



پایان نامه کارشناسی ارشد
رشته آمار ریاضی

عنوان
برآورد تابع چگالی به روش لگ-اسپلاین

استاد راهنمای
دکتر حسینعلی نیرومند

استاد مشاور
دکتر حسن دوستی

تدوین
مهدیه عرفانیان

شهریور ۱۳۸۸

یا همو

چراغ، جزبه نور چراغ، سوان دید.



بسمه تعالیٰ .

مشخصات رساله/پایان نامه تحصیلی دانشجویان .

دانشگاه فردوسی مشهد

دانشگاه فردوسی مشهد

عنوان رساله/پایان نامه: برآورد تابع چگالی به روش لگ-اسپلاین

نام نویسنده: مهدیه عرفانیان

نام استاد(ان) راهنما: دکتر حسینعلی نیرومند

نام استاد(ان) مشاور: دکتر حسن دوستی

رشته تحصیلی: آمار ریاضی	گروه: آمار	دانشکده : علوم ریاضی
تاریخ دفاع: ۱۳۸۸/۶/۲۸		تاریخ تصویب:
تعداد صفحات: ۱۴۸	دکتری ○	مقطع تحصیلی: کارشناسی ارشد ●

چکیده رساله/پایان نامه :

در این پایان نامه پس از مرور تعدادی از روش های ناپارامتری برآورد چگالی که تاکنون ارائه شده اند، به معرفی اسپلاین ها، B-اسپلاین ها و روش بهینه سازی نیوتون-رافسون پرداخته ایم. مدل های خطی گسترش یافته را بیان و حالت خاص آن، یعنی مدل لگ-اسپلاین برای تابع چگالی را معرفی کرده ایم. در روش لگ-اسپلاین، ابتدا فرض می کنیم تابع چگالی جامعه به شکل تبدیل لجستیک است و سپس برآورد پارامترهای مجھول مدل را با استفاده از برآوردگر درست نمایی ماکزیمم به دست می آوریم که از روش نیوتون-رافسون محاسبه می شوند. به منظور به دست آوردن تعداد بهینه توابع پایه، از دو روش انتخاب گام به گام گره ها و ارائه برآوردگرهای سازگار ضعیف و قوی برای تعداد توابع پایه استفاده می کنیم. روش گام به گام از ترکیب دو روش اضافه کردن گره ها و حذف کردن گره ها تشکیل شده است که در آن ها به ترتیب آماره های رائو و والد به کار گرفته شده اند. در روش ارائه برآوردگر سازگار، تعداد توابع پایه را به صورت یک پارامتر مجھول در نظر می گیریم و سپس برآوردگرهايی را برای این پارامتر با استفاده از یک جانشینی برای فاصله کوبیک-لیبلر به دست می آوریم. اثبات می شود که بر حسب نوع محک اطلاعی که در به دست آوردن این برآوردگرها استفاده می کنیم، برآوردگرها سازگار ضعیف یا قوی خواهند بود. در انتهای با استفاده از بسته لگ-اسپلاین موجود در نرم افزار آماری R به شبیه سازی پرداخته ایم و روش های ارائه شده برای یافتن تعداد بهینه توابع پایه را مقایسه کرده ایم.

کلید واژه:

B-اسپلاین، اسپلاین، انتخاب گام به گام گره ها، برآورد ناپارامتری تابع چگالی، برآورد درست نمایی ماکزیمم، برآوردگرهای سازگار ضعیف و قوی، روش بهینه سازی نیوتون-رافسون، محک اطلاع، مدل های خطی گسترش یافته.

تاریخ:

امضای استاد راهنما:



بسمه تعالى .

