

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

Ronald J. Williams
CSG220
Spring 2007

Containing many slides from the Andrew Moore tutorial of the same name.

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

Probably Approximately Correct (PAC) Learning

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

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

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
 - $X_1 \wedge X_2 \wedge X_2 \wedge X_4 \dots \wedge X_m$
- Question: if there are 3 attributes, what is the complete set of hypotheses in f ?

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

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
 - $X_1 \wedge X_2 \wedge X_2 \wedge x_4 \dots \wedge X_m$
- 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	$X_1 \wedge X_2$	$X_1 \wedge X_3$	$X_1 \wedge X_2 \wedge X_3$

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

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
 - $X_1 \wedge X_2 \wedge X_2 \wedge x_4 \dots \wedge X_m$
- Question: if there are m attributes, how many hypotheses in f ?

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

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
 - $X_1 \wedge X_2 \wedge X_2 \wedge x_4 \dots \wedge X_m$
- Question: if there are m attributes, how many hypotheses in f ? ($H = 2^m$)

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

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 ?

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

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	$\sim X_2$

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

PAC-learning: Slide 8

And-Literals Machine

- Equivalent to what we've called pure conjunctive concept descriptions when the attributes are Boolean
- E.g. $X1 \wedge \sim X3 \wedge X19$ is equivalent to $(X1 = \text{true}) \wedge (X3 = \text{false}) \wedge (X19 = \text{true})$

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

PAC-learning: Slide 9

And-Literals Machine

- $f(x, h)$ consists of all logical sentences about $X1, X2 \dots X_m$ or their negations that contain only logical ands.
- Example hypotheses:
 - $X1 \wedge \sim X3 \wedge X19$
 - $X3 \wedge \sim X18$
 - $\sim X7$
 - $X1 \wedge X2 \wedge \sim X3 \wedge \dots \wedge X_m$
- Question: if there are m attributes, what is the size of the complete set of hypotheses in f ?

True		True
True		$X2$
True		$\sim X2$
$X1$		True
$X1$	\wedge	$X2$
$X1$	\wedge	$\sim X2$
$\sim X1$		True
$\sim X1$	\wedge	$X2$
$\sim X1$	\wedge	$\sim X2$

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

PAC-learning: Slide 10

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	$\sim X_2$

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

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 ?

X_1	X_2	X_3	X_4	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

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

PAC-learning: Slide 12

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

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

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

PAC-learning: Slide 13

A Game

- We specify f , the machine
- Nature chooses hidden 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^{est} for which the training set error is 0

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

PAC-learning: Slide 14

Test Error Rate

- We specify f , the machine
- Nature chooses hidden 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^{est} for which the training set error is 0
- For each hypothesis h ,
- Say h is consistent if h has zero training set error: $\text{TRAINERR}(h) = 0$
- Define $\text{TESTERR}(h)$
 - = Fraction of test points that h will classify incorrectly
 - = $P(h \text{ classifies a random test point incorrectly})$
- Say h is bad if $\text{TESTERR}(h) > \epsilon$
- Otherwise, say h is **approximately correct**

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

PAC-learning: Slide 15

Test Error Rate

- We specify f , the machine
- Nature chooses hidden 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^{est} for which the training set error is 0
- For each hypothesis h ,
- Say h is consistent if h has zero training set error: $\text{TRAINERR}(h) = 0$
- Define $\text{TESTERR}(h)$
 - = Fraction of test points that h will classify incorrectly
 - = $P(h \text{ classifies a random test point incorrectly})$
- Say h is bad if $\text{TESTERR}(h) > \epsilon$
- Otherwise, say h is **approximately correct**

Let's consider a worst-case scenario: Among all consistent hypotheses, if any one is bad, then there's a danger that that's somehow the one we end up learning.

How probable is it that there is even one such consistent yet bad hypothesis?

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

$$\leq P(\exists h \mid h \text{ is consistent} \wedge h \text{ is bad})$$

$$= P \left(\begin{array}{l} (h_1 \text{ is consistent} \wedge h_1 \text{ is bad}) \vee \\ (h_2 \text{ is consistent} \wedge h_2 \text{ is bad}) \vee \\ \vdots \\ (h_H \text{ is consistent} \wedge h_H \text{ is bad}) \end{array} \right)$$

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

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

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

PAC-learning: Slide 16

Bounding the probability of learning a bad hypothesis

- What is $P(h_i \text{ is consistent} \mid h_i \text{ is bad})$?
- Note that if h_i is a bad hypothesis, then the probability it classifies any single training example correctly is $\leq 1-\epsilon$.
- Then, using the i.i.d. assumption, the probability it classifies all R training examples correctly is $\leq (1-\epsilon)^R$.
- Therefore we have shown that
$$P(h_i \text{ is consistent} \mid h_i \text{ is bad}) \leq (1 - \epsilon)^R$$
for any i .

