



מבוא למערכות לומדות (236756)

סמסטר חורף תשפ"ב – 10 בפברואר 2022

מרצה: ד"ר יונתן בלינקוב

מבחן מסכם מועד א'

הנחיות הבחינה:

- **משך הבחינה:** 3 שעות.
- **חומר עזר:** המבחן בחומר סגור (ללא ספרים, מחברות, דפי נוסחאות).
- אין צורך במחשבון.
- מותר לכתוב בעט או בעיפרון, כל עוד הכתב קריא וברור.
- מותר לענות בעברית או באנגלית.
- יש לכתוב את התשובות **על גבי שאלון זה** בכתב יד קריא. תשובה בכתב יד לא קריא – לא תיבדק.
- במבחן 16 עמודים ממוספרים סה"כ, כולל עמוד שער זה שמספרו 1 ושלושה עמודי טיוטה בסוף הגליון.
- נא לכתוב רק את המבוקש ולצרף הסברים קצרים עפ"י ההנחיות.
- **בתום המבחן יש להגיש את שאלון זה בלבד.**

מבנה הבחינה:

- **חלק א' [75 נק']:** 4 שאלות פתוחות.
- **חלק ב' [25 נק']:** 5 שאלות סגורות (אמריקאיות) [כל אחת 5 נק'].

בהצלחה!

חלק א' – שאלות פתוחות [75 נק']

שאלה 1 [16 נק']

חוקרת מהטכניון עובדת על בעיית סיווג בינארי כלשהי. ברשותה dataset שבו $m = 150$ דוגמאות שונות (distinct).

החוקרת הריצה שלושה מודלים, ולכל מודל ביצעה hyperparameter tuning:

(a) מודל: kNN (נק' נחשבת שכנה של עצמה), היפרפרמטר: מספר השכנים k , טווח: 1 עד 97.

(b) מודל: עץ החלטה, היפרפרמטר: עומק מירבי, טווח: 1 עד 40.

(c) מודל: Kernel SVM, היפרפרמטר: חזק הרגולריזציה λ , טווח: 10^{-3} עד 10^7 .

overfit → under
complex → simple

under → over
simple → complex

how

hard → soft
complex → simple
overfit → underfit

לכל מודל, היא ציירה גרף של דיוק האימון ודיוק ההכללה (בעזרת 5-fold cross validation) בציר y כפונקציה של ערך ההיפרפרמטר בציר x (הערכים גדלים משמאל לימין).

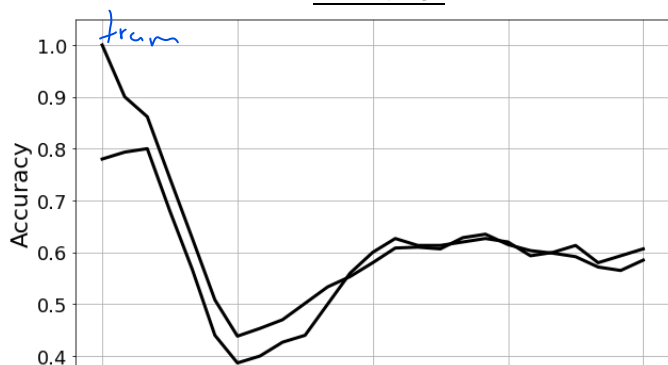
בעקבות תקלה, הכיתוב על ציר x נמחק מכל הגרפים.

בנוסף, הגרפים נשמרו בטעות בשחור לבן, כך שלא ניתן להבדיל בקלות בין דיוק האימון לדיוק ההכללה.

א. [8 נק'] במקום המתאים מתחת כל תרשים, כתבו את האות שמתאימה למודל ולהיפרפרמטר שיצרו אותו.

הערה: השאלה לא מוגדרת היטב ובדיעבד היו שתי תשובות נכונות לשאלה זו.

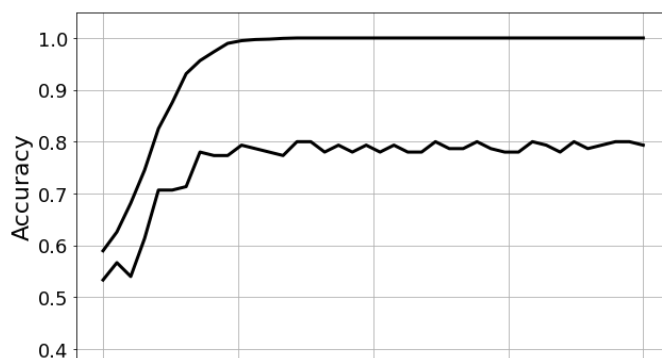
תרשים ב'



