

התאמת מודלים לנתוני תאומים עם מידע חסר

עפרא עראידה

עבודת גמר מחקרית (תיזה) המוגשת כמילוי חלק מהדרישות
לקבלת התואר "מוסמך האוניברסיטה".

אוניברסיטת חיפה

הפקולטה למדעי החברה

החוג לסטטיסטיקה

נובמבר 2014

התאמת מודלים לנתוני תאומים עם מידע חסר

מאת: עפרא עראידה

בהנחיית: ד"ר, יאיר גולדברג.

עבודת גמר מחקרית (תזה) המוגשת כמילוי חלק מהדרישות

לקבלת התואר "מוסמך האוניברסיטה".

אוניברסיטת חיפה

הפקולטה למדעי החברה

החוג לסטטיסטיקה

נובמבר 2014

_____ תאריך:

_____ מאושר על ידי:
(מנחה העבודה)

_____ תאריך:

_____ מאושר על ידי:
(יו"ר הוועדה החוגית לתואר שני)

תוכן עניינים

iii תקציר	
iv רשימת איורים	
1 הקדמה	1
2 מודל ACE	2
4 פונקציית הנראות	3
4 3.1 נראות מירבית Maximum Likelihood	
4 3.2 הנראות עבור המודל המלא	
5 3.3 הנראות עבור המודל החסר	
7 Expectation-Maximization Algorithm (EM) אלגוריתם	4
7 4.1 הקדמה	
8 4.2 הפעלת אלגוריתם EM ואלגוריתם ECM	
11 4.2.1 סיכום האלגוריתם	
14 סימולציות	5
14 5.1 הפעלת הסימולציות	
14 5.1.1 התרחיש הראשון	
18 5.1.2 התרחיש השני	
22 5.2 תוצאות הסימולציות	
23 סיכום	6
24 ביבליוגרפיה	7
25 נספחים	8
25 8.1 אומדי נראות מירבית עבור המודל המלא	
28 8.2 הרחבה- אלגוריתם EM	
30 8.3 אמידה כללית לווקטור המשקולות π	
31 8.4 תרשימי Boxplot עבור הפרמטר μ	

התאמת מודלים לנתוני תאומים עם מידע חסר

עפרא עראידה

תקציר

מחקרי תאומים הם כלי חשוב לחקירת השפעת התורשה והסביבה על שונות בין בני האדם. מחקרים רבים שמתייחסים לנתוני תאומים מניחים שהסיווג של התאומים נתון, כלומר המידע עבור אם זוג התאומים זהים או לא זהים הוא נתון. במאגרי נתונים רבים, המידע עבור הסיווג לא נתון באופן מלא. התאמת מודלים תחת ההנחה שאין נתונים חסרים, ואמידת הפרמטרים למודלים אלה יכולה להביא לאומדים מוטים. ניתוח והסתמכות על האומדים המוטים עלול להוביל למסקנות מוטעות. לפיכך צריך למצוא שיטות לטיפול בבעיה זו.

ישנם חוקרים שמנסים לטפל בבעיה של המידע החסר עבור סוג התאומים במתן משקולות קבועות מראש לסיכוי של זוג תאומים להיות זהים או לא זהים. בעבודה הזו אנו מציעים מודל תערובת של התפלגויות לניתוח נתוני תאומים, שהמידע עבור הסיווג של התאומים לא בהכרח ידוע באופן מלא. להבדיל משיטות קודמות, שבהם משלימים את המידע החסר בסיכויים ומשקולות קבועים מראש, אנו מציעים אמדי נראות מירבית לבעיה זו. בעבודה זו מציאת אומדי הנראות המירבית מהווה אתגר חישובי, מאחר שהמקסום מתבצע תחת אילוצים הנובעים ממבנה מטריצת השונות. למציאת אומדי הנראות המירבית אנו נעזרים באלגוריתם ECM, Expectation-Conditional Maximization, למימוש האלגוריתם פיתחנו קוד בשפת R המשתמש בשיטות נומריות. בעזרת קוד זה הרצנו סימולציות לבחינת המודל שהוצע. התוצאות שהתקבלו העידו שהשימוש במודל שהצגנו, מודל תערובת של התפלגויות עם אלגוריתם ECM למציאת אומדי הנראות המירבית, הוא מודל שמספק תוצאות טובות יותר מאשר השיטות הקודמות להשלמת המידע החסר.

רשימת איורים

15	1 איור	1
16	2 איור	2
17	3 איור	3
19	4 איור	4
20	5 איור	5
21	6 איור	6
31	7 איור	7

1 הקדמה

מודלים המסבירים את השונות בין בני האדם כפונקציה של הסביבה המשפחתית, התורשה, והסביבה הכללית חשובים בתחומים רבים, כגון הנדסה גנטית, פסיכולוגיה, ביולוגיה, רפואה, ביטוח, ועוד. הבנה טובה יותר של הגורמים לשונות האלה תאפשר להתאים טיפולים רפואיים טובים יותר בהקשר הרפואי, תאפשר מציאת פתרונות לבעיות פסיכולוגיות בהקשר הפסיכולוגי, וכן בהקשר ביטוחי, תאפשר התאמת פרמיה מדויקת יותר למבוטחים. אחת הדרכים המקובלות להסברת השונות היא על ידי התאמת מודלים פרמטריים. התאמת מודלים פרמטריים היא פעולה שבה מנסים להתאים התפלגות ידועה עם מספר סופי של פרמטרים שמתארת את הנתונים בצורה סבירה. לדוגמא, התאמת רגרסיה לינארית, או השימוש במודלים ליניאריים מוכללים (GLM).

חשוב לציין כי בניתוח נתונים בכל מיני תחומים, כמו בתחום האקטואריה וכן בתחום הרפואי, נתקלים פעמים רבות בבעיות של מידע חסר. לדוגמא: גיל המבוטח, הקשר המשפחתי, ועוד. השפעת הגיל על השינויים במצב הבריאותי היא משמעותית, על כן צריך לחפש דרכים שמתייחסות לנתון זה אפילו אם הוא חסר. הקשר המשפחתי והתורשה הגנטית אף הם משפיעים על המצב הבריאותי ועל הסיכויים להתפתחות מחלות. לכן באנליזה סטטיסטית צריך להתחשב בנתונים החסרים, כי התעלמות מהמידע החסר יכולה להביא לאמידה מוטת.

מחקרים רבים שמתעסקים במודלים של שונות בין בני האדם, מתעסקים בנתונים עבור תאומים. מחקרי תאומים (Twin Study) חושפים את החשיבות של השפעות הסביבה והגנטיקה (התורשה). תאומים הינם מקור חשוב להתבוננות ומחקר, מאחר והם מאפשרים למידה עבור ההשפעה של הסביבה ושל הגנים. אנו נתמקד במחקרי תאומים שבהם נחקר את השפעת הסביבה והגנטיקה על השונות בין בני האדם.

הרעיון הבסיסי שמשתמשים בו על מנת לחקור את השפעת הסביבה לעומת הגנטיקה במחקרי תאומים מסתמך על סוג התאומים: תאומים זהים או מונוזיגוטיים (MZ) ותאומים לא זהים או דיזיגוטיים (DZ). במאגרי נתונים שבהם יש נתוני משפחות מנסים לחפש את התאומים לצורך מחקר תאומים. בעיה המתעוררת במאגרים אלה, שהמידע עבור סוג התאומים לרוב לא קיים. אפשר למצוא במאגרי נתונים מידע על תאריך לידה שמאפשר איתור התאומים. אפשר למצוא מידע על המין, שזה מספק לנו מידע עבור הסוג לחלק מהתאומים, אבל עדיין אנו נתקלים בבעיית מידע חסר עבור חלק מהתצפיות. המטרה בעבודה זו היא לפתח שיטות להתאמת מודלים עם נתונים חסרים עבור מחקרי תאומים, ולהציע אומדי נראות מירבית המתחשבים במידע החסר.

2 מודל ACE

מודל ACE בודק את ההשפעה היחסית של גנים וסביבה על השונות בין בני האדם. במודל ACE אנו מסתכלים על השונות הכוללת כמורכבת משלושה גורמים: השונות האדיטיבית הגנטית (A- Additive), השונות של הסביבה המשותפת (C- The common shared environment), והשונות של הסביבה הייחודית (E- The unique environment). כלומר מסתכלים על השונות הכוללת כ-

$V = a^2 + c^2 + e^2$ כאשר a^2 היא השונות הגנטית, c^2 היא השונות עבור הסביבה המשותפת, ו- e^2 היא השונות עבור הסביבה הייחודית. מודל ACE, כמו מודלים רבים במדעי החברה, שואף להסביר את השונות של האוכלוסייה כמורכבת מרכיבים שונים. באפן דומה לטכניקת ניתוח שונות (ANOVA), מודל ACE נותן אפשרות לעשות השוואה וניתוח שוניות בתוך הקבוצות וגם בין הקבוצות.

מחקרי תאומים, לעומת מחקרים שעוסקים בגנטיקה והתנהגות בין קשרי משפחה כלליים, מאפשרים אמידה פשוטה יותר למרכיבי השונות השונים. זה נובע כי תאומים זהים או מונוזיגוטיים (MZ) חולקים כמעט 100% מהגנים שלהם, לעומת תאומים לא זהים או דיזיגוטיים (DZ) שחולקים רק כ- 50% מהגנים שלהם. נוכל לדעת מהמידע הזה על ההבדל בקורלציה בין זוג תאומים זהים לבין זוג תאומים לא זהים. החלק של השונות הגנטית, a^2 , בקורלציה בין תאומים זהים הוא 1, לעומת החלק של a^2 בקורלציה בין תאומים לא זהים הוא 0.5. החלק של השונות המוסברת על ידי הסביבה המשותפת, c^2 , בקורלציה הוא 1 לשני סוגי התאומים DZ ו- MZ, והמרכיב e^2 בין כל זוג תאומים לא משפיע על הקורלציה. מאחר ויש בידינו את המידע הזה עבור תאומים, אמידת הפרמטרים תהיה פשוטה יותר מאשר אמידת פרמטרים של שוניות במדגם שמכיל קשרים משפחתיים כלליים יותר [11].

תחת הנחה שכל זוג תאומים חיו באותה סביבה, הבדלים בין התאומים מאפשרים לנו לאמוד את מרכיב השונות עבור a^2 בין קבוצות של תאומים זהים ותאומים לא זהים. מהקשרים האלה נוכל לבנות משוואות ולאמוד מספר פרמטרים נמוך מאשר עבור המקרה של קשר משפחתי כללי יותר. אם נסמן את הקורלציה בין תאומים זהים ותאומים לא זהים ב- r_{MZ} ו- r_{DZ} , בהתאמה, נוכל להסתכל על המשוואות הבאות כאשר אנו מניחים כי $a^2 + c^2 + e^2 = 1$:

$$\begin{aligned} a^2 &= 2r_{MZ} - 2r_{DZ} \\ c^2 &= 2r_{DZ} - r_{MZ} \\ e^2 &= 1 - r_{MZ} \end{aligned}$$

כלומר נוכל להסתכל על מטריצת השונות עבור תאומים זהים באופן הבא:

$$\Sigma_{MZ} = \sigma^2 \begin{pmatrix} 1 & a^2 + c^2 \\ a^2 + c^2 & 1 \end{pmatrix}$$

ועבור תאומים לא זהים:

$$\Sigma_{DZ} = \sigma^2 \begin{pmatrix} 1 & 0.5a^2 + c^2 \\ 0.5a^2 + c^2 & 1 \end{pmatrix}$$

מחקרי תאומים רבים מניחים כי המידע על סוג התאומים ידוע. אך במציאות במאגרי מידע, האינפורמציה של סוג התאומים לא תמיד נתונה ולעיתים קשה לבדיקה, במיוחד לאחר שהנתונים כבר נאספו.

התעלמות מהמידע החסר הזה יכולה להביא לאמדים מוטים [2]. לכן אנו מנסים להתאים מודל שמתייחס לנתונים של תאומים גם בהיעדר המידע של הסיווג לכל תאום נתון.

בעיה מסוג זה ניתנת לפתרון מתמטי בעזרת תערובת של התפלגויות [10], [7]. מודל תערובת של התפלגויות הוא מודל הסתברותי פרמטרי לייצוג שילוב של תת-אוכלוסיות לאוכלוסייה כללית. מודל זה שימושי כאשר יש לנו מדגם שבו לא יודעים לאיזו תת-אוכלוסייה כל תצפית שייכת. במקרה הספציפי שלנו, כאשר אנו מניחים כי הנתונים באים מהתפלגות נורמלית, ניתן להתאים מודל (GMM Gaussian Mixture Model), תערובת של גאוסיאנים.

מודל GMM הינו מודל הסתברותי פרמטרי, שפונקציית הצפיפות בו מיוצגת כסכום משוקלל של צפיפויות גאוסיאניות. אפשר להתייחס לכל מודל כזה, כמודל עם מידע חסר: המידע החסר הוא מאיזה התפלגות באה התצפית. נניח מדגם של n זוגות תאומים $x_i = (x_{i1}, x_{i2})$, $1 \leq i \leq n$ כל זוג מפולג נורמלי, אך הפרמטרים של ההתפלגות תלויים בתת-האוכלוסייה של תאומים זהים ותאומים לא זהים. נסמן $j = 1$ את המקרה של תאומים זהים, ו- $j = 2$ את המקרה של תאומים לא זהים.

