

# Python 기반 머신러닝/딥러닝 실습

(전북대학교/데이터 청년 캠퍼스 교육과정 – Machine Learning)

2021. 7. 22 ~ 7. 30

(총 28 시간, 4시간/일 x 7일)

정 준 수 Ph.D

# 과정 목표

복잡한 데이터 구조 패턴을 기계(컴퓨터)로 하여금 스스로 학습하게 하는 인공지능 기계학습 기술을 활용하여 현업의 빅데이터를 지도학습이나 자율학습형태로 분석하고 성능을 평가하여 그 결과를 실제 업무에 적용

1. 분석 방법론, 빅데이터 분석 사례, 파이썬을 이용한 빅데이터 분석 환경 구성
2. 지도학습 모델, 회귀분석, 시계열 분석, 예측 모델 개발, HMM, SVM
3. 지도학습 모델, 분류 분석, 인공신경망, 앙상블모델(XGBoost, LGBM)
4. 의사결정나무, SVM, 자율학습 모델, 데이터 군집 모델 개발, K-means Clustering
5. 빅데이터 모델 평가, 혼돈 매트릭스, AUC, RMSE, F-measure
6. 시각화 도구, 시각화 스토리 텔링, 분석 정보 시각화하기, 시각화 보정

**Deep learning is not like pure mathematics. It is a heavily experimental field, so it's important to be a strong practitioner, not just a theoretician.**

**딥 러닝은 순수한 수학과는 다릅니다. 매우 실험적인 분야이기 때문에 이론가가 아닌 실무자가 되는 것이 중요합니다.**

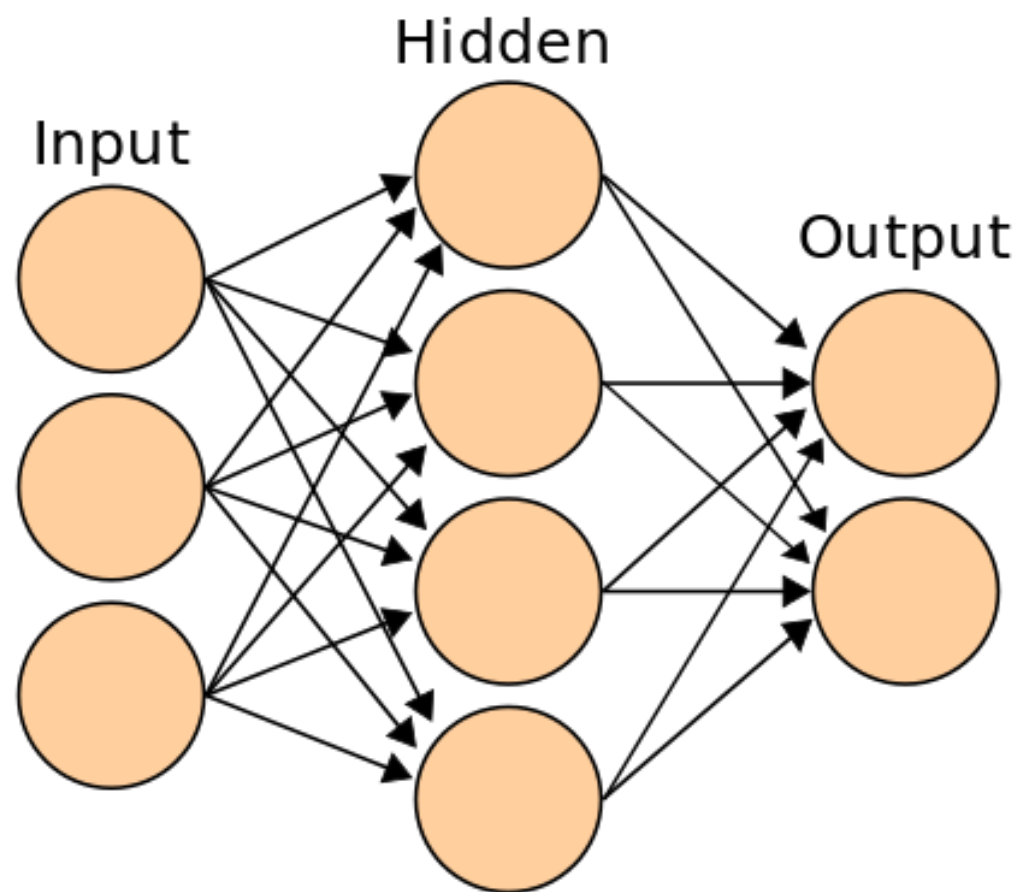
### 머신러닝이란?

인간이 다양한 경험과 시행착오를 통해 지식을 배우는 것처럼, [컴퓨터](#)에게 충분히 많은 데이터를 주고, 거기에서 일반적인 패턴을 찾아내게 하는 방법을 말합니다. (머신러닝의 대표적인 [알고리즘](#)이 딥러닝입니다.)

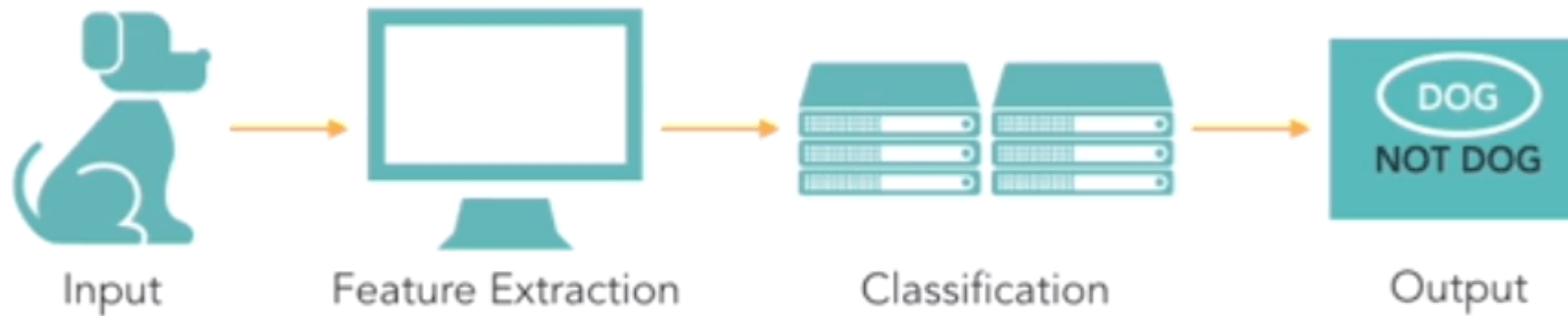
### 딥러닝이란?

머신러닝의 대표적인 학습법이 다음장 그림처럼 여러 층을 거쳐 점점 추상화 단계로 접어드는 알고리즘 형태인 '딥 러닝(Deep Learning)'입니다.

패턴을 찾기 위해선, 수많은 데이터가 있어야 되겠지요. 패턴을 견고하게 만드는 일종의 학습용 훈련데이터 말입니다.



## TRADITIONAL MACHINE LEARNING



## DEEP LEARNING

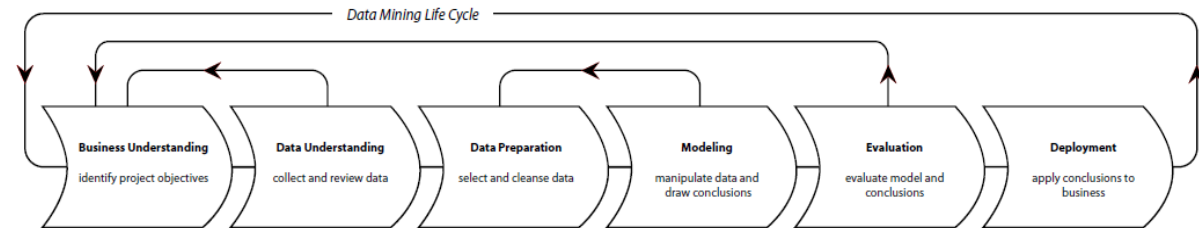


# CRISP-DM (Cross Industry Standard Process for Data Mining)

CRISP-DM(Cross Industry Standard Process for Data Mining)은 데이터 마이닝 전문가가 사용하는 일반적인 접근 방식을 설명한 가장 널리 사용되는 공개 표준 분석 모델입니다.



