

## אלגוריתמים הסתברותיים

האלגוריתם האינטואיטיבי ביותר הוא להחזיר את קלאס הרוב, או במילים אחרות להחזיר את הקלאס הסביר ביותר לדאטה אימון. היום נלמד אלגוריתמים הסתברותיים – אלגוריתמים שמשמשים בטכניקות הסתברותיות כדי לסווג את הדגימה החדשה לקלאס.

## חזרה קצרה בהסתברות:

- מרחב מדגם – הוא קבוצה של מאורעות events, שזו רשימה של כל התוצאות האפשריות של המאורע. מרחב הדגימות יכול להיות גם רציף (גבהים וכו')
- מאורע – תת קבוצה של מרחב המדגם. זריקה של שתי קוביות שהסכום על הקוביות הוא 7.
- משתנה מקרי – פונקציה ממרחב המדגם ומחזירה ערך רציף. למשל סכום של שתי קוביות. עליו ניתן לשאול מה ההסתברות שאיקס שווה לתוצאה מסוימת. למשל מה ההסתברות שאיקס = 1 עבור איקס המתאר את סכום על שתי הקוביות (ההסתברות היא 0).
- על משתנים מקריים אפשר לשאול מה התוחלת שמשומנת ב-mu ביוונית (ה-expected value),

$$E[X] = \sum_x xp(x)$$

עבור משתנה בדיד התוחלת הינה

$$E[X] = \int_{-\infty}^{\infty} xf(x)dx$$

ועבור משתנה רציף כאשר f היא פונקציית הצפיפות.

- שונות – מסומנת על ידי סיגמה בריבוע. מוגדרת להיות:  $\sigma^2 = \text{var}(X) = E[(x - \mu)^2]$

$$\sigma = \sqrt{\text{var}(X)} = \sqrt{E[(x - \mu)^2]}$$

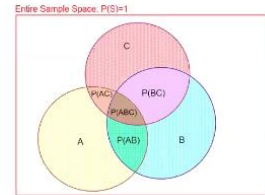
סטיית תקן – היא שורש השונות

$$P(A \cup B) = ?$$

$$P(A) + P(B) - P(A \cap B)$$

$$P(A \cup B \cup C) = ?$$

$$\begin{aligned} P((A \cup B) \cup C) &= \\ P(A \cup B) + P(C) - P((A \cup B) \cap C) &= \\ P(A) + P(B) - P(A \cap B) + P(C) - P((A \cap C) \cup (B \cap C)) &= \\ P(A) + P(B) - P(A \cap B) + P(C) - P(A \cap C) - P(B \cap C) + P(A \cap B \cap C) &= \\ P(A) + P(B) + P(C) - P(A \cap B) - P(A \cap C) - P(B \cap C) + P(A \cap B \cap C) \end{aligned}$$



## חזרה קצרה על הסתברות מותנית

### Conditional probability:

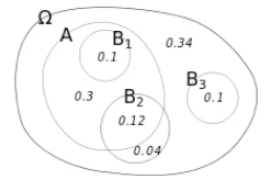
$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$P(A|B_1) = ?$$

$$\frac{P(A \cap B_1)}{P(B_1)} = \frac{0.1}{0.1} = 1$$

$$P(A|B_2) = ?$$

$$\frac{P(A \cap B_2)}{P(B_2)} = \frac{0.12}{0.16} = 0.75$$



### דוגמה עבור הסתברות מותנית:

נתון כי ההסתברות לעבור את המבחן היא  $P(\text{pass}) = 90\%$  ולהיכשל  $P(\text{fail}) = 10\%$ . בנוסף ידוע לנו כי:

- $P(\text{Learn for the test}|\text{Pass}) = 90\%$
- $P(\text{Didn't learn}|\text{Pass}) = 10\%$
- $P(\text{Learn for the test}|\text{Fail}) = 5\%$
- $P(\text{Didn't learn}|\text{Fail}) = 95\%$

- $P(\text{Pass} \cap \text{Learn for the test}) = P(\text{Pass}) \times P(\text{Learn for the test}|\text{Pass}) = 90\% \times 90\% = 81\%$
- $P(\text{Pass} \cap \text{Didn't learn}) = P(\text{Pass}) \times P(\text{Didn't learn}|\text{Pass}) = 90\% \times 10\% = 9\%$
- $P(\text{Fail} \cap \text{Learn for the test}) = P(\text{Fail}) \times P(\text{Learn for the test}|\text{Fail}) = 10\% \times 5\% = 0.5\%$
- $P(\text{Fail} \cap \text{Didn't learn}) = P(\text{Fail}) \times P(\text{Didn't learn}|\text{Fail}) = 10\% \times 95\% = 9.5\%$
- $P(\text{Learn for the test}) = P(\text{Pass} \cap \text{Learn for the test}) + P(\text{Fail} \cap \text{Learn for the test}) = 81\% + 0.5\% = 81.5\%$