כך ש- x_i שייך להתפלגות $N(\mu_j, \Sigma_j)$ בהסתברות π_j .
 כלומר בהסתברות π_j , $j = 1, 2$

$$x_i \sim N(\mu_j, \Sigma_j)$$

עם הצפיפות:

$$f_j(x_i) = \frac{1}{(2\pi)^n |\Sigma_j|^{1/2}} \exp\left(-\frac{1}{2}(x_i - \mu_j)^t \cdot \Sigma_j^{-1} \cdot (x_i - \mu_j)\right)$$

עבור $j = 1$ התאומים זהים (MZ), והשונוות והתוחלת הם (על פי מודל ACE):

$$\mu_1 \equiv \mu, \quad \Sigma_{MZ} = \sigma^2 \begin{pmatrix} 1 & a^2 + c^2 \\ a^2 + c^2 & 1 \end{pmatrix}$$

עבור $j = 2$ התאומים לא זהים (DZ), והשונוות והתוחלת הם (על פי מודל ACE):

$$\mu_2 \equiv \mu, \quad \Sigma_{DZ} = \sigma^2 \begin{pmatrix} 1 & 0.5a^2 + c^2 \\ 0.5a^2 + c^2 & 1 \end{pmatrix}$$

$$.a^2 + c^2 + e^2 = 1$$

בפרק הבא נציג את פונקציית הנראות המתאימה לנתונים כאלה, עם ההנחה שסוג התאומים נתון. פונקציה זו תהיה הנראות המלאה בנוסף נציג את הנראות עבור המצב שאין בידינו את כל המידע של סוג התאומים. פונקציה זו תהיה הנראות עבור מידע חסר.

3 פונקציות הנראות

3.1 נראות מירבית (Maximum Likelihood)

גישת הנראות המירבית היא אחת הגישות הסטטיסטיות השימושיות ביותר לאמידת פרמטרים. בשיטה זו מתאימים מודל, ואומדים את הפרמטרים על פי מדגם מקרי של נתונים. הגישה של נראות מירבית מחפשת את הפרמטר שמסביר את הנתונים באופן הטוב ביותר מתוך כל הפרמטרים האפשריים. כלומר, אומד הנראות המירבית הוא הפרמטר שנותן את ההסתברות הגבוהה ביותר (המקסימלית) לקבלת המדגם במודל.

נניח כי יש לנו מודל פרמטרי, כאשר $f(\vec{x}|\theta)$ היא פונקציית הצפיפות או ההסתברות לנתונים \vec{x} , עם ווקטור הפרמטרים θ , השייכים לקבוצת הפרמטרים Θ . נגדיר את פונקציית הנראות באופן הבא:

הגדרה 3.1 פונקציית הנראות (the likelihood function): בהינתן ערכים של נתונים $\vec{x} = (x_1, \dots, x_n)^t$, פונקציית הנראות $L(\theta|\vec{x})$ מוגדרת להיות הצפיפות (או ההסתברות) המשותפת: $L(\theta; x_1, x_2, \dots, x_n) = f_{x_1, x_2, \dots, x_n}(x_1, \dots, x_n; \theta)$ בפרט, אם x_1, \dots, x_n הינו מדגם מקרי מצפיפות (או הסתברות) $f(\cdot; \theta)$ אזי פונקציית הנראות הינה:

$$L(\theta; x_1, x_2, \dots, x_n) = f_{x_1, x_2, \dots, x_n}(x_1, \dots, x_n; \theta) = \prod_{i=1}^n f(x_i; \theta)$$

הגדרה 3.2 פונקציית לוג הנראות (the loglikelihood function): פונקציית לוג הנראות, $l(\theta|\vec{x})$, מוגדרת להיות הלוגריתם הטבעי של פונקציית הנראות:

$$l(\theta|\vec{x}) = \ln(L(\theta; x_1, x_2, \dots, x_n))$$

הגדרה 3.3 אומד נראות מירבית: (A maximum likelihood (ML) estimator of θ): אומד הנראות המירבית הוא הערך של θ שממקסם את פונקציית הנראות, או באופן שקול, ממקסם את פונקציית הלוג נראות.

בפרק הקודם דנו בשני סוגי נתוני תאומים: נתונים בהם מופיע סוג התאומים באופן מלא, ונתונים בהם סוג התאומים חסר לפחות לחלק מהתצפיות. בתתי הפרקים הבאים נציג קודם את הנראות המלאה עבור נתונים עם מידע מלא, ולאחר מכן נציג את הנראות עבור המידע החסר.

3.2 הנראות עבור המודל המלא

נתון מדגם של n זוגות של תאומים $x_i = (x_{i1}, x_{i2})$, $1 \leq i \leq n$ כך שכל זוג מפולג נורמלי. ידוע לנו הסוג של כל זוג תאומים. נניח כי מספר התאומים הזחים הוא k . לכן הנראות עבור התאומים הזחים היא:

$$\begin{aligned} L_{MZ}(x_1, \dots, x_k; \Sigma_{MZ}, \mu) &= \prod_{i=1}^k f_1(x_i; \Sigma_{MZ}, \mu) \\ &= \prod_{i=1}^k \frac{1}{(2\pi)^{1/2} |\Sigma_{MZ}|^{1/2}} \exp\left(-\frac{1}{2}(x_i - \mu)^t \cdot \Sigma_{MZ}^{-1} \cdot (x_i - \mu)\right) \end{aligned}$$

ועבור התאומים הלא זהים:

$$\begin{aligned} L_{DZ}(x_{k+1}, \dots, x_n; \Sigma_{DZ}, \mu) &= \prod_{i=k+1}^n f_2(x_i; \Sigma_{DZ}, \mu) \\ &= \prod_{i=k+1}^n \frac{1}{(2\pi)^{1/2} |\Sigma_{DZ}|^{1/2}} \exp\left(-\frac{1}{2}(x_i - \mu)^t \cdot \Sigma_{DZ}^{-1} \cdot (x_i - \mu)\right) \end{aligned}$$

בשימוש באי תלות בין הנתונים נוכל לקבל את הנראות עבור כל המדגם:

$$L(\vec{x}; \{\mu, a^2, c^2, \sigma^2\}) = \prod_{i=1}^k f_1(x_i; \Sigma_{MZ}, \mu) \cdot \prod_{i=k+1}^n f_2(x_i; \Sigma_{DZ}, \mu)$$

ולוג הנראות:

$$l(\vec{x}; \{\mu, a^2, c^2, \sigma^2\}) = \sum_{i=1}^k \log(f_1(x_i; \Sigma_{MZ}, \mu)) + \sum_{i=k+1}^n \log(f_2(x_i; \Sigma_{DZ}, \mu))$$

אמידת הפרמטרים במצב הזה היא יחסית פשוטה אף שאין ביטוי סגור. על מנת למקסם את $l(\vec{x}; \{\mu, a^2, c^2, \sigma^2\})$ יש לגזור את הפונקציה עבור כל פרמטר ולהשוות לאפס (ראה פיתוח בנספח 8.1).

3.3 הנראות עבור המודל החסר

נציג את פונקציית הנראות עבור המודל תחת ההנחה שאין בידינו את כל המידע עבור סוג התאומים. כפי שהצגנו בפרק קודם ננסה להסתכל על הבעיה הזו עם המידע החסר, כתערובת של התפלגויות. נסמן את קבוצת הפרמטרים ב- Θ . עבור כל זוג תאומים פונקציית הנראות נתונה על ידי:

$$L_i(\vec{x}; \{\mu, a^2, c^2, \sigma^2, \pi\}) = \pi_1 f_1(x_i, \theta) + \pi_2 f_2(x_i, \theta) = \sum_{j=1}^2 \pi_j f_j(x_i, \theta)$$

כאשר נסמן: $\theta = \{\mu, a^2, c^2, \sigma^2, \pi\}$.

מתקיים גם $\sum_{j=1}^2 \pi_j = 1$ כלומר: $\pi_2 = 1 - \pi_1$, לכן אפשר לכתוב:

$$L_i(x_i; \{\mu, a^2, c^2, \sigma^2, \pi_1\}) = \pi_1 f_1(x_i, \theta) + (1 - \pi_1) f_2(x_i, \theta)$$

כאשר $f_1(x_i, \theta)$ ו- $f_2(x_i, \theta)$ תלויות רק בפרמטרים $\{\mu, a^2, c^2, \sigma^2\}$. עבור תצפית של זוג תאומים זהים (MZ), מתקיים: $\pi_1 = 1$ ולכן $\pi_2 = 0$

$$L_{MZ}(x_i; \{\mu, a^2, c^2, \sigma^2\}) = f_1(x_i, \theta)$$

ועבור תצפית של זוג תאומים לא זהים (DZ), מתקיים: $\pi_1 = 0$ ולכן $\pi_2 = 1$

$$L_{DZ}(x_i; \{\mu, a^2, c^2, \sigma^2\}) = f_2(x_i, \theta)$$

עבור תצפית של זוג תאומים לא ידוע מה הסוג (NA) נקבל תערובת:

$$L_{NA}(x_i; \{\mu, a^2, c^2, \sigma^2, \pi_1\}) = \pi_1 f_1(x_i, \theta) + (1 - \pi_1) f_2(x_i, \theta)$$

כאשר הפרמטרים הם: $\theta = \{\mu, a^2, c^2, \sigma^2, \pi_1\}$

לכן הנראות של כל המדגם:

$$L(\vec{x}; \{\mu, a^2, c^2, \sigma^2, \pi_1\}) = \prod_{i=1}^n (\pi_1 f_1(x_i, \theta) + (1 - \pi_1) f_2(x_i, \theta))$$

ולוג הנראות:

$$l(\vec{x}; \{\mu, a^2, c^2, \sigma^2, \pi_1\}) = \sum_{i=1}^n \log(\pi_1 f_1(x_i, \theta) + (1 - \pi_1) f_2(x_i, \theta))$$

המטרה למצוא אומדים לפרמטרים על פי נתוני המדגם. מחפשים את האומד שממקסם את הפונקציה $L(\Theta)$. פונקציה כזו היא פונקציה לא לינארית בפרמטרים. מציאת אומדי נראות מירבית לא מתאפשרת כפתרון אנליטי. לכן נצטרך לפתח שיטות לפתרון בעיה מסוג זה. בפרק הבא נציג אלגוריתם Expectation–Maximization (EM) ואת ההרחבה Expectation–ECM, בפרק (Conditional Maximization) של אלגוריתם זה. שיעזרו לנו בפתרון בעיית מציאת אומד הנראות המירבית.

4 אלגוריתם Expectation–Maximization Algorithm (EM)

4.1 הקדמה

אלגוריתם EM [6] הוא שיטה איטרטיבית, השימושית בין השאר למציאת אומדי נראות מירבית עבור מידע חסר [1]. הרעיון העומד מאחורי EM, הוא לקחת את הנתונים הנצפים, להניח שיש לנו את הנתונים החסרים ואז להסתכל על כל הנתונים ביחד. נניח מדגם $\{x_1, x_2, \dots, x_n\}$, מדגם נצפה שאינו בהכרח מלא כלומר ישנו מידע חסר לגבי נתונים אלה, נניח $\{y_1, y_2, \dots, y_n\}$. במקרה שלנו $x_i = (x_{i1}, x_{i2})$ הוא ערכים עבור זוג של תאומים ו- y_i הוא סוג התאומים. באלגוריתם EM מניחים כי ההתפלגות של המדגם הנצפה ידועה, כלומר $P(X|Y, \theta)$ היא ידועה. במצב זה יש לאמוד את θ מתוך קבוצת הפרמטרים Θ , תוך התחשבות בנתונים החסרים. הרעיון שעומד מאחורי שיטת EM הוא להסתכל במדגם (X, Y) (המידע המלא) ולהשתמש בקשר:

$$P((X, Y)|\theta) = P(Y|X, \theta) \cdot P(X|\theta)$$

ממנו נוכל לקבל את המשוואה:

$$\log [P(X|\theta)] = \log [P((X, Y)|\theta)] - \log [P(Y|X, \theta)] \quad (1)$$

המטרה למצוא את הפרמטרים שממקסמים את הפונקציה הזו. על מנת למקסם את הפונקציה נסתכל בתוחלת של אגף ימין של (1), וננסה לחפש את המקסימום שלה (להרחבה, ראה נספח 8.2). כדי למצוא את המקסימום, נגדיר את הפונקציה הבאה:

$$Q(\theta, \theta^{(m)}) = E_{(Y|X, \theta^{(m)})} [\log(L((X, Y)|\theta)) | X, \theta^{(m)}] = \int \log(L((X, Y)|\theta)) \cdot f(Y|X, \theta^{(m)}) dY$$

שבעזרתה מבצעים את האלגוריתם. האלגוריתם מתבצע בשני צעדים והם:
צעד ה-E: חישוב תוחלת הנראות כלומר $Q(\theta, \theta^{(m)})$ כפונקציה של θ .
צעד ה-M: מציאת אומד של θ שממקסם את $Q(\theta, \theta^{(m)})$, כלומר, מחפשים $\theta^{(m+1)}$ כך ש-

$$Q(\theta^{(m+1)}, \theta^{(m)}) \geq Q(\theta, \theta^{(m)})$$

האלגוריתם מתבצע באופן רקורסיבי בתחילה בוחרים $\theta^{(0)}$, ועל סמך $\theta^{(0)}$ מבצעים את שני הצעדים של האלגוריתם ומחשבים את $\theta^{(1)}$. ממשיכים כך, משתמשים ב- $\theta^{(m)}$ לאמידת $\theta^{(m+1)}$ עד להתכנסות. אלגוריתם EM כפי שהוגדר לא תמיד עובד במצבים שיש אילוצים על הפרמטרים (להרחבה ראה [5]). המודל שאנו מציגים, כפי שהוגדר בפרק 2, מניח כי התוחלות בשני המצבים עבור תאומים זהים ולא זהים הן שוות וגם כי ישנו קשר בין השונויות. אפשר להסתכל על הקשרים האלה באופן הבא:

$$\mu_1 = \mu_2 = \mu \quad ; \quad \Sigma_{MZ} \succ \Sigma_{DZ}$$

הקשרים האלה הינם אילוצים שנצטרך להתחשב בהם בעת חיפוש הפתרון, לכן שימוש באלגוריתם EM נאיבי הוא לא מספיק. הסיבה לכך היא שצעד ה-M הופך למסובך מדי, ולא ניתן לקבל נוסחה

סגורה לפרמטרים. לכן, נצטרך להשתמש באחת ההרחבות של אלגוריתם EM שנקראת אלגוריתם Expectation-Conditional Maximization, ECM. אלגוריתם ECM הינו מחלקה של אלגוריתמים שהוצגו על ידי החוקרים Meng ו-Rubin [9]. אלגוריתם ECM מחליף את צעד ה-M "המסובך" בצעדים יותר פשוטים שבהם מעדכנים את הפרמטרים באופן הדרגתי שנקראים צעדי CM-Conditional Maximization, או בעברית מקסום מותנה. זאת אומרת, במקום לחשב את ווקטור הפרמטרים שממקסם את הפונקציה בצעד ה-M בשלב אחד, מחלקים את הצעד הזה ל-S צעדים. בכל צעד מחשבים פרמטר אחד שממקסם את הפונקציה בהינתן שאר הפרמטרים, כלומר מקבעים את כל הפרמטרים פרט לאחד, ומחפשים את המקסימום עבור אותו הפרמטר. במקרה שלנו, נחלק את צעד ה-M לשני צעדים, כלומר $S = 2$. בצעד ראשון, נעדכן את התוחלת והמשקולות, לאחר מכן נעדכן את השונות על סמך העדכון שנעשה לתוחלת וכך הלאה. בצעד ה-M של אלגוריתם זה, לאחר העדכון הראשון, חוזרים ומחשבים את הפונקציה שבצעד ה-E על סמך העדכון הראשון שנעשה. כלומר צעד ה-M מכיל בתוכו את צעד ה-E. לכן אלגוריתם ECM מצריך חישוב אחד נוסף יותר מאלגוריתם EM, ולפיכך התכנסותו יכולה להיות איטית יותר מהתכנסות אלגוריתם EM (להרחבה על התכנסות שני האלגוריתמים ראה [13], [8]).

