

Formalizing Statistical Beliefs in Hypothesis Testing Using Program Logic

Yusuke Kawamoto^{1,2}, Tetsuya Sato³, Kohei Suenaga⁴

¹AIST, Japan

²PRESTO, JST, Japan

³Tokyo Institute of Technology, Japan

⁴Kyoto University, Japan

Abstract

We propose a new approach to formally describing the requirement for statistical inference and checking whether the statistical method is appropriately used in a program. Specifically, we define belief Hoare logic (BHL) for formalizing and reasoning about the statistical beliefs acquired via hypothesis testing. This logic is equipped with axiom schemas for hypothesis tests and rules for multiple tests that can be instantiated to a variety of concrete tests. To the best of our knowledge, this is the first attempt to introduce a program logic with epistemic modal operators that can specify the preconditions for hypothesis tests to be applied appropriately.

1 Introduction

Statistical inference has been widely used to derive and justify scientific knowledge in a variety of academic disciplines, from natural sciences to social sciences. This has significantly increased the importance of statistics, but also brought concerns about the inappropriate procedure and the incorrect interpretation of statistics in scientific research. In fact, previous studies have pointed out that many research articles in biomedical science contain severe errors in the application and interpretation of statistical inference (Lang and Altman 2014). Furthermore, large proportions of these errors have been reported for basic statistical methods, possibly performed by researchers who can use only elementary techniques. For example, the concept of *statistical significance*, evaluated using *p-values*, has been commonly misused and misinterpreted (Wasserstein and Lazar 2016).

One of the main issues behind these human errors is that the logical aspects of statistical inference are described informally or implicitly using natural languages, and handled manually by analysts who may not fully understand the statistical methods. In particular, this makes them overlook some assumptions necessary for statistical methods, hence choosing inappropriate methods. Nevertheless, to our knowledge, no prior work on formal methods has specified the preconditions for statistical inference or verified the choice of statistical techniques.

In this paper, we propose a method for formalizing and reasoning about statistical inference using symbolic logic. Specifically, we introduce *belief Hoare logic (BHL)* to formalize the statistical beliefs acquired via hypothesis tests,

and to prevent errors in the choice of hypothesis tests by describing their preconditions explicitly. This is the first step to build a framework for formalizing and verifying the validity of empirical science on the basis of formal methods.

Contributions. Our main contributions are as follows:

- We propose a new approach to formalizing and reasoning about statistical inference in a program. In particular, this approach formalizes and checks the requirement for statistical methods to be used appropriately.
- We define an epistemic logic to express statistical beliefs obtained by hypothesis tests on datasets. Specifically, we formalize a statistical belief on a hypothesis φ as the knowledge that either φ holds or the sampled dataset is unlikelly far from the population. Then we introduce a Kripke semantics to define the interpretation of the logic.
- Using this epistemic logic, we construct belief Hoare logic (BHL) for formalizing and reasoning about the statistical inference based on hypothesis testing. Specifically, we define axiom schemas and rules for hypothesis tests that can be instantiated to a variety of concrete tests. In particular, BHL does not rely on a specific philosophy of statistics but can deal with both the frequentist and the Bayesian statistics by introducing corresponding axioms.
- We show that BHL is useful to reason about practical issues concerning statistical inference, such as the multiple comparison problem and *p*-value hacking.
- We provide a whole picture of the justification of statistical belief acquired via hypothesis tests inside and outside BHL. For instance, we discuss the empirical conditions and the epistemic aspects of statistical inference.

To the best of our knowledge, this is the first attempt to introduce a program logic that can specify the preconditions for hypothesis tests to be applied appropriately.

Related Work. The *Hoare logic* (Winskel 1993) is one of the program logic for an imperative programming language. This program logic is then extended and adapted so that it can handle various types of programs and assertions, including heap-manipulating programs (Reynolds 2002), hybrid systems (Suenaga and Hasuo 2011), and probabilistic programs (den Hartog and de Vink 2002). Atkinson and Carbin propose an extension of Hoare logic with epistemic assertions (Atkinson and Carbin 2020). In their work, an

epistemic assertion is used to reason about the belief of a program about a partially observable environment, whereas their logic does not deal with a statistical belief arising from statistical tests conducted in a program. To the best of our knowledge, ours is the first program logic that formalizes the concept of statistical beliefs in hypothesis testing.

Epistemic logic (von Wright 1951) is a branch of logic for reasoning about knowledge and belief (Fagin et al. 1995a; Halpern 2003), and is used to specify and verify a variety of knowledge properties in systems, e.g., authentication (Burrows, Abadi, and Needham 1990) and anonymity (Syverson and Stubblebine 1999; Garcia et al. 2005). Many previous works incorporate certain notions of degrees of belief and confidence (Huber and Schmidt-Petri 2008), but not in the sense of the statistical significance in hypothesis testing.

The first attempt to express statistical properties using modal logic is the work on statistical epistemic logic (StatEL) (Kawamoto 2019; Kawamoto 2020). They introduce a belief modality weaker than S5, and interpret it in a Kripke model with an accessibility relation defined in terms of a statistical distance between possible worlds. Unlike this work, however, StatEL cannot describe the procedures of statistical methods or reason about the appropriateness of inference.

From a broader perspective, many studies formalize and reason about programs based on knowledge (Fagin et al. 1995b) and beliefs (Laverny and Lang 2005), including belief updates in programs. For example, the situation calculus is extended to deal with the probabilistic degrees of beliefs in programs with noisy acting and sensing (Belle and Levesque 2015). However, no prior work has studied belief-based programs involving statistical hypothesis testing.

2 Preliminaries

In this section, we introduce notations used in this paper and recall background on statistical hypothesis testing.

Let \mathbb{N} , \mathbb{R} , $\mathbb{R}_{\geq 0}$ be the sets of non-negative integers, real numbers, and non-negative real numbers, respectively. Let $[0, 1]$ be the set of non-negative real numbers not greater than 1. We denote the dimension of a vector x by $\text{size}(x)$, and the set of all probability distributions over a set S by $\mathbb{D}S$.

Statistical Hypothesis Testing. *Statistical hypothesis testing* is a method of statistical inference about an unknown population x (the collection of items of interest) on the basis of a dataset y sampled from x . In a hypothesis test, an *alternative hypothesis* φ_1 is a proposition that we wish to prove about the population x , and a *null hypothesis* φ_0 is a proposition that contradicts φ_1 . The goal of the hypothesis test is to determine whether we *accept* the alternative hypothesis φ_1 by *rejecting* the null hypothesis φ_0 .

In a hypothesis test, we calculate a *test statistic* $t(y)$ from a dataset y , and see whether the $t(y)$ value contradicts the assumption that the null hypothesis φ_0 is true. Specifically, we calculate the *p-value*, showing the degree of likelihood of obtaining $t(y)$ when the null hypothesis φ_0 is true. If the *p-value* is smaller than a threshold (e.g., 0.05), we regard the dataset y is unlikely to be sampled from the population satisfying the null hypothesis φ_0 , hence we reject φ_0 and accept the alternative hypothesis φ_1 .

A hypothesis test is based on a *statistical model* $P(\xi, \theta)$ with unknown true parameters ξ , known parameters θ , and (assumed) probability distributions of the parameters in ξ .

Example 1 (*Z-test for two population means*). *As an illustrating example, we present the two-tailed Z-test for means of two populations with a known and equal variance. We introduce its statistical model as two normal distributions $N(\mu_{\text{pp}1}, \sigma^2)$ and $N(\mu_{\text{pp}2}, \sigma^2)$ with a known variance σ^2 and unknown true means $\mu_{\text{pp}1}, \mu_{\text{pp}2}$. Let y_1 and y_2 be two given datasets where each data value was sampled independently from $N(\mu_{\text{pp}1}, \sigma^2)$ and $N(\mu_{\text{pp}2}, \sigma^2)$, respectively.*

In the Z-test, we set the alternative hypothesis $\varphi_1 \stackrel{\text{def}}{=} (\mu_{\text{pp}1} \neq \mu_{\text{pp}2})$ and the null hypothesis $\varphi_0 \stackrel{\text{def}}{=} (\mu_{\text{pp}1} = \mu_{\text{pp}2})$. We calculate the Z-test statistic $t(y_1, y_2) = \frac{\text{mean}(y_1) - \text{mean}(y_2)}{\sigma \sqrt{1/\text{size}(y_1) + 1/\text{size}(y_2)}}$ where for $b = 1, 2$, $\text{size}(y_b)$ is the sample size of y_b and $\text{mean}(y_b)$ is the mean of all data in y_b . Then the p-value is defined by:

$$\Pr_{(d_1, d_2) \sim N(\mu_{\text{pp}1}, \sigma^2) \times N(\mu_{\text{pp}1}, \sigma^2)} [|t(d_1, d_2)| > |t(y_1, y_2)|]$$

under the null hypothesis φ_0 . When the p-value is small enough, the datasets y_1, y_2 are unlikely to be sampled from the same distribution, i.e., the null hypothesis $\mu_{\text{pp}1} = \mu_{\text{pp}2}$ is unlikely to hold. Hence, in the Z-test, if the p-value is smaller than a certain threshold (e.g., 0.05), we reject the null hypothesis φ_0 and accept the alternative hypothesis φ_1 .

When we have prior knowledge of $\mu_{\text{pp}1} \geq \mu_{\text{pp}2}$ (resp. $\mu_{\text{pp}1} \leq \mu_{\text{pp}2}$), then we apply the upper-tailed (resp. lower-tailed) Z-test with the alternative hypothesis $\mu_{\text{pp}1} > \mu_{\text{pp}2}$ (resp. $\mu_{\text{pp}1} < \mu_{\text{pp}2}$), and with the p-value $\Pr[t(d_1, d_2) > t(y_1, y_2)]$ (resp. $\Pr[t(d_1, d_2) < t(y_1, y_2)]$).

3 Illustrating Example

Throughout the paper, we use the following simple illustrating example to explain the basic ideas on our framework.

Example 2 (*Comparison tests on drugs*). *Let us consider three drugs 1, 2, 3 that may decrease blood pressure. To compare the efficacy of these drugs, we perform experiments and obtain a set y_i of the reduced values of blood pressure after taking drug i . Then we apply hypothesis tests on the dataset $y = (y_1, y_2, y_3)$. Let x_i be the true population from which the data values in y_i are sampled.*

Suppose that drug 1 is composed of drugs 2 and 3, and we want to know whether drug 1 has better efficacy than both drugs 2 and 3. Then we take the following procedure:

- *We first compare drugs 1 and 2 concerning the average decreases in blood pressure. We apply a two-tailed Z-test to see whether the means of the true populations x_1 and x_2 are different, i.e., $\text{mean}(x_1) \neq \text{mean}(x_2)$. In this test, the alternative hypothesis is $\varphi_{12} \stackrel{\text{def}}{=} (\text{mean}(x_1) \neq \text{mean}(x_2))$, and the null hypothesis is $\neg\varphi_{12} \equiv (\text{mean}(x_1) = \text{mean}(x_2))$.*
- *Let α_{ij} be the p-value when only comparing drugs i and j .*
- *If $\alpha_{12} \geq 0.05$, the Z-test does not reject the null hypothesis $\neg\varphi_{12}$ and concludes that the efficacy of drugs 1 and 2 may be the same. Then we are not interested in drug 1 any more, and skip the comparison with drug 3.*

- If $\alpha_{12} < 0.05$, the Z -test rejects the null hypothesis $\neg\varphi_{12}$ and concludes that the alternative hypothesis φ_{12} is true. Then we apply another Z -test to check whether the alternative hypothesis $\varphi_{13} \stackrel{\text{def}}{=} (\text{mean}(x_1) \neq \text{mean}(x_3))$ is true.
- Finally, we calculate the p -value of the combined tests, with the conjunctive alternative hypothesis $\varphi_{12} \wedge \varphi_{13}$.

