# Active Statistical Query Learning

Maria-Florina Balcan
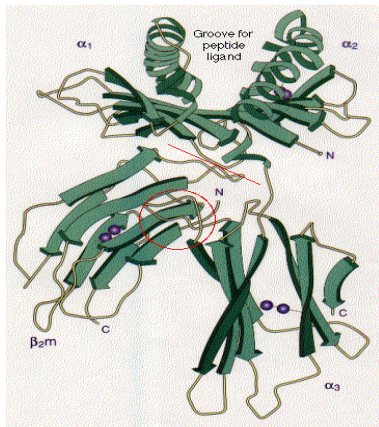
Vitaly Feldman, IBM Research
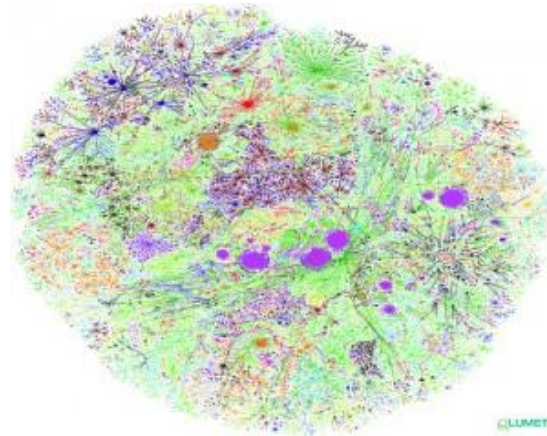
# 2-Minute Version

Modern applications: massive amounts of raw data.

Only a tiny fraction can be annotated by human experts.



Protein sequences



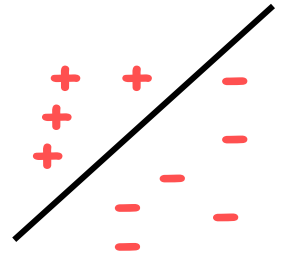Billions of webpages



Images

Active learning: leverage available data, minimize need for expert intervention.

