



מדעים דיגיטליים להיי-טק

1

תקציר מנהלים

תחילה קיבלנו שני קבצי דאטה, כאשר אחד הוא אימון והשני מבחן. ביצענו על שניהם עיבוד מקדים, כאשר בעיקר חקרנו את קובץ האימון ויישמנו על קובץ המבחן. לאחר מכן, אימנו את קובץ האימון בארבעה סוגי מודלים שונים כאשר פיצלנו את דאטה האימון ל-`validation` ו-`train`. לבסוף, בחרנו את המודל שהשיג את הביצועים הטובים ביותר לשם הפרדיקציה.

הסתכלות ראשונית על המידע

ראשית, ייבאנו את כל החבילות הדרושות עבור ביצוע המודלים והפרויקט. לאחר מכן, הסתכלנו ממבט-על על המידע שקיבלנו, בדקנו קורלציות בין כל המשתנים (**ראו נספח 1**) - מתוך כך ראינו כי משתנה `PageValues` יחסית מתואם עם `purchase`. בנוסף, ראינו כי עמודה `D` מתואמת בקשר שלילי חזק מאוד עם משתנה `purchase`. ניתן לראות גם כי משתני ה-`duration` pages וה-`num of pages` למיניהם די מתואמים אחד עם השני.

רצינו לראות איך העמודות המספריות מתפלגות, ושמנו לב כי יש הרבה יותר אי-רכישות מאשר רכישות (**ראו נספח 2**). ראינו גם כי יש הרבה משתנים אשר מתפלגים דומה להתפלגות חי-בריבוע כלומר בתחילת ההתפלגות ישנם הרבה ערכים שקרובים לאפס וככל שהערך עולה מספר התצפיות יורד.

ארגון המידע

התמודדות עם עמודות מספריות

1. תחילה, בנינו פונקציה אחת שמותאמת לכל עמודות `page` למיניהם בשביל מילוי ערכים חסרים על פי השלבים הבאים:

- א. בעמודות ה-`duration` החסרות מילאנו את עמודות ה-`num` שמתואמות להם באפס. לדוגמא אם ישנה שורה של ערך חסר בעמודה '`product_page_duration`', אז בעמודה '`num_of_product_pages`' מילאנו אפס.
- ב. הסרנו מהערכים שבעמודות ה-`duration` את המילה '`minutes`' כדי להפוך אותה לעמודה מספרית.
- ג. הנחנו כי אם בעמודת `num` יש ערך אפס, משמע המשתמש לא היה באף עמוד, אז בעמודת `duration` שמתאימה לה יהיה ערך אפס.
- ד. הנחה נוספת שהנחנו היא שסכום עמודות ה-`duration` של שלל סוגי הדפים, צריך להיות שווה ל-`total_duration`.
- ה. בעקבות הנחה 4, ובנוסף שהחלטנו למלא את עמודת '`product_page_duration`' אחרונה, מילאנו את הערכים החסרים שלה כהפרש בין `total_duration` ל-'`info_page_duration`' ו-'`admin_page_duration`'.
- ו. כעת, מילאנו ערכים חסרים בעמודות ה-`duration` לפי חציון ה-`duration` שמותאם ל-`num` שלו. לדוגמא, עבור תצפית שעמודת `num_of_product = 5`, ועמודת '`product_page_duration`' בעלת ערך חסר. ערך ה- '`product_page_duration`' תמולא על פי החציון של '`product_page_duration`' כאשר `num_of_product = 5`.
- ז. יצרנו רשימה של חציוני `duration`, אשר המיקום ברשימה הוא גם הערך שיהיה בעמודת ה-`num` של העמוד. לדוגמא, במקום האפס ברשימה נמצא החציון של עמודת ה-`duration` כאשר `num=0`, במקום האחד ברשימה נמצא החציון של עמודת ה-`duration` כאשר `num=1` וכך הלאה.
- ח. הנחנו כי כאשר עמודת ה-`duration` שווה לאפס, המשתמש לא גלש בעמודים מסוג זה ולכן מספר העמודים יהיה אפס.
- ט. בהמשך מילאנו את הערכים החסרים שבעמודות ה-`num` בעזרת הרשימה שיצרנו בסעיף 7. ראשית, בדקנו איפה הערך הספציפי שאנחנו עובדים איתו נמצא ברשימת החציונים. כלומר אם אנחנו עובדים עם ערך `duration` ששווה 4.7, אז נצמצם את הבדיקה שלנו לחציוני `duration` בין 4-5. שנית, חישבנו את הערך הממוצע של החציוני `duration` המתאימים ובדקנו האם ערך ה-`duration` הספציפי שלנו

מעל/מתחת הממוצע. אם הוא מעל, שייכנו אותו לחציון החמט הגבוה יותר ואם הוא מתחת אז שייכנו אותו לחציון החמט הנמוך יותר.

למשל, בדוגמא עם ערך duration ששווה 4.7, ערך זה ממוקם בין החציונים שבמיקומים 3-4, ומאחר ו-4.7 גדול מ-4.55 (ממוצע של 3.1 ו-6) נשייך ערך זה לחציון שבמיקום 4, כלומר נמלא את הערך החסר ב-4 (האינדקס המותאם לו).

Num of pages/index	0	1	2	3	4	5	6	7	8
Median duration	0	2	2.3	3.1	6	9.5	11.2	20.4	40.7

י. מקרה קצה שהבחנו בו היה כי כאשר ערך הduration הספציפי שלנו קטן מכל החציונים, בפרט ערך זה קטן מהחציון של כל התצפיות כאשר num=1 ולכן נמלא ערך num זה באחד. נזכיר כי אנחנו מניחים שרק אם ה-duration אז נמלא את החמט באפס.

יא. מקרה קצה נוסף- לאחר סיום מילוי הערכים, מי שנשאר לו ערך חסר אז הduration שלו היה גבוה מידי כדי להיות בין חציוני הרשימה ולכן מילאנו את החמט שלו בחמט המקסימלי שברשימה.

2. כעת, לאחר שסיימנו להריץ את הפונקציה ומילאנו את הערכים החסרים, גילינו כי יש 3 מקרי קצה ספציפיים מאוד ומילאנו אותם לפי ממוצע הערכים הקרובים אליהם.

3. בעמודת total_duration מילאנו את הערכים החסרים לפי הנחה בסעיף 4 בסכום הבא: 'admin_page_duration'+ 'product_page_duration'+ 'info_page_duration'.

4. בעמודות 'B', 'device', 'PageValues', 'ExitRates', 'BounceRates' מילאנו את הערכים החסרים לפי החציון שלהם.

5. שמנו לב כי בעמודה 'D' יש המון ערכים חסרים. רצינו להציג זאת ויזואלית (ראו נספח 3) ולכן יצרנו משתמש דמי שיגיד האם קיים/לא קיים ערך D בתצפית. קל לראות מהגרף שיצרנו כי יש יותר מידי ערכים חסרים. מהסיבה שכל מילוי אפשרי של העמודה יכול להפריע לחיזוי, בחרנו להסיר עמודה זו מהדאטה.

6. לגבי עמודה 'id', אנו מניחים כי אין קשר בינה לבין חיזוי הרכישה. אנו מניחים זאת מהסיבה שעמודה זו ממוספרת לפי סדר כרונולוגי של מספר השורות, כלומר אין לה קשר לרכישה/לא רכישה ולכן בחרנו להסיר אותה מהדאטה.

