

Lecture 17 SVM (Part II) and Online Learning

Lei Li, Yu-Xiang Wang

(some slides from my convex optimization class,
originally taught by Ryan Tibshirani in CMU)

Recap: Support Vector Machines

Given $y \in \{-1, 1\}^n$, $X \in \mathbb{R}^{n \times p}$ having rows x_1, \dots, x_n , recall the **support vector machine** or SVM problem:

$$\min_{\beta, \beta_0, \xi} \quad \frac{1}{2} \|\beta\|_2^2 + C \sum_{i=1}^n \xi_i$$

subject to $\xi_i \geq 0, i = 1, \dots, n$

$$y_i (\underline{x_i^T \beta + \beta_0}) \geq 1 - \underline{\xi_i}, i = 1, \dots, n$$

This is a quadratic program

Recap: Lagrange dual problem

Given a minimization problem

$$\begin{aligned} \min_x \quad & f(x) \\ \text{subject to} \quad & h_i(x) \leq 0, \quad i = 1, \dots, m \\ & \ell_j(x) = 0, \quad j = 1, \dots, r \end{aligned}$$

we defined the **Lagrangian**:

$$L(x, u, v) = \underbrace{f(x)}_{\text{original function}} + \sum_{i=1}^m \underbrace{u_i h_i(x)}_{\text{term involving } h_i(x)} + \sum_{j=1}^r \underbrace{v_j \ell_j(x)}_{\text{term involving } \ell_j(x)}$$

and **Lagrange dual function**:

$$\underbrace{g(u, v)}_{\text{dual function}} = \min_x L(x, u, v)$$

Recap: Lagrange dual problem

The subsequent **dual problem** is:

$$\begin{aligned} \max_{u,v} \quad & \underline{g(u, v)} \\ \text{subject to} \quad & u \geq 0 \end{aligned}$$

Important properties:

- Dual problem is always convex, i.e., g is always concave (even if primal problem is not convex)
- The primal and dual optimal values, f^* and g^* , always satisfy weak duality: $f^* \geq g^*$
- Slater's condition: for convex primal, if there is an x such that

$$h_1(x) < 0, \dots, h_m(x) < 0 \quad \text{and} \quad \ell_1(x) = 0, \dots, \ell_r(x) = 0$$

then **strong duality** holds: $f^* = g^*$. Can be further refined to strict inequalities over the nonaffine h_i , $i = 1, \dots, m$

Recap: Deriving the dual of SVM

Introducing dual variables $v, w \geq 0$, we form the Lagrangian:

$$\underbrace{L(\beta, \beta_0, \xi, v, w)}_{X} = \frac{1}{2} \|\beta\|_2^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n v_i \xi_i + \sum_{i=1}^n w_i (1 - \xi_i - y_i (\beta^T x_i + \beta_0))$$

$$0 = \nabla_{\beta} L = \beta + \sum_{i=1}^n w_i (-y_i) \vec{x}_i$$

$$\nabla_{\beta_0} L = \sum_{i=1}^n -w_i y_i$$

$$\nabla_{\xi_i} L = C - v_i - w_i$$

$$g(w, v) = \begin{cases} -\infty & \text{if } \text{not satisfied} \\ -\infty & \text{next slide} \end{cases}$$

$$\xi_i (C - v_i - w_i)$$

$$\left. \begin{array}{l} \sum_{i=1}^n w_i y_i = 0 \\ C - v_i - w_i = 0 \end{array} \right\} \text{Constraints on } w, v$$

Recap: Dual SVM

Minimizing over β, β_0, ξ gives Lagrange dual function:

$$g(v, w) = \begin{cases} -\frac{1}{2}w^T \tilde{X} \tilde{X}^T w + 1^T w & \text{if } w = C1 - v, w^T y = 0 \\ -\infty & \text{otherwise} \end{cases}$$

where $\tilde{X} = \text{diag}(y)X$. Thus SVM dual problem, eliminating slack variable v , becomes

$$\begin{aligned} \max_w \quad & -\frac{1}{2}w^T \tilde{X} \tilde{X}^T w + 1^T w \\ \text{subject to} \quad & 0 \leq w \leq C1, w^T y = 0 \end{aligned} \quad |$$

