The Likelihood Principle

Conor Mayo-Wilson

Philosophy of Statistics June 30th, 2014

REVIEW

Last Month:

- Models of experiments: Classical/Frequentist vs. Bayesian
- Classical/frequentist hypothesis tests and criticisms

TODAY

Today

- Why some take objections to hypothesis testing motivate use of interval estimates (esp. confidence intervals).
- The Likelihood Principle: Why its proof attacks classical methodology in general, and not just hypothesis tests.

To review the differences between two traditions, here's an example.

Hypothesis Testing - Toy Example

• Suppose a coin factory produces two types of coins, one with bias $\frac{1}{4}$ and the other with bias $\frac{3}{4}$.

•
$$\Theta = \{\frac{1}{4}, \frac{3}{4}\}.$$

• Suppose your **null-hypothesis** Θ_0 is that the coin has bias $\frac{1}{4}$.

Hypothesis Testing - Toy Example

Suppose you have flipped the coin ten-thousand times, and observed 2575 heads.

- $\Omega = \{0, 1\}^{10^4}$.
- $X_n: \Omega \to \mathbb{R}$ is the map

$$\omega = \langle \omega_1, \omega_2, \dots \omega_{10^4} \rangle \mapsto \omega_n.$$

• Here's what you know about the experimental outcome $\omega \in \Omega$:

$$\sum_{n\leq 10^4} X_n(\omega) = 2575$$

Objections to Classical Hypothesis Testing

Objection 1: "Statistically significant" does not mean practically significant.

EXAMPLE

Mayo and Spanos' Answer





Answer: Test the null $\mu = \mu_0$ vs. the alternative $\mu > \mu_1$, where μ_1 is a "practically significant" deviation.

VS.

Mayo and Spanos [2011]

EXAMPLE

Example:

- In the toy example, suppose, unbeknownst to you, the coin factory produces coins with bias $\frac{1}{4} + \epsilon$ for some very small $\epsilon > 0$.
- You wouldn't care if you had one such coin or a coin with bias precisely ¹/₄.
- You are flipping one such coin.
- Then, under any reasonable test, the probability of rejecting the null hypothesis $\Theta_0 = \{\frac{1}{4}\}$ at any significance level $\alpha > 0$ approaches one as the number of flips approaches infinity.

Example:

- In the toy example, the probability of obtaining 2575 or more heads is about 4%, which is significant at the .05 level.
- However, unless you own a casino, you probably do not care if your coin has bias .2575 or bias .25.

Objection 2: With a large enough sample, the smallest, insignificant deviation from a point null will be rejected.

Mayo and Spanos' Answer





Answer:

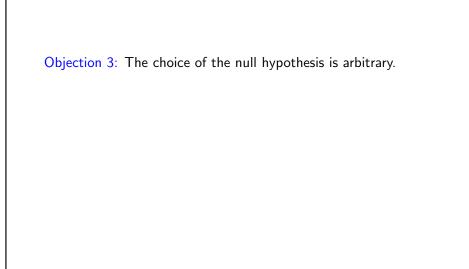
- If the null hypothesis is false, it should be rejected. That doesn't entail accepting the alternative hypothesis.
- My Note: The "change the hypothesis" response also could be used here: Test the null μ ∈ (μ₀ − ε, μ₀ + ε).

Mayo and Spanos [2011]

EXAMPLE

Example:

- If your null hypothesis is Θ₀ = {1/4}, then the probability of obtaining 2575 or more heads is about 4%, which is significant at the .05 level. So you reject the null.
- If your null hypothesis is $\Theta_0 = \{\frac{3}{4}\}$, of obtaining 2575 or fewer heads is approximately zero, which is significant at the .05 level. So you reject the null.



Mayo and Spanos' Answer





Answer: The choice of the null hypothesis can only cause erroneous inferences if one fallaciously infers that evidence again the null hypothesis is evidence for some specific alternative.

VS.

Mayo and Spanos [2011, pp. 175]

Objection 4: Statistically insignificant results are taken as evidence that the null hypothesis is true.

