

Klingon, Elvish, German, English, and Math

In Linearity, we are learning about Singular Value Decomposition and its applications. Inspired by the paper “Singular Vectors’ Subtle Secrets”, we decided to extend a cryptographical analysis of letter frequency in English to other languages. Klingon is a language invented for the science-fiction television show Star Trek, and Sindarin is an Elvish dialect invented by J.R.R. Tolkien for his fantasy world of Middle-Earth. Both languages have been developed to a degree where people can speak and write them fluently. We are interested to see how the letter patterns of these invented languages compare to those of natural languages, such as German and English.

Singular Value Decomposition

Singular Value Decomposition is a method of decomposing a matrix of any shape into a series of rank one matrices (matrices with only one linearly independent column.) The sum of the rank one matrices is equal to the original matrix.

The first step is finding three specific matrices, U, Σ, and V that decompose the given matrix A such that $A = UΣV^T$. This is very similar to the $A = PDP^{-1}$ decomposition, where P is the horizontal concatenation of the eigenvectors of A, and D is a diagonal matrix of the eigenvalues associated with each eigenvector. Since it is not possible to find the eigenvectors of a nonsquare matrix, instead the eigenvectors of square matrices AA^T and A^TA are found.

The nonzero eigenvalues of AA^T and A^TA are the same. Their associated eigenvectors u_i and v_i are related by $u_i = Av_i$. The columns of U, an NxN matrix, are the eigenvectors of AA^T in order of decreasing associated eigenvalues, and the columns of V, an MxM matrix, are the associated eigenvectors of A^TA . The singular values of A are the square roots of the eigenvalues of AA^T and A^TA . Σ is a diagonal matrix composed of the singular values of A and ‘filled in’ with zeros to have size MxN. The singular value $σ_i = \sqrt{λ_i}$. The best rank 1 approximation of A can then be calculated as $σ_1u_1v_1^T$. The best rank 2 approximation is $σ_1u_1v_1^T + σ_2u_2v_2^T$, and so on.

