

## پیش‌گفتار

برای مقایسه کردن توزیع‌ها آزمون مبتنی بر گشتاورها مانند آزمون میانگین‌ها یا واریانس‌ها و یا چندک‌ها مفید هستند و به خاطر مقایسه کردن یک تعداد از گشتاورها، اگرچه در وقت صرفه جویی می‌شود اما این نوع آزمون ناسازگار است. همچنین بکاربردن آزمون‌های پارامتری برای تساوی دو تابع چگالی این ضرورت را ایجاد می‌کند که توزیع صفر مشخص باشد و در غیر این صورت نتایج اشتباهی بدست خواهد آمد. بنابراین با وجود این شرایط، آزمون‌های ناپارامتری وجود دارند که می‌توانند کار مقایسه دو تابع توزیع بدون مشکلات ذکر شده در بالا انجام دهد. آزمون‌های مبتنی بر هسته برای آزمون تساوی توابع چگالی در حالت متغیرهای پیوسته بررسی شده که به احمد-ون بل (۱۹۷۴)، مامن (۱۹۹۲)، فن و جنسی (۱۹۹۳)، لی (۱۹۹۶) می‌توان مراجعه کرد. همچنین روش هسته مبتنی بر فراوانی یک برآورد سازگار برای تابع چگالی احتمال توام با داده‌های آمیخته پیوسته و گسسته است که با استفاده از این برآوردگر می‌توانیم به آسانی آزمون مبتنی بر هسته برای تساوی دو تابع چگالی مجهول بسازیم. در این روش هموارسازی روی متغیر گسسته انجام نمی‌شود در صورتی که ایتچیسون و ایتکن (۱۹۷۶)، هال (۱۹۸۱)، گراند و هال (۱۹۹۳)، اسکات (۱۹۹۲)، سمینف (۱۹۹۶)، لی و راسین (۲۰۰۳) و هال و همکاران (۲۰۰۴ و ۲۰۰۷) روی هموارسازی متغیر گسسته و مزیت‌های آن بحث کردند. بنابراین آزمون‌های ناپارامتری مورد بررسی قرار گرفتند که متغیر گسسته را نیز هموار کنند که این کار توسط لی و همکاران در (۲۰۰۹) انجام شد و در این پایان‌نامه این روش را بررسی می‌کنیم. اساس این کار بر اساس برآورد هسته تابع چگالی احتمال می‌باشد که در آن هموارسازی هم روی متغیرهای پیوسته و هم روی متغیرهای گسسته انجام می‌شود. روش برآورد هسته ناپارامتری برای برآورد توزیع توام روی داده‌های چندگانه دودویی اولین بار توسط ایتچیسون و ایتکن در (۱۹۷۶) پیشنهاد شد و همچنین روش مبتنی بر درست‌نمایی برای انتخاب پارامترهای هموارسازی پیشنهاد کردند. امادر شبیه‌سازی نشان داده می‌شود که این برآوردگر برای داده‌های توام پیوسته و گسسته که متغیر پیوسته آن از توزیع دم پهن آمده باشد ناسازگار است که دلیل آن استفاده از روش درست‌نمایی برای انتخاب پارامتر هموارسازی است. لی و راسین در (۲۰۰۳) روش

هسته ناپارامتری با داده های آمیخته پیوسته و گسسته پیشنهاد کردند که پارامتر هموارسازی آن از طریق کمترین مربعات اعتبار مقطعی بدست می آید و مشکل برآوردگر ایتچیسون و ایتکن را نداشت، در این پایان نامه این برآوردگر را بررسی می کنیم. این پایان نامه شامل ۴ فصل می باشد که خلاصه مطالب هر فصل به شرح زیر است.

• در فصل ۱، روش های مختلف برآورد ناپارامتری تابع چگالی احتمال با داده های پیوسته و در حالت تک متغیره را بیان می کنیم.

• در فصل ۲، برآورد تابع چگالی احتمال با داده های توام پیوسته و گسسته را بیان می کنیم.

• در فصل ۳، آزمون ناپارامتری تساوی دو تابع چگالی احتمال با داده های توام پیوسته و گسسته را بیان می کنیم.

• در فصل ۴، برآورد ناپارامتری از تابع چگالی احتمال شرطی را بیان و آزمون ناپارامتری برای تساوی دو تابع چگالی احتمال شرطی را مطرح می کنیم.

در این پایان نامه مطالبی که با (\*) نشان داده شده، مطالبی است که در کتاب نامه آمده است ولی به اثبات پرداخته است یا برنامه شبیه سازی موجود نبوده است.

