

Foundations of Machine Learning

Lecture 7

Mehryar Mohri
Courant Institute and Google Research
mohri@cims.nyu.edu

On-Line Learning

Motivation

■ PAC learning:

- distribution fixed over time (training and test).
- IID assumption.

■ On-line learning:

- no distributional assumption.
- worst-case analysis (adversarial).
- mixed training and test.
- Performance measure: mistake model, regret.

This Lecture

- Prediction with expert advice
- Linear classification

General On-Line Setting

- For $t=1$ to T do
 - receive instance $x_t \in X$.
 - predict $\hat{y}_t \in Y$.
 - receive label $y_t \in Y$.
 - incur loss $L(\hat{y}_t, y_t)$.
- **Classification:** $Y = \{0, 1\}$, $L(y, y') = |y' - y|$.
- **Regression:** $Y \subseteq \mathbb{R}$, $L(y, y') = (y' - y)^2$.
- **Objective:** minimize total loss $\sum_{t=1}^T L(\hat{y}_t, y_t)$.

Prediction with Expert Advice

- For $t=1$ to T do
 - receive instance $x_t \in X$ and **advice** $y_{t,i} \in Y, i \in [1, N]$.
 - predict $\hat{y}_t \in Y$.
 - receive label $y_t \in Y$.
 - incur loss $L(\hat{y}_t, y_t)$.
- **Objective:** minimize regret, i.e., difference of total loss incurred and that of best expert.

$$\text{Regret}(T) = \sum_{t=1}^T L(\hat{y}_t, y_t) - \min_{i=1}^N \sum_{t=1}^T L(\hat{y}_{t,i}, y_t).$$

Mistake Bound Model

- **Definition:** the maximum number of mistakes a learning algorithm L makes to learn c is defined by

$$M_L(c) = \max_{x_1, \dots, x_T} |\text{mistakes}(L, c)|.$$

- **Definition:** for any concept class C the maximum number of mistakes a learning algorithm L makes is

$$M_L(C) = \max_{c \in C} M_L(c).$$

A **mistake bound** is a bound M on $M_L(C)$.

Halving Algorithm

see (Mitchell, 1997)

HALVING(H)

```
1   $H_1 \leftarrow H$ 
2  for  $t \leftarrow 1$  to  $T$  do
3      RECEIVE( $x_t$ )
4       $\hat{y}_t \leftarrow \text{MAJORITYVOTE}(H_t, x_t)$ 
5      RECEIVE( $y_t$ )
6      if  $\hat{y}_t \neq y_t$  then
7           $H_{t+1} \leftarrow \{c \in H_t : c(x_t) = y_t\}$ 
8  return  $H_{T+1}$ 
```


Halving Algorithm - Bound

(Littlestone, 1988)

■ **Theorem:** Let H be a finite hypothesis set, then

$$M_{Halving(H)} \leq \log_2 |H|.$$

■ **Proof:** At each mistake, the hypothesis set is reduced at least by half.

VC Dimension Lower Bound

(Littlestone, 1988)

- **Theorem:** Let $\text{opt}(H)$ be the optimal mistake bound for H . Then,

$$\text{VCdim}(H) \leq \text{opt}(H) \leq M_{\text{Halving}}(H) \leq \log_2 |H|.$$

- **Proof:** for a fully shattered set, form a complete binary tree of the mistakes with height $\text{VCdim}(H)$.

Weighted Majority Algorithm

(Littlestone and Warmuth, 1988)

WEIGHTED-MAJORITY(N experts) $\triangleright y_t, y_{t,i} \in \{0, 1\}.$

1 **for** $i \leftarrow 1$ **to** N **do** $\beta \in [0, 1).$
2 $w_{1,i} \leftarrow 1$
3 **for** $t \leftarrow 1$ **to** T **do**
4 RECEIVE(x_t)
5 $\hat{y}_t \leftarrow 1_{\sum_{y_{t,i}=1}^N w_t \geq \sum_{y_{t,i}=0}^N w_t}$ \triangleright weighted majority vote
6 RECEIVE(y_t)
7 **if** $\hat{y}_t \neq y_t$ **then**
8 **for** $i \leftarrow 1$ **to** N **do**
9 **if** $(y_{t,i} \neq y_t)$ **then**
10 $w_{t+1,i} \leftarrow \beta w_{t,i}$
11 **else** $w_{t+1,i} \leftarrow w_{t,i}$
12 **return** w_{T+1}

Weighted Majority - Bound

- **Theorem:** Let m_t be the number of mistakes made by the WM algorithm till time t and m_t^* that of the best expert. Then, for all t ,

$$m_t \leq \frac{\log N + m_t^* \log \frac{1}{\beta}}{\log \frac{2}{1+\beta}}.$$

- Thus, $m_t \leq O(\log N) + \text{constant} \times \text{best expert}$.
- Realizable case: $m_t \leq O(\log N)$.
- Halving algorithm: $\beta = 0$.

