

## **OBETA Script: ETL Implementation**

<b>1. Univariate Analysis</b>	1.1 Overview of Each Variable and Define Key Variables with Statistical Description
<b>2. Staging Phase</b>	2.1 Create Tables in Database (MySQL)
<b>3. Transformation Phase</b>	3.1 Import Product Data and Pick Data from MySQL 3.2 Merge the Product and Pick Tables 3.3 Identify Dirty Data: <i>duplicated, missing, inconsistent values and outliers</i> 3.4 Clean the Data
<b>4. Production DB</b>	4.1 Create Dimensions Based on ER Model 4.2 Load Dimensions Tables Back to Production Database 4.3 Connect Production Database with Tableau
<b>5. KPIs calculation</b>	5.1 Update KPIs
<b>6. Data Visualization</b>	6.1 Create Dashboards Based on BI and KPIs 6.2 Create Tableau Story for Final Presentation

## 1.1 OVERVIEW OF EACH VARIABLE AND DEFINE KEY VARIABLES WITH STATISTICAL DESCRIPTION

Column Name	Variable Type	Key Variable	Statistical Description
<b>product_id</b>	Nominal	✓	Count: 33493460 Unique Count: 97211 Top: 109910 Freq: 170062
<b>description</b>	Nominal		
<b>product group</b>	Nominal	✓	Count: 33493460 Unique Count : 15 Top: 31_install. Freq: 8,760,367
<b>product_id</b>	Nominal		
<b>warehouse section</b>	Nominal	✓	Count:33493460 Unique Count : 5 Top: SHL Freq: 14421610
<b>origin</b> (46:store, 48:customer)	Dichotomous	✓	Count:33493460 Unique Count : 2 Top: “48” Freq: 217458068
<b>order_no</b>	Nominal	✓	Count: 33493460 Unique Count : 9344863 Top: 201604055714 Freq: 371
<b>position</b>	Ordinal		
<b>pick volume</b>	Metric	✓	Count: 33493460 Mean: 45.93 Median: 5 Mode: 1 SD: 127.42 Min: -200 Max: 1160
<b>qty_unit</b>	Nominal		
<b>date</b>	Ordinal		

## 2.1 CREATE TABLES IN DATABASE (MySQL)

Purpose	Query
Create product table in “OBETA” schema	Creation of Product data Table CREATE TABLE product_data (productid INT, product_name VARCHAR (50) NULL, product_group VARCHAR (100) NULL);
Create pick table in “OBETA” schema	Create schemas

	<pre>CREATE SCHEMA `staging_data` ; CREATE TABLE product_pick (productid VARCHAR (100) NULL, warehouse_area VARCHAR (100) NULL, order_origin VARCHAR (100) NULL, order_number VARCHAR (100) NULL, position_in_order VARCHAR (100) NULL, volume_of_order INT NULL, quantity_unit VARCHAR (50) NULL, date_of_order DATETIME);</pre>
<b>Import product.csv into MySQL</b>	<pre>LOAD DATA local INFILE "C:\\ProgramData\\MySQL\\MySQL Server 8.0\\Uploads\\product_data.csv" INTO TABLE product_data CHARACTER SET latin1 FIELDS TERMINATED BY ',' ENCLOSED BY " LINES TERMINATED BY '\n' IGNORE 1 ROWS; set global local_infile=1;</pre>
<b>Import pick.csv into MySQL</b>	<pre>LOAD DATA local INFILE "C:\\ProgramData\\MySQL\\MySQL Server 8.0\\Uploads\\pick_data.csv" INTO TABLE pick_data FIELDS TERMINATED BY ',' ENCLOSED BY "" LINES TERMINATED BY '\n'; set global local_infile=1;</pre>

### 3.1 IMPORT PRODUC DATA AND PICK DATA FROM MySQL

Purpose	Platform/Libraries	Query
<b>Connect Rstudio to MySQL (Database) and Copy Datasets</b>	Rstudio/DBI, RMySQL, dplyr	<pre># connect to staging data Schema pro_con&lt;- dbConnect(RMySQL::MySQL(),</pre>

