

Easy Learning Theory for Highly Scalable Algorithms

Claudio Gentile

DICOM

Universita' dell'Insubria, Italy

claudio.gentile@uninsubria.it

August, 2005

Content of this tutorial/1

Part 1

Intro to on-line learning problems, methods, relative loss bounds:

- On-line learning setting, examples
- Learning with expert advice (Bayes voting), relative loss bounds
- On-line learning linear-threshold functions
- Extensions: tracking the target, label-efficiency
- ~~Learning regression functions~~

focus on
BINARY
classification



Content of this tutorial/2

Part 2

Intro to statistical Pattern Recognition problem, martingales, on-line to batch conversions:

- The statistical Pattern Recognition problem
- Types of error bounds (data/algorithm - independent/dependent)
- Reduction on-line pointwise \rightarrow i.i.d.
 - ~~Expectation analysis~~
 - Data-dependent analysis
- Final comments and conclusions

(Worst-case) on-Line Learning

[L88,A88]

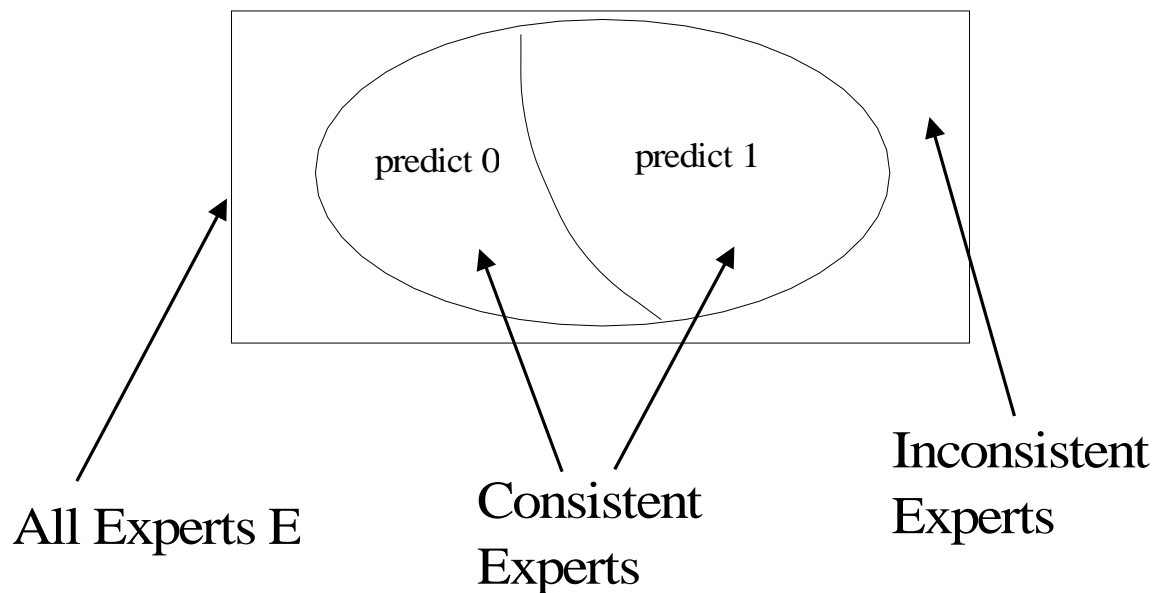
	E_1	E_2	E_3	...	E_n	experts		
						pred.	true lab.	loss
Day 1	1	1	-1	...	-1	-1	1	1
Day 2	1	-1	1	...	-1	1	-1	1
Day 3	-1	1	1	...	1	1	1	0
Day t	$z_{t,1}$	$z_{t,2}$	$z_{t,3}$...	$z_{t,n}$	\hat{y}_t	y_t	$\frac{1}{2} y_t - \hat{y}_t $

On-line protocol

For $t = 1, \dots, T$ do:	Get vector	$\mathbf{z}_t \in \{-1, 1\}^n$
	Predict	$\hat{y}_t \in \{-1, 1\}$
	Get label	$y_t \in \{-1, 1\}$
	Incur loss	$\frac{1}{2} y_t - \hat{y}_t \in \{0, 1\}$

Halving Algorithm

[BF72]



- Predicts with majority
- If mistake is made then number of consistent Experts is (at least) halved

A run of the Halving Algorithm (HA)

E_1	E_2	E_3	E_4	E_5	E_6	E_7	E_8	<i>majo rity</i>	<i>true label</i>	loss
1	1	-1	-1	1	1	-1	-1	1	-1	1
		-1	1			1	1	1	1	0
			1			-1	-1	-1	1	1
			↑							
			consistent							

\forall sequence with k consistent experts (out of n)

HA makes $m \leq \log_2(n/k)$ mistakes: $n/2^m \geq k$

Learning with expert advice/1

What if no expert E_i is consistent?

Sequence of examples $S = (z_1, y_1), \dots, (z_T, y_T)$

- $L_A(S)$ be total loss of alg. A on sequence S
- $L_i(S)$ be total loss of i -th expert E_i on S

Want bounds of the form:

$$\forall S : L_A(S) \leq a \min_i L_i(S) + b \log(n)$$

where a, b are constants

Bounds loss of algorithm **relative to** loss of best expert

Learning with expert advice/2

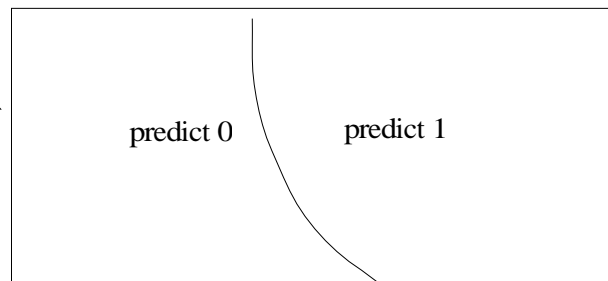
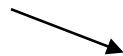
Can't wipe out experts!

Keep one weight per expert

The Weighted Majority Algorithm

[L89a,LW94]

All Experts E
vote with
their weight



- Predicts with larger side
- Weights of wrong experts are slashed by a factor $\beta \in [0, 1)$

Learning with expert advice/3

Number of mistakes of the WM algorithm/1

$L_{i,t}$ = # of mistakes of E_i before time step t

$w_{t,i}$ = $\beta^{L_{i,t}}$ weight of E_i at beginning of time step t

W_t = $\sum_{i=1}^n w_{i,t}$ total weight at time step t

Minority $\leq \frac{1}{2} W_t$

Majority $\geq \frac{1}{2} W_t$

If no mistake then minority multiplied by β : $W_{t+1} \leq 1 W_t$

If mistake then majority multiplied by β :

$$W_{t+1} \leq 1 \text{ Minority} + \beta \text{ Majority} \leq \frac{1 + \beta}{2} W_t$$

Learning with expert advice/3

Number of mistakes of the WM algorithm/2

Hence: W_{T+1} total final weight $\leq \left(\frac{1+\beta}{2}\right)^{L_{WMA}(S)} W_1$

$$W_{T+1} = \sum_{j=1}^n w_{T+1,j} = \sum_{j=1}^n \beta^{L_j(S)} \geq \beta^{L_i(S)}$$

We got: $\left(\frac{1+\beta}{2}\right)^{L_{WMA}(S)} \underbrace{W_1}_n \geq \beta^{L_i(S)}$

Solving for $L_{WMA}(S)$: $L_{WMA}(S) \leq \frac{\ln 1/\beta}{\ln \frac{2}{1+\beta}} L_i(S) + \frac{1}{\ln \frac{2}{1+\beta}} \ln n$

$$L_{WMA}(S) \leq \underbrace{2.63}_a \min_i L_i(S) + \underbrace{2.63}_b \ln n$$

$$\beta = 1/e$$

Learning with expert advice/4

Slightly more general: absolute loss/1

Alg. A keeps exponential weights

and predicts with the **weighted average**

[KW99]

$$w_{t,i} = \beta^{L_{i,t}}, \quad L_{i,t} = \sum_{j < t} \frac{1}{2} |z_{i,j} - y_j|$$

$$v_{t,i} = w_{t,i} / \sum_{i=1}^n w_{t,i} \quad \text{normalized weights}$$

$$\hat{y}_t = \mathbf{v}_t^\top \mathbf{z}_t,$$

where $z_{t,i} \in [-1, +1]$

$$\text{Absolute loss of } A: L_A(S) = \sum_{t=1}^T \frac{1}{2} |\hat{y}_t - y_t|$$

Learning with expert advice/4

Slightly more general: absolute loss/2

Relative loss bound

\forall sequences $S = (\mathbf{z}_1, \mathbf{y}_1), \dots, (\mathbf{z}_T, \mathbf{y}_T)$, $\mathbf{z}_t \in [-1, 1]^n$,

$\mathbf{y}_t \in [-1, 1]$

$$L_A(S) \leq \min_i \underbrace{(1 + \eta)}_a L_i(S) + \underbrace{\frac{1 + \eta}{\eta}}_b \ln(n), \quad \eta = \log 1/\beta$$

Notice $a > 1$ (constant η)

