

אמידת מספר עותקים מהמודל

כעת נבחר שיש לם אומדן לפונקציה $\hat{f}(seq_k)$ ואת מספרי הפרגמנטים הנצפים $\gamma_k, \dots, \gamma_1$, ואנוחנו רוצים לאמוד את מספר העותקים a_k .

הפונקציה $\hat{f}(g_{c_k})$ כפי שאמדנו בכיתה לא תקיים $\{ \hat{f}(g_{c_k}) \} = 1$ אלא $avg_k \{ \hat{f}(g_{c_k}) \}$ אלא $avg_k \{ f(seq_k) \} = avg_k \{ \lambda_k \}$

ראשית נחליץ את a_k מהמודל:

$$a_k = \frac{\lambda_k}{\hat{f}(g_{c_k})} \cdot \frac{1}{N} \cdot \frac{1}{\eta_k},$$

וכעת נשתמש באומדן ההצבה כאשר נציב γ_k במקום λ_k ואת $\hat{f}(g_c)$ במקום $\hat{f}(g_{c_k}) \cdot N$.

$$\hat{a}_k = \frac{\gamma_k}{\hat{f}(g_{c_k})}.$$

הערות

- ניתן להכפיל בעוד קבוע M^1 שידאג שהחציון של האומדים יהיה 1:
- ייתכן שיש עוד גורמים שתלויים ברצף הגנומי שמשפיעים על λ_k , כמו לדוגמה הסיכוי למפות פרמגנט, ואולי דרכי מדידה יותר עדינות ממספר הג' בתא. כדי לסמן מצב כזה, ניתן להשתמש ב $f(seq_k)$ בתיאור המודל במקום $f(g_{c_k})$. ככל שנכניס יותר גורמים ל $f(seq_k)$, כך המשמעות של η_k תשתנה ובתקווה השונות שלו תקטן.

M^1 מגולם גם את $1/N$ וגם את הממוצע של $1/\hat{f}(seq_k)$

מעבדה לסטטיסטיקה: אמידת מספר עותקים מדגימה בודדת (א)

אנחנו רוצים להשתמש בפונקציית ה GC שלמדנו ובערכים הנצפים γ_k כדי לאמוד את מספר-העותקים (Copy number) של הגנום.

בסיכום הזה נציג שיטת אמידה שדורשת רק אמידה של פונקציית ההשפעה של GC, לעיתים, קוראים לשיטות אמידה אלו "יתיקון", משום שהן מתקנות את ההטיה הנוצרת מאפקט הGC.

מודל המתאר קשר בין מספר עותקים למספר פרגמנטים

נשתמש ב $k=1...K$ לתאר את מספר התא, וב γ_k לתאר את הכיסוי הנצפה בתא.

נסמן את התוחלת של התא k ב λ_k

$$\lambda_k = E[\gamma_k].$$

המודל עבור התוחלת צריך לסמן את הרכיבים שמרכיבים את התוחלת, ולאפשר לנו לתאר איך היינו מייצרים נתונים כאלו. נשים לב ל3 רכיבים שאנחנו יודעים שמשפיעים על הכיסוי:

- כמות העותקים הוא מספר. בכפולות שלמות של חצי $(2, 1 \frac{1}{2}, 1, \frac{1}{2}, 0, \dots)$. תוחלת הכיסוי צריכה להיות פרופורציונלית למספר העותקים.
- סה"כ כמות הפרגמנטים (ע"פ כל התאים) יכולה להשתנות מניסוי לניסוי, ואינה קשורה למספר העותקים אלא למכונות הריצוף. נסמנה ב J .
- לכמות הGC בתא (או גזרות כמו כמות הGC בתתי התאים) יש השפעה על הכיסוי. תיאורנו תלוי וז בעזרת מודלים של רג'ס'יה.

אם כן, נוז לתאר את התוחלת כפונקציה כפליית של שלושה גורמים בלתי תלויים.

