TalkingData AdTracking Fraud Detection Challenge



CONTENTS









01

Competition Overview

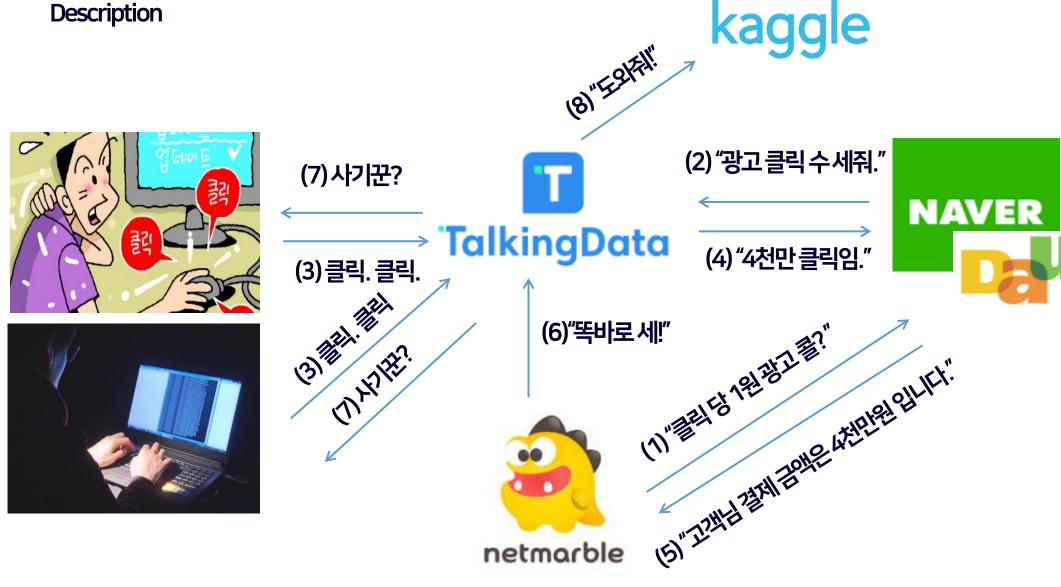
01. Competition Overview

Description

- 온라인 광고 회사는 클릭만으로도 광고 비용이 상승 할 수 있다.
- 사기성 클릭이 전체 클릭의 90% 이상이 잠재적인 사기성 클릭이다.
- 기존에는 클릭 수가 많은 추적하여 IP와 device를 이용한 블랙리스트를 이용한 후발적 조치만 하였다.
- 좀 더 빠르게 사기꾼을 색출하여 데이터 낭비를 방지하고 싶다.

01. Competition Overview

Description



netmarble

01. Competition Overview

Evaluation

Submissions are evaluated on area under the ROC curve between the predicted probability and the observed target.

Submission File

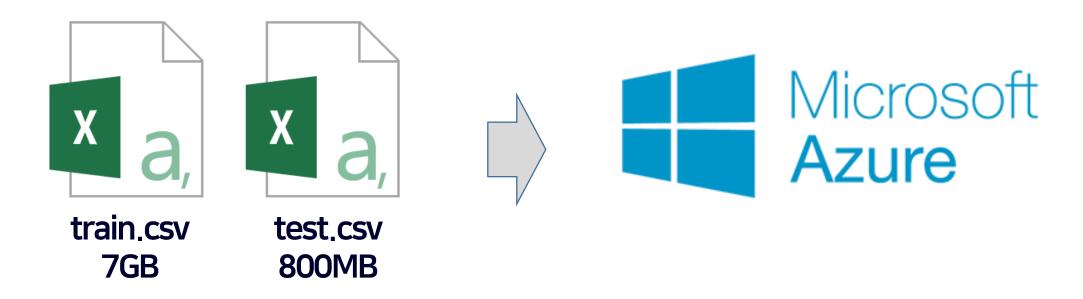
For each |click_id in the test set, you must predict a probability for the target | is_attributed variable. The file should contain a header and have the following format:

```
click_id,is_attributed
1,0.003
2,0.001
3,0.000
etc,
```

02

Analysis environment

1. 대용량 데이터 작업환경



대용량 데이터! 개인 노트북으로 분석 불가 Azure 가상머신을 사용해보자!