Check: Slater's condition is satisfied, and we have strong duality.
Further, from study of SVMs, might recall that at optimality

$$\beta = \tilde{X}^T w \quad \text{Score}(x) = \beta^T x + \beta_0$$

This is not a coincidence, as we'll later via the KKT conditions

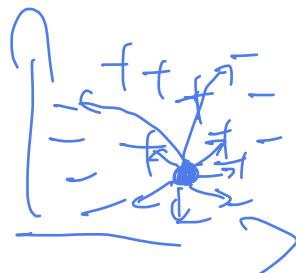
“Kernel trick” in SVM

$$\tilde{X} = \text{diag}(y) X$$

(y_1, y_2, \dots, y_d) | \tilde{X}

- The dual SVM depends only on inner products

$$\begin{cases} Ax \\ k(x, x) > 0 \\ A^T A = 0 \\ x^T A x \geq 0 \forall x \end{cases}$$



$$\max_w -\frac{1}{2} w^T \tilde{X} \tilde{X}^T w + 1^T w$$

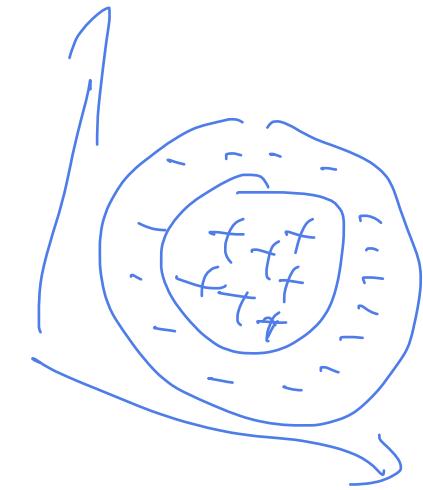
subject to $0 \leq w \leq C, w^T y = 0$

$$w^T \tilde{X} \tilde{X}^T w = \sum_{i=1}^n \sum_{j=1}^n w_i w_j y_i y_j x_i^T x_j$$

$$\phi(x_i) \in \mathcal{H}$$

example: $\phi(x_i) = \exp(-\|x_i - \cdot\|^2)$

$$k(x_i, x_j) = \langle \phi(x_i), \phi(x_j) \rangle = \exp(-\|x_i - x_j\|^2)$$



- How to make predictions?

$$\text{Sign}\left(\tilde{x}^T \tilde{X}^T w + \beta_0\right) = \text{Sign}\left(\underbrace{\phi(x)^T}_{\tilde{x}^T} \cdot \underbrace{[\phi(x_1) \cdot y_1, \phi(x_2) \cdot y_2, \dots, \phi(x_n) \cdot y_n]}_{\tilde{X}^T} - w^T\right) + \beta_0$$

$$= \text{Sign}\left(\underbrace{[k(x, x_1), k(x, x_2), \dots, k(x, x_n)]}_{\tilde{x}^T} \cdot \underbrace{[y_1, y_2, \dots, y_n]}_{\tilde{X}^T} - w^T\right) + \beta_0$$

$$= \left(\sum_{i=1}^n w_i \cdot k(x, x_i) + \beta_0\right) + \beta_0$$

This lecture

- KKT conditions
 - SVM as an example
- Online Learning

Optimality conditions: the conditions that characterizes the optimal solutions

- What you learned in high school

$$\min_{x \in \mathbb{R}} x^2 - 4x + 9 = f(x)$$

$x \in [0, 1]$ $f(x^*) = 0$ if $f''(x) \geq 0$

- Slight generalization: For convex and differentiable objective function

$$\min_{\underbrace{x \in \mathbb{R}^d}_{\text{---}}} f(x) \quad \nabla f(x) = 0 \quad \nabla^2 f(x) \succeq 0$$

Does not handle non-differentiable functions, does not handle constraints.

