

**Predictive inference:
From Bayesian inference
to Imprecise Probability**

Jean-Marc Bernard
University Paris Descartes
CNRS UMR 8069

Third SIPTA School on
Imprecise Probabilities

Montpellier, France
7 July 2008

INTRODUCTION

The “Bag of marbles” example

□ “Bag of marbles” problems (Walley, 1996)

- “I have ... a closed bag of coloured marbles. I intend to shake the bag, to reach into it and to draw out one marble. What is the probability that I will draw a red marble?”
- “Suppose that we draw a sequence of marbles whose colours are (in order):

blue, green, blue, blue, green, red.

What conclusions can you reach about the probability of drawing a red marble on a future trial?”

□ Two problems of predictive inference

- Prior prediction, before observing any item
- Posterior prediction, after observing n items

□ Inference from a state of prior ignorance about the proportions of the various colours

Categorical data (1)

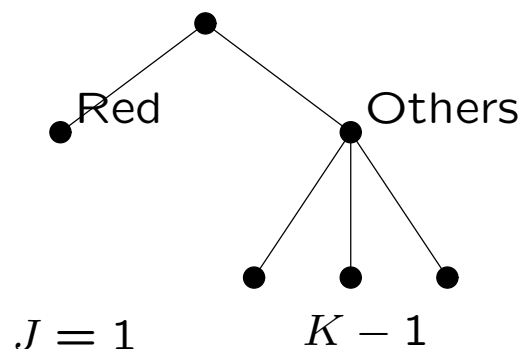
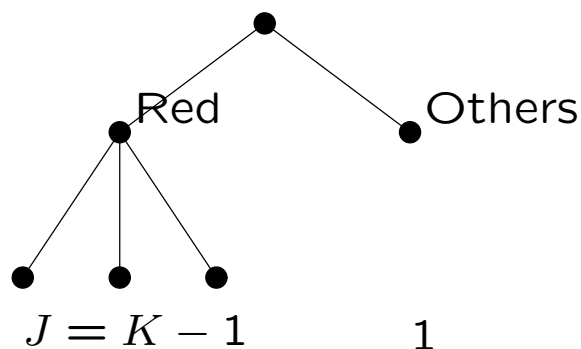
□ Categories

- Set of K of categories or types

$$C = \{c_1, \dots, c_K\}$$

- Categories c_k are exclusive and exhaustive
- Possible to add an extra category: “other colours”, “other types”

□ Categorisation is partly arbitrary



Categorical data (2)

□ Data

- Set, or sequence, I of n observations, items, individuals, *etc.*
- For each individual $i \in I$, we observe the corresponding category

$$\begin{aligned} I &\rightarrow C = \{c_1, \dots, c_K\} \\ i &\mapsto c_k \end{aligned}$$

- Observed composition, in **counts**:

$$\mathbf{a} = (a_1, \dots, a_K)$$

with $\sum_k a_k = n$

- Observed composition, in **frequencies**:

$$\mathbf{f} = (f_1, \dots, f_K) = \frac{\mathbf{a}}{n}$$

with $\sum_k f_k = 1$

□ **Compositions:** order considered as not important

Statistical inference problems (1)

□ Inference about what?

- **Predictive inference**: About future counts or frequencies in n' future observations

$$\mathbf{a}' = (a'_1, \dots, a'_K)$$

$$\mathbf{f}' = (f'_1, \dots, f'_K) = \mathbf{a}'/n'$$

$n' \geq 1$ Predictive inference (general)

$n' = 1$ Immediate prediction

- **Parametric inference**: About true/parent counts or frequencies (**parameters**) in population of

... size $N < \infty$

$$\mathbf{A} = (A_1, \dots, A_K)$$

$$\boldsymbol{\theta} = (\theta_1, \dots, \theta_K) = \mathbf{A}/N$$

... size $N = \infty$

$$\boldsymbol{\theta} = (\theta_1, \dots, \theta_K) \quad \sum_k \theta_k = 1$$

Statistical inference problems (2)

Prior vs. posterior inferences

□ Prior inferences

- $n = 0$ (no data yet)
- Unconditional
- Describes prior uncertainty about f' or θ
- Issue: formalize prior ignorance

□ Posterior inferences

- $n \geq 1$ (data a are available)
- Conditional on a
- Describes what can be inferred about f' or θ from the prior state + the knowledge of a

Relating past & future data (1)

Random sampling

□ Random sampling

- Population with a fixed, but unknown, true composition in frequencies

$$\theta = (\theta_1, \dots, \theta_K)$$

- Data (observed & future): random samples from the **same** population
- Ensures that the data are representative of the population *w.r.t.* C

□ Finite/infinite population

- Multiple-hypergeometric (N finite)
- Multinomial ($N = \infty$)

□ Stopping rule

- Fixed n
- Fixed a_k , “negative” sampling
- More complex stopping rules

□ These elements define a **sampling model**

Relating past & future data (2) Exchangeability

□ Exchangeability

- Consider any sequence S of $n^* = n + n'$ observations,

$$S = (c_1, \dots, c_n, c_{n+1}, \dots, c_{n^*})$$

having composition

$$\mathbf{a}^* = (a_1^*, \dots, a_K^*)$$

- Assumption of order-invariance, or permutation-invariance

$$\forall S, P(S | \mathbf{a}^*) = \text{constant}$$

□ Equivalence with MHyp sampling

Induced $P(\mathbf{a} | \mathbf{a}^*)$ is the same as if data with counts \mathbf{a} were obtained from random sampling from a population having counts $\mathbf{a}^* = \mathbf{a} + \mathbf{a}'$

□ **Direct link:** No need to invoke unknown parameters θ of a larger population

A statistical challenge

□ Model prior ignorance

- Model prior ignorance about θ , or a and a^*
- Arbitrariness of C and K , both may vary as data items are observed
- Model prior ignorance about both the set C and the number K of categories

□ Make reasonable posterior inferences

from such a state of prior ignorance

- Idea of “objective” methods: “let the data speak for themselves”
- Frequentist methods
- Objective Bayesian methods

□ “Reasonable”: Several desirable principles

Desirable principles / properties (1)

□ Prior ignorance

- **Symmetry (SP)**: Prior uncertainty should be invariant *w.r.t.* permutations of categories
- **Embedding pcple (EP)**: Prior uncertainty should not depend on refinements or coarsenings of categories

□ Independence from irrelevant information of posterior inferences

- **Stopping rule pcple (SRP)**: Inferences should not depend on the stopping rule, *i.e.* on data that might have occurred but have actually not
- **Likelihood pcple (LP)**: Inferences should depend on the data through the likelihood function only
- **Representation invariance (RIP)**: Posterior inferences should not depend on refinements or coarsenings of categories

Desirable principles / properties (2)

- **Reasonable account of uncertainty** in prior and posterior inferences

- **Consistency requirements** when considering several inferences
 - **Avoiding sure loss (ASL)**: Probabilistic assessments, when interpreted as betting dispositions, should not jointly lead to a sure loss
 - **Coherence (CP)**: Stronger property of consistency of all probabilistic assessments

- **Frequentist interpretation(s)**
 - **Repeated sampling principle (RSP)**: Probabilities should have an interpretation as relative frequencies in the long run

- **See Walley, 1996; 2002**

Methods for statistical inference: Frequentist approach

□ Frequentists methods

- Based upon **sampling model only** e.g. $a|\theta$
- Probabilities can be assimilated to long-run frequencies
- Significance tests, confidence limits and intervals (Fisher, Neyman & Pearson)

□ Difficulties of frequentist methods

- Depend on the stopping rule. Hence do not obey SRP, nor LP
- Not conditional on observed data; May have relevant subsets
- For multidimensional parameters' space: ad-hoc and/or asymptotic solutions to the problem of nuisance parameters

Methods for statistical inference: Objective Bayesian approach (1)

□ Bayesian methods

- Two ingredients: **sampling model + prior**
- Conjugate priors: **Dirichlet** for multinomial data, **Dirichlet-multinomial** for multiple-hypergeometric data
- Depend on the sampling model through the likelihood function only

□ Objective Bayesian methods

- Data analysis goal: let the data say what they have to say about unknown parameters
- Priors formalizing “prior ignorance”
- objective Bayesian: “non-informative” priors, *etc.* (e.g. **Kass, Wasserman, 1996**)
- Exact or approximate frequentist reinterpretations: “matching priors” (e.g. **Datta, Ghosh, 1995**)

Methods for statistical inference: Objective Bayesian approach (2)

□ **Difficulties of Bayesian methods** for categorical data

Several priors proposed for prior ignorance, but none satisfies all desirable principles.

- Inferences often depend on C and/or K
- Some solutions violate LP (Jeffreys, 1946)
- Some solutions can generate incoherent inferences (Berger, Bernardo, 1992)
- If $K = 2$, uncertainty about next observation (case $n' = 1$) is the same whether $a_1 = a_2 = 0$ (prior) or $a_1 = a_2 = 100$ (posterior)

$$P(a' = (1, 0)) = P(a' = (1, 0) | a)$$

□ **Only approximate agreement** between frequentist methods and objective Bayesian methods, for categorical data

The IDM in brief

□ **Model for parametric inference** for categorical data

Proposed by Walley (1996), generalizes the IBM (Walley, 1991).

Inference from data $\mathbf{a} = (a_1, \dots, a_K)$, categorized in K categories C , with unknown chances $\boldsymbol{\theta} = (\theta_1, \dots, \theta_K)$.

□ **Imprecise probability model**

Prior uncertainty about $\boldsymbol{\theta}$ expressed by a set of Dirichlet's.

Posterior uncertainty about $\boldsymbol{\theta}|\mathbf{a}$ then described by a set of updated Dirichlet's.

Generalizes Bayesian inference, where prior/ posterior uncertainty is described by a *single* Dirichlet.

□ **Imprecise U&L probabilities**, interpreted as reasonable betting rates *for* or *against* an event.

□ **Models prior ignorance** about $\boldsymbol{\theta}$, K and C

□ **Satisfies desirable principles** for inferences from prior ignorance, contrarily to alternative frequentist and objective Bayesian approaches.

The IDMM in brief

□ **Model for predictive inference** for categorical data

Proposed by Walley, Bernard (1999), also partly studied in (Walley, 1996).

