מימוש של מסווג Gaussian Naïve Bayes

הקדמה

בחלק זה של הפרויקט תאמנו מסווג מסוג GNB, שהוא פונקציה $f\colon \mathcal{R}^D\mapsto \{1,2,...,C\}$ המשייכת לכל וקטור במרחב בחלק זה של הפרויקט תאמנו מסווג מסוג GNB, שהוא פונקציה $f\colon \mathcal{R}^D\mapsto \{1,2,...,C\}$ שישמש אתכם לאימון המסווגים המאפיינים $\mathbf{z}\in \mathcal{R}^D$ אחד מ- $\mathbf{z}\in \mathcal{R}^D$ ערכים $\mathbf{z}\in \mathcal{R}^D$ המיינים (סה"כ 784 פיקסלים, כלומר בעבודה זאת מורכב מאוסף דוגמאות, כשכל \mathbf{z} הינו מטריצה ריבועית במימד 28x28 (סה"כ 784 פיקסלים, כלומר בעבודה זאת (D=784 פריט לבוש, ו- \mathbf{z} התיוג של כל דוגמא מייצג את סוג הפריט ע"פ הפירוט להלן (ראו איור 1)

Label 0: T-shirt/top Label 1: Trouser Label 2: Coat Label 3: Sandal Label 4: Ankle boot

אינו כי מסווג מסוג GNB מתבסס על שערוך $maximum\ likelihood$

$$\begin{array}{l} \forall \quad c=1,2,...,C: \\ \\ P(x|y=c,\mu_c,\Sigma_c) = \, \mathcal{N}(x|\mu_c,\Sigma_c) \quad \text{(where } \Sigma_c = \begin{pmatrix} \sigma_{c,1}^2 & 0 & 0 \\ 0 & \ddots & 0 \\ 0 & 0 & \sigma_{c,D}^2 \end{pmatrix} \text{ is diagonal)} \\ \\ P(y=c) = \, \pi_c \\ \\ 0 \leq \pi_c \leq 1 \qquad \text{(and} \qquad \sum_{c=1}^C \pi_c = 1 \text{)} \end{array}$$

וכי בהינתן שיערוכי הפרמטרים של המודל $\Theta = \left\{\pi_{c,ML}, \mu_{c,ML}, \Sigma_{c,ML}^{\square}
ight\}_{c=1}^{C}$ ניתן להשתמש בחוק בייס באופן הבא

$$P(y = c|x, \Theta) = \frac{P(x|y = c, \mu_{c,ML}, \Sigma_{c,ML}^{\square})P(y = c)}{\sum_{c'=1}^{C} P(x|y = c', \mu_{c',ML}, \Sigma_{c',ML}^{\square})P(y = c')}$$

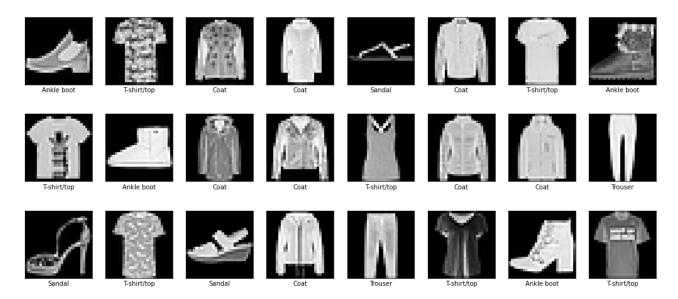
$$= \frac{\mathcal{N}(x|\mu_{c,ML}, \Sigma_{c,ML}^{\square})\pi_{c,ML}}{\sum_{c'=1}^{C} \mathcal{N}(x|\mu_{c',ML}, \Sigma_{c',ML}^{\square})\pi_{c',ML}}$$

$$\propto \mathcal{N}(x|\mu_{c,ML}, \Sigma_{c,ML}^{\square})\pi_{c,ML} \quad \forall \quad c = 1,2,...,C$$

ולסווג כל דוגמא x ע"פ הכלל

(1)
$$\hat{y}(x) = f(x) = \underset{c=1,2,\dots,C}{\operatorname{argmax}} P(y = c | x, \Theta)$$

= $\underset{c=1,2,\dots,C}{\operatorname{argmax}} \mathcal{N}(x | \mu_{c,ML} \Sigma_{c,ML}^{\square}) \pi_{c,ML}$



איור 1: דוגמאות לסוגי פרטי הלבוש.

