

ISOPERIMETRIC SETS OF ENGLISH WORDS

GB,TC,MC,MD,BF,WG,XH,SK,TL,PM,NM,AN,EP,ZT,AU,EZ

ABSTRACT. We consider a space of English words, define distance and boundary, and find subsets of small volumes of minimum or maximum perimeter.

1. INTRODUCTION

The isoperimetric problem is among the oldest in mathematics. It asks for the least-perimeter way to enclose a given volume. In this paper, we consider a space of English words and seek subsets of small volumes of minimum or maximum perimeter.

While there has been much study of numerical properties of the English language [4] and of the isoperimetric problem on discrete spaces (see e.g. [2] and [1]), as far as we know, our specific focus is new.

We use a standard list of the 3000 most common English words [3], with a notion of distance and perimeter (Sect. 2). We begin with the study of singletons (word sets with volume 1).

1.1. **Two-letter words.** For our space of two-letter words, the singletons of least perimeter are {pc}, {tv}, {hi}, and {up}, with boundaries {pm}, {to}, {he}, and {us}, respectively (Prop. 3.1). The singletons of greatest perimeter (7) are {as}, {ie}, and {me} (Prop. 3.2).

1.2. **Three-letter words.** For our space of three-letter words, there are 15 isolated singletons with no boundary. They are

{ask}, {ceo}, {dna}, {egg}, {etc}, {eye}, {fly}, {ice}, {its}, {mom}, {mrs}, {off},
{oil}, {via}, {you}

(Prop. 4.1). Note for example that {ice} is isolated only because *ace* is not on the list of 3000 words. The singleton with maximum perimeter (12) is {set}. It has boundary

$$\partial\{\text{set}\} = \{\text{bet, get, jet, let, net, pet, sea, see, sex, sit, wet, yet}\}$$

(Prop. 4.2).

Date: September 3, 2022.

1.3. **All words.** For the space of all words on our list, there are 1782 singletons with no boundary out of the 3000 words (Prop. 5.1). The singleton with maximum perimeter (16) is {ear} (Prop. 5.3). It has boundary

{bar, car, eat, era, far, war, bear, dear, earn, fear, gear, hear, near, tear, wear, year}.

1.4. **Doubletons (volume 2).** For the space of two-letter words, the doubleton with minimum perimeter (1) is {pc, pm}, with boundary {am} (Prop. 6.2). The doubleton with greatest perimeter (14) in the space of two-letter words is {as, ie}, which has boundary

$$\partial\{as, ie\} = \{ad, ah, am, at, be, he, if, in, it, me, ms, us, vs, we\}$$

(Prop. 6.3). In the universe of three-letter words, every doubleton with minimum perimeter (0) is formed in one of two ways: the doubleton is made of two isolated singletons or two singletons whose boundaries are each other (Prop. 6.4). The doubletons with maximum perimeter (21) in the universe of three-letter words are {set, cap}, {set, gay}, {set, lay}, {set, may} (Prop. 6.5). In the universe of all words, the sets of volume two with minimum perimeter (0) include the $\binom{1782}{2}$ pairs of isolated words (Prop. 6.6). In the universe of all words, the doubleton with maximum perimeter (31) is {bet, ear} (Prop. 6.7).

1.5. **Proofs.** Our methods include computation and elementary deductions therefrom. Specifically, computation was used to calculate word sets of given volume and word length and generate tables of word sets and their boundaries and perimeters. The code was written in Python and the results were confirmed with a different Python program implemented independently.

1.6. **Graph theory.** Our study could be phrased in graph theory, as words being the vertices of a graph, and words distance 1 apart being connected by an edge. However, we do not use many ideas from graph theory in this paper, leaving them as open questions.

Acknowledgements. This paper is a product of the 2022 MathPath breakout on Metric Spaces, advised by Frank Morgan. The authors would like to thank Dr. Morgan for his advice and guidance throughout the breakout and the MathPath staff for their support.

2. DEFINITIONS

Definition 2.1 (Space W). We define the space W as the 3000 most common English words as provided by ef.edu [3]. Let W_n denote the set of all n -letter words. The list generally does not include plurals or conjugations — e.g. *is* is not found in the 3000 words. We ignore punctuation and capitalization; for example, *n't*, the contraction for *not*, is considered to be a two-letter English word *nt*, *e-mail* is considered to be a

five-letter word *email*, and abbreviations such as *Mr* and *Ms* are considered to be the words *mr* and *ms* respectively.

