# 1 Online SSL

- Complete `online_ssl_update_centroids` using the pseudocode 1.

- Complete `online_ssl_compute_solution` following the pseudocode 2

---

**Algorithm 1** Incremental $k$-centers (simplified)

---

1: **Input:** an unlabeled $x_t$, a list of centroids $C_{t-1}$, a list of multiplicities $v_{t-1}$, taboo list $b$ containing the labeled centroids.
2: **if** $(|C_{t-1}| = k)$ **then**
3:     $c_1, c_2 \leftarrow$ two closest centroids such that at least one of them is not in $b$.
4:     // Decide which centroid is $c_{\text{rep}}$, that will represent both $c_1$ and $c_2$, and which centroid is $c_{\text{add}}$, that will represent the new point $x_t$.
5:     **if** $c_1$ in b **then**
6:         $c_{\text{rep}} \leftarrow c_1$
7:         $c_{\text{add}} \leftarrow c_2$
8:     **else if** $c_2$ in b **then**
9:         $c_{\text{rep}} \leftarrow c_2$
10:         $c_{\text{add}} \leftarrow c_1$
11:     **else if** $v_{t-1}(c_2) \leq v_{t-1}(c_1)$ **then**
12:         $c_{\text{rep}} \leftarrow c_1$
13:         $c_{\text{add}} \leftarrow c_2$
14:     **else**
15:         $c_{\text{rep}} \leftarrow c_2$
16:         $c_{\text{add}} \leftarrow c_1$
17:     **end if**
18:     $v_t \leftarrow v_{t-1}$
19:     $v_t(c_{\text{rep}}) \leftarrow v_t(c_{\text{rep}}) + v_t(c_{\text{add}})$
20:     $c_{\text{add}} \leftarrow x_t$
21:     $v_t(c_{\text{add}}) = 1$
22: **else**
23:     $C_t \leftarrow C_{t-1}.\text{append}(x_t)$
24:     $v_t \leftarrow v_{t-1}.\text{append}(1)$
25: **end if**

---

**Algorithm 2** Online HFS with Graph Quantization

---

1: **Input:** $t$, a list of centroids $C_t$, a list of multiplicities $v_t$ and labels $y$.
2: $V \leftarrow \text{diag}(v_t)$
3: $[\widetilde{W_q}]_{ij} \leftarrow$ weight between centroids $i$ and $j$.
4: Compute the Laplacian $L$ of the graph represented by $W_q = V\widetilde{W_q}V$
5: // Infer labels using hard-HFS.
6: $\widehat{y_t} \leftarrow \text{hardHFS}(L, y)$
7: // Remark: with the preceding construction of the centroids, $x_t$ is always present in the reduced graph and does not share the centroid with any other node.

---

Some practical considerations:

- The labeled nodes are fundamentally different from unlabeled ones. Because of this, it is always a good idea to keep them separate, and never merge them in a centroid. In the implementation this is accomplished with a taboo list $b$ that keeps track of nodes that cannot be merged together.

- In streaming applications, it is not always possible to stop execution to partition the centroids, and it is often preferable to pay a small price at every step to keep execution smooth. In our case, the centroids are updated at every step.

- Whenever a new node arrives, and we have too many centroids, we choose the two closest centroids $c_{\text{add}}$ and $c_{\text{rep}}$. $c_{\text{add}}$ will forget the old centroid and will point to the new sample that just arrived, and $c_{\text{rep}}$ will take care of representing all nodes that belonged to $c_{\text{add}}$.

# References

[1] Moses CHARIKAR, Chandra CHEKURI, Tomas FEDER, and Rajeev MOTWANI. Incremental clustering and dynamic information retrieval. *SIAM journal on computing*, 33(6):1417–1440, 2004.