

An Introduction to Boosting

Yoav Freund

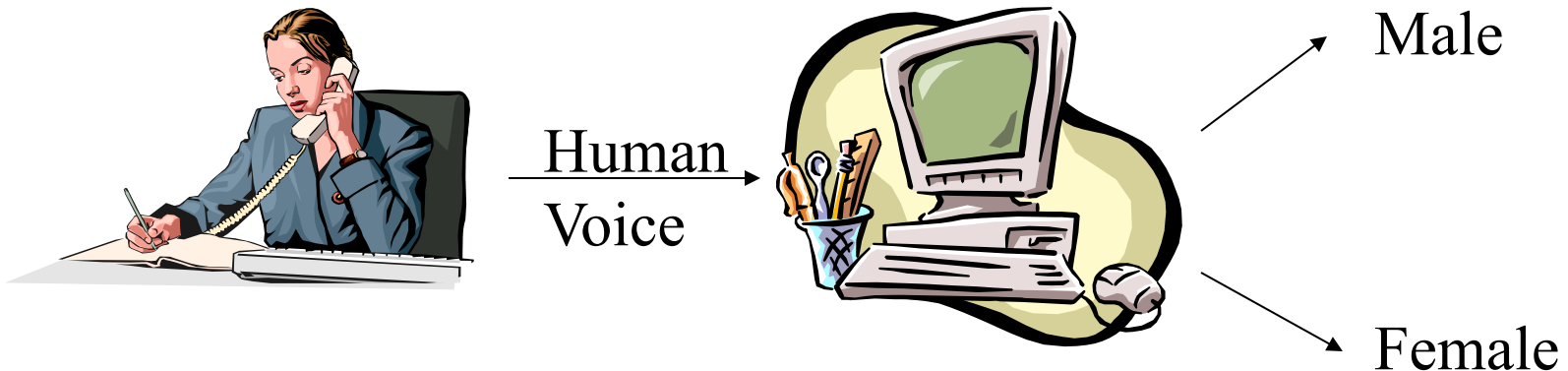
Banter Inc.

Plan of talk

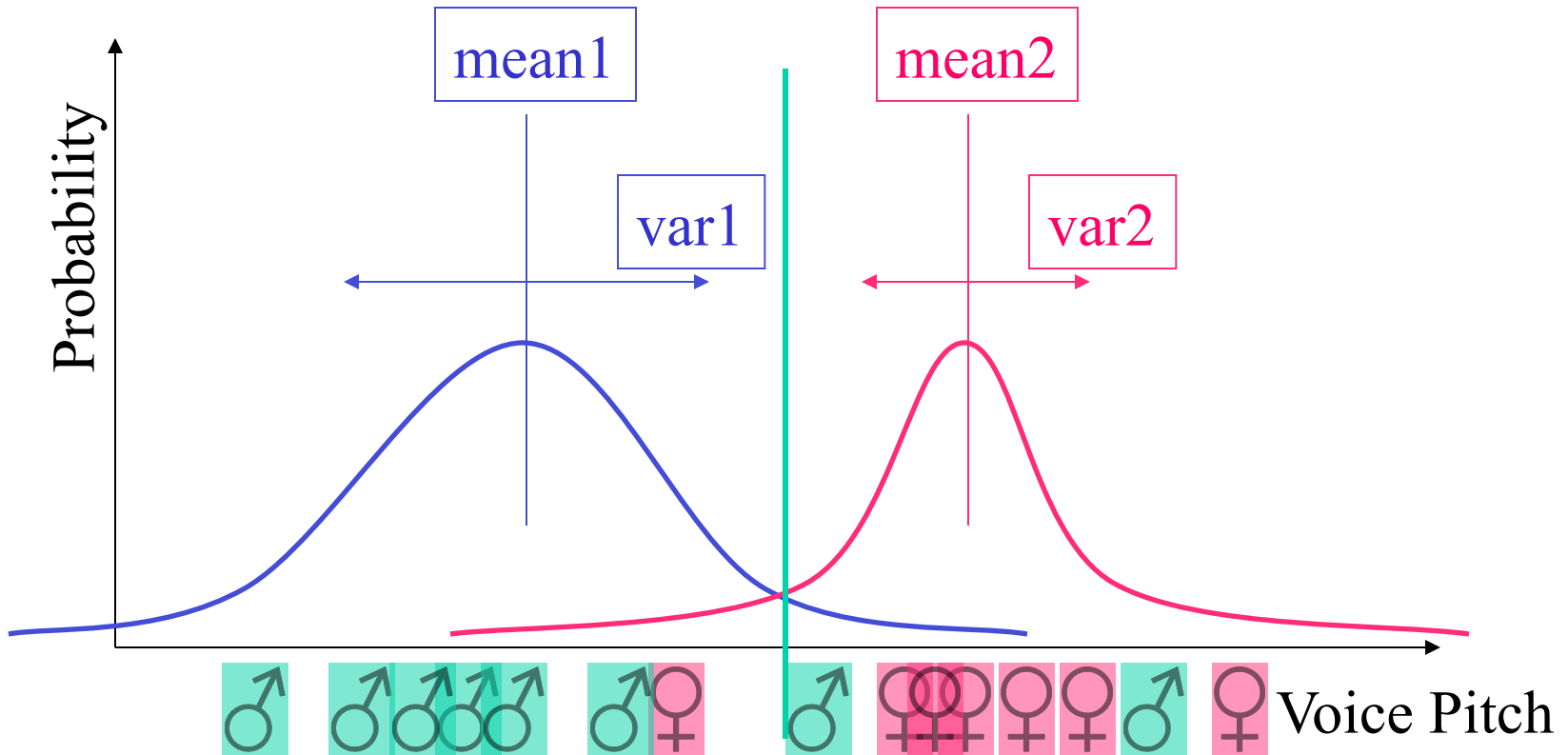
- Generative vs. non-generative modeling
- Boosting
- Alternating decision trees
- Boosting and over-fitting
- Applications

Toy Example

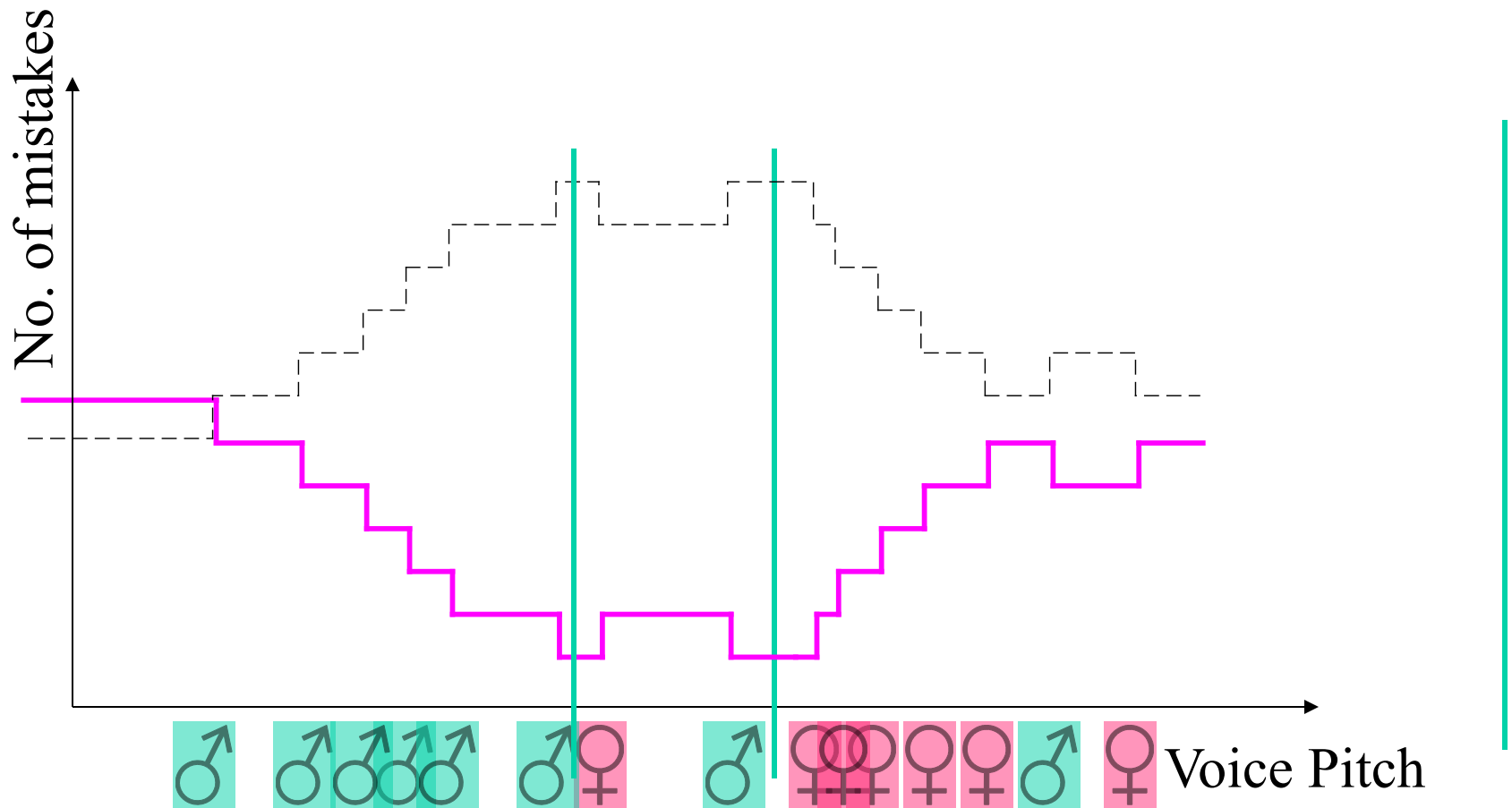
- Computer receives telephone call
- Measures Pitch of voice
- Decides gender of caller



