# Applying Cluster Analysis on SOT Oil and Gas Production Data

By

Uma Bala Devarakonda
2-2012-203

A thesis submitted to the Faculty of

ENGINEERING AND INFORMATION TECHNOLOGY

In Partial Fulfillment of the Requirements for the

MASTER'S DEGREE
in
INFORMATION TECHNOLOGY

SWISS GERMAN UNIVERSITY
EduTown BSD City
Tangerang 15339
Indonesia

January 2014

# Applying Cluster Analysis on SOT Oil and Gas Production Data

By

Uma Bala Devarakonda
2-2012-203

A thesis submitted to the Faculty of

ENGINEERING AND INFORMATION TECHNOLOGY

In Partial Fulfillment of the Requirements for the

MASTER'S DEGREE
in
INFORMATION TECHNOLOGY

SWISS GERMAN UNIVERSITY
EduTown BSD City
Tangerang 15339
Indonesia

**Revision after the Thesis defense on 18-02-2014**

## STATEMENT BY THE AUTHOR

I hereby declare that this submission is my own work and to the best of my knowledge, it contains no material previously published or written by another person, nor material which to a substantial extent has been accepted for the award of any other degree or diploma at any educational institution, except where due acknowledgement is made in the thesis.

**Uma Bala Devarakonda**

_____

Student                                                                        Date

Approved by:

**Dr. Ir. Mohammad A. Amin Soetomo, M.Sc.**

_____

Thesis Advisor                                                            Date

**Dr. Harya Widiputra**

_____

Thesis Co-Advisor                                                      Date

**Dr. Ir. Gembong Baskoro, M.Sc.**

Dean of Engineering and Information Technology            Date
Faculty

# ABSTRACT

## Applying Cluster analysis on SOT Oil and Gas Production Data

By

Uma Bala Devarakonda
Dr. Ir. Mohammad A. Amin Soetomo, M.Sc.
Dr. Harya Widiputra

SWISS GERMAN UNIVERISTY

SKK Migas has implemented SOT (Sistem Operasi Terpadu) to allow SKK Migas system to retrieve specific information from PSC Contractor data source within an integrated online system. It is expected that by the completion of SOT program, there can be huge transactional data that is exchanged on daily basis, and without Data Mining it will be extremely difficult to perform any data analysis. The main purpose of this research is to find out if we can apply clustering technique on SOT production data collected over a period of time. The entire process of the data mining analysis has been carried out using the CRISP-DM methodology. This research found out that it is possible to find certain interesting clusters applying clustering algorithm on the production data. These clusters can be interpreted as clusters of PSC contractors/oil fields/wells clustered based on their various efficiency ratios.

*Keywords: Clustering, Data mining, SOT production data, Efficiency ratio, CRISP-DM.*

Uma Bala Devarakonda

# ACKNOWLEDGEMENTS

# TABLE OF CONTENTS

Page

Uma Bala Devarakonda

## LIST OF FIGURES

Figures                                                              Page

## LIST OF TABLES