

Supplementary Material for:  
“On the Power and Limits of Distance-Based  
Learning”

Periklis A. Papakonstantinou      Jia Xu      Guang Yang

May 24, 2016

This supplementary material document contains:

1. Full proofs for every statement made in the main body.
2. Additional statements that introduce more details and more techniques for analyzing bi-Lipschitz learning settings.
3. Empirical results from the Machine Translation domain that indicate an underlying geometric structure in human and machine translation maps. The empirical results also show these properties for two “experts”: (1) Google translate and (1) Bing Translate.

## Contents

<b>1</b>	<b>Roadmap to the Proofs</b>	<b>3</b>
<b>2</b>	<b>Choosing between two in the Hamming Cube</b>	<b>3</b>
2.1	Extending to sparse almost uniform distribution. . . . .	6
2.2	Strong bounds for all target concepts and random queries . . .	8
<b>3</b>	<b>Choosing between two in Edit Spaces</b>	<b>11</b>
3.1	Proof of Theorem 2 . . . . .	13
3.2	Proof of Theorem 3: extending to all target concepts and random queries . . . . .	18
3.3	Proof of Main Lemma . . . . .	20

<b>4</b>	<b>Experimental results</b>	<b>23</b>
4.1	Data sets . . . . .	24
4.1.1	Bi-Lipschitz condition (low-distortion) . . . . .	25
4.1.2	Metric separability condition . . . . .	26
4.1.3	Dense neighboring condition . . . . .	28

# 1 Roadmap to the Proofs

The cylindric example is listed on Section 3 of the main paper. This example is artificial but more accessible than the natural settings in Section 2 and Section 3. In this example, we assert that a selector  $C$  constructed with merely metric information is always be worse than an optimal one.

Section 2 contains the bounds for the hamming space  $(\text{GF}[2]^n, \ell_1)$ , which is a previously very well-studied metric structure. In this setting we also achieve a lower bound for the generalization error of the selector  $C$ . Section 2.1 considers the extension when the input is not exactly uniformly distributed, and Section 2.2 extends to the case when the target concept is arbitrarily chosen and the learner uses random queries to construct the selector  $C$ . Lower bounds for  $C$  are proved in both extensions.

Then, we investigate more involved metric structures equipped with edit distance. A similar lower bound as in Section 2 is given as well as the extension to arbitrary target concept and the learner with random queries. Section 3.3 contains the proof of Lemma 7 which is used in Section 3.1 and Section 3.2.

Finally, we support the naturalness of our setting with experimental evidence in Section 4. We put the experimental results only in the appendix because we want to stay focus on our in-depth results on lower bounds.

# 2 Choosing between two in the Hamming Cube

The previous cylindric example is mainly for expository reasons. In this section, we investigate the hamming cube example defined over  $(\text{GF}[2]^n, \ell_1)$ , where the metric structure is more natural and well-studied.

**Theorem 1.** *For every  $\rho \in \mathbb{Z}^{\geq 0}$  and  $q(n)$ , there is a concept class  $\mathfrak{C}$ , two experts  $A_0, A_1 : \mathcal{X} \rightarrow \mathcal{Y}$ , where  $\mathcal{X} = \mathcal{Y} = (\{0, 1\}^{2n}, \ell_1)$ , such that for  $A_0, A_1$  and for every  $r \in \mathfrak{C}$  the bi-Lipschitz condition holds with constant 2, and*

- for every  $r \in \mathfrak{C}$ ,  $\mathcal{R}(A_0) + \mathcal{R}(A_1) = \rho$ ,
- if the input follows uniform distribution  $D = U_{2n}$ , then for every generic selector  $C = C_{A_0, A_1, [r]}$  constructed from  $q$  queries, there is a target concept  $r \in \mathfrak{C}$  satisfying  $\mathcal{R}(C) > \left(\frac{1}{2} - \frac{n^{\rho+1}q(n)}{2^{n+1}}\right) \rho$ .

**Corollary 2.** *Given the conditions of Theorem 1 and  $\rho = \left\lfloor \frac{n - \log q(n)}{\log n} - 3 \right\rfloor$ , there is  $\frac{n^{\rho+1}q(n)}{2^n} \leq \frac{1}{n^2}$ , and hence  $\mathcal{R}(C) > \left(\frac{1}{2} - \frac{1}{2n^2}\right) \rho \approx \rho/2$ . For example,  $q(n)$  can be  $\text{poly}(n)$  or even  $n^{\log n}$ .*

*Proof of Theorem 1.* We first introduce the example and then verify the bi-Lipschitz conditions. Finally, we analyze the generalization error of the experts and the generic selector.

**Construction (hamming space example)** For  $x \in \{0, 1\}^n, y \in \{0, 1\}^n$ , let  $A_0$  and  $A_1$  be defined as

$$\begin{aligned} A_0(x, y) &= (x, y) \\ A_1(x, y) &= (x, y + m(\rho)) \end{aligned} \quad (1)$$

where  $m : \{0, 1, \dots, n\} \rightarrow \{0, 1\}^n$  is a mask defined as  $m(t) = (\underbrace{0, \dots, 0}_{n-t}, \underbrace{1, 1, \dots, 1}_t)$ ,

and the computation of  $y + m(\rho)$  is bitwise over  $\{0, 1\}^n$ .

Note that  $q(n)$  upper bounds the number of queries used to construct the selector  $C$ . Let  $Q = \{q_i \mid 1 \leq i \leq q(n)\}$  denote the set of all queries, where  $q_i = (q'_i, q''_i) \in \{0, 1\}^{2n}$  for every  $q_i \in Q$  and  $q'_i, q''_i \in \{0, 1\}^n$ . For completeness, let  $Q' = \{q'_i \mid \exists q_i \in Q, q_i = (q'_i, q''_i)\}$  and  $Q'' = \{q''_i \mid \exists q_i \in Q, q_i = (q'_i, q''_i)\}$ .

We also construct the concept class  $\mathfrak{C} = \{r_0\} \cup \{r_Q \mid Q \subseteq \{0, 1\}^{2n}, |Q| \leq q(n)\}$ , where  $r_0$  denotes the identity mapping, and for every  $Q$ ,  $r_Q$  is defined as

$$r_Q(x, y) = (x, y + m(w_Q(x)))$$

where  $m(t) = (\underbrace{0, \dots, 0}_{n-t}, \underbrace{1, 1, \dots, 1}_t)$  as in (1), and  $w_Q(x) = \min \left\{ \min_{q'_i \in Q'} \{d(x, q'_i)\}, \rho \right\}$

for  $q'_i$  being the first half of  $q_i$  as before. By definition,  $r_Q(x, y) = r_0(x, y)$  immediately holds if  $x \in Q'$ , and in particular for every  $(x, y) \in Q$ .

**Metric properties of  $A_0, A_1$  and  $\mathfrak{C}$**  Bi-Lipschitz condition is trivial for  $A_0, A_1$  and  $r_0$ , since both their contraction and expansion are exactly 1. For  $r_Q \in \mathfrak{C}$ , the expansion of  $r_Q$  is bounded as follows,

$$\begin{aligned} & d(r_Q(x, y), r_Q(x', y')) \\ &= |x - x'| + |y + m(w_Q(x)) - y' - m(w_Q(x'))| \\ &\leq |x - x'| + |y - y'| + |w_Q(x) - w_Q(x')| \\ &\leq |x - x'| + |y - y'| + |x - x'| \\ &\leq 2d((x, y), (x', y')) \end{aligned}$$

On the other hand, noticing  $r_Q(r_Q(x, y)) = (x, y)$ , the contraction is also bounded by 2 since  $d((x, y), (x', y')) \leq 2d(r_Q(x, y), r_Q(x', y'))$ . Therefore,  $r_Q$  has bi-Lipschitz constant at most 2.

**Generalization error of  $A_0, A_1$**  For  $r = r_0$ , trivially  $\mathcal{R}(A_0) = 0$ ,  $\mathcal{R}(A_1) = \rho$ , and the conclusion follows. We prove that for arbitrary input distribution and every  $r_Q \in \mathfrak{C}$ ,  $\mathcal{R}(A_0) + \mathcal{R}(A_1) = \rho$ .

$$\begin{aligned}
\mathcal{R}(A_0) + \mathcal{R}(A_1) &= \mathbb{E} \left[ d(A_0(x, y), r_Q(x, y)) + d(A_1(x, y), r_Q(x, y)) \right] \\
&= \mathbb{E} \left[ d\left( (x, y), (x, y + m(w_Q(x))) \right) + d\left( (x, y + m(\rho)), (x, y + m(w_Q(x))) \right) \right] \\
&= \mathbb{E} \left[ \|m(w_Q(x))\|_1 + \|m(\rho) - m(w_Q(x))\|_1 \right] \\
&= \mathbb{E} \left[ w_Q(x) + (\rho - w_Q(x)) \right] = \mathbb{E}[\rho] = \rho
\end{aligned}$$

**Lower bound for  $\mathcal{R}(C)$**  We lower bound the generalization error  $\mathcal{R}(C)$ . The intuition is that as long as  $Q$  contains all the queries used to construct  $C$ ,  $C$  cannot distinguish two concepts  $r_0$  and  $r_Q$ , neither can it be close to both of them. Thus,  $\mathcal{R}(C)$  is lower bounded.

