# CONTENTS

# LIST OF TABLES

# LIST OF FIGURES

x

xi

# LIST OF ABBREVIATIONS

| | |
|---|---|
| ANN | Artificial Neural Networks |
| BDK | Bi-Directional Kohonen Network |
| BMU | Best Matching Unit |
| %CC | Percentage of Correctly Classified |
| CPN | Counter Propagation Network |
| EDA | Exploratory Data Analysis |
| EDXRF | Energy Dispersive X-Ray Fluorescence |
| GC-MS | Gas Chromatography Mass Spectroscopy |
| GSOM | Growing Self-Organizing Map |
| HSOM | Hierarchical Self-Organizing Map |
| KDML 105 | *Oriza sativa* L. cv. Khao Dawk Mali 105 |
| *k*-NN | *k*-Nearest Neighbors |
| LDA | Linear Discriminant Analysis |
| %MS | Percentage of Model Stability |
| MQE | Mean Quantization Error |
| %MQE | Percentage of Mean Quantization Error |
| MSOMs | Multiple Self-Organizing Maps |
| N | North Region of Thailand |
| NE | Northeast Region of Thailand |
| NIR | Near Infrared |
| OM | Organic Matter |
| %OM | Percentage Organic Matter |
| %PA | Percentage Predictive Ability |
| PCA | Principal Component Analysis |

| | |
|---|---|
| PCs | Principal Components |
| PLS | Partial Least Squares |
| PLS-DA | Partial Least Squares Discriminant Analysis |
| PT 1 | *Oriza sativa* L. cv. Pathumthani 1 |
| QDA | Quadratic Discriminant Analysis |
| QE | Quantization Error |
| RBFN | Radial Basis Function Network |
| RMSECV | Root Mean Square Error of Cross Validation |
| SKN | Supervised Kohonen Network |
| SIMCA | Soft Independent Modelling of Class Analogy |
| SOM | Self-Organizing map |
| SQA | Soil Quality Assessment |
| SQIs | Soil Quality Indicators |
| SSOM | Single Self-Organizing map |
| U-matrix | Unified distance matrix |
| USDA | United States Department of Agriculture |
| WRB | World Reference Based for Soil Resources |
| XYF | *XY*-Fused Network |

xiii

# LIST OF SYMBOLS

| | |
|---|---|
| *A* | Number of principal components |
| $c_g$ | Centroid of a class membership $g$ |
| $E$ | Residual |
| $F_{corrected}$ | Number of time that a sample was correctly classified |
| $F_{picked}$ | Number of time that the sample was picked as a test sample |
| *I* | Number of samples |
| *J* | Number of parameters |
| $P$ | Loadings |
| *P* | Number of rows of initial map |
| *Q* | Number of columns of initial map |
| *T* | Iterations |
| $T$ | Scores |
| $s_{ij}$ | Standardized of sample $i$th of parameter $j$th |
| $s_{(\mathbf{w}_k, \mathbf{c}_g)}$ | Dissimilarity between a centroid of class $g$ and a map unit $w_k$ |
| $X$ | Data matrix |
| $x_{ij}$ | Sample $i$th of parameter $j$th |
| $\bar{x}_j$ | Mean for variable $j$ calculated over all $I$ samples |
| $x_r$ | A random sample |
| $W$ | Weight matrix |
| $w_k$ | Weight vector of sample $k$ |

xiv