

실습 2

Logistic Regression with Doc2vec

목표(sklearn의 linear_model사용) logistic regression

- MBTI CSV데이터를 사용하여 75%:25% 의 train, test set을 만들고
- 테스트셋의 classification accuracy 80% 이상 성능 목표
- doc2vec 사용하여 size=100, min_count=2, iter=55 으로 학습
- 중간중간 단어들의 벡터화가 어떻게 진행되는지 확인하여 이해하자.
- 제공 : 기본전처리코드, csv파일 제공

● 결과:

```
In [23]: sum(preds == test_y) / len(test_y)
```



```
Out [23]: 0.814200092208391
```

입력

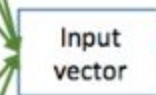
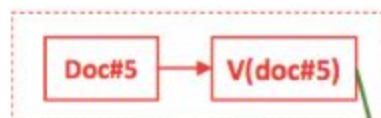
doc2vec은 word2vec의 확장이기 때문에 사용 패턴이 유사하다. 표현의 크기, 슬라이딩 윈도우의 크기, 작업자 수 또는 word2vec 모델로 변경할 수 있는 거의 모든 다른 매개변수를 쉽게 조정할 수 있다.

이 규칙의 한가지 예외는 모델에서 사용한 교육방법과 관련된 매개변수이다. word2vec 아키텍처에서는 두개의 중요 알고리즘이 있는데 “continuous bag of words”(CBOW)와 “skip-gram”(sg) 입니다. doc2vec 아키텍처에서는 “distributed memory”(dm)과 “distributed bag of words”(dbow)이다. 분산메모리 모델이 논문에서 훨씬 좋은 수행을 보였기에 이를 기본으로 사용한다. 원한다면 dm=0 플래그를 설정하여 dbow모델을 강제 실행할 수 있다.

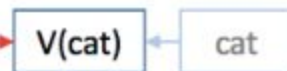
doc2vec의 입력은 **LabeledSentence** 객체의 **iterator**이다. 각 객체는 하나의 문장을 나타내며 단어 목록과 레이블 목록으로 구성된다.

```
sentence = LabeledSentence(words=[u'some', u'words', u'here'], labels=[u'SENT_1'])
```

그런다음 알고리즘은 **sentence iterator**를 두 번 수행한다.: 한번은 **vocab**를 빌드하고, 한번은 입력 데이터에서 모델을 학습하고, 각 단어 및 데이터 세트의 각 레이블에 대한 벡터 표현을 학습한다. 이 아키텍처에서는 문장당 둘 이상의 레이블을 사용할 수 있지만 가장 많이 사용되는 경우는 문장의 고유 식별자인 문장당 하나의 레이블을 갖는 것이다. 트레이닝 데이터로 다음의 클래스를 사용하여 한 줄에 한 문장의 파일에 대해 이러한 종류의 사용 사례를 구현할 수 있다.



Input vector
 v_I



Output vector
 v_o



→ Averaging
→ Softmax regression

$$\exp(v_I^T v_y) = \exp\left(\left(\frac{v_{x1} + \dots + v_{x4}}{4}\right)^T v_y\right)$$

기본 전처리 코드 제공

```
def read_corpus(fname, tokens_only=False):
    with smart_open.smart_open(fname, encoding="iso-8859-1") as f:
        for i, line in enumerate(f):
            if tokens_only:
                yield gensim.utils.simple_preprocess(line)
            else:
                # For training data, add tags
                if i==0:
                    pass
                else:
                    yield gensim.models.doc2vec.TaggedDocument(gensim.utils.simple_preprocess(line)[1:],
                                                                gensim.utils.simple_preprocess(line)[0])
```

word2vec , doc2vec 참고 사이트

https://lovit.github.io/nlp/representation/2018/03/26/word_doc_embedding/

doc2vec의 함수들을 알아 볼 수 있다.

<https://radimrehurek.com/gensim/models/doc2vec.html>

모델 **generate** 예시코드. 이후 **train**과정도 필요하다.

```
gensim.models.doc2vec.Doc2Vec
```

두가지의 과정 **build vocab**과 **train**과정이 수행되어야 한다.

```
model.build_vocab(train_corpus)
```

```
model.train(train_corpus, total_examples=model.corpus_count, epochs=model.iter)
```



```
LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,  
                    intercept_scaling=1, max_iter=100, multi_class='ovr', n_jobs=1,  
                    penalty='l2', random_state=None, solver='liblinear', tol=0.0001,  
                    verbose=0, warm_start=False)
```

sklearn에 있는 linear모델 패키지를 이용하여 LogisticRegression을
generate시킨다.

- 최종 예측 **predict** 값이다

```
sum(preds == test_y) / len(test_y)
```

```
0.8123559243891194
```

- 79~81사이로 백터가 랜덤으로 **generate**되면서
조금씩의 오차는 생길수 있다.
- 추가옵션) 임의의 **test** 파일을 만들어 **predict**해보자