

Online Classification — Part 2

Source:

Online Learning and Online Convex Optimization, by
Shai Shalev-Shwartz, Found. and Trends in ML, Vol. 4,
No. 2, 2011 (Chapters 1, 3).

Setting: feature space X . Classifiers: \mathcal{H}
Label space Y

$h \in \mathcal{H} : h: X \mapsto Y$ i.e., $x \in X$
 $y \in Y$

$h(x) = y$, a classifier.

Typically, $X = \mathbb{R}^d$, $Y = \{0, 1\}$ (or $\{-1, +1\}$)

→ binary classification

Online problem — evolves in discrete time steps,
 $t = 1, 2, \dots, T$.

Infinite Hypotheses : Unrealizable Case

($|H| = \infty$, but $L\dim(H) < \infty$).

Unrealizable case \Rightarrow Player takes randomized actions.

Setting: $|H|$ potentially infinite; $p_t \in [0, 1]$,
 $y_t \in \{0, 1\}$. Goal is to minimize expected regret
w.r.t. best hypothesis in hindsight. $L\dim(H) < \infty$

We will show that $R_T \leq O(\sqrt{L\dim(H) \cdot \ln(T)} \cdot T)$.

Idea: (to deal with $|H| = \infty$, but $L\dim(H) < \infty$)

① Previously with finite $|H|$, each hypothesis was an expert. Now, construct a finite number of experts such that:

For each $h \in H$, and any sequence of features $\{x_1, x_2, \dots, x_T\}$, \exists an expert s.t. the experts prediction \equiv prediction from $h(\cdot)$.
Assume $T \geq L\dim(H)$.

Specifically, for each $L \leq L\dim(H)$, and any sequence of time indices: $1 \leq i_1 < i_2 < i_3 \dots < i_L \leq T$,

we will construct an expert.

- ② Then use weighted majority (i.e., Exp4 will full information over this group of experts).
-

Operation of an Expert

Expert(L, i_1, i_2, \dots, i_L)

Input: $\mathcal{H}, L, i_1, i_2, \dots, i_L$.

Initial: $V_1 = \mathcal{H}$.

For each $t=1, 2, \dots, T$:

- Observe x_t

- Define: $V_t^{(0)} = \{h \in V_t : h(x_t) = 0\}$
 $V_t^{(1)} = \{h \in V_t : h(x_t) = 1\}$

$$\tilde{y}_t = \underset{r \in \{0, 1\}}{\operatorname{argmax}} L \dim(V_t^{(r)})$$

(if $V_t^{(0)} = V_t^{(1)}$, then set $\tilde{y}_t = 0$).

- For $t \in \{i_1, i_2, \dots, i_L\}$, set $\hat{y}_t = |1 - \tilde{y}_t|$

(i.e., the
opposite label)

For $t \notin \{i_1, i_2, \dots, i_L\}$, set $\hat{y}_t = \tilde{y}_t$.

- Update $v_{t+1} = v_t^{(\hat{y}_t)}$

Intuition: Recall, that for each x_t , π_{SOA} (i.e., the Standard Optimal Algorithm) predicts the label corresponding to the larger Ldim of the remaining set of hypothesis (ie instead of majority, go with larger Ldim).

This expert above ASSUMES SOA makes a mistake at times $\{i_1, i_2, \dots, i_L\}$. Thus, it flips the label predicted by SOA at those times, but follows SOA at all other times.

Note: # of experts is finite. Specifically, the number of experts E is given by

$$E = \sum_{L=0}^{\text{Ldim}(\mathcal{H})} \underbrace{\binom{T}{L}}_{\text{choose } L \text{ time-slots to}} \leq \left(\frac{e^T}{\text{Ldim}(\mathcal{H})} \right)^{\text{Ldim}(\mathcal{H})}$$

flip SOA's recommendation \rightarrow each such choice defines an expert.

