

Generation of simple structured IR functions by genetic algorithm without stagnation

Kulunchakov A. S.^a, Strijov V. V.^b

^a*Moscow Institute of Physics and Technology*

^b*Computing Centre of the Russian Academy of Sciences*

Abstract

This paper investigates an approach to construct new ranking models for Information Retrieval. The IR ranking model depends on the document description. It includes the term frequency and document frequency. The model ranks documents upon a user request. The quality of the model is defined by the difference between the documents, which experts assess as relative to the request, and the ranked ones. To boost the model quality a modified genetic algorithm was developed. It generates models as superpositions of primitive functions and selects the best according to the quality criterion. The main impact of the research is the new technique to avoid stagnation and to control structural complexity of the consequently generated models. To solve problems of stagnation and complexity, a new criterion of model selection was introduced. It uses structural metric and penalty functions, which are defined in space of generated superpositions. To show that the newly discovered models outperform the other state-of-the-art IR scoring models the authors perform a computational experiment on TREC datasets. It shows that the resulted algorithm is significantly faster than the exhaustive one. It constructs better ranking models according to the MAP criterion. The obtained models are much simpler than the models, which were constructed with alternative approaches. The proposed technique is significant for developing the information retrieval systems based on expert assessments of the query-document relevance.

Keywords: information retrieval, evolutionary stagnation, ranking function, genetic programming, overfitting

Email addresses: kulu-andrej@yandex.com (Kulunchakov A. S.), strijov@gmail.com (Strijov V. V.)

1. Introduction

Information retrieval is finding relevant documents, which satisfy an information need, from within large collections [1]. To retrieve documents relevant to a query, one needs a rank estimation procedure called *ranking model*. It is defined on pairs *document-query*. For each pair it returns relevance of the *document* to the *query*. [2] defines IR ranking models as functions of two basic features of these pairs: term frequency (*tf*) and document frequency (*idf*). In this paper ranking models are constructed considered as mathematical functions defined on tf-idf features. Instead of enlarging the set of features to provide better performance [3], current paper use the same tf-idf features to make further comparison consistent.

The ranking models in [4, 5, 6, 7, 8] are derived on some theoretical assumptions. These assumptions This allow to build ranking models without an IR collection, but these assumptions are not often met. For example, the derived ranking models are not optimal according to mean average precision [1] on TREC collections [2]. Moreover, the quality of these models significantly differs on various the collections [2].

High-performing ranking models are also discovered by automatic procedures. The paper [2] exhaustively explores a set of IR ranking models represented as superpositions of expert-given grammar elements. The grammar is an expert-given set of primitive mathematical functions, where variables are *tf-idf* features [9]. The exhaustive algorithm explores the set of superpositions, which consists of at most 8 grammar elements. The best explored ranking functions in [2] are better in average on TREC collections than ones in [4, 5, 6, 7, 8]. Moreover, these functions are guaranteed to have simple structure. However, this algorithm has high computational complexity [2]. Therefore, an exploration of more complex superpositions is an intractable problem.

Another approaches to improve IR expert systems include various genetic algorithms: search for an optimal document indexing [10, 11], clustering documents according to their relevance to queries [12, 13], tuning parameters of queries [14, 15], facilitate automatic topic selections [16], search for key words in documents [17] and optimal coefficients of a linear superposition of ranking models [18, 19]. Genetic algorithms are applied to select features in image retrieval and classification [20]. Genetic algorithms are used to generate ranking functions represented as superpositions of grammar elements [21, 22, 23]. These procedures significantly extend the set of ranking superpositions considered in [2]. However, the *basic* algorithms in [21, 22] produce superpositions with significant structural complexity after 30-40 iterations of mutations and crossovers [23]. The basic

32 algorithms do not control the structural complexity of generated superpositions and do not solve a
 33 problem of evolutionary stagnation, when a population stops to change.

Strengths	Weaknesses
[Fan et al. 2004, 2000], [22, 21]	
Large feasible set of ranking functions Fast convergence to a local optimum	Complicated final superpositions Does not provide global optimum in the feasible set of functions Have not been tested on different datasets to show uniform improvement on them
[Goswami 2014], [2]	
Provides global optimum with respect to the feasible set Compact final ranking functions Have been tested on different datasets and uniform improvement over existing approaches was shown	Small feasible set of ranking functions
[BM25]	
Theoretically justified Simple and compact explicit expression	Is not uniformly good over different datasets
[The proposed model generation algorithm]	
Large feasible set of ranking functions Fast convergence to a local optimum Compact final ranking functions Have been tested on different datasets to show uniform improvement on them	Does not provide global optimum in the feasible set of functions

Table 1: Comparison of CPU time required by structural metrics

34 The problem of evolutionary stagnation appears when a majority of stored superpositions have
 35 similar structure and high quality. Next crossover operations constructs superpositions, which are
 36 similar to the stored ones. The mutation operation constructs a superposition, which is unlikely

37 to have as high quality as the stored superpositions. This superposition highly probably will be
 38 eliminated. Therefore the population will pass to the next iteration without changes. The genetic
 39 algorithm stops actual generation.

40 To outperform the ranking functions found in [2], one needs to extend the set of superposi-
 41 tions considered there. To perform it, a modified genetic algorithm is proposed. First, it detects
 42 evolutionary stagnation and replaces the worst stored superpositions with random ones. This de-
 43 tection is implemented with a structural metric on superpositions. Regularizers solve the problem
 44 of overfitting. They penalize the excessive structural complexity of superpositions. The paper an-
 45 alyzes various pairs regularizer-metric and chooses the pair providing a selection of better ranking
 46 superpositions. All strengths and weakness of compared approaches are summarized in Table 1.

47 The paper [2] uses TREC collections to test ranking functions. To make the comparison
 48 of approaches consistent, the present paper also use these collections. The collection TREC-7
 49 (trec.nist.gov) is used as the train dataset to evaluate quality of generated superpositions. The
 50 collections TREC-5, TREC-6, TREC-8 are used as test datasets to test selected superpositions.

51 2. Problem statement

There given a collection C consisting of documents $\{d_i\}_{i=1}^{|C|}$ and queries $Q = \{q_j\}_{j=1}^{|Q|}$. For each
 query $q \in Q$ some documents C_q from C are ranked by experts. These ranks g are binary

$$g : Q \times C_q \rightarrow \mathbb{Y} = \{0, 1\},$$

52 where 1 corresponds to relevant documents and 0 to irrelevant.

