

به نام خدا

بهاره کاوسی نژاد – 99431217

کارگاه Dadmatools درس NLP

- Informal2Formal:

مشاهده می شود که در مثال اول "استاده" و "همش" و در مثال دوم "شمرون"، "معلمون" و "گنده" به حالت رسمی تبدیل نشده اند.

Testing Informal2Formal

```
✓ 1s [19] print(translator.translate('این استاد چرا انگلیسی حرف می زنه همش'))
```

این استاد است چرا انگلیسی حرف می زند همه او

```
✓ 1s print(translator.translate('دروازه شمرون به خونه بود که معلمون اونجا به مغازه گنده داشت'))
```

دروازه شمرو هستند یک خونه بود که معلمون آنجا یک مغازه گنده داشت

- Named Entity Recognition:

در مثال اول "رفت" و در مثال دوم تمامی کلمات اشتباه تشخیص داده شده اند.

```
✓ 30s [25] import dadmatools.pipeline.language as language
```

```
# as tokenizer is the default tool, it will be loaded even without calling
pips = 'ner'
nlp = language.Pipeline(pips)
text = 'محمد محمدزاده گرشاسبی به استاندارد سن دیگو رفت'
doc = nlp(text)
doc
```

Loading pretrained XLM-Roberta, this may take a while...
Model fa_tokenizer exists in cache/dadmatools/fa_tokenizer.pt
Loading tokenizer for persian
Loading multi-word expander for persian
Loading NER tagger for persian
=====

Active language: persian
=====

```
{'sentences': [{ 'id': 1,
  'tokens': [{ 'id': 1, 'text': 'محمد', 'ner': 'B-PER'},
    { 'id': 2, 'text': 'محمدزاده', 'ner': 'I-PER'},
    { 'id': 3, 'text': 'گرشاسبی', 'ner': 'E-PER'},
    { 'id': 4, 'text': 'به', 'ner': 'O'},
    { 'id': 5, 'text': 'استاندارد', 'ner': 'B-ORG'},
    { 'id': 6, 'text': 'سن', 'ner': 'I-ORG'},
    { 'id': 7, 'text': 'دیگو', 'ner': 'E-ORG'},
    { 'id': 8, 'text': 'رفت', 'ner': 'O'}] },
  'lang': 'persian'}
```

```

▶ # as tokenizer is the default tool, it will be loaded even without calling
pips = 'ner'
nlp = language.Pipeline(pips)
text = 'اي پيك نسيم صبا بار دگر گر به سر کوی دوست بگذري'
doc = nlp(text)
doc

```

```

⇨ Loading pretrained XLM-Roberta, this may take a while...
Model fa_tokenizer exists in cache/dadmatools/fa_tokenizer.pt
Loading tokenizer for persian
Loading multi-word expander for persian
Loading NER tagger for persian
=====
Active language: persian
=====
{'sentences': [{ 'id': 1,
  'tokens': [{ 'id': 1, 'text': 'اي', 'ner': 'O'},
    { 'id': 2, 'text': 'پيك', 'ner': 'O'},
    { 'id': 3, 'text': 'نسيم', 'ner': 'O'},
    { 'id': 4, 'text': 'صبا', 'ner': 'O'},
    { 'id': 5, 'text': 'بار', 'ner': 'O'},
    { 'id': 6, 'text': 'دگر', 'ner': 'O'},
    { 'id': 7, 'text': 'گر', 'ner': 'O'},
    { 'id': 8, 'text': 'به', 'ner': 'O'},
    { 'id': 9, 'text': 'سر', 'ner': 'O'},
    { 'id': 10, 'text': 'کوی', 'ner': 'O'},
    { 'id': 11, 'text': 'دوست', 'ner': 'O'},
    { 'id': 12, 'text': 'بگذري', 'ner': 'O'}]},
  'lang': 'persian'}

```

- Part of Speech Tagging:

در این مثال "بگذری" اسم شناخته شده است درحالیکه فعل است.

