

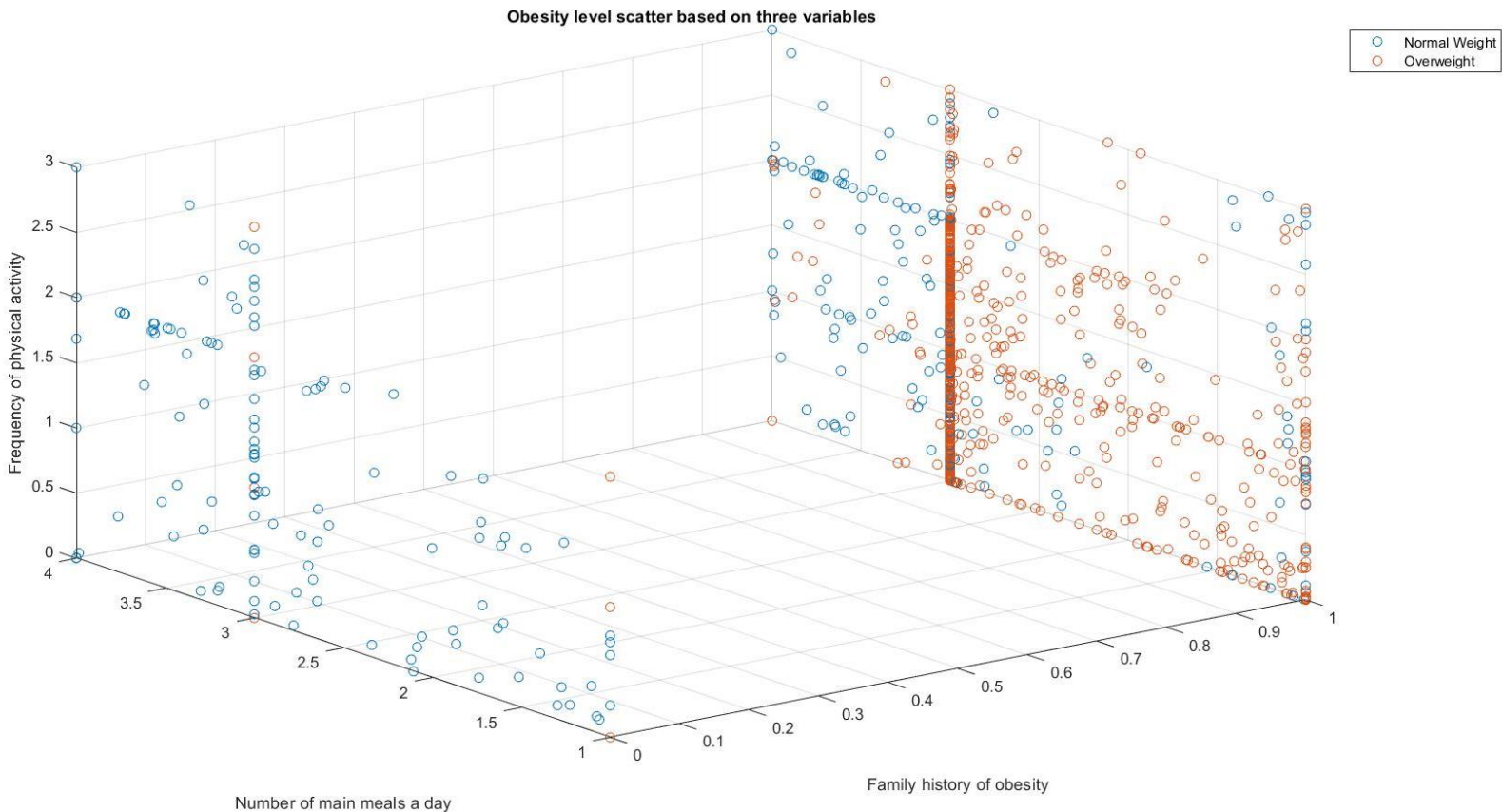
תרגיל מסכם – קורס גישות חישוביות מס' 27-431-01 (הרצת הקוד 'mainScript.m' תריץ את כל התרגיל – צריך מטלב 2019/2020) *הדאטה אצלנו עובר רנדומליזציה, לכן הגרפים והתוצאות יצאו מעט שונים בכל הרצה*****

1) המאגר מכיל משתנים שונים הקשורים לאורח החיים של הנבדק. יש מגוון משתנים מסוגים שונים לדוגמה גיל כמשתנה רציף, מין כמשתנה בינארי, שימוש באמצעי תחבורה כמשתנה בדיד, עישון כמשתנה בדיד וכו'. סך הכל יש כשישה עשר משתנים מסוגים שונים, המשתנים הללו מטרם היא לחזות משתנה אחד שהוא רמת Obesity. אנחנו בחרנו להכניס ללמידת המודל את כל המשתנים למעט – גיל, מין, גובה ומשקל. רמת obesity המשתנה אותו אנו מנבאים, הינו משתנה בינארי כאשר 0 משמעותו משקל תקין ו1 משקל עודף. לדעתנו המשתנים הרלוונטיים ביותר לניסוי הינם – היסטוריה משפחתית של השמנת יתר, מספר ארוחות עיקריות ביום ותדירות פעילות גופנית.

