

Modified Sequential Algorithmic Schema를 이용한

디지털 사진의 효율적인 분류

이상린*

*부산대학교

Modified Sequential Algorithm schema for Efficient Digital Image retrieval

Sang-lyn Lee*

*Pusan National University

E-mail : sanglynn@Gmail.com

요 약

이 논문에서는 수정된 Sequential Algorithmic Schema를 이용해서 여러 장소를 이동하면서 찍은 디지털 이미지를 효율적으로 분류할 수 있는 방법을 제안한다. 제안하는 방법은 이웃 패턴들과 특징 정보의 연속성, 유사성을 가지며 들어오는 입력 패턴에 대해 기존의 모든 군집과 유사도를 비교하는 방법이 아니라 이전 군집의 정보와 유사도를 비교하여 군집에 포함시키거나 동적으로 군집을 생성하는 효율적인 군집화 방법이다. 제안한 방법은 실험을 통해서 기존의 군집화 기법에 성능 및 속도의 효율성을 증명하였다.

키워드

Sequential Algorithmic Schema, Clustering, 디지털 사진

1. 서 론

현대사회는 멀티미디어의 급격한 발전에 따라 엄청난 양의 디지털 멀티미디어 데이터들이 생성되고 저장된다. 따라서 과거와 달리 텍스트기반의 멀티미디어 패턴 분류로는 패턴을 군집화하기에 텍스트 정보의 모호성과 부족한 정보로 군집화된 패턴이 자신의 군집과 관계없이 잘못된 군집화 결과가 발생한다. 그리하여 최근 대용량의 멀티미디어 데이터베이스에 대한 내용기반의 데이터 분류에 대한 많은 연구가 진행되어 지고 있다. 내용기반 데이터 분류를 통해서, 특히 이미지 패턴의 경우, 기존의 텍스트기반 이미지 검색에 비해 이미지를 분석하여 색상, 질감, 모양, 공간 정보 등의 특징을 추출하고 특징을 비교하는 식으로 유사도를 비교한다.[1]

전형적인 군집화 알고리즘들은 입력되는 패턴이 생성되어 있는 기존의 모든 군집과 비교함으로써 최적의 군집화 결과를 반환한다. 그러나 위와 같은 과정에서 대용량의 데이터베이스의 경우 많은 군집과 데이터베이스의 모든 패턴을 서로 비교함으로써 계산 량이 데이터베이스의 크기에 비례하여 많은 시간적인 비용이 발생한다. 특히 데이터베이스의 각 패턴이 가지는 특징들이 자신

의 이웃패턴들과 특징정보의 유사성을 가지는 연속된 이미지, 동영상의 경우 전형적인 군집화 알고리즘을 사용할 때 불필요한 비교 연산이 이루어진다. 예를 들어 A 장소에서 찍은 사진들은 같은 군집을 이루어 분류가 되어야한다. 어떤 이미지의 특징 정보(색깔, 공간정보 등)가 아무리 A 장소 사진들과 유사할지라도 사용자는 A 장소에서 찍지 않은 이미지들은 같은 군집에 포함되길 기대하지 않는다.

본 논문에서는 기존의 Sequential Algorithmic schema를 수정한 MSAS(Modified Sequential Algorithmic schema)를 이용하여 이전 군집과 이웃 패턴만 비교함으로써 여러 장소를 이동하면서 찍은 디지털 이미지, 연속된 이미지 등을 효율적으로 분류할 수 있는 방법을 제안한다.

본 논문의 구성은 2장에서는 본 논문과 관련된 있는 SAS에 대해서 기술 하였고 3장에서는 MSAS 기법을 이용한 군집화 기법에 대해서, 4장에서는 여러 장소를 이동하면서 찍은 디지털 사진 그룹에 대한 실험과 결과에 대해서 기술하고 있다. 마지막에는 결론과 향후 연구 방향에 대해 기술한다.

II. 관련 연구

SAS(Sequential Algorithmic Schema)는 입력 될 패턴들에 대해서 군집할 개수를 미리 정하지 않고 동적으로 군집을 생성하며 군집화 하는 기법이다. 입력되는 패턴 x_i 에 대해서 기존에 생성된 군집 C_k 중에 가장 유사한 군집을 선택, vigilance test를 통해 사용자가 미리 정한 vigilance parameter 값 Θ 보다 크고 사용자가 미리 정해놓은 최대허용 클러스터 q 보다 현재까지 생성된 클러스터 총 개수 m이 작을 때 새로운 군집을 생성하고 작으면 선택된 군집에 입력된 패턴을 포함시키고 선택된 군집의 대푯값 m_{C_k} 을 업데이트하는 기법이다.[2]

SAS의 의사코드(pseudo code)는 아래 표 1과 같다.