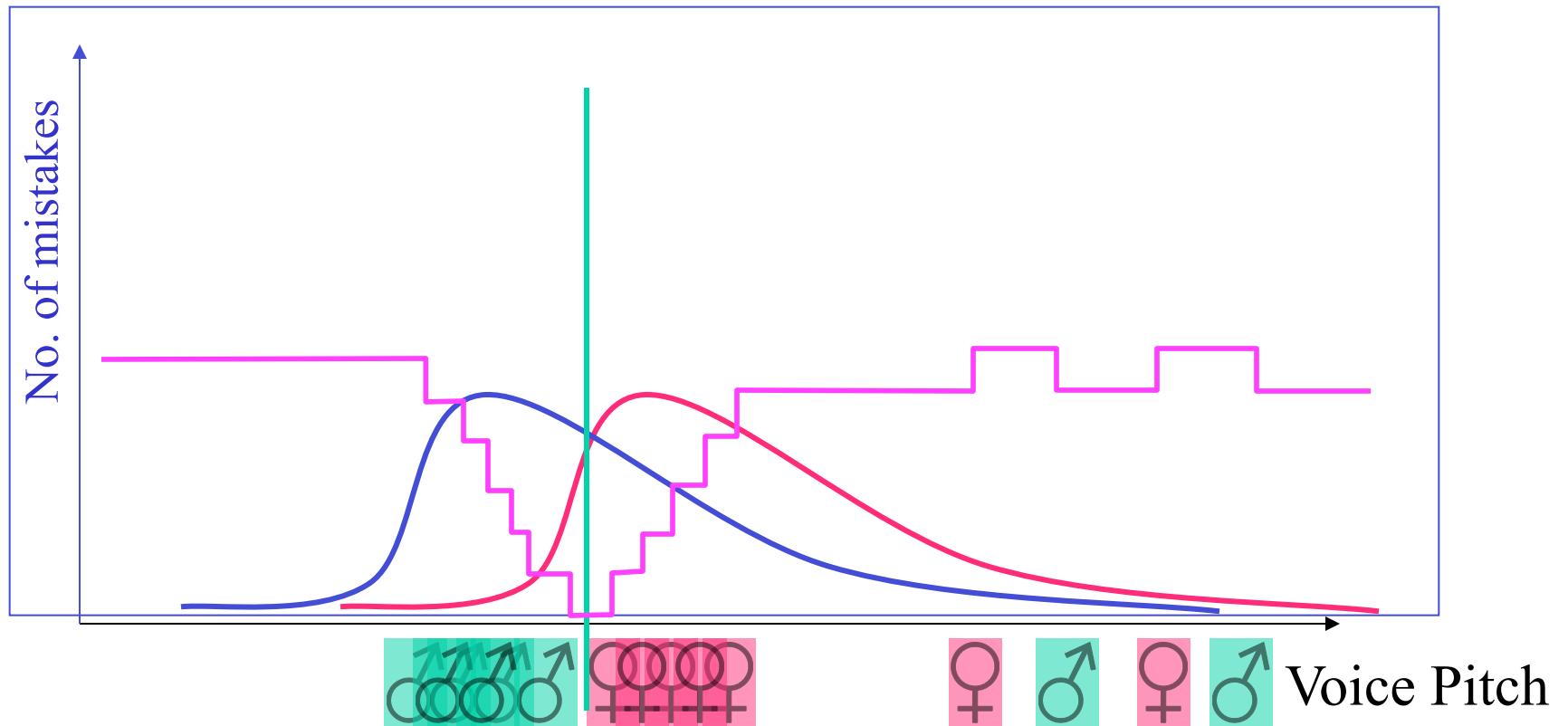
Generative modeling



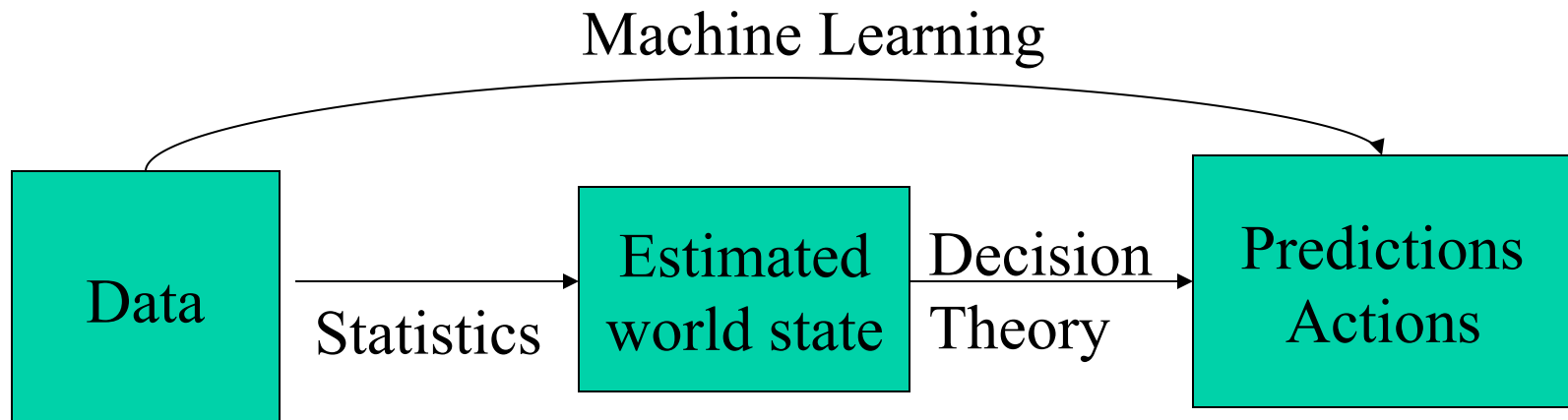
Discriminative approach



Ill-behaved data



Traditional Statistics vs. Machine Learning

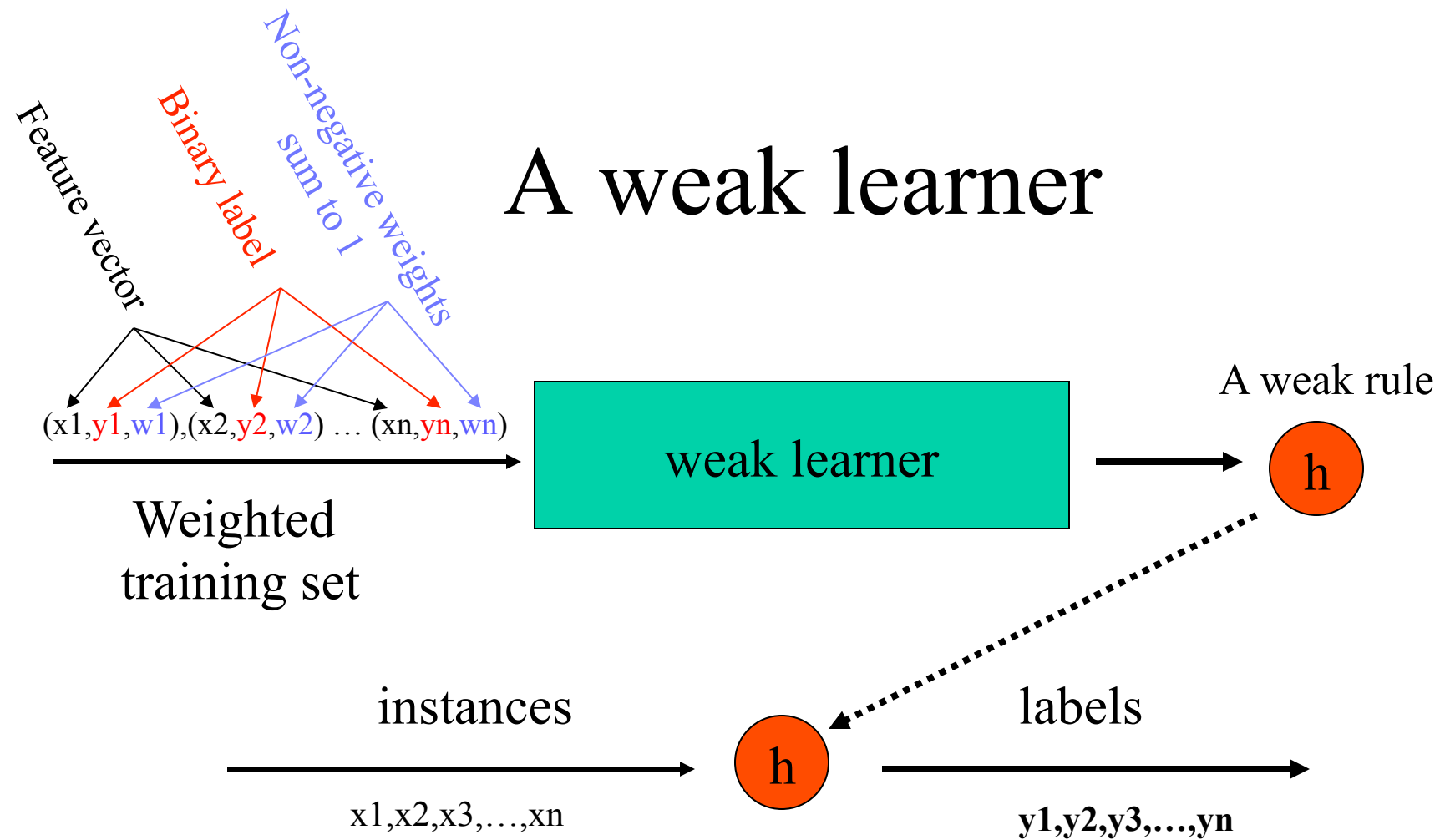


Comparison of methodologies

Model	Generative	Discriminative
Goal	Probability estimates	Classification rule
Performance measure	Likelihood	Misclassification rate
Mismatch problems	Outliers	Misclassifications

Boosting

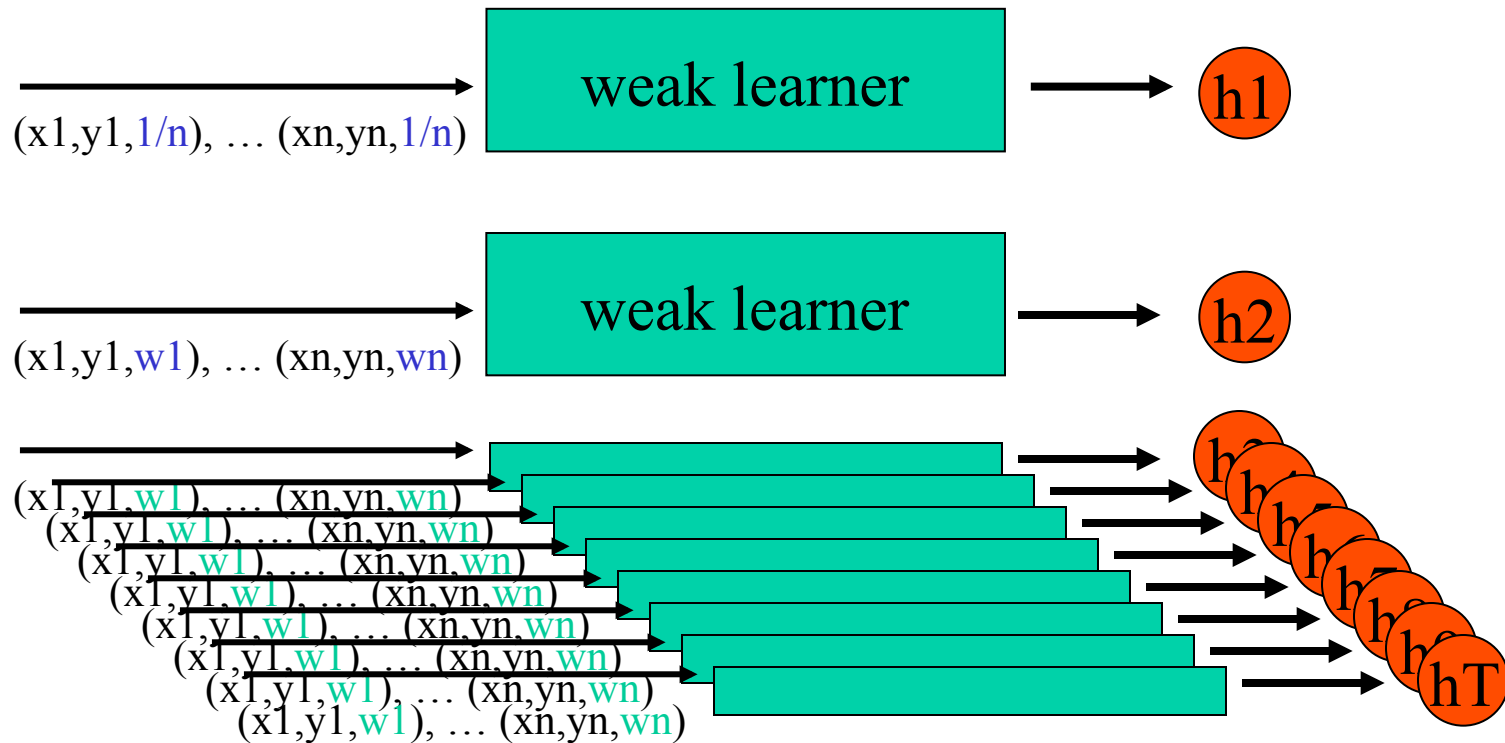
A weak learner



The weak requirement:

$$\frac{\sum_{i: y_i \neq \hat{y}_i} w_i}{\sum_{i=1}^n w_i} < \frac{1}{2} - \gamma$$

