### Advanced Machine Learning

Deep Boosting



COURANT INSTITUTE & GOOGLE RESEARCH

# Outline

- Model selection.
- Deep boosting.
  - theory.
  - algorithm.
  - experiments.

# **Model Selection**

- Problem: how to select hypothesis set H?
  - *H* too complex, no gen. bound, overfitting.
  - H too simple, gen. bound, but underfitting.

balance between estimation and approx. errors.



# **Estimation and Approximation**

General equality: for any  $h \in H$ , best in class

$$R(h) - R^* = \underbrace{[R(h) - R(h^*)]}_{\text{estimation}} + \underbrace{[R(h^*) - R^*]}_{\text{approximation}}.$$

Approximation: not a random variable, only depends on *H*.

Estimation: only term we can hope to bound; for ERM, bounded by two times gen. bound:

$$R(h_{\text{ERM}}) - R(h^*) = R(h_{\text{ERM}}) - \widehat{R}(h_{\text{ERM}}) + \widehat{R}(h_{\text{ERM}}) - R(h^*)$$
$$\leq R(h_{\text{ERM}}) - \widehat{R}(h_{\text{ERM}}) + \widehat{R}(h^*) - R(h^*)$$
$$\leq 2 \sup_{h \in H} |R(h) - \widehat{R}(h)|.$$



### SRM Guarantee

### Definitions:

- $H_{k(h)}$  simplest hypothesis set containing h.
- $f^*$  the hypothesis returned by SRM:

$$f^* = \underset{h \in H_k, k \ge 1}{\operatorname{argmin}} \widehat{R}_S(h) + R_m(H_k) + \sqrt{\frac{\log k}{m}} = F_k(h).$$

**Theorem:** for any  $\delta > 0$ , with probability at least  $1 - \delta$ ,

$$R(f^*) \le R(h^*) + 2\Re_m(H_{k(h^*)}) + \sqrt{\frac{\log k(h^*)}{m}} + \sqrt{\frac{2\log \frac{3}{\delta}}{m}}$$

# Proof

General bound for all  $h \in H$ :

$$\begin{aligned} \Pr\left[\sup_{h\in H} R(h) - F_{k(h)}(h) > \epsilon\right] \\ &= \Pr\left[\sup_{k\geq 1} \sup_{h\in H_k} R(h) - F_k(h) > \epsilon\right] \\ &\leq \sum_{k=1}^{\infty} \Pr\left[\sup_{h\in H_k} R(h) - F_k(h) > \epsilon\right] \\ &= \sum_{k=1}^{\infty} \Pr\left[\sup_{h\in H_k} R(h) - \hat{R}_S(h) - \Re_m(H_k) > \epsilon + \sqrt{\frac{\log k}{m}}\right] \\ &\leq \sum_{k=1}^{\infty} \exp\left(-2m\left[\epsilon + \sqrt{\frac{\log k}{m}}\right]^2\right) \\ &\leq \sum_{k=1}^{\infty} e^{-2m\epsilon^2} e^{-2\log k} \\ &= e^{-2m\epsilon^2} \sum_{k=1}^{\infty} \frac{1}{k^2} = \frac{\pi^2}{6} e^{-2m\epsilon^2} \leq 2e^{-2m\epsilon^2}. \end{aligned}$$

Advanced Machine Learning - Mohri@

### Proof

Using the union bound and the bound just derived gives:

$$\begin{split} \Pr\left[R(f^*) - R(h^*) - 2\mathfrak{R}_m(H_{k(h^*)}) - \sqrt{\frac{\log k(h^*)}{m}} > \epsilon\right] \\ &\leq \Pr\left[R(f^*) - F_{k(f^*)}(f^*) > \frac{\epsilon}{2}\right] + \Pr\left[F_{k(f^*)}(f^*) - R(h^*) - 2\mathfrak{R}_m(H_{k(h^*)}) - \sqrt{\frac{\log k(h^*)}{m}} > \frac{\epsilon}{2}\right] \\ &\leq 2e^{-\frac{m\epsilon^2}{2}} + \Pr\left[F_{k(h^*)}(h^*) - R(h^*) - 2\mathfrak{R}_m(H_{k(h^*)}) - \sqrt{\frac{\log k(h^*)}{m}} > \frac{\epsilon}{2}\right] \\ &= 2e^{-\frac{m\epsilon^2}{2}} + \Pr\left[\widehat{R}_S(h^*) - R(h^*) - \mathfrak{R}_m(H_{k(h^*)}) > \frac{\epsilon}{2}\right] \\ &= 2e^{-\frac{m\epsilon^2}{2}} + e^{-\frac{m\epsilon^2}{2}} = 3e^{-\frac{m\epsilon^2}{2}}. \end{split}$$

