

A Proof of Lemma 4

The following lemma states that the with high probability, the ratio $\frac{N_{\text{far}}(x)}{N(x)} \rightarrow 0$ as n approaches ∞ . Throughout this section, $B(x, r)$ denotes the closed Euclidean r -ball about x .

Lemma 4. *Let $x \sim \mathcal{D}_X$. Then, for all $\delta > 0$, $\frac{1}{2} < s < 1$, the hash table T calculated by Algorithm 1 satisfies:*

$$\begin{aligned} \mathbb{P}\left(N_{\text{far}}(x) > \delta N(x)\right) &\leq \\ &2 \exp(-0.09n^{\frac{1-s}{2}}) + \left(\frac{1.6}{n^{1-s}\sqrt{d}^d}\right)^{\frac{1}{d+1}} \\ &+ \sqrt{\frac{1}{n^{\frac{1-s}{2}}}} + \exp\left(-\frac{1}{8}n^{s-\frac{1}{2}}\right) \\ &+ 2^{-\frac{\delta}{2}n^{s-\frac{1}{2}}}, \end{aligned}$$

where the probability is over $S_n, x \sim \mathcal{D}^{n+1}$ and the choice of the function g_n .

Proof. Fix $\delta > 0$, $\varepsilon = \frac{r_n}{\sqrt{d}}$, and let C_1, \dots, C_t be a partition of $[0, 1]^d$ into $t = (\frac{1}{\varepsilon})^d$ boxes of length ε . Notice that for any x, x' in the same box, we have $\|x - x'\| \leq \sqrt{d}\varepsilon$. Put $k = n^s$ and define the random variable $L_{\varepsilon, k}(S_n) = \sum_{i: |C_i \cap S_n| < k} \mathbb{P}(C_i)$, and note that it is precisely the k -missing mass (defined in (2)) associated with the distribution $P = (\mathbb{P}(C_1), \dots, \mathbb{P}(C_t))$. By Theorem 2(a), we have $\mathbb{E}[L_{\varepsilon, k}(S_n)] \leq \frac{1.6kt}{n}$. By the law of total probability,

$$\begin{aligned} \mathbb{P}\left(N_{\text{far}}(x) > \delta N(x)\right) &\leq \mathbb{P}\left(L_{\varepsilon, m}(S_n) > \frac{1.6}{\varepsilon^d n^{1-s}} + \gamma\right) \\ &+ \mathbb{P}\left(N_{\text{far}}(x) > \delta N(x) \mid \right. \\ &\quad \left. L_{\varepsilon, m}(S_n) \leq \frac{1.6}{\varepsilon^d n^{1-s}} + \gamma\right). \end{aligned} \quad (18)$$

For the first term in (18), we apply Theorem 2(b):

$$\begin{aligned} \mathbb{P}\left(L_{\varepsilon, m}(S_n) > \frac{1.6}{\varepsilon^d n^{1-s}} + \gamma\right) &\leq \mathbb{P}\left(L_{\varepsilon, m}(S_n) > \mathbb{E}[L_{\varepsilon, m}(S_n)] + \gamma\right) \\ &\leq 2 \exp(-0.09n^{1-s}\gamma^2). \end{aligned}$$

For the second term in (18), we have

$$\begin{aligned} \mathbb{P}\left(N_{\text{far}}(x) > \delta N(x) \mid L_{\varepsilon, m}(S_n) \leq \frac{1.6}{\varepsilon^d n^{1-s}} + \gamma\right) &\leq \mathbb{P}\left(|B(x, r_n) \cap S_n| < n^s \mid \right. \\ &\quad \left. L_{\varepsilon, m}(S_n) \leq \frac{1.6}{\varepsilon^d n^{1-s}} + \gamma\right) \\ &+ \mathbb{P}\left(N_{\text{far}}(x) > \delta N(x), |B(x, r_n) \cap S_n| \geq n^s \mid \right. \\ &\quad \left. L_{\varepsilon, m}(S_n) \leq \frac{1.6}{\varepsilon^d n^{1-s}} + \gamma\right) \\ &= (*) + (**). \end{aligned} \quad (19)$$

