# Assignment Submission Report

## Data Mining & Warehousing Assignment 1

IT13011130 – Gunathilaka D. D. T. M.(08/22/2016)

The following report describes the procedure that has been followed while implementing data mining model for Assignment 1 of Data Mining & Warehousing module.

*Declaration*

I hear by declare that following task is my own creation and it does not violate any constraints of the assignment.

Gunathilaka D. D. T. M. (IT13011130)

## Table of Contents

## Introduction

## Case Study

In this assignment it has selected a data set which summarizes a heterogeneous set of features about articles published by Mashable in a period of two years. (www.mashable.com). The respective data set consisted, 39797 records with 60 columns. Each row is identified by the URL of the article and the dependent variable (shares) is continuously valued. The goal is to predict the number of shares in social networks (popularity). In other words, we need to find what kind of articles are attracting eyes from general population using the other 58 independent attributes. In the implemented solution it trains a model to predict the number of shares. Therefore "Regression" is used as the data mining technique.

## Background Study for Regression

Regression is a "Predictive" data-mining technique that falls under "Supervised Learning "category. In Oracle there are 2 algorithms to implement regression based data mining models. Both algorithms are particularly suited for mining data sets that have very high dimensionality (many attributes), including transactional and unstructured data. In the model build (training) process, a regression algorithm estimates the value of the target as a function of the predictors for each case in the build data. These relationships between predictors and target are summarized in a model, which can then be applied to a different data set in which the target values are unknown.

### Generalized Linear Models (GLM)

GLM is a popular statistical technique for linear modeling. Oracle Data Mining implements GLM for regression and for binary classification.

GLM provides extensive coefficient statistics and model statistics, as well as row diagnostics. GLM also supports confidence bounds.

### Support Vector Machines (SVM)

SVM is a powerful, state-of-the-art algorithm for linear and nonlinear regression. Oracle Data Mining implements SVM for regression and other mining functions.

SVM regression supports two kernels: The Gaussian kernel for nonlinear regression, and the linear kernel for linear regression. SVM also supports active learning.

Furthermore, Root Mean Squared Error and the Mean Absolute Error are commonly used statistics for evaluating the overall quality of a regression model.

### Root Mean Squared Error

The Root Mean Squared Error (RMSE) is the square root of the average squared distance of a data point from the fitted line.

This SQL expression calculates the RMSE.

*SQRT(AVG((predicted_value - actual_value) \* (predicted_value - actual_value)))*

### *Mean Absolute Error*

The Mean Absolute Error (MAE) is the average of the absolute value of the residuals (error). The MAE is very similar to the RMSE but is less sensitive to large errors.

This SQL expression calculates the MAE.

*AVG(ABS(predicted_value - actual_value))*

# *Procedure*

## Data Cleaning

Before starting the implementation, data preprocessed by adjusting data type and excluding URL as predictor a column. "URL" is also the only unique field in the data set which can be used as case id. Unnecessary data needs to be removed from the original data set, since they might cause to increase the error-rate of predictions. Therefore, rows with empty data were removed from original data set.

## Data Integration

The original data set was in .CSV format. Once it has been cleansed as next step; "ONLINENEWSPOPULARITY" table was created under Oracle user account which has data mining privileges. Since some attribute data had around 12 decimal points "NUMBER (38,20)" used for their column type to preserve data.

## Implementation

There are 2 ways to implement the models.

