

Bayesian Model Specification

1: Foundations

David Draper

*Department of Applied Mathematics and Statistics
University of California, Santa Cruz*

and (1 Jul–31 Dec 2013) *eBay Research Labs*

{draper@ams.ucsc.edu, dadraper@ebay.com}
www.ams.ucsc.edu/~draper

SHORT COURSE (DAY 3)
UNIVERSITY OF READING (UK)

29 Nov 2013

© 2013 David Draper (all rights reserved)

An Example, to Fix Ideas

Case Study 1. (Krnjajić, Kottas, Draper [KKD] 2008): *In-home geriatric assessment (IHGA)*. In an **experiment** conducted in the **1980s** (Hendriksen et al. 1984), **572 elderly people, representative of \mathcal{P}** = {all **non-institutionalized elderly people in Denmark**}, were **randomized, 287** to a **control (C)** group (who received **standard health care**) and **285** to a **treatment (T)** group (who received **standard care plus IHGA**: a kind of **preventive medicine** in which each person's **medical and social needs** were assessed and acted upon **individually**).

One **important outcome** was the **number of hospitalizations** during the **two-year** life of the study:

Group	Number of Hospitalizations				n	Mean	SD
	0	1	...	k			
Control	n_{C0}	n_{C1}	...	n_{Ck}	$n_C = 287$	\bar{y}_C	s_C
Treatment	n_{T0}	n_{T1}	...	n_{Tk}	$n_T = 285$	\bar{y}_T	s_T

Let μ_C and μ_T be the **mean hospitalization rates** (per two years) in \mathcal{P} under the C and T **conditions**, respectively.

Here are **four statistical questions** that **arose from this study**:

The Four Principal Statistical Activities

Q₁: Was the **mean number of hospitalizations per two years** in the IHGA group **different from** that in **control** by an **amount** that was **large** in **practical** terms? [**description** involving $\left(\frac{\bar{y}_T - \bar{y}_C}{\bar{y}_C}\right)$]

Q₂: Did IHGA (**causally**) **change** the **mean number of hospitalizations per two years** by an **amount** that was **large** in **statistical** terms? [**inference** about $\left(\frac{\mu_T - \mu_C}{\mu_C}\right)$]

Q₃: On the **basis of this study**, how **accurately** can You **predict** the **total decrease in hospitalizations** over a period of N years if **IHGA** were **implemented throughout Denmark**? [**prediction**]

Q₄: On the **basis of this study**, is the **decision to implement IHGA** throughout Denmark **optimal** from a **cost-benefit** point of view? [**decision-making**]

These questions **encompass** almost all of the **discipline of statistics**: **describing a data set D** , **generalizing outward inferentially from D** , **predicting new data D^*** , and helping people **make decisions** in the **presence of uncertainty** (I include **sampling/experimental design** under **decision-making**; **omitted**: data **quality assurance (QA)**, ...).

An Axiomatization of Statistics

- 1 (definition) **Statistics** is the study of **uncertainty**: how to **measure it well**, and how to **make good choices** in the face of it.
- 2 (definition) **Uncertainty** is a state of **incomplete information** about something of interest to **You** (Good, 1950: a **generic person** wishing to **reason sensibly** in the presence of **uncertainty**).
- 3 (axiom) (**Your uncertainty** about) “**Something of interest to You**” can always be **expressed** in terms of **propositions**: **true/false** statements A, B, \dots

Examples: You may be **uncertain** about the **truth status** of

- $A =$ (**Hillary Clinton** will be **elected U.S. President** in **2016**), or
- $B =$ (the **in-hospital mortality rate** for patients at **hospital H** admitted in **calendar 2010** with a principal diagnosis of **heart attack** was **between 5% and 25%**).

- 4 (implication) It follows from 1–3 that **statistics** concerns **Your information** (**NOT Your beliefs**) about A, B, \dots

Axiomatization (continued)

5 (axiom) But **Your information** cannot be **assessed** in a **vacuum**: all such **assessments** must be made **relative to (conditional on)** Your **background assumptions** and **judgments** about **how the world works** vis à vis A, B, \dots .