**Graduate Studies Thesis\ Dissertation Information
Ferdowsi University of Mashhad**

Title of Thesis\ Dissertation: Estimation of the density function by the Logspline method

Author: Mahdiyeh Erfaniyan

Supervisor(s): Dr. Hosseinali Niroumand

Advisor(s): Dr. Hassan Doosti

Faculty: Mathematical Sciences	Department: Statistics	Specialization: Mathematical Statistics
---------------------------------------	-------------------------------	--

Approval Date:	Defence Date: 22-August-2009
-----------------------	-------------------------------------

M.Sc. <input checked="" type="radio"/>	Ph.D. <input type="radio"/>	Number of Pages: 148
---	------------------------------------	-----------------------------

Abstract:

In this dissertation, after a review of some nonparametric density estimation methods which has been provided until now, we have introduced splines, B-splines and the Newton-Raphson optimization method. We have proposed the extended linear models and a special case of that, called the Logspline model of the density function. In the Logspline method, we first suppose that the density function of the population can be written in the form of the logistic transformation. We estimate the unknown parameters of model by the maximum likelihood estimators calculated by the Newton-Raphson method. In order to find the optimum number of basis functions, we use two methods which are the stepwise knots selection and providing the weakly and strongly consistent estimators for the number of basis functions. In the stepwise knots selection method, we use a combination of the knot addition and the knot deletion in which the Rao and Wald statistics are applied, respectively. For providing the weakly and strongly consistent estimators, we treat the number of basis functions as an unknown parameter and find some estimators for it using a proxy of the Kullback-Liebler distance. It is proved that these estimators are weakly and strongly consistent estimators in the order of the information criteria we have used for finding them. At the end, we have performed some simulation using the Logspline package written in R statistical language for comparing the two introduced method for finding the number of basis functions.

Signature of Supervisor:	Key Words:
---------------------------------	-------------------

B-splines, Extended Linear Models, Maximum Likelihood Estimation, Nonparametric Density Estimation, Splines, Stepwise Knot Selection, the Newton-Raphson Optimization Method, Information Criteria, Weakly and Strongly Consistent Estimators.

Date:

تعدیم به

او که نشانم داد فرشگان خدا فقط در آسمان ها نیستند

هستند فرشگانی که در زمینند و مثل ما از حائلند

او که خود فرشته ای زمینی است،

تعدیم به مادرم ۰۰۰

فهرست مطالب

۴	پیش‌گفتار
۸	فصل اول مروری بر معمول ترین روش‌های ناپارامتری برآورد تابع چگالی
۸	۱-۱ مقدمه
۸	۲-۱ چرا برآورد چگالی؟
۱۱	۳-۱ برآورد ناپارامتری چگالی
۱۱	۱-۳-۱ رویکرد بافت‌نگار موضعی
۱۳	۲-۳-۱ اقتباس صوری $\hat{f}_1(x)$
۱۴	۳-۳-۱ برآوردگر هسته‌ی روزنبلات-پارزن
۱۷	۴-۳-۱ برآوردگر نزدیک‌ترین همسایگی
۱۹	۵-۳-۱ برآوردگرهای پهنای نوار-متغیر
۲۱	۶-۳-۱ برآوردگرهای درست‌نمایی جریمه‌شده
۲۵	۷-۳-۱ برآوردگرهای لگاریتم درست‌نمایی موضعی
۲۸	۸-۳-۱ روش لگ-اسپلاین
۳۲	فصل دوم معرفی اسپلاین‌ها و دیگر مفاهیم اساسی
۳۲	۱-۲ مقدمه
۳۳	۲-۲ اسپلاین چیست؟
۳۵	۱-۱-۲ چندجمله‌ای‌ها
۳۵	۲-۲-۲ چندجمله‌ای‌های تکه‌ای
۳۶	۳-۱-۲ اسپلاین‌ها
۳۸	۴-۲-۲ B -اسپلاین‌ها: تعریف
۳۹	B -اسپلاین‌های مرتبه‌ی ۱
۴۰	B -اسپلاین‌های مرتبه‌ی ۲
۴۲	B -اسپلاین‌های مرتبه‌ی k

