

הנדון: ניתוח של מערך נתונים של יעילות שינה

הקדמה

מסמך זה מציג ניתוח מקיף של מערך נתונים של יעילות השינה, תוך התמקדות בטכניקות עיבוד מקדים ובמודלים של סיווג, המיושמים כדי לחזות את עמודות המגדר. מערך הנתונים מכיל מידע על דפוסי שינה של אנשים, כולל משך השינה, יעילות השינה, שעת השינה, זמן ההתעוררות ושלבי שינה שונים. מטרת מחקר זה הייתה לחקור את מערך הנתונים, לעבד מראש את הנתונים, להפעיל מודלים של סיווג ולהשוות בין התוצאות כדי לזהות את המודל היעיל ביותר ולייצר מסווג שיוכל לחזות את עמודות המגדר בנתונים חדשים המותאמים לנתונים הקיימים לנו במערך הנתונים באופן המדויק ביותר. מסמך זה מספק סקירה מפורטת של מערך הנתונים, מתאר את שלבי העיבוד המקדים ומודלי הסיווג שהופעלו, ומציג השוואה מקיפה של התוצאות.

מבוא

מערך הנתונים כולל מספר עמודות משמעותיות התורמות להבנת דפוסי יעילות השינה. אלו כוללים:

1. משך השינה- משך הזמן הכולל שכל נבדק ישן בשעות. הוא מספק תובנות לגבי משך השינה, שהוא חיוני להערכת איכות השינה.
2. שעת שינה וזמן השכמה- מאפיינים אלה מציינים את לוח השינה הממוצע של האדם, ונותנים מידע על הרגלי שינה והשגרה.
3. שלבי שינה- מערך הנתונים כולל מידע על הזמן המושקע בשלבי שינה שונים, כגון שנת REM, שינה עמוקה ושינה קלה. שלבים אלו מספקים תובנות לגבי איכות ועומק השינה שחווים בני אדם.

שינה מורכבת מארבעה שלבים-

- א. שנת חלום REM
 - ב. שינה עמוקה
 - ג. שינה קלה
 - ד. שינה שטחית
4. התעוררויות- מספר הפעמים שכל נבדק מתעורר במהלך הלילה. תכונה זו עוזרת להבין הפרעות והפרעות שינה שעשויות להשפיע על יעילות השינה.
 5. צריכת קפאין ואלכוהול- מערך הנתונים מתעד את צריכת הנבדק של קפאין ואלכוהול ב-24 השעות שלפני השינה. ידוע כי גורמים אלו משפיעים על איכות ויעילות השינה.
 6. סטטוס עישון- מידע האם הנבדק מעשן או לא, המאפשר לבחון קשרים פוטנציאליים בין עישון ויעילות שינה.
 7. תדירות פעילות גופנית- תכונה זו מתארת את תדירות הפעילות הגופנית המדווחת על ידי כל נבדק. פעילות גופנית עשויה להיקשר לאיכות שינה טובה יותר, ועמודה זו מאפשרת לחקור את השפעתה על יעילות השינה.

8. יעילות השינה- תכונת יעילות השינה היא מדד לשיעור זמן השהות במיטה שבפועל בשינה, עמודה זו באה לידי ביטוי ביחס ישיר לעמודות השונות בסט הנתונים ובאמצעותן (כלל העמודות למעט עמודת המגדר) ננסה לייצר מודל החוזה באופן הטוב ביותר את עמודת המגדר.
9. מגדר- עמודה המתארת אם מדובר בגבר או באישה- **התכונה שאותה נרצה לחזות.**

ניתוח נתונים חקרני מקיף- EDA

מטרת ה-EDA הייתה להשיג תובנות על מערך הנתונים, לזהות דפוסים ולהבין את הקשרים בין משתנים שונים. חלק זה מתאר את הצעדים שבוצעו במהלך ה-EDA ומספק ניתוח של הממצאים.

