



# *Speed dating*

문 기 태   최 윤 슬   김 나 형   안 수 빈

1

**데이터 소개  
&  
분석목적**

2

**데이터 탐색  
&  
정제**

3

**모형 구축  
&  
설명**

4

**분석결과  
&  
결론**



***CONTENTS***

# **[ 데이터 소개 ]**

---

## Speed Dating?

독신 여자,남자가 **애인을 찾을 수 있도록**  
짧은 시간동안 여러 사람들을 돌아가며  
만나보는 일종의 소개팅



# **[ 분석 목적 ]**



매칭 성공에 영향을 주는 요인들을 찾기

종속변수 **match**를 맞춰라 !!

# **[ 데이터 탐색 및 정제 ]**

---

196 variables

6702 obs

iid	id	gender	idg	condtn	wave	round	position	positin1	order	partner	pid	match	int_corr	samerace
70	5	1	10	1	3	10	6	NA	5	5	60	0	0.45	0
351	11	0	21	2	14	18	11	11	7	10	370	0	0.15	0
345	5	0	9	2	14	18	17	17	13	8	368	1	0.32	0
351	11	0	21	2	14	18	11	11	3	17	377	0	-0.40	1
488	8	1	16	2	19	15	5	14	8	4	469	0	0.31	1
359	19	0	37	2	14	18	15	15	1	14	374	0	-0.07	0
8	8	0	15	1	1	10	10	NA	2	10	20	1	-0.21	1
145	4	0	7	2	7	16	14	14	12	3	160	0	-0.18	0
370	10	1	20	2	14	20	12	12	8	12	352	0	0.16	0
286	14	1	28	2	11	21	16	16	3	20	271	0	-0.22	1
382	4	0	7	2	15	19	17	17	9	3	399	0	0.24	1
301	8	0	15	2	12	14	3	3	4	10	317	0	0.67	1
171	14	1	28	2	7	16	3	1	3	6	147	0	0.54	0
523	15	0	29	2	21	22	5	5	15	7	537	0	0.14	1
177	4	0	7	1	8	10	4	4	9	7	190	0	-0.03	1
474	9	0	17	2	19	15	8	9	1	7	487	0	-0.21	0
173	16	1	32	2	7	16	10	16	11	2	143	0	-0.03	1
162	5	1	10	2	7	16	7	14	10	15	156	0	0.45	0
106	13	1	26	2	4	18	3	NA	1	14	89	0	0.38	0



# 주요 변수 소개

변수 이름	변수 설명	범주 개수 / 범위
<i>round</i>	그 날 스피드 데이트에서 만난 파트너 수	5~22
<i>wave</i>	스피드 데이트 조 번호	1~22
<i>pf_o_attr</i>	6개 항목에 파트너가 사전에 매긴 점수	1~10
<i>from</i>	태어나서 자란 나라	270개
<i>career_c</i>	직업 코드	17개
<i>exphappy</i>	소개팅에 기대하는 정도	1~10점
<i>like / like_o</i>	이 사람이 마음에 든 정도 (_o는 상대방이 나에게 매긴 점수)	1~10점

# 주요 변수 소개

변수 이름	변수 설명	범주 개수 / 범위
<i>match_es</i>	예상 매칭 성공 수	5~22
<i>satis_2</i>	만난 사람에 대한 만족 정도	1~10
<i>length</i>	미팅의 전반적인 만족 정도	3
<i>numdat_2</i>	그 날 스피드 데이트 수의 적당한 정도	3
<i>int_corr</i>	사전조사가 파트너와 비슷한 정도	-1~1
<i>samerace</i>	인종이 같은지 여부	0/1
<i>prob / prob_o</i>	상대방이 나를 선택할 확률의 정도 (_o는 상대방이 나에게 매긴 점수)	1~10점