## Phases



**Determine Business Objectives**  
Background  
Business Objectives  
Business Success Criteria  
(Log and Report Process)

**Assess Situation**  
Inventory of Resources,  
Requirements, Assumptions,  
and Constraints  
Risks and Contingencies  
Terminology  
Costs and Benefits  
(Log and Report Process)

**Determine Data Mining Goals**  
Data Mining Goals  
Data Mining Success Criteria  
(Log and Report Process)

**Produce Project Plan**  
Project Plan  
Initial Assessment of Tools and  
Techniques  
(Log and Report Process)

**Collect Initial Data**  
Initial Data Collection Report  
(Log and Report Process)

**Describe Data**  
Data Description Report  
(Log and Report Process)

**Explore Data**  
Data Exploration Report  
(Log and Report Process)

**Verify Data Quality**  
Data Quality Report  
(Log and Report Process)

**Data Set**  
Data Set Description  
(Log and Report Process)

**Select Data**  
Rationale for Inclusion/  
Exclusion  
(Log and Report Process)

**Clean Data**  
Data Cleaning Report  
(Log and Report Process)

**Construct Data**  
Derived Attributes  
Generated Records  
(Log and Report Process)

**Integrate Data**  
Merged Data  
(Log and Report Process)

**Format Data**  
Reformatted Data  
(Log and Report Process)

**Select Modeling Technique**  
Modeling Technique  
Modeling Assumptions  
(Log and Report Process)

**Generate Test Design**  
Test Design  
(Log and Report Process)

**Build Model Parameter Settings**  
Models  
Model Description  
(Log and Report Process)

**Assess Model**  
Model Assessment  
Revised Parameter  
(Log and Report Process)

**Evaluate Results**  
Align Assessment of Data  
Mining Results with  
Business Success Criteria  
(Log and Report Process)

**Approved Models**  
Review Process  
Review of Process  
(Log and Report Process)

**Determine Next Steps**  
List of Possible Actions  
Decision  
(Log and Report Process)

**Plan Deployment**  
Deployment Plan  
(Log and Report Process)

**Plan Monitoring and Maintenance**  
Monitoring and  
Maintenance Plan  
(Log and Report Process)

**Produce Final Report**  
Final Report  
Final Presentation  
(Log and Report Process)

**Review Project**  
Experience  
Documentation  
(Log and Report Process)

## a visual guide to CRISP-DM methodology

SOURCE CRISP-DM 1.0  
<http://www.crisp-dm.org/download.htm>

DESIGN Nicole Leaper  
<http://www.nicoleleaper.com>



# Machine Learning에서 “예측”이란?

이전에 본적 없는 새로운 데이터에 대한 정확한 출력 예측



# Predictive Analytics

1 무엇을 예측하는가?

2 무엇을 할 것인가?

# Machine Learning - Training

	Input			Output
Example 1	0	0	1	0
Example 2	1	1	1	1
Example 3	1	0	1	1
Example 4	0	1	1	0

New situation	1	0	0	?
---------------	---	---	---	---

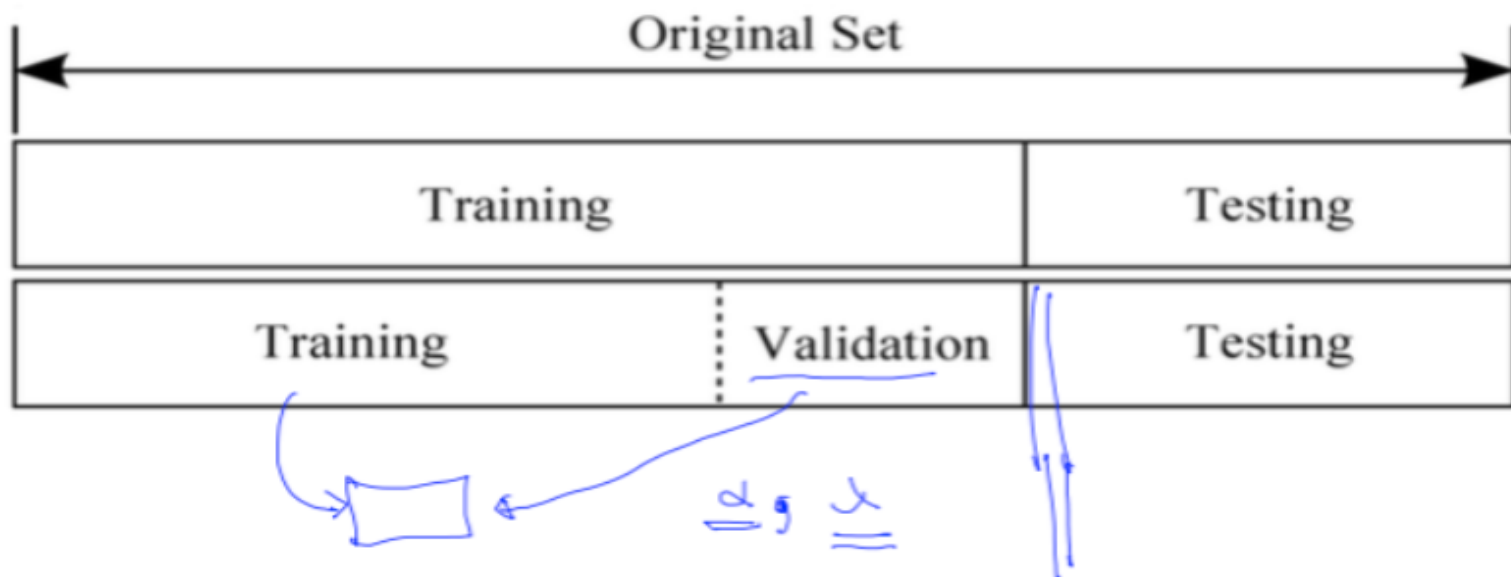


```
from numpy import exp, array, random, dot
training_set_inputs = array([[0, 0, 1], [1, 1, 1], [1, 0, 1], [0, 1, 1]])
training_set_outputs = array([[0, 1, 1, 0]]).T
random.seed(1)
synaptic_weights = 2 * random.random((3, 1)) - 1
for iteration in range(10000):
    output = 1 / (1 + exp(-(dot(training_set_inputs, synaptic_weights))))
    synaptic_weights += dot(training_set_inputs.T, (training_set_outputs - output) * output * (1 - output))
print (1 / (1 + exp(-(dot(array([1, 0, 0]), synaptic_weights)))))
```

출처: <https://medium.com/technology-invention-and-more/how-to-build-a-simple-neural-network-in-9-lines-of-python-code-cc8f23647ca1>

# 머신러닝 학습 데이터의 구성

## Training, validation and test sets



데이터 영역별 구성도다. 데이터가 충분하다면 validation set을 구성하지 않을 이유가 없다. 쉽게 보면 test를 한번 더 한다고 생각해도 좋을 것이다. 머신러닝의 교과서에 해당하는 mnist 예제는 train 55,000개, validation 5,000개, test 10,000개의 dataset으로 구성되어 있다.

## 예측 분석 응용: 이탈 모델링으로 고객 이탈 방지하기

### 1 무엇을 예측하는가?

어느 고객이 떠나갈 것인가.

### 2 무엇을 할 것인가?

떠날 위기에 있는 고객들을 타겟으로 한 고객 유지 마케팅을 수행한다.

# 예측 분석 응용: 부동산 담보대출 채권 가치 추산

## 1 무엇을 예측하는가?

부동산 담보대출 고객 중에서 누가 향후 90일 내에 조기 상환할 것인가.