שלב 1: חקר נתונים ותצפיות ראשוניות

עם טעינת מערך הנתונים, העמודות נבדקו כדי להבין את התכונות הקיימות. מערך הנתונים כלל מידע רב התורם להבנת דפוסי יעילות השינה- פירוט העמודות מצורף לעיל. עמודות אלו מילאו תפקיד חיוני בהבנת דפוסי יעילות השינה.

שלב 2: ערכים חסרים

ערכים חסרים נבדקו כדי לקבוע את מידת החסר במערך הנתונים. המשתנה `missing_values` שימש כדי לחשב את סך הערכים החסרים, ו `missing_data`-נוצר כדי להציג את אחוז הערכים החסרים בכל עמודה. ניתוח זה סיפק תובנות לגבי שלמות מערך הנתונים והדגיש פערים פוטנציאליים שיש לטפל בהם במהלך העיבוד המקדים. גילינו שיש מעט ערכים חסרים (הדאטא מכילה 452 שורות על 16 עמודות), סה"כ 4 עמודות עם ערכים חסרים:

20 בעמודה המייצגת התעוררויות, 25 בעמודה המייצגת צריכת קפאין, 14 בעמודה המייצגת צריכת אלכוהול, 6 בעמודה המייצגת תדירות פעילות גופנית- החלטנו לטפל בערכים אלו בתהליך של ניקוי הנתונים ע"י השלמת ערכים ממוצעים שכן גילינו כי מדובר על שורות בעלות מרבית הערכים האחרים מלאים.

שלב 3: חלוקת הנתונים וחריגים

זוהו עמודות מספריות, והתפלגויותיהן הוצגו באמצעות תרשים קופסאות. זה איפשר הבנה טובה יותר של חלוקת הנתונים, הנטייה המרכזית והנוכחות של חריגים. תצפית ראויה לציון הייתה נוכחותם של ערכים חריגים בעמודות כגון "אחוז שינה קלה" ו"צריכת קפאין".

שלב 4: טיפול בחריגים

כדי להבטיח את שלמות הנתונים, הוחלפו חריגים בערכי אפס. גישה זו ננקטה כדי להימנע מהטיית תהליך הניתוח והמודלים עקב ערכים קיצוניים שכן ייתכן והוזנו בטעות. באופן ספציפי, ערכי חריגים הוחלפו ב-None בעמודות "אחוז שינה קלה" ו"צריכת קפאין". שלב זה הקל על ניתוח מדויק יותר של מערך הנתונים.

שלב 5: התפלגות מגדרית

התפלגות המגדר במערך הנתונים נבדקה באמצעות גרפי עמודות. היא סיפקה סקירה כללית של ייצוג הנבדקים גברים ונשים, ואיפשרה הבנה טובה יותר של הרכב המגדר של מערך הנתונים- גילינו כי הגרף שלנו אינו סובל מהתאמת יתר (Overfitting), שכן התפלגות הגברים ונשים שוויונית- 228 גברים, 224 נשים.

שלב 6: קשר בין משתנים

דיאגרמות פיזור הופעלו כדי לחקור את הקשרים בין משתנים שונים. יש לציין כי נוצרו דיאגרמות פיזור כדי לבחון את הקשר בין "משך השינה" ו"יעילות השינה", וכן בין "משך השינה" ל"מגדר". הדמיות אלו סיפקו

תובנות לגבי מתאמים ודפוסים אפשריים בין משך השינה, יעילות השינה ומגדר, בניתוח הגרפים ניתן לומר כי אינו קיים מתאם בין נתונים אלו.

שלב 7: ניתוח ספציפי למגדר

ניתוח נוסף נערך על ידי הפרדת מערך הנתונים על סמך מגדר. עמודות הגרף שימשו להשוואת התפלגות משך השינה ויעילות השינה בין גברים ונשים. ניתוח ספציפי למגדר איפשר הבנה מעמיקה יותר של ההבדלים בדפוס השינה וביעילות המבוססת על מגדר.