표 1. Sequential Algorithmic Schema

| |
|---|
| <p>입력패턴 집합 $X = \{x_1, x_2, \dots, x_N\}$ N : 입력패턴 수 m : 현재까지 생성된 클러스터 개수 q : 최대 허용 클러스터 개수 C_k : k 번째 클러스터 n_{C_k} : k번째 클러스터에 포함된 패턴의 개수 m_{C_k} : k 번째 클러스터의 대푯값 Θ : vigilance parameter</p> <p>m = 1 $C_m = \{x_1\}$ For i = 2 to N Find $C_k : d(x_i, C_k) = \min_{1 \leq j \leq m} d(x_i, C_j)$ if $(d(x_i, C_k) > \Theta) \text{ AND } (m < q)$ then m = m+1 $C_m = \{x_i\}$ else $C_k = C_k \cup \{x_i\}$ $m_{C_k}^{new} = \frac{(n_{C_k}^{new} - 1)m_{C_k}^{old} + x}{n_{C_k}^{new}}$ End[if] End[For]</p> |
|---|

여기서 $d(x, C)$ 는 입력 패턴과 클러스터 C와의 유사도 차이로써 식 (1) 과 같이 클러스터 C의 대푯값과 입력패턴간의 거리로 계산된다.

$$d(x, C) = d(x, m_C) \tag{1}$$

III. MSAS(Modified Sequential Algorithmic Schema)

이장에서는 이웃패턴과 내용정보의 유사성 및 연속성을 가지고 입력되는 패턴들을 효율적으로 군집화 하기 위해서 기존의 방식인 SAS 개선한 MSAS(Modified Sequential Algorithmic Schema)를 제안한다. 기존의 SAS 기법의 수정해 입력되는 패턴 x_i 와 그 이웃 패턴들이 생성한 가장 최근에 생성된 군집 C_m 와의 유사도를 계산하고 C_m 에 속하는 이웃 패턴 x_{i-1} 중에서 입력 패턴과 가까운 것과의 거리를 계산한다. 계산된 두 값에서 입력패턴과 가장 가까운 값을 선택하여 vigilance test를 한다. 사용자가 미리 정한 임계값 (vigilance parameter) 값 Θ 보다 선택된 거리가 작다면 입력 패턴은 C_m 에 포함되고 C_m 의 대푯값 m_{C_m} 을 갱신시킨다. 그렇지 않으면 최대 허용 군집 개수 q 보다 현재 까지 생성된 군집의 개수 m 이 작은 경우에 한에서 새로운 군집을 생성하고 입력패턴이 새로운 패턴에 포함되고 대푯값이 갱신된다. MSAS의 의사코드(pseudo code)는 아래 표 2와 같다.

표 2. 수정된 Sequential Algorithmic Schema

| |
|--|
| <p>입력패턴 집합 $X = \{x_1, x_2, \dots, x_N\}$ N : 입력패턴 수 m : 현재까지 생성된 클러스터 개수 q : 최대 허용 클러스터 개수 C_m : m 번째 클러스터 n_{C_m} : m번째 클러스터에 포함된 패턴의 개수 m_{C_m} : 최근에 생성된 m 번째 클러스터의 대푯값 Θ : vigilance parameter</p> <p>m = 1 $C_m = \{x_1\}$ For i = 2 to N if $\min(d(x_i, C_m), d(x_i, x_{i-1})) < \Theta$ $(x_{i-1} \in C_m)$ then $C_m = C_m \cup \{x_i\}$ $m_{C_m}^{new} = \frac{(n_{C_m}^{new} - 1)m_{C_m}^{old} + x}{n_{C_m}^{new}}$ else if (m < q) then</p> |
|--|

```

        m = m+1
        Cm = {xi}
    End{if}
End{if}
End{For}
    
```

여기서 $d(x, C)$ 는 입력 패턴과 클러스터 C와의 유사도 차이로써 식 (2) 과 같이 클러스터 C의 대푯값과 입력패턴간의 거리로 계산된다.

$$d(x, C) = d(x, m_C) \quad (2)$$

기존의 Sequential Algorithmic schema 기법 또는 전형적인 군집화 기법을 사용할 경우 그림 1 과 같이 이웃 패턴과 연속성을 가지는 패턴일지라도 이웃 패턴과의 특징적 연속성, 유사성을 무시하고 모든 패턴 혹은 모든 군집과 비교를 하게 된다. 따라서 많은 시간적 비용이 발생하게 된다. 제안한 기법을 이용할 경우 그림1의 (a)부터 (c) 까지 순차적으로 비교를 하게 되며 물리적 거리가 먼 패턴(장소가 다른 곳에서 찍은 디지털 사진)의 경우 비교하지 않을 확률이 높아지게 됨에 따라 연산 량이 줄어든다.