2) נעשה בקוד – חלוקה ל 20% סט מבחן ו80% סט אימון. (פונקציית 'prepareData.m')

3) בחרנו במשתנים של היסטוריה משפחתית של השמנת יתר, מספר ארוחות עיקריות ביום ותדירות פעילות גופנית. בגרף ניתן לראות כי היסטוריה משפחתית מצליחה להפריד בצורה טובה יחסית בין משקל תקין למשקל עודף. בנוסף רואים כי רוב האנשים בעלי משקל עודף אוכל בין 1 ל3 ארוחות עיקריות ביום, בעוד משקל תקין מפוזרים יחסית על הגרף אבל ניתן לראות ריכוז טיפה יותר גדול בין 3 ל4 ארוחות. ניתן להסיק מהגרף כי שלושה משתנים אינם מספיקים לסיווג של דאטה זה, לכן יש צורך להשתמש בעוד משתנים על מנת לסווג את הדאטה טוב יותר.

גרף scatter – (פונקציית 'plotScatter.m')



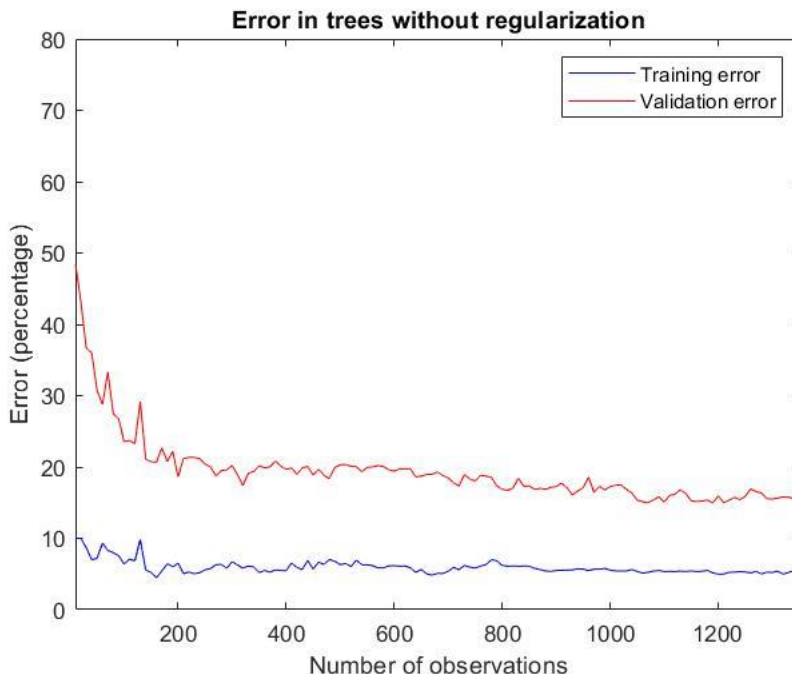
שיטת קלסיפיקציה : עצי החלטה – פונקציה 'decisionTreeModel.m'

4) לדעתנו פונקציית השגיאה המתאימה בשיטה זו הינה מיסקלסיפיקציה, בה על כל נקודה שסווגה לא נכון יש עלייה של נקודה אחת בשגיאה. פונקציית שגיאה זו עלולה לגרום למודל לעשות overfitting, היות ותמיד תהיה שגיאה אם יהיו נקודות שמסווגות לא נכון. לכן, בחרנו היפרפרמטרים מתאימים אשר יגבילו את overfitting של המודל. בחרנו Accuracy כמדד להצלחת המודל, הוא מציג את אחוז הדוגמאות אשר מסווגות נכון במודל. רשימת ההיפרפרמטרים שבחרנו וההסבר:

- **אלגוריתם פרדיקציה של הפיצול הטוב ביותר** – 'PullLeft'
- מדד להערכת פיצול של ענף – Gini הוא המדד שבחרנו. הוא בודק מהו הפיצול שיניב לנו את החלוקה הטובה ביותר של הדאטה.
- **מספר מקסימלי של פיצולים** – השארנו את ברירת המחדל שהיא $[\# \text{ of observations} - 1]$.
- **מספר מינימלי של פרדיקציות leaf** – השארנו את ברירת המחדל שהיא מינימום תצפית אחת לעלה.
- **מספר מינימלי של פרדיקציות branch nodes** – השארנו את ברירת המחדל שהיא מינימום 10 תצפיות לענף.

הסיבה שהשארנו את אפשרויות ברירת המחדל של מטלב היא שרצינו שהמודל שלנו לא יעבור רגולריזציה ויוכל להתאמן על האימון ללא הפרעה. ברירות המחדל של מטלב הינם המספרים המינימלים על מנת שלא תעשה רגולריזציה והגבלה של למידת המודל (כאשר נדרשנו להוסיף רגולריזציה הגדרנו מספרים אלו על מנת להקשות על למידת המודל).

חילקנו את סט האימון לחמש סאבסטים (5 fold) ועשינו קרוס וולידציה. על מנת לחשב את שגיאת האימון השתמשנו בפונקציה מובנית של מטלב שנקראת 'resubloss'. חישוב של שגיאת הוולידציה גם הוא נעשה בעזרת פונקציה מובנית של מטלב הנקראת 'loss'. על מנת למצוא הבדלים בין הלמידה של חמש הסאבסטים חישבנו שונות של השגיאות. קיבלנו שונות שהיא 0.04 לשגיאת האימון ו0.08 לשגיאת הוולידציה. זה אומר כי הלמידה של הסאבסטים דומה מאוד. לכן על מנת להציג את גרף שגיאת האימון ושגיאת הוולידציה, עשינו ממוצע לשגיאות האימון של כל חמשת הסאבסטים וכך גם לשגיאת הוולידציה.