Inference about future data $\mathbf{a}' = (a'_1, \dots, a'_K)$ from observed data $\mathbf{a} = (a_1, \dots, a_K)$, categorized in K categories C .

□ **Two alternative, equivalent views**

- A predictive model derived from the parametric IDM
- A model of its own, modeling only observables: available data \mathbf{a} and future data \mathbf{a}'

□ **Imprecise probability model**

Prior uncertainty about \mathbf{a} expressed by a set of Dirichlet-multinomial distributions.

Posterior uncertainty about $\mathbf{a}'|\mathbf{a}$ then described by a set of updated Dirichlet-multinomial distributions.

□ **Models prior ignorance** about \mathbf{a} , K and C

Outline

1. Introduction
2. Bayesian approach to inference
3. Important distributions
4. Objective Bayesian models
5. From Bayesian to imprecise probability models
6. Definition of the IDM & the IDMM
7. Predictive inferences from the IDMM
8. The rule of succession
9. Conclusions

References

THE BAYESIAN APPROACH

Bayesian inference

□ Focus on the Bayesian approach since

- Bayesian, precise: a single Dirichlet prior on θ yields a single Dirichlet posterior on $\theta|a$ (PDM)
- IP-model: a prior set of Dirichlet's yields a posterior set of Dirichlet's (IDM)

□ ... and for predictive inferences since

- Bayesian, precise: a single Dirichlet-Multinomial (*DiMn*) prior on a^* yields a single *DiMn* posterior on $a'|a$ (PDMM)
- IP-model: a prior set of *DiMn*'s yields a posterior set of *DiMn*'s (IDMM)

□ Goal

- Sketch Bayesian approach to inference
- Specifically: objective Bayesian models
- Indicate shortcomings of these models

Three sampling models

□ Multinomial data

- Random sampling
- Infinite population, $N = \infty$
- Data have a multinomial (Mn) likelihood

□ Multiple-hypergeometric data

- Random sampling
- Finite population, $N < \infty$
- Data have a multiple-hypergeometric ($MHyp$) likelihood

□ Exchangeable data

- Data a generated by an exchangeable process with counts $a^* = a + a'$
- Data have a $MHyp$ likelihood too

□ Hypotheses

- Set C , and number of categories, K , are considered as known and fixed

Inference from multinomial data

□ Multinomial data

- Elements of population are categorized in K categories from set $C = \{c_1, \dots, c_K\}$.
- Unknown true chances $\theta = (\theta_1, \dots, \theta_K)$, with $\theta_k \geq 0$ and $\sum_k \theta_k = 1$, i.e. $\theta \in \Theta = \mathcal{S}(1, K)$.
- Data are a random sample of size n from the population, yielding counts $\mathbf{a} = (a_1, \dots, a_K)$, with $\sum_k a_k = n$.

□ Multinomial sampling distribution

$$P(\mathbf{a}|\theta) = \binom{n}{\mathbf{a}} \theta_1^{a_1} \dots \theta_K^{a_K}$$

When seen as a function of θ , leads to the **likelihood function**

$$L(\theta|\mathbf{a}) \propto \theta_1^{a_1} \dots \theta_K^{a_K}$$

□ **Same likelihood** is obtained from observing \mathbf{a} , for a variety of stopping rules: n fixed, a_k fixed, etc.

Bayesian inference (1): a learning model

□ General scheme

$$\left\{ \begin{array}{c} \text{Prior } P(\theta) \\ + \\ \text{Sampling } P(a|\theta) \end{array} \right. \longrightarrow \left\{ \begin{array}{c} \text{Posterior } P(\theta|a) \\ + \\ \text{Prior predictive } P(a) \end{array} \right.$$

□ Iterative process

$$\left\{ \begin{array}{c} \text{Prior}' } P(\theta|a) \\ + \\ \text{Sampl.}' } P(a'|\theta, a) \end{array} \right. \longrightarrow \left\{ \begin{array}{c} \text{Posterior}' } P(\theta|a', a) \\ + \\ \text{Post. pred. } P(a'|a) \end{array} \right.$$

□ Learning model about

- unknown chances: $P(\theta)$ updated to $P(\theta|a)$
- future data: $P(a)$ updated to $P(a'|a)$

Bayesian inference (2)

□ Continuous parameters space

Since the parameters space, Θ , is continuous, probabilities on θ , $P(\theta)$ and $P(\theta|a)$, are defined via densities, denoted $h(\theta)$ and $h(\theta|a)$

□ Bayes' theorem (or rule)

$$\begin{aligned} h(\theta|a) &= \frac{h(\theta) P(a|\theta)}{\int_{\Theta} h(\theta) P(a|\theta) d\theta} \\ &= \frac{h(\theta) L(\theta|a)}{\int_{\Theta} h(\theta) L(\theta|a) d\theta} \end{aligned}$$

□ **Likelihood principle** satisfied if prior $h(\theta)$ is chosen independently of $P(a|\theta)$

□ Conjugate inference

- Prior $h(\theta)$ and posterior $h(\theta|a)$ are from the same family
- For multinomial likelihood: **Dirichlet** family

Dirichlet prior for θ

□ Dirichlet prior

Prior uncertainty about θ is expressed by

$$\theta \sim \text{Diri}(\alpha)$$

with prior strengths

$$\alpha = (\alpha_1, \dots, \alpha_K)$$

such that $\alpha_k > 0$, $\sum_k \alpha_k = s$

□ Dirichlet distribution

Density defined for any $\theta \in \Theta$, with $\Theta = \mathcal{S}(1, K)$

$$h(\theta) = \frac{\Gamma(s)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_K)} \theta_1^{\alpha_1-1} \cdots \theta_K^{\alpha_K-1}$$

□ Generalisation of the Beta distribution

$$(\theta_1, 1 - \theta_1) \sim \text{Diri}(\alpha_1, \alpha_2) \iff \theta_1 \sim \text{Beta}(\alpha_1, \alpha_2)$$

Alternative parameterization

□ Dirichlet prior on θ

$$\theta \sim \text{Diri}(\alpha)$$

□ Alternative parameterization in terms of s , the total prior strength, and the relative prior strengths

$$t = (t_1, \dots, t_K) = \alpha/s$$

with $t_k > 0$, $\sum_k t_k = 1$, i.e. $t \in \mathcal{S}^*(1, K)$

Hence,

$$\theta \sim \text{Diri}(st)$$

□ Prior expectation of θ_k

$$E(\theta_k) = t_k$$

□ Interpretation

- t determines the center of the distribution
- s determines its dispersion / concentration

Dirichlet posterior for $\theta|a$

□ Dirichlet posterior

Posterior uncertainty about $\theta|a$ is expressed by

$$\begin{aligned}\theta|a &\sim \text{Diri}(a + \alpha) \\ &\sim \text{Diri}(a + st)\end{aligned}$$

Parameters/strengths of the Dirichlet play a role of **counters**: the prior strength α_k is incremented by the observed count a_k to give the posterior strength $a_k + \alpha_k$

□ Posterior expectation of θ_k

$$\begin{aligned}E(\theta_k|a) &= \frac{a_k + \alpha_k}{n + s} \\ &= \frac{nf_k + st_k}{n + s}\end{aligned}$$

i.e. a weighted average of prior expectation, t_k , and observed frequency, f_k , with weights s and n

Prior predictive distribution

□ From Bayes theorem

$$h(\theta|a) = \frac{h(\theta) P(a|\theta)}{\int_{\Theta} h(\theta) P(a|\theta) d\theta}$$

□ Prior predictive distribution on a

$$\begin{aligned} P(a) &= \int_{\Theta} h(\theta) P(a|\theta) d\theta \\ &= \frac{h(\theta) P(a|\theta)}{h(\theta|a)} \end{aligned}$$

which yields

$$P(a) = \frac{\prod_k \binom{a_k + \alpha_k - 1}{a_k}}{\binom{n + s - 1}{n}}$$

with $\binom{m+x-1}{m} = \frac{\Gamma(m+x)}{m!\Gamma(x)}$, for any positive integer $m \geq 0$, and any real $x > 0$

□ Dirichlet-multinomial distribution

$$a \sim \text{DiMn}(n; \alpha)$$

Posterior predictive distribution

- Similarly, from Bayes theorem

$$\begin{aligned} P(a'|a) &= \frac{h(\theta|a) P(a'|\theta, a)}{h(\theta|a', a)} \\ &= \frac{h(\theta|a) P(a'|\theta)}{h(\theta|a' + a)} \end{aligned}$$

which yields

$$P(a'|a) = \frac{\prod_k \binom{a'_k + a_k + \alpha_k - 1}{a'_k}}{\binom{n' + n + s - 1}{n'}}$$

- Dirichlet-multinomial posterior

$$a'|a \sim \text{DiMn}(n'; a + \alpha)$$

- Interpretation in terms of “counters”

Here too, prior strengths α are updated into posterior strengths $a + \alpha$

Equivalence of 3 models for predictive inference

□ Multinomial + Dirichlet model

$$\left\{ \begin{array}{l} \theta \sim \text{Diri (Prior)} \\ a|\theta \sim \text{Mn (Samp.)} \\ a'|\theta, a \sim \text{Mn (Samp.)} \end{array} \right. \longrightarrow \left\{ \begin{array}{l} a \sim \text{DiMn} \\ + \\ a'|a \sim \text{DiMn} \end{array} \right.$$

□ M.-Hypergeometric + DiMn model

$$\left\{ \begin{array}{l} A \sim \text{DiMn (Prior)} \\ a|A \sim \text{MHyp (Samp.)} \\ a'|A, a \sim \text{MHyp (Samp.)} \end{array} \right. \longrightarrow \left\{ \begin{array}{l} a \sim \text{DiMn} \\ + \\ a'|a \sim \text{DiMn} \end{array} \right.$$

□ Exchangeability + DiMn model

$$\left\{ \begin{array}{l} a^* \sim \text{DiMn (Prior)} \\ a|a^* \sim \text{MHyp (Samp.)} \\ a'|a^*, a \sim \text{MHyp (Samp.)} \end{array} \right. \longrightarrow \left\{ \begin{array}{l} a \sim \text{DiMn} \\ + \\ a'|a \sim \text{DiMn} \end{array} \right.$$

