



دانشکده علوم ریاضی

پایان نامه جهت اخذ درجه کارشناسی ارشد
در رشته آمار ریاضی

عنوان :

رفتار مجانبی بر آوردگر ناپارامتری تابع رگرسیون در مدل
سانسور تصادفی با داده‌های α -آمیخته

استاد راهنما :

جناب آقای دکتر وحید فکور

استاد مشاور :

جناب آقای دکتر حسن دوستی

نگارش :

سمیه باوفا

پاییز ۱۳۹۰

فهرست مطالب

۴	پیش‌گفتار
۱۰	۱ تعاریف و مفاهیم مورد نیاز
۱۱	۱.۱ برآورد تابع چگالی احتمال
۱۲	۱.۱.۱ برآورد گر هیستوگرام
۱۳	۲.۱.۱ برآورد گر باکس کار
۱۴	۳.۱.۱ برآورد گر هسته‌ای تابع چگالی
۱۷	۲.۱ تحلیل بقا و سانسور
۲۲	۳.۱ وابستگی ضعیف و متغیرهای تصادفی آمیزنده
۲۷	۴.۱ رگرسیون ناپارامتری
۳۱	۵.۱ سایر تعاریف، قضایا و نابرابری‌های مورد نیاز
۴۰	۲ رفتار مجانبی برآورد گر هسته‌ای تابع رگرسیون در یک مدل سانسور تصادفی

۴۱	مقدمه	۱.۲
۴۲	ساختار برآوردگر	۲.۲
۴۴	پذیره‌ها و قضایا	۳.۲
۴۵	پذیره‌ها: ۱.۳.۲	
۴۸	اثبات قضایا	۴.۲
۶۱	شبه‌سازی	۵.۲

۳ برآوردگر هسته‌ای تابع رگرسیون در یک مدل سانسور تصادفی با داده‌های α -

۶۵	آمیزنده	
۶۶	مقدمه	۱.۳
۶۷	ساختار برآوردگر	۲.۳
۶۸	پذیره‌ها و نرخ همگرایی برآوردگر	۳.۳
۷۱	لم‌های مورد نیاز	۴.۳
۸۱	شبه‌سازی	۵.۳

۸۵ واژه‌نامه‌ی فارسی به انگلیسی

۹۰ واژه‌نامه‌ی انگلیسی به فارسی

۹۷ کتاب‌نامه

پیش‌گفتار

وقتی که صحبت از تحلیل رگرسیونی می‌شود، معمولاً منظور برازش یک الگوی ریاضی به داده‌ها به عنوان الگوی وابستگی متغیرها، بررسی نمودار باقیمانده‌ها به عنوان انحرافات از الگو، برآورد و آزمون فرض درباره پارامترها، پیش‌بینی و مانند آن می‌باشد. تحلیل رگرسیونی یکی از ابزارهای آماری است که در سطح بسیار وسیعی کاربرد دارد.

در رگرسیون پارامتری لازم است یک شکل از قبل تعیین شده مانند چندجمله‌ای را برای منحنی رگرسیون فرض کنیم. اما برای بسیاری از مجموعه داده‌ها در مورد شکل ارتباطی بین متغیرها دانش درستی نداریم.

در این مواقع بهتر است به جای تحلیل یک الگوی پارامتری به داده‌ها از فنونی استفاده کنیم که به داده‌ها اجازه می‌دهند رفتار الگوی خود را بروز دهند. رگرسیون ناپارامتری مجموعه‌ای از فنون برای برآورد منحنی رگرسیون است بدون آن که درباره شکل تابعی منحنی رگرسیون فرض‌های قوی بگذاریم. ریشه رگرسیون ناپارامتری را باید در برآوردگر ناپارامتری تابع چگالی احتمال جستجو کرد.

برآورد ناپارامتری تابع چگالی احتمال، با استفاده از مجموعه‌ی متناهی از مشاهدات هم‌توزیع، یکی از مسائل بنیادی و مهم در استنباط آماری است که نقش بسیار مهمی در تشخیص مدل و الگوی احتمال

جامعه دارد. روزن بلات و پارزن پیشگامان نظریه‌ی برآوردگرهای هسته‌ی چگالی می‌باشند. پس از آن دو، مسأله برآورد تابع چگالی مورد توجه بسیاری از محققان بوده است.

