

An Oracle White Paper
February 2012

Oracle Data Mining 11g Release 2

Competing on In-Database Analytics

Disclaimer

The following is intended to outline our general product direction. It is intended for information purposes only, and may not be incorporated into any contract. It is not a commitment to deliver any material, code, or functionality, and should not be relied upon in making purchasing decisions. The development, release, and timing of any features or functionality described for Oracle's products remains at the sole discretion of Oracle.

Executive Overview	1
Oracle Recognized as a Leader in Data Mining	2
In-Database Data Mining	2
Key Benefits of Oracle Data Mining	4
Introduction	5
Oracle Data Mining	5
Data Mining Deep Dive	7
Exadata and Oracle Data Mining	9
Oracle Data Mining for Data Analysts	16
Oracle Data Mining for Applications Developers	18
Competing on In-Database Analytics	20
Beyond a Tool; Enabling Predictive Applications	21
Spend Less	24
Eliminate Redundant Data and Traditional Analytical Servers	25
Conclusion	25

"Some companies have built their very businesses on their ability to collect, analyze, and act on data. ...Although numerous organizations are embracing analytics, only a handful have achieved this level of proficiency. But analytics competitors are the leaders in their varied fields—consumer products finance, retail, and travel and entertainment among them."

—Tom Davenport, author, *Competing on Analytics*

Executive Overview

Oracle Data Mining provides powerful data mining functionality within the Oracle Database. It enables you to discover new insights hidden in your data and to leverage your investment in Oracle Database technology. With Oracle Data Mining, you can build and apply predictive models that help you target your best customers, develop detailed customer profiles, and find and prevent fraud. Oracle Data Mining, a component of the Oracle Advanced Analytics Option, helps your company better *compete on analytics*.

Oracle Recognized as a Leader in Data Mining

According to a February 2010 report from independent analyst firm Forrester Research, Oracle is a leader in predictive analytics and data mining (PA/DM). [“The Forrester Wave™: Predictive Analytics And Data Mining Solutions, Q1 2010,”](#) written by Senior Analyst James G. Kobiulus, states that “Oracle provides a PA/DM solution portfolio that is built into its own widely adopted DBMS, DW, data integration, and BI platforms, with a wide range of prepackaged predictive applications, and it provides a powerful assortment of algorithms for mining complex structured and unstructured information types.”

Excerpting further from the Forrester Wave, “Oracle focuses on in-database mining in the Oracle Database, on integration of Oracle Data Mining into the kernel of that database, and on leveraging that technology in Oracle’s branded applications. Oracle’s key current offering differentiators are DBMS-integrated data-preparation tools; semistructured and unstructured information integration; text analytics; strategy maps; ensemble modeling; champion-challenger modeling; sentiment analysis; social network analysis; and support for a competitive range of statistical algorithms and variable selection techniques.”

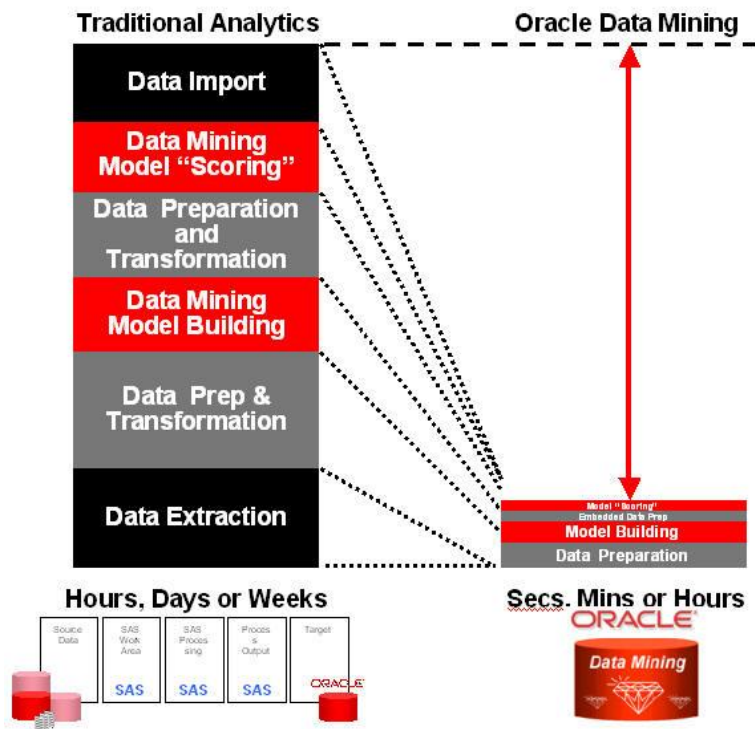
In-Database Data Mining

With Oracle Data Mining, a component of the Oracle Advanced Analytics Option, everything occurs in the Oracle Database—in a single, secure, scalable platform for advanced business intelligence. Oracle Data Mining represents a breakthrough in business intelligence. In contrast to traditional statistical software that requires data extraction to separate servers, which may be insecure and costly to maintain, Oracle Data Mining embeds a wide-range of mining functions inside the database—where the data is stored. Coupled with the power of SQL, Oracle Data Mining eliminates data movement and duplication, maintains security and minimizes latency time from raw data to valuable information.

Oracle Data Mining enables you to:

- Leverage your data to discover patterns and valuable new insights
- Build and apply predictive models and embed them into dashboards and applications
- Save money. Oracle Data Mining costs significantly less than traditional statistical software. A feature of the Oracle Database, ODM reduces the total cost of ownership.

In-Database Data Mining



Oracle Data Mining eliminates data movement, data duplication and security exposures.

Key Benefits of Oracle Data Mining

Oracle Data Mining, running natively inside the SQL kernel of the Oracle Database provides users with the following benefits:

- Mines data inside Oracle Database, an industry leader in performance and reliability
- Eliminates data extraction and movement
- Provides a platform for analytics-driven database applications
- Provides increased security by leveraging database security options
- Delivers lowest total cost of ownership (TCO) compared to traditional data mining vendors
- Leverages 30+ years of experience of ever advancing Oracle Database technology

Oracle Data Mining enables you to go beyond standard BI and OLAP tools that answer questions like: “Who are my top customers?” “What products have sold the most?” and “Where are costs the highest?” Data mining automatically sifts through data and reveals patterns and insights that help you run your business better. In today’s competitive marketplace, your company must manage its most valuable asset — its data. Moreover, your company must exploit its data for competitive advantage. If you don’t, your competitors will. With Oracle Data Mining, you can implement strategies to:

- Understand and target select customer groups
- Develop detailed customer profiles for building marketing campaigns
- Anticipate and prevent customer churn and attrition
- Identify promising cross-sell and up-sell opportunities
- Detect noncompliance and potential fraud
- Discover new clusters or customer segments
- Perform market-basket analysis to find frequently co-occurring items

"Simply put, data mining is used to discover [hidden] patterns and relationships in your data in order to help you make better business decisions."

-- Herb Edelstein, Two Crows Corporation

Introduction

Traditional business intelligence (BI) reporting tools report on what has happened in the past. OLAP provides rapid drill-through for more detailed information, roll up to aggregate information, and sometimes aggregate level forecasting. With a good BI or OLAP tool, a good analyst and enough time, you could eventually find the information you want. —But eventually can mean a very long time. According to the infinite monkey theorem, a monkey hitting keys at random will eventually type the complete works of Shakespeare.

