Active Finetuning: Exploiting Annotation Budget in the Pretraining-Finetuning Paradigm

Yichen Xie¹, Han Lu², Junchi Yan², Xiaokang Yang², Masayoshi Tomizuka¹, Wei Zhan¹ ¹ University of California, Berkeley ² Shanghai Jiao Tong University {yichen_xie, tomizuka, wzhan}@berkeley.edu, {sjtu_luhan, yanjunchi, xkyang}@sjtu.edu.cn

Introduction

- recent success of deep learning heavily relies on abundant training data.
- the annotation of large-scale datasets often requires intensive human labor.
- inspires a popular pretraining-finetuning paradigm where models are pretrained on a large amount of data in an unsupervised manner and finetuned on a small labeled subset

Introduction

 A possible explanation comes from the batch-selection strategy of most current active learning methods. Starting from a random initial set, this strategy repeats the model training and data selection processes multiple times until the annotation budget runs out. Despite their success in from-scratch training, it does not fit this pretraining-finetuning paradigm well due to the typically low annotation budget, where too few samples in each batch lead to harmful bias inside the selection process.

Motivation

 To fill in this gap in the pretraining-finetuning paradigm,we formulate a new task called active finetuning, concentrating on the sample selection for supervised finetuning. In this paper, a novel method, ActiveFT, is proposed to deal with this task.
 Starting from purely unlabeled data, ActiveFT fetches a proper data subset for supervised finetuning in a negligible time. Without any redundant heuristics, we directly bring close the distributions between the selected subset and the entire unlabeled pool while ensuring the diversity of the selected subset. This goal is achieved by continuous optimization in the high-dimensional feature space, which is mapped with the pretrained model.



$$f(\cdot; w_0) : \mathcal{X} \to \mathbb{R}^C$$

$$\mathcal{P}^{u} = \{\mathbf{x}_{i}\}_{i \in [N]} \sim p_{u}$$
$$\mathcal{S} = \{s_{j} \in [N]\}_{j \in [B]}$$
$$\mathcal{P}^{u}_{\mathcal{S}} = \{\mathbf{x}_{s_{j}}\}_{j \in [B]} \subset \mathcal{P}^{u}$$
$$\mathcal{P}^{l}_{\mathcal{S}} = \{\mathbf{x}_{s_{j}}, \mathbf{y}_{s_{j}}\}_{j \in [B]}$$

Figure 1. **Pretraining-Finetuning Paradigm:** We focus on the selection strategy of a small subset from a large unlabeled data pool for annotation, named as active finetuning task, which is under-explored for a long time.

The goal of active finetuning is to find the sampling strategy minimizing the expected model error

$$S_{opt} = \arg\min_{\mathcal{S}} \mathop{E}_{\mathbf{x}, \mathbf{y} \in \mathcal{X} \times \mathcal{Y}} [error(f(\mathbf{x}; w_{\mathcal{S}}), \mathbf{y})] \quad (1)$$

The difference with traditional active learning :

- 1) We have access to the pretrained model f, which will be finetuned, before data selection.
- 2)The selected samples are applied to the finetuning of the pretrained model f instead of from-scratch training.
- 3) The sampled subset size is relatively small, less than 10% in most cases.
- 4) We have no access to any labels such as a random initial labeled set before data selection.

Select Samples:

- 1) bringing close the distributions between the selected subset \mathcal{P}_{S}^{u} and the original pool $\mathcal{P}^{u} \sim p_{u}$.
- 2) maintaining the diversity of $\mathcal{P}_{\mathcal{S}}^{u}$

$$S_{opt} = \arg\min_{\mathcal{S}} E_{\mathbf{x}, \mathbf{y} \in \mathcal{X} \times \mathcal{Y}} [error(f(\mathbf{x}; w_{\mathcal{S}}), \mathbf{y})] \quad (1)$$

$$S_{opt} = \arg\min_{S} D(p_{f_u}, p_{f_s}) - \lambda R(\mathcal{F}_{S}^{u})$$
(2)

$$\theta_{\mathcal{S},opt} = \arg\min_{\theta_{\mathcal{S}}} D(p_{f_u}, p_{\theta_{\mathcal{S}}}) - \lambda R(\theta_{\mathcal{S}}) \ s.t. \ ||\theta_{\mathcal{S}}^j||_2 = 1 \qquad \theta_{\mathcal{S}} = \{\theta_{\mathcal{S}}^j\}_{j \in [B]}$$
(3)

$$p_{\theta_{\mathcal{S}}}(\mathbf{f}) = \sum_{j=1}^{B} \phi_{j} p(\mathbf{f} | \theta_{\mathcal{S}}^{j})$$
(4)
$$p(\mathbf{f} | \theta_{\mathcal{S}}^{j}) = \frac{\exp(sim(\mathbf{f}, \theta_{\mathcal{S}}^{j})/\tau)}{Z_{j}}$$
(5)
$$c_{i} = \arg\max_{j \in [B]} sim(\mathbf{f}_{i}, \theta_{\mathcal{S}}^{j})$$
(6)