C/A

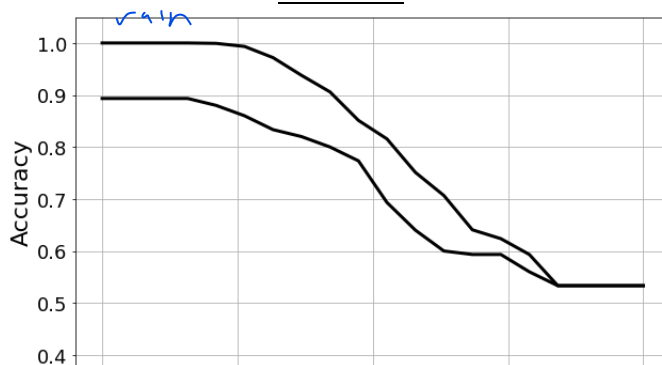
A

תרשים א'



B

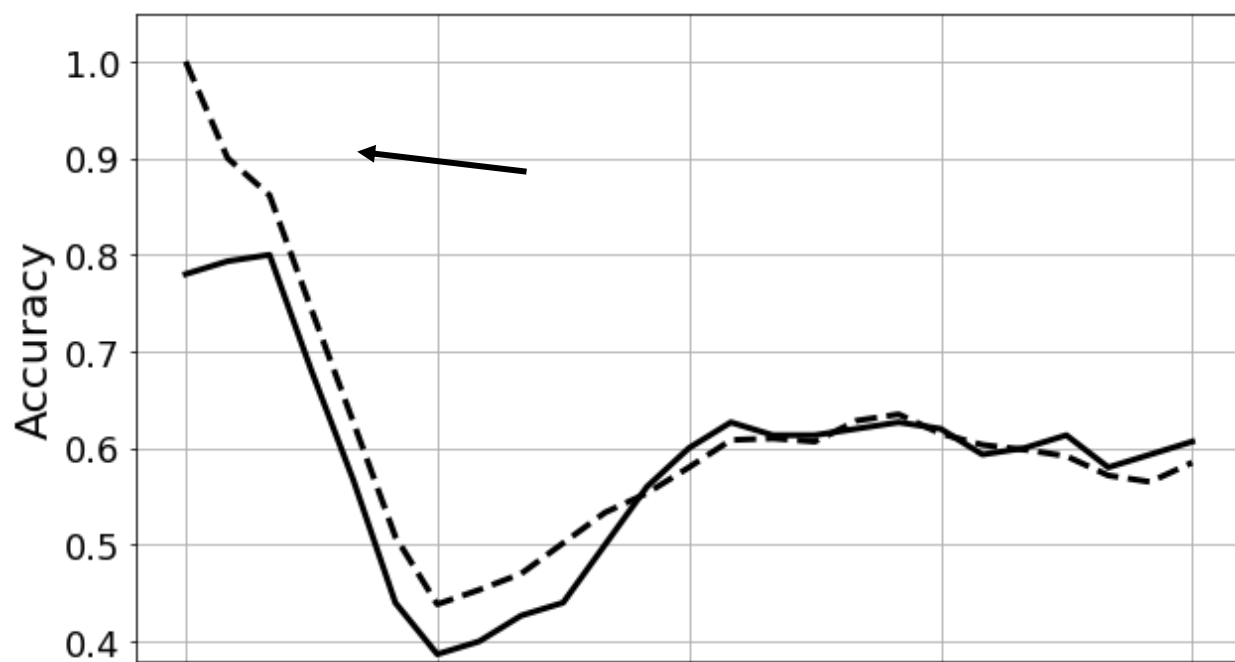
תרשים ג'



A/k

C

ב. [8 נק'] להלן תרשים ב' מוגדל.



הסתכלו על העקומה המקווקות שבתרשים (מסומנת בחץ).

האם העקומה מתארת את דיוק האימון או את דיוק ההכללה?

הסבירו בקצרה. התבססו על התרשים ועל מאפייני המודל שיצר את עקומה זו (מבין שלושת המודלים).

תשובה תמציתית:

אימון

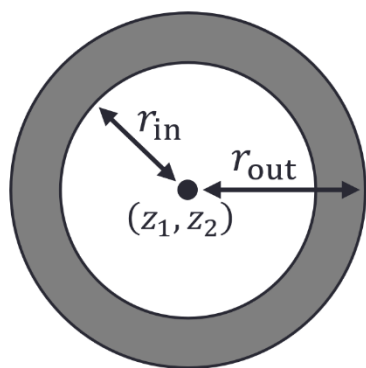
שאלה 2 – PAC learning [21 נק']

תהי \mathcal{H} מחלקת היפותזות של Bagels/donuts בדו-ממד:

$$\mathcal{H} = \{h_\theta: \mathbb{R}^2 \rightarrow \pm 1 \mid \theta = (z_1, z_2, r_{\text{out}}, r_{\text{in}}), r_{\text{out}} > r_{\text{in}} \geq 0\}$$

