

(۱) پیش پردازش

عملیات پیش پردازش شامل سه بخش است که در فایل preprocess.py قابل مشاهده است.

tokenizer(۱/۱)

در این بخش تمرکز روی جدا کردن کلمات از

یکدیگر و تولید توکن ها میباشد. در ابتدا چک میشود

که اگر آخر یک کلمه باشیم(به وسیله اسپیس یا خط جدید

یا...) کلمه را بعنوان توکن در نظر بگیریم و به سراغ

کلمه بعدی برویم.

سپس کلماتی مانند ((می خواهم)) بررسی میشوند

که بصورت یک کلمه ذخیره شوند و نه دو کلمه ((می)) و

((خواهم)).

در ادامه به بررسی مخفف ها میپردازیم و آنها را

به کلمه اصلی خود تبدیل میکنیم.

و در آخر ایمیل و آیدی را پردازش میکنیم.

```
42
43 def tokenize(input:str)->list[str]:
44     output = []
45     word = []
46     for i in range(len(input)):
47         if input[i] in space_chars:
48
49             if len(word) < 1:
50                 continue
51             #check verbs
52             if "".join(word) in special_subwords_before:
53                 word.append(special_space)
54                 continue
55             #check abbreviations
56             if "".join(word) in abbrs:
57                 output.append(abbrs["".join(word)])
58                 word = []
59                 continue
60             #check id and email
61             if '@' in word:
62                 if len(word) == 1:
63                     temp = output.pop()
64                     output.append("@ " + "".join(temp))
65                     word = []
66                     continue
67                 elif word[-1] == '@':
68                     output.append("@ " + "".join(word[:-1]))
69                     word = []
70                     continue
71
72             output.append(" ".join(word))
73             word = []
74         else:
75             word.append(input[i])
76     output.append(" ".join(word))
77     return output
```

normalizer(۱/۲)

```
def normalize(input:list[str])>->list[str]:
    output = []
    for i in range(len(input)):
        if '@' not in input[i]:
            k=0
            number_flag = False
            word = []
            for j in range(len(input[i])):
                if input[i][j] in should_change_to:
                    word.append(should_change_to[input[i][j]])
                else:
                    word.append(input[i][j])

                if word[j] in to_be_deleted_chars or (number_flag and word[j] not in numbers):
                    output.append("".join(word[k:j]))
                    k = j + 1
                    number_flag = False
                elif word[j] in numbers and not number_flag:
                    output.append("".join(word[k:j]))
                    k = j + 1
                    number_flag = True

            if len(word[k:]) > 0:
                output.append("".join(word[k:]))

            if input[i] in special_subwords_after:
                temp = output.pop()
                temp2 = output.pop()
                output.append(temp2 + "|" + temp)
                continue

        else:
            output.append(input[i])

    return output
```

در این بخش به پردازش توکن های تولید شده میپردازیم. ابتدا چک میکنیم که اگر ایمیل یا آیدی بود در آن تغییری ایجاد نکنیم.

سپس به وسیله یک دیکشنری از حروفی که باید تغییر کنند حروف یکسان با شکل های متفاوت را یکسان میکنیم. (مثلا آ و ا)

سپس به بررسی حروفی که باید حذف شوند و همچنین پردازش اعداد میپردازیم.

در آخر کلماتی مانند زیباتر را بصورت یک کلمه در نظر میگیریم مه دو کلمه جدا از هم ((زیبا)) و ((تر)).

stemming(۱/۳)

```
def stemming(input:list[str])>->list[str]:
    output = []
    for i in input:
        temp = stem.stem(i)
        if len(temp) < 3:
            temp = i
        output.append(lem.lemmatize(temp, 'V'))
    return output
```

در این بخش توسط دو کتابخانه کلمات ریشه یابی میشوند تا کلمات هم ریشه مانند هم در نظر گرفته شوند.

۲) ساخت ایندکس

در این بخش توسط سه استراکچر به ساخت inverted index میپردازیم.

positional posting(۲/۱)

```
class PositionalPosting:
    def __init__(self, doc_id: int, tf : int, positions):
        self.doc_id = doc_id
        self.tf = tf
        self.positions = positions

    def getDocID(self):
        return self.doc_id

    def getTF(self):
        self.tf = len(self.positions)
        return self.tf

    def getPositions(self):
        return self.positions

class PostingsList:
    def __init__(self):
        self.df = 0
        self.list = []
        self.tf = 0
```

این استراکچر برای نگاه داشتن پستینگ مکانی برای یک توکن درون یک داکيومنت میباشد که آیدی داک مورد نظر و عداد تکرار توکن در داک و همچنین مکان های تکرار توکن در داکيومنت را نگه میدارد. همچنین شامل سه فانکشن برای گرفتن این اطلاعات میباشد.

PosingList(۲/۲)

```
class PostingsList:
    def __init__(self):
        self.df = 0
        self.list = []
        self.tf = 0

    def getPostings(self):
        return self.list

    def getTf(self):
        frequency = 0
        for posting in self.list:
            frequency += posting.getTF()

        self.tf = frequency
        return frequency

    def getDF(self):
        return self.df

    def addPosting(self, p: PositionalPosting):
        self.list.append(p)

    def addDF(self):
        self.df += 1

    def getAllDocIdTF(self):
        res = []
        for p in self.list:
            res.append((p.getDocID(), p.getTF()))
        return res
```

این استراکچر برای نگه داشتن تمام پستینگ های مکانی یک توکن میباشد که شامل اطلاعاتی مانند تعداد داکيومنت های حاوی آن توکن(df) و تعداد تکرار توکن در کل داکيومنت ها (tf) و لیستی از پستینگ های مکانی آن میباشد.

