

제 2회 KAIST-POSTECH-UNIST

# 데이터 사이언스 경진대회 설명회

---

공동주최



주관



후원



# Contents

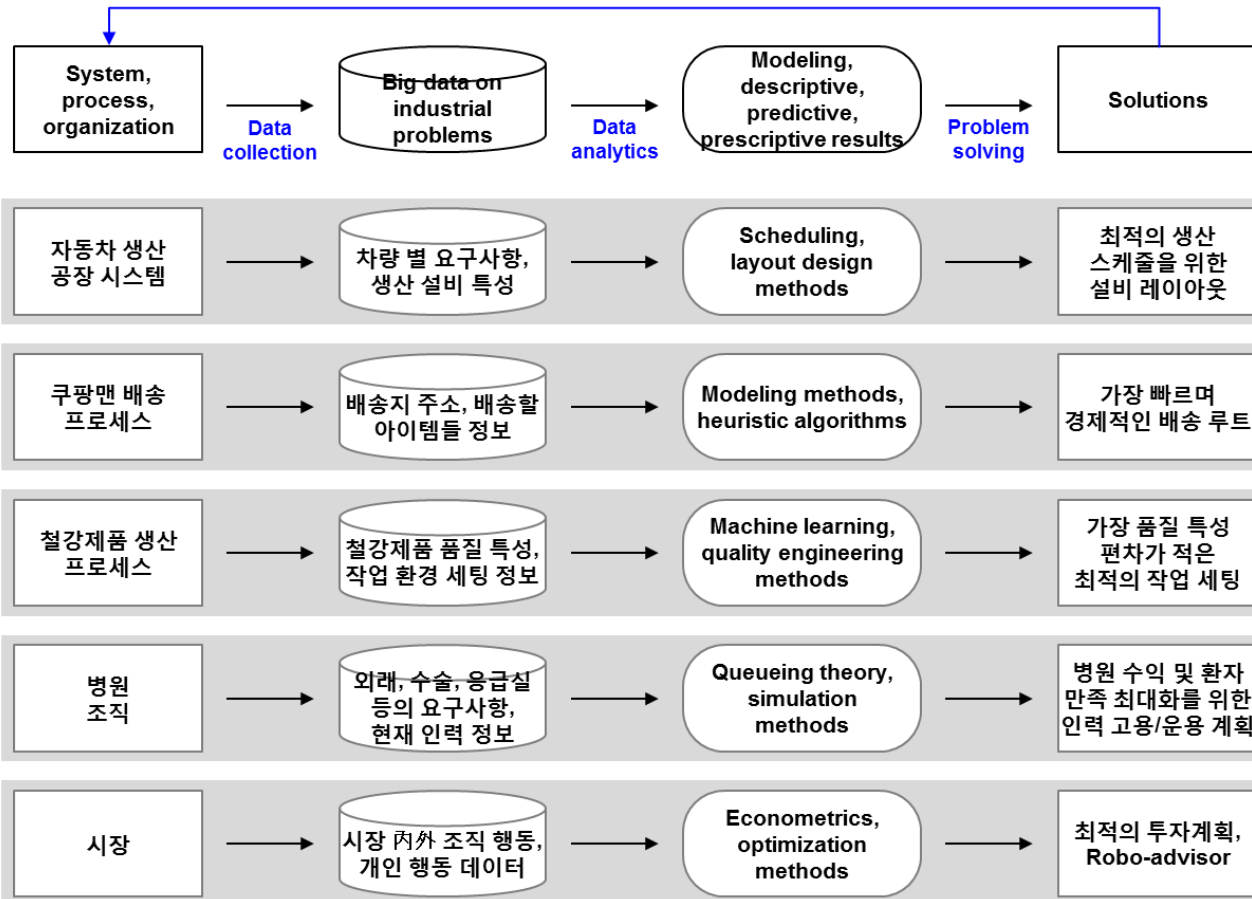
---

- 0. 경진대회 설명
- 1. 문제 배경 및 데이터 설명
- 2. Task 설명

# 산업공학 : 산업 시스템의 설계와 운영

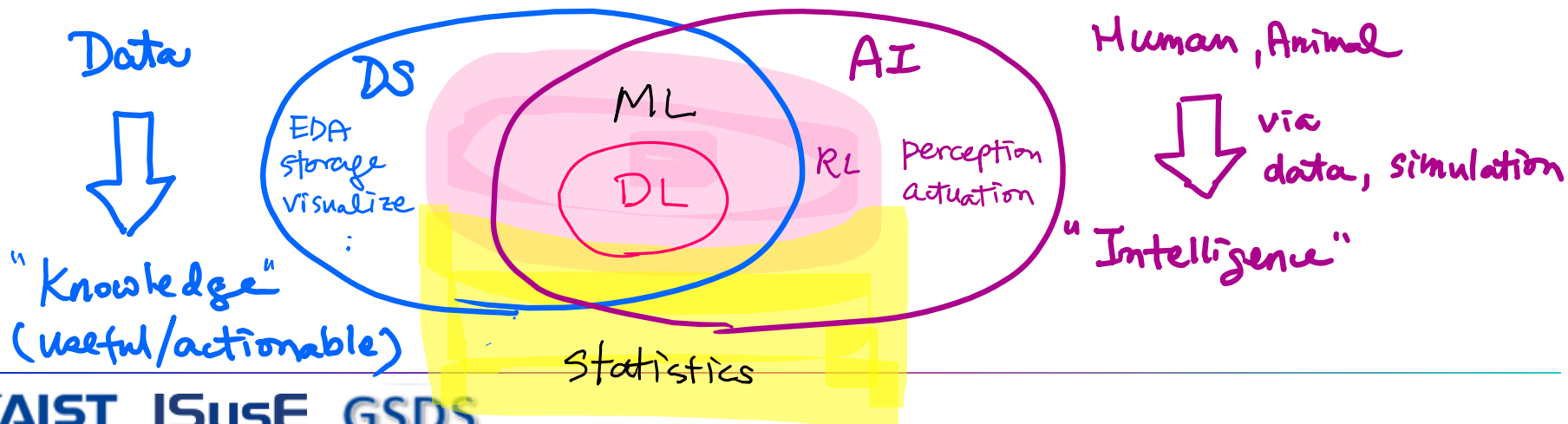
Domain knowledge → 연역적 추론 기반 모델 → 최적  
 Big data → 귀납적 추론 기반 모델 → 의사결정

Optimization and optimal decision making



# Data Science?

- **Data science** is an interdisciplinary field (\* mainly from Wikipedia with some modifications)
  - that uses scientific methods, processes, algorithms, and systems
    - math, statistics, programming, analytics, AI, ML, ...
  - **to extract knowledge and insights** from noisy, structured and unstructured data
    - inductive knowledge (vs deductive knowledge)
  - to apply knowledge from data across a broad range of application domains
- **Data science is related to**