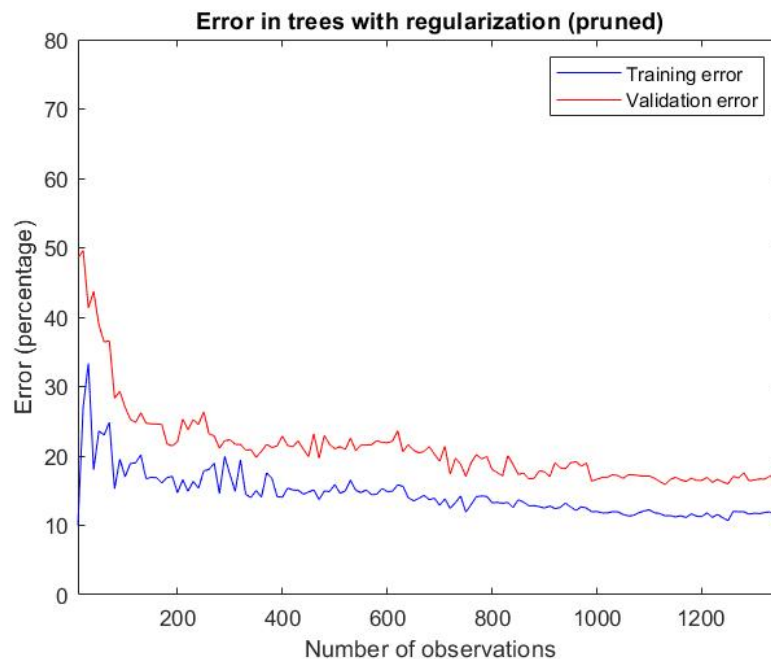


ניתן לראות בגרף השגיאה משמאל –

שגיאת האימון נשארת נמוכה בין 5% ל10% אחוז במהלך הלמידה, היא אינה מושפעת ממספר התצפיות שמהן המודל למד. בלמידה על כל סט האימון שגיאת האימון הינה 5%. לעומת זאת שגיאת הוולידציה גבוהה יחסית ויורדת מהר בין 0 ל200 תצפיות. מ200 תצפיות והלאה ישנה ירידה מתונה מאוד עד לא מורגשת של שגיאת הוולידציה, כך שבלמידה על כל סט האימון שגיאת הוולידציה הינה 16%.

(5) עשינו רגולריזציה בעזרת קיצוץ הענפים של העץ שנוצר. קיצוץ נעשה לענפים בעלי impurity הכי גבוהה. מה שאומר שמדד gini שלהם הכי קרוב ל1. מדד gini קרוב ל1 מצביע על ענף שיוצר חלוקה לא טובה של הדאטה ולכן אין צורך בחלוקה זו, כי היא אינה מוסיפה הרבה מידע לעץ. בקוד הוספנו קריטריון לpruning שהוא impurity. בחרנו רמת קיצוץ שהיא 75% הרמת הקיצוץ האפשרית הגבוהה ביותר. לדוגמה עץ בעל 20 שכבות ורמת הקיצוץ הגבוהה ביותר היא 15, אז רמת הקיצוץ שאנו בחרנו היא בערך 11.

ניתן לראות בגרף כי רגולריזציה גרמה לשגיאת האימון להיות גבוהה יותר, מה שאומר שאכן הצלחנו להגביל את למידת המודל לסט האימון באופן ספציפי מדי. אולם שגיאת הוולידציה עלתה מעט גם כן. זה מצביע על כך שהרגולריזציה שהחלנו על העץ הינה חזקה מדי, מה שגורם למודל לunderfitting זאת אומרת להיות פחות טוב לוולידציה גם כן. לאחר למידה של המודל על כל התצפיות הגענו לשגיאת אימון של 11% (לעומת 5% ללא רגולריזציה) ולשגיאת וולידציה של 17% (לעומת 16% ללא רגולריזציה). גם הפעם הירידה המשמעותית ביותר בשגיאת האימון והוולידציה נעשית בין 0 ל200 תצפיות. גם כאן חישבנו שונות על מנת לבדוק שהשונות אכן קטנה. השונות לשגיאת האימון שהתקבלה היא 0.06, ואילו השונות לשגיאת הוולידציה היא 0.1.



(6) השתמשנו בפונקציה של מטלב על מנת למצוא את ההיפרפרמטרים הטובים ביותר בשביל כל מודל מבין חמשת המודלים של הקרוס וולידציה. ההיפרפרמטרים שבחנו הינם – מספר מקסימלי של פיצולים, מספר תצפיות מינימלי בעלה, קריטריון לפיצול. השתמשנו בrandom search. מודל העץ בעל ההיפרפרמטרים האופטימליים ביותר ושגיאת הוולידציה הקטנה ביותר נבחר בשביל ההמשך בשאלות 9 ו10.

טבלה של ההיפרפרמטרים האופטימליים שהתקבלו בשביל כל fold:

מספר מודל / fold	מינימום תצפיות בעלה	מקסימום פיצולים	קריטריון לפיצול
1	7	875	gini
2	1	171	deviance
3	5	278	deviance
4	1	968	deviance
5	1	122	gini