Data Mining is now possible due to advances in computer science and machine learning. Data Mining delivers new algorithms that can automatically sift deep into your data at the individual record level to discover patterns, relationships, factors, clusters, associations, profiles, and predictions—that were previously “hidden”.

Oracle Data Mining, a collection of machine learning algorithms embedded in the Oracle Database, allows you to discover new insights, segments and associations, to make more accurate predictions, to find key variables, to detect anomalies, and to extract more information from your data. For example, by analyzing your best customers, ODM can discover profiles and embed predictive analytics in applications that identify customers who are likely to be your best customers. These customers may not represent your most valuable customers today, but they match profiles of your current best customers. With ODM you can apply predictive models to generate reports and dashboards that reveal the most promising customers for your marketing and sales departments, or real-time predictions for call center personnel. Knowing the “strategic value” of your customers — which customers are likely to become profitable customers in the future and which are not, or predicting which customers are likely to churn or likely to respond to a marketing offer — and integrating this information at just the right time into your operations is the key to successfully competing on analytics.

Oracle Data Mining

Oracle Data Mining, a component of the Oracle Advanced Analytics Option, delivers a wide range of cutting edge machine learning algorithms inside the Oracle Database. Since Oracle Data Mining functions reside natively in the Oracle Database kernel, they deliver unparalleled performance, scalability and security. The data and data mining functions never leave the database to deliver a comprehensive in-database processing solution.

Oracle Data Mining Release 11g Release 2 supports twelve in-database mining algorithms that address classification, regression, association rules, clustering, attribute importance, and feature selection problems. Working with Oracle Text (which uses standard SQL to index, search, and analyze text and documents stored in the Oracle database, in files, and on the web), many ODM mining functions can mine both structured and unstructured (text) data.

ODM provides PL/SQL and SQL application programming interfaces (APIs) for model building and model scoring functions. An optional Oracle Data Miner graphical user interface (GUI) that is available for download from OTN is available for data analysts who want to use a point and click GUI. The Oracle Spreadsheet Add-In for Predictive Analytics implements a predictive analytics PL/SQL package within Microsoft Excel. The Add-In enables automated predictive analytics functionality within a spreadsheet environment. The ODM graphical user interface(s) and APIs provide an analytical platform for data analysts and application developers to deliver data mining's results to BI dashboards and enterprise applications.

Ranges of Business Intelligence

Now let's describe what data mining is and how it both differs from and complements other business intelligence (BI) products — query and reporting, OLAP, and statistical tools. Let's also look at some common definitions of business intelligence tools.

BI, OLAP and Data Mining

BI	OLAP	Data Mining
Extraction of detailed and roll up data	Summaries, trends and forecasts	Knowledge discovery of hidden patterns
<i>"Information"</i>	<i>"Analysis"</i>	<i>"Insight & Prediction"</i>
Who purchased mutual funds in the last 3 years?	What is the average income of mutual fund buyers, by region, by year?	Who will buy a mutual fund in the next 6 months and why?

Figure 1. BI, OLAP, Statistics and Data Mining. Oracle Data Mining differs from query and reporting (BI), and OLAP tools by automatically discovering new information that was previously hidden in the data and the ability to make predictions.

BI and query and reporting tools help you to get information out of your database or data warehouse. These tools are good at answering questions such as “Who purchased a mutual fund in the past 3 years?”

OLAP tools go beyond basic BI and allow users to rapidly and interactively drill-down for more detail, comparisons, summaries and forecasts. OLAP is good at drill-downs into the details to find, for example, “What is the average income of mutual fund buyers by year by region?”

Statistical tools are used to draw conclusions from representative samples taken from larger amounts of data. Statistical tools are useful for finding patterns and correlations in “small to medium” amounts of data, but when the amount of data begins to overwhelm the tool, traditional statistical techniques struggle. Because statistical tools cannot analyze all the data, they force data analysts to use representative samples of the data and to eliminate input variables from the analysis. However, eliminating input variables and using small samples of data, makes you throw away valuable information.

Oracle Data Mining runs in the kernel of the Oracle Database and doesn’t suffer from the same limitations. Oracle Data Mining uses machine-learning techniques developed in the last decade to automatically find patterns and relationships hidden in the data. Oracle Data mining goes deep into the data and finds patterns from the data. Oracle Data Mining is good at providing detailed insights and making individual predictions, such as “Who is likely to buy a mutual fund in the next six months and why?”

Data Mining Deep Dive

Oracle Data Mining provides a collection of data mining algorithms that are designed to address a wide-range of business and technical problems. Different algorithms are good at different types of analysis. Oracle Data Mining supports classification, regression, clustering, associations, attribute importance and feature extraction problems. (For a full listing of algorithms, see Oracle Data Mining 11g Release 2 Algorithms table on page 9.)

Sequence for Determining Necessary Data. *Wrong:* Catalog everything you have, and decide what data is important.

Right: Work backward from the solution, define the problem explicitly, and map out the data needed to populate the investigation and models.

—James Taylor with Neil Raden, authors, *Smart (Enough) Systems*

The Data Mining Process

Before we start mining any data, we need to define the problem we want to solve and, most importantly, gather the right data to help us find the solution. If we don’t have the right data, we need to get it. If data mining is not properly approached, there is the possibility of “garbage in—garbage out”. To be effective in data mining, you will typically follow a four-step process:

Defining the Business Problem

This is the most important step. In this step, a domain expert determines how to translate an abstract business objective such as “How can I sell more of my product to customers?” into a more tangible and useful data mining problem statement such as “Which customers are most likely to purchase product A?” To build a model that predicts who is most likely to buy product A, we first must acquire data that describes the customers who have purchased product A in the past. Then we can begin to prepare the data for mining.

Gathering and Preparing the Data

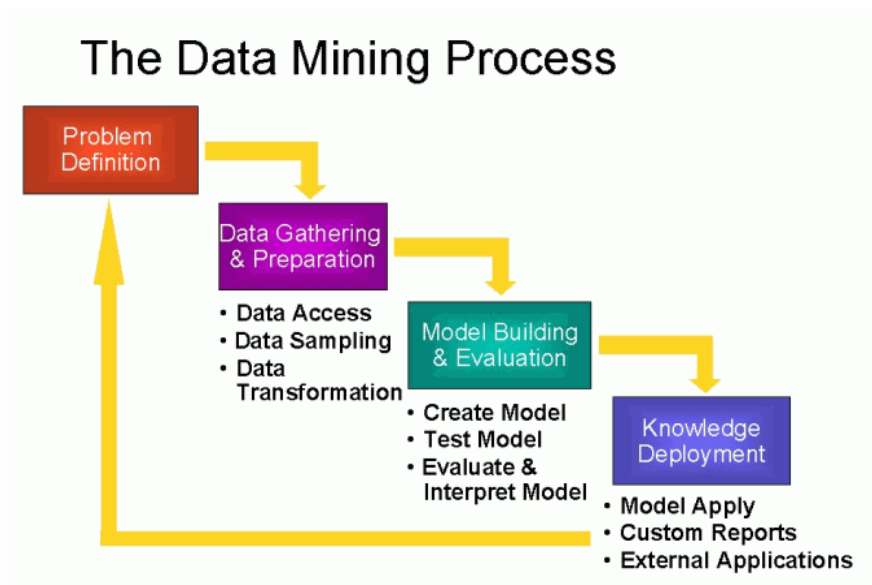
Now we take a closer look at our data and determine what additional data may be necessary to properly address our business problem. We often begin by working with a reasonable sample of the data. For example, we might examine several hundred of the many thousands, or even millions, of cases by looking at statistical summaries and histograms. We may perform some data transformations to attempt to tease the hidden information closer to the surface for mining. For example, we might transform a “Date_of_Birth” field into an “AGE” field, and we might derive new field such as “No_Times_Amt_Exceeds_N” from existing fields. The power of SQL simplifies this process.

