

آمار بسامدگرا به مثابه نظریه استنباط استقرایی*

دیبورا جی. مایو^{*}، دی. آر. کاکس^{**}

ترجمهٔ محمدرضا مشکانی

۱. آمار و فلسفهٔ استقرایی

۱.۱ فلسفهٔ آمار چیست؟

این پژوهش‌های بسیار عام که با مجادلات دیرپایی فلسفهٔ علم درهم تنیده‌اند به تبیین این موضوع که چرا رشتة آمار، خواه به تصریح یا به تلویح، گرایش به ورود به حیطهٔ فلسفه دارد، کمک می‌کنند. برخی حتی ممکن است آمار را به نوعی «فلسفهٔ علم کاربردی» (فیشر [۱۰]، کمپتون [۱۳])، و نظریهٔ آمار را به نوعی «فلسفهٔ استنباط استقرایی کاربردی» تلقی کنند. چنان‌که لین [۱۵] تأکید کرده است، نیمن کار خود را نه تنها کاری در آمار بلکه کاری در فلسفهٔ استقرایی نیز تلقی می‌کرد. پرسشی اساسی که در «فلسفهٔ استقرایی» هم در آمار و هم در فلسفهٔ مطروح است این است که: ماهیت و نقش مفهومهای احتمالاتی، روشها، و مدلها در انجام استنباط، هنگامی که با داده‌های محدود، عدم حتمیت و خطای سروکار داریم، چیست؟

در این فرضی که برای شرکت در جلسهٔ راجع به فلسفهٔ آمار در «دومین سمپوزیوم لین» به دست آورده‌ایم، توصیهٔ نیمن ([۲۲]، ص ۱۷) را به عنوان مبنای بحث خود برمی‌گزینیم و نظریهٔ آمار را اساساً «نظریهٔ بسامدگرای استنباط استقرایی» تلقی می‌کنیم. آنگاه این پرسش پیش می‌آید که چه تصور (یا تصویراتی) از استنباط استقرایی این گزینش را مجاز می‌سازد. برای بررسی اینکه آیا این روایت، تنها یا حتی رضایت‌بخش‌ترین روایت از استنباط استقرایی است یا نه، جالب است بررسی شود که چه مقدار پیشرفت به‌سوی روایتی از استنباط استقرایی، در مقابل رفتار استقرایی، ممکن است از آمار بسامدگرا (با عطف توجه به آزمون‌کردن و روش‌های مربوط) عاید شود. این روشها اغلب برای هدفهای استنباطی، یعنی کسب آگاهی از جنبه‌های سازوکار زیربنایی مولد داده‌ها، به‌کار می‌روند، و در صورتی که شفافیت بیشتری در مورد نتایج استنباط خطاها فرضی در استنباط وجود داشته باشد، می‌توان از سیاری از سردرگمی‌ها و انتقادها (مثلًاً در این مورد که آیا باید نزخ خطاهای تعدیل شوند یا نه و چرا باید تعدیل شوند) جلوگیری کرد.

با پذیرش نظرات فیشر ([۱۰]، دیدگاه لین [۱۵] دربارهٔ نیمن، و پویر [۲۶] راجع به استقرا به عنوان یک‌شیوهٔ موضوع، نقش آزمونهای معنی‌داری را در پرکردن شکافهای استقرایی موجود در استنباط استقرایی فرضی سنتی

مبانی فلسفی آمار را می‌توان به مثابهٔ بررسی مسائلهای شناخت‌شناسانه، مفهومی، و منطقی دربارهٔ کاربرد و تفسیر روش‌های آماری، به‌معنای کلی، تلقی کرد. همچون دیگر حوزه‌های فلسفهٔ علم، پژوهش در علم آمار تا حد زیادی بدون دغدغه دربارهٔ «مبانی فلسفی» به‌پیش می‌رود. با وجود این، حتی در کار حرفه‌ای آمار، بحث و جدلها دربارهٔ رویکردهای متفاوت به تحلیل آماری ممکن است بر مسائل کلی ماهیت استنباط آماری-استقرایی اثر بگذارند و از آنها متأثر شوند، و بدین ترتیب با مباحث بنیادی و فلسفی سروکار داشته باشند. حتی کسانی که عمدتاً با کاربردها سروکار دارند اغلب علاقه‌مند به شناسایی آن اصلهای کلی‌اند که زیربنای شیوه‌های مورد استفاده و توجیه‌کننده آنها هستند و کاربران به‌دلایل کمایش عملی به ارزش آنها واقف شده‌اند. دست‌کم در یک سطح از تحلیل، آمارشناسان و فیلسوفان علم بسیاری پژوهش‌های مشابه را مطرح می‌کنند.

• چه چیزی باید مشاهده شود و از داده‌های حاصل چه چیزهایی را به طور موجه می‌توان استنباط کرد؟

• داده‌ها تا چه حد مدلی را تأیید می‌کنند یا به آن می‌برازند؟

• آزمون خوب چیست؟

• آیا شکست در رد فرض H شاهدی بر «تأیید» H است؟

• چگونه می‌توان تعیین کرد که یک بی‌هنگاری ظاهری، واقعی است یا نه؟

• چگونه می‌توان تقصیر یک بی‌هنگاری را به درستی به عامل آن نسبت داد؟

• اگر نگاه‌کردن به داده‌ها بر فرضی که باید امتحان شود تأثیر بگذارد، آیا این امر به رابطهٔ بین داده‌ها و فرض ربطی دارد؟

• چگونه می‌توان رابطه‌های کاذب را از نظمهای واقعی متمایز ساخت؟

• چگونه می‌توان یک تبیین علی و یک فرض را توجیه و آزمون کرد؟

• چگونه می‌توان اختلاف بین داده‌های موجود و ادعاهای نظری را به شیوهٔ

قابل اعتمادی رفع کرد؟

آزمون شده است، و این کاملاً متمایز است با اینکه درجه‌ای از احتمال، باور، یا تأیید به آن اعطای شده باشد؛ این همان «ظرافت‌بخشی به استقراء» است. آن دو همچنین کلاً این دیدگاه مشترک را داشتند که برای آنکه داده‌ها فرضی مانند H را «تأیید» یا «پشتیبانی» کنند، آن فرض باید در معرض آزمونی قرار گیرد و در صورت نادرست بودن با احتمال یا توان زیاد رد شود. اما برخلاف ارتباط نزدیک اندیشه‌ها، به‌ظاهر هیچ ارجاعی به پویر در نوشه‌های نیمن وجود ندارد (لین [۱۶]، ص ۳) و ارجاعهای پویر به نیمن اندک و بهمندر ذی‌ربطاند. افزون بر آن، چون پویر انکار می‌کرد که هیچ‌گونه ادعای استقرایی قابل توجیه باشد، فلسفه‌اش وی را مجبور ساخت انکار کند که حتی روشنی که از آن پشتیبانی می‌کرد (حدس و ابطال) قابل اعتماد است. پویر روش ساخت که اگرچه H ممکن است درست باشد، وی پشتیبانی را حداکثر به عنوان گزارشی از عملکرد گذشتة H تلقی می‌کند: این امر هیچ ادعایی را درباره قابلیت اعتماد آن در کاربردهای آینده ایجاب نمی‌کند. بر عکس، ویزگی بارز آمار بسامدی آن است که بتوان احتمال آن را که اگر فرضی وجود دارد باشد آزمون آن را رد کند، ارزیابی و کنترل کرد. این احتمالها از فرمول بندی فرازند مولد داده‌ها برحسب یک مدل آماری ناشی می‌شوند.

نیمن در سراسر تحقیقات خود بر اهمیت یک مدل احتمالاتی برای سامانه تحت بررسی تأکید می‌کند و آمار بسامدی را به مثابه مدل بندی پدیده پایداری بسامد های نسبی تابیج «آزمایه‌های» مکرر توصیف می‌کند، در عین حال که می‌بیند امکانات دیگری برای مدل بندی پدیده‌های روان‌شناسخی مرتبط با شدتهای باور، یا مرتبط با آمادگی برای شرط‌بندی روی مبلغهای مشخص، و غیره وجود دارند، در این موارد از کارناب [۲]، دوفینتی [۸]، و سه‌ویج [۲۷] نقل قول می‌کند. بهویژه نیمن از دیدگاه استنباط «بسامدی» مورد قبول کارناب به‌خاطر نادیده گرفتن نقش اساسی مدل تصادفی پدیده مورد بررسی انتقاد می‌کند. دستاورد آماری کارناب [۲] در ارتباط با فلسفه استقرایی همان دستاورد کنیز [۱۴]، و دارای تأثیری مستقیم بر کاربردهای آماری است (جفریز [۱۲]).

۳. استقراء و استنباط فرضی-قیاسی

هر چند ممکن است گمان کنیم که «استنباط فرضی-قیاسی» استقراء را «ظرافت می‌بخشد»، در واقع استنباطهای استقرایی اند که در سراسر آزمونهای تجربی انجام می‌شوند. می‌توان ملاحظه کرد که ایده‌های آزمون آماری این شکافهای استقرایی را بر می‌کنند: اگر فرض قطعی باشد می‌توانیم تابعی ذی‌ربط از داده‌ها بیاییم که مقدارش: (i) ویزگی ذی‌ربط تحت آزمون را نمایش دهد و (ii) بتوان آن را با آن فرض پیشگویی کرد. مقدار تابع را حساب می‌کنیم و ملاحظه می‌کنیم که داده‌ها با پیشگویی موافق اند یا نیستند. اگر داده‌ها با پیشگویی در تعارض باشند، آنگاه یا فرض در معرض خطاست یا نوعی عامل کمیکی یا پس زمینه‌ای دیگر را می‌توان در این بی‌هنگاری مقصراً شناخت (مسئله دوئم^۱).

ملحوظات آماری به دو طریق وارد بحث می‌شوند. اگر H یک فرض آماری باشد، آنگاه معمولاً هیچ برآمدی اکیداً آن را نقض نمی‌کند. در نظر گرفتن داده‌ها به صورت ناسازگار با H صرفاً به این دلیل که آنها بسیار نامحتمل اند، مشکلات عمده‌ای در بردارد؛ تک تک برآمدها در صورتی که به تفصیل تشریح

1. Duhem's problem

مورد بحث قرار می‌دهیم. هدف ما شناسایی یک اصل اساسی درباره شواهد است که بهوسیله آن بتوان احتمال خطاهای فرضی را برای انجام استنباط استقرایی از داده‌های خاص بهکار بست، و نیز بررسی این نکته که این اصل چگونه می‌تواند موارد زیر را هدایت و توجیه کند: (الف) کاربردها و تقسیرهای متفاوت از سطوحهای استنباط آماری در آزمون انواع مختلف فرضهای صفر، و (ب) چه موقع و چرا «اثرهای گزینش» باید در آزمون آماری داده‌وابسته به حساب آورده شوند.

۲.۱ نقش احتمال در استقراء بسامدگر

ویزگی تعریف‌کننده استنباط استقرایی آن است که مقدمات گزاره‌های مبتنی بر شواهد ممکن است درست باشند ولی نتیجه استنباط شده بدون تناقضی منطقی نادرست باشد: یعنی نتیجه «فراتر از شواهد است». احتمال در بدست آوردن این گونه استنباطهای فراتر از شواهد به‌طور طبیعی ظاهر می‌شود. اما بیش از یک راه برای استفاده از احتمال وجود دارد. پیرسون ([۲۴]، ص ۲۲۸): دو سنت فلسفی متمایز برای بهکار بستن احتمال در استنباط را جمع‌بندی کرده است:

«در نظر یکی از دو مکتب، درجه اطمینان به یک گزاره، یعنی کیمی که با ماهیت و میزان شواهد تغییر می‌کند، مفهومی پایه‌ای است که باید به آن مقیاس عددی نسبت داد.» مکتب دیگر به ارتباط بین آگاهی از بسامد نسبی رج‌دادن رده‌ای خاص از پیشامدها در سلسه‌ای از تکرارها با زندگی عادی و با بسیاری از شاخه‌های علوم توجه دارد، و بر آن است که «از طریق پیوند احتمال با بسامد نسبی است که احتمال صریح‌ترین معنا را در ذهن انسان پیدا می‌کند.»

در استقرای بسامدی، بهر صورتی که باشد، احتمال به روای دوم بهکار گرفته می‌شود. مثلاً در آزمون معنی‌داری به احتمال متول می‌شویم تا نسبت مواردی را که در آنها فرض صفری مانند H در یک نمونه‌گیری مکرر درازمدت مفروض رد خواهد شد، یعنی احتمال خطأ، را مشخص سازیم. این تفاوت در نقش احتمال، متناظر با تفاوت در شکل استنباطی است که مقتضی انگاشته می‌شود: کاربرد قبلی احتمال معمولاً با این دیدگاه که توصیف احتمالاتی استقراء شامل کمی‌سازی درجه پشتیبانی یا تأیید موجود در ادعاهای فرضهایست، رابطه تنگاتگ دارد.

