# STATISTICA 9
## NEW FEATURES

**Even Faster!**

## Performance Improvements

*STATISTICA 9* offers many new, unique features and added functionality to the entire *STATISTICA* product line. Many low-level optimizations have been implemented that greatly improve the overall performance of *STATISTICA*. These improvements can be found in both the 32-bit version and the new, native 64-bit version.

Because of the new technology and a variety of optimizations introduced to the main computational kernel, not only the 64-bit version but also the 32-bit version of *STATISTICA 9* is significantly faster in most operations when compared to *STATISTICA 8* (which already was one of the fastest analytic applications on the market).

### Multi-Core CPU Support

A highly optimized support for "multithreading" that takes advantage of multi-core processors has been implemented in most computationally intensive analytic and data management procedures of *STATISTICA 9,* resulting in dramatic speed gains. This technology will benefit not only power users who have access to server-class computers equipped with multiple processors, but it produces dramatic speed improvements even on inexpensive, commonly used dual-core CPU machines (which account for the majority of computers produced since 2008). These improvements can be found in both the 32-bit and 64-bit versions of *STATISTICA*.

### Native 64-Bit

*STATISTICA* 64-bit runs in native 64-bit mode and takes full advantage of the 64-bit operating systems, allowing for even better memory management and performance. Consequently, *STATISTICA 9* 64-bit can process designs of an extremely large size and further improves performance, more than doubling the speed of some computationally intensive analyses when compared to *STATISTICA 9* 32-bit. *STATISTICA* 64-bit is especially useful in data mining and other operations that use extremely large data sets and iterative, computationally demanding applications.

### Support for Larger Data Sets

Further optimizations have been implemented to more efficiently support data sets of extreme size in various modules that require iterative processing of the entire data set.

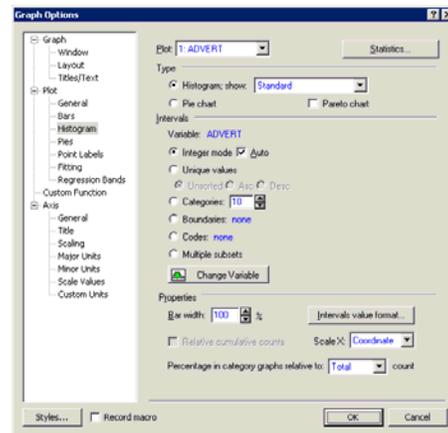## User Interface Enhancements

### Enhanced User Experience

Application navigation is simpler and more intuitive with the addition of an Office 2007-style ribbon bar as an alternative to the classic, pull-down menu based user interface. Frequently used functionality is quickly visible, and related functionality is easily found.



Note that the classic menus/toolbars will continue to be available for compatibility and to support standard customization options, and you can switch between the two interfaces at any time.

To display the classic menus/toolbars, click *Menus* on the Quick Access toolbar in the upper-left corner of the ribbon bar. To display the *STATISTICA* ribbon bar, select *Ribbon Bar* from the *View* menu.

Other improvements in the user interface can be seen with options that control the visual appearance of graphs, the behavior of spreadsheets, and the overall responsiveness of *STATISTICA*. In the more complex (multi-layered) dialogs, it is now easier to find and explore your choices with new tree controls for application options (*Options* dialog) and graph options (*Graph Options* dialog).



A quick way to specify subsets of cases for analyses is to select the *Enable Selection Conditions* check box in the *Spreadsheet Case Selection Conditions* dialog. This new option enables the visual display of selected cases, identifying them with a light green background by default. This enables you to see easily which cases will be used for the analysis.
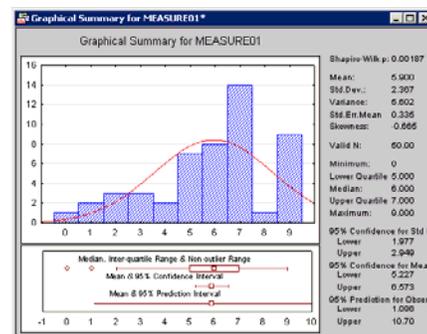
Spreadsheet cell-display enhancements apply to text that is too long to



