

A formula for partial and conditional infinite exchangeability

P.G.L. Porta Mana 

Kavli Institute, Trondheim [<pgl@portamana.org>](mailto:pgl@portamana.org)

14 May 2020; updated 2 June 2020

[draft] A formula is given for conditionally, infinitely exchangeable probability distributions.

0 Introduction and notation

De Finetti's theorem for the representation of infinitely exchangeable probability distributions yields some of the formulae, derivable from the probability calculus, with the richest practical and philosophical consequences.

In this work I present three results; each may interest a different audience.

The first result, derived in § 1, is a reformulation of partial exchangeability and of its representation theorem in a slightly unfamiliar form. Though not remarkable, this reformulation gives some insights into the connection between partial and conditional exchangeability, and their connection to regression.

The second result, derived in § 2, is an integral representation for joint (predictive) probability distributions with particular symmetries: some of the conditional distributions obtained from them satisfy partial exchangeability. This integral representation has the usual de Finetti form and, remarkably, its density must factorize in a specific way. This factorization is implied by, and implies, the conditional-exchangeability symmetries.

The third result, § 3, brings together exchangeability and Bayesian networks. A Bayesian network expresses assumptions of conditional dependence and independence. Exchangeability expresses assumptions about the mutual informational relevance of a set of phenomena or experiments, deemed to be 'similar'. When we combine these two kinds

of assumptions we obtain (predictive) probability distributions with a particular integral representation: a mixture of copies of the same Bayesian network. This result has some analogies with de Finetti's usual representation, which is a mixture of independent probability distributions ('i.i.d.').

Each of these three results builds on the preceding one(s). The final section discusses them further and hints at possible applications for them.

The remainder of this section introduces some notation and summarizes de Finetti's theorem for full exchangeability. It can be skimmed through by readers familiar with exchangeability theorems, just to grasp the notation I use.

0.1 Notation and summary of representation for full exchangeability

For the details about exchangeable distributions I refer to Bernardo & Smith¹, Diaconis & Freedman², and Dawid's³ review.

Our domain of discourse consists of a countably infinite set of atomic statements (in the logical sense)

$$\{X_i = x_i \mid i \in \mathbf{N}, \forall i \ x_i \in \mathfrak{X}\} \quad (1)$$

where \mathfrak{X} is a finite set. For each i the statements $\{X_i = x \mid x \in \mathfrak{X}\}$ are assumed mutually exclusive on information I . (The theorem holds for any set of statements with these properties, even if the statements are not of the form ' $X = x$ '.)

A probability distribution over these atomic statements is called fully (infinitely) exchangeable if

for every N , every set $\{i_1, \dots, i_N\} \subset \mathbf{N}$, and every permutation π thereof,

$$P(X_{i_1} = x_{i_1}, X_{i_2} = x_{i_2}, \dots, X_{i_N} = x_{i_N} \mid I) = P(X_{i_1} = x_{\pi(i_1)}, X_{i_2} = x_{\pi(i_2)}, \dots, X_{i_N} = x_{\pi(i_N)} \mid I), \quad (2)$$

and if all such probabilities are consistently related by marginalization. This property is equivalent to declaring the empirical frequencies of the values x to be sufficient statistics. The commas between statements denote logical conjunction (' \wedge '), so the order of the statements is immaterial.

¹ Bernardo & Smith 2000 §§ 4.3, 4.6. ² Diaconis & Freedman 1980a,b. ³ Dawid 2013.

In the following I let $\{1, 2, \dots, N\}$ denote any subset of \mathbf{N} , to avoid a proliferation of subscripts. They should be read as $\{i_1, i_2, \dots, i_N\}$.

Denote by $f_x := (f_x)$ a normalized distribution over the values $x \in \mathfrak{X}$. The set of all such distributions is a simplex of dimension $|\mathfrak{X}| - 1$.

For each $x \in \mathfrak{X}$, denote by F_x the empirical relative frequency of x in the set $\{x_1, \dots, x_N\}$:

$$NF_x := \sum_i \delta(x, x_i), \quad x \in \mathfrak{X}. \quad (3)$$

De Finetti's theorem states that a fully exchangeable distribution can be written as follows:

$$P(X_1 = x_1, \dots, X_N = x_N | I) =$$

$$\int \prod_i f_{x_i} p(f_x | I) df_x \equiv \int \prod_x f_x^{NF_x} p(f_x | I) df_x, \quad (4)$$

where the integral is over the simplex of distributions $\{f_x\}$.

