

# TIBCO Statistica®

## Stepwise Model Builder Formula Guide

*Version 14.0*

*December 2020*



# Contents

---

<b>TIBCO Documentation and Support Services .....</b>	<b>3</b>
<b>Stepwise Model Builder Overview .....</b>	<b>5</b>
<b>Notation .....</b>	<b>6</b>
<b>Model .....</b>	<b>7</b>
<b>Estimation .....</b>	<b>8</b>
<b>Statistics .....</b>	<b>9</b>
<b>Legal and Third-Party Notices .....</b>	<b>12</b>

# TIBCO Documentation and Support Services

---

## How to Access TIBCO Documentation

Documentation for TIBCO products is available on the TIBCO Product Documentation website, mainly in HTML and PDF formats.

The TIBCO Product Documentation website is updated frequently and is more current than any other documentation included with the product. To access the latest documentation, visit <https://docs.tibco.com>.

## Product-Specific Documentation

Documentation for TIBCO Statistica® is available on the [TIBCO Statistica® Product Documentation](#) page.

The following documents for this product can be found on the TIBCO Documentation site:

- *TIBCO Statistica® Release Notes*
- *TIBCO Statistica® Installation*
- *TIBCO Statistica® Quick Reference*
- *TIBCO Statistica® Product Traceability*
- *TIBCO Statistica® Configuration for Windows Server 2019*
- *TIBCO Statistica® Data Entry Administration*
- *TIBCO Statistica® Server Administrator's Guide*
- *TIBCO Statistica® Options Configuration*
- *TIBCO Statistica® R Integration: Features and Options*
- *TIBCO Statistica® Security Guide*
- *TIBCO Statistica® User's Guide*
- *TIBCO Statistica® Enterprise Manager Guide*
- *TIBCO Statistica® Statistica Object Model Guide*
- *TIBCO Statistica® Logistic Regression Formula Guide*
- *TIBCO Statistica® Stability Analysis Formula Guide*
- *TIBCO Statistica® Stepwise Model Builder Formula Guide*
- *TIBCO Statistica® Weight of Evidence Formula Guide*

The following documents for TIBCO® Data Science for TIBCO Spotfire® Analyst can be found on the [TIBCO Documentation site](#):

- *TIBCO® Data Science for TIBCO Spotfire® Analyst Release Notes*
- *TIBCO® Data Science for TIBCO Spotfire® Analyst User Guide*

The following documents for TIBCO® Data Science Service for TIBCO Spotfire® can be found on the [TIBCO Documentation site](#):

- *TIBCO® Data Science Service for TIBCO Spotfire® Release Notes*
- *TIBCO® Data Science Service for TIBCO Spotfire® Installation and Configuration Guide*
- *TIBCO® Data Science Service for TIBCO Spotfire® User Guide*

## How to Contact TIBCO Support

You can contact TIBCO Support in the following ways:

- For an overview of TIBCO Support, visit <http://www.tibco.com/services/support>.
- For accessing the Support Knowledge Base and getting personalized content about products you are interested in, visit the TIBCO Support portal at <https://support.tibco.com>.
- For creating a Support case, you must have a valid maintenance or support contract with TIBCO. You also need a user name and password to log in to <https://support.tibco.com>. If you do not have a user name, you can request one by clicking Register on the website.

## How to Join TIBCO Community

TIBCO Community is the official channel for TIBCO customers, partners, and employee subject matter experts to share and access their collective experience. TIBCO Community offers access to Q&A forums, product wikis, and best practices. It also offers access to extensions, adapters, solution accelerators, and tools that extend and enable customers to gain full value from TIBCO products. In addition, users can submit and vote on feature requests from within the [TIBCO Ideas Portal](#). For a free registration, go to <https://community.tibco.com>.

## Stepwise Model Builder Overview

---

The Stepwise Model Builder computes the marginal predictor statistics given a current model.

Specifically, the variables listed in the Marginal results table are entered one at a time into a logistic regression model containing the predictors listed in the Model results table. This enables the analyst to evaluate the unique contribution of each predictor candidate not in the equation. The model is estimated after recoding all Bad code values to 1 and Good code values to 0.

## Notation

The following notation is used throughout this model section of this document:

Notation	Description
$n$	Number of observed cases
$p$	Number of parameters
$y$	$n \times 1$ vector with $y_i$ being the observed value of the $i$ th case of the chosen dichotomous good/bad dependent variable
$x$	$n \times p$ matrix with $x_{ij}$ being the observed value of the $i$ th case of the $j$ th parameter
$\beta$	$p \times 1$ vector with $\beta_j$ being the coefficient for the $j$ th parameter
$w$	$n \times 1$ vector with $w_i$ being the weight for the $i$ th case.
$l$	Likelihood function
$L$	Log likelihood function
$I$	Information matrix

## Model

---

The Stepwise Model Builder uses the linear logistic model.