# 2-Minute Version
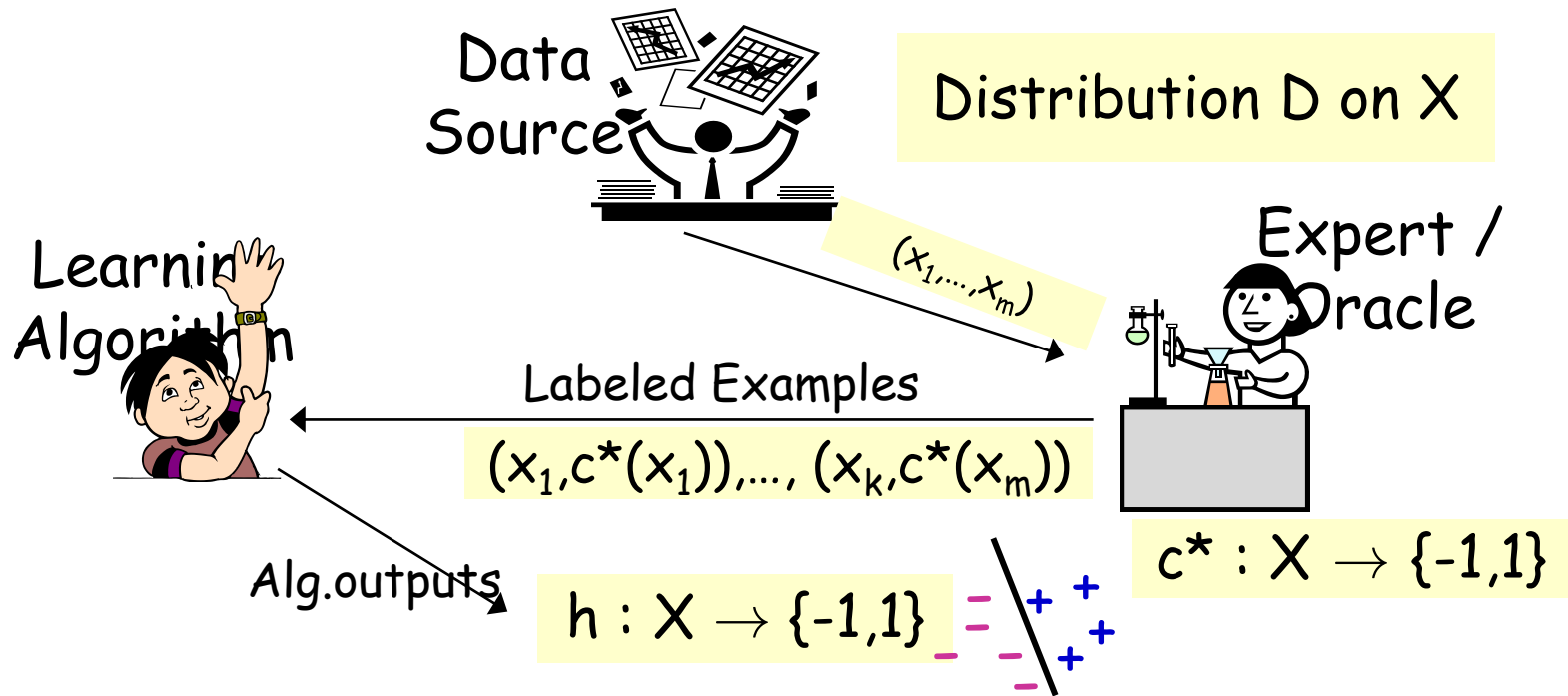
Model for designing Statistical Active Learning (AL) algos

- Poly time statistical AL algos → poly time algos tolerant to random classification noise.

- thresholds, rectangles, lin. separators.

- Naturally lead to differentially private AL algorithms.

# Outline of the talk

- Passive Learning. Statistical Query Learning

- Active Learning

- Active Statistical Query Learning

# Statistical / PAC learning model

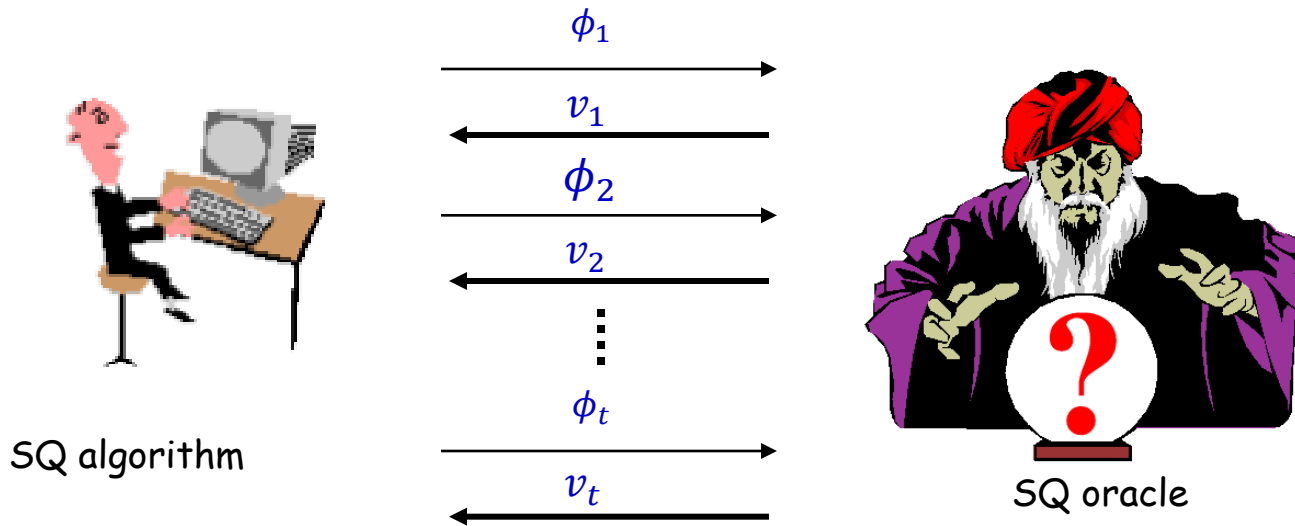

- Algo sees $(x_1, c^*(x_1)), \ldots, (x_k, c^*(x_m))$, $x_i$ i.i.d. from $D$; $c^* \in C$ [induce P]

- Do optimization over S, find hypothesis $h \in C$.

- Goal: h has small error over $D$.

$$\text{err}(h) = \Pr_{x \in D}(h(x) \neq c^*(x))$$

- PAC model: poly time algo.

# Statistical Query (SQ) Model [Kearns 93]

- Only statistical properties (not individual examples).