4.2 הפעלת אלגוריתם EM ואלגוריתם ECM

בבעיה שלנו אנו חושבים על המודל כתערובת של התפלגויות. בפרט, יש לנו את המדגם הנצפה $\{x_1, x_2, \dots, x_n\}$, כך ש- $x_i = (x_{i1}, x_{i2})$, $1 \leq i \leq n$ וישנו מידע חסר לגבי נתונים אלה שהוא הסוג של כל אחד מהתאומים $\{y_1, y_2, \dots, y_n\}$. עבור המידע המלא, נסמן את האינדיקטור z_{ij} באופן הבא:

$$z_{ij} = \begin{cases} 1 & x_i \sim N(\mu_j, \Sigma_j) \\ 0 & otherwise \end{cases} \text{ עבור } 1 \leq j \leq 2$$

כלומר x_i זוג תאומים שייך להתפלגות $N(\mu_1, \Sigma_1)$ בהסתברות π_1 , ושייך להתפלגות $N(\mu_2, \Sigma_2)$ בהסתברות $\pi_2 = 1 - \pi_1$, כאשר: $\Sigma_1 \equiv \Sigma_{MZ}$ ו- $\Sigma_2 \equiv \Sigma_{DZ}$. ההתפלגות של המדגם הנצפה ידועה, כלומר $P(X|Y, \theta)$ הינה התפלגות נורמלית. במצב זה יש לאמוד את θ מקבוצת הפרמטרים $\Theta = \{(\mu, a^2, \sigma^2, c^2, \pi_1)\}$, תחת האילוצים $\mu_1 = \mu_2$ ו- $\Sigma_{MZ} \succ \Sigma_{DZ}$, ותוך התחשבות בנתונים החסרים.

נפעיל את האלגוריתם:

בונים את הפונקציה:

$$\begin{aligned}
Q(\theta, \theta^{(m)}) &= E_{Z|X, \theta^{(m)}} [\log (L(Z, X|\Theta)|X, \theta^{(m+1)})] \\
&= E_{Z|X, \theta^{(m)}} \left[\log \prod_{i=1}^n (L(\theta, x_i, z_i)) \right] \\
&= \sum_{i=1}^n \sum_{j=1}^2 P(z_{ij}|x_i, \theta^{(m)}) \log (\pi_j f_j(x_i, \theta)) \\
&= \sum_{i=1}^n \sum_{j=1}^2 P(z_{ij}|x_i, \theta^{(m)}) \left[\log \pi_j - \frac{1}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_j| - \frac{1}{2} (x_i - \mu_j)^t \cdot \Sigma_j^{-1} \cdot (x_i - \mu_j) \right]
\end{aligned} \tag{2}$$

כאשר על פי נוסחת בייס:

$$P(z_{ij}|x_i, \theta^{(m)}) = \frac{\pi_j^{(m)} f_j(x_i|\theta^{(m)})}{\pi_1^{(m)} f_1(x_i|\theta^{(m)}) + \pi_2^{(m)} f_2(x_i|\theta^{(m)})}$$

$p_{ij}^{(m)} = P(z_{ij}|x_i, \theta^{(m)})$: נסמן:
 כלומר נקבל כי,

$$p_{i1}^{(m)} = P(z_{i1}|y_i, \theta^{(m)}) = \frac{\pi_1^{(m)} f_1(x_i|\theta^{(m)})}{\pi_1^{(m)} f_1(x_i|\theta^{(m)}) + (1 - \pi_1^{(m)}) f_2(x_i|\theta^{(m)})} \tag{3}$$

וגם,

$$p_{i2}^{(m)} = P(z_{i2}|y_i, \theta^{(m)}) = \frac{\pi_2^{(m)} f_2(x_i|\theta^{(m)})}{\pi_1^{(m)} f_1(x_i|\theta^{(m)}) + (1 - \pi_1^{(m)}) f_2(x_i|\theta^{(m)})} = 1 - p_{i1}^{(m)}$$

לכן הפונקציה (2),

$$\begin{aligned}
Q(\theta, \theta^{(m)}) &= \sum_{i=1}^n p_{i1}^{(m)} \left[\log \pi_1 - \frac{1}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_1| - \frac{1}{2} (x_i - \mu)^t \cdot \Sigma_1^{-1} \cdot (x_i - \mu) \right] \\
&\quad + \sum_{i=1}^n p_{i2}^{(m)} \left[\log \pi_2 - \frac{1}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_2| - \frac{1}{2} (x_i - \mu)^t \cdot \Sigma_2^{-1} \cdot (x_i - \mu) \right] \\
&= \sum_{i=1}^n p_{i1}^{(m)} \left[\log \pi_1 - \frac{1}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_1| - \frac{1}{2} (x_i - \mu)^t \cdot \Sigma_1^{-1} \cdot (x_i - \mu) \right] \\
&\quad + \sum_{i=1}^n (1 - p_{i1}^{(m)}) \left[\log(1 - \pi_1) - \frac{1}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_2| - \frac{1}{2} (x_i - \mu)^t \cdot \Sigma_2^{-1} \cdot (x_i - \mu) \right]
\end{aligned} \tag{4}$$

נאמוד את הפרמטרים על ידי גזירה של הפונקציה $Q(\theta, \theta^{(m)})$:

נתחיל מהפרמטר של המשקולות π_1 שלא תלוי בתוחלת ולא בשונות. נגזור את (4) לפי π_1

ונקבל:

$$d \frac{\Theta(\theta, \theta^{(m)})}{d\pi_1} = \sum_{i=1}^n p_{i1}^{(m)} \cdot \frac{1}{\pi_1} - \sum_{i=1}^n (1 - p_{i1}^{(m)}) \cdot \frac{1}{1 - \pi_1}$$

(לדרך כללית יותר לקבלת האומד למשקולות ראה נספח 8.3)

נשווה לאפס ונחלץ את π_1 :

$$\begin{aligned} 0 &= \sum_{i=1}^n p_{i1}^{(m)} \cdot \frac{1}{\pi_1} - \sum_{i=1}^n (1 - p_{i1}^{(m)}) \cdot \frac{1}{1 - \pi_1} \\ \Leftrightarrow 0 &= (1 - \pi_1) \cdot \sum_{i=1}^n p_{i1}^{(m)} - \pi_1 \cdot \sum_{i=1}^n (1 - p_{i1}^{(m)}) \\ \Leftrightarrow 0 &= \sum_{i=1}^n p_{i1}^{(m)} - \pi_1 \cdot \sum_{i=1}^n p_{i1}^{(m)} - \pi_1 n + \pi_1 \cdot \sum_{i=1}^n p_{i1}^{(m)} \\ \Leftrightarrow \pi_1 n &= \sum_{i=1}^n p_{i1}^{(m)} \end{aligned}$$

$$\pi_1^{(m+1)} = \frac{1}{n} \sum_{i=1}^n p_{i1}^{(m)}$$

כלומר:

$$\begin{aligned} \pi_1^{(m+1)} &= \frac{1}{n} \sum_{i=1}^n p_{i1}^{(m)} \\ &= \frac{1}{n} \sum_{i=1}^n \frac{\pi_1^{(m)} f_1(x_i | \theta^{(m)})}{\sum_{r=1}^2 \pi_r^{(m)} f_r(x_i | \theta^{(m)})} \\ &= \frac{1}{n} \sum_{i=1}^n \frac{\pi_1^{(m)} f_1(x_i | \theta^{(m)})}{\pi_1^{(m)} f_1(x_i | \theta^{(m)}) + (1 - \pi_1^{(m)}) f_2(x_i | \theta^{(m)})} \end{aligned}$$

עד עכשיו קיבלנו נוסחה לחישוב המשקולות. נעבור מכאן לחישוב התוחלת תוך התחשבות באילוץ. נגזור את (4) לפי $\vec{\mu}$ ונקבל:

$$d \frac{Q(\theta, \theta^{(m)})}{d\mu} = \sum_{i=1}^n \sum_{j=1}^2 P(z_{ij} | x_i, \theta^{(m)}) \frac{1}{2} \cdot (-2) \Sigma_j^{-1} (x_i - \mu)$$

נשווה לאפס ונחלץ $\vec{\mu}$:

$$0 = \sum_{i=1}^n \sum_{j=1}^2 P(z_{ij}|x_i, \theta^{(m)}) \Sigma_j^{-1} x_i - \sum_{i=1}^n \sum_{j=1}^2 P(z_{ij}|x_i, \theta^{(m)}) \Sigma_j^{-1} \mu$$

נחשב:

$$\sum_{i=1}^n \sum_{j=1}^2 P(z_{ij}|x_i, \theta^{(m)}) \Sigma_j^{-1} x_i = \sum_{i=1}^n \sum_{j=1}^2 P(z_{ij}|x_i, \theta^{(m)}) \Sigma_j^{-1} \mu$$

נקבל:

$$\mu^{(m+1)} = \left(\sum_{i=1}^n \sum_{j=1}^2 P(z_{ij}|x_i, \theta^{(m)}) \Sigma_j^{-1} \right)^{-1} \left(\sum_{i=1}^n \sum_{j=1}^2 P(z_{ij}|x_i, \theta^{(m)}) \Sigma_j^{-1} x_i \right)$$

ולכן נסכם:

$$\begin{aligned} \mu^{(m+1)} &= \left(\sum_{i=1}^n \sum_{j=1}^2 P(z_{ij}|x_i, \theta^{(m)}) \Sigma_j^{-1} \right)^{-1} \left(\sum_{i=1}^n \sum_{j=1}^2 P(z_{ij}|x_i, \theta^{(m)}) \Sigma_j^{-1} x_i \right) \\ &= \frac{\sum_{i=1}^n \sum_{j=1}^2 P(z_{ij}|x_i, \theta^{(m)}) \Sigma_j^{-1} x_i}{\sum_{i=1}^n \sum_{j=1}^2 P(z_{ij}|x_i, \theta^{(m)}) \Sigma_j^{-1}} \end{aligned}$$

קיבלנו נוסחה סגורה לתוחלת.

עבור המקרה שלנו מציאת אומד לכל אחד מהפרמטרים: a^2, c^2 ו- σ^2 , או לשונות באופן כללי, אינה מתאפשרת מבחינה אנליטית. אם נגזור את הפונקציה (4) לפי הפרמטרים של השונות לא נוכל לחלץ את הפרמטרים ולקבל נוסחה סגורה בגלל התלות בין שתי השונות Σ_1, Σ_2 . לכן נשתמש בשיטות נומריות של הנראות. לקבלת אומדים ל- a^2, c^2 ו- σ^2 .