  - **data mining**: old name for DS, narrower than DS
  - **statistics**: key foundation of DS
  - **machine learning**: core technology for DS
  - **AI**: overlaps with DS, but a little different objective



# 데이터사이언스 경진대회

## ■ 본 경진대회의 목적

- 기존에 많은 경진대회 : '예측'에 중점
- 산업공학 관점에서 본 데이터사이언스 = **예측 + 의사결정**
  - '예측' 결과를 이용한 '의사결정' → '문제 해결'
  - '산업/사회 시스템'의 설계와 운영 문제해결 → 산업공학

**산업공학** : 산업과 사회 시스템의 문제해결을 위한  
**최적 의사결정 + 고성능 예측모델 = 데이터 사이언스**



# 제2회 경진대회



## 제2회 KAIST POSTECH UNIST 데이터사이언스 경진대회

10.24 - 12.31

**산업공학 : 산업과 사회 시스템의 문제해결을 위한  
최적 의사결정 + 고성능 예측모델 = 데이터 사이언스**

(예비) 산업공학도를 모두 함께,  
데이터 사이언스로 산업 / 사회의 문제를 해결해보자!

\* 문제에 대한 상세내용은 대회 홈페이지 참조

### 대회일정

접수  
**10.24 - 11.9.**

대회일정  
**10.24 - 12.31.**

- 문제 공개 및 접수 시작 : 10.24(월)
- 예선작 중간 제출마감 : 11.26(토)
- 예선작 중간 평가 : 12.1(목)
- 예선작 최종 제출 마감 : 12.17(토)
- 예선작 최종 평가 및 결과발표 : 12.22(목)
- 본선 진출팀 최종 평가발표 : 12.29(목)
- 수상작 시상 : 23년 1월초

