

[Active Learning by Feature Mixing]

<Problem Definition>

- labeled data $D^l = \{(x_i, y_i)\}_{i=0}^M$
- unlabeled data $D^u \xrightarrow{\text{extract}} B$ instances for query

Learner : $f = f_c \odot f_e$ parameterised by $\theta = \{\theta_e, \theta_c\}$

$$f_e : X \rightarrow \mathbb{R}^D$$

↑ backbone which encodes the input to a D-dimensional representation in a latent space

$$Z = f_e(x; \theta_e) \leftarrow \text{encoded input to the latent space}$$

$$f_c : \mathbb{R}^D \rightarrow \mathbb{R}^K$$

↑ classifier (ex. multi-layer perceptron (MLP)) that maps the instances from their representations to their corresponding logits which can be converted to class likelihoods by

$$p(y|z; \theta) = \text{Softmax}(f_c(z; \theta)) \leftarrow \text{class likelihoods of instances mapped from their representations.}$$

✓ Optimise parameters end-to-end by minimizing cross-entropy loss over the labelled set:

$$\mathbb{E}_{(x,y) \sim D^l} [\ell(f_c \odot f_e(x; \theta), y)]$$

ℓ : categorical cross entropy

✓ The prediction of the label (i.e. pseudo-label) for an unseen instance:

$$y_z^* = \underset{y}{\operatorname{argmax}} \ f_c^y(z; \theta_c)$$
$$\begin{cases} z = f_e(x; \theta_e) \\ f_c^y : \text{logit output for class } y \end{cases}$$

- ✓ The logit of the predicted label: $f_c^*(z) := f_c^{y^*_z}(z)$
- ✓ The set of representations of the unlabelled data: $Z^u = \{f_e(x), \forall x \in D^u\}$
- ✓ The set of representations of the labelled data: $Z^l = \{f_e(x), \forall x \in D^l\}$
- z^* : the average representation of labelled samples per class
(= anchor)
- Z^* : The set of anchors for all classes.
(= representatives of the labelled instances)

<Feature Mixing>

Intuition: The model's incorrect prediction is mainly due to novel "features" in the input that are not recognisable.

→ We use interpolation as a way to explore novel features in the vicinity of each unlabeled point.

Interpolation between the representations of the unlabelled and labelled instances,
 ~~\underline{z}_x~~ \underline{z}^u and \underline{z}^* , respectively:

$$\underline{\tilde{z}}_\alpha = \alpha \underline{z}^* + (1 - \alpha) \underline{z}^u \quad \alpha \in [0, 1]^D : \text{interpolation ratio}$$

$$\underline{z} \sim p(\underline{z} | \underline{z}^u, \underline{z}^*, \alpha) \equiv \alpha \underline{z}^* + (1 - \alpha) \underline{z}^u, \quad \underline{z}^* \sim \underline{Z}^* \quad \leftarrow \begin{array}{l} \text{can be seen as sampling a new instance} \\ \text{without explicitly modelling the joint} \\ \text{probability of the labelled and unlabelled} \\ \text{instances.} \end{array}$$

Consider how model's prediction changes:

- { ① The change in pseudo-label (i.e. y^*) for the unlabelled instance
- ② The loss incurred with the interpolation

The model's loss for predicting the pseudo-label of an unlabelled instance at its interpolation with a labeled one:

$$l(f_c(\tilde{z}_\alpha), y^*) \approx l(f_c(z^u), y^*) + (\alpha(z^* - z^u))^\top \nabla_{z^u} l(f_c(z^u), y^*)$$

↑ for sufficiently small $\|\alpha\| \leq \epsilon$ is almost exact

Consequently, for the full labelled set,

by choosing the max loss from both sides we have:

$$\begin{aligned} & \max_{z^* \sim Z^*} [l(f_c(\tilde{z}_\alpha), y^*)] - l(f_c(z^u), y^*) \\ & \approx \max_{z^* \sim Z^*} [(\alpha(z^* - z^u))^\top \nabla_{z^u} l(f_c(z^u), y^*)] \end{aligned}$$

When performing interpolation, the change in the loss is proportionate to two terms:

① The difference of features of z^* and z^u proportionate to their interpolation α .

↳ determines which features are novel and how their value could be different between the labelled and unlabelled instance.

② The gradient of the loss with respect to the unlabelled instance.

↳ determines the sensitivity of the model to those features.

∴ If the features of the labelled and unlabelled instances are completely different but the model is reasonably consistent, there is ultimately no change in the loss, and hence those features are not considered novel to the model.

α : input specific

· determines the features to be selected

< Optimising the Interpolation Parameter α >

The objective for choosing α

$$\alpha^* = \underset{\|\alpha\| \leq \epsilon}{\operatorname{argmax}} (\alpha(z^* - z^u))^T \nabla_{z^u} l(f_c(z^u), y^*)$$

ϵ : hyperparameter governing the magnitude of the mixing

chooses the hardest case of α for each unlabeled instance and anchor

Approximate solution to this optimization from using dual norm formulation :

$$\alpha^* \approx \epsilon \frac{\|(z^* - z^u)\|_2 \nabla_{z^u} l(f_c(z^u), y^*)}{\|\nabla_{z^u} l(f_c(z^u), y^*)\|_2} \odot (z^* - z^u)$$

\odot : element-wise division

< Candidate Selection >

It is reasonable to choose instances to be queried whose loss substantially change with interpolation

↳ instances for which the model's prediction change and has novel features.
= placed close to the decision boundary in the latent space.

Candidate set: $I = \{z^u \in Z^u \mid \exists z^* \in Z^*, f_c^*(z^*) \neq y_{z^u}^*\}$

- The size of selected set I could potentially be larger than the budget B .
- Seek diverse samples since most instances in I could be chosen from the same region. (i.e. they might share the same novel features.)

K-Means → Cluster the instances in I into B groups based on their feature similarities and further choose the closest samples to the centre of each cluster to be labelled by oracle.

<Algorithm>

D^l : Initial labelled set

D^u : Unlabelled pool

B : Labelling budget at each round

ϵ : mixing parameter

for $i=1$ to max-rounds do :

Train the model f using the labelled data D^l

Initialize Z^* on the representations of D^l

$I = \{\}$

for $x^u \in D^u$ do

$$z^u = f_e(x^u)$$

for $z^* \in Z^*$ do $\# Z^*$: The set of anchors for all classes
(representatives of labelled instances)

Calculate α^* using ϵ

$$\tilde{z} = \alpha^* z^* + (1 - \alpha^*) z^u \quad \# \text{ interpolation}$$

$$\text{if } \underset{y}{\operatorname{argmax}}(f_c^y(z_u)) \neq \underset{y}{\operatorname{argmax}}(f_c^y(\tilde{z}))$$

then

$$I = I \cup (x^u, z^u)$$

Break

Cluster the samples in I into B clusters.

Select samples at the centre of each cluster (C)

$$Y^{\text{new}} = \text{Query}(C)$$

$$D^l = D^l \cup (C, Y^{\text{new}}), \quad D^u = D^u \setminus C$$

Z_{unlabeled}
from set of
all samples
in unlabeled
instances
Interpolation?
Total?