

מבוא ללמידה חישובית – מבחן מועד ב' סמסטר א' תשע"ה (2014/5)

בית הספר למדעי המחשב, אוניברסיטת תל-אביב

מרצים: פרופ' ליאור וולף, פרופ' ערן הלפרין
מתרגל: רגב שוייגר

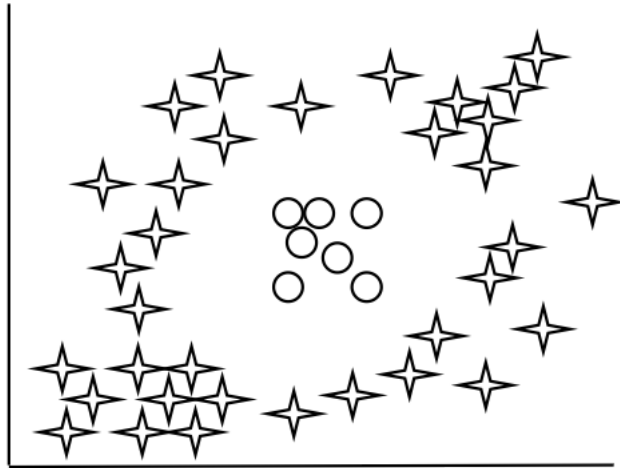
4.9.2015

הוראות:

1. מומלץ לקרוא את כל ההנחיות והשאלות בתחילת המבחן לפני תחילת כתיבת התשובות.
2. משך הבחינה – **שלוש שעות**. לא תינתן כל הארכה נוספת.
3. חומר עזר מותר: דף נוסחאות בגודל A4.
4. יש לענות על השאלות במקום המיועד לכך בטופס השאלון (טופס זה). מחברות הבחינה לא ייקראו, ותשמנה כטיוטה בלבד.
5. יש למלא בכל דף של השאלון מספר ת.ז. ומספר מחברת.
6. במבחן 6 שאלות:
 - הניקוד לכל שאלה מופיע ליד מספר השאלה.
 - יש לענות תשובות ברורות, ענייניות ותמציתיות.
7. מותר להשתמש בכל טענה שהוכחה בכיתה (בהרצאה, בתרגול, או בתרגיל בית) בתנאי שמצטטים אותה במדויק. טענות אחרות (כאלה שהוכחו בספר, בהרצאות מהסמסטר הקודם, וכו') יש להוכיח.
8. אם לא נאמר אחרת, יש להניח שדגימות במדגם נוצרות באופן בלתי תלוי ומאותה התפלגות (i.i.d).

1 שאלה 1 - 15 נקודות

נתון המדגם הדו-מימדי הבא, בו הנקודות מסווגות לשתי מחלקות:

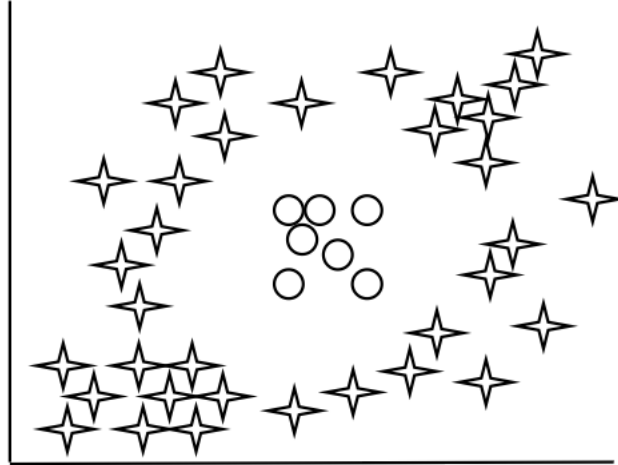


עבור כל אחד מהאלגוריתמים הבאים, קבע/י האם ניתן להריץ אותו עד לקבלת מסווג עם שגיאת למידה אפס, על המדגם הנתון. אם כן, צייר/י קו הפרדה מתאים למסווג המתקבל. אם לא, הסבר/הסבירי מדוע.

