

تقدير الامكان الأعظم شبه المعلمي لنموذج نسبة كثافة الاحتمال لعينتين بالتطبيق على بيانات مرضي الالتهاب الكبدي الوبائي

أ.د. فاطمة علي عبدالعاطى

أستاذ الإحصاء التطبيقى

كلية التجارة - جامعة المنصورة

أ.د. محمد توفيق البلقينى

أستاذ الإحصاء الإكتوارى والتأمين

كلية التجارة - جامعة المنصورة

لمياء على عبدالخالق البرعى
مدرس مساعد - قسم الإحصاء والاساليب الكمية
كلية الادارة - جامعة الدلتا للعلوم والتكنولوجيا

تقدير الامكان الأعظم شبه المعلمي لنموذج نسبة الكثافة الاحتمالي لعينتين بالتطبيق على بيانات مرضي الالتهاب الكبدي الوبائي

أ.د. محمد توفيق البليغى
أستاذ الإحصاء الإكتوارى والتأمين
كلية التجارة، جامعة المنصورة

لمياء علي عبدالخالق البرعى
مدرس مساعد - بقسم الإحصاء والاساليب الكمية
كلية الادارة - جامعة العطا للعلوم والتكنولوجيا

المستخلص

تعد الطرق الحالية لتحليل البيانات المراقبة (censoring data) لنموذج نسبة الكثافة شبه المعلمي (DRM) قليلة، وذلك بسبب صعوبة الحصول على مقدر الامكان الكامل من خلال البيانات المراقبة بسبب صعوبة الإثباتات النظرية للخصائص المتطرفة، خاصة اذا كانت دالة الامكان تتطلب عدد لا نهائي لمعلمات الابعاد.

يهدف هذا البحث إلى تقديم طرقة الامكان الأعظم شبه المعلمية (SML) لتقدير معلمات نموذج نسبة الكثافة لعينتين متضمنة بيانات مراقبة من اليمين (Right Censoring Data)، بالإضافة على بيانات تتبع التوزيع الأسوي والتوزيع الطبيعي اللوغاريتمي ومن خلال هذه الدراسة سوف يتم استخدام خوارزمية E-M لتقدير دالة الامكان الأعظم - لتيسير ايجاد المقدار المقترن - و يتم اجراء التقدير لكل من المكونات المعلمية واللامعلمية للنموذج.

سوف يتم تطبيق هذا النموذج على مجموعة من البيانات المراقبة والتي تتضمن بيانات عن مرضي الالتهاب الكبدي الوبائي (فيروس C) والتي تم الحصول عليها من احدى المراكز الطبية المتخصصه. وثبتت نتائج هذه الدراسة أن النموذج المقترن يتمتع بتقديرات ذات جودة عالية وأن النموذج ملائم لطبيعة البيانات.

الكلمات المفتاحية :

نموذج نسبة الكثافة، خوارزمية M-E، تدبر الامكان الاعظم شبه المعلمي، الامكان التجربى، البيانات المراقبة لليمين.

Abstract:

The methods currently exist for analyzing the conventional right-censored data under the semi-parametric density ratio model are considered very limited. This is mainly due to the difficulty of obtaining the full-likelihood estimator with right-censored data, the considerably complex computational algorithms and the sophisticated theoretical properties that need to be proven, especially when the likelihood function involves infinite dimensional parameters

This research aims to find semi-parametric maximum likelihood (SML) method by using the density ratio model (DRM) for two samples which include Censoring Data. Through this study the E-M algorithm will be used to estimate the maximum likelihood estimator and to facilitate finding the proposed estimator, and the estimation will be made for both parametric and nonparametric components of this model.

This model will be applied on a group of Censoring Data that include data of Hepatitis C patients obtained from a specialized medical center. the results of this study proved that the proposed model provides better fit regarding the treatment of this type of data.

Key Words: Density ratio model; EM algorithm; empirical likelihood; right-censored data; semi-parametric maximum likelihood estimation

١. المقدمة

ظهرت في العقود الماضية الكثير من الكتابات التي تركز على دراسة واحد من أهم النماذج شبه المعلمية وهو نموذج نسبة كثافة الاحتمال، فهو نموذج بحظي بالاهتمام لأنه يحتوي على العديد من المميزات؛ مثل: أنه حل بديل لمشاكل المقارنة بين عينتين أو أكثر، حيث أنه يربط توزيعين من خلال هيكل شبه المعلمي من حل الافتراضات التقليدية لبعض مشاكل تعدد العينات، ويأخذ العديد من امور التوزيعات التقليدية والنماذج كحالة خاصة.

حيث أشار (Qin & Zhang 1997) الى التكافؤ بين نموذج نسبة الكثافة ونموذج الانحدار اللوجستي لبيانات مراقبة، واقتصر (Qin 1998) نموذج نسبة الكثافة يتضمن نموذج العينات المتحجزة وأظهر ايضا العلاقة بين نموذج نسبة الكثافة ونموذج المخاطر النسبية (Cox 1972). درس العديد من الباحثين نموذج نسبة الكثافة من جوانب عديدة مثل (Zhang 2002)& Qin(2003) & Chen& Lin (2013) & Wang& Zhang (2014)

ومن هنا جاء الاهتمام بنموذج نسبة الكثافة وأهميته في ايجاد مقدراته وتطبيقاته علي البيانات المراقبة . وفي إطار ذلك تتكون هذه الدراسة من خمسة مباحث هي : المبحث الثاني يعرض الاستعراض المرجعي لنموذج نسبة الكثافة وتطبيقاتها المختلفة والمبحث الثالث تعريف نموذج نسبة الكثافة ومواصفاته، يتضمن المبحث الرابع تقدير معلمات النموذج باستخدام طريقة الامكان الاعظم. كما تم تقدير معلمات النموذج باستخدام البيانات المراقبة للبيان في المبحث الخامس، أيضا تم تطبيق التقدير المقترن لنموذج نسبة الكثافة على بيانات حقيقة لاثبات كفاءته في المبحث السادس. ولخيرا تم عرض ماتم التوصل اليه من نتائج في هذا الدراسة في المبحث السابع.

