# Learning with Statistical Queries

Yuufen

PAC model ✓

What about errors in the labels?

One model for this:

Random Classification Noise Model

each label is flipped with prob. $\eta$.

$$(x, \ell(x)) \quad \ell(x) = \begin{cases} f(x) & w.p. \ 1-\eta \\ 1-f(x) & w.p. \ \eta \end{cases}$$

Q. Can we still learn the underlying concept?

E.g. suppose we are learning an <u>OR</u>.
of Boolean variables.

Let $p_i = \Pr(f(x)=0 \text{ and } X_i = 1)$

we can let output hypothesis $h$ be
all variables for which $p_i = 0$

(any variable with $p_i = 0$ is in the true OR)

(any variable with $p_i = 0$ is ... ...

and no variables for which $p_i > \frac{\varepsilon}{n}$.

So can we estimate each $p_i \pm \frac{\varepsilon}{2n}$ (say)?

$$p_i = \underbrace{Pr\left(f(x)=0 \middle/ X_i = 1\right)}_{q_i} \cdot Pr(X_i = 1)$$

↖ independent of label noise

$$Pr_{\text{noise } \eta}\left(l(x)=0 \middle/ X_i = 1\right) = (1-\eta)\, q_i + \eta\,(1-q_i)$$
$$= \eta + q_i(1-2\eta)$$

So if we have LHS, we can subtract $\eta$, divide by $(1-2\eta)$.

suffices to approximate to within $\pm \frac{\varepsilon}{2n}(1-2\eta)$.
which we can do from samples!

---

Statistical Query Model.

Can ask for expectations of bonded
limitations $\varphi$ $(x, l(x))$ up to additive

functions of $(x', \ell(x))$ up to additive error.

E.g. $\mathbb{E}(X_i / \ell(x) = 1)$

$\mathbb{E}(X_i (1-X_j))$

$$\chi : X \times \{0, 1\} \longrightarrow [0,1] \;,\; \tau > 0$$

↑ example    ↑ label

SQ oracle responds with $\mathbb{E}_D(\chi) \pm \tau$.

$\chi$: poly time computable    $\tau \geq \dfrac{1}{poly}$

Many (almost all) known learning algorithms can be implemented with SQ.

e.g gradient descent.    $L(w) = \mathbb{E}_{x,y}(\ell(x,y,w))$

$$\nabla_w L(w) = \nabla_w \mathbb{E}_{x,y}(\ell(x,y,w))$$

$$= \mathbb{E}_{...} / \nabla_{.}(\ell(x,y,w))$$

$$SQ. \longrightarrow = \mathbb{E}_{x,y}\left(\nabla_W(\ell(x,y,w))\right)$$

---

**Thm.** SQ-learnable $\Rightarrow$ PAC learnable with Random Classification noise.

---

Estimate $P_S\left(\chi\left(x, f(x)\right) = 1\right)$

Let

$\underline{CLEAN} = \{x : \chi(x,0) = \chi(x,1)\}$

$\underline{NOISY} = \{x : \chi(x,0) \neq \chi(x,1)\}$

example
$\chi(x,\ell) = 1$
if $X_i = 1$ and $\ell = 0$.

$P_r\left(\chi(x,f(x)) = 1\right) = P_r\left(\chi(x,f(x)) = 1 \text{ and } x \in CLEAN\right)$
$\qquad\qquad + P_r\left(\chi(x,f(x)) = 1 \text{ and } x \in NOISY\right)$

Can. estimate $P_r(x \in CLEAN)$ — not affected by noise.
and hence first term.

Also $P_r(x \in NOISY)$

Also $\Pr(x \in NOISY)$

and $\Pr_{\eta}(\chi(x, f(x)) = 1 \mid x \in NOISY) \} q$

$= (1-\eta)p + \eta(1-p) = \eta + p(1-2\eta)$

where $p = \Pr(\chi(x, f(x)) = 1 \mid x \in NOISY)$

Hence $p = \dfrac{q - \eta}{1 - 2\eta}$.

Need to estimate $p$ to within $\pm \tau$

So need to estimate $q$ to within $\pm \tau(1-2\eta)$

---

Some functions are hard to learn with SQ.

e.g. parity!

Suppose the hypothesis class is the set of parities. How many $SQ(\tau)$ do we need?

Thm weak learning PARITY needs $2^{\Omega(n)}$ SQs.

Thm. weak learning PARITY needs 2 sys.

Pf. Recall that $\langle \chi_S, \chi_T \rangle_D = \begin{cases} 1 & S=T \\ 0 & o.w. \end{cases}$

Take an SQ $\chi(x, \ell(x))$

it is function fun $\{-1,1\}^n \times \{-1,1\}$
$$\longrightarrow [-1,1]$$

Let us extend $\{\chi_S\}$ to a basis for $\{-1,1\}^{n+1}$

$$\chi_S(x,\ell) = \chi_S(x)$$

$$h_S(x,\ell) = \ell(x) \cdot \chi_S(x)$$

These are $2^{n+1}$ functions.

check: $\langle h_S, \chi_T \rangle = 0$

$$\langle h_S, h_T \rangle = \begin{cases} 1 & S=T \\ 0 & o.w. \end{cases}$$

So any query $g$ can be written as

So any query $g$ can be written as

$$g(x, \ell(x)) = \sum_S \alpha_S h_S(x, \ell) + \underbrace{\sum_S \hat{f}_S \chi_S(x)}_{\substack{\text{independent} \\ \text{of } \ell(x)}}$$

$\therefore$

$$\mathbb{E}_D\Big(g(x), \ell(x)\Big) = \sum_S \alpha_S \, \mathbb{E}_D\big(\ell(x) \chi_S(x)\big) + g_0$$

Now $\ell(x) \doteq \chi_T$ for some $T$. So the first term is $0$ except when $S = T$, when we get

$$= \alpha_T$$

$\# T$ for which $\alpha_T \geq \tau$ is $\leq \dfrac{1}{\tau^2}$.

$\left(\text{since } \sum_S \alpha_S^2 + \sum_S \hat{f}_S^2 \leq 1\right)$

So if the SQ oracle responds "0", at most $\dfrac{1}{\tau^2}$ are eliminated.

at least $2^n \cdot \tau^2$ queries.

$\therefore$ need at least $2^v \cdot v^v$ queries.

---

More generally,

$$SQ\text{-}dim(\mathcal{H}, \gamma) = \max \{ d : \exists \ d \text{ concepts} \ h_1, h_2, \dots h_d \in \mathcal{H}$$

$$s.t. \ \|h_i\|^2 \leq 1$$

$$\langle h_i, h_j \rangle_D \leq \gamma \}.$$

Thm. To learn $\mathcal{H}$ in the SQ model, we need $\Omega(d\gamma)$ queries to $SQ(\sqrt{\gamma})$.