

دانشگاه رازی

دانشکده علوم
گروه آمار

پایان نامه کارشناسی ارشد آمار ریاضی

کران‌های بهبود یافته برای معیار کولیک – لیبلر بر اساس ترکیب مدل‌های رقابتی

توسط
طیبه کریمی

استاد راهنما
دکتر عبدالرضا سیاره

استاد مشاور

تیر ماه ۱۳۹۰

کلیه حقوق اعم از چاپ و تکثیر، نسخه‌برداری، ترجمه، اقتباس و ... از
این پایان‌نامه برای دانشگاه رازی محفوظ است. نقل مطالب با ذکر
مأخذ بلامانع است.

کران‌های بهبود یافته برای معیار کولبک – لیبلر بر اساس ترکیب مدل‌های رقابتی چکیده

یکی از مفاهیم بنیادی در استنباط آماری انتخاب مدل مناسب برای یک مجموعه از داده‌ها است. هنگامی که مجموعه‌ای از داده‌ها در اختیار ما قرار می‌گیرند چگالی مولد این داده‌ها یعنی تابع چگالی درست داده‌ها مجهول است. لذا با مجموعه‌ای از مدل‌های رقابتی روبرو خواهیم بود. انتخاب یک مدل قطعی از بین این مدل‌های رقابتی که براساس تعداد محدودی از مشاهدات پیشنهاد شده‌اند، به عنوان برآورده از چگالی درست جامعه موجب بروز ریسک در انتخاب مدل برای جامعه خواهد شد. به همین دلیل انتخاب مدل بهینه از بین این مدل‌های رقابتی هدف اصلی انتخاب مدل است. معیارها و آزمون‌های مختلفی برای انتخاب مدل بهینه معرفی شده‌اند. معیار واگرایی کولبک – لیبلر با کاربرد گسترده در ساختار این معیارها و آزمون‌ها مورد توجه ما است. چون معیار کولبک – لیبلر واگرا است، پیدا کردن کران‌های مناسب بالایی و پایینی برای این معیار، کمک خواهد کرد تا مدل مناسبی به عنوان برآورد مدل درست انتخاب شود. در این پایان نامه به بررسی کران‌های موجود برای معیار اطلاع کولبک – لیبلر پرداخته شده است. سپس خواص شکل نمایی این معیار مورد بررسی قرار گرفته و به کمک این خواص و ترکیب محدب مدل‌های رقابتی، این معیار بهبود داده شده است. همچنین به بررسی و مقایسه معیار اطلاع کولبک – لیبلر میانگین‌های حسابی، هارمونیک و هندسی مدل‌های رقابتی پرداخته شده است.

واژه‌های کلیدی : انتخاب مدل، ترکیب محدب، خاصیت زیر جمعی، معیار اطلاع، معیار کولبک – لیبلر، میانگین هندسی وزنی.

فهرست مندرجات

۱	۱	تعاریف و مفاهیم اولیه
۲	۱-۱	مقدمه
۳	۱-۲	تعاریف و قضایا
۵	۳-۱	مدل‌های آماری
۶	۴-۱	معیار اطلاع کولبک - لیبلر
۷	۱-۴	ویژگی‌های معیار اطلاع کولبک - لیبلر
۱۱	۱-۵	معیار اطلاع کولبک - لیبلر پارامتری
۱۲	۱-۶	معیار اطلاع آکائیک (AIC)
۱۳	۷-۱	معیارهای دیگر
۱۵	۲	کران‌هایی برای معیار اطلاع کولبک - لیبلر
۱۶	۱-۲	مقدمه

۱۶	۲-۲ کرانهایی برای معیار اطلاع کولبک - لیبلر
۲۳	۱-۲-۲ کرانهایی دیگر برای معیار کولبک - لیبلر
۲۹	۳ خواص شکل نمایی معیار اطلاع کولبک - لیبلر
۳۰	۱-۳ مقدمه
۳۰	۲-۳ میانگین هندسی وزنی و خواص آن
۳۲	۳-۳ برخی خواص $\exp[-KL(p,.)]$
۳۴	۴-۳ خاصیت‌های دیگر شکل نمایی معیار کولبک - لیبلر براساس ترکیب‌هایی از مدل‌های رقابتی
۳۷	۵-۳ شبیه سازی و مثال عددی
۲۹	۴ بهبود معیار کولبک - لیبلر براساس ترکیب مدل‌های رقابتی
۴۰	۱-۴ مقدمه
۴۰	۲-۴ ترکیب محدب مدل‌های رقابتی
۴۴	۱-۲-۴ ترکیب محدب k مدل رقابتی
۴۸	۳-۴ میانگین هندسی وزنی مدل‌های رقابتی
۴۹	۱-۳-۴ میانگین هندسی k مدل رقابتی
۴۹	۴-۴ میانگین هارمونیک مدل‌های رقابتی
۵۰	۱-۴-۴ میانگین هارمونیک k مدل رقابتی