٢. الدراسات السابقة:

بحث كلأ من (Zhang, Chen (2021) نهجاً يعتمد على تحليل المكون الاساسي الدالي (FPCA). حيث يحدد FPCA التباين والتشابه بين الدول ذات الأهمية من خلال التوسيع "الأمثل" لهذه الدول. و الحصول على تقديرات تقريرية دقيقة من خلال مجموعات خطية من عدد قليل من الدول الذاتية لنموذج DRM، وتحديد أهم الدول الذاتية لنسبة كثافة اللوغاريتم $\{ \log g_k(x)/g_0(x) : k = 0, \dots, m \}$ من خلال عينات متعددة، وأظهرت نتائج المحاكاة التي أجريت أنه عندما يتم تعميم البيانات من التوزيعات شائعة

الاستخدام، فإن تحويل DRM على أساس دالة الأساس الكافية له كفاءة كبيرة، في حين قدم كلا من Zhang, Zhu, Chen (2020) اختبار نسبة الامكان التجريبي للكميات لمودج نسبة الكثافة حيث قاما بتحديد توزيع كاي تريبيع لـ ELRT الاختبار نسبة الامكان لمودج $g_k = \exp\{\theta_k^T q(x)\} g_0(x)$ DRM ودرسوا طريقة الاستدلال للكميات عند توفر عينات متعددة من مجموعات مرتبطة. وتم اثبات ان احصائية ELRT لها توزيع كاي تريبيع تحت فرض العدم في ظل ظروف معينة، ويساعد DRM على تحسين الكفاءة الإحصائية، وأوضحت البيانات الحقيقية كفاءة الطريقة المقترنة، واقتصر كلا من Zhuang, HU, Chen(2019) على الربط بين المجتمعات المستقلة من خلال نموذج نسبة الكثافة حيث يقوم هذا النموذج بتطوير مقدر الامكان التجريبي وإثبات طبيعته المتقاربة بالإضافة إلى تحسين كفاءة التقدير من خلال دراسة أوجه الشابهة بين المجتمعات وتقديم طريقة صالحة لاختبار الفرضيات وبناء فترات ثقة بناءً على النموذج: $d F_k(x) = \exp\{\theta_k^T q(x)\} d F_0(x)$, $k = 1, \dots, m$ يحسن بشكل كبير من كفاءة التقدير وقوة الاختبار، ويؤدي إلى فترات ثقة دقيقة بشكل كافٍ. كما أوضحت البيانات الحقيقية كفاءة الطريقة المقترنة، وتوصلت دراسة كلا من Zeng, Gao, and Lin (2017) إلى مقدر الامكان الاعظم لنماذج الانحدار شبه المعلمي لبيانات مراقبة متعددة الفترات، وأظهرت الدراسة بيانات اوقات الفشل متعدد الفترات الخاصة للرقابة عندما تكون هناك أنواع متعددة من الفشل وكل وقت فشل يمكن في فترة معينة. حيث اعتمدت الدراسة على الرقت الذي يؤثر في المتغيرات بينما على اوقات الفشل متعددة المتغيرات ، من خلال النظر في قيمه واسعة من نماذج التحول شبه المعلمي مع تأثيرات عشوائية، بالإضافة على دالة الخطير التراكمي. $A_{ijk}(t) = G_k \left[\int_0^1 \exp\{\beta^T x_{ijk}^{(s)} + b_i^T z_{ijk}(s)\} d A_k(s) \right]$ وانتجت الدراسة ان المقدرات المقترنة للمعلمات ذات البعد المحدود متقاربة ومنتظمة بشكل طبيعي مع مصفوفة التغابير المحدود التي تحقق الكفاءة المزدوجة.

وبناء على ما توصلت له الدراسات السابقة وجود ندرة في تطبيق نموذج نسبة الكثافة على البيانات المراقبة تتجسد فكرة هذه الدراسة بمحاولة تقديم تقدير لمودج نسبة الكثافة باستخدام دالة الامكان الاعظم وبالتطبيق على البيانات المراقبة، كمحاولة من الباحثة لتقديم تقدير أكثر عمومية لتوفيق أكبر عدد ممكن من البيانات بمرونة أكبر.

٣ . نموذج نسبة الكثافة، ومواصفاته:

هو بديل شبه معلمي لحل مشكلة مقارنة عينتين أو أكثر حيث يعتمد تطبيقه بشكل كبير على الدالة $h(x)$ التي يفترض أنها معروفة ، حيث يؤدي التحديد الخاطئ لشكل الدالة في نموذج نسبة الكثافة إلى تغيرات متغيرة وقدان الكفاءة. فنموذج نسبة الكثافة هو نموذج انحدار شبه المعلمي يسمح بتحليل البيانات من أي عائلة أسيه دون اجراء افتراضات التوزيع المعلمي ، حيث يتميز النموذج بالقدرة في تحديد الاخطاء وكفاءة المقدر. وكما وصفه & Zhang (2021) بأنه أحد الطرق الفعالة لربط التوزيعات، حيث يمثل أحد المكونات الرئيسية في DRM في أن نسبة الكثافة اللوغاريتمية عبارة عن مجموعة خطية من دوال محددة سابقاً.

بفرض أن $(X_{(1)}, X_{(2)}, \dots, X_{(n_1)})$ عينة شوائية مستقلة لها دالة كثافة احتمال $f(x)$ ، وبفرض $(Y_{(1)}, Y_{(2)}, \dots, Y_{(n_2)})$ عينة مستقلة لها دالة كثافة احتمال $g(x)$ ولهمَا دالياً توزيع $F(X)$ ، $G(X)$ وكلما من $f(x)$ ، $g(x)$) دوال الكثافة المقابلة على التوالي، فيكون نموذج نسبة الكثافة له الشكل التالي :

$$g(x) = \exp\{(\alpha + \beta^T h(x)) f(x)\} \quad (1)$$

حيث:

$\alpha \leftarrow$ معلمة ثابتة غير معروفة للانحدار логистي

$\beta \leftarrow$ متجه معاملات المتغير x في النموذج

$h(x) \leftarrow$ دالة متجه قيم المتغيرات في النموذج

$f(x) \leftarrow$ دالة كثافة الاحتمال للمتغيرات $(X_{(1)}, X_{(2)}, \dots, X_{(n_1)})$

$g(x) \leftarrow$ دالة كثافة الاحتمال للمتغيرات $(Y_{(1)}, Y_{(2)}, \dots, Y_{(n_2)})$

ويطلق على نموذج نسبة الكثافة هذا الاسم لأنه يحدد النسبة اللوغاريتمية لدالتيين غير معروفتين لكثافة الاحتمال الخطية في بعض المعلمات (Anderson 1979) ويجد بالذكر أن فئة التوزيعات التي ينتمي إليها هي العائلة الأسيّة ، ومن ثم يمكن نمذجة مجموعة كبيرة من أنواع البيانات عن طريق نموذج نسبة الكثافة Konis ، Fokianos (2009). لقد جذب نموذج نسبة الكثافة اهتماماً كبيراً الفترة الأخيرة ، لأنه يخفف العديد من الافتراضات التقليدية في سياق حل مشاكل العينات المتعددة ، حيث يتميز DRM بالمرنة ويتضمن العديد من عائلات التوزيع المعلمية ، مثل

توزيعات اللوغاريتمات وتوزيع جاما كحالات خاصة . حيث يتم التعرف على DRM في الدراسات السابقة على أنها أداة قوية تشبه معلمية للعديد من المشاكل الاحصائية . فنموذج DRM له علاقة طبيعية في النماذج الخطية المعممة ويرتبط ارتباطاً وثيقاً بمشاكل اخذ البيانات المتغيره .