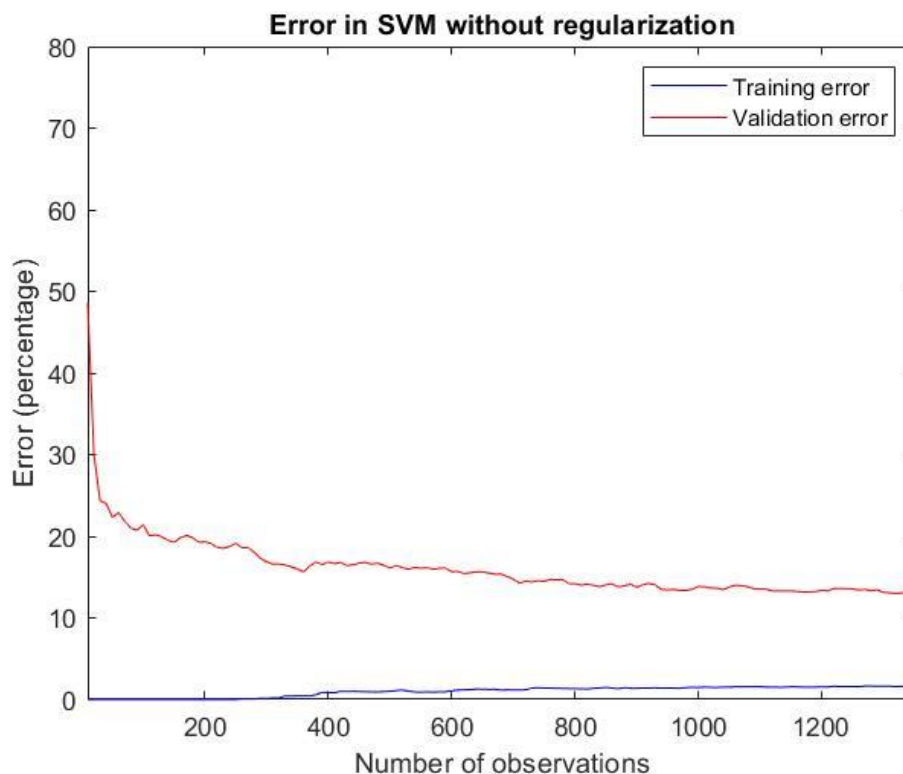
שיטת קלסיפיקציה : SVM – פונקציה 'svmModel.m'

7) 4. בסעיף זה בחרנו פונקציית שגיאה של מיסקלסיפיקציה, כאשר נקודות אשר לא מסווגות נכון מוסיפות 1 לערך השגיאה. מדד הצלחת המודל שבחרנו הינו accuracy, הוא מייצג את אחוז התצפיות אשר מסווגות נכון. ההיפרפרמטרים שבחרנו הינם כתוצאה מברירת מחדל של מטלב, אשר מחייבת להגדיר אותם. רשימת ההיפרפרמטרים שבחרנו להגדיר וההסבר:

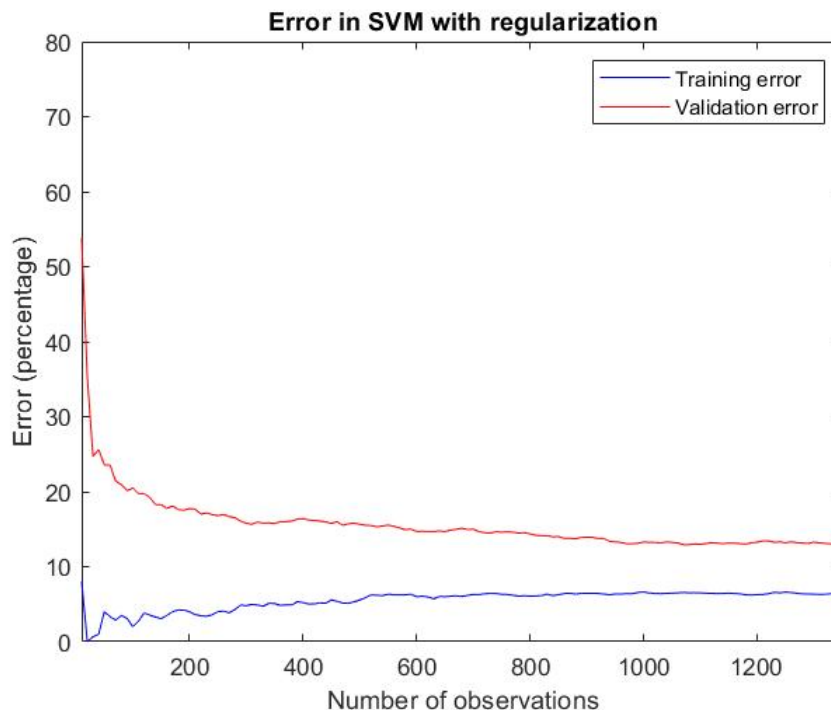
- **'BoxConstraint'** – הינו היפרפרמטר אשר מגדיר את העונש על כל נקודה שמסווגת לא נכון. היות וברירת המחדל עבור היפרפרמטר זה היא 1, ברירת מחדל זו מקשה על המודל ללמוד. לכן בחרנו להגדיל אותו ל-50 מה שמאפשר למודל ללמוד טוב ובקלות את סט האימון, מכיוון שהעונש על כל טעות הוא גדול. מה שנוצר לנו הוא מודל מאוד ספציפי לסט האימון.
- **'KernelFunction'** – השארנו לפי ברירת מחדל את הקרנל הגאוסיאני. לאחר הרצה של כמה קרנלים ובחינה של גרפי שגיאות האימון והוולידציה, ראינו כי קרנל זה הינו הטוב ביותר עבורנו למידת המודל שלנו.

על מנת לחשב את שגיאת האימון השתמשנו בפונקציה מובנית של מטלב שנקראת 'resubloss'. חישוב של שגיאת הוולידציה גם הוא נעשה בעזרת פונקציה מובנית של מטלב הנקראת 'loss'.

