

ניבוי נטייה לנטילת סיכונים

נועה בן דרור (ת.ז. 316163260)

מיכל דגן (ת.ז. 315657064)

אפריל 2021

רקע

נטייה לסיכונים הוא נושא שנחקר רבות בספרות בהקשרים שונים. הנושא הוא רחב והנטייה יכולה להיות בתחום הבריאותי, פיננסי, חברתי וכד'. ניתן למצוא ברשת שאלונים שמתבססים על דיווח עצמי של הנבדק הן בהקשר למצבים היפותטיים, והן בהקשר לפעילות ממשית יום יומית. מחקרים מראים כי קיימת קורלציה בין נתוני גיל ומין לבין הנטייה לקחת סיכונים. כך למשל, גברים נוטים לקחת סיכונים יותר מנשים, וצעירים/מתבגרים נוטים לקחת סיכונים יותר ממבוגרים (Bonem et al., 2015). בנוסף, מחקרים מראים כי למאפיין אישיותי של חיפוש אחר ריגושים יש קשר לנטילת סיכונים (Zuckerman & Kuhlman, 2000). ערכים כמו קונפורמיות משקפים אף הם נטייה לנטילת סיכונים, ובנוסף נמצא כי גם אם הנכונות ליטול סיכונים עשויה להשתנות על פני תחומים שונים בחיים, היא תמיד תהיה מושפעת מההעדפה או הרצון לקחת סיכון, וכי ההעדפה זו המבוססת על דווח עצמי נותנת תמונה מהימנה יותר מאשר התנהגות של הנבדק בניסוי (Universität Basel, 2017).

שאלת המחקר

בחרנו להשתמש באלגוריתמים שלמדנו בכיתה כדי לנבא האם בן אדם נוטה לקחת סיכונים.

איסוף הנתונים

אספנו את הנתונים בעצמנו, לשם כך בחרנו מספר שאלות המייצגות את הממצאים והשאלונים שעלו בספרות והוזכרו לעיל בהקשר של נטילת סיכונים. כך למשל שאלות כמו רכיבה על אופניים ללא קסדה או חצית כביש באור אדום משקפות ערכים של קונפורמיות, ואילו שאלות כמו: צניחה חופשית והשקעה בבורסה משקפות את הנטייה לחפש ריגושים. מאחר וכל שאלה מייצגת התנהגות בפועל או מאפיין אובייקטיבי החלטנו שהתשובות תהיינה בינאריות ואין צורך בסולם מדורג.

להלן השאלות:

- מגדר
- האם את/ה מתחת לגיל 50?
- האם את/ה רוכב/ת על אופניים/קורקינט ללא קסדה?
- האם אי פעם חצית במעבר חציה באור אדום?
- האם את/ה מעשן/ת?
- האם השתתפת בתחרות כלשהי? (ספורט, שחמט, אפייה, חידון וכו')
- האם אי פעם טיילת לבד?
- האם עשית צניחה חופשית או בנגיי?
- האם את/ה משקיע/ה בבורסה?
- האם אי פעם מילאת לוטו?
- האם תשתמש/י בגבינה שפג תוקפה לפני יומיים?
- האם אי פעם רימית במבחן?
- האם את/ה מגדיר/ה את עצמך כאדם שלוקח סיכונים?

12 השאלות הראשונות מהוות את הפיצ'רים בעזרתם נרצה לענות על שאלת המחקר, בעוד שהשאלה האחרונה - "האם את/ה מגדיר/ה את עצמך כאדם שנוטה לקחת סיכונים?" - משקפת את ההעדפה או הרצון כפי שהנשאל מדווח על עצמו והיא למעשה ה-label שיאפשר את הסווג הסופי. הציפייה שלנו היא שמענה חיובי על רוב השאלות יעיד על כך שהנבדק נוטה לקחת סיכונים. 120 נבדקים ענו על השאלון הנ"ל.

