

신경망 및 통계적 클러스터링의 성능 비교 분석

민준영¹⁾, 허문열²⁾

< 요약 >

클러스터링 알고리즘은 통계적인 방법, ISODATA 알고리즘, 신경망 클러스터링으로 크게 구분할 수 있다(김대수, 1994). 이중 신경망을 이용한 클러스터링 방법에는 Kohonen 네트워크, Carpenter와 Grossberg의 ART 네트워크가 있다. Kohonen 네트워크(1990)에는 SOFM(Self Organization Feature Maps), LVQ(Learning Vector Quantization) 알고리즘이 있으며, 이 방법은 패턴과 클러스터의 중심값과의 거리를 최소화시키는 학습 알고리즘에 따라 클러스터링을 한다는 의미에서 통계적 방법중 k-means 방법과 대응된다고 할 수 있다. 그리고 최근에 발표된 것으로 Pal et. al.(1993)의 GLVQ(Generalized Learning Vector Quantization) 알고리즘이 있다. 본 논문은 신경망을 이용한 클러스터링 알고리즘중 GLVQ 알고리즘과 이에 대응하는 통계적 클러스터링 방법에서의 k-means 방법을 비교하여 그 결과에 따른 성능을 평가한다. 평가방법으로는 Rand c 통계량으로 비교하였으며, 데이터는 Fisher의 IRIS 데이터와 필기체 숫자에서 특징추출한 패턴벡터를 이용하였다. 평가결과 신경망 클러스터링인 경우가 통계적인 방법에 비하여 Rand c 통계량의 평균이 높고 표준편차도 상대적으로 적게 나타났음을 알 수 있었다.

1. 서론

클러스터링 알고리즘은 통계적인 방법, ISODATA 알고리즘, 신경망 클러스터링으로 크게 구분할 수 있다(김대수, 1994). 이중 신경망을 이용한 클러스터링 방법에는

1) 상지대학교병설전문대학 전자계산과 교수
2) 성균관대학교 통계학과 교수

Kohonen네트워크, Carpenter와 Grossberg네트워크가 있다. Kohonen네트워크(1990)에는 SOFM(Self Organization Feature Maps), LVQ(Learning Vector Quantization)알고리즘이 있으며, 최근에 발표된 것으로 Pal et. al.(1993)의 GLVQ(Generalized Learning Vector Quantization) 알고리즘이 있다.

클러스터링의 성능을 평가하는 연구로는 각 패턴과 클러스터의 중심값에 대한 제곱오차를 계산하는 방법으로 Ismail과 Kamel이 Best-first(ABF), First-best(AFB)알고리즘의 성능을 평가 할 때 클러스터의 수를 변화시켜 가면서 k-means알고리즘과 제곱오차(sum of squared error)값을 비교하였다.

또한 Rand(1971)는 클러스터링 방법의 평가기준으로써 "Natural" 클러스터와의 비교, 원시 자료에 잡음(noise)을 추가시켜 구성된 클러스터와의 비교, 그리고 자료에 결측(missing)을 발생시켜 구성한 클러스터와 비교하는 방법을 제안하였다.

본 논문은 하드(Hard or Crisp) 클러스터링으로 국한하였을 때, 신경망을 이용한 클러스터링 알고리즘과 통계적 클러스터링방법을 비교하여 그 결과에 따른 성능을 평가한다. 비교 대상 알고리즘으로는 신경망 클러스터링에서 GLVQ알고리즘과 통계적 클러스터링 방법에서 k-means방법을 선정하였다. 클러스터링의 성능 비교방법으로는 Rand가 제안한 방법을 이용하였다.

2. 신경망클러스터링 방법

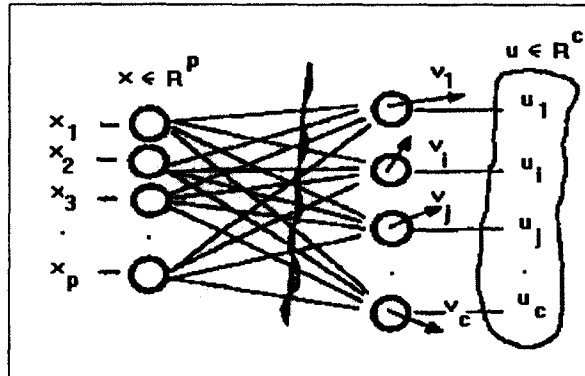
2.1 학습 벡터 양자화(Learning Vector Quantization)