۴۴	ویژگی‌های مهم B -اسپلاین‌ها	۵-۲-۲
۴۹	تقریب تابع	۳-۲
۴۹	ویژگی‌های تقریب چندجمله‌ای	۱-۳-۲
۵۰	چرا اسپلاین‌ها؟	۲-۳-۲
۵۲	چرا B -اسپلاین‌ها؟	۳-۳-۲
۵۳	فاصله‌ی یک تابع از یک فضای اسپلاین	۴-۳-۲
۵۶	تکعر	۴-۲
۵۶	حالت یک متغیره	۱-۴-۲
۵۹	حالت چندمتغیره	۲-۴-۲
۶۰	بررسی مقربه‌بودن تابع	۳-۴-۲
۶۲	بهینه‌سازی	۵-۲
۶۳	نگاهی کلی به روش نیوتون-رافسون	۱-۵-۲
۷۰	فصل سوم برآورد تابع چگالی با استفاده از لگ-اسپلاین‌ها	
۷۰	مقدمه	۱-۳
۷۱	یک مثال	۲-۳
۷۱	داده‌های درآمد	۱-۲-۳
۷۳	پیش‌زمینه	۲-۲-۳
۷۷	برآورد چگالی لگ-اسپلاین	۳-۲-۳
۸۰	روش‌شناسی لگ-اسپلاین	۳-۳
۸۰	مدل‌های خطی گسترش‌یافته	۱-۳-۳
۸۳	مدل لگ-اسپلاین	۲-۳-۳
۸۷	توابع پایه	۳-۳-۳
۹۰	جزئیات فنی	۴-۳
۹۰	مکان‌یابی اویلیه‌ی گره‌ها	۱-۴-۳
۹۲	انتگرال‌گیری عددی	۲-۴-۳
۹۵	فصل چهارم برآورد تعداد توابع پایه	
۹۵	مقدمه	۱-۴
۹۶	انتخاب گام به گام گره‌ها	۲-۴
۹۶	معرفی اجمالی	۱-۲-۴
۹۸	برآورد درست‌نمایی ماکزیمم	۲-۲-۴

۱۰۰	نیمه کردن گام.....	۳-۲-۴
۱۰۲	چگونه به یک تکرار پایان دهیم؟.....	۴-۲-۴
۱۰۳	اضافه کردن یک تابع پایه.....	۵-۲-۴
۱۰۴	حذف کردن یک تابع پایه.....	۶-۲-۴
۱۰۵	توضیح دقیق شیوه‌ی گام به گام انتخاب گره‌ها.....	۷-۲-۴
۱۰۷	برآورد سازگار برای تعداد توابع پایه با استفاده از یک جانشین برای فاصله‌ی کولبک-لیبلر.....	۳-۴
۱۰۷	معرفی اجمالی.....	۱-۳-۴
۱۰۸	شروط اضافه برای شیوه‌ی برآورد سازگار تعداد توابع پایه.....	۲-۳-۴
۱۱۴	توضیح دقیق شیوه‌ی برآورد سازگار تعداد توابع پایه با استفاده از یک جانشین برای فاصله‌ی کولبک-لیبلر.....	۳-۳-۴
۱۲۶	فصل پنجم شبیه‌سازی.....	
۱۲۶	مقدمه.....	۱-۵
۱۲۷	نتایج شبیه‌سازی.....	۲-۵
۱۳۶	بسته‌ی <i>logspine</i>	۳-۵
۱۳۸	برنامه‌های شبیه‌سازی.....	۴-۵
۱۴۴	منابع.....	
۱۴۶	واژه‌نامه‌ی انگلیسی به فارسی.....	

پیش‌گفتار

در دهه‌های گذشته شاهد تغییر باورنکردنی در نظریه‌ی آمار بوده‌ایم. پیشرفت عظیم در توان کامپیوترها، باعث به وجود آمدن شیوه‌های تابعی سازگاری شده که اکنون ابزار ضروری برای تحلیل پیشرفته‌ی داده‌ها است. حال که آماردانان از فرض‌های سخت‌گیرانه‌ای که جزء بی‌چون و چرای مدل‌های پارامتری کلاسیک هستند، رهایی یافته‌اند، از آنان انتظار می‌رود که نه تنها قادر به برگزیدن متغیرهای مهم در یک مدل باشند، بلکه بتوانند شکل تابعی وابستگی این متغیرها را نیز انتخاب کنند. برای این که بتوان در کاربرد نیز توفیق یافت، به ناچار هر شیوه‌ی جدیدی باید بین انعطاف‌پذیری و مسئله‌ی خطیر تعداد ابعاد، تعادل برقرار کند. اسپلاین‌های چندجمله‌ای، انعطاف‌پذیری مورد نیاز برای تحلیل پیشرفته‌ی داده‌ها را به ارمغان می‌آورند. در این میان یافتن برآورد چگالی جامعه‌ی مورد بحث، یکی از پرتفاضاترین مباحث، هم در زمینه‌ی آمار نظری و هم در مباحث کاربردی است.