Mayo and Spanos' Answer





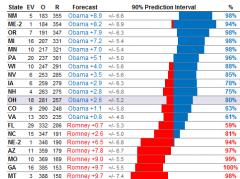
Answer:

- That's a mistake of practitioners; not of the methodology.
- In certain settings, we can quantify how "large" a variation from the null hypothesis is compatible with failure to reject. So "statistically insignificant" results can have some informative value.

Mayo and Spanos [2011]

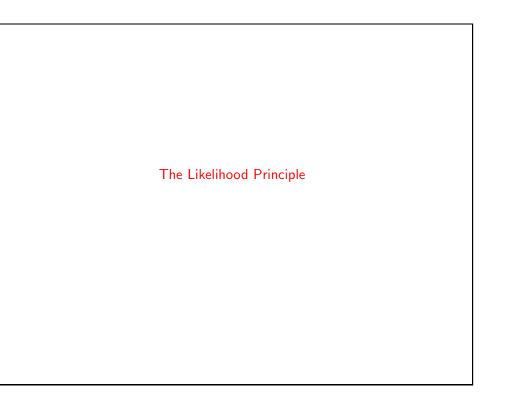
INTERVAL ESTIMATION

Competitive State Summary



Alternative Solution: Abandon hypothesis testing. Use interval estimates.

- Interval estimates quantify the effect size, not just statistical significance.
- They give a range of plausible values, not simply the verdict "reject" or "retain."



Likelihood Principle: Two data sequences (from possibly different experiments) with the same likelihood function provide the same evidence.

Comparing Evidence

- Sometimes we may wish to compare the evidence provided by two experiments.
- For example, *ceteris paribus*, if we wish to say that one study provides better evidence of climate change than another.
- Can we formalize the idea of "stronger", "similar", and "identical" evidence?

Likelihood Principle: Two data sequences (from possibly different experiments) with the same likelihood function provide the same evidence.

FORMALIZING EVIDENCE

- Birnbaum [1962] is completely agnostic about how to represent evidence.
- Let Exp_Θ be the class of experiments whose underlying set of states of the world are equal to Θ.
- Pieces of evidence about Θ are represented by elements of a set $\mathrm{EVIDENCE}_{\Theta}.$
- There is a function $Ev : Exp_{\Theta} \times \mathbb{R}^n \to \text{EVIDENCE}_{\Theta}$ with
 - Input: An experiment $\mathcal E$ and a data sequence $x \in \mathbb{R}^n$ from that experiment.
 - **Output:** A piece of evidence $Ev(\mathcal{E}, x) \in EVIDENCE_{\Theta}$.

Likelihood Principle: Two data sequences (from possibly different experiments) with the same likelihood function provide the same evidence.

Two Experiments

Consider two classical experiments with the same set of states of the worlds $\boldsymbol{\Theta}:$

- $\mathcal{E}_1 = \langle \Theta, \Omega_1, \{ P_{\theta} \}_{\theta \in \Theta}, \langle X_n \rangle_{n \in \mathbb{N}} \rangle$
- $\mathcal{E}_2 = \langle \Theta, \Omega_2, \{Q_\theta\}_{\theta \in \Theta}, \langle Y_n \rangle_{n \in \mathbb{N}} \rangle$

Two Experiments

Example:

- Experiment 1: Flip a coin of unknown bias twelve times.
- **Experiment 2:** Flip a coin of unknown bias until three heads are observed.

TWO EXPERIMENTS

The experiments may differ in any number of ways:

- They may have different sets of experimental outcomes • I.e., $\Omega_1 \neq \Omega_2$
- Different sequences of observations may be possible.
 - I.e., $\langle X_n \rangle_{n \in \mathbb{N}} \neq \langle Y_n \rangle_{n \in \mathbb{N}}$
- The probabilities of the same data sequence may differ, even if the world θ is held fixed.
 - I.e., $P_{ heta} \neq Q_{ heta}$

Likelihood Principle: Two data sequences (from possibly different experiments) with the same likelihood function provide the same evidence.

Two Experiments

Example:

- Experiment 1: Flip a coin of unknown bias twelve times.
- Experiment 2: Flip a coin of unknown bias until three heads are observed.

