# Process Mining and Simulation

## Assignment -1

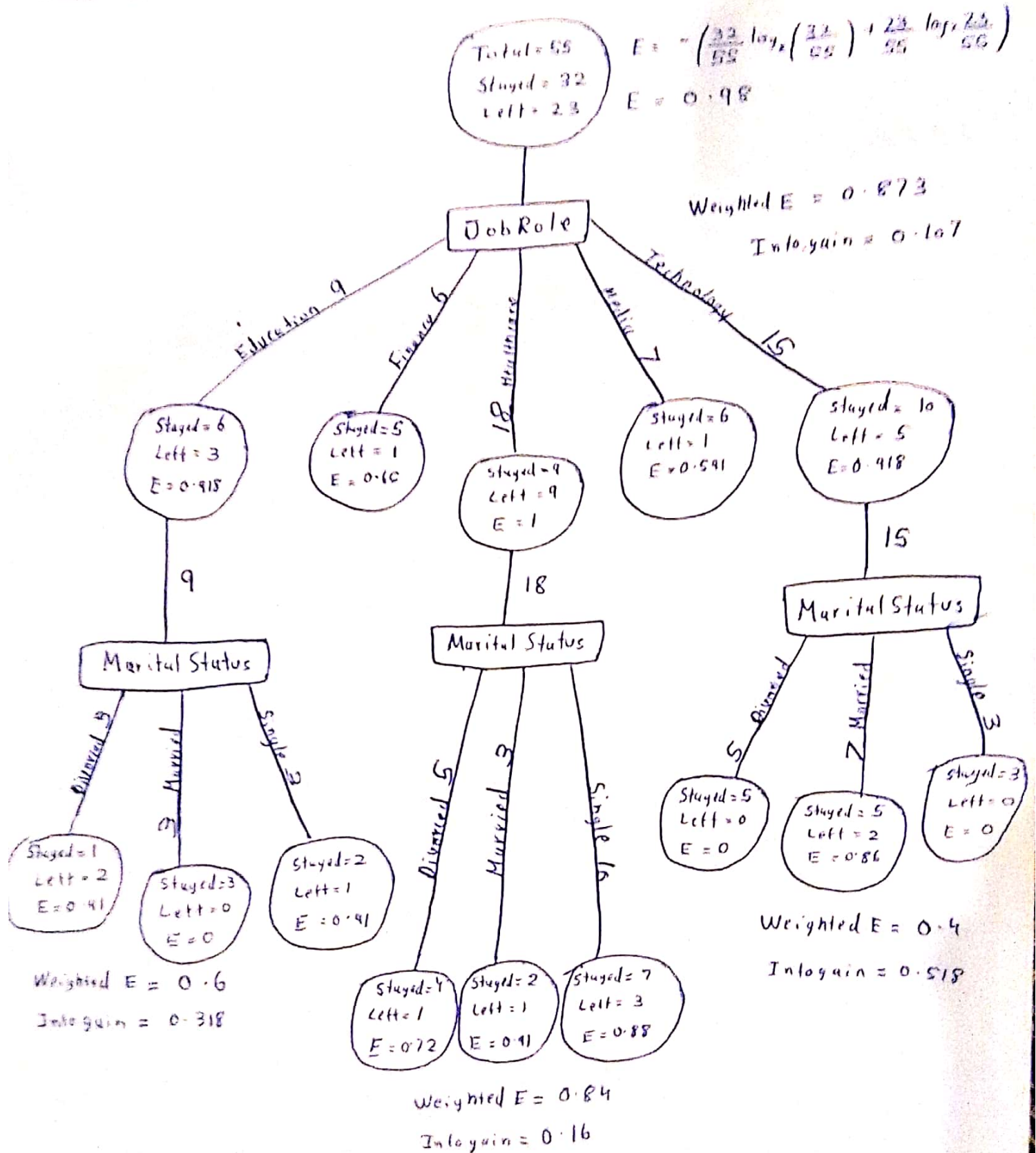Name: Abdullah Daoud

Roll No.: 22I-2626

Section: SE- E

Submitted to: Dr. Behjat Zuhaira

# Selecting Strong Attributes:

→ The rows provided to me for the dataset are 1261 - 1315

→ These rows include a total of 55 records and 24 variables.

