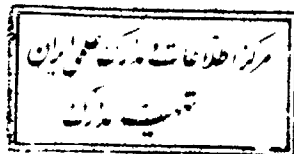


۹۲۲۸ / ۱۰ / ۸

بسم الله الرحمن الرحيم



استنباطهایی برای داده های بریده شده

بوسیله
حسین ریاض الشمس


پایان نامه

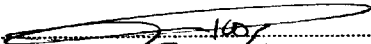
ارائه شده به دانشکده تحصیلات تکمیلی به عنوان بخشی از فعالیتهای
تحصیلی لازم برای اخذ درجه کارشناسی ارشد

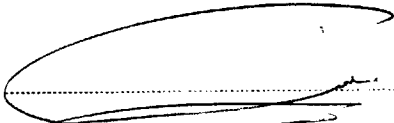
در رشته
آمار
از
دانشگاه شیراز
شیراز، ایران

۱۴۵۴۵

ارزیابی و تصویب شده توسط کمیته پایان نامه با درجه:
امضاء اعضاء کمیته پایان نامه:

.....  دکتر فریبرز حیدری، استاد یار بخش آمار (رئیس کمیته)

.....  دکتر عبدالرسول برهانی، استادیار بخش آمار

.....  دکتر احمدرضا سلطانی، اسناد بخش آمار

مهر ۷۸

۲۷۳۰۹

تقدیم به:

مادر مهربانم

۲۷۳۰۹

سیاسگزاری

با سپاس و درود خداوندی را که جان را فکرت آموخت و روان را با دم وجودش زندگی بخشید تا آدمیان ناچیز خاکی را تا به مرتبه اعلی فکرت رساند که کائنات هستی را قدرت ادراک آدمی جایگاهش باشد. در این برهه از زمان که قدرت ادراک آدمی مرز زمان و مکان را در نور دیده چه خوش است آنان را که در جهت تعقل ره پویند. و چه بزرگی، شکیبایی، پشتکار و از خود گذشتگی است افرادی چونان دکتر فریبرز حیدری را که شیوه تعقل را پایه گزارند. و زبان مرا یارای قدرتانی از چنان استاد فرزانه نباشد. همانقدر سپاس دارم که توانم که نه تنها علم آمار، که سخاوت، انسانیت، محسنات و تعقل را وارث ایشان باشم. سپاس دارم از اساتید فرزانه جناب دکتر احمد رضا سلطانی و دکتر عبدالرسول برهانی که در امر نگارش پایان نامه مرا یاری نمودند. و سپاس دارم یگانه مرد بزرگوار پروفیسور جواد بهبودیان استاد مشاور بنده. و از کلیه اساتید که در دوران تحصیل چراغ بر افروخته شان روشنی بخش راه من بوده.

چکیده

استنباط‌هایی برای داده‌های بریده شده

توسط

حسین ریاض الشمس

در سانسور تصادفی دیده شده که برآورد درست‌نمایی ماکزیمم (MLE) منحنی بقاء از فرض پارامتری بودن توزیع متغیر سانسور تاثیر نمی پذیرد. برآورد روش کاپلان مایر در ۱۹۵۸ یک برآورد MLE هم برای هر دو مدل نا پارامتری و هم مدل نیمه پارامتری می باشد. در داده های تصادفی بریده شده برآورد حدی ضربی که توسط لیندل بل در ۱۹۷۱ ارائه شده برآوردی MLE برای مدل نا پارامتری است. و آن مدل نیمه پارامتری را که در آن مکانیزم برش پارامتری باشد را در بر نمی گیرد. فرض کنید X, Y دو متغیر تصادفی مثبت باشند. متغیر تصادفی X ، (متغیر مورد بررسی) از چپ بوسیله Y ، (متغیر برش) بریده شده (یا Y از راست بوسیله X بریده شده) اگر زوج (X, Y) وقتی مشاهده شود که $X > Y$ ، چنانچه توزیع Y از یک خانواده پارامتری معلوم باشد مدل نیمه پارامتری نامیده می شود و چنانچه توزیع Y کاملاً معلوم باشد مدل ناپارامتری است. بریدگی در بسیاری از مطالعات ستاره شناسی، همه گیرشناسی، پزشکی، علوم اجتماعی، جرم شناسی و غیره، رخ می دهد. برای برآورد منحنی بقاء روشهایی از قبیل ناپارامتری، نیمه پارامتری، رگرسیون، فرایندهای مارکف ارائه شده است. بسیاری از تحقیقات بر اساس ماکزیمم کردن تابع درست‌نمایی بنا نهاده شده و مانیز به طور گسترده چنین روشهایی را بکار خواهیم برد. همچنین فرض مشخص بودن تابع توزیع متغیر برش برآوردهای بهتری از حالت نا مشخص تابع توزیع بدست می دهد. اگر متغیر برش از خانواده پارامتری باشد در روش درست‌نمایی ماکزیمم پارامترها را نیز برآورد می کنیم.

