



١٠٧٨٦

۸۷/۱۱۰۹۲۷۲

۸۷/۶/۲۸



دانشگاه شهید بهشتی

دانشکده علوم ریاضی

گروه آمار

پایان نامه کارشناسی ارشد آمار بیمه

تقریب پواسون مرکب در مدل مخاطره فردی و کران های تقریب

توسط

مریم تیموریان سفیده خوان

استاد راهنما

دکتر محمد قاسم وحیدی اصل

۱۳۸۷/۱۱/۰ - ۰

استاد مشاور

دکتر محمد ذکایی

شهریور ۱۳۸۷

۱۰۷۸۱۱

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

تقدیم به

پیشگاه لایزال یگانه اش

مادر و پدر بزرگوارم

و شارم برپایشان همه شرم است و شرم

خواهر و برادر مهربانم

و همه کسانی که همچون ستاره‌ای در آسمان زندگی ام درخشیده‌اند.

قدردانی

منت خدای را عزوجل که طاعتش موجب قربت است و به شکر اندرش مزید نعمت.
با تشکر و سپاس فراوان از استاد ارجمند جناب آقای دکتر وحیدی که با حضور
ارزشمند و راهنماییهای مؤثرشان مرا در انجام این پروژه یاری نمودند و من همسو با
علم و دانش، موفق زیستن را در مكتب ایشان آموختم. از استاد محترم جناب آقای
دکتر ذکایی که مشاوره این رساله را بر عهده داشتند صمیمانه تشکر و قدردانی
می نمایم.

از استاد محترم جناب آقای دکتر فرید روحانی و جناب آقای دکتر امیر تیمور پاینده
که رحمت مطالعه و داوری را تقبل فرمودند، کمال تشکر را دارم. از کلیه عزیزان در
پژوهشکده بیمه، به خصوص جناب آقای دکتر حسن زاده و همکاران محترم ایشان، و
مدیریت محترم شرکت بیمه آسیا که در به انجام رساندن این رساله زحمات فراوانی
متحمل شدند، صمیمانه تشکر و قدردانی می کنم.

مریم تیموریان سفیده خوان

تهران — شهریور ماه ۱۳۸۷

پیشگفتار

داشتمان خاصی را در یک شرکت بیمه در نظر بگیرید. این شرکت برای برآورد میزان ادعای کل داشتمان و محاسبه احتمال ورشکستگی و حق بیمه مربوط به هر قرارداد، به تابع توزیع مجموع ادعای کل داشتمان نیاز دارد. بنابراین یافتن تابع توزیع مجموع متغیرهای تصادفی در صنعت بیمه امری بسیار با اهمیت است. فرض کنید X_i متغیر تصادفی حقیقی مقدار ونشانگر مقدار ادعای قرارداد n ام باشد. بنابراین ما در پی یافتن توزیع متغیر تصادفی X_i $S_n = \sum_{i=1}^n X_i$ می‌باشیم. حتی در حالتی که X_i ها متغیرهای تصادفی برونولی باشند، یافتن تابع توزیع S_n برای داشتمان‌های بزرگ کاربسیار سختی است. لذا در اکثر موارد این توزیع با توزیع مناسبی تقریب زده می‌شود. توزیع‌های مناسب در طول این رساله، توزیع‌های پواسون و پواسون مرکب است. در بخش اول و دوم فصل دوم، دو مدل ممکن برای یک داشتمان توضیح داده شده است. در بخش سوم فصل دو روش استاین و کرستن برای یافتن کران بالای تقریب توزیع مجموع ادعای کل در شرایط خاص بیان شده است. در فصل سوم کران‌های خطأ با مرتبه بهتری برای تقریب مجموع ادعای کل با توزیع پواسون مرکب ارائه شده است. این کران‌ها بدون هیچ محدودیتی بر توزیع ادعاهای، برقرار می‌باشند و به کران‌های با عامل جادویی معروفند. فصل چهارم اختصاص به یک مثال کاربردی دارد که در آن کران‌های خطأ برای یک داشتمان خاص محاسبه شده است.