Regret bounds ($a = 1$) need time-changing β

[ACBG02]

Learning with expert advice/5

- Weighted Majority is just a Bayes voting scheme
- Easy to combine good experts (algorithms) so that prediction alg. is almost as good as best expert
- Bounds are logarithmic in # of experts

So far:

Learning relative to best expert/component

From now on:

Learning relative to best (thresholded) linear combination of experts/components

A more general setting

Instance	Prediction of alg A	Label	Loss of alg A
\mathbf{x}_1	\hat{y}_1	y_1	$L(y_1, \hat{y}_1)$
\vdots	\vdots	\vdots	\vdots
\mathbf{x}_t	\hat{y}_t	y_t	$L(y_t, \hat{y}_t)$
\vdots	\vdots	\vdots	\vdots
\mathbf{x}_T	\hat{y}_T	y_T	$L(y_T, \hat{y}_T)$
	Total Loss		$L_A(S)$

Sequence of examples $S = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T) \in \mathbb{R}^n \times \{-1, 1\}$

Comparison class $\{u\}$

Relative loss $L_A(S) - \inf_{\{u\}} \text{Loss}_u(S)$

Goal: Bound relative loss for arbitrary sequence S

Learning linear-threshold functions/1

Another run of the Halving Algorithm/1

Sequence of examples $S = (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_T, y_T) \in \mathbb{R}^2 \times \{-1, 1\}$

S is lin. separated by $\mathbf{u} \in \mathbb{R}^2 : \|\mathbf{u}\|_2 = 1$ with margin

$$0 < \gamma \leq y_t \mathbf{u}^\top \mathbf{x}_t \quad \forall t$$

$$R = \max_t \|\mathbf{x}_t\|_2$$

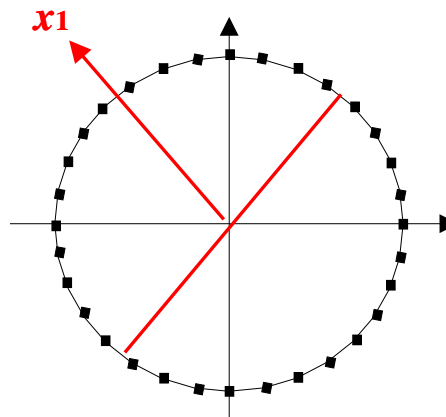
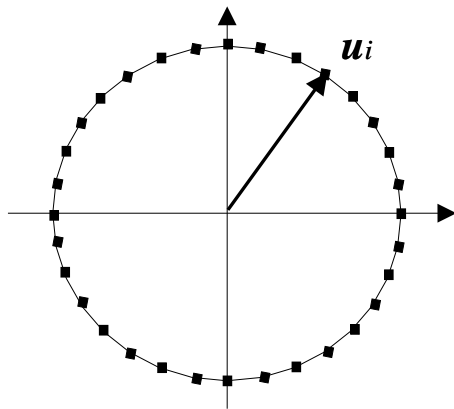
→ $\inf_{\{u\}} \text{Loss}_u(S) = 0$

Experts:

n (large) linear-threshold functions
evenly spread over unit circle

Expert i predicts $z_{it} = \text{sgn}(\mathbf{u}_i^\top \mathbf{x}_t)$

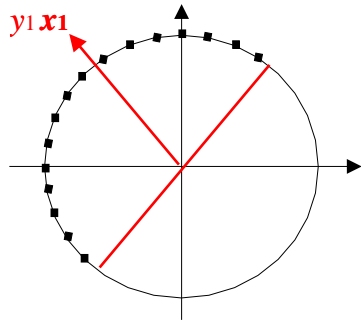
Feed experts with \mathbf{x}_1
and get expert prediction
vector \mathbf{z}_1



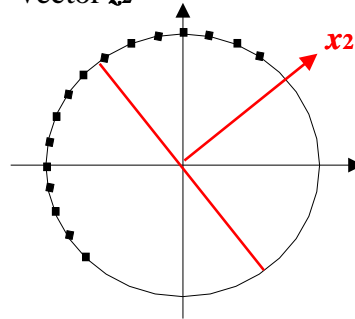
Learning linear-threshold functions/1

Another run of the Halving Algorithm/2

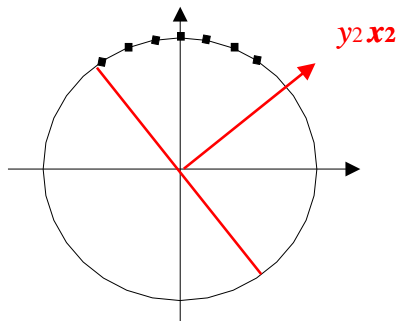
Get true label $y_1 = 1$ (mistake)
version space gets halved



Feed experts with x_2
and get expert prediction
vector z_2



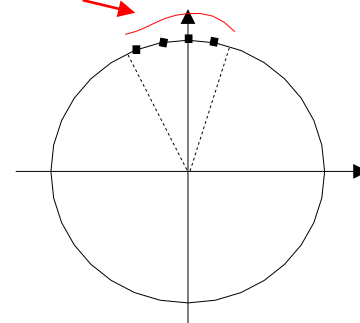
Get true label $y_2 = 1$ (mistake)
version space gets (at least) halved



...at the end

consistent experts ("margin")

...



$$m_{HA} \leq \log_2(n/k) = O(\log(R/\gamma)) \text{ for large } n$$

Learning linear-threshold functions/1

Another run of the Halving Algorithm/3

[HG02,GBNT04,...]

For d -dim vectors:

$$m_{HA} \leq \log_2 1/\mu(\text{consistent}(S)) \\ = O(d \log(R/\gamma)), \quad R = \max_t \|\mathbf{x}_t\|_2$$

Proof: $y_t \mathbf{u}^\top \mathbf{x}_t \geq \gamma$ and $\|\mathbf{u} - \mathbf{u}'\|_2 < \gamma/R$

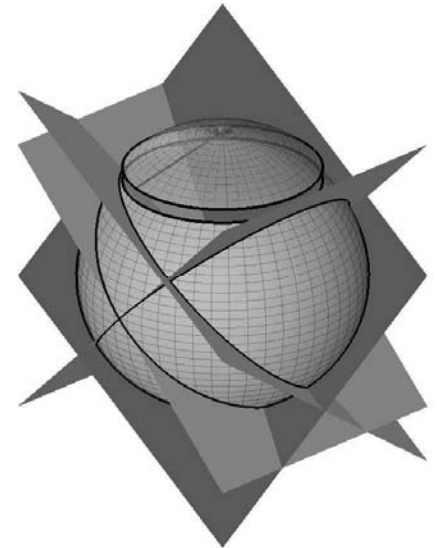
$\implies y_t (\mathbf{u}')^\top \mathbf{x}_t > 0$

$\implies \exists$ ball B of radius $\gamma/2R$: $B \subseteq \text{consistent}(S)$,

$\mu(B) = (\gamma/2R)^{d-1} \mu(\text{surface of } d\text{-dim unit sphere})$

Drawbacks: **Linear** dependence on dimension d

Looks too time-consuming (linear dependence on d ?)



Courtesy: R. Herbrich

Learning linear-threshold functions/2

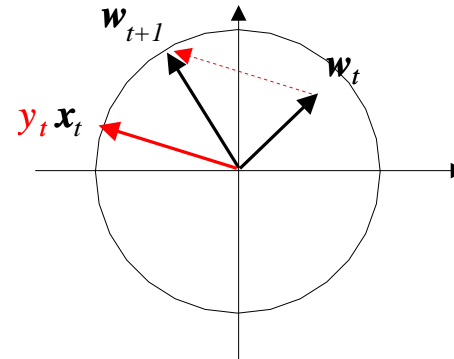
The (first-order) Perceptron algorithm

[Ro62, ...]

Keep weight vector $\mathbf{w}_t \in \mathbb{R}^n$

In trial t :

- Get instance $\mathbf{x}_t \in \mathbb{R}^n$
- Predict with $\hat{y}_t = \text{SGN}(\mathbf{w}_t^\top \mathbf{x}_t) \in \{-1, 1\}$
- Get label $y_t \in \{-1, 1\}$
- If **mistake** ($y_t \mathbf{w}_t^\top \mathbf{x}_t \leq 0$) then update $\mathbf{w}_{t+1} := \mathbf{w}_t + y_t \mathbf{x}_t$



Learning linear-threshold functions/3

Perceptron convergence theorem/1

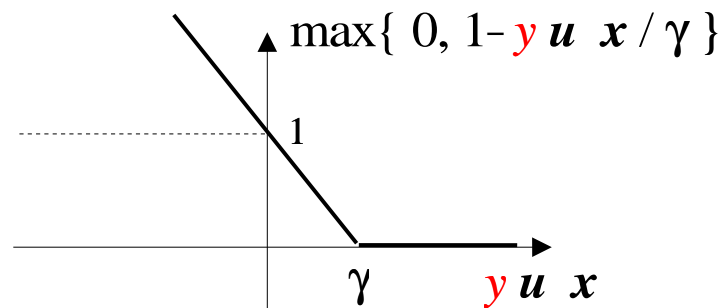
[Bl62, No62,...]

