

Top-down vs. bottom-up methods for hierarchical classification

Claudio Gentile

DICOM

Universita' dell'Insubria, Italy

`claudio.gentile@uninsubria.it`

joint work with:

N. Cesa-Bianchi

L. Zaniboni

State of the art:

Many hierarchical classification models/algorithms

[KS97, DC00, SL01, RS02, HCC03, Ro⁺05, CGZ06,...]

Top-down vs bottom-up / local vs global / on-line vs batch ...

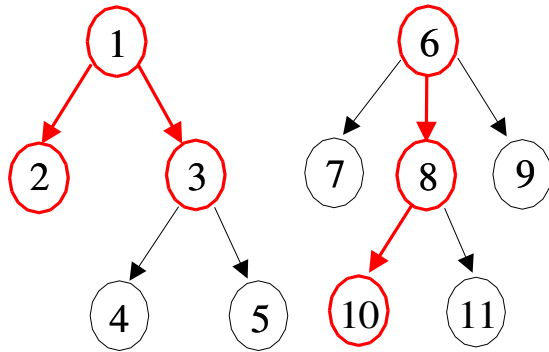
Outline:

- Hierarchical classification framework
- Bottom-up alg B-SVM:
Combines Bayes optimal (labels are **subtrees**) with SVM
- Report experimental comparison between B-SVM and local top-down H-SVM on artificial and real-world medium-size datasets
- On-line learning: model, alg, regret analysis

Hierarchical classification model/1

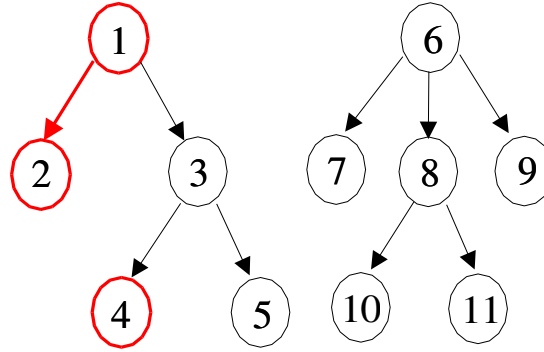
Taxonomy: tree forest ($N = 11$ nodes)

