5장 단어와 타입 임베딩

윤 예 준

INDEX

01 사전 훈련된 단어 임베딩

02 CBOW 임베딩 학습하기

03 문서 분류에 사전 훈련된 임베딩을 사용한 전이 학습

04 요약

01

사전 훈련된 임베딩

Word2Vec-GLoVe

임베딩 이산 타입과 벡터 공간의 포인트 사이에 매핑을 학습 하는 것

단어 임베딩 이산 타입이 단어일 때 밀집 벡터 표현

```
class PreTrainedEmbeddings(object):
    """ 사전 훈련된 단어 벡터 사용을 위한 래퍼 클래스 """
   def __init__(self, word_to_index, word_vectors):
       메개병수
           word_to_index (dict): 단어에서 정수로 매광
           word_vectors (numpy 배열의 리스트)
       self.word_to_index = word_to_index
       self.word_vectors = word_vectors
       self.index_to_word = {v: k for k, v in self.word_to_index.items()}
       self.index = AnnoyIndex(len(word_vectors[0]), metric='euclidean')
       print("인덱스 만드는 중!")
       for _, i in self.word_to_index.items():
          self.index.add_item(i, self.word_vectors[i])
       self.index.build(50)
       print("완료!")
   @classmethod
   def from_embeddings_file(cls, embedding_file):
       """사전 훈련된 벡터 파일에서 객체를 만듭니다
       벡터 파일은 다음과 같은 포맷입니다.
           word0 x0_0 x0_1 x0_2 x0_3 ... x0_N
           word1 x1 0 x1 1 x1 2 x1 3 ... x1 N
       매개병수.
           embedding_file (str): 파일 위치
          PretrainedEmbeddings의 인스턴스
       word_to_index = {}
       word_vectors = []
       with open(embedding_file) as fp:
           for line in fp.readlines():
              line = line.split(" ")
              word = line[0]
              vec = np.array([float(x) for x in line[1:]])
              word_to_index[word] = len(word_to_index)
              word vectors append(vec)
       return cls(word_to_index, word_vectors)
```

```
def get_embedding(self, word):
   메개변수
      word (str)
   반환값
   20 HIS (numpy.ndarray)
   return self.word_vectors[self.word_to_index[word]]
def get_closest_to_vector(self, vector, n=1):
    """벡터가 주어지면 n 개의 최근점 이웃을 반환합니다
       vector (np.ndarray): Annov 인덱스에 있는 벡터의 크기와 같아야 합니다
      n (int): 반환될 이웃의 개수
   반환값
      「str. str. ...]: 주어진 벡터와 가장 가까운 단어
          단어는 거리순으로 정렬되어 있지 않습니다
   nn_indices = self.index.get_nns_by_vector(vector, n)
   return [self.index_to_word[neighbor] for neighbor in nn_indices]
def compute_and_print_analogy(self, word1, word2, word3):
    """단어 임베임을 사용한 유추 결과를 출력합니다
   word10| word2일 때 word3은 __입니다
   이 메서드는 word1 : word2 :: word3 : word4를 출력합니다
   매개병수.
      word1 (str)
      word2 (str)
      word3 (str)
   vec1 = self.get_embedding(word1)
   vec2 = self.get_embedding(word2)
   vec3 = self.get_embedding(word3)
   # 네 번째 단어 일베일을 계산합니다
   spatial_relationship = vec2 - vec1
   vec4 = vec3 + spatial_relationship
   closest_words = self.get_closest_to_vector(vec4, n=4)
   existing_words = set([word1, word2, word3])
   closest_words = [word for word in closest_words
                     if word not in existing words)
   if len(closest_words) == 0:
      print("계산된 벡터와 가장 가까운 미웃을 찾을 수 없습니다!")
      return
   for word4 in closest_words:
```

print("{} : {} :: {} : {} .format(word1, word2, word3, word4))

언어 관계

단어 임베딩이 보여주는 언어 관계

```
embeddings.compute_and_print_analogy('man', 'he', 'woman')
man : he :: woman : she
man : he :: woman : her
embeddings.compute_and_print_analogy('fly', 'plane', 'sail')
fly : plane :: sail : ship
fly : plane :: sail : vessel
fly : plane :: sail : boat
embeddings.compute_and_print_analogy('cat', 'kitten', 'dog')
cat : kitten :: dog : puppy
cat : kitten :: dog : puppies
cat : kitten :: dog : junkyard
embeddings.compute_and_print_analogy('blue', 'color', 'dog')
blue : color :: dog : animal
blue : color :: dog : pet
blue : color :: dog : taste
blue : color :: dog : touch
```