תכנון הניסוי

בחרנו להשתמש במודל הפרספטרוני הבינארי המבצע פעולה של סיווג. הפרספטרוני מדמה רשת נוירונים מלאכותית, המבוססת על פעילות של נוירון בודד המקבל קלט, ובהתאם פעיל כתלות בסכום המשוקלל של ערכי הקלט. זהו אלגוריתם המשתיך למשפחת האלגוריתמים ל-"למידה מפוקחת". מטרת האלגוריתם היא להבדיל בין סוגים שונים של דגימות אותן הוא מקבל. כלומר, ה'מורה' מספק את דוגמאות האימון וגם את הסיווג המתאים להן, והפרספטרוני מחלק את מרחב הקלט לשני אזורים באמצעות על-מישור. אנחנו מציגים לרשת את הנקודות ואת הסיווג הנכון ו"מאמנים" אותה להפיק את הסיווג המתבקש. אלגוריתם האימון במקרה של הפרספטרוני הוא אלגוריתם איטרטיבי, שבו חוזק המשקלות (הסינפסות) משתנה באופן הדרגתי עד לקבלת הפתרון הרצוי. אם הדאטה ניתן להפרדה על ידי על-מישור, הפרספטרוני יתכנס לפתרון ויסווג בצורה נכונה את כל הדוגמאות, החל

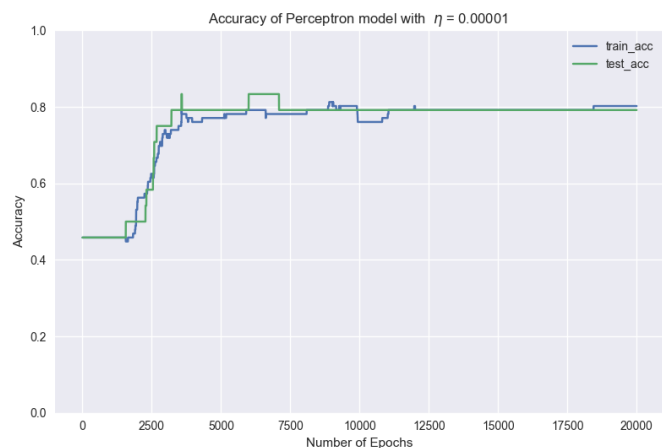
משלב כלשהו. במקרה בו הדאטה אכן ניתן להפרדה, קיים אוסף של על-מישורים שפותרים את בעיית הסיווג. אם הדאטה אינו ניתן להפרדה, לא קיים פרספטרון שסיווג את הקלט בצורה נכונה, ולכן לא נקבל תוצאה מספקת.

בחרנו לבחון את התנהגות האלגוריתם על הנתונים שלנו. מאחר ומדובר במרחב 12-מימדי לא נוכל להציג את ההפרדה באופן ויזואלי, אך נוכל לבחון את אחוז הדיוק של הפרספטרון, ככל שמספר האיטרציות גדל וחוזק המשקולות משתנה.

תוצאות

תוצאות הרצת אלגוריתם הפרספטרון על הדאטה מעידות על כך שהדאטה אינו ניתן להפרדה לינארית באופן ברור בהינתן ה-labels והפיצ'רים שבחרנו. כפי שלמדנו, לפי משפט התכנסות הפרספטרון, אם קיים פתרון אשר מתייג נכון את כל הדוגמאות (כלומר, הן linearly separable) נמצא אותו אחרי מספר סופי של עדכונים. במקרה שלנו, האלגוריתם לא התכנס ולא עצר, ולכן תוצאות הריצה של הפרספטרון בתור מסווג לבעיה שלנו לא מספקות. חילקנו את הדאטה לסט אימון וסט בדיקה (80%-20%) בגרף הבא ניתן לראות את אחוז הדיוק של אלגוריתם הפרספטרון על הדאטה, עם קצב לימוד 0.00001 לאורך מספר גדל של epochs :

ניסנו מספר קצבי לימוד, ולבסוף בחרנו בקצב לימוד יחסית קטן, שמאפשר למידה איטית, כך שלכל עדכון למעשה יש אפקט קטן.