Legal multilabel: union of paths



$$\mathbf{v}=(1, 1, 1, 0, 0, 1, 0, 1, 0, 1, 0)$$

Illegal multilabel: rest



$$\mathbf{v}=(1, 1, 0, 1, 0, 0, 0, 0, 0, 0, 0)$$

Example $(\mathbf{x}, \mathbf{v}) \in \mathbb{R}^d \times \{0, 1\}^N$

Hierarchical classification model/2

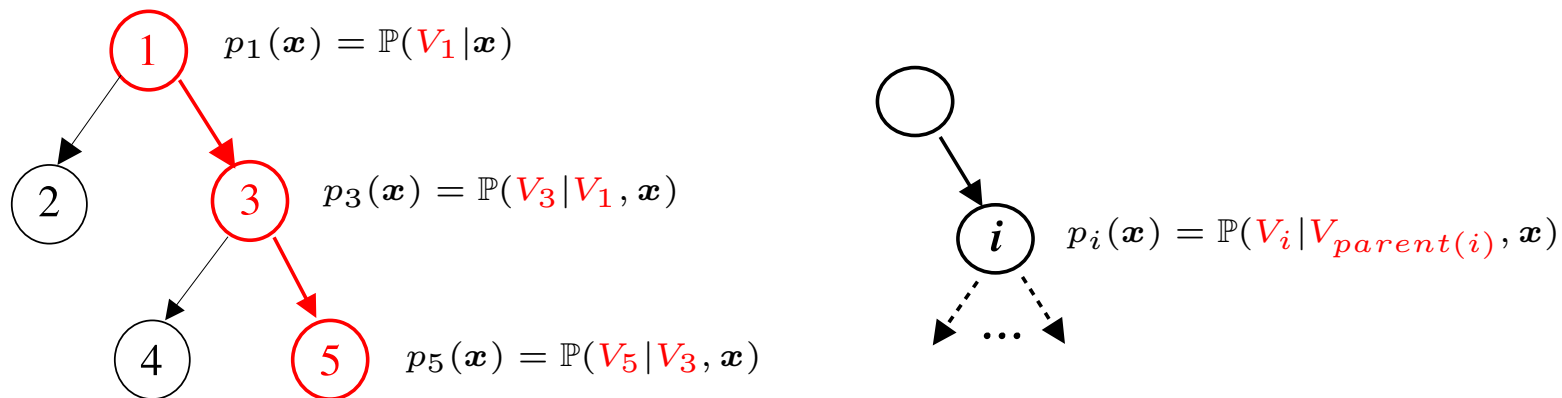
Generation of multilabels

Instance $\mathbf{x} \in \mathbb{R}^d$

Multilabel $\mathbf{V} = (V_1, V_2, \dots, V_N) \in \{0, 1\}^N$

$\mathbf{V} \mid \mathbf{x} \sim \prod_{i=1}^N \mathbb{P}(V_i \mid V_{parent(i)}, \mathbf{x})$

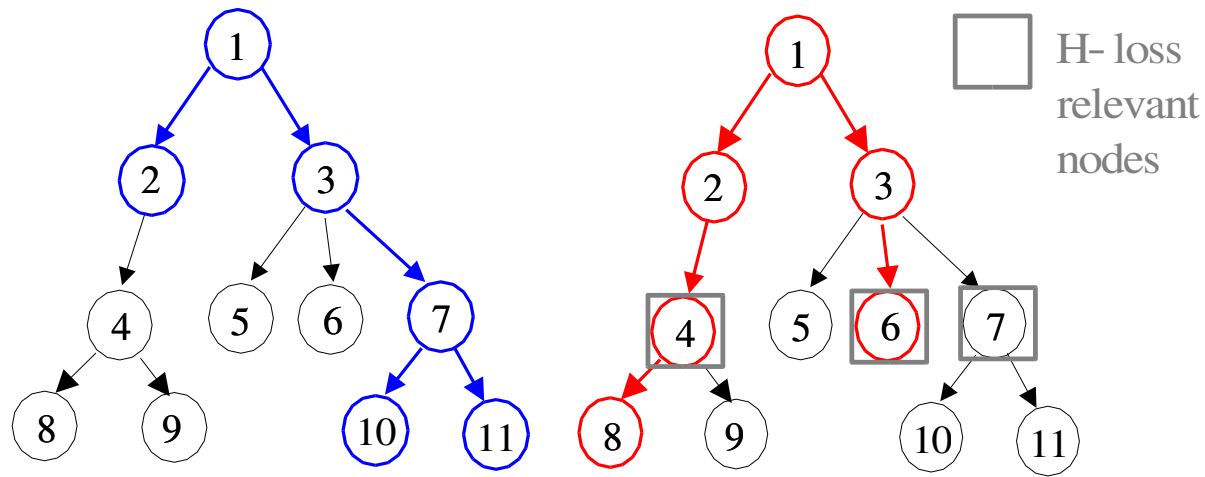
Children's labels are indep. given parent's (no correlation)



$$\mathbb{P}(V_i = 1 \mid V_{parent(i)} = 0, \mathbf{x}) = 0 \quad \forall \mathbf{x}$$

Hierarchical classification model/3

Hierarchical loss (H-loss) [CBGZ06]



prediction

label

$$\hat{\mathbf{Y}} = (\hat{y}_1, \hat{y}_2, \dots, \hat{y}_N) \quad \mathbf{V} = (V_1, V_2, \dots, V_N)$$

H-Loss:

$$\ell_H(\hat{\mathbf{Y}}, \mathbf{V}) = \sum_{i=1}^N c_i \{ \hat{y}_i \neq V_i \wedge \forall j \in \text{ancest}(i) : \hat{y}_j = V_j \}$$

