

Департамент образования и науки города Москвы  
Государственное автономное образовательное учреждение  
высшего образования города Москвы  
«Московский городской педагогический университет»  
Институт цифрового образования  
Департамент информатики управления и технологий

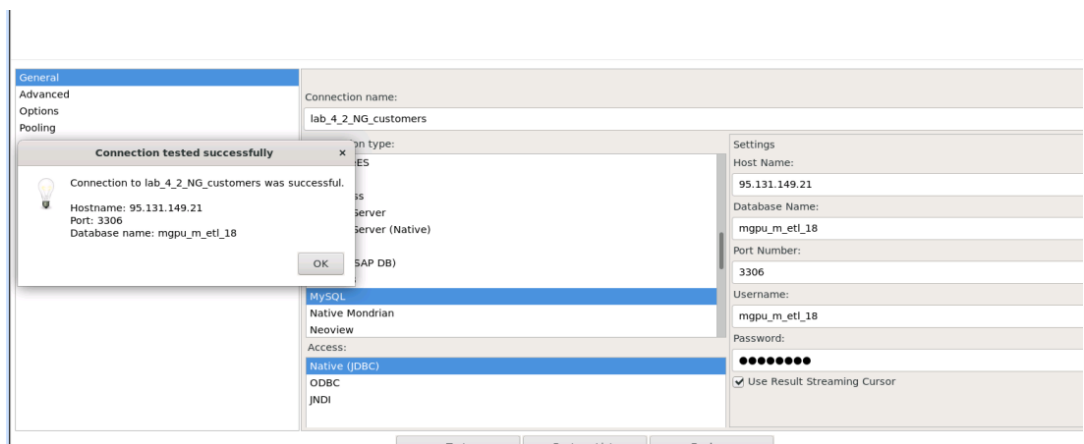
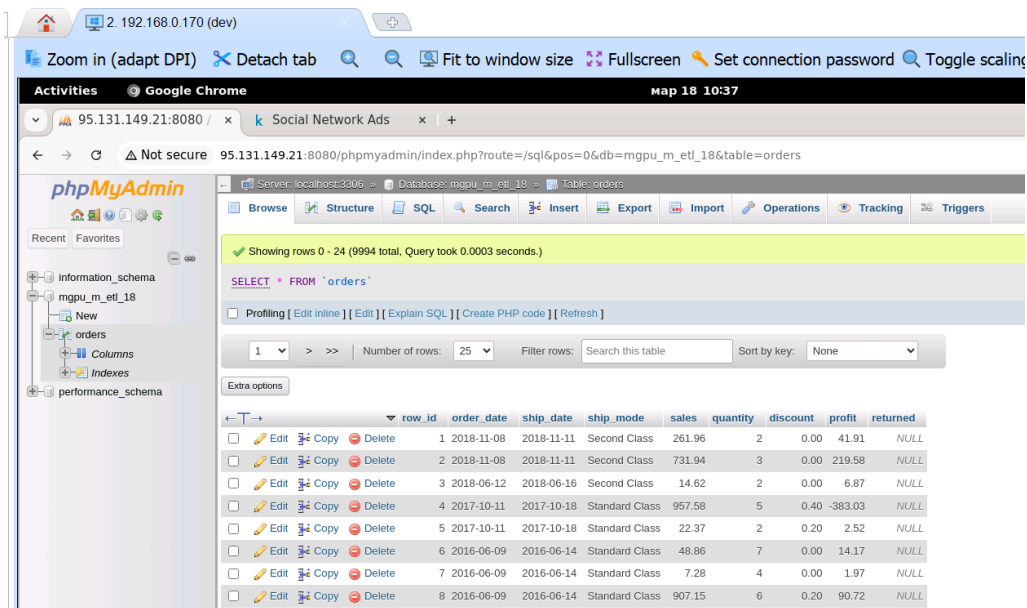
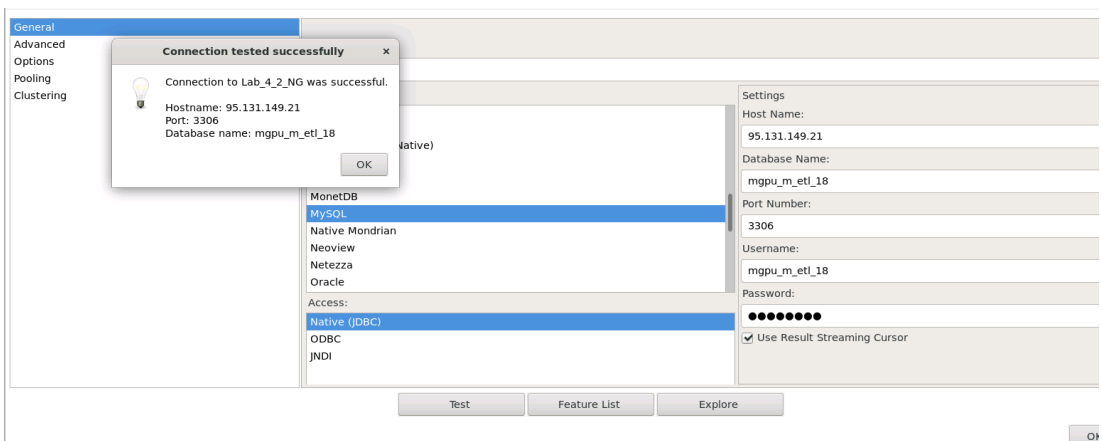
Нургалеева Гузель Рустэмовна БД-241м

