Noise-adaptive Marginbased Active Learning, and Lower Bounds

<u>Yining Wang</u>, Aarti Singh Carnegie Mellon University

## Machine Learning: the setup

- \* The machine learning problem
  - \* Each data point  $(x_i, y_i)$  consists of data  $x_i$  and label  $y_i$
  - \* Access to training data  $(x_1, y_1), \dots, (x_n, y_n)$
  - \* Goal: train classifier  $\hat{f}$  to predict *y* based on *x*
  - \* Example: Classification

$$x_i \in \mathbb{R}^d, y_i \in \{+1, -1\}$$

# Machine learning: passive vs. active

- Classical framework: passive learning
  - \* I.I.D. training data  $(x_i, y_i) \stackrel{i.i.d.}{\sim} D$
  - \* Evaluation: generalization error  $\Pr \left| y \neq \hat{f}(x) \right|$
- \* An active learning framework
  - \* Data are cheap, but labels are expensive!
  - \* Example: medical data (labels require domain knowledge)
  - \* Active learning: minimize label requests

# Active Learning

- Pool-based active learning
  - \* The learner *A* has access to unlabeled data stream  $x_1, x_2, \dots \stackrel{i.i.d.}{\sim} D$
  - \* For each  $x_i$ , the learner decides whether to query; if label requested, *A* obtains  $y_i$
  - \* Minimize number of requests, while scanning through polynomial number of unlabeled data.

# Active Learning

- \* Example: learning homogeneous linear classifier  $y_i = \operatorname{sgn}(w^{\top} x_i) + \operatorname{noise}$
- \* Basic (passive) approach: empirical risk minimization (ERM)  $\hat{w} \in \operatorname{argmin}_{\|w\|_2=1} \sum_{i=1}^n I[y_i \neq \operatorname{sgn}(w^\top x_i)]$
- \* How about active learning?

## Margin-based Active Learning

BALCAN, BRODER and ZHANG, COLT'07

- \* Data dimension *d*, query budget *T*, no. of iterations *E*
- \* At each iteration  $k \in \{1, \cdots, E\}$ 
  - \* Determine parameters  $b_{k-1}, \beta_{k-1}$
  - \* Find n = T/E samples in  $\{x \in \mathbb{R}^d : |\hat{w}_{k-1} \cdot x| \le b_{k-1}\}$
  - \* Constrained ERM:  $\hat{w}_k = \min_{\substack{\theta(w, \hat{w}_{k-1}) \le \beta_{k-1}}} L(\{x_i, y_i\}_{i=1}^n; w)$
- \* Final output:  $\hat{w}_E$

#### Tsybakov Noise Condition

\* There exist constants  $\mu > 0, \alpha \in (0, 1)$  such that  $\mu \cdot \theta(w, w^*)^{1/(1-\alpha)} \leq \operatorname{err}(w) - \operatorname{err}(w^*)$ 

\*  $\alpha \in (0,1)$  : key noise magnitude parameter in TNC

\* Which one is harder?



## Margin-based Active Learning

 Main Theorem [BBZ07]: when D is the uniform distribution, the margin-based algorithm achieves

$$\operatorname{err}(\hat{w}) - \operatorname{err}(w^*) = \widetilde{O}_P \left\{ \left(\frac{d}{T}\right)^{1/2\alpha} \right\}$$

Passive Learning:  $O((d/T)^{\frac{1-\alpha}{2\alpha}})$ 

#### **Proof outline**

BALCAN, BRODER and ZHANG, COLT'07

\* At each iteration *k*, perform *restricted* ERM over *withinmargin* data

$$\hat{w}_{k} = \underset{\theta(w, \hat{w}_{k-1}) \leq \boldsymbol{\beta}_{k-1}}{\operatorname{argmin}} \quad \widehat{\operatorname{err}}(w|S_{1}),$$
$$S_{1} = \{x : |x^{\top}\hat{w}_{k-1}| \leq \boldsymbol{b}_{k-1}\}$$