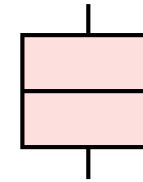
196  
variables

연속형

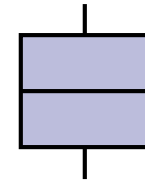
범주형

## 연속형

종속변수인 match와  
Boxplot을 그려보아  
유의미해 보이는  
변수들을 선택



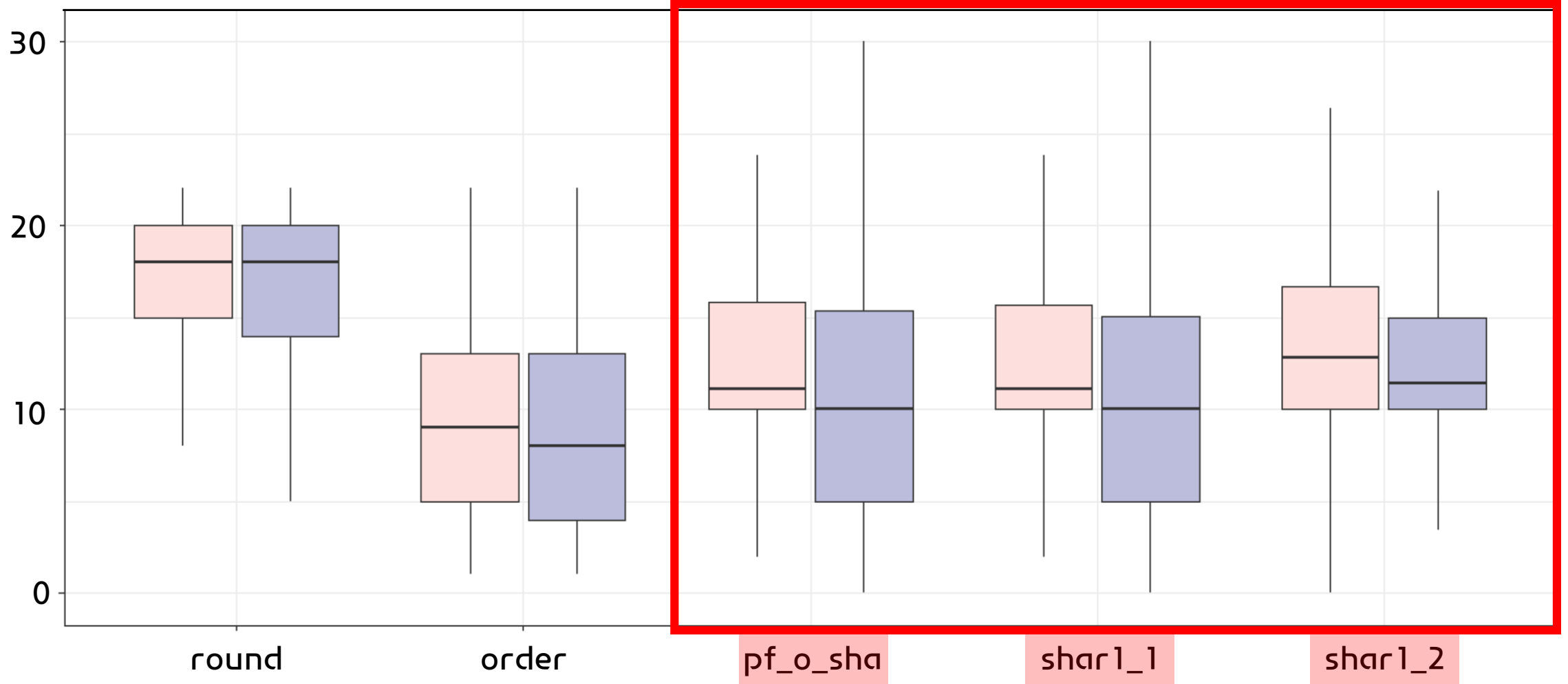
match = 0



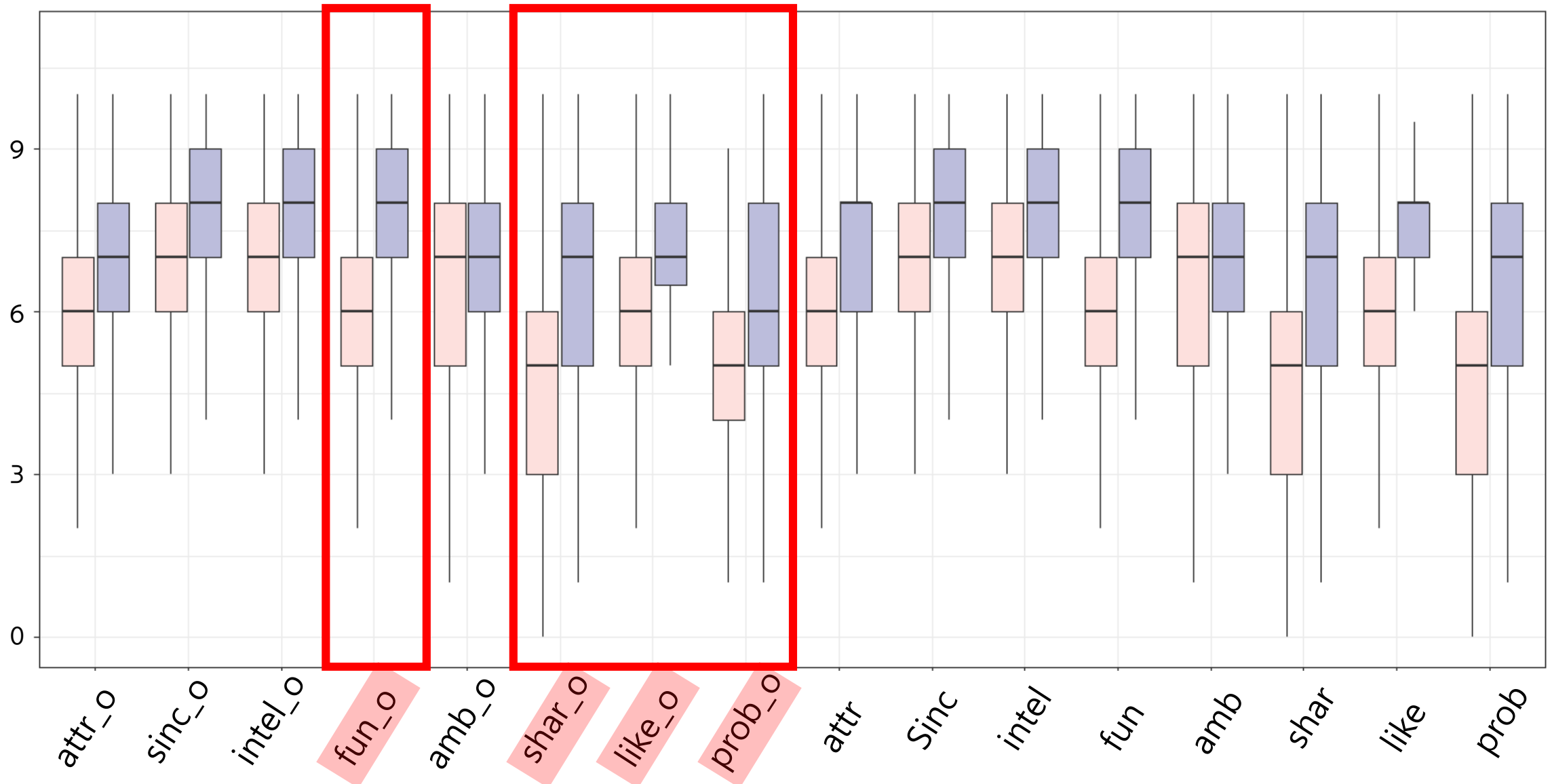
match = 1

각 변수마다 match별로  
차이 비교

# Boxplot 그리기



# Boxplot 그리기



## 중요도를 점수로 부여하는 변수들

ex) attr, attr1\_1...

✓ **NA**

→ **0으로 대체**

: 점수를 매기지 않은 것은 중요하지 않다는 것!

✓ 1부터 10까지인데  
10보다 큰 경우

→ **10으로 대체**

✓ 합산하여  
100보다 큰 경우

→ 합산한 점수로 나눠  
합이 100이 되도록  
다시 입력

- *pf\_o\_fun, pf\_o\_amb, pf\_o\_sha* 가 **NA**  
⇒ 나머지 항목들을 더해 100점이 되므로 0 할당
- *attr\_o ~ shar\_o* 가 **NA**  
⇒ 위와 같은 이유로 0 할당  
⇒ *attr\_o*가 10.5인 값, 10으로 할당
- *funl\_l, amb\_l, shar\_l* 가 **NA**  
⇒ 위와 같은 이유로 0 할당



- *like\_o, prob\_o* 가 **NA**

$$\Rightarrow attr\_o + sinc\_o + \dots + shar\_o = total\_o$$

*total\_o* 범위에 따라 *like\_o*를 채움 (10단위로 끊어서)

- *attr ~ shar*은 평가하지 않았는데 *like*값은 결측이 아닌 사람 - 4명

$\Rightarrow attr\_shar$  각각의 평균을 넣어줌

- 사전 조사 안 한 사람들 (*int\_corr*이 **NA** 값인 사람들)