Handling constraints with first-order optimality conditions

For a convex problem

$$\min_x f(x) \text{ subject to } x \in C$$

$$\forall y \in C \quad f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle$$

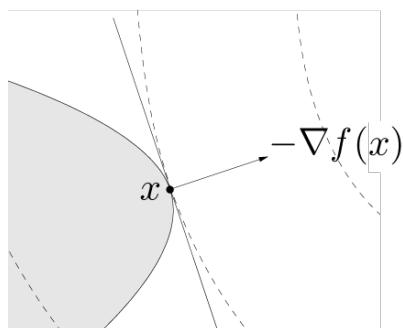
$\geq f(x)$
 $x^* = \arg \min_{x \in C} f(x)$

and differentiable f , a feasible point x is optimal if and only if

$$\nabla f(x)^T (y - x) \geq 0 \quad \text{for all } y \in C$$

$$f(x) \leq f(x^*) \\ \Rightarrow x^* \text{ is optm.}$$

This is called the **first-order condition for optimality**



In words: all feasible directions from x are aligned with gradient $\nabla f(x)$

Important special case: if $C = \mathbb{R}^n$ (unconstrained optimization), then optimality condition reduces to familiar $\nabla f(x) = 0$

Handling non-differentiable functions with “subgradient”

Recall that for convex and differentiable f ,

$$f(y) \geq f(x) + \nabla f(x)^T(y - x) \quad \text{for all } x, y$$

I.e., linear approximation always underestimates f

A subgradient of a convex function f at x is any $g \in \mathbb{R}^n$ such that

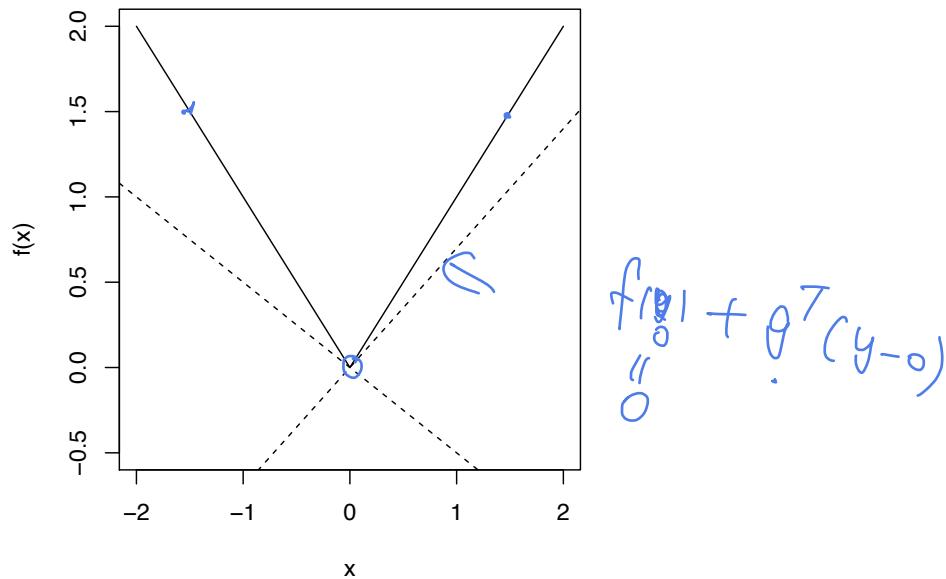
$$f(y) \geq f(x) + g^T(y - x) \quad \text{for all } y$$

- Always exists¹
- If f differentiable at x , then $g = \nabla f(x)$ uniquely
- Same definition works for nonconvex f (however, subgradients need not exist)

¹On the relative interior of $\text{dom}(f)$

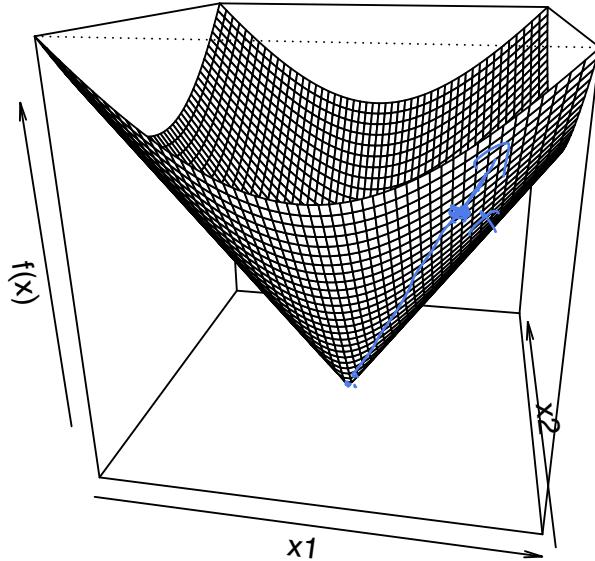
Examples of subgradients

Consider $f : \mathbb{R} \rightarrow \mathbb{R}$, $f(x) = |x|$



- For $x \neq 0$, unique subgradient $g = \text{sign}(x)$
- For $x = 0$, subgradient g is any element of $[-1, 1]$

