



# Sample Complexity Bounds for Active Learning

Paper by Sanjoy Dasgupta

Presenter: Peter Sadowski



# *Passive* PAC Learning Complexity

## ■ Based on **VC dimension**

To get error  $< \epsilon$  with probability  $\geq 1 - \delta$ :

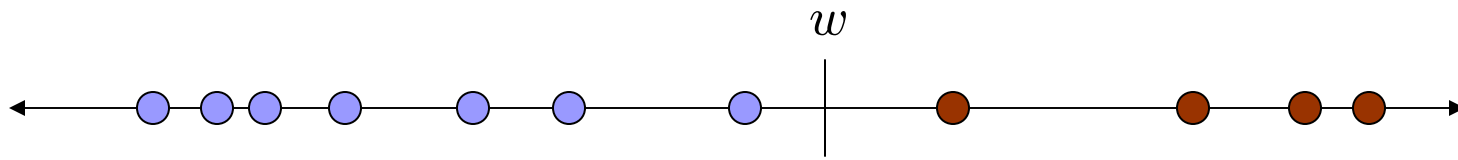
$$\text{num samples} \geq \tilde{O}\left(\frac{1}{\epsilon} (VC(H) \log(1/\delta))\right)$$

*Is there some equivalent for active learning?*

# Example: Reals in 1-D

$P$ =underlying distribution of points

$H$ =space of possible hypotheses

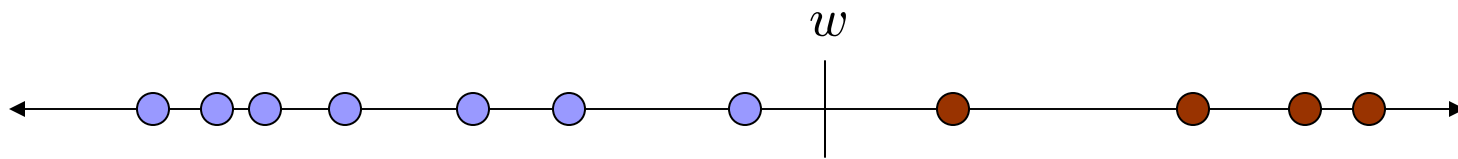


$$H = \{h_w : w \in \mathbb{R}\} \quad h_w(x) = \begin{cases} 1 & \text{if } x \geq w \\ 0 & \text{if } x < w \end{cases}$$

$O(1/\epsilon)$  random labeled examples needed from  $P$  to get error rate  $< \epsilon$

# Example: Reals in 1-D

$$h_w(x) = \begin{cases} 1 & \text{if } x \geq w \\ 0 & \text{if } x < w \end{cases}$$



Passive learning:

$O(1/\epsilon)$  random labeled examples needed from  $P$  to get error rate  $< \epsilon$

