	
26	41	management	married	tertiary	no	5883	no	no	cellular	20	nov	
27	55	blue-collar	married	primary	no	627	yes	no	unknown	5	may	

6개로 실습

age
 pdays
 duration
 marital -> maritalG
 Housing -> HousingG
 Loan -> LoanG

의사 결정 나무

변수(V):

- job
- education
- default
- balance
- contact
- day
- month
- campaign
- previous
- poutcome
- yG
- maritalG
- eduG
- housingG

종속변수(D):

y

범주(C):

독립변수(I):

- age
- pdays
- duration
- marital
- housing
- loan

☐ 첫 번째 변수 적용(F)

영향 변수(N):

확장 방법(W):

CHAID

변수를 마우스 오른쪽 단추로 클릭하여 변수 목록에서 특정 수준을 변경하십시오

확인 불여넣기(P) 재설정(R) 취소 도움말

의사 결정 나무: 출력결과

트리(T) 통계량 규칙

☒ 분류 규칙 생성(G) ✓

명령문

- ☒ SPSS Statistics
- ☐ SQL(Q)
- ☐ 단순 텍스트(M)
- ☒ 변수 및 변수값 설명 사용(U)

유형

- ☒ 케이스에 값 할당(A)
- ☐ 케이스 선택(L)
- ☒ SPSS Statistics 및 SQL 규칙에 서로게이트 포함

☐ 파일에 규칙 내보내기(X)

파일(F):

찾아보기(W)...

노드

- ☒ 모든 터미널 노드(T)
- ☐ 최고 터미널 노드(B)
- 노드 수(N):
- ☐ 지정된 케이스 퍼센트까지의 최고 터미널 노드(E)
- 퍼센트(P):
- ☐ 지수 값이 분리점 값과 크거나 같은 터미널 노드(R)
- 최소 지수 값(V):
- ☐ 모든 노드(O)

계속 취소 도움말

분류

감시됨	예측		
	no	yes	정확도(%)
no	3907	93	97.7%
yes	415	106	20.3%
전체 퍼센트	95.6%	4.4%	88.8%

