# Online Learning

Jordan Boyd-Graber
University of Colorado Boulder
LECTURE 13

Slides adapted from Mohri

- 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.

- For $t = 1$ to $T$:
  - Get instance $x_t \in X$
  - Predict $\hat{y}_t \in Y$
  - Get true label $y_t \in Y$
  - Incur loss $L(\hat{y}_t, y_t)$
- Classification: $Y = \{0, 1\}$, $L(y, y') = |y' - y|$
- Regression: $Y \subset \mathbb{R}$, $L(y, y') = (y' - y)^2$

- For $t = 1$ to $T$:
  - Get instance $x_t \in X$
  - Predict $\hat{y}_t \in Y$
  - Get true label $y_t \in Y$
  - Incur loss $L(\hat{y}_t, y_t)$
- Classification: $Y = \{0, 1\}$, $L(y, y') = |y' - y|$
- Regression: $Y \subset \mathbb{R}, L(y, y') = (y' - y)^2$
- **Objective**: Minimize total loss $\sum_t L(\hat{y}_t, y_t)$

**Plan**

Experts

Perceptron Algorithm

**Prediction with Expert Advice**

- For $t = 1$ to $T$:
    - Get instance $x_t \in X$ and advice $a_t, i \in Y, i \in [1, N]$
    - Predict $\hat{y}_t \in Y$
    - Get true label $y_t \in Y$
    - Incur loss $L(\hat{y}_t, y_t)$

**Prediction with Expert Advice**

- For $t = 1$ to $T$:
  - Get instance $x_t \in X$ and advice $a_t, i \in Y, i \in [1, N]$
  - Predict $\hat{y}_t \in Y$
  - Get true label $y_t \in Y$
  - Incur loss $L(\hat{y}_t, y_t)$
- **Objective**: Minimize regret, i.e., difference of total loss vs. best expert

$$\text{Regret}(T) = \sum_t L(\hat{y}_t, y_t) - \min_i \sum_t L(a_{t,i}, y_t) \tag{1}$$

## Mistake Bound Model

- Define the maximum number of mistakes a learning algorithm $L$ makes to learn a concept $c$ over any set of examples (until it's perfect).

$$M_L(c) = \max_{x_1,\ldots,x_T} |\text{mistakes}(L,c)| \qquad (2)$$

- For any concept class $C$, this is the max over concepts $c$.

$$M_L(C) = \max_{c \in C} M_L(c) \qquad (3)$$

**Halving Algorithm**

$H_1 \leftarrow H$;
**for** $t \leftarrow 1 \dots T$ **do**
  Receive $x_t$;
  $\hat{y}_t \leftarrow \text{Majority}(H_t, \vec{a}_t, x_t)$;
  Receive $y_t$;
  **if** $\hat{y}_t \neq y_t$ **then**
    $H_{t+1} \leftarrow \{c \in H_t : c(x_t) = y_t\}$;
**return** $H_{T+1}$
          **Algorithm 1:** The Halving Algorithm (Mitchell, 1997)

**Halving Algorithm Bound (Littlestone, 1998)**

- For a finite hypothesis set

$$M_{\text{Halving}(H)} \leq \lg |H| \tag{4}$$

- After each mistake, the hypothesis set is reduced by at least by half

**Halving Algorithm Bound (Littlestone, 1998)**

- For a finite hypothesis set

$$M_{\text{Halving}(H)} \leq \lg |H| \qquad (4)$$

- After each mistake, the hypothesis set is reduced by at least by half
- Consider the optimal mistake bound opt($H$). Then

$$\text{VC}(H) \leq \text{opt}(H) \leq M_{\text{Halving}(H)} \leq \lg |H| \qquad (5)$$

- For a fully shattered set, form a binary tree of mistakes with height VC($H$)