(constant) cost coefficients

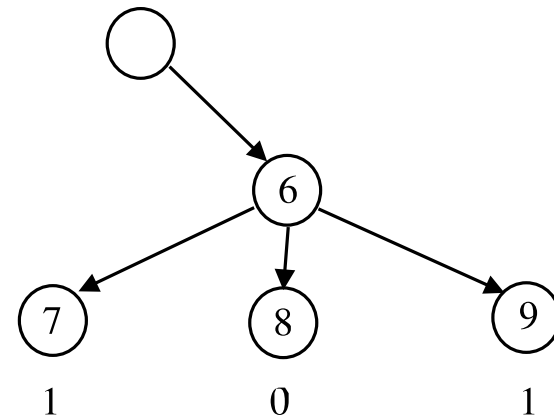
Hierarchical classification model/4

Bayes optimal classifier

$$\mathbf{y}^* = (y_1^*, \dots, y_N^*) = \operatorname{argmin}_{\mathbf{y} \in \{0,1\}^N} \mathbb{E} [\ell_H(\mathbf{y}, \mathbf{V}) \mid \mathbf{x}]$$

Bottom-up message passing

Given $p_i = p_i(\mathbf{x})$



Leaves: $y_i^* = \{p_i \geq 1/2\}$

Hierarchical classification model/4

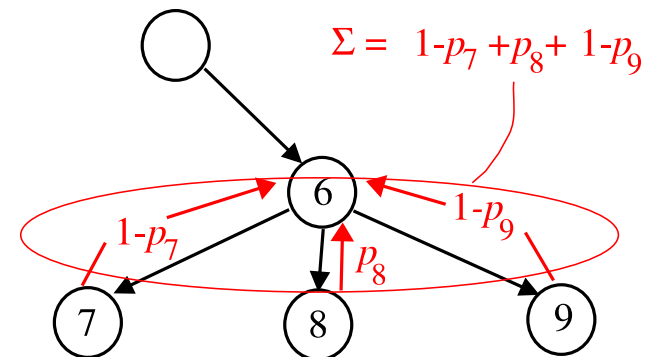
Bayes optimal classifier

$$\mathbf{y}^* = (y_1^*, \dots, y_N^*) = \operatorname{argmin}_{\mathbf{y} \in \{0,1\}^N} \mathbb{E} [\ell_H(\mathbf{y}, \mathbf{V}) \mid \mathbf{x}]$$

Bottom-up message passing

Given $p_i = p_i(\mathbf{x})$

Messages: $y_i^*(1-p_i) + (1-y_i^*)p_i$



Hierarchical classification model/4

Bayes optimal classifier

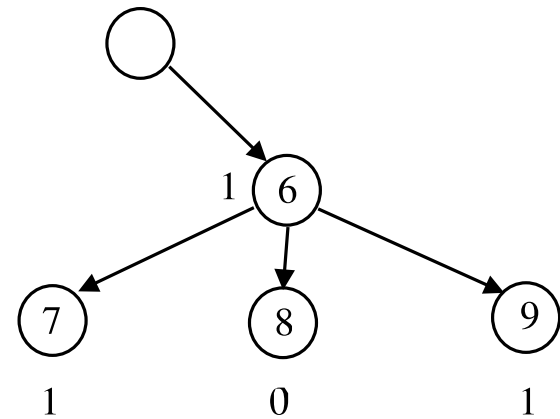
$$\mathbf{y}^* = (y_1^*, \dots, y_N^*) = \operatorname{argmin}_{\mathbf{y} \in \{0,1\}^N} \mathbb{E} [\ell_H(\mathbf{y}, \mathbf{V}) \mid \mathbf{x}]$$