\* 자세한 내용은 대회 홈페이지 참조

### 온라인 설명회

ZOOM 회의 ID : 902 953 4311

- 1차 설명회 10.26(수) 19:00
- 2차 설명회 11.2(수) 19:00

\* 설명회 참여자 전원에게 기프티콘 제공

### 참가대상

KAIST, POSTECH, UNIST의  
학부 재학생들로 구성된 팀으로 참가

- 팀 구성조건 :  
4인이하, 무학과(내과) 학생 1명 포함 필수,  
산공과 재학생 1명 포함 필수, 타전공 학생 포함 가능,  
서로 다른 학교 학생들로 구성가능
- 팀 구성에 어려울 시 개인자격으로 신청 가능  
(주최처에서 팀 배정)

### 참가방법

참가신청, 대회 플랫폼(Slack) 참여 등은  
<http://datascience-contest.com>  
(대회홈페이지) 에서 확인

### 대회 상금 및 특징

- 대상 (하나은행상) : 300만원
- 금상 3팀 x 200만원
- 금상 (루닛상), 금상 (마키나락스상), 금상 (패스트캠퍼스상)
- 은상 (ECMiner상) : 100만원
- 동상 (ECMiner상) 2팀 x 50만원
- 장려상 : 상품제공 (5팀 내외)

\*내년 수상자에게는 후원사의 도구 사용 안내와 그에 따른 특별 지원, 인턴십 선발 우선권이 있을 수 있음. 이는 후원사 상황에 따라 다르니, 상세한 내용은 대회 홈페이지 참고

### 문의처

KAIST 산업및시스템공학과 TEL.042)350-3103  
KAIST 신하용 교수 myshin@kaist.ac.kr  
POSTECH 고영명 교수 ymngko@postech.ac.kr  
UNIST 이종재 교수 yongjaelee@unist.ac.kr

# 제2회 대회 일정

대회 일정은 상황에 따라 일부 변동 될 수 있습니다.  
대회 홈페이지 및 슬랙 통해 지속적인 확인 부탁드립니다.

- 1차 온라인 설명회: 10/26, 19:00
- 2차 온라인 설명회: 11/2, 19:00
- 접수: 10/24 ~ 11/9
  - 대회 홈페이지([datascience-contest.com](https://datascience-contest.com))에서 신청
  - 참가신청 완료된 팀은 slack에 초대 후 데이터 제공 예정
- 예선 : 정량 평가
  - 예선작 중간 제출: 마감 11/26, 평가발표 12/1
  - 예선작 최종 제출: 마감 12/17, 예선 결과 발표 12/22
- 본선 : 발표 평가 (10개팀 내외)
  - 예선 정량 평가를 통해 10개 팀 내외 본선 진출팀 결정
  - 본선 진출 팀 최종 발표평가: 12/29(목)
- 수상작 시상: 2023년 1월 초



# 참가자격 및 팀구성

- KAIST, POSTECH, UNIST의 학사과정 학생으로 구성된 팀
  - 4인 이하
  - 최소 1인 이상의 새내기(무학과) 학생 포함
  - 최소 1인 이상의 산공과 학부생 포함
    - 복수 전공, 휴학생 포함 가능
  - 타전공 학생 포함 가능
  - 서로 다른 학교 학생들로 구성 가능
- 팀 구성이 어려울 시 개인 자격 (또는 부분 팀)으로 신청 가능
  - 신청 후 주최측에서 팀 매칭
- 그외 팀 구성에 있어서 특이 상황은 주최측에 문의 요망

# 대회 시상 및 상금

- 대상 (1 팀): 300만원
  - 하나은행 상
- 금상 (3 팀): 각 200만원
  - 루닛 상
  - 마키나락스 상
  - 패스트캠퍼스 상
- 은상 (1 팀): 100만원
  - ECMiner 상
- 동상 (2 팀): 각 50만원
  - ECMiner 상
- 장려상 (5팀 이내)
  - 상품

후원



Lunit

MakinaRocks



Fast campus



ECMiner

## \* 상금 이외의 특전

- (희망자에 대하여) 후원사 Tool 사용 교육
- 마키나락스 특전
  - Link™ 사용 본선진출 팀에 대해 Team building 지원 (20만원/팀)
  - 수상팀 인턴지원시 서류 평가 면제
  - 대회 참가자를 위한 Link™ 사용법 설명 Session
    - \* 11/4(금) 11am에 아래 URL로 접속
    - \* <https://meet.google.com/hpj-fxsy-ywv>
- 그 외에도 특전이 있을 경우 및 특전의 상세한 내용은 홈페이지를 통해 공지 예정

