

Letter frequency

From Wikipedia, the free encyclopedia

The most common letter in the English alphabet is E. The **frequency of letters** in text has often been studied for use in cryptography, and frequency analysis in particular. No exact letter frequency distribution underlies a given language, since all writers write slightly differently. Linotype machines assumed the letter order, from most to least common, to be **etaoin shrdlu cmfwyp vbghjk xz** based on the experience and custom of manual compositors. Likewise, Modern International Morse code encodes the most frequent letters with the shortest symbols; arranging the Morse alphabet into groups of letters that require equal amounts of time to transmit, and then sorting these groups in increasing order, yields **e it san hurdm wgvfbk opjxcz yq**. Similar ideas are used in modern data-compression techniques such as Huffman coding.

Contents

- 1 Introduction
- 2 Relative frequencies of letters in the English language
- 3 Relative frequencies of the first letters of a word in the English language
- 4 Relative frequencies of letters in other languages
- 5 See also
- 6 References
- 7 External links

Introduction

More recent analyses show that letter frequencies, like word frequencies, tend to vary, both by writer and by subject. One cannot write an essay about x-rays without using frequent Xs, and the essay will have an especially strange letter frequency if the essay is about the frequent use of x-rays to treat zebras in Qatar. Different authors have habits which can be reflected in their use of letters. Hemingway's writing style, for example, is visibly different from Faulkner's. Letter, bigram, trigram, word frequencies, word length, and sentence length can be calculated for specific authors, and used to prove or disprove authorship of texts, even for authors whose styles are not so divergent.

Accurate average letter frequencies can only be gleaned by analyzing a large amount of representative text. With the availability of modern computing and collections of large text corpora, such calculations are easily made. This Deafandblind link (http://deafandblind.com/word_frequency.htm) details examples from a variety of sources, (press reporting, religious text, scientific text and general fiction) and there are differences especially for general fiction with the position of 'h' and 'i'. The example differs from the linotype 'etaoin shrdlu' to come out as 'etaoHn lsrldu'.

Herbert S. Zim, in his classic introductory cryptography text "Codes and Secret Writing", gives the English letter frequency sequence as "ETAON RISHD LFCMU GYPWB VKJXQ Z", the most common letter pairs as "TH HE AN RE ER IN ON AT ND ST ES EN OF TE ED OR TI HI AS TO", and the most common doubled letters as "LL EE SS OO TT FF RR NN PP CC".^[1]

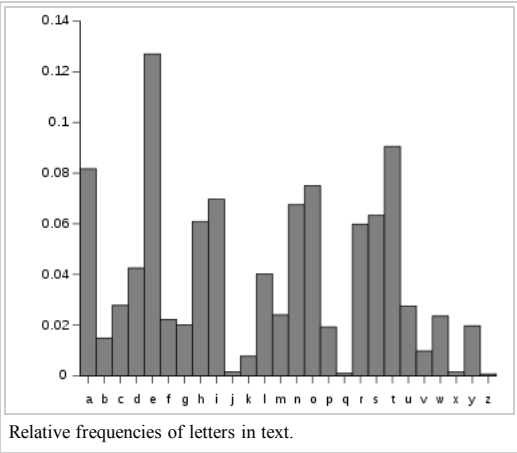
The 'top twelve' letters comprise about 80% of the total usage. The 'top eight' letters comprise about 65% of the total usage. Letter frequency as a function of rank can be fitted well by several rank functions, with the two-parameter Cocho/Beta rank function being the best.^[2] Another rank function with no adjustable free parameter also fits the letter frequency distribution reasonably well^[3] (the same function has been used to fit the amino acid frequency in protein sequences.^[4]) A spy using the VIC cipher or some other cipher based on a straddling checkerboard typically uses a mnemonic such as "a sin to err" (dropping the second "r") to remember the top 8 characters.

The use of letter frequencies and frequency analysis plays a fundamental role in cryptograms and several word puzzle games, including Hangman, Scrabble, Bananagrams, and the television game show *Wheel of Fortune*. One of the earliest description in classical literature of applying the knowledge of English letter frequency to solving a cryptogram is found in E.A. Poe's famous story *The Gold-Bug*, where the method is successfully applied to decipher a message instructing on the whereabouts of a treasure hidden by Captain Kidd.^[5]

