

Supplementary Table 1 HDI-MF-Gower

HDI-MF-Gower

Input: Dataset D , numerical features N , categorical features F , threshold ε , max iterations P

Output: Imputed dataset D^{imp}

Phase 1: Initialization

1. Identify complete samples C and incomplete samples I

2. Calculate ranges: $R_k \leftarrow \max(D[:,k]) - \min(D[:,k])$ for $k \in N$ if $R_k = 0$ then $R_k \leftarrow 1$ end if

3. Calculate and normalize weights:

$w_k^{\text{raw}} \leftarrow \text{Var}(D[:,k]) \times (1 - MR_k)$ for $k \in N$ $w_k^{\text{raw}} \leftarrow H(D[:,k]) \times (1 - MR_k)$ for $k \in F$

$w_k \leftarrow \text{normalize}(\log(1 + w_k^{\text{raw}}))$ to $[0.1, 2.0]$ for $k \in (N \cup F)$

4. Initial imputation using adaptive Gower distance:

for each $i \in I$ do

$j^* \leftarrow \text{argmin}_{j \in C} (\sum_k w_k d_k) / (\sum_k w_k)$ where:

$d_k = |x_i[k] - y_j[k]| / R_k$ for $k \in N$ (if both observed)

$d_k = I(x_i[k] \neq y_j[k])$ for $k \in F$ (if both observed)

$D^{\text{imp}}[i,k] \leftarrow D[j^*,k]$ for missing values

end for

Phase 2: Iterative Refinement

1. Sort variables by missing rate in ascending order, $p \leftarrow 0$

2. while $p < P$ do

$D^{\text{imp_old}} \leftarrow D^{\text{imp}}$

for each variable k with missing values do

Train RF_k on observed data, predict missing values

end for

Calculate: Δ_N and Δ_F (see formulas (3) and (4))

if $\Delta_N < \varepsilon$ and $\Delta_F < \varepsilon$ then break

$p \leftarrow p + 1$

end while

3. return D^{imp}
