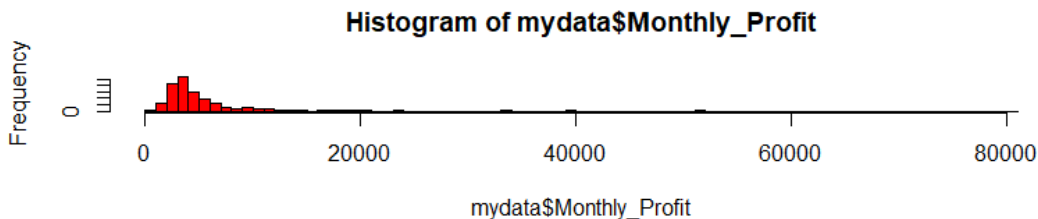


עבודה 2 מדעי הנתונים

1. ביצוע הדיסקרטיזציה:

- עבור Monthly_Profit השתמשנו ב- Equal-frequency discretization. ניתן לראות לפי ההיסטוגרמה של התפלגות הערכים שהרוב המוחלט של הערכים קטן מ-20,000 ועל כן רצינו ליצור חלוקה של שווה של ערכים בין ה-bins. בחרנו עבור משתנה זה ב-5 bins.



- עבור Spouse_Income השתמשנו בשיטה Equal-frequency discretization. הבחנו כי הערך "0" הופיע בכמות רבה של רשומות (מעל 400 רשומות), על כן החלטנו לחלק את הרשומות באופן שבו בכל אינטרוול מספר ערכים שווה. בחרנו עבור משתנה זה ב-5 bins.

- עבור Loan_Amount השתמשנו בשיטה Equal-frequency discretization. טווח הערכים במשתנה זה קטן יותר ולכן בחרנו לחלק משתנה זה ל-3 bins.

2. תצולמי העצים שהתקבלו מופיעים ברצף בדפים הבאים לפי הסדר הבא:

- Gini 50
- Gini 80
- Information 50
- Information 80

3.

Model	Accuracy
Gini30	0.7726
Information30	0.7893
Gini80	0.796
Information80	0.796

ניתן לראות כי הדיוק הטוב ביותר מתקבל כאשר משתמשים ב- $\text{minSplit}=80$ לא משנה באיזה מודל (gini או information).