Assumption 1 $\forall i \in [N], j \in [B]$, if τ is small, the following far-more-than relationship holds that

 $\exp(sim(\mathbf{f}_i, \theta_{\mathcal{S}}^{c_i})/\tau) \gg \exp(sim(\mathbf{f}_i, \theta_{\mathcal{S}}^j)/\tau), j \neq c_i$

$$p_{\theta_{\mathcal{S}}}(\mathbf{f}_{i}) \approx \phi_{c_{i}} p(\mathbf{f}_{i} | \theta_{\mathcal{S}}^{c_{i}})$$

$$= \frac{\exp(sim(\mathbf{f}_{i}, \theta_{\mathcal{S}}^{c_{i}})/\tau)}{Z_{c_{i}}/\phi_{c_{i}}}$$

$$= \frac{\exp(sim(\mathbf{f}_{i}, \theta_{\mathcal{S}}^{c_{i}})/\tau)}{\tilde{Z}_{c_{i}}}$$
(7)

$$KL(p_{f_u}|p_{\theta_{\mathcal{S}}}) = \sum_{\mathbf{f}_i \in \mathcal{F}^u} p_{f_u}(\mathbf{f}_i) \log \frac{p_{f_u}(\mathbf{f}_i)}{p_{\theta_{\mathcal{S}}}(\mathbf{f}_i)}$$
$$= \sum_{\mathbf{f}_i \in \mathcal{F}^u} \left[\log p_{f_u}(\mathbf{f}_i)\right] - \sum_{\mathbf{f}_i \in \mathcal{F}^u} \left[\log p_{\theta_{\mathcal{S}}}(\mathbf{f}_i)\right]$$
(8)

$$D(p_{f_u}, p_{\theta_{\mathcal{S}}}) = -\frac{E}{\mathbf{f}_i \in \mathcal{F}^u} [sim(\mathbf{f}_i, \theta_{\mathcal{S}}^{c_i})/\tau]$$
(9)
$$R(\theta_{\mathcal{S}}) = -\frac{E}{j \in [B]} \left[\log \sum_{k \neq j, k \in [B]} \exp\left(sim(\theta_{\mathcal{S}}^j, \theta_{\mathcal{S}}^k)/\tau\right) \right]$$
(10)

$$L = D(p_{f_u}, p_{\theta_S}) - \lambda \cdot R(\theta_S)$$

= $- \underset{\mathbf{f}_i \in \mathcal{F}^u}{E} [sim(\mathbf{f}_i, \theta_S^{c_i})/\tau] + \underset{j \in [B]}{E} \left[\log \sum_{k \neq j, k \in [B]} \exp\left(sim(\theta_S^j, \theta_S^k)/\tau\right) \right]$
(11)

$$\mathbf{f}_{s_j} = \arg \max_{\mathbf{f}_k \in \mathcal{F}^u} sim(\mathbf{f}_k, \theta_{\mathcal{S}}^j)$$
(12)



Figure 2. Parametric Model Optimization Process: By optimizing the loss in Eq. 11, each parameter θ_{S}^{j} is appealed by nearby sample features (orange in the figure, Eq. 9) and repelled by other parameters θ_{S}^{k} , $k \neq j$ (green in the figure, Eq. 10).



Figure 3. Similarity between Features and Parameters: On CIFAR100 and ImageNet, we find the Top-20 most similar parameters θ_S^j with each sample feature \mathbf{f}_i , and calculate the average exponential similarity $E_{i \in [N]}[\exp(sim(\mathbf{f}_i, \theta_S^j)/\tau]]$. Here $\theta_S = \{\theta_S^j\}_{j \in [B]}$ is randomly sampled following the distribution p_{f_u} . The model $f(\cdot; w_0)$ is DeiT-Small [42] pretrained on ImageNet [37] with DINO framework [6]. The results verify Assumption 1 that the Top-1 similarity is significantly larger than others.