The generalization distance of  $r_0$  and  $r_Q$ , under uniform distribution, is bounded as

$$\begin{aligned}
d(r_0, r_Q) &= \mathbb{E}_{(x, y) \sim U_{2n}} \left[ d(r_0(x, y), r_Q(x, y)) \right] \\
&= \frac{1}{2^{2n}} \sum_{(x, y) \in \{0, 1\}^{2n}} \|m(w_Q(x))\|_1 \\
&= \frac{1}{2^{2n}} \sum_{(x, y) \in \{0, 1\}^{2n}} w_Q(x) \\
&= \frac{1}{2^n} \sum_{x \in \{0, 1\}^n} w_Q(x) = \mathbb{E}_{x \sim U_n} [w_Q(x)]
\end{aligned}$$

Recalling that  $\mathcal{B}_z(\rho) \stackrel{\text{def}}{=} \{x \mid d(x, z) \leq \rho\}$  denotes the closed *ball of radius  $\rho$  centered at  $z$* , if  $x \notin \bigcup_{q'_i \in Q'} \mathcal{B}_{q'_i}(\rho)$ , then by definition of  $w_Q$ , it follows

$$w_Q(x) = \min \left\{ \min_{q'_i \in Q'} \{d(x, q'_i)\}, \rho \right\} = \rho$$

and hence  $r_Q(x, y) = (x, y + m(w_Q(x))) = (x, y + m(\rho)) = A_1(x, y)$  when  $x \notin \bigcup_{q'_i \in Q'} \mathcal{B}_{q'_i}(\rho)$ . Since  $|Q'| \leq |Q| \leq q(n)$  and  $|\mathcal{B}_{q'_i}(\rho)| \leq \sum_{j=0}^{\lfloor \rho \rfloor} \binom{n}{j} < n^{\rho+1}$ ,

the volume of  $\bigcup_{q'_i \in Q'} \mathcal{B}_{q'_i}(\rho)$  can be upper bounded by

$$\left| \bigcup_{q'_i \in Q'} \mathcal{B}_{q'_i}(\rho) \right| < n^{\rho+1} q(n), \quad (2)$$

Consequentially, we have

$$\begin{aligned} |\{x \in \{0, 1\}^n \mid w_Q(x) = \rho\}| &= 2^n - \left| \left\{ x \in \{0, 1\}^n \mid \min_{q'_i \in Q'} \{d(x, q'_i)\} < \rho \right\} \right| \\ &\geq 2^n - \left| \bigcup_{q'_i \in Q'} \mathcal{B}_{q'_i}(\rho) \right| > 2^n - n^{\rho+1} q(n) \end{aligned}$$

Therefore,

$$\begin{aligned} \mathbb{E}_{x \sim U_n} [w_Q(x)] &\geq \frac{|\{x \in \{0, 1\}^n \mid w_Q(x) = \rho\}|}{2^n} \cdot \rho \\ &> \frac{(2^n - n^{\rho+1} q(n)) \rho}{2^n} = \left(1 - \frac{n^{\rho+1} q(n)}{2^n}\right) \rho \end{aligned}$$

Plugging the above inequality into the bound for  $d(r_0, r_Q)$ , there is

$$d(r_0, r_Q) = \mathbb{E}_{x \sim U_n} [w_Q(x)] > \left(1 - \frac{n^{\rho+1} q(n)}{2^n}\right) \rho \quad (3)$$

For  $A_0, A_1$  as in (1), recalling that  $Q$  contains all queries made by  $\mathcal{C}$ , then the selector  $\mathcal{C}$  cannot distinguish the concept  $r_0$  from  $r_Q$ , since these two concepts coincide on every query in  $Q$ . Therefore, there is  $r \in \mathfrak{C}$  (in particular,  $r \in \{r_0, r_Q\}$ ) such that

$$\mathcal{R}(\mathcal{C}) \geq d(r_0, r_Q)/2 > \left(\frac{1}{2} - \frac{n^{\rho+1} q(n)}{2^{n+1}}\right) \rho$$

□

## 2.1 Extending to sparse almost uniform distribution.

In Theorem 1 and Theorem 1, the generalization errors are studied under the uniform distribution. Formally, to prove a PAC lower bound it is sufficient to exhibit a single example of a distribution. However, we would like to obtain even stronger bounds, aiming to understand better more realistic situations. We consider the following sparse almost uniform distribution as an extension.

**Definition 3** (sparse almost uniform distribution). A distribution  $D$  over  $\{0, 1\}^{2n}$  is called a  $k$ -sparse almost uniform distribution if for every  $(x, y) \sim D$  where  $|x| = |y|$ ,  $y$  is independent from  $x$  and distributes uniformly,  $\max_{s \in \{0, 1\}^n} \Pr_{(x, y) \sim D}[x = s] \leq 2^{-k}$ .

**Theorem 4.** For every  $\rho \in \mathbb{Z}^{\geq 0}$  and  $q(n)$ , there is a concept class  $\mathfrak{C}$  and two experts  $A_0, A_1 : \mathcal{X} \rightarrow \mathcal{Y}$  where  $\mathcal{X} = \mathcal{Y} = (\{0, 1\}^{2n}, \ell_1)$  for  $A_0, A_1$  and for every  $r \in \mathfrak{C}$  the bi-Lipschitz condition holds with constant 2, such that for every  $k$ -sparse almost uniform distribution  $D$ ,

- for every  $r \in \mathfrak{C}$ ,  $\mathcal{R}(A_0) + \mathcal{R}(A_1) = \rho$ ,
- for every selector  $C = C_{A_0, A_1, [r]}$  constructed from  $q$  queries, there is a target concept  $r \in \mathfrak{C}$  satisfying  $\mathcal{R}(C) > \left(\frac{1}{2} - \frac{n^{\rho+1}q(n)}{2^{k+1}}\right) \rho$ .

*Proof.* Let us consider exactly the same construction as in Theorem 1 (and the bi-Lipschitz condition automatically follows). For  $x \in \{0, 1\}^n, y \in \{0, 1\}^n$ ,

$$\begin{aligned} A_0(x, y) &= (x, y) \\ A_1(x, y) &= (x, y + m(\rho)) \end{aligned}$$

The concept class is  $\mathfrak{C} = \{r_0\} \cup \{r_Q \mid Q \subseteq \{0, 1\}^{2n}, |Q| \leq q(n)\}$ , where  $r_0$  denotes the identity mapping, and for every  $Q$ ,  $r_Q$  is defined as

$$r_Q(x, y) = (x, y + m(w_Q(x)))$$

The masking function  $m : \{0, 1, \dots, n\} \rightarrow \{0, 1\}^n$  is defined as before,  $m(t) = (\underbrace{0, \dots, 0}_{n-t}, \underbrace{1, 1, \dots, 1}_t)$ .  $Q = \{q_i \mid 1 \leq i \leq q(n)\}$  and  $Q' = \{q'_i \mid \exists q_i \in Q, q_i = (q'_i, q''_i)\}$ ,

$Q'' = \{q''_i \mid \exists q_i \in Q, q_i = (q'_i, q''_i)\}$ . And  $w_Q(x) = \min \left\{ \min_{q'_i \in Q'} \{d(x, q'_i)\}, \rho \right\}$  for  $q'_i$  being the first half of  $q_i \in Q$ . In particular,  $w_0(x) = 0$  for every  $x \in \{0, 1\}^n$ .

**Generalization error of  $A_0, A_1$ .** When  $r = r_0$ , it turns out  $\mathcal{R}(A_0) = 0$  and  $\mathcal{R}(A_1) = \rho$ , and the conclusion follows. For every  $r_Q \in \mathfrak{C}$ ,  $\mathcal{R}(A_0) + \mathcal{R}(A_1) = \rho$ ,

with exactly the same argument as in Theorem 1.

$$\begin{aligned}
\mathcal{R}(\mathbf{A}_0) + \mathcal{R}(\mathbf{A}_1) &= \mathbb{E}_{(x,y) \sim \mathbf{D}} [d(\mathbf{A}_0(x,y), r_Q(x,y)) + d(\mathbf{A}_1(x,y), r_Q(x,y))] \\
&= \mathbb{E}_{(x,y) \sim \mathbf{D}} \left[ d\left((x,y), (x,y + m(w_Q(x)))\right) + d\left((x,y + m(\rho)), (x,y + m(w_Q(x)))\right) \right] \\
&= \mathbb{E}_{(x,y) \sim \mathbf{D}} [\|m(w_Q(x))\|_1 + \|m(\rho) - m(w_Q(x))\|_1] \\
&= \mathbb{E}_{(x,y) \sim \mathbf{D}} [w_Q(x) + (\rho - w_Q(x))] = \mathbb{E}_{(x,y) \sim \mathbf{D}} [\rho] = \rho
\end{aligned}$$

**Upper bound for the influence of “Beacons”.** It is clear that the set  $Q$  affects the value of  $r_Q(x,y)$  only when  $x \in \bigcup_{q' \in Q'} \mathcal{B}_{q'}(\rho)$ , since otherwise  $w_Q(x) = \rho$  and  $r_Q = \mathbf{A}_1$ . Let  $\mathcal{B} = \mathcal{B}(Q, \rho) \stackrel{\text{def}}{=} \bigcup_{q' \in Q'} \mathcal{B}_{q'}(\rho)$ . We bound the probability that  $x \in \mathcal{B}$ , when  $(x,y) \sim \mathbf{D}$ . Recalling that  $|\mathcal{B}| < n^{\rho+1}q(n)$  as in (2) and  $\mathbf{D}$  is  $k$ -sparse almost uniform distribution,

$$\Pr_{(x,y) \sim \mathbf{D}} [x \in \mathcal{B}] \leq |\mathcal{B}| \cdot \max_{s \in \mathcal{B}} \Pr_{(x,y) \sim \mathbf{D}} [x = s] \leq n^{\rho+1}q(n) \cdot 2^{-k} \quad (4)$$

**Lower bound for  $\mathcal{R}(\mathbf{C})$**  Note that  $w_Q(x) = \rho$  for  $x \notin \mathcal{B}$ , thus we lower bound the distance of  $r_0$  and  $r_Q$  as follows:

$$\begin{aligned}
d(r_0, r_Q) &= \mathbb{E}_{(x,y) \sim \mathbf{D}} [d(r_0(x,y), r_Q(x,y))] \\
&= \mathbb{E}_{(x,y) \sim \mathbf{D}} [w_Q(x)] \\
&\geq \rho \cdot \Pr_{(x,y) \sim \mathbf{D}} [x \notin \mathcal{B}] \geq (1 - n^{\rho+1}q(n) \cdot 2^{-k})\rho
\end{aligned}$$

Therefore, there exists  $r \in \mathfrak{C}$  (in particular,  $r \in \{r_0, r_Q\}$ ) such that

$$\mathcal{R}(\mathbf{C}) \geq d(r_0, r_Q)/2 > \left( \frac{1}{2} - \frac{n^{\rho+1}q(n)}{2^{k+1}} \right) \rho$$

□

## 2.2 Strong bounds for all target concepts and random queries

In Theorem 1 we fix the two experts and consider a target concept chosen in an adversarial way with respect to the  $q$  queries from which  $\mathbf{C}$  is constructed.