SAME LIKELIHOOD FUNCTION

- Let x = ⟨X₁(ω),...X_m(ω)⟩ be an outcome of the first experiment.
- Let y = ⟨Y₁(ω),...Y_n(ω)⟩ be an outcome of the second experiment.
- Say x and y determine the same likelihood function if there is a constant $c_{x,y}$ such that

$$p_{\theta}(x) = c_{x,y} \cdot q_{\theta}(y)$$

for all $\theta \in \Theta$

TWO EXPERIMENTS: SAME DATA:

Outcome: Let x = y be a data sequence consisting of 12 flips, with 9 tails, and a heads on the last flip.

Two Experiments: Same Likelihood Function

Example:

- Experiment 1: Flip a coin of unknown bias twelve times.
 - Regardless of the bias $\theta \in \Theta$ of the coin:

$$P_{\theta}(x) = {12 \choose 3} \cdot \theta^3 (1-\theta)^9$$

- Experiment 2: Flip a coin of unknown bias until three heads are observed.
 - Regardless of the bias $\theta \in \Theta$ of the coin:

$$Q_ heta(x) = inom{11}{2} \cdot heta^3 (1- heta)^9$$

Note: Here, $c_{x,x} = \frac{\binom{12}{3}}{\binom{11}{2}} = \frac{220}{55} = 4.$

TWO EXPERIMENTS: SAME EVIDENCE

- Some find this intuitive: your data is your data!
- Others don't: The fact that you could have flipped to coin for a longer or shorter time in the second experiment should affect the inferences you draw.

TWO EXPERIMENTS: SAME EVIDENCE

Example: According to the Likelihood Principle, the sequence of coin flips provides the same evidence about the bias of the coin, regardless of which of the two experiments is conducted.

USELESS COUNTERFACTUALS

Here are some cases in which classical techniques make use of counterfactual dependencies that seem irrelevant.

LP AND CLASSICAL STATISTICS

Upshot: Classical statistical methods violate the LP. Why?

SAME LIKELIHOOD FUNCTION, DIFFERENT VERDICTS

Experiment 1: Recall if the coin has bias θ , then the probability of k many heads is:

$$\binom{12}{k} \cdot \theta^k (1-\theta)^{12-k}$$

When $\theta = \frac{1}{2}$:

$$heta^k (1- heta)^{12-k} = (rac{1}{2})^{12}$$

So the probability of three or fewer heads is:

$$P_{\theta}(x) = \left(\begin{pmatrix} 12\\3 \end{pmatrix} + \begin{pmatrix} 12\\2 \end{pmatrix} + \begin{pmatrix} 12\\1 \end{pmatrix} + \begin{pmatrix} 12\\0 \end{pmatrix} \right) \cdot \frac{1}{2^{12}}$$

Two Experiments: Same Likelihood Function

Example: Suppose two experimenters test the null hypothesis $\Theta_0 = \{\frac{1}{2}\}$ at the .05 level in the following two experiments.

- Experiment 1: Flip a coin of unknown bias twelve times.
- **Experiment 2:** Flip a coin of unknown bias until three heads are observed.

SAME LIKELIHOOD FUNCTION, DIFFERENT VERDICTS

Experiment 1: If you do the calculation, you'll find:

$$P_{\theta}(x) \approx 7.3\%.$$

So you would not reject the null hypothesis at the .05 level.

SAME LIKELIHOOD FUNCTION, DIFFERENT VERDICTS

Experiment 2: Flip a coin of unknown bias until three heads are observed.

The probability of needing at least k many throws is:

$$\binom{k}{2} \cdot \theta^3 (1-\theta)^{k-3}$$

When $\theta = \frac{1}{2}$, this probability is equal to the following:

$$\binom{k}{2} \cdot \frac{1}{2^k}$$

So the probability of needing 12 or more tosses is:

$$1 - \left(\binom{11}{2} \frac{1}{2^{11}} + \binom{10}{2} \frac{1}{2^{10}} + \ldots + \binom{2}{2} \frac{1}{2^3} \right)$$

SAME LIKELIHOOD FUNCTION, DIFFERENT VERDICTS

Experiment 2: If you do the calculation, you'll find this probability is about 3.3%. So you would reject the null hypothesis at the .05 level.

SAME LIKELIHOOD FUNCTION, DIFFERENT VERDICTS

Moral: Classical hypothesis testing violates the LP.