To approximate g , superpositions of grammar elements are generated. The grammar \mathfrak{G} is a
 set $\{g_1, \dots, g_m, x_w^d, y_w\}$, where each g_i stands for an mathematical function and x_w^d, y_w stand for
 variables. These variables are tf-idf features of *document-query* pair (d, q) . Feature x_w^d is a frequency
 of the word $w \in q$ in d , feature y_w is a frequency of w in C :

$$x_w^d = t_d^w \log \left(1 + \frac{l_a}{l_d} \right), \quad y_w = \frac{N_w}{|C|}, \quad (1)$$

53 where N_w is the number of documents from C containing w , t_d^w is the frequency of w in d , l_d is the
 54 number of words in d (the size of a document d), l_a is an average size of documents in C . Each
 55 superposition f of grammar elements is stored as a directed labeled tree T_f with vertices labeled
 56 by elements from \mathfrak{G} . The set of these superpositions is defined as \mathfrak{F} .

The value of f on a pair (d, q) is defined as a sum of its values on (d, w) , where w is a word from q :

$$f(d, q) = \sum_{w \in q} f(x_w^d, y_w).$$

The superposition f ranks the documents for each q . The quality of f is the mean average precision [1]

$$\text{MAP}(f, C, Q) = \frac{1}{|Q|} \sum_{q=1}^Q \text{AveP}(f, q),$$

where

$$\text{AveP}(f, q) = \frac{\sum_{k=1}^{|C_q|} (\text{Prec}(k) \times g(k))}{\sum_{k=1}^{|C_q|} \text{Rel}(k)}, \quad \text{Prec}(k) = \frac{\sum_{s=1}^k g(s)}{k},$$

57 where $g(k) \in \{0, 1\}$ is a relevance of the k -th document from C .

58 This paper aims at finding the superposition f , which maximizes the following quality function

$$f^* = \operatorname{argmax}_{f \in \mathfrak{F}} \mathcal{S}(f, C, Q), \quad \mathcal{S}(f, C, Q) = \text{MAP}(f, C, Q) - \text{R}(f), \quad (2)$$

59 where R is a regularizer controlling the structural complexity of f .

60 The exhaustive algorithm in [2] generates random ranking superpositions consisting at most of
61 8 elements of the grammar \mathfrak{G} . Let \mathfrak{F}_0 be the set of the best superpositions selected in [2]. The
62 solution f^* is compared with the superpositions from \mathfrak{F}_0 with respect to to MAP.

63 3. Generation of superpositions

IR ranking functions are superpositions of expert-given primitive functions. These superpositions are generated by the genetic algorithm. It uses an expertly given grammar \mathfrak{G} and constructs superpositions of its elements. On each iteration it **keeps stores** a population of the best selected superpositions. To update them and pass to the next iteration, it generates new superpositions with use of the stored ones. Since the superpositions are represented as trees, the algorithm applies crossover $c(f, h)$ and mutation $m(f)$ operations to the stored trees

$$c(f, h) : \mathfrak{F} \times \mathfrak{F} \rightarrow \mathfrak{F}, \quad m(f) : \mathfrak{F} \rightarrow \mathfrak{F},$$

64 **Definition 1.** Crossover operation $c(f, h) : \mathfrak{F} \times \mathfrak{F} \rightarrow \mathfrak{F}$ produces a new superpositions from given f
65 and h . This operation represents f and h as trees, uniformly selected a subtree for each of them
66 and swaps these subtrees.

Algorithm 1 Basic genetic algorithm

Require: grammar \mathfrak{G} , required value α of MAP

Ensure: superposition f of elements from G with $\text{MAP} \leq \alpha$;

create a set of initial, random superpositions \mathfrak{M}_0 ,

repeat

- crossover random pairs of stored superpositions \mathfrak{M} ,
- mutate random superpositions from the population \mathfrak{M} ,
- consider these generated superpositions and the ones stored in \mathfrak{M} . Select the best of them according to MAP,
- store the best generated superpositions in the population \mathfrak{M} and pass it to the next iteration,

until the required value of MAP is reached;

Here is an example of crossover on two superpositions, where randomly selected subtrees are in bold.

$$f(x, y) = \exp(x) + \mathbf{\ln(\mathbf{xy})}, \quad h(x, y) = \sqrt{x} + (\mathbf{x+y})$$

67

↓

$$f'(x, y) = \exp(x) + (\mathbf{x+y}), \quad h'(x, y) = \sqrt{x} + \mathbf{\ln(\mathbf{x \cdot y})},$$

68 The new superpositions are formed by swapping of these subtrees.

69 **Definition 2.** *Mutation $m(f)$ uniformly selects a subtree from f and replace it with another random*
70 *superposition. Mutation produces one new superposition.*

71 Here is an example of mutation on a superposition

$$f(x, y) = \sqrt{x} + \mathbf{\ln(\mathbf{x \cdot y})} \rightarrow f'(x, y) = \sqrt{x} + \mathbf{\exp(\mathbf{y})}.$$

72 **Definition 3.** *Size $|T|$ of a tree T is the number of its vertices.*

73 Restrict the size of substituting tree. If mutation replaces a subtree T with a tree T' , then bound
74 the size of T' by $c|T|$, where c is a constant. This restriction allows us to explore the set \mathfrak{F} more

75 gradually. The reason is to prevent the algorithm from instantaneous moving toward complicated
76 superpositions if the stored population consists mainly of simple structured superpositions. Now
77 the genetic algorithm is described in Algorithm 1. It will be referred as *basic genetic algorithm*.

78 4. Metric properties of basic genetic algorithm

To analyze the genetic algorithm, introduce a structural metric $\mu(T, T')$. It is defined on pairs of directed labeled trees. Therefore, it is defined on pairs of elements from \mathfrak{F} as well.

$$\mu(f, f') = \mu(T_f, T_{f'}).$$

79 This structural metric satisfies the following conditions

- 80 1) $\mu(f, f) = 0$, $\mu(f, f') > 0$ if $f \neq f'$ (non-negativity),
- 81 2) $\mu(f, f') = \mu(f', f)$ (symmetry),
- 82 3) $\mu(f, f') \leq \mu(f, f'') + \mu(f'', f')$ (triangle inequality). enumerate

For $r > 0$ define the r -neighborhood $U_r(f)$ of superposition f as a set of superpositions in \mathfrak{F} that are at distance less than r from f

$$U_r(f) = \{f' \in \mathfrak{F} : \mu(f, f') < r\}.$$

To associate the structural distance between superpositions with a distance on their values, introduce an extra condition. Claim that the functions, lying in one structural neighborhood, should rank the documents mainly similarly. Define a distance function η on the ranks of IR ranking functions:

$$\eta(f, f') = \frac{1}{|C|(|C| - 1)} \sum_{d_j, d_k \in C} [f(d_j) < f(d_k)][f'(d_j) > f'(d_k)],$$

where $[A]$ is the indicator of event A . It is related with Kendall rank correlation coefficient by the equation:

$$\tau(f, f') = 1 - 2\eta(f, f').$$

83 The function η is the normalized number of inversions necessary to transform one list with
84 ranks to the other. Therefore $\eta(f, f')$ is a distance on the values of the superpositions. Call the
85 neighborhood $V_r(f) = \{f' : \eta(f, f') < r\}$ the value-neighborhood.

