



STeP-1: Standard Text preparation for Persian language

کلیه حقوق این نرم افزار متعلق به آزمایشگاه پردازش زبان طبیعی دانشگاه شهید بهشتی می باشد.
لطفا در صورت استفاده از مجموعه STeP-1 به مقاله زیر ارجاع دهید.

Shamsfard, M., Jafari, H.S., Ilbeygi, M., (2010). STeP-1: A Set of Fundamental Tools for Persian Text Processing, *LREC 2010 - 8th Language Resources and Evaluation Conference*, 19-21 May, Malta.

راهنمای استفاده از توابع مربوط به ابزار STeP-1

در این سامانه دو تابع زیر برای استفاده از این ابزار ایجاد شده است:

- Tokenize (قطعه بند)
- Stem (ریشه یاب)

خروجی این توابع به فرمت جیسون (json) بوده و دارای ساختار زیر است :

```
{  
  string message;  
  bool result;  
  object Data;  
}
```

Message حاوی پیام سیستم می باشد. در صورتی که تابع با موفقیت به اجرا برسد، پیام زیر را نمایش می دهد: می گیرد "باموفقیت انجام شد" و مقدار Result برابر true خواهد بود و مقدار فیلد دیتا شامل اطلاعات مربوط به نتیجه اجرا تابع می باشد. و اگر اجرای تابع موفق نباشد مقدار فیلد Message شامل پیام خطا و فیلد result برابر false خواهد بود.

راهنمای استفاده از برنامه قطعه بند

سیستم قطعه بند، متن را به کلمات و جملات تشکیل دهنده اش تجزیه می کند. در این سیستم فاصله ها و نیم فاصله ها بین کلمات فارسی تصحیح می شود. همچنین این سیستم، متن را تا حدی بر اساس اصول نگارشی فرهنگستان زبان و ادب فارسی ویرایش می کند. در جدول زیر نمونه هایی از ورودی و خروجی برنامه آورده شده است.

ورودی	خروجی
اورادر مدرسه اش خوش حال دیدم !	او را در، اورادر] مدرسه اش خوشحال دیدم!
در 24/5 سالگی شرکت IGF را تأسیس کرد .	در 24/5 سالگی شرکت IGF [را تأسیس، را تأسیس] کرد.
کتابها تر شده اند.	[کتابها تر، کتابها تر] شده اند.
جام می بیافتاد .	جام می بیافتاد.
از تومیپرسند به کجا می روی ؟	[از تو می پرسند، از تومیپرسند] به کجا می روی؟

ساعت 2:30 دقیقه، به صرافی جی.اف.سی رفته و 45,000 تومان بابت دفترچه‌ی نشان پرداختم.	ساعت 2:30 دقیقه، به صرافی جی.اف.سی رفته و 45,000 تومان بابت دفترچه‌ی نشان پرداختم.
---	--

این تابع متن ورودی را گرفته (inputText) و متن اصلاح شده، کلمات، تمام جملات متن و کلمات
برچسب خورده را به عنوان خروجی برنامه برگردانده می‌شود.

تابع دارای ساختار زیر می‌باشد.

http://nlp.sbu.ac.ir:40101/Api/Client	آدرس پایه
Tokenize	نام تابع
Post	نوع فراخوانی

تابع دارای پارامترهای ورودی زیر است:

شرح	نوع پارامتر	نام پارامتر
متن ورودی	String	inputString
اگر این پارامتر برابر true قرار داده شود، متن داخل double quote اصلاح نمی‌شود.	Boolean	dontCorrectInsideQuote
اگر این پارامتر برابر true قرار داده شود، تنها علائم نگارشی در متن صحیح می‌شوند.	Boolean	correctDelimOnly
اگر این پارامتر برابر true قرار داده شود، در مواردی که کلمه تشخیص داده نمی‌شود، کلمه وارد تابع segmentor می‌شود. این تابع کلمه را به اجزای با معنی تجزیه می‌کند. مانند مثال 1 در جدول بالا که "اورادر" را در خروجی به صورت [او را در، اورادر] نوشته است.	Boolean	goToSegmentor
اگر این پارامتر برابر true قرار داده شود، در موارد مبهم تنها اولین حالت پیدا شده توسط برنامه، نشان داده می‌شود. مثلا در مورد "اورادر" فقط "او را در" در خروجی نمایش داده می‌شود.	Boolean	firstCaseOnly
اگر این پارامتر برابر true قرار داده شود، کلمات مرکب تشخیص داده شده توسط برنامه، در خروجی اصلاح می‌شوند	Boolean	correctCompoundWords

