

למידה חישובית – שיעור 5

נרצה להמשיך את הדיון שהתחלנו בתרגול להערכת ההסתברות שקיבלנו.

כיצד נעריך את ההסתברות?

מודל פרמטרי

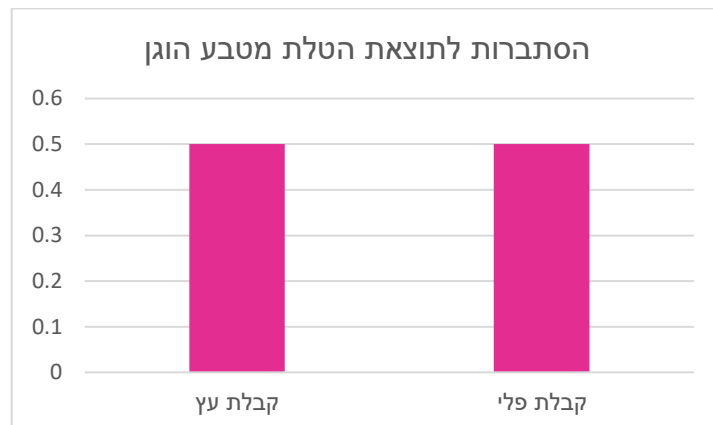
אם אנו יודעים את סוג ההתפלגות, או שאנו יכולים לנחש אותה, נוכל לעריך את הפרמטרים של ההתפלגות.

לכל Class נעריך את פרמטרי ההתפלגות לפי train dataset, ובאמצעותם נחשב את Likelihood של כל Class. לאחר שיש בידנו את ערך Likelihood נבצע קלסיפיקציה לפי ערך ההסתברות הגדול ביותר שהתקבל מההתפלגות הנורמלית

משפט הגבול המרכזי

לכל הסתברות עם ערך μ ושונות σ^2 , ההתפלגות ההמוצעת (והסכום) מתוכה היא התפלגות נורמלית עם ממוצע μ ושונות $\frac{\sigma^2}{n}$.

חשוב לזכור שלא כל התפלגות היא נורמלית, דוגמה מובהקת לכך היא הטלת מטבע הוגן. אם נבצע 10,000 הטלות הוגנות למטבע, נקבל את הגרף הבא:



נשים לב שההתפלגות שהתקבלה היא יוניפורמית, אך היא לא נראית התפלגות נורמלית. הסיבה היא כמות התוצאות האפשריות למבחן זה. מכיוון שיש לי רק 2 תוצאות אפשריות, ההתפלגות אינה נורמלית. אם ההתפלגות אינה נורמלית – לא נשתמש בכלי זה. נלמד להתמודד עם זה בשיעור הבא.

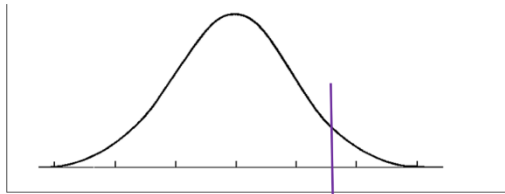
נוכל לשנות את שאלת המבחן להטלת מטבע הוגן 10 פעמים – ולבצע את המבחן הזה 10,000 פעמים, ובכך לקבל התפלגות נורמלית.

כלומר אנו מסיקים שההתפלגות הינה נורמלית (או לא) לפי המבחן עצמו ולא לפי כמות ההטלות (לפי ה-10 ולא לפי ה-10,000).

פונקציית ההתפלגות המצטברת (cumulative distribution function)

זוהי פונקציה המקבלת משתנה מקרי x ומחשבת את ההסתברות ש: $P(X \leq x)$, כאשר X הוא משתנה מקרי המקבל ערך הקטן או שווה ל- x . זהו מודל יוריסטי שמתיימר להיות קרוב לטבע. המודל הנורמלי עוזר לנו להכליל את ההתפלגות.

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x \exp\left(-\frac{x^2}{2}\right) dx$$



$\Phi(x)$ is the area under the standard Gauss curve to the left of x

הערה: נוכל להשתמש בהתפלגויות שאינן נורמליות, לדוגמה: פואסוניות. אנו לא נלמד על כך.

Normal Distribution Parameters

נרצה להתאים לכל attribute את ההתפלגות הנורמלית שלו. על מנת לעשות זאת, נחפש עבורו את μ ו- σ ונבצע קלסיפיקציה על ידי החישוב:

$$P(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

סטנדרטיזציה

נרצה לנרמל את הערכים שקיבלנו. נוריד מהמשתנה את התוחלת שלו, ונחלק בשונות.

$$U = \frac{(X - \mu)}{\sigma}$$

אינטואיציה לסטנדרטיזציה:

כשאנו מחסירים ממשתנה מקרי X , את התוחלת שלו, נקבל 0:

$$X - \mu \rightarrow E(X - \mu) = E(X - E(X)) = E(x) - E(E(x)) = E(x) - E(X) = 0$$

לאחר מכן נכפול $\frac{1}{\sigma}$ כלומר מחלקים ב-*standard deviation* (סטיית תקן)

$$Y = aX \rightarrow E(Y) = aE(X)$$

$$V(Y) = a^2 V(X)$$

$$\text{std}(Y) = a \cdot \text{std}(X) = \frac{1}{\sigma} \cdot \sigma(X) = 1$$

כלומר תהליך הסטנדרטיזציה עוזר לנו להגיע לכך ש: $\mu = 0$, $\sigma = 1$.