Note that these Z -tests assume that the distribution of each population x_i is a normal distribution with a variance σ^2 .

Overview of the Framework. In our framework, we describe a procedure of statistical tests as a program using a programming language (Section 6); in Example 2, we denote the Z -test program comparing drugs i with j by C_{ij} , and the whole procedure by:

$$C_{\text{drug}} \stackrel{\text{def}}{=} C_{12}; \text{ if } \alpha_{12} < 0.05 \text{ then } C_{13} \text{ else skip} \quad (1)$$

Then we use an assertion logic (Section 5) to describe the requirement for the hypothesis tests as a *precondition* formula. In Example 2, the precondition is given by:

$$\psi_{\text{pre}} \stackrel{\text{def}}{=} \bigwedge_{i=1,2,3} \psi_i \wedge \bigwedge_{(i,j)=(1,2),(1,3), (2,1),(3,1)} \mathbf{P}(\text{mean}(x_i) < \text{mean}(x_j))$$

$$\text{where } \psi_i \stackrel{\text{def}}{=} (x_i \approx N(\mu_i, \sigma) \wedge y_i \overset{\sim}{\leftarrow}_{n_i} x_i).$$

In this formula, ψ_i represents that the true population x_i follows a normal distribution $N(\mu_i, \sigma)$ (with an unknown true mean μ_i), and that a set y_i of n_i data is sampled from x_i . $\mathbf{P}(\text{mean}(x_1) < \text{mean}(x_2))$ and $\mathbf{P}(\text{mean}(x_1) > \text{mean}(x_2))$ represent that both the lower-tail and upper-tail are possible, hence the test C_{12} should be two-tailed. Remark that our assertion logic never deals with quantifiers over variables.

The statistical belief we want to acquire is specified as a *postcondition* formula. In Example 2, the postcondition is:

$$\varphi_{\text{post}} \stackrel{\text{def}}{=} (\mathbf{K}_y^{\leq 0.05} \varphi_{12} \rightarrow \mathbf{K}_y^{\leq \min(\alpha_{12}, \alpha_{13})} (\varphi_{12} \wedge \varphi_{13})). \quad (2)$$

Intuitively, by testing on a dataset y , when we believe φ_{12} with a p -value $\alpha \leq 0.05$, we believe the combined hypothesis $\varphi_{12} \wedge \varphi_{13}$ with a p -value at most $\min(\alpha_{12}, \alpha_{13})$.

Finally, we combine all the above and describe the whole statistical inference as a *judgment*. In Example 2, we write:

$$\Gamma \vdash \{\psi_{\text{pre}}\} C_{\text{drug}} \{\varphi_{\text{post}}\} \quad (3)$$

By proving this judgment using rules in BHL (Section 7), we conclude that the statistical inference is appropriate.

We remark that the p -value can be larger for a different purpose of testing. Suppose that in Example 2, drug 1 was a new drug and we wanted to find out it had better efficacy than *at least one* of drugs 2 and 3. Then the procedure is:

$$C_{\text{multi}} \stackrel{\text{def}}{=} C_{12} \parallel C_{13}, \quad (4)$$

and the alternative hypothesis is $\varphi_{12} \vee \varphi_{13}$ with a p -value larger than α_{12} and α_{13} (at most $\alpha_{12} + \alpha_{13}$). This is the *multiple comparisons problem* (Bretz, Hothorn, and Westfall 2010), arising when the combined alternative hypothesis is in *disjunctive* form. We explain more details in Section 7.

4 Model

In this section, we introduce a Kripke model for describing statistical properties and formally define hypothesis tests.

4.1 Variables, Data, and Actions

We introduce a set Var of variables comprised of two disjoint sets of *invisible variables* and of *observable variables*: $\text{Var} = \text{Var}_{\text{inv}} \cup \text{Var}_{\text{obs}}$. We can directly observe the values of the latter, but not those of the former. We use x as an invisible variable denoting a population, and y as an observable variable denoting a dataset sampled from the population.

We write \mathcal{O} for the set of all data values we deal with, including the Boolean values, integers, real numbers, and lists of data. A *dataset* is a list of lists of data. In particular, we deal with a list of real vectors as a dataset. Then the vectors range over $\mathcal{X} = \mathbb{R}^l$ for an $l \in \mathbb{N}$, a distribution over the population has type $\mathbb{D}\mathcal{X}$, and a dataset has type list \mathcal{X} .

We write $d \sim D^n$ for the sampling of a set d of n data from a population D where all these data are independent and identically distributed (i.i.d.). Let Smpl be a set of i.i.d. samplings of datasets from populations (e.g., $d \sim D^n$), and Cmd be a set of program commands¹ (e.g., $v := e$ and skip). Then we define an *action* as a sampling of a dataset or a program command; i.e., $\text{Act} = \text{Smpl} \cup \text{Cmd}$.

4.2 States and Possible Worlds

A *state* is a pair (m, a) consisting of the current assignment $m : \text{Var} \rightarrow \mathcal{O}$ of data values to variables, and the action $a \in \text{Act}$ that has been executed in the last transition.

A *possible world* w is a sequence of the states $(w[0], w[1], \dots, w[k-1])$ where $w[i]$ represents the i -th state. The length k is denoted by $\text{len}(w)$. $w[0]$ and $w[k-1]$ are called the *initial state* and the *current state*, respectively. Since a possible world records all updates of data values, it can be used to model the updates of knowledge and beliefs as with previous works on epistemic logic (Fagin et al. 1995a).

The *observation* of a state $w[i] = (m, a)$ is defined by $\text{obs}(w[i]) = (m_{\text{o}}, a)$ with an assignment $m_{\text{o}} : \text{Var}_{\text{obs}} \rightarrow \mathcal{O}$ such that $m_{\text{o}}(v) = m(v)$ for all $v \in \text{Var}_{\text{obs}}$. The *observation* of a possible world w is defined by $\text{obs}(w) = (\text{obs}(w[0]), \dots, \text{obs}(w[k-1]))$. We abuse notations and denote by $w : \text{Var} \rightarrow \mathcal{O}$ the assignment of data values to variables in the current state of a possible world w .

4.3 Kripke Model

We introduce a *Kripke model with labeled transitions* where two kinds of relations $\overset{a}{\rightarrow}$ and \mathcal{R} may relate possible worlds.

A *transition relation* $w \overset{a}{\rightarrow} w'$ represents a transition from a world w to another w' by performing an action a . An *observability relation* $w\mathcal{R}w'$ represents that two possible worlds w and w' have the same observation, i.e., $\text{obs}(w) = \text{obs}(w')$. In Section 5, this relation is used to model the knowledge in the conventional Hintikka-style.

Definition 1 (Kripke model). We define a *Kripke model* as a tuple $\mathfrak{M} = (\mathcal{W}, (\overset{a}{\rightarrow})_{a \in \text{Act}}, \mathcal{R}, (V_w)_{w \in \mathcal{W}})$ consisting of:

¹In Section 6, we instantiate Cmd with a concrete example of commands used in a simple programming language.

- a non-empty set \mathcal{W} of possible worlds;
- for each $a \in \text{Act}$, a transition relation $\xrightarrow{a} \subseteq \mathcal{W} \times \mathcal{W}$;
- an observability relation $\mathcal{R} = \{(w, w') \in \mathcal{W} \times \mathcal{W} \mid \text{obs}(w) = \text{obs}(w')\}$;
- for each $w \in \mathcal{W}$, a valuation $V_w : \text{Pred} \rightarrow \mathcal{P}(\mathcal{O}^k)$ that maps a k -ary predicate to a set of k -tuples of data;

We assume $w \xrightarrow{a} w'$ implies $w'[\text{len}(w') - 1] = (m, a)$ for some m . We also assume that each world in a model has the same sets Var_{inv} and Var_{obs} of variables.

4.4 Formulation of Hypothesis Testing

Next, we formalize hypothesis tests. We consider a *basic test type* $s \in \{\text{L}, \text{U}, \text{T}\}$ each representing a *lower-tailed*, *upper-tailed*, and *two-tailed* test. A *hypothesis test* is a tuple $A_{\varphi_0}^{(s)} = (\varphi_0, t, D_{t, \varphi_0}, \preceq_t^{(s)}, P(\xi, \theta))$ consisting of:

- φ_0 is an assertion, called a null hypothesis;
- t is a function that maps a dataset $d \in \text{list } \mathcal{X}$ to its test statistic $t(d)$, usually with $\text{range}(t) = \mathbb{R}^k$ for a $k \geq 1$;
- $D_{t, \varphi_0} \in \mathbb{D}(\text{range}(t))$ is a probability distribution of the test statistic when the null hypothesis φ_0 is true;
- $\preceq_t^{(s)} \in \text{range}(t) \times \text{range}(t)$ is a *likeliness relation* where for a test type s and for values d and d' of the test statistic, $d \preceq_t^{(s)} d'$ represents that d is at most as likely as d' . For brevity, we often omit t and $^{(s)}$ to write $\preceq^{(s)}$ and \preceq .
- $P(\xi, \theta)$ is a statistical model with unknown parameters ξ and known parameters θ that characterizes the population.

Then we define a *hypothesis testing* over a set Φ of possible hypotheses by $A^{(s)} = (A_{\varphi}^{(s)})_{\varphi \in \Phi}$. We denote by \mathcal{A} a finite set of hypothesis testings we consider. For brevity, we sometimes omit $^{(s)}$ to write A_{φ} and A . We also often omit the statistical model $P(\xi, \theta)$ from the description of $A_{\varphi_0}^{(s)}$.