be displayed in a cell at the current column width. The improvements in version 9 are two-fold: with the default wrapping settings, if the adjacent cells are empty, the text will now extend into those adjacent cells. Secondly, you can now *Wrap text* within the spreadsheet data, allowing lengthy text to be displayed on multiple lines in a spreadsheet cell.
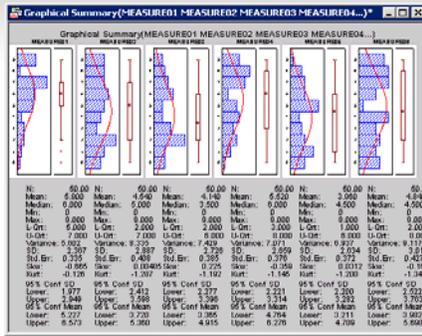
## Graphics Enhancements

### Visual Exploration of Data

One of *STATISTICA*'s strengths is visual exploration of data. StatSoft continues to build upon this strength with new visualizations for exploring data distri-bution, comparing variables, and creating color maps.



Two new compound graph options are available from *Descriptive statistics* (*Basic Statistics* module): *Graphs 2* and *Graphical comparative summary display*. Use *Graphs 2* to explore data distribution for one variable.

Use *Graphical comparative summary display* to compare up to six variables in one graph. The histogram and box plot included use the same scale per variable.

Color maps have been added to *Correlation matrices*. For example, in the display shown below, it is easy to see that winning in football is



highly correlated to playing on the home field, score, rushing, etc. (in this example, color coding is based on the statistical significance level of the correlation coefficients, as shown in the legend displayed in the title area of the table).



A *Variability plot* is used to evaluate the variability of one factor within several other organizing factors and for visual breakdown and data mining analyses. Now the *Variability plot* can be used to create user-defined graphs.

As with other user-defined graphs, this can save time, i.e., you can save common settings and then apply them to another data set.

## Marker Defaults

For improved data visualization, the default background for outline style markers is now transparent so that you can better see when multiple points are clustered together on a scatterplot.

## Graph Categories

In previous versions, the number of graph categories was limited to 255; this limit has been increased to 1,000; e.g., it is now possible to create a multiple box plot with up to 1,000 boxes.

## Graph Customization
## Macro Recording Options

Macro recording is a powerful option, and the macro object model architecture of *STATISTICA* with scripting capabilities is easily available to anyone. No technical knowledge of *STATISTICA* Visual Basic (SVB) is required. Common tasks can be automated or controlled for regulated work environments.

Macro recording enables you to access programmatically almost every aspect and virtually every detail of the functionality of the program. Even the most complex analyses and graphs can be recorded into SVB macro programs. They can later be run repeatedly or edited and used as building blocks of other applications.

Now this functionality has been expanded to include the recording of custom graph options. For example, suppose you create a line plot and want to change the line color and thickness. You open the *Graph Options* dialog and then modify the line options in the *Plot General* options pane. Select the *Record Macro* check box at the bottom of the *Graph Options* dialog, and click the *OK* button to generate a macro and produce your changes.

This powerful option also applies when recording master macros. You can now start a master macro recording session, run analyses, and customize the resulting graphs. When you re-run the master macro, the analyses will be replicated, along with the graphics customizations.

# Statistics Enhancements

## Distributions and Simulation

We are pleased to announce the beta release of *Distributions & Simulation*. This module enables users to automatically fit a large number of distributions for continuous and categorical variables to lists of variables. Standard distributions are available (normal, half-normal, log-normal, Weibull, etc.), but also included are specialized and general distributions (Johnson, Gaussian Mixture, Generalized Pareto, Generalized Extreme Value), and *STATISTICA* automatically ranks the quality of the fit for each selected distribution and variable.

In addition, the distributions fit to the list of selected variables and the covariance between the selected variables can be saved for deployment. The *Distributions & Simulation* module uses this deployment information to generate simulated data sets that not only faithfully reproduce the respective distributions, but also the covariances between variables. In short, in addition to facilitating efficient distribution fitting to large numbers of variables, this module enables users to fit general multivariate distributions, and simulate from those distributions, using cutting edge simulation techniques (e.g., Latin-Hypercube simulation).

