

Introduction To Machine Learning

פרויקט סופי -דו"ח מסכם:

מגישות:

לילך בהן, ת"ז: 209492487

נועה לוי, ת"ז: 209226687

24 מספר קבוצה:



כדי להתחיל את התהליך, טענו את קובץ ה- train ובדקנו איך הדאטה נראה, הסתכלנו על נתונים סטטיסטיים וקיבלנו תמונת מצב ראשונית, ובדקנו איזה סוגי פיצ'רים הדאטה מכיל. זיהינו 3 סוגי פיצ'רים:

- פיצ'רים קטגוריאליים- בתוכם פיצ'רים עם מספר רב של ערכים מסוג 'file_type_trid' (object ו- 'c').
- פיצ'רים בינאריים- קטגוריאליים המכילים ערכים 0 או 1 (למשל 'has_relocations', 'has_debug').
- פיצ'רים נומריים- פיצ'רים עם ערכים רציפים. (למשל 'size', 'vsize').

לכל קבוצת פיצ'רים התאמנו את השלבים בניתוח הראשוני של הדאטה ובעיבוד המקדים. ולמען הנוחות נאזרכת כל קבוצה בשמות שהוגדרו כאן.

בנוסף הסרנו את העמודה '256sha' שמכילה מזהה ייחודי לכל קובץ, כלומר ערך אחר שונה בכל שורה. בשלב האימון היא לא רלוונטית.

זיהוי התנהגות קורלטיבית בין פיצ'רים- יצרנו 3 מטריצות קורלציה

1. מטריצת קורלציה כוללת על ה- train סט - שמנו לב שיש קורלציה גבוהה בין המשתנים 'size', 'numstrings' ובין 'MZ', 'size' ובין 'numstrings' ל 'MZ'. מכאן הסקנו שאולי בהמשך נוכל להוריד חלק מהמשתנים כדי להפחית מימדיות מבלי לפגוע בתחזית. (1.נ).
2. שתי מטריצות קורלציה נוספות כאשר הנתונים מופרדים לפי label: אחת על ה- train כאשר ה- label = 1, ואחת נוספת בה label = 0 (נ. 3 ואז 2 בהתאמה).

ראינו שכאשר label = 0, הקורלציות בין הפיצ'רים גבוהות משמעותית מאשר ב- label = 1. כלומר יש קורלציה גבוהה יותר בין הפיצ'רים כאשר הקובץ הוא לא זדוני.

התפלגות של הפיצ'רים הבינאריים בכל label- בדקנו בכל משתנה בינארי איך הוא מתפלג (נ. 4-8) לכל label ומצאנו שלכל המשתנים הבינאריים, חוץ מ 'has_relocations' (נ. 5), יש הבדל משמעותי בין ה- labels. המשמעות של ההבדל הגדול יכולה להיות שאלו פיצ'רים איכותיים עבור התחזית. זיהינו מספר פיצ'רים כאלה- ביניהם: has_resources, has_tls.

התפלגות של הפיצ'רים הנומריים בכל label (נ. 9-19) השתמשנו ב- boxplot כדי לבחון את הפיצ'רים ובכך הנחנו שהם מתפלגים נורמלית, אף שאנחנו לא יודעות זאת בשלב זה. בחינה של הפלוטים מאפשרת לנו להסיק מספר מסקנות. ניתן לראות כי יש הבדל בהתפלגות בין כל label בפיצ'רים avlength, MZ, לעומת זאת בפיצ'ר A ניתן לראות כי אין הבדל משמעותי. בנוסף ניתן לראות שרובם מכילים כמות גדולה של outliers.

חלק שני- עיבוד מקדים – preprocessing

עבור כל קבוצה של משתנים ביצענו את התהליכים הבאים:

טיפול בערכים חסרים: None בפיצ'רים הקטגוריאליים file_type_trid ו- C החלפנו None במחרוזת 'NONE_VAL'. בפיצ'רים הקטגוריאליים הבינאריים החלפנו את הערכים החסרים ב- 0. בפיצ'רים הנומריים החלפנו ערכים חסרים בחציון.

התמודדות עם outliers: עדיף להסיר ערכים חריגים, ובכך אומנם נעלה מעט את ההטייה (bias) אך ה- variance ירד בצורה משמעותית ונקבל דאטה שניתן לעבוד איתו בצורה יותר נכונה.

בפיצ'רים קטגוריאליים- ב- file_type_trid - זיהינו שיש ערכים שחוזרים מספר רב של פעמים, ויש נתונים המופיעים מספר נמוך של פעמים או פעם אחת. במקרה זה החלטנו שנשמור את הערכים השכיחים ואילו ערכים 'נדירים' נאחד לערך 'other'. בחירת מספר הערכים שנשמור נעשתה מהתבוננות ב- barplot שמציג בכמה שורות כל ערך מופיע (נ. 20). בפיצ'ר זה בחרנו לשמור את 30 הערכים השכיחים ביותר, (וב 'C'- 21) שמרנו על 6 הערכים השכיחים ביותר.

בפיצ'רים נומריים- מהתבוננות בהתפלגות של הפיצ'רים הנומריים (נ. 22) זיהינו כי ההתפלגות שלהם היא skewed distribution, התפלגות הדומה לנורמלית, אך אחד הזנבות ארוך משמעותית מהאחר. לאחר שזיהינו מאפיין זה בנינו

פונקציה להתמודדות עם outliers המיועדת להתפלגות מסוג זה באמצעות (capping & Inter-Quartile Range (IQR). (קביעת גבול עליון ותחתון במבוסס על ההתפלגות של הפיצ'ר (אחוז 75 ואחוז 25) ושינוי הערך בחריגים החורגים מהגבול לערך הגבול).

לאחר שהשתמשנו בשיטה זו על כל הפיצ'רים ראינו שהיא הייתה אגרסיבית מדי עבור חלק מהם, והם איבדו הרבה מידע. לכן השתמשנו בשיטה נוספת להתמודדות עם outliers - percentile method. כך חילקנו את הפיצ'רים הנומריים לשלוש קבוצות ובכל אחת עשינו הסרת outliers מותאמת:

עבור תת-קבוצה 1 numeric או משתמשים ב-IQR - capping & Inter-Quartile Range המתאים להתפלגות המוטה. עבור תתי קבוצות 2 numeric ו-3 numeric או משתמשים בשיטת percentile, כאשר לכל קבוצה הזנו אחוז שונה. (שיטת האחוזון היא קביעת אחוזון שווה בשני הצדדים כדי לזהות חריגים, והחלפתם בערך הגבול). בנספחים ניתן לראות boxplot ו-kde של כל אחד מהפיצ'רים הנומריים לפני ואחרי הטיפול בערכים חריגים. לפני (נ.26-37) אחרי (נ.33-53) (נ.38) ניתן לראות את ההיסטוגרמות של כולם אחרי.

הערכת איכות פיצ'רים ויצירת פיצ'רים חדשים והקטנת מימדים:

פיצ'רים קטגוריאליים - בפיצ'רים קטגוריאליים המכילים ערכים שאינם מספריים (מילים) היינו צריכות להחליט על שיטה לקידוד (encoding) המשתנים-

השיטה הראשונה שניסו היא יצירת עמודות dummies, בפועל שיטה שיוצרת עמודות כמספר הערכים השונים בפיצ'ר. כפי שצינו קודם בעמודות האלה (C, file_type_trid) יש מספר רב של ערכים שונים, גם לאחר ההתמודדות עם ה-outliers. לכן שימוש בה מוסיף כ-35 עמודות ומגדיל משמעותית את מימדות הדאטה ופוגע בתחזית. לכן עברנו לשימוש בשתי השיטות האחרות.

השיטה השנייה לקידוד הקטגוריאליים היא פשוט להיפטר מהם! ניסינו גם שיטה זו, (הדאטה לאחר הסרת הקטגוריאליים נקרא (X_train_drop_cat, X_test_drop_cat), בהמשך נתייחס לתוצאות המודלים בהרצות עם דאטהסט שלא מכיל את המשתנים הקטגוריאליים.

השיטה השלישית שניסו היא ordinal encoding. בשיטה זו מחליפים את הערכים המילוליים בערכים מספריים. מתוך הבנה שבהחלפה של הערכים המספריים אנחנו יוצרות יחס סדר בין הערכים השונים, בחרנו להחליף כל ערך במספר הפעמים שהוא מופיע בסט ה-train. כך ערכים שכיחים יותר יקבלו מספרים גדולים יותר. יחס הסדר שנוצר מבוסס על שכיחות הערך ומספר משהו על הדאטה המקורי.

בנוסף בנינו פיצ'ר חדש 'interactions_file_type_trid_C' המכיל את השילובים השונים בין ערכי שני הפיצ'רים. גם אותו קודדנו. בהמשך השתמשנו במבחן chi square על המשתנים הקטגוריאליים כדי להבין את האיכות של כל אחד מהם לתחזית, וראינו כי C_count_encoded ו-interactions_file_type_trid_C_count_encoded משמעותית פחות טובים מהאחרים ולכן ויתרנו עליו. בנוסף, יצרנו פונקציה הבוחרת את 10 המשתנים הנומריים הטובים ביותר בהתבסס על מבחן ANOVA F-value. בעקבות השימוש בפונקציה הוסרו ['A', 'imports', 'printables', 'numstrings'] באמצעות הפונקציה selectkbest של sklearn עם מבחן anova f-value.

פיצ'רים בינאריים - לאחר מספר ניסיונות הגענו למסקנה שהסרה של פיצ'ר מסוג זה פוגעת בתחזית.

עשינו סטנדרטיזציה על שני הסטים באמצעות StandardScaler של sklearn. החשיבות של שלב זה היא עבור הרגרסיה הלוגיסטית ו-PCA בהם לא נרצה שהבדל בסדרי הגודל של הערכים ישפיעו על התוצאות.

בסופו של דבר הגענו לשלב הרצת המודלים עם 4 סטים של דאטה:

- דאטה שמכיל רק את הפיצ'רים הבינאריים והנומריים, שגם עברו סינון ע"י selectkbest עם מבחן anova f-value
- דאטה שמכיל את כל סוגי הפיצ'רים לאחר שהקטגוריאליים והנומריים עברו סינון לפי השיטות שתיארנו.

- דאטה שמכיל את כל הפיצ'רים ללא סינון.
- וסט נוסף שמכיל את כל הנתונים ללא סינון ועם pca.

מימדיות גדולה בעייתית מכיוון שהיא מובילה לשונות גדולה (variance) וזה יגרום ל over fitting שיפגע בתחזית על דאטה חדש (!test). בנוסף כאשר יש הרבה מימדים, התצפיות יותר קשות לסיווג כי כל התצפיות נהיות במרחק שווה מהאחרות ואז לא ניתן לסווג בצורה שמציגה את המציאות. יש מספר סימנים מעידים למימדיות גבוהה: זמן ריצה גבוה (אך יכול להיגרם גם בגלל אלגוריתמים מורכבים או כי יש הרבה דאטה), אם קשה למצוא דפוסים ויחסים משמעותיים בנתונים (כמו שהסברנו לעיל שהמרחק בין הנקודות נהיה גדול וזהו ולכן יותר קשה לסווג למחלקות), אם התחזית לא טובה (מימדיות גדולה עלולה להביא ל over fitting והמודל שהתאמן על סט הדאטה ילמד גם את הרעש ולכן יביא ביצועים לא טובים בסט הבדיקה). ההשפעה על המודל מפורטת בהמשך.

חלק שלישי: בחירת היפרפרמטרים והרצת מודלים. (בנספח 78 יש הסבר על bayes שניסינו ובסוף לא השתמשנו בו.)

בכל המודלים למעט knn השתמשנו ב GridsearchCV, הסבר על השימוש הספציפי בכל מודל נמצא בנספחים *77. נתאר את ההיפרפרמטרים שבחרנו עבור כל מודל:

מודל רגרסיה לוגיסטית:

'penalty': הוספנו למודל כפרמטר את 'penalty' כדי לבצע רגולריזציה על מנת למנוע\ לצמצם overfitting. בדרך זו הוספנו מעט bias (הטייה) למודל על מנת לצמצם משמעותית את השונות. בחרנו להשתמש ב ridge וגם ב lasso. הידוע גם כ 'lasso'. בחרנו לא להוסיף את 'elasticnet' מכיוון שהוא ממילא מוסיף גם את l וגם את 2, אך יש solvers שלא עובדים איתו.

'C': הוספנו למודל את C הבאים: [10,100,200,300,400] כאשר ככל ש C גדל, הרגולריזציה פחות משמעותית. כלומר ככל ש C גדל, השונות גדלה, וההטייה קטנה. C חשוב לרגולריזציה כי הוא מגביל את המשקולות של הרגולריזציה.

'solver': השתמשנו ב 'liblinear' וב 'lbfgs' מכיוון ש 'liblinear' תומך ב l וב (lasso&ridge,2), אך יש solvers שלא עובדים שמחזק שהוא lbfgs. 'liblinear' אבל ה 'liblinear'.

לרגרסיה לוגיסטית יש עוד המון היפרפרמטרים שניתן להוסיף, אך משיקולי זמן ריצה בחרנו רק את אלה.

הגיע לתוצאה הטובה ביותר עם דאטה סט ללא משתנים קטגוריאליים ועם בחירה של המשתנים הטובים ביותר מבין הבינאריים.

מודל KNN:

'n_neighbors' - מספר השכנים הקרובים ביותר. בחנו השפעה של מס' k שונים ובחרנו את הטוב ביותר. בשאר ההיפרפרמטרים האפשריים במודל זה לקחנו את ה- default שעבדה היטב. ככל ש- k גדל, כך גדלה ההטייה. עבור k גדול מדי, קשה לקבל "תמונה חדה" על הסיווג האמיתי. לעומת זאת, עם הקטנת K השונות גדלה, ובא קטן מדי קיים סיכוי ל overfitting (ראה דוגמא k=1)

הגיע לתוצאה הטובה ביותר עם דאטה סט בו בחרנו את המשתנים הטובים ביותר.

מודל multi layer perception: משיקולי זמן ריצה ויעילות בבחירת ההיפרפרמטרים, בחרנו בבאים:

'activation' מייצג את האקטיבציה לשכבות הנסתר (לדוגמה relu, logistic וכו'). כאשר נבחר באקטיבציה לא לינארית כמו סיגמואיד נוכל לייצג את המציאות בצורה יותר מדויקת ובכך להפחית את ההטייה, ומנגד צריך להיזהר שהמודל לא מתאמן עם אקטיבציות אלו יותר מדי זמן ומגדיל את השונות ויוצר overfitting (במידה וכן, נשנה אקטיבציה\ נעניש בעזרת ההיפרפרמטר. 'alpha').

'solver': בדקנו את שלושת lbfgs, 'adam', 'sgd' solvers: ובחרנו לוותר על 'lbfgs' הוא לא תרם לטיב התוצאות אך הכפיל ביותר מפי 2 את זמן הריצה.

verbose": ברירת המחדל היא False ושינוי ל True מכיוון שרצינו לצפות בהתקדמות.

max_iter": ברירת המחדל היא 200 ושינוי ל 1000 מכיוון שראינו הבדלים משמעותיים בתוצאות כאשר הרשנו לאלגוריתם לרוץ מספר רב יותר של איטרציות עד ההתכנסות, ועם זאת זמן האלגוריתם היה סביר. בכל איטרציה המודל מתאים את הפרמטרים כדי להוריד את training loss ולשפר את הביצוע. ככל ש' max_iter' גדל, כך המודל מותאם יותר לסט האימון, והטיית המודל קטנה, אך חשוב להיזהר כי לפעמים מאימון יתר יכול להיגרם overfitting (ולכן השונות גדלה מכל איטרציה).

'early_stopping': היפר פרמטר זה מאפשר למודל להפסיק לרוץ לפני כמות האיטרציות שנקבעו, במידה וציון ה validation לא משתפר ברמה משמעותית. בזכות בחירה זו זמני הריצה התקצרו משמעותית, ונמנע מצב של overfitting כי המודל לא ממשיך לרוץ על סט האימון ללא צורך. אמנם קיים סיכון שההטייה תגדל כי המודל עלול לא להגיע לתוצאה המיטבית, אך מדיקה שביצענו על ריצת המודל עם ההיפר פרמטר ובלעדיו, ההבדל זניח למדי.

'hidden_layer_sizes', 'alpha', 'learning_rate_init': השתמשנו בברירת המחדל לאחר שניסו מספרים אחרים שהפיקו תוצאות פחות טובות והגדילו משמעותית את זמן הריצה.

'hidden_layer_sizes' מייצג את ארכיטקטורת של רשת הנוירונים, 'alpha' מייצג את כמות הרגולריזציה (כדי למנוע overfitting), 'learning_rate_init' (המספר הראשוני לקצב הלמידה).

הגיע לתוצאה הטובה ביותר עם דאטה סט בו יש את המשתנים, ועבר pca.

מודל - Random Forest

'n_estimators' - הוא מספר עצי ההחלטה. ככל שיש יותר עצים, כך השונות תקטן אך ההטייה לא תשתנה ולכן נרצה שהיפר פרמטר זה יהיה כמה שיותר גדול, במגבלות הזמן ריצה האפשרי.

'max_depth' - מייצג את העומק המקסימלי של העצים ב-randomforest. עומק גדול מדי עלול ליצור overfitting מכיוון שעץ עמוק מדי עלול ללמוד גם דברים שלא צריכים להיכלל כמו רעש \ outliers. עם זאת לא נרצה עומק קטן מדי, אז הסיווג יסתמך על פחות מידע ותגדל ההטייה.

'criterion' - ה' log_loss' לקחה זמן רב מדי ולא הביאה תוצאות טובות, ולכן בדקנו את האופציות האחרות. אין השפעה משמעותית על ההטייה והשונות- אלו בעיקר דרכים שונות ליצירת כל חלוקה בעץ.

הגיע לתוצאה הטובה ביותר עם דאטה סט בו יש את כל המשתנים.

חלק רביעי – הערכת מודלים:

confusion matrix: עשינו לכל אחת מהדאטה סטים שיצרנו confusion matrix *71-75*, *79 שמסביר לעומק ולא רק מצרך פלט*. באופן כללי אצלנו ה TP מייצג שחזינו שהקובץ דדוני והוא באמת דדוני, FP מייצג שחזינו שהקובץ דדוני אבל הוא למעשה לא דדוני, FN מייצג שחזינו שהקובץ לא דדוני אך הוא דדוני, TN מייצג שחזינו שהקובץ לא דדוני והוא באמת לא דדוני. ניתן לראות שבכל האופציות קיבלנו מספרים יחסית דומים: ה TP וה TN מהווים את רוב המסה כלומר לרוב אם חזינו שמשנהו הוא דדוני הוא באמת דדוני, ואם חזינו שהוא לא דדוני אז הוא בד"כ לא דדוני. ולכן בכל המדדים שבדקנו (recall, precision, accuracy) התוצאה תמיד גדולה מ 0.88 (הסבר על כל מדד נמצא בנספחים), ומעידה על כך שביצועי המודל טובים.

ROC AUC K FOLD (נ.54-76) כל מודל הערכנו בעזרת K-Fold cross validation. ובנינו פלט ROC על כל K-Fold וקיבלנו בצורה ויזואלית אישור לטיב המודל שלנו, וגם תוצאות מספריות שהראו שכל המודלים הביאו תחזית טובה, בכולם למעט הרגרסיה הליניארית עם תוצאה מעל 0.9 כלומר המודלים מסווגים בצורה טובה את מרבית הדוגמאות. המודל היחיד שמסווג בצורה מעט פחות טובה הוא מודל הרגרסיה הלוגיסטית שאת ביצועיו לא הצלחנו לשפר תחת זמן ריצה סביר, יתכן שהוא לא מתאים לדאטה שלנו.

ניתן לראות כי לכל מודל יש תוצאה טובה ביותר עבור סידור שונה של הדאטה. כך למשל knn מגיע לשיא (0.945 auc) עם דאטה סט בו בחרנו את המשתנים הקטגוריאליים והנומריים בטובים ביותר ולעומתו random forest הגיע לשיא (0.976 auc) עם דאטה סט בו לא ביצענו בחירה של משתנים.

ניתן לראות כי בשלושת המודלים: רגרסיה לוגיסטית, mlp, knn הורדת המימדיות שיפרה את התוצאות והם הגיעו לתוצאה הטובה ביותר עם דאטה סט שעבר הורדת מימדיות. לעומת זאת random forest הגיע לתוצאה הטובה ביותר שלו כאשר לא הורדנו מימדים, שהיא גם התוצאה הכי טובה בהשוואה לכל המודלים (עם ובלי הורדת מימדים). מכאן אנחנו מסיקות שלכל מודל מתאים דרך אחרת מבחינת הורדת מימדיות ובחירת משתנים.

לא היה פערי ביצועים משמעותיים בין train validation - ניתן לראות זאת בפלטים של roc curve שהפעלנו על הולדייטשן, ובגלל חוסר הפערי ביצועים ניתן להניח שהמודל לא overfitted. הכחול מבוסס על validation set, הצבעוני לוקח מהחלק של k-fold. (דוגמא בנספחים 81 82)

בחירת המודל הסופי- לאחר שקלול התוצאות בחרנו להשתמש ב- random forest ולסדר את הדאטה כך שהוא מכיל גם את הפיצ'רים הקטגוריאליים שקודדנו באמצעות ordinal encoding וגם את הפיצ'רים הנומריים והבינאריים ללא שום סינון של פיצ'רים. זה המודל שהביא לנו את ה- auc הגבוה ביותר וגם את הroc הטוב ביותר ולכן בחרנו בו.

ביצירת המודל השתמשנו בכלים להתמודדות עם משתנים קטגוריאליים וכלים להתמודדות עם **outliers** אותם לא למדנו בכיתה ועליהם הרחבנו בתיאור העיבוד המקדים. תודה מעומק הלב לאתר **kaggle**, התנ"ך החדש שלנו, שבלעדיו כל זה לא היה אפשרי. קישורים למאמרים ספציפיים באתר נמצאים בנספחים.

התרומה של כל אחת:

את כל הפרויקט עשינו בשיתוף פעולה וייעוץ תמידי אחת עם השנייה כאשר לכל אחת היה תחומי אחריות עליהן היא שמה יותר דגש. את החלק הראשון והרביעי עשינו ביחד, על החלק השני והחמישי לילך היתה אחראית, על החלק השלישי ובתיבת הדו"ח נועה היתה אחראית.

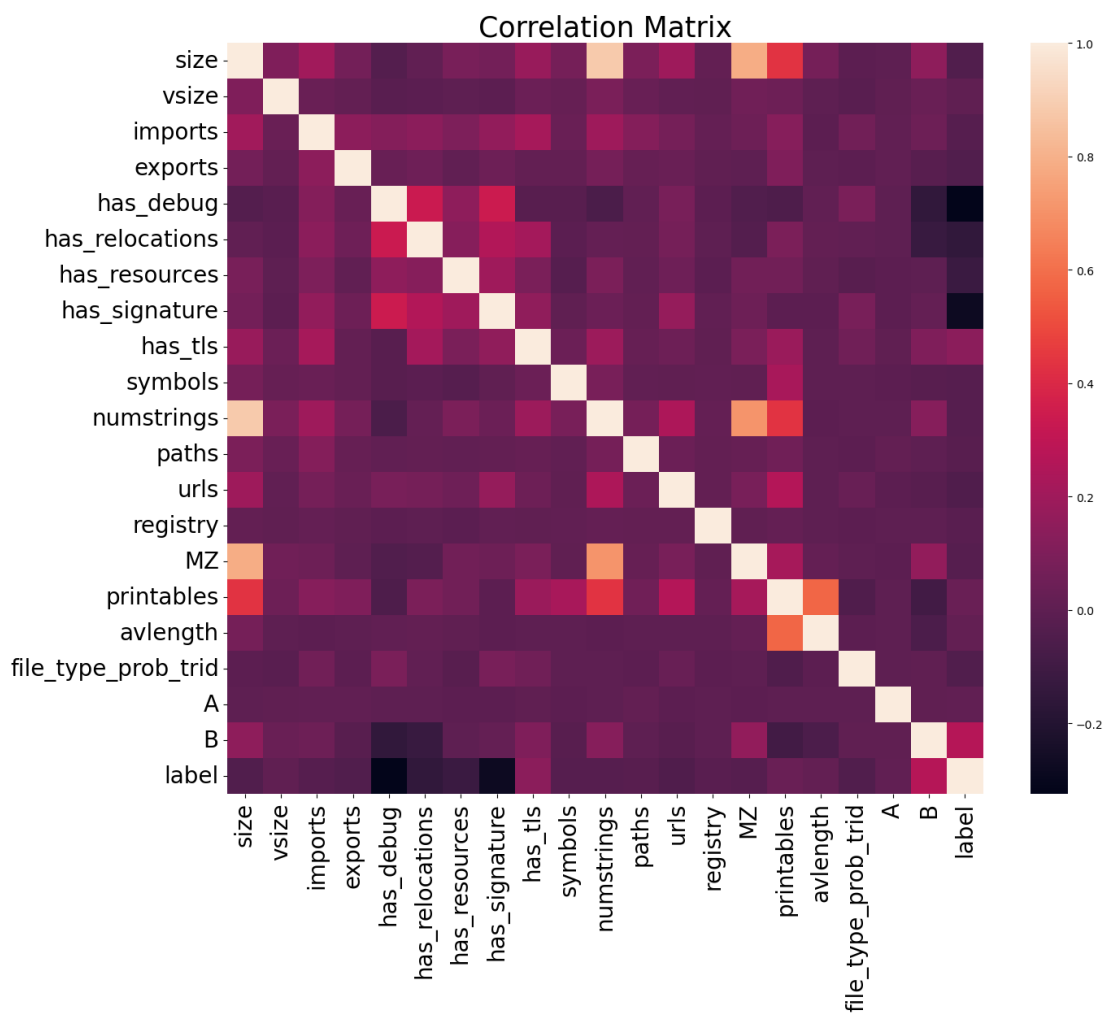
<https://www.kaggle.com/code/aimack/how-to-handle-outliers>

[ROC and AUC | Kaggle](#)

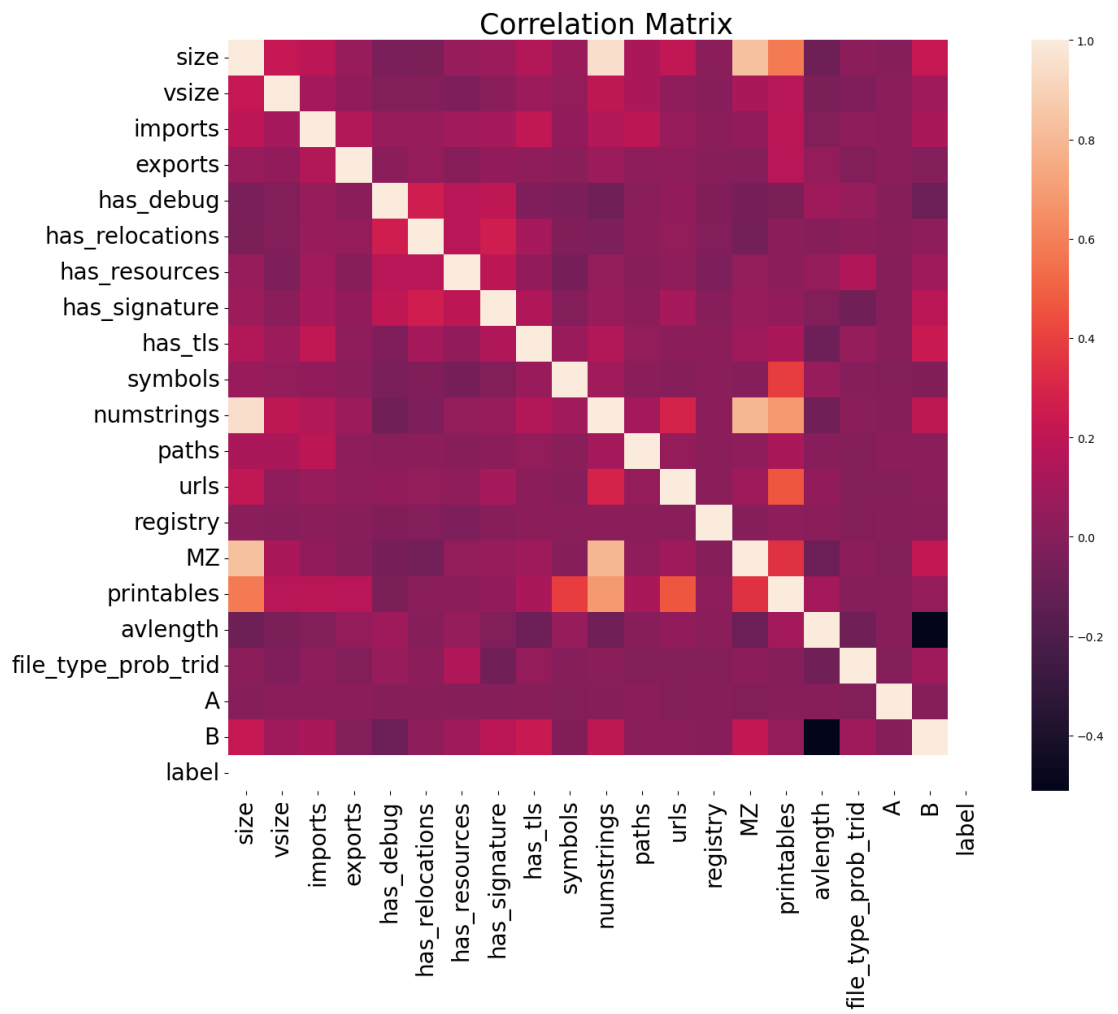
[Comprehensive Guide on Feature Selection | Kaggle](#)

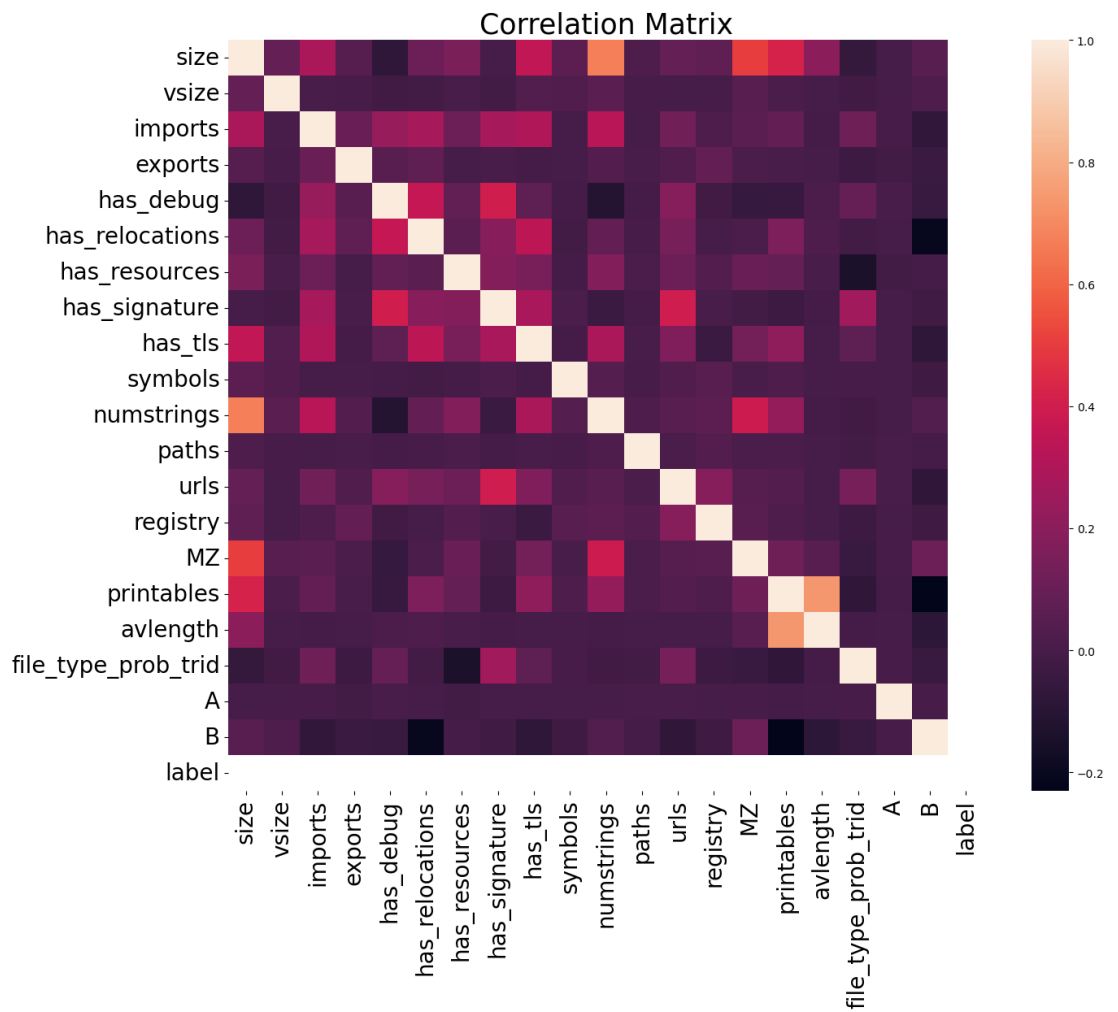
כל הפלטים של המחברת ממוספרים:

1.