단어 임베딩이 보여주는 편향

```
embeddings.compute_and_print_analogy('man', 'doctor', 'woman')
man : doctor :: woman : nurse
```

```
embeddings.compute_and_print_analogy('man', 'king', 'woman')
man : king :: woman : queen
```

02

CBOW 임베딩 학습하기

프랑켄슈타인 데이터셋

전처리

- 1. 원시 데이터로 시작
- 2. NLTK의 Punk 토큰분할기 이용하여 개별 문장으로 분할
- 3. 각 문장을 소문자로 변환 및 구두점 제거 후 공백으로 문자열 분할하여 토큰 리스트 추출

```
args = Namespace(
    raw_dataset_txt="data/books/frankenstein.txt",
    window_size=5,
    train_proportion=0.7,
    val_proportion=0.15,
    test_proportion=0.15,
    output_munged_csv="data/books/frankenstein_with_splits.csv",
    seed=1337
)

# Split the raw text book into sentences
tokenizer = nltk.data.load('tokenizers/punkt/english.pickle')
with open(args.raw_dataset_txt) as fp:
    book = fp.read()
sentences = tokenizer.tokenize(book)
```

```
# Clean sentences
def preprocess_text(text):
    text = ' '.join(word.lower() for word in text.split(" "))

text = re.sub(r"([.,!?])", r" \1 ", text)

text = re.sub(r"[^a-zA-Z.,!?]+", r" ", text)

return text
```

프랑켄슈타인 데이터셋

```
한 pitied frankenstein my pity amounted to horror i abhorred myself
원도 #1
i pitied frankenstein my pity amounted to horror i abhorred myself
원도 #2
i pitied frankenstein my pity amounted to horror i abhorred myself
원도 #3
i pitied frankenstein my pity amounted to horror i abhorred myself
원도 #4
i pitied frankenstein my pity amounted to horror i abhorred myself
```

```
1 # Global vars
2 MASK_TOKEN = "<MASK>"
 1 # Create windows
 2 flatten = lambda outer list: [item for inner list in outer list for item in inner list]
 3 windows = flatten([list(nltk.ngrams([MASK TOKEN] * args.window size + sentence.split(' ') + \
        [MASK_TOKEN] * args.window_size, args.window_size * 2 + 1)) \
        for sentence in tqdm_notebook(cleaned_sentences)])
 7 # Create cbow data
8 data = []
9 for window in tqdm_notebook(windows):
       target_token = window[args.window_size]
       context = []
        for i, token in enumerate(window):
           if token == MASK_TOKEN or i == args.window_size:
14
               continue
15
           else:
16
               context.append(token)
       data.append([' '.join(token for token in context), target token])
18
19
20 # Convert to dataframe
21 cbow data = pd.DataFrame(data, columns=["context", "target"])
```

```
# Create split data
n = len(cbow_data)
def get_split(row_num):
    if row_num <= n*args.train_proportion:
        return 'train'
elif (row_num > n*args.train_proportion) and (row_num <= n*args.train_proportion + n*args.val_proportion):
    return 'val'
else:
    return 'test'
cbow_data['split']= cbow_data.apply(lambda row: get_split(row.name), axis=1)</pre>
```

	context	target	split
0	, or the modern prometheus	frankenstein	train
1	frankenstein or the modern prometheus by		train
2	frankenstein , the modern prometheus by mary	or	train
3	frankenstein , or modern prometheus by mary wo	the	train
4	frankenstein , or the prometheus by mary wolls	modern	train

