

# 植被 Data에 對한 Information과 Entropy 理論의 實用研究

朴 勝 太  
(全北大 師大 生物教育科)

## A Study on the Presentation of Idea in Information and Entropy Theory in Vegetation Data

Park, Seung Tai

(Dept. of Biology Educ., Jeonbug Nat'l. Univ.)

### ABSTRACT

This study is concerned with some methods and applications, used as a basis on information and entropy analysis of vegetation data.

These methods are adopted for the evaluating the effect of sampling intensity on information, which represents the departure of observed variable from standard component. Classes on the data matrix are calculated by using marginal dispersion array for rank and weighting information program. Finally the information and entropy are computed by applying seven options.

On the application of vegetation studies, two models for cluster analysis and analysis of concentration are explained in detail. Cluster analysis is based on use of equivocation information and Rajski's metrics. The analysis of concentration utilizes coherence coefficient being transformed values, which has been adjusted from blocks and entropy values.

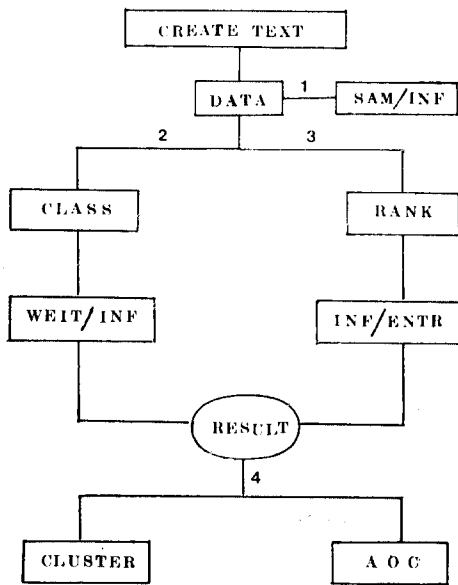
The relationship between three vegetation clusters and four stands of Naejangsan data is highly significant in 79% of total variance. Cluster A relatively tends to prefer north side, and cluster C south side.

### 結 論

生態學의 複雜한 資料(multidimensional data)에 대하여 information理論을 適用하는 것은 modeling을 하거나 統計的인 資料의 單純化(simplification) 및 資料 變數에 의한 傾向性(trend)을 찾는 것과는 상당한 차이가 있다(Legendre and Legendre, 1983; Pielou, 1977; Gauch, 1982; Orloci, 1978).

生態學의 研究에서 種의 複雜한 關係, 層化(stratification) 및 季節的인 變化등에 대한 時間的 또는 空間的인 變化를 定量的으로 叙述할 때 研究條件에 따라 個體數 또는 生體量을 나타내는 evenness와 richness에 의해 類型分析(pattern analysis), 區分分析(discrimination

\* 이 論文은 85年 尖端科學 技術分野 研究費 支援에 의한 것임.



**Fig. 1.** Flowchart of information/entropy analysis. SAM/INF=sampling information; CLASS=dispersion array program; RANK=rank program; WEIT/INF=weighting information program; INF/ENTR=information and entropy program; AOC=analysis of concentration.

data sampling의 效率의 計算을 위하여 最適標集의 intensity( $n$ )를 정하고 이에 따른 quadrat數와 크기를 결정할 수 있게하고 ② weighting information을 할때 資料의 階級(class)을 임의로 조절하여 값을 算出하고 ③ 방대한 資料에 대하여 RANK program으로 rank를 정하여 줄일 수 있게 하고 ④ 計算된 information 값을 이용하여 結果를 解析할 수 있도록 clustering 方法과 集中度分析(AOC)으로 information/entropy값을 이용하므로서 生態學 資料分析 및 解析을 용이하게 할 수 있게 한 것이다.

方 法

**Information理論** Information은 意思傳達理論을 기초로 하여 발달되었으며 Fisher(1948)는 logarithm 函數를 이용하여 確率을 information ( $I$ )으로 나타 냈다.

$$I(E) = -\ln P(E) : E=1, 2, \dots, s \dots \dots \dots (1)$$

이때 事象( $E$ )은  $s$ 번 試行할 수 있으며 information 값은 항상 +값(positive)이며 다음과 같이 變形할 수 있다.

$$I(E) = -2 \sum_{E=1}^s \ln \Pi(E) \dots \dots \dots (2)$$

또한 Kullback(1959)은 카이제승( $\chi^2$ )變量으로 이를 변형시켜 information과  $\chi^2$ 關係를 究明하였다.

$$\chi^2 = 2I \dots \dots \dots (3)$$

Information은 entropy와 divergence로 區分할 수 있는데(Orloci, 1978), entropy는 어떤 分

analysis) 및 豫見分析(predictive analysis)을 하게 된다(Greig-Smith, 1957; Orloci, 1972; MacArthur, 1965; Feoli *et al.*, 1984; Shannon, 1948).

Information理論은 意思傳達理論(communication theory)을 基礎로 하여 發達된 것으로 entropy 分析과 divergence 分析으로 區分되며, Renyi(1961)는 entropy에 대하여 disorder量으로  $H^\alpha$ 를 써서  $\alpha$ 의 次數에 따라 Shannon의 entropy函數와 Simpson指數를 算出할 수 있게 했으며 Kullback(1959)은 divergence에 대하여  $2(\text{Information}) = \chi^2$ 로 나타내어 最小識別統計(minimum discrimination information statistics)의 基礎를 마련 하므로써 資料에서 期待値와 觀側値와의 關係에 대한 divergence를 산출하여 分散分析(dispersion analysis), 多樣性分析(diversity analysis) 및 豫見分析(predictive analysis) 등을 가능하게 했다(Fisher, 1948; Margalef, 1958; McIntosh, 1967; Orloci, 1972).