1. Using query
2. Using SQL developer GUI

### *1st approach*

Settings for the selected miner models are stored in a table as key-value pairs. Only the column names have to be same. (setting name, setting value)

```
CREATE TABLE miner_model_settings(setting_name VARCHAR2(30),setting_value VARCHAR2(4000));
/
```

As next step for this example, the settings of the best regression algorithm for the dataset(GLM) with feature selection with 1.3345% predictive confidence were stored in above created table.

```
BEGIN
INSERT INTO MINER_MODEL_SETTINGS VALUES(
DBMS_DATA_MINING.ALGO_NAME,DBMS_DATA_MINING.ALGO_GENERALIZED_LINEAR_MODEL
);
INSERT INTO MINER_MODEL_SETTINGS VALUES(
DBMS_DATA_MINING.GLMS_RIDGE_REGRESSION,DBMS_DATA_MINING.GLMS_RIDGE_REG_DISABLE
);
INSERT INTO MINER_MODEL_SETTINGS VALUES(
DBMS_DATA_MINING.GLMS_FTR_SELECTION,DBMS_DATA_MINING.GLMS_FTR_SELECTION_ENABLE
);
INSERT INTO MINER_MODEL_SETTINGS VALUES(
```

```
DBMS_DATA_MINING.PREP_AUTO,DBMS_DATA_MINING.PREP_AUTO_ON
);
INSERT INTO MINER_MODEL_SETTINGS
VALUES(DBMS_DATA_MINING.GLMS_DIAGNOSTICS_TABLE_NAME,'ROW_DIAGNOSTIC_STATISTICS');
COMMIT;
END;
/
```

Following query is executed to build and train the model. In below query 'case_id_column_name' is a unique key which used to identify each record separately while 'target_column_name' is the target that is required to predict by the model.
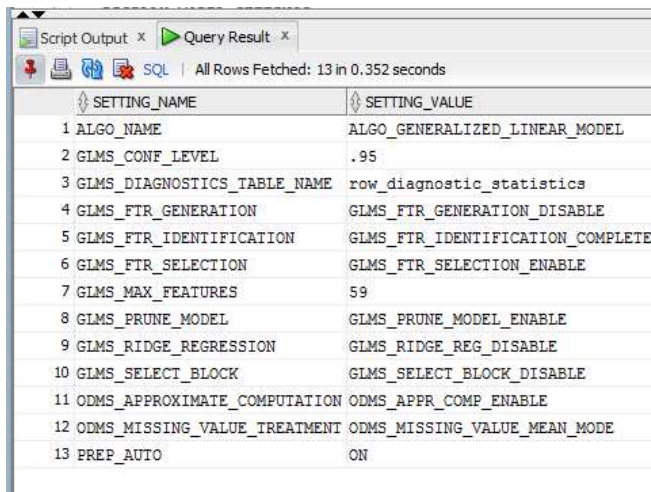
```
BEGIN
DBMS_DATA_MINING.CREATE_MODEL
(
MODEL_NAME => 'FINAL_GLM_MODEL',
MINING_FUNCTION => DBMS_DATA_MINING.REGRESSION,
DATA_TABLE_NAME    => 'ONLINENEWSPOPULARITY',
CASE_ID_COLUMN_NAME => 'URL',
TARGET_COLUMN_NAME  => 'SHARES',
SETTINGS_TABLE_NAME => 'MINER_MODEL_SETTINGS'
);
COMMIT;
END;
/
```

When PL/SQL procedure successfully completed, in order to test the created model "ONLINENEWSPOPULARITY" table is used.  Therefore, following query is executed to test the implemented model.

```
SELECT URL,SHARES, PREDICTION(FINAL_GLM_MODEL USING *) PRED FROM ONLINENEWSPOPULARITY;
```

## Model Settings:

```
COLUMN SETTING_NAME FORMAT A30
COLUMN SETTING_VALUE FORMAT A30
SELECT SETTING_NAME, SETTING_VALUE  FROM USER_MINING_MODEL_SETTINGS WHERE MODEL_NAME =
FINAL_GLM_MODEL' ORDER BY SETTING_NAME;
```