- Algo asks: "what is prob. a (labeled) example has property $\phi$? Pls. tell me up to additive error $\tau$."
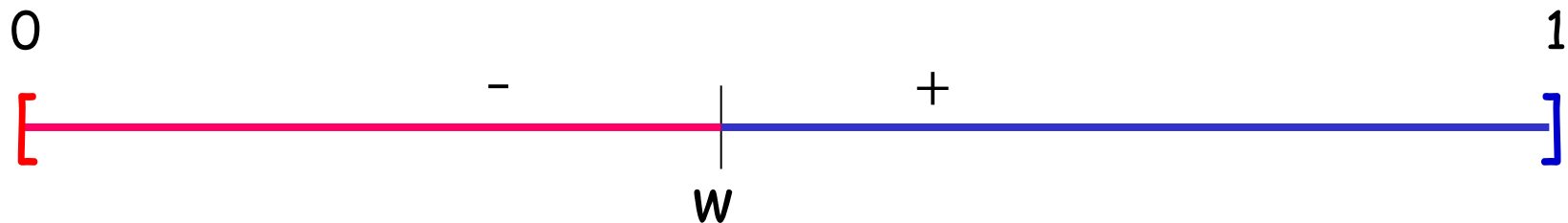


$\phi_1$

$v_1$

$\phi_2$

$v_2$

$\phi_t$

$v_t$

SQ algorithm

SQ oracle

$$|v_i - E_D[\phi_i(x, c^*(x))]| \le \tau_i$$

$\tau_i$ is tolerance of the query

Must output h of error $\le \varepsilon$.

# Simple Example: Threshold Fns

0                                                                    1

[————————————— – ——————————|——————— + ——————————]

W

If D uniform

- Ask p = Pr(x is positive) up to tolerance $\epsilon$.

  [use query $\phi(x,l) = \frac{1+l}{2}$;   $\tau = \epsilon$]
- Output 1-p.

In general,

- Ask p = Pr(x is positive) up to tolerance $\epsilon/2$.
- Ask unlabeled SQs (binary search) to find z s.t.
  
  Pr(x $\in$ (z,1]) $\in$ [p – $\epsilon/2$, p+ $\epsilon/2$].
- Output z.

# Properties of SQ model

- Can simulate SQ algos from random examples.

  [Result of query $\phi, \tau$ whp $1 - \delta$ from empirical expectation of $O\left(\frac{1}{\tau^2} \log\left(\frac{1}{\delta}\right)\right)$ random examples.]

- Can automatically convert to work in presence of random classification noise!

  - Many ML algorithms have SQ analogues.

    - E.g, Perceptron, BFKV'96,DV'06 for linear separators.

- Can be made differentially private [BDMN'05]!

# Classic Paradigm Insufficient Nowadays

Modern applications: massive amounts of raw data.

Only a tiny fraction can be annotated by human experts.

Protein sequences          Billions of webpages          Images

# Active Learning: Major Area in Modern ML

Data Source

Expert / Oracle

Learning Algorithm

Unlabeled examples

Request for the Label of an Example

A Label for that Example

Request for the Label of an Example

A Label for that Example

Algorithm outputs a classifier

- Learner can choose specific examples to be labeled.

- Goal: use fewer labeled examples.

  - Need to pick informative examples to be labeled.

# Provable Guarantees, Active Learning

- Canonical theoretical example [CAL92, Dasgupta04]

$$- \qquad\qquad\qquad\qquad +$$

W

## Active Algorithm

- Sample with $1/\varepsilon$ unlabeled examples; do binary search.

$$-\quad- \qquad\qquad +$$

Passive supervised: $\Omega(1/\varepsilon)$ labels to find an $\varepsilon$-accurate threshold.

Active: only $O(\log 1/\varepsilon)$ labels.  Exponential improvement.

# Lots of exciting activity in recent years

- Very general "disagreement based" algos [query pts from region of disagreement, throw out hyp. when statistically confi    ey are suboptimal]

  - First analyzed in [Balcan, Beyg            ford'06].

  - [Hanneke07, Dasgu            eoni'07, Wang'09, Fridman'09, Koltchinskii10            BeygelzimerHsuLangfordZhang'10, Hsu'10           ...]

**Computationally inefficient**

- Algos for specific (noise free) case    .g., linear separators.

  - QBC [Freund et al., '97]

  - Active Pe            gupta, Kalai, Monteleoni'05]

  - Margin Bas      AL [Balcan BroderZhang'07] [BalcanLong'13]

**Very specific**

Open: poly time, noise tolerant AL algos.

**This work**: framework for designing poly time AL algos tolerant to random classification noise that satisfy DP naturally.

# Active Statistical Query Model

Instead of access to random examples, algo only gets active estimates of statistical properties.

Query $(\chi, \phi)$, $\chi: X \rightarrow [0,1]$ filter [prob. of querying label of x]

"What is prob. a labeled example from $P_{|\chi}$ has property $\phi$?.