# 문의처

- 홈페이지
  - <https://datascience-contest.com>
  - 대회 자료
  - 참가신청
- 슬랙채널
  - 참가신청자에게 공개
- 담당 교수
  - KAIST 신하용 교수 ([hyshin@kaist.ac.kr](mailto:hyshin@kaist.ac.kr)),  
박찬영 교수 ([cy.park@kaist.ac.kr](mailto:cy.park@kaist.ac.kr)), 민승기 교수 ([skmin@kaist.ac.kr](mailto:skmin@kaist.ac.kr))
  - POSTECH 고영명 교수 ([youngko@postech.ac.kr](mailto:youngko@postech.ac.kr))
  - UNIST 이용재 교수 ([yongjaelee@units.ac.kr](mailto:yongjaelee@units.ac.kr))

## Part 1

# Data Description

# Dataset: Overview



## ▪ Provided by Hana Bank

- Usage record of Hana 1Q App (mobile banking app)
- for 8 months (2022.01.01 ~ 2022.08.26, 238 days)

## ▪ Data description

- Input features (raw) :  $3 + 238 \times 3 = 717$  columns per user

feature name	meaning	remark
gender	Male(0) or female(1)	binary
age_code	age band (1:<20, 2~13:20~79, 14:>=80)	
region_code	region code (1~18)	categorical
c20220101~ c20220826	Number of logins for each day	time series (cardinal)
t20220101~ t20220826	Number of logins with money transfer	time series (cardinal)
s20220101~ s20220826	Duration of staying with the app	time series (cardinal)

- Label = {1: small business owner, 0: general (non-business)}
  - The proportion of the small business owner (label=1) is approximately 6.55%

# Dataset: Example

## ■ Example

– Input features

Gender	age_code	region_code	c20220101	...	s20220101	...	t20220101	...
2	4	4	4		121		2	
1	8	1	1		15		-	

- YYYYMMDD : 20220101 ~ 20220826
- cYYYYMMDD : # logins on YYYYMMDD
- sYYYYMMDD : Duration of stay on YYYYMMDD
- tYYYYMMDD : # logins with money transfer on YYYYMMDD

– Label

Small business owner
1
0

- 1: A small business owner
- 0: Not a small business owner

# Dataset: Files

- Dataset for modeling ( $D_{\text{model}}$ ): input features + label
  - 800,000 user records
  - can be used for developing your model
    - X\_model.csv: Input features
    - Y\_model: Label
- Dataset for competition ( $D_{\text{exam}}$ ) : input features only
  - 200,000 user records
    - X\_exam.csv: Input features
- File to submit : prediction and decision results for  $D_{\text{exam}}$ 
  - submission.csv: prediction and decision results
    - a sample file will be given
    - $\{p_i, a_i, b_i\}$  in the same order as X\_exam.csv