Example 3 (The likeliness relation for Z -test). *The two-tailed Z -test for two populations (Example 1) is denoted by $A_{\varphi_0} = (\varphi_0, t, N(0, 1), \preceq_t^{(\text{T})}, N(\mu_{\text{pp1}}, \sigma^2) \times N(\mu_{\text{pp2}}, \sigma^2))$. The likeliness relation $d \preceq_t^{(\text{T})} d'$ expresses $|d| \geq |d'|$. When the null hypothesis φ_0 is true, the test statistic $t(y_1, y_2)$ follows the standard normal distribution $N(0, 1)$, hence*

$$\Pr[t(y_1, y_2) \preceq_t^{(\text{T})} 1.96] = \Pr[|t(y_1, y_2)| \geq 1.96] = 0.05.$$

In contrast, for the upper-tailed (lower-tailed) test, with alternative hypothesis $\varphi_{\text{U}} \stackrel{\text{def}}{=} (\mu_{\text{pp1}} > \mu_{\text{pp2}})$ (resp. $\varphi_{\text{L}} \stackrel{\text{def}}{=} (\mu_{\text{pp1}} < \mu_{\text{pp2}})$), the likeliness relation $d \preceq_t^{(\text{U})} d'$ (resp. $d \preceq_t^{(\text{L})} d'$) is defined by $d \geq d'$ (resp. $d \leq d'$).

Next we define the *disjunctive combination* of two hypothesis tests $A_{\varphi_b} = (\varphi_b, t_b, D_{t_b, \varphi_b}, \preceq_{t_b}^{(s_b)}, P_b)$ for $b = 1, 2$ by $A_{\varphi_1 \vee \varphi_2} = (\varphi_1 \vee \varphi_2, t, D_{t, (\varphi_1, \varphi_2)}, \preceq_t^{(s_1, s_2)}, P)$ where $t(y_1, y_2) = (t_1(y_1), t_2(y_2))$, $D_{t, (\varphi_1, \varphi_2)} = D_{t_1, \varphi_1} \times D_{t_2, \varphi_2}$, $(d_1, d_2) \preceq_t^{(s_1, s_2)} (d'_1, d'_2)$ iff either $d_1 \preceq_{t_1}^{(s_1)} d'_1$ or $d_2 \preceq_{t_2}^{(s_2)} d'_2$, and $P = P_1 \times P_2$. Similarly, we define the *conjunctive combination* by $A_{\varphi_1 \wedge \varphi_2} = (\varphi_1 \wedge \varphi_2, t, D_{t, (\varphi_1, \varphi_2)}, \preceq_t^{(s_1, s_2)}, P)$ where $(d_1, d_2) \preceq_t^{(s_1, s_2)} (d'_1, d'_2)$ iff $d_1 \preceq_{t_1}^{(s_1)} d'_1$ and $d_2 \preceq_{t_2}^{(s_2)} d'_2$.

5 Assertion Language

Next we define an *assertion logic* that can express epistemic properties including knowledge and statistical beliefs.

5.1 Syntax of the Assertion Logic

We introduce two kinds of epistemic modality \mathbf{K} and $\mathbf{K}_{y, A}^{< \epsilon}$. Intuitively, a *knowledge* $\mathbf{K}\varphi$ expresses that we know φ , and this has been studied in a lot of previous work on epistemic logic. In contrast, a *statistical belief* $\mathbf{K}_{y, A}^{< \epsilon} \varphi$ expresses that we believe a hypothesis φ based on a statistical test A on an observed dataset y with a certain error level (p -value) at most ϵ . We formalize this as the knowledge that either the hypothesis φ holds or the observed dataset y is unluckily far from the population (from which y is sampled).

Formally, for a set Var of variables and a set Pred of predicates, the set Fml of formulas are defined by:

$$\varphi ::= \eta(x_1, \dots, x_n) \mid \neg\varphi \mid \varphi \vee \varphi \mid \mathbf{K}\varphi$$

where $\eta \in \text{Pred}$ and $x_1, \dots, x_n \in \text{Var}$. The formulas have no quantifiers over variables. We denote the set of all observable (resp. invisible) variables occurring in a formula φ by $\text{fv}^{\text{obs}}(\varphi)$ (resp. $\text{fv}^{\text{inv}}(\varphi)$). Let $\text{fv}(\varphi) = \text{fv}^{\text{obs}}(\varphi) \cup \text{fv}^{\text{inv}}(\varphi)$.

As syntax sugar, we use *conjunction* \wedge , *implication* \rightarrow , and *epistemic possibility* \mathbf{P} , defined as usual by: $\varphi_0 \wedge \varphi_1 \stackrel{\text{def}}{=} \neg(\neg\varphi_0 \vee \neg\varphi_1)$, $\varphi_0 \rightarrow \varphi_1 \stackrel{\text{def}}{=} \neg\varphi_0 \vee \varphi_1$, and $\mathbf{P}\varphi \stackrel{\text{def}}{=} \neg\mathbf{K}\neg\varphi$.

We introduce three predicates for statistical inference:

- $x \approx P$ represents that a population x follows a probability distribution that we assume in a statistical model P .
- For $x = (x_1, \dots, x_k)$ and $n = (n_1, \dots, n_k)$, $y \overset{\leftarrow}{\sim}_n x$ represents that for each $i = 1, 2, \dots, k$, a dataset y_i is obtained by sampling n_i data from the population x_i .
- For $\boxtimes \in \{=, \leq, \geq, <, >\}$ and $\epsilon \in [0, 1]$, $\tau_A^{\boxtimes \epsilon}(y)$ represents that the observation of a dataset y is unlikely to occur (with exception $\boxtimes \epsilon$) according to a hypothesis test A , and the dataset y is used in no other hypothesis tests.

We introduce a set $\text{Pred}_{\mathcal{G}}$ of *global predicate*, whose interpretations are identical in every possible worlds. We assume $\approx \in \text{Pred}_{\mathcal{G}}$ and $\text{Pred} = \text{Pred}_{\mathcal{G}} \cup \{\tau_A^{\boxtimes \epsilon}, \overset{\leftarrow}{\sim}\}$.

As syntax sugar, we introduce the *statistical belief modality* $\mathbf{K}_{y, A}^{\boxtimes \epsilon}$ such that for a formula φ representing a hypothesis,

$$\mathbf{K}_{y, A}^{\boxtimes \epsilon} \varphi \stackrel{\text{def}}{=} \mathbf{K}(\varphi \vee \tau_{A-\varphi}^{\boxtimes \epsilon}(y))$$

where $A_{\neg\varphi}$ is a test with a null hypothesis $\neg\varphi$. Then we define the *statistical possibility* $\mathbf{P}_{y, A}^{\boxtimes \epsilon}$ by $\mathbf{P}_{y, A}^{\boxtimes \epsilon} \varphi \stackrel{\text{def}}{=} \neg\mathbf{K}_{y, A}^{\boxtimes \epsilon} \neg\varphi$. For brevity, we write $\mathbf{K}_{y, A}^{\epsilon}$ instead of $\mathbf{K}_{y, A}^{< \epsilon}$. For a finite set \mathcal{A} of hypothesis testings, we write $\mathbf{K}_{y, A}^{\boxtimes \epsilon} \varphi \stackrel{\text{def}}{=} \bigvee_{A \in \mathcal{A}} \mathbf{K}_{y, A}^{\boxtimes \epsilon} \varphi$ and $\mathbf{P}_{y, A}^{\boxtimes \epsilon} \varphi \stackrel{\text{def}}{=} \bigvee_{A \in \mathcal{A}} \mathbf{P}_{y, A}^{\boxtimes \epsilon} \varphi$. A formula ψ is $\tau_A^{\boxtimes \epsilon}$ -free if $\tau_A^{\boxtimes \epsilon}$, $\mathbf{K}_{y, A}^{\boxtimes \epsilon}$, $\mathbf{P}_{y, A}^{\boxtimes \epsilon}$, \mathbf{K}_y , \mathbf{P}_y do not occur in ψ .

5.2 Semantics of the Assertion Logic

In this section, we define semantics for the assertion logic.

We define the interpretation of formulas in a world w in a Kripke model \mathfrak{M} (Definition 1) by:

$$\begin{aligned} \mathfrak{M}, w \models \eta(x_1, \dots, x_k) &\text{ iff } (w(x_1), \dots, w(x_k)) \in V_w(\eta) \\ \mathfrak{M}, w \models \neg\varphi &\text{ iff } \mathfrak{M}, w \not\models \varphi \\ \mathfrak{M}, w \models \varphi \vee \varphi' &\text{ iff either } \mathfrak{M}, w \models \varphi \text{ or } \mathfrak{M}, w \models \varphi' \\ \mathfrak{M}, w \models \mathbf{K}\varphi &\text{ iff for all } w' \in \mathcal{W}, (w, w') \in \mathcal{R} \\ &\text{ implies } \mathfrak{M}, w' \models \varphi. \end{aligned}$$

\mathfrak{M} is sometimes omitted when it is clear from the context.

Next we define the interpretation of predicates. Each global predicate has the same interpretation in all worlds; i.e., for any $\eta \in \text{Pred}_{\mathcal{G}}$ and $w, w' \in \mathcal{W}$, $V_w(\eta) = V_{w'}(\eta)$. Let $A_\varphi = (\varphi, t, D_{t,\varphi}, \preceq^{(s)})$ be a hypothesis test. Recall that the distribution over the population has type $\mathbb{D}\mathcal{X}$, and that $\preceq^{(s)}$ is the likeliness relation (Section 4.4). In a world w , we interpret the predicates by:

$$\begin{aligned} V_w(\approx) &= \{(X, D) \in \mathbb{D}\mathcal{X} \times \mathbb{D}\mathcal{X} \mid X = D\} \\ V_w(\leftrightarrow) &= \left\{ (d, D, n) \in (\text{list } \mathcal{X}) \times \mathbb{D}\mathcal{X} \times \mathbb{N} \mid \begin{array}{l} \text{There is an } i \in \mathbb{N} \\ \text{s.t. } w[i] \xrightarrow{d \sim D^n} w[i+1] \end{array} \right\} \\ V_w(\tau_{A_\varphi}^{\leq \epsilon}) &= \left\{ o \in \mathcal{O} \mid \begin{array}{l} \Pr_{d \sim D_{t,\varphi}} [d \preceq^{(s)} t(o)] < \epsilon \text{ and} \\ w \text{ has no transition where } o \text{ is used} \\ \text{in other hypothesis tests than } A_\varphi. \end{array} \right\}. \end{aligned}$$

Intuitively, $V_w(\tau_{A_\varphi}^{\leq \epsilon})^2$ is the set of all dataset that reject a null hypothesis φ . More specifically, the p -value $\Pr_{d \sim D_{t,\varphi}} [d \preceq^{(s)} t(o)]$ is the probability that a data d is at most as likely as the test statistic $t(o)$ when it is sampled from the distribution $D_{t,\varphi}$ in the world w ; e.g., $\Pr_{d \sim N(0,1)} [d \preceq^{(T)} 1.96] = \Pr_{d \sim N(0,1)} [|d| \geq 1.96] = 0.05$. By definition, we have:

$$\mathfrak{M}, w \models \tau_{A_\varphi}^{\leq \epsilon}(y) \text{ iff } \begin{array}{l} \Pr_{d \sim D_{t,\varphi}} [d \preceq^{(s)} t(w(y))] < \epsilon \text{ and} \\ w(y) \text{ is not used in other tests than } A_\varphi. \end{array}$$

This represents that in the possible world w , the observation of a dataset y is unlikely to occur (except with probability ϵ) according to the hypothesis test A_φ where the test statistic follows the distribution $D_{t,\varphi}$ in the world w .