Consider $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $f(x) = \|x\|_2 = \sqrt{\sum_{i=1}^n x_i^2}$



- For $x \neq 0$, unique subgradient $g = \underline{x/\|x\|_2}$
- For $x = 0$, subgradient g is any element of $\{z : \|z\|_2 \leq 1\}$

Subdifferential

Set of all subgradients of convex f is called the **subdifferential**:

$$\underline{\partial f(x)} = \{g \in \mathbb{R}^n : g \text{ is a subgradient of } f \text{ at } x\}$$

- Nonempty (only for convex f)
- $\partial f(x)$ is closed and convex (even for nonconvex f)
- If f is differentiable at x , then $\partial f(x) = \{\nabla f(x)\}$
- If $\partial f(x) = \{g\}$, then f is differentiable at x and $\nabla f(x) = g$

First order optimality condition with subgradient

$$\begin{cases} \min f_0(x) \\ \text{s.t. } x \in C \end{cases}$$

$$f(x) = f_0(x) + I_C(x)$$

For any f (convex or not),

$$f(x^*) = \underbrace{\min_x f(x)}_{\text{ }} \iff 0 \in \partial f(x^*)$$

i.e., x^* is a minimizer if and only if 0 is a subgradient of f at x^* .

This is called the **subgradient optimality condition**

Why? Easy: $g = 0$ being a subgradient means that for all y

$$f(y) \geq f(x^*) + 0^T(y - x^*) = f(x^*)$$

Note the implication for a convex and differentiable function f ,
with $\partial f(x) = \{\nabla f(x)\}$

Karush-Kuhn-Tucker conditions

Given general problem

$$\begin{aligned} \min_x \quad & f(x) \\ \text{subject to} \quad & \begin{cases} h_i(x) \leq 0, \quad i = 1, \dots, m \\ \ell_j(x) = 0, \quad j = 1, \dots, r \end{cases} \end{aligned}$$

The **Karush-Kuhn-Tucker conditions** or **KKT conditions** are:

- $0 \in \partial \left(f(x) + \sum_{i=1}^m u_i h_i(x) + \sum_{j=1}^r v_j \ell_j(x) \right)$ (stationarity) *X~~argm~~, X L(x_u)*
- $u_i \cdot h_i(x) = 0$ for all i (complementary slackness)
- $h_i(x) \leq 0, \ell_j(x) = 0$ for all i, j (primal feasibility)
- $u_i \geq 0$ for all i (dual feasibility)

Necessity

Let x^* and u^*, v^* be primal and dual solutions with zero duality gap (strong duality holds, e.g., under Slater's condition). Then

$$\begin{aligned} f(x^*) &= g(u^*, v^*) \\ &= \min_x f(x) + \sum_{i=1}^m u_i^* h_i(x) + \sum_{j=1}^r v_j^* \ell_j(x) \\ &\stackrel{\text{optimality}}{\leq} f(x^*) + \sum_{i=1}^m u_i^* h_i(x^*) + \sum_{j=1}^r v_j^* \ell_j(x^*) \\ &\stackrel{\text{complementary slackness}}{=} f(x^*) \end{aligned}$$

In other words, all these inequalities are actually equalities

Two things to learn from this:

- The point x^* minimizes $L(x, u^*, v^*)$ over $x \in \mathbb{R}^n$. Hence the subdifferential of $L(x, u^*, v^*)$ must contain 0 at $x = x^*$ —this is exactly the **stationarity** condition
- We must have $\sum_{i=1}^m \underline{u_i^* h_i(x^*)} = 0$, and since each term here is ≤ 0 , this implies $u_i^* h_i(x^*) = 0$ for every i —this is exactly **complementary slackness**

Primal and dual feasibility hold by virtue of optimality. Therefore:

If x^* and u^*, v^* are primal and dual solutions, with zero duality gap, then x^*, u^*, v^* satisfy the KKT conditions

(Note that this statement assumes nothing a priori about convexity of our problem, i.e., of f, h_i, ℓ_j)

Sufficiency

If there exists x^*, u^*, v^* that satisfy the KKT conditions, then

$$\begin{aligned} g(u^*, v^*) &= f(x^*) + \sum_{i=1}^m u_i^* h_i(x^*) + \sum_{j=1}^r v_j^* \ell_j(x^*) \\ &= f(x^*) \end{aligned}$$

min L(x^{}, u^{*}, v^{*})*

where the first equality holds from stationarity, and the second holds from complementary slackness

Therefore the duality gap is zero (and x^* and u^*, v^* are primal and dual feasible) so x^* and u^*, v^* are primal and dual optimal. Hence, we've shown:

If x^* and u^*, v^* satisfy the KKT conditions, then x^* and u^*, v^* are primal and dual solutions

Putting it together

In summary, KKT conditions:

- always sufficient
- necessary under strong duality



Putting it together:

For a problem with strong duality (e.g., assume Slater's condition: convex problem and there exists x strictly satisfying non-affine inequality constraints),

x^* and u^*, v^* are primal and dual solutions

\iff x^* and u^*, v^* satisfy the KKT conditions

(Warning, concerning the stationarity condition: for a differentiable function f , we cannot use $\partial f(x) = \{\nabla f(x)\}$ unless f is convex!
There are other versions of KKT conditions that deal with local optima.)

Example: support vector machines

Given $y \in \{-1, 1\}^n$, and $X \in \mathbb{R}^{n \times p}$, the support vector machine problem is:

$$\begin{aligned} \min_{\beta, \beta_0, \xi} \quad & \frac{1}{2} \|\beta\|_2^2 + C \sum_{i=1}^n \xi_i \\ \text{subject to} \quad & \xi_i \geq 0, \quad i = 1, \dots, n \\ & y_i(x_i^T \beta + \beta_0) \geq 1 - \xi_i, \quad i = 1, \dots, n \end{aligned}$$

$$\boxed{\beta = X^T w^*}$$

Introduce dual variables $v, w \geq 0$. KKT stationarity condition:

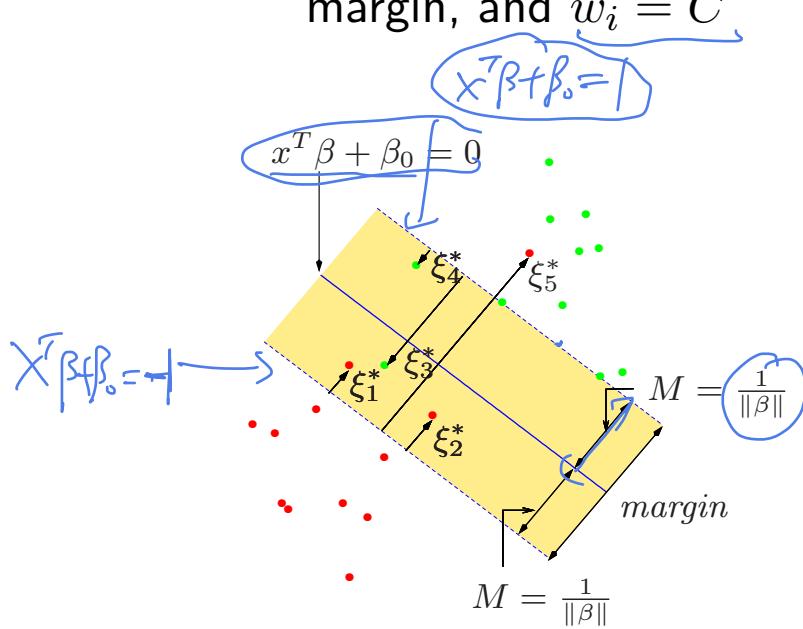
$$0 = \sum_{i=1}^n w_i y_i, \quad \boxed{\beta = \sum_{i=1}^n w_i y_i x_i}, \quad w = C1 - v$$