Bottom-up message passing

Given $p_i = p_i(\mathbf{x})$

Nodes: $y_i^* = \{p_i \geq 1/(2-\Sigma)\}$

If $y_i^* = 0$ then $y_j^* = 0 \quad \forall j \in \text{subtree}(i)$



Hierarchical classification model/4

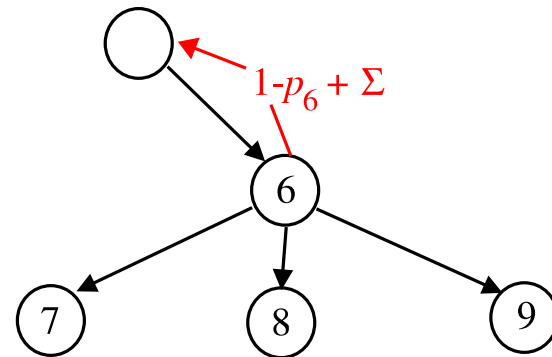
Bayes optimal classifier

$$\mathbf{y}^* = (y_1^*, \dots, y_N^*) = \operatorname{argmin}_{\mathbf{y} \in \{0,1\}^N} \mathbb{E} [\ell_H(\mathbf{y}, \mathbf{V}) \mid \mathbf{x}]$$

Bottom-up message passing

Given $p_i = p_i(\mathbf{x})$

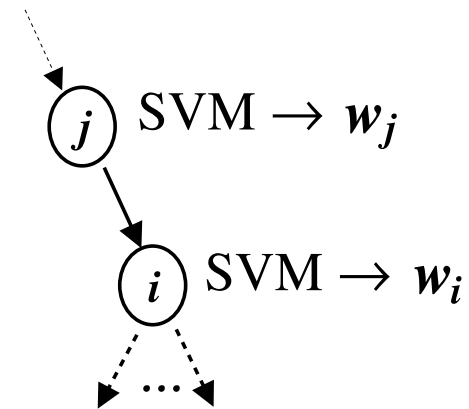
Messages: $y_i^*(1-p_i) + (1-y_i^*) p_i + \Sigma$



B-SVM algorithm (bottom-up)/1

SVM sitting at each node

$$\mathbf{w}_i \leftarrow \text{SVM}(\{(\mathbf{x}, \mathbf{v}) : v_{parent(i)} = 1\})$$



$$p_i(\mathbf{x}) \simeq \hat{p}_i(\mathbf{x}) = \frac{1}{(1 + e^{\alpha_i \mathbf{w}_i \cdot \mathbf{x}})}$$

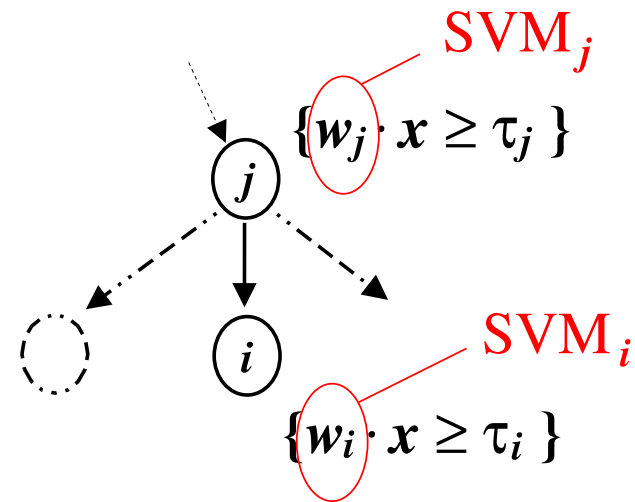
Platt's method [Pl99] to fit parameters $\alpha_1, \dots, \alpha_N$ via cross-validation on training set

Then **bottom-up** label assignment

B-SVM algorithm (bottom-up)/2

Results in SVM with
modified thresholds $\tau_j > 0$
 $\{\mathbf{w}_j \cdot \mathbf{x} \geq \tau_j\}$

$$\hat{p}_i(\mathbf{x}) = \frac{1}{(1 + e^{\alpha_i \mathbf{w}_i \cdot \mathbf{x}})}$$



τ_j depends on behavior of children i :

- $\tau_j \gg 0$ when $\hat{p}_i(\mathbf{x}) \approx 1/2$
- $\tau_j \approx 0$ when $|\hat{p}_i(\mathbf{x}) - 1/2| \approx 1/2$