כאשר היפותזה בודדת מוגדרת באופן הבא:

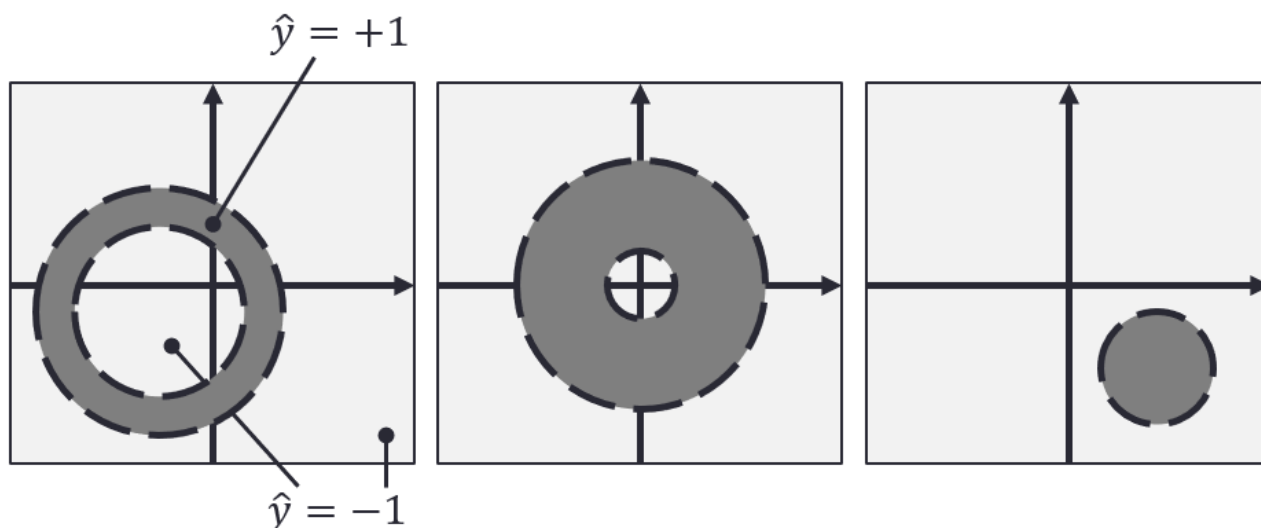
$$h_\theta(x) = \begin{cases} +1, & r_{\text{out}} \geq \sqrt{(x_1 - z_1)^2 + (x_2 - z_2)^2} \geq r_{\text{in}} \\ -1, & \text{otherwise} \end{cases}$$



דגשים לגבי כל היפותזה $h_\theta \in \mathcal{H}$:

- המרכזים של המעגלים משותפים ולא בהכרח בראשית הצירים.
- הרדיוס של המעגל הפנימי יכול להיות אפס.
- השטח שבתוך ה-donut לא יכול להיות אפס.
- האזור בין שני המעגלים מסווג כחיובי, והאזורים האחרים כשליליים.
- מדובר אך ורק במעגלים ולא באליפסות.

דוגמה לשלוש היפותזות מתוך \mathcal{H} :



א. [3 נק'] להלן ההגדרה של "ניתן". השמטנו מההגדרה את הכמתים.

