# Trees

# Previous class

- Regression-based Forecasting

Predictive Models

- Supervised
  - Regression
    - ➤ **_k_-Nearest Neighbor**
    - ➤ **Linear Regression**
    - ➤ **Regression Trees**
    - ➤ Neural Networks
    - ➤ Ensembles
    - ➤ ……
  - Classification
    - ➤ **_k_-Nearest Neighbor**
    - ➤ Naïve Bayes
    - ➤ **Logistic Regression**
    - ➤ **Classification Trees**
    - ➤ Neural Networks
    - ➤ Discriminant Analysis
    - ➤ Ensembles
    - ➤ ……
  - Time Series Forecasting
    - ➤ **Regression-based**
    - ➤ Smoothing methods
    - ➤ ……
- Unsupervised
  - Segmentation
    - ➤ **Clustering**
    - ➤ ……

# Trees

- Flexible data-driven method

- Used for

  - Classification ( called Classification Tree)

  - Regression (called Regression Tree)

- Transparent

- Easy interpretation

- Doesn't require enormous effort

- Method

  - **Recursive Partitioning** : Separating records into subgroups by creating splits on predictors

# Recursive Partitioning

- Outcome variable Y

- Predictor variables $X_1, X_2, X_3, \cdots\cdots X_p$

- Recursive Partitioning

  - ➢ Divides the p-dimensional space of predictors into non-overlapping multidimensional rectangles

- Accomplished recursively

  - ➢ Operating on the results of prior division

- Idea is to divide the entire variable-space up into rectangles such that each rectangle is as **homogeneous** or **pure**

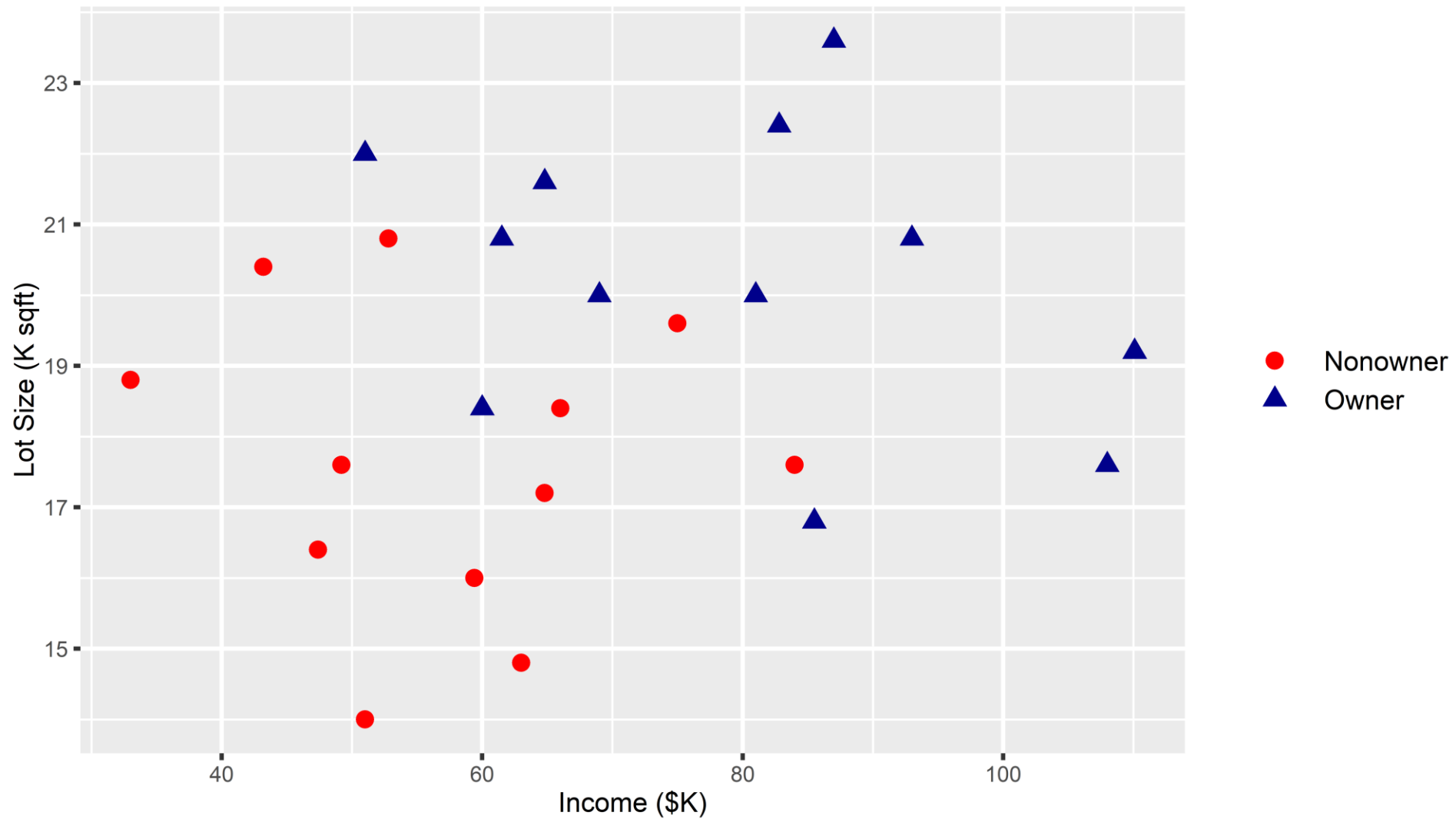- **Homogeneous** or **Pure** meaning containing records mostly of one class

# Data on Riding Mowers

- Riding-mower manufacturer would like to find a way of classifying families in a city into an **owner** or **non-owner**