Since $r_n = \sqrt{d}\varepsilon$, we have $\{|B(x, r_n) \cap S_n| < n^s\} \implies \{|C(x) \cap S_n| < n^s\}$, where $C(x)$ is the ε -length box containing x . Thus,

$$(*) \leq \frac{1.6}{\varepsilon^d n^{1-s}} + \gamma.$$

We are left to bound the second term in (19)

$$\begin{aligned} (**) &\leq \mathbb{P}\left(N_{\text{close}}(x) < \frac{1}{2}n^{s-\frac{1}{2}} \mid \right. \\ &\quad \left. L_{\varepsilon, m}(S_n) \leq \frac{1.6}{\varepsilon^d n^{1-s}} + \gamma, |B(x, r_n) \cap S_n| > n^s\right) \\ &+ \mathbb{P}\left(N_{\text{far}}(x) > \delta N(x), N_{\text{close}}(x) \geq \frac{1}{2}n^{s-\frac{1}{2}} \mid \right. \\ &\quad \left. L_{\varepsilon, m}(S_n) \leq \frac{1.6}{\varepsilon^d n^{1-s}} + \gamma, |B(x, r_n) \cap S_n| > n^s\right) \\ &= (***) + (****). \end{aligned} \quad (20)$$

Since the algorithm set $m_n = \lfloor \frac{\log n}{2 \log(\frac{1}{p_1})} \rfloor$, we have

$$\begin{aligned} \mathbb{E}\left[N_{\text{close}}(x) \mid |B(x, r) \cap S_n| > n^s\right] &\geq p_1^{m_n} n^s \\ &\geq p_1^{\frac{\log n}{2 \log(\frac{1}{p_1})}} n^s \\ &\geq (2 \log p_1)^{\frac{1}{2 \log(\frac{1}{p_1})} \log n} n^s \\ &\geq n^{s-\frac{1}{2}}. \end{aligned}$$

Let $Z \sim \text{Bin}(n^s, p_1^{m_n})$. We have $\mathbb{E}\left[N_{\text{close}}(x) \mid |B(x, r_n) \cap S_n| > n^s\right] \geq \mathbb{E}[Z] = n^{s-\frac{1}{2}}$. In addition, for each $x' \in A_{\text{close}}(x)$ we have $\mathbb{P}(g_n(x) = g_n(x')) \geq p_1^{m_n}$, and thus, invoking the Chernoff bound,

$$\begin{aligned} (***) &\leq \\ &\mathbb{P}(Z < \frac{1}{2}n^{s-\frac{1}{2}}) \\ &= \mathbb{P}(Z < \frac{1}{2}\mathbb{E}[Z]) \\ &\leq \exp(-\frac{1}{8}\mathbb{E}[Z]) \\ &\leq \exp(-\frac{1}{8}n^{s-\frac{1}{2}}). \end{aligned}$$

The last term we have to bound is the second term in (20). Notice that

$$\begin{aligned} \{N_{\text{far}}(x) > \delta N(x), N_{\text{close}}(x) \geq \frac{1}{2}n^{s-\frac{1}{2}}\} \\ \implies \{N_{\text{far}}(x) > \frac{\delta}{2}n^{s-\frac{1}{2}}\}. \end{aligned}$$

In addition, since $p_1^2 > p_2$, we have

$$\begin{aligned} \mathbb{E}[N_{\text{far}}(x)] &\leq p_2^{m_n} n \leq p_1^{2m_n} n \leq p_1^{2\left(\frac{\log n}{2 \log \frac{1}{p_1}} - 1\right)} n \\ &= p_1^{-2} = O(1). \end{aligned}$$