In the first integral form, the product is over the set of instances $1, \dots, N$. In the second, equivalent integral form, the product is over the set of values x . This form shows that the empirical frequency distribution (F_x) is a sufficient statistic; it also hint at the important role played in the theorem by the relative entropy of (F_x) with respect to (f_x).

For enough large N , the probability of observing an empirical frequency distribution F_x within a small volume v centred around the distribution f_x is approximately given by the density $p(f_x | I) df_x$:

$$P(F_x \in v | N \text{ large}, I) \approx p(f_x | I) v. \quad (5)$$

For this reason the parameter f_x can be interpreted as a long-run frequency distribution⁴. I will therefore call it so sometimes, but without the intention to force such interpretation on you.

1 Partial exchangeability: alternative form

In de Finetti's theorem for partially exchangeable distributions, the set $\{X_i = x_i\}$ of § 0.1 is divided into two or more categories represented by subsets $\{Y_j = y_j\}, \{Z_k = z_k\}, \dots$. Partial exchangeability of the distribution

⁴ But this *long run* is a misleading guide to current affairs. *In the long run* we are all dead' (Keynes 2013 § 3.I p. 65).

for such statements means that permutations are allowed within each subset but not necessarily across subsets. The usual representation in this case, after a suitable re-indexing $\{1, 2, \dots\} \mapsto \{1', 2', \dots, 1'', 2'', \dots\}$, has the form

$$P(Y_{1'} = y_{1'}, Y_{2'} = y_{2'}, \dots, Z_{1''} = z_{1''}, Z_{2''} = z_{2''}, \dots | J) = \iint \prod_j g_{y_j} \prod_k h_{z_k} p(g, h | J) dg dh, \quad (6)$$

with distinct normalized distributions g, h for each category. If the density $p(g, h | J) dg dh$ is diagonal, that is, if it contains a term $\delta(g - h)$, the fully exchangeable form (4) is recovered.

A little reflection shows that if we know the quantities X_i to belong to category Y in instances $i = 1', 2', \dots$, and to category Z in instances $i = 1'', 2'', \dots$, then (a) there are some other quantities C_i that allows us to distinguish the two categories, and (b) the values of these quantities are *known* for all instances.

Let us say, for example, that the quantities X_i are the results of animal treatments, with values 'S'uccess and 'F'ailure. Y refers to the results for treatments on Yaks, and Z on Zebras. If we write

$$P(Y_3 = S, Z_5 = F | J) = 0.2,$$

then we must already know that animal number 3 is a yak: $C_3 = Y$, and animal number 5 is a zebra: $C_5 = Z$. This is clear from our very notation, otherwise we would not have known whether to use the symbol Y or Z for those instances. This information is evidently implicit in our background information J .

We now make the category information more explicit. A slightly different definition of partial exchangeability is thus obtained, with a slightly different form of its representation theorem.

Besides the statements $\{X_i = x_i\}$, we introduce an additional, set of atomic statements

$$\{C_i = c_i \mid i \in \mathbf{N}, \forall i \ c_i \in \mathfrak{C}\}. \quad (7)$$

For each i the statements $\{C_i = c \mid c \in \mathfrak{C}\}$ are mutually exclusive on information I .

These statements allow us to identify each instance $1, 2, \dots$ as belonging to one or another category out of the finite set \mathfrak{C} .

A probability distribution over the $X_i = x_i$ atomic statements is called partially exchangeable if

for every N , every set of indices $\{1, \dots, N\} \subset \mathbf{N}$, and every permutation π thereof such that $\pi(i) = j \Rightarrow c_i = c_j$,

$$P(X_1 = x_1, \dots, X_N = x_N \mid C_1 = c_1, \dots, C_N = c_N, I) = \\ P(X_{i_1} = x_{\pi(1)}, \dots, X_{i_N} = x_{\pi(N)} \mid C_1 = c_1, \dots, C_N = c_N, I). \quad (8)$$

that is, the only allowed permutations are those *which exchange indices having the same c value*, that is, belonging to the same category.

