

Note to other teachers and users of these slides. Andrew would be delighted if you found this source material useful in giving your own lectures. Feel free to use these slides verbatim, or to modify them to fit your own needs. PowerPoint originals are available. If you make use of a significant portion of these slides in your own lecture, please include this message, or the following link to the source repository of Andrew's tutorials: <http://www.cs.cmu.edu/~awm/tutorials> . Comments and corrections gratefully received.

# PAC-learning

**Andrew W. Moore**  
**Associate Professor**  
**School of Computer Science**  
**Carnegie Mellon University**

[www.cs.cmu.edu/~awm](http://www.cs.cmu.edu/~awm)  
[awm@cs.cmu.edu](mailto:awm@cs.cmu.edu)  
412-268-7599

Copyright © 2001, Andrew W. Moore

Nov 30th, 2001

## Probably Approximately Correct (PAC) Learning

- Imagine we're doing classification with categorical inputs.
- All inputs and outputs are binary.
- Data is noiseless.
- There's a machine  $f(x, h)$  which has  $H$  possible settings (a.k.a. hypotheses), called  $h_1, h_2 \dots h_H$ .

Copyright © 2001, Andrew W. Moore

PAC-learning: Slide 2

## Example of a machine

- $f(x,h)$  consists of all logical sentences about  $X_1, X_2 \dots X_m$  that contain only logical ands.
- Example hypotheses:
  - $X_1 \wedge X_3 \wedge X_{19}$
  - $X_3 \wedge X_{18}$
  - $X_7$
  - $X_1 \wedge X_2 \wedge X_3 \wedge \dots \wedge X_m$
- Question: if there are 3 attributes, what is the complete set of hypotheses in  $f$ ?

Copyright © 2001, Andrew W. Moore

PAC-learning: Slide 3

## Example of a machine

- $f(x,h)$  consists of all logical sentences about  $X_1, X_2 \dots X_m$  that contain only logical ands.
- Example hypotheses:
  - $X_1 \wedge X_3 \wedge X_{19}$
  - $X_3 \wedge X_{18}$
  - $X_7$
  - $X_1 \wedge X_2 \wedge X_3 \wedge \dots \wedge X_m$
- Question: if there are 3 attributes, what is the complete set of hypotheses in  $f$ ? ( $H = 8$ )

True	$X_2$	$X_3$	$X_2 \wedge X_3$
$X_1$	$X_1 \wedge X_2$	$X_1 \wedge X_3$	$X_1 \wedge X_2 \wedge X_3$

Copyright © 2001, Andrew W. Moore

PAC-learning: Slide 4

## And-Positive-Literals Machine

- $f(x,h)$  consists of all logical sentences about  $X_1, X_2 \dots X_m$  that contain only logical ands.
- Example hypotheses:
  - $X_1 \wedge X_3 \wedge X_{19}$
  - $X_3 \wedge X_{18}$
  - $X_7$
  - $X_1 \wedge X_2 \wedge X_3 \wedge \dots \wedge X_m$
- Question: if there are  $m$  attributes, how many hypotheses in  $f$ ?

Copyright © 2001, Andrew W. Moore

PAC-learning: Slide 5

## And-Positive-Literals Machine

- $f(x,h)$  consists of all logical sentences about  $X_1, X_2 \dots X_m$  that contain only logical ands.
- Example hypotheses:
  - $X_1 \wedge X_3 \wedge X_{19}$
  - $X_3 \wedge X_{18}$
  - $X_7$
  - $X_1 \wedge X_2 \wedge X_3 \wedge \dots \wedge X_m$
- Question: if there are  $m$  attributes, how many hypotheses in  $f$ ? ( $H = 2^m$ )

Copyright © 2001, Andrew W. Moore

PAC-learning: Slide 6

## And-Literals Machine

- $f(x,h)$  consists of all logical sentences about  $X_1, X_2 \dots X_m$  or their negations that contain only logical ands.
- Example hypotheses:
  - $X_1 \wedge \sim X_3 \wedge X_{19}$
  - $X_3 \wedge \sim X_{18}$
  - $\sim X_7$
  - $X_1 \wedge X_2 \wedge \sim X_3 \wedge \dots \wedge X_m$
- Question: if there are 2 attributes, what is the complete set of hypotheses in  $f$ ?