- Attributes

  - Income : Income of the household in thousand of dollars

  - Lot Size : Lot size in thousand of square foot
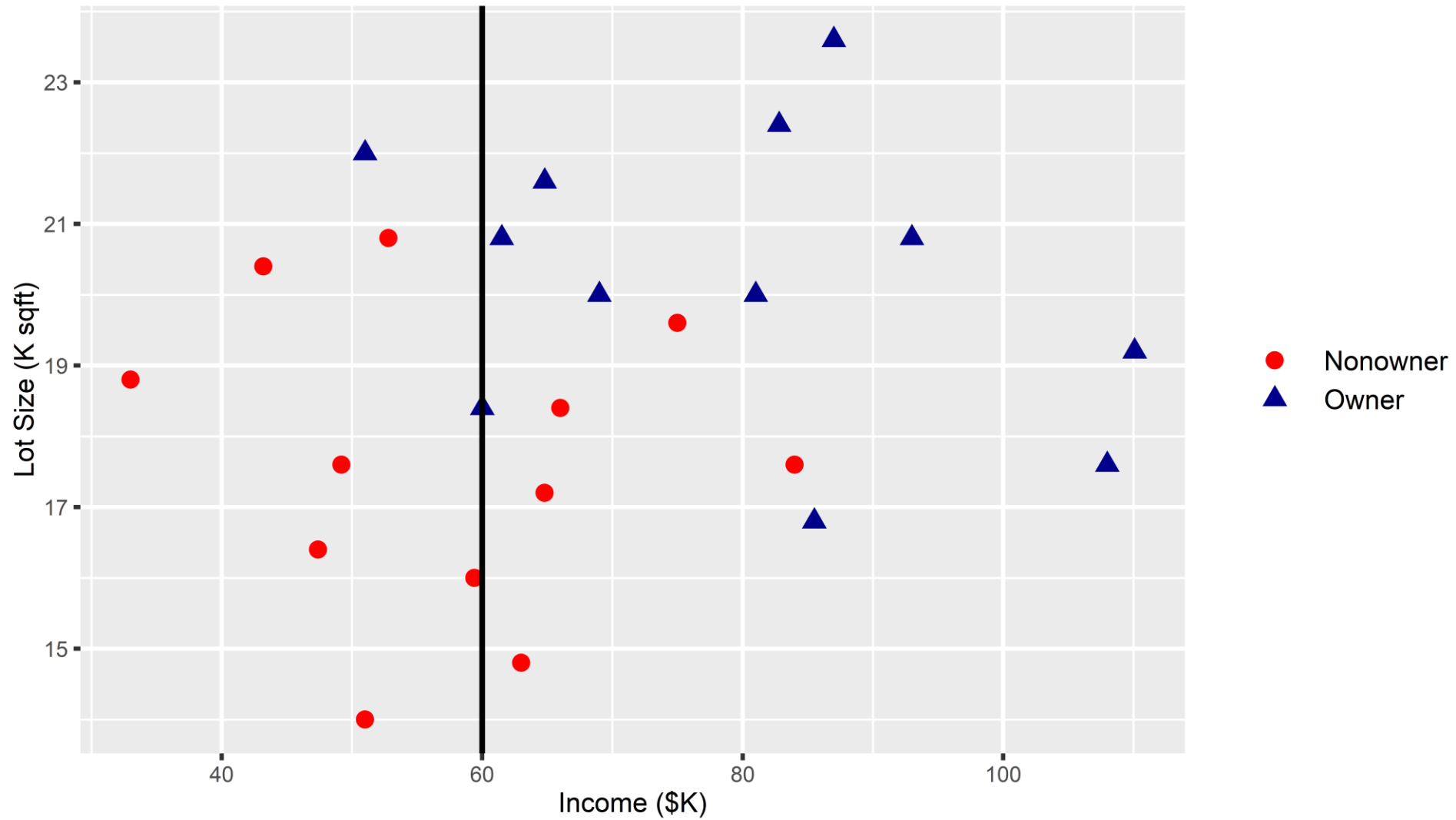
  - Ownership : Owner or Non-owner

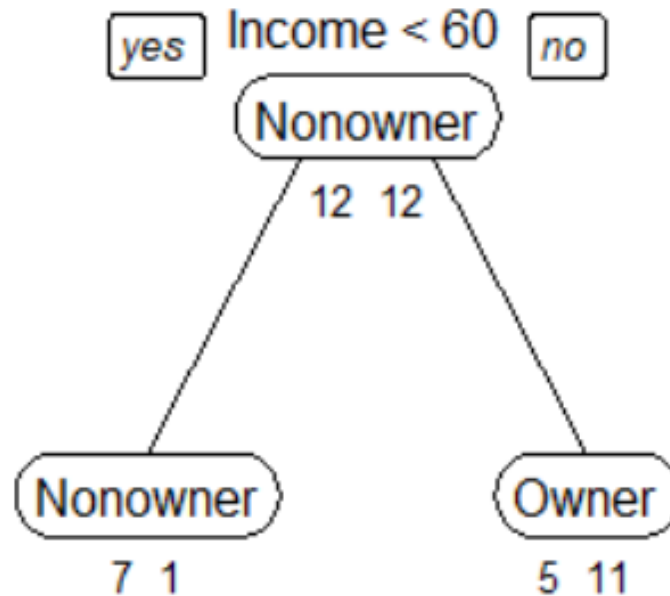