התמודדות עם עמודות קטגוריאליות

קודם כל, ביצענו כמה גרפים (ראו נספחים 4,5) כדי לזהות את הקטגוריות הנפוצות בחלק מהעמודות. לאחר מכן, יצרנו פונקציה שמותאמת לכל העמודות הקטגוריאליות שבה היא מחליפה את ערכיהם למספרים על פי השכיחות שלהם ואז ממלאה את הערכים החסרים בחציון העמודה. כאשר הערך הכי נפוץ מקבל את הערך אחד, אחריו שתיים והלאה. לדוגמא בעמודת ה'internet_browser' ראינו כי הדפדפן chrome_89 הכי נפוץ, ולכן הערך החדש שלו יהיה אחד.

שמנו לב כי בעמודת 'closeness_to_holiday', העמודה מספרית אך לא רציפה ולכן התייחסנו אליה כאל עמודה קטגוריאלית. רצינו למלא את הערכים החסרים לפי החציון של העמודה בהתאם לחודש של התצפית, כי הנחנו שהקרבה לחג תלויה בחודש שבו אנחנו נמצאים. כאשר בדקנו את החציונים של כל חודש, גילינו כי חציון העמודה בכל חודש הוא אפס, ולכן מילאנו את הערכים החסרים באפס (לא עשינו זאת על פי ממוצע משום שהיינו מקבלים ערכים לא שלמים).

בעמודת 'Region' הבנו כי למספרים יש משמעות, לכן לא היינו צריכים לעשות להם טרנספורמציה ומילאנו את ערכיהם בחציון העמודה.

בעמודת 'Weekend' ראינו כי עמודה זו היא מסוג בוליאני, אז החלפנו את ערכיה לאפס ואחד כאשר אפס מייצג False ואחד מייצג True ומילאנו את הערכים החסרים עם חציון העמודה.

התמודדות עם ערכים קיצוניים

הבחנו כי עמודה 'B' היא היחידה שמתפלגת דומה להתפלגות נורמלית (ראו **נספח 8**), ולכן בעמודה זו בחרנו להשתמש בשיטת IQR אשר מורידה ערכים קיצוניים משני קצוות ההתפלגות.

בנוסף, ניתן לראות מהגרפים (ראו **נספחים 9,10**) והשונויות שיצרנו, כי לעמודות 'ExitRates' ו-'BounceRates' יש שונות נמוכה, והרבה מידע בקצוות התפלגות הנתונים. השערותנו היא במידה ונוריד מהן ערכים קיצוניים תפגע יכולת החיזוי של המודלים.

בדיקה ראשונית להשערת העמודות החשובות לחיזוי

באמצעות **נספחים 11-16** הצגנו את ההשערות שלנו לחיזוי העמודות החשובות בהן נרצה לשים דגש.

בנספח 11 הבחנו בהתפלגות הדומה של שתי העמודות: 'product_page_duration' ו-'total_duration'. **בנספח 12** מוצגות כמות הרכישות אל מול אי הרכישות בכל חודש. נקודה מעניינת היא כי החודש עם מספר הרכישות הגבוה ביותר, אינו החודש עם אי ביצוע הרכישות הגבוה ביותר. **בנספח 13**, ניתן לראות את הקשר בין 'ExitRates' ל-'BounceRates' אל מול רכישות. כאשר ערכי שתי העמודות הללו קטן מ-0.035 יש יותר סיכוי שתבוצע רכישה. **בנספח 14**, מתואר הקשר בין 'Weekend' אל מול רכישות, ניתן לראות כי בזמן סופ"ש יש יותר סיכוי שתבוצע רכישה. **בנספח 15**, מתואר הקשר בין 'PageValues' ל-'total_duration' אל מול רכישות. כאשר ערכי שתי העמודות עולה (בעיקר PageValues) אז יש יותר סיכוי שתבוצע רכישה. **בנספח 16**, בדקנו את התפלגות עמודה 'B' בגרף מסוג היסטוגרמה, שאיפשר להסיק כי ההתפלגות של העמודה נורמלית בקירוב.

הורדת מימדים

בחרנו להשתמש בשתי שיטות להורדת מימדים. האחת, שיטת PCA – שאינה תלויה בפרמטרים. נרמלנו את הנתונים ולאחר מכן ביצענו את השיטה פעם אחת לפני ביצוע המודלים. השנייה, שיטת **Forward Selection** שכן תלויה בפרמטרים.

לאחר מכן, בנינו פונקציה כללית אשר תכריע איזה שיטה הפחתת מימדים טובה יותר עבור מודל-PCA/Forward Selection תוך התחשבות במדד MSE.

ROC K-Fold

פונקציה זו מקבלת מערך דו מימדי עם כל הפיצ'רים, מערך חד מימדי של ה-labels ואת המסווג שתלוי במודל. תחילה חילקנו את הדאטה ל-10 חלקים. אימנו את הדאטה בלולאה ויצרנו Roc Curve לכל פיצול דאטה. חישבנו את הממוצע של הציונים ואת סטיית התקן שלהם והצגנו זאת בגרף.

מודלים

רגרסיה לוגיסטית

תחילה, ביצענו את שיטת הורדת המימדים-Forward Selection, (ראו **נספח 17**). בהמשך, בדקנו איזה שיטת הורדת מימדים תביא ל-MSE קטנה יותר. ראינו כי ה-MSE בשיטת ה-Forward Selection הייתה נמוכה יותר. אחרי כן, פיצלנו את הדאטה לtrain ולvalidation כך שהvalidation מהווה 20% מהדאטה, ויצרנו פונקציה שבחרת את ה-C האופטימלי ברגרסיה (ראו **נספח 18**), שמייצגת את אחד חלקי כוח ההענשה. לאחר מכן, יצרנו מילון של היפר-פרמטרים ובעזרת GridSearchCV לקחנו את ההיפר-פרמטרים האופטימליים.

ההיפר-פרמטרים שבחרנו לבדוק הם: **(1) Penalty** - פרמטר זה אחראי על סוג ההענשה של הרגרסיה. **(2) solver** - פרמטר זה אחראי על סוג האלגוריתם שדרכו אנחנו מגיעים למשקולות האופטימליות. **(3) Max_iter** - פרמטר זה אחראי על מקסימום מספר האיטרציות. **(4) Fit_intercept** - פרמטר זה בודק האם חותך המשוואה שייך אליה או שזה רעש.

הפרמטרים שנבחרו ע"י GridSearchCV (בהתאמה): (1) סוג ההענשה שנבחר הוא l1 כלומר סוג הענשה לפי שיטת lasso. (2) סוג האלגוריתם שדרכו אנחנו מגיעים למשקולות האופטימליות הוא מסוג 'saga'. (3) מקסימום האיטרציות = 100 (4) False = Fit_intercept.

$C = 0.1$ (היישום בוצע ע"י הפונקציה שבחרת את ה-C האופטימלי ברגרסיה).

לאחר מכן, בנינו פונקציה אשר מציגה את ה confusion matrix (ראו נספח 19,20).

1. ערכי מטריצת ה **train** הם: $TN=6874$, $TP=496$, $FP=144$, $FN=808$. כלומר, צדקנו כאשר סיווגנו לאי-רכישה ב 6874 תצפיות וצדקנו כאשר סיווגנו לרכישה ב 496 תצפיות ולכן **דיוק מטריצה האימון שלנו הוא 0.88**.
2. ערכי מטריצת ה **validation** הם: $TN=1745$, $TP=105$, $FP=34$, $FN=197$. כלומר, צדקנו כאשר סיווגנו לאי-רכישה ב 1745 תצפיות וצדקנו כאשר סיווגנו לרכישה ב 105 תצפיות ולכן **דיוק מטריצה ה-validation שלנו הוא 0.88**.