LVQ는 입력패턴 X 를 c 개의 출력층으로 클러스터링 하는 네트워크로 SOFM네트워크의 형태와 다른 점은 SOFM은 출력층이 2차원 정방격자인 반면 LVQ의 출력층은 1차원이라는 점에서 차이가 있다. 즉, p 차원의 패턴 N 개를 c 개의 클러스터로 구분을 지으려고 할 때, p 개의 노드를 갖는 입력층과 c 개의 노드를 갖는 1차원의 출력층으로 구성된다. LVQ네트워크의 구성은 (그림 1)과 같다. 여기서 입력층은 $X = \{x_1, x_2, x_3, \dots, x_N\} \subset R^p$ 으로 구성된 패턴벡터가 입력되며, c 개의 출력층은 $V = \{v_1, v_2, v_3, \dots, v_c\}$, $v_r \in R^p$ ($1 \leq r \leq c$)로 구성되고 여기서 출력층에 연결된 연결강도 v_r 이 클러스터의 중심값이 된다. 학습은 x_i 와 v_r 의 거리를 계산하여 그 거리가 최소가 되는 i 번째 출력노드가 승자노드가 되며 이 승자노드에 연결된 연결강도 v_i 에 대해서만 식 (2.3)에 나와있는 학습규칙에 의하여 학습을 하게 된다.

$$v_{i,t} = v_{i,t-1} + \alpha_t (x_k - v_{i,t-1}) \quad (2.1)$$

여기서, $v_{i,t}$ 는 반복 t 시점에서 승자노드와 연결된 연결강도 벡터.

x_k 는 k 번째 패턴벡터. ($1 \leq k \leq N$)



(그림 1) Kohonen의 LVQ 네트워크 : p 차원의 패턴을 c 개의 클러스터로 구분하는 네트워크

(그림 1)에서 u 는 $c \times N$ 행렬로 이루어지며, 패턴 $x_k(1 \leq k \leq N)$ 가 c 개의 출력노드 중 i 번째 노드인 u_i 에 연결되어 있는 연결강도 v_i 와의 거리가 최소가 된다는것을 의미하는 것으로

$$\|x_k - v_i\| \leq \|x_k - v_j\|, \quad (1 \leq k \leq N, 1 \leq j < c, j \neq i)$$

의 조건을 만족하면 $u(i,k)=1$, 만족하지 않으면 $u(j,k)=0$ 의 값을 갖고, 최종적으로 학습이 완료되고 나면 각 패턴은 u 행렬의 요소 중에서 1의 값을 갖는 패턴이 u_i 에 속하게 된다. 학습은 $t-1$ 시점에서의 연결강도와 t 시점에서의 연결강도의 차이가 허용 오차 안에 들어올 때까지 계속되거나 또는 t 가 최대 반복횟수에 도달할 때 까지 계속 된다.

2.2 Generalized 학습 벡터 양자화(GLVQ)

LVQ알고리즘은 초기 연결강도 $v_{i,0}$ 에 의해서 클러스터링 결과가 많은 영향을 받기 때문에 Pal et. al.은 학습을 승자노드 뿐 만 아니라 승자가 아닌 노드도 같이 학습을 시키는 GLVQ클러스터링 알고리즘을 제안하였다. GLVQ네트워크에서의 학습규칙은 입력 패턴 x 와 출력노드간의 거리에 대한 가중값을 준 손실함수(loss function) L_x 를 식 (2.4)와 같이 정의하고 이 손실함수를 최소화시키는 방법으로 학습규칙을 유도하였다(1993).

$$L_x = \sum_{i=1}^c g_{ir} \| \mathbf{x} - \mathbf{v}_r \|^2 \quad (2.2)$$

여기서, $g_{ir} = \begin{cases} 1 & , \text{승자노드일 경우} \\ \frac{1}{\sum_{j=1}^c \| \mathbf{x} - \mathbf{v}_j \|^2} & , \text{승자노드가 아닌 경우} \end{cases}$

로써 i 번째 패턴에 대한 가중값.

여기서 승자노드일 경우에는 L_x 를 승자노드와 연결된 연결강도 \mathbf{v}_i 로 미분한 $\Delta_{\mathbf{v}_i} L_x$ 를 계산하였고, 승자노드가 아닐 경우에는 승자노드가 아닌 노드에 연결된 연결강도 \mathbf{v}_j 로 미분한 $\Delta_{\mathbf{v}_j} L_x$ 를 계산함으로써 식 (2.5)와 식 (2.6)으로 유도하였으며, 이 학습규칙에 의하여 학습을 한다. (유도식은 참고문헌 [13]참조)

$$\mathbf{v}_{i,t} = \mathbf{v}_{i,t-1} + \alpha_t (\mathbf{x}_k - \mathbf{v}_{i,t-1}) \frac{D^2 - D + \| \mathbf{x}_k - \mathbf{v}_{i,t-1} \|^2}{D^2}$$

