

Adaptive Region-Based Active Learning

Corinna Cortes Google Research New York, NY

Giulia DeSalvo Google Research New York, NY Claudio Gentile Google Research New York, NY Mehryar MohriNingshan ZhangGoogle Research & CourantNew York UniversityNew York, NYNew York, NY

Introduction

- **ARBAL** (Adaptive Region-Based Active Learning)
- 1. Adaptively partitions the input space into a finite number of regions
- 2. Subsequently seeks a distinct predictor for each region
- both phases actively requesting labels.

prove theoretical guarantees for both the generalization error and the label complexity

	Algorithm	Algorithm 1 ARBAL $(\mathcal{H}, \tau, \kappa, (\gamma_t)_{t \in t})$	[T]				
		$ \begin{array}{c} \hline K \leftarrow 1, \mathfrak{X}_1 \leftarrow \mathfrak{X}, \mathfrak{H}_1 \leftarrow \mathfrak{H} \\ \textbf{for } t \in [T] \textbf{ do} \end{array} \end{array} $					
	IWAL-query	Observe x_t ; set $k_t \leftarrow k$ such that x	$x_t \in \mathfrak{X}_k$				
		$p_t \leftarrow \max_{h,h' \in \mathcal{H}_{k_t}, y \in \mathcal{Y}} \ell(h(x_t), y_t) - $	$\ell(h'(x_t),y_t)$				
		$Q_t \leftarrow BERNOULLI(p_t)$					
		if $Q_t = 1$ then					
		$y_t \leftarrow \text{LABEL}(x_t)$					
		end if					
		if $t \leq \tau$ and $K < \kappa$ then					
		$\mathfrak{X}_l, \mathfrak{X}_r \leftarrow \operatorname{Split}(\mathfrak{X}_{k_t}, \gamma_t)$	#split phase				
		if split then					
-	The same hypothesis space <i>H</i>	$K \leftarrow K + 1, \mathfrak{X}_{k_t} \leftarrow \mathfrak{X}_l, \mathfrak{X}_K$	$\leftarrow \mathfrak{X}_r$				
		$\mathcal{H}_{K} \leftarrow \mathcal{H}, \mathcal{H}_{k_{t}} \leftarrow \mathcal{H}$					
		end if					
		else					
		$\mathcal{H}_{k_t} \leftarrow UPDATE(\mathcal{H}_{k_t})$	#IWAL phase				
		end if					
	end for						
	return $h_T \leftarrow \sum_{k=1}^{K} 1_{x \in \mathcal{X}_k} h_{k,T}$						

SPLIT phase

SPLIT splits a region if and only if the best-in-class error is likely to improve by a strictly positive amount

Algorithm 2 SPLIT (\mathfrak{X}_k, γ)

$$\begin{array}{ll} \text{for } d \in [D] \text{ and } c \in \mathbb{R} \text{ do} & \text{the splitting parameters } (d, c) \\ (\mathfrak{X}_l, \mathfrak{X}_r) \leftarrow \operatorname{REGSPLIT}(\mathfrak{X}_k, d, c) \\ \gamma_{d,c} \leftarrow \mathfrak{p}_k \Big[L_{k,t}(\widehat{h}_{k,t}) - L_{k,t}(\widehat{h}_{lr,t}) - \sqrt{\frac{2\sigma_T}{T_{k,t}}} \Big] \\ \text{end for} \\ (d^*, c^*) \leftarrow \operatorname{argmax}_{d \in [D], c \in \mathbb{R}} \gamma_{d,c} & L_{k,t}(h) = \frac{1}{T_{k,t}} \sum_{s \in [t], x_s \in \mathfrak{X}_k} \frac{Q_s}{p_s} \ell(h(x_s), y_s) \\ (d^*, c^*) \leftarrow \operatorname{argmax}_{d \in [D], c \in \mathbb{R}} \gamma_{d,c} & L_{k,t}(h) = \frac{1}{T_{k,t}} \sum_{s \in [t], x_s \in \mathfrak{X}_k} \frac{Q_s}{p_s} \ell(h(x_s), y_s) \\ (f_{\gamma d^*, c^*} \geq \gamma \text{ then} & \\ \mathfrak{X}_l^* \leftarrow \{x \in \mathfrak{X}_k : x[d^*] \leq c^*\} & \text{ \# split} \\ \mathfrak{X}_r^* \leftarrow \{x \in \mathfrak{X}_k : x[d^*] > c^*\} & \\ return \mathfrak{X}_l^*, \mathfrak{X}_r^* & \sigma_T = \kappa D \log \left[\frac{8T^3 |\mathcal{H}|^3 \kappa D}{\delta} \right] \\ else & \\ return \emptyset & \text{ \# no split} \\ end \text{ if} \end{array}$$