Because of the close relationship between confidence intervals and hypothesis tests, so will confidence intervals.

So what?

Birnbaum's Theorem: LP is entailed by two principles nearly universally endorse by classical statisticians: (i) the conditionality principle, and (ii) the sufficiency principle.

MIXED EXPERIMENTS

- Suppose you want to measure the charge of an electron.
- There are two experiments that you might conduct, and you think both are equally reliable.
- You flip a coin to decide which experiment to conduct.
- Randomly choosing among a collection of experiments is called a mixed experiment.

Conditionality Principle

- Making this precise requires giving a formal definition of a mixed experiment, which I won't do.
- Important:
 - $\bullet\,$ The randomizing device should not provide additional information about Θ :
 - I.e., The probability that the randomizing device takes particular values is not a function of $\theta \in \Theta$.
 - E.g., Don't flip the coin of unknown bias to choose which of the two experiments above to conduct. The extra coin flip gives you additional information!
 - Birnbaum is likewise agnostic about how to formalize the notion of evidence here.

Conditionality Principle

Conditionality Principle: If you conduct a mixed experiment in which you observe x after conducting the component experiment \mathcal{E} , then your evidence is the same as if you had observed x after conducting \mathcal{E} without randomizing among

SUFFICIENCY PRINCIPLE

Sufficiency: All evidence from an experiment concerning Θ in the experiment is available in a sufficient statistic.

SUFFICIENCY PRINCIPLE

- Note: A sufficient statistic may provide no evidence for any θ ∈ Θ. For instance, suppose you flip a coin to determine whether or not it will rain one year from today in Munich.
- **Major Lemma:** The set of likelihoods is a sufficient statistic, i.e., the function:

 $T: x \mapsto \langle P_{\theta}(x) \rangle_{\theta \in \Theta}$

Proof of Birnbaum's Theorem:

- Consider the mixed experiment \mathcal{E}^* obtained by flipping a coin with bias $\frac{1}{c_{x,y}+1}$ to determine whether to conduct \mathcal{E}_1 or \mathcal{E}_2 .
- Consider the probability of observing $\langle Heads, x \rangle$ and $\langle Tails, y \rangle$ in the mixed experiment.

Proof of Birnbaum's Theorem:

- Let x and y be data from experiments *E*₁ and *E*₂ determining the same likelihood function, with constant c_{x,y} ≥ 0.
- We want to show $Ev(\mathcal{E}_1, x) = Ev(\mathcal{E}_2, y)$.

Proof of Birnbaum's Theorem: In the mixed experiment \mathcal{E}^* :

$$p_{\theta}^{*}(\langle \text{Heads}, x \rangle) = \frac{1}{c_{x,y} + 1} \cdot p_{\theta}(x)$$
$$= \frac{1}{c_{x,y} + 1} \cdot c_{x,y} \cdot q_{\theta}(y)$$
$$= (1 - \frac{1}{c_{x,y} + 1}) \cdot q_{\theta}(y)$$
$$= p *_{\theta} (\langle \text{Tails}, y \rangle)$$

Proof of Birnbaum's Theorem:

• ${\rm SUFFICIENCY}$ and the lemma \Rightarrow

$$Ev(\mathcal{E}^*, \langle H, x \rangle) = Ev(\mathcal{E}^*, \langle T, y \rangle)$$

 $\bullet \ {\rm CONDITIONALITY} \Rightarrow$

$$Ev(\mathcal{E}^*, \langle H, x \rangle) = Ev(\mathcal{E}_1, x)$$

• CONDITIONALITY \Rightarrow

$$Ev(\mathcal{E}^*, \langle T, y \rangle) = Ev(\mathcal{E}_2, y)$$

• Transitivity of equality then entails

$$Ev(\mathcal{E}_1, x) = Ev(\mathcal{E}_2, y)$$

which is what we wanted to show.

References I

Birnbaum, A. (1962). On the foundations of statistical inference. *Journal* of the American Statistical Association, 57(298):269–306.

Mayo, D. and Spanos, A. (2011). Error statistics. In Forster, M. and Bandyopadhyay, P. S., editors, *Philosophy of statistics*, number 7 in Handbook of the Philosophy of Science, pages 153—198.