Pls. tell me up to additive error $\tau$ if $\mathrm{E}_D[\chi(x)] \geq \tau_f$ "

If $\mathrm{E}_D[\chi(x)] \geq \tau_f$ then $\left| v - \boldsymbol{E}_{P_{|\chi}}[\phi(x, f(x))] \right| \leq \tau$

$\tau_f$ filter tolerance; $\tau$ tolerance of query $(\chi, \phi)$

Algo gets an estimate of the prob that $\phi$ satisfied cond. on **x** satisfying $\chi$.
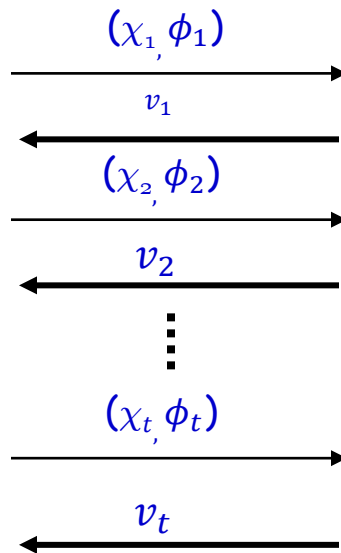
# Active Statistical Query Model

Query $(\chi, \phi)$, $\chi: X \to [0,1]$ filter [prob. of querying label of x]

"What is the prob. a labeled example from $P_{|\chi}$ has property $\phi$?.

Pls. tell me up to additive error $\tau$ if $E_D[\chi(x)] \geq \tau_f$ "



$(\chi_1, \phi_1)$

$v_1$

$(\chi_2, \phi_2)$

$v_2$

$(\chi_t, \phi_t)$

$v_t$

Active SQ algorithm

Active SQ oracle

# Simulating Active Statistical Queries

Query $(\chi, \phi)$ , $\tau_f$ , $\tau$ : if $\mathrm{E}_D[\chi(x)] \geq \tau_f$ then $\left| v - \boldsymbol{E}_{p_{|\chi}}[\phi(x, f(x))] \right| \leq \tau$

**Fact** Can be simulated with $\frac{1}{\tau^2} \log\left(\frac{1}{\delta}\right)$ labeled examples and $\frac{1}{\tau_f} \frac{1}{\tau^2} \log\left(\frac{1}{\delta}\right)$ unlabeled samples.

Design algo with $\tau$ large [only $\tau_f$ small], much less labeled data.
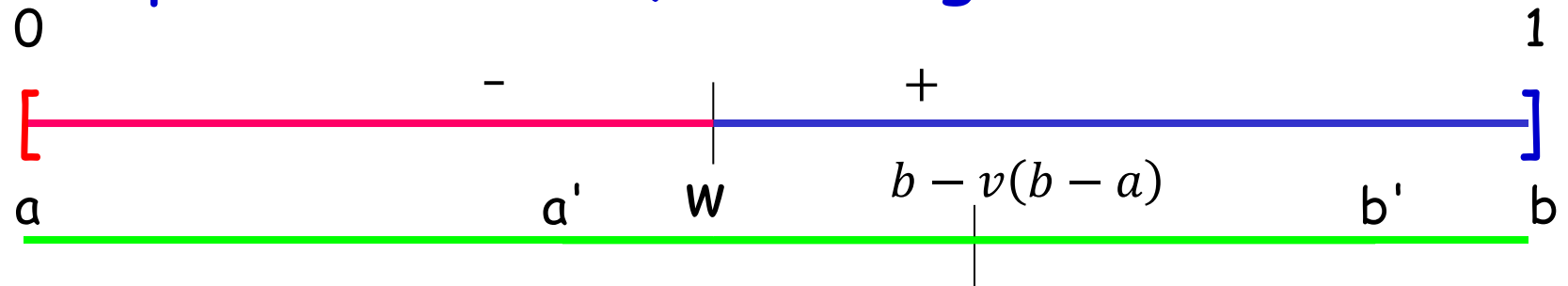
Notes:

1. Generalizes SQ model ($\chi = 1, \tau_f = 1$).

2. Since $E_{P_{|\chi}}[\phi(x, l)] = \frac{E_P[\phi(x,l)\, \chi(x)]}{E_P[\phi(x,l)]}$, can use 2 passive SQs.

   Need to estimate $E_P[\phi(x, l)\, \chi(x)]$ within $\tau E_P[\phi(x, l)]$

   Too much labeled data.