# Remarks

#### SRM bound:

- similar to learning bound when  $k(h^*)$  is known!
- can be extended if approximation error assumed to be small or zero.
- if *H* contains the Bayes classifier, only finitely many hypothesis sets need to be considered in practice.
- restriction: *H* decomposed as countable union of families with converging Rademacher complexity.
- Issues: (1) SRM typically computationally intractable;
   (2) how should we choose H<sub>k</sub>s?

# Voted Risk Minimization

#### Ideas:

- no selection of specific  $H_k$ .
- instead, use all  $H_k$ s:  $h = \sum_{k=1}^p \alpha_k h_k$ ,  $h_k \in H_k$ ,  $\boldsymbol{\alpha} \in \Delta$ .
- hypothesis-dependent penalty:

$$\sum_{k=1}^{p} \alpha_k \Re_m(H_k).$$



# Outline

- Model selection.
- Deep boosting.
  - theory.
  - algorithm.
  - experiments.

# Ensemble Methods in ML

- Combining several base classifiers to create a more accurate one.
  - Bagging (Breiman 1996).
  - AdaBoost (Freund and Schapire 1997).
  - Stacking (Smyth and Wolpert 1999).
  - Bayesian averaging (MacKay 1996).
  - Other averaging schemes e.g., (Freund et al. 2004).
- Often very effective in practice.
- Benefit of favorable learning guarantees.

### **Convex Combinations**

- Base classifier set *H*.
  - boosting stumps.
  - decision trees with limited depth or number of leaves.
- Ensemble combinations: convex hull of base classifier set.  $\operatorname{conv}(H) = \left\{ \sum_{t=1}^{T} \alpha_t h_t \colon \alpha_t \ge 0; \sum_{t=1}^{T} \alpha_t \le 1; \forall t, h_t \in H \right\}.$

### **Ensembles - Margin Bound**

(Bartlett and Mendelson, 2002; Koltchinskii and Panchencko, 2002)

Theorem: let H be a family of real-valued functions. Fix  $\rho > 0$ . Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$ , the following holds for all  $f = \sum_{t=1}^{T} \alpha_t h_t \in \operatorname{conv}(H)$ :

$$R(f) \le \widehat{R}_{S,\rho}(f) + \frac{2}{\rho} \Re_m(H) + \sqrt{\frac{\log \frac{1}{\delta}}{2m}},$$

• where 
$$\widehat{R}_{S,\rho}(f) = \frac{1}{m} \sum_{i=1}^m \mathbb{1}_{y_i f(x_i) \le \rho}$$
.

### Questions

Can we use a much richer or deeper base classifier set?

- richer families needed for difficult tasks in speech and image processing.
- but generalization bound indicates risk of overfitting.

### AdaBoost

(Freund and Schapire, 1997)

Description: coordinate descent applied to

$$F(\alpha) = \sum_{i=1}^{m} e^{-y_i f(x_i)} = \sum_{i=1}^{m} \exp\left(-y_i \sum_{t=1}^{T} \alpha_t h_t(x_i)\right).$$

- Guarantees: ensemble margin bound.
  - but AdaBoost does not maximize the margin!
  - some margin maximizing algorithms such as arc-gv are outperformed by AdaBoost! (Reyzin and Schapire, 2006)

# Suspicions

- Complexity of hypotheses used:
  - arc-gv tends to use deeper decision trees to achieve a larger margin.
- Notion of margin:
  - minimal margin perhaps not the appropriate notion.
  - margin distribution is key.



can we shed more light on these questions?



- Main question: how can we design ensemble algorithms that can succeed even with very deep decision trees or other complex sets?
  - theory.
  - algorithms.
  - experimental results.

### Base Classifier Set H

Decomposition in terms of sub-families or their union.



### **Ensemble Family**

Non-negative linear ensembles  $\mathcal{F} = \operatorname{conv}(\cup_{k=1}^{p} H_k)$ :



### Ideas

- Use hypotheses drawn from  $H_k$ s with larger ks but allocate more weight to hypotheses drawn from smaller ks.
  - how can we determine quantitatively the amounts of mixture weights apportioned to different families?
  - can we provide learning guarantees guiding these choices?