CBOW_Data

Vocabulary

```
class Vocabulary(object):
    """ 매핑을 위해 텍스트를 처리하고 어휘 사전을 만드는 클래스 """
                                                                                       38
                                                                                       39
    def __init__(self, token_to_idx=None, mask_token="<MASK>", add_unk=True, unk_token="<UNK>"):
                                                                                       41
       매개변수:
           token_to_idx (dict): 기존 토큰-인덱스 매핑 딕셔너리
                                                                                       42
           mask_token (str): Vocabulary에 추가할 MASK 토큰.
                                                                                       43
               모델 파라미터를 업데이트하는데 사용하지 않는 위치를 나타냅니다.
                                                                                       44
           add_unk (bool): UNK 토큰을 추가할지 지정하는 플래그
                                                                                       45
           unk_token (str): Vocabulary에 추가할 UNK 토큰
                                                                                       46
                                                                                       47
        if token to idx is None:
                                                                                       48
           token to idx = {}
                                                                                       49
        self. token to idx = token to idx
                                                                                       50
                                                                                       51
       self._idx_to_token = {idx: token
                                                                                       52
                           for token, idx in self. token to idx.items()}
        self. add unk = add unk
                                                                                       54
        self. unk token = unk token
                                                                                       55
       self. mask token = mask token
                                                                                       56
                                                                                       57
        self.mask index = self.add token(self. mask token)
        self.unk index = -1
                                                                                       58
        if add unk:
                                                                                       59
           self.unk_index = self.add_token(unk_token)
                                                                                       60
                                                                                       61
    def to serializable(self):
        """ 직렬화할 수 있는 딕셔너리를 반환합니다 """
                                                                                       62
        return {'token to idx': self. token to idx,
                                                                                       63
               'add unk': self. add unk,
                                                                                       64
               'unk token': self. unk token,
                                                                                       65
               'mask token': self. mask token}
                                                                                       66
```

```
@classmethod
                                                      69
def from serializable(cls, contents):
   """ 직렬화된 딕셔너리에서 Vocabulary 객체를 만듭니다 """
   return cls(**contents)
def add token(self, token):
                                                      75
   """ 토큰을 기반으로 매핑 딕셔너리를 업데이트합니다
   매개변수:
                                                      78
       token (str): Vocabularv에 추가할 토큰
                                                      79
                                                      80
                                                      81
       index (int): 토큰에 상용하는 정수
                                                      82
                                                      83
   if token in self. token to idx:
                                                      84
       index = self. token to idx[token]
                                                      85
                                                      86
                                                      87
       index = len(self. token to idx)
       self. token to idx[token] = index
       self. idx to token[index] = token
   return index
                                                      91
                                                      92
                                                      93
def add many(self, tokens):
                                                      94
   """ 토큰 리스트를 Vocabulary에 추가합니다.
                                                      95
                                                      96
   매개변수:
                                                      97
       tokens (list): 문자열 토큰 리스트
                                                      98
   반환값:
                                                      99
      indices (list): 토큰 리스트에 상용되는 인덱스 리스트
   return [self.add token(token) for token in tokens]
                                                     103
```

74

76

88

89

90

```
def lookup token(self, token):
   """ 토콘에 대응하는 인덱스를 추출합니다.
   토콘이 없으면 UNK 인덱스를 반환합니다.
   매개변수:
       token (str): 찾을 토큰
   반환값:
       index (int): 토콘에 해당하는 인덱스
       UNK 토큰을 사용하려면 (Vocabulary에 추가하기 위해)
       'unk index'가 0보다 커야 합니다.
   if self.unk index >= 0:
       return self._token_to_idx.get(token, self.unk_index)
       return self. token to idx[token]
def lookup_index(self, index):
   """ 인덱스에 해당하는 토큰을 반환합니다.
   매개변수:
       index (int): 찾을 인덱스
       token (str): 인텍스에 해당하는 토큰
   메러:
       KeyError: 인덱스가 Vocabulary에 없을 때 발생합니다.
   if index not in self. idx to token:
       raise KeyError("the index (%d) is not in the Vocabulary" % index)
   return self. idx to token[index]
def str (self):
   return "<Vocabulary(size=%d)>" % len(self)
def len (self):
   return len(self. token to idx)
```

Vectorizer

```
1 class CBOWVectorizer(object):
       """ 머휘 사전을 생성하고 관리합니다 """
       def __init__(self, cbow_vocab):
          매개변수:
              cbow_vocab (Vocabulary): 단어를 정수에 매핑합니다
          self.cbow_vocab = cbow_vocab
 9
       def vectorize(self, context, vector_length=-1):
10
          매개변수:
              context (str): 공백으로 나누어진 단어 문자열
              vector length (int): 인덱스 벡터의 길이 매개변수
          indices = [self.cbow vocab.lookup token(token) for token in context.split(' ')]
18
          if vector length < 0:
              vector_length = len(indices)
20
          out_vector = np.zeros(vector_length, dtype=np.int64)
          out vector[:len(indices)] = indices
          out vector[len(indices):] = self.cbow vocab.mask index
24
          return out vector
26
```

```
27
       @classmethod
28
       def from_dataframe(cls, cbow_df):
           """데이터셋 데이터프레임에서 Vectorizer 객체를 만듭니다
29
30
31
           매개변수::
               cbow df (pandas.DataFrame): 타깃 데이터셋
32
           반환값:
33
34
               CBOWVectorizer 객체
35
           cbow_vocab = Vocabulary()
36
           for index, row in cbow df.iterrows():
37
               for token in row.context.split(' '):
38
                  cbow vocab.add token(token)
39
               cbow vocab.add token(row.target)
40
41
42
           return cls(cbow vocab)
43
44
       @classmethod
45
       def from_serializable(cls, contents):
46
           cbow vocab = \
47
               Vocabulary.from_serializable(contents['cbow_vocab'])
           return cls(cbow vocab=cbow vocab)
48
49
50
       def to serializable(self):
           return {'cbow_vocab': self.cbow_vocab.to_serializable()}
51
```

