



Neighborhood-Regularized Self-Training for Learning with Few Labels

method

task

Advisor : Jia-Ling, Koh

Speaker : Ting-I, Weng

Source : AACL'23

Date : 2023/12/19



Outline

- Introduction
- Method
- Experiment
- Conclusion

Task



text



AGnews



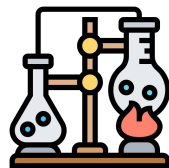
world



business



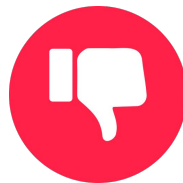
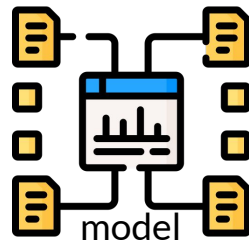
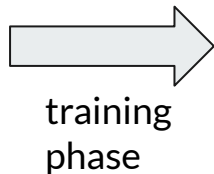
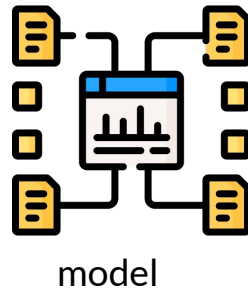
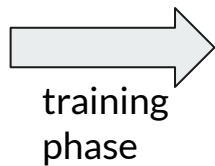
sports

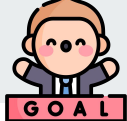


Sci/Tech

classification

Problem - few label data





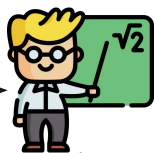
Hope to use unlabeled data to assist training

Self-training

few label data



init model



teacher model



unlabeled data

strategy



pseudo labeled data

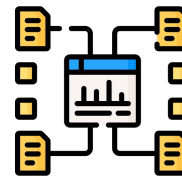


training dataset

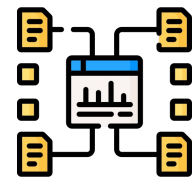


student model

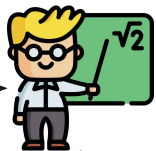
final output



Challenge of self-training



few label data



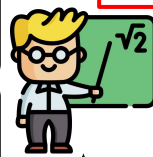
teacher model

init model



unlabeled data

strategy



over-confident



training dataset



biased



training instability

student model



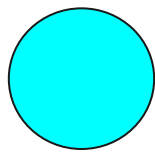
final output



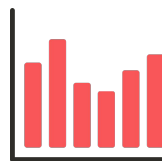
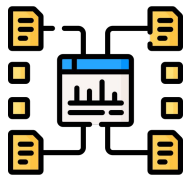
samples with similar labels tend to share similar representations

Hypothesize

labeled data

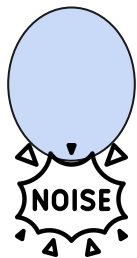


label : 1

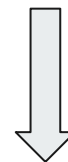
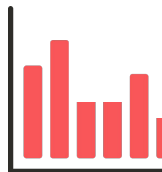
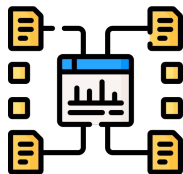


very similar

pseudo labeled data with noise



label : 1



pseudo labeled may be correct

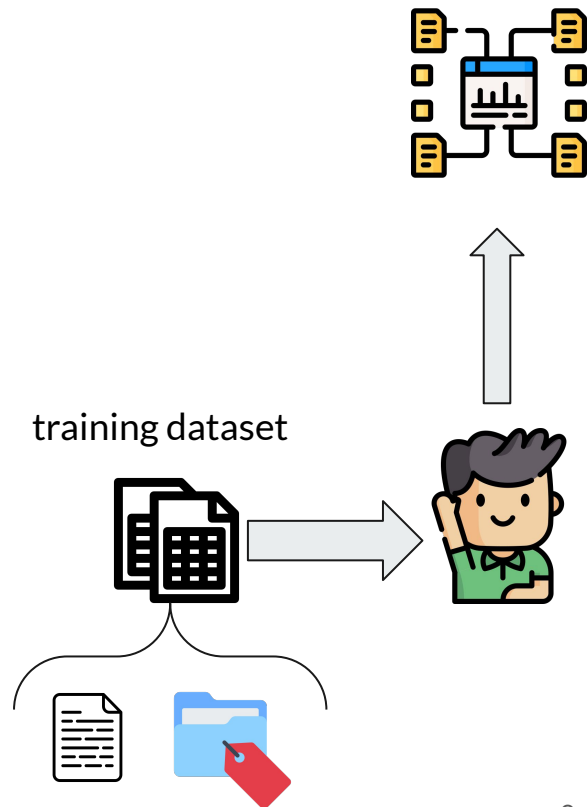
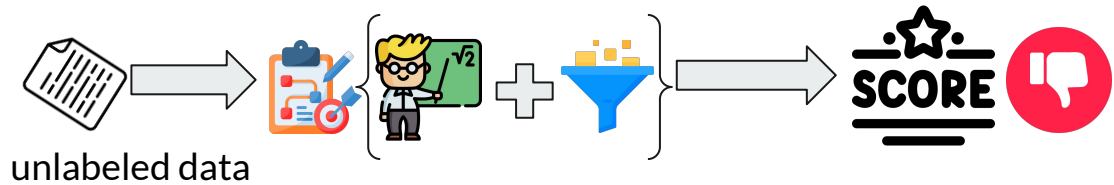
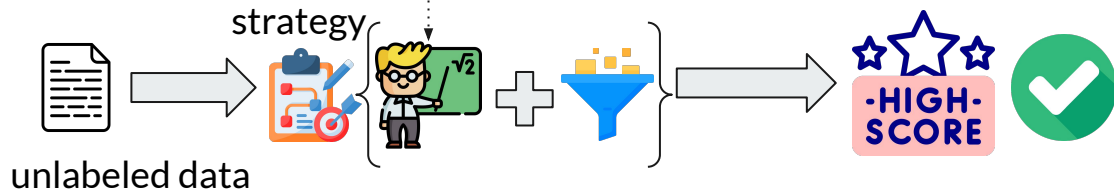
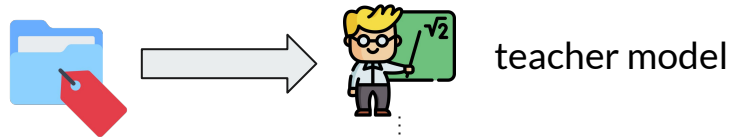




