

e-mail : kimgogo02@cs.yonsei.ac.kr

The gene prediction method considering stages of cancer, obtained by integrating gene expression, genetic interaction data and document

Jungrim Kim, Yunku Yeu, Sanghyun Park*
Dept. of Computer Science, Yonsei University

가 , 가 .

가 , (genetic interaction) , (heterogeneous) (disease)- (normal) 가

(stage) .

1. 가 , 가 .

(Genome project) 가 가 .

(Drug repositioning) .

. DNA [1] .

(Pearson Correlation Coefficient) 가 가 가 (differential expression) . [2, 3] ,

2. 2.1

가 BioGRID[4] BioGRID Release 3.2.104

(relation) 가

* : , e-mail: sanghyun@cs.yonsei.ac.kr 13,993 127,688

2013 () GEO (Gene Expression Omnibus)[5] (2012R1A2A1A01010775). GSE21815[6]

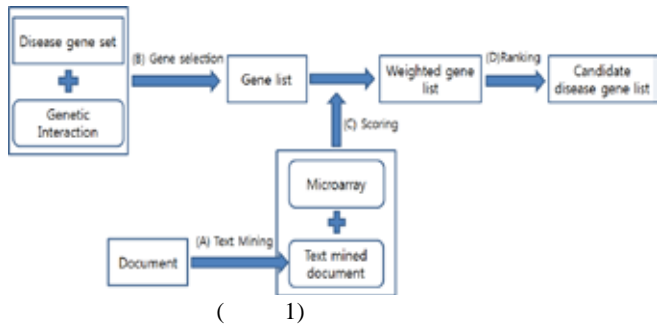
가 132
 9, 1
 12, 2 27, 3
 16, 4 11
 (gene (scoring)
 symbol) 가 (row)

< 1> KEGG NCI

HGNC (HUGO Gene Nomenclature Committee)[7]
 41,321
 PubMed[10] (colorectal cancer)
 157,740
 GeneCards version 3.10 [8,9] (colorectal
 cancer) 3,272
 2,024

KEGG	NCI
Oncogenes genes CTNNB1 (β – catenin), KRAS	Tumor suppressor genes APC, TP53, NEK4(STK11), PTEN, BMPR1A, SMAD4
Tumor suppressor genes APC, DCC, TGFBR2, SMAD2, SMAD4, BAX, TP53	Repair/stability genes MYH(MUTYH), MLH1, MSH2, MSH6, PMS2, EPCAM
DNA repair genes MLH1, MSH2, MSH3, MSH6	

2.2



(text mining), (disease gene set)
 (B) (gene selection),
 (training set) (C) (scoring), (D)
 (ranking)

KEGG (Kyoto Encyclopedia of Genes and
 Genomes)[11] NCI(National Cancer Institution) [12]
 NCI KEGG
 (cross validation)

2.2.1

, POS(Part of Speech) tagger[13]

HGNC
 (frequency)

2.2.2

KEGG NCI 1

2.2.3

(threshold)

Score_{PCC_i}

가

Score_{diff_i}

Score_{diff_i}

$$Score_{PCC_i} = \frac{\sum_{stage=1,2,3,4} P(i)}{N_{i,PCC}}$$

$$P(i) = \begin{cases} \text{if } PCC_{i,stage} > \delta_{PCC}, P(i) = PCC_{i,stage} \\ \text{else}, P(i) = 0 \end{cases} \quad (1)$$

$$Score_{diff_i} = \frac{\sum_{stage=1,2,3,4} D(i)}{N_{i,diff}}$$

$$D(i) = \begin{cases} \text{if } diff_{i,stage} > \delta_{diff}, D(i) = diff_{i,stage} \\ \text{else}, D(i) = 0 \end{cases} \quad (2)$$

score_{PCC_i} : i
 score_{diff_i} : i
 PCC_{i,stage} : stage i
 diff_{i,stage} : stage i
 δ_{PCC} :
 δ_{diff} :
 N_{i,PCC} : P(i) 0 stage
 N_{i,diff} : D(i) 0 stage

가 , , 100
 , 1 , 2 , 3 , 4 , 7
 KEGG KEGG
 6

Score_{freq_i}

$$Score_{freq_i} = \log_{10}(freq_i + 1) \quad (3)$$

score_{freq_i} : i
 freq_i : i

Score_{PCC_i} Score_{diff_i}, Score_{freq_i}
 (min-max normalization)

가

$$Score_i = w_1 * Score_{freq_i} + w_2 * Score_{PCC_i} + w_3 * Score_{diff_i} \quad (4)$$

score_i = i
 w₁ : 가
 w₂ : 가
 w₃ : 가

2.2.4

3.

