Intro 0000 Consistency

Implications 00000

Consistent Estimators and the Argument from Inductive Risk

Corey Dethier

Samuel C. Fletcher

Minnesota Center for Philosophy of Science University of Minnesota, Twin Cities

April 7, 2023



NSF Grant No. 2042366

Consistency

Implications 00000

The classical debate

Statistics

Intro

0000



"As is well known, the acceptance or rejection of such a hypothesis presupposes that a certain level of significance or level of confidence or critical region be selected." (Rudner 1953, 3)

"the activity proper to the scientist is the assignment of probabilities (with respect to currently available evidence) to the hypotheses which, on the usual view, he simply accepts or rejects." (Jeffrey 1956, 237)





 Strictly speaking, Rudner's argument only applies to a (strict) NP testing framework. Fisherian/hybrid approaches don't employ acceptance levels or regions.



- Strictly speaking, Rudner's argument only applies to a (strict) NP testing framework. Fisherian/hybrid approaches don't employ acceptance levels or regions.
- Modern rejoinders to Jeffrey—e.g., Douglas (2000), Steele (2013)—have focused on other ways that values can enter into the testing process or the need to communicate results.



- Strictly speaking, Rudner's argument only applies to a (strict) NP testing framework. Fisherian/hybrid approaches don't employ acceptance levels or regions.
- Modern rejoinders to Jeffrey—e.g., Douglas (2000), Steele (2013)—have focused on other ways that values can enter into the testing process or the need to communicate results.
- It is thus an open question whether the scientist qua (classical) statistician must make value judgments.

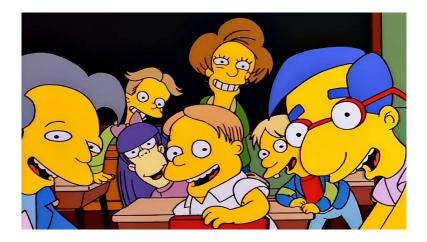
Stat 000

Intro 00●0

tistics

Consistency 00000 Implication 00000 References

Obligatory Simpsons reference



◆□▶ ◆□▶ ◆臣▶ ◆臣▶ ─臣 ─ のへで





Focusing on the use of estimators:

• The choice of estimators is open to *exactly* analogous conditions as the choice of acceptance level.

Focusing on the use of estimators:

- The choice of estimators is open to *exactly* analogous conditions as the choice of acceptance level.
- All permissible estimators will eventually converge on the truth a feature shared by Rudner's original example.

Focusing on the use of estimators:

- The choice of estimators is open to *exactly* analogous conditions as the choice of acceptance level.
- All permissible estimators will eventually converge on the truth a feature shared by Rudner's original example.
- Which calls into question whether Rudner's example (a) can or
 (b) should play the role often assigned to it in the literature.

Statistics 000000

Consistency

Implications

Estimators and inductive risk

An **estimator** is a kind of test statistic: it's a rule for deriving a "best guess" (the "estimate") for a quantity of interest ("estimand") from the sample. E.g.:

If the quantity of interest is the population mean, the sample mean, median, or mode serve as the estimator.

ション ふゆ アメリア メリア しょうくしゃ

Similarly for higher moments of the distribution, such as variance or skew.

An **estimator** is a kind of test statistic: it's a rule for deriving a "best guess" (the "estimate") for a quantity of interest ("estimand") from the sample. E.g.:

If the quantity of interest is the population mean, the sample mean, median, or mode serve as the estimator.

ション ふゆ アメリア メリア しょうくしゃ

Similarly for higher moments of the distribution, such as variance or skew.



Generally, prefer the estimator with the smallest (expected) loss.

Data 3.23 1.92 4.28 2.57 4.20 2.97 3.87 2.60 2.62 3.72 2.72 2.72 2.60 2.

Loss function

Estimator Estimate

Data are estimates for equilibrium climate sensitivity taken from Tokarska et al. (2020). Notably, there's discussion in the literature about how to estimate *variance* in this case (see Annan and Hargreaves 2011; Dethier 2022).



Generally, prefer the estimator with the smallest (expected) loss.

 Data
 3.23
 1.92
 4.28
 2.57
 4.20
 2.97
 3.87
 2.60

 2.62
 3.72
 2.72
 2.72
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.

Loss functionEstimatorEstimateMean absolute error $(\Sigma_x(|Y-x|)/n)$ Median2.97

Data are estimates for equilibrium climate sensitivity taken from Tokarska et al. (2020). Notably, there's discussion in the literature about how to estimate *variance* in this case (see Annan and Hargreaves 2011; Dethier 2022).

Statistics

Generally, prefer the estimator with the smallest (expected) loss.

