

«ELFS: LABEL-FREE Coreset selection with proxy Training Dynamics»

١ ELFS utilizes deep clustering to estimate training dynamics based data difficulty scores without ground truth labels.

٢ Pseudo-labels introduce a distribution shift in the data difficulty scores, and we propose a simple but effective double-end pruning method to mitigate bias on calculated scores.

Outcomes of training dynamics; Classification

accuracy vs. number of iterations (Proxy-based)

and also

* For the second challenge, we observe that directly applying existing selection methods like CCS results in a performance gap when compared to coresets selection with ground truth labels.

رسوں کے تاثر میں ہے جو دادھنکی و تجزیہ حسبی (یا درستہ تجزیہ) کر پائے جاوے
مختصر صورتی انتخاب کرنے کے لائق دادھنکی و تجزیہ کرنے کے لئے

* This distribution shift leads traditional sampling methods to select more easy examples, resulting in inferior performance.

→ To mitigate this distribution shift issue, we propose a simple yet effective double-end pruning method for the pseudo-label-based difficulty scores.

* This double-end pruning significantly reduces the number of selected easy examples and improves the coresset selection performance.

بيان
پنجشنبہ

۲۶

26 Oct 2023

ربيع الثاني ١٤٤٥

* CCS introduces a method that combines hard example pruning with stratified sampling to jointly consider data difficulty and coverage, significantly improving coreset selection performance.

نیو ٹکنالوجیز

وفات حضرت مصطفیٰ (ص) (۲۰۱ هـ)



بيان
جمعہ

27 Oct ■

ربيع الثاني ۱۱

* Deep clustering is a technique that combines deep learning methods to group data into clusters without explicit supervision.

* we utilize Teacher Ensemble-weighted pointwise mutual Information (TE MI) to generate pseudo-labels for the unlabeled dataset.

- * calculates the K nearest neighbor set N_x .
- * After getting the nearest neighbor sets for all examples,
- * TEMI trains an ensemble of classification heads to assign pseudo labels.
- * TEMI uses an ensemble of student head $h_s(\cdot)$ and a teacher head $h_t(\cdot)$ that share the same architecture but differ in their updates.
- * Each instance x is processed through both types of heads producing probabilistic classifications $q_s(c|x)$ and $q_t(c|x)$ respectively, where $c \in \{1, \dots, C\}$ denotes the class label.

* The student head's parameters θ_s are updated via backpropagation, while the teacher head's parameters θ_t are updated via an exponential moving average (EMA) of the student's parameters, ensuring more stable target distributions.

Loss Function

* Given H clustering heads, TEMI aims to minimize the following accumulated weighted pointwise mutual information (PMI) objective to encourage the model to align instances that are likely to belong to the same cluster while penalizing those that are less likely to be related:

$$* L_{TEMI}(x) = -\frac{1}{2H} \sum_{h=1}^H \sum_{x' \in N_x} w_h(x, x') \cdot (pmi^h(x, x') + pmi^h(x', x))$$

$$* pmi(x, x') = \log \left(\sum_{c=1}^C \frac{(q_s(c|x) q_t(c|x'))^\beta}{q_t(c)} \right)$$

\hookrightarrow pmi pointwise mutual information \leftarrow pmi درجات حرارة

١٢ | بيان
3 Nov | ١٨ | ربيع الثاني
 $w(x, x') = \sum_{c=1}^C q_s(c|x) q_t(c|x')$

\hookrightarrow درجات حرارة بين حبيبات كثافة ملائمة في نفس cluster

* TEMI assigns the pseudo label \tilde{y} by aggregating these predictions and selecting the class with the highest combined probability:

$$\tilde{y} = \arg \max_{c \in \{1, \dots, C\}} \frac{1}{H} \sum_{h=1}^H q_h(c|x)$$

* Pseudo-labels generated by deep clustering make it feasible to calculate training dynamics scores with supervised training.

* However, the clustering-based pseudo-labeling can introduce many label noises, which causes the distribution shift of difficulty scores. ?

برای این کار داده است که ابتدا داده های سخت را حذف و در نهاد می شوند و سپس داده های آسان را حذف می شوند

۱۸ | آبان | ۷ Nov 2023
۲۲ | ۱۴۴۵ | دین الشانی

- * Double-end pruning consists of two steps: 1 Prune β hard examples first, where β is a hyperparameter.
- 2 Continue pruning easy examples until the budget is met.
- * A common practice is to perform a grid search to select the optimal β .
- * We show that pseudo-labels generated in the deep clustering step provide a good estimation for this β search: After generating the pseudo-labeled dataset, we split it into 90% for training and 10% for validation.
- * We use the validation set to determine the optimal β . Once the optimal β is identified, ECLFS selects coresets from the

١٨ | آبان پنجشنبه ۹ Nov 2023
۱۴۴۵، ربیع الثانی ۲۴

- * entire training dataset. Our evaluation shows that β grid search with pseudo-labels provides a good estimation for the optimal β .
- * ELFS guarantees that ELFS does not use any ground truth labels in the coresnet selection process.
- * We report ELFS results using DINO and SwAV as the pretrained vision model to extract embeddings.

۱۹ | آبان جمعه ۱۰ Nov ۱۴۴۵، ربیع الثانی ۲۵

- * we use the area under the margin (AUM) as the training dynamics metric for all datasets.

روز جهانی علم در خدمت صلح و توسعه