

# Povzetek

Pričujoče diplomsko delo obravnava v rudarjenje večjezičnih besedil. Na začetku predstavi osnove rudarjenja teksta s poudarkom na predstavitvi dokumentov, algoritem za iskanje po bazi dokumentov ter na kratko tudi algoritem za avtomatsko kategorizacijo dokumentov. Za predstavitev uporabi vrečo besed, ter preuči njene prednosti in slabosti. Iskanja dokumentov se loti z metodo najbližjih sosedov, kategorizacije pa z metodo podpornih vektorjev.

V nadaljevanju je podrobno predstavljena kanonična korelacijska analiza (KKA). To je metoda, ki za par slučajnih vektorjev poišče smeri, vzdolž katerih sta slučajna vektorja visoko korelirana. Podana je klasična definicija KKA ter njena posplošitev na metodo jeder. Ker ima klasična KKA težave pri delu z visoko dimenzionalnimi podatki se vpelje regularizacijo. Predstavljen je tudi numerični postopek za reševanje KKA.

Za konec je predstavljena še uporaba KKA pri rudarjenju večjezičnih besedil. Prikazani so rezultati opravljenih poskusov na večjezičnih bazah in prototip iskalnika po večjezičnih dokumentih. Le-ta se nahaja tudi na priloženi zgoščenki.

**Math. Subj. Class. (MSC 2000):** 62H20, 62H30, 62P99, 65F15, 68U15

**Ključne besede:**

rudarjenje teksta, večjezično iskanje, vreča besed, kanonična korelacijska analiza, metoda podpornih vektorjev, metode jeder

**Keywords:**

text mining, cross-lingual information retrieval, bag of words, canonical correlation analysis, support vector machine, kernel methods

# Literatura

- [1] L.W.BEINEKE, R.J.WILSON, *Selected topics in graph theory*, Academic Press Inc. (London) LTD., 1978.
- [2] P.J.CAMERON, A.G.CHETWYND, J.J.WATKINS, *Decomposition of snarks*, *J.Graph Th.*, 11, 1987, 13 - 19.
- [3] M.GARDNER, *Mathematical games*, Scientific American 234, No.4 (April 1976), 126 - 130, ter 234 No.9 (September 1976), 210 - 211.
- [4] M.K.GOLDBERG, *Construction of class 2 graphs with maximum vertex degree 3*, *J.Comb. Theory Ser.B* 31 (1981), 282 - 291.
- [5] D.A.HOLTON, J. SHEEHAN, *The Petersen Graph*, Austral. Math. Soc. Lecture Series 7, Cambridge University Press, 1993.
- [6] R.ISAACS, *Infinite families of non-trivial trivalent graph, which are not Tait-colorable*, *Amer. Math. Monthly* 82 (1975), 221 - 239.
- [7] R.ISAACS, *Loupekhine's snarks: a bifamily of not Tait-colorable graphs*, Tehnical report No.263, Dept. of Math. Sciences, John Hopkins University, 1976.
- [8] F.JAEGER, *Sur l'indice chromatique du graphe representatif des aretes d'un graphe regulier*, *Disc.Math.*, 9, 1974, 161 - 172.
- [9] R.KALINOWSKI, Z.SKUPIEN, *Large Isaccs' graphs are maximally non-hamiltonian-connected*, *Disc.math.*, 82, 1990, 101 - 107.
- [10] B.MOHAR, A.VODOPIVEC, *On polyhedral embeddings of cubic graphs*, članek.
- [11] A.ORBANIČ, T.PISANSKI, M RANDIČ, B SERVATIUS, *Blanuša double*, *Mathematical Communications* 9 (2004), 91 - 103.

- 
- [12] R.C.READ, R.J.WILSON, *An Atlas of graphs*, Oxford University Press, 1998, 276 - 281.
- [13] A.VODOPIVEC, *On embeddings of snarks in the torus*, članek.
- [14] DOUGLAS B.WEST, *Introduction to Graph Theory, second edition*, Prentice-Hall, Inc. Upper Saddle River, NJ 07458, 2001.
- [15] R.J.WILSON, J.J.WATKINS, *Uvod v teorijo grafov*, Društvo matematikov, fizikov in astronomov Slovenije, 1997, 296 - 307.
- [16] INTERNET:
- <http://mathworld.wolfram.com/Snark.html>
  - <http://en.wikipedia.org>
  - <http://www.csse.uwa.edu.au/gordon/remote/cubics/index.html>