6 (axiom) These **assumptions** and **judgments**, which are themselves a form of **information**, can always be **expressed** in a **set \mathcal{B}** of **background propositions**, all of which **You believe** to be **true**.

Examples of \mathcal{B} :

- In the **IHGA study**, based on the **experimental design**, \mathcal{B} would include the **propositions**

(**Subjects were representative of [like a random sample from] \mathcal{P}**),

(**Subjects were randomized** into one of two groups, **treatment (standard care + IHGA)** or **control (standard care)**).

7 (definition) Call the **“something of interest to You”** θ ; in **applications** θ is often a **vector** (or **matrix**, or **array**) of **real numbers**, but **in principle** it could be **almost anything** (a **function**,

Axiomatization (continued)

an **image** of the **surface of Mars**, a **phylogenetic tree**, ...).

IHGA example: $\theta =$ mean relative decrease $\left(\frac{\mu_T - \mu_C}{\mu_C}\right)$ in hospitalization rate in \mathcal{P} .

8 (axiom) There will typically be an **information source (data set)** D that You judge to be **relevant** to **decreasing** Your uncertainty about θ ; in **applications** D is often again a **vector** (or **matrix**, or **array**) of **real numbers**, but in **principle** it too could be **almost anything** (a **movie**, the **words** in a **book**, ...).

9 (implication) The **presence** of D creates a **dichotomy**:

- **Your information** about θ **{internal, external}** to D .

(People often talk about a **different dichotomy**: **Your information** about θ **{before, after}** D arrives (**prior, posterior**), but **temporal considerations** are actually **irrelevant**.)

10 (implication) It follows from **1-9** that **statistics** concerns itself principally with **five things** (omitted: **description, data QA**, ...):

- (1) **Quantifying Your information** about θ **internal** to D (given \mathcal{B}), and doing so **well** (this term is **not yet defined**);

Foundational Question

(2) **Quantifying Your information** about θ **external** to D (given \mathcal{B}),
and doing so **well**;

(3) **Combining** these two **information sources** (and doing so **well**) to
create a **summary** of **Your uncertainty** about θ (given \mathcal{B}) that includes
all available information You judge to be **relevant** (this is **inference**);

and using **all Your information** about θ (given \mathcal{B}) to make

(4) **Predictions** about **future** data values D^* and

(5) **Decisions** about how to **act sensibly**, even though **Your
information** about θ may be **incomplete**.

Foundational question: How should these tasks be **accomplished**?

This question has **two parts**: **probability** and **statistics**.

The **probability foundations** (addressed here **first**), have an **interesting
and unfortunate history**, in which **much** of the **20th century** will (in
my view) be **seen** in the **21st century** to have been a **series** of **missed
scientific opportunities**.

Theory of Probability: Kolmogorov

From the **1650s (Fermat, Pascal)** through the **18th century (Bayes, Laplace)** to the period **1860–1930 (Venn, Boole, von Mises)**, **three different approaches** for how to think about **uncertainty quantification** — **classical, Bayesian**, and **frequentist probability** — were put forward in an **intuitive** way, but no one ever tried to prove a **theorem** of the form **{given these premises, there's only one sensible way to quantify uncertainty}** until **Kolmogorov, de Finetti, and RT Cox**.

— **Kolmogorov (1933)**: following (and **rigorizing**) **Venn, Boole** and **von Mises**, **probability** is a **function** on (possibly **some of**) the **subsets** of a **sample space Ω** of **uncertain possibilities**, **constrained** to obey some **reasonable axioms**; this is **excellent, as far as it goes**, but **many types of uncertainty cannot (uniquely, comfortably) be fit into this framework** (examples follow).

Kolmogorov was trying to **make precise** the **intuitive notion** of **repeatedly choosing a point at random** in a **Venn diagram** and asking **how frequently** the point falls **inside a specified set**, i.e., his **concept of probability** had a **repeated-sampling, frequentist** character:

Frequentist Probability: Kolmogorov

“The basis for the applicability of the results of the mathematical theory of probability to real ‘random phenomena’ must depend on some form of the frequency concept of probability, the unavoidable nature of which has been established by von Mises in a spirited manner.”