فهرست مطالب

صفحه	عنوان
۱	فصل اول: مقدمه
۱	۱-۱- مقدمه و خلاصه
۱	۲-۱- تاریخچه
۱	۳-۱- تعاریف
۲	۴-۱- طرح مسئله
۴	فصل دوم: مدل‌های ریاضی
۴	۱-۲- مقدمه و خلاصه
۵	۲-۲- توابع توزیع، مدل ناپارامتری
۱۱	۳-۲- مدل نیمه پارامتری
۱۲	۴-۲- نمونه‌های تصادفی از داده‌های بریده شده
۱۲	۵-۲- تابع درست‌نمایی، مدل ناپارامتری
۱۴	۶-۲- تابع درست‌نمایی، مدل نیمه پارامتری
۱۶	۷-۲- پیش‌زمینه‌های ریاضی
۱۸	فصل سوم: برآورد
۱۸	۱-۳- مقدمه و خلاصه
۱۸	۲-۳- برآورد درست‌نمایی ماکزیمم نیمه پارامتری
۲۰	۳-۳- خواص حدی
۲۵	۴-۳- تحلیل اطلاعات
۲۹	۵-۳- مدل نمایی
۳۲	۶-۳- مدل وایبل
۳۴	فصل چهارم: فاصله اطمینان
۳۴	۱-۴- مقدمه و خلاصه
۳۴	۲-۴- روش نسبت درست‌نمایی حاشیه‌ای (MLR)
۳۵	۳-۴- فاصله اطمینان برای $F(x)$

صفحه

عنوان

۳۹ ۴-۴- فاصله اطمینان برای احتمال برش

۴۱ ۵-۴- روش درستمایی نیمه پارامتری (SLR)

۴۲ ۶-۴- فاصله اطمینان درستمایی نیمه پارامتری برای $F(x)$

۵۳ ۷-۴- استنباط در مورد احتمال برش

فهرست مراجع

صفحه چکیده و صفحه عنوان به زبان انگلیسی

۱- فصل اول

بریدگی

۱-۱- مقدمه و خلاصه

تجزیه و تحلیل داده های بریده شده در سالهای اخیر یک موضوع کاربردی همراه با پشتوانه قوی ریاضی به صورت یک ابزار بوسیله آماردانان کاربردی در آمده. البته مباحثی که جدیدتر از آن تحت عنوان سانسور تاخیری بوسیله کوپاس و حیدری ۱۹۹۷ مطرح شده است بریدگی را نیز در بر می گیرد. هدف این فصل معرفی بریدگی و داده های بریده شده است. همچنین طرح مسائلی که در آنالیز بقاء مطرح است.

۱-۲- تاریخچه

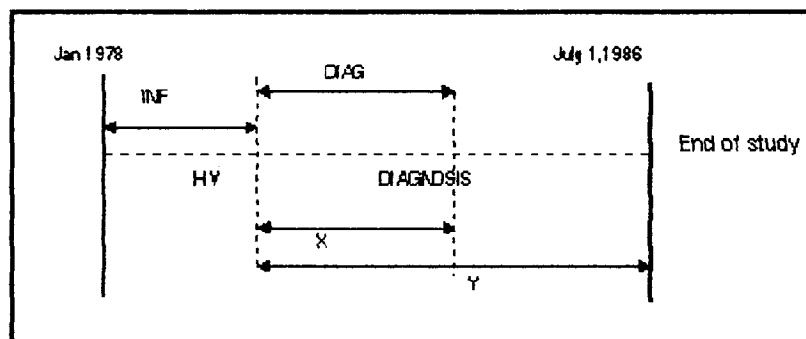
بریدگی اولین بار توسط لیندن بل (۱۹۷۱) در یک بررسی نجومی معرفی شد که البته می توان آنرا به عنوان روش ناپارامتری در نظر گرفت. اولین بحث تئوری بر اساس نسبت درستمایی ناپارامتری را آوین (۱۹۹۰، ۱۹۸۸) ارائه کرده، و از آن به بعد نظرات را متوجه خود نمود. برآورد ناپارامتری درستتایی ماکزیمم برای بقاء بوسیله گانگ لی (۱۹۹۵) بحث شده و برآورد احتمال بقاء و احتمال برش در حالت نیمه پارامتری بوسیله مئی چانگ وانگ (۱۹۸۹) و فاصله اطمینان برای احتمال بقاء و احتمال برش بوسیله گانگ لی و همکاران (۱۹۹۷) بدست آمده.