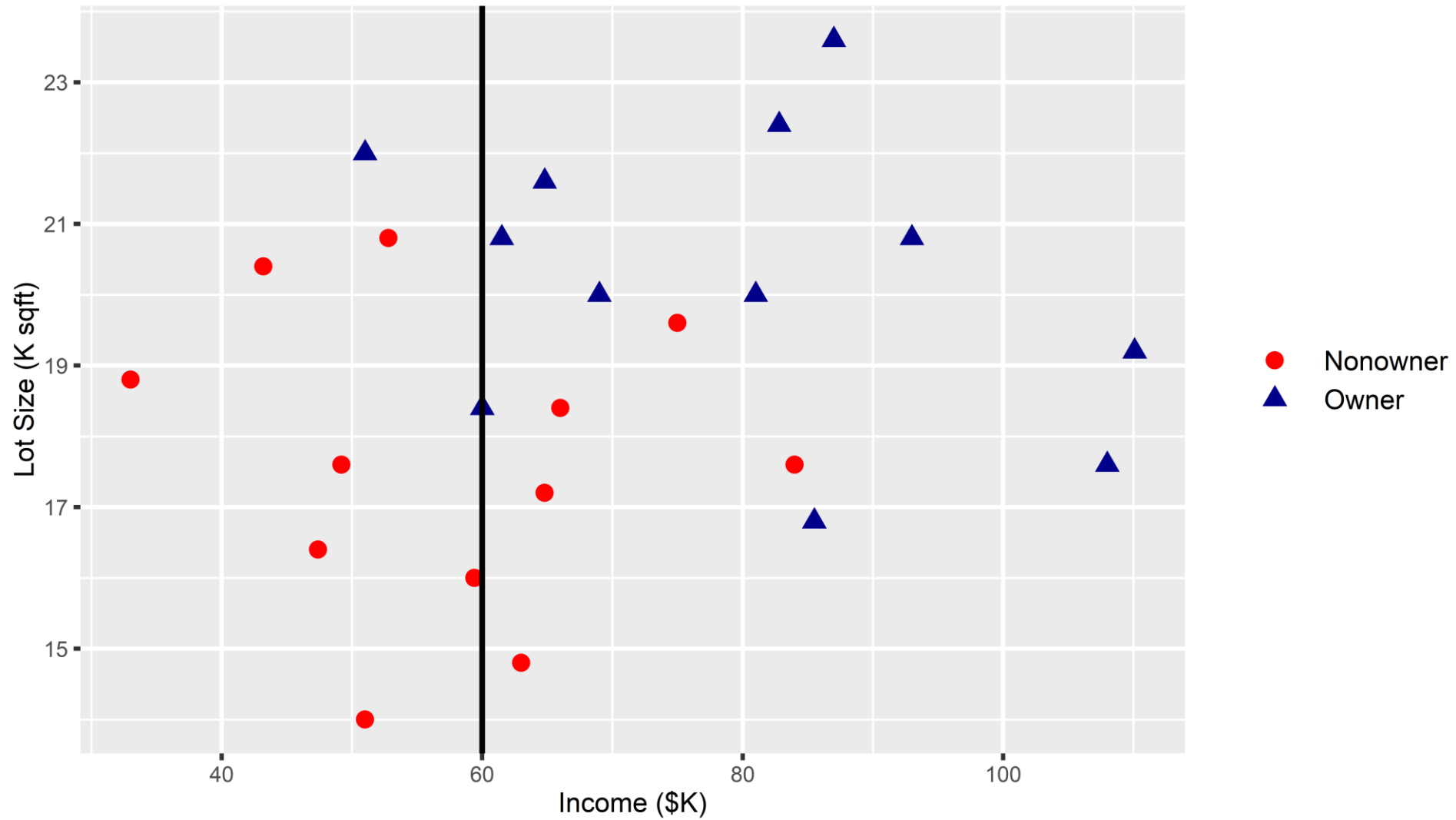| Income | Lot_Size | Ownership |
|--------|----------|-----------|
| 60     | 18.4     | Owner     |
| 85.5   | 16.8     | Owner     |
| 64.8   | 21.6     | Owner     |
| 61.5   | 20.8     | Owner     |