$\Rightarrow$  사전 조사가 결과에 중요하다고 생각하여 삭제

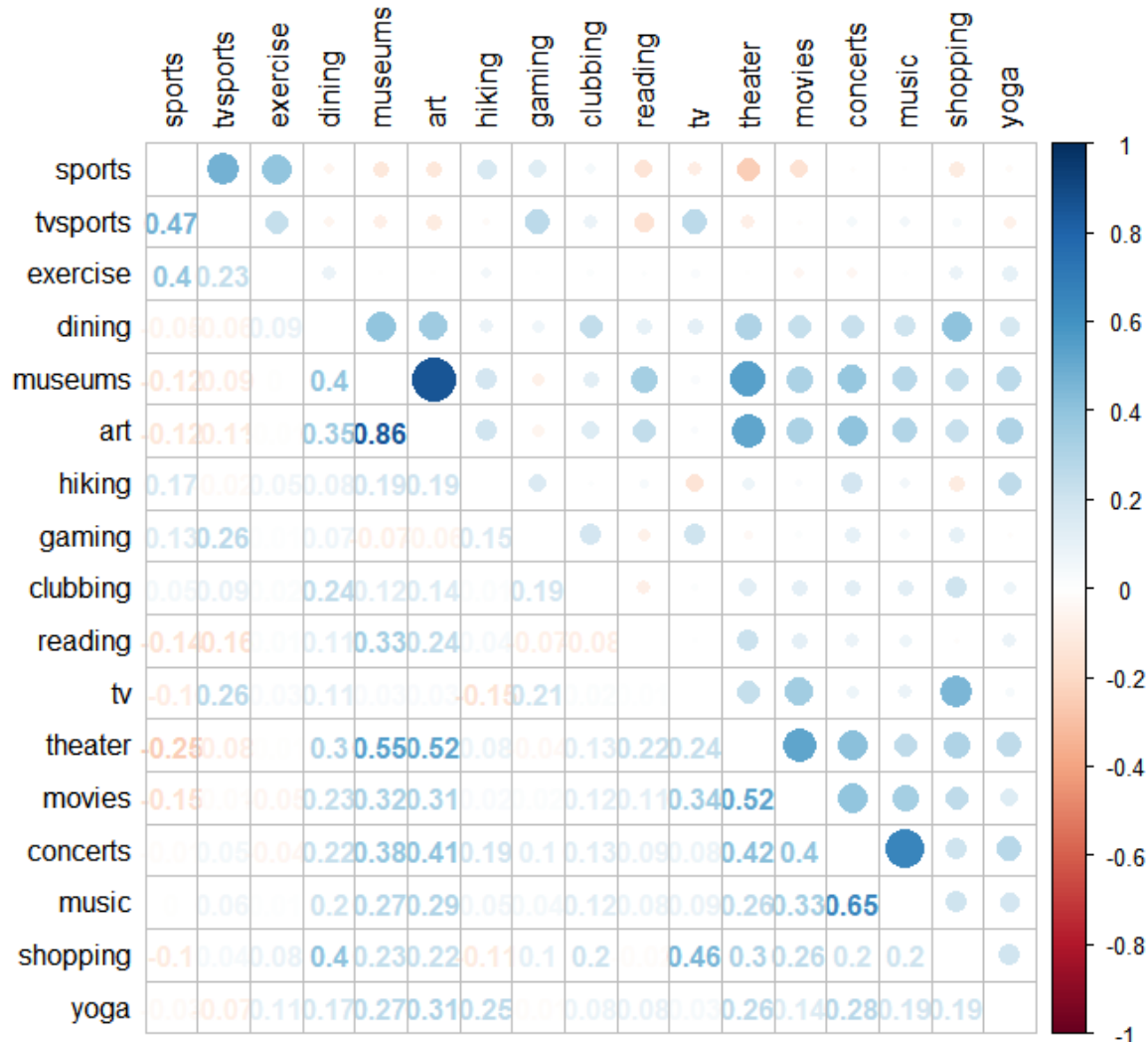
## 범주형

1

서로 상관계수가  
높게 나오거나 의미없는  
범주가 많은 변수통합

2

종속변수인 match와  
도수분포표를 그리고  
카이제곱 테스트를 하여  
유의미한 변수만 선택



## 1. 취미 변수

: 상관계수 높은 변수끼리 묶기

*museum art music reading  
theater movies concert*  
= **cult** : 문화생활

*dining shopping clubbing*  
= **outside** : 외출활동

•  
•  
•

## 2. *career\_c* 가 **NA** 인 사람들

⇒ *career*를 보고 범주를 확인해 넣어줌

<i>career</i>		<i>career_c</i>
Lawyer & Law	➡	1
Economist	➡	2
Tech Professional	➡	5
NA	➡	15

## 3. *dec* / *dec\_o* 변수 **제거**

⇒ 제거하지 않으면  
*match*가 그대로 예측이 됨

⇒ *dec* X *dec\_o* = *match*

## 4. 카이제곱 검정을 통한 변수 제거

예 ) *numdat\_2*

<i>match</i> <i>numdat_2</i>	<i>match</i> X	<i>match</i> 0
1	455	95
2	2728	509
3	2259	484



p-value = 0.1284 > 0.05

***numdat\_2*** 제거

- *match\_es* 가 **NA**

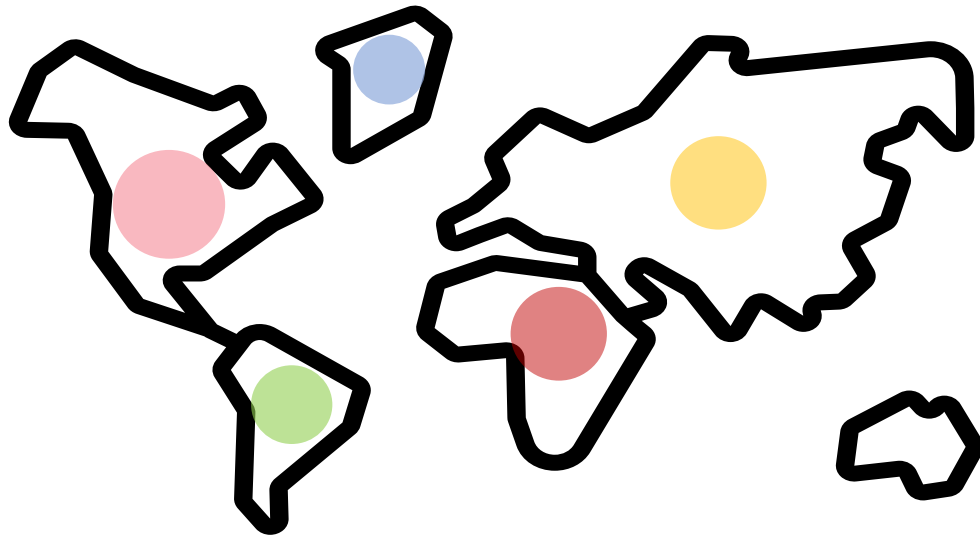
⇒ 만난 사람의 수 (*round*)와 밀접한 관련이 있을 것이라 생각,  
*table*을 그린 후 **최빈값**을 넣는다.