Model Building and Evaluation

Now we are ready to build models that sift through the data to discover patterns. Generally, we will build several models, each one using different mining parameters, before we find the best or most useful model(s).

Knowledge Deployment

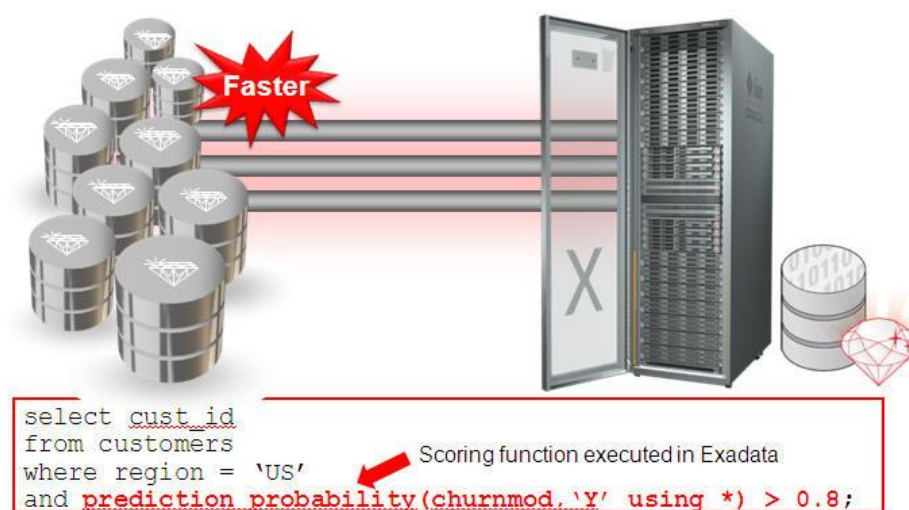
Once ODM has built a model that models relationships found in the data, we will deploy it so that users, such as managers, call center representatives, and executives, can apply it to find new insights and generate predictions. ODM’s embedded data mining algorithms eliminate any need to move (rewrite) the models to the data in the database or to extract huge volumes of unscored records for scoring using a predictive model that resides outside of the database. Oracle Data Mining provides the ideal platform for building and deploying advanced business intelligence applications.



The data mining process involves a series of steps to define a business problem, gather and prepare the data, build and evaluate mining models, and apply the models and disseminate the new information.

Exadata and Oracle Data Mining

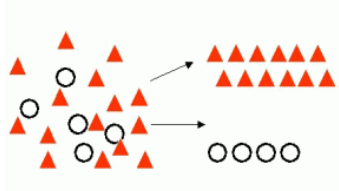

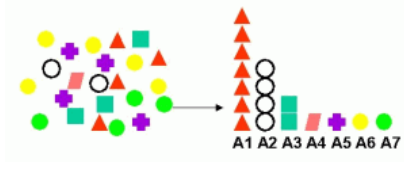
Oracle Exadata is a family of high performance storage software and hardware products that can improve data warehouse query performance by a factor of 10X or more. Oracle Data Mining scoring functions in Oracle Database 11g Release 2 score in the storage layer and permit very large data sets to be mined very quickly, thus further increasing the competitive advantage already gained from Oracle's in-database analytics.

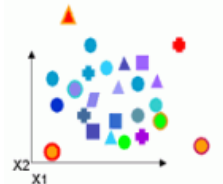
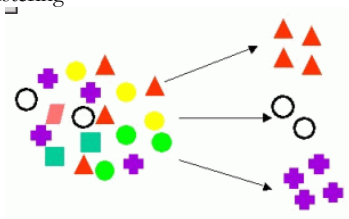
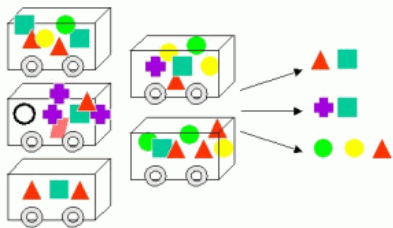
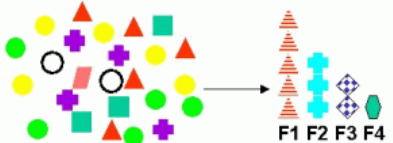


A Wide Range of Cutting Edge Algorithms

Oracle Data Mining provides a broad suite of data mining techniques and algorithms to solve many types of business and technical problems:

ORACLE DATA MINING 11G RELEASE 2 ALGORITHMS

TECHNIQUE	APPLICABILITY	ALGORITHMS
<p>Classification</p> 	<p>Classification techniques use historical data to build models that can be used to classify new data and make predictions about class membership (e.g. 0 or 1) or class value (numerical value).</p>	<p>Logistic Regression (GLM)—classic statistical technique inside the Oracle Database. Supports thousands of input variables, text and transactional data</p> <p>Naive Bayes—Fast, simple, commonly applicable</p> <p>Support Vector Machine—Cutting edge algorithm that supports binary and multi-class problems. SVMs also excel at handling shallow, yet wide, and nested data problems e.g. transaction data or gene expression data analysis. Supports text mining use cases.</p> <p>Decision Tree—Popular algorithm useful for many classification problems that that can help explain the model's logic using human-readable “If... Then...” rules.</p>
<p>Regression</p> 	<p>Technique for predicting a continuous numerical outcome such as customer lifetime value, house value, process yield rates.</p>	<p>Multiple Regression (GLM)—classic statistical technique available inside the Oracle Database Supports thousands of input variables, text and transactional data</p> <p>Support Vector Machine — Cutting edge algorithm that supports regression problems. Supports text mining and transactional data use cases.</p>
<p>Attribute Importance</p> 	<p>Ranks attributes according to strength of relationship with target attribute, for example finding factors associated with people who respond to an offer.</p>	<p>Minimum Description Length—Considers each attribute as a simple predictive model of the target class. Attribute Importance algorithm finds the attributes that have the most influence on a target attribute.</p>
<p>Anomaly Detection</p>	<p>Identifies unusual or</p>	<p>One-Class Support Vector</p>

	<p>suspicious cases based on deviation from the norm. Common examples include health care fraud, expense report fraud, and tax compliance.</p>	<p>Machine —Unsupervised learning technique that trains on “normal cases” to build a model. Then when applied, it flags unusual cases with the probability that they are not “normal”.</p>
<p>Clustering</p> 	<p>Useful for exploring data and finding natural groupings within the data. Members of a cluster are more like each other than they are like members of a different cluster. Common examples include finding new customer segments, and life sciences discovery.</p>	<p>Enhanced K-Means—Supports text mining, hierarchical clustering, distance based</p> <p>Orthogonal Partitioning Clustering—Hierarchical clustering, density based</p>
<p>Association</p> 	<p>Finds rules associated with frequently co-occurring items, used for market basket analysis, cross-sell, and root cause analysis. Useful for product bundling, in-store placement, and defect analysis.</p>	<p>Apriori—Industry standard for market basket analysis.</p>
<p>Feature Extraction</p> 	<p>Produces new attributes as linear combination of existing attributes. Applicable for text data, latent semantic analysis, data compression, data decomposition and projection, and pattern recognition.</p>	<p>Non-negative Matrix Factorization (NMF)— Creates new attributes that represent the same information using fewer attributes</p>

Supervised Learning Algorithms

Most data mining algorithms can be separated into supervised learning and unsupervised learning data mining techniques. Supervised learning requires the data analyst to identify a target attribute or dependent variable with examples of the possible classes (e.g., 0/1, Yes/No, High/Med, Low, etc.). The supervised-learning technique then sifts through data trying to find patterns and relationships among the independent attributes (predictors) that can help separate the different classes of the dependent attribute.