These methods have proven useful in various domains such as modern DOE, reliability engineering, and risk modeling.

StatSoft welcomes your comments or observations regarding this new addition to the selection of analytic modules in *STATISTICA*. Please send your input to beta@statsoft.com.

## Basic Statistics

Several new basic statistics have been added. The computation of Welch's F statistic to test for equality of means when the variances are unequal is now available on the *ANOVA & tests* tab, located in the *Breakdown & one-way ANOVA* analysis results dialog.

Confidence interval estimates for the difference between means are now available on the *Options* tab of the *T-Test for Independent Samples by Variables* dialog and the *Advanced* tab of the *T-Test for Dependent Samples* dialog.

## General Optimization

The *General Optimization* module, which is part of *STATISTICA Process Optimization*, is a unique, powerful, open-architecture product that enables users to optimize arbitrary functions of virtually any complexity using Simplex, Genetic Algorithm, or Grid-Search methods. This module (released in beta version) has applications in virtually all domains in which there is a need to find the best parameters that control specific processes to achieve optimal results according to user-specified criteria (e.g., process industries, business, finance, science). The function to be optimized can be specified in a *STATISTICA* Visual Basic (SVB) function or a set of formulas. This new module was specifically designed to make repeated invocation of other *STATISTICA* (or other, e.g., R) functions referenced in the optimization function very efficient. Therefore, optimization problems that involve multiple data mining prediction models (e.g., complex cost models) or simulation (for stochastic optimization, or for optimizing for multivariate process capability) can now be easily set up and solved efficiently.

## Generalized Linear/Nonlinear Models

There have been multiple improvements to the *Generalized Linear/Nonlinear Models (GLZ)* module.

The Tweedie distribution is now an available distribution for dependent variables. This distribution is useful for modeling insurance claim amounts.

Odds ratios and their confidence intervals for parameters are now computed for appropriate analyses. Using over-parameterized coding and assuming a binomial distribution, a second spreadsheet containing odds ratios and their confidence intervals is produced when the parameter estimates results spreadsheet is requested.

Users can now generate a Receiver Operating Characteristic (ROC) Curve from within the *GLZ* results dialog. An ROC Curve can be used to help assess the goodness of fit for a binary predictor variable.

### Nonlinear Estimation

Changes have been made for *User-specified regression, least squares* and *User-specified regression, custom loss function*. While creating the estimated function, there is an option to review the variables. Now variables can be reviewed, selected, and inserted into the function via the *Review vars* button.

### Principal Components & Classification Analysis

The *Results* dialog - *Descriptives* tab contains a *2D scatterplots* button. You can now select multiple variables and generate all the scatterplots with one click.
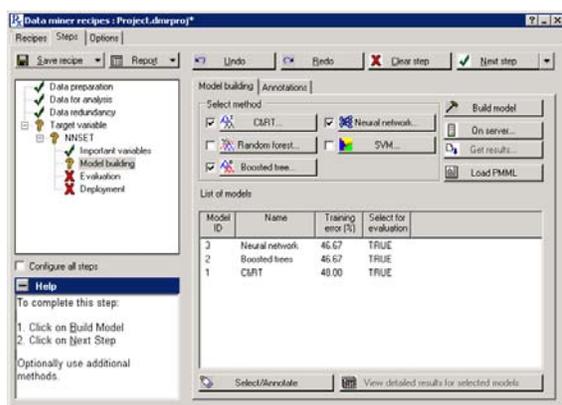
### Quality Control

Quality Control Charts now allow for one-sided control limits, including X, X-bar, MA, EWMA, Individuals with Moving Range, and CuSum, which can be deployed to *STATISTICA Enterprise*.

## Data Mining Enhancements

### Data Miner Recipes

*STATISTICA 9* features the long awaited final release of *Data Miner Recipes (DMR)*, previously available only as a beta release. DMR is an easy and flexible step-by-step data mining guide, and it is now



available to all of StatSoft's data miner customers. Novice data miners can quickly clean and analyze data, while advanced users can work more efficiently and have one more option to automate routine tasks. DMR explores the data and makes default decisions for you. You can easily modify these defaults as needed and save them for repeated use.