Lemma (3.7 in monograph): \mathcal{H} a hypothesis class with $\text{Ldim}(\mathcal{H}) < \infty$. $\{x_1, x_2, \dots, x_T\}$ a sequence of features. For any $h \in \mathcal{H}$, \exists an expert $\text{Expert}(L, i_1, i_2, \dots, i_L)$ s.t. for this sequence of features, this expert predicts $y_t = h(x_t)$ (i.e., the $\text{Expert}(L, i_1, \dots, i_L)$ and h agree on labels for all times $t = 1, 2, \dots, T$.

Proof:

- Fix $h \in \mathcal{H}, \{x_1, x_2, \dots, x_T\}$.
- Consider the sequence $\{(x_1, h(x_1)), (x_2, h(x_2)), \dots, (x_T, h(x_T))\}$. Run SOA with this sequence. Suppose that SOA makes L mistakes, at times i_1, i_2, \dots, i_L . (Recall that $L \leq \text{Ldim}(\mathcal{H})$ for the SOA algorithm). Thus, SOA predicts:
 $\{(x_1, h(x_1)), \dots, (x_{i_1}, 1-h(x_{i_1})), \dots, (x_k, h(x_k)), \dots$

$$\dots, (x_{i_L}, \textcolor{red}{1-h(x_{i_L})}), \dots, (x_T, \textcolor{green}{h(x_T)}) \}$$

Here, green times are correct labels by SOA (and thus match h), and red times correspond to incorrect labels by SOA (and thus flip h).

- Pick the $\text{Expert}(L, i_1, i_2, \dots, i_L)$. By construction, this expert flips SOA at red times, and thus matches h at all times.

□

Corollary (3.8): Let $\{(x_i, y_i), \dots, (x_T, y_T)\}$ be a labeled sequence, \mathcal{H} a hypothesis class with $\text{Ldim}(\mathcal{H}) < \infty$. Then, $\exists \text{ Expert}(L, i_1, i_2, \dots, i_L)$ that makes no more mistakes than the best $h \in \mathcal{H}$ for this labeled sequence, i.e;

of mistakes by $\text{Expert}(L, i_1, i_2, \dots, i_L)$

$$\leq \min_{h \in \mathcal{H}} \sum_{t=1}^T |h(x_t) - y_t|.$$

Proof: Immediately follows from Lemma above

□

Theorem (§-6 in monograph): For any hypothesis class \mathcal{H}
 s.t. $Ldim(\mathcal{H}) < \infty$, \exists algorithm π s.t. for
 any $h \in \mathcal{H}$, and labeled sequence $\{(x_1, y_1), \dots, (x_T, y_T)\}$,

$$\sum_{t=1}^T |p_t - y_t| = \sum_{t=1}^T |h(x_t) - y_t|$$

$$\leq O\left(\sqrt{Ldim(\mathcal{H}) \cdot (\ln T) \cdot T}\right).$$

Proof: Use weighted majority algorithm with the finite set of experts constructed above.

$$\text{Note } E = \left(\frac{e^T}{Ldim(\mathcal{H})} \right)^{Ldim(\mathcal{H})}$$

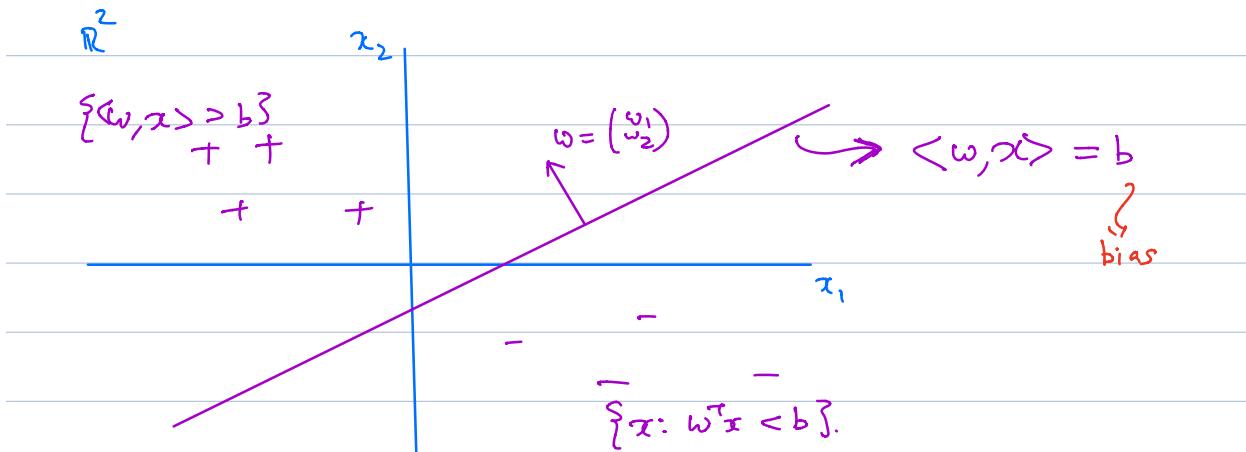
Corollary 3.8 + Weighted majority proof \Rightarrow Result. \square