ניתן להסיק מה confusion matrix כי ברוב המקרים צדקנו כשחזינו אי-רכישה. בנוסף, יש לשים לב שהריבוע שמתייחס לכך שחזינו נכונה רכישה אינו כהה, משום שאנחנו מתייחסים במטריצה לגודל אבסולוטי ולא יחסי וזאת בהתאם לכך שאין הרבה רכישות בדאטה. מפני ש TP גדול משמעותית מ FP אנחנו מסיקים כי המודל שלנו פעל בצורה טובה.

ניתן לראות מנספחים 21,22 כי ממוצע ה ROC-KFOLD לרגרסיה הלוגיסטית הוא 0.89 על ה **train** ו 0.88 על ה **validation**.

מתוך ביצועי מודל אלו, אנו מסיקים כי המודל שלנו הוא לא **overfitted**, משום שה $train_score$ גדול מה $validation_score$. ניסינו להימנע מ **overfitting** כשהורדנו את מימדיות הדאטה, כלומר הקטנו את מספר הפיצ'רים תוך התחשבות במדד MSE ושימוש ב cp מallow בשיטת **forward selection**. יש לציין כי מספר הפיצ'רים שנבחרו במודל הוא 9 שזהו מספר די נמוך, ולכן השונות של המודל יחסית נמוכה. בנוסף השתמשנו ברגולריזציה ושיטת ענישה אשר גם אמורים לעזור בהפחתת ה **overfitting**.

מסווג נאיב ביים

אנו מניחים כי התצפיות הן בלתי תלויות אחת בשנייה. תחילה, ביצענו את שיטת הורדת המימדים - **Forward Selection**, (ראו נספח 23) ולאחר מכן בדקנו איזה שיטת הורדת מימדים תביא ל-MSE קטנה יותר. מצאנו כי ה-MSE בשיטת ה-**Forward Selection** הייתה נמוכה יותר. בהמשך פיצלנו את הדאטה ל **train** ול **validation** כך שה **validation** מהווה 20% מהדאטה.

ניתן להבחין מנספחים 25,24 כי ממוצע ה ROC-KFOLD לביים הוא 0.85 על ה **train** ו 0.83 על ה **validation**.

לדעתנו מודל זה לא **overfitted** משום שזהו מודל אשר חוזה לפי הסתברות, ובנוסף ה $train_score$ גדול מה $validation_score$. בנושא סיבוכיות זמן הריצה, הבחנו כי מודל זה מחשב את ההסתברות של כל פיצ'ר, לכן ככל שיש יותר פיצ'רים הסיבוכיות עולה. כדי להקטין את בעיית החישוביות של המודל הורדנו את מימדיות הדאטה, ובכך הקטנו את מספר הפיצ'רים תוך התחשבות במדד MSE ושימוש ב cp מallow בשיטת **forward selection**. יש לציין כי מספר הפיצ'רים שנבחרו במודל הוא 4 שזהו מספר נמוך, ולכן הסיבוכיות של המודל יחסית נמוכה.

Random Forest

בחרנו שלא להשתמש בשיטת הורדת מימדים משום שבמודל זה, ככל שיש יותר עצים וככל שהם יותר עמוקים המודל יחזה יותר טוב. בהתאם לכך, לא ראינו בעייתיות שיכולה להיווצר ממספר הפיצ'רים. פיצלנו את הדאטה ל **train** ול **validation** כך שה **validation** מהווה 20% מהדאטה. בהמשך, יצרנו מילון של היפר-פרמטרים ובעזרת GridSearchCV לקחנו את ההיפר-פרמטרים האופטימליים.

ההיפר-פרמטרים שבחרנו לבדוק הם: (1) **n_estimators** - מספר העצים ביער, (2) **criterion** - סוג הפונקציה שמודדת את השגיאה בחיזוי, (3) **Max_depth** - מקסימום עומק עץ, (4) **Min_samples_leaf** - מינימום תצפיות בעלה, (5) **Min_samples_split** - מספר תצפיות מינימלי שצריך לבצע עבורן עוד פיצול בעץ

הפרמטרים שנבחרו ע"י GridSearchCV (בהתאמה): (1) **n_estimators=200**, (2) **criterion='entropy'**, כלומר הוא מודד את המידע המעיד על אי-סדר הפיצ'רים עם הlabel. (3) **Max_depth=200**, (4) **Min_samples_leaf=5**, (5) **Min_samples_split=3**.

לאחר מכן, הצגנו בגרף את חשיבות הפיצ'רים (ראו נספח 26). ניתן לראות כי בפער משמעותי הפיצ'ר 'PageValues' החשוב ביותר לביצוע הפרדיקציה כפי ששיערנו (נספח 15). ממצא זה הגיוני משום שה'PageValues' זוהי עמודה אשר מייצגת את ערך הדף, וכאשר ערך הדף גבוה גדל הסיכוי לביצוע רכישה. עמודה חשובה נוספת אשר חזינו זוהי עמודה 'product_page_duration' (נספח 11). ניתן להסיק כי היא חשובה משום שככל שמשך השהייה גדלה בעמוד מסוג זה, כך הלקוח יותר מעוניין במוצר והסיכוי לרכישה גדל (במידה והלקוח אהב את המוצר כמובן). עמודה זו מתואמת חזק עם עמודת 'total_duration', ולכן גם 'total_duration' עמודה די חשובה.

ניתן לראות מנספחים 28,27 כי ממוצע הROC-KFOLD לבייס הוא 0.93 על train ו-0.90 על validation.

מתוך ביצועי מודל אלו, אנו מסיקים כי המודל שלנו הוא לא overfitted, משום שהtrain_score גדול מהvalidation_score.

Neural Networks

במודל זה בחרנו שלא להשתמש בשיטת הורדת מימדים. PCA בחר לקחת 18 משתנים, וזה רק הבדל של 2 משתנים מהדאטה ללא הורדת מימדים ולכן השיטה לא יעילה. בחרנו שלא להשתמש בForward selection מטעמי סיבוכיות הריצה. לאחר מכן פיצלנו את הדאטה לtrain ולvalidation כך שהvalidation מהווה 20% מהדאטה. כעת, יצרנו מילון של היפר-פרמטרים ובעזרת GridSearchCV לקחנו את ההיפר-פרמטרים האופטימליים.

ההיפר-פרמטרים שבחרנו לבדוק הם: (1) **Hidden_layer_sizes** - משתנה זה קובע את מספר השכבות ומספר הפרמטרים בכל שכבה, (2) **Solver** - שיטת אופטימיזציה למשקלים, (3) **Activation** - סוג הפונקציה שנמצאת בשכבות.

הפרמטרים שנבחרו ע"י GridSearchCV (בהתאמה): (1) **Hidden_layer_sizes=3** - שכבות של 100 משתנים, (2) **Solver='adam'**, (3) **Activation='logistic'**.

ניתן לראות מנספחים 29,30 כי ממוצע הROC-KFOLD לבייס הוא 0.90 על train ו-0.89 על validation.

מתוך ביצועי מודל אלו, אנו מסיקים כי המודל שלנו הוא לא overfitted, משום שהtrain_score גדול מהvalidation_score.