Arbitrary sequence $S = (\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_T, \mathbf{y}_T) \in \mathbb{R}^n \times \{-1, 1\}$

$$\# \text{ of mistakes} \leq \inf_{\gamma > 0, \|\mathbf{u}\|_2 = 1} \left(\underbrace{D_\gamma(\mathbf{u}; S)}_{\text{"loss" of } \mathbf{u}} + \frac{\sqrt{\sum_{t \in \mathcal{M}} \|\mathbf{x}_t\|_2^2}}{\gamma} \right),$$

\mathcal{M} is set of mistaken trials t ,

$$D_\gamma(\mathbf{u}; S) = \sum_{t \in \mathcal{M}} \max\{0, 1 - \mathbf{y}_t \mathbf{u}^\top \mathbf{x}_t / \gamma\}$$



Learning linear-threshold functions/3

Perceptron convergence theorem/2

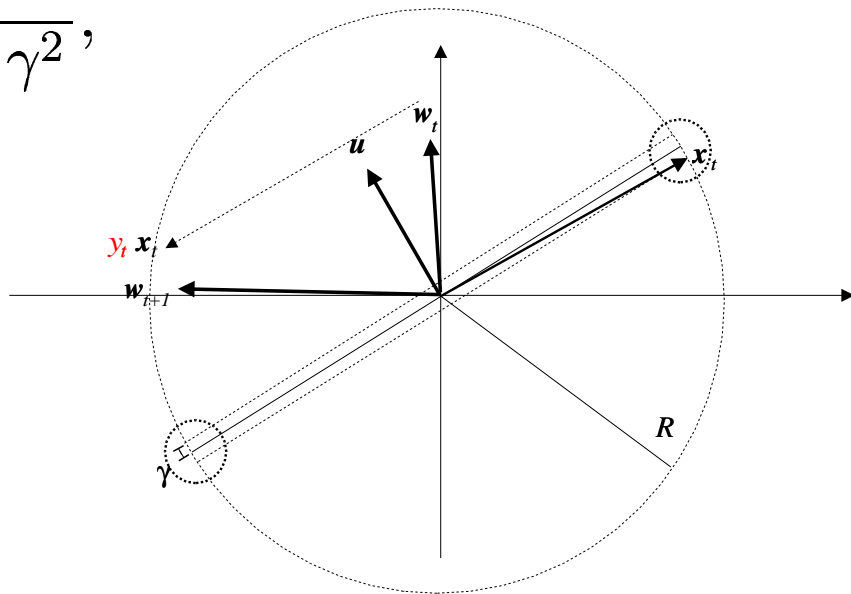
When S is separated by $\mathbf{u} : \|\mathbf{u}\|_2 = 1$ with margin

$$\gamma \leq \mathbf{y}_t \mathbf{u}^\top \mathbf{x}_t \quad \forall t$$

gets

$$\# \text{ of mistakes} \leq \frac{R^2}{\gamma^2},$$

$$\|\mathbf{x}_t\| \leq R$$



Pointwise bound:

Depends on radius R and margin γ

Learning linear-threshold functions/4

The second-order Perceptron algorithm

[CBCG05]

Keep weight vector $\mathbf{w}_t \in \mathbb{R}^n$ and matrix S_t

In trial t :

- Get instance $\mathbf{x}_t \in \mathbb{R}^n$
- Predict with $\hat{y}_t = \text{SGN}(\mathbf{w}_t^\top (aI + S_t)^{-1} \mathbf{x}_t) \in \{-1, 1\}$
positive parameter
- Get label $y_t \in \{-1, 1\}$
- **If mistake then** update
 - $\mathbf{w}_{t+1} := \mathbf{w}_t + y_t \hat{\mathbf{x}}_t$
 - $S_{t+1} = S_t + \hat{\mathbf{x}}_t \hat{\mathbf{x}}_t^\top$, $\hat{\mathbf{x}}_t = \mathbf{x}_t / \|\mathbf{x}_t\|$

Turns to first-order when $a \rightarrow \infty$

Learning linear-threshold functions/5

Second-order convergence theorem

[Ge04]

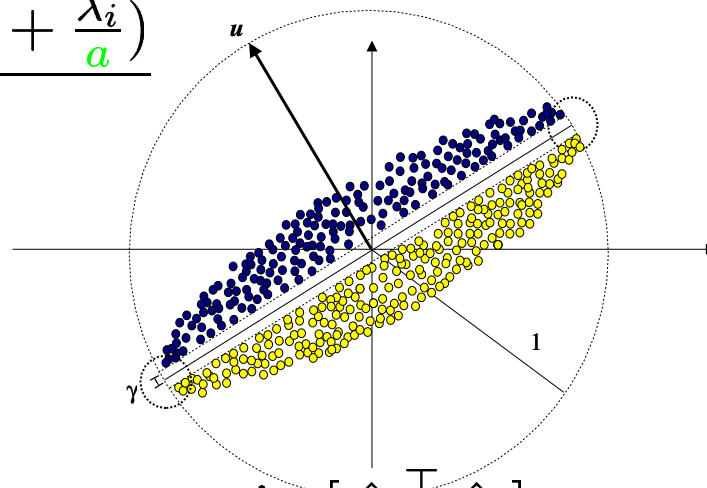
When $S = (\hat{\mathbf{x}}_1, y_1), \dots, (\hat{\mathbf{x}}_T, y_T) \in \mathbb{R}^n \times \{-1, 1\}$
is separated by \mathbf{u} with margin $\gamma \leq y_t \mathbf{u}^\top \hat{\mathbf{x}}_t$, $\|\hat{\mathbf{x}}_t\| \leq 1 \forall t$
gets

$$\# \text{ of mistakes} \leq \frac{a + \sum_{i=1}^n \ln(1 + \frac{\lambda_i}{a})}{\gamma}$$

More complicated bound
in the nonseparable case

Pointwise bound:

Depends on **eigenstructure** $\{\lambda_i\}$ of Gram matrix $[\hat{\mathbf{x}}_j^\top \hat{\mathbf{x}}_k]_{j,k \in \mathcal{M}}$
and **linearly** on inverse margin γ



Learning linear-threshold functions/6

Kernel Perceptron

[FS98,...]

Keep pool of "support vectors" \mathcal{M}_t

In trial t :

- Get instance $\mathbf{x}_t \in \mathbb{R}^n$
- Predict with $\hat{y}_t = \text{SGN}(\sum_{i \in \mathcal{M}_t} y_i K(\mathbf{x}_i, \mathbf{x}_t)) \in \{-1, 1\}$
- Get label $y_t \in \{-1, 1\}$
- **If mistake then** update $\mathcal{M}_{t+1} := \mathcal{M}_t \cup \{t\}$

Learning linear-threshold functions/7

Kernel Perceptron convergence theorem/1

Arbitrary sequence $S = (\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_T, \mathbf{y}_T) \in \mathbb{R}^n \times \{-1, 1\}$

$$\# \text{ of mist.} \leq \inf_{\gamma > 0, f \in H_K, \|f\|=1} \left(\underbrace{D_\gamma(f; S)}_{\text{"loss" of } f} + \frac{\sqrt{\sum_{t \in \mathcal{M}} K(\mathbf{x}_t, \mathbf{x}_t)}}{\gamma} \right)$$

$$H_K = \{f(\cdot) = \sum_{t=1}^T \alpha_t K(\mathbf{x}_t, \cdot) : \alpha_t \in \mathbb{R}\},$$

\mathcal{M} is set of mistaken trials t ,

$$D_\gamma(f; S) = \sum_{t \in \mathcal{M}} \max\{0, 1 - \mathbf{y}_t f(\mathbf{x}_t) / \gamma\}$$

Separable case:

$$\# \text{ of mistakes} \leq R^2 / \gamma^2, \quad K(\mathbf{x}_t, \mathbf{x}_t) \leq R^2$$

Learning linear-threshold functions/8

Kernel Second-order Perceptron

[CBCG05]

Keep pool of "support vectors" \mathcal{M}_t

In trial t :

- Get instance $\mathbf{x}_t \in \mathbb{R}^n$

$$\begin{bmatrix} \vdots \\ y_i \\ \vdots \end{bmatrix} \quad i \in M_t$$

- Predict with $\hat{y}_t =$

$$\text{SGN} \left(\mathbf{y}_t^\top \left(a I + \underbrace{[\hat{K}(\mathbf{x}_i, \mathbf{x}_j)]_{i,j \in \mathcal{M}_t}}_{\text{current Gram matrix}} \right)^{-1} \mathbf{v}_t \right) \in \{-1, 1\},$$

$$\begin{bmatrix} \vdots \\ \hat{K}(\mathbf{x}_i, \mathbf{x}_i) \\ \vdots \end{bmatrix} \quad i \in M$$

- Get label $y_t \in \{-1, 1\}$
- If **mistake** then update $\mathcal{M}_{t+1} := M_t \cup \{t\}$

