

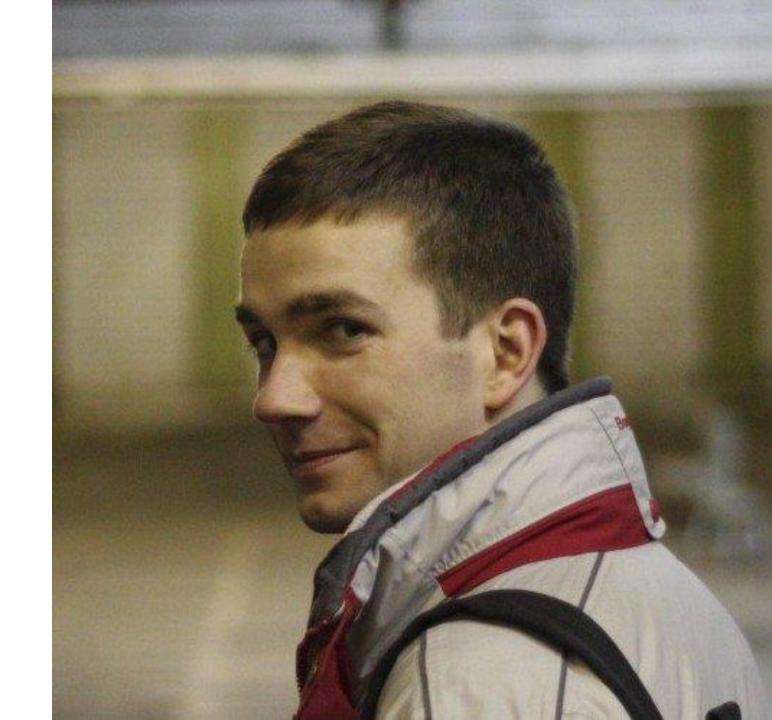
#### ЗАНЯТИЕ 13

## ОБЪЕДИНЕНИЕ ДАТАФРЕЙМОВ

### КОНСТАНТИН БАШЕВОЙ

Яндекс

Habr @kpi\_maker



#### Объединение таблиц это всегда весело

- Чем merge отличается от join и их типы в pandas
- Дубликаты как верный спутник объединений
- Оптимизация хранения данных с помощью join

## PANDAS И БОЛЬШИЕ ФАЙЛЫ

#### Количество уникальных ID

```
user_id
   9
```

Данные отсортированы

#### Количество уникальных ID

```
user id
```

- Читаем файл построчно
- Текущее VS прошлое значение
- Смена значения = пользователь

#### Количество уникальных ID

```
user id
```

- Читаем файл построчно
- Текущее VS прошлое значение
- Смена значения = пользователь
- В памяти 3 числа
- Размер файла не имеет значения

### СКВОЗНАЯ АНАЛИТИКА

#### Склеить лог визитов и лог покупок

	user_id	source
0	11	ad
1	22	yandex
2	55	email
3	11	google
4	77	ad

	user_id	category
0	11	Спорт
1	22	Авто
2	55	Дача
3	11	Дети
4	99	Авто

#### Склеить лог визитов и лог покупок

	user_id	source
0	11	ad
1	22	yandex
2	55	email
3	11	google
4	77	ad

	user_id	source	category
0	22	yandex	Авто
1	55	email	Дача
+ eщe user_id = 11 и 77			

	user_id	category
0	11	Спорт
1	22	Авто
2	55	Дача
3	11	Дети
4	99	Авто

## ПЕРВЫЕ ПРОБЛЕМЫ

### Нет однозначного соответствия

	user_id	source
0	11	ad
1	22	yandex
2	55	email
3	11	google
4	77	ad

	user_id	category
0	11	Спорт
1	22	Авто
2	55	Дача
3	11	Дети
4	99	Авто

#### Сумма визитов и покупок для user\_id

```
visits_grouped = visits.groupby('user_id').count()
visits_grouped.rename(columns={'source': 'visits'}, inplace=True)
visits_grouped
```

#### visits

#### user\_id

11	2
----	---

**22** 1

**55** 1

**77** 1

#### Сумма визитов и покупок для user\_id

category	Авто	Дача	Спорт
user_id			
11	0	0	2
22	1	0	0
55	0	1	0
99	1	0	0

#### Сумма визитов и покупок для user\_id

```
purchases_pivot = purchases.pivot_table(index='user_id', columns='category', values='user_id', aggfunc='size', fill_value=0)

category Abto Дача Спорт

user_id

11 0 0 2

22 1 0 0

55 0 1 0

99 1 0 0
```

C count не работает

## ТИПЫ ОБЪЕДИНЕНИЙ В PANDAS

#### Join - по индексу, merge - по столбцам

	visits
user_id	
11	2
22	1
55	1
77	1

category	Авто	Дача	Спорт
user_id			
11	0	0	2
22	1	0	0
55	0	1	0
99	1	0	0

user\_id сейчас индекс – используем join

#### В таком варианте merge

```
visits.groupby('user_id').count().reset_index()
```

	user_id	source
0	11	2
1	22	1
2	55	1
3	77	1

user\_id сейчас столбец

#### В таком варианте merge



category	user_id	Авто	Дача	Спорт
0	11	0	0	2
1	22	1	0	0
2	55	0	1	0
3	99	1	0	0

user\_id сейчас столбец

#### Все параметры по умолчанию

visits\_grouped.join(purchases\_pivot)

	visits	Авто	Дача	Спорт
user_id				
11	2	0.0	0.0	2.0
22	1	1.0	0.0	0.0
55	1	0.0	1.0	0.0
77	1	NaN	NaN	NaN

