# A Boosting Tutorial

**Rob Schapire**
Princeton University

www.cs.princeton.edu/~schapire

# Example: "How May I Help You?"

- **goal:** automatically categorize type of call requested by phone customer (Collect, CallingCard, PersonToPerson, etc.)

  - `yes I'd like to place a collect call long distance please`  (Collect)
  - `operator I need to make a call but I need to bill it to my office` (ThirdNumber)
  - `yes I'd like to place a call on my master card please`  (CallingCard)
  - `I just called a number in sioux city and I musta rang the wrong number because I got the wrong party and I would like to have that taken off of my bill`  (BillingCredit)

- **observation:**

  - **easy** to find "rules of thumb" that are "often" correct
    - e.g.: "IF 'card' occurs in utterance
      THEN predict 'CallingCard' "
  - **hard** to find **single** highly accurate prediction rule

# The Boosting Approach

- devise computer program for deriving rough rules of thumb
- apply procedure to subset of examples
- obtain rule of thumb
- apply to 2nd subset of examples
- obtain 2nd rule of thumb
- repeat $T$ times

## Details

- how to choose examples on each round?
  - concentrate on "hardest" examples
    (those most often misclassified by previous rules of thumb)
- how to combine rules of thumb into single prediction rule?
  - take (weighted) majority vote of rules of thumb

# Boosting

- boosting = general method of converting rough rules of thumb into highly accurate prediction rule
- technically:
  - assume given "weak" learning algorithm that can consistently find classifiers ("rules of thumb") at least slightly better than random, say, accuracy $\geq 55\%$ (in two-class setting)
  - given sufficient data, a boosting algorithm can provably construct single classifier with very high accuracy, say, 99%

## Outline of Tutorial

- brief background

- basic algorithm and core theory

- other ways of understanding boosting

- experiments, applications and extensions

# Brief Background

# The Boosting Problem

- "strong" PAC algorithm
  - for any distribution
  - $\forall \epsilon > 0, \delta > 0$
  - given polynomially many random examples
  - finds classifier with error $\leq \epsilon$ with probability $\geq 1 - \delta$

- "weak" PAC algorithm
  - same, but only for $\epsilon \geq \frac{1}{2} - \gamma$

- [Kearns & Valiant '88]:
  - does weak learnability imply strong learnability?

## Early Boosting Algorithms

- [Schapire '89]:
  - first provable boosting algorithm
    - call weak learner three times on three modified distributions
    - get slight boost in accuracy
    - apply recursively

- [Freund '90]:
  - "optimal" algorithm that "boosts by majority"

- [Drucker, Schapire & Simard '92]:
  - first experiments using boosting
  - limited by practical drawbacks

# AdaBoost

- [Freund & Schapire '95]:

  - introduced "AdaBoost" algorithm
  - strong practical advantages over previous boosting algorithms

- ## experiments and applications using AdaBoost:

| | | |
|---|---|---|
| [Drucker & Cortes '96] | [Abney, Schapire & Singer '99] | [Tieu & Viola '00] |
| [Jackson & Craven '96] | [Haruno, Shirai & Ooyama '99] | [Walker, Rambow & Rogati '01] |
| [Freund & Schapire '96] | [Cohen & Singer' 99] | [Rochery, Schapire, Rahim & Gupta '01] |
| [Quinlan '96] | [Dietterich '00] | [Merler, Furlanello, Larcher & Sboner '01] |
| [Breiman '96] | [Schapire & Singer '00] | [Di Fabbrizio, Dutton, Gupta et al. '02] |
| [Maclin & Opitz '97] | [Collins '00] | [Qu, Adam, Yasui et al. '02] |
| [Bauer & Kohavi '97] | [Escudero, Màrquez & Rigau '00] | [Tur, Schapire & Hakkani-Tür '03] |
| [Schwenk & Bengio '98] | [Iyer, Lewis, Schapire et al. '00] | [Viola & Jones '04] |
| [Schapire, Singer & Singhal '98] | [Onoda, Rätsch & Müller '00] | [Middendorf, Kundaje, Wiggins et al. '04] |
| | | ⋮ |

- ## continuing development of theory and algorithms:

| | | |
|---|---|---|
| [Breiman '98, '99] | [Duffy & Helmbold '99, '02] | [Koltchinskii, Panchenko & Lozano '01] |
| [Schapire, Freund, Bartlett & Lee '98] | [Freund & Mason '99] | [Collins, Schapire & Singer '02] |
| [Grove & Schuurmans '98] | [Ridgeway, Madigan & Richardson '99] | [Demiriz, Bennett & Shawe-Taylor '02] |
| [Mason, Bartlett & Baxter '98] | [Kivinen & Warmuth '99] | [Lebanon & Lafferty '02] |
| [Schapire & Singer '99] | [Friedman, Hastie & Tibshirani '00] | [Wyner '02] |
| [Cohen & Singer '99] | [Rätsch, Onoda & Müller '00] | [Rudin, Daubechies & Schapire '03] |
| [Freund & Mason '99] | [Rätsch, Warmuth, Mika et al. '00] | [Jiang '04] |
| [Domingo & Watanabe '99] | [Allwein, Schapire & Singer '00] | [Lugosi & Vayatis '04] |
| [Mason, Baxter, Bartlett & Frean '99, '00] | [Friedman '01] | [Zhang '04] |
| | | ⋮ |

# Basic Algorithm and Core Theory

# A Formal Description of Boosting

- given <u>training set</u> $(x_1, y_1), \ldots, (x_m, y_m)$

- $y_i \in \{-1, +1\}$ correct label of instance $x_i \in X$

- for $t = 1, \ldots, T$:
  - construct distribution $D_t$ on $\{1, \ldots, m\}$
  - find <u>weak classifier</u> ("rule of thumb")
  
    $$h_t : X \rightarrow \{-1, +1\}$$
    
    with small <u>error</u> $\epsilon_t$ on $D_t$:
    
    $$\epsilon_t = \Pr_{D_t}[h_t(x_i) \neq y_i]$$

- output <u>final classifier</u> $H_{\text{final}}$

# AdaBoost

- constructing $D_t$:
  - $D_1(i) = 1/m$
  - given $D_t$ and $h_t$:

$$D_{t+1}(i) = \frac{D_t(i)}{Z_t} \times \begin{cases} e^{-\alpha_t} & \text{if } y_i = h_t(x_i) \\ e^{\alpha_t} & \text{if } y_i \neq h_t(x_i) \end{cases}$$

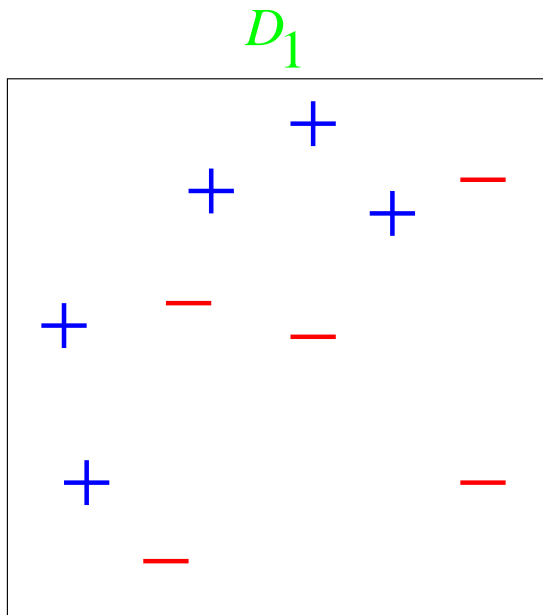$$= \frac{D_t(i)}{Z_t} \exp(-\alpha_t \, y_i \, h_t(x_i))$$

