

Midterm Exam

Due: 2025-4-10 11:59 | 4 Questions, 100+10 Pts

Name: Jiacong Fang, ID: 2200017849

Problem 1 (10')

For today's dinner, you are making a decision between two restaurants, *Yannan* and *Shaoyuan*. You know that one of them is a strictly better choice under all circumstances, and everyone prefers it to the other. However, you are the only person in the world who has no idea *which* is the better one. You decide to ask some folks on `treehole.pku.edu.cn` to give scores for the two restaurants. According to the following responses, which one will you choose?

- Alice gives a score of 84 for *Yannan*.
- Bob gives a score of 3.92 for *Shaoyuan*.
- Carol gives a score of 99999 for *Yannan*.
- Dave gives a score of 666 for *Shaoyuan*.
- Eve gives a score of $-\pi/6$ for *Yannan*.
- ...

At first glance, it seems impossible that one can do better than random guessing. However, we are going to show that this is not the case.

Formally, consider two score sequences $(x_1, \dots, x_n), (y_1, \dots, y_n) \in D^n$, where $D \subseteq \mathbb{R}$, such that either of the two cases is true:

- Case 0: $x_i < y_i, \forall i \in [n]$.
- Case 1: $x_i > y_i, \forall i \in [n]$.

We need to design a potentially randomized algorithm \mathcal{A} , whose output is either 0 or 1, to distinguish the two cases. Notably, we are restricted to “half-blind” algorithms: For each $i \in [n]$, algorithm \mathcal{A} can only query one of the two values of $\{x_i, y_i\}$; once it decided to query one value, the other one in the i 'th pair becomes inaccessible.

- a. (3') Suppose $n = 1$ and $D = \mathbb{R}$. Design a half-blind algorithm \mathcal{A}_1 such that for any real numbers x, y ,

$$\Pr[\mathcal{A}_1 \text{ is correct} \mid x_1 = x, y_1 = y] > \frac{1}{2}.$$

- b. (7') Suppose $D = [0, 100]$, $|x_i - y_i| \geq 1$. Design a half-blind algorithm \mathcal{A} for sufficiently large n , and show that it has $1 - 1/2^{p(n)}$ success probability given any feasible input, where $p(n) = \Omega(n)$.

Direction: For each i , your half-blind algorithms can decide freely which one of $\{x_i, y_i\}$ to observe. However, it can only observe one of the two values. You should give proofs for the error bounds.

Answer. a). Consider the following algorithm \mathcal{A}_1 :

Algorithm 1 Half-bind Algorithm \mathcal{A}_1

Input: $x_1, y_1 \in \mathbb{R}$

Output: Case 1 ($x_1 > y_1$) or Case 0 ($x_1 < y_1$).

```

1: Randomly query  $x_1$  or  $y_1$  with equal probability
2: if  $x_1$  is queried then
3:   Randomly sample  $z_x \sim \mathcal{N}(0, 1)$ 
4:   if  $x_1 > z_x$  then
5:     Output Case 1 ( $x_1 > y_1$ ).
6:   else
7:     Output Case 0 ( $x_1 < y_1$ ).
8: else
9:   Randomly sample  $z_y \sim \mathcal{N}(0, 1)$ 
10:  if  $y_1 > z_y$  then
11:    Output Case 0 ( $x_1 < y_1$ ).
12:  else
13:    Output Case 1 ( $x_1 > y_1$ ).
```

For any real numbers x, y , without loss of generality, we assume $x > y$ here.

$$\begin{aligned} \Pr[\mathcal{A}_1 \text{ is correct} \mid x_1 = x, y_1 = y] &= \frac{1}{2} \left(\Pr_{z_x \sim \mathcal{N}(0,1)}[x > z_x] + \Pr_{z_y \sim \mathcal{N}(0,1)}[y < z_y] \right) \\ &= \frac{1}{2} (\Phi(x) + 1 - \Phi(y)) > \frac{1}{2}. \end{aligned}$$

where $\Phi(x)$ is the cumulative distribution function of the standard normal distribution. The last inequality holds because $\Phi(x) - \Phi(y) > 0$ when $x > y$.