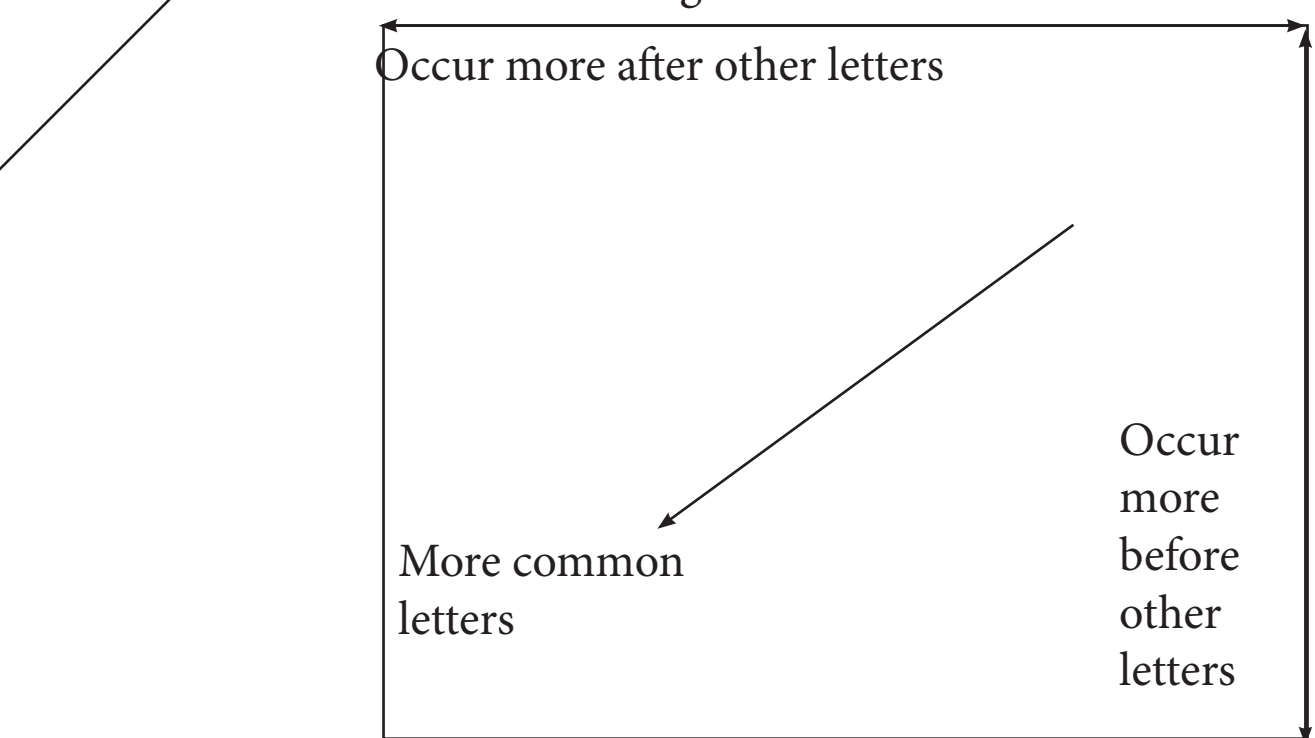
How We Use It

We use a Python script to generate an adjacency matrix (such as the one below for English) showing how frequently one letter follows another in a large sample text. This matrix is conveniently square.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z
1	0	628	889	787	5	244	650	11	1337	0	319	1080	664	666	16	815	0	3443	1640	2338	210	322	355	3	2681	12
2	335	13	0	4	2758	0	0	0	97	80	0	844	30	1	400	0	158	35	27	799	1	0	0	0	0	447
3	3056	0	1334	1	735	0	0	3734	420	0	187	356	0	0	2347	0	1037	101	2	1460	313	0	0	0	0	4
4	385	0	1	87	1474	17	96	1	1033	3	0	113	123	33	653	0	0	187	35	1	351	48	16	0	0	4
5	2453	33	1306	277	5475	356	224	60	791	21	23	2134	677	3692	65	325	124	8506	1536	1502	10	351	85	450	93	2
6	533	0	0	0	788	256	0	0	1556	0	0	351	0	0	2368	0	0	150	1	1222	275	0	0	0	0	5
7	394	1	0	4	629	0	31	740	361	0	0	167	2	47	475	0	0	634	3	16	199	0	0	0	0	4
8	144	29	0	0	100	5	0	0	1800	0	0	23	50	7	1059	0	0	133	3	126	128	0	0	0	0	13
9	282	244	3065	425	951	559	687	1	1	0	353	1029	641	1006	1701	78	0	651	1630	3338	101	555	0	134	0	28
10	33	0	0	0	152	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	342	0	0	0	0	0
11	7	0	0	0	428	26	0	0	2388	0	0	156	2	329	5	0	0	0	0	0	7	0	0	0	0	0
12	1375	0	26	114	1354	14	33	0	1646	0	318	1091	82	21	944	4	0	53	61	273	187	77	110	0	0	159
13	3408	98	0	29	1432	45	0	0	1070	0	0	13	181	14	1029	502	0	547	176	11	878	1	0	0	0	255
14	440	97	1091	876	1497	171	355	10	641	48	72	311	0	743	2206	17	39	61	100	1653	156	164	94	384	93	1
15	62	234	165	201	62	2872	91	59	266	91	136	499	1161	2726	906	413	4	2238	752	744	3620	548	672	2	53	2
16	898	1	0	2	1832	0	9	59	278	0	0	868	5	1	732	556	0	1812	8	449	121	0	2	0	0	22
17	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
18	1018	11	234	258	3070	120	93	73	1476	0	58	365	255	273	1806	70	0	347	692	892	319	250	73	0	0	124
19	750	38	242	15	2151	47	180	1344	1455	0	33	79	168	26	1179	487	0	1	695	1395	918	0	72	0	0	13
20	1405	0	78	1	1054	26	0	1000	2038	0	1	477	90	20	1648	1	0	635	33	829	478	0	429	1	23	0
21	331	148	165	139	170	54	538	0	470	0	0	1020	163	349	14	430	0	1000	1367	406	0	2	0	1	0	3
22	266	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
23	2336	2	0	4	4440	4	0	1392	1470	0	9	47	1	19	360	0	0	63	7	1	0	0	0	0	0	0
24	33	0	0	0	51	0	0	2	74	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
25	13	38	0	0	2	241	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
26	93	0	0	0	0	23	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

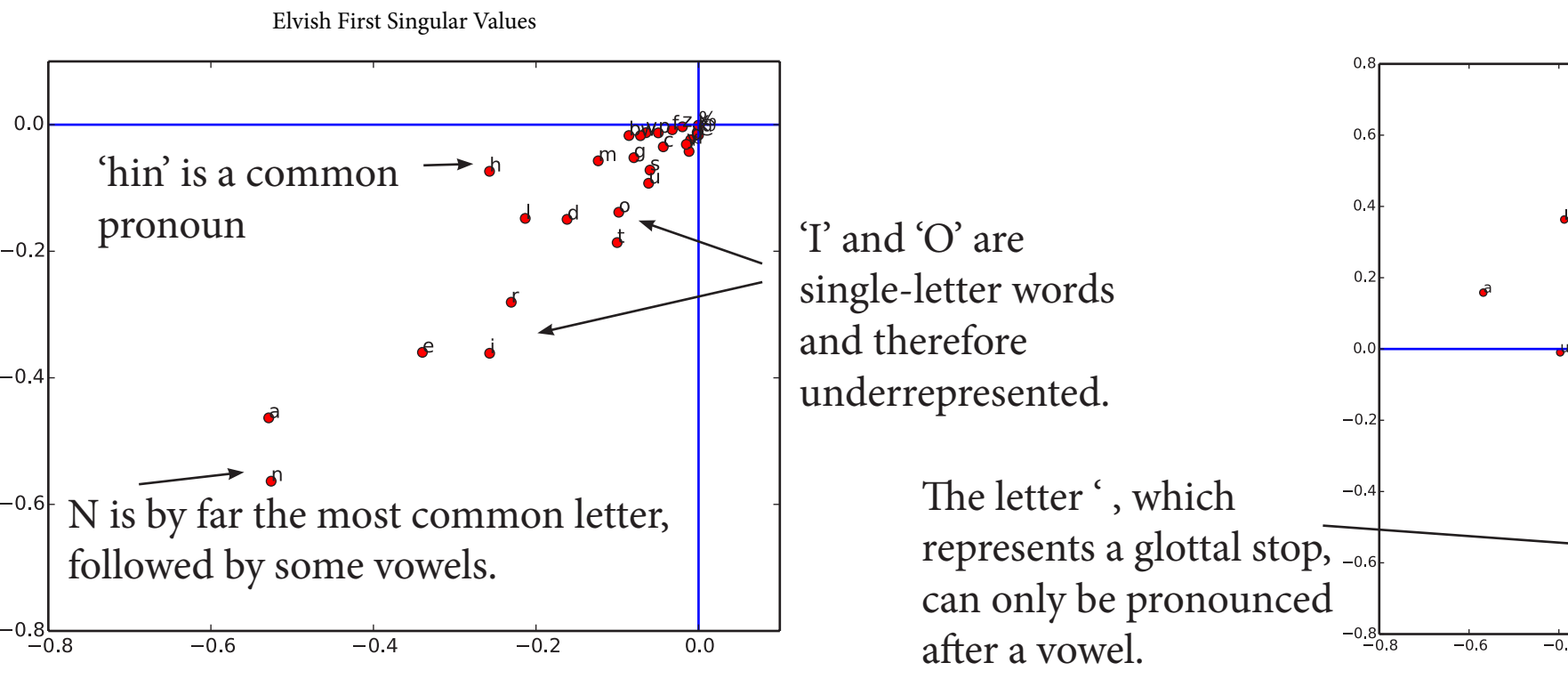
We then find U, Σ, and V for the adjacency matrix. The coordinates for a letter on the first graph are (u_{1n}, v_{1n}) , where u_1 is the first column vector in U, v_1 is the first column vector in V, and n is the index of the letter in the alphabet (so A would be 0, B would be 1, etc.) The second graph is the same except for using (u_{2n}, v_{2n}) , from the second columns of U and V.

