# HANDBOOK OF
# NATURAL LANGUAGE PROCESSING
## SECOND EDITION

Chapman & Hall/CRC
Machine Learning & Pattern Recognition Series

SERIES EDITORS

**Ralf Herbrich and Thore Graepel**
Microsoft Research Ltd.
Cambridge, UK

**AIMS AND SCOPE**

This series reflects the latest advances and applications in machine learning and pattern recognition through the publication of a broad range of reference works, textbooks, and handbooks. The inclusion of concrete examples, applications, and methods is highly encouraged. The scope of the series includes, but is not limited to, titles in the areas of machine learning, pattern recognition, computational intelligence, robotics, computational/statistical learning theory, natural language processing, computer vision, game AI, game theory, neural networks, computational neuroscience, and other relevant topics, such as machine learning applied to bioinformatics or cognitive science, which might be proposed by potential contributors.

**PUBLISHED TITLES**

MACHINE LEARNING: An Algorithmic Perspective
*Stephen Marsland*

HANDBOOK OF NATURAL LANGUAGE PROCESSING,
Second Edition
*Nitin Indurkhya and Fred J. Damerau*

# HANDBOOK OF
# NATURAL LANGUAGE PROCESSING

## SECOND EDITION

Edited by
NITIN INDURKHYA
FRED J. DAMERAU

**Visit the Taylor & Francis Web site at**
**http://www.taylorandfrancis.com**

**and the CRC Press Web site at**
**http://www.crcpress.com**

*To Fred Damerau*
*born December 25, 1931; died January 27, 2009*

Some enduring publications:

Damerau, F. 1964. A technique for computer detection and correction of spelling errors. *Commun. ACM* 7, 3 (Mar. 1964), 171–176.

Damerau, F. 1971. *Markov Models and Linguistic Theory: An Experimental Study of a Model for English*. The Hague, the Netherlands: Mouton.

Damerau, F. 1985. Problems and some solutions in customization of natural language database front ends. *ACM Trans. Inf. Syst.* 3, 2 (Apr. 1985), 165–184.

Apté, C., Damerau, F., and Weiss, S. 1994. Automated learning of decision rules for text categorization. *ACM Trans. Inf. Syst.* 12, 3 (Jul. 1994), 233–251.

Weiss, S., Indurkhya, N., Zhang, T., and Damerau, F. 2005. *Text Mining: Predictive Methods for Analyzing Unstructured Information*. New York: Springer.

# Contents

## PART I    Classical Approaches

## PART II    Empirical and Statistical Approaches

## PART III   Applications

# List of Figures