Weighted Majority - Proof

■ **Potential:** $\Phi_t = \sum_{i=1}^N w_{t,i}$.

■ **Upper bound:** after each error,

$$\Phi_{t+1} \leq \left[1/2 + 1/2 \beta\right] \Phi_t = \left[\frac{1+\beta}{2}\right] \Phi_t.$$

Thus, $\Phi_t \leq \left[\frac{1+\beta}{2}\right]^{m_t} N.$

■ **Lower bound:** for any expert i , $\Phi_t \geq w_{t,i} = \beta^{m_{t,i}}$.

■ **Comparison:** $\beta^{m_t^*} \leq \left[\frac{1+\beta}{2}\right]^{m_t} N$

$$\Rightarrow m_t^* \log \beta \leq \log N + m_t \log \left[\frac{1+\beta}{2}\right]$$

$$\Rightarrow m_t \log \left[\frac{2}{1+\beta}\right] \leq \log N + m_t^* \log \frac{1}{\beta}.$$

Weighted Majority - Notes

- **Advantage:** remarkable bound requiring no assumption.
- **Disadvantage:** no deterministic algorithm can achieve a regret $R_T = o(T)$ with the binary loss.
 - better guarantee with randomized WM.
 - better guarantee for WM with convex losses.

Exponential Weighted Average

■ Algorithm:

total loss incurred by
expert i up to time t

- weight update: $w_{t+1,i} \leftarrow w_{t,i} e^{-\eta L(\hat{y}_{t,i}, y_t)} = e^{-\eta L_{t,i}}$.
- prediction: $\hat{y}_t = \frac{\sum_{i=1}^N w_{t,i} y_{t,i}}{\sum_{i=1}^N w_{t,i}}$.

■ **Theorem:** assume that L is convex in its first argument and takes values in $[0, 1]$. Then, for any $\eta > 0$ and any sequence $y_1, \dots, y_T \in Y$, the regret at T satisfies

$$\text{Regret}(T) \leq \frac{\log N}{\eta} + \frac{\eta T}{8}.$$

For $\eta = \sqrt{8 \log N / T}$,

$$\text{Regret}(T) \leq \sqrt{(T/2) \log N}.$$

Exponential Weighted Avg - Proof

■ **Potential:** $\Phi_t = \log \sum_{i=1}^N w_{t,i}$.

■ **Upper bound:**

$$\begin{aligned}\Phi_t - \Phi_{t-1} &= \log \frac{\sum_{i=1}^N w_{t-1,i} e^{-\eta L(\hat{y}_{t,i}, y_t)}}{\sum_{i=1}^N w_{t-1,i}} \\&= \log \left(\mathbb{E}_{w_{t-1}} [e^{-\eta L(\hat{y}_{t,i}, y_t)}] \right) \\&= \log \left(\mathbb{E}_{w_{t-1}} \left[\exp \left(-\eta \left(L(\hat{y}_{t,i}, y_t) - \mathbb{E}_{w_{t-1}} [L(\hat{y}_{t,i}, y_t)] \right) - \eta \mathbb{E}_{w_{t-1}} [L(\hat{y}_{t,i}, y_t)] \right) \right] \right) \\&\leq -\eta \mathbb{E}_{w_{t-1}} [L(\hat{y}_{t,i}, y_t)] + \frac{\eta^2}{8} \quad (\text{Hoeffding's ineq.}) \\&\leq -\eta L(\mathbb{E}_{w_{t-1}} [\hat{y}_{t,i}], y_t) + \frac{\eta^2}{8} \quad (\text{convexity of first arg. of } L) \\&= -\eta L(\hat{y}_t, y_t) + \frac{\eta^2}{8}.\end{aligned}$$

Exponential Weighted Avg - Proof

■ **Upper bound:** summing up the inequalities yields

$$\Phi_T - \Phi_0 \leq -\eta \sum_{t=1}^T L(\hat{y}_t, y_t) + \frac{\eta^2 T}{8}.$$

■ **Lower bound:**

$$\begin{aligned} \Phi_T - \Phi_0 &= \log \sum_{i=1}^N e^{-\eta L_{T,i}} - \log N \geq \log \max_{i=1}^N e^{-\eta L_{T,i}} - \log N \\ &= -\eta \min_{i=1}^N L_{T,i} - \log N. \end{aligned}$$

■ **Comparison:**

$$\begin{aligned} -\eta \min_{i=1}^N L_{T,i} - \log N &\leq -\eta \sum_{t=1}^T L(\hat{y}_t, y_t) + \frac{\eta^2 T}{8} \\ \Rightarrow \sum_{t=1}^T L(\hat{y}_t, y_t) - \min_{i=1}^N L_{T,i} &\leq \frac{\log N}{\eta} + \frac{\eta T}{8}. \end{aligned}$$

Exponential Weighted Avg - Notes

- **Advantage:** bound on regret per bound is of the form $\frac{R_T}{T} = O\left(\sqrt{\frac{\log(N)}{T}}\right)$.
- **Disadvantage:** choice of η requires knowledge of horizon T .