Elvish shows much less separation between beginning and ending letters. Its distribution is more natural than Klingon's, but not as varied as English or German.



Elvish has a strong consonant-vowel pattern like English, but since its most common letter, ‘N’, is a consonant, this method of analysis has placed consonants in the second quadrant and vowels in the fourth.

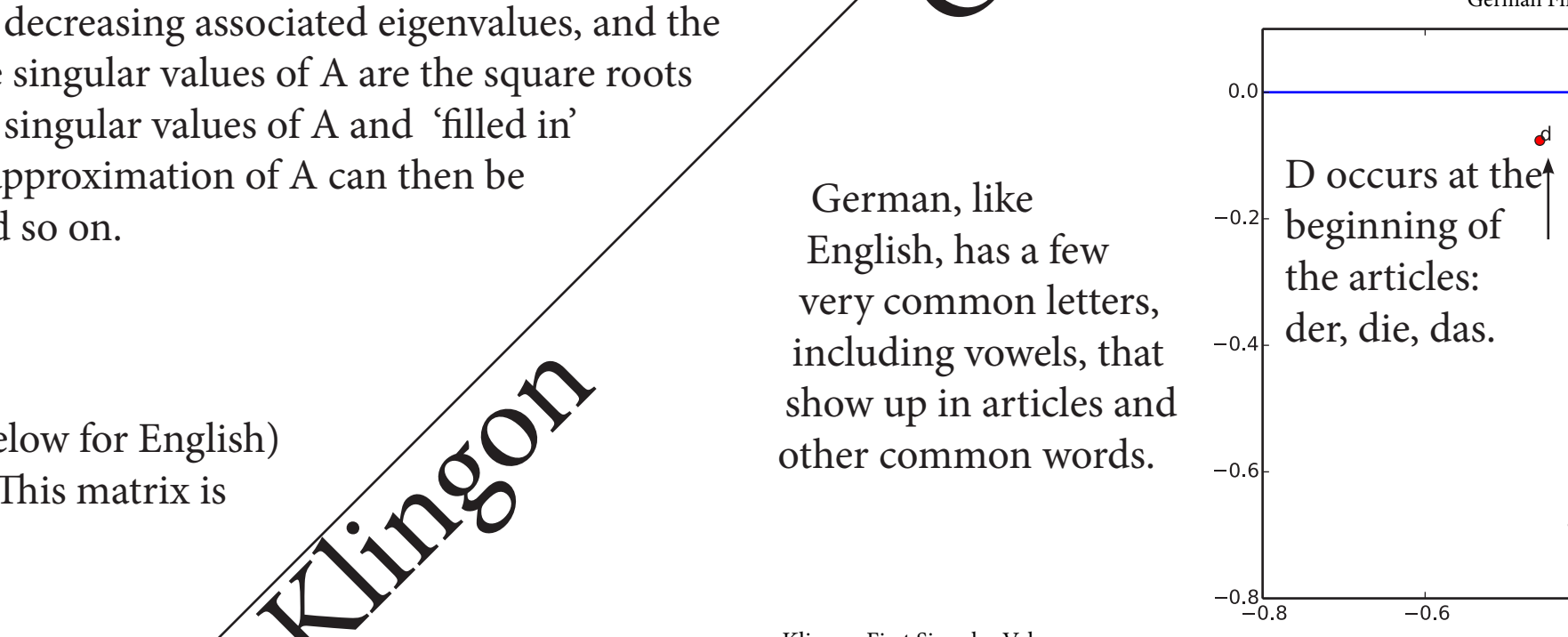
The letters ‘th’ followed by a vowel occurs in the middle of many words, accounting for h being an outlier in this quadrant. Since the graph is flipped, this is actually the same place as English’s ‘H’.



Klingon has the most artificial separation between vowels and consonants, representing the strict rules used to create it.

The most common ending letter is ‘ followed by D.