⋮
⋮

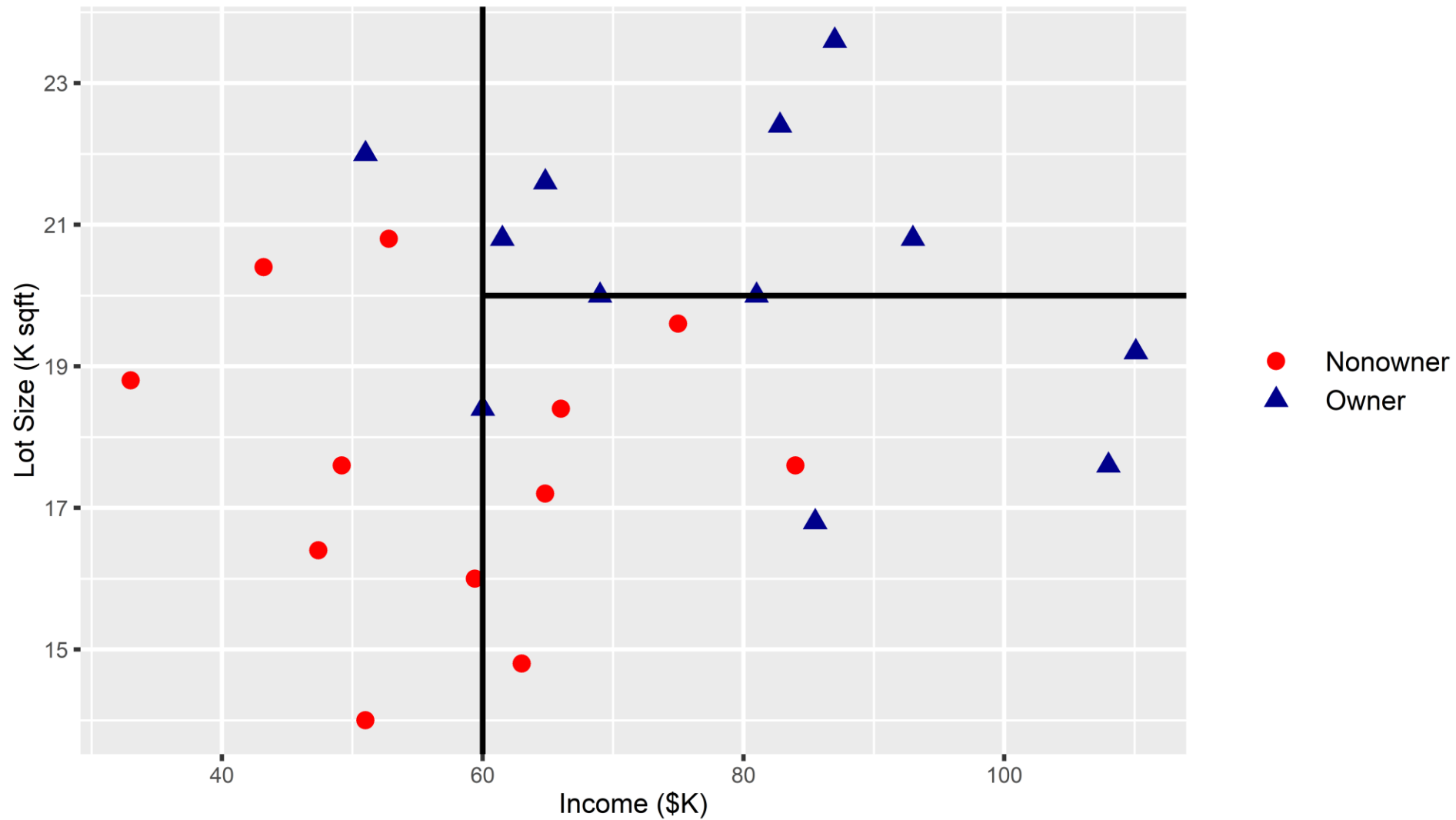Scatter plot

© Pradeep Pendem

# First split at Income = 60
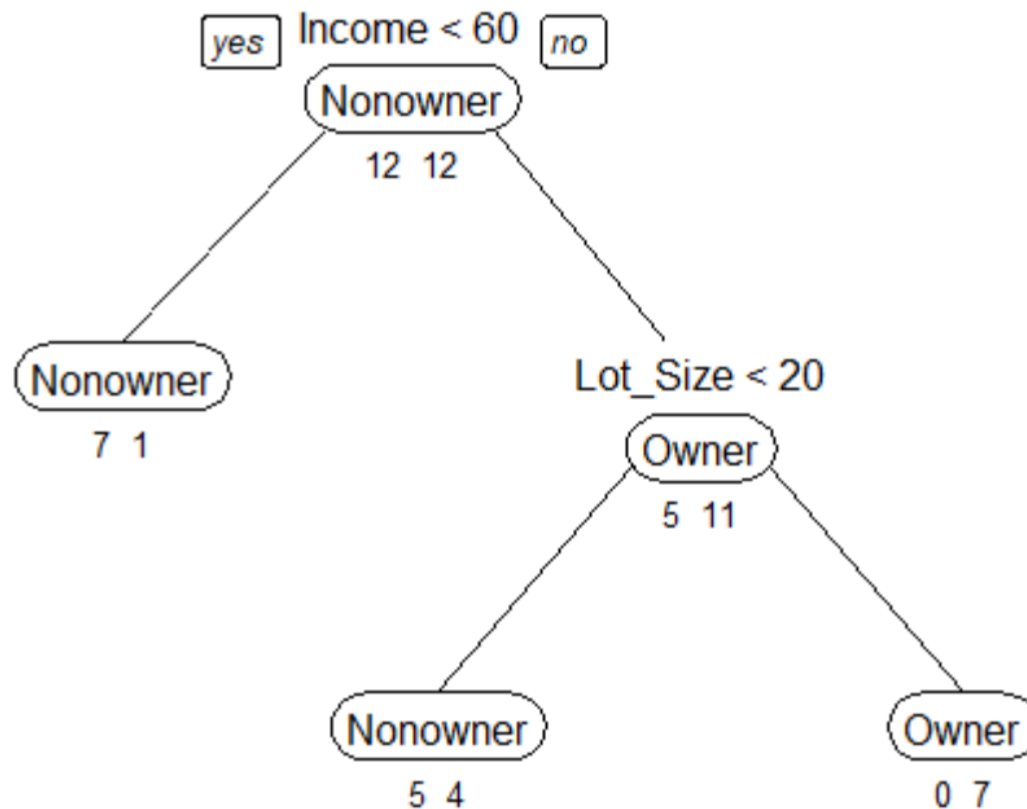