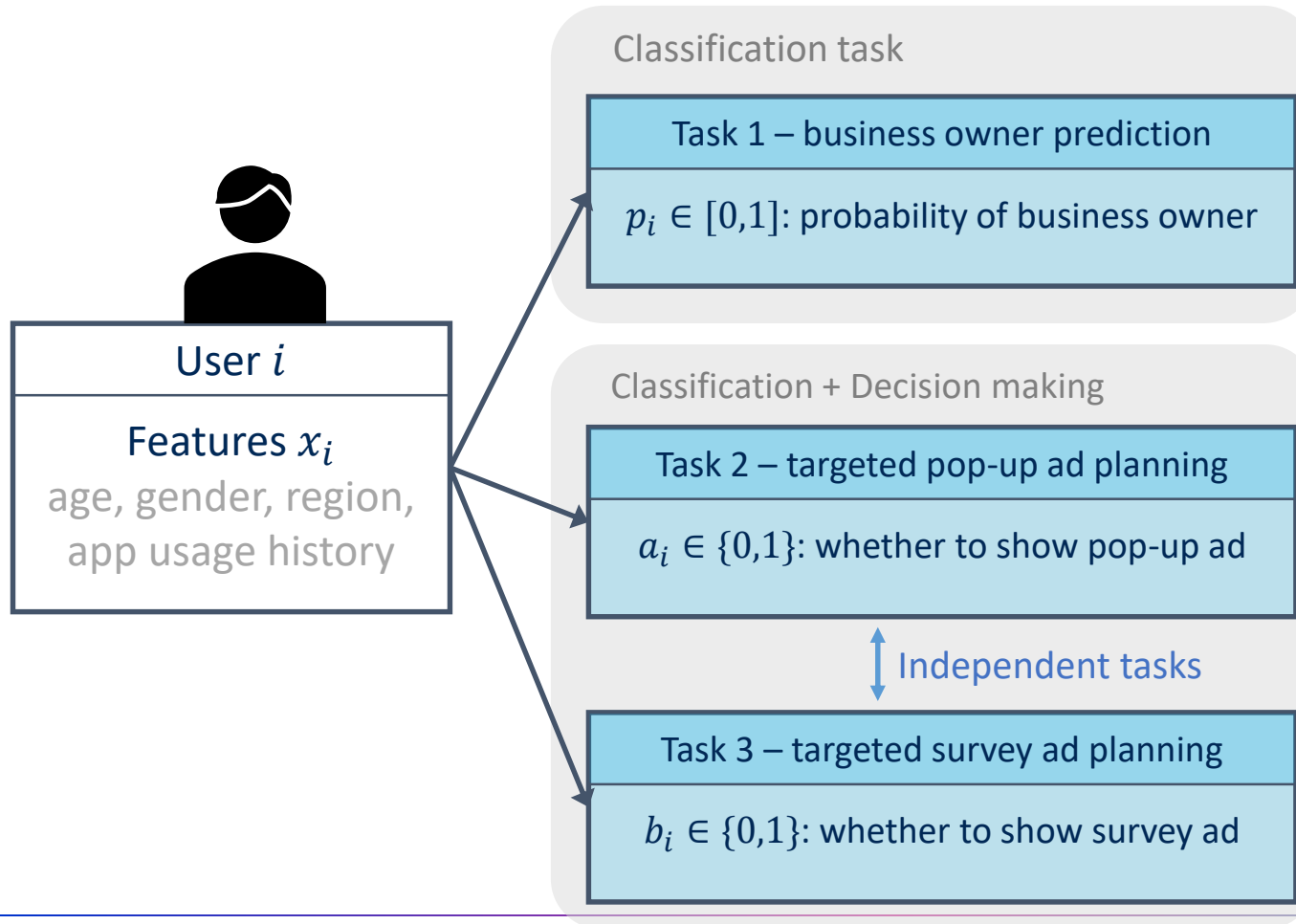
## Part 2

# Task Description



# Overview

- For each user  $i$  in the exam dataset (200,000 users), you have to submit three numbers  $p_i, a_i, b_i$

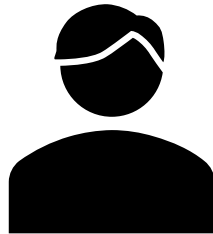


submission.csv

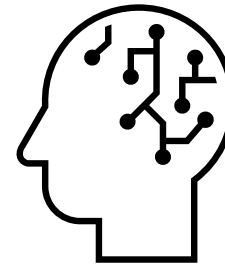
Task 1 $p_i$	Task 2 $a_i$	Task 3 $b_i$
business_prob	popup	survey
0.61845	1	0
0.97296	1	1
0.78073	1	0
0.02992	0	0
0.25522	0	0
0.92836	1	1
0.89172	1	0
0.63428	1	0
0.78102	1	0
0.93564	1	1
0.2145	0	0
0.97235	1	1
0.64393	1	0
0.36104	0	0
0.51847	1	0
0.53315	1	0
0.43127	0	0
0.28834	0	0

# Task 1. Small Business Owner Prediction

- **Context:** the company wants to predict which users are small business owners based on user profiles & app usage history



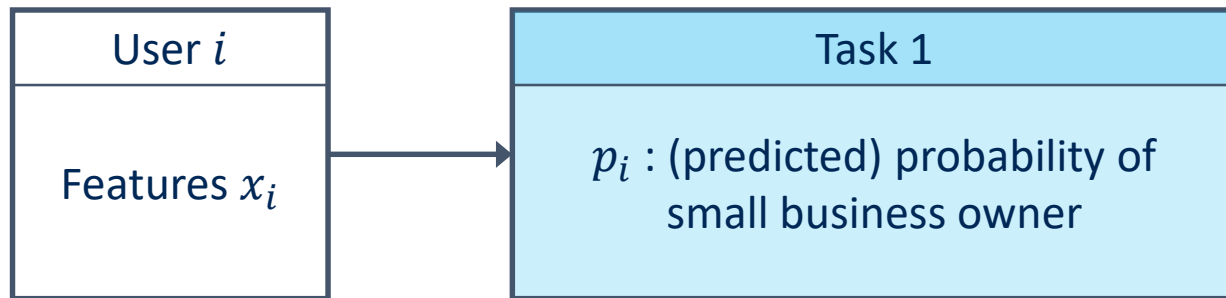
- Gender: male
- Age: in 20s
- On Jan 1<sup>st</sup>,
  - logged in 2 times
  - transferred money once
  - used app for 10 mins
- On Jan 2<sup>nd</sup>, no activity
- ...



