

IGAWorks

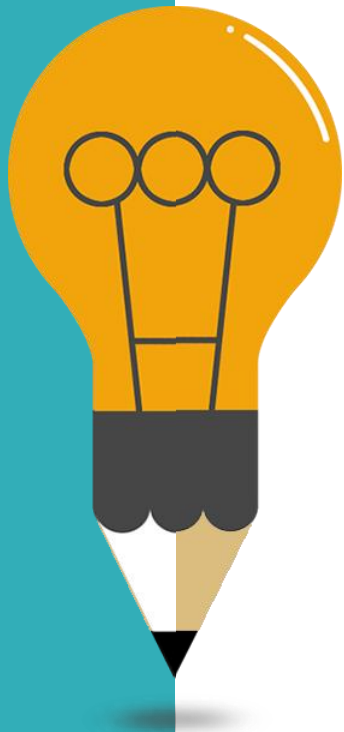
BIG DATA COMPETITION

CTR Prediction

팀명: 한강수

팀원: 김연강, 이승한, 차수만

목차

**01****Introduction**

데이터 소개 & EDA

02**Data Preprocessing**

변수 선택

03**Modeling**

DeepFM (Deep Factorization Machine)

04**Result**



1. Introduction

데이터 소개 & EDA

1. train, test

학습/평가기간 노출 로그 데이터

(train: 5,500,000 rows,
test: 550,000 rows)

2. audience_profile

오디언스 관련 정보 모음

(10,000,001 rows)

1-1. 데이터 소개

1. Train.csv & Test.csv

구분	변수명	변수 설명	카테고리 개수
광고 데이터	ssp_id	SSP 아이디	17
	campaign_id	캠페인 아이디	196
	adset_id	광고 아이디	924
	placement_type	광고 타입	4
	media_id	미디어 아이디	6010
	media_name	미디어 한글 이름	7031
	media_bundle	미디어 앱명	6504
	media_domain	미디어 도메인	162
	publisher_id	매체사 아이디	4062
	publisher_name	매체사 이름	1249
	advertisement_id	광고주 아이디	31

1. 데이터 소개

1. Train.csv & Test.csv

구분	변수명	변수 설명	카테고리 개수
기기 데이터	device_ifa	기기 구별 아이디	1,869,137
	device_os	기기 OS	2
	device_os_version	기기 OS 버전	125
	device_model	기기 모델명	1,664
	device_carrier	기기 통신사	536
	device_make	기기 제조사	299
	device_connection_type	기기 연결방식	8
	device_language	기기 언어	33
	device_country	기기 국가	1
	device_region	기기 지역	148
	device_city	기기 도시	1,326
기타 데이터	click	클릭 여부	2
	event_datetime	로그 발생 시간	5,478,966

1. 데이터 소개

2. audience_profile.csv

구분	변수명	변수 설명	카테고리 개수
고객 데이터	device_ifa	기기 구별 아이디	1,869,137
	age	연령 (추정)	12
	gender	성별 (추정)	2
	marry	기혼여부 (추정)	2
	install_pack	설치된 앱 정보	767,235
	cate_code	IGAW 카테고리별 등급	767,235
	predicted_house_price	자산 가격 (추정)	3466

1-1. 데이터 소개

Problem

Train과 Test의 bid_id가, audience_profile에 전부 있는 것은 아님

1-1. 데이터 소개

Problem

Train과 Test의 bid_id가, audience_profile에 전부 있는 것은 아님



Solution

223만

277만

audience_profile 정보가 **있는** 사람과 **없는** 사람으로 구분

[Data 1] **train** (기존과 동일)

[Data 2] **train & audience_profile**를 merge

1-1. 데이터 소개

[Data 1] train/test

train 데이터 전체

(5,500,000 rows)



audience 정보가 없는 test 데이터

(366,391 rows)

No_aud_merged (5,866,391 rows)

[Data 2] train/test & audience_profile

audience 정보가 있는 train 데이터만

(2,232,520 rows)



audience 정보가 있는 test 데이터만

(183,609 rows)

aud_merged (2,416,129 rows)

1-1. 데이터 소개

[Data 1] **train/test**

train 데이터 전체

(5,508,570 rows)
ap가 **없**는 test data는 **Data1**로 만든 Model로,

ap가 **있**는 test data는 **Data2**로 만든 Model로 !

audience 정보가 **없**는 test 데이터

(366,391 rows)

[Data 2] **train/test & audience_profile**

audience 정보가 있는 train 데이터만

audience 정보가 **있**는 test 데이터만

(audience_profile : ap)

(183,609 rows)

No_aud_merged (5,866,391 rows)

aud_merged (2,416,129 rows)

1-2. EDA

GOAL : 유저가 특정 광고를 클릭할 지 여부를 파악! ($y = \text{click}$)