86 Introduce a condition for μ to detect evolutionary stagnation of the genetic algorithm enumerate

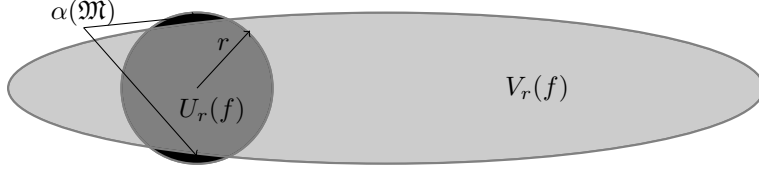


Figure 1: Illustration of supposed relation between $U_r(f)$ and $V_r(f)$.

4)

$$\alpha(\mathfrak{M}) = \nu\left([\mu(f, f') \leq \alpha_1] \Rightarrow [\eta(f, f') \leq \alpha_2] \mid f, f' \in \mathfrak{M}\right) \geq 1 - \varepsilon, \quad (3)$$

87 where $\alpha_1, \alpha_2, \varepsilon$ are some constants and $\nu(A)$ is the frequency of event A .

88 It claims that structurally similar functions rank documents mainly similarly. Figure 1 shows sup-
 89 posed relation between structural neighborhood $U_r(f)$ and value-neighborhood $V_r(f)$. Condition (3)
 90 states that the area of the black region on Figure 1 should be relatively small.

91 Let f_{opt} be a superposition of high quality according to \mathcal{S} . If μ satisfies condition (3), then the
 92 superpositions in the neighborhood $U_r(f_{\text{opt}})$ will also have high quality. Suppose that $f_{\text{opt}} \neq f^*$ (2).
 93 It means that the optimal ranking superposition f^* is not found yet. If all superpositions of a stored
 94 population \mathfrak{M}_i lie in $U_r(f_{\text{opt}})$, then they will rarely leave $U_r(f_{\text{opt}})$ on the next iterations, since
 95 crossovers produce superpositions mainly from $U_r(f_{\text{opt}})$ and mutations produce superpositions
 96 mainly of lower quality. Therefore, the optimal function f^* will frequently become unreachable for
 97 the genetic algorithm, as consequence of this evolutionary stagnation.

98 **Definition 4.** *Evolutionary stagnation is a situation in a genetic algorithm, when stored superpo-
 99 sitions are pairwise similar. The generated algorithm stops generation of principally new superpo-
 100 sitions and the population mainly does not change from iteration to iteration.*

Definition 5. *Radius $r(\mathfrak{M})$ of a population \mathfrak{M} is the minimum size of r -neighborhood with center
 in $f \in \mathfrak{M}$, which accommodates \mathfrak{M} . It shows how are the functions from \mathfrak{M} scattered across the
 set \mathfrak{F} .*

$$r(\mathfrak{M}) = \operatorname{argmin}_{r>0} \{\exists f \in \mathfrak{F} \forall f' \in \mathfrak{M} : f' \in U_r(f)\} = \min_{f \in \mathfrak{M}} \max_{f' \in \mathfrak{M}} \{\mu f', f\}. \quad (4)$$

101 Detect evolutionary stagnation with structural metric μ . Lets consider a population \mathfrak{M} stored by
 102 the genetic algorithm. If the genetic algorithm stagnates, then $r(\mathfrak{M})$ is relatively small. Oppositely,
 103 if the population is diverse, then the $r(\mathfrak{M})$ is big. Therefore evolutionary stagnation could be
 104 detected with the radius $r(\mathfrak{M})$. However, it is an intractable problem to find the exact value
 105 of $r(\mathfrak{M})$. Therefore, propose an empirical estimation of this radius.

106 **Definition 6.** *Structural complexity $|f|$ of superposition f is the number of grammar elements,*
 107 *which f consists of.*

Definition 7. *Empirical radius $r_e(\mathfrak{M})$ of is a normalized average distance between superpositions*
in \mathfrak{M} .

$$r_e(\mathfrak{M}) = \frac{\sum_{f, f' \in \mathfrak{M}} \mu_i(f, f')}{|\mathfrak{M}| \sum_{f \in \mathfrak{M}} |f_j|}. \quad (5)$$

108 This estimation is used to detect evolutionary stagnation of the genetic algorithm. If $r_e(\mathfrak{M})$ is
 109 less than a threshold $r(\mathfrak{M}) < \text{Thresh}$, eliminate the worst superpositions from \mathfrak{M} and replace them
 110 with random superpositions of the same structural complexity. This procedure increases the radius
 111 of \mathfrak{M} and diversifies it. Therefore, the present aim of this paper is to select a proper structural
 112 metric μ , which satisfies all mentioned conditions.

113 5. Structural metrics

114 Each ranking superposition $f \in \mathfrak{F}$ is represented as directed tree T_f , which vertices are labeled
 115 by elements from grammar \mathfrak{G} . Structural metrics are defined on pairs of such trees. It automatically
 116 defines them on pairs of superpositions. This paper analyzes three metrics.

117 5.1. Similarity according to an isomorphism

118 The first structural metric μ_1 uses a definition of common subgraph of two graphs [24].

119 **Definition 8.** *Two graphs G_1 and G_2 are called isomorphic if there is an edge-preserving bijection*
 120 *between their vertex sets. The edge-preserving property states that two vertices are adjacent iff their*
 121 *images are adjacent.*

122 **Definition 9.** *Two trees T_i, T_j have a common subtree T if each of them has a subtree isomorphic*
123 *to T .*

124 **Definition 10.** *A size $|T|$ of a tree T is the number of its vertices.*

125 **Definition 11.** *The largest common subtree T_{ij} of two directed labeled trees T_i and T_j is the tree*
126 *of the largest size among all common subtrees of T_i and T_j .*

The distance between T_i and T_j is calculated by the following formula

$$\mu_1(T_i, T_j) = |T_i| + |T_j| - 2|T_{ij}|.$$

127 The paper [24] defines μ_1 likewise on pairs of graphs and proves that μ_1 satisfies 1-3 conditions if
128 the graph size is defined as the number of its edges. For a tree the number of its vertices is equal
129 to the number of its edges plus 1. Therefore, the results mentioned in [24] are applicable for our
130 case and μ_1 satisfies 1-3 conditions. The last 4th condition is checked empirically.

131 5.2. Similarity according to edit distance

132 As before, a superposition is represented by a directed labeled tree. Represent a tree as a string
133 of characters. This string is constructed as a sequence of labels of vertices written in pre-order [25].

134 Now define a structural metric μ_2 on pairs of character strings. It automatically defines the
135 structural metric on pairs of superpositions. As the arities of functions from \mathfrak{G} are known, each
136 superposition could be reconstructed from its string representation. Therefore, there is no two
137 character strings corresponding to one superposition of primitive functions. The structural metric μ_2
138 is called a Levenshtein distance.