b). For any $i \in [n]$, we first define the following algorithm \mathcal{A}_i :

Algorithm 2 Half-bind Algorithm \mathcal{A}_i

Input: $x_i, y_i \in [0, 100]$

Output: Case 1 ($x_i > y_i$) or Case 0 ($x_i < y_i$).

```

1: Randomly query  $x_i$  or  $y_i$  with equal probability
2: if  $x_i$  is queried then
3:   Randomly sample  $z_x \sim \mathcal{U}(0, 100)$  ▷ Uniform distribution over  $[0, 100]$ 
4:   if  $x_i > z_x$  then
5:     Output Case 1 ( $x_i > y_i$ ).
6:   else
7:     Output Case 0 ( $x_i < y_i$ ).
8: else
9:   Randomly sample  $z_y \sim \mathcal{U}(0, 100)$ 
10:  if  $y_i > z_y$  then
11:    Output Case 0 ( $x_i < y_i$ ).
12:  else
13:    Output Case 1 ( $x_i > y_i$ ).
```

Then we can design the following algorithm \mathcal{A} :

Algorithm 3 Half-bind Algorithm \mathcal{A}

Input: $(x_1, \dots, x_n), (y_1, \dots, y_n) \in [0, 100]^n$

Output: Case 1 ($x_i > y_i$), $\forall i \in [n]$ or Case 0 ($x_i < y_i$), $\forall i \in [n]$.

```

1: num1 ← 0, num0 ← 0 ▷ Count of Case 1 and Case 0
2: for  $i = 1$  to  $n$  do
3:   Run Algorithm  $\mathcal{A}_i$  on  $(x_i, y_i)$ 
4:   if Algorithm  $\mathcal{A}_i$  outputs Case 1 then
5:     num1 ← num1 + 1
6:   else
7:     num0 ← num0 + 1
8:   if num1 >  $n/2$  then
9:     Output Case 1  $(x_i > y_i), \forall i \in [n]$ .
10: else
11:   Output Case 0  $(x_i < y_i), \forall i \in [n]$ .

```

Error Bounds Analysis: Without loss of generality, we assume $x_i > y_i$ for all $i \in [n]$ and $n = 2k + 1$ for $k \in \mathbb{N}^*$. (When n is even, we can simply ignore x_n and y_n and only run \mathcal{A} on the first $n - 1$ pairs, it won't affect the error bounds to much.)

For $i \in [n]$, let random variable $X_i := \mathbf{1}[\mathcal{A}_i \text{ outputs } x_i > y_i]$. We first analyse Algorithm \mathcal{A}_i , similar to the analysis in part a, we have:

$$\begin{aligned} \Pr[\mathcal{A}_i \text{ is correct} \mid x_i, y_i] &= \frac{1}{2} \left(\Pr_{z_x \sim \mathcal{U}(0,100)}[x_i > z_x] + \Pr_{z_y \sim \mathcal{U}(0,100)}[y_i < z_y] \right) \\ &= \frac{1}{2} \left(1 + \frac{|x_i - y_i|}{100} \right) \geq \frac{1}{2} + \frac{1}{200} = \frac{101}{200}. \end{aligned}$$

Thus, we have $\mathbb{E}[X_i] \geq 101/200$. Therefore, the probability that \mathcal{A} makes a mistake can be bounded as follows:

$$\begin{aligned} \Pr[\mathcal{A} \text{ is wrong}] &= \Pr\left[\sum_{i=1}^n X_i \leq \frac{n}{2}\right] = \Pr\left[\frac{1}{n} \sum_{i=1}^n X_i - \frac{101}{200} \leq -\frac{1}{200}\right] \\ &\leq \Pr\left[\frac{1}{n} \sum_{i=1}^n X_i - \mathbb{E}[X_i] \leq -\frac{1}{200}\right] \\ &\leq \exp\left(-2n \cdot \frac{1}{200^2}\right) \quad (\text{Chernoff Bound}) \\ &\leq \frac{1}{2^{p(n)}} \quad (\text{for sufficiently large } n, \text{ where } p(n) = \Omega(n)). \end{aligned}$$

Then we can conclude that:

$$\Pr[\mathcal{A} \text{ is correct}] = 1 - \Pr[\mathcal{A} \text{ is wrong}] \geq 1 - \frac{1}{2^{p(n)}}. \quad (1)$$

where $p(n) = \Omega(n)$.