אם X נורמלי, אז $U \sim (0,1)$ עם צפיפות נתונה של:

$$p(u) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u^2}{2}\right)$$

אם אנו משתמשים במודל נורמלי, ננסה למצוא את המודל הנורמלי הטוב ביותר עבור $P(x|A_i)$, לכל i .
לכל attribute שיש לי, אני אחשב:

$$P(x|A_i) = \frac{1}{\sigma_i \sqrt{2\pi}} e^{-\frac{(x-\mu_i)^2}{2\sigma_i^2}}$$

נחשב את השונות וסטיית התקן:

$$\mu_i = \frac{1}{|A_i|} \sum_{x \in A_i} x$$

$$\sigma_i = \sqrt{\frac{1}{|A_i|} \sum_{x \in A_i} (x - \mu_i)^2}$$

את התהליך הזה נבצע עבור לכל instance כמספר attributes שיש ברשותי. (בדוגמה בכיתה של האיריסים, נחשב את התוחלת והשונות 12 פעמים - 3 instances ו-4 attributes, ואז נחשב את $P(x|A_i)$ עבור ערכים אלה)

התפלגות נורמלית רב מימדית

נאמר ש $X \sim N(\mu, S)$, כאשר ל- X יש התפלגות גאוסיאנית רב מימדית:

$$p(x) = \frac{1}{(2\pi)^{d/2} |S|^{1/2}} \exp\left[-\frac{1}{2} (x - \mu)^t S^{-1} (x - \mu)\right]$$

כאשר:

- μ וקטור הממוצעים.
- S האלכסון הראשי והמשולש העליון של מטריצה covariation.

$$S = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1d} \\ \sigma_{21} & \sigma_{22} & \dots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{d1} & \sigma_{d2} & \dots & \sigma_{dd} \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1d} \\ \sigma_{21} & \sigma_2^2 & \dots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{d1} & \sigma_{d2} & \dots & \sigma_d^2 \end{bmatrix}$$

טענה (משתנים נורמליים רב מימדיים): X_1, X_2, \dots, X_d משתנים מקריים ממימד d מתפלגים נורמלית אם \underline{m} כל צירוף לינארי $a_1 X_1 + a_2 X_2 + \dots + a_d X_d$ יניב משתנה מקרי גאוסיוני חד מימדי.

Naïve Bayes Classifier

גישה נאיבית (ומכאן שמו) המבוססת על חוק בייס ועל ההנחה הנאיבית שאין תלות בין ה-attributes בהינתן ה-class, וכן כל ה-instances בלתי תלויים זה בזה. בעקבות ההנחה אנו צריכים לצפות רק ב-instances המכילים features מסוימים ולא צירופים של features עבור instances שונים. כלומר, נקבל:

$$P(x_1, x_2, \dots, x_d | A_i) = \prod_{j=1}^d P(x_j | A_i)$$

נסווג instance עם תכונה \vec{x} כ:

$$\begin{aligned} \operatorname{argmax}_i P(A_i) P(\vec{x} | A_i) = \\ \operatorname{argmax}_i P(A_i) \prod_{j=1}^d P(x_j | A_i) \end{aligned}$$

הצעד הראשון בקלסיפיקציה באמצעות נאיב בייס הינו שערך ההתפלגות המותנית לכל feature ולכל class.

הגישה הנאיבית פחות מדויקת, אך היא דורשת פחות data ללמוד ממנו. בנוסף השיטה הנאיבית מהירה יותר, בגלל שאנו מבצעים הנחת אי תלות שחוסכת לנו הרבה חישובים.

הגישה הנאיבית מספקת לנו שיערוך לא מספיק טוב וכן לעיתים (קרובות) ההנחה של אי התלות אינה נכונה, אך בפועל התוצאות שיתקבלו משימוש ב-Naive Bayes מספיק טובות בשביל להשוות בין ה- argmax (המקום שבו מתקבל המקסימום) של ההסתברות:

$$\operatorname{argmax}_j \hat{P}(A_j) \prod_i \hat{P}(x_i | A_j) = \operatorname{argmax}_j \hat{P}(A_j) \hat{P}(x_1, \dots, x_n | A_j)$$

זה בעצם מה שאנו מחפשים.

Laplace Estimation

אם ניתקל במידע חסר שגורר מצב בו ההסתברות של המאורע היא 0 (כי יש לי bin ריקים עבור מקרים מסוימים, אבל זה לא אומר בהכרח שהם לא קיימים), נשתמש בשיערוך Laplace:

$$\hat{P}(x_i | A_j) = \frac{n_{ij} + 1}{n_j + |V_i|}$$

כאשר:

- n_j - מספר ה-instances עם class A_j
- n_{ij} - מספר ה-instances עם class A_j , עם ערך x_i attribute
- $|V_i|$ - מספר הערכים האפשריים של ה-attribute x