همچنین شامل فانکشن هایی برای دریافت این اطلاعات و همچنین افزودن پستینگ جدید و افزودن به df و دریافت tf های هم داکيومنت میشود.

InvertedIndex(۲/۳)

در این استراکچر عملیات ساخت و ذخیره سازی inverted index و همچنین تولید بردار های داکيومنت ها و ذخیره آنها و ذخیره لیست champion انجام می شود.

این بخش شامل سه دیکشنری برای تمام PostingList ها و لیست champion و همچنین بردار های نشان دهنده داکيومنت ها در فضای برداری میشود. همچنین تعداد تمام داکيومنت ها هم برای محاسبه امتیاز نگه داری میشود.

فانکشن هایی برای دریافت اطلاعات مختلف و همچنین افزودن پستینگ لیست جدید وجود دارند.

علاوه بر این سه فانکشن برای پاک کردن k

توکن پرتکرار استفاده شده است.

دو فانکشن هم برای سیو و لود کردن مدل استفاده شده است تا هر بار نیاز به تولید دوباره ایندکس نباشد.

در آخر هم دو فانکشن addScore و

CreateVers برای تولید بردار های داکيومنت ها تعریف شده اند.

```
class InvertedIndex:
    def __init__(self, docs_num : int):
        self.dict = {}
        self.docs_num= docs_num
        self.champions = {}
        self.vertices = {}

    def getPostingList(self, term: str) -> PostingsList:
        if term not in self.dict:
            return None
        return self.dict[term]

    def getChampionPostingList(self, term: str) -> PostingsList:
        if term not in self.champions:
            return None
        return self.champions[term]

    def getdocsNum(self):
        return self.docs_num

    def addPosting(self, term: str, doc_id: int, tf: int, positions = None):
        if term not in self.dict:
            self.dict[term] = PostingsList()

        self.dict[term].addDF()
        self.dict[term].addPosting(PositionalPosting(doc_id, tf, positions))

        if tf >= 3:
            if term not in self.champions:
                self.champions[term] = PostingsList()
            self.champions[term].addDF()
            self.champions[term].addPosting(PositionalPosting(doc_id, tf, positions))
```

```
def deleteRepeatedWords(self, k: int):
    heap = []
    for key in self.dict:
        heapq.heappush(heap, (-self.dict[key].getTf(), key))
    remove_list = []
    for _ in range(k):
        item = heapq.heappop(heap)
        self.delete(item[1])
        remove_list.append(item)

    heap = []
    for key in self.champions:
        heapq.heappush(heap, (-self.champions[key].getTf(), key))
    champion_remove_list = []
    for _ in range(k):
        item = heapq.heappop(heap)
        self.deleteChamp(item[1])
        champion_remove_list.append(item)

    return [(-key, value) for key, value in remove_list] , [(-key, value) for key, value in champion_remove_list]

def delete(self, term):
    self.dict.pop(term)

def deleteChamp(self, term):
    self.champions.pop(term)

def save(obj, path, mode):
    with open(path, mode) as file:
        pickle.dump(obj, file)
    print(f"{obj} saved.")
```

```

def load(path, mode):
    with open(path, mode) as file:
        obj = pickle.load(file)
    return obj

def addScore(self, token: str):
    postings = self.getPostingList(token)
    id_tf_pairs = postings.getAllDocIdTF()
    df = postings.getDF()
    for (id, tf) in id_tf_pairs:
        try:
            self.vertices[str(id)][token] = (1 + log10(tf)) * log10(self.getdocsNum() / df)
        except:
            pass

def createVers(self, data):
    for i in range(len(data)):
        doc_id = str(i)
        self.vertices[doc_id] = {}
    for i in self.dict:
        self.addScore(i)

```

```

f create_index(data, delete):
    inverted = InvertedIndex(len(data))
    for i in range(len(data)):
        content = data[str(i)]['content']
        tokens = preprocess.preprocess(content)
        terms = {}
        for j in range(len(tokens)):
            if tokens[j] not in terms:
                terms[tokens[j]] = [0, []]

            terms[tokens[j]][0] += 1
            terms[tokens[j]][1].append(j + 1)

        for key in terms:
            inverted.addPosting(key, i + 1, terms[key][0], terms[key][1])

        if i % 1000 == 0:
            print("docs processed: ", i)
    print("processing finished\ndeleting for inverted...")
    inverted_deletes = inverted.deleteRepeatedWords(delete)
    print("normal deletes = ", inverted_deletes[0])
    print("champion deletes = ", inverted_deletes[1])
    print("creating vers...")
    inverted.createVers(data)
    print("vers finished")

    return inverted

```

فانکشن create_index

برای تولید کردن ایندکس تعریف

شده است که در فایل main با

فراخوانی این تابع ایندکس ساخته

و سپس ذخیره میشود.