Doubling Trick

- **Idea:** divide time into periods $[2^k, 2^{k+1} - 1]$ of length 2^k with $k = 0, \dots, n$, $T \geq 2^n - 1$, and choose $\eta_k = \sqrt{\frac{8 \log N}{2^k}}$ in each period.
- **Theorem:** with the same assumptions as before, for any T , the following holds:

$$\text{Regret}(T) \leq \frac{\sqrt{2}}{\sqrt{2} - 1} \sqrt{(T/2) \log N} + \sqrt{\log N/2}.$$

Doubling Trick - Proof

■ By the previous theorem, for any $I_k = [2^k, 2^{k+1} - 1]$,

$$L_{I_k} - \min_{i=1}^N L_{I_k, i} \leq \sqrt{2^k / 2 \log N}.$$

Thus,
$$\begin{aligned} L_T = \sum_{k=0}^n L_{I_k} &\leq \sum_{k=0}^n \min_{i=1}^N L_{I_k, i} + \sum_{k=0}^n \sqrt{2^k (\log N) / 2} \\ &\leq \min_{i=1}^N L_{T, i} + \sum_{k=0}^n 2^{\frac{k}{2}} \sqrt{(\log N) / 2}. \end{aligned}$$

with

$$\sum_{i=0}^n 2^{\frac{k}{2}} = \frac{\sqrt{2}^{n+1} - 1}{\sqrt{2} - 1} = \frac{2^{(n+1)/2} - 1}{\sqrt{2} - 1} \leq \frac{\sqrt{2}\sqrt{T+1} - 1}{\sqrt{2} - 1} \leq \frac{\sqrt{2}(\sqrt{T} + 1) - 1}{\sqrt{2} - 1} \leq \frac{\sqrt{2}\sqrt{T}}{\sqrt{2} - 1} + 1.$$

Notes

- Doubling trick used in a variety of other contexts and proofs.
- More general method, learning parameter function of time: $\eta_t = \sqrt{(8 \log N)/t}$. Constant factor improvement:

$$\text{Regret}(T) \leq 2\sqrt{(T/2) \log N} + \sqrt{(1/8) \log N}.$$

This Lecture

- Prediction with expert advice
- Linear classification

Perceptron Algorithm

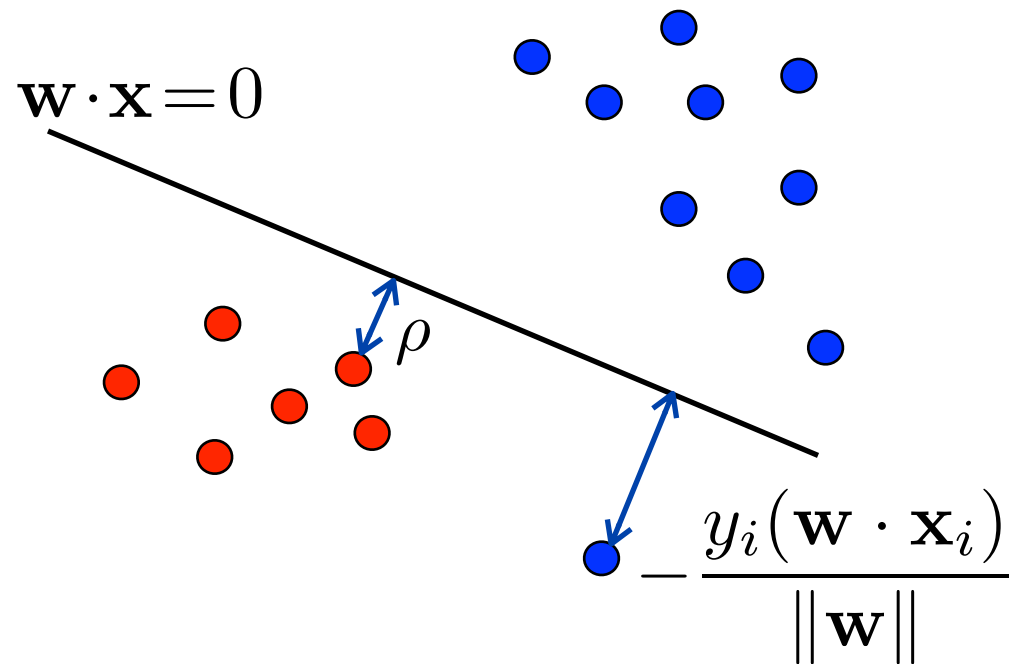
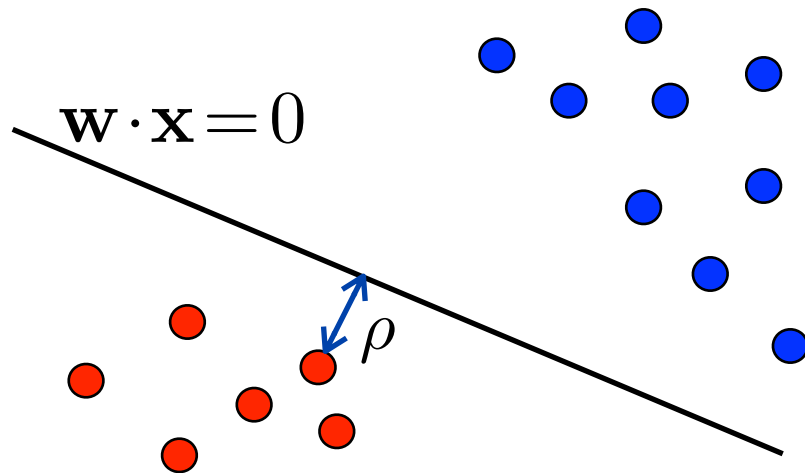
(Rosenblatt, 1958)