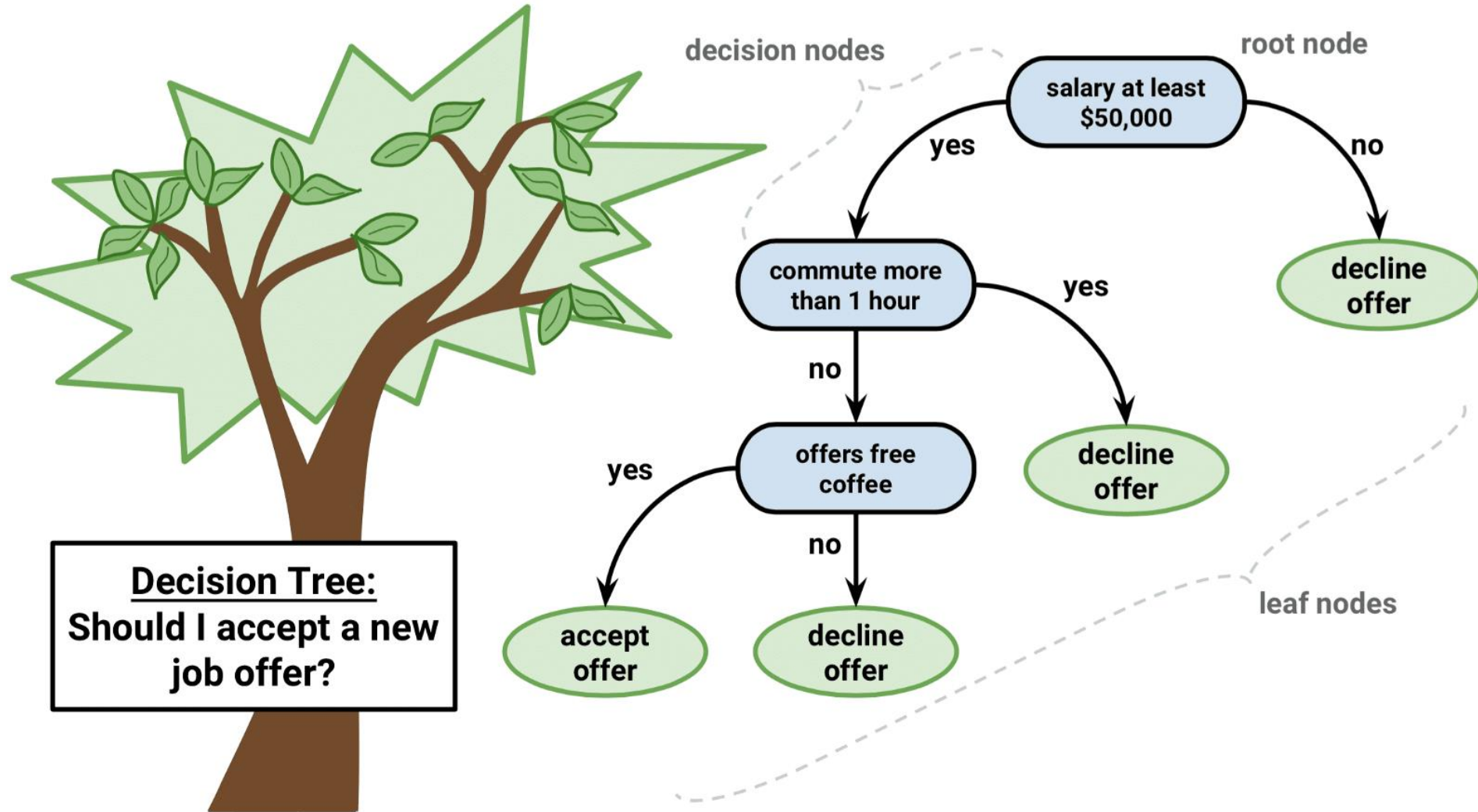
## 2 무엇을 할 것인가?

부동산 담보대출 채권의 가치를 계산한 후 다른 은행에 팔아 넘길지 여부를 결정한다.

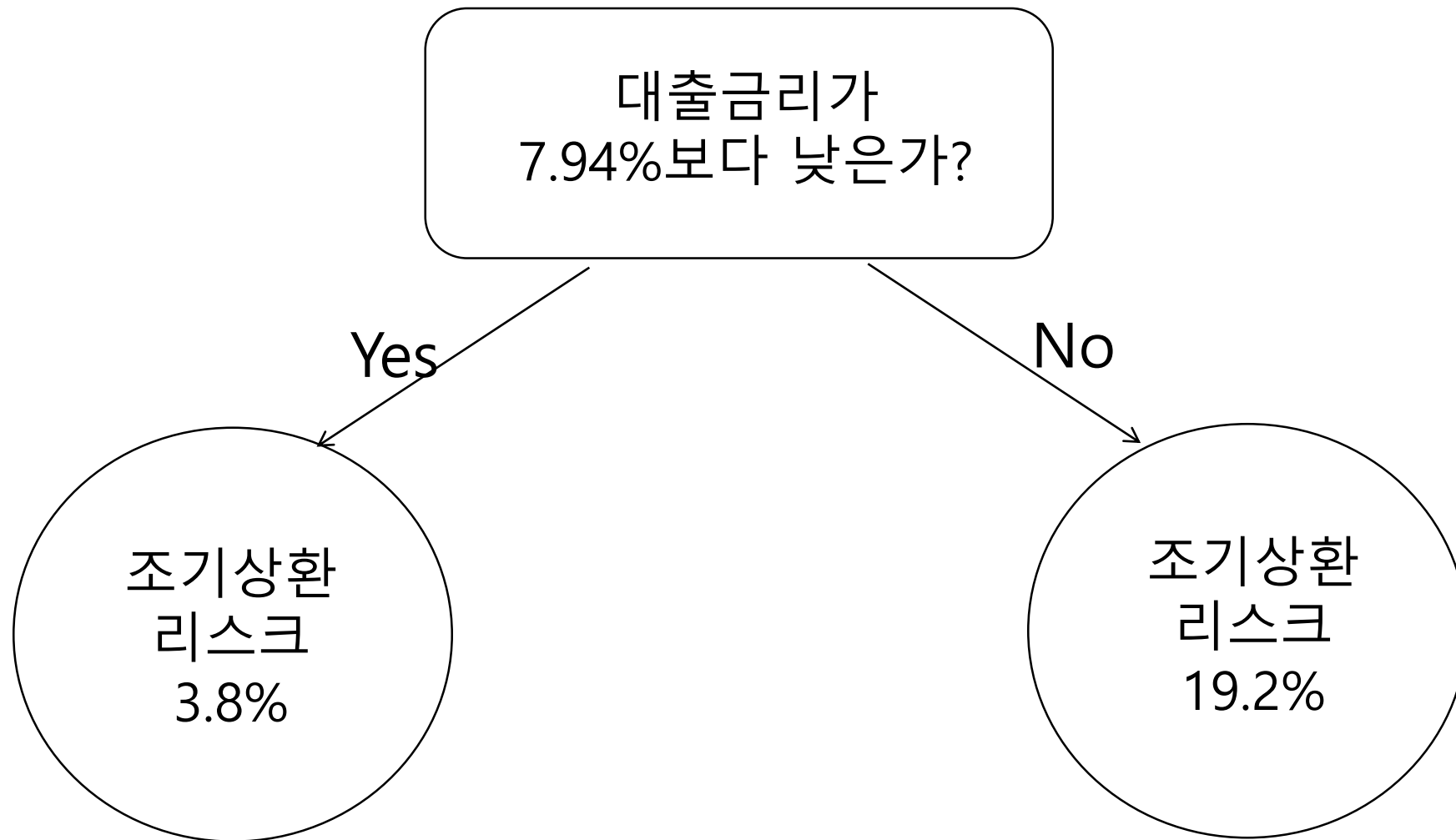
# 나무에서 돈이 자란다!

