

گزارش پروژه پایانی بازیابی - صبا سبحان ۹۹۳۱۰۹۶

۱- ساخت شاخص مکانی

۱.۱ مجموعه داده:

```
def load_docs():
    docs = {}
    contents = []
    urls = []
    with open("IR_data_news_12k.json", 'r') as file:
        docs = json.load(file)
        for key in docs.keys():
            idx = str(key)
            docs[idx] = {'title': docs[idx]['title'],
                        'content': docs[idx]['content'],
                        'url': docs[idx]['url'],
                        }
            contents.append(docs[idx]['content'])
    return docs, contents, urls
```

```
docs, contents, urls = load_docs()
docs['1']
```

```
{
  'title': 'اعلام برنامه نشست خبری گل محمدی/ مجیدی هم باید به محل تمرین پرسپولیس برود!',
  'content': '\n'
  سوم لیگ برتر (جام خلیج فارس) با\xa0مدیریت سازمان لیگ و معاونی\xa0باشگاه میزبان (پرسپولیس) در ورزشگاه شهید کاظمی به شرح زیر\xa0برگزار
  می\xa0شود: چهارشنبه ۲۵ اسفند ساعت ۱۲ فرهاد مجیدی سرمربی استقلال ساعت ۱۳:۳۰ پچی گل محمدی سرمربی پرسپولیس\xa0مسابقه دو تیم روز پنجشنبه در ورزشگاه
  آزادی برگزار می\xa0شود. \xa0به گزارش خبرنگاری فارس، پیش از این باشگاه استقلال اعلام کرده بود قرار است نشست خبری فرهاد مجیدی، سرمربی این تیم از
  ساعت ۱۵ در سازمان لیگ برگزار شود. انتهای پیام /'\n\n',
  'url': 'https://www.farsnews.ir/news/14001224000865'
}
```

این کد فایل JSON شامل داده‌های خبری را بارگذاری می‌کند، اطلاعات هر خبر (عنوان، محتوا و لینک) را در یک دیکشنری به نام docs ذخیره می‌کند. همچنین محتوای تمام اخبار را در لیست contents قرار می‌دهد. خروجی تابع، دیکشنری docs و لیست محتواهاست به همراه لینک خبر.

۲.۱ پیش پردازش اسناد

۱. بارگذاری مخفف‌ها (load_abbreviations)

```
def load_abbreviations(file_path):
    """Loads abbreviations from the file into a dictionary."""
    abbreviations = {}
    with open(file_path, 'r', encoding='utf-8') as f:
        for line in f:
            try:
                entry = eval(line.strip())
                abbreviations.update(entry)
            except Exception as e:
                print(f"Error parsing line: {line.strip()} - {e}")
    return abbreviations
```

- فایل abbreviations.txt که شامل مخفف‌ها و معانی آن‌ها به صورت خط‌به‌خط است، خوانده می‌شود.
- هر خط به یک دیکشنری تبدیل شده و به دیکشنری کلی abbreviations اضافه می‌شود.

2. جایگزینی مخفف‌ها (expand_abbreviations)

```
def expand_abbreviations(text, abbreviations):
    """Expands abbreviations using the loaded abbreviations dictionary."""
    for short_form, expansion in abbreviations.items():
        pattern = fr"\b{re.escape(short_form)}\b"
        text = re.sub(pattern, expansion, text)
    return text
```

- با استفاده از دیکشنری abbreviations، تمامی مخفف‌ها در متن ورودی پیدا شده و با معانی کاملشان جایگزین می‌شوند.

3. اصلاح فاصله‌ها (spacing_correction)

```
def spacing_correction(text):
    """Applies spacing corrections."""
    for pattern, repl in SPACING_PATTERNS:
        text = re.sub(pattern, repl, text)
    return text
```

- الگوهای مشخصی برای تصحیح فاصله‌ها تعریف شده است.
 - مثال: تبدیل می خواهم به می‌خواهم یا اتصال پسوندهای فارسی مثل ها، ترین به کلمه اصلی با نیم‌فاصله.
- متن بررسی شده و تمامی فاصله‌های اشتباه اصلاح می‌شوند.

- هدف: ایجاد متنی خواناتر و سازگارتر برای پردازش.

4. حذف علائم اضافی (remove_punct)

```
def remove_punct(text):
    """Removes diacritics and extra punctuation."""
    text = re.sub(EXTRA_PUNCT_PATTERN, "", text)
    return text
```

- علائم نگارشی و دیگر کاراکترهای غیرضروری (مثل نقطه، ویرگول، نقل قول‌ها، ...) از متن حذف می‌شوند.
- این کار با استفاده از یک الگوی Regex انجام می‌شود که تمامی این کاراکترها را شناسایی و حذف می‌کند.

5. نرمال‌سازی متن (Normalizer)

- متن با استفاده از مازول هضم نرمال‌سازی می‌شود.
- تغییرات شامل موارد زیر است:
 - تبدیل کاف و ی عربی به فارسی.
 - حذف نیم‌فاصله‌های اضافه.
 - حذف علائم تزیینی یا غیرضروری.
- هدف: یکنواخت‌سازی ساختار متن.

6. توکنیزه کردن (word_tokenize)

- متن نهایی به توکن‌های مجزا (کلمات) تقسیم می‌شود.
- این مرحله، خروجی نهایی را به شکل لیستی از کلمات آماده می‌کند.
- هدف: آماده‌سازی متن برای مراحل بعدی تحلیل (مثل محاسبه فراوانی کلمات یا ساخت شاخص معکوس).

7. تابع اصلی (preprocess_single_text)

```
def preprocess_single_text(text, expand_abbr=True, do_spacing_corr=True,
                           remove_diacritics=True):

    if expand_abbr:
        text = expand_abbreviations(text, abbreviations)

    text = normalizer.normalize(text)

    if do_spacing_corr:
        text = spacing_correction(text)

    if remove_diacritics:
        text = remove_punct(text)

    tokens = word_tokenize(text)
    return tokens
```

- این تابع مراحل بالا را با توجه به ورودی‌های کاربر اجرا می‌کند:
 - جایگزینی مخفف‌ها (در صورت نیاز).
 - اصلاح فاصله‌ها.
 - حذف علائم نگارشی.
 - نرمال‌سازی متن.
 - در نهایت توکنیزه کردن متن.
- خروجی: لیستی از توکن‌های پردازش‌شده

توضیح نرمالایزر Hazm:

استفاده از WordTokenizer:

جداسازی کلمات با استفاده از فاصله، نیم‌فاصله، تب و خطوط جدید.

شناسایی و ترکیب کلمات با پیشوندها و پسوندها (مثل "می‌روم").

توکنایز صحیح ایمیل‌ها، آیدی‌ها و اعداد.

```
self.email_pattern = re.compile(
    "{,r"[a-zA-Z0-9._+-]+@[a-zA-Z0-9-]+\.[A-Za-z]{2,}"
```

```
("(+[_self.id_pattern = re.compile(r"(?![\w._]) (@[\w
```

- WordTokenizer با قابلیت join_verb_parts=False برای شناسایی و جداسازی صحیح کلمات استفاده می‌شود. این ابزار:
 - فاصله‌ها و نیم‌فاصله‌ها را مدیریت می‌کند.
 - پیشوندها (مثل "می" و "نمی") را به درستی شناسایی می‌کند.

```

17 self._correct_spacing or self._decrease_repeated_chars.
    self.tokenizer = WordTokenizer(join_verb_parts=False)
    self.words = self.tokenizer.words

```

ترکیب پیشوندها و پسوندها:

- در تابع token_spacing، کلماتی که با نیم‌فاصله یا پسوند ترکیب می‌شوند (مثل "کتاب‌ها") شناسایی و ترکیب می‌شوند

```

result = []
for t, token in enumerate(tokens):
    joined = False

    if result:
        token_pair = result[-1] + "ZWNJ" + token
        if (
            token_pair in self.verbs
            or token_pair in self.words
            and self.words[token_pair][0] > 0
        ):
            joined = True

        if (
            t < len(tokens) - 1

```

مدیریت ایمیل‌ها، آیدی‌ها و اعداد:

- با استفاده از regex، ایمیل‌ها و اعداد شناسایی و به صورت توکن مستقل در نظر گرفته می‌شوند.