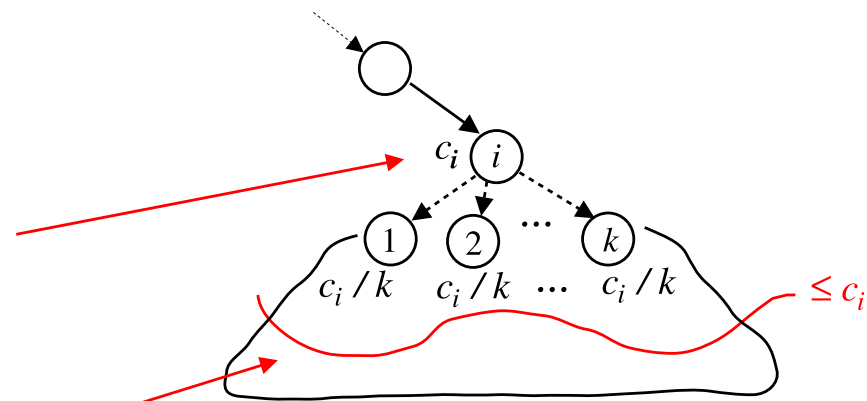
B-SVM algorithm (bottom-up)/3: Choice of H-loss coefficients c_i

B-SVM training

Cost of node i

\geq

Cost of subtree rooted at i



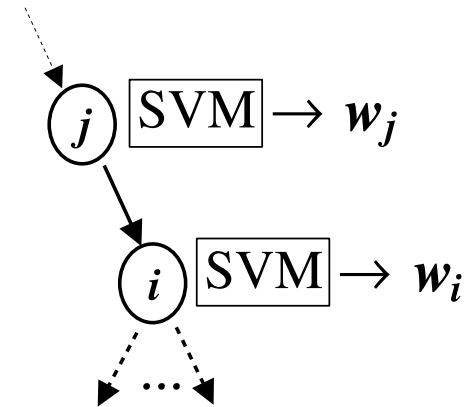
Removes bias towards short paths labelling

H-SVM algorithm (top-down) [HCC03,CGZ06], ...

SVM sitting at each node

$$\mathbf{w}_i \leftarrow \text{SVM}(\{(\mathbf{x}, \mathbf{v}) : v_{\text{parent}(i)} = 1\})$$

(same training as B-SVM)



Top-down label assignment

$$\hat{y}_i = \begin{cases} \{\mathbf{w}_i^\top \mathbf{x} \geq 0\} & \text{if } i \text{ is root} \\ \{\mathbf{w}_i^\top \mathbf{x} \geq 0\} & \text{if } i \text{ is not root \& } \hat{y}_j = 1 \\ 0 & \text{otherwise} \end{cases}$$

Independent of c_i !

Experimental results/1: Datasets

- Reuters Corpus Volume 1 (**RCV1**), first 100,000 docs: 101 nodes, 4 trees, height 3
Avg # of paths/label: 1.5
5 ordered chunks
- **OHSUMED** corpus of medical abstracts, subtree rooted in "Quality of Health Care": 55,503 docs, 94 nodes, height 4
Avg # of paths/label: 1.53
5 random splits 40,000/15,503
- Two synthetic datasets: 40,000 examples, 3 complete ternary trees, 39 nodes, height 2
Avg # of paths/label: **SYNTH1**: 2.66; **SYNTH2**: 1.28
4 ordered chunks

Experimental results/2: Avg test error (H-loss)

DATASET	SYNTH1	SYNTH2
H-SVM	1.454 (± 0.007)	0.350 (± 0.007)
B-SVM	1.269 (± 0.008)	0.322 (± 0.003)

DATASET	RCV1	OHSUMED
H-SVM	0.716 (± 0.024)	1.171 (± 0.005)
B-SVM	0.712 (± 0.023)	1.158 (± 0.005)

Experimental results/3:

Test error (H-loss) per chunk

RCV1					
CHUNK #	1	2	3	4	5
H-SVM	0.702	0.727	0.727	0.741	0.682
B-SVM	0.701	0.727	0.712	0.741	0.681

OHSUMED					
CHUNK #	1	2	3	4	5
H-SVM	1.176	1.163	1.171	1.170	1.171
B-SVM	1.164	1.152	1.158	1.157	1.161

Experimental results/4: Breakdown on levels

DEPTH #	SYNTH1		SYNTH2	
	H-SVM	B-SVM	H-SVM	B-SVM
0	0.432	0.445	0.123	0.121
1	0.572	0.519	0.184	0.166
2	0.757	0.678	0.187	0.183

DEPTH #	RCV1		OHSUMED	
	H-SVM	B-SVM	H-SVM	B-SVM
0	0.190	0.192	1.010	1.018
1	0.471	0.474	0.316	0.297
2	0.131	0.126	0.160	0.152
3	0.097	0.093	0.099	0.101
4			0.413	0.389

Experimental results/5: Comments

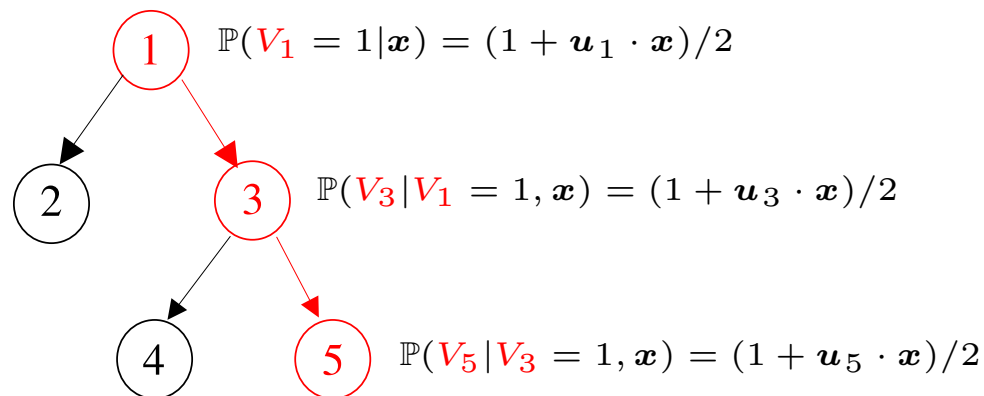
- B-SVM beats H-SVM in all cases:
 - clear on **SYNTH1** and **SYNTH2**
 - significant on **OHSUMED**
 - marginal on **RCV1**
- Chunk-wise (**RCV1**, **OHSUMED**):
B-SVM beats H-SVM on all chunks
- Across levels:
B-SVM beats H-SVM at deeper levels

On-line hierarchical classification [CGZ06]/1

Parametric model

Instance $\mathbf{x} \in \mathbb{R}^d : \|\mathbf{x}\| = 1$

$$\mathbb{P}(V_i = 1 \mid V_{parent(i)} = 1, \mathbf{x}) = (1 + \mathbf{u}_i \cdot \mathbf{x})/2 \in [0, 1]$$