4.2.1 סיכום האלגוריתם

נתון המדגם $\vec{x} = (x_1, \dots, x_n)$, שמכיל n זוגות של תאומים. נרצה להתאים לו מודל של תערובת התפלגויות, ולאמוד את הפרמטרים של הנראות. כלומר יש לנו (x_1, \dots, x_n) זוגות של תאומים. מסתכלים על זה כמטריצה $X_{n \times 2}$, (כל תצפית הינה זוג של תאומים). משלימים את הנתונים האלה בוקטור $Y_{n \times 1}$, כך שווקטור זה מייצג את הסוג של התאומים. y_i שווה ל- DZ, MZ או NA . משתמשים באלגוריתם ECM:

צעד ה-E:

מחשבים את הפונקציה:

$$\begin{aligned} Q(\theta, \theta^{(m)}) &= \\ &= \sum_{i=1}^n p_{ij}^{(m)} \left[\log \pi_1 - \frac{1}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_1| - \frac{1}{2} (x_i - \mu)^t \cdot \Sigma_1^{-1} \cdot (x_i - \mu) \right] \\ &+ \sum_{i=1}^n (1 - p_{ij}^{(m)}) \left[\log(1 - \pi_1) - \frac{1}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_2| - \frac{1}{2} (x_i - \mu)^t \cdot \Sigma_2^{-1} \cdot (x_i - \mu) \right] \end{aligned}$$

כאשר:

$$p_{i1}^{(m)} = P(z_{i1}|x_i, \theta^{(m)}) = \frac{\pi_1^{(m)} f_1(x_i|\theta^{(m)})}{\pi_1^{(m)} f_1(x_i|\theta^{(m)}) + (1-\pi_1^{(m)}) f_2(x_i|\theta^{(m)})}$$

צעד CM הראשון:

מחשבים את המשקולות:

$$\pi_1^{(m+1)} = \frac{1}{n} \sum_{i=1}^n \frac{\pi_1^{(m)} f_1(x_i|\theta^{(m)})}{\pi_1^{(m)} f_1(x_i|\theta^{(m)}) + (1-\pi_1^{(m)}) f_2(x_i|\theta^{(m)})}$$

ואת התוחלת:

$$\mu^{(m+1)} = \frac{\sum_{i=1}^n \sum_{j=1}^2 P(z_{ij}|x_i, \theta^{(m)}) \left(\Sigma_j^{(m)}\right)^{-1} x_i}{\sum_{i=1}^n \sum_{j=1}^2 P(z_{ij}|x_i, \theta^{(m)}) \left(\Sigma_j^{(m)}\right)^{-1}}$$

צעד CM השני:

נחזור ונעדכן את הפונקציה בצעד ה- E , כלומר נחשב את (3) ו-(4) עם העידכון החדש של התוחלת והמשקולות. ונקבל:

$$p_{i1}^{(m+\frac{1}{2})} = P(z_{i1}|x_i, \theta^{(m)}) = \frac{\pi_1^{(m+1)} f_1(x_i|\theta^{(m+\frac{1}{2})})}{\pi_1^{(m+1)} f_1(x_i|\theta^{(m+\frac{1}{2})}) + \pi_2^{(m+1)} f_2(x_i|\theta^{(m+\frac{1}{2})})}$$

$$Q\left(\theta, \theta^{(m+\frac{1}{2})}\right) = \sum_{i=1}^n p_{i1}^{(m+\frac{1}{2})} \left[\log \pi_1 - \frac{1}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_1| - \frac{1}{2} (x_i - \mu)^t \cdot \Sigma_1^{-1} \cdot (x_i - \mu) \right] + \sum_{i=1}^n \left(1 - p_{i1}^{(m+\frac{1}{2})}\right) \left[\log(1 - \pi_1) - \frac{1}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_2| - \frac{1}{2} (x_i - \mu)^t \cdot \Sigma_2^{-1} \cdot (x_i - \mu) \right]$$

חישוב הפרמטרים $(a^2)^{(m+1)}$, $(c^2)^{(m+1)}$, $(\sigma^2)^{(m+1)}$ שממקסמים את הפונקציה $Q\left(\theta, \theta^{(m+\frac{1}{2})}\right)$ מתבצע על ידי שיטה נומרית.

ולכן הלולאה שנבנה תסתכם בשלבים הבאים:

(1) בוחרים תנאים התחלתיים $\mu^{(0)}$, $\pi_1^{(0)}$, $\Sigma_j^{(0)}$, כאשר $\Sigma_j^{(0)}$ תלויה ב- $a^{2(0)}$, $c^{2(0)}$, $\sigma^{2(0)}$.

(2) בונים את הפונקציה: $Q\left(\theta, \theta^{(m)}\right)$ ומחשבים:

$$p_{ij}^{(m)} = P(z_{ij}|x_i, \theta^{(m)}) = \frac{\pi_j^{(m)} f_j(x_i|\theta^{(m)})}{\sum_{r=1}^k \pi_r^{(m)} f_r(x_i|\theta^{(m)})}$$

(3) מעדכנים את הפרמטר $\vec{\pi}$:

$$\pi_j^{(m+1)} = \frac{1}{n} \sum_{i=1}^n p_{ij}^{(m)}$$

(4) מעדכנים את התוחלת $\vec{\mu}$:

$$\mu^{(m+1)} = \frac{\sum_i \sum_j p_{ij}^{(m)} (\Sigma_j^{(m)})^{-1} x_i}{\sum_i \sum_j p_{ij}^{(m)} (\Sigma_j^{(m)})^{-1}}$$

(5) חוזרים ומעדכנים את p_{ij} :

$$p_{ij}^{(m+\frac{1}{2})} = P(z_{ij}|x_i, \Sigma_j^{(m)}, \mu^{(m+1)}, \pi_j^{(m+1)}) = \frac{\pi_j^{(m+1)} f_j(x_i|\Sigma_j^{(m)}, \mu^{(m+1)})}{\sum_{r=1}^2 \pi_r^{(m+1)} f_r(x_i|\Sigma_r^{(m)}, \mu^{(m+1)})}$$

(6) ממקסמים את הפונקציה $Q(\theta, \theta^{(m+\frac{1}{2})})$ בעזרת שיטות נומריות (למשל פונקציית Optim ב-R) ומקבלים את האומדים: $(a^2)^{(m+1)}, (c^2)^{(m+1)}, (\sigma^2)^{(m+1)}$.

(7) חוזרים לשלב 2 עד לקירוב הרצוי.

5.1 הפעלת הסימולציות

לבדיקת המודל שהצענו, ביצענו מספר סימולציות. כל הסימולציות והחישובים לבדיקה התבצעו בתוכנת R. לחישוב הפונקציה נעזרנו בחבילת `mvtnorm` [3]. נציג את הסימולציות עבור שני תרחישים.

5.1.1 התרחיש הראשון

בתרחיש הראשון הגרלנו מדגם אקראי של נתונים, כך שבממוצע חצי מהתאומים יהיו זהים והחצי השני לא זהים, באופן הבא: הגרלנו 30% תאומים זהים ו- 30% תאומים לא זהים והשאר לא ידוע (NA). עבור אלה הלא ידועים הנחנו כי 50% הם זהים (MZ) ו- 50% הם לא זהים (DZ). ביצענו מספר סימולציות עבור תרחיש זה: בכל פעם הגרלנו את הנתונים מהתפלגות שבה קבענו ערכים שונים לפרמטרים: μ, a^2, c^2 ו- σ^2 . לכל ערך של ווקטור הפרמטרים שקבענו הגרלנו נתונים עבור גדלי המדגם הבאים:

$n = 50, 100, 200, 400, 800$. עבור כל n הרצנו את הסימולציה 100 פעמים עבור שלוש שיטות. בשיטה הראשונה הנחנו עבור חלק מהנתונים שסוג התאומים נתון וידוע, ועבור חלק מהנתונים שהסוג לא ידוע. לכן השתמשנו בשיטה של המידע החסר שהוצעה בתת-פרק 4.2 שבו אמדנו את חמשת הפרמטרים: μ, σ^2, a^2, c^2 ו- π_1 . השיטה השנייה היא השיטה שבה עבור התאומים שלא ידוע הסוג שלהם הנחנו כי $\pi_1 = 0.5$, כלומר הנחנו משקולות נתונות מראש, ואמדנו את הפרמטרים μ, a^2, c^2 ו- σ^2 . השיטה השלישית היא אמידת נראות מירבית עבור המודל המלא. בשיטה זו, בניגוד לשתי השיטות הקודמות המידע הכיל את סוג התאומים עבור כל התצפיות. השיטה השלישית נועדה לצורך השוואה.

באיור 1 אפשר לראות את התוצאות בתרשימי BOXPLOTS עבור הפרמטרים:

$a^2 = 0.4, c^2 = 0.3, \sigma^2 = 2, \mu = 5$. ניתן לראות שהאומד של σ^2 הופך למדויק יותר ככל שמגדילים את גודל המדגם. האומדים לפרמטרים a^2 ו- c^2 פחות קרובים מאשר האומד של σ^2 , אבל ניתן לראות שהשיטה הראשונה שבה לקחנו את הנראות עבור המידע החסר (הקופסא השמאלית האדומה) טובה יותר מהשיטה שבה הנחנו $\pi_1 = 0.5$ (הקופסא האמצעית הכחולה). האומד של μ יצא קרוב מאוד וכמעט זהה בכל התרחישים¹.

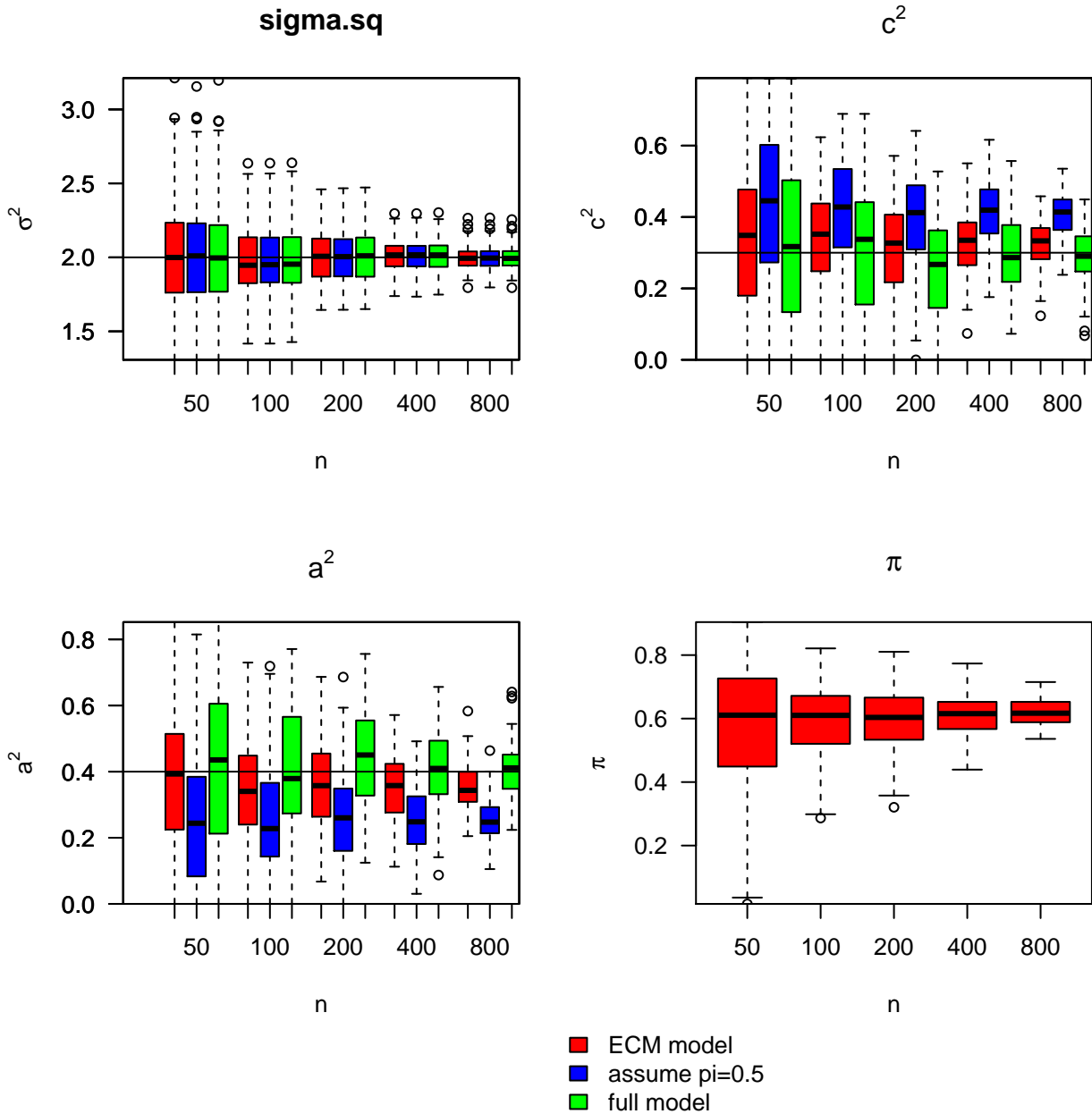
הרצנו עוד מספר פעמים את התרחיש הזה, בכל פעם שינינו את הערכים של פרמטרים. באיור 2 התוצאות מופיעות בתרשימי BOXPLOTS עבור הפרמטרים:

$a^2 = 0.3, c^2 = 0.3, \sigma^2 = 2, \mu = 5$. גם בסימולציה זו רואים שהאומד של σ^2 מתקרב ככל שמגדילים את גודל המדגם. האומדים לפרמטרים a^2 ו- c^2 במצב שבו $a^2 = c^2$ טובים יותר ממה שראינו בסימולציה הקודמת באיור 1.

באיור 3 רואים את ההרצה לאחר שינוי נוסף בערך של σ^2 . הפרמטרים שנבחרו כעת הם: $a^2 = 0.6, c^2 = 0.2, \sigma^2 = 1.33, \mu = 5$. שוב נוכל לראות שהאומד של σ^2 כמעט זהה בשלושת השיטות. אך כאן רואים שהאומדים לפרמטרים a^2 ו- c^2 פחות טובים ממצבים אחרים, אבל עדיין יותר קרובים בשיטה הראשונה שהוצעה בעבודה מאשר מהשיטה השנייה.

