

From AdaBoost to LDM

Zhi-Hua Zhou

<http://cs.nju.edu.cn/zhouzh/>

Email: zhouzh@nju.edu.cn

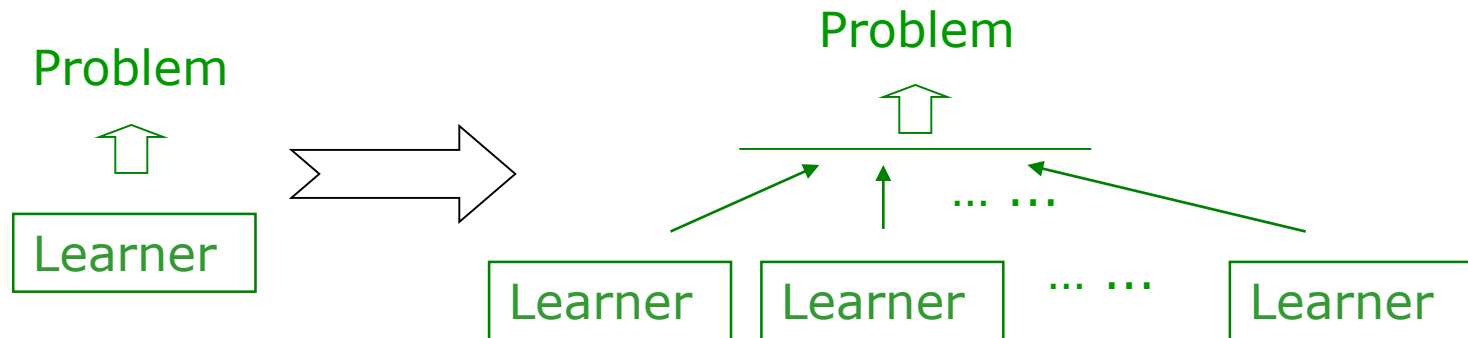
LAMDA Group

National Key Laboratory for Novel Software Technology,
Nanjing University, China



Ensemble learning (集成学习)

“Ensemble methods” is a machine learning paradigm where multiple (homogenous/heterogeneous) individual learners are trained for the same problem
e.g. neural network ensemble, decision tree ensemble, etc.



The more **accurate** and **diverse** the component learners, the better the ensemble

Great success of ensemble methods

- ❑ KDDCup'07: 1st place for "... Decision Forests and ..."
- ❑ KDDCup'08: 1st place of Challenge1 for a method using Bagging; 1st place of Challenge2 for "... Using an Ensemble Method "
- ❑ KDDCup'09: 1st place of Fast Track for "Ensemble ... "; 2nd place of Fast Track for "... bagging ... boosting tree models ...", 1st place of Slow Track for "Boosting ... "; 2nd place of Slow Track for "Stochastic Gradient Boosting"
- ❑ KDDCup'10: 1st place for "... Classifier ensembling"; 2nd place for "... Gradient Boosting machines ... "

Great success of ensemble methods (cont')

- **KDDCup'11**: 1st place of Track 1 for "A Linear Ensemble ..."; 2nd place of Track 1 for "Collaborative filtering Ensemble"; 1st place of Track 2 for "Ensemble ..."; 2nd place of Track 2 for "Linear combination of ..."
- **KDDCup'12**: 1st place of Track 1 for "Combining... Additive Forest..."; 1st place of Track 2 for "A Two-stage Ensemble of..."
- **KDDCup'13**: 1st place of Track 1 for "Weighted Average Ensemble"; 2nd place of Track 1 for "Gradient Boosting Machine"; 1st place of Track 2 for "Ensemble the Predictions"

Great success of ensemble methods (cont')

- ❑ KDDCup'14: 1st place for “ensemble of GBM, ExtraTrees, Random Forest...” and “the weighted average” ; 2nd place for “use both R and Python GBMs”; 3rd place for “gradient boosting machines... random forests” and “the weighted average of...”

- ❑ Netflix Prize:
 - ✓ 2007 Progress Prize Winner: Ensemble
 - ✓ 2008 Progress Prize Winner: Ensemble
 - ✓ 2009 \$1 Million Grand Prize Winner:

Ensemble !!

Many effective ensemble methods

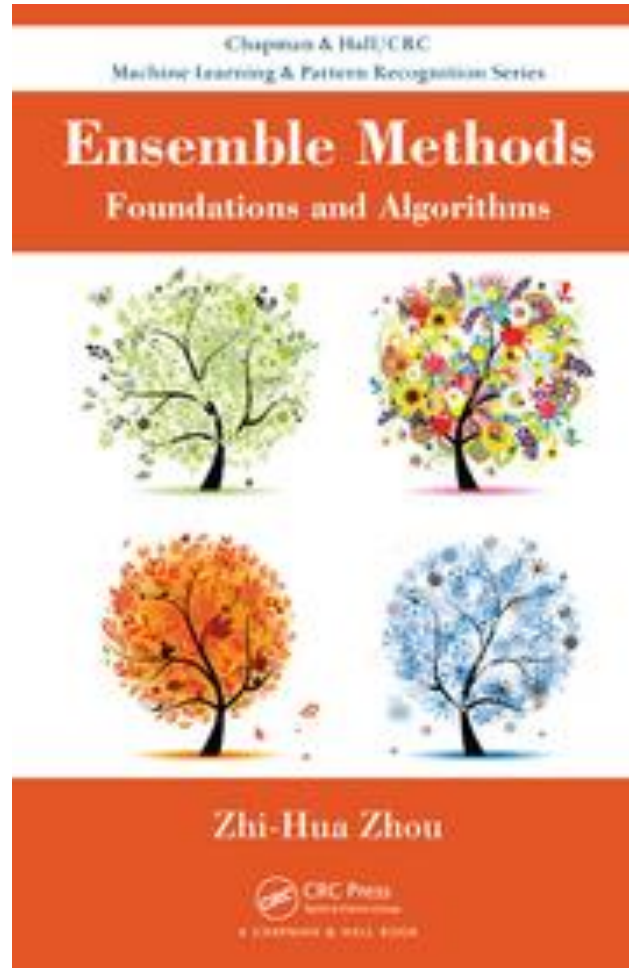
■ Sequential methods

- **AdaBoost** [Freund & Schapire, JCSS97]
- Arc-x4 [Breiman, AnnStat98]
- LPBoost [Demiriz, Bennett, Shawe-Taylor, MLJ06]
-

■ Parallel methods

- Bagging [Breiman, MLJ96]
- Random Subspace [Ho, TPAMI98]
- Random Forests [Breiman, MLJ01]
-

More about ensemble methods



Z.-H. Zhou.
Ensemble Methods:
Foundations and Algorithms,
Boca Raton, FL: Chapman &
Hall/CRC, Jun. 2012.
(ISBN 978-1-439-830031)

Special focus of this talk: AdaBoost

Significant advantageous:

- Very accurate prediction
- Very simple ("*just 10 lines of code*" as Schapire said)
- Wide and successful applications
- Sound theoretical foundation
-



Gödel Prize (2003)

Freund & Schapire, A decision theoretic generalization of on-line learning and an application to Boosting. *Journal of Computer and System Sciences*, 1997, 55: 119-139.

Outline

- **A long march of margin theory for Boosting**
(a theory breakthrough)

- **LDM: Large margin Distribution Machine**
(an algorithmic paradigm)

This talk will show how a theory result on Boosting gives born to a powerful general learning paradigm

The born of AdaBoost

An open problem [Kearns & Valiant, STOC'89]:
“weakly learnable” $\stackrel{?}{=}$ “strongly learnable”

a problem is *learnable* or *strongly learnable* if there exists an algorithm that outputs a learner h in polynomial time such that for all $0 < \delta, \epsilon \leq 0.5$, $P(\mathbb{E}_{\mathbf{x} \sim \mathcal{D}}[\mathbb{I}[h(\mathbf{x}) \neq f(\mathbf{x})]] < \epsilon) \geq 1 - \delta$

a problem is *weakly learnable* if there exists an algorithm that outputs a learner with error $0.5 - 1/p$ where p is a polynomial in problem size and other parameters

In other words, whether a “weak” learning algorithm that works just slightly better than random guess can be “boosted” into an arbitrarily accurate “strong” learning algorithm

The born of AdaBoost (con't)

- Amazingly, in 1990 Schapire proves that the answer is “yes”. More importantly, the proof is a construction!

This is the first Boosting algorithm

- In 1993, Freund presents a scheme of combining weak learners by majority voting in Phd thesis at UC Santa Cruz

However, these algorithms are not practical

- Later, at AT&T Bell Labs, Freund & Schapire published **the 1997 journal paper** (the work was reported in EuroCOLT'95), **which proposed the AdaBoost algorithm**, a practical algorithm

The AdaBoost algorithm

Given: $(x_1, y_1), \dots, (x_m, y_m)$ where $x_i \in X, y_i \in Y = \{-1, +1\}$

Initialize $D_1(i) = 1/m$.

For $t = 1, \dots, T$:

- Train base learner using distribution D_t .
- Get base classifier $h_t : X \rightarrow \mathbb{R}$.
- Choose $\alpha_t \in \mathbb{R}$.
- Update:

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

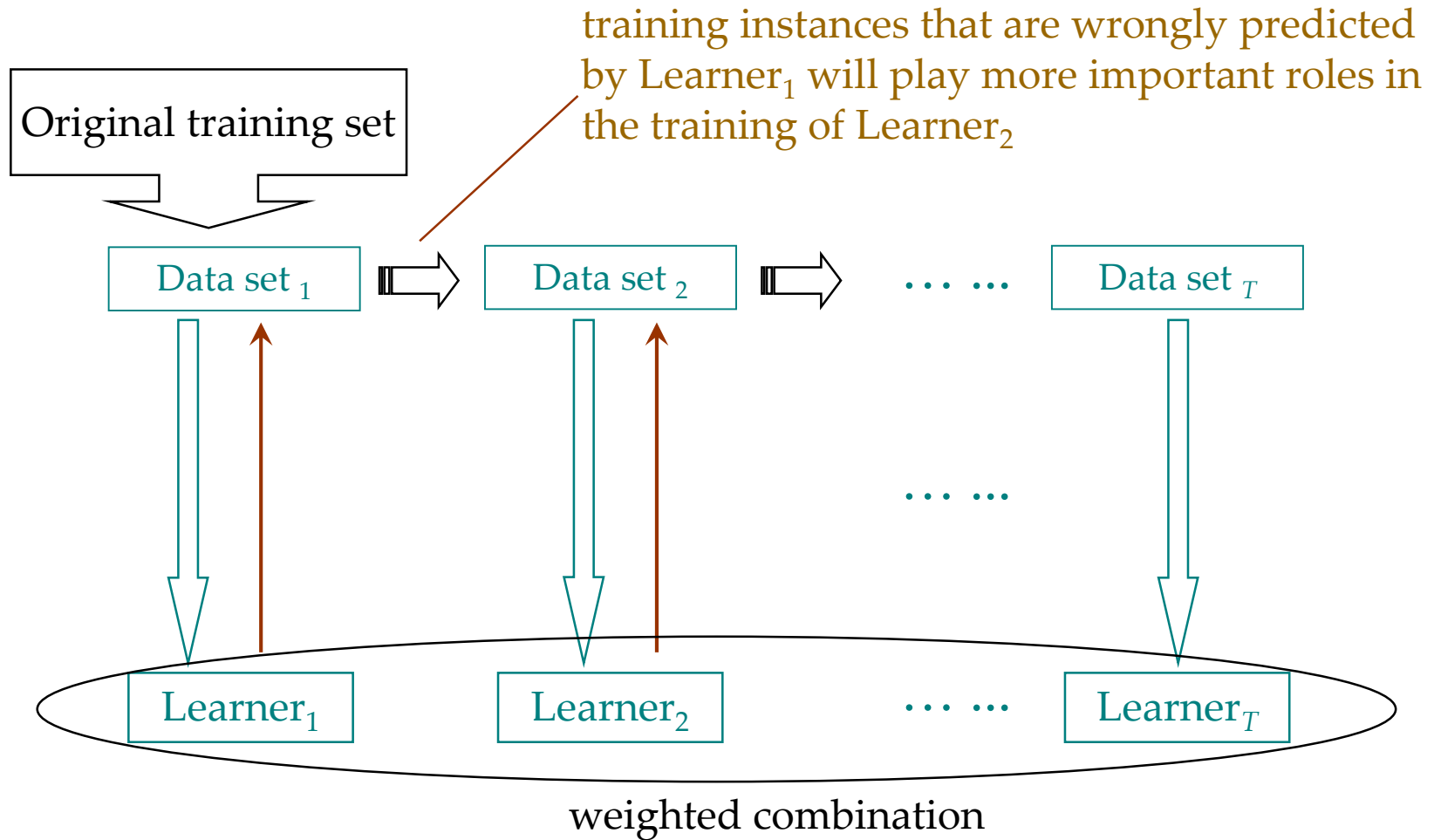
where Z_t is a normalization).

the weights of incorrectly classified examples are increased such that the base learner is forced to focus on the “hard” examples in the training set

Output the final classifier:

$$H(x) = \text{sign} \left(\sum_{t=1}^T \alpha_t h_t(x) \right).$$

A flowchart illustration



Why AdaBoost high impact?

First, it is simple yet effective

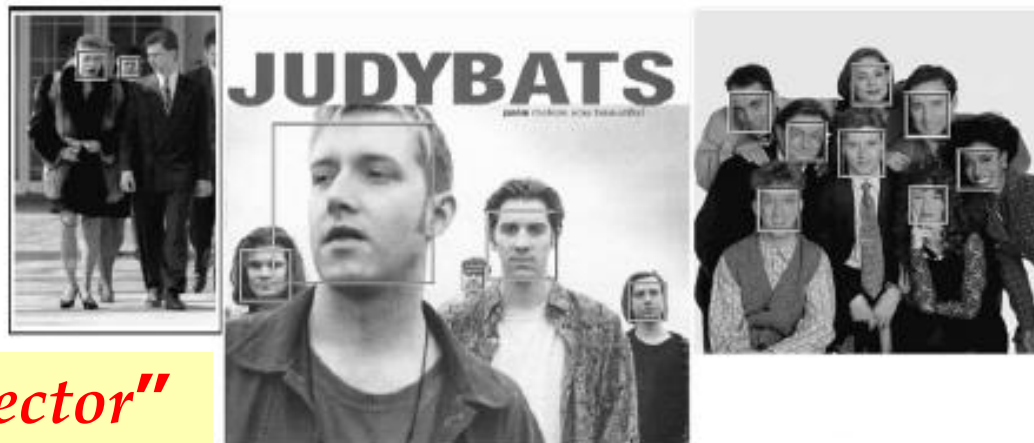
can be applied to almost all tasks where one wants to apply machine learning techniques

For example, in computer vision, the **Viola-Jones detector**
AdaBoost using harr-like features in a cascade structure



in average, only 8 features needed to be evaluated per image

The Viola-Jones detector



“the first real-time face detector”

Comparable accuracy, but **15 times faster** than state-of-the-art of face detectors (at that time)



Longuet-Higgins Prize (2011)

Viola & Jones, Rapid object detection using a Boosted cascade of simple features. CVPR, 2001.

Why AdaBoost high impact? (con't)

Second, it generates the Boosting Family of algorithms

A general boosting procedure

Input: Sample distribution \mathcal{D} ;
Base learning algorithm \mathcal{L} ;
Number of learning rounds T .

Process:

1. $\mathcal{D}_1 = \mathcal{D}$. % Initialize distribution
2. **for** $t = 1, \dots, T$:
3. $h_t = \mathcal{L}(\mathcal{D}_t)$; % Train a weak learner from distribution \mathcal{D}_t
4. $\epsilon_t = P_{\mathbf{x} \sim \mathcal{D}_t}(h_t(\mathbf{x}) \neq f(\mathbf{x}))$; % Evaluate the error of h_t
5. $\mathcal{D}_{t+1} = \text{Adjust_Distribution}(\mathcal{D}_t, \epsilon_t)$
6. **end**

Output: $H(\mathbf{x}) = \text{Combine_Outputs}(\{h_1(\mathbf{x}), \dots, h_t(\mathbf{x})\})$

A lot of Boosting algorithms:

AdaBoost.M1, AdaBoost.MR, FilterBoost, GentleBoost, GradientBoost, MadaBoost, LogitBoost, LPBoost, MultiBoost, RealBoost, RobustBoost, ...

Why AdaBoost high impact? (con't)

Third, there are sound theoretical results

Freund & Schapire [JCSS97] proved that the training error of AdaBoost is bounded by:

$$\epsilon = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \mathbb{I}[H(\mathbf{x}) \neq f(\mathbf{x})] \leq 2^T \prod_{t=1}^T \sqrt{\epsilon_t(1 - \epsilon_t)} \leq e^{-2 \sum_{t=1}^T \gamma_t^2}$$

where $\gamma_t = 0.5 - \epsilon_t$

Thus, if each base classifier is slightly better than random such that $\gamma_t \geq \gamma$ for some $\gamma > 0$, then **the training error drops exponentially fast** in T because the above bound is at most $e^{-2T\gamma^2}$

Generalization bound

Freund & Schapire [JCSS97] proved that the generalization error of AdaBoost is bounded by:

$$\epsilon_{\mathcal{D}} \leq \epsilon_D + \tilde{O} \left(\sqrt{\frac{dT}{m}} \right)$$

with probability at least $1 - \delta$, where d is the **VC-dimension** of base learners, m is the number of training instances, T is the number of learning rounds and $\tilde{O}(\cdot)$ is used instead of $O(\cdot)$ to hide logarithmic terms and constant factors.

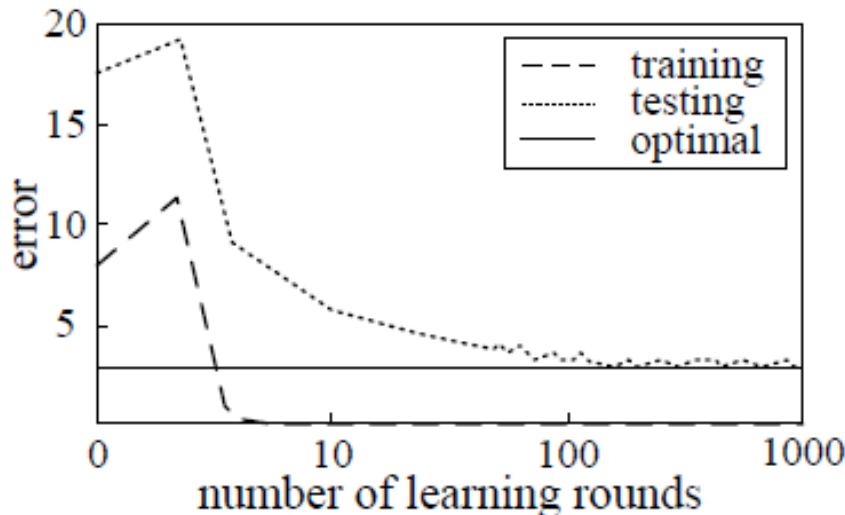
It implies that AdaBoost will **overfit** if T is large

Overfit (过拟合): The trained model fits the training data too much such that it can exaggerate minor fluctuations in the training data, leading to poor generalization performance

The Mystery

However, AdaBoost often does not overfit in real practice

A typical performance plot of AdaBoost on real data



Seems contradict with the **Occam's Razor**

Knowing the reason may inspire new methodology for algorithm design

Understanding why AdaBoost seems resistant to overfitting is the most fascinating fundamental theoretical issue

Major theoretical efforts

□ Margin Theory

Started from [Schapire, Freund, Bartlett & Lee, Boosting the margin: A new explanation for the effectiveness of voting methods. Annals of Statistics, 26(5):1651–1686, 1998]

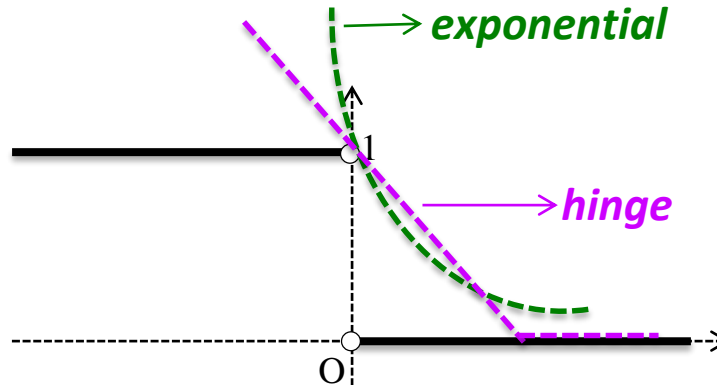
□ Statistical View

Started from [Friedman, Hastie & Tibshirani. Additive logistic regression: A statistical view of boosting (with discussions). Annals of Statistics, 28(2):337–407, 2000]

Intuition of the statistical view

In binary classification, we want to optimize the 0/1-loss

Because it is non-smooth, non-convex, ..., in statistical learning usually we instead optimize a **surrogate loss**



The key step of the AdaBoost algorithm seems closely related to the exponential loss:

$$D_{t+1}(i) = \frac{D_t(i)}{Z_t} e^{-\alpha_t y_i h_t(\mathbf{x}_i)} \quad \alpha_t = \frac{1}{2} \ln \left(\frac{1-\epsilon_t}{\epsilon_t} \right)$$

Statistical view of AdaBoost

Friedman, Hastie & Tibshirani [Ann. Stat. 2000] showed that if we consider the **additive model** $H(\mathbf{x}) = \sum_{t=1}^T \alpha_t h_t(\mathbf{x}_t)$, take a logistic function and estimate probability via

$$P(f(\mathbf{x}) = 1 \mid \mathbf{x}) = \frac{e^{H(\mathbf{x})}}{e^{H(\mathbf{x})} + e^{-H(\mathbf{x})}}$$

then AdaBoost algorithm is a Newton-like procedure optimizing the exponential loss function and the log loss function (negative log-likelihood)

$$\ell_{\log}(h \mid \mathcal{D}) = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \left[\ln \left(1 + e^{-2f(\mathbf{x})h(\mathbf{x})} \right) \right]$$

That is, **AdaBoost can be viewed as a stage-wise estimation procedure for fitting an additive logistic regression model**

Implications of the statistical view

As alternatives, one can fit the additive logistic regression model by optimizing the log loss function via other procedures, leading to many variants

- e.g., LogitBoost [Friedman, Hastie & Tibshirani, Ann. Stat. 2000]
- LPBoost [Demiriz, Bennett & Shawe-Taylor, MLJ 2002]
- L2Boost [Bühlmann & Yu, JASA 2003]
- RegBoost [Lugosi & Vayatis, Ann. Stat. 2004], etc.