# First split at Income = 60
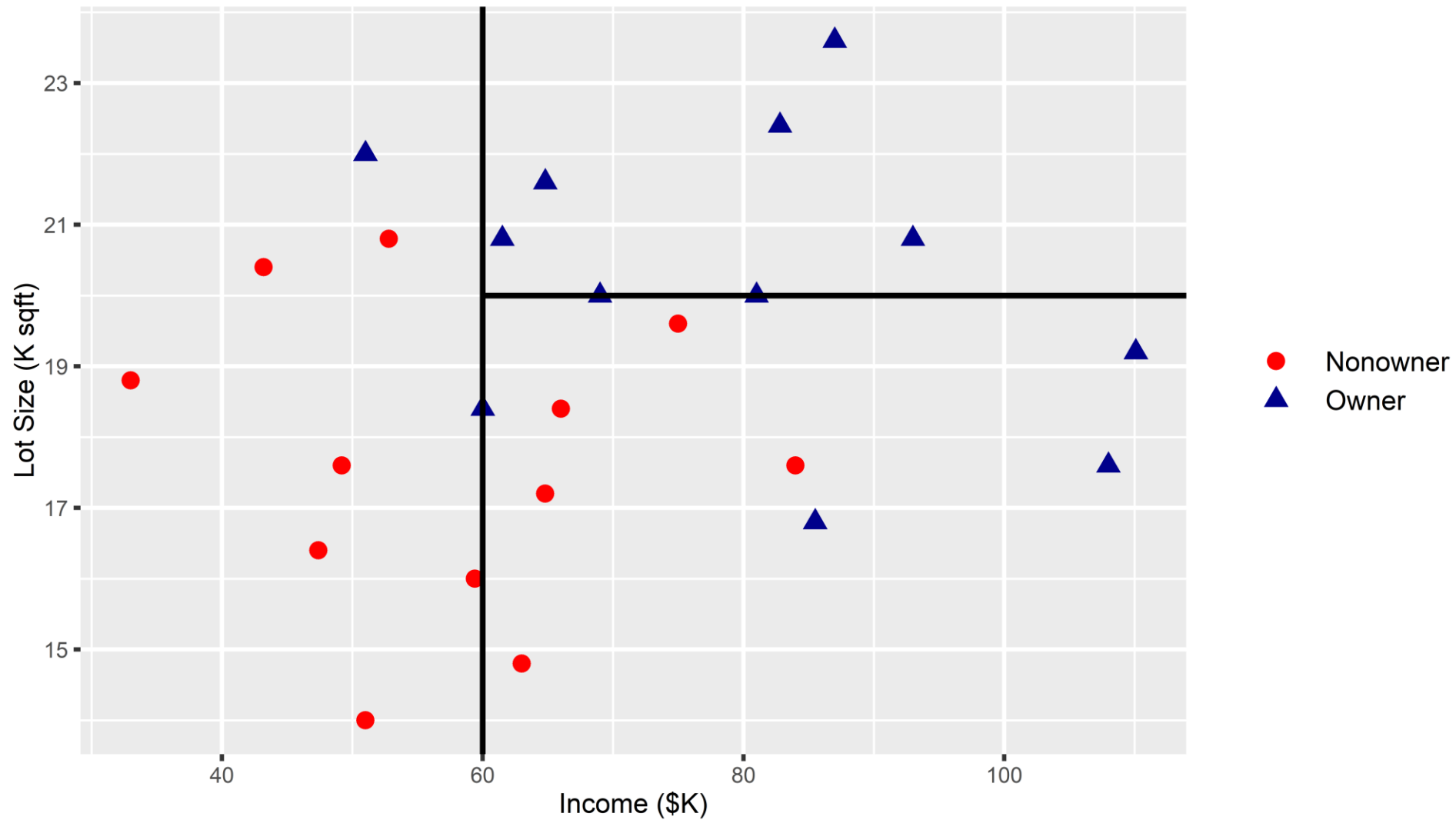
# First split at Income = 60
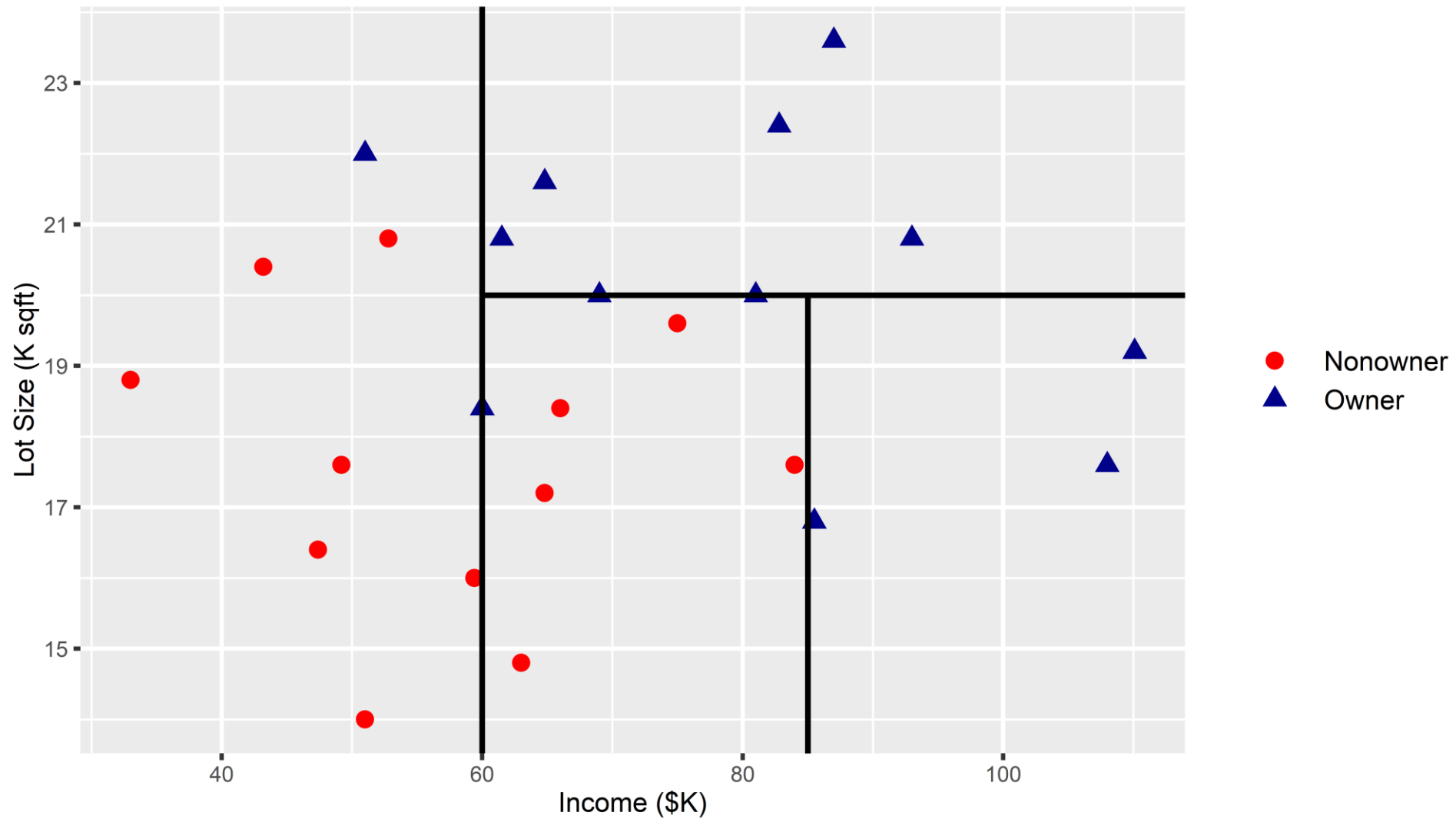
# Second split at Lot