Let us rewrite the representation formula (6) accordingly.

For each category $c \in \mathfrak{C}$, introduce a normalized distribution $\{f_{x|c} \mid x \in \mathfrak{X}\}$ over the values x . As the notation suggests, it can be considered as a *conditional* distribution over x given c . Denote (with some abuse of the symbols) by $f_{x|c} := (f_{x|c})$ the set of all such conditional distributions. This set is the Cartesian product of $|\mathfrak{C}|$ simplices, each of dimension $|\mathfrak{X}| - 1$.

Denote by $F_{x,c}$ the empirical, joint relative frequency of the pair of values (x, c) occurring in the set of pairs $\{(x_1, c_1), \dots, (x_N, c_N)\}$:

$$NF_{x,c} := \sum_i \delta(x, x_i) \delta(c, c_i), \quad x \in \mathfrak{X}, c \in \mathfrak{C}. \quad (9)$$

Thus $NF_{x,c}$ is the total number of times value x appears among the pairs with $c_i = c$.

De Finetti's theorem states that the partially exchangeable distribution (8) can be written as follows:

$$P(X_1 = x_1, \dots, X_N = x_N \mid C_1 = c_1, \dots, C_N = c_N, I) = \\ \int \prod_{c,x} f_{x|c}^{NF_{x,c}} P(f_{x|c} \mid I) df_{x|c}. \quad (10)$$

Scrutiny of this formula shows that this form is equivalent to the more familiar representation (6). The integral contains one product of $f_{x|c}$ terms for every category c . In each such product, $f_{x_i|c}$ terms are multiplied together for those i such that $c_i = c$. There are exactly $NF_{x,c}$ such terms.

This alternative formulation of partial exchangeability shows that this symmetry could also be called *conditional* exchangeability instead. The role of conditional distributions is clear in the representation (10). In the following I will use the term 'conditional' instead of 'partial' to emphasize this.

2 Representation for joint distributions with conditional exchangeability symmetries

Suppose that we would assign a conditionally (i.e., ‘partially’: see previous section) exchangeable distribution of probability to the statements $\{X_i = x_i\}$, if we knew the true $\{C_i = c_i\}$. But we do not know the latter. What kind of properties does the joint probability distribution of these statements have? And the marginal distribution for $\{X_i = x_i\}$?

The joint probability distribution can be rewritten

$$\begin{aligned} P(X_1 = x_1, C_1 = c_1, \dots, X_N = x_N, C_N = c_N | I) = \\ P(X_1 = x_1, \dots, X_N = x_N | C_1 = c_1, \dots, C_N = c_N, I) \times \\ P(C_1 = c_1, \dots, C_N = c_N | I), \quad (11) \end{aligned}$$

where the first factor, conditionally exchangeable, can be represented by the integral of eq. (10).

Let us suppose that our uncertainty about the statements $\{C_i = c_i\}$ is expressed by a fully exchangeable marginal probability distribution. An integral representation analogous to (4) then holds:

$$P(C_1 = c_1, \dots, C_N = c_N | I) = \int \prod_c f_{,c}^{NF_{,c}} p(f_{,c} | I) df_{,c}, \quad (12)$$

where $f_{,c} := \{f_{,c}\}$, and $F_{,c} := \sum_x F_{x,c}$ is the marginal empirical distribution for the c values.

We can now replace the integral representations (10) and (12) into the product (11). The products within their integrals can be combined considering that

$$f_{,c}^{NF_{,c}} = f_{,c}^{N \sum_x F_{x,c}} = \prod_x f_{,c}^{NF_{x,c}}. \quad (13)$$

We obtain

$$\begin{aligned} P(X_1 = x_1, C_1 = c_1, \dots, X_N = x_N, C_N = c_N | I) = \\ \int \prod_{c,x} f_{x,c}^{NF_{x,c}} p(f_{xc} | I) df_{xc} \quad (14a) \end{aligned}$$

$$\text{with } \boxed{p(f_{xc} | I) df_{xc} = p(f_{x|c} | I) p(f_{,c} | I) df_{x|c} df_{,c}} \quad (14b)$$

where we have defined $f_{x,c} = f_{x|c} f_{,c}$ and $f_{xc} := (f_{x,c})$. Note that f_{xc} does behave like a join distribution.