סיכום

בפרויקט זה הרצנו ארבעה סוגי מודלים: רגרסיה לוגיסטית, Neural Networks, מסווג נאיב ביס, Random Forest. בהתאם לביצועים של כל מודל, לזמן הסיבוכיות ול-overfitting, בחרנו במודל Random Forest השיג את הביצועים הטובים ביותר לפי ROC K-FOLD. יחסית למורכבות המודל ולביצועים שלו, זמן הריצה תקין ולכן בחרנו בו בעבור הפרדיקציה.

נספחים

נספח אחראיות כל שותף ותרומתו לעבודה

תחילה, נדגיש כי את הפרויקט ביצענו בסביבת עבודה Visual Studio Code תוך כדי שיתוף פעולה מלא אשר התאפשר באפליקציית ה-Zoom, כך שחלקים נרחבים מכתובת הקוד עשינו בצורה משותפת ובזמן אמת. עם זאת, חלקים קטנים חולקו בינינו כדי למזער את העומס:

ניצן כהן:

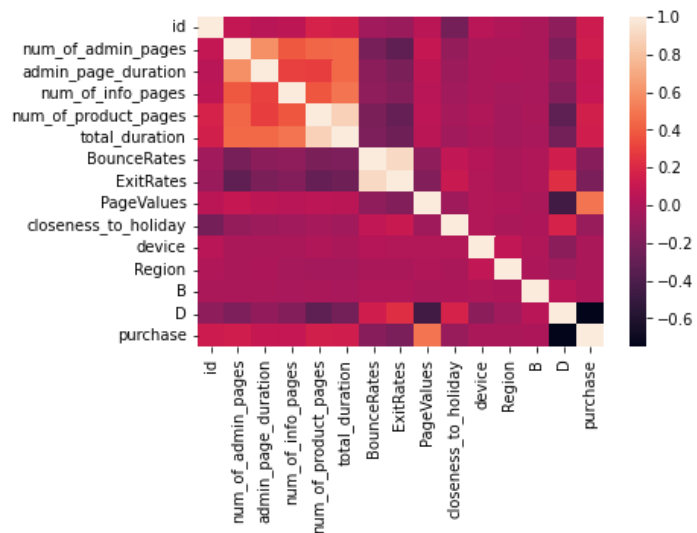
1. ניצן החליט על הצגת הוויזואליזציה שבחלק ביצוע המודלים.
2. ניצן רשם את ההערות שבחלק markdown בחלק ביצוע המודלים.

ינאי איתם ברדוש:

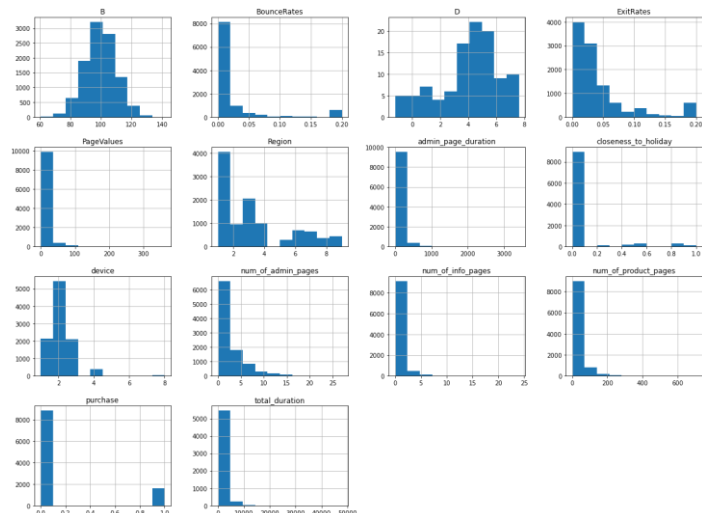
1. ינאי החליט על הצגת הוויזואליזציה שבחלק העיבוד המקדים.
2. ינאי רשם את ההערות שבחלק markdown בחלק העיבוד המקדים.

נספחי ויזואליזציה

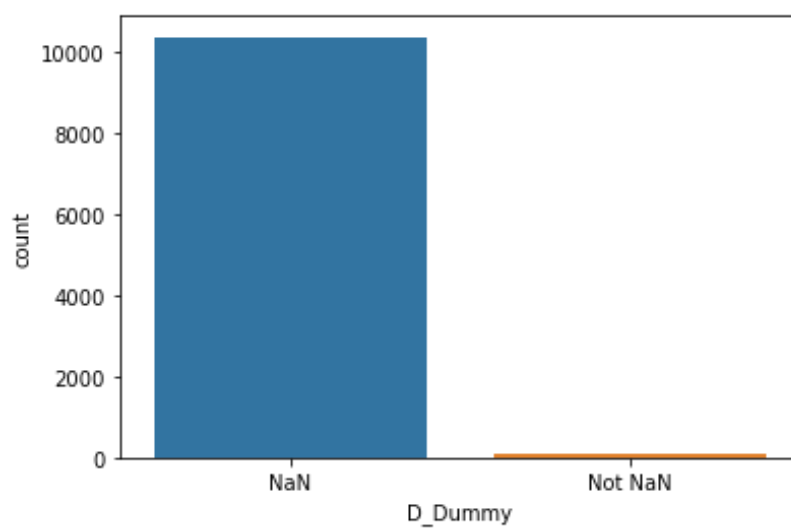
1.



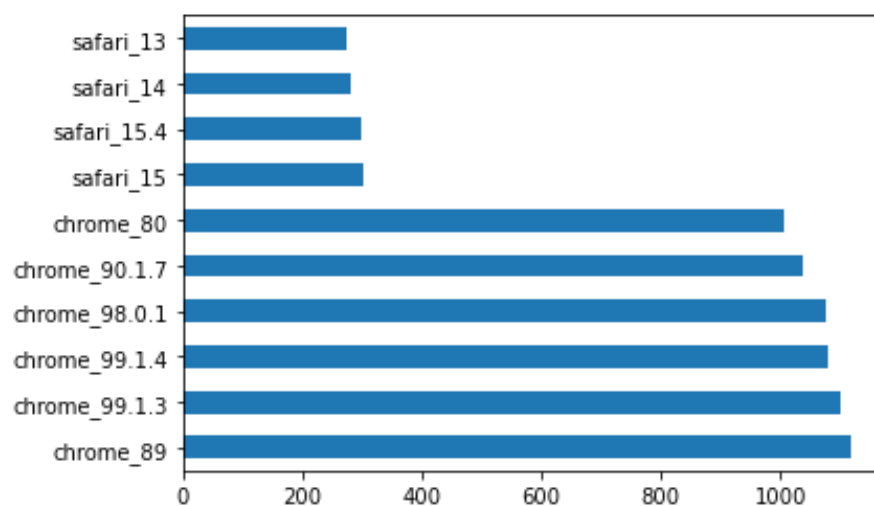
2.