Bayesian answers to inference (1)

Parametric problems

□ **Prior uncertainty:** $P(\theta)$

□ **Posterior uncertainty:** $P(\theta|\mathbf{a})$

For drawing all inferences, from observed data to unknown parameters

□ **Inferences** about θ

- Expectations, $E(\theta_k|\mathbf{a})$; Variances, $Var(\theta_k|\mathbf{a})$; etc.
- Any event about θ : $P(\theta \in \Theta^* | \mathbf{a})$

□ **Inferences** about real-valued $\lambda = g(\theta)$

- Marginal distribution function: $h(\lambda|\mathbf{a})$
- Expectation, variance: $E(\lambda|\mathbf{a})$, $Var(\lambda|\mathbf{a})$
- Cdf: $F_\lambda(u) = P(\lambda < u|\mathbf{a}) = \int_{-\infty}^u h(\lambda|\mathbf{a}) d\lambda$
- Credibility intervals: $P(\lambda \in [u_1; u_2] | \mathbf{a})$
- Any event about λ

Bayesian answers to inference (2)

Predictive problems

□ **Prior uncertainty:** $P(a)$ or $P(f)$

□ **Posterior uncertainty:** $P(a'|a)$ or $P(f'|a)$

For drawing all inferences, from observed data to future data

□ **Inferences** about f'

- Expectations, $E(f'_k|a)$; Variances, $Var(f'_k|a)$; etc.
- Any event about f' : $P(f' \in \Theta^* | a)$

□ **Inferences** about real-valued $\lambda = g(f')$

- Marginal distribution function: $P(\lambda|a)$
- Expectation, variance: $E(\lambda|a)$, $Var(\lambda|a)$
- Cdf: $F_\lambda(u) = P(\lambda < u|a) = \sum_{\lambda < u} P(\lambda|a)$
- Credibility intervals: $P(\lambda \in [u_1; u_2] | a)$
- Any event about λ

IMPORTANT DISTRIBUTIONS

Relevant distributions

□ Parametric inference on infinite population

- Dirichlet (*Diri*), any K
- Beta (*Beta*), $K = 2$

□ Predictive inference on future n' data

- Dirichlet-Multinomial (*DiMn*), any K
- Beta-Binomial (*BeBi*), $K = 2$

□ Links

	n'	$n' \rightarrow \infty$
$K = 2$	<i>BeBi</i>	<i>Beta</i>
K	<i>DiMn</i>	<i>Diri</i>

Beta distribution

□ **Consider** the variable

$$\theta \in [0, 1]$$

and the hyper-parameters

$$\alpha_1 > 0, \alpha_2 > 0$$

or $s = \alpha_1 + \alpha_2$, $t_1 = \alpha_1/s$, $t_2 = \alpha_2/s$,
with $s > 0$, $t_1 > 0$, $t_2 > 0$, $t_1 + t_2 = 1$

□ **Beta density**

$$\theta \sim \text{Beta}(\alpha_1, \alpha_2) = \text{Beta}(st_1, st_2)$$

$$h(\theta) = \frac{\Gamma(s)}{\Gamma(\alpha_1)\Gamma(\alpha_2)} \theta^{\alpha_1-1} (1-\theta)^{\alpha_2-1}$$
$$\propto \theta_1^{\alpha_1-1} (1-\theta)^{\alpha_2-1}$$

□ **Expectation and variance**

$$E(\theta) = \alpha_1/s = t_1$$

$$\text{Var}(\theta) = \frac{\alpha_1\alpha_2}{s^2(s+1)} = \frac{t_1t_2}{s+1}$$

Dirichlet distribution

□ Consider

$$\boldsymbol{\theta} = (\theta_1, \dots, \theta_K) \quad \boldsymbol{\theta} \in \Theta = \mathcal{S}(1, K)$$

$$\boldsymbol{t} = (t_1, \dots, t_K) \quad \boldsymbol{t} \in \mathcal{T} = \mathcal{S}^*(1, K)$$

and $s > 0$, or $\boldsymbol{\alpha} = \boldsymbol{st}$, $\alpha_k > 0$

□ Dirichlet density

$$\boldsymbol{\theta} \sim \text{Diri}(\boldsymbol{\alpha}) = \text{Diri}(\boldsymbol{st})$$

$$h(\boldsymbol{\theta}) = \frac{\Gamma(s)}{\prod_k \Gamma(\alpha_k)} \theta_1^{\alpha_1-1} \dots \theta_K^{\alpha_K-1}$$
$$\propto \theta_1^{\alpha_1-1} \dots \theta_K^{\alpha_K-1}$$

□ Generalization of Beta distribution ($K = 2$)

$$(\theta_1, \theta_2) \sim \text{Diri}(\alpha_1, \alpha_2) \iff \theta_1 \sim \text{Beta}(\alpha_1, \alpha_2)$$

□ Basic properties

- $E(\theta_k) = t_k$
- s determines dispersion of distribution

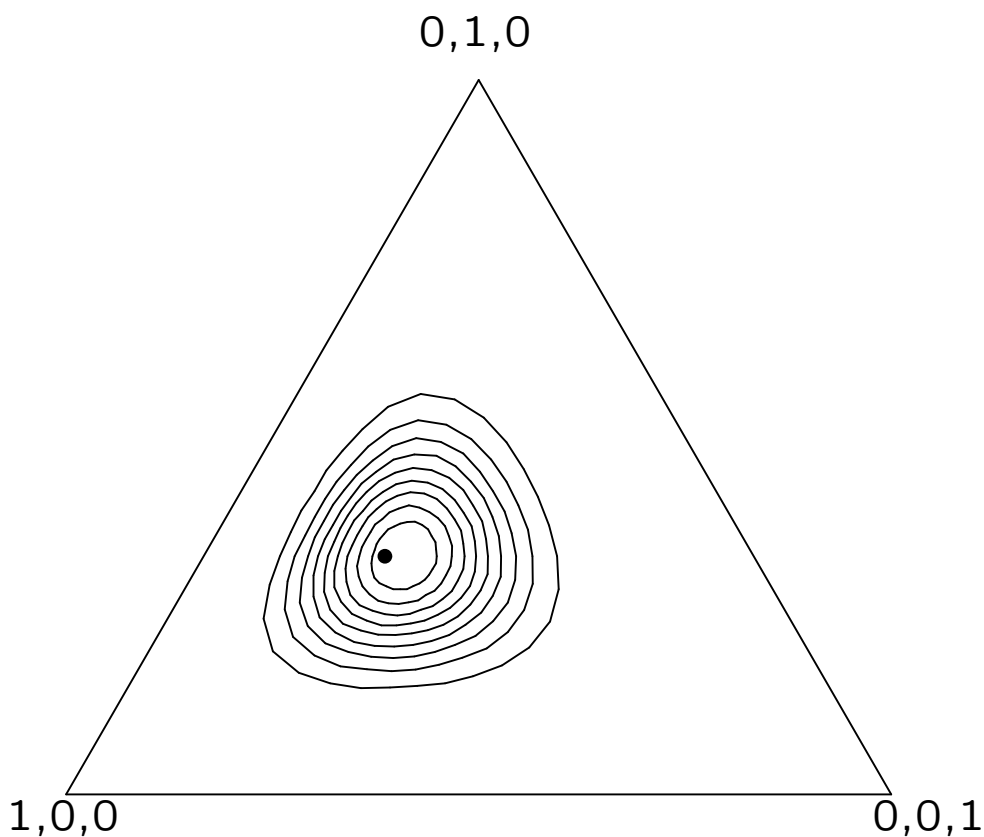
Examples of Dirichlet's

□ Example 1

$Diri(1, 1, \dots, 1)$ is uniform on $\mathcal{S}(1, K)$

□ Example 2

$(\theta_1, \theta_2, \theta_3) \sim Diri(10, 8, 6)$



(Highest density contours: [100%, 90%, ..., 10%])

Properties of the Dirichlet

General properties given on an example.

Assume $(\theta_1, \dots, \theta_5) \sim \text{Diri}(\alpha_1, \dots, \alpha_5)$. Then,

□ Pooling property

$$(\theta_1, \theta_{234}, \theta_5) \sim \text{Diri}(\alpha_1, \alpha_{234}, \alpha_5),$$

where pooling categories amounts to add corresponding chances, $\theta_{234} = \theta_2 + \theta_3 + \theta_4$, and strengths, $\alpha_{234} = \alpha_2 + \alpha_3 + \alpha_4$.