+ (1.5 , 2.5인 사람 -> 각각 2,3 / 3.4인 사람 -> 3)

- *match\_ess* = *match\_es* / *round*

⇒ 예상 매칭 성공 수를 뜻하는 *match\_es*는 *wave*별로 인원이 다르기 때문에 상이.  
인원을 뜻하는 *round*로 나눠서 척도를 맞추준다.

- **from** 출신 국가 및 지역 변수



270개의 범주  
→ 대륙 기준 6개의 범주  
(미국은 동부/서부)

- 1) 94115 -> *San Francisco*
- 2) 미국 -> 동부, 서부로 나눔 (*East/West*)
- 3) 유럽 (*Europe*)
- 4) 북아메리카 (*North America*)
- 5) 아시아 (*Asia*)
- 6) 남아메리카 (*South America*)





iid condtn round  
pid order pf\_o\_sha  
⋮

196개의 변수

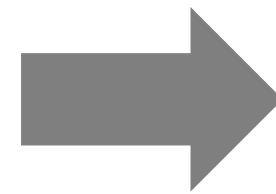
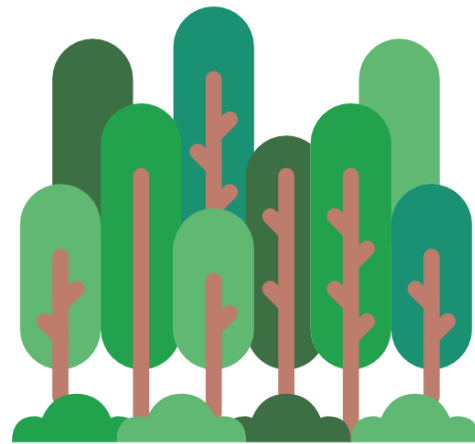
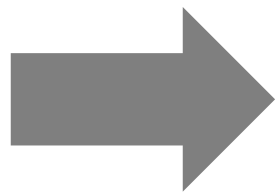


~~iid~~ condtn round  
~~pid~~ order pf\_o\_sha  
⋮

30개의 변수

# [ 모형 구축 ]

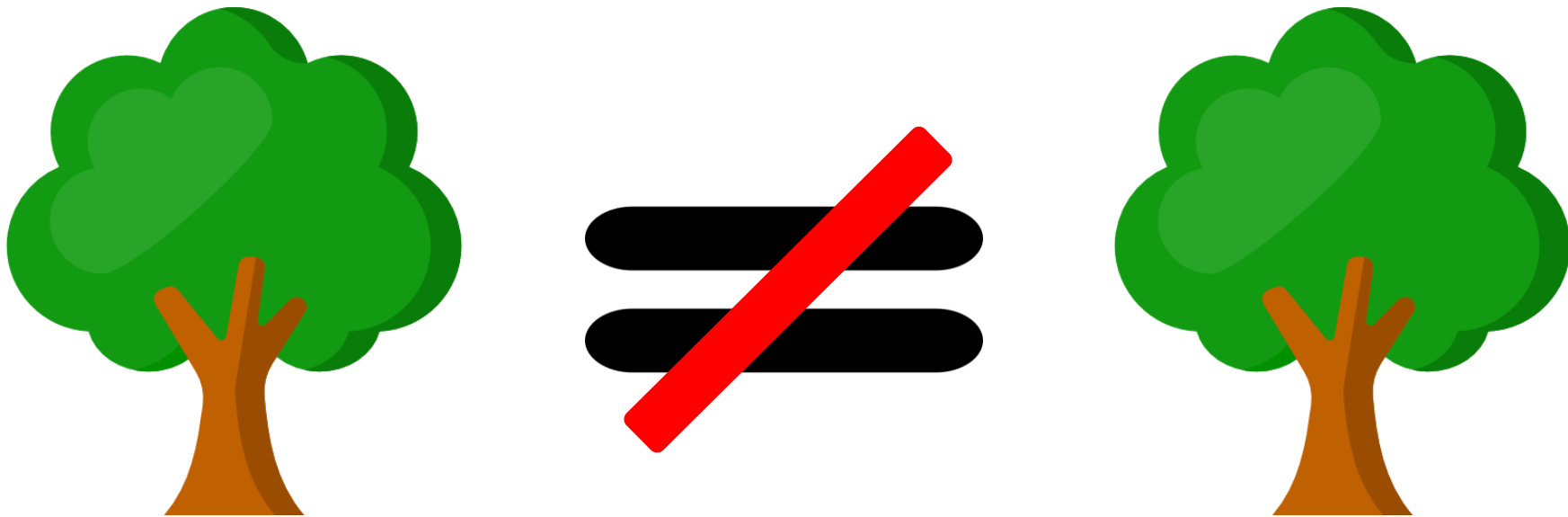
# 분석 1. 랜덤포레스트



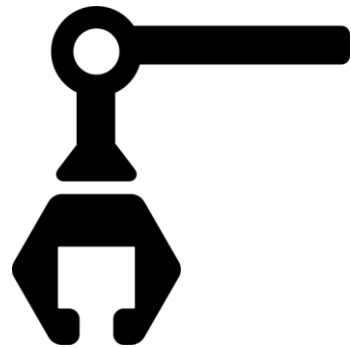
분산 감소  
편향 유지

각각의 결정

# 분석 1. 랜덤포레스트

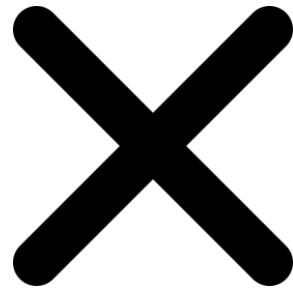


# 분석 1. 랜덤포레스트

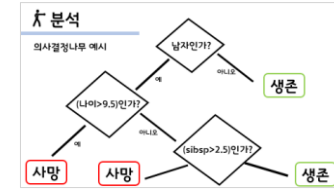


gender  
condtn  
round  
order  
pf\_o\_sha  
.....