در این پایان‌نامه، به معرفی یک روش ناپارامتری برای برآورد چگالی، با عنوان لگ-اسپلاین می‌پردازیم. این روش از دو ابزار اساسی در برآوردهای استفاده می‌کند که عبارتند از اسپلاین‌ها و برآورد درست‌نمایی ماکریم. اسپلاین‌ها که در اصل یکی از ابزارهای معرفی شده در ریاضیات کاربردی هستند، اینکه کاربرد فراوانی در زمینه‌ی علوم دیگر، به خصوص مهندسی یافته‌اند. اسپلاین‌ها، یکی از مفاهیمی هستند که در کنار موجک‌ها، سری‌های فوریه و دیگر توابع پایه در

کار تقریب توابع مورد استفاده قرار می‌گیرند و خواص مطلوب آن‌ها در این امر، بر همگان واضح است. لگ-اسپلاین از مفهوم دیگری به نام برآورده درست‌نمایی ماکزیمم بهره می‌برد. شاید بتوان گفت این روش برآورده، ملموس‌ترین روش برآورده‌یابی است، زیرا در آمار جزو اوّلین روش‌هایی است که برای برآورده از آن یاد می‌شود و هم چنین بررسی‌های بسیاری برای یافتن خواص آن انجام شده است. روش درست‌نمایی ماکزیمم، ویژگی مهم دیگری دارد و آن هم این که در ارائه‌ی طیف گسترده‌ای از آزمون فرضیه‌های آماری کلاسیک به کار می‌رود. این روش، به خاطر مفهوم ملموس و قابل درک آن، حتی در بین دیگر رشته‌های علوم نیز روشنی پرکاربرد و پرطرفدار است؛ چنان که تا کنون کارهای تحقیقی زیادی بر اساس روش درست‌نمایی ماکزیمم انجام یافته و حتی در انجام آزمون فرضیه‌ها نیز، بر حسب زمینه‌ی مورد لزوم، انواع سازگاری‌افته‌ای از این روش ارائه شده است، مانند انواع درست‌نمایی شرطی، حاشیه‌ای، برآورده‌شده، نیمرخ، جریمه‌شده، تجربی و ... بنابراین لگ-اسپلاین را به خاطر استفاده از این دو ابزار اساسی، می‌توان از یک سو روشی سازگار و کارآمد و از سویی دیگر، روشی ملموس و قابل فهم، حتی برای کاربران دیگر رشته‌ها، دانست.

در فصل اوّل این پایان‌نامه، مختصری از معروف‌ترین روش‌های ناپارامتری برآورده چگالی را بیان و سپس در فصل دوم، مفاهیم مورد نیاز از قبیل اسپلاین‌ها و روش‌های عددی به کار رفته در این شیوه را معرفی خواهیم کرد. سپس در فصل سوم مدل لگ-اسپلاین را تعریف و در فصل چهارم دو روش را برای یافتن تعداد توابع پایه ارائه می‌کنیم. در پایان و در فصل پنجم، بر اساس یافته‌های فصول پیشین، به شبیه‌سازی می‌پردازیم.

فرصت موجود، مغتنم است تا از استاد راهنمای محترم خود، جناب آقای دکتر حسینعلی نیرومند تشکر به عمل آورم، که در طول انجام این پایاننامه، از مراتب بالای علمی ایشان، و هم چنین از فضایل والای اخلاقیشان بهره‌های بسیار بردم. همین طور از استاد مشاور این کار، جناب آقای دکتر حسن دوستی سپاسگزارم، که با صبوری بسیار، مرا در مراحل مختلف تدوین این کار، راهنمایی کردند.

