**ENEE 420**
**FALL 2007**
**COMMUNICATIONS SYSTEMS**

**DATA COMPRESSION:**

**Finite sources** ─────────────────────────────

Let $\mathcal{X}$ denote a finite set, hereafter called the alphabet, and we refer to an element $x$ of $\mathcal{X}$ as a symbol. A probability mass function (pmf) $\boldsymbol{p} = (p(x),\ x \in \mathcal{X})$ on $\mathcal{X}$ is any collection of scalars indexed by $\mathcal{X}$ such that

$$0 < p(x) \le 1,\ x \in \mathcal{X} \quad \text{with} \quad \sum_{x \in \mathcal{X}} p(x) = 1.$$

A source is simply a pair $(\mathcal{X}, \boldsymbol{p})$ where $\mathcal{X}$ is a finite alphabet and $\boldsymbol{p}$ is a pmf on $\mathcal{X}$. It is sometimes convenient to refer to such a source by the notation $X = (\mathcal{X}, \boldsymbol{p})$ where the $\mathcal{X}$-valued random variable $X : \Omega \to \mathcal{X}$ is defined on some probability space $(\Omega, \mathcal{F}, \mathbb{P})$ such that

$$p(x) = \mathbb{P}\left[X = x\right], \quad x \in \mathcal{X}.$$

In short, we can think of $p(x)$ as the likelihood that the source generates symbol $x$.

**Divergence** ─────────────────────────────
The divergence between the pmfs $\boldsymbol{p}$ and $\boldsymbol{q}$ on $\mathcal{X}$ is defined by

$$D(\boldsymbol{p}\|\boldsymbol{q}) := -\sum_{x \in \mathcal{X}} p(x) \log_2 \left(\frac{q(x)}{p(x)}\right).$$

The basic bound

$$D(\boldsymbol{p}\|\boldsymbol{q}) \ge 0$$

holds with equality if and only if $\boldsymbol{p} = \boldsymbol{q}$.

**Entropy** ─────────────────────────────
For pmf $\boldsymbol{p}$ on $\mathcal{X}$, its (binary) entropy is defined by

$$H_2(\boldsymbol{p}) := -\sum_{x \in \mathcal{X}} p(x) \log_2 p(x).$$

The basic bounds

$$0 \leq H_2(\boldsymbol{p}) \leq \log_2 |\mathcal{X}|$$

holds, and we have

1. The lower bound is achieved if and only if the pmf $\boldsymbol{p}$ is degenerate, i.e.,

$$H_2(\boldsymbol{p}) = 0 \quad \text{if and only if } p(x) = 1 \quad \text{for some } x \in \mathcal{X};$$

2. The upper bound is achieved if and only if the pmf $\boldsymbol{p}$ is the uniform pmf on $\mathcal{X}$, i.e.,

$$H_2(\boldsymbol{p}) = \log_2 |\mathcal{X}| \quad \text{if and only if } p(x) = \frac{1}{|\mathcal{X}|}, \quad x \in \mathcal{X}.$$

**Compression codes**

Let $\mathcal{B}^\star$ denote the collection of all binary words, i.e.,

$$\mathcal{B}^\star = \cup_{n=1}^\infty \{0, 1\}^n.$$

A binary compression code, hereafter simply a code, for a $\mathcal{X}$-valued source is any mapping

$$C : \mathcal{X} \to \mathcal{B}^\star.$$

For each $x$ in $\mathcal{X}$, $C(x)$ is known as the codeword associated with $x$ under $C$. It is customary to refer to the collection $\{C(x), \ x \in \mathcal{X}\}$ of all codewords as the codebook for $C$, and to identify it with $C$.

Some terminology: A code $C : \mathcal{X} \to \mathcal{B}^\star$ is said to be

1. non-singular if $C(x) \neq C(y)$ for any pair of distinct symbols $x, y$ in $\mathcal{X}$;

2. uniquely decipherable if the equality

$$C(x_1) \ldots C(x_n) = C(y_1) \ldots C(y_m)$$

for some $x_1, \ldots, x_n, y_1, \ldots, y_m$ in $\mathcal{X}$ implies

$$n = m \quad \text{and} \quad x_j = y_j, \ j = 1, \ldots, n.$$

   3. prefix (or to have the prefix property) if for any symbol $x$ in $\mathcal{X}$, no prefix of $C(x)$ is a codeword for some other symbol in $\mathcal{X}$.

Prefix codes are also known as instantaneous codes. We denote the collection of all prefix codes by $\mathcal{C}_{\mathrm{Pref}}$.

**Length of codes** _____

Given a code $C : \mathcal{X} \to \mathcal{B}^\star$, let $\ell_C(x)$ denote the length of the binary codeword $C(x)$ associated with the symbol $x$ in $\mathcal{X}$. Given a source $X = (\mathcal{X}, \boldsymbol{p})$, the expected codeword length of a code $C : \mathcal{X} \to \mathcal{B}^\star$ is given by

$$
\begin{aligned}
L(C; \boldsymbol{p}) &:= \mathbb{E}\left[\ell_C(X)\right] \\
&= \sum_{x \in \mathcal{X}} \ell_C(x) p(x).
\end{aligned}
$$
(1)

**Kraft Inequality** _____
For any prefix code $C : \mathcal{X} \to \mathcal{B}^\star$, we have

$$
\sum_{x \in \mathcal{X}} 2^{-\ell_C(x)} \le 1.
$$

Conversely, for any collection $(\ell(x), \ x \in \mathcal{X})$ of positive integers such that

$$
\sum_{x \in \mathcal{X}} 2^{-\ell(x)} \le 1,
$$

there exists a prefix code $C : \mathcal{X} \to \mathcal{B}^\star$ such that

$$
\ell_C(x) = \ell(x), \quad x \in \mathcal{X}.
$$

**Shannon encoding** _____
Set

$$
\ell_{\mathrm{SH}}(x) = \lceil \log_2 \frac{1}{p(x)} \rceil, \quad x \in \mathcal{X}.
$$

Since $2^{\log_2 t} = t$ for all $t > 0$, we find

$$
\begin{aligned}
\sum_{x \in \mathcal{X}} 2^{-\ell_{\mathrm{SH}}(x)} &\le \sum_{x \in \mathcal{X}} 2^{-\log_2 \frac{1}{p(x)}} \\
&= \sum_{x \in \mathcal{X}} 2^{\log_2 p(x)} \\
&= \sum_{x \in \mathcal{X}} p(x) = 1,
\end{aligned}
$$
(2)

and there exists a prefix code $C_{\mathrm{SH}} : \mathcal{X} \to \mathcal{B}^\star$ such that

(3) $$\ell_{C_{\mathrm{SH}}}(x) = \ell_{\mathrm{SH}}(x), \quad x \in \mathcal{X}.$$

Any code satsifying (3) is known as Shannon encoding.

Note that

$$
\begin{aligned}
L(C_{\mathrm{SH}}; \boldsymbol{p}) &= \sum_{x \in \mathcal{X}} p(x) \ell_{\mathrm{SH}}(x) \\
&\leq \sum_{x \in \mathcal{X}} p(x) \left( \log_2 \frac{1}{p(x)} + 1 \right) \\
&= -\sum_{x \in \mathcal{X}} p(x) \log_2 p(x) + \sum_{x \in \mathcal{X}} p(x) \\
&= H_2(\boldsymbol{p}) + 1.
\end{aligned}
$$

(4)

Shannon encoding comes from within one bit of source entropy!

**Average code length and entropy** ———————————————
Consider a prefix code $C : \mathcal{X} \to \mathcal{B}^\star$. Introduce the pmf $\boldsymbol{q}_C$ on $\mathcal{X}$ given by

$$q_C(x) = \frac{2^{-\ell_C(x)}}{\Sigma(C)}, \quad x \in \mathcal{X}$$

where

$$\Sigma(C) = \sum_{x \in \mathcal{X}} 2^{-\ell_C(x)}.$$

We have

(5) $$L(C; \boldsymbol{p}) - H_2(\boldsymbol{p}) = D(\boldsymbol{p} \| \boldsymbol{q}_C) + \log_2 \left( \frac{1}{\Sigma(C)} \right)$$

so that

$$L(C; \boldsymbol{p}) \geq H_2(\boldsymbol{p})$$

since $D(\boldsymbol{p} \| \boldsymbol{q}_C) \geq 0$ and $\Sigma(C) \leq 1$ by Kraft inequality. Equality holds if and only if $D(\boldsymbol{p} \| \boldsymbol{q}_C) = 0$ and $\Sigma(C) = 1$. In other words, equality holds if and only if there exists positive integers $(n(x), \ x \in \mathcal{X})$ such that

$$p(x) = 2^{-n(x)}, \quad x \in \mathcal{X}.$$

**A proof of (5)** ———————————————

$$
\begin{aligned}
L(C; \boldsymbol{p}) &= \sum_{x \in \mathcal{X}} \ell_C(x) p(x) \\
&= -\sum_{x \in \mathcal{X}} p(x) \log_2 \left( 2^{-\ell_C(x)} \right) \\
&= -\sum_{x \in \mathcal{X}} p(x) \log_2 \left( \frac{2^{-\ell_C(x)}}{\Sigma(C)} \cdot \Sigma(C) \right) \\
&= -\sum_{x \in \mathcal{X}} p(x) \log_2 \left( \frac{q_C(x)}{p(x)} \cdot p(x) \Sigma(C) \right) \\
&= -\sum_{x \in \mathcal{X}} p(x) \left( \log_2 \left( \frac{q_C(x)}{p(x)} \right) + \log_2 p(x) + \log_2 \Sigma(C) \right) \\
&= -\sum_{x \in \mathcal{X}} p(x) \log_2 \left( \frac{q_C(x)}{p(x)} \right) - \sum_{x \in \mathcal{X}} p(x) \log_2 p(x) - \log_2 \Sigma(C).
\end{aligned}
$$

**Source coding Theorem (Shannon 1948)**
The bounds
(6)
$$ H_2(\boldsymbol{p}) \leq L_{\min}(\boldsymbol{p}) \leq H_2(\boldsymbol{p}) + 1 $$

hold where
$$ L_{\min}(\boldsymbol{p}) := \min \left( L(C; \boldsymbol{p}) : \ C \in \mathcal{C}_{\mathrm{Pref}} \right). $$

Moreover,
$$ L_{\min}(\boldsymbol{p}) = H_2(\boldsymbol{p}) $$

if and only if there exists positive integers $(n(x), \ x \in \mathcal{X})$

$$ p(x) = 2^{-n(x)}, \quad x \in \mathcal{X}. $$

**The more likely the symbol, the shorter its description**
Consider a prefix code $C : \mathcal{X} \to \mathcal{B}^\star$. Define a new code $C' : \mathcal{X} \to \mathcal{B}^\star$ as follows:
Pick distinct $x$ and $y$ in $\mathcal{X}$, and set

$$
C'(z) = \begin{cases}
C(z) & \text{if } z \neq x, y \\
C(y) & \text{if } z = x \\
C(x) & \text{if } z = y
\end{cases}
$$

Obviously,

$$
\ell_{C'}(z) = \begin{cases} \ell_C(z) & \text{if } z \neq x, y \\ \ell_C(y) & \text{if } z = x \\ \ell_C(x) & \text{if } z = y \end{cases}
$$

so that

$$
\begin{aligned}
L(C; \boldsymbol{p}) - L(C'; \boldsymbol{p}) &= \sum_{z \in \mathcal{X}} \ell_C(z) p(z) - \sum_{z \in \mathcal{X}} \ell_{C'}(z) p(z) \\
&= (\ell_C(x) p(x) + \ell_C(y) p(y)) - (\ell_C(y) p(x) + \ell_C(x) p(y)) \\
&= (\ell_C(x) - \ell_C(y)) p(x) + (\ell_C(y) - \ell_C(x)) p(y) \\
&= (\ell_C(x) - \ell_C(y)) (p(x) - p(y)).
\end{aligned}
$$

In short, if $p(y) < p(x)$, then $L(C; \boldsymbol{p}) \leq L(C'; \boldsymbol{p})$ if and only if $\ell_C(x) \leq \ell_C(y)$.

**Reduction step behind Huffman encoding** ⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯

Consider a code $C : \mathcal{X} \rightarrow \mathcal{B}^\star$ with the following property: There exist distinct symbols $x$ and $y$ in $\mathcal{X}$ such that their codewords differ only in their last bit, i.e., for some $\ell = 1, 2, \ldots$, we have

$$
C(x) = (b_1, \ldots, b_\ell, 1) \quad \text{and} \quad C(y) = (b_1, \ldots, b_\ell, 0)
$$

with $b_1, \ldots, b_\ell$ in $\{0, 1\}$.

With the source $X = (\mathcal{X}, \boldsymbol{p})$, we associate a new source $X' = (\mathcal{X}', \boldsymbol{p}')$ as follows: The new alphabet $\mathcal{X}'$ is obtained by combining the two symbols $x$ and $y$, i.e.,

$$
\mathcal{X}' := (\mathcal{X} - \{x, y\}) \cup \{\star\}
$$

where $\star$ denotes the new symbol obtained by combining $x$ and $y$. Next, the pmf $\boldsymbol{p}'$ on $\mathcal{X}'$ is naturally derived from $\boldsymbol{p}$, namely

$$
p'(z) = \begin{cases} p(z) & \text{if } z \neq x, y \\ p(x) + p(y) & \text{if } z = \star. \end{cases}
$$

With $C$ we now associate a new code $C' : \mathcal{X}' \rightarrow \mathcal{B}^\star$ for this new source $X' = (\mathcal{X}', \boldsymbol{p}')$ given by

$$
C'(z) = \begin{cases} C(z) & \text{if } z \neq x, y \\ (b_1, \ldots, b_\ell) & \text{if } z = \star. \end{cases}
$$

Therefore,

$$\ell_{C'}(z) = \begin{cases} \ell_C(z) & \text{if } z \neq x, y \\ \\ \ell & \text{if } z = \star. \end{cases}$$

With these definitions,

$$
\begin{aligned}
L(C', \boldsymbol{p}') &= \sum_{z \in \mathcal{X}'} \ell_{C'}(z) p'(z) \\
&= \sum_{z \in \mathcal{X} - \{x,y\}} \ell_{C'}(z) p'(z) + \ell_{C'}(\star) p'(\star) \\
&= \sum_{z \in \mathcal{X} - \{x,y\}} \ell_C(z) p(z) + \ell \, (p(x) + p(y)) \\
&= \sum_{z \in \mathcal{X} - \{x,y\}} \ell_C(z) p(z) + \ell p(x) + \ell p(y) \\
&= \sum_{z \in \mathcal{X} - \{x,y\}} \ell_C(z) p(z) + (\ell_C(x) - 1) \, p(x) + (\ell_C(y) - 1) \, p(y) \\
&= \sum_{z \in \mathcal{X}} \ell_C(z) p(z) - (p(x) + p(y)) \,.
\end{aligned}
$$

In short,

(7) $$L(C', \boldsymbol{p}') = L(C, \boldsymbol{p}) - (p(x) + p(y)) \,.$$

**Properties of optimal prefix codes** ⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯⎯

For notational convenience, assume that the symbols in the alphabet $\mathcal{X}$ are relabeled so that

$$p(M) \leq p(M-1) \leq \ldots \leq p(2) \leq p(1)$$

with $|\mathcal{X}| = M$.

1. If a (prefix) code $C : \mathcal{X} \to \mathcal{B}^\star$ is optimal, then necessarily

$$\ell_C(1) \leq \ell_C(2) \leq \ldots \leq \ell_C(M-1) \leq \ell_C(M)$$

2. If the prefix code $C : \mathcal{X} \to \mathcal{B}^\star$ is optimal, then necessarily

$$\ell_C(M-1) = \ell_C(M)$$

3. The optimal prefix code $C : \mathcal{X} \to \mathcal{B}^\star$ can always be selected so that $C(M - 1)$ and $C(M)$ differ only in the last bit, i.e., if $C(M - 1) = (a_1, \ldots, a_\ell)$ and $C(M) = (b_1, \ldots, b_\ell)$ where $\ell = \ell_C(M - 1) = \ell_C(M)$, then

$$a_k = b_k, \quad k = 1, \ldots, \ell - 1.$$