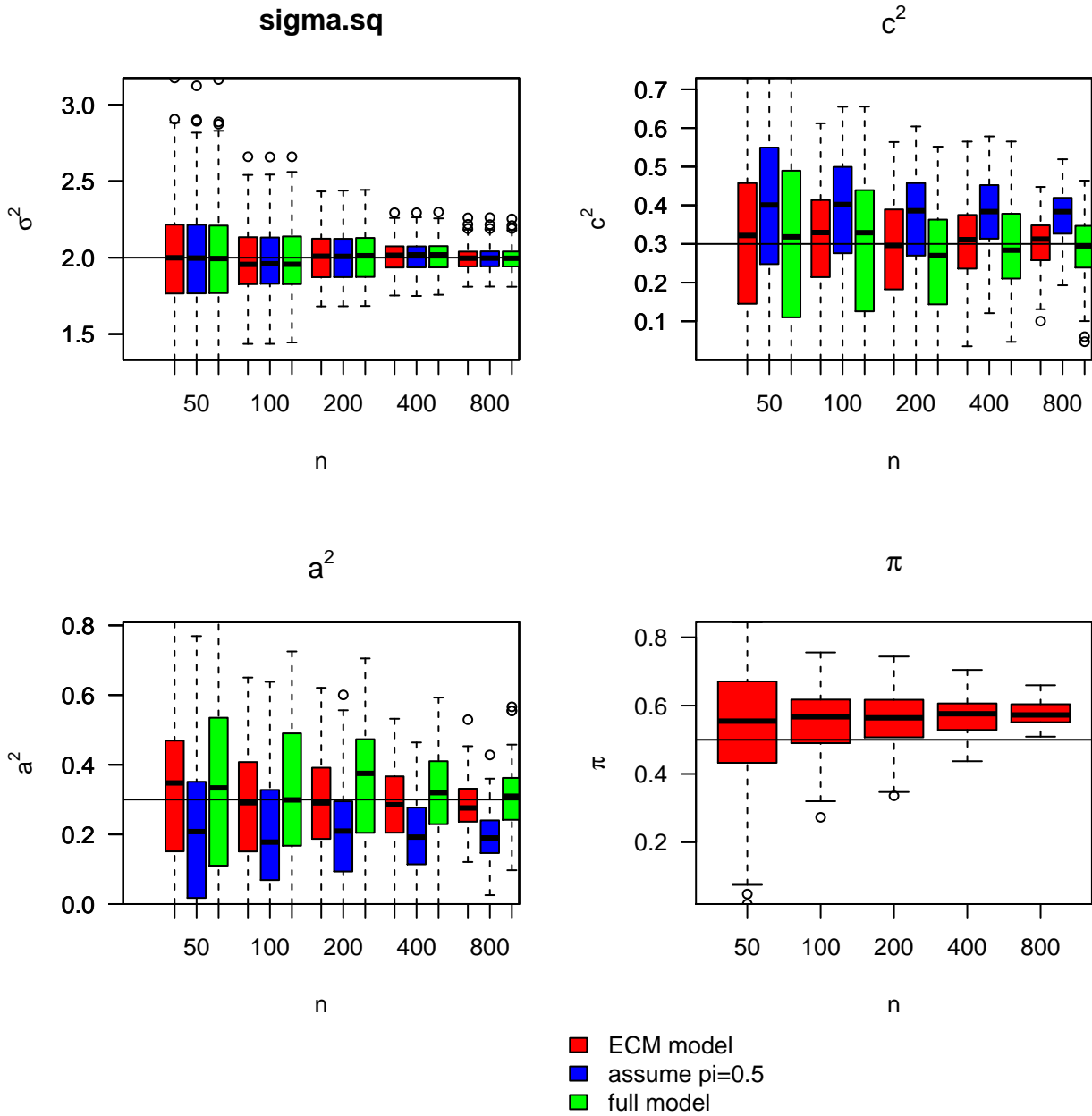
¹לתרשימי BOXPLOT של μ ראה נספח 8.4

איור 1:



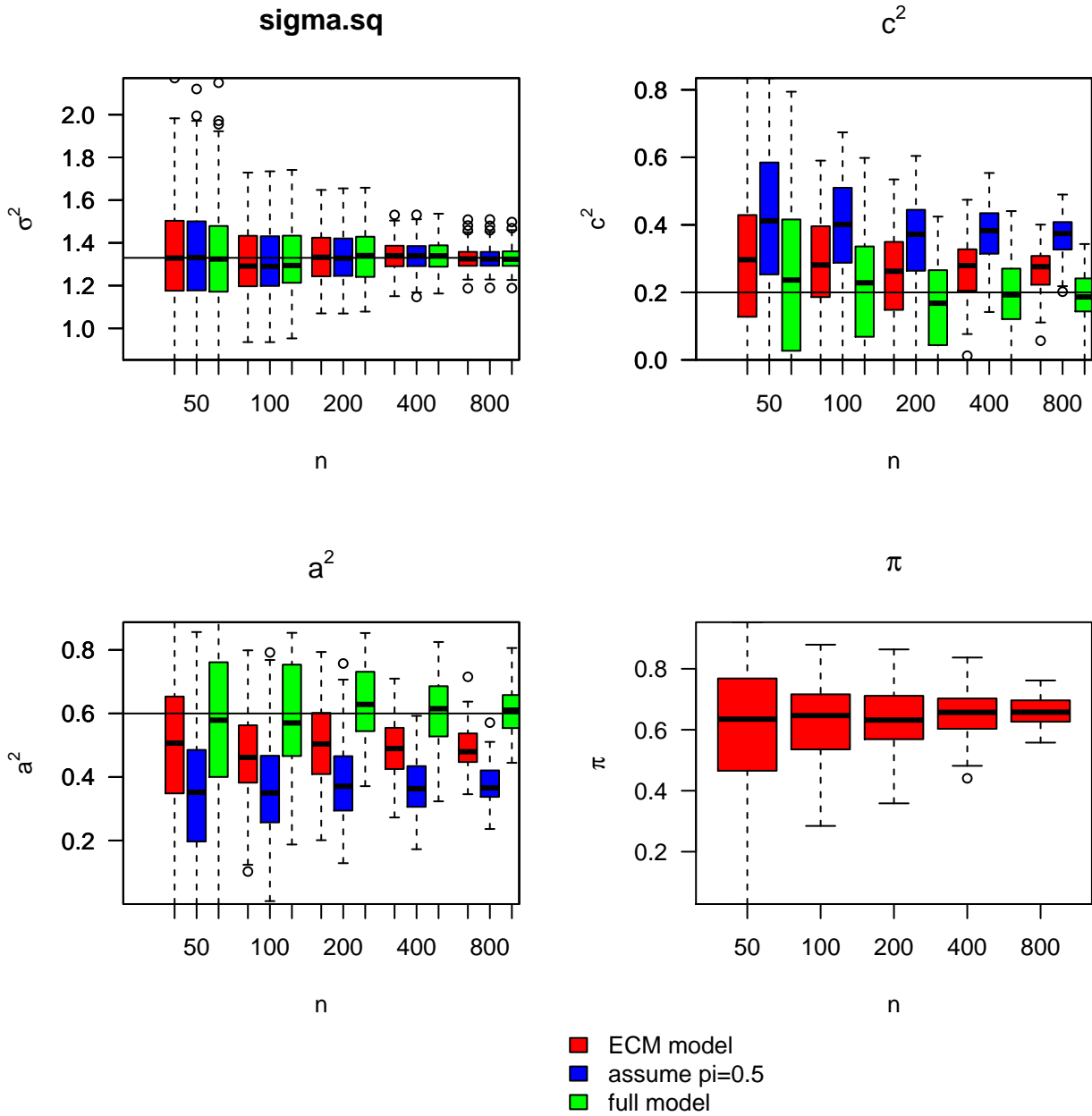
הפרמטרים שנבחרו הם $\mu=5$, $\sigma^2 = 2$, $c^2 = 0.3$, $a^2 = 0.4$. הקופסא האדומה השמאלית מתארת את השיטה הראשונה שבה הפעלנו את המודל החסר ואמדנו את π . הקופסא הכחולה האמצעית היא השיטה השנייה, והקופסא הירוקה הימנית היא השיטה השלישית שבה לקחנו את המודל המלא להשוואה.

איור 2:



הפרמטרים שנבחרו הם $a^2 = 0.3, c^2 = 0.3, \sigma^2 = 2, \mu=5$. הקופסא האדומה השמאלית מתארת את השיטה הראשונה שבה הפעלנו את המודל החסר ואמדנו את π . הקופסא הכחולה האמצעית היא השיטה השנייה, והקופסא הירוקה הימנית היא השיטה השלישית שבה לקחנו את המודל המלא להשוואה.

איור 3:



הפרמטרים שנבחרו הם $a^2 = 0.6$, $c^2 = 0.2$, $\sigma^2 = 2$, $\mu = 5$.
 הקופסא האדומה השמאלית מתארת את השיטה הראשונה שבה הפעלנו את המודל החסר ואמדנו את π .
 הקופסא הכחולה האמצעית היא השיטה השנייה, והקופסא הירוקה הימנית היא השיטה השלישית שבה לקחנו את המודל המלא להשוואה.

בתרחיש השני יצרנו מדגם אקראי של נתונים, כפי שנעשה בתרחיש הראשון. אך כאן שינינו את אופן קביעת התאומים הזהים והלא זהים. בעוד שבתרחיש הראשון לקחנו כחצי מהנתונים להיות תאומים זהים והחצי השני להיות תאומים לא זהים. כאן שינינו את האחוז הזה באופן הבא. בהתחלה הגרלנו אחוז מסוים של תאומים זהים ואחוז של תאומים לא זהים והשאר לא ידוע (NA). עבור אלה הלא ידועים הגרלנו אחוז מסוים שיהיו זהים (MZ) והשאר יהיו לא זהים (DZ). שוב הגרלנו את הנתונים מהתפלגות שבה בחרנו ערכים שונים לפרמטרים: μ, a^2, c^2 ו- σ^2 . ולקחנו ערכים שונים של n :

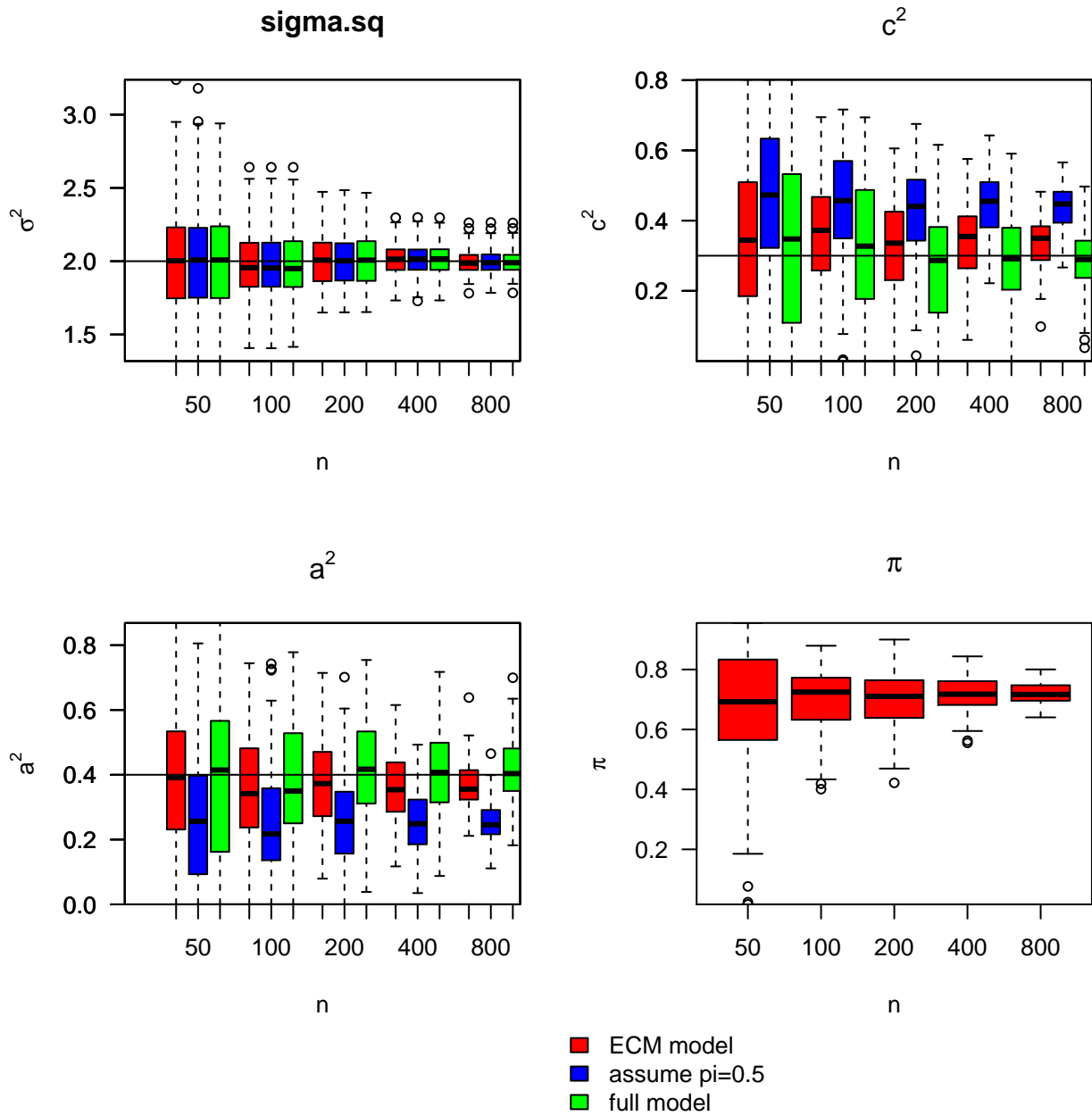
$n = 50, 100, 200, 400, 800$. עבור כל n הרצנו את הסימולציה 100 פעמים עבור שלושה שיטות שנזכרו בתרחיש 1. ניתן לראות מהתוצאות כי האומד של σ^2 הוא קרוב וכמעט זהה ככל שמגדילים את גודל המדגם. האומדים לפרמטרים a^2 ו- c^2 פחות קרובים אבל עדיין השיטה הראשונה שבו לקחנו את הנראות עבור המידע החסר (הקופסא האדומה השמאלית) טובה יותר מהשיטה השנייה שבה מניחים כי המשקולות ידועות מראש (הקופסא הכחולה האמצעית).