- over-confident
- biased
- training instability

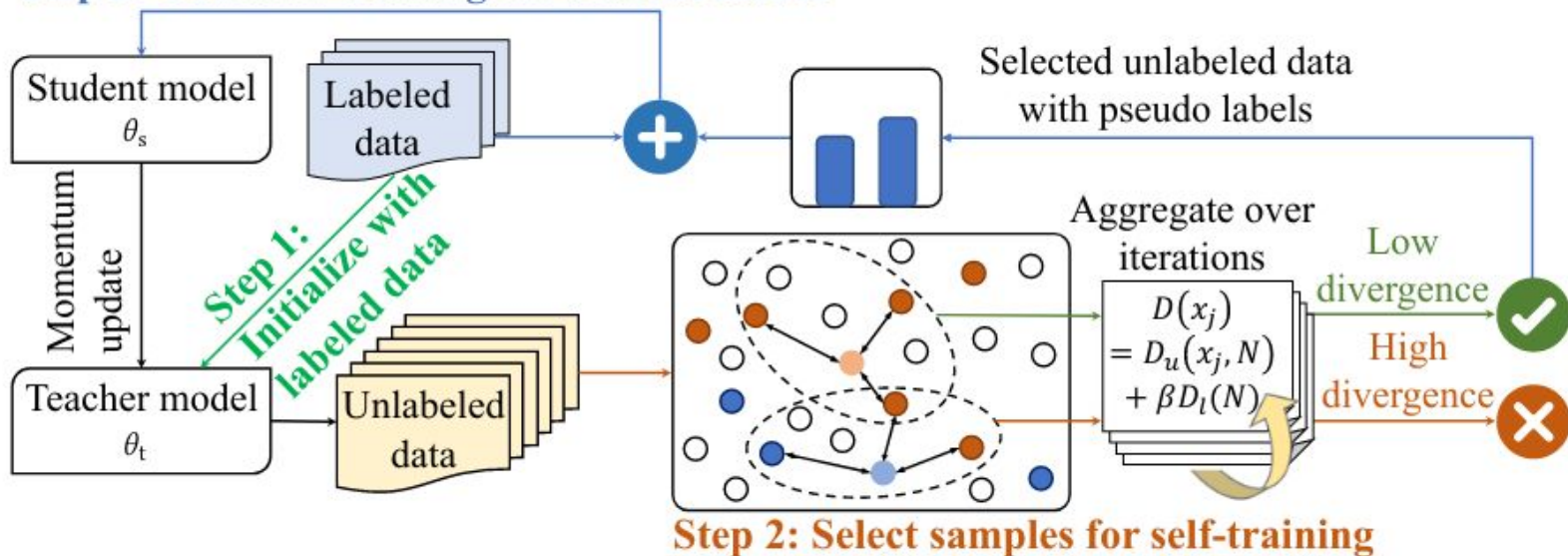
Solution

few label data



NeST(Neighborhoo Regularized Self-Training)

Step 3: Continue training the student model





Outline

- Introduction
- **Method**
- Experiment
- Conclusion

class:

● black : 0

● green : 1

● X_l : labeled data

○ black : 5

○ green : 5

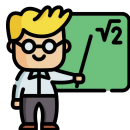
samples with similar labels tend to share similar representations

Neighborhood-Regularized Sample Selection

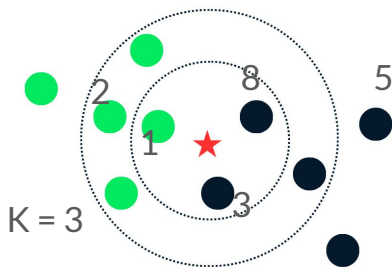
$$\mathcal{N}_j = \{x_i \mid x_i \in \mathcal{X}_l \cap \text{KNN}(v_j, \mathcal{X}_l, k)\}$$



unlabeled
data



embedding



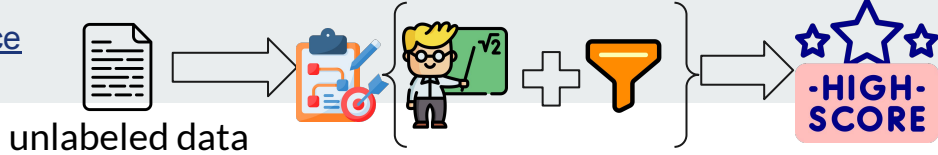
Number close to unlabeled data

$$N = \{X_1, X_3, X_8\}$$

- black : 2
- green : 1



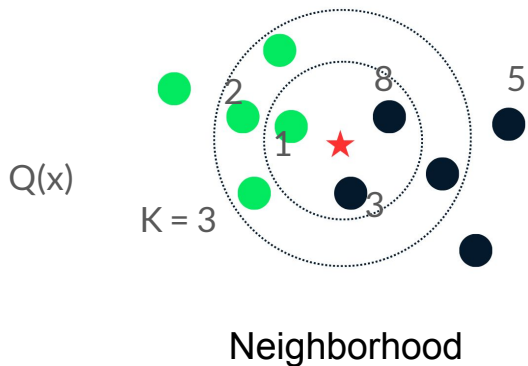
Neighborhood



Divergence-based Sample Selection

$P(x)$  unlabeled data

$$D_{\text{KL}}(Q \parallel P) = \sum_{x \in \mathcal{X}} Q(x) \ln \left(\frac{Q(x)}{P(x)} \right)$$



$$\mathcal{D}(x_j) = \mathcal{D}_u(x_j, \mathcal{N}) + \beta \mathcal{D}_1(\mathcal{N})$$

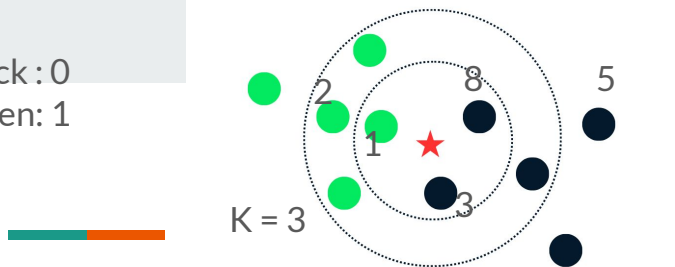


unlabeled data

$$\mathcal{N} = \{X_1, X_3, X_8\}$$

class:


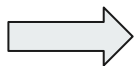
black : 0
green : 1



Unlabeled Divergence D_u

$$D_u(x_j, \mathcal{N}) = \sum_{(x_i, y_i) \in \mathcal{N}} d(f(x_j; \theta_t), y_i)$$

$$\begin{aligned} D_u(x_1^u, N_1) &= d(f(x_1^u; \theta_t), y_1) \\ &\quad + d(f(x_1^u; \theta_t), y_2) \\ &\quad + d(f(x_1^u; \theta_t), y_3) \end{aligned}$$



 $f(x_1^u; \theta_t) = [0.2, 0.8]$

unlabeled data

$$y_i \begin{cases} y_1 = [0, 1] \quad \bullet \\ y_2 = [1, 0] \quad \bullet \\ y_3 = [1, 0] \quad \bullet \end{cases}$$

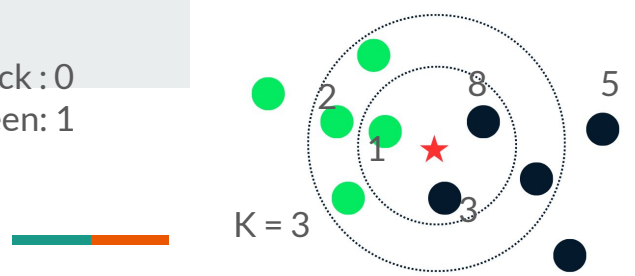
$$\begin{aligned} &= 0 \log_e \frac{0}{0.2} + 1 \log_e \frac{1}{0.8} \\ &\quad + 1 \log_e \frac{1}{0.2} + 0 \log_e \frac{0}{0.8} \\ &\quad + 1 \log_e \frac{1}{0.2} + 0 \log_e \frac{0}{0.8} \end{aligned}$$

$$= 0.2231 + 1.6094 + 1.6094 = 3.442$$

$$N = \{X_1, X_3, X_8\}$$

class:
black : 0
green : 1

- black : 2
- green : 1



Labeled Divergence DL

$$\mathcal{D}_1(\mathcal{N}) = \sum_{(x_i, y_i) \in \mathcal{N}} d(\bar{y}, y_i)$$

$$\bar{y} = \sum \frac{y_i}{|N|} = \frac{[0, 1] + [1, 0] + [1, 0]}{3} = \left[\frac{2}{3}, \frac{1}{3}\right]$$

$$y_i \begin{cases} y_1 = [0, 1] \text{ ●} \\ y_2 = [1, 0] \text{ ●} \\ y_3 = [1, 0] \text{ ●} \end{cases}$$

$$\begin{aligned} D_l(N_1) &= \sum d(\bar{y}, y_i) \\ &= d\left(\left[\frac{2}{3}, \frac{1}{3}\right], [0, 1]\right) \\ &\quad + d\left(\left[\frac{2}{3}, \frac{1}{3}\right], [1, 0]\right) \\ &\quad + d\left(\left[\frac{2}{3}, \frac{1}{3}\right], [1, 0]\right) \\ &= 1.9125 \end{aligned}$$



