Active Learning

Maria-Florina Balcan 16/11/2015

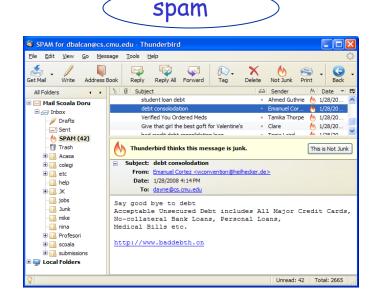
Supervised Learning

• E.g., which emails are spam and which are important.

Not spam athesis - Mozilla Thunderbird File Edit View Go Message Tools Help 🙈 Get Mail 🔹 🣝 Write 🌉 Address Book 🏽 🕙 Tag 🖜 X → Teligil Subject

Subject

Teligil From Date Congrats on the dissertation award! Michelle Leah Goodstein 10/12/2009 9:56 ... SCS Dissertation Award & ACM Dissertatio... Randy Bryant 10/14/2009 10:0... reimbursemrent · 10/15/2009 9:34 ... Re: Congrats on the dissertation award! Maria Florina Balcan reneated-eq. Re: SCS Dissertation Award & ACM Dissert Doru-Cristian Balcan 10/15/2009 12:4 review seminars-gatech sindofrii archive junk Xdelete students 10/14/2009 10:00 PM subject SCS Dissertation Award & ACM Dissertation Nominees students-diverse atalks-accross-gatech cc Catherine Copetas <copetas@cs.cmu.edu> 🗅 talks-campus atalks-gatech Nina: talks-outside teaching You might have already seen this announcement, but I would like to tech-report personally congratulate you for your outstanding dissertation. I would like to invite you to return to CMU to give a distinguished lecture theory-group sometime in the winter of 2010. Catherine Copetas will work out the theory-talks timing for you. You'll get to use the new Rashid Auditorium---a big improvement over Wean 7500. Tong total-diverse-gatech Best of wishes to you at Georgia Tech upcoming-trips Randy Downloading 26 of 29 in thesi Unread: 0 Total: 29



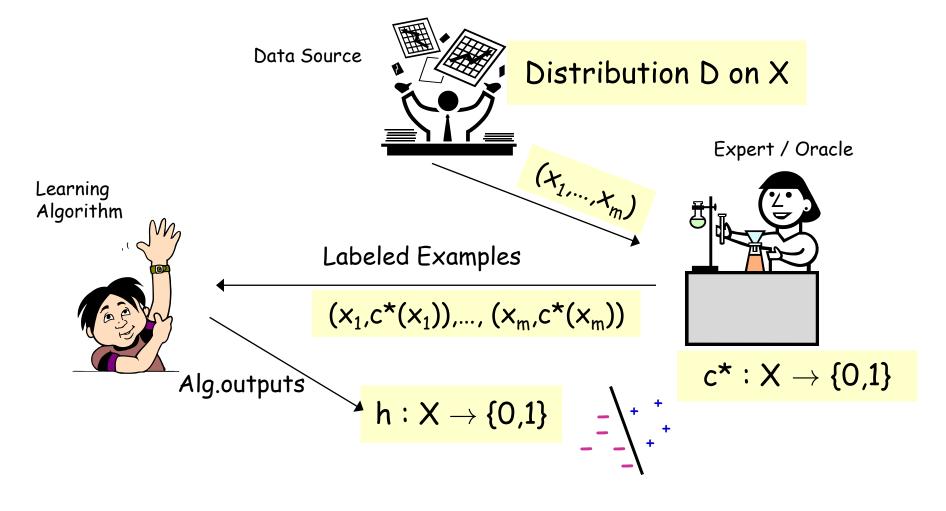
E.g., classify objects as chairs vs non chairs.