فاصله‌گذاری صحیح

نکات:

- شناسایی و اصلاح فاصله‌ها برای پیشوندها (مثل "می")، پسوندها (مثل "ها") و علائم نگارشی.

در hazm:

- اصلاح فاصله‌ها در تابع correct_spacing با الگوهای زیر انجام می‌شود:

حذف فاصله‌های اضافی:

```
self.extra_space_patterns = [  
    (r" {2,}", " "), # remove extra spaces  
    (r"\n{3,}", "\n\n"), # remove extra newlines  
    (r"\u200c{2,}", "\u200c"), # remove extra ZWNJs  
    (r"\u200c{1,} ", " "), # remove unneded ZWNJs before space  
    (r" \u200c{1,}", " "), # remove unneded ZWNJs after space  
    (r"\b\u200c*\B", ""), # remove unneded ZWNJs at the beginning of words  
    (r"\B\u200c*\b", ""), # remove unneded ZWNJs at the end of words  
    (r"[\r\n]", ""), # remove keshide, carriage returns  
]
```

اضافه کردن نیم‌فاصله به پیشوندها و پسوندها:

```
self.affix_spacing_patterns = [  
    (" ی ([ ^])"r, " یr\1), # space ی fix  
    (" (می|نمی)( |^)"r, r"\1\2\u200c"), # می, نمی put zwnj after  
    # تر, تری, ترین, گری, ها, های  
    (" تر, تری, ترین, گری, ها, های"r, r"\1\2\u200c"), # تر, تری, ترین, گری, ها, های put zwnj before  
    (  
        # ...  
    )  
]
```

اصلاح فاصله‌گذاری علائم نگارشی:

```
punc_after + "])", r"\1"), # remove space before + "])")
```

اضافه کردن فاصله بعد از اعداد:

```
, (r"(\d+)([اآبجچخدذرزژسشصضطظعفقکگلمنوهی])", r"\1 \2",
```

تعویض یونیکد

نکات:

- تبدیل کاراکترهای خاص (مثل "ي" و "ك") به فرم استاندارد فارسی.
- تبدیل " " به "بسم الله الرحمن الرحيم".

در hazm:

- تبدیل یونیکدها با استفاده از maketrans و جایگزینی خاص انجام می‌شو

تبدیل "ی" و "ک" عربی به "ی" و "ک"

```
(translations = maketrans(self.translation_src, self.translation_dst)
(text = text.translate(translations
```

جایگزینی کاراکترهای خاص:

```
] = self.replacements
, ("□", "بسم الله الرحمن الرحيم")
, ("ريال", "ريال")
, ("صلى", "صلى")
, ("الله", "الله")
, ("اكبر", "اكبر")
, ("محمد", "محمد")
, ("صلعم", "صلعم")
, ("رسول", "رسول")
, ("عليه", "عليه")
, ("وسلم", "وسلم")
, ("لا", "لا")
```

حذف کاراکترهای اضافی

حذف اعراب:

```
:if self._remove_diacritics
] = self.diacritics_patterns
remove FATHATAN, DAMMATAN, KASRATAN, FATHA, DAMMA, KASRA, SHADDA, #
SUKUN
, ("", "[u064b\u064c\u064d\u064e\u064f\u0650\u0651\u0652\]")
```

حذف علائم اضافی:

```
:if self._remove_specials_chars
```

```

] = self.specials_chars_patterns

Remove almost all arabic unicode superscript and subscript #
characters in the ranges of 00600-06FF, 08A0-08FF, FB50-FDFF, and FE70-FEFF

)

```

تبدیل اعداد انگلیسی به فارسی

```

:if self._persian_number

"self.number_translation_src = "0123456789%.۱۲۳۴۵۶۷۸۹

"self.number_translation_dst = ".۱۲۳۴۵۶۷۸۹٪.۱۲۳۴۵۶۷۸۹

```

جداسازی "می" و "نمی"

- اضافه کردن نیم فاصله بین پیشوند "می" و فعل.

در hazm:

- در تابع `seperate_mi`، پیشوندها با regex شناسایی و اصلاح می‌شوند:

```

: def seperate_mi(self: "Normalizer", text: str) -> str

```

توضیح توکنایزر در hazm:

```

def tokenize(self: "WordTokenizer", text: str) -> List[str]:

```

توکن‌های متن را استخراج می‌کند.

Examples:

```

>>> tokenizer = WordTokenizer()
(tokenizer.tokenize('این جمله (خیلی) پیچیده نیست!!!')) <<<
['این', 'جمله', '(', 'خیلی', ')', 'پیچیده', 'نیست', '!!!']
>>> tokenizer = WordTokenizer(join_verb_parts=False)

```

متد `tokenize` متن ورودی را به واحدهای کوچک‌تر به نام توکن تقسیم می‌کند. این واحدها شامل کلمات، علائم نگارشی، ایموجی‌ها، اعداد و سایر عناصر متن است. عملکرد آن به شرح زیر است:

ابتدا اگر قابلیت تشخیص مخفف‌ها فعال باشد، مخفف‌ها شناسایی شده و به‌طور موقت با مقادیر یکتا جایگزین می‌شوند تا در فرآیند شکستن متن، به بخش‌های جدا تقسیم نشوند. در پایان مخفف‌ها به حالت اولیه برمی‌گردند.

ایموجی‌های موجود در متن شناسایی شده و با یک فاصله از سایر بخش‌های متن جدا می‌شوند تا به‌عنوان توکن مستقل پردازش شوند.

عناصر خاص مثل ایمیل‌ها، لینک‌ها، آیدی‌ها (مانند @username)، هشتک‌ها، اعداد صحیح و اعشاری شناسایی شده و با مقادیر جایگزین مانند NUM، TAG، ID، LINK، EMAIL و NUMF جایگزین می‌شوند. این جایگزینی برای کاهش پیچیدگی متن انجام می‌شود.

علائم نگارشی مانند نقطه، ویرگول و علامت سوال از کلمات مجاور جدا شده و به توکن‌های مستقل تبدیل می‌شوند.

متن پس از این پردازش‌ها بر اساس فاصله‌ها شکسته شده و هر بخش به‌عنوان یک توکن جداگانه در لیست ذخیره می‌شود.

اگر قابلیت ترکیب افعال چندبخشی فعال باشد، افعالی مثل "خواهد رفت" یا "گفته شده است" با علامت _ به هم متصل می‌شوند و به‌عنوان یک توکن واحد برگردانده می‌شوند.

در نهایت، مخفف‌هایی که به‌طور موقت جایگزین شده بودند، به حالت اولیه خود بازمی‌گردند و لیستی از توکن‌های پردازش‌شده شامل کلمات، علائم نگارشی، ایموجی‌ها، و عناصر خاص آماده استفاده خواهد بود.

:Test 1

Original: من می‌خواهم که نمی‌روم به این مکان. همچنین کتاب‌هایم ترمیم شدند.

Processed Tokens: ['من', 'می', 'خواهم', 'که', 'نمی', 'روم', 'به', 'این', 'مکان', 'همچنین', 'کتاب', 'هایم', 'ترمیم', 'شدند']

:Test 2

Original: کاف و یای عربی به شکل فارسی تبدیل می‌شود. بسم الله الرحمن الرحيم.

Processed Tokens: ['کاف', 'و', 'یای', 'عربی', 'به', 'شکل', 'فارسی', 'تبدیل', 'می', 'شود', 'بسم', 'الله', 'الرحمن', 'الرحيم']

:Test 3

Original: إِنَّ اللَّهَ غَفُورٌ رَحِيمٌ.

Processed Tokens: ['إِنَّ', 'اللَّهُ', 'غَفُورٌ', 'رَحِيمٌ']

:Test 4

Original: سلام! این یک جمله است... یا شاید جمله‌ای دیگر؟!

Processed Tokens: ['سلام', 'این', 'یک', 'جمله', 'است', 'یا', 'شاید', 'جمله\200c\ای', 'دیگر']

:Test 5

Original: شماره تماس 1234567890 به شماره ۱۲۳۴۵۶۷۸۹۰ تبدیل شود.