승자노드일 경우의 학습 규칙 (2.3)

$$\mathbf{v}_{r,t} = \mathbf{v}_{r,t-1} + \alpha_t (\mathbf{x}_k - \mathbf{v}_{r,t-1}) \frac{\| \mathbf{x}_k - \mathbf{v}_{i,t-1} \|^2}{D^2}$$

승자노드가 아닐 경우의 학습 규칙 (2.4)

$$\text{여기서, } D = \sum_{i=1}^c \| \mathbf{x} - \mathbf{v}_i \|^2, \quad 1 \leq k \leq N, \quad 1 \leq r \leq c$$

\mathbf{v}_i 는 승자노드와 연결된 연결강도 벡터

GLVQ의 특징으로는 연결강도의 초기값에 영향을 받지 않고 클러스터링을 할 때 오분류되는 패턴의 수가 적어지며, 자료의 구분이 확실한 경우에 학습이 완료된 후 각 클러스터의 중심값과 실제 중심값과의 오차도 매우 근소하게 나타난다는 점에서

LVQ보다는 더 좋은 알고리즘으로 평가되고 있다. GLVQ는 x_i 의 승자노드에 연결된 연결강도 v_i 뿐만 아니라 승자가 아닌 다른 노드의 연결강도까지도 조정을 해 주기 때문에 반복이 계속되면서 v_i 가 아닌 다른 노드들도 승자가 될 수 있고, 반복이 완료된 후의 모든 연결강도는 입력된 패턴을 각 클러스터별로 포함하는 중심값이 될 수 있다.

실제로 Pal et. al.은 GLVQ를 제안한 논문에서 Fisher의 IRIS데이터를 가지고 반복 횟수와 학습률을 변화시켜 가면서 LVQ와 GLVQ로 클러스터링 한 결과를 비교하였는데 150개의 IRIS데이터중 GLVQ는 오분류된 패턴이 17개로 일정한 반면에 LVQ인 경우에는 100개에서 17개의 오분류 패턴을 나타낸 결과를 제시하였다.

3. 성능평가 기준

3.1 평가 기준

클러스터링의 결과를 비교하는 방법으로 N 개의 패턴 $X = \{x_1, x_2, \dots, x_N\}$ 이 있다고 가정한다.

이 패턴을 Y 라는 방법으로 클러스터링 한 결과를 $Y = \{y_1, y_2, \dots, y_{k1}\}$ 라 하고, Y' 이라고 하는 방법으로 클러스터링 한 결과를 $Y' = \{y'_1, y'_2, \dots, y'_{k2}\}$ 라고 했을 때 Y 와 Y' 에 있는 패턴 중 두개의 패턴 x_i, x_j ($1 \leq i \leq N, 1 \leq j \leq N, i \neq j$)을 추출한 다음 이 두개의 패턴이 Y 와 Y' 에서 모두 한 클러스터 안에 속해 있을 경우에는 Y 와 Y' 은 x_i 와 x_j 를 유사한 패턴으로 간주한다. 또한 서로 다른 클러스터에 속해 있을 경우에는 Y 와 Y' 은 x_i 와 x_j 를 서로 유사하지 않은 것으로 간주하기 때문에 이 두가지 경우가 Y 와 Y' 의 유사성을 결정하는 척도가 된다. 그러나 x_i, x_j 가 Y 에서는 같은 클러스터에 속해 있어서 서로 유사하다는 결과로 나왔으나 Y' 에서는 서로 다른 클러스터에 분리되어 속해 있어서 서로 유사하지 않다는 결과가 나왔거나 아니면 이와 반대의 결과가 나왔다면 Y 와 Y' 의 클러스터링 방법에는 차이가 있다고 할 수 있다. 따라서 N 개의 패턴중 두개의 패턴을 추출한 경우의 수를 모두 비교한다면 Y 와 Y' 의 클러스터링의 유사성을 비교할 수가 있다. 즉, N 개의 패턴에서 두개의 패턴 x_i 와 x_j 를 추출하는 경우의 수는 ${}_N C_2$ 이고, x_i, x_j 가 서로 같은 클러스터에 있는 경우의 수와 서로 다른 클러스터에 분리되어 포함된 경우의 수를 더한 값의 비를 계산한다면 Y 와 Y' 의 유사성을 비교할 수가 있다. 이를 식으로 표현하면 식 (3.1)과 같다 (Rand, 1971).

$$c(Y, Y') = \frac{\left[\binom{N}{2} - \frac{1}{2} \left(\sum_i \left(\sum_j n_{ij} \right)^2 + \sum_j \left(\sum_i n_{ij} \right)^2 \right) - \sum_i \sum_j n_{ij}^2 \right]}{\binom{N}{2}} \quad (3.1)$$