1.train

1. 날짜 데이터

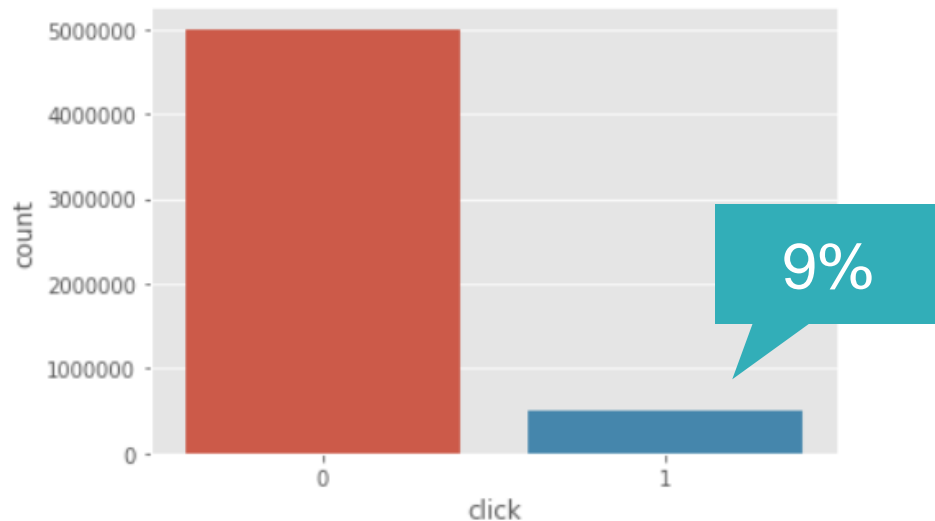
2. 광고 데이터

3. 기기 데이터

4. 기타 데이터



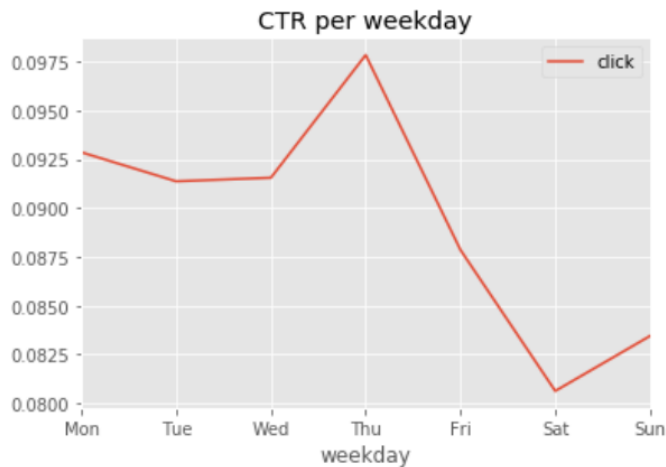
2. audience_profile



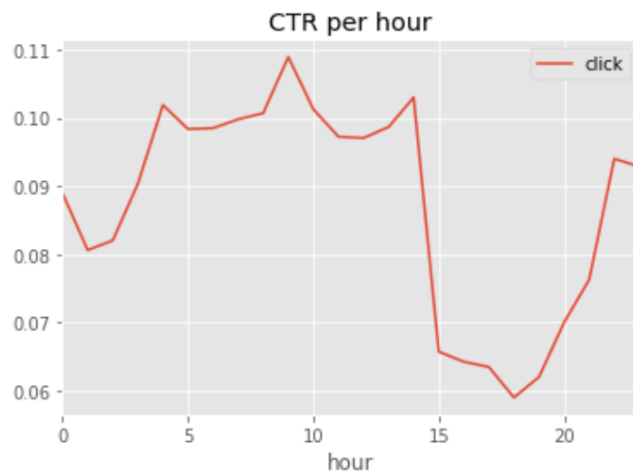
Imbalanced Dataset

1-2. EDA (1) train

1. 날짜 데이터



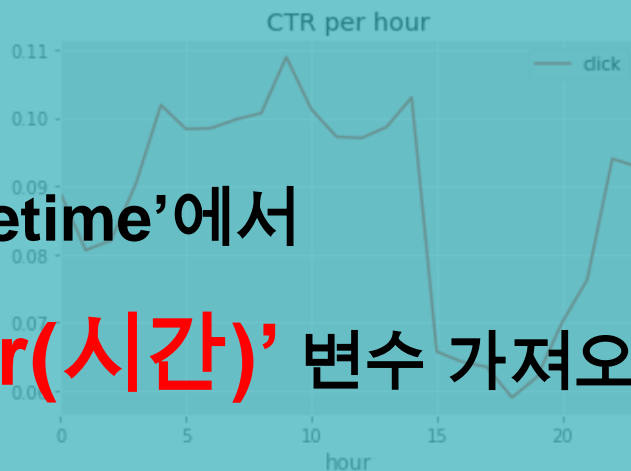
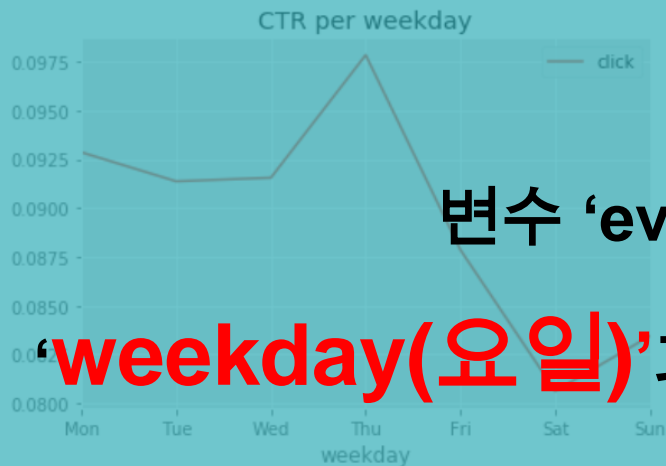
요일 별로 다른 click 비율



시간대(0~23시) 별로 다른 click 비율

1-2. EDA (1) train

1. 날짜 데이터



변수 'event_datetime'에서

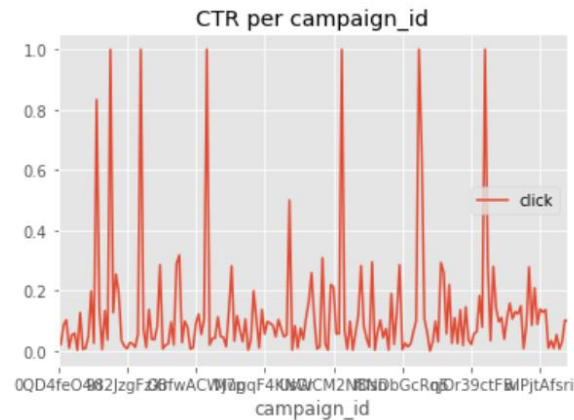
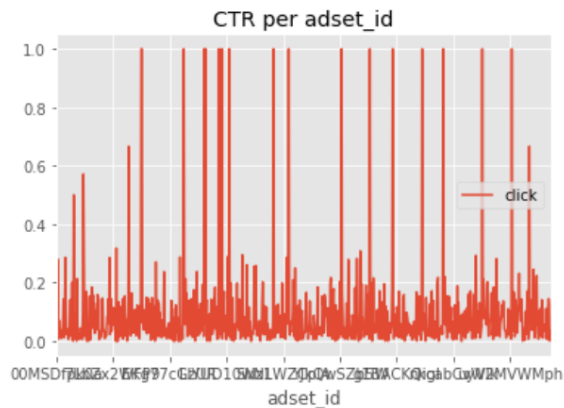
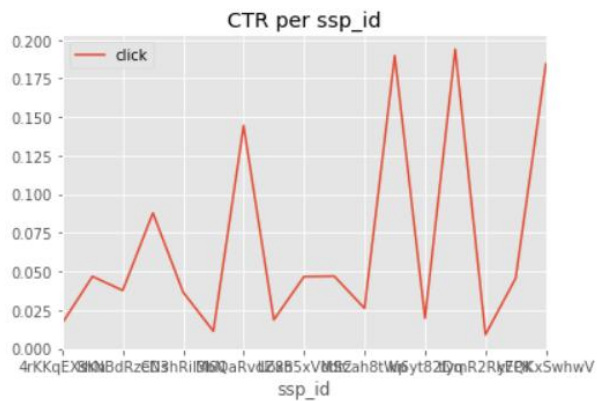
‘weekday(요일)’과 ‘hour(시간)’ 변수 가져오기

요일 별로 다른 click 비율

시간대(0~23시) 별로 다른 click 비율

1-2. EDA (1) train

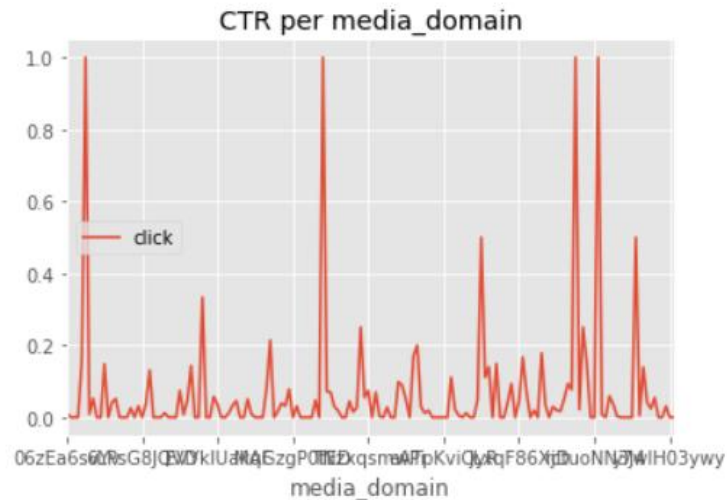
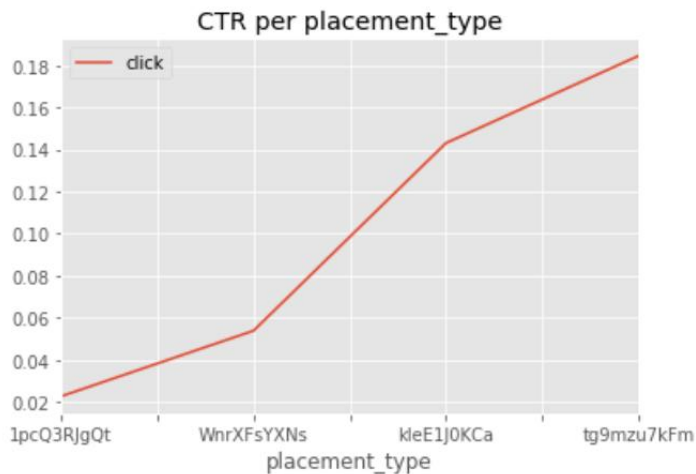
2. 광고 데이터