For example, let's say that we want to build a predictive model that can help our Marketing and Sales departments focus on people who are most likely interested in purchasing a new car. The target attribute will be a column that designates whether each customer has purchased a car—for example, a “1” for yes and a “0” for no. The supervised data mining algorithm sifts through the data searching for patterns and builds a data mining model that captures the relationships found in the data. Typically, for supervised learning, the data is separated into two parts — one for model training and another hold out sample for model testing and model evaluation. Because we already know the outcome — who purchased a car and who hasn't — we can apply our ODM predictive model to our hold out sample to evaluate the model's accuracy and make decisions about the usefulness of the model. ODM models with acceptable prediction capability can have high economic value. Binary and multi-class classification problems represent a majority of common business challenges addressed through Oracle Data Mining, including database marketing, response and sales offers, fraud detection, profitability prediction, customer profiling, credit rating, churn anticipation, inventory requirements, failure anticipation, and many others. Oracle Data Mining also provides utilities for evaluating models in terms of model accuracy and “lift” — or the incremental advantage of the predictive model over the naïve guess.

Naïve Bayes

Naïve Bayes (NB) is a supervised-learning technique for classification and prediction supported by Oracle Data Mining. The Naive Bayes algorithm is based on conditional probabilities. It uses Bayes' Theorem, a formula that calculates a probability by counting the frequency of values and combinations of values in the historical data. Bayes' Theorem finds the probability of an event occurring given the probability of another event that has already occurred. If B represents the dependent event and A represents the prior event, Bayes' theorem can be stated as follows.

Bayes' Theorem: $\text{Prob}(B \text{ given } A) = \text{Prob}(A \text{ and } B) / \text{Prob}(A)$

To calculate the probability of B given A, the algorithm counts the number of cases where A and B occur together and divides it by the number of cases where A occurs alone.

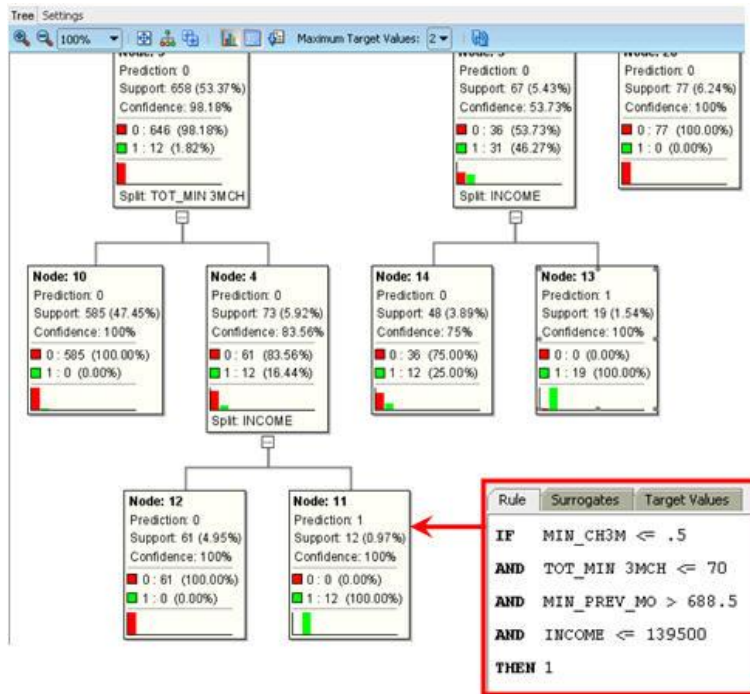
After ODM builds a NB model, the model can be used to make predictions. Application developers can integrate ODM models to classify and predict for a variety of purposes, such as:

- Identify customers likely to purchase a certain product or to respond to a marketing campaign
- Identify customers most likely to spend greater than \$3,000
- Identify customers likely to churn

NB affords fast model building and scoring and can be used for both binary and multi-class classification problems. NB cross-validation, supported as an optional way to run NB, permits the user to test model accuracy on the same data that was used to build the model, rather than building the model on one portion of the data and testing it on a different portion. Not having to hold aside a portion of the data for testing is especially useful if the amount of build data is relatively small.

Decision Trees

Oracle Data Mining supports the popular Classification Tree algorithm. The ODM Decision Tree model contains complete information about each node, including Confidence, Support, and Splitting Criterion. The full Rule for each node can be displayed, and in addition, a surrogate attribute is supplied for each node, to be used as a substitute when applying the model to a case with missing values.



Support Vector Machines

ODM's Support Vector Machines (SVM) algorithm supports binary and multi-class classification, prediction, and regression models, that is, prediction of a continuous target attribute. SVMs are particularly good at discovering patterns hidden in problems that have a very large number of independent attributes, yet have only a very limited number of data records or observations.

SVM models can be used to analyze genomic data with only 100 patients who have thousands of gene expression measurements for each patient. SVMs can build models that predict disease treatment outcome based on genetic profiles.

Generalized Linear Models (Logistic and Multiple Regression)

ODM 11g Release 2 adds support for the multipurpose classical statistical algorithm, Generalized Linear Models (GLM). ODM supports, as two mining functions: classification (binary Logistic Regression) and regression (Multivariate Linear Regression). GLM is a parametric modeling technique. Parametric models make assumptions about the distribution of the data. When the

assumptions are met, parametric models can be more efficient than non-parametric models. Oracle Data Mining's GLM implementation provides extensive model quality diagnostics and predictions with confidence bounds.

Oracle Data Mining supports ridge regression for both regression and classification mining functions. ODM's GLM automatically uses ridge if it detects singularity (exact multicollinearity) in the data. ODM supports GLM with the added capability to handle many hundreds to thousands of input attributes. Traditional external statistical software packages typically are limited to 10-30 input attributes.

Attribute Importance

Oracle Data Mining's Attribute Importance algorithm helps to identify the attributes that have the greatest influence on a target attribute. Often, knowing which attributes are most influential helps you to better understand and manage your business and can help simplify modeling activities. Additionally, these attributes can indicate the types of data that you may wish to add to your data to augment your models.

Attribute Importance can be used to find the process attributes most relevant to predicting the quality of a manufactured part, the factors associated with churn, or the genes most likely related to being involved in the treatment of a particular disease.

Unsupervised Learning Algorithms

In unsupervised learning, the user does not specify a target attribute for the algorithm. Unsupervised learning techniques, such as associations and clustering algorithms, make no assumptions about a target field. Instead, they allow the data mining algorithm to find associations and clusters in the data independent of any a priori defined business objective.