چکیده

تابع توزیع مجموع متغیرهای تصادفی نقشی مهم در مدیریت داشتمان‌ها و در نتیجه آمار بیمه دارد. به طور مثال اگر مدل یک داشتمان خاص مدل مخاطرهٔ فردی و هدف، یافتن مقدار ادعای کل این داشتمان در یک دورهٔ زمانی معین باشد، آنگاه یافتن تابع توزیع مجموع متغیرهای تصادفی امری مهم است. اما از طرفی یافتن تابع توزیع مجموع متغیرهای تصادفی در حالت کلی کاری بسیار مشکل و اغلب تقریباً ناممکن است. بنابراین تابع توزیع مجموع متغیرهای تصادفی را تقریب می‌زنیم. از جملهٔ توزیع‌های مناسب برای تقریب‌زدن، توزیع پواسون مرکب است. بدیهی است که هر تقریبی خطابی نیز دارد. بنابراین برای بهبود تقریب توزیع مجموع متغیرهای تصادفی، یافتن کران بالای تقریب اهمیت پیدا می‌کند. تاکنون روشهای متفاوتی برای یافتن این کران خطا ارائه شده است. در این رساله کران‌های بالای متفاوت و سه روش مهمی که از طریق آنها کران‌های بالایی با مرتبهٔ بهتر به دست می‌آیند، مورد بحث و بررسی قرار گرفته‌اند. توجه ما بیشتر به کران‌های بالایی است که هیچ محدودیت و شرطی برای توزیع متغیرهای تصادفی ندارند و از مرتبهٔ بهتری نسبت به کران‌های آرائه شده پیشین باشد. در انتها نیز کران‌های مزبور برای یک داشتمان واقعی محاسبه و مقایسه شده‌اند.

کلمات کلیدی: تقریب پواسون مرکب، مدل مخاطرهٔ فردی، عامل جادویی، فاصلهٔ تغییرات کل.

فهرست مندرجات

۱	تعاریف و مفاهیم اولیه	۱
۱	۱ مقدمه	۱.۱
۱	۲ نمادهای ۰ و ۱	۲.۱
۳	۳ قضیه رادون-نیکودیم	۳.۱
۴	۴ فاصله تغییرات کل	۴.۱
۱۳	۲ تعیین توزیع مجموع مبالغ ادعاهای خسارت	۲
۱۳	۱۳ مقدمه	۱.۲
۱۳	۲۰۲ مدل مخاطرهٔ فردی	۲.۲
۱۵	۳۰۲ مدل مخاطرهٔ جمعی	۳.۲
۱۸	۴۰۲ تقریب مدل مخاطرهٔ فردی	۴.۲
	۱۰۴.۲ روش استاین برای کران‌های بالای تقریب مجموع	
	متغیرهای تصادفی برنولی با توزیع پواسون	
	۲۱	۲۱

۲۰.۲ روش کرستن برای کرانهای بالای تقریب مجموع بردارهای تصادفی برنولی با توزیع پواسون چند متغیره	۳۲
۳۰.۲ روش کرستن برای کرانهای بالای تقریب مجموع متغیرهای تصادفی با توزیع پواسون مرکب	۴۰
۳ کرانهای بالا برای تقریب با توزیع پواسون مرکب در حالت پیوسته و k-بعدی	
۴۸	
۱.۳ مقدمه	۴۸
۲۰.۳ کرانهای بالا برای تقریب مجموع متغیرهای تصادفی با توزیع پواسون مرکب	۴۹
۳۰.۳ مثال	۶۱
۴۰.۳ نتیجه گیری	۶۲
۴ کران بالا برای داشتمانی شامل بیمه‌های اشیاء	
۱.۴ مقدمه	۶۳
۲۰.۴ تعاریف	۶۳
۱.۰.۴ بیمه اتومبیل	۶۳
۲۰.۴ بیمه باربری	۶۵
۳۰.۴ بیمه آتش‌سوزی	۶۵
۳۰.۴ توزیع خسارت انواع بیمه‌نامه‌ها	۶۶
۱.۰.۴ بیمه بدنی اتومبیل	۶۶
۲۰.۴ بیمه شخص ثالث	۶۸