unlabeled data



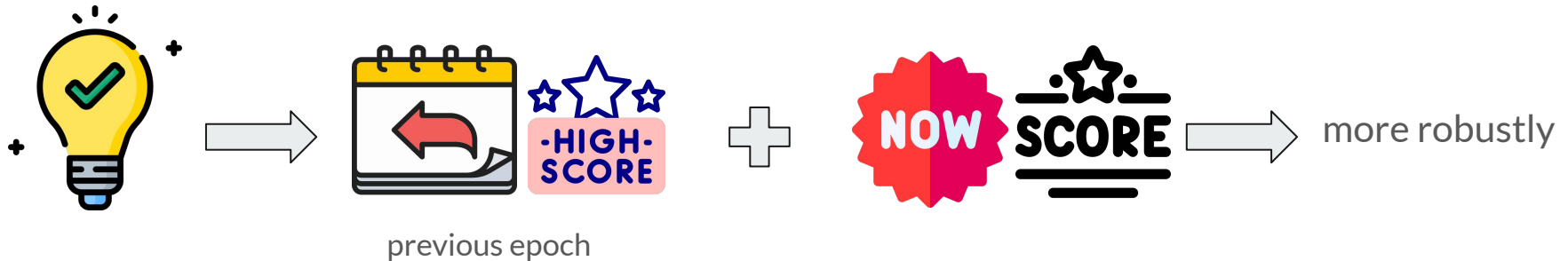
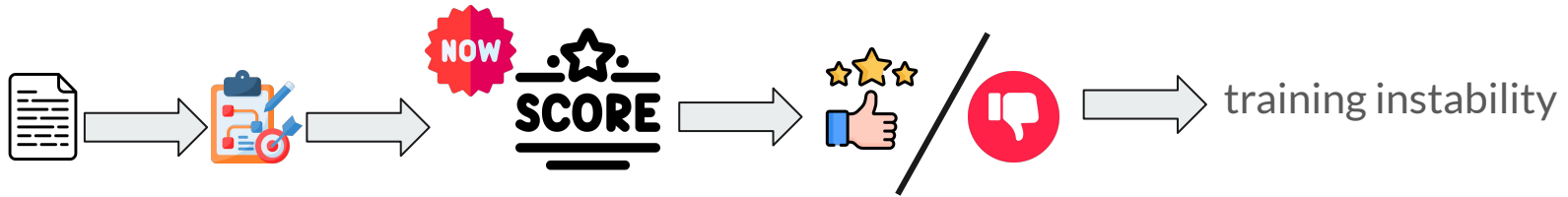
Divergence-based Sample Selection

$$D(x_j) = \mathcal{D}_u(x_j, \mathcal{N}) + \beta \mathcal{D}_1(\mathcal{N})$$

$$\begin{aligned} D(x_1^u) &= 3.442 + \beta * 1.9125 \\ &= 3.442 + 0.1 * 1.9125 = 3.6332 \end{aligned}$$




Aggregation of Predictions from Different Iterations



Aggregation of Predictions from Different Iterations

$$\mu^{(t)}(x_j) = (1 - m) \times \mu^{(t-1)}(x_j) + m \times \mathcal{D}^{(t)}(x_j)$$

 $\mu^{(t)}(x_j)$ is associated with **unlabeled data** (document icon) and **previous** (calendar icon with red arrow) $\mu^{(t-1)}(x_j)$ is associated with **HIGH SCORE** (stars icon) $\mathcal{D}^{(t)}(x_j)$ is associated with **NOW SCORE** (red starburst icon)

$$\mu^{(t)}(x_j) = (1 - m) \times \mu^{(t-1)}(x_j) + m \times (\mathcal{D}^{(t)}(x_j))$$

more robustly



$m = 0.6$



$t = 1$

Aggregation of Predictions from Different Iterations

$$\mu^1(x_1^u) = (1 - m) * \mu^{1-1}(x_1^u) + m * (D^1(x_1^u)) = 0.6 * 3.442 = 2.0652$$

$$\mu^1(x_2^u) = (1 - m) * \mu^{1-1}(x_2^u) + m * (D^1(x_2^u)) = 0.6 * 1.9125 = 1.1475$$

$$\mu^1(x_3^u) = (1 - m) * \mu^{1-1}(x_3^u) + m * (D^1(x_3^u)) = 0.6 * 0.9 = 0.54$$

label	$D^{t=1}(x_j)$	$D^{t=2}(x_j)$	$D^{t=3}(x_j)$	
 x_1^u	3.442	2.5	1.5	
x_2^u	1.9125	1.6	2.7	
x_3^u	0.9	0.7	0.5	

t = 2





	μ^1
x_1^u	2.0652
x_2^u	1.1475
x_3^u	0.54

m = 0.6

$$\mu^2(x_1^u) = (1 - m) * \mu^{2-1}(x_1^u) + m * (D^2(x_1^u)) = 0.4 * 2.0652 + 0.6 * 2.5 = 2.326$$

$$\mu^2(x_2^u) = (1 - m) * \mu^{2-1}(x_2^u) + m * (D^2(x_2^u)) = 0.4 * 1.1475 + 0.6 * 1.6 = 1.149$$

$$\mu^2(x_3^u) = (1 - m) * \mu^{2-1}(x_3^u) + m * (D^2(x_3^u)) = 0.4 * 0.54 + 0.6 * 0.6243 = 0.5905$$

	label	$D^{t=1}(x_j)$	$D^{t=2}(x_j)$	$D^{t=3}(x_j)$	 SCORE
	x_1^u	3.442	2.5	1.5	
	x_2^u	1.9125	1.6	2.7	
	x_3^u	0.9	0.7	0.5	

Robust Aggregation of Predictions from Different Iterations

x_1^u

model gives **inconsistent predictions** in different iterations



→ hurt model

x_3^u

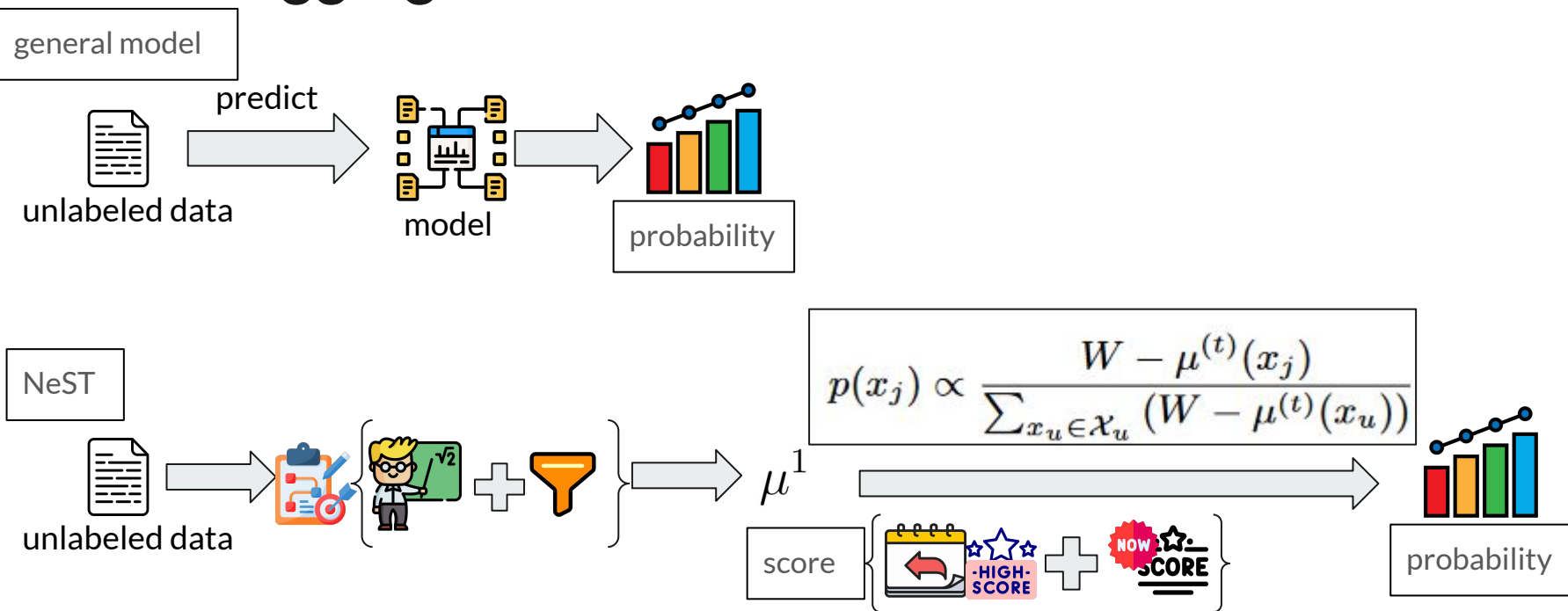
model output **consistently low scores** in different iterations



