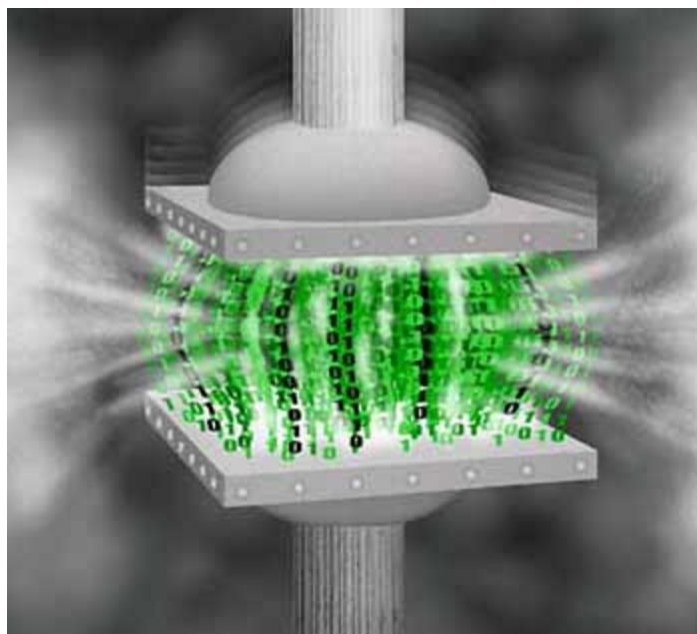


第2章 数据无损压缩





目录

2.0 简介

2.1 信息论基础知识

信息量、熵、信源编码、决策量、编码冗余

2.2 统计编码

2.2.1 香农-范诺编码

2.2.2 霍夫曼编码

2.2.3 算术编码

2.3 游程编码

2.4 词典编码

2.4.1 词典编码的思想

2.4.2 LZ77算法

2.4.3 LZSS算法

2.4.4 LZ78算法

2.4.5 LZW算法



目录

2.0 简介

2.1 信息论基础知识

信息量、熵、信源编码、决策量、编码冗余

2.2 统计编码

2.2.1 香农-范诺编码

2.2.2 霍夫曼编码

2.2.3 算术编码

2.3 游程编码

2.4 词典编码

2.4.1 词典编码的思想

2.4.2 LZ77算法

2.4.3 LZSS算法

2.4.4 LZ78算法

2.4.5 LZW算法



什么是压缩？

■ 压缩：

减少表达一个数字信号（例如声音、图像、视频、或其他信号）所用比特数，同时

- 能够完全重建出原始数据（无损压缩） 或
- 恢复出质量可被接受的原始数据（有损压缩）

为什么要压缩？

Telephony (220-3400 Hz)	$8000 \text{ samples/s} \times 12 \text{ b/sample} = 96 \text{ kbps}$
Wideband speech (50-7000 Hz)	$16,000 \text{ samples/s} \times 14 \text{ b/sample} = 224 \text{ kbps}$
Wideband audio (20-20,000 Hz)	$44,100 \text{ samples/s} \times 2 \text{ channels} \times 16 \text{ b/sample} = 1.412 \text{ Mbps}$
Color images	$512 \times 512 \text{ pixels} \times 24 \text{ bpp} = 6.3 \text{ Mbits} = 786 \text{ Kbytes}$
Video (CCIR601)	$720 \times 480 \text{ gray pixels} \times 8 \text{ bpp} \times 30 \text{ frames/s} + 2 \times 360 \times 480 \text{ chroma pixels} \times 8 \text{ bpp} \times 30 \text{ frames/s} = 166 \text{ Mbps}$
HDTV	$1280 \times 720 \text{ color pixels} \times 24 \text{ bpp} \times 60 \text{ frames/s} = 1.3 \text{ Gbps}$

CCIR601: Without compression we could store

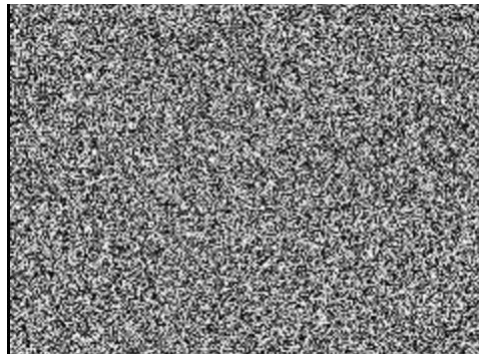
- 31 s of video on a CD-ROM (650 Mbytes)
- 4 minutes of video on a DVD-5 (40 Gbits)

HDTV Terrestrial Broadcasting:

- 19.3 Mbits/s (6 MHz channel)
- need compression ratio ~ 70

为什么信号能被压缩？

- **数据冗余：**统计冗余或数据中出现的结构（空域内、时域内、频域内）



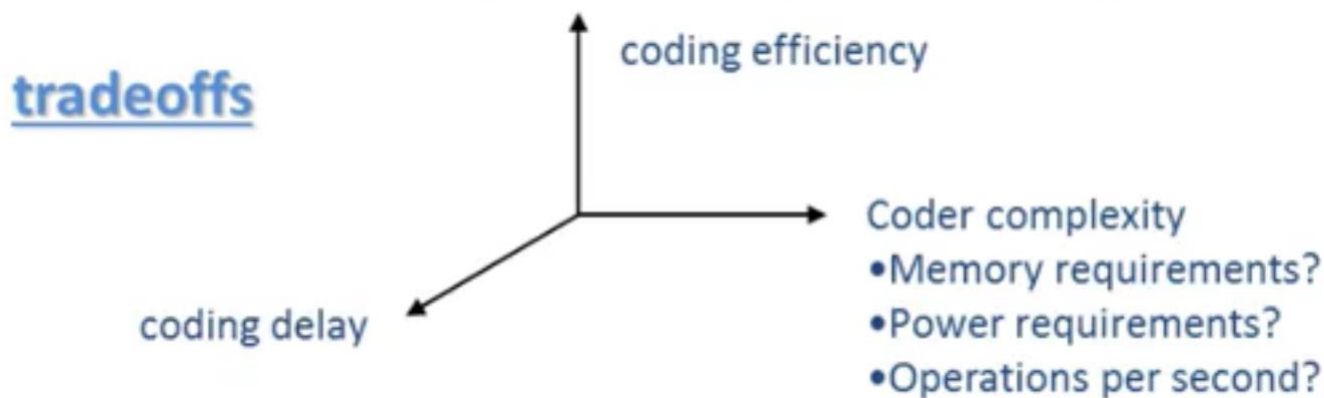
- **视听冗余：**感知上多余的信息
 - 听觉系统的敏感度有限
 - 视觉系统的敏感度有限

■ 三种多媒体数据类型

- 文字 (text)数据——无损压缩
 - 根据数据本身的冗余(Based on data redundancy)
- 声音(audio)数据——有损压缩
 - 根据数据本身的冗余(Based on data redundancy)
 - 根据人的听觉系统特性(Based on human hearing system)
- 图像(image)/视频(video) 数据——有损压缩
 - 根据数据本身的冗余(Based on data redundancy)
 - 根据人的视觉系统特性(Based on human visual system)

无损压缩

- **目标：** 由压缩数据能精确重建原始数据，即过程可逆。
- **压缩比由信源的熵来决定的**
 - 无损压缩的压缩比相对较低
- **应用举例：** 文字、医学图像、有损编码器



■ 数据无损压缩的方法

- 香农范诺编码(Shannon-Fano Coding)
- 霍夫曼编码(Huffman coding)
- 算术编码(arithmetic coding)
- 游程编码(run-length coding)
- 词典编码(dictionary coding)



目录

2.0 简介

2.1 信息论基础知识

决策量、信息量、熵、信源
编码、编码冗余

2.2 统计编码

2.2.1 香农-范诺编码

2.2.2 霍夫曼编码

2.2.3 算术编码

2.3 游程编码

2.4 词典编码

2.4.1 词典编码的思想

2.4.2 LZ77算法

2.4.3 LZSS算法

2.4.4 LZ78算法

2.4.5 LZW算法

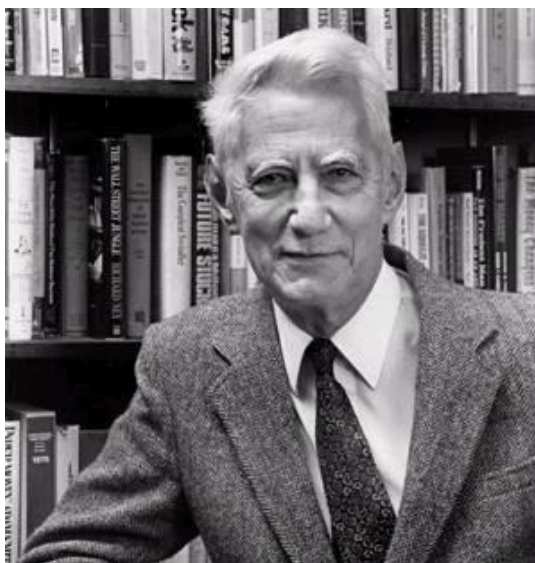


2.1 信息论基础知识

- 数据无损压缩的理论——信息论(information theory)
 - 1948年创建的数学理论的一个分支学科，研究信息的编码、传输和存储
 - Claude Shannon (香农)发表的“A Mathematical Theory of Communication”论文题目，提议用二进制数据对信息进行编码
 - 最初只应用于通信工程领域，后来扩展到包括计算在内的其他多个领域，如信息的存储、信息的检索等。在通信方面，主要研究数据量、传输速率、信道容量、传输正确率等问题。

2.1 信息论基础知识

- 信息论之父 Claude Elwood Shannon (1916-2001)



➤ In 1948:

- *A mathematical theory of communication*, landmark, climax

(Shannon理论的一个重要贡献: 熵的概念)



2.1 信息论基础知识

- 信息产生的过程可以看成是一个信源发出一串随机得从有限符号集里挑选的符号
 - 例如：手写英文；n-bit图像
- 最简单情况：离散无记忆信源（**discrete memoryless source, DMS**）
 - 信源相继产生的符号在统计上是独立的,且服从同一分布

熵 (Entropy)

■ 自信息量(self information):

$$I(s_i) = \log \frac{1}{p_i} = -\log p_i$$

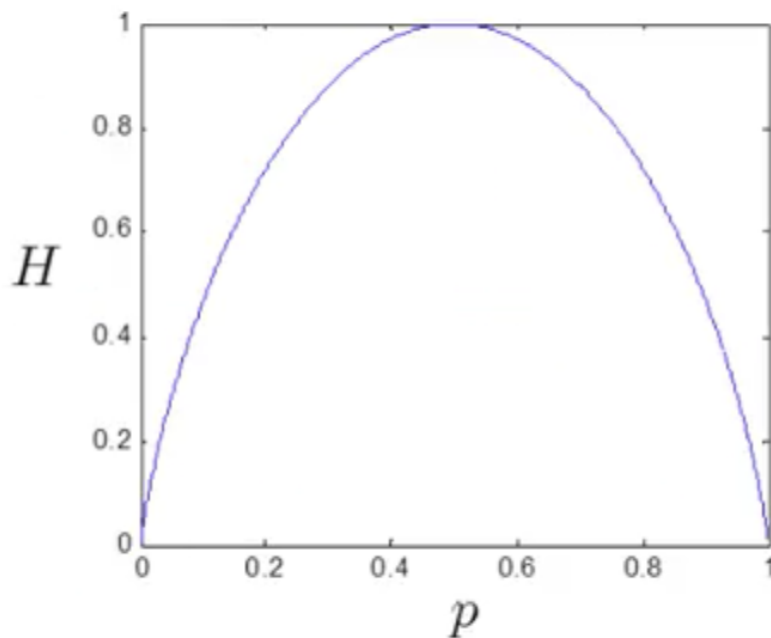
- 一个概率较小的事件发生能提供更多的信息量
- 多个独立事件作为一个事件产生的信息量相当于各个事件自信息量的总和