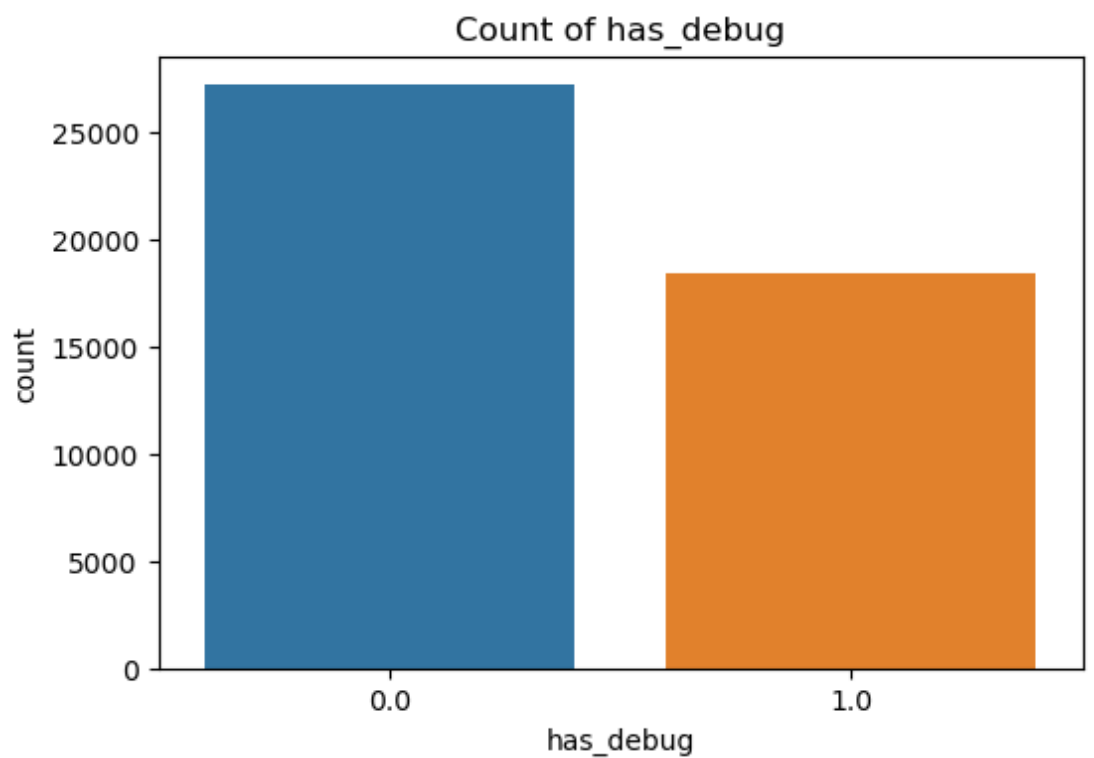


2.

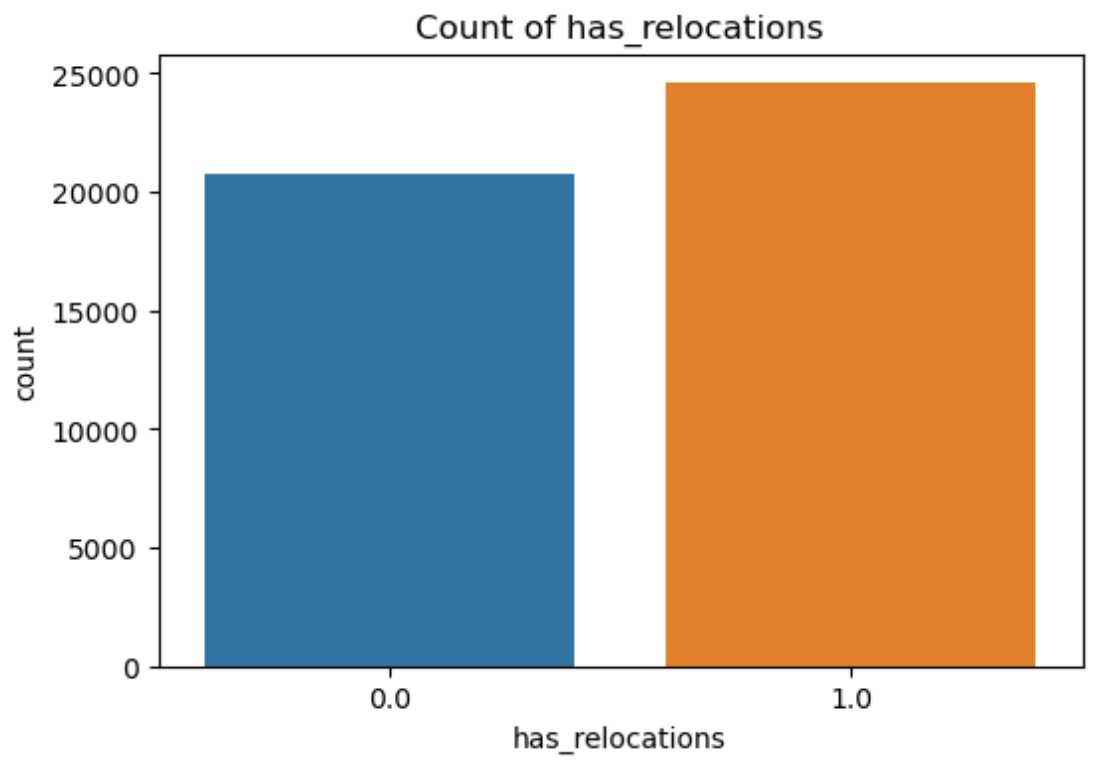




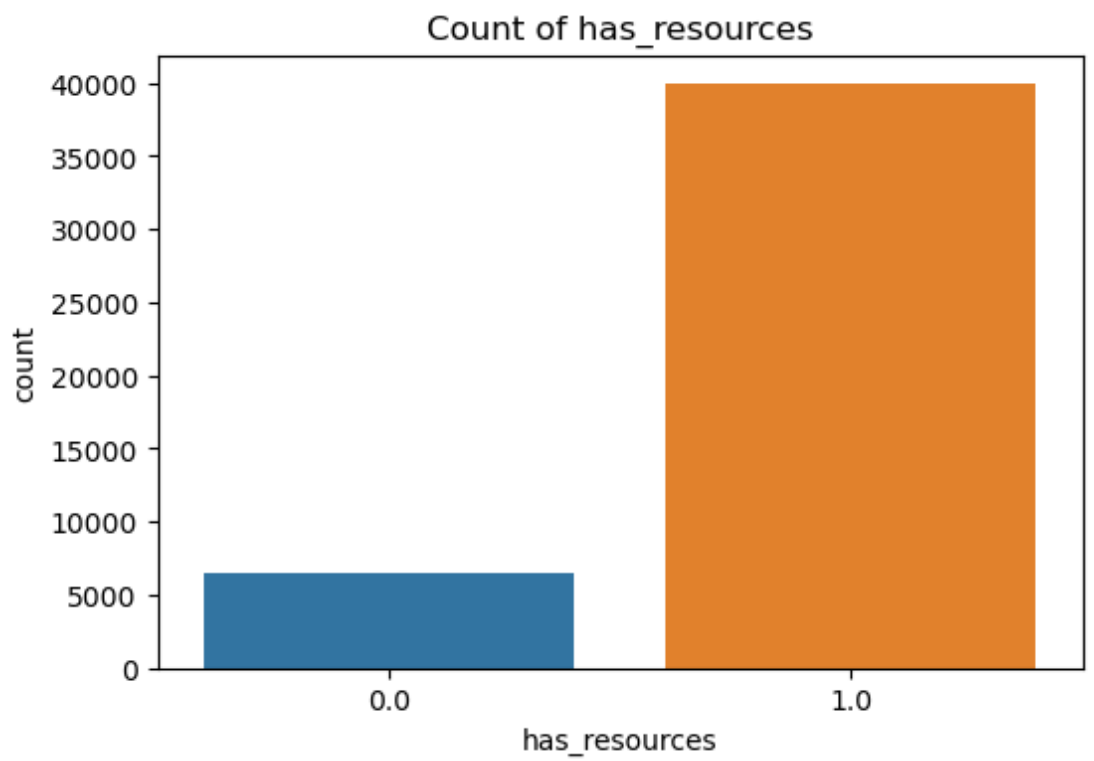
4.



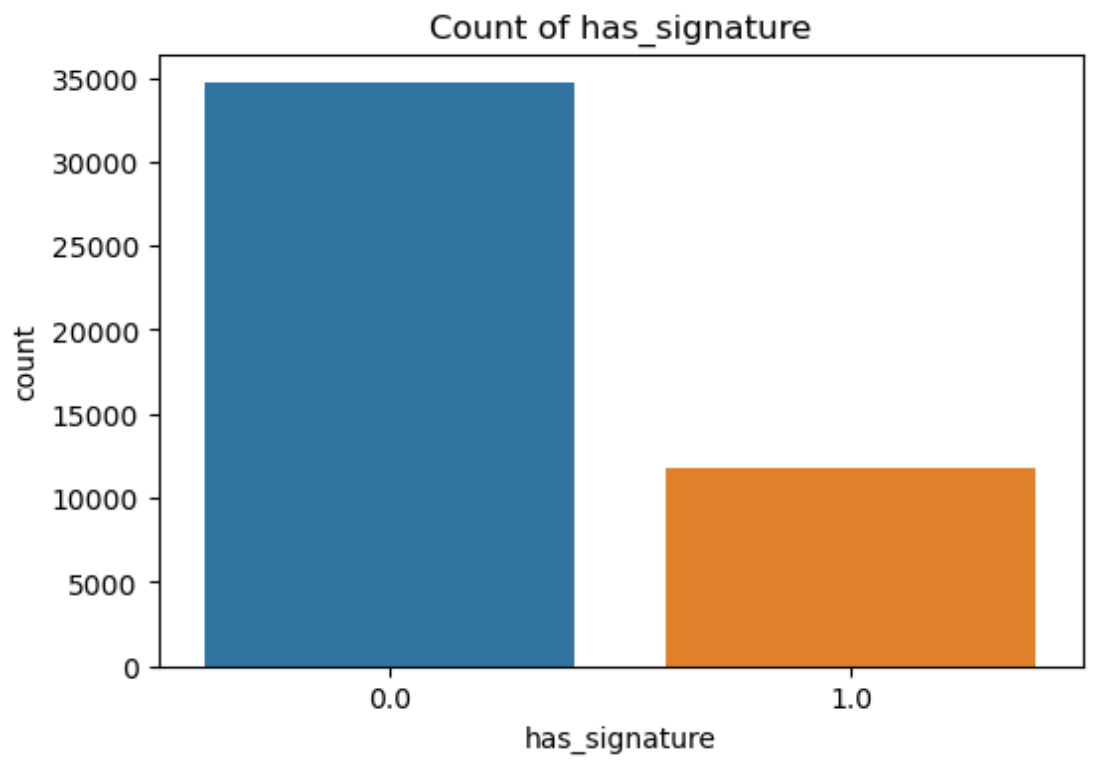
5.



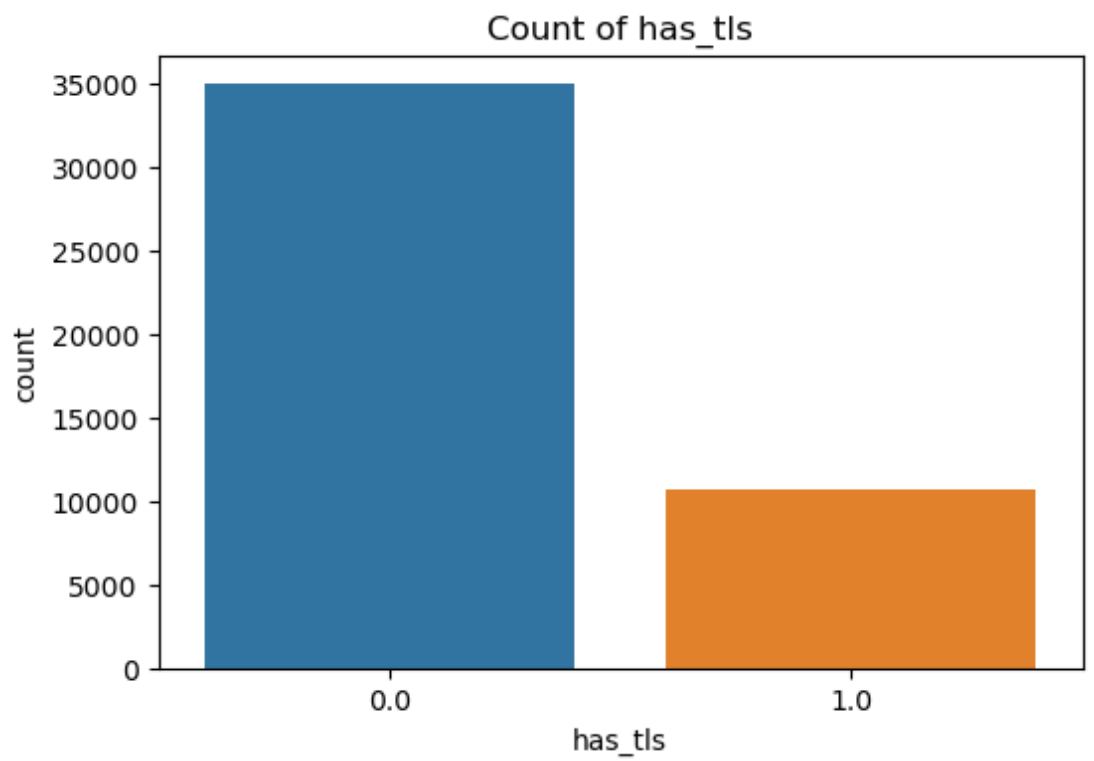
6.



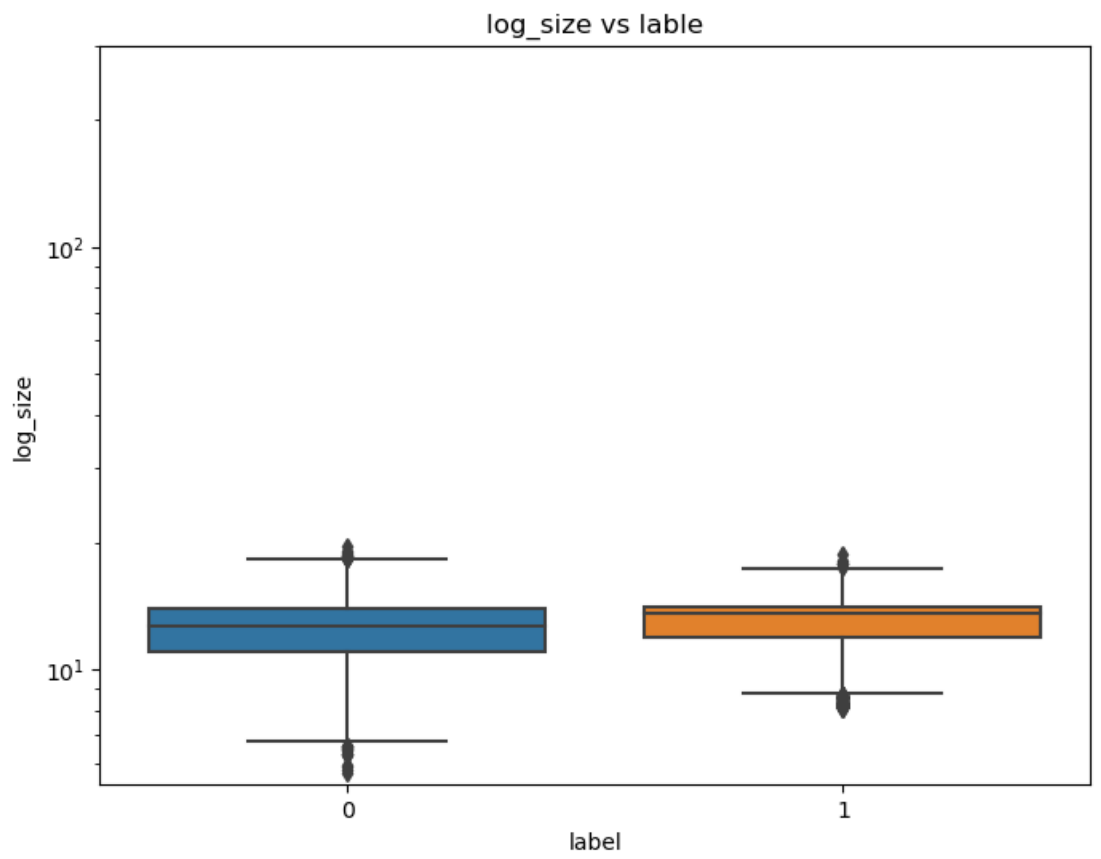
7.



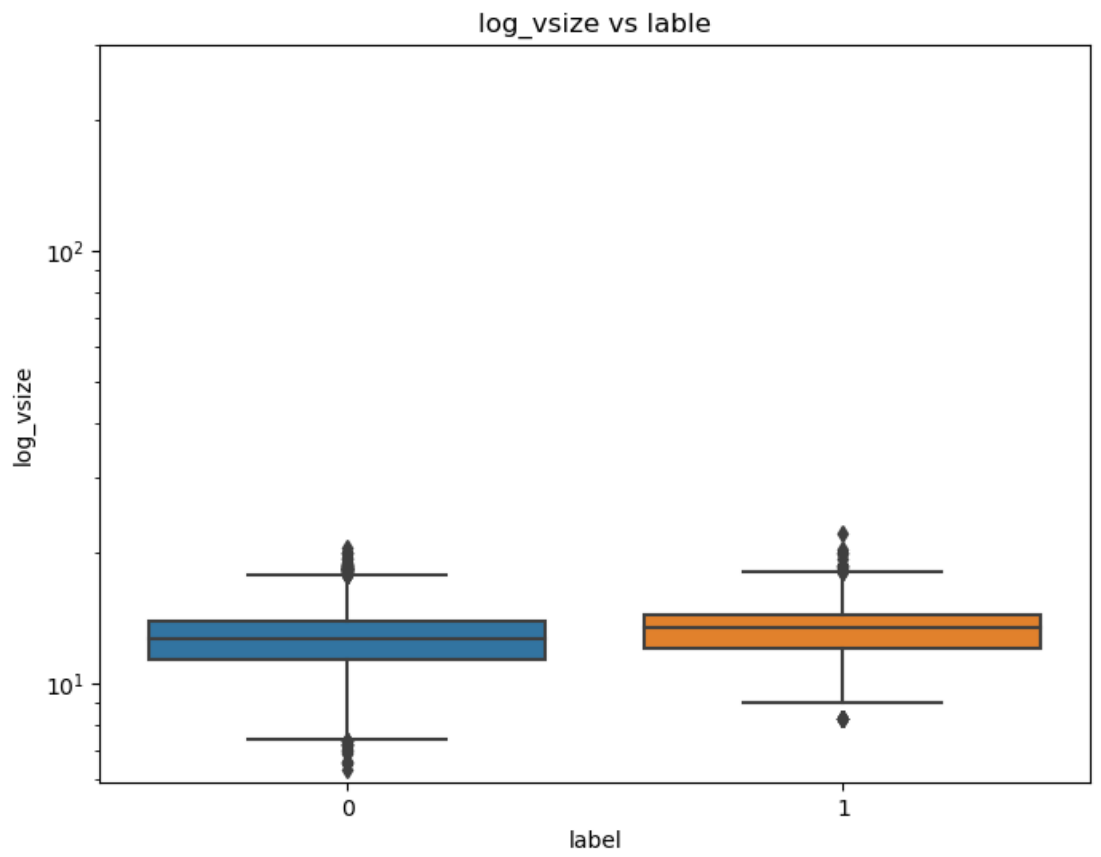
8.



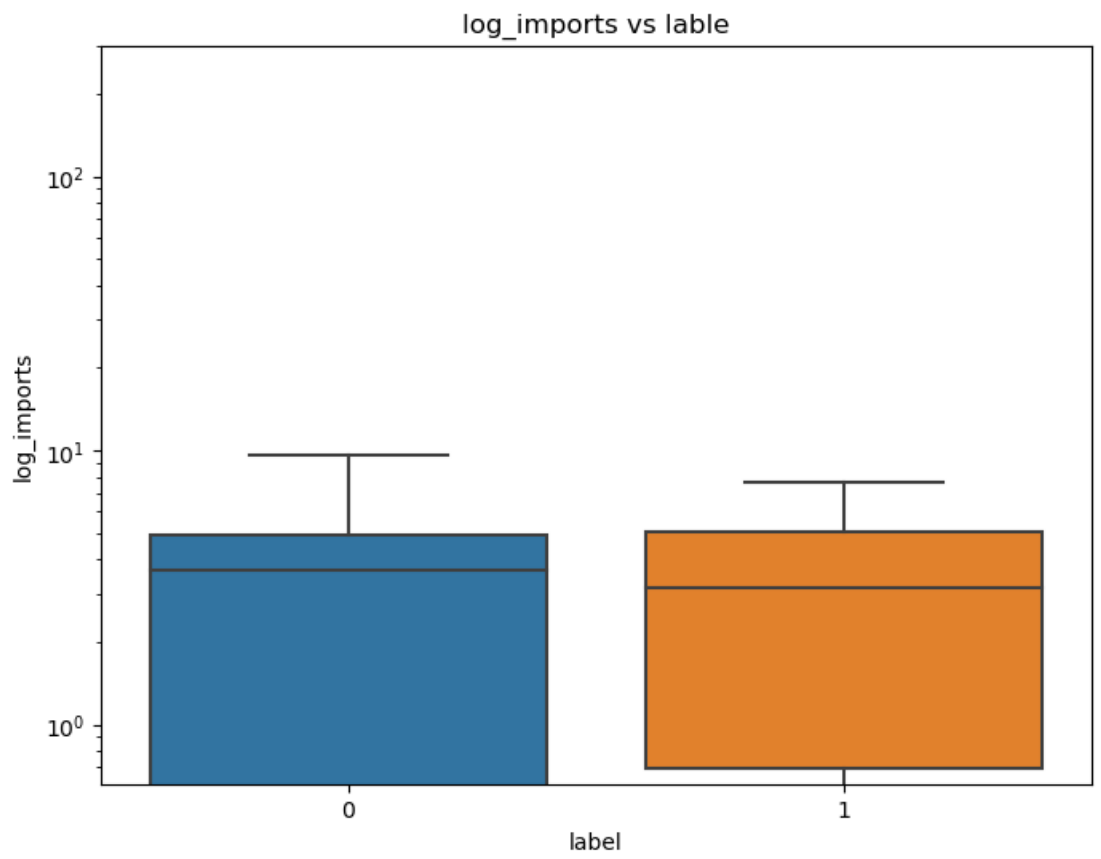
9.



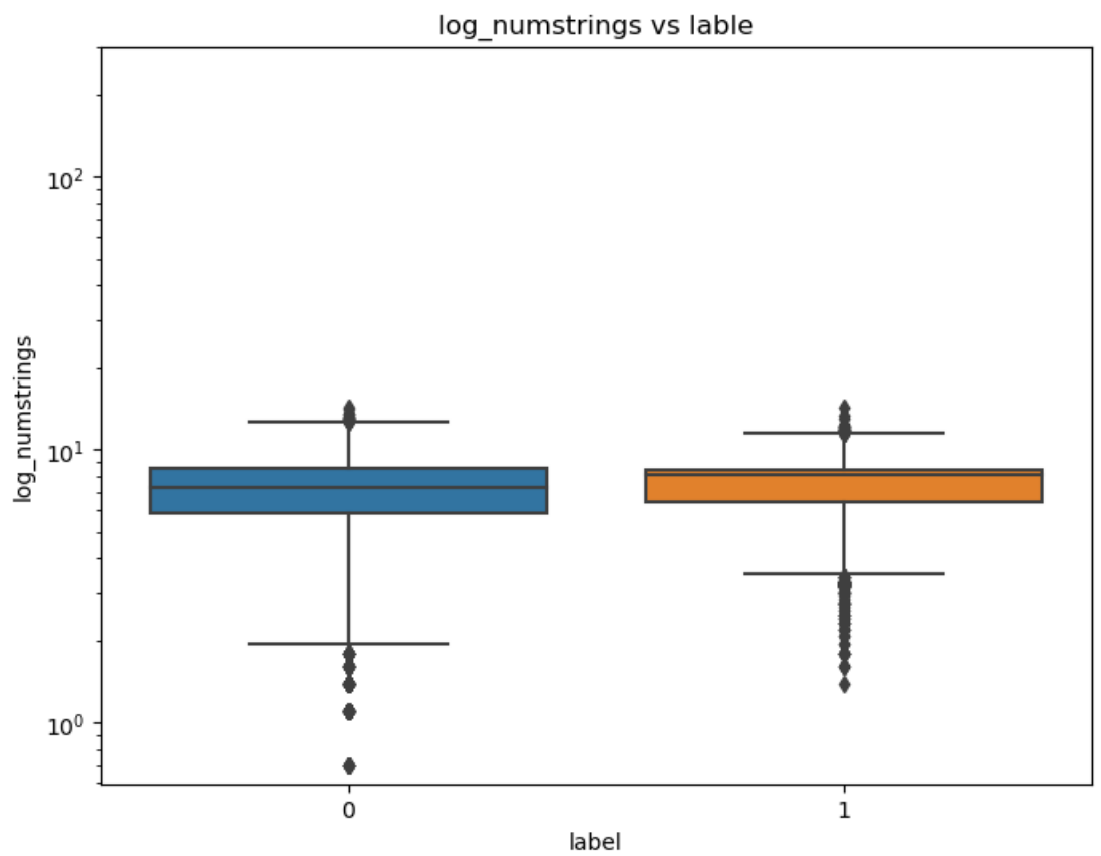
10.



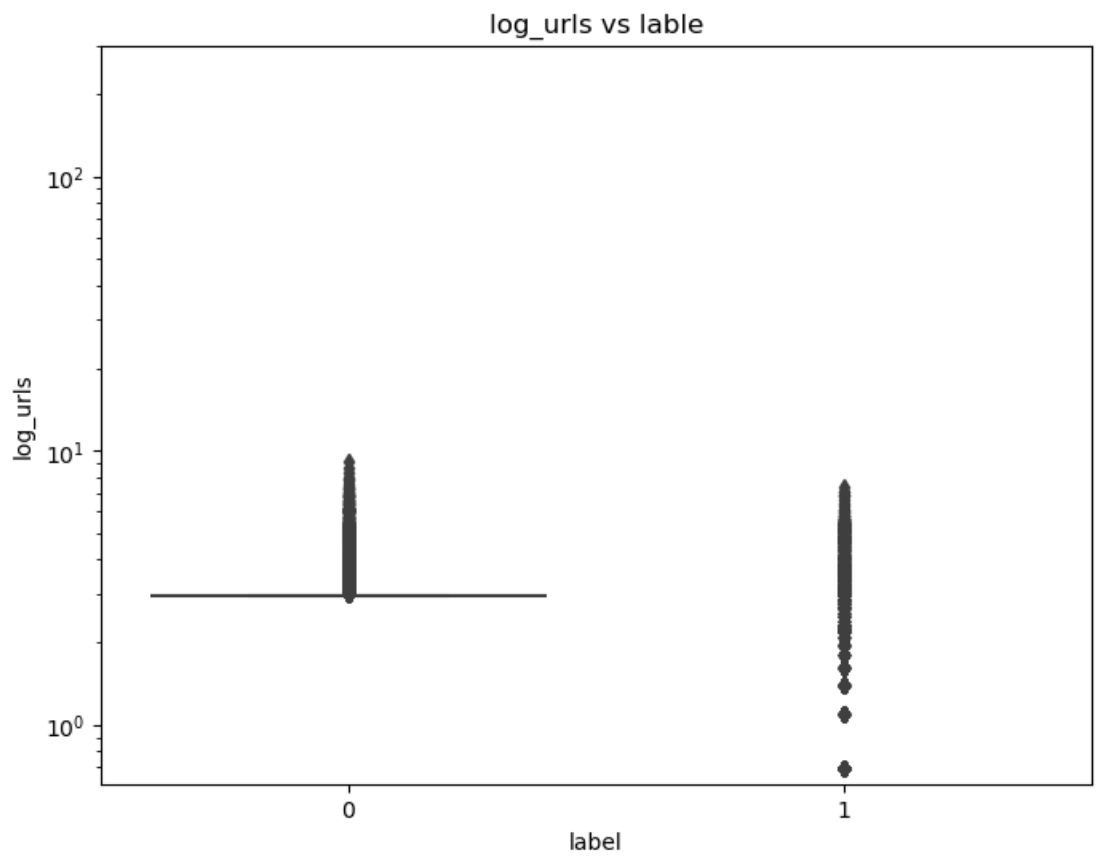
11.