Processed Tokens: ['شماره', 'تماس', '۱۲۳۴۵۶۷۸۹۰', 'به', 'شماره', '۱۲۳۴۵۶۷۸۹۰', 'تبدیل', 'شود']

:Test 6

Original: در سال‌های گذشته کتاب‌های زیادی خواندم. سال‌ها می‌گذرد و کتاب‌ها تکرار می‌شوند.

Processed Tokens: ['در', 'سال\200c\های', 'گذشته', 'کتاب\200c\های', 'زیادی', 'خواندم', 'سال\200c\u200c\ها', 'می\200c\گذرد', 'و', 'کتاب\200c\u200c\ها', 'تکرار', 'می\200c\شوند']

:Test 7

Original: بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ و .

Processed Tokens: ['بسم\200c\الله', 'الرحمن', 'الرحيم', 'و', 'بسم', 'الله', 'الرحمن', 'الرحيم']

:Test 8

Original: آیا؟! و ... را باید حذف کنیم؟

Processed Tokens: ['آیا', 'و', 'را', 'باید', 'حذف', 'کنیم']

:Test 9

Original: نمی‌شود "می‌خواهم" را به دو توکن تبدیل کرد.

1. **نرمال سازی و توکن سازی:** متن ورودی با استفاده از تابع `preprocess_single_text` نرمال و به توکن ها شکسته می شود.
2. **ریشه یابی:** با استفاده از تابع `stem_tokens`، تمامی توکن ها ریشه یابی می شوند.
3. **حذف توکن های پرتکرار:** اگر لیستی از توکن های پرتکرار (`top_k_tokens`) فراهم شود، این توکن ها از لیست نهایی حذف می شوند.

```
def simple_preprocess(content, top_k_tokens=None):  
    """  
    Simplified preprocessing pipeline for a single document:  
    - Normalize and tokenize the content.  
    - Stem tokens.  
    - Optionally remove tokens present in `top_k_tokens`.  
    """  
    # Step 1: Normalize and tokenize the content  
    tokens = preprocess_single_text(content)  
  
    # Step 2: Stem tokens  
    tokens = stem_tokens(tokens)  
  
    # Step 3: Remove top-K frequent tokens if provided  
    if top_k_tokens:  
        tokens = [token for token in tokens if token not in top_k_tokens]  
  
    return tokens
```

تابع `preprocess_all_docs`

این تابع همه اسناد را پیش پردازش می کند و مراحل زیر را انجام می دهد:

1. **پردازش تک تک اسناد:**
 - برای هر سند، متن نرمال شده، توکن سازی شده و توکن ها ریشه یابی می شوند.
 - لیست توکن های پردازش شده برای هر سند ذخیره می شود و همه توکن ها در لیستی ترکیبی جمع آوری می شوند.
2. **محاسبه پرتکرارترین توکن ها:**
 - با استفاده از `compute_top_k_frequent`، `k` توکن پرتکرار محاسبه می شود.
 - فقط خود توکن ها (بدون تعداد تکرار) برای حذف استفاده می شوند.
3. **حذف توکن های پرتکرار از هر سند:**
 - توکن های پرتکرار از لیست توکن های هر سند حذف می شوند.
 - محتوای هر سند با لیست فیلتر شده به روزرسانی می شود.

```
def preprocess_all_docs(docs, top_k=50):
    """
    Preprocesses all documents while returning the same structure as required:
    - Normalizes, tokenizes, and stems content for each document.
    - Computes top-K frequent tokens across all documents.
    - Removes top-K frequent tokens from each document's content.
    """
    combined_tokens = []
    all_tokens = {}

    for doc_id, doc_data in docs.items():
        content = doc_data['content']
        tokens = simple_preprocess(content)
        all_tokens[doc_id] = tokens
        combined_tokens.extend(tokens)

    top_k_tokens_with_counts = compute_top_k_frequent(combined_tokens, top_k)
    top_k_tokens = set(top_k_tokens_with_counts.keys()) # Extract tokens only

    for doc_id in all_tokens:
        filtered_tokens = [token for token in all_tokens[doc_id] if token not in top_k_tokens]
        docs[doc_id]['content'] = filtered_tokens

    return docs, top_k_tokens_with_counts
```

خروجی نهایی

این تابع:

1. اسناد پیش پردازش شده را برمی گرداند که در آن ها توکن های پرتکرار حذف شده اند.
2. لیستی از پرتکرارترین توکن ها به همراه تعداد تکرارشان را ارائه می دهد.

```
docs, contents, _ = load_docs()

pre_processed_docs, top_k_tokens_with_counts = preprocess_all_docs(docs, top_k=20)

print("\nTop-K Frequent Tokens:")
for token, count in top_k_tokens_with_counts.items():
    print(f"Token: {token}, Count: {count}")

print("\nProcessed Document Tokens:")
for doc_id, doc_data in list(pre_processed_docs.items())[:5]:
    print(f"Doc {doc_id} final tokens:", doc_data['content'])
```

Executed at 2025.01.11 02:31:50 in 1m 34s 625ms

Top-K Frequent Tokens:

Count: 234908	Token: و
Count: 165135	Token: در
Count: 136045	Token: به
Count: 92977	Token: از
Count: 83094	Token: این
Count: 75480	Token: که
Count: 69233	Token: با
Count: 68677	Token: را
Count: 48513	Token: اس
Count: 31829	Token: ...

Processed Document Tokens:

Doc 0 final tokens: [گزار, خبرگزار, فارس, کنفدراسیون, فوتبال, آسیا, AFC, نامه, رسم, فدراسیون, فوتبال, ایر, باشگاه, گیت, پسند, ز, قرعه\200c\کشی, ج, باشگاه, فوتسال, آسیا, رسما, اعل, اساس, ۲۵, فروردین\200c\ماه, ۱۴۰۱, مراسم, قرعه\200c\کشی, ج, باشگاه, فوتسال, آسیا, مالز, برگزار, می\200c\شود, باشگاه, گیت, پسند, یعنی, قهر, فوتسال, ایر, سال, ۱۴۰۰, مسابق, راه, پیدا, کرده-اس, پ, گیت, پسند, تجربه, ۳, دوره, حضور, ج, باشگاه, فوتسال, آسیا, داشته, هر, سه, دوره, فینال, مسابق, راه, پیدا, کرده, عنوان, قهرمان, دو, مق, دوم, بدس, آورده-اس, انت, بی, /]

Doc 1 final tokens: [گزار, خبرگزار, فارس, پیدا, حمید, مرجان, حاشیه, برای, مجامید, روز, جم, خبرنگار, رابط, روز

```
pre_processed_docs['1']
```

Executed at 2025.01.11 02:31:50 in 4ms

```
حضور تماشاگران در لیگ برتر فوتبال تابع نظر فدراسیون و سازمان لیگ است,
'content': [گزار,
'خبرگزار',
'فارس',
'سید',
'حمید',
'سجاد',
'حاشیه',
'مراس',
'گرامید',
'روز',
]
```

توضیح تابع `stemming hazm`:

۱. متد `__init__`

در این متد، یک لیست از پسوندهای رایج فارسی مانند "ات"، "ان"، "ترین"، "ها"، "یی"، "م" و غیره در متغیر `self.ends` ذخیره می‌شود. این پسوندها برای حذف از انتهای کلمات در متد `stem` استفاده می‌شوند.

۲. متد `stem`

این متد برای یافتن ریشه کلمه طراحی شده است. ابتدا بررسی می‌کند که آیا کلمه به یکی از پسوندهای موجود در لیست `self.ends` ختم می‌شود یا خیر. اگر پسوندی یافت شود، آن را از انتهای کلمه حذف می‌کند. سپس، اگر کلمه به "ه" ختم شود، آن را به "ه" تغییر می‌دهد. در نهایت، ریشه کلمه بازگردانده می‌شود.

```
class Stemmer(StemmerI):
    """این کلاس شامل توابعی برای ریشه‌یابی کلمات است."""

    def __init__(self: "Stemmer") -> None:
        self.ends = [
            "ات",
            "ان",
            "ترین",
            "تر",
            "م",
            "ت",
            "ش",
            "یی",
            "ی",
            "ها",
            "ا",
            "ا",
            "ZWNJ",
        ]

    def stem(self: "Stemmer", word: str) -> str:
        """ریشه کلمه را پیدا می‌کند.
```

۳.۱ ساخت شاخص مکانی:

```
def _add_new_token_posting(token_dict, token, doc_id, position):
    """
    Initialize a new posting entry for a token that has not yet been seen.
    """
    token_dict[token] = {
        'frequency': 1,
        'docs': {
            doc_id: {
                'positions': [position],
                'number_of_token': 1
            }
        }
    }
```

این تابع برای اضافه کردن یک توکن جدید به دیکشنری پستینگ‌ها استفاده می‌شود.