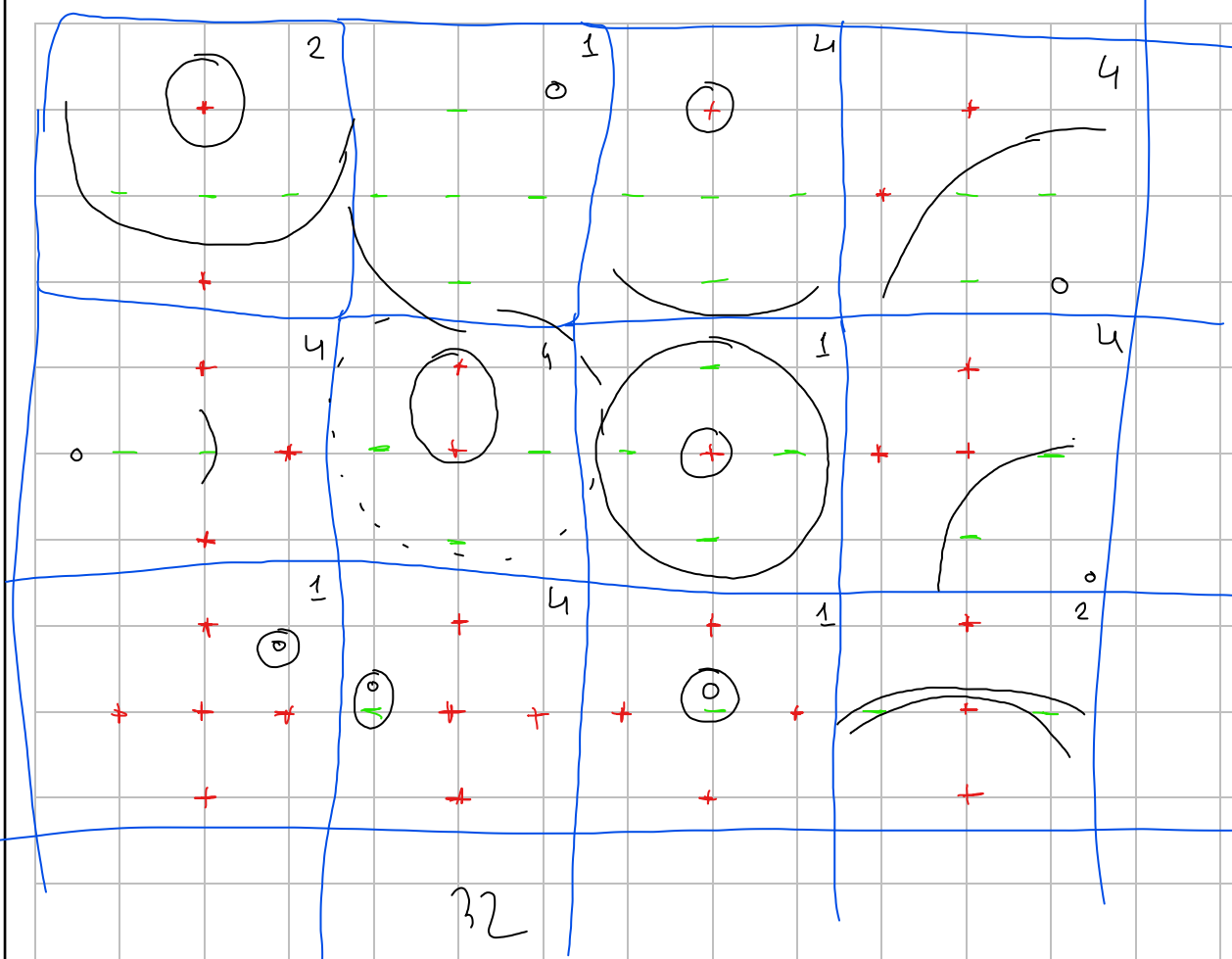
השלימו את שלושת הכמתים החסרים. בכל מקום כתבו האם חסר בהגדרה \forall או \exists .

$$\mathcal{H} \text{ shatters } \mathcal{C} \Leftrightarrow \underbrace{\forall}_{\text{השלימו}} y_1, \dots, y_{|\mathcal{C}|} \in \mathcal{Y}: \underbrace{\exists}_{\text{השלימו}} h \in \mathcal{H}: \underbrace{\forall}_{\text{השלימו}} x_i \in \mathcal{C}: h(x_i) = y_i$$

ב. [13 נק'] כתבו את החסם התחתון ההדוק ביותר שתוכלו למצוא לממד ה-VC: $VCdim(\mathcal{H}) \geq$ 5

הוכיחו את החסם התחתון שכתבתם (אין להוכיח שוויון). יש לכתוב הסבר מילולי תמציתי ולצרף תרשימים נדרשים.

הוכחה (לרשותכם דפי טיוטה בסוף הגיליון):