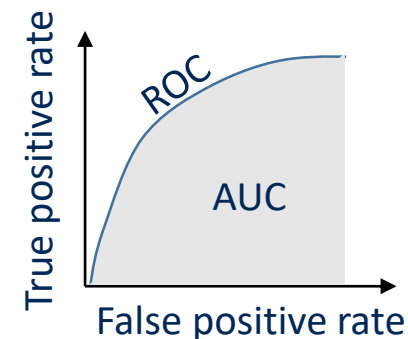
This user will be  
a small business owner  
with 60% probability

# Task 1. Small Business Owner Prediction

- **Task:** for each user in the exam dataset, predict the probability that the user is a small business owner
  - **Prediction:**  $p_i \in [0,1]$

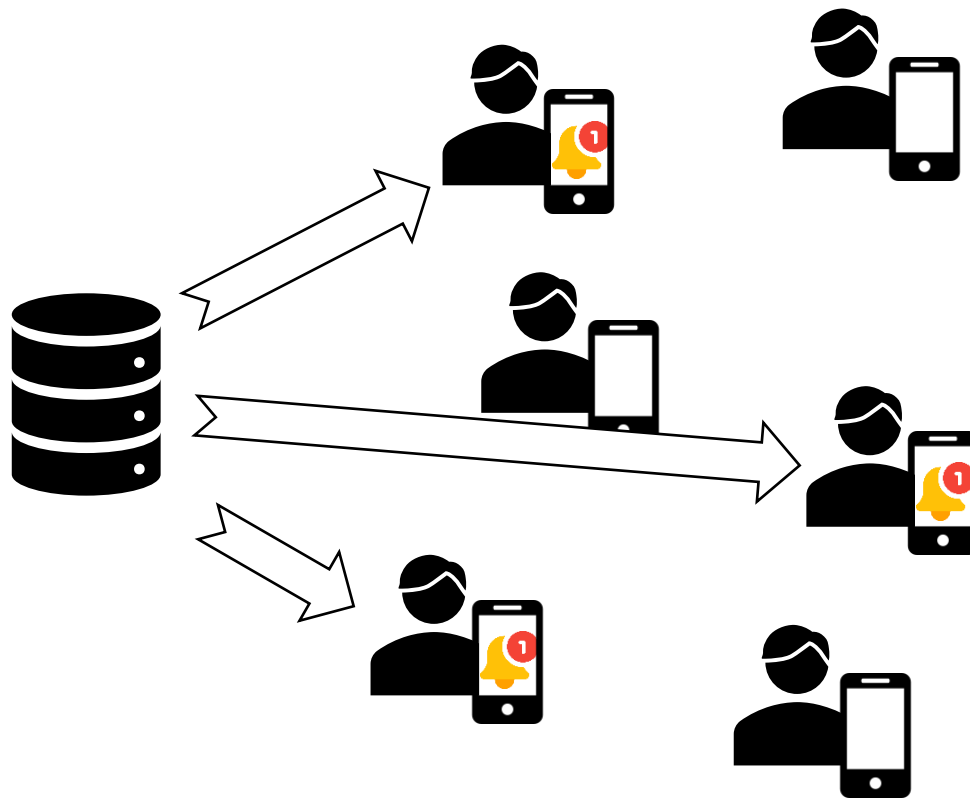


- **Evaluation:** AUC value of ROC curve



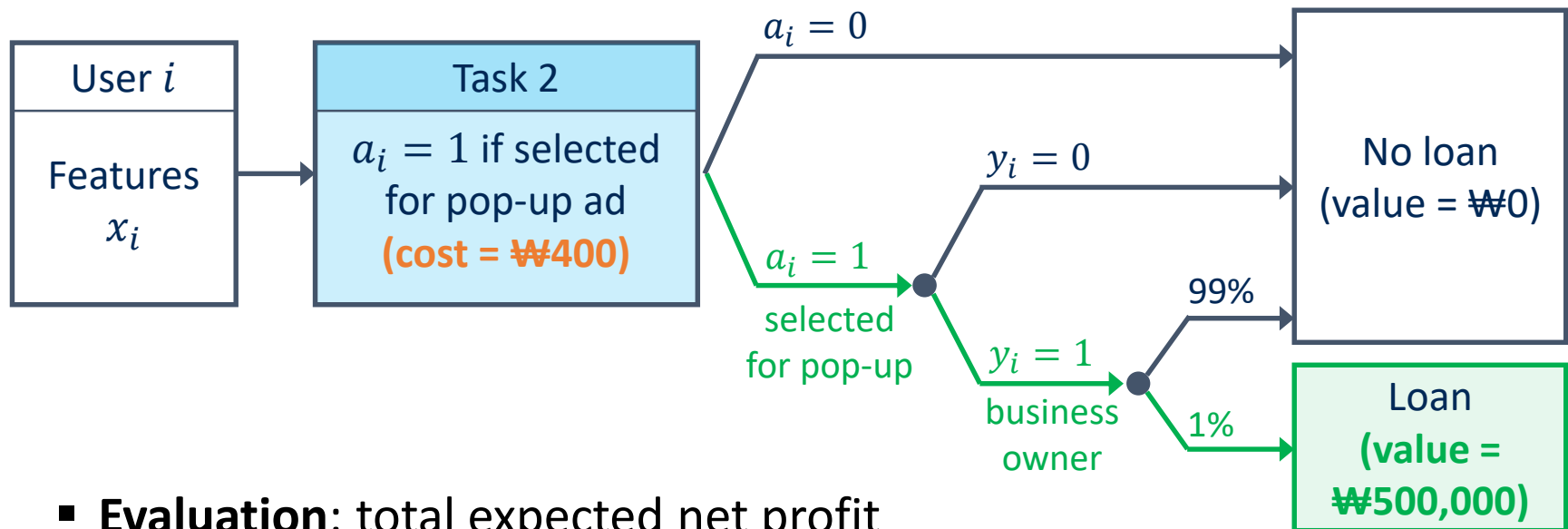
## Task 2. Pop-up Ad Planning

- **Context:** the company is planning to advertise the loan service by sending *pop-up notifications* to selected users