۵۲ ۵-۴ شبیه سازی و مثال عددی

۵۶ ۶-۴ نتیجه گیری

فصل ١

تعاريف و مفاهيم أوليه

مدل‌های آماری نقش بسیار مهمی در تحلیل و آنالیز داده‌ها دارند. در دنیای واقعی با رویدادهای پیچیده‌ای روبرو هستیم که برای تحلیل روند آنها و پیش‌بینی رفتار آینده باید بتوانیم الگوی آنها را بیابیم. در علم آمار این الگوها را با مدل‌های آماری نشان می‌دهند. مدل‌های آماری برای نشان دادن ساختارهای تصادفی، پیش‌بینی رفتار آینده و استخراج اطلاعات مهم از داده‌ها مورد استفاده قرار می‌گیرند.

آکائیک^۱ (۱۹۷۴ و ۱۹۸۵) این نقطه نظر را مطرح کرد که مدل سازی آماری فقط پیدا کردن مدلی نیست که رفتار داده‌های مشاهده شده را شرح دهد، بلکه هدف اصلی آن پیش‌بینی تا حد امکان خوب آینده فرایند تحت بررسی است. یک نقطه نظر دیگر در مدل‌های آماری استخراج اطلاعات از داده‌ها است. در بسیاری از استنباط‌های آماری فرض را براین می‌گذارند که داده‌ها از یک مدل مشخص پیروی می‌کنند. بر همین اساس یک مدل آماری پارامتری را پیشنهاد کرده و به دنبال برآورده مدل درست مشاهدات هستند. اما در دنیای واقعی مشاهدات و رخدادها تحت تاثیر متغیرهای زیادی هستند که بعضی از آنها قابل شناسایی نیستند. پس پیدا کردن مدلی که به طور کاملاً دقیق توزیع مشاهدات را نشان دهد کار بسیار پیچیده و گاهی غیر ممکن است. لذا سعی می‌شود مدلی مورد استفاده قرار گیرد که در میان مدل‌های رفاقتی که برای برآورد مدل درست پیشنهاد شده‌اند، مدل بهینه باشد.

یک مدل خوب به صورت یکتا قابل تعیین کردن نیست و می‌تواند بسته به نظر آماردان یا اطلاعات موجود متفاوت باشد. به عبارت دیگر هدف مدل سازی آماری ساختن یک مدل یکتا نیست بلکه ساختن یک مدل خوب تحت شرایط موجود است.

معیارهای متفاوتی برای انتخاب مدل معرفی شده‌اند، یکی از معیارهایی که کاربرد فراوانی در این زمینه دارد معیار اطلاع کولبک – لیبلر^۲ است. اطلاع کولبک – لیبلر در سال ۱۹۵۱ توسط سالمون کولبک^۳ و ریچارد لیبلر^۴ برای سنجش میزان نزدیکی مدل انتخابی به مدل درست معرفی شد. از آنجا که این معیار کاربرد گسترده‌ای در انتخاب مدل دارد، انتخاب حدود بالایی و پایینی مناسب برای این معیار در انتخاب یک مدل مناسب برای داده‌ها به ما کمک می‌کند.

در این پایان نامه به بررسی کران‌هایی برای اطلاع کولبک – لیبلر بخصوص کران‌های بالایی برای آن پرداخته شده و سپس برخی خواص نمایی این معیار بررسی شده است. همچنین این معیار را در حالتی

^۱ Akaike

^۲ Kullback-Leibler Information

^۳ Solomon Kullback

^۴ Richard Leibler

که مدل‌های رقابتی ترکیب شده‌اند مورد بررسی قرار داده و مقایسه‌ای بین ترکیب‌های مختلف به عمل آمده است.