Consistency

 Data
 3.23
 1.92
 4.28
 2.57
 4.20
 2.97
 3.87
 2.60

 2.62
 3.72
 2.72
 2.72
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.

Loss function

Estimator Estimate

ション ふゆ アメリア メリア しょうくしゃ

Implications

References

 $\begin{array}{lll} \mbox{Mean absolute error } (\Sigma_x(|Y-x|)/n) & \mbox{Median} & 2.97 \\ \mbox{Mean squared error } (\Sigma_x(Y-x)^2/n) & \mbox{Mean} & 3.15 \\ \end{array}$

Data are estimates for equilibrium climate sensitivity taken from Tokarska et al. (2020). Notably, there's discussion in the literature about how to estimate *variance* in this case (see Annan and Hargreaves 2011; Dethier 2022).

Statistics

Generally, prefer the estimator with the smallest (expected) loss.

Consistency

 Data
 3.23
 1.92
 4.28
 2.57
 4.20
 2.97
 3.87
 2.60

 2.62
 3.72
 2.72
 2.72
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.60
 2.

Loss function

Estimator Estimate

Implications

Mean absolute error $(\Sigma_x(Y-x)/n)$	Median	2.97
Mean squared error $(\Sigma_x(Y-x)^2/n)$	Mean	3.15
Mean quartic error $(\Sigma_x(Y-x)^4/n)$	unnamed	3.19

Data are estimates for equilibrium climate sensitivity taken from Tokarska et al. (2020). Notably, there's discussion in the literature about how to estimate *variance* in this case (see Annan and Hargreaves 2011; Dethier 2022).

References



These different functions represent different attitudes towards error.

MAE: errors of 1, 2, and 4 count for 1, 2, and 4.MSE: errors of 1, 2, and 4 count for 1, 4, and 16.MQE: errors of 1, 2, and 4 count for 1, 16, and 256.

Or: which estimator you should use depends on how you weight small vs. large errors.



To carry out the most basic hypothesis test:

• Calculate the value of the test statistic *Z*:

 $\frac{\text{estimator of the mean} - \text{hypothesized mean}}{\text{standard deviation}/\sqrt{\text{sample size}}}$

- Calculate the probability of observing a value greater than the result using a normal distribution.
- The resulting probability is the *p*-value, which quantifies how well the evidence "fits" the hypothesis.



To carry out the most basic hypothesis test:

• Calculate the value of the test statistic *Z*:

 $\frac{\text{estimator of the mean} - \text{hypothesized mean}}{\text{standard deviation}/\sqrt{\text{sample size}}}$

- Calculate the probability of observing a value greater than the result using a normal distribution.
- The resulting probability is the *p*-value, which quantifies how well the evidence "fits" the hypothesis.

The examples just given generalizes in a fairly trivial way to (almost all of?) the rest of inferential statistics, either by way of

- relying on estimators—i.e., loss functions—to select an estimate (as here); or
- relying explicitly or implicitly on loss functions in other ways (e.g., least-squares algorithms in regression).

ション ふゆ アメリア メリア しょうくしゃ

So the same value-laden choices about loss functions are ubiquitous in classical statistics.

Intro Statistics Consistency Consistency

Rudner's argument:

- (P1) Scientists qua scientists must choose an acceptance level.
- (P2) The choice of acceptance level requires weighting different errors.

- (P3) Weighting different errors requires making value judgments.
- \therefore (C) Scientists qua scientists must make value judgments.

 Intro
 Statistics
 Consistency
 Implications

 Octob
 Octob
 Octob
 Octob
 Octob

 The parallel with Rudner's argument

References

▲ロト ▲周ト ▲ヨト ▲ヨト - ヨ - の々ぐ

Our argument:

- (P1) Scientists qua classical statisticians must choose estimators.
- (P2) The choice of estimator requires weighting different errors.
- (P3) Weighting different errors requires making value judgments.
- \therefore (C) Scientists qua classical statisticians must make value judgments.

Statistics Consistency 000000 00000

Implications

Consistent estimators and inductive risk



In most contexts, any permissible estimator is consistent.

The main definition of consistency is that the estimator X_n "converges in probability" with the target θ :

$$\forall \epsilon > 0, \lim_{n \to \infty} \Pr(|X_n - \theta| > \epsilon) = 0$$

Or, more simply:

 $\lim_{n\to\infty} X_n(\theta) = \theta$



In the limit, the differences between the *estimates* generated by permissible estimators (almost surely) disappear.

Or: no matter how substantial the divergence of values, sufficient data will (almost surely) wash values out of the estimate.

▲ロ ▶ ▲ □ ▶ ▲ □ ▶ ▲ □ ▶ ▲ □ ▶ ● ○ ○ ○

And it has the same effect on Rudner's original example!