בגרף – ניתן לראות שבתחילת אימון המודל שגיאת הוולידציה גבוהה מאוד וככל שנוספות יותר תצפיות כך שגיאת הוולידציה יורדת (ירידה חדה עד ל-200 תצפיות ולאחר מכן ירידה מינורית) ושגיאת האימון עולה, אך באופן מינורי יחסית. לאחר למידה על כל התצפיות שגיאת האימון הגיעה ל-1% ושגיאת הוולידציה ל-12.5%. מאחר ובדקנו את השונות בין 5 המודלים של קרס וולידציה ולא נמצאה שונות משמעותית ביניהם, בחרנו לייצג את הגרף הממוצע של כולם. השונות של שגיאת האימון הייתה 0.009 ושל שגיאת הוולידציה הייתה 0.128.



5. רגולריזציה נעשתה בעזרת שינוי היפרפרמטר BoxConstraint ל1. מה שאומר שההענשה על נקודות שלא סווגו נכון הינה קטנה ומתונה יותר. דבר זה מאפשר למודל שלנו ללמוד את סט האימון בפחות ספציפיות, מכיוון שהוא מאפשר לו לעשות טעויות. כמו שניתן לראות מהגרף, הרגולריזציה גרמה אכן לשגיאת האימון לעלות מכיוון שכעת המודל פחות ספציפי לסט האימון, בעוד ששגיאת הולידציה נשארה יחסית זהה למודל בלי רגולריזציה. בין 0 ל200 תצפיות ניתן לראות את הירידה הגדולה ביותר בשגיאת הוולידציה (אבל לא בשגיאת האימון). לאחר מכן שגיאת הוולידציה והאימון נשארות יחסית קבועות, שגיאת אימון בין 2.5% ל5% ושגיאת הוולידציה בין 15% ל12.5%. העובדה ששגיאת האימון עלתה ושגיאת הוולידציה נשארה ללא שינוי, יכולה להצביע על כך שהרגולריזציה אשר ביצענו אינה מספיק חזקה והמודל עדיין יחסית ספציפי לסט האימון. גם כאן חישבנו שונות על מנת לוודא שהיא אכן זניחה וניתן להציג את חמשת המודלים כממוצע בגרף אחד. השונות של שגיאת האימון שהתקבלה היא 0.0019, השונות של שגיאת הוולידציה הינה 0.065.



6. השתמשנו בפונקציה של מטלב על מנת למצוא את ההיפרפרמטרים הטובים ביותר בשביל כל מודל מבין חמשת המודלים של הקרוס וולידציה. ההיפרפרמטרים שבחנו הם - 'BoxConstraint', 'KernelScale', 'KernelFunction'. השתמשנו בrandom search. מודל הSVM בעל ההיפרפרמטרים האופטימליים ביותר ושגיאת הוולידציה הקטנה ביותר נבחר בשביל ההמשך בשאלות 9 ו10.

טבלה של ההיפרפרמטרים האופטימליים שהתקבלו בשביל כל fold:

Kernel function	Kernel Scale	Box constraint	fold / מודל
גאוסיאני	1.55	50.45	1
פולינומיאלי	none	16.78	2
פולינומיאלי	none	0.0039	3
פולינומיאלי	none	0.0086	4
גאוסיאני	0.522	51.31	5

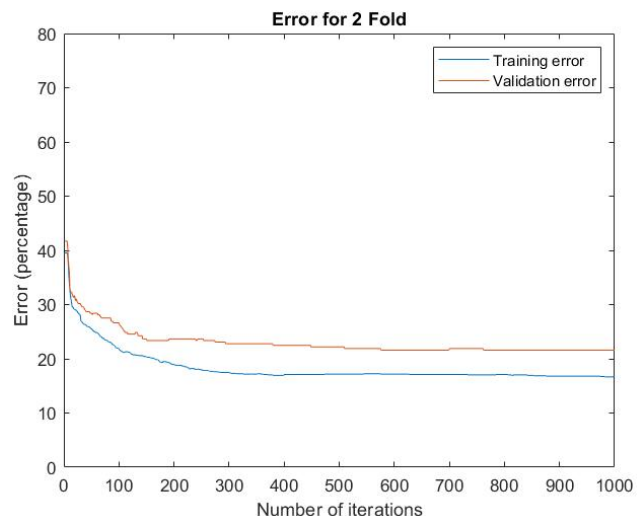
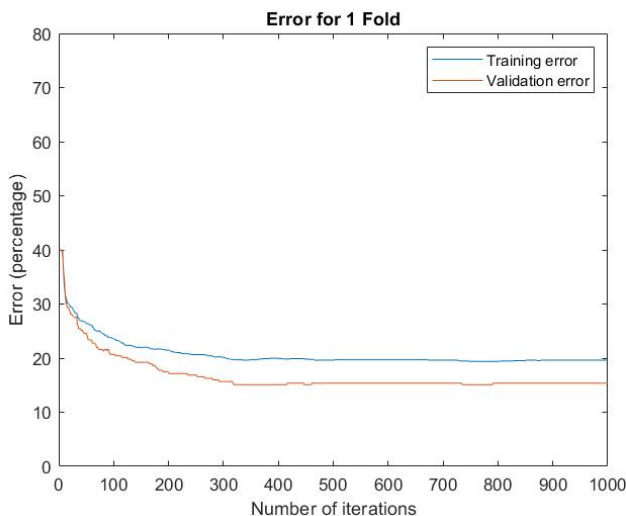
שיטת קלסיפיקציה : רגרסיה לוגיסטית – פונקציות, 'logisticRegressionModel.m', 'convertData.m', 'cost.m', 'costRegu.m', 'gridSearch.m', 'predTest.m', 'sigmoid.m'

