# Trustworthy Machine Learning Statistical Query

Sangdon Park

POSTECH

#### Statistical Queries (SQ)

Efficient Noise-Tolerant Learning From Statistical Queries

> Michael Kearns \* AT&T Laboratories – Research Florham Park, New Jersey

> > June 15, 1998

#### 1 Introduction

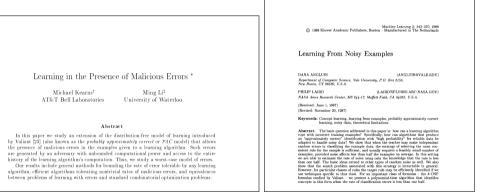
In this paper, we study the extension of Valiant's learning model [23] in which the positive or negative classification label provided with each random example may be corrupted by random noise. This extension was first examined in the learning theory literature by Angluin and Laird [1], which formalized the simplest type of white label noise and them sought algorithms tolerating the higher statement of a number of theorem. In this case, in additional state of a number of theorem at the classification noise model has become a common parading for experimental machine learning theory literature learning theory literature and the learning theory literature and the classification noise model has become a common parading for experimental machine learning research.

Anginia and Laird provided an algorithm for learning boelen conjunctions that tolerates an solise rate approaching the information-theoretic barrier of 1/2. Subsequently, there have been some isolated instances of efficient noise-tolerant algorithms [0, 27, 29], but little work on characterizing which classes can be efficiently learned in the presence of noise, and no general transformations of Valiant model algorithms into noise-tolerant algorithms. The primary conibution of the present paper is in making significant progress in both of these areas.

We identify and formalize an apparently rather weak sufficient condition on learning algorithms in Valant's model that permits the immediate derivation of noise-toferant learning algorithms. More precisely, we define a natural rearrietion on Valant model algorithms that allows them to be reliably and mains. This allows us to obtain efficient noise-toferant learning algorithms for perickally very cooped, and for which an efficient learning algorithms for min the state of the state

<sup>\*</sup>Author's address: AT&T Laboratories - Research, Room A235, 180 Park Avenue, Florham Park, New Jersey 07932. Electronic mail address: mkearns@research.att.com.

### Learning with Noise



#### • SQ generalizes learning with random classification noise.

- TL;DR: a generalized version of PAC learning for designing classification noise-tolerant PAC learning algorithms.
  - ▶ PAC learning: access to  $\mathsf{EX}(h^*, \mathcal{D})$  for a labeled example
  - ▶ PAC learning with random classification noise: access to  $\mathsf{EX}^{\eta}(h^*, \mathcal{D})$  for a labeled example
  - ▶ SQ: access to  $STAT(h^*, D)$  for a statistic of labeled examples

- TL;DR: a generalized version of PAC learning for designing classification noise-tolerant PAC learning algorithms.
  - ▶ PAC learning: access to  $\mathsf{EX}(h^*, \mathcal{D})$  for a labeled example
  - ▶ PAC learning with random classification noise: access to  $\mathsf{EX}^{\eta}(h^*, \mathcal{D})$  for a labeled example
  - $\blacktriangleright$  SQ: access to  $\mathsf{STAT}(h^*,\mathcal{D})$  for a statistic of labeled examples
- Main difference: access to the estimate of statistics from multiple samples, instead of one sample

- TL;DR: a generalized version of PAC learning for designing classification noise-tolerant PAC learning algorithms.
  - ▶ PAC learning: access to  $\mathsf{EX}(h^*, \mathcal{D})$  for a labeled example
  - ▶ PAC learning with random classification noise: access to  $\mathsf{EX}^{\eta}(h^*, \mathcal{D})$  for a labeled example
  - $\blacktriangleright$  SQ: access to  $\mathsf{STAT}(h^*,\mathcal{D})$  for a statistic of labeled examples
- Main difference: access to the estimate of statistics from multiple samples, instead of one sample
- Recall adversarial examples and noisy labels

- TL;DR: a generalized version of PAC learning for designing classification noise-tolerant PAC learning algorithms.
  - ▶ PAC learning: access to  $\mathsf{EX}(h^*, \mathcal{D})$  for a labeled example
  - PAC learning with random classification noise: access to  $\mathsf{EX}^\eta(h^*, \mathcal{D})$  for a labeled example
  - $\blacktriangleright$  SQ: access to  $\mathsf{STAT}(h^*,\mathcal{D})$  for a statistic of labeled examples
- Main difference: access to the estimate of statistics from multiple samples, instead of one sample
- Recall adversarial examples and noisy labels
- Why we have to learn? maybe useful for differential privacy and unlearning

# PAC Learning with Random Classification Noise

#### Definition (simplified definition)

An algorithm  $\mathcal{A}$  is a PAC-learning algorithm for  $\mathcal{H}$  with random classification noise if for any  $\varepsilon > 0$ ,  $\delta > 0$ ,  $0 \le \eta \le \frac{1}{2}$ ,  $h^* \in \mathcal{H}$ , and  $\mathcal{D}$  separable by  $h^*$ , and for some minimum sample size n' (which depends on  $\varepsilon, \delta, \eta, \mathcal{D}$ ), the following holds with any sample size  $n \ge n'$ :

 $\mathbb{P}\left\{L(\mathcal{A}(\mathcal{S})) \le \varepsilon\right\} \ge 1 - \delta,$ 

where  $\mathcal{S} \coloneqq ((x_1, y_1), \dots, (x_n, y_n))$  and  $(x_i, y_i) \sim \mathsf{EX}^{\eta}(h^*, \mathcal{D})$ .