همان طور که گفتیم ریشه رگرسیون ناپارامتری را باید در برآوردگر ناپارامتری تابع چگالی احتمال جستجو کرد. اولین نوشته‌ی رسمی در زمینه‌ی رگرسیون ناپارامتری در سال ۱۹۶۴ منتشر شد. در آن سال نادارایا و واتسون به طور جداگانه با ایده گرفتن از برآوردگر هسته برای چگالی احتمال، برآوردگری را برای منحنی رگرسیون ارائه دادند که بعدها به نام این افراد نام‌گذاری شد. نادارایا در مقاله‌ی خود شرایط کافی برای سازگاری این برآوردگر را بیان کرده و به توزیع مجانبی نرمال آن نیز اشاره نموده است.

از طرفی در اغلب مسائل آمار کاربردی، مشاهدات در دست بررسی، از جامعه‌ای نمونه‌گیری شده‌اند که توزیع آن با توزیع مورد مطالعه‌ی محقق متفاوت است. چنین مشاهداتی، نمونه‌های اریب نامیده می‌شوند، زیرا با استفاده از این داده‌ها در استنباط در مورد توزیع جامعه، اریبی به وجود می‌آورند. داده‌های سانسور از مشهورترین نمونه‌های اریب می‌باشند.

در سانسور، برخی از مشاهدات به طور کامل مشخص نمی‌شوند و فقط این اطلاع در اختیار محقق قرار می‌گیرد که این مشاهدات عضوی از یک مجموعه می‌باشند. به عنوان مثال فرض کنید متغیر X نشان‌دهنده‌ی زمان رخ دادن یک پیشامد باشد. برای پیشامدی که تا زمان t رخ نمی‌دهد، مقدار X نامشخص است و فقط این اطلاع در مورد آن وجود دارد که مقدار این متغیر در بازه‌ی (t, ∞) قرار دارد. به عبارت دیگر یک مشاهده در مجموعه‌ی A سانسور می‌شود اگر به جای مشاهده‌ی X فقط $X \in A$ مشخص باشد. انتخاب‌های متفاوت از مجموعه‌ی A گونه‌های مختلفی از سانسور را به وجود می‌آورند. یکی از مهمترین انواع سانسور، سانسور تصادفی است.

بنابراین وجود سانسور در مشاهدات را نمی‌توان نادیده گرفت.

بران (۱۹۸۱) یک کلاس از برآوردگرهای ناپارامتری تابع رگرسیون در مدل سانسور شده را معرفی کرد. بعدها داباروسکا (۱۹۸۷)، کاربونز (۱۹۹۵) و کهلر (۲۰۰۲) نیز به معرفی برآوردگرهای مختلف تابع رگرسیون در شرایط سانسور، همراه با بیان ویژگی‌های آن‌ها پرداختند.

بدیهی است که پذیرفتن استقلال بین متغیرهای تصادفی همیشه امکان‌پذیر نیست و در برخی رویدادهای واقعی، نوعی از وابستگی بین متغیرهای تصادفی دیده می‌شود. وابستگی α -آمیزنده، یکی از ضعیف‌ترین وابستگی‌های موجود می‌باشد که تاکنون مورد توجه بسیاری قرار گرفته است. این نوع وابستگی توسط روزن بلات (۱۹۵۶) معرفی شد. وی قضیه‌ی حد مرکزی را برای متغیرهای تصادفی α -آمیزنده ثابت نمود. از آن زمان تاکنون توجه به متغیرهای تصادفی α -آمیزنده، از اهمیت بسیاری برخوردار بوده است. کاربرد این متغیرها در نظریه‌ی آمار و احتمال بسیار گسترده می‌باشد.