where $Z_t$ = normalization constant

$$\alpha_t = \tfrac{1}{2} \ln\left(\frac{1 - \epsilon_t}{\epsilon_t}\right) > 0$$
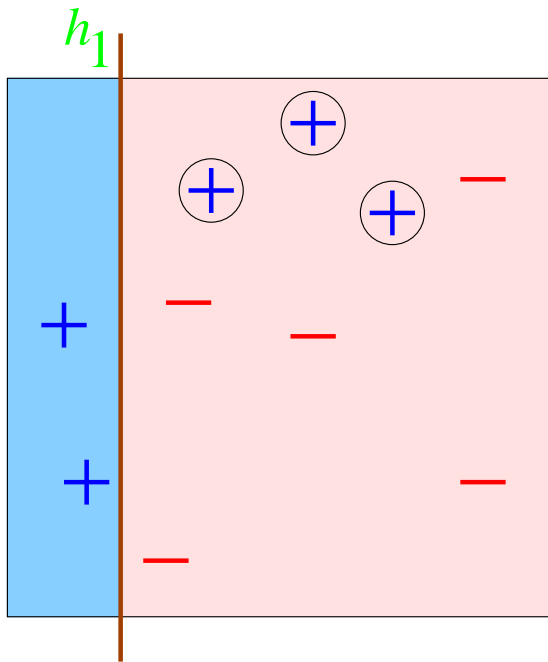
- final classifier:
  - $H_{\text{final}}(x) = \text{sign}\left(\sum\limits_t \alpha_t h_t(x)\right)$
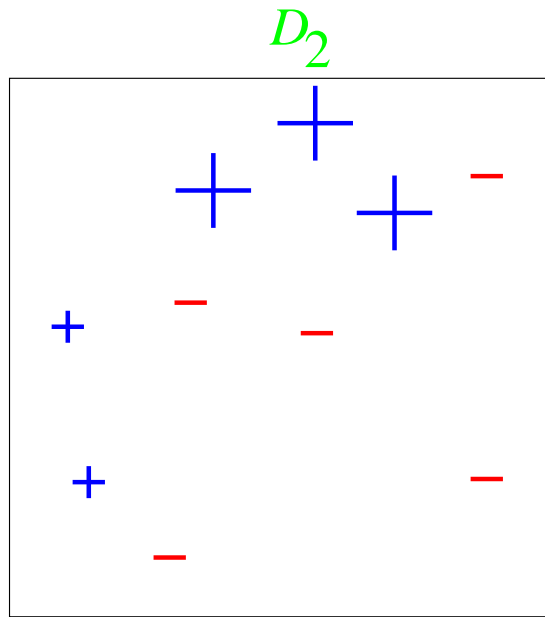
# Toy Example

$D_1$



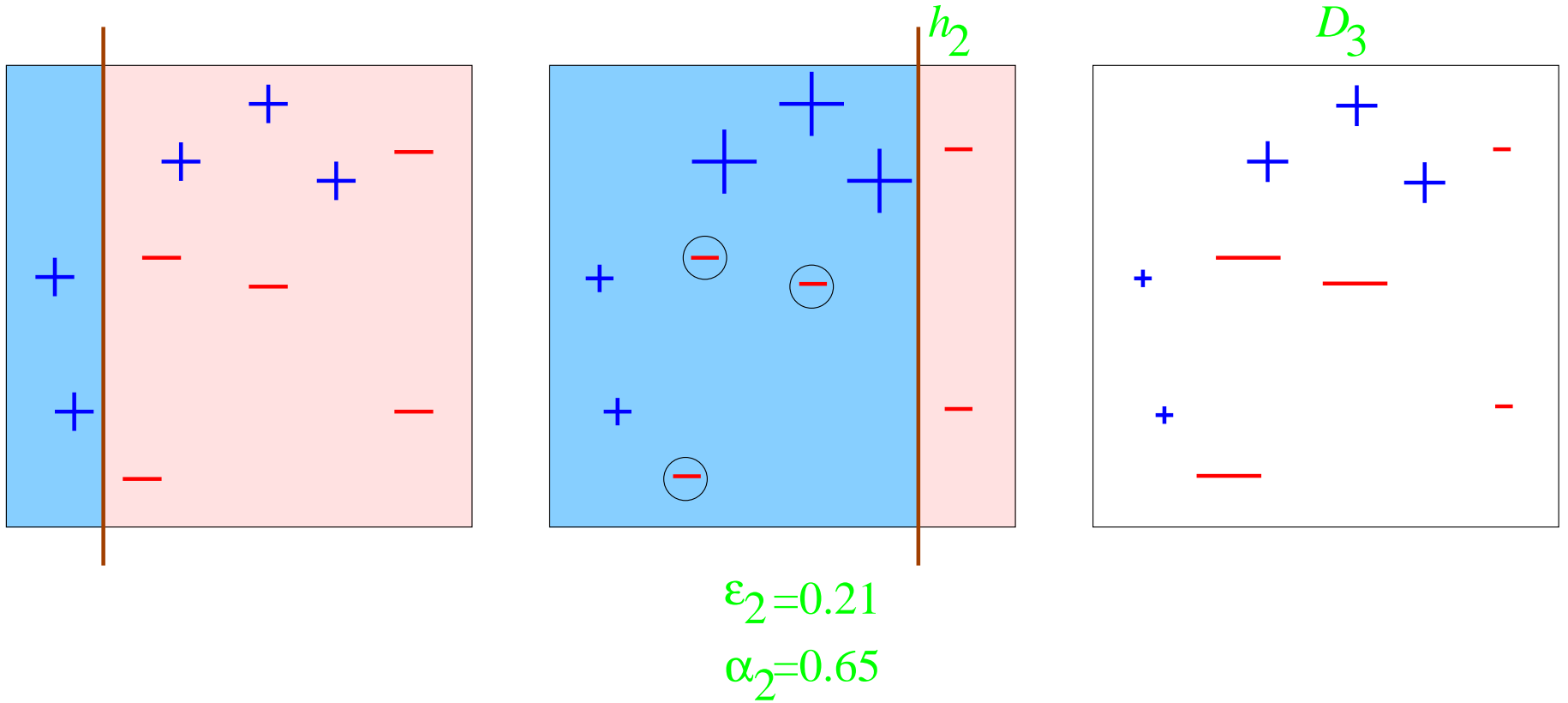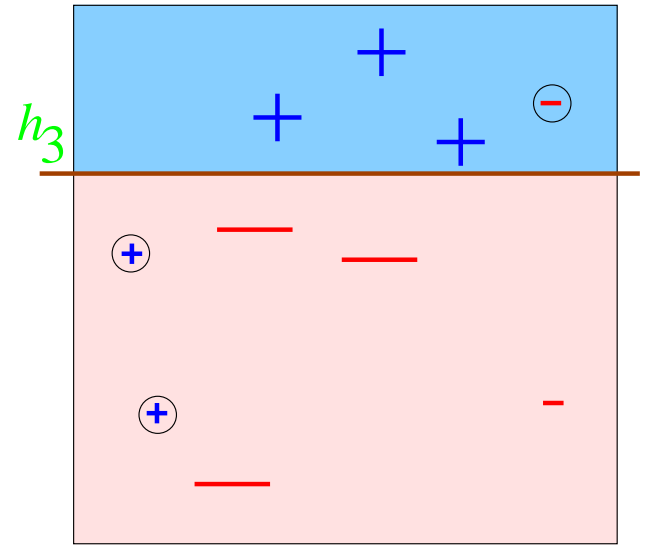weak classifiers = vertical or horizontal half-planes

# Round 1



$h_1$

$D_2$

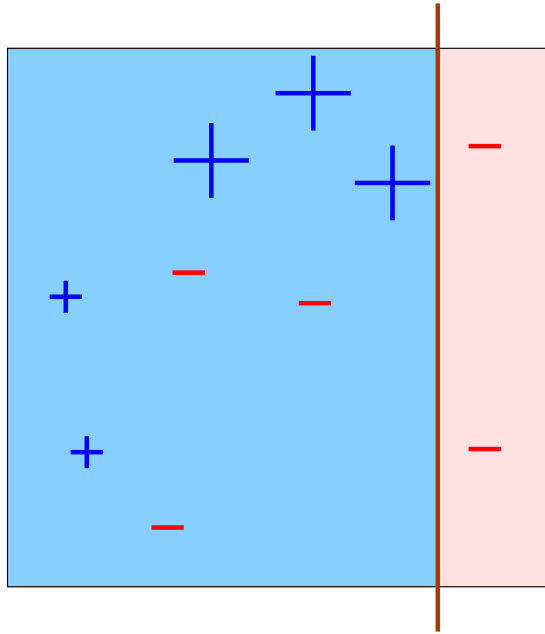$\varepsilon_1 = 0.30$

$\alpha_1 = 0.42$

# Round 2



$h_2$

$D_3$

$\varepsilon_2 = 0.21$

$\alpha_2 = 0.65$

# Round 3



$h_3$

$\varepsilon_3=0.14$

$\alpha_3=0.92$

# Final Classifier

$$H_{final} = \text{sign}\left( 0.42 \quad\boxed{\phantom{XXXX}}\quad + 0.65 \quad\boxed{\phantom{XXXX}}\quad + 0.92 \quad\boxed{\phantom{XXXX}} \right)$$