۲/۴) کلمات پرتکرار حذف شده

```
normal deletes = [(219489, 'و'), (164431, 'رد'), (133643, 'هب'), (131956, 'تسه#'), (92960, 'زا'), (838, 'نی'), (58166, 'ش#دش'), (59098, 'ار'), (68895, 'اب'), (69065, 'نک#درک'), (70005, 'هک'), (76187, 'نی'), (43344, 'شاپ#دوب'), (41085, 'راد#تشاد'), (30993, 'ارب'), (23648, 'کی'), (22371, 'مه'), (22261, 'وگ#تفگ'), (17525, 'دوخ'), (18866, 'ام'), (19767, 'میت'), (22020, 'هد#داد'), (22237, 'روشک'), (16969, 'یو'), (15486, 'اسا'), (15616, 'ات'), (15881, 'دیاب'), (16203, 'رب'), (16681, 'زاب'), (16876, 'ن'), (1, 'رازگ'), (13790, 'سراف'), (14232, 'سلجم'), (14555, 'ریم#درم'), (15052, 'ریا'), (15189, 'ام'), (5235, 'مایپ'), (11715, 'یاهتن'), (12250, 'رازگربخ'), (12413, 'سییر'), (12779, 'راک'), (12804, 'لود'), (3375, 'رارق'), (9497, 'هکنی'), (10055, 'يلم'), (10151, 'ور#تفر'), (10691, 'لاس'), (11201, 'نکیزاب'), (11373, 'ود'), (9124, 'ه اوخ#تساوخ'), (9416, 'و')
champion deletes = [(217244, 'و'), (162665, 'رد'), (131464, 'هب'), (128105, 'تسه#'), (88577, 'زا'), (7, 'نی'), (52792, 'ش#دش'), (52989, 'ار'), (63399, 'اب'), (63655, 'نک#درک'), (64697, 'هک'), (71306, 'نی'), (37646, 'شاپ#دوب'), (35296, 'راد#تشاد'), (24695, 'ارب'), (18404, 'میت'), (18063, 'مه'), (1793, 'دوخ'), (12923, 'زاب'), (13290, 'هد#داد'), (16074, 'ام'), (17070, 'روشک'), (17859, 'کی'), (8, 'وگ#تفگ'), (11877, 'اسا'), (11936, 'ریم#درم'), (12098, 'دیاب'), (12301, 'سلجم'), (12531, 'ن'), (8, 'یو'), (9323, 'ات'), (10024, 'سییر'), (10035, 'رب'), (10622, 'ام'), (10938, 'لود'), (11335, 'ریا'), (598, 'هکنی'), (6475, 'يلم'), (7002, 'لاس'), (7378, 'بالقنا'), (7747, 'راک'), (8572, 'نکیزاب'), (598, 'وناق'), (5529, 'اروش'), (5667, 'ه اگشاب'), (5773, 'نم'), (5812, 'روهج'), (6004, 'بوخ'), (6049, 'ور#تفر'), (5147, 'ن')]
```

کلمات حذف شده لیست اصلی:

و-۲۱۹۴۸۹	در-۱۶۴۴۳۱	به-۱۳۳۶۴۳	هست (و انواع آن)-۱۳۱۹۵۶
از-۹۲۹۶۰	این-۸۳۸۳۸	که-۷۶۱۸۷	
کرد (و انواع آن)-۷۰۰۵	با-۶۹۰۶۵	را-۶۸۸۹۵	
شد (و انواع آن)-۵۹۰۹۸	است-۵۸۱۶۶	بود-۴۳۳۴۴	داشت-۴۱۰۸۵

و...

کلمات حذف شده لیست champion:

مانند حالت اصلی است بغیر از اواخر که فرق میکند.

QueryProcessor(۳)

این بخش که در فایل

querryrunner.py قرار گرفته است

وظیفه پردازش کویری کاربر و برگرداندن

نتیجه را دارد که توسط فایل runner.py

فراخوانی میشود.

دو فانکشن اصلی آن findNormal و

FindChampion هستند که ابتدا از لیست

Champion برای جستجو استفاده میشود و در

صورتی که به وسیله آن تعداد داکيومنت دلخواه

به دست نیاید سراغ normal میرویم.