Implementation as a Learning Model

Algorithm 1: Pseudo-code for ActiveFT

Input: Unlabeled data pool $\{\mathbf{x}_i\}_{i \in [N]}$, pretrained model $f(\cdot; w_0)$, annotation budget B, iteration number T for optimization **Output:** Optimal selection strategy $S = \{s_j \in [N]\}_{j \in [B]}$ 1 for $i \in [N]$ do 2 $\lfloor \mathbf{f}_i = f(\mathbf{x}_i; w_0)$ /* Construct $\mathcal{F}^u = \{\mathbf{f}_i\}_{i \in [N]}$ based on \mathcal{P}^u , normalized to $||\mathbf{f}_i||_2 = 1$ */ 3 Uniformly random sample $\{s_j^0 \in [N]\}_{j \in [B]}$, and initialize $\theta_S^j = \mathbf{f}_{s_j^0}$ /* Initialize the parameter $\theta_S = \{\theta_S^j\}_{j \in [B]}$ based on \mathcal{F}^u */

4 for $iter \in [T]$ do Calculate the similarity between $\{\mathbf{f}_i\}_{i \in [N]}$ and 5 $\{\theta_{\mathcal{S}}^{j}\}_{j\in[B]}$: $Sim_{i,j} = \mathbf{f}_{i}^{\top}\theta_{\mathcal{S}}^{j}/\tau$ $MaxSim_i = \max_{i \in [B]} Sim_{i,j} = Sim_{i,c_i}$ 6 /* The Top-1 similarity between \mathbf{f}_i and $\theta_{S}^{j}, j \in [B]$ */ Calculate the similarity between θ_{S}^{j} and 7 $\theta_{S}^{k}, k \neq j$ for regularization: $RegSim_{j,k} = \exp(\theta_{\mathcal{S}}^{j} \theta_{\mathcal{S}}^{k}/\tau), k \neq j$ $Loss = -\frac{1}{N} \sum_{i \in [N]} MaxSim_i +$ 8 $\frac{1}{B}\sum_{j\in[B]}\log\left(\sum_{k\neq j}RegSim_{j,k}\right)$ /* Calculate the loss function in Eq. 11 $\theta_{\mathcal{S}} = \theta_{\mathcal{S}} - lr \cdot \bigtriangledown_{\theta_{\mathcal{S}}} Loss$ 9 /* Optimize the parameter through gradient descent */ $\theta_{S}^{j} = \theta_{S}^{j} / ||\theta_{S}^{j}||_{2}, j \in [B]$ 10 /* Normalize the parameters to ensure $||\theta_{S}^{j}||_{2} = 1$ */ 11 Find \mathbf{f}_{s_j} closest to $\theta_{\mathcal{S}}^j$: $s_j = \arg \max_{k \in [N]} \mathbf{f}_k^\top \theta_{\mathcal{S}}^j$ for each $j \in [B]$ 12 Return the selection strategy $S = \{s_i\}_{i \in [B]}$



Figure 3. Similarity between Features and Parameters: On CIFAR100 and ImageNet, we find the Top-20 most similar parameters $\theta_{\mathcal{S}}^{j}$ with each sample feature \mathbf{f}_{i} , and calculate the average exponential similarity $E_{i \in [N]} [\exp(sim(\mathbf{f}_{i}, \theta_{\mathcal{S}}^{j})/\tau]]$. Here $\theta_{\mathcal{S}} = \{\theta_{\mathcal{S}}^{j}\}_{j \in [B]}$ is randomly sampled following the distribution $p_{f_{u}}$. The model $f(\cdot; w_{0})$ is DeiT-Small [42] pretrained on ImageNet [37] with DINO framework [6]. The results verify Assumption 1 that the Top-1 similarity is significantly larger than others.

Table 1. Image Classification Results: Experiments are conducted on natural images with different sampling ratios. We report the mean and std over three trials. Explanation of N/A results ("-") is in our supplementary materials.

