

Active Learning and Optimized Information Gathering

Lecture 8 – Active Learning

CS 101.2

Andreas Krause

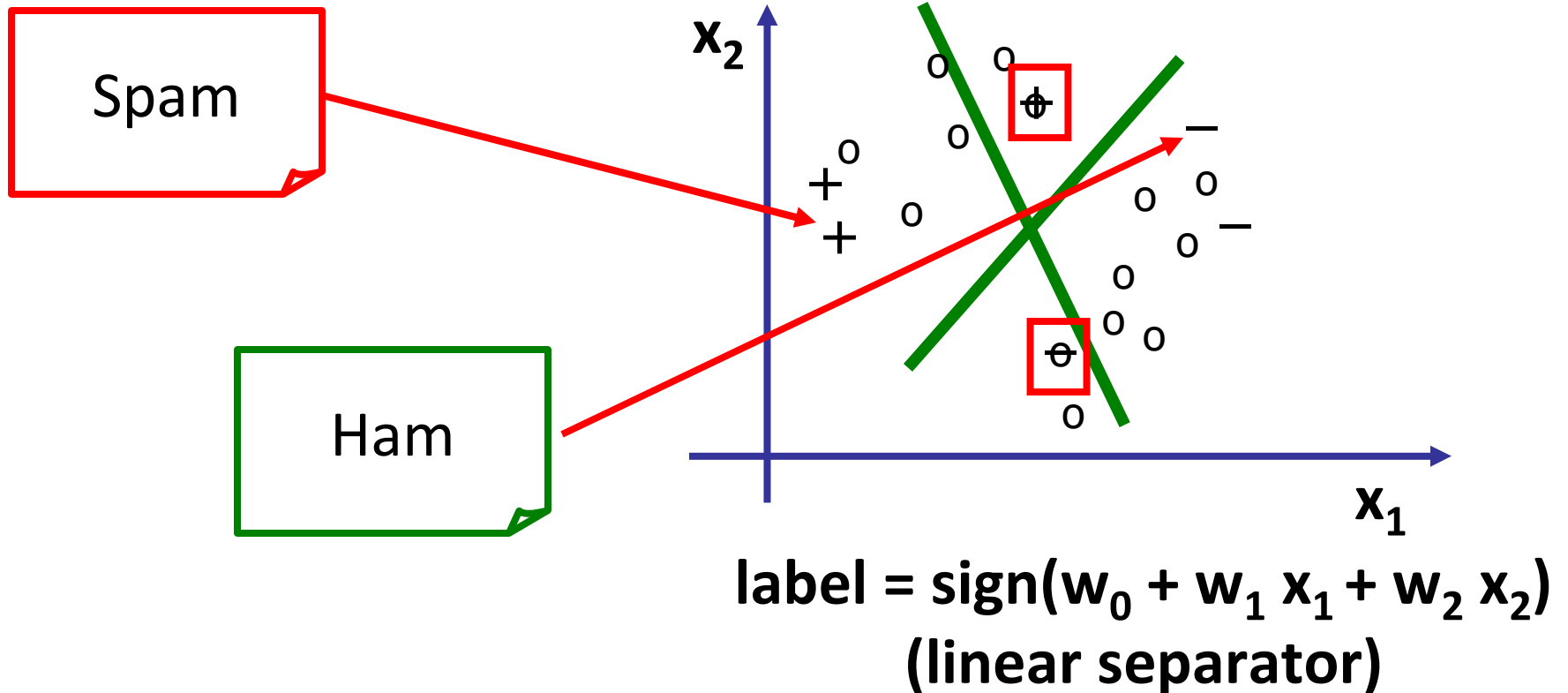
Announcements

- Homework 1: Due today
- Office hours
 - Come to office hours before your presentation!
 - Andreas: **Monday 3pm-4:30pm**, 260 Jorgensen
 - Ryan: Wednesday 4:00-6:00pm, 109 Moore

Outline

- Background in learning theory
- Sample complexity
- Key challenges
- Heuristics for active learning
- Principled algorithms for active learning

Spam or Ham?



