



CNN을 이용한 텍스트 분류

Sang Yup Lee

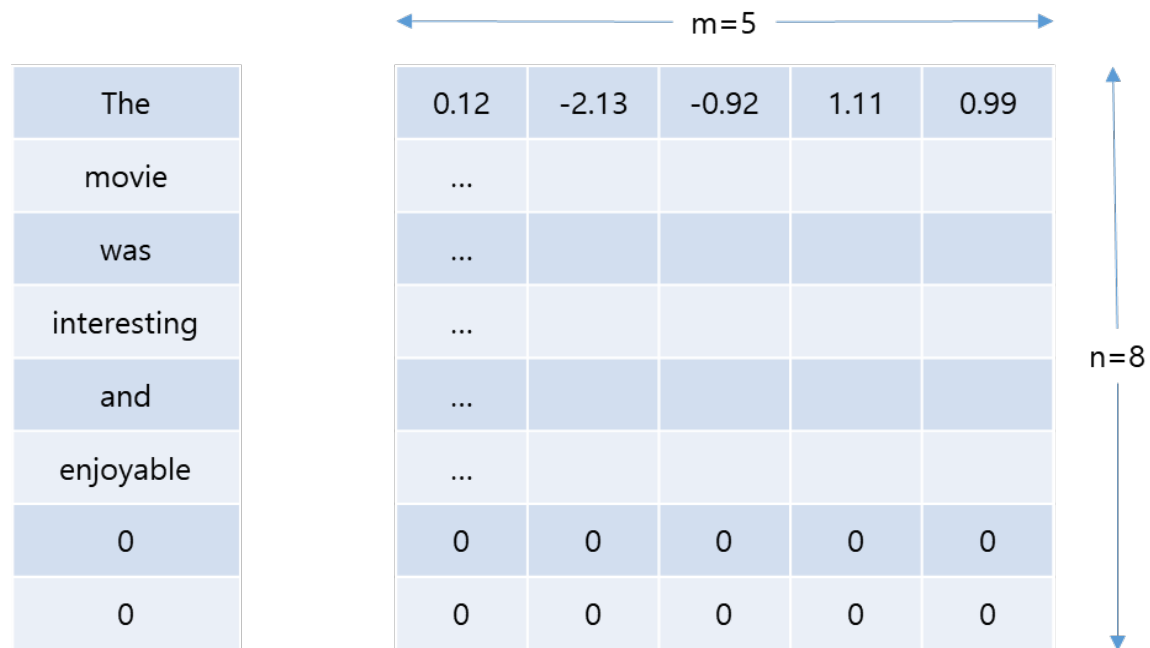


CNN을 이용한 텍스트 분류

- 순서 1
 - 문서를 3D 형태, 즉 (n, m, c) 로 표현 (즉, 하나의 이미지와 같은 형태로 표현)
 - n 은 문서를 표현할 때 사용하는 최대 단어 수
 - m 은 한 단어를 표현하는 임베딩 벡터의 차원
 - c 는 이미지 데이터에서의 채널수로 문서에서는 1