# Example: Active SQ Learning of Thresholds

0                                                                          1



$a$                           $a'$        $w$    $b - v(b-a)$              $b'$    $b$

**Passive SQ:**  Ask query $\phi(x,l) = \frac{1+l}{2}$ with tolerance $\epsilon$ ; so $1/\epsilon^2$ labels

**Active SQ:**  Key : localize/filter and use only constant tolerance

Assume $w \in [a,b]$

Ask query $\phi(x,l) = \frac{1+l}{2}$ ; $\chi(x) = I_{x \in [a,b]}$; $\tau = \frac{1}{4}, \tau_f = b - a$; get $v$.

Know $|v - E[\phi(x,l)| \, x \ \in \ [a,b]| \leq \frac{1}{4}$; $|E[\phi(x,l)| \, x \ \in \ [a,b]| = \frac{b-w}{b-a}$

So $w \in \left[ b - \left(v + \frac{1}{4}\right)(b-a), b - \left(v - \frac{1}{4}\right)(b-a) \right]$ twice smaller than $[a,b]$

Only $\log\left(\frac{1}{\epsilon}\right)$ rounds, and $\log\left(\frac{1}{\epsilon}\right)\log\left(\frac{\log\left(\frac{1}{\epsilon}\right)}{\delta}\right)$ labeled examples

# Noise Tolerance

**Fact**  Query $(\chi, \phi)$, $\tau_f$, $\tau$

Under RCN given access to $P^\eta$ estimate $E_{P_{|\chi}}[\phi(x, l)]$ within $\tau$ using $\frac{1}{\tau^2}\frac{1}{(1-2\eta)^2}\log\left(\frac{1}{\delta}\right)$ labeles and $\frac{1}{\tau_f}\frac{1}{\tau^2}\frac{1}{(1-2\eta)^2}\log\left(\frac{1}{\delta}\right)$ unlabeled examples.

[Active SQs can be simulated from examples corrupted with RCN noise.]

**Key points**:  Break into part affected by noise, and part unaffected; estimate each within $\frac{\tau}{2}$

$$\phi(x, l) = \frac{\phi(x, 1) - \phi(x, -1)}{2} l + \frac{\phi(x, 1) + \phi(x, -1)}{2}$$

$$E_{P^\eta_{|\chi}}\left[\frac{\phi(x, 1) - \phi(x, -1)}{2} l\right] = (1 - 2\eta)E_{P_{|\chi}}\left[\frac{\phi(x, 1) - \phi(x, -1)}{2} l\right]$$

sufficient to estimate it within $(1 - 2\eta)\tau/2$

# Active SQ Learning of Linear Separators

Run a passive SQ algo to get $w_0$ with err($w_0$)<C.
**iterate** k = 2, ..., s

- let $\mu_k$ be the indicator fnc of being within $\gamma_{k-1}$ of $w_{k-1}$.

- Let $\chi_k = \frac{\sum_{\{i \leq k\}} \mu_i}{k}$

- Run **passive SQ** over $D_{|\chi_k}$ to output $w_k$ of error $\frac{c}{k}$ over $D_{|\chi_k}$.

[passive SQs over $D_{|\chi_k}$ implemented as active SQs with $\tau_f = C\epsilon$]
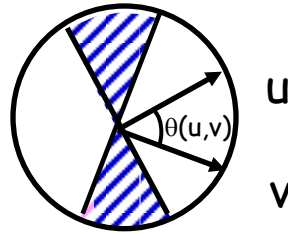
**Theorem** D log-concave in $R^d$.
If $\gamma_k = O\left(\frac{c}{2^k}\right)$ then after $s = \log\left(\frac{1}{\epsilon}\right)$ iterations err($w_s$) $\leq \epsilon$

Total number of labeled examples is poly$\left(d, \log\left(\frac{1}{\epsilon}\right)\right)$

# Linear Separators, Log-Concave Distributions

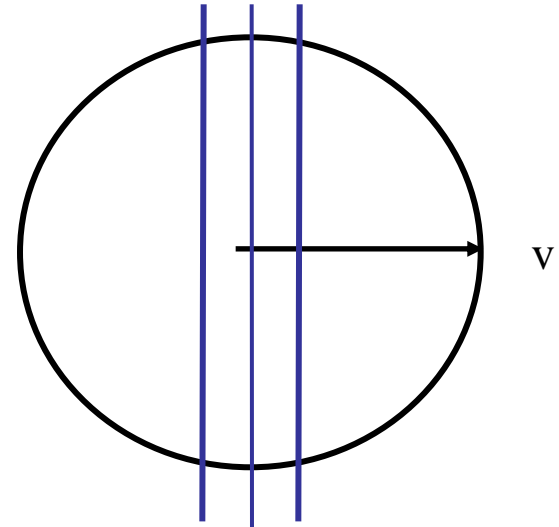Fact 1 $\quad d(u, v) \approx \dfrac{\theta(u,v)}{\pi}$