The boxed equality above comes from the one-one correspondence between the variables $(f_{x|c}, f_c)$ and f_{xc} , so that the product of density functions for $f_{x|c}$ and f_c is just a specific case of a density function for f_{xc} , apart from a Jacobian factor.

The integral expression (14) is the representation of a fully exchangeable predictive distribution. Thus the joint distribution for the set of *pairs* of statements $\{(X_i = x_i, C_i = c_i)\}$ is fully exchangeable.

The noteworthy feature of the integral expression (14) for the fully exchangeable predictive distribution is that *the density for the joint distribution f_{xc} is factorizable into the product of a density for the conditional distribution $f_{x|c}$ and a density for the marginal distribution f_c* . This factorization expresses the conditional (i.e., partial) exchangeability for the statements $\{X_i = x_i\}$ given the $\{C_i = c_i\}$.

It is easy to show that the reverse also holds: if the density of an integral representation is factorizable as in (14b), then the corresponding probability distributions enjoy a symmetry of conditional exchangeability.

The factorization (14b) is not trivial. With a change of variables the following identities hold:

$$\begin{aligned} p(f_{xc} | I) df_{xc} &\equiv p[(f_{x|c}, f_c) | I] df_{x|c} df_c \equiv \\ & p(f_{x|c} | f_c, I) p(f_c | I) df_{x|c} df_c . \end{aligned} \quad (15)$$

The factorization condition is thus equivalent to conditional independence:

$$p(f_{x|c} | f_c, I) df_{x|c} = p(f_{x|c} | I) df_{x|c} . \quad (16)$$

3 Exchangeable Bayesian networks

3.1 Graphical representation of conditional exchangeability

The conditional (i.e., ‘partial’: see § 1) exchangeability (10) of the probabilities of the X-statements given the C-statements, the full exchangeability (12) of the probabilities of the latter, and the final integral representation (14) for their join probability distribution, can be all together expressed in the guise of a Bayesian network:

The two nodes represent the two sets of statements. The arrow from the C-node to the X-node represents the conditional exchangeability of the probabilities for the latter set of statements conditional on the former

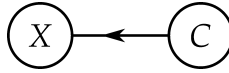


Figure 1

set. The absence of incoming arrows to the C-node represents the full exchangeability for its probability distribution. The final integral representation for the joint probability of the full network has a factorizable density, with one factor per node.

Using the reasoning and integral representations of §§ 1–2 it is indeed possible to generalize these rules to more complex networks of statements. It can be calculated that the following network, for example: has the representation

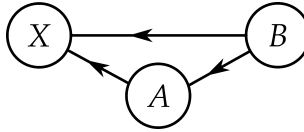


Figure 2

$$P(X_1 = x_1, A_1 = a_1, B_1 = b_1, \dots, X_N = x_N, A_N = a_N, B_N = b_N | I) = \int \prod_{x,a,b} \frac{f_{x,a,b}}{f_{x|a,b} f_{a|b} f_{,,b}}^{N_{F_{x,a,b}}} \underbrace{p(f_{x|ab} | I) p(f_{a|b} | I) p(f_b | I) df_{x|ab} df_{a|b} df_b}_{=p(f_{xab}|I) df_{xab}} \quad (17)$$

with $f_{,,b} := \sum_{x,a} f_{x,a,b}$ and so on.

3.2 Representing additional independence assumptions

In the last example the joint probability density for all the statements is decomposed in full accord with the product rule:

$$P(\{X_i = x_i\}, \{A_i = a_i\}, \{B_i = b_i\} | I) \equiv P(\{X_i = x_i\} | \{A_i = a_i\}, \{B_i = b_i\}, I) \times P(\{A_i = a_i\} | \{B_i = b_i\}, I) \times P(\{B_i = b_i\} | I), \quad (18)$$

which is an identity in the probability-calculus. In other words, no special properties of conditional independence hold.

