

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

Кузьмина Дарья Юрьевна БД-241м

Практическая работа

Лабораторная работа 4.1. Установка и настройка ETL-инструмента.

Вариант 11

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

Руководитель дисциплины:
Босенко Т.М., доцент департамента
информатики, управления и технологий,
доктор экономических наук

Москва
2025

Содержание

Введение	2
Основная часть.....	4
Заключение	9

Введение

Создание конвейеров данных

Цель работы: изучение основных принципов работы с ETL-инструментами на примере Pentaho [Data](#) Integration (PDI), настройка конвейера обработки данных, фильтрация и замена значений в Excel-файле, а также выгрузка обработанных данных в базу данных MySQL/PostgreSQL.

Условие выполнения работы:

Работа выполняется в образе **Ubuntu 22.04 (.ova) для VirtualBox 7.0** https://disk.yandex.ru/d/gagWU_zn1erR8g, в котором предварительно установлены все необходимые компоненты для работы с Pentaho [Data](#) Integration, либо проводится установка окружения в ОС Linux.

Задачи

- **Настроить среду для работы с Pentaho [Data](#) Integration (PDI):**

Запуск виртуальной машины с **Ubuntu 22.04** в **VirtualBox**.

Проверка установки Java и WebKitGTK.

Развертывание Pentaho [Data](#) Integration.

- **Создать ETL-конвейер:**

Загрузить данные из **CSV-файла**.

Очистить, преобразовать и отфильтровать данные.

Выполнить замену значений.

Выгрузить обработанные данные в **MySQL** или **PostgreSQL**.

- **Проверить корректность обработки:**

Выполнить SQL-запросы для проверки результата.

Подготовить отчет с описанием проделанных шагов.

Инструменты и технологии

Для выполнения лабораторной работы используются следующие инструменты:

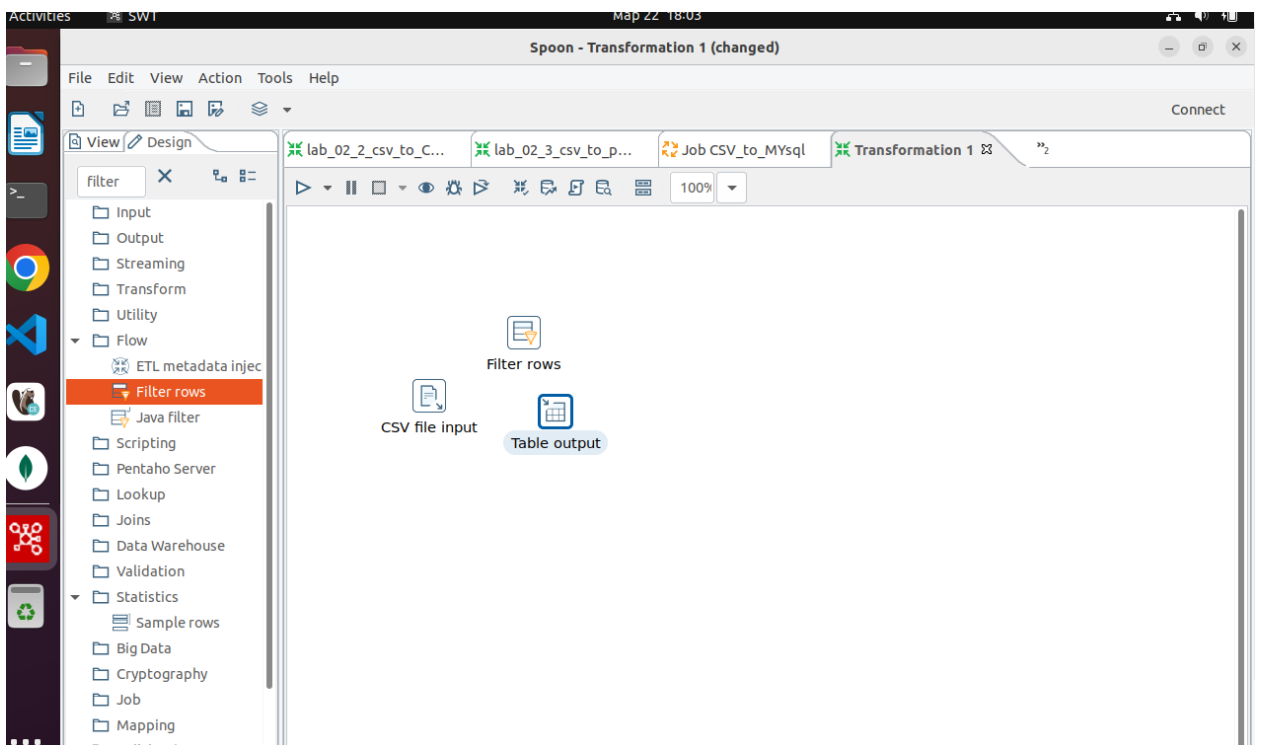
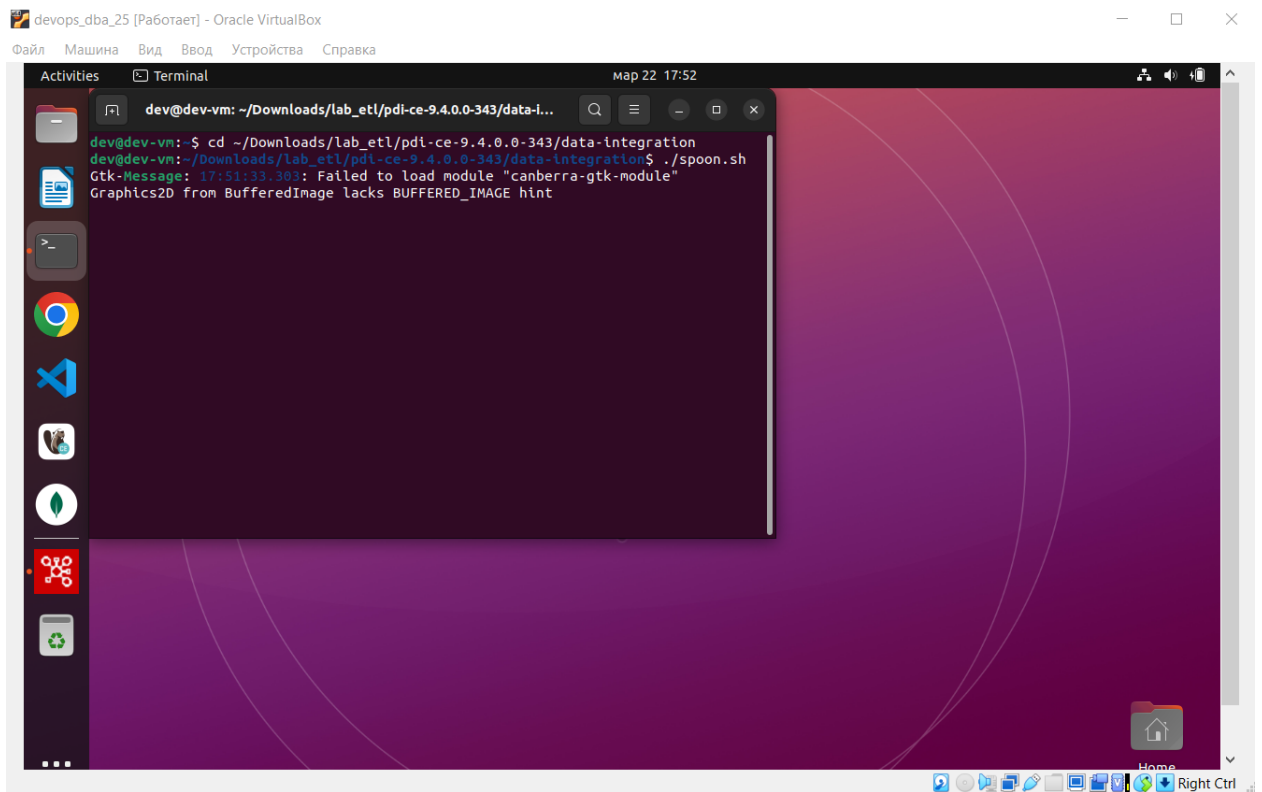
Компонент	Описание
Ubuntu 22.04 (.ova)	Образ операционной системы для VirtualBox 7.0
VirtualBox 7.0	Виртуализация среды
Pentaho Data Integration 9.4	ETL-инструмент для работы с данными
MySQL/PostgreSQL	База данных для хранения обработанных данных
CSV-файлы	Исходные данные для обработки
Java 11	Необходима для работы Pentaho
libwebkitgtk-1.0-0	Библиотека для корректного запуска Spoon
SQL	Язык запросов для работы с базами данных