Size = 20
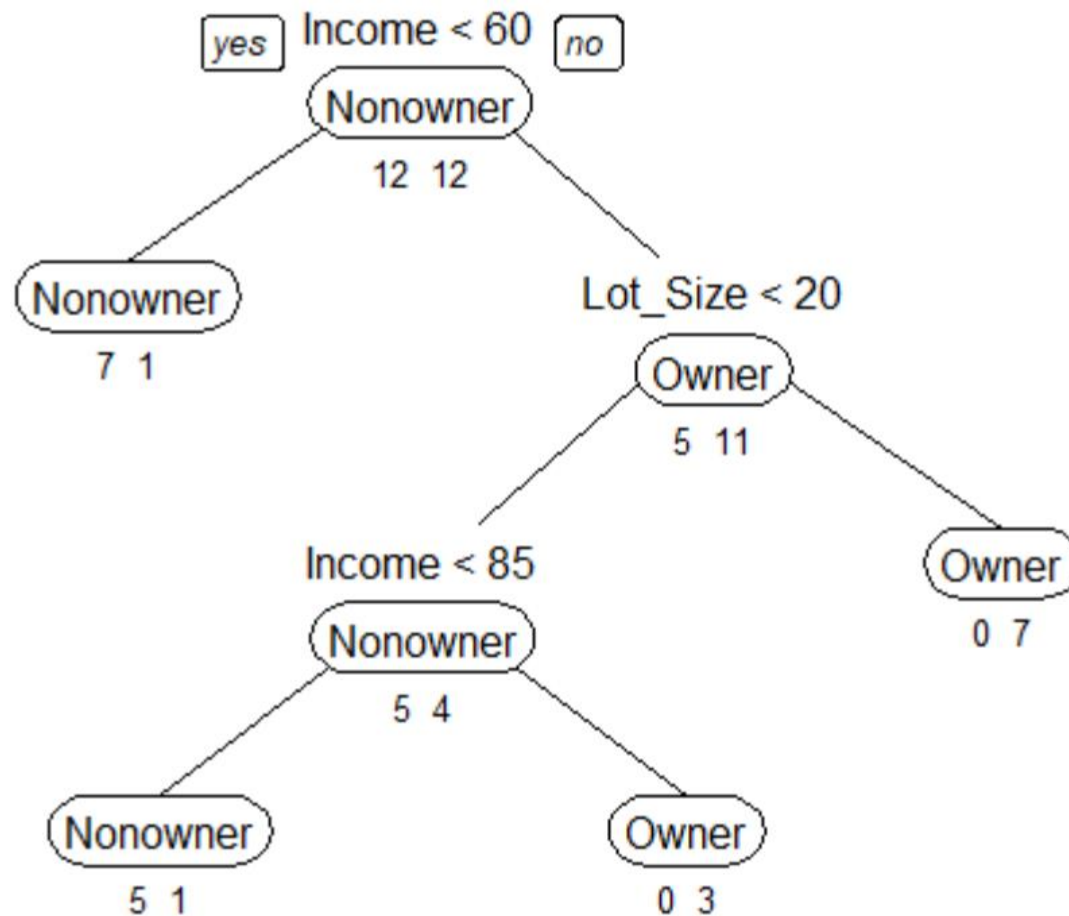
# Second split at Lot Size = 20
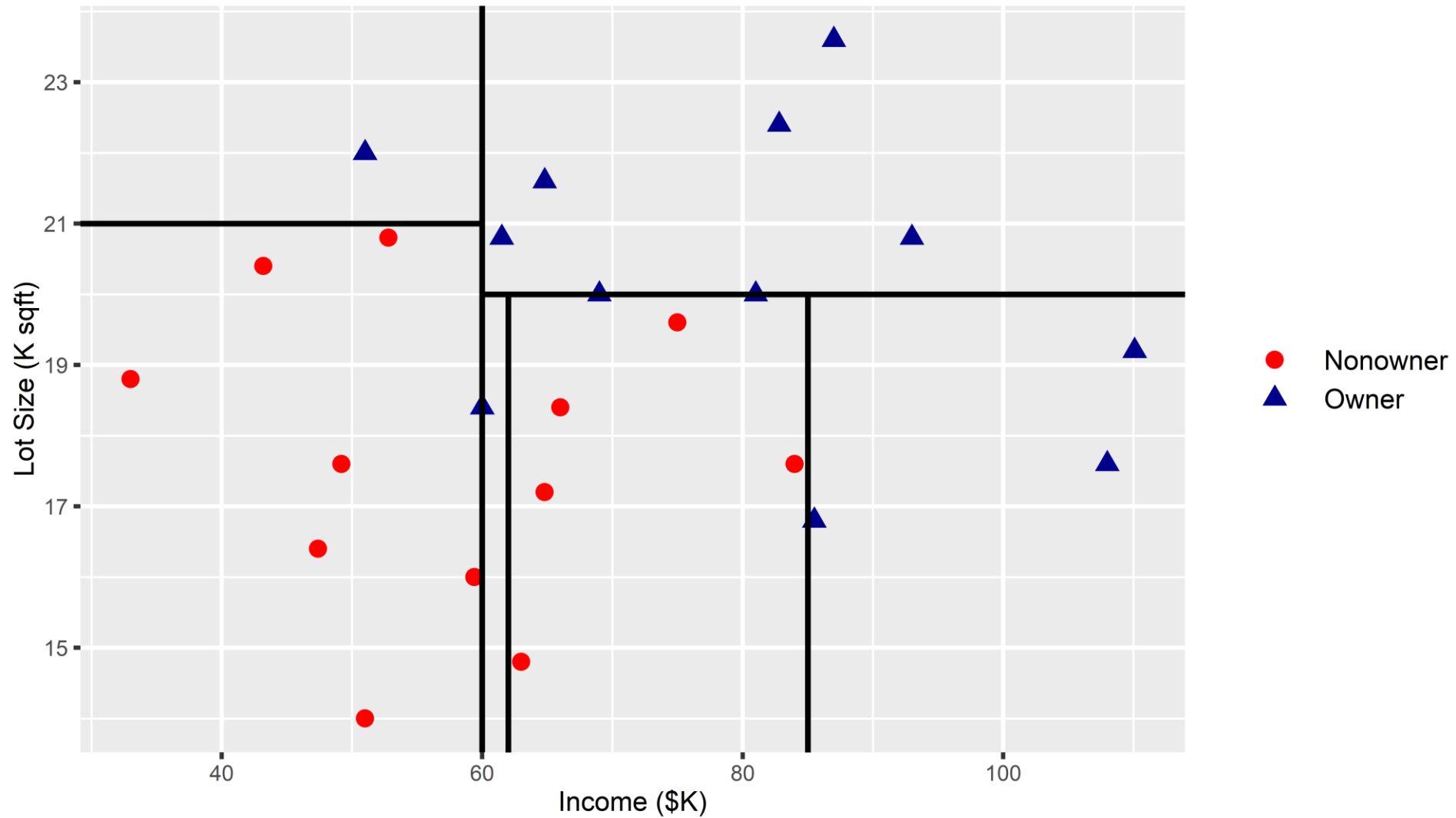
# Second split at Lot Size = 20

