

## Support Vector Machine을 이용한 유해 이미지 분류

송철환<sup>0</sup> 유성준

세종대학교 컴퓨터 소프트웨어 공학과

peternara<sup>0</sup>@naver.com, sjyoo@sejong.ac.kr

### Adult Image Filtering using Support Vector Machine

Chull Hwan Song<sup>0</sup>, Seong Joon Yoo

School of Computer Engineering, Sejong University, 98 Gunja, Gwangjin

Seoul, Korea 143-747

#### 요 약

본 논문은 인터넷의 대표적인 문제점중의 하나인 Adult Image 분류 연구에 대해 기술한다. 특히 우리는 이러한 Adult Image를 분류하기 위한 Data Set을 5가지 타입으로 구성한다. 이러한 각 Image에 대해 Color, Gradient, Edge Direction 특성의 Feature들을 추출하고 이를 Histogram으로 구성한다. 이렇게 구성된 Histogram을 Support Vector Machine에 적용하여 Adult Image를 분류한다. 그 결과, 우리는 8250개의 Test Set에 대하여 Recall(96.53%), Precision(97.33%), False Positive(2.96%), F-Measure(96.93%)의 성능 결과를 보여준다.

#### 1. 서 론

최근 인터넷의 빠른 성장에 힘입어 많은 사람들은 어느 곳에서든 인터넷을 쉽게 접근하여 이용한다. 또한 대부분의 인터넷 사용자들은 각 개인의 선호도에 따라 그 유용한 정보를 얻는다. 반면에 많은 문제점을 내포하고 있는 것도 사실이다. 그 대표적인 문제점으로 인식하고 있는 것들이 Spam Mail과 Adult Multimedia Data이다. 특히 Adult Multimedia Data의 경우에 미성년자들이 인터넷을 통하여 접근한다면 심각한 문제가 아닐 수 없다. 따라서 본 논문은 인터넷의 심각한 부작용 중의 하나인 Adult Image 분류 연구에 대해 기술한다. 대부분의 인터넷 사용자들은 각 개인의 선호도에 따라 그 유용한 정보를 얻는다. 반면에 많은 문제점을 내포하고 있는 것도 사실이다. 그 대표적인 문제점으로 인식하고 있는 것들이 Spam Mail과 Adult Multimedia Data이다. 특히 Adult Multimedia Data의 경우에 미성년자들이 인터넷을 통하여 접근한다면 심각한 문제가 아닐 수 없다. 따라서 본 논문은 인터넷의 심각한 부작용 중의 하나인 Adult Image 분류 연구에 대해 기술한다. 그리고 우리의 연구는 먼저 Adult Image를 분류하고자 할 때 나타나는 문제점들을 정의한다. 즉, Adult Image는 Skin Color에 기반하여 분류하게 된다. 이러한 방법은 많은 문제점을 가지고 있는데 특히, Web 상에 빈번하게 나타나는 얼굴 이미지들이 많이 존재한다. 이를 위해 우리는 Face Detection방법을 적용하거나 Segmentation알고리즘들을 적용한 후, Skin Color를 검사하는 방법으로 그 문제점들을 해결한다. 그 후, Blob 이미지들에 대한 다양한 Image Feature를 추출하여 분류 알고리즘을 Support Vector Machine에 적용

하여 Adult Image를 분류한다.

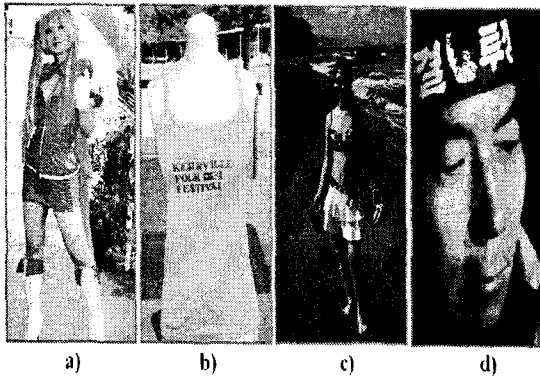
#### 2. 이전 연구

지금까지 Adult Image Classification 연구는 다양한 형태의 이미지 Feature를 추출하고 그 Feature를 다양한 형태의 분류 알고리즘에 적용하여 진행되어 왔다.

