

## 1.1 רקע

פרויקט זה עוסק בפיתוח מערכת מתקדמת לזיהוי הונאות ופעילות בלתי חוקית בעסקאות ביטקוין. המערכת מתבססת על ניתוח מעמיק של גרף העסקאות של הביטקוין, תוך שימוש בטכניקות של למידת מכונה ואנליזת נתונים מתקדמת. הפרויקט נשען על מערך נתונים ייחודי שמופה לישויות אמיתיות, המאפשר זיהוי מדויק של פעילות חשודה.

## 1.2 מטרות הפרויקט

- זיהוי מדויק של עסקאות ביטקוין בלתי חוקיות, כולל הונאות, תוכנות זדוניות, מימון טרור, וסכמות פונזי.
- סיווג אוטומטי של עסקאות לקטגוריות חוקיות ובלתי חוקיות.
- שיפור משמעותי בדיוק הזיהוי של פעילות בלתי חוקית בהשוואה לשיטות הקיימות בתעשייה.
- יצירת כלי יעיל לרשויות אכיפת החוק ומוסדות פיננסיים למאבק בפשיעה פיננסית בעולם הקריפטו.

## 1.3 יעדים מדידים

- השגת דיוק (accuracy) של לפחות 97% בסיווג עסקאות.
- השגת רגישות (recall) של לפחות 80% עבור זיהוי עסקאות בלתי חוקיות.
- שיפור ה-F1 score ל-0.8 ומעלה, המשקף איזון טוב בין דיוק לרגישות.
- הפחתת שיעור האזעקות השווא (false positives) ל-5% או פחות.

## 2. ארכיטקטורת המערכת

### 2.1 מבנה כללי

המערכת בנויה ממספר מודולים מרכזיים המשתלבים יחד ליצירת פתרון מקיף:

1. \*\*מודול טעינת נתונים\*\*<sup>1</sup>: אחראי על קריאה וארגון של נתוני העסקאות מקבצי ה-CSV.
2. \*\*מודול הנדסת תכונות\*\*<sup>2</sup>: מבצע עיבוד ראשוני והעשרה של הנתונים.
3. \*\*מודול אימון מודל\*\*<sup>3</sup>: אחראי על בניית ואימון מודל הלמידה העמוקה.
4. \*\*מודול זיהוי אנומליות\*\*<sup>4</sup>: מיישם אלגוריתמים לזיהוי עסקאות חריגות.
5. \*\*מודול ויזואליזציה\*\*<sup>5</sup>: מייצר גרפים ותרשימים להצגת התוצאות.
6. \*\*מודול הערכת ביצועים\*\*<sup>6</sup>: מחשב ומציג מדדי ביצוע של המודל.

## 2.2 זרימת מידע

1. נתוני העסקאות נטענים ממקורות חיצוניים.
2. הנתונים עוברים תהליכי ניקוי, נרמול והעשרה.
3. המודל מאומן על סט האימון ומכיל על סט הוולידציה.
4. עסקאות חדשות מועברות דרך המודל לזיהוי אנומליות.
5. תוצאות הניתוח מוצגות באמצעות ממשק המשתמש וכלי הוויזואליזציה.

## 3. מפרט טכני

### 3.1 טכנולוגיות ושפות תכנות

- \*\*שפת תכנות ראשית\*\* Python 3.x
- \*\*ספריות מרכזיות\*\*:
- NumPy ו-Pandas לעיבוד נתונים
- Scikit-learn למודלים של למידת מכונה
- TensorFlow או PyTorch לרשתות נוירונים עמוקות
- NetworkX לניתוח גרפים
- Matplotlib ו-Seaborn לויזואליזציה

### 3.2 מבנה הנתונים

- \*\*מקור הנתונים\*\* :גרף עסקאות ביטקוין מה-Elliptic Data Set
- \*\*גודל המערכת\*\* : 203,769 צמתים (עסקאות) ו-234,355 קשתות
- \*\*מאפיינים\*\* : 166 מאפיינים לכל צומת, כולל מידע מקומי ומצרפי
- \*\*תיוג\*\* : 2% מהצמתים מסווגים כבלתי חוקיים, 21% כחוקיים, והשאר לא מסווגים