בשתי העקומות אחוז הדיוק נמצא במגמת עלייה עד סף מסוים, והחל ממנו נותר יחסית קבוע, מה שעולה בקנה אחד עם העובדה שהאלגוריתם לא הגיע להתכנסות מבחינת עדכון המשקולות. בנוסף, התלכדות העקומות מעידה על כך שסט הבדיקה מיוצג בצורה טובה ע"י סט האימון, וכי הפרספטרון מצליח פחות או יותר להכליל גם על דוגמאות שהוא לא ראה.

לאחר שהרצנו את אלגוריתם הפרספטרון, וראינו שהמדגם אינו פריד לינארית, בחרנו להשתמש גם באלגוריתם ה-SVM (Support Vector Machine), מכונת וקטורים תומכים.

מטרת הפרספטרון היא למצוא מפריד לינארי כלשהו, אך אם הדוגמאות ניתנות להפרדה יהיו הרבה מפרידים אפשריים. לאלגוריתם אין דרך לקבוע מי מהם "עדיף", וההתכנסות תהיה שרירותית לאחד מהם (כתלות בתנאי ההתחלה וסדר הצגת הדוגמאות). יחד עם זאת, יש מצבים שבהם ברור לנו שמפריד אחד עדיף על פני השאר.

לפי עקרון מקסימום שוליים (max-margin), נרצה למצוא על-מישור מפריד שהמרחק המינימלי שלו מהנקודות הנתונות במדגם, הוא הגדול ביותר האפשרי (מבין העל-מישורים שמפרידים את המדגם באופן נכון). ב-SVM המסווג מסתכל על ה"דמיון" בין הנקודה שקיבל לבין מספר קטן יחסית של נקודות קריטיות שהיו במדגם, המהוות את הוקטורים התומכים. הסיווג של הנקודה החדשה ייקבע לפי "מידת הדמיון" שלה לנקודות הקריטיות.

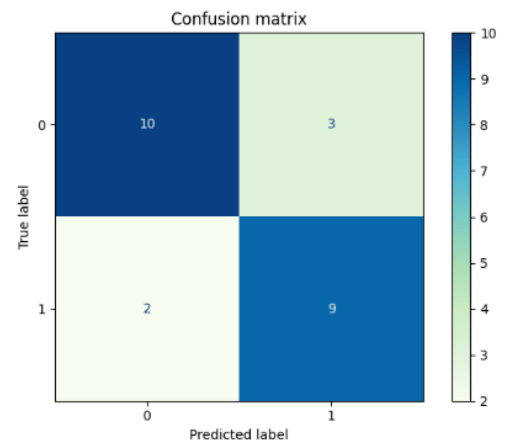
השתמשנו במימוש של soft-SVM דרך ספריית sklearn בפייתון. Soft-SVM משמש למקרה שבו הדאטה לא ניתן להפרדה על ידי על-מישור – וזה בדיוק המקרה שלנו, לפי מה שראינו לאחר הרצת הפרספטרון. אלגוריתם ה-soft-SVM

SVM מכליל את האלגוריתם המקורי בכך שהוא מרשה מספר קטן של טעויות בסיווג, ועדיין מנסה למקסם את השוליים ביחס לרוב הנקודות שכן מתוייגות נכון. מדובר בוריאציה לבעיה המקורית שבה מרשים מספר קטן של טעויות בסיווג (מה שמאפשר להסיר את ההנחה שאכן בכלל קיים מסווג לינארי עבור מדגם האימון) כדי לשמור על שוליים גדולים עבור רוב שאר הנקודות. בנוסף, SVM אינו מוגבל רק לסיווג לינארי, ויכול לבצע גם סיווג לא לינארי באמצעות הוספת kernel trick שבו בדרך כלל הקלט ממופה למרחב במימד גבוה יותר.