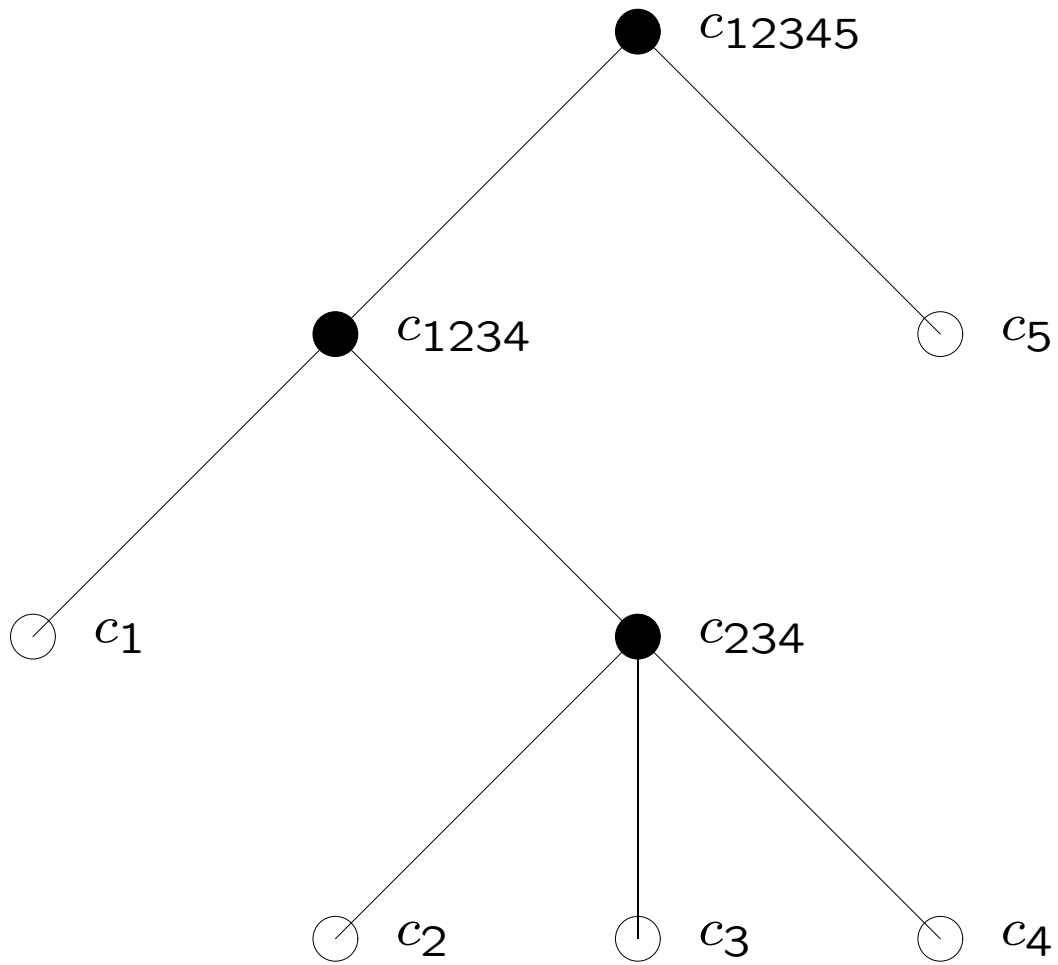
□ Restriction property

$$(\theta_2^{234}, \theta_3^{234}, \theta_4^{234}) \sim \text{Diri}(\alpha_2, \alpha_3, \alpha_4),$$

where $\theta_2^{234} = \theta_2 / \theta_{234}$, etc., are conditional chances.

□ Generalizes to any tree underlying the set C .

Tree representation of categories



Beta-Binomial distribution (1)

□ Notation

$$(a_1, a_2) \sim \text{BeBi}(n; \alpha_1, \alpha_2)$$

for a_1 and a_2 positive integers, with $a_1 + a_2 = n$
and $\alpha_1 > 0$ and $\alpha_2 > 0$, with $\alpha_1 + \alpha_2 = s$

□ Probability distribution function

$$\begin{aligned} P(a_1, a_2) &= \frac{\binom{a_1 + \alpha_1 - 1}{a_1} \binom{a_2 + \alpha_2 - 1}{a_2}}{\binom{n + s - 1}{n}} \\ &= \frac{\Gamma(a_1 + \alpha_1)}{a_1! \Gamma(\alpha_1)} \frac{\Gamma(a_2 + \alpha_2)}{a_2! \Gamma(\alpha_2)} \frac{n! \Gamma(s)}{\Gamma(n + s)} \\ &= \binom{n}{a_1} \frac{\alpha_1^{[a_1]} \alpha_2^{[a_2]}}{s^{[n]}} \end{aligned}$$

Beta-Binomial distribution (2)

□ **Expectation & variance** of a_1 and $f_1 = a_1/n$

$$E(a_1) = n \frac{\alpha_1}{s} = nt_1$$

$$E(f_1) = t_1$$

$$\text{Var}(f_1) = \frac{t_1(1-t_1)}{s+1} \times \frac{n+s}{n}$$

where $t_1 = \alpha_1/s$, $1 - t_1 = t_2 = \alpha_2/s$

□ **Convergence** of distribution of f_1

$$t_1 \rightarrow \text{Beta}(\alpha_1, \alpha_2)$$

when $n \rightarrow \infty$

Dirichlet-Multinomial distribution

□ Notation

$$\mathbf{a} \sim \text{DiMn}(n; \boldsymbol{\alpha})$$

for $\mathbf{a} = (a_1, \dots, a_K)$, a_k positive ints, $\sum_k a_k = n$
and $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_K)$, $\alpha_k > 0$, $\sum_k \alpha_k = s$

□ Probability distribution function

$$\begin{aligned} P(\mathbf{a}) &= \frac{\prod_k \binom{a_k + \alpha_k - 1}{a_k}}{\binom{n + s - 1}{n}} \\ &= \frac{n! \Gamma(s)}{\Gamma(n + s)} \prod_k \frac{\Gamma(a_k + \alpha_k)}{a_k! \Gamma(\alpha_k)} \\ &= \binom{n}{\mathbf{a}} \frac{\prod_k \alpha_k^{[a_k]}}{s^{[n]}} \end{aligned}$$

Mathematical functions or coefficients

□ Binomial coefficient

$$\binom{n}{a} = \frac{n!}{a!(n-a)!}$$

for n, a integers, $n \geq a$

□ Multinomial coefficients

$$\binom{n}{a} = \frac{n!}{a_1! \cdots a_k!}$$

for $a = (a_1, \dots, a_k)$ integers, $\sum_k a_k = n$

□ Generalized binomial coefficients

$$\binom{m+x-1}{m} = \frac{\Gamma(m+x)}{m! \Gamma(x)}$$

for integer $m \geq 0$, and real $x > 0$

□ Ascending factorial (from Appell ?)

$$x^{[m]} = x(x+1)\cdots(x+m-1), \quad x^{[0]} = 1$$

for integer $m \geq 0$, and real x

OBJECTIVE BAYESIAN MODELS

Objective Bayesian models

□ Priors proposed for objective inference

Idea: α expressing prior ignorance about θ or a^*
(Kass & Wasserman, 1996; Bernard, 1996)

□ For direct (Mn or MHyp) sampling

Almost all proposed solutions for fixed n are **sym-metric** Dirichlet priors, *i.e.* $t_k = 1/K$:

- Haldane (1948): $\alpha_k = 0$ ($s = 0$)
- Perks (1947): $\alpha_k = \frac{1}{K}$ ($s = 1$)
- Jeffreys (1946): $\alpha_k = \frac{1}{2}$ ($s = K/2$)
- Bayes-Laplace, uniform: $\alpha_k = 1$ ($s = K$)
- Berger-Bernardo reference priors

□ For negative (Mn or MHyp) sampling

Some proposed solutions for fixed a_k are *non-symmetric* Dirichlet priors

Which principles are satisfied? (1)

□ Prior ignorance

- **Symmetry (SP)**. Yes: for all usual symmetric priors with $t_k = 1/K$. No: for some priors proposed for negative-sampling.
- **Embedding Principle (EP)**. Yes: for Haldane's prior. No: for all other priors

□ Internal consistency

- **Coherence (CP)**, including ASL. Yes: if prior is proper. No: for Haldane's improper prior.

□ Frequentist interpretation

- **Repeated sampling principle (RSP)**. No in general. Yes asymptotically. Exact or conservative agreement for some procedures.

Which principles are satisfied? (2)

□ Invariance, Independence from irrelevant information

- **Likelihood principle (LP)**, including SRP. Yes, if prior ($P(\theta)$ or $P(a^*)$) chosen independently of sampling model ($P(a|\theta)$ or $P(a|a^*)$). No, for Jeffreys' or Berger-Bernardo's priors
- **Representation invariance (RIP)**. Yes: Haldane. No: all other priors
- **Invariance by reparameterisation**. Yes, for Jeffreys' or Berger-Bernardo's priors

□ Difficulties of objective Bayesian approach

None of these solutions simultaneously satisfies all desirable principles for inferences from prior ignorance

Focus on Haldane's prior

□ Satisfies most principles

- Satisfies most of the principles: symmetry, LP, EP and RIP
- Incoherent because of impropriety, but can be extended to a coherent model (Walley, 1991)

□ But

- Improper prior
- Improper posterior if some $a_k = 0$
- Too data-glued:
If $a_k = n = 1$, essentially says that $\theta_k = 1$, or that $a'_k = n'$, with probability 1.
If $a_k = 0$, essentially says that $\theta_k = 0$, or that $a'_k = 0$ for any n' , with probability 1.
- Doesn't give a reasonable account of uncertainty.

□ Limit case of the ID(M)M

**FROM PRECISE
BAYESIAN MODELS
TO AN IMPRECISE
PROBABILITY MODEL**

Precise Bayesian Dirichlet model

□ Elements of a (precise) standard Bayesian model

- Prior distribution: $P(\theta)$, $\theta \in \Theta$
- Sampling distribution: $P(a|\theta)$, $a \in \mathcal{A}$, $\theta \in \Theta$
- Posterior distribution: $P(\theta|a)$, $\theta \in \Theta$, $a \in \mathcal{A}$, obtained by Bayes' theorem

□ Elements of a precise Dirichlet model

- Dirichlet $P(\theta)$
- Multinomial $P(a|\theta)$
- Dirichlet $P(\theta|a)$

Probability vs. Prevision (1)

□ Three distributions

$$P(\theta) \quad P(a|\theta) \quad P(\theta|a)$$

These are probability distributions, which allocate a mass probability (or a probability density) to any event relative to θ and/or a .

□ From probability of events to previsions of gambles

Since each one is a precise model, each defines a unique linear prevision for each possible gamble. So, each $P(\cdot)$ or $P(\cdot|\cdot)$ can be assimilated to a linear prevision

□ Domains of these linear previsions

Here, we always consider **all possible gambles**, so these linear previsions are each defined on the linear space of all gambles (on their respective domains).

Probability vs. Prevision (2)

Remarks

□ Remark on terms used

- Random quantity = Gamble
- Expectation = Prevision

□ Previsions of gambles are more fundamental than probabilities of events

- Precise world:

Previsions \iff Probabilities

- Imprecise world:

Previsions \implies Probabilities

□ See (de Finetti, 1974-75; Walley, 1991)

Coherence of a standard Bayesian model

□ Coherence of these linear previsions

- If prior is proper, then $P(\theta)$ is coherent
- $P(a|\theta)$ always coherent
- If prior is proper, then posterior is proper, and hence $P(\theta|a)$ is coherent

□ Joint coherence (Walley, 1991, Thm. 7.7.2)

- The linear previsions, $P(\theta)$, $P(a|\theta)$ and $P(\theta|a)$ are jointly coherent
- This is assured by **generalized Bayes' rule**, which reduces to Bayes' rule/theorem in the case of linear previsions.