۱-۳- تعاریف

فرض کنید X, Y دو متغیر تصادفی مثبت باشند که به ترتیب نشان دهنده طول عمر (X) و زمان برش (Y) مربوط به یک موضوع مورد بررسی باشند. زمان طول عمر X را میگوییم از چپ بوسیله Y بریده شده (بنابراین Y از راست بوسیله X بریده شده) اگر زوج (X, Y) تنها وقتی قابل مشاهده باشند که $X > Y$. و یا میگوییم X از راست توسط Y بریده شده اگر زوج (X, Y) تنها وقتی قابل مشاهده باشند که $X < Y$.

بریدگی در بسیاری از موضوعها ممکن است رخ دهد مانند: ستاره شناسی، همه گیر شناسی، قابلیت اعتماد، مطالعات پزشکی، علوم اجتماعی، جرم شناسی، اقتصاد و غیره. به عنوان نمونه در مطالعات پزشکی وقتی می خواهیم طول زمان بقاء را بعد از

حمله یک ویروس به بدن بررسی کنیم، اگر X بیانگر زمان بین حمله ویروس به بدن تا مرگ باشد و زمان دنبال کردن بیماری Y واحد زمانی بعد از شروع حمله ویروس شروع شده باشد در اینجا به وضوح X از چپ بوسیله Y بریده شده است ($X > Y$). به عنوان مثال داده های مربوط به انتقال خون که به وسیله مرکز کنترل بیماریها در آتلانتا توسط کالب فلیش و لاولس در ۱۹۸۹ ارائه شده را در نظر بگیرید. متغیر مورد بررسی زمان کمون می باشد. متغیر مورد بررسی زمان کمون یعنی طول زمان از حمله ویروس HIV تا زمان تشخیص بیماری (X) میباشد. متغیر برش (Y) زمان از حمله ویروس به بدن تا پایان دوره مطالعه یعنی (اژولای ۱۹۸۶) میباشد. به انگاره زیر توجه کنید.



$$X = \text{DIAG}$$

$$\text{INF} + Y = 103 \Rightarrow Y = 103 - \text{INF}$$

در اینجا داده هایی قابل مشاهده هستند که در آنها $X < Y$ یعنی X زمان کمون از راست بوسیله Y زمان برش بریده شده.

۱-۴- طرح مسئله

اکنون فرض کنید G, F به ترتیب توابع توزیع غیر شرطی Y, X باشند. تابع

$$R(x) = 1 - F(x) \quad (1-1)$$

در آنالیز بقاء به عنوان تابع بقاء و در قابلیت اعتماد به عنوان تابع قابلیت اعتماد یک سیستم شناخته می شود و برآورد این تابع همیشه به عنوان مسئله اصلی مطرح است. در واقع این تابع در آنالیز بقاء برابر احتمال اینست که طول عمر یک مولفه بعد از شروع بیماری بیشتر از x واحد زمانی باشد و در قابلیت اعتماد برابر احتمال اینست که سیستم بیشتر از یک زمان معین x کار کند.

همچنین تابع

$$v = p(X < Y) \quad (2-1)$$

احتمال برش نام دارد. معیار v در کاربرد نقش مهمی جهت دانستن درصد موارد مشاهده نشده دارد.

در فصلهای آتی برآوردهایی برای v , $R(x)$ یا $F(x)$ ، همچنین فاصله اطمینان برای این معیارها ارائه خواهد شد.

