Random partitions and other combinatorial objects

Peter McCullagh

Department of Statistics University of Chicago

Hotelling Lecture II UNC Chapel Hill December 3 2008

◆□▶ ◆□▶ ▲□▶ ▲□▶ □ のQ@

Outline

Random partitions

Ewens partition process

Applications Gauss-Ewens process Illustration Classification Illustration

Trees

Gibbs fragmentation trees

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

Partitions

 $[n] = \{1, \ldots, n\}$ a finite set A partition B of the set [n] = [6] is (i) a set of disjoint non-empty subsets $b \subset [n]$ called blocks... e.g. $B = \{\{2, 4, 6\}, \{1, 3\}, \{5\}\} \equiv 246|13|5 \equiv 13|246|5$ (ii) an equivalence relaton $B: [n] \times [n] \rightarrow \{0, 1\}$ s.t. B(i, j) = 1 if $i \sim j$ belong to the same block (iii) a symmetric Boolean matrix

#B: number of elements (no. of blocks)

for $b \in B$, #b > 0 is the number of elements (#b > 0) Integer partition $\nu(B) = 1 + 2 + 3$ associated with B = 246|13|5

The set \mathcal{E}_n of partitions of [n]

 $\begin{array}{l} \mathcal{E}_1: 1 \\ \mathcal{E}_2: 12, \quad 1|2 \\ \mathcal{E}_3: 123, \quad 12|3, \quad 13|2, \quad 23|1, \quad 1|2|3| \\ \mathcal{E}_4: 1234, \quad 123|4[4], \quad 12|34[3], \quad 12|3|4[6], \quad 1|2|3|4 \\ \mathcal{E}_5: 12345, \quad 1234|5[5], \quad 123|45[10], \quad 123|4|5[10], \quad 12|34|5[15], \\ 12|3|4|5[10], \quad 1|2|3|4|5 \\ \#\mathcal{E}_n: 1, 2, 5, \quad 15, \quad 52, \quad 203, \quad 877, \quad 4140, \quad 21147, \quad 115975, \quad 678570 \end{array}$

Permutation map $\pi: [n] \rightarrow [n]$ also acts $\pi^*: \mathcal{E}_n \rightarrow \mathcal{E}_n$ partition type $\nu(B)$ is maximal invariant Deletion map: $D_n: \mathcal{E}_n \rightarrow \mathcal{E}_{n-1}$ (onto) $D_4: 1234, 123|4 \rightarrow 123$ $12|3|4, 12|34, 124|3 \rightarrow 12|3$ (three types) $1|2|3|4, 1|2|34, 1|24|3, 14|2|3 \rightarrow 1|2|3$

Same as removal of last row and column from matrix

 \mathcal{E} represents the sets $\{\mathcal{E}_n\}$ with permutation and deletion maps

Probability distributions on partitions

 P_n a probability distribution on \mathcal{E}_n Finitely exchangeable if $\nu(B) = \nu(B')$ implies $P_n(B) = P_n(B')$ Examples:

\mathcal{E}_3	123	12 3	13 2	23 1	1 2 3	
P_3	1/3	1/6	1/6	1/6	1/6	
P'_3	1/6	1/6	1/6	1/6	1/3	
\mathcal{E}_4	1234	123 4	13 24	23	1 4	1 2 3 4
P_4	1/4	1/12	1/24	1	/24	1/24
P'_4	1/10	1/15	1/30	1	/15	2/15

Compatibility:

 $\begin{array}{l} P_4(1234\cup 123|4)=1/4+1/12=1/3=P_3(123)\\ P_4(12|3|4\cup 12|34\cup 124|3)=2/24+1/12=1/6=P_3(12|3)\\ P_4(1|2|3|4\cup 1|2|34[3])=1/24+3/24=1/6=P_3(1|2|3)\\ P_3 \text{ is the marginal distribution of } P_4 \text{ under deletion}\\ P_3' \text{ is the marginal distribution of } P_4' \text{ under deletion} \end{array}$

Exchangeable partition process

An *exchangeable partition process* is a sequence P_n such that each P_n is finitely exchangeable $P_n(B)$ depends only on block sizes $\nu(B)$ P_n is the marginal distribution of P_{n+1} .

Kolmogorov compatibility condition:

$$P_n(B) = \sum_{B': D_{n+1}B'=B} P_{n+1}(B')$$

Conditional distribution

$$P_{n+1}(B' | B \in \mathcal{E}_n) = egin{cases} P_{n+1}(B')/P_n(B) & D_{n+1}B' = B \ 0 & ext{otherwise.} \end{cases}$$

Kingman's paintbox characterization

The Ewens partition process

Ewens distribution with parameter $\lambda > 0$

$$P_n(B) = rac{\Gamma(\lambda)\lambda^{\#B}}{\Gamma(n+\lambda)} \prod_{b\in B} \Gamma(\#b)$$

Conditional distributions

$$P_{n+1}(u_{n+1} \mapsto b \mid B) = \frac{P_{n+1}(B')}{P_n(B)} = \begin{cases} \#b/(n+\lambda) & b \in B\\ \lambda/(n+\lambda) & b = \emptyset \end{cases}$$

(Pitman's CRP description)

Induced distribution on integer partitions $\nu(B) = 1^{\nu_1} 2^{\nu_2} \cdots n^{\nu_n}$

$$Q_n(\nu) = \frac{\Gamma(\lambda)\lambda^{\nu}}{\Gamma(n+\lambda)} \prod_{j=1}^n ((j-1)!)^{\nu_j} \times \frac{n!}{\prod (j!)^{\nu_j} \nu_j!}$$

▲□▶▲□▶▲□▶▲□▶ □ のへで

No deletion operation for integer partitions Hence no process on integer partitions

Permutation process

Exponential family of distributions on permutations $\sigma: [n] \rightarrow [n]$

$$P_{n}(\sigma; \lambda) = \lambda^{\#\sigma} / M_{n}(\lambda), \quad \#\sigma = \#\text{cycles}$$
$$M_{n}(\lambda) = \sum_{\lambda} \lambda^{r} S_{n,r}$$
$$= \lambda(\lambda + 1) \cdots (\lambda + n - 1) = \Gamma(n + \lambda) / \Gamma(\lambda)$$
$$P_{n}(\sigma; \lambda) = \frac{\Gamma(\lambda) \lambda^{\#\sigma}}{\Gamma(n + \lambda)}$$

(ロ) (同) (三) (三) (三) (○) (○)

Exponential family with canonical statistic $\#\sigma$ Cumulant function $K(\lambda) = \log M(\lambda)$ determines the mean, variance,... of $\#\sigma$ (Goes back to Euler)

In what sense is this an exchangeable process?

Permutations $\{\Pi_n\}$ as a projective system

Projective system with respect to sub-sampling

A sub-sample of size *m* taken from [*n*] (not random) is an ordered subset $\varphi_1, \ldots, \varphi_m$ distinct in [*n*] a 1–1 map $\varphi : [m] \rightarrow [n]$

sub sample $\varphi : [m] \to [n] \longrightarrow$ deletion $\varphi^* : \Pi_n \to \Pi_m$ φ^* in reverse direction on permutations by conjugation $\varphi^* \sigma = \varphi^{-1} \sigma \varphi$ if m = nby deletion from cycle representation if $m \le n$ $(i, j, \ldots)(\ldots) \mapsto (\varphi^{-1}(i), \varphi^{-1}(j), \ldots)(\ldots)$ delete if $\varphi^{-1}(\{j\}) = \emptyset$

 $(\varphi\psi)^* = \psi^*\varphi^*$ (composition in reverse order)

Exponential family distributions are compatible with these maps

Ewens permutation process

Ewens distribution on permutations Π_n

$$P_n(\sigma;\lambda) = \frac{\Gamma(\lambda)\lambda^{\#\sigma}}{\Gamma(n+\lambda)}$$

Induced distribution on partitions (cycles ignoring cyclic order)

$$P_n(B;\lambda) = \frac{\Gamma(\lambda)\lambda^{\#B}}{\Gamma(n+\lambda)} \prod_{b \in B} \Gamma(\#b)$$

Conditional distribution given Π_n

$$\begin{aligned} P_{n+1}(n+1 \mapsto (i,n+1,\sigma(i),\ldots) \,|\, \sigma) &= 1/(n+\lambda) \quad 1 \le i \le n \\ P_{n+1}(n+1 \mapsto (n+1) \,|\, \sigma) &= \lambda/(n+\lambda) \quad \text{new cycle} \end{aligned}$$

Defines an infinite exchangeable random permutation $\#\sigma$ is approximately $Po(\lambda \log n)$ Note difference between permutation and a ranking

Other interpertations of the Ewens process

Conditional Poisson interpretation X_1, X_2, \ldots independent Poisson variables $X_j \sim Po(\lambda/j)$ as multiplicities $1^{X_1} 2^{X_2} 3^{X_3} \cdots$ in integer partition

Conditional distribution of X_1, \ldots, X_n given $\sum jX_j = n$

$$p(x) \propto e^{-\lambda \sum 1/j} \frac{\lambda^{x_{\bullet}}}{\prod j^{x_j} x_j!}$$

is exactly the Ewens partition

Negative binomial model for the number of species Fisher 1943; Good 1953; Mosteller & Wallace 197?; Efron & Thisted 1976

Kendall's (1975) family-size process (Kelly's book)

Prime factorization of large integers (Billingsley 1972; Donnelly & Grimmett)

Partition induced by Dirichlet process

Characterization of the Ewens distribution

Why is the Ewens distribution ubiquitous?

- (i) Exchangeability: $B \sim P_n$ implies $B^{\pi} = \pi^{-1} B \pi \sim P_n$
- (ii) Restriction to subsets $[m] \subset [n]$ if $B \sim P_n$, the restriction is $B[m] \sim P_m$ (process property)

(iii) Self-similarity (lack of memory) Given that $B \le b|b'$, conditional distn is $B \sim P_{\#b} \times P_{\#b'}$ (Aldous, 199?)

A D F A 同 F A E F A E F A Q A

Leading to a theory of Markov fragmentation trees... by recursive partitioning...

The Gauss-Ewens cluster process

Cluster process has following parts:

(i) An index set \mathbb{N}

(ii) A random sequence Y_1, Y_2, \ldots with $Y_i \in S$, $(S = \mathbb{R}^d)$

(iii) A random partition B of \mathbb{N} (not a partition of S)

(iv) Finite-dimensional distributions such as

$$P_n(B) = \frac{\lambda^{\#B} \Gamma(\lambda)}{\Gamma(n+\lambda)} \prod_{b \in B} \Gamma(\#b)$$

$$Y[n] \mid B[\mathbb{N}] \sim N(\mathbf{1}\mu, I_n \otimes \Sigma + B[n] \otimes \Sigma')$$

Note $B \equiv B[\mathbb{N}]$ (no interference)

Parameters: $\mu \in \mathcal{R}^d$, $\lambda > 0$;

 Σ, Σ' within- and between-cluster covariance matrices

(日) (日) (日) (日) (日) (日) (日)

Exchangeability of cluster process

Observation for a finite sample $[n] \subset \mathbb{N}$ is (Y[n], B[n])

Observation space is $S^n \times \mathcal{E}_n$

Permutation π : $[n] \rightarrow [n]$ acts on observation space $(Y[n], B[n]) \mapsto (Y^{\pi}, B^{\pi})$ by composition $B^{\pi}(i, j) = B(\pi_i, \pi_j)$

Restriction $\varphi : [m] \to [n]$ acts on observation spaces $(Y[n], B[n]) \mapsto (Y^{\varphi}, B^{\varphi})$ by composition $Y^{\varphi}(i) = Y(\varphi(i)), \qquad B^{\varphi}(i, j) = B(\varphi_i, \varphi_j)$

(i) Distribution *Q_n* on *E_n × Sⁿ* unaffected by permutation π of [*n*]
(ii) *Q_m* on *E_m × S^m* is the marginal distn of *Q_n*

Hence there exists an infinite random clustering (Y, B)

More conventional version

Gauss-Ewens cluster model is more or less equivalent to

 $\eta \sim IIDN(0, \Sigma'), \quad \epsilon \sim IIDN(0, \Sigma) \text{ independent}$ $ext{tbl}_i = 1 + \max(ext{tbl}[i - 1]) \quad \text{w.p. } \lambda/(i - 1 + \lambda);$ $ext{else one of (tbl_1, ..., tbl_{i-1}) with equal prob}$ $Y_i = \mu + \epsilon_i + \eta_{ ext{tbl}(i)}$

with $(Y_1, \text{tbl}_1), \dots, (Y_n, \text{tbl}_n)$ observed ...except that this version is not exchangeable

Can fix this by forgetting/ignoring table numbers i.e. by defining B = outer(tbl, tbl, "==") and saying that (*Y*[*n*], *B*[*n*]) is observed.





Statistical classification (aka supervised learning)

Feature $Y_u \in S$ in feature space and class $t_u \in C$

Training sample u_1, \ldots, u_n : observed features (Y_1, \ldots, Y_n) and types (t_1, \ldots, t_n) and $Y_{n+1} = y^*$, to which class does u_{n+1} belong?

Deterministic interpretation: (forced choice of one $t^* \in C$) Statistical classification: probability distribution on CEnormous literature going back to Fisher (1936)

Logistic classification model (more recent; 1970s?)

$$\log \operatorname{pr}(t_u = r \mid y_u = y) = \frac{e^{\beta_r' y}}{\sum_r e^{\beta_s' y}}$$

for $r \in C$, independently for distinct units

Cluster models for classification w/o classes

Problem: No set C of classes in a cluster process (Y, B)Observation (Y, B)[n] in training sample u_1, \ldots, u_n How can we assign new unit u_{n+1} to classes?

Conditional distribution

$$\operatorname{pr}(u_{n+1} \mapsto b \mid (Y, B)[n], y^*) = \begin{cases} f(\ldots) & b \in B \\ \ldots & \text{otherwise.} \end{cases}$$

Blocks of *B* are the classes!

Also need parameter estimates (at least $\lambda, \theta = \Sigma' \Sigma^{-1}$)

Lack of C is a big advantage! possibility of assigning u_{n+1} to a previously unseen class

Explicit calculation of conditional distribution

Simplification $\Sigma' = \theta \Sigma$ in $S = \mathcal{R}^d$

$$\mathsf{pr}(u' \mapsto b \mid ...) \propto \begin{cases} \# b \, \phi_d(y(u') - \tilde{\mu}_b; \tilde{\Sigma}_b) & b \in B \\ \lambda \, \phi_d(y(u'); \Sigma(1 + \theta)) & b = \emptyset \end{cases}$$

$$\tilde{\mu}_b = (\mu + n_b \theta \bar{y}_b)/(1 + n_b \theta), \quad \tilde{\Sigma}_b = \Sigma (1 + \theta/(1 + n_b \theta))$$

Typical values $\theta \ge 5$ and $n_b \ge 5$ so $\tilde{\mu}_b/\bar{y}_b = n_b\theta/(1+n_b\theta) \ge 0.96$

(similar to Fisher discriminant model, but with shrinkage)

(日) (日) (日) (日) (日) (日) (日)

Tree version with classes and sub-classes

Block having maximum conditional probability



< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Trees (rooted and leaf-labelled)



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - 釣��

Set of rooted leaf-labelled trees

Each tree T is a positive definite matrix

Skeleton tree sk(T) is the set of edges (also called topology) $sk(T) = \{abcde, ace, bd, ce, a, b, c, d, e\}$ No. of edges $\leq 2n - 1$.

◆□▶ ◆□▶ ◆□▶ ◆□▶ ● ● ● ●

Structure of T_n

Symmetric, non-negative, $T_{ij} \ge \min(T_{ik}, T_{jk})$

(i) Closed under component-wise monotone transformation $T \mapsto g(T)$ with $g(0) \ge 0$

(ii) closed under component-wise scalar multiplication

(iii) Define
$$(AB)_{ij} = \max_k \{\min(A_{ik}, B_{kj})\}$$

then $T^2 = T$ if and only if $T \in \mathcal{T}_n$.

(iv) if A is non-negative, then $\lim_{n\to\infty} A^n = T$ exists

(v) T_n is the union of intersecting manifolds of dimension 2n-1. How many? $1 \cdot 3 \cdots (2n-3)$

(vi) Contrast unrooted trees: $U_{ij} = T_{ii} + T_{jj} - 2T_{ij}$ $U_{ij} + U_{kl} \le \max\{U_{ik} + U_{jl}, U_{il} + U_{jk}\}$ (Buneman)

Exchangeable fragmentation trees

Gibbs fragmentation trees:

Exchangeable random tree:

(i) distribution P_n on T_n invariant under permutation (ii) P_n is marginal distribution of P_{n+1}

Markov property:

splits independent of waiting times edge lengths independent exponential Branches behave independently following split Branch of size r distributed as P_r (self-similarity)

(日) (日) (日) (日) (日) (日) (日)

Binary splits

Kolmogorov consistency for the Gibbs skeletal tree

Gibbs skeletal tree T: a random collection of subsets of [*n*] satisfying certain tree conditions

$$p_n(T) = K_n^{-1} \prod_{b \in T} \psi(\#b)$$
(1)

 $\psi(n), n = 1, 2, 3, \dots$ Gibbs potential function

Kolmogorov consistency condition P_n is marginal distn of P_{n+1}

Can take $\psi_1 = \psi_2 = 1$ w.l.o.g. Consistency implies

$$\psi_{n+1} = \frac{\psi_n(1+(n-1)\gamma)}{2+\psi_n+(n-2)\gamma}$$

for some $\gamma > 0$. $\psi_3 = (1 + \gamma)/3$ determines the entire sequence

Leaf deletion for n = 3



Gibbs fragmentation trees

The class of exchangeable homogeneous Markovian trees one-dimensional family (Aldous's beta-splitting rules $\beta > -2$) (Bertoin, Le Gall, Berestyki, McC, Pitman, Winkel,)

One member ($\beta = -1$) Waiting time exponential with rate $\varphi(n) = 1 + 1/2 + \dots + 1/(n-1)$ mean waiting time $1/\varphi(n) \simeq 1/\log(n)$ ($\varphi(1) = 0$) splitting distribution $n \mapsto r + s$

$$p_n(r,s) = \frac{n}{2 r s \varphi(n)}$$
 $1 \le r \le n-1$, $r+s=n$

Exercises connected with Gibbs trees

(i) Beginning with [n] at t = 0, find the partition at time t

(ii) Description of behaviour as $n \to \infty$

e.g. size of largest block at time t

(iii) Time to complete fragmentation: $T = 2 \log(n) + O_p(1)$ ($\beta = -1$)

Density: $f(t) \propto (1 - e^{-t/2})^{n-2} e^{-t}$

(iv) Distn of time until u_1 is isolated (leaf height)

Density: $\sum_{r=1}^{n} {\binom{n-1}{r-1}} (-1)^r \varphi_r e^{-\varphi_r t}$

(v) Expected fragmentation rate given [n] at time 0

$$\lambda_n(t) = \sum_{r=1}^n \binom{n}{r} (-1)^r \varphi_r e^{-\varphi_r t}$$

(vi) Analogous theory for unrooted trees

(vii) Applications of Gibbs trees in statistical models