CNN을 이용한 텍스트 분류

- 예) $n=8, m=5, c=1$
 - 문서: The movie was interesting and enjoyable



CNN for text classification

CNN을 이용한 텍스트 분류

■ 순서 2

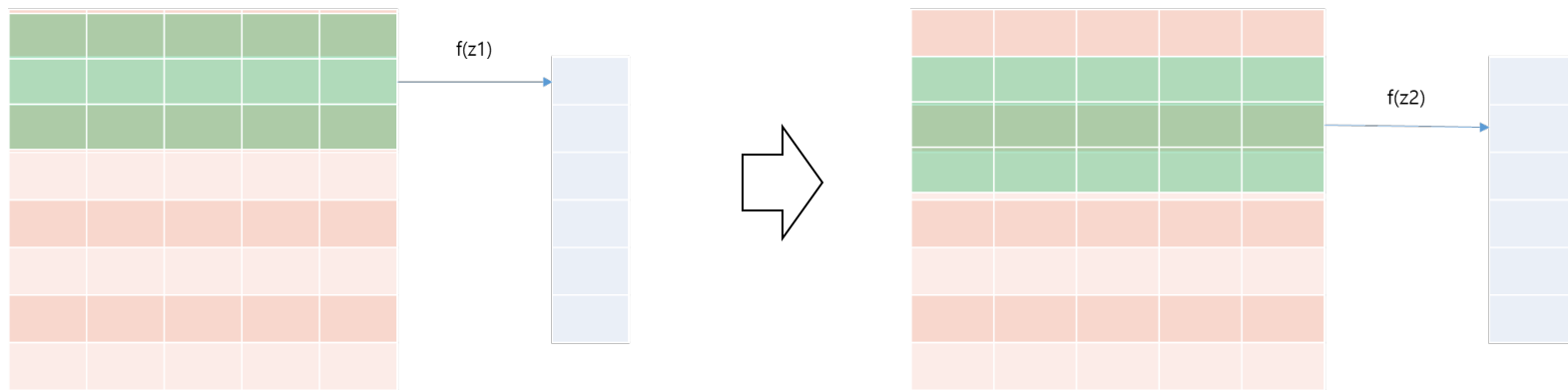
■ 합성곱 필터 적용하기

- 텍스트의 경우, $n \times m \times c$ 의 문서에 대해서 $k \times m \times c$ 형태의 필터를 적용
 - 여기에서 $k < m$
 - 앞의 예에서는 $k \times 5 \times 1$ 를 적용
 - $k=3$ 인 경우

$w_{1,1}$	$w_{1,2}$	$w_{1,3}$	$w_{1,4}$	$w_{1,5}$
$w_{2,1}$	$w_{2,2}$	$w_{2,3}$	$w_{2,4}$	$w_{2,5}$
$w_{3,1}$	$w_{3,2}$	$w_{3,3}$	$w_{3,4}$	$w_{3,5}$

CNN을 이용한 텍스트 분류

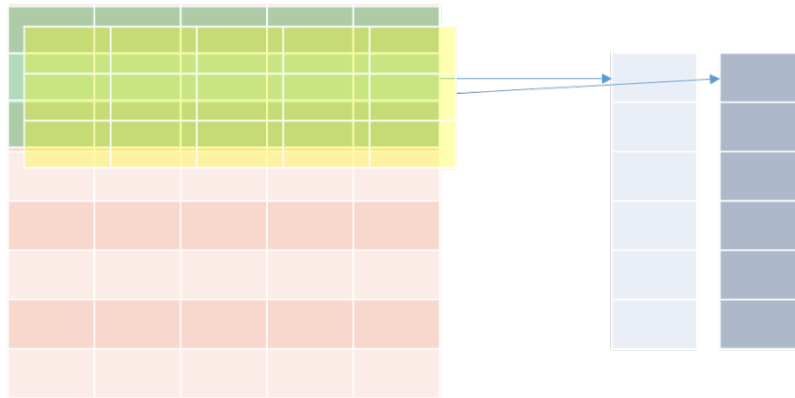
- 필터의 적용
 - 문서의 경우는 보통 $\text{stride}=1$



결과물은 $(n-k+1)$ 크기의 1D array

CNN을 이용한 텍스트 분류

- 필터의 적용 (cont'd)
 - 크기가 같은 필터 2개를 적용하는 경우



$n \times m$ 형태의 문서 array에 $k \times m$ 형태의 필터를 h 개 적용하는 경우에 생기는 결과물은 $(n-k+1) \times h$



CNN을 이용한 텍스트 분류

- 그 다음 순서
 - 이미지에서와 마찬가지로 pooling layer 추가
 - Flattening
 - Fully connected layer
 - 출력층 with softmax()
- Python coding
 - See "CNN_imdb_example.ipynb"