ssp_id, adset_id, campaign_id 별로 다른 click 비율

1-2. EDA (1) train

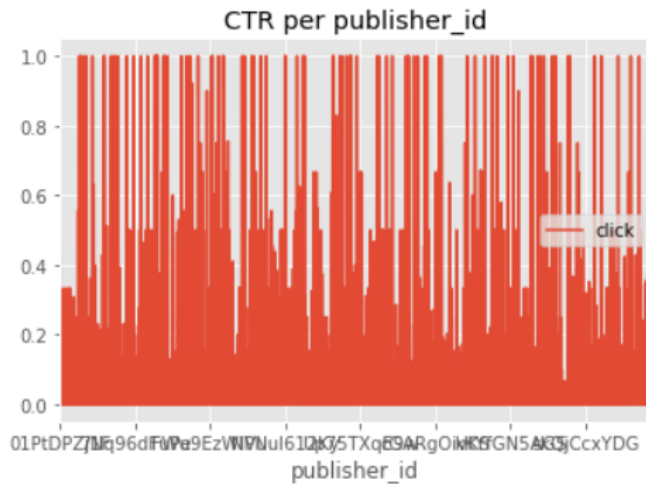
2. 광고 데이터



placement_type, media_domain 별로 다른 click 비율

1-2. EDA (1) train

2. 광고 데이터



publisher_name, publisher_id 별로 다른 click 비율

1-2. EDA (1) train

2. 광고 데이터

media

media_id	media_name	media_bundle	media_domain
lyDyyhXBnW	dAWR8DOmzo	V8yCzbCKcB	pjBT3sDGbH

publisher

publisher_id	publisher_name
YOSTF4U4h4	tXBXkBEsgT

다중공산성 문제?

1-2. EDA (1) train

2. 광고 데이터

media

media_id	media_name	media_bundle	media_domain
lyDyyhXBnW	dAWR8DOmzo	V8yCzbCKcB	pjBT3sDGbH

publisher

publisher_id	publisher_name
YOSTF4U4h4	tXBXkBEsgT



BUT 특정 변수를 **버리지 않는 것**이 더 나은 결과를 가져왔음
+ **다른 변수와의 상호작용** 고려 시, 다른 영향을 가질 수도

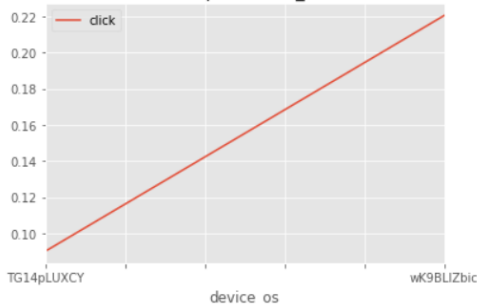
다중공산성 문제?

전부 사용!

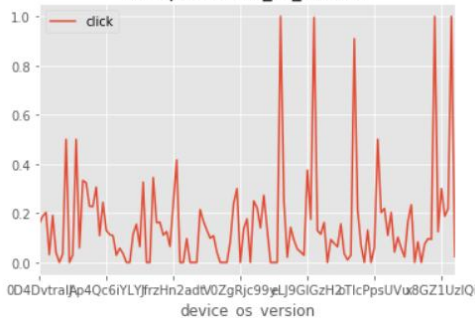
1-2. EDA (1) train

3. 기기 데이터

CTR per device_os



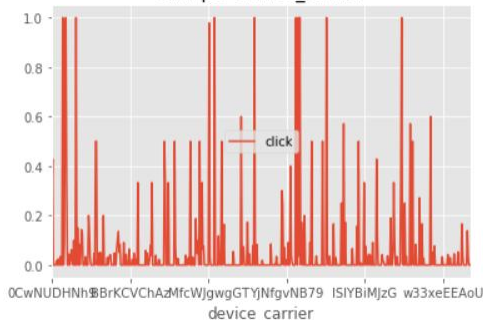
CTR per device_os_version



CTR per device_model



CTR per device_carrier

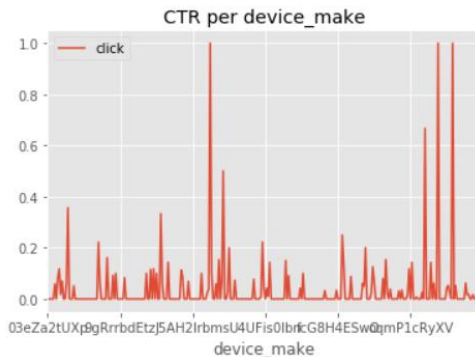
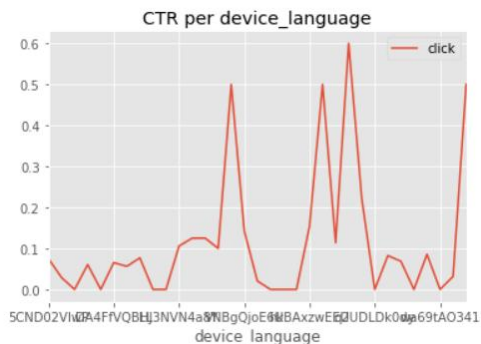
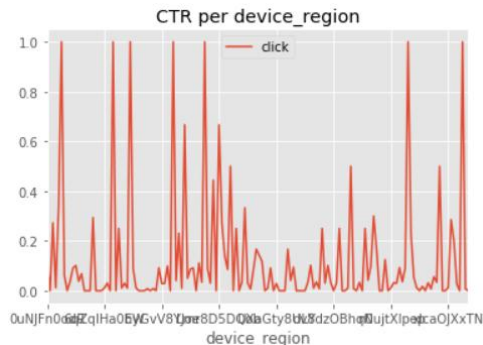
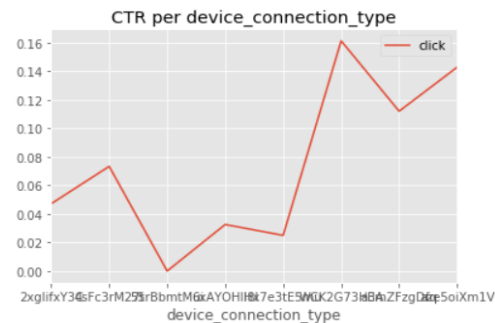


각 종 기기(device) 정보 별로 다른 click률!

(device_os , device_os_version,
device_model, device_carrier)

1-2. EDA (1) train

3. 기기 데이터

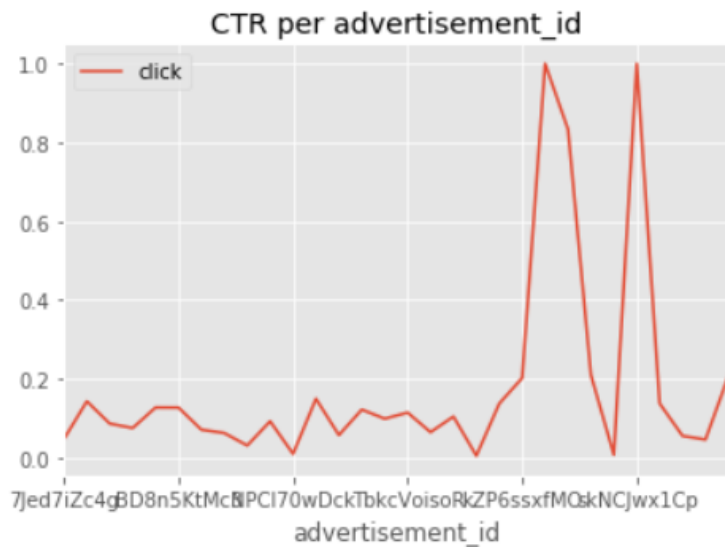


각 종 기기(device) 정보 별로 다른 click률!