## Weighted Majority (Littlestone and Warmuth, 1998)

**for** $t \leftarrow 1 \ldots N$ **do**
$\quad \mid \quad w_{1,i} \leftarrow 1$;
**for** $t \leftarrow 1 \ldots T$ **do**
$\quad \mid \quad$ Receive $x_t$;
$\quad \mid \quad \hat{y}_t \leftarrow \mathbb{1} \left[ \sum_{y_{t,i}=1} w_t \geq \sum_{y_{t,i}=0} w_t \right]$;
$\quad \mid \quad$ Receive $y_t$;
$\quad \mid \quad$ **if** $\hat{y}_t \neq y_t$ **then**
$\quad \mid \quad \quad \mid \quad$ **for** $t \leftarrow 1 \ldots N$ **do**
$\quad \mid \quad \quad \mid \quad \quad \mid \quad$ **if** $\hat{y}_t \neq y_t$ **then**
$\quad \mid \quad \quad \mid \quad \quad \mid \quad \quad \mid \quad w_{t+1,i} \leftarrow \beta w_{t,i}$;
$\quad \mid \quad \quad \mid \quad \quad \mid \quad$ **else**
$\quad \mid \quad \quad \mid \quad \quad \mid \quad \quad \mid \quad w_{t+1,i} \leftarrow w_{t,i}$
**return** $w_{T+1}$

- Weights for every expert
- Classifications in favor of side with higher total weight ($y \in \{0, 1\}$)
- Experts that are wrong get their weights decreased ($\beta \in [0, 1]$)
- If you're right, you stay unchanged

**Weighted Majority (Littlestone and Warmuth, 1998)**

**for** $t \leftarrow 1 \ldots N$ **do**
   $w_{1,i} \leftarrow 1$;
**for** $t \leftarrow 1 \ldots T$ **do**
   Receive $x_t$;
   $\hat{y}_t \leftarrow \mathbb{1}\left[\sum_{y_{t,i}=1} w_t \geq \sum_{y_{t,i}=0} w_t\right]$;
   Receive $y_t$;
   **if** $\hat{y}_t \neq y_t$ **then**
      **for** $t \leftarrow 1 \ldots N$ **do**
         **if** $\hat{y}_t \neq y_t$ **then**
            $w_{t+1,i} \leftarrow \beta w_{t,i}$;
         **else**
            $w_{t+1,i} \leftarrow w_{t,i}$
**return** $w_{T+1}$

- Weights for every expert
- Classifications in favor of side with higher total weight ($y \in \{0, 1\}$)
- Experts that are wrong get their weights decreased ($\beta \in [0, 1]$)
- If you're right, you stay unchanged

## Weighted Majority (Littlestone and Warmuth, 1998)

**for** $t \leftarrow 1 \dots N$ **do**
$\quad \mid \quad w_{1,i} \leftarrow 1;$
**for** $t \leftarrow 1 \dots T$ **do**
$\quad \mid \quad$ Receive $x_t$;
$\quad \mid \quad \hat{y}_t \leftarrow \mathbb{1}\left[\sum_{y_{t,i}=1} w_t \geq \sum_{y_{t,i}=0} w_t\right];$
$\quad \mid \quad$ Receive $y_t$;
$\quad \mid \quad$ **if** $\hat{y}_t \neq y_t$ **then**
$\quad \mid \quad \quad \mid \quad$ **for** $t \leftarrow 1 \dots N$ **do**
$\quad \mid \quad \quad \mid \quad \quad \mid \quad$ **if** $\hat{y}_t \neq y_t$ **then**
$\quad \mid \quad \quad \mid \quad \quad \mid \quad \quad \mid \quad w_{t+1,i} \leftarrow \beta w_{t,i};$
$\quad \mid \quad \quad \mid \quad \quad \mid \quad$ **else**
$\quad \mid \quad \quad \mid \quad \quad \mid \quad \quad \mid \quad w_{t+1,i} \leftarrow w_{t,i}$
**return** $w_{T+1}$

- Weights for every expert
- Classifications in favor of side with higher total weight ($y \in \{0, 1\}$)
- Experts that are wrong get their weights decreased ($\beta \in [0, 1]$)
- If you're right, you stay unchanged