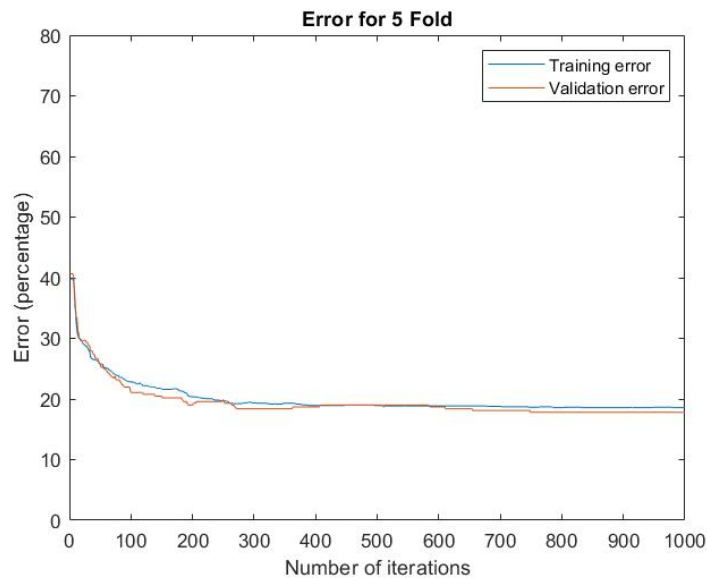
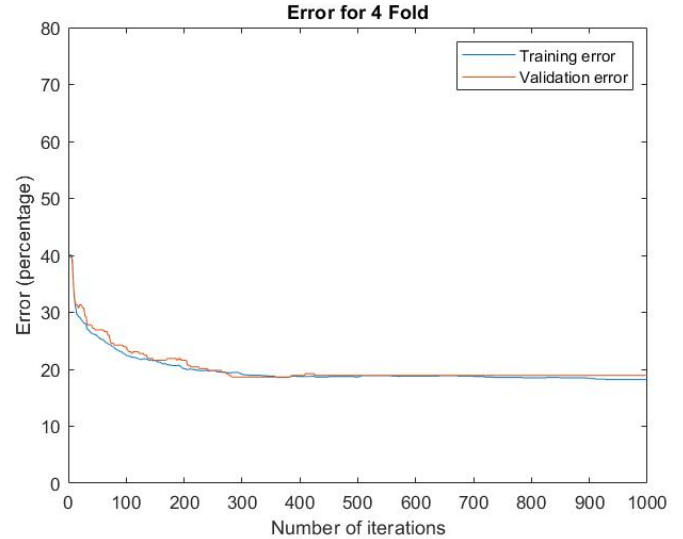
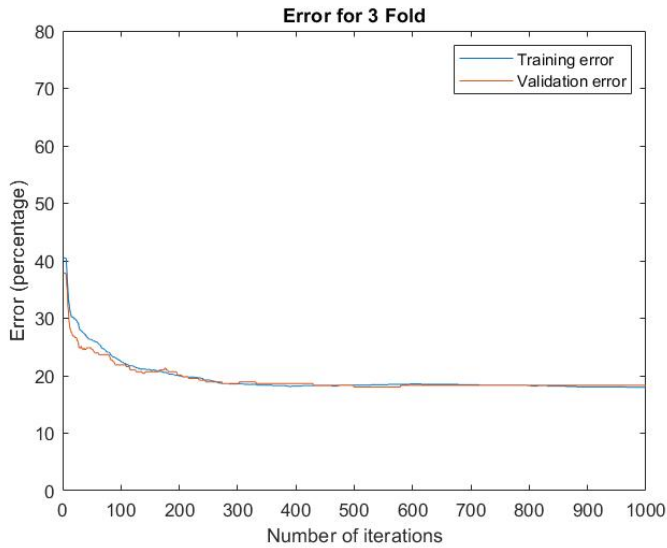
8) 4. רגרסיה לוגיסטית הינה פונקציה מסוג supervised learning כלומר אנחנו מספקים למודל את X והמטרה שלנו הוא לבצע פרדיקציה על Y . זוהי שיטת רגרסיה שנוכל להשתמש בה כאשר נרצה להתגבר על בעיית הערכים הקיצוניים הקיימת לפעמים בשימוש בפונקציה ליניארית. היתרון שלה זה שהערך נע בין 0 ל1, ולכן נוכל להגיד שערכים שהם מעל חצי הם 1 ומה שקטן מחצי הוא 0. גם שאנחנו מסווגים נכון את המודל עדיין יש לפונקציה ערך ונוכל להמשיך לתקן את המשקולות למרות שהאלגוריתם צודק, ולקבל כך מודל יותר טוב. בחרנו את פונקציית השגיאה כ $Cost(h_\theta(x), y) = -y \log(h_\theta(x)) - (1 - y) \log(1 - h_\theta(x))$ ואת מדד ההצלחת המודל Accuracy, שהוא אחוז התצפיות אשר מסווגות נכון.