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

PAC-learning: Slide 17

Bounding the prob. of a bad hypothesis

- Thus

$$\begin{aligned} P(\text{we learn a bad } h) &\leq \sum_{i=1}^H P(h_i \text{ is consistent} \mid h_i \text{ is bad}) \\ &\leq \sum_{i=1}^H (1-\epsilon)^R \\ &= H(1-\epsilon)^R \end{aligned}$$

- We can combine this with the fact that $1-\epsilon \leq e^{-\epsilon}$ to conclude

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

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

PAC-learning: Slide 18

Probably Approximately Correct

- Suppose we want the probability to be at least $1-\delta$ that the h we learn is not bad.
- A sufficient condition is that

$$\delta \geq H e^{-\epsilon R}$$

- If H , R , δ , and ϵ satisfy this relationship, then with probability $\geq 1-\delta$ we are assured that the test error rate of the h we learn is $\leq \epsilon$.
- The h we learn is **probably** (with probability $\geq 1-\delta$) **approximately** (with error rate $\leq \epsilon$) **correct**.

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

PAC-learning: Slide 19

PAC Learning

Two ways to use a sufficient condition like

$$\delta \geq H e^{-\epsilon R}$$

1. Given that we've found a consistent hypothesis h^{est} for a training set of size R , how confident are we that its test error rate is no worse than some given ϵ ? **Like confidence intervals in statistical parameter estimation theory.**

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

PAC-learning: Slide 20

PAC Learning

Two ways to use a sufficient condition like

$$\delta \geq H e^{-\epsilon R}$$

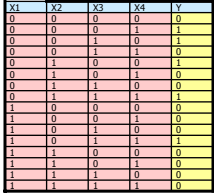
1. Given that we've found a consistent hypothesis h^{est} for a training set of size R , how confident are we that its test error rate is no worse than some given ϵ ? Like confidence intervals in statistical parameter estimation theory.
2. Sample complexity: Given δ and ϵ , how large must R be to guarantee that, with probability at least $1 - \delta$, h^{est} has a test error rate no worse than ϵ ? Get an answer by solving for R :

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

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

PAC-learning: Slide 21

PAC in action

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

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

PAC-learning: Slide 22

Extensions to PAC Analysis

- What if our learner does not produce a hypothesis with $\text{TRAINERR}(h) = 0$ (perhaps because of noisy data or limited representational power)? More generally, say h is a bad hypothesis if $\text{TESTERR}(h) > \text{TRAINERR}(h) + \epsilon$.
- In this case it turns out that the corresponding probability of learning a bad hypothesis is bounded by

$$He^{-2\epsilon^2 R}$$

- Thus to guarantee with probability at least $1-\delta$ that $\text{TESTERR}(h) \leq \text{TRAINERR}(h) + \epsilon$, it is sufficient to have a training set of size

$$R \geq \frac{1}{2\epsilon^2} \left(\ln H + \ln \frac{1}{\delta} \right)$$

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

PAC-learning: Slide 23

Extensions to PAC Analysis

- What if our hypothesis space is infinite?
- E.g.
 - perceptrons
 - multilayer neural networks
 - support vector machines
- In this case the bounds we've given are useless.
- Can we still bound the probability that $\text{TESTERR}(h) \leq \text{TRAINERR}(h) + \epsilon$ for given ϵ ?

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

PAC-learning: Slide 24

Extensions to PAC Analysis

- What if our hypothesis space is infinite?
- E.g.
 - perceptrons
 - multilayer neural networks
 - support vector machines
- In this case the bounds we've given are useless.
- Can we still bound the probability that $\text{TESTERR}(h) \leq \text{TRAINERR}(h) + \epsilon$ for given ϵ ?
- Perhaps surprisingly, the answer is YES, at least in many situations
- **Magic words: VC (Vapnik-Chervonenkis) dimension**

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

PAC-learning: Slide 25

Remarks

- This form of analysis makes no assumption about the underlying distribution of examples – just assumes same one used for both training and testing. Therefore valid for *any* distribution.
 - **Distribution free.**
- The lower bounds we've computed on the sample complexity are sufficient but not necessary for PAC-learning. But there are corresponding results providing lower bounds on the number of training examples necessary for PAC-learning with certain distributions.

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

PAC-learning: Slide 26

Remarks

- The underlying randomness in this theory is based on the randomness in the training sample
- The bounds derived from this theory are very conservative, for several reasons:
 - designed to handle any distribution of examples, including worst-case
 - derivation in PAC case, for example, based on bounding the prob. that there is *any* h that is both consistent and bad – when we select one, it could easily be better than this worst-case one

Questions to test your understanding of our PAC analysis:

1. What can be said about the *best-case* consistent hypothesis?
2. Can you see how to easily make a very, very slight improvement in the bound we derived on the probability of learning a bad h ?

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

PAC-learning: Slide 27

What you should know

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

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

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

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

- Understand examples of deducing H for various machines

Originals © 2001, Andrew W. Moore, Modifications © 2003, Ronald J. Williams

PAC-learning: Slide 28

What you should know

- Understand the generalization to nonzero training error, where having this many records is sufficient to guarantee with high probability that $\text{TESTERR}(h)$ is not much worse than $\text{TRAINERR}(h)$ when learning a machine with H hypotheses:

$$R \geq \frac{1}{2\epsilon^2} \left(\ln H + \ln \frac{1}{\delta} \right)$$