Testing Part of Speech Tagging

```

✓ [43] # as tokenizer is the default tool, it will be loaded even without calling
pips = 'pos'
nlp = language.Pipeline(pips)
text = 'اي پيك نسيم صبا بار دگر گر به سر کوی دوست بگذري'
doc = nlp(text)
doc

```

```

⇨ {'head': 8,
  'deprel': 'nmod:poss'},
{'id': 5,
  'text': 'بار',
  'upos': 'NOUN',
  'xpos': 'N_SING',
  'feats': 'Number=Sing',
  'head': 8,
  'deprel': 'nmod:poss'},
{'id': 6,
  'text': 'دگر',
  'upos': 'NOUN',
  'xpos': 'N_SING',
  'feats': 'Number=Sing',
  'head': 8,
  'deprel': 'nmod:poss'},
{'id': 7,
  'text': 'گر',
  'upos': 'NOUN',
  'xpos': 'N_SING',
  'feats': 'Number=Sing',
  'head': 8,
  'deprel': 'nmod:poss'}

```

✓ 26s

```

    { 'id': 9,
      'text': 'سر',
      'upos': 'NOUN',
      'xpos': 'N_SING',
      'feats': 'Number=Sing',
      'head': 8,
      'deprel': 'obj'},
    { 'id': 10,
      'text': 'کوی',
      'upos': 'NOUN',
      'xpos': 'N_SING',
      'feats': 'Number=Sing',
      'head': 8,
      'deprel': 'nmod:poss'},
    { 'id': 11,
      'text': 'دوست',
      'upos': 'NOUN',
      'xpos': 'N_SING',
      'feats': 'Number=Sing',
      'head': 8,
      'deprel': 'nsubj'},
    { 'id': 12,
      'text': 'بگنری',
      'upos': 'NOUN',
      'xpos': 'N_SING',
      'feats': 'Number=Sing',
      'head': 8,
      'deprel': 'root'}}],
    'lang': 'persian'}

```

- Dependency Parsing:

```
✓ [28] # as tokenizer is the default tool, it will be loaded even without calling
27s pips = 'dep'
nlp = language.Pipeline(pips)
```

```
↳ Loading pretrained XLM-Roberta, this may take a while...
Model fa_tokenizer exists in cache/dadmatools/fa_tokenizer.pt
Loading tokenizer for persian
Loading tagger for persian
Loading multi-word expander for persian
=====
Active language: persian
=====
```

```
✓ [29] text = 'مرده آن است که نبرند نامش به نکویی'
1s doc = nlp(text)
doc
```

```
↳ {'sentences': [{ 'id': 1,
  'tokens': [{ 'id': 1,
    'text': 'مرده',
    'upos': 'NOUN',
    'xpos': 'N_SING',
    'feats': 'Number=Sing',
    'head': 0,
    'deprel': 'root'},
    ...
    'upos': 'SCONJ',
    'xpos': 'CON',
    'head': 2,
    'deprel': 'mark'},
  { 'id': 5,
    'text': 'نبرند',
    'upos': 'NOUN',
    'xpos': 'N_SING',
    'feats': 'Number=Sing',
    'head': 2,
    'deprel': 'ccomp'},
  { 'id': 6,
    'text': 'نامش',
    'upos': 'NOUN',
    'xpos': 'N_SING',
    'feats': 'Number=Sing',
    'head': 2,
    'deprel': 'nsubj'},
  { 'id': 7,
    'text': 'به',
    'upos': 'NOUN',
    'xpos': 'N_SING',
    'feats': 'Number=Sing',
    'head': 2,
    'deprel': 'advmod'}]}
```

- Kasreh Ezafe Detection:

در این مثال سپر کسره اضافه ندارد اما تشخیص داده شده است.

Testing Kasreh Ezafe Detection

✓ 24s [30] # as tokenizer is the default tool, it will be loaded even without calling
pips = 'kasreh'
nlp = language.Pipeline(pips)

➞ Loading pretrained XLM-Roberta, this may take a while...
Model fa_tokenizer exists in cache/dadmatools/fa_tokenizer.pt
Loading tokenizer for persian
Loading multi-word expander for persian
Loading Kasreh tagger for persian
=====

Active language: persian
=====

✓ 0s ▶ text = 'سپر عقب ماشین جلویی خورد به سپر جلویی ماشین عقبی'
doc = nlp(text)
doc