معصومه اکبری لاکه

شهریور ۱۳۸۹

# فهرست مندرجات

۶	۱	روش‌های مختلف برآورد ناپارامتری تابع چگالی احتمال برای متغیرهای پیوسته یک بعدی
۷	۱-۱	برآوردگر بافت نگار
۷	۱-۱-۱	برآوردگر بافت نگار
۸	۱-۱-۲	ویژگی‌های بافت نگار
۱۱	۱-۱-۳	انتخاب پارامتر هموارسازی، $h$
۱۲	۲-۱	برآوردگر فراوانی چند ضلعی
۱۲	۱-۲-۱	ویژگی‌هایی از برآوردگر فراوانی چند ضلعی
۱۳	۲-۲-۱	انتخاب پارامتر هموارسازی، $h$
۱۴	۳-۱	تغییر دادن پارامتر هموارسازی، $h$
۱۶	۴-۱	برآوردگر هسته چگالی
۱۶	۱-۴-۱	دلیل پیدایش برآوردگر هسته چگالی
۱۷	۲-۴-۱	چند ویژگی برآوردگر هسته

۲۰	انتخاب پارامتر هموارسازی $h$ در برآوردگر هسته	۳-۴-۱
۲۴	مشکلات برآوردگر هسته چگالی	۵-۱
۲۴	اریبی مرز	۱-۵-۱
۲۵	عدم تغییر پذیری موضعی در هموار کردن	۲-۵-۱
۲۷	همواری قله‌ها و دره‌ها	۳-۵-۱
۲۸	تعدیل‌ها و بهبودها برای برآوردگر هسته چگالی	۶-۱
۲۸	هسته‌های مرزی	۱-۶-۱
۳۰	تغییر دادن پارامتر هموارسازی، $h$	۲-۶-۱
۳۴	هسته‌های مرتبه بالاتر	۷-۱
۳۵	برآورد مبتنی بر تبدیل	۸-۱
۳۷	برآورد تابع چگالی احتمال با داده‌های آمیخته پیوسته و گسسته	۲
۳۸	مقدمه	۱-۲
۳۸	برآورد تابع جرم احتمال متغیرهای تصادفی چندگانه گسسته	۲-۲
۴۱	برآورد تابع چگالی احتمال با داده‌های آمیخته پیوسته و گسسته	۳-۲
۴۵	تعمیم متغیر تصادفی چندگانه کلی	۴-۲

۴۵	.....	بخش کامپیوتری مسئله	۵-۲
۴۶	.....	نتایج شبیه سازی مونت کارلو	۶-۲
۴۶	.....	عملکرد نمونه متناهی، توزیع های مستقل و هم توزیع (*)	۱-۶-۲
		عملکرد نمونه متناهی، کمترین مربعات اعتبار مقطعی در برابر	۲-۶-۲
		ماکزیمم دستنمایی اعتبار مقطعی با توزیع های دم پهن (*)	۴۷
۴۹	.....	پیوست A	۷-۲
۵۷	.....	پیوست B	۸-۲
۷۶		آزمون ناپارامتری تساوی دو تابع چگالی احتمال غیرشرطی	۳
۷۷	.....	مقدمه	۱-۳
۷۸		آزمون ناپارامتری برای تساوی توابع چگالی غیرشرطی با داده توام پیوسته و گسسته	۲-۳
۷۸	.....	آزمون تساوی دو تابع چگالی	۱-۲-۳
۸۸	.....	انتخاب پارامترهای هموارسازی	۳-۳
۹۵	.....	انواع دیگر تابع آزمون تساوی توابع چگالی	۴-۳
۹۸	.....	شیوه بوت استرپ	۵-۳

- ۶-۳ شیب سازى مونت كارلو . . . . . ۹۹
- ۱-۶-۳ آزمون تساوى توابع چگالى غير شرطى با داده هاى آميخته پيوسته  
و گسسته . . . . . ۱۰۰
- ۲-۶-۳ آزمون تساوى توابع چگالى غير شرطى تحت فرض مقابل توابع  
چگالى با تغييرات آهسته يا چگالى هاى هموار(\*) . . . . . ۱۰۲
- ۳-۶-۳ آزمون تساوى توابع چگالى غير شرطى تحت فرضيه مقابل توابع  
چگالى با تغييرات سريع يا چگالى هاى ناهموار(\*) . . . . . ۱۰۳
- ۴ برآورد ناپارامترى تابع چگالى شرطى و آزمون ناپارامترى تساوى دو تابع چگالى شرطى ۱۰۴
- ۱-۴ مقدمه . . . . . ۱۰۵
- ۲-۴ برآوردگر ناپارامترى تابع چگالى شرطى در حالت پيوسته و تک متغيره . . . . . ۱۰۶
- ۱-۲-۴ برآوردگر يك مرحله اى . . . . . ۱۰۶
- ۲-۲-۴ برآوردگر دو مرحله اى . . . . . ۱۰۶
- ۳-۴ انتخاب پارامتر هموارسازى در دو برآوردگر بالا . . . . . ۱۰۸
- ۴-۴ برآورد تابع چگالى احتمال شرطى با داده آميخته پيوسته و گسسته . . . . . ۱۰۹
- ۱-۴-۴ بررسى قاعده اعتبار مقطعى . . . . . ۱۱۰
- ۵-۴ آزمون ناپارامترى تساوى دو تابع چگالى شرطى با داده هاى آميخته پيوسته و  
گسسته . . . . . ۱۱۳

- ۱-۵-۴ آماره آزمون  $J_n$  برای آزمون ناپارامتری تساوی توابع چگالی شرطی ۱۱۴
- ۲-۵-۴ آماره آزمون  $J_{n,\lambda}$  برای آزمون ناپارامتری تساوی توابع چگالی شرطی ۱۱۶
- ۶-۴ شبیه سازی مونت کارلو ۱۱۷