Binary Classification over $x_t \in \mathbb{R}^d$ where $Ldim(\mathcal{H}) = \infty$

Perceptron and Winnow Algorithms

Setting: $y \in \{-1, +1\}$ (instead of $\{0, 1\}$ as before).

$$X = \mathbb{R}^d, D: p_t \in \{-1, +1\}.$$



Note: We will set $b=0$ in the sequel. Bias can easily be addressed by considering the above in $(d+1)$ dimensions with $x_{d+1}=1$, and $w_{d+1}=-b$.

Setting: In round $t=1, 2, \dots, T$, player observes $x_t \in \mathbb{R}^d$.

- Player has $w_t \in \mathbb{R}^d$ (weight vector), and updates it sequentially. Specifically, player predicts

$$p_t = \text{Sign}(\langle w_t, x_t \rangle) \in \{-1, +1\}$$

- Observes label $y_t \in \{-1, +1\}$

$$\cdot \ell(p_t, y_t) = \chi_{\{p_t \neq y_t\}} = |p_t - y_t|$$

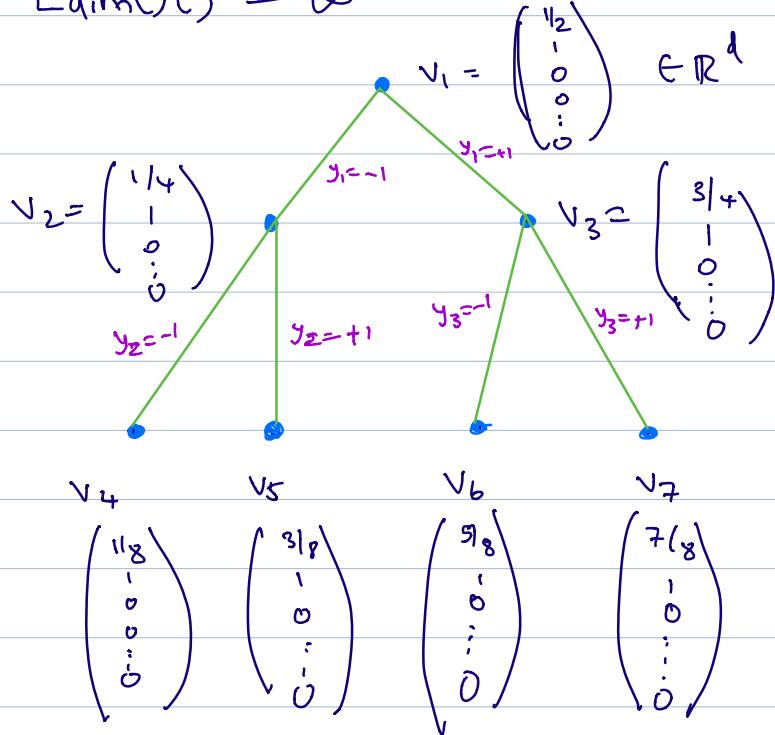
In this case, since $w_t \in \mathbb{R}^d$,

$$\mathcal{H} = \{w \in \mathbb{R}^d : p_t = \text{Sign}(\langle w, x \rangle)\}$$

i.e., $h \in \mathcal{H}$, with $h_w(x) = \text{sign}(\langle w, x \rangle)$.
 (set of half spaces in \mathbb{R}^d).