#### **Proof outline**

- \* Key fact: if  $\theta(\hat{w}_{k-1}, w^*) \leq \beta_{k-1}$  and  $b_k = \tilde{\Theta}(\beta_k/\sqrt{d})$  then  $\operatorname{err}(\hat{w}_k) - \operatorname{err}(w^*) = \tilde{O}\left(\beta_{k-1}\sqrt{d/T}\right)$
- Proof idea: decompose the excess error into two terms

$$\underbrace{\left[\operatorname{err}(\hat{w}_{k}|S_{1}) - \operatorname{err}(w^{*}|S_{1})\right]}_{\tilde{O}(\sqrt{d/T})} \underbrace{\Pr[x \in S_{1}]}_{\tilde{O}(b_{k-1}\sqrt{d})}$$

$$\underbrace{\operatorname{Forr}(\hat{w}_{k}|S^{c}) - \operatorname{err}(w^{*}|S^{c})]}_{\operatorname{Pr}[x \in S^{c}]} - \underbrace{\tilde{O}(\operatorname{tan}\beta_{k-1})}_{\operatorname{Pr}[x \in S^{c}]} = \underbrace{\tilde{O}(\operatorname{tan}\beta_{k-1}$$

 $\left[\operatorname{err}(\hat{w}_k | S_1^c) - \operatorname{err}(w^* | S_1^c)\right] \Pr[x \in S_1^c] = O(\tan \beta_{k-1})$ 

Must ensure w\* is always within reach!

$$\beta_k = 2^{\alpha - 1} \beta_{k - 2}$$

### Problem

- \* What if  $\alpha$  is not known? How to set key parameters  $b_k, \beta_k$
- \* If the true parameter is  $\alpha$  but the algorithm is run with  $\alpha' > \alpha$
- \* The convergence is  $\alpha'$  instead of  $\alpha$  !

## Noise-adaptive Algorithm

Agnostic parameter settings

$$E = \frac{1}{2} \log T, \beta_k = 2^{-k} \pi, b_k = \frac{2\beta_k}{\sqrt{d}} \sqrt{2E}$$

- Main analysis: two-phase behaviors
  - \* *"Tipping point"*:  $k^* \in \{1, \cdots, E\}$ , depending on  $\alpha$
  - \* *Phase I:*  $k \leq k^*$ , we have that  $\theta(\hat{w}_k, w^*) \leq \beta_k$
  - \* *Phase II:*  $k > k^*$ , we have that  $\operatorname{err}(\hat{w}_{k+1}) - \operatorname{err}(\hat{w}_k) \le \beta_k \cdot \widetilde{O}(\sqrt{d/T})$

### Noise-Adaptive Analysis

\* Main theorem: for all  $\alpha \in (0, 1/2)$ 

$$\operatorname{err}(\hat{w}) - \operatorname{err}(w^*) = \widetilde{O}_P \left\{ \left(\frac{d}{T}\right)^{1/2\alpha} \right\}$$

- \* Matching the upper bound in [BBZ07]
- \* ... and also a *lower bound* (this paper)

#### Lower Bound

- \* Is there any active learning algorithm that can do better than the  $\tilde{O}_P((d/T)^{1/2\alpha})$  sample complexity?
- \* In general, no [Henneke, 2015]. But the data distribution *D* is quite contrived in the negative example.
- \* We show that  $\tilde{O}_P((d/T)^{1/2\alpha})$  is tight even if *D* is as simple as the uniform distribution over unit sphere.

### Lower Bound

- \* The "Membership Query Synthesis" (QS) setting
  - \* The algorithm A picks an *arbitrary* data point  $x_i$
  - \* The algorithm receives its label  $y_i$
  - \* Repeat the procedure *T* times, with *T* the budget
- \* QS is more powerful than pool-based setting when *D* has density bounded away from below.
- \* We prove lower bounds for the QS setting, which implies lower bounds in the pool-based setting.

