

How to Choose a Data Mining Suite

The choice of a data mining suite is not an easy task. This article provides a brief outline of some considerations that could affect your decision. Contrary to common opinion, the best tool suite for you may not be the most advanced tool; it may not be the one with the most data mining algorithms, nor the one that gives the greatest accuracy in prediction. More important than all of these things is identifying the tool suite that is:

- Easy to use.
- Provides acceptable accuracy (even though not the highest accuracy available)
- Able to perform all the common tasks in a data mining project.

Ease of Use

Some traditional (and heavily advertised) data mining tools may provide a rich variety of data processing and modeling capabilities, but require a legion of “priests” to use them. Often, these “priests” of data mining are developed only after many years of practice and travel up the “learning curve” of the tool’s capabilities. Rather than be very procedural (programmed with a scripting language), the user interface to data mining technology should be like the interface to automobile technology. The great success of the automobile is because it brings the benefits of sophisticated engineering technology down to the level of use appropriate to the common man and woman. You don’t need to be an expert in internal combustion technology, or understand the complex relationships between gear ratios and acceleration to use a car effectively. All you have to do is get behind the wheel, turn on the ignition, step on the gas, and turn the wheel or use the brakes at the appropriate times, and voila’, you are an expert user of the automobile! Using a data mining tool should be like that. You might be surprised to learn that several modern data mining tool suites approach that ease of use.

Accuracy: How High?

Suppose you could buy a tool for a fraction of the cost of the priestly tool (maybe 20% or less), which permitted ordinary business analysts and statisticians to create models that were 80% as good as the priest could create. Would you choose to buy the priestly tool, or the 80% tool? I think for most companies, the answer is, “of course, we want the 80% tool”. For this case as well, several data mining tool suites available today can provide this functionality. For other purposes, the best tool might be the most accurate tool.

Ability to Perform All Common Data Mining Tasks

Most data miners will tell you that 70-90% of the time required to perform a data mining project is spent in data preparation for modeling. Reasons for this include:

- Most data mining algorithms require clean and complete data records as input. No data mining tool in the world can analyze data that does not exist (missing data in some fields).

- Most data in commercial databases was collected from transactional systems to serve query and reporting purposes, not analytical purposes.
- Most data in commercial databases are rather “dirty”. That is, databases often contain inappropriate data, training data, improperly input data, or just plain garbage data. Even if the data appears to be clean, historical data records may reflect changes in coding and aggregation rules at various times in the past, which must be reconciled. In addition, data formats may not be consistent across databases used as data sources. Finally, data may require transformation to different ranges or different expressions (letters changed to numbers), or new variables may be needed that are combinations of existing variables. A good data mining suite will provide tools for performing all of these operations. Some data mining suites are better at it than others.

This article reviews five of the most useful and powerful data mining suites available today, STATISTICA-Data Miner, SPSS-Clementine, Affinium Model, Insightful Miner, and KXEN. We can use these tools to illustrate how you can evaluate how suitable a data mining tool suite is for your use. Let’s cut to the chase right in the beginning.

There is no best tool overall.

Are you surprised? Well, competition in the marketplace almost guarantees this. If a particular tool is successful enough to make it into the “mainstream” of data mining use, it must serve well at least a moderate segment of business needs. Each tool suite has its strengths and weaknesses; each tool suite may be the best for particular needs in particular companies. Each of the five tool suites will be reviewed and classified according to their best uses. From this evaluation, you can gain enough information to take the first step in the choice of the data mining tool suite that is right for you.

The first step is to look at the features and functions of the data mining tool suite. While this only first step in the decision-making process, it may not be the most important consideration for you. Table 1 shows a weighted comparison of the features and functions of the tool suites. You will notice by comparing the relatively moderate cost and the weighted score across all features and functions, that STATISTICA Data Miner is the clear winner. This *does not* mean that this tool is best for you. Your needs may not require (or your budget may not permit) the rich variety of capabilities provided by STATISTICA Data Miner; Insightful Miner (with its great ease of use and affordability) may be just the right tool for you, regardless of its relatively low score in Table 1. Or, you may want a fully automatic data mining engine that can generate models of the very highest accuracy, to which you are willing to submit data in the suitable format. If so, KXEN is the right tool suite for you, providing the cost is acceptable. Clementine and Affinium Model tool suites provide intermediate solutions between those of KXEN and Insightful Miner, in terms of functionality and cost.