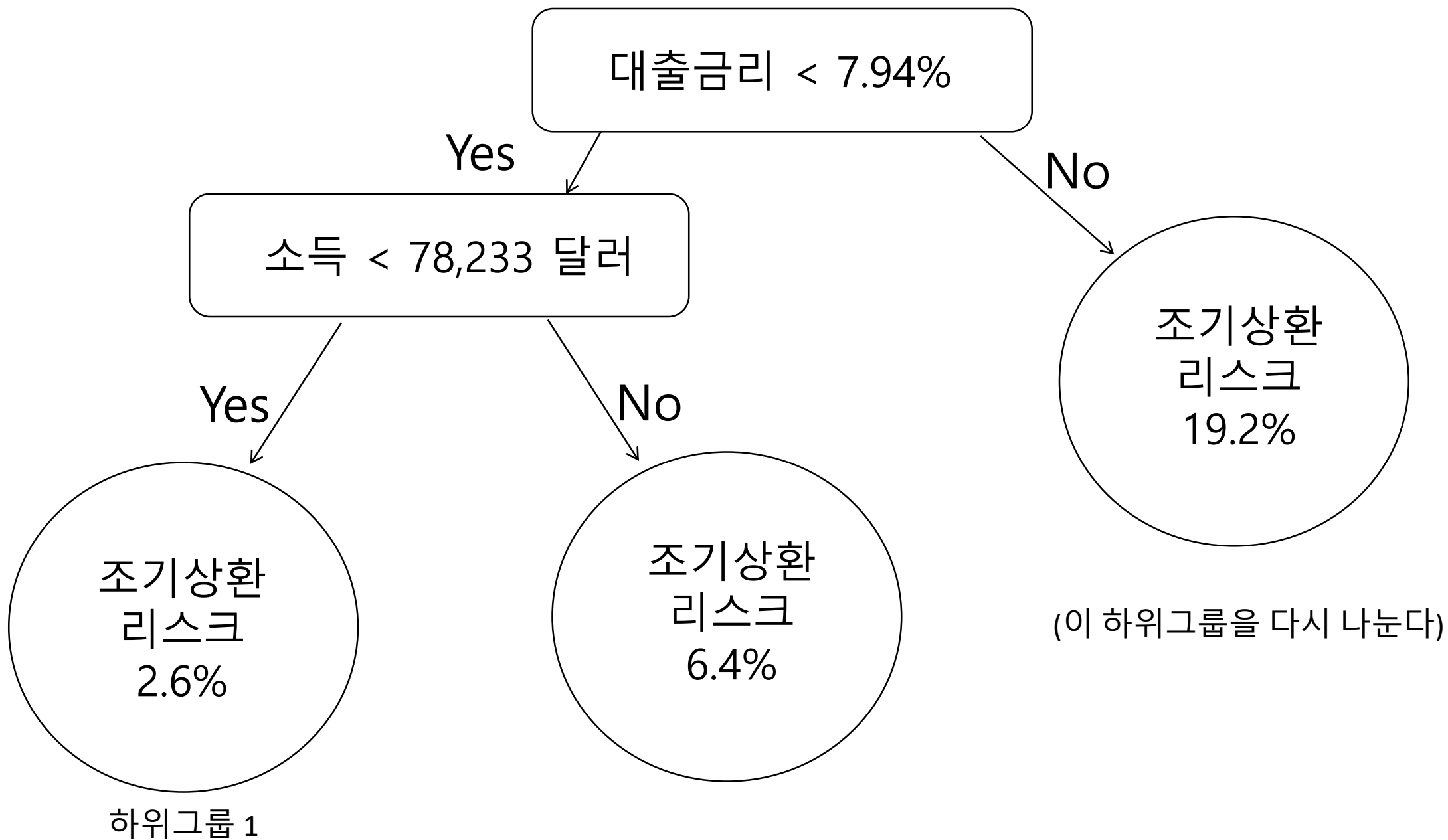
체이스 은행 예측모델은 부동산 담보대출 중에서 실제로  
조기상환된 대출 건들 중 74%를 정확하게 인식해 내어서  
부동산 담보대출 포트폴리오를 성공적으로 관리함.