## فصل ۱

# روش‌های مختلف برآورد ناپارامتری تابع چگالی احتمال برای متغیرهای پیوسته یک بعدي

۱-۱ برآوردگر بافت نگار

۲-۱ برآوردگر فراوانی چند ضلعی

۳-۱ تغییر دادن پارامتر هموارسازی  $h$

۴-۱ برآوردگر هسته چگالی

۵-۱ مشکلات برآوردگر هسته چگالی

۶-۱ تعدیل‌ها و بهبودها در برآوردگر هسته چگالی احتمال

۷-۱ هسته‌های مرتبه بالاتر

۸-۱ برآورد مبتنی بر تبدیل

## ۱-۱ برآوردگر بافت نگار

یک مفهوم اساسی در آنالیز داده تک متغیره بحث تابع چگالی احتمال است. اگر  $X$  متغیر تصادفی با تابع چگالی احتمال  $f(x)$  باشد با استفاده از آن می‌توانیم توزیع  $X$  را توصیف کنیم و همچنین احتمال‌هایی را با استفاده از روابط زیر بدست آوریم

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

یک انگیزه برای ساختن برآورد ناپارامتری تابع چگالی احتمال با استفاده از تعریف تابع چگالی احتمال است

$$f(x) = \frac{d}{dx} F(x) = \lim_{h \rightarrow 0} \frac{F(x+h) - F(x)}{h}$$

### ۱-۱-۱ برآوردگر بافت نگار

اگر  $x_1, \dots, x_n$  یک نمونه تصادفی به حجم  $n$  از تابع چگالی احتمال  $f$  باشد و  $F(x)$  را با تابع توزیع تجربی،  $\hat{F}(x) = \frac{\#\{x_i \leq x\}}{n}$ ، جایگذاری کنیم برآوردگر بافت نگار بصورت زیر خواهد بود.

$$\hat{f}(x) = \frac{(\#\{x_i \leq b_{j+1}\} - \#\{x_i \leq b_j\})/n}{h} \quad x \in (b_j, b_{j+1}]$$

یا

$$\hat{f}(x) = \frac{n_j}{nh} \quad x \in (b_j, b_{j+1}] \quad (1-1)$$

$b_1, \dots, b_n$  بیانگر باندها است و  $n_j$  تعداد مشاهداتی است که در باندها قرار می‌گیرد و

$$h = b_{j+1} - b_j \text{ است.}$$

بدون شک رایج‌ترین برآوردگر چگالی تک متغیره، بافت نگارها هستند زیرا محاسبه آسان و ساختمان و تفسیری ساده دارند و نیاز به ابزار پیشرفته گرافیکی هم نیست.

### ۲-۱-۱ ویژگی‌های بافت نگار

از رابطه (۱-۱) واضح است که ویژگی‌های برآوردگر بافت نگار به  $h$  یا به عبارت دیگر به تعداد باندها وابسته است. در شکل ۱-۱ نمودارهای  $c, b, a$  بافت نگارهایی با  $20$  باند هستند که هر سه بطور تصادفی از توزیع نرمال استاندارد به حجم  $100$  تولید شده‌اند در حالی که نمودارهای  $d, e, f$  بافت نگارهایی با  $4$  باند هستند و چگالی نرمال واقعی بر روی هر نمودار نیز رسم شده است.

نمودارهای  $c, b, a$  ناهمواری زیادی را نشان می‌دهد. در این نمودارها اگرچه ارتفاع باندها از چگالی واقعی پیروی می‌کند ولی آنها بطور قابل توجهی از یک نمودار به نمودار دیگر تغییر می‌کنند یعنی برآوردگر حاصل اریبی کم و تغییرپذیری زیاد دارد. در مقایسه، نمودارهای  $d, e, f$  همواری زیادی را نشان می‌دهند. در این نمودارها ارتفاع باندها بطور معمولی از یک نمودار به نمودار دیگر ثابت است اما آنها خیلی خوب از چگالی واقعی پیروی نمی‌کنند یعنی برآوردگر تغییرپذیری کم و اریبی بزرگ دارد. یک روش برای ارزیابی  $\hat{f}(x)$  از طریق اندازه تفاوت آن با  $f(x)$  است و ساده ترین اندازه هم مربعات خطاست که بصورت زیر می‌باشد.