۷۰	بیمه آتش سوزی	۳.۳.۴
۷۲	بیمه باربری	۴.۳.۴
۷۴	یافتن کران بالا برای تقریب توزیع ادعای کل	۴.۴
۷۷	برنامه کامپیوتری به زبان Matlab	A
۸۱	الگوریتم بازگشتیتابع $g_{\lambda,A}$	B
۸۴	واژه‌نامه	C
۸۵	نامنامه	D
۸۷	مراجع	E

فصل ۱

تعریف و مفاهیم اولیه

۱.۱ مقدمه

این فصل شامل تعریف‌ها و مفاهیم اولیه مورد نیاز می‌باشد. در بخش اول، تعریف‌ها و قضیه‌های مهم در مورد نمادهای لانداو تشریح شده است. این نمادها به خصوص در ساده‌نویسی گزاره‌هایی که شامل مفهوم حد هستند، به کار می‌روند. در بخش دوم این فصل قضیه رادون-نیکودیم بیان شده است که در فصل‌های بعدی از آن استفاده خواهد شد. در بخش سوم نیز توضیح مختصری درباره فاصله تغییرات کل که به عنوان معیاری برای اندازه‌گیری خطای تقریب استفاده می‌شود و نیز لم‌هایی که برای اثبات قضیه‌ها در فصل‌های بعد به کار خواهد رفت، آرائه شده است.

۲.۱ نمادهای O و o

نمادهای O و o که به افتخار لانداو، ریاضیدان آلمانی، به نمادهای لانداو مشهور‌اند، در سال ۱۸۹۲ توسط پل گوستاو هانریش باخمن مطرح گردید.

تعریف ۱.۲.۱ فرض کنید $\mathcal{R} \rightarrow \mathcal{N} : f, g \mapsto \mathcal{N}$. در این صورت f از مرتبهٔ حداقل g است هرگاه ثابت مثبتی مانند C و عدد صحیح مثبتی مانند n موجود باشند به‌طوری که برای هر $n \geq n_0$ ،

$$|f(n)| \leq C|g(n)|. \quad (1.2.1)$$

این مفهوم با نماد $f = O(g)$ نشان داده می‌شود.

عدد C یکتا نیست؛ زیرا رابطه $(1.20.1)$ برای هر عدد صحیح بزرگتر از C نیز برقرار است. تعریف کلی تری نیز برای O به صورت زیر آنند شده است.

تعریف $2.20.1$ فرض کنید g و f دو تابع حقیقی مقدار باشند که در همسایگی نقطه x_0 تعریف شده‌اند. گوییم f از مرتبه O تابع g است، وقتی $x \rightarrow x_0$ و $f = O(g)$ و می‌نویسیم $f = O(g)$ هرگاه یک همسایگی V از x_0 و ثابت مثبتی M وجود داشته باشد به‌طوری که برای هر $x \in V$ که $x \neq x_0$

$$|f(x)| \leq M|g(x)| \quad (2.20.1)$$

اگر g در نزدیکی x_0 صفر نشود، آنگاه رابطه بالا معادل است با شرط

$$\limsup_{x \rightarrow x_0} \left| \frac{f(x)}{g(x)} \right| < \infty.$$

نماد o کوچک مفهوم مهم دیگری است که به صورت زیر تعریف می‌شود.

تعریف $3.20.1$ فرض کنید g و f دو تابع حقیقی مقدار باشند که در همسایگی نقطه x_0 تعریف شده‌اند. $(f = o(g))$ وقتی $x \rightarrow x_0$ هرگاه برای هر $\epsilon > 0$ همسایگی $(\epsilon) = V(\epsilon)$ از x_0 وجود داشته باشد به‌طوری که برای هر $x \in V$ که $x \neq x_0$ $|f(x)| \leq \epsilon |g(x)|$. اگر g در همسایگی x_0 صفر نشود، در این صورت شرط $f = o(g)$ معادل است با

$$\lim_{x \rightarrow x_0} \frac{f(x)}{g(x)} = 0.$$

بدیهی است که برای $x \rightarrow x_0$ ، اگر $f = O(g)$ آنگاه $f = o(g)$ اماعکس این رابطه برقرار نیست.