برخی طرفداران رویکرد بسامدگر، که ترجیح می‌دهند اصطلاح «رفتار استقرایی» را بهکار بینند، با این نظر موافق‌اند که نقش احتمال را در آمار بسامدی توصیف کنند. در این حال، استدلال‌کننده استقرایی «تصمیم می‌گیرد که نتیجه را استنباط کند»، و احتمال مخاطره خطای مربوط را کمی می‌سازد. این اندیشه که یک نقش احتمال در علم عبارت است از مشخص کردن «مخاطره‌آمیزی» یا تأمین شواهد یا شدت آزمونهایی که فرضها در معرض آنها قرار داده می‌شوند، یادآور فلسفه کارل پویر [۲۶] است. بهویژه، لین ([۱۶]، ص ۳۲) شباهت زمانی و مفهومی اندیشه‌های پویر و نیمن را درباره «ظرافت بخشی» به موضوع استقراء با درنظر گرفتن یک فرایند آزمون فرض به جای استدلال استقرایی، خاطر نشان کرده است.

این مطلب درست است که پویر و نیمن رویکردهایی مشابه مبتنی بر این اندیشه داشته‌اند که می‌توانیم از فرضی سخن بگوییم که به یک معنا خوب

این توزیع احتمال را باید نمایشی غالباً انتزاعی و یقیناً آرمانی شده از فرایند زیربنایی مولد داده‌ها تلقی کرد. سپس فرضی درباره توزیع احتمال داریم، که گاهی به آن فرض تحت آزمون می‌گویند اما اغلب آن را فرض صفر می‌نامند و با H_0 نشان می‌دهند. بعداً شماری از انواع کاملاً متفاوت فرضهای صفر را بیان خواهیم کرد، اما فعلاً بین آنهایی که گاهی ساده خوانده می‌شوند، به این معنی که توزیع Y را به طور کامل (اصولاً به طور عددی) مشخص می‌سازند و آنهایی که گاهی مرکب نامیده می‌شوند، به این معنی که برخی جنبه‌های معین را کاملاً مشخص می‌کنند و دیگر جنبه‌ها را نامشخص باقی می‌گذارند، تمایز قابل می‌شوند.

مثالی که از بسیاری جهات ابتدایی ترین مثال، هر چند مثالی تکراری، است آن است که Y متشکل از n مؤلفه مستقل و هم توزیع است که به طور نرمال با میانگین μ و انحراف معیار امکاناً مجھول σ توزیع شده‌اند. اگر مقدار σ معلوم و برابر با، مثلاً σ ، باشد و فرض صفر آن باشد که μ برابر با μ ، یعنی ثابتی مفروض است، فرضی ساده حاصل می‌شود. یک فرض مرکب در همان چارچوب بالا ممکن است شامل σ مجھول باشد و باز هم مقدار μ را مشخص کند.

توجه کنید که در این فرمول بندی لازم است جنبه‌ای مجھول از توزیع، نوعاً یک یا چند پارامتر مجھول، دقیقاً مشخص شود. فرض اینکه، مثلاً $\mu \leq \mu$ ، یک فرمول بندی پذیرفتی برای فرض صفر در آزمون فیشری نیست؛ در حالی که این صورت کلی تر فرض صفر در فرمول بندی‌های نیمن-پیرسون مجاز است.

هدف فوری عبارت است از آزمودن سازگاری داده‌های خاص تحت تحلیل با H_0 ، از جنبه‌ای که قرار است مشخص شود. برای انجام این کار تابعی مانند $t(y) = t$ از داده‌ها می‌باشیم، که آماره آزمون نامیده می‌شود، به‌گونه‌ای که

- هر چه مقدار t بزرگ‌تر باشد داده‌ها با H_0 ناسازگارترند؛
- وقتی H_0 درست باشد، متغیر تصادفی متناظر $(Y = t)$ دارای یک توزیع احتمال (به طور عددی) معلوم است.

این دو شرط با شرطهای قطعی متناظر هم‌تازند. برای ارزیابی اینکه یک ناهمخوانی واقعی با H_0 (یا انحراف بازآفریدنی از H_0) وجود دارد یا نه، به‌اصطلاح p -مقدار متناظر با هر t را به‌صورت

$$p = p(t) = P(T \geq t; H_0)$$

تعییف می‌کنیم، که شاخصی از ناهمخوانی با H_0 از جنبه مورد آزمون تلقی می‌شود. دست‌کم در فرمول بندی آغازین، فرضهای مقابل در لایه‌های زیرین پنهان‌اند اما از جنبه احتمالاتی به طور صریح فرمول بندی نمی‌شوند؛ همچنین پیش‌پیش هیچ پریشی درباره تعیین یک مقدار آستانه‌ای از پیش تعیین شده و «رد» H_0 اگر و تنها اگر $\alpha \leq p$ ، مطرح نیست. افزون بر آن، برای توجیه آزمونها صرفاً به رفتار درازمدت متول نخواهیم شد بلکه یک ملاک عقلی استنباطی یا مبتنی بر شواهد را شناسایی خواهیم کرد. اینک در این باره بیشتر توضیح می‌دهیم.

شوند ممکن است احتمال وقوع خیلی کمی داشته باشند. اما اساساً طبق نظر پوپر ([۲۶]، صص ۸۶، ۲۰۳) موضوع این است که آیا برامدی که امکاناً بی‌亨جارت باشد نمایشگر اثری نظام‌مند و بازآفریدنی است یا نه.

تمرکز نظریه پوپر بر ابطال به عنوان هدف آزمونها، و ابطال به مثابه ملاک تعريف‌کننده برای یک نظریه یا فرض علمی، به روشنی و قویاً یادآور تفکر فیشر است. در حالی که شواهدی با تأثیر مستقیم عملًا غایب باشد، دیدگاه‌های پوپر با این بیان فیشر ([۹]، ص ۱۶) موافق است که می‌توان گفت هر آزمایش تنها برای آن است که شناس رکدن فرض صفر را به واقعیات اعطای کند. اما، چون موضع پوپر این است که هیچ‌گاه دلیلی برای استنباط درباره قابلیت اعتماد نخواهیم داشت، وی قبول ندارد که هرگز بتوانیم دلایلی برای استنباط راجع به انحرافهای بازآفریدنی داشته باشیم.

مزیت چارچوب آماری نوین آن است که احتمال‌ها از تعريف مدل احتمال برای نمایش پدیده مورد نظر ناشی می‌شوند. اگر پوپر از ایده‌های آزمون آماری که تقریباً در همان زمانها ارائه شدند بهره گرفته بود، امکان داشت توصیف خود از ابطال را مستدل سازد.

موضوع دوم مربوط به این مسأله است که وقتی داده‌ها با پیشگویی «توازن» دارند چگونه باید استدل کرد. استدل بر این اساس که H_0 مستلزم داده‌های y است، و اینکه y مشاهده شده است، برای رسیدن به این استنباط که H_0 درست است، البته از لحاظ قیاسی نامعتبر است. مسأله‌ای اساسی در توصیف استقرایی عبارت است از اینکه بتوان این استنباط H_0 را به معنای تضمین کرد. اما، مسأله کلاسیک، حتی در حالتهای قطعی، آن است که بسیاری فرضهای رقیب (برخی خواهند گفت به تعداد بینهایت از آنها) نیز y را پیشگویی خواهند کرد، و از این‌رو مانند H_0 قابل قبول خواهند بود. برای آنکه آزمونی گواه آور [فرام آورنده شواهد] باشد، پیشگویی از روی H_0 باید چیزی باشد که اگر H_0 نادرست و رقیبان مهم H_0 درست باشند، در همان حال این پیشگویی به معنایی پسیار شگفت‌انگیز باشد و به‌آسانی قابل توجیه نباشد. اینک بحث می‌کنیم که چگونه شکافهای موجود در آزمون استقرایی را می‌توان با یک نوع شیوه آماری خاص، یعنی آزمون معنی‌داری، پر کرد.

۲. آزمونهای معنی‌داری آماری

اگرچه آزمون معنی‌داری آماری طی بیش از ۵۰ سال با مباحثاتی همراه بوده و در بالاتر سوئتفاهم‌های موجود در نوشتگان آماری فرو رفته است، شماری از ویژگی‌های اساسی دیدگاه استقرایی بسامدی را که ما در نظر داریم، به شکل ساده‌ای نشان می‌دهد. برای مثال، رک. [۲۱] و [۱۶]. ما بحث را تا حد امکان با عناصر اصلی آزمون معنی‌داری آغاز می‌کنیم: روایت ما از این موضوع قویاً مرتبط، اتا از جهاتی متفاوت، با هر دو رویکرد فیشری و نیمن-پیرسونی، دست‌کم در فرمول بندی معمولی آنهاست.

۱.۲ تذکرات کلی و تعريف

فرض می‌کنیم داده‌هایی تجربی در دست داریم که به طور دسته جمعی با y نمایانه شده‌اند و آنها را مقدارهای مشاهده شده متغیر تصادفی Y درنظر می‌گیریم. داده‌های y را تنها تا آنجا ممکن است درباره توزیع احتمال Y به‌گونه‌ای که با مدل آماری ذی‌ربط تعریف شده، فراهم می‌سازند.

کارگزارشی است که به نوع استدلالهای نابی که دانشمندان به قصد کسب دانش قابل اعتماد و درک یک رشتۀ علمی بنا می‌نهند نزدیک می‌شود. در میان این پیچیدگی، استدلال آزمون معنی‌داری مفهومی نسبتاً سرراست از ارزیابی شواهد ناهمخوان با H . را در زمینه‌ای آماری بازسازی تاباند، مفهومی که شاید پوپر در ذهن خود داشت اما فاقد ابزارهایی برای تحقق بخشیدن به آن بود. ایده اساسی در اینجا آن است که احتمالهای خطای را می‌توان برای ارزیابی «مخاطره‌آمیز بودن» پیشگویی‌هایی که لازم است H در آنها صدق کند به کار بست، این کار با سنجش قابلیت اعتماد آزمون در تعایین قاتل شدن بین اینکه فرایند واقعی مولد داده‌ها با فرایند توصیف شده در H همخوان است یا نه، انجام می‌شود. آگاهی از این ظرفیت گواه‌آوری امکان می‌دهد که تعیین کنیم آیا شواهد قوی مبنی بر ناهمخوانی وجود دارند یا نه. برای تعیین اینکه شواهدی H وجود دارند یا نه، استدلال برایه اصل بسامدگرایانه زیر بنا نهاده می‌شود.

شواهد بسامدگرایانه (ش ب) (i). داده‌های y شواهدی (قوی) علیه H اند، یعنی شواهدی (قوی) مبنی بر وجود اختلاف با H . اند اگر و تنها اگر هنگامی که H توصیفی درست از سازوکار مولد y باشد، آنگاه با احتمال زیاد، این سازوکار نتیجه‌ای بدده که کمتر از آنچه به وسیله y نشان داده می‌شود، ناهمخوان باشد.

یک فرع ش ب آن است که y شواهدی (قوی) علیه H نیست اگر احتمال نتیجه‌ای ناهمخوان تر خیلی کم نباشد، حتی اگر H درست باشد. یعنی، اگر احتمالی نسبتاً زیاد برای نتیجه‌ای ناهمخوان تر وجود داشته باشد، حتی اگر H درست باشد، آنگاه H از جنبه مورد آزمون با y همخوانی دارد.

مطلوبی که قدری بحث‌انگیزتر است تفسیر شکست در یافتن p -مقداری کوچک است؛ اما برایه صورت فوق از ش ب می‌توان شرح و تفسیری رسا برای آن عرضه کرد.

۳.۲ شکست و تأیید

مشکل درنظرگرفتن مقداری متوسط از p به مثبته شاهدی به نفع H ، آن است که حتی اگر رقیبان H که جدا با H نقاوت دارند درست باشند، همخوانی بین H و y ممکن است رخ دهد. این موضوع بهویژه وقتی حاد است که اندازه داده‌ها محدود باشد. اما، گاهی می‌توانیم شواهدی بهنفع H بیاییم، که به مثبته ادعایی درک می‌شود مبنی براینکه اختلافی، عیی، یا خطایی در بین نیست، و این کار را می‌توانیم به وسیله آزمونهایی انجام دهیم که با احتمال زیاد، در صورت وجود اختلاف، آن را نشان می‌دادند. همان‌قدر که نام نیمن با فنون تصمیم‌گونه خودکار همراه است، دست‌کم در عمل، هم او و هم پیرسون، انتخاب مقتضی احتمال خطاهای را بازتاب‌دهنده بافتار ویژه مورد نظر تلقی می‌کردند. (نیمن [۲۳]، پیرسون [۲۴]).