$$SE(x)^1 = [\hat{f}(x) - f(x)]^2$$

$$MSE^2 = E_f[f(\hat{x}) - f(x)]^2$$

$$ISE^3 = \int_{-\infty}^{\infty} [\hat{f}(u) - f(u)]^2 dx$$

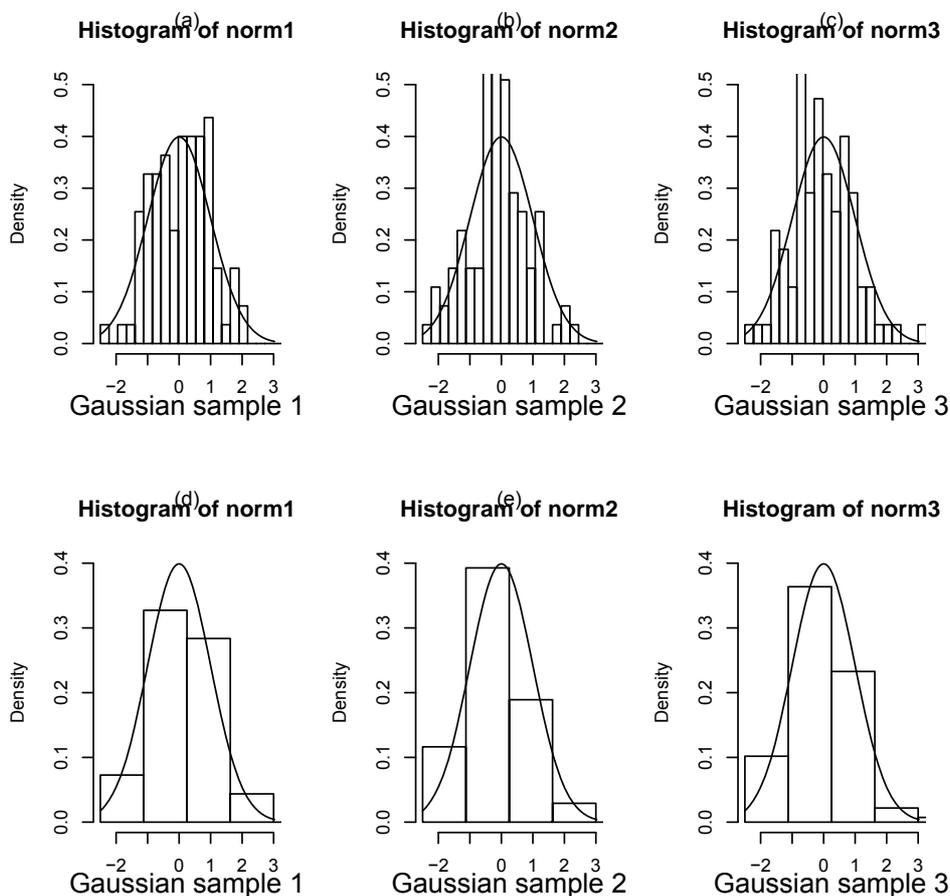
---

square error<sup>1</sup>

mean square error<sup>2</sup>

integral square error<sup>3</sup>

فصل ۱ روش‌های مختلف برآورد ناپارامتری تابع چگالی احتمال برای متغیرهای پیوسته یک بعدی ۹



شکل ۱-۱: برآوردگر بافت نگار با داده‌های تصادفی نرمال

چگالی  $f(x)$  را هموار گوئیم هرگاه  $f'(x)$  مطلقاً پیوسته و توان دوم آن انتگرال پذیر باشد در اینصورت می‌توان نشان داد که  $MSE$  مجانبی برابر است با

$$Bias[\hat{f}(x)] = E_f[\hat{f}(x)] - f(x) = \frac{1}{2} f''(x) [h - 2(x - b_j)] + O(h^2) \quad x \in (b_j, b_{j+1}]$$

$$var(\hat{f}(x)) = \frac{f(x)}{nh} + O(n^{-1})$$

فصل ۱ روش‌های مختلف برآورد ناپارامتری تابع چگالی احتمال برای متغیرهای پیوسته یک بعدی ۱۰

$$\begin{aligned}MSE(\hat{f}(x)) &= \text{var}(\hat{f}(x)) + [\text{Bias}(\hat{f}(x))]^2 \\ &= \frac{f(x)}{nh} + \frac{[f'(x)]^2}{4} [h - 2(x - b_j)]^2 + O(n^{-1}) + O(h^2) \quad (2-1)\end{aligned}$$

سرانجام با انتگرال گرفتن بر روی همه باندها  $MISE$  به این صورت است

$$MISE = \frac{1}{nh} + \frac{h^2 R(f'(x))}{12} + O(n^{-1}) + O(h^2) \quad (3-1)$$

که  $R(g) = \int g^2(u) du$  بیانگر است.

مبادله اریبی و واریانس که قبلاً توجه کردیم در رابطه (۳-۱) نیز دیده می‌شود.  $h$  بطور مستقیم وابسته به انتگرال توان دوم اریبی  $(\frac{h^2 R(f'(x))}{12})$  هست و رابطه معکوس با انتگرال واریانس  $(\frac{1}{nh})$  دارد. یعنی باندهای تنگ تر منجر به برآوردگری می‌شود که اریبی آن کم تر است اما تغییر پذیری زیاد دارد در نتیجه ساختن باندها با پهنای بیش تر اریبی را زیاد می‌کند.

به آسانی قابل محاسبه است  $h$  که  $MISE$  را مینیمم می‌کند برابر است با

$$h_O = \left[ \frac{6}{R(f'(x))} \right]^{\frac{1}{5}} n^{-1/5} \quad (4-1)$$

با جایگذاری این  $h$  بهینه در رابطه (۳-۱) داریم

$$AMISE_O = \left[ \frac{9R(f'(x))}{16} \right]^{\frac{1}{5}} n^{-\frac{2}{5}} \quad (5-1)$$

رابطه (۵-۱) نشان می‌دهد که سرعت همگرایی برآوردگر بافت نگار تابعی از  $n^{-\frac{2}{5}}$  است. در ضمن  $AMISE$  مجانبی بودن  $MISE$  را نشان می‌دهد.