* **Example:** You’re about to roll a **pair of dice** and **You regard** this dice-rolling as **fair**, by which You mean that **(in Your judgment)** all $6^2 = 36$ **elemental outcomes** in $\Omega = \{(1, 1), (1, 2), \dots, (6, 6)\}$ are **equally probable**; then the **Kolmogorov probability of snake eyes** $((1, 1))$ **exists** and is **unique** (from Your **fairness judgment**), namely $\frac{1}{36}$; but

* **Example:** You’re a **doctor**; a **new patient** presents **saying** that he **may** be **HIV positive**; what’s the **Kolmogorov probability** that he is?

What’s Ω ? **This patient** is not the result of a **uniquely-specifiable repeatable “random” process**, he’s just a guy who **walked into Your doctor’s office**, and — throughout the **repetitions** of whatever **repeatable phenomenon** anyone might **imagine** — his **HIV status** is **not fluctuating “randomly”**: he’s either **HIV positive** or he’s **not**.

Theory of Probability: de Finetti

The **closest** You can come to making **Kolmogorov's** approach work here is to **imagine** the set Ω of **all people** **{similar to this patient in all relevant ways}** and ask **how often** You'd get an **HIV-positive person** if You **repeatedly chose** one person **at random** from Ω , but to **make this operational** You have to **specify** what You mean by **"similar to, in all relevant ways,"** and if You **try** to do this You'll notice that it's **not possible** to do so **uniquely** (in such a way that **all other reasonable people** would **unanimously agree** with You).

— **de Finetti** (1937): rigorizing **Bayes**, **probability** is a **quantification** of **betting odds** about the **truth** of a **proposition**, constrained to obey **axioms** guaranteeing **coherence** (absence of **internal contradictions**); this is **more general** than **Kolmogorov** — in fact, it's **as general as You can get**: any **statement** about **sets** can be **expressed** in terms of **propositions** — but **betting odds** are **not fundamental to science**.

de Finetti made **many important contributions** — in particular, his concept of **exchangeability** (see **AMS 206** or **206B**) is **crucial** in **Bayesian modeling** — but (in **my view**) science is about **information**, not **betting**.

— **RT Cox** (1946): following **Laplace**, **probability** is a **quantification of information** about the **truth** of one or more **propositions**, constrained to obey **axioms** guaranteeing **internal logical consistency**; this is both **fundamental to science** and **as general as You can get**.

Cox's goal was to identify what **basic rules** $p(A|B)$ — the **plausibility (weight of evidence)** in favor of (the **truth** of) A given B — should follow so that $p(A|B)$ behaves **sensibly**, where A and B are **propositions** with B **assumed** by You to be **true** and the truth status of A **unknown** to You.

He did this by **identifying** a set of **principles** making **operational** the word **“sensible”** (Jaynes, 2003):

- Suppose You're **willing** to represent **degrees of plausibility** by **real numbers** (i.e., $p(A|B)$ is a function from propositions A and B to \mathbb{R});
 - You insist that **Your reasoning** be **logically consistent**:
 - If a **plausibility assessment** can be arrived at in **more than one way**, then **every possible way** must lead to the **same value**.

Cox's Principles and Axioms

- You always take into account **all of the evidence** You judge to be **relevant** to the **plausibility assessment** under consideration (this is the **Bayesian** version of **objectivity**).
- You always represent **equivalent states of information** by **equivalent plausibility assignments**.

From these **principles** Cox derived a set of **axioms**:

- The **plausibility** of a **proposition** determines the **plausibility** of the proposition's **negation**; each **decreases** as the other **increases**.
 - The **plausibility** of the **conjunction** $AB = (A \text{ and } B)$ of **two propositions** A, B **depends** only on the **plausibility** of B and that of $\{A \text{ given that } B \text{ is true}\}$ (or **equivalently** the **plausibility** of A and that of $\{B \text{ given that } A \text{ is true}\}$).
 - Suppose AB is **equivalent** to CD ; then if You acquire **new information** A and later acquire **further new information** B , and **update** all **plausibilities** each time, the **updated plausibilities** will be the **same** as if You had **first acquired new information** C and **then acquired further new information** D .