→ model more certain

	μ^1	μ^2	μ^3 suppose	μ^4 suppose	
x_1^u	2.0652	2.326	0.9	2.2	
x_2^u	1.1475	1.149	1.2	1.05	
x_3^u	0.54	0.5905	0.4	0.3	

Robust Aggregation of Predictions from Different Iterations



t = 1

Robust Aggregation of Predictions from Different Iterations

$$p(x_j) \propto \frac{W - \mu^{(t)}(x_j)}{\sum_{x_u \in \mathcal{X}_u} (W - \mu^{(t)}(x_u))}$$

$W = \max_{x \in \mathcal{X}_u} (\mu^{(t)}(x))$ is the normalizing factor

$$W = \max_x (\mu^t(x)) = 2.0652$$

$$p(x_1) = \frac{2.0652 - 2.0652}{0 + [2.0652 - 1.1475] + [2.0652 - 0.54]} = 0$$

$$p(x_2) = \frac{2.0652 - 1.1475}{0 + [2.0652 - 1.1475] + [2.0652 - 0.54]} = 0.3756$$

$$\star \star \star \text{👍} p(x_3) = \frac{2.0652 - 0.54}{0 + [2.0652 - 1.1475] + [2.0652 - 0.54]} = 0.6243$$

	μ^1	μ^2
x_1^u	2.0652	2.326
x_2^u	1.1475	1.149
x_3^u	0.54	0.5905

t = 2

Robust Aggregation of Predictions from Different Iterations

$$p(x_j) \propto \frac{W - \mu^{(t)}(x_j)}{\sum_{x_u \in \mathcal{X}_u} (W - \mu^{(t)}(x_u))}$$

$W = \max_{x \in \mathcal{X}_u} (\mu^{(t)}(x))$ is the normalizing factor

$$W = \max_x (\mu^t(x)) = 2.236$$

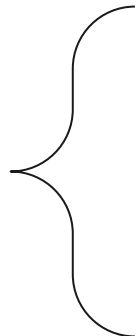
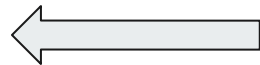
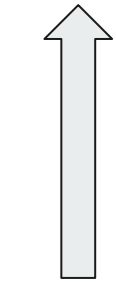
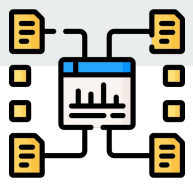
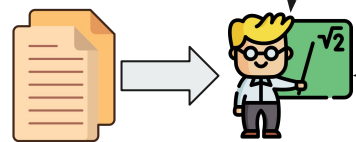
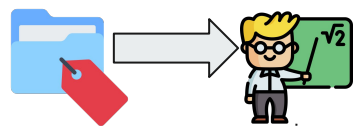
$$p(x_1) = \frac{2.326 - 2.326}{[2.326 - 1.149] + [2.326 - 0.5905]} = 0$$

$$p(x_2) = \frac{2.326 - 1.149}{[2.326 - 1.149] + [2.326 - 0.5905]} = 0.4041$$

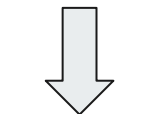
$$\star \star \star \text{👍} p(x_3) = \frac{2.326 - 0.5905}{[2.326 - 1.149] + [2.326 - 0.5905]} = 0.5958$$

	μ^1	μ^2
x_1^u	2.0652	2.326
x_2^u	1.1475	1.149
x_3^u	0.54	0.5905

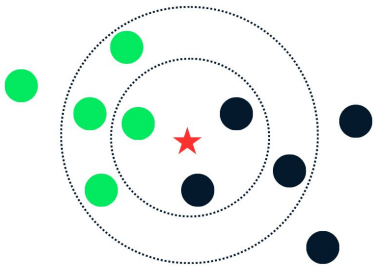
Model process



unlabeled data



embedding



$$N = \{X_1, X_3, X_8\}$$





Outline

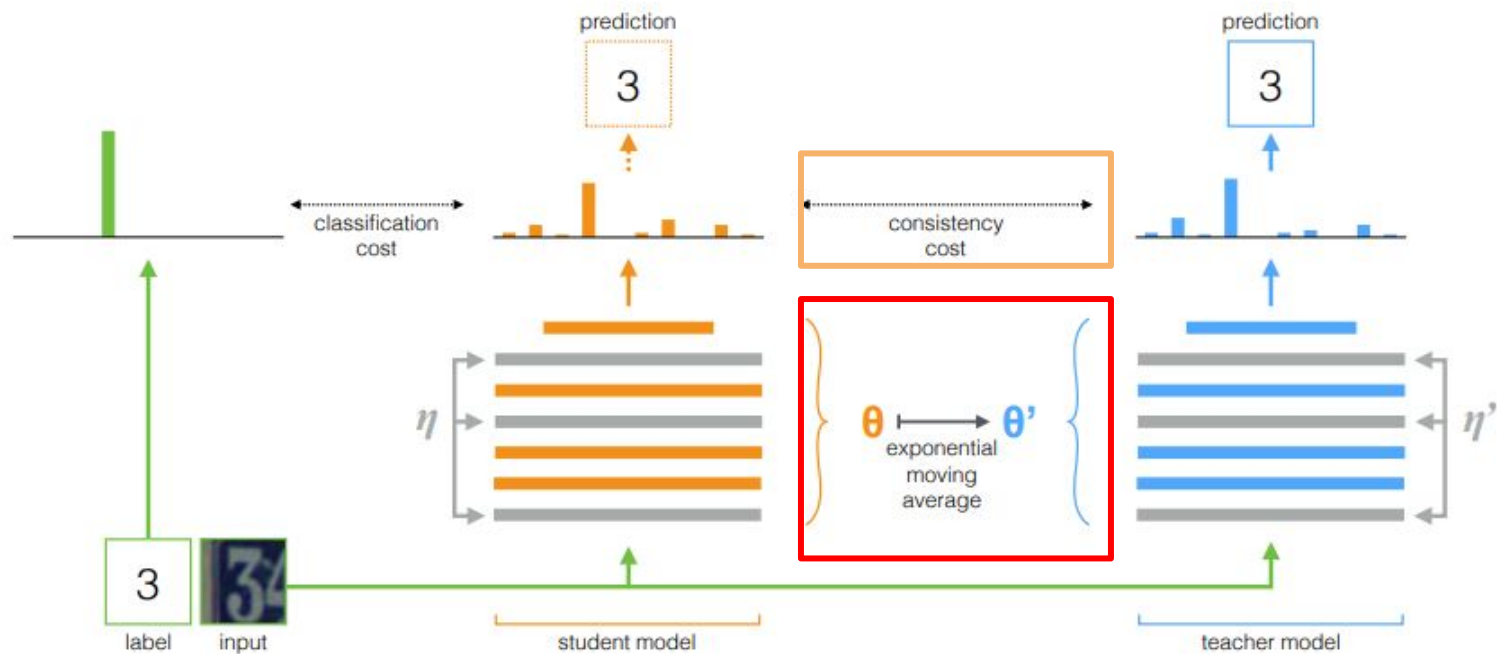
- Introduction
- Method
- **Experiment**
- Conclusion