Note: The error bound can be analysed by **Median Trick**. Each run of \mathcal{A}_i is a Bernoulli trial with success probability p_i greater than $101/200$. Then

$$\Pr[\mathcal{A}_i \text{ success}] \geq \frac{101}{200} \implies \Pr\left[\#\{\mathcal{A}_i \text{ success in } \mathcal{A}\} \leq \frac{n}{2}\right] \leq \left(\frac{101}{200}\right)^{n/2}$$

Thus, (1) holds for $p(n) = \Omega(n)$ with sufficiently large n .

Q.E.D.

◁

Problem 2 (15') (Knapsack Counting)

Recall the 0-1 Knapsack problem you have encountered before. We are given n objects of weight $0 \leq a_1 \leq \dots \leq a_n \leq b$ where b is the total volume of the knapsack, and we are asked to find a 0-1 assignment x_1, \dots, x_n such that $\sum_{i=1}^n a_i x_i \leq b$. In the counting version, we will compute the number of valid assignments, the set of which is denoted as S , instead of finding the optimal solution. The key idea here is similar to #DNF; that is, sampling from a restricted space. Suppose b is much larger than n^2 .

We will use the *scaling method*. Define $a'_i := \lfloor n^2 a_i / b \rfloor$, then $0 \leq a'_i \leq n^2$. Note the floor operator may cause error when trying to recover a_i from a'_i , but the total error is no more than $n(b/n^2) = b/n$. Let $S' := \{x : \sum_{i=1}^n a'_i x_i \leq n^2\}$ be the set of solutions to the rescaled problem.

a. (2') Prove that $S \subseteq S'$.

Apparently, the converse does not necessarily hold. However, we claim that any $x \in S'$ can be shoehorned into S by setting at most one of the x_i from 1 to 0.

b. (5') Prove the claim above, which implies $|S'| \leq (n+1)|S|$.

In other words, with probability at least $1/(n+1)$, a uniformly random element of S' is in S . So the remaining task here is to compute $|S'|$ and sample from S' . Let $C(k, m) := \left| \left\{ x : \sum_{i=1}^k a'_i x_i \leq m \right\} \right|$.

c. (3') Design an algorithm based on dynamic programming that computes $C(i, j)$ for all $0 \leq i \leq n$ and $0 \leq j \leq n^2$. Your algorithm should take $O(n^3)$ time.

Definitely, we can fetch the value of $|S'|$ from $C(n, n^2)$. But in fact, we can also use this table to sample uniformly from S' .

d. (5') Based on the values of $C(\cdot, \cdot)$, design an algorithm which outputs a uniform sample of S' . Your algorithm should run in $O(n)$ time (excluding the time cost by the random generator).

Answer. a). For any $x = (x_1, x_2, \dots, x_n) \in S$, we have

$$\sum_{i=1}^n a_i x_i \leq b \iff \sum_{i=1}^n \frac{n^2 a_i}{b} x_i \leq n^2$$

And notice that

$$\sum_{i=1}^n a'_i x_i = \sum_{i=1}^n \left\lfloor \frac{n^2 a_i}{b} \right\rfloor x_i \leq \sum_{i=1}^n \frac{n^2 a_i}{b} x_i \leq n^2$$

Then we have $x \in S'$ which implies $S \subseteq S'$.

b). Define k be s.t. $a_j \leq b/n$ for $j \leq k$ and either $k = n$ or $a_{k+1} > b/n$. Let $C = \{0, 1\}^k \times \{0\}^{n-k}$. If $x \in C$ then

$$\sum_{i=1}^n a_i x_i \leq \sum_{i=1}^k a_i x_i \leq \frac{bk}{n} \leq n \implies C \subseteq S$$

For any $x \in S' \setminus S$, then there exists $p = p(x) \in \mathbb{N}^*$, s.t. $x_p = 1$ and $p \notin [k]$ (Otherwise $x \in C \subseteq S$ which is a contradiction). Then the sheohorn(or a mapping) from S' to S is defined as:

$$f : S' \rightarrow \{0, 1\}^n$$

$$f(x) \mapsto \begin{cases} x & \text{if } x \in S \\ y & \text{if } x \in S' \setminus S \end{cases}$$

where $y_j = x_j, \forall j \neq p$ and $y_p = 0$.

