



Contextual Active Model Selection

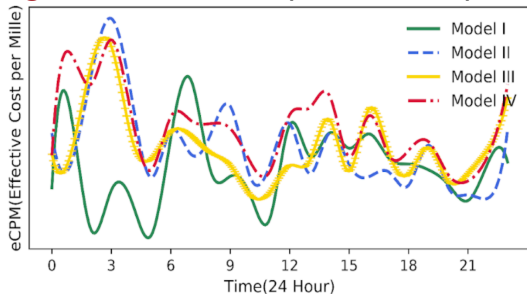
Xuefeng Liu¹ Fangfang Xia² Rick L. Stevens^{1,2} Yuxin Chen¹

¹University of Chicago

²Argonne National Laboratory

Motivation

- As **pre-trained models** become increasingly prevalent in a variety of real-world machine learning applications, there is a growing demand for **label-efficient** approaches for model selection
- **No single** pre-trained model achieves the best performance for every context
- Model performance depends on the **context**
- **Cost-sensitive** to evaluate and access the models / labels
- Online **streaming data** instead of a pool of data points



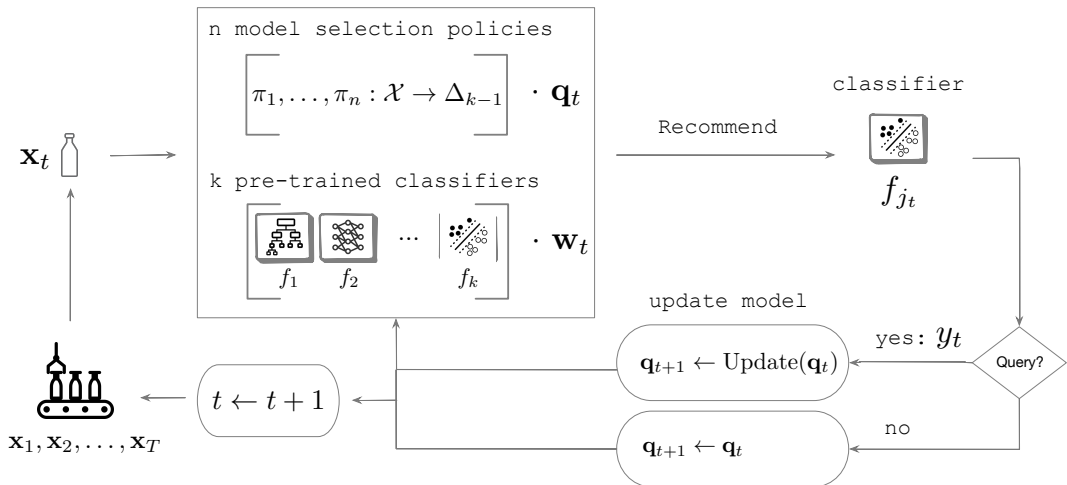
Research question

- How to select data-adaptive models when facing **heterogeneous data stream**?
- How to make it **labeling efficient**?

- We want a robust **cost-effective online-learning** algorithms that
 - effectively identify best model selection policy
 - works under limited labeling resources
 - adaptive to arbitrary data streams

We formally define it as the Contextual Active Model Selection (CAMS) problem and propose a novel algorithm, also named CAMS, to effectively address it.

Learning Protocol



Learning Protocol

Algorithm Contextual Active Model Selection Protocol

- 1: Given a set of classifiers \mathcal{F} and model selection policies Π
- 2: **for** $t = 1, 2, \dots, T$ **do**
- 3: The learner receives a data instance $\mathbf{x}_t \in \mathcal{X}$ as the context for the current round
- 4: Compute the predicted label $\hat{y}_{t,j} = f_j(\mathbf{x}_t)$ for pre-trained classifier indexed by $j \in [k]$
- 5: The learner identifies a model/classifier f_{j_t} and makes a prediction \hat{y}_{t,j_t} for the instance \mathbf{x}_t based on previous observations.
- 6: **if** The learner decide to query **then**
- 7: The learner incurs a **query cost**
- 8: The learner observes true label y_t and receives a (full) 0-1 loss vector $\ell_t = \mathbb{I}_{\{\hat{y}_t \neq y_t\}}$
- 9: The learner can then use the queried labels to adjust its model selection criterion for future rounds.