Now, we consider the extension when  $\mathcal{C}$  is constructed from random queries. The following theorem can be realized as a dual to Theorem 1.

**Theorem 5.** *For every  $\rho \in \mathbb{Z}^{\geq 0}$ , polynomial  $q(n)$ , and every target concept  $r$  with bi-Lipschitz constant  $c$ , there exists a set  $\mathcal{A}$  of experts defined over  $\mathcal{X} \rightarrow \mathcal{Y}$  where  $\mathcal{X} = \mathcal{Y} = (\{0, 1\}^{2n}, \ell_1)$  and all experts in  $\mathcal{A}$  have bi-Lipschitz constant  $\leq 2c$ , such that when the instances follow the uniform distribution  $\mathcal{D} = U_{2n}$ , then for every  $q$ -queries generic selector  $\mathcal{C} = \mathcal{C}_{\mathcal{A}_0, \mathcal{A}_1, [r]}$ , there are two experts  $\mathcal{A}_0, \mathcal{A}_1 \in \mathcal{A}$  such that  $\min \{\mathcal{R}(\mathcal{A}_0), \mathcal{R}(\mathcal{A}_1)\} = 0$  but*

$$\mathcal{R}(\mathcal{C}) > \left( \frac{1}{2} - \frac{n^{\rho+1}q(n)}{2^{n+1}} \right) \frac{\rho}{c}$$

*Proof.* We begin with the construction of  $\mathcal{A}$  with respect to  $r$ . Then, we verify its metric properties, and finally we will analyze the lower bound.

**Construction** Recalling that  $q(n)$  upper bounds the number of queries used to construct the selector  $\mathcal{C}$ , let  $Q = \{q_i | 1 \leq i \leq q(n)\}$  denote the set of all queries, where  $q_i = (q'_i, q''_i) \in \{0, 1\}^{2n}$  for every  $q_i \in Q$  and  $q'_i, q''_i \in \{0, 1\}^n$ . For completeness, let  $Q' = \{q'_i | \exists q_i \in Q, q_i = (q'_i, q''_i)\}$  and  $Q'' = \{q''_i | \exists q_i \in Q, q_i = (q'_i, q''_i)\}$ .

For  $x \in \{0, 1\}^n, y \in \{0, 1\}^n$ , define  $\mathcal{A}_Q$  as

$$\mathcal{A}_Q(x, y) = r(x, y + m(w_Q(x))) \tag{5}$$

where  $m : \{0, 1, \dots, n\} \rightarrow \{0, 1\}^n$  is a mask defined as  $m(t) = (\underbrace{0, \dots, 0}_{n-t}, \underbrace{1, 1, \dots, 1}_t)$ ,

and  $w_Q(x) = \min \{d(x, Q'), \rho\} = \min \left\{ \min_{q'_i \in Q'} \{d(x, q'_i)\}, \rho \right\}$ , and the addition is bitwise over  $\{0, 1\}^n$ .

By definition,  $w_Q(x) = 0$  and immediately  $\mathcal{A}_Q(x, y) = r(x, y)$  holds if  $x \in Q'$ , and in particular when  $(x, y) \in Q$ . Intuitively,  $\mathcal{A}_Q$  is close to  $r$  on input  $(x, y)$  only when  $x$  is close to some element in the set  $Q'$ . When  $x$  is far away from  $Q'$ , i.e.  $w_Q(x) = \rho$ , the distance of  $\mathcal{A}_Q$  and  $r$  is bounded by  $d(\mathcal{A}_Q(x, y), r(x, y)) \geq d(y + m(w_Q(x)), y)/c = \rho/c$  and similarly  $d(\mathcal{A}_Q(x, y), r(x, y)) \leq \rho c$ .

Then, we construct the class  $\mathcal{A} = \{r\} \cup \left\{ \mathcal{A}_Q | Q \subseteq \{0, 1\}^{2n}, |Q| \leq q(n) \right\}$

**Metric properties of  $\mathcal{A}$**  Recalling that  $r$  has bi-Lipschitz constant  $c$ , we consider the bi-Lipschitz condition of  $\mathcal{A}_Q \in \mathcal{A}$ . For the expansion of  $\mathcal{A}_Q$ , we

have

$$\begin{aligned}
& d(\mathbf{A}_Q(x, y), \mathbf{A}_Q(x', y')) \\
&= d(\mathbf{r}(x, y + m(w_Q(x))), \mathbf{r}(x', y' + m(w_Q(x')))) \\
&\leq c \cdot d((x, y + m(w_Q(x))), (x', y' + m(w_Q(x')))) \\
&= c (\|x - x'\|_1 + \|y + m(w_Q(x)) - y' - m(w_Q(x'))\|_1) \\
&\leq c \cdot (d((x, y), (x', y')) + \|w_Q(x) - w_Q(x')\|_1) \\
&\leq c \cdot d((x, y), (x', y')) + c \cdot d(x, x') \\
&\leq 2c \cdot d((x, y), (x', y'))
\end{aligned}$$

On the other hand, the contraction of  $\mathbf{A}_Q$  is also bounded by  $2c$ ,

$$\begin{aligned}
& d(\mathbf{A}_Q(x, y), \mathbf{A}_Q(x', y')) \\
&= d(\mathbf{r}(x, y + m(w_Q(x))), \mathbf{r}(x', y' + m(w_Q(x')))) \\
&\geq c^{-1} \cdot d((x, y + m(w_Q(x))), (x', y' + m(w_Q(x')))) \\
&= c^{-1} (\|x - x'\|_1 + \|y + m(w_Q(x)) - y' - m(w_Q(x'))\|_1) \\
&\geq c^{-1} \cdot \max \{ \|x - x'\|_1, \|x - x'\|_1 + \|y - y'\|_1 - \|w_Q(x) - w_Q(x')\|_1 \} \\
&\geq c^{-1} \cdot \max \{ \|x - x'\|_1, \|y - y'\|_1 \} \\
&\geq c^{-1} \cdot \frac{\|x - x'\|_1 + \|y - y'\|_1}{2} \\
&= \frac{d((x, y), (x', y'))}{2c}
\end{aligned}$$

Thus, the bi-Lipschitz condition follows and  $\mathbf{A}_Q$  has bi-Lipschitz constant  $2c$ .

**Lower bound for  $\mathcal{R}(\mathbf{C})$**  Let  $Q$  be the set of all queried points made in the construction of  $\mathbf{C}$ . Randomly set  $b \in \{0, 1\}$ , then let  $\mathbf{A}_b = \mathbf{r}$  and  $\mathbf{A}_{1-b} = \mathbf{A}_Q$ . Since  $\mathbf{A}_b(x, y) = \mathbf{A}_{1-b}(x, y)$  for every  $(x, y) \in Q$ , and moreover  $\mathbf{A}_b(x, y) = \mathbf{A}_{1-b}(x, y + m(w_Q(x)))$ ,  $\mathbf{A}_0$  and  $\mathbf{A}_1$  are symmetric and hence indistinguishable to  $\mathbf{C}$ . Therefore  $\mathcal{R}(\mathbf{C}) \geq (\mathcal{R}(\mathbf{A}_0) + \mathcal{R}(\mathbf{A}_1))/2 = \mathcal{R}(\mathbf{A}_Q)/2$ . By the same calculation as in the proof of Theorem 1, we get the lower bound  $\mathcal{R}(\mathbf{A}_Q) > \left(1 - \frac{n^{\rho+1}q(n)}{2^n}\right) \rho/c$  and as a result,

$$\mathcal{R}(\mathbf{C}) > \left(\frac{1}{2} - \frac{n^{\rho+1}q(n)}{2^{n+1}}\right) \frac{\rho}{c}$$

□

### 3 Choosing between two in Edit Spaces

Natural metric spaces lack the nice properties of the previous cylindrical example. Our main result regards edit spaces. The technical obstacle is to give an explicit construction which deals with isoperimetric phenomena on spaces whose underlying graphs have high expansion. Roughly speaking this prevents constructing examples by simply gluing together independent copies (as we did for the cylinders).

In this section, we first introduce Theorem 6 and Theorem 9. We also exposit Corollary 8 for a better understanding of the theorem. Then, we provide the proof of Theorem 6 in Section 3.1. Section 3.2 contains the proof of Theorem 9 and Section 3.3 provides the proof to Lemma 7.

