

Machine learning

Naïve Bayes – part 2

- * Smoothing
 - * Continuous features
-

Exercise V

פיתוח:

ד"ר יהונתן שלר

משה פרידמן

מהתפלגות לפונקציית צפיפות

מושגים - תזכורת

משתנה מקרי: פונקציה המתאימה כל אירוע אפשרי במרחב הסתברות לערך מספרי. אצלינו – מאפיין. דוגמאות:

- ❖ מ"מ בדיד: זריקת מטבע אקראית. נוצר מ"מ בדיד בינארי. התאמת צד מטבע לערך 0, וצדו השני לערך 1
- ❖ מ"מ רציף: גובהו של אדם שנבחר באקראי הוא גם כן משתנה מקרי.

מרחב המדגם Ω : קבוצת כל התוצאות האפשריות בניסוי. אצלינו – אוסף הערכים המאפיינים האפשריים. דוגמאות:

- ❖ זריקת מטבע אקראית. מרחב המדגם: $\{0,1\}$
- ❖ טמפרטורה של מים. מרחב המדגם $[0,100]$

מאורע / תצפית: תוצאה נצפת מסוימת בניסוי מסוים. אצלינו – ערך מאפיין של דוגמה ב-dataset. דוגמאות:

- ❖ התוצאה 3 בזריקת קוביה
- ❖ גובה 1.72 של סטודנט

מושגים

הסתברות מאורע: מידת הסבירות שמאורע מסוים יתרחש.

❖ ההסתברות של מאורע יכולה לקבל ערך מספרי שבין 0 ל-1.

פונקציית צפיפות הסתברות (של משתנה מקרי) **[PDF]**: פונקציה המתארת את צפיפות המשתנה בכל נקודה במרחב המדגם.

❖ במ"מ בדיד - הצפיפות בנקודה מסוימת היא בעצם ההסתברות של המאורע (פונקצית המסה). סך כל הערכים שבפונקציית הצפיפות $= 1$

❖ במ"מ רציף - פונקציית הצפיפות לא שווה להסתברות של קיום אירוע. אפשר לראות את ה-PDF במ"מ רציף כסבירות היחסית שערך שייך להסתברות. ערכיו אי שליליים, אך לא מוגבלים ל-1 (כמו במ"מ בדיד).

פונקציית ההתפלגות המצטברת (של משתנה מקרי) **[CDF]**: פונקציה הקובעת את ההסתברות למאורעות $X \leq a$, (לכל a ממשי).

❖ נדרשת עבור מ"מ רציף

ההתפלגות (של משתנה מקרי): קובעת מהי פונקציית הצפיפות (ומהי ההסתברות של כל מאורע).

❖ במשתנה מקרי בדיד בעל אוכלוסיה סופית (או במדגם train-set) נחשב את ההסתברות (הצפיפות) בנקודה מסוימת כמספר המופעים של האירוע לחלק לסך כמות האירועים.

❖ במשתנה מקרי רציף, נמדדת בד"כ כפונקציה של הממוצע וסטיית התקן (דוגמאות בהמשך).

מושגים - תזכורת

התפלגות אחידה: התפלגות בה הצפיפות (סבירות) לכל מאורע היא זהה.

❖ **התפלגות אחידה בדידה:** ההסתברות שווה ל-1 חלקי מספר הערכים האפשריים במרחב המדגם. לדוגמה: הסתברות 1/6 לקבלת הערך 4 בקובייה הוגנת.

❖ **התפלגות אחידה רציפה:** לדוגמה: נניח ש- X מתפלג באופן אחיד בקטע $[0, 1]$. אז פונקציית ההתפלגות המצטברת שלו:

$$F(x) = \begin{cases} 0 & : x < 0 \\ x & : 0 \leq x < 1 \\ 1 & : x \geq 1. \end{cases}$$

התפלגות ברנולי: מ"מ בדיד בינארי, עם מרחב המדגם: $\{0, 1\}$. 1 – מסמן הצלחה ו-0 מסמן כישלון. אם סיכוי ההצלחה הוא p , סיכוי הכישלון הוא $q=1-p$.

מושגים - תזכורת

תוחלת: מייצגת תוצאה "צפויה" (Expected) של ניסוי זהה החוזר על עצמו פעמים רבות.

❖ עבור משתנה מקרי בדיד: $\mu = E[X] = \sum_{x \in A} P(X = x)x$ ה: התוחלת של קובייה הוגנת $3.5 =$

❖ עבור משתנה מקרי רציף: $\mu = \int x f(x) dx$

$$\text{Var}(X) = \mathbb{E}((X - \mu)^2) = \mathbb{E}(X^2) - \mu^2 = \mathbb{E}(X^2) - (\mathbb{E}(X))^2$$

שונות: מדד לפיזור ערכים באוכלוסייה נתונה ביחס לתוחלת שלה.