	Wt	CLEM \$50K	STAT \$15K	IM \$15K	AM \$50K	KXEN \$20K
Price, Single-User						
1. Data Prep tools	8					
Data Exploration tools	2	5	10	6	5	4
Data Cleaning tools	2	5	6	5	3	0
Data Transf.tools	2	7	9	4	0	3
Data splitting tools	2	5	8	5	3	0
2. Algorithm Richness	12					
Clustering Methods	2	2	10	8	8	0
Neural Nets	2	3	9	4	7	8
Regression methods	2	6	10	7	6	
Classification methods	2	10	8	6	7	8
Time-Series	2	6	10	0	0	
Other	2		10	0	0	10
3. Tunability of Algorithms	6					
Parameter Search?	2	7	0	0	8	8
Facilitation of Boosting?	1	0	8	0	2	0
Facilitation of Bagging?	1	0	7	0	1	0
4. Degree of Automation	6					
Data Preparation	2	2	2	2	2	10
Model building	2	3	3	2	8	10
Recoding/standardization	2	3	4	2	0	10
5. Scalability	4					
Parallel Support?	2	2	8	0	0	0
IDB capability for scoring?	2	0	10	0	0	0
6. Model Portability	6					
C/C++	2	8	10	0	0	10
SQL	2	10	10	0	0	
XML	1	0	10	0	0	10
PMML	1	0	10	0	0	
7. Web Enablement	7	0	10	0	0	0
8. Stat. Functionality	5					
Level of Integration	2	1	9	5	0	0
Procedural or Menu	3	3	10	10	10	0
9. Graphical/Reporting	8					
Correlation Matrix	2	0	9	9	6	0
Scatterplots/matrix	1	0	9	3	7	0
Line plots	2	8	9	0	7	0

Reports	1	8	5	0	5	0
Total Graphics types	2	3	10	3	2	0
10. GUI Intuitiveness	10					
Integration with DM Proc.	2	8	8	2	8	7
East of building models	2	4	9	4	9	10
Help or Wizards	2	5	5	4	5	5
Complexity of streams	2	7	9	7	0	0
Visual progr. capability	2	10	10	10	0	0
11. Expandability	3					
Custom nodes	1	1	6	0	0	0
C-Program integration	1	0	8	0	0	0
Add Algorithms?	1	0	8	0	0	0
12. Scripting/Macros	4					
Standard Script Lang.	1	6	10	0	7	0
Help/Syntax examples	1	8	7	0	7	0
Operators/Functions	1	9	9	0	9	0
Modeling & Scoring	1	0	8	0	8	0
13 DM Process	6					
Explicit in design	3	9	9	3	7	6
CRISP/SEMMA, or ?	3	10	8	4	4	4
14. Data Preparation	7					
Recoding	1	3	10	2	0	10
New variable derivation	2	8	10	5	8	10
Standardization	1	0	10	0	0	10
Dataset splitting	1	8	10	7	0	10
Aggregation capability	1	10	10	9	0	0
Temporal Abstractions		0	10			10
15. Metadata re-mapping		0	0	0	10	0
16. Model Eval./Reporting	8					
Lift charts	2	0	8	7	8	10
KS statistic	1	0	0	0	0	0
RMS, MAE, etc.	1	0	7	0	0	7
Cross-validation	1	0	8	0	0	0
Variable sens. reports	1	8	8	6	8	10
Pred. accuracy reports	1	8	8	8	8	8
Confusion matrix	1	0	8	10	7	0
17. Data Import/Output	8					
ASCII						
Fixed	1	10	9	9	10	10

Free form	2	10	9	9	10	10
Choice of separators	1	10	9	9	10	10
Excel	1	5	9	9	10	0
Database struct.(ODBC)	1	7	9	9	10	8
Metadata capture	1	10	9	9	10	10
Date formats		10	7	5		10
18. Data Exploration	7					
Missing value reports	1	8	9	9	10	8
Missing value imputation	2	6	9	0	9	7
Descriptive stat reports	2	4	10	5	0	1
Binning	1	0	7	8	0	0
Sampling	1	7	8	7	0	0
19. DM Project Mgt.	3					
Project-level organization	1	1	5	1	8	1
Multiple modeling	1	8	9	8	0	8
Meta-modeling	1	6	9	5	0	8
20. Dimensionality Reduct.	8					
CHAID	2	0	9	0	5	0
PCA	2	9	9	9	5	7
Correlation matrix	2	0	9	8	0	0
Other	2	0	10	0	0	0
	4					
21. Drilldown	1	0	10	2	2	0
SQL/internal DB	1	0	8	0	0	0
SQL/external DB	1	2	8	7	0	8
Canned reports.	1	0	3	0	5	0
22. Client-Server capability	6	8	10	0	5	10
23. Product maturity	5	10	7	2	10	8
WEIGHTED TOTALS		641	1147	474	565	562

Table 1. Weighted scores for the ability of five data mining tool suites to perform common data mining tasks.

Performance

The performance of each tool suite was tested with a neural net (when possible). Prediction accuracies are shown in Table 2. For data sets with a binary target variable, percent accuracies are listed for the 1's and 0's in the form of 85/56. For the Abalone dataset with an integer target variable, overall accuracy and mean absolute error (MAE) are listed.