תוצאות

הרצנו את אלגוריתם ה-Soft-SVM על הדאטה שלנו. גם כאן לא נוכל להציג את הסיווג באופן ויזואלי. ולכן נתייחס לאחוז הדיוק של המודל. עבור סט אימון בגודל 96, וסט בדיקה בגודל 24 (20%-80%), קיבלנו שאחוז הדיוק, היחס בין הסיווגים הנכונים לבין סך כל הסיווגים של סט הבדיקה, הינו 0.79. בנוסף, נתייחס למדד נוסף שהוא ה-F1 score, העושה ממוצע משוקלל של ה-precision וה-recall. Recall הינו היחס בין דוגמאות חיוביות שזוהו נכונה (true positive) לבין כל הדוגמאות החיוביות שהמודל זיהה. Precision הינו היחס של תצפיות חיוביות שהמודל זיהה נכון מכל התצפיות שהמודל זיהה שהן חיוביות.

	precision	recall	f1-score	support
0	0.83	0.77	0.80	13
1	0.75	0.82	0.78	11
accuracy			0.79	24
macro avg	0.79	0.79	0.79	24
weighted avg	0.80	0.79	0.79	24



ב-confusion matrix ניתן לראות ש-10 דוגמאות (המהוות 77%, מתוך 13) מתוך כלל דוגמאות ה-test, שהסיווג שלהן הוא בן אדם שאינו נוטל סיכונים - סווגו נכון ע"י המודל. וכן ש-9 דוגמאות (המהוות 82%, מתוך 11) שצריכות להיות מסווגות כנוטל סיכונים סווגו נכון. כמו כן, ניתן לראות שאחוז הדיוק וציון ה-F1 score שהתקבלו גבוהים יחסית.

כעת, היינו רוצות להציג את סט הנתונים שהתקבל בגרף, כדי שנוכל לראות בעיניים את חלוקת הנבדקים שלנו לקבוצת לוקחי סיכונים וקבוצת נמנעי סיכונים. אך, נשים לב כי סט הנתונים שלנו הוא 12-מימדי, ולכן לא נוכל להציג אותו ויזואלית כפי שהוא. לכן, כדי שנוכל ליצור ויזואליזציה של הדאטה, החלטנו להשתמש באלגוריתם PCA.

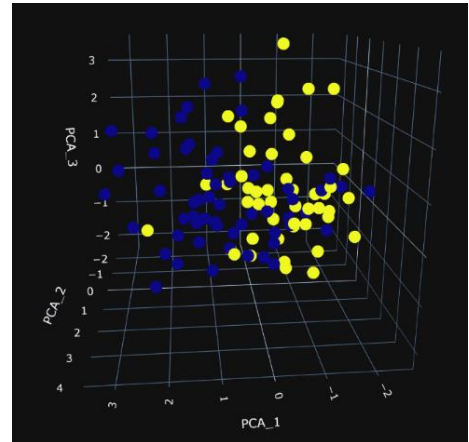
PCA היא שיטה סטטיסטית אשר ממירה את סט הנתונים שלנו ממימד p למרחב לינארי חדש ממימד d, הקטן מ-p. העיקרון המרכזי בשיטה זו הוא שהשגיאה הריבועית הכוללת של שחזור הנתונים היא מינימלית. למעשה, השיטה לא רק מבצעת מינימיזציה של השגיאה הריבועית, אלא היא מבצעת גם מקסימיזציה של השונות בסט הנתונים שלנו. כלומר, ההנחה המרכזית היא שכאשר נוריד מימד נרצה שהמימדים החדשים שלנו ידגישו את ההבדלים בין הדגימות (השונות). נזכיר כיצד עובד אלגוריתם ה-PCA:

1. נבחר היפר-פרמטר d להוריד מימד אליו (אנחנו בחרנו d=3). כלומר, אנחנו נעבור ממרחב 12 מימדי למרחב תלת מימדי.
2. מרכז הדאטה
3. נחשב את מטריצת הקורלציה C
4. נמצא את הערכים העצמיים של C, ונמין את הוקטורים העצמיים בסדר יורד לפי הערכים העצמיים שלהם
5. ניקח את d הוקטורים העצמיים שמתאימים ל-d הערכים העצמיים הגדולים ביותר
6. מטילים את הדאטה ממרחב 12 מימדי למרחב תלת מימדי שנפרש על ידי d הוקטורים שמצאנו

תוצאות

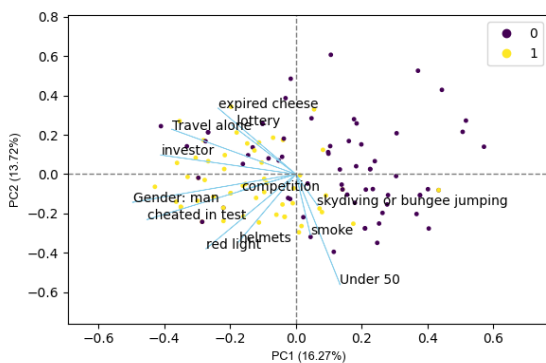
להלן התוצאות שקיבלנו לאחר הרצת PCA על הדאטה, והורדה של הנתונים ממרחב 12 מימדי למרחב תלת מימדי:

לפנינו 120 נקודות במרחב תלת מימדי, כאשר כל נקודה מייצגת נבדק. בצהוב מסומנים הנבדקים שהעידו על עצמם כנוטלי סיכונים, ובכחול מסומנים הנבדקים שהעידו על עצמם כנמנעי סיכונים. במצב אידיאלי, לפי ההשערה שלנו היינו מצפים לראות חלוקה ברורה, כלומר מישור מפריד בין שתי הקבוצות. לפי התוצאות שקיבלנו שתואמות גם את מה שעלה מהמסווגים, נראה כי קיימת חלוקה כלשהי, אם כי לא חד משמעית.

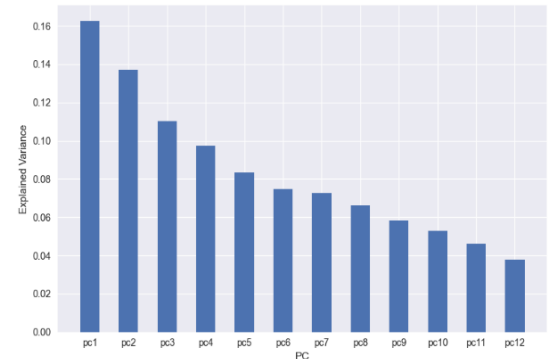


ההיסטוגרמה להלן מציגה את השונות המוסברת עבור כל PC. ניתן לראות כי ל-PC1 יש 0.162 שונות מוסברת, ל-PC2 יש 0.137 שונות מוסברת ול-PC3 יש 0.11 שונות מוסברת.

בטבלה, מתוארים הוקטורים העצמיים הראשונים עם הערכים הגדולים ביותר, ושלושת המקדמים הגדולים ביותר בערך מוחלט (כדי לשקף את המתאם הגבוה, בין אם שלילי או חיובי) של כל אחד מהם המייצגים את השאלה לצידם. ה-biplot מציג את הדאטה המוטל על המרחב שנפרש ע"י PC1 ו-PC2 והוקטורים המתאימים למקדמים, ומשקף ויזואלית את הטבלה. בנוסף, הוא מאפשר לנו לראות את הוקטורים שמייצגים את המשתנים שביניהם ישנה קורלציה חיובית חזקה, למשל: בין מגדר לבין האם אי פעם רימית במבחן.