באיור 4 אפשר לראות את תרשימי ה-BOXPLOTS של תוצאות הנתונים שהוגרלו באופן הבא: 30% תאומים זהים ו-30% תאומים לא זהים והשאר לא ידוע (NA). עבור אלה הלא ידועים הנחנו כי 90% הם זהים (MZ) ו-10% הם לא זהים (DZ). עבור הפרמטרים: $a^2 = 0.4, c^2 = 0.3, \sigma^2 = 2, \mu = 5$. באיור 5 אפשר לראות את תרשימי ה-BOXPLOTS של תוצאות הנתונים שהוגרלו באופן הבא: 30% תאומים זהים ו-30% תאומים לא זהים והשאר לא ידוע (NA). עבור אלה הלא ידועים הנחנו כעת כי 10% הם זהים (MZ) ו-90% הם לא זהים (DZ). עבור הפרמטרים:

$$a^2 = 0.3, c^2 = 0.3, \sigma^2 = 2, \mu = 5$$

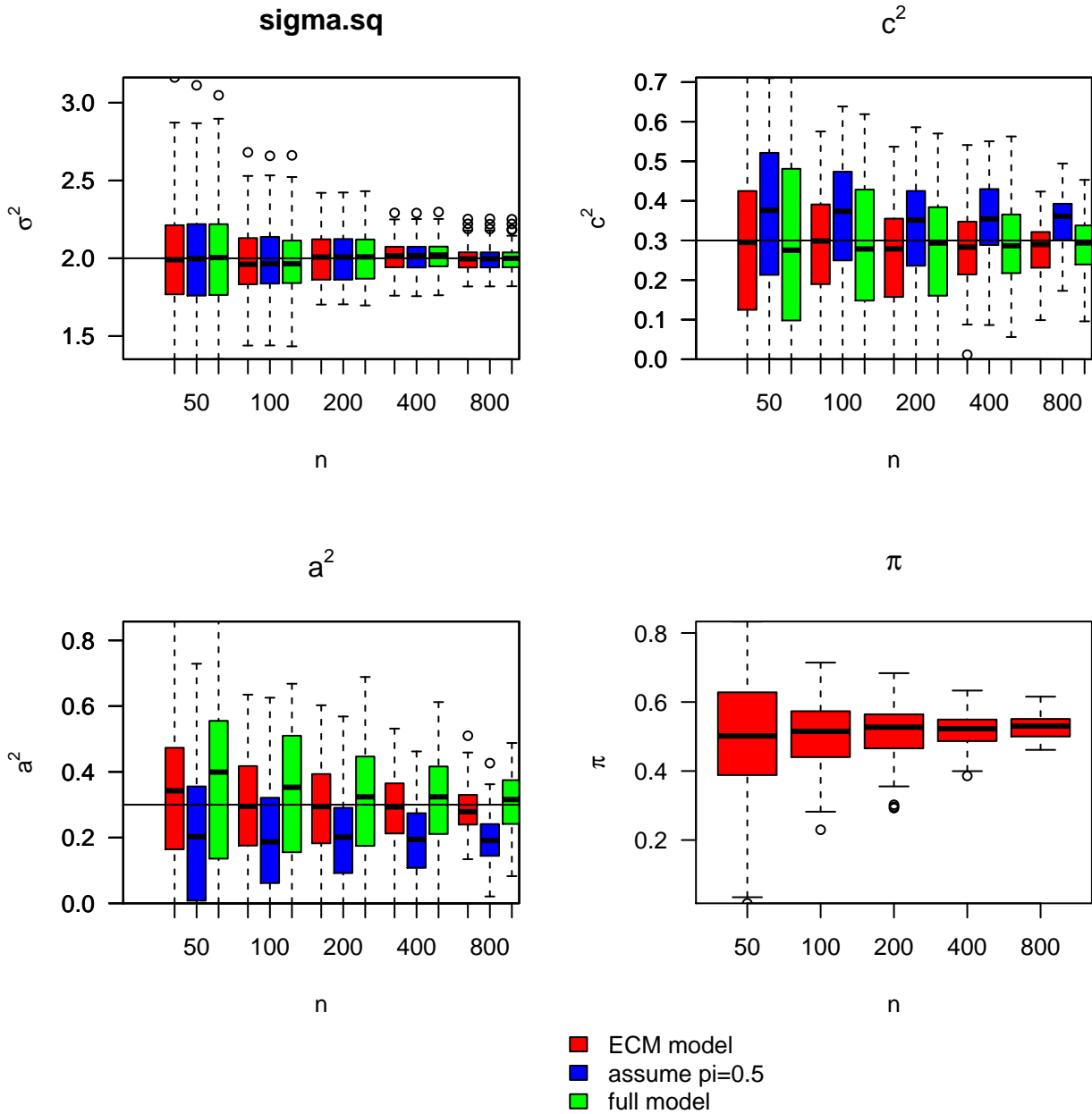
באיור 6 מוצגים תרשימי ה-BOXPLOTS של תוצאות הנתונים שהוגרלו באופן הבא: 40% תאומים זהים ו-30% תאומים לא זהים והשאר לא ידוע (NA). עבור אלה הלא ידועים הנחנו כעת כי 15% הם זהים (MZ) ו-85% הם לא זהים (DZ). עבור הפרמטרים: $a^2 = 0.6, c^2 = 0.2, \sigma^2 = 1.33, \mu = 5$.

איור 4:



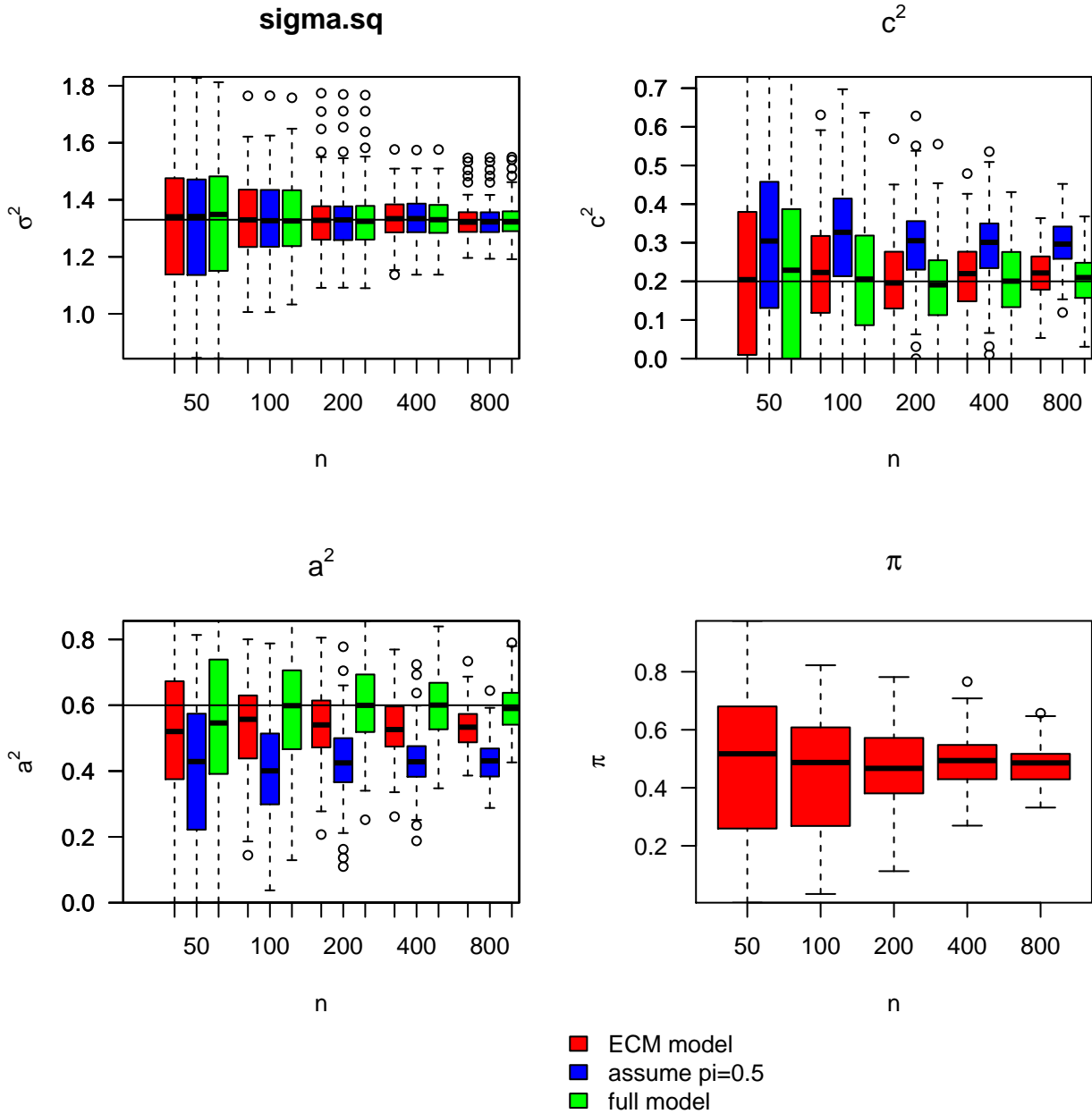
הפרמטרים שנבחרו הם $a^2 = 0.4$, $c^2 = 0.3$, $\sigma^2 = 2$, $\mu = 5$. הקופסא האדומה השמאלית מתארת את השיטה הראשונה שבה הפעלנו את המודל החסר ואמדנו את π . הקופסא הכחולה האמצעית היא השיטה השנייה, והקופסא הירוקה הימנית היא השיטה השלישית שבה לקחנו את המודל המלא להשוואה.

איור 5:



הפרמטרים שנבחרו הם $\mu = 5$, $\sigma^2 = 2$, $c^2 = 0.3$, $a^2 = 0.3$. הקופסא האדומה השמאלית מתארת את השיטה הראשונה שבה הפעלנו את המודל החסר ואמדנו את π . הקופסא הכחולה האמצעית היא השיטה השנייה, והקופסא הירוקה הימנית היא השיטה השלישית שבה לקחנו את המודל המלא להשוואה.

איור 6:



הפרמטרים שנבחרו הם $\mu = 5$, $\sigma^2 = 1.33$, $c^2 = 0.2$, $a^2 = 0.6$. הקופסא האדומה השמאלית מתארת את השיטה הראשונה שבה הפעלנו את המודל החסר ואמדנו את π . הקופסא הכחולה האמצעית היא השיטה השנייה, והקופסא הימנית שבה לקחנו את המודל המלא להשוואה.

5.2 תוצאות הסימולציות

על פי תרשימי ה-BOXPLOTS, התוצאות שהתקבלו מהרצת הסימולציות מעידות שהשיטה, שהוצעה בתת-פרק 4.2, למציאת אומדי הנראות המירבית עבור מודל עם מידע חסר תוך השימוש במודל של תערובת התפלגויות עבור נתוני תאומים עם מידע חסר, סיפקה תוצאות יותר טובות מהמצב שבו מניחים משקולות ידועות מראש, למשל בחירת $\pi_1 = 0.5$ לכל התצפיות.

בתרחיש הראשון, בכל המצבים שהרצנו, קיבלנו שהאומדים של μ ו- σ^2 הם די קרובים ככל שמגדילים את גודל המדגם. לעומת זאת האומדים לפרמטרים a^2 ו- c^2 פחות קרובים אבל עדיין השיטה הראשונה טובה יותר מאשר השיטה השנייה שבה לקחנו משקולות נתונות. אפשר גם לראות שכאשר בחרנו ערכים שונים עבור הפרמטרים a^2 ו- c^2 . ככל שההבדל בין a^2 ל- c^2 היה קטן יותר, התוצאות היו יותר קרובות בשיטה שהצענו, לערכים האמיתיים מהתוצאות של המודל "הקלאסי", שהוא המודל שבו מניחים את המשקולות כידועות מראש.

בתרחיש השני שבו שינינו את המשקולות ההתחלתיות, אנו רואים שוב כי האומדים של התוחלת μ ו- σ^2 מתקרבים ומתכנסים לערכים האמיתיים ככל שמגדילים את n . בשלושת השיטות. האומדים של הפרמטרים a^2 ו- c^2 הם טובים יותר מאשר עבור השיטה השנייה של משקולות קבועות מראש. אנו יכולים לזהות כי האומדים של π_1 הם פחות טובים: בתרחיש הראשון הערך האמיתי של π_1 הוא חצי, אך על פי התוצאות שקיבלנו האומדים לא כל כך תואמים לערך האמיתי. גם בתרחיש השני, באיור 6, אפשר לראות כי האומדים של π_1 אינם במגמת התכנסות. אפשרות לבדיקת נקודה זו יכולה להיות בהגדלת מספר התצפיות, ובדיקת הסימולציה מחדש. למרות שבחלק מההרצות יצא כי האומדים של המודל המוצע בעבודה זו לא ממש תואמים לערכים האמיתיים, עדיין המודל הזה מספק תוצאות קרובות יותר לערכים האמיתיים משיטה המקובלת של קביעת משקולות מראש.

בעבודה זו הצגנו מודל של תערובת התפלגויות, שמתאים לנתוני תאומים עם מידע סיווג חסר, ומתחשב באילוצים בין הפרמטרים. כמו כן הצענו דרך לאמידת הפרמטרים של המודל בשיטת הנראות המירבית.

בפרק 2 הוצג מודל של ACE, מודל שהוא שימושי לחקירת שונות בין בני אדם באופן כללי. על פי מודל ACE לנתוני תאומים, ניתן היה לראות כי ישנו קשר בין השונות של תאומים זהים ולא זהים, וכמו כן יש קשר בין התוחלות. הקשרים האלה נחשבים לאילוצים שיש לקחת אותם בחשבון, בעת חישוב ומציאת האומדים. המודל שהתאמנו לנתוני התאומים הוא מודל תערובת GMM. לאחר התאמת המודל, בפרק 3, הוצגה שיטת הנראות המירבית לאמידת הפרמטרים, עם המידע המלא והמידע החסר. ניתן היה לראות כי פונקציית הנראות הינה פונקציה לא לינארית בפרמטרים, שמצריכה שיטות פתרון לא טריוויאליות למציאת האומדים. בפרק 4 הוצגו אלגוריתם EM והרחבה שלו, אלגוריתם ECM, לאמידת הפרמטרים, תוך התחשבות בשני דברים: האילוצים, והמידע החסר. לאחר הצגת האלגוריתמים אמדנו את הפרמטרים. לתוחלת ולמשקולות התקבלה נוסחה סגורה, אך למרכיבי השונות לא התאפשר לקבל נוסחה סגורה, לכן היה צורך לשימוש בשיטות נומריות על מנת לאמוד את מרכיבי השונות. בפרק 5 בדקנו את המודל שהצגנו על ידי סימולציות. ולאחר מכן דנו בתוצאות שהתקבלו. הצלחנו להראות, בעזרת הסימולציות שבוצעו בתוכנת R, כי המודל שהצגנו הוא מודל טוב יותר מאשר המודל הקלאסי, המודל שבו מניחים את המשקולות כידועות מראש.

