

מעבדה לסטטיסטיקה: מודל פוואסון עם פרמטר משתנה

כזכור Y_i הוא מספר הפרגמנטים המתחילים בבסיס ה- i , $i=1, \dots, J$, ואילו X_j , $j=1, \dots, J$ משתנים מקריים המצינים לאן נפל הפרגמנט ה- j . בשבועיים החולפים הראינו שהנחת האחידות לא מתארת היטב את הנתונים שלנו, ואנחנו צריכים מודל גמיש יותר. מודל הפוואסון המשתנה אכן גמיש בהרבה, ומהווה את המודל הבסיסי לתופעות דיסקרטיות רבות.

מודל פוואסוני לבסיס בודד

ע"פ מודל זה, נמשיך להניח שכל הפרגמנטים מתפלגים אחיד (הנחה 1), ונזכיר את הסימון

$$p_i = P(X_1 = i) = P(X_2 = i) = \dots = P(X_J = i)$$

עבור ההסתברות שפרגמנט כלשהו יתחיל בבסיס ה- i (הסתברות זהה לכל פרמנט).

כמו כן, נניח שהפרגמנטים אינם תלויים (ההנחה כנראה איננה בדיוק נכונה, אך ניתן להאמין שהתלות קטנה). בפרט, משתני האינדיקטור $1(X_j = i)$ הינם משתני ברנולי בלתי תלויים ושווי התפלגות, עם הסתברות הצלחה p_i . מכאן, ש $Y_i = \sum_{j=1}^J 1(X_j = i)$ מתפלג בינומית:

$$Y_i \sim \text{Binomial}(J, p_i),$$

ובקירוב פוואסוני

$$Y_i \sim \text{Poisson}(J \cdot p_i).$$

מודל פוואסוני לתא

היתרון הגדול שמאפשר הקירוב הפוואסוני זה שקל מאוד לחבר משתנים פוואסונים בלתי תלויים. תזכורת:

א. Y_1, Y_2 משתנים פוואסונים בלתי תלויים אז $Y_1 + Y_2 \sim \text{Poisson}(E[Y_1] + E[Y_2])$.

ב. Y_1, \dots, Y_n משתנים פוואסונים בלתי תלויים אז $Y_1 + \dots + Y_n \sim \text{Poisson}(\sum_{i=1}^n E[Y_i])$.

עבור תאים בגודל b , נסמן בעזרת \tilde{Y}_k את מספר הפרגמנטים בתא ה- k , כלומר:

$$\tilde{Y}_k = Y_{bk-b+1} + \dots + Y_{bk}$$

אם כן, מהנחת חוסר התלות:

$$\tilde{Y}_k \sim \text{Poisson}(\lambda_k), \quad \lambda_k = J \cdot (p_{bk-b+1} + \dots + p_{bk}).$$

הערות:

1. בהמשך נזניח את הטילדה ב \tilde{Y}_k ונסמן במקומו Y_k , כאשר האינדקס (k) וההקשר יבהירו

שמדובר בתא ולא בבסיס יחיד.

2. כזכור, השונות של פואסון שווה לתוחלת. כלומר $\text{Var}(Y_k) = E[Y_k] = \lambda_k$. בפרט, ככל שהתוחלת גבוהה יותר גם השונות גדלה. עם זאת, כדאי לזכור שהמשמעות של המשוואה היא שהפיזור (סטיית התקן, או ממוצע המרחקים) גדלה בסדר גודל של שורש התוחלת.

3. פירוק השונות: כאשר מסתכלים על השונות הנצפית בין ערכי Y_k השונים, נשים לב שישנם שני מרכיבים לשונות. ראשית, התוחלות של כל Y_k אחרת, ושנית, כל Y_k מתפלג סביב תוחלתו. למעשה, השונות הנצפית היא סכום של שני רכיבים אלו:

$$\begin{aligned}\text{Var}[Y_k] &= E_{k \in K}[\text{Var}_{\text{Pois}}[Y_k]] + \text{Var}_{k \in K}[E_{\text{Pois}}[Y_k]] \\ &= E_{k \in K}[\lambda_k] + \text{Var}_{k \in K}[\lambda_k] = \bar{\lambda} + \frac{1}{K}(\lambda_k - \bar{\lambda})^2.\end{aligned}$$

כאשר $\bar{\lambda} = \frac{1}{K} \sum \lambda_k = J \cdot b/I$ היא התוחלת הממוצעת לתא בגודל b .

4. אפילו אם p_1, \dots, p_I שונים אחד מהשני, זה לא מספיק כדי לראות הבדלים גדולים בין ה λ_k ות. הסיבה להבדלים הגדולים שאנחנו בין ה λ_k היא שההסתברויות בסביבה הקרובה דומות אחת לשניה, ולכן לא מתבטלות כאשר נסכמים יחדיו. אם לא היה מבנה מרחבי ובאופן מקרי היינו מאחדים תאים, רוב ההבדלים היו נעלמים.

5. לפחות בשלב זה ננסה לחזות את λ_k ישירות ע"י התאמת קו רגרסיה מ GC_k , אולם המודל מזמין גם לנסות לחזות את ה λ_k ישירות. אני מקווה שנראה עוד מזה בשבועות האחרונים של המעבדה.