 $:\mathcal{D}$ עבור הדאטה, $f(\mathbf{x})$ המסווג (accuracy) כמו כן הגדרנו

$$P_c(f, \mathcal{D}) = \frac{1}{N} \sum_{n=1}^{N} \mathbb{I}(y_n == f(x_n))$$

כאשר $\mathbb{I}(A)$ היא פונקציית האינדיקטור

$$\mathbb{I}(A) = \begin{cases} 1 & \text{if } A \text{ is TRUE} \\ & \text{iii} \\ 0 & \text{otherwise} \end{cases}$$

שימו לב:

- בפתרונכם ניתן להשתמש בקוד שכתבתם במהלך הסמסטר בתרגילי הקידוד ובמטלות הבית.
- בנוסף למחברת שתגישו, יש לצרף מסמך מלווה עם התשובות לחלקים העיוניים במטלה ופירוט תוצאות ניסוייכם. בהירות ושלמות המסמך ישפיעו על הציון הסופי.
- 1. (1 נקודות) טענו את הדאטה שקיבלתם, והכינו אותו לאימון המסווג ע"י ייצוג כל תמונה כוקטור במרחב \mathcal{R}^{784} . סדרת האימון מכילה 24,000 דוגמאות, כל אחת כוללת מטריצה במימד 28x28 (המאפיינים 24,000 והתיוגים של כל אחת מהדוגמאות. אוסף כל התמונות נתון לכם כמטריצה תלת-ממדית במימד 28x28x24,000, והתיוגים של כל הדוגמאות כוקטור באורך 24,000.
- א. טענו את הדאטה הנתון בקובץ TrainData.pkl, הכולל את מטריצת התמונות ואת וקטור התיוגים של סדרת המבחן.

- ב. הציגו בתמונות נפרדות דוגמה אחת מכל אחד מהתיוגים הקיימים בדאטה, והוסיפו לתמונה את תיוגה ע"פ הדאטה.
- עם ערכי כל הפיקסלים של הדוגמא ה-n-ית באורה ה-n-ית בה מכילה וקטור באורך 24,000x784, כאשר השורה ה-n-ית בה מכילה וקטור באורך x_n (המתקבל לאחר שכל השורות בה שורשרו זו לזו). np.reshape(...)
 - 2. (34 נקודות) ממשו את אלגוריתם האימון של מסווג GNB עבור מידע רב-ממדי כפי שלמדנו בקורס
 - של כל אחד maximum likelihood א. כתבו פונקציה המקבלת את סדרת האימון, ומחשבת את שערוך $\Theta = \left\{\pi_{c,ML}, \mu_{c,ML}, \Sigma_{c,ML}^{\square}\right\}_{c=1}^{c}$ מפרמטרי המודל
- **הערה ו:** ערכי הפיקסלים בתמונות שקיבלתם מוגדרים כמשתנה מסוג int8. בכדי להימנע מבעיות נומריות בשערוכי הפרמטרים מומלץ ראשית לשמור אותם בפורמט אחר דוגמת float64.

(np.cov אין להשתמש בפונקציות כגון להשתמש בפונקציות כגון enp.cov הערה וו: בתשובתכם יש לבצע חישוב מלא של השערוכים

ב. כתבו פונקציה נוספת המקבלת את כל פרמטרי המודל ששערכתם ודוגמאות לא מתויגות, ומחזירה את סיווג כל אחת מהתמונות ע"פ

$$\hat{y}(x) = \underset{c=1,2,\dots,C}{\operatorname{argmax}} P(y = c | x, \Theta)$$

- 784 יצגו כל אחד מהוקטורים הממוצעים ששערכתם כמטריצה ריבועית (כלומר יצגו כל וקטור באורך 784 כמטריצה במימד 28x28) והציגו אותו כתמונה. הוסיפו לתמונה את התיוג הרלוונטי ודונו בקצרה בתוצאות שקיבלתם.
 - 3. (5 נקודות) חשבו והציגו עבור סדרת האימון את
 - $P_c(f,\mathcal{D})$ א. דיוק המסווג
 - ב. מטריצת הערבול (confusion matrix)

דונו בתוצאות שקיבלתם, תוך התייחסות מפורשת לכל אחד מהמדדים הנ"ל.

4. (10 נקודות) חיזרו על שאלה 3 עבור סדרת המבחן הנתונה בקובץ TestData.pkl. דונו <u>בהרחבה</u> בתוצאות שקיבלתם, תוך התייחסות <u>מפורטת</u> לדומה ולשונה בין הביצועים שמדדתם עבור כל אחת מהסדרות.

בהצלחה!