- اگر توکن جدید باشد:
 - تعداد تکرار آن (frequency) برابر 1 تنظیم می‌شود.
 - در دیکشنری docs، شناسه سند (doc_id) ذخیره می‌شود.
 - موقعیت توکن در متن (لیستی از موقعیت‌ها) به همراه تعداد تکرار در همان سند نگهداری می‌شود.

```
def _update_existing_token_posting(token_dict, token, doc_id, position):
    """
    Update an existing token entry with a new position in an existing or new document.
    """
    token_dict[token]['frequency'] += 1

    if doc_id in token_dict[token]['docs']:
        token_dict[token]['docs'][doc_id]['positions'].append(position)
        token_dict[token]['docs'][doc_id]['number_of_token'] += 1
    else:
        token_dict[token]['docs'][doc_id] = {
            'positions': [position],
            'number_of_token': 1
        }
```

.تابع update_existing_token_posting

این تابع برای به‌روزرسانی توکنی که قبلاً وجود دارد استفاده می‌شود.

- تعداد تکرار توکن (frequency) در کل دیکشنری افزایش می‌یابد.
- اگر توکن در همان سند (doc_id) باشد:
 - موقعیت جدید آن به لیست موقعیت‌ها اضافه می‌شود.

- تعداد تکرار توکن در همان سند به روزرسانی می‌شود.
- اگر توکن در سند جدید باشد:
 - ورودی جدیدی برای آن سند اضافه می‌شود که شامل موقعیت و تعداد تکرار است.

```
def _build_postings_dict(Docs):
    """
    Build a postings dictionary from the input Docs, collecting token frequency
    and positions per document.
    """
    token_dict = {}
    for doc_id, doc_content in Docs.items():
        for position, token in enumerate(doc_content['content']):
            if token in token_dict:
                _update_existing_token_posting(token_dict, token, doc_id, position)
            else:
                _add_new_token_posting(token_dict, token, doc_id, position)
    return token_dict
```

3. تابع `build_postings_dict`

این تابع یک دیکشنری پستینگ برای تمامی اسناد می‌سازد.

- اسناد (Docs) به صورت `doc_id` و محتوای آن بررسی می‌شوند.
- تمامی توکن‌های هر سند پیمایش شده و:
 - اگر توکن جدید باشد، به کمک تابع `add_new_token_posting` اضافه می‌شود.
 - اگر توکن قبلاً وجود داشته باشد، به کمک تابع `update_existing_token_posting` به روزرسانی می‌شود.

خروجی:

یک دیکشنری که شامل توکن‌ها، تعداد تکرار آن‌ها و اطلاعات مربوط به هر سند است.

```
def _calculate_tf_idf(token_dict, total_docs):
    """
    Calculate TF-IDF for each token in each document.
    """
    for term, term_data in token_dict.items():
        term_docs = term_data['docs']
        n_t = len(term_docs) # Number of documents containing this term

        for doc_id, doc_info in term_docs.items():
            tf = doc_info['number_of_token']
            # TF-IDF = log10(N / n_t) * (1 + log10(tf))
            tf_idf_value = (np.log10(total_docs / n_t)) * (1 + np.log10(tf))
            doc_info['tf_idf'] = tf_idf_value
```

این تابع مقدار **TF-IDF** را برای هر توکن در هر سند محاسبه می‌کند:

- برای هر توکن:
 - تعداد اسنادی که شامل این توکن هستند (n_t) محاسبه می‌شود.
 - سپس برای هر سند (doc_id) شامل توکن:
 - مقدار TF (تعداد تکرار توکن در سند) محاسبه می‌شود.
 - مقدار TF-IDF با فرمول زیر محاسبه شده و ذخیره می‌شود:

```
# TF-IDF = log10(N / n_t) * (1 + log10(tf))
```

تابع `build_champions_list_and_docs_vectors_`

این تابع دو وظیفه دارد:

1. ساخت لیست قهرمانان (Champions List):

- توکن‌ها براساس تعداد تکرارشان در هر سند مرتب می‌شوند.
- اگر تعداد اسناد بیش از مقدار `champ_len` باشد، بالاترین `champ_len` سند در لیست قهرمانان ذخیره می‌شود.
- این لیست شامل اطلاعات `number_of_token` و `TF-IDF` برای هر سند است.

```

def _build_champions_list_and_docs_vectors(token_dict, champ_len):
    docs_vectors = {}

    for term, term_data in token_dict.items():
        term_docs = term_data['docs']

        # Sort documents by number_of_token (descending) for the champion list
        sorted_term_docs = sorted(
            term_docs,
            key=lambda d: term_docs[d]['number_of_token'],
            reverse=True
        )

        # Build champion list (take top `champ_len` if needed)
        champions_list = {}
        for doc_id in sorted_term_docs:
            champions_list[doc_id] = {
                'number_of_token': term_docs[doc_id]['number_of_token'],
                'tf_idf': term_docs[doc_id]['tf_idf']
            }

        if champ_len < len(term_docs):
            champions_list = dict(list(champions_list.items())[:champ_len])

        token_dict[term]['champions_list'] = champions_list

```

2. ساخت docs_vectors:

○ بردار توکن‌های هر سند شامل اطلاعات TF-IDF و TF در قالب دیکشنری ذخیره می‌شود.

```

# Populate docs_vectors
for doc_id, doc_info in term_docs.items():
    if doc_id not in docs_vectors:
        docs_vectors[doc_id] = {}
    docs_vectors[doc_id][term] = {
        'tf_idf': doc_info['tf_idf'],
        'tf': doc_info['number_of_token']
    }

return docs_vectors

```

تابع اصلی Postings_List

این تابع تمامی مراحل بالا را مدیریت می‌کند:

1. ایجاد لیست پستینگ:

- ابتدا دیکشنری پستینگ با استفاده از تابع `build_postings_dict_` ساخته می‌شود.

2. محاسبه TF-IDF:

- برای تمامی توکن‌ها و اسناد مقدار TF-IDF محاسبه می‌شود.

3. ساخت لیست قهرمانان و بردارهای اسناد:

- با استفاده از تابع `build_champions_list_and_docs_vectors_`، لیست قهرمانان و بردارهای سند ایجاد می‌شود.

خروجی:

- دیکشنری پستینگ با اطلاعات کامل هر توکن.
- دیکشنری `docs_vectors` که بردارهای توکن هر سند را شامل می‌شود.

```
def Postings_List(Docs, champ_len):  
    """  
    Orchestrates the creation of the postings list (token_dict) and docs_vectors.  
    Steps:  
    1) Build an initial postings dictionary with token frequencies/positions.  
    2) Compute TF-IDF for each token in each document.  
    3) Build the champions lists and docs_vectors.  
    """  
    # Step 1: Build the core postings dictionary  
    token_dict = _build_postings_dict(Docs)  
  
    # Step 2: Calculate TF-IDF  
    total_docs = len(Docs)  
    _calculate_tf_idf(token_dict, total_docs)  
  
    # Step 3: Create champions lists and docs_vectors  
    docs_vectors = _build_champions_list_and_docs_vectors(token_dict, champ_len)  
  
    return token_dict, docs_vectors
```

```
dictionary, docs_vectors = Postings_List(pre_processed_docs, 20)
```