שלב 8: ניתוח רב משתנים

כדי לחקור קשרים מורכבים בין משתנים, יושמו טכניקות ניתוח רב משתנים. לדוגמה, נוצרה מטריצת פיזור כדי להמחיש את הקשר בין יעילות שינה, מצב עישון ואחוז שנת REM. דיאגרמת קו נוספת נוצרה כדי לבחון את הקשר בין גיל, יעילות שינה ומגדר. הדמיות אלו סיפקו תובנות לגבי משחק הגומלין של משתנים מרובים והשפעתם על יעילות השינה. לדוגמה ניתן לראות בדיאגרמת הקו שלאורך חיי אדם (הן גבר והן אישה) ישנה יעילות שינה זהה באופן יחסי הן אצל גברים והן אצל נשים. כמו כן, ע"פ נתונים אלו נראה שאצל גברים בשנות ה-20 המוקדמות יעילות השינה נמוכה מנשים ואילו בשנות ה-30 המוקדמות יעילות השינה של נשים נמוכה יחסית לשל גברים (גם בשנות ה-50 המאוחרות אצל גברים יעילות השינה נמוכה יותר).

שלב 9: השפעת עישון ואלכוהול

דיאגרמת קופסאות נוצלו כדי לנתח את ההשפעה של מצב העישון וצריכת אלכוהול על יעילות השינה. ניתוח זה נועד לזהות קשרים פוטנציאליים בין גורמי אורח חיים אלו ואיכות השינה. בדיאגרמת השינה המתארת קשר בין צריכת אלכוהול ליעילות שעות שינה, אכן ניתן לזהות שאי צריכת אלכוהול מייצרת את יעילות השינה הגבוהה ביותר, אמנם ניתן לומר ע"פ נתונים אלו שהקשר ביניהם אינו לינארי דבר המקשה עלינו להגיע למסקנה מובהקת בנושא הקשר בין המשתנים. ע"פ התרשים המקשר בין יעילות השינה לסטאטוס מעשן ולא מעשן, ניתן לומר שיעילות השינה אצל לא מעשנים אכן טובה יותר.

סיכום EDA-

תהליך ה-EDA שנערך על מערך הנתונים של יעילות השינה חשף תובנות חשובות לגבי דפוס השינה, יעילות השינה והגורמים המשפיעים עליהם. באמצעות הדמיות וניתוח סטטיסטי, זוהו דפוסים ונחקרו קשרים בין משתנים. החלפת חריגים בערכים אפסיים הבטיחה את שלמות הנתונים במהלך הניתוח והמודלים הבאים. הניתוח הספציפי למגדר שופך אור על הבדלי דפוס השינה בפן המגדרי.