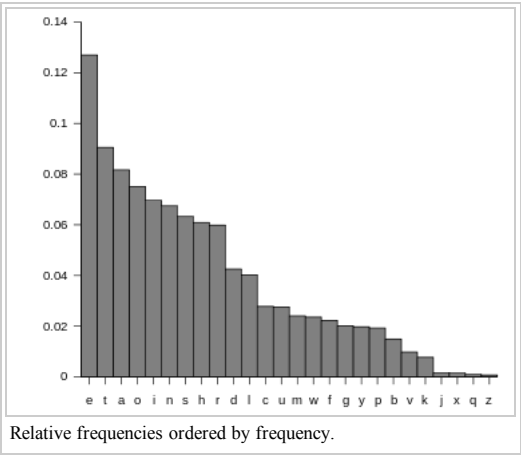
Letter frequencies had a strong effect on the design of some keyboard layouts. The most-frequent letters are on the bottom row of the Blickensderfer typewriter. The most-frequent letters are on the home row of the Dvorak Simplified Keyboard.

Relative frequencies of letters in the English language

Analysis of entries in the Concise Oxford dictionary is published by the compilers.^[6] The table below is taken from Pavel Mička's website, which cites Robert Lewand's *Cryptological Mathematics*.^[7]



Letter	Relative frequency in the English language	
a	8.167%	
b	1.492%	
c	2.782%	
d	4.253%	
e	12.702%	
f	2.228%	
g	2.015%	
h	6.094%	
i	6.966%	
j	0.153%	
k	0.772%	
l	4.025%	
m	2.406%	
n	6.749%	
o	7.507%	
p	1.929%	
q	0.095%	
r	5.987%	
s	6.327%	
t	9.056%	
u	2.758%	
v	0.978%	
w	2.360%	
x	0.150%	
y	1.974%	
z	0.074%	



This table differs slightly from others, such as Cornell University Math Explorer’s Project, which produced a table after measuring 40,000 words.^[8]

In English, the space is slightly more frequent than the top letter (e) ^[9] and the non-alphabetic characters (digits, punctuation, etc.) collectively occupy the fourth position, between *t* and *a*.^[10]

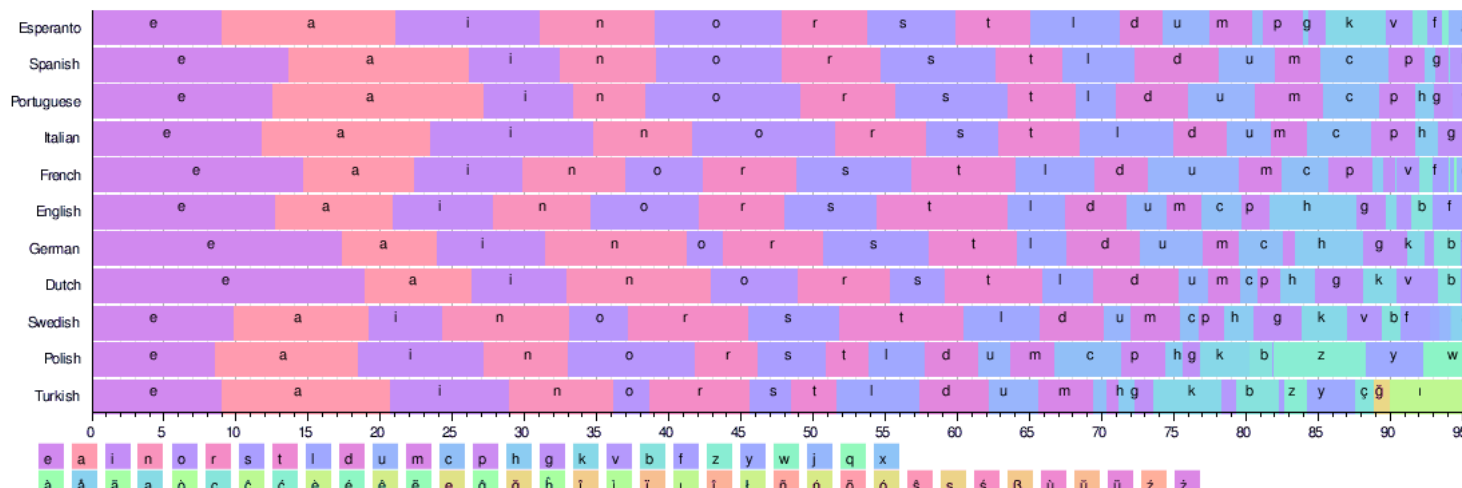
Relative frequencies of the first letters of a word in the English language