Definition 2.2 (Distance). Given words w_1 and w_2 in $Z \subset W$, we define the distance from w_1 to w_2 in Z as the minimum number of letter changes it takes to change w_1 into w_2 . A letter change is defined as inserting a letter, deleting a letter, replacing a letter, or swapping two adjacent letters. Each word along the path from w_1 to w_2 must be another member of the set Z .

Definition 2.3 (Boundary). The *boundary* in $Z \subset W$ of a set $S \subset Z$ is the set of words in $Z - S$ which are distance one from at least one word in S . We denote the boundary of S as ∂S . The number of elements of the boundary is called the *perimeter*. A word w is *isolated* if $\{w\}$ has an empty boundary.

3. TWO-LETTER WORDS

Propositions 3.1 and 3.2 provide, in the universe W_2 of two-letter words, the sets of volume 1 of minimum and maximum perimeter. They follow from the following Table 3.1, which shows two-letter words and their boundaries, along with their boundary size (perimeter).

TABLE 3.1. Singletons and their boundaries in the universe W_2 of two-letter words.

Word Set	Boundary	Boundary Size
{hi}	{he}	1
{pc}	{pm}	1
{tv}	{to}	1
{up}	{us}	1
{by}	{be, my}	2
{pm}	{am, pc}	2
	⋮	
{at}	{ad, ah, am, as, it, nt}	6
{ms}	{as, me, mr, my, us, vs}	6
{no}	{do, go, nt, on, so, to}	6
{on}	{in, no, of, oh, ok, or}	6
{as}	{ad, ah, am, at, ms, us, vs}	7
{ie}	{be, he, if, in, it, me, we}	7
{me}	{be, he, ie, mr, ms, my, we}	7

Proposition 3.1. *In the universe W_2 of two-letter words, the singletons of minimum perimeter (1) are {pc}, {tv}, {hi}, and {up}. Their respective boundaries are {pm}, {to}, {he}, and {us}.*

Proposition 3.2. *In the universe W_2 of two-letter words, the singletons of maximum perimeter (7) are {as}, {ie}, and {me}.*

The boundary of {as} is

$$\partial\{as\} = \{ad, ah, am, at, ms, us, vs\}.$$

The boundary of {ie} is

$$\partial\{ie\} = \{be, he, if, in, it, me, we\}.$$

The boundary of {me} is

$$\partial\{me\} = \{be, he, ie, mr, ms, my, we\}.$$

Proof. Both propositions follow immediately from Table 3.1. ■

4. THREE-LETTER WORDS

Propositions [4.1](#) and [4.2](#) provide in the universe W_3 of three-letter words the sets of volume 1 of minimum and maximum boundary. They follow from the following [Table 4.1](#), which shows all three-letter words and their boundaries, along with their boundary size.

TABLE 4.1. Singletons and their boundaries in the universe W_3 of three-letter words.

Word Set	Boundary	Boundary Size
{ask}	{}	0
{ceo}	{}	0
{dna}	{}	0
{egg}	{}	0
{etc}	{}	0
{eye}	{}	0
{fly}	{}	0
{ice}	{}	0
{its}	{}	0
{mom}	{}	0
{mrs}	{}	0
{off}	{}	0
{oil}	{}	0
{via}	{}	0
{you}	{}	0
	:	
{cap}	{can, car, cat, cop, cup, gap, lap, map, tap}	9
{gay}	{day, gap, gas, guy, lay, may, pay, say, way}	9
{jet}	{bet, get, jew, let, net, pet, set, wet, yet}	9
{lay}	{day, gay, lab, lap, law, may, pay, say, way}	9
{may}	{day, gay, lay, mad, man, map, pay, say, way}	9
{yet}	{bet, get, jet, let, net, pet, set, wet, yes}	9
{let}	{bet, get, jet, leg, lot, net, pet, set, wet, yet}	10
{bet}	{bed, bit, but, get, jet, let, net, pet, set, wet, yet}	11
{net}	{bet, get, jet, let, new, not, nut, pet, set, wet, yet}	11
{pet}	{bet, get, jet, let, net, per, pot, put, set, wet, yet}	11
{set}	{bet, get, jet, let, net, pet, sea, see, sex, sit, wet, yet}	12