**Theorem 6.** *For every integer  $\rho \leq n$  and  $q$ , there is a concept class  $\mathfrak{C}$ , two experts  $A_0, A_1 : \mathcal{X} \rightarrow \mathcal{Y}$  (same experts for every concept), where  $\mathcal{X} = \mathcal{Y} = (\{0, 1\}^{2n}, \text{ed})$ , such that the following holds true for uniform  $D = U_{2n}$ . For every learner with expert advice that makes  $q$  many queries there exists  $r \in \mathfrak{C}$  such that for the constructed selector  $\mathcal{C}$  we have  $\mathcal{R}(\mathcal{C}) > \left(\frac{1}{2} - \frac{(4n+4\rho)^\rho q(n)}{2^{n+1}}\right) \Omega(\sqrt{\rho})$ , while the optimal selector has generalization error  $\leq \frac{(4n+4\rho)^\rho q(n)}{2^n} \cdot \rho$ . Furthermore, for  $A_0, A_1$  and for every  $r \in \mathfrak{C}$  the bi-Lipschitz condition holds with constant less than 5.*

*Proof sketch of Theorem 6.* Consider two experts  $A_0, A_1$  and the concept class  $\mathfrak{C}$ :

For  $x \in \{0, 1\}^n, y \in \{0, 1\}^n$ ,

$$\begin{aligned} A_0(x, y) &= (x, y) \\ A_1(x, y) &= (x, y + m(\rho)) \end{aligned} \tag{6}$$

where the computation of  $y + m(\rho)$  is bitwise over  $\{0, 1\}^n$ , and  $m : \{0, 1, \dots, n\} \rightarrow \{0, 1\}^n$  is a mask defined as  $m(t) = \underbrace{(0, \dots, 0)}_{n-t}, \underbrace{(1, 1, \dots, 1)}_t$ . Let

$$\mathfrak{C} = \{r_0\} \cup \left\{ r_Q \mid Q \subseteq \{0, 1\}^{2n}, |Q| \leq q(n) \right\}$$

where  $r_0 = A_0$  is identity and  $r_Q(x, y) = (x, y + m(w_Q(x)))$ , where  $w_Q(x) = \min \left\{ \min_{q'_i \in Q'} \{ \text{ed}(x, q'_i) \}, \rho \right\}$ , for  $Q' = \{q'_i \mid \exists q_i \in Q, q_i = (q'_i, q''_i)\}$ .

**Generalization error of  $A_0, A_1$ .** For every  $r \in \mathfrak{C}$ ,  $\mathcal{R}(A_0) + \mathcal{R}(A_1) \geq \mathbb{E} [\text{ed}(y, y + m(\rho))] = \Omega(\sqrt{\rho})$  by the following lemma. This lemma is proved

by calculating the probability that  $\rho$  uniformly chosen bits have at least  $(\rho + \sqrt{\rho})/2$  many 0's, and the full proof is deferred to Section 3.3.

**Lemma 7.** *For every  $\rho \leq n$ , and function  $m$  as in (6) we have*

$$\mathbb{E}_{y \sim U_n} [\text{ed}(y, y + m(\rho))] \geq \Omega(\sqrt{\rho})$$

**Upper bound for the influence of “Beacons”.** Suppose  $Q$  consists of all queries instances, which are the “beacons”. By definition,  $r_Q(x, y) = r_0(x, y)$  only when  $x \in Q'$ , and the distance of  $r_0$  and  $r_Q$  gradually grows to  $\rho$  as the input moves away from “beacons”. Since each beacon influences a ball of radius  $\rho$  which has size at most  $(4(n + \rho))^\rho$ , the influence of all beacons can be upper bounded by the union of balls.

$$\left| \bigcup_{q'_i \in Q'} \mathcal{B}_{q'_i}(\rho) \right| < (4n + 4\rho)^\rho q(n) \quad (7)$$

**Lower bound for  $\mathcal{R}(C)$**  By inequality (7) and Lemma 7,

$$\begin{aligned} \text{ed}(r_0, r_Q) &= \mathbb{E}_{(x, y) \sim U_{2n}} [\text{ed}(r_0(x, y), r_Q(x, y))] \\ &= \frac{1}{2^{2n}} \sum_{(x, y) \in \{0, 1\}^{2n}} \text{ed}((x, y), (x, y + m(w_Q(x)))) \\ &\geq \frac{1}{2^n} \sum_{\substack{x \in \{0, 1\}^n \\ w_Q(x) = \rho}} \left( \frac{1}{2^n} \sum_{y \in \{0, 1\}^n} \text{ed}((x, y), (x, y + m(w_Q(x)))) \right) \\ &> \frac{2^n - (4n + 4\rho)^\rho q(n)}{2^n} \cdot \mathbb{E}_{y \sim U_n} [\text{ed}(y, y + m(\rho))] \geq \left( 1 - \frac{(4n + 4\rho)^\rho q(n)}{2^n} \right) \Omega(\sqrt{\rho}) \end{aligned}$$

Since  $C$  cannot distinguish  $r_0$  from  $r_Q$ , there exists  $r \in \mathfrak{C}$  (in particular,  $r \in \{r_0, r_Q\}$ ) such that

$$\mathcal{R}(C) \geq \text{ed}(r_0, r_Q)/2 > \left( \frac{1}{2} - \frac{(4n + 4\rho)^\rho q(n)}{2^{n+1}} \right) \Omega(\sqrt{\rho})$$

The proof for the bi-Lipschitz condition (see the appendix) is technical but not difficult, which can be obtained by analyzing the optimal matching function.  $\square$

**Corollary 8.** *Given the conditions of Theorem 6 for  $q = q(n) \leq n^{\log n}$  and  $\rho \leq \left\lfloor \frac{n - \log q(n)}{3 + \log n} - 3 \right\rfloor = \frac{n}{3 + \log n} - \log n - O(1)$ , then  $\frac{(4n+4\rho)^\rho q(n)}{2^n} < \frac{1}{n^3}$ . Thus, the optimal selector has generalization error  $\min \{\mathcal{R}(A_0), \mathcal{R}(A_1)\} \leq \frac{(4n+4\rho)^\rho q(n)}{2^n} \cdot \rho < \frac{\rho}{n^3}$ , while the generalization error  $\mathcal{R}(C) > \left(\frac{1}{2} - \frac{1}{2n^3}\right) \Omega(\sqrt{\rho}) = \Omega(\sqrt{\rho})$  for every generic selector from  $q$  queries.*

Theorem 9 shares significant parts in its proof from that of Theorem 6. It also has a similar statement as the above one, where instead of a generic selector we have a *random samples generic selector* and the lower bound for  $\mathcal{R}(C)$  holds for every concept  $r \in \mathcal{C}$ .

**Theorem 9.** *For every  $\rho \geq 0$ ,  $q(n) \geq 0$ , and target concept  $r$  with distortion constant  $c$ , there exists a set  $\mathcal{A}$  of bi-Lipschitz experts defined over  $\mathcal{X} \rightarrow \mathcal{Y}$  for the edit space  $\mathcal{X} = \mathcal{Y} = (\{0, 1\}^{2n}, \text{ed})$ , such that for uniformly chosen input  $(x, y) \sim D = U_{2n}$  the following holds true. For every learner with black-box access that makes  $q$ -many queries to  $r$  and  $q$ -many queries to the two experts, there exist two experts  $A_0, A_1 \in \mathcal{A}$  with distortion  $\leq 5c$  and the optimal selector has 0 generalization error. Then, such a learner constructs  $C$  where  $\mathcal{R}(C) \geq \left(\frac{1}{2} - \frac{(4n+4\rho)^\rho q(n)}{2^{n+1}}\right) \cdot \Omega\left(\frac{\sqrt{\rho}}{c}\right)$ . In particular,  $\mathcal{R}(C) \geq \Omega\left(\frac{\sqrt{\rho}}{c}\right)$  when  $\rho \leq \left\lfloor \frac{n - \log q(n)}{3 + \log n} - 1 \right\rfloor$ .*

### 3.1 Proof of Theorem 2

*Proof of Theorem 6.* Similar to the proof of Theorem 1, we first introduce the example, then analyze the generalization error of the experts and the generic selector, and finally verify the bi-Lipschitz conditions.

**Construction (family of examples)** For  $x \in \{0, 1\}^n, y \in \{0, 1\}^n$ , let  $A_0$  and  $A_1$  be defined as follows

$$\begin{aligned} A_0(x, y) &= (x, y) \\ A_1(x, y) &= (x, y + m(\rho)) \end{aligned} \tag{8}$$

where  $m : \{0, 1, \dots, n\} \rightarrow \{0, 1\}^n$  is a mask  $m(t) = (\underbrace{0, \dots, 0}_{n-t}, \underbrace{1, 1, \dots, 1}_t)$ ,

and the computation of  $y + m(\rho)$  is bitwise over  $\{0, 1\}^n$ .

Note that  $q(n)$  upper bounds the number of queries made in the construction of the selector  $C$ . Let  $Q = \{q_i \mid 1 \leq i \leq q(n)\}$  denote the set of all queries, where  $q_i = (q'_i, q''_i) \in \{0, 1\}^{2n}$  for every  $q_i \in Q$  and  $q'_i, q''_i \in$

$\{0, 1\}^n$ . For completeness, let  $Q' = \{q'_i | \exists q_i \in Q, q_i = (q'_i, q''_i)\}$  and  $Q'' = \{q''_i | \exists q_i \in Q, q_i = (q'_i, q''_i)\}$ .

Then we construct the concept class  $\mathfrak{C} = \{r_0\} \cup \{r_Q | Q \subseteq \{0, 1\}^{2n}, |Q| \leq q(n)\}$ .  $r_0$  denotes the identity mapping, and for every  $Q$ ,  $r_Q$  is defined as

$$r_Q(x, y) = (x, y + m(w_Q(x)))$$

where  $m(t) = (\underbrace{0, \dots, 0}_{n-t}, \underbrace{1, 1, \dots, 1}_t)$  and  $w_Q(x) = \min \left\{ \min_{q'_i \in Q'} \{\text{ed}(x, q'_i)\}, \rho \right\}$

as before. By definition,  $r_Q(x, y) = r_0(x, y)$  if and only if  $x \in Q'$ , which holds in particular for every  $(x, y) \in Q$ .