Learning linear-threshold functions/9

Kernel Second-order convergence theorem

When $S = (\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_T, \mathbf{y}_T) \in \mathbb{R}^n \times \{-1, 1\}$
is separated by $f(\cdot) = \sum_{t=1}^T \alpha_t \hat{K}(\mathbf{x}_t, \cdot)$, $\alpha_t \in \mathbb{R}$,
with margin $\gamma \leq \mathbf{y}_t f(\mathbf{x}_t) \forall t$

gets

$$\# \text{ of mist.} \leq \frac{a + \sum_i \ln(1 + \frac{\lambda_i}{a})}{\gamma},$$

λ_i is i -th eigenvalue of (normalized) kernel Gram matrix
 $[\hat{K}(\mathbf{x}_i, \mathbf{x}_j)]_{i,j \in \mathcal{M}}$,

\mathcal{M} is set of mistaken trials

Learning linear-threshold functions/10

Second-order Perceptron: computational aspects

Time per trial

Primal formulation

compute $(aI + S_t)^{-1}$ based on $(aI + S_{t-1})^{-1}$

$O(n^2)$ extra time per trial

Dual formulation

compute $\left(aI + [\hat{G}(\mathbf{x}_i, \mathbf{x}_j)]_{i,j \in \mathcal{M}_t} \right)^{-1}$

based on $\left(aI + [\hat{G}(\mathbf{x}_i, \mathbf{x}_j)]_{i,j \in \mathcal{M}_{t-1}} \right)^{-1}$

$O(|\mathcal{M}_t|^2)$ extra inner products (kernel evaluations) per trial

of mistakes so far

Learning linear-threshold functions/11

Additive algorithms

An **additive** algorithm (e.g. first/second-order Perceptron):

- Relies on linear algebra
- Is rotation invariant (depends on data via angles)
- Can be **easily** kernelized ($\mathbf{x}_i^\top \mathbf{x}_j \rightarrow K(\mathbf{x}_i, \mathbf{x}_j)$)
- Has **no** bias for axes-parallel directions (no feature selection)

Learning linear-threshold functions/12

Nonadditive algorithms

- **No** linear algebra
- **No** rotation invariance
- **Harder** to kernelize
- Bias for sparse solutions (built-in feature selection)

Example: p -norm algorithms

Learning linear-threshold functions/13

p -norm algs

[GLS01, GL99, Ge03]

Keep weight vector $\mathbf{w}_t \in \mathbb{R}^n$

In trial t :

$$\mathbf{f}(\cdot) = \nabla \frac{1}{2} \|\cdot\|_p^2, p \geq 2$$

- Get instance $\mathbf{x}_t \in \mathbb{R}^n$
- Predict $\hat{y}_t = \text{SGN}(\mathbf{f}(\mathbf{w}_t)^\top \mathbf{x}_t) \in \{-1, 1\}$
- Get label $y_t \in \{-1, 1\}$
- If **mistake** then update $\mathbf{w}_{t+1} := \mathbf{w}_t + y_t \mathbf{x}_t$

Notice:

- $p = 2$ gets (first-order) Perceptron
- $p = O(\ln n)$ gets Weighted Majority/Winnow [L88, LW94]
- $2 < p < O(\ln n)$ interpolates between the two extremes

Learning linear-threshold functions/14

p -norm Perceptron convergence theorem/1

[GLS01, GL99, Ge03]

Arbitrary sequence $S = (\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_T, \mathbf{y}_T) \in \mathbb{R}^n \times \{-1, 1\}$

$$\# \text{ mistakes} \leq \inf_{\gamma > 0, \|\mathbf{u}\|_q = 1} \left(\underbrace{D_\gamma(\mathbf{u}; S)}_{\text{"loss" of } \mathbf{u}} + \frac{\sqrt{(p-1) \sum_{t \in \mathcal{M}} \|\mathbf{x}_t\|_p^2}}{\gamma} \right)$$

\mathcal{M} is set of mistaken trials t ,

$$D_\gamma(\mathbf{u}; S) = \sum_{t \in \mathcal{M}} \max\{0, 1 - \mathbf{y}_t \mathbf{u}^\top \mathbf{x}_t / \gamma\}$$

Learning linear-threshold functions/14

p -norm Perceptron convergence theorem/2

When S is separated by $\mathbf{u} : \|\mathbf{u}\|_q = 1$ with margin

$$\gamma \leq \mathbf{y}_t \mathbf{u}^\top \mathbf{x}_t \quad \forall t \quad \underbrace{(1/p + 1/q = 1)}_{\text{dual norms}}$$

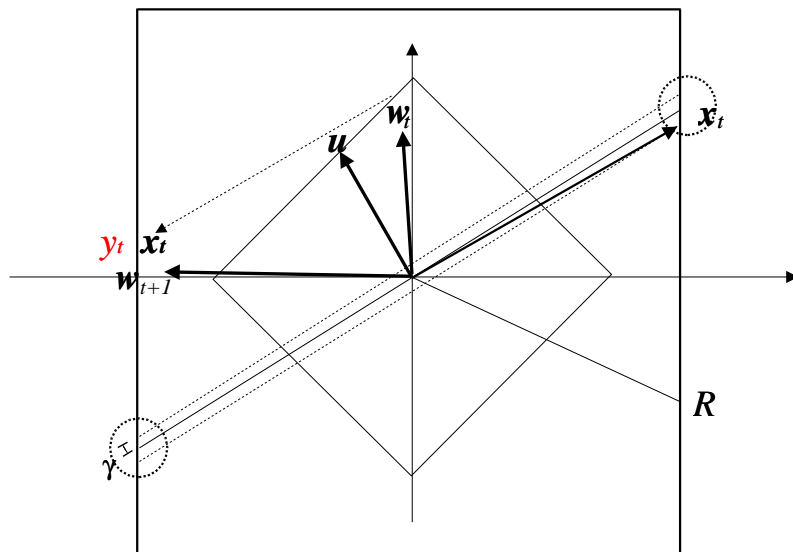
gets

$$\# \text{ of mistakes} \leq (p - 1) \frac{R^2}{\gamma^2}$$

$$\|\mathbf{x}_t\| \leq R$$

Pointwise bound:

Depends on p -norm radius R
and (q -norm) margin γ



Learning linear-threshold functions/15

p -norm algorithms with kernels/1 (wild slide) [Ge04]

$$\begin{array}{r} \mathbf{x} = \begin{array}{c} x_1 \\ x_2 \\ \vdots \\ x_n \end{array} \Rightarrow \Phi(\mathbf{x}) = \begin{array}{c} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_n \\ x_1 x_2 \\ \vdots \\ x_1 x_2 \dots x_n \end{array} \end{array} \quad \begin{array}{l} K(\mathbf{x}, \mathbf{y}) = \Phi(\mathbf{x})^\top \Phi(\mathbf{y}) \\ \\ = \prod_{i=1}^n (1 + x_i y_i) \\ \\ \text{(Simple poly kernel)} \end{array}$$

Learning linear-threshold functions/15

p -norm algorithms with kernels/2 (wild slide)

p -norm hypothesis: $\mathbf{w} = \sum_{i \in \mathcal{M}} y_i \Phi(\mathbf{x}_i)$

$$p\text{-norm margin: } = \mathbf{f}(\mathbf{w})^\top \Phi(\mathbf{x}) \qquad \mathbf{f}(\mathbf{w}) = \mathbf{w}^{p-1}$$

$$= \underbrace{\left(\sum_{i \in \mathcal{M}} y_i \Phi(\mathbf{x}_i)^\top \right)^{p-1}}_{\text{expand!}} \Phi(\mathbf{x})$$

Then expand polynomial and use $\Phi(\mathbf{x})\Phi(\mathbf{y}) = \Phi(\mathbf{x}\mathbf{y})$

Learning linear-threshold functions/15

p -norm algorithms with kernels/3 (wild slide)

Example: $p = 4$, $f(\mathbf{w}) = \mathbf{w}^3$

$$\mathbf{w} = y_1 \Phi(\mathbf{x}_1) + y_2 \Phi(\mathbf{x}_2)$$

follow pattern $(a + b)^3 = a^3 + 3a^2b + 3ab^2 + b^3$

$$f(\mathbf{w}) =$$

$$y_1^3 \Phi^3(\mathbf{x}_1) + 3y_1^2 y_2 \Phi^2(\mathbf{x}_1) \Phi(\mathbf{x}_2) + 3y_1 y_2^2 \Phi(\mathbf{x}_1) \Phi^2(\mathbf{x}_2) +$$

$$y_2^3 \Phi^3(\mathbf{x}_2) =$$

$$y_1 \Phi(\mathbf{x}_1^3) + 3y_2 \Phi(\mathbf{x}_1^2 \mathbf{x}_2) + 3y_1 \Phi(\mathbf{x}_1 \mathbf{x}_2^2) + y_2 \Phi(\mathbf{x}_2^3) =$$

$$y_1 \Phi(\mathbf{x}_1^3) + 3y_2 \Phi(\mathbf{x}_1^2 \mathbf{x}_2) + 3y_1 \Phi(\mathbf{x}_1 \mathbf{x}_2^2) + y_2 \Phi(\mathbf{x}_2^3)$$

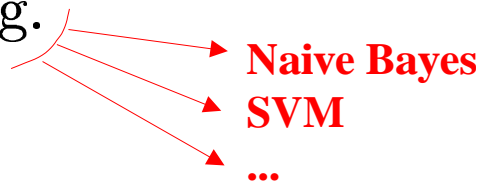
Then p -norm margin $f(\mathbf{w})^\top \Phi(\mathbf{x}) =$

$$y_1 \underbrace{K(\mathbf{x}_1^3, \mathbf{x})}_{SV} + 3y_2 \underbrace{K(\mathbf{x}_1^2 \mathbf{x}_2, \mathbf{x})}_{SV} + 3y_1 \underbrace{K(\mathbf{x}_1 \mathbf{x}_2^2, \mathbf{x})}_{SV} + y_2 \underbrace{K(\mathbf{x}_2^3, \mathbf{x})}_{SV}$$

Batch algorithm run on-line

$$S = (\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_T, \mathbf{y}_T) \in \mathbb{R}^n \times \{-1, +1\}$$

A is generic batch classification alg.