از جناب آقای دکتر مهدی عمامی و جناب آقای دکتر محمد آرشی برای قبول داوری این کار، بسیار متشرکرم.

از دوست و هم‌کلاسی عزیز خود، خانم حکیمه مریبی هروی، به خاطر کمک‌های بی‌دریغشان در کار شبیه‌سازی این پایاننامه عرض تشکر فراوان دارم.

مهدیه عرفانیان

Email: mahdiyeh.erfaniyan@gmail.com

فصل اول

مروی بر معمول نرین (وش های نایارا هندی برآورده چگالی

۱-۱ مقدمه

در این فصل مروری کوتاه بر روش‌های ناپارامتری معمول برآورد چگالی یک متغیر تصادفی خواهیم داشت. در بعضی موارد، لازم است بیان کنیم که اصولاً چرا به برآورد چگالی علاقه‌مند هستیم. یک دلیل مهم آن است که لزوم برآورد چگالی، در عمل آن قدر به حد کافی احساس می‌شود که طیف گستره‌ای از تحقیقات آماری به آن اختصاص یافته است. بر این اساس، در ابتدا به بیان دلایل ارائه‌ی برآورد چگالی می‌پردازیم و در ادامه، روش‌های ناپارامتری مختلف برآورد چگالی را، به طور کوتاه، معرفی می‌کنیم.

۲-۱ چرا برآورد چگالی؟

به طور کلی، می‌توان سه زمینه را بیان کرد که در آن‌ها به برآورد چگالی احتیاج است. اوّل آن که، برآورد چگالی می‌تواند برای شناسایی پدیده‌های الگودار مهم باشد. برای مثال، اگر بدانیم متغیری

که تحت مطالعه است دم‌های پهن یا قله‌های بلند دارد، هر مدل از داده‌ها که متناظر با چنین متغیری باشد، باید بتواند یک چگالی با چنین ویژگی‌هایی تولید کند. دوم آن که، برای برآوردگری که در یک مطالعه استفاده می‌شود، اغلب ممکن است یک تحلیل مونت کارلو^۱ لازم باشد. به طور سنتی، تنها تعداد کمی از گشتاورهای برآوردگر، ثبت یا آماره آزمون‌هایی مانند کولموگوروف-سمیرنوف^۲ برای ارزیابی انحراف از توزیع نرمال فراهم می‌شوند. اما برآوردهای ناپارامتری چگالی تصویر کاملی از توزیع برآوردگر را مهیا می‌کنند و بنابراین راه مطلوب‌تری برای خلاصه‌کردن خروجی آزمایش مونت کارلو به نظر می‌رسند. در نهایت آن که، گاهی اوقات، مسئله‌ی موجود آن است که برآوردگرهای پارامتری، توزیع مجانبی دارند که به یک چگالی محاسبه‌شده در یک نقطه‌ی مشخص بستگی دارد. برای مثال، میانه‌ی X دارای واریانس $0.25n^{-1}f^{-2}(0)$ است که در آن $f(0)$ چگالی X در نقطه‌ی $x=0$ است. بنابراین هر آماره آزمونی که شامل میانه باشد، به برآورد مقدار $f(0)$ احتیاج دارد.

حال فرض کنید $f(x) = f$ نشان‌دهنده‌ی تابع چگالی پیوسته‌ی متغیر تصادفی X در هر نقطه‌ی x باشد و x_1, \dots, x_n مشاهدات حاصل از f باشند. دو روش کلی برای برآورد f پیشنهاد شده است:

(الف) **برآوردهای پارامتری:** روش‌های پارامتری شکل مخصوصی برای f در نظر

می‌گیرند، مثلاً چگالی نرمال،

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

¹ Monte Carlo analysis

² Kolmogorov-Smirnov

که در آن میانگین μ و واریانس σ^2 پارامترهای f هستند. برآورد f می‌تواند به صورت زیر

نوشته شود:

$$\hat{f}(x) = \frac{1}{\hat{\sigma}\sqrt{2\pi}} \exp\left[-\frac{1}{2}\left(\frac{x - \hat{\mu}}{\hat{\sigma}}\right)^2\right],$$