We now prove that for any $\mathbf{x} \in S' \setminus S$, $\mathbf{y} = f(\mathbf{x}) \in S$. Let $\delta_i = a_i - a'_i$, then we have

$$\begin{aligned} \sum_{i=1}^n a_i y_i &= \frac{b}{n^2} \sum_{i=1}^n (a'_i + \delta_i) y_i = \frac{b}{n^2} \left(\sum_{i=1}^n a'_i y_i + \sum_{i=1}^n \delta_i y_i \right) \\ &= \frac{b}{n^2} \left(\sum_{i=1}^n a'_i x_i - a'_p + \sum_{i=1}^n \delta_i y_i \right) \\ &\leq \frac{b}{n^2} (n^2 - n + n) = b. \end{aligned}$$

The last inequality holds because $\delta \in [0, 1)$ and $p \notin [k] \implies a_p > b/n \implies a'_p \geq n$. Therefore, $\mathbf{y} \in S$ which implies $|S'| \leq (n+1)|S|$.

c). The dynamic programming algorithm is as follows:

Algorithm 4 Dynamic Programming for $|S'|$

Input: $a'_1, a'_2, \dots, a'_n \in [0, n^2]$

Output: $C(i, j)$ for all $0 \leq i \leq n$ and $0 \leq j \leq n^2$.

1: $\forall 0 \leq j \leq n^2$, initialize $C(\cdot, \cdot)$ with:

$$C(1, j) = \begin{cases} 2 & \text{if } j \geq a'_1 \\ 1 & \text{otherwise} \end{cases}$$

2: **for** $i = 2$ to n **do**

▷ Compute the table C using following recursion.

3: **for** $j = 0$ to n^2 **do**

4: **if** $j < a'_i$ **then**

5: $C(i, j) = C(i-1, j)$

6: **else**

7: $C(i, j) = C(i-1, j) + C(i-1, j - a'_i)$

The time complexity of the algorithm is $O(n^3)$.

d). We can determine a uniform sample from S' by tracing back probabilistically from $C(n, n^2)$ as following:

Algorithm 5 Uniform sampling from S'

Input: $C(i, j)$ for all $0 \leq i \leq n$ and $0 \leq j \leq n^2$.

Output: A uniform sample from S' .

1: $j \leftarrow n^2$

2: **for** $i = n$ to 2 **do**

3: $x_i \leftarrow 0$ with prob. $C(i-1, j)/C(i, j)$ else $x_i \leftarrow 1$ with the remaining prob.

4: **if** $x_i = 1$ **then**

5: $j \leftarrow j - a'_i$

▷ Update the remaining capacity

6: **else**

7: $j \leftarrow j$

8: **if** $C(1, j) = 2$ **then**

▷ Handle the special case when $i = 1$

9: $x_1 \leftarrow 0$ and $x_1 \leftarrow 1$ with equal prob.

10: **else**

11: $x_1 \leftarrow 0$

12: **Output** (x_1, x_2, \dots, x_n) .

The time complexity of the algorithm is $O(n)$ because each iteration of the for loop takes $O(1)$ time, and the loop runs n times.

Q.E.D.

◁

Problem 3 (50') (Low Diameter Decomposition)

In this problem, you'll learn Low Diameter Decomposition, and use it to give approximate algorithm for some fundamental graph problem like Tree Embedding, All Pair Shortest Path (APSP).

In this problem, the graph can be seen as weighted graph with all the weight being positive integer.