DataLoader

```
1 class CBOWDataset(Dataset):
       def init (self, cbow df, vectorizer):
           매개변수:
               cbow df (pandas.DataFrame): 데이터셋
               vectorizer (CBOWVectorizer): 데미터셋에서 만든 CBOWVectorizer 액체
           self.cbow_df = cbow_df
           self. vectorizer = vectorizer
           measure len = lambda context: len(context.split(" "))
           self. max seq length = max(map(measure len, cbow df.context))
14
           self.train df = self.cbow df[self.cbow df.split=='train']
           self.train size = len(self.train df)
           self.val_df = self.cbow_df[self.cbow_df.split=='val']
18
           self.validation_size = len(self.val_df)
19
20
           self.test df = self.cbow df[self.cbow df.split=='test']
           self.test size = len(self.test df)
           self. lookup dict = {'train': (self.train df, self.train size),
                                'val': (self.val df, self.validation size),
                                'test': (self.test_df, self.test_size)}
26
           self.set split('train')
```

```
def load_dataset_and_make_vectorizer(cls, cbow_csv):
          """데이터셋을 로드하고 처음부터 새로운 Vectorizer 만들기
          매개변수:
              cbow csv (str): 데미터셋의 위치
              CBOWDataset의 인스턴스
          cbow_df = pd.read_csv(cbow_csv)
          train cbow df = cbow df[cbow df.split=='train']
          return cls(cbow df, CBOWVectorizer.from dataframe(train cbow df))
       @classmethod
       def load_dataset_and_load_vectorizer(cls, cbow_csv, vectorizer_filepath):
          """ 데이터셋을 로드하고 새로운 CBOWVectorizer 객체를 만듭니다.
          캐시된 CBOWVectorizer 객체를 재사용할 때 사용합니다.
          매개변수:
48
              cbow_csv (str): 데미터셋의 위치
              vectorizer filepath (str): CBOWVectorizer 액체의 저장 위치
49
50
          반환값:
              CBOWVectorizer의 인스턴스
52
          cbow df = pd.read csv(cbow csv)
54
          vectorizer = cls.load vectorizer only(vectorizer filepath)
55
          return cls(cbow_df, vectorizer)
        @staticmethod
58
        def load_vectorizer_only(vectorizer_filepath):
           """파일에서 CBOWVectorizer 객체를 로드하는 정적 메서드
59
60
61
62
               vectorizer filepath (str): 직렬화된 CBOWVectorizer 액체의 위치
           반환값:
               CBOWVectorizer의 인스턴스
64
           with open(vectorizer_filepath) as fp:
               return CBOWVectorizer.from_serializable(json.load(fp))
68
        def save vectorizer(self, vectorizer filepath):
70
           """CBOWVectorizer 액체를 ison 형태로 디스크에 저장합니다
               vectorizer filepath (str): CBOWVectorizer 액체의 저장 위치
           with open(vectorizer_filepath, "w") as fp:
               json.dump(self._vectorizer.to_serializable(), fp)
        def get_vectorizer(self):
            """ 벡터 변환 객체를 반환합니다 """
           return self. vectorizer
        def set_split(self, split="train"):
           """ 데이터프레임에 있는 열을 사용해 분할 세트를 선택합니다 """
           self. target split = split
           self._target_df, self._target_size = self._lookup_dict[split]
86
        def __len__(self):
           return self._target_size
```

```
90
       def __getitem__(self, index):
            """파이토치 데미터센의 주요 진입 메셔트
91
92
93
           매개변수:
94
               index (int): 데이터 포인트의 인덱스
95
              데이터 포인트의 특성(x_data)과 레이블(y_target)로 이루어진 딕셔너리
96
97
98
           row = self. target df.iloc[index]
99
100
           context vector = \
101
               self. vectorizer.vectorize(row.context, self. max seq length)
102
           target_index = self._vectorizer.cbow_vocab.lookup_token(row.target)
103
104
           return {'x data': context vector,
105
                   'y target': target index}
106
107
       def get num batches(self, batch size):
108
           """배치 크기가 주머지면 데이터셋으로 만들 수 있는 배치 개수를 반환합니다
109
110
           매개변수:
               batch size (int)
           반환값:
               배치 개수
114
           return len(self) // batch size
    def generate batches(dataset, batch size, shuffle=True,
                       drop_last=True, device="cpu"):
        파이토치 DataLoader를 감싸고 있는 제너레이터 함수.
       걱 텐서를 지정된 장치로 이동합니다.
       dataloader = DataLoader(dataset=dataset, batch size=batch size,
124
                             shuffle=shuffle, drop last=drop last)
126
       for data_dict in dataloader:
           out data dict = {}
128
           for name, tensor in data dict.items():
129
               out data dict[name] = data dict[name].to(device)
130
           yield out_data_dict
```