여기서, n_{ij} 는 Y와 Y'으로 구성된 클러스터링 결과 중 Y의 i 번째 클러스터와 Y'의 j 번째 클러스터 안에 모두 포함되어 있는 동일한 패턴의 수를 의미한다. 식 (3.1)에서 $c=0$ 이면 Y와 Y'의 결과는 전혀 유사하지 않은 것이며, $c=1$ 일 경우에는 두개의 클러스터의 결과가 동일한 것임을 나타낸다. 위 식에서 표현된 c 통계량의 특징으로는 첫째, c 통계량은 두 클러스터의 유사성을 나타내고, 둘째, $(1-c)$ 는 두 클러스터의 비유사성을 의미하며, 셋째, X가 어느 특정 분포를 갖고 있을 경우에 c 는 확률변수이다.

3.2 데이터 선정

3.2.1 데이터

본 논문에서 이용한 자료는 대학교 재학생인 100명의 학생을 무작위로 추출하여 0-9의 숫자 중 하나씩 쓰게 하여 수집된 100개의 패턴을 스캐너로 그 이미지를 입력 받아서 특징추출(feature extraction)을 하였고, 이 특징추출한 자료를 이용하여 클러스터링을 하였다. 비교분석의 편의상 0-9까지의 숫자를 쓸 때 각 숫자마다 학생 10명씩을 할당하여 쓰게 하였다.

3.2.2 특징추출

16×16 격자에 숫자에 대하여 이진화(binanzed)된 원시 자료를 받고, 이 숫자를 좌, 우, 상, 하로 여백을 제거한 다음 이 패턴을 10×10으로 스케일링을 한다. 그 다음 좌, 우, 상, 하에 한 행 또는 열을 추가시켜 12×12로 만든다. 이 12×12로 스케일링한 자료를 좌측 상단에 있는 격자부터 한 격자씩 탐색해 가면서 '1'이 있을 경우에 그 격자를 중심으로 좌, 우, 상, 하, 좌상, 좌하, 우상, 우하의 방향에 '1'이 있으면 1을 증가시키는 방법이다. 따라서 모든 방향에 '1'이 있을 경우에는 최고 9의 값을 갖는다. 12×12의 격자가 모두 탐색되고 나면 (그림 2c)와 같은 자료를 얻을 수 있다. 이 자료를 가로, 세로 각각 3개씩의 격자를 묶어서 하나의 블록으로 만들고, 이 블록 안의 값을 모두 더한 다음에 10으로 나누어 준 자료가 입력패턴벡터가 된다.

(그림 2)은 100개의 패턴 중 숫자 '0'에 대한 특징추출 과정에서 얻은 자료의 예를 보여주고 있다.

```

. . . . . 1 1 1 1 .
. . . . 1 1 . . 1 .
. . . 1 1 . . . 1 1
. . 1 1 . . . . . 1
. 1 1 . . . . . . 1
. 1 . . . . . . 1 1
1 1 . . . . . . 1 1
1 . . . . . . 1 1 .
1 . . . . . . 1 1 .
1 . . . . . 1 1 . .
1 1 . . 1 1 1 . . .
. 1 1 1 1 1 . . . .
    
```

(a) 원시자료

```

. . . . . . . . . . .
. . . . . . . 1 1 1 1 . .
. . . . . . . 1 1 . . 1 . .
. . . . . 1 1 1 . . . 1 1 .
. . . . 1 1 . . . . . 1 .
. . . . 1 . . . . . . 1 1 .
. . . . 1 1 . . . . . 1 1 .
. . . . 1 . . . . . . 1 1 .
. . . . 1 . . . . . 1 1 1 . .
. . . . 1 1 . . 1 1 1 . . .
. . . . 1 1 1 1 1 . . . .
. . . . . . . . . . .
    
```

(b) 12×12 스케일링

```

0 0 0 0 0 1 2 3 3 2 1 0
0 0 0 0 1 3 4 4 4 3 2 0
0 0 1 2 4 5 5 4 5 5 4 1
0 1 3 4 5 4 3 1 2 4 4 2
0 2 4 5 4 2 1 0 2 5 5 3
1 4 5 4 1 0 0 0 2 5 5 3
2 4 4 2 0 0 0 1 4 6 5 2
3 4 4 1 0 0 1 3 6 6 4 1
3 4 4 1 1 2 4 5 6 4 2 0
2 4 5 4 4 5 6 5 4 2 1 0
1 3 4 4 4 5 5 3 1 0 0 0
0 1 2 3 3 3 2 1 0 0 0 0
    
```

(c) 12×12 누적 자료

```

0.1 1.6 3.4 1.8
2.0 2.9 1.1 3.6
3.2 0.7 3.0 3.0
2.2 3.5 2.7 0.3
    
```

(d) 4×4 입력자료

(그림 2) 16×16 자료를 4×4로 특징추출 예 : 숫자 '0'

4. 성능평가 실험

이 장에서는 Rand비교 방법에 의해서 신경망 클러스터링과 통계적 방법의 클러스터링의 성능을 평가한다. Rand방법에 의한 클러스터링의 성능 평가는 최적의 클러스터 개수를 결정하는 것은 배제되었기 때문에 본 논문에서 이용한 자료 중 IRIS데이터는 클러스터의 수를 3개로 하고, 숫자자료는 클러스터의 수를 10개로 하여 비교하였다.