The boosting process



Final rule: $\text{Sign}[a_1 h_1 + a_2 h_2 + \dots + a_T h_T]$

Adaboost

- Binary labels $y = -1, +1$
- $\text{margin}(\mathbf{x}, y) = y [\mathbf{S}_t \mathbf{a}_t \mathbf{h}_t(\mathbf{x})]$
- $P(\mathbf{x}, y) = (1/Z) \exp(-\text{margin}(\mathbf{x}, y))$
- Given \mathbf{h}_t , we choose \mathbf{a}_t to minimize
$$\sum_{(\mathbf{x}, y)} \exp(-\text{margin}(\mathbf{x}, y))$$

Main property of adaboost

- If advantages of weak rules over random guessing are: g_1, g_2, \dots, g_T then in-sample error of final rule is at most

$$\exp \left(-2 \sum_{t=1}^T \gamma_t^2 \right)$$

(w.r.t. the initial weights)

Adaboost as gradient descent

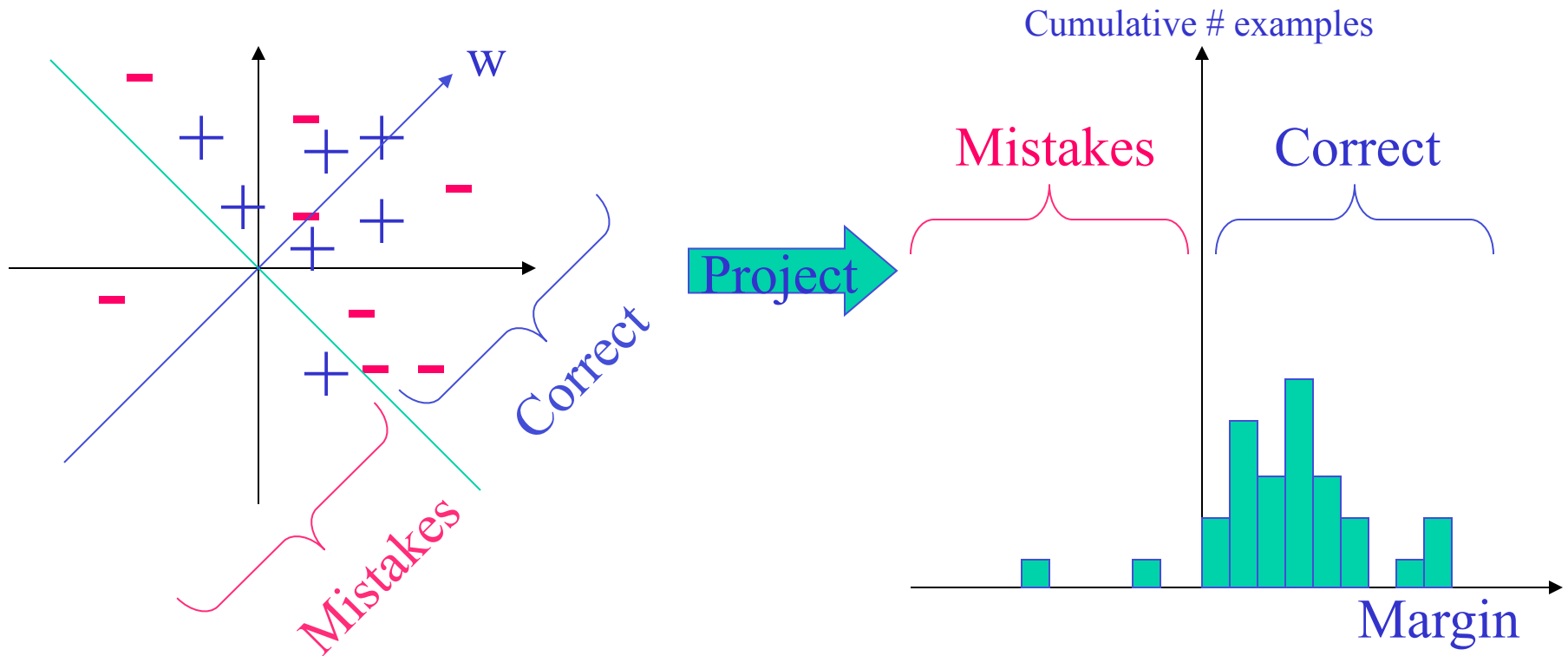
- **Discriminator class:** a linear discriminator in the space of “weak hypotheses”
- **Original goal:** find hyper plane with smallest number of mistakes
 - Known to be an **NP-hard** problem (no algorithm that runs in time polynomial in **d**, where **d is the dimension** of the space)
- **Computational method:** Use exponential loss as a surrogate, perform gradient descent.

Margins view

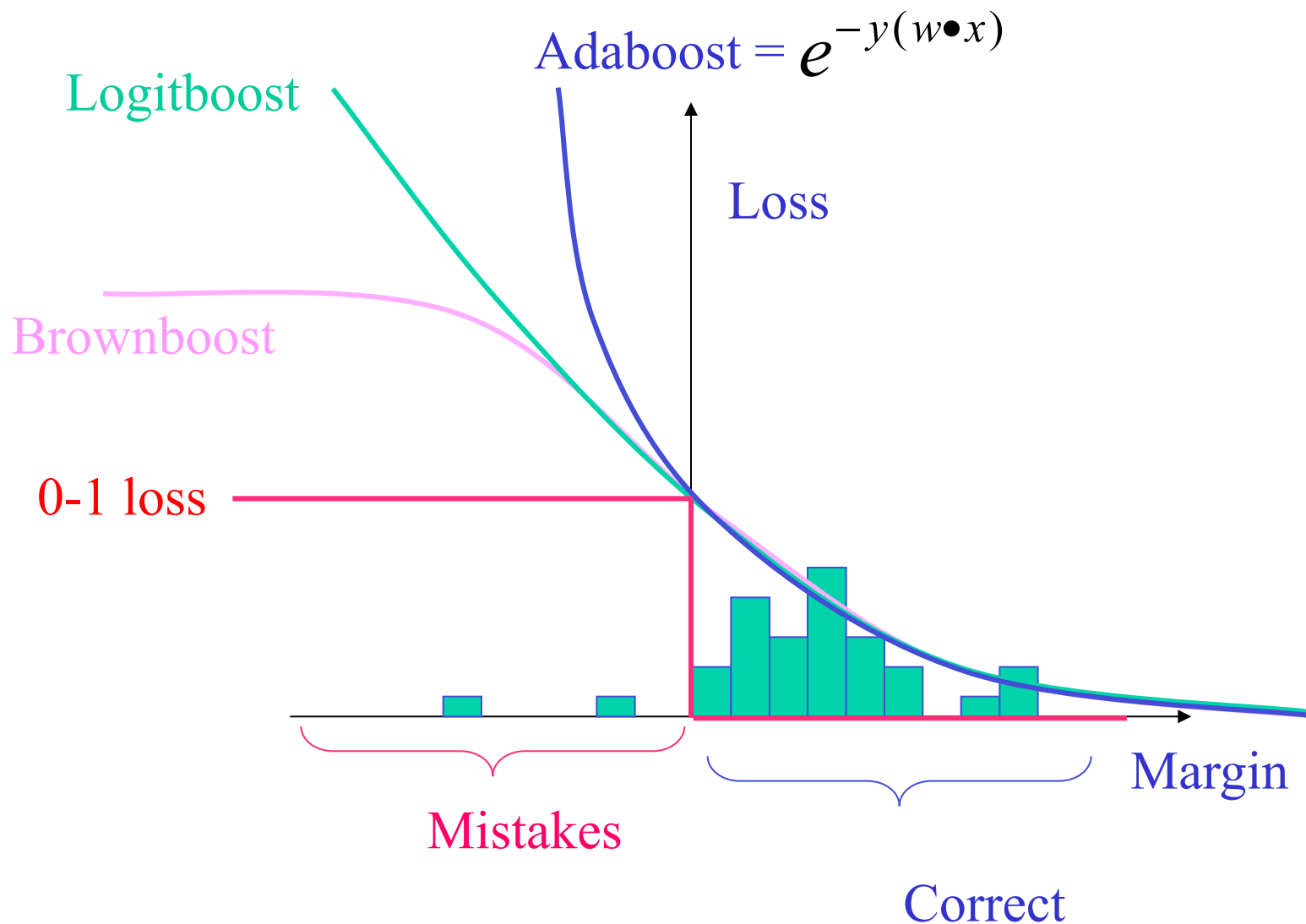
$$x, w \in R^n; y \in \{-1, +1\}$$

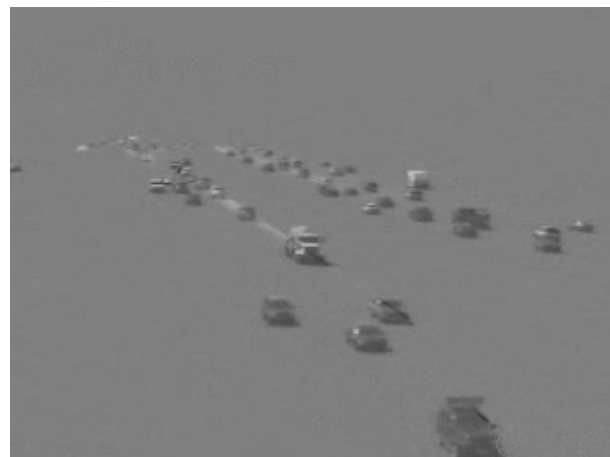
$$\text{Prediction} = \text{sign}(w \bullet x)$$