Class of coherent models

- **One privileged way** of constructing coherent imprecise posterior probabilities

“... is to form the lower envelopes of a class of standard Bayesian priors and the corresponding class of standard Bayesian posteriors”

(Walley, 1991, p. 397)

- **Lower envelope theorem** (id., Thm. 7.1.6)

The lower envelope of a class of separately coherent lower previsions, is a coherent lower prevision.

- **Class of Bayesian models** (id., Thm. 7.8.1):

Suppose that $P_\gamma(\cdot)$, $P_\gamma(\cdot|\Theta)$ and $P_\gamma(\cdot|\mathcal{A})$ constitute a standard Bayesian model, for every $\gamma \in \Gamma$. Then their lower envelopes, $\underline{P}(\cdot)$, $\underline{P}(\cdot|\Theta)$ and $\underline{P}(\cdot|\mathcal{A})$ are coherent.

Towards the IDM & the IDMM

□ Building an Imprecise Dirichlet model

- **Class** of Dirichlet priors
- A **single precise** Mn sampling model
- Update each prior, using Bayes' theorem
- **Class** of Dirichlet posteriors
- Form the associated posterior lower prevision

□ ... or an Imprecise Dirichlet-multinomial model

- **Class** of Dirichlet-multinomial priors
- A **single precise** $MHyp$ sampling model
- Update each prior, using Bayes' theorem
- **Class** of Dirichlet-multinomial posteriors
- Form the associated posterior lower prevision

The IDM & IDMM

Class of priors for the IDM & the IDMM

□ **Models proposed** by Walley (1996) for the IDM, and by Walley, Bernard (1999) for the IDMM.

□ **Which prior class?**

Choosing a *Diri* or a *DiMn* prior amounts to choosing prior strengths

$$\begin{aligned}\alpha &= (\alpha_1, \dots, \alpha_K) \\ &= s t \\ &= s (t_1, \dots, t_K)\end{aligned}$$

In the IDM or the IDMM

- Fix the total prior strength s
- Let t take all possible values in $\mathcal{T} = \mathcal{S}^*(1, K)$

□ **Yielding which properties?**

- Nice properties for modeling prior ignorance
- Satisfy several desirable principles

Prior and posterior IDM

□ Prior IDM

The prior IDM(s) is defined as the set \mathcal{M}_0 of all Dirichlet distributions on θ with a fixed total prior strength $s > 0$:

$$\mathcal{M}_0 = \{Diri(st) : t \in \mathcal{T} = \mathcal{S}^*(1, K)\}$$

□ Posterior IDM

Posterior uncertainty about θ , conditional on a , is expressed by the set

$$\mathcal{M}_n = \{Diri(a + st) : t \in \mathcal{T} = \mathcal{S}^*(1, K)\}.$$

□ Updating

Each Dirichlet distribution on θ in the set \mathcal{M}_0 is updated into another Dirichlet on $\theta|a$ in the set \mathcal{M}_n , using Bayes' theorem.

This procedure guarantees the **coherence** of inferences (Walley, 1991, Thm. 7.8.1).

Prior and posterior IDMM

□ Prior IDMM

The prior IDMM(s) is defined as the set \mathcal{M}_0 of all Dirichlet-Multinomial distributions on \mathbf{a}^* with a fixed total prior strength $s > 0$:

$$\mathcal{M}_0 = \{DiMn(n^*; st) : t \in \mathcal{T} = \mathcal{S}^*(1, K)\}$$

□ Posterior IDMM

Posterior uncertainty about \mathbf{a}' , conditional on \mathbf{a} , is expressed by the set

$$\mathcal{M}_n = \{DiMn(n'; \mathbf{a} + st) : t \in \mathcal{T} = \mathcal{S}^*(1, K)\}.$$

□ Updating

Similarly, each *DiMn* distribution on \mathbf{a}^* in the set \mathcal{M}_0 is updated into another *DiMn* on $\mathbf{a}'|\mathbf{a}$ in the set \mathcal{M}_n .

□ Counts / frequencies

Prior on \mathbf{a}^* or \mathbf{f}^* , posterior on $\mathbf{a}'|\mathbf{a}$ or $\mathbf{f}'|\mathbf{a}$.

Drawing inferences from the IDM or IDMM

□ Events, indicator functions

- Compute lower & upper (L&U) probabilities of events of interest
- Substantial conclusion if lower probability is sufficiently large

□ Random quantities

- Compute L&U cumulative distribution functions (cdf)
- Compute L&U expectations
- Compute L&U variances
- Compute L&U credible limits
- Compute (conservative) credible interval having a fixed (e.g. 0.95) lower probability

□ Optimization problems:

minimizing and maximizing

L&U probabilities of an event

□ Prior L&U probabilities

Consider an event B relative to f' , and $P_{st}(B)$ the prior probability obtained from the distribution $DiMn(n'; st)$ in \mathcal{M}_0 .

Prior uncertainty about B is expressed by

$$\underline{P}(B) \text{ and } \overline{P}(B),$$

obtained by min-/maximization of $P_{st}(B)$ w.r.t. $t \in \mathcal{S}^*(1, K)$.

□ Posterior L&U probabilities

Denote $P_{st}(B|\mathbf{a})$ the posterior probability of B obtained from the prior $DiMn(n'; st)$ in \mathcal{M}_0 , i.e. the posterior $DiMn(n'; \mathbf{a} + st)$ in \mathcal{M}_n .

Posterior uncertainty about B is expressed by

$$\underline{P}(B|\mathbf{a}) \text{ and } \overline{P}(B|\mathbf{a}),$$

obtained by min-/maximization of $P_{st}(B|\mathbf{a})$ w.r.t. $t \in \mathcal{S}^*(1, K)$.

Posterior inferences about $\lambda = g(f')$

- **Derived parameter of interest** (real-valued)

$$\lambda = g(f') = \begin{cases} f'_k \\ \sum_k y_k f'_k \\ f'_i / f'_j \\ \text{etc.} \end{cases}$$

Inferences about λ can be summarized by

- **L&U expectations**

$$\underline{E}(\lambda|\mathbf{a}) \quad \text{and} \quad \overline{E}(\lambda|\mathbf{a}),$$

obtained by min-/maximization of $E_{st}(\lambda|\mathbf{a})$ *w.r.t.* $t \in \mathcal{S}^*(1, K)$,

- **L&U cumulative distribution functions (cdf)**

$$\underline{F}_\lambda(u|\mathbf{a}) = \underline{P}(\lambda \leq u|\mathbf{a})$$

$$\overline{F}_\lambda(u|\mathbf{a}) = \overline{P}(\lambda \leq u|\mathbf{a})$$

obtained by min-/maximization of $P_{st}(\lambda \leq u|\mathbf{a})$ *w.r.t.* $t \in \mathcal{S}^*(1, K)$,

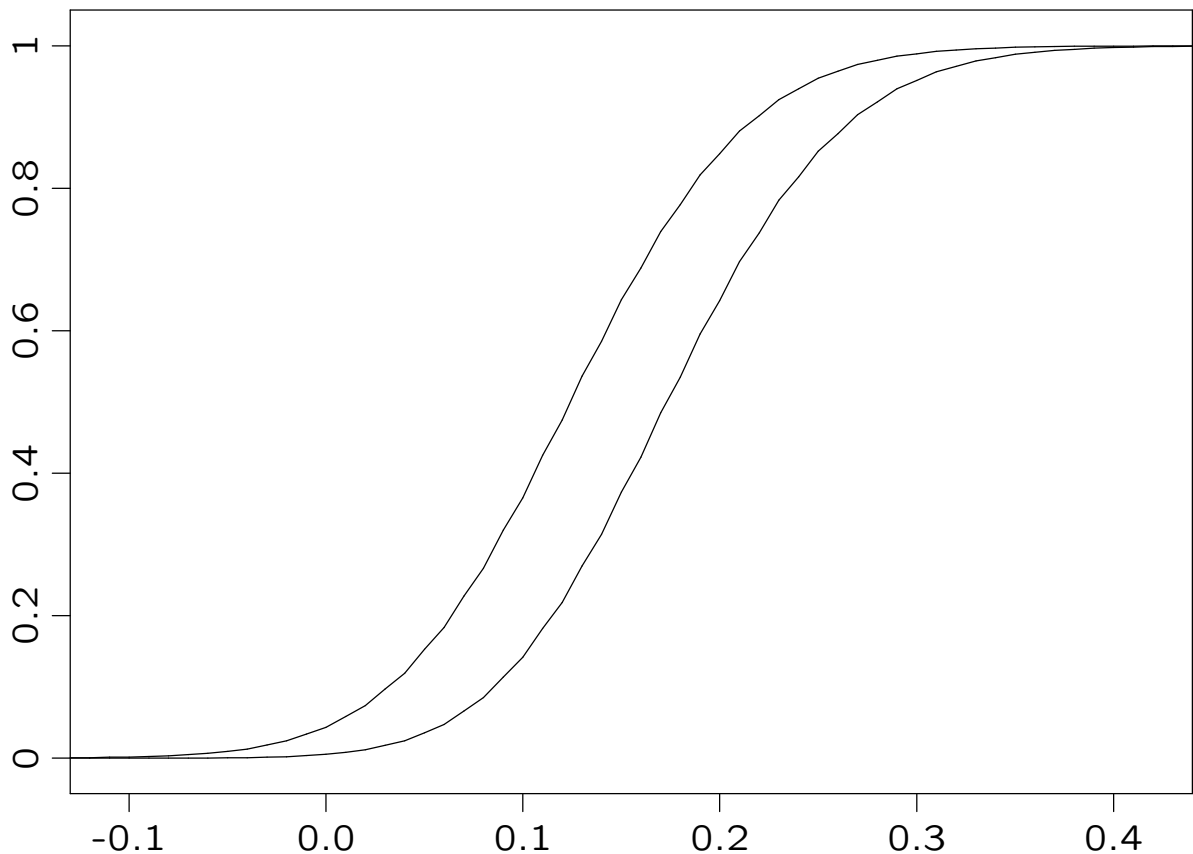
Example of L&U cdf's

□ Example from Walley, Bernard (1999)

Data $a = (2, 12, 46, 6, 0)$ with $n = 66$ and $K = 5$.
Prediction for $n' = 384$ (i.e. $n^* = 450$), on

$$\begin{aligned}\lambda = g(f^*) &= 2f_1^* + f_2^* - f_4^* - 2f_5^* \\ &= \frac{384}{450}g(f') + \frac{66}{450}g(f)\end{aligned}$$

□ L&U cdf's of λ



Optimization problems

□ Set or convex combinations?