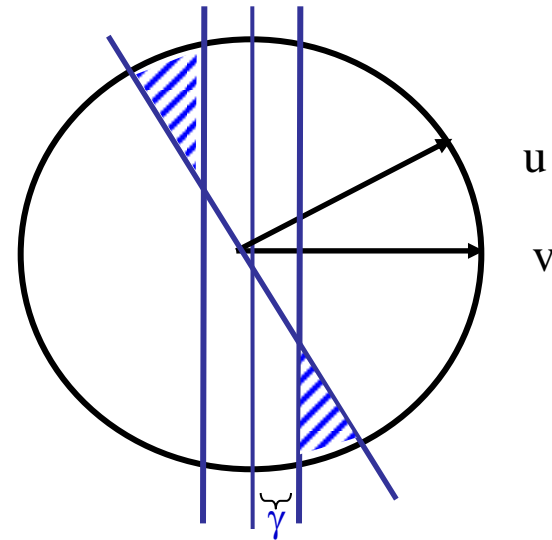
Fact 2

$$\Pr_x\left[\,|v \cdot x| \leq \gamma\,\right] \leq \gamma.$$

# Linear Separators, Log-Concave Distributions

Fact 3   If   $\theta(u,v) = \beta$ and   $\gamma = C\beta$

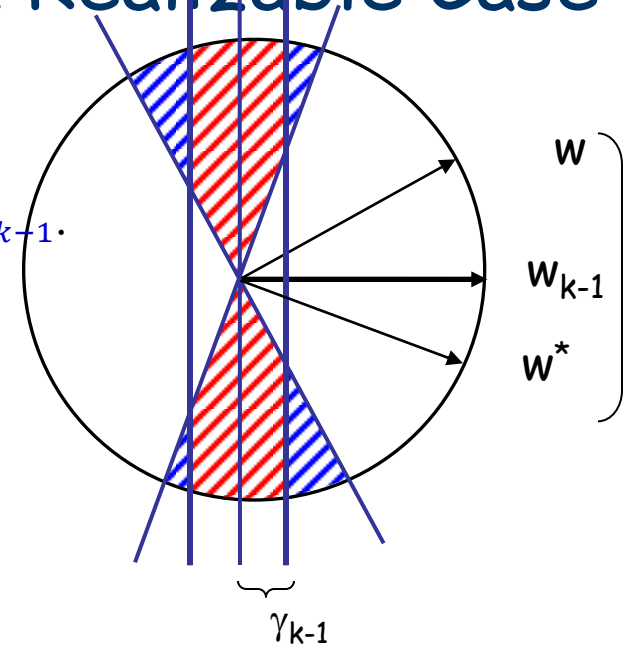$$\Pr_x\left[(u \cdot x)(v \cdot x) < 0, |v \cdot x| \geq \gamma\right] \leq \frac{\beta}{4}.$$

# Margin Based Active-Learning, Realizable Case

Run a passive SQ algo to get $w_0$ with err($w_0$)<C.

**iterate** k = 2, ..., s

- let $\mu_k$ indicator fnc of being within margin $\gamma_{k-1}$ of $w_{k-1}$.

- Let $\chi_k = \frac{\sum_{\{i \leq k\}} \mu_i}{k}$

- Run passive SQ over $D_{|\chi_k}$ with error $\frac{c}{k}$ and filter

  tolerance $C\epsilon$ to obtain $w_k$



$w$

$w_{k-1}$

$w^*$

$\gamma_{k-1}$

---

## Proof Idea

Induction: all w s.t. $\mathrm{err}_{D_{|\mu_i}}(w) \leq C, i \leq k$, $\mathrm{err}_D(w) \leq \frac{1}{2^k}$. So, $\mathrm{err}_D(w_k) \leq \frac{1}{2^k}$

Let w s. t. $\mathrm{err}_{D_{|\mu_i}}(w) \leq C$ for $i \leq k$.

$$\mathrm{err}(w) = \Pr(w \text{ errs on } x, |w_{k-1} \cdot x| \geq \gamma_{k-1}) \quad +$$

$$\Pr(w \text{ errs on } x, |w_{k-1} \cdot x| \leq \gamma_{k-1})$$

# Proof Idea

By induction $\mathrm{err}_D(w) \leq \frac{1}{2^{k-1}}$, so $\theta(w, w^*) \leq 2^{-k+1}$