1. עצי החלטה, כאשר ה-classifiers מתוכם האלגוריתם בוחר הם כל ה-decision stumps המקבילים לצירים (כלומר, עבור כל מימד i , וסף a , ניתן לבחור classifier שמסווג על פי תוצאת השוואה $x_i < a$).

לא יכול להגיע לשגיאה אפס. הסבר:

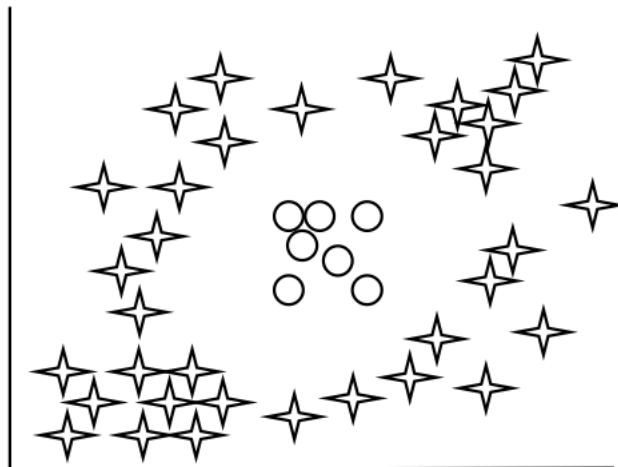
יכול להגיע לשגיאה אפס. קו הפרדה של המסווג:



.2 .3-Nearest Neighbor

לא יכול להגיע לשגיאה אפס. הסבר:

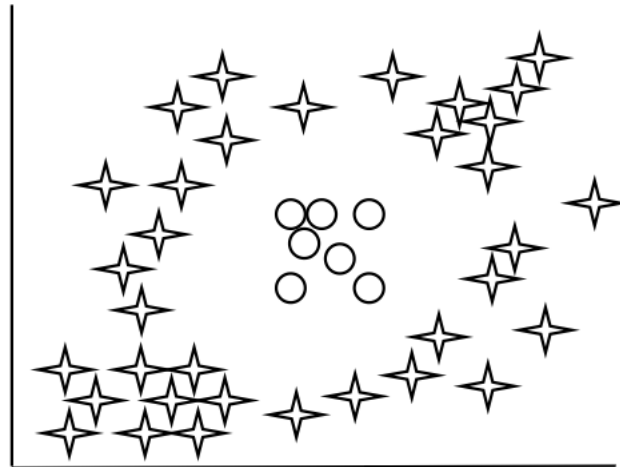
יכול להגיע לשגיאה אפס. קו ההפרדה של המסווג:



3. SVM עם Polynomial kernel, $d=2$.

לא יכול להגיע לשגיאה אפס. הסבר:

יכול להגיע לשגיאה אפס. קו ההפרדה של המסווג:



2 שאלה 2 - 15 נקודות

1. מפעילים את אלגוריתם PCA, על מדגם אימון מסוים, לקבלת מערכת צירים בעלת k וקטורים (כאשר אנחנו כוללים שלב מירכוז מקדים, בו אנו מזיזים את כל נקודות המדגם באותו וקטור, כך שממוצע כל קואורדינטה הוא 0). כעת מפעילים על כל נקודה במדגם טרנספורמציה נתונה T , ועל המדגם החדש מפעילים PCA בשנית לקבלת k וקטורים נוספים. האם שתי מערכות הצירים זהות (עד כדי החלפת סימן אפשרית של כל וקטור עצמי)? נמק/י את תשובתך בקצרה.

(א) $T(x) = x + x_0$, כאשר x_0 הוא וקטור.

מערכות זהות

לא בהכרח זהות

הסבר: כן, כיוון ששלב המירכוז יבטל את השפעת הטרנספורמציה.

(ב) $T(x) = ax$, כאשר $a \neq 0$ הוא סקלר.

מערכות זהות

לא בהכרח זהות

הסבר: כן, מכיוון שהוקטורים העצמיים של XX^T ושל $(aX)(aX)^T$ זהים. אינטואיטיבית, מתיחה של כל הצירים בשיעור אחיד לא תשנה אז זהות תת-המרחב המתאים ביותר.