성장방법: CHAID

종속변수: y

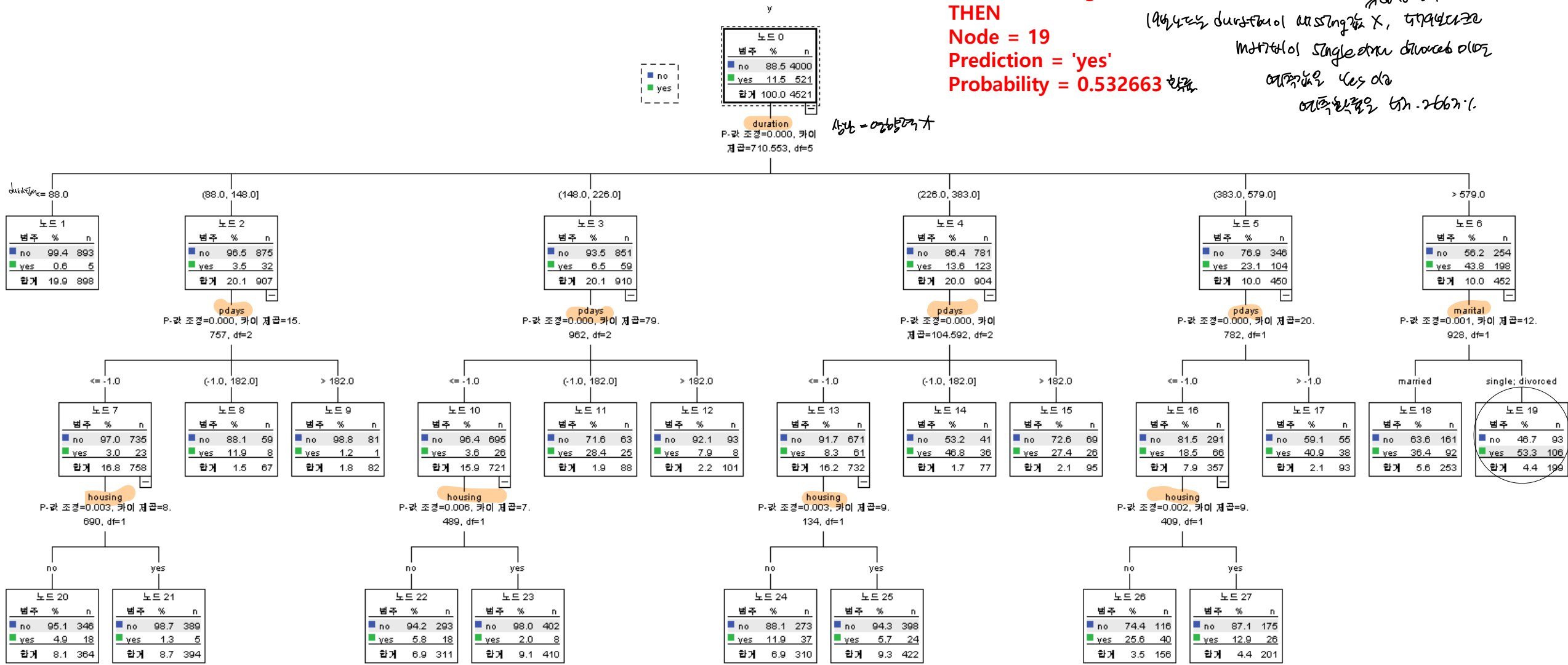
1-2주

$$\frac{93}{447} \approx 11.2\%$$

/* Node 19 */.
IF (duration NOT MISSING AND (duration > 579)) AND
(marital = "single" OR marital = "divorced")
THEN
Node = 19
Prediction = 'yes'
Probability = 0.532663

duration이 missing한 X, 그러면
 marital이 single or divorced 이고
 그러면 yes 다
 그러면 확률은 53.2663%.