ويفرض ان لدينا متغيرين يقاء عشوائين (X^0, Y^0) مستقلين ولهمما دالتي كثافة احتمال $(g(x), f(x))$ ودالتي توزيع $(G(x), F(x))$ غير معروفيين ، حيث تعتبر نموذج نسبة الكثافة شبه المعلمي للدلائل $(g(x), f(x))$ هو :

$$g(x) = \exp(\alpha + \beta^T h(x))f(x) \quad (2)$$

عندما $\alpha = \beta = 0$ تؤدي الى $\alpha = 0$ ، كذلك ، حيث $h(x) = -\log \int e^{B^T h(x)} f(x) dx$ لها العديد من الاشكال المستخدمة في التوزيعات اكثراً استخاداماً : $h(x) = x$ ، $h(x) = \log(x)$ ، $h(x) = (x, x^2)^T$ ومن بين هذه الاشكال $(h(x) = x)$ يناسب الكثير من التوزيعات المعروفة ، مثل توزيعن أسيين لها متوسطين مختلفين أو توزيعن طبيعين لها نفس التابع ولكن المتوسطات مختلفة بينما $(h(x) = (x, x^2)^T)$ يناسب توزيعن طبيعين مختلفين المتوسط والتباين ، $(h(x) = \log(x))$ تناسب توزيعن جاما لها نفس المعلمة او توزيعن لوغاريتم الطبيعي لها نفس قيمة المعلمة .

في تحليل البقاء على قيد الحياة نهتم دائماً بالاستدلال الاحصائي للبيانات المراقبة ، بفرض (X^0, Y^0) نخضعان للمراقبة العشوائية بواسطة (V, U) ، حيث (V, U) متغيرات مراقبة مستقلة لها توزيع غير معروف $. H_u(x), H_v(x)$

في هذا البحث يتم اشتقاق مقدر الامكان الاعظم شبه المعلمي (SML) $\theta = (\alpha, \beta^T)^T$ ودالة التوزيع $F(x)$ لنموذج نسبة الكثافة (1) . حيث تلعب خوارزمية E-M (Vardi 1989) دور هام جداً في تنفيذ المقدر المقترن .

٤ . طريقة الامكان الاعظم لتقدير معلمات النموذج :

يتم التقدير بطريقة الامكان التجاري هي طريقة غير معلمية للاستدلال الإحصائي . يسمح لمحلل البيانات باستخدام طرق الامكان ، دون الحاجة إلى افتراض أن البيانات تأتي من مجموعة معروفة من التوزيعات . حيث أن مزاياها في أنها تجمع بين موثوقية الطرق الالمعلمية ومرنة وفعالية نهج الامكان وتنحصر بتقدير المعلم من نموذج نسبة الكثافة .

بفرض ان لدينا بيانات بقاء على قيد الحياة كالتالي: $X_i = X_i^0$, $\delta_i = 1$, $Y_j = Y_j^0$, $\eta_j = 1$ حيث : $i = 1, 2, \dots, n_1$, $j = 1, 2, \dots, n_2$ وبالتالي فان دالة الامكان الاعظم الكاملة هي :

$$\begin{aligned} L(\alpha, \beta, F) &= \prod_{i=1}^{n_1} f(X_{1i}^0) \prod_{j=1}^{n_2} g(Y_j^0) \\ &= \prod_{i=1}^{n_1} f(X_i^0) \prod_{j=1}^{n_2} f(Y_j^0) \exp(\alpha + \beta^T h(X)) \quad (3) \end{aligned}$$

وبفرض أن $b_{1k} = \sum_{i=1}^{n_1} I(X_i^0 = t_k)$, $b_{2k} = \sum_{j=1}^{n_2} I(Y_j^0 = t_k)$ عند الحدث

t_k ($k = 1, \dots, K$) في مجموعتي المتغيرات، حيث تلاحظ القيمة المميزة للمشاهدات بين المجموعات متزايدة حيث $t_1 < t_2 < \dots < t_K$ حيث يلاحظ أن القيم لوقت t_k في المشاهدات بين المجموعات تكون موجودة بترتيب تصاعدي ، وبفرض : $b_k = b_{1k} + b_{2k} = n_k$ وبالتالي فإن: $\sum_{k=1}^K b_{1k} = n_1$; $\sum_{k=1}^K b_{2k} = n_2$ وأيضاً $\sum_{k=1}^K b_k = n$ ، فان دالة الامكان الكاملة يمكن كتابتها كما يلي :

$$\begin{aligned} L(\alpha, \beta, F) &= \prod_{k=1}^K f(t_k)^{b_{1k}} g(t_k)^{b_{2k}} \\ &= \prod_{k=1}^K f(t_k)^{b_{1k}} \left(\exp(\alpha + \beta^T h(t_k)) \cdot f(t_k) \right)^{b_{2k}} \\ &= \prod_{k=1}^K f(t_k)^{b_{1k}} (\exp(\alpha + \beta^T h(t_k)) b_{2k} \cdot f(t_k))^{b_{2k}} \\ &= \prod_{k=1}^K f(t_k)^{b_k} \exp\{(\alpha + \beta^T h(t_k)) b_{2k}\} \end{aligned}$$

و تكون دالة لوغاریتم الامكان هي :

$$l(\alpha, \beta, F) = \sum_{k=1}^K \{b_k \log f(t_k) + b_{2k} (\alpha + \beta^T h(t_k))\} \quad (4)$$

وكما ذكر (Vardi 1989) الى تعليم (4) لأنها كافية وللنظر في دالة التوزيع المنفصلة F التي لديها الكتل الإيجابية تقترن على هذه الن نقاط $\{t_1, t_2, \dots, t_k\}$ فقط وحيث أن :

$$f(t_k) = p_k \quad , k = 1, \dots, K, \quad p = (p_1, \dots, p_K)$$

فكرين دالة لوغاریتم الامکان بالنسبة للمعلمات (α, β, p) كما يلي :

$$l(\alpha, \beta, p) = \sum_{k=1}^K \{ b_k \log p_k + b_{2k} (\alpha + \beta^T h(t_k)) \} \quad (5)$$