139 **Definition 12.** *The Levenshtein distance between two character strings is the minimum number of*
140 *single-character edits (insertions, deletions and rewritings) required to change one string into the*
141 *other.*

142 Each edit distance satisfies the conditions 1-3. The metric μ_2 also satisfies them in the case
143 when it is defined on pairs of superpositions, because the string representation is bijective. The last
144 4th condition is checked empirically.

145 The third structural metric μ_3 is a Levenshtein distance defined on pairs of directed labeled
146 trees.

147 **Definition 13.** *The Levenshtein distance between two trees is the minimum number of edits (edge*
148 *insertions, edge deletions and vertex relabeling) required to change one tree into the other.*

149 The structural metric μ_3 satisfies the metric axioms [26]. The last 4th condition is checked
150 empirically.

151 6. Regularizers

152 To approximate noisy data accurately, the genetic algorithm generate complex superpositions
153 after some iterations. To prevent this overfitting, it should control the structural complexity of
154 superpositions by a regularizer. The regularizer restricts a set $\mathfrak{F}' \subset \mathfrak{F}$ of superpositions reachable
155 by the genetic algorithm. Search for a regularizer, which makes the set \mathfrak{F}' sufficiently rich to
156 find there a proper approximating superposition and sufficiently small to avoid overfitting of the
157 algorithm. Lets consider the structural parameters of directed labeled trees

- 158 1) The size of a tree, see Definition 3.
- 159 2) The number of leaves in a tree.
- 160 3) The height of a tree.

161 A restriction of these parameters makes complex superpositions unreachable for the genetic algo-
162 rithm. This paper analyzes three regularizers built on these structural parameters. To penalize
163 accurate superpositions less, all of these regularizers are proportional to MAP.

164 1) $R_1(f) = p \cdot \text{MAP}(f) \cdot I(|f| < \text{CT})$,

165 where CT is a threshold for the structural complexity, p is a penalty parameter. The reg-
166 ularizer R_1 penalizes those superpositions, which have structural complexity larger than the
167 threshold CT.

168 2) $R_2(f) = p \cdot \text{MAP}(f) \cdot I(|f| \geq \text{CT}) \cdot (|f| - \text{CT})$,

169 where C is a positive parameter. The regularizer R_2 penalizes the superpositions having struc-
170 tural complexity larger than the threshold CT . And the more complex a superposition, the
171 higher the penalty.

172 3) $R_3(f) = p \cdot \text{MAP}(f) \cdot |f|^* \cdot \log(|f| + 1)$,

173 The regularizer R_3 treats a structural complexity of a superposition as the number of leaves $|f|^*$
174 of its tree multiplied by the estimation $\log(|f| + 1)$ of its height.

175 All parameters from the definitions should be set empirically. To set them one needs to follow the
176 principle mentioned above: the set \mathfrak{F}' should be sufficiently rich to find there a proper approximating
177 superposition and sufficiently small to avoid overfitting of the genetic algorithm.

178 Select proper structural metric and regularizer to modify the basic genetic algorithm. The
179 modified version solves the problems of overfitting and evolutionary stagnation. This version is
180 described in Algorithm 2.

181 7. Computational experiment

182 The main goal of this paper is to generate superpositions outperforming the ones from \mathfrak{F}_0
183 selected in [2]. These functions, in turn, outperform known ranking models BM25, LGD, LM_{DIR} .
184 Therefore, if the modified genetic algorithm succeeds in outperforming functions from \mathfrak{F}_0 , it will
185 also outperform BM25, LGD, LM_{DIR} as well. Now describe the data used to estimate the quality
186 of the generated superpositions.

187 *Data.* Authors in [2] estimate the quality ranking functions on TRECs. To make the comparison
188 with \mathcal{F}_0 consistent, use TRECS as well. Perform the computational experiment on Trec-5, Trec-
189 6, Trec-7, Trec-8 (trec.nist.gov). The Text REtrieval Conference (TREC), co-sponsored by the
190 National Institute of Standards and Technology (NIST) and U.S. Department of Defense, was
191 started in 1992 as part of the TIPSTER Text program. For each TREC, National Institute of
192 Standards and Technology (NIST) provides a test set of documents and questions. Participants
193 run their own retrieval systems on the data, and return to NIST a list of the retrieved top-ranked

Algorithm 2 Modified genetic algorithm

Require: grammar \mathcal{G} , required value α of MAP

Ensure: superposition f of elements from G with $\text{MAP} \leq \alpha$;

create a set of initial, random superpositions \mathfrak{M}_0 ,

repeat

- crossover random pairs of stored superpositions \mathfrak{M} ,
- mutate random superpositions from the population \mathfrak{M} ,
- consider these generated superpositions and the ones stored in \mathfrak{M} . Select the best of them according to the quality function \mathcal{S} (2),
- store the best superpositions in a population \mathfrak{M}' and pass it to the next iteration,
- **if** $d_e(\mathfrak{M}') < \text{Thresh}$ **then**
 - evolutionary stagnation is detected and we replace the worst superpositions from the population \mathfrak{M}' by random superpositions,
- **end if**
- $\mathfrak{M} = \mathfrak{M}'$.

until the required value of MAP is reached;

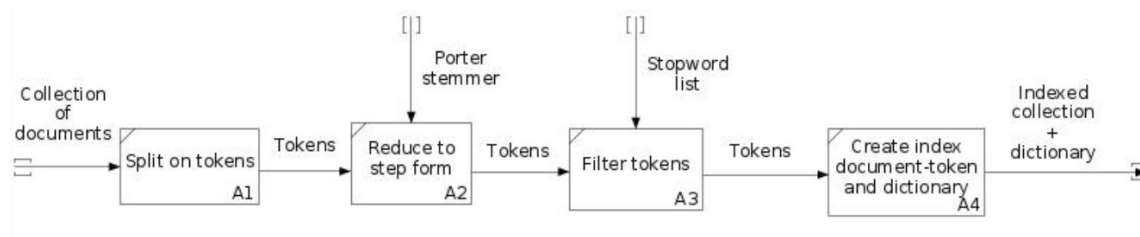


Figure 2: Scheme of data preprocessing steps.

194 documents. NIST pools the individual results, judges the retrieved documents for correctness, and
 195 evaluates the results. Thus each TREC consists of a collection of documents, user queries and
 196 judgments for a subset of a collection Each TREC is associated with this triplet. Each triplet has
 197 a collection of nearly 500 000 documents. 50 queries to the collection and 2000 judgments for each
 198 query in average. The number specified after the name $Trec_{i,j}$ denotes the year of the creation of
 199 the TREC.