The statistical view also encouraged the study of some specific statistical properties of AdaBoost

- e.g., for **consistency**: Boosting with early stopping is consistent [Zhang & Yu, Ann. Stat. 2004], Exponential and logistic loss is consistent [Zhang, Ann. Stat. 2004, Bartlett, Jordana & McAuliffea, JASA 2006], etc.

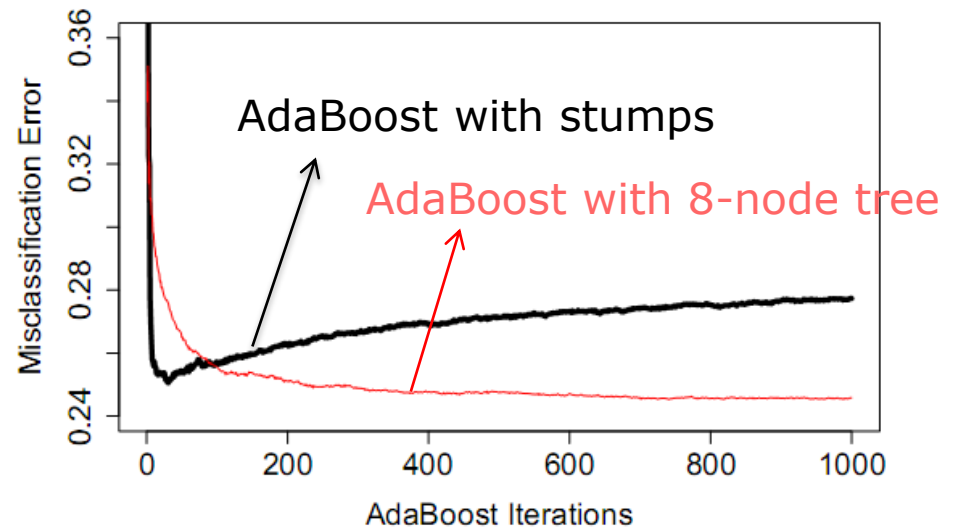
Concerns about the statistical view

However, many aspects of the statistical view have been questioned by empirical results

e.g., in a famous article [Mease & Wyner. Evidence contrary to the statistical view of boosting (with discussions). JMLR, 9:131–201, 2008] it was disclosed that:

Larger-size trees will lead to overfitting because of higher-level interaction [Friedman, Hastie & Tibshirani, Ann. Stat. 2000]

But in practice ...



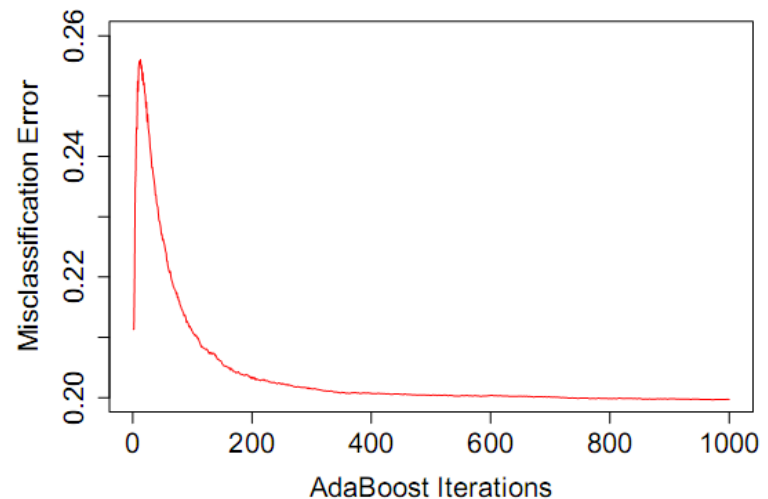
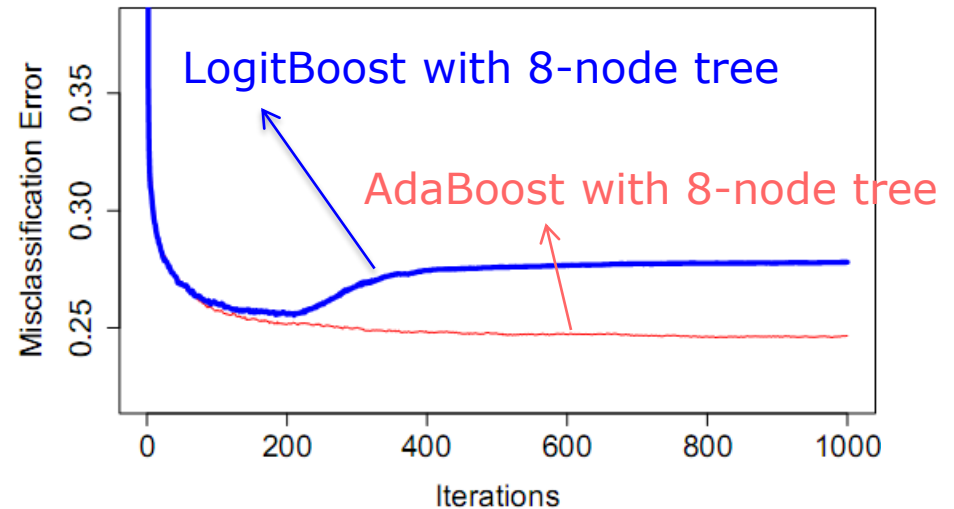
Concerns about the statistical view (con't)

LogitBoost is better than AdaBoost for noisy data
[Hastie, Tibshirani & Friedman, "The Elements of Statistical Learning", Springer 2001]

But in practice ...

Early stopping can be used to prevent overfitting [Zhang & Yu, Ann. Stat. 2004]

But in practice ...



Major theoretical efforts

□ **Margin Theory**

Started from [Schapire, Freund, Bartlett & Lee, Boosting the margin: A new explanation for the effectiveness of voting methods. Annals of Statistics, 26(5):1651–1686, 1998]

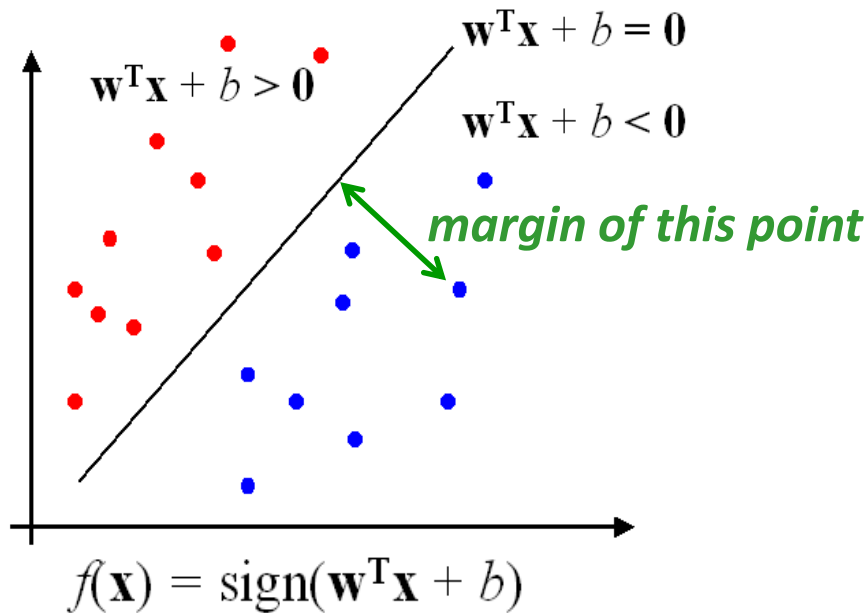
□ **Statistical View**

Started from [Friedman, Hastie & Tibshirani. Additive logistic regression: A statistical view of boosting (with discussions). Annals of Statistics]

**The biggest issue:
The statistical view did not explain why
AdaBoost is resistant to overfitting**

The “margin” (间隔)

Binary classification can be viewed as the task of separating classes in a feature space



The bigger the margin,
the higher the predictive confidence

For binary classification, the ground-truth $f(\mathbf{x}) \in \{-1, +1\}$

The margin of a single classifier h : $f(\mathbf{x})h(\mathbf{x})$

For $H(\mathbf{x}) = \sum_{t=1}^T \alpha_t h_t(\mathbf{x}_t)$
the margin is

$$f(\mathbf{x})H(\mathbf{x}) = \sum_{t=1}^T \alpha_t f(\mathbf{x})h_t(\mathbf{x})$$

and the normalized margin:

$$\frac{\sum_{t=1}^T \alpha_t f(\mathbf{x})h_t(\mathbf{x})}{\sum_{t=1}^T \alpha_t}$$

Margin explanation of AdaBoost

Based on the concept of margin, Schapire et al. [1998] proved that, given any threshold $\theta > 0$ of margin over the training data D , with probability at least $1 - \delta$, the generalization error of the ensemble $\epsilon_{\mathcal{D}} = P_{\mathbf{x} \sim \mathcal{D}}(f(\mathbf{x}) \neq H(\mathbf{x}))$ is bounded by

$$\begin{aligned} \epsilon_{\mathcal{D}} &\leq P_{\mathbf{x} \sim \mathcal{D}}(f(\mathbf{x})H(\mathbf{x}) \leq \theta) + \tilde{O} \left(\sqrt{\frac{d}{m\theta^2} + \ln \frac{1}{\delta}} \right) \\ &\leq 2^T \prod_{t=1}^T \sqrt{\epsilon_t^{1-\theta} (1 - \epsilon_t)^{1+\theta}} + \tilde{O} \left(\sqrt{\frac{d}{m\theta^2} + \ln \frac{1}{\delta}} \right) \end{aligned}$$