Clustering

Oracle Data Mining provides two algorithms, Enhanced k-Means and Orthogonal Partitioning Clustering (O-Cluster), for identifying naturally occurring groupings within a data population. ODM's Enhanced k-Means (EKM) and O-Cluster algorithms support identifying naturally occurring groupings within the data population. ODM's EKM algorithm supports hierarchical clusters, handles numeric and categorical attributes and will cut the population into the user specified number of clusters.

ODM's O-Cluster algorithm handles both numeric and categorical attributes and will automatically select the best cluster definitions. In both cases, ODM provides cluster detail information, cluster rules, cluster centroid values, and can be used to "score" a population on their cluster membership. For example, Enhanced k-Means Clustering can be used to find new customer segments or to reveal subgroups within a diseased population.

Association Rules (Market Basket Analysis)

ODM's Association Rules (AR) finds co-occurring items or events within the data. Often called “market basket analysis”, AR counts the number of combinations of every possible pair, triplet, quadruplet, etc., of items to find patterns. Association Rules represent the findings in the form of antecedents and consequents. An AR rule, among many rules found, might be “Given Antecedents Milk, Bread, and Jelly, then Consequent Butter is also expected with Confidence 78% and Support 12%. Translated in simpler English, this means that if you find a market basket having the first three items, there is a strong chance (78% confidence) that you will also find the fourth item and this combination is found in 12% of all the market baskets studied. The associations or “rules” thus discovered are useful in designing special promotions, product bundles, and store displays.

AR can be used to find which manufactured parts and equipment settings are associated with failure events, what patient and drug attributes are associated with which outcomes or which items or products a person who has purchased item A most likely to buy?

Anomaly Detection

Release 2 of Oracle Data Mining 10g introduced support for a new mining application—anomaly detection, that is, the detection of “rare cases” when very few or even no examples of the rare case are available. Oracle Data Mining can “classify” data into “normal” and “abnormal” even if only one class is known. ODM uses a special case of the Support Vector Machines algorithm to create a model of known cases. When the model is applied to the general population, cases that don't fit the profile are flagged as anomalies (that is, abnormal or suspicious). ODM's anomaly detection algorithm is extremely powerful in finding truly rare occurrences when you have a lot of data but need to find needles in the haystacks.

Feature Extraction

ODM's Nonnegative Matrix Factorization (NMF) is useful for reducing a large dataset into representative attributes. Similar in high level concept to Principal Components Analysis (PCA), but able to handle much larger amounts of attributes and create new features in an additive nature, NMF is a powerful, cutting-edge data mining algorithm that can be used for a variety of use cases.

NMF can be used to reduce large amounts of data, e.g., text data, into smaller, more sparse representations that reduce the dimensionality of the data, i.e., the same information can be preserved using far fewer variables. The output of NMF models can be analyzed using supervised learning techniques such as SVMs or unsupervised learning techniques such as clustering techniques. Oracle Data Mining uses NMF and SVM algorithms to mine unstructured text data.

Text Mining and Unstructured Data

Oracle Data Mining provides a single unified analytic server platform capable of mining both structured, that is, data organized in rows and columns, and unstructured data. ODM can mine unstructured data, that is, “text” as a text attribute that can be combined with other structured data, for example, age, height, and weight to build classification, prediction, and clustering models. ODM could add, for example, a physician’s notes to the structured “clinical” data to extract more information and build better data mining models.

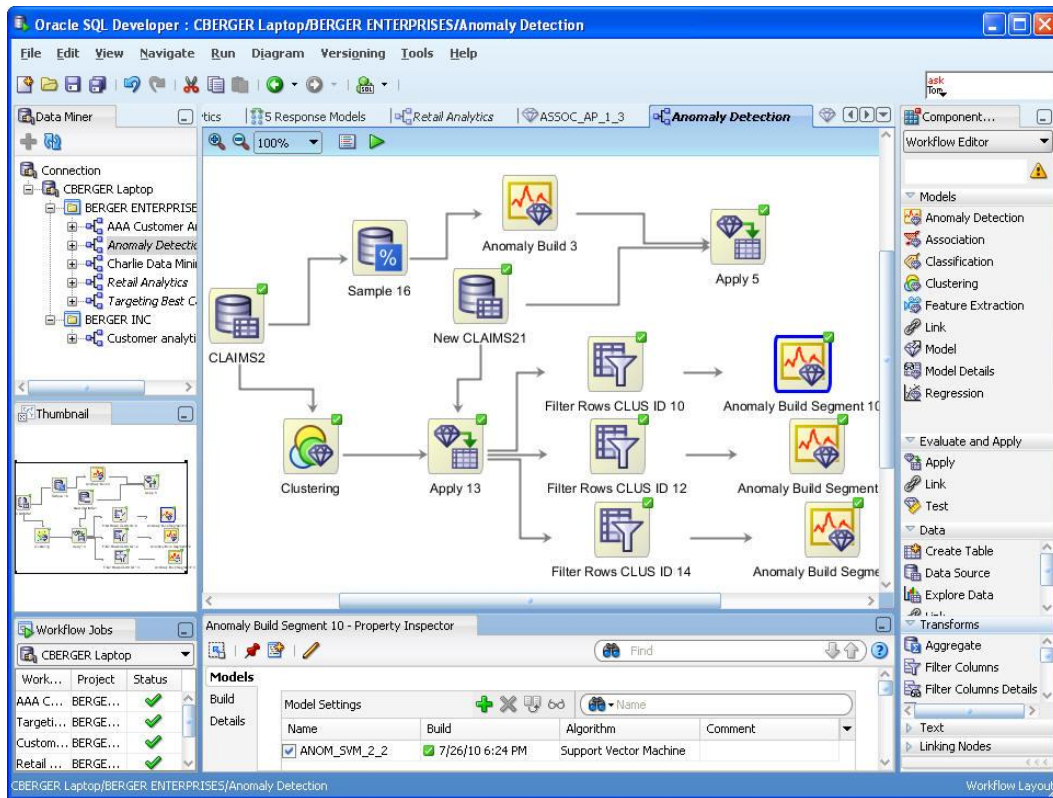
This ability to combine structured data with unstructured data opens new opportunities for mining data. For example, law enforcement personnel can build models that predict criminal behavior based on age, number of previous offenses, income, and so forth, and combine a police officer’s notes about the person to build more accurate models that take advantage of all available information.

Additionally, ODM’s ability to mine unstructured data is used within Oracle Text to classify and cluster text documents stored on the Database, e.g. Medline. Oracle Data Mining’s NMF and SVM models can be used with Oracle Text to build advanced document classification and clustering models.