که در آن μ و σ به راحتی از داده‌ها به صورت زیر برآورد می‌شوند:

$$\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 \quad \text{و} \quad \hat{\mu} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i.$$

(ب) برآوردهای ناپارامتری: یکی از نوافص روش پارامتری نیازمند بودن آن به دانستن قید

چگالی واقعی f است. در روش ناپارامتری، $f(x)$ مستقیماً و بدون در نظر گرفتن شکل آن برآورد

می‌شود. بافت‌نگار، یکی از روش‌های برآورد ناپارامتری است و یکی از قدیمی‌ترین روش‌های

برآورد چگالی به شمار می‌رود (وان رایزن^۱ (۱۹۷۳) و اسکات^۲ (۱۹۷۹)). اما با وجودی که

بافت‌نگار روش مفیدی برای برآورد چگالی است، مشکلاتی مانند ناپیوستگی و ناهموار^۳ بودن را

با خود دارد. علاوه بر این، برای استفاده در دو متغیر یا بیش‌تر بسیار پیچیده است. به خاطر این

نوافص، در سه دهه‌ی گذشته، برآوردهای ناپارامتری مختلفی به کمک برآوردهای هموار^۴ ($f(x)$)

ارائه شده‌اند.

¹ Van Ryzin

² Scott

³ rough

⁴ smooth

۱-۳ براورد ناپارامتری چگالی

هیچ روش منحصر به فردی برای انجام برآورد ناپارامتری چگالی وجود ندارد. ما در این فصل، به معرفی هشت روش برآورد ناپارامتری تابع چگالی می‌پردازیم.

۱-۳-۱ رویکرد بافت‌نگار موضعی^۱

برای فهم برخی تکنیک‌های برآورد چگالی که بعداً مورد بحث قرار می‌گیرند، کار را با شرایطی آغاز می‌کنیم که در آن X یک متغیر تصادفی گستته است. فرض کنید یکی از مقادیری که این متغیر می‌تواند بگیرد، x باشد و هدف ما برآورد $f(x)$ از داده‌های x_i که $i = 1, \dots, n$ است.

برآورد $f(x)$ در حالت گستته، لزوماً، برآورد نسبت مقادیر x در جامعه‌ی X است. با توجه به داده‌های x_1, \dots, x_n یک برآورد ساده و معروف، عبارت از نسبت نمونه یعنی

$$\hat{f}_1(x) = n^* / n$$

اکنون، حالتی را در نظر می‌گیریم که در آن X یک متغیر تصادفی پیوسته است، یعنی احتمال برابر بودن x با x_i مساوی صفر است و $f(x)$ لازم است با میانگین‌گیری از x_i ‌ها برابر برآورد شود که در یک بازه در اطراف x ، مثلاً $x \pm h/2$ قرار دارد که در آن h عرض بازه است.

بنابراین براوردگر تجربی چگالی می‌تواند به صورت

$$\hat{f}_1(x) = (nh)^{-1} \sum_{i=1}^n I\left(x - \frac{h}{2} \leq x_i \leq x + \frac{h}{2}\right)$$

\mathcal{A} درست باشد، و در غیر صورت برابر صفر است. به صورت دیگر، می‌توانیم بنویسیم

¹ local histogram approach

$$\begin{aligned}\hat{f}_1(x) &= \frac{1}{nh} \sum_{i=1}^n I\left(-1/2 \leq \frac{x_i - x}{h} \leq 1/2\right) \\ &= \frac{1}{nh} \sum_{i=1}^n I(-1/2 \leq \psi_i \leq 1/2),\end{aligned}\quad (1-3-1)$$