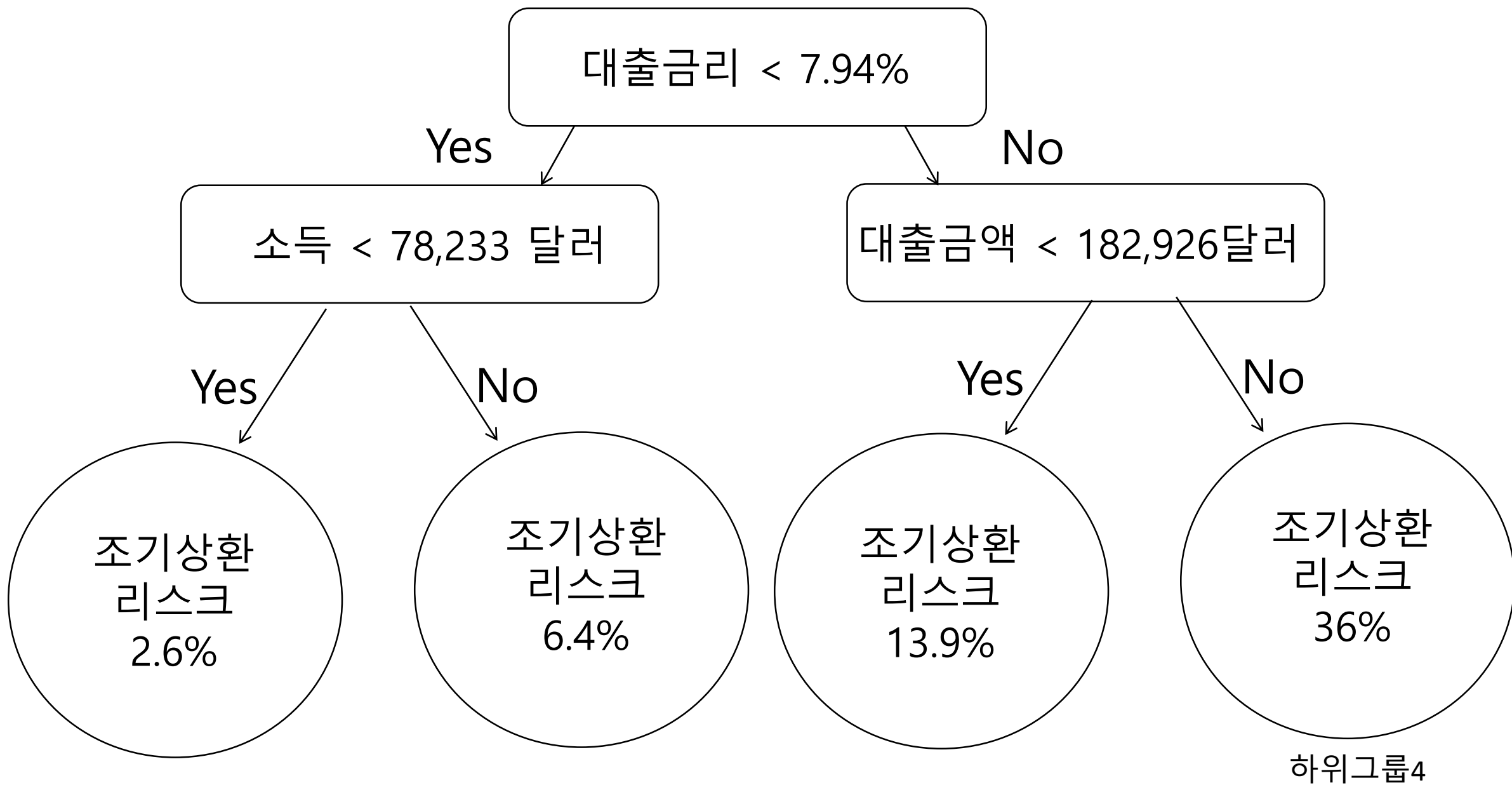
# Decision Tree

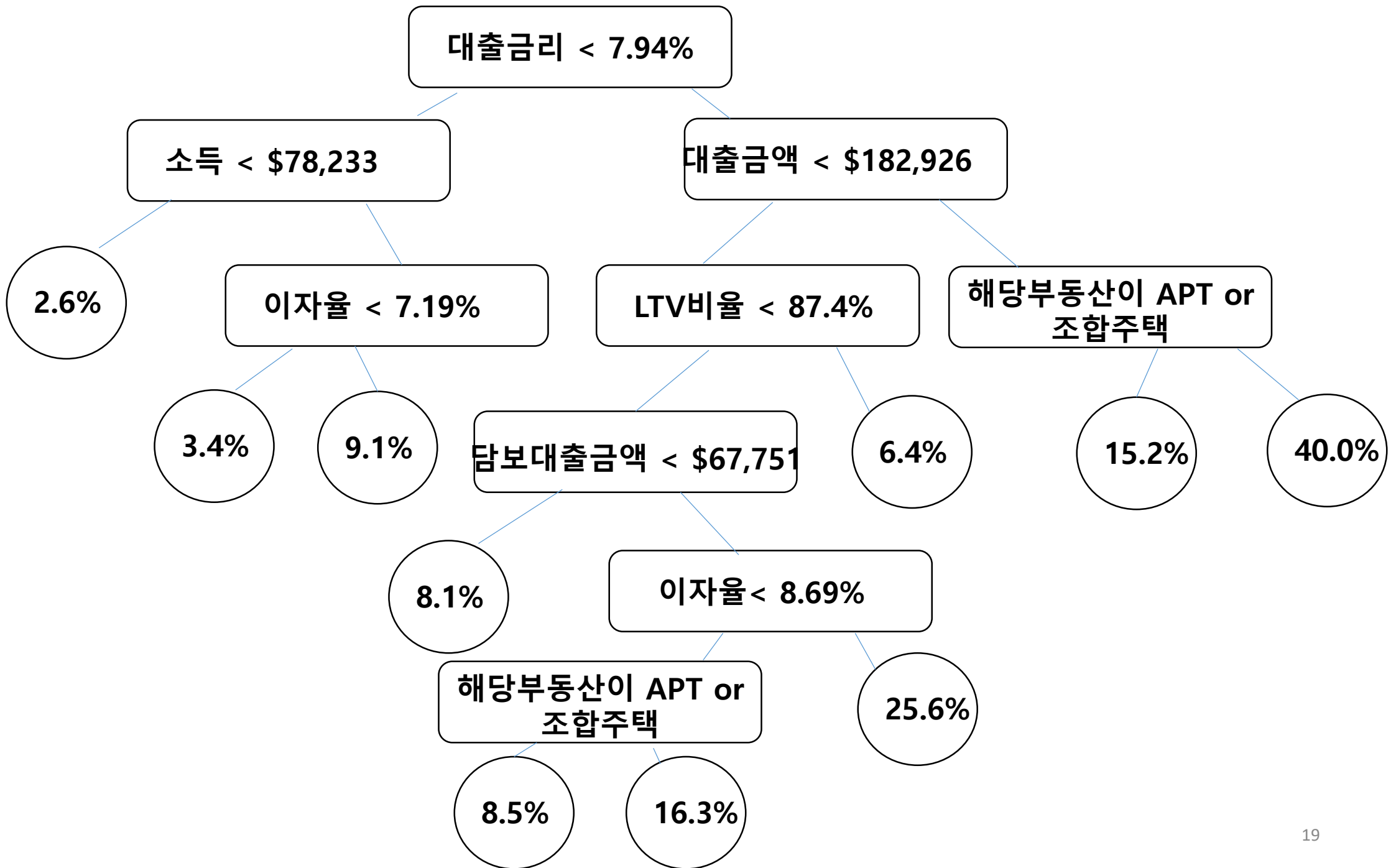












**만약(IF):**

부동산 담보대출 금액이 67,751 달러와 같거나 그보다 더 많고 182,926 달러보다 작다.

**그리고(AND):**

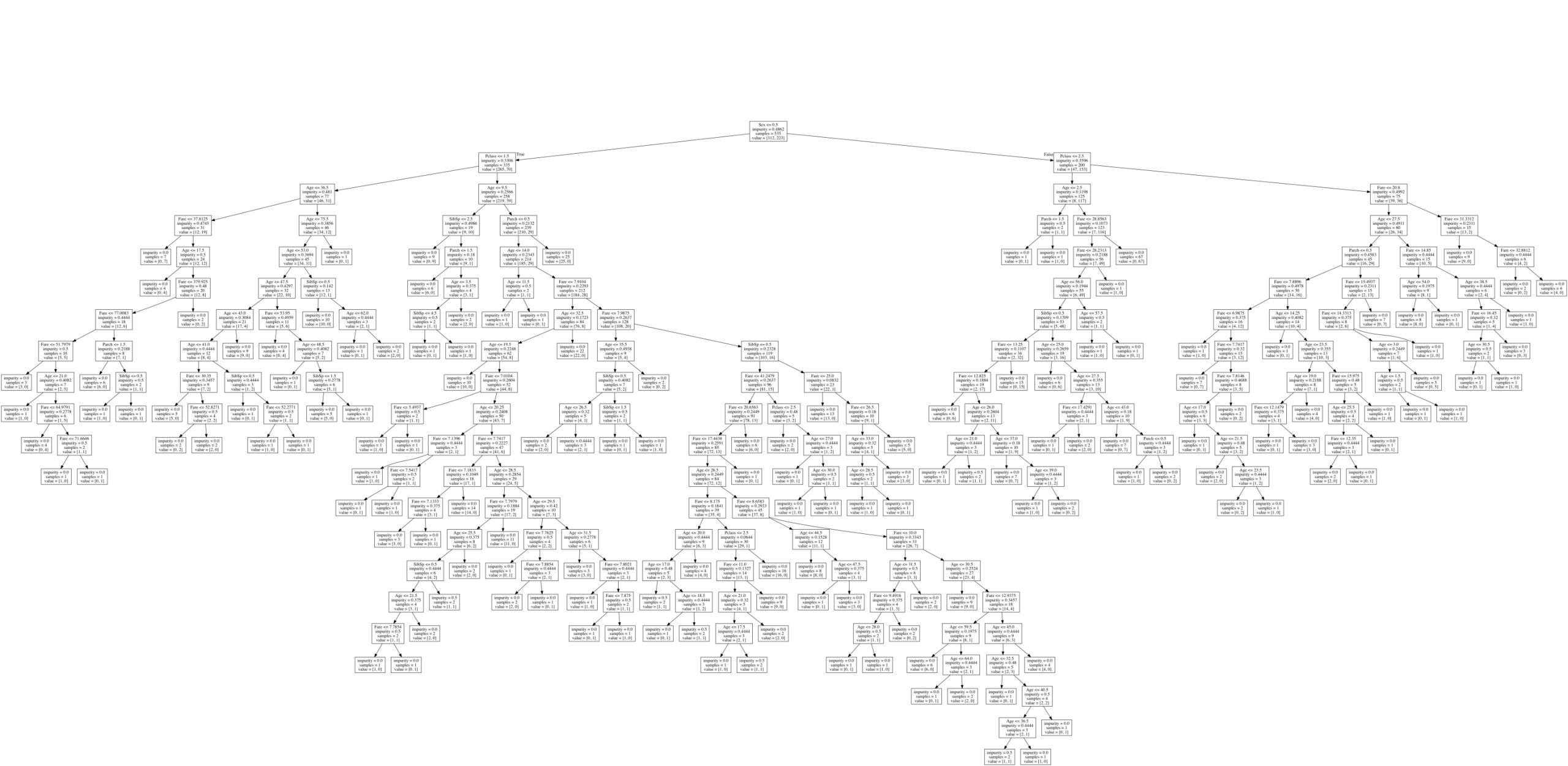
이자율이 8.69%와 같거나 그보다 더 높다.