בשביל תהליך הלמידה הראשוני של המודל בחרנו את ההיפרפרמטרים הבאים:

- **איטרציות** – הגדרנו כ1000. כלומר המודל יבצע אלף חזרות כדי להגיע לרמת למידה מספקת.
- **אלפא (קבוע למידה)** – הגדרנו כ0.1. בחרנו אותו כך שהמודל מצד אחד יוכל לבצע למידה טובה של המשקולות ב"צעדים" יחסית קטנים, אבל מצד שני לא קטן מידי אשר גורם לזמן למידה ארוך מידי.
- **למדה (קבוע רגולריזציה)** – בשביל סעיף זה לא הכנסנו את למדה לפונקציית הcost של למידת המודל. אולם כאשר הוספנו רגולריזציה בסעיף 5, קבענו כי קבוע הרגולריזציה יהיה 0.5. בחרנו מקדם הענשה זה על מנת שלא ייצור הענשה גדולה מידי על המשקולות שתפריע ללמידה, אבל מספיק גדול על מנת לייצר הענשה משמעותית אשר תשפר את למידת המודל.

על פי הגרפים של שגיאות הוולידציה והאימון, ניתן לראות כי אכן יש הבדל בין הצלחת חמשת המודלים (5 fold). לכן ברגרסיה הלוגיסטית הצגנו את המודלים בגרפים נפרדים בשונה מעצי החלטה וSVM. בכל המודלים ניתן לראות כי הירידה הגדולה ביותר בשגיאת האימון והוולידציה נעשית בין 0 ל200 איטרציות, לאחר מכן השגיאות נשארות יחסית קבועות. בכל המודלים שגיאת האימון מגיעה ב1000 איטרציות ל20%. שגיאת הוולידציה לעומת זאת מגיעה ל15% במודל הראשון, ל22% במודל השני, ל20% במודל השלישי, הרביעי והחמישי. במודל הראשון שגיאת הוולידציה קטנה יותר משגיאת האימון, זה יכול לבבוע מכמה סיבות. ייתכן כי בסט האימון הספציפי לfold זה היו הרבה דוגמאות קשות לסיווג, או שבסט הוולידציה הזה היו דוגמאות קלות לסיווג. במודל השני שגיאת האימון קטנה יותר משגיאת הוולידציה, מה שאומר שהמודל התאמן יותר מידי טוב על סט האימון. במודל השלישי, הרביעי והחמישי שגיאות האימון והוולידציה חופפות מה שיכול להצביע על כך שהמודל למד בצורה המיטבית על סט האימון והצליח להכליל את הלמידה על סט הוולידציה.

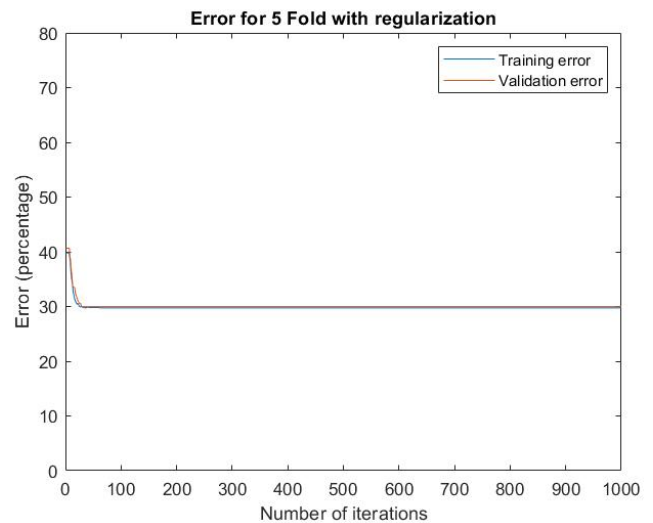
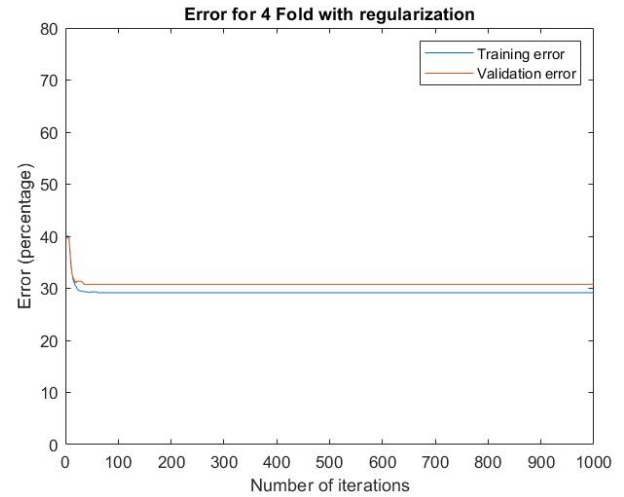
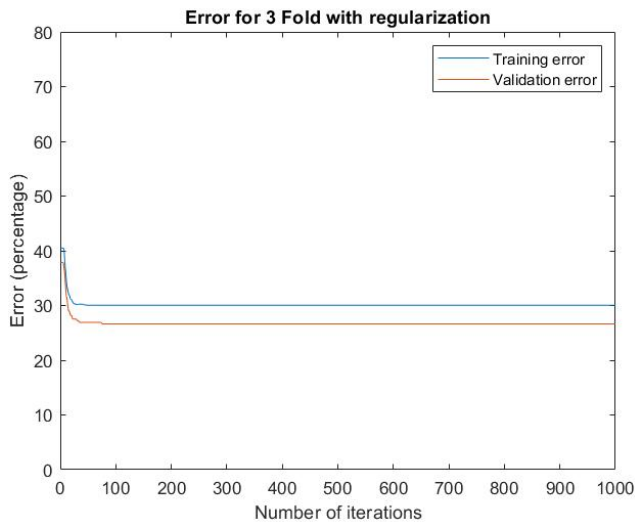
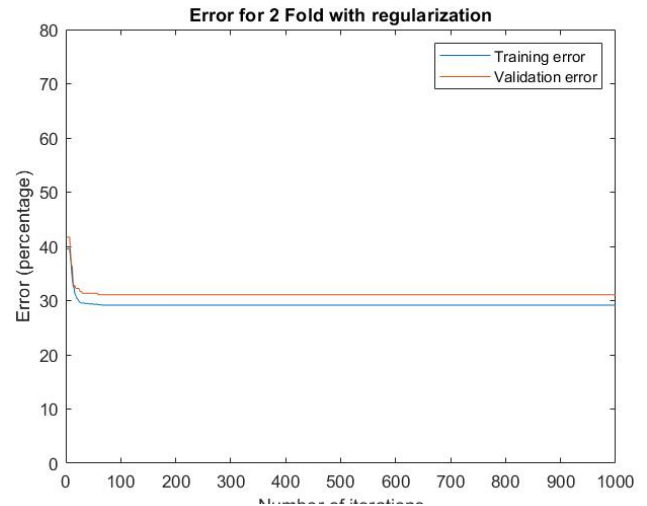
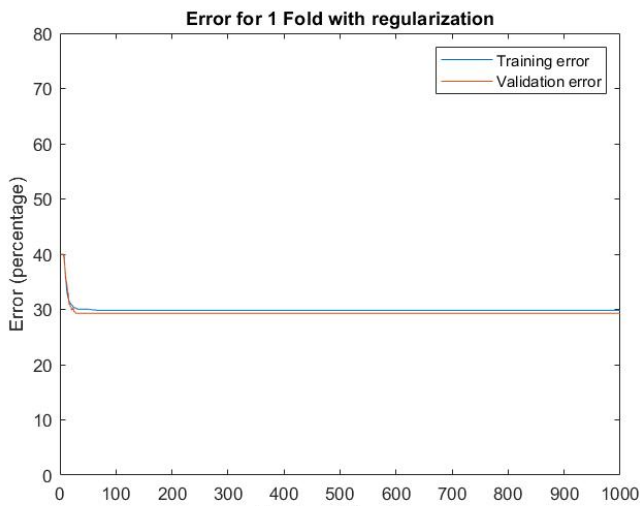