# Third split at Income = 85

# Third split at Income = 85

# Exhaustive splits

# Fully grown tree



Income < 60

| yes | | no |

Nonowner
12  12

Lot_Size < 21

Nonowner
7  1

Nonowner        Owner
7  0            0  1

Lot_Size < 20

Owner
5  11

Income < 85

Nonowner
5  4

Owner
0  7

Income >= 62

Nonowner
5  1

Owner
0  3

Nonowner        Owner
5  0            0  1

Decision node

Terminal node

- Size of the tree is calculated by the number of terminal modes = 6
- Number of splits = Size of the tree − 1 = 6 − 1 = 5

# Measures of Impurity
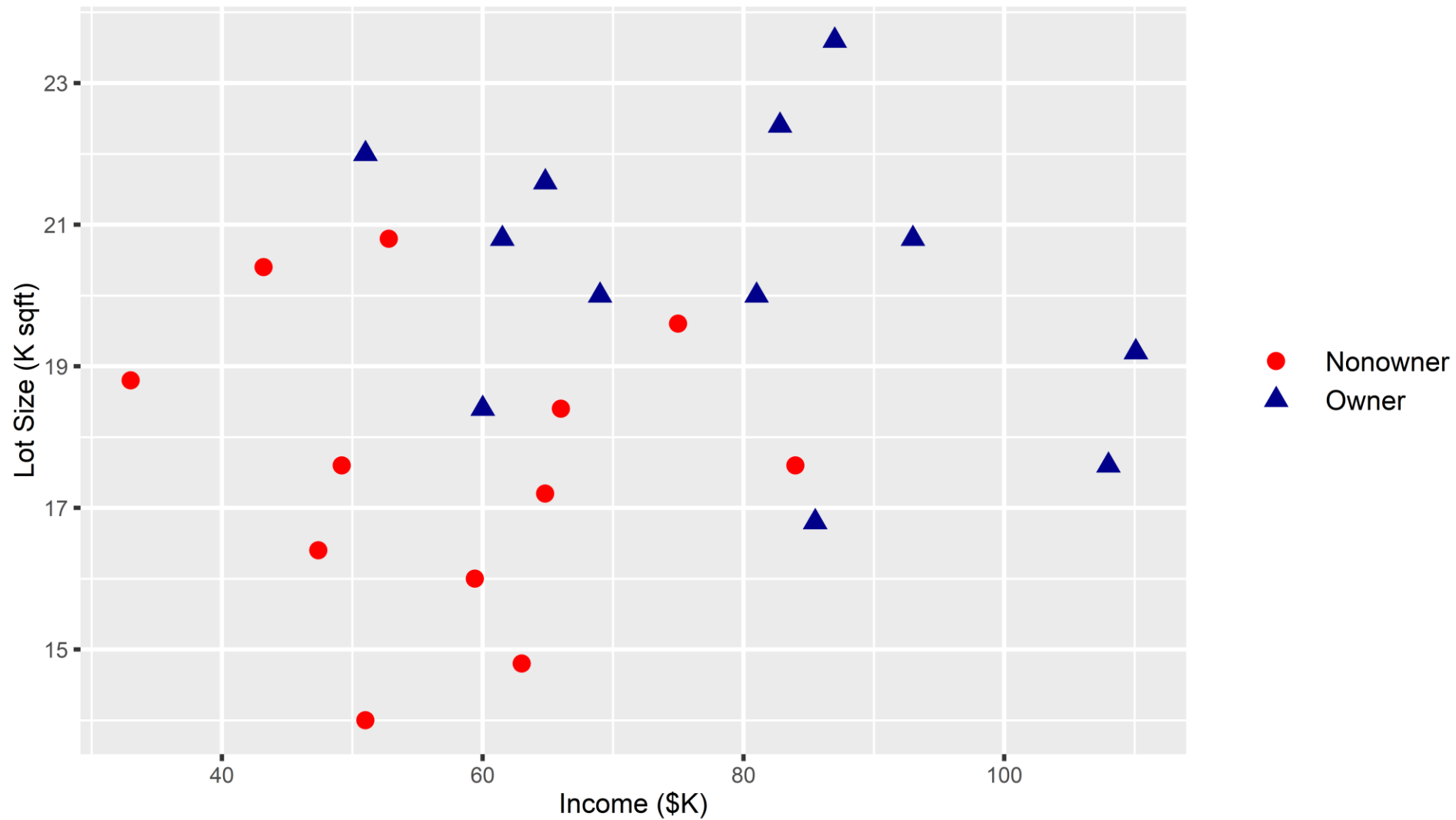
- Two popular measures

    ➢ Gini Index

    ➢ Entropy measure

- Gini impurity index for a rectangle A is given by

$$I(A) = 1 - \sum_{k=1}^{m} p_k^2$$

- $p_k$ is the proportion of records in rectangle A that belong to class $k$

- Measure takes values between 0 (when all records belong to same class) and $m-1/m$ (when all m classes are equally represented)
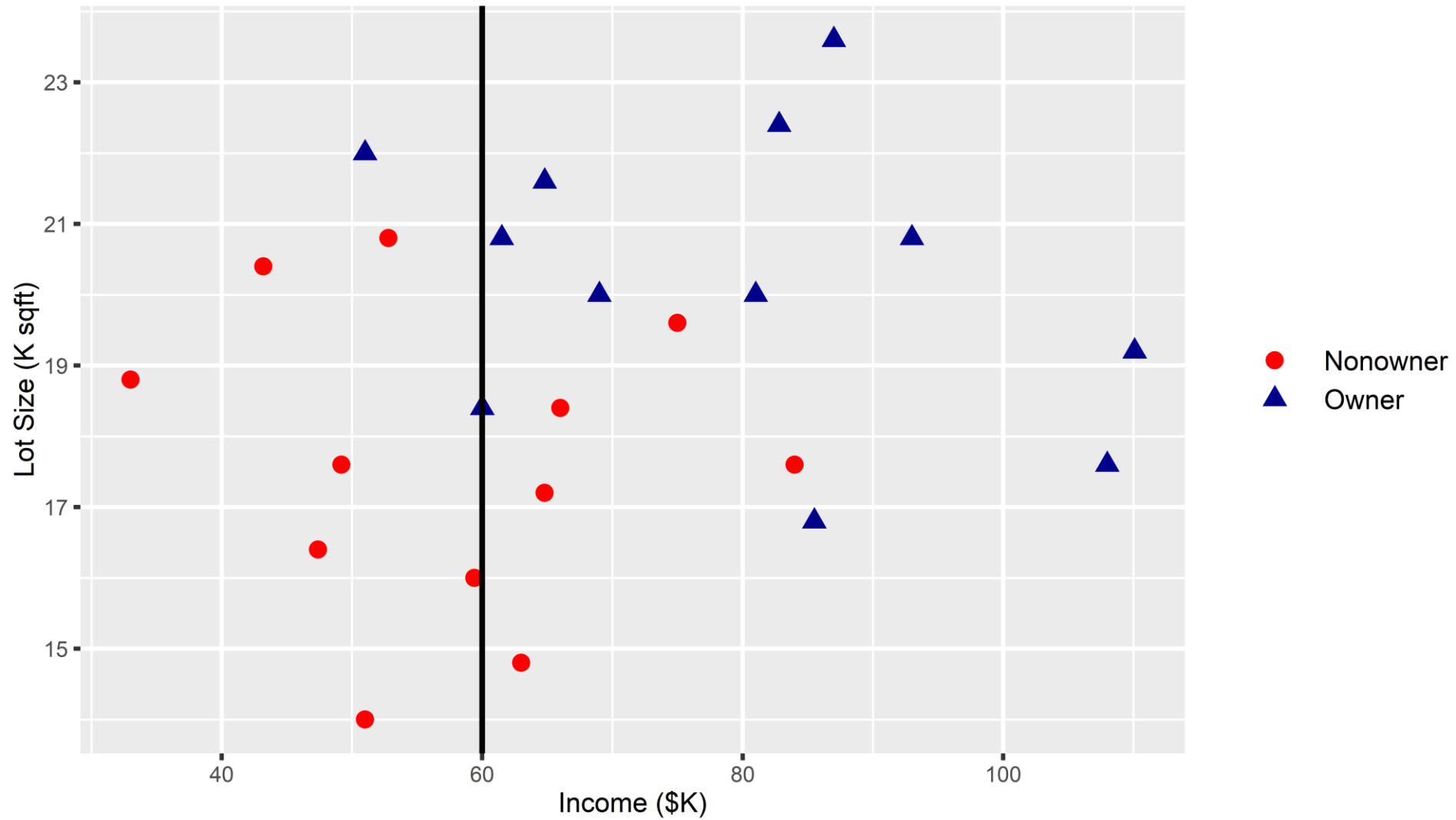
# Gini Index with no split

# Gini Index : No-split

- Gini impurity index for no-split

$$= 1 - \sum_{k=1}^{m} p_k^2$$

$$= 1 - \left( \frac{1}{2^2} + \frac{1}{2^2} \right) = \frac{1}{2}$$

- $p_k$ is the proportion of records in rectangle A that belong to class $k$

- Measure takes values between 0 (when all records belong to same class) and $^{m-1}/_m$ (when all m classes are equally represented)

