# Hierarchical Dynamic Labeling (HDL)

## Semi-supervised learning



datasets → labeled / unlabeled

labeled → □ ← unlabeled
pesudo-label ↙

$$□ : \begin{cases} \text{predict network ( classifier )} \\ \quad\quad\quad or \\ \text{Represent network (feature extractor)} \end{cases}$$ ✓ Generalization ability ↑

## k-NN label clusterability ( unsupervised learning )



embedding space

Similar embedding → assemable
[metric learning ( clip . Dino )]

## Defined :

labeled dataset: $D = \{(x_n, y_n)\}_{n \in [N]}$ $[N] = \{1, 2, \cdots, N\}$ $y_n \in \{1, 2, \cdots, C\}$

unlabeled dataset: $D' = \{x'_m\}_{m \in [M]}$ $y'_m$ : potential label

$$\begin{cases} \forall m \in [M], y'_m \in [C] \\ \exists m \in [M], y'_m \notin [C] \end{cases} \xrightarrow{closed-set} \forall m \in [M], y'_m \in \{1, 2, \cdots, C\}$$

## Define 1 : label clusterability $(k, \delta_k)$

Dataset $D$ 、 image encoder $f(\cdot)$

→ embeddings $\mathcal{X} = f(D)$

$\forall x \in \mathcal{X}$, $P(x$ and $k$-nearest neighbors → Same Class$) \geq 1 - \delta_k$

↘ $D → (k, \delta_k)$   $\delta_k$: $D$ violate clusterability probability

(1) $\delta_k = 0$ , $D → k$-NN label clusterability

(2) $k\uparrow → \delta_k \uparrow$

(3) Representational ability ↓ → $\delta_k \uparrow$

(2)   $k\uparrow$ Risk↑ $\delta_k\uparrow$  (3) 

---

## Iterative algorithm process :

① Ideal label clustrability ( paper 4.1 (KNN-DV) )

Ideal : Strict clustering → lables

labeled embeddings $D = \{(x_n, y_n)\}_{n \in [N]}$

unlabeled embeddings $D' = \{x'_m\}_{m \in [M]}$

Since ideal So $D → k$-NN label clustrability

$x'_m, m \in [N] \longrightarrow y'_m, m \in [M]$

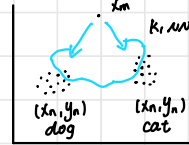$x'_m$ $k_1$ NN from $D = x'_{m_1}, x'_{m_2} \cdots x'_{mk_1}$

Corresponding labels: $y'_{m_1}, y'_{m_2} \cdots y'_{mk_1}$

label convert : integer → one-hot

Vote : $y'_m = \underset{i \in [C]}{\arg\max}\ (\frac{1}{k_1} \sum_{j}^{k_1} y'_{mj})[i],\ y'_{mj} \in R^c$

Traversal $\begin{cases} 1: i \in [C] \quad C: \text{number of labels} \\ 2: \sum_{j}^{k_1} y'_{mj} : k_1 \text{ nearest neighbors labels from } D \\ \quad\quad\quad\quad \text{one-hot form addition} \end{cases}$

### Geometric perspective :



| Classes | integer | one-hot |
|---------|---------|---------|
| dog | 0 | 10 |
| cat | 1 | 01 |

$([1,0] \times 2 + [0,1] \times 3) / 5$

$\Rightarrow \arg\max ([\frac{2}{5}, \frac{3}{5}])$

$\Rightarrow 1 \Rightarrow cat = y'_m$
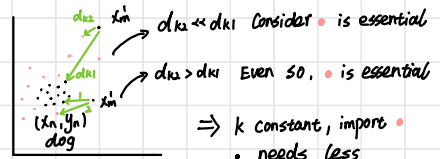
## Weakness :

1: Rely on highly clustering !

2: Step by step import unlabeled point

Ignore the relationship between them !

$y_n \in [C]$ , $y'_m \in [C]$   Same of label space

$\Rightarrow$ D mix D' follow clustrability in embedding space



$d_{k_2} \ll d_{k_1}$ Consider ● is essential

$d_{k_2} > d_{k_1}$ Even so, ● is essential

$\Rightarrow k$ constant, import ●

● needs less

● unlabeled point

We can proved clustrability be allowed ↓ Since ● needs less

→ Continous ...

---

| ResNet 50 | ResNet 50 (pre-trained) | clip ( vit-B31) |
|-----------|-------------------------|-----------------|
| dog k=7 | dog k=7 | dog k=7 |

we assume $x_m'$ authentic label is dog
$\Rightarrow$ how • ensure $x_m'$ labeling correctly when clustrability $\downarrow$
Simultaneous analysis of near and far points $x_{m1}'$, $x_{m2}'$



(k₁=5)

dog        Cat
$x_{m1}'$
$x_{m2}'$

• unlabeled point
( But no others)

(k₁=5)

dog        Cat
$x_{m1}'$
$x_{m2}'$

• unlabeled point
( have others help)

$\Longrightarrow$ I think • as fault-tolerant when clustrability $\downarrow$
② Unlabeled embedding help each others (paper 4.2 (KNN))
May $x_{m1}'$, $x_{m2}'$ ⋯ $x_{mk}'$ is unlabeled embedding
$\exists j \in [k] \rightarrow x_{mj}' \in D' \rightarrow$ how we know its label
$\Rightarrow$ Depth-First Search (DFS) for tree structure