Executed at 2025.01.11 02:32:03 in 12s 559ms

```
docs_vectors['8332']
```

Executed at 2025.01.11 02:32:03 in 2ms

```
,{'tf_idf': 3.003141219509179, 'tf': 2}: 'برف'
,{'tf_idf': 0.6759912344460222, 'tf': 1}: 'پارلمان'
,{'tf_idf': 1.9860604755388058, 'tf': 1}: 'دام'
,{'tf_idf': 2.5549521036141134, 'tf': 1}: 'شر'
,{'tf_idf': 3.38746101632035, 'tf': 1}: 'غلام\U200cرضا'
,{'tf_idf': 1.8216131976468322, 'tf': 1}: 'شهرساز'
,{'tf_idf': 1.928068528561119, 'tf': 1}: 'انتخابیه'
,{'tf_idf': 2.3788608445584325, 'tf': 1}: 'دامدار'
,{'tf_idf': 2.8559820992780947, 'tf': 1}: 'گلگاه'
,{'tf_idf': 3.132188511217044, 'tf': 1}: 'علوفه'
,{'tf_idf': 6.036154116557285, 'tf': 3}: 'بهشهر'
,{'tf_idf': 4.086431020656368, 'tf': 1}: 'نکاء'
```

```
dictionary['فهرست'] ['champions_list']
```

Executed at 2025.01.11 02:32:03 in 4ms

```
'821': {'number_of_token': 25, 'tf_idf': 0.005047725200337222},
'2388': {'number_of_token': 24, 'tf_idf': 0.0050104056913280555},
'7744': {'number_of_token': 24, 'tf_idf': 0.0050104056913280555},
'525': {'number_of_token': 19, 'tf_idf': 0.00479683475620343},
'8680': {'number_of_token': 16, 'tf_idf': 0.004639729076116687},
'10183': {'number_of_token': 16, 'tf_idf': 0.004639729076116687},
'10025': {'number_of_token': 15, 'tf_idf': 0.00458072789429455},
'2831': {'number_of_token': 14, 'tf_idf': 0.00451765454084172},
'8897': {'number_of_token': 14, 'tf_idf': 0.00451765454084172},
'9687': {'number_of_token': 14, 'tf_idf': 0.00451765454084172},
'1684': {'number_of_token': 13, 'tf_idf': 0.004449904957144393},
'2098': {'number_of_token': 13, 'tf_idf': 0.004449904957144393}
```

```
dictionary['کیمسیون']
```

Executed at 2025.01.11 03:50:50 in 70ms

```
{'frequency': 6298,
 'docs': {'13': {'positions': [79, 91],
                  'number_of_token': 2,
                  'tf_idf': 1.084825715800724},
          '89': {'positions': [301],
                  'number_of_token': 1,
                  'tf_idf': 0.8338206800889957},
          '100': {'positions': [149],
                  'number_of_token': 1,
                  'tf_idf': 0.8338206800889957},
          '106': {'positions': [34],
                  'number_of_token': 1
```

```
'champions_list': {'11882': {'number_of_token': 55,
    'tf_idf': 2.2849710814446},
    '7506': {'number_of_token': 48, 'tf_idf': 2.2356743920311724},
    '9970': {'number_of_token': 46, 'tf_idf': 2.2202625461250296},
    '7401': {'number_of_token': 42, 'tf_idf': 2.187319507362553},
    '10042': {'number_of_token': 35, 'tf_idf': 2.121296546932828},
    '11449': {'number_of_token': 29, 'tf_idf': 2.053198373257889},
    '9500': {'number_of_token': 28, 'tf_idf': 2.040490973979017},
    '11143': {'number_of_token': 28, 'tf_idf': 2.040490973979017},
    '11163': {'number_of_token': 26, 'tf_idf': 2.0136547194018264},
```

ذخیره دیکشنری:

این کد یک دیکشنری را به دو قسمت مساوی تقسیم می‌کند. ابتدا تعداد کل کلیدهای دیکشنری محاسبه می‌شود و سپس کلیدها به دو بخش تقسیم می‌شوند: نیمه اول شامل اولین نصف کلیدها و نیمه دوم شامل بقیه کلیدها است. هر بخش به صورت جداگانه در فایل‌های first_half.json و second_half.json ذخیره می‌شود. داده‌ها در این فایل‌ها با فرمت JSON و به صورت خوانا ذخیره می‌شوند.

```
half_length = len(dictionary) // 2
first_half = {key: dictionary[key] for key in list(dictionary.keys())[:half_length]}
second_half = {key: dictionary[key] for key in list(dictionary.keys())[half_length:]}

with open('first_half.json', "w", encoding="utf-8") as first_file:
    json.dump(first_half, first_file, indent=4)

with open('second_half.json', "w", encoding="utf-8") as second_file:
    json.dump(second_half, second_file, indent=4)
```

۲. پاسخ‌دهی به پرسمان در فضای برداری

تابع vector_length

این تابع طول بردار وزن‌های TF-IDF یک سند را محاسبه می‌کند. این طول در محاسبه شباهت کوساین استفاده می‌شود.

- ورودی: دیکشنری‌ای که شامل وزن‌های TF-IDF کلمات یک سند است.
- محاسبات:
 - مربع تمام وزن‌های TF-IDF محاسبه و جمع زده می‌شود.
 - جذر جمع این مقادیر به عنوان طول بردار برگردانده می‌شود.
- خروجی: طول بردار.

```
def vector_length(vector_dict):
    length = math.sqrt(sum(tf_idf_value['tf_idf'] ** 2 for tf_idf_value in vector_dict.values()))
    return length
```

تابع get_query_tokens

عبارت جستجو را نرمال سازی و توکنایز می کند.

```
def get_query_tokens(query):
    """
    Preprocesses the query into tokens.
    Returns a list of tokens.
    """
    return simple_preprocess(query)
```

تابع get_query_tokens_count

تعداد دفعات تکرار هر کلمه در عبارت جستجو را محاسبه می کند.

```
def get_query_tokens_count(query_tokens):
    """
    Counts how many times each token appears in the query.
    Returns a dictionary of token -> frequency.
    """
    return dict(collections.Counter(query_tokens))
```

وزن دهی به کلمات عبارت جستجو

تابع compute_query_weight

این تابع وزن TF-IDF هر کلمه از عبارت جستجو را محاسبه می کند.

• ورودی ها:

- term: کلمه ای که باید وزن آن محاسبه شود.
- query_tokens_count: تعداد دفعات تکرار کلمه در عبارت جستجو.
- dictionary: دیکشنری حاوی اطلاعات مربوط به کلمات و اسناد.
- champion_list: مشخص می کند که آیا از لیست منتخب (champion list) استفاده شود یا کل اسناد.
- total_number_of_docs: تعداد کل اسناد.

• منطق:

- ابتدا بررسی می کند که آیا کلمه در دیکشنری وجود دارد یا نه. اگر کلمه موجود نباشد، مقدار 0 و یک دیکشنری خالی بازگردانده می شود.
- وزن TF-IDF برای کلمه محاسبه می شود.
- لیست اسناد مرتبط (لیست کامل یا منتخب) بازیابی می شود.

- خروجی: وزن TF-IDF کلمه و لیست اسناد حاوی آن

```
def compute_query_weight(term, query_tokens_count, dictionary, champion_list, total_number_of_docs):
    """
    Computes the TF-IDF weight of a query term.
    Returns w_tq (the query-term weight) and the term_docs to iterate over.
    """
    if term not in dictionary:
        return 0, {}

    # Retrieve the appropriate document list (champion list or full list)
    term_docs = (dictionary[term]['champions_list']
                 if champion_list
                 else dictionary[term]['docs'])

    w_tq = calculate_tf_idf(query_tokens_count[term],
                           total_number_of_docs,
                           len(term_docs))
    return w_tq, term_docs
```

به روزرسانی امتیازات اسناد

محاسبه طول بردار سند

تابع `vector_length`

طول بردار وزن های TF-IDF یک سند را محاسبه می کند که برای نرمال سازی شباهت کسینوسی استفاده می شود.

- منطق:
 - مربع تمام مقادیر TF-IDF جمع زده می شود.
 - جذر مجموع این مقادیر طول بردار را می دهد.
- خروجی: طول بردار.

```
def vector_length(vector_dict):
    length = math.sqrt(sum(tf_idf_value['tf_idf'] ** 2 for tf_idf_value in vector_dict.values()))
    return length
```

پیش پردازش عبارت جستجو

تابع `get_query_tokens`

عبارت جستجو را نرمال سازی و توکنایز می کند و کلمات مجزا را استخراج می کند.