PERCEPTRON(\mathbf{w}_0)

```
1   $\mathbf{w}_1 \leftarrow \mathbf{w}_0$        $\triangleright$  typically  $\mathbf{w}_0 = \mathbf{0}$ 
2  for  $t \leftarrow 1$  to  $T$  do
3      RECEIVE( $\mathbf{x}_t$ )
4       $\hat{y}_t \leftarrow \text{sgn}(\mathbf{w}_t \cdot \mathbf{x}_t)$ 
5      RECEIVE( $y_t$ )
6      if ( $\hat{y}_t \neq y_t$ ) then
7           $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t + y_t \mathbf{x}_t$      $\triangleright$  more generally  $\eta y_t \mathbf{x}_t, \eta > 0$ 
8      else  $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t$ 
9  return  $\mathbf{w}_{T+1}$ 
```

Separating Hyperplane

■ Margin and errors



Perceptron = Stochastic Gradient Descent

- **Objective function:** convex but not differentiable.

$$F(\mathbf{w}) = \frac{1}{T} \sum_{t=1}^T \max \left(0, -y_t(\mathbf{w} \cdot \mathbf{x}_t) \right) = \mathbb{E}_{\mathbf{x} \sim \hat{D}} [f(\mathbf{w}, \mathbf{x})]$$

with $f(\mathbf{w}, \mathbf{x}) = \max \left(0, -y(\mathbf{w} \cdot \mathbf{x}) \right)$.

- **Stochastic gradient:** for each \mathbf{x}_t , the update is

$$\mathbf{w}_{t+1} \leftarrow \begin{cases} \mathbf{w}_t - \eta \nabla_{\mathbf{w}} f(\mathbf{w}_t, \mathbf{x}_t) & \text{if differentiable} \\ \mathbf{w}_t & \text{otherwise,} \end{cases}$$

where $\eta > 0$ is a learning rate parameter.

- **Here:**
$$\mathbf{w}_{t+1} \leftarrow \begin{cases} \mathbf{w}_t + \eta y_t \mathbf{x}_t & \text{if } y_t(\mathbf{w}_t \cdot \mathbf{x}_t) < 0 \\ \mathbf{w}_t & \text{otherwise.} \end{cases}$$

Perceptron Algorithm - Bound

(Novikoff, 1962)

- **Theorem:** Assume that $\|x_t\| \leq R$ for all $t \in [1, T]$ and that for some $\rho > 0$ and $\mathbf{v} \in \mathbb{R}^N$, for all $t \in [1, T]$,

$$\rho \leq \frac{y_t(\mathbf{v} \cdot \mathbf{x}_t)}{\|\mathbf{v}\|}.$$

Then, the number of mistakes made by the perceptron algorithm is bounded by R^2 / ρ^2 .

- **Proof:** Let I be the set of t s at which there is an update and let M be the total number of updates.

- Summing up the assumption inequalities gives:

$$\begin{aligned} M\rho &\leq \frac{\mathbf{v} \cdot \sum_{t \in I} y_t \mathbf{x}_t}{\|\mathbf{v}\|} \\ &= \frac{\mathbf{v} \cdot \sum_{t \in I} (\mathbf{w}_{t+1} - \mathbf{w}_t)}{\|\mathbf{v}\|} \quad (\text{definition of updates}) \\ &= \frac{\mathbf{v} \cdot \mathbf{w}_{T+1}}{\|\mathbf{v}\|} \\ &\leq \|\mathbf{w}_{T+1}\| \quad (\text{Cauchy-Schwarz ineq.}) \\ &= \|\mathbf{w}_{t_m} + y_{t_m} \mathbf{x}_{t_m}\| \quad (t_m \text{ largest } t \text{ in } I) \\ &= \left[\|\mathbf{w}_{t_m}\|^2 + \|\mathbf{x}_{t_m}\|^2 + 2 \underbrace{y_{t_m} \mathbf{w}_{t_m} \cdot \mathbf{x}_{t_m}}_{\leq 0} \right]^{1/2} \\ &\leq \left[\|\mathbf{w}_{t_m}\|^2 + R^2 \right]^{1/2} \\ &\leq \left[MR^2 \right]^{1/2} = \sqrt{M}R. \quad (\text{applying the same to previous } ts \text{ in } I) \end{aligned}$$