Then the statistical belief modality $\mathbf{K}_{y,A}^{\leq \epsilon}$ is interpreted as:

$$\begin{aligned} \mathfrak{M}, w \models \mathbf{K}_{y,A}^{\leq \epsilon} \varphi & \\ \text{iff } \mathfrak{M}, w \models \mathbf{K}(\varphi \vee \tau_{A_\varphi}^{\leq \epsilon}(y)) & \\ \text{iff for all } w', (w, w') \in \mathcal{R} \text{ implies } \mathfrak{M}, w' \models \neg\varphi \rightarrow \tau_{A_\varphi}^{\leq \epsilon}(y). & \end{aligned}$$

Intuitively, $\mathbf{K}_{y,A}^{\leq \epsilon} \varphi$ represents a belief that an alternative hypothesis φ on the population is true. More specifically, $w' \models \neg\varphi \rightarrow \tau_{A_\varphi}^{\leq \epsilon}(y)$ means that if we consider a possible world w' where the null hypothesis $\neg\varphi$ is true, then a hypothesis test A_φ would conclude that the observation of a dataset y is unlikely to occur (with exceptions at most ϵ), i.e., $\tau_{A_\varphi}^{\leq \epsilon}(y)$ holds in w' . We discuss the implication of our formalization of $\mathbf{K}_{y,A}^{\leq \epsilon}$ in Section 5.3.

Although the modality \mathbf{K} expresses the knowledge in terms of S5, the syntax sugar $\mathbf{K}_{y,A}^{\leq \epsilon} \varphi$ represents belief instead of knowledge. This is because φ can be false when

²We also define the interpretation of $\tau_{A_\varphi}^{\geq \epsilon}(y)$ analogously.

$\tau_{A_\varphi}^{\leq \epsilon}(y)$ holds (i.e., we may have a false belief on φ when the sampled dataset y is unluckily far from the population).

Example 4 (Statistical belief in Z -tests). Recall again the two-tailed Z -test for two population means in Example 1.

The alternative hypothesis is $\varphi_1 \stackrel{\text{def}}{=} (\mu_{\text{pp1}} \neq \mu_{\text{pp2}})$, and the null hypothesis is $\varphi_0 \stackrel{\text{def}}{=} (\mu_{\text{pp1}} = \mu_{\text{pp2}})$. We denote this Z -test by $A_{\varphi_0} = (\varphi_0, t, N(0, 1), \preceq^{(T)})$.

Suppose that in a world w , we sample two datasets $w(y_1)$ and $w(y_2)$ respectively from two populations $w(x_1)$ and $w(x_2)$, and calculate the Z -test statistic $t(w(y_1), w(y_2))$ defined in Example 1. When the null hypothesis φ_0 is true, $t(w(y_1), w(y_2))$ follows the distribution $N(0, 1)$.

If $t(w(y_1), w(y_2)) = 3$, $\Pr_{d \sim N(0,1)} [d \preceq^{(T)} t(w(y_1), w(y_2))] < 0.05$. Then the null hypothesis φ_0 is rejected, and we obtain the statistical belief that the alternative hypothesis φ_1 is true with the significance level 0.05, i.e., $w \models \mathbf{K}_{y,A}^{\leq 0.05} \varphi_1$. In contrast, if $t(w(y_1), w(y_2)) = 1.8$, then $w \models \neg \mathbf{K}_{y,A}^{\leq 0.05} \varphi_1$ because $\Pr_{d \sim N(0,1)} [d \preceq^{(T)} t(w(y_1), w(y_2))] > 0.05$.

5.3 Remarks on the Formalization

Implication The universe \mathcal{W} of the model \mathfrak{M} is assumed to include all possible worlds we can imagine. If there exists no possible world satisfying the null hypothesis $\neg\varphi$ in the model \mathfrak{M} , then φ is satisfied in all worlds in \mathfrak{M} , hence so are $\mathbf{K}\varphi$ and $\mathbf{K}_{y,A}^{\leq \epsilon} \varphi$. This reflects the fact that if we cannot imagine a possible world where $\neg\varphi$ is true, then we already know that φ is true without executing hypothesis tests.

Type II error The *type II error rate* is the probability that the hypothesis test A does not reject the null hypothesis φ_{null} when φ_{null} is false. Assume that the true population satisfies a hypothesis ξ in the world w . Let y' be a dataset such that the p -value (type I error rate) α of the test A is 0.05; i.e., $w \models \mathbf{K}_{y',A}^{0.05} \neg\varphi_{\text{null}}$. Then the type II error rate β when $\alpha = 0.05$ is given by $w \models \mathbf{K}_{y',A}^\beta \neg\xi$.

5.4 Properties of Statistical Beliefs

The statistical possibility $\mathbf{P}_{y,A}^{\leq \epsilon} \varphi$ means that we think a null hypothesis φ may be true after a hypothesis test A did not reject φ with a significance level ϵ . Formally, we have:

$$\begin{aligned} \mathfrak{M}, w \models \mathbf{P}_{y,A}^{\leq \epsilon} \varphi & \\ \text{iff there is a } w' \text{ s.t. } (w, w') \in \mathcal{R} \text{ and } \mathfrak{M}, w' \not\models \neg\varphi \vee \tau_{A_\varphi}^{\leq \epsilon}(y) & \\ \text{iff } \mathfrak{M}, w \models \mathbf{P}(\varphi \wedge \neg\tau_{A_\varphi}^{\leq \epsilon}(y)). & \end{aligned}$$

Now we show basic properties of statistical beliefs.

Proposition 1 (Properties of $\mathbf{K}_{y,A}^{\leq \epsilon}$). Let φ be a formula, A be a hypothesis test, $y \in \text{Var}_{\text{obs}}$, and $\epsilon \in \mathbb{R}_{\geq 0}$.

1. Knowledge is also regarded as belief: $\models \mathbf{K}\varphi \rightarrow \mathbf{K}_{y,A}^{\leq \epsilon} \varphi$.
2. If we believe φ based on a test A , then we know this statistical belief; i.e., $\models \mathbf{K}_{y,A}^{\leq \epsilon} \varphi \rightarrow \mathbf{K} \mathbf{K}_{y,A}^{\leq \epsilon} \varphi$.
3. If we failed to reject φ and think it possible, then we know this possibility; i.e., $\models \mathbf{P}_{y,A}^{\leq \epsilon} \varphi \rightarrow \mathbf{K} \mathbf{P}_{y,A}^{\leq \epsilon} \varphi$.
4. If $\epsilon \leq \epsilon'$, $\models \mathbf{K}_{y,A}^{\leq \epsilon} \varphi \rightarrow \mathbf{K}_{y,A}^{\leq \epsilon'} \varphi$ and $\models \mathbf{P}_{y,A}^{\leq \epsilon'} \varphi \rightarrow \mathbf{P}_{y,A}^{\leq \epsilon} \varphi$.

5. $\mathbf{K}_{y,A}^{<\epsilon} \varphi$ may represent a false belief. The alternative hypothesis φ we believe may be false, i.e., the rejected null hypothesis $\neg\varphi$ may be true: $\epsilon > 0$ iff $\not\models \mathbf{K}_{y,A}^{<\epsilon} \varphi \rightarrow \varphi$.
6. $\models \mathbf{K}_{y,A}^{<\epsilon} \varphi \rightarrow \mathbf{K}_y^{<\epsilon} \varphi$.

The proofs are straightforward from definitions. See the full version of this paper for the proofs.

6 A Simple Programming Language

We introduce an imperative programming language Prog.

6.1 Syntax of Prog

We define the syntax of Prog by the following BNF:

$$\begin{aligned}
 T &::= \text{bool} \mid \text{int} \mid \text{real} \mid T \times T \mid \text{list}(T) && \text{(Types)} \\
 e &::= v \mid f(e_1, \dots, e_k) && \text{(Terms)} \\
 c &::= \text{skip} \mid v := e && \text{(Commands)} \\
 C &::= c \mid C; C \mid C \parallel C \mid \text{if } e \text{ then } C \text{ else } C \mid \text{loop } e \text{ do } C && \text{(Programs)}
 \end{aligned}$$

where v is an element of Var_{obs} , f is a (built-in) function symbol, and constants are dealt as functions with arity 0. Notice that a program can handle only observable variables.

T represents *types*. A type is either `bool` for Boolean values, `int` for integers, `real` for real numbers, $T_1 \times T_2$ for pairs consisting of a value of type T_1 and a value of type T_2 , or `list`(T) for lists of values of type T . e represents *expressions* that evaluate to values. An expression is either a variable v or a function call $f(e_1, \dots, e_k)$; the latter is typically a call to a function that computes a test statistic. c and C represent *commands* and *programs* respectively. We give their intuitive explanation as follows.

- `skip` does nothing.
- $v := e$ updates v with the result of an evaluation of e .
- $C_1; C_2$ executes C_1 and then C_2 .
- $C_1 \parallel C_2$ executes C_1 and C_2 in parallel that may share some data.
- `if` e `then` C_1 `else` C_2 executes C_1 if e evaluates to `true`; executes C_2 if e evaluates to `false`.
- `loop` e `do` C iteratively executes C as long as e evaluates to `true`.

For instance, the example program in Section 7.5 conforms to the programming language Prog.

Hereafter we assume that all programs are well-typed although we do not explicitly mention the types. Checking this condition for our language can be done by adapting a standard type-checking algorithm to our setting.

We write $\text{upd}(C)$ for the set of variables that may be updated by executing C : $\text{upd}(\text{skip}) = \emptyset$, $\text{upd}(v := e) = \{v\}$, $\text{upd}(C_1; C_2) = \text{upd}(C_1 \parallel C_2) = \text{upd}(\text{if } e \text{ then } C_1 \text{ else } C_2) = \text{upd}(C_1) \cup \text{upd}(C_2)$, and $\text{upd}(\text{loop } e \text{ do } C) = \text{upd}(C)$.

Then we impose the following restriction to every occurrence of $C_1 \parallel C_2$: $\text{upd}(C_1) \cap \text{Var}(C_2) = \text{upd}(C_2) \cap \text{Var}(C_1) = \emptyset$. This restriction is to ensure that an execution of C_1 does not interfere with that of C_2 .

6.2 Semantics of Prog

We define the semantics of Prog over a Kripke model \mathfrak{M} with labeled transitions given in Section 4.3. The semantics is based on the standard structural operational semantics (e.g. (Nielsen and Nielson 2007)).

For a possible world $w \in \mathcal{W}$ and $n = \text{len}(w)$, we write

$$w = w[0], w[1], \dots, w[n-2], (m, a)$$

where (m, a) is the current state $w[n-1]$ with an assignment $m: \text{Var} \rightarrow \mathcal{O}$ and an action a in the last transition in \mathfrak{M} .