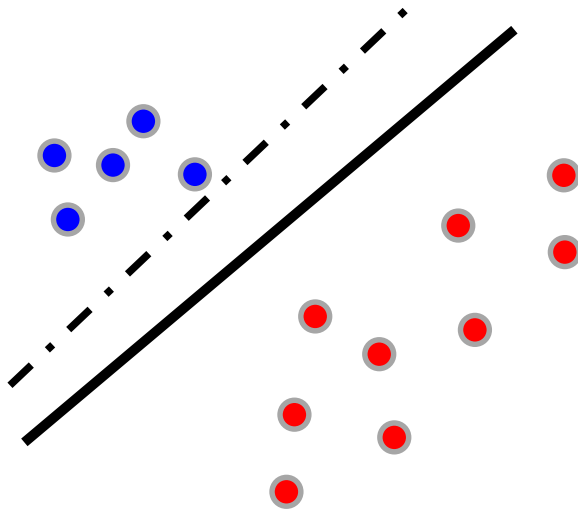
This bound implies that, when other variables are fixed, the larger the margin over the training data, the smaller the generalization error

Margin explanation of AdaBoost (con't)

Why AdaBoost tends to be resistant to overfitting?

the margin theory answers:

Because it is able to increase the ensemble margin even after the training error reaches zero



This explanation is quite intuitive

It receives good support in empirical study

The minimum margin bound

Schapire et al.'s bound depends heavily on the smallest margin, because $P_{\mathbf{x} \sim D}(f(\mathbf{x})H(\mathbf{x}) \leq \theta)$ will be small if the smallest margin is large

Thus, by considering the minimum margin:

$$\varrho = \min_{\mathbf{x} \in D} f(\mathbf{x})H(\mathbf{x})$$

Breiman [Neural Comp. 1999] proved a generalization bound, which is tighter than Schapire et al.'s bound

The two generalization bounds

Theorem 1. (Schapire et al., 1998) For any $\delta > 0$ and $\theta > 0$, with probability at least $1 - \delta$ over the random choice of sample S with size m , every voting classifier $f \in \mathcal{C}(\mathcal{H})$ satisfies the following bound:

$$\Pr_D[yf(x) < 0] \leq \Pr_S[yf(x) \leq \theta] + O\left(\frac{1}{\sqrt{m}} \left(\frac{\ln m \ln |\mathcal{H}|}{\theta^2} + \ln \frac{1}{\delta}\right)^{1/2}\right).$$

$O(\sqrt{\log m / m})$

Theorem 2. (Breiman, 1999) If

$$\theta = \hat{y}_1 f(\hat{x}_1) > 4\sqrt{\frac{2}{|\mathcal{H}|}} \text{ and } R = \frac{32 \ln 2 |\mathcal{H}|}{m\theta^2} \leq 2m,$$

then, for any $\delta > 0$, with probability at least $1 - \delta$ over the random choice of sample S with size m , every voting classifier $f \in \mathcal{C}(\mathcal{H})$ satisfies the following bound:

$$\Pr_D[yf(x) < 0] \leq R \left(\ln(2m) + \ln \frac{1}{R} + 1 \right) + \frac{1}{m} \ln \frac{|\mathcal{H}|}{\delta}.$$

$O(\log m / m)$

The doubt about margin theory

Breiman [Neural Comp. 1999] designed a variant of AdaBoost, the arc-gv algorithm, which directly maximizes the minimum margin

the margin theory would appear to predict that arc-gv should perform better than AdaBoost

However, experiments show that, comparing with AdaBoost:

- arc-gv does produce **uniformly larger minimum margin**
- **the test error increases drastically** in almost every case

Thus, Breiman convincingly concluded that **the margin theory was in serious doubt**. This almost sentenced the margin theory to death

7 years later ...

Reyzin & Schapire [ICML'06 best paper] found that, amazingly, Breiman had not controlled model complexity well in exps

Breiman controlled the model complexity by using decision trees with a fixed number of leaves

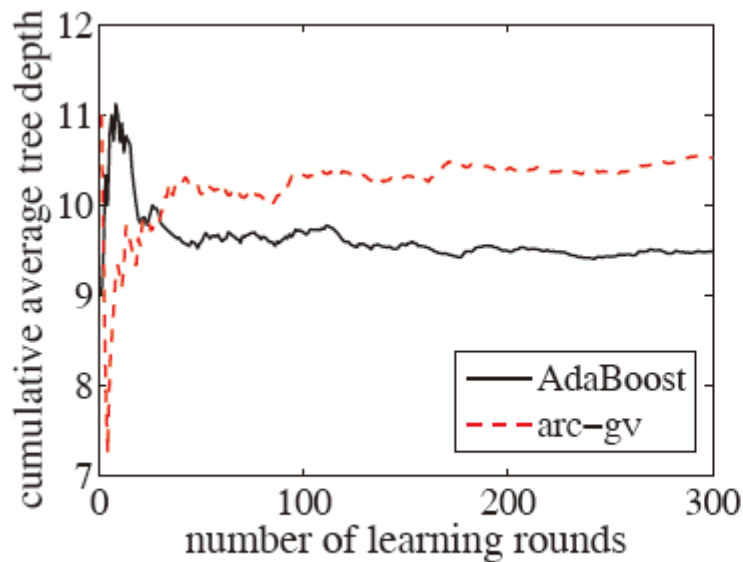
Reyzin & Schapire found that, the trees of arc-gv are generally "deeper" than the trees of AdaBoost

Reyzin & Schapire repeated Breiman's exps using decision stumps with two leaves: arc-gv is with larger minimum margin, but worse margin distribution

R&S claimed that the minimum margin is not crucial, and the *average* or *median margin* is crucial

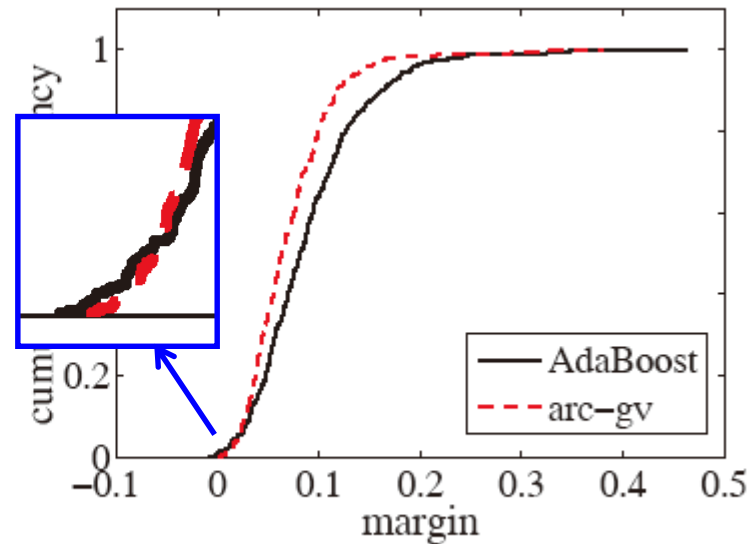
Experimental results

*Tree depth
using fixed number of leaves*



(a)

*Margin distribution
using decision stumps*



(b)

FIGURE 2.8: (a) Tree depth and (b) margin distribution of AdaBoost against arc-gv on the UCI *clean1* data set.

Margin theory survive?

Not necessarily ...

Breiman's minimum margin bound is tighter

To claim margin distribution is more crucial, we need a margin distribution bound which is even tighter

Equilibrium margin (Emargin) bound

Theorem 3. (Wang et al., 2011) If $8 < |\mathcal{H}| < \infty$, then for any $\delta > 0$, with probability at least $1 - \delta$ over the random choice of the training set S of size $m > 1$, every voting classifier $f \in \mathcal{C}(\mathcal{H})$ such that

$$q_0 = \Pr_S \left[yf(x) \leq \sqrt{8/|\mathcal{H}|} \right] < 1 \quad (3)$$

satisfies the following bound:

$$\Pr_D[yf(x) < 0] \leq \frac{\ln |\mathcal{H}|}{m} + \inf_{q \in \{q_0, q_0 + \frac{1}{m}, \dots, 1\}} KL^{-1}(q; u[\hat{\theta}(q)]),$$

where

$$u[\hat{\theta}(q)] = \frac{1}{m} \left(\frac{8 \ln |\mathcal{H}|}{\hat{\theta}^2(q)} \ln \frac{2m^2}{\ln |\mathcal{H}|} + \ln |\mathcal{H}| + \ln \frac{m}{\delta} \right)$$

and $\hat{\theta}(q) = \sup \{ \theta \in (\sqrt{8/|\mathcal{H}|}, 1] : \Pr_S[yf(x) \leq \theta] \leq q \}$. Also, the Emargin is given by $\theta^* \in \arg \inf_{q \in \{q_0, q_0 + \frac{1}{m}, \dots, 1\}} KL^{-1}(q; u[\hat{\theta}(q)])$.

Proved to be tighter than Breiman's bound

$O(\log m / m)$

- Considered factors different from Schapire et al. and Breiman's bounds
- No intuition to optimize

The k th margin bound

Given a sample S of size m , we define the k th margin as the k th smallest margin over sample S , i.e., the k th smallest value in $\{y_i f(x_i), i \in [m]\}$

Theorem 4. For any $\delta > 0$ and $k \in [m]$, if $\theta = \hat{y}_k f(\hat{x}_k) > \sqrt{8/|\mathcal{H}|}$, then with probability at least $1 - \delta$ over the random choice of sample with size m , every voting classifier $f \in \mathcal{C}(\mathcal{H})$ satisfies the following bound:

$$\Pr_D[yf(x) < 0] \leq \frac{\ln |\mathcal{H}|}{m} + KL^{-1}\left(\frac{k-1}{m}; \frac{q}{m}\right), \quad (4)$$

where

$$q = \frac{8 \ln(2|\mathcal{H}|)}{\theta^2} \ln \frac{2m^2}{\ln |\mathcal{H}|} + \ln |\mathcal{H}| + \ln \frac{m}{\delta}.$$

The minimum margin bound and Emargin bound are special cases of the k th margin bound, both are single-margin bound (not margin distribution bound)

Finally, our margin distribution bound

Theorem 8. *For any $\delta > 0$, with probability at least $1 - \delta$ over the random choice of sample S with size $m \geq 5$, every voting classifier $f \in \mathcal{C}(\mathcal{H})$ satisfies the following bound:*

$$\Pr_D[yf(x) < 0] \leq \frac{2}{m} + \inf_{\theta \in (0,1]} \left[\Pr_S[yf(x) < \theta] + \frac{7\mu + 3\sqrt{3\mu}}{3m} + \sqrt{\frac{3\mu}{m} \Pr_S[yf(x) < \theta]} \right]$$

where

$$\mu = \frac{8}{\theta^2} \ln m \ln(2|\mathcal{H}|) + \ln \frac{2|\mathcal{H}|}{\delta}.$$

$O(\log m / m)$

- ✓ Uniformly tighter than Breiman's as well as Schapire et al.' bounds
- ✓ Considers the same factors as Schapire et al. and Breiman

thus, defends the margin theory against Breiman's doubt

New insight?

