# Tips, Tricks, and What's New in IBM SPSS Modeler

Smarter software for a smarter planet

# Agenda, Goals for the Session, and Opening Remarks

- Welcome

- The Agenda for Today's Session includes General Buckets
  Addressing:
  - Naming
  - GUI – the User Interface
  - Sources – Accessing Data
  - Data Preparation
  - Modeling
  - Integration
  - Output

- The Goal for Today's Session
  - To take away at least one tip, trick, or fact about what's new
    that will improve the efficiency and/or effectiveness of your
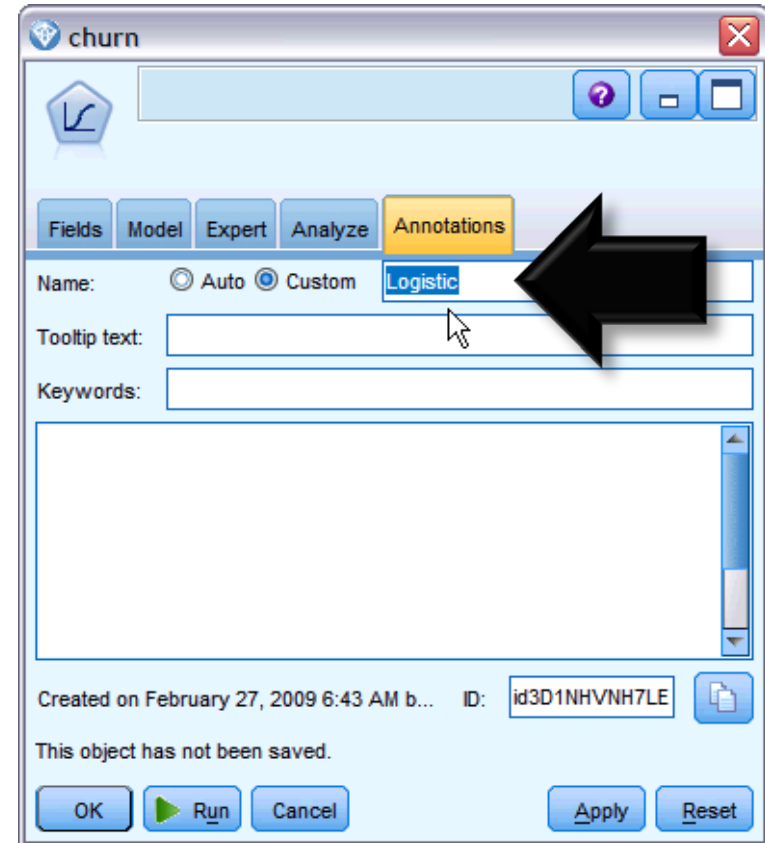    data mining / predictive analytics efforts

# Intro: Clarifying the Name of IBM SPSS Modeler

- Version 12: "Clementine "
  - Long-time users refer to Modeler as "Clementine" or "Clem"

- Version 13: "PASW Modeler"
  - "PASW" was an umbrella term used to signify that Modeler was one facet of a Predictive Analytics Software portfolio solution (along side Data Collection, Statistics, Collaboration & Deployment Services...)

- Version 14: "IBM SPSS Modeler"
  - With the acquisition, IBM well-understood the equity of the SPSS brand name



IBM® SPSS® Modeler 14.2

IBM® SPSS® Modeler
Version 14.2

Licensed Materials - Property of IBM Corp. © Copyright IBM Corporation and its licensors 1994, 2011. IBM, IBM logo, ibm.com, and SPSS are trademarks or registered trademarks of International Business Machines Corp., registered in many jurisdictions worldwide. A current list of IBM trademarks is available on the Web at www.ibm.com/legal/copytrade.shtml. Java and all Java-based trademarks and logos are trademarks or registered trademarks of Oracle and/or its affiliates. Other product and service names might be trademarks of IBM or other companies. This Program is licensed under the terms of the license agreement accompanying the Program. This license agreement may be either located in a Program directory folder or library identified as "License" or "Non_IBM_License", if applicable, or provided as a printed license agreement. Please read the agreement carefully before using the Program. By using the Program you agree to these terms.