(device_connection_type, device_region,
Device_language, device_make)

1-2. EDA (1) train

4. 기타 데이터

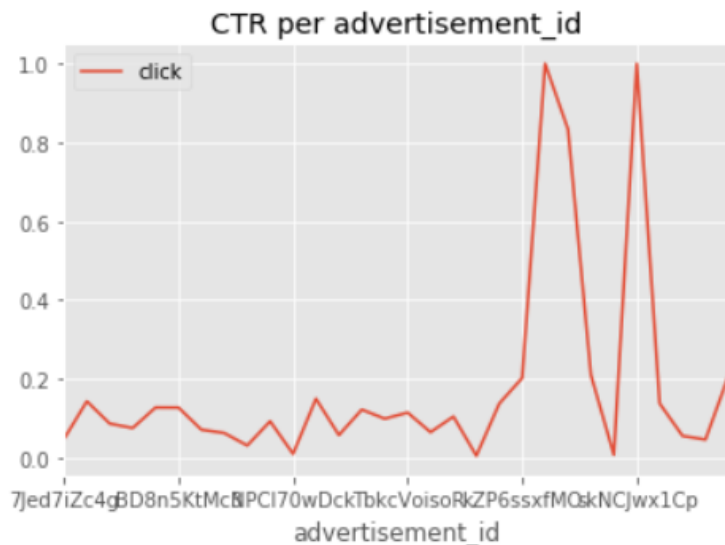


광고(advertisement) 정보 별로 다른 click률!

(adset_id & advertisement_id)

1-2. EDA (1) train

4. 기타 데이터



광고(advertisement) 정보 별로 다른 click률!

(adset_id & advertisement_id)



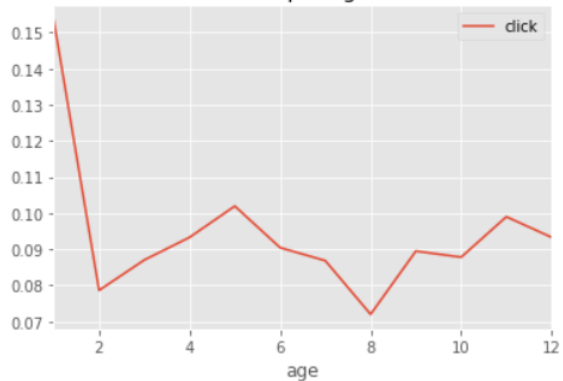
두 변수 중 하나만 사용하는 것 보다,

둘 다 사용했을 때 더 좋은 결과!

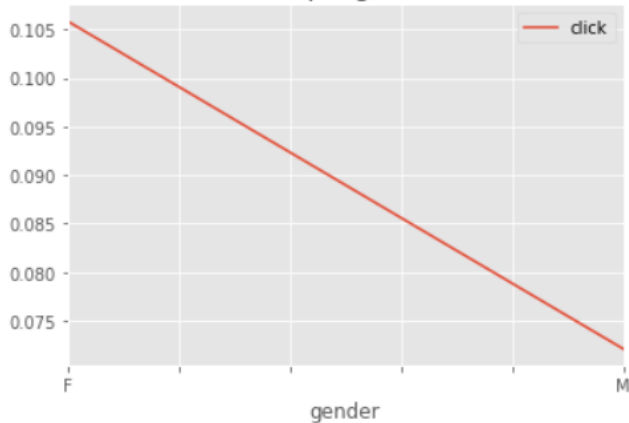
1-2. EDA (2) audience

1. 개인 정보

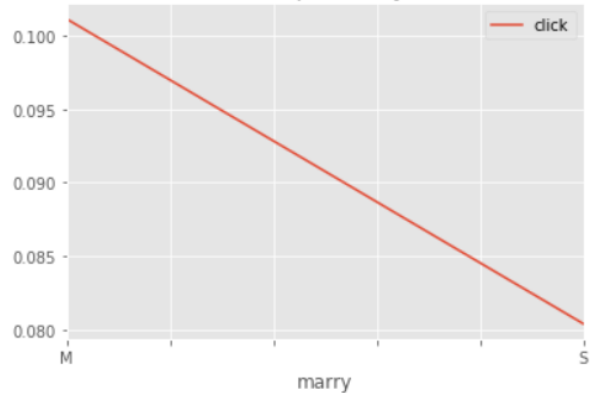
CTR per age



CTR per gender



CTR per marry



나이(age), 성별(gender), 결혼 여부(marry)에 따라 click률에 큰 차이가 있음

1-2. EDA (2) audience

2. 앱 설치 정보

install_pack

```
train_merged_aud['install_pack'].value_counts()
```

```
p20150g, p32857g, p12957g, p17595g, p16464g, p25493g, p25785g, p12920g, p4320g, p9102g, p21616g, p26644g, p26662g, p26674g, p6935g, p9951g, p10883g, p4786g, p29512g, p29051g, p35286g, p22119g, p10176g, p6450g, p6933g, p15194g, p14938g, p36404g, p22125g, p3879g, p26892g, p16531g, p15365g, p16466g, p18585g, p30097g, p20142g, p47955g, p10831g, p28081g, p3481g, p25784g, p21103g, p21742g, p27007g, p10954g, p25026g, p13001g, p22114g, p12968g, p12389g, p41431g, p38837g, p21512g, p7098g, p12912g, p16460g, p8458g, p12346g, p17586g, p15448g, p18604g, p6818g, p21154g, p22121g, p10674g, p22766g, p12948g, p20149g, p25499g, p14329g, p23424g, p22786g, p13413g, p16534g, p21525g, p26661g, p25787g, p26696g, p7583g, p25782g, p10178g, p26622g, p12950g, p25478g, p48173g, p26700g, p26916g, p18202g, p12928g, p37605g, p25504g, p6071g, p25792g, p29688g, p22106g, p17600g, p11116g, p26669g, p21747g, p26621g, p34971g, p33379g, p18578g, p26677g, p23840g, p26699g, p12965g, p30448g, p31442g, p25483g, p23418g, p6637g, p5095g, p20156g, p22126g, p21812g, p38244g, p16489g, p21472g, p25789g, p44448g, p15296g, p23256g, p43404g, p17533g, p22120g, p25485g, p6029g, p15192g, p31070g, p22105g, p41852g, p16517g, p6531g, p3470g, p39923g, p39241g, p34342g, p20147g, p26660g
```

```
mean = np.mean(pack_psudeo_length100/pack_length100)
mean
```

7.853241826853476

(평균적으로 7.85개의 앱을 설치)

Name: install_pack, Length: 767235, dtype: int64

[PROBLEM] 너무 많은 ‘설치한 앱’ 종류 & 개인 별로 상이한 ‘설치한 앱의 종류’