در هر دو فانکشن ابتدا کلماتی از کویری

که در دیکشنری نیستند حذف میشوند و امتیاز

باقی کلمات محاسبه میشود. و سپس داکيومنت

ها به صورت برداری دریافت میشوند و به

وسیله ی فانکشن دیگری شباهت کوسینوسی

بین کویریو بردار های داکيومنت ها به دست

می آید و داکيومنت ها با امتیاز بیشتر برگردانده

میشوند.

```
6 class QueryProcessor:
7     def __init__(self, index: index.InvertedIndex):
8         self.index = index
9
10    def findNormal(self, query: str, k: int):
11        tokens = preprocess.preprocess(query.strip())
12        remove_list = []
13        for i in range(len(tokens)):
14            pl = self.index.getPostingList(tokens[i])
15            if pl is None:
16                remove_list.append(i)
17        for i in remove_list:
18            tokens.pop(i)
19
20        query_scores = {}
21        query_frequency = {}
22        for i in tokens:
23            if i not in query_frequency:
24                query_frequency[i] = 0
25            query_frequency[i] += 1
26
27        for i in query_frequency:
28            query_scores[i] = 1 + log10(query_frequency[i])
29
30        doc_scores = {}
31        for i in tokens:
32            scores = self.getScoreList(i)
33            for id in scores:
34                if id not in doc_scores:
35                    doc_scores[id] = {}
36                doc_scores[id][i] = scores[id]
37
38        heap = []
39        for docID, doc_vector in doc_scores.items():
40            doc_score = self.getSimilarity(query_scores, doc_vector)
41            heapq.heappush(heap, (-doc_score, docID))
42
43        if len(heap) < k:
44            k = len(heap)
45
46        result = []
47        for _ in range(k):
48            neg_score, id = heapq.heappop(heap)
49            result.append((id, -neg_score))
50
51        return result
52
53    def findChampion(self, query: str, k: int):
54        tokens = preprocess.preprocess(query.strip())
55        remove_list = []
56        for i in range(len(tokens)):
57            pl = self.index.getChampionPostingList(tokens[i])
58            if pl is None:
59                remove_list.append(tokens[i])
60        for i in remove_list:
61            tokens.remove(i)
62
63        query_scores = {}
64        query_frequency = {}
65        for i in tokens:
66            if i not in query_frequency:
67                query_frequency[i] = 0
68            query_frequency[i] += 1
69
70        for i in query_frequency:
71            query_scores[i] = 1.0 + log10(query_frequency[i])
72
73        doc_scores = {}
74        for i in tokens:
75            scores = self.getScoreListChampion(i)
76            for id in scores:
77                if id not in doc_scores:
78                    doc_scores[id] = {}
79                doc_scores[id][i] = scores[id]
80
81        heap = []
82        for docID, doc_vector in doc_scores.items():
83            doc_score = self.getSimilarity(query_scores, doc_vector)
84            heapq.heappush(heap, (-doc_score, docID))
85
86        if len(heap) < k:
87            result = []
88            for _ in range(len(heap)):
89                neg_score, id = heapq.heappop(heap)
90                result.append((id, -neg_score))
```



```

89         result.append((id, -neg_score))
90     normal_res = self.findNormal(query, 2 * k)
91     for i in range(len(normal_res)):
92         if normal_res[i] not in result:
93             result.append(normal_res[i])
94     else:
95         result = []
96         for _ in range(k):
97             neg_score, id = heapq.heappop(heap)
98             result.append((id, -neg_score))
99     return result
100
101 def getScoreList(self, token: str):
102     scores = self.index.vertices
103     postings = self.index.getPostingList(token)
104     result = {}
105     for i in postings.list:
106         result[str(i.getDocID())] = scores[str(i.getDocID())][token]
107     return result
108
109 def getScoreListChampion(self, token: str):
110     scores = self.index.vertices
111     postings = self.index.getChampionPostingList(token)
112     result = {}
113     for i in postings.list:
114         result[str(i.getDocID())] = scores[str(i.getDocID())][token]
115     return result
116
117 def getSimilarity(self, vector1: dict, vector2: dict):
118     squared_sum1 = 0.0
119     squared_sum2 = 0.0
120
121     for t in vector1:
122         squared_sum1 += vector1[t] ** 2
123     for t in vector2:
124         squared_sum2 += vector2[t] ** 2
125
126     squared_sum1 = sqrt(squared_sum1)
127     squared_sum2 = sqrt(squared_sum2)
128
129     score = 0.0
130     for i in vector1:
131         if i in vector2:
132             score += (vector1[i] / squared_sum1) * (vector2[i] / squared_sum2)

```

```

    for i in vector1:
        if i in vector2:
            score += (vector1[i] / squared_sum1) * (vector2[i] / squared_sum2)
    return score

```

۴) کویری ها

با اجرای فایل runner.py میتوان کویری زد و نتیجه را دریافت کرد.

در زیر چند کویری نشان داده شده است.

۴/۱) شکایت:

نتایج کویری یک کلمه ای ساده ((شکایت)) را میتوان در زیر مشاهده کرد. به دلیل اینکه کویری یک کلمه ای است امتیاز tfidf آن به وسیله نرمالسازی محاسبه شده است، برای اسنادی که حداقل یکبار آنرا دارند ۱ میشود. همچنین برای هر داک مکان هایی که شکایت آمده است نشان داده شده است. نتایج بدست آمده مرتبط هستند و در حوزه موارد قضایی و شکایت میباشند.

```
Enter query : شکایت
time : 0.0011661052703857422
doc_id: 10824
Score: 1.0
Title: دش رشتنم سلجم تابومم هرابرد نابهگن یاروش تارظن تایئزج
URL: https://www.farsnews.ir/news/14000829000328/دش-رشتنم-سلجم-تابومم-ه-رابرد-ن-ابهگن-ی-اروش-تارظن-تایئزج
In doc:

ایس ای گدی سروتی اکشلب اقی ای افیق ع ارم ری اسرد
=====
ارن آل ع ج ح تی اکشد روم لم عل اروتسد عضو ع ر م ه چ ن ا ن ج
=====
و د در گ ر ا رض ابج و م تی اکشد روم هم ا ن ی ی آ ه کن ی ا _ : د در گیم
=====
ه تخ ان ش ع ر ش ف الخ ه سفن ی ق تی اکشد روم ه بومم ا ذ ل دش ن ز ا ر ح ا
=====
نوم ز آ ی ر ا ز گ ر بی ه گ آ و تی اکشد روم ی غ الب ا زوج م _ : د در گیم
=====
ز ا ر گ ا ل ک اضری ل ع ی اق آ تی اکش صومخ ر دن ابه گ ن ی ا روشم ر ح م ی ا ه ق ف
=====
```

doc_id: 11406
Score: 1.0
Title: دنک ققوتم ار راگنربخ کی زا دوختی اکش تفتن ترازو
URL: <https://www.farsnews.ir/news/14000816000080/>
In doc:

کی ز اتفتن ترازوتی اکش تفتن ترازو
=====

=====

دن اه تفرگر ارقی ای امتی اکش دروم تفتن ترازو وسزا
=====

=====

رد امساره دشد ای تی اکش تفتن ترازو تمیم و معر اکفا اب
=====

=====

doc_id: 11524
Score: 1.0
Title: نوناق تمرح طفیحات نیوچشناد زا هناردپ تیامح /؟دوب هج سلجم شرازگ هب بالقن اربهر رکذت یارجام
URL: <https://www.farsnews.ir/news/14000812000656/>
In doc:

تقوچی ام اء درکده اوختی اکش درقنی ام ادق ازا هک
=====

doc_id: 11741
Score: 1.0
Title: دوش حیحمیت یکشزپ میهنخت یرایتسد 49 هرود نومزآ غیرات
URL: <https://www.farsnews.ir/news/14000809000238/>
In doc:

48 نومزآ ابیطیترم تی اکش در 45 بعشز اهدش
=====

=====

فلکم هکلبه دشارج ال ام زالتی اکش فطرطی ارب اهن تن، هدش
=====

=====

doc_id: 11830
Score: 1.0
Title: ؟دوب هج ااهیسروب یارجام / هدنورپ کی گرزب نافلختم یریگن اهج و ین اهور
URL: <https://www.farsnews.ir/news/14000726000151/>
In doc:

ترازو ااهیسروب ین ایوچشناد تی اکش* د ادربخ رادل کشم و
=====

=====

630 هبومیم ام اء دن کتی اکشه ن اختر ازون یازادنه اوخی ئافق
=====

=====

ئی اضیق کاحم رد ااهیسرویتی اکشه زاج ایگنه رقبالقن ایل ای اروش
=====

=====

Score: 1.0
Title: دنک کمک ام لایتوف هب دن اوت یم چیچوکسا /دوب نازلل مد اخیزیزع یارب هم ه زا رتشیب تسایر یلدننص :یژیزع
URL: <https://www.farsnews.ir/news/14001207000863/> یم-چیچوکسا-دوب-ن-ازلل-م-د-اخیزیزع-ی-ارب-هم-ه-زا-رتشیب-تسایر-یلدننص-یژیزع
In doc:

رگی دکیزا اهم یترم تسیم تی اکشون اری ال ایتوفریخ ای شاوچ
=====

دن تنگی یم یاه رودرد .دن کی م تی اکش، دروخی م تسپک شه کی م یتره
=====

یل اجنم یارب، دن رادتی اکشم ه زا هم ه کن اری ا
=====

doc_id: 1342
Score: 1.0
Title: ه دش رانکوب تسرپرس و یزیریت ه اگشاپ ره م بیج ع یارج ام /داب آ م رخ ریخ هب یزاس نیشام نکیز اب ل ا قتن ا رد ل ا جنج
URL: <https://www.farsnews.ir/news/14001206000543/> ره م-بیج-ع-ی-ارج-ام-د-اب-آ-م-رخ-ریخ-ه-ب-یزاس-نیشام-نکیز-اب-ل-ا-قتن-ا-رد-ل-ا-ج-ن-ج
In doc:

هیرگا ام ادسیرنتی اکش هیر اکا اتدد رگرین کی ز اب
=====

ل ا ورقبط ارن ام دوختی اکش ا در فاعطق، دن شا به تشا در ارسا
=====

doc_id: 139
Score: 1.0
Title: ن اناوچ گیل رد ین ایت ه بیخاش هب یگدی سر تباب یدج ام زا سیر اف جیلخ ل ایتوف میت یبرم ری دقت
URL: <https://www.farsnews.ir/news/14001222000297/> رد-ین-ایت-ه-بیخاش-ه-ب-یگدی-سر-تباب-ی-دج-ام-زا-سیر-اف-جیلخ-ل-ایتوف-میت-یبرم-ری-دقت
In doc:

یدج ام د اشریم روتسد ایتی اکش نی ایپرد .درکض ارتع ا
=====

جیلخ ن اناوچ ل ایتوف میت یختی اکشد روم ردسیر افش زور اگن ریخ
=====

تیتز اسپ، مید رکتی اکشل ایتوفن ویسا ا ردفه بیز ابه جیتن
=====

ن ویسا ا ردفتمیرپرسی دج ام د اشریم تی اکش تیتز اسپ، مید رکتی اکش
=====

doc_id: 1584
Score: 1.0
Title: سگع+دش ی رطق هن اسیر ه ژوس ل القیخسا تس اوخرد هب CAS در تسبد
URL: <https://www.farsnews.ir/news/14001203000704/> سگع-دش-ی-رطق-هن-اسیر-ه-ژوس-ل-القیخسا-تس-اوخرد-هب-CAS-در-تسبد
In doc:

```
doc_id: 1737
Score: 1.0
Title: دن‌آه‌دادن ار مم‌کج -زونه -ام ا -دن‌تفری‌ذپ -ار مم‌سایر -این‌د لک :یع‌اس
URL: https://www.farsnews.ir/news/14001201000714/
In doc:
```

نم :درکن اش‌ن‌رط‌اخ ،دن‌آه‌درکتی اکش‌رفن 15دن‌تفگی‌م‌ه‌کنی‌ان‌ایب

ه‌ک‌ذ‌وب‌ه‌دش‌تبیثتی اکش‌گی‌آه‌نت‌م‌آه‌دی‌دن‌ار‌رفن

رد :درکن اش‌ن‌رط‌اخ ،دن‌آه‌درکتی اکش‌رگی‌دد‌ار‌ف‌ار‌وط‌ج‌ود‌ن‌د‌وب

نم ،تس‌ین‌ع‌فن‌ی‌ذه‌درکتی اکش‌ه‌کی‌رفن‌ه‌ک‌د‌ش‌ر‌ط‌م

Enter query :

۴/۲) کریسمس

نتایج کریسمس به عنوان کویری تک کلمه ای سخت را در زیر مشاهده میکنید. مانند بخش قبل بدلیل استفاده از tfidf و نرمالسازی امتیاز همه ۱ شده است. همچنین زمان آن از کویری ((شکایت)) کمتر شده است زیرا کلمه نادرتری است و هنگام index elimination داک های بیشتری حذف میشوند.