### Other Enhancements

New shortcut keys have been added for connecting data miner nodes in the workspace, and further drag/drop functionality was also added. Previous versions of *STATISTICA* required you to click on the second

node to complete a connection. Now you can just drop the connection arrow on top of the second node to complete a connection.

Large classification and regression tree displays are now scrollable.

The *MARSplines* results spreadsheet for outcome, where the variable name used to be truncated to 8 characters, now displays the whole variable name for the independent variable.

The *General Optimization* module (see the description, above), which we are now releasing as a beta version, is included in *STATISTICA Process Optimization*. This module enables users to optimize arbitrary functions using Simplex, Genetic Algorithm, or Grid-Search methods. The function to be optimized can be specified in a *STATISTICA Visual Basic (SVB)* formula or program.

This new module was specifically designed to easily call any *STATISTICA* (or other) functions. Therefore, optimization problems that involve multiple data mining prediction models and simulation can now be easily set up and solved efficiently.

## Other System Enhancements
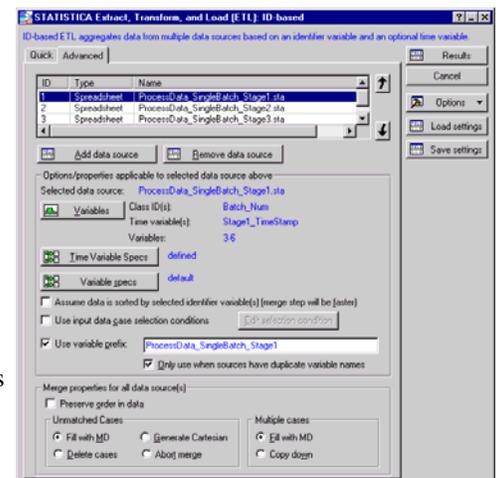
### Deployment of *STATISTICA* Models to SAS

If your company has a SAS system installed, you can now build your models in *STATISTICA* and then deploy them in the SAS environment. This new functionality offers several advantages. For example, (1) it enables you to deploy predictive data mining models in the SAS corporate environment even if you do not license SAS Enterprise Miner, and (2) it enables you to create predictive models that are not supported by and cannot be created in the SAS system (even if you do have a SAS Enterprise Miner license).

### *STATISTICA* ETL (Extract, Transform, and Load)

The new *STATISTICA Extract, Transform, and Load (ETL)* product offers powerful tools to align, merge, and intelligently combine data from databases and submit them to the powerful *STATISTICA* data processing capabilities for data filtering, aggregation, alignment, and analyses.



ID-Based *STATISTICA ETL* can be used to align data from disparate sources by batch number and time interval, and/or by one or more ID fields. Merging many-to-one data sets is a common scenario in many process and manufacturing industries.

For example, using *STATISTICA ETL* it is easy to pre-process batch manufacturing data to constant batch length (solve the "unequal batch length" problem for *PLS/PCA* model-based quality control with *STATISTICA* or on-line real-time *STATISTICA Enterprise*). Other typical applications include processing data collected at different time intervals that need to be aligned either by aggregating the values for the variables collected at the higher frequency, or by replicating the values for the variables collected at the lower frequency (e.g., when analyzing historical process data describing the performance of a

furnace, to align combustion parameters collected at one-minute intervals with fuel quality data collected daily). In version 9, the product has been further updated for improved performance and scalability and more detailed reporting (e.g., to report the actual intervals in the output results).

Plus, you can now use these *ETL* options directly from within the *Multivariate Statistical Process Control (MSPC)* module.

### STATISTICA Enterprise (Reports and More)

*STATISTICA Enterprise*'s reporting interface allows you to create reports in HTML, PDF, and RTF formats. Now the system can automatically email these reports when they are run.

*STATISTICA Enterprise* also offers a number of other enhancements in version 9. There is a new database connection role for users who are not system administrators. When using QC analysis configurations, users can define not only two-sided but also one-sided control limits.