Script Output ×  ▶ Query Result ×

SQL | All Rows Fetched: 13 in 0.352 seconds

|   | SETTING_NAME | SETTING_VALUE |
|---|---|---|
| 1 | ALGO_NAME | ALGO_GENERALIZED_LINEAR_MODEL |
| 2 | GLMS_CONF_LEVEL | .95 |
| 3 | GLMS_DIAGNOSTICS_TABLE_NAME | row_diagnostic_statistics |
| 4 | GLMS_FTR_GENERATION | GLMS_FTR_GENERATION_DISABLE |
| 5 | GLMS_FTR_IDENTIFICATION | GLMS_FTR_IDENTIFICATION_COMPLETE |
| 6 | GLMS_FTR_SELECTION | GLMS_FTR_SELECTION_ENABLE |
| 7 | GLMS_MAX_FEATURES | 59 |
| 8 | GLMS_PRUNE_MODEL | GLMS_PRUNE_MODEL_ENABLE |
| 9 | GLMS_RIDGE_REGRESSION | GLMS_RIDGE_REG_DISABLE |
| 10 | GLMS_SELECT_BLOCK | GLMS_SELECT_BLOCK_DISABLE |
| 11 | ODMS_APPROXIMATE_COMPUTATION | ODMS_APPR_COMP_ENABLE |
| 12 | ODMS_MISSING_VALUE_TREATMENT | ODMS_MISSING_VALUE_MEAN_MODE |
| 13 | PREP_AUTO | ON |

## *Model Signature(Attributes):*

*COLUMN ATTRIBUTE_NAME FORMAT A40*
*COLUMN ATTRIBUTE_TYPE FORMAT A20*
*SELECT ATTRIBUTE_NAME, ATTRIBUTE_TYPE  FROM USER_MINING_MODEL_ATTRIBUTES WHERE MODEL_NAME =*
*'FINAL_GLM_MODEL'ORDER BY ATTRIBUTE_NAME;*

| | ATTRIBUTE_NAME | ATTRIBUTE_TYPE |
|---|---|---|
| 1 | AVERAGE_TOKEN_LENGTH | NUMERICAL |
| 2 | AVG_NEGATIVE_POLARITY | NUMERICAL |
| 3 | DATA_CHANNEL_IS_ENTERTAINMENT | NUMERICAL |
| 4 | KW_AVG_AVG | NUMERICAL |
| 5 | KW_MAX_AVG | NUMERICAL |
| 6 | KW_MIN_AVG | NUMERICAL |
| 7 | LDA_02 | NUMERICAL |
| 8 | LDA_03 | NUMERICAL |
| 9 | NUM_HREFS | NUMERICAL |
| 10 | N_TOKENS_TITLE | NUMERICAL |
| 11 | SELF_REFERENCE_MIN_SHARES | NUMERICAL |
| 12 | SHARES | NUMERICAL |
| 13 | TIMEDELTA | NUMERICAL |

## *Global statistics:*

*SELECT *  FROM TABLE(DBMS_DATA_MINING.GET_MODEL_DETAILS_GLOBAL('FINAL_GLM_MODEL'))ORDER BY*
*GLOBAL_DETAIL_NAME;*

| GLOBAL_DETAIL_NAME | GLOBAL_DETAIL_VALUE |
|---|---|
| ADJUSTED_R_SQUARE | 0.0232352706043227 |
| AIC | 632415.191258281 |
| COEFF_VAR | 327.051975283298 |
| CORRECTED_TOTAL_DF | 33952 |
| CORRECTED_TOT_SS | 4267394094335.55 |
| DEPENDENT_MEAN | 3387.87129266928 |
| ERROR_DF | 33940 |
| ERROR_MEAN_SQUARE | 122768615.627308 |
| ERROR_SUM_SQUARES | 4166766814390.85 |
| F_VALUE | 68.3041559050875 |
| GMSEP | 122815640.724438 |
| HOCKING_SP | 3617.33155447445 |
| J_P | 122815621.547349 |
| MODEL_CONVERGED | 1 |
| MODEL_DF | 12 |
| MODEL_F_P_VALUE | 0 |
| MODEL_MEAN_SQUARE | 8385606662.05944 |
| MODEL_SUM_SQUARES | 100627279944.713 |
| NUM_PARAMS | 13 |
| NUM_ROWS | 33953 |
| ROOT_MEAN_SQ | 11080.0999827307 |
| R_SQ | 0.0235804984775775 |
| SBIC | 632524.816780708 |
| TERMINATION | 0 |
| VALID_COVARIANCE_MATRIX | 1 |