먼저 Jiao et al.[1] 연구에서는 Adult Image에 대하여 Color Coherence와 Color Histogram을 구성하고 이를 Support Vector Machine에 적용하여 분류하였다. 이 연구의 실험 결과는 2400개의 Test Image에 대하여 약 89.3% 비율로 Adult Image를 분류하였고 이때의 False Positive는 약 9.4%의 결과를 보여주었다.

Zheng et al.[2]은 Hybrid Approach 방법을 이용하여 분류한다. 즉, Face Detection과 Skin Color Detection에 기반한 Semantic High-Level Feature를 구성하고 이를 Adult Image에 적용하여 분류한다. //이 연구의 실험 결과는 2196개의 Test Image에 대하여 약 92% 비율로 분류하였고 이때의 False Positive는 약 3.96%의 결과를 보여주었다.

Yoo et al.[3]은 여러MPEG-7 Visual Descriptor 알고리즘들을 적용하여 이미지에 대한 각 Feature 들을 추출한 후, 이미지들의 Feature에 대한 유사성 측정 방법으로 분류하였다. 그리고 전체 2800개의 Test Image에 대한 실험 결과로써 약 99.2% 비율(Recall)로 Adult Image를 분류하였고 False Positive는 약 23%의 결과를 보여주었다.



[그림 1] Skin Color Region 기반으로 하여 분류할 때 문제점을 가지고 있는 사진. a) 여자가 입고 있는 옷이 Skin Color를 포함하고 있는 사진, b)는 사람이 포함되지 않고 Skin Color 영역을 포함하고 있는 사진, c)는 해변에서 수영복 입고 있는 사진, d)는 얼굴 영역이 대부분 사진 지역을 차지하고 있는 사진

Kim et al.[4]은 [3]보다 다양한 MPEG-7 Visual Descriptor를 이용하여 이미지 Feature를 추출한 후, Neural Network에 적용하여 Adult Image를 분류하였다. 이 연구의 결과는 5397개의 Test Set에 대하여 Color Structure Visual Descriptor를 적용할 때 가장 좋은 성능을 보여 주었고 그 때의 결과가 약 94.69%의 Recall과 약 2.55%의 False Positive 비율을 나타내었다.

### 3. 유해 이미지 분류 문제

Adult Image Classification의 대부분의 연구는 이미지 안에 나타나는 Skin 영역에 기반한다. 그러나 이러한 Skin Color 영역 중심의 연구는 많은 문제점을 내포하고 있다. 즉, 이미지에 나타나는 스킨 영역이 많을 경우에 Adult Image 일 수 있지만, 그렇지 않을 경우도 있다. [그림 1]은 이에 대한 예제를 보여준다. [그림 1]과 같이 Skin Color 영역이 사진의 대부분의 영역을 차지하고 있다면 Adult Image로 분류될 가능성이 다분하다.

따라서 이러한 특수한 사진에 대해서는 다양한 방법을 사용하여 분류해야 한다. 예를 들어 [그림 1]의 d)경우는 [2]의 연구에서 언급한 것처럼 Face Detection 알고리즘을 이용하여 그 Face Detection 지역이 사진 영역의 일정 크기 이상을 가진다면 Adult Image가 아니라고 판단한다. 우리는 이를 OpenCV의 Face Detection 라이브러리를 이용하여 해결한다. 한편 [그림 1]의 a)~b) 사진의 경우에는 매우 감지하기 어렵다. 따라서 우리는 이를 최대한 해결하기 위해서 Image Segmentation Algorithm을 사용한 다음 각 Segment Image들에 대한 Skin Color 영역의 크기를 검사하여 일정 기준(Threshold) 이하의 경우는 제외 시켜서 최대한의 잡음을 제거한다.