PERFORMANCE	CLEM	STAT	IM	AM	KXEN
Failures dataset					
Prediction Accuracy	69.1/ 83.4%	81.2/80.8%	74.8/88.7%	81.8/96.8%	90.6/82.5
Model used	Neural Net	Neural Net	Neural Net	Neural Net	SVM
Thrombin dataset					
Prediction Accuracy	30.5/100%	1:32% 0:100%	31.9/100%	31.9/100%	75/0%
Model used	Neural Net	Neural Net	Neural Net	Neural Net	SVM
Abalone dataset					
Prediction Accuracy	94.4% MAE = 0.14	64.6% MAE = 1.45	59% MAE 1.68 (Error)		88% MAE 1.95
Model used	Neural Net	Neural Net	Neural Net		SVM

Table 2. Typical performance of the five tool suites.

Tool Suite Comparisons

SPSS-Clementine

Clementine has been around for a long time. The tool suite is very mature and has a very faithful following, particularly in Europe (it was developed in England). Clementine was the first data mining suite to use the graphical programming approach used previously by the scientific programming tools Stella, I-Think, MathCad and MatLab in the 1980's.

Pros:

- Good variety of data mining algorithms
- Very powerful optimal parameter search routines built into many of the data mining algorithms (automatic trials of different parameter sets)
- Very powerful combination of the Type node and Quality node for data quality checks and missing value imputation
- Power meta-learning models can be built, in which the results of one modeling algorithm can be easily streamed as input to another modeling algorithm
- Powerful (but proprietary) internal scripting language (CLEM) for creating complex variable processing
- Moderately easy to use

Cons:

- Relatively little descriptive statistical or parametric statistical analysis capabilities are available directly in the tool (although SPSS nodes can be used for input from and output to the SPSS StatPak)
- Relatively poor descriptive or output graphics forms
- Model export for scoring outside the tool suite (or to perform deployment calculations at a faster speed than the interpreter based internal tool) must be done via an optional Publisher product (\$25,000).

STATISTICA Data Miner

STATISTICA Data Miner (like KXEN) is a tool in a class by itself. This uniqueness is defined primarily by in terms of the many things it does well, and the completeness in facilitating all tasks of the data mining project. Other tools may be easier to use (e.g. Insightful-Miner) or employ more automation (Affinium Model or KXEN), but no data mining suite available today provides more tools for performing data mining projects. This tool suite is my personal favorite.

Pros:

- Provides the richest combination of parametric statistical and machine learning data mining algorithms
- Relatively easy to use graphical programming user interface
- Provides tools for all common data mining tasks
- Highly flexible tools for model output
- Powerful tools for reduction of dimensionality
- Scalability (STATISTICA Data Miner can rapidly process larger data sets both in terms of their dimensionality and the overall size than the other products)
- Powerful customization options based on the industry standard VB language

Cons:

- Lift charts are not easily available for evaluation of neural net models
- Training in statistical analysis is best for properly interpreting the results of the parametric statistical algorithms

KXEN

KXEN is one of two tool suites that provide an implementation of a Support Vector Machine (SVM). STATISTICA Data Miner is the other. For this algorithm, the input data must be transformed to the range of -1 to $+1$ (done automatically by KXEN) and projected into the feature space. An SVM finds the solution of minimum error by identifying the “maximal separable hyperplane” through a cloud of data points in feature-space (a theoretical space with a number of dimensions equal to the number of predictor variables). The error minimization approach of a neural net is an iterative search for the set of predictor variable weights that produces the least amount of error between actual and predicted values. The problem with this kind of a search is that it is quite possible for the search to end at a solution that represents only a “local minimum” for total error across only a region of data points, rather than a “global minimum” total error across all data points. An SVM *induces* the solution of minimum error across all the data points in the feature-space, rather than *deduces* it by an iterative procedure.

Pros:

- KXEN is clearly the most accurate data mining tool available today
- Various combinations and transforms of existing variables are automatically created and included in the analysis as derived predictor variables.
- This tool is almost fully automatic! For performing analyses on appropriately cleansed and formatted datasets, this tool comes closest to the ideal of the automobile user interface. KXEN is the *only* data mining tool available today that can be so easily embedded into your data processing stream to use as a data mining engine. In fact, KXEN means “Knowledge Extraction Engine”; and it lives up to its name.

Cons:

- A clean data set must be submitted to the Consistent Coder of KXEN in the form of one record per entity to be modeled.