מסקנות לגבי טכניקות עיבוד מקדים ומודל חלוקה לאשכולות

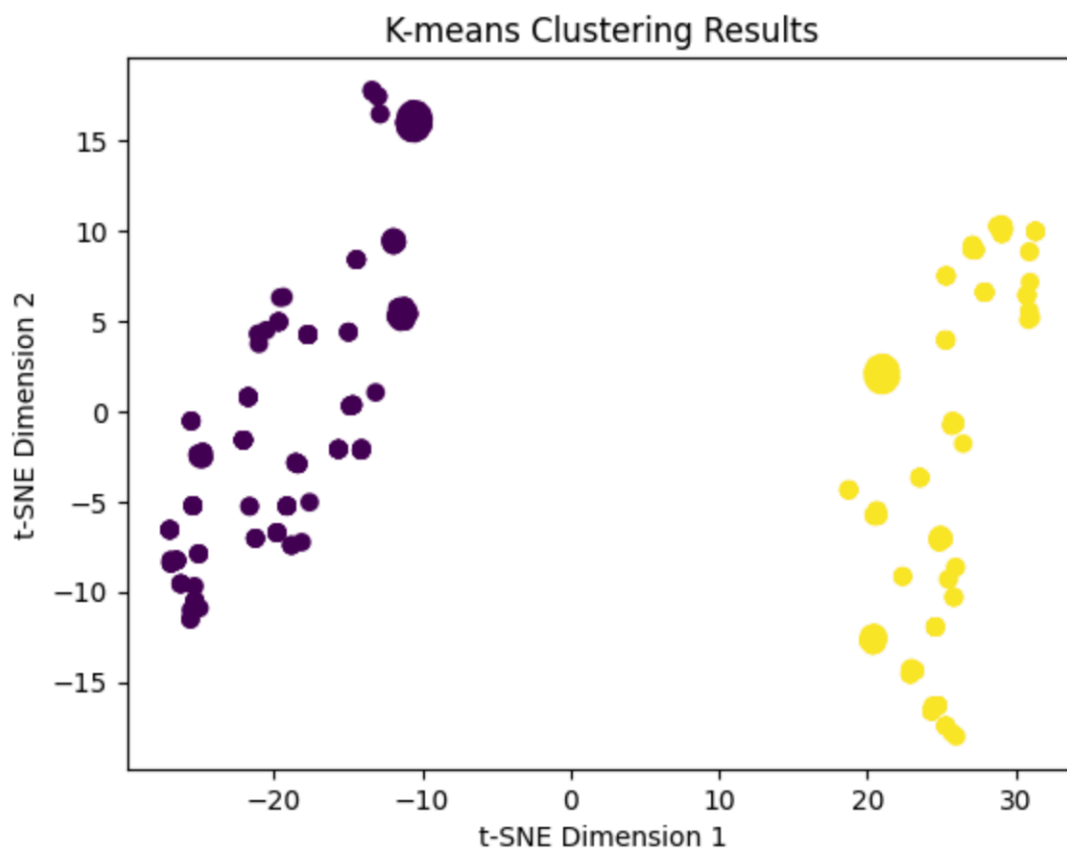
בפרויקט, עשינו שימוש בטכניקות שונות של עיבוד מקדים של נתונים ומודלים של למידת מכונה כדי לנתח את מערך הנתונים. **המטרה העיקרית** הייתה לקבוע את **טכניקת העיבוד המקדים והשילוב המודל הטוב ביותר** שמניב את מדדי הביצועים הגבוהים ביותר.

טכניקות עיבוד מקדים של נתונים כללו קנה מידה מינימלי-מקסימלי, קנה מידה של ציון z, קנה מידה עשירי ונרמול לפי טכניקות וקטור יחידה. טכניקות אלו יושמו כדי לנרמל את התכונות המספריות של מערך הנתונים לפני אימון המודלים.

כדי להעריך את הביצועים, נעשה שימוש במספר מדדים. מדדים אלו סיפקו לנו תובנות לגבי הדיוק הכולל, דיוק החיזוי החיובי, שיעור החיובי האמיתי ואיזון בין דיוק לזכירה.

כמו כן, בדקנו שילובים שונים של טכניקות ומודלים של עיבוד מקדים, והתוצאות נותחו בקפידה כדי לקבל החלטות מושכלות לגבי השלבים הבאים.

כחלק מתהליך ניסיון שיפור המודל, החלטנו להוסיף סיווג לפי אשכולות של מגדר באמצעות מודל האשכולות K-means.



הגרף מציג ייצוג דו מימדי של הנתונים באמצעות אלגוריתם t-SNE (t-Distributed Stochastic Neighbor Embedding), שהיא טכניקת הפחתת מימדיות. בגרף זה, הצירים אינם מייצגים משתנים ספציפיים ממערך הנתונים. במקום זאת, הם צירים שרירותיים שנוצרו על ידי אלגוריתם t-SNE כדי להמחיש את הנתונים במרחב ממדי מצומצם. אלגוריתם t-SNE הופך נתונים בעלי מימד גבוה למרחב בעל מימד נמוך תוך שמירה

על המבנה המקומי והיחסים בין נקודות הנתונים. הוא משמש לזיהוי אשכולות או דפוסים שעשויים להתקיים בתוך הנתונים.