**그리고(AND):**

부동산 자산가치 대비 대출금액의 비율이 87.4% 보다 작다.

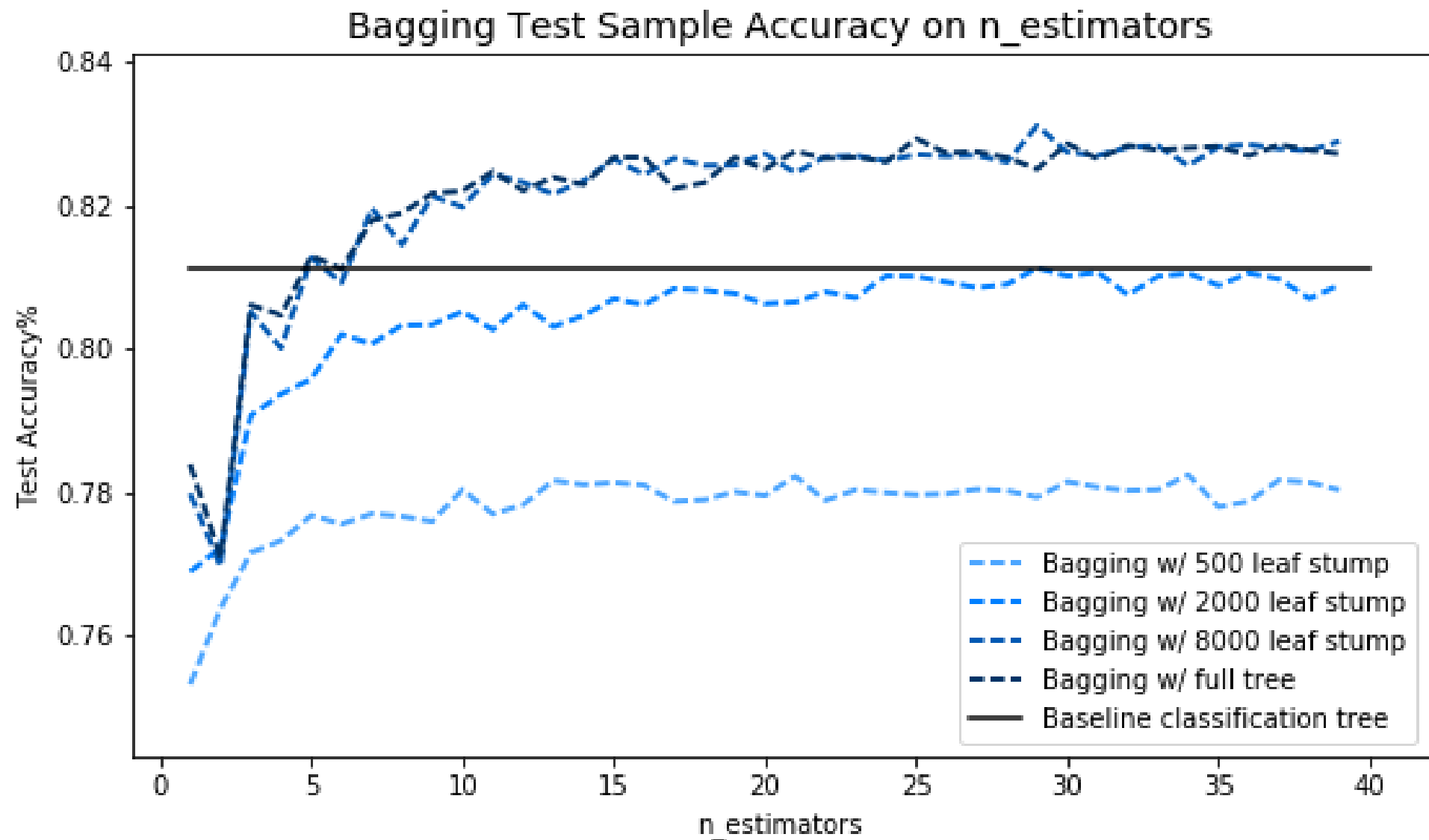
**그러면(THEN):**

조기상환 확률은 25.6% 이다.



## Decision Tree: 20% 기준 다양한 향상도 성과 기록

Decision Tree(Leafs)	20% 기준 향상도
4개의 그룹	2.5
10개의 그룹	3.0
39개의 그룹	3.0





보·이·는·대·로·**밀**·지·마·라!

대/반/전/ 음악추리쇼



Mnet tvN 공동 방송

10월 22일 목요일 | 밤 9시 40분

6명 중에는 음치도 있고 노래를 잘하는 실력자(정상)도 있다.

번호 [ 1, 2, 3, 4, 5, 6 ]

정답 : [음치, 음치, 음치, 음치, 정상, 정상] 가 있다.

누가 음치인지 겉모습만 보고 맞춰야 한다.

감으로 예측을 한다.

예측 : [음치, 음치, 정상, 정상, 정상, 정상]