The prior & posterior sets, \mathcal{M}_0 and \mathcal{M}_n , of *Diri* or *DiMn* distributions, are used to define lower previsions $\underline{P}(\cdot)$ (by taking lower envelopes). Each $\underline{P}(\cdot)$ is equivalent to the class of its dominating linear previsions, which contains also all convex combinations of these *Diri* or *DiMn* distributions.

□ Optimization of $\mathbf{E}_{st}(\lambda)$ or $\mathbf{E}_{st}(\lambda|\mathbf{a})$

Since $E(\cdot)$ is linear, only requires optimization on the original set of Dirichlet's, \mathcal{M}_0 or \mathcal{M}_n .

□ Optimization of $\mathbf{F}_{st,\lambda}(u)$ or $\mathbf{F}_{st,\lambda}(u|\mathbf{a})$

Similarly, since $F(\cdot)$ is the probability of the event $(\lambda \leq u)$ (*i.e.* the expectation of the corresponding indicator function), optimization only requires the original set \mathcal{M}_0 or \mathcal{M}_n .

□ Optimization attained

- often by corners for $\mathbf{t} \in \mathcal{T}$, *i.e.* when some $t_k \rightarrow 1$, and all others tend to 0,
- but, not always

Inferences about θ_k from the IDM

□ **Prior L&U expectations and cdf's**

Expectations

$$\underline{E}(\theta_k) = 0 \quad \text{and} \quad \overline{E}(\theta_k) = 1$$

Cdf's

$$\underline{P}(\theta_k \leq u) = P(\text{Beta}(s, 0) \leq u)$$

$$\overline{P}(\theta_k \leq u) = P(\text{Beta}(0, s) \leq u)$$

□ **Posterior L&U expectations and cdf's**

Expectations

$$\underline{E}(\theta_k | \mathbf{a}) = \frac{a_k}{n + s} \quad \text{and} \quad \overline{E}(\theta_k | \mathbf{a}) = \frac{a_k + s}{n + s}$$

Cdf's

$$\underline{P}(\theta_k \leq u | \mathbf{a}) = P(\text{Beta}(a_k + s, n - a_k) \leq u)$$

$$\overline{P}(\theta_k \leq u | \mathbf{a}) = P(\text{Beta}(a_k, n - a_k + s) \leq u)$$

□ **Optimization** attained for $t_k \rightarrow 0$ or $t_k \rightarrow 1$.

Equivalent to:

Haldane + s extreme observations.

Extreme ID(M)M's (1)

□ Ignorance vs. Near-ignorance

- Ignorance in the IP theory: vacuous probabilistic statements
- Complete ignorance: ignorance about **all** gambles and events
- Near-ignorance: ignorance about **some** gambles and/or events

□ Two extremes

- $s \rightarrow 0$: Haldane's model, precise
- $s \rightarrow \infty$: vacuous model, maximally imprecise

□ Haldane's model: $s \rightarrow 0$

- Unreasonable account of prior uncertainty
- Inferences over-confident with extreme data
- **You learn too quickly!**

Extreme ID(M)M's (2)

□ **Vacuous model:** $s \rightarrow \infty$

- The $IDM(s_{sup})$ contains all IDM's with $s \leq s_{sup}$, i.e. all $Diri_{st}$, $s \leq s_{sup}$, $t \in \mathcal{T}$. At the limit, the $IDM(s_{sup} \rightarrow \infty)$ contains all Dirichlet's
- Hence, the $IDM(s_{sup} \rightarrow \infty)$ contains all mixtures (convex combinations) of Dirichlet's
- But, any distribution on Θ can be approximated by a finite convex mixture of Dirichlet's. So, the $IDM(s_{sup} \rightarrow \infty)$, contains **all** distributions on Θ
- Leads to **vacuous statements for any gamble**, and for both **prior and posterior** inferences
- **You never learn anything!**

□ **Conclusions**

- $s \rightarrow 0$: **Too precise!**
- $s \rightarrow \infty$: **Too imprecise!**

Hyperparameter s

□ Interpretations of s

- Determines the degree of imprecision in *posterior* inferences; the larger s , the more cautious inferences are
- s as a number of additional **unknown** observations

□ Hyperparameter s must be small

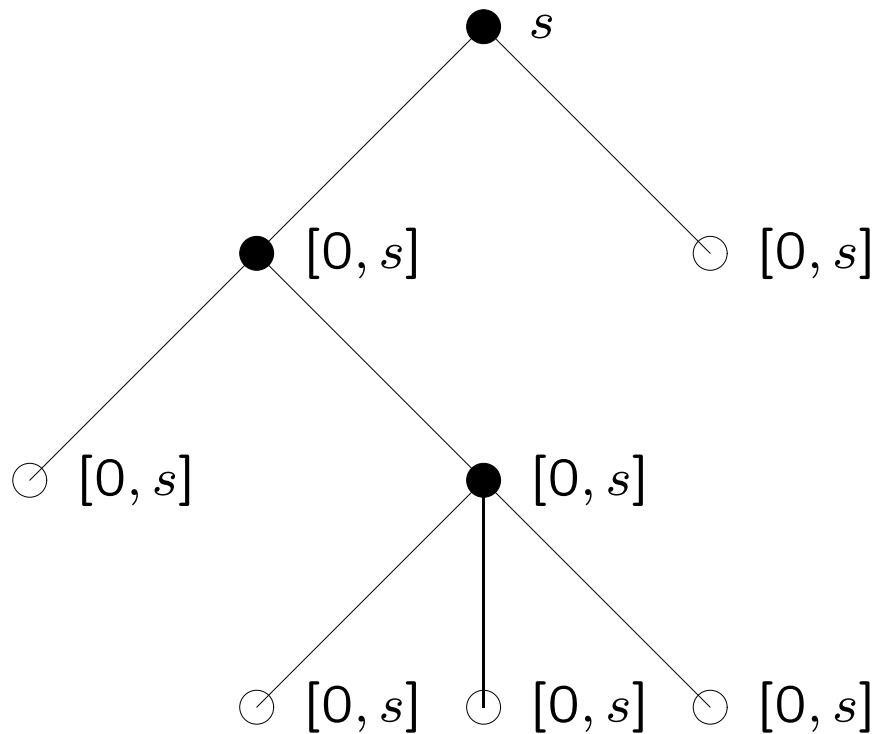
- If too high, inferences are too weak

□ Hyperparameter s must be large enough to

- Encompass objective Bayesian inferences:
Haldane: $s > 0$; Perks: $s \geq 1$
Other solutions? Problem: $s \geq K/2$ or $\geq K$
- Encompass frequentist inferences

□ Suggested values: $s = 1$ or $s = 2$

Why does the ID(M)M satisfy the EP and RIP?



- *Diri* or *DiMn* distributions compatible with any tree. But, under a PDM or PDMM, total prior strength s scatters when moving down the tree
- In the IDM or IDMM, all allocations of s to the nodes are possible (due to imprecision)
- Each sub-tree **inherits** the same s

**PREDICTIVE INFERENCE
FROM THE ID(M)M**

Bayesian inference (recall)

□ Apply Bayes' theorem once

$$\left\{ \begin{array}{c} \text{Prior } P(\theta) \\ + \\ \text{Sampling } P(a|\theta) \end{array} \right. \longrightarrow \left\{ \begin{array}{c} \text{Posterior } P(\theta|a) \\ + \\ \text{Prior predictive } P(a) \end{array} \right.$$

□ Apply Bayes' theorem a second time

$$\left\{ \begin{array}{c} \text{Prior}' } P(\theta|a) \\ + \\ \text{Sampl.}' } P(a'|\theta, a) \end{array} \right. \longrightarrow \left\{ \begin{array}{c} \text{Posterior}' } P(\theta|a', a) \\ + \\ \text{Post. pred. } P(a'|a) \end{array} \right.$$

□ Learning model about

- unknown chances: $P(\theta)$ updated to $P(\theta|a)$
- future data: $P(a)$ updated to $P(a'|a)$

Bayesian prediction from a single $\text{Diri}(\alpha)$ prior

□ Dirichlet-multinomial prior

$$a \sim \text{DiMn}(n; \alpha)$$

$$\begin{aligned} P(a) &= \prod_k \binom{a_k + \alpha_k - 1}{a_k} / \binom{n + s - 1}{n} \\ &= \binom{n}{a} \frac{\alpha_1^{[a_1]} \dots \alpha_K^{[a_K]}}{s^{[n]}} \end{aligned}$$

□ Dirichlet-multinomial posterior

$$a'|a \sim \text{DiMn}(n'; a + \alpha)$$

$$\begin{aligned} P(a'|a) &= \prod_k \binom{a'_k + a_k + \alpha_k - 1}{a'_k} / \binom{n' + n + s - 1}{n'} \\ &= \binom{n'}{a'} \frac{(a_1 + \alpha_1)^{[a'_1]} \dots (a_K + \alpha_K)^{[a'_K]}}{(n + s)^{[n']}} \end{aligned}$$

Beta-binomial marginals under a single Dir(α) prior

□ **Beta-binomial marginal prior** for a_k

$$a_k \sim \text{BeBi}(n; \alpha_k, s - \alpha_k)$$

$$\begin{aligned} P(a_k) &= \frac{\binom{a_k + \alpha_k - 1}{a_k} \binom{n - a_k + s - \alpha_k - 1}{n - a_k}}{\binom{n + s - 1}{n}} \\ &= \binom{n}{a_k} \frac{\alpha_k^{[a_k]} (s - \alpha_k)^{[n - a_k]}}{s^{[n]}} \end{aligned}$$