❖ מ"מ בדיד, אם האוכלוסייה בגודל N :
$$\text{Var}(X) = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2 = \left(\frac{1}{N} \sum_{i=1}^N x_i^2 \right) - \mu^2$$

❖ מ"מ רציף:
$$\text{Var}(X) = \sigma^2 = \int (x - \mu)^2 f(x) dx = \int x^2 f(x) dx - \mu^2$$

סטיית תקן: שורש השונות.

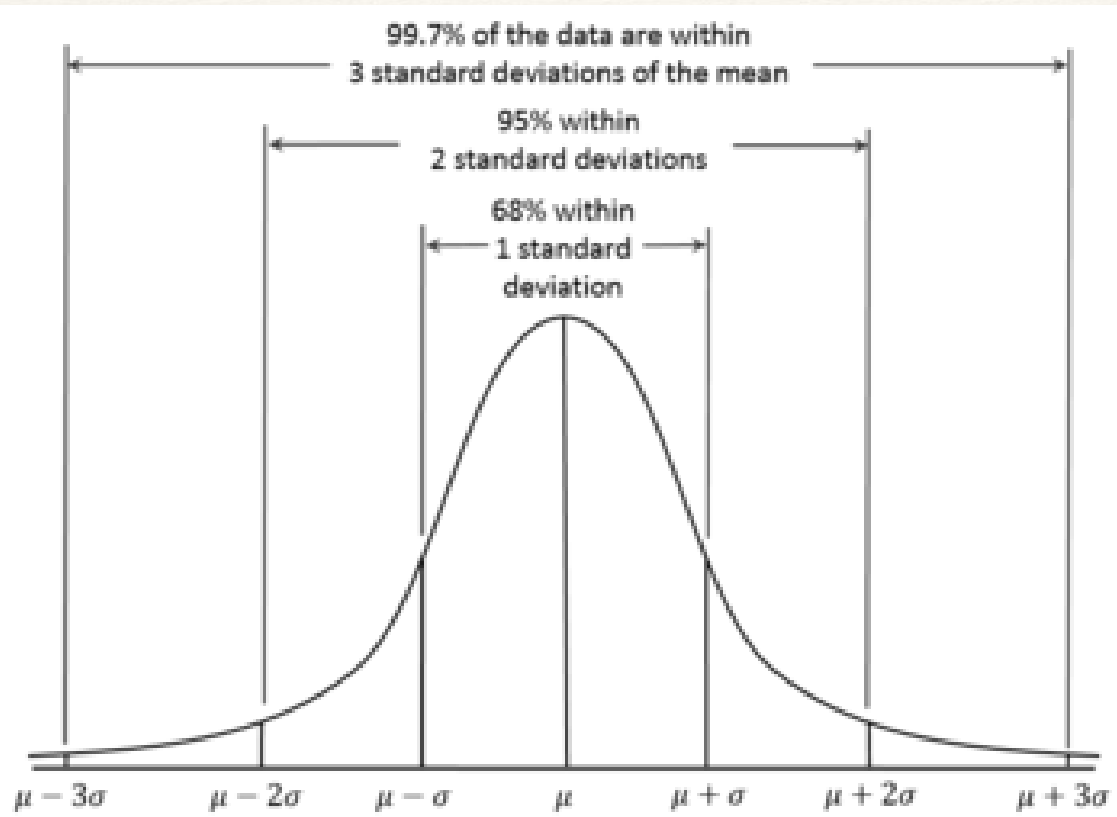
התפלגות נורמלית

התפלגות נורמלית: נקראת גם גאוסיאן (Gaussian) או עקומת פעמון.

❖ פונקציית צפיפות סמטרית.

התפלגות z : תת קבוצה של התפלגות נורמלית בו התוחלת/הממוצע $=0$ וסטיית התקן $=1$.

❖ כל התפלגות נורמלית ניתן להפוך להתפלגות z



התפלגות במדגם

מדגם (sample): מדגם הוא קבוצת פרטים, המהווה מודל לאוכלוסייה, שאליה היא שייכת. אצלינו – ה-train-set.

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

ממוצע במדגם:

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

סטיית התקן במדגם:

התפלגות t: התפלגות המבוססת על מידע שנאסף במדגם.

❖ התפלגות t שואפת להתפלגות z, כאשר גודל המדגם שואף לאינסוף.

חזרה להתפלגות נורמלית

התפלגות נורמלית: נקראת גם גאוסיאן (Gaussian) או עקומת פעמון.