در اینجا دو موضوع متفاوت مطرح است. یکی اینکه آیا مقداری خاص از p باید به مثبته آستانه‌ای در هر کاربرد به کار رود یا نه. این همان شیوه‌ای است که اگر نه در همه ولی در اکثر شرحهای رسمی از نظریه نیمن-پیرسون بیان می‌شود. موضوع دوم آن است که آیا کنترل نز خطاهای درازمدت، توجیهی برای آزمونهای بسامدگرایی را آیا توجیه نهایی آزمونها در نقش آنها

۲.۲ رفتار استقرایی در برابر استنباط استقرایی

این استدلال را می‌توان روایتی آماری از استدلال موسوم به برهان خلف درنظر گرفت که در منطق قیاسی معتبر است. در این روش، انکار یک فرض H از ترکیب گزاره H مستلزم E است با این اطلاع که نادرست است، استنباط می‌شود. چون اگر H درست می‌بود احتمال زیاد ($p - 1$) وجود می‌داشت که نتیجه‌ای کمتر معنی دارخ دهد، می‌توانیم انتخاب p -مقدارهای کوچک را که درست حساب شده باشند، به عنوان شواهدی علیه H توجیه کنیم. چرا؟ به دو دلیل عده:

اولاً چنین قاعده‌ای باعث می‌شود که در درازمدت، وقتی H درست باشد، نرخهای خطای (یعنی، رکوردن‌های نادرست) کوچک باشند، که استدلالی رفتارگرایانه است. همان‌گونه با نگرش مبتنی بر ارزیابی خطای H به آمار، می‌توانیم به هر مقدار خاص مثلاً p ، تفسیر فرضی زیر را نسبت دهیم: فرض کنید که قرار است داده‌ها را تنها شواهد قاطع علیه H قلمداد کنیم. در این صورت، در تکرارهای فرضی، H به نسبت درازمدت p از مواردی که در آنها واقعاً درست است، رد خواهد شد. اما، آگاهی از این احتمالهای خطای فرضی را می‌توان برای حمایت از توجیهی متمایز درنظر گرفت.

بدین‌گونه است که چنین قاعده‌ای راهی را برای تعیین اینکه مجموعه‌ای از داده‌های خاص شواهدی از ناهمخوانی با H است یا نه، پیش پایی می‌گذارد.

بهویژه، یک p -مقدار کوچک، به شرط آنکه درست حساب شود، شواهدی از اختلاف با H را از حیث جنبه مورد امتحان به دست می‌دهد، در حالی که p -مقداری که کوچک نباشد شواهدی از همخوانی یا سازگاری با H را در اختیار می‌گذارد (که در اینجا بین این حالت و شواهد مثبت به نفع H آن‌گونه که در زیر در بخش ۳-۲ بررسی می‌شود، باید فرق گذاشت). در کاربردها توجه مانوعاً به این است که آیا p در دامنه‌ای مانند $1 - p \geq 0.05$ واقع است یا نه، که می‌توان آن را به مثبته همخوانی معقول با H از جنبه مورد آزمون تلقی کرد، یا اینکه p نزدیک به عدددهای قراردادی از قبیل $0.01, 0.05, 0.1$ را است یا نه. شیوه معمول در اکثر کاربردها آن است که مقدار مشاهده شده p را به صورت تقریبی اعلام کنند. مقدار کوچک p بیانگر یکی از موارد زیر است:

(i) H نادرست است (اختلافی با H وجود دارد) یا (ii) اساس آزمون آماری معیوب است، اغلب از این نظر که خطاهای واقعی کم باور دشده‌اند، مثلاً به دلیل مفروضات نامعتبر استقلال، یا (iii) نقش شناس فوق العاده بوده است.

بخشی از هدف طرح پژوهشی خوب و انتخاب مناسب روش تحلیل این است که با کسب اطمینان از اینکه ارزیابیهای خطای مناسب‌اند، از شق (ii) بالا پرهیز شود.

به هیچ وجه قصد آن نیست که بگوییم آزمون معنی داری نواعاً تنها تحلیلی خواهد بود که گزارش می‌شود. در واقع، یک اصل بنیادی مفهوم یادگیری استقرایی که بیشترین سازگاری را با فلسفه بسامدگرایی دارد آن است که استنباط استقرایی لازم می‌نماید که استدلالها و استنباطهای قاطع از طریق به هم پیوستن چند خرد نتیجه متفاوت انجام شود. هر چند به خاطر پیچیدگی داستان، توصیف شسته رفته‌ای از آن، مثلاً چنانکه گویی تک‌الگوریتمی یگانه می‌تواند کل استنباط استقرایی را نمایش دهد، بسیار مشکل است، حاصل

می‌شود. برای تشخیص تفاوت بین این توجیه «شواهدی» استدلال آزمونهای معنی‌داری، و توجیه «رفتارگرایی»، ممکن است بررسی یک مثال غیررسمی از کاربرد این استدلال در «حالت خاص» سودمند باشد. مثلاً فرض کنید که افزایش وزن بهوسیله روش‌های کاملاً کالبیده^۱ و پایدار، در صورت امکان با استفاده از چند ابزار اندازه‌گیری و مشاهده‌گر، اندازه گرفته می‌شوند و نتایج طی دوره زمانی موردنظر آزمون، تغییر ناچیزی را نشان می‌دهند. این امر را می‌توان به مثابه دلیلی برای این استنباط تلقی کرد که افزایش وزن افراد در محدوده‌ای که توسط حساسیت ترازووها تعیین می‌شود چشم پوشیدنی است. چرا؟

هر چند شخص با پیروی از چنین شیوه‌ای در درازمدت به ندرت افزایش وزنها را نادرست گزارش خواهد کرد، این تنها دلیل منطقی برای این استنباط خاص نیست. بلکه توجیه آن است که خواص احتمالاتی خطای شیوه توزین، آنچه را که واقعاً در این مورد خاص مطرح است نشان می‌دهد. (باید بین این تفسیر و تفسیر شواهدی نظریه نیمن-پرسون که برناوم [۱] آن را پیشنهاد کرده، و داده‌وابسته است، فرق گذاشت.)

آزمون معنی‌داری یک ابزار اندازه‌گیری برای همخوانی با فرضی ویژه است که، همچون ابزارهای اندازه‌گیری به طور عام، با عملکردش در کاربردهای مکرر، در این مورد نوعاً به طور نظری یا با شبیه‌سازی، کالبیده می‌شود. درست همانند کاربست ابزارهای اندازه‌گیری، که در موردی خاص به کار برده می‌شوند، ویژگی‌های عملکرد را برای انجام استنباط درباره جنبه‌های چیز خاصی که اندازه‌گیری می‌شود، جنبه‌هایی که ابزار اندازه‌گیری به طور مقتضی قادر به آشکارسازی آنهاست، به کار می‌گیریم.

البته برای آنکه این وضع برقرار باشد، محاسبات درازمدت احتمالاتی باید برای مورد تحت بررسی مناسب و عملی باشند. اجرای این موضوع در نظریه آمار در بحث‌های استنباط شرطی، انتخاب توزیع مقتضی برای تعیین مقدار p ، ظاهر می‌شود. به نظر می‌رسد مشکلاتی که این امر به همراه دارد بیشتر فنی اند تا مفهومی، و در اینجا به آنها نمی‌پردازیم جز آنکه تذکر دهیم که اقدام به کاربست (یا کوشش برای به کار بستن) ش ب می‌تواند به مشخص کردن آزمون مناسب کمک کند.

۳. انواع فرض صفر و استنباط‌های استقرایی متناظر آنها

در تحلیل آماری داده‌های علمی و فن‌شناختی، تقریباً همیشه اطلاعاتی بیرونی وجود دارند که باید وارد بحث شوند تا از آنچه داده‌ها درباره پرسش اولیه مورد نظر نشان می‌دهند، نتیجه‌گیری شود. معمولاً این ملاحظات پس زمینه‌ای نه از طریق تخصیص احتمال بلکه از راه شناسایی پرسشی که باید مطرح شود، طراحی بررسی، تفسیر نتایج آماری، و ربط دادن آن استنباط‌ها به استنباط‌های علمی اولیه و کاربست آنها برای گسترش و پشتیبانی نظریه زیربنایی وارد بحث می‌شوند. برای آنکه ابزارها به گونه‌ای عاری از مغایطه و طبق منظور به کار روند باید داوری‌هایی درباره اینکه چه چیزی مناسب و آگاهی بخش است در اختیار گذاشته شود. با وجود این، کاربست‌های نظاممندی وجود دارند که می‌توان آنها را متناظر با انواع آزمون و انواع فرضهای صفر، تشریح کرد.

در تفسیر شواهد موجود در موارد خاص نهفته است یا نه. در شرحی که در اینجا عرضه می‌شود، مقدار تحقیق‌یافته p ، دست کم به تقریب، گزارش می‌شود، و شرح مبتنی بر «بی‌ذیر-ردکن» صرفاً شرحی فرضی است که به p تفسیری عملیاتی ببخشد. مشهور است که پیرسون [۲۴] از تفسیر رفتارگرای تنگ‌نظرانه پشتیبانی نمی‌کرد (مايو [۱۷]). نیمن نیز، دست کم در بحث خود با کارناب (نیمن [۲۳])، به نظر می‌رسد که به تمایزی بین تفسیرهای رفتاری و استنباطی اشاره دارد.

در این بحث، نیمن در کوششی برای روشن ساختن ماهیت آمار بسامدگرای نگران اصطلاح «درجۀ تایید» بود که کارناب به کار برده است. نیمن در چارچوب مثالی که در آن یک آزمون بهینه نتوانسته بود H_0 را «رد کند»، به بررسی این مطلب پرداخت که آیا این آزمون H_0 را «تایید» کرده است یا نه. وی تذکر داد که این امر بستگی به معنای واژه‌هایی نظری «تایید» و «اطمینان» دارد و در بافتاری که در آن H_0 «ردنشده» بوده است تلقی آن به مثابه تایید H_0 در صورتی که این آزمون در واقع شانس اندکی برای کشف اختلافی مهم با H_0 داشته باشد، حتی اگر آن اختلاف وجود داشته باشد، «خطرنک» است. از سوی دیگر، اگر این آزمون توان قابل توجهی می‌داشت که آن اختلاف را کشف کند وضعیت «از بیخ و بن متفاوت» می‌بود.

نیمن بر مغایطه‌ای مرتبط با «نتایج منفی» تأکید می‌کند به این مضمون که اگر داده‌های \mathcal{D} یک نتیجه آزمونی به دست دهنده از لحاظ آماری تفاوت معنی‌داری با H_0 (مثلاً فرض صفر «بی اثر» بودن) نداشته باشد، و با این حال آزمون مذکور احتمال کوچکی برای رد H_0 ، حتی زمانی که اختلاف جدی در بین است، داشته باشد آنگاه \mathcal{D} شواهد خوبی برای استنباط اینکه H_0 از سوی \mathcal{D} تایید می‌شود، نیست. طبق این استدلال، تها اگر شانس آنکه آزمون به درستی توانسته باشد اختلافی را کشف کند زیاد باشد، می‌توان از نبود اختلاف مطمئن بود.

نیمن این وضعیت را با تفسیرهای مقتضی برای رفتار استقرایی مقایسه می‌کند. در اینجا تایید و اطمینان را می‌توان برای توصیف انتخاب عمل به کار H_0 بست، مثلاً در خودداری از اعلام یک کشف یا تصمیمی مبنی بر اینکه H_0 را رضایت‌بخش در نظر بگیریم، منطق کار، منطق رفتارگرای عملگرایانه درباره کنترل خطاهای در درازمدت است. این تمایز بدان معنی است که حتی در نظر نیمن، برای شواهد مربوط به تصمیم‌گیری ممکن است ملاکی تمایز از شواهد تشریح نکرد. ما این نظر را در پیش می‌نہیم که اصل شواهدی مورد نیاز اقتصاسی از ش ب (i) برای حالتی است که p -مقدار کوچک نیست:

ش ب (ii). مقدار متوسطی از p -مقدار شاهدی از نبود اختلاف بین δ و H_0 است تها اگر احتمال زیادی وجود داشته باشد که آزمون، بازشی بدتر به H_0 (یعنی، p -مقدار کوچک‌تری) به دست می‌داد در صورتی که اختلافی مانند δ وجود می‌داشت. ش ب (ii) مخصوصاً در بافتار فرضهای «نشانده شده» (زیر) پیش می‌آید.

آنچه نوع استدلال فرضی را به حالت مورد نظر مربوط می‌سازد صرفاً یا عمده‌تاً مقدار پایین نرخ خطاهای درازمدت وابسته به استفاده از این ابزار (یا آزمون) با این روال نیست، بلکه چیزی است که آن نرخهای خطای درباره منع یا پدیده مولد داده‌ها آشکار می‌سازند. محاسبات مبتنی بر خطای اطمینان خاطری می‌بخشند که از تفسیرهای نادرست شواهد در حالت خاص پرهیز

۱.۳ انواع فرض صفر

اکنون شماری از انواع فرض صفر را تشریح می‌کنیم. در بحث زیر، بحثی را که کاکس ([۴]، [۵]) و کاکس و هینکلی ([۶] عرضه کرده‌اند گسترش می‌دهیم. هدف ما در اینجا آن نیست که راهنمایی برای مجموعه همه بافتارهایی که پژوهشگر ممکن است با آنها روبرو شود، ارائه کنیم، بلکه برخی تفسیرهای متفاوت از نتایج آزمون p -مقندرهای مربوط را توضیح می‌دهیم. در بخش ۴.۳، تفسیر ژرفتر استباطهای استقرایی متناظر را که، به نظر ما، با استدلال p -مقداری مجازند (و مجاز نیستند)، بررسی می‌کنیم.