□ **Beta-binomial marginal posterior** for a'_k

$$a'_k | \mathbf{a} \sim \text{BeBi}(n'; a_k + \alpha_k, n - a_k + s - \alpha_k)$$

$$\begin{aligned} P(a'_k | \mathbf{a}) &= \frac{\binom{a'_k + a_k + \alpha_k - 1}{a'_k} \binom{n' - a'_k + n - a_k + s - \alpha_k - 1}{n' - a'_k}}{\binom{n' + n + s - 1}{n'}} \\ &= \binom{n'}{a'_k} \frac{(a_k + \alpha_k)^{[a'_k]} (n - a_k + s - \alpha_k)^{[n' - a'_k]}}{(n + s)^{[n']}} \end{aligned}$$

Prior predictive distribution under the IDMM

□ **Prior prediction** about a and $f = a/n$

Prior uncertainty about a is described by a set of *DiMn* distributions:

$$\mathcal{M}_0 = \{DiMn(n; st) : t \in \mathcal{S}^*(1, K)\}$$

□ **Vacuous L&U prior expectations** of a_k and f_k

$$\begin{array}{ll} \underline{E}(a_k) = 0 & \overline{E}(a_k) = n \\ \underline{E}(f_k) = 0 & \overline{E}(f_k) = 1 \end{array}$$

obtained as $t_k \rightarrow 0$ and $t_k \rightarrow 1$ respectively

□ **Vacuous L&U prior cdf's** of a_k

(Notation: $F_k(u) = P(a_k \leq u)$, for $u = 0, \dots, n$)

$$\begin{array}{ll} \underline{F}_k(u) = 0 & \text{if } 0 \leq u < n \\ \overline{F}_k(u) = 1 & \text{if } 0 \leq u \leq n \end{array}$$

obtained as $t_k \rightarrow 1$ and $t_k \rightarrow 0$ respectively

Posterior predictive distribution under the IDMM (1)

□ **Posterior prediction** about $a'|a$ and $f'|a$

Posterior uncertainty about a' , conditional on a , is described by the corresponding set of updated *DiMn* distributions:

$$\mathcal{M}_n = \{DiMn(n'; a + st) : t \in \mathcal{S}^*(1, K)\}$$

□ **L&U posterior expectations** of a'_k and f'_k

$$\underline{E}(a'_k|a) = n' \frac{a_k}{n + s} \quad \bar{E}(a'_k|a) = n' \frac{a_k + s}{n + s}$$

$$\underline{E}(f'_k|a) = \frac{a_k}{n + s} \quad \bar{E}(f'_k|a) = \frac{a_k + s}{n + s}$$

obtained as $t_k \rightarrow 0$ and $t_k \rightarrow 1$ respectively

Posterior predictive distribution under the IDMM (2)

- **L&U posterior cdf's** of a'_k
 (Notation: $F_k(u|\mathbf{a}) = P(a'_k \leq u|\mathbf{a})$, for $u = 0, \dots, n'$)

$$\underline{F}_k(u|\mathbf{a}) = \sum_{a'_k=0}^u \frac{\binom{a'_k + a_k + s - 1}{a'_k} \binom{n' - a'_k + n - a_k - 1}{n' - a'_k}}{\binom{n' + n + s - 1}{n'}}$$

$$\overline{F}_k(u|\mathbf{a}) = \sum_{a'_k=0}^u \frac{\binom{a'_k + a_k - 1}{a'_k} \binom{n' - a'_k + n - a_k + s - 1}{n' - a'_k}}{\binom{n' + n + s - 1}{n'}}$$

obtained as $t_k \rightarrow 1$ and $t_k \rightarrow 0$ respectively

- **L&U posterior exp. & cdf's** are obtained using either

$$\text{BeBi}(n'; a_k, n - a_k + s)$$

or $\text{BeBi}(n'; a_k + s, n - a_k)$

Pooling categories

□ **Pooling** categories c_k and c_l into c_j

$$\begin{aligned}a_j &= a_k + a_l \\a'_j &= a'_k + a'_l \\ \alpha_j &= \alpha_k + \alpha_l\end{aligned}$$

□ **Then**

- Each $DiMn_K$, prior or posterior, is transformed into a $DiMn_{K-1}$ where c_j replaces c_k and c_l , with all absolute strengths obtained by summation.
- Recursively, for any pooling in $J < K$ categories, the **DiMn form** and the **value of s** are both preserved.

□ **Thus, in the IDMM,**

L&U prior and posterior probabilities for any event involving pooled counts with $J < K$ categories are invariant whether we

- Pool first, then apply IDMM(s)
- Apply IDMM(s) first, then pool

Properties & principles

□ Prior ignorance about C and K

- Symmetry in the K categories
- Embedding principle (EP) satisfied, due to the pooling property

□ Prior near-ignorance about a & f

- Near-ignorance properties: L&U exp. $E(a_k)$, $E(f_k)$ and cdf's $F_{a_k}(\cdot)$, $F_{f_k}(\cdot)$ are vacuous
- Many other events, or derived parameters, have vacuous prior probabilities, or previsions
- But not all, unless $s \rightarrow \infty$

□ Posterior inferences

- Satisfy coherence (CP)
- Satisfy the likelihood principle (LP)
- Representation invariance (RIP) is satisfied, for the same reason as EP is

Frequentist prediction

□ “Bayesian and confidence limits for predictions” (Thatcher, 1964)

- Considers binomial or hypergeometric data ($K = 2$), $\mathbf{a} = (a_1, n - a_1)$.
- Studies the prediction about n' future observations, $\mathbf{a}' = (a'_1, n' - a'_1)$.
- Derives lower and upper **confidence limits** (frequentist) for a'_1 .
- Compares these confidence limits to **credibility limits** (Bayesian) from a Beta prior.

□ Main result

- Upper confidence and credibility limits for a'_1 coincide *iff* the prior is $Beta(\alpha_1 = 1, \alpha_2 = 0)$.
- Lower confidence and credibility limits for a'_1 coincide *iff* the prior is $Beta(\alpha_1 = 0, \alpha_2 = 1)$.

□ IDM with $s = 1$!

These two *Beta* priors are the most extreme priors under the IDM with $s = 1$

Towards the IDMM? (Thatcher, 1964)

□ A “difficulty”

“... is there a prior distribution such that both the upper and lower Bayesian limits always coincide with confidence limits? ... In fact there are not such distributions.” (Thatcher, 1964, p. 184)

□ Reconciling frequentist and Bayesian

“... we shall consider whether these difficulties can be overcome by a more general approach to the prediction problem: in fact, by ceasing to restrict ourselves to a single set of confidence limits or a single prior distribution.” (Thatcher, 1964, p. 187)

THE RULE OF SUCCESSION

Rule of succession problem

□ **Problem** $P(a'|a)$ for $n' = 1$

- Prediction about the next observation
- Also called **immediate prediction**

□ **A solution to it**

- Called a **rule of succession**
- So many rules for such an (apparently) simple problem!

□ **Highly debated problem**

- Very early problem in Statistics
- **Laplace** computing the probability that the sun will rise tomorrow

□ **Two types of problems / solutions**

- Prior rule, before observing any data
- Posterior rule, after observing some data

The “Bag of marbles” example

□ “Bag of marbles” problems (Walley, 1996)

- “I have ... a closed bag of coloured marbles. I intend to shake the bag, to reach into it and to draw out one marble. What is the probability that I will draw a red marble?”
- “Suppose that we draw a sequence of marbles whose colours are (in order):

blue, green, blue, blue, green, red.

What conclusions can you reach about the probability of drawing a red marble on a future trial?”

□ Two problems of predictive inference

- Prior prediction, before observing any item
- Posterior prediction, after observing n items

□ Inference from a state of prior ignorance about the proportions of the various colours

Notation

□ Event, elementary or combined

Let B_j be the event that the next observation is of type c_j , where c_j is a subset of C with J elements

$$1 \leq J \leq K$$

If $J = 1$, then $c_j = c_k$ is an elementary category
If $J > 1$, then c_j is a combined category

□ Define

The observed count and frequency of c_j

$$a_j = \sum_{k \in j} a_k \quad f_j = \sum_{k \in j} f_k$$

The prior strength, and relative strength, of c_j from a $Diri(\alpha)$ prior

$$\alpha_j = \sum_{k \in j} \alpha_k \quad t_j = \sum_{k \in j} t_k$$

Rule of succession under a PDMM

□ Bayesian rule of succession

The rule of succession obtained from a PDMM, with hyper-parameters $\alpha = st$, is

$$\begin{aligned} P(B_j|\mathbf{a}) &= \frac{a_j + \alpha_j}{n + s} \\ &= \frac{nf_j + st_j}{n + s} \end{aligned}$$

The prior prediction, obtained for $n = a_j = 0$, is

$$P(B_j) = t_j$$

□ Generally

Denoting $f'_j = \sum_{k \in j} f'_k$, the future frequencies in n' data, and possibly $\theta_j = \sum_{k \in j} \theta_k$, the population frequencies, then

$$\begin{aligned} P(B_j) &= E(f'_j) = E(\theta_j) \\ P(B_j|\mathbf{a}) &= E(f'_j|\mathbf{a}) = E(\theta_j|\mathbf{a}) \end{aligned}$$

Prior rule of succession under the IDMM

□ Prior rule of succession

The L&U prior probabilities of B_j are vacuous:

$$\underline{P}(B_j) = 0 \quad \text{and} \quad \overline{P}(B_j) = 1,$$

obtained as $t_j \rightarrow 0$ and $t_j \rightarrow 1$ respectively

□ Prior ignorance

Prior imprecision is maximal, L&U probabilities are vacuous:

$$\Delta(B_j) = \overline{P}(B_j) - \underline{P}(B_j) = 1$$

irrespectively of s

Posterior rule of succession under the IDMM

□ Posterior rule of succession

After data \mathbf{a} have been observed, the posterior L&U probabilities of event B_j are

$$\underline{P}(B_j|\mathbf{a}) = \frac{a_j}{n + s} \quad \text{and} \quad \overline{P}(B_j|\mathbf{a}) = \frac{a_j + s}{n + s},$$

obtained as $t_j \rightarrow 0$ and $t_j \rightarrow 1$ respectively

□ Posterior imprecision

$$\Delta(B_j|\mathbf{a}) = \overline{P}(B_j|\mathbf{a}) - \underline{P}(B_j|\mathbf{a}) = \frac{s}{n + s}$$

□ L&U probabilities and f_j

The interval always contains $f_j = a_j/n$. The L&U probabilities both converge to f_j as n increases.