## Weighted Majority (Littlestone and Warmuth, 1998)

**for** $t \leftarrow 1 \ldots N$ **do**
$\quad | \quad w_{1,i} \leftarrow 1;$
**for** $t \leftarrow 1 \ldots T$ **do**
$\quad |$ Receive $x_t$;
$\quad | \quad \hat{y}_t \leftarrow \mathbb{1}\left[\sum_{y_{t,i}=1} w_t \geq \sum_{y_{t,i}=0} w_t\right];$
$\quad |$ Receive $y_t$;
$\quad |$ **if** $\hat{y}_t \neq y_t$ **then**
$\quad | \quad |$ **for** $t \leftarrow 1 \ldots N$ **do**
$\quad | \quad | \quad |$ **if** $\hat{y}_t \neq y_t$ **then**
$\quad | \quad | \quad | \quad | \quad w_{t+1,i} \leftarrow \beta w_{t,i};$
$\quad | \quad | \quad |$ **else**
$\quad | \quad | \quad | \quad | \quad w_{t+1,i} \leftarrow w_{t,i}$
**return** $w_{T+1}$

- Weights for every expert
- Classifications in favor of side with higher total weight ($y \in \{0, 1\}$)
- Experts that are wrong get their weights decreased ($\beta \in [0, 1]$)
- If you're right, you stay unchanged

**Weighted Majority**

- Let $m_t$ be the number of mistakes made by WM until time $t$
- Let $m_t^*$ be the best expert's mistakes until time $t$

$$m_t \leq \frac{\log N + m_t^* \log \frac{1}{\beta}}{\log \frac{2}{1+\beta}} \tag{6}$$

- Thus, mistake bound is $O(\log N)$ plus the best expert
- Halving algorithm $\beta = 0$

**Proof: Potential Function**

- Potential function is the sum of all weights

$$\Phi_t \equiv \sum_i w_{t,i} \tag{7}$$

- We'll create sandwich of upper and lower bounds

**Proof: Potential Function**

- Potential function is the sum of all weights

$$\Phi_t \equiv \sum_i w_{t,i} \tag{7}$$

- We'll create sandwich of upper and lower bounds
- For any expert $i$, we have lower bound

$$\Phi_t \geq w_{t,i} = \beta^{m_{t,i}} \tag{8}$$

## Proof: Potential Function

- Potential function is the sum of all weights

$$\Phi_t \equiv \sum_i w_{t,i} \tag{7}$$

- We'll create sandwich of upper and lower bounds
- For any expert $i$, we have lower bound

$$\Phi_t \geq w_{t,i} = \beta^{m_{t,i}} \tag{8}$$

Weights are nonnegative, so $\sum_i w_{t,i} \geq w_{t,i}$

**Proof: Potential Function**

- Potential function is the sum of all weights

$$\Phi_t \equiv \sum_i w_{t,i} \tag{7}$$

- We'll create sandwich of upper and lower bounds
- For any expert $i$, we have lower bound

$$\Phi_t \geq w_{t,i} = \beta^{m_{t,i}} \tag{8}$$

Each error multiplicatively reduces weight by $\beta$

**Proof: Potential Function (Upper Bound)**

- If an algorithm makes an error at round $t$

$$\Phi_{t+1} \leq \frac{\Phi_t}{2} + \frac{\beta\Phi_t}{2} \qquad (9)$$

**Proof: Potential Function (Upper Bound)**

- If an algorithm makes an error at round $t$

$$\Phi_{t+1} \leq \frac{\Phi_t}{2} + \frac{\beta \Phi_t}{2} \qquad (9)$$

Half (at most) of the experts by weight were right

**Proof: Potential Function (Upper Bound)**

- If an algorithm makes an error at round $t$

$$\Phi_{t+1} \leq \frac{\Phi_t}{2} + \frac{\beta \Phi_t}{2} \tag{9}$$

Half (at least) of the experts by weight were wrong

**Proof: Potential Function (Upper Bound)**

- If an algorithm makes an error at round $t$

$$\Phi_{t+1} \leq \frac{\Phi_t}{2} + \frac{\beta \Phi_t}{2} = \left[\frac{1+\beta}{2}\right] \Phi_t \qquad (9)$$

**Proof: Potential Function (Upper Bound)**

- If an algorithm makes an error at round $t$

$$\Phi_{t+1} \leq \frac{\Phi_t}{2} + \frac{\beta \Phi_t}{2} = \left[\frac{1+\beta}{2}\right] \Phi_t \qquad (9)$$