Proposition 4.1. *In the universe W_3 of three-letter words, the singletons of words with minimum perimeter (0) are $\{\text{ask}\}, \{\text{ceo}\}, \{\text{dna}\}, \{\text{egg}\}, \{\text{etc}\}, \{\text{eye}\}, \{\text{fly}\}, \{\text{ice}\}, \{\text{its}\}, \{\text{mom}\}, \{\text{mrs}\}, \{\text{off}\}, \{\text{oil}\}, \{\text{via}\},$ and $\{\text{you}\}$. There are 15 of these singletons, all of which have empty boundary.*

Proposition 4.2. *In the universe W_3 of three-letter words, the set of volume 1 and maximum perimeter (12) is $\{\text{set}\}$. Its boundary is*

$$\partial\{\text{set}\} = \{\text{bet}, \text{get}, \text{jet}, \text{let}, \text{net}, \text{pet}, \text{sea}, \text{see}, \text{sex}, \text{sit}, \text{wet}, \text{yet}\}.$$

Proof. Both propositions follow immediately from Table 4.1. ■

5. ALL WORDS

Propositions 5.1 and 5.3 provide in the universe W of all words the singletons of minimum and maximum boundary. The following propositions follow from a computer search.

Proposition 5.1. *In the universe W of all words, there are 1782 singletons with empty boundary.*

Proof. This follows directly from a computer search. ■

Corollary 5.2. *For $V \leq 1782$, the minimum perimeter of a set with volume V is 0.*

Proof. For a given V , we take the union of V arbitrary singletons that have perimeter 0. This set will also have perimeter 0 since none of its elements are connected to other words. Thus, we have a set with volume V that has no perimeter, proving our corollary. ■

Proposition 5.3. *In the universe W of all words, the subset with volume 1 and maximum perimeter (16) is $\{\text{ear}\}$. The boundary of $\{\text{ear}\}$ is:*

$$\{\text{bar}, \text{car}, \text{eat}, \text{era}, \text{far}, \text{war}, \text{bear}, \text{dear}, \text{earn}, \text{fear}, \text{gear}, \text{hear}, \text{near}, \text{tear}, \text{wear}, \text{year}\}.$$

Proof. This follows directly from a computer search. ■

6. VOLUME TWO

Earlier sections dealt with singletons, but now we will consider word sets of volume two. We will start with two-letter words. But first, a useful lemma.

Lemma 6.1. *The boundary of a doubleton D in any universe $Z \subset W$ is the union of the boundaries of the singleton subsets, minus the singletons themselves, i.e.*

$$\partial\{w_1, w_2\} = \partial\{w_1\} \cup \partial\{w_2\} - \{w_1, w_2\},$$

where ∂ is the boundary of a set.

Proof. If a word is in the boundary of D , it is within distance 1 of a word of D , and consequently within the union of the boundaries of those word singletons, while by definition not in D itself. On the other hand, if a word is in that union, it is within distance 1 of w_1 or w_2 , so if in addition it is not in D , it is in the boundary of D . ■

Proposition 6.2. *In the universe W_2 of two-letter words, the set of volume 2 of minimum perimeter (1) is $\{pm, pc\}$, with boundary $\{am\}$.*

Proof. Note from Table 3.1 that for every singleton in W_2 , the perimeter is at least 1. Consider a set of volume 2 $\{w_1, w_2\}$ with perimeter at most 1. Neither of the singletons $\{w_1\}$ or $\{w_2\}$ can have perimeter greater than 2, or by Lemma 6.1, $\{w_1, w_2\}$ would have perimeter greater than 1. If both singletons had perimeter 1, by Table 3.1, the perimeter would be 2 because the boundaries of all singletons with perimeter 1 are disjoint. Therefore, w_1 is *by* or *pm* since they are the only singletons with perimeter 2. If w_1 is *by*, w_2 must be *my* or *be* because they are the only words in the boundary of $\{by\}$. But w_2 cannot be *my* or *be*, because those singletons have perimeter greater than 2. Therefore, w_1 must be *pm*. This means that w_2 must be *am* or *pc*, the words in the boundary of $\{pm\}$. However, $\{am\}$ has boundary greater than 2, so w_2 must be *pc*. The boundary of $\{pm, pc\}$ is $\{am\}$, so it is the unique doubleton with minimum perimeter 1. ■

