

بسم الله الرحمن الرحيم

پردازش زبانهای طبیعی

جلسه ۵ و ۶ (عمده جلسه ۶ از روی نوتبوک)

احسان الدین عسگری

اسفند ۱۴۰۲



<http://language.ml/>

دانشکده مهندسی کامپیوتر

آزمایشگاه پردازش هوشمند متن و زبان و علوم انسانی محاسباتی

asgari@berkeley.edu

مراحل پیش پردازش

توکنایزیشن

دیگران هم بکنند آن چه مسیح می کرد

فیض روح القدس اگر باز مدد فرماید

نرمالایزیشن

دیگران هم بکنند آن چه مسیحا می کرد

فیض روح القدس از باز مدد فرماید

ظاهر

مراحل پیش پردازش

دیگران هم بکنند آن چه مسیح می کرد	فیض روح القدس اگر باز مدد فرماید	توکنایزیشن
دیگران هم بکنند آن چه مسیح می کرد	فیض روح القدس اگر باز مدد فرماید	نرمالایزیشن
دیگران هم بکنند آن چه مسیحا می کرد	فیض روح القدس از باز مدد فرماید	ظاهر

تقلیل فرم (stemming)

- خرد کردن لغت به فرم کوچک
 - حذف پسوندها و پیشوندها
 - کلمات با ریشه یکسان به یک کلمه واحد
- خانوار - خانواده

ریشه‌یابی (lemmatization)

– بردن کلمه به ریشه

– افعال

– اسامی

با استفاده از دیکشنری یا پارسر ساخت لغت

– گرد

– مرد

?

– قوی

– شست

کاربرد ادات سخن، جزء کلام، یا Part of speech

- محدود کردن معنا (مرد؟ فعل یا اسم)
- چگونگی تلفظ (قوی؟ صفت یا اسم)
- ریشه یابی

آنچه در POS نهفته

- معنا
 - قواعد نحوی
 - ریشه‌یابی
 - و ترکیبی از آنها
- چه چالش‌هایی؟

مراحل پیش پردازش

```
In [9]: tagger = POSTagger(model='resources/postagger.model')
```

```
for x, y in tagger.tag(tokens):  
    print(F"{x} ----> {y} ")
```

```
Ne <--- فیض  
Ne <--- روح  
N <--- القدس  
CONJ <--- ار  
ADV <--- باز  
N <--- مدد  
V <--- فرماید  
PRO <--- دیگران  
CONJ <--- هم  
V <--- بکنند  
PRO <--- آن  
DET <--- چه  
N <--- مسیحا  
V <--- می‌کرد
```

ادوات سخن

دیگران هم بکنند آن چه مسیح می‌کرد

توکنایزیشن فیض روح القدس اگر باز مدد فرماید

دیگران هم بکنند آن چه مسیح می‌کرد

نرمالایزیشن فیض روح القدس اگر باز مدد فرماید

دیگران هم بکنند آن چه مسیحا می‌کرد

ظاهر فیض رُوح القدس ار باز مدد فرماید

مراحل پیش پردازش

ادات سخن

```
In [9]: tagger = POSTagger(model='resources/postagger.model')

for x, y in tagger.tag(tokens):
    print(f"{x} ----> {y} ")
```

```
Ne <--- فیض
Ne <--- روح
N <--- القدس
CONJ <--- ار
ADV <--- باز
N <--- مدد
V <--- فرماید
PRO <--- دیگران
CONJ <--- هم
V <--- بکنند
PRO <--- آن
DET <--- چه
N <--- مسیحا
V <--- میگرد
```

ریشه یابی و تقلیل فرم

```
In [8]: tokens = word_tokenize(normalized)

stemmer = Stemmer()
lemmatizer = Lemmatizer()

for tok in tokens:
    print(f"{tok} ----stem--> ", stemmer.stem(tok))
    print(f"{tok} ----lemma--> ", lemmatizer.lemmatize(tok))
```

```
فیض <--stem--> فیض
فیض <--lemma--> فیض
روح <--stem--> روح
روح <--lemma--> روح
القدس <--stem--> القدس
القدس <--lemma--> القدس
ار <--stem--> ار
ار <--lemma--> ار
باز <--stem--> باز
باز <--lemma--> باز
مدد <--stem--> مدد
مدد <--lemma--> مدد
فرماید <--stem--> فرماید
فرماید <--lemma--> فرماید
دیگران <--stem--> دیگران
دیگران <--lemma--> دیگران
هم <--stem--> هم
هم <--lemma--> هم
بکنند <--stem--> بکنند
بکنند <--lemma--> بکنند
آن <--stem--> آن
آن <--lemma--> آن
چه <--stem--> چه
چه <--lemma--> چه
مسیحا <--stem--> مسیحا
مسیحا <--lemma--> مسیحا
میگرد <--stem--> میگرد
میگرد <--lemma--> میگرد
```

توکنایزیشن

فیض روح القدس اگر باز مدد فرماید

دیگران هم بکنند آن چه مسیح می کرد

نرمالایزیشن

فیض روح القدس اگر باز مدد فرماید

دیگران هم بکنند آن چه مسیح می کرد

ظاهر

فیض رُوح القدس ار باز مدد فرماید

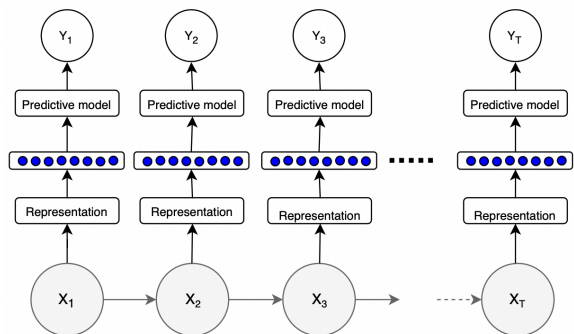
دیگران هم بکنند آن چه مسیحا می کرد

موجودیت‌های نامدار

- **PER** (Person شخص): «مولا امیرالمومنین علیه السلام»
- **LOC** (Location مکان): «حافظیه شیراز»
- **ORG** (Organization موسسه): «دانشگاه صنعتی شریف»
- **GPE** (Geo-Political Entity موقعیت ژئوپولیتیک): «جمهوری اسلامی ایران»

◦ معمولاً چند کلمه‌ای است.

- شامل موارد متعدد دیگری نیز می‌شود:
- زمان، تاریخ، قیمت



چالشهای موجودیت‌های نامدار

◦ دشواریهای تقطیع

◦ ابهام موجودیت‌های نامدار

– ایران واکسن بیشتری تولید کرد.

– مسابقات جام ملت‌های آسیا در ایران برگزار خواهد شد.

BIO تگ‌های

شیخ	B- PER
بهای	I-PER
معمار	O
و	O
دانشمند	O
برجسته	O
ایران	B-GPE
در	O
...	...

دیتاست‌های فارسی

پیما (۷۱۴۵ جمله)

Label	#
Organization	16964
Money	2037
Location	8782
Date	4259
Time	732
Person	7675
Percent	699

آرمان (۷۶۸۲ جمله)

Label	#
Organization	30108
Location	12924
Facility	4458
Event	7557
Product	4389
Person	15645