Theorem 9. For any $\delta > 0$, with probability at least $1 - \delta$ over the random choice of sample S with size $m \geq 5$, every voting classifier $f \in \mathcal{C}(\mathcal{H})$ satisfies the following bound:

$$\Pr_D[yf(x) < 0] \leq \frac{1}{m^{50}} + \inf_{\theta \in (0,1]} \left[\Pr_S[yf(x) < \theta] + m^{-2/(1-E_S^2[yf(x)]+\theta/9)} \right]$$

related to average margin

$$+ \frac{3\sqrt{\mu}}{m^{3/2}} + \frac{7\mu}{3m} + \sqrt{\frac{3\mu}{m} \hat{\mathcal{I}}(\theta)}$$

$O(\log m / m)$

where

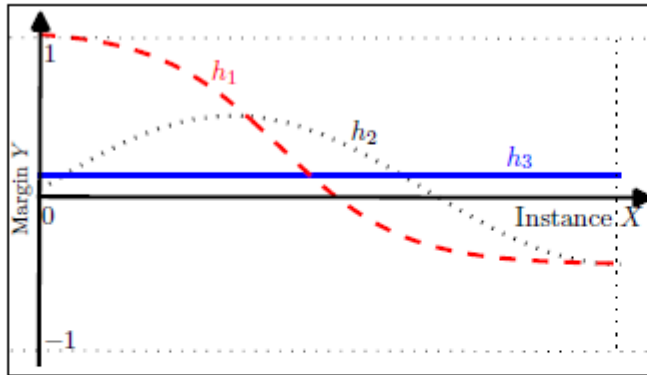
$$\mu = 144 \ln m \ln(2|\mathcal{H}|)/\theta^2 + \ln(2|\mathcal{H}|/\delta),$$

$$\hat{\mathcal{I}}(\theta) = \Pr_S[yf(x) < \theta] \Pr_S[yf(x) \geq 2\theta/3].$$

related to margin variance

We should pay attention to not only the average margin, but also the margin variance !

In practice



Margin variance really important

Figure from [Gao & Zhou, AIJ 2013]

Figure 1: Each curve represents a voting classifier. The X -axis and Y -axis denote example and margin, respectively, and uniform distribution is assumed on the example space. The voting classifiers h_1 , h_2 and h_3 have the same average margin but with different generalization error rates: $1/2$, $1/3$ and 0 .

[Shivaswamy & Jebara, NIPS 2011] tried to design new boosting algorithms by maximizing average margin and minimizing margin variance simultaneously, and the results are encouraging

Long march of margin theory for AdaBoost

- 1989, [Kearns & Valiant], [open problem](#)
- 1990, [Schapire], [proof by construction](#), the first Boosting algorithm
- 1993, [Freund], [another impractical boosting algorithm by voting](#)
- 1995/97, [Freund & Schapire], [AdaBoost](#)

- 1998, [Schapire, Freund, Bartlett & Lee], [Margin theory](#)
- 1999, [Breiman], [serious doubt by minimum margin bound](#)
- 2006, [Reyzin & Schapire], [finding the model complexity issue in exps, emphasizing the importance of margin distribution](#)
- 2008, [Wang, Sugiyama, Yang, Zhou & Feng], [Emargin bound, believed to be a margin distribution bound](#)
- 2013, [Gao & Zhou], [a real margin distribution bound, shedding new insight ; margin theory defended](#)

Outline

- **A long march of margin theory for Boosting**
(a theory breakthrough)

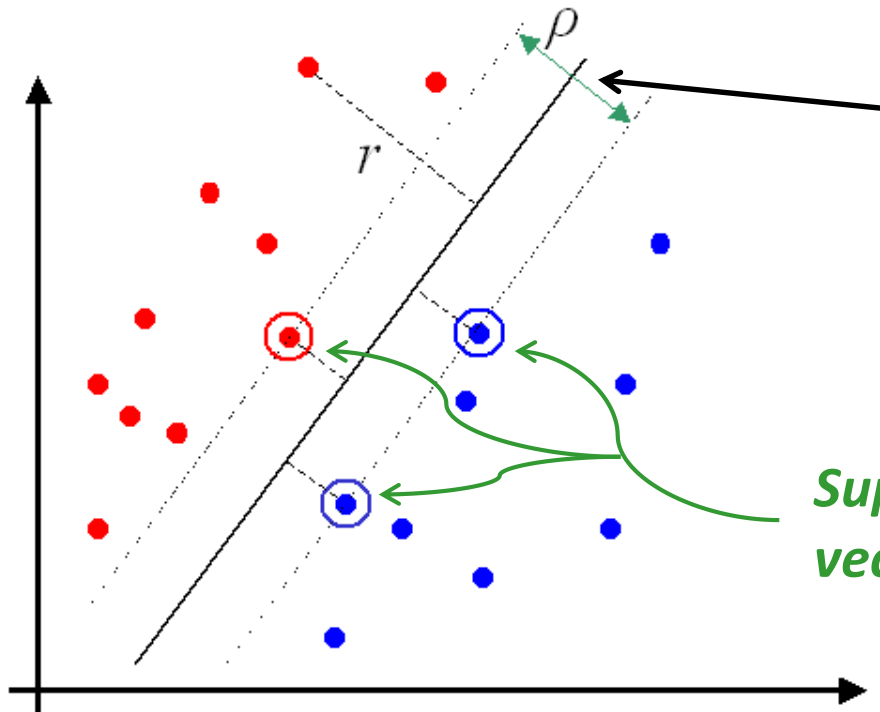
- **LDM: Large margin Distribution Machine**
(an algorithmic paradigm)

Question

Do we have direct and principled way to:

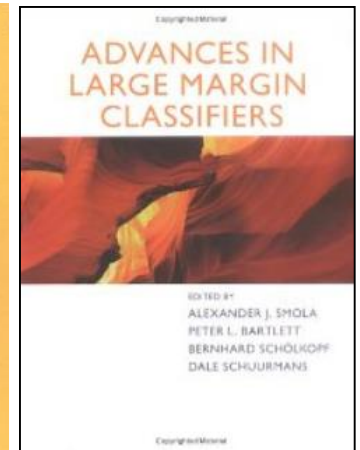
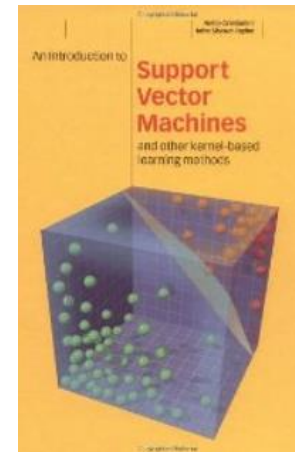
maximize the margin mean and minimize the margin variance simultaneously?