**Data Engineering. Практическая работа 4-2.**  
**Вариант задания: 18**

Направление подготовки/специальность  
38.04.05 - Бизнес-информатика  
Бизнес-аналитика и большие данные  
(очная форма обучения)

Москва  
2024

**Цель работы:** получить практические навыки создания ETL-процесса для загрузки данных из CSV-файла в базу данных MySQL с использованием Pentaho Data Integration.



phpMyAdmin interface showing the structure of the 'customers' table in the 'mgpu\_m\_etl\_18' database. The table has 9 columns: id, customer\_id, customer\_name, segment, country, city, state, postal\_code, and region. The 'id' column is the primary key.

#	Name	Type	Collation	Attributes	Null	Default	Comments	Extra	Action
1	id	int			No	None		AUTO_INCREMENT	Change Drop More
2	customer_id	varchar(20)	utf8mb4_0900_ai_ci		No	None			Change Drop More
3	customer_name	varchar(100)	utf8mb4_0900_ai_ci		Yes	NULL			Change Drop More
4	segment	varchar(50)	utf8mb4_0900_ai_ci		Yes	NULL			Change Drop More
5	country	varchar(100)	utf8mb4_0900_ai_ci		Yes	NULL			Change Drop More
6	city	varchar(100)	utf8mb4_0900_ai_ci		Yes	NULL			Change Drop More
7	state	varchar(100)	utf8mb4_0900_ai_ci		Yes	NULL			Change Drop More
8	postal_code	varchar(20)	utf8mb4_0900_ai_ci		Yes	NULL			Change Drop More
9	region	varchar(50)	utf8mb4_0900_ai_ci		Yes	NULL			Change Drop More

phpMyAdmin interface showing the data of the 'customers' table. The table contains 9 rows of customer data.