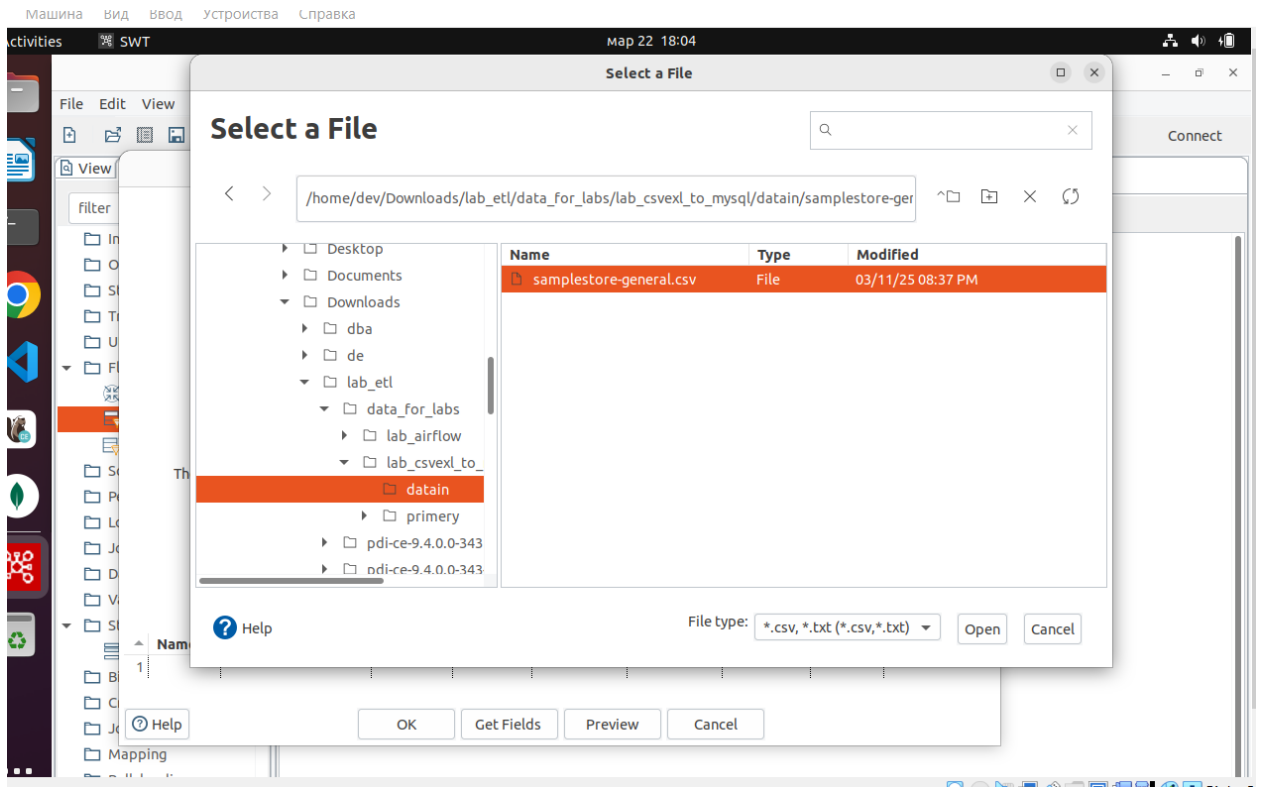
Формат предоставления отчета по лабораторной работе:

1. Скачать CSV-датасет с Kaggle.
2. Создать Pentaho ETL-конвейер с фильтрацией, заменой и обработкой данных.
3. Выгрузить данные в MySQL/PostgreSQL.
4. Подготовить отчет с:
 - Описанием процесса;
 - Скриншотами Spoon;
 - SQL-запросами для проверки.
5. Загрузить отчет, CSV-датасет (или ссылку на github) и **lab_4_01.ktr** в LMS.