# Analyzing the training error

- Theorem:
  - write $\epsilon_t$ as $1/2 - \gamma_t$
  - then

$$\text{training error}(H_{\text{final}}) \leq \prod_t \left[ 2\sqrt{\epsilon_t(1 - \epsilon_t)} \right]$$

$$= \prod_t \sqrt{1 - 4\gamma_t^2}$$

$$\leq \exp\left( -2 \sum_t \gamma_t^2 \right)$$

- so: if $\forall t: \gamma_t \geq \gamma > 0$
  - then $\text{training error}(H_{\text{final}}) \leq e^{-2\gamma^2 T}$
- AdaBoost is adaptive:
  - does not need to know $\gamma$ or $T$ a priori
  - can exploit $\gamma_t \gg \gamma$

## Proof

- let $f(x) = \sum_t \alpha_t h_t(x) \Rightarrow H_{\text{final}}(x) = \text{sign}(f(x))$

- *Step 1*: unwrapping recurrence:

$$D_{\text{final}}(i) = \frac{1}{m} \frac{\exp\left(-y_i \sum_t \alpha_t h_t(x_i)\right)}{\prod_t Z_t}$$

$$= \frac{1}{m} \frac{\exp\left(-y_i f(x_i)\right)}{\prod_t Z_t}$$

# Proof (cont.)

- *Step 2*: training error$(H_{\text{final}}) \leq \prod_t Z_t$

- Proof:
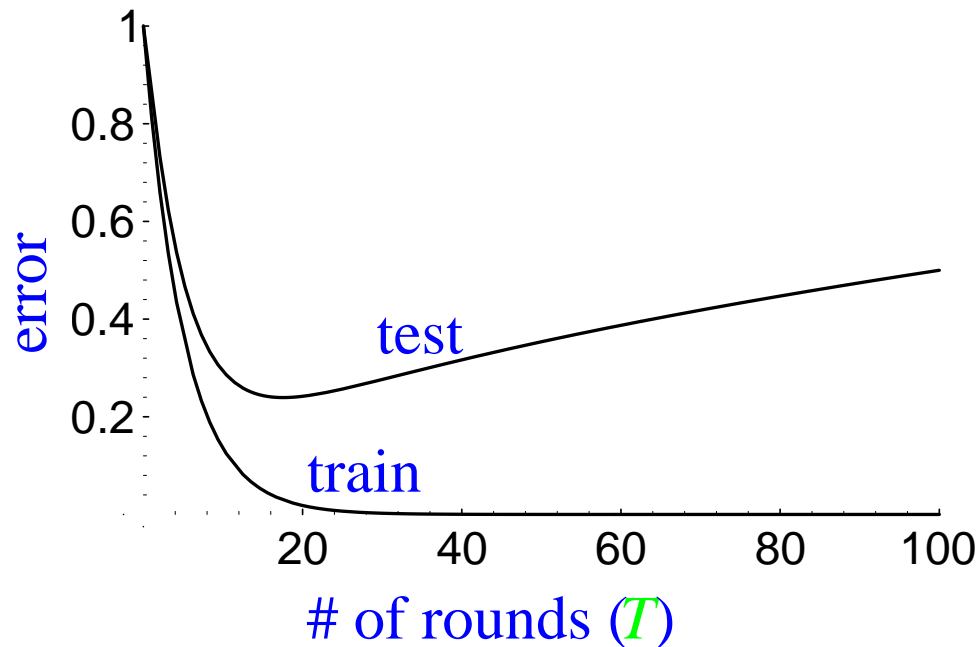
$$\text{training error}(H_{\text{final}}) = \frac{1}{m} \sum_i \begin{cases} 1 & \text{if } y_i \neq H_{\text{final}}(x_i) \\ 0 & \text{else} \end{cases}$$

$$= \frac{1}{m} \sum_i \begin{cases} 1 & \text{if } y_i f(x_i) \leq 0 \\ 0 & \text{else} \end{cases}$$

$$\leq \frac{1}{m} \sum_i \exp(-y_i f(x_i))$$

$$= \sum_i D_{\text{final}}(i) \prod_t Z_t$$

$$= \prod_t Z_t$$

## Proof (cont.)

- *Step 3*: $Z_t = 2\sqrt{\epsilon_t(1 - \epsilon_t)}$

- Proof:

$$
\begin{aligned}
Z_t &= \sum_i D_t(i) \exp(-\alpha_t \, y_i \, h_t(x_i)) \\
&= \sum_{i:y_i \neq h_t(x_i)} D_t(i) e^{\alpha_t} + \sum_{i:y_i = h_t(x_i)} D_t(i) e^{-\alpha_t} \\
&= \epsilon_t \, e^{\alpha_t} + (1 - \epsilon_t) \, e^{-\alpha_t} \\
&= 2\sqrt{\epsilon_t(1 - \epsilon_t)}
\end{aligned}
$$

# How Will Test Error Behave? (A First Guess)



expect:

- training error to continue to drop (or reach zero)
- test error to increase when $H_{\text{final}}$ becomes "too complex"
  - "Occam's razor"
  - overfitting
    - hard to know when to stop training

# Actual Typical Run



(boosting C4.5 on "letter" dataset)

- test error does <u>not</u> increase, even after 1000 rounds
  - (total size $> 2{,}000{,}000$ nodes)

- test error continues to drop even after training error is zero!

|  | # rounds | | |
| --- | --- | --- | --- |
|  | 5 | 100 | 1000 |
| train error | 0.0 | 0.0 | 0.0 |
| test error | 8.4 | 3.3 | 3.1 |

- Occam's razor <u>wrongly</u> predicts "simpler" rule is better

# A Better Story: Theory of Margins

[with Freund, Bartlett & Lee]

- <u>key idea</u>:
  - training error only measures whether classifications are right or wrong
  - should also consider <u>confidence</u> of classifications

- recall: $H_{\text{final}}$ is weighted majority vote of weak classifiers

- measure confidence by <u>margin</u> = strength of the vote
     = (fraction voting correctly) − (fraction voting incorrectly)

high conf. incorrect        low conf.        high conf. correct

$H_{\text{final}}$      $H_{\text{final}}$

−1    incorrect    0    correct    +1

# Empirical Evidence: The Margin Distribution

- margin distribution
  = cumulative distribution of margins of training examples



| | # rounds | | |
|---|---|---|---|
| | 5 | 100 | 1000 |
| train error | 0.0 | 0.0 | 0.0 |
| test error | 8.4 | 3.3 | 3.1 |
| % margins $\leq 0.5$ | 7.7 | 0.0 | 0.0 |
| minimum margin | 0.14 | 0.52 | 0.55 |

# Theoretical Evidence: Analyzing Boosting Using Margins

- Theorem: large margins $\Rightarrow$ better bound on generalization error (independent of number of rounds)
  - proof idea: if all margins are large, then can approximate final classifier by a much smaller classifier (just as polls can predict not-too-close election)

- Theorem: boosting tends to increase margins of training examples (given weak learning assumption)
  - proof idea: similar to training error proof

- so:
  although final classifier is getting larger,
  margins are likely to be increasing,
  so final classifier actually getting close to a simpler classifier,
  driving down the test error

## More Technically...

- with high probability, $\forall \theta > 0$ :

$$\text{generalization error} \leq \hat{\Pr}[\text{margin} \leq \theta] + \tilde{O}\left(\frac{\sqrt{d/m}}{\theta}\right)$$

  where
  - $m = \#$ training examples
  - $d = $ "complexity" of weak classifiers
- $\hat{\Pr}[\text{margin} \leq \theta] \to 0$ exponentially fast (in $T$) if $\gamma_t > \theta$ ($\forall t$)

# Other Ways of Understanding AdaBoost

# Game Theory

- game defined by matrix $\mathbf{M}$:

|  | Rock | Paper | Scissors |
|---|---|---|---|
| Rock | 1/2 | 1 | 0 |
| Paper | 0 | 1/2 | 1 |
| Scissors | 1 | 0 | 1/2 |

- row player chooses row $i$

- column player chooses column $j$ (simultaneously)

- row player's goal: minimize loss $\mathbf{M}(i, j)$

- usually allow randomized play:
  - players choose distributions $\mathbf{P}$ and $\mathbf{Q}$ over rows and columns

- learner's (expected) loss

$$= \sum_{i,j} \mathbf{P}(i)\mathbf{M}(i, j)\mathbf{Q}(j)$$

$$= \mathbf{P}^\mathrm{T}\mathbf{M}\mathbf{Q} \equiv \mathbf{M}(\mathbf{P}, \mathbf{Q})$$

# The Minmax Theorem

- von Neumann's minmax theorem:

$$\min_{\mathbf{P}} \max_{\mathbf{Q}} \mathbf{M}(\mathbf{P}, \mathbf{Q}) = \max_{\mathbf{Q}} \min_{\mathbf{P}} \mathbf{M}(\mathbf{P}, \mathbf{Q})$$
$$= v$$
$$= \text{"value" of game } \mathbf{M}$$