# Learning Guarantee

(Cortes, MM, and Syed, 2014)

Theorem: Fix  $\rho > 0$ . Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$ , the following holds for all  $f = \sum_{t=1}^{T} \alpha_t h_t \in \mathcal{F}$ :

$$R(f) \le \widehat{R}_{S,\rho}(f) + \frac{4}{\rho} \sum_{t=1}^{T} \alpha_t \Re_m(H_{k_t}) + \widetilde{O}\left(\sqrt{\frac{\log p}{\rho^2 m}}\right)$$

### Consequences

- Complexity term with explicit dependency on mixture weights.
  - quantitative guide for controlling weights assigned to more complex sub-families.
  - bound can be used to inspire, or directly define an ensemble algorithm.

### Set-Up

- $H_1, \ldots, H_p$ : disjoint sub-families of functions taking values in [-1, +1].
- Further assumption (not necessary): symmetric subfamilies, i.e.  $h \in H_k \Leftrightarrow -h \in H_k$ .
- Notation:

• 
$$r_j = \mathfrak{R}_m(H_{k_j})$$
 with  $h_j \in H_{k_j}$ .

### Derivation

Learning bound suggests seeking  $\alpha \ge 0$  with  $\sum_{t=1}^{T} \alpha_t \le 1$  to minimize

$$\frac{1}{m} \sum_{i=1}^{m} 1_{y_i \sum_{t=1}^{T} \alpha_t h_t(x_i) \le \rho} + \frac{4}{\rho} \sum_{t=1}^{T} \alpha_t r_t.$$

# **Convex Surrogates**

- Let  $u \mapsto \Phi(-u)$  be a decreasing convex function upper bounding  $u \mapsto 1_{u \leq 0}$ , with  $\Phi$  differentiable.
- Two principal choices:
  - Exponential loss:  $\Phi(-u) = \exp(-u)$ .
  - Logistic loss:  $\Phi(-u) = \log_2(1 + \exp(-u))$ .

### Optimization Problem (Cortes, MM, and Syed, 2014)

Moving the constraint to the objective and using the fact that the sub-families are symmetric leads to:

$$\min_{\boldsymbol{\alpha}\in\mathbb{R}^N}\frac{1}{m}\sum_{i=1}^m \Phi\left(1-y_i\sum_{j=1}^N\alpha_jh_j(x_i)\right) + \sum_{t=1}^N(\lambda r_j + \beta)|\alpha_j|,$$

where  $\lambda, \beta \geq 0$ , and for each hypothesis, keep either *h* or *-h*.

# DeepBoost Algorithm

- Coordinate descent applied to convex objective.
  - non-differentiable function.
  - definition of maximum coordinate descent.



# **Direction & Step**

Maximum direction: definition based on the error

$$\epsilon_{t,j} = \frac{1}{2} \Big[ 1 - \mathop{\mathrm{E}}_{i \sim \mathcal{D}_t} [y_i h_j(x_i)] \Big],$$

where  $D_t$  is the distribution over sample at iteration t.

- Step:
  - closed-form expressions for exponential and logistic losses.
  - general case: line search.