본 研究에서는 Fig. 1에서와 같이 information/entropy理論을 基礎로 하여 ① Data

布에 대한 disorder를 logarithm 함수로 표시했고, Renyi(1961)는 disorder ( $H^\alpha$ )를 다음과 같이 나타냈다.

$$H^\alpha = \frac{\ln \sum_{E=1}^S P(E)^\alpha}{1-\alpha} \dots\dots\dots(4)$$

위의 식 (4)에서  $\alpha=0$ 일때  $H^0=\ln S$ 로  $H_{max}$ 가 되며,  $\alpha \rightarrow 1$  일때 L'Hopital 法則을 이용하여 풀면  $H^1 = -\sum_{E=1}^S P(E)\ln P(E)$ 가 되는데 이는 Shannon index이며,  $\alpha=2$ 일때  $H^2 = -\ln \sum_{E=1}^S P(E)^2$ 로 Simpson index이다.

또한 entropy는 joint entropy ( $I(\mathbf{a}, \mathbf{b})$ ), mutual entropy ( $I(\mathbf{a}; \mathbf{b})$ ) 및 equivocation entropy ( $E(\mathbf{a}; \mathbf{b})$ )등으로 區分되는데 이들의 計算은 다음 program에서 자세히 설명된다.

Divergence는 觀側値와 期待値에 대한 確率이나 頻度를 이용하여 Kullback(1959)는 다음과 같이 나타냈다.

$$\chi^2 = 2I(\mathbf{P}; \mathbf{P}^0) \dots\dots\dots(5)$$

위의 식 (5)에서 觀側値의 確率은  $P = \{P(1)P(2)\dots\dots P(s)\}$ 이고 期待値의 確率은  $P^0 = \{P^0(1) P^0(2)\dots\dots P^0(s)\}$ 이며 divergence를 disorder ( $H$ )로 나타낼 때 다음과 같다.

$$H(\mathbf{P}; \mathbf{P}^0) = H(\mathbf{P}^0) - H(\mathbf{P}) \\ = \sum_{E=1}^S P(E)\ln P(E)/P^0(E) \dots\dots\dots(6)$$

또한 頻度( $F$ )를 이용할 때에도 觀側頻度( $\mathbf{F}$ )와 期待頻度( $\mathbf{F}^0$ ) 計算할 수 있다.

$$\chi^2 = 2I(\mathbf{F}; \mathbf{F}^0) \dots\dots\dots(7)$$

이때  $\mathbf{F} = n\mathbf{P}$ 이며  $\mathbf{F}^0 = n\mathbf{P}^0$ 로서  $n$ 은 標集크기이며 期待頻度( $\mathbf{F}^0$ )의 計算은 對應되는 두 頻度の 平均을 이용하여 算出한다.

그러나 觀側頻度( $f_{ij}$  또는  $f_{hj}$ )만을 이용하여 對應되는 두 quadrat間의 divergence를 計算할 경우는 다음과 같다.

$$I(\mathbf{A}; \mathbf{B}) = I(\text{quadrat1}; \text{quadrat2}) \\ = \sum_h \sum_j f_{hj} \ln 2f_{ij} / (f_{h1} + f_{h2}) \\ h=1, \dots, s; j=1, 2 \dots\dots\dots(8)$$

이때  $\mathbf{A}$ 와  $\mathbf{B}$ 는 quadrat이며 group으로 區分된 quadrat에서도 計算할 수 있다.

또한 種과 quadrat의 information 計算은 다음과 같이 할 수 있다.

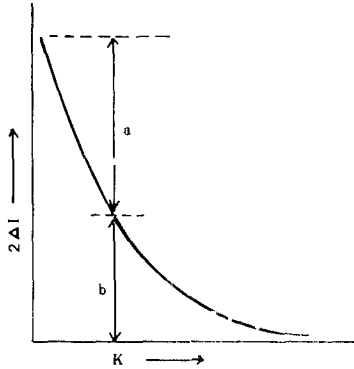
$$I(\text{species}; \text{quadrat}) = \sum_h \sum_j f_{hj} \ln f_{hj} f_{..} / (f_{h.} f_{.j}) \\ h=1, \dots, s; j=1, 2 \dots\dots\dots(9)$$

위의 식(9)에서  $f_{..} = \sum_h \sum_j f_{hj}$ 이고  $f_{h.} = \sum_j f_{hj}$ 이며  $f_{.j} = \sum_h f_{hj}$ 이다.

이때 식(8)과 식(9)의 차이 값을 divergence ( $D$ )로 이용한다.

$$D(1, 2) = I(\mathbf{A}; \mathbf{B}) - I(\text{species}; \text{quadrat}) \\ = \sum_j f_{.j} \ln 2f_{.j} / f_{..} \\ j=1, 2 \dots\dots\dots(10)$$

본 研究에서는 information/entropy 理論만을 이용하여 植被(vegetation) 資料에 적용할 수 있도록 program 했으며 entropy 값으로 clustering 하거나 集中度(concentration)를 算出하여 解析할 수 있게 했다.



**Fig. 2.** Relationship of sampling effort and information. The vertical scale indicates information loss due to not counting entries beyond a maximum K. The symbol a indicates the gain of information and b information loss.

c인 資料에서 觀測된 要素인  $f_{ij}$  및  $(n - f_{ij})$ 와 期待된 要素  $k$  및  $(n - k)$ 의 거리 계산은  $f_{ij} \leq k$  일때만 가능하다.  $n$ 은 한개 quadrat내의 最小數(intensity)이며  $k$ 는 quadrat 最適數로 intensity에 따라 quadrat 크기를 결정하게 된다(Orloci, 1970).

② Class program(CLASS); 植被資料를 分析할 때 array하는 方法이 4가지가 있다. R-array는  $r \times c$  matrix 資料에서  $r$ 變數(predictor)의 分散配列(dispersion array)을 말하며, Q-array는  $c$ 變數(object)를 分散配列하는 것이다. Predictive array는  $r \times c$  matrix에서 resemblance를 산출하여  $r \times r$  또는  $c \times c$  array하는 것이고 diversity array는 R 또는 Q-array하는 方法을 이용하여 多樣度를 산출하는 array法이다(Feoli *et al.*, 1984).

array 할때 class의 限界를 임의로 정하고(본 연구에서는 6개 class; Table 4와5 참고) marginal dispersion을 산출하여 weighting information을 구하거나 RANK program을 운영할 때 이용했다.