Oracle Data Mining for Data Analysts

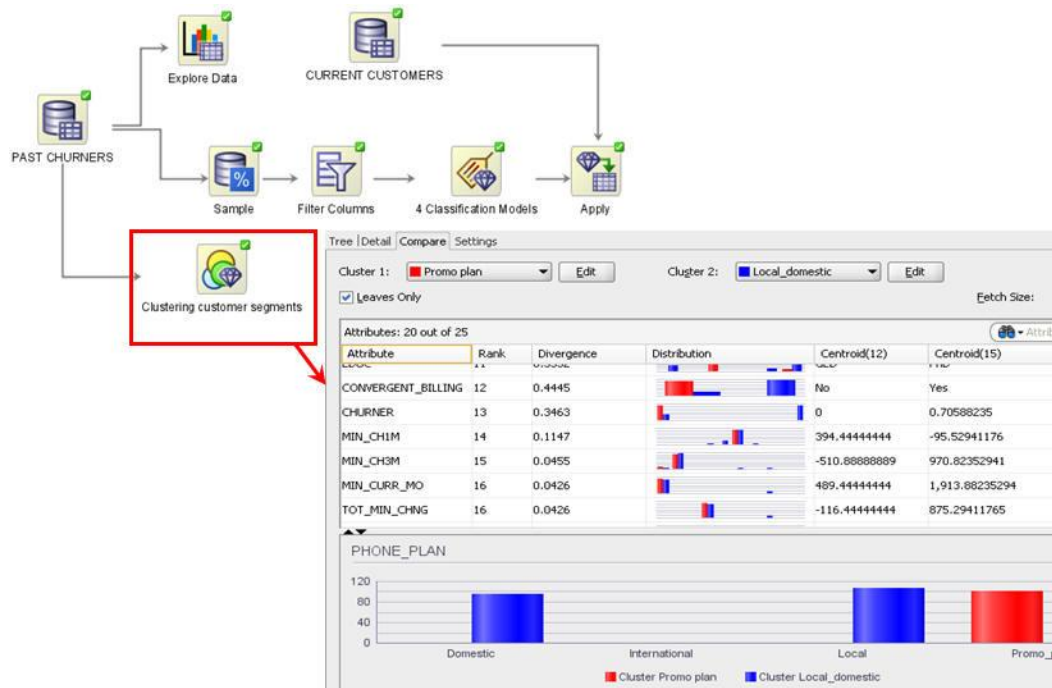
Oracle Data Mining empowers you to analyze your data like never before. With Oracle Data Mining you can build more models faster, and deploy those models more quickly and easily. With Oracle Data Mining, you can mine more sources and types of data, scale to high dimensional and large volumes of data, use state of the art algorithms, and build, evaluate and deploy models immediately throughout the enterprise. Oracle Data Mining provides disruptive technology that changes the way you extract information from data and reduces your dependence on traditional statistical environments. With Oracle Data Mining, you can now mine the data you want to mine.

Oracle Data Mining provides a graphical user interface, Oracle Data Miner, an extension to SQL Developer Release 3.1, with “nodes” for data exploration, data preparation, data mining, model evaluation, and model scoring process. Intelligent defaults are provided to help the data analyst to successfully mine their data. Power users may optionally change advanced settings or use the APIs.



Oracle Data Miner, available for download from the Oracle Technology Network (OTN) makes data mining easy

Oracle Data Miner can export work flows, generate PL/SQL scripts and SQL code snippets to accelerate transforming the data mining methodology into an integrated data mining/BI application.



Clustering analysis to discover customer segments based on behavior, demographics, plans, equipment, etc.

Oracle Data Mining for Applications Developers

Oracle Data Mining empowers IT departments, and application vendors by making data mining a natural extension of the Oracle Database. Many companies already rely on Oracle Database to store and manage their data. Now, companies can leverage in-database analytics and perform data mining right where the data exists: in the database. In-database mining, as provided by Oracle Data Mining, keeps not only the data in the database where it is safe and secure, but also keeps the data mining models and results there as well, where they are immediately usable in SQL queries and applications. The Oracle Database with Oracle Data Mining, Oracle R Enterprise, Oracle Business Intelligence and Oracle Applications provide a complete platform for data management, data analysis and end user applications.

With ODM's SQL and PL/SQL APIs, you can integrate data mining into existing business processes, workflows, and applications. Oracle Data Mining provides numerous sample code examples that can be used as starting points for quickly building data mining applications. Additionally, the Oracle Data Miner graphical user interface generates PL/SQL code that helps you to operationalize your mining activities.

```

drop table CLAIMS_SET;
exec dbms_data_mining.drop_model('CLAIMSMODEL');
create table CLAIMS_SET (setting_name varchar2(30), setting_value varchar2(4000));
insert into CLAIMS_SET values ('ALGO_NAME','ALGO_SUPPORT_VECTOR_MACHINES');
insert into CLAIMS_SET values ('PREP_AUTO','ON');
commit;

begin
dbms_data_mining.create_model('CLAIMSMODEL', 'CLASSIFICATION',
'CLAIMS', 'POLICYNUMBER', null, 'CLAIMS_SET');
end;
/
-- Top 5 most suspicious fraud policy holder claims
select * from
(select POLICYNUMBER, round(prob_fraud*100,2) percent_fraud,
rank() over (order by prob_fraud desc) rnk from
(select POLICYNUMBER, prediction_probability(CLAIMSMODEL, '0' using *) prob_fraud
from CLAIMS
where PASTNUMBEROFCLAIMS in ('2 to 4', 'more than 4'))
where rnk <= 5
order by percent_fraud desc;

```

<u>Mining Function</u>	<u>Algorithm</u>	<u>Sample Program</u>
Anomaly Detection	One-Class Support Vector Machine	dmsvodem.sql
Association Rules	Apriori	dmardemo.sql
Attribute Importance	Minimum Descriptor Length	dmaidemo.sql
Classification	Decision Tree	dmdtdemo.sql
Classification	Decision Tree (cross validation)	dmdtxvlddemo.sql
Classification	Logistic Regression	dmglcdem.sql
Classification	Naive Bayes	dmnbdemo.sql
Classification	Support Vector Machine	dmsvodem.sql
Clustering	k-Means	dmkndemo.sql
Clustering	O-Cluster	dmocdemo.sql
Feature Extraction	Non-Negative Matrix Factorization	dmrmdemo.sql
Regression	Linear Regression	dmglrdem.sql
Regression	Support Vector Machine	dmsvrtem.sql
Text Mining	Text transformation using Oracle Text	dmtxtfe.sql
Text Mining	Non-Negative Matrix Factorization	dmtxtnmf.sql
Text Mining	Support Vector Machine (Classification)	dmtxtsvm.sql

The PL/SQL Sample Programs provide examples of mini-solutions and use cases for Oracle Data Mining and are an excellent starting point when developing an ODM Application

Competing on In-Database Analytics

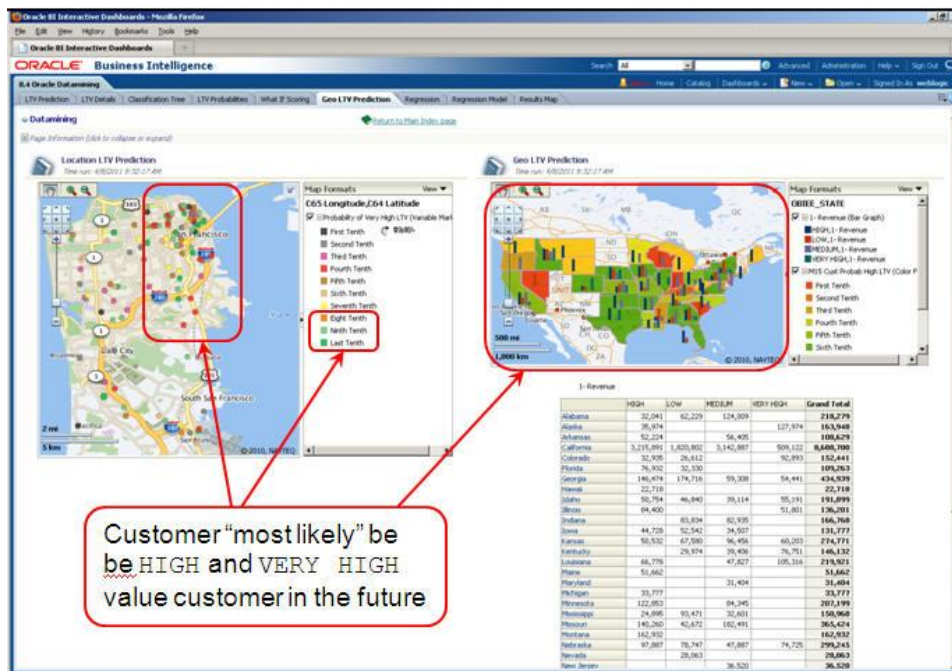
“Any company can generate simple descriptive statistics about aspects of its business—average revenue per employee, for example, or average order size. But analytics competitors look well beyond basic statistics. These companies use predictive modeling to identify the most profitable customers—plus those with the greatest profit potential and the ones most likely to cancel their accounts”