وطبقاً لطريقة كلاً من Qin & Lawless (1994) نعن نعزم (5) طبقاً للشروط :

$$\sum_{k=1}^K p_k = 1, \quad 0 \leq p_k \leq 1, \quad k = 1, \dots, K$$

$$\sum_{k=1}^K \exp\{\alpha + \beta^T h(t_k)\} p_k = 1 \quad (6)$$

ما يضمن أن كلاً من $g(x)$ هي دوال التوزيع وباستخدام طريقة مضاعف لاجرانج Lagrange تتحقق القيمة العظمى للقيم $l(\alpha, \beta, p)$ عند :

$$\tilde{p}_k = \frac{b_k}{n_1 + n_2 \exp(\alpha + \beta^T h(t_k))}$$

حيث $k = 1, 2, \dots, K$ للثوابت (α, β) ، وبالتالي تجاهل العناصر الثابتة ، تعرف دالة لوغاریتم الامکان $\mathcal{L}(\alpha, \beta)$ هي :

$$\begin{aligned} l(\alpha, \beta) &= \sum_{k=1}^K \{ b_k \log p_k + b_{2k} (\alpha + \beta^T h(t_k)) \} \\ &= \sum_{k=1}^K \left\{ b_k \log \frac{b_k}{n_1 + n_2 \exp(\alpha + \beta^T h(t_k))} + b_{2k} (\alpha + \beta^T h(t_k)) \right\} \\ &= \sum_{k=1}^K \left\{ b_k \log b_k - b_k \log(n_1 + n_2 \exp(\alpha + \beta^T h(t_k))) + \sum_{k=1}^K b_{2k} (\alpha \right. \\ &\quad \left. + \beta^T h(t_k)) \right\} \\ &= - \sum_{k=1}^K \left\{ b_k \log(n_1 + n_2 \exp(\alpha + \beta^T h(t_k))) + \sum_{k=1}^K b_{2k} (\alpha + \beta^T h(t_k)) \right\} \end{aligned}$$

وبحساب التفاضل لدالة لوغاریتم الامکان بالنسبة لكل من (α, β) والحصول على مقدارتهم $(\tilde{\alpha}, \tilde{\beta})$ فان :

$$l(\alpha, \beta) = \sum_{k=1}^K b_{2k} (\alpha + \beta^T h(t_k)) - \sum_{k=1}^K \{ b_k \log(n_1 + n_2 \exp(\alpha + \beta^T h(t_k))) \} \quad (7)$$

$$\frac{\partial l}{\partial \alpha} = \sum_{k=1}^K b_{2k} - \sum_{k=1}^K \frac{b_k n_2 \exp(\alpha + \beta^T h(t_k))}{n_1 + n_2 \exp(\alpha + \beta^T h(t_k))}$$

$$\frac{\partial l}{\partial \alpha} = n_2 - \sum_{k=1}^K \frac{b_k n_2 \exp(\alpha + \beta^T h(t_k))}{n_1 + n_2 \exp(\alpha + \beta^T h(t_k))} = 0$$

$$\frac{\partial l}{\partial \beta} = \sum_{k=1}^K b_{2k} h(t_k) - \sum_{k=1}^K \frac{b_k n_2 h(t_k) \exp(\alpha + \beta^T h(t_k))}{n_1 + n_2 \exp(\alpha + \beta^T h(t_k))} = 0$$

ولذلك نحصل على :

$$\tilde{p}_k = \frac{b_k}{n_1 + n_2 \exp(\tilde{\alpha} + \tilde{\beta}^T h(t_k))}$$

حيث : $k = 1, 2, \dots, K$ وبالاعتماد على مقدار الامكان الاعظم شب المعلمي $SML(\tilde{\alpha}, \tilde{\beta})$ و \tilde{p}_k فان تغير كل من $G(x), F(x)$ هو :

$$\tilde{F}(x) = \sum_{k=1}^K \tilde{p}_k I(t_k \leq x)$$

$$\tilde{G}(x) = \sum_{k=1}^K \tilde{p}_k \exp\left\{(\tilde{\alpha} + \tilde{\beta}^T h(t_k))\right\} I(t_k \leq x) \quad (8)$$

٥. تقيير المعلمات للبيانات المراقبة من اليمين: The right censored data

في هذا الجزء سنقوم بتطبيق المقدر السابق الى حالة البيانات المراقبة من اليمين ولبيانات بقاء مراقبة من اليمين تكون دالة الامكان الكاملة هي :

$$L(\alpha, \beta, F) = \prod_{i=1}^{n_1} f(X_i)^{\delta_i} (1 - F(X_i))^{1-\delta_i} \cdot \prod_{j=1}^{n_2} g(Y_j)^{\eta_j} (1 - G(Y_j))^{1-\eta_j}$$

$$= \prod_{i=1}^{n_1} f_1(X_i)^{\delta_i} (1 - F_1(X_i))^{1-\delta_i} \cdot \prod_{j=1}^{n_2} f_2(Y_j)^{\eta_j} \exp \left\{ (\alpha + \beta^T h(Y_j)) \eta_j \right\} (1 - F_2(Y_j))^{1-\eta_j}$$

ويفرض $(t_k < \dots < t_1)$ تدل على قيم متميزة بما في ذلك الخاضعة للمراقبة وغير الخاضعة للمراقبة حيث $K = 1, \dots, k$

$$r_{1k} = \sum_{i=1}^{n_1} I(X_i = t_k, \delta_i = 1), \quad \xi_{1k} = \sum_{i=1}^{n_1} I(X_i = t_k, \delta_i = 0)$$

$$r_{2k} = \sum_{j=1}^{n_2} I(Y_j = t_k, \eta_j = 1), \quad \xi_{2k} = \sum_{j=1}^{n_2} I(Y_j = t_k, \eta_j = 0)$$

و تكون المضاعفات لكل المشاهدات المراقبة وغير المراقبة في الوقت t في المجموعتين علي التوالى ، ثم يمكن توضيب دالة الامكان الكاملة كما يلى:

$$L(\alpha, \beta, F) = \prod_{k=1}^K f_0(t_k)^{r_{1k}+r_{2k}} \exp\{(\alpha + \beta^T h(t_k))r_{2k}\} (1 - F_0(t_k))^{\xi_{1k}} (1 - F_1(t_k))^{\xi_{2k}}$$

وتكون دالة الامكان اللوغاريتمية هي كما يلي :

$$l(\alpha, \beta, F) = \sum_{k=1}^K \{ (r_{1k} + r_{2k}) \log f_0(t_k) + r_{2k} (\alpha + \beta^T h(t_k)) + F_0(t_k) + \xi_{2k} \log(1 - F_0(t_k)) \} \quad (9)$$