برخی از ویژگی‌های مفید نمادهای o و O در قضیه زیر و بدون اثبات آمده است.

قضیه $1.20.1$ ویژگی‌های زیر برای نمادهای مرتبه برقرار است.

• اگر C_1 و C_2 مقادیر ثابتی باشند، آنگاه $O(g) + O(g) = O(g)$

• اگر C_1 و C_2 مقادیر ثابتی باشند، آنگاه $O(g) + o(g) = o(g)$

$O(O(g)) = O(g)$ •

$O(o(g)) = o(O(g)) = o(o(g)) = o(g)$ •

$O(f)O(g) = O(fg)$ •

$O(f)o(g) = o(fg)$ •

تابع g در رابطه‌های $f = o(g)$ یا $f = O(g)$ یا $f = \frac{f(x)}{g(x)}$ و قتی $x \rightarrow x$ ، دارای کران متناهی است و رابطه $f = o(g)$ نشانگر این است که وقتی $x \rightarrow x$ ، $\frac{f(x)}{g(x)} \rightarrow 0$ ؛ اما هیچ اطلاعی از سرعت همگرایی آن به دست نمی‌دهد. برای مثال $\sin x = O(\lambda x)$ و $\sin x = O(x)$ وقتی $x \rightarrow \infty$ هر دو درست هستند. همچنین گزاره‌های $o(1) = o\left(\frac{1}{x}\right)$ و $o(1) = o\left(\frac{1}{1+x^2}\right)$ وقتی $x \rightarrow \infty$ نیز درست‌اند. بنابراین اطلاعی که نمادهای مرتبه به ما می‌دهند، کاملاً دقیق نبوده و از طرفی همین ویژگی مبهم بودن این امکان را به ما می‌دهد که آن‌ها را در شرایط مختلفی به کار ببریم. در پایان ذکر این نکته ضروری است که از نماد O به طور وسیعی در مطالعه رشد توابع استفاده می‌شود. به این معنی که $O(g) = f$ فقط دلالت بر این دارد که تابع f رشدی سریعتر از g ندارد. در واقع این نماد، یک کران بالا برای f به ازای مقادیر بزرگ x ارائه می‌دهد.

۳.۱ قضیه رادون-نیکودیم

در فصل سوم از قضیه رادون-نیکودیم استفاده خواهیم کرد. لذا در این قسمت آن را بدون اثبات بیان می‌کنیم. برای اثبات این قضیه می‌توان به کتاب آش (۲۰۰۰) مراجعه کرد. قبل از بیان قضیه ابتدا تعریف استفاده شده در قضیه را ذکر می‌کنیم. فرض می‌کنیم (Ω, \mathcal{F}) فضای اندازه‌پذیر بوده و μ اندازه‌ای سیگما-متناهی و λ اندازه‌ای علامتدار بر \mathcal{F} باشند.

تعریف ۱.۳.۱ اندازه λ را نسبت به اندازه μ مطلقاً پیوسته گوییم اگر به ازای هر $A \in \mathcal{F}$ که $\mu(A) = 0$ ، آنگاه $\lambda(A) = 0$. این مفهوم با نماد $\ll \lambda$ نمایش داده می‌شود.

۶-

قضیه ۱.۳.۱ فرض کنید اندازه μ نسبت به λ مطلقاً پیوسته باشد. بنابراین تابع اندازه‌پذیر بورل مانند $g: \Omega \rightarrow \mathbb{R}$ ، وجود دارد به طوری که برای هر مجموعه $A \in \mathcal{F}$

$$\lambda(A) = \int_A g d\mu.$$

اگر h تابع اندازه‌پذیر دیگری با همین خاصیت باشد، آنگاه تقریباً همه جا، نسبت به اندازه μ ، $g = h$.