۱. فرضهای صفر شناخته شده. در این مسأله‌ها نه تنها یک مدل احتمال برای فرض صفر فرمول‌بندی می‌شود، بلکه مدل‌های دیگری نیز فرمول‌بندی می‌شوند که امکانهای دیگری را نمایش می‌دهند که در آنها فرض صفر نادرست است، و بنابراین، عمولاً امکانهایی را نمایش می‌دهند که اگر وجود داشته باشند مایلیم آنها را کشف کنیم. در بین شماری از وضعیت‌های ممکن، در معمول‌ترین آنها خانواده‌ای پارامتری از توزیعها وجود دارد که با پارامتر مجهول مانند θ اندیس‌گذاری می‌شود؛ θ به مؤلفه‌های $(\phi, \lambda) = \theta$ افزار می‌شود به‌گونه‌ای که فرض صفر عبارت باشد از $\phi = \phi_0$ ، و λ یک پارامتر مزاحم مجهول است و دستکم در بحث آغازین، ϕ تک‌بعدی است. توجه معطوف به فرضهای مقابل $\phi > \phi_0$ است.

فرمول‌بندی بالا این مزیت فنی را دارد که تا حد زیادی آماره آزمون مناسب (y_t) را با این شرط که حساس‌ترین آزمون ممکن با توجه به داده‌های در دست باشد، تعیین می‌کند.

دو روایت متفاوت از فرمول‌بندی بالا وجود دارد. در یکی، خانواده کامل، یک فرمول‌بندی موقتی است که نه چندان به عنوان پایه‌ای ممکن برای تفسیر نهایی بلکه بیشتر به عنوان ابزاری برای تعیین یک آماره آزمون مناسب در نظر گرفته می‌شود. یک مثال، استفاده از مدلی درجه دوم برای آزمون‌کردن رسانی یک رابطه خطی است؛ به طور کلی رگرسیونهای چندجمله‌ای پایه‌ای ضعیف برای تحلیل نهایی اند اما برای کشف انحرافهای کوچک از یک صورت مفروض، مناسب و قابل تفسیرند. در حالت دوم، خانواده مذکور پایه‌ای محکم برای تفسیر است. بازه‌های اطمینان برای ϕ تفسیری معقول دارند.

یک امکان دیگر، که خیلی به ندرت پیش می‌آید، آن است که یک فرض صفر ساده و یک فرض مقابل ساده تک وجود دارد، یعنی تنها دو توزیع ممکن تحت بررسی‌اند. اگر این دو فرض برپایه‌ای برابر در نظر گرفته شوند، بهتر است این تحلیل به عنوان یک تحلیل تشخیصی فرضی یا واقعی در نظر گرفته شود یعنی تعیین اینکه کدام یک از دو امکان (یا به طور کلی ترکدام یک از چند امکان محدود) مناسب‌تر است اگر با امکانهای مختلف برپایه‌ای از لحاظ مفهومی برابر رفتار شود.

در این حالت دو رویکرد کلی وجود دارند. یکی کاربست نسبت درستنمایی به عنوان شاخصی از برازش نسبی است، که امکان دارد همراه با کاربردی از قضیه بیز باشد. دیگری، که با رویکرد احتمال خطای همخوانی بیشتری دارد، آن است که هر مدل را به نوبت به مثابه فرض صفر و دیگری را به عنوان فرض مقابل می‌گیرد و منجر به ارزیابی این امر می‌شود که آیا داده‌ها همخوان با هر دو فرض، یکی از آنها، یا هیچ یک از آنهاست. با کاربست ش ب در این حالت، وقتی که در چارچوب نیمن-پرسون درآمده باشد، اساساً یک تفسیر نتیجه می‌شود.

این سه حالت را به ترتیب خانواده‌ای صوری از فرضهای مقابل، خانواده‌ای خوش‌بینی از فرضهای مقابل، و خانواده‌ای از امکانهای گسته می‌نامیم.
۲. تقسیم فرضهای صفر. در بسیاری از موارد، مخصوصاً در کاربردهای فن‌شناسخی اما نه تنها در آنها، کانون توجه ما مقایسه دو یا چند شرط، فرایند یا تیمار است بدون هیچ دلیل خاصی برای اینکه انتظار داشته باشیم آنها دقیقاً یا تقریباً همانند باشند، مثلاً دارویی جدید در مقایسه با دارویی استاندارد ممکن است نرخهای بقا را افزایش یا کاهش دهد.

در واقع، شخص دو آزمون را ترکیب می‌کند، نخست برای امتحان این امکان که مثلاً $\mu_1 < \mu_2$ ، دیگری برای $\mu_1 > \mu_2$. در این حالت، این آزمون دوطرفه هر دو آزمون یک طرفه را، هر یک با سطح معنی داری خود، ترکیب می‌کند. سطح معنی داری به‌دلیل «اثر گزینش»، دو برابر p کوچک‌تر است (کاکس و هینکلی ([۶]، ص ۱۰۶). در بخش ۴ به این موضوع برمی‌گردیم. در این حال فرض صفر مبنی بر تفاوت صفر، وضعیت‌های ممکن را به دو ناحیه کیفیتاً متفاوت نسبت به ویژگی مورد آزمون تقسیم می‌کند، وضعیت‌هایی که در آنها یکی از تیمارها نسبت به دیگری برتر است و وضعیت‌هایی که در آنها همان تیمار بدتر است.

۳. فرضهای صفر مبنی بر نبود ساختار. در بسیاری از پژوهش‌های کتابیش تجربی در رشته‌هایی که پایه نظری محکمی ندارند، داده‌ها به امید یافتن ساختاری، اغلب به صورت واستگی‌هایی بین ویژگی‌های ناشناخته، گردآوری می‌شوند. این کار در همه‌گیرشناصی به صورت آزمونهایی از عاملهای خطر بالقوه برای بیماری‌ای که علت ناشناخته دارد، در می‌آید.

۴. فرضهای صفر مبتنی بر رسانی مدل. حتی در حالت کاملاً نشانده شده که در آن خانواده‌ای کامل از توزیعهای تحت بررسی وجود دارند و این خانواده‌ها به طور بالقوه آنقدر غنی اند که خواه فرض صفر درست یا غلط باشد داده‌ها را تبیین می‌کنند، این امکان وجود دارد که اختلافی مهم با مدل وجود داشته باشد که گسترش، اصلاح، یا جانشینی کردن کامل مدل به کار رفته برای تفسیر را توجیه کند. در بسیاری از رشته‌ها، مدل‌های آغازین به کار رفته برای تفسیر کاملاً آزمایشی اند؛ در دیگر رشته‌ها، به خصوص در برخی زمینه‌های فیزیک، مدل‌ها پایه‌ای محکم در نظریه و آزمایش‌های گسترده دارند. اما در همه موارد اگرچه به صورت غیررسمی، امکان مشخص کردن نادرست مدل را باید پذیرفت.

در این صورت انتخابی نه چندان آسان بین یک آماره آزمون نسبتاً خاص‌نگر که نسبت به انواع خاص نارسانی مدل (جهت‌های انحراف خاص) حساس است، و آزمونهای موسوم به کلی نگر که در برآرایه ماهیت انحرافها مفروضات قوی ندارند، مطرح است. روشن است که آزمونهای اخیر نسبت به فرضهای مقابل خاص ناحساس، و اغلب فوق العاده ناحساس، خواهند بود. این دونوع به طور کلی با آزمونهای خرى دو با درجه‌های آزادی کم و زیاد متناظرند. برای آزمون خاص‌نگر می‌توانیم یا آماره آزمونی مناسب، یا تقریباً به طور هم‌ارز خانواده‌ای تصویری از فرضهای مقابل را برگزینیم. مثلاً برای امتحان توافق n مشاهده متنقل با یک توزیع پواسون در واقع می‌توانیم توافق واریانس نمونه با میانگین نمونه را به وسیله یک آزمون پراکنش خى دو (یا معادل دقیق آن) آزمون کنیم یا توزیع پواسون را مثلاً در یک خانواده دوچمله‌ای منفی بشناسیم. ۵. فرضهای صفر پرمایه. در برخی بافتارهای خاص، نتایج مبتنی بر فرض صفر ممکن است شواهدی جدی برای ادعاهای علمی در بافتارهایی باشند

نقاطی از نمونه که به اندازه مقادیر مشاهده شده یا کمتر از آنها محتمل اند. هر چند این کار اغلب به نتایج معقول می‌انجامد در اینجا این مسیر را دنبال خواهیم کرد.

۳.۳ استنباطهای استقرایی برپایه برامدهای آزمونها

استدلال براساس آزمون معنی‌داری چگونه استنباطهای استقرایی با ارزیابیهای شواهدی را در حالتی‌گوناگون تأیید می‌کند؟ تفسیر عملیاتی فرضی p -مقدار روشن است اما پیامدهای ژرفتر مربوط به یکی از دو مقدار متوسط یا کوچک p -مقدار چیست‌اند؟ این پیامدها بشدت به (i) نوع فرض صفر، و (ii) ماهیت انحراف یا فرض مقابل تحت برسی، وابسته‌اند و همچنین (iii) به اینکه آیا به تفسیر مجموعه‌هایی خاص از داده‌ها علاقه‌مندیم، چنانکه در اکثر کارهای آماری تفصیلی پیش می‌آید، یا اینکه مدلی کلی برای تحلیل و تفسیر در یک رشته علمی را برسی می‌کنیم. حالت اخیر به فرمول‌بندی سنتی نیمن-پیرسون مبنی بر تثبیت یک سطح بحرانی و به معنایی، پذیرفتن H_0 اگر $\alpha > p$ و ردکردن H_0 در غیر این صورت، نزدیک است. برخی نوافع آشنای کاربست طبق عادت یا مکانیکی p -مقدار را برسی می‌کنیم.

۴.۳ کاربست عادت‌وار p -مقدارها

تصور کنید که $5^{\circ}\text{ر} = \alpha$ گرفته می‌شود و نتایج در صورتی به مقاله‌ای قابل انتشار می‌انجامد که تنها به ازای p ذی‌ربط، داده‌ها $5^{\circ}\text{ر} < p$ را نتیجه دهند. منطق کار از نوع منطق رفتارگرایی است که پیش‌تر توضیح داده شد. اما در اکثریت عمدۀ بحث‌های آماری در این زمینه، که سابقه‌اش به کار تیتس [۳۲] و قبل از آن برمی‌گردد، چنین رویکردی محکوم می‌شود، هم به واسطه این نگرانی که شیوه‌های مکانیکی، خودکار، و بدون تفکر را شویق می‌کند، و هم به خاطر تأکید این رویکرد بر باورهای ذی‌ربط علاوه بر آزمون فرضها. در واقع چند مجله علمی در برخی رشته‌ها عملًا استفاده از p -مقدار را منوع کرده‌اند. در دیگر رشته‌ها، مانند تعدادی از زمینه‌های همه‌گیرشناختی، معمول آن است که بر بازه‌های اطمینان ۹۵ درصدی تأکید شود که با جریان اصلی بحث آماری هماهنگ است. البته، این کار شخص را از ایناز به ارائه توضیح مقتضی بسامدگرایانه درباره کاربست و تفسیر سطحهای اطمینان، که در اینجا ارائه نمی‌کنیم، نمی‌رهاند (رک. بخش ۶.۳).

با وجود این، کاربست نسبتاً مکانیکی p -مقدارها، در عین حال که در معرض ریختن است، خیلی دور از عمل رایج در برخی رشته‌ها نیست؛ و این روش به عنوان یک ابزار غربالگری، با به رسمیت شناختن امکان خطأ، و کاستن از امکان انتشار نتایج گمراهنکننده، به کار گرفته می‌شود. یک نقش مشابه آزمونها در کار بنگاه‌های ناظر، به ویژه FDA (بنگاه دارویی فدرال آمریکا)، دیده می‌شود. ملزم ساختن بررسیها به نشان دادن مقداری از p که کوچک‌تر از یک سطح تخصیص یافته به وسیله آزمونی از پیش تعیین شده باشد، ممکن است باعث انعطاف‌ناپذیری شود، و انتخاب سطح بحرانی دلخواه باشد. با این حال، چنین شیوه‌هایی دارای محسن بی‌طرفی و استقلال نسبی از دستکاریهای نامعقول‌اند. هر چند پایین‌تری یه یک p -مقدار ثابت ممکن است باعث سوگیری نوشتگان به سمت نتیجه‌گیری‌های مثبت شود، اطمینانی خواهد شد از برخی خواص درازمدت معلوم و مطلوب را به همراه می‌آورد. ملاحظه خواهد شد که این خواص به ویژه برای مثال ۳ بخش ۲.۴ مناسب‌اند.

که شایسته قرار گرفتن در یک رسته پنجم‌اند. در اینجا، نظریه‌ای مانند T که برای آن شواهد نظری و/یا تجربی قابل ملاحظه‌ای وجود دارند، پیشگویی می‌کند که H_0 دست کم با تقریب بسیار خوبی، وضعیت واقعی است.