Claim: $\text{Ldim}(\mathcal{H}) = \infty$

Proof:



Now consider $\tilde{\mathcal{H}} \subseteq \mathcal{H}$, where $h \in \tilde{\mathcal{H}}$ parameterized by $w = (-1, a, 0, \dots, 0)$, $a \in [0, 1]$. Observe that the infinite tree above can be shattered by the hypotheses in $\tilde{\mathcal{H}}$. Thus, $\text{Ldim}(\mathcal{H}) = \infty$ ④.

Note: Convexification Using Surrogate Loss Functions

Ref: Sec. 2.1.2, monograph + Sec. 3.3.1 (Perception).

Above setting has $\text{Ldim}(\mathcal{H}) = \infty$, so it is not immediately clear how to use randomization. This note discusses an alternate approach, namely, surrogate convex loss functions (from sec. 2.1.2). The high-level idea:

- ① Observe that the loss function here is $\chi_{\{p_t \neq y_t\}}$
 $= |p_t - y_t|$, where $p_t = \text{sign}(\langle w_t, x_t \rangle)$, i.e.,
the loss function at each time:

$$l_t(w, (x, y)) = \chi_{\{y \langle w, x \rangle \leq 0\}}$$

$\uparrow \quad \uparrow$
iff signs are opposite, the
product is $-w$, and $l_t(\cdot) = 1$.

(2) Upper bound $l_t(\cdot)$ by a surrogate convex loss function, i.e., construct a sequence of functions $\{f_t(\omega), t=1, 2, \dots\}$ s.t.

- $f_t(\omega)$ is convex, and
- $\forall \omega, t, f_t(\omega) \geq l(\omega, (x_t, y_t))$

(3) Using online convex optimization (in full feedback setting), show that online gradient descent with $f_t(\cdot)$ has low regret.

(4) Use upper bound property in (2) to derive a mistake bound for the $l_t(\cdot)$ loss function.

The Perceptron Algorithm

$$\text{Loss } l_t(\omega, (x, y)) = \chi_{\{y < \omega \cdot x \leq 0\}}$$

actual loss

$$\text{Define } \tilde{f}_t(\omega) = [1 - y_t \langle \omega, x_t \rangle]_+,$$

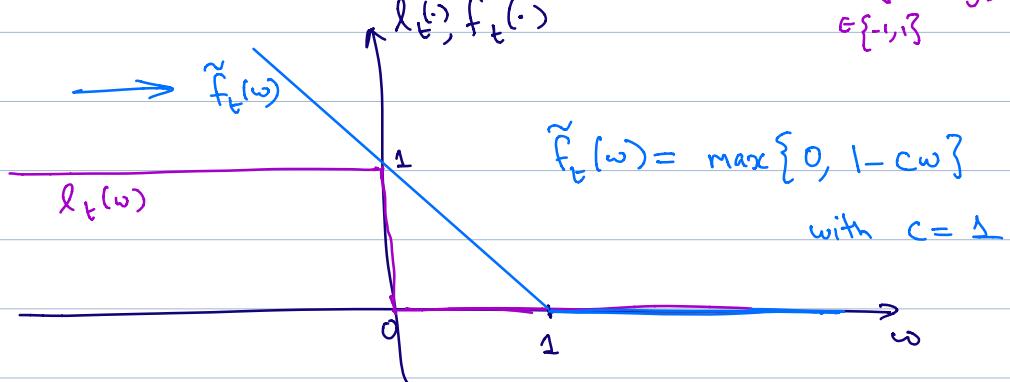
where $[r]_+ = \max\{r, 0\}$

hinge loss

e.g. if $\omega \in \mathbb{R}$,

(for some fixed $(x, y) \in \mathbb{R}^2$ with $xy > 0$)
 $\in \{-1, 1\}$

hinge loss.