Dataset

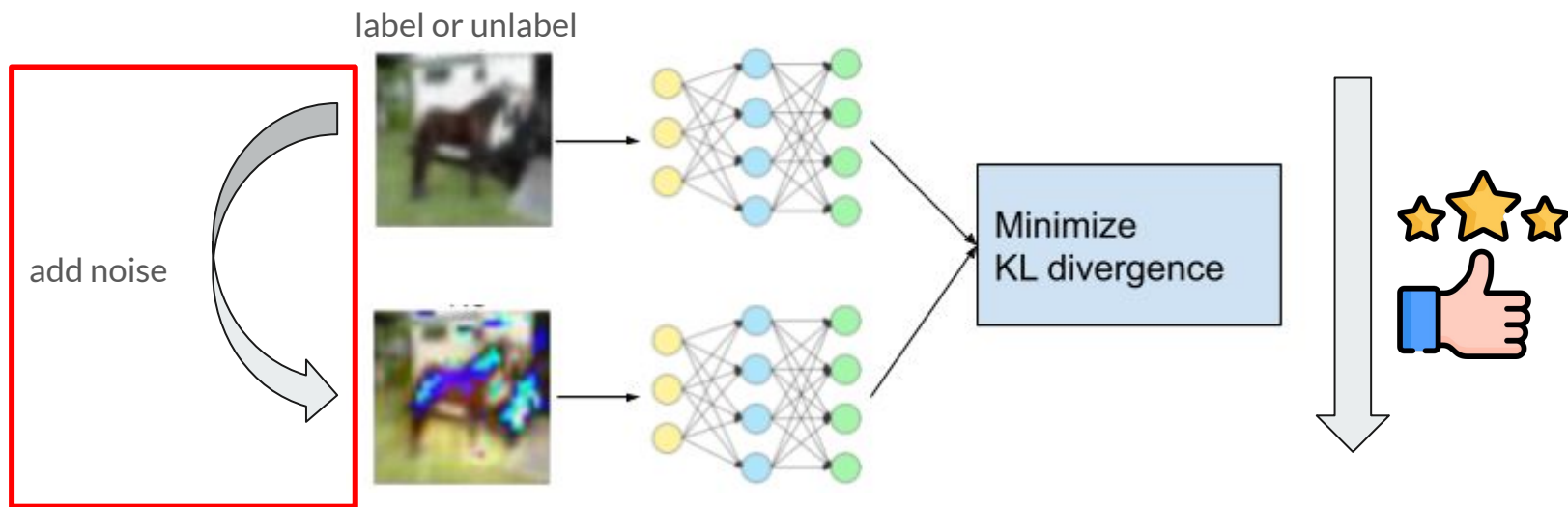
Dataset	Domain	Task	# Train / Test	# Class	Metric
Elec	Reviews	Sentiment Analysis	25K / 25K	2	Acc.
AG News	News	Topic Classification	120K / 7.6K	4	Acc.
NYT	News	Topic Classification	30K / 3.0K	9	Acc.
Chemprot	Chemical	Relation Classification	12K / 1.6K	10	F1

Dataset	Elec	AG News	NYT	Chemprot
description	Amazon shopping review	collection of news	New York Times	Contains PubMed abstracts containing chemical-protein interactions annotated by experts
category	positive、negative	World、Sports、Business、Sci/Tech	science、sports、music...	upregulator上調劑、downregulator下調劑 agonist激動劑 antagonist拮抗劑

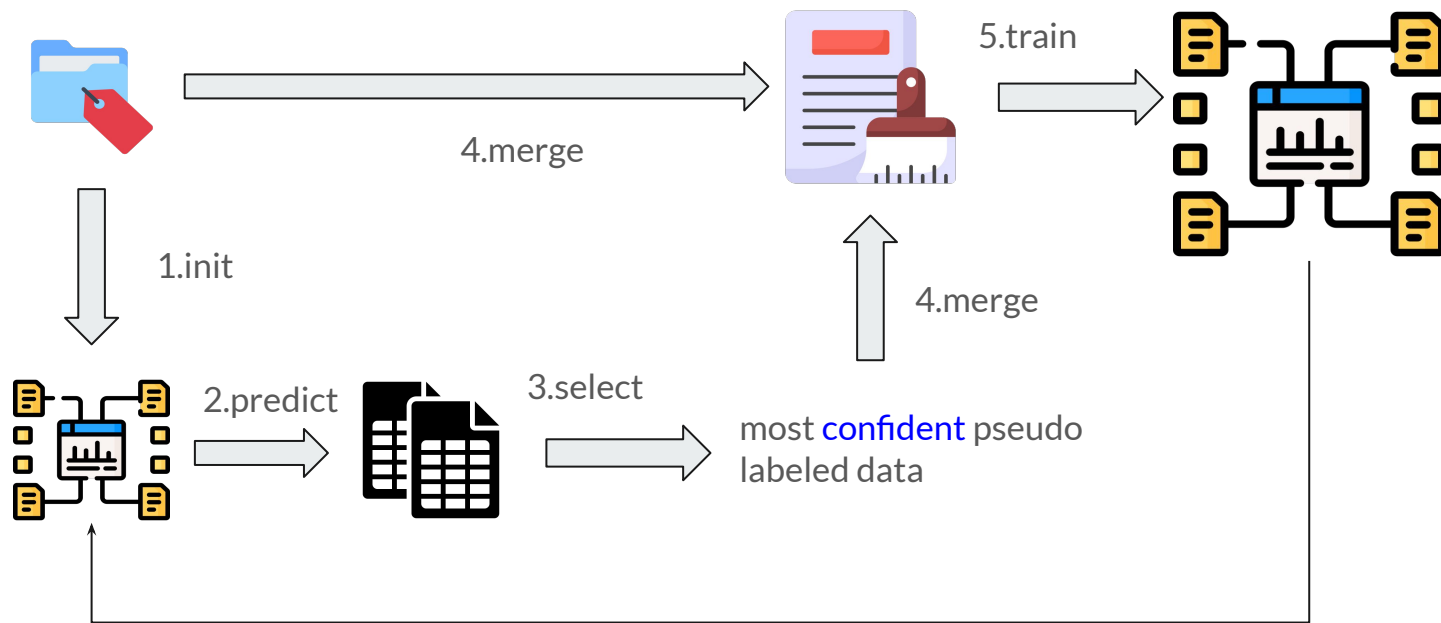
Baseline - Mean-Teacher



Baseline - Virtual Adversarial Training(VAT)

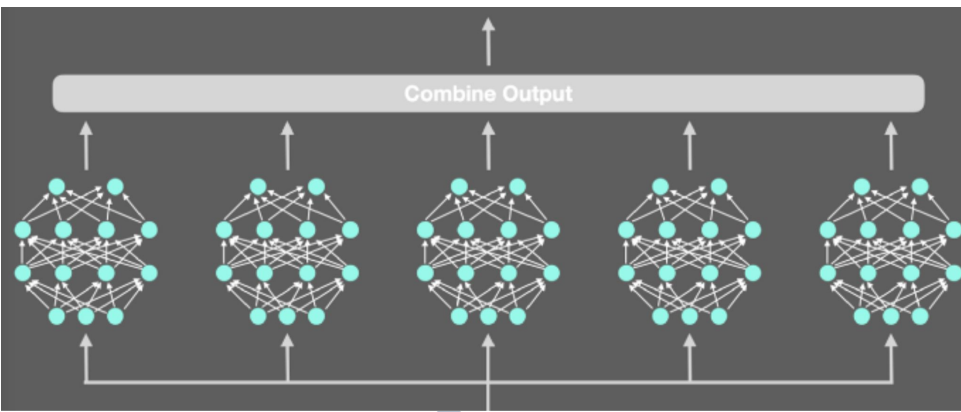


Baseline - Self-training(ST)