③ Rank program (RANK); information에 의한 變數의 rank를 정하는 理論은 다양하나(Orloci, 1976) 본 연구에서는 equivocation information ( $E(a;b)$ )을 이용하여 rank를 결정하고 資料變數를 줄였다.

④ Weighting information program(WEIT/INF); Orloci(1978)의 program을 이용하였다.

⑤ Information entropy program(INF/ENTROPY); 生態學 資料分析에 대하여 Fig. 3과 같이 venn diagram에서 7가지 information을 계산할 수 있도록 program하고 각각을 선택할 수 있게했다.

① information a;  $I(a)$

$$I(a) = -\sum_j f_{hj} \ln f_{hj} / f_h$$

$$= \sum_j f_{hj} \ln f_h - \sum_j f_{hj} \ln f_{hj}$$

$$j=1, 2, \dots, S_h \dots \dots \dots (12)$$

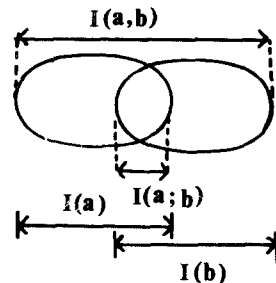
② information b;  $I(b)$

**Program ①** Sampling information program(SAM/INF); 標集方法은 研究戰略과 研究條件에 따라 quadrat 크기와 수를 한정해야하며 특히 quadrat 수는 結果에 상당한 영향을 하게 된다(Greig-Smith, 1957). 본 연구에서는 Fig. 2와 같이 sampling effort에 따라 information 得失에 차이가 있으므로 quadrat의 적합한 數와 quadrat 크기에 한정된 intensity( $n$ )의 결정은 information loss ( $2\Delta I$ )가 最小일 때를 다음과 같은 식으로 정했다.

$$2\Delta I = 2 \sum_{i=1}^r \sum_{j=1}^c f_{ij} \ln f_{ij} / k + (n - f_{ij}) \ln (n - f_{ij}) / (n - k)$$

$$i=1, 2, \dots, r; j=1, 2, \dots, c \dots (11)$$

위의 식(11)에서 species는  $r$ , quadrat은



**Fig. 3.** Venn diagram.

$$\begin{aligned}
 I(\mathbf{b}) &= -\sum_j f_{ij} \ln f_{ij} / f_{i\cdot} \\
 &= \sum_j f_{ij} \ln f_{i\cdot} - \sum_j f_{ij} \ln f_{ij} \\
 & \quad j=1, 2, \dots, S_i \dots \dots \dots (13)
 \end{aligned}$$

③ joint information;  $I(\mathbf{a}, \mathbf{b})$

$$\begin{aligned}
 I(\mathbf{a}, \mathbf{b}) &= -\sum_j \sum_k f_{hj, ik} \ln f_{hj, ik} / f_{h\cdot, i\cdot} \\
 &= \sum_j \sum_k f_{hj, ik} \ln f_{h\cdot, i\cdot} - \sum_j \sum_k f_{hj, ik} \ln f_{hj, ik} \\
 & \quad j=1, \dots, S_h; k=1, \dots, S_i \dots \dots \dots (14)
 \end{aligned}$$

④ mutual information (interaction information);  $I(\mathbf{a}; \mathbf{b})$

$$\begin{aligned}
 I(\mathbf{a}; \mathbf{b}) &= I(\mathbf{a}) + I(\mathbf{b}) - I(\mathbf{a}, \mathbf{b}) \\
 &= \sum_j \sum_k f_{hj, ik} \ln f_{hj, ik} f_{h\cdot, i\cdot} / f_{hj, i\cdot} f_{h\cdot, ik} \\
 & \quad j=1, \dots, S_h; k=1, \dots, S_i \dots \dots \dots (15)
 \end{aligned}$$

⑤ equivocation information;  $E(\mathbf{a}; \mathbf{b})$

$$E(\mathbf{a}; \mathbf{b}) = I(\mathbf{a}, \mathbf{b}) - I(\mathbf{a}; \mathbf{b}) \dots \dots \dots (16)$$

⑥ Rajski metrics;  $d(\mathbf{a}; \mathbf{b})$

$$d(\mathbf{a}; \mathbf{b}) = \frac{E(\mathbf{a}; \mathbf{b})}{I(\mathbf{a}, \mathbf{b})} \dots \dots \dots (17)$$

⑦ coherence coefficient;  $r(\mathbf{a}; \mathbf{b})$

$$r(\mathbf{a}; \mathbf{b}) = |1 - d^2(\mathbf{a}; \mathbf{b})|^{1/2} \dots \dots \dots (18)$$

모든 program은 APPLE BASIC, DOS 3.3을 이용할 수 있게 제작 되었다.

이때의  $\mathbf{a}, \mathbf{b}$ 는 대립되는 變數로서 빈도나 밀도를 나타낸다.

**Clustering (CLUSTER)** Information에 의한 Clustering은 equivocation information과 Rajski metrics를 이용하여 sum of square(SS) algorithm으로 fusion하였고 cluster의 결정은 information SS값을 임의로 선정하여 구분 하였다. Clustering 效率은 區分된 cluster內的 SS(W) information을 計算하고 全體 cluster SS(T) information에서 빼면 cluster間的 SS(B) information이 된다. 이를 이용하여  $\frac{B}{T} \times 100(\%)$ 로 效率을 計算 했다(Sneath and Sokal, 1973).

**集中度分析(AOC)** 植被內에 分布된 種과 環境과의 關係를 正準分析法(canonical analysis)을 이용하여 傾向을 추출하는 方法인데 본 研究에서는 coherence coefficient( $r(\mathbf{a}; \mathbf{b})$ ) 값 ( $I_{pp}$ )을 이용하여 資料를 다음과 같이 adjust ( $A_{pq}$ ) 했고  $N$ 은 식(20)과 같이 補正값이다.

$$A_{pq} = I_{pp} R_{pq} / N \dots \dots \dots (19)$$

$$N = pq / |p + q|^{1/2} \dots \dots \dots (20)$$

위의 식(19)에서  $R_{pq}$ 는 raw data이며  $p=1, \dots, r$ 로서 data matrix에서 種을 나타내며  $q=1, \dots, c$ 는 quadrat를 나타낸다. Clustering 結果에 따른 種의 cluster와 調査地所를 정리하여 分割表를 Block으로 만들어서 이에 대한 集中度分析은 AOC program으로 내장산 資料 (Table 2) 만을 分析 하였다(Park, 1984).