در فصل اول تعاریف، قضایا و نامساوی‌های مورد نیاز فصل‌های بعد و معیار اطلاع کولبک – لیبلر معرفی می‌شوند. در فصل دوم کران‌های پایینی و بالایی برای اطلاع کولبک – لیبلر آورده شده است. در فصل سوم خواص شکل نمایی اطلاع کولبک – لیبلر بیان می‌شود و در فصل چهارم ترکیب محدب مدل‌های رقابتی را به عنوان مدل پیشنهادی معرفی کرده و اطلاع کولبک – لیبلر در این حالت بررسی می‌شود و مقایسه‌ای بین ترکیب‌های مختلف مدل‌های رقابتی به عمل خواهد آمد.

۱-۲ تعاریف و قضایا

در این بخش برخی تعاریف، قضایا و نامساوی‌های مورد نیاز فصل‌های بعدی بیان شده است.

تعريف ۱.۱ میانگین هارمونیک^۵

فرض کنید $a = (a_1, a_2, \dots, a_n)$ برداری از اعداد حقیقی مثبت باشد، آنگاه میانگین هارمونیک به صورت

$$H(a_1, a_2, \dots, a_n) = \frac{n}{\sum_{i=1}^n \frac{1}{a_i}}$$

تعریف می‌شود.

تعريف ۲.۱ میانگین هندسی وزنی^۶

فرض کنید $p = (p_1, p_2, \dots, p_n)$ توزیع احتمال و $a = (a_1, a_2, \dots, a_n)$ برداری از اعداد حقیقی نامنفی باشند، میانگین هندسی وزنی که با نماد $G_n(p, a)$ نشان داده می‌شود، به صورت

$$G_n(p, a) = \prod_{i=1}^n a_i^{p_i}$$

تعریف می‌شود.

^۵ Harmonic Mean

^۶ Weighted Geometric Mean

تعريف ۳.۱ تابع محدب^۷

تعاریف مختلفی از توابع محدب وجود دارد که به سه نوع از این تعاریف اشاره می‌کنیم:

الف) تابع پیوسته^(.) g که دامنه و برد آن اعداد حقیقی است، محدب نامیده می‌شود اگر برای هر

$x \in \mathbb{R}$ ، خطی مانند^(.) $l(x)$ وجود داشته باشد که از نقطه^(.) $(x_0, g(x_0))$ بگذرد و رابطه

$$l(x) \leq g(x), \quad \forall x \in \mathbb{R} \quad (1.2.1)$$

برقرار باشد.

ب) برای مقادیر حقیقی تابع اکیداً محدب^(.) g که روی فاصله حقیقی I تعریف می‌شود نامساوی

$$g(b) - g(a) \geq g'(a)(b - a) \quad (2.2.1)$$

برای هر $a, b \in I$ برقرار است. هرگاه $b = a$ باشد در نامساوی فوق تساوی رخ می‌دهد.

ج) تابع پیوسته^(.) g اکیداً محدب است اگر مشتق دوم آن همواره مثبت باشد.

قضیه ۱.۱ . نامساوی جنسن^۸

اگر X متغیری تصادفی با میانگین $E[X]$ و^(.) g تابعی محدب باشد، آنگاه

$$E[g(X)] \geq g(E[X]). \quad (3.2.1)$$

برهان. اگردر تعریف الف (۱.۳) فرض کنیم $x_0 = E[X]$ با توجه به تحدب و پیوستگی^(.) g خطی

مانند^(.) $l(x) = a + bx$ وجود دارد به طوری که

$$l(E[X]) = g(E[X]) \quad (4.2.1)$$

لذا طبق رابطه (۱.۲.۱)

^۷ Convex Function

^۸ Jensen Inequality

$$l(x) = a + bx \leq g(x), \quad \forall x \in \mathbb{R} \quad (5.2.1)$$

از طرفی

$$E[l(X)] = E[a + bX] = a + bE[X] = l(E[X]) \quad (6.2.1)$$

لذا طبق رابطه فوق و رابطه های (4.2.1) و (5.2.1) نامساوی

$$g(E[X]) = l(E[X]) = E[l(X)] \leq E[g(X)] \quad (7.2.1)$$

نتیجه می شود.