### 3.3 אלגוריתמים ומודלים

- \*\*מודל בסיסי\*\* : Random Forest Classifier
- \*\*מודלים מתקדמים לשקילה\*\* :
- Graph Convolutional Networks (GCN)
- Graph Attention Networks (GAT)
- Temporal Graph Networks (TGN) לניתוח דינמיקה לאורך זמן
- \*\*טכניקות לטיפול בחוסר איזון\*\* :
- SMOTE (Synthetic Minority Over-sampling Technique)

- Class weighting
- Ensemble methods with balanced bagging

## 4. מודולים ופונקציונליות

### 4.1 מודול טעינת נתונים

- \*\*פונקציונליות\*\*:
- קריאת קבצי CSV של עסקאות, קשרים ומאפיינים
- טיפול בערכים חסרים וחריגים
- הפרדת הנתונים לסטים של אימון (60%), וולידציה (20%) ובדיקה (20%)
- \*\*יישום\*\*:
- שימוש ב-Pandas לקריאה יעילה של קבצים גדולים
- מימוש פונקציות לניקוי וטרנספורמציה של נתונים

### 4.2 מודול הנדסת תכונות

- \*\*פונקציונליות\*\*:
- נרמול וסטנדרטיזציה של מאפיינים מספריים
- קידוד מאפיינים קטגוריים
- יצירת מאפיינים מצרפיים מבוססי גרף (כגון דרגת צומת, מרכזיות)
- בחירת מאפיינים רלוונטיים באמצעות שיטות כמו PCA או feature importance
- \*\*יישום\*\*:
- שימוש ב-Scikit-learn לטרנספורמציות סטנדרטיות
- פיתוח פונקציות מותאמות ליצירת מאפייני גרף באמצעות NetworkX

### 4.3 מודול אימון מודל

- \*\*פונקציונליות\*\*:
- אימון מודל Random Forest כבסיס
- יישום וכוונון מודלים מבוססי גרף (GCN, GAT)
- ביצוע cross-validation לבחירת היפר-פרמטרים אופטימליים
- שימוש בטכניקות לטיפול בחוסר איזון בנתונים
- \*\*יישום\*\*:
- שימוש ב-Scikit-learn עבור Random Forest

- יישום מודלים מבוססי גרף באמצעות PyTorch Geometric או DGL
- פיתוח מנגנון לכוונן אוטומטי של היפר-פרמטרים

#### 4.4 מודול זיהוי אנומליות

- \*\*פונקציונליות\*\*:
- חישוב ציוני אנומליה לעסקאות חדשות
- יישום אלגוריתמים לזיהוי אנומליות מבוססי גרף
- קביעת סף דינמי לסיווג עסקאות כחשודות
- \*\*יישום\*\*:
- פיתוח אלגוריתם המשלב את תוצאות המודל עם ניתוח טופולוגי של הגרף
- יישום שיטות לקביעת סף אדפטיבי המתעדכן עם הזמן

#### 4.5 מודול ויזואליזציה

- \*\*פונקציונליות\*\*:
- יצירת גרפים של חשיבות מאפיינים
- ויזואליזציית t-SNE או UMAP של העסקאות במרחב דו-ממדי
- גרף התפלגות ציוני אנומליה
- ויזואליזציה של תת-גרפים חשודים
- \*\*יישום\*\*:
- שימוש ב-Matplotlib ו-Seaborn ליצירת גרפים סטטיסטיים
- פיתוח ויזואליזציות אינטראקטיביות באמצעות Plotly או Bokeh

#### 4.6 מודול הערכת ביצועים

- \*\*פונקציונליות\*\*:
- חישוב מדדי ביצוע מקיפים: דיוק, רגישות, AUC-PR, AUC-ROC, F1 score
- יצירת מטריצת בלבול
- ניתוח שגיאות וזיהוי מקרים קשים
- \*\*יישום\*\*:
- שימוש בפונקציות מ-Scikit-learn לחישוב מדדים
- פיתוח דוחות ביצועים אוטומטיים