**資料標集** 캐나다 Sifton Bog 中心地域에 인접된 곳에 300 區域을 정하고 random으로 (0.5×0.5)m 크기의 quadrat 30個를 設置하고 出現한 植物 頻度를 조사하여 37種×30 quadrat의 data를 얻었다(Orloci, 1970) 이를 RANK program으로 20種×30 quadrat로 줄여서 資料로 이용 했다(Table 1).

Table 1. Sifton Bog data (Orlaci, 1970)

Species	Quadrat number																															
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30		
1. <i>Sphagnum fuscum</i>	12	13	10	23	14	0	13	0	13	24	21	17	9	18	19	19	16	23	22	20	8	9	5	17	0	12	0	0	13	21		
2. <i>Sphagnum palustre</i>	1	18	16	4	5	0	8	0	5	5	10	14	3	14	4	4	9	12	18	17	21	5	0	18	0	0	11	2				
3. <i>Vaccinium oxycoccus</i>	2	4	4	7	15	0	8	0	10	17	12	13	5	10	1	9	0	9	4	1	1	7	13	15	0	3	0	0	10			
4. <i>Rhamnus frangula</i>	1	0	8	1	2	11	0	14	1	3	3	0	2	0	16	0	3	0	0	6	3	0	1	12	9	21	1	0	0			
5. <i>Rhynchospora alba</i>	11	0	0	11	11	0	0	0	8	4	14	8	2	2	0	15	0	20	3	0	0	0	0	19	0	13	0	0	19			
6. <i>Chamaedaphne calyculata</i>	25	21	22	25	19	0	23	0	20	22	18	25	21	25	1	20	16	20	25	16	18	19	23	22	0	5	0	0	25	25		
7. <i>Andromeda glaucophylla</i>	0	0	0	3	8	0	7	0	6	12	6	0	0	0	4	0	11	1	2	0	0	0	7	0	0	0	0	0	0	0		
8. <i>Drosera rotundifolia</i>	5	1	2	2	9	0	3	0	5	12	6	1	1	3	0	4	1	6	5	0	1	2	2	1	0	0	0	1	2			
9. <i>Sphagnum recurvum</i>	6	0	0	0	0	2	0	5	0	5	0	2	4	0	2	0	0	0	2	0	0	2	0	0	2	0	6	0	5	0		
10. <i>Hypericum verginianum</i>	0	0	0	3	0	0	0	0	5	9	0	3	0	3	0	0	3	0	0	0	0	7	12	0	0	0	0	0	1	0		
11. <i>Thelypteris palustris</i>	0	0	0	0	3	0	5	0	0	1	0	0	0	1	0	0	0	0	0	0	0	8	16	1	0	0	0	1	0	2	0	
12. <i>Sphagnum riparium</i>	10	2	1	0	0	0	0	4	1	0	1	6	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
13. <i>Picea mariana</i>	0	5	1	1	12	1	4	0	0	5	0	2	1	1	3	3	1	2	4	0	0	1	3	0	4	0	0	0	2			
14. <i>Acer rubrum</i>	4	7	1	3	1	5	2	5	0	3	3	4	5	3	9	3	4	2	1	1	1	0	1	10	4	7	3	5	0	0		
15. <i>Carex sp</i>	1	0	0	0	3	0	0	0	16	3	0	2	5	0	0	5	0	2	1	0	0	0	0	0	0	0	0	0	0	3		
16. <i>Pogonia ophioglossoides</i>	0	0	0	0	6	0	0	0	0	12	0	3	0	0	0	0	0	0	0	0	0	1	1	2	0	0	0	0	0	0		
17. <i>Impatiens biflora</i>	0	0	0	0	13	0	16	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	4	0	5	22	0
18. <i>Solanum dulcamara</i>	0	0	0	0	4	0	9	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	18	0	3	6	0	0	
19. <i>Sphagnum capillaceum</i>	0	0	0	1	9	0	5	0	0	0	4	0	0	0	0	0	0	0	0	0	0	0	0	0	3	5	0	0	0	0	0	
20. <i>Sarracenia purpurea</i>	0	0	0	0	4	0	1	0	1	2	0	0	0	0	0	0	0	0	0	0	1	3	1	0	0	3	0	0	0	0	0	



또한 우리나라에서 內臟山을 國立公園으로 지정할 때 조사된 자료를 토대로(Park, 1974) 1985년 4월부터 10월까지 내장산 일대에서 4個地域으로 區分하여 전당대 부근의 6個地所, 선인봉 부근의 9個地所, 까치봉 부근에서 8個地所와 원적암 부근의 9個地所를 택하여 (15×15)m 크기의 quadrat를 각 地所에 5~7개씩 설치하여 出現되는 植物 頻度を 조사하여 78種×52 quadrat의 資料를 얻어 RANK program으로 이를 24種×32 quadrat로 줄여서 分析에 이용하였다(Table 2).

## 結 果

**Sampling information** 資料標集에서 information loss가 最小일때의  $k$  (quadrat 數)와 intensity ( $n$ )을 정했다. 본 研究에서  $n=25$ 일때로 캐나다의 Sifton Bog에서는 0.25 m<sup>2</sup> (0.5 m×0.5m)의 크기의 quadrat를 이용했고 내장산에서는 225 m<sup>2</sup>(15m×15m) 크기의 quadrat에서 植物의 頻度を 이용했다.

Fig. 4에서와 같이 Sifton Bog 資料(Table 1)에서는  $k=5$  일때 information loss가 32%,  $k=20$ 에서 1%였으며 내장산 자료(Table 2)에서는  $k=5$  일때 information loss가 32%,  $k=20$  일때는 0.6%였다.