تذکر: برای اثبات نامساوی (7.2.1) از این ویژگی امید ریاضی، که اگر $f(x) \leq h(x)$ آنگاه

□ استفاده شده است. $E[f(X)] \leq E[h(X)]$

۱-۳ مدل‌های آماری

فرض می‌کنیم $\{X_1, X_2, \dots, X_n\}$ یک نمونه تصادفی *i.i.d* از توزیعی با تابع چگالی $h(x)$ باشد. تابع چگالی درست یا مدل درست مشاهدات است. این توزیع مجهول است و برآورد آن یکی از مسائل اساسی در آمار است. با توجه به ساختار داده‌ها، مدلی مانند $F(x)$ برای تقریب یا برآورد توزیع درست معرفی می‌شود. تابع چگالی احتمال مرتبط با این مدل با $f(x)$ نشان داده می‌شود. اگر هدف برآورد مدل درست باشد معمولاً مدل پارامتری $F_\theta = \{f(x; \theta), \theta \in \Theta \subset \mathbb{R}^p\}$ به عنوان یک مدل رقبتی انتخاب خواهد شد. با برآورد پارامترهای مدل رقبتی، برآوردهای برای مدل درست به دست می‌آید. مدل‌های آماری نسبت به هم می‌توانند آشیانه‌ای^۹، غیر آشیانه‌ای^{۱۰} و متداخل^{۱۱} باشند.

تعریف ۴.۱ مدل‌های غیر آشیانه‌ای.

فرض کنید دو مدل $G_\gamma = \{g(x; \gamma), \gamma \in \Gamma \subset \mathbb{R}^q\}$ و $F_\theta = \{f(x; \theta), \theta \in \Theta \subset \mathbb{R}^p\}$ مدل‌های رقبتی

^۹ Nested

^{۱۰} Non - nested

^{۱۱} Overlap

برای چگالی درست باشند، این دو مدل را نسبت به هم غیرآشیانه‌ای گوییم اگر و تنها اگر $F_\theta \cap G_\gamma = \emptyset$.

تعريف ۵.۱ مدل‌های آشیانه‌ای.

گوییم مدل F_θ در مدل G_γ آشیانه دارد هرگاه

تعريف ۶.۱ مدل‌های متداخل.

دو مدل F_θ و G_γ را متداخل گوییم اگر و تنها اگر:

$$F_\theta \cap G_\gamma \neq \emptyset \quad (1)$$

$$G_\gamma \not\subseteq F_\theta \text{ و } F_\theta \not\subseteq G_\gamma \quad (2)$$

مدل‌های آماری را می‌توان در دو دسته طبقه‌بندی کرد: مدل‌های خوب – توصیف شده^{۱۲} و مدل‌های بد – توصیف شده^{۱۳}.

تعريف ۷.۱ مدل خوب – توصیف شده.

مدل $F_\theta = \{f(x; \theta), \theta \in \Theta \subset \mathbb{R}^p\}$ را خوب – توصیف شده گوییم هرگاه یک $\theta \in \Theta$ وجود داشته باشد که

$h(x) = f(x; \theta_0)$. به عبارت دیگر $h \in F_\theta$.

تعريف ۸.۱ مدل بد – توصیف شده.

مدل F_θ را بد – توصیف شده گوییم هرگاه $\theta \in \Theta$ وجود نداشته باشد که $h \in F_\theta$. به عبارت دیگر

برای هر $\theta \in \Theta$ $h(x) \neq f(x; \theta)$.

۱-۴ معیار اطلاع کولبک – لیبلر

معیار اطلاع کولبک – لیبلر یک معیار برای ارزیابی میزان نزدیکی مدل‌های آماری به توزیع درست جامعه است. معیار اطلاع KL در واقع میزان ریسک احتمالی انتخاب یک مدل رقابتی به جای توزیع درست است. هرچه این معیار کوچکتر باشد مدل انتخابی مدلی مناسب‌تر خواهد بود.

فرض می‌کنیم که $f(x)$ یک مدل آماری ساخته شده بر اساس مشاهدات باشد. می‌خواهیم نزدیکی این مدل را به تابع چگالی (تابع جرم احتمال) درست $h(x)$ بر اساس معیار اطلاع KL ارزیابی کنیم. این معیار به صورت

^{۱۲} Well-Specified

^{۱۳} Miss-Specified

$$KL(h, f) = E_h \left[\log \left\{ \frac{h(x)}{f(x)} \right\} \right]$$

است که در آن E_h امید ریاضی تحت تابع چگالی یاتابع جرم احتمال h است.
اگر h یک تابع چگالی پیوسته باشد این معیار به صورت

$$KL(h, f) = \int_{-\infty}^{+\infty} \log \left\{ \frac{h(x)}{f(x)} \right\} h(x) dx \quad (8.4.1)$$