Since for each $x' \in A_{\text{far}}(x)$ we have $\mathbb{P}(g_n(x) = g_n(x')) \leq p_2^{m_n}$, if we let $Z \sim \text{Bin}(n, p_2^{m_n})$ then, by Chernoff's bound,

$$(\ast \ast \ast) \leq \mathbb{P}(Z > \frac{\delta}{2}n^{s-\frac{1}{2}}) \leq 2^{\frac{\delta}{2}n^{s-\frac{1}{2}}}.$$

For $s > \frac{1}{2}$ and large enough n s.t. $2e\mathbb{E}[N_{\text{far}}(x)] \leq 2e\mathbb{E}[Z] \leq 2e \leq \frac{\delta}{2}n^{s-\frac{1}{2}}$. □

Finally, setting $\gamma = \sqrt{\frac{1}{n^{\frac{1-s}{2}}}}$, $r_n = \left(\frac{1.6\sqrt{d}^{d+2}}{n^{1-s}}\right)^{\frac{1}{d+1}}$ we conclude our proof. □

B Proof of Lemma 5

Here we show that with high probability, the variable $N(x) \rightarrow \infty$. Namely, the number of sample points at each bucket is increasing as n goes to ∞ .

Lemma 5. *Let $x \sim \mathcal{D}_X$ be a test point. Then, for all $M > 0$, $\frac{1}{2} < s < 1$ the hash table calculated by Algorithm 1 satisfies:*

$$\begin{aligned} \mathbb{P}(N(x) < M) &\leq \\ &\exp\left(-\frac{n^{s-\frac{1}{2}}}{2} + M\right) + 2\exp\left(-0.09n^{\frac{1-s}{2}}\right) \\ &+ \left(\frac{1.6}{n^{1-s}\sqrt{d}^d}\right)^{\frac{1}{d+1}} + \sqrt{\frac{1}{n^{\frac{1-s}{2}}}}. \end{aligned}$$

Where, again, the probability is over $S_n, x \sim \mathcal{D}^{n+1}$ and the choice of the function g_n .

Proof. Fix $M > 0$. Similar to Lemma 4, we have

$$\begin{aligned} \mathbb{P}(N(x) < M) &\leq \\ &\mathbb{P}\left(N(x) < M \mid \right. \\ &\quad \left. L_{\varepsilon, m}(S_n) \leq \frac{1.6}{\varepsilon^d n^{1-s}}, |B(x, r_n) \cap S_n| > n^s\right) \\ &+ 2\exp\left(-0.09n^{\frac{1-s}{2}}\right) + \frac{1.6\sqrt{d}^d}{r_n^d n^{1-s}} + \sqrt{\frac{1}{n^{\frac{1-s}{2}}}}. \end{aligned} \tag{21}$$

We only have to bound the first term in (21). Observe that

$$\{N(x) < M\} \implies \{N_{\text{close}}(x) < M\}$$

and that $\mathbb{E}[N_{\text{close}}(x) \mid |B(x, r_n) \cap S_n| > n^s] \geq \mathbb{E}[Z] = p_1^{m_n} n^s = n^{s-\frac{1}{2}}$. Now for $Z \sim \text{Bin}(n^s, p_1^{m_n})$, if we let $\xi = 1 - \frac{M}{\mathbb{E}[Z]}$, then by Chernoff's bound we have,

$$\begin{aligned} \mathbb{P}\left(N_{\text{close}}(x) < M \mid \right. \\ &\quad \left. L_{\varepsilon, m}(S_n) \leq \frac{1.6}{\varepsilon^d n^{1-s}}, |B(x, r_n) \cap S_n| > n^s\right) \\ &\leq \mathbb{P}(Z < (1-\xi)\mathbb{E}[Z]) \\ &\leq \exp\left(-\frac{\xi^2}{2}\mathbb{E}[Z]\right) \\ &\leq \exp\left(-\frac{n^{s-\frac{1}{2}}}{2} + M\right). \end{aligned}$$