Original :
loss fn

$c\omega \rightarrow \{-1, +1\}$
replaced by

$(c\omega \rightarrow \text{sign}(y \langle \omega, x \rangle))$

Observe: ① $\tilde{f}_t(\omega)$ convex in ω .

$$\begin{aligned} \textcircled{2} \quad \tilde{f}_t(\omega) &= [1 - y_t \langle \omega, x_t \rangle]_+ \geq x_{\{y_t \langle \omega, x_t \rangle \leq 0\}} \\ &= l_t(\omega, (x_t, y_t)). \end{aligned}$$

Now, define

$$f_t(\omega) = \begin{cases} \tilde{f}_t(\omega) & \text{if mistake made at time } t \\ 0 & \text{if no mistake} \end{cases}$$

Note that $f_t(\omega)$ satisfies ① $f_t(\omega)$ convex in ω for each t , and ② $f_t(\omega) \geq \ell_t(\omega, (x_t, y_t))$. $\forall \omega, t$.

$\{f_t(\omega)\}$ is the sequence of convex surrogate loss functions.

Algorithm — Online Gradient Descent on ω_t using surrogate functions.

$$\omega_1 = 0$$

$$(\omega_i \in \mathbb{R}^d)$$

$$\omega_{t+1} = \omega_t - \gamma z_t, \quad z_t \in \partial f_t(\omega_t).$$

\hookrightarrow subgradient.

Here: no mistake $\Leftrightarrow y_t \langle \omega_t, x_t \rangle > 0$
 $\Leftrightarrow f_t(\cdot) = 0$: Set $z_t = 0$

mistake: $z_t = -y_t x_t \in \partial f_t(\omega_t)$.

$$\therefore \omega_{t+1} = \begin{cases} \omega_t + \gamma y_t x_t & \text{if } y_t \langle \omega_t, x_t \rangle \leq 0 \\ \omega_t & \text{if } y_t \langle \omega_t, x_t \rangle > 0 \end{cases}$$

Let $W \subseteq \{1, 2, 3, \dots, T\}$ be the time slots where mistakes were made, i.e.,

$$s \in W \Leftrightarrow \underbrace{\text{sign}(\langle w_t, x_t \rangle)}_{P_t} \neq y_t$$

Observe also that:

$$\begin{aligned} P_t &= \text{sign}(\langle w_t, x_t \rangle) = \text{sign}\left(\eta \sum_{\substack{i \in W: \\ i < t}} y_i \langle x_i, x_t \rangle\right) \xrightarrow{\eta > 0} \\ &= \text{sign}\left(\sum_{\{i \in W: i < t\}} y_i \langle x_i, x_t \rangle\right) \end{aligned}$$

i.e., invariant to η .

\therefore Perception Algorithm:

(single neuron with threshold activation)

$$P_t = \text{sign}(\langle w_t, x_t \rangle)$$

$$w_{t+1} = \begin{cases} w_t + y_t x_t & \text{if } y_t \langle w_t, x_t \rangle \leq 0 \\ w_t & \text{if } y_t \langle w_t, x_t \rangle > 0 \end{cases}$$

$\eta > 0$ is immaterial. Choose this appropriately for analysis bounds. Algo. does not change.

Analysis :

Use OGD (online gradient descent) bound for convex functions. This is analogous to the bounds we studied for online linear optimization, see Chap 2 of monograph for details. Specifically, it can be shown that:

$$\sum_{t=1}^T f_t(w_t) - \sum_{t=1}^T f_t(u) \geq |w| \leq \frac{1}{2\eta} \|u\|_2^2 + \frac{\eta}{2} |w| R^2,$$

by property of surrogate loss.

where $R = \max_{1 \leq t \leq T} \|x_t\|_2$.

$$\therefore \text{Choose } \eta = \|u\|_2 / R \sqrt{|w|}$$

$$\Rightarrow |w| - \sum_{t=1}^T f_t(u) - R \|u\|_2 \sqrt{|w|} \leq 0$$



Theorem (3.9 in monograph): With sequence $\{(x_1, y_1), \dots, (x_T, y_T)\}$,
 and $R = \max_{1 \leq t \leq T} \|x_t\|_2$, and $|W| = \# \text{ of}$

mistakes until time T . Then for any $u \in \mathbb{R}^d$,

$$|W| \leq \sum_{t=1}^T f_t(u) + R \|u\|_2 \sqrt{\sum_{t=1}^T f_t(u)} + R^2 \|u\|_2^2$$