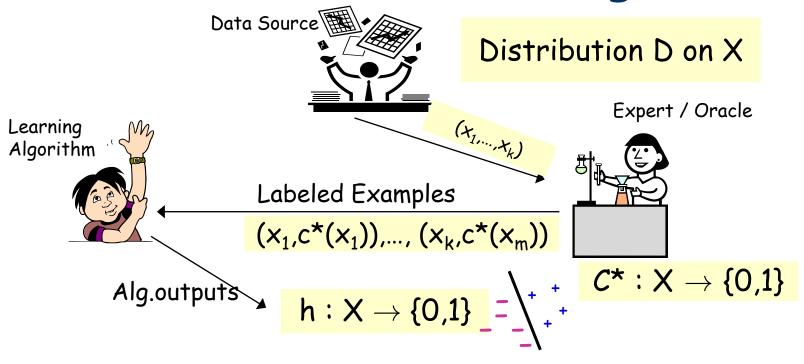


chair

Statistical / PAC learning model



Statistical / PAC learning model



- Algo sees $(x_1,c^*(x_1)),...,(x_k,c^*(x_m)), x_i$ i.i.d. from D
- Do optimization over S, find hypothesis $h \in C$.
- Goal: h has small error over D.

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

· c* in C, realizable case; else agnostic

Two Main Aspects in Classic Machine Learning

Algorithm Design. How to optimize?

Automatically generate rules that do well on observed data.

E.g., Boosting, SVM, etc.

Generalization Guarantees, Sample Complexity

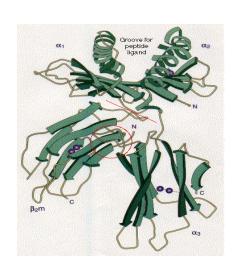
Confidence for rule effectiveness on future data.

- Realizable: $0\left(\frac{1}{\epsilon}\left(VC\dim(C)\log\left(\frac{1}{\epsilon}\right) + \log\left(\frac{1}{\delta}\right)\right)\right)$
- Agnostic replace ε with ε^2 .

Classic Fully Supervised Learning 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

Modern ML: New Learning Approaches

Modern applications: massive amounts of raw data.

Techniques that best utilize data, minimizing need for expert/human intervention.

Paradigms where there has been great progress.

· Semi-supervised Learning, (Inter)active Learning.







Active Learning

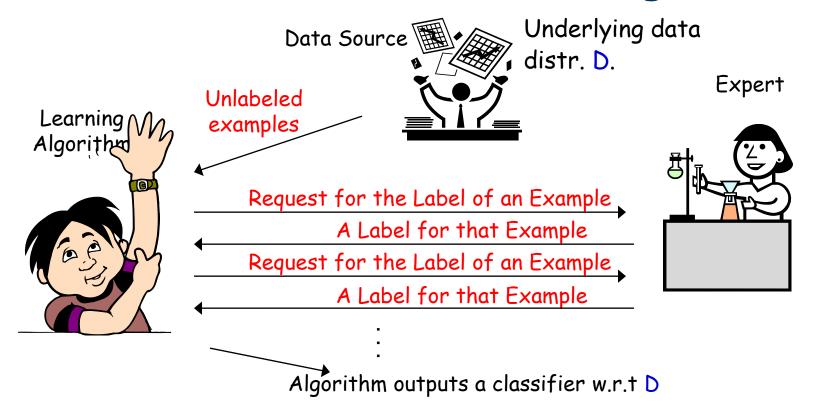
Additional resources:

- Two faces of active learning. Sanjoy Dasgupta. 2011.
- Active Learning. Balcan-Urner. Encyclopedia of Algorithms. 2015
- Theory of Active Learning. Hanneke. 2014

Additional resources:

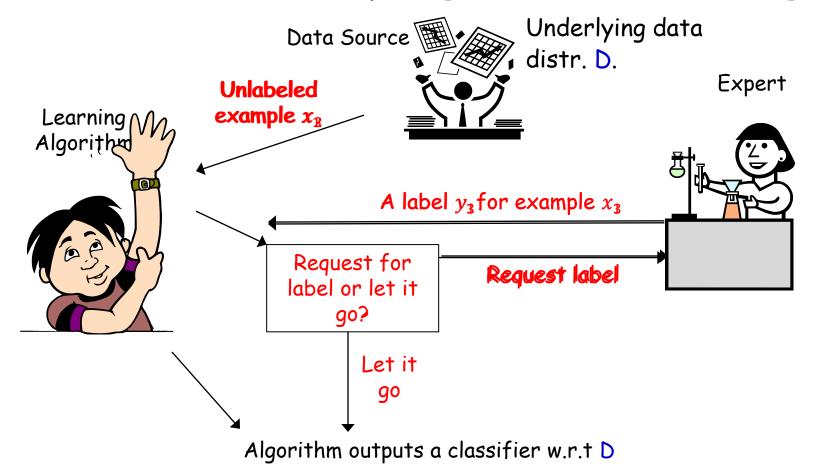
Active Learning. Bur Settles. 2012.

Batch Active Learning



- Learner can choose specific examples to be labeled.
- Goal: use fewer labeled examples [pick informative examples to be labeled].

Selective Sampling Active Learning



- Selective sampling AL (Online AL): stream of unlabeled examples, when each arrives make a decision to ask for label or not.
- · Goal: use fewer labeled examples [pick informative examples to be labeled].

What Makes a Good Active Learning Algorithm?

- Guaranteed to output a relatively good classifier for most learning problems.
- Doesn't make too many label requests.

Hopefully a lot less than passive learning and SSL.

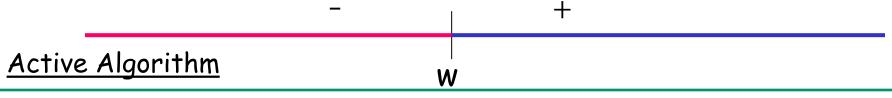
 Need to choose the label requests carefully, to get informative labels.

Can adaptive querying really do better than passive/random sampling?

- YES! (sometimes)
- We often need far fewer labels for active learning than for passive.
- This is predicted by theory and has been observed in practice.

Can adaptive querying help? [CAL92, Dasgupta04]

• Threshold fns on the real line: $h_w(x) = 1(x \ge w)$, $C = \{h_w : w \in R\}$



- Get N unlabeled examples
- How can we recover the correct labels with $\ll N$ queries?
- Do binary search! Just need O(log N) labels!



- Output a classifier consistent with the N inferred labels.
- $N = O(1/\epsilon)$ we are guaranteed to get a classifier of error $\leq \epsilon$.

<u>Passive supervised</u>: $\Omega(1/\epsilon)$ labels to find an ϵ -accurate threshold.

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

Uncertainty sampling in SVMs common and quite useful in practice. E.g., [Tong & Koller, ICML 2000; Jain, Vijayanarasimhan & Grauman, NIPS 2010; Schohon Cohn, ICML 2000]

Active SVM Algorithm

- At any time during the alg., we have a "current guess" \mathbf{w}_t of the separator: the max-margin separator of all labeled points so far.
- Request the label of the example closest to the current separator.

Active SVM seems to be quite useful in practice.

[Tong & Koller, ICML 2000; Jain, Vijayanarasimhan & Grauman, NIPS 2010]

Algorithm (batch version)

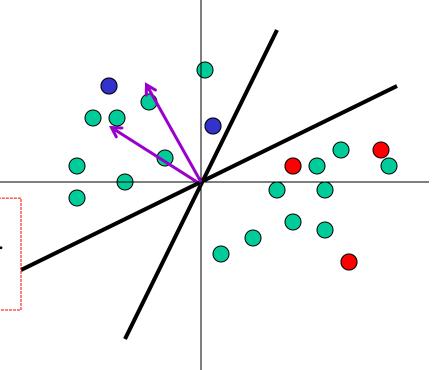
Input $S_u = \{x_1, ..., x_{m_u}\}$ drawn i.i.d from the underlying source D

Start: query for the labels of a few random x_i s.

For $t = 1, \ldots,$

- Find w_t the max-margin separator of all labeled points so far.
- Request the label of the example closest to the current separator: minimizing $|x_i \cdot w_t|$.

