# Measure of Impurity: Entropy

- Entropy at a given node t:

$$Entropy(t) = -\sum_j p(j \mid t)\log p(j \mid t)$$

(NOTE: $p(j/t)$ is the relative frequency of class j at node t).

- ◆ Maximum ($\log n_c$) when records are equally distributed among all classes implying least information
- ◆ Minimum (0.0) when all records belong to one class, implying most information

- – Entropy based computations are quite similar to the GINI index computations

# Computing Entropy of a Single Node

$$Entropy(t) = -\sum_j p(j \mid t)\log_2 p(j \mid t)$$

| | |
|---|---|
| C1 | **0** |
| C2 | **6** |

P(C1) = 0/6 = 0     P(C2) = 6/6 = 1

Entropy = – 0 log 0 – 1 log 1 = – 0 – 0 = 0

| | |
|---|---|
| C1 | **1** |
| C2 | **5** |

P(C1) = 1/6        P(C2) = 5/6

Entropy = – (1/6) $\log_2$ (1/6) – (5/6) $\log_2$ (1/6) = 0.65

| | |
|---|---|
| C1 | **2** |
| C2 | **4** |

P(C1) = 2/6        P(C2) = 4/6

Entropy = – (2/6) $\log_2$ (2/6) – (4/6) $\log_2$ (4/6) = 0.92

# Computing Information Gain After Splitting

l Information Gain:

$$GAIN_{split} = Entropy(p) - \left( \sum_{i=1}^{k} \frac{n_i}{n} Entropy(i) \right)$$
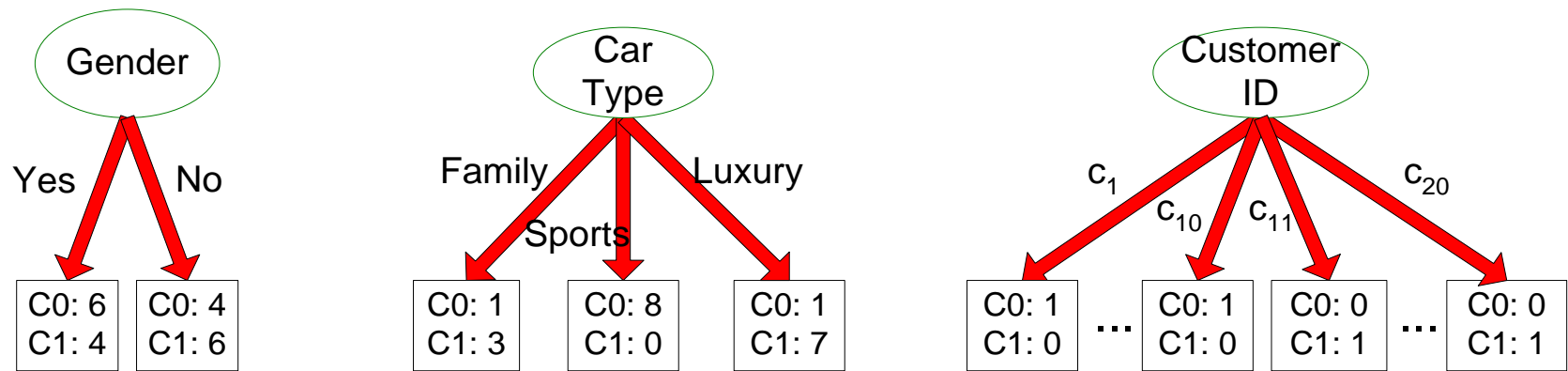
Parent Node, p is split into k partitions;

$n_i$ is number of records in partition i

– Choose the split that achieves most reduction (maximizes GAIN)

– Used in the ID3 and C4.5 decision tree algorithms

# Problem with large number of partitions

☐ Node impurity measures tend to prefer splits that result in large number of partitions, each being small but pure



Gender
Yes / No
C0: 6 / C1: 4    C0: 4 / C1: 6

Car Type
Family / Sports / Luxury
C0: 1 / C1: 3    C0: 8 / C1: 0    C0: 1 / C1: 7

Customer ID
$c_1$ ... $c_{10}$ $c_{11}$ ... $c_{20}$
C0: 1 / C1: 0    C0: 1 / C1: 0    C0: 0 / C1: 1    C0: 0 / C1: 1

– Customer ID has highest information gain because entropy for all the children is zero

# Gain Ratio

l  Gain Ratio:

$$GainRATIO_{split} = \frac{GAIN_{Split}}{SplitINFO}$$   $$SplitINFO = -\sum_{i=1}^{k} \frac{n_i}{n} \log \frac{n_i}{n}$$

   Parent Node, p is split into k partitions

   $n_i$ is the number of records in partition i

– Adjusts Information Gain by the entropy of the partitioning (SplitINFO).

   ◆ Higher entropy partitioning (large number of small partitions) is penalized!

– Used in C4.5 algorithm

– Designed to overcome the disadvantage of Information Gain

# Gain Ratio

l Gain Ratio:

$$GainRATIO_{split} = \frac{GAIN_{Split}}{SplitINFO} \qquad SplitINFO = -\sum_{i=1}^{k} \frac{n_i}{n} \log \frac{n_i}{n}$$

Parent Node, p is split into k partitions

$n_i$ is the number of records in partition i

| CarType | | |
|---|---|---|
| **Family** | **Sports** | **Luxury** |

| | Family | Sports | Luxury |
|---|---|---|---|
| **C1** | 1 | 8 | 1 |
| **C2** | 3 | 0 | 7 |
| **Gini** | | 0.163 | |

**SplitINFO = 1.52**

| CarType | |
|---|---|
| **{Sports, Luxury}** | **{Family}** |

| | {Sports, Luxury} | {Family} |
|---|---|---|
| **C1** | 9 | 1 |
| **C2** | 7 | 3 |
| **Gini** | | 0.468 |

**SplitINFO = 0.72**

| CarType | |
|---|---|
| **{Sports}** | **{Family, Luxury}** |

| | {Sports} | {Family, Luxury} |
|---|---|---|
| **C1** | 8 | 2 |
| **C2** | 0 | 10 |
| **Gini** | | 0.167 |

**SplitINFO = 0.97**