Further, if $\exists u$ st. $y_t \langle u, x_t \rangle \geq 1 \forall t$ (also
 called separability with large margin)

$$|W| \leq R^2 \|u\|^2.$$

Proof: follows from $\textcircled{1}$. Also $x - b\sqrt{x} - c \leq 0$

$$\Downarrow$$

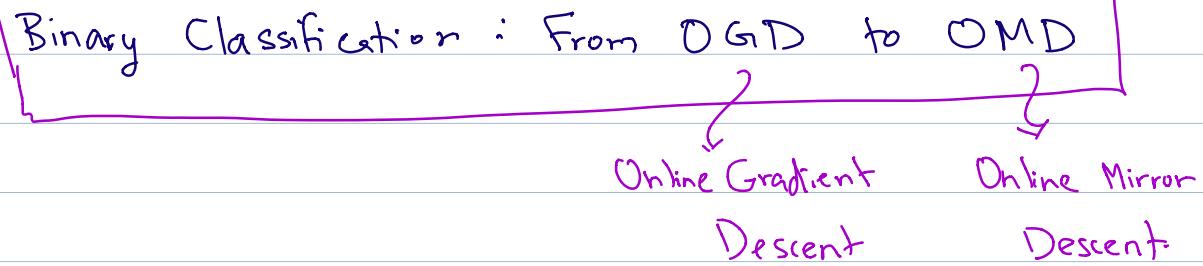
$$x \leq c + b^2 + b\sqrt{c}.$$

Last claim using algebra from $\textcircled{2}$.



Note: It can happen that $\exists u$ with no mistakes
 (but does not satisfy large margin condition), but the
 perceptron algorithm can make many mistakes, because
 regret wrt $\sum_t f_t(u)$ (and not the actual loss
 function).

Thus, bound is useful only when $\exists u$ s.t. $\sum_t f_t(u)$ is small.



Previously: By changing regularizer and generalizing gradient descent to mirror descent, substantial improvement in regret in certain settings (mirror descent over simplex with KL divergence as regularizer).

This section: Similar plan: Interpret Perceptron as Online Gradient Descent. Winnow (algo. to be described below) is an Online Mirror Descent algo.

Note : These algorithms have a much longer history. Precursor to Exp3 (and boosting in the case of full information) is the Weighted Majority Algorithm (Littlestone and Warmuth, 1994). Winnow predates even that (late 80s). Perceptron goes back to (Rosenblatt, 1958).

The unaltered problem: $\mathcal{X} = \{0, 1\}^d$, i.e.,

$x \in \mathcal{X}$ is a binary feature vector. We are in a setting where we know that only $k \leq d$ of these coordinates are useful; rest are spurious.

If any of the useful coordinates equals ' 1 ', the label is ' 1 '.

Thus, we need to find a $w \in \{0, 1\}^d$ (a binary weight vector) that is k -sparse, and the associated classifier is given by:

$$h(x) = \text{sign}(2 \langle w, x \rangle - 1).$$

$$\begin{array}{l} \left\{ \begin{array}{l} \langle w, x \rangle = 0 \\ \Rightarrow h(x) = -1 \end{array} \right. \\ \left\{ \begin{array}{l} \langle w, x \rangle \geq 1 \\ \Rightarrow h(x) = +1 \end{array} \right. \end{array} \leftarrow \begin{cases} +1 & \langle w, x \rangle > 0 \\ -1 & \langle w, x \rangle \leq 0 \end{cases}$$

$$\mathcal{H} = \{x \mapsto \text{sign}(2 \langle w, x \rangle - 1), \|w\|_1 = k, w \in \{0, 1\}^d\}$$

Convex Relaxation + Surrogate Hinge Loss

1. Domain Relaxation: $w \in \mathbb{R}_+^d = \{x \in \mathbb{R}^d : x_i \geq 0\}$.
(positive orthant).

2. Surrogate Loss Function: Exact loss function is

$$l(\omega, (x, y)) = \chi_{\{y(2\langle \omega, x \rangle - 1) \leq 0\}}.$$

Let $f_t(\omega) = \begin{cases} 0 & \text{if no mistake by algo. at time } t \text{ (i.e., } y = \text{sign}(2\langle \omega, x \rangle - 1)\text{),} \\ [1 - y_t(2\langle \omega, x_t \rangle - 1)]_+ & \text{if mistake at time } t. \end{cases}$

hinge loss.