Mtry= 사용할 변수

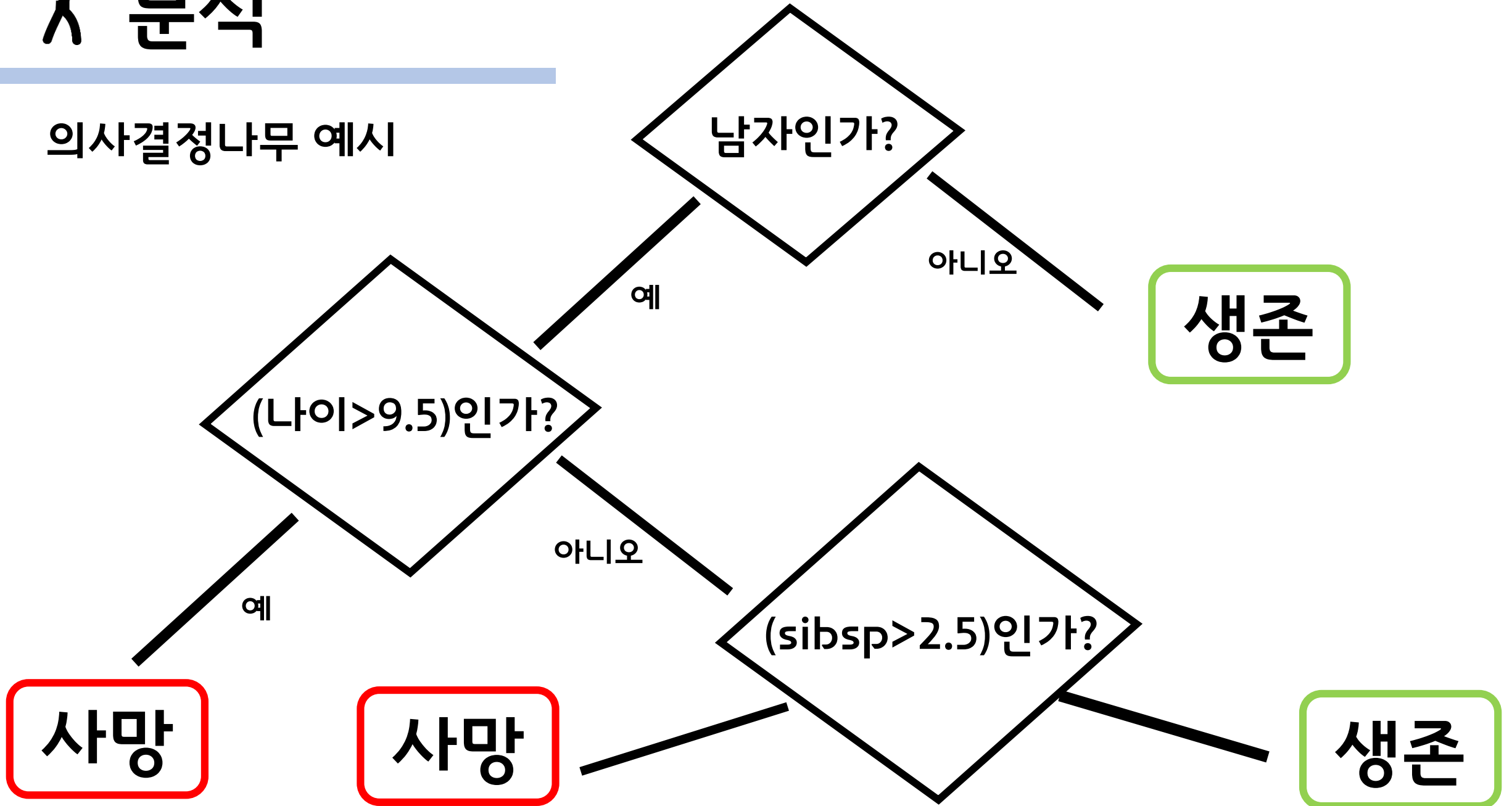


ntree

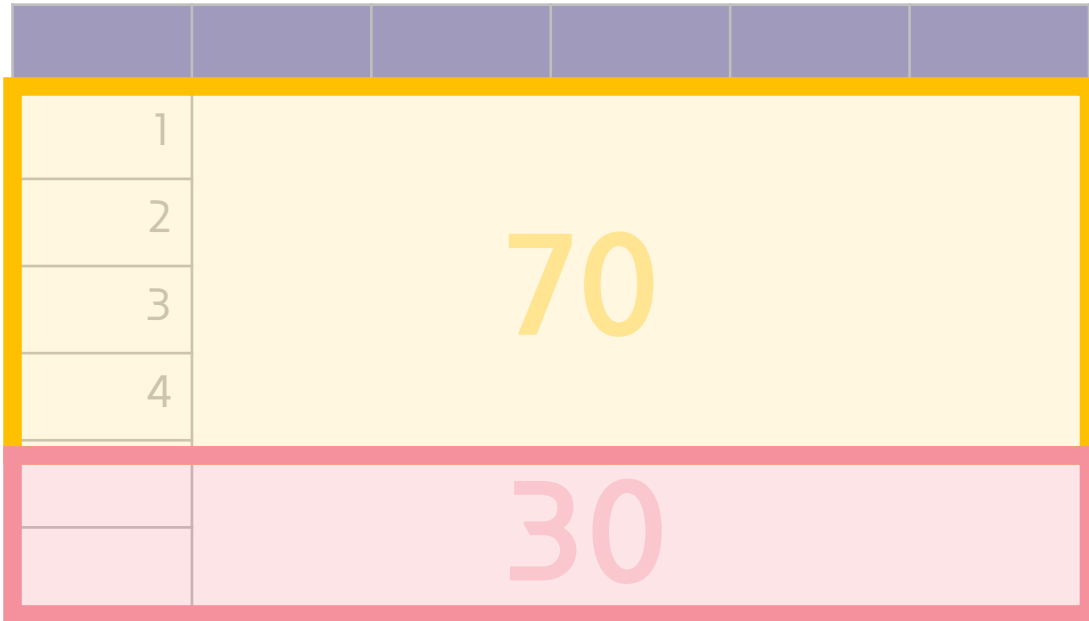


# 분석

의사결정나무 예시



# 분석 1. 랜덤포레스트

