

New entrants are currently prohibited. You must wait until after the deadline (6.7 days) to download the data or submit.



\$25,000 • 959 teams

## Outbrain Click Prediction

Merger and Entry Deadline

Wed 5 Oct 2016

Wed 18 Jan 2017 (6.7 days to go)

### Dashboard

Home

Data

Make a submission

### Information

Description

Evaluation

Rules

Prizes

Timeline

### Forum

### Kernels

New Script

New Notebook

### Leaderboard

### My Submissions

### Public Leaderboard

1. Three Data Points
2. code monkey
3. brain-afk
4. Neuron
5. FG Knight
6. Andrii Cherednychenko
7. CV
8. Sangxia
9. rokh
10. Igor Pasechnik

### 1,451 Kernels

Evaluate map score without groupby in python  
14 Votes / 16 days ago / Python

Outbrain EDA  
191 Votes / 3 months ago / Python

[Competition Details](#) » [Get the Data](#) » [Make a submission](#)

### Data Files

File Name	Available Formats
documents_categories.csv	<a href="#">.zip (32.34 mb)</a>
clicks_test.csv	<a href="#">.zip (135.43 mb)</a>
documents_meta.csv	<a href="#">.zip (15.51 mb)</a>
documents_entities.csv	<a href="#">.zip (125.67 mb)</a>
promoted_content.csv	<a href="#">.zip (2.52 mb)</a>
sample_submission.csv	<a href="#">.zip (99.57 mb)</a>
documents_topics.csv	<a href="#">.zip (120.91 mb)</a>
clicks_train.csv	<a href="#">.zip (389.75 mb)</a>
events.csv	<a href="#">.zip (477.74 mb)</a>
page_views.csv	<a href="#">.zip (29.71 gb)</a>
page_views_sample.csv	<a href="#">.zip (148.51 mb)</a>

The dataset for this challenge contains a sample of users' page views and clicks, as observed on multiple publisher sites in the United States between 14-June-2016 and 28-June-2016. Each viewed page or clicked recommendation is further accompanied by some semantic attributes of those documents. For full details, see data specifications below.

The dataset contains numerous sets of content recommendations served to a specific user in a specific context. Each context (i.e. a set of recommendations) is given a display\_id. In each such set, the user has clicked on at least one recommendation. The identities of the clicked recommendations in the test set are not revealed. Your task is to rank the recommendations in each group by decreasing predicted likelihood of being clicked.

Unveiling page\_views.csv with PySpark  
87 Votes / 2 months ago / Python

pypy implementation of fm(with adam)  
6 Votes / 5 days ago / Python

FTRL Starter (with leakage vars)  
33 Votes / 2 months ago / Python

pandas is cool! LB: 0.63714  
25 Votes / 2 months ago / Python

## Forum (116 topics)

FFM Input  
4 minutes ago

Best single model performance  
14 minutes ago

let's break 0.7  
3 hours ago

Why pandas read both strings and integers value for platforms in events.csv  
5 hours ago

Future publish\_time in documents\_meta.csv file  
8 hours ago

Benchmark 0.65251 using BTB with 95GB page\_views.csv  
12 hours ago

teams

players

entries

As a warning, this is a very large relational dataset. While most of the tables are small enough to fit in memory, the page views log (page\_views.csv) is over 2 billion rows and 100GB uncompressed. We have also uploaded a sample version of this file with the first 10,000,000 rows. The MD5 checksum of page\_views.csv.zip is 3742c116bab4030e0a7ea1c0be623bd9.

## Data Fields

Each user in the dataset is represented by a unique id (uuid). A person can view a document (document\_id), which is simply a web page with content (e.g. a news article). On each document, a set of ads (ad\_id) are displayed. Each ad belongs to a campaign (campaign\_id) run by an advertiser (advertiser\_id). You are also provided metadata about the document, such as which entities are mentioned, a taxonomy of categories, the topics mentioned, and the publisher.

## File Descriptions

**page\_views.csv** is a the log of users visiting documents. To save disk space, the timestamps in the entire dataset are relative to the first time in the dataset. If you wish to recover the actual epoch time of the visit, add 1465876799998 to the timestamp.

- uuid
- document\_id
- timestamp (ms since 1970-01-01 - 1465876799998)
- platform (desktop = 1, mobile = 2, tablet =3)
- geo\_location (country>state>DMA)
- traffic\_source (internal = 1, search = 2, social = 3)

**clicks\_train.csv** is the training set, showing which of a set of ads was clicked.

- display\_id
- ad\_id
- clicked (1 if clicked, 0 otherwise)

**clicks\_test.csv** is the same as clicks\_train.csv, except it does not have the clicked ad. This is the file you should use to predict. Each display\_id has only one clicked ad. Note that test set contains display\_ids from the entire dataset timeframe. Additionally, the public/private sampling for the competition is uniformly random, not based on time. These sampling choices were intentional, in spite of the possibility that participants can look ahead in time.

**sample\_submission.csv** shows the correct submission format.

**events.csv** provides information on the display\_id context. It covers both the train and test set.

- display\_id
- uuid
- document\_id
- timestamp
- platform
- geo\_location

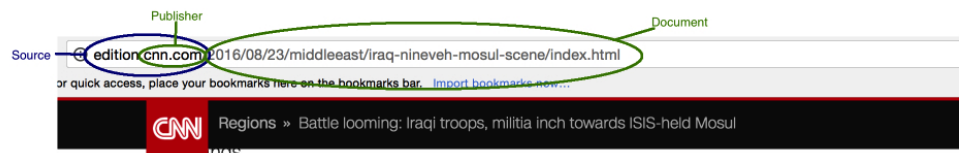
**promoted\_content.csv** provides details on the ads.

- ad\_id
- document\_id
- campaign\_id
- advertiser\_id

**documents\_meta.csv** provides details on the documents.

- document\_id
- source\_id (the part of the site on which the document is displayed, e.g. edition.cnn.com)
- publisher\_id
- publish\_time

**documents\_topics.csv**, **documents\_entities.csv**, and **documents\_categories.csv** all provide information about the content in a document, as well as Outbrain's confidence in each respective relationship. For example, an **entity\_id** can represent a person, organization, or location. The rows in **documents\_entities.csv** give the confidence that the given entity was referred to in the document.









Promoted  
Content  
Set

"I am so happy for them," the man said. "But I am heartbroken myself. My parents were not able to come with me. I don't know how I am going to get them out."



Promoted  
Content  
Item

**Paid Content** Recommended by **Outbrain**

 Mapping the Startup Nation: The 12 most popular Tech Hubs in... Viola Notes	 First time in Israel: Business degrees in Ramat Gan and New... Israel News	 The most addictive game of the year! Play with 15 million Players... Forge Of Empires
 How to Avoid Everyday Pain Landmines Womens Health	 How One Brand is Disrupting the \$63 Billion Makeup Industry The Huffington Post	 Find out what special ingredient makes this omelette so tasty HomeMadebyYou

## Privacy Reminder

Outbrain is releasing 2 Billion page views and 16,900,000 clicks of 700 Million unique users, across 560 sites. The data is anonymized. Please remember that participants are prohibited from de-anonymizing or reverse engineering data or combining the data with other publicly available information. Outbrain does not collect or hold PII (personally identifiable information), and the user identifiers we are releasing here are obscured. To protect its publisher partners, Outbrain is not releasing URLs of viewed or

clicked stories, but rather anonymized document and site identifiers. The task at hand is click prediction, and by downloading the dataset, participants agree to use the data for that task alone, and will not attempt to reverse engineer the mapping from document, site, and user identifiers to URLs, site names or actual users.