■ DMS的熵 (entropy): 信源符号的信息量的平均值

$$H(S) = \sum_{i=1}^n p_i I(s_i) = -\sum_{i=1}^n p_i \log_2 p_i$$

熵-示例

■ 离散无记忆信源的熵



$$H = -p \log_2 p - (1 - p) \log_2 (1 - p)$$



信源的熵

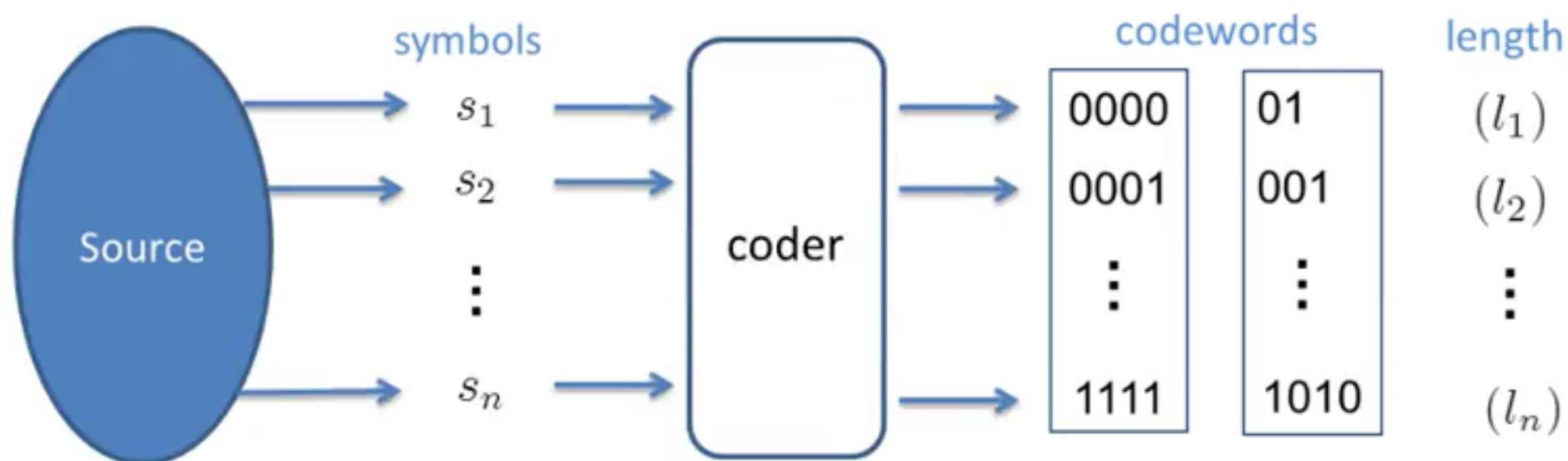
■ 英语的熵

- DMS, 相同概率: $H(S) = 4.70\text{bits/letter}$
- DMS, 实际概率: $H(S) = 4.14 \text{ bits/letter}$

■ 利用数据的结构

- 例子1: 考虑数据: 1234323456
- 例子2: 考虑数据: 334533334545

编码



Alphabet: $A = \{s_1, s_2, \dots, s_n\}$

$\downarrow \quad \downarrow \quad \downarrow$

$p_1 \quad p_2 \quad p_n$

平均码字长度: $l_{avg} = \sum_{i=1}^n l_i p_i$

信源编码定理

- S 为一个符号集大小为 n 、熵为 $H(S)$ 的信源。考虑把 N 个信源符号编成二进制的码字。对于任意的 $\delta > 0$ ，都可以通过选择足够大的 N 来构建码字，使得原始符号的平均比特数 l_{avg} 满足：

$$H(S) \leq l_{avg} \leq H(S) + \delta$$

例子

Symbols	Probability	Code 1	Code 2	Code 3	Code 4
s_1	1/2	0	0	0	0
s_2	1/4	0	1	10	01
s_3	1/8	1	00	110	011
s_4	1/8	10	11	111	0111
Average length		1.125	1.25	1.75	1.875

Entropy $H(S)=1.75$ bits/symbol

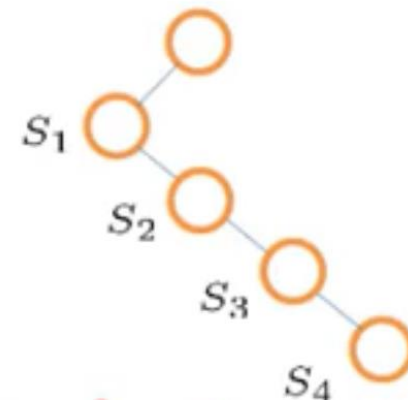
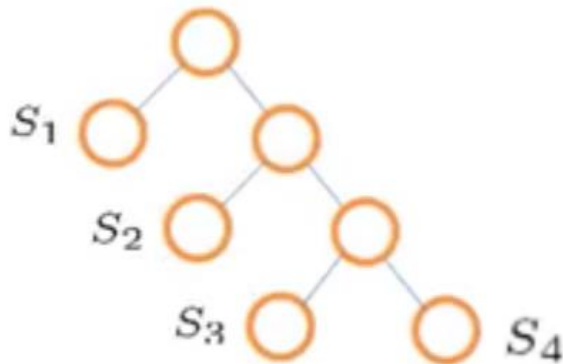
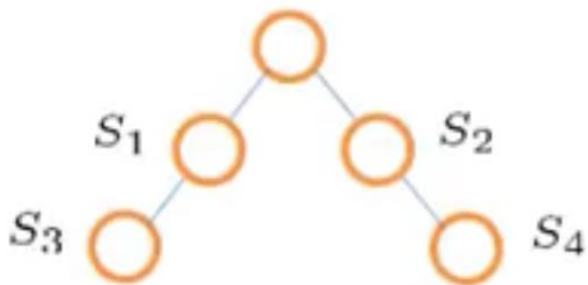
- 可唯一解码：任何一段有限码字序列对应于唯一一个信息序列
- 前缀码字：每个码字都不是其他码字的前缀

码字的二叉树表示

Code 2
0
1
00
11

Code 3
0
10
110
111

Code 4
0
01
011
0111



2.1 信息论基础知识

■ 决策量(decision content)

- 决策量是有限字符集大小的对数值
- 在数学上表示为

$H_0 = \log(n)$ 其中, n 是有限字符集的字符数

- 决策量的单位由对数的底数决定
 - Bit或 Sh (Shannon): 用于以2为底的对数
 - Nat (natural unit): 用于以e为底的对数
 - Hart (hartley): 用于以10为底的对数

2.1 信息论基础知识

■ 数据的冗余量

在信息论中，数据的冗余量(R)定义为决策量(H_0)超过熵(H)的量，数学上表示为

$$R = H_0 - H$$

例如，令 $\{a, b, c\}$ 为 3 个事件构成的数据集，它们出现的概率分别为 $p(a) = 0.5$,

$p(b) = 0.25$, $p(c) = 0.25$ ，这个数据集的冗余量则为，

$$R = H_0 - H = \log_2 3 - \left[-\sum_{i=1}^n p(x_i) \log_2 p(x_i) \right] = 1.58 - 1.50 = 0.08 \quad (\text{Sh})$$



目录

2.0 简介

2.1 信息论基础知识

信息量、熵、信源编码、决策量、编码冗余

2.2 统计编码

2.2.1 香农-范诺编码

2.2.2 霍夫曼编码

2.2.3 算术编码

2.3 游程编码

2.4 词典编码

2.4.1 词典编码的思想

2.4.2 LZ77算法

2.4.3 LZSS算法

2.4.4 LZ78算法

2.4.5 LZW算法



2.2 统计编码

■ 统计编码

- 给已知统计信息的符号分配代码的数据无损压缩方法

■ 编码方法

- 香农-范诺编码
- 霍夫曼编码
- 算术编码

■ 编码特性

- 香农-范诺编码和霍夫曼编码的原理相同，都是根据符号集中各个符号出现的频繁程度来编码，出现次数越多的符号，给它分配的代码位数越少
- 算术编码使用0和1之间的实数的间隔长度代表概率大小，概率越大间隔越长，编码效率可接近于熵



2.2.1 香农-范诺(Shannon-Fano)编码

- 最早阐述和实现这种编码的是Shannon(1948年)和Fano(1949年), 采用从上到下的方法进行编码。
- **Shannon-Fano编码步骤:**
 - (1) 按照符号出现的频度或概率排序
 - (2) 使用递归方法分成两个部分, 每一部分具有近似相同的次数

2.2.1 香农-范诺编码(例1)

■ 香农-范诺编码举例

- 有一幅40个像素组成的灰度图像，灰度共有5级，分别用符号A，B，C，D和E表示，如表格所示
 - (1) 计算该图像可能获得的压缩比的理论值
 - (2) 对5个符号进行编码
 - (3) 计算该图像可能获得的压缩比的实际值

符号在图像中出现的数目

符号	A	B	C	D	E
出现的次数	15	7	7	6	5
出现的概率	15/40	7/40	7/40	6/40	5/40

2.2.1 香农-范诺编码(例1)

(1) 压缩比的理论值

- 按照常规的编码方法，表示5个符号最少需要3位编码这幅图像总共需要120位。

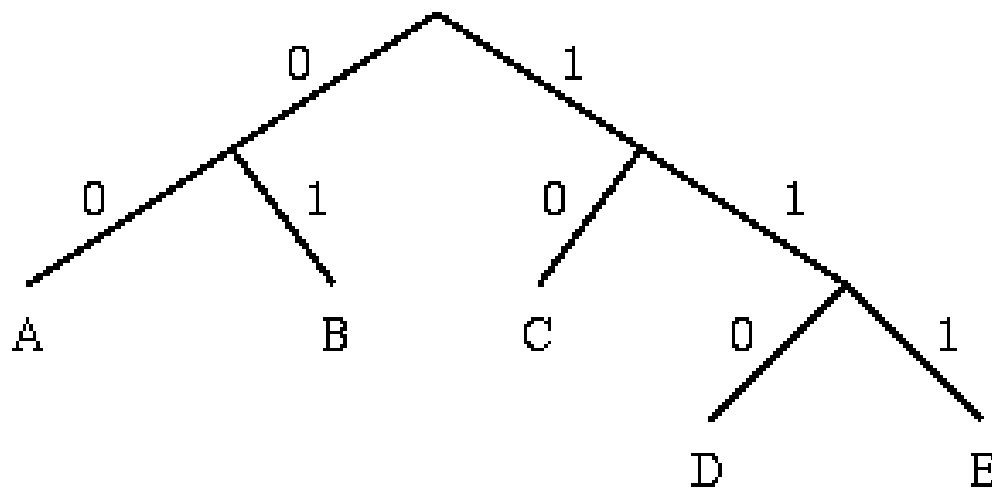
- 按照香农理论，这幅图像的熵为

$$\begin{aligned} H(X) &= -\sum_{i=1}^n p(x_i) \log_2 p(x_i) \\ &= -p(A) \log_2(p(A)) - p(B) \log_2(p(B)) - \dots - p(E) \log_2(p(E)) \\ &= (15/40) \log_2(40/15) + (7/40) \log_2(40/7) + \dots + (5/40) \log_2(40/5) \approx 2.196 \text{ bits} \end{aligned}$$

- 40个像素只需用87.84位
理论压缩比为120:87.84 \approx 1.37:1
即3:2.196 \approx 1.37

2.2.1 香农-范诺编码(例1)

(2) Shannon-Fano编码



香农-范诺算法编码举例

2.2.1 香农-范诺编码(例1)

符号	出现的次数(P_i)	$\log_2(1/P_i)$	分配的代码	需要的位数
A	15 (0.375)	1.4150	00	30
B	7 (0.175)	2.5145	01	14
C	7 (0.175)	2.5145	10	14
D	6 (0.150)	2.7369	110	18
E	5 (0.125)	3.0000	111	15

(3) 压缩比的实际值

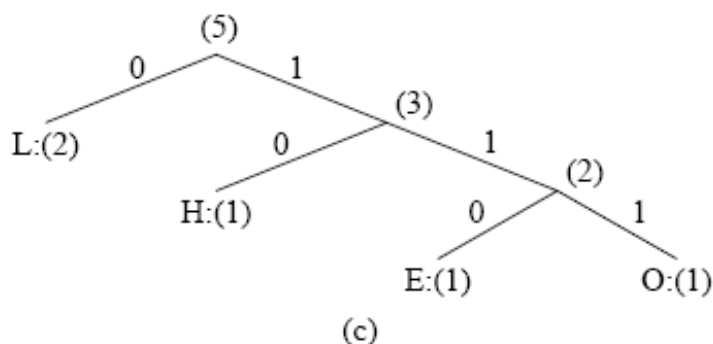
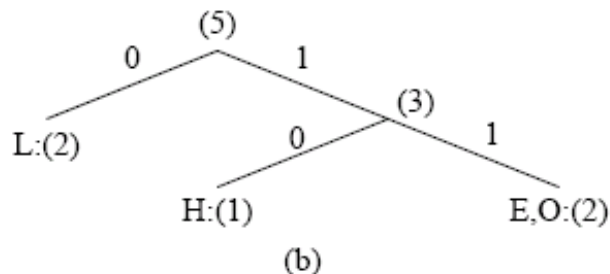
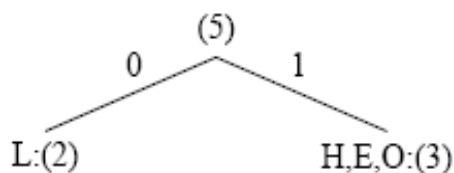
- 按照这种方法进行编码需要的总位数为
 $30+14+14+18+15=91$ ，实际的压缩比为 $120:91 \approx 1.32$
: 1