## Task 2. Pop-up Ad Planning

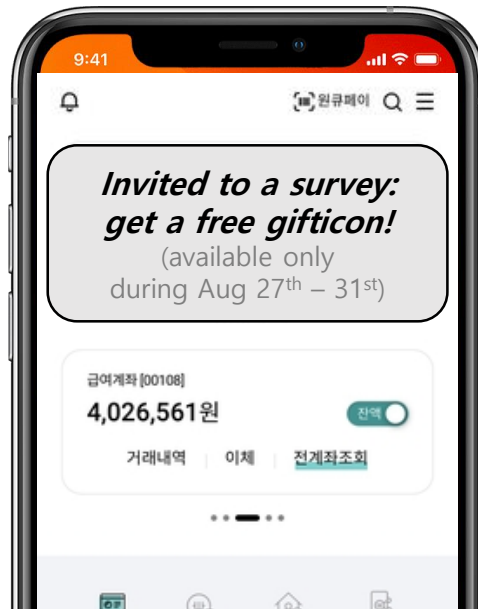
- **Task:** for each user in the exam dataset, decide whether to send a pop-up notification
  - **Decision:**  $a_i = 0$  or  $1$
  - **Cost:** ₩400 per user (pop-up add setup cost)
  - **Value:** 1% of small business owners who viewed pop-up ad will use the loan service, yielding value of ₩500,000



- **Evaluation:** total expected net profit  
= total expected value – total cost