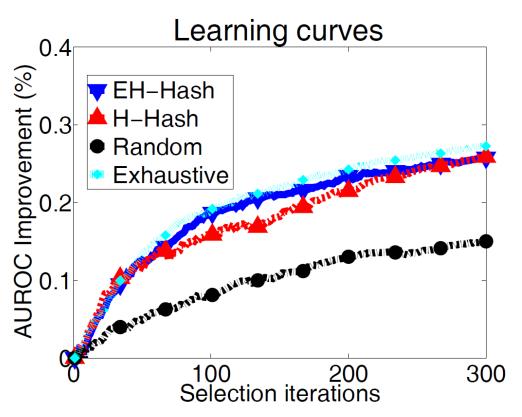
(highest uncertainty)



Active SVM seems to be quite useful in practice.

E.g., Jain, Vijayanarasimhan & Grauman, NIPS 2010

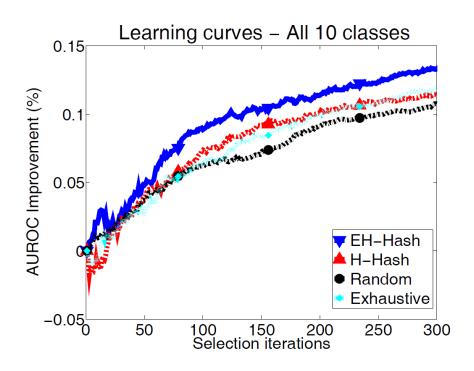
Newsgroups dataset (20.000 documents from 20 categories)



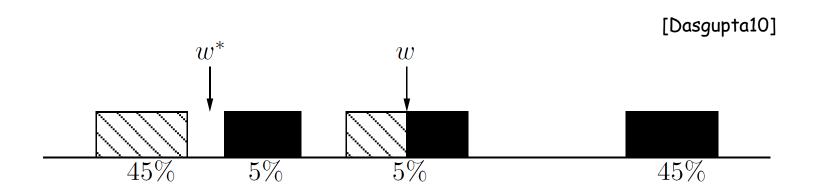
Active SVM seems to be quite useful in practice.

E.g., Jain, Vijayanarasimhan & Grauman, NIPS 2010

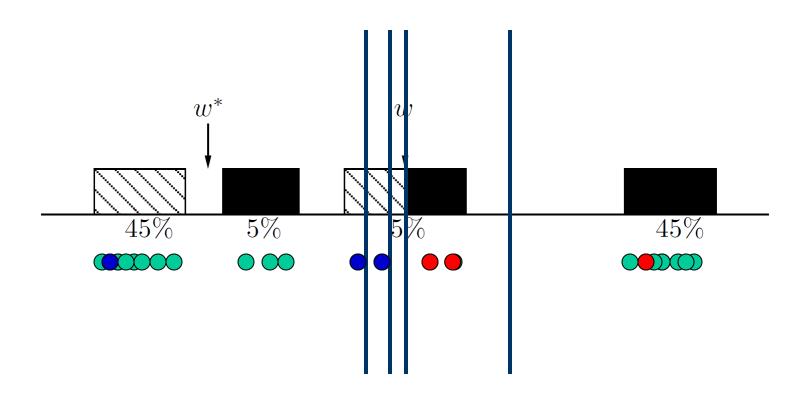
CIFAR-10 image dataset (60.000 images from 10 categories)



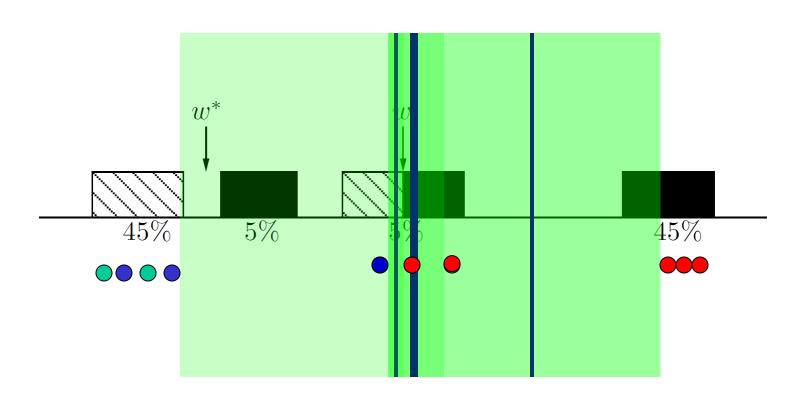
- Works sometimes....
- However, we need to be very very careful!!!
 - Myopic, greedy technique can suffer from sampling bias.
 - A bias created because of the querying strategy; as time goes on the sample is less and less representative of the true data source.