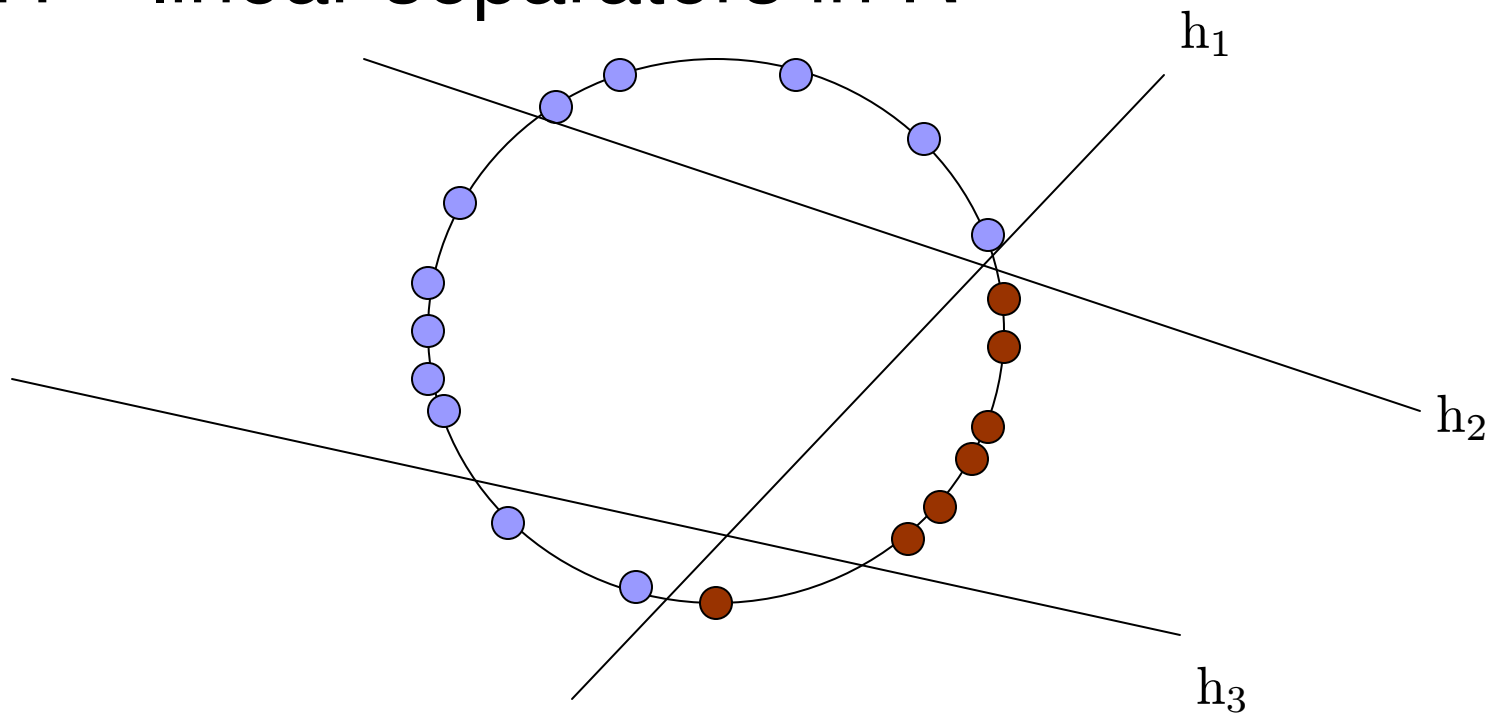
Active learning (Binary Search):

$O(\log 1/\epsilon)$  examples needed to get error  $< \epsilon$

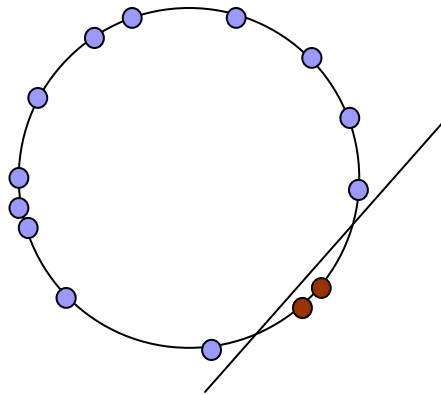
*Active learning gives us an exponential improvement!*

# Example 2: Points on a Circle

- $P$  = some density on circle perimeter
- $H$  = linear separators in  $\mathbb{R}^2$



## Example 2: Points on a Circle



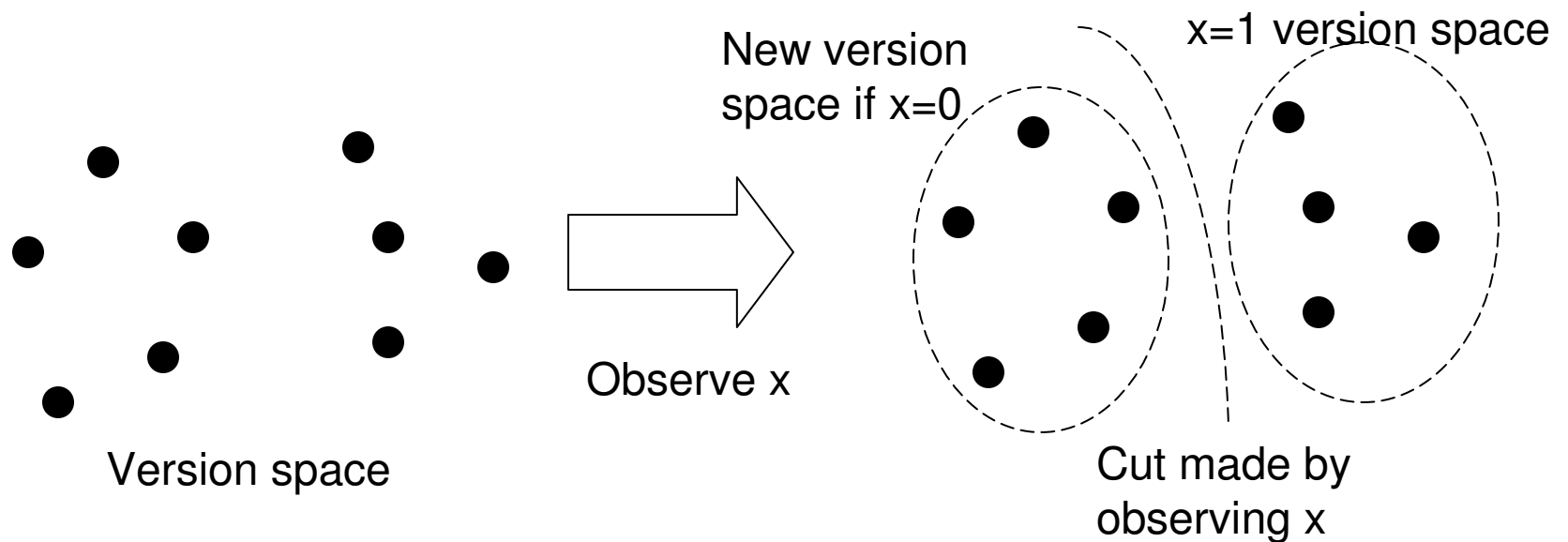
Worst case: small  $\epsilon$  slice of the circle is different