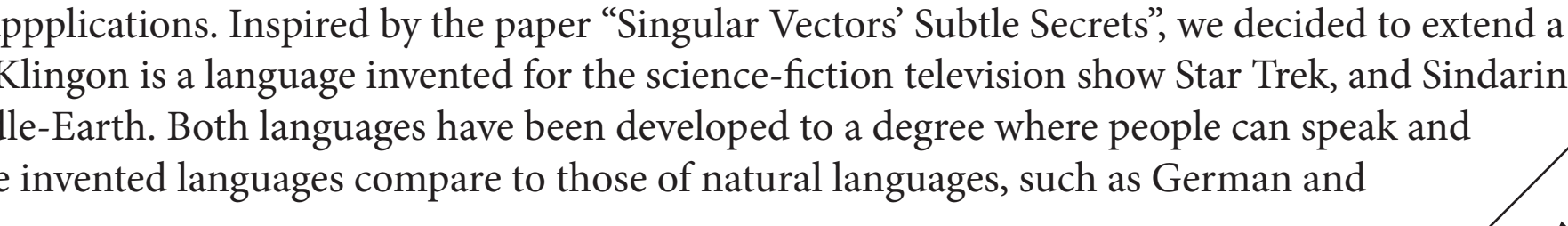
A word never ends with a vowel.



German, like English, has a few very common letters, including vowels, that show up in articles and other common words.

D occurs at the beginning of the articles: der, die, das.

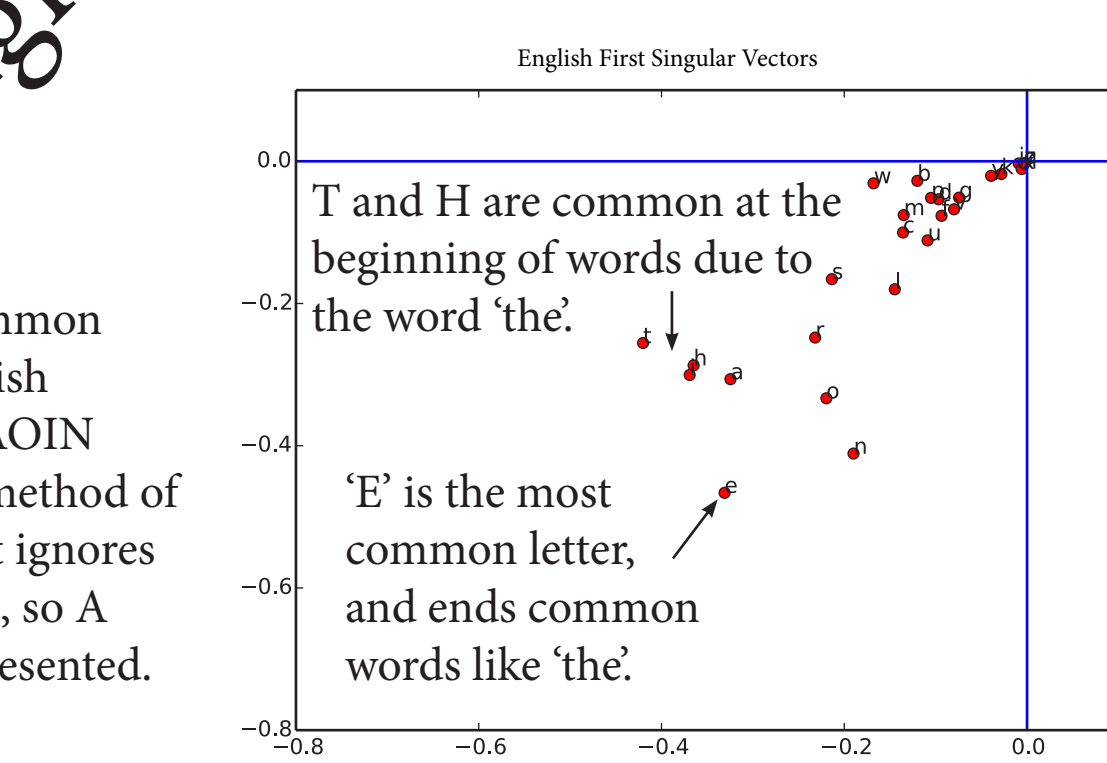
E is more common than in English, especially at the end of words.



The most common letters in English are usually ETAOIN SHRDLU. Our method of analyzing the text ignores single-letter words, so A and I are underrepresented.

T and H are common at the beginning of words due to the word ‘the’.

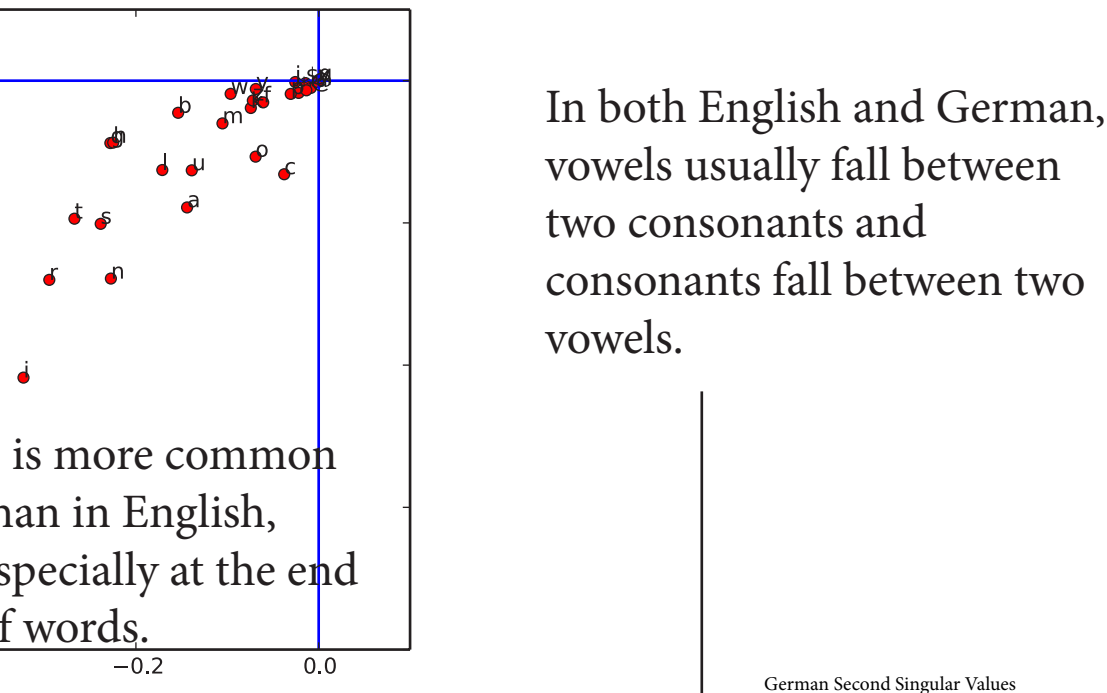
‘E’ is the most common letter, and ends common words like ‘the’.



