



ACM RecSys Challenge 2016

Решение команды Avito (7 место)

Василий Лексин, Андрей Остапец

Avito.ru

9 июля 2016 г.



- **Users** — данные по пользователям сайта XING.com: `jobroles`, `career_level`, `discipline`, `industry`, `location`, `experience`, `education`.
- **Items** — данные по вакансиям: `title`, `career_level`, `discipline`, `industry`, `location`, `employment_type`, `tags`, `created_at`, `active_during_test`.
- **Impressions** — какие вакансии и когда рекомендовались пользователю.
- **Interactions** — по каким вакансиям и когда пользователь совершал действия: `clicked`, `bookmarked`, `replied`, `deleted`.



Интервал дат: 2015-08-19 – 2015-11-09

### Impressions

- 201M уникальных user-item-week сочетаний.
- 2.7M уникальных пользователей.
- 846K уникальных вакансий.

### Interactions

- 8.8M events: **clicked** – 7.2M, **deleted** – 1.0M, **replied** – 422K, **bookmarked** – 206K.
- 785K уникальных пользователей.
- 1.03M уникальных вакансий.
- 2.8M из 6.9M (40%) пар user-item содержатся в impressions.



150K пользователей, которым нужно порекомендовать, из них:

- у 39.7K (26.5%) нет событий;
- у 59.5K (39.6%) менее 1 события;
- у 70.6K (47.1%) менее 2-х событий.

327K активных вакансий, из них:

- у 129K (39.5%) нет событий;
- у 164K (50.1%) менее 1 события;
- у 188K (57.6%) менее 2-х событий.



Для каждого из 150K целевых пользователей отобрать по 30 вакансий, на которые пользователь наиболее вероятно совершил одно из действий: clicked, bookmarked or replied в следующую после 2015-11-09 неделю.

$$score = \sum_{all\_users} 20(P2 + P4 + R + S) + 10(P6 + P20),$$

где  $P_k$  – precision at  $k$ ;  $R$  – recall;  $S$  – user success (хотя бы одно положительное действие на любой из 30 вакансий).

# Решение команды

## Пример пользовательской сессии



CV					
experience_n	experience	discipline_id_user	country_user	region_user	jobroles
5 or more entr	10-15 year	Sales & Commerce	Germany	Bavaria	['962959', '283291', '502342']
SESSIONS					
created_at	impression	discipline_id_item	country_item	region_item	tags
09-01 1:27	1	Production & Manufact	Germany	Baden-Württer	['620383', '1118975']
09-01 1:27	1	Other Disciplines	Germany	Hamburg	['4572761', '3543754', '19689']
09-01 1:27	1	Other Disciplines	Germany	Berlin	['18091']
09-02 0:21	1	Health, Medical & Social	non_dach	not specified	['165415', '1986087', '258579']
09-04 20:46	0	IT & Software Development	Germany	Brandenburg	['655030']
09-08 20:16	1	Other Disciplines	Germany	Hamburg	['2915824', '4035399', '15676']
09-08 22:40	1	Sales & Commerce	Germany	not specified	['3418410', '3413328']
09-08 22:41	1	IT & Software Development	non_dach	not specified	['3408137']
09-09 22:58	1	Administration	Germany	Berlin	['4141254', '1118975']
09-09 23:00	0	Production & Manufact	Germany	Lower Saxony	['4454260', '502342']
09-09 23:01	0	Other Disciplines	Germany	Lower Saxony	['1567693', '568776']
09-09 23:02	1	Health, Medical & Social	non_dach	not specified	['1567693']
09-09 23:08	0	Finance, Accounting & Tax	Austria	not specified	['2865345', '3294368']
09-09 23:08	0	Production & Manufact	Germany	Lower Saxony	['494116']
09-09 23:08	0	Health, Medical & Social	Switzerland	not specified	['128836', '1836819']
09-09 23:09	0	Production & Manufact	Austria	not specified	['2846960', '76751', '4227194']
09-09 23:09	0	Engineering & Technology	Germany	Berlin	['4141254']
09-10 0:59	1	Production & Manufact	Austria	not specified	['1118975', '3478136']
09-10 0:59	1	Engineering & Technology	Germany	Bavaria	['128836', '76887']
09-10 0:59	1	Engineering & Technology	Germany	North Rhine-W	['119117', '3705605', '34781']
09-10 1:00	0	Other Disciplines	Austria	not specified	['2915824', '4035399', '15676']
09-10 1:00	1	Teaching, R&D	Germany	North Rhine-W	['1986087']
09-10 1:00	0	Other Disciplines	Germany	not specified	['2573697', '4035399', '44843']
09-10 1:00	0	Engineering & Technology	Germany	Baden-Württer	['2140778', '3241763']
09-10 1:00	1	Management & Corporate	Germany	Bavaria	['494116', '1119117', '238737']
09-10 1:00	1	Teaching, R&D	Germany	Bremen	['494116', '3376322']
09-10 1:01	0	Production & Manufact	Germany	Lower Saxony	['4454260', '502342']
09-10 1:01	0	Other Disciplines	Germany	North Rhine-W	['494116', '3408137']
09-10 1:02	0	Production & Manufact	Germany	Bavaria	['502342']
09-10 1:02	0	Production & Manufact	Germany	Schleswig-Hols	['398644', '2865345']
09-10 1:02	1	Administration	Germany	Berlin	['4141254', '1118975']
09-10 1:02	0	Production & Manufact	Germany	Bavaria	['4013443']
10-16 22:33	1	Other Disciplines	Switzerland	not specified	['1986087', '1971244']