가 w₁ = 1, w₂ = 1, w₃ = 2
 δ_{PCC} = 0.5, δ_{diff} = 0.2

. NCI KEGG

NCI

NCI

7

< 2> NCI

	Microarray	Document	Microarray + GI + Document
Top 30	1	2	2
Top 50	1	3	5
Top 100	1	3	7

2

가 30

< 3> KEGG

	Microarray	Document	Microarray + GI + Document
Top 30	0	0	3
Top 50	0	1	4
Top 100	0	3	5

3

30 , 50 , 100

Score_{PCC_i} Score_{diff_i}

GeneCards

가

. 4 , 10 , 30 ,
 50 , 100

< 4>

Score_{PCC_i} Score_{diff_i}

	Score _{PCC_i}	Score _{diff_i}
Top 10	9	6
Top 30	26	23
Top 50	40	35
Top 100	70	59

4.

가

(preprocess)

- [1] Duggan. D. J, Bittner. M, Chen Y, Meltzer. P, Trent. J. M,
“Expression profiling using cDNA microarrays,” Nature Genetics Supplement, vol. 21, pp.10-14, 1999.
- [2] Eunji Shin, Yongmi Yoon, Jaegyeon Ahn, Sanghyun Park, “TC-VGC: A Tumor Classification System using variations in Genes’ Correlation”, Comput. Methods Programs in Biomed, vol. 104, pp.87-101, 2011
- [3] Liu X, Liu ZP, Zhao XM, et al “Identifying disease genes and module biomarkers by differential interactions” J Am Med Inform Assoc, vol. 19, pp.241-248, 2012
- [4] Chatr-Aryamontri A, Breitkreutz BJ, Heinicke S, Boucher L, Winter A, Stark C, Nixon J, Ramage L, Kolas N, O'Donnell L, Reguly T, Breitkreutz A, Sellam A, Chen D, Chang C, Rust JM, Livstone MS, Oughtred R, Dolinski K, Tyers M. “The BioGRID Interaction Database: 2013 update” Nucleic Acids Res. 2012 Nov 30
- [5] Edgar R, Domrachev M, Lash AE. “Gene Expression Omnibus: NCBI gene expression and hybridization array data repository”, Nucleic Acids Res. 2002 Jan 1;30(1):207-10
- [6] "Mori M, Mimori K, Yokobori T et al., 2011, “Gene expression profiles in 132 laser microdissected colorectal cancer tissues” (data accessible at NCBI GEO database (Edgar et al., 2002), accession GSE21815)."
- [7] HUGO Gene Nomenclature Committee at the European Bioinformatics Institute, http://www.genenames.org/cgi-bin/hgnc_downloads/
- [8] GeneCards <http://www.genecards.org/>
- [9] Belinky, F, Bahir, I, Stelzer, G, Zimmerman, S, Rosen, N, Nativ, N, Dalah, I, Iny Stein, T, Rappaport, N, Mituyama, M, Safran, M and Lancet, D. “Non-redundant compendium of human ncRNA genes in GeneCards”, Bioinformatics 29, 2: 255-61 (2013)
- [10] PubMed, <http://www.ncbi.nlm.nih.gov/pubmed>
- [11] KEGG, <http://www.kegg.jp/kegg/kegg2.html>
- [12] National Cancer Institute, “Genetics of Colorectal Cancer”, <http://www.cancer.gov/cancertopics/pdq/genetics/colorectal/HealthProfessional/page2>, 2013. 7. 25
- [13] Stanford Log-linear Part-Of-Speech Tagger, <http://nlp.stanford.edu/>