1-2. EDA (2) audience

2. 앱 설치 정보

install_pack

```
train_merged_aud['install_pack'].value_counts()
```

```
p20150g,p32857g,p12957g,p17595g,p16464g,p25493g,p25785g,p12920g,p4320g,p9102g,p21616g,p26644g,p26662g,p26674g,p6935g,p9951g,p10983g,p4786g,p29512g,p23051g,p35286g,p22119g,p10176g,p6450g,p6933g,p15194g,p14938g,p36404g,p22125g,p3879g,p26992g,p16531g,p15365g,p16466g,p18585g,p30097g,p20142g,p47855g,p10631g,p28061g,p3481g,p26784g,p21103g,p21742g,p27007g,p10954g,p25026g,p19001g,p22114g,p12969g,p12899g,p41481g,p38837g,p21512g,p7098g,p12912g,p16460g,p8459g,p12345g,p17596g,p16448g,p18604g,p6818g,p21154g,p22121g,p10674g,p22769g,p12948g,p20149g,p25499g,p14329g,p23424g,p22786g,p13413g,p16534g,p21525g,p26661g,p25787g,p26696g,p7589g,p25782g,p10178g,p26622g,p12950g,p25478g,p48173g,p26700g,p26916g,p18202g,p12928g,p37905g,p25504g,p6071g,p25782g,p29688g,p22106g,p17600g,p11116g,p26668g,p21747g,p26621g,p34971g,p33379g,p18578g,p26677g,p23840g,p26689g,p12965g,p30448g,p31442g,p25483g,p23418g,p6637g,p5095g,p20156g,p22126g,p21812g,p38244g,p16489g,p21472g,p25789g,p44448g,p15296g,p23259g,p43404g,p17533g,p22120g,p25495g,p6029g,p15192g,p31070g,p22105g,p41652g,p16517g,p6531g,p3470g,p3923g,p39241g,p34324g,p20147g,p26660g
```

```
mean = np.mean(pack_psudeo_length100/pack_length100)
mean
```

7.853241826853476

이렇게 다양한 categorical 변수들을,

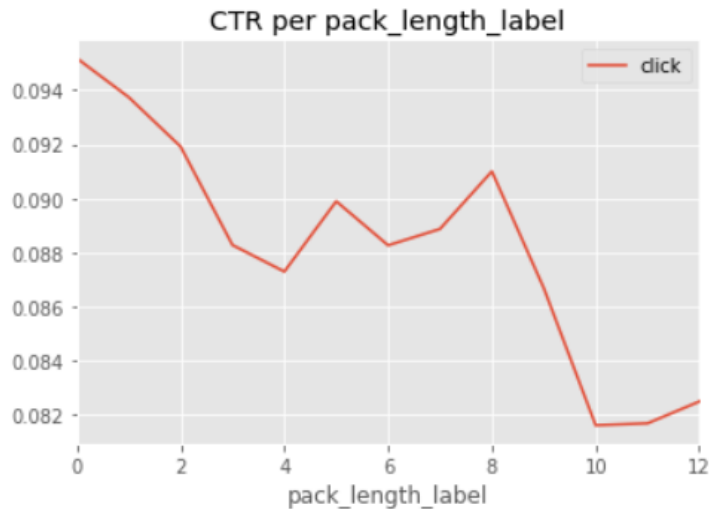
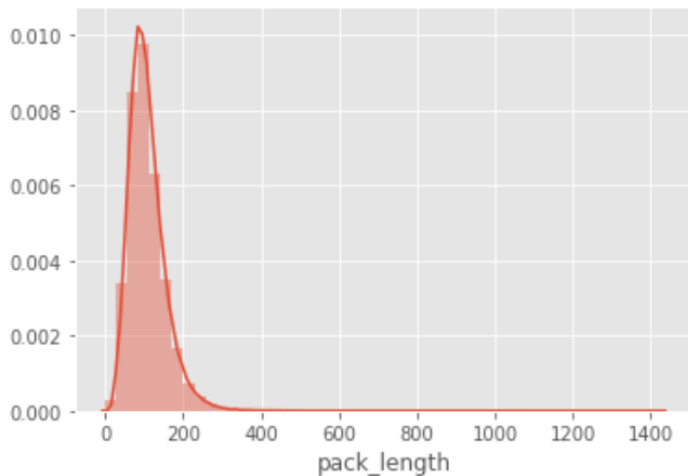
(평균적으로 7.85개의 앱을 설치)

“설치한 앱의 개수(pack_length)”라는 하나의 numeric 변수로!

1-2. EDA (2) audience

2. 앱 설치 정보

새로운 변수, “**pack_length**” (설치한 앱의 개수)



Click률에 있어서 **유의미한 차이를 보이는 pack_length**

Data Overview Summary

너무 많은 변수들 !

(버리는 변수들 거의 없음 & category들 간의 grouping도 안함)

+ 심지어 categorical variables들의 dummy화까지 하게 되면?

1-2. EDA (2) audience

Data Overview Summary

너무 많은 변수들 !

(버리는 변수들 거의 없음 & category들 간의 grouping도 안함)

+ 심지어 categorical variables들의 dummy화까지 하게 되면?

NO WORRY! Pre-trained Model을 사용!



2. Data Preprocessing

2. Data Preprocessing

Data2 기준으로 소개

(Data1도 audience_profile의 feature 없다는 점 제외하고 동일)

1. Data Merge (공통 column인 Bid_id를 기준으로)

[Train / Test]

Bid_id	특징1	특징2
A
B
C



[Audience_Profile]

Bid_id	특징3	특징4
A
C
F

2. Data Preprocessing

Data2 기준으로 소개

(Data1도 audience_profile의 feature 없다는 점 제외하고 동일)

1. Data Merge (공통 column인 Bid_id를 기준으로)

[Train / Test]

Bid_id	특징1	특징2
A
B
C



[Audience_Profile]

Bid_id	특징3	특징4
A
C
F



[Data 2 (aud_merged)]

Bid_id	특징1	특징2	특징3	특징4
A	
C	

2. Data Preprocessing

Data2 기준으로 소개

(Data1도 audience_profile의 feature 없다는 점 제외하고 동일)

1. Data Merge (공통 column인 Bid_id를 기준으로)

[Train / Test]

Bid_id	특징1	특징2
A
B
C



[Audience_Profile]

Bid_id	특징3	특징4
A
C
F



[Data 2 (aud_merged)]

Bid_id	특징1	특징2	특징3	특징4
A	
C	

그대로 사용

[Data 1 (no_aud_merged)]

Bid_id	특징1	특징2
A
B
C

2. 변수 추가 및 제거

변수 추가

1. hour(로그 발생 시간(시)) - 접속한 시간대 (대낮에 길가면서 핸드폰 할 때 vs 자기 직전에 핸드폰 할 때)
2. dayofweek(로그 발생 요일) - 홍보하는 기업에 따라 맞춤형 요일이 있을 수도!

변수 제거

1. device_country - 모든 data가 동일한 값을 가짐
2. device_ifa - 너무 많은 unique 값!
3. event_datetime - 위에서 hour & min을 뽑아서 사용함
4. predicted_house_price - 추가 시 오히려 성능 저하 (영향을 미치지 않는 요소 or 추정치의 신뢰 문제)
5. cate_code
6. install_pack