- **Notes:**
 - bound independent of dimension and tight.
 - convergence can be slow for small margin, it can be in $\Omega(2^N)$.
 - among the many variants: **voted perceptron algorithm**. Predict according to

$$\text{sign}\left(\left(\sum_{t \in I} c_t \mathbf{w}_t\right) \cdot \mathbf{x}\right),$$

where c_t is the number of iterations \mathbf{w}_t survives.

- $\{x_t : t \in I\}$ are the **support vectors** for the perceptron algorithm.
- non-separable case: **does not converge**.

Perceptron - Leave-One-Out Analysis

■ **Theorem:** Let h_S be the hypothesis returned by the perceptron algorithm for sample $S = (x_1, \dots, x_T) \sim D$ and let $M(S)$ be the number of updates defining h_S . Then,

$$\mathbb{E}_{S \sim D^m} [R(h_S)] \leq \mathbb{E}_{S \sim D^{m+1}} \left[\frac{\min(M(S), R_{m+1}^2 / \rho_{m+1}^2)}{m+1} \right].$$

■ **Proof:** Let $S \sim D^{m+1}$ be a sample linearly separable and let $\mathbf{x} \in S$. If $h_{S-\{\mathbf{x}\}}$ misclassifies \mathbf{x} , then \mathbf{x} must be a ‘support vector’ for h_S (update at \mathbf{x}). Thus,

$$\hat{R}_{\text{loo}}(\text{perceptron}) \leq \frac{M(S)}{m+1}.$$

SVMs - Leave-One-Out Analysis

(Vapnik, 1995)

- **Theorem:** let h_S be the optimal hyperplane for a sample S and let $N_{SV}(S)$ be the number of support vectors defining h_S . Then,

$$\mathbb{E}_{S \sim D^m} [R(h_S)] \leq \mathbb{E}_{S \sim D^{m+1}} \left[\frac{\min(N_{SV}(S), R_{m+1}^2 / \rho_{m+1}^2)}{m+1} \right].$$

- **Proof:** one part proven in lecture 4. The other part due to $\alpha_i \geq 1/R_{m+1}^2$ for \mathbf{x}_i misclassified by SVMs.

Comparison

- Bounds on expected error, not high probability statements.
- Leave-one-out bounds not sufficient to distinguish SVMs and perceptron algorithm. Note however:
 - same maximum margin ρ_{m+1} can be used in both.
 - but different radius R_{m+1} of support vectors.
- Difference: margin distribution.

Non-Separable Case - LI Bound

(MM and Rostamizadeh, 2013)

■ **Theorem:** let I denote the set of rounds at which the Perceptron algorithm makes an update when processing $\mathbf{x}_1, \dots, \mathbf{x}_T$ and let $M_T = |I|$. Then,

$$M_T \leq \inf_{\rho > 0, \|u\|_2 \leq 1} \sum_{t \in I} \left(1 - \frac{y_t(\mathbf{u} \cdot \mathbf{x}_t)}{\rho} \right)_+ + \frac{\sqrt{\sum_{t \in I} \|\mathbf{x}_t\|^2}}{\rho}.$$

● when $\|\mathbf{x}_t\| \leq R$ for all $t \in I$, this implies

$$M_T \leq \inf_{\rho > 0, \|u\|_2 \leq 1} \left(\frac{R}{\rho} + \sqrt{\|\mathbf{L}_\rho(\mathbf{u})\|_1} \right)^2,$$

$$\text{where } \mathbf{L}_\rho(\mathbf{u}) = \left[\left(1 - \frac{y_t(\mathbf{u} \cdot \mathbf{x}_t)}{\rho} \right)_+ \right]_{t \in I}.$$

- **Proof:** for any t , $1 - \frac{y_t(\mathbf{u} \cdot \mathbf{x}_t)}{\rho} \leq \left(1 - \frac{y_t(\mathbf{u} \cdot \mathbf{x}_t)}{\rho}\right)_+$, summing up these inequalities for $t \in I$ yields:

$$M_T \leq \sum_{t \in I} \left(1 - \frac{y_t(\mathbf{u} \cdot \mathbf{x}_t)}{\rho}\right)_+ + \sum_{t \in I} \frac{y_t(\mathbf{u} \cdot \mathbf{x}_t)}{\rho}.$$