→ The Response variable for this dataset is Attrition which has 2 possible values Stayed and Left

→ Strong attributes are those predictor variables that have the most effect on the Response variable

→ To identify the strong attributes from the given dataset, we need to calculate the information gained after a split based on specific variables.

→ To make the process intuitive, we will set the threshold of Information Gain to be 0.03 meaning that those variables whose information gain is 0.03 or above after one split will be considered as strong attributes.

→ After performing rigorous calculations, here are the strong attributes for the 1261-1315 dataset

→ Job Role = 0.107 information gain

→ Marital Status = 0.102 information gain

→ Company Size = 0.0852 information gain

→ Job Level = 0.05 information gain

→ Monthly Income = 0.05 information gain

→ Based on these strong attributes, 3 decision trees will be created

# Tree 1: JobRole → Marital Status



**Root node:** Total = 55, Stayed = 32, Left = 23

$$E = -\left(\frac{32}{55}\log_2\left(\frac{32}{55}\right) + \frac{23}{55}\log_2\frac{23}{55}\right)$$

$$E = 0.98$$

Weighted E = 0.873

Info gain = 0.107

**JobRole** branches:

- **Education 9** → Stayed = 6, Left = 3, E = 0.915
- **Finance 6** → Stayed = 5, Left = 1, E = 0.16
- **Healthcare 18** → Stayed = 9, Left = 9, E = 1
- **Media 7** → Stayed = 6, Left = 1, E = 0.591
- **Technology 15** → Stayed = 10, Left = 5, E = 0.918

**Education → 9 → Marital Status:**

- **Divorced 3** → Stayed = 1, Left = 2, E = 0.41
- **Married 3** → Stayed = 3, Left = 0, E = 0
- **Single 3** → Stayed = 2, Left = 1, E = 0.91

Weighted E = 0.6

Info gain = 0.318

**Healthcare → 18 → Marital Status:**

- **Divorced 6** → Stayed = 4, Left = 1, E = 0.72
- **Married 3** → Stayed = 2, Left = 1, E = 0.91
- **Single 16** → Stayed = 7, Left = 3, E = 0.88

Weighted E = 0.84

Info gain = 0.16

**Technology → 15 → Marital Status:**

- **Divorced 5** → Stayed = 5, Left = 0, E = 0
- **Married 7** → Stayed = 5, Left = 2, E = 0.86
- **Single 3** → Stayed = 3, Left = 0, E = 0

Weighted E = 0.4

Info gain = 0.518

→ Depth of tree is 2

→ To calculate information gain of the whole tree, we will first find out the weighted Average Entropy of the whole tree and then subtract it from the root node's Entropy

Weighted Average Entropy $= 3/55 \times 0.41 + 3/55 \times 0.41 + 6/55 \times 0.65$

$$+ 5/55 \times 0.72 + 3/55 \times 0.41 + 10/55 \times 0.88$$

$$+ 7/55 \times 0.541 + 7/55 \times 0.86$$

$$= 0.629$$

Information Gain $= 0.98 - 0.629 = 0.351$

$\rightarrow$ So information gain of Tree1 $= 0.351$

(P.T.O)