200 7.1. Data processing

201 As TREC collections are large, calculations of the variables x_w^d and y_w (1) are computationally
 202 expensive. To speed up the calculations, one should perform data preprocessing. Terrier IR Plat-
 203 form v3.6 (terrier.org) perform necessary steps for this preprocessing. It provides flexible processing
 204 of terms through a pipeline of components (stopwords removing, stemmers, etc.). The platform
 205 indexes a collection of documents. The preprocessing steps include stemming using Porter stem-
 206 mer and removing stop-words using the stopword list. Second, Terrier performs a query expansion
 207 techniques and retrieves required documents efficiently. It processes the data stored in Trec5-8
 208 and returns the matrices of features x_w^d and y_w for each word $w \in q$ and each document from the
 209 collection having this word.

210 The algorithm of primary data preprocessing makes the following steps, see Figure 2.

- 211 1. Split documents on tokens. Reduce each token to its stem form by Porter stemmer [4].
- 212 2. Filter the set of stemmed tokens is according to the stopwords list.

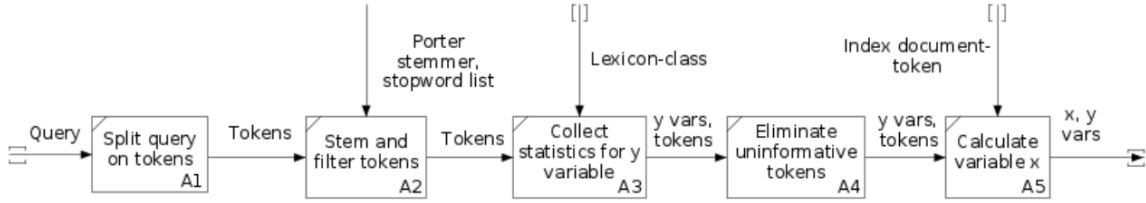


Figure 3: Scheme of query processing steps.

- 213 3. The collection is represented as an index *document-token*.
 214 4. Create a *lexicon-class*, which represents the list of terms (dictionary) in the index.

215 After the preliminary steps are performed, one can calculate the variables x_w^d and y_w for each
 216 query q , see Figure 3.

- 217 1. Split q on tokens. Process each token by the stemmer and filter the resulted set by the
 218 stopword list.
 219 2. *Lexicon-class* collects statistics about the tokens. It calculates the feature y_w .
 220 3. Eliminate tokens with high value of y_w as uninformative.
 221 4. For each token the platform retrieves the information about its second feature x_w^d from the
 222 index.

223 The described scheme is used by the modified genetic algorithm to estimate the quality of a
 224 superposition. Now describe the system performing this modified genetic algorithm. This system
 225 generates superpositions of primitive functions.

226 *7.2. Generation system*

227 Algorithm 2 gives the description of the modified genetic algorithm used for generation of ranking
 228 superpositions. These superpositions are constructed from the elements of $\mathfrak{G} = \{x_w^d, y_w, +, -, \times, \div, \log, \exp, \sqrt{\cdot}\}$.

229 On each iteration the algorithm stores 20 best generated superpositions. To create new super-
230 positions, it performs 10 crossovers and 10 mutations on the stored ones. Then it selects 20
231 best according to (2) and pass to the next iteration. This paper terminates the generation af-
232 ter 300 iterations. The selected superpositions are compared with the ones from \mathfrak{F}_0 To use
233 this algorithm, one must select proper regularizer and structural metric. The code for this sys-
234 tem is found in [https://github.com/KuluAndrej/Generation-of-simple-structured-IR-functions-by-](https://github.com/KuluAndrej/Generation-of-simple-structured-IR-functions-by-genetic-algorithm-without-stagnation)
235 [genetic-algorithm-without-stagnation](https://github.com/KuluAndrej/Generation-of-simple-structured-IR-functions-by-genetic-algorithm-without-stagnation).

236 7.3. Selection of regularizer and structural metric

237 This paper analyzes three metrics and three regularizers defined above with respect to the genetic
238 algorithm. There are 9 combinations of these metrics and regularizers. Selects the pair, which
239 provides better generation of superpositions both in terms of structural diversity and prediction
240 accuracy. The selected pair is used by the modified genetic algorithm to generate an optimal ranking
241 superposition.

242 Table 2 shows a computational efficiency of calculation of different metrics with respect to
243 different regularizers. There are 9 possible pairs metric-regularizer. The modified genetic algorithm
244 is launched 100 times for each pair. The CPU time required to calculate all values of a metric
245 is averaged over these 100 launches and 300 iterations for each launch. Table 2 shows that μ_2
246 is uniformly easiest to calculate. At the same time, μ_1 is uniformly hardest to calculate. This
247 efficiency is considered in the selection. Now analyze the pairs with respect to the generation of
248 superpositions.

249 First, analyze the modified genetic algorithm without regularizers. All measured values are aver-
250 aged over 100 launches, see Figure 4. On the last 300-th iteration the average structural complexity
251 of superpositions in the population is more than 40. Figure 4 shows slow trend to evolutionary stag-
252 nation. The reason is that structural complexity of the generated superpositions grows dramatically
253 with the iteration number. It makes the stored superpositions sufficiently diverse. Therefore during
254 the whole evolution the empirical diameter d_e of the stored population is large. However, the gener-
255 ated superpositions are significantly overfitted and should be penalized for the excessive structural

256 complexity.

Regularizer	μ_1	μ_2	μ_3
R ₁	11.52	1.84	4.54
R ₂	6.7876	0.9347	1.5666
R ₃	7.63	1.05	1.87

Table 2: Comparison of CPU time required by structural metrics

257 Now let us analyze 3 metrics with presence of a regularizer. For each pair metric-regularizer
258 plot the empirical diameter d_e depending on the number of iteration. Figures 5, 6, 7 also shows the
259 average structural complexity l_a of stored superpositions. It allows to make inferences about the
260 presence of overfitting.

261 Note that the empirical diameter $d(\mathfrak{M})$ calculated with μ_1 remains approximately unchanged
262 during the whole evolution, see Figures 5, 6, 7. This particular feature does not allow to detect
263 evolutionary stagnation in proper time. The actual start of evolutionary stagnation can not be
264 denoted with μ_1 . Moreover, calculation of μ_1 is computationally inefficient comparing with μ_2
265 and μ_3 , see Table 2. These reasons lead to elimination of μ_1 from the further analysis.

266 Two other metrics μ_2 and μ_3 provide almost equal values of $d(\mathfrak{M})$, see Figures 5, 6, 7. The
267 relative difference in these values is under 5% for all variants of used regularizer. Therefore, without
268 loss of generality, select the structural metric μ_2 as more efficiently calculated, see Table 2.