بیشتر کویری های برگردانده شده مرتبط هستند ولی برخی از آنها ارتباط کمی دارند.

```
Enter query : سم‌سیرک
time : 0.000982522964477539
doc_id: 6117
Score: 1.0
Title: سکع+سم‌سیرک رط‌اخ ه‌ب‌آه‌یرم‌م‌اب‌یل‌ا‌غ‌ت‌ر‌پ‌درم‌ف‌ال‌ت‌خ‌ا‌ت‌ف‌ر‌گ‌ب‌ی‌ل‌«ر‌وت‌ا‌ن‌ک‌د»‌ش‌و‌ریک
URL: https://www.farsnews.ir/news/14001005000165/
In doc:
```

قی‌وع‌ت‌ه‌ب‌ار‌ود‌را‌سم‌سیرک‌ت‌الی‌ط‌ع‌ت‌لی‌لد‌ه‌ب‌یل‌و‌دن‌ک

ی‌آه‌ت‌ن‌ا‌.ت‌س‌ا‌ه‌ت‌ف‌ر‌گ‌ال‌اب‌سم‌سیرک‌ت‌الی‌ط‌ع‌ت‌رط‌اخ‌ه‌ب‌نی‌ف‌ر‌ط‌نی‌ب

```
doc_id: 5431
Score: 1.0
Title: ت‌ار‌ا‌خ‌ت‌ف‌ا‌ت‌ش‌گ‌ز‌اب‌دی‌رد‌ام‌ل‌ا‌ئر‌ه‌ب‌ن‌ادی‌ز‌ع‌ت‌ق‌و
URL: https://www.farsnews.ir/news/14001014000983/
In doc:
```

ل‌ا‌ئر‌یه‌ار‌ز‌ر‌ب‌ون‌ی‌ت‌ن‌ر‌ول‌ف‌سم‌سیرک‌ه‌یده‌ن‌ا‌ون‌ع‌ه‌ب‌ه‌ی‌ون‌ا‌ژ‌4

doc_id: 5483
Score: 1.0
Title: دددرگ یم رب سیراپ هب ینام ز هج یسم
URL: <https://www.farsnews.ir/news/14001014000228/> دددرگ-یم رب-سیراپ-هب-ینام ز-هج-یس
In doc:

الحجیم سوریو نی ا هب سم سیرک م ای ا رد ناهج نکیزاب نیرتهب

doc_id: 5921
Score: 1.0
Title: درادن مه نیرمت نیمز کی ناردنزام زبیسرس ناتسا/میتسه یزوریپ هب م وکحم نام یلس دجسم تن لباقم یم اهلا
URL: <https://www.farsnews.ir/news/14001007000852/> زبیسرس-ناتسا-میتسه-یزوریپ-هب-م وکحم-نام یلس-دجسم-تن-لباقم-یم اهلا
In doc:

ناوریپ همه هب ار سم سیرک :تشاد راهظا نام یلس دجسم تن

doc_id: 5926
Score: 1.0
Title: سکع+ دنادرگرب لیصحت لحم هب ار کدوک 4، ناجنسفر سم یجراخ مچاهم

doc_id: 5926
Score: 1.0
Title: سکع+ دنادرگرب لیصحت لحم هب ار کدوک 4، ناجنسفر سم یجراخ مچاهم
URL: <https://www.farsnews.ir/news/14001007000809/> سکع-دنادرگرب-لیصحت-لحم-هب-ار-کدوک-4-ناجنسفر-سم-یجراخ-مچاهم
In doc:

doc_id: 5933
Score: 1.0
Title: سکع+یواژ هب الویدراوگ سم سیرک هیده ؛یی این ایسا هراتس
URL: <https://www.farsnews.ir/news/14001007000739/> سکع-یواژ-هب-الویدراوگ-سم سیرک-هیده-یی این ایسا-هراتس
In doc:

یواژ هب الویدراوگ پپ سم سیرک هیده هب سروت ،جرط نی ا

doc_id: 6120
Score: 1.0
Title: دیدج لاس هناتسآ رد لایتوف هگشزرو رد راتشک
URL: <https://www.farsnews.ir/news/14001005000143/> دیدج-لاس-هناتسآ-رد-لایتوف-هگشزرو-رد-راتشک
In doc:

لایتوف نیمز کی رد سم سیرک نشج م اگنه رد رفن جنب

=====

doc_id: 6148
Score: 1.0
Title: ؟دوش یم رترب گیل یه ار رادرس؛لساکوین ای نوتروا/نوم زآ بذج یارب نویل هب سیلگن ا میت هبرض
URL: <https://www.farsnews.ir/news/14001004000708/>یه ار رادرس؛لساکوین-ای نوتروا-نوم زآ-بذج-یارب-نویل-هب-سیلگن ا میت-هبرض/
In doc:

یم رطن هب .دنسیونب سمسیرگ یارب نالوط رایبب تسرهف دن دوب
=====

doc_id: 6697
Score: 1.0
Title: ده دیم ناشن ار وا توفیک چیچوکسا جیاتن /تشاذگ ریثات یلیخ ام تنه ذیور شوریگ /دوب رتهب یسلج هب نم نادرگرب تفگ یمراط :شجین اه ج
URL: <https://www.farsnews.ir/news/14000927000153/>ام-تنه ذیور-شوریگ-دوب-رتهب-یسلج-هب-نم-نادرگرب-تفگ-یمراط-شجین اه ج/
In doc:

تسا توافقتم سیلگن ا رد سمسیرگ .تسا دیدج قافتا کی هک
=====

doc_id: 9413
Score: 1.0
Title: دوشوم دانگدوب سمسیرگ هلاکاتوآ اب یغید یلمتولقا ناگدتیلمن هلاک :نوم سمسیرگ شلاع

۴/۳) لیگ برتر:

نتایج کویری ((لیگ برتر)) به عنوان کویری چند کلمه ای ساده را در زیر مشاهده میکنید.
مشاهده میشود که دیگر تمام امتیاز ها ۱ نداشتند زیرا کویری چند کلمه ای است و داک ها بر اساس امتیاز و شباهت به کویری مرتب شده اند.