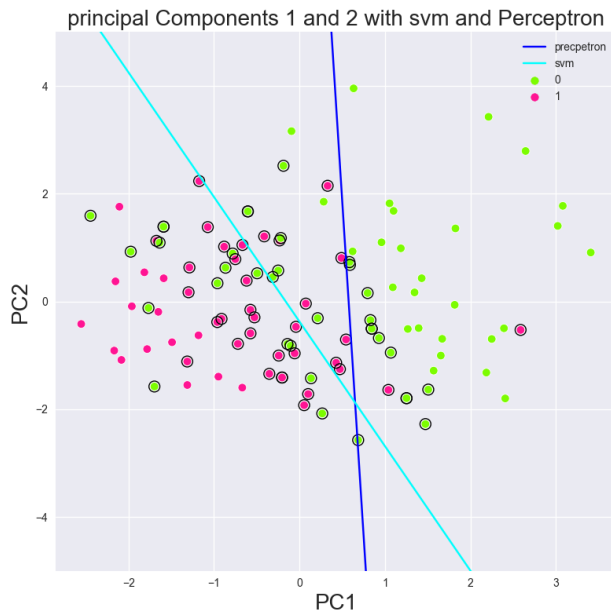
PC2	PC1
0.55 האם את/ה מתחת לגיל 150?	0.49 מגדר
0.37 האם אי פעם חצית במעבר חציה באור אדום?	0.44 האם אי פעם רימית במבחן?
0.34 האם את/ה רוכבת/על אופניים/קורקינט ללא קסדה?	0.40 האם את/ה משקיע/ה בבורסה?



לאחר שבחנו את שלושת האלגוריתמים על הדאטה שלנו, נרצה לראות ויזואליזציה של המישורים המפרידים בהרצת פרספטרון ו-SVM. אך, מאחר שאנחנו נמצאים במרחב 12 מימדי, כדי לעשות זאת, בחרנו להשתמש ב-PCA.

כפי שהראנו בהתחלה, ה-PCA החזיר עבורנו את הוקטורים העצמיים המתאימים לערכים העצמיים של מטריצת הקורלציה. כעת, נבחר בשני הוקטורים העצמיים המתאימים לערכים העצמיים הגדולים ביותר, כפי שראינו בגרף לעיל, שני הרכיבים הראשונים תורמים לכ-30% מכלל השונות. נטיל את הדאטה שלנו, שנמצא במרחב 12-מימדי, למרחב דו מימדי, שנפרש על ידי שני הוקטורים העצמיים שבחרנו. נרץ את אלגוריתם הפרספטרון ואלגוריתם ה-SVM על הדאטה החדש, המוטל. בגלל שאנחנו נמצאים במרחב דו מימדי, נוכל לראות בעיניים את המישורים המפרידים שמתקבלים מהרצת האלגוריתמים.

ניתן לראות ויזואלית שלא ניתן ליצור הפרדה מוחלטת בין שתי המחלקות. גם כאן, הפרספטרון לא התכנס, אך ניתן לראות שיחסית הגיע להפרדה מסוימת. הנקודות המוקפות בעיגול מהוות את הוקטורים התומכים – היושבים "על השוליים" של המפריד שיצר ה-SVM.