[그림 1]의 c) 같은 경우는 그 의미적으로 모호한 문제이므로 우리는 이러한 사진은 Adult Image로 포함 시킨다.

### 4. 이미지 특징추출

#### 4.1 Color Histogram based on HSV Color Space

Image의 대표적인 Feature는 Color, Texture, Shape 등이 있다. 이들 요소 중 가장 중요한 요소는 당연히 Color이다. 인간은 사물을 인식하는데 있어서 대부분을 Color에 의존하기 때문이다. 우리는 Hue-Saturation-Value (HSV) Color Space를 이용한다. HSV Color Space는 RGB 보다 Illumination에 민감하지 않고 HS(색상 요소)를 쉽게 분리할 수 있다. Cox et al.[5]의 연구에서는 HSV를 이용하여 11개의 Bin을 구성하였다. 우리는 이를 HSV 36Bin으로 재구성한다. 그리고 이를 Color Feature로 선택한다.

#### 4.2. Histogram of Canny Edge Detection using Sobel Filter

Color는 Image Classification을 하기 위한 중요 요소 중 하나이지만 이미지 안의 다른 오브젝트가 같은 색깔로 이루어질 때 Color만을 사용하여 이미지 분류에 이용한다면 문제가 발생할 것이다. 따라서 우리는 Edge Detection 알고리즘에서 획득할 수 있는 Gradient 및 Edge Direction요소들을 추출하여 이용한다. 대표적인 Edge Detection 알고리즘의 하나인 Laplacian Edge Detection을 비롯한 대부분의 Edge Detection 알고리즘은 이미지의 잡음에 민감하여 작은 잡음이 이미지 안에 존재할 경우 Edge로 간주하여 추출되는 경우가 발생한다. 우리는 보다 덜 민감한 알고리즘인 Canny Edge Detection[6] 알고리즘을 이용하여 Gradient와 Edge Direction 요소를 추출하여 Histogram을 구성한다. Canny Edge Detection 알고리즘은 다음과 같은 4 Step의 과정을 갖는다.

1. Gaussian Smoothing Filtering: 원 이미지의 잡음을 제거해주는 역할을 한다.
2. Edge Strength: 이미지의 Gradient를 가져오기 위해 Edge Strength를 찾는다. 이때 Sobel, Priwitt, Roberts 과 같은 Mask를 적용할 수 있는데 우리는 Edge를 상대적으로 잘 감지하는 Sobel Mask를 이용한다. 즉, [그림 2]를 식(1)에 적용하여 Gradient를 구한다. 이때 Gx는 수평적 요소, Gy는 수직적 요소를 나타낸다.

$$|G| = |Gx| + |Gy| \dots\dots\dots(1)$$

3. Edge Direction Finding: Gx와 Gy를 입력 값으로 하여 식(2)에 적용하면 Edge Direction을 구할 수 있다.

$$\theta = \arctan (Gy / Gx) \dots\dots\dots(2)$$

-1	0	+1
-2	0	+2
-1	0	+1

**Gx**

+1	+2	+1
0	0	0
-1	-2	-1

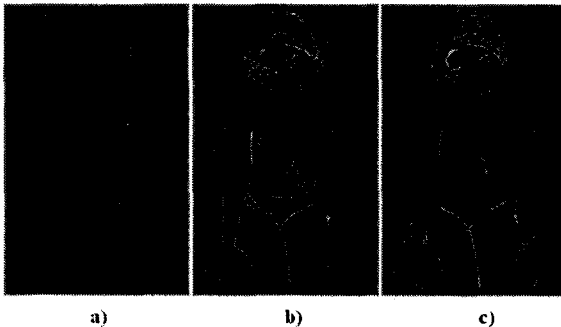
**Gy**

[그림 2] Sobel Mask

그래서, 우리는 0 degrees (in the horizontal direction), 45 degrees (along the positive diagonal), 90 degrees (in the vertical direction), 135 degrees (along the negative diagonal)를 기준으로 4 가지 방향을 구할 수 있다.