- خروجی: لیستی از کلمات جستجو.

تابع `get_query_tokens_count`

تعداد تکرار هر کلمه در عبارت جستجو را محاسبه می‌کند.

- **خروجی:** دیکشنری‌ای شامل هر کلمه و تعداد تکرار آن.

```
def get_query_tokens(query):  
    """  
    Preprocesses the query into tokens.  
    Returns a list of tokens.  
    """  
    return simple_preprocess(query)
```

وزن دهی به کلمات جستجو

تابع `compute_query_weight`

وزن TF-IDF هر کلمه از عبارت جستجو را محاسبه می‌کند و لیست اسناد حاوی آن را برمی‌گرداند.

- **منطق:**
 - اگر کلمه در دیکشنری وجود نداشته باشد، مقدار 0 و یک دیکشنری خالی بازگردانده می‌شود.
 - لیست اسناد مرتبط با کلمه از دیکشنری (لیست کامل یا لیست منتخب) دریافت می‌شود.
 - وزن TF-IDF کلمه محاسبه می‌شود.
- **خروجی:** وزن TF-IDF کلمه و لیست اسناد مرتبط.

```
def compute_query_weight(term, query_tokens_count, dictionary, champion_list, total_number_of_docs):  
    """  
    Computes the TF-IDF weight of a query term.  
    Returns w_tq (the query-term weight) and the term_docs to iterate over.  
    """  
    if term not in dictionary:  
        return 0, {}  
  
    # Retrieve the appropriate document list (champion list or full list)  
    term_docs = (dictionary[term]['champions_list']  
                 if champion_list  
                 else dictionary[term]['docs'])  
  
    w_tq = calculate_tf_idf(query_tokens_count[term],  
                           total_number_of_docs,  
                           len(term_docs))  
    return w_tq, term_docs
```

به‌روزرسانی امتیازات اسناد

تابع `update_doc_scores`

امتیازات کسینوسی را برای اسنادی که کلمه مورد نظر را دارند، به روزرسانی می‌کند.

- **منطق:**

- وزن TF-IDF کلمه در سند (w_{td}) در وزن کلمه در عبارت جستجو (w_{tq}) ضرب می‌شود.
- اگر سند قبلاً امتیاز داشته باشد، به آن اضافه می‌شود؛ در غیر این صورت، مقدار اولیه تنظیم می‌شود.

- **خروجی:** به روزرسانی امتیازات در دیکشنری `cosine_scores`.

```
def update_doc_scores(term_docs, w_tq, cosine_scores):  
    """  
    Updates the cosine scores for each document that contains the term.  
    """  
    for doc in term_docs:  
        w_td = term_docs[doc]['tf_idf']  
        doc_id = int(doc)  
  
        # Update cosines similarity  
        if doc_id in cosine_scores:  
            cosine_scores[doc_id] += w_td * w_tq  
        else:  
            cosine_scores[doc_id] = w_td * w_tq
```

محاسبه نهایی امتیازات کسینوسی

تابع `finalize_cosine_scores`

امتیاز نهایی شباهت کسینوسی را برای هر سند محاسبه می‌کند.

- **منطق:**

- امتیاز هر سند بر طول بردار آن تقسیم می‌شود تا نرمال‌سازی انجام شود.

- **خروجی:** امتیازات نهایی نرمال‌شده.

```
def finalize_cosine_scores(cosine_scores):  
    """  
    Divides each document's cosine score by its vector length.  
    """  
    for doc_number in cosine_scores:  
        cosine_scores[doc_number] /= vector_length(docs_vectors[str(doc_number)])
```

مرتب‌سازی امتیازات

تابع `sort_scores`

امتیازات اسناد را به ترتیب نزولی مرتب می‌کند.

- **خروجی:** لیستی از جفت‌های (doc_id, score) مرتب‌شده بر اساس امتیاز.

```
def sort_scores(scores_dict):  
    """  
    Sorts scores in descending order by value.  
    Returns a list of (doc_id, score) tuples.  
    """  
    return sorted(scores_dict.items(), key=lambda x: x[1], reverse=True)
```

تابع query_scoring

این تابع تمام مراحل بالا را ترکیب می‌کند و امتیازات نهایی را برای عبارت جستجو محاسبه می‌کند.

- **مراحل:**

1. عبارت جستجو پیش‌پردازش می‌شود و تعداد تکرار کلمات آن محاسبه می‌شود.
2. برای هر کلمه در عبارت جستجو:
 - وزن TF-IDF محاسبه می‌شود.
 - امتیازات کسینوسی برای اسناد مرتبط به‌روزرسانی می‌شود.
3. امتیازات کسینوسی نهایی نرمال‌سازی می‌شوند.
4. اسناد بر اساس امتیاز کسینوسی مرتب می‌شوند.
5. تاپ k سند برتر برگردانده می‌شود.

- **خروجی:** لیستی از k سند برتر بر اساس شباهت کسینوسی

```
def query_scoring(query, total_number_of_docs, dictionary, k, champion_list=False):
    # Initialize scores
    cosine_scores = {}

    # Preprocess query
    query_tokens = get_query_tokens(query)
    query_tokens_count = get_query_tokens_count(query_tokens)
    query_terms_num = sum(query_tokens_count.values())

    print(query_tokens_count)

    # Compute and update scores for each term in the query
    for term in query_tokens_count:
        w_tq, term_docs = compute_query_weight(term,
                                                query_tokens_count,
                                                dictionary,
                                                champion_list,
                                                total_number_of_docs)

        if w_tq != 0:
            update_doc_scores(term_docs, w_tq, cosine_scores)

    # Finalize scores
    finalize_cosine_scores(cosine_scores)

    # Sort and retrieve top k results
    sorted_doc_cosine = sort_scores(cosine_scores)

    return sorted_doc_cosine[:k]
```

```
1 r1, r2 = query_search('کریسمس', result_numbers = 5, champion_list = True)
Executed at 2025.01.11 03:51:12 in 5ms
```

```
{1: 'کریسمس'}
=== Cosine Scores ===
----- Results -----
Rank: 1 | ID: 5933
Title : ستاره اسپانیایی؛ هدیه کریسمس گواردیولا به ژاوی+عکس
URL : https://www.farsnews.ir/news/14001007000739
-----
Rank: 2 | ID: 6117
Title : کیروش «دیکتاتور» لقب گرفت/ اختلاف مرد پرتغالی با مصری‌ها به خاطر کریسمس+عکس
URL : https://www.farsnews.ir/news/14001007000739
-----
Rank: 3 | ID: 6128
```

الف) ساده:کلمه اخبار

{اخبار: 1}

=== Cosine Scores ===

----- Results -----

Rank: 1 | ID: 4082

Title : انتصاب های جدید در باشگاه پرسپولیس رد شد

URL :

https://www.farsnews.ir/news/14001101000140/انتصاب-های-جدید-در-باشگاه-پرسپولیس-رد-

شد

Rank: 2 | ID: 3781

Title : گفت‌وگوی فارس با مسئول فدراسیون عراق درباره جزئیات کرونایی‌ها قبل از بازی با ایران

URL :

<https://www.farsnews.ir/news/14001104000419>/گفت‌وگوی-فارس-با-مسئول-فدراسیون-عراق-در

باره-جزئیات-کرونایی‌ها-قبل-از

Rank: 3 | ID: 5749

Title : پیش بینی مفسر مشهور فوتبال روسیه از آینده آزمون+عکس

URL :

<https://www.farsnews.ir/news/14001010000074>/پیش-بینی-مفسر-مشهور-فوتبال-روسیه-از-آیند

ه-آزمون-عکس

Rank: 4 | ID: 10436

Title : مقتدایی: درباره خلل یا آسیب در روند نگهداری برخی از اسناد سازمان اوقاف شفاف سازی شود

URL :

<https://www.farsnews.ir/news/14000909000220>/مقتدایی-درباره-خلل-یا-آسیب-در-روند-نگهداری

-برخی-از-اسناد-سازمان-اوقاف

Rank: 5 | ID: 5330

Title : از تکذیب توافق آزمون با لیون تا خوشحالی روس‌ها+عکس

URL :

<https://www.farsnews.ir/news/14001016000067>/از-تکذیب-توافق-آزمون-با-لیون-تا-خوشحالی-رو

س‌ها-عکس



انتصاب های جدید در باشگاه پرسپولیس رد شد

باشگاه پرسپولیس با انتشار بیانیه‌ای انتصاب‌های جدید در این باشگاه را تکذیب کرد.