□ Rule independent from C , K and J

Rule of succession and imprecision

□ Degree of imprecision about B_j

- Prior state: imprecision is maximal

$$\Delta(B_j) = 1$$

- Posterior state:

$$\Delta(B_j|a) = \frac{s}{n + s}$$

□ Interpretation of s

Hyper-parameter s controls how fast imprecision diminishes with n : s is the number of observations necessary to halve imprecision about B_j .

Objective Bayesian models

□ Bayesian rule of succession

The rule of succession obtained from a single symmetric $DiMn$ distribution, $DiMn(n'; \alpha)$ with $n' = 1$ and $\alpha_k = s/K$, is

$$P(B_j|\mathbf{a}) = \frac{a_j + \alpha_j}{n + s} = \frac{nf_j + s\frac{J}{K}}{n + s}$$

□ Objective Bayesian rules: $P(B_j|\mathbf{a}) =$

Haldane	a_j/n
Perks	$(a_j + J/K)/(n + 1)$
Jeffreys	$(a_j + J/2)/(n + K/2)$
Bayes	$(a_j + J)/(n + K)$

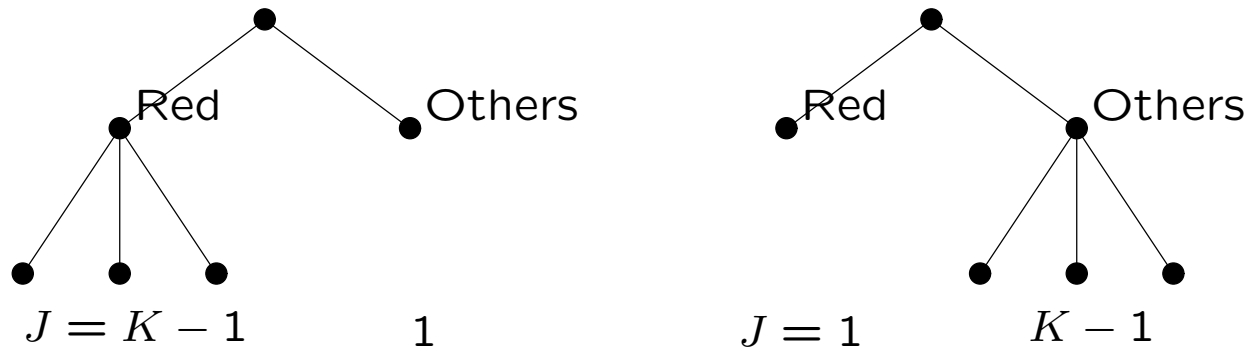
□ Dependence on K and J except Haldane

□ Particular case $J = 1, K = 2$

If $a_j = n/2$, i.e. $f_j = 1/2$, each Bayesian rule leads to $P(B_j|\mathbf{a}) = 1/2$, whether $n = 0$, or $n = 10$, 100 or 1000.

Categorization arbitrariness

□ Arbitrariness of C , i.e. of J and K



Most extremes cases obtained as $K \rightarrow \infty$

□ Bayesian rules

Yield **intervals** when arbitrariness is introduced

Bayes-Laplace	$[0; 1]$,	IDM($s \rightarrow \infty$)
Jeffreys	$[0; 1]$,	IDM($s \rightarrow \infty$)
Perks	$[\frac{a_j}{n+1}; \frac{a_j+1}{n+1}]$,	IDM($s = 1$)
Haldane	$[\frac{a_j}{n}; \frac{a_j}{n}]$,	IDM($s \rightarrow 0$)

Frequentist rule of succession

□ “Bayesian and confidence limits for prediction” (Thatcher, 1964)

- Studies the particular case of immediate prediction

□ Main result (reminder)

- Upper confidence and credibility limits for a'_1 coincide *iff* the prior is $Beta(\alpha_1 = 1, \alpha_2 = 0)$.
- Lower confidence and credibility limits for a'_1 coincide *iff* the prior is $Beta(\alpha_1 = 0, \alpha_2 = 1)$.

□ Frequentist “rule of succession”

When reinterpreted as Bayesian rules of succession, the lower and upper confidence limits respectively correspond to:

$$P(B_j|\mathbf{a}) = \frac{a_j}{n+1} \quad \text{and} \quad P(B_j|\mathbf{a}) = \frac{a_j + 1}{n+1}$$

i.e. to the IDM interval for $s = 1$.

CONCLUSIONS

Comments on predictive inference

□ **Predictive approach is more fundamental**
(see, Geisser, 1993)

- Finite population & data
- Models observables only, not hypothetical parameters
- Relies on the exchangeability assumption only.
- Pearson (1920) considered predictive inference as “the fundamental problem of practical statistics”

□ **Predictive approach is more natural,**

□ **For the IDMM, in particular**

- Gives the IDM as a limiting case as $n' \rightarrow \infty$
- Covers sampling with replacement from a finite population

Why using a set of Dirichlet's Walley (1996, p. 7)

□ About Dirichlet's

- (a) Dirichlet prior distributions are **mathematically tractable** because ... they generate Dirichlet posterior distributions;
- (b) when categories are combined, Dirichlet distributions **transform to other Dirichlet** distributions (this is the crucial property which ensures that the **RIP** is satisfied);
- (c) sets of Dirichlet distributions are **very rich**, because they produce the same inferences as their convex hull and any prior distribution can be approximated by a finite mixture of Dirichlet distributions;
- (d) the most common **Bayesian models** for prior ignorance about θ are Dirichlet distributions.

□ **Same arguments hold** for $DiMn$ distributions

Links between IDM and IDMM

□ Parametric and predictive inference

In general, in both precise Bayesian models and in the ID(M)M,

- $\theta, \theta|a$ yield $f, f'|a$ (from Bayes' theorem)
- $f, f'|a$ yield $\theta, \theta|a$ (as $n' \rightarrow \infty$)

□ Equivalence between IDM and IDMM

- The IDM and the IDMM are equivalent, if we assume that n' can tend to infinity
- Any IDMM statement about f' which is independent of n' is also a valid IDM statement about θ

□ Two views of the IDMM

- The IDMM is the predictive side of the IDM
- The IDMM is a model of its own

Fundamental properties of the ID(M)M

□ Principles

Satisfies several desirable principles for prior ignorance: SP, EP, RIP, LP, SRP, coherence.

□ ID(M)M vs. Bayesian and frequentist

- Answers several difficulties of alternative approaches
- Provides means to reconcile frequentist and objective Bayesian approaches (Walley, 2002)

□ Generality

More general than for multinomial data. Valid under a general hypothesis of exchangeability between observed and future data. (Walley, Bernard, 1999).

□ Degree of imprecision and n

Degree of imprecision in posterior inferences enables one to distinguish between: (a) prior uncertainty still dominates, (b) there is substantial information in the data.

The two cases can occur within the **same** data set.

REFERENCES

- Berger, J. O., Bernardo, J. M. (1992), "Ordered Group Reference Priors with Application to the Multinomial Problem", *Biometrika*, 79, no. 1, pp. 2–5–37.
- Bernard, J.-M. (1996), "Bayesian Interpretation of Frequentist Procedures for a Bernoulli Process", *The American Statistician*, 50, no. 1, 7–13.
- Bernard, J.-M. (1997), "Bayesian Analysis of Tree-Structured Categorized Data", *Revue Internationale de Systémique*, Vol. 11 no. 1, pp. 1–1–29.
- Fang, K. T., Kotz, S., Ng, K. W. (1990), *Symmetric Multivariate and Related Distributions*, New-York: Chapman and Hall.
- Geisser, S. (1993), *Predictive Inference: An Introduction*, Monographs on Statistics and Applied Probability 55, New-York: Chapman & Hall.
- Haldane, J. B. S. (1948), "The Precision of Observed Values of Small Frequencies," *Biometrika*, 35, pp. 2–97–300.
- Hutter, M. (2003). Robust Estimators under the Imprecise Dirichlet Model. In Bernard J.-M., Seidenfeld T., Zaffalon M. (Eds), *ISIPTA'03: Proceedings of the Third International Symposium on Imprecise Probabilities and Their Applications, Lugano, Switzerland*, Proceedings in Informatics 18, Waterloo, Ontario, Canada: Carleton Scientific, pp. 274–289.
- Jeffreys, H. (1946), "An Invariant Form for the Prior Probability in Estimation Problems", *Proceedings of the Royal Society of London, Ser. A*, 186, pp. 4–53–461.
- Jeffreys, H. (1961), *Theory of Probability*, 3rd ed., Oxford: Clarendon Press.

- Kass, R. E., Wasserman, L. (1996), "The Selection of Prior Distributions by Formal Rules", *Journal of the American Statistical Association*, Vol. 91, no. 435, pp. 1–343–1370.
- Perks, F. J. A. (1947), "Some Observations on Inverse Probability Including a New Indifference Rule (with discussion)", *Journal of the Institute of Actuaries*, 73, pp. 2–85–334.
- Thatcher, A. R. (1964), "Relationships Between Bayesian and Confidence Limits for Predictions" (with discussion), *Journal of the Royal Statistical Society*, Ser. B, 26, pp. 1–76–210.
- Walley, P. (1991), *Statistical Reasoning with Imprecise Probabilities*, Monographs on Statistics and Applied Probability 42, London: Chapman & Hall.
- Walley, P. (1996), "Inferences from Multinomial Data: Learning about a Bag of Marbles", *Journal of the Royal Statistical Society*, Ser. B., 58 no. 1, pp. 3–57.
- Walley, P. (2002). Reconciling frequentist properties with the likelihood principle. *Journal of Statistical Planning and Inference*, Vol. 105, no. 1, pp. 35–65.
- Walley, P., Bernard, J.-M. (1999), "Imprecise Probabilistic Prediction for Categorical Data", Technical Report CAF-9901, Laboratoire Cognition et Activités Finalisées, University Paris 8, Saint-Denis, France, January 1999.