For the assignment m of the current state of w , we define the evaluation $\llbracket e \rrbracket_m$ of a term e inductively by $\llbracket v \rrbracket_m = m(v)$ and $\llbracket f(e_1, \dots, e_k) \rrbracket_m = \llbracket f \rrbracket(\llbracket e_1 \rrbracket_m, \dots, \llbracket e_k \rrbracket_m)$.

As in Figure 1, we define a binary relation

$$\longrightarrow \subseteq (\text{Prog} \times \mathcal{W}) \times ((\text{Prog} \times \mathcal{W}) \cup \mathcal{W})$$

that relates a pair $\langle C, w \rangle$ consisting of a program C and a possible world w to its next step of execution. If C is terminated, the next step will be a possible world w' , otherwise the execution continues to the $\langle C', w' \rangle$.

$$\begin{aligned}
 \langle \text{skip}, w \rangle &\longrightarrow w; (m, \text{skip}), \\
 \langle v := e, w \rangle &\longrightarrow w; (m[v \mapsto \llbracket e \rrbracket_m], v := e), \\
 \langle C_1, w \rangle &\longrightarrow w' \\
 \hline
 \langle C_1; C_2, w \rangle &\longrightarrow \langle C_2, w' \rangle \\
 \langle C_1, w \rangle &\longrightarrow \langle C'_1, w' \rangle \\
 \hline
 \langle C_1; C_2, w \rangle &\longrightarrow \langle C'_1; C_2, w' \rangle \\
 \langle \text{if } e \text{ then } C_1 \text{ else } C_2, w \rangle &\longrightarrow \begin{cases} \langle C_1, w \rangle & \llbracket e \rrbracket_m = \top \\ \langle C_2, w \rangle & \llbracket e \rrbracket_m = \perp \end{cases} \\
 \langle \text{loop } e \text{ do } C, w \rangle &\longrightarrow \begin{cases} \langle C; \text{loop } e \text{ do } C, w \rangle & \llbracket e \rrbracket_m = \top \\ w & \llbracket e \rrbracket_m = \perp \end{cases} \\
 \langle C_1, w \rangle &\longrightarrow \langle C'_1, w' \rangle \\
 \hline
 \langle C_1 \parallel C_2, w \rangle &\longrightarrow \langle C'_1 \parallel C_2, w' \rangle \\
 \langle C_1, w \rangle &\longrightarrow w' \\
 \hline
 \langle C_1 \parallel C_2, w \rangle &\longrightarrow \langle C_2, w' \rangle \\
 \langle C_2, w \rangle &\longrightarrow \langle C'_2, w' \rangle \\
 \hline
 \langle C_1 \parallel C_2, w \rangle &\longrightarrow \langle C_1 \parallel C'_2, w' \rangle \\
 \langle C_2, w \rangle &\longrightarrow w' \\
 \hline
 \langle C_1 \parallel C_2, w \rangle &\longrightarrow \langle C_1, w' \rangle
 \end{aligned}$$

Figure 1: Rules of execution of programs.

Remark that the semantics of a program contains the trace of commands executed in it. Hence, even if programs finally have the same result, their semantics may be different.

$$\begin{aligned}
 \langle v := v + 1, ([v \mapsto 1], a) \rangle &\longrightarrow ([v \mapsto 1], a), ([v \mapsto 2], v := 1 + 1) \\
 \langle v := 2 * v, ([v \mapsto 1], a) \rangle &\longrightarrow ([v \mapsto 1], a), ([v \mapsto 2], v := 2 * v)
 \end{aligned}$$

We define the semantic relation $\llbracket C \rrbracket \subseteq \mathcal{W} \times \mathcal{W}$ by

$$\llbracket C \rrbracket(w) = \{w' \mid \langle C, w \rangle \longrightarrow^* w'\}$$

where \longrightarrow^* is the transitive closure of \longrightarrow .

Remark on Parallel Compositions Since parallel compositions are nondeterministic, $w' \in \llbracket C \rrbracket(w)$ may not be unique. However, the resulting world w' are essentially the same, because for parallel composition $C_1 \parallel C_2$, the world $w' \in \llbracket C_1 \parallel C_2 \rrbracket(w)$ is convertible to a pair of $w_1 \in \llbracket C_1 \rrbracket(w)$ and $w_2 \in \llbracket C_2 \rrbracket(w)$ and vice versa.

If we have $\langle C_b, w \rangle \xrightarrow{*} w; u_b$ for $b = 1, 2$, then by $\text{upd}(C_b) \cap \text{Var}(C_{3-b}) = \emptyset$, we obtain a sequence u' such that $\langle C_1 \parallel C_2, w \rangle \xrightarrow{*} w; u'$ by combining u_1 and u_2 .

Conversely, if $\langle C_1 \parallel C_2, w \rangle \xrightarrow{*} w; u'$, we can decompose u' into u_1 and u_2 such that $\langle C_b, w \rangle \xrightarrow{*} w; u_b$ for $b = 1, 2$ (for detail, see the full version of this paper).

Then, for any pair of τ^∞ -free assertions φ_1 and φ_2 satisfying $\text{upd}(C_b) \cap \text{fv}(\varphi_{3-b}) = \emptyset$ for $b = 1, 2$, for any possible world $w' \in \llbracket C_1 \parallel C_2 \rrbracket(w)$, we can write $w' = w; u'$, and decompose u' into u_1 and u_2 as above. We then obtain $w' \models \varphi_1 \wedge \varphi_2$ iff $w; u_1 \models \varphi_1$ and $w; u_2 \models \varphi_2$.

Procedures of Hypothesis Testing Finally, we present the interpretation of a program $f_{A_{\varphi_0}^{(s)}}$ for a hypothesis test

$A_{\varphi_0}^{(s)} = (\varphi_0, t, D_{t, \varphi_0}, \preceq^{(s)})$ with a null hypothesis φ_0 , a test statistic t , and a test type s . For a dataset y and an assignment m , $\llbracket f_{A_{\varphi_0}^{(s)}}(y) \rrbracket_m$ represents the p -value:

$$\llbracket f_{A_{\varphi_0}^{(s)}}(y) \rrbracket_m = \Pr_{d \sim D_{t, \varphi_0}} [d \preceq^{(s)} t(m(y))], \quad (5)$$

which is the probability that a data d is at most as likely as the test statistic $t(m(y))$ when it is sampled from D_{t, φ_0} in the world where the null hypothesis φ_0 is true.

7 Belief Hoare Logic for Hypothesis Testing

We introduce *belief Hoare logic* (BHL) for formalizing and reasoning about statistical inference using hypothesis tests. Then we describe the reasoning about the *multiple comparison problem* and *p-value hacking* using BHL.

7.1 Hoare Triples

An *environment* is defined as a pair $\Gamma = (\Gamma^{\text{inv}}, \Gamma^{\text{obs}})$ consisting of an *invisible environment* Γ^{inv} and an *observable environment* Γ^{obs} that assign types to invisible variables and to observable variables, respectively. We denote by $\text{Var}(\Gamma)$ the set of all variables occurring in an environment Γ , and by Env the set of all possible environments.

A *judgment* is of the form $\Gamma \vdash \{\psi\} C \{\varphi\}$ where $\Gamma \in \text{Env}$, $\psi, \varphi \in \text{Fml}$, and $C \in \text{Prog}$. Intuitively, this represents that when the *precondition* ψ is satisfied, executing the program C results in satisfying the *postcondition* φ .

We say that a judgment $\Gamma \vdash \{\psi\} C \{\varphi\}$ is *valid* iff for any model \mathfrak{M} and any possible world w , if $\mathfrak{M}, w \models \psi$, then $\mathfrak{M}, w' \models \varphi$ for all $w' \in \llbracket C \rrbracket(w)$.

We write $\Gamma \models \varphi$ if $\mathfrak{M}, w \models \varphi$ for any model \mathfrak{M} and any world w that respects the type information in Γ (i.e., the type of $w(v)$ being $\Gamma(v)$ for any variable $v \in \text{Var}$).

7.2 Inference Rules

Next, we present the inference rules for belief Hoare logic (BHL). The rules are classified into those for basic command constructs (Figure 2) and for hypothesis testing constructs

$$\begin{array}{c} \Gamma \vdash \{\psi\} \text{skip} \{\varphi\} \quad (\text{SKIP}) \\ \frac{\Gamma(x) = \Gamma(y)}{\Gamma \vdash \{\varphi[x \mapsto y]\} x := y \{\varphi\}} \quad (\text{UPDVAR}) \\ \frac{\Gamma \vdash \{\psi\} C_1 \{\psi'\} \quad \Gamma \vdash \{\psi'\} C_2 \{\varphi\}}{\Gamma \vdash \{\psi\} C_1; C_2 \{\varphi\}} \quad (\text{SEQ}) \\ \frac{\Gamma \vdash \{\psi\} C \{\varphi\}}{\Gamma \vdash \{\psi\} \text{skip} \parallel C \{\varphi\}} \quad (\text{PAR-SKIPADDL}) \\ \frac{\Gamma \vdash \{\psi\} C \{\varphi\}}{\Gamma \vdash \{\psi\} C \parallel \text{skip} \{\varphi\}} \quad (\text{PAR-SKIPADDR}) \\ \frac{\Gamma \vdash \{\psi\} \text{skip} \parallel C \{\varphi\}}{\Gamma \vdash \{\psi\} C \{\varphi\}} \quad (\text{PAR-SKIPRML}) \\ \frac{\Gamma \vdash \{\psi\} C \parallel \text{skip} \{\varphi\}}{\Gamma \vdash \{\psi\} C \{\varphi\}} \quad (\text{PAR-SKIPRMR}) \\ \frac{\Gamma \vdash \{\psi \wedge e\} C_1 \{\varphi\} \quad \Gamma \vdash \{\psi \wedge \neg e\} C_2 \{\varphi\}}{\Gamma \vdash \{\psi\} \text{if } e \text{ then } C_1 \text{ else } C_2 \{\varphi\}} \quad (\text{IF}) \\ \frac{\Gamma \vdash \{\psi \wedge e\} C \{\psi\}}{\Gamma \vdash \{\psi\} \text{loop } e \text{ do } C \{\psi \wedge \neg e\}} \quad (\text{LOOP}) \\ \frac{\Gamma \models \psi \rightarrow \psi' \quad \Gamma \vdash \{\psi'\} C \{\varphi'\} \quad \Gamma \models \varphi' \rightarrow \varphi}{\Gamma \vdash \{\psi\} C \{\varphi\}} \quad (\text{CONSEQ}) \end{array}$$

Figure 2: Rules for basic constructs for commands.

(Figure 3). The latter includes axiom schemas that can be instantiated to a variety of concrete hypothesis test methods; we present such instantiation in Section 7.3.

The rules for basic constructs in Figure 2 are standard; the readers are referred to a standard textbook on Hoare logic (Winskel 1993) for details. We remark the features:

- In the rules IF and LOOP, the guard condition e is a Boolean expression implicitly used as a logical predicate in the preconditions and the postconditions as usual.
- The rule CONSEQ is used to weaken the precondition and strengthen the postcondition of a triple. The relation $\Gamma \models \varphi$ defined above is used in this rule.