# Actual Values

1

0

Predicted Values

1

TRUE POSITIVE

You're pregnant

FALSE POSITIVE

You're pregnant

TYPE 1 ERROR

FALSE NEGATIVE

You're not pregnant

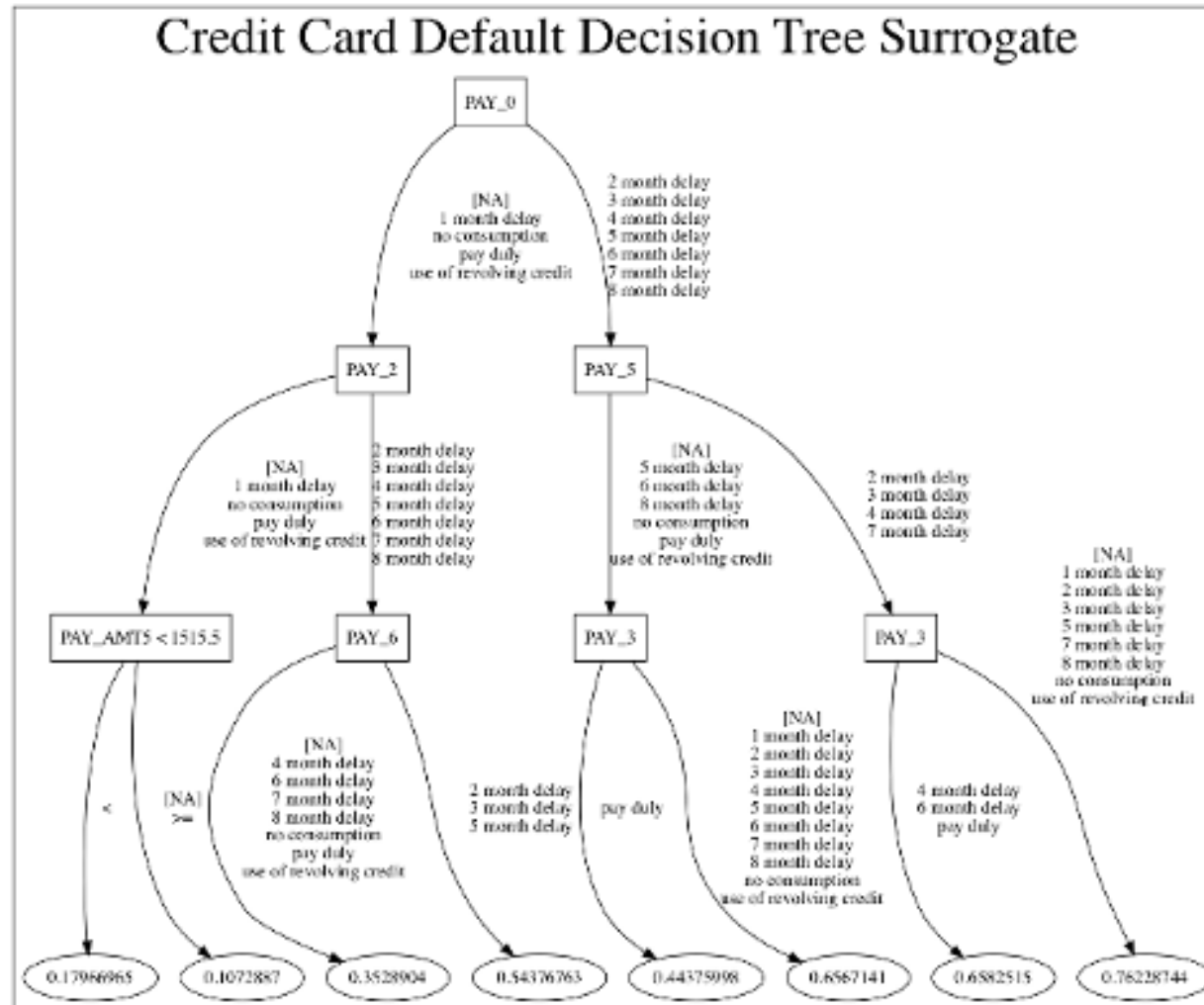
TYPE 2 ERROR

TRUE NEGATIVE

You're not pregnant

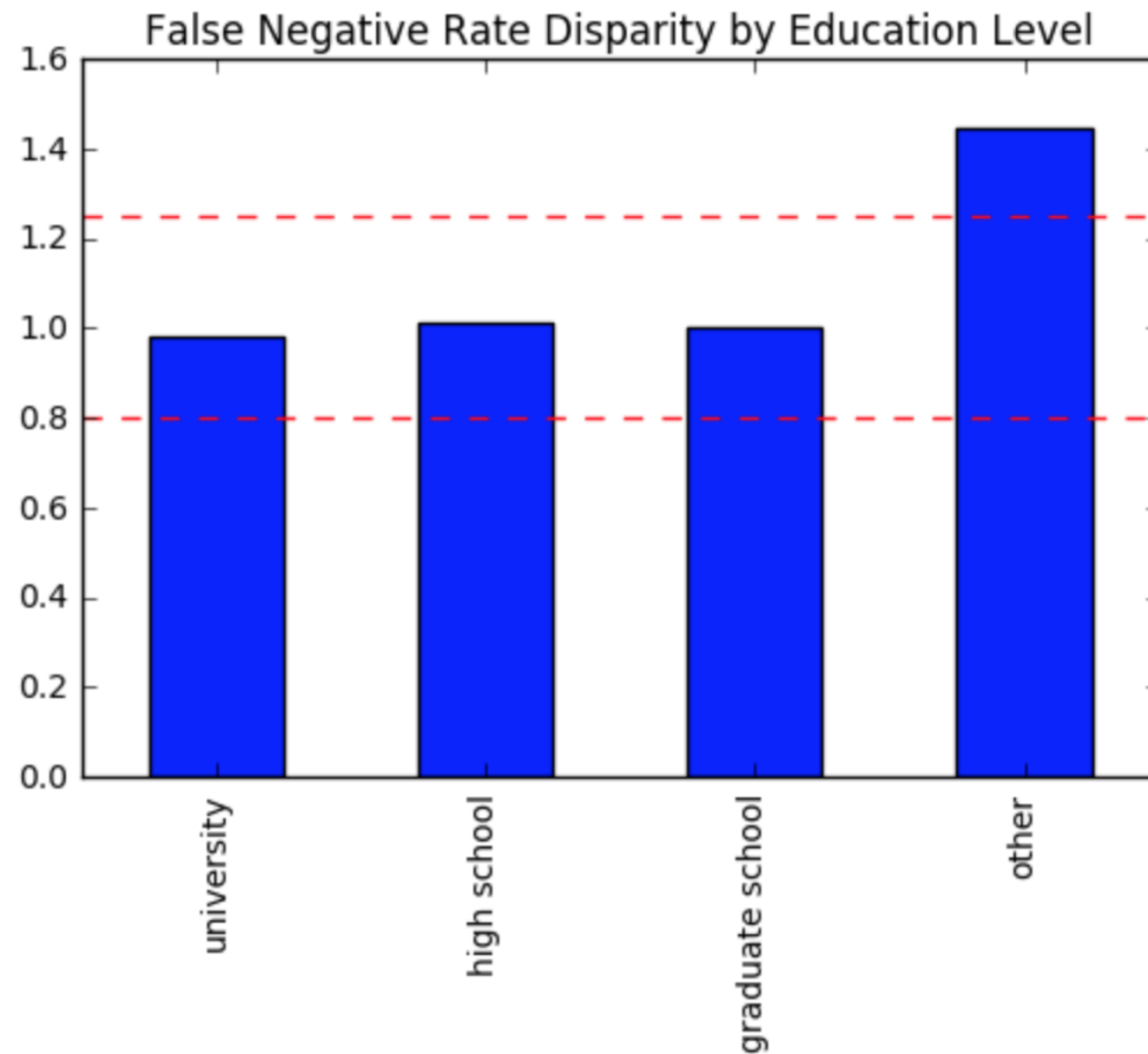
0

# 머신러닝 결과 분석(Decision Tree)



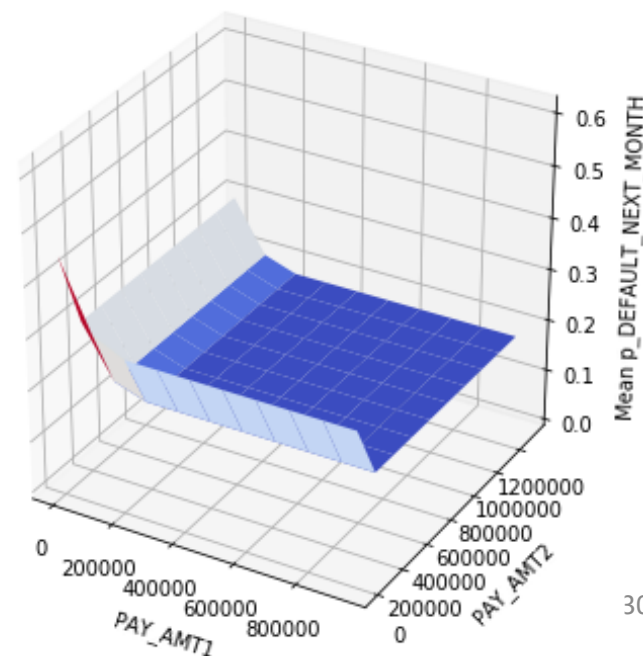
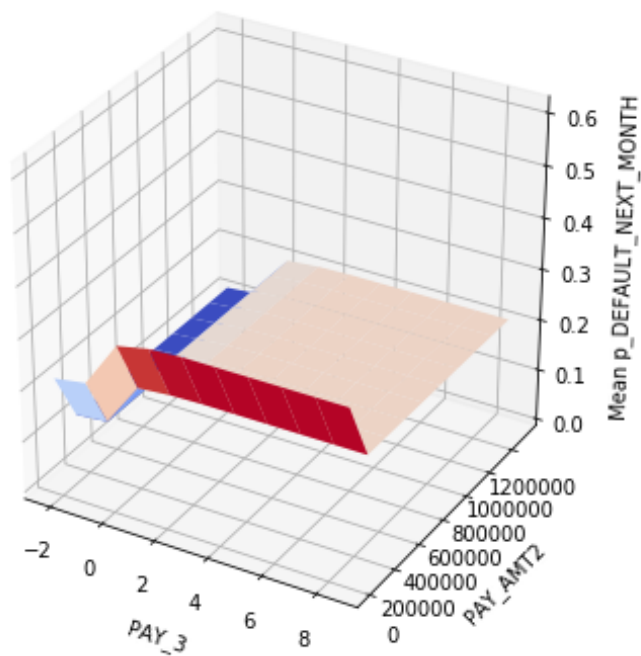
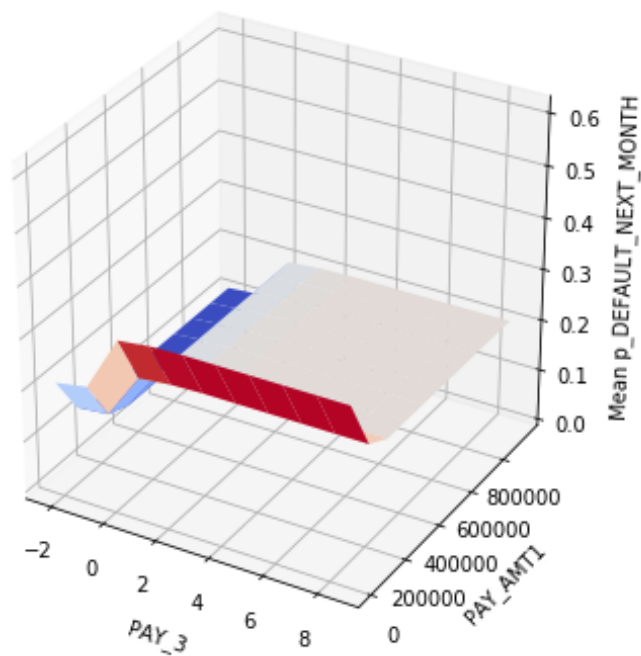
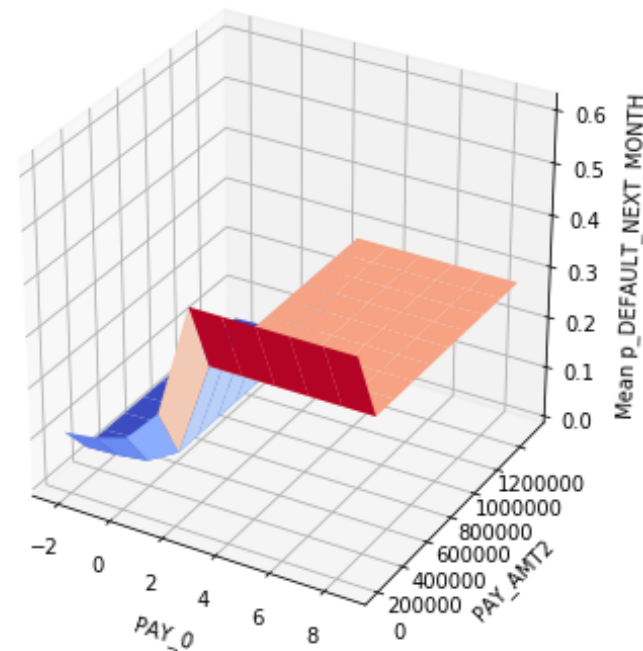
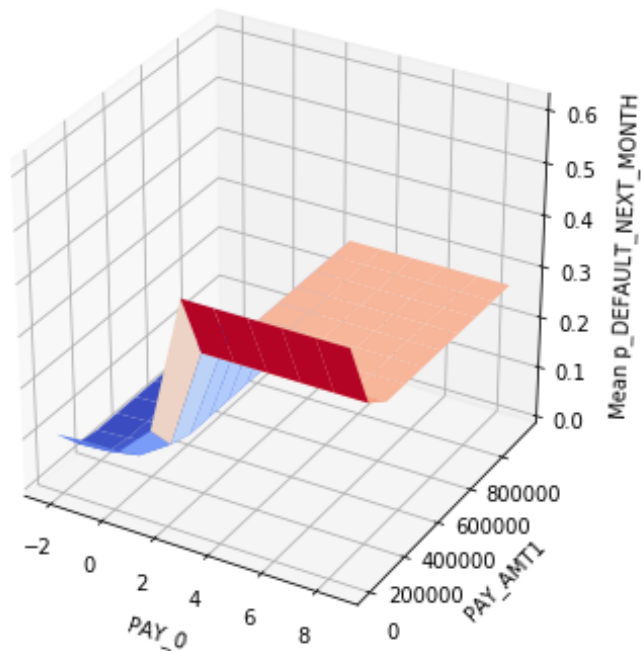
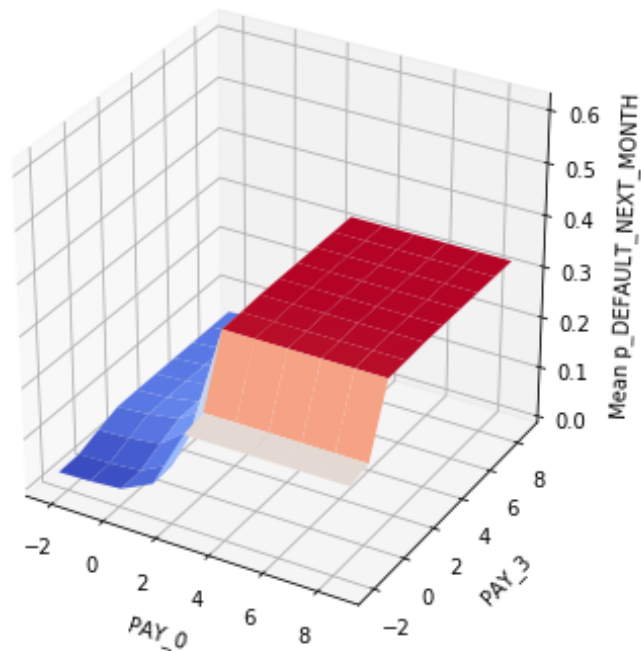
Error Metrics for PAY\_0

	Prevalence	Accuracy	True Positive Rate	Precision	Specificity	Negative Predicted Value	False Positive Rate	False Discovery Rate	False Negative Rate	False Omissions Rate
PAY_0										
-2	0.049	0.857	0.3	0.119	0.885	0.961	0.115	0.881	0.7	0.039