وبعد التركيز على دالة التوزيع المنفصلة F التي لها كتلها إيجابية محصورة على النطاق $\{t_k, t_1, t_2, \dots\}$ للحصول على، نعطي (9) وبالمثل تدل على:

$$p_k = f_0(t_k), \quad q_k = f_1(t_k) = \exp\{(\alpha + \beta^T h(t_k))\} f_0(t_k)$$

$$k = 1, \dots, K \quad , p = (p_1, \dots, \dots, p_K)$$

فإن دالة لوغاریتم الامکان بالنسبة (α, β, p) هي :

$$l(\alpha, \beta, p) = \sum_{k=1}^K \left\{ (r_{1k} + r_{2k}) \log p_k + r_{2k} (\alpha + \beta^T h(t_k)) + \xi_{1k} \log (\sum_{l=k+1}^K p_l) + \xi_{2k} \log (\sum_{l=k+1}^K \exp(\alpha + \beta^T h(t_l)p_l)) \right\} \quad (10)$$

ولتعطيم (10) تخطي للشروط:

$$\sum_{k=1}^K p_k = 1, \quad 0 \leq p_k \leq 1, \quad k = 1, \dots, K$$

$$\sum_{k=1}^K \exp\{\alpha + \beta^T h(t_k)\} p_k = 1 \quad (11)$$

ويمكنا الحصول على مقدر SML α, β, p حيث نعرض مقدر SML بشكل فريد في ظل بعض الشروط المنتظمة ، حيث تطبيق (10) بطريقة مباشرة مهمه صعبه ، وبالتالي نجأ هنا إلى خوارزمية حسابية ممكنه لتنفيذها بسهولة وهي خوارزمية E-M ، التي اقترحها المؤلف Vardi (1989) للحصول على مقدر الامكان الاعظم غير المعلملي لدوال البقاء على قيد الحياة مع بيانات مراقبة مضايقة الان نحن بحاجه الي حساب المتوقع من كل تكرار وفي خطوة M حيث تقوم بتحديث مقدر لهذه المعلمات ، ويمكنا ان بندا بالتكرار :

$p^{(0)} = (p_1^{(0)}, \dots, p_K^{(0)})$, $\alpha^{(0)}, \beta^{(0)}$ التي تحقق الشرط في (11) ولتبسيط نحن دائمًا نضع القيمة الاولية $\alpha^{(0)} = 0, \beta^{(0)} = 0, p_1^{(0)} = \dots = p_K^{(0)} = \frac{1}{K}$ تم عملية التكرار وفقا للإجراءات التالية :

أولاً : خطوة التوقع E-step :

مع n من المشاهدات و المعلمات المقدرة $(\alpha^{(m-1)}, \beta^{(m-1)}, p^{(m-1)})$ ، فإن التكرار الشرطي المتوقع للمجموعة الاولى هو كالتالي:

$$\begin{aligned} b_{1k}^{(m)} &= \sum_{i=1}^{n_1} E[I(X_i^0 = t_k, \delta_i = 1) | (X_i, \delta_i), \alpha^{(m-1)}, \beta^{(m-1)}, p^{(m-1)}] \\ &\quad + \sum_{i=1}^{n_1} E[I(X_i^0 = t_k, \delta_i = 0) | (X_i, \delta_i), \alpha^{(m-1)}, \beta^{(m-1)}, p^{(m-1)}] \\ &= r_{1k} + \sum_{i=1}^{n_1} E[I(X_i^0 = t_k, \delta_i = 0) | (X_i, \delta_i), \alpha^{(m-1)}, \beta^{(m-1)}, p^{(m-1)}] \end{aligned}$$

حيث:

$$\begin{aligned}
 & E[I(X_i^0 = t_k, \delta_i = 0) | (X_i, \delta_i), \alpha^{(m-1)}, \beta^{(m-1)}, p^{(m-1)}] \\
 &= I(X_i^0 \leq t_k, \delta_i = 0) \frac{P(X_i^0 = t_k, U_i = X_i)}{\sum_{l:t_l \geq X_i} P(X_i^0 = t_l, U_i = X_i)} \\
 &= I(X_i \leq t_k, \delta_i = 0) \frac{f(t_k) dH_u(X_i)}{\sum_{l:t_l \geq X_i} f(t_l) dH_u(X_i)} \\
 &= p_K^{(m-1)} \sum_{j=1}^k I(X_i = t_j, \delta_i = 0) \frac{1}{\sum_{l:t_l \geq X_i} p_l^{(m-1)}}
 \end{aligned}$$

وتأتي المعادلة الثانية من فرض المراقبة المستقلة ، نحصل على :

$$\begin{aligned}
 b_{1k}^{(m)} &= r_{1k} + \sum_{j=1}^k I(X_i = t_j, \delta_i = 0) \frac{1}{\sum_{l:t_l \geq X_i} p_l^{(m-1)}} + \sum_{i=1}^{n_1} p_K^{(m-1)} \\
 &= r_{1k} + p_K^{(m-1)} \sum_{j=1}^k \sum_{i=1}^{n_1} \frac{I(X_i = t_j, \delta_i = 0)}{\sum_{l=j}^k p_l^{(m-1)}} \\
 &= r_{1k} + p_K^{(m-1)} \sum_{j=1}^k \frac{\xi_{1j}}{\sum_{l=j}^k p_l^{(m-1)}}
 \end{aligned}$$

وبالمثل فإن التكرار الشرطي المتوقع للمجموعة الثانية هو :

$$b_{2k}^{(m)} = r_{2k} + q_K^{(m-1)} \sum_{j=1}^k \frac{\xi_{2j}}{\sum_{l=j}^k q_l^{(m-1)}}$$

حيث:

$$q_k^{(m-1)} = \exp\{\alpha^{(m-1)} + \beta^{(m-1)} h(t_k)\} p_k^{(m-1)}$$

وفيما يلي: نقوم بتحديث المعندرات باستخدام مجموعة البيانات الكاملة الجديدة:

$$\{b_{1k}^{(m)}, b_{2k}^{(m)}, t_k; k = 1, \dots, K\}$$

ثانياً : خطوة التعظيم M-step

وفقاً للمعادلة (11) فإن دالة لوغاريتم الامكان لا هي كما يلي:

$$l^m(\alpha, \beta, p) = \sum_{k=1}^K \{(b_{1k}^{(m)} + b_{2k}^{(m)}) \log p_k + b_{2k}^{(m)}(\alpha + \beta^T h(t_k))\}$$