- 1. 대용량 데이터 작업환경
- Q) 어떤 가상머신을 구현해야 할까?
 - 요구사항
 - 1. 큰데이터를 처리할 작업공간
 - 2. 작업공간의 공유
 - 3. R언어 사용의 편리성

• Rstudio Server 웹에서 가상머신 IP주소로 간편하게 Rstudio에 접근 가능 Linux환경필요 (Windows 미지원)







Linux 필요

Linux 기반 OS "Ubuntu" 선택

- 2. Azure로 가상머신 구현하기
- @) Azure에서 Ubuntu 구현하기
 - 1. Linux에 접근하기 위해 PuTTY 라는 프로그램을 설치
 - 2. PuTTYgen 실행후키생성
 - 3. Azure에서 Ubuntu 선택 후 SSH 공개 키에 키 입력

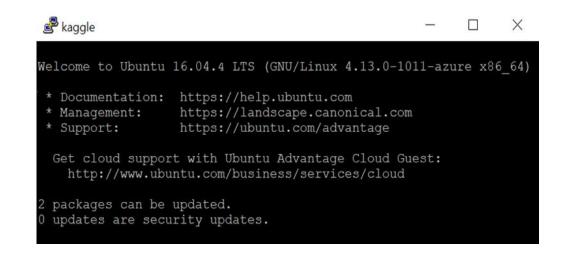
···완성!



PuTTygen SSH Key



- 2. Azure로 가상머신 구현하기
- @) Ubuntu에 R 환경 구현하기
 - 1. PuTTy를 실행하여 가상머신 IP 로 Ubuntu 실행
 - 2. Linux 명령어로 R과 Rstudio Server설치
 - 3. Azure에서 Rstudio Server의 포트 8787 열어주기
 - 4. 인터넷 주소창에 가상머신 IP:8787로 R server 접근 후 Linux ID로 서버 로그인





3. Ubuntu에 R작업환경 구축하기

커널의 여러 기법들을 사용하려 하였으나, 우분투의 기본 지원 R버전이 낮아서 문제가 발생 -> 우분투의 Sources.list에 R repository 추가해서 R버전 업데이트

You need to add R's repository to your system:

- 1. Use your favorite text editor (I'm using gedit as an example) to open /etc/apt/sources.list:
 sudo -H gedit /etc/apt/sources.list
- 2. Add this line to the file (if this is slow, use another mirror. You may also want to change precise into the codename for your Ubuntu version --- e.g., trusty for 14.04):

deb http://cran.rstudio.com/bin/linux/ubuntu precise/



추가적으로, 패키지 설치를 위하여 Git과 Cmake를 설치 -> LightGBM 등의 모델링 패키지 설치가능

```
library(devtools)
options(devtools.install.args = "--no-multiarch") # if yo
install_github("Microsoft/LightGBM", subdir = "R-package"
```

- 3. Ubuntu에 R작업환경 구축하기
- ...였지만 끝난 것이 아니었다. 문제발생!
- 문제1) 가상머신이 체험판이라 7G데이터를 분석하기에는 메모리가 부족하다..
- 방안)
- 1. data.frame이 아닌 data.table의 fread()로 데이터 불러오기.
- 2. 필요 없는 변수 없애기. (attributed_time)
- 3. 전체 데이터를 작은 데이터로 줄여서 사용하기.
- 4. 코드 중간중간에 가비지 컬렉션 gc() 자주 사용해서 메모리 관리하기.
- 문제2) EDA 시각화에 시간이 오래 걸린다.
- 커널 EDA를 참고하고 train_sample데이터로 확인해보기
- => 비록 모델링 자체에 1억 8천건의 데이터를 전부 사용하지는 못하고 있지만 데이터를 불러오고 모델링 할 수 있는 환경 구축 성공!