Definition 1 (Low Diameter Decomposition (LDD)). Given an undirected graph $G = (V, E)$, a **Low Diameter Decomposition (LDD)** scheme with approximation factor β and diameter bound D is a randomized algorithm that partitions V into disjoint clusters V_1, V_2, \dots, V_k satisfying:

1. **Bounded Diameter:** For each V_i , the diameter of $G[V_i] \leq D$.
2. **Separation Probability:** For any $x, y \in V$,

$$\Pr [x \text{ and } y \text{ lie in different clusters}] \leq \beta \cdot \frac{d_G(x, y)}{D},$$

where $d_G(x, y)$ denotes the shortest-path distance between x and y in G .

Remarks:

- The **diameter** of $G[V_i]$ is $\max_{u, v \in V_i} d_G(u, v)$.
 - β balances cluster tightness (D) and separation likelihood. Lower β implies better decomposition quality.
- a. (10') Prove that the following algorithm gives a LDD with approximation factor $\beta = O(\log n)$, with probability $\geq 1 - n^{-1}$.

Algorithm 6 Low Diameter Decomposition (LDD)

Input: Undirected graph $G = (V, E)$, target diameter $D > 0$

Output: Vertex partition $\{V_1, \dots, V_k\}$ with $\text{diam}(G[V_i]) \leq D$

- 1: Initialize $\text{marked}[v] \leftarrow \text{FALSE}$ for all $v \in V$
 - 2: **while** $\exists v \in V$ with $\neg \text{marked}[v]$ **do**
 - 3: Select arbitrary unmarked vertex $v_0 \in V$
 - 4: Sample $R_{v_0} \sim \text{Geometric}(p)$ with $p = \min(1, \frac{4 \log_e n}{D})$.
 - 5: Compute $B \leftarrow \{u \in V \mid \neg \text{marked}[u] \wedge d_G(v_0, u) \leq R_{v_0}\}$
 - 6: Create cluster $C \leftarrow B$
 - 7: $\text{marked}[u] \leftarrow \text{TRUE}$ for all $u \in C$
 - 8: Add C to output partition
 - 9: **return** the computed clustering
-

Remarks: The naive way to run LDD takes time $O(n^3)$, since one need to run Dijkstra for n times. However, there are ways to do LDD in $\tilde{O}(n^2)$ time. You will get **10 Bonus Point** if you can find out.

We will use this tool to solve approximate APSP problem. First, we will introduce Low Stretch Tree.

Definition 2. A randomized low-stretch tree of stretch α for a graph $G = (V, E)$ is a probability distribution \mathcal{D} over spanning trees of G s.t.

1. $d_G(x, y) \leq d_T(x, y)$, for all T in the support \mathcal{D} .
2. $\mathbb{E}_{T \sim \mathcal{D}}[d_T(x, y)] \leq \alpha \cdot d_G(x, y)$, $\forall x, y \in V$

Theorem 3. For any metric space $M = (V, d)$, there exists an efficiently sampleable α_B -stretch spanning tree distribution \mathcal{D}_B , where

$$\alpha_B = O(\log n \log \Delta_M)$$

Δ_M is defined as $\max_{x,y} d(x, y)$, we assume $\forall x \neq y, d(x, y) \geq 1$.

We will prove theorem 3 by the following algorithm.

Algorithm 7 Low Stretch Tree Construction, $LST(M, \delta)$

Input: Metric space $M = (V, d)$, target diameter $D = 2^\delta$

Output: Spanning tree T with low stretch. **Invariant:** $\text{diameter}(M) \leq 2^\delta$

- 1: **if** $|V| = 1$ **then**
 - 2: **return** trivial tree containing the single point
 - 3: Partition V into clusters $C_1, \dots, C_t \leftarrow \text{LDD}(M, D/2)$
 - 4: **for** $j = 1$ **to** t **do**
 - 5: Let M_j be M restricted to C_j
 - 6: Recursively build $T_j \leftarrow LST(M_j, \delta - 1)$
 - 7: Connect roots r_2, \dots, r_t to r_1 with edges of length 2^δ
 - 8: **return** final tree T rooted at r_1
-

Lemma 4. If the random tree T returned by some call $LST(M, \delta)$ has root r , then

1. every vertex x in T has distance $d_T(x, r) \leq 2^{\delta+1}$
2. the expected distance between any $x, y \in T$ has $\mathbb{E}[d_T(x, y)] \leq 8\delta\beta d(x, y)$. (Recall that β is the approximate factor in LDD)

If we can prove Lemma 4, then one can see from Problem a that if $\beta = O(\log n)$, then with high probability, $d_T(x, y) \geq d(x, y)$ for any x, y , thus theorem 3 can be proved.

b. (20') Prove the Lemma 4.

c. (10') Prove the following theorem.

