

Mathematical Pseudocode – Anomaly-Based Duplicate Detection: A Probabilistic Approach

Andreas Obermeier^[0000-0002-0048-8208]

University of Ulm, 89069 Ulm, Germany
andreas.obermeier@uni-ulm.de

1 Interval Method

algorithm interval-method **is**

input: set of feature vectors $\{\zeta_{i,j} \in \mathbb{R}^f\}$ in the dataset to be analyzed
set of feature vectors $\{\hat{\zeta}_{i,j} \in \mathbb{R}^f\}$ in duplicate-free training data
multidimensional interval $I_h \subset \mathbb{R}^f$
number of samples N

output: duplicate probability $P(D|\zeta_{i,j} \in I_h)$ for pairs of records in I_h

$m \leftarrow |\{\zeta_{i,j} \in \mathbb{R}^f\}|$ // sample size = size of the dataset to be analyzed
 $q_h \leftarrow |\{\zeta_{i,j} | \zeta_{i,j} \in I_h\}|$ // count of feature vectors in I_h

for $n=1$ **to** N :
 $\{\hat{\zeta}_{i,j}^s\} \leftarrow$ sample of m feature vectors from $\{\hat{\zeta}_{i,j} \in \mathbb{R}^f\}$
 $\hat{q}_{h,n} \leftarrow |\{\hat{\zeta}_{i,j}^s | \hat{\zeta}_{i,j}^s \in I_h\}|$ // count of feature vectors in I_h in sample n

for $k=1$ **to** m :
 $\hat{P}(\hat{q}_h = k) \leftarrow \frac{1}{N} \sum_{n=1}^N \mathbb{I}_{\{\hat{q}_{h,n}=k\}}$ // probability of k counts in training data

$\hat{E}(\hat{q}_h) \leftarrow \sum_{k=1}^m k \cdot \hat{P}(\hat{q}_h = k)$ // expected value in duplicate-free data

if $q_h > \hat{E}(\hat{q}_h)$: // test for anomaly
 $P(D|\zeta_{i,j} \in I_h) \leftarrow \frac{\sum_{\hat{q}_h^D=0}^{q_h} \hat{P}(\hat{q}_h = q_h - \hat{q}_h^D) \cdot \hat{q}_h^D}{q_h}$ // estimated duplicate probability
 return $P(D|\zeta_{i,j} \in I_h)$

else:
 return 0 // no anomaly, duplicate probability zero

2 Kernel Density Estimation Method

algorithm kde-method **is**

input: set of feature vectors $\{\zeta_l \in \mathbb{R}^f\}$ in the dataset to be analyzed with
 $(1 \leq l \leq \text{dataset size})$ (for example, $\zeta_1 = \zeta_{1,2}$)
feature vector $\zeta_{i,j} \in \{\zeta_l\}$ of pair of records in question
set of feature vectors $\{\hat{\zeta}_l \in \mathbb{R}^f\}$ in duplicate-free training data
kernel function K for kernel density estimation
number of samples N

output: duplicate probability $P(D|\zeta_{i,j})$ for a specific pair of records with
feature vector $\zeta_{i,j}$

function KDE($x, \{x_i\}$): // function for Kernel Density Estimation
 $L \leftarrow |\{x_i\}|$
return $\frac{1}{L} \sum_{i=1}^L K(x - x_i)$
end function

$m \leftarrow |\{\zeta_{i,j} \in \mathbb{R}^f\}|$ // sample size = size of the dataset to be analyzed
 $q_{i,j} \leftarrow \text{KDE}(\zeta_{i,j}, \{\zeta_l\})$ // estimated density at $\zeta_{i,j}$ in the dataset to be analyzed

for $n=1$ **to** N :
 $\{\hat{\zeta}_l^s\} \leftarrow$ sample of m feature vectors from $\{\hat{\zeta}_l\}$
 $\hat{q}_{i,j}^{(n)} \leftarrow \text{KDE}(\zeta_{i,j}, \{\hat{\zeta}_l^s\})$ // estimated density at $\zeta_{i,j}$ in sample n

$\{\hat{q}_{i,j}^{(n)}\} \leftarrow \{\hat{q}_{i,j}^{(1)}, \hat{q}_{i,j}^{(2)}, \dots, \hat{q}_{i,j}^{(N)}\}$ // set of sampled densities at $\zeta_{i,j}$

$\hat{E}(\hat{q}_{i,j}) \leftarrow \int_{-\infty}^{\infty} \hat{q}_{i,j} \cdot \text{KDE}(\hat{q}_{i,j}, \{\hat{q}_{i,j}^{(n)}\}) d\hat{q}_{i,j}$ // expected value in duplicate-free data

if $q_{i,j} > \hat{E}(\hat{q}_{i,j})$: // test for anomaly
 $P(D|\zeta_{i,j}) \leftarrow \frac{\int_0^{q_{i,j}} \text{KDE}(q_{i,j} - q_{i,j}^D, \{\hat{q}_{i,j}^{(n)}\}) \cdot q_{i,j}^D dq_{i,j}^D}{q_{i,j}}$ // estimated duplicate probability
return $P(D|\zeta_{i,j})$

else:
return 0 // no anomaly, duplicate probability zero