Complementary slackness: $\boxed{\nabla_{\beta} L(\beta, w, v) = 0}$

$$v_i \xi_i = 0, \quad w_i (1 - \xi_i - y_i(x_i^T \beta + \beta_0)) = 0, \quad i = 1, \dots, n$$

Hence at optimality we have $\beta = \sum_{i=1}^n w_i y_i x_i$, and w_i is nonzero only if $y_i(x_i^T \beta + \beta_0) = 1 - \xi_i$. Such points i are called the **support points**

- For support point i , if $\xi_i = 0$, then x_i lies on edge of margin, and $w_i \in (0, C]$;
- For support point i , if $\xi_i \neq 0$, then x_i lies on wrong side of margin, and $w_i = C$



$$\text{distance between two hyperplanes} = \frac{1}{\|\beta\|}$$

KKT conditions do not really give us a way to find solution, but gives a better understanding

In fact, we can use this to screen away non-support points before performing optimization

Checkpoint: KKT conditions and SVM

- A generalized set of conditions that characterizes the optimal solutions
 - Stationarity, complementary slackness, primal / dual feasibility
 - Always sufficient for optimality
 - Necessary when we have strong duality
- Complementary slackness implies
 - SVM dual solutions are sparse!
 - The number of “support vector”s is small

This lecture

- KKT conditions
 - SVM as an example
- Online Learning

Batch

Recap: Statistical Learning Setting

$$(x_1, y_1) \dots (x_n, y_n) \stackrel{iid}{\sim} D \quad X \in \mathcal{X} = \mathbb{R}^d \text{ or } \{0,1\}^d$$
$$Y \in \mathcal{Y} = \{0,1\}$$

\mathcal{H} hypothesis class $h \in \mathcal{H} \quad h: \mathcal{X} \rightarrow \mathcal{Y}$

Realizable case $\exists h^* \in \mathcal{H}$ s.t. w.p.1 $h^*(x) = y$
 $(x, y) \in D$

Goal of learning: find $h \in \mathcal{H}$ s.t.

$$\underbrace{\text{err}(h)}_{D} = \mathbb{E}_{(x,y) \in D} \left[\underbrace{\mathbb{I}_{\{h^*(x) \neq h(x)\}}}_{\text{I}} \right] \xrightarrow{n \rightarrow \infty} 0$$

function of data

(Adversarial) Online Learning Setting

- Data points show up sequentially (non-iid), learner makes online predictions

x_i , chosen by nature

$h_i \leftarrow h(x_i)$, predict $\hat{y}_i = h(x_i)$

y_i is revealed by nature,

⋮

$h_t \in \{h_1, \dots, h_M\}$, x_t , predict $\hat{y}_t = h_t(x_t)$, receive y_t

$\text{loss}_t = \mathbb{I}(\hat{y}_t \neq y_t)$

- Performance metric: Mistake bounds M

If Alg A satisfies # of mistakes A makes $\leq M$

for all seq. of $(x_1, h^*(x_1)), (x_2, h^*(x_2)), \dots, (x_T, h^*(x_T))$

Then Alg A has a ~~mistake~~ mistake bound of M .

Algorithm A “Consistency”

1. $V_1 = \mathcal{H}$

2. for $t=1, 2, 3, \dots$

Receive X_t , pick any $h \in V_t$

prediction $\hat{y}_t = h(x_t)$

Receive $y_t = h^*(x_t)$

update $V_{t+1} = \{h \in V_t \mid h(x_t) = y_t\}$

Check: $\forall h \in V_{t+1}, h^*(x_i) = y_i$
 $\quad \quad \quad \text{for } \forall i = 1, 2, \dots, t$

Each mistake, we can eliminate at least 1 hypothesis
trivial upper bound

$$|V_t(\text{consistency})| \leq |\mathcal{H}| - 1$$

Example: $\mathcal{X} = \{1, 2, \dots, |\mathcal{H}|\}$

$$\mathcal{H} = \{h_1, \dots, h_{|\mathcal{H}|}\}$$

$$h_i(x) = \begin{cases} 0 & \text{if } x < i \\ 1 & \text{otherwise} \end{cases}$$

$$x_1 = 1, x_2 = 2, \dots, x_{|\mathcal{H}|} = |\mathcal{H}|, \dots$$

$$y_1 = 0, y_2 = 0, \dots, y_{|\mathcal{H}|} = 1$$

$$h^* = h_{|\mathcal{H}|}$$

Predict $h_1, h_2, \dots, h_{|\mathcal{H}|-1}, h_{|\mathcal{H}|}$

ht classifiers (x_i, y_i)

but $h_t(x_t) \neq y_t$

Algorithm B “Halving”