(ג) $T(x) = Dx$, כאשר D היא מטריצה אלכסונית, ללא 0 על האלכסון.

מערכות זהות

לא בהכרח זהות

הסבר: לא בהכרח, מכיוון שהוקטורים העצמיים של XX^T ושל $(DX)(DX)^T$ לא בהכרח זהים. אינטואיטיבית, באמצעות D ניתן למתוח ולכווץ כל קואורדינטה בנפרד ולתת לה חשיבות עודפת או פחותה, ובכך להשפיע על תת-המרחב הקרוב ביותר אל המדגם.

(ד) $T(x) = Ux$, כאשר U היא מטריצה אורתונורמלית.

מערכות זהות

□ לא בהכרח זהות

הסבר: לא בהכרח, מכיוון שהוקטורים העצמיים של XX^T ושל $(UX)(UX)^T$ לא בהכרח זהים. אפשר באמצעות U לשובב את מערכת הצירים ובכך לשנות אותה.

2. מריצים אלגוריתם EM ללמידת הפרמטרים של מודל Gaussian Mixture Model על מדגם נתון. האלגוריתם מאותחל כך: וקטור ההתפלגות המגדיר את ההסתברויות לבחור בכל גאוסיאן מאותחל להתפלגות האחידה. השונות ההתחלתית של של גאוסיאן היא 1, והתוחלת ההתחלתית של הגאוסיאן ה- i היא הנקודה ה- i במדגם (יש פחות צבירים מאשר נקודות). לוקחים נקודה קבועה z ומחשבים את ההסתברות האפוסטריורית שלה להשתייך לכל אחד מהצבירים, לאחר הרצת האלגוריתם. כעת מפעילים על כל נקודה במדגם טרנספורמציה נתונה T , ועל המדגם החדש מפעילים את האלגוריתם בשנית. כעת מחשבים את ההסתברות האפוסטריורית של $T(z)$ להשתייך לכל אחד מהצבירים. בשתי הפעמים מריצים את האלגוריתם אותו מספר איטרציות.

האם וקטור ההסתברויות עבור z זהה לוקטור ההסתברויות עבור $T(z)$?

(א) $T(x) = x + x_0$, כאשר x_0 הוא וקטור.

בהכרח זהה

לא בהכרח זהה

הסבר: כן, כיוון שהמקום היחיד בו זה יבוא לידי ביטוי הוא באומדן התוחלות, שיוזז גם הוא, אבל כל שאר החישובים יישארו זהים. אפשר להחליף בכל החישובים $T(\mu) \rightarrow \mu$ והכל ישאר זהה.

(ב) $T(x) = ax$, כאשר $a \neq 0$ הוא סקלר.

בהכרח זהה

לא בהכרח זהה

הסבר: לא בהכרח זהה. כדי שיהיה זהה, יש לשנות גם את השונות ההתחלתית.

(ג) $T(x) = Dx$, כאשר D היא מטריצה אלכסונית, ללא 0 על האלכסון.

בהכרח זהה

לא בהכרח זהה

הסבר: לא.

(ד) $T(x) = Ux$, כאשר U היא מטריצה אורתונורמלית.

בהכרח זהה

תעודת זהות:
מספר מחברת:

□ לא בהכרח זהה

הסבר: כן, סיבוב אחיד לכל לא ישנה את ההסתברויות. אפשר להחליף בכל החישובים $\mu \rightarrow U\mu$ והכל ישאר זהה.