269 The first regularizer R₁ is too strict, see Figure 5. The algorithm falls into evolutionary stag-
270 nation on the first iterations, because the set of reachable superpositions \mathfrak{F}' is small. The similar
271 situation is observed for the second regularizer R₂, see Figure 6. The algorithm does not immediately
272 fall into evolutionary stagnation. The stored superpositions are updated up to the 300-th iteration.
273 However, the empirical diameter $d(\mathfrak{M})$ significantly decreases after 30-40 iterations, see Figure 6.
274 It means that although the stored superpositions are being updated throughout the evolution, they
275 have mainly similar structures. These reasons lead us to the use of the third regularizer R₃. The
276 value of the empirical diameter $d(\mathfrak{M})$ decreases smoothly with R₃, see Figure 7. It allows to have
277 enough iterations to learn the structure of optimal superposition and detect evolutionary stagna-

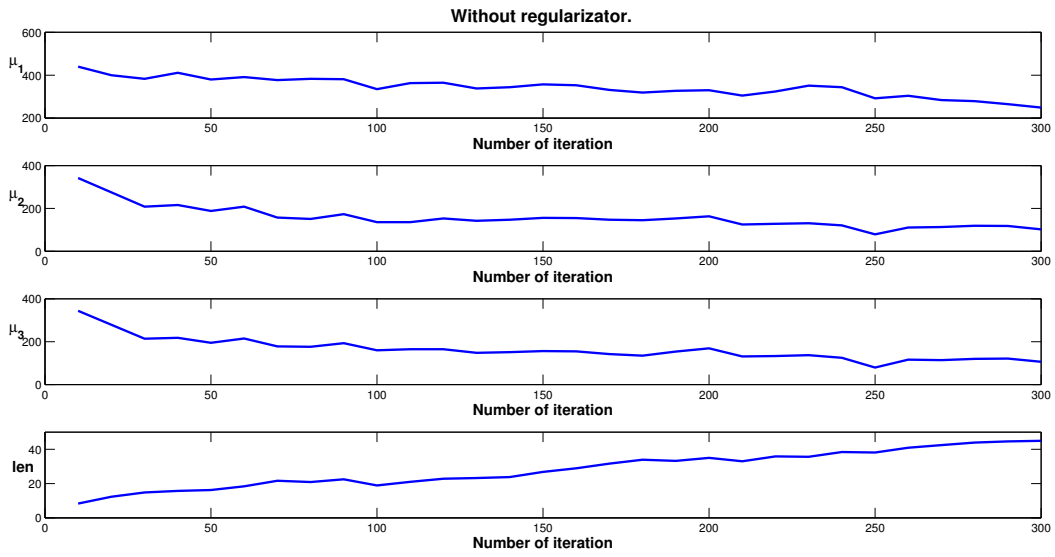


Figure 4: Dynamics of $d(\mathfrak{M})$ and l_a when no regularizer is used.

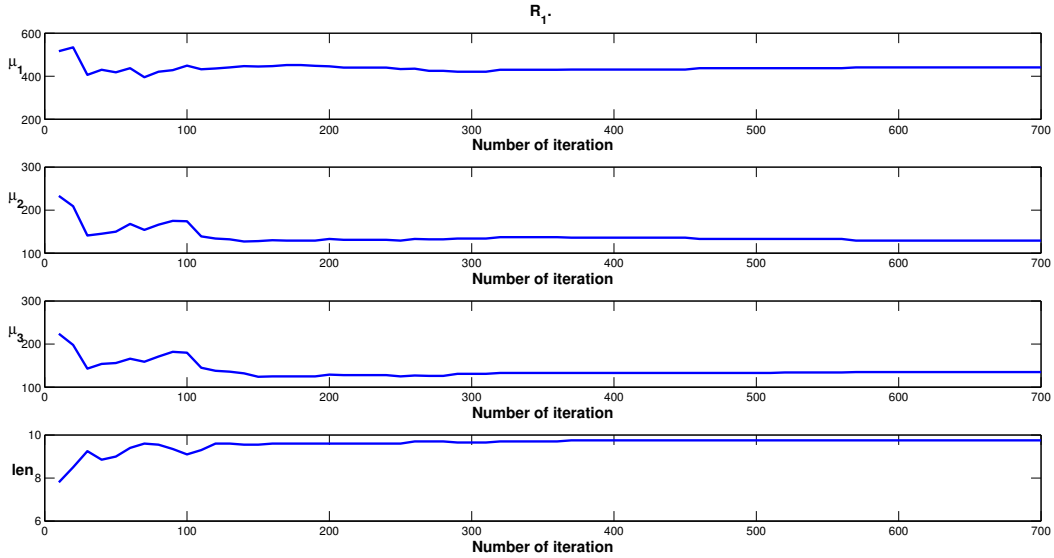


Figure 5: Dynamics of $d(\mathfrak{M})$ and l_a when the regularizer R_1 is used.

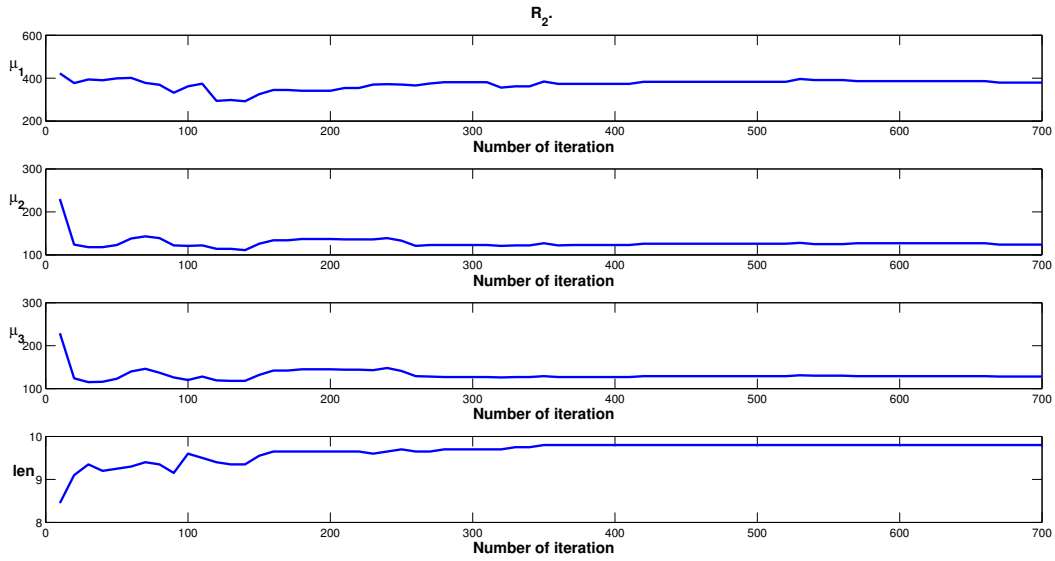


Figure 6: Dynamics of $d(\mathfrak{M})$ and l_a when the regularizer R_2 is used.