Copyright © 2001, Andrew W. Moore

PAC-learning: Slide 7

## And-Literals Machine

- $f(x,h)$  consists of all logical sentences about  $X_1, X_2 \dots X_m$  or their negations that contain only logical ands.
- Example hypotheses:
  - $X_1 \wedge \sim X_3 \wedge X_{19}$
  - $X_3 \wedge \sim X_{18}$
  - $\sim X_7$
  - $X_1 \wedge X_2 \wedge \sim X_3 \wedge \dots \wedge X_m$
- Question: if there are 2 attributes, what is the complete set of hypotheses in  $f$ ? ( $H = 9$ )

True		True
True		$X_2$
True		$\sim X_2$
$X_1$		True
$X_1$	$\wedge$	$X_2$
$X_1$	$\wedge$	$\sim X_2$
$\sim X_1$		True
$\sim X_1$	$\wedge$	$X_2$
$\sim X_1$	$\wedge$	$\sim X_2$

Copyright © 2001, Andrew W. Moore

PAC-learning: Slide 8

## And-Literals Machine

- $f(x,h)$  consists of all logical sentences about  $X_1, X_2 \dots X_m$  or their negations that contain only logical ands.
- Example hypotheses:
  - $X_1 \wedge \sim X_3 \wedge X_{19}$
  - $X_3 \wedge \sim X_{18}$
  - $\sim X_7$
  - $X_1 \wedge X_2 \wedge \sim X_3 \wedge \dots \wedge X_m$
- Question: if there are  $m$  attributes, what is the size of the complete set of hypotheses in  $f$ ?

True		True
True		$X_2$
True		$\sim X_2$
$X_1$		True
$X_1$	$\wedge$	$X_2$
$X_1$	$\wedge$	$\sim X_2$
$\sim X_1$		True
$\sim X_1$	$\wedge$	$X_2$
$\sim X_1$	$\wedge$	$\sim X_2$

Copyright © 2001, Andrew W. Moore

PAC-learning: Slide 9

## And-Literals Machine


- $f(x,h)$  consists of all logical sentences about  $X_1, X_2 \dots X_m$  or their negations that contain only logical ands.
- Example hypotheses:
  - $X_1 \wedge \sim X_3 \wedge X_{19}$
  - $X_3 \wedge \sim X_{18}$
  - $\sim X_7$
  - $X_1 \wedge X_2 \wedge \sim X_3 \wedge \dots \wedge X_m$
- Question: if there are  $m$  attributes, what is the size of the complete set of hypotheses in  $f$ ? ( $H = 3^m$ )

True		True
True		$X_2$
True		$\sim X_2$
$X_1$		True
$X_1$	$\wedge$	$X_2$
$X_1$	$\wedge$	$\sim X_2$
$\sim X_1$		True
$\sim X_1$	$\wedge$	$X_2$
$\sim X_1$	$\wedge$	$\sim X_2$

Copyright © 2001, Andrew W. Moore

PAC-learning: Slide 10

## Lookup Table Machine


- $f(x,h)$  consists of all truth tables mapping combinations of input attributes to true and false
- Example hypothesis: 
- Question: if there are  $m$  attributes, what is the size of the complete set of hypotheses in  $f$ ?

x1	x2	x3	x4	Y
0	0	0	0	0
0	0	0	1	1
0	0	1	0	1
0	0	1	1	0
0	1	0	0	1
0	1	0	1	0
0	1	1	0	0
0	1	1	1	1
1	0	0	0	0
1	0	0	1	0
1	0	1	0	0
1	0	1	1	1
1	1	0	0	0
1	1	0	1	0
1	1	1	0	0
1	1	1	1	0

Copyright © 2001, Andrew W. Moore

PAC-learning: Slide 11

## Lookup Table Machine

- $f(x,h)$  consists of all truth tables mapping combinations of input attributes to true and false
- Example hypothesis: 
- Question: if there are  $m$  attributes, what is the size of the complete set of hypotheses in  $f$ ?

x1	x2	x3	x4	Y
0	0	0	0	0
0	0	0	1	1
0	0	1	0	1
0	0	1	1	0
0	1	0	0	1
0	1	0	1	0
0	1	1	0	0
0	1	1	1	1
1	0	0	0	0
1	0	0	1	0
1	0	1	0	0
1	0	1	1	1
1	1	0	0	0
1	1	0	1	0
1	1	1	0	0
1	1	1	1	0

$$H = 2^{2^m}$$

Copyright © 2001, Andrew W. Moore

PAC-learning: Slide 12

## A Game

- We specify  $f$ , the machine
- Nature chooses hidden random hypothesis  $h^*$
- Nature randomly generates  $R$  datapoints
  - How is a datapoint generated?
    1. Vector of inputs  $\mathbf{x}_k = (x_{k1}, x_{k2}, \dots, x_{km})$  is drawn from a fixed unknown distrib:  $D$
    2. The corresponding output  $y_k = f(\mathbf{x}_k, h^*)$
- We learn an approximation of  $h^*$  by choosing some  $h^{\text{est}}$  for which the training set error is 0