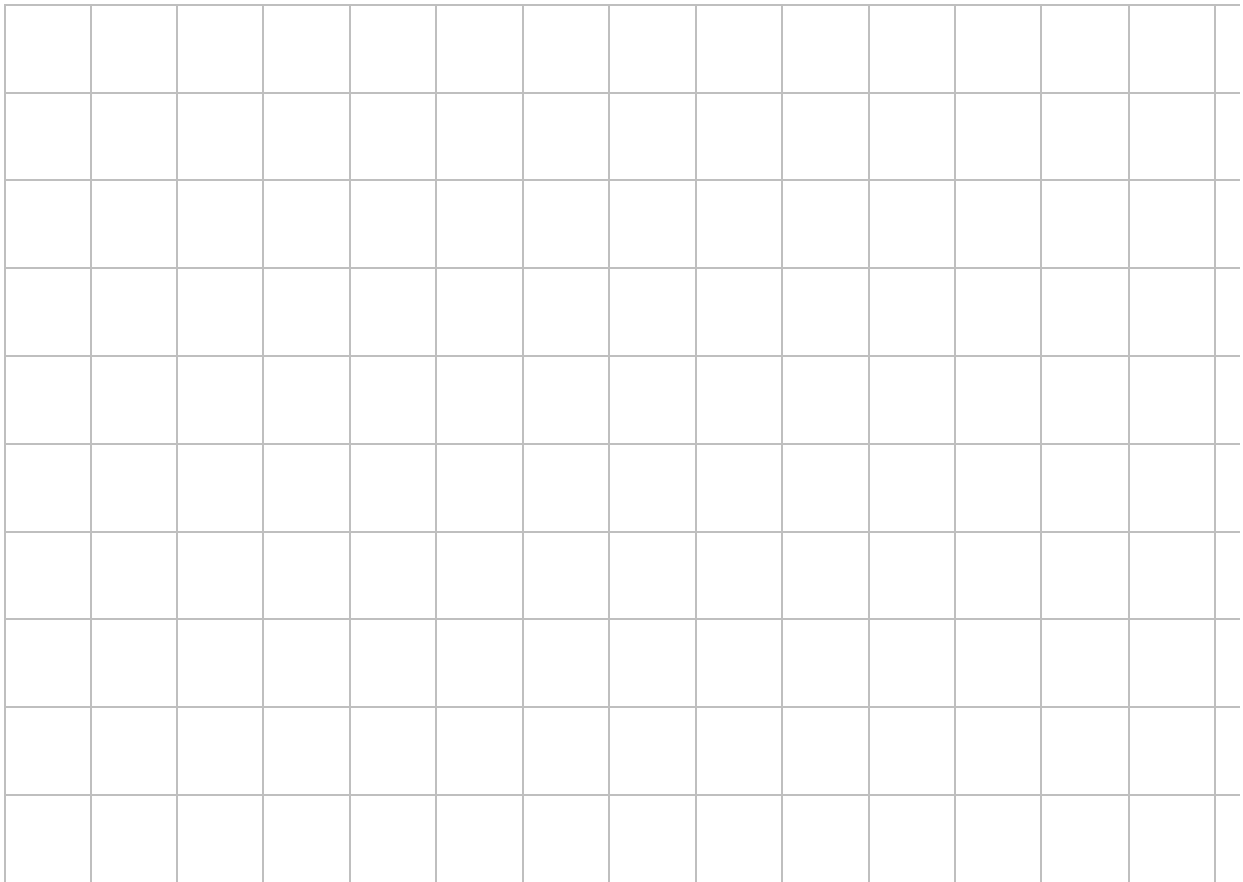
5
10
10
5
1
1
2
20
(2) = 10
24
28
3
1 ✓
2 ✓
3 ✓
5.2.3 = 10
7.2

כאמור הפעמי בין ה + δ - אב אחרו עוקרין
3 מ"נ ש במקום + כתבתי -
וב מקום - כרבוני +

32 קומה סה"כ

כאמור 12 מקום "שאלה"

המשך סעיף ב':



$\frac{1}{2}$

$$\frac{10}{20}$$

$$\frac{5}{10}$$

ג. [5 נק'] חוקרת וחוקר רוצים לאמן מודל סיווג בינארי.

החוקרת משתמשת במחלקת ההיפותזות \mathcal{H} שהגדרנו.

החוקר משתמש במחלקת היפותזות של donuts שמרכזם בראשית הצירים, משמע:

$$\mathcal{H}' = \{h_\theta \mid \theta = (0, 0, r_{\text{out}}, r_{\text{in}}), r_{\text{out}} > r_{\text{in}} \geq 0\} \subset \mathcal{H}$$

מי צפוי להזדקק לפחות דוגמאות בתהליך הלמידה ע"מ להבטיח (במונחי PAC) שגיאת הכללה $\epsilon = 0.1$? נמקו בקצרה.

תשובה תמציתית:

$$VC(\mathcal{H}') \leq VC(\mathcal{H})$$

$$\frac{VC(\mathcal{H})}{\epsilon} \leq \frac{VC(\mathcal{H}') + \log(\frac{1}{\epsilon})}{\epsilon} \Rightarrow VC(\mathcal{H}') \leq VC(\mathcal{H})$$

החוקר, יצחק, שמואל, אהרון

שאלה 3 – רגרסיה ליניארית ו-Generative models [21 נק']

נתון דאטה $S = \{(x_i, y_i)\}_{i=1}^m$ שהגיע ממודל ליניארי $y_i = \mathbf{w}^\top x_i + \varepsilon_i$ עם רעש אקראי מפילוג i.i.d. נורמלי: $\varepsilon_i \sim \mathcal{N}(0, 1)$.

שימו לב: הדוגמאות $x_i \in \mathbb{R}^d$ והתייגים $y_i \in \mathbb{R}$ נתונים. וקטור המשקלים $\mathbf{w} \in \mathbb{R}^d$ לא ידוע ואותו אנו רוצים ללמוד.