- Labels are expensive (need to ask expert)
- **Which labels should we obtain to maximize classification accuracy?**

Recap: Concept learning

- Set X of instances, with distribution P_X
- True concept $c: X \rightarrow \{0,1\}$
- Data set $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$, $x_i \sim P_X$, $y_i = c(x_i)$
- Hypothesis $h: X \rightarrow \{0,1\}$ from $H = \{h_1, \dots, h_n, \dots\}$
- Assume $c \in H$ (c also called “target hypothesis”)

- $\text{error}_{\text{true}}(h) = E_X |c(x) - h(x)|$
- $\text{error}_{\text{train}}(h) = (1/n) \sum_i |c(x_i) - h(x_i)|$
 y_i

If n large enough, $\text{error}_{\text{true}}(h) \approx \text{error}_{\text{train}}(h)$ for all h

Recap: PAC Bounds

How many samples n do we need to get error $\leq \epsilon$ with probability $1-\delta$?

No noise: $n \geq \frac{1}{\epsilon} (\log |H| + \log \frac{1}{\delta})$

Noise: $n \geq \frac{1}{\epsilon^2} (\log |H| + \log \frac{1}{\delta})$

Requires that data is i.i.d.!

Today: Mainly no-noise case (more next week)

Statistical passive/active learning protocol

Data source P_x (produces inputs x_i)



Active learner assembles
data set $D_n = \{(x_1, y_1), \dots, (x_n, y_n)\}$
by selectively obtaining labels



Learner outputs hypothesis h



$$\text{error}_{\text{true}}(h) = E_{x \sim P}[h(x) \neq c(x)]$$

Data set NOT sampled i.i.d.!!

Example: Uncertainty sampling

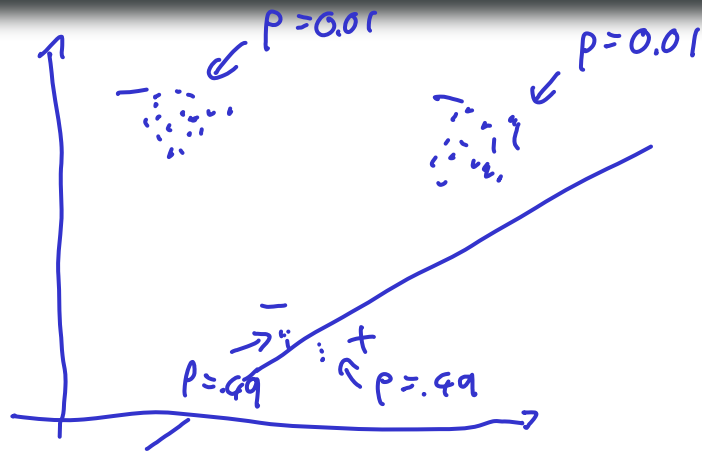
- Budget of m labels
- Draw n unlabeled examples $m = \frac{n}{2}$
- Repeat until we've picked m labels
 - Assign each unlabeled data an “uncertainty score”
 - Greedily pick the most uncertain example
- One of the most commonly used class of heuristics!

Uncertainty sampling for linear separators



$$u(x) = \frac{1}{\text{distance to separator}}$$

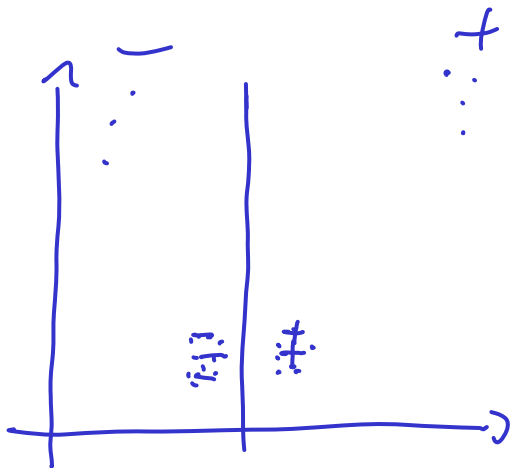
Active learning bias



Uncertainty sampling
with $m = \frac{n}{2}$

Even if
 $n \rightarrow \infty$
 $m = \frac{n}{2}$


\exists hypothesis h
consistent with labels
we've seen
 $\text{error}(h) \geq 0.01$



Active learning bias

- If we can pick at most $m = n/2$ labels, with overwhelmingly high probability, US pick points such that there remains a hypothesis with error $> \epsilon$!!!
- With standard passive learning, error $\rightarrow 0$ as $n \rightarrow \infty$

