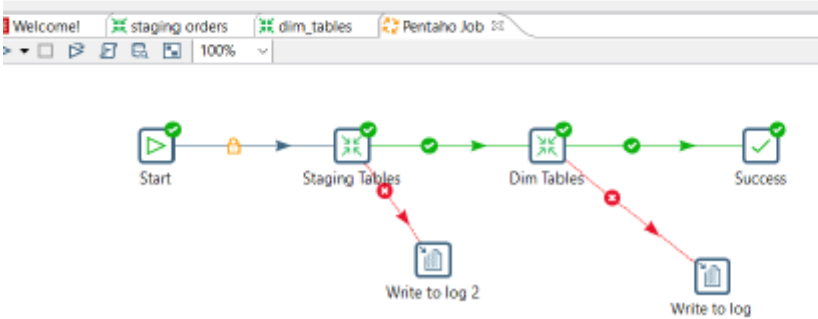
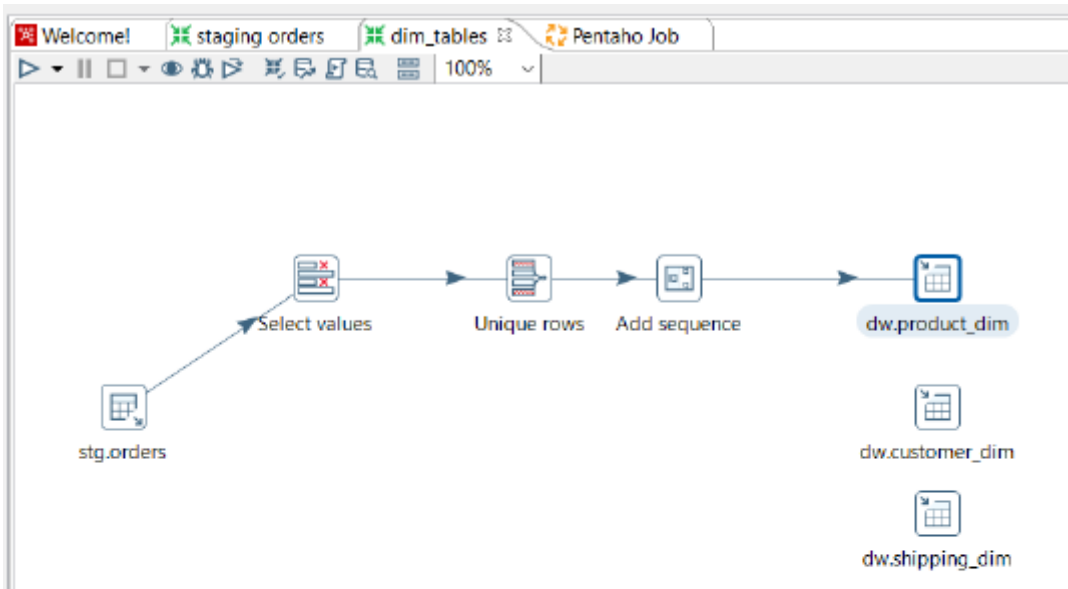


Лабораторная работа 4.1 Соловьев Михаил БД-231м

- 1. Скачать отсюда и запустить Pentaho DI. Pentaho DI требует установку Java 8. Попробуйте скачать архив и распаковать его. Необходимо запустить spoon.sh для Linux/Mac и spoon.bat для Windows. Видео по установке Pentaho DI на примере Windows 10 +
- 2. Скачать примеры Pentaho jobs для Staging и Dimension Tables. +