4.1 "Natural" 클러스터와의 비교

본 절에서는 "Natural" 클러스터를 알고 있다는 가정 하에 비교하는 방법이기 때문에 0-9의 숫자 패턴 자료를 이용하여 3.1절의 c 통계량을 계산하였다.

0-9까지의 숫자 자료를 비교하는 데 있어서는 실제 Y 를 자료의 성질에 따라 구분되는 기준 클러스터로 총 10개의 패턴을 각 숫자별로 구분하여 구성한 클러스터이다. 본 논문에서는 분석의 편의상 각 숫자별로 10명의 필기자를 할당하여 쓰게 하였으므로 10개의 숫자 '0'은 클러스터 0에, 10개의 숫자 '1'은 클러스터 1에 포함시키고, 마지막으로 숫자 '9'는 클러스터 9에 포함시킨 10개의 클러스터 집합이다.

$$Y_{NUM} = \{ (x_{0,0}, x_{1,0}, x_{2,0}, \dots, x_{9,0}), (x_{0,1}, x_{1,1}, x_{2,1}, \dots, x_{9,1}), \dots, \dots, (x_{0,9}, x_{1,9}, x_{2,9}, \dots, x_{9,9}) \}$$

여기서, $x_{i,j} \in R^p$: 각 숫자별 i 번째($0 \leq i \leq 9$) 필기자가 쓴 숫자 j .

$$Y'_{NUM} = \{ (x'_{0,1}, x'_{0,2}, \dots, x'_{0,n_1}), (x'_{1,1}, x'_{1,2}, \dots, x'_{1,n_2}), \dots, \dots, (x'_{9,1}, x'_{9,2}, \dots, x'_{9,n_9}) \}$$

$$\text{여기서, } \sum_{i=1}^9 n_i = 100,$$

$x'_{i,j}$: 0-9숫자자료를 어느 특정 클러스터링방법으로

클러스터를 구성하였을 때 i 번째 클러스터에 포함된 j 번째 패턴 벡터.

Y' 은 GLVQ와 k -means 클러스터링방법으로 구성된 10개의 클러스터로써 $Y'_{NUM, GLVQ}$, $Y'_{NUM, k\text{-means}}$ 을 각각 GLVQ, k -means 알고리즘으로 구성된 집합이라고 한다.

두 자료를 GLVQ에 의해서 클러스터링을 하였을 경우 학습률 α_0 는 0.4, 반복횟수는 2500번으로 하여 클러스터링을 하였다.

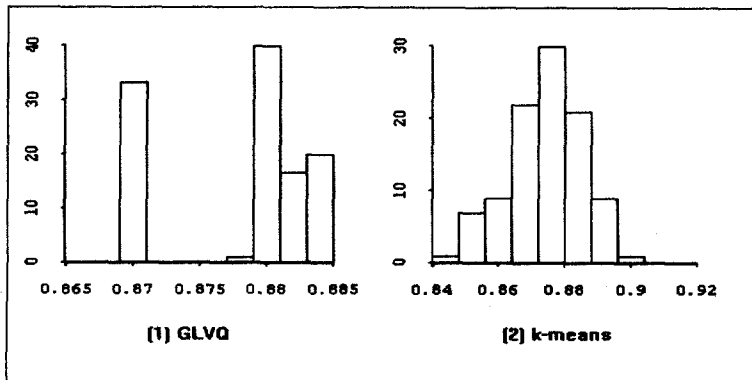
<표 1>은 3.2.2절에서 특징추출한 자료를 각각의 방법으로 클러스터링 하여 "Natural" 클러스터와 비교한 결과인 데 이때 GLVQ 방법과 k -means 방법은 100번 시행하여 c 통계량의 평균과 표준편차를 비교한 결과이다. (그림 3)은 각 특징추출방법

에 의하여 얻은 자료를 GLVQ와 k-means알고리즘으로 클러스터링 결과에 대한 c통계량의 분포를 히스토그램으로 표현한 것이다.