➞ {'sentences': [{ 'id': 1,
 'tokens': [{ 'id': 1, 'text': 'سپر', 'kasreh': 'S-kasreh'},
 { 'id': 2, 'text': 'عقب', 'kasreh': '0'},
 { 'id': 3, 'text': 'ماشین', 'kasreh': '0'},
 { 'id': 4, 'text': 'جلویی', 'kasreh': '0'},
 { 'id': 5, 'text': 'خورد', 'kasreh': '0'},
 { 'id': 6, 'text': 'به', 'kasreh': '0'},
 { 'id': 7, 'text': 'سپر', 'kasreh': 'S-kasreh'},
 { 'id': 8, 'text': 'جلویی', 'kasreh': '0'},
 { 'id': 9, 'text': 'ماشین', 'kasreh': '0'},
 { 'id': 10, 'text': 'عقبی', 'kasreh': '0'}]}],
 'lang': 'persian'}

- Spell Checker:

املاي کلمات "آمیتیس" و "معطلی" تصحيح نشده است.

```

✓ 35s [32] # as tokenizer is the default tool, it will be loaded even without calling
      pips = 'spellchecker'
      nlp = language.Pipeline(pips)

↳ Loading pretrained XLM-Roberta, this may take a while...
Model fa_tokenizer exists in cache/dadmatools/fa_tokenizer.pt
Loading tokenizer for persian
Loading multi-word expander for persian
/usr/local/lib/python3.10/dist-packages/huggingface_hub/file_download.py:1132: FutureWarning: `resume_download`
  warnings.warn(
=====
Active language: persian
=====

✓ 1s [33] text = 'آمیٹیکس پس از مدتی معتلی بلخره رسید'
      doc = nlp(text)
      doc

↳ 1it [00:00, 2.55it/s]
{'spellchecker': {'original': 'آمیٹیکس پس از مدتی معتلی بلخره رسید',
                  'corrected': 'آمیٹیکس پس از مدتی معتلی بالآخره رسید',
                  'checked_words': [('مدتی', 'مدتی'), ('بلخره', 'بالآخره'), ('رسید', 'رسید')]},
 'sentences': [{'id': 1,
                  'tokens': [{'id': 1, 'text': 'آمیٹیکس'},
                             {'id': 2, 'text': 'پس'},
                             {'id': 3, 'text': 'از'},
                             {'id': 4, 'text': 'مدتی'},
                             {'id': 5, 'text': 'معتلی'},
                             {'id': 6, 'text': 'بلخره'},
                             {'id': 7, 'text': 'رسید'}]}],
 'lang': 'persian'}

```

- Normalizer
- Tokenizer:
- Lemmatizer:

"آرد#" ریشه درستی نیست.

```
✓ [36] pips = 'lem'
24s nlp = language.Pipeline(pips)
```

```
⇌ Loading pretrained XLM-Roberta, this may take a while...
Model fa_tokenizer exists in cache/dadmatools/fa_tokenizer.pt
Loading tokenizer for persian
Loading tagger for persian
Loading multi-word expander for persian
Loading lemmatizer for persian
=====
Active language: persian
=====
```

```
✓ [38] text = 'گردآوري كرديد'
0s doc = nlp(text)
doc
```

```
⇌ {'sentences': [{'id': 1,
  'tokens': [{'id': 1,
    'text': 'گردآوري',
    'upos': 'NOUN',
    'xpos': 'N_SING',
    'feats': 'Number=Sing',
    'head': 2,
    'deprel': 'root',
    'lemma': 'گردآوري'}],
  {'id': 2,
    'text': 'كرديد',
    'upos': 'VERB',
    'xpos': 'V_PA',
    'head': 0,
    'deprel': 'root',
    'lemma': 'آرد#آرد'}]},
  'lang': 'persian'}
```

- Sentiment Analysis:

این کنایه معنی منفی ندارد.

Testing Sentiment Analysis

```
[34] pips = 'sent'  
nlp = language.Pipeline(pips)
```

⇒ Loading pretrained XLM-Roberta, this may take a while...
Model fa_tokenizer exists in cache/dadmatools/fa_tokenizer.pt
Loading tokenizer for persian
Loading multi-word expander for persian
=====

Active language: persian
=====

```
[35] text = 'آب توی دلش تکان نمی خورد'  
doc = nlp(text)  
doc
```

⇒ {'sentences': [{'id': 1,
 'tokens': [{'id': 1, 'text': 'آب'},
 {'id': 2, 'text': 'توی'},
 {'id': 3, 'text': 'دلش'},
 {'id': 4, 'text': 'تکان'},
 {'id': 5, 'text': 'نمی'},
 {'id': 6, 'text': 'خورد'}]}],
 'lang': 'persian',
 'sentiment': [{'label': 'negative', 'score': 0.527174711227417}]}