חשוב לציין כי בעבודה הזו, הנתונים הסתמכו רק על נתוני תאומים והתאמת מודל תערובת שמתייחס לנתונים אלה ספציפית, תוך טיפול בשתי הבעיות: האילוצים בין הפרמטרים, והמידע החסר עבור סוג של כל זוג תאומים. כעבודה עתידית אפשר לקחת את המודל שהותאם כאן בכיוון נוסף, למשל להתאים מודל לינארי מוכלל (GLM), לאמידת כל מיני הסתברויות וסיכויים. למשל, אמידת סיכויים להתפתחות מחלה מסוימת על סמך ידע נתון על אחד מהתאומים. בנוסף אפשר להרחיב את מה שנעשה בעבודה על ידי בחירת נתונים עם קשרים משפחתיים כלליים יותר, כלומר קשרים משפחתיים כדוגמת: אב ובנו, אחים ואחיות, בני דודים וכו'.

References

- [1] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, 39(1):1–38, 1977.
- [2] Mark M Fredrickson. Ace in the hole: Correlation and classical twin studies. Unpublished manuscript, 2009.
- [3] A. Genz, F. Bretz, T. Miwa, X. Mi, F. Leisch, F. Scheipl, and T. Hothorn. mvtnorm: Multivariate normal and t distributions. 2012. URL <http://CRAN.R-project.org/package=mvtnorm>. R package version 0.9-9992.
- [4] D. R. Hunter and K. Lange. A tutorial on MM algorithms. *The American Statistician*, 58(1):30–37, 2004.
- [5] R. J. A. Little and D. B. Rubin. *Statistical Analysis with Missing Data (Wiley Series in Probability and Statistics)*. Wiley, second edition, 2002.
- [6] G. J. MacLachlan and T. Krishnan. *The EM algorithm and extensions*. Wiley, 1997.
- [7] J. J. McArdle and C. A. Prescott. Mixed-effects variance components models for biometric family analyses. *Behavior genetics*, 35(5):631–652, 2005.
- [8] X. L. Meng. On the rate of convergence of the ECM algorithm. *The Annals of Statistics*, 22(1):326–339, 1994.
- [9] X. L. Meng and D. B. Rubin. Maximum likelihood estimation via the ECM algorithm: A general framework. *Biometrika*, 80(2):267–278, 1993.
- [10] M. C. Neale. A finite mixture distribution model for data collected from twins. *Twin Research and Human Genetics*, 6(03):235–239, 2003.
- [11] M. C. Neale and L. Cardon. *Methodology for Genetic Studies of Twins and Families*. Springer, 1992.
- [12] F. Vaida. Parameter convergence for EM and MM algorithms. *Statistica Sinica*, 15(3):831, 2005.
- [13] C. F. J. Wu. On the convergence properties of the EM algorithm. *The Annals of statistics*, 11(1):95–103, 1983.

8.1 אומדי נראות מירבית עבור המודל המלא

נחשב את אומדי הנראות המירבית עבור המודל המלא. כפי שראינו בתת פרק 3.2 פונקציית לוג הנראות עבור המודל המלא היא:

$$l(\vec{x}; \{\mu, a^2, c^2, \sigma^2\}) = \sum_{i=1}^k \log(f_1(x_i; \Sigma_{MZ}, \mu)) + \sum_{i=k+1}^n \log(f_2(x_i; \Sigma_{DZ}, \mu))$$

כלומר:

$$l(\vec{x}; \{\mu, a^2, c^2, \sigma^2\}) = \sum_{i=1}^k \left(-\frac{1}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_{MZ}| - \frac{1}{2} (x_i - \mu)^t \cdot \Sigma_{MZ}^{-1} \cdot (x_i - \mu) \right) \quad (5)$$

$$+ \sum_{i=k+1}^n \left(-\frac{1}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_{DZ}| - \frac{1}{2} (x_i - \mu)^t \cdot \Sigma_{DZ}^{-1} \cdot (x_i - \mu) \right)$$

נסמן: $\Sigma_1 = \Sigma_{MZ}$ ו- $\Sigma_2 = \Sigma_{DZ}$
 למציאת אומד נראות מירבית עבור μ , נגזור את משוואה (5) לפי μ .
 ונקבל:

$$\frac{\partial l(\vec{x}; \{\mu, a^2, c^2, \sigma^2\})}{\partial \mu} = \sum_{i=1}^k (\Sigma_1^{-1} (x_i - \mu)) + \sum_{i=k+1}^n (\Sigma_2^{-1} \cdot (x_i - \mu))$$

$$= \sum_{i=1}^k \left(\Sigma_1^{-1} \begin{pmatrix} x_{i1} - \mu \\ x_{i2} - \mu \end{pmatrix} \right) + \sum_{i=k+1}^n \left(\Sigma_2^{-1} \cdot \begin{pmatrix} x_{i1} - \mu \\ x_{i2} - \mu \end{pmatrix} \right)$$

$$= \sum_{i=1}^k \left(\Sigma_1^{-1} \begin{pmatrix} x_{i1} \\ x_{i2} \end{pmatrix} \right) - k \Sigma_1^{-1} \begin{pmatrix} 1 \\ 1 \end{pmatrix} \mu + \sum_{i=k+1}^n \left(\Sigma_2^{-1} \begin{bmatrix} x_{i1} \\ x_{i2} \end{bmatrix} \right)$$

$$- (n - k) \Sigma_2^{-1} \begin{pmatrix} 1 \\ 1 \end{pmatrix} \mu$$

נשווה לאפס למציאת μ :

$$\frac{\partial l(\vec{x}; \{\mu, a^2, c^2, \sigma^2\})}{\partial \mu} = 0$$

$$k \Sigma_1^{-1} \begin{pmatrix} 1 \\ 1 \end{pmatrix} \mu + (n - k) \Sigma_2^{-1} \begin{pmatrix} 1 \\ 1 \end{pmatrix} \mu = \sum_{i=1}^k \left(\Sigma_1^{-1} \begin{pmatrix} x_{i1} \\ x_{i2} \end{pmatrix} \right) + \sum_{i=k+1}^n \left(\Sigma_2^{-1} \begin{pmatrix} x_{i1} \\ x_{i2} \end{pmatrix} \right)$$

נחלץ את μ :

$$\hat{\mu} = \left(\left[(n - k) \Sigma_2^{-1} \begin{pmatrix} 1 \\ 1 \end{pmatrix} + k \Sigma_1^{-1} \begin{pmatrix} 1 \\ 1 \end{pmatrix} \right]^{-1} \right)^t \left(\sum_{i=1}^k \left(\Sigma_1^{-1} \begin{pmatrix} x_{i1} \\ x_{i2} \end{pmatrix} \right) + \sum_{i=k+1}^n \left(\Sigma_2^{-1} \begin{pmatrix} x_{i1} \\ x_{i2} \end{pmatrix} \right) \right)$$

עכשיו נסתכל בשאר הפרמטרים. קודם נסדר את הפונקציה (5).
 על פי מודל ACE:

$$\Sigma_2 = \Sigma_{DZ} = \sigma^2 \begin{pmatrix} 1 & 0.5a^2 + c^2 \\ 0.5a^2 + c^2 & 1 \end{pmatrix} \quad \Sigma_1 = \Sigma_{MZ} = \sigma^2 \begin{pmatrix} 1 & a^2 + c^2 \\ a^2 + c^2 & 1 \end{pmatrix}$$

כל תצפית: $x_i = (x_{i1}, x_{i2})$
 נסמן

$$\begin{aligned} \rho_1 &= a^2 + c^2 \\ \rho_2 &= 0.5a^2 + c^2 \end{aligned}$$

נציב ב- (5) ונרשום את החלק התלוי בפרמטרים $\tau = \{\rho_1, \rho_2, \sigma^2\}$

$$\begin{aligned} \zeta(\tau) &= \sum_{i=1}^k \left(-\frac{1}{2} \log(\sigma^4(1 - \rho_1^2)) - \frac{1}{2} \cdot \frac{1}{\sigma^4(1 - \rho_1^2)} \left[(x_{i1} - \mu)^2 - 2\rho_1 \cdot (x_{i1} - \mu) \cdot (x_{i2} - \mu) + (x_{i2} - \mu)^2 \right] \right) \\ &+ \sum_{i=k+1}^n \left(-\frac{1}{2} \log(\sigma^4(1 - \rho_2^2)) - \frac{1}{2} \cdot \frac{1}{\sigma^4(1 - \rho_2^2)} \left[(x_{i1} - \mu)^2 - 2\rho_2 \cdot (x_{i1} - \mu) \cdot (x_{i2} - \mu) + (x_{i2} - \mu)^2 \right] \right) \end{aligned} \quad (6)$$

נסמן: $v = \frac{1}{\sigma^4}$
 נגזור את (6) לפי v :

$$\begin{aligned} \frac{\partial \zeta(\{\mu, \rho_1, \rho_2, \sigma^2\})}{\partial v} &= \sum_{i=1}^k \left(\frac{1}{2} \cdot \frac{1}{v} - \frac{1}{2} \cdot \frac{1}{(1 - \rho_1^2)} \left[(x_{i1} - \mu)^2 - 2\rho_1 \cdot (x_{i1} - \mu) \cdot (x_{i2} - \mu) + (x_{i1} - \mu)^2 \right] \right) \\ &+ \sum_{i=k+1}^n \left(\frac{1}{2} \cdot \frac{1}{v} - \frac{1}{2} \cdot \frac{1}{(1 - \rho_2^2)} \left[(x_{i1} - \mu)^2 - 2\rho_2 \cdot (x_{i1} - \mu) \cdot (x_{i2} - \mu) + (x_{i2} - \mu)^2 \right] \right) \end{aligned}$$

נשווה לאפס ונקבל:

$$\begin{aligned} 0 &= \frac{k}{v} - \sum_{i=1}^k \frac{1}{2(1 - \rho_1^2)} \left[(x_{i1} - \mu)^2 - 2\rho_1 \cdot (x_{i1} - \mu) \cdot (x_{i2} - \mu) + (x_{i1} - \mu)^2 \right] \\ &+ \frac{n - k}{v} - \sum_{i=k+1}^n \frac{1}{2(1 - \rho_2^2)} \left[(x_{i1} - \mu)^2 - 2\rho_2 \cdot (x_{i1} - \mu) \cdot (x_{i2} - \mu) + (x_{i2} - \mu)^2 \right] \end{aligned}$$

נסדר:

$$\begin{aligned} \frac{n}{v} &= \sum_{i=1}^k \frac{1}{2(1 - \rho_1^2)} \left[(x_{i1} - \mu)^2 - 2\rho_1 \cdot (x_{i1} - \mu) \cdot (x_{i2} - \mu) + (x_{i1} - \mu)^2 \right] \\ &+ \sum_{i=k+1}^n \frac{1}{2(1 - \rho_2^2)} \left[(x_{i1} - \mu)^2 - 2\rho_2 \cdot (x_{i1} - \mu) \cdot (x_{i2} - \mu) + (x_{i2} - \mu)^2 \right] \end{aligned}$$

לכן:

$$\frac{1}{v} = \frac{1}{n} \left(\sum_{i=1}^k \frac{1}{2(1-\rho_1^2)} [(x_{i1} - \mu)^2 - 2\rho_1 \cdot (x_{i1} - \mu) \cdot (x_{i2} - \mu) + (x_{i1} - \mu)^2] \right. \\ \left. + \sum_{i=k+1}^n \frac{1}{2(1-\rho_2^2)} [(x_{i1} - \mu)^2 - 2\rho_2 \cdot (x_{i1} - \mu) \cdot (x_{i2} - \mu) + (x_{i2} - \mu)^2] \right)$$

קיבלנו:

$$\hat{\sigma}^4 = \frac{1}{n} \left(\sum_{i=1}^k \frac{1}{2(1-\rho_1^2)} [(x_{i1} - \mu)^2 - 2\rho_1 \cdot (x_{i1} - \mu) \cdot (x_{i2} - \mu) + (x_{i1} - \mu)^2] \right. \\ \left. + \sum_{i=k+1}^n \frac{1}{2(1-\rho_2^2)} [(x_{i1} - \mu)^2 - 2\rho_2 \cdot (x_{i1} - \mu) \cdot (x_{i2} - \mu) + (x_{i2} - \mu)^2] \right)$$

נסכים:

$$\hat{\sigma}^2 = \left[\frac{1}{n} \left(\sum_{i=1}^k \frac{1}{2(1-\rho_1^2)} [(x_{i1} - \mu)^2 - 2\rho_1 \cdot (x_{i1} - \mu) \cdot (x_{i2} - \mu) + (x_{i1} - \mu)^2] \right. \right. \\ \left. \left. + \sum_{i=k+1}^n \frac{1}{2(1-\rho_2^2)} [(x_{i1} - \mu)^2 - 2\rho_2 \cdot (x_{i1} - \mu) \cdot (x_{i2} - \mu) + (x_{i2} - \mu)^2] \right) \right]^{\frac{1}{2}}$$

נשאר לנו לאמוד את הפרמטרים ρ_1, ρ_2 . לפרמטרים האלה לא נוכל לחשב את האומדים בצורה אנליטית כי הפונקציה (6) מסובכת בפרמטרים אלה כך שלא נוכל לקבל נוסחה סגורה. השיטה החילופית היא השיטה הנומרית. נוכל לאמוד את הפרמטרים בעזרת שיטות נומריות כגון ניוטון רפסון או אחרות.