تابع g تعریف شده در قضیه ۱.۳.۱ را مشتق رادون-نیکودیم یا چگالی λ نسبت به μ نامیده و با $\frac{d\lambda}{d\mu}$ نشان می‌دهیم. اگر μ اندازه لبگ باشد، آنگاه تابع g را چگالی λ می‌نامیم.

۴.۱ فاصله تغییرات کل

برای بررسی همگرایی در توزیع، از فاصله بین توزیع‌ها استفاده می‌شود. فاصله تغییرات کل برای دو متغیر تصادفی X و Y به صورت زیر تعریف می‌شود.

فرض کنید متغیرهای تصادفی X و Y بر فضای احتمال (Ω, \mathcal{F}, P) تعریف شده باشند و \mathcal{B} ، سیگما میدان بورل (سیگما میدان روی خط حقیقی) باشد. در این صورت فاصله تغییرات کل بین توزیع‌های این دو متغیر تصادفی برابر است با

$$d_{TV}(X, Y) = \sup_{A \in \mathcal{B}} |P(X \in A) - P(Y \in A)| \quad (۳.۴.۱)$$

اگر متغیرهای تصادفی X و Y به ترتیب دارای توابع توزیع F_X و F_Y باشند، آنگاه فاصله تغییرات کل بین توزیع‌های این دو متغیر تصادفی اغلب به صورت $d_{TV}(F_X, F_Y)$ نیز نشان داده می‌شود.

برای فاصله تغییرات کل تعریف‌های متفاوتی بیان شده است. این تعریف‌ها معادل باهم بوده و هر یک در شرایط مناسب استفاده می‌شود. در ادامه به یکی از این تعریف‌ها که در فصل سوم استفاده خواهد شد، اشاره می‌کنیم.

فرض کنید P_1 و P_2 دو اندازه احتمال بر فضای اندازه‌پذیر (Ω, \mathcal{F}) و نسبت به اندازه μ مطلقاً پیوسته باشند. بنابراین طبق قضیه رادون-نیکودیم به ترتیب دارای توابع چگالی f_1 و f_2

نسبت به اندازه μ خواهند بود. فاصله تغییرات کل میان این دو اندازه برابر است با

$$d_{TV}(P_1, P_2) = \frac{1}{\mu} \int |f_1 - f_2| d\mu.$$

در ادامه ثابت می‌کنیم d_{TV} یک متر است.

لم ۱.۴.۱ اگر d_{TV} فاصله تغییرات کل و Z, Y, X متغیرهای تصادفی باشند، آنگاه

- i) $d_{TV} \geq 0$.
- ii) $d_{TV}(X, Y) = d_{TV}(Y, X)$.
- iii) $d_{TV}(X, Y) \leq d_{TV}(X, Z) + d_{TV}(Z, Y)$.
- iv) $d_{TV}(X, Y) = 0 \iff P_X(A) = P_Y(A) \quad \forall A \in \mathcal{B}$.

برهان : با توجه به تعریف فاصله تغییرات کل و خواص قدرمطلق، روابط (i) و (ii) بدیهی‌اند. برای اثبات (iii) با استفاده از نابرابری مثلثی برای متغیرهای تصادفی Z, Y, X داریم :

$$\begin{aligned} |P(X \in A) - P(Y \in A)| &\leq |P(X \in A) - P(Z \in A)| \\ &+ |P(Z \in A) - P(Y \in A)| \end{aligned}$$

حال اگر از طرفین نابرابری بالا نسبت به $A \in \mathcal{B}$ سوپرتم بگیریم، خواهیم داشت:

$$d_{TV}(X, Y) \leq d_{TV}(X, Z) + d_{TV}(Z, Y).$$

برای اثبات حکم (iv)،

$$d_{TV}(X, Y) = 0 \iff \sup_{A \in \mathcal{B}} |P(X \in A) - P(Y \in A)| = 0$$

اگر و تنها اگر به ازای هر $A \in \mathcal{B}$ داشته باشیم $P_X(A) = P_Y(A)$.