-1	0.117	0.805	0.383	0.267	0.861	0.913	0.139	0.733	0.617	0.087
0	0.05	0.864	0.345	0.143	0.891	0.963	0.109	0.857	0.655	0.037
1	0.822	0.457	0.368	0.93	0.871	0.229	0.129	0.07	0.632	0.771
2	1	0.709	0.709	1	0.5	0	0.5	0	0.291	1
3	1	0.748	0.748	1	0.5	0	0.5	0	0.252	1
4	1	0.571	0.571	1	0.5	0	0.5	0	0.429	1
5	1	0.444	0.444	1	0.5	0	0.5	0	0.556	1
6	1	0.25	0.25	1	0.5	0	0.5	0	0.75	1
7	1	0.5	0.5	1	0.5	0	0.5	0	0.5	1
8	1	0.75	0.75	1	0.5	0	0.5	0	0.25	1





Mean p\_DEFAULT\_NEXT\_MONTH for ['PAY\_0', 'PAY\_3', 'PAY\_AMT1', 'PAY\_AMT2']



# 강사 소개

정 준 수 / Ph.D ( heinem@naver.com )

- 前) 삼성전자 연구원
- 前) 삼성의료원 (삼성생명과학연구소)
- 前) 삼성SDS (정보기술연구소)
- 現) (사)한국인공지능협회, AI, 머신러닝 강의
- 現) 한국소프트웨어산업협회, AI, 머신러닝 강의
- 現) 서울디지털재단, AI 자문위원
- 現) 한성대학교 교수(겸)
- 전문분야: 시각 모델링, 머신러닝(ML), RPA
- <https://github.com/JSJeong-me/>