03

EDA

Data fields

- lp : 클릭한 ip 주소
- App : 마케팅용앱 id
- Device : 사용자의 핸드폰 기종 id
- Os: 사용자 전화의 os 버전 id

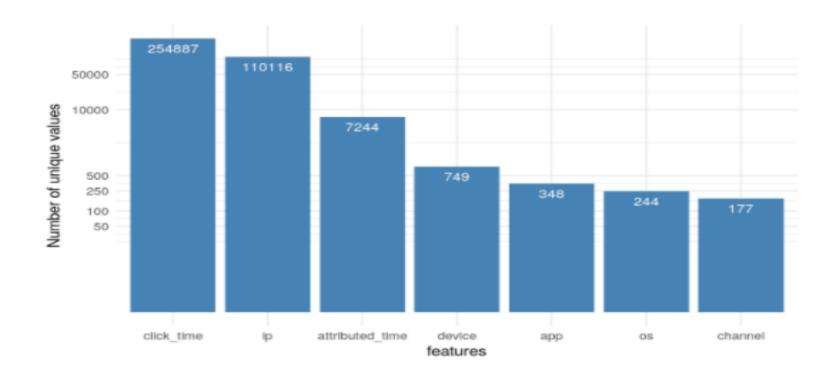
- Channel: 모바일 광고 게시자의 채널 id
- Click_time : 클릭시각
- Attributed_time : 앱 다운로드 시간
- ls_attributed : 1 :다운로드 / 0 : 사기

Data fields

1억 8천 4백만……

- 관측치 수: 1억 8천4백9십만 3천8백9십 ->데이터가 많아서 학습에 오랜 시간이 걸린다. 혹은 불가능하다. (샘플링)
- 변수 class : 비 식별화 되어있는 변수들. ->시간빼고모두 범주형 변수라고 보는 것이 적절하다.

Unique(train)



가장 많은 unique 변수를 가지고 있는 ip가 11만개 (샘플 300만개 중에서)
 -> 상당히 많은 중복이 발생하였고, 각 변수의 분포를 알아볼 필요가 있다.

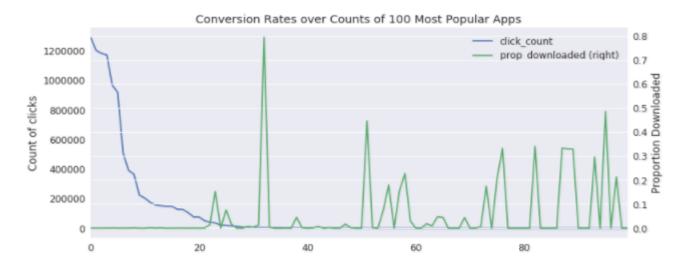
Count feature를 이용한 분포 그래프



1차자료와 2차자료

Conversions by App

Check 100 most popular apps by click count:



2017년 중국 스마트폰 시장 점유율 (단위:%)

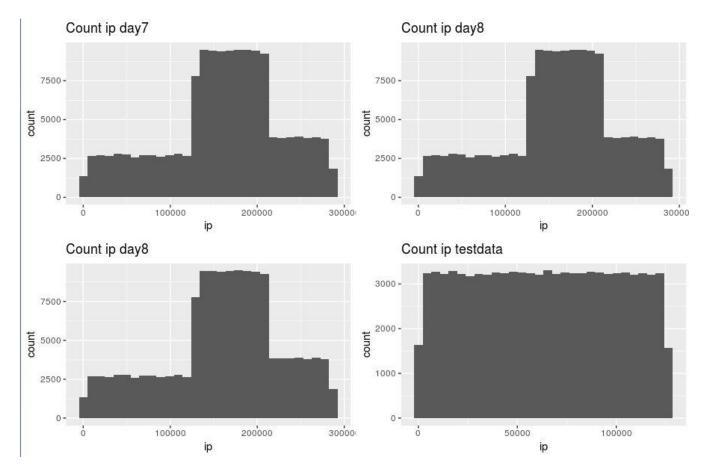
스마트폰 제조사	1분기	2분기	3분기	4분기
화웨이	18.9	21	19	19.6
오포	18.7	19.1	18.8	18.4
비보	16.8	14.6	15.2	16.6
샤오미	7.1	12.9	12.8	12.5
애플	7.7	5.2	7.2	13.3
지오니	5.5	5.3	4.9	3.6
메이쥬	4.5	4	3.4	2.4
삼성전자	3.1	2.7	2	1.7
ZTE	3.1	2.7	2	1.5
시아오라지아오	1.9	1.8	1.5	1.1
기타	12.6	10.6	13.2	9.2