Cox's Theorem

From these **axioms** Cox proved a **theorem** showing that **uncertainty quantification** about **propositions** behaves in **one and only one way**:

Theorem: If You accept **Cox's axioms**, then to be **logically consistent** You **must** quantify uncertainty as follows:

- Your **plausibility operator** $pl(A|B)$ — for **propositions** A and B — can be referred to as Your **probability** $P(A|B)$ that A is true, **given** that You regard B as true, and $0 \leq P(A|B) \leq 1$, with **certain truth** of A (given B) represented by **1** and **certain falsehood** by **0**.

- **(normalization)** $P(A|B) + P(\bar{A}|B) = 1$, where $\bar{A} = (\text{not } A)$.

- **(the product rule):**

$$P(AB|C) = P(A|C) \cdot P(B|A C) = P(B|C) \cdot P(A|B C).$$

The **proof** (see, e.g., Jaynes (2003)) involves deriving two **functional equations** $F[F(x, y), z] = F[x, F(y, z)]$ and $x S \left[\frac{S(y)}{x} \right] = y S \left[\frac{S(x)}{y} \right]$ that $pl(A|B)$ must satisfy and then **solving** those equations.

A number of **important corollaries** arise from **Cox's Theorem**:

Optimal Reasoning Under Uncertainty

- (the sum rule):

$$P(A \text{ or } B|C) \equiv P(A + B|C) = P(A|C) + P(B|C) - P(AB|C).$$

- **Extensions** of the **product** and **sum rules** to an **arbitrary finite number** of **propositions** are **easy**, e.g.,

$$P(ABC|D) = P(A|D) \cdot P(B|AD) \cdot P(C|ABD) \text{ and}$$

$$P(A + B + C|D) = P(A|D) + P(B|D) + P(C|D) - P(AB|D) \\ - P(AC|D) - P(BC|D) + P(ABC|D).$$

- This **framework** (obviously) covers **optimal reasoning** about **uncertain quantities** θ taking on a **finite** number of **possible values**; less obviously, it **also handles** (equally well) situations in which the **set** Θ of **possible values** of θ has **infinitely** many elements.

— **Example:** You're studying **quality of care** at the **17 Kaiser Permanente (KP) northern California hospitals** in **2003–7**, before the era of **electronic medical records**; during that time there was a **population** \mathcal{P} of $N = 8,561$ **patients** at these facilities with a **primary admission diagnosis** of **heart attack**.

Inference About a Population Parameter

You take a **simple random sample** of $n = 112$ of these admissions and **record** whether or not each patient had an **unplanned transfer to the intensive care unit (ICU)**, observing $s = 4$ who did; θ is the **proportion** of such **unplanned transfers** in all of \mathcal{P} ; here $\Theta = \{\frac{0}{N}, \frac{1}{N}, \dots, \frac{N}{N}\}$, which can be **conveniently approximated** by $\Theta' = [0, 1]$.

Prior to 2003, the **proportion** of such **unplanned transfers** for **heart attack patients** at **KP** in the **northern California region** was about $q = 0.07$, so **interest** focuses on $P(A|D\mathcal{B})$, where A is the **proposition** ($\theta \leq q$), D is the **proposition** ($s = 4$), and \mathcal{B} includes (among other things) **details** about the **sampling experiment** (e.g., ($n = 112$)).

In this **setup** θ is **usually called** a **(population) parameter**, and is **not itself the result of any sampling experiment** (random or otherwise); for this reason, it's **not possible** to **(directly) quantify uncertainty** about θ from the **Kolmogorov (set-theoretic)** point of view, but it makes **perfect sense** to do so from the **RT Cox (propositional)** point of view.

Optimal Reasoning About a Continuous θ

You could now **more generally** define a function $F_{(\theta|D\mathcal{B})}(q) = P(\theta \leq q|D\mathcal{B})$ and call it the **cumulative distribution function (CDF)** **for (not of)** $(\theta|D\mathcal{B})$, which is **shorthand** for the **CDF** for **Your uncertainty about θ** given D and \mathcal{B} .