۲- فصل دوم مدلهای ریاضی

۲-۱- مقدمه و خلاصه

فرض کنید زمان عمر X بوسیله Y از چپ بریده شده باشد. در حقیقت زوج مشاهده هم توزیع و مستقل $(X_1, Y_1), (X_2, Y_2), \dots, (X_n, Y_n)$ از تابع درستی نامی شرطی $L((X, Y) | X > Y)$ بدست آمده است. یعنی توزیع شرطی (X, Y) به شرط $X > Y$. مسئله برآورد احتمال برش $v = P(X < Y)$ و تابع بقاء $P(X \geq x)$ (ویا به عوض آن تابع توزیع $P(X \leq x)$) میباشد. فرض کنید G, F_0 به ترتیب توابع توزیع Y, X باشند. همانطور که وانگ در ۱۹۸۹ بیان داشته F_0 با شرط

$$\inf\{t; G(t) > 0\} > \inf\{t; F_0(t) > 0\}$$

قابل برآورد نیست. در این حالت ما در حقیقت تابع توزیع شرطی زیر را برآورد می کنیم.

$$P(X < x | X > c) = \frac{P(X < x \cap X > c)}{P(X > c)} = \frac{F_0(x) - F_0(c)}{1 - F_0(c)} \quad x \geq c \quad (1-2)$$

که در آن $c \geq \inf\{t; G(t) > 0\}$.

توجه کنید که $F = F_0$ اگر G, F_0 دارای محمل یکسان باشند و

$$c = \inf\{t; G(t) > 0\} \quad (2-2)$$

در عمل اغلب c را برابر مقدار زیر می گیریم.

$$c = \min\{X_i, i=1 \dots n\} \quad (3-2)$$

برای جلوگیری از بحثهای جزئی فرض می کنیم $P(X=Y)=0$ و $P(X>Y)>0$.

وودروف در ۱۹۸۵ در بررسی یک مسئله نجومی با استفاده از توزیعهای تجربی مربوط به X_1, X_2, \dots, X_n و Y_1, Y_2, \dots, Y_n در حالی که G, F کاملاً مجهول فرض شده بودند مدلهای ناپارامتری برای برآورد G, F ارائه نمود.

ولی ما فرض را بر معلوم بودن G می گذاریم چرا که در آنالیز بقاء همانطور که در فصلهای آینده خواهیم دید می توانیم توزیع متغیر Y (G) یعنی زمان حمله بیماری تا تشخیص بیماری را بدست و یا در مثال بخش ۱-۲ توزیع زمان از حمله بیماری تا پایان مطالعه کامل مشخص است. و یک چنین فرض اضافی منجر به

برآوردهای بهتری خواهد شد تا زمانی که G کاملاً ناشناخته است. چرا که در این حالت از اطلاعات موجود در متغیر برش Y استفاده می کنیم. در بسیاری موارد توزیع برش G متعلق به یک خانواده پارامتری

$$k = \{G(t, \theta), \theta \in \Theta \subset \mathbb{R}^q\}$$

می باشد. که در آن θ می تواند پارامتر q بعدی متعلق به Θ باشد. و برای هر θ ، $G(t, \theta)$ یک تابع توزیع شناخته شده است. این حالت مدل نیمه پارامتری نامیده می شود. مدل نیمه پارامتری اولین بار بوسیله وانگ در ۱۹۸۹ معرفی شده و از آن برای برآورد درستمایی ماکزیمم نیمه پارامتری نام برده می شود. برآوردهایی که بوسیله ماکزیمم کردن تابع درستمایی بدست می آیند توجه زیادی را به خود جلب کرده اند که ما نیز از این امکانات به طور وسیع استفاده خواهیم کرد. به طور کلی در مورد تجزیه و تحلیل داده های بریده شده از مدل های رگرسیونی، پارامتری، نیمه پارامتری و نا پارامتری استفاده شده است.

۲-۲- توابع توزیع، مدل نا پارامتری

فرض کنید X متغیر تصادفی طول عمر و Y متغیر برش به ترتیب با توزیعهای $F(x)$ ، $G(y)$ باشند. و $G(y)$ به پارامتر بستگی نداشته باشد. تعاریف:

$$F(x) = P(X \leq x)$$

$$G(y) = P(Y \leq y)$$

$$\alpha = P(Y < X) \quad (۴-۲)$$

$$v = P(X < Y) \quad (۵-۲)$$

$$H_*(x, y) = P(X \leq x, Y \leq y | Y < X) \quad (۶-۲)$$

$$F_*(x) = P(X \leq x | Y < X) \quad (۷-۲)$$

$$G_*(y) = P(Y \leq y | Y < X) \quad (۸-۲)$$

F_* ، G_* به ترتیب توابع توزیع کناری شرطی Y, X به شرط $Y < X$ میباشند. و $H_*(x, y)$ تابع توزیع شرطی (X, Y) به شرط $X > Y$ است.