This model has a dichotomous dependent variable, and in this module the outcome is labeled as either having a Good or Bad outcome. Included in the Stepwise Model Builder dialog box are options for Dependent Variable (Y), Good code and Bad code. For the model, the dependent variable is assumed to have a probability  $\pi$ , where for the  $i$ th case:

$$\pi_i = \frac{\exp(\eta_i)}{1 + \exp(\eta_i)}$$

or

$$\ln\left(\frac{\pi_i}{1 - \pi_i}\right) = \eta_i = \mathbf{X}_i' \boldsymbol{\beta}$$

For  $n$  observations,  $y_1$  through  $y_n$ , with probabilities  $\pi_1$  through  $\pi_n$  and case weights  $w_1$  through  $w_n$ , the likelihood function is:

$$l = \prod_{i=1}^n \pi_i^{w_i y_i} (1 - \pi_i)^{w_i (1 - y_i)}$$

The logarithm of  $l$  is:

$$L = \ln(l) = \sum_{i=1}^n (w_i y_i \ln(\pi_i) + w_i (1 - y_i) \ln(1 - \pi_i))$$

The derivative of  $L$  with respect to  $\beta_j$  is:

$$L_{x_j}^* = \frac{\partial L}{\partial \beta_j} = \sum_{i=1}^n w_i (y_i - \pi_i) x_{ij}$$

## Estimation

---

Maximum Likelihood Estimation is achieved through the Fisher Scoring algorithm.

$$\beta_{k+1} = \beta_k + I^{-1}S$$

$$I_{pxp} = -E_{\beta} \left[ \frac{\partial^2 l(\beta)}{\partial \beta^2} \right]$$

$$S = \frac{\partial l(\beta)}{\partial \beta}$$



# Statistics

## Estimated Variance Covariance Matrix

The estimated covariance matrix is the inverse of the information matrix (negative of the expected Hessian) evaluated at the MLE values of the parameters.

$$I(\theta_{MLE})^{-1}$$

$$I(\theta_{MLE}) = -E_{\theta}[H(\theta)]$$

$$H(\theta) = \frac{\partial^2 l(\theta)}{\partial \theta^2}$$

## Estimated Correlation Matrix

The estimated correlation matrix is the standardized version of the covariance matrix, that is, all entries are divided by the product of the standard deviations.

## Gini Coefficient

$$G = 2 \left( \frac{\text{Number of Bads}}{N} \right) \left( \frac{\text{Number of Goods}}{N} \right)$$

Notation:

$N$  = Total number of observations

## Hosmer-Lemeshow (HL) Goodness of Fit Statistic

$$H = \sum_{g=1}^n \frac{(O_g - E_g)^2}{N_g \pi_g (1 - \pi_g)}$$

Notation:  $O_g$  = Observed event  $g$

$E_g$  = Expected event  $g$

$N_g$  = Observations of event  $g$

$\pi_g$  = Predicted risk for the  $g^{th}$  risk decile group

$n$  = number of groups



The Hosmer-Lemeshow statistic is asymptotically distributed and follows a  $\chi^2$  distribution with  $n-2$  degrees of freedom.

## Kolmogorov-Smirnov (KS) test

For all Good observations, predicted probability of Bad is computed, that is the relative frequency of bad cases in the bin a Good observation is placed. This process is repeated for all Bad observations. The KS test is then completed with the Good/Bad indicator as the group variable and the predicted probability of Bad as the response.

$$Z = \max_j |D_j| \sqrt{\frac{n_1 n_2}{n_1 + n_2}}$$

Significance level (p) approximation is based on the formula:

$$p = 2 \sum_{i=1}^{\infty} (-1)^{i-1} e^{-2i^2 \left( \frac{KS \sqrt{\frac{n_1}{n_1+n_2} + 0.12 + 0.11}}{\sqrt{\frac{n_1}{n_1+n_2}}} \right)^2}$$

**Lift Value**

$$\text{Lift} = \frac{\text{Result Predicted by Model}}{\text{Result Predicted with No Model}}$$

**ROC - Area Under Curve (AUC)**

$$\text{AUC} = \frac{G + 1}{2}$$

Notation:  $G$  = Gini coefficient

**ROC - Sensitivity**

$$\text{Sensitivity} = \frac{\text{number of true positives}}{\text{number of true positives} + \text{number of false negatives}}$$

Note: If a test has two outcomes, positive and negative:

- True Positive ☐ Both the observed and predicted response is positive.
- False Positive ☐ Predicted response is positive but the observed response is negative.
- True Negative ☐ Both the observed and predicted response is negative.
- False Negative ☐ Predicted response is negative but the observed response is positive.

**ROC - Specificity**

$$\text{Specificity} = \frac{\text{number of true negatives}}{\text{number of true negatives} + \text{number of false positives}}$$

**Somers' D**

If ties are present:

$$d = \frac{(n_c - n_d)}{t}$$

If ties are not present:

$$d = 2c - 1 \text{ where } c = (n_c + 0.5(t - n_c - n_d)) / t$$

**Wald Statistic**

For continuous variables:

$$W_i = \left( \frac{\beta_i}{SE_{\beta_i}} \right)^2$$

For categorical variables:

If  $\beta_i$  is a vector of MLEs associated with  $m-1$  dummy coded variables, and  $C$  is the asymptotic covariance matrix for  $\beta_i$ , the Wald statistic is calculated as:

$$W_i = \beta_i' C^{-1} \beta_i$$



Asymptotically distributed as a  $\chi^2$  distribution with degrees of freedom equal to the number of parameters estimated and is analogous to the t-test in linear regression.