- in words:
  - $v = \min \max$ means:
    - row player has strategy $\mathbf{P}^*$
      such that $\forall$ column strategy $\mathbf{Q}$
      loss $\mathbf{M}(\mathbf{P}^*, \mathbf{Q}) \leq v$
  - $v = \max \min$ means:
    - this is optimal in sense that
      column player has strategy $\mathbf{Q}^*$
      such that $\forall$ row strategy $\mathbf{P}$
      loss $\mathbf{M}(\mathbf{P}, \mathbf{Q}^*) \geq v$

# The Boosting Game

- let $\{\overline{h}_1, \ldots, \overline{h}_N\}$ = space of <u>all</u> weak classifiers
- row player $\leftrightarrow$ booster
- column player $\leftrightarrow$ weak learner
- matrix **M**:
  - row $\leftrightarrow$ example $(x_i, y_i)$
  - column $\leftrightarrow$ weak classifier $\overline{h}$
  - $\mathbf{M}(i,j) = \begin{cases} 1 & \text{if } y_i = \overline{h}_j(x_i) \\ 0 & \text{else} \end{cases}$

weak learner

$\overline{h}_1 \quad\cdots\quad \overline{h}_j \quad\cdots\quad \overline{h}_N$

booster

$x_1\, y_1$

$x_i\, y_i \qquad\qquad M(i,j)$

$x_m y_m$

# Boosting and the Minmax Theorem

- if:
  - $\forall$ distributions over examples
    $\exists h$ with accuracy $\geq \frac{1}{2} - \gamma$
- then:
  - $\min_{\mathbf{P}} \max_{h} \mathbf{M}(\mathbf{P}, h) \geq \frac{1}{2} - \gamma$
- by minmax theorem:
  - $\max_{\mathbf{Q}} \min_{i} \mathbf{M}(i, \mathbf{Q}) \geq \frac{1}{2} - \gamma > \frac{1}{2}$
- which means:
  - $\exists$ weighted majority of classifiers which correctly classifies
    all examples with positive margin $(2\gamma)$
- optimal margin $\leftrightarrow$ "value" of game

# AdaBoost and Game Theory

- AdaBoost is special case of general algorithm for solving games through repeated play

- can show
  - distribution over examples converges to (approximate) minmax strategy for boosting game
  - weights on weak classifiers converge to (approximate) maxmin strategy

- different instantiation of game-playing algorithm gives on-line learning algorithms (such as weighted majority algorithm)

# AdaBoost and Exponential Loss

- many (most?) learning algorithms minimize a "loss" function
  - e.g. least squares regression
- training error proof shows AdaBoost actually minimizes

$$\prod_t Z_t = \frac{1}{m} \sum_i \exp(-y_i f(x_i))$$

where $f(x) = \sum_t \alpha_t h_t(x)$

- on each round, AdaBoost <u>greedily</u> chooses $\alpha_t$ and $h_t$ to minimize loss

- exponential loss is an upper bound on 0-1 (classification) loss

- AdaBoost <u>provably</u> minimizes exponential loss



$yf(x)$

# Coordinate Descent

[Breiman]

- $\{\overline{h}_1, \ldots, \overline{h}_N\}$ = space of <u>all</u> weak classifiers
- want to find $\lambda_1, \ldots, \lambda_N$ to minimize

$$L(\lambda_1, \ldots, \lambda_N) = \sum_i \exp\left(-y_i \sum_j \lambda_j \overline{h}_j(x_i)\right)$$

- AdaBoost is actually doing <u>coordinate descent</u> on this optimization problem:
  - initially, all $\lambda_j = 0$
  - each round: choose one coordinate $\lambda_j$ (corresponding to $h_t$) and update (increment by $\alpha_t$)
  - choose update causing biggest decrease in loss
- powerful technique for minimizing over huge space of functions

# Functional Gradient Descent

- want to minimize

$$L(f) = L(f(x_1), \ldots, f(x_m)) = \sum_i \exp(-y_i f(x_i))$$

- say have current estimate $\overline{f}$ and want to improve

- to do gradient descent, would like update

$$\overline{f} \leftarrow \overline{f} - \eta \nabla_f L(\overline{f})$$

- but update restricted in class of weak classifiers

- so choose $h_t$ "closest" to $\nabla_f L(\overline{f})$

- equivalent to AdaBoost

## Benefits of Model Fitting View

- immediate generalization to other loss functions
  - e.g. squared error for regression
  - e.g. logistic regression (by only changing one line of AdaBoost)

- sensible approach for converting output of boosting into conditional probability estimates

- caveat: wrong to view AdaBoost as just an algorithm for minimizing exponential loss
  - other algorithms for minimizing same loss will (provably) give very poor performance
  - thus, this loss function cannot explain why AdaBoost "works"

# Estimating Conditional Probabilities

- often want to estimate <u>probability</u> that $y = +1$ given $x$

- AdaBoost minimizes (empirical version of):

$$\mathrm{E}_{x,y}\left[e^{-yf(x)}\right] = \mathrm{E}_x\left[\mathrm{P}\left[y = +1 | x\right] e^{-f(x)} + \mathrm{P}\left[y = -1 | x\right] e^{-f(x)}\right]$$

  where $x, y$ random from true distribution

- over <u>all</u> $f$, minimized when

$$f(x) = \frac{1}{2} \cdot \ln\left(\frac{\mathrm{P}\left[y = +1 | x\right]}{\mathrm{P}\left[y = -1 | x\right]}\right)$$

  or

$$\mathrm{P}\left[y = +1 | x\right] = \frac{1}{1 + e^{-2f(x)}}$$

- so, to convert $f$ output by AdaBoost to probability estimate, use same formula

# Calibration Curve



- order examples by $f$ value output by AdaBoost

- break into bins of size $r$

- for each bin, plot a point:
  - $x$-value: average estimated probability of examples in bin
  - $y$-value: actual fraction of positive examples in bin

# Other Ways to Think about AdaBoost

- dynamical systems

- statistical consistency

- maximum entropy

# Experiments, Applications and Extensions

## Practical Advantages of AdaBoost

- fast

- simple and easy to program

- no parameters to tune (except $T$)

- flexible — can combine with any learning algorithm

- no prior knowledge needed about weak learner

- provably effective, provided can consistently find rough rules of thumb

  → shift in mind set — goal now is merely to find classifiers barely better than random guessing

- versatile

  - can use with data that is textual, numeric, discrete, etc.
  - has been extended to learning problems well beyond binary classification

## Caveats

- performance of AdaBoost depends on <u>data</u> and <u>weak learner</u>

- consistent with theory, AdaBoost can <u>fail</u> if
  - weak classifiers too complex
  
    $\rightarrow$ overfitting
  - weak classifiers too weak ($\gamma_t \rightarrow 0$ too quickly)
  
    $\rightarrow$ underfitting
  
    $\rightarrow$ low margins $\rightarrow$ overfitting

- empirically, AdaBoost seems especially susceptible to uniform noise

# UCI Experiments

- tested AdaBoost on UCI benchmarks

- used:
    - C4.5 (Quinlan's decision tree algorithm)
    - "decision stumps": very simple rules of thumb that test on single attributes

```
         eye color = brown ?                    height > 5 feet ?
           yes           no                       yes          no
      predict       predict                 predict       predict
        +1             -1                      -1            +1
```

boosting Stumps

boosting C4.5

# Multiclass Problems

- say $y \in Y = \{1, \ldots, k\}$
- direct approach (AdaBoost.M1):

$$h_t : X \to Y$$

$$D_{t+1}(i) = \frac{D_t(i)}{Z_t} \cdot \begin{cases} e^{-\alpha_t} & \text{if } y_i = h_t(x_i) \\ e^{\alpha_t} & \text{if } y_i \neq h_t(x_i) \end{cases}$$

$$H_{\text{final}}(x) = \arg\max_{y \in Y} \sum_{t:h_t(x)=y} \alpha_t$$

- can prove same bound on error <u>if</u> $\forall t : \epsilon_t \leq 1/2$
  - in practice, not usually a problem for "strong" weak learners (e.g., C4.5)
  - significant problem for "weak" weak learners (e.g., decision stumps)
- instead, reduce to binary

# Reducing Multiclass to Binary

- say possible labels are $\{a, b, c, d, e\}$

- each training example replaced by five $\{-1, +1\}$-labeled examples:

$$x \ , \ c \ \longrightarrow \begin{cases} (x, a) \ , \ -1 \\ (x, b) \ , \ -1 \\ (x, c) \ , \ +1 \\ (x, d) \ , \ -1 \\ (x, e) \ , \ -1 \end{cases}$$