5. בעת צירפנו את קבוע הרגולריזציה (למדה) לפונקציית ה-cost. מטרת הרגולריזציה היא להעניש משקולות גדולות ולגרום לכך שהרשת תהיה פשוטה יותר. נקבע קבוע רגולריזציה המוכפל בפונקציית השיאה- אם נבחר קבוע קטן מידי אז לא תהיה מספיק הענשה של המשקולות ונשאר במצב של overfitting, אך אם נבחר קבוע גדול מידי הענשת המשקולות תהיה כבדה ונגרום למצב של underfitting.

בכל המודלים ניתן לראות כי הוספת רגולריזציה הפריעה ללמידה של המודל. ניתן לראות ששיאות האימון ושיאות הוולידציה עלו יחד ל-30%, בהשוואה ל-20% שגיאה ללא רגולריזציה. אנו משערים כי ללמידה של המודלים, מה שיכול להיות גרם למודל לעשות Underfitting גם לאימון ולגם לוולידציה.

גרפים של פונקציית השיאה עם רגולריזציה -



6. ביצענו grid search עבור כל אחד מחמשת המודלים.

חיפשנו את ההיפרפרמטרים האופטימליים הבאים – מספר

איטרציות, קבוע למידה (אלפא), קבוע רגורלריזציה (למדה).

מסעיף זה לקחנו את המודל הטוב ביותר בעל

ההיפרפרמטרים האופטימליים וטעות הוולידציה הקטנה ביותר והשתמשנו בו לשאלות 9 ו-10.

טבלה של ההיפרפרמטרים האופטימליים שהתקבלו בשביל כל fold:

מספר מודל / fold	מספר איטרציות	קבוע למידה (אלפא)	קבוע רגולריזציה (למדה)
1	700	0.1	1.4
2	700	0.1	1.4
3	700	0.1	1.4
4	700	0.1	1.4
5	700	0.1	1.4

(9) עשינו זאת בפונקציה **'ensembleModelandAccuracy.m'**. את הגרף של הדיוק נציג ונסביר בשאלה הבאה (10), זאת על מנת לצרף להשוואת הדיוק את המודל של ensemble אותו בנינו בשאלה 10.

(10) עשינו זאת בפונקציה **'ensembleModelandAccuracy.m'**. חשוב לציין כי במודל ensemble לא הייתה אפשרות של בחירה רנדומלית בין הפרדיקציה של המודלים השונים. זאת מכיוון שישנם שלושה מודלים ואילו הסיווג הינו בינארי, לכן תמיד יהיו שני מודלים שיעשו את אותה הפרדיקציה. לדוגמה מודל עץ – 0, מודל SVM – 1 ומודל רגרסיה לוגיסטית – $1 \leftarrow$ כך ensemble יעשה פרדיקציה של 1.

ניתן לראות מגרף הדיוק כי המודל בעל אחוז הדיוק הגבוהה ביותר הינו מודל SVM (80%), לאחר מכן למודל עץ ההחלטה יש את הדיוק הבינוני (75%) ולבסוף מודל הרגרסיה הלוגיסטית בעל אחוז הדיוק הנמוך ביותר (70%). ניתן לראות כי מודל ensemble הינו בעל אחוז דיוק בין SVM לעץ ההחלטה (77%), תוצאה זו הגיונית מכיוון שהשתמשנו בכלל של "הרוב קובע" על מנת לייצור את מודל זה.