Method Overview

1: Input: Models \mathcal{F} , policies Π , #rounds T , budget b	
2: Initialize loss $\tilde{L}_0 \leftarrow 0$; query cost $C_0 \leftarrow 0$	
3: Set $\Pi^* \leftarrow \Pi \cup \{\pi_1^{\text{const}}, \dots, \pi_k^{\text{const}}\}$	
4: for $t = 1, 2, \dots, T$ do	
5: Receive \mathbf{x}_t	
6: $\eta_t \leftarrow \text{SETRATE}(t, \mathbf{x}_t, \Pi^*)$	} Contextual Model Selection
7: Set $q_{t,i} \propto \exp(-\eta_t \tilde{L}_{t-1,i}) \forall i \in \Pi^* $	
8: $j_t \leftarrow \text{RECOMMEND}(\mathbf{x}_t, \mathbf{q}_t)$	
9: Output $\hat{y}_{t,j_t} \sim f_{t,j_t}$ as the prediction for \mathbf{x}_t	
10: Compute $z_t = \max\{\delta_0^t, \mathfrak{E}(\hat{\mathbf{y}}_t, \mathbf{w}_t)\}$	
11: Sample $U_t \sim \text{Ber}(z_t)$	} Active Queries
12: if $U_t = 1$ and $C_t \leq b$ then	
13: Query the label y_t	
14: $C_t \leftarrow C_{t-1} + 1$	
15: Compute $\ell_t: \ell_{t,j} = \mathbb{I}\{\hat{y}_{t,j} \neq y_t\}, \forall j \in [\mathcal{F}]$	} Model Updates
16: Estimate model loss: $\hat{\ell}_{t,j} = \frac{\ell_{t,j}}{z_t}, \forall j \in [\mathcal{F}]$	
17: Update $\tilde{\ell}_t: \tilde{\ell}_{t,i} \leftarrow \langle \pi_i(\mathbf{x}_t), \hat{\ell}_{t,j} \rangle, \forall i \in [\Pi^*]$	
18: $\tilde{L}_t = \tilde{L}_{t-1} + \hat{\ell}_t$	
19: else	
20: $\tilde{L}_t = \tilde{L}_{t-1}$	
21: $C_t \leftarrow C_{t-1}$	
22: end if	
23: end for	
	21: procedure SETRATE(t, \mathbf{x}_t, m)
	22: if STOCHASTIC then
	23: $\eta_t = \sqrt{\frac{\ln m}{t}}$
	24: end if
	25: if ADVERSARIAL then
	26: Set ρ_t as in adversarial setting section
	27: $\eta_t = \sqrt{\frac{1}{\sqrt{t}} + \frac{\rho_t}{c^2 \ln c}} \cdot \sqrt{\frac{\ln m}{T}}$
	28: end if
	29: return η_t
	30: end procedure
	29: procedure RECOMMEND($\mathbf{x}_t, \mathbf{q}_t$)
	30: if STOCHASTIC then
	31: $\mathbf{w}_t = \sum_{i \in \Pi^* } q_{t,i} \pi_i(\mathbf{x}_t)$
	32: $j_t \leftarrow \text{maxind}(\mathbf{w}_t)$
	33: end if
	34: if ADVERSARIAL then
	35: $i_t \sim \mathbf{q}_t$
	36: $j_t \sim \pi_{i_t}(\mathbf{x}_t)$
	37: end if
	38: return j_t
	39: end procedure

We denote by $\bar{\ell}_t^y := \langle \mathbf{w}_t, \mathbb{I}\{\hat{\mathbf{y}}_t \neq y\} \rangle$ as the expected loss if the true label is y , where $\mathbf{w}_t = \pi_{\text{maxind}(\mathbf{q}_t)}(\mathbf{x}_t)$ and $\text{maxind}(\mathbf{w}) := \arg \max_{j: w_j \in \mathbf{w}} w_j$.

Comparison against Related Work

Algorithm \ Setup	Online bagging bagging	Hedge online learning	EXP3 bandit	EXP4 contextual bandits	Query by Committee active learning	ModelPicker model selection	CAMS (ours)
model selection	×	✓	✓	✓	×	✓	✓
full-information	✓	✓	×	×	✓	✓	✓
active queries	×	×	×	×	✓	✓	✓
context-aware	×	×	×	✓	×	×	✓

Theoretical Guarantees

- **Pseudo-regret** for **stochastic** setting

$$\bar{\mathcal{R}}_T(\mathcal{A}) = \mathbb{E}[L_T^{\mathcal{A}}] - T \min_{i \in [\Pi^*]} \mu_i, \quad (1)$$

where μ_i represents the expected loss of policy i if recommending the most probable model $\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{\mathbf{x}_t, y_t} [\ell_{t, \max \text{ind}(\pi_i(\mathbf{x}_t))}]$.

