

K-GRAM INDEX

USED TO APPROXIMATE EDIT DISTANCE
(WHICH IS TOO EXPENSIVE) IN THE CASE OF
MORE THAN ONE SPELL ERROR (WHICH THE
1 ERROR CORRECTION APPROACH DON'T COVER, OFC).

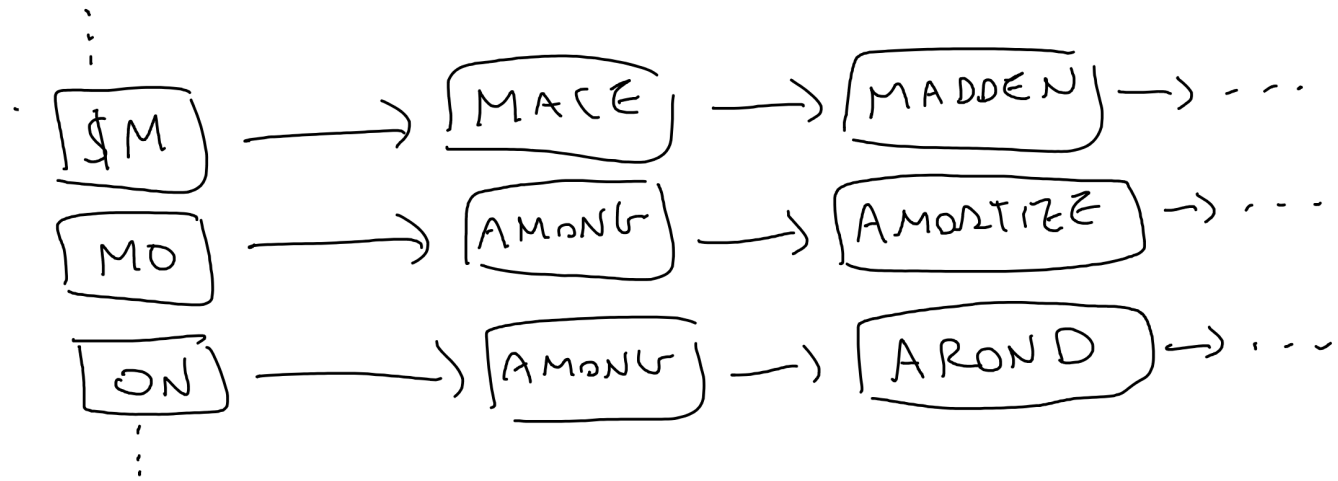
THE K-GRAM INDEX IS AN INVERTED INDEX CONTAINING,
FOR EACH K-GRAM, ALL TERMS INCLUDING THAT
K-GRAM.

Let's see how to build it and use it for our problem

(1) First off, we have to build the
K-gram index using the terms from our
lexicon. Let's assume that each term
is anticipated by $K-1$ special char $\$$,
in order to ensure that the number of
K-grams is equal to the length of
the string generating them.

GENERATE ALL POSSIBLE OVERLAPPING KGRAM
FOR EACH TERM AND THEN USE
THEM TO BUILD THE INVERTED LISTS

EXAMPLE OF K-GRAM INDEX ($K=2$)



NOTE: I shown only the K-GRAMS of the term "MON", BUT THE K-GRAM $\boxed{\$M}$ IS GENERATED BY ALL STRINGS STARTING WITH M

NOTE: GIVEN A TERM OF SIZE $|t|$, THE NUMBER OF RESULTING KGRAMS IS $|t|$

$\$CASA \Rightarrow \C, CA, AS, SA

(2) NOW WE CAN CHECK FOR ERRORS IN A QUERY USING THE K-GRAM INDEX.

ENUMERATE ALL OVERLAPPING KGRAMS IN A QUERY

AND SEARCH THEM IN THE INDEX, THEN

KEEP TRACK OF HOW MANY K-GRAMS ARE IN COMMON BETWEEN Q AND THE TERMS OF THE LEXICON

(3) SELECT THE POSSIBLE CORRECT TERMS BY THRESHOLDING ON THE K-GRAMS MATCHING THE MOST

- SINCE OUR QUERY TERM IS OF LENGTH $|Q|$
- GIVEN e THE MAXIMUM AMOUNT OF ALLOWED ERRORS, $e * K$ K-GRAMS OF OUR QUERY TERM MIGHT BE DIFFERENT FROM THE K-GRAMS OF A TERM IN THE LEXICON.