Wish list for active learning

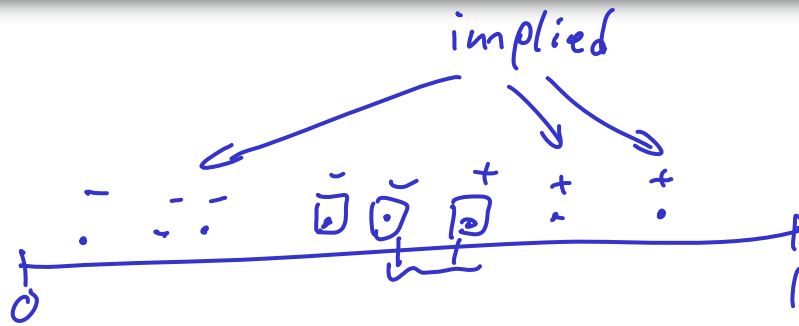
- Minimum requirement
 - **Consistency:** Generalization error should go to 0 asymptotically
- We'd like more than that: 
 - **Fallback guarantee:** Convergence rate of error of active learning “at least as good” as passive learning
- What we're **really** after
 - **Rate improvement:** Error of active learning decreases much faster than for passive learning

From passive to active

- Passive PAC learning
 1. Collect data set D of $n \geq \frac{1}{\epsilon} (\log |H| + \log 1/\delta)$ data points and their labels i.i.d. from P_X
 2. Output consistent hypothesis h
 3. With probability at least $1-\delta$, $\text{error}_{\text{true}}(h) \leq \epsilon$

- Key idea
 - Sample **n unlabeled** data points $D_X = \{x_1, \dots, x_n\}$ i.i.d.
 - Actively query labels until all hypotheses consistent with these labels **agree** on the labels of **all unlabeled** data

Why might this work?



Formalization: “Relevant” hypothesis

- Data set $D = \{(x_1, y_1), \dots, (x_n, y_n)\}$, Hypothesis space H
- Input data: $D_x = \{x_1, \dots, x_n\}$
- Relevant hypothesis
 $H'(D_x) = H' = \text{Restriction of } H \text{ on } D_x$
- Formally:
 $H' = \{h' : D_x \rightarrow \{0, 1\} \mid \exists h \in H \text{ s.t. } \forall x \in D_x: h'(x) = h(x)\}$

Example: Threshold functions



$$H = \{h(x) = [x \geq t] \text{ for some } t \in [0, 1]\}$$

$$h(x) = \begin{cases} 1 & \text{if } x \geq t \\ 0 & \text{otherwise} \end{cases}$$

$$H' = \{h_1, h_2, h_3\}$$

	x_1	x_2
h_1	-	-
h_2	-	+
h_3	+	+

Version space

- Input data $D_x = \{x_1, \dots, x_n\}$
- Partially labeled: Have $L = \{(x_{i_1}, y_{i_1}), \dots, (x_{i_m}, y_{i_m})\}$
- The (relevant) **version space** is the set of all **relevant hypotheses consistent** with the labels L
- Formally:

$$V(H', L) = V(L) = V = \{h \in H' : \forall (x, y) \in L : h(x) = y\}$$

- Why useful?
 - Partial labels L imply all remaining labels for $D_x \Leftrightarrow |V| = 1$

Version space

- Input data $D_X = \{x_1, \dots, x_n\}$
- Partially labeled: Have $L = \{(x_{i_1}, y_{i_1}), \dots, (x_{i_m}, y_{i_m})\}$
- The (relevant) **version space** is the set of all **relevant hypotheses consistent** with the labels L

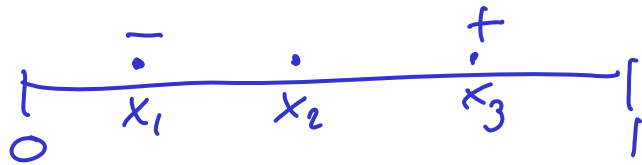
- Formally:

$$V(D_X, L) = V = \{h' \in H'(D_X) : h'(x_{i_j}) = y_{i_j} \text{ for } 1 \leq j \leq m\}$$

- Why useful?