دوکن (۱۹۹۴)، بوسک (۱۹۹۸) و لیچر (۲۰۰۱) به بررسی ویژگی‌های تابع رگرسیون تحت شرایط α -آمیزنده پرداختند.

گوسوم و اولدسعید (۲۰۰۹) به معرفی برآوردگرهای نادارایا-واتسون تحت شرایط سانسور و نیز α -آمیزنده پرداختند و ویژگی‌های مجانبی آن را بیان کردند.

ما نیز با توجه به نتایج به دست آمده، به بررسی رفتار مجانبی برآوردگرهای تعریف شده توسط آن‌ها در این مدل می‌پردازیم و برای نشان دادن رفتار برآوردگر در نمونه‌های با حجم کم، به مطالعه‌ی شبیه‌سازی انجام شده می‌پردازیم.

در فصل اول پایان‌نامه، تعاریف و مفاهیم موردنیاز در فصل‌های آتی، ارائه خواهد شد. در این فصل، ابتدا مروری بر برآوردگر ناپارامتری تابع چگالی احتمال داریم. پس از آن، با مفهوم متغیرهای سانسور تصادفی آشنا می‌شویم که پایه و اساس ورود به فصل دوم می‌باشد. در ادامه، تعریف وابستگی ضعیف و

متغیرهای تصادفی آمیزنده مورد مطالعه قرار می‌گیرد. سپس به معرفی برخی از انواع رگرسیون ناپارامتری، قضایا و نابرابری‌های لازم می‌پردازیم.

در فصل دوم با برآوردگر هسته‌ای تابع رگرسیون در مدل سانسور تصادفی در حالت یک متغیره آشنا می‌شویم و سازگاری قوی را همراه با یک نرخ برای برآوردگر معرفی شده و نرمال بودن مجانبی آن نشان می‌دهیم. همچنین با توجه به نتایج کسب شده، یک بازه‌ی اطمینان برای برآوردگر معرفی شده ارائه می‌کنیم. در نهایت به منظور روشن ساختن نتایج و نشان دادن ویژگی نرمال بودن مجانبی برآوردگر در نمونه‌های متناهی، شبیه‌سازی را گنجانده‌ایم.

در فصل سوم به بررسی برآوردگر هسته‌ای تابع رگرسیون با شرط α -آمیزنده بودن متغیرهای تصادفی مورد بررسی در حالت چندمتغیره می‌پردازیم و همچنین سازگاری قوی را با نرخ برای برآوردگر معرفی شده نشان می‌دهیم.

لازم به ذکر است که در طول این پایان‌نامه همواره سعی بر آن بوده است که تمامی قضایا و نتایج موردنیاز در همین مجموعه آورده شود تا خواننده بتواند بدون نیاز به مراجع، برهان قضایای اصلی را دنبال کند. همچنین برای جلوگیری از حجیم شدن این گردایه، از آوردن برخی برهان‌ها اجتناب کرده و تنها به ذکر نام مرجع بسنده کرده‌ایم. علی‌رغم این که نمادهای مربوط به هر فصل، در همان قسمت معرفی می‌شوند، در ادامه فهرستی از نمادهای استفاده شده در کل پایان‌نامه ارائه می‌شود تا سبب سهولت در امر آشنایی با نمادها شود. به منظور حفظ یکپارچگی و انسجام در نوشتار پایان‌نامه و نیز جلوگیری از ترکیب کلمات انگلیسی با متن فارسی، واژه‌نامه‌ای را در انتها گردآوری نمودیم تا خوانندگان بتوانند به صورت یکجا و آنی، برگردان بیشتر واژه‌های مطرح شده در پایان‌نامه را مشاهده نمایند. در برگردان علائم، اسامی و ترجمه کلمات کلیدی از کتاب واژه‌ها و اصطلاحات آماری چاپ مرکز آمار ایران استفاده گردیده است.