1. $V_1 = \mathcal{H}$

2. for $t=1, 2, 3, \dots$

Receive x_t

Predict $\hat{y}_t = \text{Vote}_{h \in V_t} (h(x_t)) = \arg\max_{r \in \{0, 1\}} \left| \left\{ h \in V_t \mid h(x_t) = r \right\} \right|$

Receive $y_t = h^*(x_t)$

Update $V_{t+1} = \left\{ h \in V_t \mid h(x_t) = y_t \right\}$

Claim: $M(\text{halving}) \leq \log_2(|\mathcal{H}|)$

Proof: for each mistake,
at least $\frac{|V_t|}{2}$ hypotheses
are wrong.

$$|V_{t+1}| \leq |V_t| \cdot \frac{1}{2}$$

$$1 \leq |V_{t+1}| \leq |\mathcal{H}| \cdot 2^{-M}$$

$$2^M \leq |\mathcal{H}| \Rightarrow M \leq \log_2 |\mathcal{H}| \quad \square$$

Now let's get rid of "Realizability". The setting is called "Agnostic learning"

Compete v.s. the best $h \in \mathcal{H}$ in indsight

$$\text{Regret} = \sum_{t=1}^T \mathbb{1}(h_t(x_t) \neq y_t) - \min_{h \in \mathcal{H}} \sum_{t=1}^T \mathbb{1}(h(x_t) \neq y_t)$$

(x_t, y_t) chosen by adversary

as $T \rightarrow \infty$

If $\text{Regret}(T) = o(T)$

$$\frac{1}{T} \text{Regret}(T) = o(1)$$

Example: Stock forecasting

n experts

	Exp 1 (Sigi)	Exp 2 (Esha)	Exp 3 (Lei)	Exp 4 (Raffles the cat)	Outcome
Day 1	Down	Up	Up	Down	Down
Day 2	Up	Up	Down	Down	Down
Day 3	Up	Down	Up	Up	Up
Weighted Majority	1	$\frac{1}{n}$	$\frac{1}{2}$	1	Day 1 Day 2
	$\frac{1}{2}$	$\frac{1}{4}$	$\frac{1}{2}$	1	30

Alg C Weighted Majority

$$\text{predict } \hat{y}_t = \underset{\text{Weighted Output}}{\underbrace{\underset{h \in H}{\sum} w_h^t \cdot h(x_t)}} \geq \frac{1}{\text{Total weight}} \underset{\text{Total weight}}{\underbrace{\sum_{h \in H} w_h^t}}$$

Recent y_t : discount $w_h^{t+1} = w_h^t \cdot \frac{1}{2}$ for h that made a mistake

How do we fix “weighted majority”? Instead of discounting by $1/2$, let’s try discounting by $1-\epsilon$

- Following the same analysis

Fact: For all $0 \leq x \leq 0.5$

$$-x - x^2 \leq \log(1 - x) \leq -x$$

Algorithm D: Randomized Weighted Majority

Analysis of RWM

From mistake bounds to loss minimization

- Loss function
- Regret
- The “Hedge” Algorithm:

Checkpoint: Online Learning

- Learning with expert advice
 - A summary of regret bound: # mistakes - Oracle # of mistakes

	Consistency	Halving	Weighted Majority	Randomized WM
Realizable setting	$\min(T, \mathcal{H})$	$\min(T, \log \mathcal{H})$	$\min(T, \log \mathcal{H})$	$\min(T, \log \mathcal{H})$
Agnostic setting	n.a.	n.a.	$(1 + \epsilon)m + \log \mathcal{H} / \epsilon$	$\sqrt{m \log \mathcal{H} } = O(\sqrt{T \log \mathcal{H} })$

Next lecture

- Online Learning (Part II)
 - Online Gradient Descent
- Reinforcement Learning
 - Markov Decision Processes