

XY절단알고리즘과 중첩신경망에 의한 문서화상의 영역분할방법

리광철, 고진혁

문서화상으로부터 검색을 위한 정보들을 추출하고 분석하고 저장하기 위해서는 먼저 문서화상을 본문구역과 비본문구역에 해당하는 영역들로 분할하여야 한다.

지금까지 문서화상의 영역분할을 위한 많은 방법[1, 2]들이 제안되었지만 대부분의 방법들은 문서화상의 내용이나 특징에 대한 사전지식을 필요로 하거나 문서화상의 질에 의존하는 결함이 있었다.

본문에서는 XY절단알고리즘과 중첩신경망을 리용하여 문서화상의 영역들을 유형별로 분할하기 위한 한가지 방법을 제안하였다.

1. XY절단알고리즘과 삽화검출에 의한 문서화상의 영역분할

문서화상에는 본문, 그림, 표 등 각이한 형식의 영역들이 존재한다. 그리하여 먼저 XY절단알고리즘[3]을 리용하여 문서화상에 대한 초기분할을 진행한다.

XY알고리즘은 반복적인 top-down알고리즘으로서 문서화상의 매 화소들을 수평 및 수직방향으로 투영하여 얻어지는 투영히스토그램에 기초하여 문서화상영역을 반복적으로 분할해나간다.

그러나 분할된 영역들에는 본문으로 된 영역뿐아니라 삽화와 본문이 함께 있는 영역들도 있다.

이로부터 다음의 알고리즘에 의하여 매 분할된 영역들에서 삽화영역을 검출한다.

본문영역과 삽화영역을 구별하기 위한 특징으로서 자체상관행렬을 리용한다.

문서화상을 크기가 $n \times n$ 인 블록들로 분할하고 매 블록에서 자체상관행렬을 다음의 식에 따라 계산한다.

$$C(k, l) = \sum_{y=\max(0, l)}^{n-1+\min(0, l)} \sum_{x=\max(0, k)}^{n-1+\min(0, k)} I(x, y) \cdot I(x+k, y+l) \quad (1)$$

여기서 l 과 k 는 $[-n/2, n/2]$ 에서 정의되며 블록의 크기 n 은 본문블록의 반복적인 패턴들이 강조되도록 설정한다.

식 (1)에 의한 자체상관계산은 많은 계산량을 필요로 하며 효율적이지 못하다.

자체상관함수의 푸리에변환은 에너지스펙트럼과 같고 결국 자체상관은 에너지스펙트럼의 역푸리에변환과 같다.

그러므로 블록의 자체상관을 고속푸리에변환의 두단계 계산을 통하여 효과적으로 계산할수 있다.

블록에 대한 자체상관은 무늬에서 의미가 있는 방향을 추정하는데 리용된다.

그러므로 $C(k, l)$ 을 극자리표계로 변환하고 다음의 식에 따라 방향히스토그램을 구한다.

$$w(\theta) = \sum_{r \in (0, n/2]} C(r \cos \theta, r \sin \theta) \quad (2)$$

자체상관행렬은 정의에 의하여 중심에 관하여 대칭이기때문에 $[0^\circ, 180^\circ)$ 범위내에서의 방향히스토그램만 고려한다.

매 블록들을 방향히스토그램에 의하여 부호화하면 본문블록들은 0° 와 180° 근방에서 봉우리가 나타나며 반대로 삽화블록들은 다중양상형태로 나타난다.

최종적으로 방향히스토그램과 자체상관행렬을 수평 및 수직방향으로 투영하여 얻어지는 사영히스토그램을 연결하여 블록에 대한 특징벡토르를 구성하고 SVM분류기를 리용하여 본문영역과 삽화영역을 판별한다.

검출된 삽화영역들을 제거하고 다시 XY절단알고리즘을 적용하여 최종적인 분할영역들을 얻는다.

2. 중첩신경망에 의한 문서화상의 영역류형분류

논문에서는 문서화상에 포함되는 영역의 류형을 크게 6가지 즉 본문, 그림, 그래프, 표, 수학적식, 악보로 하였다.

영역류형의 분류는 중첩신경망을 리용하여 진행한다.

중첩신경망의 입력은 하나의 영역으로부터 추출한 크기가 $n \times n$ 인 블록이며 출력은 위의 6가지 류형중 어느 하나의 류형이다.

영역의 류형분류를 위한 중첩신경망의 구성을 표에 보여주었다.