- There are no data preparation tools to help you put the data in this form (although many data preparation steps required by parametric statistical or neural nets/decision trees are not necessary with SVMs).
- No coincidence (or “confusion”) matrix is available for binary output, from which precision and recall values can be calculated. You can create one, if you can determine the correct threshold to use on the decimal output to convert it to binary predicted values

Insightful Miner

This tool suite may be the best one available for a company that would like to use ordinary business analysts to do relatively simple data mining projects. Insightful Miner is a good data mining tool for use with S-Plus systems, because the entire library of S-Plus functions are available for use with it! In addition, it provides a rich assortment of data mining and statistical data mining algorithms (but not nearly as rich as does STATISTICA Data Miner). For almost all common steps in the data mining project, Insightful Miner clearly gives the best bang for the buck.

Pros:

- Excellent tools for data import/export, data exploration and data cleansing tasks, and reduction of dimensionality prior to modeling
- Even though it does not employ a graphical programming interface, it is relatively easy to use by non data miners
- The most complete general purpose data mining suite available, and it is relatively inexpensive

Cons:

- A relatively low level of automation
- No scripting interface for coding of complex problems
- Recoding must be done via an expression language in the Create Columns node
- No model exporting capabilities

Affinium Model

This tool is the easiest to use response modeling product on the market, even easier than STATISTICA- Data Miner. This is the best package for use by the non-data miner/statistician, for whom the lack of a rich statistical and graphical backbone is not a problem. The automatic operation of the modeling engine shields the user from many data mining operations that must be manually performed by users of other packages, including choice of algorithms. The user has only to choose the level of analysis from quick to extensive, and the tool automatically creates models from a small to a large number of algorithms and parameter settings, while saving the current best model. Four different modeling applications: Response modeler, Cross-Seller, Customer Segmenter, and Customer Valuator are actually very similar in function, and differ only the terms used in creating the model. This seems like just a repackaging of the same thing, but that is part of the appeal of this tool. In the mind of the non-data miner, perception is most of

the problem in using data mining tools for different purposes. The Modeling button brings up a list of modeling options (Quick-Intermediate-Extensive) that cycles through an increasingly large number of modeling algorithms and associated parameters to find the optimum model for each data set. Optionally, you can select one particular model type, or no model. Following that choice, the modeling process in Affinium Model is automatic.

Pros:

- The menu items are arranged from left to right showing modeling application, data import, modeling, model reports, scoring, scoring reports, and variable name editing.
- The data is imported into an internal spreadsheet, like in STATISTICA-Data Miner, but the only manipulation of the data is permitted through the Edit button (only for variable name changes) and via a drop-down menu with a data quality option (for missing data reports and imputation).
- New variables can be derived in the spreadsheet with a rich set of macro functions.
- Interpretation of the model results is very intuitive. The user has the choice of viewing a brief report, a detailed report, a lift curve, or a variable sensitivity report.

Cons:

- No data exploration tools
- The biggest potential drawback with this product is the almost complete lack of data preparation functions. Input data must be properly prepared in other tools before import into Affinium Model.
- There is no evidence that data is standardized (conversion of the ranges of all variables to a common scale) before submission to the modeling algorithm.

How can you tell which tool would be best for you? One way is to match the tool to the data scenarios for which you plan to use it. The following are several common data scenarios that you might encounter in your business.

Data Scenarios

The choice of the proper data mining suite for your use may depend on the data environment in which you would like to use it. Here are some data scenarios and choices of appropriate tools to use.

Scenario #1. If the company has access to (or is willing to hire) people with statistical expertise, the best tool will be one that statisticians understand and can use effectively:

- STATISTICA Data Miner
- SPSS Clementine (in conjunction with SPSS Stat package)

Scenario #2. If data preparation must be done by hand inside the data mining package, then the best tools would include:

- STATISTICA Data Miner
- Insightful Miner
- SPSS-Clementine

Affinium Model would not be a good choice here, because relatively few data preparation operations are supported in the tool.

Scenario #3. If the company wants to do data mining modeling with lower-level business analysts, then the best tool will have a relatively high degree of automation:

- Affinium Model
- KXEN
- Insightful Miner

Scenario #4. If the company has its own in-house analytical tools that require some enhancement to provide data mining capability, then the best data mining tool will be one that is easily embedded into their existing systems:

- KXEN
For example, NCR's Warehouse Miner could be coupled with KXEN to provide access to very fast data mining solutions within the Teradata database system. KXEN can be used to find the optimum support vector for a solution, in which some of the important variables in the solution may include non-intuitive constructs created by the KXEN Consistent-Coder prior to modeling.
- Affinium Model

About the Author

Dr. Nisbet has over 30 years experience in complex systems analysis and modeling as a Research Professor (University of California, Santa Barbara) and as consultant in data mining sciences. While at NCR Corporation and Torrent Systems, he pioneered the design and development of configurable data mining applications for retail sales forecasting, and Churn, Propensity-to-buy, and Customer Acquisition in Telecommunications and Insurance.