### 5. ממשק משתמש

## **5.1 ממשק קונסולה**

- הצגת התקדמות האימון בזמן אמת
- דיווח על מדדי ביצוע עיקריים
- אפשרות להפעלת ניתוחים שונים דרך פקודות

## **5.2 ממשק גרפי**

- דשבורד אינטראקטיבי להצגת תוצאות וסטטיסטיקות
- כלים לחקירת עסקאות ספציפיות
- ויזואליזציות דינמיות של גרף העסקאות

## **6. דרישות אבטחה ופרטיות**

### **6.1 אנונימיזציה של נתונים**

- שמירה על אנונימיות של פרטי העסקאות
- הצפנת מידע רגיש

### **6.2 אבטחת מידע**

- הגנה על הנתונים ותוצאות המודל
- מערכת הרשאות לגישה למידע רגיש

## **7. ביצועים ומדדים**

### **7.1 זמני עיבוד**

- אימון המודל הראשוני: פחות מ-24 שעות
- זיהוי אנומליות בזמן אמת: פחות מ-5 שניות לעסקה
- עדכון המודל: פחות מ-6 שעות

### **7.2 דיוק ורגישות**

- דיוק כללי: שאיפה ל-97% ומעלה
- רגישות לעסקאות בלתי חוקיות: מטרה של 80% ומעלה
- F1 score: יעד של 0.8 ומעלה

### 7.3 יכולת הרחבה

- תמיכה בניתוח של לפחות 1 מיליון עסקאות ביום
- יכולת לעדכן את המודל עם נתונים חדשים באופן שוטף

### 8. סיכום

מערכת זיהוי הונאות הקריפטו המוצעת מציגה גישה חדשנית ומקיפה לזיהוי פעילות בלתי חוקית ברשת הביטקוין. באמצעות שילוב של טכניקות למידת מכונה מתקדמות, ניתוח גרפים, ומערכ נתונים עשיר, המערכת מבטיחה לספק כלי יעיל ומדויק למאבק בפשיעה פיננסית בעולם הקריפטו.

היתרונות העיקריים של המערכת כוללים:

1. דיוק גבוה בזיהוי עסקאות חשודות
2. יכולת לנתח כמויות גדולות של נתונים בזמן אמת
3. גמישות בהתאמה לדפוסי הונאה חדשים
4. ויזואליזציות מתקדמות להבנה מעמיקה של דפוסי פעילות
5. שילוב של ידע מתחום הפיננסים עם טכנולוגיות מתקדמות של בינה מלאכותית

האתגרים העיקריים שהמערכת מתמודדת איתם:

1. חוסר איזון משמעותי בנתונים בין עסקאות חוקיות ובלתי חוקיות
2. דינמיות גבוהה של דפוסי הונאה המשתנים במהירות
3. צורך באיזון בין דיוק גבוה לבין מיעוט אזעקות שווא
4. שמירה על פרטיות ואבטחת מידע תוך כדי ניתוח מעמיק

המערכת המוצעת מתוכננת להתמודד עם אתגרים אלו באמצעות:

- שימוש בטכניקות מתקדמות לטיפול בחוסר איזון בנתונים
- יכולת למידה מתמשכת והתאמה לדפוסים חדשים
- שילוב של מודלים סטטיסטיים עם למידת מכונה מבוססת גרפים
- מנגנוני אבטחה ואנונימיזציה מובנים

לסיכום, מערכת זיהוי הונאות הקריפטו המוצעת מהווה צעד משמעותי קדימה ביכולת לזהות ולמנוע פעילות פיננסית בלתי חוקית בעולם המטבעות הדיגיטליים. עם יישום מוצלח, היא צפויה לתרום משמעותית לאבטחה ולאמינות של מערכת הביטקוין ומטבעות קריפטוגרפיים אחרים.