2.2.1 香农-范诺编码 (例2)

■ 香农范诺算法

假设对单词HELLO中字符进行编码，每个字符出现的频率如右：

符号	H	E	L	O
数目	1	1	2	1



Symbol	Count	$\log_2 \frac{1}{p_i}$	Code	# of bits used
L	2	1.32	0	2
H	1	2.32	10	2
E	1	2.32	110	3
O	1	2.32	111	3
TOTAL number of bits:				10



2.2.2 霍夫曼(Huffman) 编码

- 出现概率大的信源符号赋予短码字
- 两个最少出现的符号会有相同长度的码字

- **Huffman编码步骤:**
 - (1) 按照概率对符号进行排序;
 - (2) 概率最小的两个符号组成一个新符号, 概率等于对应符号的概率和 (进行这一步, 符号集大小减1);
 - (3) 对剩下的符号重复这个过程;
 - (4) 从根节点顺着树枝到叶子, 从编码过程对应的二叉树上得到码字。



2.2.2 霍夫曼编码 (例1)

■ 霍夫曼编码举例1

- 现有一个由5个不同符号组成的30个符号的字符串：

BABACACADADABBCBABEBEDDABEEEBB

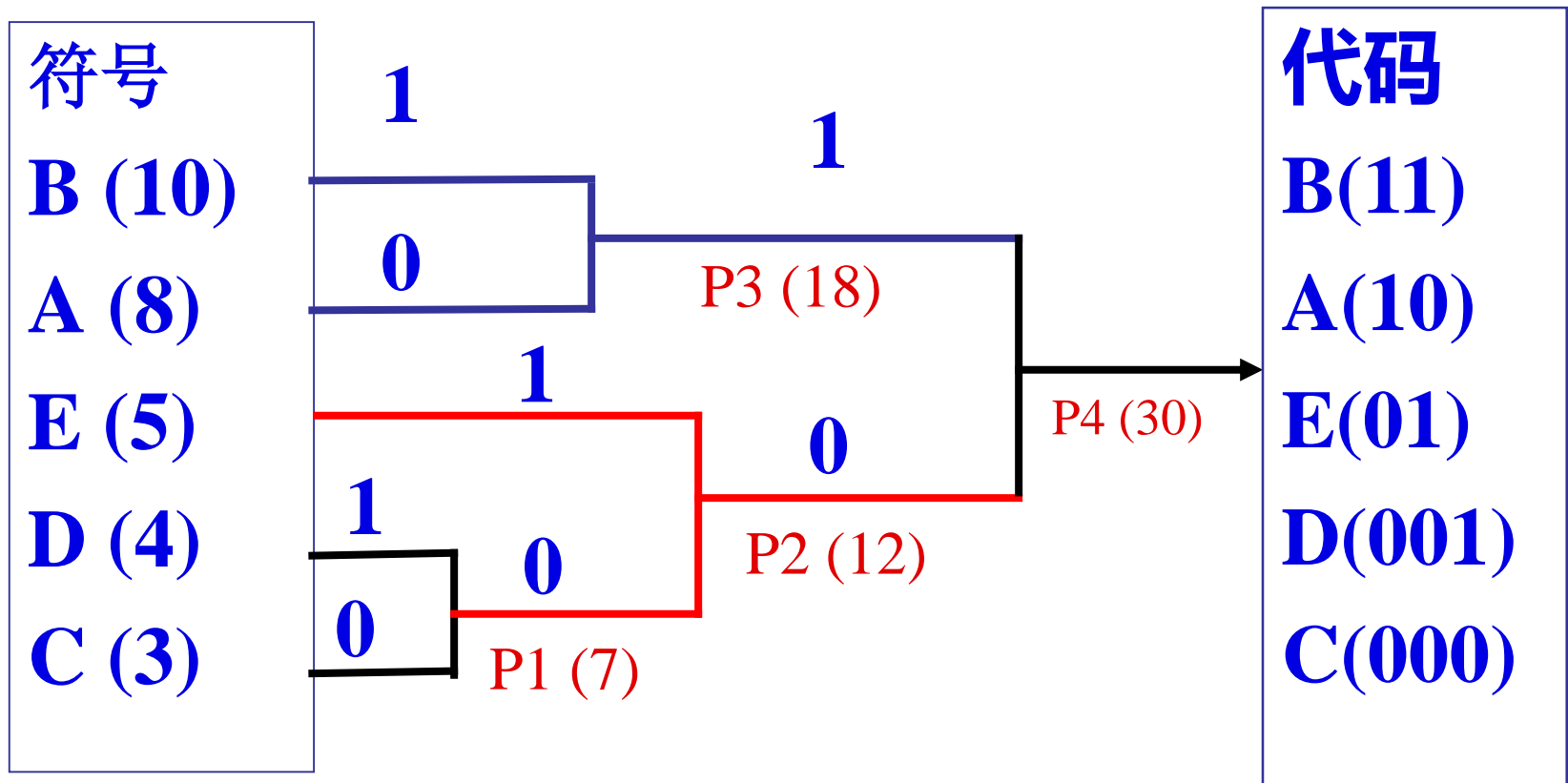
- 计算
 - (1) 该字符串的霍夫曼码
 - (2) 该字符串的熵
 - (3) 该字符串的平均码长
 - (4) 编码前后的压缩比

2.2.2 霍夫曼编码 (例1)

符号出现的概率

符号	出现的次数	$\log_2(1/p_i)$	分配的代码	需要的位数
B	10	1.585	?	
A	8	1.907	?	
C	3	3.322	?	
D	4	2.907	?	
E	5	2.585	?	
合计	30			

2.2.2 霍夫曼编码 (例1)





2.2.2 霍夫曼编码 (例1)

5个符号的代码

符号	出现的次数	$\log_2(1/p_i)$	分配的代码	需要的位数
B	10	1.585	11	20
A	8	1.907	10	16
C	3	3.322	010	9
D	4	2.907	011	12
E	5	2.585	00	10
合计	30	1.0		67

30个字符组成的字符串需要67位

2.2.2 霍夫曼编码 (例1)

(2) 计算该字符串的熵

$$H(X) = \sum_{i=1}^n p(x_i) I(x_i) = - \sum_{i=1}^n p(x_i) \log_2 p(x_i)$$

其中, $X = \{x_1, \dots, x_n\}$ 是事件 $x_i (i = 1, 2, \dots, n)$ 的集合,

并满足 $\sum_{i=1}^n p(x_i) = 1$

➤ $H(S)$

$$\begin{aligned} &= (8/30) \times \log_2(30/8) + (10/30) \times \log_2(30/10) + \\ &\quad (3/30) \times \log_2(30/3) + (4/30) \times \log_2(30/4) + \\ &\quad (5/30) \times \log_2(30/5) \\ &= [30 \lg 30 - (8 \times \lg 8 + 10 \times \lg 10 + 3 \times \lg 3 + 4 \\ &\quad \times \lg 4 + 5 \lg 5)] / (30 \times \log_2 2) \\ &= (44.3136 - 24.5592) / 9.0308 = 2.1874 \text{ (Sh)} \end{aligned}$$

2.2.2 霍夫曼编码 (例1)

(3) 计算该字符串的平均码长

➤ 平均码长: $l = \sum_{i=1}^N l_i p(l_i)$

$$\begin{aligned} &= (2 \times 8 + 2 \times 10 + 3 \times 3 + 3 \times 4 + 2 \times 5) / 30 \\ &= 2.233 \text{ 位/符号} \end{aligned}$$

平均码长: $67/30=2.233$ 位

(4) 计算编码前后的压缩比

- 编码前: 5个符号需3位, 30个字符, 需要90位
- 编码后: 共67位

压缩比: $90/67=1.34:1$



2.2.2 霍夫曼编码 (例2)

■ 霍夫曼编码举例

例：设有7个符号： $a_1, a_2, a_3, a_4, a_5, a_6, a_7$
它们出现概率是：0.2, 0.19, 0.18, 0.17, 0.15, 0.1, 0.01

2.2.2 Huffman编码 (分析)

每个编码均非其它码的前缀，因此唯一可译

($a_1=10, a_2=11, a_3=000, a_4=001, a_5=010, a_6=0110, a_7=0111$)

11 010 10 001 10 11 10 0111

a_2 a_5 a_1 a_4 a_1 a_2 a_1 a_7

1. 简单易实现
2. 编码效率较高
3. 必须预先知道信源的统计特性



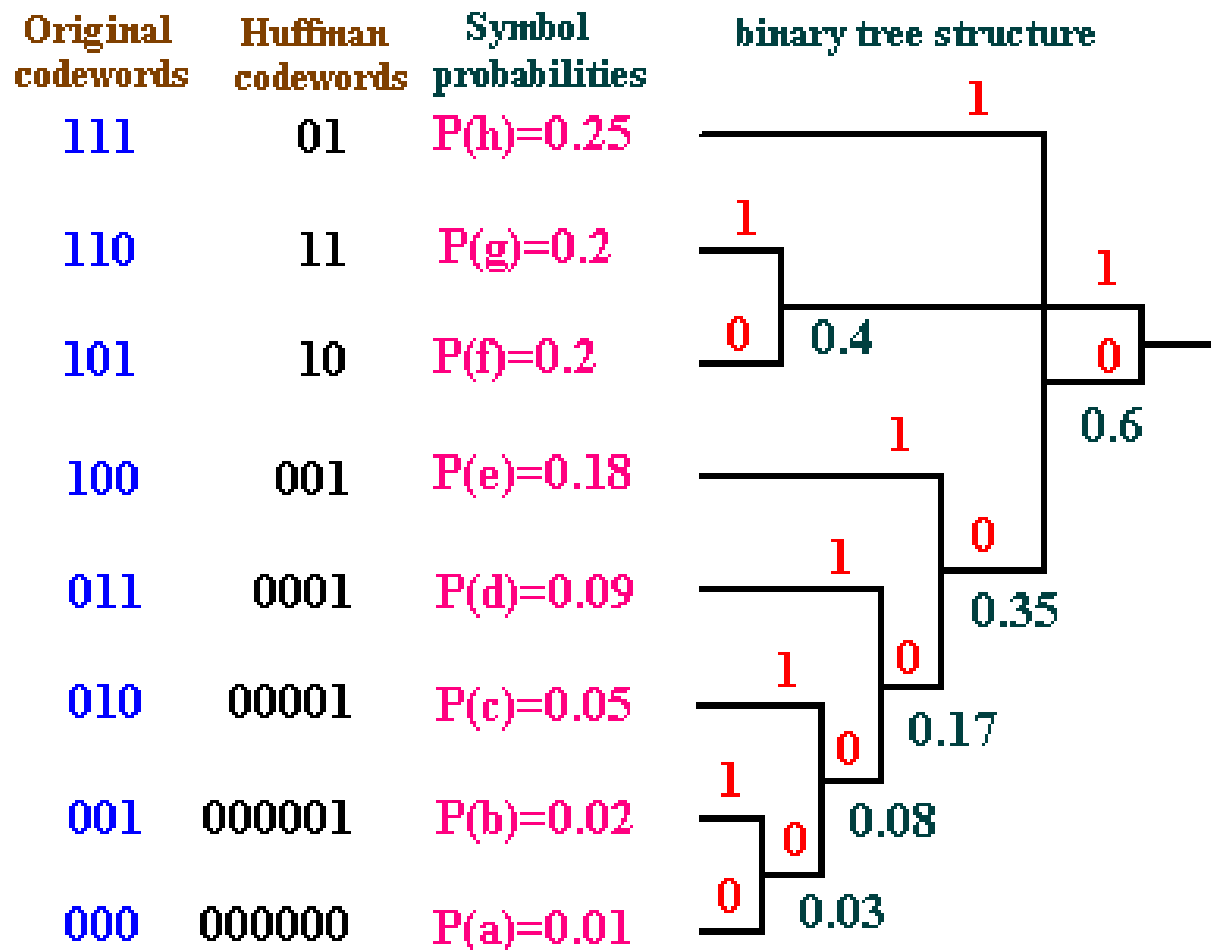
2.2.2 Huffman编码(例3)