میانگین انتگرال مربع خطا

### ۳-۱-۱ انتخاب پارامتر هموارسازی، $h$

معادله (۴-۱) یک روش برای انتخاب کردن  $h$  در برآوردگر بافت نگار با پهنای ثابت است اما این روش به چگالی  $f$  بستگی دارد و ساده ترین روش برای انتخاب  $h$ ، آن است که یک چگالی مرجع  $f$  انتخاب کنیم و سپس در رابطه (۵-۱) جایگذاری کنیم.

اگر چگالی مرجع را نرمال انتخاب کنیم  $h$  که  $AMISE$  را مینیمم می‌کند به صورت زیر است

$$h_O = 3.491 \sigma n^{-\frac{1}{5}} \quad (6-1)$$

در این روش نیاز است که  $\sigma$  را برآورد کنیم.

مانع اصلی در انتخاب  $h$  بر اساس (۶-۱) این است که توجیه نظری برای این روش وقتی که چگالی اصلی نرمال نباشد، وجود ندارد.

اگر اندازه  $ISE$  را دوباره در نظر بگیریم

$$\begin{aligned} ISE &= \int [f(u) - \hat{f}(u)]^2 du = \int f^2(u) du + \int \hat{f}^2(u) du - 2 \int \hat{f}(u) f(u) du \\ &= R(f) + R(\hat{f}(u)) - 2 \int \hat{f}(u) f(u) du \end{aligned} \quad (7-1)$$

$R(f)$  به  $\hat{f}$  بستگی ندارد پس در انتخاب  $h$ ، تأثیری ندارد و جمله آخر (۷-۱) به ظاهر  $-2E(\hat{f}(u))$  هست که باید برآورد شود.

یک روش برای انجام کار قاعده اعتبار مقطعی<sup>۵</sup> است.

اگر  $\hat{f}_{-i}(x_i)$  برآوردگر بافت نگار  $f$  با حذف مشاهده  $i$  ام باشد آنگاه

$$E(\hat{f}_{-i}(x_i)) = E(\hat{f}(x))$$

یا

$$E\left(\frac{1}{n} \sum_{i=1}^n \hat{f}_{-i}(x_i)\right) = E(\hat{f}(x)) = \int \hat{f}(u) f(u) du$$

## فصل ۱ روش‌های مختلف برآورد ناپارامتری تابع چگالی احتمال برای متغیرهای پیوسته یک بعدی ۱۲

بنابراین یک روش برای انتخاب  $h$ ، قاعده اعتبار مقطعی است.

$$CV = R(\hat{f}) - \frac{2}{n} \sum_{i=1}^n \hat{f}_{-i}(x_i)$$

### ۲-۱ برآوردگر فراوانی چند ضلعی

#### ۱-۲-۱ ویژگی‌هایی از برآوردگر فراوانی چند ضلعی

اگرچه بافت نگارها مفید هستند اما بدلیل ناهمواریشان بدنبال برآوردگرهای هموارتری هستیم. یک روش برای ساختن برآوردگر هموارتر این است که نقاط میانی بازه‌ها را با خط‌های مستقیم به هم وصل کنیم.

این برآوردگر را برآوردگر فراوانی چند ضلعی<sup>۶</sup> می‌نامیم. اگرچه پیوسته است اما مشتقات آن در نقاط میانی باندها تعریف نشده است.

اگر باندها با طول ثابت  $h$  را در نظر بگیریم و  $\{b_0, \dots, b_{k+1}\}$  ابتدای باندها را نشان دهد و  $C_j$  ها را بصورت زیر تعریف کنیم

$$\begin{aligned} C_j &= \frac{b_j + b_{j+1}}{2} & j = 1, \dots, k \\ C_{k+1} &= b_{k+1} + \frac{h}{2} \\ C_0 &= b_1 - \frac{h}{2} \end{aligned}$$

آنگاه برآوردگر فراوانی چند ضلعی بصورت زیر تعریف می‌شود

$$\hat{f}(x) = \frac{1}{nh} [n_j c_{j+1} - n_{j+1} c_j + (n_{j+1} - n_j)x], \quad x \in [c_j, c_{j+1}] \quad (8-1)$$

که در آن  $n_0 = n_{k+1} \equiv 0$  است.

آنالیز مجانبی نشان می‌دهد که این بهبود ظاهری می‌تواند بهبود دقت برآوردگر را در برداشته باشد.