מה ההסתברות שלמדת למבחן: 0.815 להלן החישובים הדרושים לפתרון

$$P(\text{Pass}|\text{Learn for the test}) = \frac{P(\text{Pass} \cap \text{Learn for the test})}{P(\text{Learn for the test})} = \frac{81\%}{81.5\%} = 99\%$$

$$P(\text{Fail}|\text{Learn for the test}) = \frac{P(\text{Fail} \cap \text{Learn for the test})}{P(\text{Learn for the test})} = \frac{0.5\%}{81.5\%} = 1\%$$

$$P(\text{Pass}|\text{Didn't learn}) = \frac{P(\text{Pass} \cap \text{Didn't learn})}{P(\text{Didn't learn})} = \frac{9\%}{18.5\%} = 49\%$$

$$P(\text{Fail}|\text{Didn't learn}) = \frac{P(\text{Fail} \cap \text{Didn't learn})}{P(\text{Didn't learn})} = \frac{9.5\%}{18.5\%} = 51\%$$

#### Independent events

- If  $P(A \cap B) = P(A)P(B)$  then A & B are independent
- From conditional probability we get:

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

$$\downarrow$$

$$P(A \cap B) = P(A|B)P(B)$$

- If A & B are independent:

$$P(A)P(B) = P(A \cap B) = P(A|B)P(B)$$

$$P(A) = P(A|B)$$

\* And also  $P(B) = P(B|A)$

#### מתבקש לשאול, ולהלן החישובים

- מה ההסתברות שסטודנט עבר בהינתן שלמד?
- מה ההסתברות שסטודנט נכשל בהינתן שלמד?
- מה ההסתברות שסטודנט עבר בהינתן כי לא למד?
- ומה ההסתברות שסטודנט נכשל בהינתן כי לא למד?

#### נגזיר מאורעות בלתי תלויים:

אם ההסתברות לחיתוך שווה למכפלת ההסתברויות אזי המאורעות בלתי תלויים.  
המסקנה: זה שאחד מהם קרה, לא גורע ולא מוסיף למאורע השני.

כעת, נוכל להתקדם למושגים שיותר רלוונטיים עבור Bayes.

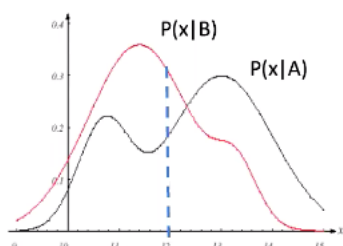
#### Prior Probability / הסתברות פריורית

זו היא ההסתברות של המקרים בדאטה שיש לנו, האלגוריתם הנאיבי שתואר בעמוד הקודם משתמש בהסתברות prior. כלומר, החלוקה הפנימית של מאורעות, למשל 40% בנים ו-60% בנות.

#### Likelihood Probability / הסתברות הלייקליהוד

זו היא ההסתברות שמותנית על ידי קלאס: ההסתברות של דגימה x בהינתן

הקלאס. למשל עבור דגימה  $x=12$  ושני קלאסים אפשריים A ו-B: האדומה היא ההסתברות של B והשחורה היא A. נחשב את ערך ההסתברות בשני הגרפים עבור  $x=12$  ונראה כי ההסתברות יותר גבוהה עבור קלאס B ולכן נסווג B. אם נחזור לדוגמת עבר/נכשל, ההסתברות זו היא כמו לשאול מה ההסתברות שמישהו למד למבחן בהינתן שעבר. אבל, אנחנו רוצים לדעת מה ההסתברות לעבור או להיכשל בהינתן שלמדנו. אז אנחנו צריכים דרך לעבור מה-likelihood probability ל-posterior probability.



If  $x=12$ , we'll predict B, because  $P(x|B) > P(x|A)$

#### Posterior Probability and Bayes rule / הסתברות הפוסטריור וחוק בייס

חוק בייס מניב עבורנו את ההסתברות הפוסטריור = מה ההסתברות לקלאס, בהינתן האינסטנס.