2. Data Preprocessing

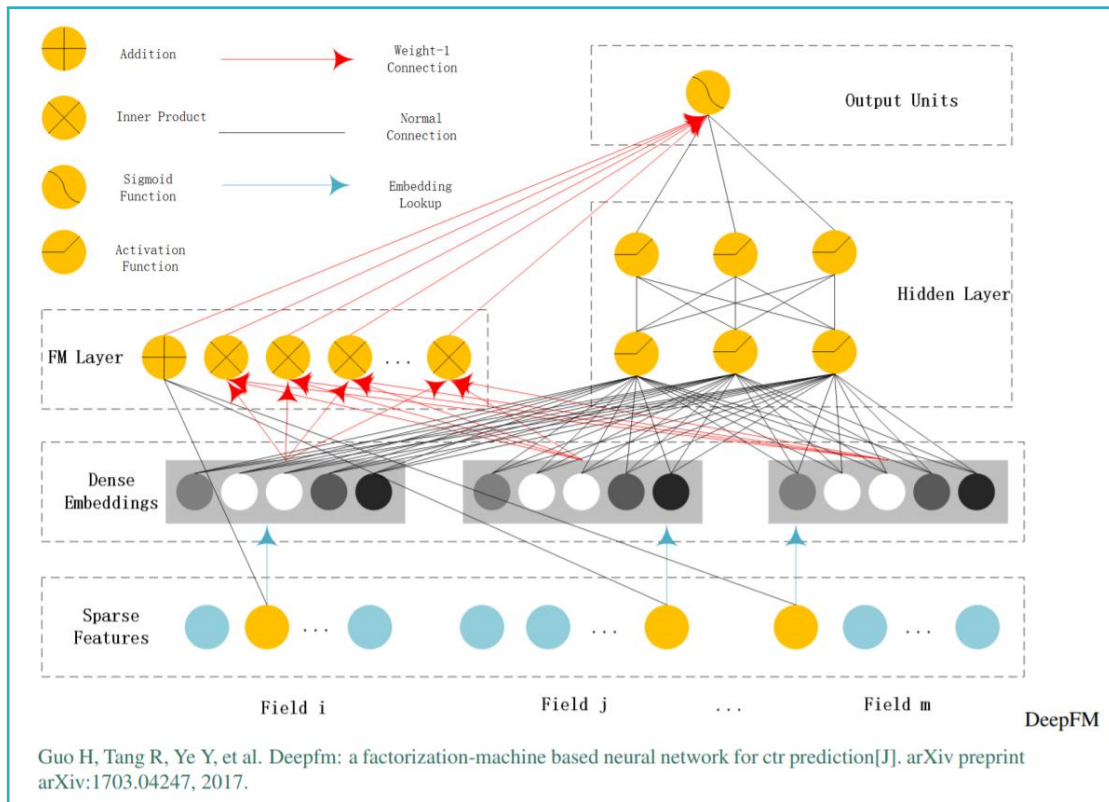
3. Label Encoding

(추후에 사용할 DeepFM이라는 모델을 위해서 필요한 전처리 과정)



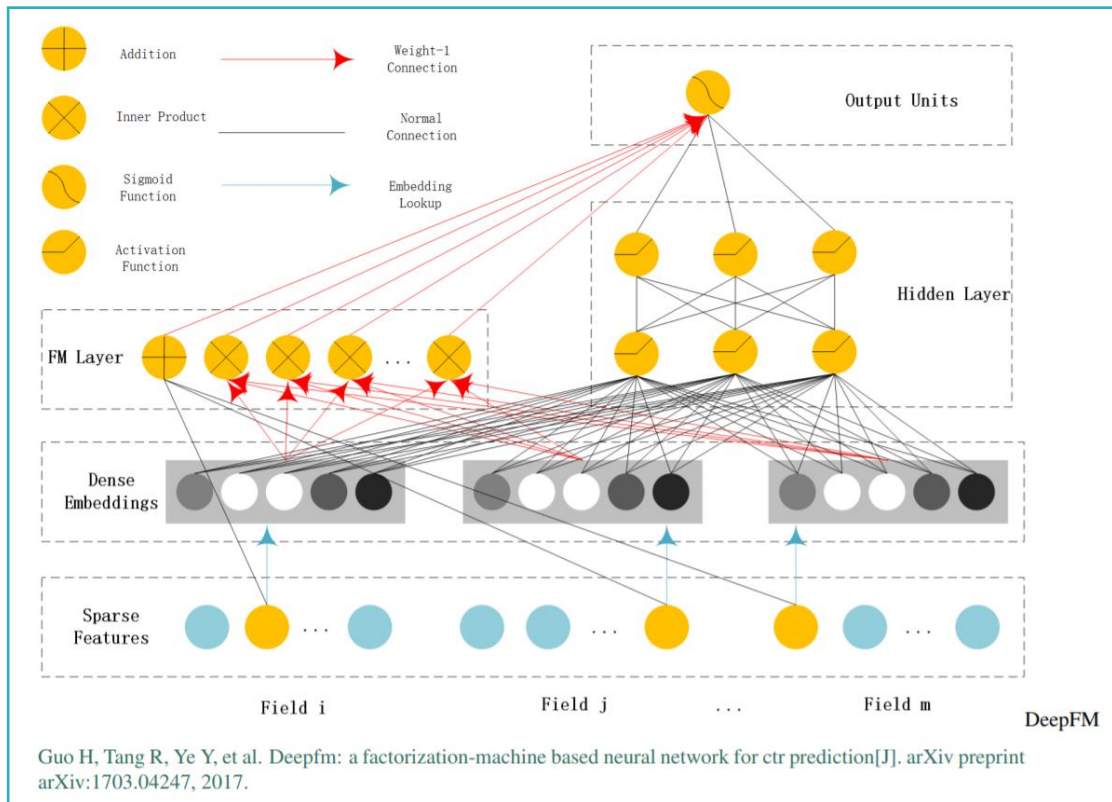
	adset_id	advertisement_id	age	bid_id	campaign_id	cate_code
0	109	19	7	aAEDD9Aelv	3	537223
1	109	19	7	120KZBpPEp	3	537223
2	117	19	7	AMFiNF3X7r	3	537223
3	187	18	7	Mza3hx3DOX	110	537223
4	117	19	7	4GbWwwNnJZ	3	537223

3. Modeling



Deep FM

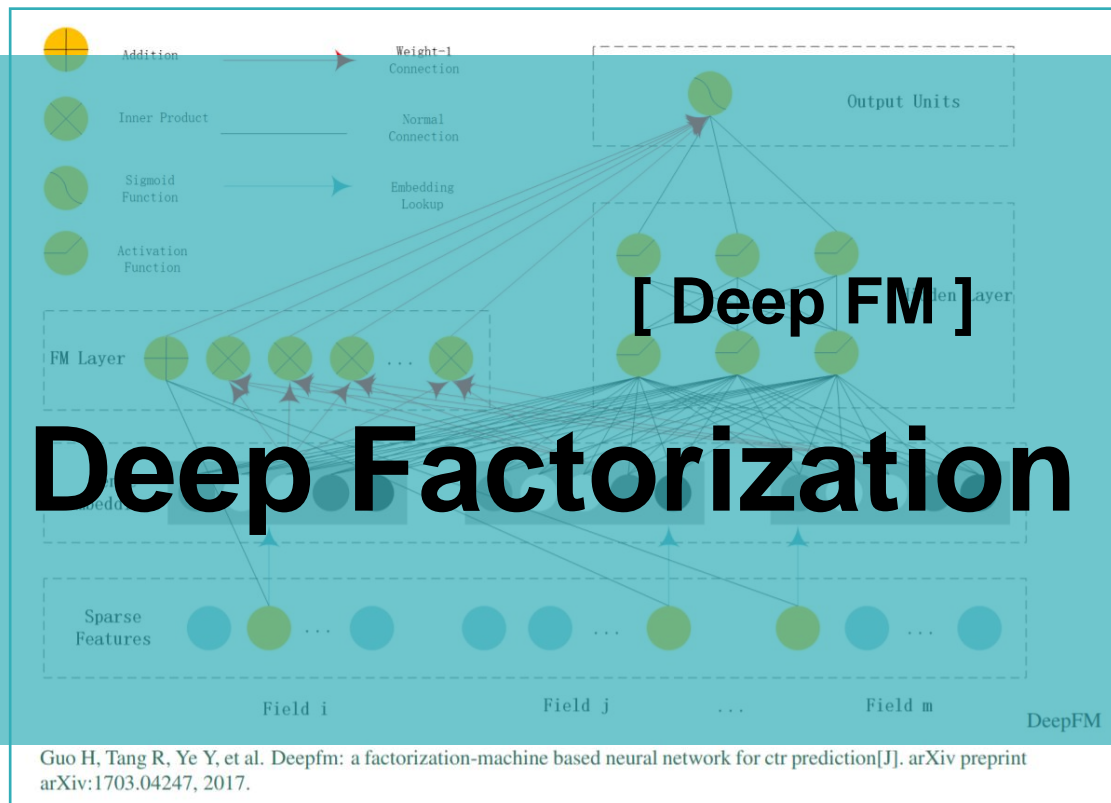
3. Modeling



Deep FM

?