- Passive learning:  $O(1/\epsilon)$
- Active learning:  $O(1/\epsilon)$

*No improvement!*

# Active Learning Abstracted

- Goal: Narrow down the **version space**, (hypotheses that fit with known labels)
- Idea: Think of hypotheses as points



# Shrinking the Version Space

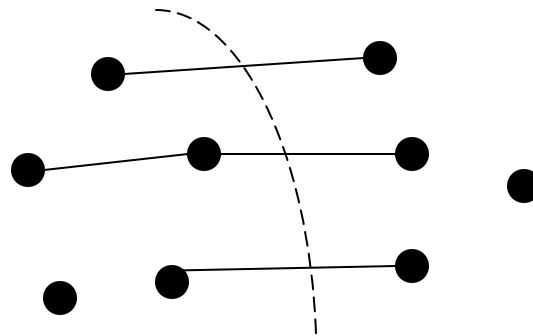
- Define distance between hypotheses:

$$d(h, h') = P\{x: h(x) \neq h'(x)\}$$

- Ignore distances less than  $\epsilon$

$$Q = H \times H$$

$$Q_\epsilon = \{(h, h') \in Q : d(h, h') > \epsilon\}$$

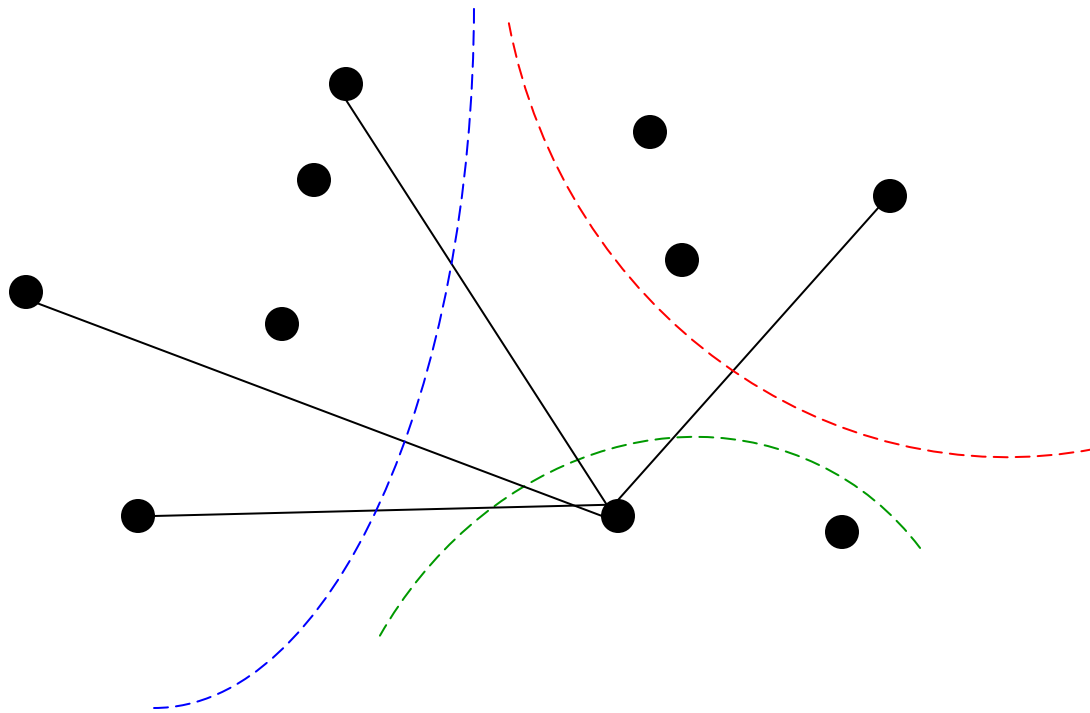


A good cut!