- Works sometimes....
- However, we need to be very very careful!!!



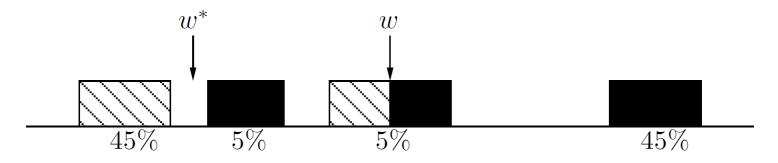
- Works sometimes....
- However, we need to be very very careful!!!



- Works sometimes....
- However, we need to be very very careful!!!
 - Myopic, greedy technique can suffer from sampling bias.
 - Bias created because of the querying strategy; as time goes on the sample is less and less representative of the true source.
 - Observed in practice too!!!!



 Main tension: want to choose informative points, but also want to guarantee that the classifier we output does well on true random examples from the underlying distribution.



Safe Active Learning Schemes

Disagreement Based Active Learning Hypothesis Space Search

[CAL92] [BBL06]

[Hanneke'07, DHM'07, Wang'09, Fridman'09, Kolt10, BHW'08, BHLZ'10, H'10, Ailon'12, ...]

Version Spaces

- X feature/instance space; distr. D over X; c^* target fnc
- Fix hypothesis space H.

```
\label{eq:Definition} \begin{array}{ll} \textbf{Definition} & \text{Assume realizable case: } c^* \in H. \\ \textbf{Given a set of labeled examples } (x_1,y_1), ..., (x_{m_l},y_{m_l}), y_i = c^*(x_i) \\ \textbf{Version space of H: part of H consistent with labels so far.} \\ \textbf{I.e., } h \in VS(H) \text{ iff } h(x_i) = c^*(x_i) \ \forall i \in \{1, ..., m_l\}. \end{array}
```

Version Spaces

- X feature/instance space; distr. D over X; c^* target fnc
- Fix hypothesis space H.

Definition Assume realizable case: $c^* \in H$.

Given a set of labeled examples (x_1, y_1) , ..., (x_{m_1}, y_{m_1}) , $y_i = c^*(x_i)$

Version space of H: part of H consistent with labels so far.

E.g.,: data lies on circle in R², H = homogeneous linear seps.

region of disagreement in data space

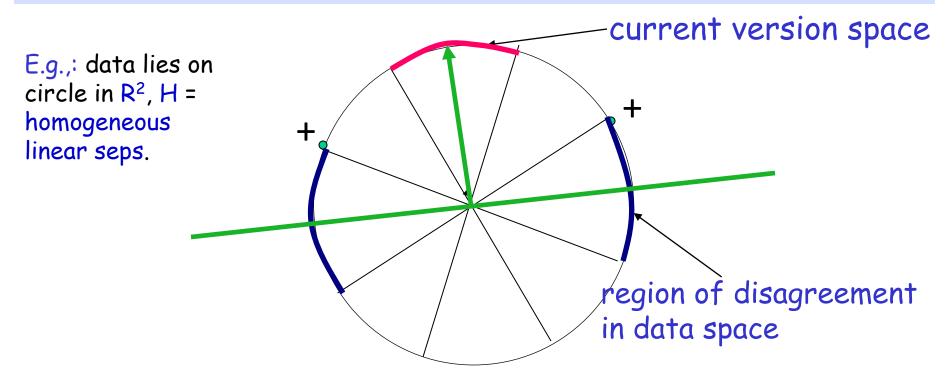
Version Spaces. Region of Disagreement

Definition (CAL'92)

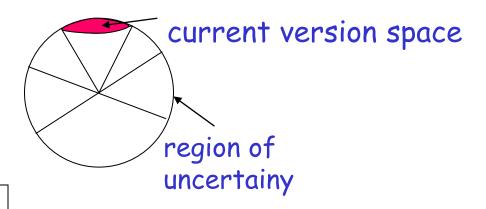
Version space: part of H consistent with labels so far.