In trial t :

- train A on prefix $S_{t-1} = (\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_{t-1}, \mathbf{y}_{t-1})$
- Get $h_t(\mathbf{x}) = A_{S_{t-1}}(\mathbf{x})$
- Mistake if $h_t(\mathbf{x}_t) \neq \mathbf{y}_t$

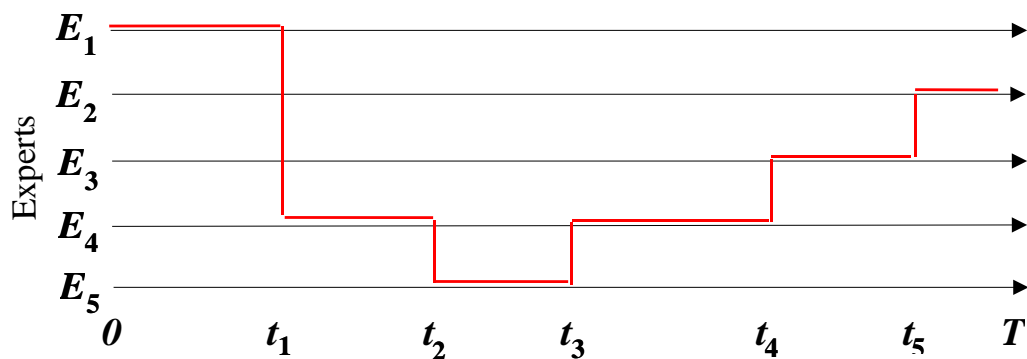
Can count # mistakes on sequence S

Tracking the target:

Motivation and relative loss

[LW94,HW98,AW98]

Previous methods yield good results when *single* target (expert/vector) is good on whole sequence S



Non-stationary tasks:
good targets
change with time

Given $S = (\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_T, \mathbf{y}_T)$, want to bound

$$L_A(S) - a \min_{\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{u}_T} \underbrace{\left(\text{Loss}_{\mathbf{u}_0, \mathbf{u}_1, \dots, \mathbf{u}_T}(S) + \sum_{t=1}^T \text{dist}(\mathbf{u}_t, \mathbf{u}_{t+1}) \right)}_{\text{loss of "shifting" target}}$$

Tracking the target: Examples

- Experts framework:

\mathbf{u}_t are the n unit vectors, $\mathbf{x}_t \in [-1, +1]^n$

$$Loss_{\mathbf{u}_0, \mathbf{u}_1, \dots, \mathbf{u}_T}(S) = \sum_{t=1}^T \frac{1}{2} |\mathbf{u}_t^\top \mathbf{x}_t - y_t| \quad (\text{mistakes})$$

$$dist(\mathbf{u}_t, \mathbf{u}_{t+1}) = \|\mathbf{u}_{t+1} - \mathbf{u}_t\| = \begin{cases} 1 & \text{if } \mathbf{u}_{t+1} \neq \mathbf{u}_t \\ 0 & \text{otherw.} \end{cases} \quad (\text{shifts})$$

- Linear-threshold functions:

$\mathbf{u}_t, \mathbf{x}_t \in \mathbb{R}^n$

$$Loss_{\mathbf{u}_0, \mathbf{u}_1, \dots, \mathbf{u}_T}(S) = \sum_{t=1}^T \max\{0, 1 - y_t \mathbf{u}_t^\top \mathbf{x}_t / \gamma\}$$

$$dist(\mathbf{u}_t, \mathbf{u}_{t+1}) = \|\mathbf{u}_{t+1} - \mathbf{u}_t\|$$

Tracking the target: Algorithms/1

- Experts framework:

Basic expert algs. (exponential update) prevent experts to recover when they start performing better

Modification to Weighted Majority (fixed share alg.):

- $w'_{t,i} = \beta^{L_{i,t+1}}$
- Each expert sends fraction $\frac{\alpha}{n-1}$, $\alpha > 0$,
of total weight to other $n - 1$ experts $\rightarrow \mathbf{w}_{t+1}$
- $v_{t+1,i} = w_{t+1,i} / \sum_{j=1}^n w_{t+1,j}$
- $\hat{\mathbf{y}}_{t+1} = \mathbf{v}_{t+1}^\top \mathbf{x}_{t+1}$, $x_{t+1,i} \in [-1, +1]$

Tracking the target: Algorithms/2

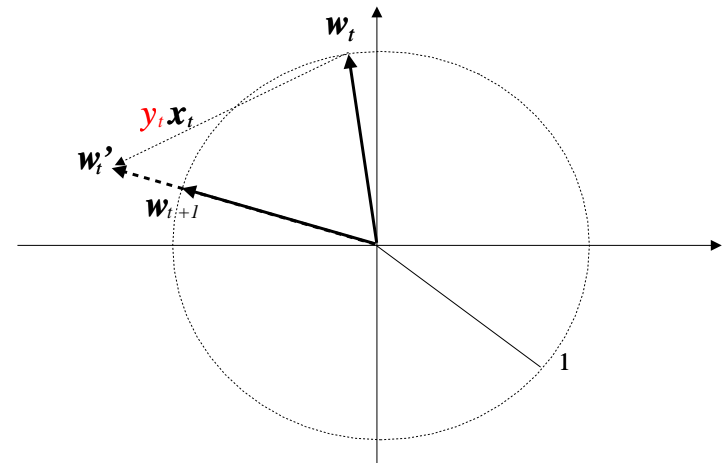
- Linear-threshold functions:

Same story ...

Modification to linear-threshold algs. (projection-based)

Simplest example:

- $\mathbf{w}'_t = \mathbf{w}_t + \eta y_t \mathbf{x}_t$
- *projects* \mathbf{w}'_t onto convex set (e.g. unit ball)
 $\rightarrow \mathbf{w}_{t+1}$
- $\hat{y}_{t+1} = \text{SGN}(\mathbf{w}_{t+1}^\top \mathbf{x}_{t+1})$



Tracking the target: Bounds

[HW98,AW98,...]

- Experts framework:

A is share update alg. with $\alpha = \frac{k}{T-1}$

Bits to encode
shift times

Bits to encode experts
active at each shift

$$L_A(S) \leq a \min_{P(k)} L_{P(k)}(S) + b \left(k \log(T/K) + k \log n \right)$$

- Linear-threshold functions:

Simple case: compare to sequence of linear separators

$$\mathbf{u}_1, \dots, \mathbf{u}_T : 0 < \gamma \leq \mathbf{y}_s \mathbf{u}_t^\top \mathbf{x}_s \quad s, t = 1, \dots, T$$

A is projection-based Perceptron with $\eta = \frac{\gamma}{2R^2}$, $\|\mathbf{x}_t\| \leq R$

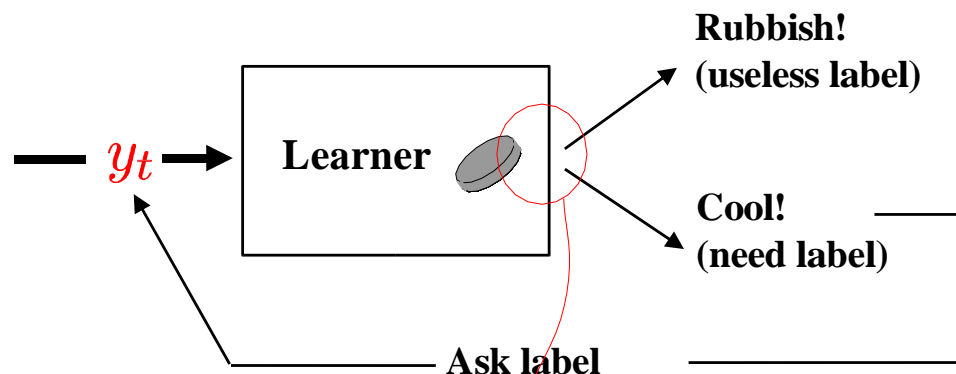
$$\# \text{ of mistakes} \leq \frac{4R^2}{\gamma^2} \left(\frac{1}{2} + \sum_{t=1}^T \|\mathbf{u}_{t+1} - \mathbf{u}_t\| \right)$$

Can be generalized in several ways

Label-efficient prediction

[ACL90,FSST97,CGZ04...]