Consistency and acceptance levels

Let z_{α} indicate the "critical value": if the hypothesis is true, the probability of observing $|Z| > z_{\alpha}$ for a given sample size is α .

The accept and reject criteria are then:

Accept: if $|z| \leq z_{\alpha}$ Reject: $|z| > z_{\alpha}$

Claim: If the underlying estimator is consistent, as $n \to \infty$, the probability of accepting a true hypothesis goes to 1 and the probability of accepting a false one goes to 0 for any $z_{\alpha} \in (0, \infty)$.

Technically, you *could* define the acceptance region without reference to a test statistic. The resulting tests are not "completely consistent" (Andrews 1986) – which obviates the point of using a consistent estimator in the test.



Recall that Z is defined as follows:

$$Z = \frac{\text{estimator of the mean} - \text{hypothesized mean}}{\text{standard deviation}/\sqrt{\text{sample size}}} = \frac{X_n - \theta_0}{\sigma/\sqrt{n}}$$

▲□▶ ▲圖▶ ▲ 臣▶ ▲ 臣▶ ― 臣 … のへで

Consistency (recall): $\forall \epsilon > 0$, $\lim_{n \to \infty} Pr(|X_n - \theta| > \epsilon) = 0$.



Recall that Z is defined as follows:

$$Z = \frac{\text{estimator of the mean} - \text{hypothesized mean}}{\text{standard deviation}/\sqrt{\text{sample size}}} = \frac{X_n - \theta_0}{\sigma/\sqrt{n}}$$

Consistency (recall): $\forall \epsilon > 0$, $\lim_{n \to \infty} Pr(|X_n - \theta| > \epsilon) = 0$.

If the hypothesis is true $(\theta_0 = \theta)$: $\forall z_{\alpha} > 0$, $\lim_{n \to \infty} Pr(|Z| \leq z_{\alpha}) = 1$

▲ロト ▲冊 ト ▲ 臣 ト ▲ 臣 ト ○ ○ ○ ○ ○



Recall that Z is defined as follows:

$$Z = \frac{\text{estimator of the mean} - \text{hypothesized mean}}{\text{standard deviation}/\sqrt{\text{sample size}}} = \frac{X_n - \theta_0}{\sigma/\sqrt{n}}$$

Consistency (recall): $\forall \epsilon > 0$, $\lim_{n \to \infty} Pr(|X_n - \theta| > \epsilon) = 0$.

If the hypothesis is **true** $(\theta_0 = \theta)$: $\forall z_{\alpha} > 0$, $\lim_{n \to \infty} Pr(|Z| \leq z_{\alpha}) = 1$

If the hypothesis is false $(\theta_0 \neq \theta)$: $\forall z_{\alpha} > 0$, $\lim_{n \to \infty} Pr(|Z| \leq z_{\alpha}) = 0$

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ ・ う く つ ・



▲ロト ▲周ト ▲ヨト ▲ヨト - ヨ - の々ぐ

Different choices of how to prioritize errors w.r.t. either

- Ioss functions / estimators
- acceptance levels / critical values

converge on the same conclusion / decision in the limit.



Different choices of how to prioritize errors w.r.t. either

- Ioss functions / estimators
- 2 acceptance levels / critical values

converge on the same conclusion / decision in the limit.

When we focus on specific cases rather than on the full range of consistent estimators, we can be more specific.

E.g., when i.i.d. sampling from a normal distribution, the sample variance s^2 converges on the population variance σ^2 as $n \rightarrow 30$ for most practical purposes.

Statistics 000000 Consistency 00000 Implications

What have we learned about inductive risk?



"I will argue that non-epistemic values are a required part of the internal aspects of scientific reasoning for cases where inductive risk includes risk of nonepistemic consequences." (Douglas 2000, 559)

Implications

Consistency

It's inaccurate to describe Douglas's *argument* as a "revival, reiteration, or rediscovery" of Rudner's (Havstad 2022, 309).



▲ロ ▶ ▲ □ ▶ ▲ □ ▶ ▲ □ ▶ ▲ □ ▶ ● ○ ○ ○

References

Rudner's *example* of balancing false positives and negatives continues of play an extremely important role in the literature.

Particularly in motivating the claim that values can *legitimately* influence scientific reasoning in other contexts / domains.

We see this not just of Douglas (2000) but also in Biddle (2013), Brown (2013), Elliott (2022), Frank (2019), John (2015), Parker (2014), Plutynski (2017), Steel (2010), Steele (2013), Stegenga (2017), and Wilholt (2009).



These statistical examples have the special feature that the influence of values will eventually wash out.

In other words:

• Values can only influence conclusions in short/medium term

Conclusions should (with probability) converge over time (regardless of actually *reaching* the infinite limit)



These statistical examples have the special feature that the influence of values will eventually wash out.