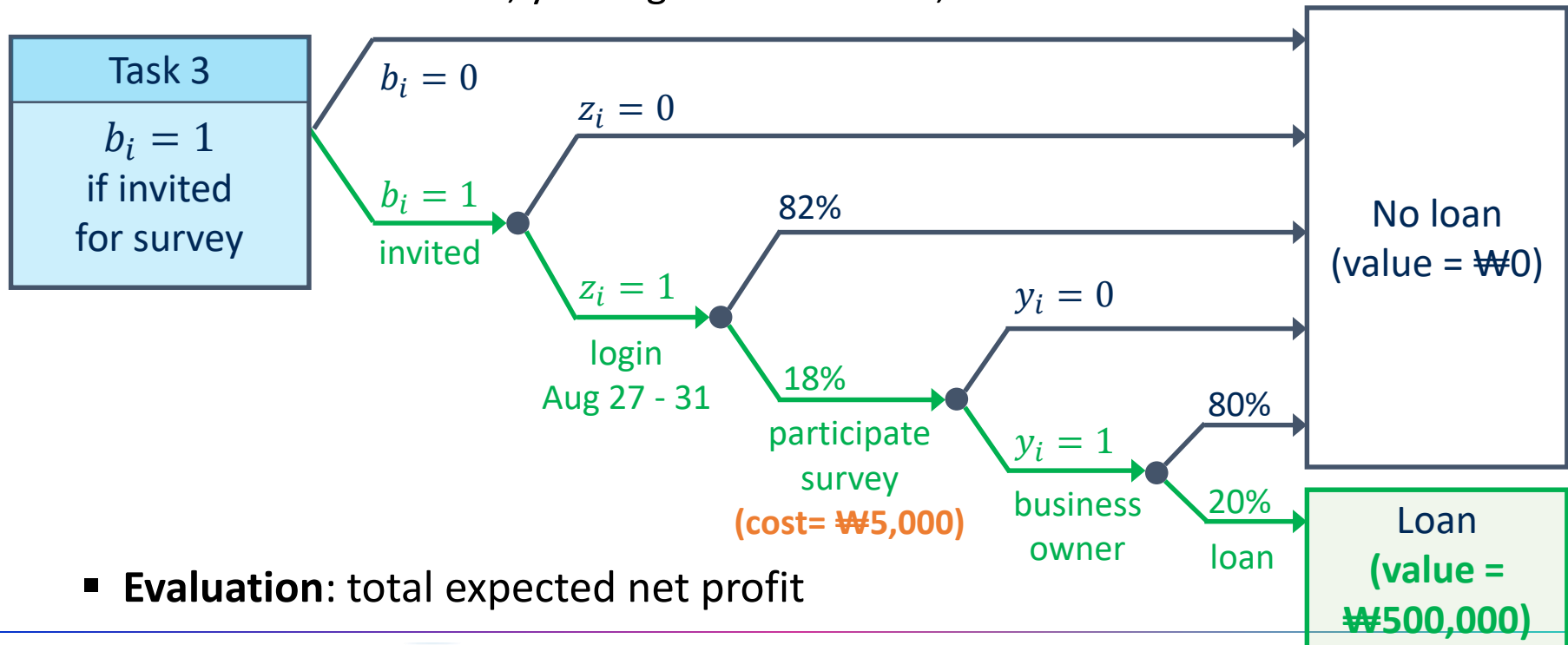
## Task 3. Survey Ad Planning (indep. of Task 2)

- **Context:** the company is planning to advertise the loan service through *in-app survey-based promotion*
  - over the next 5 days (August 27<sup>th</sup> – 31<sup>st</sup>)
  - only for selected users (up to 50,000 users)
  - gifticon for survey participants



## Task 3. Survey Ad Planning (indep. of Task 2)

- **Task:** for each user in the exam dataset, decide whether to invite him/her to the survey promotion
  - **Decision:**  $b_i = 0$  or  $1$  (up to 50,000 users: i.e.  $\sum_i b_i \leq 50000$ )
  - **Cost:** upon login, an invited user will participate the survey with 18% chance, incurring ~~₩~~5,000 cost (gifticon)
  - **Value:** upon participation, a small business owner will use the loan service with 20% chance, yielding value of ~~₩~~500,000



- **Evaluation:** total expected net profit

## Conclusion

- 대회 진행 관련 세부 사항은 아래 참조
  - 대회 홈페이지 (<https://datascience-contest.com>)
    - 참가신청 페이지 (대회 홈페이지 통해 접속 가능)  
[https://docs.google.com/forms/d/e/1FAIpQLSdwOdkCwmWmsoTQ5Q0x7r\\_XVobrRS3j4kSGWVE8v8g4SlraFA/viewform](https://docs.google.com/forms/d/e/1FAIpQLSdwOdkCwmWmsoTQ5Q0x7r_XVobrRS3j4kSGWVE8v8g4SlraFA/viewform)
  - 오픈카톡방 (대회 참여 전, <https://open.kakao.com/o/gaALZVGe>)
  - 슬랙 (대회 참여 후)
- 질문사항은 오픈카톡방과 슬랙 통해서 편하게 물어봐주세요

# Thank you for your attention!