Schemas for Single Hypothesis Testing In Figure 3 the axiom schemas TWO-T, LOW-T, and UP-T correspond to two-tailed, lower-tailed, and upper-tailed tests, respectively.

In these schemas, a dataset y is sampled from a population x , which follows a statistical model $P(\xi; \theta)$ with unknown true parameters $\xi \in \Phi$ and known parameters $\theta \in \Theta$.

To reason about the unknown parameter ξ , we perform a hypothesis test A_{φ_0} with a null hypothesis φ_0 . Let φ_L, φ_U , and $\varphi_T \stackrel{\text{def}}{=} \varphi_L \vee \varphi_U$ be the alternative hypotheses for the lower-tailed, upper-tailed, and two-tailed tests, respectively.

When we consider both the lower-tail φ_L and upper-tail φ_U are possible before performing the test, we express them by $\mathbf{P}\varphi_L$ and by $\mathbf{P}\varphi_U$ in the precondition. Then we apply the schema TWO-T. When the two-tailed test $f_{A_{\varphi_0}^{(\text{T})}}(y)$ returns the p -value $\alpha \in [0, 1]$, we obtain a statistical belief on the alternative hypothesis φ_T , namely, $\mathbf{K}_{y,A}^\alpha \varphi_T$. When we consider only the lower-tail φ_L (resp. upper tail φ_U) possible, we apply LOW-T (resp. UP-T) and obtain a statistical belief on φ_L (resp. φ_U), namely, $\mathbf{K}_{y,A}^\alpha \varphi_L$ (resp. $\mathbf{K}_{y,A}^\alpha \varphi_U$).

$$\begin{array}{c}
\Gamma^{\text{inv}} = \{\xi : \Xi, x : \mathbb{D}\mathcal{X}\} \cup \text{Var}_{\text{inv}}(\{\psi', \varphi_L, \varphi_U\}), \quad \Gamma^{\text{obs}} = \{\theta : \Theta, n : \mathbb{N}, y : \text{list } \mathcal{X}, \alpha : [0, 1]\} \cup \text{Var}_{\text{obs}}(\{\psi', \varphi_L, \varphi_U\}), \\
\alpha \notin \text{fv}(\{\varphi_L, \varphi_U\}), \quad \psi' : \tau^{\text{pd}}\text{-free}, \quad \psi \stackrel{\text{def}}{=} (x \approx P(\xi, \theta) \wedge y \stackrel{\leftarrow}{\sim}_n x \wedge \psi') \\
\hline
(\Gamma^{\text{inv}}, \Gamma^{\text{obs}}) \vdash \{\psi[\alpha \mapsto f_{A_{\varphi_0}^{(\top)}}(y)] \wedge \mathbf{P}\varphi_L \wedge \mathbf{P}\varphi_U \mid \alpha := f_{A_{\varphi_0}^{(\top)}}(y) \mid \psi \wedge \mathbf{K}_{y,A}^{\alpha} \varphi_{\top}\} \quad (\text{Two-T}) \\
\\
\Gamma^{\text{inv}} = \{\xi : \Xi, x : \mathbb{D}\mathcal{X}\} \cup \text{Var}_{\text{inv}}(\{\psi', \varphi_L, \varphi_U\}), \quad \Gamma^{\text{obs}} = \{\theta : \Theta, n : \mathbb{N}, y : \text{list } \mathcal{X}, \alpha : [0, 1]\} \cup \text{Var}_{\text{obs}}(\{\psi', \varphi_L, \varphi_U\}), \\
\alpha \notin \text{fv}(\{\varphi_L, \varphi_U\}), \quad \psi' : \tau^{\text{pd}}\text{-free}, \quad \psi \stackrel{\text{def}}{=} (x \approx P(\xi, \theta) \wedge y \stackrel{\leftarrow}{\sim}_n x \wedge \psi') \\
\hline
(\Gamma^{\text{inv}}, \Gamma^{\text{obs}}) \vdash \{\psi[\alpha \mapsto f_{A_{\varphi_0}^{(\downarrow)}}(y)] \wedge \mathbf{P}\varphi_L \wedge \neg \mathbf{P}\varphi_U \mid \alpha := f_{A_{\varphi_0}^{(\downarrow)}}(y) \mid \psi \wedge \mathbf{K}_{y,A}^{\alpha} \varphi_L\} \quad (\text{Low-T}) \\
\\
\Gamma^{\text{inv}} = \{\xi : \Xi, x : \mathbb{D}\mathcal{X}\} \cup \text{Var}_{\text{inv}}(\{\psi', \varphi_L, \varphi_U\}), \quad \Gamma^{\text{obs}} = \{\theta : \Theta, n : \mathbb{N}, y : \text{list } \mathcal{X}, \alpha : [0, 1]\} \cup \text{Var}_{\text{obs}}(\{\psi', \varphi_L, \varphi_U\}), \\
\alpha \notin \text{fv}(\{\varphi_L, \varphi_U\}), \quad \psi' : \tau^{\text{pd}}\text{-free}, \quad \psi \stackrel{\text{def}}{=} (x \approx P(\xi, \theta) \wedge y \stackrel{\leftarrow}{\sim}_n x \wedge \psi') \\
\hline
(\Gamma^{\text{inv}}, \Gamma^{\text{obs}}) \vdash \{\psi[\alpha \mapsto f_{A_{\varphi_0}^{(\cup)}}(y)] \wedge \neg \mathbf{P}\varphi_L \wedge \mathbf{P}\varphi_U \mid \alpha := f_{A_{\varphi_0}^{(\cup)}}(y) \mid \psi \wedge \mathbf{K}_{y,A}^{\alpha} \varphi_U\} \quad (\text{UP-T}) \\
\\
\text{For } b = 1, 2, \Gamma_b^{\text{inv}} = \text{fv}^{\text{inv}}(\{\psi_b, \psi'_b, \varphi_b\}), \quad \Gamma_b^{\text{obs}} = \{y_b : \text{list } \mathcal{X}, \alpha_b : [0, 1]\} \cup \text{fv}^{\text{obs}}(\{\psi_b, \psi'_b, \varphi_b\}) \cup \text{upd}(C_b), \\
\text{upd}(C_b) \cap (\text{fv}(\{\psi'_{3-b}, \varphi_{3-b}\}) \cup \{y_1, y_2\}) = \emptyset, \quad \psi'_b : \tau^{\text{pd}}\text{-free}, \quad \alpha_1, \alpha_2 \notin \text{fv}(\{\varphi_1, \varphi_2\}) \\
(\Gamma_1^{\text{inv}}, \Gamma_1^{\text{obs}}) \vdash \{\psi_1\} C_1 \mid \{\psi'_1 \wedge \mathbf{K}_{y_1}^{\alpha_1} \varphi_1\} \quad (\Gamma_2^{\text{inv}}, \Gamma_2^{\text{obs}}) \vdash \{\psi_2\} C_2 \mid \{\psi'_2 \wedge \mathbf{K}_{y_2}^{\alpha_2} \varphi_2\} \\
\hline
(\Gamma_1^{\text{inv}} \cup \Gamma_2^{\text{inv}}, \Gamma_1^{\text{obs}} \cup \Gamma_2^{\text{obs}}) \vdash \{\psi_1 \wedge \psi_2\} C_1 \parallel C_2 \mid \{\psi'_1 \wedge \psi'_2 \wedge \mathbf{K}_{(y_1, y_2)}^{\leq \alpha_1 + \alpha_2} (\varphi_1 \vee \varphi_2)\} \quad (\text{MULT-V}) \\
\\
\text{For } b = 1, 2, \Gamma_b^{\text{inv}} = \text{fv}^{\text{inv}}(\{\psi_b, \psi'_b, \varphi_b\}), \quad \Gamma_b^{\text{obs}} = \{y_b : \text{list } \mathcal{X}, \alpha_b : [0, 1]\} \cup \text{fv}^{\text{obs}}(\{\psi_b, \psi'_b, \varphi_b\}) \cup \text{upd}(C_b), \\
\text{upd}(C_b) \cap (\text{fv}(\{\psi'_{3-b}, \varphi_{3-b}\}) \cup \{y_1, y_2\}) = \emptyset, \quad \psi'_b : \tau^{\text{pd}}\text{-free}, \quad \alpha_1, \alpha_2 \notin \text{fv}(\{\varphi_1, \varphi_2\}) \\
(\Gamma_1^{\text{inv}}, \Gamma_1^{\text{obs}}) \vdash \{\psi_1\} C_1 \mid \{\psi'_1 \wedge \mathbf{K}_{y_1}^{\alpha_1} \varphi_1\} \quad (\Gamma_2^{\text{inv}}, \Gamma_2^{\text{obs}}) \vdash \{\psi_2\} C_2 \mid \{\psi'_2 \wedge \mathbf{K}_{y_2}^{\alpha_2} \varphi_2\} \\
\hline
(\Gamma_1^{\text{inv}} \cup \Gamma_2^{\text{inv}}, \Gamma_2^{\text{obs}} \cup \Gamma_2^{\text{obs}}) \vdash \{\psi_1 \wedge \psi_2\} C_1 \parallel C_2 \mid \{\psi'_1 \wedge \psi'_2 \wedge \mathbf{K}_{(y_1, y_2)}^{\leq \min(\alpha_1, \alpha_2)} (\varphi_1 \wedge \varphi_2)\} \quad (\text{MULT-}\wedge)
\end{array}$$

Figure 3: Axiom schemas and rules for hypothesis tests. TWO-T, LOW-T, and UP-T are schemas for two-tailed, lower-tailed, and upper-tailed tests, respectively. MULT-V is the rule for the Bonferroni's method, and MULT- \wedge is for the simultaneous tests without correction of α .

Rules for Multiple Tests The rule MULT-V corresponds to the multiple tests by the *Bonferroni's method*. As illustrated in Section 3, a typical example is to test whether a drug has better efficacy than *at least one* of multiple drugs.

In the rule MULT-V, we have two datasets y_1 and y_2 respectively obtained by sampling from two populations x_1 and x_2 (that may have statistical relevance). Then we apply two separate hypothesis tests on y_1 in C_1 and on y_2 in C_2 to derive the *disjunctive* alternative hypothesis $\varphi_1 \vee \varphi_2$. We denote by α_1 and α_2 the p -values of these two tests when performed separately; i.e., $\mathbf{K}_{y_1}^{\alpha_1} \varphi_1$ and $\mathbf{K}_{y_2}^{\alpha_2} \varphi_2$ are satisfied.

However, the p -value when performing both the tests simultaneously (in $C_1 \parallel C_2$) is larger than α_1 and α_2 . This is the so-called *multiple comparison problem*. By applying the Bonferroni's method, the p -value in total is bounded above by $\alpha_1 + \alpha_2$, namely, $\mathbf{K}_{(y_1, y_2)}^{\leq \alpha_1 + \alpha_2} (\varphi_1 \vee \varphi_2)$ is satisfied. This reflects that BHL does *not* derive elementary mistakes (e.g., $\mathbf{K}_{(y_1, y_2)}^{\alpha_1} \varphi_1$) where the reported p -value α_1 is lower than the actual p -value in multiple comparison.