CBOWClassifier 모델

```
class CBOWClassifier(nn.Module): # Simplified chow Model
   def __init__(self, vocabulary_size, embedding_size, padding_idx=0):
       刚개增全:
          vocabulary_size (int): 어휘 사전 크기, 일베일 개수와 예측 벡터 크기를 결정합니다
          embedding_size (int): 일베일 크기
       padding_idx (int): 기본값 0: 일베임은 이 인덱스를 사용하지 않습니다
       super(CBOWClassifier, self).__init__()
       self.embedding = nn.Embedding(num_embeddings=vocabulary_size.
                                  embedding_dim=embedding_size.
                                  padding_idx=padding_idx)
       self.fc1 = nn.Linear(in_features=embedding_size.
                         out_features=vocabularv_size)
   def forward(self, x_in, apply_softmax=False):
       """분류기의 정박향 계산
       別別増全に
          x in (torch.Tensor): 일력 데이터 텐서
              x_in.shape는 (batch, input_dim)입니다
          apply_softmax (bool): 소프트맥스 활성화 함수를 위한 플래그
              크로스-엔트로피 손실을 사용하려면 False로 지정합니다
       발활값:
          결과 텐서. tensor.shape은 (batch, output_dim)입니다
       x_{embedded_sum} = F.dropout(self.embedding(x_in).sum(dim=1), 0.3)
       v_out = self.fc1(x_embedded_sum)
       if apply_softmax:
          y_out = F.softmax(y_out, dim=1)
       return y_out
```

모델 훈련 및 모델 평가

```
def make_train_state(args):
        return {'stop_early': False,
               'early_stopping_step': 0,
               'early_stopping_best_val': 1e8,
               'learning rate': args.learning rate,
 5
               'epoch_index': 0,
               'train_loss': [],
              'train acc': [],
9
               'val_loss': [],
               'val acc': [],
10
11
               'test loss': -1,
12
               'test acc': -1,
13
               'model filename': args.model state file}
```