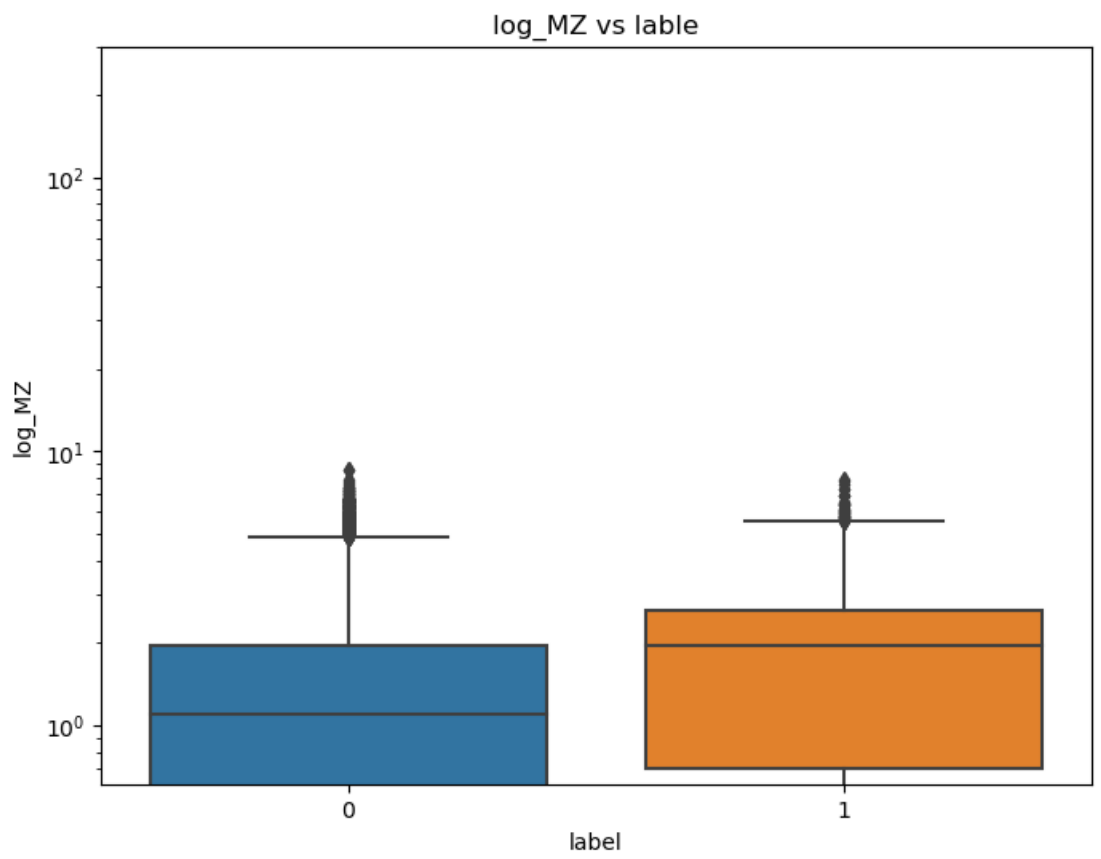
12.



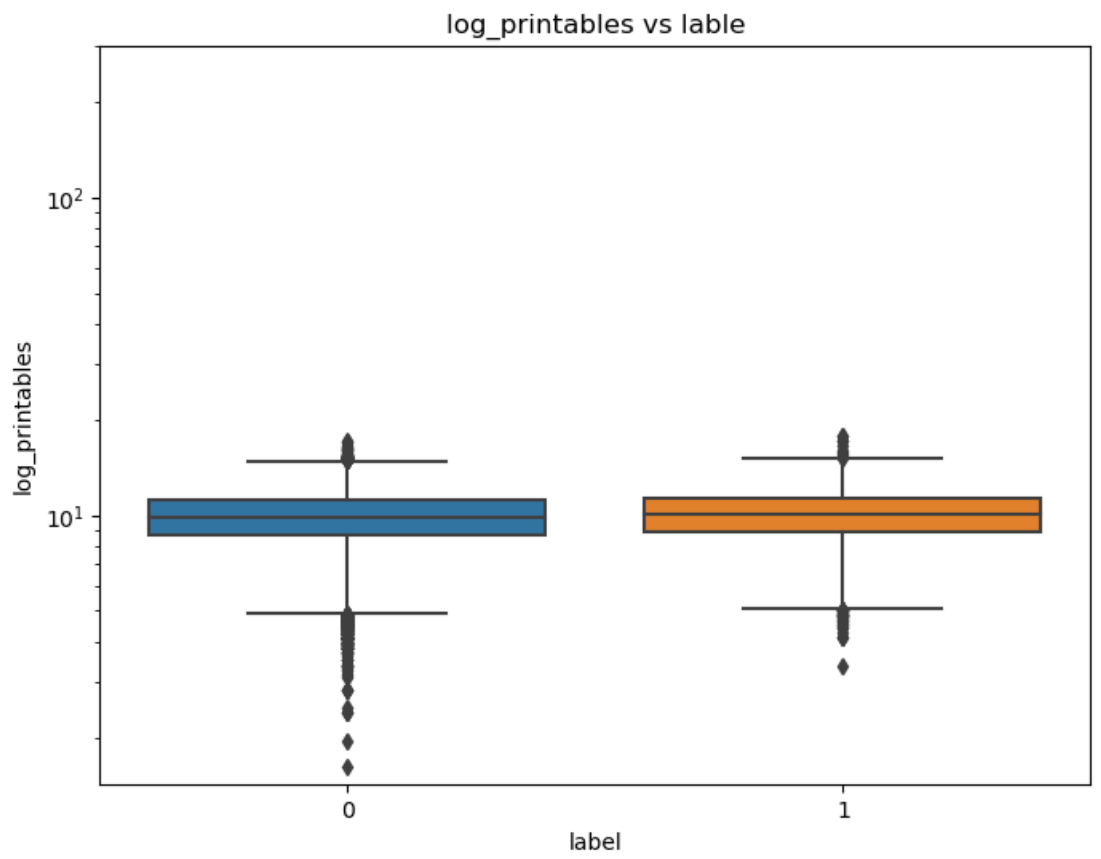
13.



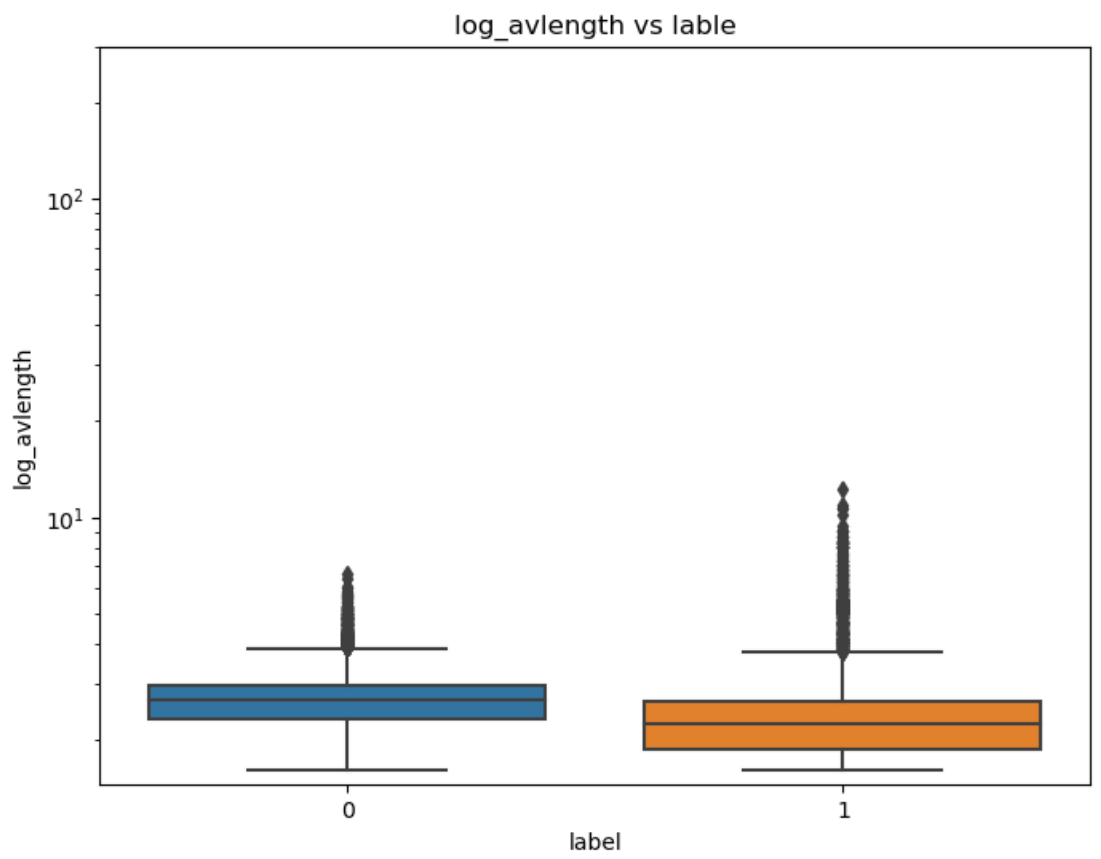
14.



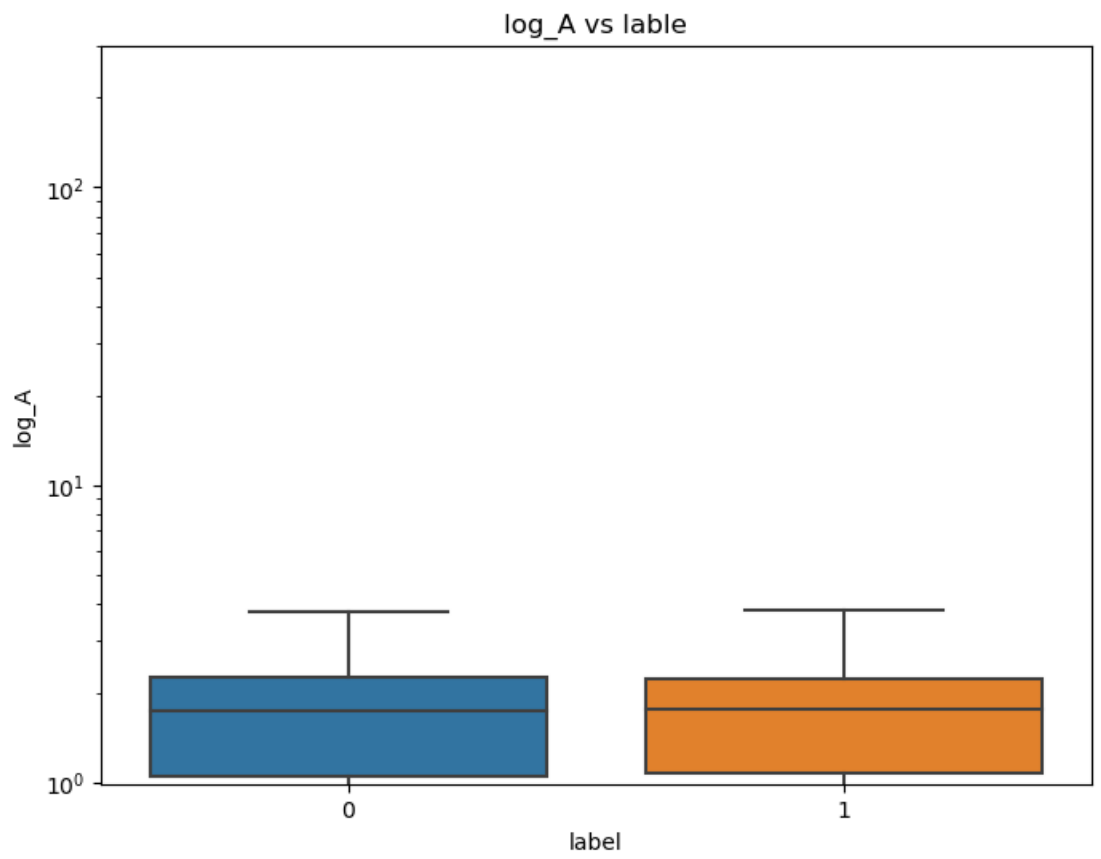
15.



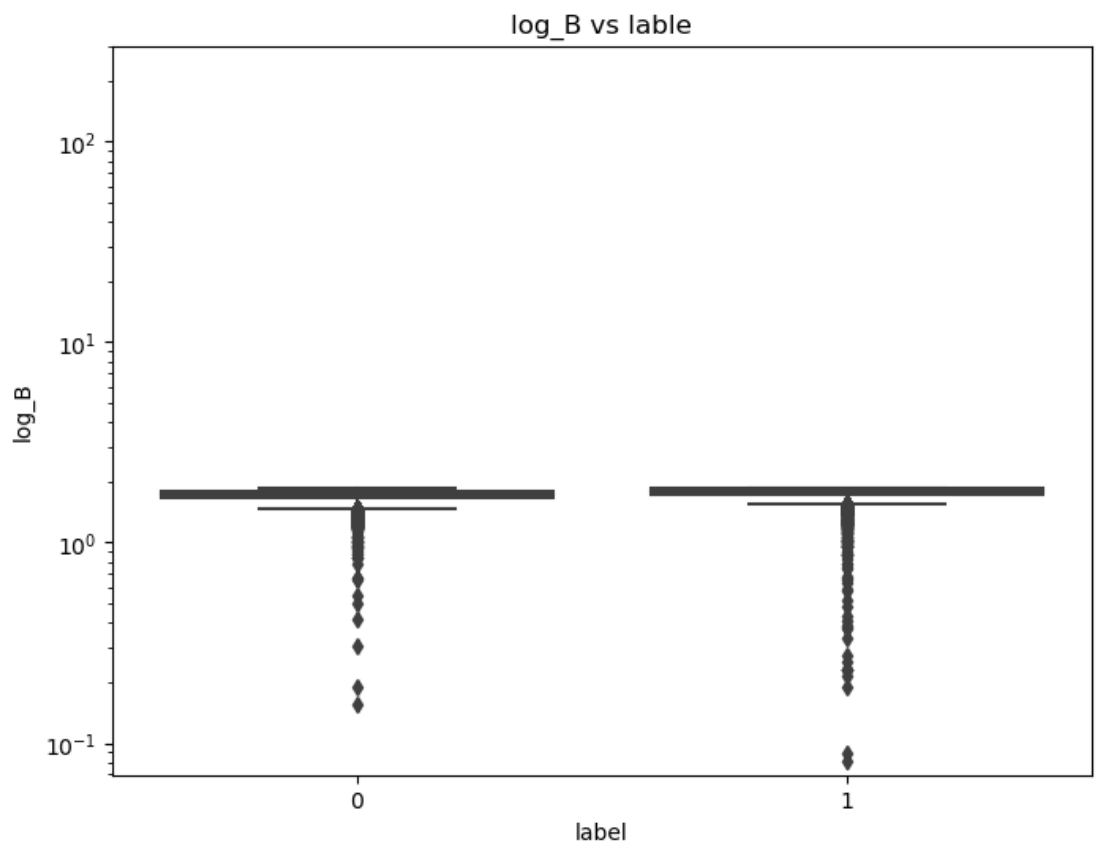
16.



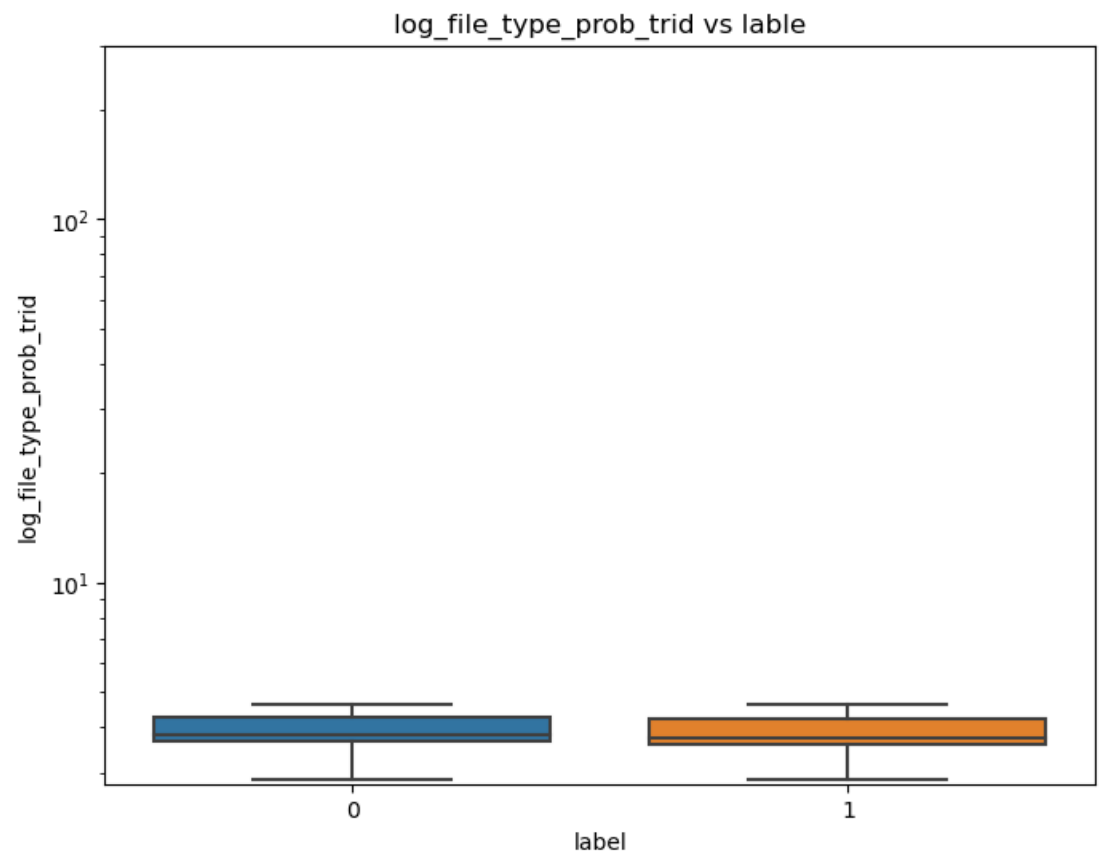
17.

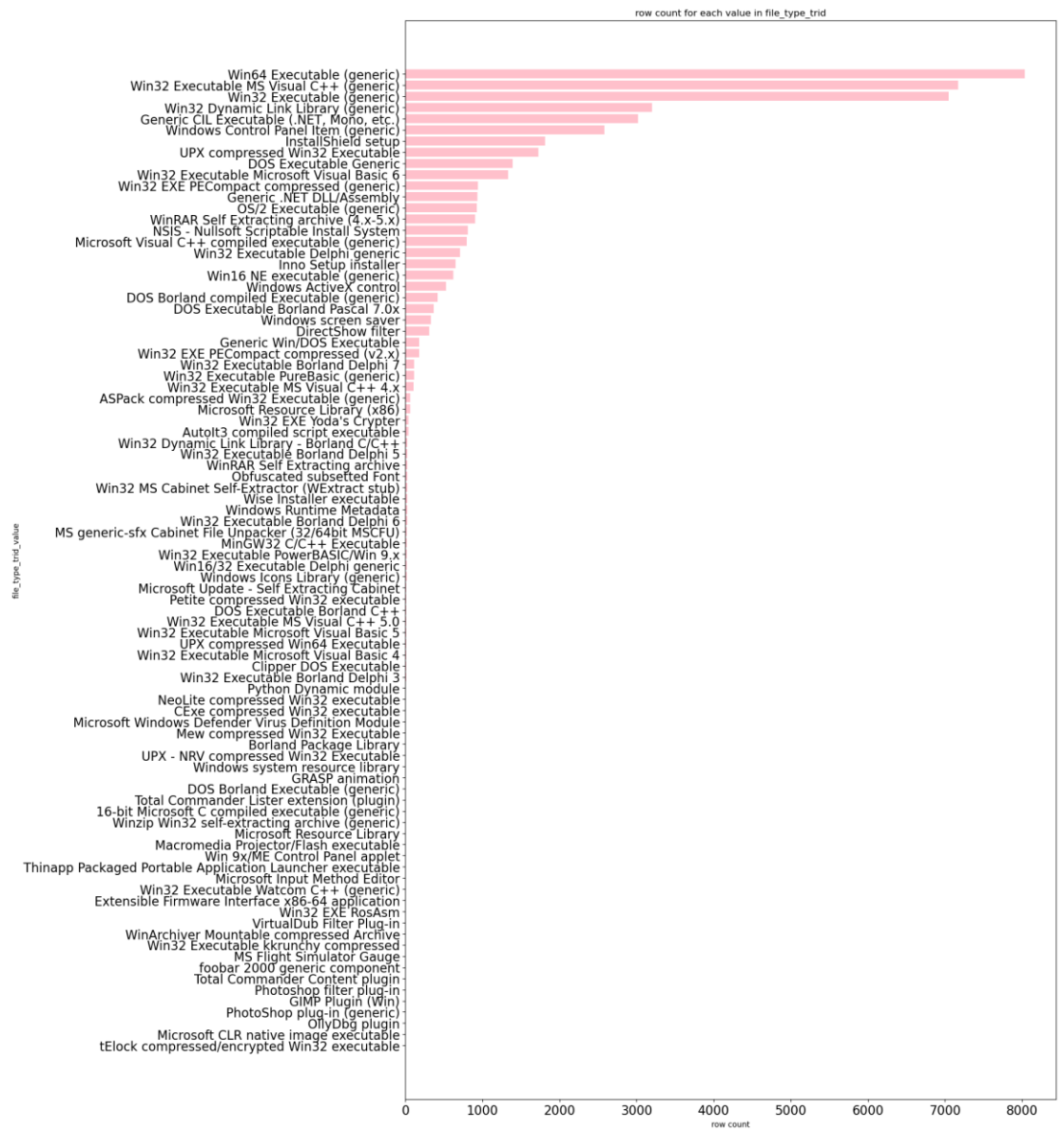


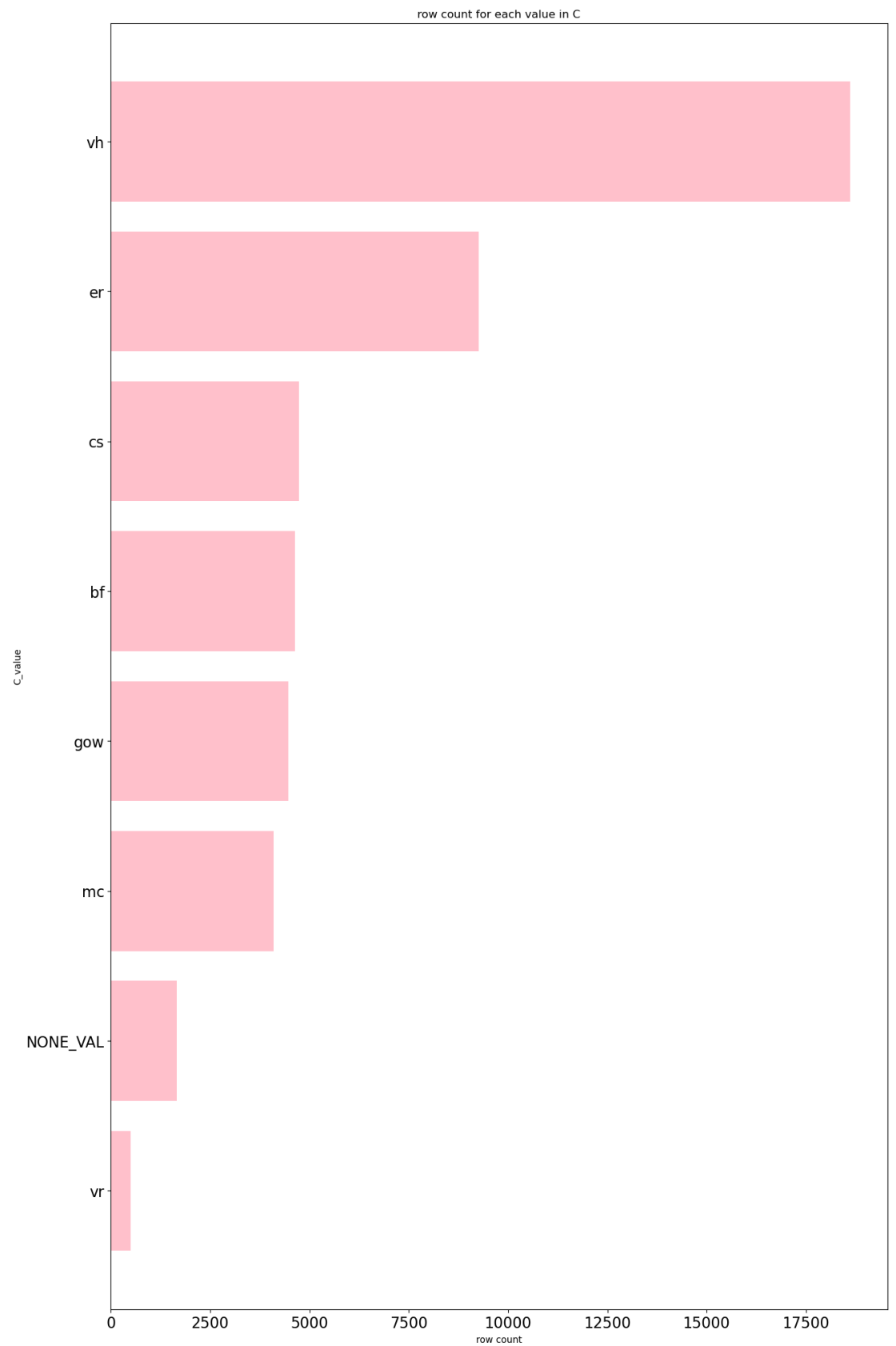
18.

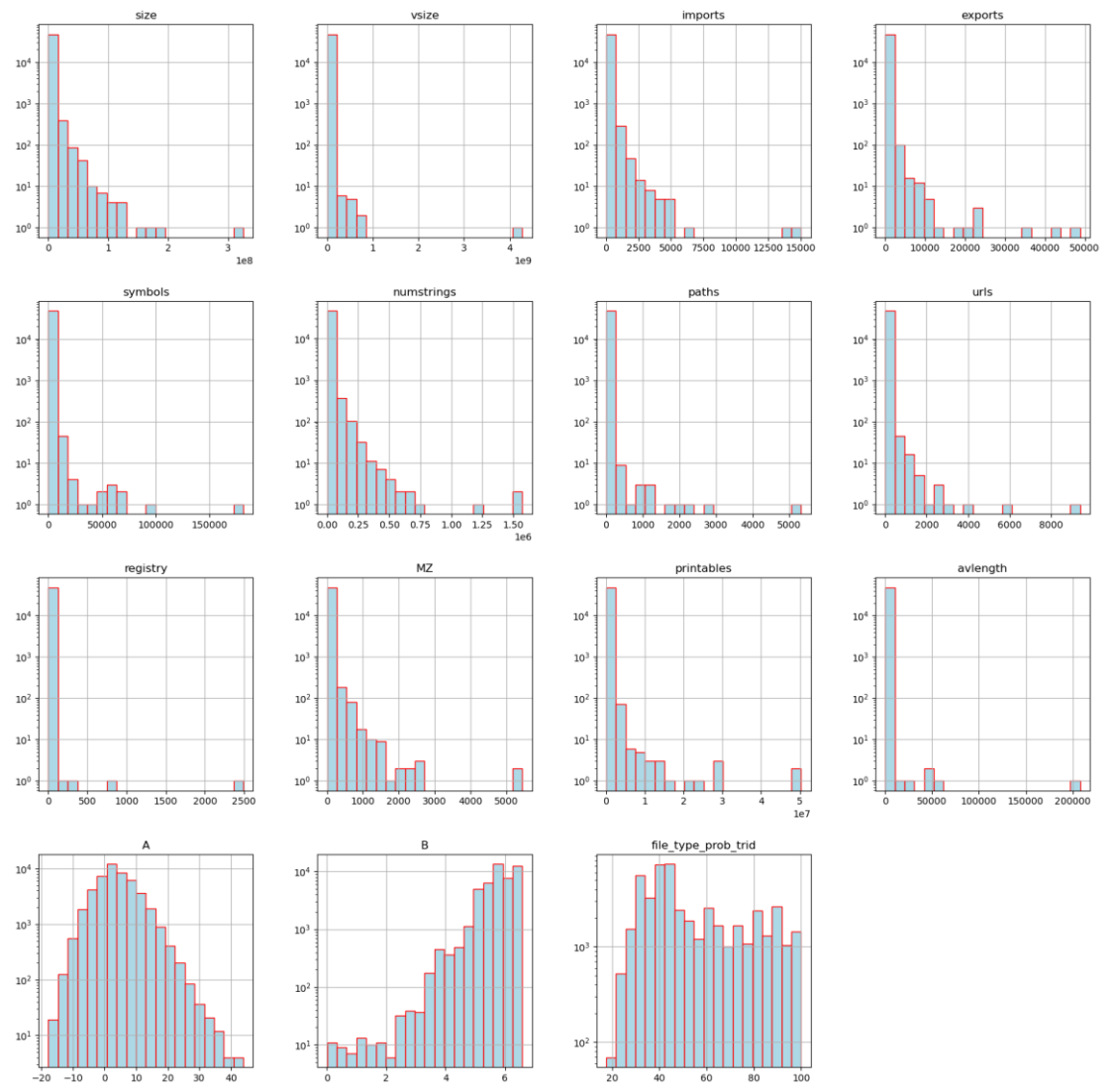


19.