Theorem 5. There's an algorithm that output $O(\log n \log \Delta)$ approximation of APSP on an undirected graph in $\tilde{O}(n^2)$ time, and success with probability $\geq 1 - \frac{1}{\text{poly}(n)}$.

(This means, the algorithm output $d'(x, y)$ for any pair x, y , and satisfies $d(x, y) \leq d'(x, y) \leq \log n \log \Delta \cdot d(x, y)$, where $d(x, y)$ is the length of shortest path between x, y .)

d. (10') Prove the following theorem.

Theorem 6. There's an algorithm that output $O(\log n)$ approximation of ASAP on an undirected graph in $\tilde{O}(n^2)$ time, and success with probability $\geq 1 - \frac{1}{\text{poly}(n)}$.

Problem 4 (25') (Sub-Network Approximation for Two-Layer Neural Networks)

Definition 7 (Two-Layer Neural Network). Consider a wide neural network $f_M(x; \theta) : \mathbb{R}^d \rightarrow \mathbb{R}$ of the form

$$f_M(x; \theta) := \sum_{j=1}^M a_j \sigma(w_j^\top x),$$

where $\theta = \{a_1, \dots, a_M, w_1, \dots, w_M\}$ denotes the parameters of the network, $a_j \in \mathbb{R} \setminus \{0\}$ are the non-zero coefficients, and $w_j \in \mathbb{R}^d$ are the weights associated with the j -th neuron. Here, $\sigma : \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$ denotes ReLU activation function that is defined as $\sigma(\cdot) = \max\{0, \cdot\}$.

Definition 8 (Parameter Norm). The norm of the parameters θ is defined as

$$\|\theta\|_{\mathcal{P}} := \sum_{j=1}^M |a_j| \|w_j\|_2,$$

where $\|w_j\|_2$ denotes the Euclidean norm of the weight vector w_j , and $\|\theta\|_{\mathcal{P}}$ is assumed to be finite.

Let $\pi = \text{Unif}(\mathbb{S}^{d-1})$ denote the uniform distribution over the unit sphere $\mathbb{S}^{d-1} \subset \mathbb{R}^d$. The goal is to approximate the network $f_M(x; \theta)$ with a sparse subset of the neurons, while controlling the error in expectation. Indeed, please prove the following result:

Theorem 9. For any integer $m \in \mathbb{N}$, there exists a sparse vector $\tilde{a} = (\tilde{a}_1, \dots, \tilde{a}_M)$ such that $\|\tilde{a}\|_0 = m$, i.e., the vector \tilde{a} contains exactly m non-zero entries, which satisfies the following approximation bound:

$$\mathbb{E}_{x \sim \pi} \left(\sum_{j=1}^M \tilde{a}_j \sigma(w_j^\top x) - \sum_{j=1}^M a_j \sigma(w_j^\top x) \right)^2 \lesssim \frac{\|\theta\|_{\mathcal{P}}^2}{dm}.$$

Remark 10. In essence, this result demonstrates that it is possible to approximate the network with a sparse representation without incurring significant loss in performance, provided the number of selected neurons m is sufficiently large relative to the network's parameter norm $\|\theta\|_{\mathcal{P}}$, and the size of the input space d .