CBOW 훈련 매개변수

테스트 손실: 7.6765486717224105; 테스트 정확도: 13.191176470588225

03

문서 분류에 사전 훈련된 임베딩을 사용한 전이 학습

AG뉴스 데이터셋

2005년 수집한 뉴스 기사 모음

스포츠, 과학/기술, 세계, 비즈니스 범주로 균등하게 분할된 뉴스 기사 12만개 사용

```
class NewsDataset(Dataset):
       @classmethod
       def load dataset and make vectorizer(cls, news csv):
           """데미터셋을 로드하고 처음부터 새로운 Vectorizer 만들기
           매개변수:
               news_csv (str): 데미터셋의 위치
           반화값:
              NewsDataset의 인스턴스
           news df = pd.read csv(news csv)
           train news df = news df[news df.split=='train']
           return cls(news_df, NewsVectorizer.from_dataframe(train_news_df))
14
      @classmethod
       def load dataset and load vectorizer(cls, news csv, vectorizer filepath):
16
          """ 데이터셋과 새로운 Vectorizer 객체를 로드합니다.
          캐시된 Vectorizer 액체를 재사용할 때 사용합니다.
18
19
          매개변수:
20
             news csv (str): 데미터셋의 위치
             vectorizer_filepath (str): Vectorizer 객체의 저장 위치
24
             NewsDataset의 민스턴스
26
          news df = pd.read csv(news csv)
          vectorizer = cls.load vectorizer only(vectorizer filepath)
28
          return cls(news_csv, vectorizer)
29
       def getitem (self, index):
30
          """파이토치 데이터셋의 주요 진입 메서드
              index (int): 데이터 포인트의 인덱스
34
35
          반환값:
              데이터 포인트의 특성(x_data)과 레이블(y_target)로 이루어진 딕셔너리
36
38
          row = self. target df.iloc[index]
39
40
          title vector = \
41
              self. vectorizer.vectorize(row.title, self. max seq length)
42
43
          category index = \
44
             self._vectorizer.category_vocab.lookup_token(row.category)
45
          return {'x_data': title_vector,
46
                  'y target': category index}
```

```
args = Namespace(
    raw dataset csv="data/ag news/news.csv",
    train proportion=0.7,
    val proportion=0.15,
    test proportion=0.15,
    output munged csv="data/ag news/news with splits.csv",
    seed=1337
# Read raw data
news = pd.read csv(args.raw dataset csv, header=0)
news.head()
   category
                                                    title
                Wall St. Bears Claw Back Into the Black (Reuters)
0 Business
   Business Carlyle Looks Toward Commercial Aerospace (Reu...
              Oil and Economy Cloud Stocks' Outlook (Reuters)
3 Business
               Iraq Halts Oil Exports from Main Southern Pipe...
4 Business
                 Oil prices soar to all-time record, posing new...
```

Vocabulary

UNK: 모델이 드물게 등장하는 단어에 대한

표현을 학습하도록 함

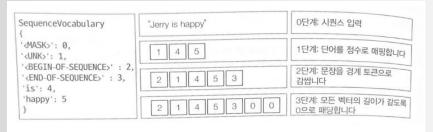
MASK: Embedding 층의 마스킹 역할 수행 및

가변 길이의 시퀸스가 있을 때 손실 계산을 도움

BEGIN-OF-SEQUENCE, END-OF0SEQUENCE: 시퀸스 경계에 관한 힌트를 신경망에 제공

```
1 class SequenceVocabulary(Vocabulary):
       def init (self, token to idx=None, unk token="<UNK>",
                   mask token="<MASK>", begin seq token="<BEGIN>",
                   end seq_token="<END>"):
           super(SequenceVocabulary, self). init (token to idx)
           self. mask token = mask token
           self. unk token = unk token
10
           self._begin_seq_token = begin_seq_token
           self. end seq token = end seq token
           self.mask index = self.add token(self. mask token)
14
           self.unk index = self.add token(self. unk token)
15
           self.begin seq index = self.add token(self. begin seq token)
16
           self.end seg index = self.add token(self. end seg token)
17
18
       def to_serializable(self):
19
           contents = super(SequenceVocabulary, self).to serializable()
20
           contents.update({'unk token': self. unk token,
                           'mask_token': self._mask_token,
                           'begin seq token': self. begin seq token,
23
                           'end seq token': self. end seq token})
24
           return contents
25
26
       def lookup token(self, token):
27
           """ 토콘에 대용하는 인덱스를 추출합니다.
           토콘미 없으면 UNK 인덱스를 반환합니다.
28
29
30
           매개변수:
31
              token (str): 찾을 토큰
32
           반환값:
33
              index (int): 토콘에 해당하는 인덱스
34
           노트:
              UNK 토큰을 사용하려면 (Vocabulary에 추가하기 위해)
35
36
               `unk index`가 0보다 커야 합니다.
38
           if self.unk index >= 0:
39
               return self._token_to_idx.get(token, self.unk_index)
40
           else:
41
              return self. token to idx[token]
```