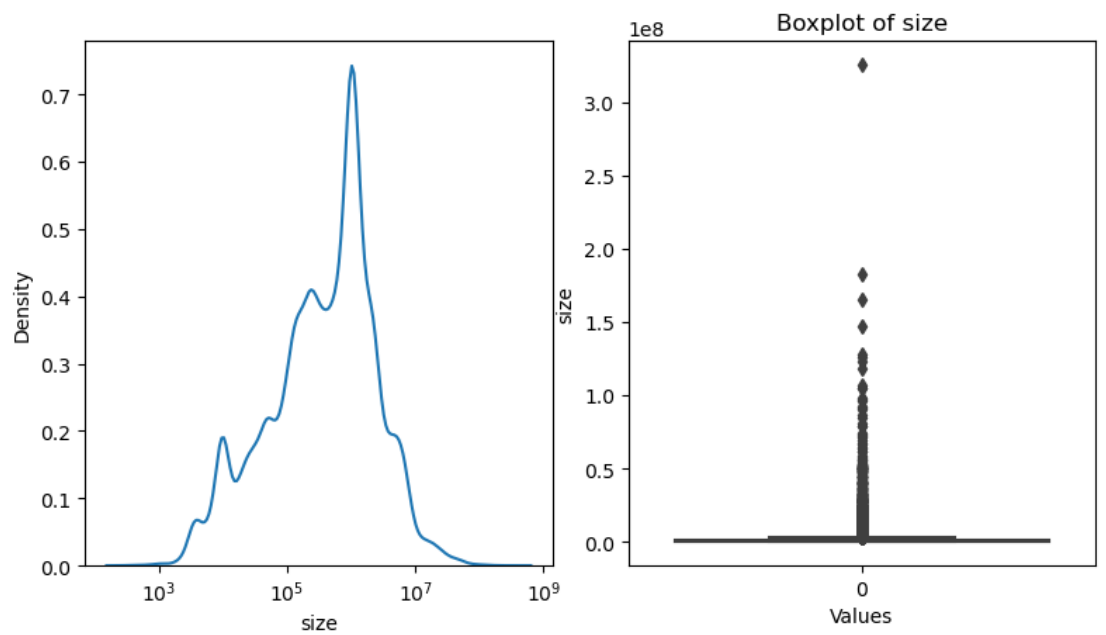






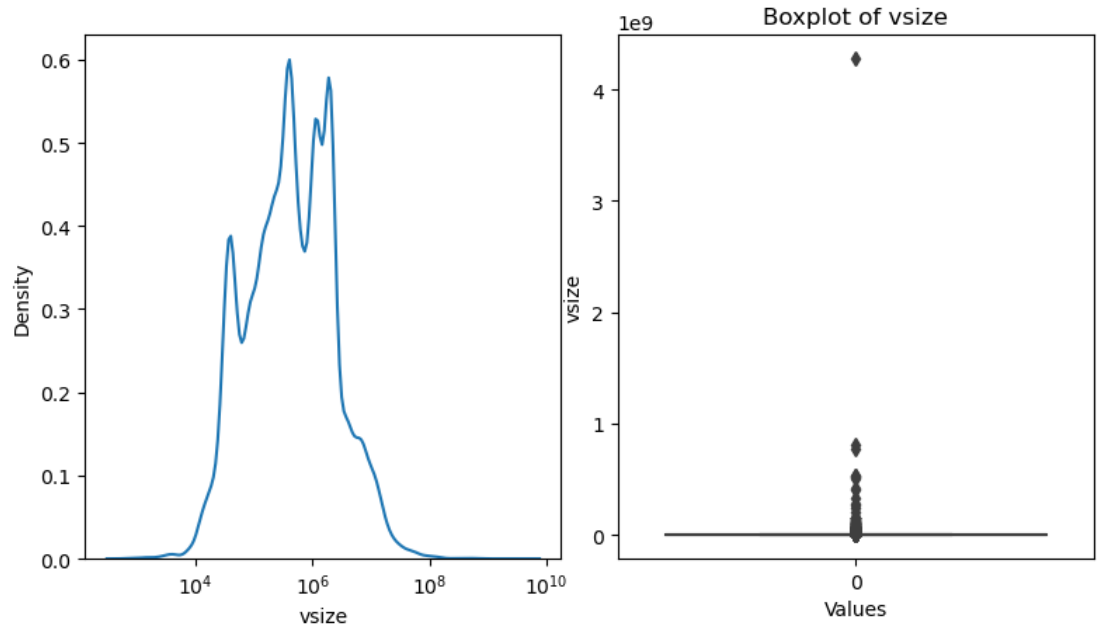
23)

Distribution before Capping



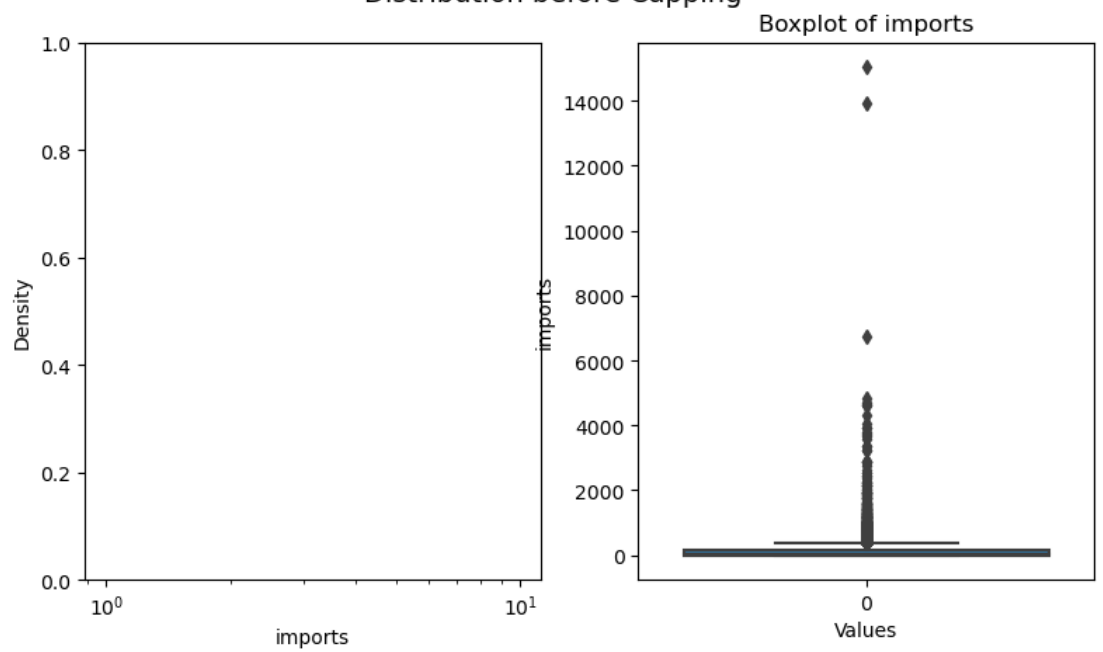
24)

Distribution before Capping



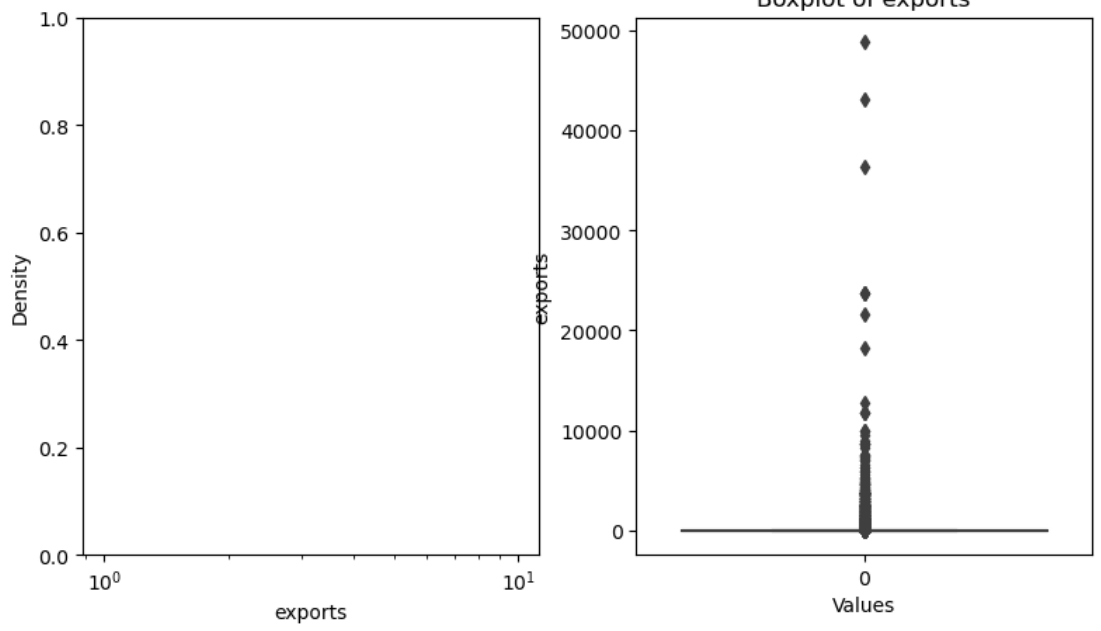
25)

Distribution before Capping



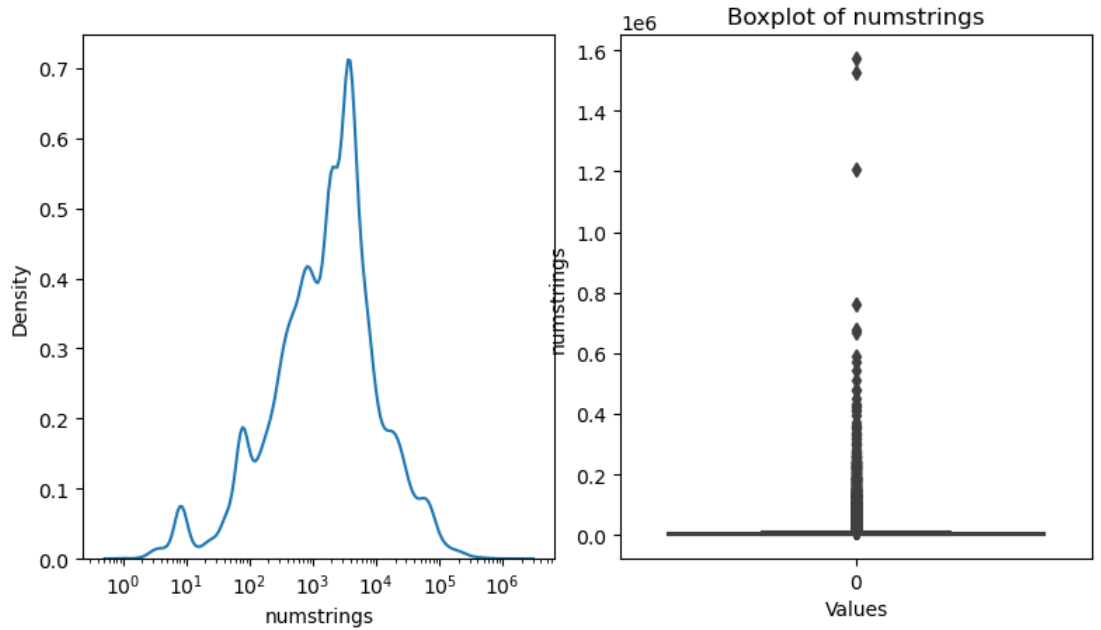
26)

Distribution before Capping



27)

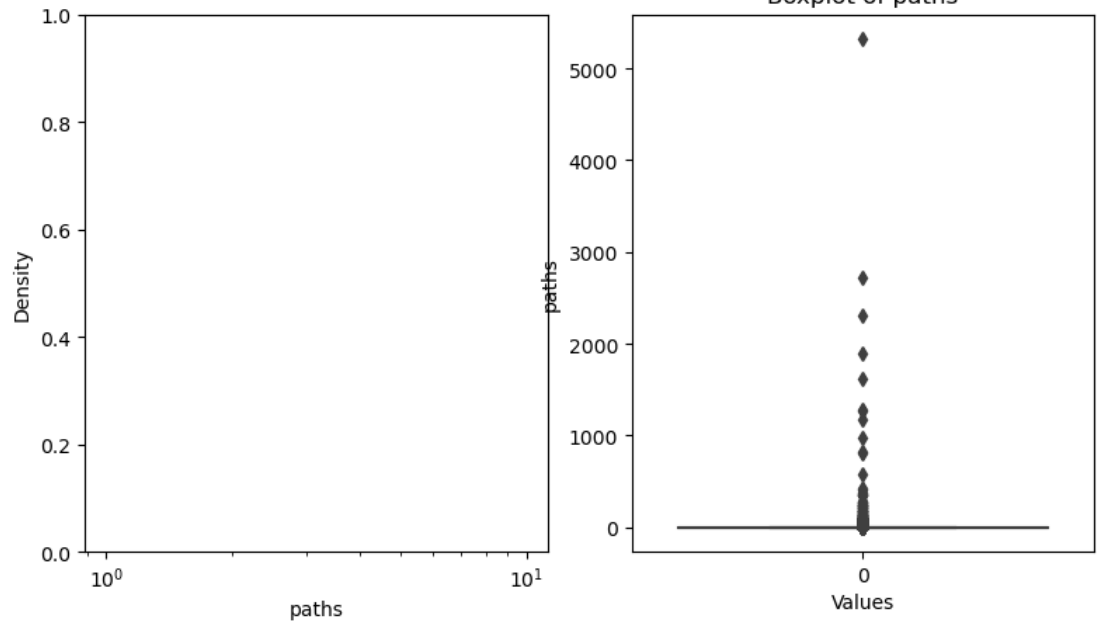
Distribution before Capping



28)

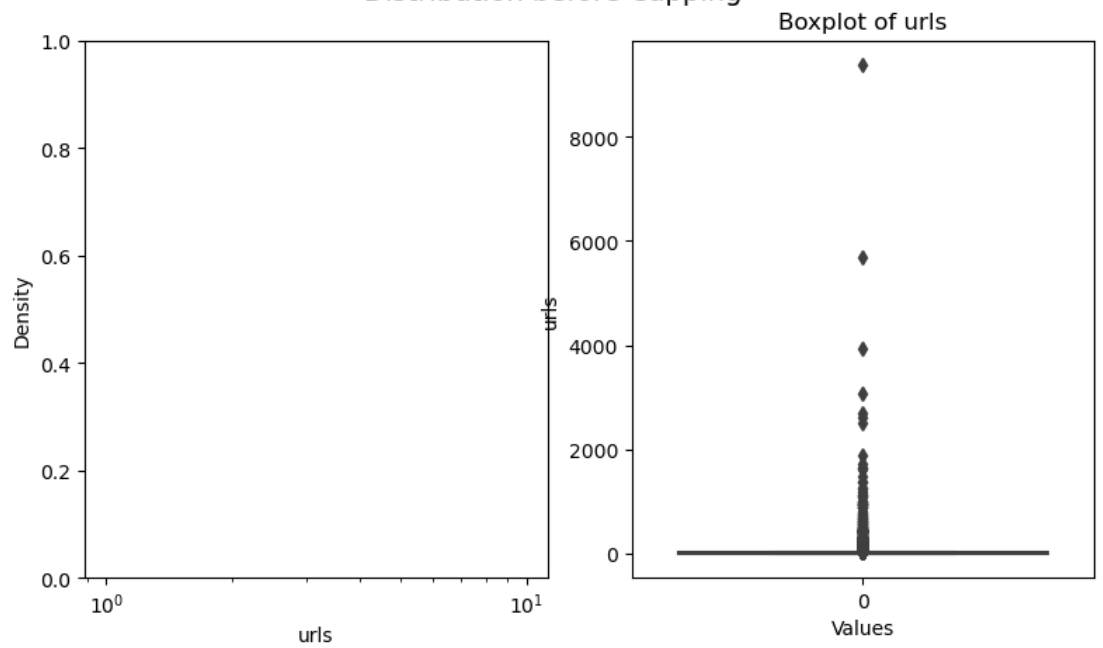
29)

Distribution before Capping



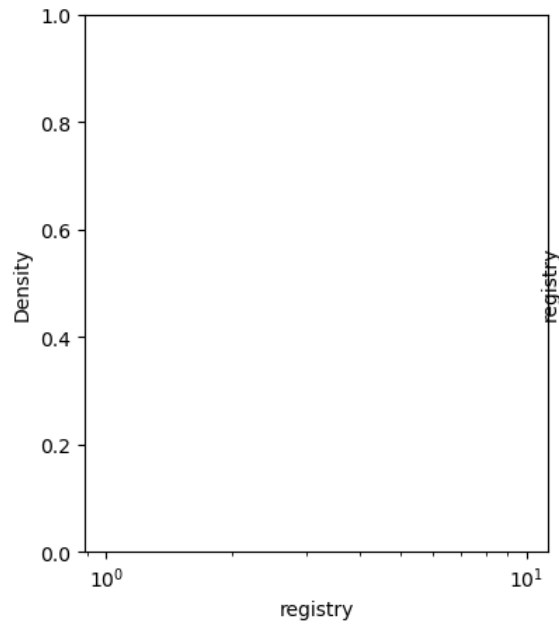
30)

Distribution before Capping

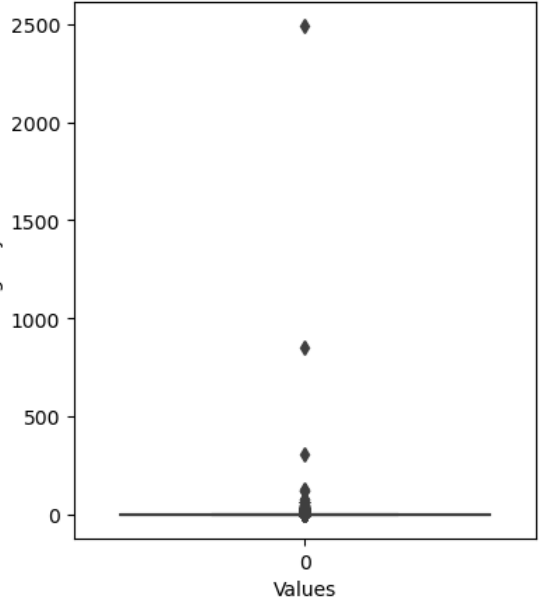


31)

Distribution before Capping

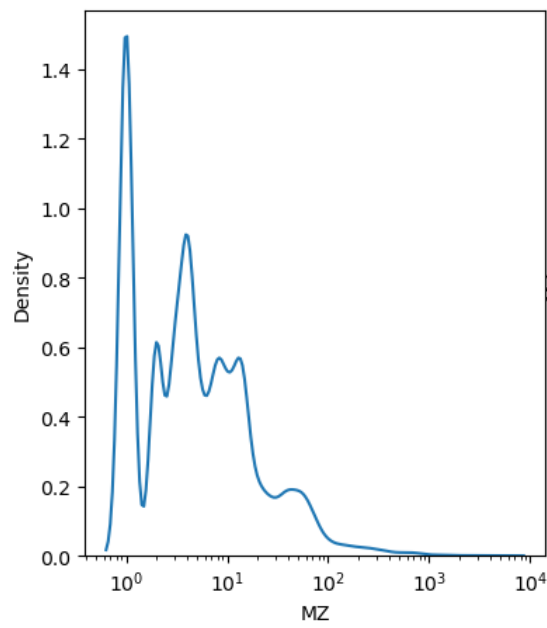


Boxplot of registry

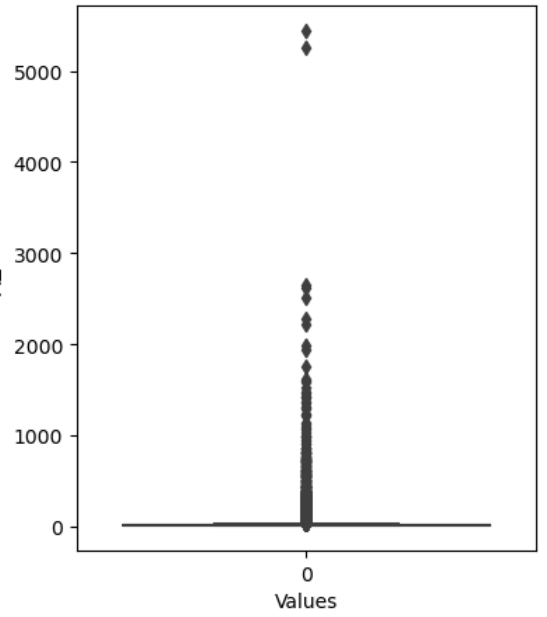


32)

Distribution before Capping

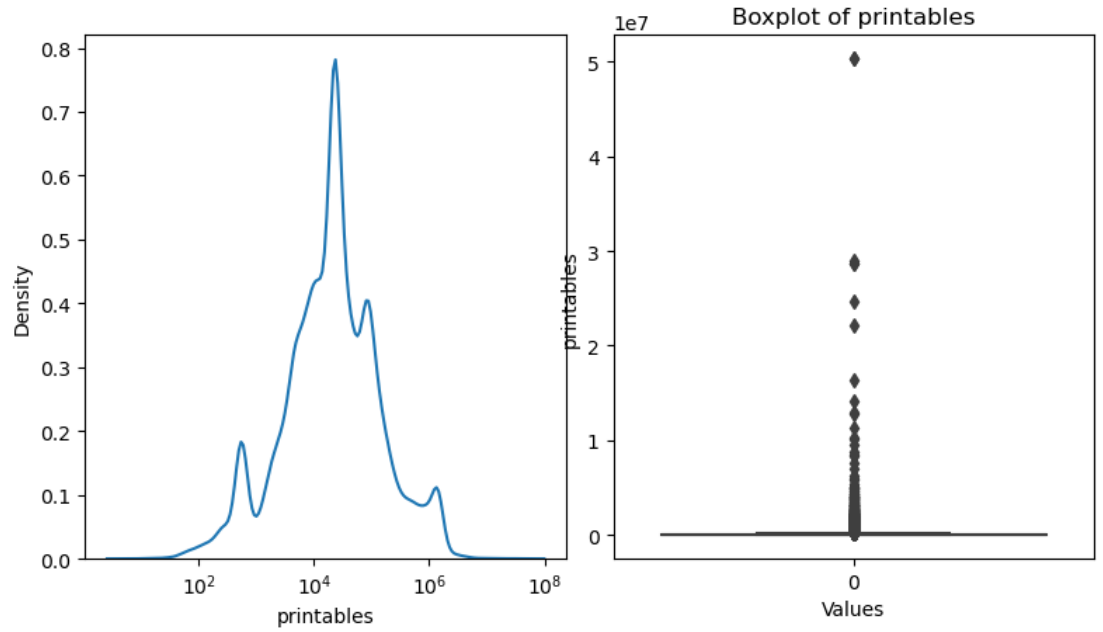


Boxplot of MZ



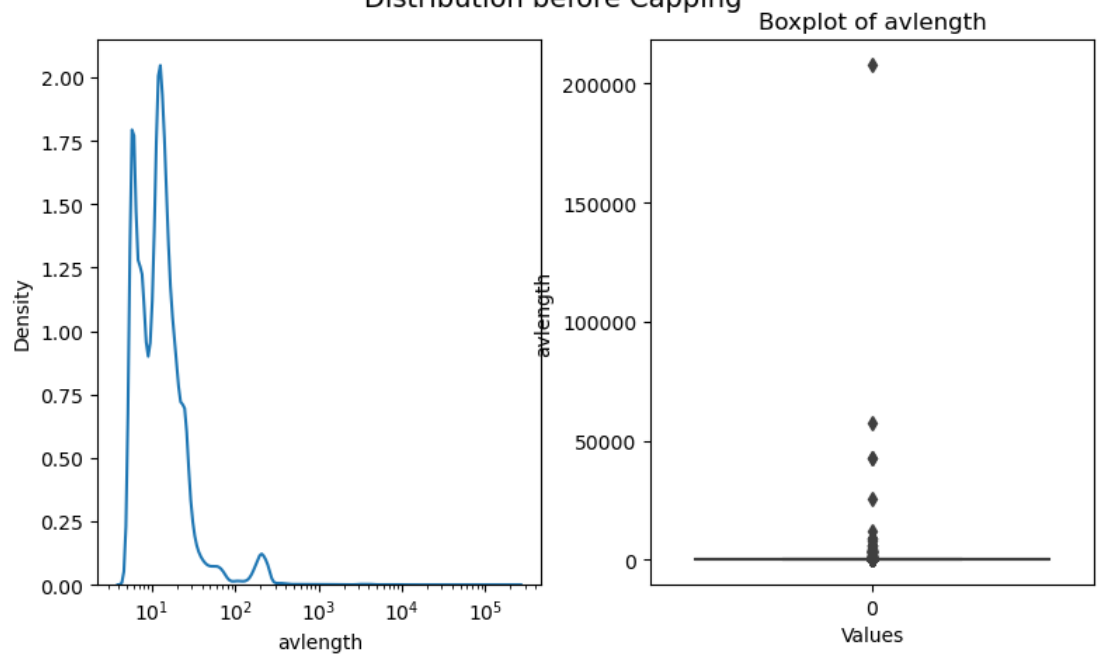
33)

Distribution before Capping

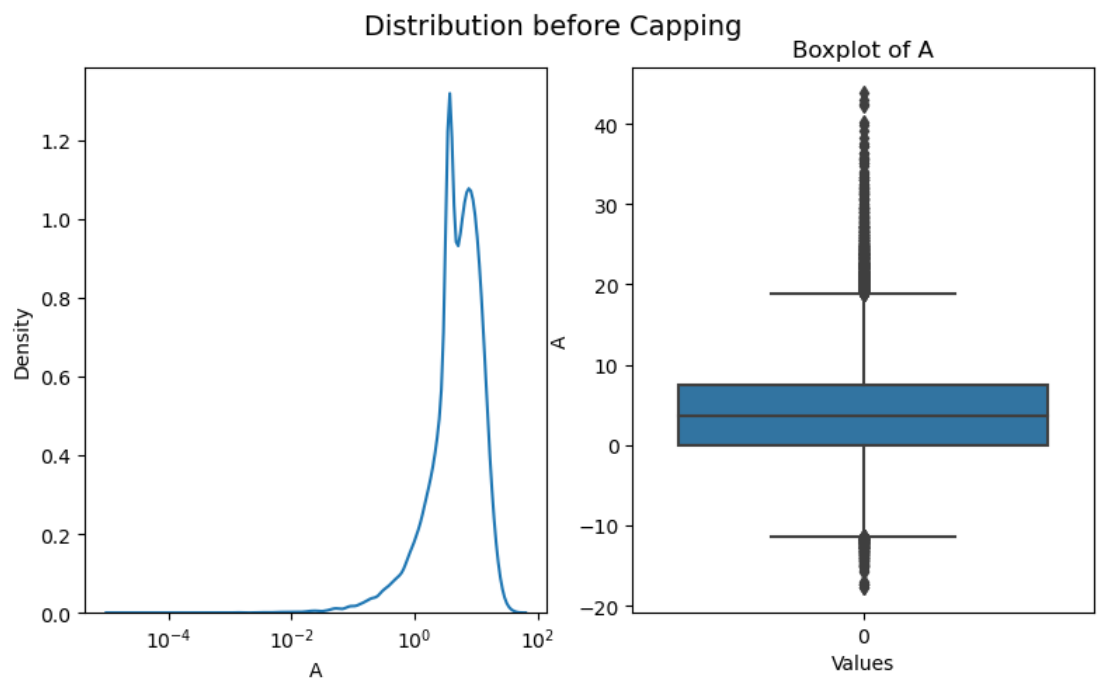


34)

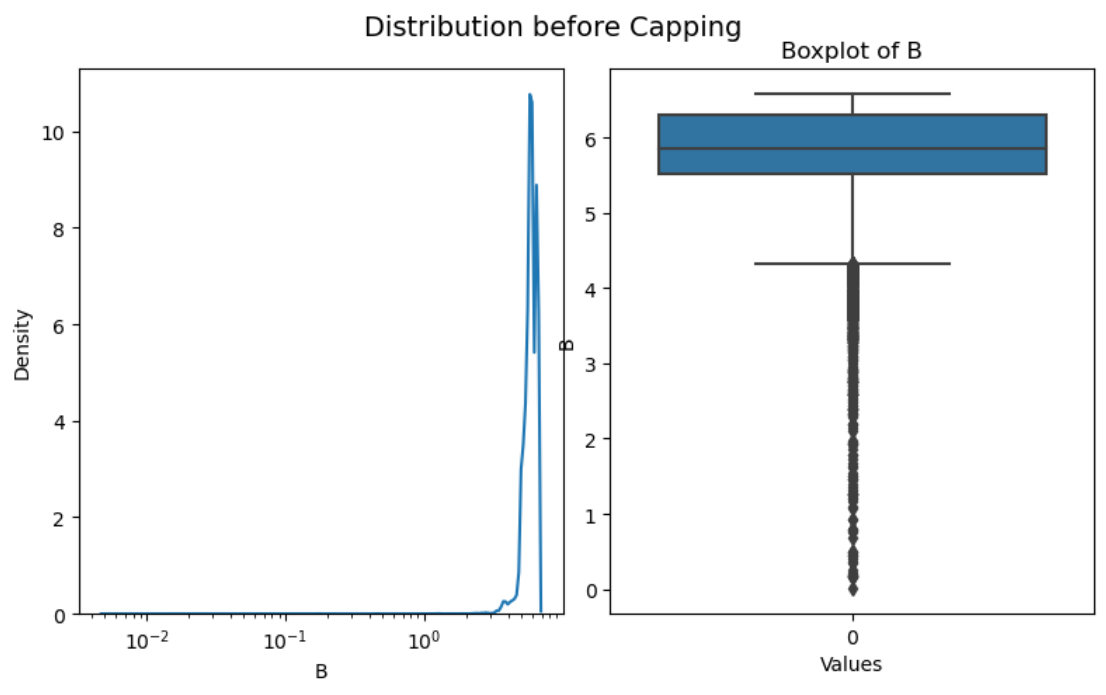
Distribution before Capping



35)

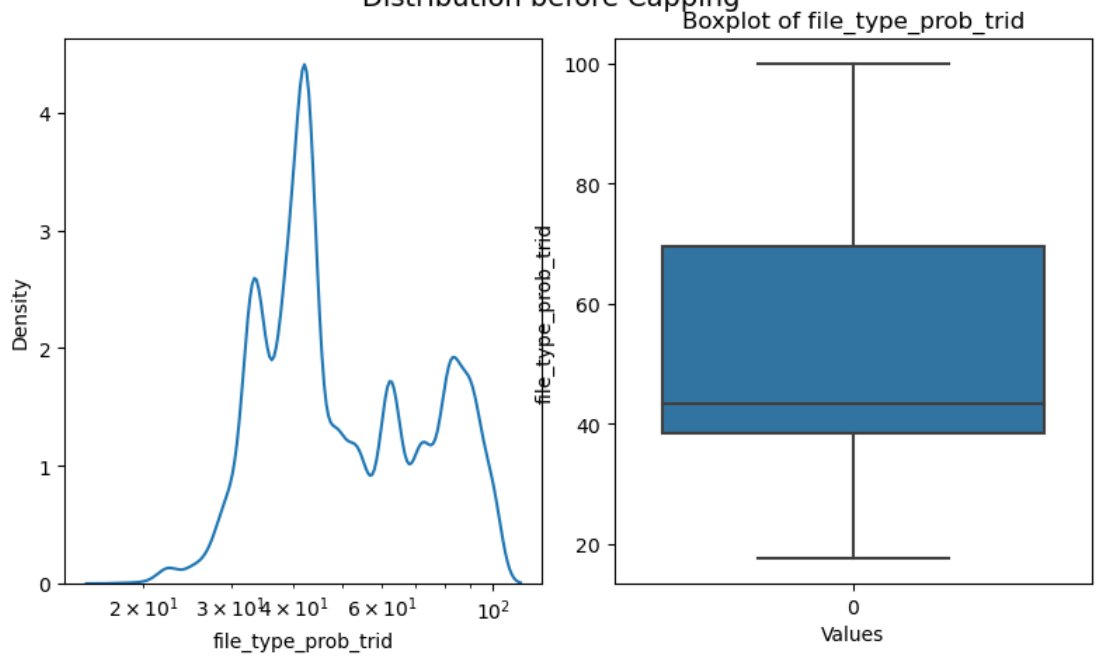


36)

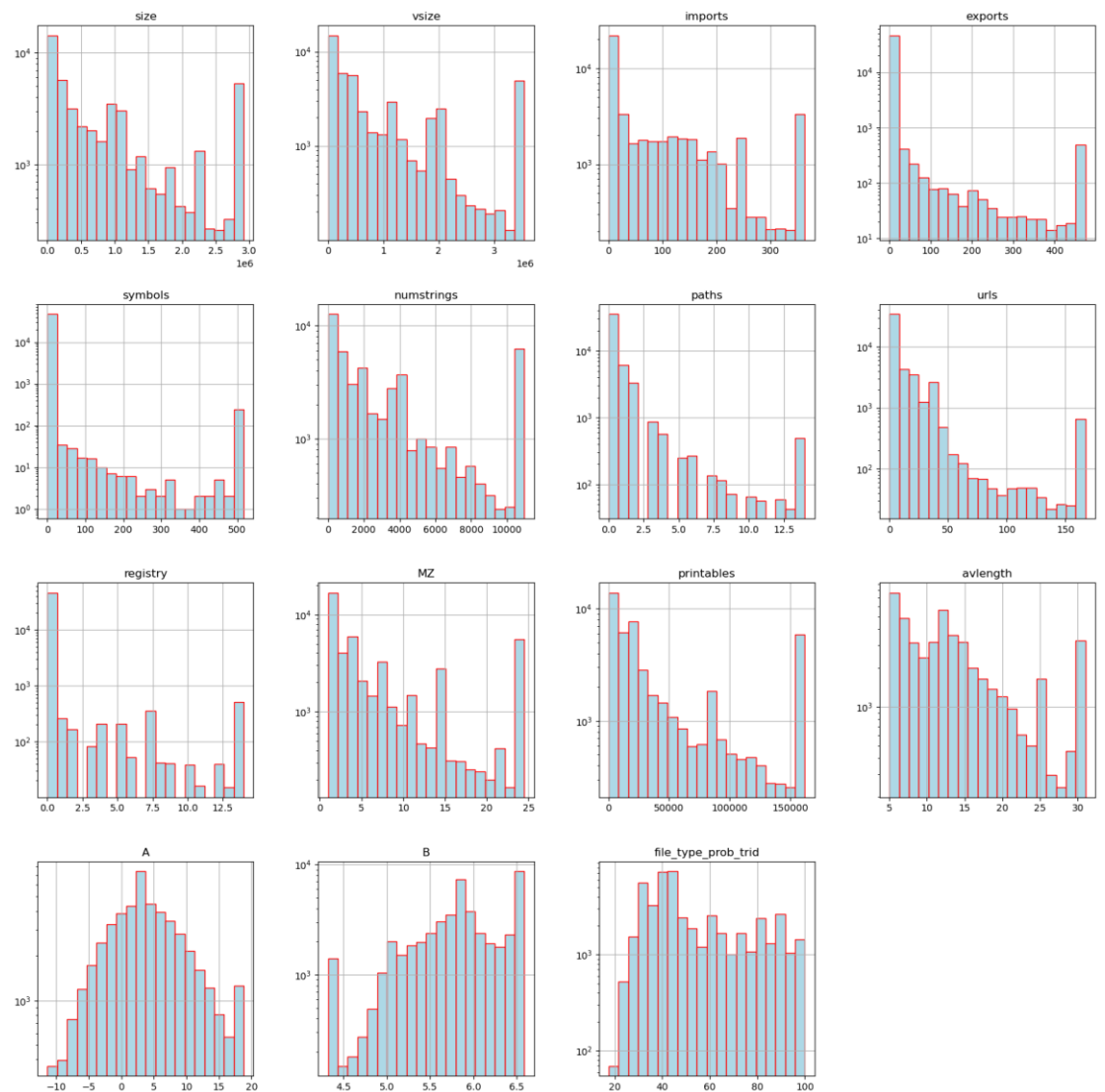


37)

