

# Model Selection

• Given a task and a large zoo of pretrained models, Model Selection (MS) selects the top few models for the best fine-tuning performance, avoiding fine-tuning of all models.

### **Feature-based Model Selection**

- Linearization Assumption
- Label-Feature Correlation (LFC)<sup>[1]</sup>  $\bullet$

$$S_{\mathbf{LFC}}(\mathbf{x}, \mathbf{y}) = (f_w(\mathbf{x}) f_w(\mathbf{x})^T) \cdot \mathbf{y} \mathbf{y}^T$$

• **PARC**<sup>[2]</sup> introduces heuristics of layer depth

$$S'_{\text{PARC}} = \frac{S_{\text{PARC}} - \mu^t}{\sigma^t} + \frac{\ell_s}{\ell_{\text{max}}}$$



Difficulty with heterogeneous models



Pareto Frontier Models are dataset dependent.

- MS can fail for ViTs when features are not normalized
- Weights may change a lot for hard/dissimilar datasets



Stanford Dogs is very close to ImageNet, its MS score is more accurate.

However, dataset like Aircrafts is much different from ImageNet, so its MS scores are much lower.

Hard to integrate more prior knowledge (e.g., capacity, dataset size) The scale is ad-hoc and it is not easy to integrate more signals or prior knowledge







# **Guided Recommendation for Model Fine-tuning**

Hao Li, Charless Fowlkes, Hao Yang, Onkar Dabeer, Zhuowen Tu, Stefano Soatto

Learning to Recommend Models

AWS AI Labs



		dataset features								model features				additional features			
		[			۱			[		)					(		
		d <sub>o</sub>	dı	d2	d3	d4	<b>d</b> 5	m <sub>o</sub>	$m_1$	<i>m</i> 2	m3	$m_4$	$m_5$	<i>s</i> 1	<b>s</b> <sub>2</sub>	у	
,	<b>x</b> 0	1	0	0	0	0	0.2	0	1	0	0	0	0.8	0.7	0.6	0.9	
2 2 6	<b>x</b> 1	0	0	1	0	0	0.3	0	0	0	1	0	0.5	0.6	0.4	0.8	
	x <sub>2</sub>	0	0	0	0	1	0.4	0	0	0	0	1	0.7	0.5	0.7	0.6	
	<b>X</b> 3	0	0	0	1	0	0.1	0	0	1	0	0	0.5	0.4	0.3	0.7	
	<b>X</b> 4	0	0	1	0	0	0.5	0	1	0	0	0	0.6	0.4	0.3	?	



### **Training History and Meta Features**



	dataset features								model features				additional features		
													[		٦
	d <sub>0</sub>	d <sub>1</sub>	d <sub>2</sub>	d <sub>3</sub>	d4	d <sub>5</sub>	<i>m</i> 0	<i>m</i> 1	<i>m</i> <sub>2</sub>	<i>m</i> 3	<i>m</i> 4	$m_5$	s <sub>1</sub>	s <sub>2</sub>	у
х <sub>0</sub>	1	0	0	0	0	0.2	0	1	0	0	0	0.8	0.7	0.6	0.9

- task difficulty: If a task can be solved with a simple model, then the task is relatively easy in comparison with other dataset.
- number of samples: a few-shot task is generally harder and often requires a strong model than a larger dataset size.
- number of classes: the task difficulty usually increase as the number of classes when the total images are fixed.

Features

Each job is a fine-tuning on a dataset with 8 HPs and the best result is obtained

### **Recommendation Models**

- Linear Regression (LR)  $S_{LR}(\mathbf{z}) = w_0 + \sum w_i z_i$
- Factorization Machines (FM)

$$S_{\text{FM}}(\mathbf{z}) = w_0 + \sum_{i=1}^{|z|} w_i z_i + \sum_{i=1}^{|z|-1} \sum_{j=i+1}^{|z|} \langle \mathbf{u}_i, \mathbf{u}_j \rangle z_i z_j$$

linear regression

feature interactions





### • We convert model selection as a model **recommendation** problem, which **learns** the model selection criteria from the past **fine-tuning history**.

### The goal is to predict performance on the target dataset for a given model.

### Embedding

- architecture family: architectures of the same family usually have similar inductive biases as they consist of similar modules.
- input size: architectures with higher resolution usually helps for downstream tasks.
- model capacity: a model with high capacity usually generalizes better with more data.
- model complexity: GMACs
- pre-trained domain: the pretrained domain matters for the downstream task performance

- MS score: it considers the feasibility of the model's initial features.
- semantic distance: semantic embedding of labels of the target task and the source task
- any features that are relevant for performance prediction

# **Experiments: More Training Data Helps!**



 Learning from the history of single dataset with **a subset** of models Evaluating unseen models on the same dataset.