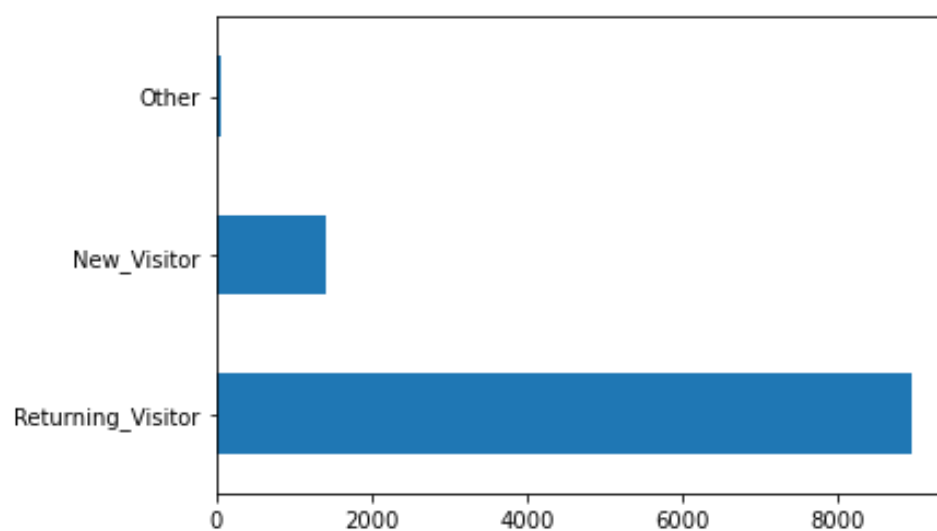
.3



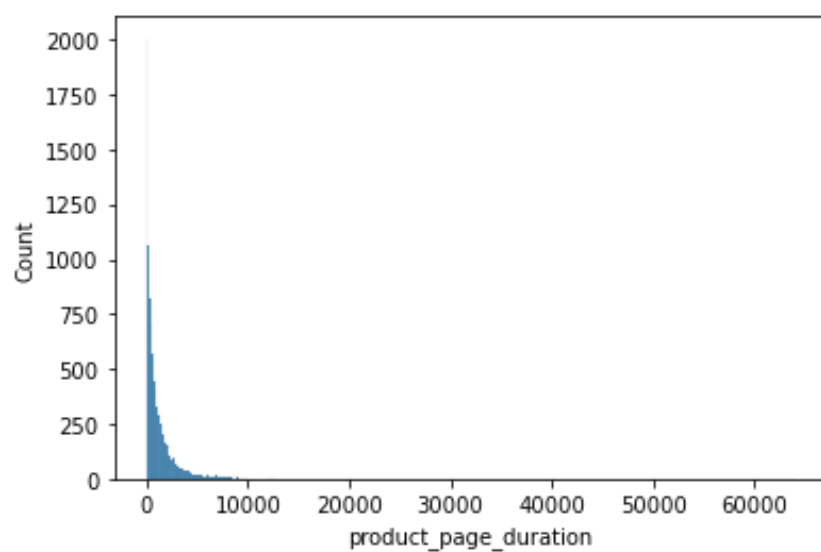
.4



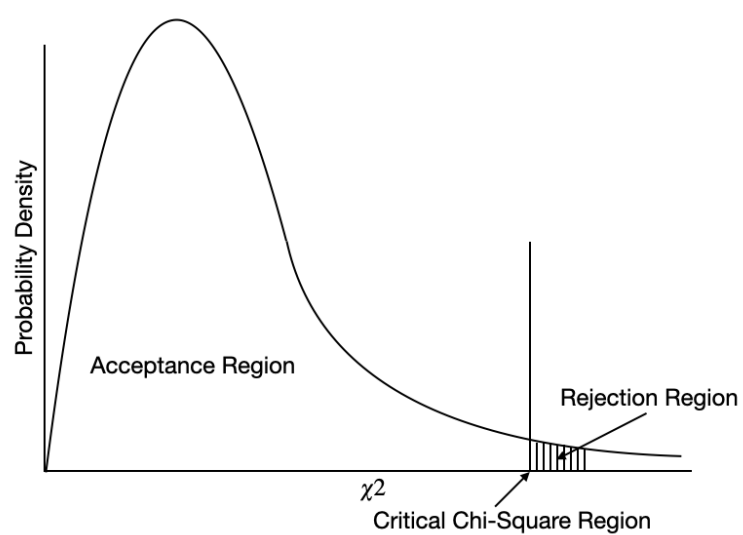
.5



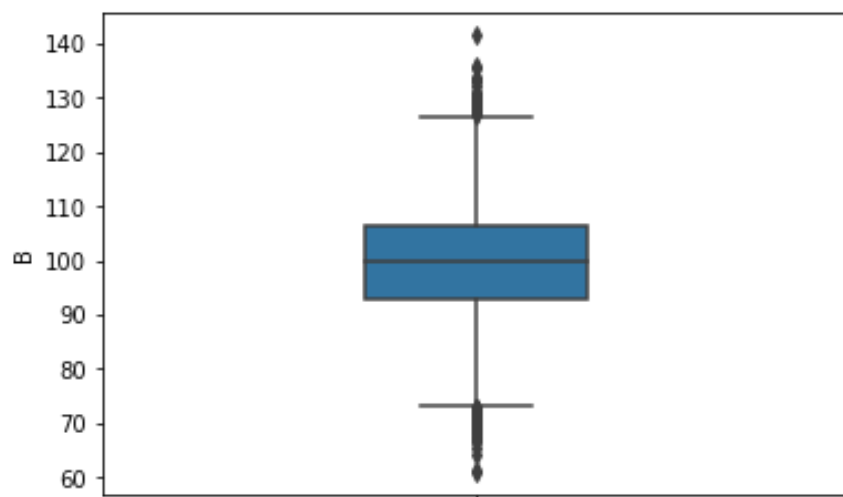
.6



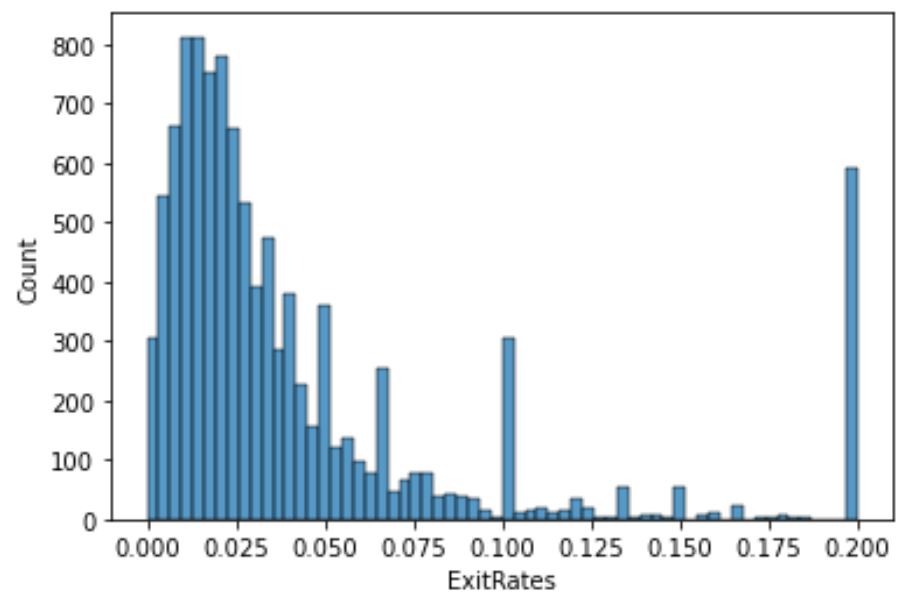
.7



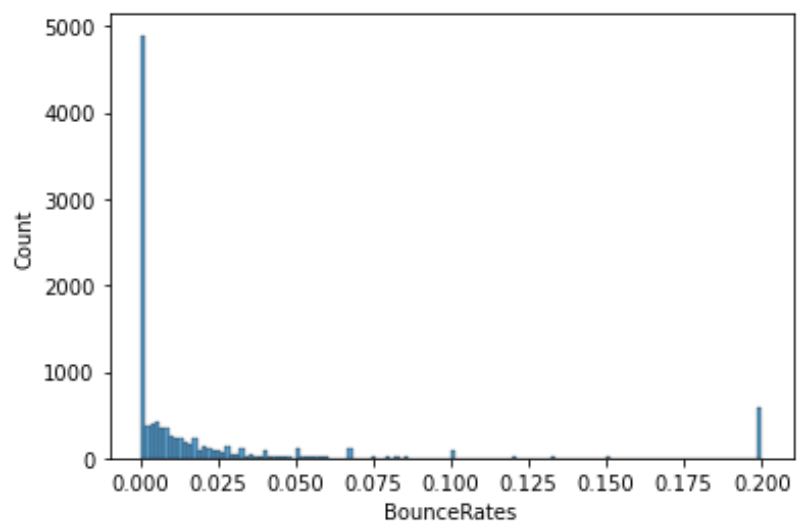