SPLIT phase

Lemma 1. With probability at least $1 - \delta/4$, for all binary trees with (at most) κ leaf nodes, the improvement in the minimal empirical error by splitting concentrates around the improvement in the best-in-class error:

$$\left| \left[R_k(h_k^*) - R_k(h_{lr}^*) \right] - \left[L_{k,t}(\widehat{h}_{k,t}) - L_{k,t}(\widehat{h}_{lr,t}) \right] \right| \leq \sqrt{\frac{2\sigma_T}{T_{k,t}}}.$$

Corollary 2. With probability at least $1 - \delta/4$, for all splits made by ARBAL, the improvement in the best-in-class error is at least γ_t , where γ_t is the threshold at the time of split.

$$\gamma_{d,c} \leftarrow \mathbf{p}_k \left[L_{k,t}(\widehat{h}_{k,t}) - L_{k,t}(\widehat{h}_{lr,t}) - \sqrt{\frac{2\sigma_T}{T_{k,t}}} \right]$$

$$\kappa D \log \left[\frac{8T^3 |\mathcal{H}|^3 \kappa D}{\delta} \right]$$

OT

$$\sigma_T = 0.01/\sqrt{T_k}$$

	Algorithm	Algorithm 1 ARBAL $(\mathcal{H}, \tau, \kappa, (\gamma_t)_{t \in [T]})$	_		
		$ \begin{array}{c} - \\ \hline K \leftarrow 1, \mathfrak{X}_1 \leftarrow \mathfrak{X}, \mathfrak{H}_1 \leftarrow \mathfrak{H} \\ \textbf{for } t \in [T] \textbf{ do} \end{array} $			
	IWAL-query	Observe x_t ; set $k_t \leftarrow k$ such that $x_t \in \mathfrak{X}_k$			
		$p_t \leftarrow \max_{h,h' \in \mathcal{H}_{k_t}, y \in \mathcal{Y}} \ell(h(x_t), y_t) - \ell(h'(x_t), y_t)$			
		$Q_t \leftarrow \text{Bernoulli}(p_t)$			
		if $Q_t = 1$ then			
		$y_t \leftarrow \text{LABEL}(x_t)$			
end if					
if $t \leq \tau$ and $K < \kappa$ then					
		$\mathfrak{X}_l, \mathfrak{X}_r \leftarrow \operatorname{SPLIT}(\mathfrak{X}_{k_t}, \gamma_t)$ #split phase	3		
		if split then			
The same hypothesis space H $\begin{array}{c} K \leftarrow K + 1, \mathfrak{X}_{k_t} \leftarrow \mathfrak{X}_l, \mathfrak{X}_K \leftarrow \mathfrak{X}_r \\ \mathfrak{H}_K \leftarrow \mathfrak{H}, \mathfrak{H}_{k_t} \leftarrow \mathfrak{H} \end{array}$					
		end if			
else					
		$\mathcal{H}_{k_t} \leftarrow \text{UPDATE}(\mathcal{H}_{k_t}) $ #IWAL phase			
		end if			
		end for return $\hat{h}_T \leftarrow \sum_{k=1}^K \mathbb{1}_{x \in \mathcal{X}_k} \hat{h}_{k,T}$			

SPLIT phase



ARBAL maintains throughout the split phase the original hypothesis space H

■ The shrinkage of *H* only takes place during the IWAL phase

IWAL phase

■ IWAL (Importance Weighted Active Learning)[1]

Based on the largest possible disagreement among the current set of hypotheses on the current input: flips a coin $Q_t \in \{0, 1\}$ with bias $p_t = p(x_t)$

$$p_{t} = \max_{h,h' \in \mathcal{H}_{t}, y \in \mathcal{Y}} \ell(h(x_{t}), y) - \ell(h'(x_{t}), y)$$
$$\hat{h}_{T} = \operatorname{argmin}_{h \in \mathcal{H}_{T}} \sum_{t=1}^{T} Q_{t} \ell(h(x_{t}), y_{t}) / p_{t}$$
$$\mathcal{H}_{k,t} = \left\{ h \in \mathcal{H}_{k,t-1} : L_{k,t}(h) \leq \min_{h \in \mathcal{H}_{k,t-1}} L_{k,t}(h) + \sqrt{\frac{8\sigma_{T}}{T_{k,t}}} \right\}$$

[1]Beygelzimer, A., Dasgupta, S., and Langford, J. Importance weighted active learning. In *Proceedings of ICML*, pp. 49–56. ACM, 2009.