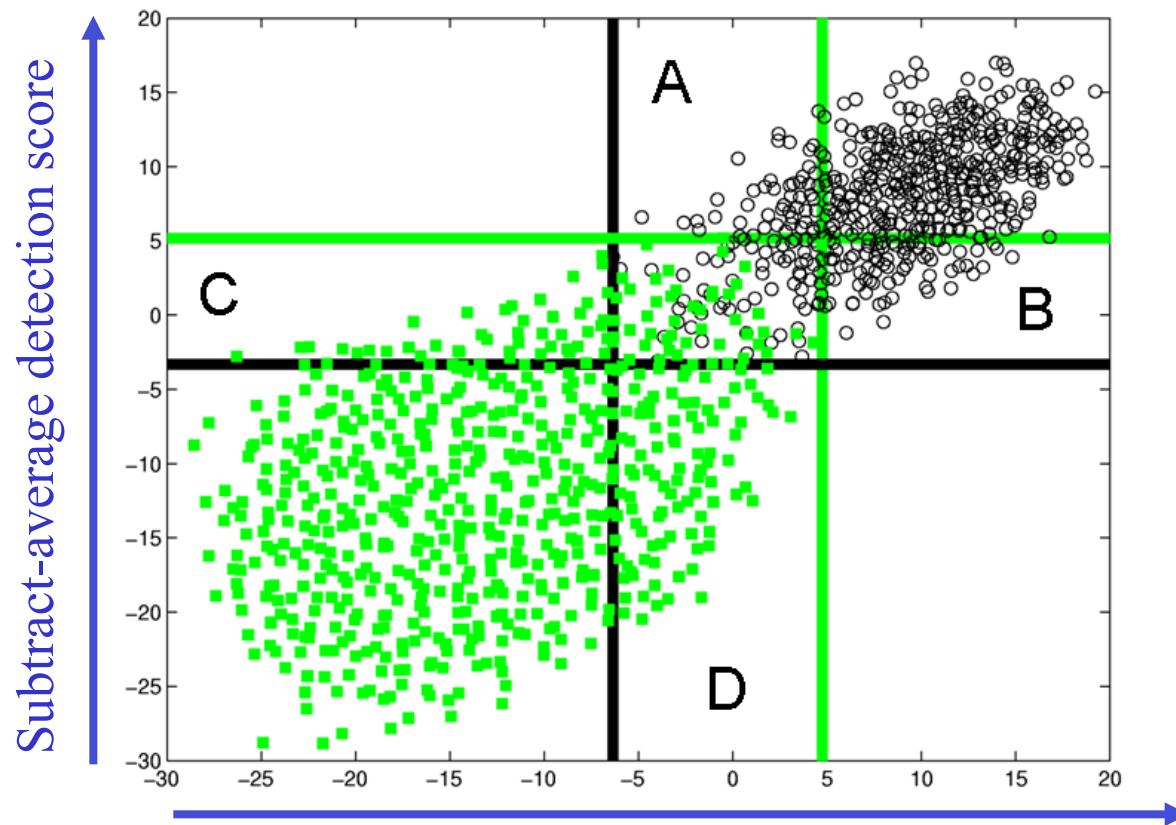
$$\text{Margin} = \frac{y(w \bullet x)}{\|w\| \cdot \|x\|}$$



Adaboost et al.







Grey-scale detection score

One coordinate at a time

- Adaboost performs **gradient descent** on exponential loss
- Adds one coordinate (“**weak learner**”) at each iteration.
- Weak learning in **binary classification** = slightly better than random guessing.
- Weak learning in regression – unclear.
- Uses **example-weights** to communicate the gradient direction to the weak learner
- Solves a **computational** problem

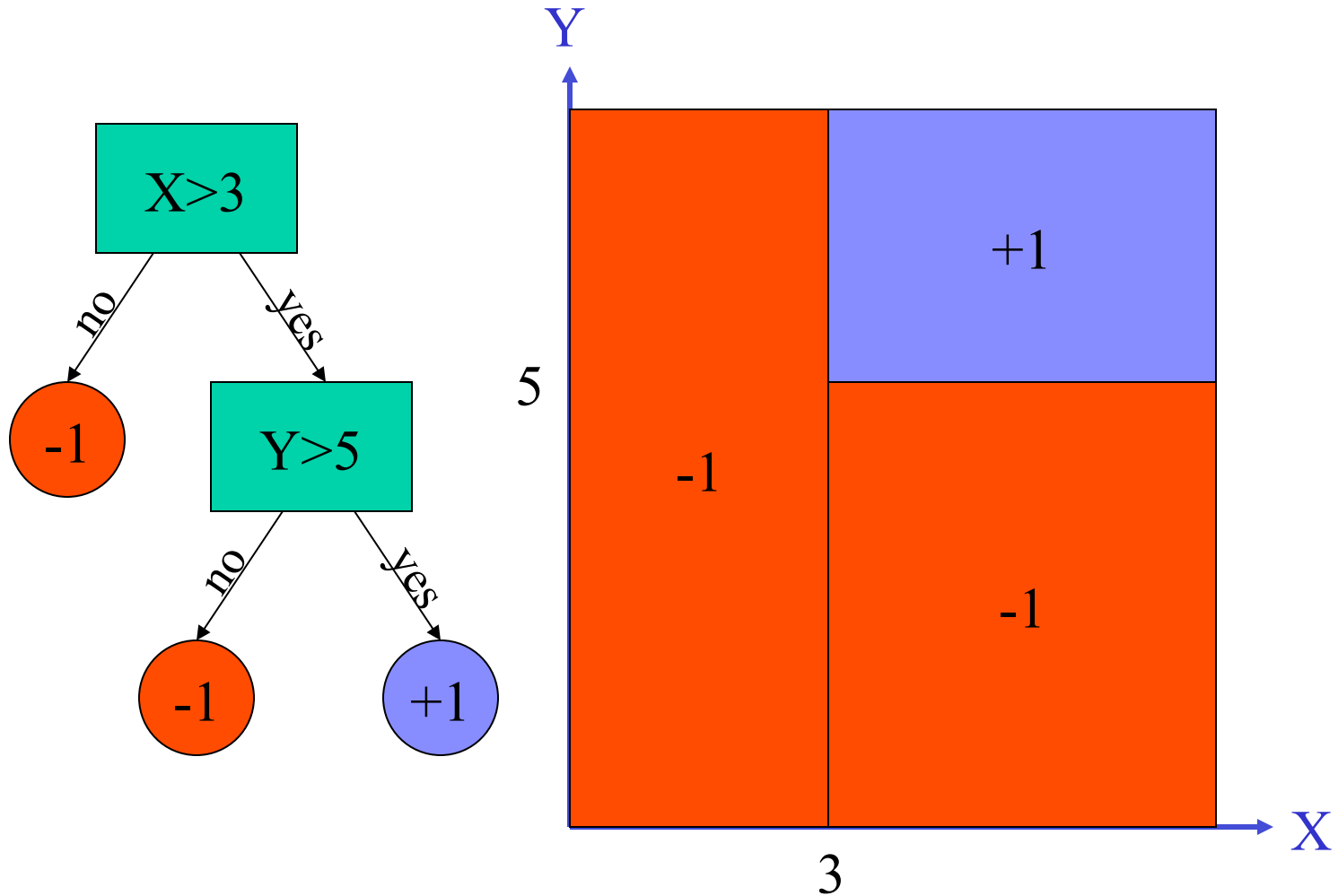
What is a good weak learner?

- The set of weak rules (features) should be **flexible enough to be (weakly) correlated** with most conceivable relations between feature vector and label.
- **Small enough to allow exhaustive search** for the minimal weighted training error.
- **Small enough to avoid over-fitting.**
- Should be able to **calculate predicted label very efficiently.**
- Rules can be “**specialists**” – predict only on a small subset of the input space and **abstain from predicting** on the rest (output 0).

Alternating Trees

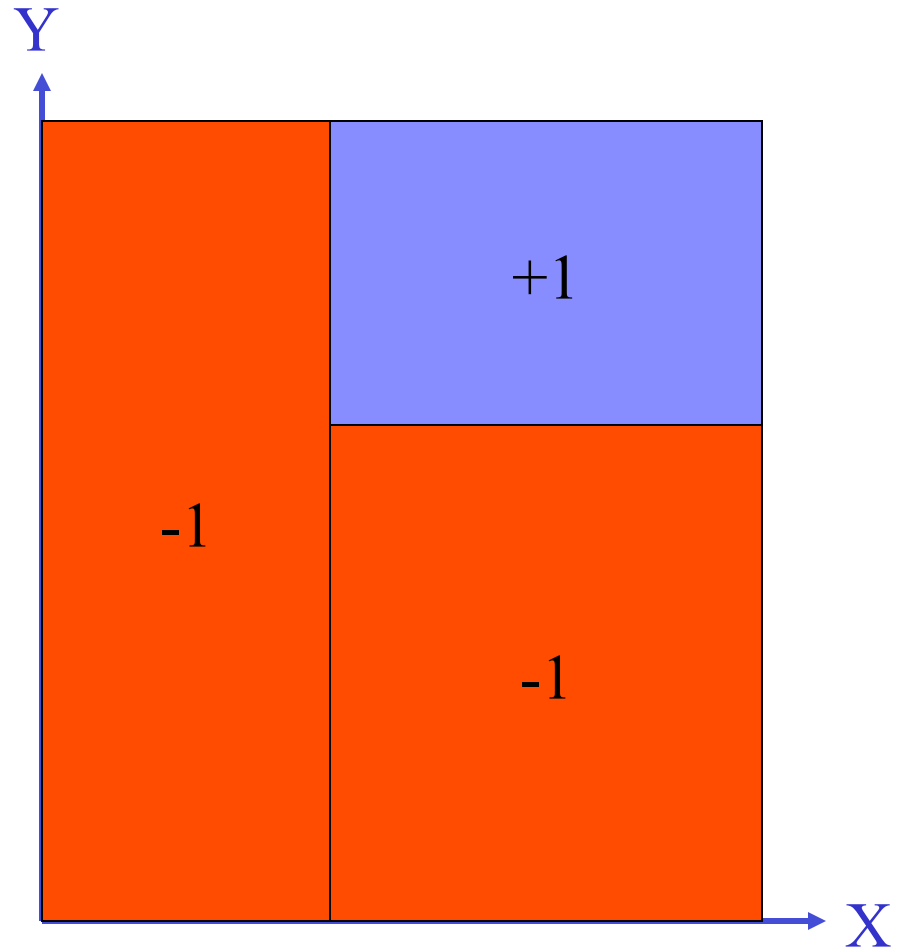
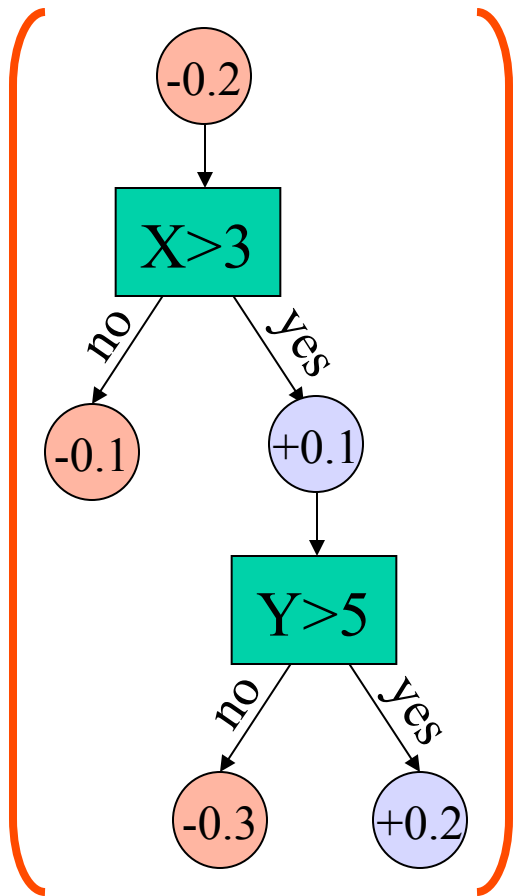
Joint work with Llew Mason