#### Tsybakov's Main Theorem

TSYBAKOV and ZAIATS, Introduction to Nonparametric Estimation

- \* Let  $\mathcal{F}_0 = \{f_0, \dots, f_M\}$  be a set of models. Suppose
  - \* Separation:  $D(f_j, f_k) \ge 2\rho, \forall j, k \in \{1, \dots, M\}, j \neq k$ \* Closeness:  $\frac{1}{M} \sum_{j=1}^{M} \operatorname{KL}(P_{f_j} \| P_{f_0}) \le \gamma \log M$
  - \* <u>Regularity</u>:  $P_{f_j} \ll P_{f_0}, \forall j \in \{1, \cdots, M\}$

Then the following bound holds

$$\inf_{\hat{f}} \sup_{f \in \mathcal{F}_0} \Pr_f \left[ D(\hat{f}, f) \ge \rho \right] \ge \frac{\sqrt{M}}{1 + \sqrt{M}} \left( 1 - 2\gamma - 2\sqrt{\frac{\gamma}{\log M}} \right)$$

### Negative Example Construction

- \* <u>Separation</u>:  $D(f_j, f_k) \ge 2\rho, \forall j, k \in \{1, \cdots, M\}, j \neq k$ 
  - \* Find hypothesis class  $\mathcal{W} = \{w_1, \cdots, w_m\}$  such that  $t \le \theta(w_i, w_j) \le 6.5t, \quad \forall i \ne j$
  - \* ... can be done for all  $t \in (0, 1/4)$ , using constant weight coding

\* ... can guarantee that  $\log |\mathcal{W}| = \Omega(d)$ 

#### Negative Example Construction



## Negative Example Construction

#### Lower Bound

TSYBAKOV and ZAIATS, Introduction to Nonparametric Estimation

- \* Let  $\mathcal{F}_0 = \{f_0, \cdots, f_M\}$  be a set of models. Suppose
  - \* <u>Separation</u>:  $D(f_j, f_k) \ge 2\rho, \forall j, k \in \{1, \dots, M\}, j \neq k$ \* <u>Closeness</u>:  $\frac{1}{M} \sum_{j=1}^{M} \operatorname{KL}(P_{f_j} \| P_{f_0}) \le \gamma \log M$ \* <u>Popularity</u>  $P_{\mathcal{A}} \ll P_{\mathcal{A}} \forall i \in \{1, \dots, M\}$
  - \* <u>Regularity</u>:  $P_{f_j} \ll P_{f_0}, \forall j \in \{1, \cdots, M\}$
- \* Take  $\rho = \Theta(t) = \Theta((d/T)^{(1-\alpha)/2\alpha})$   $\log M = \Theta(d)$
- \* We have that  $\inf_{\hat{w}} \sup_{w^*} \Pr\left[\theta(\hat{w}, w^*) \ge \frac{t}{2}\right] = \Omega(1)$

#### Lower Bound

\* Suppose *D* has density bounded away from below and fix  $\mu > 0, \alpha \in (0, 1)$ . Let  $\mathcal{P}_{Y|X}$  be class of distributions satisfying  $(\mu, \alpha)$ -TNC. Then we have that

$$\inf_{A} \sup_{P \in \mathcal{P}_{Y|X}} \mathbb{E}_{P}\left[\operatorname{err}(\hat{w}) - \operatorname{err}(w^{*})\right] \geq \Omega\left[\left(\frac{d}{T}\right)^{1/2\alpha}\right]$$

## Extension: "Proactive" learning

- Suppose there are *m* different users (labelers) who share the same classifier *w*<sup>\*</sup> but with different TNC parameters *α*<sub>1</sub>, · · · , *α*<sub>m</sub>
- \* The TNC parameters are *not* known.
- At each iteration, the algorithm picks a data point x and also a user *j*, and observes *f*(*x*;*j*)
- \* The goal is to estimate the Bayes classifier  $w^*$

## Extension: "Proactive" learning

- \* Algorithm framework:
  - \* Operate in  $E = O(\log T)$  iterations.
  - At each iteration, use conventional *Bandit* algorithms to address exploration-exploitation tradeoff
- \* Key property: search space  $\{\beta_k\}$  and margin  $\{b_k\}$  does *not* depend on unknown TNC parameters.
- \* Many interesting extensions: what if multiple labelers can be involved each time?

## Thanks! Questions?