# First split at Income = 60

# Gini Index at First split of Income = 60



- Gini impurity index for <u>left</u> rectangle

$$= 1 - \sum_{k=1}^{m} p_k^2 = 1 - \left( \frac{7^2}{8^2} + \frac{1^2}{8^2} \right) = 0.219$$

- Gini impurity index for <u>right</u> rectangle

$$= 1 - \sum_{k=1}^{m} p_k^2 = 1 - \left( \frac{11^2}{16^2} + \frac{5^2}{16^2} \right) = 0.430$$

- Weighted Index $= \frac{8}{24} * 0.219 + \frac{16}{24} * 0.430 = 0.359$

- No split to first split $\rightarrow$ 0.5 to 0.359

- This is the lowest drop that can be expected, hence the choice 60

# Measures of Impurity

- Entropy

- Entropy for a rectangle A is given by

$$\text{entropy(A)} = -\sum_{k=1}^{m} p_k * \log_2 p_k$$

- $p_k$ is the proportion of records in rectangle A that belong to class k

- Measure takes values between 0 (when all records belong to same class) and $\log_2(m)$(when all m classes are equally represented)

# Today's class mandatory steps

- Create a folder name "**n. classification_tree**" within the folder

  "**oba_455_555_ddpm_r/rproject**"

- Download "**classification_tree_code.R**", and all **csv** files from canvas

- Place all downloaded files in

  "**oba_455_555_ddpm_r /rproject/ n. classification_tree**"

- Open RStudio project

- Open "**classification_tree_code.R**" file within RStudio

# Example : Acceptance of Personal Loan

- Response : Bank customer accepting a loan (1) or not (0)

- Predictors (X)

  - Age, Experience, Income, Family Size, Education

  - Spending on Credit cards

  - Mortgage, Securities account
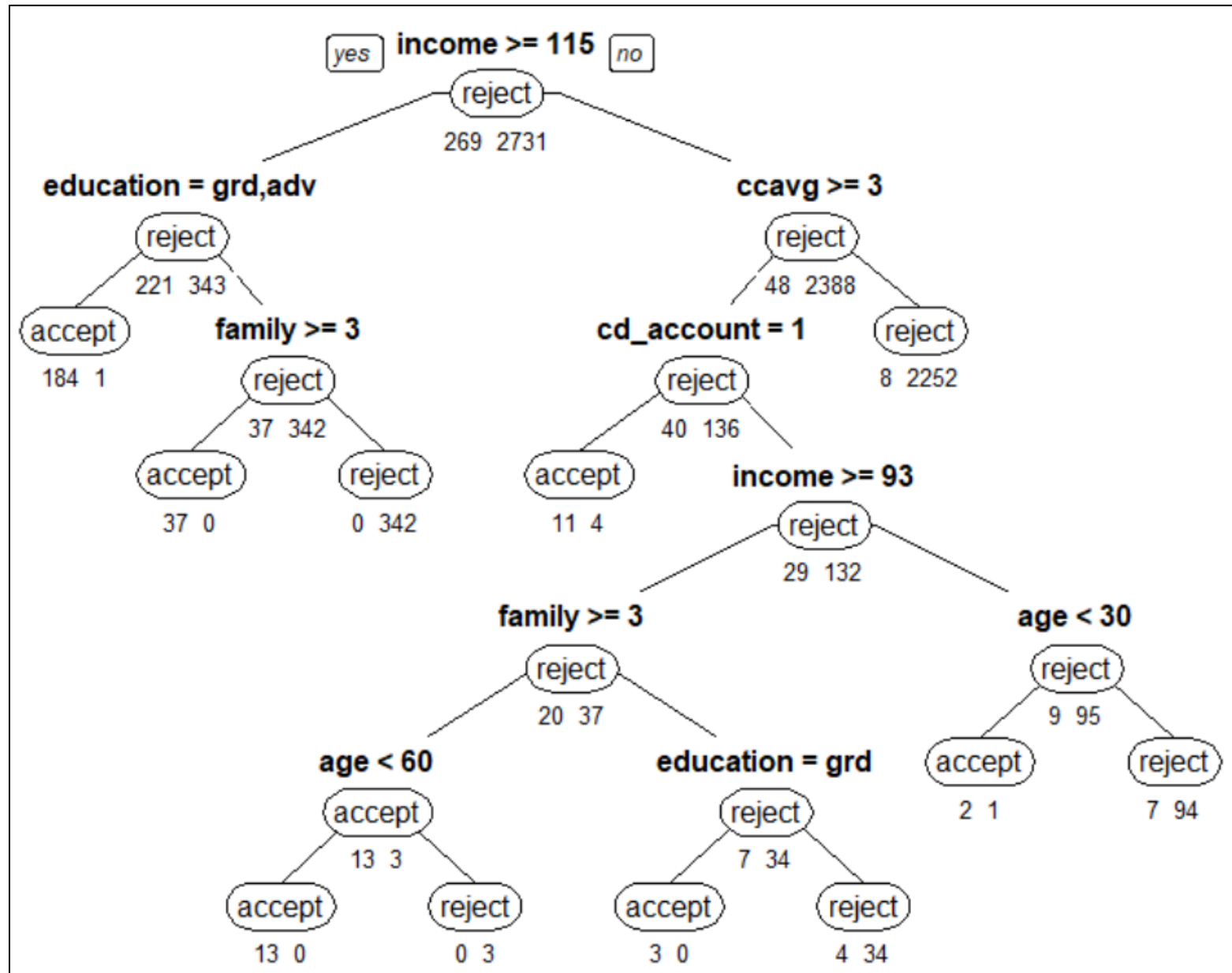
  - Online banking

  - ……

# Pruning a tree

- Step 1: Set the seed, create train and validation data

- Step 2: Run a tree with options cp = 0.00001, minsplit = 5 or 10, xval = 5 or 10

- Step 3: Plot the cp or relative error

- Step 4: Find "nsplit" value and its associated cp value from "size of the tree"

- Step 5: Prune the tree with the optimal cp

- Step 6: Predict the loan status for validation data.

- Step 7 : Develop confusion matrix and accuracy measures

# Pruning – Key variables

- Complexity parameter (cp)

    ➢ Any split that does not improve the fit by cp is not attempted

    ➢ Saves computing time by pruning off splits that are not worthwhile

- minsplit

    ➢ minimum number of observations that must exist in a node in order for a split to be attempted.

# Pruned tree

# Next class

- Regression Tree

# Thank You