تعریف می شود. اگر h یک تابع احتمال باشد آنگاه

$$KL(h, f) = \sum_{i=1}^n f(x_i) \log \left\{ \frac{h(x_i)}{f(x_i)} \right\}. \quad (9.4.1)$$

۱-۴-۱ ویژگی های معیار اطلاع کولبک - لیبلر

معیار اطلاع کولبک - لیبلر به عنوان معیار واگرایی بین دو توزیع احتمال معرفی شد. این معیار یک متر بر روی فضای توزیع های احتمال نیست زیرا:

- ۱) اطلاع کولبک - لیبلر متقارن نیست یعنی

$$KL(h, f) \neq KL(f, h)$$

۲) این معیار در نامساوی مثلثی^{۱۴} صدق نمی کند.

$$KL(h, f) \not\leq KL(h, g) + KL(g, f)$$

اطلاع کولبک - لیبلر دارای ویژگی های زیر است:

$$KL(h, f) \geq 0. \quad (1)$$

$$\cdot h(x) = f(x) \text{ اگر و تنها اگر } KL(h, f) = 0. \quad (2)$$

برهان.

برای بررسی دو ویژگی بالا، تابع $k(t) = \log t - t + 1$ را که برای $t > 0$ تعریف شده است در نظر می گیریم. در این صورت مشتق $(t)^k$ برابر است با

$$k'(t) = \frac{1}{t} - 1,$$

^{۱۴} Triangle Inequality

که $k'(1) = 0$ و همچنین

$$k''(t) = -\frac{1}{t^2} < 0,$$

در نتیجه $k(t)$ ماکسیمم مقدار خود را در نقطه یک اختیار می کند، $0 = k(1)$.

بنابراین

$$k(t) \leq k(1) = 0, \quad \forall t > 0$$

لذا

$$\log(t) \leq t - 1, \quad \forall t > 0$$

و تساوی زمانی برقرار است اگر و فقط اگر $t = 1$ باشد.

اثبات ویژگی ۱:

الف) برای حالت پیوسته با جایگذاری $t = \frac{f(x)}{h(x)}$ داریم:

$$\log\left(\frac{f(x)}{h(x)}\right) \leq \left(\frac{f(x)}{h(x)} - 1\right),$$

با ضرب $h(x)$ در دو طرف نامساوی فوق و انتگرال گرفتن خواهیم داشت:

$$\begin{aligned} \int_{\mathbb{R}} \log\left\{\frac{f(x)}{h(x)}\right\} h(x) dx &\leq \int_{\mathbb{R}} \left\{\frac{f(x)}{h(x)} - 1\right\} h(x) dx \\ &= \int_{\mathbb{R}} \{f(x) - h(x)\} dx \\ &= \int_{\mathbb{R}} f(x) dx - \int_{\mathbb{R}} h(x) dx \end{aligned} \quad (10.4.1)$$

چون انتگرال هر تابع چگالی روی فضای متغیر تصادفی X برابر یک است در نتیجه طرف دوم رابطه

فوق صفر خواهد شد، بنابراین:

$$\int_{\mathbb{R}} \log\left\{\frac{f(x)}{h(x)}\right\} h(x) dx \leq 0,$$

و

$$\begin{aligned} KL(h, f) &= \int_{\mathbb{R}} \log\left\{\frac{h(x)}{f(x)}\right\} h(x) dx \\ &= - \int_{\mathbb{R}} \log\left\{\frac{f(x)}{h(x)}\right\} h(x) dx \geq 0 \end{aligned} \quad (11.4.1)$$

λ

ب) در حالت گستته تابع جرم احتمال به جای تابع چگالی و سیگما به جای انتگرال را می‌توان در نظر گرفت.

اثبات ویرگی ۲:

همانطور که گفته شد در نامساوی $t - \log(t)$ تساوی برقرار است اگر و تنها اگر $t = 1$ باشد. با توجه با اینکه در اثبات ویرگی ۱، $t = \frac{f(x)}{h(x)}$ را فرض کردیم ویرگی ۲ برقرار است.

براساس ویرگی‌های فوق در می‌یابیم که هر چه معیار کولبک – لیبلر کوچک‌تر باشد مدل مورد نظر به چگالی درست مشاهدات نزدیک‌تر است و تنها در صورتی برابر صفر خواهد شد که مدل رقابتی معادل چگالی درست مشاهدات باشد.