- **Expected regret** for **adversarial** setting

$$\mathcal{R}_T(\mathcal{A}) = \mathbb{E}[L_T^{\mathcal{A}}] - \min_{i \in [\Pi^*]} \sum_{t=1}^T \tilde{\ell}_{t,i}, \quad (2)$$

where $\tilde{\ell}_{t,i}$ represents the expected loss of policy i if **randomizing the model** recommendation at t , $\tilde{\ell}_{t,i} := \langle \pi_i(\mathbf{x}_t), \ell_t \rangle$

Theoretical Guarantees

- **Query complexity** and **tight regret bound** under
 - Stochastic data streams
 - Adversarial data streams
 - Finite policy / model classes

Algorithm	Regret	Query Complexity
Exp3	$2\sqrt{Tk \log k}$	-
Exp4	$\sqrt{2Tk \log G}$	-
Model Picker _{stochastic}	$62 \max_i \Delta_i k / (\lambda^2 \log k)$ $\lambda = \min_{j \in [k] \setminus \{i^*\}} \Delta_j^2 / \theta_j$	$\sqrt{2T \log k} (1 + 4 \frac{c}{\Delta})$
Model Picker _{adversarial}	$2\sqrt{2T \log k}$	$5\sqrt{T \log k} + 2L_{T,*}$
CAMS _{stochastic}	$\left(\frac{\ln \frac{ \Pi^* - 1}{\gamma} + \sqrt{\ln \Pi^* \cdot 2b^2 \ln \frac{2}{\delta}}}{\sqrt{\ln \Pi^* \Delta}} \right)^2$	$\left(\left(\frac{\ln \frac{ \Pi^* - 1}{\gamma} + \sqrt{\ln \Pi^* \cdot 2b^2 \ln \frac{2}{\delta}}}{\sqrt{\ln \Pi^* \Delta}} \right)^2 + T \mu_{i^*} \right) \frac{\ln T}{c \ln c}$
CAMS _{adversarial}	$2c \sqrt{\ln c / \max\{\rho_T, \sqrt{1/T}\}} \cdot \sqrt{T \log \Pi^* }$	$O \left(\left(\sqrt{\frac{T \log \Pi^* }{\max\{\rho_T, \sqrt{1/T}\}}} + \tilde{L}_{T,*} \right) (\ln T) \right)$

Experiments

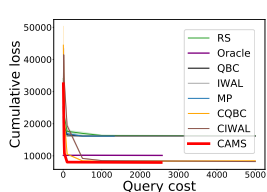
dataset	classification	total instances	test set	stream size	classifier	policy
CIFAR10	10	60000	10000	10000	80	85
DRIFT	6	13910	3060	3000	10	11
VERTEBRAL	3	310	127	80	6	17
HIV	2	40000	4113	4000	4	20
CovType	55	580000	100000	100000	6	17

- Baselines

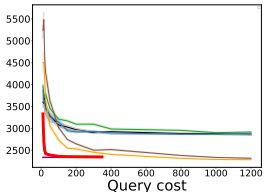
context-free: Random Sampling (RS, query the instance label with fixed probability), Query by Committee (QBC, committee-based sampling), Importance Weighted Active Learning (IWAL, Calculate query probability based on labeling disagreements of surviving classifiers), Model Picker (MP, employ variance based active sampling)

contextual: Contextual-QBC (CQBC), Contextual-IWAL (CIWAL), Oracle

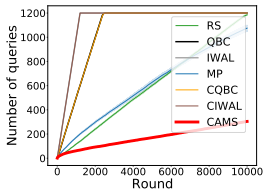
Experiments



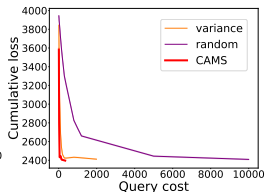
(a) COVTYPE



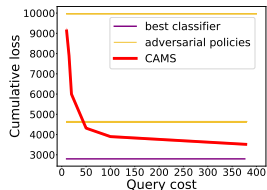
(b) CIFAR10



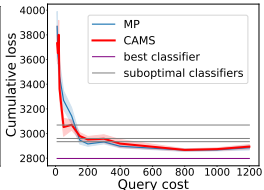
(c) Query complexity study, maintaining label efficiency



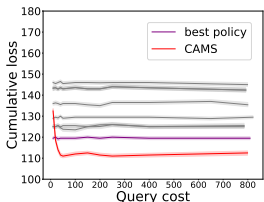
(d) Query strategy ablation



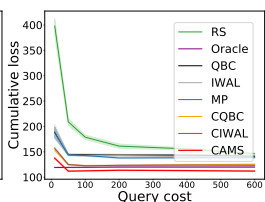
(e) Robustly recover from malicious advices



(f) Reduction to context-free setting when no experts are available



(g) CAMS vs best policy (HIV), effectively converge to the best expert



(h) Adjust probability (HIV), possibly outperforming all the experts including oracle