	id	customer_id	customer_name	segment	country	city	state	postal_code	region	
<input type="checkbox"/>	Edit Copy Delete	1	CC-12670	Craig Carreira	Consumer	United States	Chicago	Illinois	60610	Central
<input type="checkbox"/>	Edit Copy Delete	2	SO-20335	Sean O'Donnell	Consumer	United States	Fort Lauderdale	Florida	33311	South
<input type="checkbox"/>	Edit Copy Delete	3	BS-11590	Brendan Sweed	Corporate	United States	Columbus	Indiana	47201	Central
<input type="checkbox"/>	Edit Copy Delete	4	RF-19840	Roy Franz	Consumer	United States	Chesapeake	Virginia	23320	South
<input type="checkbox"/>	Edit Copy Delete	5	DR-12880	Dan Reichenbach	Corporate	United States	Inglewood	California	90301	West
<input type="checkbox"/>	Edit Copy Delete	6	JE-15745	Joel Eaton	Consumer	United States	Newark	Ohio	43055	East
<input type="checkbox"/>	Edit Copy Delete	7	SJ-20215	Sarah Jordon	Consumer	United States	Columbia	Tennessee	38401	South
<input type="checkbox"/>	Edit Copy Delete	8	MM-18055	Michelle Moray	Consumer	United States	Aurora	Colorado	80013	West
<input type="checkbox"/>	Edit Copy Delete	9	AC-10450	Amy Cox	Consumer	United States	Seattle	Washington	98105	West

phpMyAdmin interface showing the connection settings for the 'lab\_4\_2\_products\_NG' database. A dialog box indicates that the connection to 'lab\_4\_2\_products\_NG' was successful.

Connection name: lab\_4\_2\_products\_NG

Settings:

- Host Name: 95.131.149.21
- Database Name: mgpu\_m\_etl\_18
- Port Number: 3306
- Username: mgpu\_m\_etl\_18
- Password: [Redacted]
- Use Result Streaming Cursor: ☒

phpMyAdmin interface showing the data of the 'products' table. The table contains 7 rows of product data.

	id	product_id	category	sub_category	product_name	person	
<input type="checkbox"/>	Edit Copy Delete	1	OFF-AP-10002578	Office Supplies	Appliances	Fellowes Premier Superior Surge Suppressor, 10-Out...	Chuck Magee
<input type="checkbox"/>	Edit Copy Delete	2	OFF-PA-10000575	Office Supplies	Paper	Wirebound Message Books, Four 2 3/4 x 5 White Form...	Chuck Magee
<input type="checkbox"/>	Edit Copy Delete	3	TEC-MA-10002790	Technology	Machines	NeatDesk Desktop Scanner & Digital Filing System	Kelly Williams
<input type="checkbox"/>	Edit Copy Delete	4	OFF-AR-10000255	Office Supplies	Art	Newell 328	Kelly Williams
<input type="checkbox"/>	Edit Copy Delete	5	TEC-PH-10001061	Technology	Phones	Apple iPhone 5C	Cassandra Brandon
<input type="checkbox"/>	Edit Copy Delete	6	OFF-AR-10003179	Office Supplies	Art	Dixon Ticonderoga Core-Lock Colored Pencils	Anna Andreadi
<input type="checkbox"/>	Edit Copy Delete	7	OFF-AP-10003040	Office Supplies	Appliances	Fellowes 8 Outlet Superior Workstation Surge Prote...	Anna Andreadi

И вариант с использованием job

The workflow diagram illustrates a process starting with 'Start' and 'Set variables'. It then branches into two paths: one leading to 'File exists' and 'HTTP', and another leading to 'Write to log 2' and 'Abort job'. Both paths converge at 'Write to log' and 'Success'. The 'File exists' path includes 'Wait for', 'csv to orders', 'csv to customers', and 'csv\_to\_products' steps. The 'HTTP' path includes 'Wait for' and 'csv to orders' steps. The 'Write to log 2' path includes 'Abort job' and 'Write to log' steps. The 'Write to log' path includes 'csv to customers' and 'csv\_to\_products' steps.

**Execution Results**

Logging History Job metrics Metrics

2025/03/18 23:32:56 - Job CSV\_to\_Mysql\_wo\_PSQL - Finished job entry [HTTP ] (result=[true])  
2025/03/18 23:32:56 - Job CSV\_to\_Mysql\_wo\_PSQL - Finished job entry [File exists] (result=[true])  
2025/03/18 23:32:56 - Job CSV\_to\_Mysql\_wo\_PSQL - Finished job entry [Set variables] (result=[true])  
2025/03/18 23:32:56 - Job CSV\_to\_Mysql\_wo\_PSQL - Job execution finished  
2025/03/18 23:32:56 - Spoon - Job has ended.