Additional Details

OK

# GUI: Comments, Custom Name Attribute, and SuperNodes Provide Communication, Security

- **Comments**
  - A sticky notes-like feature that allows users to document thoughts and share details directly on the canvas

- **Custom Name** in the Attributes Tab
  - Specifies the name assigned to the model that is created when the node is executed.
    - Auto - generates the model name automatically
    - Custom - allows you to specify a custom name for the model created

- **SuperNodes**
  - Creating a SuperNode "shrinks" the stream by encapsulating several nodes into one

# GUI: Other Interface Shortcuts that Speed the Development and Execution of SPSS Modeler Streams

- Build streams quickly by double-clicking
  - Double-clicking a node on the palette will add and connect it to the current stream

- Use key combinations to select all downstream nodes
  - Pressing Ctrl-Q and Ctrl-W will toggle the selection of <u>all</u> nodes downstream.

- Use shortcut keys to connect and disconnect nodes
  - When a node is selected in the canvas, pressing:
    - F2 will begin a connection
    - Tab will move to the desired node
    - Shift-spacebar will complete the connection
    - F3 will disconnect all inputs and outputs to the selected node

- Customize the Nodes Palette tab with your favorite nodes
  - From the Tools menu, choose Manage Palettes to open a dialog box for adding, removing, or moving the nodes shown on the Nodes Palette.

# Sources: Importing, Exporting to Cognos BI Enables a Distribution of Predictive Analytics based on "One Version of the Truth"

- **IBM Cognos BI Source Node**
  - Bring Cognos BI database data into your data mining session
  - Combine business intelligence features with the predictive analytics capabilities
  - Import relational, dimensionally-modeled relational and OLAP data

- **IBM Cognos BI Export Node**
  - Export data from a stream to Cognos BI (UTF-8 format)
  - Enables Cognos BI to make use of transformed or scored data from Modeler
    - Specifically, saved to a Cognos BI server and distributed to Cognos users

# Sources: Save Time Connecting to a Preferred Database by Specifying a Default Database

- **Database Connection** - For database import and export users can specify a default connection and store details of different connections for display. A few highlights:
  - Data Sources
    - Lists the available data sources
  - Connections
    - Shows currently connected databases
  - Default
    - Optionally choose one connection as the default causing
    - This setting can be edited if desired.
  - Preset
    - Indicates (with a * character) whether preset values have been specified for the database connection

# Sources: Options in Sub-Dialog Boxes also help Save Time when Connecting to a Databases

- **Database Export Node**
  - The "Schema" sub-dialog box helps improve efficiency when exporting to an InfoSphere Warehouse database
  - The "Advanced" sub-dialog box enables technical details for exporting results to a database

# Data Prep: Reduce Time Spent Preparing Data for Analysis with Automated Data Preparation

- **Automated Data Preparation (ADP) Node**
  - Analyzes data
  - Identifies possible fixes
  - Screens out fields that are problematic or not likely to be useful
  - Derives new attributes when appropriate
  - Improves performance through intelligent screening techniques

# Data Prep: Enhancements to the Sample Node include Data Mining that is more Efficient and Effective

- **Sample Node**
  - Oracle or IBM DB2 database
    - Block-level sampling can be more efficient with random percentage sampling when performing in-database mining
  - Enhanced support for SQL generation (14.2)
    - Support for SQL generation in the Sample node when using Simple sampling has been enhanced for DB2, Netezza, and Teradata databases