If $F_{(\theta|D\mathcal{B})}(q)$ turns out to be **continuous** and **differentiable** in q (I haven't **said yet** how to **calculate F**), it will be **convenient** to write

$$F_{(\theta|D\mathcal{B})}(b) - F_{(\theta|D\mathcal{B})}(a) = P(a < \theta \leq b|D\mathcal{B}) = \int_a^b p_{(\theta|D\mathcal{B})}(q) dq, \quad (1)$$

where the **(partial) derivative** $p_{(\theta|D\mathcal{B})}(q)$ of $F_{(\theta|D\mathcal{B})}$ with respect to q can be called the **density** **for (not of)** **(Your uncertainty about θ)** given D and \mathcal{B} .

In a **small abuse of notation** it's **common** to **write** $F(\theta|D\mathcal{B})$ and $p(\theta|D\mathcal{B})$ instead of $F_{(\theta|D\mathcal{B})}(q)$ and $p_{(\theta|D\mathcal{B})}(q)$ (respectively), letting the **argument θ** of $F(\cdot|D\mathcal{B})$ and $p(\cdot|D\mathcal{B})$ serve as a **reminder** of the **uncertain quantity** in question.

Ontology and Epistemology

NB In the **Kolmogorov approach** a **random variable** X is a **function** from Ω to some **outcome space** O , and if $O = \mathfrak{R}$ You'll often find it **useful to summarize** X 's **behavior** through the **CDF** **of** X :
 $F_X(x) = P(\text{the set of } \omega \in \Omega \text{ such that } X(\omega) \leq x)$, usually written in **propositional-style shorthand** as $F_X(x) = P(X \leq x)$.

In the **RT Cox approach**, there are **no random variables**; there are **uncertain things** θ whose **uncertainty** (when $\Theta = \mathfrak{R}^k$, for integer $1 \leq k < \infty$) can **usefully** be **summarized** with **CDFs** and **densities**.

Jaynes (2003) makes a **worthwhile distinction**: the **statements**

There is noise in the room.

The room is noisy.

seem **quite similar** but are **in fact quite different**: the former is **ontological** (asserting the **physical existence** of something), whereas the latter is **epistemological** (expressing the **personal perception** of the **individual** making the **statement**).

Talking about “the **density** **of** θ ” would be to **confuse ontology** and **epistemology**;

The Mind-Projection Fallacy

Jaynes calls this confusion of **{the world}** (ontology) with **{Your uncertainty about the world}** (epistemology) the **mind-projection fallacy**, and it's clearly a **mistake worth avoiding**.

Returning to the **corollaries** of **Cox's Theorem**,

- Given the set \mathcal{B} , of **propositions** summarizing Your **background assumptions and judgments** about **how the world works** as far as θ , D and future data D^* are **concerned**:

(a) It's **natural** (and indeed **You must be prepared** in this approach) to specify **two conditional probability distributions**:

— $p(\theta|\mathcal{B})$, to quantify **all information** about θ **external** to D that You judge **relevant**; and

— $p(D|\theta\mathcal{B})$, to quantify Your **predictive uncertainty**, given θ , about the **data set D before it's arrived**.

(b) Given the **distributions** in (a), the distribution $p(\theta|D\mathcal{B})$ quantifies **all relevant information** about θ , both **internal and external** to D , and **must be computed** via **Bayes's Theorem**:

Optimal Inference, Prediction and Decision

$$p(\theta|D\mathcal{B}) = c p(\theta|\mathcal{B}) p(D|\theta\mathcal{B}), \quad \text{(inference)} \quad (2)$$

where $c > 0$ is a **normalizing constant** chosen so that the **left-hand side** of (2) **integrates** (or sums) over Θ to **1**;

(c) Your **predictive distribution** $p(D^*|D\mathcal{B})$ for future data D^* given the **observed data set** D **must be expressible** as follows:

$$p(D^*|D\mathcal{B}) = \int_{\Theta} p(D^*|\theta D\mathcal{B}) p(\theta|D\mathcal{B}) d\theta;$$

often there's **no information** about D^* contained in D if θ is known, in which case this expression **simplifies** to

$$p(D^*|D\mathcal{B}) = \int_{\Theta} p(D^*|\theta\mathcal{B}) p(\theta|D\mathcal{B}) d\theta; \quad \text{(prediction)} \quad (3)$$