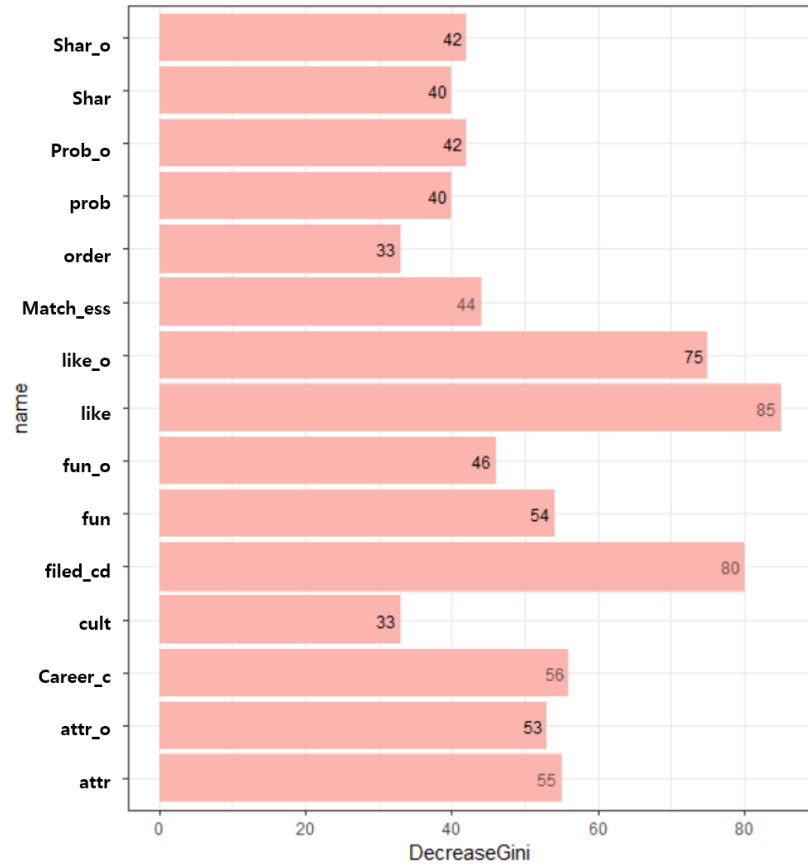


Train 데이터

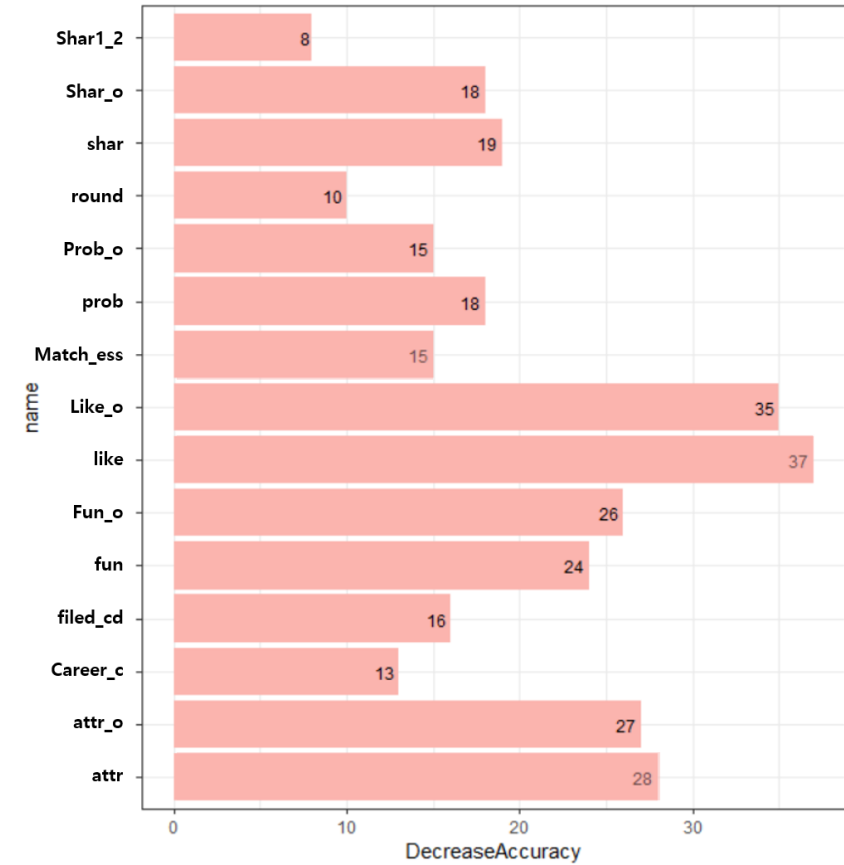
**Train70** : 모형 설계  
randomForest 사용

**Test30** : 모형 확인  
OOB 및 confusionMatrix 사용