בניתוח זה, הגרף המתקבל מספק ייצוג חזותי של נקודות הנתונים המקובצות במרחב דו מימדי, ומסייע בזיהוי ובפירוש של דפוסים או קבוצות בתוך מערך הנתונים.

חשוב להתייחס לכך שהתוספת של מודל האשכולות K-means הציגה מימד חדש לניתוח שלנו. **לאחר החלת K-means clustering והוספת תוויות האשכולות כעמודה חדשה למערך הנתונים, ראינו חלוקה ברורה ומושלמת לשני אשכולות מופרדים לחלוטין** (כפי הנראה בגרף לעיל). יתר על כן, בעת ניתוח ההתפלגות המגדרית בתוך כל אשכול, מצאנו שהאשכולות היו מופרדים היטב על סמך מגדר. זה מדגיש את החשיבות של מודל האשכולות K-means בקיבוץ נקודות נתונים דומות יחדיו וחשיפת דפוסים בסיסיים בנתונים. **ובהתאמה ישירה לכך, אכן המודל הנבחר הסופי עם טכניקת הנרמול המתאימה לכל עמודה נתן את המדדים הטובים ביותר לסט הנתונים שלנו עם תוספת העמודה שיצר מודל האשכולות** (נתייחס לכך בסעיף המודל הנבחר).

נתייחס לתוצאות המדדים השונים וננתח את כל האפשרויות והסיבות לשינויים:

1. נרמול ע"י MIN_MAX

א. עץ החלטה: הבחירה הראשונית של שימוש בקנה מידה מינימלי-מקסימלי עם מודל עץ ההחלטות הביאה לביצועים נאותים בכל המדדים. עם זאת, הבחנו שהדיוק, הריקול וציון ה F1 היו נמוכים יחסית לאפשרויות אחרות.

ב. יער אקראי: עם המעבר למודל היער האקראי, נצפה שיפור משמעותי בכל המדדים. הדיוק, וציון ה F1 עלו בצורה ניכרת, מה שמצביע על כך שאופי של מודל יער אקראי התאים יותר למערך הנתונים.

ג. naivebayes: למרות תוצאות סבירות שהושגו עם עץ ההחלטות והיער האקראי, מודל naivebayes הניב ביצועים מעט נמוכים יותר בכל המדדים.

ד. KNN: נבדקו ערכים שונים של K (3,5,7) תוך שימוש בקנה מידה מינימלי-מקסימלי. התוצאות הצביעו על כך ש K=5 הפיק את הביצועים הכוללים הטובים ביותר. היא נבחרה כאופציה המועדפת במקרה הזה.

2. נרמול ע"י Z-score

א. עץ החלטה: מעבר לקנה מידה של ציון z עם מודל עץ ההחלטות הביא לשיפור הביצועים בכל המדדים.

ב. יער אקראי: בדומה לעץ ההחלטות, יער אקראי עם קנה מידה של ציון z השיג ביצועים גבוהים יותר בהשוואה למקבילו לקנה מידה מינימלי. המדדים השתפרו באופן משמעותי, והדגישו עוד יותר את האפקטיביות של קנה מידה של ציון z בשיפור יכולת הניבוי של המודל.

ג. Naivebayes: הביצועים של Bayes הנאיביים עם קנה מידה של ציון z נשארו דומים לביצועים שלו עם קנה מידה מינימום-מקסימום.

ד. KNN: התוצאות של KNN עם קנה מידה של ציון z היו עקביות עם אפשרות קנה המידה של המינימום המקסימלי.

ה. K-means Clustering: בדומה לתרחיש הקודם, K-means Clustering עם קנה מידה Z-Score הניב ביצועים נמוכים יותר בהשוואה לדגמים אחרים.