In contrast, the rule MULT- \wedge formalizes the tests for the *conjunctive* alternative hypothesis $\varphi_1 \wedge \varphi_2$, e.g., the program C_{drug} in Example 2, which tests whether a drug has better efficacy than *both* drugs. According to statistics, this does not make the p -value higher, i.e., the p -value is at most $\min(\alpha_1, \alpha_2)$. See the full version of this paper for details.

Finally, we obtain the soundness of BHL from the validity of the rules. See the full version for the proof.

Theorem 1 (Soundness). *Every derivable judgment is valid.*

7.3 Instantiation to Concrete Test Methods

The axiom schemas for hypothesis tests in Figure 3 are instantiated with concrete examples of tests as follows. We first show the case of the two-tailed Z -test (Example 1).

Example 5 (Z -test). The axiom for the Z -test comparing means of two populations with datasets y_1, y_2 of sample sizes $\text{size}(y_1), \text{size}(y_2)$ and a known variance σ^2 is given by instantiating TWO-T with the following parameters:

$$\begin{aligned}
P_b(\mu_{\text{pp1 } b}, \sigma^2) &\stackrel{\text{def}}{=} N(\mu_{\text{pp1 } b}, \sigma^2) \quad \text{for } b = 1, 2 \\
\varphi_0 &\stackrel{\text{def}}{=} (\mu_{\text{pp1 } 1} = \mu_{\text{pp1 } 2}) \quad \varphi_{\top} \stackrel{\text{def}}{=} (\mu_{\text{pp1 } 1} \neq \mu_{\text{pp1 } 2}) \\
\leq^{(\top)} &= \{(r_1, r_2) \in \mathbb{R} \times \mathbb{R} \mid |r_1| > |r_2|\} \\
t(y_1, y_2) &= \frac{\text{mean}(y_1) - \text{mean}(y_2)}{\sigma \sqrt{1/\text{size}(y_1) + 1/\text{size}(y_2)}} \\
D_{t_0, \varphi_0} &= N(0, 1) \quad (\text{the standard normal distribution}).
\end{aligned}$$

Next we show the instantiation to the classical likelihood ratio test with a simple null hypothesis $\xi = \xi_0$ and a simple alternative hypothesis $\xi = \xi_1$, namely, in the setting of the Neyman-Pearson lemma.

Example 6 (Likelihood ratio test). The goal of (the simplest version of) the likelihood ratio test is to determine which of two candidate distributions $D_p, D_q \in \mathbb{D}\mathbb{R}$ is better to fit a dataset $y = (y_1, \dots, y_n)$ of sample size n . The test can be reformulated with the statistical model $x \approx P(\xi)$ defined by

$$P(\xi) = \begin{cases} D_q & \text{if } \xi = \xi_0 \\ D_p & \text{if } \xi = \xi_1 \end{cases}$$

and the following null and alternative hypotheses:

$$\varphi_0 \stackrel{\text{def}}{=} (\xi = \xi_0) \quad \varphi_L \stackrel{\text{def}}{=} (\xi = \xi_1).$$

The the likelihood function L for this test is given by

$$L(y|\xi_0) = \prod_{i=1}^n q(y_i) \quad L(y|\xi_1) = \prod_{i=1}^n p(y_i).$$

Then, the likelihood ratio $t(y)$ is given by

$$t(y) = \frac{L(y|\xi_0)}{L(y|\xi_1)} = \frac{\prod_{i=1}^n q(y_i)}{\prod_{i=1}^n p(y_i)}$$

where p and q are the density functions of D_p and D_q respectively. In the likelihood ratio test, for a given α and a threshold k such that $\Pr_{d_1, \dots, d_n \sim D_q} [t((d_1, \dots, d_n)) \leq k] \leq \alpha$, if we have $t(y) \leq k$, the likelihood $L(y|\xi_0)$ is too small to accept the distribution D_q . We then conclude that the other candidate D_p is better to fit y (thus this test is lower-tailed).

Conversely, the p -value of this test is defined by

$$\Pr_{d_1, \dots, d_n \sim D_q} [t((d_1, \dots, d_n)) \leq t(y)]. \quad (6)$$

The axiom for the likelihood ratio test is given by instantiating LOW-T with the above $P(\xi)$, φ_0 , φ_L , $t(y)$ and

$$\begin{aligned} \preceq^{(L)} &= \{(r_1, r_2) \in \mathbb{R} \times \mathbb{R} \mid r_1 \leq r_2\} \\ D_{t_\theta, \varphi_0} &= \frac{\prod_{i=1}^n q(D_q)}{\prod_{i=1}^n p(D_q)} \end{aligned}$$

where $p(D_q)$ and $q(D_q)$ are the probability distributions respectively defined by $p(D_q)(A) \stackrel{\text{def}}{=} D_q(p^{-1}(A))$ and $q(D_q)(A) \stackrel{\text{def}}{=} D_q(q^{-1}(A))$ for any measurable subset $A \subseteq \mathbb{R}$. Intuitively, $p(D_q)$ and $q(D_q)$ represent the probability distributions of $p(y)$ and $q(y)$ when y is sampled from D_q . By instantiating the p -value $\llbracket f_{A_{\varphi_0}^{(s)}}(y) \rrbracket$ in (5), we obtain (6).

Bayesian hypothesis test is given in an analogous way.

Example 7 (Bayesian hypothesis test). Consider the Bayesian likelihood ratio test with a dataset y of sample size n , prior distributions $D_{p'}, D_{q'} \in \mathbb{D}\mathbb{R}$ with density functions p' and q' , and posterior distributions $D_{p(z)}, D_{q(z)} \in \mathbb{D}\mathbb{R}$ with density functions $p(-|z)$ and $q(-|z)$.

The goal of this test is to determine whether the dataset y is sampled from $D_{q(z)}$ when z follows $D_{q'}$. The null hypothesis is that y is sampled from $D_{p(z)}$ when z follows $D_{p'}$.

A statistical model $x \approx P(\xi)$ of this test is defined as follows: First, we introduce the statistical models $z \approx P_0(\xi)$ and $x \approx P_1(\xi, z)$ of prior and posterior distributions by

$$P_0(\xi) = \begin{cases} D_{q'} & \text{if } \xi = \xi_0 \\ D_{p'} & \text{if } \xi = \xi_1 \end{cases} \quad P_1(\xi, z) = \begin{cases} D_{q(z)} & \text{if } \xi = \xi_0 \\ D_{p(z)} & \text{if } \xi = \xi_1 \end{cases}.$$

Next, for each $\xi = \xi_0, \xi_1$, we define the probability measure $P(\xi) \in \mathbb{D}\mathbb{R}$ by for any measurable subset $A \subseteq \mathbb{R}$,

$$(P(\xi))(A) \stackrel{\text{def}}{=} \int_{\mathbb{R}} h_{\xi, A} d\mu_\xi$$

where $\mu_\xi = P_0(\xi)$ and $h_{\xi, A}: \mathbb{R} \rightarrow \mathbb{R}$ is a measurable function defined by $h_{\xi, A}(z) = (P_1(\xi, z))(A)$.

The axiom for this test is given by instantiating LOW-T with the above model P and the following parameters:

$$\begin{aligned} \varphi_0 &\stackrel{\text{def}}{=} (\xi = \xi_0) \quad \varphi_L \stackrel{\text{def}}{=} (\xi = \xi_1) \\ \preceq^{(L)} &= \{(r_1, r_2) \in \mathbb{R} \times \mathbb{R} \mid r_1 \leq r_2\} \\ t(y) &= \frac{\int q'(z) \prod_{i=1}^n q(y_i|z) dz}{\int p'(z) \prod_{i=1}^n p(y_i|z) dz} \\ D_{t_\theta, \varphi_0} &= \frac{\int q'(z) \prod_{i=1}^n q(D_{q(z)}|z) dz}{\int p'(z) \prod_{i=1}^n p(D_{q(z)}|z) dz}. \end{aligned}$$

Unlike the (classical) likelihood ratio test, the *Bayes factor* $t(y)$ is the ratio of the following marginal likelihoods:

$$\begin{aligned} L(y|\xi_0) &= \int q'(z) \prod_{i=1}^n q(y_i|z) dz \\ L(y|\xi_1) &= \int p'(z) \prod_{i=1}^n p(y_i|z) dz. \end{aligned}$$

7.4 Reasoning About Multiple Comparison

We illustrate how BHL reasons about the multiple comparison in Example 2, where the derivation of (3) guarantees that the hypothesis tests are applied appropriately in C_{drug} in (1). (Details are shown in the Supplementary Material.)

In the derivation, we obtain the following judgments:

$$\begin{aligned} \Gamma &\vdash \{\psi_{\text{pre}}\} C_{12} \{\psi_{\text{pre}} \wedge \mathbf{K}_y^{\alpha_{12}} \varphi_{12}\} \\ \Gamma &\vdash \{\psi_{\text{pre}} \wedge \mathbf{K}_y^{\alpha_{12}} \varphi_{12} \wedge \alpha_{12} \leq 0.05\} C_{13} \{\varphi_{12}^{\text{bel}} \rightarrow \varphi^{\text{bel}}\} \\ \Gamma &\vdash \{\psi_{\text{pre}} \wedge \mathbf{K}_y^{\alpha_{12}} \varphi_{12} \wedge \alpha_{12} > 0.05\} \text{skip} \{\varphi_{12}^{\text{bel}} \rightarrow \varphi^{\text{bel}}\} \end{aligned}$$

where $\varphi_{12}^{\text{bel}} \stackrel{\text{def}}{=} \mathbf{K}_y^{\leq 0.05} \varphi_{12}$ and $\varphi^{\text{bel}} \stackrel{\text{def}}{=} \mathbf{K}_y^{\leq \alpha} (\varphi_{12} \wedge \varphi_{13})$. The second judgment is derived by the two-tailed test axiom, MULT- \wedge , and CONSEQ. The last judgment is derived from $\Gamma \models (\psi_{\text{pre}} \wedge \mathbf{K}_y^{\alpha_{12}} \varphi_{12} \wedge \alpha_{12} > 0.05) \rightarrow \neg \varphi_{12}^{\text{bel}}$ and $\Gamma \models \neg \varphi_{12}^{\text{bel}} \rightarrow (\varphi_{12}^{\text{bel}} \rightarrow \varphi^{\text{bel}})$ and CONSEQ. Applying IF to the last two judgments, we have $\Gamma \vdash \{\psi_{\text{pre}} \wedge \mathbf{K}_y^{\alpha_{12}} \varphi_{12}\}$ if $\alpha_{12} \leq 0.05$ then C_{13} else skip $\{\varphi_{12}^{\text{bel}} \rightarrow \varphi^{\text{bel}}\}$; composing it with the first judgment by applying SEQ, we have the judgment in (3).