4. Hysteresis: 위의 Step에서 구한 Edge Direction을 이용하여 이미지에 적용시키는 단계이다.

[그림 3]은 각 Mask에 따라 Canny Edge Detection 알고리즘을 적용한 결과를 보여준다. Sobel Mask를 이용할 때 좀 더 좋은 Edge를 감지하는 것을 알 수 있다.



[그림 3] 각 Mask를 적용한 Canny Edge Histogram 결과. a) Roberts Mask, b) Prewitt Mask, c) Sobel Mask

그래서, 우리는 Sobel Mask를 이용한 Canny Edge Detection에서 Gradient 요소와 Edge Direction 요소를 획득하여 Histogram을 구성한다. 좀더 자세히 설명하면 식(1)의 Gx, Gy에 대해 각각 4Bin을 구성하고 식(2)의 theta 값에 나타나는 4가지의 Edge Direction을 이용하여 4Bin을 구성한다.

### 4.3 Combined Image Feature Histogram

우리는 이제까지 36Bin HSV Color Histogram, Skin Color 1Bin, Canny Edge Detection 알고리즘을 이용한 8Bin Gradient Histogram(수직요소 4Bin+수평요소 4Bin = 8Bin) 그리고 4Bin Edge Direction Histogram 구성에 대해 알아 보았다. 우리는 이를 결합하여 사용한다.

즉, 49Bin (36+1+4+4+4) Image Feature Histogram을 구성한다. 이러한 49 크기의 Bin은 SVM에 적용하는데 있어서 큰 입력 차원에 해당되지 않는다. 즉, SVM은 분류하는데 있어서 조금 느리다는 단점이 존재하는데 49 크기의 Bin은 SVM에 많은 부담을 주지 않고 분류할 수 있게 만든다. 그리고 우리의 실험에서 37Bin Color Histogram과 49Bin Image Feature Histogram을 Support Vector Machine에 적용하여 Adult Image를 분류하고 그 성능을 비교한다.

### 5. 실험 및 결과

첫 번째, 37Bin Color Histogram을 이용하고 두 번째, 49Bin Image Feature (Combined Color + Gradient + Edge Direction) Histogram 들을 Support Vector Machine에 적용하여 Adult Image를 분류한다. // 즉, 분류기로는 Support Vector Machine[7]을 이용하고 각 Kernel에 따른 성능비율을 보여준다. 실험을 위해 우리는 LIBSVM을 이용한다. 그리고 학습과 Test는 10-fold Cross-Validation 과정을 통해 이루어진다.

실험을 위한 Data Set의 구성은 5 가지(A~E Type)로 구성된다. 첫 번째, A Type은 여성의 가슴이 포함된 사진 1700 장, 두 번째, B Type은 비키니 수영복과 이와 비슷한 복장을 포함하는 여성 사진 1701 장, 세 번째, C Type은 전신 여성 누드 사진으로 1702 장, 네 번째, D Type은 성행위가 포함된 사진으로 1703 장, 다섯 번째, E Type은 일반 Non-Adult 사진으로 1444 장을 포함한다. 또한 A~D Type은 Adult Image로 구성하고, E Type은 Adult Image가 아닌 일반 이미지(Non-Adult Image)로 구성한다. 따라서 전체 이미지 8250 장은 Adult Image를 분류하기 위한 Data Set으로 이용한다.

우리 연구에서의 성능평가는 Precision, Recall, F-Measure 그리고 False Positive 과 같은 평가 방법들을 이용하여 측정한다.