3. Decimal scaling

- א. עץ החלטה: מודל עץ ההחלטות עם קנה מידה עשרוני הפגין ביצועים דומים לקנה מידה מינימום-מקסימום, והשיג תוצאות נאות בכל המדדים. עם זאת, הוא עדיין נפל בהשוואה לקנה מידה של ציון z במונחים של דיוק כולל וציון $F1$.
- ב. יער אקראי: מודל היער האקראי עם קנה מידה עשרוני הציג ביצועים משופרים מעט בהשוואה לקנה מידה מינימלי. אמנם המדדים לא היו שונים באופן משמעותי, אך העלייה ברמת הדיוק, הדיוק, הריקול והציון $F1$ הייתה בולטת.
- ג. Naivebayes: המודל עם קנה מידה עשרוני הראה ביצועים דומים כמו בטכניקות קנה המידה האחרות.
- ד. KNN: המודל עם קנה מידה עשרוני הציג תוצאות עקביות.

4. נרמול לפי חלוקה לפחים

- א. עץ החלטה: בעת החלת נורמליזציה של חלוקה בינומית על עמודות ספציפיות, ראינו תוצאות דומות לגישת קנה המידה של המינימום המקסימלי. מודל עץ ההחלטות הציג ביצועים נאותים בכל המדדים, אך לא עלה על הביצועים של קנה מידה של ציון z .
- ב. יער אקראי: ניצול נורמליזציה של חלוקה בינומית במודל היער האקראי הוביל לשיפור ניכר בכל המדדים. הדיוק, הדיוק, ההיזכרות והציון $F1$ עלו באופן משמעותי, מה שמצביע על כך שאלגוריתם היער האקראי התאים היטב לטיפול במערך הנתונים בטכניקת נורמליזציה זו.
- ג. Naivebayes: בדומה לעץ ההחלטות, מודל Bayes הנאיבי הראה ביצועים עקביים עם נורמליזציה של חלוקה בינומית, אם כי מעט נמוך יותר בהשוואה לטכניקות נורמליזציה אחרות.
- ד. KNN: ניסינו עם ערכים שונים של K ($K=3,5,7$) באמצעות נורמליזציה של חלוקה בינומית. התוצאות הצביעו על כך ש- $K=5$ הפיק את הביצועים הכוללים הטובים ביותר, מה שהפך אותו לאופציה המועדפת עבור KNN עם טכניקת נורמליזציה זו.

כחלק מהתהליך השלם, בדקנו את ערך האנטרופיה עבור כל עמודה וניסינו להריץ שוב את המודלים השונים לאחר בחירת פיצ'רים, כלומר בחרנו עמודות חדשות בעלות ערך האנטרופיה הנמוך ביותר. הערכנו שלאחר בחירת עמודות בעלות אנטרופיה נמוכה נצליח לשפר את המודל.

עם זאת, ההשפעה של בחירת תכונה על המדדים לא הייתה כצפוי. בניגוד להשערה הראשונית, המדדים הראו למעשה ירידה בביצועים לאחר בחירת תכונה, במיוחד במונחים של דיוק. עניין זה מרמז שייתכן שהתכונות שנבחרו לא ייצגו כראוי את הדפוסים הבסיסיים בנתונים, מה שגרם לירידה ברמת הדיוק החזוי של המודל. בהשוואת המדדים שהתקבלו לאחר בחירת תכונות עם המדדים שהושגו בעת הפעלת המודלים על כל הנתונים, ניכר כי בחירת תכונות לא שיפרה את הביצועים הכוללים.

לאחר ניתוח מדוקדק, ניתן לייחס מספר סיבות לתפקוד נמוך של המודל לאחר בחירת תכונה. ראשית, ייתכן שאלגוריתם בחירת התכונות לא תפס ביעילות את התכונות המבדילות ביותר עבור משימת הסיווג. ייתכן שהגישה מבוססת האנטרופיה במקרה זה לא הייתה מתאימה למערך הנתונים, וכתוצאה מכך הוסרו תכונות חשובות. שנית, הנוכחות של תכונות מתואמות או הרגישות של אלגוריתם היער האקראי שהוכח כמודל הטוב ביותר עבור סט הנתונים שלנו לתת-קבוצת התכונות שנבחרה עשוי להשפיע לרעה על ביצועי המודל.