(الف) در یک روایت، ممکن است نتایجی به ظاهر بی‌هنگار برای وجود داشته باشد، و آزمونی طراحی می‌شود که مجال کافی برای آشکار ساختن ناهمخوانی با H_0 در صورتی که این نتایج بی‌هنگار، واقعی باشد.

(ب) در روایت دوم، نظریه‌ای رقیب مانند T^* ناهمخوانی مشخصی با H_0 را پیشگویی می‌کند و آزمون معنی‌داری برای تشخیص بین T و آن نظریه T^* (در قلمرویی که تاکنون آزمون نشده) طراحی می‌شود. به عنوان مثالی از (الف)، نظریه فیزیکی حاکی از آن است که چون کواتنوم ازرسی در میدانهای برقطیسی غیریوننده، مانند ازرسیهای حاصل از خطهای انتقال فشارقوی، بسیار کمتر از حد لازم برای شکستن یوند ملکولی است، نباید واقع شدن در معرض این گونه میدانها اثر سلطان زایی داشته باشد. بدین ترتیب، در یک آزمایش تصادفیده که در آن دو گروه موش تحت شرایط همانند هستند جز آنکه یک گروه در معرض چنین میدانی قرار داده می‌شود، فرض صفر مبنی بر اینکه نرخهای شیوع سرطان در این دو گروه یکسان‌اند ممکن است کاملاً درست باشد و در تحلیل داده‌ها در درجه اول اهمیت و توجه باشد. البته فرض صفری از این نوع کلی لازم نیست که مدلی با اثر صفر باشد؛ ممکن است اشاره به توافقی با بافت‌های تجربی کاملاً اثبات شده قبلی یا نظریه‌ای داشته باشد.

۲.۳ برخی نکته‌های کلی

در بالا اساساً آزمونهای یک‌طرف را توصیف کردیم. گسترش بحث به آزمونهای دو‌طرفه برخی مسائل را در زمینه تعریف پیش می‌آورد اما در اینجا درباره آنها بحث خواهیم کرد.

چند نوع از فرضهای صفر مستلزم مشخص کردن احتمال به طور ناقص‌اند. یعنی ممکن است تنها فرض صفر به روشنی مشخص شده باشد. می‌توان استدلال کرد که در فرمول‌بندی کامل احتمال باید همیشه هم فرض صفر و هم امکانات فرضهای مقابل شدنی را در نظر گرفت. این امر شاید از لحاظ اصولی معقول به نظر برسد اما به مثابه راهبردی که مستقیماً به کار برد سود، اغلب شدنی نیست؛ در هر حالت، مدل‌هایی که همه امکانات معقول را در برخواهند گرفت باز هم ناکامل خواهند بود و با اثرهای جانی زیان‌بخش خود حتی مسائلهای ساده را پیچیده‌تر خواهند ساخت.

اما توجه داشته باشید که در همه فرمول‌بندی‌هایی که در اینجا به کار گرفته شده‌اند نوعی مفهوم تبیین داده‌ها به صورت بدیلی برای فرض صفر از راه انتخاب آماره آزمون در بین است؛ موضوع این است که چه هنگام این انتخاب از طریق یک فرمول‌بندی احتمالاتی صریح صورت می‌گیرد. اصل کلی شواهد ش ب کمک می‌کند بفهمیم که در بافت‌هایی مشخص، حالت قبلی برای انجام یک ارزیابی شواهدی کفایت می‌کند (رک. بخش ۳.۳).

اما، گاهی استدلال می‌شود که انتخاب آماره آزمون را می‌توان برپایه توزیع احتمال به عنوان آماره آزمون، و بدین ترتیب با جمع بستن احتمالها روی همه

داد و با این کار پایه‌ای گستردگر برای نتیجه‌گیری‌ها در بارهٔ پارامتر تعريف‌کننده به دست می‌آید.

فرضهای صفر تقسیم‌کننده و فرضهای مبنی بر نبود ساختار در حالت فرضهای صفر تقسیم‌کننده، ناهمخوانی با فرض صفر (با بهکارگیری مقدار دو طرفه p) بیانگر جهت انحراف است (یعنی، اینکه کدام یک از دو تیمار برتر است؟)؛ همخوانی با H_0 حاکی از آن است که این داده‌ها شواهدی رسا حتی برای جهت اختلاف به دست نمی‌دهند. اغلب انتقادهایی شنیده می‌شود که آزمودن فرض صفری که می‌دانیم نادرست است بی‌فایده است، اما حتی اگر انتظار نداشته باشیم که دو میانگین مثلاً برابر باشند، این آزمون برای تقسیم انحرافها به انواعی که از لحاظ کیفی متفاوت‌اند آگاهی بخش است. وقتی فرض صفر مانند فرض صفر نبود ساختار است تفسیر مشابهی انجام می‌شود. مقدار متوسط p بیانگر آن است که داده‌ها حساسیت کافی برای کشف ساختار ندارند. اگر داده‌ها محدود باشند این امر ممکن است فقط هشداری علیه تفسیر بیش از حد باشد و نه شاهدی براینکه فکر کنیم در واقع ساختاری وجود ندارد. یعنی آزمون ممکن است ظرفیت اندکی برای کشف هرگونه ساختاری داشته بوده باشد. اما مقدار کوچک p شواهدی مبنی بر وجود اثری واقعی را نشان می‌دهد؛ که جستجو برای تفسیر واقعیت‌بنیاد از چنین اثری ذاتاً مستعد خطا نخواهد بود.

هنگامی که ارزیابی‌های غیررسمی در بارهٔ گواه‌آوری یا حساسیت آزمونها صورت می‌گیرد، استدلال مشابهی مصدق دارد. اگر داده‌ها آنقدر گستردۀ باشند که همخوانی با فرض صفر دال بر نبود اثری با اهمیت عملی باشد، و یک p -مقدار به قدر معقول بزرگ به دست آمده باشد، آنگاه می‌توان آن را به عنوان شواهدی مبنی بر نبود اثری با اهمیت عملی گرفت. همین‌طور، اگر داده‌ها از چنان دامنه محدودی باشند که بتوان فرض کرد که داده‌های همخوان با فرض صفر با انحرافهای واحد اهمیت علمی نیز سازگارند، آنگاه بزرگ بودن p -مقدار، استنباطی مبنی بر نبود انحرافهای مهم از فرض صفر را تضمین نمی‌کند.

فرضهای صفر مبنی بر رسایی مدل. وقتی فرضهای صفر ادعاهایی مبنی بر رسایی مدل‌اند، تفسیر نتیجه‌های آزمون بستگی به آن خواهد داشت که آیا آماره آزمون نسبتاً خاص‌نگر داریم که نسبت به انواع خاص نارسایی مدل حساس است، یا آزمونهایی به اصطلاح کلی‌نگر داریم. همخوانی با فرض صفر در حالت اول شواهدی مبنی بر نبود آن نوع انحرافی که آزمون برای کشف آن حساس است به دست می‌دهد، در صورتی که، با آزمون کلی‌نگر، این همخوانی کمتر آگاهی بخش است. در هر دو نوع آزمون، p -مقدار کوچک شواهدی از نوعی انحراف است، اما تا آنجا که مدل‌های بدیل گوناگون بتوانند تخطی مشاهده شده را توجیه کنند (یعنی، تا آنجا که این آزمون توانایی اندکی برای تمیز دادن بین آنها داشته باشد)، این داده‌ها به‌نهایی، فقط می‌توانند پیشنهادهایی موقتی در بارهٔ مدل‌های بدیلی که باید امتحان شوند، فراهم سازند.

فرضهای صفر پرمایه. در حالتهای پیشین، همخوانی با یک فرض صفر حداقل می‌توانست، برطبق توانایی آزمون در آشکارسازی اختلاف، شواهدی برای کنار گذاشتن اختلافهای با مقدار یا نوع مشخص به دست دهد. در مورد فرضهای صفر پرمایه [جدی]، چیز بیشتری می‌توان گفت. اگر فرض صفر نشان‌دهنده یک پیش‌بینی از نظریه‌ای باشد که از لحاظ قابلیت کاربرد عام تحت بررسی است، سازگاری با فرض صفر را می‌توان به عنوان

۰.۳ کاربست شواهدی-استقرایی p -مقدارها

اکنون به کاربست آزمونهای معنی‌داری روی می‌آوریم که هر چند رایج‌ترند، در عین حال، به عنوان یک ابزار تنها برای کمک به تحلیل مجموعه‌های خاص از داده‌ها، و/یا بناهادن استنباطهای استقرایی براساس داده‌ها، بحث‌انگیزترند. در این بحث پیش‌پیش فرض می‌شود که توزیع احتمال مورد استفاده برای ارزیابی p -مقدار تاحد ممکن برای داده‌های خاص تحت تحلیل مناسب است.

اصل بسامدگرای کلی برای استدلال استقرایی، ش ب، یا چیزی مانند آن، می‌تواند راهنمایی برای رسیدن به حکم مقتضی دربارهٔ شواهد یا استنباط راجع به هر نوع فرض صفر باشد. تقریباً همان‌گونه که شخص در بارهٔ تغییرات وزن بدن برایه مشخصات عملکردی ترازووهای مختلف استنباطهایی می‌کند، می‌تواند از نتایج آزمون معنی‌داری با استفاده از خواص نزخ طای آزمونها نیز استنباط به عمل آورد. این خواص ظرفیت آزمون خاص را به‌دلیل آنکه توانسته است ناسازگاریها و اختلافهای موجود در جنبه‌های تحت بررسی را آشکار کند، نشان می‌دهند، و این امر به نوبه خود ربط‌دادن p -مقدارها را به فرضهایی در بارهٔ فرایندی که به‌طور آماری مدل‌بندی شده است، مجاز می‌سازد. نتیجه می‌شود که برای ارائه یک شرح بسامدگرای رسا از استنباط باید در تأمین اطلاعات به منظور اجرای ش ب اهتمام ورزید.

فرضهای صفر نشانده شده. در حالت فرضهای صفر نشانده شده، استفاده از p -مقدارهای کوچک به عنوان شواهدی از اختلاف با فرض صفر در جهت فرض مقابل، سرراست است. اما، فرض کنید که معلوم شود داده‌ها موافق با فرض صفرند (p کوچک نباشد). ممکن است شخص، با استفاده از همان منطق آزمون معنی‌داری، این امر را، به مثابة شاهدی برای اینکه هر اختلاف با فرض صفر کوچک‌تر از δ است تلقی کند. در چنین حالتهایی، همخوانی با فرض صفر ممکن است شواهدی از نبود اختلاف با فرضهای صفر دارای اندازه‌های گوناگون فراهم سازد، چنانکه در ش ب (ii) تصریح شد.

برای استنباط نبود اختلافی به بزرگی δ با H_0 می‌توانیم (δ, β) احتمال مشاهده برازشی بذرگی به H_0 را در صورتی که $\delta = \mu + \beta\mu$ ، بیازمایم، اگر این احتمال نزدیک به یک باشد، براساس ش ب (ii)، داده‌ها شواهد خوبی برای $\delta < \mu$ دارند. بدین ترتیب، (δ, β) را می‌توان به عنوان جدیت یا شدت آزمون در سنجش اختلاف δ در نظر گرفت؛ به عبارت دیگر می‌توان گفت که $\delta < \mu < \mu$ آزمون شدیدی را گذرانده است (مايو [17]).

این کار از تفسیر تضمین نشده سازگاری با H_0 در آزمونهای غیرحساس جلوگیری می‌کند. یک چنین ارزیابی برای داده‌های خاص مناسب‌تر است تا مفهوم توان، که نسبت به یک مقدار بحرانی از پیش تعیین شده حساب می‌شود که فراتر از آن، آزمون فرض صفر را «رد کند». یعنی، توان به یک ناحیه رد از پیش مشخص شده مربوط می‌شود، نه به داده‌های خاص تحت تحلیل.

هر چند بیش حساسی کمتر محتمل است که مشکلی باشد، اگر آزمونی آنچنان حساس باشد که p -مقداری برابر یا حتی کوچک‌تر از مقدار مشاهده شده، محتمل باشد حتی وقتی $\delta < \mu$ ، آنگاه مقداری کوچک از p شاهدی برای انحراف از H_0 به اندازه بیش از δ نیست.

اگر خانواده‌ای از فرضهای مقابله‌ای وجود داشته باشد، می‌توان مجموعه‌ای از بازه‌های اطمینان برای پارامتر مجهول تعريف‌کننده H_0 به دست

را به دست می‌دهد. اگر چنین باشد، تحلیل بازه اطمینان نیز باید ذیل اصل بسامدگرای کلی ما قرار گیرد و قرار می‌گیرد. در آزمون یک طرفه $\mu = \mu$ در مقابل $\mu < \mu$ ، p -مقدار کوچک متناظر با این است که μ (صرفاً) از بازه اطمینان متناظر $(\mu - 2\sigma, \mu + 2\sigma)$ (دو طرفه) (یا $p = 1$ برای بازه اطمینان یک طرفه) حذف شود. اگر μ را μ_L یعنی کران اطمینان پایینی، قرار می‌دادیم، آنگاه نتیجه‌های ناهمخوان‌تر با احتمال زیاد $(p = 1)$ رخ می‌داد. بدین ترتیب، شش اجازه قبول این وضع را به مثابه شواهد ناسازگاری با $\mu_L = \mu$ (در جهت مثبت) می‌دهد. افزون بر آن، این استدلال مزیت درنظر گرفتن چند بازه اطمینان را در دامنه‌ای از سطحها، بهجای آنکه صرفاً گزارش شود که مقدار پارامتری مفروض درون بازه‌ای در یک سطح اطمینان ثابت قرار دارد یا خیر، نشان می‌دهد.