به صورت پیش فرض بهتر است این تابع با مقادیر زیر اجرا شود :

```
"InputString": "inputString"  
"dontCorrectInsideQuote": false,  
"correctDelimOnly" :false,  
"goToSegmentor" : false,  
"firstCaseOnly" : true,  
"correctCompoundWords":false,
```

صدا زدن تابع به صورت زیر باعث درست شدن برخی از کلمات مرکب و اشکالات تایپی که در اثر اتصال چند کلمه به هم در متن وجود دارد (مانند "اورادر") میشود. البته در این صورت سرعت اجرای برنامه نیز کندتر می‌شود.

```
"InputString": "inputString"  
"dontCorrectInsideQuote": false,  
"correctDelimOnly" :false,  
"goToSegmentor" : true,  
"firstCaseOnly" : true,  
"correctCompoundWords":true,
```

همچنین برای استفاده بهتر از سامانه بهتر است طول متن ورودی در هر فراخوانی از صدهزار کاراکتر بیشتر نباشد. چنان که اشاره شد، در صورت فراخوانی در حالت تصحیح کلمات مرکب، این محدودیت نیز با توجه به متن ورودی ممکن است بیشتر شود.

ساختار خروجی کلمات برچسب خورده Morph است، که شامل موارد زیر می باشد:
در خروجی، اولین Morph پیدا شده برای هر کلمه از متن در یک لیست برگردانده می شود. هر کدام از عناصر این لیست از نوع Morph می باشند.

```
public struct Morph  
{  
    public string word;  
    public List<string> prefixes;  
    public List<string> postfixs;  
    public string stem;  
    public string tag;  
    public string kind;  
    public string tense;  
    public int frequency;  
    public string phonetic;  
};
```

- نحوه فراخوانی:

<http://nlp.sbu.ac.ir:40101/Api/Client/Tokenize>

- نمونه ورودی:

```

{
    "InputString": "سلام",
    "dontCorrectInsideQuote": false,
    "correctDelimOnly": false,
    "goToSegmentor": false,
    "firstCaseOnly": true,
    "correctCompoundWords": true,
}

```

• نمونه خروجی:

```

1 {
2   "message": "بیا موفقیت انجام شد",
3   "result": true,
4   "data": {
5     "inputString": "سلام",
6     "tokenOutput": "سلام",
7     "taggedSubstrings": [
8       {
9         "kind": "اسم",
10        "prefixs": [],
11        "postfixs": [],
12        "word": "سلام",
13        "frequency": 270,
14        "phonetic": "salAm",
15        "num": 0,
16        "stem": "سلام",
17        "tag": "N1",
18        "tense": "",
19        "zamirMafoliChasbide": ""
20      }
21    ],
22    "sentences": [
23      "سلام"
24    ],
25    "words": [
26      "سلام"
27    ]
28  }
29 }

```

راهنمای استفاده از برنامه ریشه یاب

سیستم ریشه یاب قادر به ریشه یابی تمام کلمات تصریفی، تعدادی از کلمات اشتقاقی و تحلیل ساختاری آنهاست. این سیستم به زبان C# نوشته شده است. در جدول زیر نمونه هایی از کلماتی که این سیستم قادر به ریشه یابی است، آورده شده است.

نخواهم خواند	آینده منفی
خواندم	ساده

ساده منفی	نخواندم
ماضی استمراری	می خواندم
ماضی استمراری منفی	نمی خواندم
ماضی استمراری با داشت	داشتم می خواندم
ماضی نقلی	خوانده ام

علامت جمع "ان"	زنان
علامت جمع "گان"	پرندگان
علامت جمع "یان"	دانیان
علامت جمع "ات"	محصولات
علامت جمع "جات"	ترشیجات
علامت جمع "ون"	روحانیون
علامت جمع "ین"	معلمین

انه	ماهرانه
گی	زندگی
آگین	زهرآگین
ا	ژرفا
ه	کناره
ستان	گلستان
ین	سنگین
ینه	نقدینه
چی	درشکه چی
چه	ناوچه

این سیستم قادر به ریشه یابی افعال مرکب مانند "آماده شدیم" و کلمات مرکب مانند "کلاس بندی" نمی باشد.

تابع ریشه یاب کلمه ورودی را گرفته و لیستی از مورفولوژی های پیدا شده برای کلمه را برمی گرداند. هر کدام از عناصر این لیست از نوع Morph می باشند. Morph ساختاری است که شامل موارد زیر می باشد:

```
public struct Morph
{
    public string word;
    public List<string> prefixes;
    public List<string> postfixs;
    public string stem;
    public string tag;
    public string kind;
    public string tense;
    public int frequency;
}
```