Setting 1. MS learned with only ImageNet training history. 80% of the 409 models as sampled as training set. Our methods still get reasonable corr. scores even when models are random initialized

Methods	Features	ImageNet				
Methods	1 cutures	Pre-trained	Ran			
	$S_{ m LFC}$	$0.65\pm0.07$	0.0			
feature-based	$S_{ m LogME}$	$0.35\pm0.09$	0.04			
	$S_{\mathrm{PARC}}$	$0.83 \pm 0.04$	0.0			
	$\mathbf{d}, \mathbf{m}$	$0.53\pm0.07$	0.5			
$I \mathbf{D} (ours)$	$\mathbf{d}, \mathbf{m}, S_{ ext{LFC}}$	$0.73\pm0.06$	0.5			
LK (Ouis)	$\mathbf{d}, \mathbf{m}, S_{\text{LogME}}$	$0.55\pm0.08$	0.5			
	$\mathbf{d}, \mathbf{m}, S_{\mathrm{PARC}}$	$0.85 \pm 0.04$	0.5			
	$\mathbf{d}, \mathbf{m}$	$0.54\pm0.06$	0.5			
EM (ours)	$\mathbf{d}, \mathbf{m}, S_{ ext{LFC}}$	$0.70\pm0.12$	0.5			
	$\mathbf{d}, \mathbf{m}, S_{\text{LogME}}$	$0.55\pm0.09$	0.5			
	$\mathbf{d}, \mathbf{m}, S_{\text{PARC}}$	<b>0.84</b> $\pm$ 0.05	0.5			



### Reference

[1] A linearized framework and a new benchmark for model selection for fine-tuning, Deshpande et al, arXiv 2021 [2] Scalable Diverse Model Selection for Accessible Transfer Learning, Bolya et al, NeurIPS 2021 [3] LogME: Practical Assessment of Pre-trained Models for Transfer Learning, You et al, ICML 2021

					l					-	)
<b>z</b> <sub>2</sub>	<b>Z</b> 3	Z4	<b>z</b> 5	z <sub>6</sub>	Z <sub>7</sub>	Z <sub>8</sub>	Z <sub>9</sub>	z <sub>10</sub>	z <sub>11</sub>	z <sub>12</sub>	z <sub>13</sub>
0	0	0	0.2	0	1	0	0	0	0.8	0.7	0.6
d	( [		J	L		n	n 1		]	$\subseteq_{S_{ m LFC}}$	

<b>u</b> <sub>2</sub>	<b>u</b> <sub>3</sub>	$u_4$	$\mathbf{u}_5$	u <sub>6</sub>	<b>u</b> <sub>7</sub>	<b>u</b> <sub>8</sub>	u <sub>9</sub>	<b>u</b> <sub>10</sub>	<b>u</b> <sub>11</sub>	<b>u</b> <sub>12</sub>	<b>u</b> <sub>13</sub>
.4	.6	.7	.5	.2	.1	.2	.0	.0	.6	.7	.5
.2	.4	.5	.4	.3	.0	.0	.1	.0	.5	.6	.4
.1	.2	.1	.1	.1	.1	.0	.9	.1	.1	.3	.2
.5	.3	.4	.2	.4	.5	.1	.0	.2	.2	.2	.1





- Learning from the history of **single** dataset with **all** models Evaluating known models or
- unseen datasets



- Learning from the history of leave-oneout datasets. Evaluating known models on unseer
- tasks

**Setting 2 & 3**. Average Pearson Correlation of predicted performance and the ground-truth performance of 22 models. The ImageNet column is trained on 409 ImageNet training jobs. The LOO column denotes MS learned with the training history combining ImageNet and all other downstream jobs. MS trained with LOO outperforms over ImageNet only

Random Init.	downstream j						agenter	ny.
$0.03 \pm 0.10$	Methods	Features	19 fine-gr	ained	6 Domai	nNet	15 VT/	4B
$0.04\pm0.08$		$S_{\rm LFC}$	0.55		0.63		0.14	
$0.08\pm0.09$	feature-based MS	$S_{ m LogME}$	0.54	0.54		0.52		
$0.57 \pm 0.10$		S <sub>PARC</sub>	0.54		0.50		0.13	
$0.56 \pm 0.10$			ImageNet	LOO	ImageNet	LOO	ImageNet	LOO
$0.50 \pm 0.10$		$\mathbf{d}, \mathbf{m}$	<u>0.53</u>	<u>0.66</u>	<u>0.80</u>	<u>0.82</u>	<u>0.29</u>	<u>0.37</u>
$0.56 \pm 0.09$		$\mathbf{d}, \mathbf{m}, S_{ ext{LFC}}$	0.67	0.74	0.84	0.85	0.38	0.41
$0.57 \pm 0.11$	LR (ours)	$\mathbf{d}, \mathbf{m}, S_{LogME}$	0.54	0.65	0.81	0.84	0.30	0.36
$0.57 \pm 0.10$		$\mathbf{d}, \mathbf{m}, S_{\mathrm{PARC}}$	0.54	0.66	0.81	0.85	0.30	0.40
$0.56 \pm 0.10$		$\mathbf{d}, \mathbf{m}$	0.53	<u>0.65</u>	<u>0.81</u>	<u>0.85</u>	<u>0.35</u>	0.39
$0.50 \pm 0.10$		$\mathbf{d}, \mathbf{m}, S_{\mathrm{LFC}}$	0.64	0.74	0.82	0.87	0.39	0.41
$0.56 \pm 0.10$	FM (ours)	$\mathbf{d}, \mathbf{m}, S_{\text{LogME}}$	0.60	0.67	0.82	0.86	0.31	0.40
$0.57 \pm 0.11$		$\mathbf{d}, \mathbf{m}, S_{\mathrm{PARC}}$	0.56	0.69	0.86	0.86	0.30	0.43