سمیه باوفا

پاییز ۱۳۹۰

نمادها

E امید ریاضی
F تابع توزیع مورد بررسی
f تابع چگالی مورد بررسی
g تابع توزیع متغیر سانسور تصادفی
$I(A)$ تابع نشانگر مجموعه A
K تابع هسته
h_n پهنای باند
C_1, C_2, D ثابت عمومی
$\xrightarrow{a.s}$ همگرایی تقریباً حتمی
\xrightarrow{d} همگرایی در توزیع
\xrightarrow{P} همگرایی در احتمال
A, B, \mathcal{F} سیگما میدان
$a(n)$ ضریب آمیزنده
\mathbb{R} مجموعه‌ی اعداد حقیقی

\mathbb{N}	مجموعه‌ی اعداد طبیعی
\mathbb{Z}	مجموعه‌ی اعداد صحیح
$[n]$	بزرگ‌ترین عدد صحیح کوچک‌تر مساوی n
\emptyset	مجموعه‌ی تهی
A^c	متمم مجموعه‌ی A
Y	متغیر مورد بررسی
$Y \wedge T$	مینیمم بین Y و T
$f^{(j)}$	مشتق j -ام تابع چگالی
$K^{(j)}$	مشتق j -ام تابع هسته
var	واریانس
$BV([a, b])$	مجموعه‌ی همه‌ی توابع با تغییر کراندار روی $[a, b]$
τ_F	$\sup\{y, \bar{F}(y) > \circ\}$
a_F	$\inf\{z, F(z) > \circ\}$
$:=$	تعریف می‌کنیم

فصل ۱

تعاریف و مفاهیم مورد نیاز

۱.۱ برآورد تابع چگالی

۲.۱ تحلیل بقا و سانسور

۳.۱ وابستگی ضعیف و متغیرهای تصادفی آمیزنده

۴.۱ رگرسیون ناپارامتری

۵.۱ سایر تعاریف، قضایا و نابرابری‌های مورد نیاز

در این فصل با تعاریف و ویژگی‌های مقدماتی و برخی مفاهیم دیگر که در فصل‌های آینده مورد استفاده قرار می‌گیرند آشنا می‌شویم. به این منظور ابتدا در نخستین بخش به بررسی برآورد تابع چگالی و بعضی از انواع آن می‌پردازیم.

در بخش دوم مروری بر تحلیل بقا، تابع بقا و سانسور خواهیم داشت. سپس در بخش بعد وابستگی ضعیف و متغیرهای تصادفی آمیزنده مورد مطالعه قرار می‌گیرند. در بخش چهارم به معرفی رگرسیون ناپارامتری و برخی از انواع مهم آن می‌پردازیم و سرانجام در بخش پنجم به بیان سایر تعاریف، قضایا و نابرابری‌های لازم در برهان‌های ارائه شده، می‌پردازیم.

۱.۱ برآورد تابع چگالی احتمال

در این قسمت به بررسی برآورد تابع چگالی و انواع آن توسط یک نمونه‌ی تصادفی مانند X_1, \dots, X_n که دارای چگالی f هستند می‌پردازیم.

ابتدا فرض کنید $a < x < b$ ، برای برآورد تابع چگالی در نقطه‌ی x ، یعنی برآورد $f(x)$ ، چنانچه فاصله‌ی a تا b یعنی $h = b - a$ به اندازه کافی کوچک باشد آن‌گاه می‌توان گفت

$$f(x) \sim \frac{\int_a^b f(u) du}{b - a} \quad (1.1)$$

از طرفی اگر فاصله‌ی a تا b یعنی h به اندازه‌ی کافی بزرگ باشد به طوری که مقدار قابل ملاحظه‌ای از مشاهدات X_1, \dots, X_n در آن قرار بگیرند آن‌گاه

$$\int_a^b f(u) du = P(a < X < b)$$

را می‌توان توسط فراوانی نسبی تعداد مشاهداتی که در بازه‌ی (a, b) قرار می‌گیرند تقریب زد. فرض کنید

متغیر تصادفی Y تعداد مشاهداتی باشد که در بازه‌ی (a, b) قرار می‌گیرند. پس

$$P(a < X < b) = \int_a^b f(u)du \sim \frac{Y}{n} \quad (۲.۱)$$

توسط (۱.۱) و (۲.۱) داریم:

$$f(x) \sim \frac{\int_a^b f(u)du}{b-a} \sim \frac{Y}{n(b-a)}$$

بنابراین برآوردی برای $f(x)$ را می‌توان $\frac{Y}{n(b-a)}$ در نظر گرفت که آن را با $f_n(x)$ نمایش می‌دهیم.