In other words:

- **1** Values can only influence conclusions in short/medium term
- Conclusions should (with probability) converge over time (regardless of actually *reaching* the infinite limit)

It is an open question whether any of the extensions mentioned on the last page have similar features.

Insofar as arguments for extending legitimate values-influence depend on the analogy to the Rudner example, those arguments bear re-examination.

More positively, our discussion suggests at least one path forward on "the new demarcation question" (Holman and Wilholt 2022):

Namely, value-influence is legitimate when? if? only if? they wash out in a manner analogous to what we find in the Rudner example.

Intro	Statistics	Consistency	Implications	References
0000	000000	00000	0000●	
The end				

Thank you!!

<□▶ <□▶ <□▶ < □▶ < □▶ < □▶ = - つへぐ

- Andrews, Donald W. K. (1986). Complete Consistency: A Testing Analogue of Estimator Consistency. The Review of Economic Studies 53.2: 263–69. DOI: 10.2307/2297650.
- Annan, James D. and Julia C. Hargreaves (2011). Understanding the CMIP3 Model Ensemble. *Journal of Climate* 24: 4529–38.
- Biddle, Justin B. (2013). State of the Field: Transient Underdetermination and Values in Science. Studies in History and Philosophy of Science Part A 44.1: 124–33. DOI: 10.1016/j.shpsa.2012.09.003.
- Brown, Matt J. (2013). Values in Science Beyond Underdetermination and Inductive Risk. *Philosophy of Science* 80.5: 829–39.
- Dethier, Corey (2022). When is an Ensemble Like a Sample? 'Model-Based' Inferences in Climate Modeling. Synthese 200.52: 1–20. DOI: 10.1007/s11229-022-03477-5.
- Douglas, Heather (2000). Inductive Risk and Values in Science. Philosophy of Science 67.4: 559–79. DOI: 10.1086/392855.
- Elliott, Kevin C. (2022). Values in Science. Cambridge: Cambridge University Press.
- Frank, David M. (2019). Ethics of the Scientist qua Policy Advisor: Inductive Risk, Uncertainty, and Catastrophe in Climate Economics. Synthese 196: 3123–38. DOI: 10.1007/s11229-017-1617-3.
- Havstad, Joyce C. (2022). Sensational Science, Archaic Hominin Genetics, and Amplified Inductive Risk. Canadian Journal of Philosophy 52.3: 295–320. DOI: 10.1017/can.2021.15.
- Holman, Bennett and Torsten Wilholt (2022). The New Demarcation Problem. Studies in History and Philosophy of Science Part A 91: 211–20. DOI: 10.1016/j.shpsa.2021.11.011.
- Jeffrey, Richard (1956). Valuation and Acceptance of Scientific Hypotheses. Philosophy of Science 23.3: 237–46. DOI: 10.1086/287489.
- John, Stephen (2015). Inductive Risk and the Contexts of Communication. Synthese 192.1: 79–96. DOI: 10.1007/s11229-014-0554-7.
- Parker, Wendy S. (2014). Values and Uncertainties in Climate Prediction, Revisited. Studies in History and Philosophy of Science Part A 46: 24–30.
- Plutynski, Anya (2017). Safe or Sorry? Cancer Screening and Inductive Risk. In: Exploring Inductive Risk: Case Studies of Values in Science. Ed. by Kevin C. Elliott and Ted Richards. Oxford: Oxford University Press: 149–70.
- Rudner, Richard (1953). The Scientist qua Scientist makes Value Judgments. *Philosophy of Science* 20.1: 1–6. DOI: 10.1086/287231.
- Steel, Daniel (2010). Epistemic Values and the Argument from Inductive Risk. Philosophy of Science 77.1: 14–34. DOI: 10.1086/650206.

Intro	Statistics	Consistency	Implications	References
0000	000000	00000	00000	

- Steele, Katie (2013). The Scientist qua Policy Advisor makes Value Judgments. *Philosophy of Science* 79.5: 893–904. DOI: 10.1086/667842.
- Stegenga, Jacob (2017). Drug Regulation and the Inductive Risk Calculus. In: Exploring Inductive Risk: Case Studies of Values in Science. Ed. by Kevin C. Elliott and Ted Richards. Oxford: Oxford University Press: 17–36.
- Tokarska, Katarzyna B. et al. (2020). Past Warming Trend Constrains Future Warming in CMIP6 Models. Science Advances 6.12: 1–13.
- Wilholt, Torsten (2009). Bias and Values in Scientific Research. Studies in History and Philosophy of Science Part A 40.1: 92–101. DOI: 10.1016/j.shpsa.2008.12.005.