- predict with label receiving most (weighted) votes

## **AdaBoost.MH**

- can prove:

$$\text{training error}(H_{\text{final}}) \leq \frac{k}{2} \cdot \Pi \, Z_t$$

  - reflects fact that small number of errors in binary predictors can cause overall prediction to be incorrect

- extends immediately to <u>multi-label</u> case (more than one correct label per example)

# <u>Using Output Codes</u>

[with Allwein & Singer]

- alternative: choose "code word" for each label

|   | $\pi_1$ | $\pi_2$ | $\pi_3$ | $\pi_4$ |
|---|---|---|---|---|
| a | $-$ | $+$ | $-$ | $+$ |
| b | $-$ | $+$ | $+$ | $-$ |
| c | $+$ | $-$ | $-$ | $+$ |
| d | $+$ | $-$ | $+$ | $+$ |
| e | $-$ | $+$ | $-$ | $-$ |

- each training example mapped to one example per column

$$x \ , \ c \ \rightarrow \begin{cases} (x, \pi_1) \ , \ +1 \\ (x, \pi_2) \ , \ -1 \\ (x, \pi_3) \ , \ -1 \\ (x, \pi_4) \ , \ +1 \end{cases}$$

- to classify new example $x$:
  - evaluate classifier on $(x, \pi_1), \ldots, (x, \pi_4)$
  - choose label "most consistent" with results
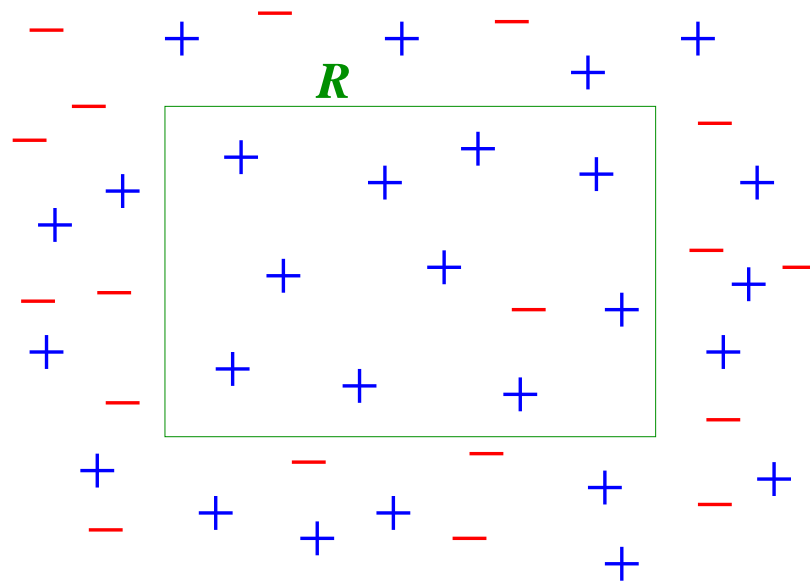
# Output Codes (cont.)

- training error bounds <u>independent</u> of # of classes

- overall prediction robust to large number of errors in binary predictors

- <u>but</u>: binary problems may be harder

# Ranking Problems

- other problems can also be handled by reducing to binary

- e.g.: want to learn to <u>rank</u> objects (say, movies) from examples

- can reduce to multiple binary questions of form:
  *"is or is not object A preferred to object B?"*

- now apply (binary) AdaBoost

# Problem with "Hard" Predictions



- ideally, want weak classifier that says:

$$h(x) = \begin{cases} +1 & \text{if } x \in R \\ \text{"}don't\ know\text{"} & \text{else} \end{cases}$$

- problem: cannot express using "hard" predictions
- if must predict $\pm 1$ outside $R$, will introduce many "bad" predictions
  - need to "clean up" on later rounds
- dramatically increases time to convergence

# Confidence-rated Predictions

- useful to allow weak classifiers to assign <u>confidences</u> to predictions

- formally, allow $h_t : X \to \mathbb{R}$

$$\text{sign}(h_t(x)) = \text{prediction}$$
$$|h_t(x)| = \text{``confidence''}$$

- use identical update:

$$D_{t+1}(i) = \frac{D_t(i)}{Z_t} \cdot \exp(-\alpha_t \, y_i \, h_t(x_i))$$

  and identical rule for combining weak classifiers

- <u>question</u>: how to choose $\alpha_t$ and $h_t$ on each round

## Confidence-rated Predictions (cont.)
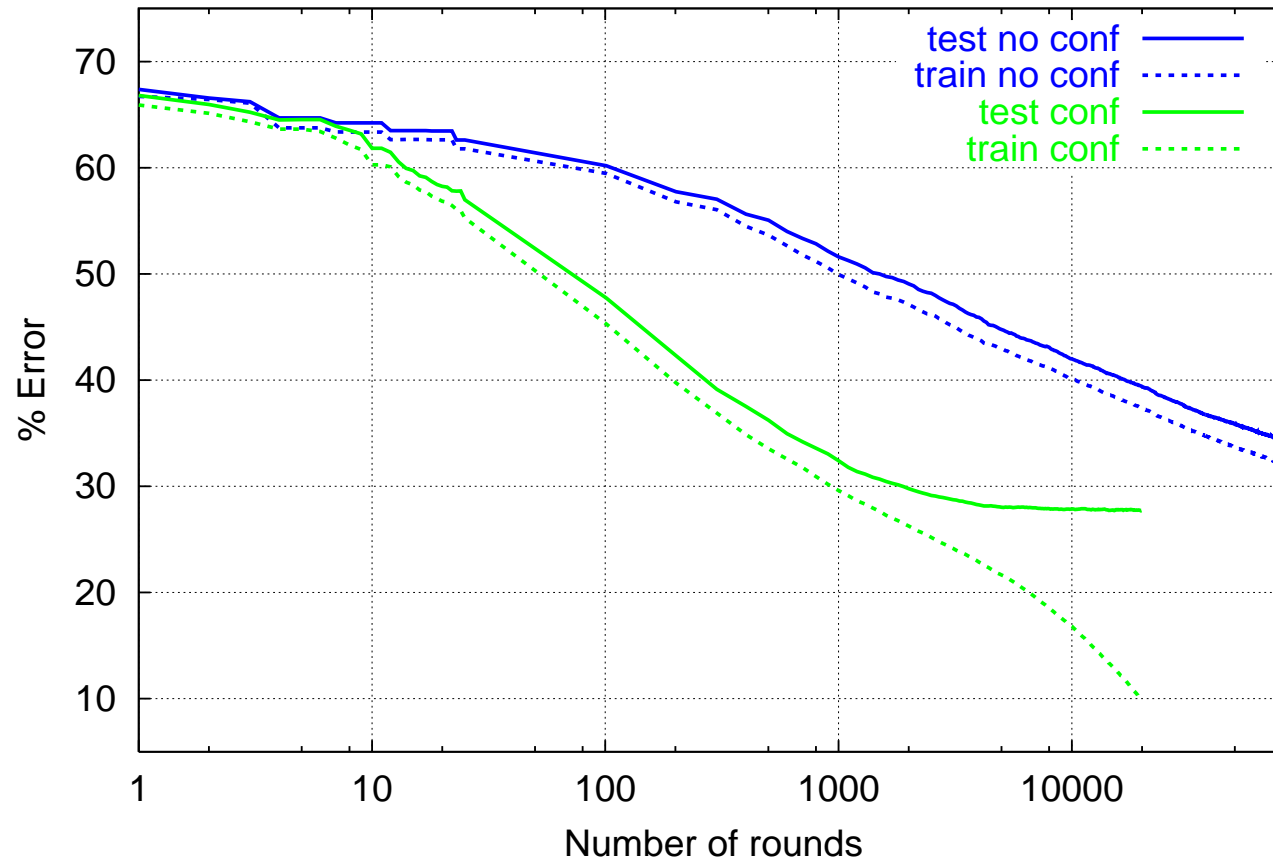
- saw earlier:

$$\text{training error}(H_{\text{final}}) \leq \prod_t Z_t = \sum_i \exp\left(-y_i \sum_t \alpha_t h_t(x_i)\right)$$

- therefore, on each round $t$, should choose $\alpha_t h_t$ to minimize:

$$Z_t = \sum_i D_t(i) \exp(-\alpha_t \ y_i \ h_t(x_i))$$

- in many cases (e.g., decision stumps), best confidence-rated weak classifier has simple form that can be found efficiently

# Confidence-rated Predictions Help a Lot



| % error | round first reached | | speedup |
|---|---|---|---|
| | conf. | no conf. | |
| 40 | 268 | 16,938 | 63.2 |
| 35 | 598 | 65,292 | 109.2 |
| 30 | 1,888 | >80,000 | – |

# Application: Boosting for Text Categorization

[with Singer]

- weak classifiers: very simple weak classifiers that test on simple patterns, namely, (sparse) $n$-grams
  - find parameter $\alpha_t$ and rule $h_t$ of given form which minimize $Z_t$
  - use efficiently implemented exhaustive search
- "How may I help you" data:
  - 7844 training examples
  - 1000 test examples
  - categories: AreaCode, AttService, BillingCredit, CallingCard, Collect, Competitor, DialForMe, Directory, HowToDial, PersonToPerson, Rate, ThirdNumber, Time, TimeCharge, Other.