כפול גורם לא ידוע η_k :

$$\lambda_k = N \cdot a_k \cdot f(g_{c_k}) \cdot \eta_k.$$

לאורך זיהוי הפרמטרים צריך להניח הנחות לגבי הממוצעים של $\eta_k, f(g_{c_k})$ ושל a_k . נבחר לדרוש מהחציון של a_k להיות 1, והממוצעים של הגורמים האחרים להיות 1:

$$avg_k \{ a_k \} = 1, \quad avg_k \{ f(seq_k) \} = 1, \quad avg_k \{ \eta_k \} = 1,$$

ואז N הקבוע שמתקן בהתאם, כלומר

$$N \approx med_k \{ \lambda_k \}$$

מעבדה לסטטיסטיקה: אמידת מספר עותקים (2)

אנחנו רוצים להשתמש בפונקציית ה GC שלמדנו ובערכים הנצפים Y_k לבין מספר-העותק של הגנום. ראשית נדבר על אמידה שאינה דורשת דוגמת control, ואח"כ נציג אמידה שמשתמשת בדוגמת control. לעיתים, קוראים לשיטות אמידה אלו "ניתוק", משום שהן מתקנות את ההטיה הנוצרת מאפקט GC.

1. אמידה ללא בדיגמה בודדת

1.1. מודל

ניח שמספר הפרגמנטים Y_k הממופים לתא ka מתפלגים $Poisson(\lambda_k)$. אנחנו בנה מודל לפרמטר התחלת λ_k שנותן שמות לתרומה של GC ומספר העותקים.

נסמן ב $\{0, \frac{1}{2}, 1, \frac{3}{2}, 2, \dots\}$ את מספר העותקים. כאשר בדרך כלל $c_k = 1$. ניח ששינוי מספר העותקים משפיע כפולית על מספר הפרגמנטים הנצפים. נגדיר לכן:

$$\lambda_k = a_k \cdot \mu_k$$

כאשר μ_k הוא הערך הצפוי של פרמנטים אם היה עותק יחיד. מזכר, שאנחנו ניסו להבין כיצד מושפע מהרצף הגנומי המרכיב את התא k, לדוגמה מכמות GC בתא או מתוך כמות GC באזורים בתא. נסמן פונקציה זו ב $f(GC_k)$. ראינו שהפונקציות אכן מנבאות חלק אך לא את כל התחלת של התא. אם כן, ניתן לכתוב:

$$\lambda_k = a_k \cdot f(GC_k) \cdot \beta_k, \quad Y_k = c_k \cdot f(GC_k) \cdot \beta_k \cdot \eta_{k,Pois}$$

כאשר β_k היא השגיאה בתחלת של התא שאינה קשורה ב GC, $\eta_{k,Pois}$ מסמן את השגיאה הנובעת מההגרלה הפואסונית. ניח ש $E_k[\beta] = 1$ ו $E_{Pois}[\eta_{k,Pois}] = 1$.

הערות:

- לשם נוחות ההצגה השמטתי את הקבוע 1 המסמן את התחלת הכללית. מבחינת הניסוי, בדי"כ כן מתעניינים בקבוע זה כדבר נפרד מפונקציית $f(GC_k)$.

1.2. אמידה בדיגמה בודדת

כעת, ניח שאמדנו את הפונקציה $f(GC_k)$ בעזרת רגרסיה, ונסמן את הפונקציה שנאמדה ב $f(GC_k)$. אם כן, ניתן לאמוד את c_k ע"י הצבה (plug-in) במשוואת התחלת:

$$E_k[E_{Pois}[Y_k]] = c_k \cdot f(GC_k) \\ a_k = \frac{E[Y_k]}{f(GC_k)}$$

ולכן, מעקרון אמידת מומנטים נקבל את האומד הבא:

$$\hat{a}_k^{(1)} = \frac{Y_k}{\hat{f}(GC_k)}$$

1.3. התפלגות ע"פ המודל

נציב את המודל Y_k ונקבל:

$$\hat{a}_k^{(1)} = \frac{a_k \cdot f(GC_k) \cdot \beta_k \cdot \eta_{k,Pois}}{\hat{f}(GC_k)} = a_k \cdot \frac{f(GC_k)}{\hat{f}(GC_k)} \cdot \beta_k \cdot \eta_{k,Pois}$$