In both English and German, vowels usually fall between two consonants and consonants fall between two vowels.

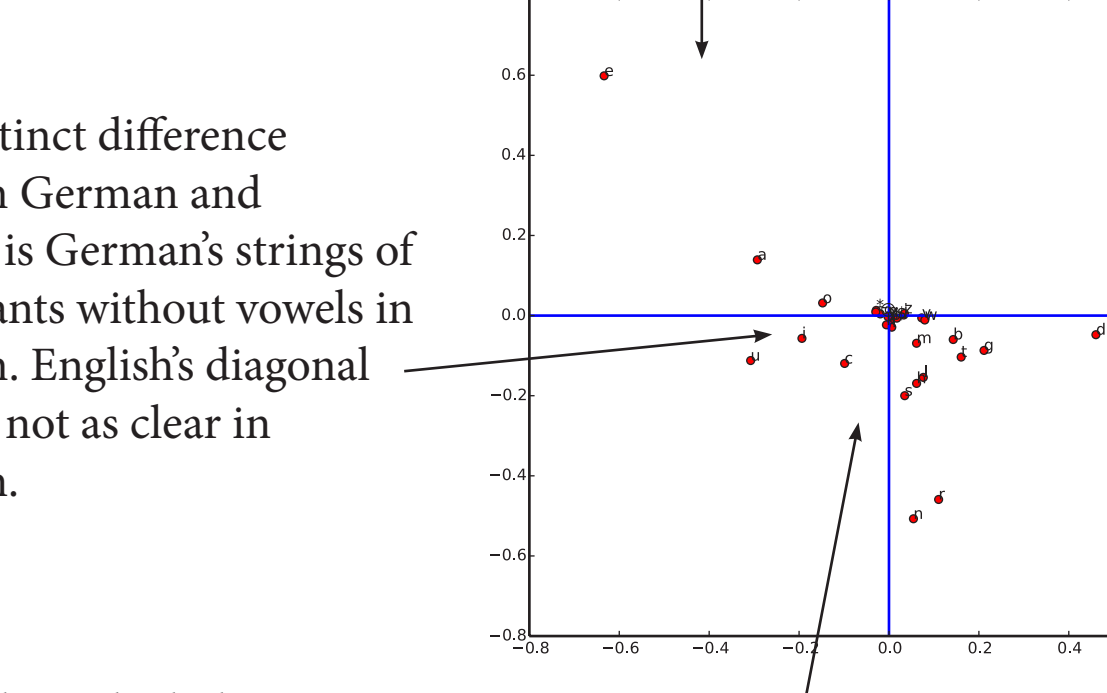
D occurs at the beginning of the articles: der, die, das.

E is more common than in English, especially at the end of words.



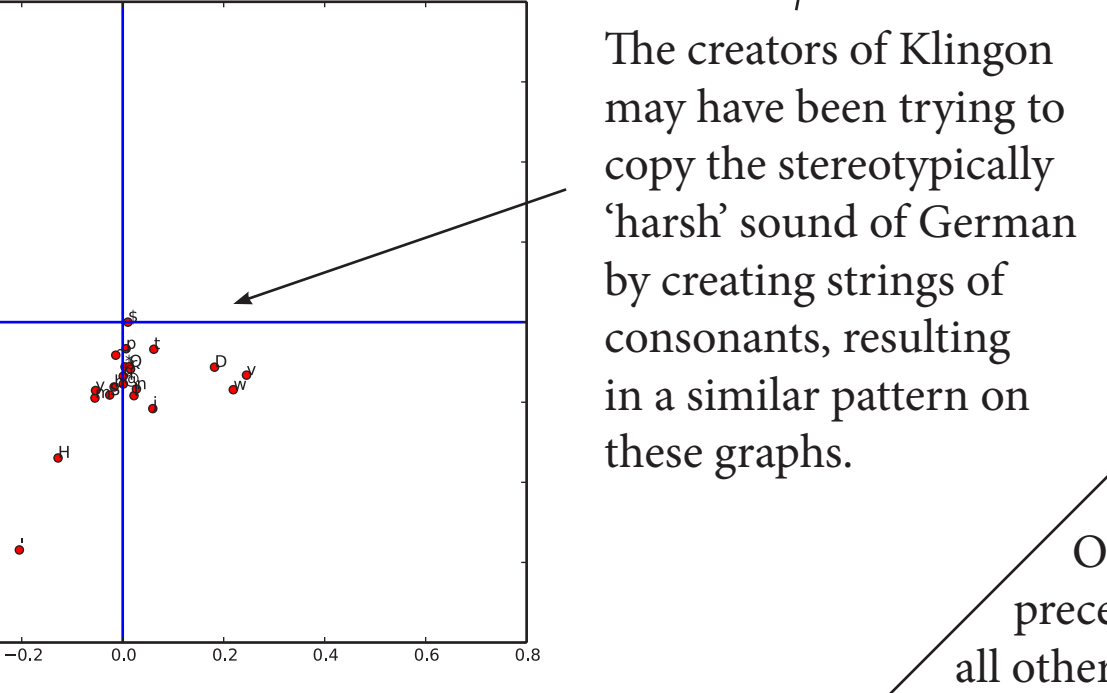
One distinct difference between German and English is German's strings of consonants without vowels in between. English's diagonal trend is not as clear in German.