وباستخدام مضاعف لاجرانج Lagrange Multiplier مرة أخرى يكون الحد الاقصي لدالة لوغاريتم الامكان التي تحفظ الشروط (11) هي كما يلي :

$$p_K^{(m)} = \frac{b_{1k}^{(m)} + b_{2k}^{(m)}}{\omega_1^{(m)} + \omega_2^{(m)} \exp(\alpha^{(m)} + \beta^{(m)T} h(t_k))}, \quad k = 1, \dots, K,$$

حيث :

$$\omega_1^{(m)} = \sum_{k=1}^K b_{1k}^{(m)}, \quad \omega_2^{(m)} = \sum_{k=1}^K b_{2k}^{(m)}$$

و α^m, β^m هي جذر معادلات النتيجة التالية :

$$\frac{\partial l^m}{\partial \alpha} = \omega_2^{(m)} - \sum_{k=1}^K \frac{(b_{1k}^{(m)} + b_{2k}^{(m)}) \omega_2^{(m)} \exp(\alpha + \beta^T h(t_k))}{\omega_1^{(m)} + \omega_2^{(m)} \exp(\alpha + \beta^T h(t_k))} = 0$$

$$\frac{\partial l^m}{\partial \beta} = \sum_{k=1}^K b_{2k}^{(m)} h(t_k) - \sum_{k=1}^K \frac{(b_{1k}^{(m)} + b_{2k}^{(m)}) \omega_2^{(m)} h(t_k) \exp(\alpha + \beta^T h(t_k))}{\omega_1^{(m)} + \omega_2^{(m)} \exp(\alpha + \beta^T h(t_k))} = 0$$

نكرر كلا من الخطوة E والخطوة M بالتناوب حتى يتقارب كلاً من :

$$\theta^m = (\alpha^{(m)}, \beta^{(m)T})^T, p_K^{(m)}, k = 1, \dots, K$$

والتقريب يمكن تطبيق عدة معايير على الدراسات العددية حيث سيتم اثناء التكرار عندما تكون $\hat{\theta}$ معياراً الفرق بين اثنين من مقدرات التكرارات المجاورة أقل من الحد الأقصى المحدد، ونحصل أخيراً على مقدر SML للبيانات المراقبة ويشار اليه:

$G(x) = (\hat{\alpha}, \hat{\beta}^T)^T$, ومقدار SML هو $\hat{P} = (\hat{P}_1, \dots, \hat{P}_K)$

$$\begin{aligned}\hat{P}(x) &= \sum_{k=1}^K \hat{p}_k I(t_k \leq x), \\ \hat{G}(x) &= \sum_{k=1}^K \hat{p}_k \exp\{\hat{\alpha} + \hat{\beta}^T h(t_k)\} I(t_k \leq x),\end{aligned}\quad (12)$$

٦. الجانب التطبيقي:

نتناول في هذا الجزء التطبيق العملي لنموذج شبه المعلمي DRM لعينتين بافتراض المراقبة من جهة اليمين على بيانات واقعية تمثل عينة من مرضى الالتهاب الكبدي الوبائي المزمن (HCV). تشتمل البيانات على عدد 379 مريضاً بالغاً يعالون من النوع الجيني 4 من فيروس الالتهاب الكبدي الوبائي المزمن (Hepatitis C) والذين تم تعيينهم لتلقي العلاج من التهاب الكبد الفيروسي مع تركيبة جرعة ثابتة من OBV / PTV / RBV (75/50 / 12.5 مجم، بمعدل فرصة مرتاحية واحدة يومياً على أن تكون الجرعة على أساس الوزن ((1000 RBV مجم / يوم) إذا كان وزن الجسم أقل من 75 كجم أو جرعة (1200 مجم / يوم) إذا كان الوزن أكثر من 75 كجم، مع أو بدون 400 مجم من SOF يومياً (لأولئك الذين فشلوا في العلاج السابق بـ daclatasvir plus sofosbuvir (DCV) plus sofosbuvir (NCCVH)؛ وذلك خلال الفترة ما بين يناير 2016 ويونيو 2018، من مركز علاج التهاب الكبد الفيروسي بالقاهرة الجديدة مركز (NCVHTC)؛ أحد مراكز العلاج المتخصصة التابعة للجنة الوطنية لمكافحة التهاب الكبد الفيروسي (NCCVH) في مصر. تم تحديد رفت لجميع المرضى المعينين

لتقي العلاج PTV / RBV / OBV لمدة 12 أسبوعاً وفقاً لما إذا كانوا لم يسبق لهم أي علاج بأدوية الالتهاب

(Treatment Experienced) (مشاهدات غير مراقبة) أو 24 أسبوعاً (Treatment Naïve)

(مشاهدات مراقبة) إذا كانوا سبق لهم العلاج بأدوية أخرى للالتهاب الكبدي الفيروسي، وذلك وفقاً للمعايير الموحدة

لبروفوكول صادر عن مركز (NCCVH). (El Kassas 2019).

من هذه البيانات سوف نعتمد على ثلاثة متغيرات للتطبيق في نموذج نسبة الكثافة شبه المعلمي بين مجموعتين

في حالة المراقبة جهة اليمين، وهذه المتغيرات هي:

• جدول (١) : وصف المتغيرات المستخدمة في البحث

التعريف	اسم المتغير باللغة العربية	اسم المتغير بالإنجليزية
ويمثل نوع المريض الذكر بالكود 1 بينما نوع المريض الأنثى بالكود 0.	النوع	Gender
المعالجة الأولى: "Treatment Naïve" وهي عبارة عن المرضى الذين لم يسبق لهم تناول أي علاج للالتهاب الكبدي الوبائي الفيروسي، وسوف يتم تكديرها بالرقم 1	المعالجة	Treatment
المعالجة الثانية: "Treatment Experienced" وهي عبارة عن المرضى الذين سبق لهم تناول علاجات أخرى للالتهاب الكبدي الوبائي الفيروسي وسوف يتم تكديرها بالرقم 0.	البيليروبين	Total BILIRUBIN

في النموذج شبه المعلمي المقترن DRM يمثل متغير النوع المجموعتين حيث أن المجموعة الأولى تمثل عينة المرضى الذكور وعددها 218 مريض بنسبة مؤوية 57.5% بينما المجموعة الثانية تمثل عينة المرضى الإناث وعددها 161 بنسبة مؤوية 42.5%. بينما يمثل متغير المعالجة نوع المراقبة من جهة البين حيث أن المعالجة تمثل الوحدات غير المراقبة وعددها 327 مشاهدة بنسبة مؤوية 86.3% بينما المعالجة تمثل الوحدات المراقبة وعددها 52 بنسبة مؤوية 13.7% Treatment Experienced

وبالتالي يكون لدينا جدول الأفقران للعلاقة بين نوع المعالجة المستخدمة ومتغير النوع كما يلي

Table

	Female	Male
Treatment Experienced	17	35
Treatment Naïve	144	183

ومن الجدول السابق فلن نسبة المراقبة في المجموعة الأولى تمثل عدد الذكور تحت المراقبة ما يعادل تقريباً نسبة 60% بينما نسبة المراقبة في المجموعة الثانية تمثل عدد الإناث تحت المراقبة أي ما يعادل تقريباً نسبة 640%. نحدد دالة مؤشر كل مجموعة (المجموعة الأولى المؤشر لها I_1 والمجموعة الثانية المؤشر لها I_2) وهذا المؤشر يأخذ قيمة 1 إذا كانت الوحدات مراقبة ويأخذ قيمة 0 إذا كانت الوحدة غير مراقبة.