- upper-bounding $\sum_{t \in I} (y_t \mathbf{u} \cdot \mathbf{x}_t)$ as in the proof for separable case shows the first inequality.
- the second inequality is obtained by solving

$$M_T \leq \|\mathbf{L}_\rho(\mathbf{u})\|_1 + \frac{R}{\rho} \sqrt{M_T},$$

which gives $\sqrt{M_T} \leq \frac{\frac{R}{\rho} + \sqrt{\frac{R^2}{\rho^2} + 4\|\mathbf{L}_\rho(\mathbf{u})\|_1}}{2}.$

Non-Separable Case - L2 Bound

(Freund and Schapire, 1998; MM and Rostamizadeh, 2013)

■ **Theorem:** let I denote the set of rounds at which the Perceptron algorithm makes an update when processing $\mathbf{x}_1, \dots, \mathbf{x}_T$ and let $M_T = |I|$. Then,

$$M_T \leq \inf_{\rho > 0, \|u\|_2 \leq 1} \left[\frac{\|\mathbf{L}_\rho(\mathbf{u})\|_2}{2} + \sqrt{\frac{\|\mathbf{L}_\rho(\mathbf{u})\|_2^2}{4} + \frac{\sqrt{\sum_{t \in I} \|\mathbf{x}_t\|^2}}{\rho}} \right]^2.$$

● when $\|\mathbf{x}_t\| \leq R$ for all $t \in I$, this implies

$$M_T \leq \inf_{\rho > 0, \|u\|_2 \leq 1} \left(\frac{R}{\rho} + \|\mathbf{L}_\rho(\mathbf{u})\|_2 \right)^2,$$

$$\text{where } \mathbf{L}_\rho(\mathbf{u}) = \left[\left(1 - \frac{y_t(\mathbf{u} \cdot \mathbf{x}_t)}{\rho} \right)_+ \right]_{t \in I}.$$

- **Proof:** Reduce problem to separable case in higher dimension. Let $l_t = \left(1 - \frac{y_t \mathbf{u} \cdot \mathbf{x}_t}{\rho}\right)_+ \mathbf{1}_{t \in I}$, for $t \in [1, T]$.
- Mapping (similar to trivial mapping):

$(N+t)$ th component

$$\mathbf{x}_t = \begin{bmatrix} x_{t,1} \\ \vdots \\ x_{t,N} \end{bmatrix} \rightarrow \mathbf{x}'_t = \begin{bmatrix} x_{t,1} \\ \vdots \\ x_{t,N} \\ 0 \\ \vdots \\ 0 \\ \Delta \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

$\mathbf{u} \rightarrow \mathbf{u}' = \begin{bmatrix} \frac{u_1}{Z} \\ \vdots \\ \frac{u_N}{Z} \\ \frac{y_1 \rho l_1}{\Delta Z} \\ \vdots \\ \frac{y_T \rho l_T}{\Delta Z} \end{bmatrix}$

$$\|\mathbf{u}'\| = 1 \implies Z = \sqrt{1 + \frac{\rho^2 \|\mathbf{L}_\rho(\mathbf{u})\|^2}{\Delta^2}}.$$

- Observe that the Perceptron algorithm makes the same predictions and makes updates at the same rounds when processing $\mathbf{x}'_1, \dots, \mathbf{x}'_T$.
- For any $t \in I$,

$$\begin{aligned}
 y_t(\mathbf{u}' \cdot \mathbf{x}'_t) &= y_t \left(\frac{\mathbf{u} \cdot \mathbf{x}_t}{Z} + \Delta \frac{y_t \rho l_t}{Z \Delta} \right) \\
 &= \frac{y_t \mathbf{u} \cdot \mathbf{x}_t}{Z} + \frac{\rho l_t}{Z} \\
 &= \frac{1}{Z} (y_t \mathbf{u} \cdot \mathbf{x}_t + [\rho - y_t(\mathbf{u} \cdot \mathbf{x}_t)]_+) \geq \frac{\rho}{Z}.
 \end{aligned}$$

- Summing up and using the proof in the separable case yields:

$$M_T \frac{\rho}{Z} \leq \sum_{t \in I} y_t(\mathbf{u}' \cdot \mathbf{x}'_t) \leq \sqrt{\sum_{t \in I} \|\mathbf{x}'_t\|^2}.$$

- The inequality can be rewritten as

$$M_T^2 \leq \left(\frac{1}{\rho^2} + \frac{\|\mathbf{L}_\rho(\mathbf{u})\|^2}{\Delta^2} \right) (r^2 + M_T \Delta^2) = \frac{r^2}{\rho^2} + \frac{r^2 \|\mathbf{L}_\rho(\mathbf{u})\|^2}{\Delta^2} + \frac{M_T \Delta^2}{\rho^2} + M_T \|\mathbf{L}_\rho(\mathbf{u})\|^2;$$

where $r = \sqrt{\sum_{t \in I} \|\mathbf{x}_t\|^2}$.