Large margin classifiers

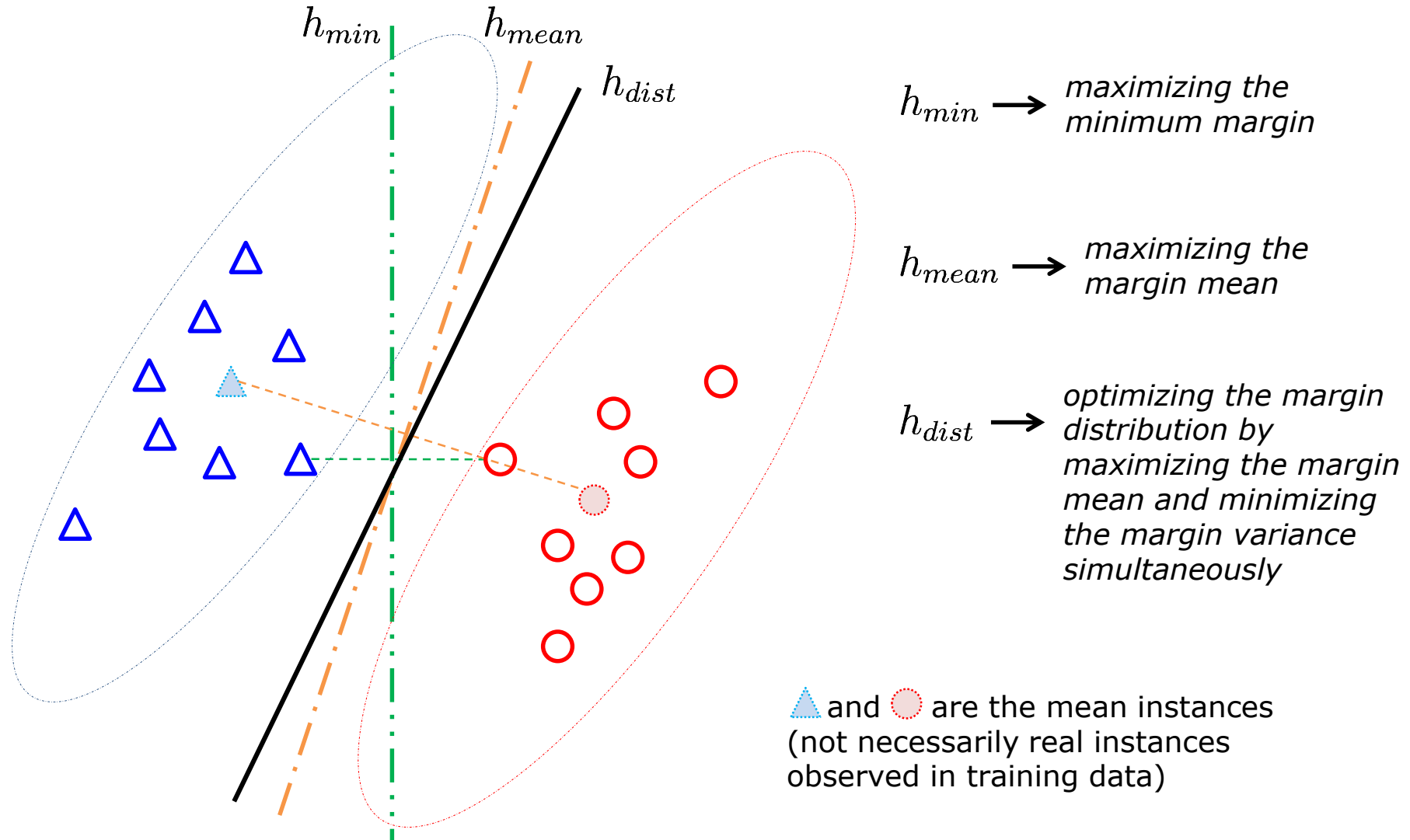


Large margin classifiers are actually trying to maximize the **minimum margin**

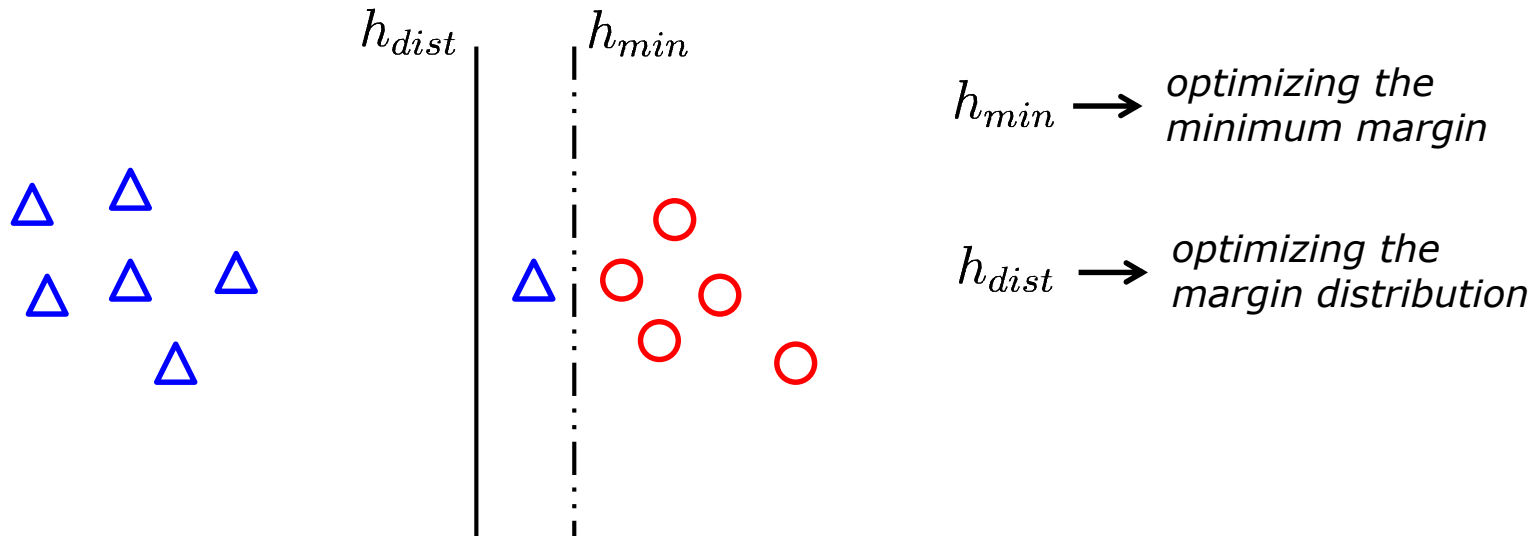
Support vectors



Big difference between “large margin” and “large margin distribution”



Another advantage of “large margin distribution learning”



Less sensitive to outliers or noisy data points

Formally

Margin: $\gamma = yf(\mathbf{x}) = y\mathbf{w}^\top \phi(\mathbf{x})$

Margin mean: $\bar{\gamma} = \frac{1}{m} \sum_{i=1}^m y_i \mathbf{w}^\top \phi(\mathbf{x}_i) = \frac{1}{m} (\mathbf{X}\mathbf{y})^\top \mathbf{w}$

Margin variance:

$$\begin{aligned} \hat{\gamma} &= \frac{1}{m^2} \sum_{i=1}^m \sum_{j=1}^m (y_i \mathbf{w}^\top \phi(\mathbf{x}_i) - y_j \mathbf{w}^\top \phi(\mathbf{x}_j))^2 \\ &= \frac{2}{m^2} (m\mathbf{w}^\top \mathbf{X}\mathbf{X}^\top \mathbf{w} - \mathbf{w}^\top \mathbf{X}\mathbf{y}\mathbf{y}^\top \mathbf{X}^\top \mathbf{w}). \end{aligned}$$

The LDM formulation

minimizing the margin variance

$$\min_{\mathbf{w}, \xi} \frac{1}{2} \mathbf{w}^\top \mathbf{w} + \frac{2\lambda_1}{m^2} (m \mathbf{w}^\top \mathbf{X} \mathbf{X}^\top \mathbf{w} - \mathbf{w}^\top \mathbf{X} \mathbf{y} \mathbf{y}^\top \mathbf{X}^\top \mathbf{w})$$

maximizing the margin mean

$$- \frac{\lambda_2}{m} (\mathbf{X} \mathbf{y})^\top \mathbf{w} + C \sum_{i=1}^m \xi_i$$

$$\text{s.t. } y_i \mathbf{w}^\top \phi(\mathbf{x}_i) \geq 1 - \xi_i, \\ \xi_i \geq 0, \quad i = 1, \dots, m.$$

*Not really needed;
 kept just to reuse
 SVM implementation*

With representer theorem, the dual is:

$$\min_{\beta} f(\beta) = \frac{1}{2} \beta^\top \mathbf{H} \beta + \left(\frac{\lambda_2}{m} \mathbf{H} \mathbf{e} - \mathbf{e} \right)^\top \beta, \\ \text{s.t. } 0 \leq \beta_i \leq C, \quad i = 1, \dots, m.$$

The optimization

$$\begin{aligned} \min_{\beta} \quad & f(\beta) = \frac{1}{2} \beta^\top \mathbf{H} \beta + \left(\frac{\lambda_2}{m} \mathbf{H} \mathbf{e} - \mathbf{e} \right)^\top \beta, \\ \text{s.t.} \quad & 0 \leq \beta_i \leq C, \quad i = 1, \dots, m. \end{aligned}$$

convex quadratic programming
decoupled box constraint

It can be solved by **coordinate descent** efficiently since a **closed-form solution** can be achieved in each iteration

$$\beta_i^{\text{new}} = \min \left(\max \left(\beta_i - \frac{[\nabla f(\beta)]_i}{h_{ii}}, 0 \right), C \right)$$

For large scale problems, we can solve the prime LDM directly by **average stochastic gradient descent (ASGD)**

Analysis

Based on leave-one-out cross-validation estimate, we can derive a bound on the expectation of error for LDM:

THEOREM 3. Let α denote the optimal solution of (18), and $E[R(\alpha)]$ be the expectation of the probability of error, then we have

$$E[R(\alpha)] \leq \frac{E[h \sum_{i \in I_1} \alpha_i + |I_2|]}{m}, \quad (19)$$

where $I_1 \equiv \{i \mid 0 < \alpha_i < C\}$, $I_2 \equiv \{i \mid \alpha_i = C\}$, $h = \max\{\text{diag}\{\mathbf{H}\}\}$ and $\mathbf{H} = \mathbf{Y}\mathbf{X}^\top \mathbf{Q}^{-1} \mathbf{X}\mathbf{Y}$.

- A similar bound also holds for SVM [Vapnik, 1995], and the only difference is: $\mathbf{H} = \mathbf{Y}\mathbf{X}^\top \mathbf{X}\mathbf{Y}$
- \mathbf{Q} encodes the information of the margin distribution with the result that the value of h for LDM is **much smaller** than SVM

Related work

- MDO [Garg & Roth, ICML 2003]
Optimize the weighted margin combination; the setting of weights is heuristic, and the objective function is non-convex
- MAMC [Pelckmans et al, NIPS 2007]
Optimize the average margin; it suffers from unequal class size and is a special case of LDM
- KM-OMD [Aioli et al, ICANN 2008]
Optimize the weighted margin combination; hard-margin only

These methods consider the margin mean in some sense, but ignore the influence of margin variance

Experimental settings

Compared methods:

LDM, MDO, MAMC, KM-OMD, SVM

Evaluation:

30 times hold-out tests, 1/2 training, 1/2 testing

Parameters:

selected by 5-fold cross validation on training data

Experimental data

20 regular scale data sets and 12 large scale data sets

Scale	Dataset	#Instance	#Feature	Dataset	#Instance	#Feature
regular	<i>promoters</i>	106	57	<i>haberman</i>	306	14
	<i>planning</i>	182	12	<i>vehicle</i>	435	16
	<i>colic</i>	188	12	<i>glass1</i>	176	166
	<i>house</i>	252	10	<i>australian</i>	890	42
	<i>heart</i>	270	9	<i>fourclass</i>	862	2
	<i>breast</i>	277	9	<i>german</i>	1,000	59
large	<i>farm-ads</i>	4,143	54,877	<i>ijcnn1</i>	141,691	22
	<i>news20</i>	19,996	1,355,191	<i>skin</i>	245,057	3
	<i>adult-a</i>	32,561	123	<i>covtype</i>	581,012	54
	<i>w8a</i>	49,749	300	<i>rcv1</i>	697,641	47,236
	<i>cod-rna</i>	59,535	8	<i>url</i>	2,396,130	3,231,961
	<i>real-sim</i>	72,309	20,958	<i>kdd2010</i>	8,407,752	20,216,830

Data size ranged from 106 to more than 8,000,000
Dimensionality ranged from 2 to more than 20,000,000

Results (regular scale data, linear kernel)