.8



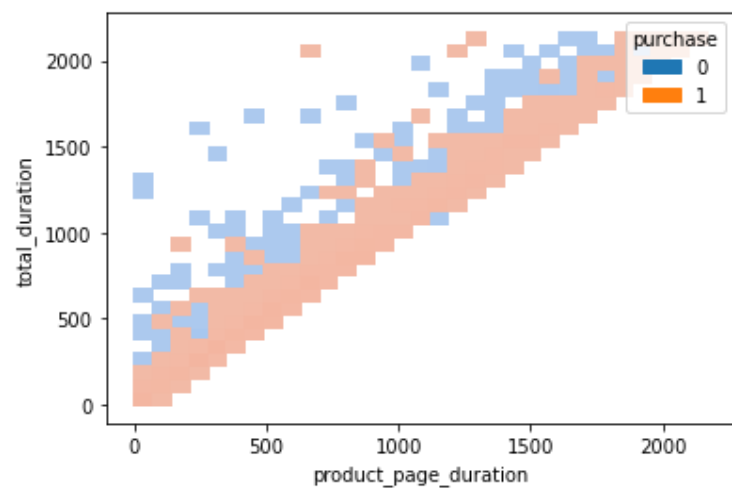
.9



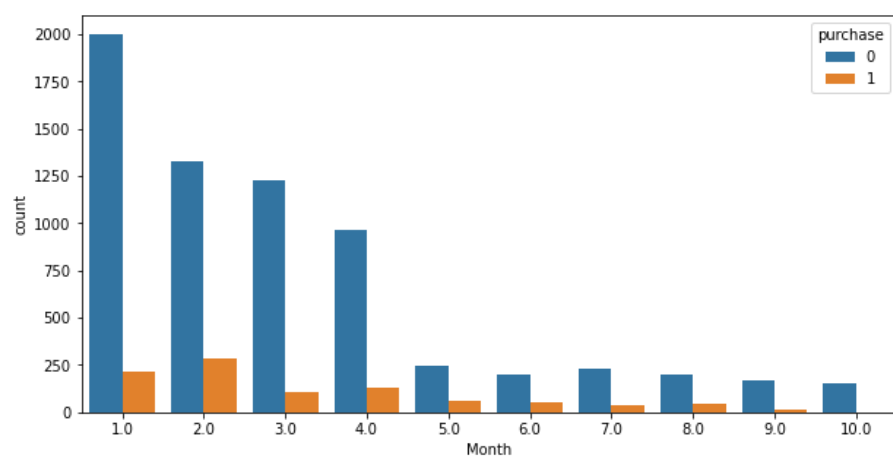
.10



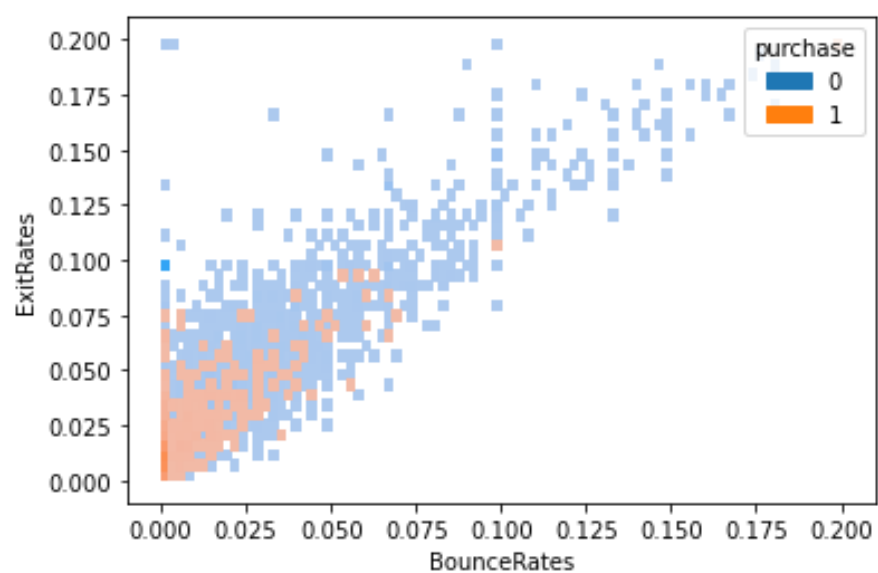
.11



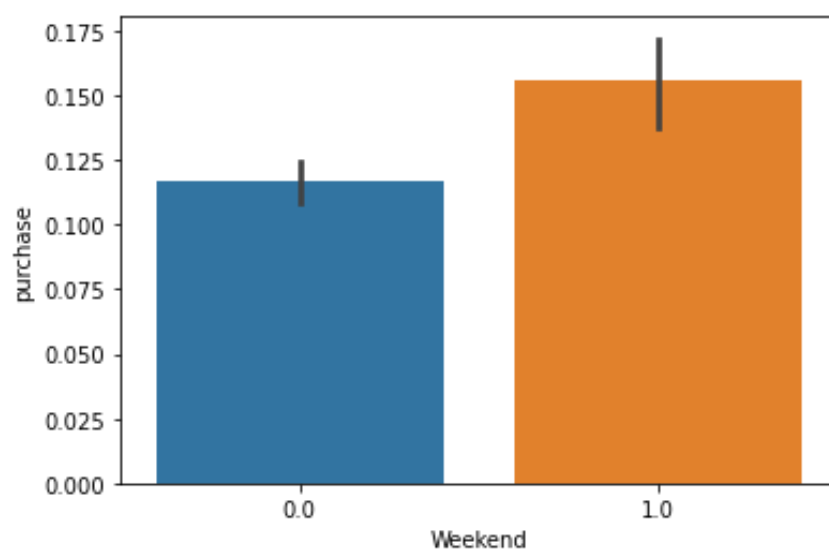
.12



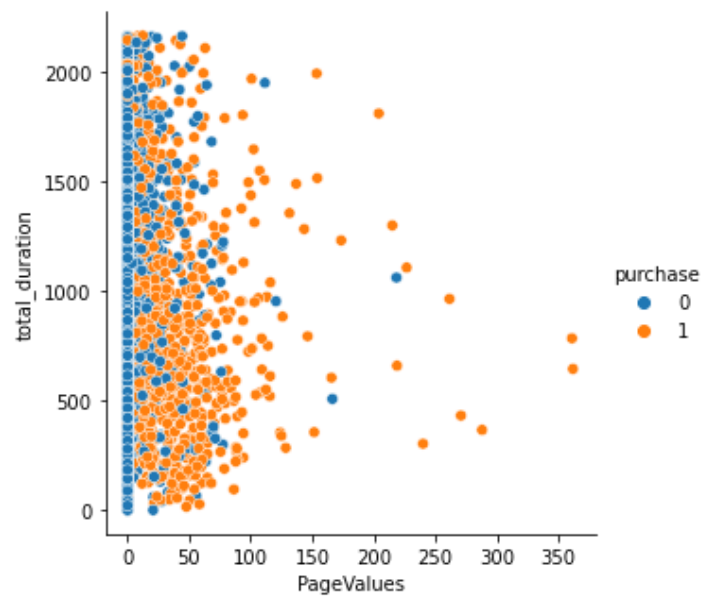
.13