که در آن $\psi_i = (x_i - x) / h$

توجه کنید که $\hat{f}(x)$ در (1-3-1) فراوانی نسبی به ازای هر واحد در بازه‌ی $(x - h/2, x + h/2)$ است که نقطه‌ی میانی آن x است. بنابراین به نظر می‌آید برآوردگر (1-3-1) تلاش می‌کند بافت‌نگاری بسازد که بر پایه‌ی مشاهدات موضعی نسبت به x بوده، در آن هر نقطه‌ی x مرکز بازه‌ی نمونه‌گیری است. عرض بازه، یعنی h ، مقداری را کنترل می‌کند که به وسیله‌ی آن، داده‌ها برای تولید برآورد (1-3-1) هموار (میانگین‌گیری) می‌شوند. هم چنین \hat{f}_1 مانند آن چه در فیکس و هاجز^۱ (۱۹۵۱) آمده، به عنوان برآوردگر «اویله»^۲ شناخته می‌شود. واضح است که تابع نشانگر یا تابع وزن $I(-1/2 < \psi_i < 1/2)$ در (1-3-1) به فاصله‌ی x_i از x بستگی دارد. اگر قدر مطلق این فاصله، کمتر از یا مساوی با $1/2$ باشد، وزن برابر ۱ و در غیر این صورت برابر صفر است. علاوه بر این موارد، تابع وزن $I(\psi) = I(-1/2 < \psi < 1/2)$

$$\begin{aligned}\int_{-\infty}^{\infty} I(\psi) d\psi &= \int_{-\infty}^{-1/2} I(\psi) d\psi + \int_{-1/2}^{1/2} I(\psi) d\psi + \int_{1/2}^{\infty} I(\psi) d\psi \\ &= \int_{-1/2}^{1/2} I(\psi) d\psi = \int_{-1/2}^{1/2} d\psi = 1.\end{aligned}\quad (2-3-1)$$

بنابراین

¹ Fix and Hodges
² naive

$$\begin{aligned} \int_{-\infty}^{\infty} \hat{f}_1(x) dx &= \frac{1}{nh} \sum_{i=1}^n \int_{-\infty}^{\infty} I\left(-\frac{1}{2} < \frac{x_i - x}{h} < \frac{1}{2}\right) dx \\ &= \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^{\infty} I\left(-\frac{1}{2} < \psi_i < \frac{1}{2}\right) d\psi_i = 1, \end{aligned} \quad (3-3-1)$$

و برآورد چگالی به این دلیل مناسب است که نامنفی و دارای انتگرال برابر با یک است. یک ویژگی (3-3-1) آن است که انتگرال روی x گرفته شده، زیرا می‌تواند فرض کند که x روی تمام برد X مقدار می‌گیرد.

۲-۳-۱ اقتباس صوری^۱

فرض کنید $F(x) = P(X \leq x)$ نشان‌دهنده‌ی تابع توزیع احتمال تجمعی X باشد. آن گاه تابع چگالی $f(x)$ به صورت زیر تعریف می‌شود:

$$\begin{aligned} f(x) &= \frac{d}{dx} F(x) = \lim_{h \rightarrow 0} \frac{F\left(x + \frac{h}{2}\right) - F\left(x - \frac{h}{2}\right)}{h} \\ &= \lim_{h \rightarrow 0} \frac{P\left(x - \frac{h}{2} < X < x + \frac{h}{2}\right)}{h}. \end{aligned} \quad (4-3-1)$$

مسئله‌ی ما برآورد $f(x)$ بر اساس x_1, \dots, x_n است. برای این منظور، فرض کنیم h تابع مثبتی از n باشد که با $\infty \rightarrow n$ به صفر میل می‌کند، و $P\left(x - \frac{h}{2} < X < x + \frac{h}{2}\right)$ را به وسیله‌ی نسبت مشاهدات x_1, \dots, x_n که در بازه‌ی $\left(x - \frac{h}{2}, x + \frac{h}{2}\right)$ قرار می‌گیرند، برآورد می‌کنیم.