■ 编码前

- $N = 8$ symbols: {a,b,c,d,e,f,g,h}
- 3 bits per symbol ($N = 2^3 = 8$)
- $p(a) = 0.01$, $p(b) = 0.02$, $p(c) = 0.05$, $p(d) = 0.09$,
 $p(e) = 0.18$, $p(f) = 0.2$, $p(g) = 0.2$, $p(h) = 0.25$

➤ 对该字符串进行Huffman编码

2.2.2 Huffman编码(例3)





2.2.2 Huffman编码 (分析)

霍夫曼码的码长虽然是可变的，但却不需要另外附加同步代码。

几个问题值得注意：

- 1.霍夫曼码没有错误保护功能；
- 2.霍夫曼码是可变长度码，因此很难随意查找或调用压缩文件中间的内容，然后再译码；
- 3.接收端需保存一个与发送端相同的霍夫曼码表。

Huffman编码-信源拓展

Symbols	Probability	Code
s_1	0.8	0
s_2	0.02	11
s_3	0.18	10

Huffman code: entropy = 0.816 bits/symbol;
average length = 1.2 bits/symbol;
redundancy = 0.384 bits/symbol, or 47% of entropy

Letter	Probability	Code
$s_1 s_1$	0.64	0
$s_1 s_2$	0.016	10101
$s_1 s_3$	0.144	11
$s_2 s_1$	0.016	101000
$s_2 s_2$	0.0004	10100101
$s_2 s_3$	0.0036	1010011
$s_3 s_1$	0.1440	100
$s_3 s_2$	0.0036	10100100
$s_3 s_3$	0.0324	1011

The extended alphabet and corresponding Huffman code; average length = 1.7516 bits/new symbol; or average length = 0.8758 bits/original symbol; redundancy = 0.06 bits/symbol, or 7% of entropy

2.2.3 算术编码

■ 算术编码思想

- Huffman编码中每个符号都用整数个bits来表示，影响编码效率。
- 若能把一串符号作为编码单位，则效率还可提高。
- 符号串的区间表示法
设符号串为： S_1, S_2, \dots, S_m
则它可以映射成为 $0 \cdots 1$ 中的唯一的一个子区间来表示

基本原理： 将编码的信息表示成实数0和1之间的一个间隔，信息越长，编码表示它的间隔就越小，表示这一间隔所需的二进制位就越多。



2.2.3 算术编码

■ 算术编码(Arithmetic Coding)

- 给已知统计信息的符号分配代码的数据无损压缩技术
- 基本思想是用0和1之间的一个数值范围表示输入流中的一串字符，而不是给输入流中的每个字符分别指定一个码字
- 实质上是为整个输入字符流分配一个“码字”，因此它的编码效率可接近于熵

2.2.3 算术编码-标签的生成

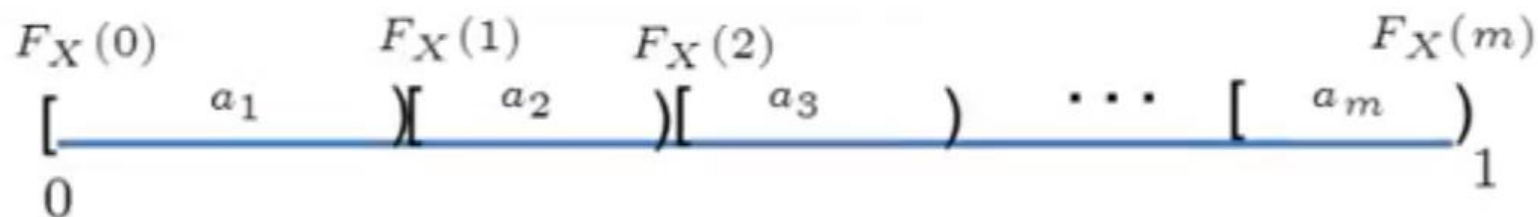
$\mathcal{A} = \{a_1, a_2, \dots, a_m\}; \quad P(a_i)$ Random variable: $X(a_i) = i$

Probability density function (pdf): $P(X = i) = P(a_i)$

Cumulative distribution function (cdf): $F_X(i) = \sum_{k=1}^i P(X = k)$

- 核心思想：随着序列中收到的符号越来越多，标签所在的区间逐渐缩小。

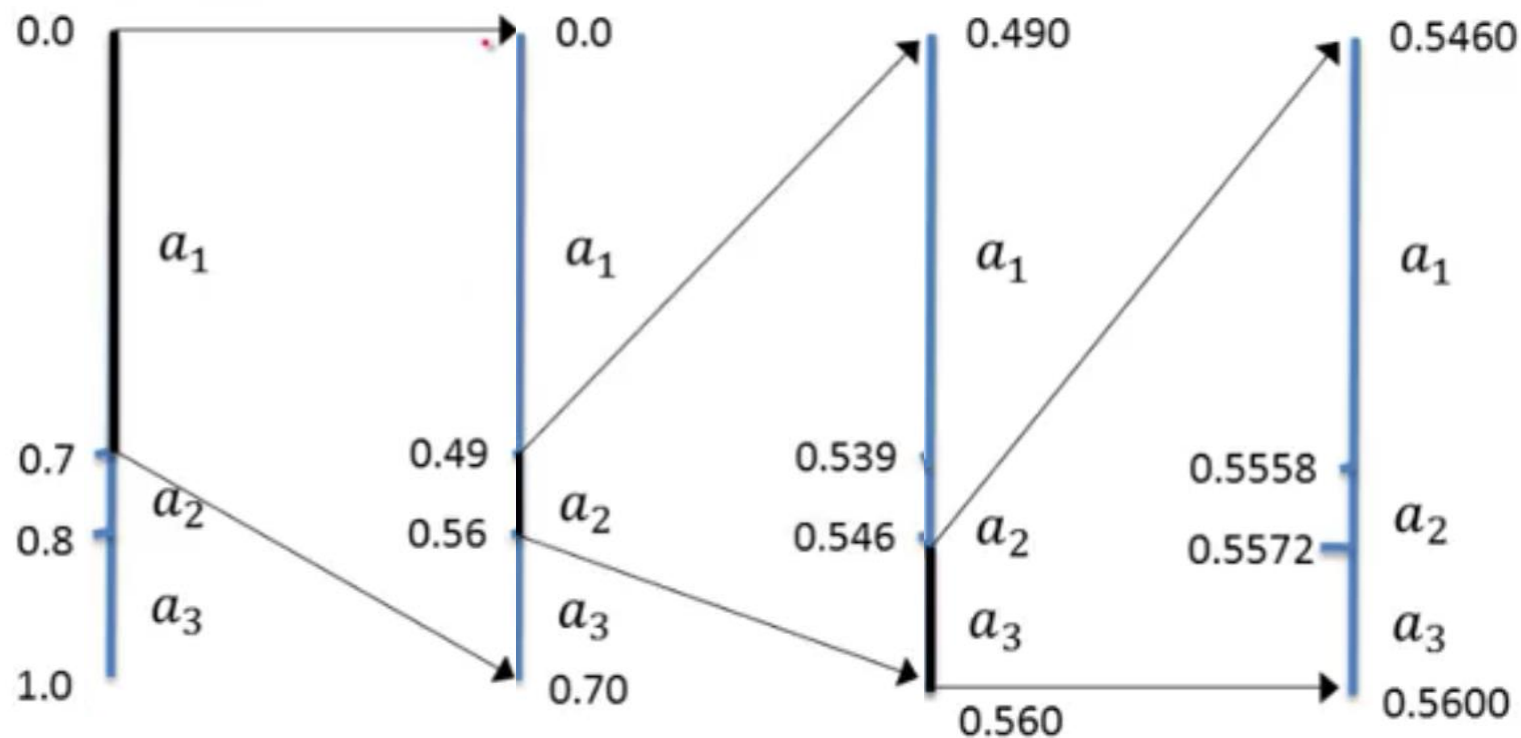
把 $[0,1)$ 单元区间划分成下面的子区间



$$a_i \in [F_X(i-1), F_X(i))$$

2.2.3 算术编码-(例1)

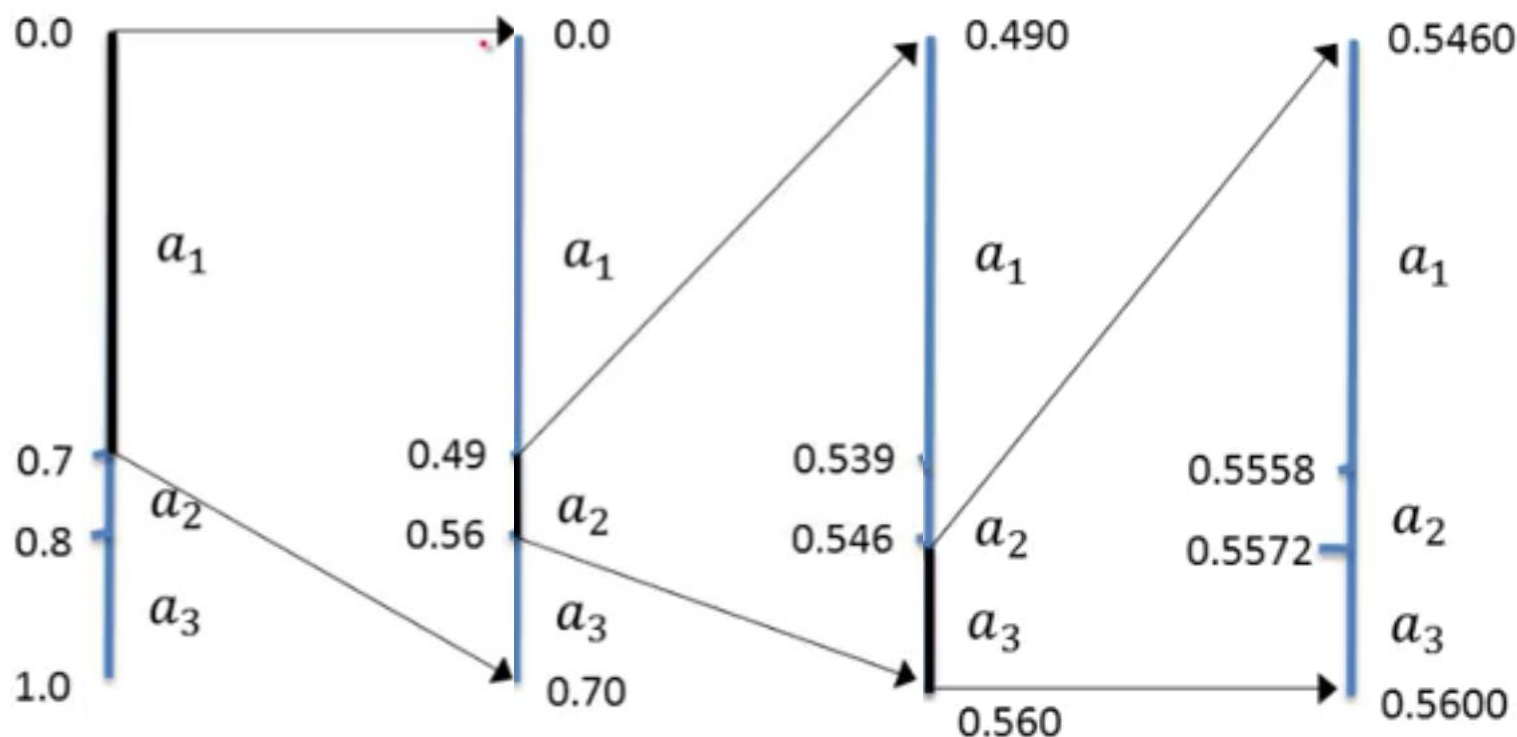
$\mathcal{A} = \{a_1, a_2, a_3\}$; $P(a_1) = 0.7, P(a_2) = 0.1, P(a_3) = 0.2$; $F_X(1) = 0.7, F_X(2) = 0.8, F_X(3) = 1$



2.2.3 算术编码-标签解码

$\mathcal{A} = \{a_1, a_2, a_3\}$; $P(a_1) = 0.7, P(a_2) = 0.1, P(a_3) = 0.2$; $F_X(1) = 0.7, F_X(2) = 0.8, F_X(3) = 1$