3. Modeling



Deep FM

Deep Factorization Machine?

Guo H, Tang R, Ye Y, et al. Deepfm: a factorization-machine based neural network for ctr prediction[J]. arXiv preprint arXiv:1703.04247, 2017.

3. Modeling

[Deep FM] Deep Factorization Machine

Factorization Machine

example

	취미	전공	Y
A	독서	경영	1
B	골프	컴과	0
C	야구	철학	1

취미 & 전공  Y

3. Modeling

Factorization Machine

Linear model

$$W_{\text{독서}} + W_{\text{경영}}$$

[Deep FM] Deep Factorization Machine

	취미	전공	Y
A	독서	경영	1
B	골프	컴과	0
C	야구	철학	1

3. Modeling

Factorization Machine

Linear model

$$W_{\text{독서}} + W_{\text{경영}}$$



Too Simple

[Deep FM] Deep Factorization Machine

	취미	전공	Y
A	독서	경영	1
B	골프	컴과	0
C	야구	철학	1

3. Modeling

[Deep FM] Deep Factorization Machine

Factorization Machine

Linear model

$W_{\text{독서}} + W_{\text{경영}}$ ➡ Too Simple

Polynomial model

$W_{\text{독서}} + W_{\text{경영}} + W_{\text{독서, 경영}}$

	취미	전공	Y
A	독서	경영	1
B	골프	컴과	0
C	야구	철학	1

3. Modeling

[Deep FM] Deep Factorization Machine

Factorization Machine

Linear model

$W_{\text{독서}} + W_{\text{경영}}$ ➡ Too Simple

Polynomial model

$W_{\text{독서}} + W_{\text{경영}} + W_{\text{독서, 경영}}$ ➡ Overfitting & Unseen Data (ex. 골프&경영)

	취미	전공	Y
A	독서	경영	1
B	골프	컴과	0
C	야구	철학	1

3. Modeling

[Deep FM] Deep Factorization Machine

Factorization Machine

Linear model

$W_{\text{독서}} + W_{\text{경영}}$ ➡ Too Simple

Polynomial model

$W_{\text{독서}} + W_{\text{경영}} + W_{\text{독서, 경영}}$ ➡ Overfitting & Unseen Data (ex. 골프&경영)

Factorization Machine

$W_{\text{독서}} + W_{\text{경영}} + V_{\text{경영}} \cdot V_{\text{독서}}$

	취미	전공	Y
A	독서	경영	1
B	골프	컴과	0
C	야구	철학	1

3. Modeling

[Deep FM] Deep Factorization Machine

Factorization Machine

	취미	전공	Y
A	독서	경영	1
B	골프	컴과	0
C	야구	철학	1

Linear model

$W_{\text{독서}} + W_{\text{경영}} \Rightarrow$ Too Simple

Polynomial model

$W_{\text{독서}} + W_{\text{경영}} + W_{\text{독서, 경영}} \Rightarrow$ Overfitting & Unseen Data (ex. 골프&경영)

Factorization Machine

$W_{\text{독서}} + W_{\text{경영}} + \underline{V_{\text{경영}} \cdot V_{\text{독서}}}$

각각의 값들이 “**latent vector**” 를 가지고 있음!

(unseen data를 다룰 수 있음! (ex. 이전에는 없던 “취미가 독서”이고 “전공이 컴과:인 사람))

Modeling

[Deep FM] Deep Factorization Machine

Factorization Machine

	취미	전공	Y
A	독서	경영	1
B	골프	컴과	0
C	야구	철학	1

[핵심] 어떻게 latent vectors들을 훈련시킬 것인가?

Polynomial model

$W_{\text{독서}} + W_{\text{경영}} + W_{\text{독서, 경영}}$



Overfitting & Unseen Data (ex. 골프&경영)

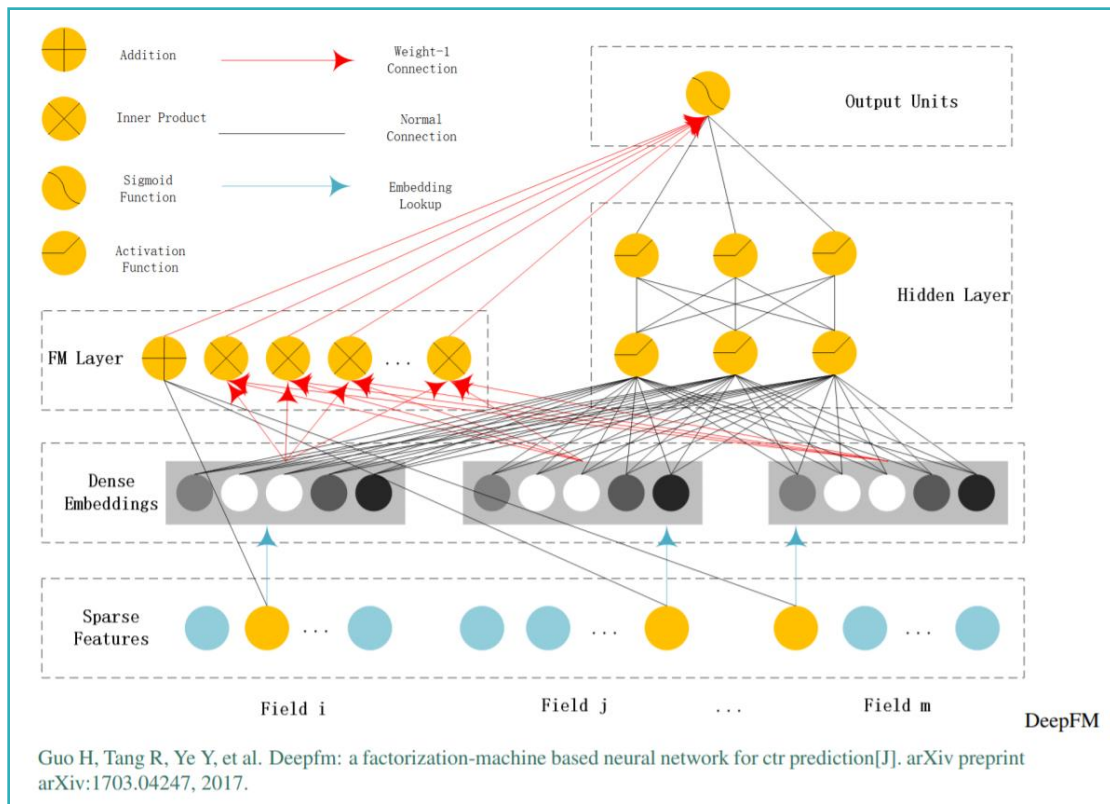
FM은 “Click Prediction” 에서 특히 잘 작동!

(각기 다른 선호도를 가지고, 방문한 사이트들과 설치 앱들이 다른 유저들)

have “latent vector” for each feature!