표. 영역의 류형분류를 위한 중첩신경망의 구성

층	핵크기	걸음크기	러파기수	입력	출력
conv1	11×11	1×1	32	(1, n, n)	(32, n, n)
mp1	2×2	2×2	-	(32, n, n)	(32, n/2, n/2)
conv2	5×5	1×1	32	(32, n/2, n/2)	(32, n/2, n/2)
mp2	2×2	2×2	-	(32, n/2, n/2)	(32, n/4, n/4)
conv3	3×3	1×1	32	(32, n/4, n/4)	(32, n/4, n/4)
mp3	2×2	2×2	-	(32, n/4, n/4)	(32, n/8, n/8)
fc1				$(32 \cdot n^2/64)$	1 024
dropout				1 024	512
fc2				512	6

어떤 영역에 대한 분류를 하기 위하여 영역안의 매 블록들을 학습된 중첩신경망에 입력하였을 때의 fc1층의 출력벡토르(1 024차원)를 얻고 영역안의 모든 블록들에 대한 1 024차원의 특징벡토르들에 대한 평균과 공분산행렬을 각각 구한다.

영역의 기하학적자리표가 (R_x, R_y, R_w, R_h) 이고 블록들의 모임 B 를 포함하고있는 영역 R 의 특징벡토르는 다음과 같다.

$$f_R = \left[\frac{R_x}{W}, \frac{R_y}{H}, \frac{R_w}{W}, \frac{R_h}{H}, \frac{R_w}{R_h}, \mu(B), \sigma(B) \right] \quad (3)$$

여기서 W, H 는 각각 페이지의 너비와 높이이며 $\mu(B), \sigma(B)$ 는 각각 B 안의 모든 블록들에 대한 평균벡토르와 공분산행렬의 왼쪽윗3각형원소들을 1차원으로 라렬한 벡토르이다.

영역에 대한 최종분류는 영역에 대한 특징벡터 f_R 와 우연숲분류기(Random Forest Classifier)[4]를 리용하여 진행한다.

3. 문서화상의 영역분할실험

실험에서 리용된 문서화상의 크기는 2 448×3 264이며 블록의 크기는 64×64로 설정하였다.

문서화상에 대하여 XY절단알고리즘, 삽화검출알고리즘, 영역류형분류알고리즘을 각각 적용하였을 때의 결과를 그림에 보여주었다.

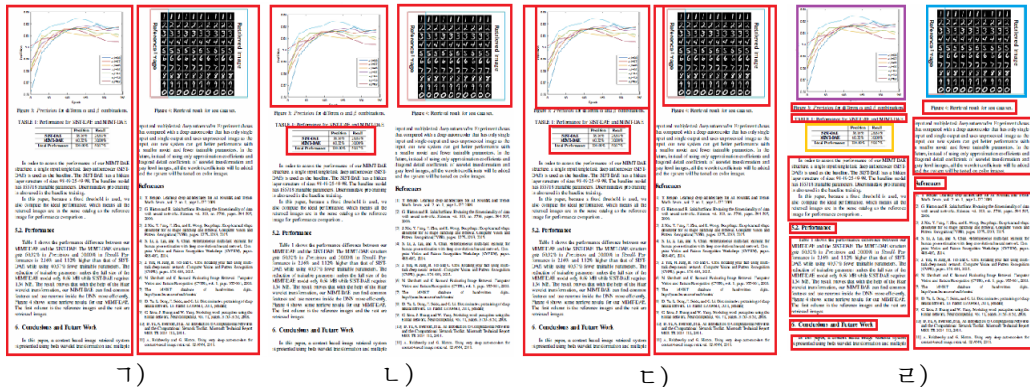


그림. 문서화상의 영역분할실험결과

- 가) XY절단알고리즘적용결과, 나) 삽화검출알고리즘적용결과,
다) XY절단알고리즘을 다시 적용한 결과, 라) 영역류형분류결과

본문의 방법에 의한 문서화상의 영역분할정확도는 약 92%로서 XY절단방법이나 백색영역분석방법에 비하여 약 20% 더 높다.

맺는 말

XY절단알고리즘과 중첩신경망을 리용하여 문서화상의 영역분할과 영역들의 류형결정을 동시에 진행하기 위한 한가지 방법을 제안하고 실험을 통하여 그 효과성을 검증하였다.

참고 문헌

- [1] Pradipta Maji, Shaswati Roy; Applied Soft Computing, 1, 1, 2015.
- [2] Y. L. Chen, B. F. Wu; Pattern Recognit, 42, 1419, 2009.
- [3] Jaekyu Ha et al.; Proceedings of the Third International Conference on Document Analysis and Recognition, 2, 952, 1995.
- [4] Archana Chaudhary et al.; Information Processing in Agriculture, 3, 4, 215, 2016.

Segmentation Method of Document Image by XY-cut Algorithm and CNN

Ri Kwang Chol, Ko Jin Hyok

In this paper, we propose a method for clustering regions of document image using XY-cut algorithm and CNN.

Keywords: XY-cut algorithm, CNN, document image, segmentation