Region of disagreement = part of data space about which there is still some uncertainty (i.e. disagreement within version space)

 $x \in X, x \in DIS(VS(H))$ iff $\exists h_1, h_2 \in VS(H), h_1(x) \neq h_2(x)$



Disagreement Based Active Learning [CAL92]



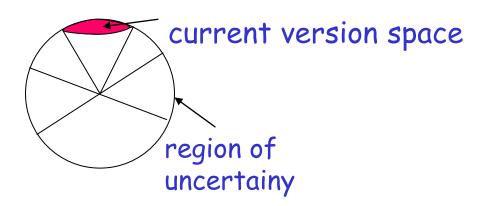
Algorithm:

Pick a few points at random from the current region of uncertainty and query their labels.

Stop when region of uncertainty is small.

Note: it is active since we do not waste labels by querying in regions of space we are certain about the labels.

Disagreement Based Active Learning [CAL92]



Algorithm:

Query for the labels of a few random x_i s.

Let H_1 be the current version space.

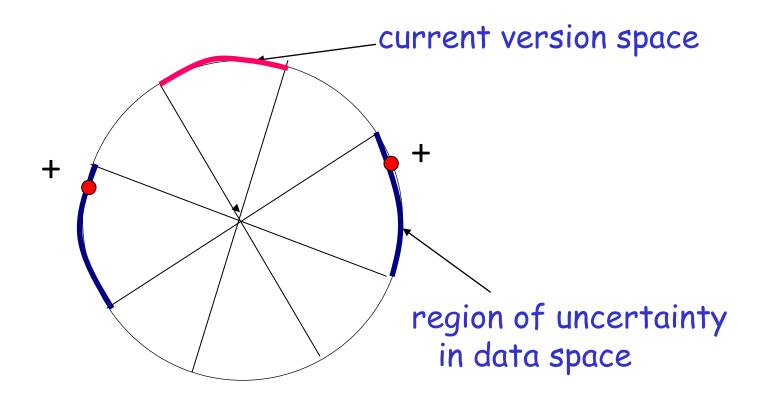
For $t = 1, \ldots,$

Pick a few points at random from the current region of disagreement $DIS(H_t)$ and query their labels.

Let H_{t+1} be the new version space.

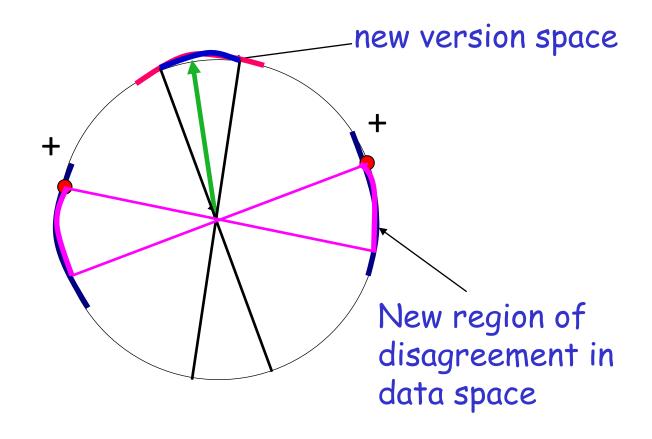
Region of uncertainty [CAL92]

- Current version space: part of C consistent with labels so far.
- "Region of uncertainty" = part of data space about which there is still some uncertainty (i.e. disagreement within version space)



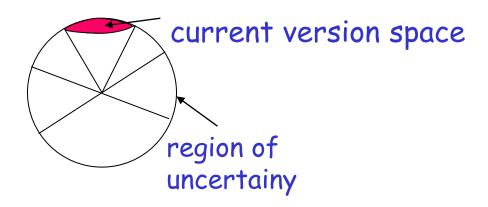
Region of uncertainty [CAL92]

- Current version space: part of C consistent with labels so far.
- "Region of uncertainty" = part of data space about which there is still some uncertainty (i.e. disagreement within version space)



How about the agnostic case where the target might not belong the H?