**Upper bounding the volume of  $\bigcup_{q'_i \in Q'} \mathcal{B}_{q'_i}(\rho)$ .** Let  $\mathcal{B}_z(\rho) \stackrel{\text{def}}{=} \{x \mid \text{ed}(x, z) \leq \rho\}$  denote the closed ball of radius  $\rho$  centered at  $z$ . For  $x \notin \bigcup_{q'_i \in Q'} \mathcal{B}_{q'_i}(\rho)$ , by definition of  $w_Q$ , it follows

$$w_Q(x) = \min \left\{ \min_{q'_i \in Q'} \{\text{ed}(x, q'_i)\}, \rho \right\} = \rho$$

To upper bound the volume of  $\mathcal{B}_z(\rho)$ , we calculate the number of distinct  $\rho$ -step operation sequences from  $z \in \{0, 1\}^n$ . For each of the  $\rho$  steps, we can perform either insertion 0/1 to one of at most  $n + \rho$  positions, or deletion and substitution in one of at most  $n + \rho - 1$  positions, or no operation. Totally, we have less than  $4(n + \rho)$  choices in each step. Thus, we bound  $\mathcal{B}_{q'_i}(\rho)$  as

$$\left| \mathcal{B}_{q'_i}(\rho) \right| < (4(n + \rho))^\rho$$

By union bound, the volume of the union of  $q(n)$  balls is upper bounded by

$$\left| \bigcup_{q'_i \in Q'} \mathcal{B}_{q'_i}(\rho) \right| < (4n + 4\rho)^\rho q(n)$$

As a consequence,

$$\begin{aligned} \left| \{x \in \{0, 1\}^n \mid w_Q(x) = \rho\} \right| &= 2^n - \left| \left\{ x \in \{0, 1\}^n \mid \min_{q'_i \in Q'} \{d(x, q'_i)\} < \rho \right\} \right| \\ &\geq 2^n - \left| \bigcup_{q'_i \in Q'} \mathcal{B}_{q'_i}(\rho) \right| > 2^n - (4n + 4\rho)^\rho q(n) \end{aligned}$$

**Generalization error of  $A_0, A_1$**  We consider the generalization error of an optimal selector for the concept class  $\mathfrak{C}$ . If  $r = r_0$ , then the generalization error is 0 since  $A_0 = r_0$ . Otherwise, if  $r = r_Q$ ,

$$\begin{aligned}
\mathcal{R}(A_1) &= \mathbb{E} [\text{ed}(A_1(x, y), r(x, y))] \\
&= \mathbb{E} \left[ \text{ed} \left( (x, y + m(\rho)), (x, y + m(w_Q(x))) \right) \right] \\
&\leq \mathbb{E} \left[ \text{ed} \left( y + m(\rho), y + m(w_Q(x)) \right) \right] \\
&\leq \Pr_{x \sim \mathcal{D}} [w_Q(x) < \rho] \mathbb{E} \left[ \text{ed} \left( y + m(\rho), y + m(w_Q(x)) \right) \mid w_Q(x) < \rho \right] \\
&\leq \frac{(4n + 4\rho)^\rho q(n)}{2^n} \cdot \rho
\end{aligned}$$

Thus, we conclude

$$\min \{ \mathcal{R}(A_0), \mathcal{R}(A_1) \} \leq \frac{(4n + 4\rho)^\rho q(n)}{2^n} \cdot \rho$$

**Lower bound for  $\mathcal{R}(C)$**  The lower bound of  $\mathcal{R}(C)$  follows the same intuition that  $C$  cannot distinguish two concepts which are far away from each other and on the other hand  $C$  can never be close to both of them.

To lower bounding the generalization distance of  $r_0$  and  $r_Q$ , we first recall that

$$|\{x \in \{0, 1\}^n \mid w_Q(x) = \rho\}| > 2^n - (4n + 4\rho)^\rho q(n)$$

Then using Lemma 7,

$$\begin{aligned}
\text{ed}(r_0, r_Q) &= \mathbb{E}_{(x, y) \sim U_{2n}} [\text{ed}(r_0(x, y), r_Q(x, y))] \\
&= \frac{1}{2^{2n}} \sum_{(x, y) \in \{0, 1\}^{2n}} \text{ed}((x, y), (x, y + m(w_Q(x)))) \\
&\geq \frac{1}{2^n} \sum_{\substack{x \in \{0, 1\}^n \\ w_Q(x) = \rho}} \left( \frac{1}{2^n} \sum_{y \in \{0, 1\}^n} \text{ed}((x, y), (x, y + m(w_Q(x)))) \right) \\
&> \frac{2^n - (4n + 4\rho)^\rho q(n)}{2^n} \mathbb{E}_{y \sim U_n} [\text{ed}(y, y + m(\rho))] \\
&\geq \left( 1 - \frac{(4n + 4\rho)^\rho q(n)}{2^n} \right) \Omega(\sqrt{\rho})
\end{aligned}$$

For  $A_0, A_1$  as in (6), recalling that  $Q$  consists of all queries made by  $C$ , then  $C$  cannot distinguish the concept  $r_0$  from  $r_Q$ , since these two concepts

coincide on every query in  $Q$ . Therefore, there exists  $r \in \mathfrak{C}$  (in particular,  $r \in \{r_0, r_Q\}$ ) such that

$$\mathcal{R}(\mathfrak{C}) \geq \text{ed}(r_0, r_Q)/2 > \left( \frac{1}{2} - \frac{(4n + 4\rho)^\rho q(n)}{2^{n+1}} \right) \Omega(\sqrt{\rho})$$

**Metric properties of  $A_0, A_1$  and  $\mathfrak{C}$**  The bi-Lipschitz condition is trivial for  $A_0$  and  $r_0$  since they are identity. Then, we prove the bi-Lipschitz constant is at most 3 for  $A_1$  and no more than 5 for every  $r_Q \in \mathfrak{C} \setminus \{r_0\}$ .

Firstly, we bound  $\text{ed}(A_1(x, y), A_1(x', y')) = \text{ed}((x, y + m(\rho)), (x', y' + m(\rho)))$  with  $\text{ed}((x, y), (x', y'))$ . Let us consider the optimal matching  $\mathcal{M}((x, y), (x', y'))$  between  $(x, y)$  and  $(x', y')$ , which maps every pair of unchanged bits when editing<sup>1</sup>  $(x, y)$  to  $(x', y')$ . Denote by  $|\mathcal{M}((x, y), (x', y'))|$  the number of matched pairs defined by the matching, then (note the input length is  $2n$ ),

$$2n - |\mathcal{M}((x, y), (x', y'))| \leq \text{ed}((x, y), (x', y')) \leq 4n - 2|\mathcal{M}((x, y), (x', y'))|$$

Similarly,

$$2n - |\mathcal{M}(A_1(x, y), A_1(x', y'))| \leq \text{ed}(A_1(x, y), A_1(x', y')) \leq 4n - 2|\mathcal{M}(A_1(x, y), A_1(x', y'))|$$

Therefore, it suffices to compare  $|\mathcal{M}((x, y), (x', y'))|$  and  $|\mathcal{M}(A_1(x, y), A_1(x', y'))|$ . Realizing the matchings as graphs with each edge corresponding to a pair of matched bits, we let  $\Gamma$  denote the cut in  $\mathcal{M}((x, y), (x', y'))$  at the position  $2n - \rho$ , i.e.  $\Gamma$  consists of all pairs containing one bit from the first  $2n - \rho$  positions and the other from the last  $\rho$  positions.

On one hand, we assert that  $\Gamma$  only contains edges in either of the following two cases: from the first  $2n - \rho$  positions of  $(x, y)$  to the last  $\rho$  of  $(x', y')$ , or from the last  $\rho$  positions of  $(x, y)$  to the first  $2n - \rho$  positions of  $(x', y')$ . Because in the process of editing, the relative order of remaining bits can never be changed. Then, it trivially follows that  $2|\Gamma| \leq \text{ed}((x, y), (x', y'))$ , since each single edge in  $\Gamma$  will cause at least two insertion/deletion operations.

On the other hand, we notice that if a pair in  $\mathcal{M}((x, y), (x', y'))$  is not in  $\Gamma$ , it can also appear in the matching between  $A_1(x, y)$  and  $A_1(x', y')$ . Thus,  $|\mathcal{M}((x, y), (x', y'))| - |\mathcal{M}(A_1(x, y), A_1(x', y'))| \leq |\Gamma|$ .

Combining above two parts,

$$|\mathcal{M}((x, y), (x', y'))| - |\mathcal{M}(A_1(x, y), A_1(x', y'))| \leq |\Gamma| \leq \text{ed}((x, y), (x', y'))/2$$

<sup>1</sup>In case there are multiple ways to edit  $(x, y)$  to  $(x', y')$  with the same number of operations, fix  $\mathcal{M}$  to any of them.



Then, recalling the lower and upper bound for  $\text{ed}((x, y), (x', y'))$  and  $\text{ed}(\mathbf{A}_1(x, y), \mathbf{A}_1(x', y'))$ , the expansion of  $\mathbf{A}_1$  is at most 3.

$$\begin{aligned} & \text{ed}(\mathbf{A}_1(x, y), \mathbf{A}_1(x', y')) \\ & \leq 4n - 2 |\mathcal{M}(\mathbf{A}_1(x, y), \mathbf{A}_1(x', y'))| \\ & \leq 4n - 2 |\mathcal{M}((x, y), (x', y'))| + \text{ed}((x, y), (x', y')) \\ & \leq 3\text{ed}((x, y), (x', y')) \end{aligned}$$

Since  $\mathbf{A}_1(\mathbf{A}_1(x, y)) = (x, y)$ , the contraction of  $\mathbf{A}_1$  follows the same bound, i.e.  $\text{ed}((x, y), (x', y')) \leq 3\text{ed}(\mathbf{A}_1(x, y), \mathbf{A}_1(x', y'))$ . To conclude,  $\mathbf{A}_1$  has bi-Lipschitz constant 3.

For  $r_Q$ , we verify its bi-Lipschitz constant is at most 5.

By definition of  $r_Q$ , we have

$$\text{ed}(r_Q(x, y), r_Q(x', y')) = \text{ed}((x, y + m(w_Q(x))), (x', y' + m(w_Q(x'))))$$

Again, let  $\mathcal{M}((x, y), (x', y'))$  be the optimal matching between  $(x, y)$  and  $(x', y')$ . And similarly define  $\mathcal{M}(r_Q(x, y), r_Q(x', y'))$ , then,

$$\text{ed}(r_Q(x, y), r_Q(x', y')) \leq 4n - 2 |\mathcal{M}(r_Q(x, y), r_Q(x', y'))|$$

Let  $\Gamma$  denote the cut in  $\mathcal{M}((x, y), (x', y'))$  at the position  $2n - w_Q(x)$  of  $(x, y)$  and the position  $2n - w_Q(x')$  of  $(x', y')$ , i.e.  $\Gamma$  consists of all edges in the matching  $\mathcal{M}((x, y), (x', y'))$  that

- either from the first  $2n - w_Q(x)$  positions of  $(x, y)$  to the last  $w_Q(x')$  positions of  $(x', y')$ ,
- or from the last  $w_Q(x)$  positions of  $(x, y)$  to the first  $2n - w_Q(x')$  positions of  $(x', y')$ .