(d) to make a sensible **decision** about which **action** a You should take in the face of Your **uncertainty** about θ , You **must be prepared to specify**

(i) the set \mathcal{A} of **feasible actions** among which You're **choosing**, and

(ii) a **utility function** $U(a, \theta)$, taking values on \mathfrak{R} and **quantifying** Your **judgments** about the **rewards** (**monetary** or **otherwise**) that would **ensue** if You chose **action** a and the **unknown** actually **took** the **value** θ ; **without loss of generality** You can take **large values** of $U(a, \theta)$ to be **better than small values**;

then the **optimal decision** is to choose the action a^* that **maximizes** the **expectation** of $U(a, \theta)$ over $p(\theta|D\mathcal{B})$:

$$a^* = \operatorname{argmax}_{a \in \mathcal{A}} E_{(\theta|D\mathcal{B})} U(a, \theta) = \operatorname{argmax}_{a \in \mathcal{A}} \int_{\Theta} U(a, \theta) p(\theta|D\mathcal{B}) d\theta. \quad (4)$$

The equation solving the **inference problem** is **traditionally** attributed to **Bayes (1764)**, although it's just an **application** of the **product rule** (page 13), which was **already in use** by **(James) Bernoulli** and **de Moivre** around **1715**, and **Laplace** made **much better use** of this equation from **1774** to **1827** than Bayes did in **1764**; nevertheless the **Laplace/Cox propositional approach** is typically referred to as **Bayesian reasoning**.

Logical Consistency \rightarrow Bayesian Reasoning Justified

Cox's Theorem is equivalent to the assertion

If You wish to **quantify Your uncertainty** about an **unknown θ** (and make **predictions** and **decisions** in the **presence** of that **uncertainty**) in a **logically internally consistent** manner (as **specified** through **Cox's axioms**), on the basis of **data D** and **background assumptions/judgments \mathcal{B}** , then You can **achieve this goal with Bayesian reasoning**, by **specifying** $p(\theta|\mathcal{B})$, $p(D|\theta\mathcal{B})$, and $\{\mathcal{A}, U(a, \theta)\}$ and **using equations (2–4)**.

This **assertion** has not rendered **Bayesian analyses ubiquitous**, although the **value of Bayesian reasoning** has become **increasingly clear** to an **increasingly large number of people** in the **last 20 years**, now that **advances in computing** have made the **routine use of equations (2–4) feasible**.

Advantages include a **unified probabilistic framework**: e.g., in my earlier **ICU example**, **Kolmogorov's non-Bayesian approach** does not permit **direct probability statements** about a **population parameter**, but **Cox's Theorem permits You** to make such statements (summarizing **all relevant available information**) in a natural way.

The Specification Burden

It's **worth noting**, however, that **there really is a theorem here**, of the form $A \rightarrow B$, from which $\bar{B} \rightarrow \bar{A}$; this **comes close to the assertion**

If You employ **non-Bayesian reasoning** then You're **open to the possibility** of **logical inconsistency**,

and indeed there have been some **embarrassing moments** in **non-Bayesian inference** over the past **100 years** (e.g., **negative estimates** for quantities that are **constrained** to be **non-negative**).

Challenges: These **corollaries** to **Cox's theorem** solve problems (3–5) above (page 7) — they leave **no ambiguity** about how to draw **inferences**, and make **predictions** and **decisions**, in the presence of **uncertainty** — but problems (1) and (2) are still **unaddressed**: to **implement** this **logically-consistent approach** in a given application, You have to **specify**

- $p(\theta|\mathcal{B})$, usually called Your **prior information** about θ (given \mathcal{B} ; this is **better understood** as a **summary of all relevant information** about θ **external** to D , rather than by appeal to any **temporal (before-after) considerations**);

The Specification Burden (continued)

- $p(D|\theta \mathcal{B})$, often referred to as Your **sampling distribution** for D given θ (and \mathcal{B} ; this is **better understood** as Your **conditional predictive distribution** for D given θ , before D has been **observed**, rather than by appeal to **other data sets that might have been observed**); and
 - the **action space** \mathcal{A} and the **utility function** $U(a, \theta)$ for **decision-making purposes**.