성별 = 여성일지 여부



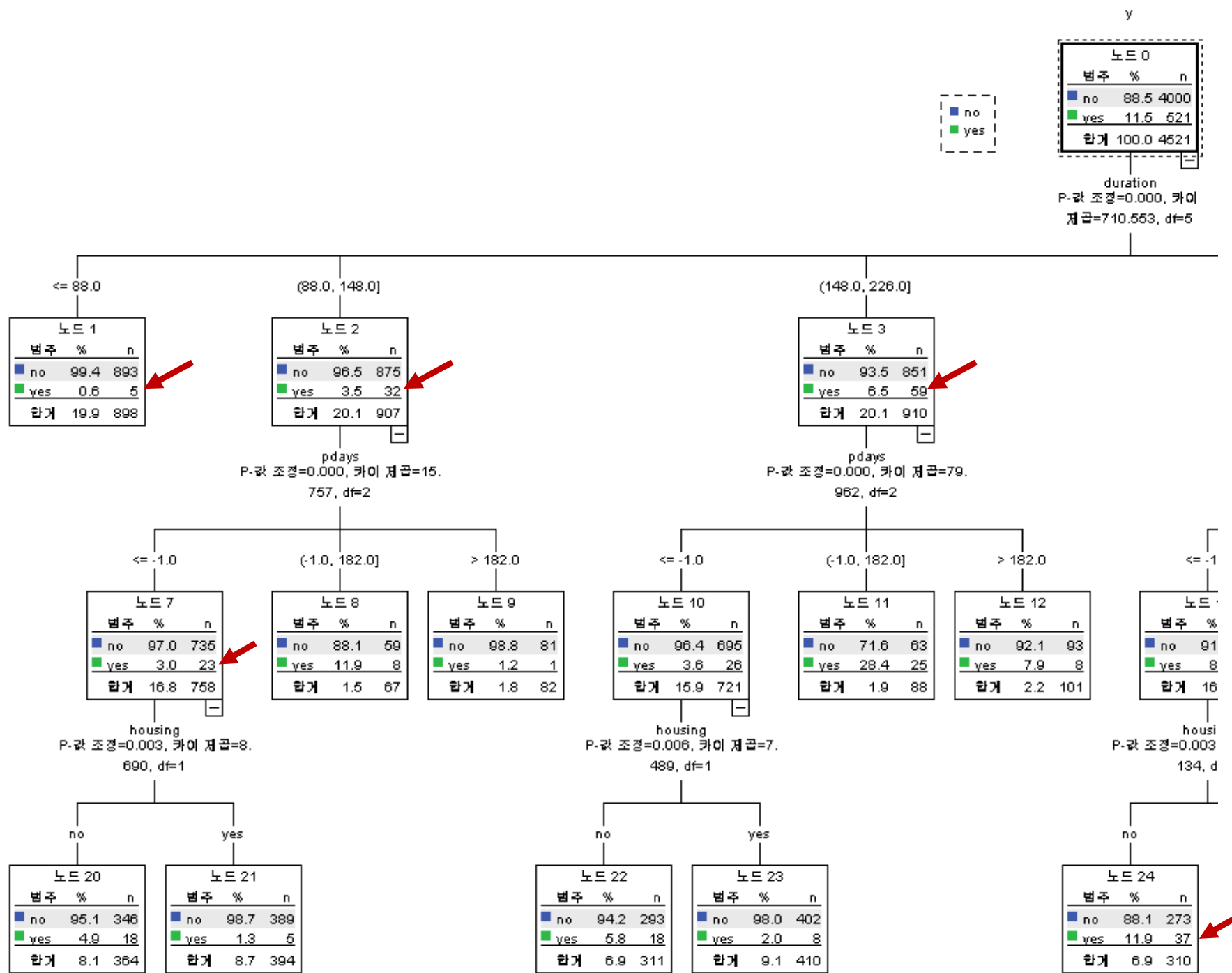
Decision Tree 이용

415 + 93 = 508

분류			
감시됨	예측		
	.00	1.00	정확도(%)
.00	3907	93	97.7%
1.00	415	106	20.3%