قضیه -- ۱۱۲-۱: چنانچه X, Y دو متغیر تصادفی مثبت و مستقل باشند آنگاه روابط زیر برقرار است.

$$\alpha = \int_0^{\infty} G(x) dF(x) \quad (9-2)$$

$$= \int_0^{\infty} (1-F(v)) dG(v) \quad (10-2)$$

$$v = \int_0^{\infty} [1-G(x)] dF(x) \quad (11-2)$$

$$F_*(x) = \alpha^{-1} \int_0^x G(u) dF(u) \quad (12-2) \text{ توزیع:}$$

$$dF_*(x) = \alpha^{-1} G(x) dF(x) \quad (13-2) \text{ چگالی:}$$

$$G_*(y) = \alpha^{-1} \int_0^y (1-F(v)) dG(v) \quad (14-2) \text{ توزیع:}$$

$$dG_*(y) = \alpha^{-1} (1-F(y)) dG(y) \quad (15-2) \text{ چگالی:}$$

$$H_*(x, y) = \alpha^{-1} \int_0^x \int_0^y I(v < u) dG(v) dF(u) \quad (16-2)$$

$$= \alpha^{-1} \int_0^y (F(x) - F(v)) dG(v)$$

$$dH_*(x, y) = \alpha^{-1} dF(x) dG(y) I(x > y) \quad (17-2)$$

البته فرمول مربوط به $G_*(y)$ جهت کامل بودن مطلب در اینجا آورده شده وگرنه از این فرمول استفاده ای نخواهیم کرد.

هر کدام از عبارتهای بالا را می توان به چندین راه ثابت کرد ما در اینجا یک راه حل ساده به منظور قابل فهم بودن مطلب و همچنین یک راه حل ریاضی دقیقتر را بیان خواهیم کرد.

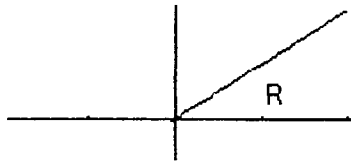
اثبات: راه حل اول.

$$\begin{aligned} \alpha = P(Y < X) &= \iint_R dF(x, y) dy dx \\ &= \iint_R dG(y) dF(x) \end{aligned}$$

که در آن R ناحیه انتگرالگیری مطابق شکل زیر است.

$$R = \{(x, y) \mid 0 \leq y \leq x, 0 \leq x < \infty\} \quad (18-2)$$

$$R = \{(x, y) \mid 0 \leq y < \infty, y \leq x < \infty\} \quad (19-2)$$



طبق ناحیه (۱۸-۲)

$$\alpha = \int_0^{\infty} \int_0^x dG(x)dF(x) = \int_0^{\infty} G(x)dF(x)$$

و یا طبق ناحیه (۱۹-۲)

$$\alpha = \int_0^{\infty} \int_y^{\infty} dF(x)dG(y) = \int_0^{\infty} (1-F(y))dG(y)$$

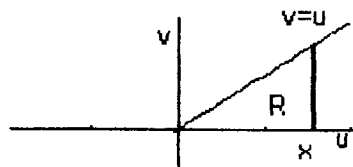
که (۹-۲) و (۱۰-۲) ثابت شده.

فرمول مربوط به v نیز به همین ترتیب با کمی تغییرات بدست می آید. می توان $G_*(y)$ و $F_*(x)$ را مستقیماً با استفاده از تعریف آنها به شرح زیر محاسبه نمود.

$$\begin{aligned} F_*(x) &= P(X \leq x \mid Y < X) = \frac{P(X \leq x, Y < X)}{P(Y < X)} \\ &= \alpha^{-1} P(X \leq x, Y < X) \\ &= \alpha^{-1} \iint_R dG(v)dF(u) \end{aligned}$$

که در آن

$$R = \{(u, v) \mid 0 \leq u \leq x, 0 \leq v \leq u\}$$



$$\begin{aligned} \Rightarrow F_*(x) &= \alpha^{-1} \int_0^x \int_0^u dG(v) dF(u) \\ &= \alpha^{-1} \int_0^x G(u) dF(u) \end{aligned}$$

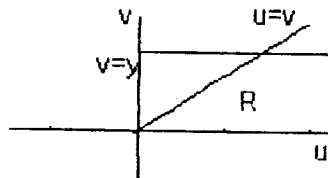
که رابطه (۱۲-۲) را ثابت می کند. رابطه (۱۳-۲) به راحتی از این رابطه نتیجه می شود.