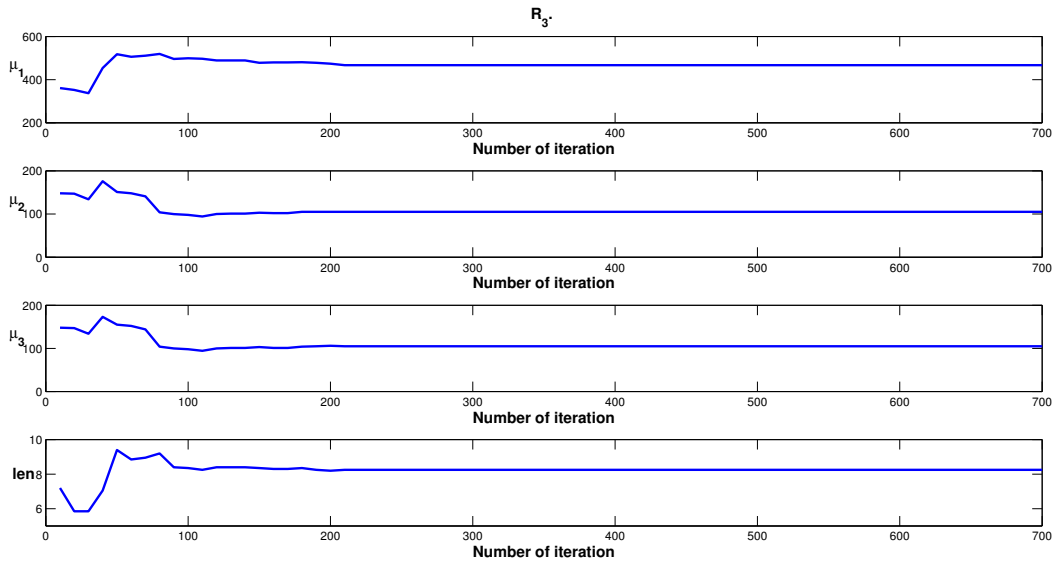


Figure 7: Dynamics of $d(\mathfrak{M})$ and l_a when the regularizer R_3 is used.

278 tion. Since the structural metric μ_2 and the regularizer R_3 are selected, the modification of the
 279 genetic algorithm is ready to generate ranking superpositions. s

280 *Generation of ranking superpositions.* Modified genetic algorithm is launched on TREC-7. The best
 281 selected superpositions are compared with ones from \mathfrak{F}_0 . The superpositions in \mathfrak{F}_0 are of simple
 282 structure and have a high quality in average on analyzed collections. Besides, these superpositions
 283 are better in average than the traditionally used ranking models BM25, LGD, LM_{DIR} . Here is the
 284 list of the best superpositions from \mathfrak{F}_0

285

2

- 286 1. $f_1 = e \sqrt{\ln\left(\frac{x}{y}\right)},$
- 287 2. $f_2 = \sqrt{\frac{\ln(x)}{\sqrt{y}}},$
- 288 3. $f_3 = \sqrt[4]{\frac{x}{y}},$
- 289 4. $f_4 = \sqrt{y + \sqrt{\frac{x}{y}}},$
- 290 5. $f_5 = \sqrt[4]{\frac{x}{y}} \cdot e^{-y/2},$
- 291 6. $f_6 = \sqrt{\sqrt{x} + \sqrt{\frac{x}{y}}}.$

292 The selection of the best superpositions is performed by the modified genetic algorithm on
 293 TREC-7. The other datasets TREC-5, TREC-6, TREC-8 serve as test datasets. After 1000
 294 iterations the modified genetic algorithm selects the following family of superpositions (for the
 295 convenience denote $\ln(x+1)$ as $\ln(x)$ and $g(x) = \ln \ln(x)$):

296

2

- 297 1. $h_1 = g\left(\frac{g(x)}{\sqrt{\ln(x) + x}}\right) - \ln(y),$
- 298 2. $h_2 = g\left(\frac{g(x)}{\sqrt{\frac{1}{2}\ln(x) + x}}\right) - \ln(y),$

$$\begin{aligned}
299 \quad & 3. \ h_3 = g \left(\ln \left(\frac{g(x)}{\sqrt{\frac{1}{2} \ln(x) + x}} \right) - \ln(y) \right), \\
300 \quad & 4. \ h_4 = g \left(\frac{g(x)}{\sqrt{g(\sqrt{x}) + x}} \right) - \ln(y), \\
301 \quad & 5. \ h_5 = g \left(\frac{g(x)}{\sqrt{\ln(x) + \ln(y)}} \right) - \ln(y), \\
302 \quad & 6. \ h_6 = g \left(\frac{g(\ln(x))}{\sqrt{\ln(x) + x}} \right) - \ln(y).
\end{aligned}$$

303 The values of MAP of the superpositions $\{h_j\}$ and $\{f_i\}$ are is presented in Table 3. The
304 superpositions from \mathfrak{F}_0 are presented in the upper half of the table. The superpositions $\{h_j\}$ are
305 presented in the lower half. The qualities of the best functions $\{f_i\}$ are bold in each column in the
306 upper half. In the lower half we bold those values, which are higher than the bold values in the
307 corresponding column in the upper half.

Superposition	TREC-5	TREC-6	TREC-7	TREC-8
Superpositions from \mathfrak{F}_0				
f_1	8.785	13.715	10.038	13.902
f_2	8.518	12.996	9.216	13.074
f_3	8.908	13.615	9.905	13.708
f_4	8.908	13.615	9.905	13.708
f_5	8.908	13.615	9.908	13.709
f_6	8.872	13.613	9.890	13.695
Family of selected superpositions				
h_1	8.965	13.693	10.600	14.403
h_2	9.472	13.723	10.650	14.402
h_3	9.558	13.786	10.631	14.376
h_4	9.226	13.713	10.5	14.374
h_5	8.862	13.388	10.439	14.359
h_6	8.104	13.483	10.421	14.355

Table 3: Comparison of the superpositions $\{h_j\}$ to $\{f_i\}$ according to the MAP criterion

308 Note that the superpositions h_1, h_2, h_3, h_4 are uniformly better than the functions from [2] on
309 all 4 datasets. The other superpositions are better in average. The modified genetic algorithm is
310 able to build effective yet simple structured superpositions, which outperform the known ones.

311 8. Conclusion

312 This paper investigates an Information Retrieval ranking function construction technique. It
313 develops a genetic algorithm, which consequently generates ranking functions. The basic version
314 of this algorithm is inclined to generate overfit ranking functions. To avoid overfitting, one must
315 control their structural complexity and solve an evolutionary stagnation problem. It is solved
316 by use of regularizers and structural metrics respectively. A regularizer, presented in the quality
317 function, controls the structural complexity of the functions. Overfit functions are penalized and
318 unlikely to pass to the next iterations in the genetic algorithm. A structural metric estimates the
319 diversity of the generated functions. If all generated functions are similar to each other, some of
320 them are replaced by random ones. It solves a problem of evolutionary stagnation. This paper
321 analyzes different regularizers and structural metrics and chooses those, which provide a better
322 generation. The modified genetic algorithm uses the selected pair metric-regularizer and generate
323 effective yet simple structured functions. These functions outperform BM25, LGD, LMDIR. and
324 the ones selected by are exhaustive approach.