Answer. We try to sample \tilde{a} and construct an unbiased estimator whose expectation is $f_M(x; \theta)$. For $j \in [M]$, we define the following probability mass function (Denote as distribution \mathcal{D}_M).

$$p_j = \frac{|a_j| \cdot \|w_j\|_2}{\|\theta\|_{\mathcal{P}}} \quad \text{and} \quad \sum_{j=1}^M p_j = 1.$$

Then we can sample m neurons from M neurons with $\{p_j\}_{j=1}^M$ independently, i.e.,

$$\text{sample } j_1, j_2, \dots, j_m \stackrel{i.i.d.}{\sim} \mathcal{D}_M.$$

Then we can define the following estimator, for $t \in [m]$,

$$Y_t(x) := \frac{\text{sign}(a_{j_t})}{\|w_{j_t}\|_2} \cdot \sigma(w_{j_t}^\top x) \quad \text{and} \quad \tilde{f}(x) := c \cdot \sum_{t=1}^m Y_t(x).$$

where c is a constant waiting to be determined to makes sure that $\tilde{f}(x)$ is unbiased. Then we can assign \tilde{a} as follows:

$$\tilde{a}_i = c \cdot \frac{\text{sign}(a_i)}{\|w_i\|_2} \text{ for } i \in \{j_t\}_{t=1}^m, \text{ and } \tilde{a}_i = 0 \text{ otherwise.}$$

Notice that j_1, j_2, \dots, j_m may be not distinct, but we can merge those terms into a single nonzero entry. Thus, here we have $\|\tilde{a}\|_0 \leq m$ (it's a weaker condition than $\|\tilde{a}\|_0 = m$ actually).

Now we determine the constant c first. Given x , we have:

$$\begin{aligned} \mathbb{E}_{j_t \sim \mathcal{D}_M} [Y_t(x)] &= \sum_{j=1}^M p_j \cdot \frac{\text{sign}(a_j)}{\|w_j\|_2} \sigma(w_j^\top x) \\ &= \sum_{j=1}^M \frac{|a_j| \|w_j\|_2}{\|\theta\|_{\mathcal{P}}} \cdot \frac{\text{sign}(a_j)}{\|w_j\|_2} \sigma(w_j^\top x) \\ &= \frac{1}{\|\theta\|_{\mathcal{P}}} \sum_{j=1}^M a_j \sigma(w_j^\top x) \quad (\text{use } \text{sign}(a_j) \cdot |a_j| = a_j) \\ &= \frac{1}{\|\theta\|_{\mathcal{P}}} f_M(x; \theta). \end{aligned}$$

Then (let $c = \|\theta\|_{\mathcal{P}}/m$)

$$\mathbb{E}_{j_t \sim \mathcal{D}_M} [\tilde{f}(x)] = c \cdot \sum_{t=1}^m \mathbb{E}_{j_t \sim \mathcal{D}_M} [Y_t(x)] = \frac{c \cdot m}{\|\theta\|_{\mathcal{P}}} f_M(x; \theta) = f_M(x; \theta).$$

Therefore, $\tilde{f}(x)$ is an unbiased estimator of $f_M(x; \theta)$ over \mathcal{D}_M .

Look back the error term we want to bound, notice that:

$$\begin{aligned} &\mathbb{E}_{x \sim \pi} \left(\sum_{j=1}^M \tilde{a}_j \sigma(w_j^\top x) - \sum_{j=1}^M a_j \sigma(w_j^\top x) \right)^2 \\ &= \mathbb{E}_{x \sim \pi} \left[\left(\tilde{f}(x) - \mathbb{E}_{j_t \sim \mathcal{D}_M} [\tilde{f}(x)] \right)^2 \right] = \mathbb{E}_{x \sim \pi} \left[\text{Var}_{j_t \sim \mathcal{D}_M} (\tilde{f}(x)) \right] \\ &= \mathbb{E}_{x \sim \pi} \left[\frac{\|\theta\|_{\mathcal{P}}^2 \cdot m}{m^2} \text{Var}_{j_t \sim \mathcal{D}_M} (Y_t(x)) \right] \quad (t \text{ is a fixed index here}) \\ &= \frac{\|\theta\|_{\mathcal{P}}^2}{m} \cdot \mathbb{E}_{x \sim \pi} \left[\mathbb{E}_{j_t \sim \mathcal{D}_M} [(Y_t(x))^2] - \left(\mathbb{E}_{j_t \sim \mathcal{D}_M} [Y_t(x)] \right)^2 \right] \\ &\leq \frac{\|\theta\|_{\mathcal{P}}^2}{m} \cdot \mathbb{E}_{x \sim \pi} \left[\mathbb{E}_{j_t \sim \mathcal{D}_M} [(Y_t(x))^2] \right] \\ &= \frac{\|\theta\|_{\mathcal{P}}^2}{m} \cdot \mathbb{E}_{j_t \sim \mathcal{D}_M} \left[\mathbb{E}_{x \sim \pi} [(Y_t(x))^2] \right]. \end{aligned}$$