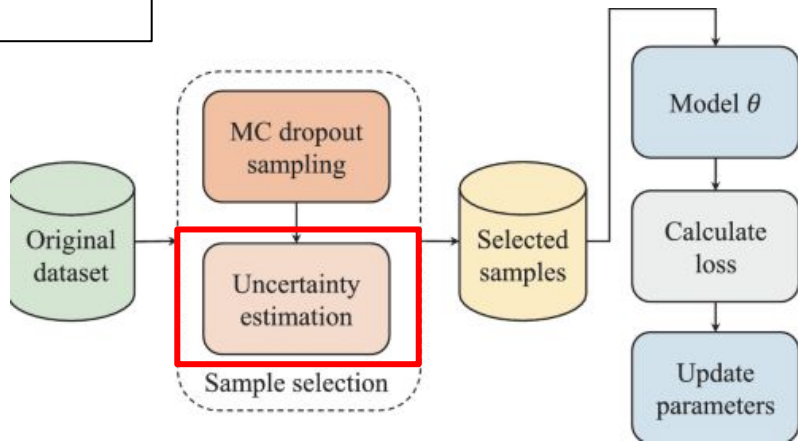


Baseline - Uncertainty-aware Self-training(UST)


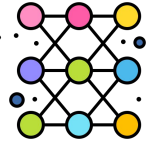
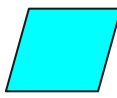
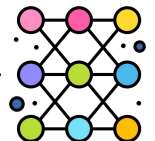
MC dropout



UST



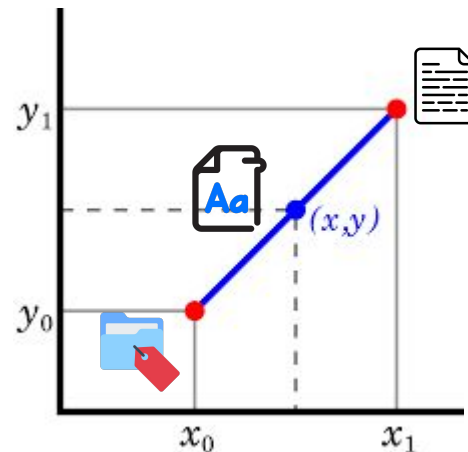
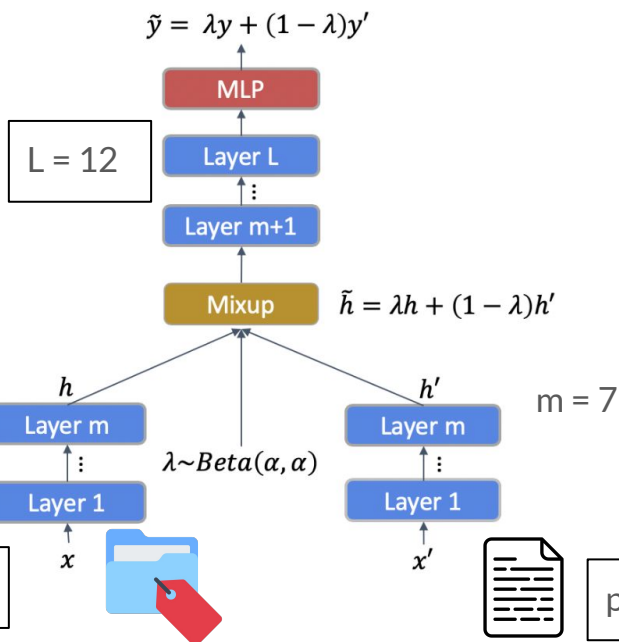
Uncertainty estimation - Entropy

	$\log P_{\theta}(y_i x)$	$P_{\theta}(y_i x)\log P_{\theta}(y_i x)$	$x_E^* = \underset{x}{\operatorname{argmax}} - \sum_i P_{\theta}(y_i x)\log P_{\theta}(y_i x)$
 →  <p>class A: 0.93 class B: 0.05 class C: 0.02</p>	<p>class A: -0.104 class B: -4.321 class C: -5.6438</p>	<p>class A: -0.09672 class B: -0.21605 class C: -0.11287</p>	<p>$-(0.09672+0.21605+0.11287) = -0.4256$</p> <p>$x_E^* = -(-0.4256) = 0.4256$</p> <p>entropy ↓ uncertain ↓</p>
 →  <p>class A: 0.55 class B: 0.35 class C: 0.1</p>	<p>class A: -0.8624 class B: -1.5145 class C: -3.3219</p>	<p>class A: -0.47432 class B: -0.53007 class C: -0.33219</p>	<p>$-(0.47432+0.53007+0.33219) = -1.33658$</p> <p>$x_E^* = -(-1.33658) = 1.33658$</p>

Baseline - MixText

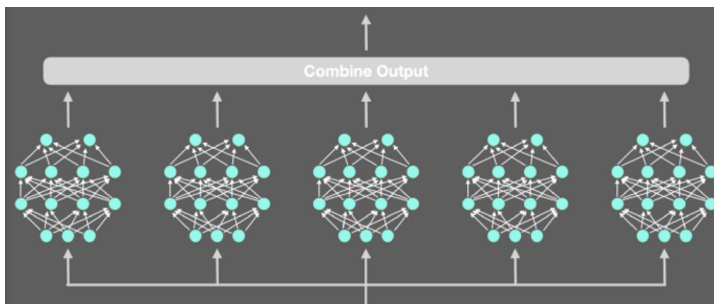
Data Augmentation by interpolating

BERT Base

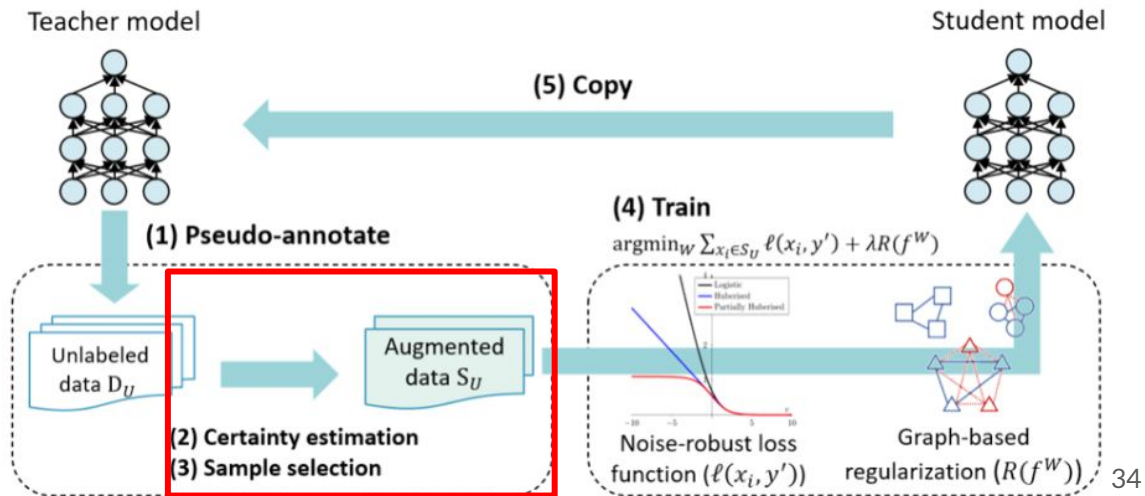


Baseline - Contrast-Enhanced Semi-supervised Text Classification(CEST)

MC dropout



CEST



Experiment

Method	AG News (Accuracy, \uparrow)		
	30	50	100
BERT	80.6 \pm 1.4		
MT	81.8 \pm 1.2		
VAT	82.1 \pm 1.2		
UDA	86.5 \pm 0.9		
MixText [†]	87.0 \pm 1.2		
<hr/>			
NeST	87.8\pm0.8		
Superv.		93.0*	

all unlabels are marked as pseudo labels and used to train the model

Method	Name	Description
MT	Mean Teacher	average model weight
VAT	Virtual Adversarial Training	add noise with unlabel
UDA	Unsupervised Data Augmentation	data augmentation with unlabel
MixText	MixText	data augmentation + interpolating with unlabel
ST	self-training	use strategy to select unlabel
UST	Uncertainty-aware Self-training	MCdropout + uncertainty to select unlabel
CEST	Contrast-Enhanced Semi-supervised	MCdropout + certainty + Graph-based Contrast
Nest	Neighborhood-Regularized Self-Training	KNN + self-training

- use all unlabel
 - data augmentation methods are more effective than BERT , e.g. UDA、MixText