حيث أن فروض الدراسة هي:

الفرض العددي H_0 : لا يوجد اختلاف بين المجموعتين (ذكور وإناث) تحت نوع المعالجة المستخدم

الفرض البديل H_1 : يوجد اختلاف بين المجموعتين (ذكور وإناث) تحت نوع المعالجة المستخدم

فإذا كانت الأحتمالية المحسوبة للاحصائية // أقل من 5% فإن ذلك دليل على قبول الفرض البديل

ورفض الفرض العددي مما يعني أن هناك اختلاف معنوي ذو دلالة إحصائية بين المجموعتين (الذكور

والإناث) تحت نوع المعالجة المستخدم وباختبار المراقبة جهة اليمين. أما إذا كانت الأحتمالية المحسوبة

للحصائية // أكبر من 5% فإن ذلك دليل على رفض الفرضية البديلة وقبول الفرض العددي بمعنى أنه

لا يوجد اختلاف بين مجموعتي (الذكور والإناث) تحت نوع المعالجة المستخدم.

في الجدول التالي نستعرض نتائج تقديرات المعالم (α, β) والخطأ المعياري وكذلك فترة الثقة لكل معلمة

عند مستوى % 95 وتقدير R_{η} في حالة أن دالة الشغل تتبع التوزيع الأسوي والتوزيع اللوغاريثمي الطبيعي

بالتطبيق على البيانات الواقعية .

جدول رقم (٢): نتائج البيانات الحقيقة (المتوسط والخطأ المعياري وفترة الثقة عند مستوى 95%)

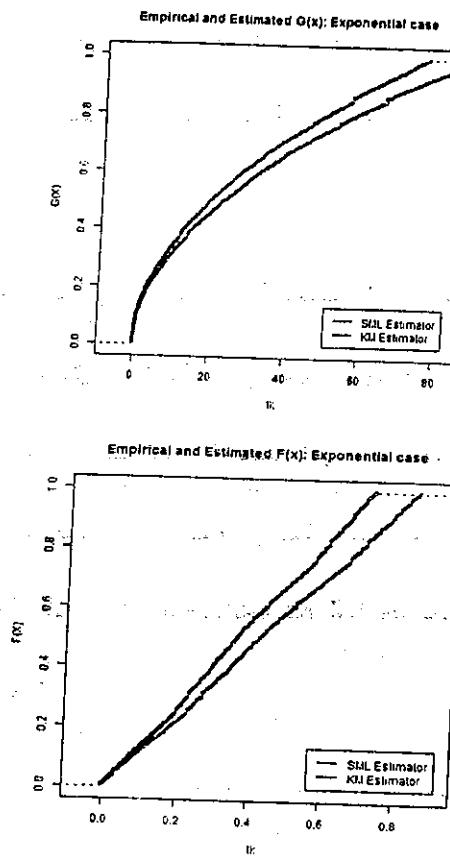
للتقديرات شبه المعلمية للمعلمتين (α, β) وكذلك تغير حجم نسبة الأمكان تبعاً لنوع التوزيع المفترض

تقدير p-value	تقدير لفترة الثقة		تقدير الخطأ المعياري		تقدير التغيير		التوزيع
	$\hat{\beta}$	$\hat{\alpha}$	$\hat{\beta}$	$\hat{\alpha}$	$\hat{\beta}$	$\hat{\alpha}$	
0.01254	(-0.6579, 0.7802)	(-0.2556, 0.9622)	0.9165	0.4308	0.2271	0.1996	الأسوي
0.01281	(-0.6613, 0.7669)	(-0.2738, 0.9679)	0.3161	0.3631	0.0533	0.3465	اللوغاريثم ال الطبيعي

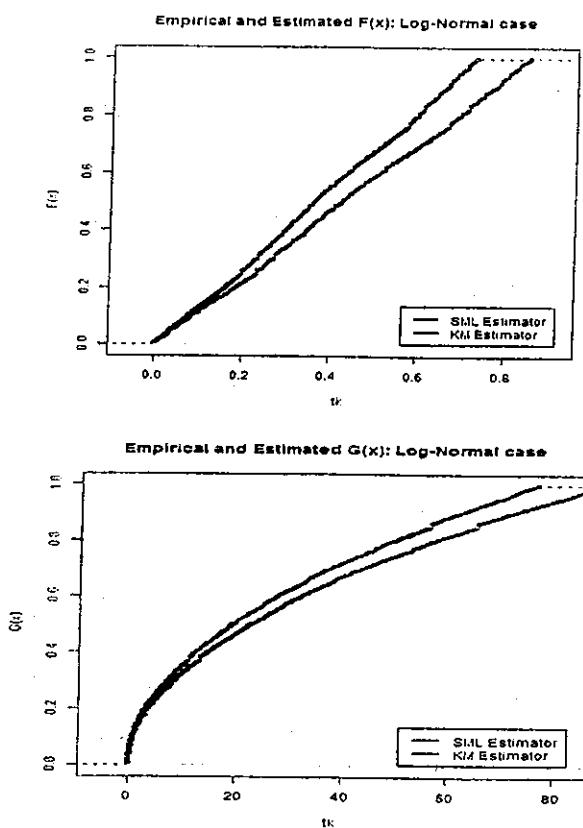
من الجدول السابق نستنتج أن هناك فرق بين المجموعتين محل التطبيق في حالة استخدام كلاً من التوزيع الأسوي أو التوزيع اللوغاريتمي الطبيعي حيث أن قيمة الأحتمالية المحسوبة $p\text{-value}$ في كلا الحالتين أقل من 5% وبالتالي لا نستطيع أن نقبل الفرض العلمي القائل بأن المجموعتين متساويتين في التأثير عليه يمكننا القول بأن المجموعتين مختلفتين في التأثير باستخدام المعالجة المقترحة في حالة مراعبة من جهة اليمين (كلا المجموعتين من الذكور والإناث).