- Initially potential function sums all weights, which start at 1

$$\Phi_1 = N \qquad (10)$$

**Proof: Potential Function (Upper Bound)**

- If an algorithm makes an error at round $t$

$$\Phi_{t+1} \leq \frac{\Phi_t}{2} + \frac{\beta \Phi_t}{2} = \left[\frac{1+\beta}{2}\right] \Phi_t \tag{9}$$

- Initially potential function sums all weights, which start at 1

$$\Phi_1 = N \tag{10}$$

- After $m_T$ mistakes after $T$ rounds

$$\Phi_T \leq \left[\frac{1+\beta}{2}\right]^{m_T} N \tag{11}$$

**Weighted Majority Proof**

- Put the two inequalities together, using the best expert

$$\beta^{m_T^*} \leq \Phi_T \leq \left[\frac{1+\beta}{2}\right]^{m_T} N \qquad (12)$$

**Weighted Majority Proof**

- Put the two inequalities together, using the best expert

$$\beta^{m_T^*} \leq \Phi_T \leq \left[\frac{1+\beta}{2}\right]^{m_T} N \tag{12}$$

- Take the log of both sides

$$m_T^* \log \beta \leq \log N + m_T \log \left[\frac{1+\beta}{2}\right] \tag{13}$$

**Weighted Majority Proof**

- Put the two inequalities together, using the best expert

$$\beta^{m_T^*} \leq \Phi_T \leq \left[\frac{1+\beta}{2}\right]^{m_T} N \tag{12}$$

- Take the log of both sides

$$m_T^* \log \beta \leq \log N + m_T \log \left[\frac{1+\beta}{2}\right] \tag{13}$$

- Solve for $m_T$

$$m_T \leq \frac{\log N + m_T^* \log \frac{1}{\beta}}{\log \left[\frac{2}{1+\beta}\right]} \tag{14}$$

## Weighted Majority Recap

- Simple algorithm
- No harsh assumptions
- Depends on best learner
- Downside: Takes a long time to do well in worst case (but okay in practice)
- Solution: Randomization

## Plan

Experts

Perceptron Algorithm

## Perceptron Algorithm

- Online algorithm for classification
- Very similar to logistic regression (but 0/1 loss)
- But what can we prove?

**Perceptron Algorithm**

$\vec{w}_1 \leftarrow \vec{0}$;
**for** $t \leftarrow 1 \ldots T$ **do**
    Receive $x_t$;
    $\hat{y}_t \leftarrow \mathsf{sgn}(\vec{w}_t \cdot \vec{x}_t)$;
    Receive $y_t$;
    **if** $\hat{y}_t \neq y_t$ **then**
        $\vec{w}_{t+1} \leftarrow \vec{w}_t + y_t \vec{x}_t$;
    **else**
        $\vec{w}_{t+1} \leftarrow w_t$;
**return** $w_{T+1}$

        **Algorithm 2:** Perceptron Algorithm (Rosenblatt, 1958)
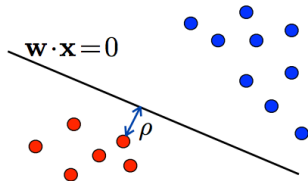
**Objective Function**

- Optimizes

$$\frac{1}{T} \sum_t \max\left(0, -y_t(\vec{w} \cdot x_t)\right) \tag{15}$$
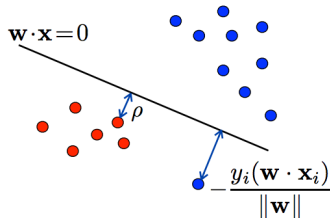
- Convex but not differentiable

## Margin and Errors



- If there's a good margin $\rho$, you'll converge quickly

## Margin and Errors



- If there's a good margin $\rho$, you'll converge quickly
- Whenever you se an error, you move the classifier to get it right
- Convergence only possible if data are separable

**How many errors does Perceptron make?**

- If your data are in a $R$ ball and there is a margin

$$\rho \leq \frac{y_t(\vec{v} \cdot \vec{x}_t)}{||v||} \tag{16}$$

 for some $\vec{v}$, then the number of mistakes is bounded by $R^2/\rho^2$

- The places where you make an error are support vectors
- Convergence can be slow for small margins

**Why study Perceptron?**

- Simple algorithm
- Bound independent of dimension and tight
- Foundation of deep learning
- Proof techniques helped usher in SVMs
- Generalizes to structured prediction