لم ۲.۴.۱ اگر d_{TV} فاصلهٔ تغییرات کل و X و Y متغیرهای تصادفی گسسته باشند، آنگاه

$$\begin{aligned} i) d_{TV}(X, Y) &\leq P(X \neq Y). \\ ii) d_{TV}(X, Y) &= \sum_{k=0}^{\infty} |P(X = k) - P(Y = k)|. \end{aligned}$$

برهان : برای اثبات اولین نابرابری مجموعه دلخواه $A \in \mathcal{B}$ ، را در نظر می‌گیریم.

$$\begin{aligned} P(X \neq Y) &\geq P(X \in A, Y \in A^c) \\ &= P(X \in A) - P(X \in A, Y \in A) \geq P(X \in A) - P(Y \in A) \end{aligned}$$

همچنین با تعویض نقش‌های X و Y

$$P(X \neq Y) \geq P(Y \in A) - P(X \in A)$$

بنابراین

$$P(X \neq Y) \geq |P(X \in A) - P(Y \in A)|.$$

حال اگر از طرفین رابطهٔ بالا نسبت به $A \in \mathcal{B}$ سوپررم بگیریم، حکم حاصل می‌شود.
برای اثبات تساوی دوم فرض کنید که $\{0, 1, 2, \dots\} \subseteq A \subseteq \{0, 1, 2, \dots\}$. در این صورت

$$\begin{aligned} \sum_{k=0}^{\infty} |P(X = k) - P(Y = k)| &= \\ &= \sum_{k \in A} |P(X = k) - P(Y = k)| + \sum_{k \in A^c} |P(X = k) - P(Y = k)| \\ &\geq |\sum_{k \in A} P(X = k) - P(Y = k)| + |\sum_{k \in A^c} P(X = k) - P(Y = k)| \\ &= |P(X \in A) - P(Y \in A)| + |P(X \in A^c) - P(Y \in A^c)| \\ &= 2|P(X \in A) - P(Y \in A)|. \end{aligned}$$

بنابراین

$$\sup_A |P_X(A) - P_Y(A)| \leq \sum_{k=0}^{\infty} |P(X = k) - P(Y = k)|. \quad (\text{۴.۴.۱})$$

حال فرض کنید $A = \{k : P(X = k) \geq P(Y = k)\}$. در این صورت

$$\begin{aligned} 2|P(X \in A) - P(Y \in A)| &= 2 \left| \sum_{k \in A} P(X = k) - P(Y = k) \right| \\ &= 2 \sum_{k \in A} [P(X = k) - P(Y = k)]. \end{aligned}$$

از طرفی

$$\begin{aligned} \sum_{k=0}^{\infty} |P(X = k) - P(Y = k)| &= \sum_{k \in A} P(X = k) - P(Y = k) + \sum_{k \in A^c} P(Y = k) - P(X = k) \\ &= 2 \sum_{k \in A} [P(X = k) - P(Y = k)]. \end{aligned}$$

بنابراین

$$2|P(X \in A) - P(Y \in A)| = \sum_{k=0}^{\infty} |P(X = k) - P(Y = k)|$$

و

$$\sup_A |P_X(A) - P_Y(A)| \geq \frac{1}{2} \sum_{k=0}^{\infty} |P(X = k) - P(Y = k)|. \quad (5.4.1)$$

از روابط (۴.۴۱) و (۵.۴۱) اثبات کامل می‌شود.

لم زیر برای یافتن کران بالای فاصله تغییرات کل مورد استفاده قرار خواهد گرفت.