כמו כן, ייתכן שבחירת התכונות לא שיפרה את התוצאות לאור כך שגילנו במהלך העבודה שהעמודה מגדר הינה עמודה לא מדויקת לחיזוי בסט הנתונים שלנו ואילו היינו בוחרות לחזות את העמודה של יעילות שעות שינה ייתכן שהתוצאות היו משתפרות.

אף על פי זאת, לאור כך שהשגנו תוצאות גבוהות במיוחד 97.8% במודל היער האקראי, עם סוגי נרמול שונים ראינו לנכון לא לשנות שוב את בחירת העמודות ולהשאיר את המודל הטוב שהושג עם סט הנתונים השלם בתוספת עמודת האשכולות החדשה.

המודל הנבחר ומסקנות

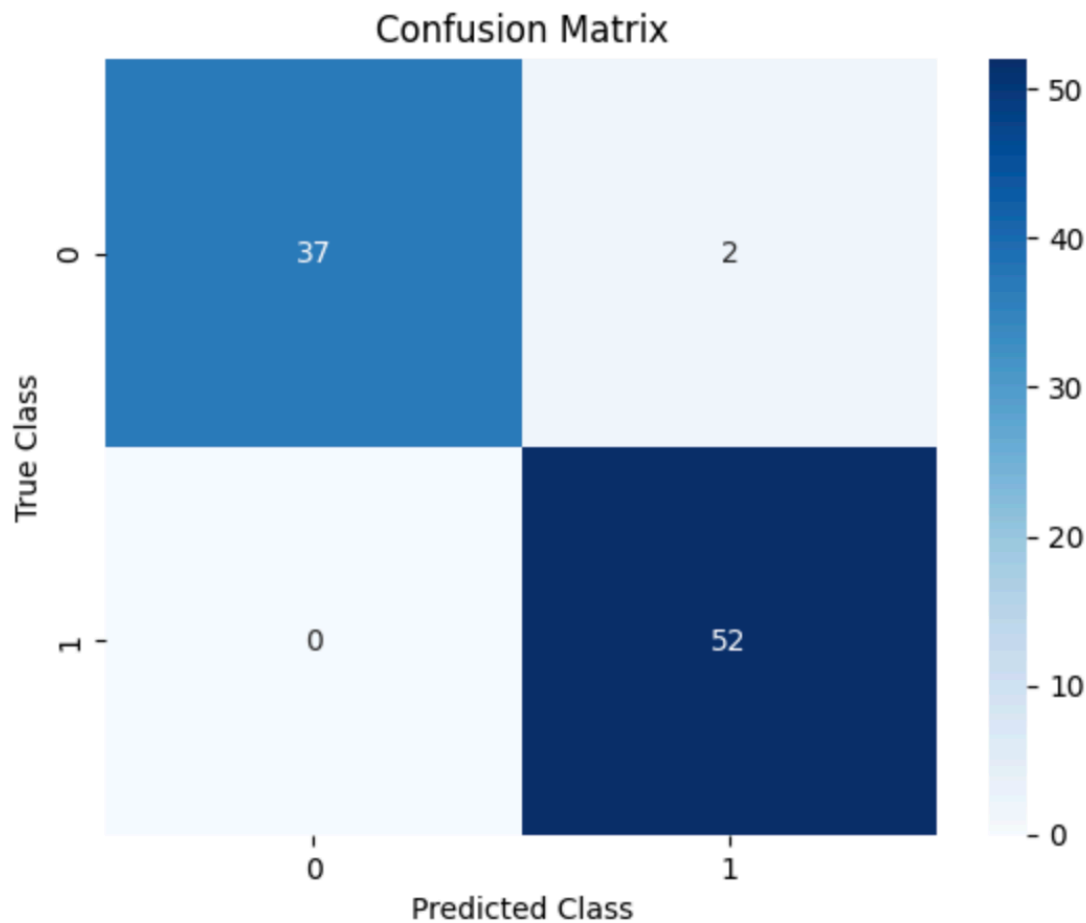
בהתבסס על הניתוח המקיף שביצענו לאורך הפרויקט, המודל הסופי שבחרנו כמודל בעל הביצועים הטובים ביותר הוא **מודל היער האקראי עם נורמליזציה מותאמת אישית לפי עמודות**. מודל זה הניב תוצאות מרשימות, והשיג **דיוק של 97.8% והפגין שיפורים משמעותיים לעומת טכניקות ומודלים מקדימים אחרים שנחקרו במהלך הפרויקט שלנו**.