The results of **implementing** this approach are

- $p(\theta|D \mathcal{B})$, often referred to as Your **posterior** distribution for θ given D (and \mathcal{B} ; as above, this is **better understood** as the **totality of Your current information** about θ , again without appeal to **temporal considerations**);
- Your **posterior predictive distribution** $p(D^*|D \mathcal{B})$ for future data D^* given the **observed data set** D ; and
 - the **optimal decision** a^* given **all available information** (and \mathcal{B}).

Theory of Applied Statistics

To summarize: **Inference** and **prediction** require You to **specify** $p(\theta|\mathcal{B})$ and $p(D|\theta\mathcal{B})$; **decision-making** requires You to **specify** the same two **ingredients** plus \mathcal{A} and $U(a, \theta)$; **how** should **this** be **done** in a **sensible** way?

Cox's Theorem and its **corollaries** provide **no constraints** on the **specification process**, apart from the requirement that **all probability distributions** be **proper** (**integrate** or **sum** to **1**).

In **my view**, in seeking **answers** to these **specification questions**, as a **profession** we're approximately where the **discipline of statistics** was in arriving at an **optimal theory of probability before Cox's work**: many people have made **ad-hoc suggestions** (**some** of them **good**), but **little formal progress** has been made.

Developing (1) **principles**, (2) **axioms** and (3) **theorems** about **optimal specification** could be regarded as creating a **Theory of Applied Statistics**, which we **need** but **do not yet have**.

Definition. Let's agree to call $\{p(\theta|\mathcal{B}), p(D|\theta\mathcal{B}), \mathcal{A}, U(a, \theta)\}$ Your **model** M for **Your uncertainty about** θ , with the **convention** that

Optimal Model Specification

when **no decision-making** is involved this **simplifies** to
 $\{p(\theta|\mathcal{B}), p(D|\theta\mathcal{B})\}$.

How should M be **specified**? Where is the **progression**

Principles \rightarrow **Axioms** \rightarrow **Theorems**

to **guide You**, the way **Cox's Theorem** settled the **foundational questions** for **probability**?

In my view this is the **central unsolved foundational problem** in **statistical inference, prediction** and **decision-making**.

Optimal model specification. Can M be **specified optimally**?

That **depends** on **what You mean** by **optimal**; here's **my definition**:

Definition. In model specification, **optimal** = {**conditioning** only on **propositions rendered true** by the **context** of the **problem** and the **design** of the **data-gathering process**, while at the **same time ensuring** that the **set of conditioning propositions** includes **all relevant problem context**}.

Repeat **Q:** Can this **optimality goal** be achieved?

A, Part 1: Yes, sometimes.

Focusing for the **moment** on **inference** and **prediction**, in some **special settings**, You really **don't have any model uncertainty at all** — sometimes the **prior distribution** and/or **likelihood/sampling distribution** arise **directly** from **problem context** (**Lecture Notes, Part 2**).

Distribution-free well-calibrated Bayesian inferential methods are **sometimes** also **available** (**Lecture Notes, Part 2A**).

References

- Cox RT (1946). Probability, frequency and reasonable expectation. *American Journal of Physics*, **14**, 1–13.
- De Finetti B (1937). La prévision: ses lois logiques, ses sources subjectives. *Annals of the Institute of Henri Poincaré*, **7**, 1–68.

References (continued)

- Hendriksen C, Lund E, Stromgard E (1984). Consequences of assessment and intervention among elderly people: a three year randomized controlled trial. *British Medical Journal*, **289**, 1522–1524.
- Jaynes ET (2003). *Probability Theory: The Logic of Science*. Cambridge UK: University Press.
- Kolmogorov A (1933). *Grundbegriffe der Wahrscheinlichkeitsrechnung*. Berlin: Springer.
- Krnjajić M, Kottas A, Draper D (2008). Parametric and nonparametric Bayesian model specification: a case study involving models for count data. *Computational Statistics and Data Analysis*, **78**, 2110–2128.