- Suppose binary classification
- $\eta$ : a noise rate
- EX $^{\eta}(h^*, \mathcal{D})$ : the noisy example oracle that randomly flips the label with  $\eta$  probability
- If  $\eta = \frac{1}{2}$ , no hope to learn.

# Setup for SQ

- $\bullet~ \mathcal{D}:$  a distribution over  $\mathcal X$
- Binary classification, i.e.,  $\mathcal{Y} \coloneqq \{0, 1\}$

# Setup for SQ

- $\bullet~ \mathcal{D}:$  a distribution over  $\mathcal X$
- Binary classification, i.e.,  $\mathcal{Y} \coloneqq \{0, 1\}$
- $\chi: \mathcal{X} \times \mathcal{Y} \rightarrow [-1, 1]$ : a map (e.g., 0-1 loss)
- $\bullet~\alpha:$  a tolerance parameter
- $(\chi, \alpha)$ : a statistical query

# Setup for SQ

- $\bullet~ \mathcal{D}:$  a distribution over  $\mathcal X$
- Binary classification, *i.e.*,  $\mathcal{Y} \coloneqq \{0, 1\}$
- $\chi: \mathcal{X} \times \mathcal{Y} \rightarrow [-1, 1]$ : a map (e.g., 0-1 loss)
- $\alpha$ : a tolerance parameter
- $(\chi, \alpha)$ : a statistical query
- STAT $(h, D) : (\chi, \alpha) \mapsto [-1, 1]$ : a statistical query oracle (*i.e.*, a data source)
- $v \in [0,1]$ : The response of  $\mathsf{STAT}(h^*,\mathcal{D})$ , where

$$\left| \mathop{\mathbb{E}}_{x \sim \mathcal{D}} \left\{ \chi(x, h^*(x)) \right\} - v \right| \le \alpha.$$

Suppose that the statistical query oracle satisfies this with probability one.

#### Definition (simplified definition)

An algorithm  $\mathcal{A}$  is a statistical query algorithm for  $\mathcal{H}$  if for any  $\varepsilon > 0$ ,  $h^* \in \mathcal{H}$ , and  $\mathcal{D}$  over  $\mathcal{X}$ , and for some minimum number of queries n' (which depends on  $\varepsilon, \delta, \mathcal{D}$ ), the following holds with any number of queries  $n \ge n'$ :

 $L(\mathcal{A}(\mathcal{S})) \leq \varepsilon,$ 

where  $\mathcal{S} \coloneqq (v_1, \dots, v_n)$  and  $v_i \sim \mathsf{STAT}(h^*, \mathcal{D})$ .

#### Definition (simplified definition)

An algorithm  $\mathcal{A}$  is a statistical query algorithm for  $\mathcal{H}$  if for any  $\varepsilon > 0$ ,  $h^* \in \mathcal{H}$ , and  $\mathcal{D}$  over  $\mathcal{X}$ , and for some minimum number of queries n' (which depends on  $\varepsilon, \delta, \mathcal{D}$ ), the following holds with any number of queries  $n \ge n'$ :

 $L(\mathcal{A}(\mathcal{S})) \leq \varepsilon,$ 

where  $\mathcal{S} \coloneqq (v_1, \ldots, v_n)$  and  $v_i \sim \mathsf{STAT}(h^*, \mathcal{D})$ .

• Suppose that the statistic is efficiently computed.

#### Definition (simplified definition)

An algorithm  $\mathcal{A}$  is a statistical query algorithm for  $\mathcal{H}$  if for any  $\varepsilon > 0$ ,  $h^* \in \mathcal{H}$ , and  $\mathcal{D}$  over  $\mathcal{X}$ , and for some minimum number of queries n' (which depends on  $\varepsilon, \delta, \mathcal{D}$ ), the following holds with any number of queries  $n \ge n'$ :

 $L(\mathcal{A}(\mathcal{S})) \leq \varepsilon,$ 

where  $\mathcal{S} \coloneqq (v_1, \ldots, v_n)$  and  $v_i \sim \mathsf{STAT}(h^*, \mathcal{D})$ .

- Suppose that the statistic is efficiently computed.
- No confidence parameter  $1 \delta$ ; it is required in the statistical query oracle.

#### Definition (simplified definition)

An algorithm  $\mathcal{A}$  is a statistical query algorithm for  $\mathcal{H}$  if for any  $\varepsilon > 0$ ,  $h^* \in \mathcal{H}$ , and  $\mathcal{D}$  over  $\mathcal{X}$ , and for some minimum number of queries n' (which depends on  $\varepsilon, \delta, \mathcal{D}$ ), the following holds with any number of queries  $n \ge n'$ :

 $L(\mathcal{A}(\mathcal{S})) \leq \varepsilon,$ 

where  $\mathcal{S} \coloneqq (v_1, \ldots, v_n)$  and  $v_i \sim \mathsf{STAT}(h^*, \mathcal{D})$ .