برای رابطه (۱۴-۲) داریم.

$$\begin{aligned} G_*(y) &= P(Y \leq y | Y < X) = \frac{P(Y \leq y, Y < X)}{P(Y < X)} \\ &= \alpha^{-1} \iint_R dF(u) dG(v) \end{aligned}$$

که در آن:

$$R = \{(u, v) \mid 0 \leq v \leq y, v \leq u < \infty\}$$



$$\begin{aligned} \Rightarrow G_*(y) &= \alpha^{-1} \int_0^y \int_v^\infty dF(u) dG(v) \\ &= \alpha^{-1} \int_0^y (1-F(v)) dG(v) \end{aligned}$$

که همان رابطه (۱۴-۲) می باشد و رابطه (۱۵-۲) به سادگی از این رابطه بدست می آید.

اکنون جهت محاسبه $H_*(x, y)$ داریم:

$$H_*(x, y) = P(X \leq x, Y \leq y | Y < X) = \frac{P(X \leq x, Y \leq y, Y < X)}{P(Y < X)} \quad (۲۰-۲)$$

$$P(X \leq x, Y \leq y, Y < X) = \iint_R dG(v) dF(u) \quad (۲۱-۲) \text{ و}$$

با استفاده از ناحیه زیر

$$R = R_1 \cup R_2$$

$$R_1 = \{(u, v) \mid 0 \leq u \leq y, 0 \leq v \leq u\}$$

$$R_2 = \{(u, v) \mid y \leq u \leq x, 0 \leq v \leq y\}$$

داریم:

$$\begin{aligned} P(X \leq x, Y \leq y, Y < X) &= \iint_{R_1} dG(v) dF(u) + \iint_{R_2} dG(v) dF(u) \\ &= \int_0^y \int_0^u dG(v) dF(u) + \int_y^x \int_0^y dG(v) dF(u) \\ &= \int_0^y G(u) dF(u) + \int_y^x G(u) dF(u) \\ &= \int_0^y G(\min(u, y)) dF(u) \quad (22-2) \end{aligned}$$

به سادگی می توان عبارت احتمال (20-2) را به فرم زیر نوشت.

$$\begin{aligned} P(X \leq x, Y \leq y, Y < X) &= \int_0^x \int_0^y I(v \leq u) dF(u) dG(v) \quad (23-2) \\ &= \int_0^y (F(x - F(v))) dG(v) \quad (24-2) \end{aligned}$$

با استفاده از رابطه (20-2) و (23-2) و (24-2) فرمول (16-2) مربوط به $H_*(x, y)$ بدست می آید.

البته توجه کنید که با داشتن H_* می توان F_* ، G_* را با توجه به رابطه زیر بدست آورد.

$$F_*(x) = P(X \leq x \mid Y < X) = \lim_{x \rightarrow \infty} H_*(x, y) = H_*(x, \infty)$$

$$G_*(y) = P(Y \leq y \mid Y < X) = \lim_{x \rightarrow \infty} H_*(x, y) = H_*(\infty, y)$$

و به راحتی روابط (12-1) و (14-2) را نتیجه خواهد داد.

اثبات: راه حل دوم.

می توان با استفاده از قضیه زیر اثبات دقیقتر بر پایه ریاضی ارائه داد.

قضیه 2-2: اگر X, Y بردارهای مستقل با توزیعهای μ, ν در R^j, R^k باشند آنگاه:

$$P[(X, Y) \in B] = \int_{R^j} P[(x, Y) \in B] \mu(dx) \quad ; \quad B \in \mathcal{R}^{j+k} \quad (25-2)$$

$$P[X \in A, (X, Y) \in B] = \int_A P[(x, Y) \in B] \mu(dx) \quad ; \quad A \in \mathcal{R}^j, B \in \mathcal{R}^{j+k} \quad (26-2)$$