# 1   Introduction

$k$-means is a clustering algorithm that partitions a set of data points into $k$ classes.

Suppose we have a set of data points $\{\mathbf{x}_i\}_{i=1}^n$ where each data point is a vector of generally continuous values $\mathbf{x}_i = [x_i, \ldots, x_p]$. Note that $k$-means works best with continuous data because it updates clusters with Euclidean distance.

An approximate approach is given by:

(i) Randomly assigning cluster centers $\hat{\mu}_1, \ldots, \hat{\mu}_M$. Generally, we these to be points in the training data.

(ii) Determine which of the clusters $R_1, \ldots, R_M$ each data point $\mathbf{x}_i$ belongs to by computing the closest cluster center $\hat{\mu}$.

(iii) Update the cluster centers $\hat{\mu}_m$ as the average of all $\mathbf{x}_i \in R_m$.

This process is iterated until we reach some max number of iterations or a convergence threshold where the change in cluster centers is deemed negligible.

# 2   Algorithm

---
**Algorithm 1** $k$-Means Clustering
---
**Input:** Data set $\mathcal{T} = \{\mathbf{x}_i\}_{i=1}^n$, max iterations $s$, stopping threshold $\epsilon$ (optional)
**Output:** Set of $k$ clusters

1: **function** $k$-MEANS($\mathcal{T}, s, \epsilon,$)
2:      $\hat{\mu}_i \leftarrow$ random $\mathbf{x}_p \in \mathcal{T}$ for $p = 1, \ldots, k$
3:      $t \leftarrow 1$
4:      **repeat**
5:          $R_j \leftarrow \emptyset$ for $j = 1, \ldots, k$
6:          **for** $\mathbf{x}_i \in \mathcal{T}$ **do**
7:              $j^* \leftarrow \arg\min_i \| x_j - \hat{\mu}_i \|^2$    (assign $\mathbf{x}_j$ to the closest cluster center)
8:              $R_{j^*} \leftarrow R_j \cup \mathbf{x}_j$
9:          **end for**
10:         **for** $i = 1, \ldots, k$ **do**
11:             $\hat{\mu}_i \leftarrow \frac{1}{|R_i|} \sum_{\mathbf{x}_j \in R_i} \mathbf{x}_j$    (update cluster centers to average of points in cluster)
12:         **end for**
13:      **until** $t = s$ or largest change in cluster center less than $\epsilon$
14: **end function**

---