Information loss는 觀側된 要素( $f_{ij}$ ,  $n-f_{ij}$ )와 期待 要素( $k$ ,  $n-k$ )의 차이를 계산한 것으로  $f_{ij} \leq k$  일때 식(11)을 이용하여 계산된 것이다. Sifton Bog의 植物 頻度は  $n=25$ 가 적절하나 내장산에서는  $n=23$ 에서 information loss가 0이 되므로 標集할 때 時間과 經濟의 一面을 考慮하여 quadrat 크기를 225 m<sup>2</sup>(15m×15m)로 調査하는 것이 効果적인 것 같다.

**Rank에 의한 資料 整理** Table 1은 37種×30 quadrat의 data를 equivocation information으로 種에 대한 rank를 정하여 information rank 값이 17.397 (10%이상)인 20種을 擇하여 20種×30 quadrat로 정리하였고 Table 2는 78種×52 quadrat 크기의 Data를 같은 方法으로 rank에 의해 information rank 값이 16.215 (10%) 이상인 24種을 擇하여 24種×32 quadrat로 정리하였다(Table 3).

R과 Q array와 weighting information 및 CLASS program을 이용하여 Table 1에서 rank가 1과 2인 *Sphagnum fuscum*과 *Sphagnum palustre*에 대하여 class 限界를 0, 1~5, 6~10, 11~15, 15~20, 21~25의 6等級으로 區分하여 R-array하고 marginal dispersion을 구하고 weighting information을 계산하였다(Table 4).

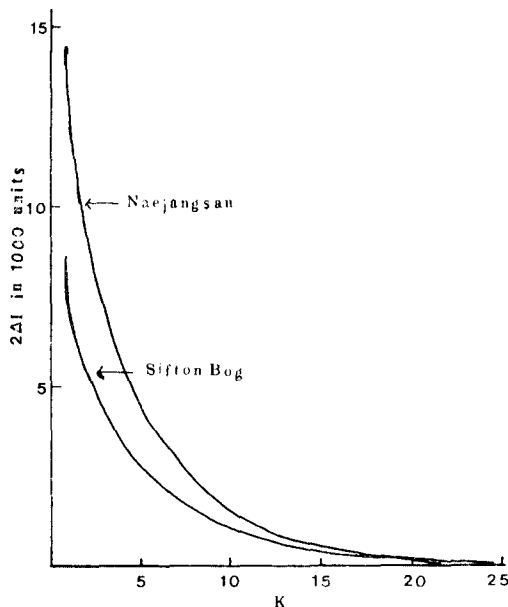


Fig. 4. Information loss curve. The sampling intensity( $n$ ) in Sifton Bog data (Table 1) represented 25 and it is 23 in Naejangsan data (Table 2).



Table 3. Determination of rank for species based on an analysis (RANK) of equivocation information

Species of Sifton Bog	Rank	Information	Species of Naejang San	Rank	Information
<i>Sphagnum fuscum</i>	1	58.4592555	<i>Lespedeza bicolor</i>	1	45.6960449
<i>Sphagnum palustre</i>	2	47.2816482	<i>Lespedeza maximowicz</i>	2	48.1873107
<i>Vaccinium oxycoccum</i>	3	43.9568073	<i>Sasa purpurascens</i>	3	46.3309865
<i>Rhamnus frangula</i>	4	43.7972082	<i>Fraxinus sieboldiana</i>	4	45.921967
<i>Rhynchospora alba</i>	5	40.437303	<i>Pueraria thunbergiana</i>	5	44.8283105
<i>Chamaedaphne calyculata</i>	6	35.8886137	<i>Cornus controversa</i>	6	44.039641
<i>Andromeda glaucophylla</i>	7	31.1129195	<i>Daphniphyllum macropodum</i>	7	43.1043711
<i>Drosera rotundifolia</i>	8	33.0778586	<i>Acer japonica</i>	8	40.8669782
<i>Sphagnum recurvum</i>	9	24.9302653	<i>Carpinus laxiflora</i>	9	38.592689
<i>Hypericum virginianum</i>	10	24.5995935	<i>Akebia quinata</i>	10	37.2623754
<i>Thelypteris palustris</i>	11	23.2824387	<i>Pinus densiflora</i>	11	36.8732983
<i>Sphagnum riparium</i>	12	23.0732848	<i>Rhododendron schlippenbachii</i>	12	36.8165513
<i>Picea mariana</i>	13	23.0657299	<i>Actinidia arguta</i>	13	34.295363
<i>Acer rubrum</i>	14	22.9426326	<i>Staphylea bumalda</i>	14	29.6390097
<i>Carex sp</i>	15	22.3462548	<i>Carpinus tschonoskii</i>	15	28.6522889
<i>Pogonia ophioglossoides</i>	16	21.8723649	<i>Celastrus orbiculatus</i>	16	25.983487
<i>Impatiens biflora</i>	17	20.1777315	<i>Hydrangea serrata</i>	17	25.1403285
<i>Solanum dulcamara</i>	18	18.7914371	<i>Zelkova serrata</i>	18	23.5001815
<i>Sphagnum capillaceum</i>	19	17.7154399	<i>Styrax japonica</i>	19	22.4934058
<i>Sarracenia purpurea</i>	20	17.3974551	<i>Smilax china</i>	20	22.1807098
			<i>Acer mono</i>	21	22.118169
			<i>Rhus javanica</i>	22	16.8102822
			<i>Benzoin obtusifolium</i>	23	16.370776
			<i>Meliosma myriantha</i>	24	16.2155282

**Table 4.** R dispersion array for a two species set

A) Raw Data by Table 1.		quadrat number																													
		1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30
1.	<i>S. fuscum</i>	12	13	10	23	14	0	13	0	13	24	21	17	9	18	19	19	16	23	22	20	8	9	5	17	0	12	0	0	13	21
2.	<i>S. palustre</i>	1	18	16	4	5	0	8	0	5	5	5	10	14	3	14	4	4	4	9	12	18	17	21	5	0	18	0	0	11	2
B) Marginal dispersion arrays																															
Class limits		0	1~5					6~10					11~15					16~20					21~25								
1.	<i>S. fuscum</i>	5	1					4					7					7					6								
2.	<i>S. palustre</i>	5	12					3					4					5					1								
C) R-dispersion array																															
		species 2						Total																							
		0	1~5			6~10			11~15			16~20			21~25																
species 1	0	5	0			0			0			0			0			5													
	1~5	0	0			0			0			0			1			1													
	6~10	0	0			1			3			0			4																
	11~15	0	3			1			1			2			7																
	16~20	0	4			1			2			0			7																
	21~25	0	5			1			0			0			6																
Total		5	12			3			4			5			1			30													