App

Device

 특정 집단 혹은 개인이 비정상적인 데이터를 만들고 있고, 이들은 일정한 Device를 사용하고 특정 App의 광고 수를 인위적으로 늘리고 있다.

About ip count feature



• ip가 고르게 분포한 형태로 Test data가 변경됨 -> ip count feature의 사용을 재고해야 한다.

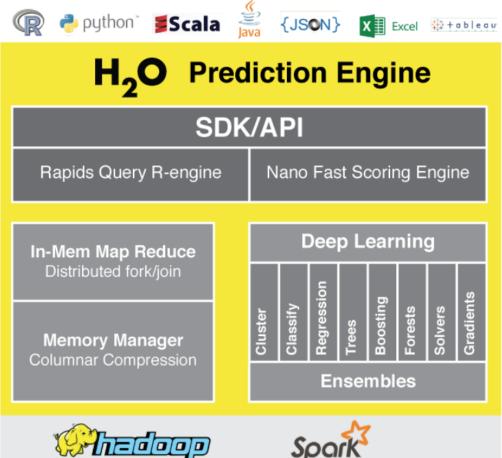
04

H2O Package

H2O패키지소개

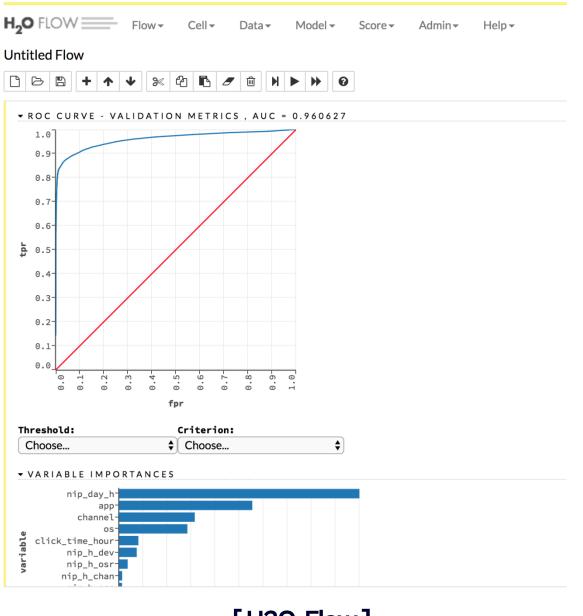
빅데이터 분석을 위한 JAVA기반의 오픈소스 머신러닝/AI 플랫폼

H2O is an open source, in-memory, distributed, fast, and scalable machine learning and predictive analytics platform that allows you to build machine learning models on big data and provides easy productionalization of those models in an enterprise environment. [http://docs.h2o.ai/h2o/latest-stable/h2o-docs/welcome.html]



H2O패키지의 장점

- **빠른속도 (**분산 병렬 처리)
- 웹 인터페이스 제공 (H2O Flow)
- 확장성이좋다.
- 모델 배포가 편리
- 사용하기 쉬움



[H2O Flow]