نیمن نظریه بازه‌های اطمینان را از ابتدا یعنی تنها با انتکای تلویحی و نه تصریحی بر کار قبلى اش با پیرسون در نظریه آزمونها ابداع کرد. این امر تا حدودی بستگی به نحوه عرضه مطلب دارد که براورده بازه‌ای را اصولاً آنقدر متفاوت از آزمون حفظ تمایز مفهومی [بین این دو] پدید آمده است یا نه. از سوی مجزا برای حفظ تمایز مفهومی [بین این دو] پدید آمده است یا نه. از سوی دیگر مزایای قابل توجهی برای درنظر گرفتن یک حد اطمینان، بازه یا ناحیه‌ای به عنوان مجموعه مقدارهای پارامتر سازگار با داده‌ها در سطحی مشخص، وجود دارد که با آزمودن هر مقدار ممکن به وسیله شیوه‌های متقابلاً همخوان تعیین می‌شود. به خصوص، این رویکرد، در سطح معنی‌داری مشخصی، به راحتی با بازه‌های اطمینانی که تهی هستند یا مرکب از همه مقادیر ممکن پارامترند مواجه می‌شود. چنین ناحیه‌تهی یا نامتناهی دال بر آن است که داده‌ها با همه مقادیر ممکن پارامترها ناسازگارند یا با همه مقادیر ممکن سازگارند. ساختن مثالهایی که در آنها این نتایج کاملاً مناسب به نظر برستند، آسان است.

۴. برخی پیچیدگیها: اثرهای گزینش

فرمول‌بندی آرمانی که در تعریف آغازین آزمون معنی‌داری به‌چشم می‌خورد علی‌الاصول با یک فرض و یک آماره آزمون شروع می‌شود، و مراحل بعدی، کسبِ داده‌ها، به‌کار بستن آزمون و بررسی برآمد است. پس شیوه فرضی مستتر در تعریف آزمون در حد قابل قبولی با آنچه انجام شده تطابق دارد؛ برآوردهای ممکن، مقدارهای ممکن متفاوت از آماره آزمون مشخص هستند. این امر امکان می‌دهد که ویژگیهای توزیع آماره آزمون در کسب اطلاع از ویژگیهای متناظر سازوکار مولد داده‌ها به‌کار آید. شیوه‌ای که عملاً از آن پیروی می‌شود به دلایل گوناگونی ممکن است متفاوت باشد و اینک یک جنبه عام آن را بررسی می‌کنیم.

اگلی اتفاق می‌افتد که یا فرض صفر یا آماره آزمون از وارسی مقدماتی داده‌ها تأثیر می‌ذیرد، به‌طوری که شیوه واقعی تولید نتیجه آزمون نهایی تغییر می‌باید. این امر به‌نوبه خود ممکن است توانایی آزمون را در کشف اختلافات با فرض صفر به‌گونه‌ای قابل اعتماد، تغییر دهد، و تعدیهایی را در احتمالهای خطای آن لازم بنماید.

تا آنجا که p به عنوان جنبه‌ای از رابطه منطقی یا ریاضی بین داده‌ها و مدل احتمالی در نظر گرفته شود چنین گزینشهای مقدماتی بی‌ربطاند. این امر برای حصول اطمینان از اینکه p -مقدارها آن منظوری را براورده سازند که

شواهدی اضافی برای این نظریه تلقی کرد، بهویژه اگر آزمون و داده‌ها به قدر کافی حساس باشند که انحرافهای عمدۀ از نظریه را کنار بگذارند. یک جنبه موضوع در اندرز فیشر (کوکان، [۳]) خلاصه شده است که می‌گوید برای کمک به اینکه بررسیهای مشاهداتی ما تفسیر علمی دقیق‌تری داشته باشند، باید نظریه‌هایمان را تفصیلی تر سازیم، که منظورش این بود که باید آزمونهای گوناگونی برای پیامدهای متفاوت نظریه طراحی کنیم، تا امتحانی جامع از استلزم‌های آن به دست آوریم. این نتیجه محدود که مجموعه‌ای از داده‌ها با نظریه همخوان است یک جزء به شواهدی می‌افزاید که اعتبار آنها از روی هم ایناشتن توانایی دشوقق دیگر ناشی می‌شود.

نخستین نوع مثال تحت این رهنمود، ممکن است چنین باشد که نتیجه‌های آشکارا بی‌هنگار برای نظریه یا فرضی مانند T وجود داشته باشند، در حالی که T بازرسی دقیق‌نظری و/یا تجربی قابل توجهی را با موقوفیت پشت سر گذاشته است. اگر نتیجه‌های آشکارا بی‌هنگار برای T واقعی باشند، انتظار می‌رود که H_0 رد شود، به‌طوری که اگر رد نشود، این نتیجه‌ها شواهد مثبت علیه واقعیت بی‌هنگاری‌اند. در دومین نوع این مورد، باز هم یک نظریه کاملاً آزمون شده T داریم، و نظریه‌ای رقیب مانند T^* تعیین می‌شود که در حوزه‌ای تاکنون آزموده نشده از حیث یک اثر با T در تضاد است. اگر پیشگویی حاصل از T را با فرض صفر یکی بگیریم، هرگونه اختلاف در جهت T^* شناس خوبی دارد که کشف گردد، به‌گونه‌ای که اگر هیچ انحراف معنی‌داری یافته نشود، این امر شاهدی به نفع T از جنبه مورد آزمون است. اگرچه نظریه نسبیت عام (GTR) در دهه ۱۹۶۰ با بی‌هنگاری‌های مواجه نبود، رقبای GTR شکست اصل هم‌ارزی ضعیف،^۱ WEP، را برای اجسام خودگرانشی سنگین، مثل سامانه زمین-ماه پیشگویی می‌کردند: این اثر که اثر نوردوت^۲ نامیده می‌شود، برای GTR برابر صفر است (که با فرض صفر یکی گرفته می‌شود) و برای رقبا ناصر خواهد بود. اندازه‌گیری زمانهای سفر رفت و برگشت بین زمین و ماه (بین سالهای ۱۹۶۹ و ۱۹۷۵) امکان داد که دربرابر وجود چنین بی‌هنگاری برای GTR تحقیق شود. پیدائشدن شواهدی علیه فرض صفر، کرانهای بالایی را برای تخطی ممکن از WEP تعیین کرد، و چون آزمونها به قدر کافی حساس بودند، این اندازه‌گیری‌ها شواهد خوبی فراهم ساختند که اثر نوردوت وجود ندارد و بدین ترتیب شواهدی به نفع فرض صفر در دست است (ویل، [۳۱]). توجه کنید که یک چنین نتیجه منفی شواهدی به نفع همه GTR (در همه حوزه‌های پیشگویی) فراهم نمی‌سازد، بلکه شواهدی به نفع درستی آن از حیث این اثر به دست می‌دهد. منطق کار چنین است: نظریه T گزاره H_0 را دستکم با تقریبی نزدیک به وضعیت واقعی پیشگویی می‌کند؛ نظریه رقیت T^* اختلاف مشخصی را با H_0 پیشگویی می‌کند، و احتساب زیاد دارد که اگر T^* صحیح باشد آزمون چنین اختلافی از T را کشف کند. پس کشف نشدن هیچ اختلافی، شاهدی بر نبود آن است.

۶.۳ بازه‌های اطمینان

چنانکه در بالا اشاره شد، در بسیاری از مسائله‌ها تدارک بازه‌های اطمینان، علی‌الاصول در دامنه‌ای از سطحهای احتمال، ثمر بخش‌ترین تحلیل فراوانی گرا

1. Weak Equivalence Principle 2. Nordvedt effect

وفرض صفر متناظر آن را گزارش می‌کند. نکته‌های اساسی عبارت‌انداز: استقلال آزمونها و گزارش نکردن نتیجه‌های حاصل از آزمونهای غیرمعنی دار. مثال ۲. روایتی بسیار آرمانی از آزمون مربوط به تطبیق DNA با نمونه‌ای مفروض، که شاید از آن یک بزهکار باشد، این است که جستجو در سراسر یک پایگاه داده‌های تطبیق‌های ممکن، هر بار با یک فرد، صورت گرفته استخان می‌شود که فرض توافق با نمونه مفروض رد می‌شود یا نه. فرض کنید که حساسیت و ویژگی هر دو خیلی زیاد هستند. یعنی، احتمالهای منفی‌های کاذب و مثبت‌های کاذب هر دو خیلی کوچک‌اند. نخستین فرد، از پایگاه داده‌ها که فرض در مورد آن رد شود، در صورتی که چنین فردی وجود داشته باشد، تطبیق واقعی اعلام می‌شود و این فرایند همانجا متوقف می‌گردد.

مثال ۳. در یک بررسی ریزایله^۱ چندهزار ژن از لحظه تجلی بالقوه، مثلاً تفاوتی بین وضعیت بیماری نوع ۱ و نوع ۲ آزمون می‌شوند. بنابراین چند هزار فرض وجود دارند که هر یک با فرض صفر مربوط در یک گام تحت بررسی قرار می‌گیرند.

مثال ۴. برای بررسی واستگی یک متغیر پاسخ یا برمدی مانند y به متغیری تبیینی مانند x ، قصد آن است که از تحلیل رگرسیون خطی y بر حسب x استفاده شود. وارسی داده‌ها حاکی از آن است که بهتر است از رگرسیون $\log y$ بر حسب x استفاده شود، مثلاً این دلیل که رابطه به حالت خطی y نزدیک‌تر است یا چون مفروضات ثانوی، مانند ثابت‌بودن واریانس خط، با تقریب بهتری، صادق‌اند.

مثال ۵. برای بررسی واستگی یک متغیر پاسخ یا برمدی مانند y به تعداد قابل توجهی از متغیرهای تبیینی بالقوه x ، شیوه داده‌واستگی گزینش متغیر بدکار گرفته می‌شود تا نمایشی به دست آید که سپس با روش‌های استاندارد برآزانده شده و فرضهای ذی‌ربط آزمون شوند.

مثال ۶. فرض کنید وارسی مقدماتی داده‌ها نوعی اثرکلآنامتنظریاً نوعی نظم را نشان دهد که در مرحله‌های آغازین به آن فکر نشده بود. در نتیجه آزمونی رسمی، این اثر «به شدت معنی دار» است. معقول است چه نتیجه‌ای گرفته شود؟

۴. نیاز به تعدیلهایی به خاطر گزینش

برای بحث عمیق درباره همه این مثالها مجال کافی در این مقاله نیست. یک موضوع اساسی ناظر به این است که کدام یک از این وضعیتها نیاز به تعدیل به خاطر آزمون کردن چندگانه یا گزینش داده‌واستگی چندگانه وجود دارد با هم باشد. تصور کلی از سرشت یک نظریه تحلیل و تفسیر بسامدگرا چگونه کمک می‌کند که به جوابها برسیم؟

پیشنهاد ما این است که این امر به روای زیر صورت گیرد: نخست باید در نظر گرفته شود که آیا بافتار مورد نظر بافتاری است که در آن موضوع اصلی مورد نظر کنترل نرخهای خطای دریک سلسه از کاربردهاست (هدف رفتاری)، یا انجام استنباط استقرایی خاص یا ارزیابی شواهد خاص (هدف انتباطی). احتمالهای خطای مربوط را برای حالت اول می‌توان تغییر داد و برای حالت اخیر نمی‌توان. دوم، لازم است دنباله نکارهای ذی‌ربطی را که باید بسامدگرا را برپایه

1. microarray

برای آنها در استنباط بسامدگرا، خواه در بافت‌های رفتاری یا شواهدی، در نظر گرفته شده است، کفایت نمی‌کند. در حدی که شخص محاسباتی مبتنی بر خطای بخواهد که به آزمون معنایی حاکی از قابلیت کاربرد در کارهای آمار بسامدگرا بیخشش، تحلیل مقدماتی و گزینش ممکن است بسیار مناسب باشد.