- Partial labels L imply all remaining labels for $D_X \Leftrightarrow |V|=1$

Example: Binary thresholds



$$V = \{h_1, h_2\}$$


	x_1	x_2	x_3
h_1	-	-	+
h_2	-	+	+

Pool-based active learning with fallback

1. Collect $n \geq \frac{1}{\epsilon} (\log |H| + \log 1/\delta)$ **unlabeled** data points D_X from P_X
2. Actively request labels L until there remains a single hypothesis $h' \in H'$ that's consistent with these labels (i.e., $|V(H', L)| = 1$)
3. Output any hypothesis $h \in H$ consistent with the obtained labels. With probability $\geq 1 - \delta$ $\text{error}_{\text{true}}(h) \leq \epsilon$

Get PAC guarantees for active learning
Bounds on #labels for fixed error ϵ
carry over from passive to active
→ **Fallback guarantee**

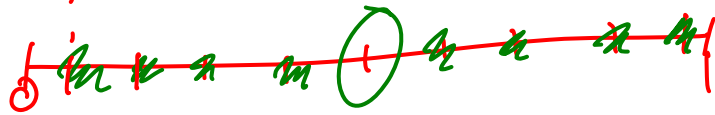
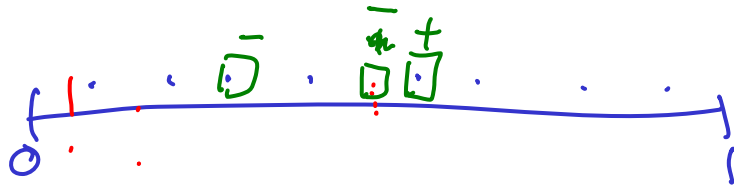
Wish list for active learning

- Minimum requirement
 - **Consistency:** Generalization error should go to 0 asymptotically
- We'd like more than that:
 - **Fallback guarantee:** Convergence rate of error of active learning “at least as good” as passive learning 
- What we're **really** after
 - **Rate improvement:** Error of active learning decreases much faster than for passive learning

Pool-based active learning with fallback

1. Collect $n \geq \frac{1}{\epsilon} (\log |H| + \log \frac{1}{\delta})$ **unlabeled** data points D_X from P_X
2. Actively request labels L until there remains a single hypothesis $h' \in H'$ that's consistent with these labels (i.e., $|V(H', L)| = 1$)
3. Output any hypothesis $h \in H$ consistent with the obtained labels. With probability $\geq 1 - \delta$ $\text{error}_{\text{true}}(h) \leq \epsilon$

Example: Threshold functions



↑
must be the
only consistent relevant hyp

⇒ only need $O(\log n)$ labels

(instead of n for passive)

Generalizing binary search [Dasgupta '04]

- Want to shrink the version space (number of consistent hypotheses) as quickly as possible.

• *Halve labels L*

- General (greedy) approach:

- For each unlabeled instance x_i compute

$$v_{i,1} = V(H', L \cup \{(x_i, 1)\}) \leftarrow \text{"hallucinate" label 1}$$

$$v_{i,0} = V(H', L \cup \{(x_i, 0)\}) \leftarrow \text{"hallucinate" label 0}$$

- $v_i = \min \{v_{i,1}, v_{i,0}\}$

- Obtain label y_i for x_i where $i = \operatorname{argmax}_j \{v_j\}$

Ideal case

$$L_0 = \{\emptyset\} \subset L_1 \subset L_2 \subset \dots \subset L_m \quad |L_m| = n$$

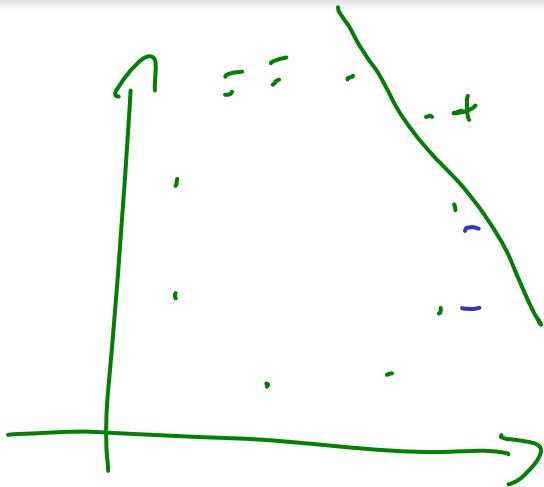
$$\text{Ideal case: } \forall i \exists x : |V(L_i \cup \{(x_i, 1)\})| = |V(L_i \cup \{(x_i, 0)\})|$$