Основная часть

Pentaho запускаем

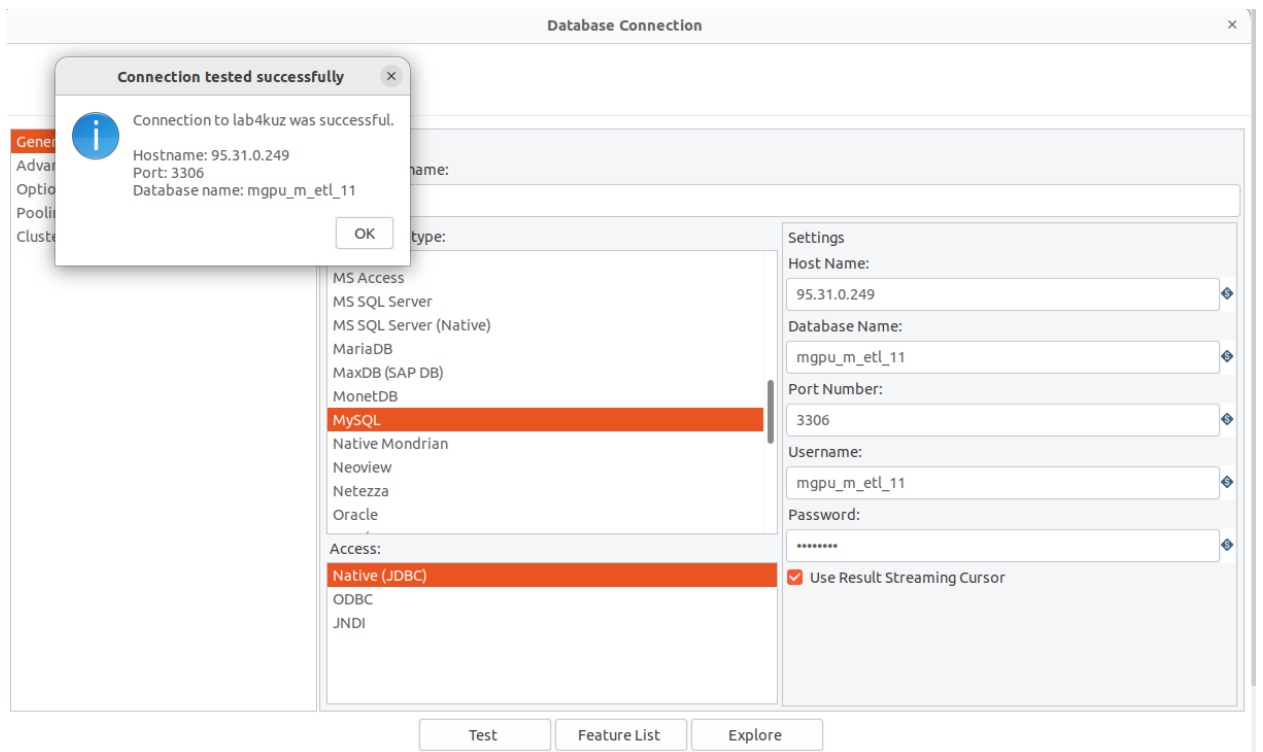




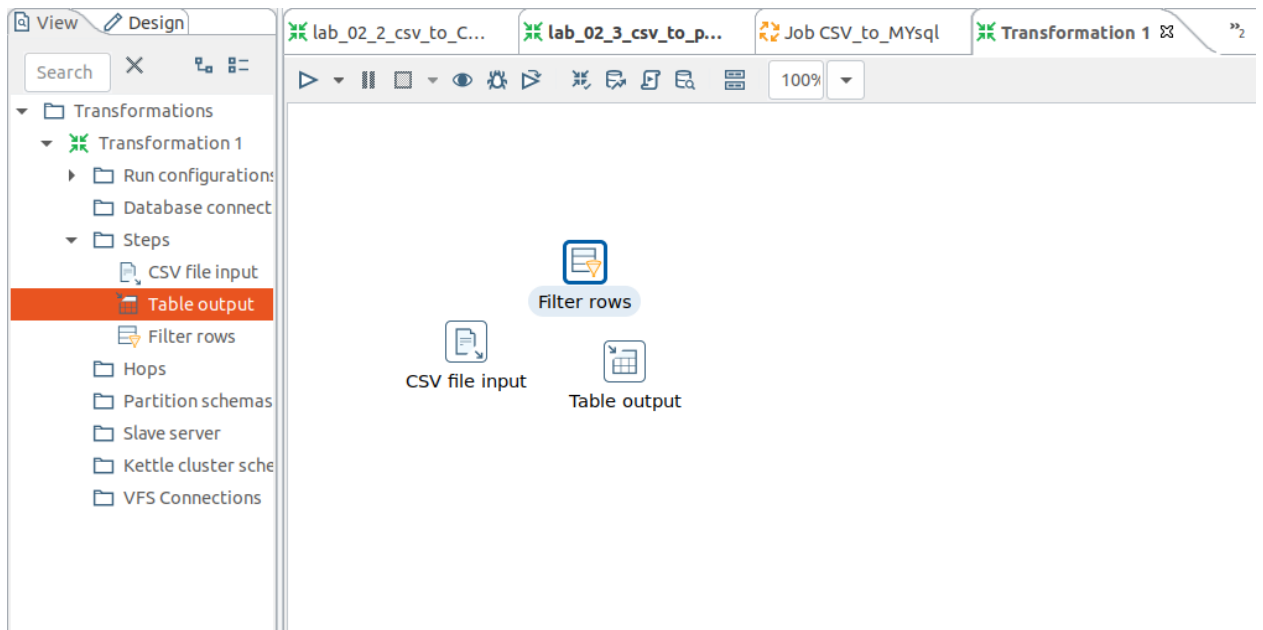
Выбираем файл

Выбираем разделитель

Протестировали и сделали первую трансформацию.