به گزارش خبرگزاری فارس و به نقل از سایت باشگاه پرسپولیس، در روزهای اخیر موج شایعات و اخبار غیر موثق در ارتباط با انتصاب‌ها در حوزه مدیریت، معاونت و مش منتشر شده است. باشگاه پرسپولیس، ضمن رد این اخبار، تاکید می‌کند؛ برنامه‌ریزی با توجه به اولویت‌هایی که وجود دارد، حل و فصل مسائل باشگاه و تیم است. بر ه اساس هیچ گونه تغییر و انتصابی تاکنون صورت نگرفته است. اخبار رسمی باشگاه پرسپولیس، از طریق رسانه‌های رسمی باشگاه اطلاع‌رسانی خواهد شد. انتهای پیام/

گفت‌وگوی فارس با مسئول فدراسیون عراق درباره جزئیات کرونایی‌ها قبل از بازی با ایران

مدیر روابط عمومی فدراسیون فوتبال عراق توضیحاتی را درباره جزئیات ابتلای بازیکنان و اعضای کاروان این تیم در تهران ارائه کرد.

به گزارش خبرگزاری فارس، امروز رسانه های عراق **اخباری** متفاوت از کرونایی های تیم ملی این کشور منتشر کردند. برخی از رسانه از 6 کرونایی، برخی دیگر از 5 کرونایی و **اخباری** از ابتلای 3 بازیکن مطرح شد. به همین خاطر خبرنگار خبرگزاری فارس برای تایید رسمی این **اخبار** به سراغ محمد عماد مدیر روابط عمومی فدراسیون فوتبال عراق رفت از واقعیت این **اخبار** مطلع شود. مدیر روابط عمومی فدراسیون فوتبال عراق در این زمینه به خبرنگار فارس، گفت: **بله اخبار** مطرح شده در زمینه ابتلای چند بازیکن و عضو تیم ملی صحت دارد. ما روز گذشته در تهران تست کرونا دادیم و امروز پاسخ تست ها آمد. وی درباره تعداد دقیق کرونایی‌های تیم ملی عراق در بدو ورود به تهران، گفت: 5 تست کرونا اعضای کاروان در تهران مثبت شد. در میان این کرونایی‌ها تست دو بازیکن به نام های یاسر قاسم و مناف یونس مثبت شد. همچنین تست کرونای احمد جاسم مربی دروازه بان های عراق، غیث مهنا مدیر اداری و یاسین خبیر پزشک تیم مثبت شده است. مدیر روابط عمومی فدراسیون عراق

در این ها خبر تعداد تکرارهای کلمه اخبار زیاد است چون از لیست قهرمان استفاده شده که در آن داکيومنت ها با بیشترین tf قرار دارندو اگر همین سرچ بدون لیست قهرمان انجام میشد احتمالا داکيومنت ها کوتاه‌تر و با تعداد تکرار کمتر کلمه اخبار بود.
زیاد بودن تکرار کلمه اخبار نشانه از درست بودن کوئری است.

(ب) ساده و مرکب: باشگاه والیبال

{باشگاه: 1, 'فوتبال': 1}

=== Cosine Scores ===

----- Results -----

Rank: 1 | ID: 4238

Title : منافی: امیدواریم درویش بتواند از پس مشکلات پرسپولیس بر بیاید/ هیات مدیره جای بگو و بخند

نیست

URL :

https://www.farsnews.ir/news/14001028000893/منافی-امیدواریم-درویش-تواند-از-پس-مشکلات-پرسپولیس-بر-بیاید-هیات

ت-پرسپولیس-بر-بیاید-هیات

Rank: 2 | ID: 2642

Title : شهرداری ارومیه به دنبال دور زدن فدراسیون جهانی والیبال/ نه پول هست، نه وکیل

URL :

https://www.farsnews.ir/news/14001118000468/شهرداری-ارومیه-به-دنبال-دور-زدن-فدراسیون-جهانی-والیبال-نه-پول-هست-نه

جهانی-والیبال-نه-پول-هست-نه

Rank: 3 | ID: 5403

Title : اولیایی: امیدواریم AFC قول وزیر ورزش را بپذیرد/ واگذاری سهام استقلال و پرسپولیس به مردم

سراب است

URL :

<https://www.farsnews.ir/news/14001014001185>/اولیایی-امیدواریم-AFC-قول-وزیر-ورزش-را-پیذ
یرد-واگذاری-سهام-استقلال-و

Rank: 4 | ID: 2391

Title : عطابخش: به جز ۲-۳ بازیکن سایر نفرات مان می‌توانند در لیگ امید بازی کنند/ پیکان در بودجه
مضيقه دارد

URL :

<https://www.farsnews.ir/news/14001122000769>/عطابخش-به-جز-۲-۳-بازیکن-سایر-نفرات مان-
می‌توانند-در-لیگ-امید-بازی

Rank: 5 | ID: 534

Title : پولادگر: خیلی از باشگاه‌های ما فقط تابلوی خصوصی دارند/ ساختمان‌ها به استقلال و پرسپولیس
هبه شده

URL :

<https://www.farsnews.ir/news/14001216001105>/پولادگر--خیلی-از-باشگاه‌های-ما-فقط-تابلوی-خ
صوصی-دارند-ساختمان‌ها-به



تکلیف میزبان مسابقات والیبال جام باشگاه‌های آسیا مشخص شد

با انصراف باشگاه شهرداری ارومیه از میزبانی رقابت‌های والیبال قهرمانی باشگاه‌های مردان آسیا، این مسابقات در سالن فدراسیون در تهران برگزار خواهد شد.

به گزارش خبرگزاری فارس، رقابت‌های والیبال قهرمانی باشگاه‌های مردان آسیا از 25 اردیبهشت تا یکم خردادماه سال 1401 با حضور 11 تیم از 10 کشور قاره کهن برگزار خواهد شد. پیش از این قرار بود این رقابت‌ها به میزبانی شهرداری ارومیه برگزار شود که با توجه به محرومیت‌های مالی این باشگاه ازسوی فدراسیون جهانی والیبال و با وجود دادن چندین مهلت و فرصت فدراسیون والیبال به شهرداری ارومیه، این باشگاه نتوانست رضایت شاکیان خود از FIVB را جلب نماید. مدیرکل ورزش و جوانان استان آذربایجان غربی و مسئولین باشگاه شهرداری ارومیه با وجود تلاش‌های و مذاکرات فراوان اعلام کردند که نمی‌توانند این مشکلات را برطرف کنند و از حضور و میزبانی رقابت‌های والیبال قهرمانی باشگاه‌های مردان آسیا انصراف دادند. این در حالی بود که مسئولین فدراسیون والیبال در نظر داشتند آذربایتری را به عنوان میزبان این رقابت‌ها اعلام کنند که مسئولین این تیم هم اعلام کردند نمی‌توانند میزبانی این رقابت‌ها را بر عهده گیرند. بر این اساس و با توجه به ترافیک رویدادهای بین‌المللی والیبال ایران در سال آینده، فدراسیون والیبال تصمیم گرفت که رقابت‌های والیبال قهرمانی باشگاه‌های مردان آسیا 2022 را در سالن والیبال مجموعه ورزشی آزادی تهران برگزار کند و تهران را به عنوان شهر میزبان این رقابت‌ها به کنفدراسیون والیبال آسیا معرفی کرد.

بدین ترتیب علاوه بر قهرمان لیگ برتر والیبال سال 1400 ایران، یک تیم دیگر از ایران به‌عنوان میزبان این رقابت‌ها در مسابقات باشگاه‌های مردان آسیا حضور خواهد یافت که فدراسیون والیبال در روزهای آینده پس از پایان لیگ برتر ۱۴۰۰ در این خصوص تصمیم‌گیری خواهد کرد انتهای پیام/.