(a) (b) (c)
그림 1. 특징의 유사성을 가지는 연속된 이미지

IV. 실험 및 결과

장소를 이동하면서 직접 찍은 디지털 사진 총 300장을 가지고 제안한 기법에 대한 성능 평가를 하였다. 시스템에 입력되는 사진은 되어있다. 실험에 사용된 이미지는 일반적으로 여행에서 디지털 카메라를 이용하여 자연스럽게 찍은 데이터로써 장소에 따라 특정장소에 찍은 사진들이 내용의 유사성이 내포된 데이터로 구성되어 있다.

이미지의 특징을 생성하기 전에 선처리에 필요한 연산 량을 줄이기 위해 실험할 이미지 그림 4 의 픽셀의 크기를 256x256 픽셀로 줄였으며 또한 64컬러로 양자화 하여 이미지의 특징을 분석하였다.

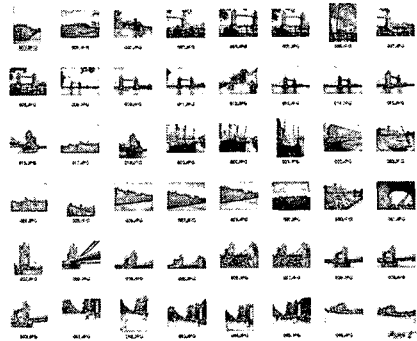


그림 2. 실험에 사용된 디지털 사진

본 논문에서는 제안한 기법인 MSAS와 SAS 기법의 성능을 분류 속도와 분류 정확도를 이용해 비교하였으며 이를 아래의 표 3 과 표 4의 실험 결과를 통해 성능 평가하였다.

표 3은 제안한 기법에 대한 분류 속도를 비교한 것으로 입력되는 디지털 사진의 양에 따라 두 기법이 분류하는데 걸리는 시간을 비교한 것이다. 시스템의 분류 속도는 한 번만 측정하게 되면 수행 시간에 오차가 생길 수 있기 때문에, 여러 번 측정하여 평균(mean, average)값을 취하였다. 군집화 임계값은 1.0을 두 패턴 간 가장 유사하다고 보았을 때 0.75로 지정 하였다. 두 패턴의 유사도가 0.75이상일 경우 같은 군집에 포함되도록 하였다. 최대 허용 그룹의 개수는 무한대로 두고 실험하였다.

표 3. 제안한 기법에 대한 분류 속도 비교 (단위:ms)

| 사진 개수 | SAS | MSAS |
|-------|---------|--------|
| 100 | 31.745 | 5.4328 |
| 150 | 67.211 | 8.958 |
| 200 | 116.103 | 15.671 |
| 250 | 147.692 | 18.490 |
| 300 | 218.622 | 19.178 |

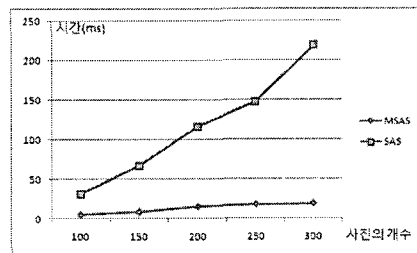


그림 3. 입력되는 사진의 개수에 따른 분류 속도 비교 그래프

그림 3의 그래프를 통해 입력되는 사진의 양이 늘어남에 따라 MSAS 기법과 SAS 기법을 이용한 분류의 연산 량을 차이를 알 수 있다. 사진의 개수가 늘어날수록 MSAS와 SAS 기법이 연산하는 연산 량의 차이가 크게 나타나는 것을 보인다. SAS 기법을 이용할 경우 동적으로 생성되는 군집의 개수가 작을 때는 연산 량이 작지만 개수가 증가할수록 연산 량이 커지게 된다. MSAS 기법의 경우 최근에 생성된 군집과 이전의 사진과 비교만 하기 때문에 군집의 개수가 늘어나더라도 연산 량에 영향을 받지 않는다.