Готово



Переходим в заготовку, проверяем зашли ли все данные

Spoon - lab_02_1_csv_orders

File Edit View Action Tools

View

Search

Examine preview data

Rows of step: CSV file input (1000 rows)

▲	Row ID	Order ID	Order Date	Ship Date	Ship Mode	Customer ID	Customer Name
1	6569	CA-2016-100678	18/04/2016	22/04/2016	Standard Class	KM-16720	Kunst Miller
2	6570	CA-2016-100678	18/04/2016	22/04/2016	Standard Class	KM-16720	Kunst Miller
3	6571	CA-2016-100678	18/04/2016	22/04/2016	Standard Class	KM-16720	Kunst Miller
4	6572	CA-2016-100678	18/04/2016	22/04/2016	Standard Class	KM-16720	Kunst Miller
5	6315	CA-2016-100762	24/11/2016	29/11/2016	Standard Class	NG-18355	Nat Gilpin
6	6316	CA-2016-100762	24/11/2016	29/11/2016	Standard Class	NG-18355	Nat Gilpin
7	6317	CA-2016-100762	24/11/2016	29/11/2016	Standard Class	NG-18355	Nat Gilpin
8	6318	CA-2016-100762	24/11/2016	29/11/2016	Standard Class	NG-18355	Nat Gilpin
9	6252	CA-2016-101147	02/12/2016	04/12/2016	First Class	MC-17575	Matt Collins
10	1575	CA-2016-101602	15/12/2016	18/12/2016	First Class	MC-18100	Mick Crebagga
11	1576	CA-2016-101602	15/12/2016	18/12/2016	First Class	MC-18100	Mick Crebagga
12	9558	CA-2016-103086	17/10/2016	19/10/2016	Second Class	EB-14170	Evan Bailliet
13	4283	CA-2016-103100	20/12/2016	23/12/2016	First Class	AB-10105	Adrian Barton
14	4284	CA-2016-103100	20/12/2016	23/12/2016	First Class	AB-10105	Adrian Barton
15	5725	CA-2016-103191	22/09/2016	27/09/2016	Standard Class	VG-21805	Vivek Grady
16	7546	CA-2016-103492	10/10/2016	15/10/2016	Standard Class	CM-12715	Craig Molinari
17	7547	CA-2016-103492	10/10/2016	15/10/2016	Standard Class	CM-12715	Craig Molinari
18	7548	CA-2016-103492	10/10/2016	15/10/2016	Standard Class	CM-12715	Craig Molinari

Close

Show Log

Help

OK

Get Fields

Preview

Cancel

SWT Map 22 21:24

Spoon - lab_02_1_csv_orders

File Edit View Action Tools Help

Connect

Examine preview data

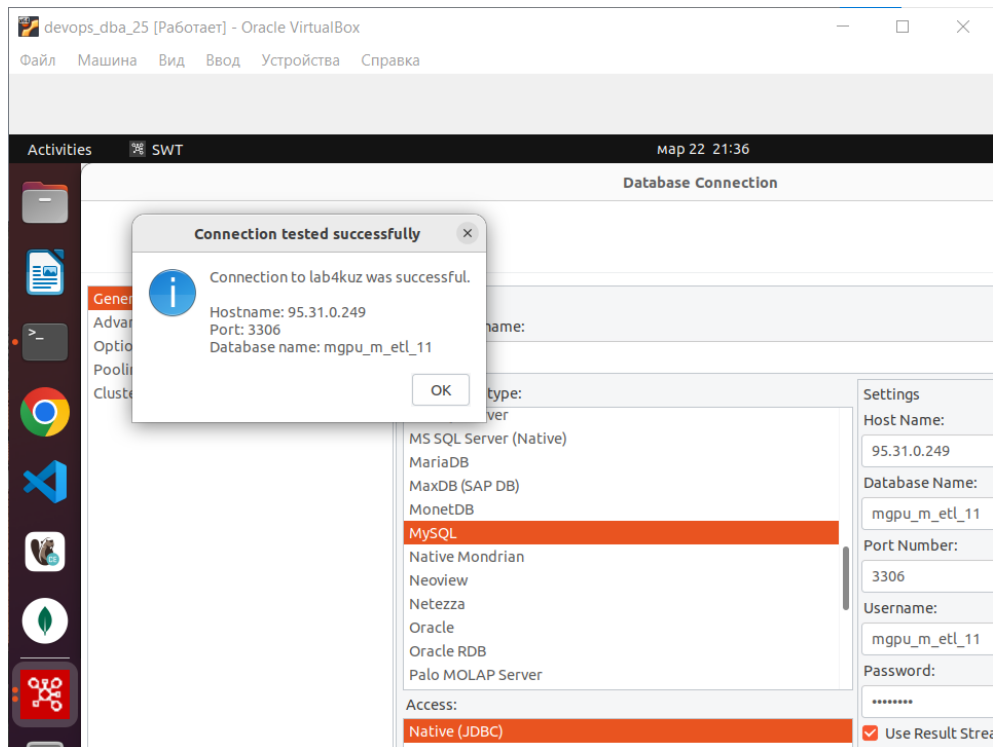
Rows of step: CSV file input (1000 rows)