# Weak Classifiers

| rnd | term | AC | AS | BC | CC | CO | CM | DM | DI | HO | PP | RA | 3N | TI | TC | OT |
|-----|------|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 1 | collect | | | | | | | | | | | | | | | |
| 2 | card | | | | | | | | | | | | | | | |
| 3 | my home | | | | | | | | | | | | | | | |
| 4 | person ? person | | | | | | | | | | | | | | | |
| 5 | code | | | | | | | | | | | | | | | |
| 6 | I | | | | | | | | | | | | | | | |

# More Weak Classifiers

| rnd | term | AC | AS | BC | CC | CO | CM | DM | DI | HO | PP | RA | 3N | TI | TC | OT |
|-----|------|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 7 | time | | | | | | | | | | | | | | | |
| 8 | wrong number | | | | | | | | | | | | | | | |
| 9 | how | | | | | | | | | | | | | | | |
| 10 | call | | | | | | | | | | | | | | | |
| 11 | seven | | | | | | | | | | | | | | | |
| 12 | trying to | | | | | | | | | | | | | | | |
| 13 | and | | | | | | | | | | | | | | | |

# More Weak Classifiers

| rnd | term | AC | AS | BC | CC | CO | CM | DM | DI | HO | PP | RA | 3N | TI | TC | OT |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 14 | third | | | | | | | | | | | | | | | |
| 15 | to | | | | | | | | | | | | | | | |
| 16 | for | | | | | | | | | | | | | | | |
| 17 | charges | | | | | | | | | | | | | | | |
| 18 | dial | | | | | | | | | | | | | | | |
| 19 | just | | | | | | | | | | | | | | | |

# Finding Outliers
examples with most weight are often outliers
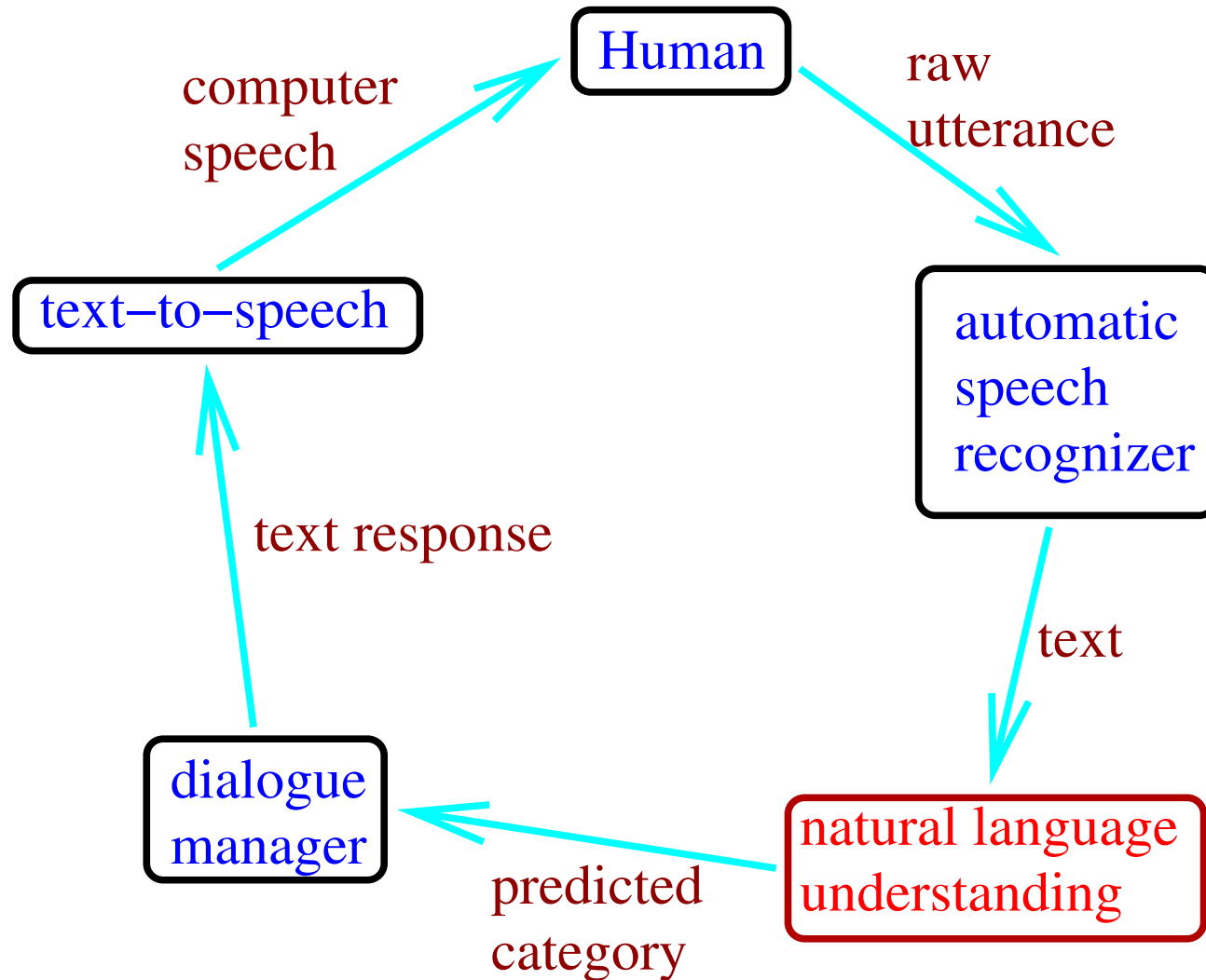(mislabeled and/or ambiguous)

- I'm trying to make a credit card call  (Collect)

- hello  (Rate)

- yes I'd like to make a long distance collect call please  (CallingCard)

- calling card please  (Collect)

- yeah I'd like to use my calling card number  (Collect)

- can I get a collect call  (CallingCard)

- yes I would like to make a long distant telephone call and have the charges billed to another number  (CallingCard DialForMe)

- yeah I can not stand it this morning I did oversea call is so bad  (BillingCredit)

- yeah special offers going on for long distance  (AttService Rate)

- mister allen please william allen  (PersonToPerson)

- yes ma'am I I'm trying to make a long distance call to a non dialable point in san miguel philippines  (AttService Other)

- yes I like to make a long distance call and charge it to my home phone that's where I'm calling at my home  (DialForMe)

## Application: Human-computer Spoken Dialogue
[with Rahim, Di Fabbrizio, Dutton, Gupta, Hollister & Riccardi]

- **application**: automatic "store front" or "help desk" for AT&T Labs' Natural Voices business

- caller can request demo, pricing information, technical support, sales agent, etc.