Dataset	SVM	MDO	MAMC	KM-OMD	LDM
<i>promoters</i>	0.723±0.071	0.713±0.067	0.520±0.096○	0.736±0.061	0.721±0.069
<i>planning-relax</i>	0.683±0.031	0.605±0.185○	0.706±0.034●	0.479±0.050○	0.706±0.034●
<i>colic</i>	0.814±0.035	0.781±0.154	0.661±0.062○	0.813±0.028	0.832±0.026●
<i>parkinsons</i>	0.846±0.038	0.732±0.270○	0.764±0.035○	0.814±0.024○	0.865±0.030●
<i>colic.ORIG</i>	0.618±0.027	0.624±0.040	0.623±0.027	0.635±0.045●	0.619±0.042
<i>sonar</i>	0.725±0.039	0.734±0.035	0.533±0.045○	0.766±0.033●	0.736±0.036
<i>vote</i>	0.934±0.022	0.587±0.435○	0.884±0.022○	0.957±0.013●	0.970±0.014●
<i>house</i>	0.942±0.015	0.943±0.015	0.883±0.029○	0.957±0.020●	0.968±0.011●
<i>heart</i>	0.799±0.029	0.826±0.026●	0.537±0.057○	0.836±0.026●	0.791±0.030
<i>breast-cancer</i>	0.717±0.033	0.710±0.031	0.706±0.027	0.696±0.031○	0.725±0.027●
<i>haberman</i>	0.734±0.030	0.728±0.029	0.738±0.020	0.667±0.040○	0.738±0.020
<i>vehicle</i>	0.959±0.012	0.956±0.012	0.566±0.160○	0.960±0.010	0.959±0.013
<i>clean1</i>	0.803±0.035	0.798±0.031	0.561±0.025○	0.821±0.027●	0.814±0.019●
<i>wdbc</i>	0.963±0.012	0.966±0.010	0.623±0.020○	0.968±0.009●	0.968±0.011●
<i>isolet</i>	0.995±0.003	0.501±0.503○	0.621±0.207○	0.995±0.003	0.997±0.002●
<i>credit-a</i>	0.861±0.014	0.862±0.013	0.596±0.063○	0.863±0.013	0.864±0.013●
<i>austra</i>	0.857±0.013	0.842±0.055	0.567±0.044○	0.858±0.013	0.859±0.015
<i>australian</i>	0.844±0.019	0.842±0.020	0.576±0.049○	0.858±0.016●	0.866±0.014●
<i>fourclass</i>	0.724±0.014	0.377±0.238○	0.641±0.020○	0.736±0.014●	0.723±0.014
<i>german</i>	0.711±0.030	0.737±0.014●	0.697±0.017○	0.729±0.017●	0.738±0.016●
Ave. accuracy	0.813	0.743	0.650	0.807	0.823
LDM: w/t/l	12/8/0	9/10/1	17/3/0	10/5/5	

bold: best

●/○:
significantly
better/worse
than SVM

**LDM is best
on 13/20 data
sets, and
best average
accuracy**

w/t/l counts:
after *t*-tests
(95% SI)

**LDM never loses to SVM
(other approaches often
lose to SVM)**

LDM outperforms others significantly

Results (regular scale data, RBF kernel)

Dataset	SVM	MDO	MAMC	KM-OMD	LDM
<i>promoters</i>	0.684±0.100	N/A	0.638±0.121○	0.701±0.085	0.715±0.074●
<i>planning-relax</i>	0.708±0.035	N/A	0.706±0.034	0.683±0.031○	0.707±0.034
<i>colic</i>	0.822±0.033	N/A	0.623±0.037○	0.825±0.024	0.841±0.018●
<i>parkinsons</i>	0.929±0.029	N/A	0.852±0.036○	0.906±0.033○	0.927±0.029
<i>colic.ORIG</i>	0.638±0.043	N/A	0.623±0.027	0.621±0.039	0.641±0.044
<i>sonar</i>	0.842±0.034	N/A	0.753±0.052○	0.821±0.051○	0.846±0.032
<i>vote</i>	0.946±0.016	N/A	0.913±0.019○	0.930±0.029○	0.968±0.013●
<i>house</i>	0.953±0.020	N/A	0.561±0.139○	0.938±0.022○	0.964±0.013●
<i>heart</i>	0.808±0.025	N/A	0.540±0.043○	0.805±0.048	0.822±0.029●
<i>breast-cancer</i>	0.729±0.030	N/A	0.706±0.027○	0.691±0.024○	0.753±0.027●
<i>haberman</i>	0.727±0.024	N/A	0.742±0.021●	0.676±0.042○	0.731±0.027
<i>vehicle</i>	0.992±0.007	N/A	0.924±0.025○	0.988±0.008○	0.993±0.006
<i>clean1</i>	0.890±0.020	N/A	0.561±0.025○	0.772±0.043○	0.891±0.024
<i>wdbc</i>	0.951±0.011	N/A	0.740±0.042○	0.941±0.040	0.961±0.010●
<i>isolet</i>	0.998±0.002	N/A	0.994±0.004○	0.995±0.003○	0.998±0.002
<i>credit-a</i>	0.858±0.014	N/A	0.542±0.032○	0.845±0.029○	0.861±0.013
<i>austra</i>	0.853±0.013	N/A	0.560±0.018○	0.854±0.017	0.857±0.014●
<i>australian</i>	0.815±0.014	N/A	0.554±0.015○	0.860±0.014●	0.854±0.016●
<i>fourclass</i>	0.998±0.003	N/A	0.791±0.014○	0.838±0.014○	0.998±0.003
<i>german</i>	0.731±0.019	N/A	0.697±0.017○	0.742±0.017●	0.743±0.016●
Ave. accuracy	0.844	N/A	0.701	0.822	0.854
LDM: w/t/l	10/10/0	N/A	18/1/1	15/5/0	

bold: best

●/○:
significantly
better/worse
than SVM

**LDM is best
on 15/20
data sets, and
best average
accuracy**

w/t/l counts:
after *t*-tests
(95% SI)

**LDM never loses to SVM
(other approaches often
lose to SVM)**

LDM outperforms others significantly

Results (large scale data, linear kernel)

No results return in 48 hours

Dataset	SVM	MDO	MAMC	KM-OMD	LDM
<i>farm-ads</i>	0.880±0.007	0.880±0.007	0.759±0.038○	N/A	0.890±0.008●
<i>news20</i>	0.954±0.002	0.948±0.002○	0.772±0.017○	N/A	0.960±0.001●
<i>adult-a</i>	0.845±0.002	0.788±0.053○	0.759±0.002○	N/A	0.846±0.003●
<i>w8a</i>	0.983±0.001	0.985±0.001●	0.971±0.001○	N/A	0.983±0.001
<i>cod-rna</i>	0.899±0.001	0.774±0.203	0.667±0.001○	N/A	0.899±0.001
<i>real-sim</i>	0.961±0.001	0.955±0.002○	0.744±0.004○	N/A	0.971±0.001●
<i>ijcnn1</i>	0.921±0.003	0.921±0.002	0.904±0.001○	N/A	0.921±0.002
<i>skin</i>	0.934±0.001	0.929±0.003○	0.792±0.000○	N/A	0.934±0.001
<i>covtype</i>	0.762±0.001	0.760±0.003○	0.628±0.002○	N/A	0.763±0.001
<i>rcv1</i>	0.969±0.000	0.959±0.000○	0.913±0.000○	N/A	0.977±0.000●
<i>url</i>	0.993±0.006	0.993±0.006	0.670±0.000○	N/A	0.993±0.006
<i>kdd2010</i>	0.852±0.001	N/A	0.853±0.000●	N/A	0.881±0.001●
Ave. accuracy	0.913	0.899	0.786	N/A	0.919
LDM: w/t/l	6/6/0	7/3/1	12/0/0	N/A	

LDM is best on 8/12 data sets, and best average accuracy

LDM never loses to SVM (other approaches often lose to SVM)

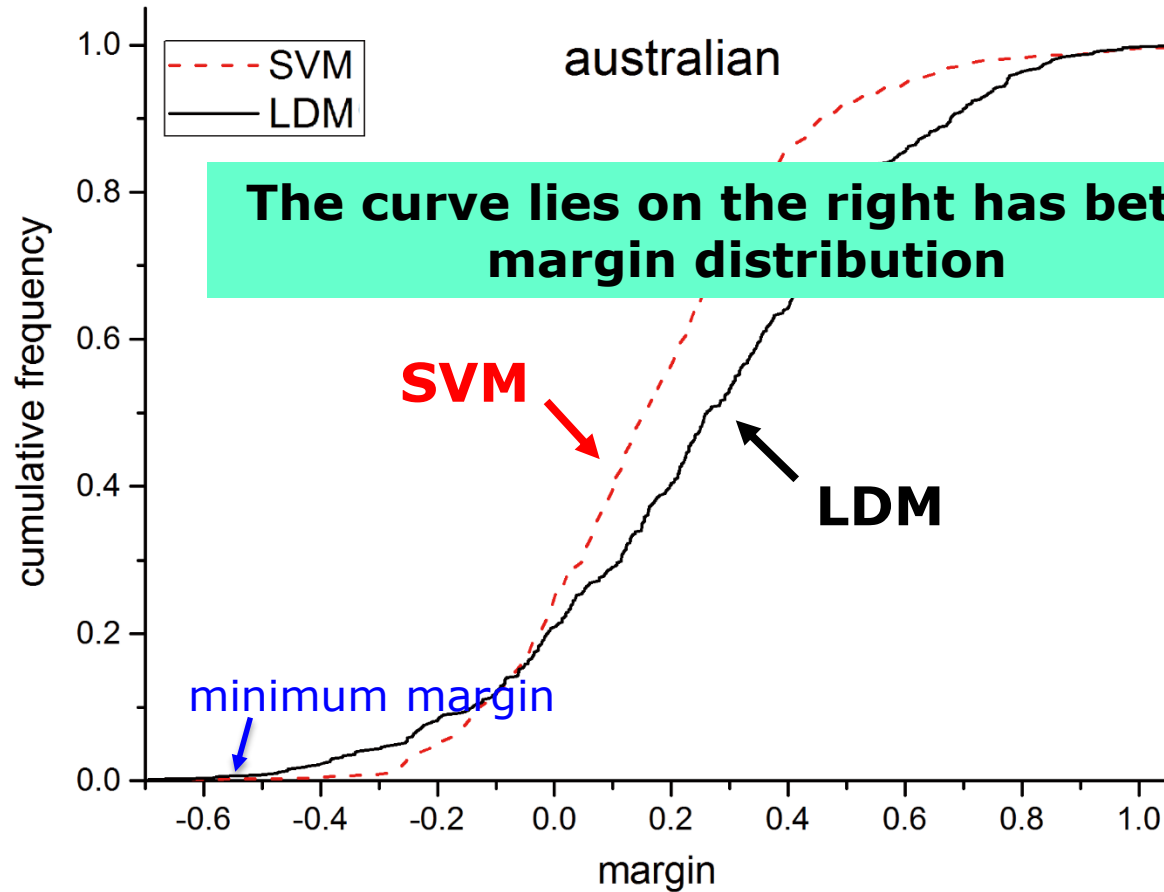
LDM outperforms others significantly

LDM always the best in empirical study (no matter data size, kernel type)

In all comparisons SVM used soft-margin version, whereas LDM used hard-margin version. LDM's superiority much larger when soft-margin used