Moreover, all edges in  $\Gamma$  must be in the same case. By counting the number of unmatched bits on both sides of the cut, we can get  $\text{ed}((x, y), (x', y')) \geq 2(|\Gamma| - |w_Q(x) - w_Q(x')|)$ .

On the other hand,  $|\mathcal{M}((x, y), (x', y'))| - |\mathcal{M}(r_Q(x, y), r_Q(x', y'))| \leq |\Gamma|$  since all edges in  $\mathcal{M}((x, y), (x', y')) \setminus \Gamma$  can be simultaneously reserved in a matching from  $r_Q(x, y)$  to  $r_Q(x', y')$ .

As a result,

$$|\mathcal{M}((x, y), (x', y'))| - |\mathcal{M}(r_Q(x, y), r_Q(x', y'))| \leq |\Gamma| \leq \text{ed}((x, y), (x', y'))/2 + |w_Q(x) - w_Q(x')|$$

Plugging it into previous lower and upper bound for  $\text{ed}((x, y), (x', y'))$  and

$\text{ed}(r_Q(x, y), r_Q(x', y'))$ , the expansion of  $r_Q$  is at most 5.

$$\begin{aligned}
& \text{ed}(r_Q(x, y), r_Q(x', y')) \\
& \leq 4n - 2 |\mathcal{M}(r_Q(x, y), r_Q(x', y'))| \\
& \leq 4n - 2 |\mathcal{M}((x, y), (x', y'))| + \text{ed}((x, y), (x', y')) + 2 |w_Q(x) - w_Q(x')| \\
& \leq 3\text{ed}((x, y), (x', y')) + 2 |w_Q(x) - w_Q(x')| \\
& \leq 3\text{ed}((x, y), (x', y')) + 2\text{ed}(x, x') \\
& \leq 5\text{ed}((x, y), (x', y'))
\end{aligned}$$

Since  $r_Q(r_Q(x, y)) = (x, y)$ , the contraction of  $r_Q$  is also bounded by 5. Thus,  $r_Q$  has bi-Lipschitz constant 5.  $\square$

### 3.2 Proof of Theorem 3: extending to all target concepts and random queries

Similarly as with the extension in Section 2.2, we consider the case where the target concept is a fixed mapping with low distortion in edit space. We show that for the fixed target concept there exists a pair of experts such that from which no selector with small generalization error can be constructed.

**Theorem 9 (informally stated).** For every  $\rho \geq 0$ ,  $q(n) \geq 0$ , and target concept  $r$  with distortion constant  $c$ , there exists a set  $\mathcal{A}$  of bi-Lipschitz experts defined over  $\mathcal{X} \rightarrow \mathcal{Y}$  for the edit space  $\mathcal{X} = \mathcal{Y} = (\{0, 1\}^{2n}, \text{ed})$ , such that for uniformly chosen input  $(x, y) \sim D = U_{2n}$  the following holds true. For every learner with black-box access that makes  $q$ -many queries to  $r$  and  $q$ -many queries to the two experts, there exist two experts  $A_0, A_1 \in \mathcal{A}$  with distortion  $\leq 5c$  and the optimal selector has 0 generalization error. Then, such a learner constructs  $C$  where  $\mathcal{R}(C) \geq \left(\frac{1}{2} - \frac{(4n+4\rho)^\rho q(n)}{2^{n+1}}\right) \cdot \Omega\left(\frac{\sqrt{\rho}}{c}\right)$ . In particular,  $\mathcal{R}(C) \geq \Omega\left(\frac{\sqrt{\rho}}{c}\right)$  when  $\rho \leq \left\lfloor \frac{n - \log q(n)}{3 + \log n} - 1 \right\rfloor$ .

*Proof.* We begin with the construction of  $\mathcal{A}$ . Then, we discuss the metric property and the generalization error of experts  $A_Q \in \mathcal{A} \setminus \{r\}$ . Finally, we prove the lower bound for  $\mathcal{R}(C)$ .

**Construction** Recalling that  $q(n)$  upper bounds the number of queries made by the selector  $C$ , let  $Q = \{q_i | 1 \leq i \leq q(n)\} \subseteq \{0, 1\}^{2n}$  denote the set of all queries, and define  $Q' = \{q'_i | \exists q_i \in Q, q_i = (q'_i, q''_i)\} \subseteq \{0, 1\}^n$  to be the set of first half of elements in  $Q$ .

For  $x \in \{0, 1\}^n, y \in \{0, 1\}^n$ , define  $A_Q$  as in (5)

$$A_Q(x, y) = r(x, y + m(w_Q(x)))$$

where  $m : \{0, 1, \dots, n\} \rightarrow \{0, 1\}^n$  is a mask defined as  $m(t) = (\underbrace{0, \dots, 0}_{n-t}, \underbrace{1, 1, \dots, 1}_t)$ ,

and  $w_Q(x) = \min \{d(x, Q'), \rho\} = \min \left\{ \min_{q'_i \in Q'} \{d(x, q'_i)\}, \rho \right\}$ . When  $x \in Q'$ , in particular when  $(x, y) \in Q$ ,  $w_Q(x) = 0$  and immediately  $A_Q(x, y) = r(x, y)$  follows by definition.

The expert class is constructed as  $\mathcal{A} = \{r\} \cup \left\{ A_Q \mid Q \subseteq \{0, 1\}^{2n}, |Q| \leq q(n) \right\}$

**Metric properties of  $A_Q \in \mathcal{A}$**  For the bi-Lipschitz condition of elements in  $\mathcal{A}$ , we need to bound  $\text{ed}(A_Q(x, y), A_Q(x', y'))$ , with  $\text{ed}((x, y), (x', y'))$ . Recalling that  $r$  has bi-Lipschitz constant  $c$  and  $\text{ed}(A_Q(x, y), A_Q(x', y')) = \text{ed}(r(x, y + m(w_Q(x))), r(x', y' + m(w_Q(x'))))$ ,

$$\frac{1}{c} \leq \frac{\text{ed}(r(x, y + m(w_Q(x))), r(x', y' + m(w_Q(x'))))}{\text{ed}((x, y + m(w_Q(x))), (x', y' + m(w_Q(x'))))} \leq c$$

Following the bi-Lipschitz proof of  $r_Q$  in Theorem 6, we get the upper bound for expansion

$$\text{ed}((x, y + m(w_Q(x))), (x', y' + m(w_Q(x')))) \leq 5\text{ed}((x, y), (x', y'))$$

and symmetrically for contraction

$$\text{ed}((x, y), (x', y')) \leq 5\text{ed}((x, y + m(w_Q(x))), (x', y' + m(w_Q(x'))))$$

Therefore, the bi-Lipschitz constant of  $A_Q$  is at most  $5c$ .

**Generalization error of  $A_Q \in \mathcal{A}$**  Before lower bounding the generalization error  $\mathcal{R}(A_Q)$ , we upper bound the fraction of elements close to  $Q'$ . Let  $B = \{x \in \{0, 1\}^n \mid d(x, Q') < \rho\}$  be the set of all elements close to  $Q'$ , then  $B = \bigcup_{q'_i \in Q'} \mathcal{B}_{q'_i}(\rho)$  as in (7) and the volume is upper bounded by

$$|B| = \left| \bigcup_{q'_i \in Q'} \mathcal{B}_{q'_i}(\rho) \right| < (4n + 4\rho)^p q(n)$$

The generalization error of  $A_Q$  is  $\Omega\left(\frac{\sqrt{\rho}}{c}\right)$  as long as  $(x, y) \sim U_{2n}$ . Recalling Lemma 7 that  $\mathbb{E}[\text{ed}(y, y + m(\rho))] = \Omega(\sqrt{\rho})$ , we have

$$\begin{aligned}
\mathcal{R}(A_Q) &= \mathbb{E}[\text{ed}(A_Q(x, y), r(x, y))] \\
&= \mathbb{E}\left[\text{ed}\left(r(x, y + m(w_Q(x))), r(x, y)\right)\right] \\
&> \mathbb{E}\left[\text{ed}\left(r(x, y + m(\rho)), r(x, y)\right) \mid x \notin B\right] \cdot \Pr_{x \sim U_n}[x \notin B] \\
&\geq \mathbb{E}\left[\frac{1}{c} \cdot \text{ed}\left((x, y + m(\rho)), (x, y)\right)\right] \cdot \left(1 - \frac{(4n + 4\rho)^\rho q(n)}{2^n}\right) \\
&\geq \frac{1}{c} \cdot \mathbb{E}[\text{ed}(y, y + m(\rho))] \cdot \left(1 - \frac{(4n + 4\rho)^\rho q(n)}{2^n}\right) \\
&\geq \left(1 - \frac{(4n + 4\rho)^\rho q(n)}{2^n}\right) \cdot \frac{\Omega(\sqrt{\rho})}{c}
\end{aligned}$$

and when  $\rho \leq \left\lfloor \frac{n - \log q(n)}{3 + \log n} - 1 \right\rfloor$ ,

$$\mathcal{R}(A_Q) = \Omega\left(\frac{\sqrt{\rho}}{c}\right)$$

**Lower bound for  $\mathcal{R}(C)$**  Let  $Q$  be the set of all queried points made by  $C$ . Uniformly at random select  $b \in \{0, 1\}$ , then let  $A_b = r$  and  $A_{1-b} = A_Q$ . Trivially we have

$$\min\{\mathcal{R}(A_0), \mathcal{R}(A_1)\} = 0$$

Note that  $A_b(x, y) = A_{1-b}(x, y)$  when  $(x, y) \in Q$ , and moreover  $A_b(x, y) = A_{1-b}(x, y + m(w_Q(x)))$ . Therefore,  $A_0$  and  $A_1$  are indistinguishable to  $C$ , and hence

$$\mathcal{R}(C) \geq \frac{\mathcal{R}(A_0) + \mathcal{R}(A_1)}{2} = \mathcal{R}(A_Q)/2 = \left(1 - \frac{(4n + 4\rho)^\rho q(n)}{2^n}\right) \cdot \Omega\left(\frac{\sqrt{\rho}}{c}\right)$$

□

### 3.3 Proof of Main Lemma

**Lemma 10.** *For every  $\rho \leq n$ , we have the following lower bound for the expected edit distance between a random string and its complement.*

$$\mathbb{E}_{y \sim U_n}[\text{ed}(y, y + m(\rho))] \geq \Omega(\sqrt{\rho})$$

*Proof of Lemma 10.* We lower bound the expected edit distance by calculating the probability the  $y$  has  $\sqrt{\rho}$  more 0's than 1's in its last  $\rho$  bits.