Vectorizer



```
1 class NewsVectorizer(object):
       """ 머휘 사전을 생성하고 관리합니다 """
       def __init__(self, title_vocab, category_vocab):
           self.title vocab = title vocab
           self.category_vocab = category_vocab
       def vectorize(self, title, vector_length=-1):
           매개변수:
 9
              title (str): 공백으로 나누어진 단어 문자열
10
               vector length (int): 인덱스 벡터의 길이 매개변수
11
           반환값:
               벡터로 변환된 제목 (numpy.array)
13
14
           indices = [self.title vocab.begin seq index]
16
           indices.extend(self.title vocab.lookup token(token)
17
                         for token in title.split(" "))
           indices.append(self.title_vocab.end_seq_index)
18
19
20
           if vector length < 0:
21
               vector length = len(indices)
           out vector = np.zeros(vector length, dtype=np.int64)
24
           out_vector[:len(indices)] = indices
25
           out_vector[len(indices):] = self.title_vocab.mask_index
26
           return out vector
28
```

```
29
       @classmethod
30
       def from_dataframe(cls, news_df, cutoff=25):
           """데이터셋 데이터프레임에서 Vectorizer 객체를 만듭니다
31
32
           매개변수:
33
34
               news df (pandas.DataFrame): 타깃 데이터셋
               cutoff (int): Vocabulary에 포함할 빈도 임곗값
35
36
           반환값:
               NewsVectorizer 객체
37
38
39
           category vocab = Vocabulary()
40
           for category in sorted(set(news df.category)):
41
               category vocab.add token(category)
42
43
           word counts = Counter()
44
           for title in news df.title:
45
               for token in title.split(" "):
46
                   if token not in string.punctuation:
47
                       word counts[token] += 1
48
49
           title_vocab = SequenceVocabulary()
50
           for word, word count in word counts.items():
51
               if word count >= cutoff:
52
                   title vocab.add token(word)
53
54
           return cls(title vocab, category vocab)
55
56
       @classmethod
57
       def from serializable(cls, contents):
58
           title vocab = \
59
               SequenceVocabulary.from serializable(contents['title vocab'])
60
           category vocab = \
61
               Vocabulary.from serializable(contents['category vocab'])
62
63
           return cls(title vocab=title vocab, category vocab=category vocab)
64
65
       def to_serializable(self):
66
           return {'title vocab': self.title vocab.to serializable(),
                   'category vocab': self.category vocab.to serializable()}
67
```

NewsClassifier 모델

```
11 def load glove from file(glove filepath):
       """Glove 임베딩 로드
       매개변수:
14
           glove filepath (str): 임베딩 파일 경로
           word to index (dict), embeddings (numpy.ndarary)
20
       word to index = {}
       embeddings = []
       with open(glove_filepath, "r") as fp:
          for index, line in enumerate(fp):
              line = line.split(" ") # each line: word num1 num2 ...
              word to index[line[0]] = index # word = Line[0]
              embedding_i = np.array([float(val) for val in line[1:]])
              embeddings.append(embedding i)
28
       return word to index, np.stack(embeddings)
       make embedding matrix(glove filepath, words):
       특정 단대 진합에 대한 일베딩 행력을 만듭니다.
       매개변수:
34
           glove filepath (str): 임베딩 파일 경로
          words (list): 단메 리스트
38
       word to idx, glove embeddings = load glove from file(glove filepath)
       embedding size = glove embeddings.shape[1]
39
       final_embeddings = np.zeros((len(words), embedding_size))
       for i, word in enumerate(words):
44
           if word in word to idx:
               final_embeddings[i, :] = glove_embeddings[word_to_idx[word]]
46
47
               embedding i = torch.ones(1, embedding size)
              torch.nn.init.xavier uniform (embedding i)
49
              final_embeddings[i, :] = embedding_i
       return final embeddings
```

```
1 class NewsClassifier(nn.Module):
        def __init__(self, embedding_size, num_embeddings, num_channels,
                     hidden dim, num classes, dropout p,
                    pretrained_embeddings=None, padding_idx=0):
               embedding_size (int): 임베딩 벡터의 크기
               num embeddings (int): 임베딩 벡터의 개수
               num channels (int): 합성곱 커널 개수
               hidden_dim (int): 은닉 차원 크기
10
               num classes (int): 클래스 개수
               dropout_p (float): 드롭마웃 확률
               pretrained embeddings (numpy.array): 사전에 훈련된 단어 임베딩
                   기본값은 None
               padding_idx (int): 패딩 인덱스
          super(NewsClassifier, self), init ()
18
19
          if pretrained embeddings is None:
20
              self.emb = nn.Embedding(embedding dim=embedding size,
                                   num_embeddings=num_embeddings,
                                   padding_idx=padding_idx)
24
              pretrained embeddings = torch.from numpy(pretrained embeddings).float()
              self.emb = nn.Embedding(embedding dim-embedding size,
                                   num embeddings=num embeddings,
                                   padding idx-padding idx.
29
                                   _weight-pretrained_embeddings)
32
           self.convnet = nn.Sequential(
               nn.Conv1d(in channels=embedding size,
34
                      out_channels=num_channels, kernel_size=3),
35
               nn.Conv1d(in channels=num channels, out channels=num channels,
                      kernel size=3, stride=2).
38
39
               nn.Conv1d(in_channels=num_channels, out_channels=num_channels,
40
                      kernel_size=3, stride=2),
41
42
               nn.Conv1d(in_channels=num_channels, out_channels=num_channels,
                      kernel size=3),
44
45
46
           self, dropout p = dropout p
48
           self.fc1 = nn.Linear(num channels, hidden dim)
49
           self.fc2 = nn.Linear(hidden dim, num classes)
```

```
def forward(self, x in, apply softmax=False):
          """분류기의 정방향 계산
52
          매개변수:
54
55
             x in (torch.Tensor): 입력 데이터 텐서
56
                 x_in.shape는 (batch, dataset._max_seq_length)입니다.
             apply softmax (bool): 소프트맥스 활성화 함수를 위한 플래그
                 크로스-엔트로피 손실을 사용하려면 False로 지정합니다.
58
          반환값:
59
              결과 텐서. tensor.shape은 (batch, num_classes)입니다.
60
61
62
          # 일베일을 적용하고 특성과 채널 차원을 바꿉니다
63
64
          x_embedded = self.emb(x_in).permute(0, 2, 1)
65
66
          features = self.convnet(x embedded)
67
          # 평균 값을 계산하여 부가적인 차원을 제거합니다
68
          remaining size = features.size(dim=2)
69
70
          features = F.avg pool1d(features, remaining size).squeeze(dim=2)
          features = F.dropout(features, p=self. dropout p)
          # MLP 분류기
          intermediate vector = F.relu(F.dropout(self.fc1(features), p=self._dropout_p))
74
75
          prediction vector = self.fc2(intermediate vector)
76
          if apply softmax:
78
              prediction vector = F.softmax(prediction vector, dim=1)
79
80
          return prediction vector
```