در واقع $\{\frac{h(X)}{f(X)}\}$ تابع زیان حاصل از انتخاب مدل $f(x)$ به جای مدل درست $h(x)$ است و امید ریاضی آن نسبت به چگالی درست $h(x)$ ریسک کولبک – لیبلر خواهد بود. این ریسک وقتی به صفر نزدیک‌تر می‌شود که مدل ساخته شده تقریباً معادل چگالی درست مشاهدات باشد. لذا اطلاع کولبک – لیبلر، ریسک کولبک – لیبلر نیز نامیده می‌شود.

در مثال‌های زیر معیار کولبک – لیبلر برای چند مدل پیوسته و گستته بدست آمده است.

معیار کولبک – لیبلر برای حالت پیوسته

مثال ۱.۱ فرض کنید مدل درست h و مدل رقابتی f به ترتیب دارای توزیع $N(\mu, \sigma^2)$ و $N(\zeta, \tau^2)$ در نظر گرفته شده باشند. برای توزیع $N(\mu, \sigma^2)$ نظر گرفته شده باشند. برای توزیع $N(\zeta, \tau^2)$

$$f(x) = (2\pi\sigma^2)^{-1/2} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$$

لذا

$$\begin{aligned} E_h[\log f(X)] &= E_h\left[-\frac{1}{2}\log(2\pi\sigma^2) - \frac{(X-\mu)^2}{2\sigma^2}\right], \\ &= -\frac{1}{2}\log(2\pi\sigma^2) - \frac{\tau^2 + (\zeta-\mu)^2}{2\sigma^2}. \end{aligned} \quad (12.4.1)$$

و

$$E_h[\log h(X)] = -\frac{1}{2}\log(2\pi\tau^2) - \frac{1}{2}. \quad (13.4.1)$$

در نتیجه اطلاع کولبک – لیبلر به صورت

$$\begin{aligned}
KL(h, f) &= E_h[\log h(X)] - E_h[\log f(X)] \\
&= \frac{1}{2} \left\{ \log\left(\frac{\sigma^2}{\tau^2}\right) + \frac{\tau^2 + (\zeta - \mu)^2}{\sigma^2} - 1 \right\}
\end{aligned} \tag{14.4.1}$$

خواهد بود.

مثال ۲.۱ فرض کنید که مدل درست جامعه توزیع لابلانس با تابع چگالی $h(x) = \frac{1}{\sqrt{2\pi}} \exp(-|x|)$ و مدل رقلابتی $f(x)$ نرمال با میانگین μ و واریانس σ^2 باشد. در این حالت داریم:

$$\begin{aligned}
E_h[\log h(X)] &= -\log 2 - \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} |x| e^{-|x|} dx \\
&= -\log 2 - \int_0^{+\infty} x e^{-x} dx \\
&= -\log 2 - 1
\end{aligned} \tag{15.4.1}$$

$$\begin{aligned}
E_h[\log f(X)] &= -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} \int_{-\infty}^{+\infty} (x - \mu)^2 e^{-|x|} dx \\
&= -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (4 + 2\mu^2)
\end{aligned} \tag{16.4.1}$$

اطلاع کولبک – لیبلر برای مدل $f(x)$ تحت چگالی درست $h(x)$ برابر است با:

$$\begin{aligned}
KL(h, f) &= E_h[\log h(X)] - E_h[\log f(X)] \\
&= \frac{1}{2} \log(2\pi\sigma^2) + \frac{2 + \mu^2}{2\sigma^2} - \log 2 - 1
\end{aligned} \tag{17.4.1}$$

مثال ۳.۱ معیار کولبک – لیبلر برای حالت گسسته در پرتاب یک تاس سالم، احتمال آمدن هر یک از خالها برابر $\frac{1}{6}$ است. لذا تابع جرم احتمال درست به

صورت $\{ \frac{1}{7}, \frac{1}{7}, \frac{1}{7}, \frac{1}{7}, \frac{1}{7}, \frac{1}{7} \} = p$ است. دوتابع جرم احتمال برای مشاهده خالهای این تاس به شکل

زیر پیشنهاد شده است:

$$q_1 = \{ 0/20, 0/12, 0/18, 0/12, 0/20, 0/18 \}$$

و

$$q_2 = \{ 0/18, 0/12, 0/14, 0/19, 0/22, 0/15 \}$$

می خواهیم بررسی کنیم کدام مدل برای توصیف نتایج حاصل از پرتاب یک تاس سالم مناسب‌تر است.