Consider the case when  $y$  has at least  $(\rho + \sqrt{\rho})/2$  bits equal to 0 in the last  $\rho$  bits, then the number of 0's in last  $\rho$  bits of  $y + m(\rho)$  is at most  $(\rho - \sqrt{\rho})/2$ . As a result,

$$\text{ed}(y, y + m(\rho)) \geq |\# \text{ of 0's in } y - \# \text{ of 0's in } y + m(\rho)| \geq \sqrt{\rho}$$

Then, we bound the probability of the above case. In particular, we lower bound it by the probability that in the last  $\rho$  bits of  $y$  there are at least  $(\rho + \sqrt{\rho})/2$  but at most  $\rho/2 + \sqrt{\rho}$  bits equal to 0.

$$\begin{aligned} & \Pr_{y \sim U_n} \left[ (\rho + \sqrt{\rho})/2 \leq (\# \text{ of 0's in last } \rho \text{ bits of } y) \leq \rho/2 + \sqrt{\rho} \right] \\ &= \sum_{i=\rho/2+\sqrt{\rho}/2}^{\rho/2+\sqrt{\rho}} \frac{\binom{\rho}{i} \cdot 2^{n-\rho}}{2^n} \\ &\geq \frac{1}{2^\rho} \cdot \left(\frac{\sqrt{\rho}}{2} + 1\right) \cdot \binom{\rho}{\rho/2 + \sqrt{\rho}} \\ &\geq \frac{\sqrt{\rho}}{2^{\rho+1}} \cdot \frac{\rho!}{(\rho/2 + \sqrt{\rho})!(\rho/2 - \sqrt{\rho})!} \end{aligned} \tag{9}$$

Recalling that by Stirling's Formula (the version due to Robbins 1955),

$$\sqrt{2\pi n}^{n+1/2} \exp\left(-n + \frac{1}{12n+1}\right) \leq n! \leq \sqrt{2\pi n}^{n+1/2} \exp\left(-n + \frac{1}{12n}\right)$$

we derive the lower bound for  $\binom{\rho}{\rho/2+\sqrt{\rho}}$  as

$$\begin{aligned}
& \binom{\rho}{\rho/2+\sqrt{\rho}} \\
&= \frac{\rho!}{(\rho/2+\sqrt{\rho})!(\rho/2-\sqrt{\rho})!} \\
&\geq \frac{\sqrt{2\pi}\rho^{\rho+1/2}\exp(-\rho+\frac{1}{12\rho+1})}{(\sqrt{2\pi}(\rho/2+\sqrt{\rho})^{\rho/2+\sqrt{\rho}+1/2}\exp(-(\rho/2+\sqrt{\rho})+1))} \\
&\quad \frac{1}{(\sqrt{2\pi}(\rho/2-\sqrt{\rho})^{\rho/2-\sqrt{\rho}+1/2}\exp(-(\rho/2-\sqrt{\rho})+1))} \\
&= \frac{\rho^{\rho+1/2}\exp\left(-\rho+\frac{1}{12\rho+1}+(\rho/2+\sqrt{\rho})+(\rho/2-\sqrt{\rho})-2\right)}{\sqrt{2\pi}(\rho/2+\sqrt{\rho})^{\rho/2+\sqrt{\rho}+1/2}(\rho/2-\sqrt{\rho})^{\rho/2-\sqrt{\rho}+1/2}} \\
&\geq \frac{\rho^{\rho+1/2}\exp\left(\frac{1}{12\rho+1}-2\right)}{\sqrt{2\pi}\left(\frac{\rho}{2}+\frac{2\rho}{\rho+1}\right)^{\rho+1}} \tag{10} \\
&\geq \frac{\rho^{\rho+1/2}\exp\left(\frac{1}{12\rho+1}-2\right)}{\sqrt{2\pi}(\rho/2+2)^{\rho+1}} \\
&= \frac{1}{\sqrt{2\pi\rho}} \cdot \left(\frac{\rho}{\rho/2+2}\right)^{\rho+1} \exp\left(\frac{1}{12\rho+1}-2\right) \\
&\geq \frac{2^{\rho+1}}{\sqrt{2\pi\rho} \cdot e^2} \cdot \left(\frac{\rho/2}{\rho/2+2}\right)^{\rho+1} = \frac{2^{\rho+1}}{\sqrt{2\pi\rho} \cdot e^2} \cdot \left(1-\frac{1}{\rho/4+1}\right)^{\rho+1} \\
&\geq \frac{2^{\rho+1}}{\sqrt{2\pi\rho} \cdot e^2} \cdot \left(1-\frac{1}{\rho/4+1}\right)^{4(\rho/4+1)} \\
&> \frac{2^{\rho+1}}{\sqrt{2\pi\rho} \cdot e^2} \cdot \left(\frac{1}{5}\right)^5 \tag{11}
\end{aligned}$$

The inequality (10) holds since

$$\begin{aligned}
& (\rho/2+\sqrt{\rho})^{\rho/2+\sqrt{\rho}+1/2} \cdot (\rho/2-\sqrt{\rho})^{\rho/2-\sqrt{\rho}+1/2} \\
&\leq \left(\frac{(\rho/2+\sqrt{\rho})(\rho/2+\sqrt{\rho}+1/2)+(\rho/2-\sqrt{\rho})(\rho/2-\sqrt{\rho}+1/2)}{\rho/2+\sqrt{\rho}+1/2+\rho/2-\sqrt{\rho}+1/2}\right)^{\rho/2+\sqrt{\rho}+1/2+\rho/2-\sqrt{\rho}+1/2} \\
&= \left(\frac{\rho^2/2+5\rho/2}{\rho+1}\right)^{\rho+1} = \left(\frac{\rho}{2}+\frac{2\rho}{\rho+1}\right)^{\rho+1}
\end{aligned}$$

And the inequality (11) holds for  $\rho \geq 1$ , since the function  $(1 - \frac{1}{x})^x$  is monotonically increasing and converges to  $e^{-1}$  when  $x > 1$ , therefore

$$\left(1 - \frac{1}{\rho/4 + 1}\right)^{\rho/4 + 1} \geq \left(\frac{1}{5}\right)^{5/4}$$

To conclude, we have the lower bound of  $\binom{\rho}{\rho/2 + \sqrt{\rho}}$ ,

$$\binom{\rho}{\rho/2 + \sqrt{\rho}} > \frac{2^{\rho+1}}{\sqrt{2\pi\rho} \cdot e^2 \cdot 5^5} \quad (12)$$

In fact, this lower bound can be further improved to  $\frac{2^{\rho+1}}{\sqrt{2\pi\rho \cdot e^6}}$  when  $\rho$  is sufficiently large.

Then, plugging (12) into (9),

$$\begin{aligned} & \Pr_{y \sim U_n} \left[ \rho/2 + \sqrt{\rho}/2 \leq (\# \text{ of } 0\text{'s in last } \rho \text{ bits of } y) \leq \rho/2 + \sqrt{\rho} \right] \\ & \geq \frac{\sqrt{\rho}}{2^{\rho+1}} \cdot \binom{\rho}{\rho/2 + \sqrt{\rho}} \\ & \geq \frac{\sqrt{\rho}}{2^{\rho+1}} \cdot \frac{2^{\rho+1}}{\sqrt{2\pi\rho} \cdot e^2 \cdot 5^5} = \frac{1}{\sqrt{2\pi} \cdot e^2 \cdot 5^5} = \Omega(1) \end{aligned} \quad (13)$$

Therefore,

$$\begin{aligned} & \mathbb{E}_{y \sim U_n} [\text{ed}(y, y + m(\rho))] \\ & \geq \mathbb{E}_{y \sim U_n} \left[ \text{ed}(y, y + m(\rho)) \mid \rho/2 + \sqrt{\rho}/2 \leq (\# \text{ of } 0\text{'s in last } \rho \text{ bits of } y) \leq \rho/2 + \sqrt{\rho} \right] \\ & \quad \cdot \Pr_{y \sim U_n} \left[ \rho/2 + \sqrt{\rho}/2 \leq (\# \text{ of } 0\text{'s in last } \rho \text{ bits of } y) \leq \rho/2 + \sqrt{\rho} \right] \\ & \geq \sqrt{\rho} \cdot \Omega(1) = \Omega(\sqrt{\rho}) \end{aligned}$$

The expected edit distance between  $y$  and  $y + m(\rho)$  is  $\Omega(\sqrt{\rho})$ . □

## 4 Experimental results

This work presents analytical results. This is the first theoretical study on metric embeddings between metric spaces of combinatorial nature. The hope is that it will provide machinery in areas such as in image processing, speech recognition, data mining, and natural language processing (NLP).

In this section, we take one of the mainstream applications nowadays in NLP, machine translation, and verify that the setting where our theoretical developments happen are consistent with this setting.