כדי לחשב את הפוסטריור נצטרך להשתמש בהסתברות הלייקליהוד והפריור (במונה) ובהסתברות המכנה (אווידנס).

$$P(A|x) = \frac{P(x|A)P(A)}{P(x)}$$

- We will classify A if

$$P(A|x) = \frac{P(x|A)P(A)}{P(x)} > \frac{P(x|B)P(B)}{P(x)} = P(B|x)$$

$$P(x|A)P(A) > P(x|B)P(B)$$

- Note that  $P(x)$  is removed from both denominators simply because it is the same

מסווג שישווג A אם מתקיים  $P(A|x) > P(B|x)$ , הוא מסווג שמביא למקסימום את ההסתברות הפוסטריור – MAP Classifier. הקלסיפיקציה תלויה בפריור והלייקליהוד, האוידנס לא יעניין אותנו מפני שהוא במכנה והוא שווה בכל החישובים שלנו.

#### שגיאה מינימלית במסווג מסוג MAP:

But, we classify B only if  $P(B|x) > P(A|x)$ , and therefore the probability of the error is minimal

$$P(\text{error}|x) = \min[P(A|x), P(B|x)]$$

יש פה קאץ', זה לא שאנחנו ממזערים את הטעות הגלובלית העולמית של הפרדיקציה, אלא בהתייחס להסתברויות שלמדנו מהדאטה הטעות היא מינימלית. כלומר אם ההסתברויות של הדאטה אכן מייצגות את העולם האמיתי, אכן השגיאה היא מינימלית ביחס לעולם האמיתי.

0-1 loss (the simplest one):

$$\lambda_{ij} = \lambda(\text{Choose } A_i | A_j) = \begin{cases} 1, & \text{if } i \neq j \\ 0, & \text{if } i = j \end{cases}$$

הגדרת פונקציית ה-loss: אם בחרתי את  $A_i$  בהינתן שהקלאס הנכון היה  $A_j$  ניתן דוגמה

של 0/1 loss

**בהינתן פונקציית ה-loss נגדיר את הסיכון ה-risk:** זהו בעצם הסיכון לבחור קלאס מסוים, למשל  $A_i$  בהינתן האינסטנס, זהו ההפסד כפול ההסתברות של הקלאסים. אנחנו עושים סכום על כמות הקלאסים  $k$ , של הלוס על כל אחד של הפוסטריורים של הקלאסים.

$$R(\text{Choose } A_i|x) = \sum_{j=1}^k \lambda_{ij} P(A_j|x) = \sum_{j \neq i} P(A_j|x) = 1 - P(A_i|x)$$

דוגמה לסיכון במקרה שפונקציית הלוס הינה  $0/1$ -loss: נוכל להגדיר פונקציות לוס הרבה יותר מורכבות.

מסווג שרוצה למזער את הסיכון יבחר  $A_i$  כך ש:  $P(A_i|x) > P(A_j|x) \forall j \neq i$

$$g_i(x) = P(A_i|x) = \frac{P(x|A_i)P(A_i)}{\sum_{j=1}^k P(x|A_j)P(A_j)} = P(x)$$

**הסתברות הפוסטריור** המלאה (לפני שנשמיט את  $P(x)$ ) הינה על פי ביים:

$$g_i(x) = P(x|A_i)P(A_i)$$

ולאחר שנשמיט את  $P(x)$  (המכנה) יתקיים:

### בייס עבור קלאסים רבים

$$g_i(x) = \ln(P(x|A_i)P(A_i)) = \ln(P(x|A_i)) + \ln(P(A_i))$$

**כדי להפוך את תהליך הסיווג ליותר יעיל נוכל להשתמש ב- $\ln()$ :**

שימוש ב- $\ln$  עוזר להפחית את ההכפלות במספרים נמוכים (בין 0-1) (חוקי לוגריתמים, הופכת מכפלות לסכומים) וכן עוזר להתמודד יותר טוב עם פונקציית הצפיפות הנורמלית ( $e^{-f(x)}$ ). אנחנו יכולים להשתמש ב- $\ln$  מפני שהיא פונקציה מונוטונית עולה.