❖ פונקציית צפיפות סמטרית.

התפלגות z : תת קבוצה של התפלגות נורמלית בו התוחלת/הממוצע $=0$ וסטיית התקן $=1$.

❖ כל התפלגות נורמלית ניתן להפוך להתפלגות z

❖ מהתפלגות נורמלית להתפלגות z : $z\text{-val} = \frac{(x-\mu)}{\sigma}$

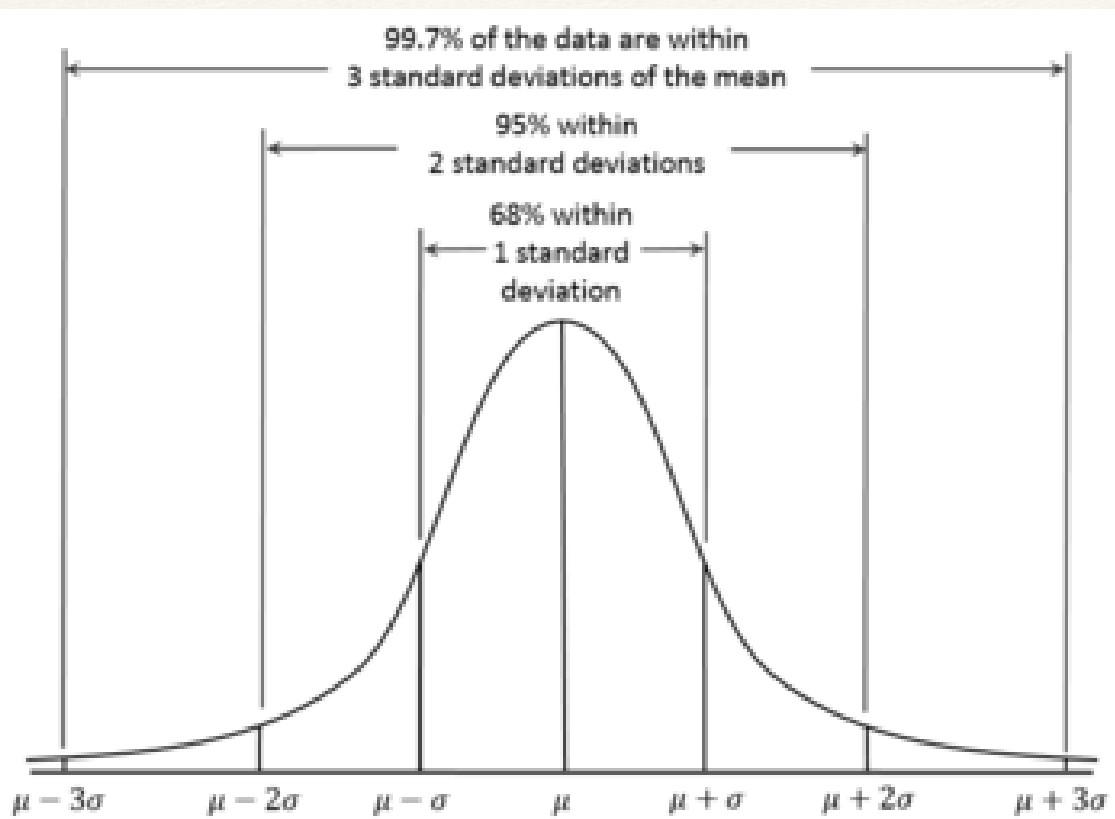
ערך מאורע

תוחלת/ממוצע

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

פונקציית צפיפות:

סטיית תקן



הערה חשובה: אנחנו נחשב סטיית תקן במדגם ואת ה- $t\text{-val}$

MAP = Maximum a posteriori (estimation) - תזכורת

❖ בהינתן ווקטור לסיווג $(x_1, x_2, x_3, \dots, x_n)$, נעריך את ההסתברות עבור כל סיווג c_i השייך לקבוצה C ונבחר את הסיווג עם ההסתברות הגבוהה ביותר.

$$P(c_1 | x_1, x_2, x_3, \dots, x_n) \quad \diamond$$

$$P(c_2 | x_1, x_2, x_3, \dots, x_n) \quad \diamond$$

$$P(c_3 | x_1, x_2, x_3, \dots, x_n) \quad \diamond$$

... \diamond

$$h_{MAP} = \arg \max_{c \in C} P(c | X)$$

כלומר, נבחר את הקטגוריה c , המקיימת

חוק בייס והנחת חוסר התלות - תזכורת

Class prior

Likelihood probability

$$P(c | x_1, x_2, \dots, x_D) = \frac{P(c)P(x_1, x_2, \dots, x_D | c)}{P(x_1, x_2, \dots, x_D)}$$