You can now add an existing Dashboard to a Dashboard definition via the *Dashboard Admin - Add Task*, using the *Add Dashboard* button.

Changes have also been made to the audit logging facilities to always record audit log times in universal coordinated time. The conversion is made to the local time zone when the audit log is displayed. This enables users in different time zones to have an accurate view of what changed when.

### STATISTICA Web Data Entry

*STATISTICA Web Data Entry* is a product and a development system enabling users to configure data entry scenarios for data entry via Web browsers. *STATISTICA Web Data Entry* builds on, and seamlessly integrates with, the configuration objects in *STATISTICA Enterprise.*

### WebSTATISTICA Licensing

*WebSTATISTICA* has always been licensed per-processor. Now the licensing system has been extended to support splitting the license into separate *WebSTATISTICA* instances on separate servers. For example, a site licensed for eight CPUs can either deploy this license on a single eight-CPU server or on two separate four-CPU servers pointing at the same license file.

### Microsoft Installer (MSI) Support

With the release of version 9, we have changed the *STATISTICA* application installation platform to Microsoft Installer (MSI) instead of InstallShield installer that we used in version 8 and before.

When installing interactively, the user experience is similar to that from version 8, but the dialogs are more attractive.

However, the real benefit from using MSI is in how the *STATISTICA* installer can now be integrated into other installation packages and enterprise installation tools. The MSI allows for a totally "silent" installation, where all the information the user enters (CDKEY, serial number, netID, install code, and user registration information) can be passed either directly on the command line as command-line parameters or in a parameter file that the command line references.

There are three main use cases for the MSI installer:

1. Integrating our installer with other installers.

2. Working with enterprise-level deployment solutions for single-user installations.

3. Working with enterprise-level deployment solutions for concurrent workstation installations.

For more information on these new installation options, please contact your local StatSoft office.

## Expanded Interfaces for Developers

In version 9 of *STATISTICA*, a variety of enhancements have been added for developers and system integrators.

1. A new lightweight *STATISTICA* Spreadsheet library is now distributed with *STATISTICA*. It is available at no cost to third party developers who intend to read or write *STATISTICA* data files. It is multi-threaded and has a separate multi-threaded library for .NET access.

2. Graphs have a new event interface: OnGroupingSelect. This event is used for graphs that are categorized or aggregated in some fashion, for instance, a histogram, a box whisker plot, or categorized graphs. When the user selects items in the plots, the OnGroupingSelect event is fired to provide information about what groups/categories that selection represents. This new interface now allows applications using *STATISTICA* graphs to implement drill-down capability.

3. *STATISTICA* will no longer by default include all macro references in a newly created macro. Instead, each individual module will add its specific references when a macro is recorded. Suppressing all the references makes macros start up more quickly, but you could run into unresolved references if you copy/paste code from one macro into another. Therefore, the program will now check when you are copying/pasting between macros that have a different list of references, and offer to copy the additional references.

4. A new command-line parameter /MacroArgument has been added that can be used in conjunction with the /RunMacro argument. This will allow you to pass a parameter to the macro being run, which the macro can access with the GetScriptArgument call.

5. *STATISTICA* offers a comprehensive set of integration options to use with procedures written in R, which is a highly extensible programming language and environment for statistical computing (http://www.r-project.org). All recent versions of R (up to 2.10.1, the most current version as of this writing) can be executed within *STATISTICA*. R results can be displayed in native *STATISTICA* Spreadsheets and Graphs. A variety of integration options are offered including execution of R code on *STATISTICA* Servers. See the *Integration Options and Features to Leverage Specialized R Functionality in STATISTICA and WebSTATISTICA Solutions* white paper (http://www.statsoft.com/products/webserver.htm) for more details.

6. ANSI-92 SQL JOIN syntax is now supported by *STATISTICA Query*. Newer versions of SQL Server will require these types of joins to be used. By default, this option is not selected. You can set the option per query or for all queries.

Please contact StatSoft at **918-749-1119** or info@statsoft.com if you have any questions about your *STATISTICA 9* upgrade.