Observe:

- ① $f_t(\omega)$ is convex in ω for $\forall t=1, 2, \dots, T$
 - ② $f_t(\omega) \geq l(\omega, (x_t, y_t))$
-

Goal

Develop regret bound on number of mistakes with algo. over time horizon T .

Caveat: Algo. uses weight vector $\omega \in \mathbb{R}_+^d$. However regret with respect to best weight vector u s.t. $\|u\|_2 = k$, $u \in \{0, 1\}^d$. In other words, the algo. has more "power" than the genie's best in hindsight.

If realizable, assume $\exists u \in \{0, 1\}^d$ s.t. $\|u\|_1 = k$
 and $f_t(u) = 0 \ \forall t$.

Base line Algo: Perceptron

$$x \rightarrow \phi(x) = \begin{pmatrix} 2x \\ -1 \end{pmatrix} \in \mathbb{R}^{d+1}$$

$\underbrace{x}_{\in \{0, 1\}^d}$ $\underbrace{\phi(x)}_{\text{modified feature vector}}$

Use Perceptron with $P_t = \text{sign}(\langle w, \phi(x) \rangle)$.

Then, with realizability assumption $\left(\exists u \in \{0, 1\}^d \text{ s.t. } \|u\|_1 = k, f_t(u) = 0 \right)$,

$$\begin{aligned} \text{we have: } R^2 \|u\|_2^2 &= (\max \|x\|_2)^2 (\|u\|_2^2) \\ &= (4d+1) k \end{aligned}$$

\therefore From Thm 3.9 (pp. 16 in these notes), we have

$$\|w\| \leq R^2 \|u\|_2 = (4d+1) k.$$

We will see with Winnow (a exponential weighting

algorithm) that $|W| \leq 8k \ln(d)$, an thus an exponential improvement due to OMD instead of OGD.

Winnow Algorithm \equiv Unnormalized Exp Weighting.

Input: $\eta > 0$ (set $\eta = 1/4$ for bounds).

Initial: $w_1 = \left(\frac{1}{d} \ \frac{1}{d} \ \dots \ \frac{1}{d}\right)^T$, $w \in \mathbb{R}_+^d$.

For each $t=1, 2, \dots, T$:

- Observe $x_t \in \{0, 1\}^d$
- Predict $p_t = \text{sign}(2 \langle w_t, x_t \rangle - 1)$
- If $y_t (2 \langle w_t, x_t \rangle - 1) > 0$, → correct prediction

Set $w_{t+1} = w_t$ (no change)

else if $y_t (2 \langle w_t, x_t \rangle - 1) \leq 0$,

$$w_{t+1}(i) = w_t(i) e^{-\eta \underbrace{2y_t x_t(i)}_{\in \{-1, 1\}}}, \quad i=1, 2, \dots, d.$$

Theorem (3.10 in monograph) $\{(x_1, y_1), \dots, (x_T, y_T)\}$

labeled samples, and $\eta \leq 1/2$. Let $\mathcal{W} \subseteq \{1, 2, \dots, T\}$ time steps where mistakes occur with Winnow. Let

$$f_t(w) = \chi_{\{t \in \mathcal{W}\}} \cdot [1 - y_t(2 \langle w_t, x_t \rangle - 1)]_+.$$

Then, for any $u \in \{0, 1\}^d$, we have:

$$|\mathcal{W}| \leq \sum_{t=1}^T f_t(w_t) \leq \frac{1}{(1-2\eta)} \left(\sum_{t=1}^T f_t(u) + \frac{k \ln(d)}{\eta} \right)$$