First Letter of a word frequencies:^[11]

v	1.628%	0.846%	1.138%	1.665%	1.904%	2.097%	0.977%	2.415%	0	1.854%	2.332%	2.437%
w	0.074%	1.921%	0.017%	0.037%	0	0.033%	0.016%	0.142%	6.313%	1.821%	0.069%	0
x	0.427%	0.034%	0.215%	0.253%	0	0	0.007%	0.159%	0	0.036%	0.028%	0.046%
y	0.128%	0.039%	1.008%	0.006%	0	0.020%	3.371%	0.708%	3.206%	0.035%	0.698%	0.900%
z	0.326%	1.134%	0.517%	0.470%	0.494%	1.181%	1.497%	0.070%	5.852%	1.374%	0.034%	0
à	0.486%	0	0	0.072%	0	0.635%	0	0	0	0	0	0
â	0.051%	0	0	0.562%	0	0	0	0	0	0	0	0
á	0	0	0.502%	0.118%	0	0	0	0	0	0	0	1.799%
â	0	0	0	0	0	0	0	1.338%	0	0	1.190%	0
ä	0	0.447%	0	0	0	0	0	1.797%	0	0	0	0
ã	0	0	0	0.733%	0	0	0	0	0	0	0	0
ä	0	0	0	0	0	0	0	0	0.699%	0	0	0
æ	0	0	0	0	0	0	0	0	0	0	0.872%	0.867%
œ	0.018%	0	0	0	0	0	0	0	0	0	0	0
ç	0.085%	0	0	0.530%	0	0	0.825%	0	0	0	0	0
ê	0	0	0	0	0.657%	0	0	0	0	0	0	0
é	0	0	0	0	0	0	0	0	0.743%	0	0	0
ö	0	0	0	0	0	0	0	0	0	0	0	4.393%
è	0.271%	0	0	0	0	0.263%	0	0	0	0	0	0
é	1.504%	0	0.433%	0.337%	0	0	0	0	0	0	0	0.647%
ê	0.225%	0	0	0.450%	0	0	0	0	0	0	0	0
ë	0.001%	0	0	0	0	0	0	0	0	0	0	0
ę	0	0	0	0	0	0	0	0	1.035%	0	0	0
ĝ	0	0	0	0	0.691%	0	0	0	0	0	0	0
ğ	0	0	0	0	0	0	1.129%	0	0	0	0	0
ĥ	0	0	0	0	0.022%	0	0	0	0	0	0	0
î	0.045%	0	0	0	0	0	0	0	0	0	0	0
ì	0	0	0	0	0	0.030%	0	0	0	0	0	0
í	0	0	0.725%	0.132%	0	0	0	0	0	0	0	1.570%
ï	0.005%	0	0	0	0	0	0	0	0	0	0	0
ı	0	0	0	0	0	0	5.199%*	0	0	0	0	0
ĵ	0	0	0	0	0.055%	0	0	0	0	0	0	0
l	0	0	0	0	0	0	0	0	2.109%	0	0	0
ñ	0	0	0.311%	0	0	0	0	0	0	0	0	0
ń	0	0	0	0	0	0	0	0	0.362%	0	0	0
ò	0	0	0	0	0	0.002%	0	0	0	0	0	0
ö	0	0.573%	0	0	0	0	0.270%	1.305%	0	0	0	0.777%
ô	0.023%	0	0	0.635%	0	0	0	0	0	0	0	0
ó	0	0	0.827%	0.296%	0	0	0	0	1.141%	0	0	0.994%
ø	0	0	0	0	0	0	0	0	0	0	0.939%	0
ș	0	0	0	0	0.385%	0	0	0	0	0	0	0
ş	0	0	0	0	0	0	1.938%	0	0	0	0	0
ś	0	0	0	0	0	0	0	0	0.814%	0	0	0
ß	0	0.307%	0	0	0	0	0	0	0	0	0	0
þ	0	0	0	0	0	0	0	0	0	0	0	1.455%
ù	0.058%	0	0	0	0	0.166%	0	0	0	0	0	0
ú	0	0	0.168%	0.207%	0	0	0	0	0	0	0	0.613%
ü	0	0	0	0	0.520%	0	0	0	0	0	0	0
ü	0	0.995%	0.012%	0.026%	0	0	1.992%	0	0	0	0	0
ý	0	0	0	0	0	0	0	0	0	0	0	0.228%
ź	0	0	0	0	0	0	0	0	0.078%	0	0	0
ž	0	0	0	0	0	0	0	0	0.706%	0	0	0

*See Dotted and dotless I

The figure below illustrates the frequency distributions of the 26 most common Latin letters across some languages.