$x_m'$ $y_m'$ (?)        $y_m' = \underset{i \in [c]}{\operatorname{argmax}} (\frac{1}{k} \overset{k}{\underset{j}{\sum}} y_{mj}')[i]$

$x_{m1}'$ $y_{m1}'$ ⋯ $x_{mi}'$ $y_{mi}'$ (?) ⋯ $x_{mk}'$ $y_{mk}'$

$y_{mi}' = \underset{j \in [c]}{\operatorname{argmax}} (\frac{1}{k} \overset{k}{\underset{a}{\sum}} y_{mia}')[j]$

$x_{mi1}'$ $y_{mi1}'$ ⋯ $x_{mia}'$ $y_{mia}'$ (?) ⋯ $x_{mik}'$ $y_{mik}'$

⎧ Exhaustively $\rightarrow$ find all labels
⎨        or
⎩ proportional threshold $\rightarrow$ partial labels compute manager

HDL : ↗ Dynamically increase determined labels

weakness :    ⎧ ↻ → loop search
fall into      ⎨
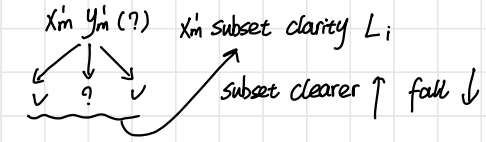               ⎩ ↻ multi branch repeated search

③ Hierarchical Dynamic labeling (paper 4.3 HDL)
   Try best to avoid fall : (core)
   (1) Prioritize labeling clearer ingredients ( P(fall) minimum )
   (2) Dynamically update D ( D'↓ D↑) ⇒ P(the rest, fall) ↓

Self-optimal ⎫
Overall-optimal ⎭

$D = \{(x_n, y_n)\}_{n \in [N]}$  $D' = \{x_m'\}_{m \in [M]}$
$x_m'$ $y_m'$ (?)    $x_m'$ subset clarity $L_i$
                    ↗ subset clearer ↑ fall ↓

$D'$, kNN, belong to $D$, is $L_1, L_2, \cdots L_m$
clearest : $L_m = \max\{L_1, L_2, \cdots L_m\}$, $m \in [M]$
$L_m$ ⎧ one ⇒ unlabeled embedding
     ⎨
     ⎩ multi   $x_{max,1}', \cdots, x_{max,s}'$, $s \leq M$
       (S)

(2) overall optimal

$x_{max,1}', \cdots, x_{max,s}'$  $s \leq M$ ( Self-optimal )
   arbitrary select ?
Assume labeling $x_{max,i}'$, $i \in [s] \rightarrow x_{max,i}'$ from $D'$ to $D$
Then the rest of self-optimal's $l$ may change
$[l_{i,1}, \cdots, l_{i,i-1}, l_{i,i+1}, \cdots l_{i,s}] \in R^{s-1}$
                                        ↘ expand

$$l = \begin{bmatrix} l_{1,2} \cdots l_{1,i} & l_{1,i+1} \cdots l_{1,s} \\ \vdots \ddots \vdots & \vdots \vdots \\ l_{i,1} \cdots l_{i,i-1} & l_i l_{i+1} \cdots l_{i,s} \\ \vdots \ddots \vdots & \vdots \ddots \vdots \\ l_{s,1} \cdots l_{s,i-1} & l_{s,i} \cdots l_{s,i-1} \end{bmatrix} \in R^{s \times (s-1)}$$

$l[i]$ : labeling the i-th, the rest of $D' \rightarrow l$
$\underset{i \in [s]}{\operatorname{argmax}} \operatorname{sum}(l[i]) \rightarrow$ Rest-optimal

self-optimal + Rest-optimal = overall-optimal

④ Adaptive k (paper 4.4)
balance ⎧ → kNN : k↑ clustrability ↓
        ⎩ → majority voting : k↑ robust↑ Reliable↑
Ref [43] (2022 ICML) :
" Detecting corrupted labels without training a model to predict "
                                    → Continous ...

2022 ICML: (lower bound)

$P(\text{Vote is correct} \mid k) \geq (1-\delta_k) \cdot I_{1-e}(k+1-k', k'+1)$

$(1-\delta_k)$: Quality of features

Based on it

$\Rightarrow P(k) \geq \mu_k \cdot I_{1-e}(k+1-k', k'+1)$

$\begin{cases} \mu_k: P(\text{clustrability}) \\ \quad e: P(\text{error labels in embedding set}) \\ \quad k': [(k+1)/2] - 1 \\ \quad I_{1-e}(k+1-k', k'+1) \cdot \text{regularized incomplete beta function} \end{cases}$

$\mu_k = (1-\delta_k) = $ Quality of embedding

$I_{1-e}(k+1-k', k'+1) = \dfrac{(k+1)!}{(k-k')! \, k'!} \displaystyle\int_0^{1-e} t^{k-k'}(1-t)^{k'} dt$

$k\uparrow \Rightarrow \begin{cases} I_{1-e}(k+1-k', k'+1) \uparrow \\ \\ \text{clustrability} \downarrow \Rightarrow \mu_k \downarrow \end{cases}$

Qutification $\mu_k$:

From labeled embedding space
randomly select instances to statistics

$k = 5$



estimate $\mu_k$

dog        cat

(backbone: clip)