Tag=0.55 , 3 symbols



模仿编码过程



2.2.3 算术编码

■ 算术编码算法描述

可用伪代码描述如下：

set Low to 0

set High to 1

while there are input symbols do

 take a symbol

$\text{CodeRange} = \text{High} - \text{Low}$

$\text{High} = \text{Low} + \text{CodeRange} * \text{HighRange}(\text{symbol})$

$\text{Low} = \text{Low} + \text{CodeRange} * \text{LowRange}(\text{symbol})$

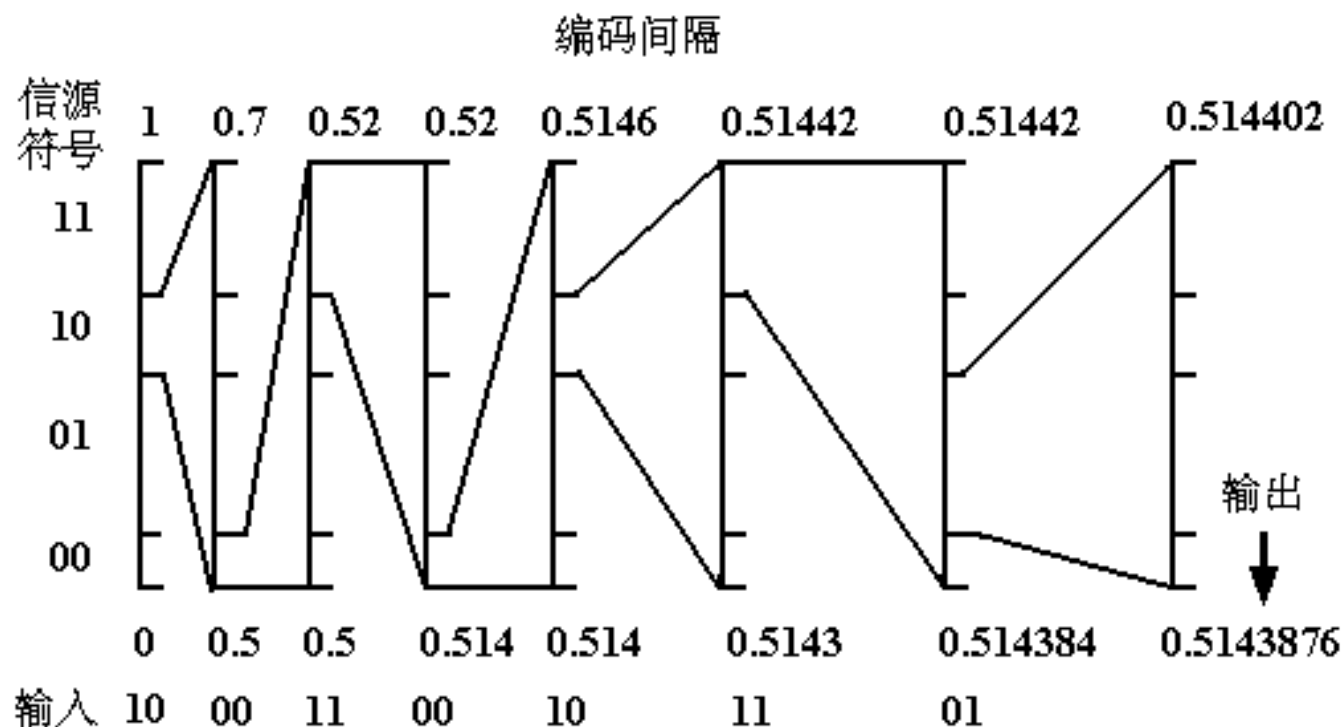
end of while

output Low

2.2.3 算术编码 (例2)

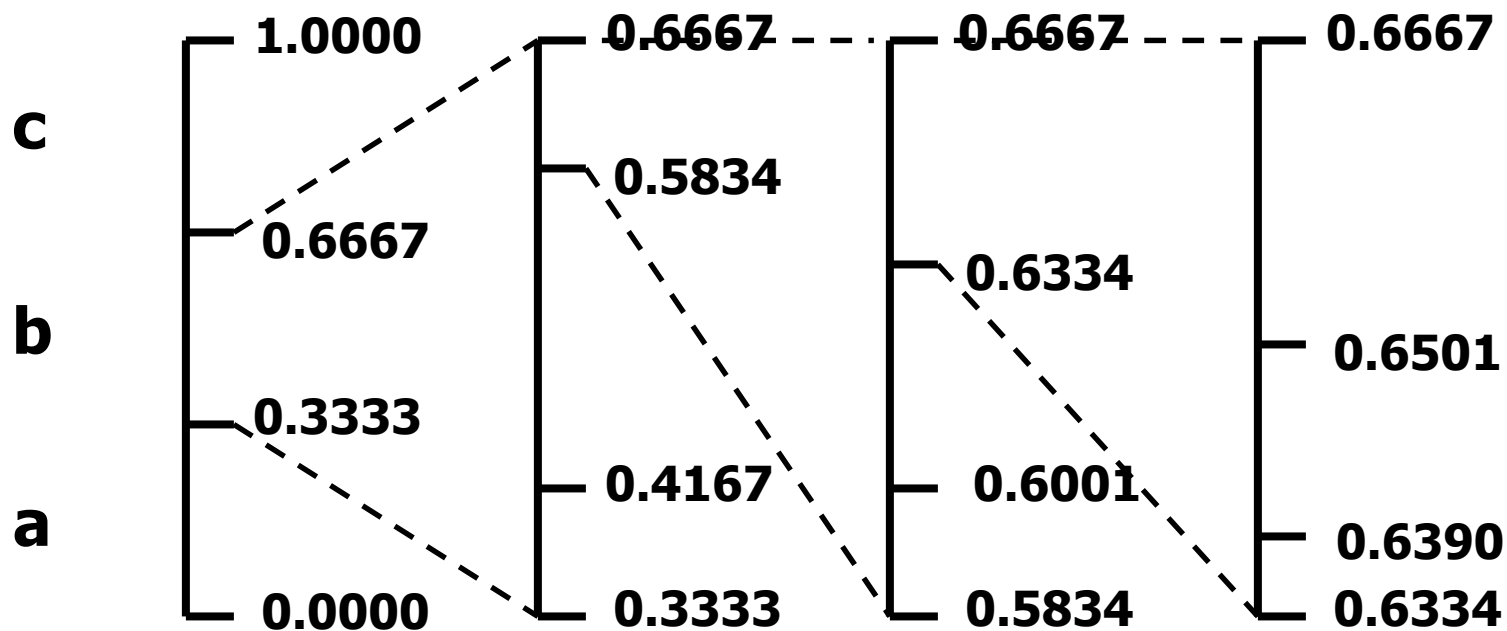
假设信源符号为{00, 01, 10, 11}, 这些符号的概率分别为{ 0.1, 0.4, 0.2, 0.3 }, 根据这些概率可把间隔[0, 1)分成4个子间隔: [0, 0.1), [0.1, 0.5), [0.5, 0.7), [0.7, 1), 二进制消息序列的输入为: 10 00 11 00 10 11 01

符号	00	01	10	11
概率	0.1	0.4	0.2	0.3
初始编码间隔	[0, 0.1)	[0.1, 0.5)	[0.5, 0.7)	[0.7, 1)



2.2.3 算术编码 (补充: 自适应算术编码)

输入序列为: bcc



c	1/3	1/4	2/5	3/6
b	1/3	2/4	2/5	2/6
a	1/3	1/4	1/5	1/6

2.2.3 算术编码

■ 算术编码分析



子区间的大小对应于符号串中各符号概率的乘积大小

子区间的位置对应着 符号串中符号的排列顺序

子区间的头尾均用二进制数表示

该符号串的编码是满足下列条件的、最短的数F:

$$\text{头} < = F < \text{尾}$$

可见:

不同的串, 对应着不同的子区间, 选择子区间中任意的一个数作为该符号串的代码, 则必有不同的编码

选择 子区间中**最短的数**作为 该符号串的代码

一般而言, 子区间越宽, 码长越短; 子区间越窄, 码长越长



2.2.3 算术编码

■ 算术编码注意点

1. 由于实际计算机精度不可能无限长，运算中溢出是明显的问题，但多数机器都有16位、32位或者64位的精度，因此可使用比例缩放法解决。
2. 算术编码器对消息只产生一个码字，这个码字是在 $[0, 1)$ 中的一个实数，因此译码器在接受到表示这个实数的所有位之前不能进行译码。
3. 算术编码也是一种对错误很敏感的编码方法，如果有一位发生错误就会导致整个消息译错。

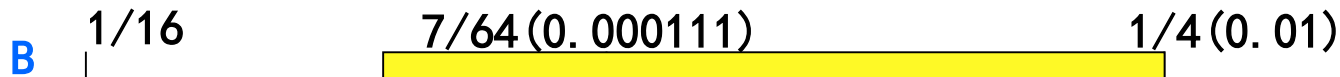
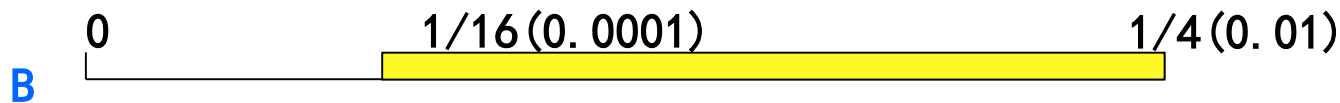
2.2.3 算术编码 (思考)

信源符号及其概率如下:

符号	A	B
概率	0.25	0.75

求输入串ABBA的算术二进制编码(赋予大概率为右区间), 并求最短编码。

2.2.3 算术编码 (答案)



ABBA符号串对应的区间为 $[7/64, 37/256)$ ，即对应的二进制区间为 $[0.000111, 0.00100101)$ ，取最短二进制编码为：0.001



目录

2.0 简介

2.1 信息论基础知识

信息量、熵、信源编码、决策量、编码冗余

2.2 统计编码

2.2.1 香农-范诺编码

2.2.2 霍夫曼编码

2.2.3 算术编码

2.3 游程编码

2.4 词典编码

2.4.1 词典编码的思想

2.4.2 LZ77算法

2.4.3 LZSS算法

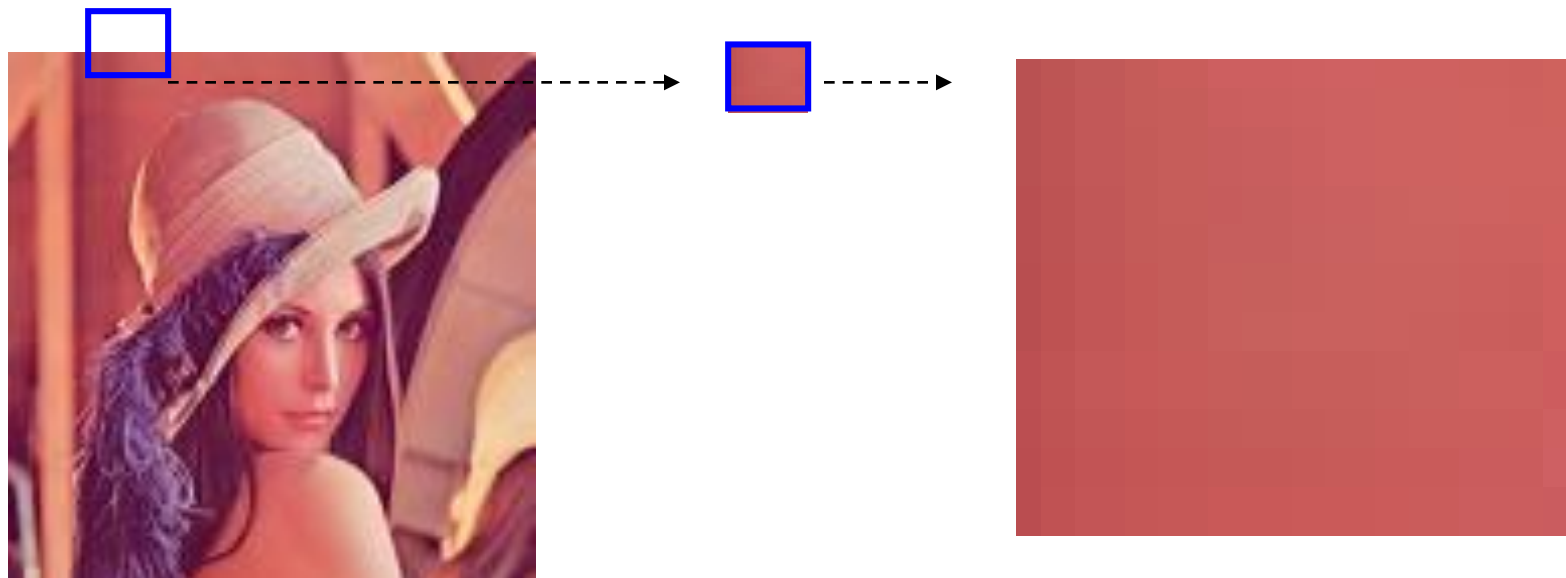
2.4.4 LZ78算法

2.4.5 LZW算法

2.3 RLE(Run Length Encoding)编码

RLE编码

- 利用连续数据单元有相同数值的特点对数据进行压缩





2.3 RLE(Run Length Encoding)编码

■ 行程编码:

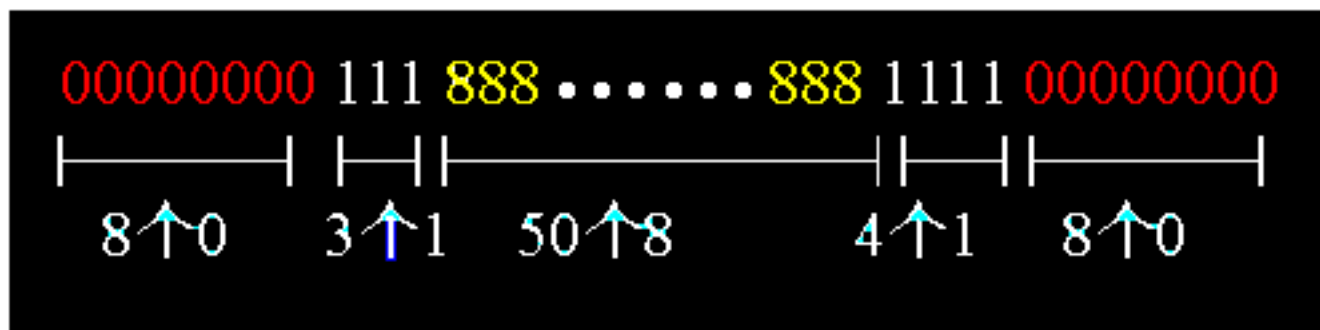
- 不需存储每个像素的颜色值，而仅存储一个像素的颜色值以及具有相同颜色的像素数目。

■ 行程长度:

- 具有相同颜色并且是连续的像素数目。

2.3 RLE编码(例)

- 假定一幅灰度图像，第n行的像素值为：



用RLE编码方法得到的代码为：

80315084180。

压缩比 $\approx 7:1$



2.3 RLE编码 (特点)

- **RLE的压缩比取决于图像本身特点。**
 - 如果图像中具有相同颜色的图像块越大，压缩比就越高。反之，压缩比越小。
- **译码时按照与编码时采用的相同规则进行，还原后得到的数据与压缩前的数据完全相同。**



2.3 RLE编码适用范围

- **适用于:**

- 计算机生成的图像、黑白二值图像

- **不适用于:**

- 颜色丰富的自然图像（除非需和其他的压缩编码技术联合应用）、文字

2.3 RLE编码 (适用场景)

- 行程编码多用于黑白二值图像的压缩中。

例如

0000000011111111111100000111111

8

12

5

6

被转化为一系列黑串和白串长度的编码:

8C56 (0x)



2.3 RLE编码

■ 对于文本压缩:

- RLE的效果不好
- 例如: HoHai->H1o1h1a1i1
- 除非经常遇到
 - “Ahhhhhhhhhhhhhhhhhhhh!”,
 - “Noooooooooooo!”.
 - “Goooooooooooooooooooooooooooooogle”



目录

2.0 简介

2.1 信息论基础知识

信息量、熵、信源编码、决策量、编码冗余

2.2 统计编码

2.2.1 香农-范诺编码

2.2.2 霍夫曼编码

2.2.3 算术编码

2.3 游程编码

2.4 词典编码

2.4.1 词典编码的思想

2.4.2 LZ77算法

2.4.3 LZSS算法

2.4.4 LZ78算法

2.4.5 LZW算法

板凳宽，扁担长。
扁担没有板凳宽，
板凳没有扁担长。





2.4.1 词典编码

- 在很多应用中，信源的输出包含了很多重复的样式（**recurring patterns**）。
- 对于这种信源，可以用一个列表、或者词典来记录重复出现的样式。
- 词典编码有静态和自适应两种。
 - 大多自适应的方法来根源于Ziv和Lempel在1977年(LZ77)和1978年(LZ78)的两篇文章。

静态词典编码

$$\mathcal{A} = \{a, b, c, d, r\}$$

Code	Entry
000	a
001	b
010	c
011	d
100	r
101	ab
110	ac
111	ad

a b r a c a d a b r a



2.4.1 词典编码

■ 通用编码技术

- 不需要知道统计特性

■ 利用数据本身包含许多重复字符串

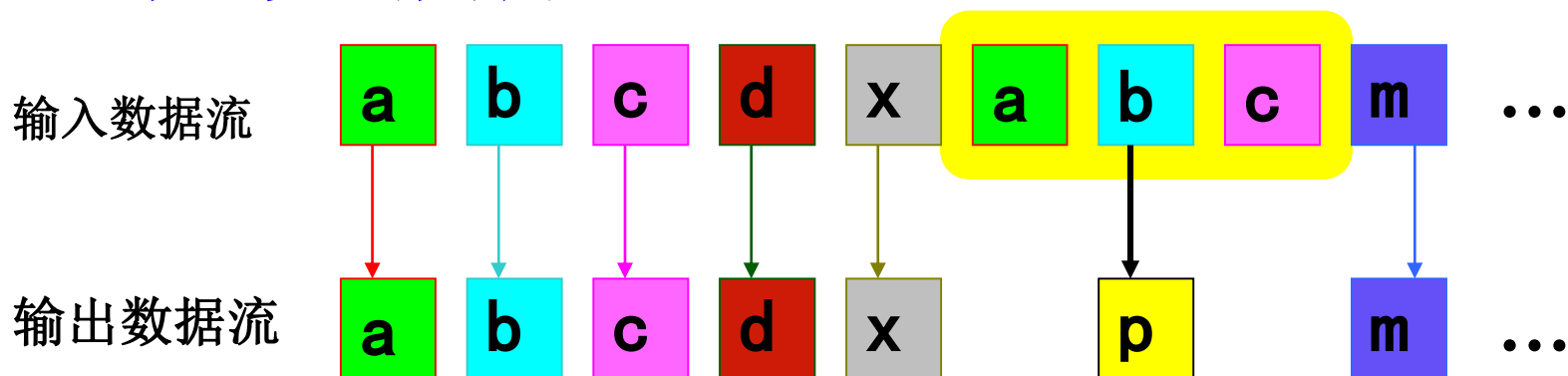
例：“吃葡萄不吐葡萄皮,不吃葡萄倒吐葡萄皮”，
用简单代号代替这些字符串,实现压缩。

■ 实用的词典编码算法的核心

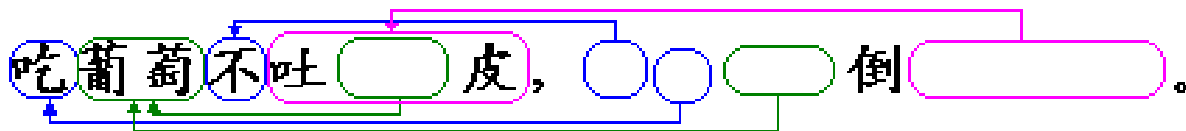
- 如何动态地形成词典，以及如何选择输出格式以减小冗余

2.4.1 词典编码(分类)

第一类词典编码思想



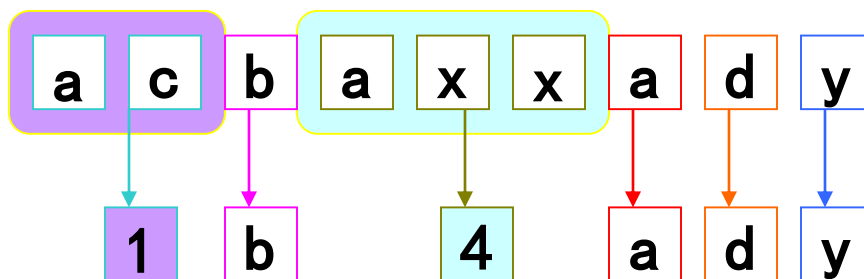
第一类词典法的想法是企图**查找**正在压缩的字符序列是否在**以前输入**的数据中出现过，然后用已经出现过的字符串替代重复的部分，它的输出仅仅是指向早期出现过的字符串的“指针”。



2.4.1 词典编码(分类)

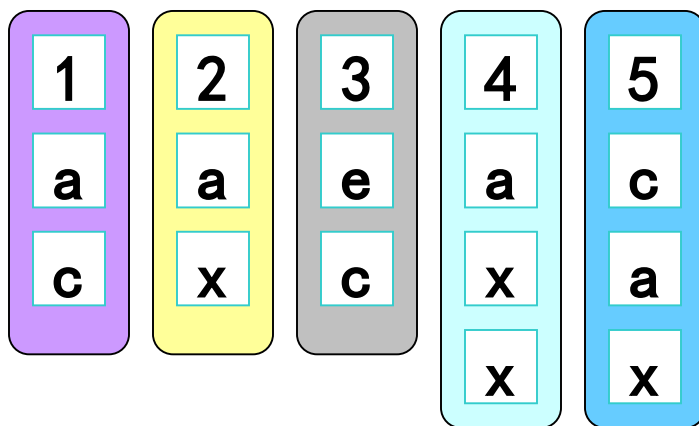
第二类词典编码思想

输入
数据流



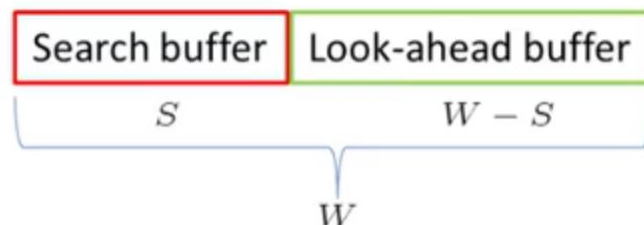
输出
数据流

编码
词典



... 第二类算法的想法是企图从输入的数据中**创建一个“短语词典”** (dictionary of the phrases), 这种短语可以是任意字符的组合。编码数据过程中当遇到已经在词典中出现的“短语”时, 编码器就输出这个词典中的短语的“索引号”, 而不是短语本身。

2.4.2 第一类词典编码-LZ77



- 输出: Encode: $\langle o, l, c \rangle$
<位移, 匹配长度, 下个字符的码字>

- 窗口前移 $l+1$ 个字符

Bits needed (using fixed length codes)

$$\lceil \log_2 S \rceil + \lceil \log_2 W \rceil + \lceil \log_2 |\mathcal{A}| \rceil$$

2.4.2 LZ77 (例1)



Example: Encoding

c a b r a c a d a b r a r r a d

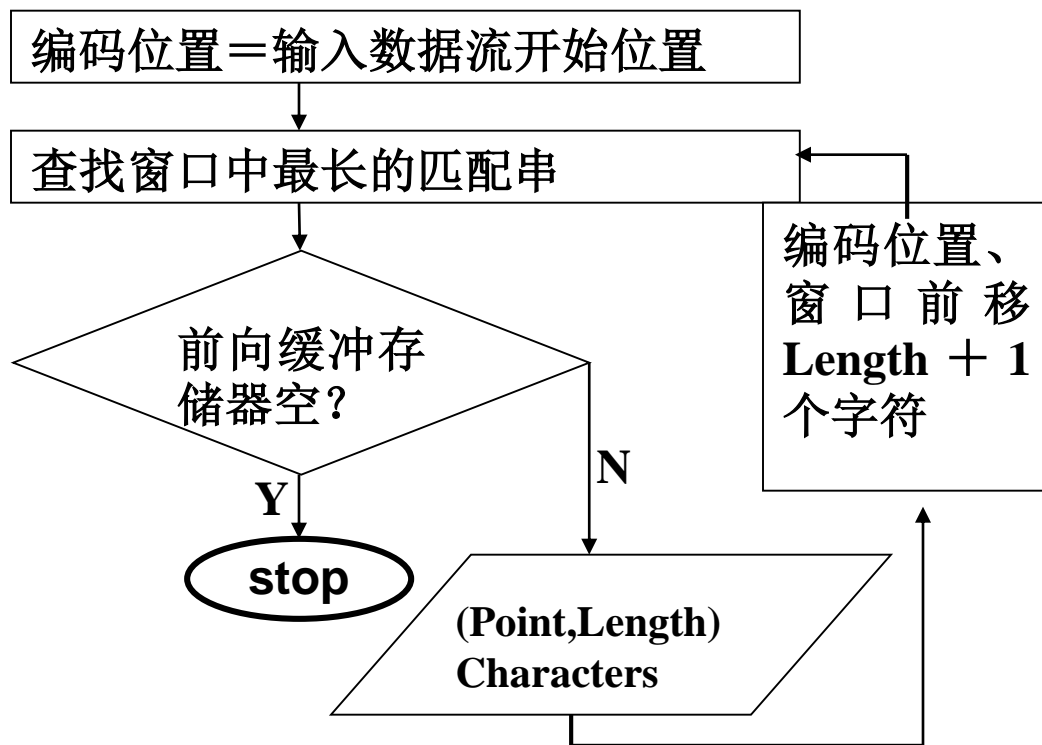
c a b r a c a d a b r a r r a d
 $\langle o, l, c \rangle$