**Wald Statistic - Standard Error**

The standard error (SE) is the square root of the  $i^{\text{th}}$  diagonal entry of the inverse information matrix.

**Wald Statistic Confidence Interval**

$$100 (1 - \alpha)\% \text{ CI for } \beta_i = \hat{\beta}_i \pm z_{1-\alpha/2} SE_{\beta_i}$$

Notation:  $z_{1-\alpha/2}$  = 100 (1 -  $\alpha/2$ )th percentile of the standard normal distribution

$\hat{\beta}_i$  = Estimate of parameter  $\beta_i$

$SE_{\beta_i}$  = Standard error estimate of  $\hat{\beta}_i$

## Legal and Third-Party Notices

---

SOME TIBCO SOFTWARE EMBEDS OR BUNDLES OTHER TIBCO SOFTWARE. USE OF SUCH EMBEDDED OR BUNDLED TIBCO SOFTWARE IS SOLELY TO ENABLE THE FUNCTIONALITY (OR PROVIDE LIMITED ADD-ON FUNCTIONALITY) OF THE LICENSED TIBCO SOFTWARE. THE EMBEDDED OR BUNDLED SOFTWARE IS NOT LICENSED TO BE USED OR ACCESSED BY ANY OTHER TIBCO SOFTWARE OR FOR ANY OTHER PURPOSE.

USE OF TIBCO SOFTWARE AND THIS DOCUMENT IS SUBJECT TO THE TERMS AND CONDITIONS OF A LICENSE AGREEMENT FOUND IN EITHER A SEPARATELY EXECUTED SOFTWARE LICENSE AGREEMENT, OR, IF THERE IS NO SUCH SEPARATE AGREEMENT, THE CLICKWRAP END USER LICENSE AGREEMENT WHICH IS DISPLAYED DURING DOWNLOAD OR INSTALLATION OF THE SOFTWARE (AND WHICH IS DUPLICATED IN THE LICENSE FILE) OR IF THERE IS NO SUCH SOFTWARE LICENSE AGREEMENT OR CLICKWRAP END USER LICENSE AGREEMENT, THE LICENSE(S) LOCATED IN THE "LICENSE" FILE(S) OF THE SOFTWARE. USE OF THIS DOCUMENT IS SUBJECT TO THOSE TERMS AND CONDITIONS, AND YOUR USE HEREOF SHALL CONSTITUTE ACCEPTANCE OF AND AN AGREEMENT TO BE BOUND BY THE SAME.

This document is subject to U.S. and international copyright laws and treaties. No part of this document may be reproduced in any form without the written authorization of TIBCO Software Inc.

TIBCO, the TIBCO logo, the TIBCO O logo, Statistica, Spotfire, Process Tree Viewer, Process Data Explorer, Predictive Claims Flow, Making the World More Productive, Live Score, Electronic Statistics Textbook, Decisioning Platform, Data Health Check, and Better Decisioning are either registered trademarks or trademarks of TIBCO Software Inc. and/or its subsidiaries in the United States and/or other countries.

Java and all Java based trademarks and logos are trademarks or registered trademarks of Oracle and/or its affiliates.

All other product and company names and marks mentioned in this document are the property of their respective owners and are mentioned for identification purposes only.

This software may be available on multiple operating systems. However, not all operating system platforms for a specific software version are released at the same time. Please see the readme.txt file for the availability of this software version on a specific operating system platform.

THIS DOCUMENT IS PROVIDED "AS IS" WITHOUT WARRANTY OF ANY KIND, EITHER EXPRESS OR IMPLIED, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE, OR NON-INFRINGEMENT.

THIS DOCUMENT COULD INCLUDE TECHNICAL INACCURACIES OR TYPOGRAPHICAL ERRORS. CHANGES ARE PERIODICALLY ADDED TO THE INFORMATION HEREIN; THESE CHANGES WILL BE INCORPORATED IN NEW EDITIONS OF THIS DOCUMENT. TIBCO SOFTWARE INC. MAY MAKE IMPROVEMENTS AND/OR CHANGES IN THE PRODUCT(S) AND/OR THE PROGRAM(S) DESCRIBED IN THIS DOCUMENT AT ANY TIME.

THE CONTENTS OF THIS DOCUMENT MAY BE MODIFIED AND/OR QUALIFIED, DIRECTLY OR INDIRECTLY, BY OTHER DOCUMENTATION WHICH ACCOMPANIES THIS SOFTWARE, INCLUDING BUT NOT LIMITED TO ANY RELEASE NOTES AND "READ ME" FILES.

This and other products of TIBCO Software Inc. may be covered by registered patents. Please refer to TIBCO's Virtual Patent Marking document (<https://www.tibco.com/patents>) for details.

Copyright © 1995-2020. TIBCO Software Inc. All Rights Reserved.