		<pre> dbname="staging_data", host='localhost', port=3306, user='root', password='mysql')  dbListTables(pro_con, 'product_data')  #call product data from mysql prod_d &lt;- tbl(pro_con, "product_data")  #copy product data to R global env't prod_df &lt;- collect(prod_d)  #call pick data from mysql pick_d &lt;- tbl(pro_con, "pick_data")  #copy pick data to R global env't pick_df &lt;- collect(pick_d) </pre>
--	--	---

### 3.2 MERGE THE PRODUCT AND PICK TABLES WITH RSTUDIO

Purpose	Query
Join the Product and Pick datasets	<code>joined_table &lt;- inner_join (pick_df, prod_df, by = "product_id")</code>

### 3.3 IDENTIFY DIRTY DATA WITH RSTUDIO

Purpose	Query
Find duplicated rows	<pre> #Create a DF and flag for duplicated rows joined_df_dupli &lt;- joined_df %&gt;%   group_by_all() %&gt;% </pre>

	<pre>mutate(is_duplicate = n() &gt; 1)  table(joined_df_dupli\$is_duplicate)  # filter duplicate rows only_dupli_rows &lt;- joined_df %&gt;%   group_by_all() %&gt;%   filter(n() &gt; 1)</pre>
<b>Find missing values</b>	<pre>missing_values &lt;- joined_df %&gt;%   summarise_all(~ sum(is.na(.)))  #Find missing strings in all columns missing_strings &lt;- joined_df %&gt;%   summarise_all(~ sum(. == ""))</pre>
<b>Find duplicated order IDs (more than one order year)</b>	<pre>#Create a order year column joined_table\$order_year &lt;- year(joined_table\$date_of_order)  #filter out order IDs which duplicate themselves in more than one order year duplicated_orderid &lt;- joined_table %&gt;%   group_by(order_number) %&gt;%   filter(n_distinct(order_year) &gt; 1) %&gt;%   ungroup() table()</pre>
<b>Find rows with zero order volumes</b>	<pre>#Filter rows with zero order volumes zero_values &lt;- filter(joined-table, volume_of_order == 0)</pre>
<b>Find rows with negative order volumes</b>	<pre>#Filter rows with negative order volumes negative_values &lt;- filter(joined-table, volume_of_order &lt; 0)  table(zero_values)</pre>
<b>Detect outliers</b>	<pre># Create outlier flag function outlier.create.flag &lt;- function(x) {hist(x)</pre>

	<pre> summary(x)  # using here the z values flag&lt;- ifelse(test = scale(x) &gt; 3   scale(x)&lt; -3, yes = 1, no = 0) # number of outlier table(flag)  # return the vector representing the flag variable flag}  # Create outlier column for joined data and add to the dataframe join_non_zero\$volume.outlier.flag&lt;- outlier.create.flag(join_non_zero\$volume_of_order)  table(join_non_zero\$volume.outlier.flag)  #create only outlier table on R  only_outlier_table &lt;- filter(joined_df, volume.outlier.flag ==1) </pre>
--	---

### 3.4 CLEAN THE DATA WITH RSTUDIO

Purpose	Query
Remove duplicated rows	<code>join_df_unique &lt;- joined_df %&gt;% distinct()</code>
Remove rows with zero value	<code>join_non_zero &lt;- join_df_unique %&gt;% filter(volume_of_order != "0")</code>
Create unique order ID for duplicated values	<code>#Merge year of order and order ID join_non_zero\$orderid_unique&lt;- do.call(paste, c(join_non_zero [column_merge], sep=""))</code>
Create a filtered table to exclude outliers	<code>join_clean_df &lt;- filter(join_non_zero, volume.outlier.flag ==0)</code>