$$\Rightarrow |V(L_{i+1})| \leq \frac{1}{2} |V(L_i)|$$

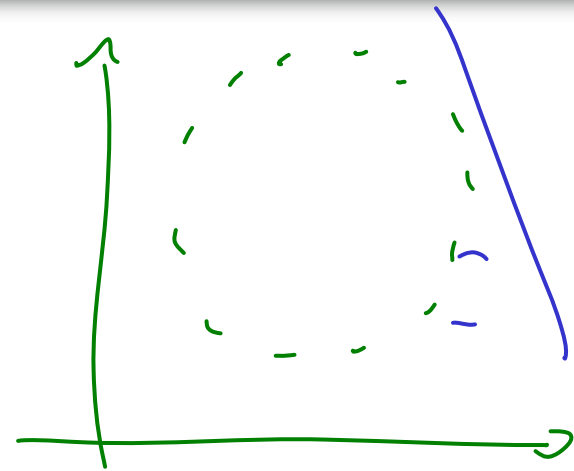
$$\Rightarrow |V(L_m)| \leq \frac{1}{2^m} |V(L_0)| \leq 1$$

$$\Rightarrow m = O(\log n)$$

Is it always possible to half the version space?

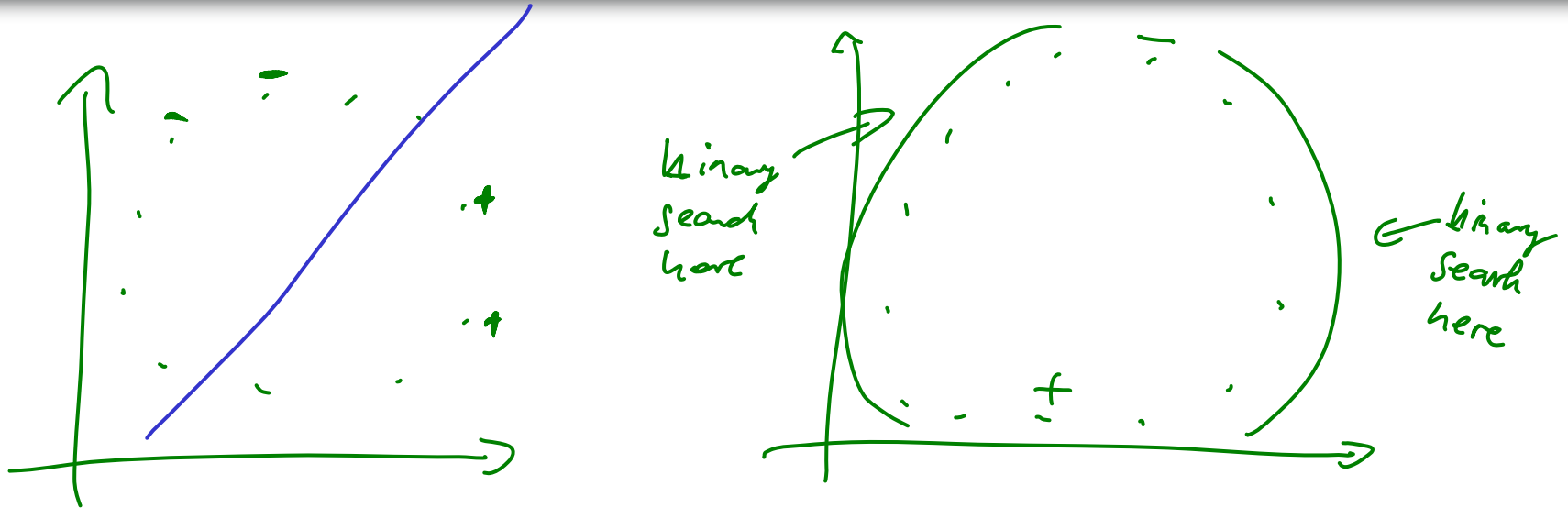


need to
distinguish
from



need to see all n data points i

Typical case much more benign



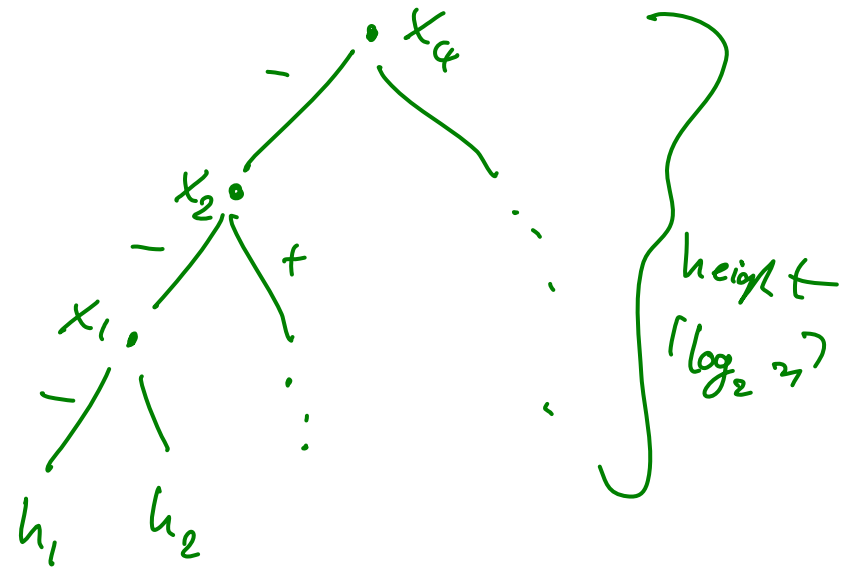
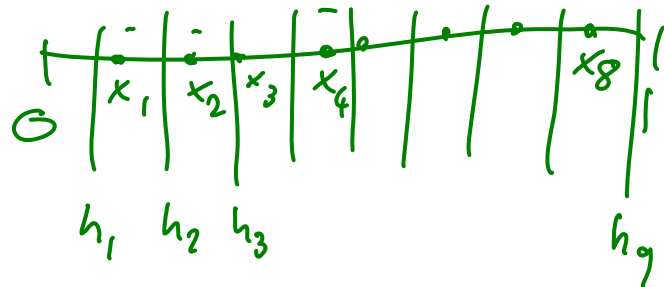
Typically, $O(\log n)$ samples suffice

Query trees

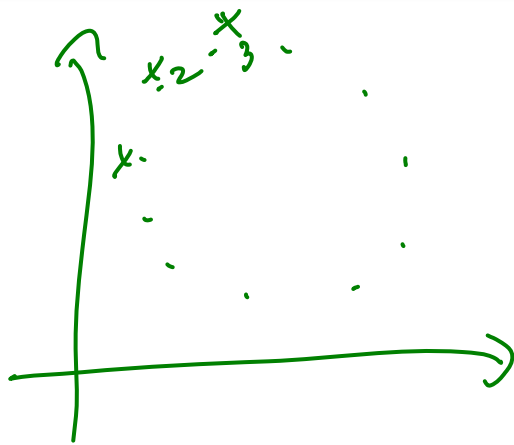
- A **query tree** is a rooted, labeled tree on the relevant hypothesis H'
- Each node is labeled with an input $x \in D_x$
- Each edge is labeled with $\{0,1\}$
- Each path from root to hypothesis $h' \in H'$ is a labeling L such that $V(D_x, L) = \{h'\}$

- **Want query trees of minimum height**

Example: Threshold functions



Example: linear separators (2D)



Number of labels needed to identify hypothesis

- Depends on target hypothesis!