<표 1> "Natural" 클러스터 :
GLVQ와 k-means클러스터링에 대한 c통계량의 비교 결과

| 클러스터링 방법 | | 4×4 |
|----------|---------|----------|
| c의 평균 | GLVQ | 0.878133 |
| | k-means | 0.873909 |
| c의 표준편차 | GLVQ | 0.000035 |
| | k-means | 0.000126 |

주) 4×4 : 16×16 자료를 4×4로 특징추출한 자료



(그림 3) "Natural" 클러스터와의 비교 : 16×16 자료를 4×4로
특징추출한 자료에 대한 c통계량의 히스토그램

4.2 자료에 오류가 있는 클러스터와의 비교

본 절에서는 자료의 오류 또는 잡음이 들어 왔을 때 어느 클러스터링 방법이 잡음의 영향에 민감하게 반응하는가 하는 정도로써 그 성능평가를 하는 방법이다.

IRIS데이터를 이용할 경우에는 150개의 자료의 각 변수에 난수를 발생시킨 오류를 더하여 클러스터링을 비교하였다. 0-9까지의 숫자 자료에서도 역시 Y를 4.1절에서 각 클러스터링 방법으로 구성한 Y'집합이라고 하고, Y'은 원시 자료에 10%의 잡음을 추가하여 3.2.2절의 특징추출 방법을 이용하여 100개의 패턴을 새로 구성하여 GLVQ와 k-means방법으로 클러스터링을 하였다. 여기서 10%의 잡음은 원시 자료에서 '1'의 수를 모두 합한 총수의 10%에 해당하는 수(=n이라고 한다)를 의미하며 이

'1'을 원래 패턴에 임의의 셀 위치에 n 개 추가하였다.

$Y_{NUM} = Y'_{Natural\ NUM}$: "Natural"클러스터를 비교하기 위하여 각 클러스터링방법에 의해서 구성된 클러스터의 집합.

$$Y'_{NUM} = \{ (x'_{0,1}, x'_{0,2}, \dots, x'_{0,n_1}), (x'_{1,1}, x'_{1,2}, \dots, x'_{1,n_2}), \dots, \dots, (x'_{9,1}, x'_{9,2}, \dots, x'_{9,n_9}) \}$$

여기서, $\sum_{i=0}^9 n_i = 100,$

x' : 10%의 잡음을 추가하여 특징추출한 패턴 벡터.

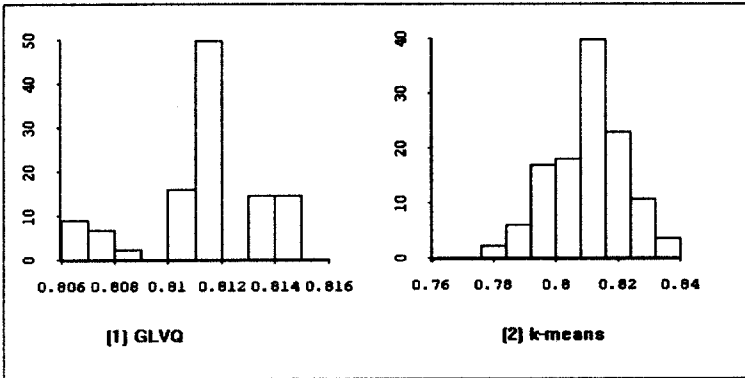
IRIS데이터와 각 특징추출 별로 얻은 자료에 대한 비교 결과는 <표 2>와 같으며, c 통계량에 대한 히스토그램은 (그림 4)에 나와 있다.

<표 2> 자료에 오류를 추가 :

신경망과 통계적 클러스터링에 대한 c 통계량의 비교

| 클러스터링 방법 | | 4×4 |
|----------|---------|----------|
| c의 평균 | GLVQ | 0.811303 |
| | k-means | 0.810192 |
| c의 표준편차 | GLVQ | 0.000005 |
| | k-means | 0.000140 |

주) 4×4 : 16×16 자료를 4×4로 특징추출한 자료



(그림 4) 자료에 오류가 있는 클러스터의 비교 :

16×16 자료를 4×4로 특징추출한 자료에 대한 c통계량의 히스토그램

4.3 자료에 결측이 있는 클러스터와의 비교

이 방법은 자료를 우선 N_1 로 클러스터링을 한 결과와, (N_1+N_2) 개로 클러스터링 한 결과가 일치하는가를 비교하는 것으로 비교 대상은 처음 클러스터링을 한 N_1 개의 자료이다. 즉, N_1 개로 클러스터링을 한 결과와 (N_1+N_2) 개로 클러스터링을 했을 때 이 중 N_1 개의 자료가 처음 클러스터링 한 결과와 얼마나 일치하는가를 비교하는 것이다.