Based on these tables, the 'etaoin shrldu'-equivalent results for each language is as follows:

- French: 'esait nrulo'; (Indo-European: Romance; traditionally, 'esartinulop' is used, in part for its ease of pronunciation^[24])
- Spanish: 'eaosr nidlt'; (Indo-European: Romance)
- Portuguese: 'aeosr indmt' (Indo-European: Romance)
- Italian: 'eaion lrtsc'; (Indo-European: Romance)
- Esperanto: 'aieon lsrtk' (artificial language – influenced by Indo-European languages, Romance, Germanic mostly)
- German: 'enisr atdhu'; (Indo-European: Germanic)
- Swedish: 'eantr isldo'; (Indo-European: Germanic)
- Turkish: 'aeinr ldkmu'; (Altaic: Turkic)
- Dutch: 'enati rods'l'; (Indo-European: Germanic)^[20]
- Polish: 'aoien wszrd'; (Indo-European: Slavic)

All these languages use a basically similar 25+ character alphabet.

See also

- Corpus linguistics
- ETAOIN SHRDLU
- RSTLNE (Wheel of Fortune)
- Frequency analysis (cryptanalysis)
- Linotype machine
- Most common words in English
- Scrabble
- Arabic Letter Frequency

References

- ¹ ^ Zim, Herbert Spencer. (1961). *Codes & Secret Writing: Authorized Abridgement*. Scholastic Book Services. OCLC 317853773 (<http://www.worldcat.org/oclc/317853773>).
- ² ^ Li, Wentian; Miramontes, Pedro (2011). "Fitting ranked English and Spanish letter frequency distribution in US and Mexican presidential speeches". *Journal of Quantitative Linguistics* **18** (4): 359. doi:10.1080/09296174.2011.608606 (<http://dx.doi.org/10.1080/09296174.2011.608606>).
- ³ ^ Gusein-Zade, S.M. (1988). "Frequency distribution of letters in the Russian language". *Probl. Peredachi Inf.* **24** (4): 102–7.
- ⁴ ^ Gamow, George; Ycas, Martynas (1955). "Statistical correlation of protein and ribonucleic acid composition" (<http://www.pnas.org/content/41/12/1011.full.pdf>). *Proc. Natl. Acad. Sci.* **41** (12): 1011–19. doi:10.1073/pnas.41.12.1011 (<http://dx.doi.org/10.1073/pnas.41.12.1011>). PMC 528190 (<http://www.ncbi.nlm.nih.gov/pmc/articles/PMC528190>).
- ⁵ ^ Poe, Edgar Allan. "The works of Edgar Allan Poe in five volumes" (http://www.gutenberg.org/catalog/world/readfile?fk_files=1977099). Project Gutenberg.
- ⁶ ^ "What is the frequency of the letters of the alphabet in English?" (<http://oxforddictionaries.com/words/what-is-the-frequency-of-the-letters-of-the-alphabet-in-english>). *Oxford Dictionary*. Oxford University Press. Retrieved 29 December 2012.
- ⁷ ^ Mička, Pavel. "Letter frequency (English)" (<http://en.algoritmy.net/article/40379/Letter-frequency-English>). Algoritmy.net.
- ⁸ ^ <http://www.math.cornell.edu/~mec/2003-2004/cryptography/subs/frequencies.html>
- ⁹ ^ Statistical Distributions of English Text (<http://www.data-compression.com/english.html>)
- ¹⁰ ^ Lee, E. Stewart. "Essays about Computer Security" (<http://www.cl.cam.ac.uk/~mgk25/lee-essays.pdf>) (PDF). University of Cambridge Computer Laboratory. p. 181.
- ¹¹ ^ Calculated from "Project Gutenberg Selections" available from the NLTK Corpora (http://nltk.googlecode.com/svn/trunk/nltk_data/index.xml)
- ¹² ^ "CorpusDeThomasTempé" (<http://gpl.insa-lyon.fr/Dvorak-Fr/CorpusDeThomasTemp%C3%A9>). Retrieved 2007-06-15.
- ¹³ ^ Beutelspacher, Albrecht (2005). *Kryptologie* (7 ed.). Wiesbaden: Vieweg. p. 10. ISBN 3-8348-0014-7.
- ¹⁴ ^ Pratt, Fletcher (1942). *Secret and Urgent: the Story of Codes and Ciphers*. Garden City, N.Y.: Blue Ribbon Books. pp. 254–5. OCLC 795065 (<http://www.worldcat.org/oclc/795065>).
- ¹⁵ ^ "Frequência da ocorrência de letras no Português" (<http://www.numaboia.com/criptografia/criptoanalise/310-Frequencia-no-Portugues>). Retrieved 2009-06-16.
- ¹⁶ ^ "La Ofitecoj de la Esperantaj Literoj" (<http://lingvakritiko.com/2007/09/13/literofitecoj-kaj-tabelvortofitecoj/>). Retrieved 2007-09-14.
- ¹⁷ ^ Singh, Simon; Galli, Stefano (1999). *Codici e Segreti* (in Italian). Milano: Rizzoli. ISBN 978-8-817-86213-4. OCLC 535461359 (<http://www.worldcat.org/oclc/535461359>).
- ¹⁸ ^ "Practical Cryptography" (<http://practicalcryptography.com/cryptanalysis/letter-frequencies-various-languages/swedish-letter-frequencies/>). Retrieved 2013-10-30.
- ¹⁹ ^ Wstep do kryptologii (<http://www.korzen.org/wsiziz/wstep%20do%20kryptologii/wyklady/Wykad%202.doc>), counting [space] 17.2%, [dot point] 0.9%, [comma] 0.9% and [semicolon] 0.5%