חוק בייס:

a posteriori probability

Feature (predictor) priors

בגלל הנחת חוסר התלות בין המאפיינים:

$$P(x_1, x_2, \dots, x_D | c) = P(x_1 | c)P(x_2 | c)P(x_3 | c) \dots P(x_D | c) = \prod_{i=1}^D P(x_i | c)$$

תרגיל 8 - סימולציית סיווג – מ"מ רצוף

מסוג Gaussian Naïve Bayes עבור מ"מ רציף - תזכורת

- Train Naïve Bayes (examples)

for each value y_k

estimate* $\pi_k \equiv P(Y = y_k)$

for each attribute X_i estimate $P(X_i | Y = y_k)$

- class conditional mean μ_{ik} , standard deviation σ_{ik}

- Classify (X^{new})

$$Y^{new} \leftarrow \arg \max_{y_k} P(Y = y_k) \prod_i P(X_i^{new} | Y = y_k)$$

$$Y^{new} \leftarrow \arg \max_{y_k} \pi_k \prod_i \mathcal{N}(X_i^{new}; \mu_{ik}, \sigma_{ik})$$

* probabilities must sum to 1, so need estimate only n-1 parameters...

עבור Gaussian Naïve Bayes, נשתמש
בפונקציית צפיפות:

$$P(X_i = x | Y = y_k) = \frac{1}{\sigma_{ik} \sqrt{2\pi}} e^{-\frac{(x - \mu_{ik})^2}{2\sigma_{ik}^2}}$$

תרגיל 8 - סימולציית סיווג – מ"מ רציף

outlook	temperature	humidity	windy	play
sunny	85	85	false	no
sunny	80	90	true	no
overcast	83	86	false	yes
rainy	70	96	false	yes
rainy	68	80	false	yes
rainy	65	70	true	no
overcast	64	65	true	yes
sunny	72	95	false	no
sunny	69	70	false	yes
rainy	75	80	false	yes
sunny	75	70	true	yes
overcast	72	90	true	yes
overcast	81	75	false	yes
rainy	71	91	true	no

❖ נתונות 14 דוגמאות של
טניסאי שהיה צריך
להחליט אם לשחק טניס
ביום מסוים.

❖ בנו מסווג שיחזה האם
האדם ישחק טניס ביום
עם התנאים הבאים:

outlook=overcast,
temperature=60,
humidity=62,
windy=false.

תרגיל 8 - סימולציית סיווג – מ"מ רציף - פתרון

outlook	temperature	humidity	windy	play
sunny	85	85	false	no
sunny	80	90	true	no
overcast	83	86	false	yes
rainy	70	96	false	yes
rainy	68	80	false	yes
rainy	65	70	true	no
overcast	64	65	true	yes
sunny	72	95	false	no
sunny	69	70	false	yes
rainy	75	80	false	yes
sunny	75	70	true	yes
overcast	72	90	true	yes
overcast	81	75	false	yes
rainy	71	91	true	no

שלב האימון:

א. חישוב הסתברויות priors של המחלקות:

- $p(\text{yes}) = 9 / (9+5) = 0.643$
- $p(\text{no}) = 5 / 14 = 0.357$

תרגיל 8 - סימולציית סיווג – מ"מ רציף - פתרון

outlook	temperature	humidity	windy	play
sunny	85	85	false	no
sunny	80	90	true	no
overcast	83	86	false	yes
rainy	70	96	false	yes
rainy	68	80	false	yes
rainy	65	70	true	no
overcast	64	65	true	yes
sunny	72	95	false	no
sunny	69	70	false	yes
rainy	75	80	false	yes
sunny	75	70	true	yes
overcast	72	90	true	yes
overcast	81	75	false	yes
rainy	71	91	true	no

שלב האימון:

א. חישוב הסתברויות priors של המחלקות:

- $p(\text{yes}) = 0.643$
- $p(\text{no}) = 0.357$

תרגיל 8 -

סימולציית סיווג -

מ"מ רצוף - פתרון

הווקטור החדש:

outlook=overcast,
temperature=60,
humidity=62,
windy=false.

$$P(X = x | Y = y) = P(X = x, Y = y) / P(Y = y)$$

שלב האימון (נחשב רק עבור הערכים שבווקטור החדש):