با استفاده از تعریف اطلاع کولبک – لیبلر برای حالت گسسته داریم:

$$KL(p, q_1) = \sum_{i=1}^7 p_i \log\left\{\frac{p_i}{q_{1i}}\right\} = 0/023$$

و

$$KL(p, q_2) = \sum_{i=1}^7 p_i \log\left\{\frac{p_i}{q_{2i}}\right\} = 0/020$$

نامنفی بودن اطلاع کولبک – لیبلر در این مثال به وضوح دیده می‌شود. همچنین با توجه به مقادیر بدست آمده، احتمال‌های متناظر با مدل q_2 به مدل درست نزدیک‌تر است. بنابراین مدل q_2 به مدل q_1 ترجیح داده می‌شود.

۱-۵ معیار اطلاع کولبک – لیبلر پارامتری

معیار اطلاع کولبک – لیبلر در حالت پارامتری به صورت

$$KL(h, f_\theta) = E_h \left[\log \left(\frac{h(X)}{f(X; \theta)} \right) \right]$$

است. این معیار می‌تواند به صورت زیر تجزیه شود:

$$KL(h, f_\theta) = E_h [\log h(X)] - E_h [\log f(X; \theta)]$$

به وضوح دیده می‌شود که معیار اطلاع کولبک – لیبلر تابعی از پارامتر θ است. چون جمله اول سمت راست رابطه فوق یک مقدار ثابت است و فقط به چگالی درست و مجھول h بستگی دارد، بنابراین برای مقایسه مدل‌های رقابتی کافیست عبارت دوم سمت راست را که جمله مرتبط معیار اطلاع کولبک – لیبلر نامیده می‌شود، محاسبه کرد. از آنجایی که در بحث انتخاب مدل علاقمند به پیدا کردن θ ای

هستیم که معیار KL را مینیمم کند، کافیست θ ای بی را پیدا کنیم که جمله مرتبط را ماقسیم کند. توجه کنید جمله مرتبط در واقع امید لگاریتم درستنمایی است. بنابراین در مقایسه مدل‌های رقابتی مدلی را بر سایر مدل‌ها ترجیح می‌دهیم که امید لگاریتم درستنمایی بزرگتری داشته باشد.

تعريف ۹.۱ مقدار شبه پارامتر درست^{۱۵}

فرض کنید $F_\theta = \{f(x; \theta), \theta \in \Theta \subset \mathbb{R}^p\}$ مدل پیشنهاد شده باشد. پارامتر θ^* ای که معیار کولبک – لیبلر را مینیمم کند، مقدار شبه پارامتر درست نامیده می‌شود.

با توجه به مجھول بودن چگالی درست جامعه h و پارامتر θ ، مقدار امید لگاریتم درستنمایی یک کمیت مجھول است بنابراین باید برآورد شود تا بتوان از آن برای انتخاب مدل استفاده کرد.

۱-۶ معیار اطلاع آکائیک (AIC)

فرض کنید $\{X_1, X_2, \dots, X_n\}$ یک نمونه تصادفی از چگالی درست و مجھول h و یک خانواده F_θ از مدل‌های رقابتی برای برآورد h در نظر گرفته شده باشند. با استفاده از روش درستنمایی ماقسیم، برآورد پارامترها را پیدا کرده و چگالی $f(x; \hat{\theta}_n)$ تشکیل داده می‌شود. در اینجا هدف محاسبه و بررسی خوبی یا بدی مدل برآورد شده $f(z; \hat{\theta}_n)$ از لحاظ پیش‌بینی برای آینده است که در آن Z به عنوان مشاهده آینده از چگالی h و مستقل از X_i ها است. میزان نزدیکی $f(z; \hat{\theta}_n)$ به توزیع درست مشاهدات از اطلاع KL به صورت

$$\begin{aligned} KL(h(z), f(z; \hat{\theta}_n)) &= E_h \left[\log \left(\frac{h(Z)}{f(Z; \hat{\theta}_n)} \right) \right] \\ &= E_h [\log h(Z)] - E_h [\log f(Z; \hat{\theta}_n)], \end{aligned}$$