3 שאלה 3 - 20 נקודות

1. ברגרסיה לינארית סטנדרטית, נתון המודל הבא: יהיו וקטורים $\mathbf{x}_1, \dots, \mathbf{x}_p \in \mathbb{R}^n$ נתונים, ותהא X המטריצה שעמודותיה הם $\mathbf{x}_1, \dots, \mathbf{x}_p$. יהי וקטור של משתנים מקריים שעבורו לכל $i = 1, \dots, n$, $e_i \sim N(0, \sigma^2)$. כאשר σ^2 היא שונות קבועה ונתונה. יהיו וקטור פרמטרים, ונגדיר:

$$\begin{aligned}\mathbf{y} &= X\mathbf{a} + \mathbf{e} \\ &= a_1 \cdot \mathbf{x}_1 + \dots + a_p \cdot \mathbf{x}_p + \mathbf{e}\end{aligned}$$

מהי ההתפלגות הרב מימדית של הוקטור \mathbf{y} ?

$$\begin{aligned}\mathbf{e} &\sim N(\mathbf{0}_n, \sigma^2 I_n) \Rightarrow \\ \mathbf{y} = X\mathbf{a} + \mathbf{e} &\sim N(X\mathbf{a}, \sigma^2 I_n)\end{aligned}$$

2. כעת נכליל את המודל. בנוסף להגדרות סעיף 1, יהיו $\mathbf{z}_1, \dots, \mathbf{z}_q \in \mathbb{R}^n$ וקטורים נתונים, ותהא Z המטריצה שעמודותיה הם $\mathbf{z}_1, \dots, \mathbf{z}_q$. יהיו $b_1, \dots, b_q \sim N(0, \tau^2)$ משתנים מקריים, ונסמן $\mathbf{b} = (b_1, \dots, b_q)^T$. השונות τ^2 נתונה. כעת:

$$\begin{aligned}\mathbf{y} &= X\mathbf{a} + Z\mathbf{b} + \mathbf{e} \\ &= a_1 \cdot \mathbf{x}_1 + \dots + a_p \cdot \mathbf{x}_p + b_1 \cdot \mathbf{z}_1 + \dots + b_q \cdot \mathbf{z}_q + \mathbf{e}\end{aligned}$$

מודל כזה נקרא Linear Mixed Model (LMM). כעת מהי ההתפלגות של \mathbf{y} ?

$$\begin{aligned}\mathbf{b} &\sim N(\mathbf{0}_q, \tau^2 I_q) \Rightarrow \\ Z\mathbf{b} &\sim N(\mathbf{0}_n, Z \cdot \tau^2 I_q \cdot Z^T) = N(\mathbf{0}_n, \tau^2 Z Z^T) \Rightarrow \\ \mathbf{y} = X\mathbf{a} + Z\mathbf{b} + \mathbf{e} &\sim N(X\mathbf{a}, \sigma^2 I_n + \tau^2 Z Z^T)\end{aligned}$$

4 שאלה 4 - 20 נקודות

באלגוריתם SVM, נתון לנו מדגם מסווג (\mathbf{x}_n, y_n) , כאשר $\mathbf{x}_n \in \mathbb{R}^d$ ו- $y_n \in \{+1, -1\}$.
אנו פותרים את הבעיה הבאה:

$$\min_{\mathbf{w}, b} \frac{1}{2} \mathbf{w}^T \mathbf{w}$$
$$\text{s.t. } y_n(\mathbf{w}^T \mathbf{x}_n + b) \geq 1 \quad \forall n = 1, \dots, N$$

כאשר $\mathbf{w} \in \mathbb{R}^d$ הוא וקטור המשקולות, $b \in \mathbb{R}$ הוא קבוע ההזזה. כזכור, את הבעיה הזו ניתן להמיר לבעיה דואלית ולפתור אותה באופן שקול. נסמן ב- α_i את כופלי לגרנז' של הלגרנז'יאן בבעיה הדואלית.

1. הראה/י כי מתקיימות שתי המשוואות:

$$\sum_{n=1}^N \alpha_n y_n = 0, \quad \mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$$

ראינו בשיעור. הלגרנז'יאן הוא

$$L(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{n=1}^N \alpha_n (y_n(\mathbf{w}^T \mathbf{x}_n + b) - 1)$$

גזירה נותנת:

$$\frac{d}{db} L = - \sum_{n=1}^N \alpha_n y_n = 0 \implies \sum_{n=1}^N \alpha_n y_n = 0$$

$$\nabla_{\mathbf{w}} L = \mathbf{w} - \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n = 0 \implies \mathbf{w} = \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$$

2. נניח כי עבור i מסוים, \mathbf{x}_i הוא support vector. הראי כיצד ניתן לחשב את b כפונקציה של $\mathbf{x}_n, y_n, \alpha_n$.