### Pseudocode

 $DEEPBOOST(S = ((x_1, y_1), ..., (x_m, y_m)))$ for  $i \leftarrow 1$  to m do 1  $D_1(i) \leftarrow \frac{1}{m}$ 2 for  $t \leftarrow 1$  to T do 3 for  $j \leftarrow 1$  to N do 4 if  $(\alpha_{t-1,j} \neq 0)$  then 5  $d_j \leftarrow \left(\epsilon_{t,j} - \frac{1}{2}\right) + \operatorname{sgn}(\alpha_{t-1,j}) \frac{\lambda_j m}{2S_t}$  $\Lambda_i = \lambda r_i + \beta.$ 6 elseif  $\left(\left|\epsilon_{t,j}-\frac{1}{2}\right| \leq \frac{\Delta_j m}{2S_j}\right)$  then 7  $d_i \leftarrow 0$ 8 else  $d_j \leftarrow (\epsilon_{t,j} - \frac{1}{2}) - \operatorname{sgn}(\epsilon_{t,j} - \frac{1}{2}) \frac{\Lambda_j m}{2S_i}$ 9  $k \leftarrow \operatorname{argmax} |d_i|$ 10 $j \in [1,N]$ 11  $\epsilon_t \leftarrow \epsilon_{t,k}$ if  $\left( \left| (1 - \epsilon_t) e^{\alpha_{t-1,k}} - \epsilon_t e^{-\alpha_{t-1,k}} \right| \leq \frac{\Lambda_k m}{S_t} \right)$  then 12  $\begin{array}{l} \eta_t \leftarrow -\alpha_{t-1,k} \\ \text{elseif}\left((1-\epsilon_t)e^{\alpha_{t-1,k}} - \epsilon_i e^{-\alpha_{t-1,k}} > \frac{\Lambda_k m}{S_t}\right) \text{ then} \end{array}$ 13 14  $\eta_t \leftarrow \log \left[ -\frac{\Lambda_k m}{2\epsilon_t S_t} + \sqrt{\left[\frac{\Lambda_k m}{2\epsilon_t S_t}\right]^2 + \frac{1-\epsilon_t}{\epsilon_t}} \right]$ 15 else  $\eta_t \leftarrow \log \left[ + \frac{\Lambda_k m}{2\epsilon_t S_t} + \sqrt{\left[\frac{\Lambda_k m}{2\epsilon_t S_t}\right]^2 + \frac{1-\epsilon_t}{\epsilon_t}} \right]$ 16 17  $\boldsymbol{\alpha}_t \leftarrow \boldsymbol{\alpha}_{t-1} + \eta_t \mathbf{e}_k$  $S_{t+1} \leftarrow \sum_{i=1}^{m} \Phi' \left( 1 - y_i \sum_{j=1}^{N} \alpha_{t,j} h_j(x_i) \right)$ 18 19 for  $i \leftarrow 1$  to m do  $D_{t+1}(i) \leftarrow rac{\Phi'\left(1-y_i\sum_{j=1}^N lpha_{i,j}h_j(x_i)
ight)}{S_{t+1}}$ 20  $f \leftarrow \sum_{j=1}^{N} \alpha_{t,j} h_j$ 21return f 22

# **Connections with Previous Work**

- For  $\lambda=\beta=0$  , DeepBoost coincides with
  - AdaBoost (Freund and Schapire 1997), run with union of subfamilies, for the exponential loss.
  - additive Logistic Regression (Friedman et al., 1998), run with union of sub-families, for the logistic loss.
- For  $\lambda = 0$  and  $\beta \neq 0$ , DeepBoost coincides with
  - L1-regularized AdaBoost (Raetsch, Mika, and Warmuth 2001), for the exponential loss.
  - coincides with L1-regularized Logistic Regression (Duchi and Singer 2009), for the logistic loss.

# Rad. Complexity Estimates

Benefit of data-dependent analysis:

- empirical estimates of each  $\mathfrak{R}_m(H_k)$ .
- example: for kernel function  $K_k$ ,

$$\widehat{\mathfrak{R}}_{S}(H_{k}) \leq \frac{\sqrt{\operatorname{Tr}[\mathbf{K}_{k}]}}{m}$$

alternatively, upper bounds in terms of growth functions,

$$\Re_m(H_k) \le \sqrt{\frac{2\log \Pi_{H_k}(m)}{m}}$$

# Experiments (1)

Family of base classifiers defined by boosting stumps:

- boosting stumps  $H_1^{\text{stumps}}$  (threshold functions).
  - in dimension d ,  $\Pi_{H_1^{\mathrm{stumps}}}(m) \leq 2md$  , thus

$$\Re_m(H_1^{\text{stumps}}) \le \sqrt{\frac{2\log(2md)}{m}}$$

- decision trees of depth 2,  $H_2^{\rm stumps}$ , with the same question at the internal nodes of depth 1.
  - in dimension d ,  $\Pi_{H_2^{\rm stumps}}(m) \leq (2m)^2 \frac{d(d-1)}{2}$  , thus

$$\mathfrak{R}_m(H_2^{\text{stumps}}) \le \sqrt{\frac{2\log(2m^2d(d-1))}{m}}$$

# Experiments (1)

- Base classifier set:  $H_1^{\text{stumps}} \cup H_2^{\text{stumps}}$ .
- Data sets:
  - same UCI Irvine data sets as (Breiman 1999) and (Reyzin and Schapire 2006).
  - OCR data sets used by (Reyzin and Schapire 2006): ocr17, ocr49.
  - MNIST data sets: ocr17-mnist, ocr49-mnist.
- Experiments with exponential loss.
- Comparison with AdaBoost and AdaBoost-L1.