0-9까지의 숫자 자료는 100명의 필기자에서 얻은 100개(= N_1)의 패턴을 GLVQ와 k-means알고리즘으로 구성한 10개의 클러스터의 집합이라고 한다.

$$Y_{NUM} = \{ (x_{0,1}, x_{0,2}, \dots, x_{0,n_1}), (x_{2,1}, x_{2,2}, \dots, x_{2,n_2}), \dots, (x_{9,1}, x_{9,2}, \dots, x_{9,n_9}) \}$$

여기서, $\sum_{i=0}^9 n_i = 100$

Y'은 Y를 구성한 패턴에 30(= N_2)명의 필기자를 추가하여 총 130명의 필기자에게

서 얻은 130개(=N₁+N₂)의 패턴을 GLVQ와 k-means방법에 의하여 10개의 클러스터로 클러스터링을 한 집합이다.

$$Y'_{NUM} = \{ (x'_{0,1}, x'_{0,2}, \dots, x'_{0,n'}), (x'_{1,1}, x'_{1,2}, \dots, x'_{1,n'}), \dots, \dots, (x'_{9,1}, x'_{9,2}, \dots, x'_{9,n'}) \}$$

여기서, $\sum_{i=0}^9 n'_i = 130,$

Y'_{NUM} 는 Y_{NUM} 중 Y_{NUM} 를 구성한 자료만을 추출하여 구성한 집합이라고 했을 때 $c(Y_{NUM}, Y'_{NUM})$ 를 계산하여 비교한다.

$$Y''_{NUM} = \{ (x''_{0,1}, x''_{0,2}, \dots, x''_{0,n''}), (x''_{1,1}, x''_{1,2}, \dots, x''_{1,n''}), \dots, \dots, (x''_{9,1}, x''_{9,2}, \dots, x''_{9,n''}) \}$$

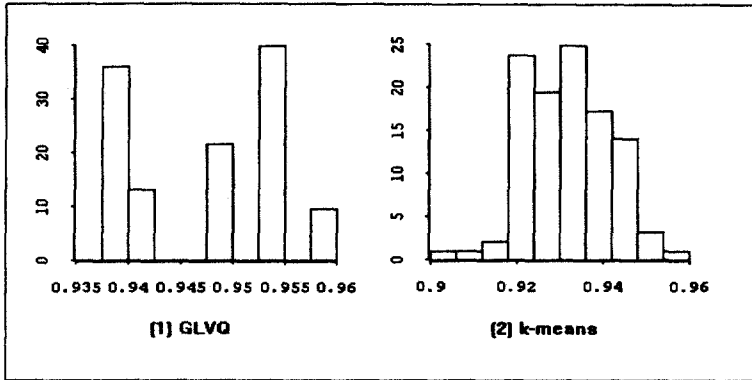
여기서, $x'' \in Y_{NUM}$
 $\sum_{i=0}^9 n''_i = 100,$

결측이 있는 자료를 클러스터링한 결과는 <표 3>과 같고 c통계량의 분포는 (그림 5)에 나와있는 히스토그램과 같다.

<표 3> 결측자료 :
 신경망과 통계적 클러스터링에 대한 c통계량의 비교

| 클러스터링 방법 | | 4×4 |
|----------|---------|----------|
| c의 평균 | GLVQ | 0.947161 |
| | k-means | 0.931827 |
| c의 표준편차 | GLVQ | 0.000045 |
| | k-means | 0.000100 |

주) 4×4 : 16×16 자료를 4×4로 특징추출한 자료



(그림 5) 자료에 결측이 있는 클러스터의 비교 :
 16×16 자료를 4×4로 특징추출한 자료에 대한 c통계량의 히스토그램

5. 결론

본 논문은 클러스터링 방법 중 신경망에 의한 방법과 통계적인 방법을 비교하여 성능을 평가하는 데 그 목적이 있다. IRIS데이터와 필기체 숫자에 대한 자료를 신경망에 의한 클러스터링과 통계적인 방법에 의하여 클러스터링을 한 후 Rand c통계량을 적용하여 비교한 결과에 국한시켰을 때 신경망클러스터링의 GLVQ알고리즘이 통계적 방법의 k-means방법보다는 매 반복시행을 할 때마다 오분류되는 패턴이 약간 적게 나타났다. 또한 매 반복 시행을 할 때마다 오분류되는 패턴의 편차도 상대적으로 크게 나타나지 않았다.