Decision Trees

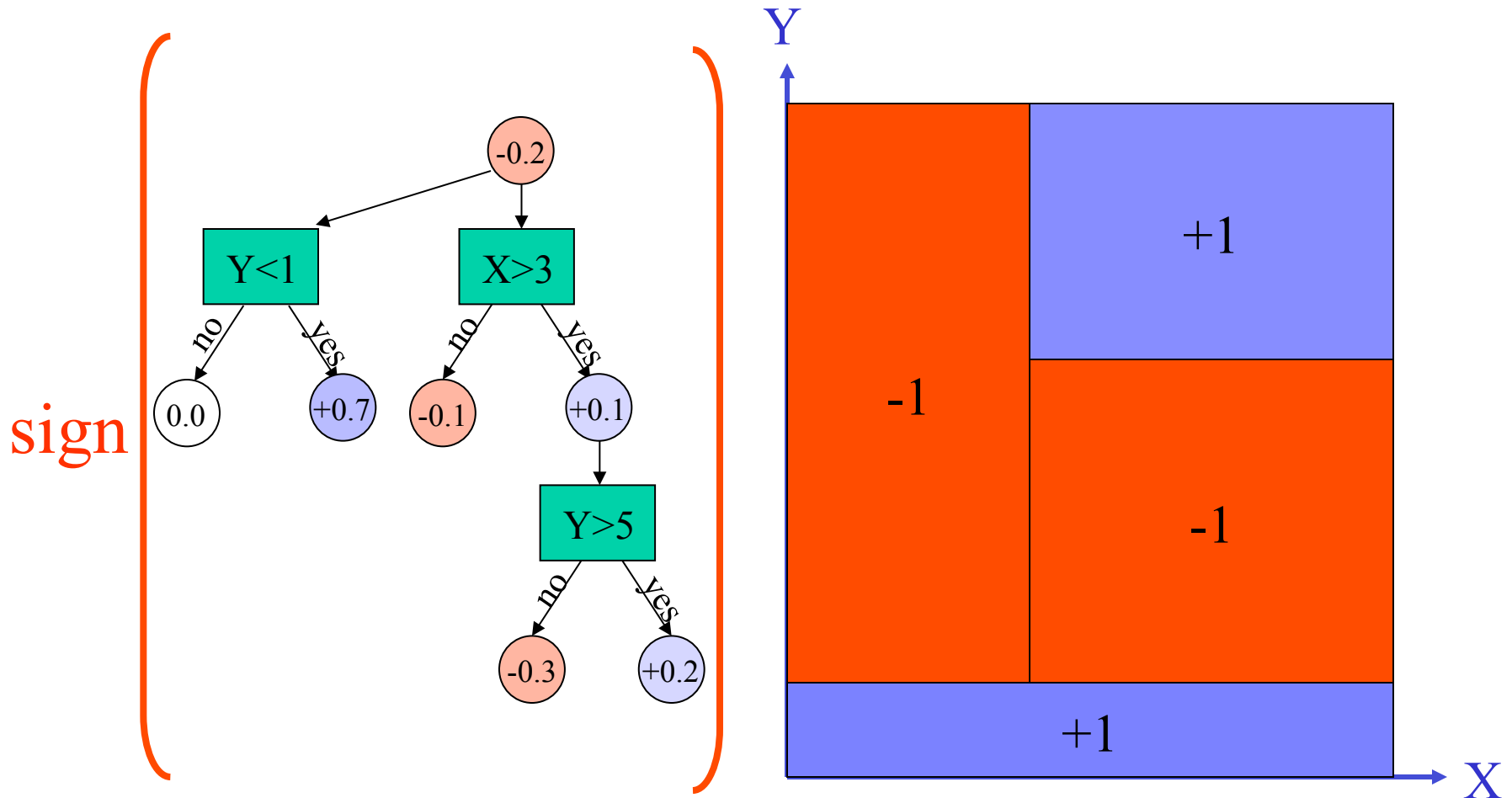


Decision tree as a sum

sign



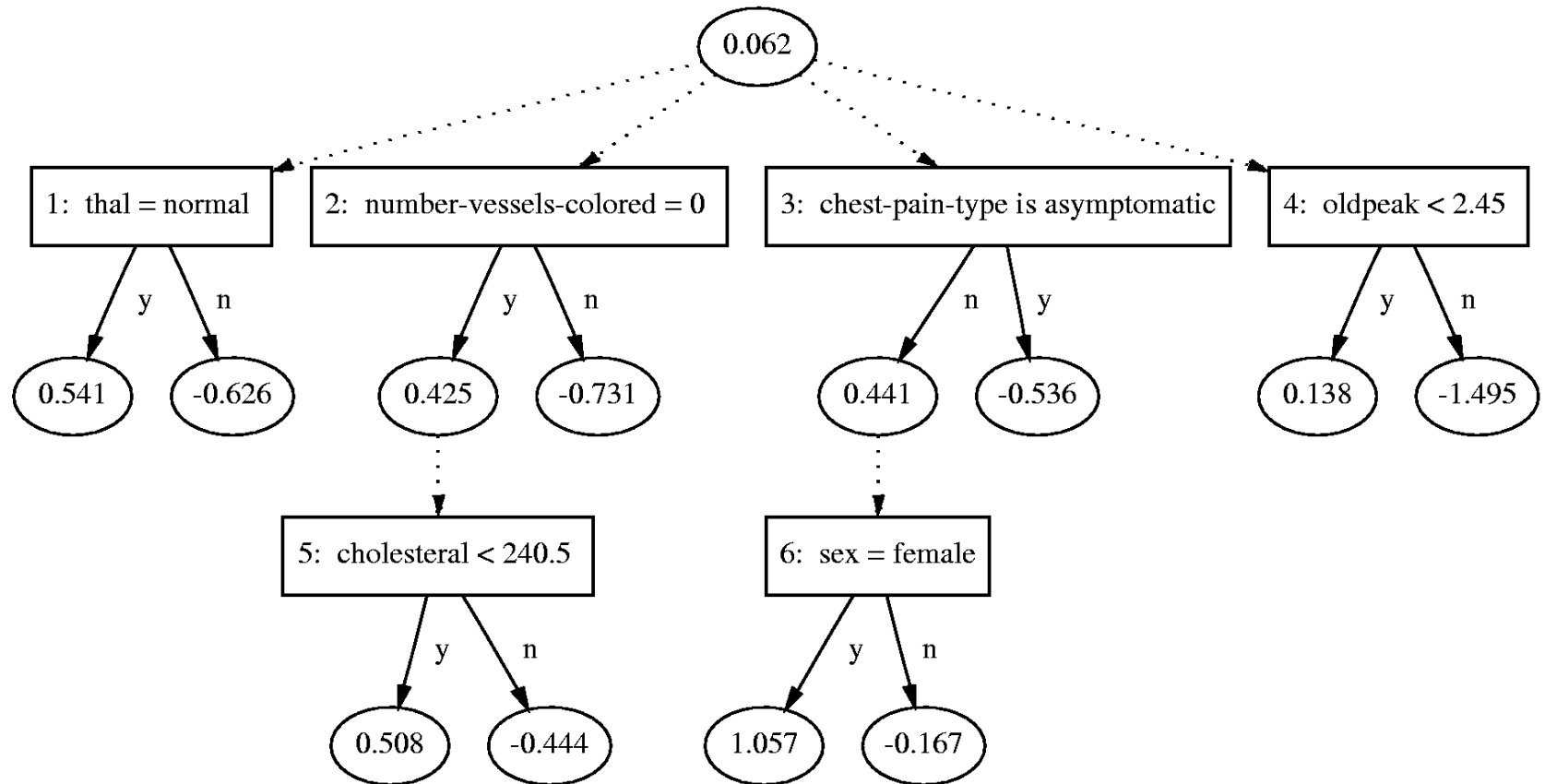
An alternating decision tree



Example: Medical Diagnostics

- **Cleve** dataset from UC Irvine database.
- Heart disease diagnostics (+1=healthy, -1=sick)
- 13 features from tests (real valued and discrete).
- 303 instances.

Adtree for Cleveland heart-disease diagnostics problem



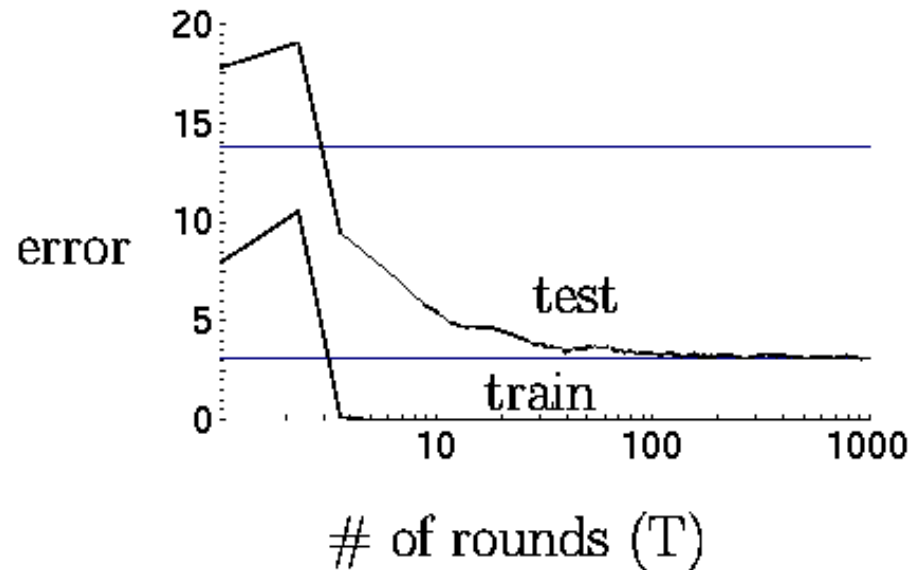
Cross-validated accuracy

Learning algorithm	Number of splits	Average test error	Test error variance
ADtree	6	17.0%	0.6%
C5.0	27	27.2%	0.5%
C5.0 + boosting	446	20.2%	0.5%
Boost Stumps	16	16.5%	0.8%