Copyright © 2001, Andrew W. Moore

PAC-learning: Slide 13

## Test Error Rate

- We specify  $f$ , the machine
- Nature chooses hidden random hypothesis  $h^*$
- Nature randomly generates  $R$  datapoints
  - How is a datapoint generated?
    1. Vector of inputs  $\mathbf{x}_k = (x_{k1}, x_{k2}, \dots, x_{km})$  is drawn from a fixed unknown distrib:  $D$
    2. The corresponding output  $y_k = f(\mathbf{x}_k, h^*)$
- We learn an approximation of  $h^*$  by choosing some  $h^{\text{est}}$  for which the training set error is 0
- For each hypothesis  $h$ ,
- Say  $h$  is Correctly Classified (CCd) if  $h$  has zero training set error
- Define  $\text{TESTERR}(h)$ 
  - = Fraction of test points that  $h$  will classify correctly
  - =  $P(h \text{ classifies a random test point correctly})$
- Say  $h$  is BAD if  $\text{TESTERR}(h) > \epsilon$

Copyright © 2001, Andrew W. Moore

PAC-learning: Slide 14

# Test Error Rate

- We specify  $f$ , the machine
- Nature chooses hidden random hypothesis  $h^*$
- Nature randomly generates  $R$  datapoints
  - How is a datapoint generated?
    1. Vector of inputs  $\mathbf{x}_k = (x_{k1}, x_{k2}, \dots, x_{km})$  is drawn from a fixed unknown distrib:  $D$
    2. The corresponding output  $y_k = f(\mathbf{x}_k, h^*)$
- We learn an approximation of  $h^*$  by choosing some  $h^{\text{est}}$  for which the training set error is 0
- For each hypothesis  $h$ ,
- Say  $h$  is Correctly Classified (CCd) if  $h$  has zero training set error
- Define  $\text{TESTERR}(h)$ 
  - = Fraction of test points that  $h$  will classify correctly
  - =  $P(h \text{ classifies a random test point correctly})$
- Say  $h$  is BAD if  $\text{TESTERR}(h) > \epsilon$

$$P(h \text{ is CCd} \mid h \text{ is bad}) = P(\forall k \in \text{Training Set}, f(x_k, h) = y_k \mid h \text{ is bad}) \leq (1 - \epsilon)^R$$

Copyright © 2001, Andrew W. Moore

PAC-learning: Slide 15

# Test Error Rate

- We specify  $f$ , the machine
- Nature chooses hidden random hypothesis  $h^*$
- Nature randomly generates  $R$  datapoints
  - How is a datapoint generated?
    1. Vector of inputs  $\mathbf{x}_k = (x_{k1}, x_{k2}, \dots, x_{km})$  is drawn from a fixed unknown distrib:  $D$
    2. The corresponding output  $y_k = f(\mathbf{x}_k, h^*)$
- We learn an approximation of  $h^*$  by choosing some  $h^{\text{est}}$  for which the **training set error is 0**
- For each hypothesis  $h$ ,
- Say  $h$  is Correctly Classified (CCd) if  $h$  has zero training set error
- Define  $\text{TESTERR}(h)$ 
  - = Fraction of test points that  $h$  will classify correctly
  - =  $P(h \text{ classifies a random test point correctly})$
- Say  $h$  is **BAD** if  $\text{TESTERR}(h) > \epsilon$

$$P(h \text{ is CCd} \mid h \text{ is bad}) = P(\forall k \in \text{Training Set}, f(x_k, h) = y_k \mid h \text{ is bad}) \leq (1 - \epsilon)^R$$

$$P(\text{we learn a bad } h) \leq$$

$$P\left(\begin{array}{c} \text{the set of CCd } h\text{'s} \\ \text{contains a bad } h \end{array}\right) =$$

$$P(\exists h. h \text{ is CCd} \wedge h \text{ is bad}) =$$

$$P\left(\begin{array}{c} (h_1 \text{ is CCd} \wedge h_1 \text{ is bad}) \vee \\ (h_2 \text{ is CCd} \wedge h_2 \text{ is bad}) \vee \\ \vdots \\ (h_H \text{ is CCd} \wedge h_H \text{ is bad}) \end{array}\right) \leq$$

$$\sum_{i=1}^H P(h_i \text{ is CCd} \wedge h_i \text{ is bad}) \leq \sum_{i=1}^H P(h_i \text{ is CCd} \mid h_i \text{ is bad}) =$$

$$H \times P(h_i \text{ is CCd} \mid h_i \text{ is bad}) \leq H(1 - \epsilon)^R$$