در ادامه به تعریف بعضی از انواع برآوردگرهای تابع چگالی می‌پردازیم.

تابع چگالی را می‌توان به وسیله‌ی روش‌های مختلفی برآورد کرد، از جمله برآوردگر هیستوگرام،

برآوردگر باکس کار، برآوردگر هسته‌ی تابع چگالی و ... که در ادامه به بعضی از آن‌ها اشاره می‌کنیم.

۱.۱.۱ برآوردگر هیستوگرام

ساده‌ترین برآوردگر تابع چگالی برآوردگر هیستوگرام است. فرض کنید f در فاصله‌ی دلخواهی تعریف

شده باشد. بدون از دست دادن کلیت فرض می‌کنیم که f روی فاصله‌ی $[0, 1]$ تعریف شده باشد. اگر m

یک عدد صحیح باشد بازه‌هایی را به صورت زیر تعریف می‌کنیم

$$B_1 = \left[0, \frac{1}{m}\right], B_2 = \left[\frac{1}{m}, \frac{2}{m}\right], \dots, B_m = \left[\frac{m-1}{m}, 1\right]$$

$h = \frac{1}{m}$ را پهنای باند تعریف می‌کنیم. فرض کنیم Y_j ، تعداد مشاهدات که در بازه‌ی B_j قرار می‌گیرند،

باشد. قرار می‌دهیم

$$\hat{P}_j = \frac{Y_j}{n} \quad , \quad P_j = \int_{B_j} f(u)du$$

تعریف ۱.۱.۱ برآوردگر هیستوگرام تابع چگالی عبارت است از

$$\hat{f}_n(x) = \sum_{j=1}^m \frac{\hat{P}_j}{h} I(x \in B_j)$$

توجه کنید که برای هر $x \in B_j$ و h کوچک

$$E(\hat{f}_n(x)) = \frac{E(\hat{P}_j)}{h} = \frac{P_j}{h} = \frac{\int_{B_j} f(u) du}{h} \approx \frac{f(x)h}{h} = f(x)$$

بنابراین برآوردگر هیستوگرام تابع چگالی تقریباً نارایب می‌باشد.

□

۲.۱.۱ برآوردگر باکس کار

فرض کنید Y_j در قسمت قبل را با $K_n(x - \frac{h}{4}, x + \frac{h}{4})$ تعویض کنیم که $K_n(x - \frac{h}{4}, x + \frac{h}{4})$ برابر است با

تعداد مشاهدات در بازه‌ی $(x - \frac{h}{4}, x + \frac{h}{4})$. با توجه به این که

$$x - \frac{h}{4} < X_i < x + \frac{h}{4} \iff \left| \frac{x - X_i}{h} \right| < \frac{1}{4}$$

و با در نظر گرفتن تابع

$$e(t) = \begin{cases} 1 & |t| < \frac{1}{4} \\ 0 & |t| > \frac{1}{4} \end{cases}$$

به تعریف زیر دست می‌یابیم.

تعریف ۲.۱.۱ برآوردگر تابع چگالی باکس کار عبارت است از

$$\hat{f}_n(x) = \frac{K_n(x - \frac{h}{4}, x + \frac{h}{4})}{nh} = \frac{1}{n} \sum_{i=1}^n \frac{1}{h} e\left(\frac{x - X_i}{h}\right)$$

□

۳.۱.۱ برآوردگر هسته‌ای تابع چگالی

به علت هموار نبودن هیستوگرام‌ها، اکنون برآوردگرهای هسته‌ای تابع چگالی که هموارتر می‌باشند و سریع‌تر به یک چگالی صحیح همگرايند، را بیان می‌کنیم. قبل از تعریف و بیان ویژگی‌های برآوردگرهای هسته‌ای به تعریف هسته و بیان انواع آن می‌پردازیم.