<sup>۶</sup>frequency polygon

## فصل ۱ روش‌های مختلف برآورد ناپارامتری تابع چگالی احتمال برای متغیرهای پیوسته یک بعدی ۱۳

فرض کنید مشتق دوم  $f$  مطلقاً پیوسته و  $R(f')$  و  $R(f'')$  متناهی باشند در این صورت  $MISE$  خواهد شد

$$MISE = \frac{2}{3nh} + \frac{49h^4 R(f'')}{2880} + O(n^{-1}) + O(h^6) \quad (9-1)$$

جمله اول (واریانس) به طور معکوس با  $h$  و جمله دوم (توان دوم اریبی) بطور مستقیم با  $h$  تغییر می‌کند. تفاوت قابل ملاحظه رابطه (۹-۱) نسبت به رابطه (۳-۱) این است که توان دوم اریبی از  $O(h^2)$  به  $O(h^4)$  تبدیل شده است. یعنی اگر  $h$  را کوچک کنیم اریبی با سرعت بیشتری کم می‌شود. با مینیمم کردن (۹-۱)  $h$  بهینه به صورت زیر بدست می‌آید

$$h_O = 2 \left[ \frac{15}{49} R(f'') \right]^{\frac{1}{5}} n^{-\frac{1}{5}} \quad (10-1)$$

پس

$$AMISE_O = \frac{5}{12} \left[ \frac{49 R(f'')}{15} \right]^{\frac{1}{5}} n^{-\frac{4}{5}}$$

مشاهده می‌کنیم  $AMISE$  با سرعت  $O(n^{-\frac{4}{5}})$  به صفر میل می‌کند که در مقایسه با  $AMISE$  در برآوردگر بافت نگار بیشتر شده است. مقایسه کردن (۱۰-۱) و (۴-۱) نشان می‌دهد که  $h$  بهینه برای برآوردگر فراوانی چند ضلعی فرم متفاوت از بافت نگار دارد و بطور مجانبی بزرگتر می‌شود.

### ۲-۲-۱ انتخاب پارامتر هموارسازی، $h$

ساده ترین روش برای انتخاب  $h$  برای برآوردگر فراوانی چند ضلعی جایگذاری یک فرم ویژه از  $f$  در (۱۰-۱) هست. اگر  $f$  رانرمال بگیریم،  $h$  بهینه به این صورت می‌شود

$$h_O = 2.15 \sigma n^{-\frac{1}{5}}$$

### ۳-۱ تغییر دادن پارامتر هموارسازی، $h$

تا اینجا برآوردگرهای بافت نگار و فراوانی چند ضلعی براساس  $h$  ثابت بررسی شدند. می دانیم که  $h$  مبادله اریبی و واریانس در هر نقطه را کنترل می کند وهم چنین این مبادله به ویژگی های موضعی چگالی وابسته است.

معادله (۱-۲) نشان می دهد که ویژگی های موضعی چگالی در هر نقطه  $x$ ، دقت بافت نگار در  $x$  را تعیین می کند. با توجه به آنچه گفته شد  $h$  در مناطقی با چگالی زیاد باید بزرگتر باشد تا جمله اول یعنی واریانس در  $MSE$  کاهش یابد و باید بطور عکس نسبت به  $|f'(x)|$  تغییر کند تا جمله دوم اریبی مینیمم شود.

تعریف بافت نگار که  $h$  آن به طور موضعی تغییر می کند به صورت زیر است

$$\hat{f}(x) = \frac{n_j}{n(b_{j+1} - b_j)} \quad x \in (b_j, b_{j+1}] \quad (1-11)$$

حال مسئله، پیدا کردن باندهای مناسب  $b_1, \dots, b_{k+1}$  برای برآوردگراست.

$MISE$  مجانبی برآوردگر در باندی که شامل  $x$  هست فرم زیر را دارد

$$AMISE(x) = \frac{2f(x)}{3nh} + \frac{49h^4 f''(x)^2}{2880}$$

در نتیجه  $h$  موضعی بهینه عبارت است از

$$h_x = 2 \left[ \frac{15f(x)}{49f''(x)^2} \right]^{1/5} n^{-1/5}$$

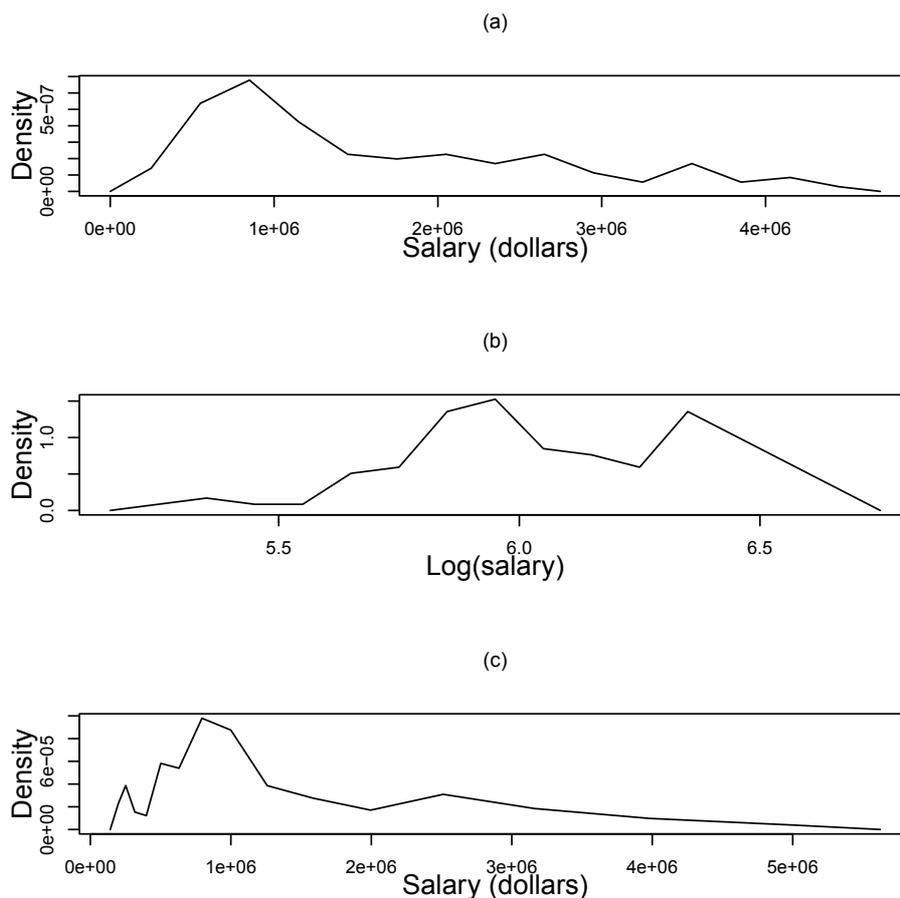
برآوردگر فراوانی چند ضلعی با  $h$  موضعی به صورت زیر تعریف می شود