A² Agnostic Active Learner [BBL'06]



Algorithm:

Let
$$H_1 = H$$
.

For
$$t = 1,,$$

- Pick a few points at random from the current region of disagreement $DIS(H_t)$ and query their labels.
- Throw out hypothesis if you are statistically confident they are suboptimal.

When Active Learning Helps. Agnostic case

 A^2 the first algorithm which is robust to noise.

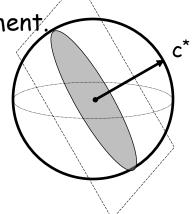
[Balcan, Beygelzimer, Langford, ICML'06] [Balcan, Beygelzimer, Langford, JCSS'08]

"Region of disagreement" style: Pick a few points at random from the current region of disagreement, query their labels, throw out hypothesis if you are statistically confident they are suboptimal.

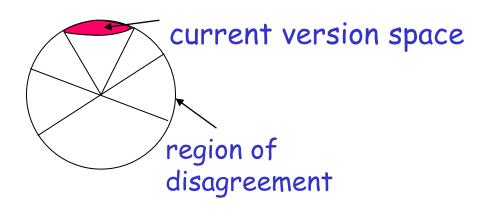
Guarantees for A² [BBL'06,'08]:

- Fall-back & exponential improvements.
 - C thresholds, low noise, exponential improvement;
 - · C homogeneous linear separators in Rd,
 - D uniform, low noise, only $d^2 \log (1/\epsilon)$ labels.

A lot of subsequent work.



A² Agnostic Active Learner [BBL'06]



Algorithm:

Let $H_1 = H$.

Careful use of generalization bounds; Avoid the selection bias!!!! /

For t = 1,,

- Pick a few points at random from the current region of disagreement $DIS(H_t)$ and query their labels.
- Throw out hypothesis if you are statistically confident they are suboptimal.

General guarantees for A² Agnostic Active Learner

"Disagreement based": Pick a few points at random from the current region of uncertainty, query their labels, throw out hypothesis if you are statistically confident they are suboptimal. [BBL'06]
How quickly the region of disagreement

collapses as we get closer and closer to optimal classifier

Guarantees for A² [Hanneke'07]:

Disagreement coefficient
$$\theta_{c^*} = \sup_{r \geq \eta + \epsilon} \frac{\Pr(DIS(B(c^*, r)))}{r}$$

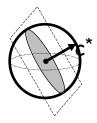
Theorem

$$m = \left(1 + \frac{\eta^2}{\epsilon^2}\right) VCdim(C)\theta_{c^*}^2 \log(\frac{1}{\epsilon})$$

labels are sufficient s.t. with prob. $\geq 1-\delta$ output h with $err(h) \leq \eta + \epsilon$.

Realizable case: $m = VCdim(C)\theta_{c^*}\log(\frac{1}{\epsilon})$

Linear Separators, uniform distr.: $\theta_{c^*} = \sqrt{d}$



Disagreement Based Active Learning

"Disagreement based" algos: query points from current region of disagreement, throw out hypotheses when statistically confident they are suboptimal.

- Generic (any class), adversarial label noise.
- Computationally efficient for classes of small VC-dimension

Still, could be suboptimal in label complex & computationally inefficient in general.

Lots of subsequent work trying to make is more efficient computationally and more aggressive too: [HannekeO7, DasguptaHsuMontleoni'07, Wang'09, Fridman'09, Koltchinskii10, BHW'08, BeygelzimerHsuLangfordZhang'10, Hsu'10, Ailon'12, ...]