Example: Decoding

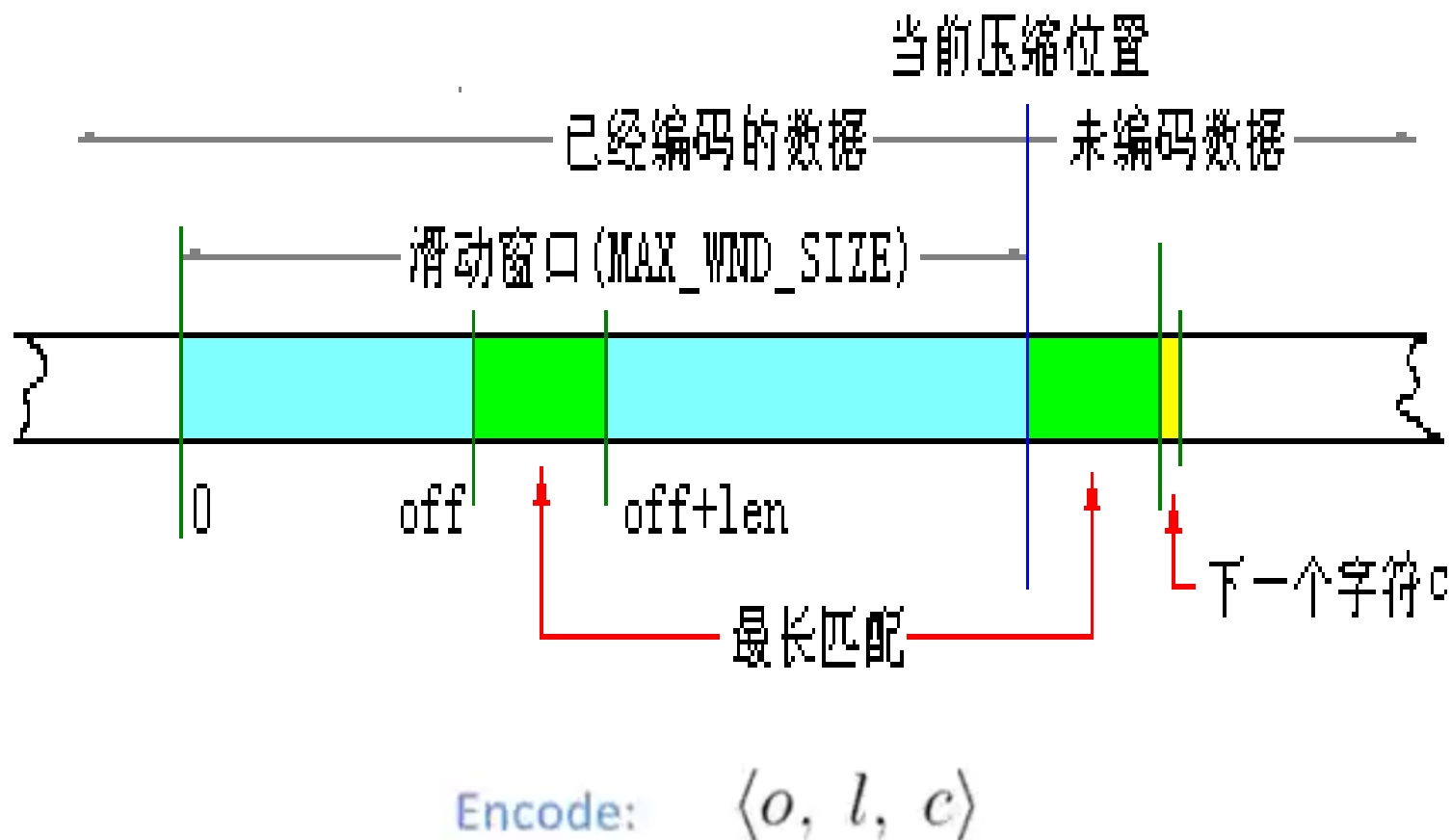
c a b r a c a

2.4.2 第一类词典编码-LZ77

核心：查找从前向缓冲存储器开始的最长的匹配串



2.4.2 第一类词典编码-LZ77



2.4.2 第一类词典编码-LZ77(例2)

位置	1	2	3	4	5	6	7	8	9	10
字符	A	A	B	C	B	B	A	B	C	A

步骤	位置	匹配串	字符	输出
1	1	--	A	(0,0)A
2	2	A	B	(1,1)B
3	4	--	C	(0,0)C
4	5	B	B	(2,1)B
5	7	ABC	A	(5,3)A

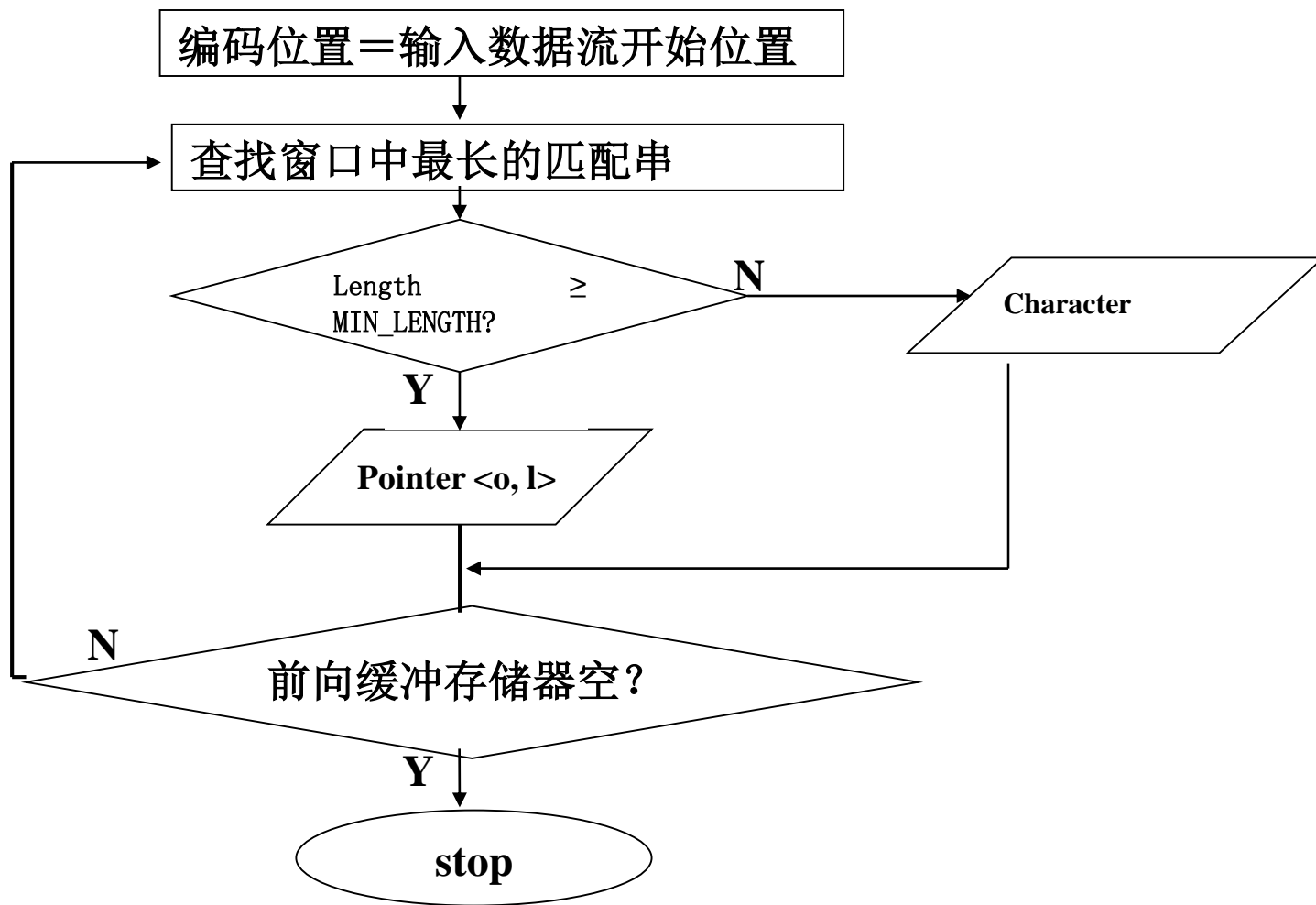


2.4.3 第一类词典编码-LZSS

LZ77通过输出真实字符解决了在窗口中出现没有匹配串的问题，但这个解决方案包含有冗余信息。冗余信息表现在两个方面，一是空指针，二是编码器可能输出额外的字符，这种字符是指可能包含在下一个匹配串中的字符。

LZSS算法以比较有效的方法解决这个问题，它的思想是如果匹配串的长度比指针本身的长度长就输出指针，否则就输出真实字符。另外要输出额外的标志位区分是指针还是字符。

2.4.3 第一类词典编码-LZSS



2.4.3 第一类词典编码-LZSS (例)

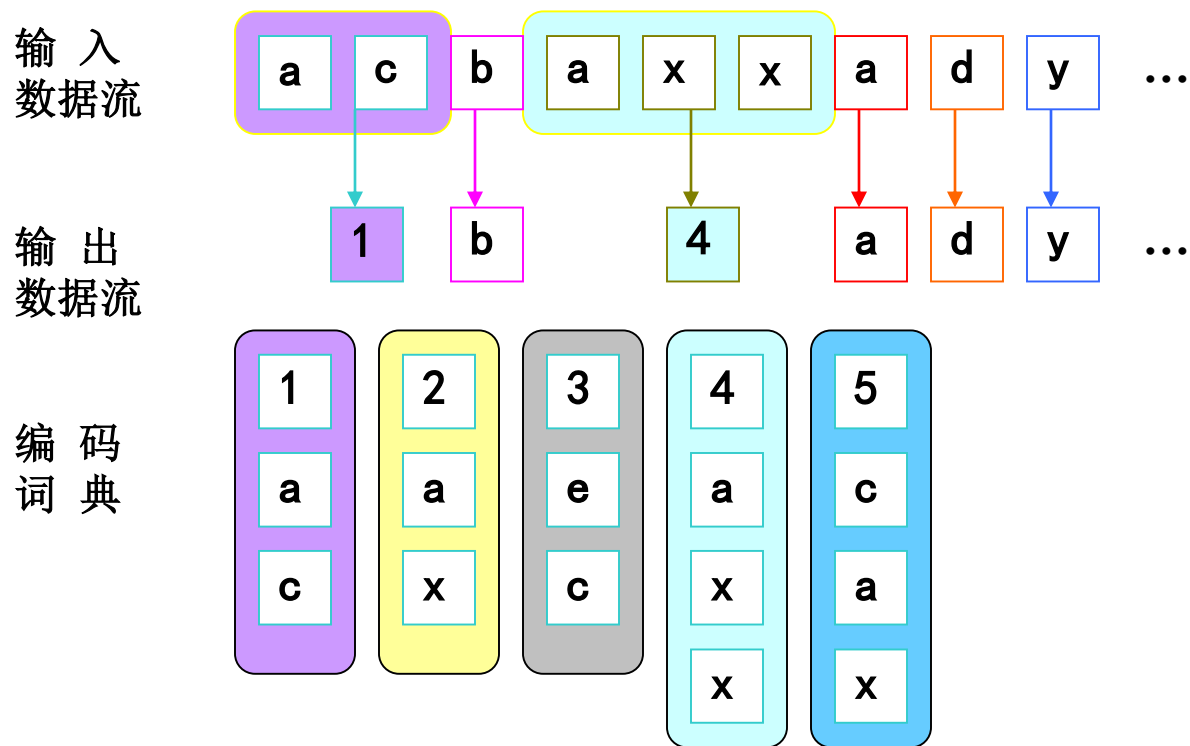
位置	1	2	3	4	5	6	7	8	9	10	11
字符	A	A	B	B	C	B	B	A	A	B	C

Length \geq MIN_LENGTH?