In contrast, the program $C_{12} \parallel C_{13}$ in (4) shows the multiple comparison problem. Since the alternative hypothesis $\varphi_{12} \vee \varphi_{13}$ is disjunctive, we apply the rule MULT- \vee to drive the belief $\mathbf{K}_{(y, y)}^{\leq \alpha_{12} + \alpha_{13}} (\varphi_{12} \vee \varphi_{13})$, with a p -value (larger than α_{12} and α_{13}) at most $\alpha_{12} + \alpha_{13}$.

7.5 Reasoning About p -Value Hacking

We informally describe how our framework can be applied to reason about a program for *p-value hacking*, i.e., a scientifically malignant technique to obtain a low p -value. The following program c_{pHack} is an example of p -value hacking that conducts a hypothesis test on different datasets y_1 and y_2 , and ignores the experiment showing a higher p -value to report only a lower p -value:

$$\begin{aligned} (\alpha_1 := f_{A_{\varphi_0}^{(\top)}}(y_1)) \parallel (\alpha_2 := f_{A_{\varphi_0}^{(\top)}}(y_2)); \\ \text{if } \alpha_1 < \alpha_2 \text{ then } \alpha := \alpha_1 \text{ else } \alpha := \alpha_2 \end{aligned}$$

We write φ_{alt} for the alternative hypothesis of this test.

For the reported p -value α to be an actual p -value, $\mathbf{K}_{(y_1, y_2)}^{\leq \alpha} \varphi_{\text{alt}}$ needs to hold as a postcondition of c_{pHack} . Then

$(\alpha_1 < \alpha_2 \rightarrow \mathbf{K}_{(y_1, y_2)}^{\leq \alpha_1} \varphi_{\text{alt}}) \wedge (\alpha_1 \geq \alpha_2 \rightarrow \mathbf{K}_{(y_1, y_2)}^{\leq \alpha_2} \varphi_{\text{alt}})$ must hold at the end of the first line of c_{pHack} due to the rules UPDVAR and IF. However, this is not implied by the postconditions of MULT- \vee or MULT- \wedge . Since $\mathbf{K}_{(y_1, y_2)}^{\leq \alpha_1 + \alpha_2} (\varphi_1 \vee \varphi_2)$ is a postcondition of c_{pHack} , the total p -value $\alpha_1 + \alpha_2$ should be reported without ignoring any experiments.

8 Discussion

In this section, we provide a whole picture of the justification of statistical belief inside and outside BHL.

A statistical belief derived in a program relies on the following three issues: (i) the validity of hypothesis test methods themselves, (ii) the satisfaction of the empirical conditions required for the hypothesis tests, and (iii) the appropriate usage of hypothesis tests in the program. In our framework, these are respectively addressed by (a) the validity of axioms and rules, (b) the (manual) confirmation of the preconditions in a judgment, and (c) the proof for the judgment.

8.1 Validity of Hypothesis Test Methods

The validity of hypothesis test methods is not ensured by mathematics alone. The philosophy of statistics has a long history of argument on the proper interpretation of hypothesis testing. One of the most notable examples is the argument between the frequentist and the Bayesian statistics, which still has many issues to be discussed (Sober 2008).

We also remark that statistical methods occasionally involve some approximation of numerical values. However, we may not always confirm the validity of approximation rigorously, i.e., whether the approximation is valid to the specific situation we apply the statistical methods.

For these reasons, we do not attempt to formalize the “justification” for hypothesis test methods within BHL, and left them for future work. Instead, we define axiom schemas for hypothesis tests that are commonly used in practice and explained in textbooks, e.g., (Kanji 2006). Then we focus on the logical aspects of the appropriate usage of hypothesis tests, which has been a long-standing, practical concern but has not been formalized using a symbolic logic before.

One of the advantages of this approach is that we do not adhere to a specific philosophy of statistics, but can model both the frequentist and the Bayesian statistics by introducing an axiom/rule corresponding to each hypothesis test.

8.2 Clarification of Empirical Conditions

The hypothesis test methods usually assume some empirical conditions on the unknown population from which the dataset is sampled. Typically, many parametric tests require that the population follows a normal distribution. For instance, the Z -test in Example 5 assumes that the population follows a normal distribution with known variance, but this cannot be rigorously confirmed or justified in general.

In some cases, such conditions on the unknown population may be confirmed approximately or partially by some exploratory observations on the sampled data and by prior knowledge of some properties on the population (outside the statistical inference). As far as we know, there has been no

general method for justifying such empirical conditions rigorously. Thus, the formal justification of those conditions themselves would require further research in statistics.

In the present paper, the empirical conditions on the unknown population remain to be assumptions from the viewpoint of formal logic. Hence, we describe empirical conditions as the preconditions of a judgment in BHL.

Explicit specification of the preconditions would be useful for non-experts to prevent errors in the choice of statistical methods. Furthermore, when we formalize empirical science in future work, it would be crucial to clarify the empirical conditions that justify scientific conclusions.

8.3 Epistemic Aspects of Statistical Inference

One of our contributions is to show that epistemic logic is useful to formalize statistical inference. Although the outcome of a hypothesis test is a *knowledge* determined by the test action, it may form a false *belief*; i.e., a rejected null hypothesis may be true, and a retained one may be false. Hence the formalization of statistical inference deals with both truth and beliefs, for which epistemic logic is suitable.

The key to formalizing statistical beliefs is to introduce a Kripke semantics with a possible world w_0 where a null hypothesis is true (Section 5.2). This possible world w_0 may not be the real world where we actually apply the hypothesis test on an observed dataset.

In the Kripke model, a transition between states is used to model the update of statistical beliefs by a hypothesis test. Since the world records the executions of all tests, BHL does not allow for hiding some tests to manipulate the statistics (e.g., p -value in multiple comparison and p -value hacking).

Furthermore, the choice of two-tailed or one-tailed tests requires describing a prior belief using the possibility modality \mathbf{P} . Without this modality, we cannot express the belief that both lower-tail and upper-tail are possible before applying the test, since this belief may not be true.

Without this Kripke semantics, the formalization of hypothesis testing would deal with only the purely mathematical propositions satisfied in the possible world w_0 where a null hypothesis is true, hence could not reason about the appropriate usage of hypothesis tests in the real world.

9 Conclusion

In this work, we proposed a new approach to formalizing and reasoning about statistical inference in programs. Specifically, we introduced belief Hoare logic (BHL) for formalizing and checking the requirement for hypothesis tests to be employed appropriately. Then we showed that BHL is useful to reason about practical issues in statistics. We also discussed a whole picture of the justification of statistical inference. We emphasize that this is the first attempt to introduce a program logic for the appropriate application of hypothesis tests.

In ongoing and future work, we are extending our framework to other kinds of statistical methods. We plan to investigate the relative completeness of BHL and to develop a verification tool based on this framework. Another possible research would be to study justification logic for statistics.

Acknowledgments

The authors are supported by ERATO HASUO Metamathematics for Systems Design Project (No. JPMJER1603), JST. In particular, we thank Ichiro Hasuo for providing the opportunity for us to meet and collaborate in that project. Yusuke Kawamoto is supported by JST, PRESTO Grant Number JPMJPR2022, Japan, and by JSPS KAKENHI Grant Number 21K12028, Japan. Tetsuya Sato is supported by JSPS KAKENHI Grant Number 20K19775, Japan. Kohei Suenaga is supported by JST CREST Grant Number JPMJCR2012, Japan.

References

- Atkinson, E., and Carbin, M. 2020. Programming and reasoning with partial observability. *Proc. ACM Program. Lang.* 4(OOPSLA):200:1–200:28.
- Belle, V., and Levesque, H. J. 2015. ALLEGRO: belief-based programming in stochastic dynamical domains. In *Proc. IJCAI 2015*, 2762–2769. AAAI Press.
- Bretz, F.; Hothorn, T.; and Westfall, P. 2010. *Multiple Comparisons Using R*. Chapman and Hall/CRC.
- Burrows, M.; Abadi, M.; and Needham, R. M. 1990. A logic of authentication. *ACM Trans. Comput. Syst.* 8(1):18–36.
- den Hartog, J., and de Vink, E. P. 2002. Verifying probabilistic programs using a Hoare like logic. *Int. J. Found. Comput. Sci.* 13(3):315–340.
- Fagin, R.; Halpern, J.; Moses, Y.; and Vardi, M. 1995a. *Reasoning about Knowledge*. The MIT Press.
- Fagin, R.; Halpern, J. Y.; Moses, Y.; and Vardi, M. Y. 1995b. Knowledge-based programs. In *Proc. PODC 1995*, 153–163. ACM.
- Garcia, F. D.; Hasuo, I.; Pieters, W.; and van Rossum, P. 2005. Provable anonymity. In *Proc. of FMSE*, 63–72.
- Halpern, J. Y. 2003. *Reasoning about uncertainty*. The MIT press.
- Huber, F., and Schmidt-Petri, C. 2008. *Degrees of belief*, volume 342. Springer Science & Business Media.
- Kanji, G. K. 2006. *100 statistical tests*. Sage.
- Kawamoto, Y. 2019. Statistical epistemic logic. In *The Art of Modelling Computational Systems: A Journey from Logic and Concurrency to Security and Privacy*, volume 11760 of LNCS, 344–362. Springer.
- Kawamoto, Y. 2020. An epistemic approach to the formal specification of statistical machine learning. *Software and Systems Modeling* 20(2):293–310.
- Lang, T. A., and Altman, D. G. 2014. *Statistical Analyses and Methods in the Published Literature: The SAMPL Guidelines*. John Wiley & Sons, Ltd. chapter 25, 264–274.
- Laverny, N., and Lang, J. 2005. From knowledge-based programs to graded belief-based programs, part I: on-line reasoning*. *Synth.* 147(2):277–321.
- Nielson, H. R., and Nielson, F. 2007. *Semantics with Applications: An Appetizer (Undergraduate Topics in Computer Science)*. Berlin, Heidelberg: Springer-Verlag.
- Reynolds, J. C. 2002. Separation logic: A logic for shared mutable data structures. In *Proc. LICS 2002*, 55–74. IEEE Computer Society.
- Sober, E. 2008. *Evidence and evolution: The logic behind the science*. Cambridge University Press.
- Suenaga, K., and Hasuo, I. 2011. Programming with infinitesimals: A while-language for hybrid system modeling. In *Proc. ICALP 2011, Part II*, volume 6756 of LNCS, 392–403. Springer.
- Syverson, P. F., and Stubblebine, S. G. 1999. Group principals and the formalization of anonymity. In *World Congress on Formal Methods (1)*, 814–833.
- von Wright, G. H. 1951. *An Essay in Modal Logic*. Amsterdam: North-Holland Pub. Co.
- Wasserstein, R. L., and Lazar, N. A. 2016. The ASA statement on p-values: Context, process, and purpose. *The American Statistician* 70(2):129–133.
- Winskel, G. 1993. *The Formal Semantics of Programming Languages—An Introduction*. The MIT Press.