참 고 문 헌

1. 김대수 (1994), "신경망 이론과 응용 (I),(II)", 하이테크정보.
2. 이성환 (1994), "패턴인식의 원리 (I),(II)", 홍릉과학출판사.
3. 최병수 (1990), "군집분석을 위한 통계전문가 시스템의 구현", 성균관대학교 박사학위 논문.
4. Byron J.T. Morgan (1984), "Elements of Simulation", Greate Britain at the Univ. Press Cambridge.
5. D.J. Hand, F. Daly, A.D. Lunn, K.J. McConway, E. Ostrowski (1994) , "A Handbook of Small Data Sets", Chapman & Hall.
6. Geoffrey H. Ball, David J. Hall (1967), "A Clustering Technique for Summarizing Multivariate Data", Behavior Science Vol. 12, pp.153-155.

7. Helge Ritter, Thomas Martinets, Klaus Schulten (1992), "*Neural Computation and Self-Organizing Maps*", Addison-Wesley Publishing Co. Inc.
8. John Hertz, Anders Krogh, Richard G. Palmer (1991), "*Introduction to the Theory Neural Computation*", Addison-Wesley Publishing Co. Inc.
9. J. R. Slagle, C. L. Chang, R. C. T. Lee (1974), "Experiments with some Cluster Analysis Algorithm", *Pattern Recognition*, Vol. 6, pp.181-187.
10. M.A.Ismail and M.S.Kamel (1989), "Multidimensional Data Clustering Utilizing Hybrid Search Strategies", *Pattern Recognition*, Vol. 22, No. 1, pp.75-89.
11. Mark S. Aldenderfer, Roger K. Blashfield (1984), "*Cluster Analysis*", Sage Publications Inc.
12. Michael R. Anderberg (1973), "*Clustering Analysis for Applications*", Academic Press, N.Y. and London.
13. Nikhil R. Pal, James C. Bezdek, Etric C.K.Tsao (1993), "Generalized Clustering Networks and Kohonen's Self-Organizing Scheme", *IEEE Trans. on Neural Networks*, Vol. 4, No. 4, pp.549-557.
14. R. Dubes and R.K.Jain (1976), "Clustering Techniques: the user's dilemma", *Pattern Recognition*, Vol. 8, pp.247-260.
15. R.A. Jarvis, Edward A. Patrick (1973), "Clustering Using a Similarity Measure Based on Shared Near Neighbors", *IEEE Trans. on Computers*, Vol. C-22, No. 11, Nov., pp.1025-1034.
16. Richard A. Johnson, Dean W. Wichern (1992), "*Applied Multivariate Statistical Analysis*", 3rd Ed., Prentice Hall.
17. Robert Schalkkoff (1992), "*Pattern Recognition -statistical structureal and approaches*", John and Wiley & Sons, Inc.
18. Shunichi Shimoji, Sukhan Lee (1994), "Data Clustering with Entropical Sheduling", *International Joint Conference on Neural Network*, Vol. 4, pp.2423-2428.
19. Shokri Z. Selim, M.A. Ismail (1984), "k-means-Type Algorithms : A Generalized Convergence Theorem and Characterization of Local Optimality", *IEEE Trans. on Pattern Analysis and Machine Intelligence*, Vol. PAMI-6, No. 1, Jan, pp. 81-87.
20. T. Kohonen (1990), "The Self-Organizing Map", *Proc. IEEE*, Vol. 78, No. 9, pp.1464-1480.

21. Terence D. Sanger (1989), "Optimal Unsupervised Learning in a Single-Layer Linear Feedforward Neural Network", *Neural Networks*, Vol. 2, p.459.
22. William M. Rand (1971), "Objective Criteria for the Evaluation of Clustering Methods", *Journal Am. Stat. Assoc.*, Vol. 66, pp.846-850.
23. Yoh-Han Pao (1989), "*Adaptive Pattern Recognition and Neural Network*", *Addison-Wesley Publishing Co. Inc.*

Comparative Analysis of Neural and Statistical Clustering

Jun Young Min¹⁾, Mun Yul Huh²⁾

Abstract

Three popular techniques for clustering algorithms are statistical methods, ISODATA algorithm and neural networks. In these methods, neural networks include the Kohonen's SOFM(Self Organize Feature Maps) and ART(Adaptive Resonance Theory) algorithm. Kohonen's SOFM is corresponding to the k-means algorithm in statistical clustering.

This paper evaluates the clustering performance of a neural network and a statistical method. Algorithms which are used in this paper are the GLVQ(Generalized Learning vector Quantization) for a neural method and the k-means algorithm for a statistical clustering method. For comparison of two methods, we calculate the Rand's c statistics. As a result, the mean of c value obtained with the GLVQ is higher than that obtained with the k-means algorithm, while standard deviation of c value is lower. Experimental data sets were the Fisher's IRIS data and patterns extracted from handwritten numerals.

1) Sangji Junior College, Wonju

2) Department of Statistics, Sung Kyun Kwan University, 3-53, Myungryun-dong Chongri-Ku, Seoul 110-745, Korea.