### היפותזת מקסימום לייקליהוד ML-Hypothesis

- נרצה לבחור את ההיפותזה שיש לה את ההיפותזה הכי סבירה (בעלת ההסתברות המקסימלית) בהינתן הדאטה, ובכתיב הסתברות

$$P(h|D) = P(D|h)P(h)$$

נחפש אחר ההסתברות הפוסטריורית הבאה המקסימלית:

- נוכל להניח שלכל ההיפותזות במרחב ההיפותזות יש את אותו הפריור  $P(h)$ , ולכן נוכל למצוא את ההיפותזה הכי סבירה  $h$  על פי ה-

$$h_{ML} = \underset{h \in H}{\operatorname{argmax}} p(D|h)$$

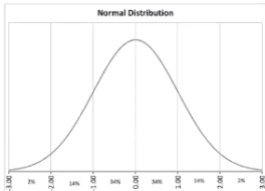
: maximum likelihood

לכן למעשה ההיפותזה שממקסמת את הלליקליהוד היא ההיפותזה בעלת ההסתברות המקסימלית בהינתן הדאטה.

- נרצה למצוא את  $h$ -ML. נוכל להניח שכל האינסטנסים הם בלתי תלויים זה בזה ועל כן מתקיים:

$$P(D|h) = \prod_i P(y_i|h)$$

- **הנחה:** אם הטעות בהיפותזה מתפלגת באופן נורמלי עם ממוצע / תוחלת 0,  $e_i \sim N(0, \sigma)$ , אזי נוכל לומר שההסתברות ש- $h(x_i) = y_i$  (שההיפותזה תדע לחזות בדיוק את ה-target value עבור אינסטנס מסויים  $x_i$ ) היא אותה ההסתברות ש- $e_i = 0$  (שאינן טעות עבור אינסטנס  $i$ ) לפי ההסתברות הנורמלית של  $e_i$ .



$$h_{ML} = \underset{h \in H}{\operatorname{argmax}} p(D|h) = \underset{h \in H}{\operatorname{argmax}} \prod_i P(y_i|h)$$

- **כעת נקבל:** (משמאל)

הנחנו אי תלות ולכן השורה הראשונה מתקיימת.

הנחנו שהטעות מתפלגת נורמלית ולכן ההסתברות לקבל את ה-target value בהינתן ההיפותזה, שווה להסתברות של הטעות ולכן נציב את נוסחת ההסתברות של הטעות המתפלגת נורמלית עבור כל אינסטנס (שורה שניה).

בשורה השלישית אנחנו מציבים את המרחק שלנו מהטעות, כי ההיפותזה נתונה.

לאחר כל זאת נוכל להפעיל  $\ln$ . שיהפוך את הפאי, המכפלה, לסכום. ונשמיט ערכים קבועים כי הם לא משנים את המקסימום.

- **לבסוף נקבל כי ההיפותזה הכי סבירה (maximum likelihood hypothesis) היא זו שמזעזעת את**

**ה-MSE = שורה אחרונה.**

$$\begin{aligned} h_{ML} &= \underset{h \in H}{\operatorname{argmax}} \prod_i \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2} \left(\frac{h(x_i) - y_i}{\sigma}\right)^2} \\ h_{ML} &= \underset{h \in H}{\operatorname{argmax}} \ln \left( \prod_i \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2} \left(\frac{h(x_i) - y_i}{\sigma}\right)^2} \right) \\ h_{ML} &= \underset{h \in H}{\operatorname{argmax}} \sum_i \ln \left( \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2} \left(\frac{h(x_i) - y_i}{\sigma}\right)^2} \right) \\ &= \underset{h \in H}{\operatorname{argmax}} \sum_i \left( \ln \left( \frac{1}{\sqrt{2\pi\sigma^2}} \right) - \frac{1}{2} \left( \frac{h(x_i) - y_i}{\sigma} \right)^2 \right) \\ &= \underset{h \in H}{\operatorname{argmax}} \sum_i -\frac{1}{2} \left( \frac{h(x_i) - y_i}{\sigma} \right)^2 \\ &= \underset{h \in H}{\operatorname{argmin}} \sum_i (h(x_i) - y_i)^2 \end{aligned}$$

## סיכום עד כה:

- Prior classifier:  $P(A) > P(B)$
- ML classifier:  $P(x|A) > P(x|B)$  – assuming  $P(A) = P(B)$
- MAP classifier:  

$$P(A|x) = P(x|A)P(A) > P(x|B)P(B) = P(B|x)$$
- \* Drooping  $P(x)$  from the denominator

## איך מחשבים/משערים את ההסתברויות? הלייקליהוד, הפוסטריו

- **Parametric estimation** – אם ידוע לנו שאנחנו יכולים לנחש את סוג ההתפלגות נוכל להעריך את הפרמטרים של ההתפלגות. למשל אם נוכל לנחש שמשנתה מקרי מסויים מתפלג נורמלית נוכל לשערך עבורו את ה-miu וה-sigma, או אם הוא מתפלג פואסונית נוכל לשערך עבורו את הלמדא.
- **Non parametric estimation** – אם אנחנו לא יכולים להניח אף סוג של התפלגות על הדאטה שלנו נשתמש בהיסטוגרמה (=ספירה) או ב-Kernel Density Estimation (שזו למעשה היסטוגרמה חלקה = smooth histogram).