전체 퍼센트	95.6%	4.4%	88.8%
성장방법: CHAID			
종속변수: YG			

415,
93,
어디서 나오나?

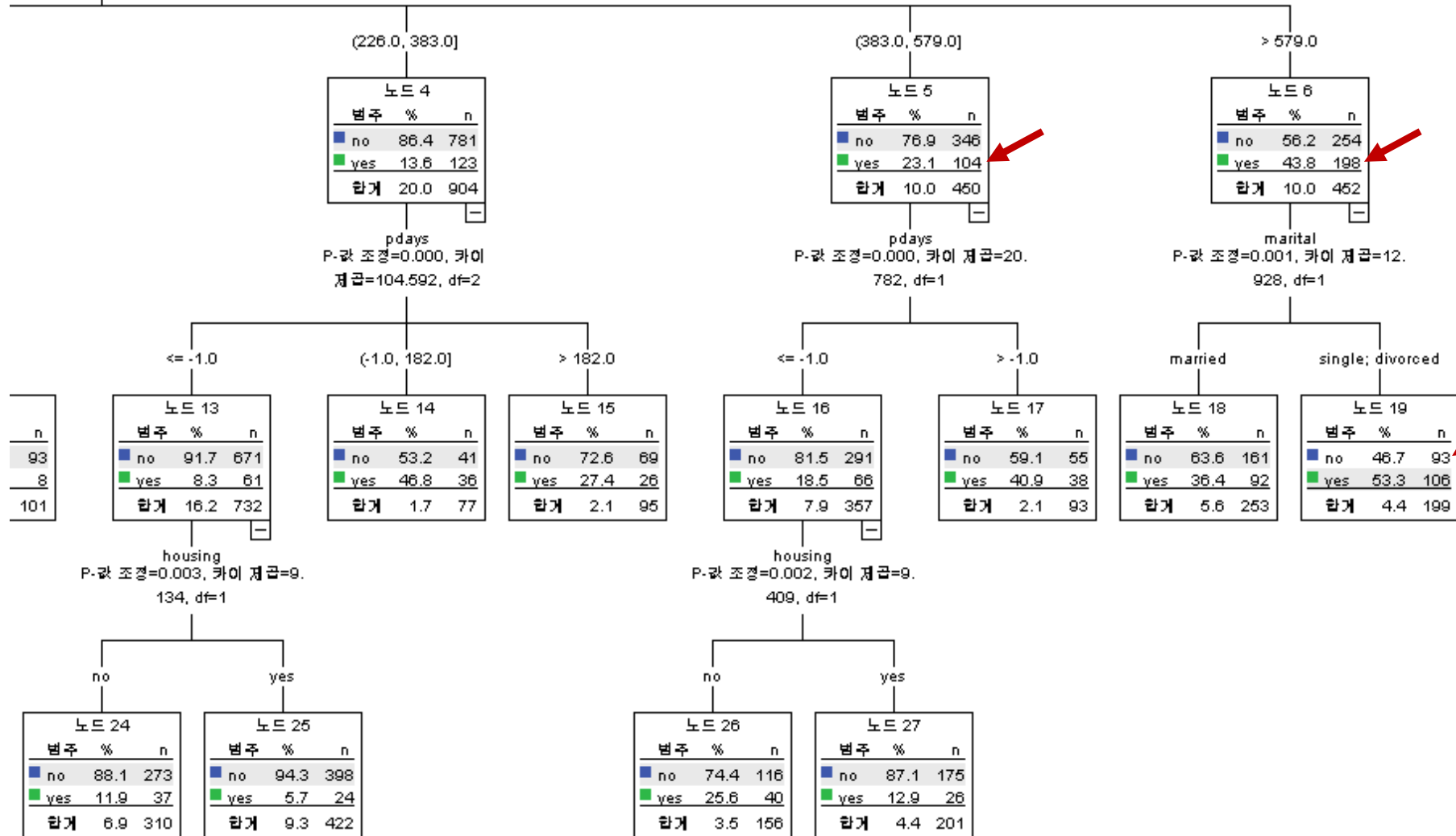
No 예측 but Yes
Yes 예측 but No.