Labels are scarce and/or expensive



margin-based
random decision

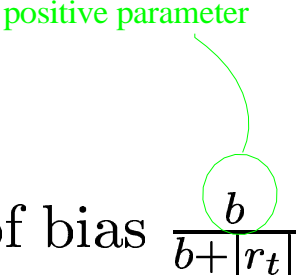
Learner decides in on-line fashion which labels are needed:
The **higher** the margin the **less** likely label y_t is requested

Label-efficient prediction: Algorithms

Simplest example: Label-efficient Perceptron

Keep weight vector $\mathbf{w}_t \in \mathbb{R}^n$ and margin $r_t \in \mathbb{R}$

In trial t :

- Get instance $\mathbf{x}_t \in \mathbb{R}^n$, set $r_t = \mathbf{w}_t^\top \mathbf{x}_t$
- Predict with $\hat{y}_t = \text{SGN}(r_t) \in \{-1, 1\}$
- Draw Bernoulli random variable $Z_t \in \{0, 1\}$ of bias $\frac{b}{b+|r_t|}$

- If $Z_t = 1$ then
 - Ask for label $y_t \in \{-1, 1\}$
 - if **mistake** ($y_t r_t \leq 0$) then update $\mathbf{w}_{t+1} = \mathbf{w}_t + y_t \mathbf{x}_t$

Else ($Z_t = 0$) $\mathbf{w}_{t+1} = \mathbf{w}_t$

Label-efficient prediction: Bounds

Simplest result:

When S is separated by $\mathbf{u} : \|\mathbf{u}\|_2 = 1$ with margin $\gamma \leq \mathbf{y}_t \mathbf{u}^\top \mathbf{x}_t \forall t$ and $\|\mathbf{x}_t\| \leq R$ gets

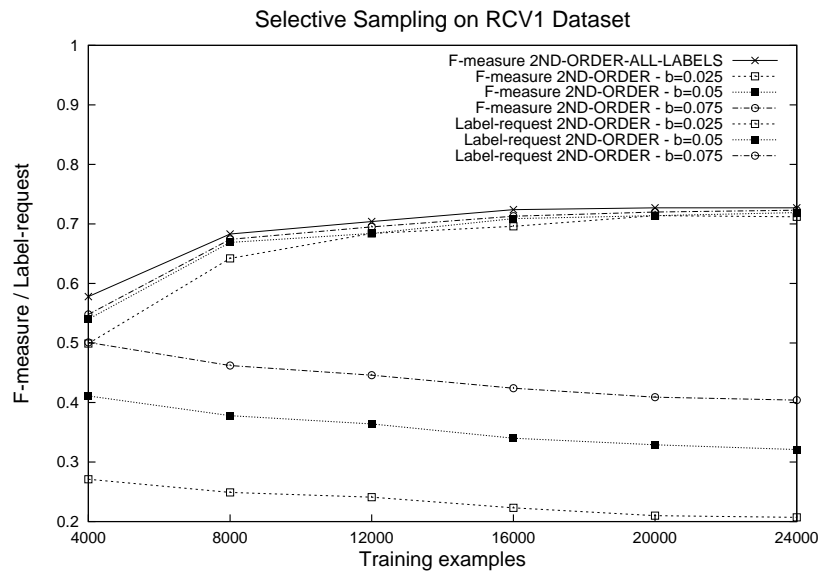
$$\text{Expected \# of labels} = \sum_{t=1}^T \mathbb{E} \left[\frac{b}{b + |r_t|} \right] \leq T,$$

and, for $b = R^2/2$,

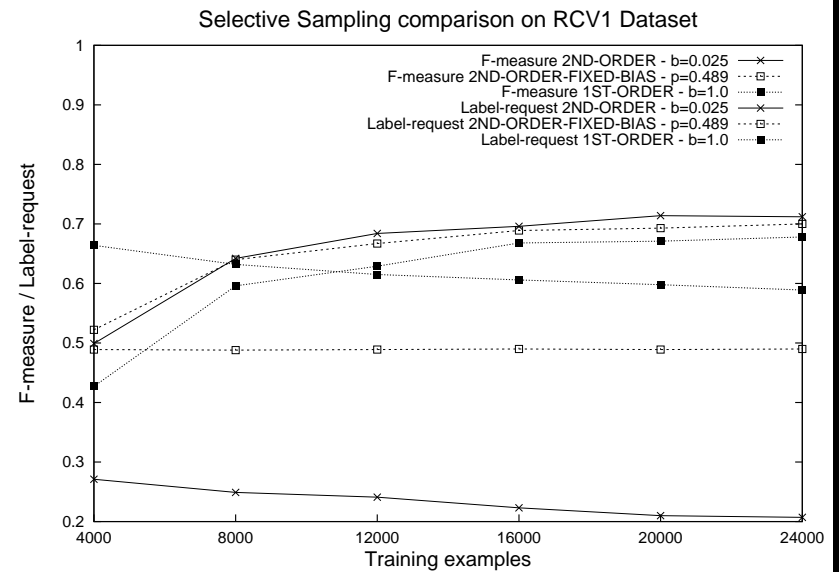
$$\text{Expected \# of mistakes} \leq \frac{R^2}{\gamma^2}$$

Can be generalized in several ways

Label-efficient prediction: Empirical evidence/1

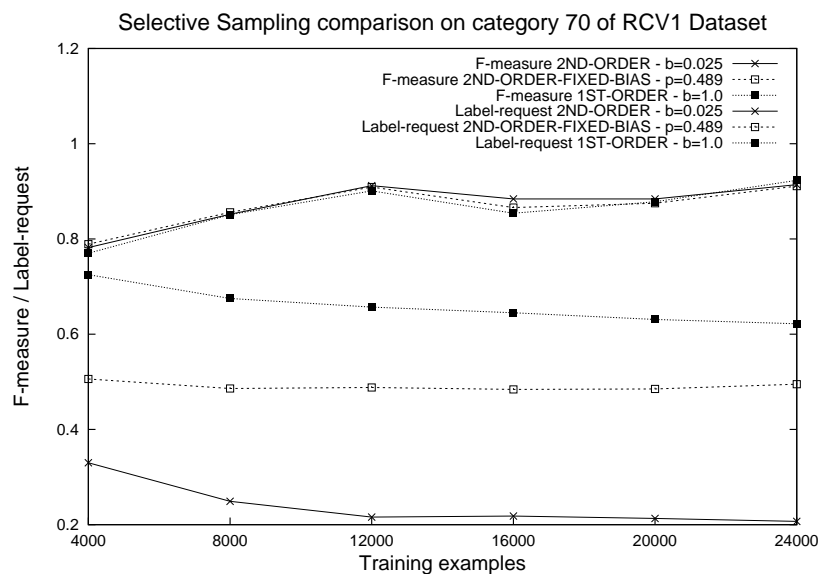


(a)

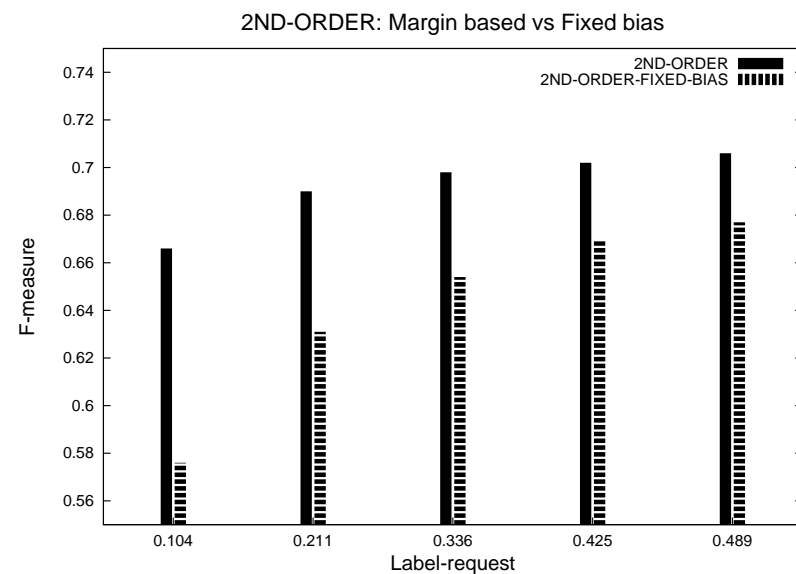


(b)

Label-efficient prediction: Empirical evidence/2



(c)



(d)

End of Part 1

Generalization bounds/1

Given

- class \mathcal{H} of ± 1 functions
- i.i.d. sequence $S = (X_1, Y_1), \dots, (X_T, Y_T)$ over $\mathbb{R}^n \times \{-1, 1\}$,