325 Future research TODO

326 9. Acknowledgements

327 This research is supported by RFBR grant 16-37-00486.

328 **References**

- 329 [1] C. D. Manning, P. Raghavan, H. Schütze, Introduction to Information Retrieval, Cambridge
330 University Press, Cambridge, UK, 2008.
331 URL <http://nlp.stanford.edu/IR-book/information-retrieval-book.html>
- 332 [2] P. Goswami, S. Moura, E. Gaussier, M.-R. Amini, F. Maes, Exploring the space of ir functions,
333 in: ECIR'14, 2014, pp. 372–384.
- 334 [3] Z. Yea, J. X. Huangb, H. Lina, Incorporating rich features to boost information retrieval
335 performance: A svm-regression based re-ranking approach, Expert Systems with Applications
336 38 (6) (2011) 75697574.
- 337 [4] M. F. Porter, Readings in information retrieval, Morgan Kaufmann Publishers Inc., San Fran-
338 cisco, CA, USA, 1997, Ch. An Algorithm for Suffix Stripping, pp. 313–316.
- 339 [5] D. Metzler, W. B. Croft, A markov random field model for term dependencies, in: Proceedings
340 of the 28th Annual International ACM SIGIR Conference on Research and Development in
341 Information Retrieval, SIGIR '05, ACM, New York, NY, USA, 2005, pp. 472–479.
- 342 [6] G. Amati, C. J. Van Rijsbergen, Probabilistic models of information retrieval based on mea-
343 suring the divergence from randomness, ACM Trans. Inf. Syst. 20 (4) (2002) 357–389.
- 344 [7] S. Clinchant, E. Gaussier, Information-based models for ad hoc ir, in: Proceedings of the 33rd
345 International ACM SIGIR Conference on Research and Development in Information Retrieval,
346 SIGIR '10, ACM, New York, NY, USA, 2010, pp. 234–241.
- 347 [8] J. M. Ponte, W. B. Croft, A language modeling approach to information retrieval, in: Proceed-
348 ings of the 21st Annual International ACM SIGIR Conference on Research and Development
349 in Information Retrieval, SIGIR '98, ACM, New York, NY, USA, 1998, pp. 275–281.
- 350 [9] G. Salton, M. J. McGill, Introduction to Modern Information Retrieval, McGraw-Hill, Inc.,
351 New York, NY, USA, 1986.
- 352 [10] M. Gordon, Probabilistic and genetic algorithms in document retrieval, Commun. ACM 31 (10)
353 (1988) 1208–1218.

- 354 [11] H. Valizadegan, R. Jin, R. Zhang, J. Mao, Learning to rank by optimizing ndcg measure, in:
355 Y. Bengio, D. Schuurmans, J. Lafferty, C. Williams, A. Culotta (Eds.), *Advances in Neural*
356 *Information Processing Systems 22*, Curran Associates, Inc., 2009, pp. 1883–1891.
- 357 [12] M. D. Gordon, User-based document clustering by redescribing subject descriptions with a
358 genetic algorithm, *JASIS* 42 (5) (1991) 311–322.
- 359 [13] V. Raghavan, B. Agarwal, Optimal determination of user-oriented clusters: An application for
360 the reproductive plan, in: *Proceedings of the Second International Conference on Genetic Al-*
361 *gorithms on Genetic Algorithms and Their Application*, L. Erlbaum Associates Inc., Hillsdale,
362 NJ, USA, 1987, pp. 241–246.
- 363 [14] J. Yang, R. Korfhage, E. M. Rasmussen, Query improvement in information retrieval using
364 genetic algorithms - A report on the experiments of the TREC project, in: *TREC*, 1992, pp.
365 31–58.
- 366 [15] F. E. Petry, B. P. Buckles, T. Sadasivan, D. H. Kraft, The use of genetic programming to build
367 queries for information retrieval., in: *International Conference on Evolutionary Computation*,
368 1994, pp. 468–473.
- 369 [16] D.-Y. Chiu, Y.-C. Pan, S.-Y. Yu, An automated knowledge structure construction approach:
370 Applying information retrieval and genetic algorithm to journal of expert systems with appli-
371 cations, *Expert Systems With Applications* 36 (5) (2009) 9438–9447. doi:10.1016/j.eswa.
372 2008.12.032.
- 373 [17] H. Chen, Machine learning for information retrieval: Neural networks, symbolic learning, and
374 genetic algorithms, *Journal of the American Society of Information Science* 46 (3) (1995) 194–
375 216.
- 376 [18] H. Billhardt, D. Borrajo, V. Maojo, Using genetic algorithms to find suboptimal retrieval
377 expert combinations., in: *SAC*, ACM, 2002, pp. 657–662.
- 378 [19] P. Pathak, M. Gordon, W. Fan, Effective information retrieval using genetic algorithms based
379 matching function adaptation, in: in *Proceedings of the 33rd Hawaii International Conference*
380 *on System Science (HICSS, 2000)*.

- 381 [20] C.-H. Lina, H.-Y. Chenb, Y.-S. Wua, Study of image retrieval and classification based on
382 adaptive features using genetic algorithm feature selection, *Expert Systems with Applications*
383 41 (2014) 66116621.
- 384 [21] W. Fan, M. D. Gordon, P. Pathak, A generic ranking function discovery framework by genetic
385 programming for information retrieval, *Inf. Process. Manage.* 40 (4) (2004) 587–602.
- 386 [22] W. Fan, M. D. Gordon, P. Pathak, Personalization of search engine services for effective re-
387 trieval and knowledge management, in: *Proceedings of the Twenty First International Con-*
388 *ference on Information Systems, ICIS '00*, Association for Information Systems, Atlanta, GA,
389 USA, 2000, pp. 20–34.
- 390 [23] J. R. Koza, *Genetic Programming: On the Programming of Computers by Means of Natural*
391 *Selection*, MIT Press, Cambridge, MA, USA, 1992.
- 392 [24] Makarov, Metric properties of the functions of distances between molecular graphs, *Journal of*
393 *Structural Chemistry* 2 (2007) 223–229.
- 394 [25] J. M. Morris, Traversing binary trees simply and cheaply., *Inf. Process. Lett.* 9 (5) (1979)
395 197–200.
396 URL <http://dblp.uni-trier.de/db/journals/ip1/ip19.html#Morris79a>
- 397 [26] K. Zhang, D. Shasha, D.: Simple fast algorithms for the editing distance between trees and
398 related problems, *SIAM J. Comput.*