## GLM statistics | Coefficient statistics | Features and their p_values:

```
SELECT * FROM TABLE (DBMS_DATA_MINING.GET_MODEL_DETAILS_GLM('FINAL_GLM_MODEL'));
SELECT * FROM TABLE(DBMS_DATA_MINING.GET_MODEL_DETAILS_GLM(MODEL_NAME =>
'FINAL_GLM_MODEL'));

SET LINE 120
COLUMN FEATURE_EXPRESSION FORMAT A53
SELECT FEATURE_EXPRESSION, COEFFICIENT, STD_ERROR, TEST_STATISTIC,
  P_VALUE, STD_COEFFICIENT, LOWER_COEFF_LIMIT, UPPER_COEFF_LIMIT
  FROM TABLE(DBMS_DATA_MINING.GET_MODEL_DETAILS_GLM('FINAL_GLM_MODEL'));

SET LIN 80
SET PAGES 20
SELECT FEATURE_EXPRESSION, COEFFICIENT, P_VALUE
  FROM TABLE(DBMS_DATA_MINING.GET_MODEL_DETAILS_GLM('FINAL_GLM_MODEL'))
  ORDER BY P_VALUE;
```



## Validation:

Root Mean Square Error - $\mathrm{Sqrt(Mean((x - x')^2))}$

Mean Absolute Error - $\mathrm{Mean(|(x - x')|)}$

```
 COLUMN RMSE FORMAT 9999.99
 COLUMN MAE FORMAT 9999.99
SELECT SQRT(AVG((A.PRED - B.SHARES) * (A.PRED - B.SHARES))) RMSE,
     AVG(ABS(A.PRED - B.SHARES)) MAE
  FROM (SELECT URL,SHARES, PREDICTION(FINAL_GLM_MODEL USING *) PRED
     FROM ONLINENEWSPOPULARITY) A,
     ONLINENEWSPOPULARITY B
  WHERE A.URL = B.URL;
```

| | RMSE | MAE |
|---|---|---|
| 1 | 11502.5333274091471849189618877368572007 | 3051.5554104560556301541847442358995056 |

## 2nd approach

In the second method of model implementation a graphical user interface provided by Oracle SQL-Developer tool is used. Once the tool is opened following steps are followed to implement Classification mining model.

### Steps

Create new project and work flow in created data miner user account under Oracle Data Miner. Name the project as "Online News Popularity" (any name) and define the workflow as "Popularity by Number of Shares" (any name)



**Figure 1 Create Project & workflow**

Next drag and drop  from the tool palette to the workflow which is under "Data" section of palette. Once the "Data Source" node is added to workflow it will open a window to select a schema. This schema would be the table which provides training data to the model. Therefore, in our case select "DMUSER. ONLINENEWSPOPULARITY" as the schema. Since we have cleaned irrelevant columns previously all the existing columns of the table are required to build the model. Therefore, click "Finish"



**Figure 2 Select Data Source**

As the next step the model node should be created. Therefore drag and drop  node from tool palette which is under "Models" section. Then link the two nodes (  )by using

 which is under "Linking Nodes" category or right click →  connect.

Next double-click or right click → Edit and open edit window of "Regress Build" (i.e. Regression) node. Then use following setting:



**Figure 3 Add, Delete, Set Target and Case ID**



**Figure 4 Oracle will evaluate all attributes to determine which are important**



**Figure 5 Advance settings for best confidence in dataset**

Double Click on the model for advanced settings:

Press Ok twice to go back to work flow.

Now right click on regress build and click on Run button.