Experiment

Method	AG News (Accuracy, ↑)		
	30	50	100

UDA	86.5±0.9
MixText [†]	87.0±1.2

} all unlabels are marked as pseudo labels and used to train the model

NeST	87.8±0.8
------	----------

Superv. 93.0*

Method	Name	Description
MT	Mean Teacher	average model weight
VAT	Virtual Adversarial Training	add noise with unlabel
UDA	Unsupervised Data Augmentation	data augmentation with unlabel
MixText	MixText	data augmentation + interpolating with unlabel
ST	self-training	use strategy to select unlabel
UST	Uncertainty-aware Self-training	MCdropout + uncertainty to select unlabel
CEST	Contrast-Enhanced Semi-supervised	MCdropout + certainty + Graph-based Contrast
Nest	Neighborhood-Regularized Self-Training	KNN + self-training

- wrong pseudo label causes model confusion

Experiment

Method	AG News (Accuracy, \uparrow)		
	30	50	100

ST	86.0 \pm 1.4	} use strategy to select unlabel
UST	86.9*	
CEST [‡]	86.5*	
NeST	87.8\pm0.8	
Superv.	93.0*	

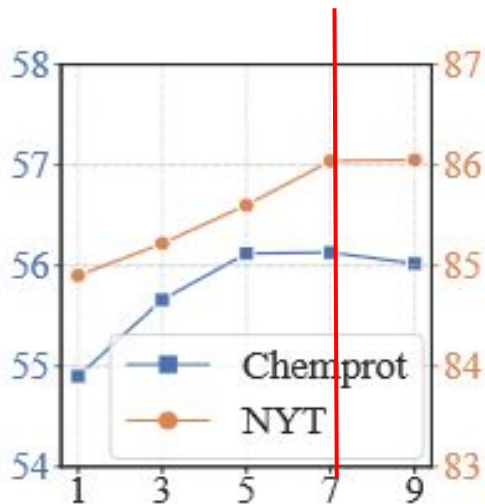
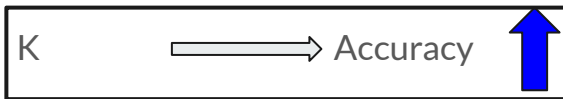
Method	Name	Description
MT	Mean Teacher	average model weight
VAT	Virtual Adversarial Training	add noise with unlabel
UDA	Unsupervised Data Augmentation	data augmentation with unlabel
MixText	MixText	data augmentation + interpolating with unlabel
ST	self-training	use strategy to select unlabel
UST	Uncertainty-aware Self-training	MCdropout + uncertainty to select unlabel
CEST	Contrast-Enhanced Semi-supervised	MCdropout + certainty + Graph-based Contrast
Nest	Neighborhood-Regularized Self-Training	KNN + self-training

- too much reliance on model predictions
- NeST is selected by aggregating the scores from the previous iteration

Experiment

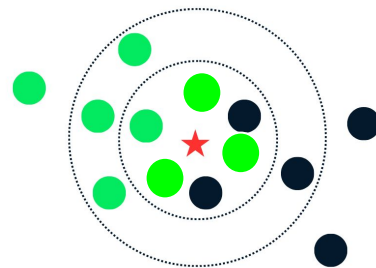
Method	AG News (Accuracy, \uparrow)			Elec (Accuracy, \uparrow)			NYT (Accuracy, \uparrow)			Chemprot (F1, \uparrow)		
	30	50	100	30	50	100	30	50	100	30	50	100
BERT	80.6 \pm 1.4	83.1 \pm 1.6	86.0 \pm 1.1	85.0 \pm 1.9	87.2 \pm 1.0	90.2 \pm 1.2	79.4 \pm 1.6	83.0 \pm 1.1	85.7 \pm 0.5	49.1 \pm 2.3	51.2 \pm 1.7	54.9 \pm 1.4
MT	81.8 \pm 1.2	83.9 \pm 1.4	86.9 \pm 1.1	87.6 \pm 0.9	88.5 \pm 1.0	91.7 \pm 0.7	80.2 \pm 1.1	83.5 \pm 1.3	86.1 \pm 1.1	50.0 \pm 0.7	54.1 \pm 0.8	56.8 \pm 0.4
VAT	82.1 \pm 1.2	85.0 \pm 0.8	87.5 \pm 0.9	87.9 \pm 0.8	89.8 \pm 0.5	91.5 \pm 0.4	80.7 \pm 0.7	84.4 \pm 0.9	86.5 \pm 0.6	50.7 \pm 0.7	53.8 \pm 0.4	57.0 \pm 0.5
UDA	86.5 \pm 0.9	87.1 \pm 1.2	87.8 \pm 1.2	89.6 \pm 1.1	91.2 \pm 0.6	92.3 \pm 1.0	—	—	—	—	—	—
MixText [†]	87.0 \pm 1.2	87.7 \pm 0.9	88.2 \pm 1.0	91.0 \pm 0.9	91.8 \pm 0.4	92.4 \pm 0.5	—	—	—	—	—	—
ST	86.0 \pm 1.4	86.9 \pm 1.0	87.8 \pm 0.6	89.6 \pm 1.2	91.4 \pm 0.4	92.1 \pm 0.5	<u>85.4\pm0.9</u>	<u>86.9\pm0.5</u>	<u>87.5\pm0.5</u>	<u>54.1\pm1.1</u>	55.3 \pm 0.7	59.3 \pm 0.5
UST	86.9*	87.4*	87.9*	90.0*	91.6*	91.9*	85.0 \pm 0.6	86.7 \pm 0.4	87.1 \pm 0.3	53.5 \pm 1.3	<u>55.7\pm0.4</u>	<u>59.5\pm0.7</u>
CEST [‡]	86.5*	87.0*	<u>88.4*</u>	<u>91.5*</u>	<u>92.1*</u>	<u>92.5*</u>	—	—	—	—	—	—
NeST	87.8\pm0.8	88.4\pm0.7	89.5\pm0.3	92.0\pm0.3	92.4\pm0.2	93.0\pm0.2	86.5\pm0.7	88.2\pm0.7	88.6\pm0.6	56.5\pm0.7	57.2\pm0.4	62.0\pm0.5
Superv.		93.0*			95.3*			93.6 \pm 0.5			82.5 \pm 0.4	

Parameter Studies



(a) k

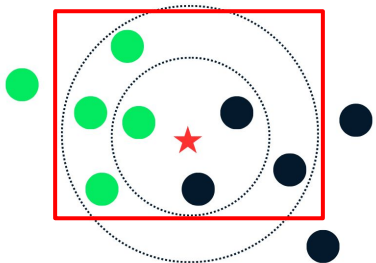
K = count of KNN



unlabel can find more labeled data



stable

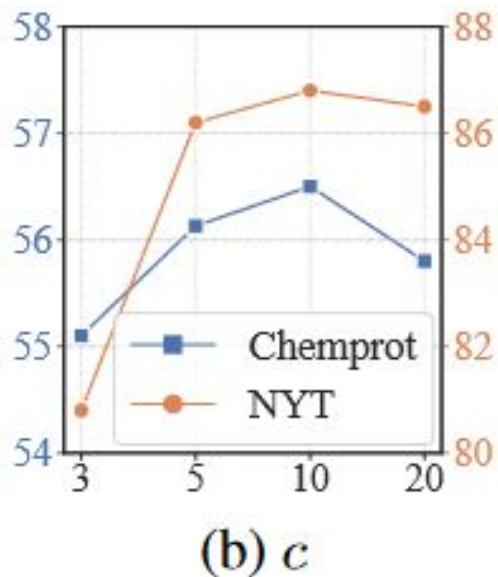


unlabel is too far away from label and has high divergence.



unstable

Parameter Studies



- if $c = 3$, labeled data = 120, $b=c|x_i| = 3 * 120 = 360$
 - the number of pseudo labels is not enough