phpMyAdmin interface showing the 'customers' table data. The table has columns: id, customer\_id, customer\_name, segment, country, city, state, postal\_code, and region. The data is displayed in a table with 3 rows.

id	customer_id	customer_name	segment	country	city	state	postal_code	region
1	CC-12670	Craig Carreira	Consumer	United States	Chicago	Illinois	60610	Central
2	SO-20335	Sean O'Donnell	Consumer	United States	Fort Lauderdale	Florida	33311	South
3	BS-11590	Brendan Sweed	Corporate	United States	Columbus	Indiana	47201	Central

Activities Google Chrome map 19 20:12

95.131.149.21:8080 / x Social Network Ads x pgAdmin 4 x +

Not secure 95.131.149.21:8080/phpmyadmin/index.php?route=/sql&pos=0&db=mgpu\_m\_etl\_18&table=orders

phpMyAdmin

Recent Favorites

information\_schema  
mgpu\_m\_etl\_18  
New  
ads\_results  
customers  
Columns  
Indexes  
orders  
products  
performance\_schema

Server: localhost:3306 » Database: mgpu\_m\_etl\_18 » Table: orders

Browse Structure SQL Search Insert Export Import Operations Tracking

Showing rows 0 - 24 (9994 total, Query took 0.0003 seconds.)

SELECT \* FROM `orders`

Profiling [ Edit inline ] [ Edit ] [ Explain SQL ] [ Create PHP code ] [ Refresh ]

1 > >> Number of rows: 25 Filter rows: Search this table Sort by key: None

Extra options

	row_id	order_date	ship_date	ship_mode	sales	quantity	discount	profit	returned	
<input type="checkbox"/>	Edit Copy Delete	1	2018-11-08	2018-11-11	Second Class	261.96	2	0.00	41.91	NULL
<input type="checkbox"/>	Edit Copy Delete	2	2018-11-08	2018-11-11	Second Class	731.94	3	0.00	219.58	NULL
<input type="checkbox"/>	Edit Copy Delete	3	2018-06-12	2018-06-16	Second Class	14.62	2	0.00	6.87	NULL
<input type="checkbox"/>	Edit Copy Delete	4	2017-10-11	2017-10-18	Standard Class	957.58	5	0.40	-383.03	NULL
<input type="checkbox"/>	Edit Copy Delete	5	2017-10-11	2017-10-18	Standard Class	22.37	2	0.20	2.52	NULL

Activities Google Chrome map 19 20:12

95.131.149.21:8080 / x Social Network Ads x pgAdmin 4 x +

Not secure 95.131.149.21:8080/phpmyadmin/index.php?route=/sql&pos=0&db=mgpu\_m\_etl\_18&table=products

phpMyAdmin

Recent Favorites

information\_schema  
mgpu\_m\_etl\_18  
New  
ads\_results  
customers  
Columns  
Indexes  
orders  
products  
performance\_schema

Server: localhost:3306 » Database: mgpu\_m\_etl\_18 » Table: products

Browse Structure SQL Search Insert Export Import Operations Tracking Triggers

Showing rows 0 - 24 (5371 total, Query took 0.0002 seconds.)

SELECT \* FROM `products`

Profiling [ Edit inline ] [ Edit ] [ Explain SQL ] [ Create PHP code ] [ Refresh ]

1 > >> Number of rows: 25 Filter rows: Search this table Sort by key: None

Extra options

	id	product_id	category	sub_category	product_name	person
<input type="checkbox"/>	Edit Copy Delete	1	OFF-AP-10002578	Office Supplies Appliances	Fellowes Premier Superior Surge Suppressor, 10-Out...	Chuck Magee
<input type="checkbox"/>	Edit Copy Delete	2	OFF-PA-10000575	Office Supplies Paper	Wirebound Message Books, Four 2 3/4 x 5 White Form...	Chuck Magee
<input type="checkbox"/>	Edit Copy Delete	3	TEC-MA-10002790	Technology Machines	NeatDesk Desktop Scanner & Digital Filing System	Kelly Williams
<input type="checkbox"/>	Edit Copy Delete	4	OFF-AR-10000255	Office Supplies Art	Newell 328	Kelly Williams
<input type="checkbox"/>	Edit Copy Delete	5	TEC-PH-10001061	Technology Phones	Apple iPhone 5C	Cassandra Brandow
<input type="checkbox"/>	Edit Copy Delete	6	OFF-AR-10003179	Office Supplies Art	Dixon Ticonderoga Core-Lock Colored Pencils	Anna Andreadi