# Tree 2: Marital Status → Job Role



Total = 55
Stayed = 32
Left = 23   E = 0.48

**Marital Status**

Divorced 13 — Weighted E = 0.878   info gain = 0.102 — Single 22

(Divorced branch)
Stayed = 10
Left = 3
E = 0.77

(Married branch)
Stayed = 14
Left = 6
E = 0.88

(Single branch)
Stayed = 8
Left = 14
E = 0.94

---

**Job Role** (Divorced, 13)

3 Education — 5 Healthcare — 5 Technology

Stayed = 1
Left = 2
E = 0.91

Stayed = 4
Left = 1
E = 0.72

Stayed = 5
Left = 0
E = 0

Weighted E = 0.486
info gain = 0.284

---

**Job Role** (Married, 20)

Education 3 — Finance — Healthcare 3 — Medium 5 — Technology

Stayed = 3
Left = 0
E = 0

Stayed = 5
Left = 2
E = 0

Stayed = 2
Left = 1
E = 0.91

Stayed = 4
Left = 1
E = 0.72

Stayed = 5
Left = 1
E = 0.81

Weighted E = 0.6175
info gain = 0.2625

---

**Job Role** (Single, × 2)

Education 2 — Finance — Healthcare — Medium — Technology 2

Stayed = 2
Left = 1
E = 0.91

Stayed = 1
Left = 3
E = 0.81

Stayed = 3
Left = 1
E = 0.98

Stayed = 5
Left = 0
E = 0

Stayed = 0
Left = 3
E = 0

weighted E = 0.671
info gain = 0.269

---

weighted Average Entropy = 3/55 × 0.91 + 5/55 × 0.72 + 3/55 × 0.91

+ 5/55 × 0.72 + 7/55 × 0.88 + 3/55 × 0.91

+ 4/55 × 0.81 + 10/55 × 0.88

= 0.60

Information Gain = 0.48 − 0.60 = 0.38

# Tree 3: Company Size → Monthly Income → Job Level

**Root node:**
Total = 55
Stayed = 32
Left = 23
E = 0.98

**Company Size**

weighted E = 0.8948
info gain = 0.0852

Branches from Company Size:
- Large 12
- Medium 34
- Small 9

**Large (12):**
Stayed = 9
Left = 3
E = 0.811

**Medium (34):**
Stayed = 21
Left = 13
E = 0.459

**Small (9):**
Stayed = 2
Left = 7
E = 0.764

weighted E = 0.44
info gain = 0.324

---

**Monthly Income (12):**
weighted E = 0.78
info gain = 0.031

Branches:
- <7336
- ≥7336

**<7336:**
Stayed = 4
Left = 2
E = 0.91

**≥7336:**
Stayed = 5
Left = 1
E = 0.65

---

**Monthly Income (Small, 9):**
Branches:
- <8092
- ≥8092

**<8092:**
Stayed = 1
Left = 2
E = 1

**≥8092:**
Stayed = 0
Left = 5
E = 0

---

**Job Level (6):**
weighted E = 0.33
info gain = 0.58

Branches:
- Entry 1
- Mid 2
- Senior 3

**Entry:**
Stayed = 0
Left = 1
E = 0

**Mid:**
Stayed = 1
Left = 1
E = 1

**Senior:**
Stayed = 3
Left = 0
E = 0

---

**Monthly Income (34):**
weighted E = 0.93
info gain = 0.029

Branches:
- <7687
- ≥7687

**<7687:**
Stayed = 12
Left = 5
E = 0.87

**≥7687:**
Stayed = 9
Left = 8
E = 0.99

---

**Job Level (17):**
Weighted E = 0.767
info gain = 0.103

Branches:
- Entry 5
- Mid 6
- Senior

**Entry:**
Stayed = 3
Left = 2
E = 0.97

**Mid:**
Stayed = 6
Left = 3
E = 0.91

**Senior:**
Stayed = 3
Left = 0
E = 0

---

**Job Level (17):**
weighted E = 0.848
info gain = 0.044

Branches:
- Entry
- Mid
- Senior

**Entry:**
Stayed = 4
Left = 4
E = 1

**Mid:**
Stayed = 3
Left = 2
E = 0.97

**Senior:**
Stayed = 3
Left = 1
E = 0.81

→ Depth of tree is 3

Weighted Average Entropy = $2/55 \times 1 + 6/55 \times 0.65 + 5/55 \times 0.97$

$+ 9/55 \times 0.91 + 8/55 \times 1 + 5/55 \times 0.97$

$+ 4/55 \times 0.811 + 7/55 \times 1$

$= 0.709$

Information Gain = $0.98 - 0.709 = 0.271$

Choosing Optimal Tree:

→ The information gains for the 3 decision trees are:

Tree 1 = 0.351

Tree 2 = 0.38

Tree 3 = 0.271

→ The most optimal tree is the one which gives the highest information gain

→ High information gain means that tree or specific node has better prediction capabilities

→ Since Tree 2 has the highest information gain of all 3 trees, therefore

Tree 2 is the most optimal Tree.