8.2 הרחבה - אלגוריתם EM

אלגוריתם EM מסתמך על הקשר (1). בהינתן קבוצת פרמטרים $\theta^{(m)}$, ועל ידי הכפלת המשוואה (1) ב- $P(Y|\theta^{(m)})$ נקבל:

$$\log [P(X|\theta)] \cdot P(Y|\theta^{(m)}) = \log [P((X, Y)|\theta)] \cdot P(Y|\theta^{(m)}) - \log [P(Y|X, \theta)] \cdot P(Y|\theta^{(m)})$$

נניח Y בדיד נסכום על כל האפשרויות לערכי Y ונקבל:

$$\log [P(X|\theta)] = \sum_y \log [P((X, Y)|\theta)] \cdot P(Y|\theta^{(m)}) - \sum_y \log [P(Y|X, \theta)] \cdot P(Y|\theta^{(m)})$$

נסמן:

$$Q(\theta|\theta^{(m)}) = \sum_y \log [P((X, Y)|\theta)] \cdot P(Y|\theta^{(m)})$$

$$H(\theta|\theta^{(m)}) = -\sum_y \log [P(Y|X, \theta)] \cdot P(Y|\theta^{(m)})$$

קיבלנו:

$$\log [P(x|\theta)] = Q(\theta|\theta^{(m)}) + H(\theta|\theta^{(m)})$$

במקום למקסם את פונקציית המטרה $\log [P(X|\theta)]$ אנו ממקסימים את הפונקציה $Q(\theta|\theta^{(m)})$. על פי אי שוויון ינסן אפשר לקבל כי:

$$H(\theta^{(m)}|\theta^{(m)}) \leq H(\theta|\theta^{(m)})$$

כלומר:

$$Q(\theta|\theta^{(m)}) + H(\theta^{(m)}|\theta^{(m)}) \leq Q(\theta|\theta^{(m)}) + H(\theta|\theta^{(m)})$$

ולכן הפונקציה הנוצרת בצעד ה- E , $Q(\theta|\theta^{(m)})$, עד הקבוע: $c = H(\theta^{(m)}|\theta^{(m)})$ נחשבת לפונקציה חילופית שמחפשים את המקסימום של פונקציה זו במקום הפונקציה המקורית.

כלומר אלגוריתם EM משתמש ברעיון של שיטת כללית יותר שנקראת שיטת MM (Minorization - Maximization), נגדיר את השיטה הזו:

MM הינה שיטה שבעזרתה בונים אלגוריתם לפתרון בעיות אופטימיזציה באופן נומרי (היא נקראת גם "אלגוריתם MM" [4]). שיטת MM מציעה פונקציה "תחליפית" לפונקציה "מסובכת" שנרצה לעשות לה אופטימיזציה. כלומר היא הופכת את הבעיה לצורה "פשוטה" יותר. במקרה שלנו אנו מתעניינים בחיפוש אומד נראות מירבית כלומר הבעיה היא בעיית מקסימיזציה, נגדיר את MM בבעיית מקסימיזציה:

בבעיות מקסימיזציה: MM = Minorization - Maximization

ההגדרה של פעולת Minorization

הגדרה 8.1 A function $g(\theta|\theta^{(m)})$ is said to minorize a real-valued function $f(\theta)$ at point $\theta^{(m)}$ provided:

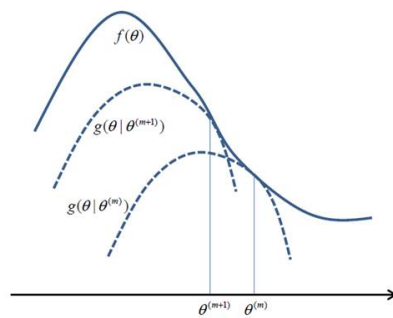
$$\forall \theta \quad g(\theta|\theta^{(m)}) \leq f(\theta), \quad g(\theta^{(m)}|\theta^{(m)}) = f(\theta^{(m)})$$

אלגוריתם MM מתבצע בשני צעדים:

1. מציאת את הפונקציה $g(\theta|\theta^{(m)})$ "המנורייר" "

$$2. \quad \theta^{(m+1)} = \underset{\theta}{\operatorname{argmax}} g(\theta|\theta^{(m)})$$

ממשיכים כך עד להתכנסות להרחבה על התכנסות של אלגוריתם MM ראה [12].
בשרטוט



אנו רואים ששני הצעדים של אלגוריתם MM מקבילים לשני הצעדים של אלגוריתם EM. במקרה של EM הפונקציה $Q(\theta|\theta^{(m)})$ היא המינורייר של הפונקציה $\log [P(x|\theta)] - H(\theta|\theta^{(m)})$. ולכן ניתן להסתכל על כל אלגוריתם EM כמקרה פרטי של אלגוריתם כללי יותר MM.

8.3 אמידה כללית לווקטור המשקולות π

החישוב שהוצג בפרק 4.2 לאמידת המשקולות π , עובד רק במצב שיש שתי אפשרויות לתערובות כלומר רק למצב הספיציפי שלנו. אפשר להסתכל על החישוב באופן אחר, שיתאים למקרה הכללי. נציג כאן את החישוב הכללי הזה.

נגזור (4) לפי ווקטור המשקולות $\vec{\pi}$ תוך התחשבות בתנאי $\sum_j \pi_j = 1$.

החישוב הכללי:

על מינת לטפל באילוץ $\sum_j \pi_j = 1$ נוסיף כופל לגרנג λ .

כלומר נוסיף: $\lambda \left(\sum_j \pi_j - 1 \right) = \lambda (\pi_1 + \pi_2 - 1)$ עכשיו נצטרך לגזור את:

$$\xi(\theta, \theta^{(m)}) = Q(\theta, \theta^{(m)}) + \lambda \left(\sum_j \pi_j - 1 \right)$$

$$d \frac{\xi(\theta, \theta^{(m)})}{d\pi} = \sum_{i=1}^n \sum_{j=1}^2 P(z_{ij}|y_i, \theta^{(m)}) \frac{1}{\pi_j} + \lambda$$

נשווה לאפס:

$$0 = \sum_{i=1}^n \sum_{j=1}^2 P(z_{ij}|y_i, \theta^{(m)}) \frac{1}{\pi_j} + \lambda$$

עבור j בודד:

$$0 = \sum_{i=1}^n P(z_{ij}|y_i, \theta^{(m)}) \frac{1}{\pi_j} + \lambda \quad (7)$$

נסכום את המשוואות לכל j , במצב הזה נסכום את שתי המשוואות, על מנת לקבל אומד ל- λ :

$$\begin{aligned} 0 &= \sum_{i=1}^n P(z_{i1}|y_i, \theta^{(m)}) \frac{1}{\pi_1} + \lambda + \sum_{i=1}^n P(z_{i2}|y_i, \theta^{(m)}) \frac{1}{\pi_2} + \lambda \\ 0 &= \frac{1}{\pi_1} \cdot \pi_1 n + \lambda + \frac{1}{\pi_2} \cdot \pi_2 n + \lambda \\ 0 &= 2n + 2\lambda \\ -n &= \lambda \end{aligned}$$

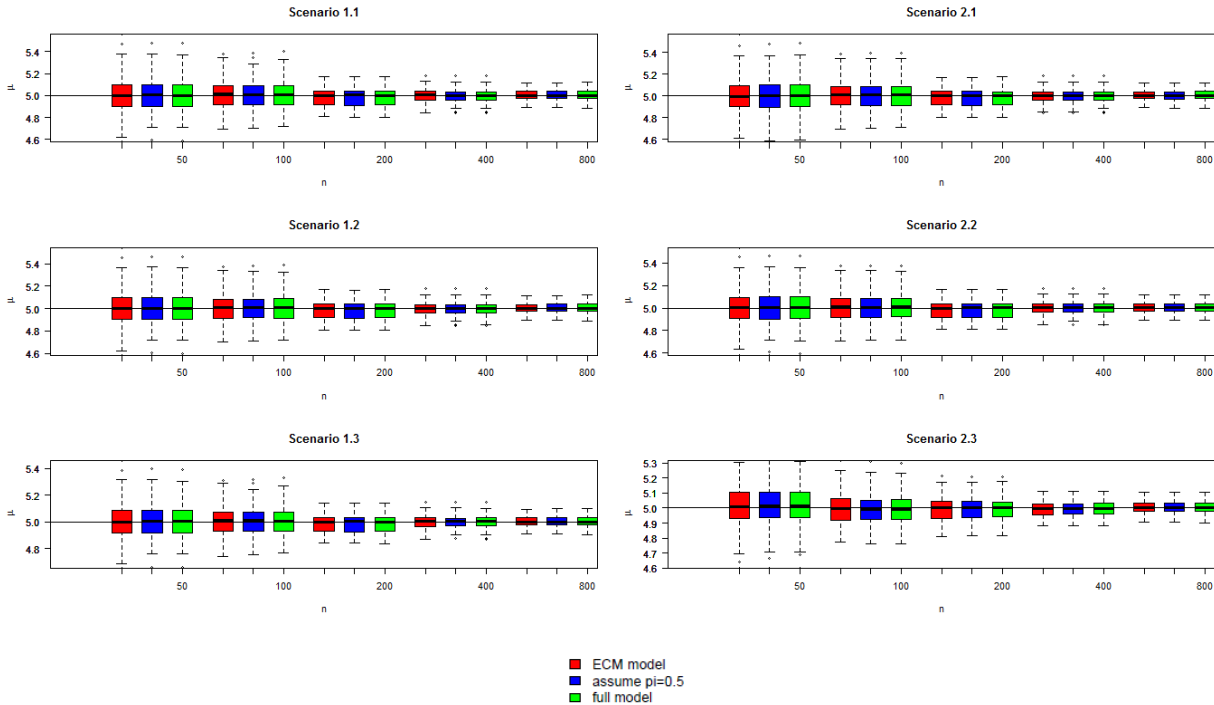
לכן אם נחזור ל-(7) ונציב את מה שקיבלנו במשוואה האחרונה נקבל:

$$\pi_j^{(m+1)} = \frac{1}{n} \sum_{i=1}^n p_{ij}^{(m)}$$

ראינו כי עבור המקרה הספיציפי שלנו קיבלנו את אותה התוצאה.

8.4 תרשימי Boxplot עבור הפרמטר μ

איור 7:



ב-Scenario 1.1 אפשר לראות את תרשימי ה-BOXPLOTS של תוצאות הנתונים שהיו בתרחיש 1 בסימולציה הראשונה שבה הפרמטרים הוגרלו באופן הבא: $a^2 = 0.4, c^2 = 0.3, \sigma^2 = 2, \mu = 5$.
 ב-Scenario 1.2 אפשר לראות את תרשימי ה-BOXPLOTS של תוצאות הנתונים שהוגרלו בתרחיש 1 בסימולציה השנייה שהוצגה, עבור הפרמטרים: $a^2 = 0.3, c^2 = 0.3, \sigma^2 = 2, \mu = 5$.
 ב-Scenario 1.3 מוצגים תרשימי ה-BOXPLOTS של תוצאות הנתונים שהוגרלו בתרחיש 1 בסימולציה השלישית שהוצגה. עבור הפרמטרים: $a^2 = 0.6, c^2 = 0.2, \sigma^2 = 1.33, \mu = 5$.
 ב-Scenario 2.1 אפשר לראות את תרשימי ה-BOXPLOTS של תרחיש 2 עבור תוצאות הנתונים שהוגרלו באופן הבא: 30% תאומים זהים ו-30% תאומים לא זהים והשאר לא ידוע (NA). עבור אלה הלא ידועים הנחנו כי 90% הם זהים (MZ) ו-10% הם לא זהים (DZ). עבור הפרמטרים: $a^2 = 0.4, c^2 = 0.3, \sigma^2 = 2, \mu = 5$.
 ב-Scenario 2.2 אפשר לראות את תרשימי ה-BOXPLOTS של תרחיש 2 עבור תוצאות הנתונים שהוגרלו באופן הבא: 30% תאומים זהים ו-30% תאומים לא זהים והשאר לא ידוע (NA). עבור אלה הלא ידועים הנחנו כעת כי 10% הם זהים (MZ) ו-90% הם לא זהים (DZ). עבור הפרמטרים: $a^2 = 0.6, c^2 = 0.2, \sigma^2 = 1.33, \mu = 5$.
 ב-Scenario 2.3 מוצגים תרשימי ה-BOXPLOTS של תרחיש 2 עבור תוצאות הנתונים שהוגרלו באופן הבא: 40% תאומים זהים ו-30% תאומים לא זהים והשאר לא ידוע (NA). עבור אלה הלא ידועים הנחנו כעת כי 15% הם זהים (MZ) ו-85% הם לא זהים (DZ). עבור הפרמטרים: $a^2 = 0.6, c^2 = 0.2, \sigma^2 = 1.33, \mu = 5$.

Fitting Models for Twins Data with Missing Data

By: Afraa Araidy

Supervised by: Dr Yair Goldberg

THESIS SUBMITTED IN PARTIAL FULFILLMENT OF THE
REQUIREMENTS FOR THE MASTER'S DEGREE

University of Haifa
Faculty of Social Science
Department of Statistics

November 2014

Fitting Models for Twins Data with Missing Data

Afraa Araidy

THESIS SUBMITTED IN PARTIAL FULFILLMENT OF THE
REQUIREMENTS FOR THE MASTER'S DEGREE

University of Haifa
Faculty of Social Science
Department of Statistics

November 2014

Fitting Models for Twins Data with Missing Data

Afraa Araidy

Abstract

Twin studies are an important tool for investigating the influence of heredity and environment on variation in human population. The majority of twin studies assume that the zygosity is given, i.e., the information for each pair of twins whether they are monozygotic (MZ) or dizygotic (DZ) is given. However, this can be missing. Matching models under the wrong assumption that there is no missing data, and estimation of the parameters of these models can lead to biased estimators. Analysis on biased estimators may lead to erroneous conclusions. For this reason we need to find methods to solve this problem.

Some researchers tried to solve the problem of missing data by providing predetermined weights for the probability of twins to be either MZ or DZ.

In this work we propose a model of mixture distributions for twin data analysis without knowing the zygosity for some of the twins. In other words, we consider a model that deals with missing data. In contrast to previous methods which complement the missing information with probability and predetermined weights, we offer a maximum likelihood estimator (MLE). Finding the MLE in this problem is challenging since this maximization include constraints which arise from the structure of the variance matrices. To find the maximum likelihood estimator we use ECM algorithm (Expectation-Conditional Maximization). For implementing of the algorithm we developed an R code using numerical methods. In this work, we ran simulations to examine the proposed method. The results showed that using the method presented here, i.e. the model of mixture distributions with ECM algorithm for finding the maximum likelihood estimator is a model that provides better results than the previous methods which complement the missing data in this context.