الأشكال البيانية التالية توضح مقارنة بين تقدير $(x)\hat{F}$ وتقدير $(x)\hat{G}$ باستخدام مقدرات الأمكان الأعظم شبه الملمي (SML) ومقدر Kolmogorov-Smirnov (KM). في شكل (١) يوضح تقدير كلاً من تقدير $(x)\hat{F}$ وتقدير $(x)\hat{G}$ في حالة توزيع الأسوي بينما شكل (٢) يوضح تقدير كلاً من تقدير $(x)\hat{F}$ وتقدير $(x)\hat{G}$ في حالة توزيع اللوغاريتمي الطبيعي، وكما هو واضح من الأشكال البيانية تقارب طريقي التقدير: الأمكان الأعظم شبه الملمي (SML) ومقدر Kolmogorov-Smirnov (KM) وهو مطابق للنتائج النظرية التي تم الحصول عليها في جدول رقم (٢).

شكل (١) : مقارنة بين تدبر $\hat{G}(x)$ وتدبر $\hat{F}(x)$ باستخدام مقدرات الأمكان الأعظم شب المعلمي (SML)
ومقدر Kolmogorov-Smirnov (KM) وذلك في حالة التوزيع الأسوي



شكل (٢): مقارنة بين تقدير $\hat{F}(x)$ وتقدير $\hat{G}(x)$ باستخدام مقدرات الأمكان الأعظم شبه المعلمي (SML) ومقدر Kolmogorov-Smirnov (KM) وذلك في حالة التوزيع اللوغاريتمي الطبيعي



٧. النتائج:

يمكن إيجاز النتائج التي تم التوصل إليها في هذا البحث في النقاط التالية:

- عند استخدام البيانات الحقيقة باستخدام البيانات المرادفة من المبنين فيوجد هناك فرق بين المجنزعين محل التطبيق في حالة استخدام التوزيع الأسني فإن قيمة p-value (0.01254) أي أنها أقل من 5%

وبالتالي فإن المجموعتين مختلفتين في التأثير باستخدام المعالجة المقترحة في حالة المراقبة من جهة اليمين.

٢- عند التطبيق على حالة توزيع الطبيعي اللوغاريتمي فإن قيمة $p\text{-value}$ أيضاً (0.1281) أقل من 5% وبالتالي فإن المجموعتين (الذكور والإناث) مختلفتين في التأثير باستخدام المعالجة المقترحة في حالة المراقبة من جهة اليمين.

٣- تقارب ذاتي الكثافة عند استخدام دالة الامكان الأعظم شبه المعلمي SML ومقرر كولوجروف سميرنوف KM في حالتي التوزيعين توزيع الأسبي وتوزيع اللوغاريتمي الطبيعي مما يعني قوة الاختبار وأن النموذج المقترض ملائم لطبيعة البيانات في حالة وجود المراقبة من اليمين.

٨ . التوصيات

بناءً على ما توصلت إليه الدراسة من نتائج، توصي هذه الدراسة بالآتي:

١- استخدام طريقة تقدير أخرى لتقدير معلمات نموذج DRM مع البيانات المراقبة المختلفة أو بيانات عاربة.

٢- تطبيق نموذج نسبة الكثافة DRM في حالات التوزيعات الأخرى لدالة الفشل مثل: توزيع جاما $.h(x) = (x, \log|x|)$ ، أو التوزيع الطبيعي $(x, x^2) = h(x)$

٣- استخدام طريقة بيز لتقدير معلمات النموذج DRM ومقارنتها بنتائج تقديرات الامكان الأعظم.

المراجع

- 1) Anderson, J. A. (1979) "Multivariate logistic compounds". *Biometrika*, 66(1):17–26.
- 2) Anderson, T. W. and Darling, D. A. (1954), "A Test of Goodness-of-Fit". *Journal of the American Statistical Association*, 49(268), 765–769.
- 3) Cai, S. (2014). "On dual empirical likelihood inference under semiparametric density ratio models in the presence of multiple samples with applications to long term monitoring of lumber quality" (Doctoral dissertation, University of British Columbia).
- 4) Chakravarti, I. M., Laha, R. G., & Roy, J. (1967). "Handbook of methods of applied statistics." Wiley Series in Probability and Mathematical Statistics (USA) eng.
- 5) Chen, J., & Liu, Y. (2013). "Quantile and quantile-function estimations under density ratio model." *The Annals of Statistics*, 41(3), 1669–1692.
- 6) Cheng, K. F., & Chu, C. K. (2004). "Semi parametric density estimation under a two-sample density ratio model." *Bernoulli*, 10(4), 583–604.
- 7) Darling, D. A. (1957), "The Kolmogorov-Smirnov, Cramer-von Mises Tests". *Ann. Math. Statist*, 28(4), 823–838
- 8) Dempster, A. P., Laird, N. M., & Rubin, D. B. (1977). "Maximum likelihood from incomplete data via the EM algorithm." *Journal of the Royal Statistical Society: Series B (Methodological)*, 39(1), 1–22.

- 9) Diao, G., & Ning, J. (2012). "Maximum likelihood estimation for semi parametric density ratio model." *The international journal of biostatistics*, 8(1).
- 10) Epstein, B., & Sobel, M. (1953)." Life testing." *Journal of the American Statistical Association*, 48(263), 486–502.
- 11) El Kassas, M., Alboraei, M., Omar, H., El Latif, Y. A., Algaber, M. A., El Tahan, A. & Doss, W. (2019). High success rates for the use of ombitasvir/paritaprevir/ritonavir containing regimens in treatment of naïve and experienced chronic hepatitis C genotype 4: Real world results. *Journal of medical virology*. DOI: 10.1002/jmv.25478, 2019;1~8
- 12) Fokianos, K. (2007). "Density ratio model selection." *Journal of Statistical Computation and Simulation*, 77(9), 805–819
- 13) Zhang, A. G., & Chen, J. (2021). "Density ratio model with data-adaptive basis function." arXiv:2103.03445.
- 14) Zeng, D., Gao, F., & Lin, D. Y. (2017). "Maximum likelihood estimation for semi parametric regression models with multivariate interval-censored data." *Biometrika*, 104(3), 505–525.
- 15) Zhang, A. G., Zhu, G., & Chen, J. (2020). "Empirical Likelihood Ratio Test on quantiles under a Density Ratio Model." arXiv:2007.10586.
- 16) Zhuang, W. W., Hu, B. Y., & Chen, J. (2019). "Semi parametric inference for the dominance index under the density ratio model." *Biometrika*, 106(1), 229–241.