Further, if $\nexists u$ s.t. $y_t(2 \langle u, x_t \rangle - 1) \geq 1 \forall t$,

then with $\eta = 1/4$:

$$|\mathcal{W}| \leq 8k \ln(d).$$

Prof: Uses the OMD bound for online convex optimization (mirror descent with KL regularizer on positive orthant, read Thm 2.23 in monograph).

$$\text{Define } z_t = \begin{cases} 2y_t x_t & \text{if } t \in \mathcal{W} \\ 0 & \text{if } t \notin \mathcal{W} \end{cases}$$

Then from Thm 2.23 (please read monograph); this is similar to several proj s we have done with KL regularizer in OCO), we have: $u \in \{0, 1\}^d$, $\|u\|_1 = k$,

$$\sum_{t=1}^T \langle w_t - u, z_t \rangle \leq 1 + \frac{\sum_{i=1}^d u_i \ln(\frac{d u_i}{e})}{2}$$

$$+ \gamma \sum_{t=1}^T w_t(i) \cdot z_t(i)^2$$

Further, $\|u\|_1 = k$, $u \in \{0, 1\}^d \Rightarrow$

$$1 + \sum_{i=1}^d u_i \ln\left(\frac{d u_i}{e}\right) \leq 1 + \ln d.$$

Also, $z_t \in \partial f_t|_{w_t}$, and f_t convex. Thus,

$$\begin{aligned} \sum_{t=1}^T (f_t(w_t) - f_t(u)) &\leq \sum_{t=1}^T \langle w_t - u, z_t \rangle \\ &\leq \frac{k \ln d}{\gamma} + \gamma \sum_{t=1}^T \sum_{i=1}^d w_t(i) z_t(i)^2 \end{aligned}$$

Next, $\forall t = 1, 2, \dots, T$, $\sum_{i=1}^d w_t(i) z_t(i)^2 \leq 2 f_t(w_t)$.

$$\hookrightarrow \textcircled{*}$$

Case(i): $t \notin \mathcal{W} \Rightarrow f_t(w_t) = 0$, $z_t(i) = 0$. $\forall i = 1, 2, \dots, d$.
 $\Rightarrow \textcircled{*}$ holds.

Case(ii): $t \in \mathcal{W}$, and $y_t = 1$. Then LHS in $\textcircled{*}$ is

$$\sum_{i=1}^d w_t(i) z_t(i)^2 = 4 \langle w_t, x_t \rangle$$

\downarrow

$$2 y_t x_t(i) \underset{y_t=1}{\hookrightarrow} \in \{0, 1\}$$

} $\Rightarrow 4 \langle w_t, x_t \rangle \leq 2$

$$t \in \mathcal{W} \Rightarrow 2 \langle w_t, x_t \rangle - 1 \leq 0$$

Also RHS in $\textcircled{*}$ = $2 f_t(w_t) \geq 2 \geq \text{LHS}$.

$\underbrace{\geq 1}_{\text{when } t \in \mathcal{W}}$

Case(iii): $t \in \mathcal{W}$, $y_t = -1$. LHS in $\textcircled{*}$ = $4 \langle w_t, x_t \rangle$.

\therefore RHS in $\textcircled{*}$ = $2 f_t(w_t) = 2 \left(2 \langle w_t, x_t \rangle \right) = \text{RHS}$.

\therefore mistake occurs

$$= [1 - (-1)2 \langle w_t, x_t \rangle - 1] \times = 2 \langle w_t, x_t \rangle$$

$$\therefore \sum_{t=1}^T f_t(w_t) \leq \frac{1}{1-2\eta} \left(\sum_{t=1}^T f_t(u) + \frac{k \ln(d)}{\eta} \right)$$

(4)

Please read monograph chap. 2 for the OGD, OMD discussion with convex functions. This is analogous to the discussions for the linear case we studied. However, ideas immediately apply because of linear approximation for a convex fn.

f_t convex \Rightarrow for any u ,

$$f_t(w_t) - f_t(u) \leq \langle w_t - u, z_t \rangle$$