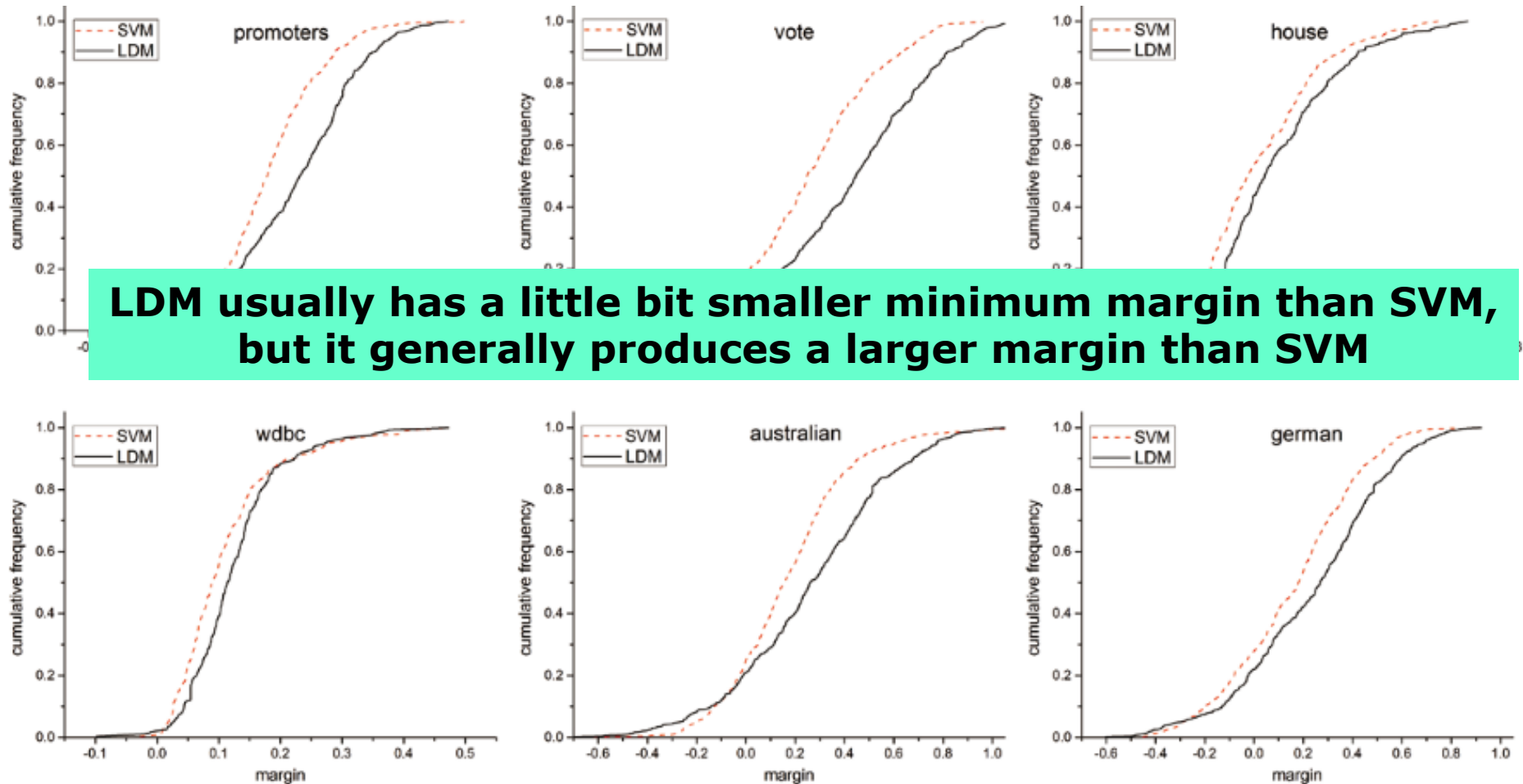
Margin distribution

Cumulative frequency (y-axis) with respect to margin (x-axis) of SVM and LDM

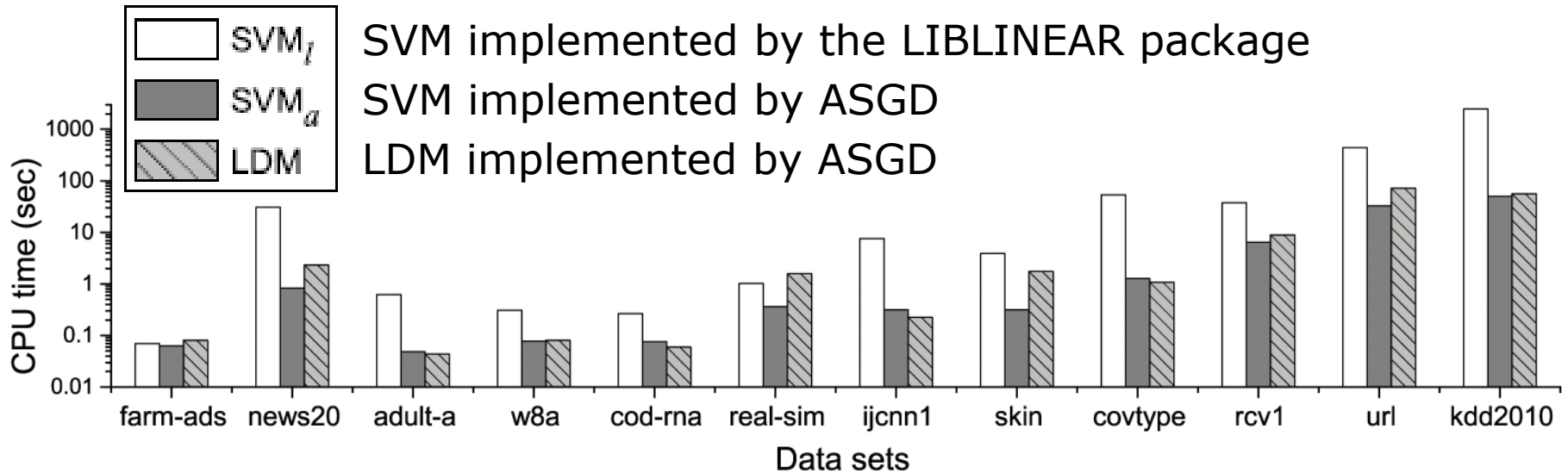


Margin distribution

Cumulative frequency (y-axis) with respect to margin (x-axis) of SVM and LDM



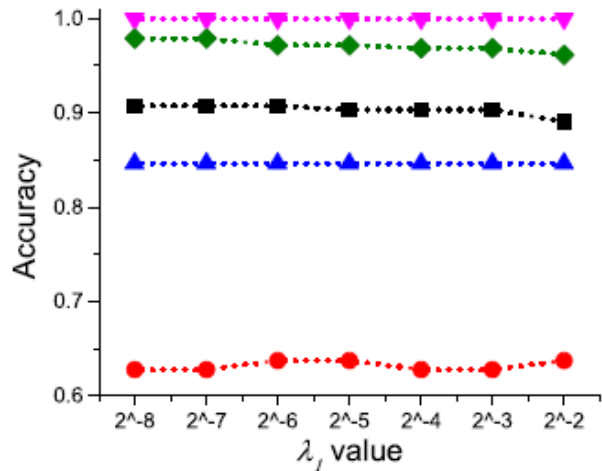
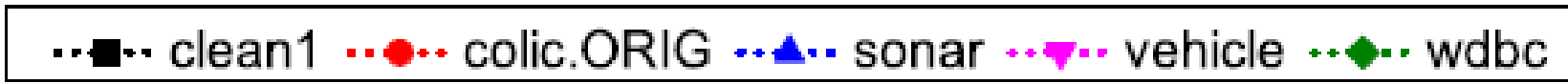
Time cost (twelve large scale data sets)



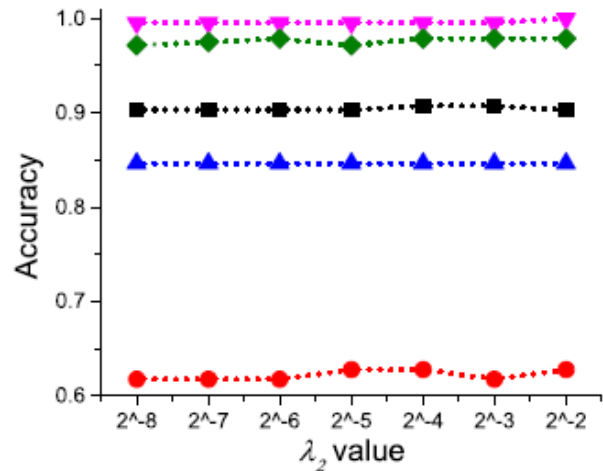
LDM quite efficient

- Comparing with common implementation of SVM (such as LIBLINEAR), LDM is faster
- Even when comparing with ASGD implementation of SVM, LDM is highly competitive: LDM is slightly slower than SVM_a on 3 data sets, but highly competitive to SVM_a on the other 9 data sets

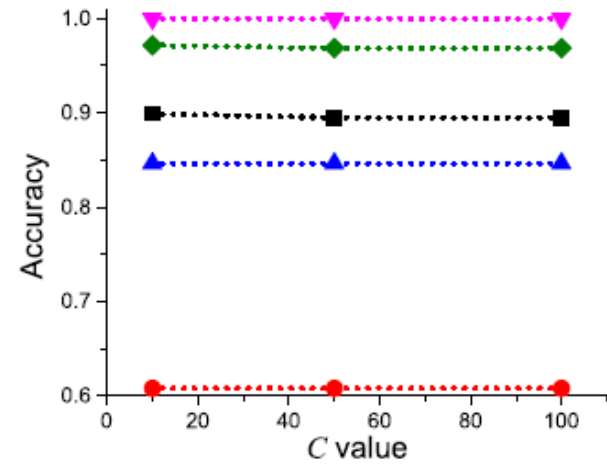
Parameter influence



(a) Influence of λ_1 on LDM



(b) Influence of λ_2 on LDM



(c) Influence of C on LDM

LDM quite insensitive to parameter settings

Joint work with my students



Wei Gao
(高 尉)

W. Gao and Z.-H. Zhou. On the doubt about margin explanation of boosting. Artificial Intelligence, 2013, 203: 1-18.

(arXiv:1009.3613, Sept.2010)

An easy-to-read article:

Z.-H. Zhou. Large margin distribution learning. ANNPR 2014, pp.1-11. (keynote article)



Teng Zhang
(张 腾)

T. Zhang and Z.-H. Zhou. Large margin distribution machine. KDD'14, pp.313-322.

(arXiv:1311.0989, Nov.2013)

Code: http://lamda.nju.edu.cn/code_LDM.ashx

Take-Home Messages

- **Why AdaBoost is less prone to overfitting?**
Margin theory stands
- **What's crucial in margin theory?**
Margin mean and margin variance, together
- **Large margin methods good?**
Large margin distribution methods better
LDM can be applied to generalize all kinds of large margin methods

Thanks!