- Значительная доля пользователей и вакансий с малым количеством событий или без событий, значит нужен гибридный подход, учитывающий контент.
- Impressions слабо меняются во времени, 40% действий пользователя совершаются из impressions, то есть наличие пары user-item в impressions хороший признак.
- Географические признаки почему-то не работают.



- Выделяем последнюю полную неделю из interactions в контроль.
- Выделяем случайных 10К пользователей из тех, что совершал действия в эту неделю.
- Отбрасываем старые вакансии (created\_at более месяца), которые никто не кликал.
- Хорошая корреляция с лидербордом.





Метрики сходства:

- Jaccard
- Cosine
- Pearson

Типы событий, на которых обучаемся:

- Interactions без deleted
- Только clicked
- Только impressions

# Решение команды

## Factorization Machines: общая идея



Feature vector $x$																			Target $y$		
$x^{(1)}$	1	0	0	...	1	0	0	0	...	0.3	0.3	0.3	0	...	13	0	0	0	0	...	5 $y^{(1)}$
$x^{(2)}$	1	0	0	...	0	1	0	0	...	0.3	0.3	0.3	0	...	14	1	0	0	0	...	3 $y^{(2)}$
$x^{(3)}$	1	0	0	...	0	0	1	0	...	0.3	0.3	0.3	0	...	16	0	1	0	0	...	1 $y^{(2)}$
$x^{(4)}$	0	1	0	...	0	0	1	0	...	0	0	0.5	0.5	...	5	0	0	0	0	...	4 $y^{(3)}$
$x^{(5)}$	0	1	0	...	0	0	0	1	...	0	0	0.5	0.5	...	8	0	0	1	0	...	5 $y^{(4)}$
$x^{(6)}$	0	0	1	...	1	0	0	0	...	0.5	0	0.5	0	...	9	0	0	0	0	...	1 $y^{(5)}$
$x^{(7)}$	0	0	1	...	0	0	1	0	...	0.5	0	0.5	0	...	12	1	0	0	0	...	5 $y^{(6)}$
	A	B	C	...	TI	NH	SW	ST	...	TI	NH	SW	ST	...	Time	TI	NH	SW	ST	...	
	User				item Movie					Other Movies rated						Last Movie rated					



### **Users** - все признаки (OneHotEncoder)

- jobroles
- career\_level, discipline\_id, industry\_id
- country, region
- experience: n\_entries\_class, years, years\_in\_current
- edu: degree, fieldofstudies

### **Items** - все признаки, кроме latitude и longitude (OneHotEncoder)

- title, tags
- career\_level, discipline\_id, industry\_id
- country, region
- employment\_type



- Количество латентных факторов
- Количество случайных негативных примеров на каждого пользователя
- Количество итераций
- Параметры регуляризации



- Документ = пользователь (все токены вакансий, просмотренных пользователем + токены из CV) или вакансия (все токены вакансии).
- Строим TFIDF на токенах пользователей.
- Обучаем Latent Semantic Indexing модель (по сути SVD).
- Сходство пользователя и вакансии аналогично сходству двух документов.