다음은 MSAS 기법의 분류 정확도를 통해 MSAS 기법의 성능을 평가하였다. 분류 정확도는 사용자가 기대하는 그룹으로 군집화 되는지 평가한다. 사용자는 장소의 이동에 따라 찍은 사진이 내용의 연속성을 가지고 분류되도록 기대한다. A 공간에서 찍은 사진들은 A공간에 찍은 사진들끼리 군집화를 이루고 B 공간에서 찍은 사진들은 B 공간에서 찍은 사진끼리 군집화를 이루기를 기대한다. 표 4가 나타내는 분류 정확도 성능평가는 다음과 같이 식3으로 나타낼 수 있다.

$$\text{분류정확도} = \frac{a}{a+b} \quad (3)$$

a는 분류된 사진 중에서 사용자가 기대하는 그룹으로 군집화된 사진 개수의 평균, b는 분류된 사진 중에서 사용자가 기대하는 그룹으로 군집화 되지 않는 사진 개수의 평균이라고 정의한다. 또한 군집화 임계값은 0.75로 두었으며 최대 허용 그룹 개수는 입력되는 사진개수만큼 두고 실험하였다.



(a) 129번째 사진 (b) 149번째 사진 (c) 151번째 사진
그림 4. 기존의 분류방식을 이용한 분류 문제점

표 4. 제안한 기법에 대한 분류 정확도

| 사진개수 | MSAS |
|------|--------|
| | 분류 정확도 |
| 30 | 73.3% |
| 50 | 74% |
| 70 | 62.9% |

그림 4의 (a) 사진과 (c) 사진 사이에 다른 장

소의 사진이 있음에도 불구하고 그림4.(a) 사진과 그림 4.(b) 사진은 같은 군집으로 겹쳐지는 결과가 나타난다. 생성된 그룹에서 사진이 포함하는 내용이나 시간적 순서를 전혀 고려하지 않고 오로지 Color Histogram을 통해 구해진 사진의 특징 정보만을 모든 군집에 대해서 유사도를 비교되기 때문에 입력되는 사진들이 장소에 따라 혹은 시간에 따라 군집화될 가능성이 MSAS에 비해 현저히 떨어진다.

장소의 이동에 따라 찍힌 순서대로 입력되는 사진을 MSAS를 통해 차례로 군집화 됨으로써 사용자가 기대하는 분류의 만족도를 표4와 같이 나타냈다. 단 표4의 결과는 사용자의 주관적 기대치에 대한 만족도를 나타낸다. 표 4의 결과를 볼 때 MSAS 기법을 이용하여 군집화할 때 사진의 개수가 늘어나도 분류의 정확도가 크게 떨어지지 않는다. 입력되는 사진에 대하여 순차적으로 군집화가 이루어지며 가장 최근에 만들어진 그룹을 제외한 과거의 군집은 다음 군집을 생성할 때 사용되지 않기 때문이다.

V. 결 론 및 향후 연구 과제

본 논문에서는 수정된 Sequential Algorithmic Schema를 이용하여 이전 군집, 이웃 패턴만 비교함으로써 여러 장소를 이동하면서 찍은 디지털 사진, 연속된 이미지 등에 대해서 효율적으로 분류할 수 있는 방법을 제안하였다. 4장에서 이루어진 실험을 통해서 본 논문이 제안한 기법을 이용하여 여러 장소에서 이동하면서 찍은 디지털 사진에 대하여 전형적인 군집화방법중 하나인 SAS 기법에 비해 속도와 군집화 되는 정확도면에서 나은 성능을 보였다. 특히 입력 패턴의 내용이 이웃 패턴과 연속성, 유사성을 띄면서 들어올 경우 기존의 군집화 방법을 수정한 기법을 통해 연산 량을 크게 줄일 수 있음을 보였다.

향후 연구로서 내용의 전개나 공간의 이동이 뚜렷한 동영상에 대하여 제안한 기법을 이용하여 영상을 몇 개의 인덱스로 나눌 수 있도록 하여 영상의 내용에 기반을 둔 동영상의 인덱스 설정을 할 수 있도록 한다.

참고문헌

- [1] 송석진, 이희봉, 김효성, 남기곤, "히스토그램 인터섹션과 오토코릴로그램을 이용한 내용기반 영상검색 시스템", 한국 신호처리, 시스템 학회 2002년도 제3권 1호, pp1-7
- [2] S. Theodoridis, K. Koutroumbas, Pattern Recognition, Academic Press, pp.433-437 2003