Mathada	CIFAR10			CIFAR100				ImageNet	
Methods	0.5%	1%	2%	1%	2%	5%	10%	1%	5%
Random	77.3±2.6	82.2±1.9	88.9±0.4	14.9±1.9	24.3±2.0	50.8±3.4	69.3±0.7	45.1±0.8	64.3±0.3
FDS	64.5±1.5	73.2±1.2	$81.4 {\pm} 0.7$	8.1±0.6	12.8 ± 0.3	16.9 ± 1.4	52.3±1.9	26.7 ± 0.6	55.5±0.1
K-Means	83.0±3.5	$85.9{\pm}0.8$	89.6±0.6	17.6±1.1	31.9±0.1	42.4±1.0	70.7±0.3	-	1
CoreSet [38]	-	81.6±0.3	88.4±0.2	-	$30.6 {\pm} 0.4$	48.3±0.5	62.9 ± 0.6	-	61.7±0.2
VAAL [39]	-	80.9 ± 0.5	$88.8{\pm}0.3$	-	24.6 ± 1.1	46.4 ± 0.8	70.1 ± 0.4	-	64.0 ± 0.3
LearnLoss [48]	-	81.6 ± 0.6	86.7 ± 0.4	-	19.2 ± 2.2	38.2 ± 2.8	65.7 ± 1.1	-	63.2 ± 0.4
TA-VAAL [21]	-	82.6±0.4	88.7 ± 0.2	-	34.7 ± 0.7	46.4 ± 1.1	66.8 ± 0.5	-	64.3±0.2
ALFA-Mix [33]	<u>1</u>	83.4±0.3	89.6±0.2	-	35.3 ± 0.8	50.4 ± 0.9	69.9±0.6	1 0	64.5±0.2
ActiveFT (ours)	85.0 ±0.4	88.2 ±0.4	90.1 ±0.2	26.1 ±2.6	40.7 ±0.9	54.6 ±2.3	71.0 ±0.5	50.1 ±0.3	65.3 ±0.1

Table 2. Semantic Segmentation Results: experiments are conducted on ADE20k with sampling ratios 5%, 10%. Results are averaged over three trials.

Sel. Ratio	Random	FDS	K-Means	ActiveFT (ours)
5%	14.54	6.74	13.62	15.37±0.11
10%	20.27	12.65	19.12	21.60 ±0.40

Table 3. **Data Selection Efficiency:** We compare the time cost to select different percentages of samples from the CIFAR100 training set.

Sel. Ratio	K-Means	CoreSet	VAAL	LearnLoss	ours
2%	16.6s	1h57m	7h52m	20m	12.6s
5%	37.0s	7h44m	12h13m	1h37m	21.9s
10%	70.2s	20h38m	36h24m	9h09m	37.3s



Figure 4. **tSNE Embeddings of CIFAR10:** We visualize the embedding of selected samples using different algorithms. Different colors denote categories, and the black dots are the 1% samples selected by our method.

Table 4. Generality on Pretraining Frameworks and Model Architectures: We examine the performance of ActiveFT on different pretraining frameworks and models on CIFAR-10.

Methods	0.5%	1%	2%				
Random	81.7	83.0	89.8				
CoreSet [38]	-	82.8	89.2				
LearnLoss [48]	-	83.6	89.2				
VAAL [39]	a .	85.1	89.3				
ActiveFT (ours)	87.6 ±0.8	88.3±0.2	90.9 ±0.2				
(b) Performance on ResNet-50 Pretrained with DINO							
Methods	0.5%	1%	2%				
Random	64.8	76.2	83.7				
CoreSet [38]	-	70.4	83.2				
LearnLoss [48]	-	71.7	81.3				
VAAL [39]	-	75.0	83.3				
ActiveFT (ours)	68.5 ±0.4	78.6 ±0.7	84.9 ±0.3				

(a) Performance on DeiT-Small Pretrained with iBOT

Table 5. **Ablation Study:** We examine the effect of two modules in our method. Experiments are conducted on CIFAR100 with pretrained DeiT-Small model.

(a) c_i Update Manner			(b) Regularization Design			
Ratio	No-Update	Update	Ratio	S1	S2	ours
2%	20.6	40.7	2%	33.1	26.8	40.7
5%	52.8	54.6	5%	51.5	46.9	54.6

Table 6. **Effect of Temperatures:** We try different temperatures in our method. Experiments are conducted on CIFAR10 with pre-trained DeiT-Small model.

Ratio	$\tau = 0.04$	$\tau = 0.07$	$\tau = 0.2$	$\tau = 0.5$
0.5%	85.6	85.0	84.1	83.5
1%	87.4	88.2	85.3	86.1
2%	90.3	90.1	89.6	89.0

Thanks