THEREFORE, AT LEAST $|Q| - e * K$
K-GRAMS OF THE QUERY TERM
MATCH WITH THOSE OF A
LEXICON'S TERM.

$$|Q| - e * K \geq \text{\# matching K-grams with a term}$$

WARNING: THIS IS ONLY AN APPROXIMATION
AND WE MIGHT WANT TO DO THE \subseteq
ANYWAY

EXAMPLE OF WARWING

GIVEN $K=3$, COMPARE $$$NOVEMBER$
(FROM LEXICON) WITH $$$DECEMBER (4)$

$$$NOVEMBER \Rightarrow $$N, $NO, NOV, OVE, VEM, \underline{EMB}, \underline{MBE}, \underline{BER}$

$$$DECEMBER \Rightarrow $$D, $DE, DEC, ECE, CEM, \underline{EMB}, \underline{MBE}, \underline{BER}$

if $e=1$:

$$|4| - e^* K = 8 - 3 = 5 \quad \text{NO}$$

if $e=2$:

$$|4| - e^* K = 8 - 6 = 2 \quad \text{OK}$$

BUT WE NEED E.D. TO SAY

THAT THE DISTANCE IS ACTUALLY 3

EXAMPLE

GIVEN $S = \{ \text{PITOM, DAD, DADDY, ZOOM} \}$

(1) BUILD THE CORRESPONDING ZGRAM INDEX

(2) SHOW THE EXECUTION OF THE 1-EDIT EDIT SEARCH ON S USING THE INDEX, GIVEN $P = \text{ATOM}$

(1) $\$PITOM, \$DAD, \$DADDY, \$ZOOM$
 1 2 3 4

$\$P \rightarrow 1$

$PI \rightarrow 1$

$IT \rightarrow 1$

$TO \rightarrow 1$

$OM \rightarrow 1, 4$

$\$D \rightarrow 2, 3$

$DA \rightarrow 2, 3$

$AD \rightarrow 2, 3$

$DD \rightarrow 3$

$DY \rightarrow 3$

$\$Z \rightarrow 4$

$ZO \rightarrow 4$

$OO \rightarrow 4$

$\$D$

$\$P$

$\$Z$



SORT
FOR
BETTER
I/Os

② \$ATOM

K grams = \$A, AT, TO, OM

Let's check the matches with the KGRAM INDEX

\$A $\rightarrow \phi$

AT $\rightarrow \phi$

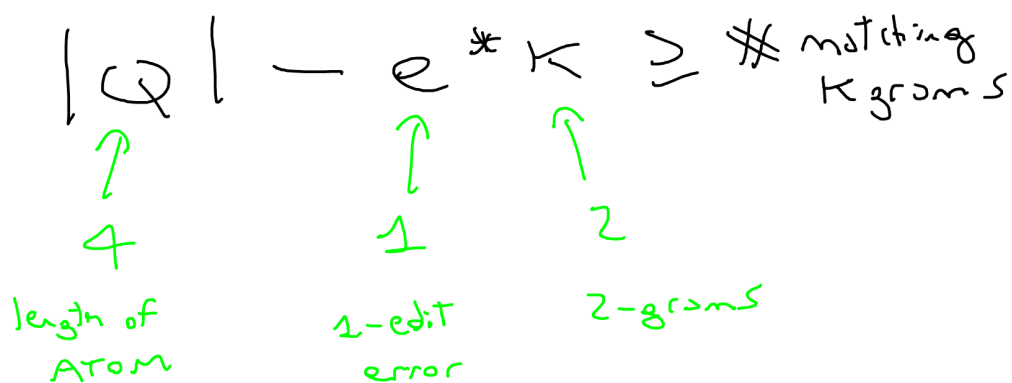
TO $\rightarrow 1$

OM $\rightarrow 1, 4$

2 CANDIDATES = PITOM, ZOOM

THE CORRECT ONE IS THE ONE
SUCH THAT

$$|Q| - e^* K \geq \# \text{ matching Kgrams}$$



4 1 2

length of ATOM 1-edit error 2-grams

$$4 - 2 = 2$$

the correct candidate
is PITOM