תזכורת: הוכחנו שתחת הנחות אלה ה-likelihood שווה ל:

$$L(\mathbf{w} | \{x_i, y_i\}_i) = \Pr(\{x_i, y_i\}_i | \mathbf{w}) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{1}{2}(\mathbf{w}^\top x_i - y_i)^2\right\}$$

א. [5 נק'] הוכיחו שתחת הנחות השאלה, בעיית ה-LS ללא רגולריזציה, משמע $\arg\min_{\mathbf{w}} \frac{1}{m} \sum_{i=1}^m (\mathbf{w}^\top x_i - y_i)^2$,

שקולה לבעיית ה-MLE הבאה: $\arg\max_{\mathbf{w}} L(\mathbf{w} | \{x_i, y_i\}_i)$.

הוכחה ע"י פיתוח תמציתי מנומק:

שימו לב: עדיין מניחים שהרעש ϵ_i מתפלג גאוסיאנית ככתוב בתחילת השאלה.

~~ב. [11 נק'] הוכיחו שתחת כלל ההנחות, בעיית LS עם רגורזיציית ℓ^1 , משמע $\min_{\mathbf{w}} \frac{1}{m} \sum_{i=1}^m (\mathbf{w}^T \mathbf{x}_i - y_i)^2 + \lambda \|\mathbf{w}\|_1$~~

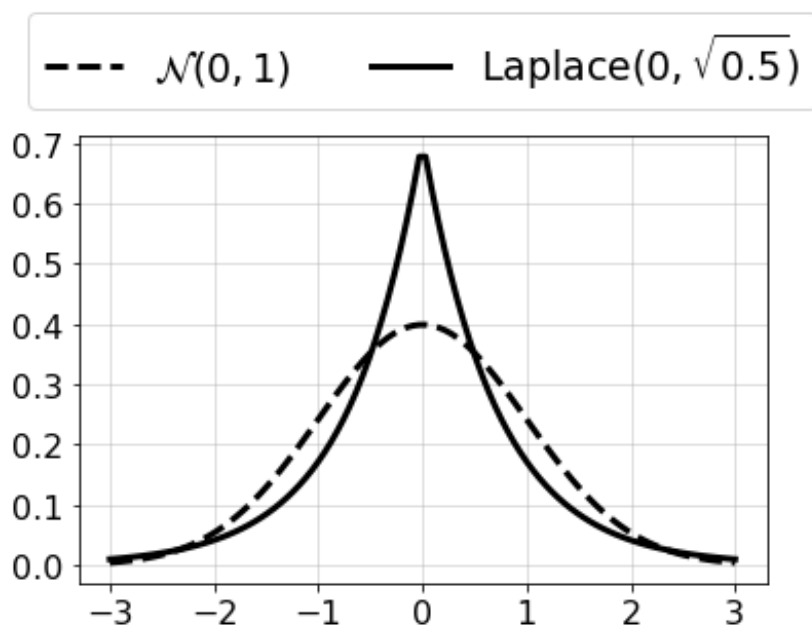
ניתן להשתמש בחישובים שכבר הוכחו בסעיף הקודם.

הוכחה ע"י פיתוח תמציתי מנומק:

[illegible]

השונויות של התפלגות לפלאס נתונה ע"י $\text{Var}[w_k] = 2b^2$.

התרשים משווה בין התפלגות לפלאס להתפלגות נורמלית שהשונויות שלהן היא 1.



תזכורת: הוכחנו בתרגול שאם מניחים $w_k \sim \mathcal{N}(0, 1)$, בעיית ה-MAP שקולה לבעיית LS עם רגולריזציה ℓ^2 .

ג. [5 נק'] מתוך הסתכלות בתרשים, מתוך התזכורת ומתוך מה שהוכחתם בסעיף הקודם, הסבירו בקצרה ובאופן אינטואיטיבי (לא פורמלי) הבדל שלמדנו בין אופי הפיתרונות שמתקבלים ע"י רגולריזציה ℓ^2 לאלה המתקבלים ע"י רגולריזציה ℓ^1 .

הסבר קצר:

שאלה 4 – Kernel SVM [17 נק']

עבור פרמטר נתון $\gamma > 0$, נגדיר את ה-Gaussian kernel לקלט חד-ממדי באופן הבא:

$$K: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}, \quad K(a, b) = \exp(-\gamma(a - b)^2)$$

א. [12 נק'] הציעו פונקציית מיפוי $\phi: \mathbb{R} \rightarrow \mathbb{R}^p$ והוכיחו בעזרתה שהפונקציה K מהווה קרנל חוקי (בחד ממד).