0-1 loss
in our case

want to compute hypothesis $\hat{H} = \hat{H}_S$ with small risk

$\text{risk}(\hat{H}) = \mathbb{E}_{X, Y} [\text{loss}(Y, \hat{H}(X))]$:

$$\mathbb{P} \left(\text{risk}(\hat{H}) \leq \inf_{h \in \mathcal{H}} \text{risk}(h) + \epsilon \right) \geq 1 - \delta$$

Generalization bounds/2: VC Uniform conv. [VC71]

Key quantity is **empirical** risk

$$\text{risk}_{\text{emp}}(h) = \frac{1}{T} \sum_{t=1}^T \text{loss}(\mathbf{Y}_t, h(\mathbf{X}_t))$$

VC-bound:

$$\mathbb{P} \left(\sup_{h \in \mathcal{H}} |\text{risk}_{\text{emp}}(h) - \text{risk}(h)| \geq c \sqrt{\frac{d + \ln 1/\delta}{T}} \right) \leq \delta$$

[VC71, Va98, ...]

constant (pointing to c) *VC-dim(H)* (pointing to d)

$\implies \hat{H} = \text{arginf}_{h \in \mathcal{H}} \text{risk}_{\text{emp}}(h)$ is s.t.

$$\mathbb{P} \left(\text{risk}(\hat{H}) \leq \inf_{h \in \mathcal{H}} \text{risk}(h) + 2c \sqrt{\frac{d + \ln 2/\delta}{T}} \right) \geq 1 - \delta$$

Generalization bounds/3:

Data-dep. uniform conv./1

[Ba98, BLM00, WSTSS99, BM02, ...]

$$\sqrt{\frac{d + \ln 2/\delta}{T}} \rightarrow C_T(S) + \sqrt{\frac{\ln 1/\delta}{T}}$$

$$C_T(S) = C_T(S, \mathcal{H})$$

is sample statistic:

- empirical VC-entropy [BLM00, WSTSS99]
- Rademacher complexity [BM02]
- Maximum discrepancy [BLM00]
- ...

Stronger than VC since $C_T(S) \approx \mathbb{E}[C_T(S)] \ll \sqrt{d/T}$

Generalization bounds/3:

Data-dep. uniform conv./2

Others (e.g., margin-based bounds for linear-threshold functions)

[AKLL02, KP02, LSM01, SFBL98, ...]

$$\mathbb{P} \left(\forall h \in \mathcal{H} : \text{risk}(h) \leq \text{risk}_{\text{emp}}(h) + C_T(h, S) + c \sqrt{\frac{\ln 1/\delta}{T}} \right) \geq 1 - \delta$$

Leave algorithmic problem of computing $h \in \mathcal{H}$ optimizing trade-off

$$\text{risk}_{\text{emp}}(h) \quad \text{vs} \quad C_T(h, S)$$

Digression: martingales/1

Coin-tossing game:

X_1, X_2, \dots, X_t are ± 1 i.i.d. variables with $\mathbb{P}(X_t = +1) = 1/2$

Gambler's strategy:

$L_t = L_t(X_1, X_2, \dots, X_t)$ is gain (or loss) at time t

L_1, L_2, \dots, L_t are **no longer** independent

Game is *fair* if

$$\mathbb{E}[L_{t+1} \mid X_1, \dots, X_t] = 0 \quad (\text{w.p.1}) \quad \forall t$$

- Sequence L_1, L_2, \dots , is *martingale difference sequence* (w.r.t. X_1, X_2, \dots)
- Partial sums $S_t = L_1 + L_2 + \dots + L_t$ is *martingale* (w.r.t. X_1, X_2, \dots)

Digression: martingales/2

Laws of large numbers

(empirical average concentrates around mean)

extend to martingales

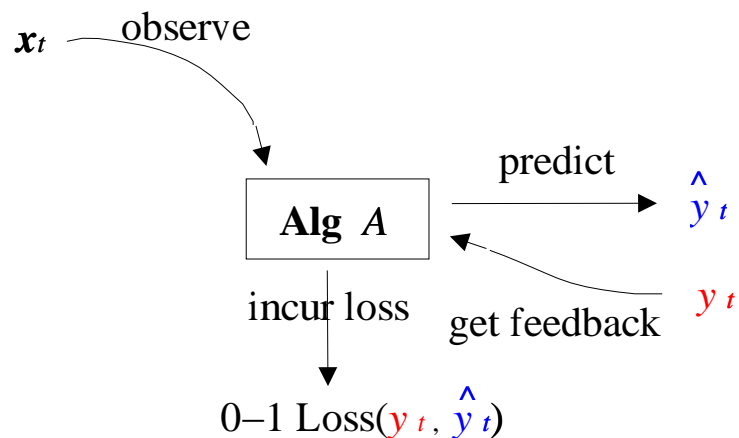
Independent variables	Dependent variables
(Zero-mean) independent r.v.	Martingale diff. sequence
Sum of (zero-mean) indep. r.v.	Martingale

(Hoeffding-Azuma) If $L_1 + L_2 + \dots + L_T$ is martingale difference sequence with bounded L_t

$$\frac{L_1 + L_2 + \dots + L_T}{T} \approx 0 \quad (\text{with high probability})$$

On-line pointwise bounds

$$S = (\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_T, \mathbf{y}_T)$$

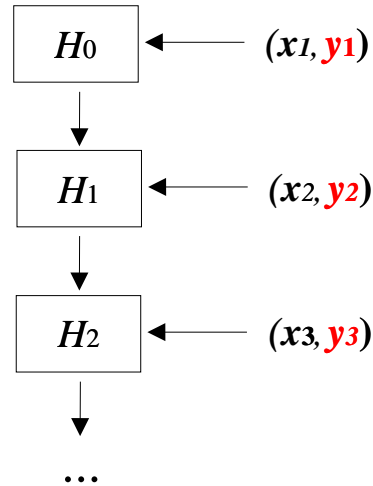


Pointwise bounds so far:

$$\text{Total \# mistakes}_A(S) \leq \text{some_function}(S)$$

n, R, γ (Halving) R, γ (1st Perc) λ_i, γ (2nd Perc) $R, \gamma(\text{dual})$ (p -norm) ...

On-line pointwise \rightarrow i.i.d. data-dependent/1



Sweep through sequence of examples S just **once!**

Get sequence of hypotheses

$$H_0, H_1, H_2, \dots, H_T : H_t = H_t((x_1, y_1), \dots, (x_t, y_t))$$

Goal: Extract one with small risk

Early ref: [L] (separate test set)

On-line pointwise \rightarrow i.i.d. data-dependent/2

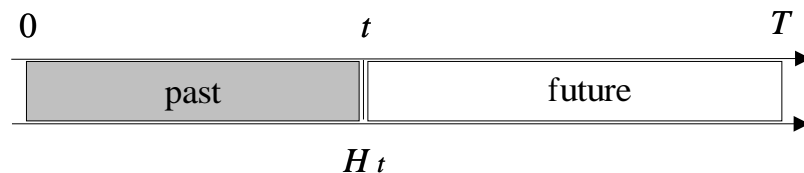
Which one?

1. **Last** one: H_T (back to uniform convergence ...)

2. **Average** one: $\bar{H} = \frac{1}{T} \sum_{t=0}^T H_t \in [0, 1]$
(convex upper bound on 0-1 loss)

3. **Best penalized** one:

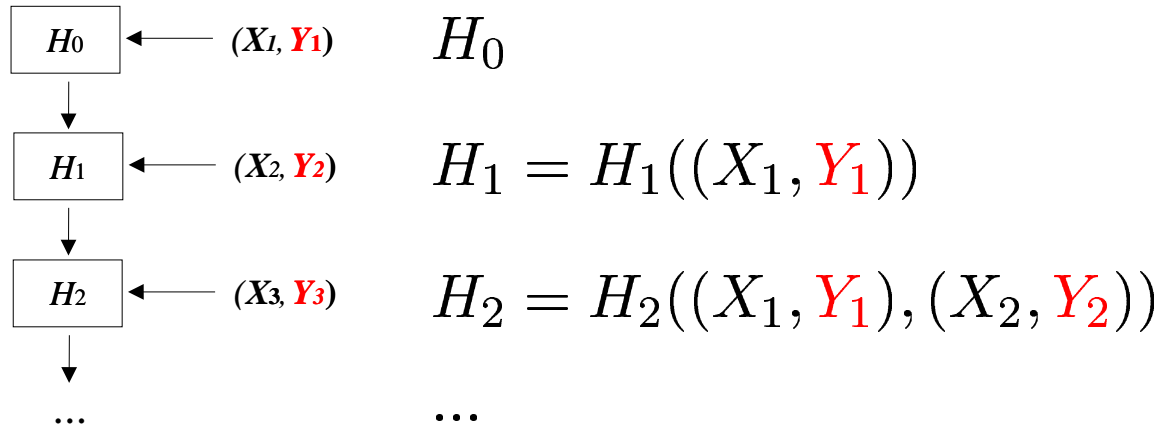
$$\text{riskemp}(H_t, t + 1) = \frac{1}{T - t} \sum_{i=t+1}^T \text{loss}(Y_i, H_t(X_i))$$



$$\hat{H} = \operatorname{argmin}_{t=0 \dots T-1} \left(\text{riskemp}(H_t, t + 1) + \underbrace{\sqrt{\frac{1}{T - t} \ln \frac{T}{\delta}}}_{\text{penalty}} \right)$$