MIN_LENGTH=2

步骤	位置	匹配串	输出
1	1	— —	A
2	2	A	A
3	3	— —	B
4	4	B	B
5	5	— —	C
6	6	BB	(3,2)
7	8	AAB	(7,3)
8	11	C	C

词典编码 (第二类词典编码思想)



- 从输入的数据中**创建一个“短语词典”** (dictionary of the phrases)。
- 编码过程中遇到已在词典中的“短语”时，输出词典中的短语的“索引号”，而不是短语本身。



词典编码(第二类词典编码历史)

■ 第二类词典编码

- J.Ziv和A.Lempel在1978年首次发表了介绍这种编码方法的文章。(LZ78)
- 在此基础上, Terry A.Welch在1984年发表了改进这种编码算法的文章, 因此称为LZW(Lempel-Ziv Welch)压缩编码, 首先在高速硬盘控制器上应用了这种算法。



2.4.4 第二类词典编码-LZ78

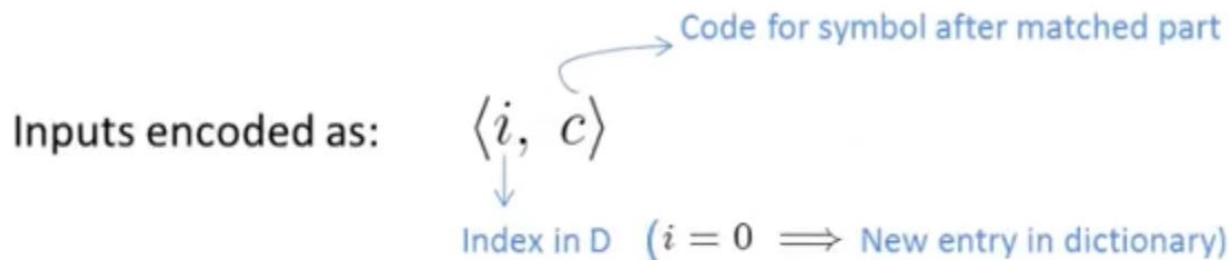
■ LZ78的编码思想:

- 不断地从字符流中提取新的字符串(String), 通俗地理解为新“词条”, 然后用“代号”也就是码字(Code word)表示这个“词条”。
- 对字符流的编码就变成了用码字(Code word)去替换字符流(Char stream), 生成码字流(Code stream), 从而达到压缩数据的目的。

■ LZ77利用滑动窗口里的内容作为字典, 找最长子串; LZ78动态构建字典

2.4.4 LZ78

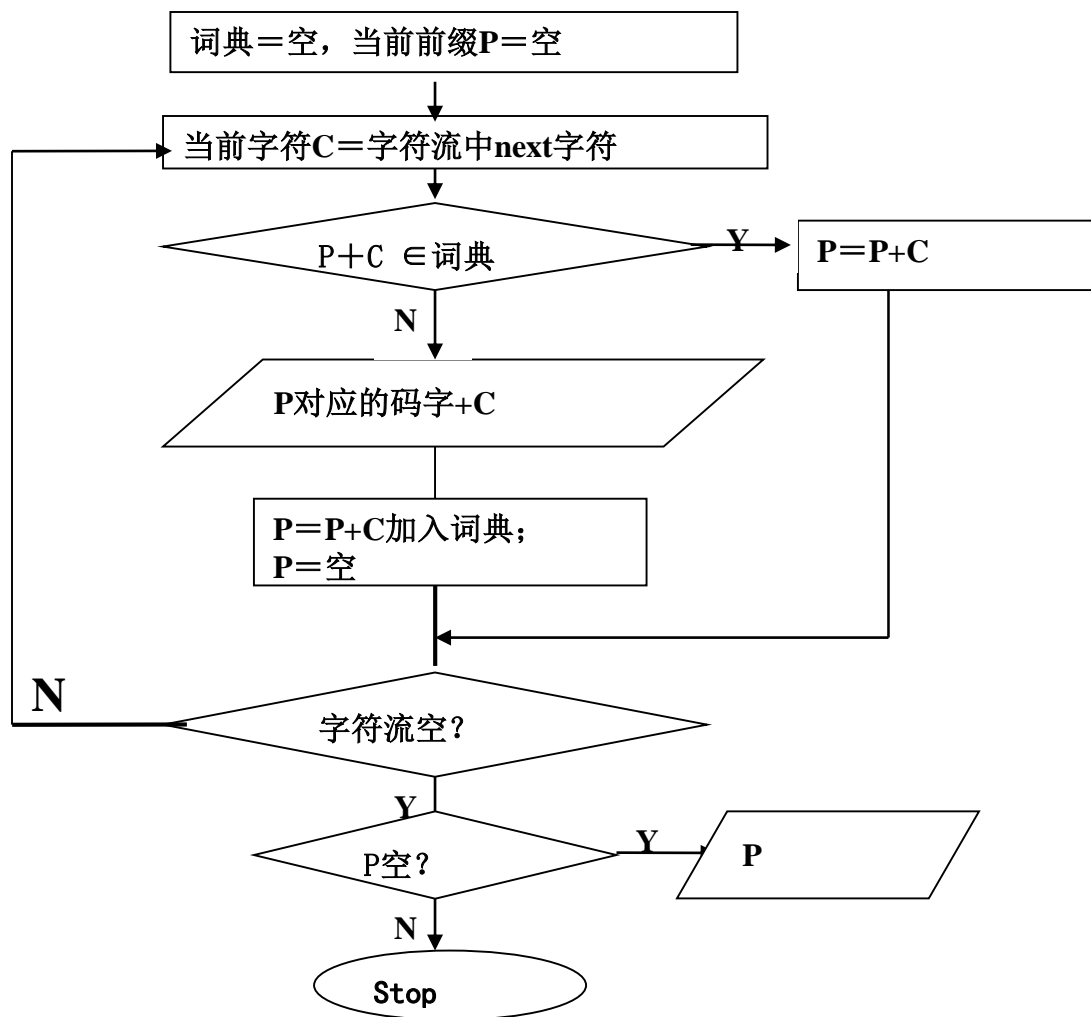
- 不依赖于搜索窗口、建立了一个专门的词典
- 编码端和解码端都需要建立这个词典，并且两个词典保持一致。



码字-字符(W,C)对，每次输出一对到码字流中，与码字W相对应的缀-字符串(String)用字符C进行扩展生成新的缀-字符串(String)，然后添加到词典中。

2.4.4 LZ78示例

位置	1	2	3	4	5	6	7	8	9
字符	A	B	B	C	B	C	A	B	A
步骤		位置		词典			输出		
1		1							
2		2							
3		3							
4		5							
5		8							



2. 4. 4 LZ78示例

b a w w a / b a w w a /

Encoder output	Index	Entry
<0,c(b)>	1	b
<0,c(a)>	2	a
<0,c(w)>	3	w
<3,c(a)>	4	wa
<0,c(/)>	5	/
<1,c(a)>	6	ba
<3,c(w)>	7	ww
<2,c(/)>	8	a/

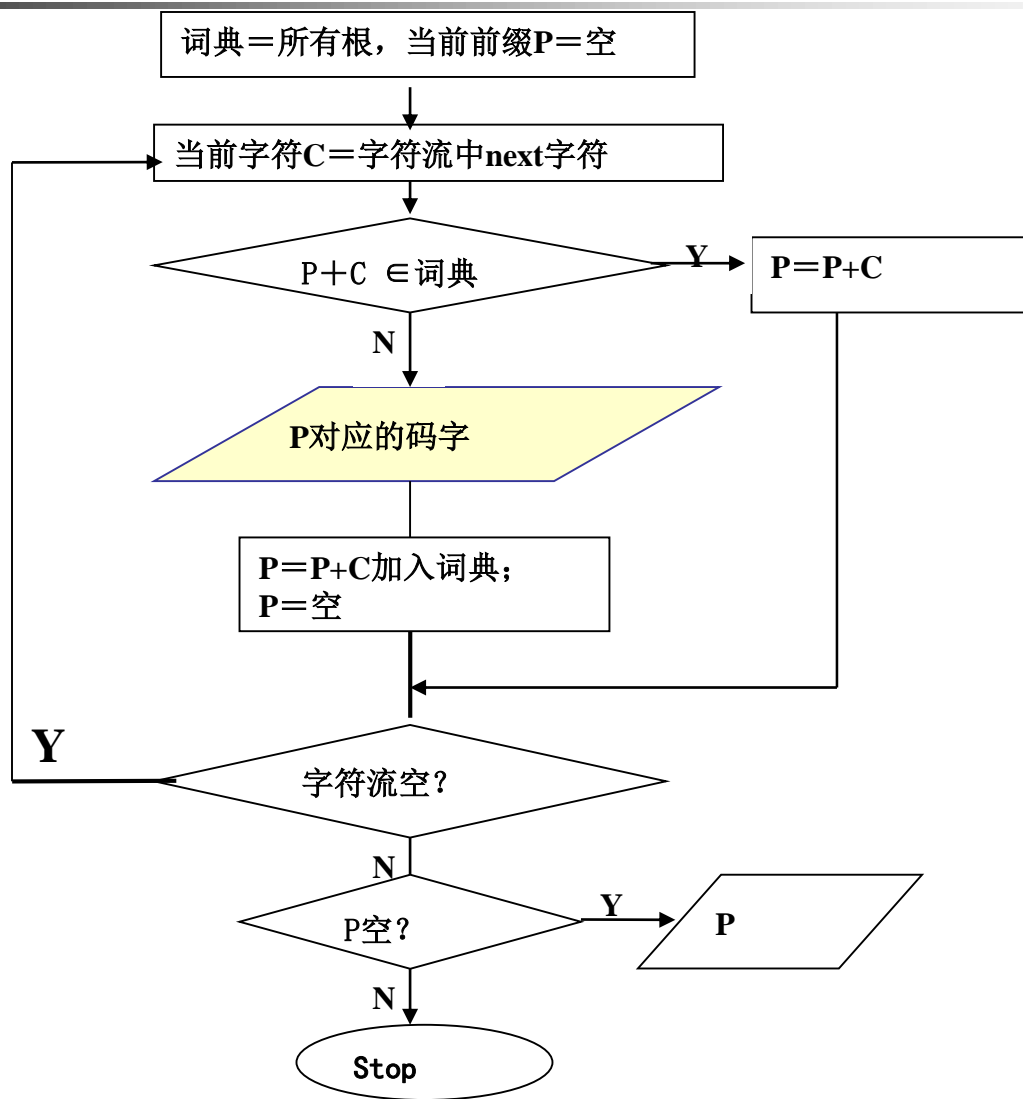


2.4.5 第二类词典编码-LZW

■ LZW算法

- 只输出代表词典中的字符串(String)的码字(code word)。开始时词典包含可能在字符流出现中的所有单个字符。
- 每个编码步骤开始时都使用一字符前缀(one-character prefix)，至少可在词典中找到长度为1的匹配串。

2.4.5 第二类词典编码-LZW



2.4.5 第二类词典编码-LZW

位置	1	2	3	4	5	6	7	8	9
字符	A	B	B	A	B	A	B	A	C

步骤	位置	词典	输出
		(1)	A
		(2)	B
		(3)	C
1	1	(4)	AB (1)
2	2	(5)	BB (2)
3	3	(6)	BA (2)
4	4	(7)	ABA (4)
5	6	(8)	ABAC (7)
6	---	---	--- (3)

LZW解码

步骤	位置	词典		输出
		(1)	A	
		(2)	B	
		(3)	C	
1	1	(4)	AB	(1)
2	2	(5)	BB	(2)
3	3	(6)	BA	(2)
4	4	(7)	ABA	(4)
5	6	(8)	ABA C	(7)
6	--	--	--	(3)



总结

- 无损压缩
- 熵编码
- 游程编码
- 词典编码



作业题：香农-范诺(Shannon- Fano)编码

现有7个待编码的符号，它们的概率如下。计算该符号集的：（1）Shannon-Fano编码；（2）平均码长；（3）“x1 x3 x6 x4” 编码结果。

Character	Probability
X1	0.1
X2	0.05
X3	0.2
X4	0.15
X5	0.15
X6	0.25
X7	0.1



作业

- 教科书第43页
 - Shannon-Fano编码，课后题 2.4, 2.5, 2.6
- 提交时间：11月16号（周五）

END

第2章 数据无损压缩