- 3. Создайте еще одну трансформацию, в которой создать sales_fact таблицу. +

id	customer_id	customer_name	ship_date	ship_mode	product_id	region	country	city	state	quantity	sales	profit	discount
1	16-15985	Steven Brown	16-01-07 00:00:00	Standard Class	QIF-PA-10000074	Central	United States	Houston	Texas	2	16,448	5,551.82	
2	16-160264	Jack Citland	16-01-07 00:00:00	First Class	QIF-AB-10000209	South	United States	Atlanta	Georgia	3	5,178	5,238	
3	16-112526	Phyllis Ober	16-01-08 00:00:00	Standard Class	QIF-BI-10004984	Central	United States	Peapack	New Jersey	2	3,54	3,487	
4	16-112526	Phyllis Ober	16-01-08 00:00:00	Standard Class	QIF-LA-10000223	Central	United States	Peapack	New Jersey	3	11,384	4,217	
5	16-112526	Phyllis Ober	16-01-08 00:00:00	Standard Class	QIF-SI-10000214	Central	United States	Peapack	New Jersey	3	27,738	14,748	
6	16-112526	Phyllis Ober	16-01-08 00:00:00	Standard Class	QIF-PA-10000000	Central	United States	Peapack	New Jersey	3	16,44	9,312	
7	16-167198	Maria Hernandez	16-01-10 00:00:00	Standard Class	QIF-AB-10000000	South	United States	Henderson	Kentucky	2	5,182	5,467	
8	16-167198	Maria Hernandez	16-01-10 00:00:00	Standard Class	QIF-AB-10000000	South	United States	Henderson	Kentucky	2	5,18	5,467	
9	16-167198	Maria Hernandez	16-01-10 00:00:00	Standard Class	QIF-AB-10000000	South	United States	Henderson	Kentucky	2	5,18	5,467	
10	16-167198	Maria Hernandez	16-01-10 00:00:00	Standard Class	QIF-AB-10000000	South	United States	Henderson	Kentucky	2	5,18	5,467	
11	16-167198	Maria Hernandez	16-01-10 00:00:00	Standard Class	QIF-AB-10000000	South	United States	Henderson	Kentucky	2	5,18	5,467	

Main options Database fields		
Fields to insert:		
#	Table field	Stream field
1	order_id	order_id
2	customer_id	customer_id
3	customer_n...	customer_na...
4	ship_date	ship_date
5	ship_mode	ship_mode
6	product_id	product_id
7	region	region
8	country	country
9	city	city
1..	state	state
1..	quantity	quantity
1..	sales	sales
1..	profit	profit
1..	discount	discount
1..	sales_PK	sales_PK

4.Выявить 8-10 подсистем в ETL Pentaho DI и написать небольшой отчет, в котором приложить print screen компонента (ETL подсистемы) и написать про его свойства. Результат сохраните в Git.



stg.orders

Table input - Используется для считывание данных их таблицы в базе данных.

Необходимо настраивать подключение к базе.



Microsoft Excel input

Microsoft Excel input - Используется для чтения данных из файлов формата Excel. Он предоставляет возможность извлечения данных из листов Excel и дальнейшей их обработки в рамках трансформации. Excel Input позволяет читать данные из файлов формата Excel (.xls или .xlsx). Также можно указать конкретный лист, из которого следует извлечь данные. И можно определить структуру данных, указав имена столбцов, их типы данных и тд.



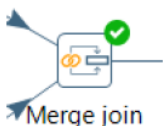
Select values

Select values - Используется для выбора (в том числе для удаления лишних столбцов), переименования и преобразования столбцов данных в рамках трансформации.



Sort rows

Sort rows - Используется для сортировки строк данных в рамках трансформации. В частности, используется перед выполнением операции слияния данных Merge Join.



Merge join

Merge join - Используется для объединения двух наборов данных из различных источников данных, таких как таблицы баз данных, файлы или другие источники данных, на основе общих значений столбцов. Можно настроить тип объединения (INNER, LEFT, RIGHT) и указать столбцы для сравнения при выполнении операции объединения. Также поддерживает обработку дубликатов.



Calculator

Calculator - Используется для проверки качества данных и выполнения различных видов валидации данных в рамках трансформации. Data Validator может:

- проверять типы данных в столбцах данных и выявлять некорректные значения, несоответствующие ожидаемым типам данных.
- проверять наличие значений в определенных столбцах данных и выявлять строки с недостающими данными.
- поддерживает проверку формата данных, таких как даты, времена, числа или текстовые строки, на соответствие определенным форматам.



Unique rows Unique rows - Используется для удаления дубликатов строк из набора данных. Этот компонент позволяет получить уникальные строки данных на основе определенных критериев и удалить повторяющиеся строки.



Group by Group by - Позволяет группировать данные и использовать агрегаты