↓

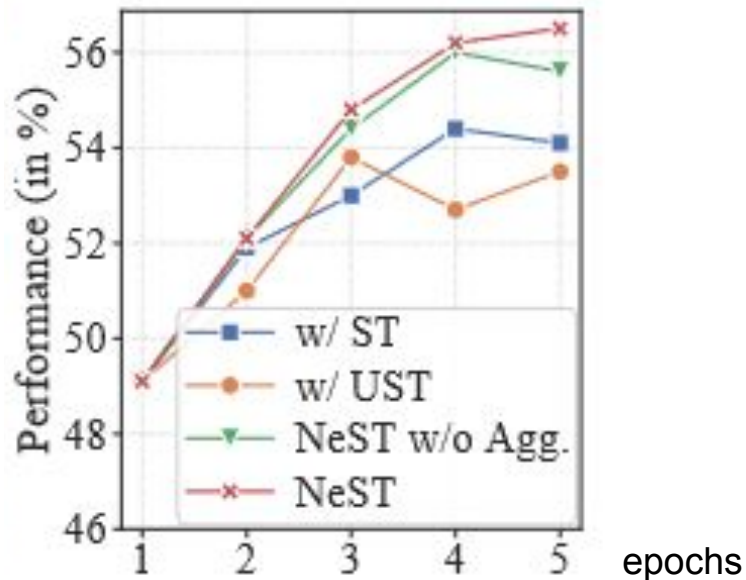
 - accuracy is not high
-
- if $c = 20$, labeled data = 400, $b=c|x_i| = 20 * 400 = 8000$
 - pseudo data selected is too messy and poor quality

↓

 - disrupt model learning

c : multiple of how many samples to select in an epoch

Ablation Studies

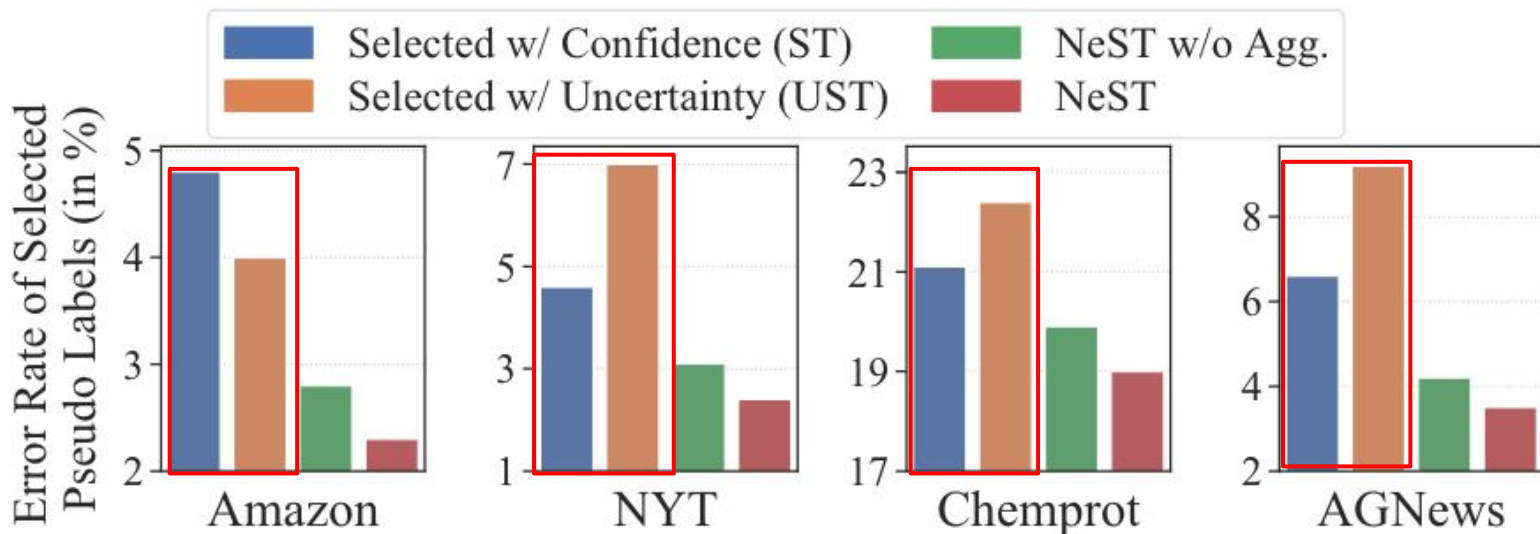


- in the early stage
 - selected pseudo labels that can help the model learn.

(c) Acc. over iters.

ST	self-training	use strategy to select unlabel
UST	Uncertainty-aware Self-training	MCdropout + uncertainty to select unlabel
Nest	Neighborhood-Regularized Self-Training	KNN + self-training

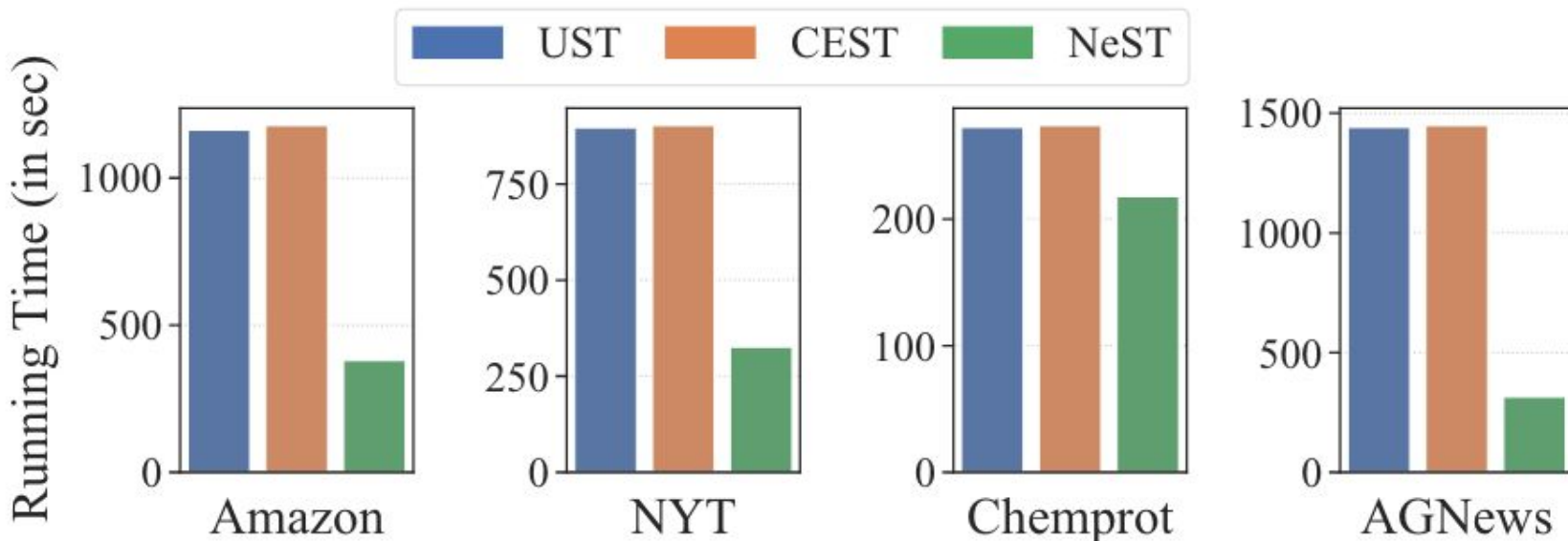
Error of Pseudo Labels





UST	Uncertainty-aware Self-training	MCdropout + uncertainty to select unlabel
CEST	Contrast-Enhanced Semi-supervised	MCdropout + certainty + Graph-based Contrast
Nest	Neighborhood-Regularized Self-Training	KNN + self-training

Running time of different methods





Outline

- Introduction
- Method
- Experiment
- **Conclusion**



Conclusion

- propose NeST to improve sample selection in self-training for robust label efficient learning
- design a neighborhood-regularized approach to select more reliable samples based on representations for self-training
- propose to aggregate the predictions on different iterations to stabilize self-training