모델 훈련 및 평가

```
args = Namespace(
   # 날짜와 경로 정보
   news_csv="data/ag_news/news_with_splits.csv",
   vectorizer file="vectorizer.json",
   model state file="model.pth",
   save_dir="model_storage/ch5/document_classification",
   # 모델 하이퍼파라미터
   glove filepath='data/glove/glove.6B.100d.txt',
   use glove=False,
   embedding size=100,
   hidden_dim=100,
   num channels=100,
   # 훈련 하이퍼파라미터
   seed=1337,
   learning rate=0.001,
   dropout p=0.1,
   batch size=128,
   num epochs=100,
   early stopping criteria=5,
   # 실행 옵션
   cuda=True,
   catch keyboard interrupt=True,
   reload from files=False,
   expand filepaths to save dir=True
```

테스트 손실: 0.6176266238093375; 테스트 정확도: 79.6819196428571

새로운 뉴스 제목의 카테고리 예측

```
def predict category(title, classifier, vectorizer, max length):
       """뉴스 제목을 기반으로 카테고리를 예측합니다
       매개변수:
          title (str): 원시 제목 문자열
          classifier (NewsClassifier): 훈련된 분류기 객체
          vectorizer (NewsVectorizer): 해당 Vectorizer
          max length (int): 최대 시퀀스 길이
 8
              노트: CNN은 입력 텐서 크기에 민감합니다.
 9
                   훈련 데이터처럼 동일한 크기를 갖도록 만듭니다.
10
       11 11 11
11
12
       title = preprocess text(title)
       vectorized title = \
13
          torch.tensor(vectorizer.vectorize(title, vector_length=max_length))
14
       result = classifier(vectorized title.unsqueeze(0), apply softmax=True)
15
16
       probability values, indices = result.max(dim=1)
17
       predicted category = vectorizer.category vocab.lookup index(indices.item())
18
19
       return {'category': predicted category,
20
              'probability': probability values.item()}
```

⁰⁴ 요약

- 1. 사전 훈련된 임베딩을 블랙박스처럼 사용하는 방법
- 2. CBOW를 포함해 데이터에서 이런 임베딩을 직접 훈련하는 방법(간략히)
- 3. 언어 모델링 측면에서 CBOW 모델을 훈련하는 방법
- 4. 문서 분류에 사전 훈련된 임베딩을 사용한 전이 학습

THE

감사합니다

END