Parameters \mathbf{u}_i

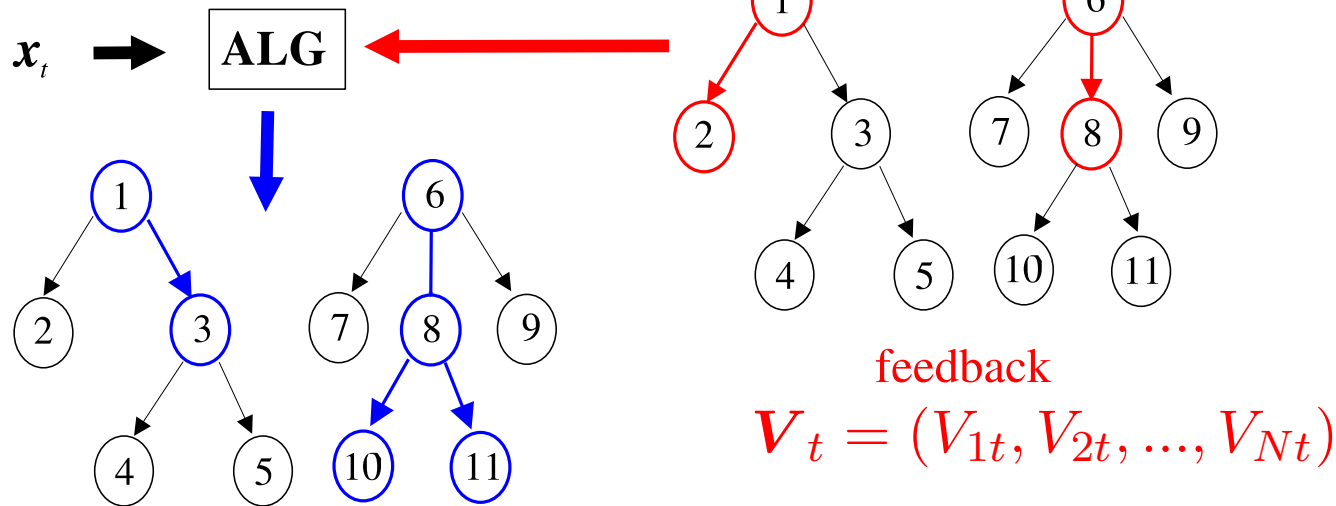
$$\|\mathbf{u}_i\| = 1, i = 1 \dots N$$

$$\mathbb{P}(V_i = 1 \mid V_{parent(i)} = 0, \mathbf{x}) = 0 \quad \forall \mathbf{x}$$

On-line hierarchical classification/2

On-line learning protocol

At time t :



$$\hat{\mathbf{Y}} = (\hat{y}_{1t}, \hat{y}_{2t}, \dots, \hat{y}_{Nt}) \in \{0, 1\}^N$$

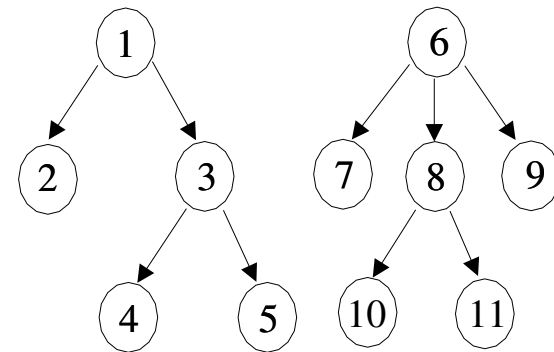
Accuracy on H-loss

On-line hierarchical classification/3

Comparison predictor and regret

Top-down comparator:

$$y_i = \begin{cases} \{\mathbf{u}_i \cdot \mathbf{x} \geq 0\} & \text{if } i \text{ is root} \\ \{\mathbf{u}_i \cdot \mathbf{x} \geq 0\} & \text{if } i \text{ is **not** root \& } \\ & y_{parent(i)} = 1 \\ 0 & \text{otherwise} \end{cases}$$



Not Bayes optimal (independent of c_i !)

Cumulative regret:

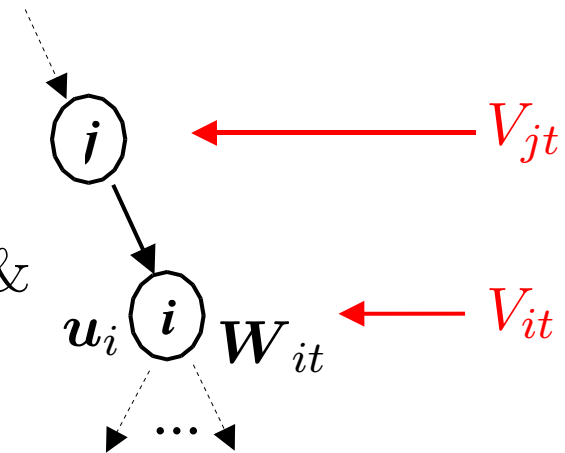
$$\sum_t \left(\mathbb{E}[\ell_H(\hat{\mathbf{Y}}_t, \mathbf{V}_t)] - \mathbb{E}[\ell_H(\mathbf{Y}_t, \mathbf{V}_t)] \right)$$