 - Binary thresholds (on n inputs D_X)
 - Optimal query tree needs $O(\log n)$ labels! 😊
 - For linear separators in 2D (on n inputs D_X)
 - For some hypotheses, even optimal tree needs n labels 😞
 - On average, optimal query tree needs $O(\log n)$ labels! 😊
- ➔ Average-case analysis of active learning

Average case query tree learning

- Query tree T
- $\text{Cost}(T) = 1/|H'| \sum_{h' \in H'} \text{depth}(h', T)$

Want $T^* = \text{argmin}_T \text{Cost}(T)$

- Superexponential number of query trees ☹️
- Finding the optimal one is hard ☹️

Greedy construction of query trees [Dasgupta '04]

Algorithm GreedyTree(D_x, L)

$V' = H'(D_x, L)$

If $V' = \{h\}$

return Leaf(h)

Else

For each unlabeled instance x_i compute

$v_{i,1} = |V'(H', L \cup \{(x_i, 1)\})|$ and $v_{i,0} = |V'(H', L \cup \{(x_i, 0)\})|$

$v_i = \min \{v_{i,1}, v_{i,0}\}$

Let $i = \operatorname{argmax}_j \{v_j\}$

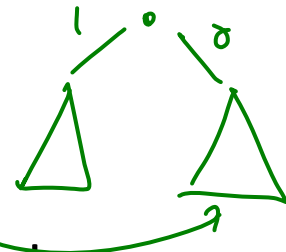
LeftSubTree = GreedyTree($D_x, L \cup \{(x_i, 1)\}$)