### **Data Statistics**

	breastcancer	ionosphere	german (numeric)
Examples	699	351	1000
Attributes	9	34	24

	diabetes	ocr17	ocr49
Examples	768	2000	2000
Attributes	8	196	196

	ocr17-mnist	ocr49-mnist
Examples	15170	13782
Attributes	400	400

# Experiments - Stumps Exp Loss

(Cortes, MM, and Syed, 2014)

Table 1. Results for boosted decision stumps and the exponential loss function.

	AdaBoost	AdaBoost				AdaBoost	AdaBoost		
breastcancer	$H_1^{\text{stimps}}$	$H_2^{\rm stumps}$	AdaBoost-L1	DeepBoost	ocr17	$H_1^{\text{starsps}}$	$H_2^{\rm stamps}$	AdaBoost-L1	DeepBoost
Error	0.0429	0.0437	0.0408	0.0373	Error	0.0085	0.008	0.0075	0.0070
(std dev)	(0.0248)	(0.0214)	(0.0223)	(0.0225)	(std dev)	0.0072	0.0054	0.0068	(0.0048)
Avg tree size	1	2	1.436	1.215	Avg tree size	1	2	1.036	1.369
Avg no. of trees	100	100	43.6	21.6	Avg no. of trees	100	100	37.8	36.9
ionosphere	AdaBoost H <sup>etimps</sup>	AdaBoost H <sup>stumys</sup> <sub>2</sub>	AdaBoost-L.	DeepBoost	ocr49	AdaBoost H <sup>starsps</sup>	AdaBoost H <sup>stumps</sup> <sub>2</sub>	AdaBoost-L1	DeepBoost
Error	0.1014	0.075	0.0708	0.0638	Error	0.0555	0.032	0.03	0.0275
(std dev)	(0.0414)	(0.0413)	(0.0331)	(0.0394)	(std dcv)	0.0167	0.0114	0.0122	(0.0095)
Avg tree size	1	2	1.392	1.168	Avg tree size	1	2	1.99	1.96
Avg no. of trees	100	100	39.35	17.45	Avg no. of trees	100	100	99.3	96
	AdaBoost	AdaBoost				AdaBoost	AdaBoost		
german	$H_1^{\text{stumps}}$	$H_2^{\text{stumps}}$	AdaBoost-L1	DeepBoost	ocr17-mnist	$H_1^{\text{stamps}}$	$H_2^{stamps}$	AdaBoost-L1	DeepBoost
Error	0.243	0.2505	0.2455	0.2395	Error	0.0056	0.0048	0.0046	0.0040
(std dev)	(0.0445)	(0.0487)	(0.0438)	(0.0462)	(std dev)	0.0017	0.0014	0.0013	(0.0014)
Avg tree size	1	2	1.54	1.76	Avg tree size	1	2	2	1.99
Avg no. of trees	100	100	54.1	76.5	Avg no. of trees	100	100	100	100
	AdaBoost	AdaBoost				AdaBoost	AdaBoost		
diabetes	$H_1^{\text{stimps}}$	$H_2^{\text{stumps}}$	AdaBoost-L1	DeepBoost	ocr49 mnist	$H_1^{sumps}$	$H_2^{sumps}$	AdaBoost-L1	DeepBoost
Error	0.253	0.260	0.254	0.253	Error	0.0414	0.0209	0.0200	0.0177
(std dev)	(0.0330)	(0.0518)	(0.04868)	(0.0510)	(std dev)	0.00539	0.00521	0.00408	(0.00438)
Avg tree size	1	2	1.9975	1.9975	Avg tree size	1	2	1.9975	1.9975
Avg no. of trees	100	100	100	100	Avg no. of trees	100	100	100	100

# Experiments (2)

Family of base classifiers defined by decision trees of depth k. For trees with at most n nodes:

$$\Re_m(\mathsf{T}_n) \le \sqrt{\frac{(4n+2)\log_2(d+2)\log(m+1)}{m}}.$$

- Base classifier set:  $\cup_{k=1}^{K} H_k^{\text{trees}}$  .
- Same data sets as with Experiments (1).
- Both exponential and logistic loss.
- Comparison with AdaBoost and AdaBoost-L1.