On-line hierarchical classification/4: Algorithm/1

- **Storage:** Keeps at node i weight vector $\mathbf{W}_{it} \simeq \mathbf{u}_i$

- **Prediction:**

$$\hat{y}_{it} = \begin{cases} \{\mathbf{W}_{it} \cdot \mathbf{x}_t \geq 0\} & \text{if } i \text{ is root} \\ \{\mathbf{W}_{it} \cdot \mathbf{x}_t \geq 0\} & \text{if } i \text{ is not root \& } \hat{y}_{jt} = 1 \\ 0 & \text{otherwise} \end{cases}$$



- **Update:** On feedback V_{1t}, \dots, V_{Nt} , \mathbf{x}_t is “passed” to node i (“ $\mathbf{x}_t \rightarrow i$ ”) iff $V_{jt} = 1$

Top-down and (once again) independent of c_i !

On-line hierarchical classification/4: Algorithm/2

\mathbf{W}_{it} is regularized least squares-like [HK,AW,CBCG,RYT, ...]:

$$\mathbf{W}_{it} = \underbrace{\left(\underbrace{I}_{\text{Identity}} + \sum_{s:\mathbf{x}_s \rightarrow i} \mathbf{x}_s \mathbf{x}_s^\top + \mathbf{x}_t \mathbf{x}_t^\top \right)^{-1}}_{\text{2nd-order structure}} \underbrace{\left(\sum_{s:\mathbf{x}_s \rightarrow i} V_{is} \mathbf{x}_s \right)}_{\text{(1st-order) Perc.-like}}$$

Variance control

- Naturally arises from parametric model:

$\mathbf{W}_{it} \cdot \mathbf{x}_t$ (almost) conditionally unbiased estimator of $\mathbf{u} \cdot \mathbf{x}_t$
 accuracy depends on $\#\{s : \mathbf{x}_s \rightarrow i\}$

- Dual form (kernels): running time quadratic in $\#\{s \leq t - 1 : \mathbf{x}_s \rightarrow i\}$

On-line hierarchical classification/5

Regret bound (pointwise)

Example sequence $(\mathbf{x}_1, \mathbf{V}_1) \dots (\mathbf{x}_T, \mathbf{V}_T) \in \mathbb{R}^d \times \{0, 1\}$

$$\sum_{\text{examples } t} \left(\mathbb{E}[\ell_H(\hat{\mathbf{Y}}_t, \mathbf{V}_t)] - \mathbb{E}[\ell_H(\mathbf{Y}_t, \mathbf{V}_t)] \right)$$

$$\leq \frac{1}{\Delta^2} \sum_{\text{nodes } i} C_i \mathbb{E} \left[\sum_{\text{dim. } j} \log(1 + \lambda_{ij}) \right]$$

$$\min_{i,t} (\mathbf{u}_i \cdot \mathbf{x}_t)^2$$

(margin)

$$\sum_{k \in \text{subtree}(i)} c_k$$

(“value” of subtree)

j -th eigenval of

$$\sum_{s: \mathbf{x}_s \rightarrow i} \mathbf{x}_s \mathbf{x}_s^\top$$

Conclusions

- Framework for hierarchical classification
- Improving baseline top-down **evaluation** (label assignment) scheme (H-SVM) by bottom-up Bayes-optimal-like
- Modular approach
- Preliminary experiments
- On-line (top-down) regret analysis of regularized least-squares alg on linear parametric model

Open questions

- Experimental:
 - Clear advantage of **bottom-up** on synthetic, less clear on real-world:
"human noise" ? Naive independence assumption ?
 - Replace SVM by better algs for probability estimation (e.g. regularized logistic regression)
- Theoretical:
 - Regret analysis w.r.t. Bayes optimal logistic model ?
 - Remove independence assumption