The inequality holds because $\mathbb{E}_{j_t \sim \mathcal{D}_M} [Y_t(x)]^2 = f_M^2(x; \theta) / \|\theta\|_2^2 \geq 0$. Then given j_t , we try to calculate

$$\begin{aligned} \mathbb{E}_{x \sim \pi} \left[\left(\sigma(w_{j_t}^\top x) \right)^2 \right] &= \int_{\mathbb{S}^{d-1}} \left(\max\{0, w_{j_t}^\top x\} \right)^2 \cdot d\pi(x) \\ &= \int_{\{w_{j_t}^\top x \geq 0\}} \left(w_{j_t}^\top x \right)^2 \cdot d\pi(x) \\ &= \frac{1}{2} \int_{\mathbb{S}^{d-1}} \left(w_{j_t}^\top x \right)^2 \cdot d\pi(x) \quad (w_{j_t}^\top x \text{ is symmetric over } \mathbb{S}^{d-1}) \\ &= \frac{1}{2} \mathbb{E}_{x \sim \pi} \left[\left(w_{j_t}^\top x \right)^2 \right] \\ &= \frac{\|w_{j_t}\|_2^2}{2} \cdot \mathbb{E}_{x \sim \pi} \left[\left(u^\top x \right)^2 \right] \quad (\text{where } u = \frac{w_{j_t}}{\|w_{j_t}\|_2}). \end{aligned}$$

We have $\mathbb{E}_{x \sim \pi} \left[(u^\top x)^2 \right] = 1/d$ (it will be proved at the end, skip it for now). Therefore,

$$\mathbb{E}_{x \sim \pi} \left[\left(\sigma(w_{j_t}^\top x) \right)^2 \right] = \frac{\|w_{j_t}\|_2^2}{2d}. \implies \mathbb{E}_{x \sim \pi} [Y_t(x)^2] = \frac{\|w_{j_t}\|_2^2}{2d} \cdot \frac{\text{sign}(a_{j_t})^2}{\|w_{j_t}\|_2^2} = \frac{1}{2d}.$$

Then we have

$$\mathbb{E}_{x \sim \pi} \left(\sum_{j=1}^M \tilde{a}_j \sigma(w_j^\top x) - \sum_{j=1}^M a_j \sigma(w_j^\top x) \right)^2 \leq \frac{\|\theta\|_{\mathcal{P}}^2}{m} \cdot \mathbb{E}_{j_t \sim \mathcal{D}_M} \left[\frac{1}{2d} \right] = \frac{\|\theta\|_{\mathcal{P}}^2}{2dm} \lesssim \frac{\|\theta\|_{\mathcal{P}}^2}{dm}.$$

At the end, we prove that

$$\mathbb{E}_{x \sim \pi} \left[(u^\top x)^2 \right] = 1/d.$$

Expand u to d -orthonormal basis of \mathbb{R}^d , denoted as $\{e_1, e_2, \dots, e_d\}$, then for any given $x \in \mathbb{S}^{d-1}$, we have

$$\begin{aligned} x &= \sum_{i=1}^d \alpha_i e_i \quad \text{where} \quad \sum_{i=1}^d \alpha_i^2 = 1, \implies \sum_{i=1}^d (e_i^\top u)^2 = \sum_{i=1}^d \alpha_i^2 = 1 \\ \implies 1 &= \mathbb{E}_{x \sim \pi} \left[\sum_{i=1}^d (e_i^\top u)^2 \right] = d \cdot \mathbb{E}_{x \sim \pi} \left[(u^\top x)^2 \right] \implies \mathbb{E}_{x \sim \pi} \left[(u^\top x)^2 \right] = \frac{1}{d}. \end{aligned}$$

Q.E.D.

◁