کویری های برگردانده شده مرتبط هستند و در ارتباط با لیگ های برتر انگلیس یا ایران هستند.

```
PS C:\Users\parsa> & C:/Users/parsa/AppData/Local/Programs/Python/Python39/python.exe g:/IR-Project/runner.py
Enter query : رترب گیل
time : 0.0029997825622558594
doc_id: 5134
Score: 0.9999968932456669
Title: دراودوو دتیانوی رتسچنم اب یظفاح ادخ ؛ م اج رفص ، وروی درایلیم کی
URL: https://www.farsnews.ir/news/14001018000474/کی-رفص-وروی-درایلیم-کی
In doc:
```

```
و دناسرن سفلگن ا رترب گیل م اج کی هب یتح ار
=====
رد ینام رهق و اپورا گیل رد ینام رهق کی هب و
=====
رد .دوب اپورا ن انام رهق گیل هب ندیسر اهنت دتیانوی هج یتن
=====
زین ار اپورا ن انام رهق گیل هیمهس هتیل .دیسر سفلگن ا یفدح
=====
نیلا و دیوب ایسدا گیل رد دتیانوی رتسچنم ینام رهق نآ هج یتن
```

```
doc_id: 4890
Score: 0.9999913339558058
Title: لمف مین ینام رهق هن اتس آ رد ن اشوپ یب آ سفل و پسرپ رد ه لزلزل ، ل القتسا رد ن شج ؛ رترب گیل م ه دزن اپ هتفه
URL: https://www.farsnews.ir/news/14001021000814/رد-ن-اشوپ-یب-آ-سفل-وپسرپ-رد-ه-لزلزل-ل-الق-تسا-رد-ن-شج-رترب-گیل-م-ه-دزن-اپ-ه-تفه
In doc:
```

```
ن ای اپ هب لایتوف رترب گیل لمف مین تاقب اسم هین شچنپ ، سراف
=====
کی دن اوت یم دنتسه هی اسم ه گیل لودج رد هک ناری لایتوف
=====
تیت مکی و تهیب گیل رد ار رگید یایی ز و
=====
م ه ن اج نسفر سم تس ا گیل ورشیپ یاه میت زا یکی
=====
درا د رایخ ا رد ار گیل قیغض عاقد طخ نیمود ه دز
=====
```


doc_id: 2453
Score: 0.9999680009022884
Title: ؟ن اچ ی الپ ای اشگه ار لاسیتوف گیل نامزاس تامیمیت
URL: https://www.farsnews.ir/news/14001119001101/نامزاس-تامیمیت-ار-لاسیتوف-گیل-ای-اشگه
In doc:

۱۳۹۸ لاس رد لاسیتوف گیل نامزاس هسپئر تایه تاباختنا ،گیل
=====

نی ا اب .دوب لاسیتوف گیل نامزاس سیسات ودب زا هسپئر
=====

باختنا اب .دش هدرپس گیل نامزاس ریبد یرباج شرآ هب
=====

ریبد مه و لاسیتوف گیل نامزاس رادن اکس مه یرباج ،لاسیتوف
=====

مه لاسیتوف نویس اردف و گیل نامزاس هم انساسا و نوناق رد .دش
=====

Score: 0.9999420218152351
Title: میگنجیم نامدوخ یبلطاج یارب هدنیآ هتفه 9 رد /تسا روسج مییت راداو ه :یعیبر
URL: https://www.farsnews.ir/news/14001220000217/نامدوخ-یبلطاج-یارب-هدنیآ-هتفه-9-رد-تسا-روسج-مییت-راداو-ه-یعیبر
In doc:

و درک یزاب یرترب گیل مییت ود اب هتشذگ یزاب
=====

مینک فذخ ار یرترب گیل مییت کی میتمیناوت یفذخ م اچ
=====

هتفه 21 رد .میشاب گیل زاسیتفگش هک دراد دوجو مییت
=====

دوب یرگید یراذگفده رترب گیل رد ام یراذگفده .میاهد درک شالط
=====

زا مینکیم یس ،رترب گیل هدنامیقاب هتفه 9 و هدش
=====

امش زا :تفگ رترب گیل رد اهنآ روضح و شمیت
=====

فتک هیحان زا دشیم گیل رد ام موس هتفه زا
=====

عل عملان ذایق دایک رد گمل نایکم نایب سیرت راذگفده شات سیرترب نایب

doc_id: 4952
Score: 0.9999420218152351
Title: دوب ده اوخ لایتوف ه دنی آ عفن هب سیل وپسرب و لالقتسا فذح/تسا باداش و لاجرس ناولم :ه دازدمح ا
URL: <https://www.farsnews.ir/news/14001021000138/> ه دازدمح ا
In doc:

د: درک ناو ن ع ، ددرگرب رترب گیل هب اهلاس زا سپ میت
=====

یتخس رایسب راک کی گیل رد ینیشنردس موات . دندن ادیم هم ه
=====

نی ا مساسح ا . دنشاب کی گیل میت نیرت هب دندنکیم شالت و
=====

گیل یاهیزاب :تفگ رترب گیل توفیک و روشک لایتوف یلک
=====

لایند یلیخ ار رترب گیل یاهیزاب :تفگ رترب گیل توفیک
=====

لایتوف تم الس ده اش یاهقرح گیل عورش لوا لاس 10 لثم
=====

یو . مینیب میتسین اوتیم ار گیل تم الس لقادح اهلاس نآ رد

doc_id: 5576
Score: 0.9999420218152351
Title: سیکع+دش هیسور گیل نکیزاب نیرت نارگ یلم میت هراتس
URL: <https://www.farsnews.ir/news/14001012000643/> هراتس
In doc:

10 یفرعم هب هیسور گیل لوا لیمف مین مامت اب
=====

مجاهم نومزآ رادرس . تخادرب گیل نی تمیق نارگ نکیزاب 10
=====

تشاد رارق روشک نی گیل نانکیزاب نیرت نارگ مود هدر
=====

هراتس . دش هیسور رترب گیل نکیزاب نیرت نارگ وروی نویلیم
=====

دش لیدبت هیسور رترب گیل نکیزاب نیرت نارگ هب شا
=====

زاغآ تینز اب هیسور گیل رد یناری هراتس ییالط نارود
=====

کی نومزآ نینچم ه . دش گیل نامرهق راب 3 دش قفوم
=====

۴,۴) جیانی، اینفانتینو:

بعنوان کویری چندکلمه ای دشوار ((جیانی اینفانتینو)) در نظر گرفته شده است که اسم خاصی است. مشاهده میشود که سرعت جوابگویی آن از کویری چمد کلمه ای ساده بیشتر است(به دلیل تاثیر index elimination) ولی امتیاز های آن کمتر از کویری ساده است زیرا تکرار کمتری دارد و امتیاز و شباهت کمتری دارد.

