**Algorithm 1** Iterative Voting Consensus

**Input:** a set of $N$ data points $\mathbf{X} = \{x_1, x_2, ..., x_N\}$
a set of $C$ clusterings $\mathbf{\Pi} = \{\pi_1, \pi_2, ..., \pi_C\}$
$K$ is a desired number of clusters

**Output:** a consensus clustering $\pi^*$ with $K$ clusters

Initialize $\pi^*$
**repeat**
  Let $P_i = \{y \mid \pi^*(y) = i\}$ be the $i^{th}$ cluster

  Compute the representation of each cluster:
  $y_{P_i} = \langle majority\{(P_i)_1\}, ..., majority\{(P_i)_C\}\rangle$,
  where $(P_i)_j$ is the set of the $j^{th}$ features of all data
  points in $P_i$

  **for** $y$ in $\mathbf{Y}$ **do**
    Re-assign $\pi^*(y) \leftarrow argmin_i D(y, y_{P_i})$, where
    $D(y, y_{P_i}) = \sum_{j=1}^{C} \mathcal{I}((y)_j \neq (y_{P_i})_j)$
  **end for**
**until** $\pi^*$ does not change

---

**Algorithm 2** Iterative Probabilistic Voting Consensus

**Input:** a set of $N$ data points $\mathbf{X} = \{x_1, x_2, ..., x_N\}$
a set of $C$ clusterings $\mathbf{\Pi} = \{\pi_1, \pi_2, ..., \pi_C\}$
$K$ is a desired number of clusters

**Output:** a consensus clustering $\pi^*$ with $K$ clusters

Initialize $\pi^*$
**repeat**
  Let $P_i = \{y \mid \pi^*(y) = i\}$ be the $i^{th}$ cluster, $n_i = |P_i|$

  **for** $y$ in $\mathbf{Y}$ **do**
    Re-assign $\pi^*(y) \leftarrow argmin_i D(y, P_i)$, where

$$D(y, P_i) = \sum_{j=1}^{C} \frac{\sum_{y' \in P_i} \mathcal{I}((y)_j \neq (y')_j)}{n_i}$$

  **end for**
**until** $\pi^*$ does not change

---

for each data point via a defined distance measure between them. Formally, the distance between a data point $y$ and a cluster of $c$ data points $\{y_1, y_2, ..., y_c\}$, is defined as

$$D(y, \{y_1, y_2, ..., y_c\}) = \sum_{j=1}^{C} \frac{\sum_{i=1}^{c} \mathcal{I}((y)_j \neq (y_i)_j)}{c}. \quad (2)$$

The pseudo-code of the IPVC algorithm is described in algorithm 2. This algorithm can be viewed as a refinement of the IVC algorithm. Particularly, the distance function takes

into account the proportion that each feature of a point differs from those of points in a cluster of the target clustering.

### 5.3. Iterative Pairwise Consensus (IPC)

This iterative algorithm utilizes the similarity matrix $S$ as defined in equation 1. In each iteration, each data point is reassigned to different clusters based on the similarity measure between the data point and the previous established clusters in the target consensus clustering. Formally, the similarity measure between a data point $x$ and a cluster of $c$ data points $\{x_1, x_2, ..., x_c\}$, is defined as

$$S(x, \{x_1, x_2, ..., x_c\}) = \frac{\sum_{i=1}^{c} S(x, x_i)}{c}. \quad (3)$$

---

**Algorithm 3** Iterative Pairwise Consensus

**Input:** a set of $N$ data points $\mathbf{X} = \{x_1, x_2, ..., x_N\}$
a set of $C$ clusterings $\mathbf{\Pi} = \{\pi_1, \pi_2, ..., \pi_C\}$
$K$ is a desired number of clusters

**Output:** a consensus clustering $\pi^*$ with $K$ clusters

Compute the $n \times n$ similarity matrix $S$, where
$S(x_i, x_j) = \frac{1}{C} \sum_{c=1}^{C} \mathcal{I}(\pi_c(x_i) = \pi_c(x_j))$

Initialize $\pi^*$
**repeat**
  Let $P_i = \{x \mid \pi^*(x) = i\}$ be the $i^{th}$ cluster of the data
  points, and $n_i = |P_i|$

  **for** $x$ in $\mathbf{X}$ **do**
    Re-assign $\pi^*(x) \leftarrow argmax_i S(x, P_i)$, where
    $S(x, P_i) = \frac{\sum_{x' \in P_i} S(x, x')}{n_i}$
  **end for**
**until** $\pi^*$ does not change

---

Different from the previous two, this algorithm is a pairwise similarity approach. We can view this algorithm as applying a variation of k-means method to the similarity matrix $S$. The pseudo-code of the IPC algorithm is described in algorithm 3.

## 6. Evaluation Criteria

Evaluating the quality of a clustering is a nontrivial and ill-posed task [16]. In supervised learning, model performance is assessed by comparing model predictions to targets. In clustering we do not have targets and usually do not know *a priori* what groupings of the data are best. This hinders discerning when one clustering is better than another, or when one clustering algorithm outperforms another. In