تعریف ۳.۱.۱ تابع نامنفی K را هسته گویند هرگاه شرایط زیر را داشته باشد:

$$1. \int K(x)dx = 1$$

$$2. \int xK(x)dx = 0$$

$$3. \int x^2 K(x)dx \neq 0 \quad (\sigma^2 > 0)$$

□

کلمه‌ی هسته به هر تابع هموار K که شرایط بالا را در بر داشته باشد اطلاق می‌شود. هسته‌ها به دو گروه تقسیم می‌شوند.

۱. هسته‌های متقارن

۲. هسته‌های نامتقارن

بعضی از انواع هسته‌های متقارن و نامتقارن را در زیر بیان می‌کنیم:

هسته‌های متقارن

هسته‌ی باکس‌کار: $K(x) = \frac{1}{\sigma} I(x)$

$$K(x) = \frac{1}{\sqrt{\pi}} e^{-x^2}: \text{هسته ی گوسی:}$$

$$K(x) = \frac{\sqrt{2}}{\pi} (1 - x^2) I(x): \text{هسته ی اپانچنیکوف:}$$

$$K(x) = \frac{\sqrt{2}}{\pi} (1 - |x|^3)^3 I(x): \text{هسته ی تریکوب:}$$

هسته های نامتقارن

$$K_{\frac{x}{b+1}, b}(t) = \frac{t^{\frac{x}{b}} e^{-\frac{t}{b}}}{b^{\frac{x}{b+1}} \Gamma(\frac{x}{b+1})}: \text{هسته ی گاما:}$$

هسته ی معکوس گوسی:

$$K_{IG}(m, d)ty = \frac{\sqrt{d}}{\sqrt{\pi} y^3} \exp\left(-\frac{d}{m} \left(\frac{d}{m} - \sqrt{2} + \frac{m}{y}\right)\right) \quad y > 0$$

یا

$$K(M, h) = \frac{1}{\sqrt{\pi} h x^3} e^{-\left(\frac{1}{\sqrt{h}}\right)(x-M)^2/xM^2} \quad x > 0$$

هسته ی عکس معکوس گوسی:

$$K_{RIG}(m, d)^{(2)} = \frac{\sqrt{d}}{\sqrt{\pi} z} \exp\left(-\frac{d}{m} \left(mz - \sqrt{2} + \frac{1}{mz}\right)\right) \quad z > 0$$

تعریف ۴.۱.۱ فرض کنید هسته ی متقارن K و عدد مثبت h که پهنای باند نامیده می شود داده شده است،

در این صورت برآورد گر چگالی هسته به صورت زیر تعریف می شود:

$$\hat{f}_n(x) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h} K\left(\frac{x - X_i}{h}\right)$$

□

معمولاً انتخاب K و انتخاب پهنای باند یعنی h مهم تر است. مثلاً ممکن است با پهنای باند کوچک یک برآورد بسیار ناهمواری به دست آید، در حالی که با پهنای باند بزرگ تر، برآورد هموارتری داشته باشیم.

تحت شرایطی برآوردگر هسته‌ی تابع چگالی، برآوردگری سازگار (ضعیف) برای f می‌باشد که این مطلب در قضیه‌ی زیر بیان شده است.

قضیه ۱.۱.۱ فرض کنید f در x پیوسته باشد و وقتی $n \rightarrow \infty$ ، $h \rightarrow 0$ و $nh \rightarrow \infty$ در این صورت داریم:

$$\hat{f}_n(x) \xrightarrow{P} f(x)$$

یعنی برآوردگر $\hat{f}_n(x)$ در احتمال به $f(x)$ میل می‌کند و یا

$$\lim_{n \rightarrow \infty} P(|\hat{f}_n(x) - f(x)| > \epsilon) = 0 \quad \forall \epsilon > 0$$

تعریف ۵.۱.۱ (برآوردگر هسته‌ای تابع چگالی در حالت چندمتغیره)