**Table 5.** Q dispersion array for a two quadrats set

A) Raw Data by Table 1.		species number																							
		1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20				
1.	quadrat 11	21	5	12	3	14	18	6	6	5	9	0	0	0	3	0	12	0	0	0	0				
2.	quadrat 24	17	5	15	1	19	22	7	1	2	0	0	0	3	10	0	2	0	0	5	0				
B) Marginal dispersion array																									
Class limits		0	1~5					6~10					11~15					16~20					21~25		
1.	quadrat 11	8	4					3					3					1					1		
2.	quadrat 24	7	7					2					1					2					1		
C) Q-dispersion array quadrat																									
		quadrat 24						Total																	
		0	1~5			6~10			11~15			16~20			21~25										
quadrat 11	0	6	2			0			0			0			0			8							
	1~5	0	3			1			0			0			4										
	6~13	1	1			1			0			0			3										
	11~15	0	1			0			1			1			3										
	16~20	0	0			0			0			0			1										
	21~25	0	0			0			0			1			1										
Total		7	7			2			1			2			1			20							





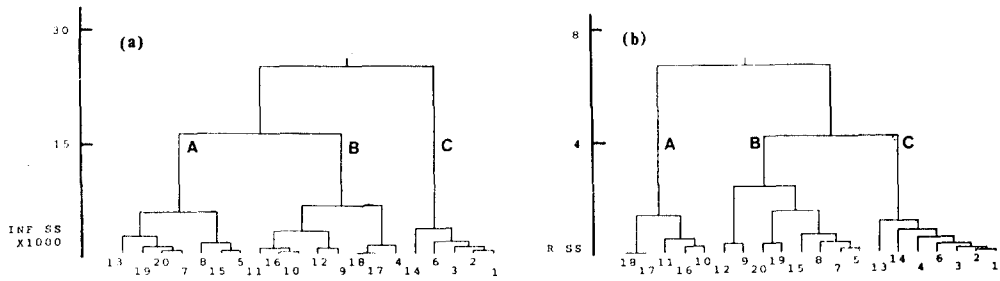


Fig. 5. Dendrograms of 20 species in Table 1. The vertical scale indicates sum of square(SS) of information. The (a) classified by basis on equivocation information and (b) by Rajski's metrics.

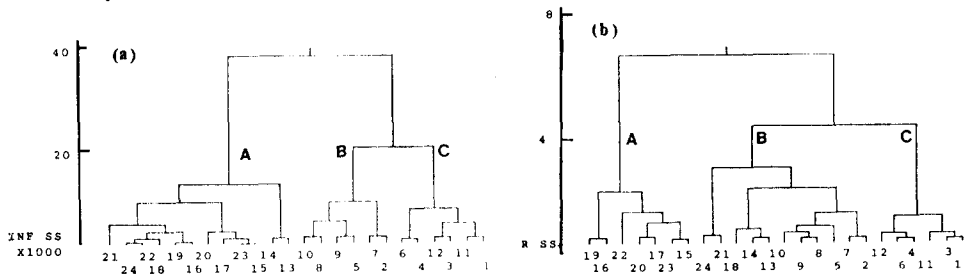


Fig. 6. Dendrograms of 24 species in Table 2. The vertical scale indicates sum of square of information. The (a) classified by basis on equivocation information and (b) on Rajski's metrics.

sum of square(SS) algorithm으로 fusion해서 Fig. 5(a,b)와 Fig. 6(a,b) 같이 dendrogram으로 나타냈다.

Fig. 5의 clustering은 Sifton Bog 資料(Table 1)의 結果로 Fig. 5(a)는 equivocation information 값으로 cluster해서 SS-information이  $15 \times 1,000$  일때 區分했고 Fig. 5(b)는 Rajski metrics 값을 이용한 것으로 SS-information이 4일때 區分한 것이다. 이때 cluster의 效率은 equivocation information으로 clustering 했을때는全體 cluster SS가 30248, cluster內的 SS가 18695이므로 cluster間的 SS는 11553이 된다. 그래서 效率은  $38.2\% (11553/30248 \times 100)$ 가 된다. 반면에 Rajski metrics에 의한 clustering에서는 전체 cluster SS가 5,789, cluster內的 SS는 4,219로써 cluster間 SS가 1,571로 效率은  $27.1\%$ 였다.

또한 equivocation information에 의한 clustering에서 cluster A에는 13, 19, 20, 7, 8, 15, 5의 7種이 포함되었고, cluster B에는 11, 16, 10, 12, 9, 18, 17, 4의 8種이 포함되었고, cluster C에는 14, 6, 3, 1, 2의 5種이 포함되었다(Fig. 5(a)).

Rajski metric에 의한 clustering 結果는 cluster A는 18, 17, 11, 16, 10의 5種이 포함되며, cluster B에는 12, 9, 20, 19, 15, 8, 7, 5,의 8種이 포함되었고, cluster C에는 15, 14, 4, 6, 3, 2, 1의 7種이 포함 되었다(Fig. 5(b)).

두 方法에 의한 clustering의 結果에서 cluster C는 거의 비슷하나 cluster A와 B는 큰 차이가 있었다.

Fig. 6는 內藏山의 資料(Table 2)의 clustering의 結果로 Fig. 6(a)는 equivocation information에 의한 것이며, Fig. 6(b)는 Rajski metrics에 의한 것이다. equivocation information에 의한 clustering 效率은全體 cluster의 SS가 43882이고, cluster內的 SS가 31989이므로

cluster間的 SS는 11894로서 效率은 27.1%(11894/43882×100)였으며 Rajski metrics에 의한 clustering 效率은全體 cluster의 SS가 8.144이고, cluster內的 SS가 6.721이므로 cluster間的 SS가 1.421이므로 區分 效率이 17.4%(1.421/8.144×100)으로 나타 났다.