Also $\theta(w_{k-1}, w^*) \leq 2^{-k+1}$

For $\gamma_k = O\left(\frac{c}{2^k}\right)$

$$\mathrm{err}(w) = \underbrace{\Pr(w \text{ errs on } x, |w_{k-1} \cdot x| \geq \gamma_{k-1})}_{\leq 1/2^{k+1}} +$$

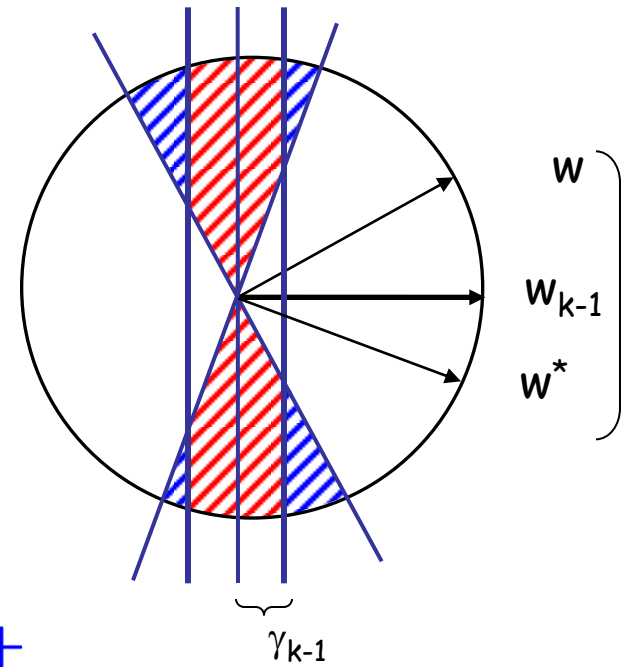$$\Pr(w \text{ errs on } x, |w_{k-1} \cdot x| \leq \gamma_{k-1})$$

# Proof Idea



By induction $\mathrm{err}_D(w) \leq \frac{1}{2^{k-1}}$, so $\theta(w, w^*) \leq 2^{-k+1}$

Also $\theta(w_{k-1}, w^*) \leq 2^{-k+1}$

For $\gamma_k = O\left(\frac{c}{2^k}\right)$

$$\leq 1/2^{k+1}$$

$$\mathrm{err}(w) = \underbrace{\Pr(w \text{ errs on } x, |w_{k-1} \cdot x| \geq \gamma_{k-1})}_{} +$$

$$\Pr(w \text{ errs on } x \mid |w_{k-1} \cdot x| \leq \gamma_{k-1}) \Pr(|w_{k-1} \cdot x| \leq \gamma_{k-1})$$
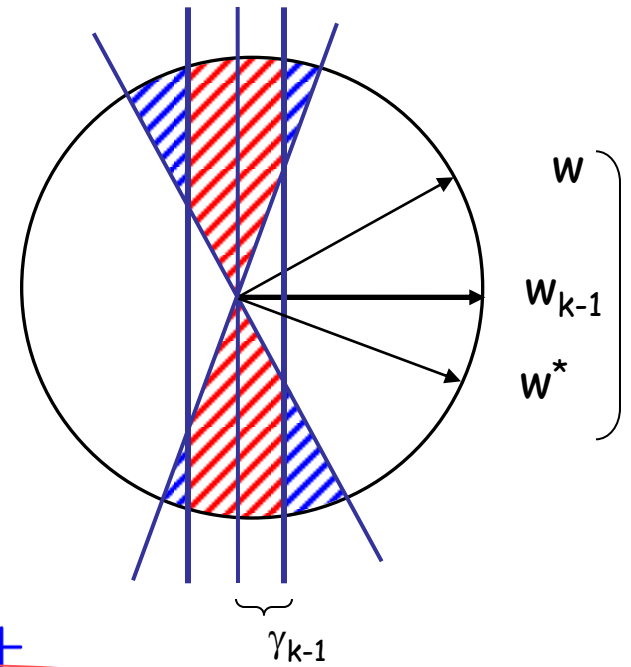
$$\leq C\gamma_{k-1}.$$

By assumption          $\leq C$

$$\leq 1/2^{k+1}$$

So $\mathrm{err}_D(w) \leq \frac{1}{2^k}$, as desired.