- interactive dialogue

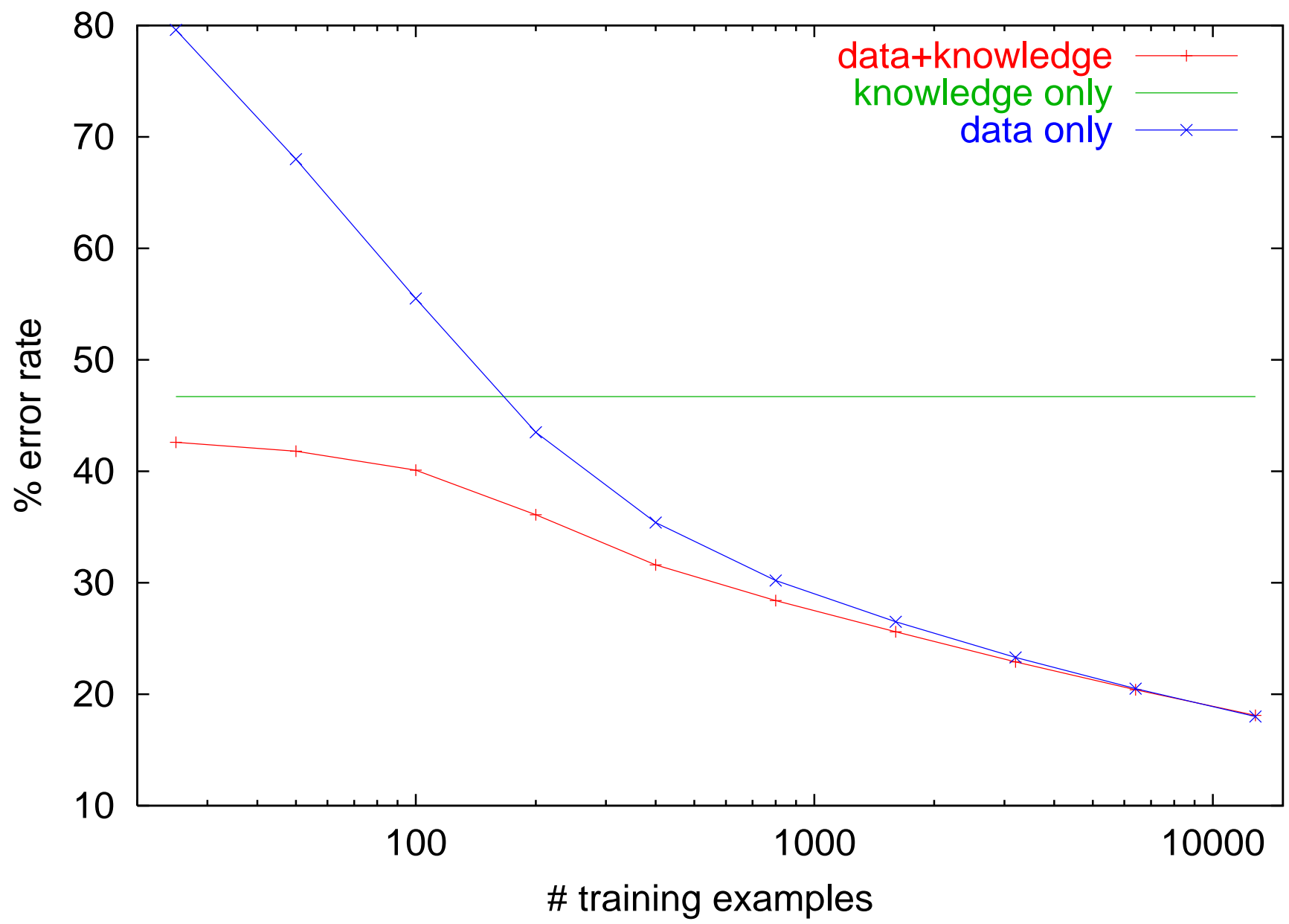- `naturalvoices.att.com`, 1-877-741-4321

## How It Works



- NLU's job: classify caller utterances into 24 categories (demo, sales rep, pricing info, yes, no, etc.)

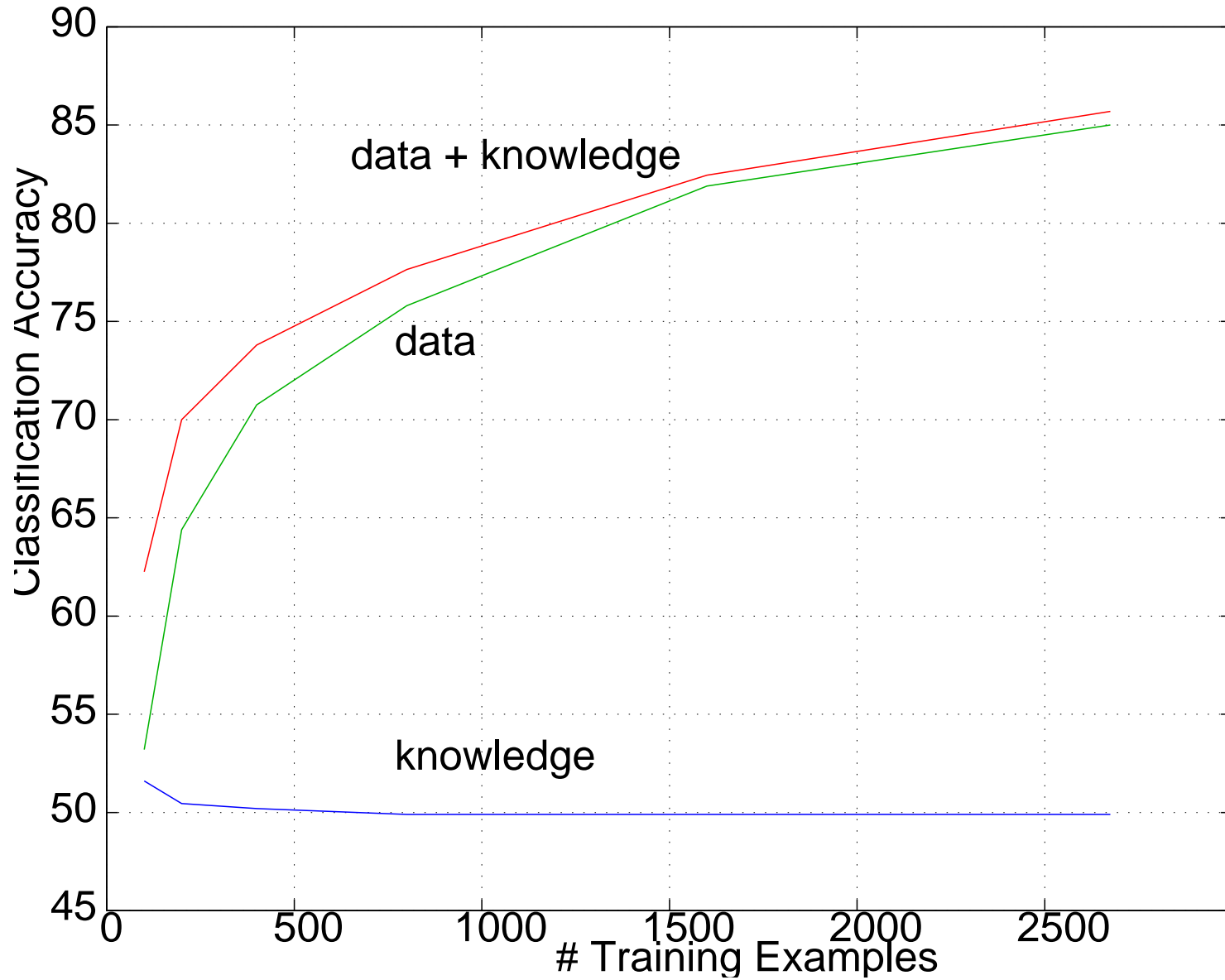## Need for Prior, Human Knowledge

[with Rochery, Rahim & Gupta]

- building NLU: standard text categorization problem

- need <u>lots of data</u>, but for cheap, <u>rapid</u> deployment, can't wait for it

- <u>bootstrapping</u> problem:
  - need labeled data to deploy
  - need to deploy to get labeled data

- <u>idea</u>: use human knowledge to compensate for insufficient data
  - modify loss function to balance <u>fit to data</u> against <u>fit to prior model</u>

# Results: AP-Titles

# Results: Helpdesk

# Problem: Labels are Expensive

- for spoken-dialogue task
  - getting examples is cheap
  - getting <u>labels</u> is expensive
    - must be annotated by humans
- how to reduce number of <u>labels</u> needed?