نکته‌کلی مورد نظر هم در نوشتگان فلسفی و هم نوشتگان آماری بسیار مورد بحث قرار گرفته است؛ در مباحث فلسفی تحت عنوانهای مانند لزوم نوبودن یا پرهیز از فرضهای مورد ویژه، و در مباحث آماری به عنوان قواعدی علیه دزدکی نگاه‌کردن به داده‌ها یا جستجوی معنی داری، و در نتیجه لزوم به حساب آوردن اثرهای گزینش، بحث شده‌اند. موضوع کلی این است که آیا وقتی فرض صفر H_0 به طریقی برای آزمون شدن ساخته یا انتخاب می‌شود که به رابطه مشاهده شده خاصی بین H_0 و y ، خواه توافق یا عدم توافق، منجر شود، آیا تأثیر شواهدی داده‌های y بر استنباط یا فرض H_0 تغییر می‌یابد یا نه. کسانی که با رویکردهای منطقی به تأیید نظر مساعد دارند می‌گویند نه (متلاً میل [۲۰]، کیزن [۱۴]، در صورتی که کسانی که به مفهوم آماری خطای بیشتری دارند می‌گویند آری (ویول [۳۰]، پیرس [۲۵]). به پیروی از این فلسفه‌ای خیر، پوپر خواستار آن بود که دانشمندان پیش‌پیش مشخص سازند چه برآمدۀای را به مثابه ابطال H_0 در نظر خواهند گرفت، شرطی که حتی خود او وادار به رد آن شد؛ کل موضوع در فلسفه لایحل باقی مانده است (مايو [۱۷]). ملاحظات آماری خطای امکان پیش‌پیش را با تدارک ملاک‌هایی برای زمانی که گزینشهای داده‌واستگی گوناگون اهمیت دارند و نحوه به حساب آوردن تأثیر آنها بر احتمالهای خطای، فراهم می‌سازد. بهویژه، اگر به دلیل آنکه آماره آزمون بزرگ است فرض صفر برای آزمون انتخاب شود، احتمال یافتن چنین ناهمخوانی یا ناهمخوانی‌های دیگر حتی تحت فرض صفر ممکن است بزرگ باشد. بدین ترتیب، بنا به ش ب (i)، شواهدی واقعی مبنی بر ناهمخوانی با فرض صفر نخواهیم داشت، و جز آنکه p -مقدار به طرزی مقتضی اصلاح شود استنباط گمراحته خواهد بود. در حدی که شخص از محاسبات مبتنی بر خطای انتظار دارد که معنایی به آزمون بیخشش، معنایی حاوی اطمینان دهی محدود در این باره که ناسازگاری ظاهری در این حالت خاص واقعی است و صرفاً معلول شناس نیست، تعديل p -مقدار لازم است.

چنین تعدیلهایی اغلب در مواردی شامل گزینشهای داده‌واستگی داشته یا در مدل گزینی یا در مدل سازی پیش می‌آیند؛ مسئله تعديل p اغلب در مواردی شامل آزمون فرضهای چندگانه رخ می‌نماید، اما مهم است که این حالتها صرفاً به این دلیل که داده‌واستگی یا آزمون فرض چندگانه وجود دارد با هم اجرا نشوند. اکنون چند حالت خاص را به اجمال شرح می‌دهیم تا نکته‌های اساسی در روایتهای مختلف را آشکار سازیم. سپس بررسی می‌کنیم که در نظر گرفتن گزینش در هر حالت لازم است یا نه.

۱.۴ مثالها

مثال ۱. پژوهشگری مثلاً ۲۰ مجموعه داده‌های مستقل دارد، که هر یک حاکی از اتری متفاوت اما بسیار مرتبط با یقیه است. پژوهشگر همه ۲۰ آزمون را انجام می‌دهد و تنها کوچک‌ترین p را، که در واقع حدود ۵٪ است،

مثال ۱ را می‌توان با یک آزمایش عاملی استاندارد که برای تحقیق اثرهای چند متغیر تبیینی به طور همزمان طراحی شده است، قیاس کرد. در اینجا شماری از پرسش‌های متمایز، هر یک با فرض مربوط و هر یک با p -مقدار مخصوص به خود، وجود دارد. اینکه پرسشها را از طریق یک مجموعه داده‌ها مورد بررسی قرار می‌هیم بدجای آنکه از طریق مجموعه‌های داده‌های جداگانه بررسی کنیم، به یک معنا یک تصادف فنی است. هر p به درستی در زمینه پرسش خودش تفسیر می‌شود. مشکلات برای استنباطهای خاص در صورتی پیش می‌آید که عملای پسیاری از پرسشها را دور بیندازیم و توجه خود را تنها بر یکی، یا به طور کلی، بر شماری اندک، بدليل آنکه دارای کوچکترین p ‌اند، متوجه کنیم. زیرا در این صورت ظرفیت آزمون را در تغییردادن نظرمان به وسیله p -مقداری که به درستی محاسبه شده، تغییر داده‌ایم، خواه شواهدی بر استنباط مورد نظر داشته باشیم یا نه.

مثال ۲ (تبیین اثری معلوم با استقرای حذفی). مثال ۲ شباهت سطحی با مثال ۱ دارد؛ یافتن یک تطبیق DNA تا حدی مثل یافتن یک انحراف از لحاظ آماری معنی دار از یک فرض صفر است: شخص درون داده‌ها جستجو می‌کند و بر یک حالت که در آن «تطبیق» با DNA بزرگار یافت می‌شود تمرکز کرده، عدم تطبیق‌ها را نادیده می‌گیرد. اگر شخص برای «جستجو» در مثال ۱ تعديل انجام می‌دهد، آیا باید به طور کلی به همان طریق در مثال ۲ تعديل انجام دهد؟ نه.

در مثال ۱ نگرانی عبارت است از استنباط کردن یک اثر «بازآفریدنی» واقعی، در حالی که چنین اثری وجود ندارد؛ در مثال ۲، اثری معلوم یا پیشامدی خاص، یعنی DNA بزرگار، وجود دارد و شیوه‌هایی قابل اعتماد به کار گرفته می‌شوند تا علت یا منبع خاص (چنانکه با نزخ پایین «تطبیق اشتباہ» بیان می‌شود) ردیابی گردد. احتمال زیاد دارد که اگر H_0 بزرگار نباشد تطبیقی با فرد H_0 بددست نیاوریم؛ پس بنابر ش ب، یافتن یک تطبیق در سطحی کیفی، شاهد خوبی بر بزرگار بودن H_0 است. افزون بر آن، هر عدم تطبیق یافت شده، بنا به تصریحات مثال، عملان آن شخص را حذف می‌کند؛ از این‌رو، هر قدر چنین نتیجه‌های منفی زیادتر یافت شوند، «تطبیق» استنباط شده محکم‌تر است؛ در صورتی که در مثال ۱ چنین نیست.

چون حداقل یک فرض صفر مبنی بر یک گناهی نادرست است، شواهد بی‌گناهی برای یک نفر، اگرچه به قدر اندک، شناسنامه‌بودن دیگری را افزایش می‌دهد. پس از آنکه شیوه نمونه‌گیری برای آزمون کردن مشخص شود، ارزیابی نزخهای خطأ حتماً میسر است. در اینجا از ارائه جزئیات خودداری می‌کنیم. موردی که شباهت کلی با این وضعیت دارد به بی‌هنگاری مدار عطارد مربوط می‌شود: چون تلاش‌های پرشمار برای ارائه نتیجه‌های بی‌هنگار را با دقت و بدون هیچگونه تعديل‌های موردی پیشگویی کند، بیش از پیش اعتباریافت.

مثال ۳ (داده‌های ریزآرایه). در تحلیل داده‌های ریزآرایه، یک فرض آغازی معقول این است که شمار بسیار بزرگی از فرضهای صفر باید آزمون شوند که کسر نسبتاً اندکی از آنها (ایکیا) نادرست‌اند، و فرض صفر فراگیری که حاکی از نبود هیچ اثر حقیقی باشد، اغلب نامعقول است. در این صورت، مسئله از نوع انتخاب جایگاه‌هایی است که در آنها اثری می‌تواند ثابت شده تلقی گردد. در اینجا، نیاز به تعديل به خاطر آزمون کردن چندگانه، عمدتاً به واسطه این نگرانی

آنها بنا نهاد شناسایی کرد. شرط کلی آن است که ناهمخوانی با یک فرض صفر را به کمک شیوه‌ای ارائه نکنیم که ناهمخوانیها را با فراوانی نسبتاً بیشتری مطرح می‌کند حتی اگر فرض صفر درست باشد. تعیین سریهای فرضی ذی‌ربطی که براساس آنها بسامد خطای یابد حساب شود، ایجاد می‌کند که طبیعت مسئله استنباط در نظر گرفته شود. به طور مشخص‌تر، موضع ویژه‌ای را که به‌منظور استنباطی قابل اعتماد در حالت خاص باید از آنها پرهیز کرد، و ظرفیت آزمون را، به عنوان یک ابزار اندازه‌گیری برای آشکارساختن وجود موضع، باید شناسایی کرد. هنگامی که هدف اصلی ما ارزیابی شواهد خاص است، ش ب برخی رهنمودها را در اختیار می‌گذاریم. به طور مشخص‌تر، مشکل وقتی رخ می‌نماید که برای گرینش یک فرض برای آزمون کردن یا تعیین مشخصات یک مدل زیربنایی، داده‌ها به طریقی به کار گرفته می‌شوند که یا از ش ب عدول می‌شود یا نمی‌توان تعیین کرد که آیا ش ب صادق است یا خیر (مايو و کروس [۱۸]).

مثال ۱ (جستجوی معنی داری آماری). شیوه آزمون بسیار متفاوت با حالتی است که در آن تک‌فرض صفری که معلوم شده از هیئت آماری معنی دار است از قبل به عنوان فرض برای آزمون تعیین شده باشد؛ شاید آن فرض H_0 ، سیزدهمین فرض صفر از بین ۲۰ فرض، باشد. در مثال ۱، نتیجه‌های ممکن عبارت از اعمالهایی که ممکن است از لحاظ آماری معنی دار باشند و معلوم شود که انحراف معنی دار آماری «محاسبه شده» ای از فرض صفر را نشان می‌دهند. بنابراین، احتمال خطای نوع ۱ احتمال یافتن دست‌کم یک چنین متفاوت معنی دار در بین ۲۰ تاست، هر چند که فرض صفر فراگیر درست است (یعنی، همه بیست تفاوت مشاهده شده به واسطه شناسن رخ داده‌اند). احتمال آنکه این شیوه باعث رده کردن خطأ آزمیز شود متفاوت با، و خیلی بزرگ‌تر از ۵٪ (و تقریباً ۶۴٪) است. راههای متفاوت، و در واقع راههای بسیار بیشتری، برای اتکاب خطای در این مثال نسبت به زمانی که یک فرض صفر از قبل مشخص شده، وجود دارد و این امر در p -مقدار تعديل شده بازتاب می‌یابد. این قدر از مطلب کاملاً معلوم است، اما آیا این امر باید بر تفسیر نتیجه در یک چارچوب استنباط استقرایی تأثیر بگذارد؟ بنابر ش ب باید تأثیر بگذارد. اما دغدغه این نیست که از اعلام اشتباہی و مکر اثرهای واقعی در یک سری پرهیز شود، نگرانی آن است که این آزمون به عنوان ابزاری برای تیزی دادن اثرهای واقعی از اثرهای شناسی در این حالت خاص به طور ضعیف عمل می‌کند. زیرا می‌دانیم که دست‌کم یک انحراف چشمگیر از این نوع، حتی اگر همه انحرافها به واسطه شناسن باشند، رایج است؛ آزمون به سختی به ما اطمینان داده است که با پرهیز از اتکاب به چنین اشتباہی در این حالت، کار خوبی انجام داده است. حتی اگر دلیل‌هایی دیگر برای باورداشتن به اصالت اثری که پیدا شده است وجود داشته باشد، انکار می‌کنیم که این آزمون به تنهایی چنین شواهدی را فراهم ساخته است.

گفته‌ایم که محاسبات بسامدگرایانه برای آن به کار می‌روند که حالت خاص را بیازمایند؛ و این امر با مشخص‌سازی ظرفیت آزمونها در عیان ساختن اشتباہ در استنباط صورت می‌گیرد؛ بر این مبنای «شیوه جستجو» ظرفیت اندکی هشداردادن در این مورد دارد که عملان اشتباہ خود بکاهیم حتی وقتی چنین کاهشی ضرورت دارد. از سوی دیگر، اگر شخص p -مقدار را تعديل کند تا نزخ خطای کلی را بازتاب دهد، باز هم آزمون ابزاری می‌شود که در خدمت این هدف است.

این نوع در تحلیلهای آماری بیچیده اهمیت دارند از این لحاظ که سلسه‌ای از انتخابها را که به طور غیررسمی مشخص شده‌اند درباره فرمول‌بندی مدلی که بهترین مدل برای تحلیل و تفسیر است می‌توان انجام داد (اسپانوس [۲۹]).