ב. חישוב ההסתברויות המותנות הבדידות:

outlook	temperature	humidity	windy	play
sunny	85	85	false	no
sunny	80	90	true	no
overcast	83	86	false	yes
rainy	70	96	false	yes
rainy	68	80	false	yes
rainy	65	70	true	no
overcast	64	65	true	yes
sunny	72	95	false	no
sunny	69	70	false	yes
rainy	75	80	false	yes
sunny	75	70	true	yes
overcast	72	90	true	yes
overcast	81	75	false	yes
rainy	71	91	true	no

$$P(\text{outlook}=\text{overcast}|\text{yes})=$$
$$P(\text{windy}=\text{false}|\text{yes})=$$

$$P(\text{outlook}=\text{overcast}|\text{no})=$$
$$P(\text{windy}=\text{false}|\text{no})=$$

outlook	temperature	humidity	windy	play
sunny	85	85	false	no
sunny	80	90	true	no
overcast	83	86	false	yes
rainy	70	96	false	yes
rainy	68	80	false	yes
rainy	65	70	true	no
overcast	64	65	true	yes
sunny	72	95	false	no
sunny	69	70	false	yes
rainy	75	80	false	yes
sunny	75	70	true	yes
overcast	72	90	true	yes
overcast	81	75	false	yes
rainy	71	91	true	no

$$P(\text{outlook}=\text{overcast}|\text{yes})=4/9 \sim 0.44$$

$$P(\text{windy}=\text{false}|\text{yes})=6/9 \sim 0.66$$

הווקטור החדש:

outlook=overcast,
temperature=60,
humidity=62,
windy=false.

תרגיל 8 -

סימולציית סיווג -

מ"מ רציף - פתרון

$$P(X = x | Y = y) = P(X = x, Y = y) / P(Y = y)$$

שלב האימון (נחשב רק עבור הערכים שבווקטור החדש):

ב. חישוב ההסתברויות המותנות הבדידות:

$$P(\text{outlook}=\text{overcast}|\text{no})=0/5=0$$

$$P(\text{windy}=\text{false}|\text{no})=2/5=0.4$$

תרגיל 8 -

סימולציית סיווג -

מ"מ רציף - פתרון

הווקטור החדש:

outlook=overcast,
temperature=60,
humidity=62,
windy=false.

$$P(\text{outlook}=\text{overcast}|\text{yes})=4/9 \sim 0.44$$
$$P(\text{windy}=\text{false}|\text{yes})=6/9 \sim 0.66$$

$$P(\text{outlook}=\text{overcast}|\text{no})=0/5=0$$
$$P(\text{windy}=\text{false}|\text{no})=2/5=0.4$$

$$P(X = x | Y = y) = \frac{n_c + mp}{n + m}$$

$$P(\text{outlook} = \text{overcast} | \text{yes}) =$$

$$P(\text{outlook} = \text{overcast} | \text{no}) =$$

$$P(X = x | Y = y) = P(X = x, Y = y) / P(Y = y)$$

שלב האימון:

ב. חישוב ההסתברויות המותנות הבדידות:

שימו לב, שיש לנו כאן הסתברות 0

- נשתמש בהחלקה עם הפרמטרים: $p=1/3, m=3$

תזכורת - Smoothing solution

- ❖ Probability estimates are adjusted or *smoothed*.
- ❖ Assumes that each feature is given a prior probability, p , that is assumed to have been previously observed in a “virtual” sample of size m .
- ❖ Usually, in the binary case p is simply assumed to be 0.5

$$P(X = x | Y = y) = \frac{n_c + mp}{n + m}$$

- ❖ n = number of training examples for which $Y = y$
- ❖ n_c = number of examples where $X=x$ and $Y=y$
- ❖ p = a prior estimation for $P(X=x | Y=y)$
- ❖ m = the equivalent sample size

תרגיל 8 -

סימולציית סיווג -

מ"מ רציף - פתרון

הווקטור החדש:

outlook=overcast,
temperature=60,
humidity=62,
windy=false.

$$P(\text{outlook}=\text{overcast}|\text{yes})=4/9 \sim 0.44$$
$$P(\text{windy}=\text{false}|\text{yes})=6/9 \sim 0.66$$

$$P(\text{outlook}=\text{overcast}|\text{no})=0/5=0$$
$$P(\text{windy}=\text{false}|\text{no})=2/5=0.4$$

$$P(X = x | Y = y) = \frac{n_c + mp}{n + m}$$

$$P(\text{outlook} = \text{overcast} | \text{yes}) = \frac{4+1}{9+3} = \frac{5}{12} = 0.4167$$

$$P(\text{outlook} = \text{overcast} | \text{no}) = \frac{0+1}{5+3} = \frac{1}{8} = 0.125$$

$$P(X = x | Y = y) = P(X = x, Y = y) / P(Y = y)$$

שלב האימון:

ב. חישוב ההסתברויות המותנות הבדידות:

שימו לב, שיש לנו כאן הסתברות 0

- נשתמש בהחלקה עם הפרמטרים: $p=1/3, m=3$

תרגיל 8 - סימולציית סיווג – מ"מ רצוף - פתרון

שלב האימון:

ג. חישוב ההסתברויות המותנות הרציפות:

(יש לחשב בנפרד עבור המחלקות השונות)

outlook	temperature	humidity	windy	play
sunny	85	85	false	no
sunny	80	90	true	no
overcast	83	86	false	yes
rainy	70	96	false	yes
rainy	68	80	false	yes
rainy	65	70	true	no
overcast	64	65	true	yes
sunny	72	95	false	no
sunny	69	70	false	yes
rainy	75	80	false	yes
sunny	75	70	true	yes
overcast	72	90	true	yes
overcast	81	75	false	yes
rainy	71	91	true	no

Temperature:

$\mu_{temp_yes} =$, $\sigma_{temp_yes} =$

$\mu_{temp_no} =$, $\sigma_{temp_no} =$

Humidity:

$\mu_{hum_yes} =$, $\sigma_{temp_yes} =$

$\mu_{hum_no} =$, $\sigma_{temp_no} =$

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

תרגיל 8 - סימולציית סיווג – מ"מ רצוף - פתרון

שלב האימון:

ג. חישוב ההסתברויות המותנות הרציפות:

(יש לחשב בנפרד עבור המחלקות השונות)

outlook	temperature	humidity	windy	play
sunny	85	85	false	no
sunny	80	90	true	no
overcast	83	86	false	yes
rainy	70	96	false	yes
rainy	68	80	false	yes
rainy	65	70	true	no
overcast	64	65	true	yes
sunny	72	95	false	no
sunny	69	70	false	yes
rainy	75	80	false	yes
sunny	75	70	true	yes
overcast	72	90	true	yes
overcast	81	75	false	yes
rainy	71	91	true	no

Temperature:

$$\mu_{\text{temp_yes}}=73, \sigma_{\text{temp_yes}}=6.2$$

$$\mu_{\text{temp_no}}=74.6, \sigma_{\text{temp_no}}=8.0$$

Humidity:

$$\mu_{\text{hum_yes}}=79.1, \sigma_{\text{hum_yes}}=10.2$$

$$\mu_{\text{hum_no}}=86.2, \sigma_{\text{hum_no}}=9.7$$

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$s = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

תרגיל 8 - סימולציית סיווג – מ"מ רציף - פתרון

$$P(X_i = x \mid Y = y_k) = \frac{1}{\sigma_{ik}\sqrt{2\pi}} e^{-\frac{(x-\mu_{ik})^2}{2\sigma_{ik}^2}}$$

$$f(\text{temperature} = 60 \mid \text{yes}) =$$

$$f(\text{temperature} = 60 \mid \text{no}) =$$

$$f(\text{humidity} = 62 \mid \text{yes}) =$$

$$f(\text{humidity} = 62 \mid \text{no}) =$$

שלב הסיווג – חישוב פונקצית
הצפיפות עבור המ"מ הרציפים:

temperature=60 humidity=62

חישובנו:

$\mu_{\text{temp_yes}}=73, \sigma_{\text{temp_yes}}=6.2$

$\mu_{\text{temp_no}}=74.6, \sigma_{\text{temp_no}}=8.0$

$\mu_{\text{hum_yes}}=79.1, \sigma_{\text{temp_yes}}=10.2$

$\mu_{\text{hum_no}}=86.2, \sigma_{\text{temp_no}}=9.7$

תרגיל 8 - סימולציית סיווג – מ"מ רציף - פתרון

$$P(X_i = x \mid Y = y_k) = \frac{1}{\sigma_{ik}\sqrt{2\pi}} e^{-\frac{(x-\mu_{ik})^2}{2\sigma_{ik}^2}}$$

$$f(\text{temperature} = 60 \mid \text{yes}) = \frac{1}{6.2\sqrt{2\pi}} e^{-\frac{(60-73)^2}{2(6.2)^2}} = 0.071$$

$$f(\text{temperature} = 60 \mid \text{no}) = \frac{1}{8\sqrt{2\pi}} e^{-\frac{(60-74.6)^2}{2(8)^2}} = 0.0094$$

$$f(\text{humidity} = 62 \mid \text{yes}) = \frac{1}{10.2\sqrt{2\pi}} e^{-\frac{(62-79.1)^2}{2(10.2)^2}} = 0.0096$$

$$f(\text{humidity} = 62 \mid \text{no}) = \frac{1}{9.7\sqrt{2\pi}} e^{-\frac{(62-86.2)^2}{2(9.7)^2}} = 0.0018$$