- Suppose that the statistic is efficiently computed.
- No confidence parameter  $1 \delta$ ; it is required in the statistical query oracle.
- If an algorithm is a SQ algorithm, then it is a tolerant algorithm for random classification noise.

# **Example: Stochastic Convex Optimization**

#### Setup:

- $\bullet~\mathcal{H}:$  a set of convex functions
- $\ell(h,z):$  convex and sub-differentiable in h
- separable assumption (i.e.,  $\mathbb{E}\{\ell(h^*,z)\}=0$  for some  $h^*\in\mathcal{H}$ )

### **Example: Stochastic Convex Optimization**

#### Setup:

- $\bullet~\mathcal{H}:$  a set of convex functions
- $\ell(h,z)$ : convex and sub-differentiable in h
- separable assumption (*i.e.*,  $\mathbb{E}\{\ell(h^*, z)\} = 0$  for some  $h^* \in \mathcal{H}$ )

#### Stochastic convex optimization (in learning):

 $\min_{h \in \mathcal{H}} \mathbb{E}_{z \sim \mathcal{D}} \left\{ \ell(h, z) \right\}$ 

• We use the mirror descent algorithm to solve this.

- A generalized version of the gradient descent algorithm.
- The mirror descent considers the "geometry" of optimization.

- A generalized version of the gradient descent algorithm.
- The mirror descent considers the "geometry" of optimization.
- Example: a proximal gradient descent algorithm
  - Minimization:

 $\min_{h \in \mathcal{H}} \mathbb{E}_{z \sim \mathcal{D}} \{\ell(h, z)\}$ 

- A generalized version of the gradient descent algorithm.
- The mirror descent considers the "geometry" of optimization.
- Example: a proximal gradient descent algorithm
  - Minimization:

$$\min_{h \in \mathcal{H}} \mathbb{E}_{z \sim \mathcal{D}} \{\ell(h, z)\}$$

Minimization with linear approximation:

$$\begin{split} h_{t+1} &\leftarrow \arg\min_{h} \left\{ \eta \underbrace{\left(\ell(h_{t}, z_{t}) + \langle \nabla \ell(h_{t}, z_{t}), h - h_{t} \rangle\right)}_{\text{linear approximation}} + \frac{1}{2} \underbrace{\|h - h_{t}\|_{2}^{2}}_{\text{regularizer}} \right\} \\ h_{t+1} &\leftarrow \arg\min_{h} \left\{ \eta \langle \nabla \ell(h_{t}, z_{t}), h \rangle + \frac{1}{2} \|h - h_{t}\|_{2} \right\} \end{split}$$

- A generalized version of the gradient descent algorithm.
- The mirror descent considers the "geometry" of optimization.
- Example: a proximal gradient descent algorithm
  - Minimization:

$$\min_{h \in \mathcal{H}} \mathbb{E}_{z \sim \mathcal{D}} \{\ell(h, z)\}$$

Minimization with linear approximation:

$$\begin{split} h_{t+1} &\leftarrow \arg\min_{h} \left\{ \eta \underbrace{\left(\ell(h_{t}, z_{t}) + \langle \nabla \ell(h_{t}, z_{t}), h - h_{t} \rangle\right)}_{\text{linear approximation}} + \frac{1}{2} \underbrace{\|h - h_{t}\|_{2}^{2}}_{\text{regularizer}} \right\} \\ h_{t+1} &\leftarrow \arg\min_{h} \left\{ \eta \langle \nabla \ell(h_{t}, z_{t}), h \rangle + \frac{1}{2} \|h - h_{t}\|_{2} \right\} \end{split}$$

> This is equivalent to the conventional gradient descent algorithm, *i.e.*,

$$\eta \nabla \ell(h_t, z_t) + (h_{t+1} - h_t) = 0 \implies h_{t+1} = h_t - \eta \nabla \ell(h_t, z_t)$$

#### Gradient Descent Algorithm with SQ

#### Algorithm (for $\|\cdot\|_2$ ):

$$\bar{h} \coloneqq \frac{1}{T} \sum_{t=1}^{T} h_t \quad \text{where} \quad h_{t+1} \leftarrow \arg\min_h \left\{ \eta \langle \bar{\boldsymbol{g}}_t, h \rangle + \frac{1}{2} \| h - h_t \|_2^2 \right\}$$

• 
$$\bar{g}_t \coloneqq \frac{1}{K} \sum_{k=1}^K \nabla \ell(f_t, z_t^{(k)})$$

•  $\mathsf{STAT}(h^*,\mathcal{D})=\bar{g}_t$ : an " $\alpha$ -good" gradient value

• The above algorithm is a SQ algorithm, i.e.,  $\mathbb{E}\{\ell(\bar{h},z)\} \leq \varepsilon$ .