# Quick Example

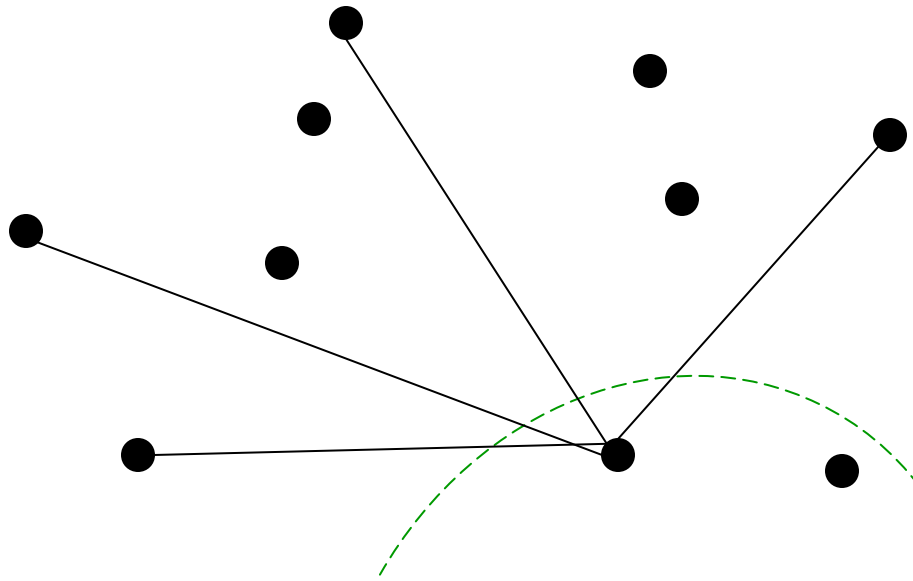
- What is the best cut?



$$Q_\epsilon = \{(h, h') \in Q : d(h, h') > \epsilon\}$$

# Quick Example

- Cut edges => shrink version space



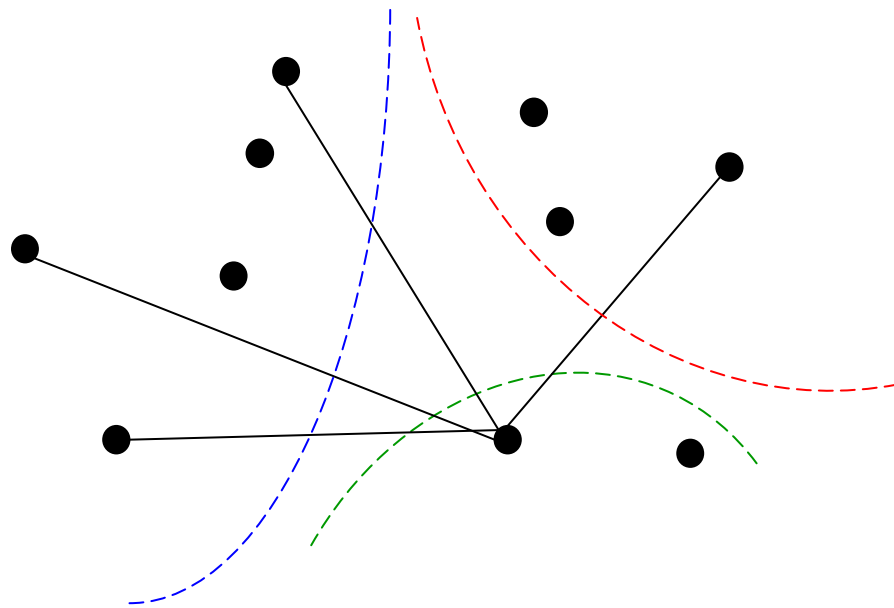
After this cut, we have a solution!

The hypotheses left are insignificantly different.