שימו לב: עליכם לבחור $p \in \mathbb{N} \cup \{\infty\}$ מתאים, סופי או אינסופי.

$$e^x = \sum_{n=0}^{\infty} \frac{x^n}{n!}$$

תזכורת: פירוק טור טיילור של e^x לכל x נתון ע"י:

תשובה:

$$e^{-\gamma(a-b)^2} = e^{-\gamma a^2} \cdot e^{2\gamma ab} \cdot e^{-\gamma b^2} = \alpha \cdot \beta \cdot \gamma$$

$$= \alpha \beta \sum_{n=0}^{\infty} \frac{(2\gamma ab)^n}{n!} = \alpha \cdot \sum_{n=0}^{\infty} \frac{(2\gamma)^{\frac{n}{2}} a^n}{\sqrt{n!}} \cdot \frac{(2\gamma)^{\frac{n}{2}} b^n}{\sqrt{n!}} \cdot \beta$$

$\phi: \mathbb{R} \rightarrow \mathbb{R}^{\infty}$

$$\phi(a) = [\phi_0(a), \phi_1(a), \dots] \quad \phi_i(a) = e^{-\gamma a^2} \cdot \frac{(2\gamma)^{\frac{i}{2}} a^i}{\sqrt{i!}}$$

$$\phi(a)^T \phi(b) = e^{-\gamma(a-b)^2} = K(a, b)$$

המשך לסעיף ב':

ב. [5 נק'] נתון dataset עם $m = 1000$ דוגמאות חד-ממדיות. נרצה לפתור את הבעיה עם ה-Gaussian kernel שהגדרנו. מבחינת יעילות, האם עדיף לפתור את ה-primal problem עם ה-feature mapping שמצאתם, או שעדיף לפתור את ה-dual problem עם פונקציית ה-kernel שהוגדרה? ענו והסבירו בקצרה.

תשובה :

א dual מסיון שלו קיצר היחסים לזמן מחייו גרמי הינו
 למחז ∞ דקטורציה היקרה וזה נשק לבסס
 ס. ליוויז צול זק כהיו (n) בזה $(\frac{1000}{2})$

חלק ב' – שאלות אמריקאיות [25 נק']

בשאלות הבאות סמנו את התשובות המתאימות (לפי ההוראות). בחלק זה אין צורך לכתוב הסברים.

א. [5 נק'] סמנו את כל התשובות הנכונות ביחס לאלגוריתמי One vs. One (1v1) ו-One vs. All (1vA).

(הניחו שיש 10 מחלקות ומעלה).

$$\binom{10}{2} \text{ vs } 10$$

a. ל-1vA סיבוכיות מקום נמוכה יותר מזו של 1v1 בזמן האימון.

b. ל-1vA סיבוכיות מקום נמוכה יותר מזו של 1v1 בזמן המבחן (לאחר שהאימון הושלם).

c. רק אחד משני האלגוריתמים ניתן למיקבול (parallelization).

d. 1v1 נוטה יותר ליצור בעיות לא מאוזנות (imbalanced).

ב. [5 נק'] סמנו את כל הטענות שמסלימות בצורה הגיונית את הטענה הבאה.

באופן כללי, ככל שה-complexity של מחלקת היפותזות עולה:

a. ה-bias עולה.

b. ה-variance עולה.

c. צריך פחות דאטה על מנת להכליל כראוי.

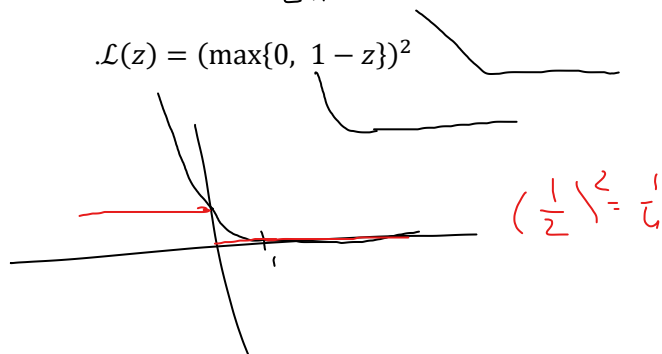
d. יש יותר נטייה ל-overfitting.

e. תהליך האימון של מסווג בודד דורש זמן רב יותר.

$$z \geq 1: z^2 - 2z + 1$$

$$z < 1: 0$$

$$\mathcal{L}(z) = (\max\{0, 1 - z\})^2$$



ג. [5 נק'] נגדיר את פונקציית ה-squared hinge loss:

סמנו את כל הטענות הנכונות ביחס לפונקציה זו.

a. הפונקציה קמורה ביחס ל- z .

b. הנגזרת של הפונקציה היא $\frac{\partial}{\partial z} \mathcal{L} = 2 - 2z$.

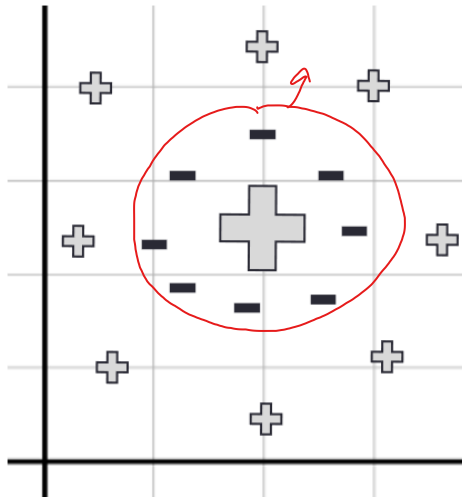
c. הפונקציה חוסמת מלמעלה את ה-0-1 loss בכל מקום.

d. הפונקציה חוסמת מלמעלה את ה-hinge loss בכל מקום.

e. עבור בעיות סיווג תחת מודל ליניארי, משמע $z = y_i \mathbf{w}^T \mathbf{x}_i$, הפונק' מעודדת margin מהמפריד.