SIM_jac	Item-based jaccard similarity
SIM_click	Item-based jaccard similarity on clicks
SIM_pearson	Item-based pearson similarity
SIM_imp	Item-based jaccard similarity on impressions
FR_f100_i25	Factorization, n_factors=100, iter=25
FR_f400_i70	Factorization, n_factors=400, iter=70
FR_f400_i50_no_side	Factorization, no side data
FR_imp	Factorization on impressions
TM	LSI topic model



Базовые модели:

$FR_0$	$SIM_0$	impression_prev	local score
1	0	0	76995
0	1	0	69622
0	0	1	104495
1	1	1	132505

- $SIM_0$  – item-based jaccard similarity
- $FR_0$  – factorization, 400 factors
- impression\_prev – был ли impression этой вакансии пользователю



«Нулевая» версия:

$FR_0$	$SIM_0$	$impression\_prev$	$local\ score$
1	2	1	134285

Первая версия:

$FR_0$	$SIM_0$	$FR_0^{8.0} * SIM_0$	$impression\_prev$	$local\ score$
1	13	8	1	138073

Вторая версия:

$FR_1$	$SIM_1$	$FR_1^{8.0} * SIM_1$	$impression\_prev$	$local\ score$
1	13	8	1	140876

$$FR_1 = FR\_f100\_i25$$

$$SIM_1 = SIM\_jac$$





Третья версия:

$FR_2$	$SIM_2$	$FR_2^{8.0} * SIM_2$	<i>impression_prev</i>	<i>local score</i>
1	13	8	1	143653

$$FR_2 = 0.5 * FR\_f100\_i25 + 0.5 * FR\_f400\_i70$$

$$SIM_2 = 0.5 * SIM\_jac + 0.5 * SIM\_click$$



### local score (145841):

$1.0 * FR_3 + 15.0 * (FR_3^{8.0} * SIM_3) + 13.0 * SIM_3 + 1.0 * impression\_prev - 0.5 * SIM\_pearson - 0.3 * FR\_f400\_i50\_no\_side + 0.5 * (FR\_imp^{2.0} * SIM\_imp)$ , где

$$FR_3 = 0.5 * FR\_f100\_i25 + 0.5 * FR\_f400\_i70$$

$$SIM_3 = SIM\_click$$

### local score (146569):

$1.0 * FR_4 + 15.0 * (FR_4^{8.0} * SIM_4) + 13.0 * SIM_4 + 1.0 * impression\_prev - 0.4 * SIM\_pearson - 0.3 * FR\_f400\_i50\_no\_side + 0.5 * (FR\_imp^{2.0} * SIM\_imp) + 0.2 * TM$ , где

$$FR_4 = 0.5 * FR\_f100\_i25 + 0.5 * FR\_f400\_i70$$

$$SIM_4 = 0.4 * SIM\_jac + 0.6 * SIM\_t = 1$$



- 1 сервер 56 ядер 256Gb
- Полное обучение + скоринг = 16 часов
- Библиотеки: graphlab, gensim



Score	Rank	Label	Time
554655	9	Topic model 100 factors	06/27/16
548366	8	Top 150 candidates from every model	06/25/16
543284	8	8 model set: 4FR + 4SIM + TM	06/24/16
537157	9	Topic model	06/23/16
530599	10	3 models set: FR + 2SIM + tuned coef.	06/23/16
497136	15	3 models set: FR + 2SIM	06/22/16
496241	1	FR with side data	03/20/16
397604	1	Impression prev feature	03/11/16
132790	1	Simple item-similarity recommender	03/10/16

# Leaderboard

Финальный рейтинг



Rank	Team	Leaderboard Score	Full Score
1	YunOS-OneSearch	681707.38	2052185.54
2	mim-solutions	675985.03	2035964.16
3	DaveXster	665592.06	2005263.73
4	PumpkinPie	622408.55	1866477.77
5	milk tea	613125.21	1846420.12
6	mdr_rec	605048.58	1823472.31
<b>7</b>	<b>Avito</b>	<b>554654.72</b>	<b>1677898.52</b>
8	recometric	556133.18	1677233.84
9	nodalpoints	555483.39	1671812.08
10	lucky_dog	542213.51	1632828.82
21	XING_TELECOM	461000.32	1397030.74

# Спасибо!

Приглашаем принять участие:

- Конкурс Avito-2016: Распознавание категории объявления (призовой фонд конкурса - 500 000 рублей)
- Data science meetup в Avito – 13 августа 2016г.