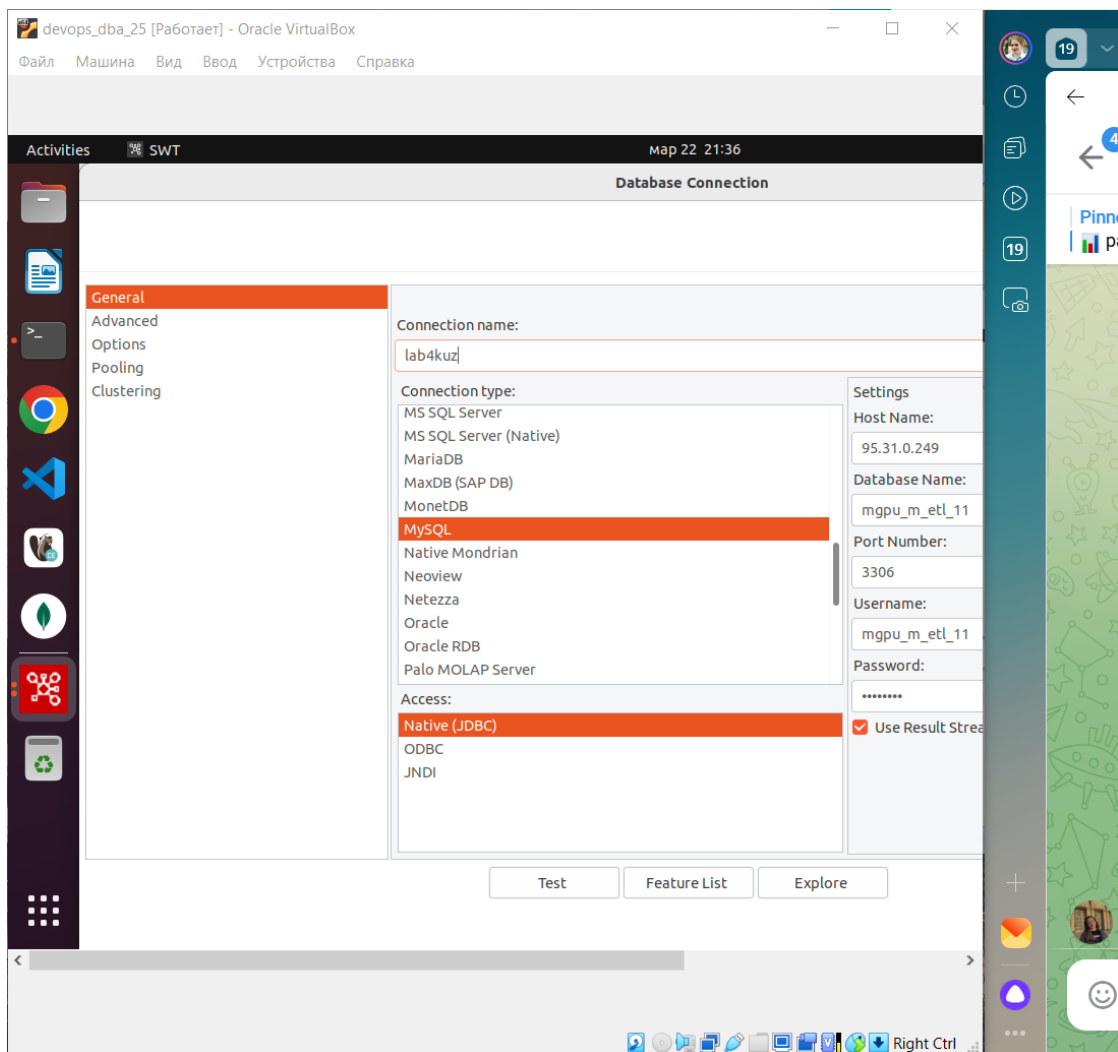
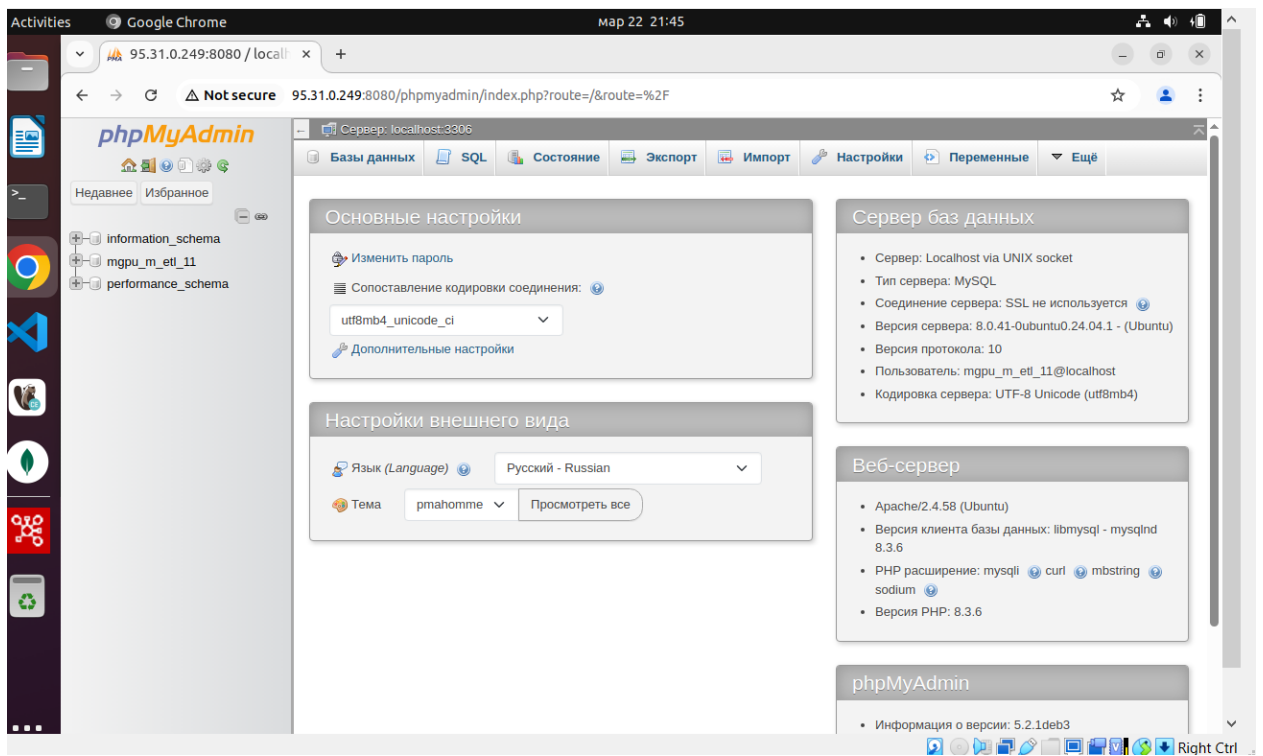
	Sales	Quantity	Discount	Profit	Returned	Person
ticks	2.7	2	0.2	1	<null>	Kelly Williams
Mission-style Design Wood Chair	317.1	3	0.3	-18.1	<null>	Kelly Williams
elopes	149.4	3	0.2	50.4	<null>	Kelly Williams
use Optical USB Trackball for PC or Mac	228	3	0.2	28.5	<null>	Kelly Williams
tric Pencil Sharpener, Blue	151.9	4	0	45.6	Yes	Kelly Williams
e Reel Labels, White, 5000/Box	196.6	2	0	96.3	Yes	Kelly Williams
	144.1	3	0	69.2	Yes	Kelly Williams
sage Book w/Frequently-Called Numbers Space, 400 Messages per Book	16	2	0	8	Yes	Kelly Williams
	2.4	1	0.8	-6.3	<null>	Kelly Williams
shboard Air Vent Car Mount Holder	40.7	3	0.2	-9.2	<null>	Kelly Williams
lding Chairs, 4/Carton	763.3	5	0.3	-21.8	<null>	Kelly Williams
Recessed Floodlight Bulbs	5.3	2	0.6	-1.6	<null>	Kelly Williams
	3.7	1	0	1.7	<null>	Kelly Williams
annual Binding System	1104	3	0	496.8	<null>	Kelly Williams
re Shelving, Black	331.5	3	0.2	-82.9	<null>	Kelly Williams
	720	6	0.2	72	<null>	Kelly Williams
Business Phone System	755.9	7	0.2	66.1	<null>	Kelly Williams
Strips	12	5	0.8	-19.2	<null>	Kelly Williams