We consider all sentences of bounded length in one human language as a finite metric space. We impose metric structure by choosing edit (or Levenshtein) distance [Lev66] to measure the distance between sentences. Other scoring functions commonly used in machine translation won't be appropriate because they fail to satisfy the symmetry and triangle inequality. We find rather intriguing the fact that edit distance is able to reveal very non-trivial geometry in human and machine translations. This human language can be a natural (perfect) human language such as French, English, but can also be a machine generated language, such as translations of a natural language by a translation engine, e.g. Google translate. The translation of sentences between two languages now becomes a metric embedding.

In our setting, the source metric space is Chinese language and the target metric space is English. The mapping is translating from a Chinese sentence into an English sentence. The experts are Bing translate [Bin] ( $A_0$ ) from Chinese to English and Google translate [Goo] ( $A_1$ ). Therefore, we have three metric mappings: (i) from Chinese natural language to English natural language; (ii) from Chinese natural language to Google translation; (iii) from Chinese natural language to Bing translation.

#### 4.1 Data sets

Experiment	Data	Original Words	Words	Vocabulary
Natural	Chinese	2752	2624	479
Language	English Reference	2755	2748	481
Google	Chinese	2752	2747	480
	Google Translation	2542	2536	444
Bing	Chinese	2752	2746	480
	Bing Translation	2554	2548	369

Table 1: IWSLT 2011 corpus statistics in three experimental settings: natural language, Google and Bing translations, on running words before and after filtering as well as vocabulary size in thousand units [K].

The data is provided from IWSLT [IWS] competition with 139K sentence pairs in Chinese and English. The details of this corpus can be found



in Table 1. We used Google and Bing online translation engine application to obtain the English machine translation results. To eliminate the effect of noise in data, such as incomplete sentences, and other sporadic translation anomalies (distortion is a worst-case quantity), we show the distortion value on selected subsets of the data disregarding sentences triggering high distortion. This happens either due to errors in the original corpus, or due to large difference in the sentence length between two languages. For filtering, we use simple greedy heuristics to compute these sets which in reality could be enough larger (finding these sets can be somewhat easily shown to be NP-hard by a reduction from INDEPENDENT SET and it is not clear if any reasonable approximation exists – recall that independent set itself is vastly inapproximable). We remove the sentence pair with the highest value on contraction or expansion and re-calculate the distortion iteratively.

#### 4.1.1 Bi-Lipschitz condition (low-distortion)

Our *metric embeddings in learning* developments in the main paper is based on the Bi-Lipschitz condition only. In this subsection, we show that translation within our setting satisfies this condition (at least within some relaxed notion of locality), thus our theory is very naturally related to language translation. The term bi-Lipschitz is used informally to refer to small distortion. The distortion is the product of *expansion* and *contraction*, where the contraction is the maximum value among each distance of two Chinese sentences divided by the distance their English translations; and the expansion is the maximum value among each distance of the two translated English sentences divided by the distance of the corresponding sentences in Chinese. Here, bi-Lipschitz means that for every two Chinese sentences, which are close in edit distance their English translations are also close, and vice-versa.

Figure 1 depicts the expansion, contraction, and distortion on the subsets of sentences selected based on their lengths, ranging from 1 to 100. Each subset is composed of 1000 sentence pairs (sampled uniformly random) with no more or less than two of the given length in average, when filtering out 10% of the sentences (which includes noise).

Shorter sentences tend to have higher distortion, because of the sensitivity of distortion on the number of words. For instance, two far related single-word Chinese idioms (”众说纷纭” and ”接下来”) are translated into an English sentence with eight and six words (“many things have been said about this .” and “now follow me , ok ? ”) , respectively, their distance is two in Chinese and fourteen in English.

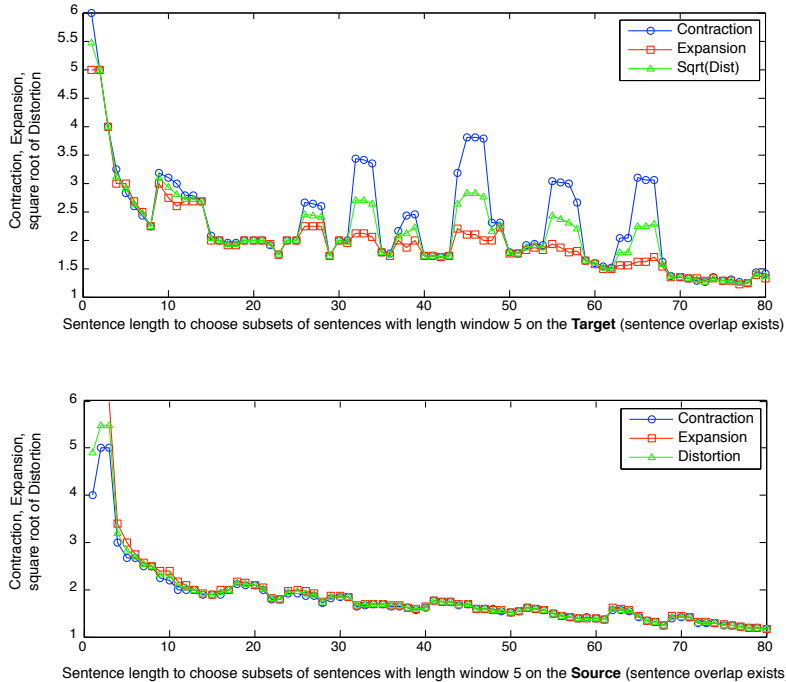


Figure 1: Distortion in sentence length range for the *natural language*. The x-axis depicts the window position  $x$  for sentence length in the window  $[x - 2, x + 2]$ . The y-axis is the distortion occurring in the corresponding window.

Our empirical study revealed two more conditions satisfied in translation setting: METRIC SEPARABILITY and DENSE NEIGHBORING. These conditions are not used in the statements of our theorems. However, we remark that *all of our constructions of the concept classes and the experts in our theorems do satisfy these conditions as well*. Further understanding in the geometry of human translation maps can lead to new classification techniques.

#### 4.1.2 Metric separability condition

We partition the translation output in Table 1 into two sets: one that contains the translated sentences generated using Bing that has smaller edit distance than those generated using Google (this is the set: “Bing-better”); the other contains all translations generated using Google which has smaller

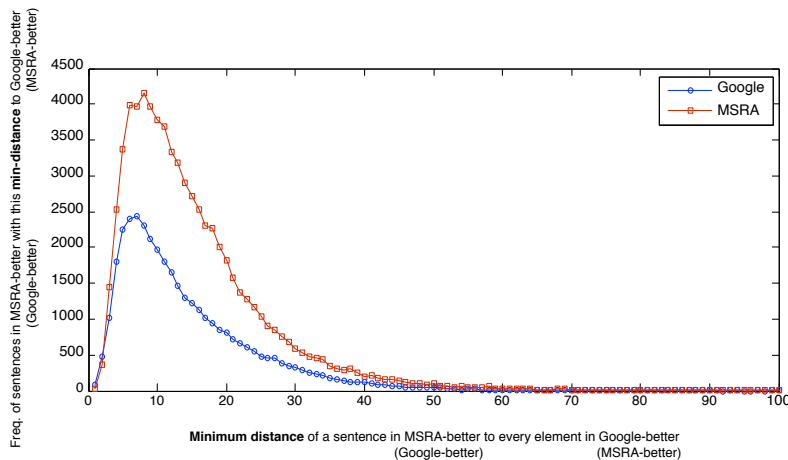


Figure 2: Frequency of sentences in Google or Bing better performed set given the minimum Levenshtein distance of this sentence from the other. Google-better: the set where Google does better than Bing; Bing-better: the set where Bing does better than Google. The x-axis is the minimum distance of Bing-better (Google-better) to Google-better (Bing-better). The y-axis is the number of sentences in the Bing-better (Google-better) set with the corresponding minimum distance from Google-better (Bing-better).

distance than the Bing one (Google-better).

There is a large subset of Bing-better that has significant distance from Google-better set (where distance between two sets is the min-distance of a point in one and a point in the other). This means that for every sentence in Bing-better, its minimum distance to all sentences in Google-better set is greater than a threshold, i.e. there is a certain distance from each sentence in Bing-better to the Google-better. Therefore Bing-better and Google-better set can be separated in a metric/geometric sense.

Figure 2 depicts the histogram of the number of sentences given its minimum distance to the target set. An important experimental finding is that the better performed translations are concentrated close to the area that has the minimum distance of 10 to the other set. The curve for Bing as in Figure 2 is greater than that for the Google, and is pushed slightly to the area with longer distances. This result indicates a strong geometric property: it is not only the case that there are large sets where Bing outperforms Google and vice-versa, but also when they do so they do it for sentences that are far away in an edit distance sense.

### 4.1.3 Dense neighboring condition

Consider one of the two systems performing better and the other performing worse for each sentence in a maximum subset of source sentences. For this subset if we translate a sentence with the worse system, then we can always find another sentence in the same subset whose translation is close to the first translation.

In our experiments, for each source sentence better translated with Google, we apply Bing to translate it as well. For this translation, we look for the closest translation generated by Google, whose source sentence is translated better by Google than by Bing with the respect to the edit distance. The results show that 10% of the worse translated sentences by Bing has a neighborhood with less than or equal to 5 in edit distance, where as a neighborhood we take a good translated sentence by Google. For edit distance 10 the ; and 50% of sentences with a edit distance of 10. Google versus Bing shows analogous result. We believe (and the tendency of the graph is such) that the percentage of sentences fulfilling this assumption will increase by increasing the corpus size.

## References

- [Bin] Bing. Bing online machine translation system.
- [Goo] Google. Google translate online machine translation system.
- [IWS] IWSLT. International workshop on spoken language translation 2011.
- [Lev66] V. I. Levenshtein. Binary codes capable of correcting deletions, insertions and reversals. In *Soviet physics doklady*, volume 10, page 707, 1966.