

Data Science Game Qualification phase: Music recommendation.

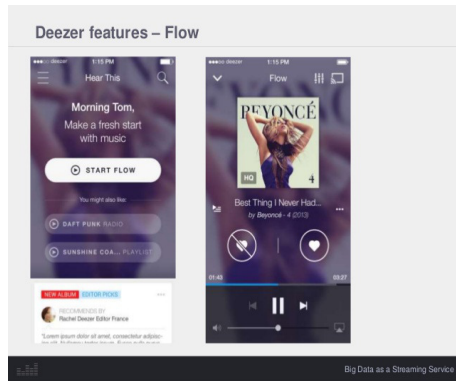
Vikulin Vsevolod Soboleva Daria Popov Nikolay Shapovalov Nikita

Lomonosov Moscow State University

10 June 2017

Goal

The goal of this challenge is to predict whether the users of the test dataset listened to the first track Flow proposed them or not.



Data

user_id	media_id	f1	f2	...	fn	y
0	211678					1
0	211678					1
0	238286					1
0	615655					0
1	211277					1
...

f1 ... fn – information about user and media.

User-specific features

- user_id – anonymized id of the user
- user_gender – gender of the user
- user_age – age of the user

Media-specific features

- `media_id` – identifier of the song listened by the user
- `media_duration` – duration of the song
- `context_type` – type of content where the song was listened: playlist, album etc.
- `release_date` – release date of the song with the format YYYYMMDD
- `artist_id` – identifier of the artist of the song
- `genre_id` – identifier of the genre of the song

Platform-specific features

- ts_listen – timestamp of the listening in UNIX time
- platform_name – type of os
- platform_family – type of device
- listen_type – if the songs was listened in a flow or not

Artist-specific features

- id – the artist's Deezer id
- name – the artist's name
- radio – true if the artist has a smartradio
- nb_fan – the number of artist's fans

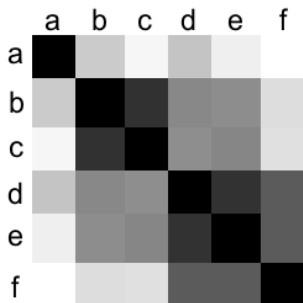
Album-specific features

- id – the Deezer album id
- genre_id – the album's first genre id
- title – the track's fulltitle
- nb_tracks – the number of album's tracks
- rating – the album's rate
- duration – the track's duration in seconds
- fans – the number of album's Fans

Distance matrix features

We create the cosine distance matrix between users.

$$distance(user_1, user_2) = 1 - CosSim = 1 - \frac{(user_1, user_2)}{||user_1|| * ||user_2||}$$



Cluster matrix and get the cluster id.

SVD features

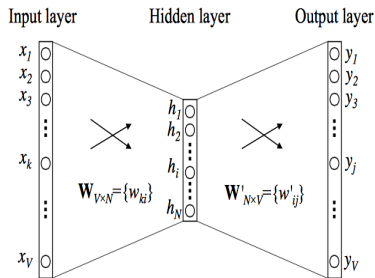
We use SVD for matrix UserMedia.

The diagram illustrates the SVD decomposition of matrix A . Matrix A is represented by a tall vertical rectangle. It is equal to the product of three matrices: U , L , and V^T . Matrix U is a tall vertical rectangle with a shaded vertical strip on its left side. Matrix L is a square with a shaded top-left corner. Matrix V^T is a square with a shaded top row. The equation is shown as $A = U L V^T$.

Cluster users-embedding matrix and get the cluster id or calculate mean of vector for each user.

User2Vec features

We use Word2Vec architecture to predict song for each user.



Get cluster id or calculate mean of user vector for each user.

Solution. Target statistic.

- $K = data.groupby(f).size()$
- $mean_y = data.groupby(f)['y'].mean()$
- $global_mean_y = data['y'].mean()$
- $\frac{mean_y * K + global_mean_y * 10}{K + 10}$

Solution. Real features statistic.

`data.groupby(f)[real_f].agg([func])`

$func \in \text{min, max, median, etc.}$

Solution. Bug for luck.

1st {

user_id	media_id	ts_listen	f1	f2	...	fn	y
0	122437314	$1477 * 10^6$					1
0	131340580	$1478 * 10^6$					0

2nd {

user_id	media_id	ts_listen	f1	f2	...	fn	y
0	131576046	$1479 * 10^6$					1
0	129011934	$1480 * 10^6$					1

$$(1st + 2nd)/2$$

→ Top1!!!