- Selecting Δ to minimize the bound gives $\Delta^2 = \frac{\rho \|\mathbf{L}_\rho(\mathbf{u})\|_2 r}{\sqrt{M_T}}$ and leads to

$$M_T^2 \leq \frac{r^2}{\rho^2} + 2 \frac{\sqrt{M_T} \|\mathbf{L}_\rho(\mathbf{u})\| r}{\rho} + M_T \|\mathbf{L}_\rho(\mathbf{u})\|^2 = \left(\frac{r}{\rho} + \sqrt{M_T} \|\mathbf{L}_\rho(\mathbf{u})\|_2 \right)^2.$$

- Solving the second-degree inequality

$$M_T - \sqrt{M_T} \|\mathbf{L}_\rho(\mathbf{u})\|_2 - \frac{r}{\rho} \leq 0$$

yields directly the first statement. The second one results from replacing r with $\sqrt{M_T} R$.

Dual Perceptron Algorithm

DUAL-PERCEPTRON(α^0)

```
1   $\alpha \leftarrow \alpha^0$        $\triangleright$  typically  $\alpha^0 = \mathbf{0}$ 
2  for  $t \leftarrow 1$  to  $T$  do
3      RECEIVE( $\mathbf{x}_t$ )
4       $\hat{y}_t \leftarrow \text{sgn}(\sum_{s=1}^T \alpha_s y_s (\mathbf{x}_s \cdot \mathbf{x}_t))$ 
5      RECEIVE( $y_t$ )
6      if ( $\hat{y}_t \neq y_t$ ) then
7           $\alpha_t \leftarrow \alpha_t + 1$ 
8  return  $\alpha$ 
```

Kernel Perceptron Algorithm

(Aizerman et al., 1964)

K PDS kernel.

KERNEL-PERCEPTRON(α^0)

```
1   $\alpha \leftarrow \alpha^0$        $\triangleright$  typically  $\alpha^0 = \mathbf{0}$ 
2  for  $t \leftarrow 1$  to  $T$  do
3      RECEIVE( $x_t$ )
4       $\hat{y}_t \leftarrow \text{sgn}(\sum_{s=1}^T \alpha_s y_s K(x_s, x_t))$ 
5      RECEIVE( $y_t$ )
6      if ( $\hat{y}_t \neq y_t$ ) then
7           $\alpha_t \leftarrow \alpha_t + 1$ 
8  return  $\alpha$ 
```

Winnow Algorithm

(Littlestone, 1988)

WINNOW(η)

```
1   $w_1 \leftarrow \mathbf{1}/N$ 
2  for  $t \leftarrow 1$  to  $T$  do
3      RECEIVE( $\mathbf{x}_t$ )
4       $\hat{y}_t \leftarrow \text{sgn}(\mathbf{w}_t \cdot \mathbf{x}_t)$   $\triangleright y_t \in \{-1, +1\}$ 
5      RECEIVE( $y_t$ )
6      if ( $\hat{y}_t \neq y_t$ ) then
7           $Z_t \leftarrow \sum_{i=1}^N w_{t,i} \exp(\eta y_t x_{t,i})$ 
8          for  $i \leftarrow 1$  to  $N$  do
9               $w_{t+1,i} \leftarrow \frac{w_{t,i} \exp(\eta y_t x_{t,i})}{Z_t}$ 
10         else  $\mathbf{w}_{t+1} \leftarrow \mathbf{w}_t$ 
11 return  $\mathbf{w}_{T+1}$ 
```


Notes

- Winnow = weighted majority:
 - for $y_{t,i} = x_{t,i} \in \{-1, +1\}$, $\text{sgn}(\mathbf{w}_t \cdot \mathbf{x}_t)$ coincides with the majority vote.
 - multiplying by e^η or $e^{-\eta}$ the weight of correct or incorrect experts, is equivalent to multiplying by $\beta = e^{-2\eta}$ the weight of incorrect ones.
- Relationships with other algorithms: e.g., boosting and Perceptron (Winnow and Perceptron can be viewed as special instances of a general family).

Winnnow Algorithm - Bound

- **Theorem:** Assume that $\|x_t\|_\infty \leq R_\infty$ for all $t \in [1, T]$ and that for some $\rho_\infty > 0$ and $\mathbf{v} \in \mathbb{R}^N$, $\mathbf{v} \geq 0$ for all $t \in [1, T]$,

$$\rho_\infty \leq \frac{y_t(\mathbf{v} \cdot \mathbf{x}_t)}{\|\mathbf{v}\|_1}.$$

Then, the number of mistakes made by the Winnnow algorithm is bounded by $2(R_\infty^2 / \rho_\infty^2) \log N$.

- **Proof:** Let I be the set of t s at which there is an update and let M be the total number of updates.