## שיערוך פרמטרי Parametric Estimation

עבור כל קלאס נשערך את הפרמטרים של ההתפלגות על פי הדאטה אימון. אם אנחנו מדברים על ההתפלגות הנורמלית, עלינו לשערך את התוחלת ואת השונות של כל קלאס: ואם נרצה לסווג לפי ההסתברות הגבוהה ביותר בהינתן ההתפלגות הנורמלית (הלייקליהוד):

$$\mu = \frac{1}{m} \sum_{k=1}^m x_k$$

$$\sigma^2 = \frac{1}{m} \sum_{k=1}^m (x_k - \mu)^2$$

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

הבעיה היא שכל זאת מתאים רק לאטריביוט/קלאס יחיד, **מה אם יש יותר אטריביוטים?** אז במקרה כזה הסתברות הלייקליהוד תחושב לפי **התפלגות נורמלית רב-ממדית**. לשם כך נצטרך את וקטור התוחלות (כל ממד יהיה התוחלת של אטריביוט מסוים) ואת מטריצת השונות המשותפת the covariance matrix. (במטריצה, השונותיות שלא באלכסון הן ההסתברויות המשותפות)

## :Multivariate normal distribution

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

$|S|$  - is the determinant of the covariance matrix  
 $S^{-1}$  - is the inverse matrix of the covariance matrix

$$P(\bar{x}|A_i) = \frac{1}{\sqrt{(2\pi)^d |S|^2}} e^{-\frac{1}{2}(\bar{x} - \bar{\mu}_i)^T S^{-1} (\bar{x} - \bar{\mu}_i)}$$

## שיערוך לא פרמטרי Non Parametric

מכיוון שבשיערוך זה איננו יודעים להניח את ההתפלגות, עלינו למצוא דרך אחרת לשערך את הפריור  $P(A_i)$  ואת הלייקליהוד  $P(x|A_i)$ . הסתברות הפריור  $P(A_i)$  ניתנת לשערך מתדירות הקלאסים בדאטה אימון (מספר המופעים של הקלאס לחלק למספר הדגימות). לגבי הלייקליהוד, ראינו בהרצאה שניתן להשתמש בהיסטוגרמה וב-Interpolation, ובשבוע הבא נראה גישה נוספת.

## התמודדות עם multiple features

כדי להעריך נכון את הלייקליהוד עבור דגימה נתונה עלינו להחזיק דאטה-סט עצום: אם בדינו d אטריביוטים בדידים, ו-k קלאסים, מספר האפשרויות של הלייקליהוד  $P(x_1, x_2, \dots, x_d|A_i)$  הוא  $k \cdot |V_1| \cdot |V_2| \dots |V_d|$ . עבור רק 2 אטריביוטים וקלאס בודד נוכל להגיע ל- $2^d$ . לכן אנחנו צריכים דרך / הנחה שתעזור לנו להתגבר על בעיה זו. לכן נוכל להניח שהאטריביוטים הם בלתי תלויים בהינתן קלאס – ברע שנגניח זאת נוכל להמיר את המשוואה למכפלת ההסתברויות של כל אטריביוט בנפרד. בדאטה מופיעים לנו כל הערכים של האטריביוטים האלה וכך הדבר ניתן לחישוב:

### הנחה זו נקראת Naïve Bayes:

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

- אנחנו מניחים שכל האטריביוטים הם בלתי תלויים בהינתן הקלאס, ונקבל:

$$V_{NB} = \underset{i}{\operatorname{argmax}} P(A_i) \prod_{j=1}^d P(x_j | A_i)$$

- כעת נוכל למצוא את ה-MAP (Maximum A-Posterior Probability):  
הקלאס שיבחר הוא זה שממקסם את הביטוי הנתון.

$$\sum_{j=1}^d |V_j|$$

- ברגע שאנו מניחים כך, אנחנו מורידים משמעותית את הגודל של הדאטה-סט שנצטרך הבדיד ל-

- ניתן להניח זאת גם במקרה הרציף, ניתן להשתמש בו גם בשיערוך הפרמטרי.