$$\hat{f}(x) = \frac{1}{n(c_{j+1} - c_j)} \left[ \frac{n_j c_{j+1}}{b_{j+1} - b_j} - \frac{n_{j+1} c_j}{b_{j+2} - b_{j+1}} + \left( \frac{n_{j+1}}{b_{j+2} - b_{j+1}} - \frac{n_j}{b_{j+1} - b_j} \right) x \right] \quad x \in [c_j, c_{j+1}] \quad (1-12)$$

یک روش ساده برای ساختن برآوردگرهای بافت نگار یا فراوانی چند ضلعی با  $h$  موضعی که در عمل هم خوب کار می کند، این است که ابتدا داده ها را به یک مقیاس دیگری تبدیل کنیم و سپس این

فصل ۱ روش‌های مختلف برآورد ناپارامتری تابع چگالی احتمال برای متغیرهای پیوسته یک بعدی ۱۵

داده‌های تبدیل شده را هموار کنیم. سپس  $h_j$  های ساخته شده را به مقیاس اولیه برمی گردانیم و سپس از رابطه‌های (۱-۱۱) و (۱-۱۲) برای برآورد می‌توانیم استفاده کنیم. البته داده‌ها باید تحت تبدیل یکنوا تغییر مقیاس دهند.



شکل ۱-۲: برآوردگر فراوانی چندضلعی برای داده حقوق بازیکنان بیسبال

شکل ۱-۲ چگونگی روش انجام کار را شرح می‌دهد. داده‌ها مربوط به حقوق ۱۱۸ بازیکن بیسبال لیگ میجر در سال ۱۹۹۳ هستند.

در نمودار  $a$ ، برآوردگر فراوانی چند ضلعی در مقیاس اصلی داده‌ها با  $h$  ثابت  $3000000$  دلار محاسبه و رسم شده است. چنانکه شکل نشان می‌دهد سمت راست توزیع بلند است یعنی چوله به

## فصل ۱ روش‌های مختلف برآورد ناپارامتری تابع چگالی احتمال برای متغیرهای پیوسته یک بعدی ۱۶

راست است. بیشترین مد حدود ۱ میلیون دلار هست. اما به دلیل ناهموازی در دم نمی‌توان ساختار دیگری در داده‌ها را تعیین کرد. بلند بودن دم، تبدیل لگاریتمی را پیشنهاد می‌کند. با این تبدیل ساختاری با عدم تاکید بلند بودن دم در شکل  $b$  قابل مشاهده است. در نمودار  $c$ ، پس از تبدیل به مقیاس اصلی (با نمایی گرفتن) برآوردگر فراوانی چند ضلعی با  $h$  موضعی رسم شده است. اگرچه در این نمودار هنوز دم سمت راست بلند است اما علاوه بر مد ۱ میلیون دلاری، مدهای کوچکتری در ۲۵۰۰۰۰ دلار و ۵۰۰۰۰۰ دلار و ۲۵۰۰۰۰۰ دلار قابل مشاهده است.

### ۴-۱ برآوردگر هسته چگالی

#### ۱-۴-۱ دلیل پیدایش برآوردگر هسته چگالی

اگرچه برآوردگرهای معرفی شده در بخش‌های قبلی مفید هستند اما دو اشکال عمده دارند، یکی اینکه هموار نیستند و دیگر آنکه به قدر کافی به ویژگی‌های محلی  $f$  حساس نیستند. برای این منظور تعریف تابع چگالی را به صورت زیر در نظر می‌گیریم

$$f(x) = \frac{dF(x)}{dx} = \lim_{h \rightarrow 0} \frac{F(x+h) - F(x-h)}{2h} \quad (13-1)$$

بافت نگار (۱۳-۱) را با تقسیم کردن خط به باندها برآورد می‌کند اما یک روش منطقی این است که بطور جداگانه در هر نقطه  $x$  تابع چگالی را برآورد کنیم. با جایگذاری تابع توزیع تجربی به جای  $F(x)$  داریم

$$\hat{f}(x) = \frac{\#\{x_i \in (x-h, x+h)\}}{2nh}$$

که می‌تواند به این صورت هم نوشته شود

$$\hat{f}(x) = \frac{1}{nh} \sum_{i=1}^n k\left(\frac{x-x_i}{h}\right), \quad (14-1)$$

## فصل ۱ روش‌های مختلف برآورد ناپارامتری تابع چگالی احتمال برای متغیرهای پیوسته یک بعدی ۱۷

که  $k$  بصورت زیر تعریف می‌شود

$$k(u) = \begin{cases} 1/2 & -1 < u \leq 1 \\ 0 & \text{otherwise.} \end{cases}$$

فرم (۱-۱۴) برآوردگر هسته چگالی با تابع هسته  $k$  می‌باشد. توجه کنید که این تابع هسته، تابع چگالی یکنواخت روی  $[-1, 1]$  هست.

