

SOUNDEX

CLASS OF HEURISTICS THAT EXPAND
A QUERY INTO PHONETIC EQUIVALENTS.
THESE ALSO ARE LANGUAGE SPECIFIC
AND MAINLY USED FOR NAMES
(CHEBYSNEV \rightarrow TCHEBJCHEFF).

IN ORDER TO DO THIS, IT COMPUTES A
REDUCED FORM OF 4 CHARS OF BOTH
THE LEXICON AND THE QUERY TERM

Why?

Even if the queried name is spelled wrong
(Ex: Herman instead of Hermann), THE REDUCED
FORM IS THE SAME

SSS

THE SOUNDEX FAMILY IS NOT USED IN IR,
BECAUSE THERE ARE BETTER ALTERNATIVES.

Typical algo

① Take in input & word

EXAMPLE: Herman

② Retain the first letter of the word

EXAMPLE: H

③ Change each occurrence of A, E, I, O, U, H, W, Y with zero 0 and topped it to the first letter

EXAMPLE H0rmond

④ Change letters to digit as follow:

• B, F, P, V \rightarrow 1

• C, G, J, K, Q, S, X, Z \rightarrow 2

• D, T \rightarrow 3

• L \rightarrow 4

• M, N \rightarrow 5

• R \rightarrow 6

EXAMPLE

H06505

⑤ Remove all pairs of consecutive equal digit

EXAMPLE: None, in this case

⑥ Remove all zeroes

Example: H6SS

⑦ PAD WITH TRAILING ZEROS AND
RETURNS THE FIRST 4 positions as

< uppercase letter > < digit > < digit > < digit >

NOTE: Hermann gives the same result

oss: You might notice that the way the
mapping is done in ⑦, the reduction
is biased towards certain nationality names