# 분석 1. 랜덤포레스트

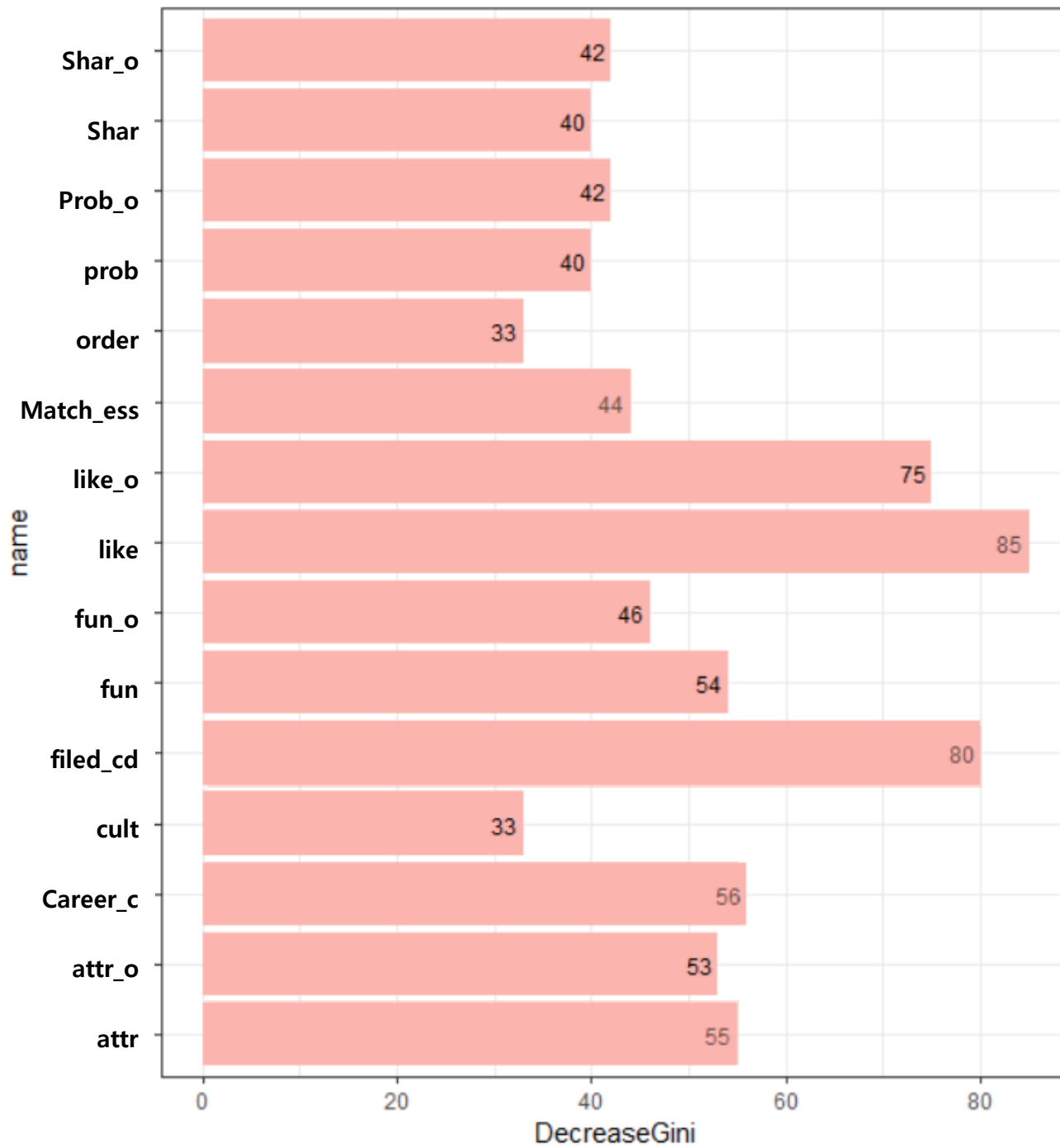


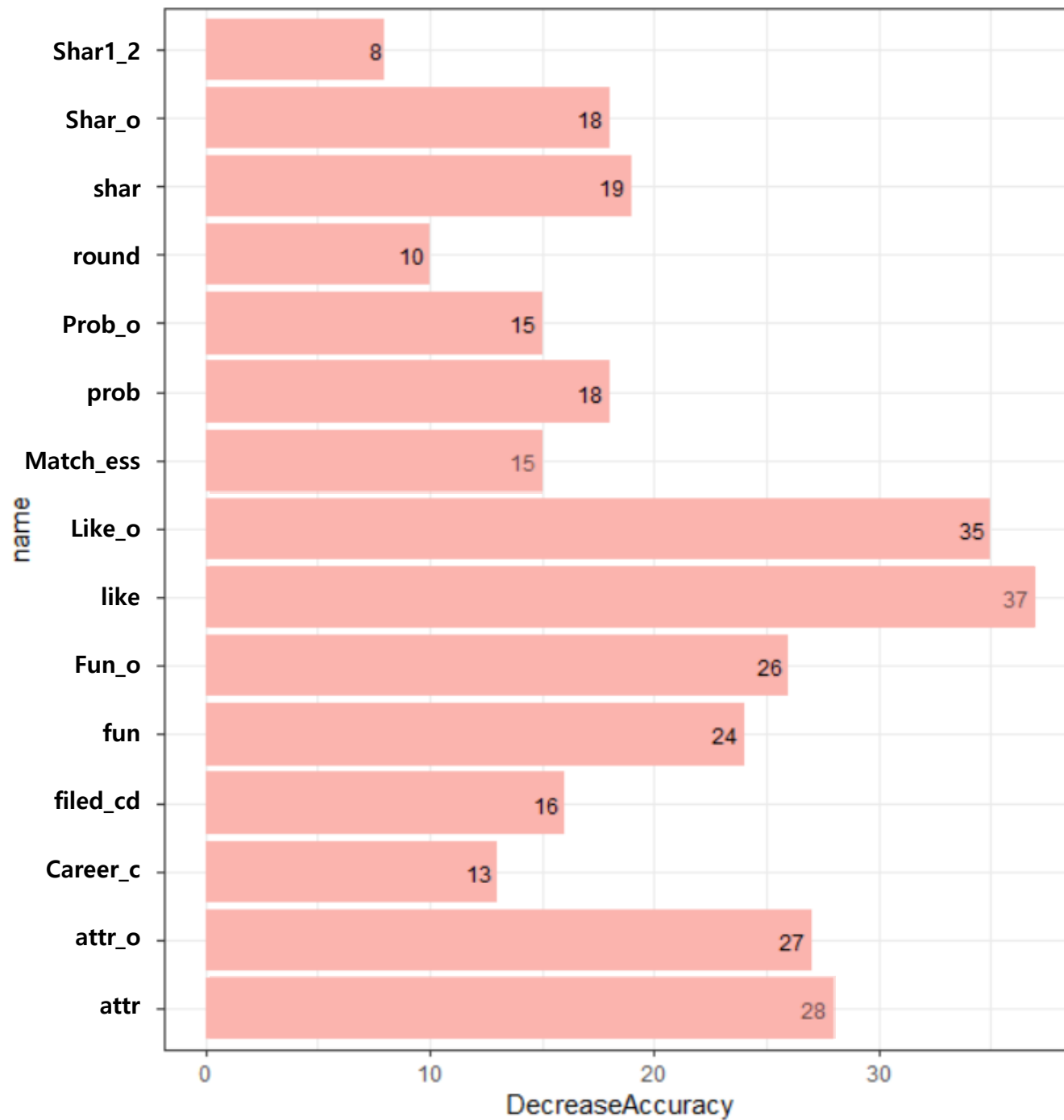
**like attr match\_ess**  
**field\_cd career\_c**