Distribution before Capping

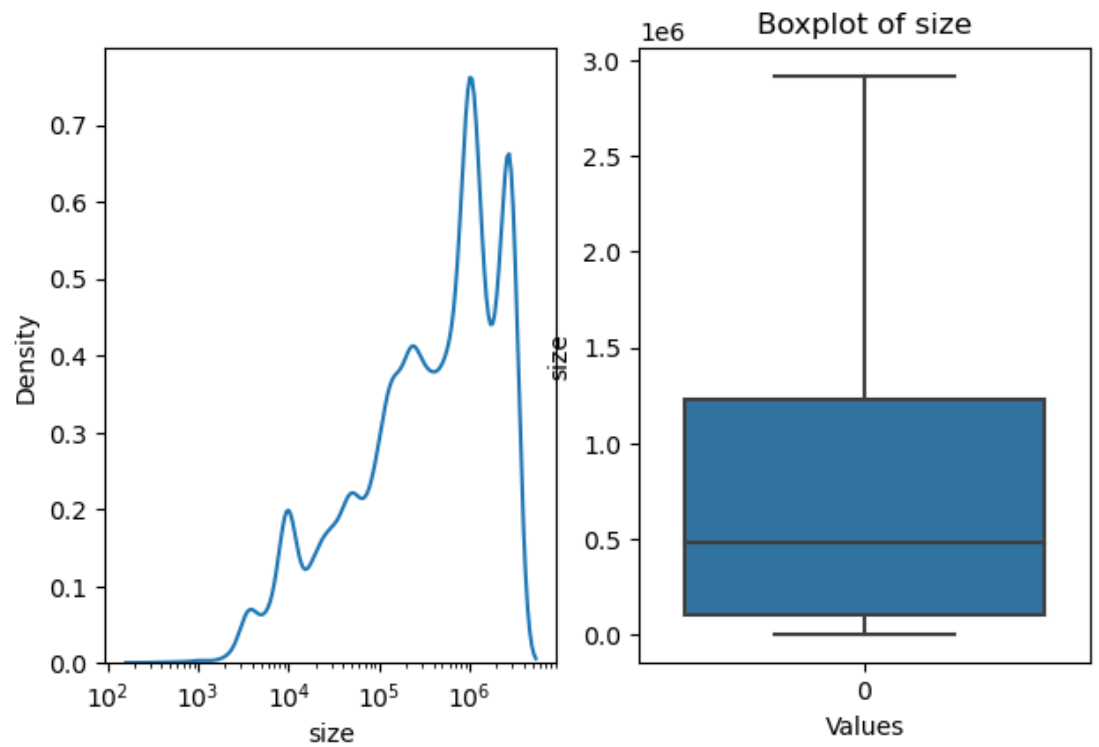


38)



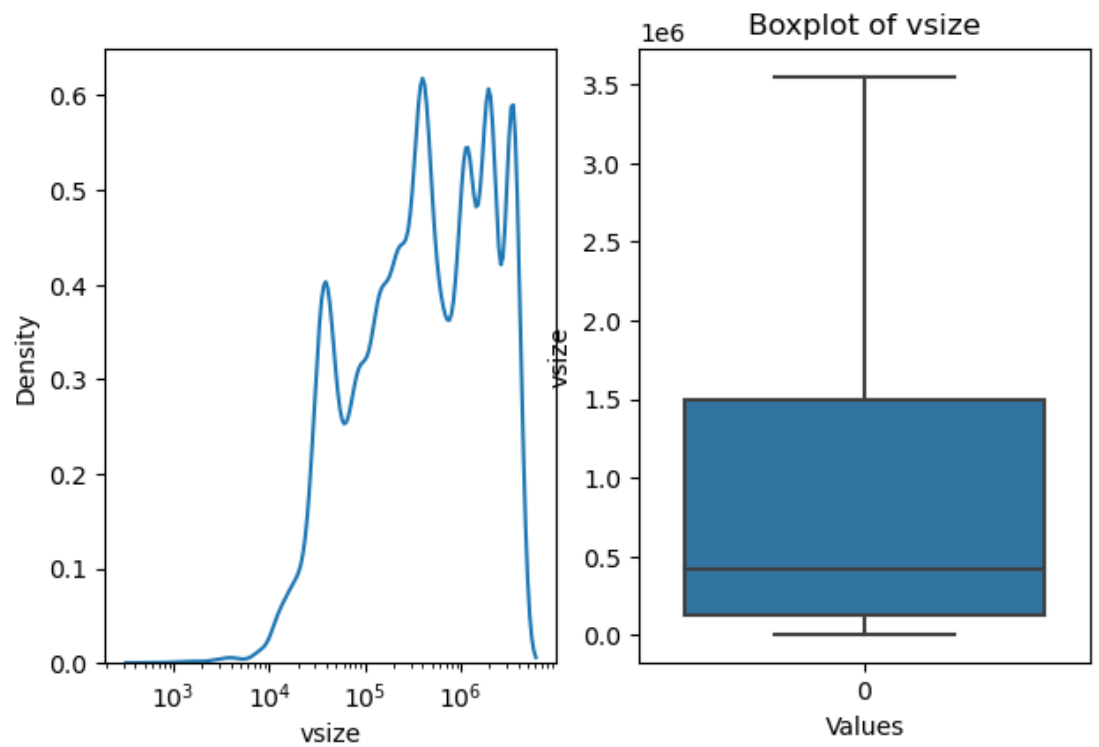
39)

Distribution after Capping



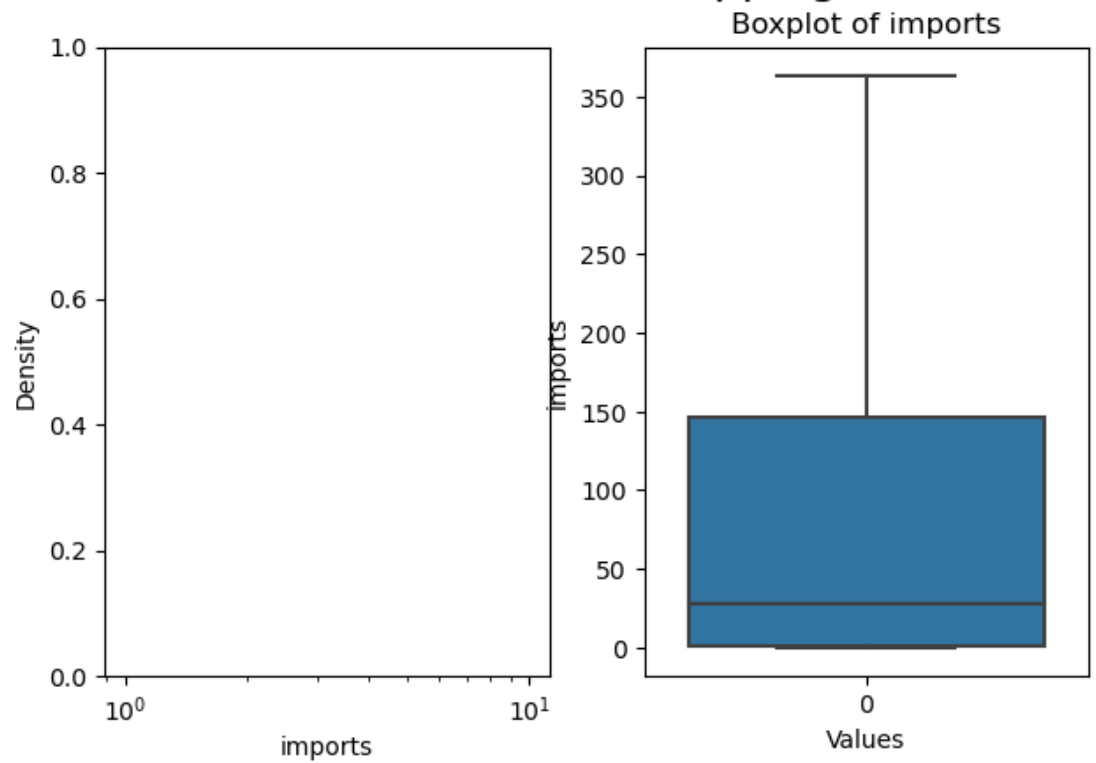
40)

Distribution after Capping



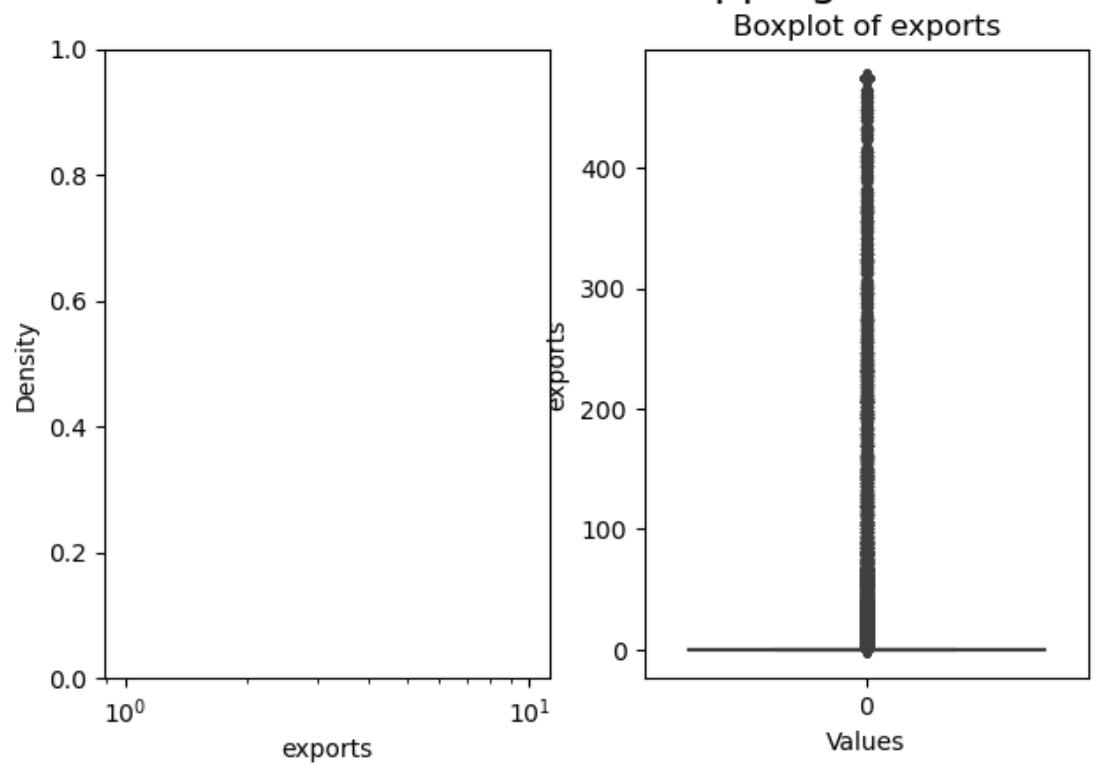
41)

Distribution after Capping



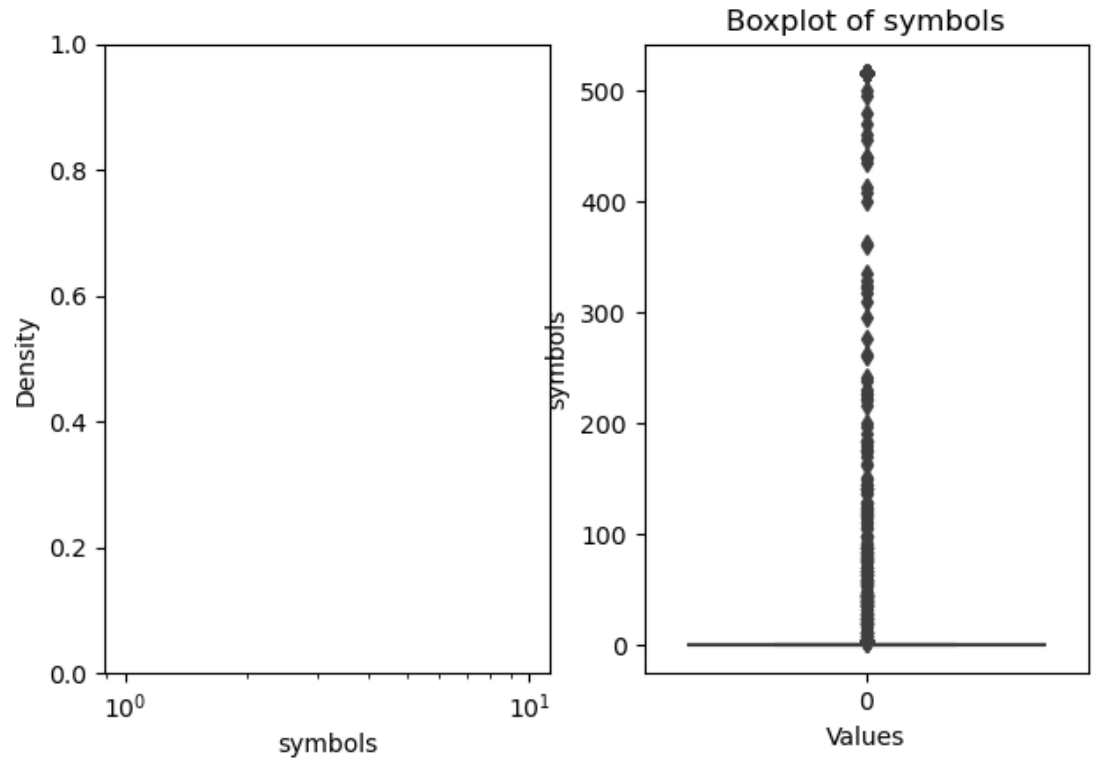
42)

Distribution after Capping



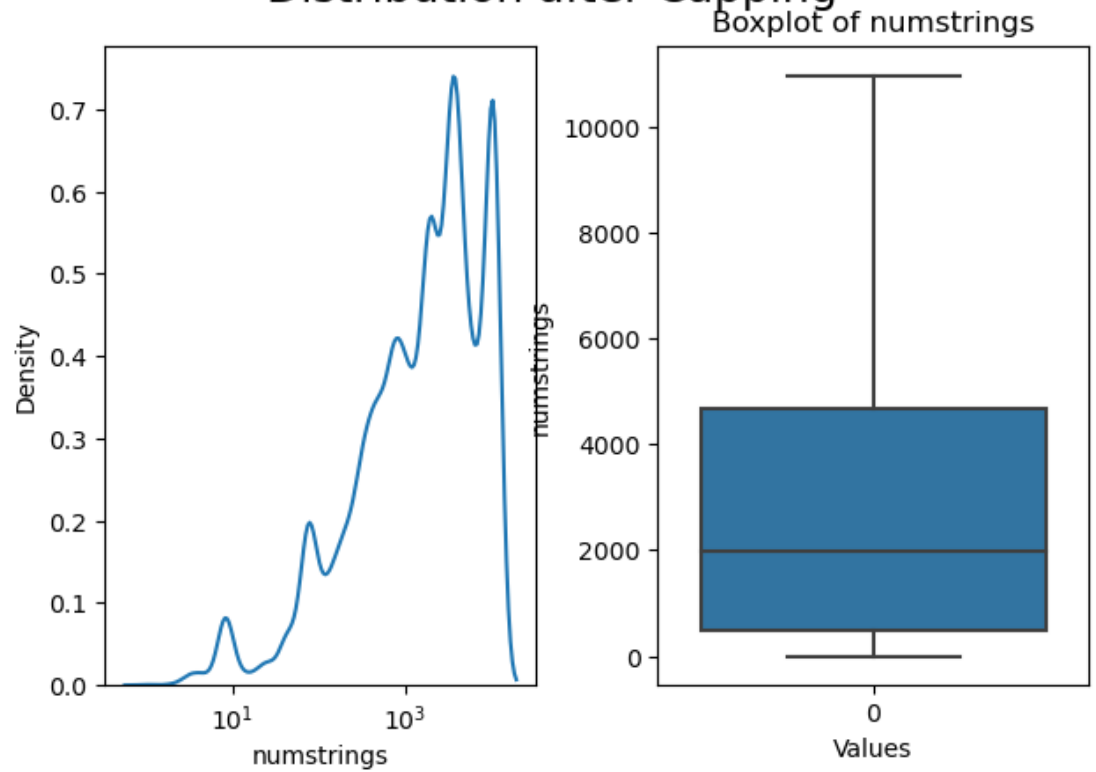
43)

Distribution after Capping



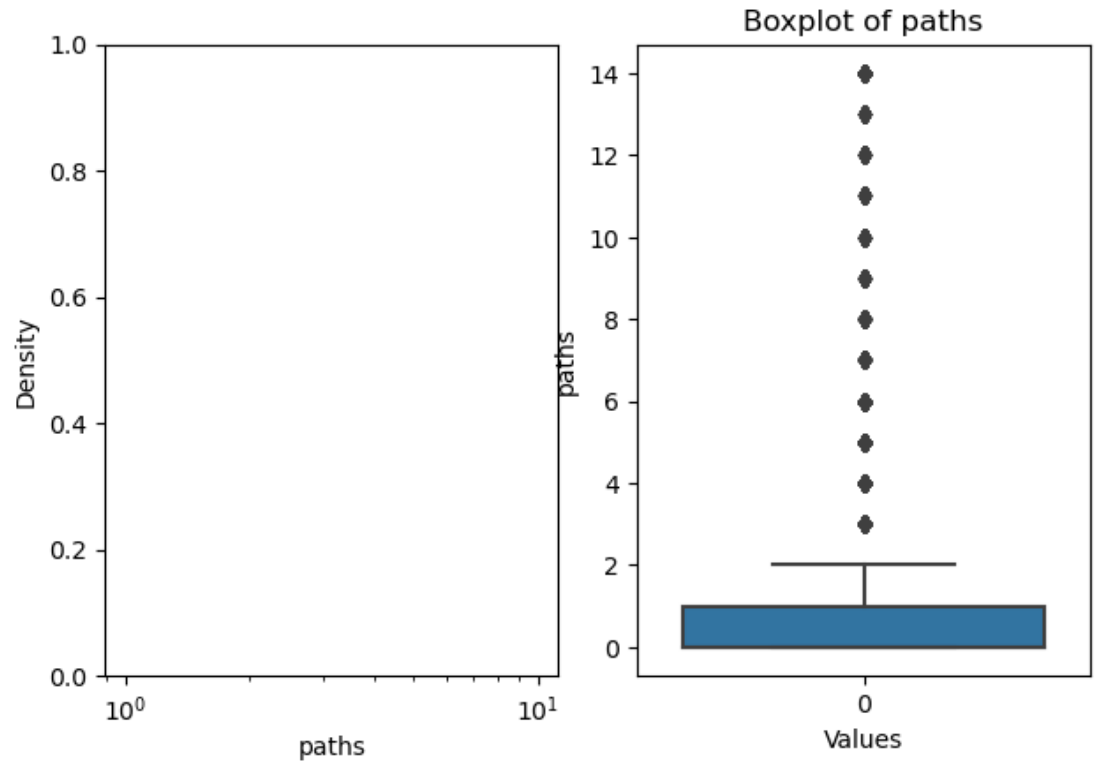
44)

Distribution after Capping



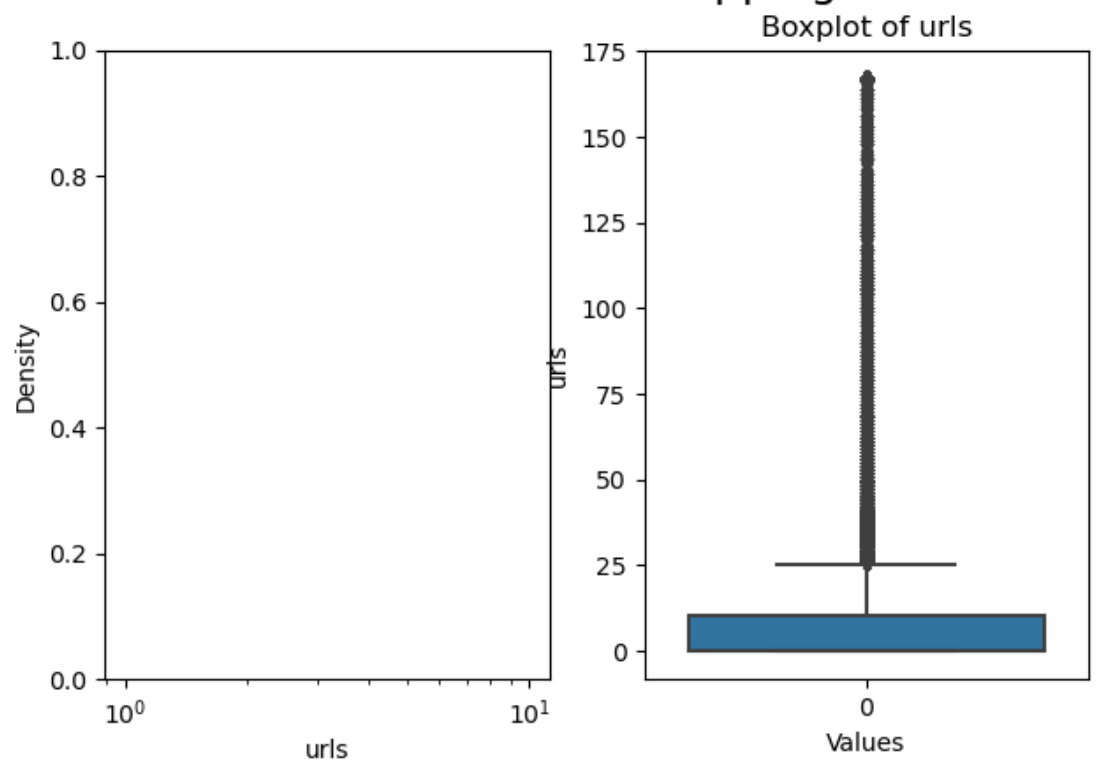
45)

Distribution after Capping



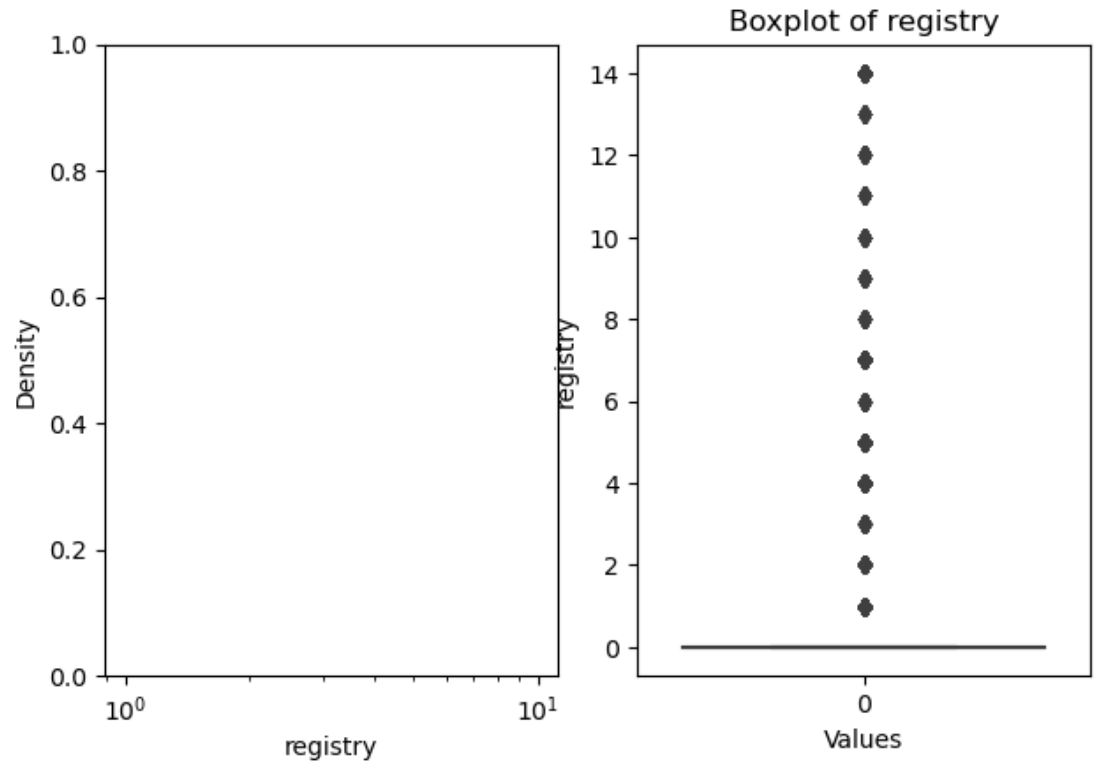
46)

Distribution after Capping



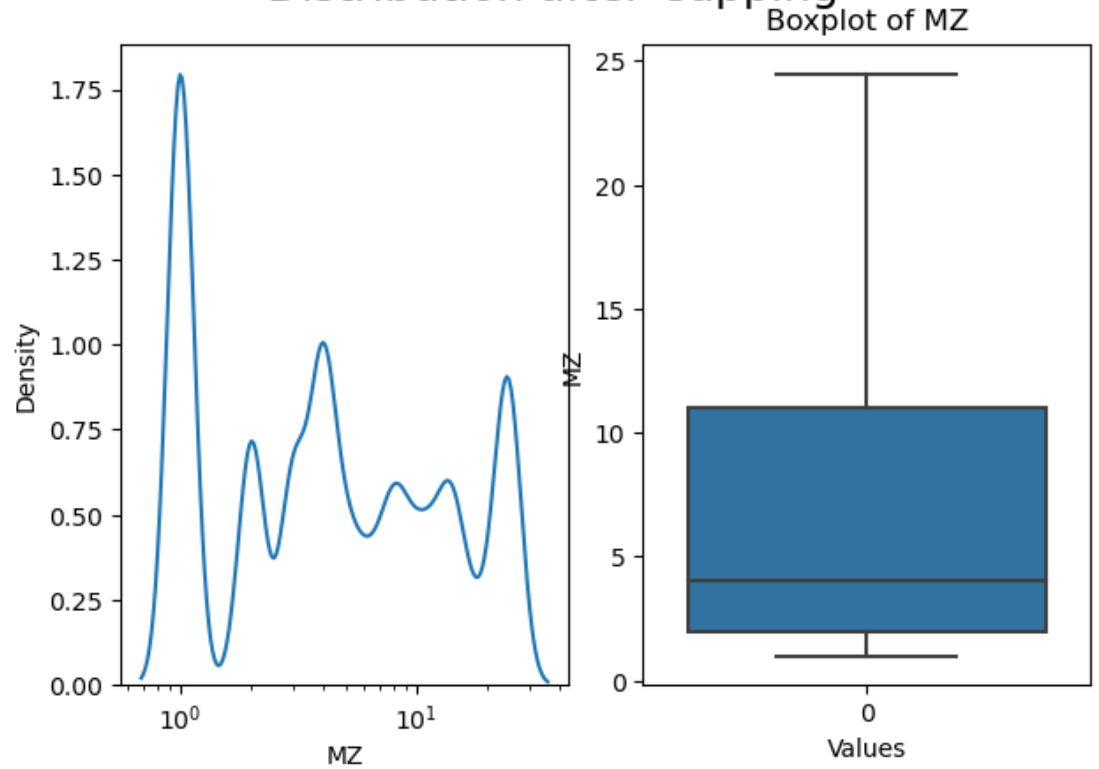
47)

Distribution after Capping



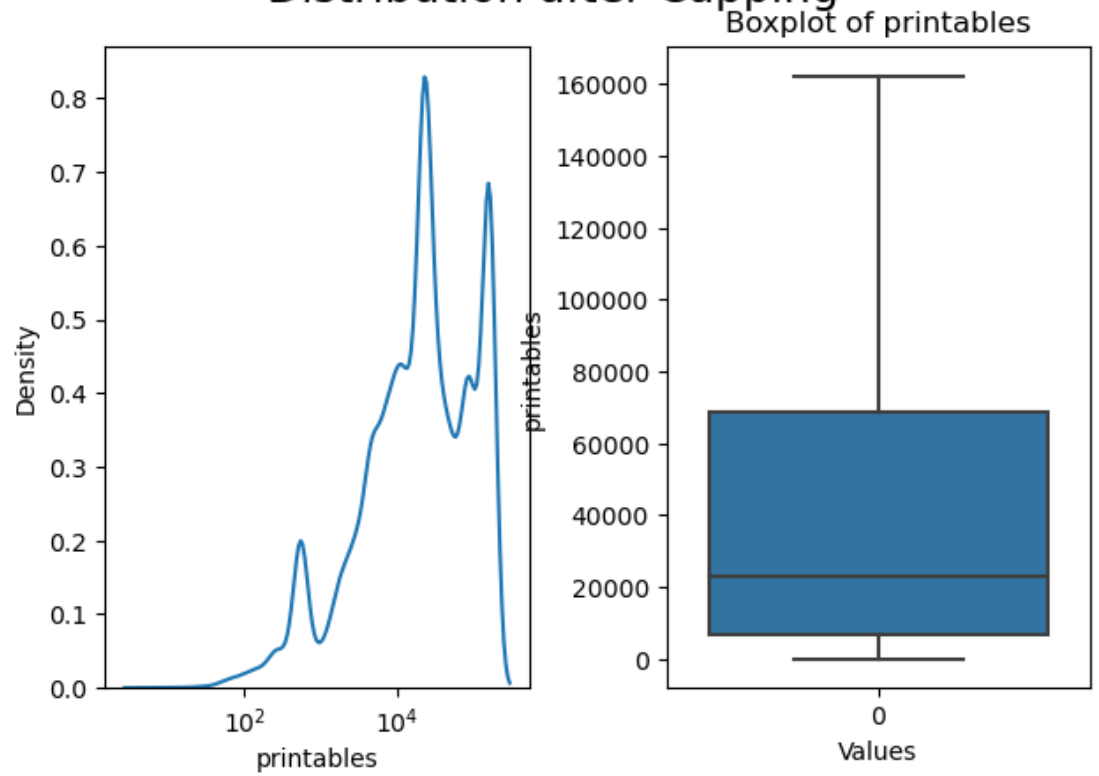
48)