The creators of Klingon may have been trying to copy the stereotypically ‘harsh’ sound of German by creating strings of consonants, resulting in a similar pattern on these graphs.



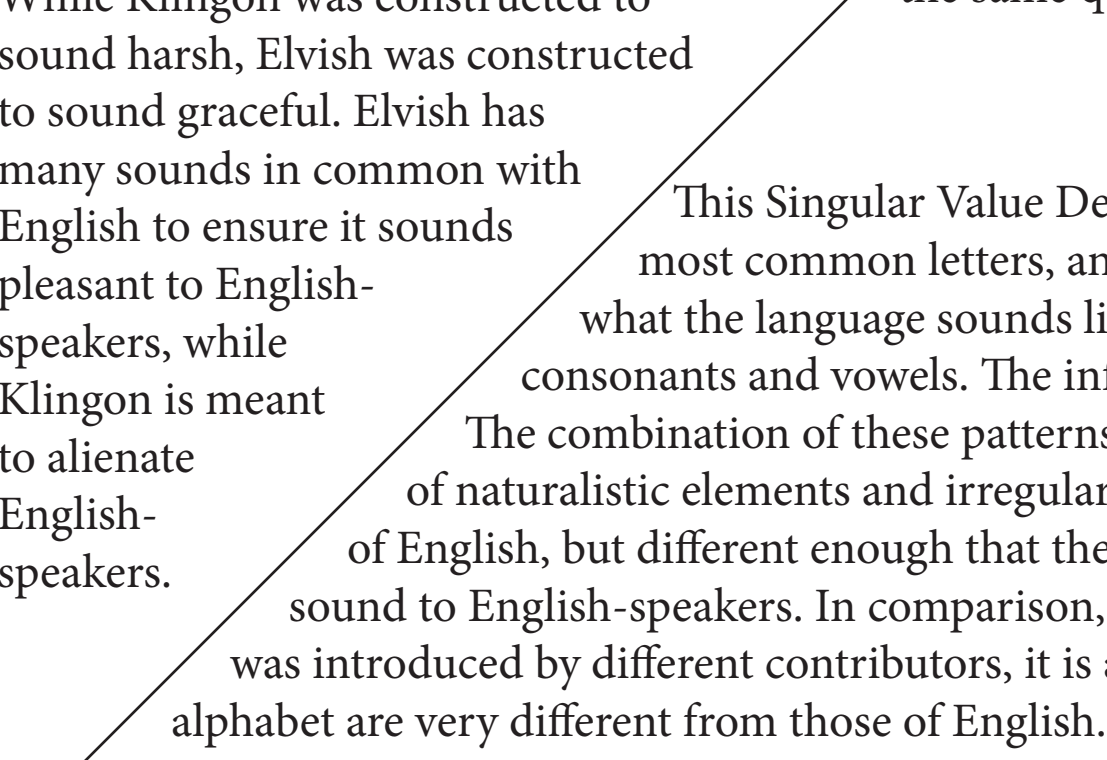
While Klingon was constructed to sound harsh, Elvish was constructed to sound graceful. Elvish has many sounds in common with English to ensure it sounds pleasant to English-speakers, while Klingon is meant to alienate English-speakers.

The letters ‘th’ followed by a vowel occurs in the middle of many words, accounting for h being an outlier in this quadrant. Since the graph is flipped, this is actually the same place as English’s ‘H’.



Elvish has a strong consonant-vowel pattern like English, but since its most common letter, ‘N’, is a consonant, this method of analysis has placed consonants in the second quadrant and vowels in the fourth.

The letters ‘th’ followed by a vowel occurs in the middle of many words, accounting for h being an outlier in this quadrant. Since the graph is flipped, this is actually the same place as English’s ‘H’.