# Active SQ Learning of Linear Separators

**Theorem** $D$ log-concave in $R^d$.

If $\gamma_k = O\left(\frac{c}{2^k}\right)$ then after $s = \log\left(\frac{1}{\epsilon}\right)$ iterations $\operatorname{err}(w_s) \leq \epsilon$

Total number of labeled examples $\operatorname{poly}\left(d, \log\left(\frac{1}{\epsilon}\right)\right)$

Label complexity:

- Round $k$, **passive SQ** over $D_{|X_k}$, get $\operatorname{err}_{D_{|X_k}}(w_k) \leq \epsilon' = \frac{c}{k}$.
- Only $\operatorname{poly}(d, \epsilon')$ passive SQs over $D_{|X_k}$ with $\tau = 1/\operatorname{poly}(d, \epsilon')$
   [can be implemented as active SQs with $\tau = 1/\operatorname{poly}(d, \epsilon'), \tau_f = C\epsilon$].

# Active Differential Privacy

Learner has full access to unlabeled portion of database $S$.

For every element of $S$ can request the label.

**Goal**: 1. Do learning while minimize # label request

2. Ensure differential privacy [modifying a record in $S$ does not affect much prob. that any $h$ is output]

| | |
|---|---|
| $x_1$ | |
| ... | |
| $x_i$ | |
| ... | |
| $x_n$ | |

# Active Differential Privacy

- A is $\alpha$-*differentially private* if for any two neighbor datasets $S$, $S'$ (differ in just one element $(x_i, y_i) \rightarrow (x_i', y_i')$).

| $x_1$ | |
|---|---|
| ... | |
| $x_i$ | $y_i$ |
| ... | |
| $x_n$ | |

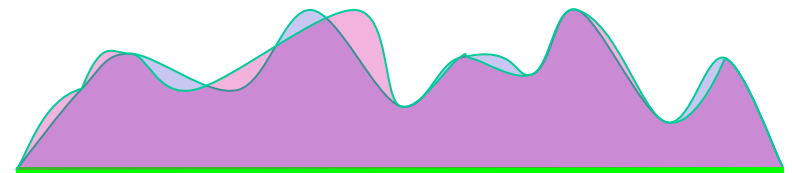| $x_1$ | |
|---|---|
| ... | |
| $x_i'$ | $y_i'$ |
| ... | |
| $x_n$ | |

For all outcomes $v$,

$$\mathrm{e}^{-\alpha} \leq \frac{\Pr(A(S) = v)}{\Pr(A(S') = v)} \leq e^{\alpha}$$

$\approx 1 + \alpha$

$\approx 1 - \alpha$

Prob. over randomness in A

# Active Differential Privacy

**Theorem** Any active SQ alg with **M** queries of tolerance $\tau$, filter-tolerance $\tau_f$, can be made to preserve $\alpha$-Diff Privacy using $O(\left(\frac{M}{\alpha\tau} + \frac{M}{\tau^2}\right)\log(M))$ label requests, $O(\left(\frac{M}{\alpha\tau\tau_f} + \frac{M}{\tau^2\tau_f}\right)\log(M))$ unlabeled examples.

Privacy cost     original             Privacy cost    original

**Implications:**

- For $\alpha \geq \tau$, privacy "for free" in terms of # of labeled requests.

- For lin. sep. & thresholds, can learn and preserve DP with much fewer label requests than **non**-private passive as long as $\alpha$ is large compared to $\epsilon$.

# Active Differential Privacy

**Theorem** Any active SQ alg with **M** queries of tolerance $\tau$, filter-tolerance $\tau_f$, can be made to preserve $\alpha$-Diff Privacy using $O(\left(\frac{M}{\alpha\tau} + \frac{M}{\tau^2}\right)\log(M))$ label requests, $O(\left(\frac{M}{\alpha\tau\tau_f} + \frac{M}{\tau^2\tau_f}\right)\log(M))$ unlabeled examples.

Privacy cost    original

Privacy cost    original

**Proof sketch:**

- Answer each query using disjoint set of $O(\left(\frac{1}{\alpha\tau\tau_f} + \frac{1}{\tau^2\tau_f}\right)\log(M))$ unlabeled exs.

- Will query the **T** examples that pass the filter and add Laplace noise.

- Sets are disjoint so suffices to satisfy $\alpha$-DP per query.

- Query sensitivity is $\frac{1}{T}$, so suffices to add $\frac{1}{\alpha T}$ Laplace noise per query.

- Sample size large enough so that whp, noise added is $\leq \tau/2$. Combine with $\tau/2$ from sample size to get whp overall error $\leq \tau$ per query.

# Discussion

Model for designing Statistical Active Learning (AL) algos

- Poly time statistical AL algos $\rightarrow$ poly time algos tolerant to random classification noise.

- Naturally lead to differentially private AL algorithms.

**Open questions**

Deal with more general types of noise [ABL'13].

Practical Implications?