چندتا از کویری ها کاملاً مرتبط هستند و در مورد فرد جستجو شده هستند ولی بیشتر آنها کاملاً مرتبط نیستند.

```
Enter query : ونیئن افنی ،ین ایج
time : 0.001157999038696289
doc_id: 3457
Score: 1.0
Title: ؟تسا هتفر چی وکسا رگم /میه دیم تایل ام شوریک یارب زونه /تسا کرش یرگوداج یلو دندرک م هتم یرگوداج هب ار ام :مداج یزیزع
URL: https://www.farsnews.ir/news/14001108000181/ یارب-زونه-تسا-کرش-یرگوداج-یلو-دندرک-م-هتم-یرگوداج-هب-ار-ام-:مداج-یزیزع
In doc:

هب افیف سییر ،ونیئن افنی ،ین ایج رفس وغل لیلد دروم رد
=====
-----

doc_id: 4033
Score: 0.9999983756898863
Title: دوش یم یسرب یی اسمش یی این ایسا ن ارایتسد تیعضو /ن ارهت هب ونیئن افنی رفس وغل لامتحا
URL: https://www.farsnews.ir/news/14001102000084/ یی اسمش-یی این ایسا-ن ارایتسد-تیعضو-ن ارهت-هب-ونیئن افنی-رفس-وغل-لامتحا
In doc:
```

```
doc_id: 3457
Score: 0.9981087778902797
Title: ؟تسا هتفر چی وکسا رگم /میه دیم تایل ام شوریک یارب زونه /تسا کرش یرگوداج یلو دندرک م هتم یرگوداج هب ار ام :مداج یزیزع
URL: https://www.farsnews.ir/news/14001108000181/ یارب-زونه-تسا-کرش-یرگوداج-یلو-دندرک-م-هتم-یرگوداج-هب-ار-ام-:مداج-یزیزع
In doc:

هب افیف سییر ،ونیئن افنی ،ین ایج رفس وغل لیلد دروم رد
=====
-----

doc_id: 6731
Score: 0.9913235060590921
Title: ین اهج م اج ه اگشزرو رد ن اگراتس نیدام ن یزاب رد ای یوده م روضح
URL: https://www.farsnews.ir/news/14000926000283/ ین اهج-م اج-ه اگشزرو-رد-ن اگراتس-نیدام-ن یزاب-رد-ای-یوده-م-روضح
In doc:

-----
-----

doc_id: 1020
Score: 0.7071067811865475
Title: نیارکوا و نیطسلف هلئسم رد افیف فیکانتم راتفر /تسایس زا شزرو یی ادج» م ان هب یلاخ و تراغش
URL: https://www.farsnews.ir/news/14001209000606/ هلئسم-رد-افیف-فیکانتم-راتفر-تسایس-زا-شزرو-یی-ادج-م-ان-هب-یلاخ-وت-راغش
```

doc_id: 1140
Score: 0.7071067811865475
Title: دش تارام-اقارغ-ن ابزیم-دادغب/ناریا هورگم هین ابزیم میرحت عفر
URL: <https://www.farsnews.ir/news/14001209000085/>
In doc:

doc_id: 127
Score: 0.7071067811865475
Title: دوشیم ناریا-ی ابقر-یزاب-رگاشامت-افیف-سئئر
URL: <https://www.farsnews.ir/news/14001222000684/>
In doc:

doc_id: 2041
Score: 0.7071067811865475
Title: هلاس-2-ین اهه-م اج-فل احم-اهتسیل ابیتوف-دمرد 75
URL: <https://www.farsnews.ir/news/14001127000839/5->
In doc:

doc_id: 4580
Score: 0.7071067811865475
Title: دی آیم ناریا هین وین افنی ا دشاین کانرطخ طیارش رگا/دوش تحاران چی وکسا هک دشن هدزیدب فرح: ل ابیتوف نویساردف یوگنخس
URL: <https://www.farsnews.ir/news/14001025000590/>
In doc:

doc_id: 4841
Score: 0.7071067811865475
Title: ینامل آتکرش کی ییوری رازه 32 بیلط و شزرو ریو داقتن ا هب ل ابیتوف نویساردف شنکاو
URL: <https://www.farsnews.ir/news/14001022000544/>
In doc:

doc_id: 4910
Score: 0.7071067811865475

doc_id: 6370
Score: 0.7071067811865475
Title: دنوشیم فذح هخرج زا دننیین شزوم آ ار VAR ام نارواد/دهدیم هم ادا شراک هب تردق اب چی وکسا: تم داخیززع
URL: <https://www.farsnews.ir/news/14001001000441/>
In doc:

doc_id: 6449
Score: 0.7071067811865475
Title: نویساردف ره راطتن ا رد رالد نویلیم 19 /رایکی لاس ود ره یناهه م اج یرازگرب یارب افیف لاس دنفرت
URL: <https://www.farsnews.ir/news/14000930000060/>
In doc: