# Information Content

Based on the previous discussion we can easily prove the following lemma.

**Lemma 9.3: Information content**

If an object occurs with frequency $p$, then the most efficient way to represent it with $\log_2(1/p)$ bits.

**Example 9.9: Information content**

- A which occurs with frequency $\frac{1}{2}$ is represented by 1-bit, B which occurs with frequency $\frac{1}{4}$ represented by 2-bits and both C and D which occurs with frequency $\frac{1}{8}$ are represented by 3 bits each.

# Entropy Calculation

We can generalize the above understanding as follows.

- If there are $m$ objects with frequencies $p_1, p_2 \ldots\ldots, p_m$, then the average number of bits (i.e. questions) that need to be examined a value, that is, entropy is the frequency of occurrence of the $i^{th}$ value multiplied by the number of bits that need to be determined, summed up values of $i$ from $1$ to $m$.

---

**Theorem 9.4: Entropy calculation**

If $p_i$ denotes the frequencies of occurrences of $m$ distinct objects, then the entropy $E$ is

$$E = \sum_{i=1}^{m} p_i \log(^1/_{p_i}) \; and \sum_{i=1}^{m} p_i = 1$$

---

**Note:**

- If all are equally likely, then $p_i = \dfrac{1}{m}$ and $E = \log_2 m$; it is the special case.

# Entropy of a Training Set

- If there are $k$ classes $c_1, c_2 \ldots \ldots, c_k$ and $p_i$ for $i = 1 \ to \ k$ denotes the number of occurrences of classes $c_i$ divided by the total number of instances (i.e., the frequency of occurrence of $c_i$) in the training set, then entropy of the training set is denoted by

$$E = -\sum_{i=1}^{m} p_i \log_2 p_i$$

Here, $E$ is measured in "bits" of information.

**Note:**

- The above formula should be summed over the non-empty classes only, that is, classes for which $p_i \neq 0$

- $E$ is always a positive quantity

- $E$ takes it's minimum value (zero) if and only if all the instances have the same class (i.e., the training set with only **one** non-empty class, for which the probability 1).

- Entropy takes its maximum value when the instances are equally distributed among $k$ possible classes. In this case, the maximum value of $E$ is $log_2 \ k$.

# Entropy of a Training Set

**Example 9.10: OPTH dataset**

Consider the OTPH data shown in the following table with total 24 instances in it.

| Age | Eye sight | Astigmatic | Use Type | Class |
|-----|-----------|------------|----------|-------|
| 1 | 1 | 1 | 1 | 3 |
| 1 | 1 | 1 | 2 | 2 |
| 1 | 1 | 2 | 1 | 3 |
| 1 | 1 | 2 | 2 | 1 |
| 1 | 2 | 1 | 1 | 3 |
| 1 | 2 | 1 | 2 | 2 |
| 1 | 2 | 2 | 1 | 3 |
| 1 | 2 | 2 | 2 | 1 |
| 2 | 1 | 1 | 1 | 3 |
| 2 | 1 | 1 | 2 | 2 |
| 2 | 1 | 2 | 1 | 3 |
| 2 | 1 | 2 | 2 | 1 |
| 2 | 2 | 1 | 1 | 3 |
| 2 | 2 | 1 | 2 | 2 |
| 2 | 2 | 2 | 1 | 3 |
| 2 | 2 | 2 | 2 | 3 |
| 3 | 1 | 1 | 1 | 3 |
| 3 | 1 | 1 | 2 | 3 |
| 3 | 1 | 2 | 1 | 3 |
| 3 | 1 | 2 | 2 | 1 |
| 3 | 2 | 1 | 1 | 3 |
| 3 | 2 | 1 | 2 | 2 |
| 3 | 2 | 2 | 1 | 3 |
| 3 | 2 | 2 | 2 | 3 |

A coded forms for all values of attributes are used to avoid the cluttering in the table.

# Entropy of a training set

Specification of the attributes are as follows.

| Age | Eye Sight | Astigmatic | Use Type |
|---|---|---|---|
| 1: Young | 1: Myopia | 1: No | 1: Frequent |
| 2: Middle-aged | 2: Hypermetropia | 2: Yes | 2: Less |
| 3: Old | | | |

**Class:    1: Contact Lens   2:Normal glass    3: Nothing**

In the OPTH database, there are 3 classes and 4 instances with class 1, 5 instances with class 2 and 15 instances with class 3. Hence, entropy $E$ of the database is:

$$E = -\frac{4}{24}\log_2\frac{4}{24} - \frac{5}{24}\log_2\frac{5}{24} - \frac{15}{24}\log_2\frac{15}{24} = 1.3261$$

**Note:**

- The entropy of a training set implies the number of yes/no questions, on the average, needed to determine an unknown test to be classified.

- It is very crucial to decide the series of questions about the value of a set of attribute, which collectively determine the classification. Sometimes it may take one question, sometimes many more.

- Decision tree induction helps us to ask such a series of questions. In other words, we can utilize entropy concept to build a better decision tree.


**How entropy can be used to build a decision tree is our next topic of discussion.**

# Decision Tree Induction Techniques

- Decision tree induction is a top-down, recursive and divide-and-conquer approach.

- The procedure is to choose an attribute and split it into from a larger training set into smaller training sets.

- Different algorithms have been proposed to take a good control over

  1. Choosing the best attribute to be split, and

  2. Splitting criteria

- Several algorithms have been proposed for the above tasks. In this lecture, we shall limit our discussions into three important of them
  - **ID3**
  - **C 4.5**
  - **CART**

# Algorithm ID3

# ID3: Decision Tree Induction Algorithms

- Quinlan [1986] introduced the ID3, a popular short form of **I**terative **D**ichotomizer 3 for decision trees from a set of training data.

- In ID3, each node corresponds to a splitting attribute and each arc is a possible value of that attribute.

- At each node, the splitting attribute is selected to be the most informative among the attributes not yet considered in the path starting from the root.

# Algorithm ID3

- In ID3, entropy is used to measure how informative a node is.

  - It is observed that splitting on any attribute has the property that average entropy of the resulting training subsets will be less than or equal to that of the previous training set.

- ID3 algorithm defines a measurement of a splitting called **Information Gain** to determine the goodness of a split.

  - The attribute with the largest value of information gain is chosen as the splitting attribute and

  - it partitions into a number of smaller training sets based on the distinct values of attribute under split.

# Defining Information Gain

- We consider the following symbols and terminologies to define information gain, which is denoted as α.

- $D \equiv$ denotes the training set at any instant

- $|D| \equiv$ denotes the size of the training set $D$

- $E(D) \equiv$ denotes the entropy of the training set $D$

- The entropy of the training set $D$

$$E(D) = -\sum_{i=1}^{k} p_i \, log_2(p_i$$

  - where the training set $D$ has $c_1, c_2, \dots, c_k$, the $k$ number of distinct classes and

  - $p_i$, $0 < p_i \leq 1$ is the probability that an arbitrary tuple in $D$ belongs to class $c_i$ ($i = 1, 2, \dots, k$).

# Defining Information Gain

- $p_i$ can be calculated as

$$p_i = \frac{|C_{i,D}|}{|D|}$$

- where $C_{i,D}$ is the set of tuples of class $c_i$ in $D$.

- Suppose, we want to partition $D$ on some attribute $A$ having $m$ distinct values $\{a_1, a_2, \ldots, a_m\}$.

- Attribute $A$ can be considered to split $D$ into $m$ partitions $\{D_1, D_2, \ldots, D_m\}$, where $D_j$ ($j = 1, 2, \ldots, m$) contains those tuples in $D$ that have outcome $a_j$ of $A$.

# Defining Information Gain

> ## Definition 9.4: **Weighted Entropy**
>
> The weighted entropy denoted as $E_A(D)$ for all partitions of $D$ with respect to $A$ is given by:
>
> $$E_A(D) = \sum_{j=1}^{m} \frac{|D_j|}{|D|} E(D_j)$$
>
> Here, the term $\frac{|D_j|}{|D|}$ denotes the weight of the $j$-th training set.
>
> More meaningfully, $E_A(D)$ is the expected information required to classify a tuple from $D$ based on the splitting of $A$.

# Defining Information Gain

- Our objective is to take $A$ on splitting to produce an exact classification (also called pure), that is, all tuples belong to one class.

- However, it is quite likely that the partitions is impure, that is, they contain tuples from two or more classes.

- In that sense, $E_A(D)$ is a measure of impurities (or purity). A lesser value of $E_A(D)$ implying more power the partitions are.

---

### Definition 9.5: **Information Gain**

Information gain, $\alpha(A, D)$ of the training set $D$ splitting on the attribute $A$ is given by

$$\alpha(A, D) = E(D) - E_A(D)$$

In other words, $\alpha(A, D)$ gives us an estimation how much would be gained by splitting on $A$. The attribute $A$ with the highest value of $\alpha$ should be chosen as the splitting attribute for $D$.

# Information Gain Calculation

**Example 9.11 : Information gain on splitting OPTH**

- Let us refer to the OPTH database discussed in Slide #48.

- Splitting on **Age** at the root level, it would give three subsets $D_1$, $D_2$ and $D_3$ as shown in the tables in the following three slides.

- The entropy $E(D_1)$, $E(D_2)$ and $E(D_3)$ of training sets $D_1$, $D_2$ and $D_3$ and corresponding weighted entropy $E_{Age}(D_1)$, $E_{Age}(D_2)$ and $E_{Age}(D_3)$ are also shown alongside.

- The Information gain $\alpha$ $(Age, OPTH)$ is then can be calculated as **0.0394**.

- Recall that entropy of OPTH data set, we have calculated as $E(OPTH)$ = **1.3261**
  *(see Slide #49)*

# Information Gain Calculation

**Example 9.11 : Information gain on splitting OPTH**

Training set: $D_1 (\text{Age} = 1)$

| Age | Eye-sight | Astigmatism | Use type | Class |
|-----|-----------|-------------|----------|-------|
| 1 | 1 | 1 | 1 | 3 |
| 1 | 1 | 1 | 2 | 2 |
| 1 | 1 | 2 | 1 | 3 |
| 1 | 1 | 2 | 2 | 1 |
| 1 | 2 | 1 | 1 | 3 |
| 1 | 2 | 1 | 2 | 2 |
| 1 | 2 | 2 | 1 | 3 |
| 1 | 2 | 2 | 2 | 1 |

$$E(D_1) = -\frac{2}{8}log_2(\frac{2}{8}) - \frac{2}{8}log_2(\frac{2}{8}) - \frac{4}{8}log_2(\frac{4}{8}) = \mathbf{1.5}$$

$$E_{Age}(D_1) = \frac{8}{24} \times 1.5 = \mathbf{0.5000}$$

# Calculating Information Gain

Training set: $D_2(\text{Age} = 2)$

| Age | Eye-sight | Astigmatism | Use type | Class |
|-----|-----------|-------------|----------|-------|
| 2 | 1 | 1 | 1 | 3 |
| 2 | 1 | 1 | 2 | 2 |
| 2 | 1 | 2 | 1 | 3 |
| 2 | 1 | 2 | 2 | 1 |
| 2 | 2 | 1 | 1 | 3 |
| 2 | 2 | 1 | 2 | 2 |
| 2 | 2 | 2 | 1 | 3 |
| 2 | 2 | 2 | 2 | 3 |

$$E(D_2) = -\frac{1}{8}log_2\left(\frac{1}{8}\right) - \frac{2}{8}log_2\left(\frac{2}{8}\right) - \frac{5}{8}log_2\left(\frac{5}{8}\right)$$
$$= 1.2988$$

$$E_{Age}(D_2) = \frac{8}{24} \times 1.2988 = 0.4329$$

# Calculating Information Gain

Training set: $D_3(\text{Age} = 3)$

| Age | Eye-sight | Astigmatism | Use type | Class |
|-----|-----------|-------------|----------|-------|
| 3 | 1 | 1 | 1 | 3 |
| 3 | 1 | 1 | 2 | 3 |
| 3 | 1 | 2 | 1 | 3 |
| 3 | 1 | 2 | 2 | 1 |
| 3 | 2 | 1 | 1 | 3 |
| 3 | 2 | 1 | 2 | 2 |
| 3 | 2 | 2 | 1 | 3 |
| 3 | 2 | 2 | 2 | 3 |

$$E(D_3) = -\frac{1}{8}log_2(\frac{1}{8}) - \frac{1}{8}log_2(\frac{1}{8})$$
$$-\frac{6}{8}log_2(\frac{6}{8}) = 1.0613$$

$$E_{Age}(D_3) = \frac{8}{24} \times 1.0613 = 0.3504$$

$$\alpha\,(Age, D) = 1.3261 - (0.5000 + 0.4329 + 0.3504) = \mathbf{0.0394}$$

# Information Gains for Different Attributes

- In the same way, we can calculate the information gains, when splitting the OPTH database on Eye-sight, Astigmatic and Use Type. The results are summarized below.

- Splitting attribute: Age

$$\alpha(Age, OPTH) = 0.0394$$

- Splitting attribute: Eye-sight

$$\alpha(Eye-sight, OPTH) = 0.0395$$

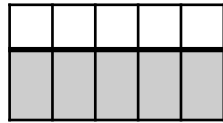- Splitting attribute: Astigmatic

$$\alpha(Astigmatic, OPTH) = 0.3770$$

- Splitting attribute: Use Type

$$\alpha(Use\ Type, OPTH) = 0.5488$$

# Decision Tree Induction : ID3 Way

- The ID3 strategy of attribute selection is to choose to split on the attribute that gives the greatest reduction in the weighted average entropy

  - The one that maximizes the value of information gain

- In the example with OPTH database, the larger values of information gain is $\alpha(\text{Use Type} , OPTH) = 0.5488$

  - Hence, the attribute should be chosen for splitting is "Use Type".

- The process of splitting on nodes is repeated for each branch of the evolving decision tree, and the final tree, which would look like is shown in the following slide and calculation is left for practice.

# Decision Tree Induction : ID3 Way

$OPTH$

**Age** ✗  **Eye-sight**  **Use Type** ✓  **Astigmatic** ✗

$\alpha = 0.395$   $\alpha = 0.0394$

$D1$  $Frequent(1)$   $D2$  $Less(2)$

$E(D_1) =?$

| Age | Eye | Ast | Use | Class |
|-----|-----|-----|-----|-------|
| 1 | 1 | 1 | 1 | 3 |
| 1 | 1 | 2 | 1 | 3 |
| 1 | 2 | 1 | 1 | 3 |
| 1 | 2 | 2 | 1 | 3 |
| 2 | 1 | 1 | 1 | 3 |
| 2 | 2 | 1 | 1 | 3 |
| 2 | 2 | 2 | 1 | 3 |
| 3 | 1 | 1 | 1 | 3 |
| 3 | 1 | 2 | 1 | 3 |
| 3 | 2 | 1 | 1 | 3 |
| 3 | 2 | 2 | 1 | 3 |

$E(D_2) =?$

| Age | Eye | Ast | Use | Class |
|-----|-----|-----|-----|-------|
| 1 | 1 | 1 | 2 | 2 |
| 1 | 1 | 2 | 2 | 1 |
| 1 | 2 | 1 | 2 | 2 |
| 1 | 2 | 2 | 2 | 1 |
| 2 | 1 | 1 | 2 | 2 |
| 2 | 1 | 2 | 2 | 1 |
| 2 | 2 | 1 | 2 | 2 |
| 3 | 1 | 1 | 2 | 3 |
| 3 | 1 | 2 | 2 | 3 |
| 3 | 2 | 1 | 2 | 2 |
| 3 | 2 | 2 | 2 | 3 |

Age  Eye-sight  Astigmatic   Age  Eye-sight  Astigmatic