محاسبه می‌شود در اینجا Z از $\hat{\theta}_n$ مستقل است زیرا تابعی از X_i ها است. با توجه به ویژگی‌های معیار اطلاع KL بزرگ بودن $E_h [\log f(Z; \hat{\theta}_n)]$ نزدیکی چگالی $f(z; \hat{\theta}_n)$ به چگالی درست h را نشان می‌دهد. چون $E_h [\log f(Z; \hat{\theta}_n)]$ کمیتی مجھول است باید برآورد شود. یک برآورد از این کمیت را می‌توان با جایگزین کردن تابع توزیع تجربی \hat{H} به جای تابع توزیع درست و مجھول H ، به صورت زیر به دست آورد:

^{۱۵} Psedo-true Value of a Parameter

$$E_{\hat{H}} \left[\log f(Z; \hat{\theta}_n) \right] = \int \log f(Z; \hat{\theta}_n) d\hat{H}(x) = \frac{1}{n} \sum_{i=1}^n \log f(x_i; \hat{\theta}_n).$$

بنابراین $E_{\hat{H}} \left[\log f(Z; \hat{\theta}_n) \right] = \frac{1}{n} \sum_{i=1}^n \log f(x_i; \hat{\theta}_n)$ دارد.

یک معیار مهم و زیربنایی در انتخاب مدل است. در واقع این معیار از نظر تئوری و عملی مورد استفاده آمارشناسان و تحلیل گران داده‌ها قرار می‌گیرد. AIC یک معیار ارزیابی برای عدم برازش مدل‌هایی است که در آن پارامترها به روش درستنمایی ماکسیمم برآورده شده‌اند و نشان می‌دهد که اریبی لگاریتم درستنمایی متناسب با پارامترهای مدل (p) است. معیار AIC به صورت زیر تعریف می‌شود:

$$AIC = -2 \sum_{i=1}^n \log f(x_i; \hat{\theta}_n) + 2p.$$

در واقع آکائیک در سال ۱۹۷۴ بیان کرد که اگر چگالی مولدهای داده‌ها به مدل پارامتری پیشنهاد شده نزدیک باشد میزان اریبی بوسیله تعداد پارامترهای پیشنهاد شده تقریب زده می‌شود. محاسبه معیار AIC به چگالی درست و مجھول h بستگی ندارد.

۱-۷ معیارهای دیگر

در این بخش معیارهای هلینگر و خی دو که در فصل ۲ برای بررسی کران‌های معیار اطلاع کولبک – لیبلر استفاده می‌شود بیان شده است.

تعريف ۱۰.۱ معیار هلینگر^{۱۶}

فرض کنید $f(x)$ و $h(x)$ برای هر $x \in \chi$ دوتابع چگالی (احتمال) باشند. آنگاه این معیار به صورت

$$D_H(h, f) = \begin{cases} \frac{1}{2} \int_{\chi} (\sqrt{h(x)} - \sqrt{f(x)})^2 dx, & \text{در حالت پیوسته} \\ \frac{1}{2} \sum_{x \in \chi} (\sqrt{h(x)} - \sqrt{f(x)})^2. & \text{در حالت گسسته} \end{cases}$$

تعریف می‌شود.

تذکر: معیار هلینگر یک متر نیست اما متقارن است.

تعريف ۱۱.۱ معیار خن دو^{۱۷}

این معیار که معیاری نامتقارن است به دو صورت تعریف می‌شود.

(الف)

$$D_{\chi^r}(h, f) = \begin{cases} \int_{\chi} \frac{h^r(x) - f^r(x)}{f(x)} dx, & \text{در حالت پیوسته} \\ \sum_{x \in \chi} \frac{h^r(x) - f^r(x)}{f(x)}. & \text{در حالت گسسته} \end{cases}$$

(ب)

$$D_{\chi^r}(h, f) = \begin{cases} \int_{\chi} \frac{(h(x) - f(x))^r}{f(x)} dx, & \text{در حالت پیوسته} \\ \sum_{x \in \chi} \frac{(h(x) - f(x))^r}{f(x)}. & \text{در حالت گسسته} \end{cases}$$

اگر تعاریف الف و ب بسط داده شوند رابطه‌های یکسانی بدست می‌آید.

در این فصل برخی تعاریف مورد نیاز فصل‌های بعد، مدل‌های آماری و معیارهای انتخاب مدل بیان

شد. در فصل بعد برخی کران‌های مهم معیار اطلاع کولبک – لیبلر مورد بررسی قرار خواهد گرفت.