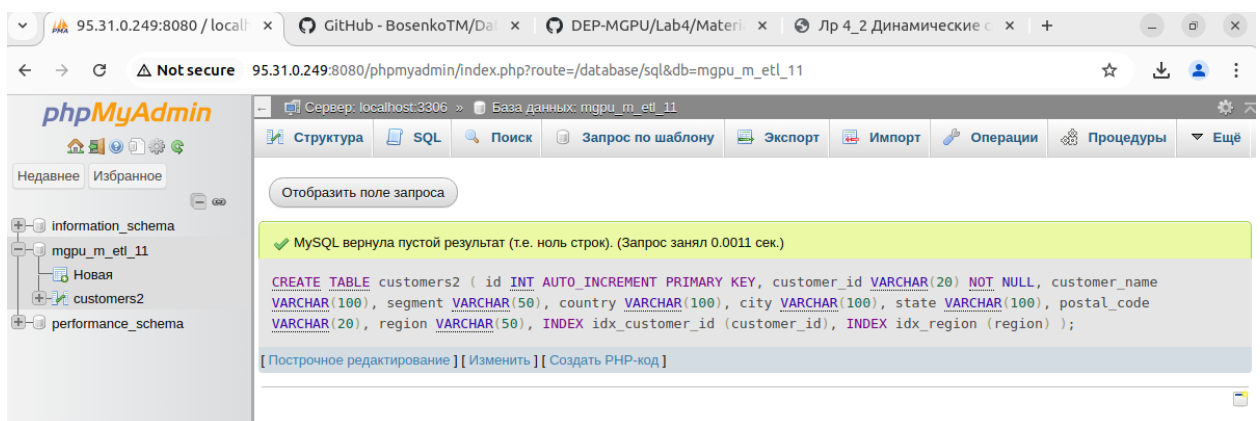
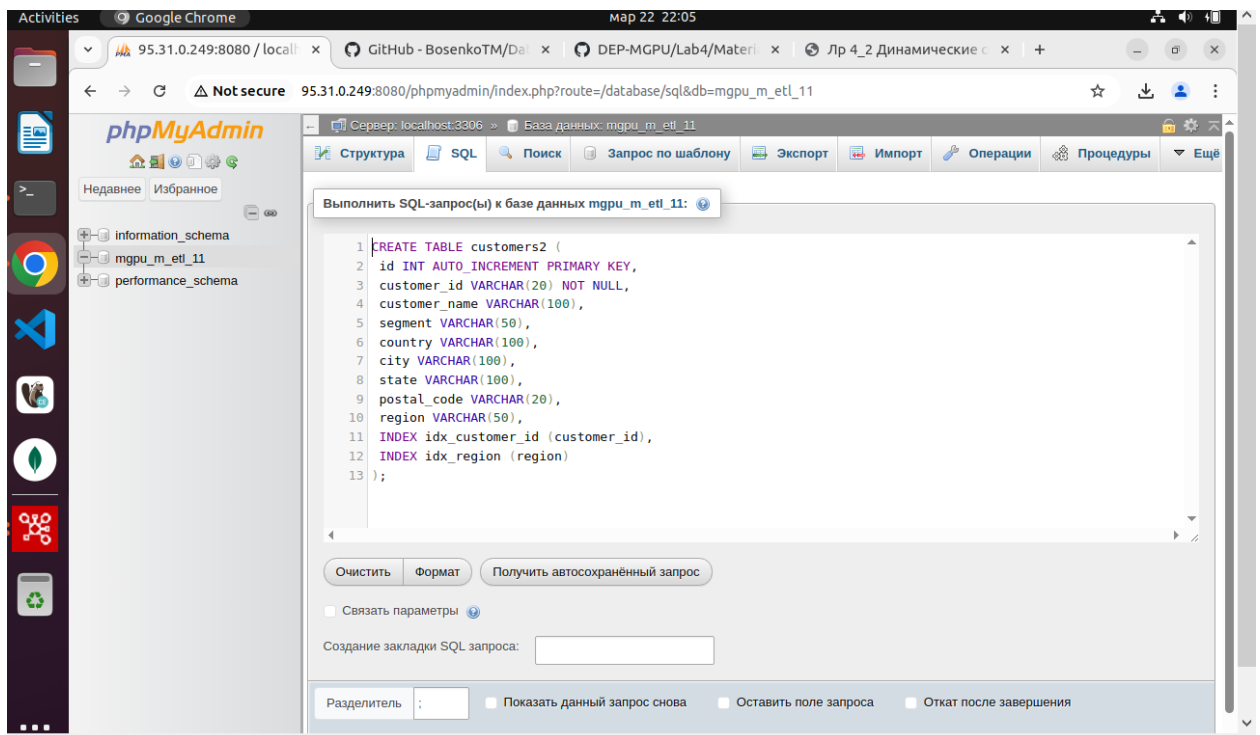
Close Show Log

Рассмотрим преобразование значение <null> в 0

Подключим датасет







Заключение

В ходе выполнения лабораторной работы были изучены основные принципы работы с ETL-инструментом Pentaho Data Integration (PDI). Была настроена среда для работы с PDI, включая установку Java, WebKitGTK и развертывание Pentaho. Создан ETL-конвейер, который включал загрузку данных из CSV-файла, их очистку, фильтрацию и преобразование. Обработанные данные были выгружены в базу данных MySQL/PostgreSQL.

- Pentaho Data Integration предоставляет удобный интерфейс для создания ETL-процессов.
- Важным этапом является предварительная очистка и трансформация данных для обеспечения их качества.
- Интеграция данных из CSV в базу данных позволяет автоматизировать процессы анализа и отчетности.