# 37th Symposium on the Interface of Computing Science and Statistics 2005 and the Annual Meeting of the Classification Society of North America

## (Interface/CSNA 2005)

**Abstracts**

**St. Louis, Missouri, USA**
**8-12 June 2005**

# TABLE OF CONTENTS

## CLUSTERING

## MASS SPECTROMETRY BASED PROTEOMICS

## DETECTING GENE-GENE INTERACTIONS

## INFERENCE IN A CASE-CONTROL STUDY WITH 100,000 SNPS

## HIGH-DIMENSIONAL BIOMEDICAL DATA

## MIXTURE MODELS

## DATABASES AND COMPUTERS

## APPLICATIONS IN MEDICINE

## BEST OF SIAM DATA MINING CONFERENCE

## MODELING ALCOHOL ABUSE

## TIME SERIES, NEURAL NETWORKS, HIERARCHICAL MODELS

## MODEL AVERAGING AND ASSESSMENT

## BIOINFORMATICS

## CLUSTER VALIDATION

## SELECTED IASC PAPERS

## MODEL-BASED/GRAPH-THEORETIC METHODS

## CLUSTERING AND CLASSIFYING TEXT

## MODEL-BASED CLUSTERING AND CLASSIFICATION

## MICROARRAYS

## APPLICATIONS

## COMPARING CLASSIFICATION METHODS

## QUANTILE REGRESSION

## SPECTRAL METHODS IN DATA ANALYSIS

## MULTIDIMENSIONAL SCALING AND ET ALIA

## BEST OF THE JOURNAL OF CLASSIFICATION

## NON-NUMERIC DATA ANALYSIS

## DISCRIMINATION AND DATA MINING

## CLUSTERING METHODS AND APPLICATIONS

## COMPUTATIONAL BIOLOGY

## GRAPH-THEORETIC PATTERN RECOGNITION

## ALGORITHMS

## <u>AUTHOR IDENTIFICATION</u>