# Active Learning

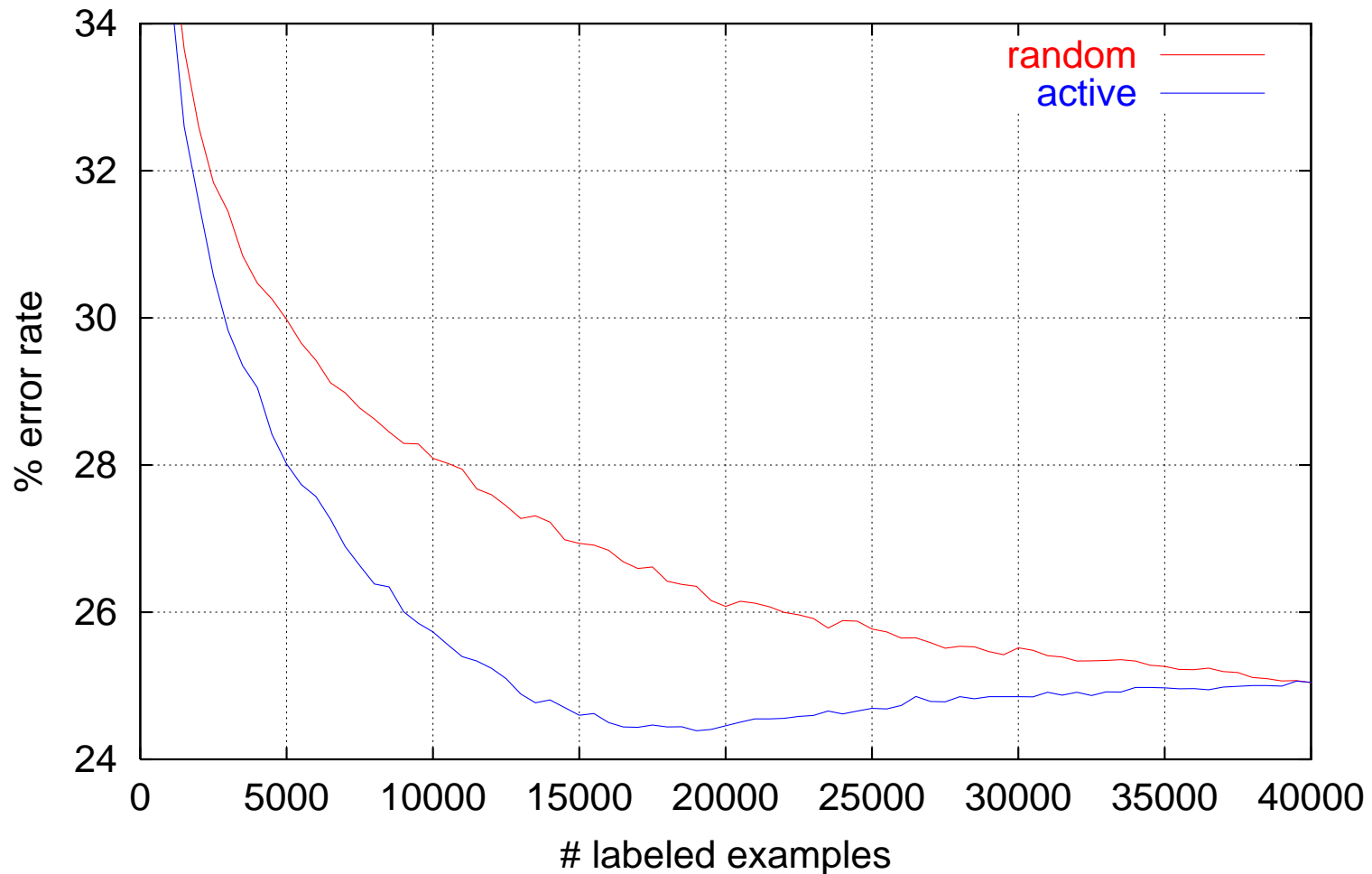- idea:
  - use <u>selective sampling</u> to choose which examples to label
  - focus on <u>least confident</u> examples                    [Lewis & Gale]

- for boosting, use (absolute) margin $|f(x)|$ as natural confidence measure

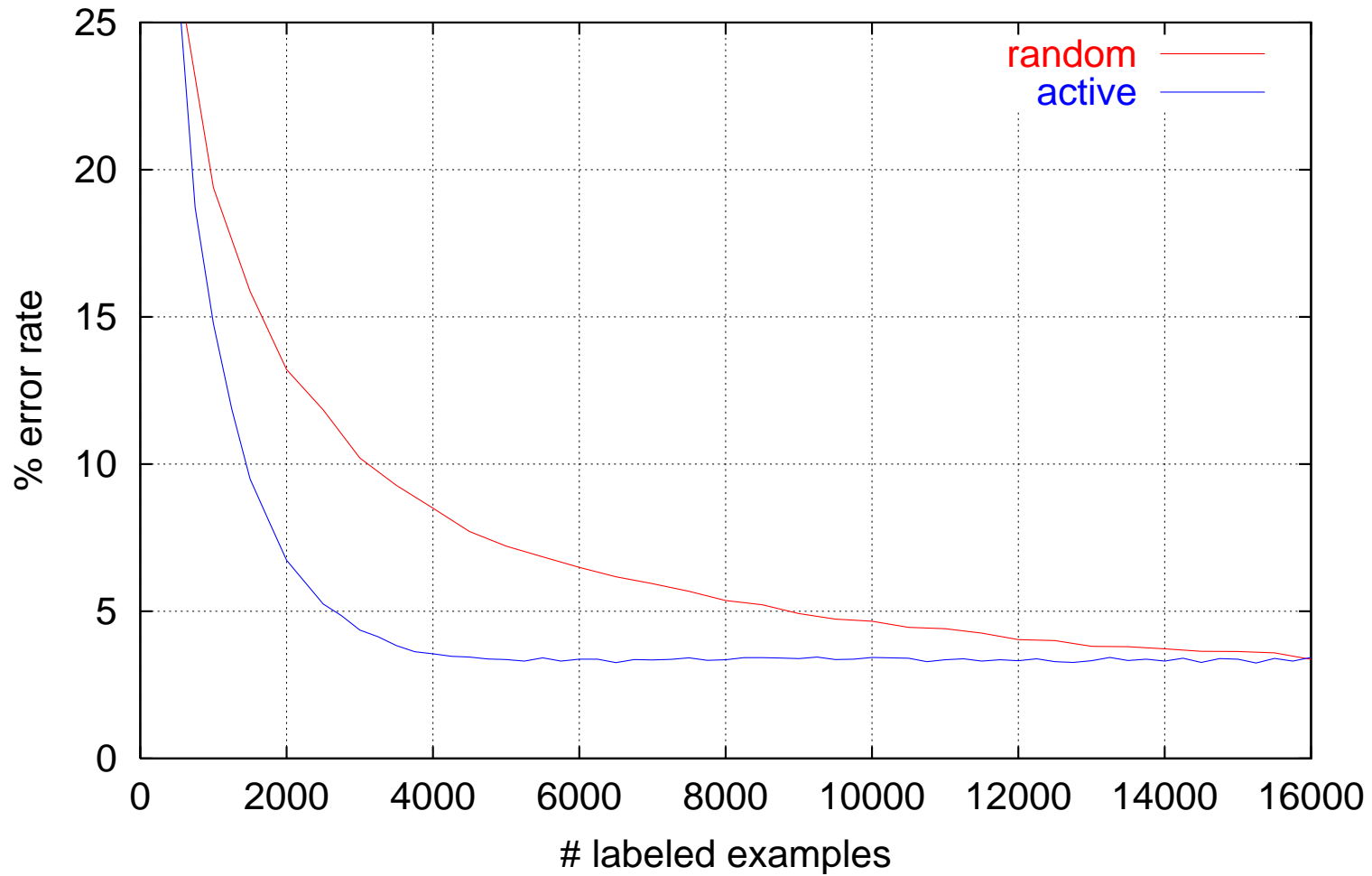[Abe & Mamitsuka]

## Labeling Scheme

- start with pool of unlabeled examples

- choose (say) 500 examples at random for labeling

- run boosting on all labeled examples
  - get combined classifier $f$

- pick (say) 250 additional examples from pool for labeling
  - choose examples with minimum $|f(x)|$

- repeat

# Results: How-May-I-Help-You?



| % error | first reached | | % label |
| | random | active | savings |
|---|---|---|---|
| 28 | 11,000 | 5,500 | 50 |
| 26 | 22,000 | 9,500 | 57 |
| 25 | 40,000 | 13,000 | 68 |

# Results: Letter



| % error | first reached | | % label savings |
|---|---|---|---|
| | random | active | |
| 10 | 3,500 | 1,500 | 57 |
| 5 | 9,000 | 2,750 | 69 |
| 4 | 13,000 | 3,500 | 73 |

# Application: Detecting Faces

[Viola & Jones]

- problem: find faces in photograph or movie

- weak classifiers: detect light/dark rectangles in image



- many clever tricks to make extremely fast and accurate

# Conclusions

- boosting is a practical tool for classification and other learning problems
  - grounded in rich theory
  - performs well experimentally
  - often (but not always!) resistant to overfitting
  - many applications and extensions

- many ways to think about boosting
  - none is entirely satisfactory by itself, but each useful in its own way
  - considerable room for further theoretical and experimental work

# References

Ron Meir and Gunnar Rätsch. An Introduction to Boosting and Leveraging.
   In *Advanced Lectures on Machine Learning (LNAI2600)*, 2003.
   http://www.boosting.org/papers/MeiRae03.pdf

Robert E. Schapire. The boosting approach to machine learning: An overview.
   In *MSRI Workshop on Nonlinear Estimation and Classification*, 2002.
   http://www.cs.princeton.edu/~schapire/boost.html