Theoretical analysis

• with a fixed γ

Theorem 3. Assume that a run of ARBAL over T rounds has split the input space into K regions. Then, for any $\delta > 0$, with probability at least $1-\delta$, the following inequality holds: $\sqrt{32K\sigma_T} = \frac{16K\sigma_T}{16K\sigma_T}$

$$R(\hat{h}_T) \le R_U + \sqrt{\frac{32K\sigma_T}{T}} + \frac{16K\sigma_T}{T}$$

where $R_U = R^* - \gamma(K-1)$ is an upper bound on the best-in-class error obtained by ARBAL. Moreover, with probability at least $1 - \delta$, the expected number of labels requested, $\tau_T = \sum_{t=1}^T \mathbb{E}_{x_t \sim \mathcal{D}_X} [p_t | \mathcal{F}_{t-1}]$, satisfies $\tau_T \leq \min\{2\theta r_0, 1\}\tau + 4\theta_{\max}(T-\tau) \left[R_U + 8\sqrt{\frac{K\sigma_T}{T-\tau}}\right]$ $+ \sqrt{32}K\sigma_T$.

$$\sigma_T = \kappa D \log \left[\frac{8T^3 |\mathcal{H}|^3 \kappa D}{\delta} \right]$$

$$\sigma_T = 0.01/\sqrt{T_k}$$

Theoretical analysis

■ with a adaptive γ_t ρ : $R_k(h_k^*) - R_k(h_{lr}^*) \ge \rho$ $\rho = 0.01$ Proposition 4. Let ARBAL be run with $\gamma_t = \rho \mathbb{P}(\chi_{k_t})/2$. Then, for any $\delta > 0$, with probability at least $1 - \delta/2$, the first split occurs before round $\left[2\sigma_T(\frac{4}{\rho}+1)^2\right]$. $\sigma_T = \kappa D \log\left[\frac{8T^3|\mathcal{H}|^3\kappa D}{\delta}\right]$

 $\min\{\mathbb{P}(\mathfrak{X}_l), \mathbb{P}(\mathfrak{X}_r)\} \geq c \mathbb{P}(\mathfrak{X}_k), \text{ with } 0 < c < 0.5 \qquad \sigma_T = 0.01/\sqrt{T_k}$ **Corollary 5.** Let ARBAL run with $\gamma_t = \mathbb{P}(\mathfrak{X}_{k_t})\rho/2$. Then, with probability at least $1 - \delta/2$, ARBAL splits more than $\min\left\{\log_{1/c}\left[\frac{\tau}{2\sigma_T\left(\frac{4}{\rho}+1\right)^2}\right], \kappa - 1\right\} \text{ times by the end of the}$ split phase.

Experiment

logistic loss function

$$\sigma_T = 0.01/\sqrt{T_k}$$

The initial hypothesis set *H* consists of 3,000 randomly drawn hyperplanes with bounded

norms.



Figure 2: Misclassification loss of ARBAL with fixed and adaptive threshold γ on held out test data vs. number of labels requested (log₁₀ scale), with $\kappa = 20$ and $\tau = 800$. The vertical lines indicate the end of the first (split) phase.

Experiment



Figure 3: Misclassification loss of ARBAL(with adaptive γ_t), ORIWAL, IWAL, and MARGIN on hold out test data vs. number of labels requested (log₁₀ scale), with $\kappa = 20$ and $\tau = 800$. The ARBAL curves are repetitions from Figure 2.

Theoretical analysis

Define the distance $\rho(f,g)$ between two hypotheses $f,g \in \mathcal{H} \text{ as } \rho(f,g) = \mathbb{E}_{(x,y)\sim \mathcal{D}} |\ell(f(x),y) - \ell(g(x),y)|$ The generalized disagreement coefficient $\theta(\mathcal{D},\mathcal{H})$ of a class of functions \mathcal{H} with respect to distribution \mathcal{D} is defined as the minimum value of θ , such that for all r > 0,

$$\mathbb{E}_{x \sim \mathcal{D}_{\mathcal{X}}} \left[\sup_{h \in \mathcal{H}, \rho(h,h^*) \le r, y \in \mathcal{Y}} \left| \ell(h(x), y) - \ell(h^*(x), y) \right| \right] \le \theta r$$

disagreement coefficient $\theta_k = \theta(\mathcal{D}_k, \mathcal{H})$

 $r_0 = \max_{h \in \mathcal{H}} \rho(h, h^*)$. Let \mathcal{F}_t denotes the σ -algebra generated by $(x_1, y_1, Q_1), \ldots, (x_t, y_t, Q_t)$.