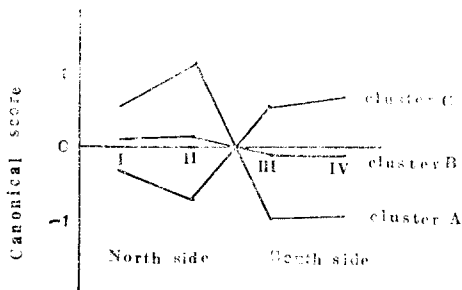
역시 내장산의 資料에 대한 cluster 區分도 3개의 cluster로서 equivocation information에 의한 clustering에서 cluster A는 21, 24, 22, 18, 19, 16, 20, 17, 23, 15, 14, 13으로 12種이 포함되며, cluster B는 10, 8, 9, 5, 7, 2으로 6種이 포함되며, cluster C에는 6, 4, 12, 3, 11, 1의 6種이 포함되었다. 반면에 Rajski metrics에 의한 clustering에서는 cluster A에는 19, 16, 22, 20, 17, 23, 15의 7種이 포함되고 cluster B에는 24, 21, 18, 14, 13, 10, 9, 8, 5, 7, 2의 11種이 포함하고 cluster C에는 12, 6, 4, 11, 3, 1의 6種이 포함 되었다.

두 方法에 의한 clustering의 結果에서도 cluster C에서는 비슷한 양상을 띠었으나 cluster A와 B는 큰 차이가 있었다. 본 研究에서는 clustering의 效率이 높은 equivocation information에 의한 clustering의 結果를 이용하여 集中度分析을 하였다.

**集中度分析(AOC)** 集中度分析을 하기 위하여 分割表는 equivocation information에 의한 cluster A(12種), B(6種) C(6種)과 北斜面 I 地所(전망대 ; 6 quadrat), II 地所(선인봉 ; 9 quadrat)와 南斜面에서 III 地所(까치봉 ; 8 quadrat), IV 地所(원적암 ; 9 quadrat)의 3×4 Block 을 정리하여 식(18)을 이용하여 coherence coefficient를 種(24種)에 대하여 算出된 값( $I_{pp}$ )

**Table 9.** Sum of Rajski's metrics and blocks among three clusters and four stands, and adjusted table corresponded with block table and information

Stands	Sum of Rajski's metrics				Block size				Adjusted value			
	I	II	III	IV	I	II	III	IV	I	II	III	IV
Cluster A	43.16	58.99	52.26	62.51	72	108	96	108	36.52	33.28	33.17	35.26
Cluster B	27.34	37.59	34.20	43.10	36	54	48	54	46.27	42.41	43.41	48.63
Cluster C	28.48	37.36	37.15	44.47	36	54	48	54	48.20	42.15	47.15	50.17



**Fig. 7.** Expected deviation resulting analysis of concentration of three clusters and four stands in Naejangsan data. Cluster A reveals the different gradient along north and south sides against cluster B.

과 Table 2 ( $R_{pp}$ )를 이용하여 block에 따라 合算한 후 adjust했다(Table 9).

集中度分析 次數의 결정은 Gittins(1979)의 理論에 따라 第一變量은 79%, 第二變量은 21%의 分析力을 나타냈는데 第一變量만을 이용하여 해석했다.

Fig. 7에 나타난 바와 같이 cluster A에 속하는 고추나무, 개서나무, 산수국, 노박덩굴, 매죽나무, 고로쇠, 생강나무등은 북사면에 높게 나타났고, 반대로 cluster C에 속하는 싸리나무, 조릿대, 쇠물푸레나무, 소나무, 층층나무 등은 남사면에 높게 나타났으나 cluster B에 속하는 굴거리나무, 참

단풍, 서나무, 칩등은 남북사면에 차이 없이 나타났다.

## 論 議

意思傳達理論(communication theory)에 의한 information理論의 適用은 새로운 理論은 아니나 여러 學問에 多樣하게 適用되어 왔다. 生物學에 이 理論을 적용한 것은 細胞生物學 部分에 Yockey(1958) 등이 이용한 것이 처음이며, 種多樣性的 定量에 Shannon指數를 Magalef(1958)가 적용한 후 Pielou(1966), Orloci(1977)등에 의해서 information 理論을 生態學에 實用 되었으며 Feoli *et al.* (1984)은 FORTRAN program을 작성하였고 Orloci(1978)는 BASIC program으로 data 分析을 할 수 있게 하였다.

본 研究에서는 Fig. 2와 같이 標集할 때 information loss가 最小일 때의 sampling intensity를 정하여 quadrat數와 크기를 限定하여 시간과 노력을 最小化할 수 있고 標集의 效果를 산출할 수 있게 하였다.

Table 1의 資料에서  $n=25$ 일 때 information loss가  $k=5$ 에서 32%였고  $k=25$ 에서 0이 됐다. 이때 quadrat의 크기는  $(0.5 \times 0.5)m$ 로 最適의 크기였다. 또한 Table 2에서는  $n$ (intensity)이 23일 때  $k=5$ 에서 information loss가 32%였으나  $k=23$ 에서 0이됐고 quadrat크기가  $(15 \times 15)m$ 로 종래의 森林生態學에서  $(10 \times 10)m$ 의 크기의 quadrat보다 내장산에서의 이 연구에서는 약간 클 때 標集效率이 높은 것으로 나타났다.

Rank를 정할 때 Orloci(1978)는 equivocation information과 mutual information을 대비했는데 效率이 좋은 equivocation information을 이용하였고 class의 限界는 임의로 택하였다.

計算된 information/entropy 값으로 clustering의 區分은 agglomerative algorithm 法으로 (Sneath and Sokal, 1973) sum of square algorithm을 이용 했으며, 이때 주로 equivocation information과 Rajski metrics를 이용했다. Table 1의 clustering에서 equivocation information으로 구분하면 效率이 38.2% 였으나 Rajski metrics에서는 27.1% 였으며, Table 2의 clustering 效率은 前者에 의해서 37.2%, 後者에 의해서 17.4%로 나타났다. 따라서 集中度分析은 equivocation information clustering 結果를 이용하여 분석했다.

같은 植物群落內에서 植物과 環境과의 關係는 類型(pattern)에 따라 環境變數에 의한 勾配分析(gradient analysis; Whittaker, 1967)과 어떤 地域內의 aggregation된 植物과 環境과의 關係를 ordination 또는 canonical analysis (Orloci, 1978; Gauch, 1982)를 이용하여 classification해서 Block間的 關係를 AOC法으로 解析하는 것이 보통이다(Feoli and Orloci, 1979; Park, 1984).

內藏山 資料(Table 2)에 대한 集中度分析(AOC)을 위해서 coherence coefficient를 산출하여  $(I_{pp})$  adjust는  $3 \times 4$  block을 clustering으로 정하고  $I_{pp}R_{pp}/N$ (식(19))로 변형하여 block內의 습을 이용하였다(Table 9). 3個의 cluster(A, B 및 C)와 北斜面(I, II地所)과 南斜面(III, IV地所)에 대한 集中度分析의 結果는 第一變量이 79%, 第二變量이 21%로 나타나 第一變量으로만 해석했다.

Fig. 7과 같이 cluster A에 포함되는 植物은 北斜面에 높은 頻度を 나타내고 cluster C에 포함되는 植物은 南斜面에 높게 나타났고 cluster B에 포함되는 植物은 南北斜面에 차이가 없이 나타났다.

Information/entropy에 의한 結果 해석은 종래의 多樣性이나 clustering으로 表現해서 資料를 해석하는 이외에도 集中度分析으로 資料內의 關係를 구체적으로 해석하는데 이용할

수 있다.

### 摘 要

Information과 entropy理論을 適用하여 캐나다 Sifton Bog 資料(Table 1)와 內藏山에서 조사된 資料(Table 2)를 分析하였다.

標集할때 information loss가 最小일 때의 sampling intensity ( $n$ )를 限定할 수 있게 했고, weighting information을 算出할때 class를 임의로 區分하여 information 값을 計算했으며, 資料가 방대할 때 equivocation information으로 rank를 정하여 資料의 크기를 줄일 수 있게 했다.

또한 information/entropy 計算을 7가지로 구분하여 선택적으로 information을 算出할 수 있게 하였다. information에 의한 結果를 解析할 때 多樣性 계산과는 다르게 entropy 값으로 clustering하고 效率을 計算하였다.

植被 資料에서 種의 cluster와 調查地所와의 關係 分析은 集中度法을 이용하여 內藏山 資料에서 3個의 種 cluster와 4個 調查地所에 대하여 Rajsiki metrics를 adjust해서 分析하였다. 이때 第一變量(79%) 만을 이용하여 解析하였다.

Cluster A에 속하는 개서나무, 노박덩굴, 고로쇠, 생강나무, 매죽나무등은 北斜面에 높게 나타났으며 cluster C에 속하는 싸리나무, 소나무, 쇠물푸레나무, 층층나무등은 南斜面에 높게 나타났으나 cluster B에 속한 참단풍, 서나무등은 남북사면에 관계없이 고르게 나타났다.

### 引 用 文 獻

- Feoli, E. and L. Orloci. (1979). Analysis of concentration and detection of underlying factors in structure tables. *Vegetatio*, **40** : 49~54.
- Feoli, E., M. Lagonegro and L. Orloci. (1984). Information analysis of vegetation data. Junk, Hague.
- Fisher, R.A. (1948). Statistical method for research workers. 10th ed. Oliver and Boyd., Edinburgh.
- Gauch, H.G. Jr. (1982). Multivariate analysis in community ecology. Cambridge Univ. Press, Cambridge.
- Greig-Smith, P. (1957). Quantitative plant ecology. Butterworth, London.
- Kullback, S. (1959). Information theory and statistics. Wiley, New York.
- MacArthur, R.H. (1965). Pattern of species diversity. *Biol. Rev.*, **40** : 510~533.
- Margalef, D.R. (1958). Information theory in ecology. *Yearbook Soc. Syst. Research*, **3** : 36~71.
- McIntosh, R.P. (1967). An index of diversity and the relation of certain concepts to diversity. *Ecology*, **48** : 392~403.
- Legendre, L. and P. Legendre. (1983). Numerical ecology. Elsevier Sci. Co., New York.
- Orloci, L. (1970). Analysis of vegetation samples based on the use of information. *J. Theor. Biol.*, **29** : 173~189.
- Orloci, L. (1972). On information analysis in phytosociology. Junk, Hague.
- Orloci, L. (1976). Ranking species by an information criterion. *J. Ecol.*, **64** : 417~419.
- Orloci, L. (1977). Ranking species based on the components of equivocation information. *Vegetatio*, **37** : 123~125.
- Orloci, L. (1978). Multivariate analysis in vegetation research. Junk, Hague.



- Park, B.K. (1974). A phytosociological study on the vegetation of National Park Mt. Naejangsan. A report of Korea. Associ. for Conservation of Nature.
- Park, S.T. (1984). An application of analysis of concentration for ecological study. Korean J. Bot., **27** : 223~231.
- Pielou, E.C. (1966). Shannon's formula as measure of species diversity. Amer. Natur., **100** : 463~465.
- Pielou, E.C. (1977). Mathematical ecology. Wiley, New York.
- Rajski, C. (1961). Entropy and metric space: Information theory. Butterworth, London.
- Renyi, A. (1961). On measures of entropy and information. 4th Berkeley Symposium on Math. Statis. and Prob. Univ. of Calif. Press, Berkeley.
- Shannon, C.E. (1948). A mathematical theory of communication. Bell system Tech., **27** : 379~423.
- Sneath, P.H.A. and V.D. Sokal. (1973). Numerical taxonomy. Freeman, San Francisco.
- Whittaker, R.H. (1967). Gradient analysis of vegetation. Biol. Rev., **42** : 207~264.
- Yockey, H.P., R.L. Platzman and H. Quastler. (1958). Information theory in biology. Pergamon, New York.

((1987年 3月 23日 接受))