فرض کنید X_1, \dots, X_n یک نمونه‌ی تصادفی d بعدی با چگالی f باشد. همچنین

$$\mathbf{X}_i = (X_{i1}, \dots, X_{id})^T$$

دوهاولز (۱۹۷۷) در حالت چندمتغیره، برآوردگر هسته‌ای تابع چگالی احتمال را به صورت زیر معرفی کرد.

$$\hat{f}(x, \mathbf{H}) = n^{-1} \sum_{i=1}^n K_{\mathbf{H}}(x - \mathbf{X}_i)$$

که در آن \mathbf{H} یک ماتریس $d \times d$ ، مثبت معین و متقارن است که ماتریس پهنای باند نامیده می‌شود و

$$K_{\mathbf{H}}(x) = |\mathbf{H}|^{-\frac{1}{d}} K(\mathbf{H}^{-\frac{1}{d}} x)$$

که در آن K یک تابع هسته d متغیره است که شرط زیر را برقرار می‌کند

$$\int_{\mathbb{R}^d} K(x_1, \dots, x_d) dx_1 \dots dx_d = 1$$

معمولاً به دو روش، با استفاده از هسته‌ی یک متغیره و متقارن K ، هسته‌ی چندمتغیره در نظر گرفته می‌شود

$$K^P(x) = \prod_{i=1}^d \kappa(x_i) \quad , \quad K^S(x) = c_{\kappa,d} \kappa\{(x^T x)^{\frac{1}{d}}\}$$

که

$$c_{\kappa,d}^{-1} \int \kappa\{(x^T x)^{\frac{1}{d}}\} dx \quad , \quad x = (x_1, \dots, x_d)$$

$K^P(X)$ را هسته‌ی حاصلضربی و $K^S(X)$ را هسته‌ی کروی نامند.

یک حالت خاص از برآوردگر هسته‌ای تابع چگالی در حالت چندمتغیره را می‌توان تابع زیر را که در

اکثر کاربردها استفاده می‌شود، در نظر گرفت (با فرض ماتریس قطری $(\mathbf{H} = \text{diag}(h_1, \dots, h_d))$)

$$\hat{f}(x; h) = n^{-1} \left(\prod_{l=1}^d h_l \right)^{-1} \sum_{i=1}^n K \left(\frac{x_1 - X_{i1}}{h_1}, \dots, \frac{x_d - X_{id}}{h_d} \right)$$

و با در نظر گرفتن پهنای باند یکسان، خواهیم داشت

$$\hat{f}(x; h) = n^{-1} h^{-d} \sum_{i=1}^n K\{(x - X_i)/h\}$$

□

۲.۱ تحلیل بقا و سانسور

تعریف ۱.۲.۱ تحلیل بقا عبارت است از گردآوری از روش‌های آماری که برای تجزیه و تحلیل داده‌های

نامنفی به کار می‌رود. معمولاً این داده‌های نامنفی (داده‌های بقا)، داده‌های مربوط به زمان شروع تا رسیدن

به یک پیشامد خاص می‌باشند.

پیشامد خاص را به عنوان مثال می توان پیشامد رسیدن به یک بیماری خاص، پیشامد رسیدن مرگ بر اثر بیماری خاص، پیشامد انتشار یک تومور و ... در نظر گرفت. در ادامه ما از نماد T برای نشان دادن متغیر تصادفی بقا استفاده می کنیم. همچنین فرض کنید متغیر تصادفی T دارای تابع چگالی $f(t)$ و تابع توزیع $F_T(t)$ باشد.

معمولاً به دو دلیل عمده زیر، تجزیه و تحلیل داده های بقا از روش های استاندارد آماری استفاده شده در آنالیز داده ها تبعیت نمی کند:

۱. داده های بقا دارای توزیع نرمال نمی باشند.

۲. دلیل دیگری که باعث می شود روش های استاندارد در مورد داده های بقا مناسب نباشد، این است

که عموماً سانسور می شوند. که در ادامه مفهوم سانسور بیان خواهد شد.