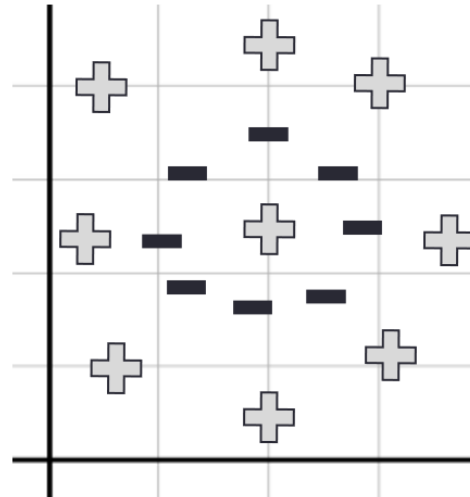
ד. [5 נק'] נתון דאטה עם תיוגים בינאריים ("+" או "-"). מריצים AdaBoost עם מסווג בסיס לא ידוע. גדלי הצורות בתרשימים מסמלים את ההסתברויות שהאלגוריתם מקצה (הסתברות גבוהה = צורה גדולה). מריצים את האלגוריתם איטרציה אחת ומקבלים את התרשימים הבאים:

ההתפלגות לאחר איטרציה אחת



(הנק' האמצעית גדלה והיתר קטנו)

ההתפלגות האחידה ההתחלתית



איזה סוג של מסווג בסיס יכול להסביר את התרשים השמאלי שהתקבל? סמנו את התשובה הנכונה.

☒ a. עץ החלטה בעומק 1 (decision stump). משמע, שורש ושני עלים.

☒ b. עץ החלטה בעומק 2. משמע, שורש, רמת ביניים ועד ארבעה עלים.

☒ c. מסווג שאומר על כל המרחב "שקר" או "אמת".

☒ d. SVM עם קרנל פולינומיאלי ממעלה 2.

☒ e. כל התשובות הקודמות לא נכונות.

$\beta \rightarrow 0 \rightarrow \text{uniform}$

$\beta \rightarrow \infty \rightarrow \text{isolat}$

ה. [5 נק'] היזכרו בפונקציית ה-Sigmoid שמשמשת כשכבה האחרונה של רשת נוירונים לסיווג ל-K מחלקות:

$$\text{softmax}(f_1(x), \dots, f_K(x); \beta) = \left[\frac{\exp\{\beta f_1(x)\}}{\sum_{i \in [K]} \exp\{\beta f_i(x)\}}, \dots, \frac{\exp\{\beta f_K(x)\}}{\sum_{i \in [K]} \exp\{\beta f_i(x)\}} \right]^T$$

בסעיף זה אנו לא מתייחסים כלל לאפשרות ש- $\beta = 0$ ומניחים של- β אותו סימן בזמן האימון ובזמן המבחן.

סמנו את כל הטענות הנכונות ביחס לפונקציה זו.

☒ a. בזמן מבחן (לאחר האימון), כאשר $\beta \rightarrow \infty$, התפלגות הפלט הולכת להתפלגות אחידה.

☒ b. כאשר משנים את β לאחר האימון בזמן המבחן, כל עוד β שומר על הסימן, אין לו השפעה על הדיוק של הרשת.

☒ c. בזמן אימון, כל עוד הפרמטר β חיובי, אין לו השפעה על מהלך האימון.

d. בזמן אימון, אם הפרמטר β שלילי, לא ניתן ללמוד את הרשת בעזרת שיטות gradient.

$\beta \rightarrow 0 \rightarrow c$

$e^0 = 1$

מסגרת נוספת (יש לציין אם מדובר בטייטה או בהמשך לתשובה אחרת):

A large rectangular box with rounded corners, containing 25 horizontal lines for writing. The box is empty, with the lines spaced evenly apart, providing a designated area for the student to provide a second answer or further explanation.

מסגרת נוספת (יש לציין אם מדובר בטייטה או בהמשך לתשובה אחרת):

A large rectangular box with rounded corners, containing 25 horizontal lines for writing. The box is empty, with the lines spaced evenly from top to bottom.

מסגרת נוספת (יש לציין אם מדובר בטייטה או בהמשך לתשובה אחרת):

A large rectangular box with rounded corners, containing 25 horizontal lines for writing. The box is empty, with no text or markings inside.