Copyright © 2001, Andrew W. Moore

PAC-learning: Slide 16



## PAC Learning

- Choose  $R$  such that with probability less than  $\delta$  we'll select a bad  $h^{est}$  (i.e. an  $h^{est}$  with 0 train error but which will have test error  $\varepsilon$  or worse)

$$P(\text{error of } h \text{ is } \geq \varepsilon) \leq \delta$$

- Probably Approximately Correct
- As we just saw, this can be achieved by choosing  $R$  such that

$$H(1 - \varepsilon)^R \leq \delta$$

- i.e.  $R$  such that

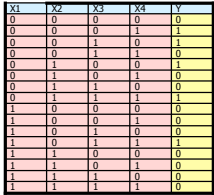
$$R \geq \frac{1}{\varepsilon} \left( \log_2 H + \log_2 \frac{1}{\delta} \right)$$

Whoah wait...  
Why is that?  
 $(1 - \varepsilon) \leq e^{-\varepsilon}$

Copyright © 2001, Andrew W. Moore

PAC-learning: Slide 17

## PAC in action

Machine	Example Hypothesis	H	R required to PAC-learn
And-positive-literals	$X3 \wedge X7 \wedge X8$	$2^m$	$\frac{1}{\varepsilon} \left( m + \log_2 \frac{1}{\delta} \right)$
And-literals	$X3 \wedge \sim X7$	$3^m$	$\frac{1}{\varepsilon} \left( (\log_2 3)m + \log_2 \frac{1}{\delta} \right)$
Lookup Table		$2^{2^m}$	$\frac{1}{\varepsilon} \left( 2^m + \log_2 \frac{1}{\delta} \right)$
And-lits or And-lits	$(X1 \wedge X5) \vee (X2 \wedge \sim X7 \wedge X8)$	$(3^m)^2 = 3^{2m}$	$\frac{1}{\varepsilon} \left( (2 \log_2 3)m + \log_2 \frac{1}{\delta} \right)$

Copyright © 2001, Andrew W. Moore

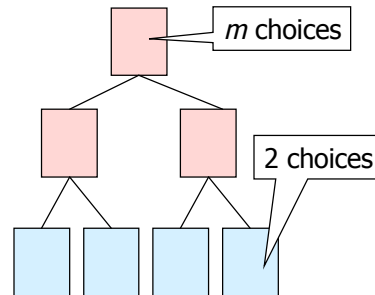
PAC-learning: Slide 18

## PAC for decision trees of depth k

- Assume  $m$  binary attributes, binary class labels
- $H_k$  = Number of decision trees of depth  $k$

Example:  $k = 2$

- $H_0 = 2$
- $H_k = (\text{\#choices of root attribute}) * (\text{\# possible left subtrees}) * (\text{\# possible right subtrees})$   
 $= m * H_{k-1} * H_{k-1}$



- Write  $L_k = \log_2 H_k$
- $L_0 = 1$
- $L_k = \log_2 m + 2L_{k-1}$
- So  $L_k = (2^k - 1)(1 + \log_2 m) + 1$
- So to PAC-learn, need

$$R \geq \frac{1}{\epsilon} \left( (2^k - 1)(1 + \log_2 m) + 1 + \log_2 \frac{1}{\delta} \right)$$

Copyright © 2001, Andrew W. Moore

PAC-learning: Slide 19

## What you should know

- Be able to understand every step in the math that gets you to

$$P(\text{we learn a bad } h) \leq H(1 - \epsilon)^R$$

- Understand that you thus need this many records to PAC-learn a machine with  $H$  different hypotheses

$$R \geq \frac{1}{\epsilon} \left( \log_2 H + \log_2 \frac{1}{\delta} \right)$$

- Understand examples of deducing  $H$  for various machines (i.e. counting hypotheses.)

Copyright © 2001, Andrew W. Moore

PAC-learning: Slide 20