Boosting and over-fitting

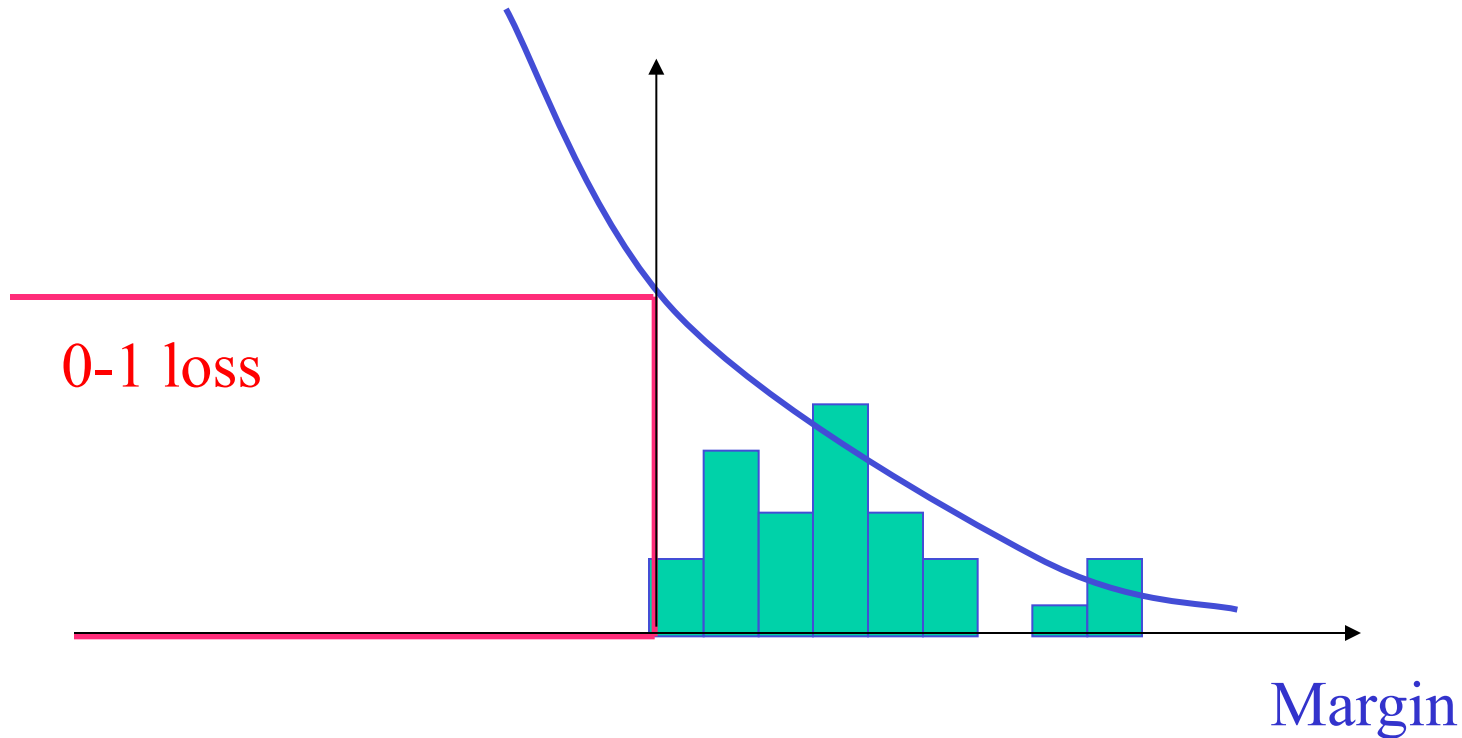
Curious phenomenon

Boosting decision trees

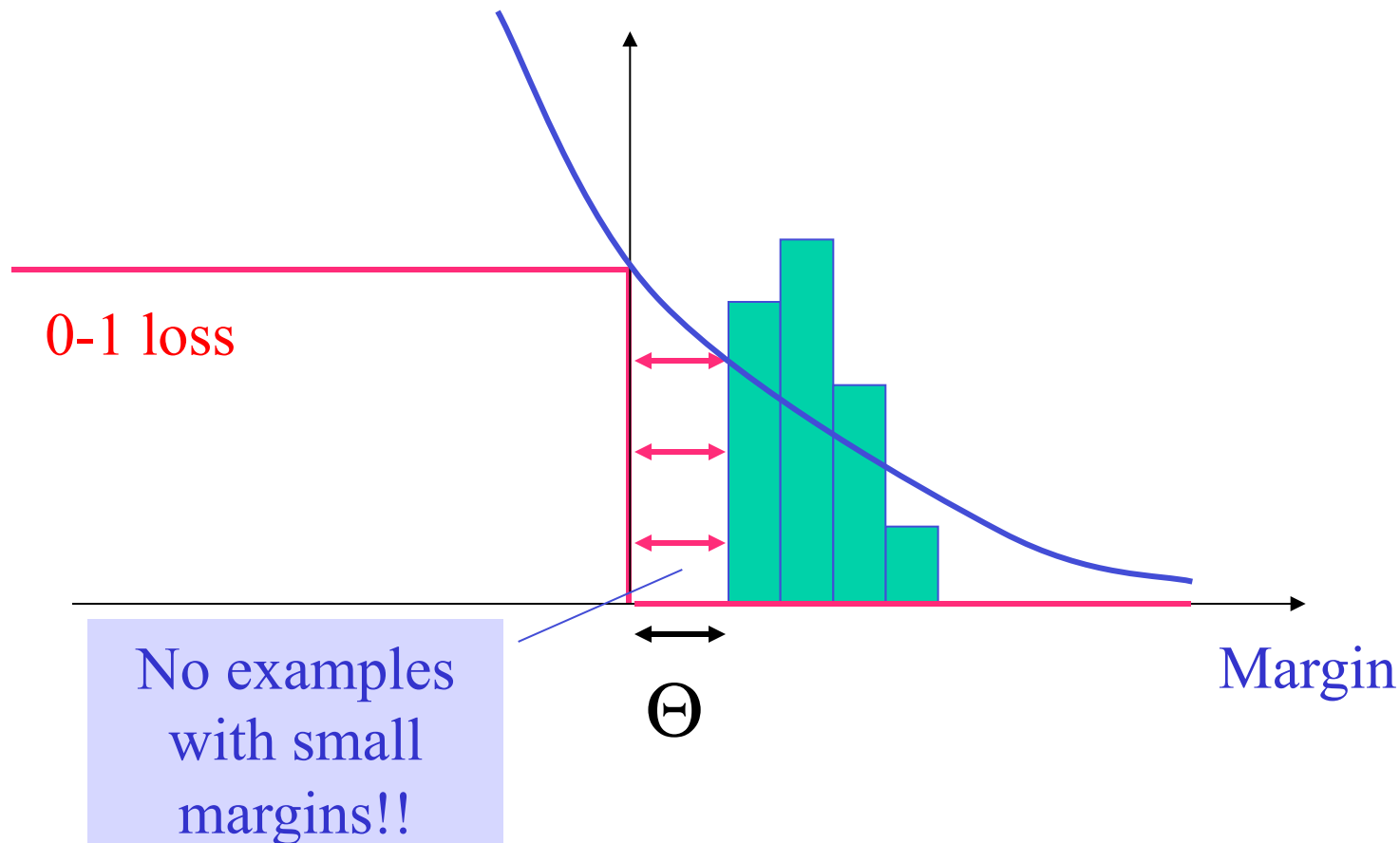


Using $<10,000$ training examples we fit $>2,000,000$ parameters

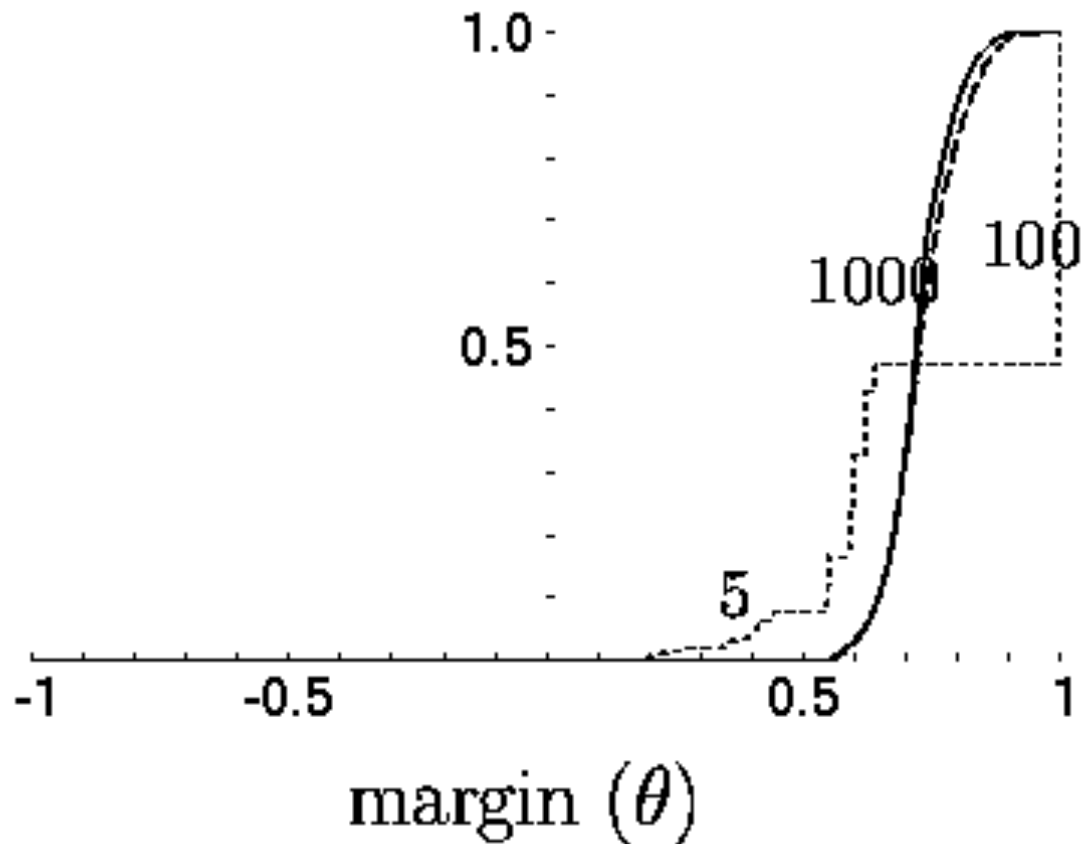
Explanation using margins



Explanation using margins



Experimental Evidence



Theorem

Schapire, Freund, Bartlett & Lee
Annals of stat. 98

For any convex combination and any threshold $\forall f \in \mathcal{C}, \forall \theta > 0$:

Probability of mistake

Fraction of training example
with small margin

$$P_{(x,y) \sim D} [\text{sign}(f(x)) \neq y] \leq P_{(x,y) \sim S} [\text{margin}_f(x, y) \leq \theta]$$