**Индивидуальное задание 18:** загрузить данные с фильтром по дате – только заказы 2016 года

Filter rows

Step name: Filter rows

Send 'true' data to step: Value mapper

Send 'false' data to step: Write to log

The condition:

```

order_date IS NOT NULL
AND
ship_date IS NOT NULL
AND
order_date >= [2016-01-01]
AND
order_date <= [2016-12-31]

```

Help OK Cancel

5/03/18 11:53:18 - Spoon - Save file as...

lab\_02\_1\_csv\_orders\_ lab\_02\_2\_csv\_to\_Cust lab\_02\_3\_csv\_to\_prod Job CSV\_to\_Mysql lab\_4\_1\_NurgaleevaG Job CSV\_to\_Mysql\_wo\_83

100%

File exists HTTP Wait for csv to orders csv to customers csv\_to\_products

Write to log 2 Abort job Write to log

### Execution Results

Logging History Job metrics Metrics

2025/03/19 20:34:21 - Job CSV\_to\_Mysql\_wo\_PSQL - Finished job entry [HTTP ] (result=[true])  
 2025/03/19 20:34:21 - Job CSV\_to\_Mysql\_wo\_PSQL - Finished job entry [File exists] (result=[true])  
 2025/03/19 20:34:21 - Job CSV\_to\_Mysql\_wo\_PSQL - Finished job entry [Set variables] (result=[true])  
 2025/03/19 20:34:21 - Job CSV\_to\_Mysql\_wo\_PSQL - Job execution finished  
 2025/03/19 20:34:21 - Spoon - Job has ended.

Server: localhost:3306 » Database: mgpu\_m\_etl\_18 » Table: orders

Browse Structure SQL Search Insert Export Import Operations Tracking Triggers

Showing rows 0 - 24 (1993 total, Query took 0.0002 seconds.)

`SELECT * FROM `orders``

☐ Profiling [ Edit inline ] [ Edit ] [ Explain SQL ] [ Create PHP code ] [ Refresh ]

1 > >> Number of rows: 25 Filter rows: Search this table Sort by key: None

Extra options

	row_id	order_date	ship_date	ship_mode	sales	quantity	discount	profit	returned
<input type="checkbox"/> Edit Copy Delete	6	2016-06-09	2016-06-14	Standard Class	48.86	7	0.00	14.17	NULL
<input type="checkbox"/> Edit Copy Delete	7	2016-06-09	2016-06-14	Standard Class	7.28	4	0.00	1.97	NULL
<input type="checkbox"/> Edit Copy Delete	8	2016-06-09	2016-06-14	Standard Class	907.15	6	0.20	90.72	NULL
<input type="checkbox"/> Edit Copy Delete	9	2016-06-09	2016-06-14	Standard Class	18.50	3	0.20	5.78	NULL
<input type="checkbox"/> Edit Copy Delete	10	2016-06-09	2016-06-14	Standard Class	114.90	5	0.00	34.47	NULL
<input type="checkbox"/> Edit Copy Delete	11	2016-06-09	2016-06-14	Standard Class	1706.18	9	0.20	85.31	NULL

Проверка:

`SELECT distinct (year(order_date)) FROM `orders`;`

☐ Profiling [ Edit inline ] [ Edit ] [ Explain SQL ] [ Create PHP code ] [ Refresh ]

☐ Show all Number of rows: 25 Filter rows: Search this table

Extra options

(year(order\_date))