لم ۳.۴.۱ برای دنباله‌های $\{X_i\}_{i=1}^n$ ، $\{Y_i\}_{i=1}^n$ ، از متغیرهای تصادفی دلخواه داریم:

$$P\left(\sum_{i=1}^n X_i \neq \sum_{i=1}^n Y_i\right) \leq \sum_{i=1}^n P(X_i \neq Y_i).$$

برهان: برای هر $i = 1, 2, \dots, n$ ، تعریف می‌کنیم:

$$A_i = \{X_i \neq Y_i\}.$$

با استفاده از نابرابری بول،

$$\begin{aligned} P\left(\sum_{i=1}^n X_i \neq \sum_{i=1}^n Y_i\right) &\leq P\left(\bigcup_{i=1}^n A_i\right) \\ &\leq \sum_{i=1}^n P(A_i) = \sum_{i=1}^n P(X_i \neq Y_i). \end{aligned}$$

لذا حکم ثابت می‌شود.

روش‌های بسیاری برای یافتن کوچکترین کران‌های بالا برای متر $d_{TV}(X, Y)$ موجود است. یکی از این روش‌ها، روش جفت‌سازی است. این روش را وانگ (۱۹۸۶) ارائه داد. فرض کنید X و Y متغیرهای تصادفی گستته بر فضای احتمال (Ω, \mathcal{F}, P) باشند. با توجه به لم (۲.۴.۱) برای یافتن کران بالای $d_{TV}(X, Y)$ می‌توان توزیع توانی برای متغیرهای تصادفی X و Y یافت به‌طوری که $P(X = Y) = 0$ ماسیم شود.

می‌دانیم که به ازای هر دو توزیع F و G توزیع‌های توانی بسیاری مانند H وجود دارند که دارای توزیع‌های حاشیه‌ای F و G باشند. به‌ویژه دو توزیع توانی H^* و H را می‌توان از توزیع‌های F و G تولید کرد به‌طوری که بهترین بیشترین و کمترین همبستگی باشند. ویت (۱۹۷۶) ثابت می‌کند که اگر X دارای توزیع بینولی و Y دارای توزیع پواسون با $E(X) = E(Y)$ باشد، آنگاه روش استفاده شده در استخراج توزیع توانی H^* که دارای بیشترین همبستگی است، موجب ماسیم شدن $P(X = Y)$ نیز خواهد شد. اما در حالت کلی توزیع توانی H^* مقدار $P(X = Y)$ را ماسیم نخواهد کرد. می‌خواهیم بدانیم که آیا روشی وجود دارد که بتوان توزیع توانی (جفت‌سازی) یافت به‌طوری که موجب برقراری تساوی $d_{TV}(X, Y) = P(X \neq Y)$ برای هر دو متغیر تصادفی گستته دلخواه X و Y شود؟ پاسخ مثبت است. برای این منظور ابتدا تعریف زیر را بیان می‌کنیم.

تعریف ۱.۴.۱ جفت‌ساز دو متغیر تصادفی گستته دلخواه X و Y را بیشین می‌نامند، هرگاه

$$d_{TV}(X, Y) = P(X \neq Y).$$

چنین جفت‌ساز بیشین را با $P^*(X \neq Y)$ نشان می‌دهیم.

برای یافتن توزیع توانی که در تعریف (۱.۴.۱) صدق کند، فرض خواهیم کرد که X و Y متغیرهای تصادفی گستته و حقیقی مقدار باشند که بر فضای احتمال (Ω, \mathcal{F}, P) تعریف شده‌اند. مجموعه B را به صورت زیر تعریف می‌کنیم:

$$B = \{x : 0 \leq P(X = x) < P(Y = x)\}.$$

برای هر $x \in B^c$ متمم مجموعه B است) تعریف می‌کنیم:

$$a_x = P(X = x) - P(Y = x)$$

$$\alpha_x = \frac{a_x}{\sum_{x \in B^c} a_x}$$

و برای هر $y \in B$

$$\begin{aligned} b_y &= P(Y = y) - P(X = y) \\ \beta_y &= \frac{b_y}{\sum_{y \in B} b_y} \end{aligned}$$

به راحتی ثابت می‌شود که برای هر x, y

$$\sum_{x \in B^c} a_x = \sum_{y \in B} b_y \quad , \quad a_x \beta_y = a_x b_y.$$

چگالی توان متفاوت X و Y را به صورت زیر تعریف می‌کیم:
برای $x \in B$

$$f(x, y) = \begin{cases} P(X = x) & y = x \\ 0 & \text{در غیر این صورت} \end{cases} \quad \text{و برای } y \in B^c$$

$$f(x, y) = \begin{cases} P(Y = y) & x = y \\ 0 & \text{در غیر این صورت} \end{cases} \quad \text{و برای } y \in B \text{ و } x \in B^c$$

$$f(x, y) = a_x \beta_y = a_x b_y. \quad (6.4.1)$$

تابع f تمام شرایط تابع چگالی را دارد است؛ زیرا اولاً مقادیر تابع f نامنفی است. ثانیاً

$$\begin{aligned} \sum f(x, y) &= \sum_{x \in B} P(X = x) + \sum_{y \in B^c} P(Y = y) + \sum_{x \in B^c} \sum_{y \in B} a_x b_y \\ &= \sum_{x \in B} P(X = x) + 1 - \sum_{y \in B} P(Y = y) + \sum_{x \in B^c} \sum_{y \in B} \frac{a_x}{\sum_{x \in B^c} a_x} b_y \\ &= \sum_{x \in B} P(X = x) + 1 - \sum_{y \in B} P(Y = y) + \sum_{y \in B} b_y = 1. \end{aligned}$$

قضیه ۱.۴.۱ اگر متغیرهای X و Y دارای تابع چگالی توان f تعریف شده به صورت
(۶.۴.۱) باشند، آنگاه

$$d_{TV}(X, Y) = P(X \neq Y).$$

برهان : با استفاده از لم ۲.۴.۱ ،

$$\begin{aligned} d_{TV} &= \frac{1}{2} \left| \sum_{x=0}^{\infty} P(Y=x) - P(X=x) \right| \\ &= \frac{1}{2} \left[\sum_{x \in B} \{P(Y=x) - P(X=x)\} + \sum_{x \in B^c} \{P(X=x) - P(Y=x)\} \right] \\ &= \frac{1}{2} \left[\sum_{x \in B} b_x + \sum_{x \in B^c} a_x \right] \end{aligned}$$

اما با توجه به اینکه $\sum_{x \in B^c} a_x = \sum_{x \in B} b_x$ ، بنابراین

$$d_{TV} = \sum_{x \in B^c} a_x. \quad (۷.۴.۱)$$

از طرفی با استفاده از قانون احتمال کل و توزیع توانم f ،

$$\begin{aligned} P(X \neq Y) &= \sum_{x=0}^{\infty} P(X=x, Y \neq x) \\ &= \sum_{x \in B} P(X=x, Y \neq x) + \sum_{x \in B^c} P(X=x, Y \neq x) \\ &= \sum_{x \in B} 0 + \sum_{x \in B^c} \left[\sum_{y \in B} a_x \beta_y + \sum_{y \in B^c, y \neq x} 0 \right] \\ &= \sum_{x \in B^c} \sum_{y \in B} [P(X=x) - P(Y=x)] \left[\frac{P(Y=y) - P(X=y)}{\sum_{y \in B} [P(Y=y) - P(X=y)]} \right] \\ &= \sum_{x \in B^c} P(X=x) - P(Y=x) = \sum_{x \in B^c} a_x. \end{aligned}$$

بنابراین $d_{TV}(X, Y) = P(X \neq Y)$

همان طور که قبل اشاره شد، روش جفت‌سازی برای اثبات قضیه زیر به کار می‌رود.

قضیه ۲.۴.۱ هرگاه $\{Y_1, Y_2, \dots, Y_n\}$ و $\{X_1, X_2, \dots, X_n\}$ دنباله‌هایی از متغیرهای تصادفی گستته و مستقل باشند، در این صورت

$$d_{TV}\left(\sum_{i=1}^n X_i, \sum_{i=1}^n Y_i\right) \leq \sum_{i=1}^n d_{TV}(X_i, Y_i).$$