[표 2] 37 Bin Color Histogram 결과

SVM Kernel	Recall (%)	Precision (%)	False Positive (%)	F-Measure (%)
Linear	90.66	97.12	5.05	93.78
RBF	93.86	96.30	2.75	95.06
Polynomial	97.01	88.95	5.01	92.81
Sigmoid	90.43	96.20	5.41	93.23

[표 1]의 경우는 37Bin Color Histogram을 입력 값으로 받는 SVM 결과와 각 SVM Kernel에 따른 평가 결과를 보여준다. 이때, SVM에 RBF Kernel을 적용할 때, Recall (93.86%), Precision (96.30%), False Positive

[표 1] 49Bin Combined Image Feature Histogram  
결과

SVM Kernel	Recall (%)	Precision (%)	False Positive (%)	F-Measure (%)
Linear	95.83	96.78	2.51	96.30
RBF	96.53	97.30	2.52	96.91
Polynomial	<b>96.53</b>	<b>97.33</b>	<b>2.48</b>	<b>96.93</b>
Sigmoid	95.75	97.02	2.96	96.38

(2.75%), F-Measure (95.06%)의 결과를 나타내어 가장 높은 성능을 보여주었다.

[표 2]의 경우는 4.3 장에서 설명하였던 49Bin 을 가지는 Histogram 을 이용한 결과를 보여준다. 특히, Table 2 의 37Bin Color Histogram 을 이용하는 것보다 좀더 나은 성능을 보여주고 있다. 즉, 가장 좋은 성능을 나타낼 경우는 Support Vector Machine 의 Polynomial Kernel, 49Bin Image Feature Histogram 을 적용할 때 Recall (96.53%), Precision (97.33%), False Positive (2.96%), F-Measure (96.93%)의 결과를 얻어내었다. 그렇지만 37Bin Color Histogram 을 적용할 때도 좋은 성능이 나타났음을 알 수 있다.

## 6. 결론

우리는 이제까지 Adult Image 분류 연구에 대해 기술하였다. 그리고 Adult Image 를 분류하기 위해 Segmentation Algorithm 과 Skin Color 추출 알고리즘을 이용하여 Skin Color Segment 이미지를 획득하였다. 이렇게 획득한 이미지에서 Color, Gradient 그리고 Edge Direction 요소에 대한 Histogram을 추출하여 Support Vector Machine에 적용하였다. 그리고 그 결과를 다른 이전 연구와 비교하였다. 이전 연구와 비교해 볼 때 우리의 연구는 더 많은 Test Set을 이용하여 평가 했는데도 불구하고 보다 높은 성능 결과를 보여주었다. 이러한 연구는 컴퓨터와 인터넷의 활성화에 따른 부작용을 해결하는데 있어서 매우 중요하다. 따라서 향후 우리는 Adult Image 분류 성능을 더욱 향상시킬 예정이다. 또한 인터넷의 다른 심각한 문제점 중의 하나인 Spam 문서 분류 연구에 대해서도 진행시킬 예정이다.

## 7. 참고 논문

[1] Jiao, F., Gao, W., Duan, L and Cui, G.: Detecting adult image using multiple features. IEEE conference, Vol.3, pp.378-383 (2001)

- [2] Zheng, Q.F., Zhang, M.J. and Wang, W.O.: A Hybrid Approach to Detect Adult Web Images. PCM 2004, LNCS 3332, pp.609-616 (2004)
- [3] Yoo, S.J., Jung, M.H. and Kang, H.B.: Composition of MPEG-7 Visual Descriptors for Detecting Adult Image on the Internet. HIS 2003, LNCS 2713, pp. 682-687 (2003)
- [4] Kim, W.I., Lee, H.K., Yoo, S.J., and Baik, W.W.: Neural Network Based Adult Image Classification. ICANN 2005, LNCS 3696, pp. 481-486 (2005)
- [5] Cox, I. J., Miller, M.L., Omohundro, S.M., Yianilos, P.N.: Target testing and the PicHunter Bayesian multimedia retrieval system in the Proceedings of the 3rd Forum on Research and Technology Advances in Digital Libraries, DL'96, pp. 66-75 (1996)
- [6] Canny, J.: A computational approach to edge detection.: IEEE Transaction Pattern Analysis Machine Intelligence 8(6), pp. 679-698 (1986)
- [7] Vanpik, V.: Statistical learning theory. Wiley, New York (1998)