2016

## SQL скрипты

-- Таблица заказов (основная информация о продажах)

```
CREATE TABLE orders (  
  row_id INT PRIMARY KEY,  
  order_date DATE,  
  ship_date DATE,  
  ship_mode VARCHAR(50),  
  sales DECIMAL(10,2),  
  quantity INT,  
  discount DECIMAL(4,2),  
  profit DECIMAL(10,2),  
  returned TINYINT(1) DEFAULT 0 -- 1 = Yes, 0 = No  
);
```

-- Таблица клиентов

```
DROP TABLE IF EXISTS customers;  
CREATE TABLE customers (  
  id INT AUTO_INCREMENT PRIMARY KEY,  
  customer_id VARCHAR(20) NOT NULL,  
  customer_name VARCHAR(100),  
  segment VARCHAR(50),  
  country VARCHAR(100),  
  city VARCHAR(100),  
  state VARCHAR(100),  
  postal_code VARCHAR(20),  
  region VARCHAR(50),  
  INDEX idx_customer_id (customer_id),  
  INDEX idx_region (region)  
);
```

-- создаем таблицу products

```
DROP TABLE IF EXISTS products;  
CREATE TABLE products (  
  id INT AUTO_INCREMENT PRIMARY KEY,  
  product_id VARCHAR(20) NOT NULL,  
  category VARCHAR(50),  
  sub_category VARCHAR(50),  
  product_name VARCHAR(255),  
  person VARCHAR(100),  
  INDEX idx_product_id (product_id),  
  INDEX idx_category (category),  
  INDEX idx_subcategory (sub_category)  
);
```

## Контрольные вопросы

### 1. Что такое динамические соединения в PDI?

Динамические соединения в PDI позволяют:

- Использовать параметры подключения из внешних источников.
- Менять настройки соединения во время выполнения.
- Обработать множество источников данных в одном процессе.

### 2. Как организовать обработку ошибок в трансформации?

Обработку ошибок можно организовать через занесение их в Log. После этого можно просмотреть журнал ошибок.

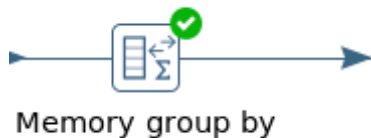
Также при варианте, в котором задействуется шаг с Log, можно предусмотреть последующим шагом остановку выполнения алгоритма.

Компоненты обработки ошибок.

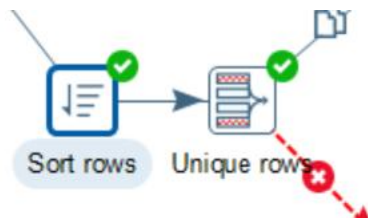
- `wrt-execution_error` - запись информации об ошибках.
- `abrt-execution_error` - прерывание выполнения при критических ошибках.

### 3. Какие методы выявления дублей существуют?

- можно удалить дубли используя группировку по всем столбцам



- через использование функции unique rows



### 4. Как настроить параметризацию подключений?

Для настройки параметризации подключений в Pentaho можно использовать именованные параметры.

Эта система позволяет параметризовать преобразования и задания.

Чтобы настроить параметры, нужно:

1. Перейти в диалоговое окно настроек преобразования или задания и найти вкладку «Параметры».
2. Установить значения для параметров, которые будут использоваться во время выполнения. Если для параметра не установлено значение, используется значение по умолчанию.
3. В диалоговом окне выполнения преобразований и заданий можно установить значение для каждого определённого именованного параметра.

### 5. Какие компоненты PDI используются для объединения данных?



При помощи компонента **Pentaho Data Integration** также известного как Kettle. Он используется для интеграции данных из разных источников.

Один из востребованных шагов в Pentaho Data Integration — объединение (join) двух потоков данных. Для этого в настройках можно выбрать тип соединения (Inner, Left Outer, Right Outer, Full Outer) и указать ключи.

Обязательное требование для шага join — входные данные должны быть отсортированы по ключевым полям, для этого в Pentaho Data Integration применяется отдельный шаг **Sorter**