20. ^ *a b* "Letterfrequenties" (<http://www.onzetaal.nl/advies/letterfreq.php>). *Genootschap OnzeTaal*. Retrieved 2009-05-17.
21. ^ "Practical Cryptography" (<http://practicalcryptography.com/cryptanalysis/letter-frequencies-various-languages/danish-letter-frequencies/>). Retrieved 2013-10-24.
22. ^ "Practical Cryptography" (<http://practicalcryptography.com/cryptanalysis/letter-frequencies-various-languages/icelandic-letter-frequencies/>). Retrieved 2013-10-24.
23. ^ "Practical Cryptography" (<http://practicalcryptography.com/cryptanalysis/letter-frequencies-various-languages/finnish-letter-frequencies/>). Retrieved 2013-10-24.
24. ^ Perec, Georges; *Alphabets*; Éditions Galilée, 1976

Notes

Some useful tables for single letter, digram, trigram, tetragram, and pentagram frequencies based on 20,000 words that take into account word-length and letter-position combinations for words 3 to 7 letters in length. The references are as follows:

1. Mayzner, M.S.; Tresselt, M.E. (1965). "Tables of single-letter and digram frequency counts for various word-length and letter-position combinations". *Psychonomic Monograph Supplements* **1** (2): 13–32. OCLC 639975358 ([//www.worldcat.org/oclc/639975358](http://www.worldcat.org/oclc/639975358)).
2. Mayzner, M.S.; Tresselt, M.E.; Wolin, B.< R.< (1965). "Tables of trigram frequency counts for various word-length and letter-position combinations". *Psychonomic Monograph Supplements* **1** (3): 33–78.
3. Mayzner, M.S.; Tresselt, M.E.; Woliin, B.< R,.. (1965). "Tables of tetragram frequency counts for various word-length and letter-position combinations". *Psychonomic Monograph Supplements* **1** (4): 79–143.
4. Mayzner, M.S.; Tresselt, M.E.Wolin, B.< R.> (1965). "Tables of pentagram frequency counts for various word-length and letter-position combinations". *Psychonomic Monograph Supplements* **1** (5): 144–190.

External links

- A site with content of *Cryptographical Mathematics* (<http://pages.central.edu/emp/LintonT/>) by Robert Edward Lewand
- Some examples of letter frequency rankings in some common languages (<http://www.bckelk.ukfsn.org/words/etaoin.html>)
- Java-Application for building letter frequencies out of a text file (http://www.imn.htwk-leipzig.de/~dborkman/offtopic/letter_frequency/letter.html)
- JavaScript Heatmap Visualization showing letter frequencies of texts on different keyboard layouts (<http://www.patrick-wied.at/projects/heatmap-keyboard/>)
- An updated version of Mayzner's work using Google books Ngrams data set (<http://norvig.com/mayzner.html>) by Peter Norvig

Retrieved from "http://en.wikipedia.org/w/index.php?title=Letter_frequency&oldid=583665039"

Categories: Cryptography | Linguistics | Linguistics terminology | Quantitative linguistics

-
- This page was last modified on 28 November 2013 at 14:00.
 - Text is available under the Creative Commons Attribution-ShareAlike License; additional terms may apply. By using this site, you agree to the Terms of Use and Privacy Policy.
- Wikipedia® is a registered trademark of the Wikimedia Foundation, Inc., a non-profit organization.