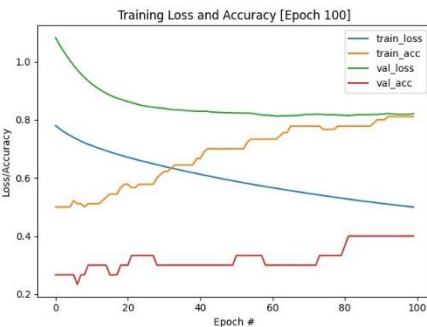


לבסוף, החלטנו להריץ מסווג לא לינארי על הדאטה שלנו, ולראות האם נצליח לקבל תוצאות טובות יותר. בחרנו ברשת נוירונים עמוקה, הגדרנו שכבה חבויה אחת עם 8 כניסות. לאחר 15 epochs קיבלנו רמת דיוק של 0.609, המהווה תוצאה פחות טובה ממה שקיוונו. בבחינה של למידת רשת הנוירונים לאורך 100 epochs עולה כי רמת הדיוק על ה-training set הולכת וגדלה, בעוד שרמת הדיוק על ה-validation set (10% מהדאטה) מגיעה לפלטו ולא משתפרת. כך גם לגבי ה-loss, בעוד שה-loss של ה-training set הולך ויורד, ה-loss של ה-validation set נותר גבוה. ניתן להסביר זאת בכך שהדאטה שלנו מורכב מכמות קטנה של דוגמאות (רק 120 נבדקים), מה שמקשה על רשת הנוירונים לבצע למידה איכותית ומהימנה ואף יכול להוביל ל-overfitting לאורך זמן.

סיכום:

אמנם, שאלת המחקר שלנו עוסקת בניבוי האם בן-אדם נוטה לקחת סיכונים או לא. אך במסגרת החיפוש אחר התשובה לשאלתנו, אחת השאלות המרכזיות שנתקלנו בה הייתה האם הדאטה שאספנו על מנת לענות על שאלת המחקר פריד לינארית, שכן התשובה לשאלה זו תאפשר לנו למצוא את האלגוריתם שייתן את המענה האופטימלי לשאלת המחקר. התחלנו ממסווג לינארי פשוט בעזרת פרספטרון ואכן, מתוצאות הריצה של אלגוריתם הפרספטרון, עלה כי הדאטה אינו פריד לינארית באופן מובהק. בהינתן התוצאות הלא מספקות הנ"ל ועל מנת לבחון את השאלה מזווית נוספת, הרצנו על הדאטה גם את אלגוריתם ה-SVM אשר תוצאותיו לא היו מושלמות (עם דיוק של 0.79). כדי לעשות ויזואליזציה של המישורים המפרידים ב-SVM ובפרספטרון בחרנו להשתמש ב-PCA. הטלנו את הדאטה למרחב דו מימדי, והרצנו את האלגוריתמים על הדאטה המוטל. גם כאן לא הגענו לתוצאות מושלמות, אך כן ראינו שיש חלוקה יחסית הגיונית של המרחב לשני חלקים. השלב הבא, היה לבחון את תוצאות הריצה של מסווג לא לינארי, בהינתן אותו הדאטה. תוצאות הלמידה של רשת נוירונים עמוקה לא אפשרו לנו לקבל ניבוי אופטימלי, מאחר וכמות הדוגמאות היתה קטנה.

בנושא איסוף הדאטה, לא ערכנו בקרה על פיזור משתנים כמו גיל ומגדר, באופן אחיד על פני הנבדקים. בנוסף, בבחירת השאלות, ייתכן שהיה מקום לדייק יותר את השאלות לתחום מסוים של נטילת סיכונים. אנו מאמינות שאיסוף דאטה ממספר רב יותר של נבדקים יאפשר בחינה של הסיווג ברשת נוירונים עמוקה והשוואה איכותית יותר אל מול תוצאות הפרספטרון וה-SVM. כמחקר המשך היינו רוצות להשוות את התוצאות שקיבלנו למסווג לא לינארי נוסף.



Bonem, E. M., Ellsworth, P. C., & Gonzalez, R. (2015). Age Differences in Risk: Perceptions, Intentions and Domains. *Journal of Behavioral Decision Making*, 28(4), 317–330.
<https://doi.org/10.1002/bdm.1848>

Zuckerman, M., & Kuhlman, D. M. (2000). Personality and Risk-Taking: Common Biosocial Factors. In *Journal of Personality* (Vol. 68).

Universität Basel. (2017, October 30). Willingness to take risks: A personality trait. ScienceDaily. Retrieved April 1, 2021 from
www.sciencedaily.com/releases/2017/10/171030095706.htm