در خلاصه کردن داده های بقا، دو تابع به نام های تابع بقا و تابع خطر نقش اساسی دارند.

□

تعریف ۲.۲.۱ تابع بقا که آن را با $S(t)$ نشان می دهیم عبارت است از

$$S(t) = P(T > t)$$

این تابع نشان دهنده ی احتمال بقای یک عنصر (فرد) تا زمان t است.

□

تعریف ۳.۲.۱ تابع خطر که آن را با $h(t)$ نشان می دهیم به صورت زیر تعریف می شود

$$h(t) = \lim_{\Delta t \rightarrow 0} \frac{P(t \leq T < t + \Delta t | T \geq t)}{\Delta t}$$

تابع فوق نشان دهنده‌ی نرخ مرگ فوری برای افراد زنده مانده تا زمان t می‌باشد. اگر T متغیر تصادفی پیوسته باشد، به سادگی داریم

$$h(t) = \frac{f(t)}{S(t)} = \frac{f(t)}{1 - F_T(t)}$$

□

همان‌طور که قبلاً بیان شد یکی از دلایل مهمی که باعث می‌شود تحلیل داده‌های بقا متمایز از تحلیل‌های استاندارد شود، سانسور می‌باشد.

تعریف ۴.۲.۱ زمان بقای یک عنصر (فرد) را سانسور شده گویند هرگاه زمان بقای آن عنصر (فرد) تا پایان دوره‌ی مطالعه نشود. به عبارت دیگر زمان بقای یک عنصر را سانسور شده گویند، هرگاه نقطه‌ی پایان پیشامد موردنظر برای آن عنصر مشاهده نشود. بنابراین داده‌ی سانسور شده یک داده‌ی ناقص است. معمولاً دلایل عمده‌ی سانسور شدن عبارت‌اند از

۱. پایان مطالعه

۲. عدم امکان پیگیری

۳. خارج شدن از مطالعه

از دیدگاهی، سانسور را می‌توان به صورت زیر تقسیم‌بندی کرد.

۱. سانسور نوع اول:

فرض کنید C یک عدد ثابت و X_1, \dots, X_n یک نمونه تصادفی از توزیع F باشند. در این سانسور

$$Y_i = \min(X_i, C) \text{ یعنی}$$

$$Y_i = \begin{cases} X_i & X_i \leq C \\ C & X_i \geq C \end{cases}$$

دیده می شوند.

۲. سانسور نوع دوم

فرض کنید $X_{(1)}, \dots, X_{(n)}$ آماره‌های ترتیبی X_1, \dots, X_n و $r < n$ باشد. در این نوع سانسور

فرض می کنیم تا مشاهده شدن r امین شکست، عناصر تحت مطالعه قرار گیرند. بنابراین مشاهدات

به صورت زیر به دست می آیند

$$Y_1 = X_{(1)}, Y_2 = X_{(2)}, \dots, Y_r = X_{(r)}, Y_{(r+1)} = X_{(r)}, \dots, Y_{(n)} = X_{(r)}$$

۳. سانسور تصادفی (نوع سوم):

متغیر تصادفی X_i توسط متغیر تصادفی Y_i از (راست) سانسور تصادفی می شود هرگاه زوج مرتب

$$(Z_i, \delta_i) \text{ مشاهده شوند که در آن } \delta_i = I(X_i \leq Y_i) \text{ و}$$

$$Z_i = \min(X_i, Y_i)$$

فرض می کنیم متغیرهای نشان دهنده‌ی زمان‌های سانسور یعنی $Y_i, 1 \leq i \leq n$ مستقل و هم توزیع با تابع

توزیع مشترک G هستند. همچنین Y_i ها مستقل از X_i ها فرض می شوند. از آنجایی که داده‌های سانسور

بیشتر در مسائل مربوط به تحلیل بقا رخ می دهند، بدون کم شدن از کلیت مسأله X_i ها و Y_i ها نامنفی فرض

می شوند. برای داده‌های سانسور شده δ_i ها برابر صفر و برای داده‌های غیر سانسور، δ_i ها برابر یک است.

□