Proposition 6.3. *In the universe W_2 of two-letter words, the set of volume 2 and maximum perimeter (14) is $\{as, ie\}$. The boundary is*

$$\partial\{as, ie\} = \{ad, ah, am, at, be, he, if, in, it, me, ms, us, vs, we\}.$$

Proof. Suppose the maximum perimeter of $\{w_1, w_2\}$ is 14 or more. By Lemma 6.1 and Table 3.1, each singleton would need perimeter 7 and no overlap with the other singleton. Only $\{as, ie\}$ satisfies these conditions. ■

The preceding two propositions may also be verified by computation, as summarized in Table 6.1.

TABLE 6.1. Doubletons and their boundaries in the universe of two-letter words.

Word Set	Boundary	Boundary Size
{pc, pm}	{am}	1
{hi, pc}	{he, pm}	2
{hi, tv}	{he, to}	2
{hi, up}	{he, us}	2
{pc, tv}	{pm, to}	2
{pc, up}	{pm, us}	2
{tv, up}	{to, us}	2
	⋮	
{as, me}	{ad, ah, am, at, be, he, ie, mr, ms, my, us, vs, we}	13
{as, no}	{ad, ah, am, at, do, go, ms, nt, on, so, to, us, vs}	13
{as, on}	{ad, ah, am, at, in, ms, no, of, oh, ok, or, us, vs}	13
{at, me}	{ad, ah, am, as, be, he, ie, it, mr, ms, my, nt, we}	13
{ie, no}	{be, do, go, he, if, in, it, me, nt, on, so, to, we}	13
{me, no}	{be, do, go, he, ie, mr, ms, my, nt, on, so, to, we}	13
{me, on}	{be, he, ie, in, mr, ms, my, no, of, oh, ok, or, we}	13
{as, ie}	{ad, ah, am, at, be, he, if, in, it, me, ms, us, vs, we}	14

Proposition 6.4. *In the universe W_3 , the 107 sets of volume 2 with minimum perimeter (0) are formed in one of two ways. 105 of the sets come from two isolated singletons. The other two come from two singletons whose boundaries are each other: {age, ago} and {all, ill}.*

Proof. By Lemma 6.1, if a doubleton has no boundary, either each singleton subset has no boundary or each singleton subset is the boundary of the other. (Note that if one is contained in the boundary of the second, then the second is contained in the boundary of the first.) Since by Proposition 4.1 there are 15 isolated words, there are

$\binom{15}{2} = 105$ pairs of them. By computation, there are only two pairs of the second type: $\{age, ago\}$ and $\{all, ill\}$ ■

Proposition 6.5. *In the universe W_3 of three-letter words, the sets of volume 2 and maximum perimeter (21) are $\{set, cap\}$, $\{set, gay\}$, $\{set, lay\}$, and $\{set, may\}$.*

Proof. Consider a doubleton $\{w_1, w_2\}$ with perimeter greater than or equal to 21. By Lemma 6.1, the only possibilities for the two singletons $\{w_1\}$, $\{w_2\}$ are the 11 of Table 4.1 with perimeter at least 9 (and at most 12 in the case of $\{set\}$). If both w_1 and w_2 end in *et*, the boundaries of the singletons have too much overlap for the perimeter of $\{w_1, w_2\}$ to reach 21. Therefore one of them must be *cap*, *gay*, *lay*, or *may*, all with perimeter 9, and the other must be *set* with perimeter 12. The resulting four doubletons indeed have perimeter 21 and are therefore the only maxima of perimeter. ■

Proposition 6.6. *In the universe W of all words, the sets of volume 2 of minimum perimeter (0) include the $\binom{1782}{2} = 1,586,871$ pairs of isolated words.*

Proof. By computation, there are 1782 isolated words. For each of the 1,586,871 pairs of these words, their set will have perimeter 0. ■

Proposition 6.7. *In the universe W of all words, the set of volume 2 and maximum perimeter (31) is $\{bet, ear\}$. The boundary is*

$\{be, bed, bit, but, get, jet, let, net, pet, set, wet, yet, beat, belt, best, bar, car, eat, era, far, war, bear, dear, earn, fear, gear, hear, near, tear, wear, year\}$.