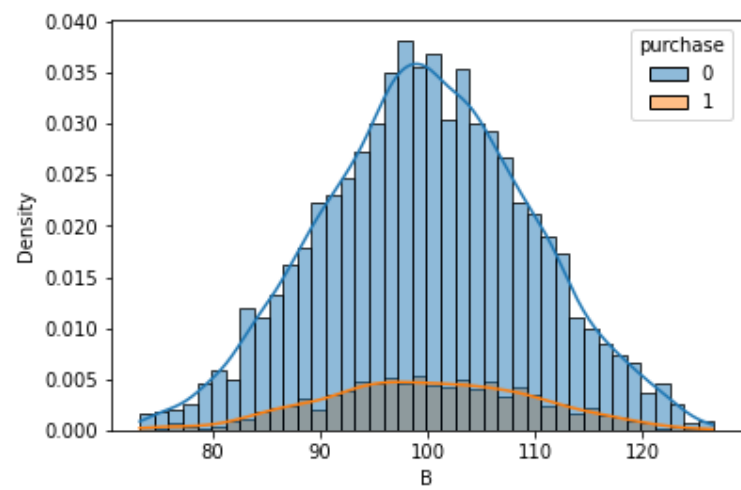
.14



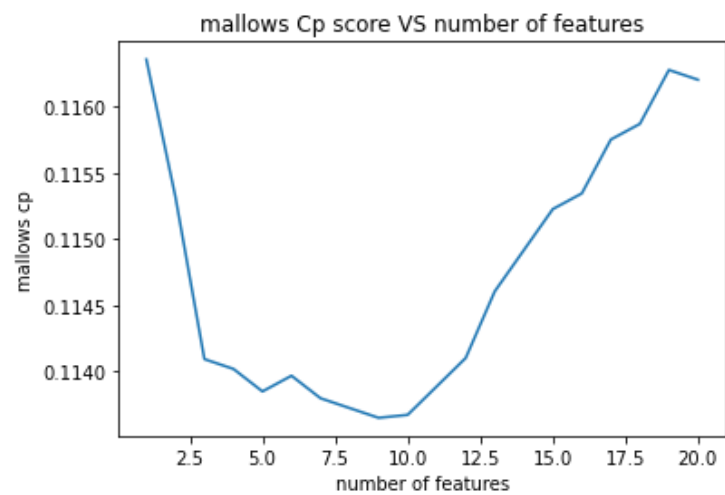
.15



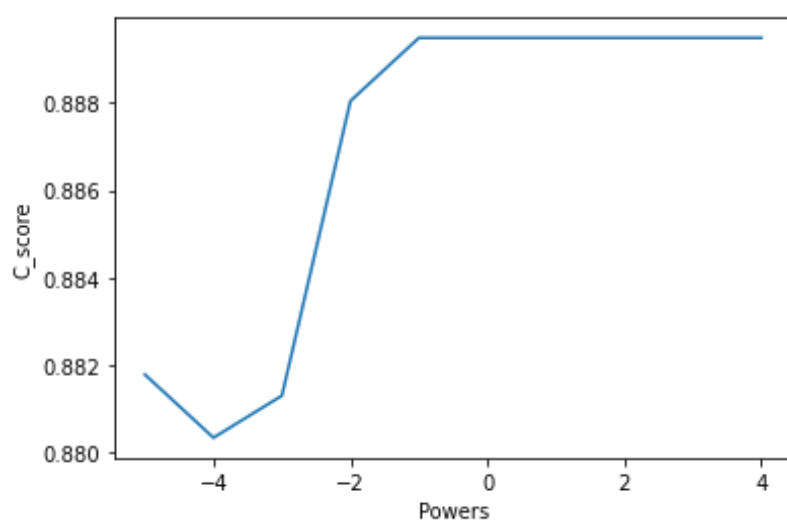
.16



.17



.18



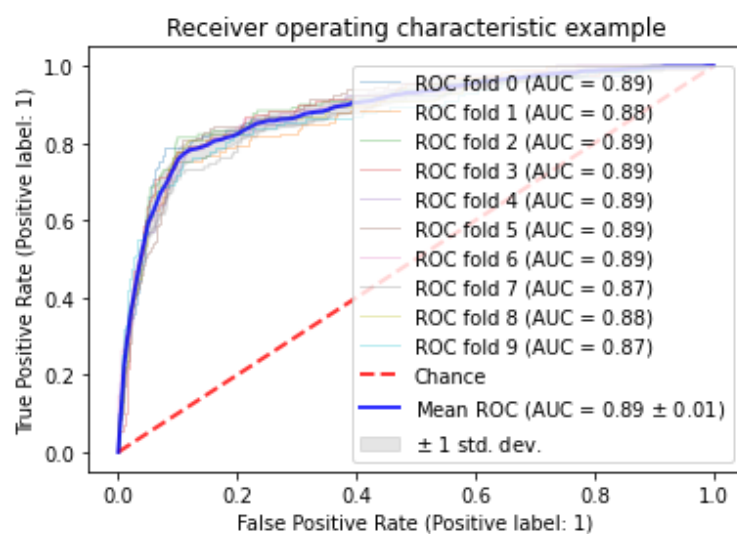
.19



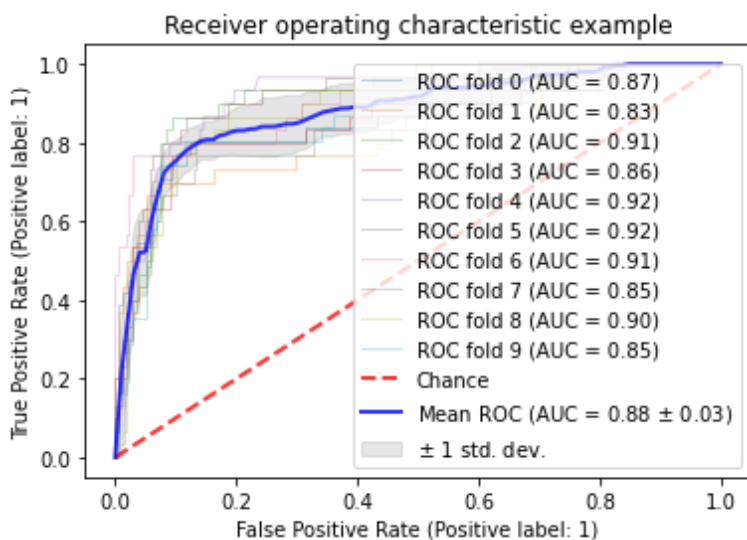
.20



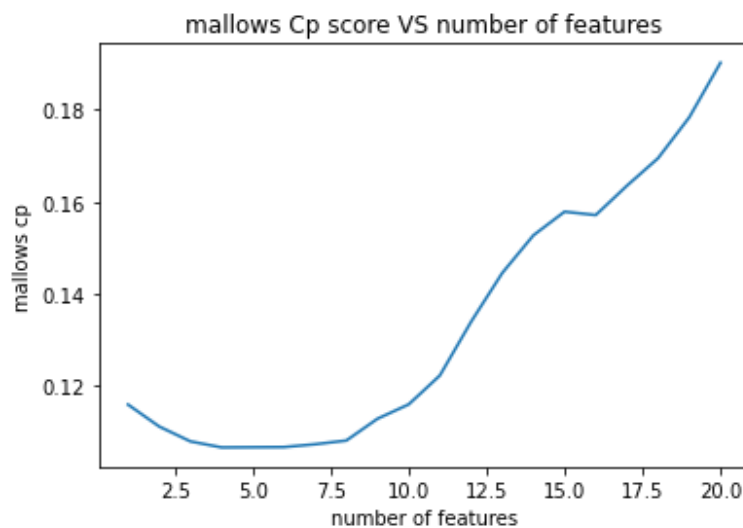