אלגוריתם היער האקראי, שיטת למידה אננסמבלית, שיחקה תפקיד מרכזי בהצלחת המודל הסופי שלנו. על ידי שילוב של עצי החלטה מרובים, אלגוריתם זה לכד ביעילות דפוסים מורכבים בתוך מערך הנתונים וערך תחזיות מדויקות. על ידי צבירה של התחזיות של עצים מרובים, מודל היער האקראי הפיק תוצאות חזקות ומהימנות.

אחד היתרונות המרכזיים של אלגוריתם Random Forest הוא היכולת שלו להתמודד עם מספר רב של תכונות. בפרויקט שלנו, זה היה **מועיל במיוחד מכיוון שהיה לנו סט מגוון של תכונות לעבוד איתם**. יתרה מכך, **הבחירה האקראית של התכונות במהלך תהליך האימון סייעה בהפחתת ההשפעה של תכונות רועשות או לא רלוונטיות בודדות**, והעלתה את כוח הניבוי הכולל של המודל. על ידי מינוף מנגנון בחירת תכונה זה, מודל ה-Random Forest זיהה והשתמש ביעילות בתכונות האינפורמטיביות ביותר, מה שהוביל לשיפור הביצועים.

במהלך תהליך שיפור המודל, גילינו שעמודות שונות במערך הנתונים שלנו הגיבו בצורה שונה לטכניקות נורמליזציה שונות. כדי לנצל את התובנה הזו, **תכננו פונקציה המיישמת שיטות נורמליזציה ספציפיות על כל עמודה, תוך התחשבות בטכניקת הנורמליזציה שהשיגה את התוצאה הטובה ביותר במהלך שלב הניסויים שלנו**. על ידי שילוב נורמליזציה מותאמת אישית זו, הצלחנו לייעל את ייצוג הנתונים ולספק למודל תשומות משמעותיות יותר וניתנות להשוואה.

המודל הסופי הפגין דיוק יוצא דופן, ריקול וציון F1 וכמו כן הפגין את יכולתו לבצע תחזיות מדויקות לסוג ביעילות מופעים בתוך מערך הנתונים. מטריצת הבלבול (טעויות) ממחישה את הביצועים של המודל, עם זיהוי נכון של 37 נשים מתוך 39 ושל 52 גברים מתוך 52 (הנתונים האלו מייצגים את הביצועים של המודל על סט הנתונים שהוקצה לבדיקה (x_test)), מה שמצביע על מהימנותו בחיזוי מקרים חיוביים.



לסיכום, באמצעות הערכה מקיפה של טכניקות ומודלים שונים של עיבוד מקדים, הפרויקט שלנו זיהה בהצלחה את מסווג היער האקראי עם נורמליזציה ספציפית לעמודה מותאמת אישית **כפתרון האופטימלי עבור מערך הנתונים הנתון (98%)**. על ידי מינוף למידת אנסמבל, בחירת תכונות וטכניקות נורמליזציה מותאמות, השגנו תוצאות יוצאות דופן וסיפקנו מודל אמין לחיזוי מגדר וניתוח יעילות השינה.

בברכה,
מתן ניצן ואסנת שבתאי