## LEFT / RIGHT JOIN



## Каждой строчке левой таблицы ищет соответствие в правой

	visits	category	Авто	Дача	Спорт
user_id		user_id			
11	2	11	0	0	2
22	1	22	1	0	0
55	1	55	0	1	0
77	1	99	1	0	0



### Каждой строчке левой таблицы ищет соответствие в правой





## Каждой строчке левой таблицы ищет соответствие в правой

	visits		category	Авто	Дача	Спорт
user_id			user_id			
11	2	<b>→</b>	11	0	0	2
22	1	<b>→</b>	22	1	0	0
55	1	<b>→</b>	55	0	1	0
77	1		99	1	0	0



## Строчки без пары левой таблицы остаются. Правой – удаляются из результата

	visits		category	Авто	Дача	Спорт			
user_id			user_id						
11	2	<b>→</b>	11	0	0	2			
22	1	<b>→</b>	22	1	0	0			
55	1	<b>→</b>	55	0	1	0			
77	1		99	1	0	0			
NaN									

все строчки левой таблицы останутся

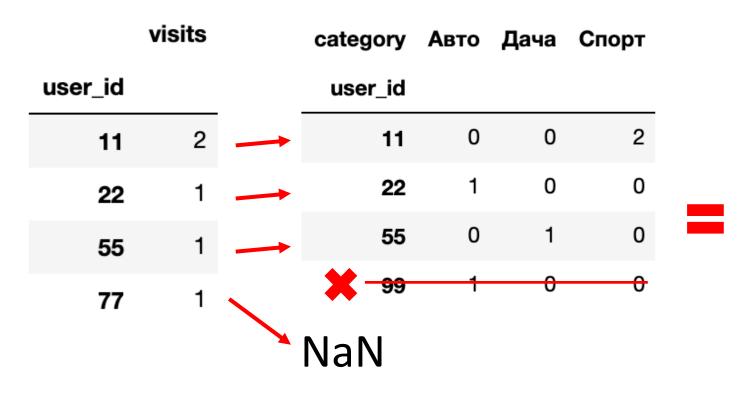
# Строчки без пары левой таблицы остаются. Правой – удаляются из результата

	visits		category	Авто	Дача	Спорт
user_id			user_id			
11	2	<b>→</b>	11	0	0	2
22	1	<b>→</b>	22	1	0	0
55	1	<b>→</b>	55	0	1	0
77	1		<del>× -99</del>	1	0	0
			NaN			

все строчки левой таблицы останутся

правой – не факт

## Строчки без пары левой таблицы остаются. Правой – удаляются из результата



	visits	Авто	Дача	Спорт
user_id				
11	2	0.0	0.0	2.0
22	1	1.0	0.0	0.0
55	1	0.0	1.0	0.0
77	1	NaN	NaN	NaN

#### Right join зеркален left join

# Не рекомендуется к использованию без особой на это необходимости

### Каждой строчке правой таблицы ищет соответствие в левой

	visits		category	Авто	Дача	Спорт	wi ai ta	~~~~!!!	200 -	oin (n	urahaa	og pirrot
user_id	I		user_id				VISIUS	_grou	jea.j		ow='ri	es_pivot, ght')
11	2	-	11	0	0	2		visits	Авто	Дача	Спорт	
22	1	<b>—</b>	22	1	0	0	 user_id					
55	. 1		55	0	1	0	11	2.0	0	0	2	
50	•			_	•	0	22	1.0	1	0	0	
<b>*</b> 77	1		99	1	0	0	55	1.0	0	1	0	
	NaN						99	NaN	1	0	0	

## INNER // OUTER JOIN

#### Оставляет строки, которые есть в обеих таблицах

		visits		category	Авто	Дача	Спорт
user	_id			user_id			
	11	2	<b>→</b>	11	0	0	2
	22	1	<b>→</b>	22	1	0	0
	55	1	<b>→</b>	55	0	1	0
×	77	1		<b>×</b> 99	1	0	0

	visits	Авто	Дача	Спорт
user_id				
11	2	0	0	2
22	1	1	0	0
55	1	0	1	0

#### Оставляет все строки

,	visits	category	Авто	Дача	Спорт	
user_id		user_id				
11	2	11	0	0	2	
22	1	22	1	0	0	
55	1	55	0	1	0	i
77	1 /	99	1	0	0	
N	laN 🔨	NaN				

	visits	Авто	Дача	Спорт
user_id				
11	2.0	0.0	0.0	2.0
22	1.0	1.0	0.0	0.0
55	1.0	0.0	1.0	0.0
77	1.0	NaN	NaN	NaN
99	NaN	1.0	0.0	0.0

## CAMOE BECEЛOE B JOIN

#### Что полезно проверять

(исходя из логики задачи)

- После LEFT-join количество строк не изменилось
- Суммы числовых столбцов не изменились
- Суммы в правой таблице тоже неплохо проверить

## JOIN И ОПТИМИЗАЦИЯ XPAHEНИЯ

#### Как сэкономить место на диске

- Длинные повторяющиеся столбцы переводим в идентификаторы
- Составляем словари
- Логи отдельно, словари отдельно