# **Experiments - Trees Exp Loss**

#### (Cortes, MM, and Syed, 2014)

DeepBoost 0.002

(0.00100)

26.061.8

breastcancer	AdaBoost	AdaBoost-L1	DeepBoost	ocr17	AdaBoost	AdaBoost-L1
Error	0.0267	0.0264	0.0243	Error	0.004	0.003
(std dev)	(0.00841)	(0.0098)	(0.00797)	(std dev)	(0.00316)	(0.00100)
Avg tree size	29.1	28.9	20.9	Avg tree size	15.0	30.4
Avg no. of trees	67.1	51.7	55.9	Avg no. of trees	88.3	65.3

ionosphere	AdaBoost	AdaBoost-L1	DeepBoost
Error	0.0661	0.0657	0.0501
(std dev)	(0.0315)	(0.0257)	(0.0316)
Avg tree size	29.8	31.4	26.1
Avg no. of trees	75.0	69.4	50.0

ocr49	AdaBoost	AdaBoost-L1	DeepBoost
Error	0.0180	0.0175	0.0175
(std dev)	(0.00555)	(0.00357)	(0.00510)
Avg tree size	30.9	62.1	30.2
Avg no. of trees	92.4	89.0	83.0

german	AdaBoost	AdaBoost-L1	DeepBoost
Error	0.239	0.239	0.234
(std dev)	(0.0165)	(0.0201)	(0.0148)
Avg tree size	3	7	16.0
Avg no. of trees	91.3	87.5	14.1

ocr17-mnist	AdaBoost	AdaBoost-L1	DeepBoost
Error	0.00471	0.00471	0.00409
(std dev)	(0.0022)	(0.0021)	(0.0021)
Avg tree size	15	33.4	22.1
Avg no. of trees	88.7	66.8	59.2

diabetes	AdaBoost	AdaBoost-L1	DeepBoost
Error	0.249	0.240	0.230
(std dev)	(0.0272)	(0.0313)	(0.0399)
Avg tree size	3	3	5.37
Avg no. of trees	45.2	28.0	19.0

ocr49-mnist	AdaBoost	AdaBoost-L1	DeepBoost
Error	0.0198	0.0197	0.0182
(std dev)	(0.00500)	(0.00512)	(0.00551)
Avg tree size	29.9	66.3	30.1
Avg no. of trees	82.4	81.1	80.9

# **Experiments - Trees Log Loss**

#### (Cortes, MM, and Syed, 2014)

breastcancer	LogReg	LogReg-L1	DeepBoost		ocr17	LogReg	LogReg-L1	DeepBoost
Error	0.0351	0.0264	0.0264		Error	0.00300	0.00400	0.00250
(std dev)	(0.0101)	(0.0120)	(0.00876)		(std dev)	(0.00100)	(0.00141)	(0.000866)
Avg tree size	15	59.9	14.0		Avg tree size	15.0	7	22.1
Avg no. of trees	65.3	16.0	23.8		Avg no. of trees	75.3	53.8	25.8
ionosphere	LogReg	LogReg-L1	DeepBoost		ocr49	LogReg	LogReg-L1	DeepBoost
Error	0.074	0.060	0.043		Error	0.0205	0.0200	0.0170
(std dev)	(0.0236)	(0.0219)	(0.0188)		(std dev)	(0.00654)	(0.00245)	(0.00361)
Avg tree size	7	30.0	18.4		Avg tree size	31.0	31.0	63.2
Avg no. of trees	44.7	25.3	29.5		Avg no. of trees	63.5	54.0	37.0
				-				
german	LogReg	LogReg-L1	DeepBoost		ocr17-mnist	m LogReg	LogReg-L1	DeepBoost
Error	0.233	0.232	0.225		Error	0.00422	0.00417	0.00399
(std dev)	(0.0114)	(0.0123)	(0.0103)		(std dev)	(0.00191)	(0.00188)	(0.00211)
Avg tree size	7	7	14.4		Avg tree size	15	15	25.9
Avg no. of trees	72.8	66.8	67.8		Avg no. of trees	71.4	55.6	27.6
diabetes	LogReg	LogReg-L1	DeepBoost		ocr49-mnist	LogReg	LogReg-L1	$\overline{\text{DeepBoost}}$
Error	0.250	0.246	0.246		Error	0.0211	0.0201	0.0201
(std dev)	(0.0374)	(0.0356)	(0.0356)		(std dev)	(0.00412)	(0.00433)	(0.00411)
Avg tree size	3	3	3		Avg tree size	28.7	33.5	72.8