Harvard Business Review, “Competing on Analytics”, by Thomas Davenport.

Successful companies must go beyond basic BI dashboards and reporting. They must harvest maximal information from their data and leverage it for competitive advantage. Traditionally, this process of extracting information from data has been left to the purview of highly technical data analyst specialists using specialized data analysis tools and extracted copies of data.

Oracle Data Mining removes this barrier, allows data analysts direct access to the data and enables them to more rapidly build, evaluate and deploy predictive analytics throughout the enterprise in dashboards and next-generation applications “*powered* by Oracle Data Mining”. With Oracle Data Mining companies can:

- Eliminate data movement and collapse information latency
- Transform their database repository into an “analytical database”
- Deliver new insights and predictive analytics throughout the enterprise



Oracle Data Mining's Predictions & probabilities available in Database for Oracle BI EE and other reporting tools

Beyond a Tool; Enabling Predictive Applications

The true value of data mining is best realized when the new insights and predictions are integrated and operationalized into existing business applications. Oracle Data Mining provides an ideal data and business intelligence IT platform to enable companies to be successful competing on analytics. With Oracle Data Mining users can deploy predictive analytics into business applications, call centers, web sites, campaign management systems, automatic teller machines (ATMs), enterprise resource management (ERM), and other operational and business planning applications. With Oracle Data Mining, special departments of advanced data analysts working in silos far away from the database are no longer needed. The Database becomes more than just a data repository—it becomes an *analytical* database that can undergird many new advanced BI use cases.

ODM eliminates the extraction of data from the database for data mining, thus significantly reducing total cost of ownership. With ODM, there is no need for multiple data storage hardware and software environments, multiple data analysis tools, and multiple support resources. With ODM, there are fewer “moving parts” resulting in a simpler, more reliable, and more efficient data management and data analysis environment.

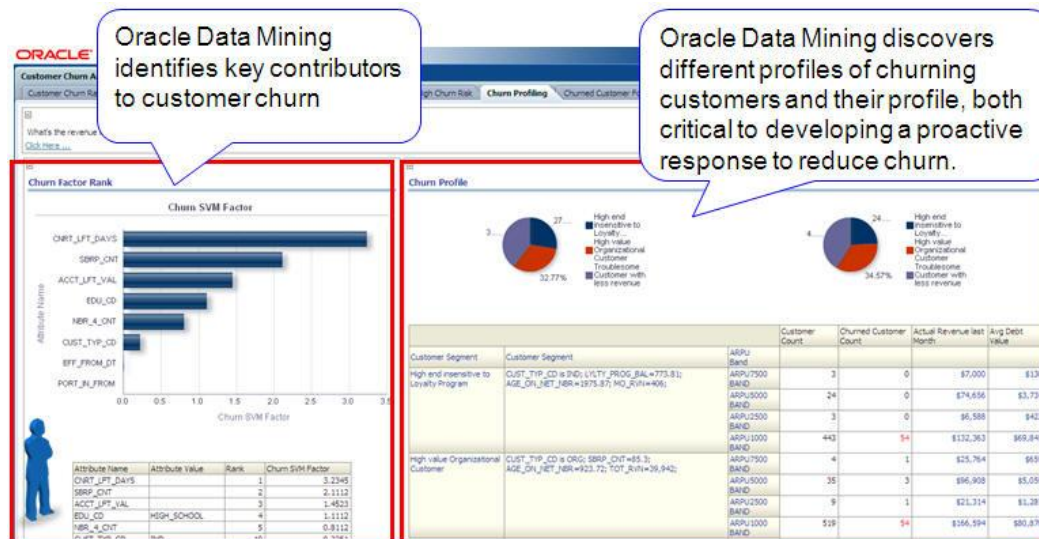
Automation of data mining tasks is facilitated by Oracle Data Mining's application programming interface (APIs). Oracle Data Mining supports SQL and PL/SQL APIs. Application programmers can control all aspects of data mining — they can expose complex settings for

advanced users or completely automate the process for business users. Programmatic control extends from data preparation through model build and to model apply for scoring of single records (on-demand) and batch scoring of large data sets. ODM model predictions may be called on demand as in the following example SQL query and/or stored in relational tables for access by other business applications:

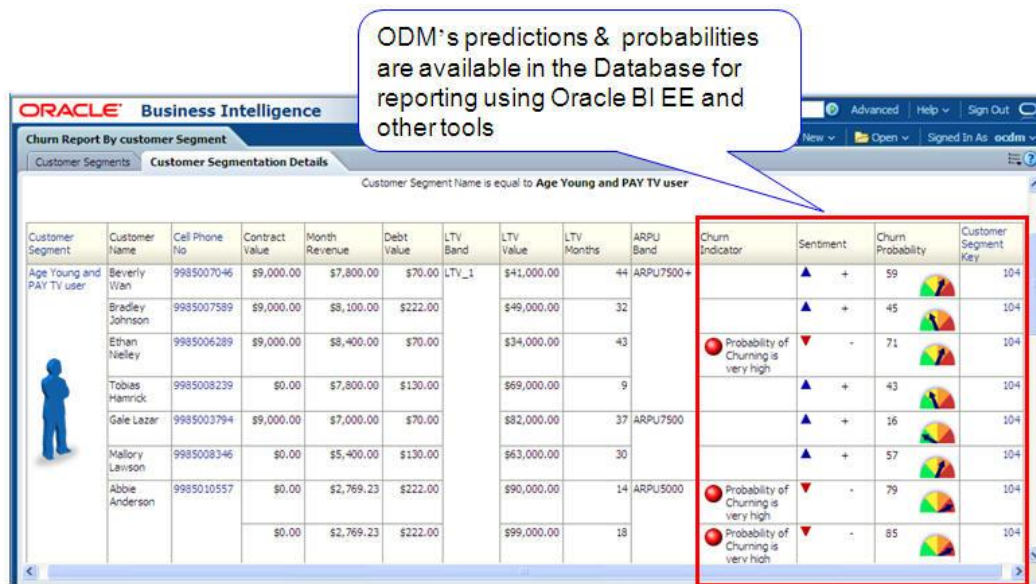
```
SELECT * from(
SELECT A.CUSTOMER_ID, A.AGE,
      MORTGAGE_AMOUNT, PREDICTION_PROBABILITY (LIKELY_RESPOND, 'YES'
      USING A.*) prob
FROM CUSTOMER_DATA.INSUR_CUST_LTV A)
WHERE prob > 0.85;
```

Oracle Data Mining's APIs provide direct, asynchronous access to ODM's functionality. Oracle Data Mining's PL/SQL and SQL based APIs enable application developers to enhance, for example, a call center application to highlight a customer's likelihood to churn or to become a profitable customer. The probability that the customer will accept the special offers can be displayed for the customer service representative as a window pop-up to provide better service the customer.