On-line pointwise \rightarrow i.i.d. data-dependent/3

Proof technique/1



Build martingale kind of process

$$L_t = L_t((X_1, Y_1), \dots, (X_t, Y_t))$$
$$= \text{loss}(Y_t, H_{t-1}(X_t)) - \text{risk}(H_{t-1})$$

$$\mathbb{E}[\text{loss}(Y_t, H_{t-1}(X_t)) \mid (X_1, Y_1), \dots, (X_{t-1}, Y_{t-1})]$$

On-line pointwise \rightarrow i.i.d. data-dependent/3 Proof technique/2

From the very definition of L_t
 L_1, L_2, \dots, L_T is bounded ($|L_t| \leq 1$)
martingale difference sequence
w.r.t. $(X_1, Y_1), \dots, (X_T, Y_T)$

Hoeffding-Azuma:

$$\frac{1}{T} \sum_{t=1}^T \left[\text{loss}(Y_t, H_{t-1}(X_t)) - \text{risk}(H_{t-1}) \right] \approx 0$$

On-line pointwise \rightarrow i.i.d. data-dependent/3
 Proof technique/3

$$\underbrace{\frac{1}{T} \sum_{t=1}^T \text{loss}(\mathbf{Y}_t, H_{t-1}(X_t))}_{\substack{M_T \\ \# \text{ of mistakes}}} \underbrace{\approx}_{(*)} \frac{1}{T} \sum_{t=1}^T \text{risk}(H_{t-1}) \left\{ \begin{array}{l} \underbrace{\geq}_{(**)} \text{risk}(\bar{H}) \\ \underbrace{\approx}_{(***)} \text{risk}(\hat{H}) \end{array} \right.$$

(*) (Hoeffding-Azuma)

[DGL96]

(**) bounded and convex (Jensen)

(***) general bounded (Chernoff-Hoeffding)

[DGL96]

On-line pointwise \rightarrow i.i.d. data-dependent/4 Simplest bounds

Convex:
$$\mathbb{P} \left(\text{risk}(\overline{H}) \leq M_T + L \sqrt{\frac{2}{T} \ln \frac{2}{\delta}} \right) \geq 1 - \delta$$

bound on range of convex loss

More general:
$$\mathbb{P} \left(\text{risk}(\hat{H}) \leq M_T + 6 \sqrt{\frac{1}{T} \ln \frac{T}{\delta}} \right) \geq 1 - \delta$$

On-line pointwise \rightarrow i.i.d. data-dependent/5
Some applications: plug and play/1

Recall bound on Halving Algorithm for separable case:

$$M_T \leq \frac{1}{T} O(d \log(R/\gamma))$$

Just plug back into

$$\mathbb{P} \left(\text{risk}(\hat{H}) \leq M_T + 6 \sqrt{\frac{1}{T} \ln \frac{T}{\delta}} \right) \geq 1 - \delta$$

Gets

$$\mathbb{P} \left(\text{risk}(\hat{H}) \leq \frac{1}{T} O(n \log(R/\gamma)) + 6 \sqrt{\frac{1}{T} \ln \frac{T}{\delta}} \right) \geq 1 - \delta$$

Similar to [HG02]

On-line pointwise \rightarrow i.i.d. data-dependent/5

Some applications: plug and play/2

Recall bound on Kernel Perceptron:

$$M_T \leq \inf_{\gamma > 0, f \in H_K, \|f\|=1} \frac{1}{T} \left(D_\gamma(f; S) + \frac{\sqrt{\sum_{t \in \mathcal{M}} K(\mathbf{x}_t, \mathbf{x}_t)}}{\gamma} \right)$$

Separable case:

$$M_T \leq \frac{1}{T} \frac{\max_{t \in \mathcal{M}} K(\mathbf{x}_t, \mathbf{x}_t)}{\gamma^2}$$

Plug back into

$$\mathbb{P} \left(\text{risk}(\hat{H}) \leq M_T + 6 \sqrt{\frac{1}{T} \ln \frac{T}{\delta}} \right) \geq 1 - \delta$$

Similar to [BM02] for SVM

On-line pointwise \rightarrow i.i.d. data-dependent/5

Some applications: plug and play/3

Recall bound on Kernel Second-order Perceptron
(separable case)

$$M_T \leq \frac{1}{T} \frac{a + \sum_i \ln(1 + \frac{\lambda_i}{a})}{\gamma},$$

Plug into

$$\mathbb{P} \left(\text{risk}(\hat{H}) \leq M_T + 6 \sqrt{\frac{1}{T} \ln \frac{T}{\delta}} \right) \geq 1 - \delta$$

Similar to [WSTSS99] for SVM

On-line pointwise \rightarrow i.i.d. data-dependent/5

Some applications: plug and play/4

Recall bound on Perceptron for shifting targets

$$M_T \leq \frac{1}{T} \frac{4R^2}{\gamma^2} \left(\frac{1}{2} + \sum_{t=1}^T \|\mathbf{u}_{t+1} - \mathbf{u}_t\| \right)$$

Again, combine with

$$\mathbb{P} \left(\text{risk}(\hat{H}) \leq M_T + 6 \sqrt{\frac{1}{T} \ln \frac{T}{\delta}} \right) \geq 1 - \delta$$

On-line pointwise \rightarrow i.i.d. data-dependent/5

Some applications: plug and play/5

Recall bounds on label-efficient Perceptron:

$$\mathbb{E}_{(*)}[M_T] \leq \frac{1}{T} \frac{R^2}{\gamma^2}$$

$$\mathbb{E}_{(*)}[\# \text{ of labels}] = \sum_{t=1}^T \mathbb{E}_{(*)} \left[\frac{b}{b + |r_t|} \right]$$

(*) = internal randomization

Get with prob $> 1 - \delta$:

$$\mathbb{E}_{(*)}[\text{risk}(\hat{H})] \leq \frac{1}{T} \frac{R^2}{\gamma^2} + 6 \sqrt{\frac{1}{T} \ln \frac{T}{\delta}}$$

and

On-line pointwise \rightarrow i.i.d. data-dependent/5

Some applications: plug and play/6

A is your favourite batch classification alg.

Run it in on-line fashion and count # of mistakes

Still get

$$\mathbb{P} \left(\text{risk}(\hat{H}) \leq M_T + 6 \sqrt{\frac{1}{T} \ln \frac{T}{\delta}} \right) \geq 1 - \delta$$

No direct mention of, e.g., “complexity of function classes”

Try it yourself with other specific algs.

On-line pointwise \rightarrow i.i.d. data-dependent/6

Remarks

These bounds:

- are algorithm-specific (**NO** uniform convergence arguments, closer in spirit to algorithmic stability/luckiness) [BE02,HeWi02,...]
- proven by **simple** large deviation on martingales
- refer to **efficient** algs (on-line, one sweep)
- are tight (I believe ...)
- Are widely applicable (in principle)

On-line pointwise \rightarrow i.i.d. data-dependent/7
Refinements/1

$$\mathbb{P} \left(\text{risk}(\hat{H}) \leq \min_{t=0 \dots T-1} \left(M_{t,T} + 6 \sqrt{\frac{1}{T-t} \ln \frac{T}{\delta}} \right) \right) \geq 1 - \delta,$$

where $M_{t,T} = \frac{1}{T-t} \sum_{i=t+1}^T \text{loss}(Y_i, H_{i-1}(X_i))$ (loss on suffix)

Basically “ $\min_{t=0 \dots T-1}$ ” replaces “ $t = 0$ ”

On-line pointwise \rightarrow i.i.d. data-dependent/7 Refinements/2

$$\mathbb{P} \left(\text{risk}(\hat{H}) \leq M_T + O \left(\frac{1}{T} \ln \frac{T}{\delta} + \sqrt{\frac{M_T}{T} \ln \frac{T}{\delta}} \right) \right) \geq 1 - \delta,$$

$$\hat{H} = \underset{t=0 \dots T-1}{\text{argmin}} \left(\text{risk}_{\text{emp}} + \underbrace{\frac{1}{T-t} \ln \frac{T}{\delta} + \sqrt{\frac{\text{risk}_{\text{emp}}}{T-t} \ln \frac{T}{\delta}}}_{\text{penalty}} \right)$$

$$\text{risk}_{\text{emp}} = \text{risk}_{\text{emp}}(H_t, t+1) = \frac{1}{T-t} \sum_{i=t+1}^T \text{loss}(\mathbf{Y}_i, H_t(X_i))$$

(Uses Bernstein-type inequalities for martingales) [F75,DZ01]

Can be combined with “Refinement/1”

Conclusions

- Pointwise bounds for on-line algorithms directly turn to (tight) data-dependent i.i.d. bounds
- Easy plug and play
- Resulting algs. are still as efficient as on-line (one epoch over training sequence)
- Simple proofs, algorithm-specific, no uniform convergence
- Can be immediately extended to regression frameworks