یک تابع هسته هموارتر به برآورد هسته چگالی هموارتر منجر خواهد شد. شکل ۱-۳ برآوردگر هسته برای داده‌های نرخ سپرده سه ماهه ۶۹ بانک ایسلند با استفاده از چگالی نرمال برای  $k$  و  $h$  های  $0.08$  و  $0.04$  و  $0.16$  رسم شده است. فقط به ازای  $h = 0.8$  منحنی بطور مطلوبی هموار شده است و یک فرم ۳ مدی با مدهای  $7/5$  و  $8$  و  $8/5$  درصد نشان می‌دهد.

### ۱-۴-۲ چند ویژگی برآوردگر هسته

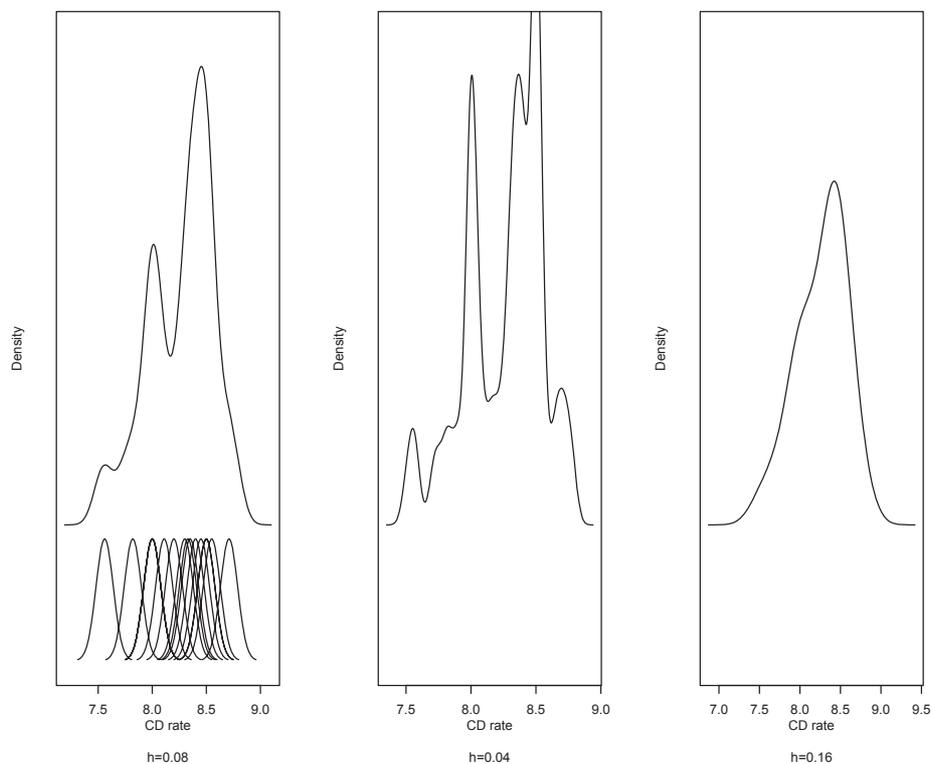
درجه هموارسازی داده‌ها تأثیر قوی روی ظاهر  $\hat{f}(x)$  دارد و از طریق تنظیم پارامتر هموارسازی  $h$  قابل کنترل است.

انتخاب پارامتر هموارسازی می‌تواند از طریق معیاری مانند  $MISE$  تعیین شود.

اگر  $\hat{f}$  برآوردگر هسته چگالی باشد و تابع هسته  $k$  در شرایط زیر صدق کند آنگاه  $MSE$ ، برآوردگر هسته بصورت زیر خواهد بود

$$\int k(u)du = 1, \quad \int uk(u)du = 0, \quad \int u^2 k(u)du = \sigma_k^2 > 0$$

فصل ۱ روش‌های مختلف برآورد ناپارامتری تابع چگالی احتمال برای متغیرهای پیوسته یک بعدی ۱۸



شکل ۱-۳: برآوردگر هسته داده‌های نرخ سپرده سه ماهه ۶۹ بانک ایسلند

فرض کنید که چگالی اصلی به قدر کافی هموار باشد یعنی  $f''$  بطور مطلق پیوسته و  $f'''$  انتگرال پذیر باشد. اگر  $h \rightarrow 0$  و  $nh \rightarrow \infty$  وقتی که  $n \rightarrow \infty$  می‌رود، سپس با بسط سری تیلور داریم