שלב הסיווג – חישוב פונקצית
הצפיפות עבור המ"מ הרציפים:

temperature=60 humidity=62

חישובנו:

$\mu_{\text{temp_yes}}=73, \sigma_{\text{temp_yes}}=6.2$

$\mu_{\text{temp_no}}=74.6, \sigma_{\text{temp_no}}=8.0$

$\mu_{\text{hum_yes}}=79.1, \sigma_{\text{temp_yes}}=10.2$

$\mu_{\text{hum_no}}=86.2, \sigma_{\text{temp_no}}=9.7$

Classify (X^{new})

$$Y^{new} \leftarrow \arg \max_{y_k} P(Y = y_k) \prod_i P(X_i^{new} | Y = y_k)$$

$$f(\text{temperature} = 60 | \text{yes}) = 0.071$$

$$f(\text{temperature} = 60 | \text{no}) = 0.0094$$

$$f(\text{humidity} = 62 | \text{yes}) = 0.0096$$

$$f(\text{humidity} = 62 | \text{no}) = 0.0018$$

$$P(\text{yes} | \text{outlook}=\text{overcast..}) =$$

הווקטור החדש:

outlook=overcast,
temperature=60,
humidity=62,
windy=false.

תרגיל 8 - סימולציית סיווג – פתרון

שלב הסיווג (המשך):

חישבנו ושלפנו את ההסתברויות המותנות ואת ה-priors:

$$- p(\text{yes}) = 0.643$$

$$- p(\text{no}) = 0.357$$

$$P(\text{outlook}=\text{overcast}|\text{no})=0.125$$

$$P(\text{windy}=\text{false}|\text{no})=0.4$$

$$P(\text{outlook}=\text{overcast}|\text{yes}) \sim 0.4167$$

$$P(\text{windy}=\text{false}|\text{yes}) \sim 0.66$$

כעת נחשב את ההסתברויות:

$$P(\text{no} | \text{outlook}=\text{overcast..}) =$$

Classify (X^{new})

$$Y^{new} \leftarrow \arg \max_{y_k} P(Y = y_k) \prod_i P(X_i^{new} | Y = y_k)$$

הווקטור החדש:

outlook=overcast,
temperature=60,
humidity=62,
windy=false.

תרגיל 8 - סימולציית סיווג – פתרון

שלב הסיווג (המשך):

$$f(\text{temperature} = 60 | \text{yes}) = 0.071$$

$$f(\text{temperature} = 60 | \text{no}) = 0.0094$$

$$f(\text{humidity} = 62 | \text{yes}) = 0.0096$$

$$f(\text{humidity} = 62 | \text{no}) = 0.0018$$

חישבנו ושלפנו את ההסתברויות המותנות ואת ה-priors:

$$- p(\text{yes}) = 0.643$$

$$- p(\text{no}) = 0.357$$

$$P(\text{outlook}=\text{overcast}|\text{no})=0.125$$

$$P(\text{windy}=\text{false}|\text{no})=0.4$$

$$P(\text{outlook}=\text{overcast}|\text{yes}) \sim 0.4167$$

$$P(\text{windy}=\text{false}|\text{yes}) \sim 0.66$$

$$\begin{aligned} P(\text{yes} | \text{outlook}=\text{overcast}..) &= p(\text{yes}) * \\ &p(\text{outlook}=\text{overcast} | \text{yes}) * .. = \\ &0.643 * 0.4167 * 0.667 * 0.071 * 0.0096 = \\ &0.00012 \end{aligned}$$