Because all results are created and stored in an open relational database, users have access to data mining results using a wide variety of business intelligence tools including Oracle Business Intelligence EE, Oracle OLAP, Oracle Reports, Oracle Portal, and Oracle Applications.



Oracle Data Mining powering the Oracle Communications Industry Data Model with embedded and automated predictive analytics.

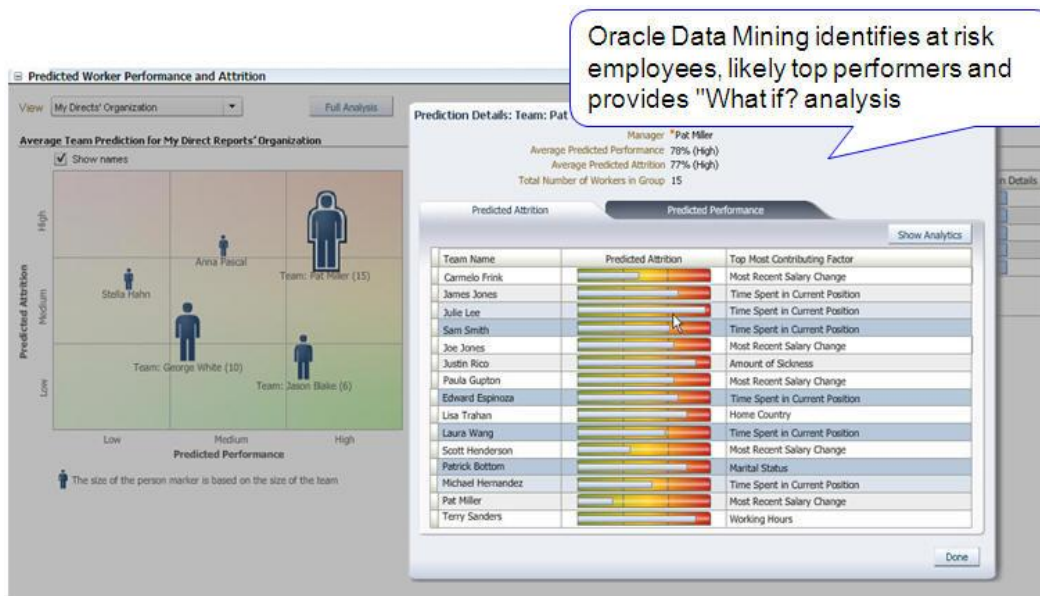


Oracle Data Mining powering the Oracle Communications Industry Data Model with embedded and automated predictive analytics.

Oracle now delivers applications and industry data models that have integrated data mining. Oracle's Human Capital Management (HCM) Fusion Application, for example, provides embedded Oracle Data Mining models for predicting employee behavior based on observed past employee behavior at that company. Oracle Data Mining's models enable the Oracle HCM Predictive Workforce product to:

- Predict and anticipate "at risk" employees who are most likely to voluntarily leave
- Identify the factors that most contribute to an employee's predicted behavior
- Perform real-time "What if? Analysis" to change those factors and see anticipated results
- Identify and predict likely top and bottom performers

Oracle's Industry Data Models delivers a standards-based data model, designed and pre-tuned for Oracle data warehouses. Oracle Retail Data Model combines market-leading retail application knowledge with the power of Oracle's Data Warehouse and Business Intelligence platforms. With pre-built Oracle Data Mining, Oracle OLAP and dimensional models, it delivers industry-specific metrics and insights you can act on immediately. With Oracle Retail Data Model, you can jump-start the design and implementation of a retail data warehouse to quickly achieve a positive ROI for your data warehousing and business intelligence project with a predictable implementation effort.



Built-in self-learning Oracle Data Mining predictive models for employee behavior.

Oracle Data Mining is part of Oracle's family of business intelligence products and features. With Oracle, the data can come from the same "single source of truth" accessed by enterprise users and protected by database security schemes. Oracle Data Mining makes it easy to build enterprise applications that automate data mining and distribute new insights within the organization.

Spend Less

Oracle Data Mining, a component of the Oracle Advanced Analytics Option, provides a cost effective alternative to expensive traditional statistical analysis software. Savings are realized in avoiding additional hardware purchases for computing and storage environments, redundant copies of the data and multiple versions of the data, duplication of personnel who perform similar functions but unnecessarily use different software packages. Whereas traditional statistical software is rented under an annual usage fee (AUF) pricing scheme, companies can reduce their overall data analysis costs by making the specialized statistical software available for only those individuals who really need it. Production analytical applications can be implemented now inside the Oracle Database for in-database analytics.

Eliminate Redundant Data and Traditional Analytical Servers

Oracle Data Mining significantly reduces the cost of data mining. It reduces the need for separate, dedicated mining servers since the data does not need to be extracted outside of the Oracle Database.

Additionally, by utilizing the same data and a “single source of truth”, Oracle Data Mining makes it easy to properly select the data and to always work with the most up to date version of the data.



Oracle Data Mining eliminates data movement, data duplication and security exposures.

Conclusion

Oracle Data Mining, a component of the Oracle Advanced Analytics Option, provides a powerful, scalable in-database data mining engine for data analysts seeking to harvest valuable new information and an industry-standard infrastructure for application developers looking to build applications that automate the discovery and deployment of predictive analytics throughout the enterprise.

Oracle Data Mining’s wide range of data mining algorithms, completely embedded in the Oracle Database, solve a wide variety of business problems and provide a powerful infrastructure for building, automating and deploying advanced enterprise business intelligence applications.

By automating, integrating, and operationalizing the discovery and distribution of predictive analytics and new business insights, companies can leverage their Oracle Database technology investment, to operate more intelligently, and most importantly, to gain competitive advantage.



Oracle White Paper— Oracle Data Mining 11g
Release 2: Competing on In-Database Analytics
February 2012
Author: Charlie Berger

Oracle Corporation
World Headquarters
500 Oracle Parkway
Redwood Shores, CA 94065
U.S.A.

Worldwide Inquiries:
Phone: +1.650.506.7000
Fax: +1.650.506.7200
oracle.com



Oracle is committed to developing practices and products that help protect the environment

Copyright © 2012, Oracle and/or its affiliates. All rights reserved. This document is provided for information purposes only and the contents hereof are subject to change without notice. This document is not warranted to be error-free, nor subject to any other warranties or conditions, whether expressed orally or implied in law, including implied warranties and conditions of merchantability or fitness for a particular purpose. We specifically disclaim any liability with respect to this document and no contractual obligations are formed either directly or indirectly by this document. This document may not be reproduced or transmitted in any form or by any means, electronic or mechanical, for any purpose, without our prior written permission.

Oracle is a registered trademark of Oracle Corporation and/or its affiliates. Other names may be trademarks of their respective owners.