אם זה וקטור תומך אז מתקיים

$$\mathbf{x}_i \cdot \mathbf{w} + b = y_i \Rightarrow b = y_i - \mathbf{x}_i \cdot \mathbf{w} = y_i - \mathbf{x}_i \cdot \sum_{n=1}^N \alpha_n y_n \mathbf{x}_n$$

3. נזכיר כי עבור i מסוים, \mathbf{x}_i הוא support vector אם $\alpha_i > 0$. השתמש/י בשני הסעיפים הקודמים כדי להוכיח כי:

$$\sum_{n=1}^N \alpha_n = \|\mathbf{w}\|^2$$

אם \mathbf{x}_i הוא לא וקטור תומך אז $\alpha_i = 0$. אז:

$$\begin{aligned} b &= y_i - \mathbf{x}_i \cdot \sum_{j=1}^N \alpha_j y_j \mathbf{x}_j \Rightarrow \\ 0 &= \sum_{i=1}^N \alpha_i y_i \cdot b = \sum_{i=1}^N \alpha_i y_i^2 - \sum_{i=1}^N \sum_{j=1}^N \alpha_i y_i \alpha_j y_j \mathbf{x}_i \cdot \mathbf{x}_j \Rightarrow \\ &= \sum_{n=1}^N \alpha_n - \|\mathbf{w}\|^2 \end{aligned}$$

5 שאלה 5 - 15 נקודות

1. נתונים $x_1, \dots, x_n \in \mathbb{R}$ הוכח/הוכיחי כי המספר μ המביא למינימום את:

$$\sum_{i=1}^n |x_i - \mu|$$

הוא החציון של המדגם.

המספר שמביא למינימום את סכום (או ממוצע) מרחקי l_1 מהמדגם הוא החציון של המדגם. ניתן גם להוכיח זאת ישירות על ידי גזירה, תוך התחשבות בנקודות אי הגזירות של הפונקציה.

2. פונקצית הצפיפות של התפלגות לפלס החד מימדית $Laplace(\mu, \sigma)$ היא

$$f(x) = \frac{1}{2\sigma} e^{-\frac{|x-\mu|}{\sigma}}$$

נתון מדגם $x_1, \dots, x_n \sim Laplace(\mu, \sigma)$ מהו אומד הנראות המירבית של הפרמטרים μ, σ ?

$$\begin{aligned} \ell(\mu, \sigma; x_1, \dots, x_n) &= \sum_{i=1}^n \log \left(\frac{1}{2\sigma} \right) - \frac{|x_i - \mu|}{\sigma} \\ &= -n \log(2\sigma) - \frac{1}{\sigma} \sum_{i=1}^n |x_i - \mu| \end{aligned}$$

אומד הנראות המירבית של μ הוא זה שמביא למקסימום את הביטוי לעיל. כאשר σ נתונה, הוא זה שמביא למינימום את סכום (או ממוצע) מרחקי l_1 מהמדגם - כלומר החציון של המדגם. כלומר,

$$\hat{\mu} = \text{median}(x_1, \dots, x_n)$$

גזירה לפי σ נותנת שהאומד שלו הוא ממוצע נורמת ה- l_1 מהחציון.

$$\begin{aligned} \frac{\partial \ell}{\partial \sigma} &= -\frac{n}{\sigma} + \frac{1}{\sigma^2} \sum_{i=1}^n |x_i - \mu| = 0 \Rightarrow \\ \hat{\sigma} &= \frac{1}{n} \sum_{i=1}^n |x_i - \hat{\mu}| \end{aligned}$$

3. נתונים $x_1, \dots, x_n \in \mathbb{R}$. בנוסף נתונות משקולות חיוביות w_1, \dots, w_n . נגדיר את $\mu_{(w_1, \dots, w_n)}(x_1, \dots, x_n)$ בתור המספר μ המביא למינימום את:

$$\sum_{i=1}^n w_i |x_i - \mu|$$

מהו $\mu_{(w_1, \dots, w_n)}(x_1, \dots, x_n)$?