بنابراین یک برآورد بدیهی و ساده‌ی $f(x)$ در (4-3-1) به صورت زیر است:

¹ formal derivation

$$\begin{aligned}
 \hat{f}_2(x) &= \frac{1}{nh} \left[\#(x_1 \dots x_n) \in \left(x - \frac{h}{2}, x + \frac{h}{2} \right) \right] \\
 &= \frac{1}{nh} \left[\# \left(\frac{x_1 - x}{h} \dots \frac{x_n - x}{h} \right) \in (-1/2, 1/2) \right] \\
 &= \hat{f}_1(x),
 \end{aligned} \tag{5-۳-۱}$$

که مانند (۱-۳-۱) است که در آن $\#\#$ نشان‌دهنده‌ی عملگر «تعداد عناصر موجود در \mathcal{A} » است. برآوردگر داده شده در (۵-۳-۱) را اوّلین بار فیکس و هاجز (۱۹۵۱) معرفی کردند.

۱-۳-۳-۱ برآوردگر هسته‌ی روزنبلات-پارزن^۱

برآوردگر چگالی که توسط تابع نشانگر معرفی شده در (۱-۳-۱) به دست می‌آید، این ویژگی را دارد که انتگرالش برابر یک است، اما نقصی هم دارد که ناهمواربودنش است. هم چنین $\hat{f}_1(x)$ یک تابع پیوسته نیست، و پرس‌هایی در نقاط $x / h \pm 2$ با مشتقات صفر در سایر نقاط دارد. این امر به برآوردها، ماهیتی پله‌ای می‌دهد، و شاید ما مجموعه‌ای از وزن‌های هموارتر را ترجیح دهیم. روزنبلات (۱۹۵۶) این مسئله را با جایگذاری یک تابع هسته‌ی مثبت K که در شرط زیر

صدق می‌کند، حل کرد:

$$\int_{-\infty}^{\infty} K(\psi) d\psi = 1. \tag{6-۳-۱}$$

برآورد کلی وی به صورت زیر است:

^۱ Rosenblatt-Parzen kernel estimator

$$\hat{f}(x) = \hat{f}_3(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x_i - x}{h}\right) = \frac{1}{nh} \sum_{i=1}^n K(\psi_i), \quad (7-3-1)$$

که در آن $\psi_i = h^{-1}(x_i - x)$ مانند (1-3-1) است و h ، عرض پنجره^۱ (هم‌چنین پارامتر هموارسازی یا پهنه‌ای نوار نیز نامیده می‌شود)، تابعی به اندازه‌ی n است و با $n \rightarrow \infty$ به صفر میل می‌کند.

تعدادی از ویژگی‌هایی که یک هسته باید دارا باشد می‌تواند از ماهیت تابع نشانگر استنباط شود. اولاً، برای مقادیر بزرگ $|\psi_i|$ (یعنی این که x_i به فاصله‌ی زیادی از x قرار می‌گیرد) باید $K(\psi_i)$ کوچک باشد، همین طور وزن‌های بسیار کوچک باید به چنین داده‌هایی در ساخت برآورد چگالی اختصاص یابند. به خصوص، چون وقتی $n \rightarrow \infty$ آن گاه $h \rightarrow 0$ ، نتیجه می‌شود که برای هر $x_i \neq x$ داریم $|\psi_i| \rightarrow \infty$ ، و بنابراین $K(-\infty) = K(\infty) = 0$ که از رابطه‌ی (1-3-6) حاصل می‌شود. این ویژگی، خاصیت صفر تابع نشانگر را دوباره به دست می‌دهد، که در آن قسمت یک با داشتن $\int_{-\infty}^{\infty} K(\psi) d\psi = 1$ نشان داده می‌شود. این امر منجر به جایگذاری یک مربع به مرکزیت x و به طول واحد با یک منحنی هموار می‌شود، هم‌چنین می‌توان آن را با مربعی به مرکزیت x به همان مساحت اما نه لزوماً با محمل کراندار جایگذاری کرد. علاوه بر این، چون این ویژگی‌ها همان ویژگی‌های تابع چگالی هستند، هسته‌ها اغلب همان توابع چگالی معروف در نظر گرفته می‌شوند، برای مثال، هسته‌ی نرمال استاندارد به صورت $K(\psi) = (2\pi)^{-1/2} \exp(-.5\psi^2)$ است. در این راستا، تابع نشانگر را می‌توان به صورت یک برآوردگر هسته با (\cdot) در نظر گرفت که روی $[2/1/2]$ دارای چگالی یکنواخت است.

^۱ band width