| Mode | Sample | Max size | Seed | DB2 OS/Z | DB2 OS/400 | DB2 Windows/UNIX | Netezza | Oracle | SQL Server | Teradata |
|---|---|---|---|---|---|---|---|---|---|---|
| Include | First | n/a | | Y | Y | Y | Y | Y | Y | Y |
| | 1-in-n | off | | Y | Y | Y | Y | Y | | Y |
| | | max | | Y | Y | Y | Y | Y | | Y |
| | Random % | off | off | | | Y | Y | Y | | Y |
| | | | on | | | Y | | Y | | |
| | | max | off | | | Y | Y | Y | | Y |
| | | | on | | | Y | | Y | | |
| Discard | First | off | | | | | Y | Y | | |
| | | max | | | | | Y | Y | | |
| | 1-in-n | off | | Y | Y | Y | Y | Y | | Y |
| | | max | | Y | Y | Y | Y | Y | | Y |
| | Random % | off | off | | | Y | Y | Y | | Y |
| | | | on | | | Y | | Y | | |
| | | max | off | | | Y | Y | Y | | Y |
| | | | on | | | Y | | Y | | |

# Data Prep: Leveraging CLEM Speeds the Data Mining Process

- Background on CLEM
  - Control Language for Expression Manipulation
  - Powerful language for analyzing and manipulating the data that flows along streams
  - Perform both simple and complex tasks
    - Simple: Deriving profit from cost and revenue data
    - Complex: Transforming web log data into a set of fields and records with usable information.

- New to Version 14.1
  - Within the CLEM language a number of date and time functions now additionally support the use of timestamp as an argument

# Modeling: In-Database Mining Offers Improved Performance and Easy Deployment



- **In-Database Mining**
  - Server supports integration with data mining and modeling tools that are available from database vendors, including:
    - Oracle Data Miner
    - IBM DB2 InfoSphere Warehouse
    - Microsoft Analysis Services
    - IBM Netezza.
  - Build, score, and store models inside the database—all from within the Modeler

- Advantages of database-native algorithms include:
  - Improved performance
  - Easy deployment to / sharing with application that access the database

- **Microsoft Time Series**



- **Microsoft Sequence Clustering**



- **Oracle Attribute Importance**

# Modeling: Using SPSS Modeler as a Front-end is a Primary Benefit of Integration with a Database (Example: Netezza)



Using Modeler as a front-end enables easy connection to Netezza data

A visual front-end coupled with SQL pushback means rapid preparation – especially for data sets with tens of millions of records

Modeler's various output nodes provide an easy way to interact with analytics running within the Netezza appliance

Data manipulation and preparation can be implemented visually without programming

Using Modeler as a front-end enables easy and rapid fine-tuning of the In-Database Analytics that is running inside the Netezza appliance

# Modeling: Automated Modeling Enables a Side-by-Side Comparison of Model Effectiveness

- **Auto Classifier, Auto Cluster, and Auto Numeric Nodes**
  - Standard a group of ensemble modeling nodes
  - Automate the building of a number of different models concurrently
  - Compare the results and choose the best model for data

  - CLEF provides the AutoModeling element to enable a model specified by the ModelBuilder element to be used by any of these ensemble nodes

# Integration: Modeler Now Fully Integrated with IBM SPSS Statistics, Modeler Premium now includes Text Analytics

- **IBM SPSS Text Analytics**
  - Advanced linguistic technologies and Natural Language Processing (NLP)
  - Rapidly process a large variety of unstructured text data, extract and organize the key concepts, and group these concepts into categories
  - Extracted concepts and categories can be combined with existing structured data