More often preceded by consonants

More often preceded by vowels

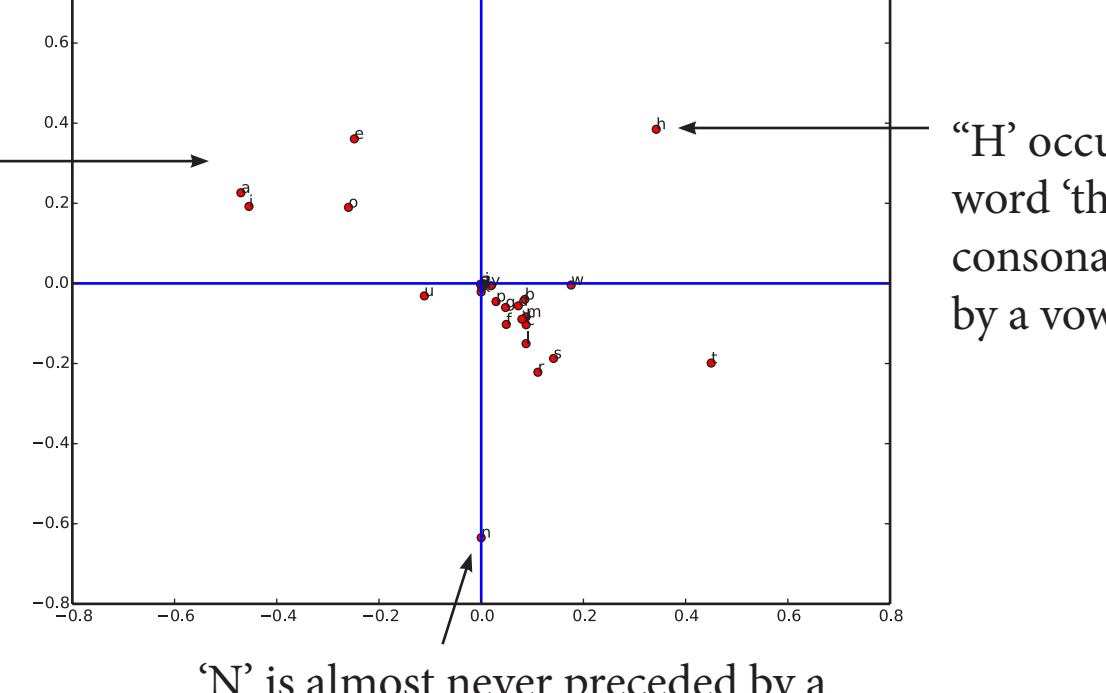
More often followed by consonants

More often followed by vowels



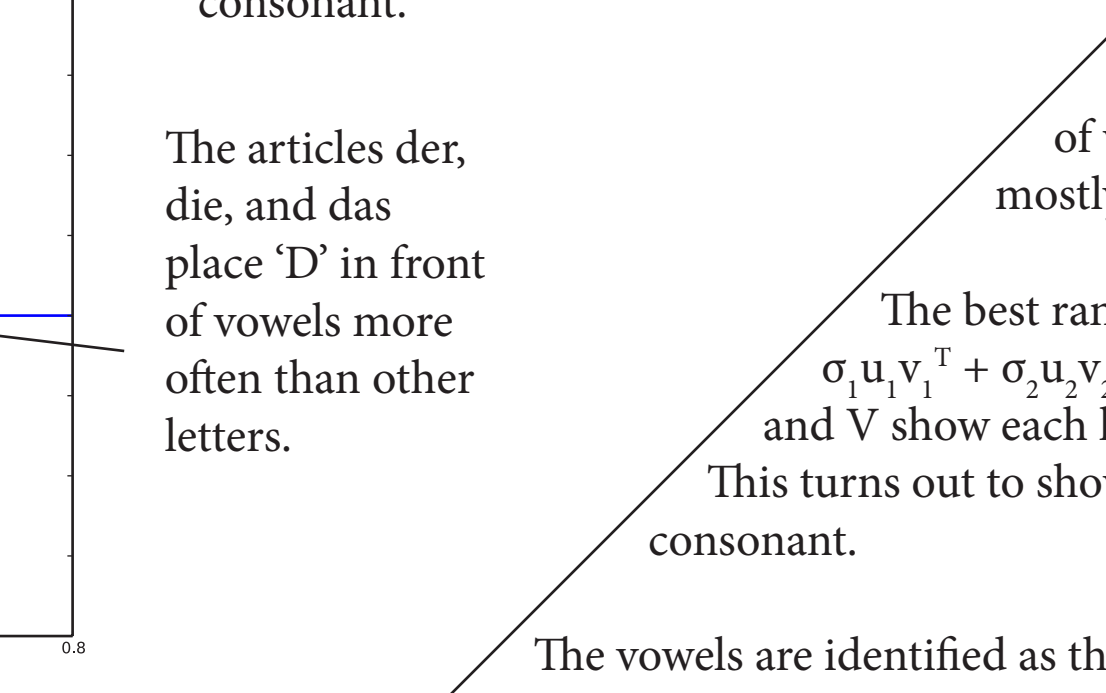
“H” occurs mostly in the word ‘the’, preceded by a consonant and followed by a vowel.

‘N’ is almost never preceded by a consonant.



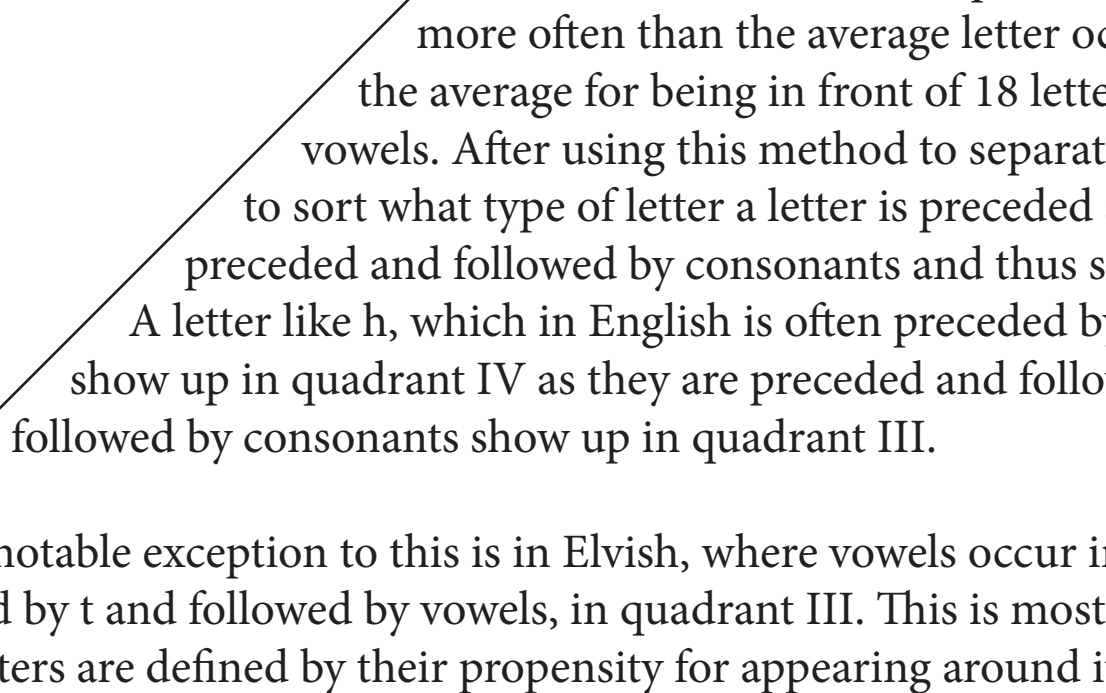
The articles der, die, and das place ‘D’ in front of vowels more often than other letters.