Avg no. of trees

79.3

61.7

46.0

45.5

45.5

Avg no. of trees

41.9

### Margin Distribution



# Multi-Class Learning Guarantee

(Kuznetsov, MM, and Syed, 2014)

Theorem: Fix $\rho > 0$ . Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$ , the following holds for all  $f = \sum_{t=1}^{T} \alpha_t h_t \in \mathcal{F}$ :

$$R(f) \le \widehat{R}_{S,\rho}(f) + \frac{8c}{\rho} \sum_{t=1}^{T} \alpha_t \Re_m(\Pi_1(H_{k_t})) + O\left(\sqrt{\frac{\log p}{\rho^2 m} \log\left[\frac{\rho^2 c^2 m}{4\log p}\right]}\right)$$

with c number of classes.

• and 
$$\Pi_1(H_k) = \{x \mapsto h(x,y) \colon y \in \mathcal{Y}, h \in H_k\}.$$

# Extension to Multi-Class

- Similar data-dependent learning guarantee proven for the multi-class setting.
  - bound depending on mixture weights and complexity of sub-families.
- Deep Boosting algorithm for multi-class:
  - similar extension taking into account the complexities of sub-families.
  - several variants depending on number of classes.
  - different possible loss functions for each variant.

# **Other Related Algorithms**

Structural Maxent models (Cortes, Kuznetsov, MM, and Syed, ICML 2015): feature functions chosen from a union of very complex families.



# **Other Related Algorithms**

Deep Cascades (DeSalvo, MM, and Syed, ALT 2015): cascade of predictors with leaf predictors and node questions selected from very rich families.



# Conclusion

- Deep Boosting: ensemble learning with increasingly complex families.
  - data-dependent theoretical analysis.
  - algorithm based on learning bound.
  - extension to multi-class.
  - ranking and other losses.
  - enhancement of many existing algorithms.
  - compares favorably to AdaBoost and additive Logistic Regression or their L1-regularized variants in experiments.

### References

- Bartlett, Peter L. and Mendelson, Shahar. Rademacher and Gaussian complexities: Risk bounds and structural results. JMLR, 3, 2002.
- Breiman, Leo. Bagging predictors. Machine Learning, 24 (2):123–140, 1996.
- Breiman, Leo. Prediction games and arcing algorithms. Neural Computation, 11(7):1493–1517, 1999.
- Corinna Cortes, Mehryar Mohri, and Umar Syed. Deep boosting. In ICML, 2014.
- Duchi, John C. and Singer, Yoram. Boosting with structural sparsity. In ICML, pp. 38, 2009.
- Yoav Freund and Robert E. Schapire. A decision-theoretic generalization of online learning and an application to boosting. Journal of Computer System Sciences, 55(1):119–139, 1997.



- Koltchinskii, Vladmir and Panchenko, Dmitry. Empirical margin distributions and bounding the generalization error of combined classifiers. Annals of Statistics, 30, 2002.
- Vitaly Kuznetsov, Mehryar Mohri, and Umar Syed. Multi-class deep boosting. In NIPS, 2014.
- Raetsch, Gunnar, Onoda, Takashi, and Mueller, Klaus- Robert. Soft margins for AdaBoost. Machine Learning, 42(3):287–320, 2001.
- Raetsch, Gunnar, Mika, Sebastian, and Warmuth, Manfred K. On the convergence of leveraging. In NIPS, pp. 487–494, 2001.

### References

- Reyzin, Lev and Schapire, Robert E. How boosting the margin can also boost classifier complexity. In ICML, pp. 753–760, 2006.
- Schapire, Robert E., Freund, Yoav, Bartlett, Peter, and Lee, Wee Sun. Boosting the margin: A new explanation for the effectiveness of voting methods. In ICML, pp. 322– 330, 1997.
- Vladimir N. Vapnik. The Nature of Statistical Learning Theory. Springer-Verlag, New York, 1995.
- Vladimir N. Vapnik and Alexey Ya.Chervonenkis. Theory of Pattern Recognition. Nauka, Moscow, 1974.