נסמן את השונות של $\eta_{k,Pois}$ ב $\beta_k \cdot \eta_{k,Pois}$, ונניח שהאמידה של f טובה ולכן $\frac{f(GC_k)}{\hat{f}(GC_k)} \approx 1$. אם כן ההטיה של $\hat{a}_k^{(1)}$ קרובה ל-0, והשונות של $\hat{a}_k^{(1)}$ בקירוב $c_k^2 \sigma^2 / f^2(GC_k)$:

$$Bias[\hat{a}_k^{(1)}] = \\ E_k E_{Pois} \left[c_k \cdot \frac{f(GC_k)}{\hat{f}(GC_k)} \cdot \beta_k \cdot \eta_{k,Pois} \cdot -a_k \right] = a_k \left(\frac{f(GC_k)}{\hat{f}(GC_k)} \cdot E_k E_{Pois}[\beta_k \cdot \eta_{k,Pois}] - 1 \right) \approx 0. \\ Var[\hat{a}_k^{(1)}] = \frac{Var[Y_k]}{\hat{f}(GC_k)^2}$$

1.4. ייצוב האומד (בעזרת epsilon)

בעזרת תוספת של $\epsilon > 0$ הן למונה והן למכנה, ניתן לשנות את האיזון בין הטיה לשונות. נגדיר את האומד

$$\hat{a}_k^{(1+\epsilon)} = \frac{Y_k + \epsilon}{\hat{f}(GC_k) + \epsilon}$$

מבחינת ההתפלגות, נשים לב שעבור $1 \neq a_k$, ההטיה של האומד $\hat{a}_k^{(1+\epsilon)}$ עם הייצוב גדלה.

מצד שני, השונות של $\hat{a}_k^{(1+\epsilon)}$ היא בקירוב $\frac{a_k^2 \sigma^2}{\hat{f}^2(GC_k + \epsilon)^2}$ ועל כן קטנה יותר מהשונות של $\hat{a}_k^{(1)}$.

2. אמידה בשתי דיגמות עם GC

כעת, נדבר על אמידה (או ניתוק) בעזרת דיגמה נוספת. במקרה כזה, נגדיר את מספר העותקים של הדיגמות הבריאה להיות תמיד 1. כמו כן, נפרק את השונות של תחלת-התא לרכיב שמשווקי לשתי הדיגמות, ולרכיב ששונה בין הדיגמות.

נסמן ב Y_k^t את דיגמות הסרטן וב Y_k^n את הדיגמות הנורמליות. ניח את המודל הבא:

$$Y_k^t = a_k \cdot f^t(GC_k) \cdot \gamma_k \cdot \delta_k^t \cdot \eta_{k,Pois}^t \\ Y_k^n = f^n(GC_k) \cdot \gamma_k \cdot \delta_k^n \cdot \eta_{k,Pois}^n$$

פצלנו את β_k^t, β_k^n לרכיב המשותף γ_k בין הדיגמות ולרכיבים הנפרדים δ_k^t, δ_k^n .

האומד של מספר העוקפים המבוסס על שתי דגימות יהיה בצורה הבאה:

$$\hat{a}_k^{(2)} = \frac{\gamma_k^t / \hat{f}^t(GC_k) + \epsilon}{\gamma_k^n / \hat{f}^n(GC_k) + \epsilon}$$

2.1. התפלגות ע"פ המודל

עם נתעלם רגע מההסתכלות

$$\hat{a}_k^{(2)} = a_k \cdot \frac{f^t(GC_k) f^n(GC_k)}{\hat{f}^t(GC_k) \hat{f}^n(GC_k)} \cdot \frac{\delta_k^t \cdot \eta_{k,Pois}^t}{\delta_k^n \cdot \eta_{k,Pois}^n} \approx a_k \cdot \frac{\delta_k^t \cdot \eta_{k,Pois}^t}{\delta_k^n \cdot \eta_{k,Pois}^n}.$$

אם כן, את הרעש ב אומד דגימה אחת $\hat{a}_k^{(1)}$ שהיה $\delta_k^t \cdot \eta_{k,Pois}^t = \gamma_k \cdot \delta_k^t$, אנחנו

מחליפים ב $\frac{\delta_k^t \cdot \eta_{k,Pois}^t}{\delta_k^n \cdot \eta_{k,Pois}^n}$.

3. תיקוני חציון

אם אנחנו מודדים את איכות האומדים על ידי השוואה לערך המקובל¹, כדאי קודם כל לדאוג שהנתונים אכן מיושרים ל1. לצורך כך, ניתן לחלק כל אחד מוקטורי האומדים בחציון שלהם, וכך נבטיח שחציון האומד יהיה 1.