The best rank two approximation of the original adjacency matrix A is $σ_1u_1v_1^T + σ_2u_2v_2^T$, so our graphs of the i^{th} values of the second vectors of U and V show each letter's deviation from the point indicated by the first graphs. This turns out to show what letters are most commonly preceded by a vowel or a consonant.



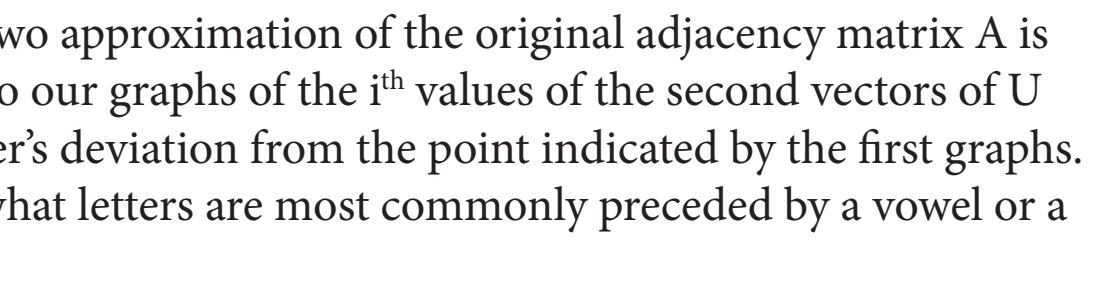
The vowels are identified as the letters that surpass the average frequency per letter for the most letters. For example, in English the letter A occurs before the letter B significantly more often than the average letter occurring before B. In fact, A has a higher frequency than the average for being in front of 18 letters, more than half. A similar pattern can be seen for all vowels. After using this method to separate vowels from consonants, the second singular value is able to sort what type of letter a letter is preceded and followed by. In English and German vowels are typically preceded and followed by consonants and thus show up in quadrant II of the secondary singular value graph. A letter like h, which in English is often preceded by t and followed by e, shows up in quadrant I. Consonants show up in quadrant IV as they are preceded and followed by vowels and letters that are preceded by vowels and followed by consonants show up in quadrant III.

One notable exception to this is in Elvish, where vowels occur in quadrant IV, consonants in quadrant II, and h, also preceded by t and followed by vowels, in quadrant III. This is most likely because the most common letter is n, a consonant, and all other letters are defined by their propensity for appearing around it. Other consonants do not appear near it, and are sorted into the same quadrant, while vowels are sorted into the quadrant opposite. This results in a graph flipped over the y=x axis.



In a first-rank approximation, all the rows are scalar multiples of v_1^T and all the columns are scalar multiples of u_1 . Therefore, the i^{th} element of u_1 increases with the frequency that the i^{th} letter in the alphabet precedes other letters, and the i^{th} element of v_1 increases with the frequency that the i^{th} letter in the alphabet follows other letters. Plotting these two values on the Cartesian plane as (u_1, v_1) shows both the relative frequency of the letters, with more common letters landing further from the origin, and the frequency of those letters beginning or ending a word, with letters occurring mostly at the beginning of words landing further from the x-axis and letters occurring mostly at the end of words landing further from the y-axis.

The best rank two approximation of the original adjacency matrix A is $σ_1u_1v_1^T + σ_2u_2v_2^T$, so our graphs of the i^{th} values of the second vectors of U and V show each letter's deviation from the point indicated by the first graphs. This turns out to show what letters are most commonly preceded by a vowel or a consonant.



Conclusion

This Singular Value Decomposition analysis revealed some fascinating patterns in these languages. In the first graphs we can easily identify the most common letters, and in the second we can see patterns of vowels and consonants. Between these two, we can see a graphical representation of what the language sounds like. German and Klingon have combinations of consonants, while English and Elvish follow a stricter pattern of alternating consonants and vowels. The influence of common articles, pronouns, and suffixes causes outliers like ‘h’ in English and Elvish. The combination of these patterns reveals key differences between the two constructed languages. J.R.R. Tolkien was known for his careful incorporation of naturalistic elements and irregularity into his constructed languages, and it shows on the graphs. Elvish has letter patterns and outliers similar to those of English, but different enough that they clearly constitute another language, just like German's. This was probably meant to give Elvish a pleasant but foreign sound to English-speakers. In comparison, Klingon's rules are very strict. The language was constructed primarily by fans of Star Trek, and while some variance was introduced by different contributors, it is a less natural and more prescriptive language. It was constructed to sound harsh and alien, so its letter patterns and even alphabet are very different from those of English. Its letter patterns more closely resemble those of German, which English-speakers tend to regard as harsh-sounding.