## 4.1 CREATE DIMENSIONS BASED ON ER MODEL

Purpose	Platform/Libraries	Query
<b>Create Warehouse Section Dimension Table</b>	Rstudio/dplyr	<pre>#create warehouse name and column join_facts &lt;- join_facts %&gt;%   mutate(section_name = case_when(     grepl("AKL", warehouse_area, ignore.case = TRUE) ~ "Automated Small-     Parts Warehouse",     grepl("HRL", warehouse_area, ignore.case = TRUE) ~ "High Bay Storage for     Pallets",     grepl("SHL", warehouse_area, ignore.case = TRUE) ~ "Shuttle house",     grepl("Kabellager", warehouse_area, ignore.case = TRUE) ~ "Cable Storage",     grepl("Manuell", warehouse_area, ignore.case = TRUE) ~ "Manual     Warehouse",     TRUE ~ "Other"))</pre>
<b>Create Product Dimension Table</b>	Rstudio	<pre>product_dim &lt;- join_facts[,c("product_id", "description", "product_group",                              "quantity_unit")] #remove duplicated rows product_dim &lt;- product_dim %&gt;% distinct()</pre>
<b>Create Order Dimension Table</b>	Rstudio	<pre>#Create Order dimension table  order_dim &lt;- join_facts[,c("orderid_unique", "order_number", "order_origin",                            "position_in_order", "date_of_order")]  #remove duplicated rows order_dim &lt;- order_dim %&gt;% distinct()</pre>
<b>Create Time Dimension Table</b>	Rstudio/lubridate	<pre>#Create Time Dimension time_dim &lt;- join_facts[,c("date_of_order", "order_year" )]  # Extract day time_dim\$day &lt;- day(time_dim\$date_of_order)</pre>

		#extract_month time_dim\$month <- month(time_dim\$date_of_order)
<b>Create Facts Table</b>	Rstudio	#Create Facts Table facts_table <- join_facts[,c("warehouse_area", "orderid_unique", "product_id", "date_of_order", "position_in_order", "volume of order", "pick efficiency")]

## 4.2 LOAD DIMENSION TABLES BACK TO PRODUCTION DATABASE

Purpose	Platform/Libraries	Query
<b>Connect Rstudio to Production Schema on MySQL</b>	Rstudio/DBI, RMySQL, dplyr	# connect to production data Schema prod_con<- dbConnect(RMySQL::MySQL(), dbname="production_data", host='localhost', port=3306, user='root', password='mysql')
<b>Copy Warehouse Section Dimension Table into MySQL's Schema Production Data</b>	Rstudio/DBI, RMySQL, dplyr	dbWriteTable(prod_con,"warehouse_section_dim",warehouse_section_dim,row.names=FALSE)
<b>Copy Product Dimension Table into MySQL's Schema Production Data</b>	Rstudio/DBI, RMySQL, dplyr	dbWriteTable(prod_con,"product_dim",product_dim,row.names=FALSE)
<b>Copy Order Dimension Table into MySQL's Schema Production Data</b>	Rstudio/DBI, RMySQL, dplyr	dbWriteTable(prod_con,"order_dim",order_dim,row.names=FALSE)
<b>Copy Time Dimension Table into MySQL's Schema Production Data</b>	Rstudio/DBI, RMySQL, dplyr	dbWriteTable(prod_con,"time_dim",time_dim,row.names=FALSE)



## 5.1 UPDATE KPIS

Purpose	Query
<b>Pick Efficiency per Product</b>	<pre>facts_table &lt;- facts_table %&gt;%   arrange(orderid_unique, date_of_order) %&gt;%   group_by(orderid_unique) %&gt;%   mutate(effi2 = as.numeric(difftime(date_of_order,     lag(date_of_order,       default = first(date_of_order)), units = "mins")))</pre>
<b>Pick Efficiency per Order</b>	<p>Create a calculated field for sum of Group Efficiency Sum  { FIXED [orderid_unique] : SUM([Pick Efficiency Kpi]) }</p> <p>create a calculated field for Average pick efficiency  WINDOW_AVG([Group Efficiency Sum])</p>