

## 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 שאיננה דורשת דוגמת control, ואח"כ נציג אמידה שמשתמשת בדוגמת control. לעיתים, קוראים לשיטות אמידה אלו "תיקון", משום שהן מתקנות את ההטיה הנוצרת מאפקט ה GC.

### 1. אמידה לתא בדגימה בודדת

#### 1.1. מודל

נניח שמספר הפרגמנטים  $Y_k$  הממופים לתא הא מתפלגים  $Poisson(\lambda_k)$ . אנחנו נבנה מודל לפרמטר התוחלת  $\lambda_k$  שנותן שמות לתרומה של ה GC ומספר העותקים.

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

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

כאשר  $\mu_k$  הוא הערך הצפוי של פרגמנטים אם היה עותק יחיד. נזכור, שאנחנו ניסינו להבין כיצד  $\mu_k$  מושפע מהרצף הגנומי המרכיב את התא ה  $k$ , לדוגמה מכמות ה GC בתא או מתוך כמות ה GC באזורים בתא. נסמן פונקציה זו ב  $f(GC_k)$ . ראינו שהפונקציות אכן מנבאות חלק אך לא את כל התוחלת 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  $\beta_k$  היא השגיאה בתוחלת של התא שאיננה קשורה ב GC, ו  $\eta_{k,Pois}$  מסמן את השגיאה הנובעת מההגרלה הפואסונית. נניח ש  $E_k[\beta] = 1$  וש  $E_{Pois}[\eta_{k,Pois}] = 1$ .

#### הערות:

א. לשם נוחות ההצגה השמטתי את הקבוע  $J$  המסמן את התוחלת הכללית. מבחינת הניסוי,

בד"כ כן מתעניינים בקבוע זה כדבר נפרד מפונקציית  $f(GC_k)$ .

### 1.2. אמידה בדגימה בודדת

כעת, נניח שאמדנו את הפונקציה  $f(GC_k)$  בעזרת רגרסיה, ונסמן את הפונקציה שנאמדה ב

$\hat{f}(GC_k)$ . אם כן, ניתן לאמוד את  $c_k$  ע"י הצבה (plug-in) במשוואת התוחלת:

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

$$a_k = \frac{E[Y_k]}{f(GC_k)}$$

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  $\beta_k \cdot \eta_{k,Pois}$  ב  $\sigma^2$ , 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  $c_k^2 \sigma^2 / \hat{f}^2(GC_k)$ :

$$Bias[\hat{a}_k^{(1)}] =$$

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

$$Var[\hat{a}_k^{(1)}] = \frac{Var[Y_k]}{\hat{f}(GC_k)^2}.$$

### 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

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

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

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

## 2. אמידה בשתי דגימות עם GC

כעת, נדבר על אמידה (או תיקון) בעזרת דגימה נוספת. במקרה כזה, נגדיר את מספר העותקים של הדגימות הבריאה להיות תמיד 1. כמו כן, נפרק את השונות של תוחלת-התא לרכיב שמשותף לשתי הדגימות, ולרכיב ששונה בין הדגימות.

נסמן ב  $Y_k^t$  את דגימות הסרטן וב  $Y_k^n$  את הדגימות הנורמליות. נניח את המודל הבא:

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

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

פיצלנו את  $\beta_k^t, \beta_k^n$  לרכיב המשותף  $\gamma_k$  בין הדגימות ולרכיבים הנפרדים  $\delta_k^t, \delta_k^n$ .

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

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

## 2.1. התפלגות ע"פ המודל

עם נתעלם רגע מהepsilon

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

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

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

## 3. תיקוני חציון

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