Distribution after Capping



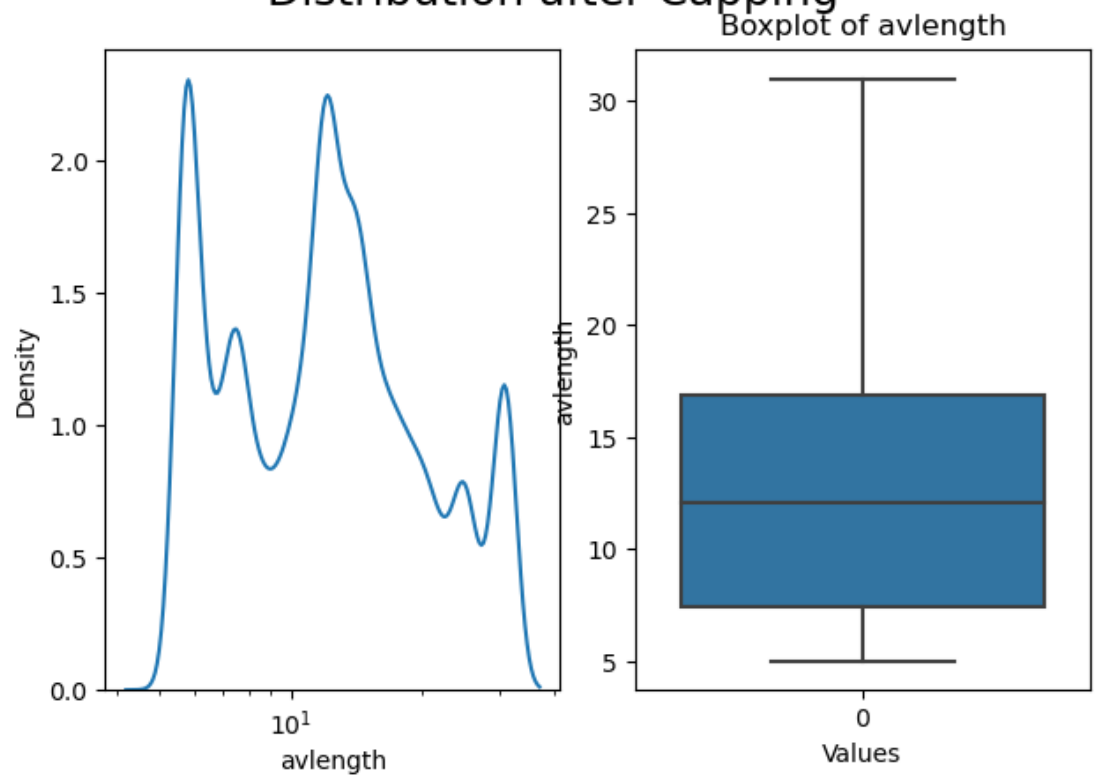
49)

Distribution after Capping



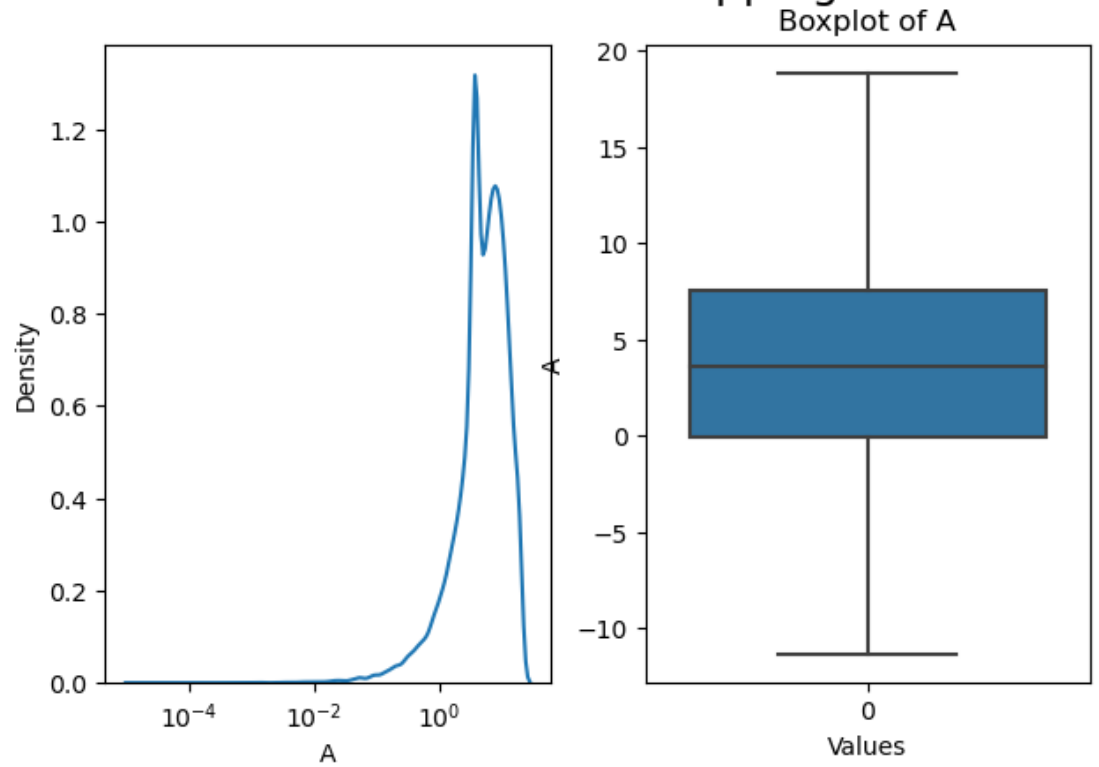
50)

Distribution after Capping



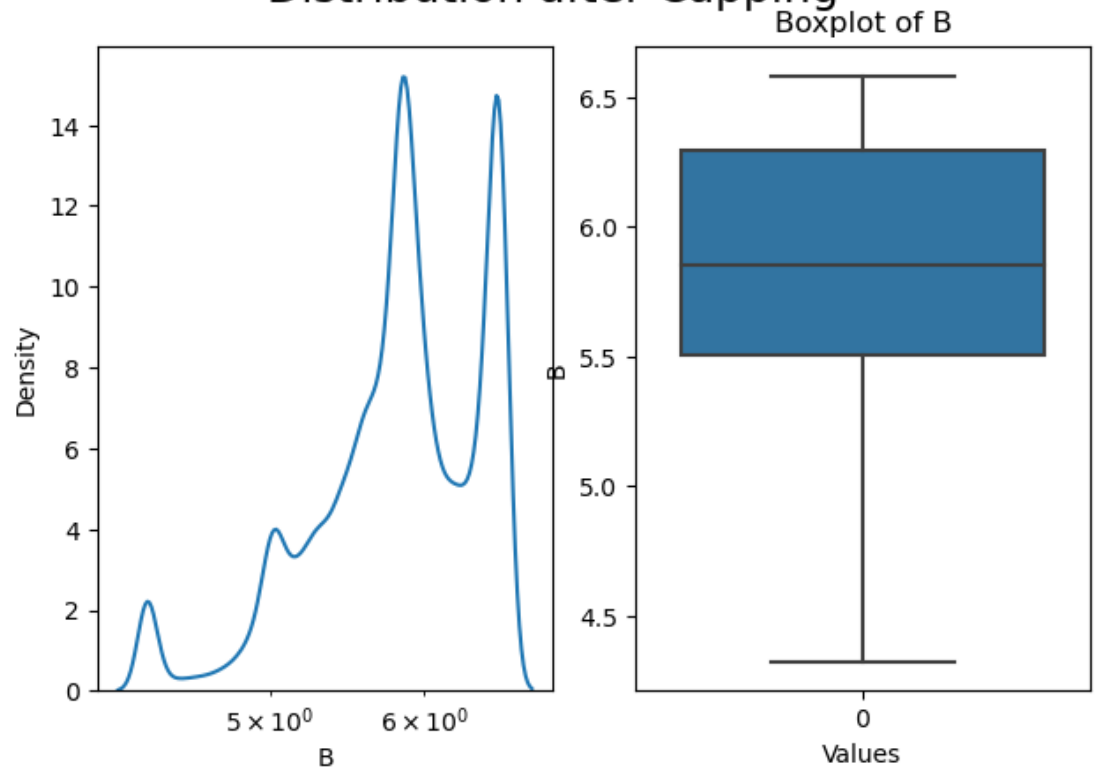
51)

Distribution after Capping



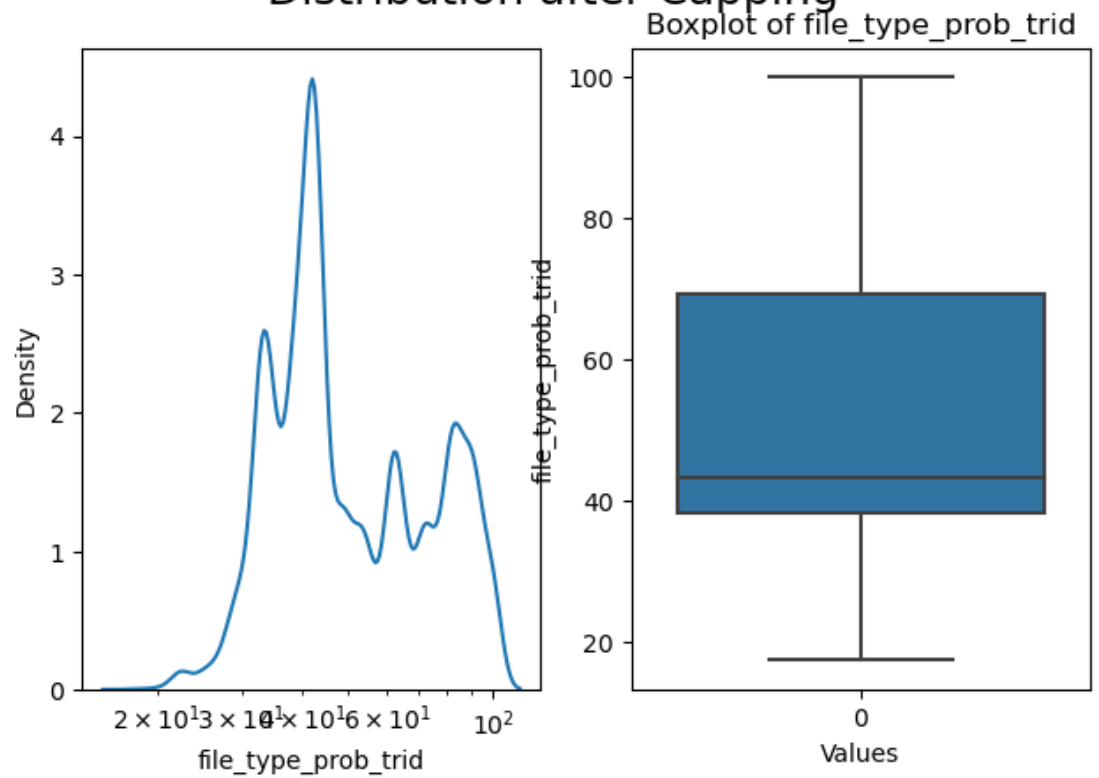
52)

Distribution after Capping



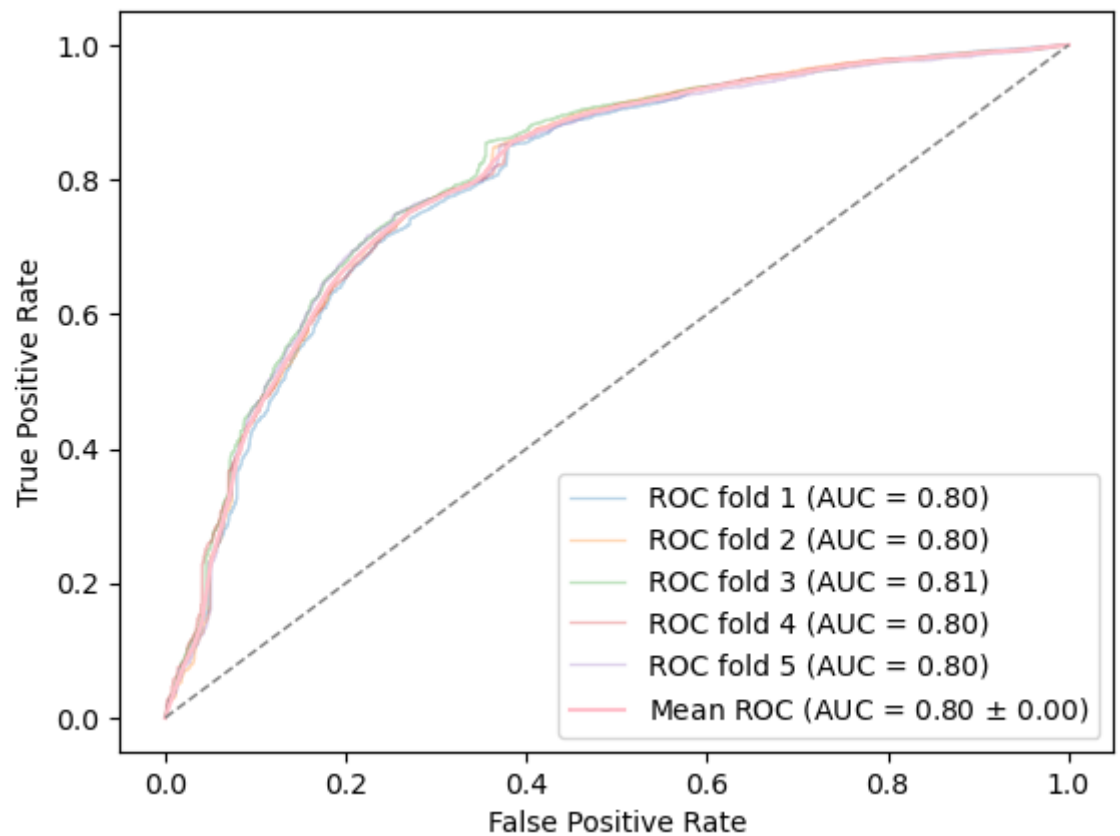
53)

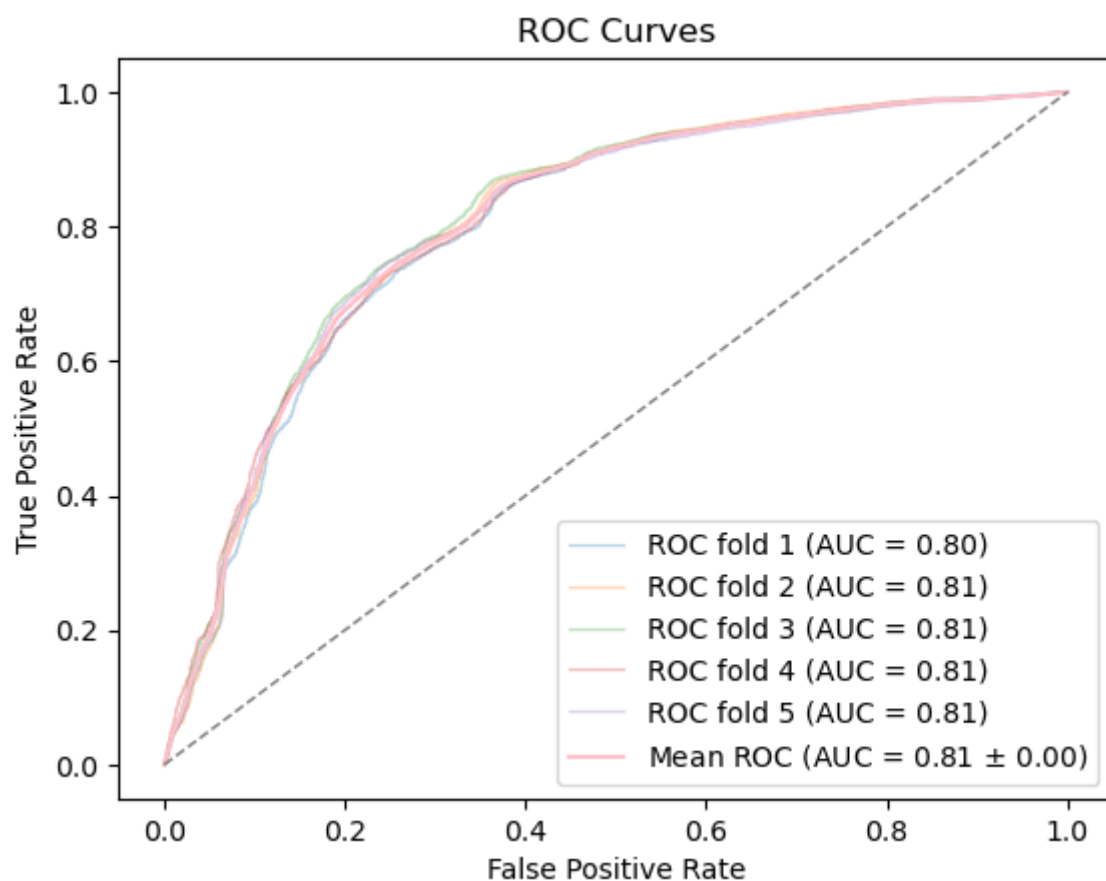
Distribution after Capping



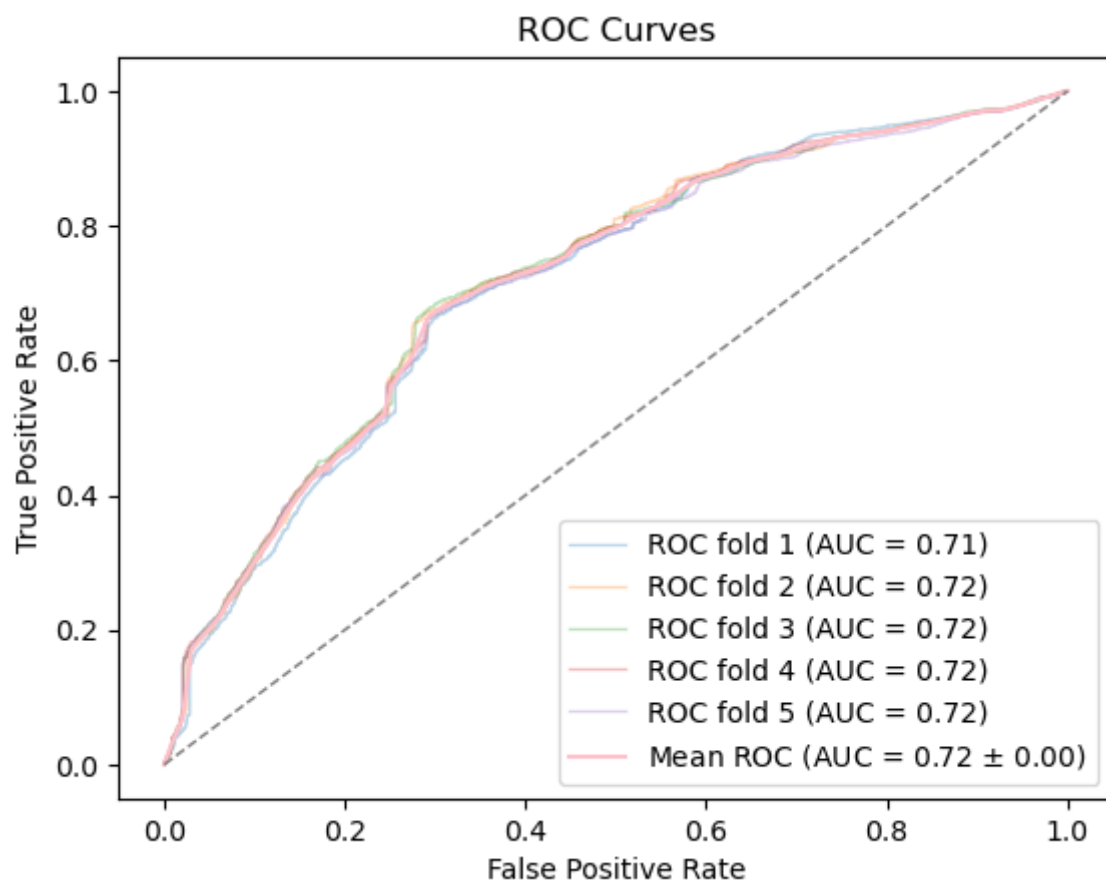
גרסיה לוגיסטית 1 (54)

ROC Curves

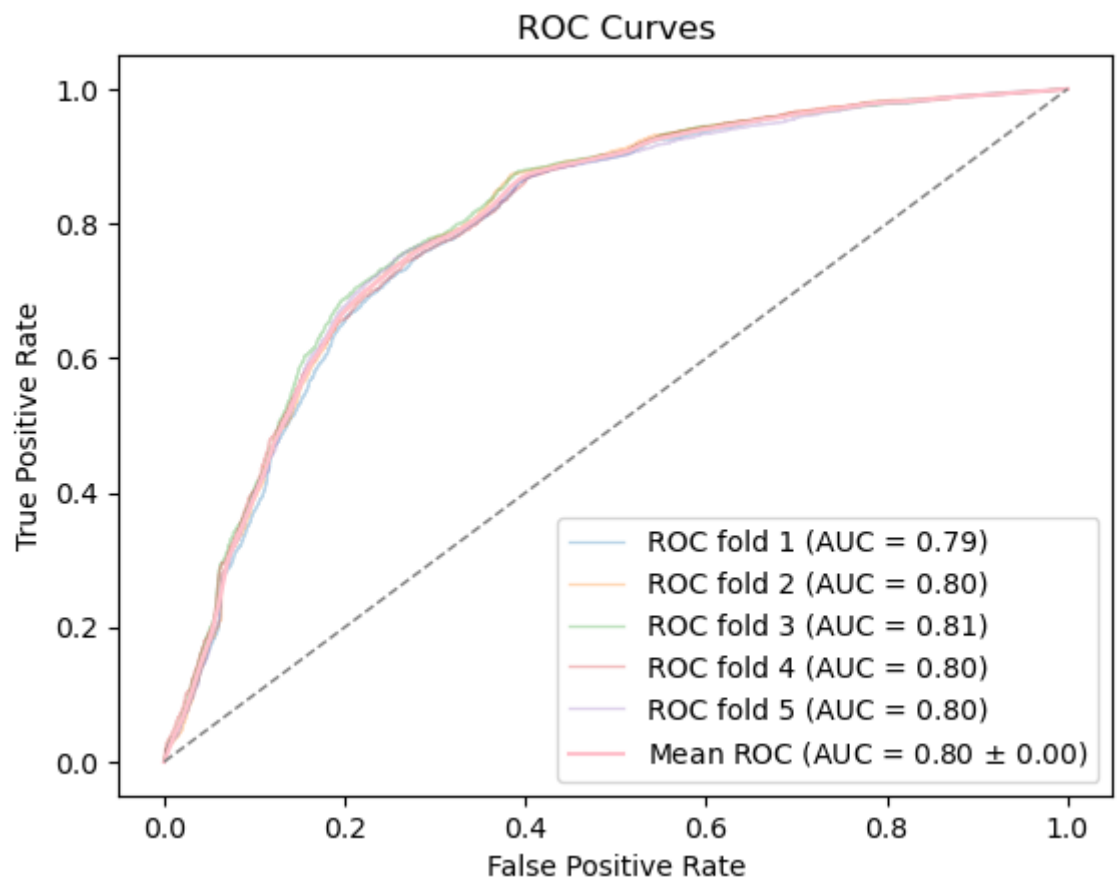




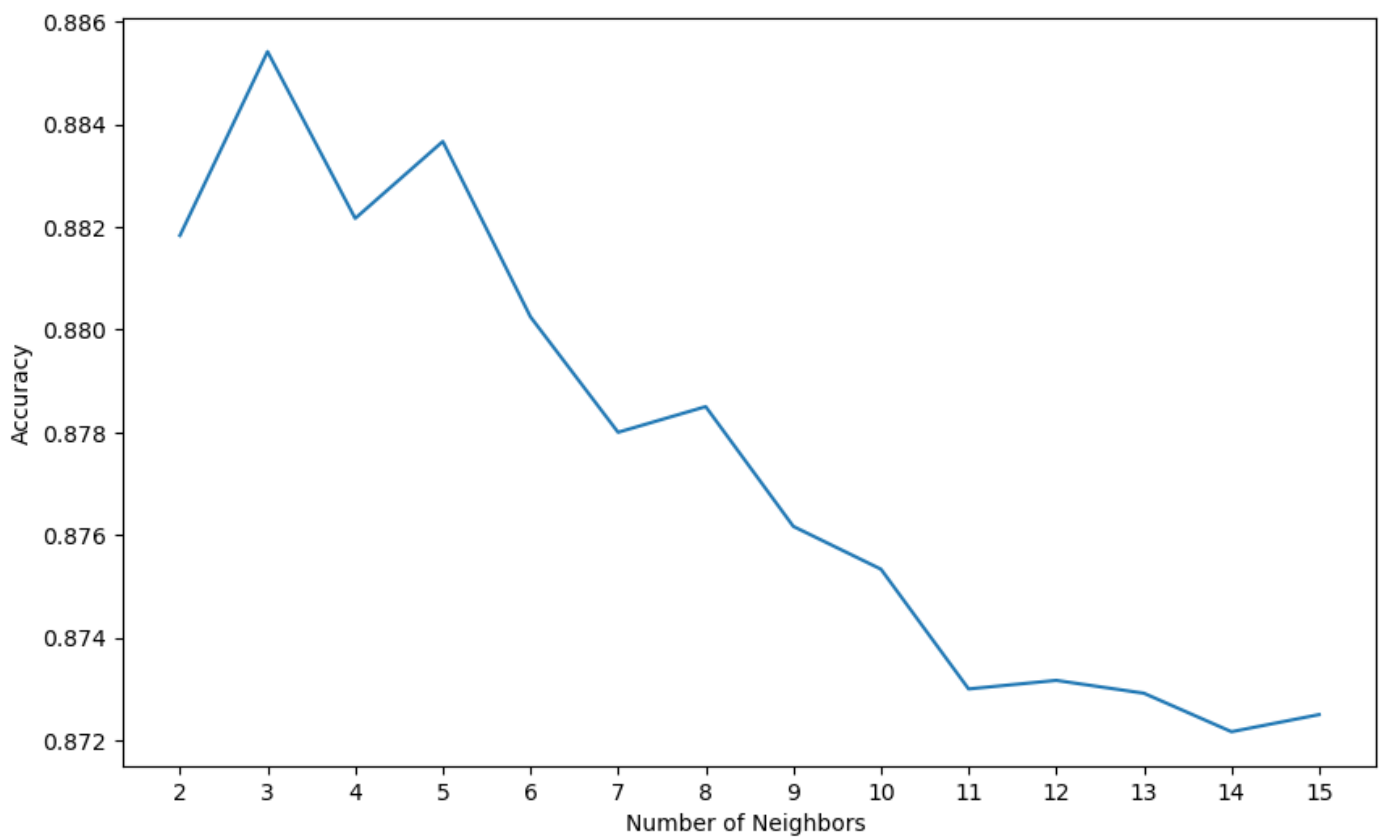
56: 3:



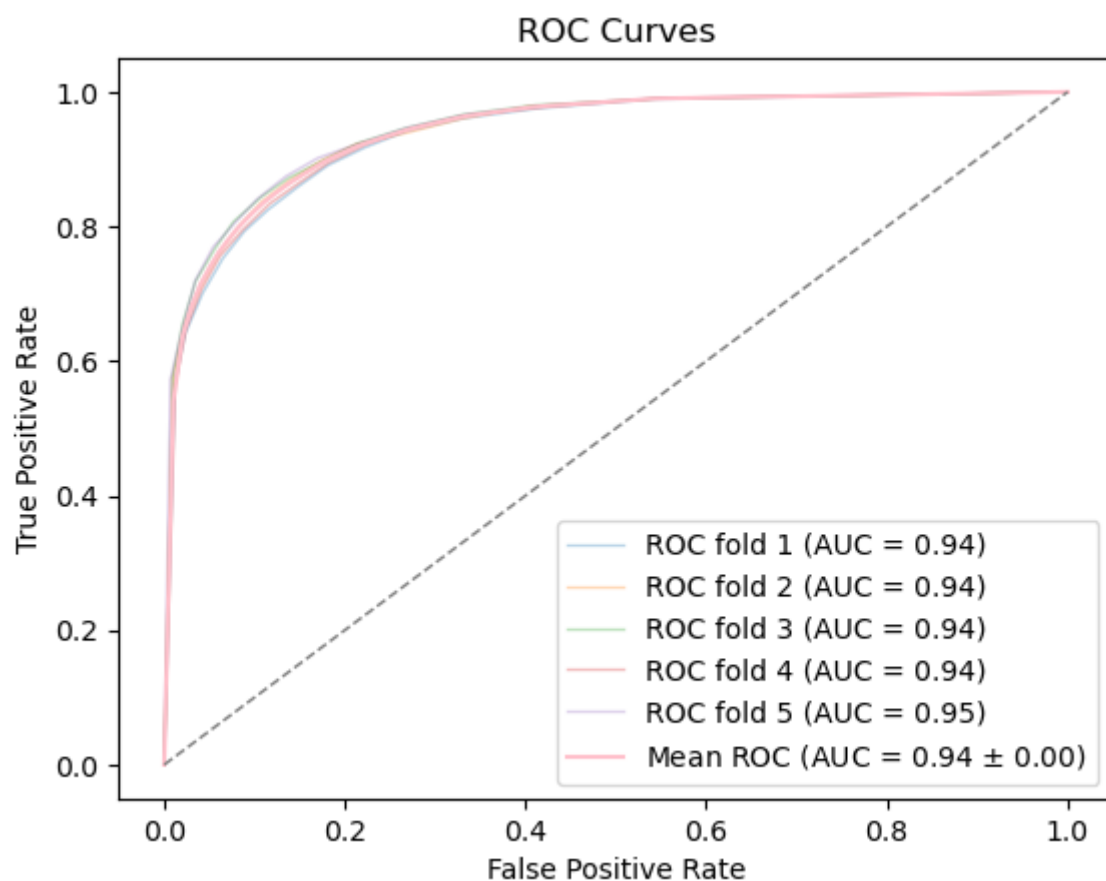
57: 4:



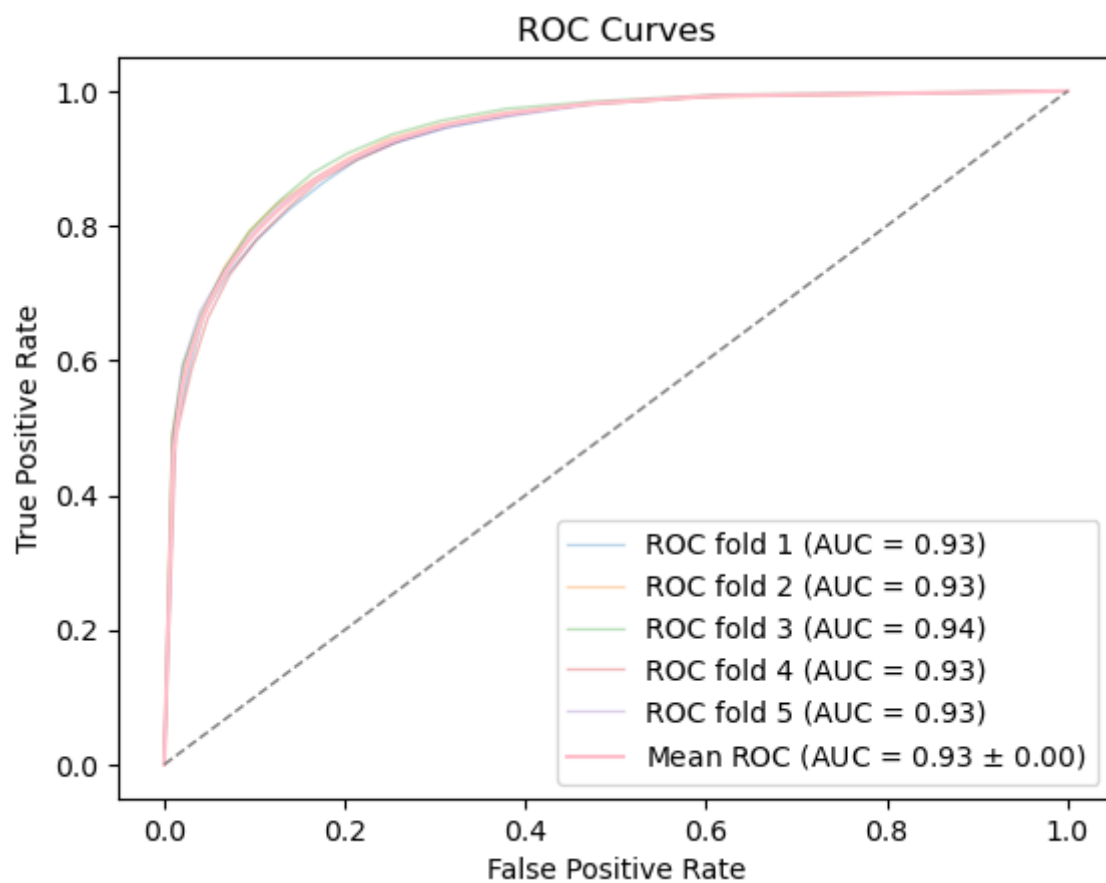
58: KNN 1:



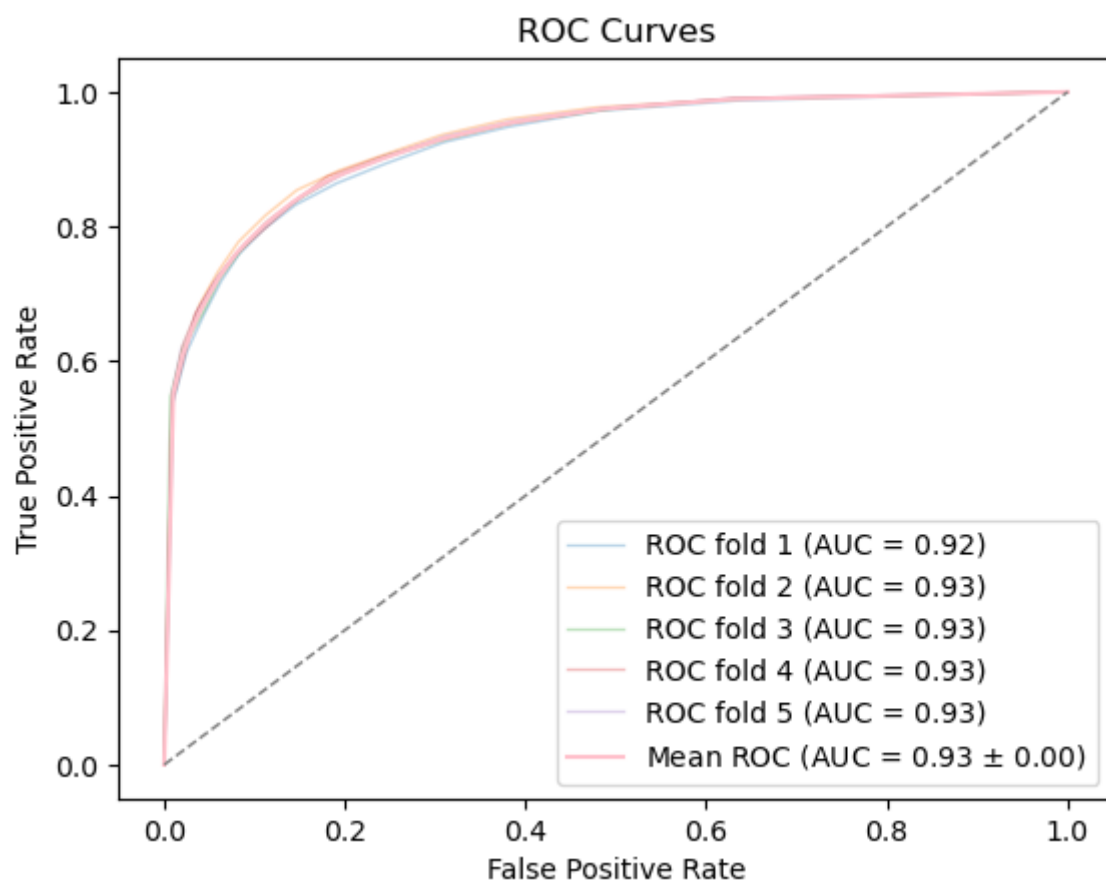
59: 1:



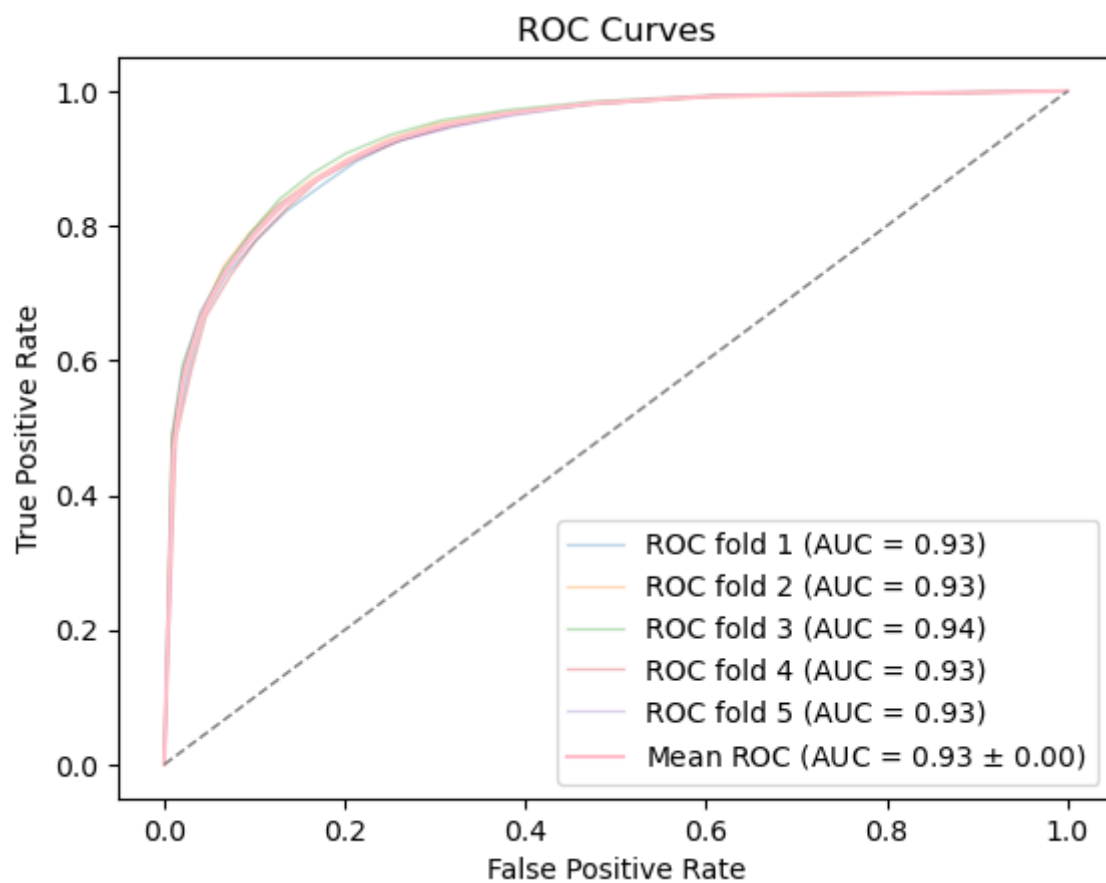
60: 2:



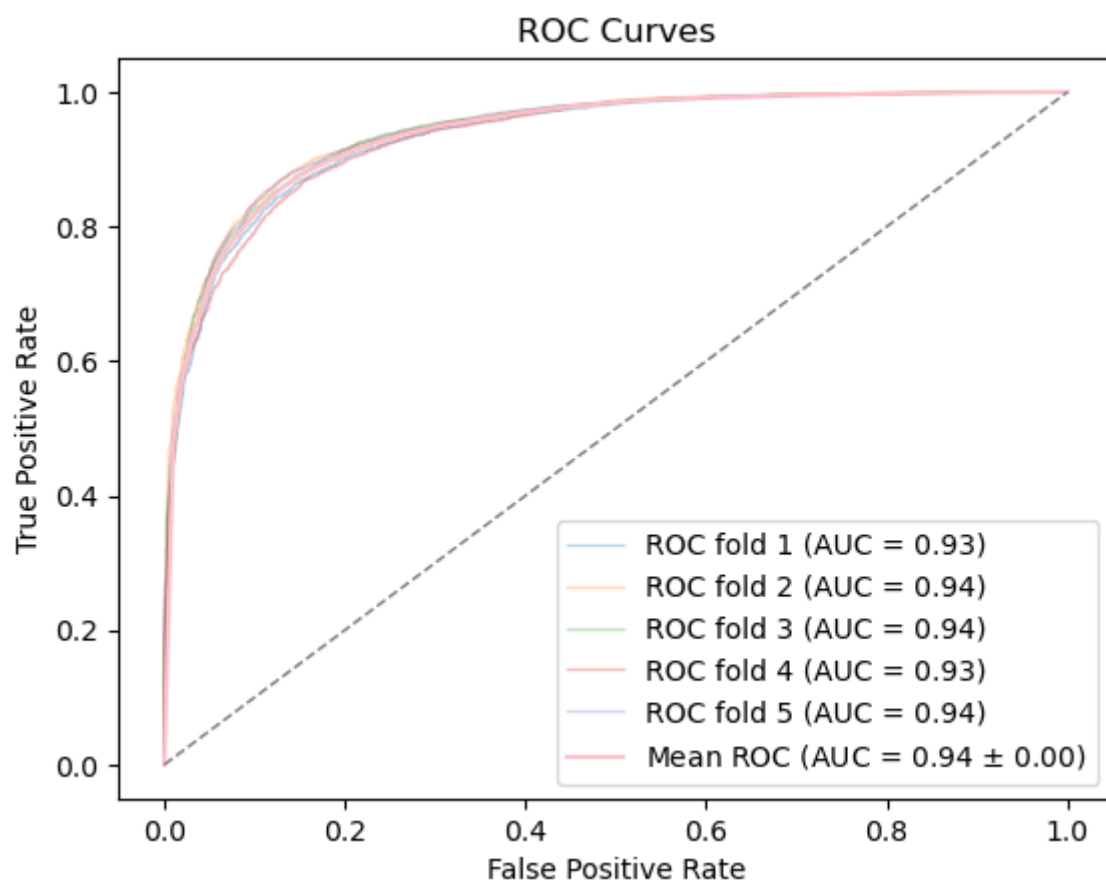
61: 3:



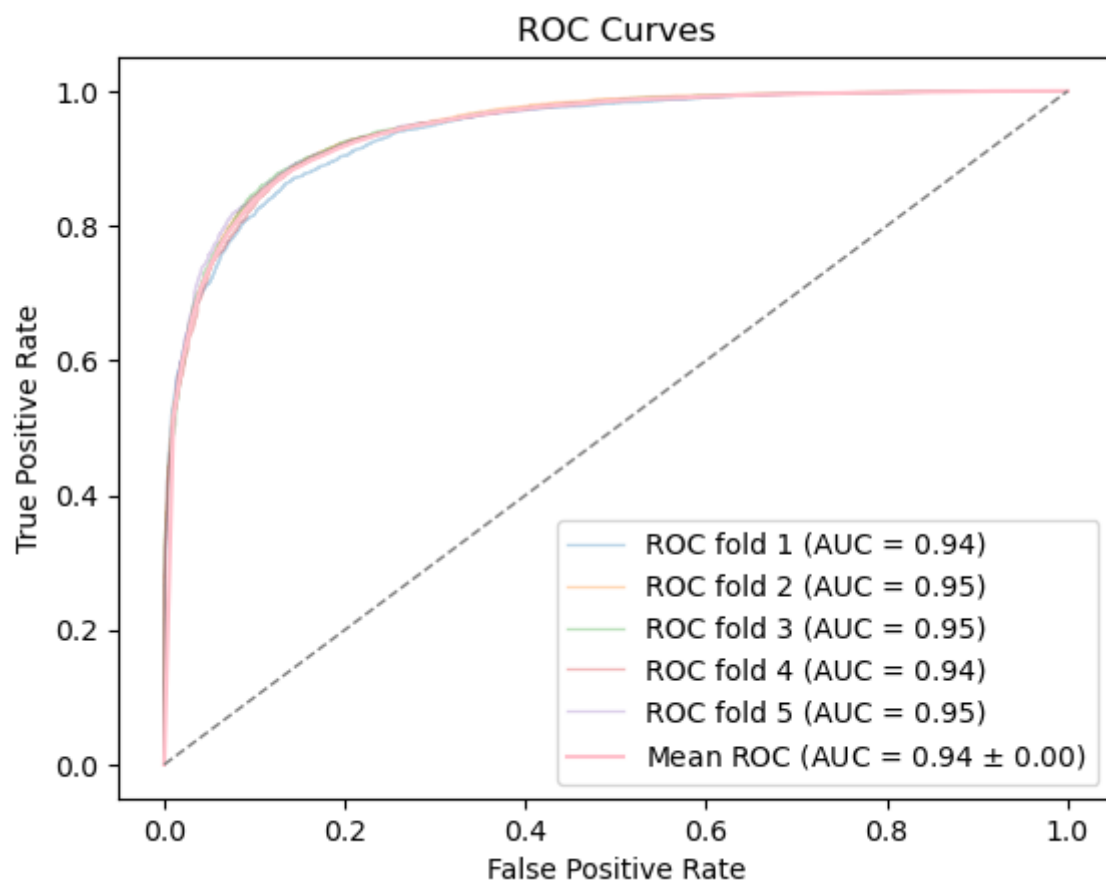
62: 4:



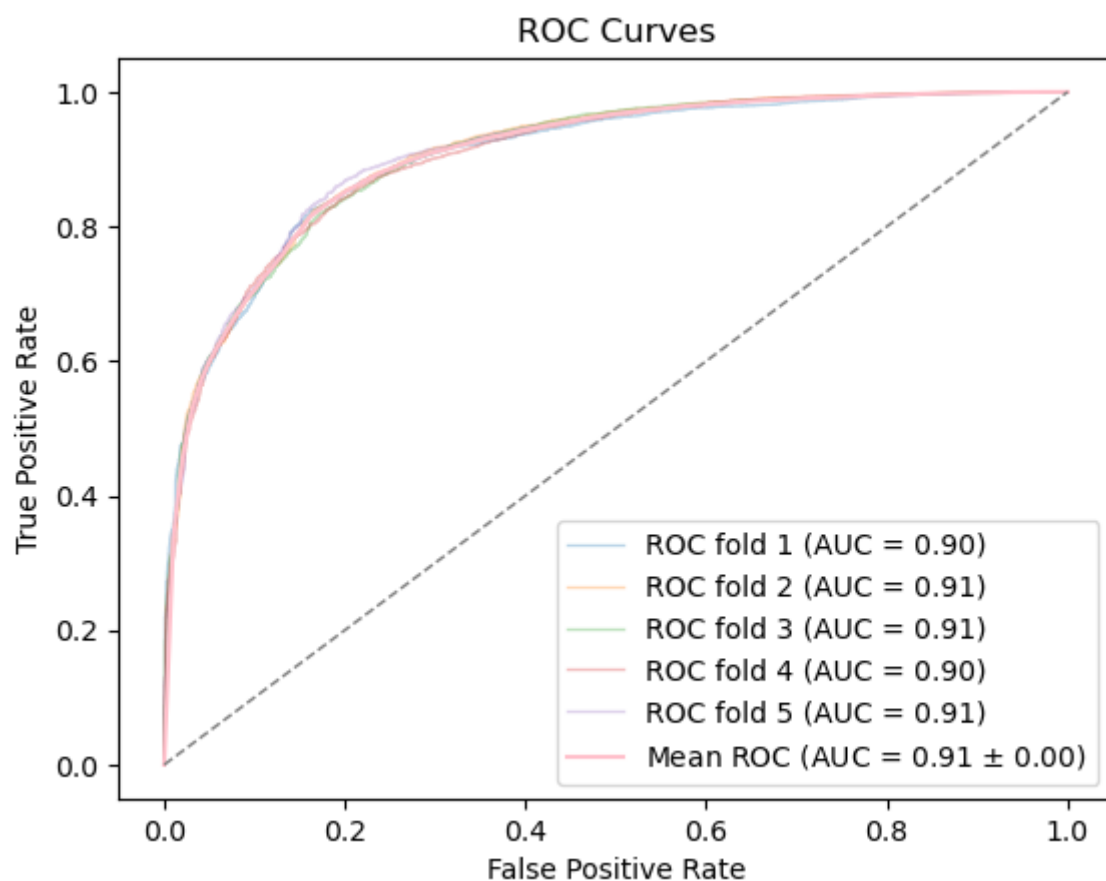
63: mlp 1



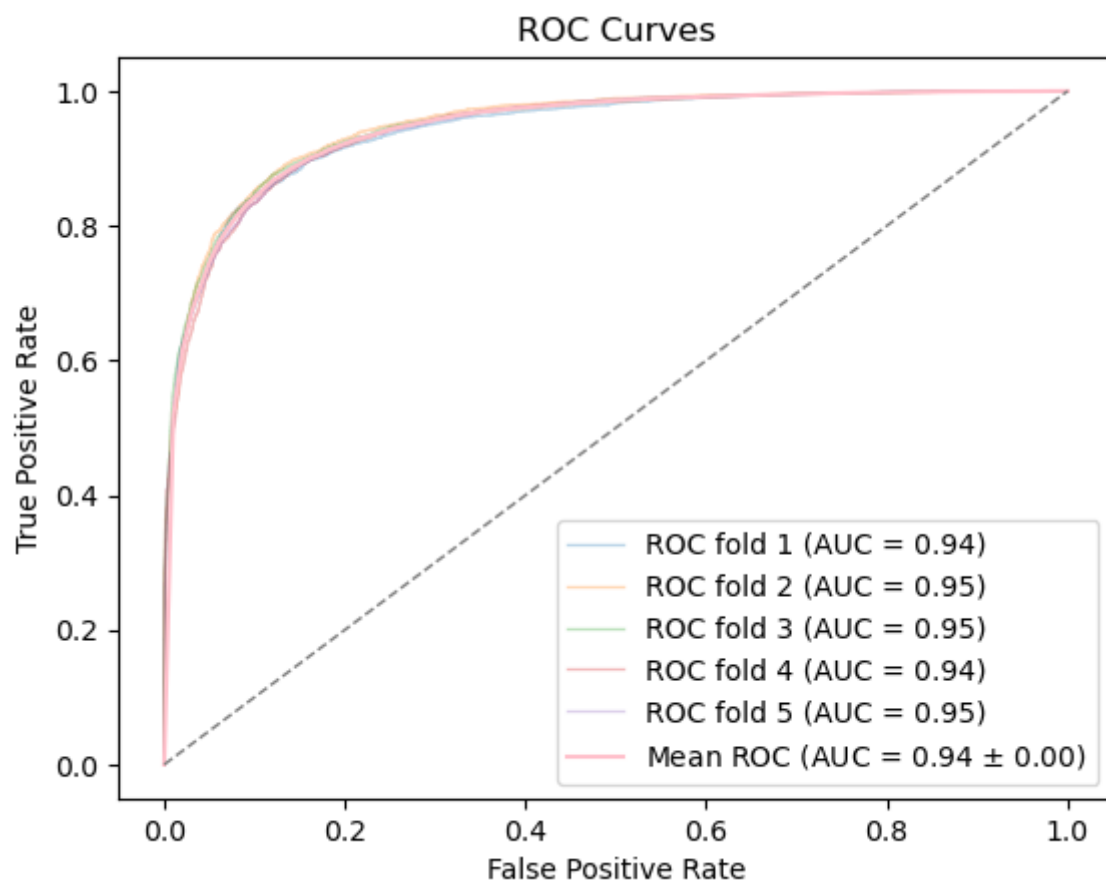
64: 2:



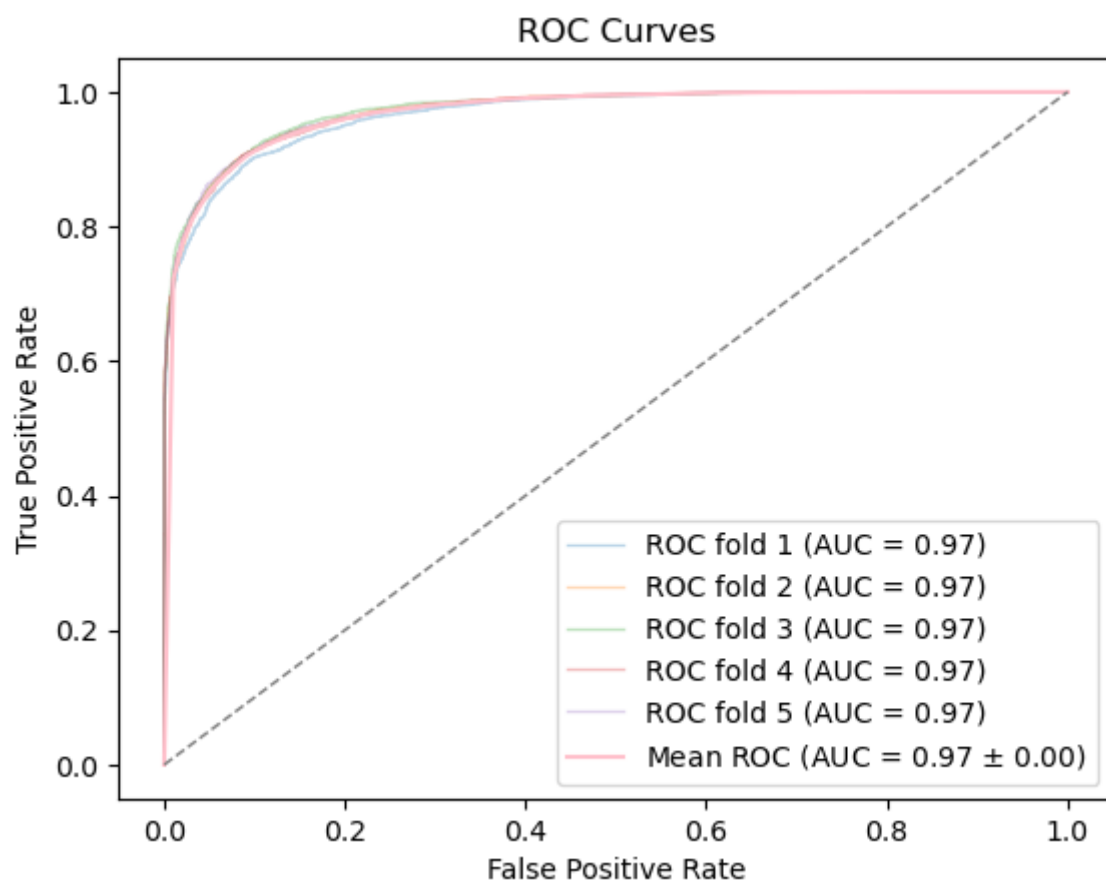
65: 3:



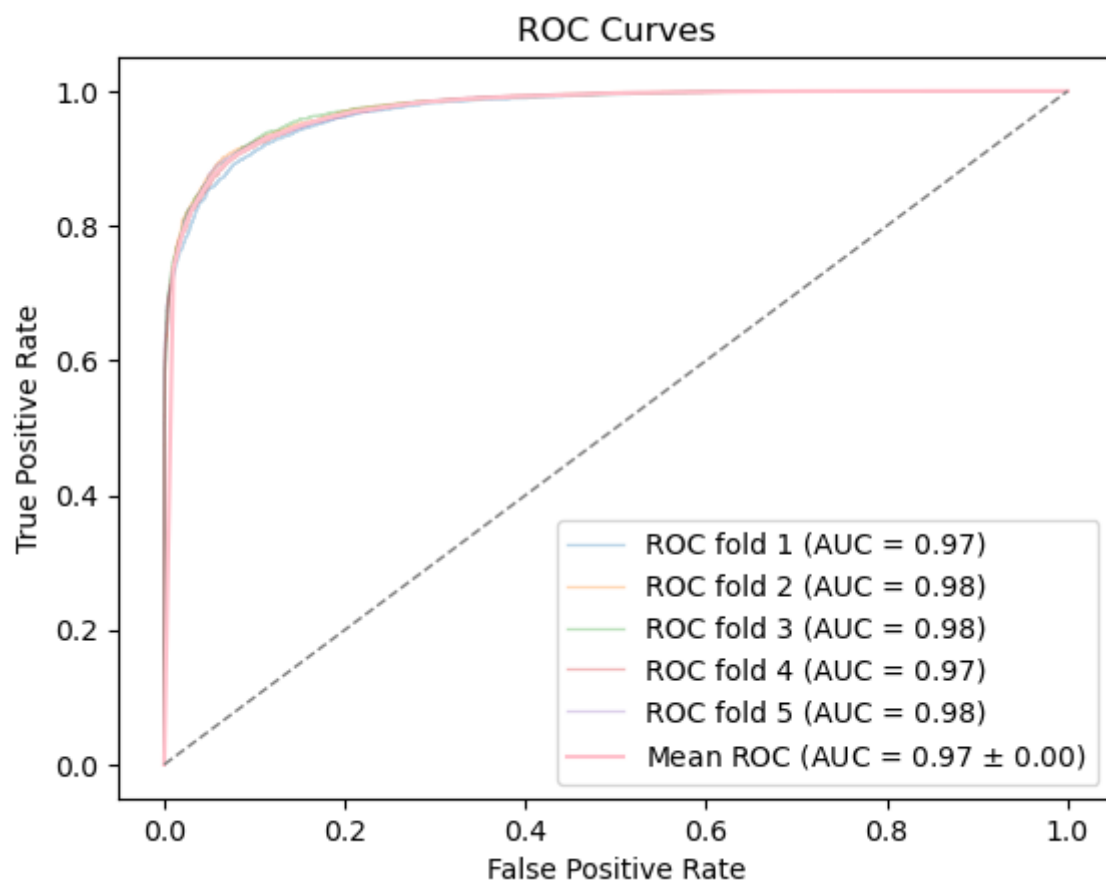
66: 4:



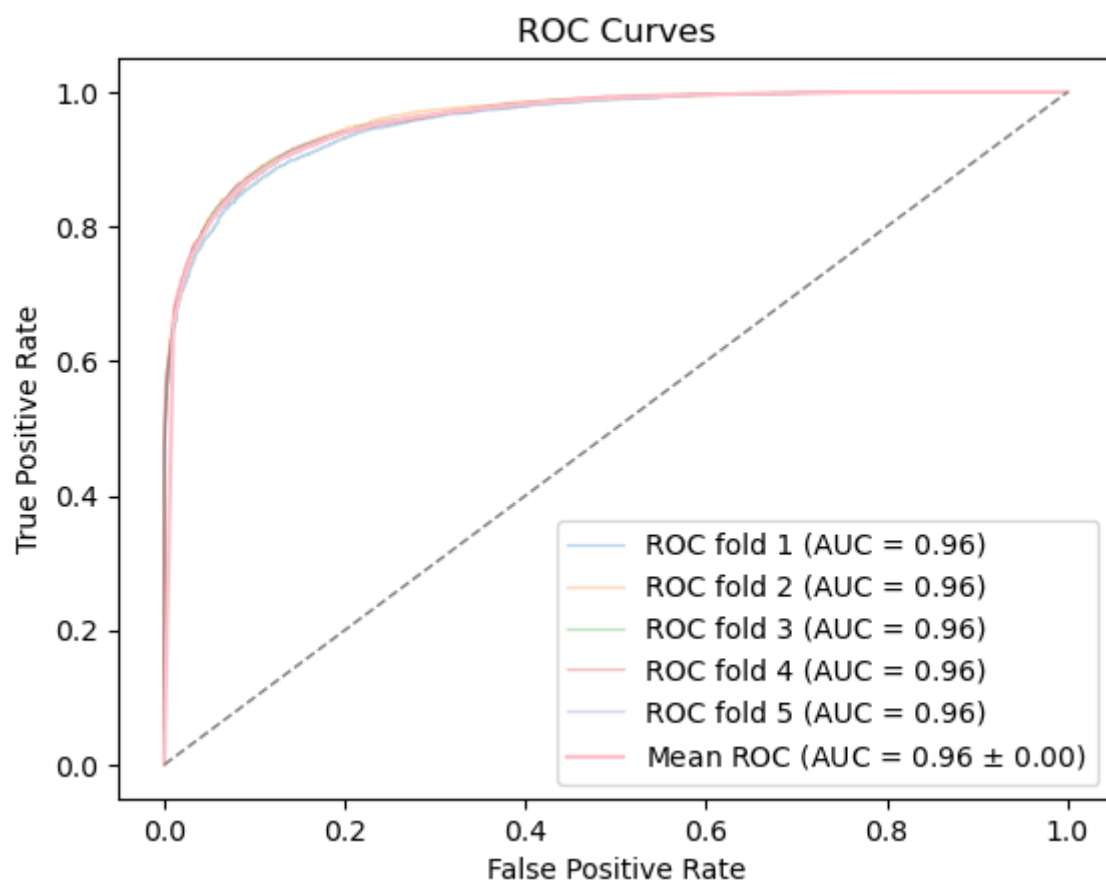
67: עץ רנדומלי 1:



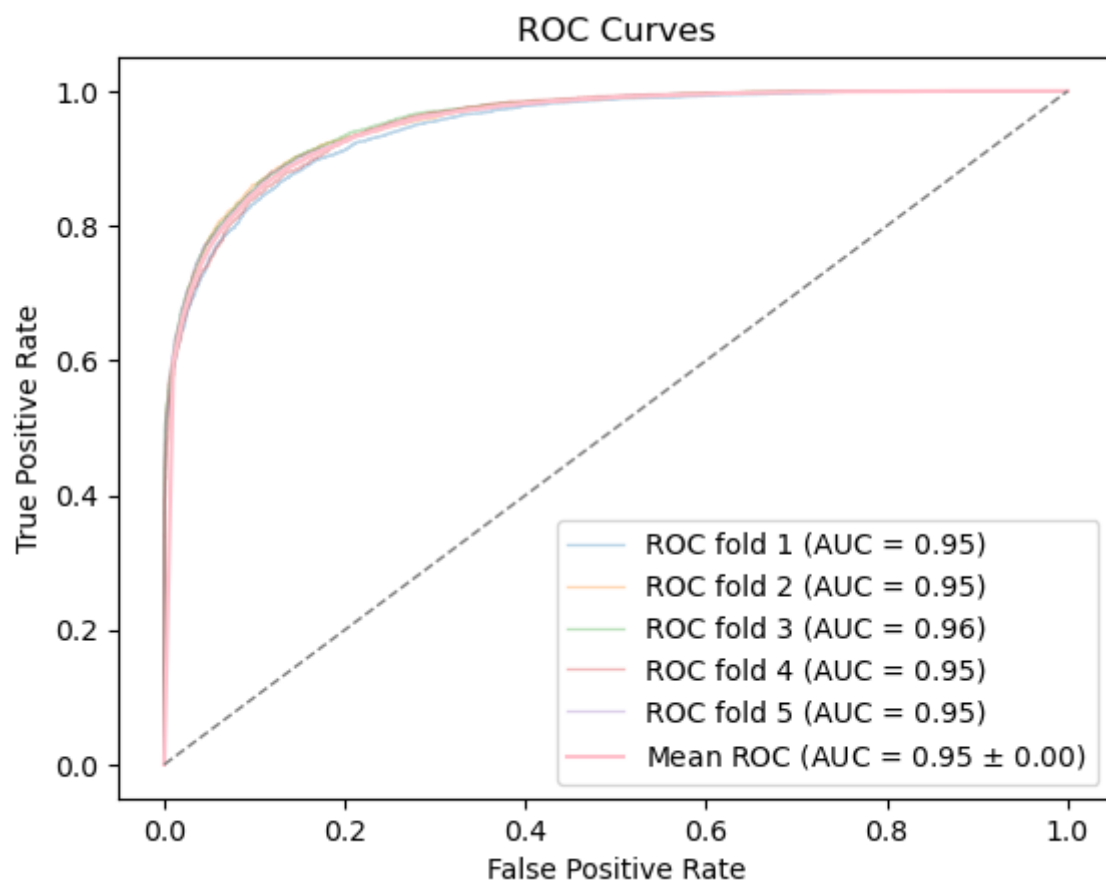
68: 2:



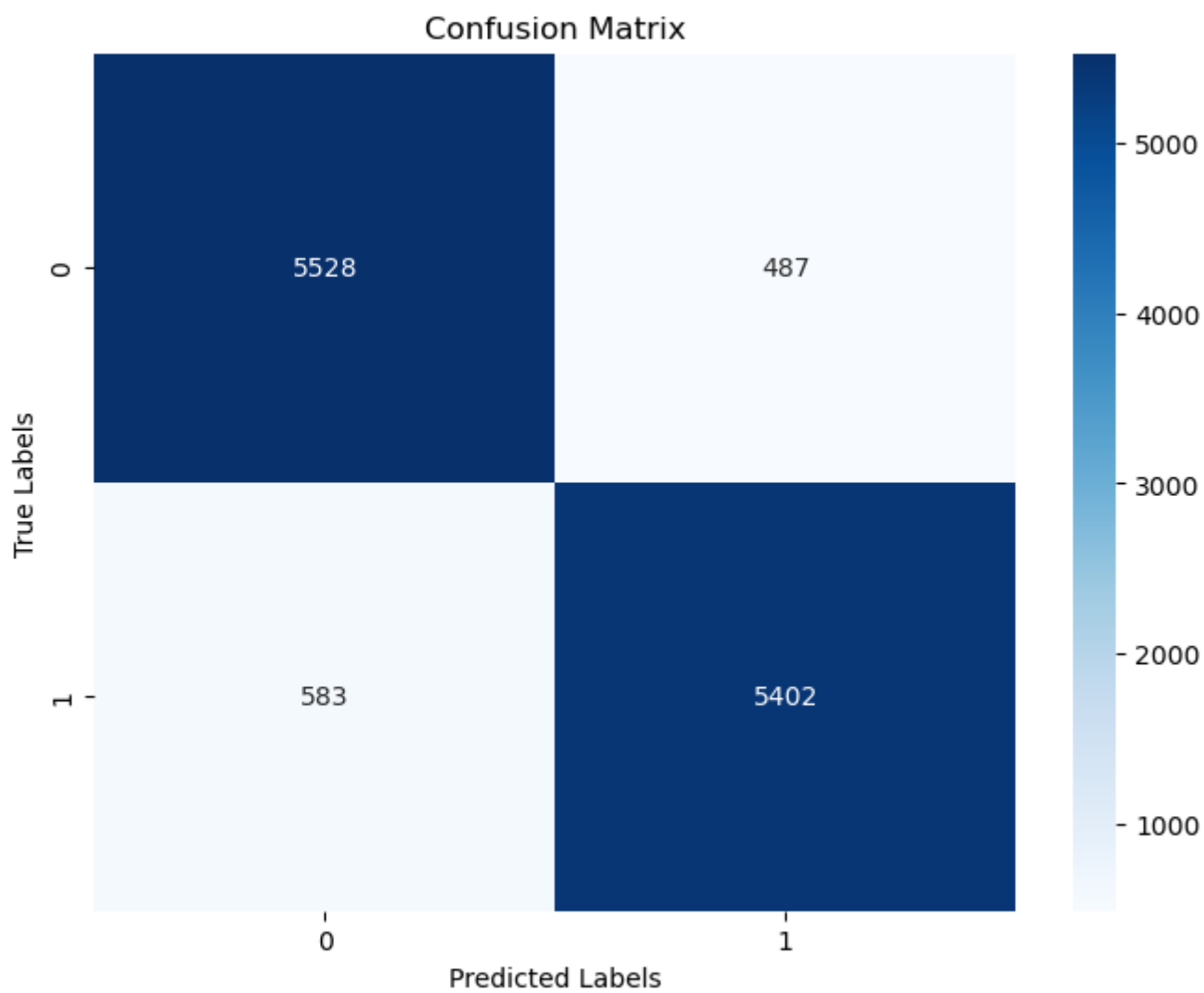
69: 3:

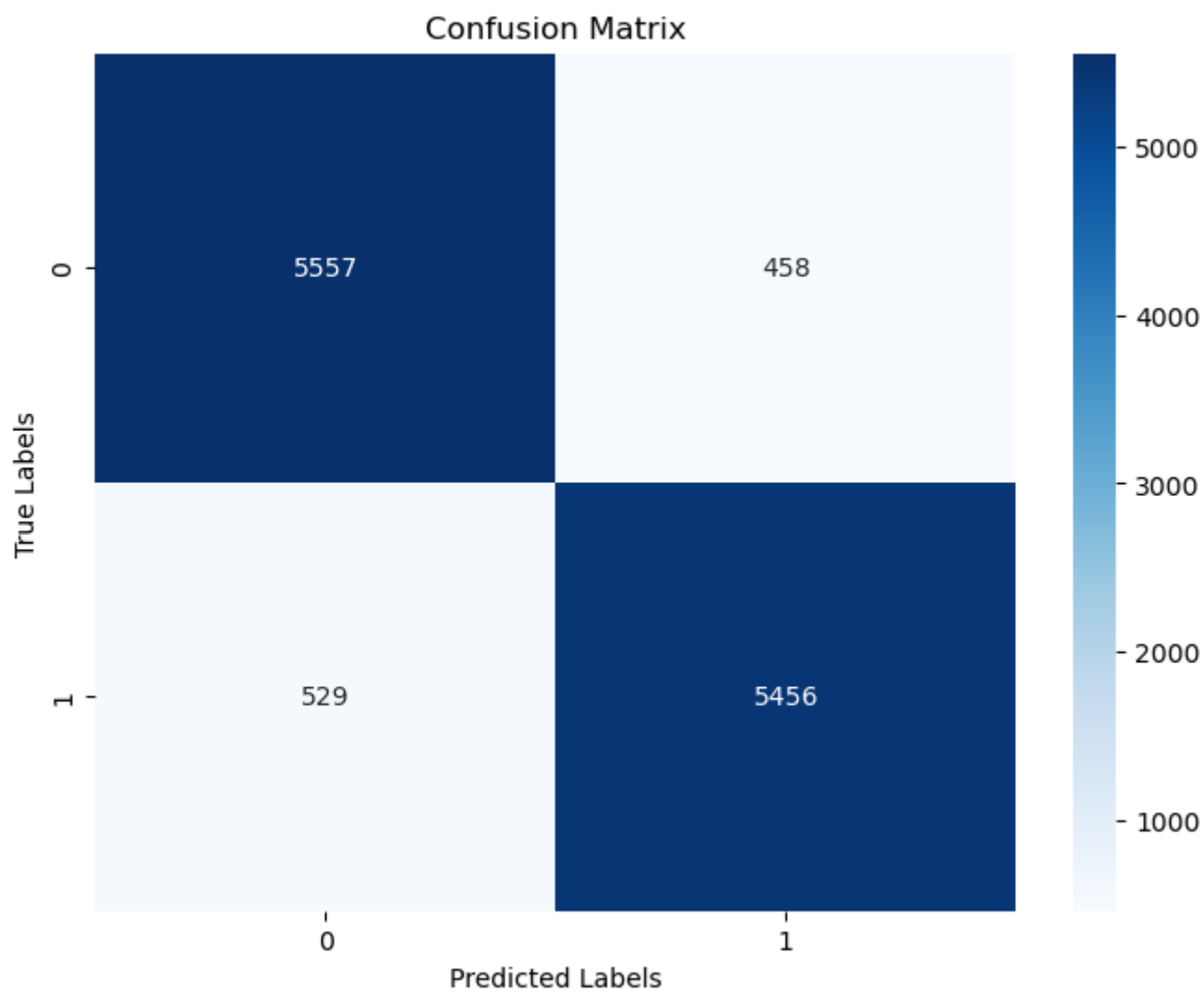


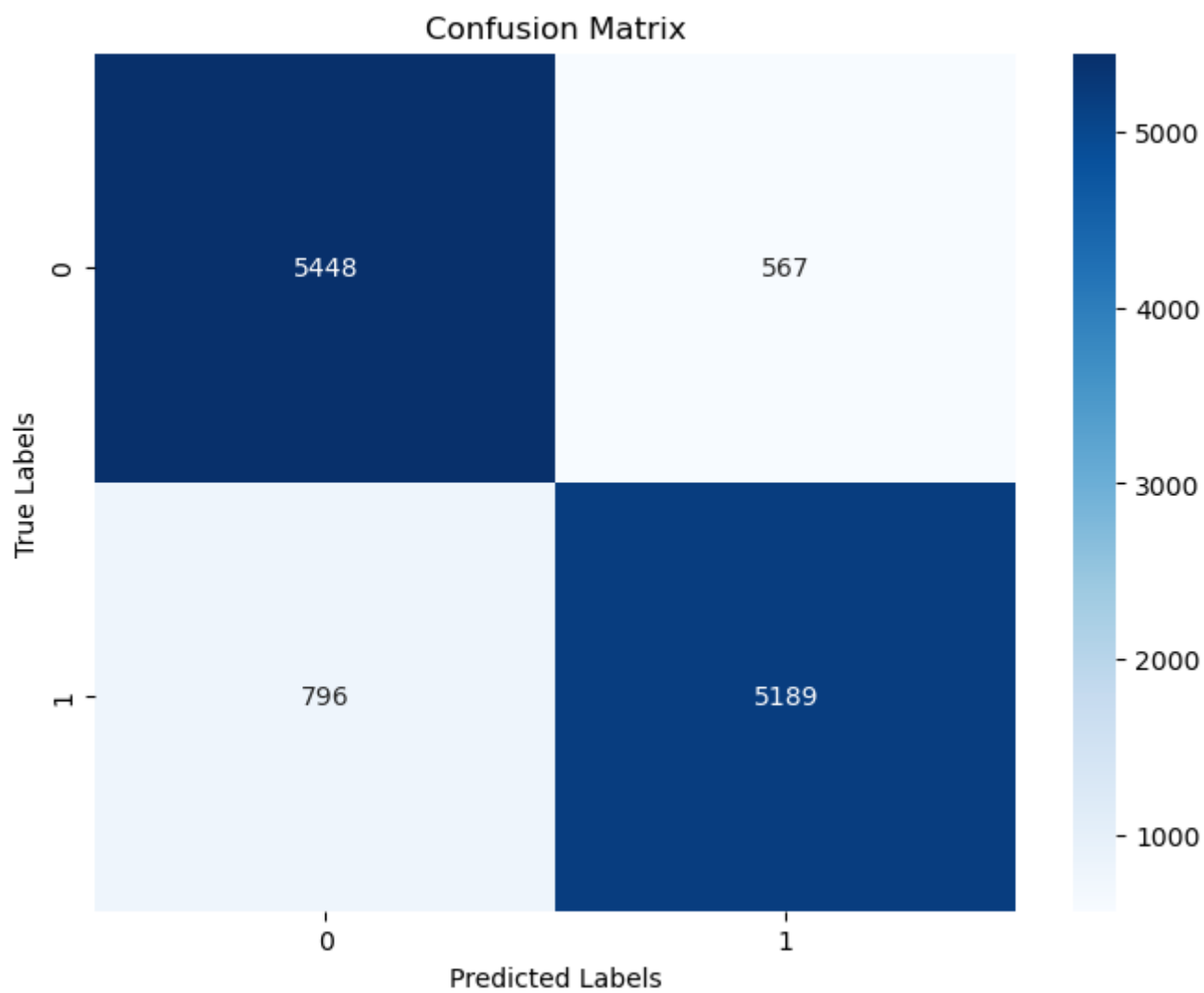
70: 4:

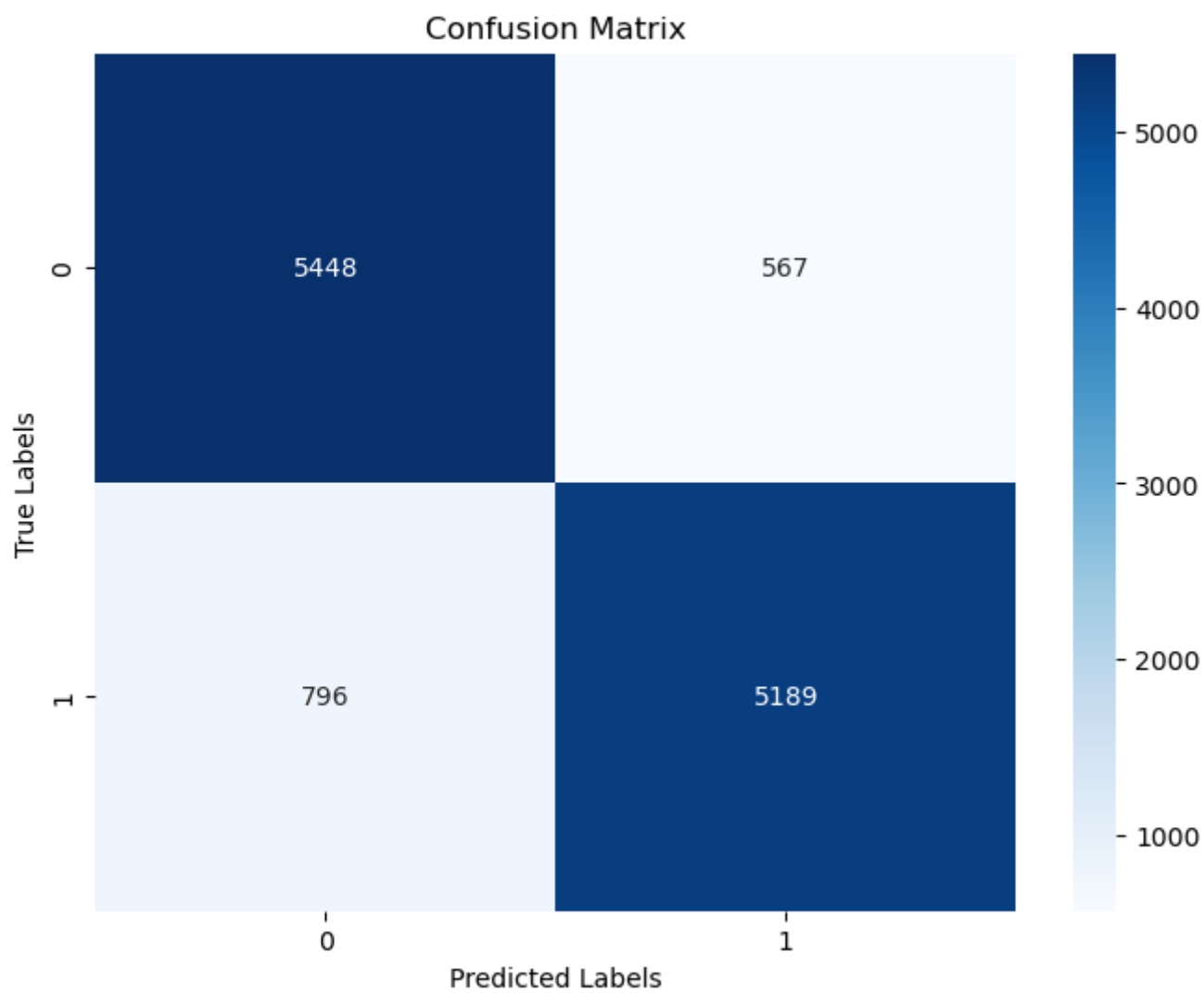


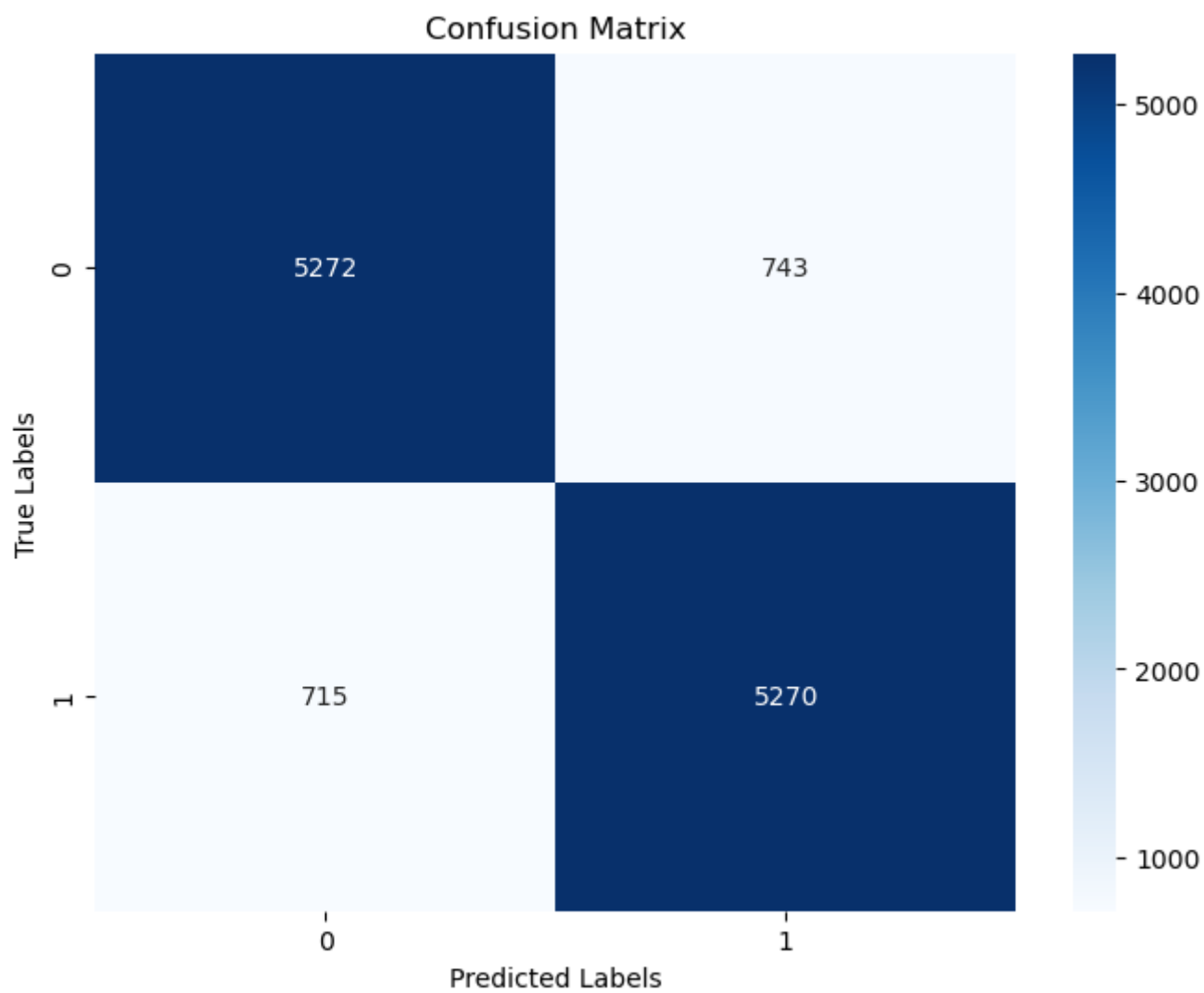
71:



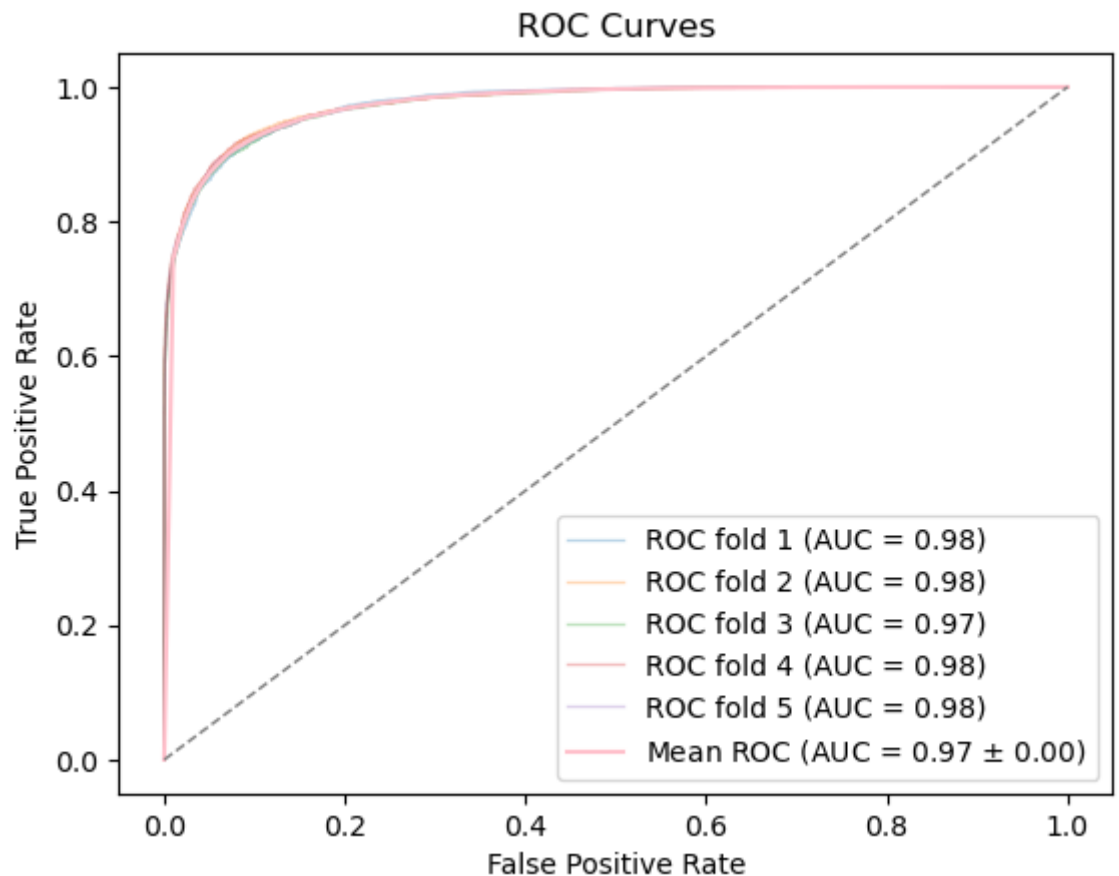








76: פייפלליין:



77*) לכל המודלים הוספנו את GridSearchCV על מנת למצוא את הקומבינציה הכי טובה של היפר פרמטרים לכל מודל. בכולם בחרנו את scoring של roc_auc והגדרנו refit=True על מנת שהפרמטרים שמצאנו כהיפר פרמטרים הטובים ביותר למודל, ייושמו על סט האימון ובנוסף שהערכת סט הולידציה תהיה על בסיס סט האימון המשודרג הנ"ל. בנוסף $n_jobs = 1$ מייצג את העבודה המקבילית שתחסוך זמני ריצה, $cv=5$ זה הברירת מחדל. CV זה cross validation שבו נשתמש במהלך ה GridSearchCV והוא מאפשר לנו להחליט איך לפצל לסט אימון וסט ולידציה לכל קומבינציה של ההיפר פרמטרים.

78*) מודל נוסף שבדקנו ובסוף לא השתמשנו בו הוא מודל **Naïve bayes** אך בסוף לא השתמשנו בו עקב תוצאות לא טובות. הפרמטרים היחידים שניתן להוסיף הם 'var_smoothing', 'priors'. השתמשנו בברירת המחדל מכיוון ש'priors' מייצג את ההסתברות הפירורית לכל מחלקה, ומכיוון ואין לנו ידע מקדים על הדאטה שיתמוך בערכים מסוימים שנשים ל'priors', העדפנו להשתמש בברירת המחדל בה האלגוריתם מחשב את ההסתברויות המקדימות של כל קבוצה על פי סט האימון, מהנחה שסט האימון מספיק גדול ויחסית מייצג את כל האוכלוסייה. הגדרת ה'prior' יכולה להשפיע על ההטייה. אם לא היינו בוחרים בברירת המחדל, האלגוריתם היה מתייחס להגדרת prior שננתנו לו ונותן משקל רב יותר למחלקה המתאימה ובכך יכולנו לשנות את ה bias. מההנחה שסט האימון מייצג נכונה את האוכלוסייה אנחנו מקטינים את ה bias שיכל להיווצר אם היינו מגדירים אקראית את ההסתברויות הפירוריות.

בנוסף השתמשנו בברירת המחדל ב'var_smoothing' פרמטר זה שולט ברמת ההחלקה של השונות. ככל שנגדיל את הערך של var_smoothing כך תגדל השונות (ותקטן ההטייה בהתאם) מכיוון שזה מוסיף ערך קבוע לשונות הפיצ'רים. כדי להימנע מ overfitting או הטייה גדולה מדי, העדפנו להשתמש בברירת המחדל.

Confusion matrix(79)

TP	FP
----	----

FN	TN
----	----

נראה ככה:

באופן כללי מדד accuracy מסביר כמה חוזו בצורה נכונה (correctly predicted) מתוך כלל הדוגמאות. כלומר בהינתן דוגמה רנדומלית, מה הסיכויים שחזה בצורה נכונה?

כאשר התשובות מוטות לצד אחד עדיף להשתמש במדד ה precision או מדד ה recall ובשמו הנוסף sensitivity.

מדד ה precision מייצג: כשחוזים שמשהו יקרה, מה הסיכויים שאפשר להיות בטוחים בשזה באמת קרה?

מדד ה recall מייצג: מתוך הדוגמאות הנכונות (True), מה הסיכויים שצדקתי?

מדרך 1 קיבלנו:

```
True Positives (TP): 5503
False Positives (FP): 512
True Negatives (TN): 5369
False Negatives (FN): 616
recall\sensitivity: 0.899329955875143
precision: 0.914879467996675
accuracy: 0.906
```

והתוצוגת confusion matrix ב89

מדרך 2:

```
True Positives (TP): 5299
False Positives (FP): 716
True Negatives (TN): 5278
False Negatives (FN): 707
recall\sensitivity: 0.8822843822843823
precision: 0.8809642560266001
accuracy: 0.8814166666666666
```

והתוצוגת confusion matrix ב92

מדרך 3:

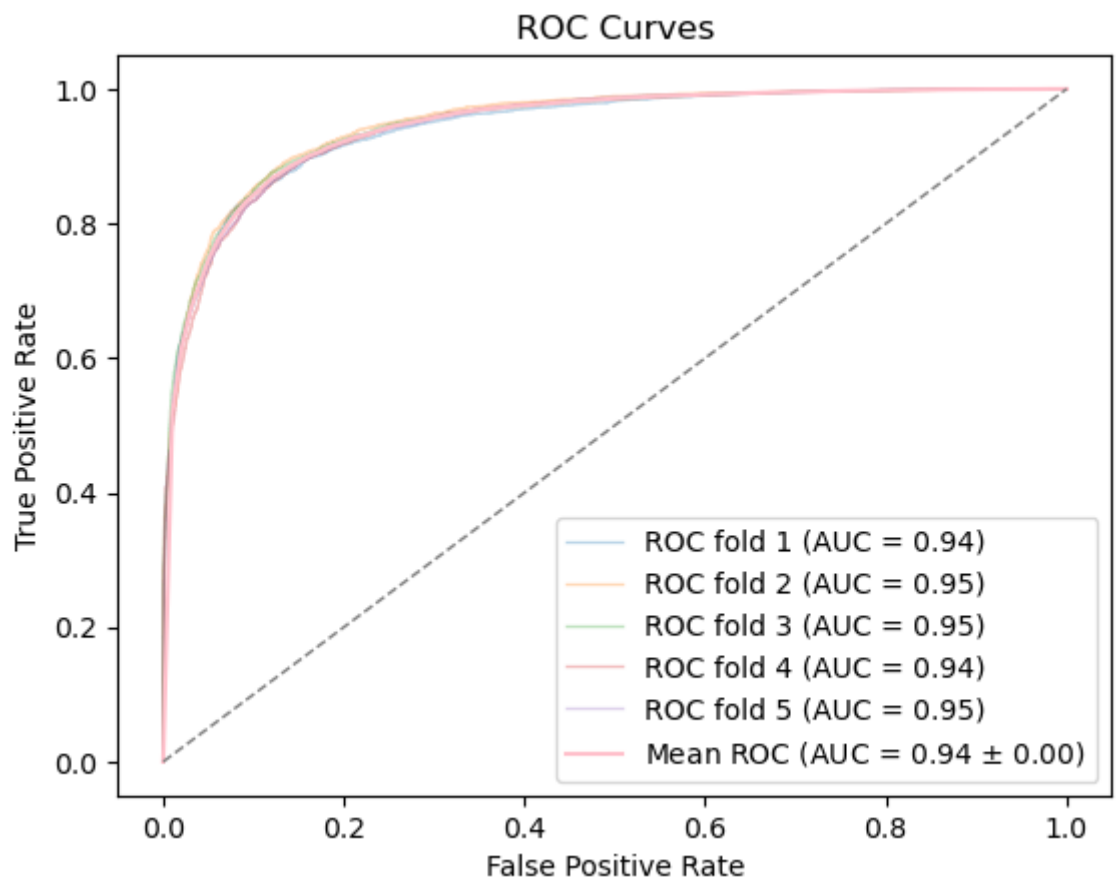
```
True Positives (TP): 5333
False Positives (FP): 682
True Negatives (TN): 5313
False Negatives (FN): 672
recall\sensitivity: 0.8880932556203164
precision: 0.886616791354946
accuracy: 0.8871666666666667
```

והתוצוגת confusion matrix ב95

(80) וקובץ ה CSV נראה ככה: (רואים פה רק חלק):

B	A	
proba	sha256	1
0.994732	023928c14	2
0.229575	6436083d4	3
0.485881	038e71f41	4
0.008172	be913ef29	5
0.37193	e78f83f0af	6
0.997935	c290de562	7
0.007651	15c686509	8
0.990066	633e2d2c1	9
0.16656	e80ed8f43	10
0.235311	a3071ca37	11
0.975499	1de60dd8a	12
0.320449	5349c2324	13
0.427857	3f9b01ac4	14
0.994283	e2c2d607c	15
0.172923	4967f1a3b	16
0.317981	540eb4eb6	17
0.325025	511312bd5	18
0.845333	7e46eac60	19

(81) הבדלים לדוגמה:



(82