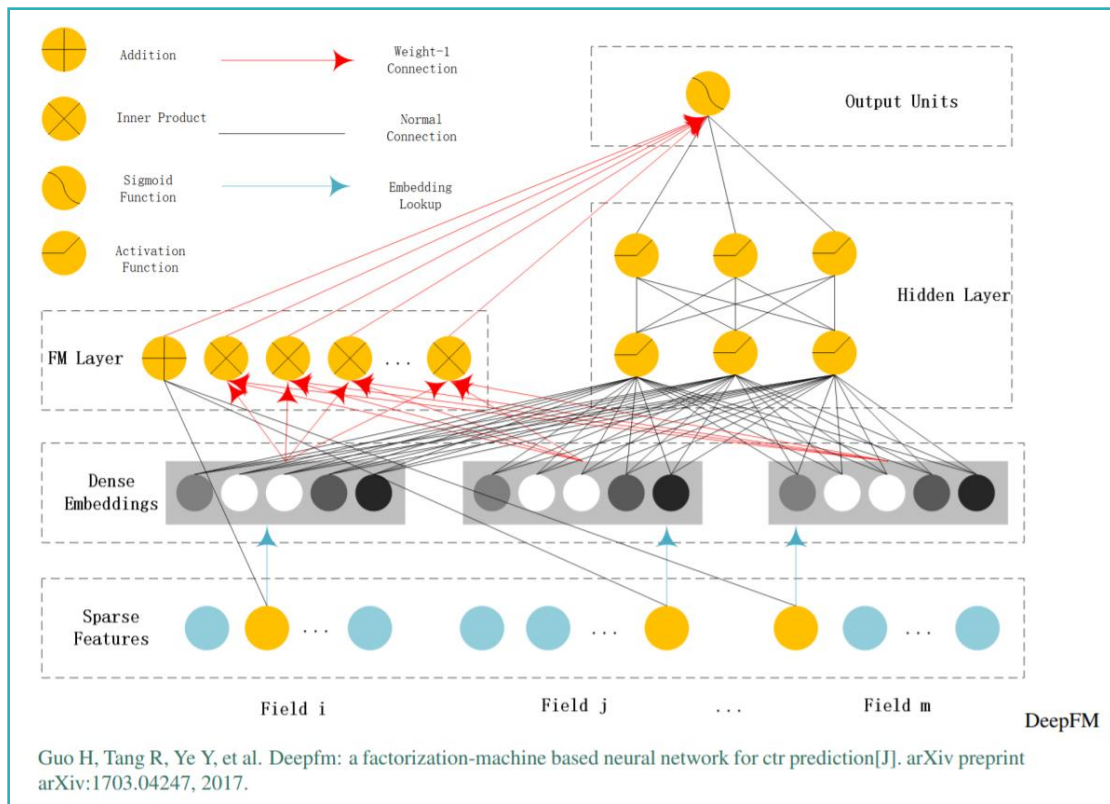
(Able to deal with unseen data! (ex. 취미: 독서, 전공: 컴과))

3. Modeling



How to train latent vectors?

3. Modeling



How to train latent vectors?

with “**Neural Network**”

Deep FM

3. Modeling



How to train latent vectors?

with “Neural Network”

why DeepFM?

Deep FM

Guo H, Tang R, Ye Y, et al. Deepfm: a factorization-machine based neural network for ctr prediction[J]. arXiv preprint arXiv:1703.04247, 2017.

3. Modeling

why DeepFM?

1. 대부분의 feature들이 categorical variable

- dummy화 할 경우, 매우 sparse해지는 문제 발생!

3. Modeling

why DeepFM?

1. 대부분의 feature들이 categorical variable
- dummy화 할 경우, 매우 sparse해지는 문제 발생!

➡ 저차원으로 embedding (4-dim)

3. Modeling

why DeepFM?

1. 대부분의 feature들이 categorical variable

- dummy화 할 경우, 매우 sparse해지는 문제 발생!

2. Unseen Data들을 다룰 수 있음

- ex) “A형”에 “럭비”가 취미인 “20kg”의 “6살”아이..

➡ 저차원으로 embedding (4-dim)

3. Modeling

why DeepFM?

1. 대부분의 feature들이 categorical variable

- dummy화 할 경우, 매우 sparse해지는 문제 발생!

➡ 저차원으로 embedding (4-dim)

2. Unseen Data들을 다룰 수 있음

- ex) “A형”에 “럭비”가 취미인 “20kg”의 “6살”아이..

➡ $V_{A형} \cdot V_{럭비} \cdot V_{한국} \cdot V_{남자}$

3. Modeling

why DeepFM?

1. 대부분의 feature들이 categorical variable

- dummy화 할 경우, 매우 sparse해지는 문제 발생!

➡ 저차원으로 embedding (4-dim)

2. Unseen Data들을 다룰 수 있음

- ex) “A형”에 “럭비”가 취미인 “한국”에 사는 “남자”

➡ $V_{A형} \cdot V_{럭비} \cdot V_{한국} \cdot V_{남자}$

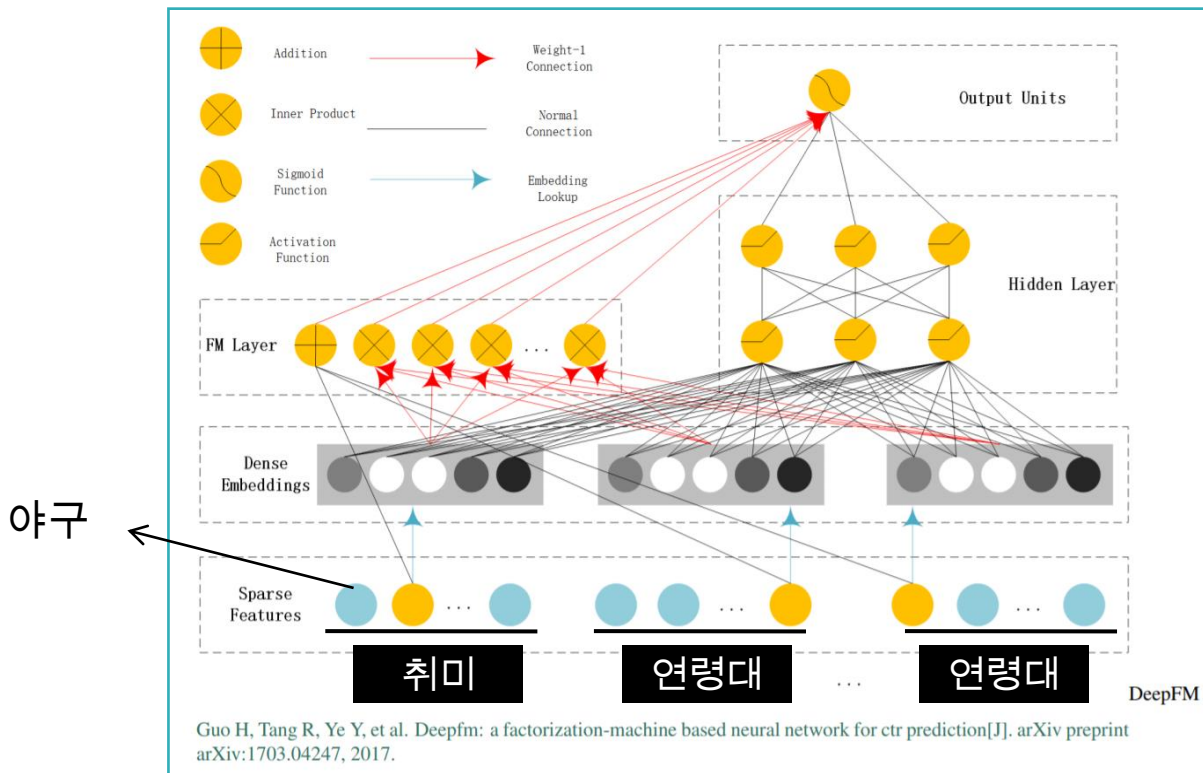
3. Pretrained Model

- 한 번 학습 이후, 새로운 데이터에 대해 빠르게 적용 가능

3. Modeling

$$y_{FM} = \langle w, x \rangle + \sum_{i=1}^d \sum_{j=i+1}^d \langle V_i, V_j \rangle x_i \cdot x_j$$

Example)



Summary



Thank you

Insert the title of your subtitle Here