مثال ۶ (اثر کلاً نامتنظر). این وضع مشکلات عمدۀ را پیش می‌آورد. در علوم آزمایشگاهی که داده‌ها با سرعت معقولی به دست می‌آیند، کوشش برای تکرار نتیجه‌گیری‌ها عملًا اجباری خواهد بود. در دیگر زمینه‌ها جستجو برای داده‌های دیگر که بر موضوع تأثیر دارند لازم می‌آید. تفسیر معنی داری آماری شدید به خودی خود بسیار مشکل خواهد بود، اساساً به این دلیل که گزینش رخ داده است و با هر نوع واقعگرایی، مشخص کردن مجموعه‌ای که در آن گزینش رخ داده است نوعاً سخت یا ناممکن است. اما ملاحظاتی که در مثالهای ۱ تا ۵ مورد بحث واقع شدند می‌توانند راهنمای باشند. مثلاً اگر وضعیت مانند وضعیت مثال ۲ (تبیین اثری معلوم) باشد منبع را می‌توان به طور قابل اعتماد، به شیوه‌ای که شواهد را تقویت کند به جای آنکه از ارزش آنها بکاهد، شناسایی کرد. در حالی شبیه به مثال ۱، یک اثرگزینش وجود دارد، اما در حد قابل قبولی روش است که مجموعه امکاناتی که در آن این گزینش رخ داده چیست، که تصحیح m -مقدار را مجاز می‌نماید. در دیگر مثالها، یک اثرگزینش وجود دارد، اما ممکن است روش نباشد که تصحیح را چگونه باید انجام داد. کوتاه سخن آنکه، بسیار غیرمعقول خواهد بود که امکان یادگیری چیزی جدید از داده‌ها در جهتی کلاً نامتنظر را انکار کنیم، اما باید بین بافتارهای مختلف به منظور کسب راهنمایی درباره اینکه چه تحلیل دیگری ممکن است لازم باشد، تمایز قائل شد.

۵. تذکرات پایانی

استدلال کردیم که احتمالهای خطأ در آزمونهای بسامدگرا را می‌توان برای ارزیابی قابلیت اعتماد یا ظرفیت آزمون در تمیزدادن این موضوع به کار برد که آیا فرایند واقعی پدیدآورنده داده‌ها با آنچه در H_0 توصیف شده همخوان است یا نیست. با آگاهی از این ظرفیت «گواه آوری» می‌توان مشخص کرد که شواهدی قوی علیه H_0 ، برپایه اصل بسامدگرای ش ب که شرح دادیم، وجود دارد یا خیر. چیزی که باعث می‌شود این نوع استدلال فرضی در مورد درستی بررسی به کار آید، نرخهای پایین خطای درازمدت در استفاده از این ابزار (یا آزمون) به این روای نیست؛ بلکه این است که آن نرخهای خطأ چه چیزی را درباره منبع یا پدیده زاینده داده‌ها آشکار می‌سازد. سعی نکرده‌ایم که رابطه بین تحلیلهای بسامدگرا و بیزی از موضوعات بسیار مشابه را مورد بحث قرار دهیم. یک اصل بینایی مفهوم استنباط استقرایی که خیلی با فلسفه بسامدگرا سازگار است، آن است که استنباط استقرایی نیاز به اقامه استدلالها و استنباطهای قاطع از طریق کنار هم گذاشتن چند خرد نتیجه متفاوت دارد؛ ملاحظاتی را برای هدایت این خرد نتیجه‌ها شرح دادیم. البته بیچیدگی موضوعها شرح شسته‌رفته‌ای از آنها را مشکل می‌سازد، شرحی که شخص می‌توانست ارائه کند اگر الگوریتمی واحد، کل استنباط استقرایی را در بر می‌گرفت. با این حال، حاصل کارگزارشی است که به نوع استدلالهایی که داشمندان بنا می‌نهند تا دانش و درک قابل اعتمادی از یک رشته علمی کسب کنند نزدیک می‌شود.

عملگرایانه که از «نوفه بسیار زیاد در شبکه» برهیز شود، موجه است. توجه عمدۀ معطوف به این است که چگونه باید نرخهای خطأ را به بهترین وجه تعديل کرد تا به مؤثرترین صورت بیانگر فرضهای زنی قابل پیگیری باشند. در این صورت، تحلیل مبتنی بر خطأ از موضوعها از طریق نزخ کشف نادرست است، یعنی از طریق نسبت درازمدت جایگاههایی که به عنوان مثبت برگزیده شده‌اند و در آنها اثری وجود ندارد. یک فرمول‌بندی بدیل از طریق مدل بیز تجربی انجام می‌شود و نتیجه‌گیری‌های حاصل از آن را می‌توان به نزخ کشف نادرست ربط داد. روش اخیر ممکن است مرجع باشد زیرا یک نزخ خطای خاص برای هر

زن انتخاب شده می‌توان یافت؛ در برخی موارد محتمل است که شواهد بسیار قوی‌تر از سایر موارد باشند و این تمایز در نزخ کشف نادرست سراسری محو می‌شود. برای ملاحظه بررسی نظاممندی از این موضوع، ر.ک. شافر [۲۸].

مثال ۴ (با تعریف آزمون). اگر آزمونها با مشخصات متفاوت اجرا شوند، و آزمونی بزرگزیده شود که معنی داری آماری کرانگین را می‌دهد، آنگاه تعديل به خاطر گزینش لازم است، هر چند که تعیین تعديل دقیق ممکن است مشکل باشد. اگر شخصی تأثیر گذاشتند نتیجه بر انتخاب مشخصات آزمون را مجاز بداند، شیوه‌ای را که منجر به m -مقدارشده تغییر می‌دهد، و این کار ممکن است غیرقابل قبول باشد. در حالی که موضوع اساسی و فرض بی‌تغییر باقی می‌مانند تعیین مشخصات دقیق مدل احتمالی با تحلیل مقدماتی داده‌ها به گونه‌ای هدایت شده است که سازوکار تصادفی‌ای را که عملًا مسئول برآمد آزمون بوده، تغییر دهد.

این وضع مانند آزمون توانایی یک تک‌تیرانداز است از این طریق که به سوی صفحه‌ای تیراندازی کند و سپس دایرة هدف به دور جایی که بیشترین تعداد اصابات به آن صورت گرفته رسم شود، و به اصطلاح اصل نشانه‌زن تکریسی به کار گرفته شود. مهارتی که شخص به ظاهر آزمون می‌کند و درباره آن استنباط انجام می‌دهد عبارت است از توانایی وی در زدن به هدف وقتی هدف داده شده و ثابت است، در حالی که آن همان مهارتی نیست که عملًا باعث به دست آمدن امتیاز زیاد شده است.

بر عکس، اگر انتخاب مشخصات نه براساس ملاحظات معنی داری آماری انحراف از فرض صفر، بلکه به آن دلیل صورت گیرید که با توجه به داده‌ها تغییراتی به‌منظور دستیابی به خطی بودن یا ثابت بودن واریانس خطأ مجاز باشد، به‌نظر می‌رسد که هیچ تعديلی به خاطر گزینش لازم نیست. بلکه دقیقاً خلاف آن: انتخاب مشخصات مناسب‌تر از لاحظه تجربی، اطمینان می‌دهد که m -مقدار حساب شده برای تفسیری از شواهد که قابل اعتماد باشد، مناسب دارد. (مایو و اسپانوس [۱۹]). این کار را می‌توان به طور رسمی‌تر با درنظر گرفتن انتخاب مشخصات به عنوان یک تحلیل ماکسیم درست‌نمایی توجیه کرد، یعنی ماکسیم‌سازی نسبت به پارامترهایی که فرض صفر مورد نظر را مشخص می‌کنند.

مثال ۵ (داده‌کاوی). این مثال شبیه مثال ۱ است، اگرچه چگونگی انجام تعديل به خاطر گزینش ممکن است روش نباشد زیرا شیوه به کار گرفته در گزینش متغیر ممکن است پر پیچ و خم باشد. در اینجا نیز برای اجتناب از مشکلات گزارش گزینشی، به جای انتخاب تها یک مدل، همه مدل‌های سازگار با داده‌ها را که در حد معقولی ساده‌اند مشخص می‌کنیم (کاکس و استل [۷]). بخشی از مشکلات اجرای یک چنین راهبردی محاسباتی اند نه مفهومی. مثالهایی از

مراجع

19. MAYOAYO, D. G. AND SPANOS, A. (2006). Severe testing as a basic concept in a Neyman-Pearson philosophy of induction. *British Journal of Philosophy of Science* **57**, 323-357. MR2249183
20. MILL, J. S. (1988). *A System of Logic*, Eighth edition. Harper and Brother, New York.
21. MORRISON, D. AND HENKEL, R. (eds.) (1970). *The Significance Test Controversy*. Aldine, Chicago.
22. NEYMAN, J. (1955). The problem of inductive inference. *Comm. Pure and Applied Maths* **8**, 13-46. MR0068145
23. NEYMAN, J. (1957). Inductive behavior as a basic concept of philosophy of science. *Int. Statist. Rev.* **25**, 7-22.
24. PEARSON, E. S. (1955). Statistical concepts in their relation to reality. *J. R. Statist. Soc. B* **17**, 204-207. MR0076234
25. PIERCE, C. S. [1931-5]. *Collected Papers*, Vols. 1-6, Hartshorne and Wiess, P. (eds.). Harvard University Press, Cambridge. MR0110632
26. POPPER, K. (1959). *The Logic of Scientific Discovery*. Basic Books, New York. MR0107593
27. SAVAGE, L. J. (1964). The foundations of statistics reconsidered. In *Studies in Subjective Probability*, Kyburg, H. E. and H. E. Smokler (eds.). Wiley, New York, 173-188. MR0179814
28. SHAFFER, J. P. (2005). This volume.
29. SPANOS, A. (2000). Revisiting data mining: 'hunting' with or without a license. *Journal of Economic Methodology* **7**, 231-264.
30. WHEWELL, W. [1874] (1967). *The Philosophy of the Inductive Sciences. Founded Upon Their History*, Second edition, Vols. 1 and 2 Reprint. Johnson Reprint, London.
31. WILL, C. (1993). *Theory and Experiment in Gravitational Physics*. Cambridge University Press. MR0778909
32. YATES, F. (1951). The influence of *Statistical Methods for Research Workers* on the development of the science of statistics. *J. Amer. Statist. Assoc.* **46**, 19-34.

- Deborah G. Mayo and D. R. Cox, "Frequentist statistics as a theory of inductive inference", *2nd Lehmann Symposium—Optimality*, IMS Lecture Notes-Monograph Series, **49** (2006) 77-97.

* دیبورا جی. مایو، بخش فلسفه و اقتصاد، ویرجینیا تک، آمریکا

mayod@vt.edu

** دی. آر. کاکس، کالج نافیلد، آکسفورد، انگلستان

david.cox@nuffield.ox.ac.uk

1. BIRNBAUM, A. (1997). The Neyman-Pearson Theory as decision theory, and as inference theory; with a criticism of the Lindley-Savage argument for Bayesian theory. *Synthese* **36**, 19-49. MR0652320
2. CARNAP, R. (1962). *Logical Foundations of Probability*. University of Chicago Press. MR0184839
3. COCHRAN, W. G. (1965). The planning of observational studies in human populations (with discussion). *J. R. Statist. Soc. A* **128**, 234-265.
4. COX, D. R. (1958). Some problems connected with statistical inference. *Ann. Math. Statist.* **29**, 357-372. MR0094890
5. COX, D. R. (1977). The role of significance tests (with discussion). *Scand. J. Statist.* **4**, 49-70. MR0448666
6. COX, D. R. AND HINKLEY, D. V. (1974). *Theoretical Statistics*. Chapman and Hall, London. MR0370837
7. COX, D. R. AND SNELL, E. J. (1974). The choice of variables in observational studies. *J. R. Statist. Soc. C* **23**, 51-59. MR0413333
8. DE FINETTI, B. (1974). *Theory of Probability*, 2 vols. English translation from Italian. Wiley, New York.
9. FISHER, R. A. (1935a). *Design of Experiments*. Oliver and Boyd, Edinburgh.
10. FISHER, R. A. (1935b). The logic of inductive inference. *J. R. Statist. Soc.* **98** 39-54.
11. GIBBONS, J. D. AND PRATT, J. W. (1975). *P-values: Interpretation and methodology*. *American Statistician* **29**, 20-25.
12. JEFFREYS, H. (1961). *Theory of Probability*, Third edition. Oxford University Press. MR0187257
13. KEMPTHORNE, O. (1976). Statistics and the philosophers. In *Foundations of Probability Theory, Statistical Inference, and Statistical Theories of Science*. Harper and Hooker (eds.) Vol. 2, 273-314. MR0488407
14. KEYNES, J. M. [1921] (1952). *A Treatise on Probability*. Reprint. St. Martin's Press, New York. MR1113699
15. LEHMANN, E. L. (1993). The Fisher and Neyman-Pearson theories of testing hypotheses: One theory or two? *J. Amer. Statist. Assoc.* **88**, 1242-1249. MR1245356
16. LEHMANN, E. L. (1995). Neyman's statistical philosophy. *Probability and Mathematical Statistics* **15**, 29-36. MR1369789
17. MAYO, D. G. (1996). *Error and the Growth of Experimental Knowledge*. University of Chicago Press.
18. MAYO, D. G. AND M. KRUSE (2001). Principles of inference and their consequences. In *Foundations of Bayesianism*, D. Cornfield and J. Williamson (eds.). Kluwer Academic Publishers, Netherlands, 381-403. MR1889643