Winnow Algorithm - Bound

■ **Potential:** $\Phi_t = \sum_{i=1}^N \frac{v_i}{\|\mathbf{v}\|} \log \frac{v_i / \|\mathbf{v}\|}{w_{t,i}}.$ (relative entropy)

■ **Upper bound:** for each t in I ,

$$\begin{aligned}\Phi_{t+1} - \Phi_t &= \sum_{i=1}^N \frac{v_i}{\|\mathbf{v}\|_1} \log \frac{w_{t,i}}{w_{t+1,i}} \\ &= \sum_{i=1}^N \frac{v_i}{\|\mathbf{v}\|_1} \log \frac{Z_t}{\exp(\eta y_t x_{t,i})} \\ &= \log Z_t - \eta \sum_{i=1}^N \frac{v_i}{\|\mathbf{v}\|_1} y_t x_{t,i} \\ &\leq \log \left[\sum_{i=1}^N w_{t,i} \exp(\eta y_t x_{t,i}) \right] - \eta \rho_\infty \\ &= \log \mathbb{E}_{\mathbf{w}_t} \left[\exp(\eta y_t \mathbf{x}_t) \right] - \eta \rho_\infty\end{aligned}$$

$$\begin{aligned}(\text{Hoeffding}) &\leq \log \left[\exp(\eta^2 (2R_\infty)^2 / 8) \right] + \eta y_t \mathbf{w}_t \cdot \mathbf{x}_t - \eta \rho_\infty \\ &\leq \eta^2 R_\infty^2 / 2 - \eta \rho_\infty.\end{aligned}$$

Winnow Algorithm - Bound

- **Upper bound:** summing up the inequalities yields

$$\Phi_{T+1} - \Phi_1 \leq M(\eta^2 R_\infty^2 / 2 - \eta \rho_\infty).$$

- **Lower bound:** note that

$$\Phi_1 = \sum_{i=1}^N \frac{v_i}{\|\mathbf{v}\|_1} \log \frac{v_i / \|\mathbf{v}\|_1}{1/N} = \log N + \sum_{i=1}^N \frac{v_i}{\|\mathbf{v}\|_1} \log \frac{v_i}{\|\mathbf{v}\|_1} \leq \log N$$

and for all t , $\Phi_t \geq 0$ (property of relative entropy).

Thus, $\Phi_{T+1} - \Phi_1 \geq 0 - \log N = -\log N$.

- **Comparison:** $-\log N \leq M(\eta^2 R_\infty^2 / 2 - \eta \rho_\infty)$. For $\eta = \frac{\rho_\infty}{R_\infty^2}$ we obtain

$$M \leq 2 \log N \frac{R_\infty^2}{\rho_\infty^2}.$$

Notes

- Comparison with perceptron bound:
 - dual norms: norms for \mathbf{x}_t and \mathbf{v} .
 - similar bounds with different norms.
 - each advantageous in different cases:
 - Winnow bound favorable when a sparse set of experts can predict well. For example, if $\mathbf{v} = \mathbf{e}_1$ and $\mathbf{x}_t \in \{\pm 1\}^N$, $\log N$ vs N .
 - Perceptron favorable in opposite situation.

Conclusion

■ On-line learning:

- wide and fast-growing literature.
- many related topics, e.g., game theory, text compression, convex optimization.
- online to batch bounds and techniques.
- online version of batch algorithms, e.g., regression algorithms (next lecture).

References

- Aizerman, M.A., Braverman, E. M., & Rozonoer, L. I. (1964). Theoretical foundations of the potential function method in pattern recognition learning. *Automation and Remote Control*, 25, 821-837.
- Nicolò Cesa-Bianchi, Alex Conconi, Claudio Gentile: On the Generalization Ability of On-Line Learning Algorithms. *IEEE Transactions on Information Theory* 50(9): 2050-2057. 2004.
- Nicolò Cesa-Bianchi and Gábor Lugosi. *Prediction, learning, and games*. Cambridge University Press, 2006.
- Yoav Freund and Robert Schapire. Large margin classification using the perceptron algorithm. In *Proceedings of COLT 1998*. ACM Press, 1998.
- Nick Littlestone. From On-Line to Batch Learning. *COLT 1989*: 269-284.
- Nick Littlestone. "Learning Quickly When Irrelevant Attributes Abound: A New Linear-threshold Algorithm" *Machine Learning* 285-318(2). 1988.

References

- Nick Littlestone, Manfred K. Warmuth: The Weighted Majority Algorithm. *FOCS* 1989: 256-261.
- Tom Mitchell. *Machine Learning*, McGraw Hill, 1997.
- Mehryar Mohri and Afshin Rostamizadeh. Perceptron Mistake Bounds. arXiv:1305.0208, 2013.
- Novikoff, A. B. (1962). On convergence proofs on perceptrons. *Symposium on the Mathematical Theory of Automata*, 12, 615-622. Polytechnic Institute of Brooklyn.
- Rosenblatt, Frank, *The Perceptron: A Probabilistic Model for Information Storage and Organization in the Brain*, Cornell Aeronautical Laboratory, Psychological Review, v65, No. 6, pp. 386-408, 1958.
- Vladimir N. Vapnik. *Statistical Learning Theory*. Wiley-Interscience, New York, 1998.