Proof. By computation, the two singletons with largest perimeter are $\{bet\}$ and $\{ear\}$, with perimeters of 15 and 16 respectively. By computation, and as referenced in Proposition 5.3, the boundary of the singleton $\{ear\}$ is

$\{bar, car, eat, era, far, war, bear, dear, earn, fear, gear, hear, near, tear, wear, year\}$.

By computation, the boundary of the singleton $\{bet\}$ is

$\{be, bed, bit, but, get, jet, let, net, pet, set, wet, yet, beat, belt, best\}$.

Since the two boundaries and $\{bet, ear\}$ are pairwise disjoint, by Lemma 6.1, the boundary of the doubleton $\{bet, ear\}$ is the union of the boundaries of the singletons, so the boundary of $\{bet, ear\}$ is as asserted, with the maximal 31 elements. ■

7. OPEN QUESTIONS

- (1) A *connected component* is a set with empty boundary and no proper subsets with empty boundary. How many connected components are there in W ? There are at least 1782, as there are 1782 singleton components.
- (2) Which connected component of W has largest volume? It likely contains the words a and i .

- (3) What is the distribution of the perimeter of word sets with given volume?
- (4) What is the largest clique in W (where every word is unit distance from every other)?
- (5) What sets have minimum and maximum fattened boundary (including words within distance two)?

REFERENCES

- [1] S. G. Bobkov and F. Götze. “Discrete isoperimetric and Poincare-type inequalities”. In: *Probability Theory and Related Fields* **114** (1997), 245–277. URL: https://www-users.cse.umn.edu/~bobko001/papers/1999_PTRF_BG.pdf.
- [2] Fan Chung. “Discrete Isoperimetric Inequalities”. In: *Surveys in Differential Geometry* **9** (2004), 53–82. URL: <https://mathweb.ucsd.edu/~fan/wp/iso.pdf>.
- [3] ef.edu. *3000 most common words in English*. 2022. URL: <https://www.ef.edu/english-resources/english-vocabulary/top-3000-words/>.
- [4] Wikipedia. *Most common words in English*. URL: https://en.wikipedia.org/wiki/Most_common_words_in_English.

SCIENCE AND ARTS ACADEMY, 1825 MINER STREET, DES PLAINES, IL 60016

RANCHO SAN JOAQUIN MIDDLE SCHOOL, 4861 MICHELSON DRIVE, IRVINE, CA 92612

HARVARD-WESTLAKE SCHOOL, 700 N FARING RD, LOS ANGELES, CA 90077

PROOF SCHOOL, 973 MISSION ST, SAN FRANCISCO, CA 94103

HORACE MANN SCHOOL, 231 W 246TH ST, THE BRONX, NY 10471

CANYON CREST ACADEMY, 5951 VILLAGE CENTER LOOP RD, SAN DIEGO, CA 92130

LEXINGTON HIGH SCHOOL, 251 WALTHAM ST, LEXINGTON, MA 02421

PROOF SCHOOL, 973 MISSION ST, SAN FRANCISCO, CA 94103

VESTAVIA HILLS HIGH SCHOOL 2235 LIME ROCK RD, VESTAVIA HILLS, AL 35216

MCCALL MIDDLE SCHOOL, 458 MAIN ST, WINCHESTER, MA 01890

10521 160 AVE, GRANDE PRAIRIE, ALBERTA

2369 FORBES AVE, SANTA CLARA, CA 95050

LIVINGSTON HIGH SCHOOL, 30 ROBERT H HARP DRIVE, LIVINGSTON, NJ 07039

LAKESIDE MIDDLE SCHOOL, 510 1ST AVE NE, SEATTLE, WA 98125

400 S KIMBALL AVE, SOUTHLAKE, TX, 76092

RANCHO SAN JOAQUIN MIDDLE SCHOOL, 4861 MICHELSON DRIVE, IRVINE, CA 92612