In this case the integral representation (17) involves an integration over a set of conditional or marginal long-run frequencies which is equivalent to the set of joint frequencies. That is, the two sets

$$\{f_{x|a,b}, f_{a|b}, f_{.,b} \mid x \in \mathfrak{X}, a \in \mathfrak{A}, b \in \mathfrak{B}\} \leftrightarrow \{f_{x,a,b} \mid x \in \mathfrak{X}, a \in \mathfrak{A}, b \in \mathfrak{B}\} \quad (19)$$

are in one-one correspondence.

The factorizability of the density $p(f_{xab} \mid I) df_{xab}$ therefore implies, and is implied by, the exchangeability symmetries of conditional probability distributions. But it does not imply additional independences.

Additional independence properties, such as

$$\begin{aligned} P(\{X_i = x_i\}, \{A_i = a_i\}, \{B_i = b_i\} \mid I) = \\ P(\{X_i = x_i\} \mid \{A_i = a_i\}, \{B_i = b_i\}, I) \times \\ P(\{A_i = a_i\} \mid I) \times P(\{B_i = b_i\} \mid I), \quad (20) \end{aligned}$$

which is not an identity of the probability-calculus, are those typically expressed by Bayesian networks.

In the integral representation under discussion, these conditional independences are expressed by *a reduction in the number of conditional or marginal long-run frequencies that are integrated over* (similarly to what happens when partial exchangeability reduces to full exchangeability; see § 1).

For example, if the independence (20) hold, a little calculation shows that the representation (17) becomes

$$\begin{aligned} P(X_1 = x_1, A_1 = a_1, B_1 = b_1, \dots, X_N = x_N, A_N = a_N, B_N = b_N \mid I) = \\ \int \prod_{x,a,b} (f_{x|a,b} f_{a|b} f_{.,b})^{N_{f_{x,a,b}}} p(f_{x|ab} \mid I) p(f_a \mid I) p(f_b \mid I) df_{x|ab} df_a df_b. \quad (21) \end{aligned}$$

We see that the integration is now over the *reduced* set of long-run distributions

$$\{f_{x|a,b}, f_{a|b}, f_{.,b} \mid x \in \mathfrak{X}, a \in \mathfrak{A}, b \in \mathfrak{B}\}. \quad (22)$$

This is a reduced set in the sense that $f_{a|b} = f_{a|}$ for all b , implying the presence of a delta term in the corresponding density.

The independence condition (20), the various conditional-exchangeability conditions, and the integral representation (21) can be expressed by the network

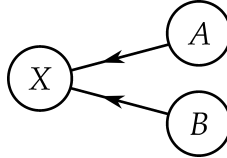


Figure 3

3.3 Generalization and connection with Bayesian networks

Let us summarize what we have found thus far. We have

0. several sets of statements, $\{X_i = x_i\}$, $\{A_i = a_i\}, \dots$. Each set is countably infinite. For each set we also have an associated set of values $x \in \mathfrak{X}, a \in \mathfrak{A}, \dots$;
- i. several assumptions of conditional independence representable in the guise of a Bayesian network, such as in eq. (20) and fig. 3;
- ii. several assumptions of conditional exchangeability for the predictive probabilities of some groups of such statements conditional on some other groups, such as in eq. (10).

These assumptions, taken together, must be sufficient to obtain the predictive, joint probability distribution for all statements through the product rule of the probability-calculus.

Then the predictive, joint probability distribution for N statements from each group has an integral representation of de Finetti form:

$$P(X_1 = x_1, A_1 = a_1, \dots, X_N = x_N, A_N = a_N, \dots \mid I) = \int \prod_{x,a,\dots} (\xi_{x,a,\dots})^{N F_{x,a,\dots}} p(\xi \mid I) d\xi \quad (23)$$

where $F_{x,a,\dots}$ are the joint empirical frequencies. This representation has these properties:

- I. the set of integration variables $\{\xi_{x,a}, \dots\}$ consists of subsets of marginal and conditional distributions: the same that would formally be associated with the Bayesian network of point [i](#). For example

$$\{\xi_{x,a,b}\} = \{\{f_{x|a,b}\}, \{f_{,a}\}, \{f_{,b}\}\}.$$

- II. the density $p(\xi | I) d\xi$ over the integration variables factorizes into a product of independent densities, one for each subset of the point above, for example

$$p(\xi | I) d\xi = p(f_{x|ab} | I) df_{x|ab} p(f_a | I) df_a p(f_b | I) df_b.$$

3.4 Discussion

The representation summarized in the preceding section combines exchangeability and Bayesian networks. I find that the combinations of these two notion is much needed.