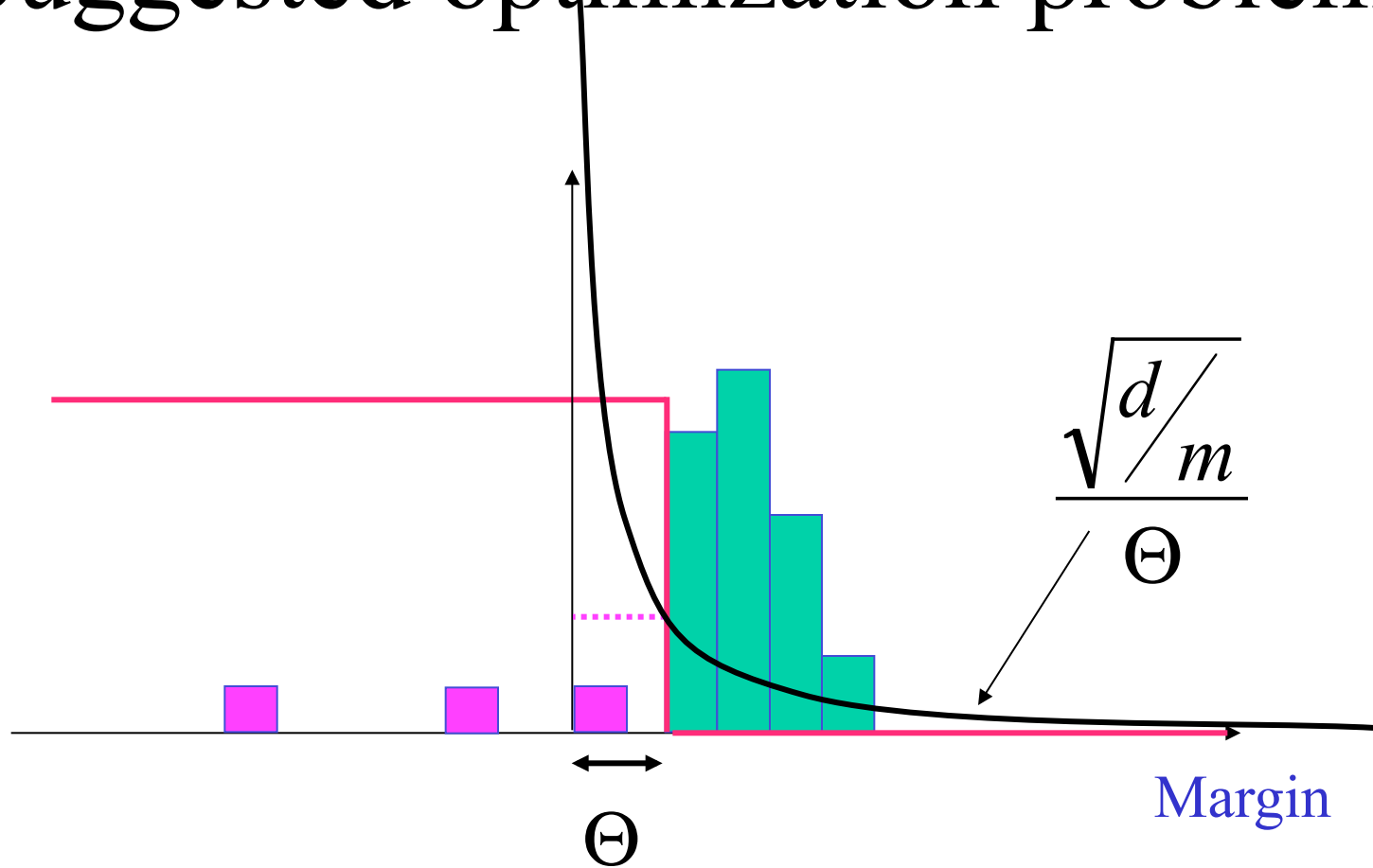
Size of training sample

$$+ \tilde{O} \left(\frac{\sqrt{d/m}}{\theta} \right)$$

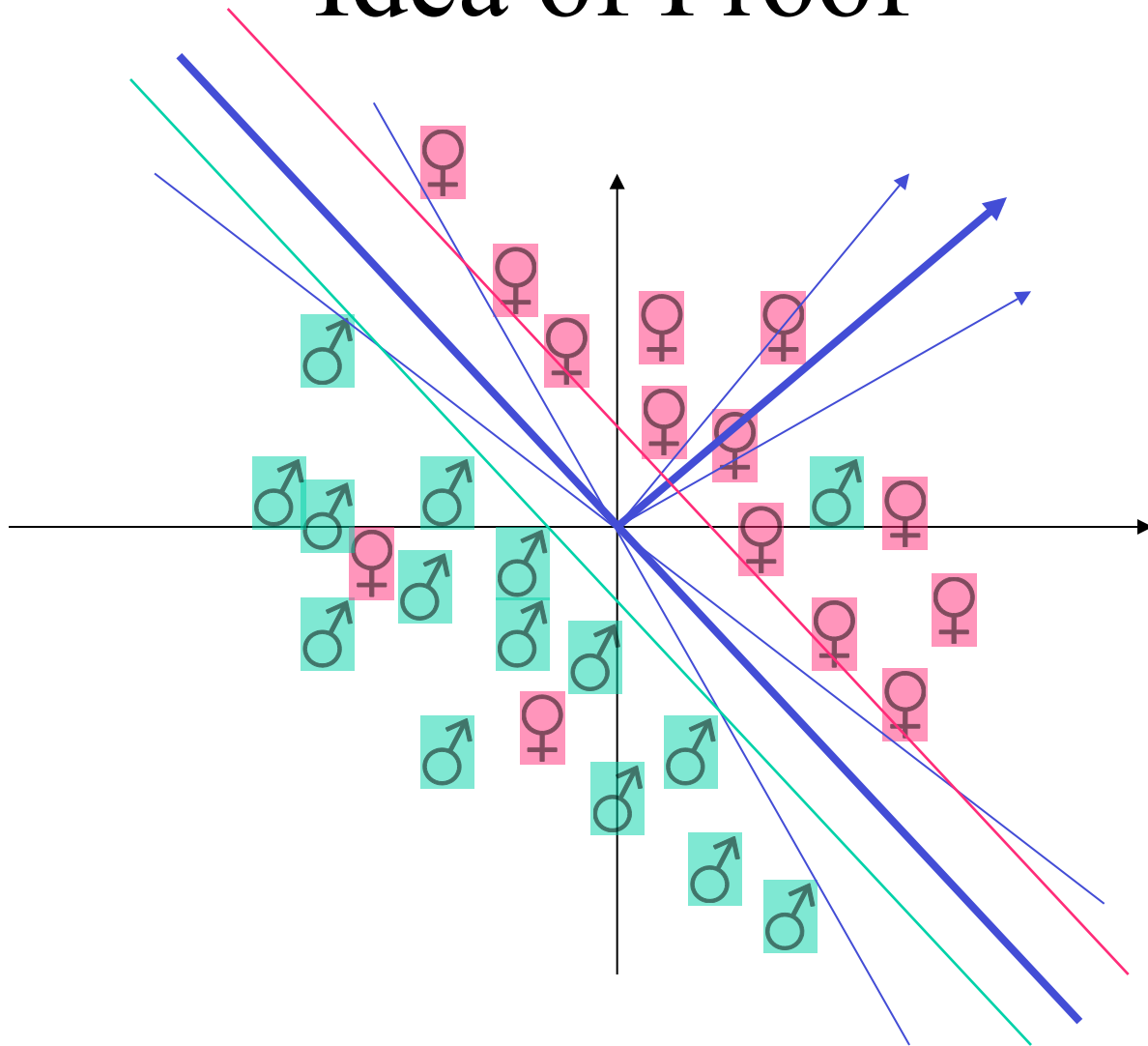
No dependence on number
of weak rules
that are combined!!!

VC dimension of weak rules

Suggested optimization problem



Idea of Proof



Applications

Applications of Boosting

- Academic research
- Applied research
- Commercial deployment

Academic research

% test error rates

Database	Other	Boosting	Error reduction
Cleveland	27.2 (DT)	16.5	39%
Promoters	22.0 (DT)	11.8	46%
Letter	13.8 (DT)	3.5	74%
Reuters 4	5.8, 6.0, 9.8	2.95	~60%
Reuters 8	11.3, 12.1, 13.4	7.4	~40%

Applied research

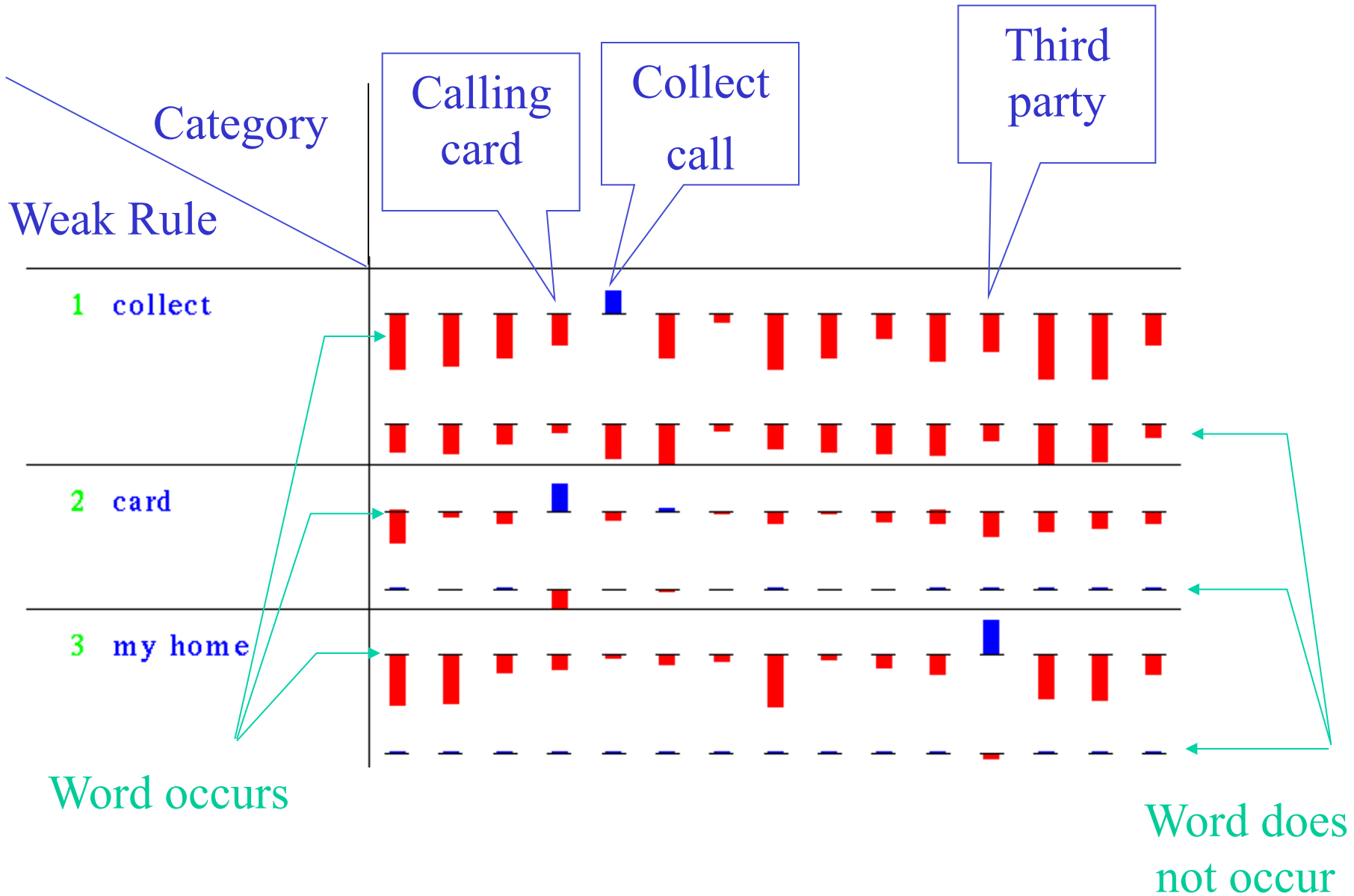
- *“AT&T, How may I help you?”*
- Classify voice requests
- Voice -> text -> category
- Fourteen categories

Area code, AT&T service, billing credit, calling card, collect, competitor, dial assistance, directory, how to dial, person to person, rate, third party, time charge ,time

Examples

- 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
➤ third party
- Yes I' d like to place a call on my master card please
➤ calling card
- 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 my bill
➤ billing credit

Weak rules generated by “boostexter”



Results

- 7844 training examples
 - hand transcribed
- 1000 test examples
 - hand / machine transcribed
- Accuracy with 20% rejected
 - Machine transcribed: 75%
 - Hand transcribed: 90%