- **IBM SPSS Statistics**
  - Fully integrated within the Modeler environment

# Output: Analysis Node can be used to Compare Models, Results can be Written to a File
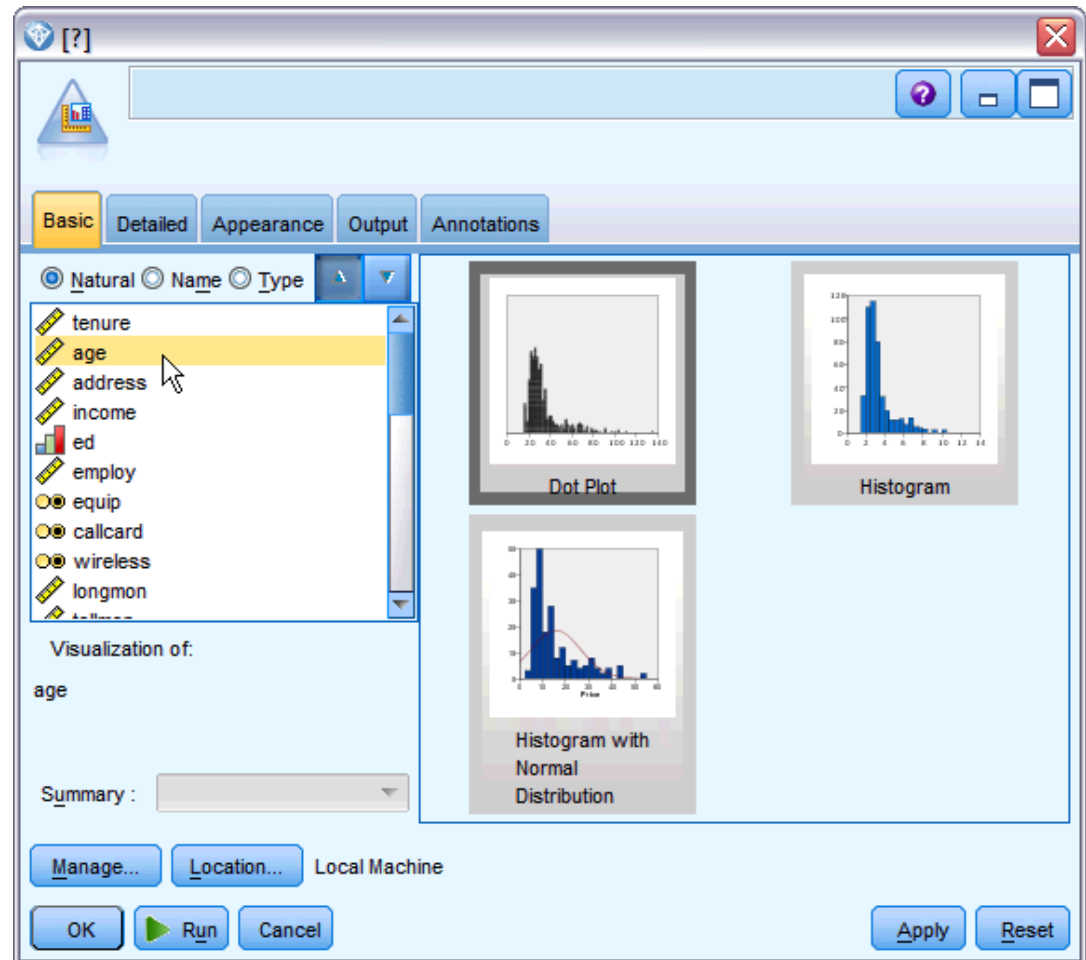
- **Analysis Node**
  - Evaluate the ability of a model to generate accurate predictions.
  - Perform various comparisons between predicted values and actual values
  - Tip: Analysis nodes can be used to compare predictive models to other predictive models
  - Upon execution a summary of the analysis results automatically added to the Analysis section on the Summary tab
  - Detailed analysis results appear on the Outputs tab of the manager window or can be written directly to a file

# Output: The Graphboard Node Presents a Variety of Different Visualizations based on Data Types

- **Graphboard Node**
  - Choose from many different graphs outputs in a single node
  - Node presents you with a choice of graph types that work for your data
  - Automatically filters out any graph types that would not work with the field choices

# Recap: Agenda, Goals for the Session, and Next Steps

- Today's Session included General Buckets Addressing:
  - Naming
  - GUI – the User Interface
  - Sources – Accessing Data
  - Data Preparation
  - Modeling
  - Integration
  - Output

- The Goal for Today's Session was to
  - Take away at least one tip, trick, or fact about what's new that will improve the efficiency and/or effectiveness of your data mining / predictive analytics efforts

- Next Steps
  - For More Information Call (800) 543-2185