**like attr match\_ess**







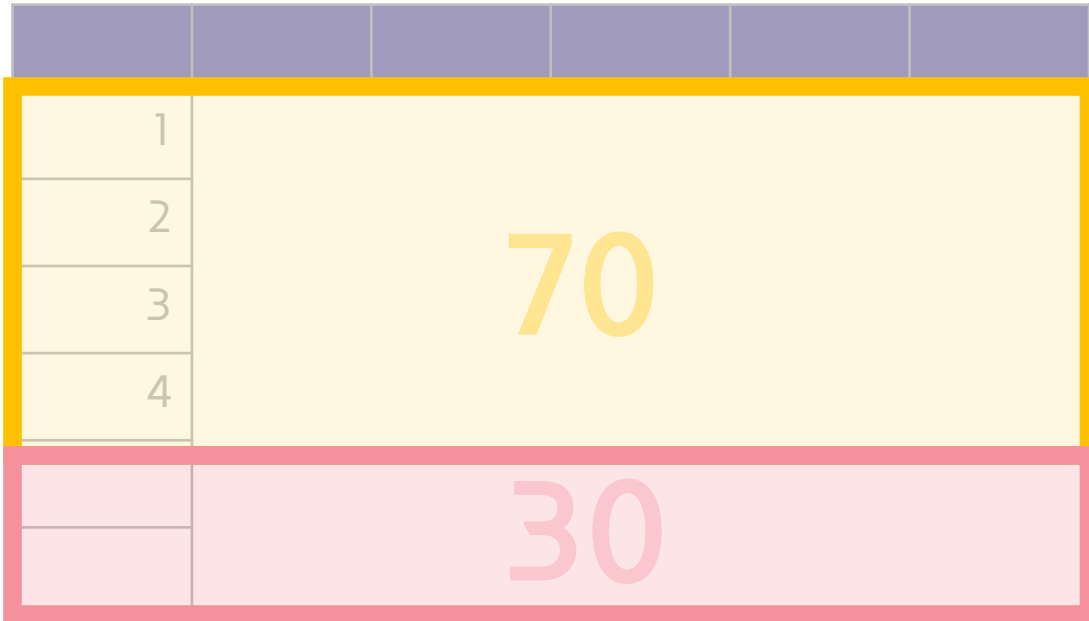
# 분석 1. 랜덤포레스트

		match	
		0	1
prediction	0	1569	240
	1	50	100

Confusion Matrix  
test set 적용

✓ accuracy = 0.855

## 분석 2. 로지스틱



Train 데이터

**Train70 : 모형 설계**

glm을 활용한  
로지스틱 회귀 모형

**Test30 : 모형 확인**

roc곡선, auc,  
confusionMatrix

**4번 반복하여 비교**

## 분석 2. 로지스틱

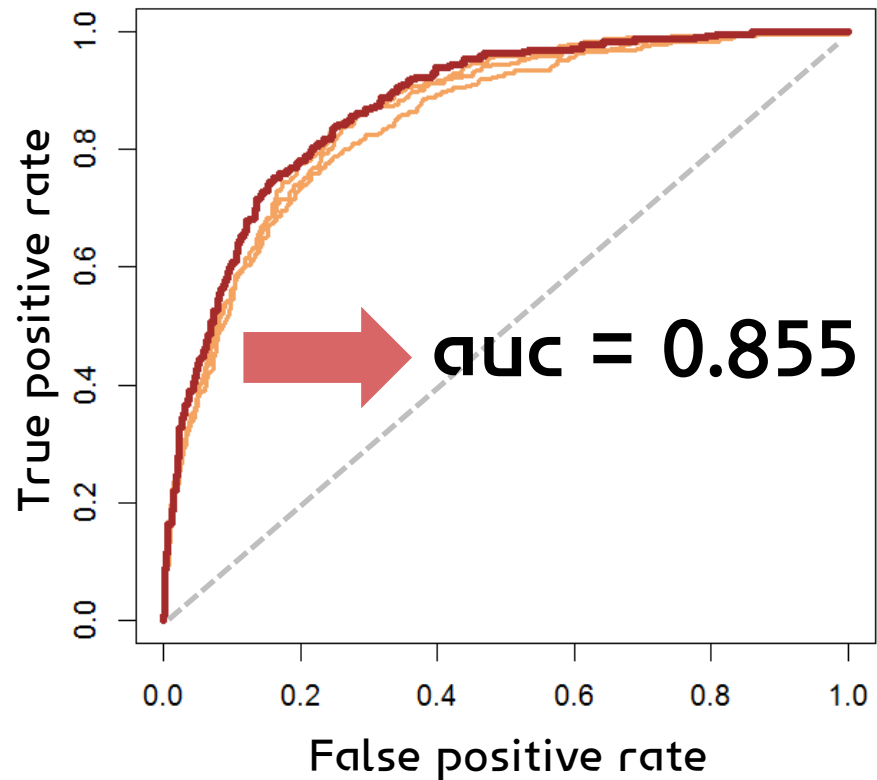
match 예측하는 로지스틱 회귀 모형 fitting  
=> 유의하지 않은 변수, 다중 공선성 문제



- 변수 내 무의미한 범주 통합
- stepwise로 변수 제거

## 분석 2. 로지스틱

ROC curve



		match	
		0	1
prediction	0	1569	203
	1	63	123

**accuracy = 0.864**

다음에 갖는 모형을 최종적으로 선택

## 분석 2. 로지스틱

	Odds ratio	Std.error	Pr(> z )
attr_o	1.38	0.03	<0.05
like	1.27	0.05	<0.05
attr	1.23	0.04	<0.05
prob_o	1.20	0.03	<0.05
....	....	....	....
shar_o	1.04	0.02	0.11
attr3_2	0.81	0.04	<0.05
likedif	0.76	0.04	<0.05
condtn2	0.31	0.22	<0.05

Summary 결과

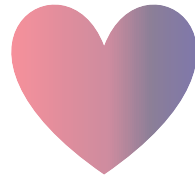
shar\_o, outside 변수 제외  
모든 변수 유의

Odds ratio를 통해  
match와의 관계 파악

# **[ 분석 결과 ]**

---





높을수록 매칭이  
성공하는 요인



shar / shar\_o

관심사를  
공유하는  
정도

match\_ess

예상  
매칭수

like

상대가  
마음에 든  
정도

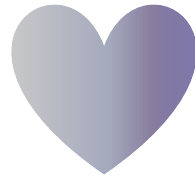
재미  
점수

fun / fun\_o

매력도  
점수

attr / attr\_o

**MATCH**



# 높을수록 매칭이 실패하는 요인

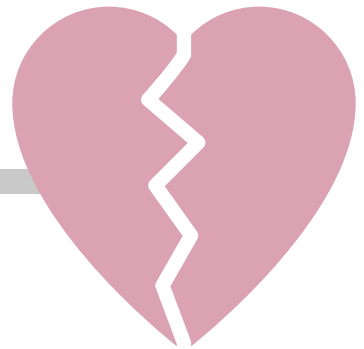
---

attr3\_2  
**자신을  
매력적으로  
생각하는 정도**

likedif  
**상대와 나의  
like점수차**

**지성**  
intel / intel\_o

**의욕적인  
정도**  
amb / amb\_o



**감사합니다**

---