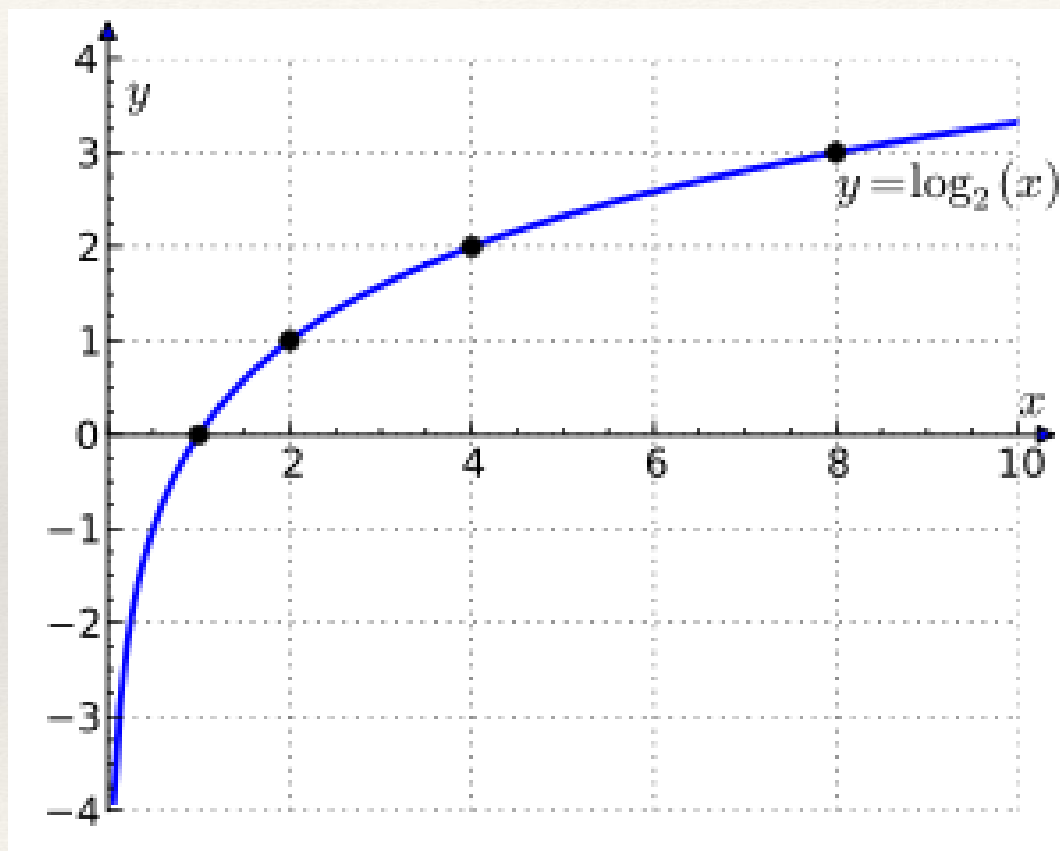
כעת נחשב את ההסתברויות:

$$\begin{aligned} P(\text{no} | \text{outlook}=\text{overcast}..) &= p(\text{no}) * \\ &p(\text{outlook}=\text{overcast} | \text{no}) * .. = \\ &0.357 * 0.125 * 0.4 * 0.0094 * 0.0018 = \\ &0.000000302 \end{aligned}$$



$P(\text{yes})$ is more probable..

פונקציית לוג



תכונות:

- פונקציית לוג של שברים תהיה שלילית, אך היא שומרת על הסדר, והיא גם מונוטונית עולה.
- $\log(x*y) = \log(x) + \log(y)$

לכן, נרצה לעבוד עם חיבור לוגים, במקום מכפלת שברים (של הסתברויות).

מדוע?

Classify (X^{new})

$$Y^{new} \leftarrow \arg \max_{y_k} P(Y = y_k) \prod_i P(X_i^{new} | Y = y_k)$$

הווקטור החדש:

outlook=overcast,
temperature=60,
humidity=62,
windy=false.

תרגיל 8 - סימולציית סיווג – פתרון

שלב הסיווג (המשך):

$$f(\text{temperature} = 60 | \text{yes}) = 0.071$$

$$f(\text{temperature} = 60 | \text{no}) = 0.0094$$

$$f(\text{humidity} = 62 | \text{yes}) = 0.0096$$

$$f(\text{humidity} = 62 | \text{no}) = 0.0018$$

חישבנו ושלפנו את ההסתברויות המותנות ואת ה-priors:

$$- p(\text{yes}) = 0.643$$

$$- p(\text{no}) = 0.357$$

$$P(\text{outlook}=\text{overcast}|\text{no})=0.125$$

$$P(\text{windy}=\text{false}|\text{no})=0.4$$

$$P(\text{outlook}=\text{overcast}|\text{yes}) \sim 0.4167$$

$$P(\text{windy}=\text{false}|\text{yes}) \sim 0.66$$

$$\begin{aligned} P(\text{yes} | \text{outlook}=\text{overcast}..) &= p(\text{yes}) * \\ &p(\text{outlook}=\text{overcast} | \text{yes}) * .. = \log(0.643) + \\ &\log(0.4167) + \log(0.667) + \log(0.071) + \log(0.0096) \\ &= -13.003 \end{aligned}$$

כעת נחשב את ההסתברויות:

$$\begin{aligned} P(\text{no} | \text{outlook}=\text{overcast}..) &= p(\text{no}) * \\ &p(\text{outlook}=\text{overcast} | \text{no}) * .. = \log(0.357) + \\ &\log(0.125) + \log(0.4) + \log(0.0094) + \log(0.0018) \\ &= -21.66 \end{aligned}$$



$P(\text{yes})$ is more probable..