.21



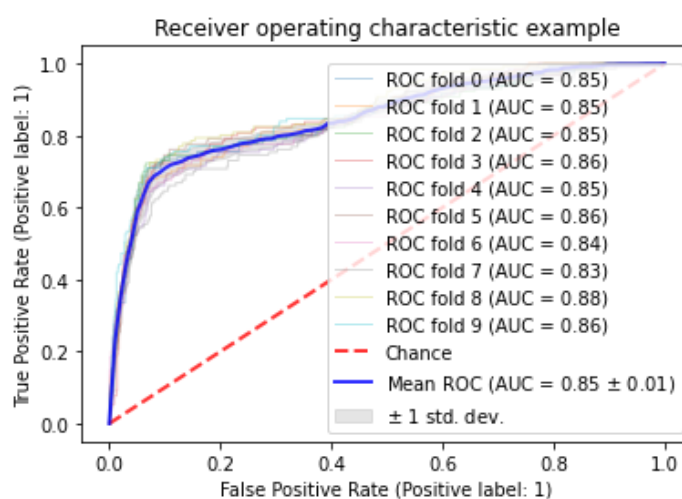
.22



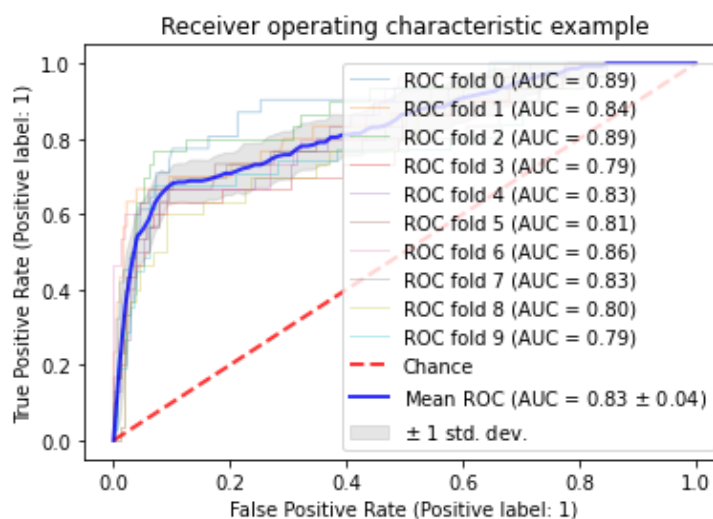
.23



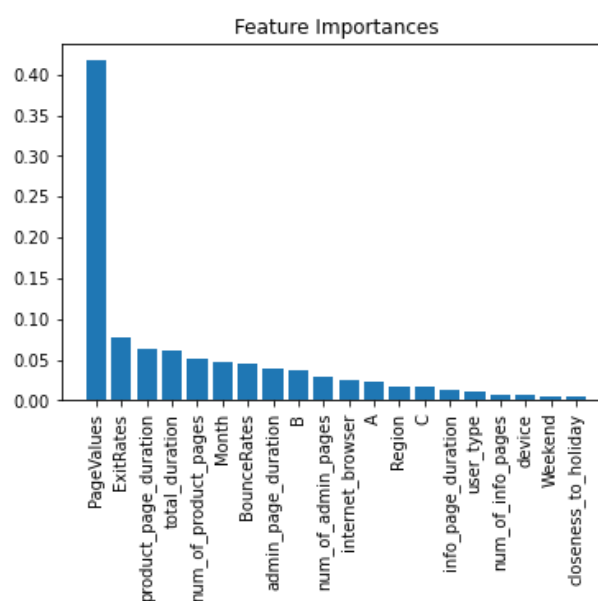
.24



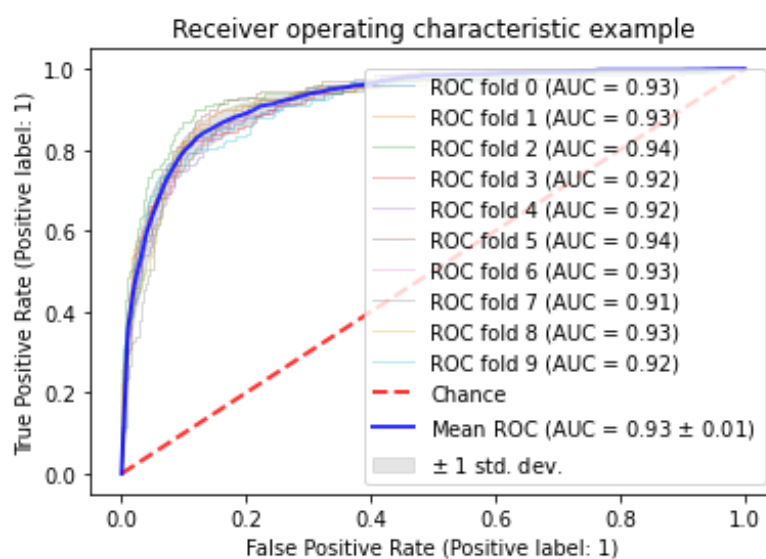
.25



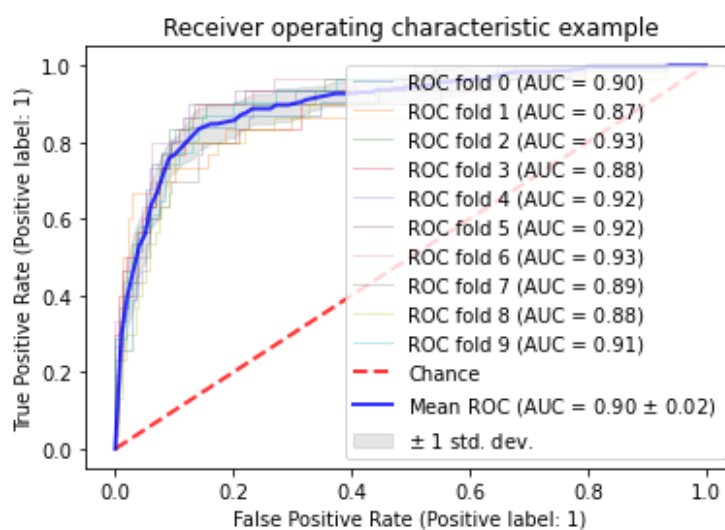
.26



.27



.28



.29