Once the model has successfully executed drag and drop "Apply" node and new "Data-Source" node from tool palette.

Next double click on newly added "Data-Source" node and "ONLINENEWSPOPULARITY" table in data source node.

In apply node:



**Figure 6 Set case ID for quicker predictions**



**Figure 7 Set Column order**



**Figure 8 Add Case Id as reference, actual shares to compare with predictions**

Next re-run the entire workflow. This will apply implemented all algorithms to the records of "ONLINENEWSPOPULARITY" table. Now models are ready to provide predictions.

In order to display predictions a separate table called "REGRESSION_RESULT" is created. This table keeps the predictions, probability of those predictions and relevant case Ids made by models.

The implemented work-flow appears as below.



**Figure 9 Selected Model for deploy**


The accuracy of the models can be visualized by selecting "Compare Test Results" option in menu which appears once right click on "Class Build" node.

Right click on apply node → deploy → select node, dependent node and children node to generate query for the model.

In model settings:

Missing value treatment set to "Delete Row" to remove rows with null value

Shares (Independent variable) and URL not used in prediction



**Figure 10 URL and Shares attributes are not used**

## *Results & Discussion*

The overall performances of implemented models are shown below.

Models

| Name | Predictive Confidence % | Mean Absolute Error | Root Mean Square Error | Mean Predicted Value | Mean Actual Value | Algorithm | Creation Date |
|---|---|---|---|---|---|---|---|
| SVM | 0.6142 | 2,840.0664 | 9,156.4713 | 3,075.1412 | 3,318.8128 | Support Vector Machine | 8/21/16 12:48 AM |
| SVM_ACTIVE | 0 | 5,830.7154 | 9,943.2918 | 7,164.4152 | 3,318.8128 | Support Vector Machine | 8/21/16 12:46 AM |
| SVM_GAUSSIAN | 0.6142 | 2,840.0664 | 9,156.4713 | 3,075.1412 | 3,318.8128 | Support Vector Machine | 8/21/16 12:48 AM |
| SVM_GAUSSIAN_ACTIVE | 0 | 5,830.7154 | 9,943.2918 | 7,164.4152 | 3,318.8128 | Support Vector Machine | 8/21/16 12:46 AM |
| SVM_LINEAR | 0.318 | 2,865.0521 | 9,183.7597 | 3,053.3908 | 3,318.8128 | Support Vector Machine | 8/21/16 12:46 AM |
| SVM_LINEAR_ACTIVE | 0 | 2,307.5263 | 9,313.4355 | 1,664.958 | 3,318.8128 | Support Vector Machine | 8/21/16 12:46 AM |

Models

| Name | Predictive Confidence % | Mean Absolute Error | Root Mean Square Error | Mean Predicted Value | Mean Actual Value | Algorithm | Creation Date |
|---|---|---|---|---|---|---|---|
| GLM_GENERATION1 | 0 | 3,327.3795 | 10,932.7822 | 3,512.3058 | 3,318.8128 | Generalized Linear Model | 8/21/16 9:31 AM |
| GLM_RIGED1 | 1.2362 | 2,991.066 | 9,099.1658 | 3,458.7539 | 3,318.8128 | Generalized Linear Model | 8/21/16 1:05 AM |
| GLM_SELECTION1 | 1.3345 | 3,000.754 | 9,090.106 | 3,466.5963 | 3,318.8128 | Generalized Linear Model | 8/21/16 1:08 AM |



According to the above result SVM with linear kernel has the lowest spread and relatively highest accuracy. But when favoring the predictive confidence, we have to pick GLM with feature selection (confidence 0.318< 1.3345) as the best model.

**Figure 11 Prediction value distribution per post**



**Figure 12 Overall Prediction value Average**

## *Integration*

The implemented models can be integrated with software components. Word form with VBA backend implemented for data insertion. Diagrams and visualization are displayed through an Excel instance but all the data manipulation done in the oracle server. Power BI used to visualize spread of training data.