H2O패키지사용하기

```
# h2o 클러스터를 R과 연결
h2o.init(nthreads = 4, max_mem_size = "18g")

# h2o 상에 데이터 올리기
training <- as.h2o(training, destination_frame = "training")
validation <- as.h2o(validation, destination_frame = "validation")
test_h2o <- as.h2o(test, destination_frame = "test_h2o")
h2o.ls() # 올라갔는지 확인

target <- "is_attributed"
features <- names(training)[!names(training) %in% target]
```
```

h2o.init() : h2o클러스터와 R 연결

(thread 갯수, 최대사용메모리 지정)

as.h2o() : 데이터를 h2o 상에 올림

(또는, h2o.importFile() 로 파일을 바로 올릴 수 있음)

\_\_\_\_

#### key

```
depth_grid_model_2
depth_grid_model_3
depth_grid_model_4
depth_grid_model_5
gbm_model
glm_model
modelmetrics_GBM_model_R_1521834389751_18@8903738402016773956_on_training@-45614059431649
```

#### H2O패키지사용하기

H2OBinomialMetrics: gbm

MSF · 0 001736085

\*\* Reported on training data. \*\*

h2o.gbm(): Gradient Boosting model

```
GBM 모델 생성
gbm_model <- h2o.gbm(x = features, y = target, training_frame = training, model_id = "gbm_model")
```{r}
                                                                                              £ ₹
summary(gbm_model)
                                                          gbm2 # AUC: 0.9612291
 Model Details:
                                                          preds <- h2o.predict(gbm2, test_h2o)</pre>
                                                          head(preds)
                                                           summary(preds$predict, exact_quantiles=TRUE)
 H2OBinomialModel: gbm
 Model Key: gbm_model
 Model Summary:
   number_of_trees number_of_internal_trees model_size_ir
                                                                    predict
                                                                    <fctr>
                                                                                                             <dbl>
   min_depth max_depth mean_depth min_leaves max_leaves
                                                                                                        0.9994166
                                                                    0
                          5.00000
                                           28
                                                                                                        0.9993592
                                                                    0
   mean_leaves
      31.62000
                                                                                                        0.9996849
                                                                    0
```

h2o.predict() 로예측

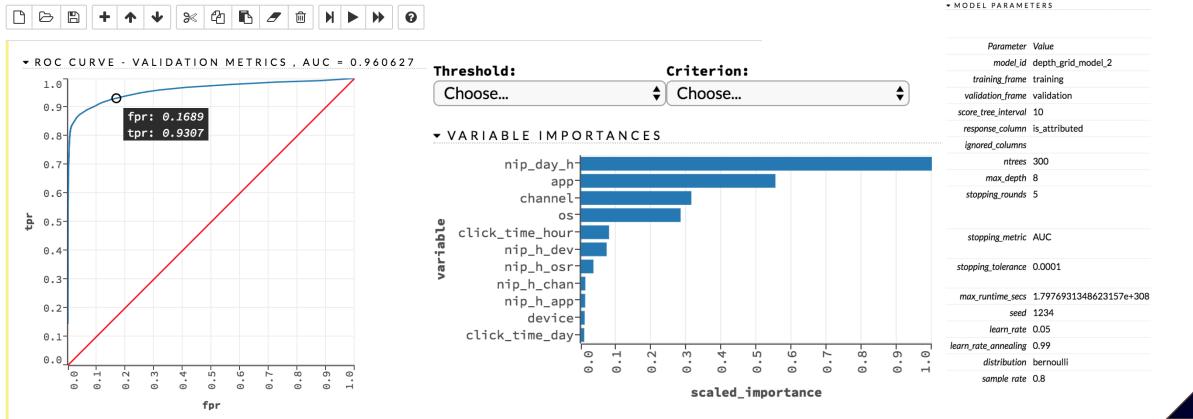
<dbl> 0.0005833523 0.0006407654 0.0003150630 0.0003288168 0.9996712 0 0.0003086371 0 0.9996914 0.0005730025 0.9994270 0

H2O패키지사용하기

h2o flow 에서 모델정보 확인 (ROC Curve, 변수중요도, 모델파라미터등)



Untitled Flow



Q&A

Thank you