# Contents

# List of Tables

# List of Figures

# List of Symbols

$H$- Hash Family

$R$- R-neighborhood

$Pr$- Probability

$h()$- Hash Function

$X$- Random Number variable

$U$- Universe of Discourse

$l_p^d$- $L_p$ in $d$-dimension

$WOS$- Weakest Outlier Score

$\Re$- Real Space

$d_{thres}$- Distance Threshold

$H(C_k)$- Within-cluster entropy of cluster $C_k$