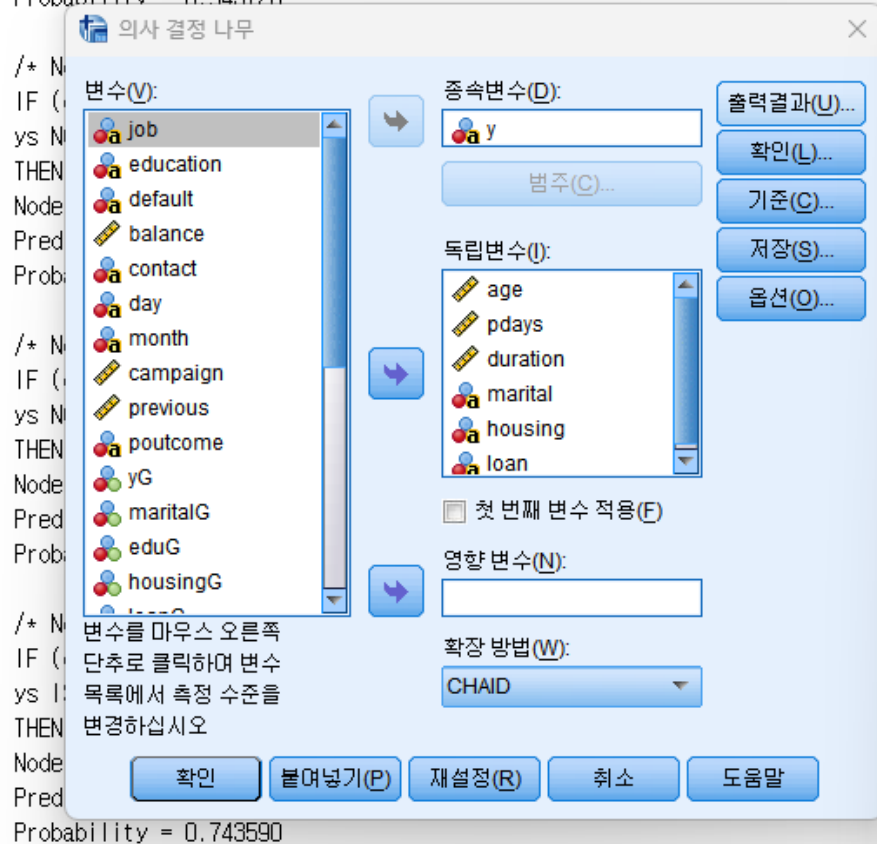
y

노드 0		
변수	%	n
no	88.5	4000
yes	11.5	521
합계	100.0	4521

duration
P-값 조정=0.000, 카이 제곱=710.553, df=5

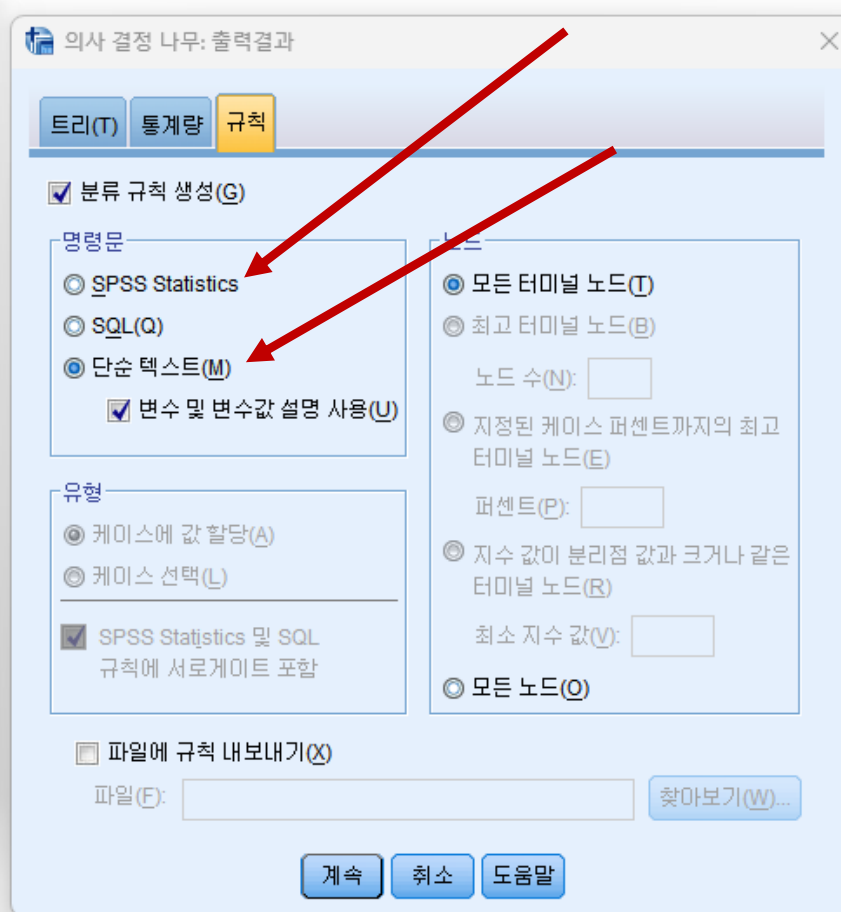


Node = 25
Prediction = 'no'
Probability = 0.943128



```
/* Node 27 */.  
IF (duration NOT MISSING AND (duration > 383 AND duration <= 579)) AND (pda  
ys IS MISSING OR (pdays <= -1)) AND (housing != "no")
```

Rules 생성 및 적용
.SPSS Statistics(적용)
.Simple Text(이해)



숙제

교재처럼 6개 변수를 예측변수로 활용하여 의사결정나무 분석을 실시하고,
그 결과를 Excel파일(*.xls)로 export(내보내기)하여 통계 결과에 대한 해석을 달기

- 1) 가장 큰 영향을 주는 변수는?
- 2) 정확도 해석하기.
- 3) simple text로 수행 후 => 노드1번, 18번, 19번, 20번, 21번 해석하기
- 4) 오류난 것 합계 계산하여 <508개>와 일치하는지 확인하기