אם גוזרים מקבלים שהנגזרת בין כל שתי נקודות היא סכום המשקולות שמתאימות לנקודות הגדולות מהנק' פחות סכום המשקולות שמתאימות לנקודות הקטנות מהנק'. באופן כללי לא ברור שיש שוויון ולכן שיש מינימום באינטרוולים בין הנקודות. לכן המינימום מתקבל באחת הנקודות המקוריות, בהן הסימן סכום המשקולות מצד ימין פחות סכום המשקולות מצד שמאל מתהפך.

6 שאלה 6 - 20 נקודות

פונקצית הצפיפות של התפלגות לפלס החד מימדית $Laplace(\mu, \sigma)$ היא

$$f(x) = \frac{1}{2\sigma} e^{-\frac{|x-\mu|}{\sigma}}$$

נתון מדגם הנלקח מתוך תערובת של התפלגויות לפלס חד מימדיות. כלומר, נתונות k התפלגויות לפלס $f_1 = Laplace(\mu_1, 1), \dots, f_k = Laplace(\mu_k, 1)$. בנוסף נתון וקטור הסתברות $p = (p_1, \dots, p_k)$ הקובע את ההסתברות לדגום מתוך כל אחת מההתפלגויות. עבור כל נקודת מדגם $i = 1, \dots, n$, תחילה מגרילים ערך z_i לפי ההתפלגות p כדי לבחור מאיזו התפלגות לפלס לדגום. ואז, הערך x_i נדגם מתוך ההתפלגות $Laplace(\mu_{z_i}, 1)$.

מהו אלגוריתם EM המתאים לאמידת הפרמטרים μ_i, p_i ?

אפשר להשתמש בתוצאות שאלה 5 ללא הוכחה.

הפיתוח דומה לאלגוריתם עבור תערובת גאוסיאניים, כאשר אומדי הנראות המירבית בשלב M -שונים. המשתנים החבויים הם Z_i . ההסתברות האפוסטריורית לקבל $z_i = j$ באיטרציה $t+1$ היא

$$a_{i,j}^{t+1} = \Pr(Z_i = j | x_1, \dots, x_n) = \frac{p_j^t f_j^t(x_i)}{\sum_{m=1}^k p_m^t f_m^t(x_i)}$$

כאשר $\mu_i^t, \sigma_i^t, p_i^t$ הם ערכי הפרמטרים המשוערכים באיטרציה t , ו- f_j^t היא התפלגות לפלס עם הפרמטרים μ_i^t, σ_i^t . בשלב M -אנו רוצים להביא למקסימום את:

$$Q = \sum_{i=1}^n \sum_{j=1}^k a_{i,j}^{t+1} \log(p_j f_j(x_i))$$

הפתרון עבור p_j זהה למקרה של תערובת גאוסיאניים:

$$p_j^{t+1} = \frac{1}{n} \sum_{i=1}^n a_{i,j}^{t+1}$$

אם ננסה להביא למקסימום את Q על פי μ_j , נקבל שכיוון ש- p_j קבועים, באופן שקול

נראה להביא למקסימום את

$$\begin{aligned} Q &\propto \sum_{i=1}^n \sum_{j=1}^k a_{i,j}^{t+1} \log(f_j(x_i)) \\ &= \sum_{j=1}^k \left(\sum_{i=1}^n a_{i,j}^{t+1} \log\left(\frac{1}{2}\right) - \sum_{i=1}^n a_{i,j}^{t+1} |x_i - \mu_j| \right) \\ &\propto \sum_{j=1}^k \sum_{i=1}^n a_{i,j}^{t+1} |x_i - \mu_j| \end{aligned}$$

כלומר יש כאן k בעיות של חציון משוקלל, והאומדנים הם בסימוני סעיף 3 שאלה 5:

$$\mu_j^{t+1} = \mu_{(a_{1,j}^{t+1}, \dots, a_{n,j}^{t+1})}(x_1, \dots, x_n)$$

תעודת זהות:
מספר מחברת:
