**Algorithm 1** Expectation-Maximization Data-centric (EM-DC) over $\mathcal{D}$

1:   **INPUT** $\mathcal{D}$: data, $k$: cluster number.
2:   **OUTPUT** $\mathcal{C}_1, \ldots, \mathcal{C}_k$: $k$-normal distributions/clusters.
3:   // Each $\mathcal{C}_j^i \sim \mathcal{N}(\mu_j^i, \Sigma_j^i)$, $P(\mathcal{C}_j^i), w_{\mathbf{x}_j}^i, X_j^i \in \mathcal{C}_j^i$.
4:   // $i$: iteration number, $j$:cluster number.
5:   // $\mu_j^i$: mean, $\Sigma_j^i$: covariance, $P(\mathcal{C}_j^i)$: prior.
6:   // $w_{\mathbf{x}_j}^i \in \mathbb{R}$: likelihood and $X_j^i \subseteq \mathcal{D}$.
7:   // Binary Search Trees (BSTs): $\mathbf{T}^i = \{\mathcal{T}_1^i, \mathcal{T}_2^i, \ldots, \mathcal{T}_k^i\}, \cup_{t=1}^k \mathcal{T}_t = \mathcal{D}$
8:   randomly construct $\mathbf{C}^0 = \{\mathcal{C}_1^0, \mathcal{C}_2^0, \ldots, \mathcal{C}_k^0\}$
9:   $i \leftarrow 0$
10:  // $\mathcal{D}_{HE}$: high expressive data. All data points are high expressive at the first iteration.
11:  $\mathcal{D}_{HE} \leftarrow \mathcal{D}$
12: **repeat**
13:     **for** $\mathbf{x} \in \mathcal{D}$ **do**
14:         // E-step:
15:         **for** $\mathcal{C}_j^i \in \mathbf{C}^i$ **do**
16:             $\mathcal{C}_j^i.w_{\mathbf{x}_j}^i \leftarrow P(\mathcal{C}_j^i \mid \mathbf{x})$
17:         **end for**
18:         // If $i = 0$, build BSTs, otherwise, update them–depends on $w_{\mathbf{x}_j}^i$.
19:         $\mathcal{C}_j^i.\mathcal{T}_j^i.insert(\mathbf{x}, w_{\mathbf{x}_j}^i) \leftarrow (\mathbf{x}, w_{\mathbf{x}_j}^i)$
20:     **end for**
21:     // $\mathcal{D}_{HE}$: A temporary variable used to store the HE data at each iteration.
22:     $\mathcal{D}_{HE}' \leftarrow \varnothing$
23:     // M-step:
24:     **for** $\mathcal{C}_j^i \in \mathbf{C}^i$ **do**
25:         $\mathcal{C}_j^{i+1}.\mu_j^i \leftarrow \Sigma_{\mathbf{x} \in \mathcal{C}_j^i.X_j^i}(\mathbf{x} \cdot \mathcal{C}_j^i.w_{\mathbf{x}_j}^i/(\Sigma\mathcal{C}_j^i.w_{\mathbf{x}_j}^i))$
26:         $\mathcal{C}_j^{i+1}.\Sigma_j^i \leftarrow \Sigma_{\mathbf{x} \in \mathcal{C}_j^i.X_j^i}(\mathcal{C}_j^i.w_{\mathbf{x}_j}^i(\mathbf{x}-\mathcal{C}_j^i.\mu_j^i)(\mathbf{x}-\mathcal{C}_j^i.\mu_j^i)^T /(\Sigma_j^i\mathcal{C}_j^i.w_{\mathbf{x}_j}^i))$
27:         $\mathcal{C}_j^{i+1}.P(\mathcal{C}_j^i) \leftarrow \Sigma(\mathcal{C}_j^i.w_{\mathbf{x}_j}^i)/|\mathcal{C}_j^i.X_j^i|$
28:         // Using BSTs to determine new HE data and storing them in $\mathcal{D}_{HE}'$
29:         // $\mathcal{C}_j^i.\mathcal{T}_j^i.flush(\Sigma)$ represents separation of HE data from BST
30:         $\mathcal{D}_{HE}' \leftarrow \mathcal{C}_j^i.\mathcal{T}_j^i.flush(\Sigma)$
31:         $\mathbf{C}^{i+1} \leftarrow \cup\{\mathcal{C}_j^{i+1}\}$
32:     **end for**
33:     $i \leftarrow i + 1$
34:     // Updating final HE data
35:     $\mathcal{D}_{HE} \leftarrow \mathcal{D}_{HE}'$
36:     // $d : \mathbb{R} \times \mathbb{R} \to \mathbb{R}$, Stopping criterion (Convergence over structure, BSTs):
37:     // Node-wise hamming distance among BSTs between two consecutive iterations.
38: **until** threshold on $d(\mathbf{C}^{i-1}, \mathbf{C}^i) \leqslant \epsilon$