RightSubTree = GreedyTree($D_x, L \cup \{(x_i, 0)\}$)

return Node x_i with children LeftSubTree (1) and

RightSubTree(0)

max. "disagreement"



Near-optimality of greedy tree [Dasgupta '04]

Theorem: Let $T^* = \operatorname{argmin}_T \operatorname{Cost}(T)$

Then GreedyTree constructs a query tree T such that

$$\operatorname{Cost}(T) = O(\log |H'|) \operatorname{Cost}(T^*)$$

Limitations of this algorithm

- Often computationally intractable
 - Finding “most-disagreeing” hypothesis is difficult
- No-noise assumption
- Will see how we can relax these assumptions in the talks next week.

Bayesian or not Bayesian?

- Greedy querying needs at most $O(\log |H'|)$ queries more than optimal query tree **on average**
- Assumes prior distribution (uniform) on hypotheses
- If our assumption is wrong, generalization bound **still holds!** (but might need more labels)

Can also do a pure Bayesian analysis:

- Query by Committee algorithm [Freund et al '97]
- Assumes that Nature draws hypotheses from **known** prior distribution

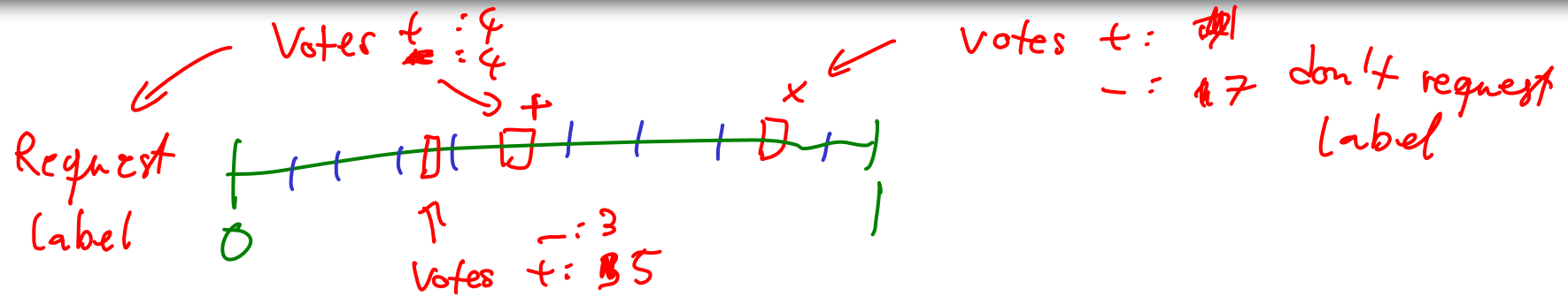
Query by Committee

- Assume prior distribution on hypotheses
- Sample a “committee” of $2k$ hypotheses drawn from the prior distribution
- Search for an input such that k “members” assign label 1, and k “members” assign 0, and query that label (“maximal disagreement”)

Theorem [Freund et al '97]

For linear separators in \mathbb{R}^d where both the coefficients w and the data X are drawn uniformly from the unit sphere, QBC requires exponentially fewer labels than passive learning to achieve same error

Example: Threshold functions



Wish list for active learning

- Minimum requirement
 - **Consistency:** Generalization error should go to 0 asymptotically
- We'd like more than that:
 - **Fallback guarantee:** Convergence rate of error of active learning "at least as good" as passive learning

What we're **really** after

- **Rate improvement:** Error of active learning decreases much faster than for passive learning

*generalized binary search
is "competitive" with optimal querying
on average if hypotheses equally likely*

*✓
using pool-based
active learning*

Beyond pool-based analysis

- Pool-based active learning just one convenient analysis technique
 - Gets around active learning bias by generalizing from a pool drawn i.i.d. at random
 - In pool-based analysis, there are examples where active learning does not outperform passive learning
- Exciting recent theoretical results show that using a more involved analysis, active learning **always** helps (asymptotically) [Balcan, Hanneke, Wortman COLT '08]
- Also other active learning paradigms
 - E.g.: Active querying (constructing rather than selecting inputs)

What you need to know

- Uncertainty sampling
- Active learning bias
- Pool-based active learning scheme
- Relevant hypotheses and version spaces
- Generalized binary search algorithm