باشگاه

1/10

ورزشی

@Sports

باشگاه‌های آسیا مشخص شد

باشگاه

باشگاه

باشگاه‌های مردان

باشگاه

باشگاه

باشگاه‌های مردان

با انصراف **باشگاه** شهرداری ارومیه از میزبانی رقابت‌های والیبال قهرمانی **باشگاه‌های** مردان آسیا، این مسابقات در سالن فدراسیون در تهران برگزار خواهد شد.

به گزارش خبرگزاری فارس، رقابت‌های والیبال قهرمانی **باشگاه** های مردان آسیا از 25 اردیبهشت تا یکم خردادماه سال 1401 با حضور 11 تیم از 10 کشور قاره کهن برگزار خواهد شد. پیش از این قرار بود این رقابت‌ها به میزبانی شهرداری ارومیه برگزار شود که با توجه به محرومیت‌های مالی این **باشگاه** ازسوی فدراسیون جهانی والیبال و با وجود دادن چندین مهلت و فرصت فدراسیون والیبال به شهرداری ارومیه، این **باشگاه** نتوانست رضایت شاکیان خود از FIVB را جلب نماید. مدیرکل ورزش و جوانان استان آذربایجان غربی و مسئولین **باشگاه** شهرداری ارومیه با وجود تلاش‌های و مذاکرات فراوان اعلام کردند که نمی‌توانند این مشکلات را برطرف کنند و از حضور و میزبانی رقابت‌های والیبال قهرمانی **باشگاه‌های** مردان آسیا انصراف دادند. این در حالی بود که مسئولین فدراسیون والیبال در نظر داشتند آذربایتری را به عنوان میزبان این رقابت‌ها اعلام کنند که مسئولین این تیم هم اعلام کردند نمی‌توانند میزبانی این رقابت‌ها را بر عهده گیرند. بر این اساس و با توجه به ترافیک رویدادهای بین‌المللی والیبال ایران در سال آینده، فدراسیون والیبال تصمیم گرفت که رقابت‌های والیبال قهرمانی **باشگاه‌های** مردان آسیا 2022 را در سالن والیبال مجموعه ورزشی آزادی تهران برگزار کند و تهران را به عنوان شهر میزبان این رقابت‌ها به کنفدراسیون والیبال آسیا معرفی کرد.

بدین ترتیب علاوه بر قهرمان لیگ برتر والیبال سال 1400 ایران، یک تیم دیگر از ایران به‌عنوان میزبان این رقابت‌ها در مسابقات **باشگاه‌های** مردان آسیا حضور خواهد یافت که فدراسیون والیبال در روزهای آینده پس از پایان لیگ برتر ۱۴۰۰ در این خصوص تصمیم‌گیری خواهد کرد انتهای پیام/.

این سند با اینکه دو کلمه کوثری را دارد ولی چون هر دو به نسبت ساده و رایج هستند و در ترکیب با کلمات دیگر هستند. برای همین در سند می توان دید که این دو کلمه در کنار هم نیامده اند و ممکن است اسناد غیرمرتبط برگردانده شود.

پ) دشوار تک کلمه ای: کریسمس
{'کریسمس': 1}

=== Cosine Scores ===

----- Results -----

Rank: 1 | ID: 5933

Title : ستاره اسپانیایی؛ هدیه کریسمس گواردیولا به ژاوی+عکس

URL :

https://www.farsnews.ir/news/14001007000739/ستاره-اسپانیایی-هدیه-کریسمس-گواردیولا-به-ژ

او-عکس

Rank: 2 | ID: 6117

Title : کی‌روش «دیکتاتور» لقب گرفت/اختلاف مرد پرتغالی با مصری‌ها به خاطر کریسمس+عکس

URL :

https://www.farsnews.ir/news/14001005000165/کی‌روش-دیکتاتور-لقب-گرفت-اختلاف-مرد-پرتغا

لی-با-مصری‌ها-به-خاطر-کریسمس

Rank: 3 | ID: 6120

Title : کشتار در ورزشگاه فوتبال در آستانه سال جدید

URL :

https://www.farsnews.ir/news/14001005000143/کشتار-در-ورزشگاه-فوتبال-در-آستانه-سال-جدید

Rank: 4 | ID: 5926

Title : مهاجم خارجی مس رفسنجان، 4 کودک را به محل تحصیل برگرداند +عکس

URL :

https://www.farsnews.ir/news/14001007000809/مهاجم-خارجی-مس-رفسنجان-4-کودک-را-به-محل-تحصیل-برگرداند-عکس

Rank: 5 | ID: 5483

Title : مسی چه زمانی به پاریس برمی گردد؟

URL : https://www.farsnews.ir/news/14001014000228/مسی-چه-زمانی-به-پاریس-برمی-گردد



تعداد تکرار با توجه به دشوار بودن کلمه پایین آمده است ولی با توجه به اینکه کریسمس کلمه ای خاص است و به شکل کلی تکرار نمیشود در برخی مواقع مانند خبر بالا ارتباط بیشتری نسبت به حالت های قبل دارد به سند زیرا رایج نیست و محدوده جست و جو کمتر است.

(ج)دشوار چند کلمه‌ای: کمیسیون کریسمس

{'کمیسیون': 1, 'کریسمس': 1}

=== Cosine Scores ===

----- Results -----

Rank: 1 | ID: 5933

Title : ستاره اسپانیایی؛ هدیه کریسمس گواردیولا به ژاوی+عکس

URL :

https://www.farsnews.ir/news/14001007000739/ستاره-اسپانیایی-هدیه-کریسمس-گواردیولا-به-ژاوی-عکس

Rank: 2 | ID: 6117

Title : کی‌روش «دیکتاتور» لقب گرفت/اختلاف مرد پرتغالی با مصری‌ها به خاطر کریسمس+عکس

URL :

https://www.farsnews.ir/news/14001005000165/کی‌روش-دیکتاتور-لقب-گرفت-اختلاف-مرد-پرتغالی-با-مصری‌ها-به-خاطر-کریسمس

Rank: 3 | ID: 6120

Title : کشتار در ورزشگاه فوتبال در آستانه سال جدید

URL :

https://www.farsnews.ir/news/14001005000143/کشتار-در-ورزشگاه-فوتبال-در-آستانه-سال-جدید

Rank: 4 | ID: 5926

Title : مهاجم خارجی مس رفسنجان، 4 کودک را به محل تحصیل برگرداند +عکس

URL :

https://www.farsnews.ir/news/14001007000809/مهاجم-خارجی-مس-رفسنجان-4-کودک-را-به-محل-تحصیل-برگرداند-عکس

Rank: 5 | ID: 5483

Title : مسی چه زمانی به پاریس برمی گردد؟

URL : https://www.farsnews.ir/news/14001014000228/مسی-چه-زمانی-به-پاریس-برمی-گردد

کریسمس چون کلمه خاصی هست امتیاز آن بالاتر رفته و به عنوان مثال در داکيومنت اول میبینیم که هیچ نشانی ای از کمیسیون وجود ندارد چون کلمه عام تری هستش. بنابراین در این حالت امکان دارد ارتباط با

کوئری پایین تر بیاید هنگامی که کوئری خاص است.

کمپیون0/0

ورزشی
@Sports

ستاره اسپانیایی؛ هدیه کریسمس گواردیولا به ژاوی+عکس

طرحی جالب توسط بلیچر ریپورت به بهانه انتقال فران تورس از منچستر سیتی به بارسلونا منتشر شده است.

به گزارش خبرگزاری فارس، فران تورس ستاره اسپانیایی باشگاه منچستر سیتی با قراردادی تا سال 2027 به باشگاه بارسلونا پیوست. انتقال این بازیکن 21 ساله اسپانیایی مورد توجه بلیچر ریپورت قرار گرفته و طرحی جالب در این باره را منتشر کرده است. در این طرح، تورس به هدیه کریسمس پپ گواردیولا به ژاوی تشبیه شده است. انتهای پیام/

B4R