A Bayesian network is a graphical representation of judgements (based for example on physical theories) about the informational relevance or irrelevance of some statements to other statements. Such relevance or irrelevance is expressed by equalities or inequalities between conditional probabilities involving those statements.

Unfortunately the use of Bayesian networks is a little ambiguous in the literature. Sometimes their application seem to refer to *individual* instances, and therefore to degrees of belief; for example, our degree of belief about the effect of a treatment upon a specific person, given the gender and health condition of that person. Sometimes their application seem to refer to whole populations or even superpopulations, and therefore to ‘long-run’ frequencies.⁵

In either case some questions arise. In the first case, the question of how our (conditional or unconditional) degrees of belief about the specific instance relate to, or are updated by, knowledge about similar instances. In the second case, the question of our uncertainty about the usually unknowable long-run frequencies, and about the empirical frequencies of future observations of small groups of individuals. These questions are really two sides of the same question.

⁵ e.g. Pearl [2009](#); Wiegerinck et al. [2013](#); the discussion by Lindley & Novick [1981](#) is more precise in this respect: they assume, for simplicity, a limit in which the two problems become numerically similar.

And these are precisely the questions addressed by exchangeability and its theorems. Exchangeability expresses judgements about ‘similarity’. More precisely, about the relevance or irrelevance of statements concerning individual events (observations, experiments, and similar) to sets of statements concerning other individual events. Such relevance or irrelevance is expressed by symmetries of conditional and unconditional probabilities involving those statements. The result is mathematical formula that quantitatively relates empirical frequencies from known observations (cf. $F_{x,a}$), long-run frequencies of superpopulations (cf. $f_{x,a}$) and degrees of beliefs about individual instances (cf. $P(X_1 = x_1 \mid \dots)$).

4 Discussion

*** first result can be useful for infinite limits, leading to regression

The result of the previous section is thus summarized: Given infinitely countable sets of statements $\{X_i = x_i\}$ and $\{C_i = c_i\}$, and assuming that

1. the marginal probability distribution for the C statements is fully exchangeable,
2. the probability distribution for the X statements is partially (or conditionally) exchangeable given the C ,

Then the joint distribution for both sets is fully exchangeable, and the density within its integral representation *factorizes* into a density for a conditional long-run frequency distribution, and a density for a marginal long-run frequency distribution, eq. (14b).

Appendix: On the ambiguous meaning of ‘ $X = x$ ’

Bibliography

- (‘de X ’ is listed under D , ‘van X ’ under V , and so on, regardless of national conventions.)
- Bernardo, J.-M., Smith, A. F. (2000): *Bayesian Theory*, repr. (Wiley, New York). First publ. 1994.
- Damien, P., Dellaportas, P., Polson, N. G., Stephens, D. A., eds. (2013): *Bayesian Theory and Applications*. (Oxford University Press, Oxford).
- Dawid, A. P. (2013): *Exchangeability and its ramifications*. In: Damien, Dellaportas, Polson, Stephens (2013), ch. 2, 19–29.
- Diaconis, P., Freedman, D. (1980a): *Finite exchangeable sequences*. *Ann. Prob.* **8**⁴, 745–764.
- (1980b): *De Finetti’s generalizations of exchangeability*. In: Jeffrey (1980), 233–249.

- Jeffrey, R. C., ed. (1980): *Studies in inductive logic and probability. Vol. II.* (University of California Press, Berkeley).
- Keynes, J. M. (2013): *A Tract on Monetary Reform*, repr. of 2nd ed. (Cambridge University Press, Cambridge). First publ. 1923.
- Lindley, D. V., Novick, M. R. (1981): *The role of exchangeability in inference.* Ann. Stat. **9**¹, 45–58.
- Pearl, J. (2009): *Causality: Models, Reasoning, and Inference*, 2nd ed. (Cambridge University Press, Cambridge). First publ. 2000.
- Wiegerinck, W., Burgers, W., Kappen, B. (2013): *Bayesian networks, introduction and practical applications.* In: **bianchinietal2013**, p. 12, 401–431.