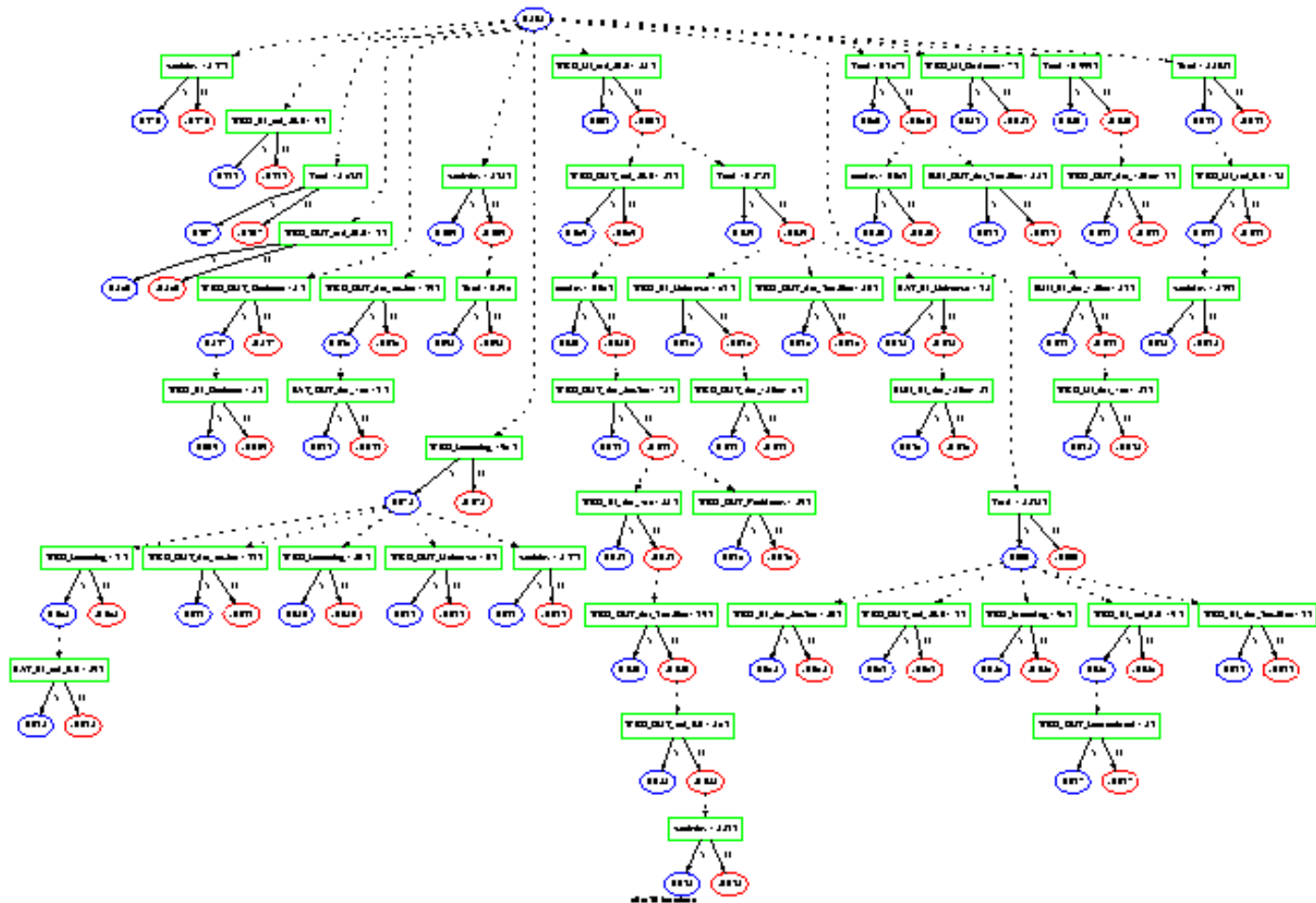
Commercial deployment

- Distinguish business/residence customers
- Using statistics from call-detail records
- Alternating decision trees
 - Similar to boosting decision trees, more flexible
 - Combines very simple rules
 - Can over-fit, cross validation used to stop

Massive datasets

- 260M calls / day
- 230M telephone numbers
- Label unknown for $\sim 30\%$
- Hancock: software for computing statistical signatures.
- 100K randomly selected training examples,
- $\sim 10K$ is enough
- Training takes about 2 hours.
- Generated classifier has to be both accurate and efficient

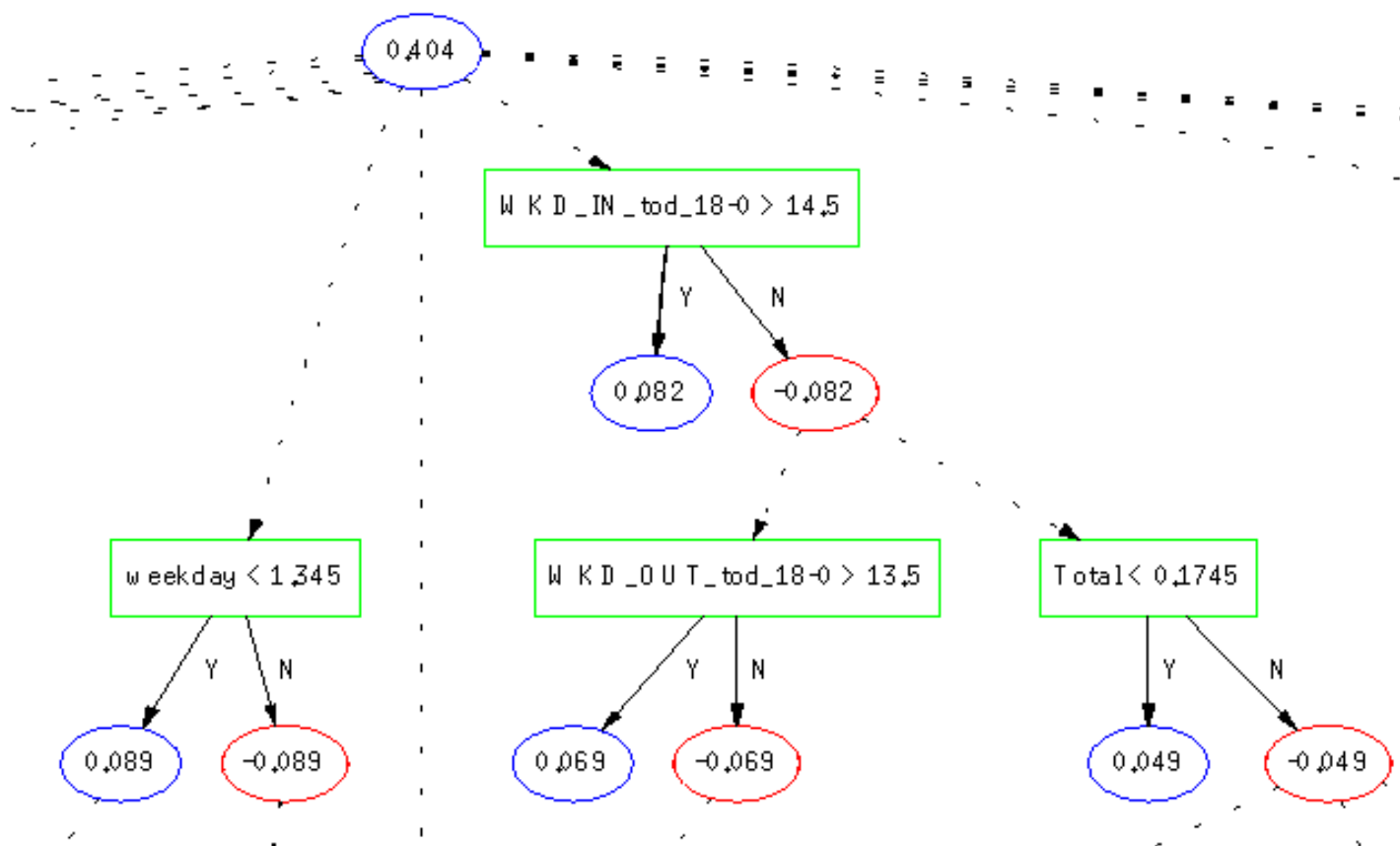
Alternating tree for “buizocity”



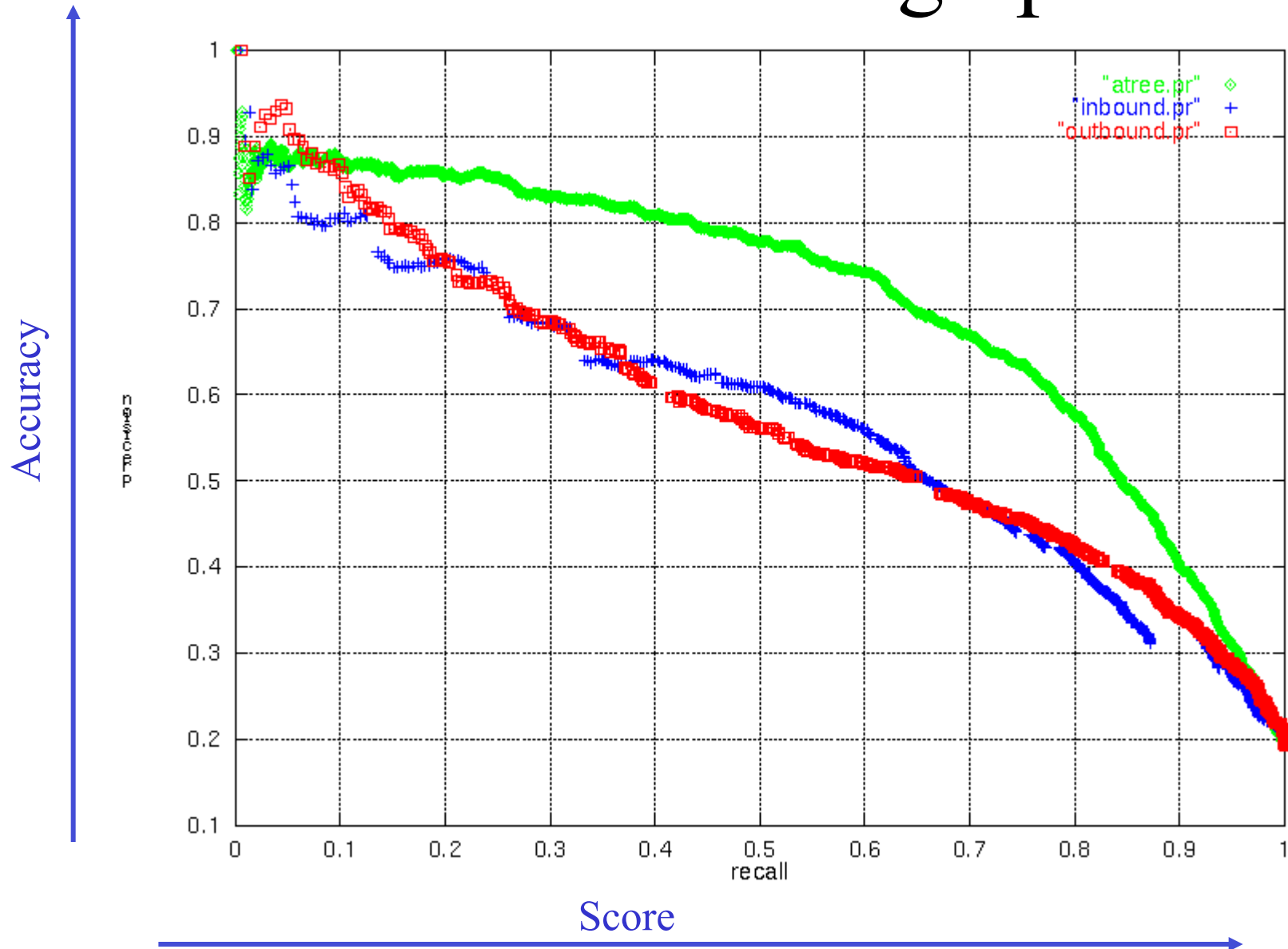
Alternating Tree (Detail)

Positive predictions \Leftrightarrow Residences

Negative predictions \Leftrightarrow Businesses



Precision/recall graphs



Business impact

- Increased coverage from 44% to 56%
- Accuracy ~94%
- Saved AT&T 15M\$ in the year 2000 in operations costs and missed opportunities.

Summary

- Boosting is a computational method for learning accurate classifiers
- Resistance to over-fit explained by margins
- Underlying explanation – large “neighborhoods” of good classifiers
- Boosting has been applied successfully to a variety of classification problems

Come talk with me!

- Yfreund@banter.com
- <http://www.cs.huji.ac.il/~yoavf>