# Quantifying “Usefulness” of Points

A point  $x \in X$  is said to  $\rho$ -split  $Q_\epsilon$

IF its label reduces the number of edges by a fraction  $\rho > 0$



1/4-split

1-split

3/4-split



# Quantifying the Difficulty of Problems

Definition:

Subset  $S$  of hypotheses is  $(\rho, \epsilon, \tau)$ -splittable if

$$P\{x : x \text{ } \rho\text{-splits } Q_\epsilon\} \geq \tau$$

”At least a fraction of  $\tau$  samples are  $\rho$ -useful in splitting  $S$ .”

$\rho$  small  $\Rightarrow$  smaller splits

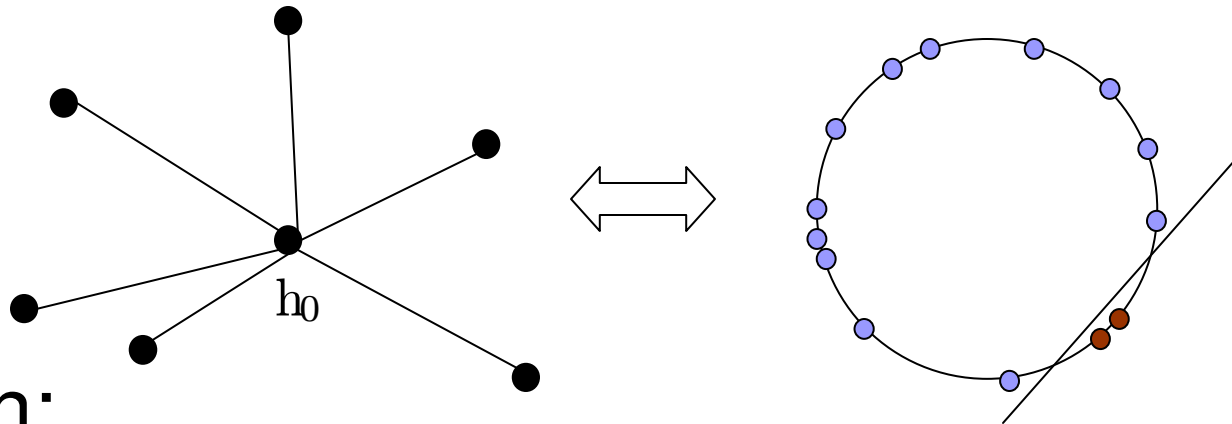
$\epsilon$  small  $\Rightarrow$  small error

$\tau$  small  $\Rightarrow$  lots of samples needed to get a good split

# Lower Bound Result

Suppose for some hypothesis space  $H$ :

- $d(h_0, h_i) > \epsilon$  for some hypotheses  $h_1, h_2, \dots, h_N$
- “disagree sets”  $\{x : h_0(x) \neq h_i(x)\}$  are disjoint



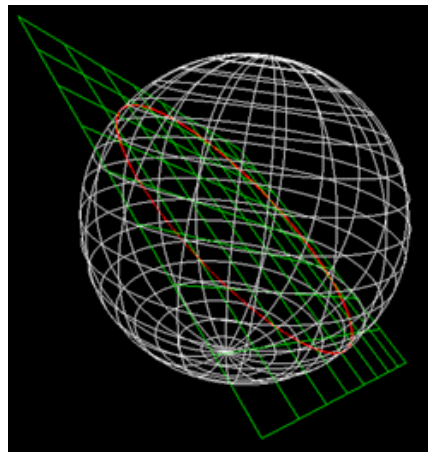
Then:

For any  $\tau$  and  $\rho > 1/N$ ,  $Q$  is not  $(\rho, \epsilon, \tau)$ -splittable.

# An Interesting Result

There is constant  $c > 0$  such that for any dimension  $d \geq 2$ , if

1.  $H$  is the class of homogeneous linear separators in  $R^d$ , and
  2.  $P$  is the uniform distribution over the surface of the unit sphere,
- then  $H$  is  $(1/4, \epsilon, c\epsilon)$ -splittable for all  $\epsilon > 0$ .



$\Rightarrow$  For any  $h \in H$ , any  $\epsilon \leq 1/(32\pi^2\sqrt{d})$ ,  
 $B(h, 4\epsilon)$  is  $\left(\frac{1}{8}, \epsilon, \Omega\left(\frac{\epsilon}{\sqrt{d}}\right)\right)$ -splittable.



# Conclusions

- Active learning not always much better than passive.
- “Splittability” is the VC dimension for active learning.
- We can use this framework to fit bounds for specific problems.