Коды Шеннона-Фано • длина кода символа обратно пропорциональна частоте встречаемости символа •7/40  $\cdot \cdot 5/40$ ce··5/40 • • 4/40

. . 2/40 . . . . 1111

Одно из применений сбалансированных бинарных деревьев (как строили на рк ака деревья сжатия) - поиск

• Более часто искомые данные - выше. По очевидным причинам: чаще ищут - быстрее выдавать

Асимптотическое достиженеие оптимальности при увеличении числа

Схема

var-var,

var-var,

block-var

block-var

Самые частые данные - выше. Данные встречаются чаще -> их чаще будут искать

сообщений (больше исходник - лучше сжатие)

2) Сжатый текст = конкатенация кодовых слов

1) Каждое сообщение - уникальное кодовое слово

Условия оптимальности

..3/40

*а,* - сообщение

Как строить деревья поиска? 2 подхода:

Статистические методы сжатия данных

Код

1975

о слова

Шеннона-Фано, 1948

Универсальные коды,

Арифметические коды,

Хаффмана, 1952

 $p(a_i)$  - вероятность сообщения  $x_i$  - кодовое слово  $-\lg p(x_i)$  - длина кодового слова  $-\lg p(x_i) + 1$  - длина кодового слова, если  $H \le S \le H+1$  - средняя длина кодового слова Статические коды Хаффмана .25 .25 .25 -33 .42 .58 1.0 .25 .20 .20 .22 .33 42 42 .15 .18 .20 -22 .25 Œ .12 .15 .18 .20 44 .12 .10 .15 45 .10 .10 44 .08 (a) 457

**(b)** 

delta

0100

0101

01100

01101

01110

01111

F(N)

11

011

0011

1011

00011

10011

01011

000011

0010011

00101011

1

0

0

1

0

1

0

0

0

0

1

Cumulative

probability

. 2

.6

.7

.9

1.0

Cumulative

probability

.05

.125

.225

.35

Range

[0,.2)

[.2,.6)

[.6,.7)

[.7,.9)

[.9, 1.0)

Range

[0,.05)

[.05, .125)

[.225,.35)

[.125,.225)

1

0

0

0

0

1

0

0

0

2

1

1

0

0

0

0

1

1

3

. 2

.1

.2

.1

.05

.1

.075

.125

Probability

00100000

001010000

001010001

0011000000

годного сообщения збыточности (Галлагер, 1978) я граница избыточности *а,* - сообщение  $p(a_i)$  - вероятность сообщения  $x_i$  - Кодовое слово  $-\lg p(x_i)$  - длина кодового слова  $H \leq S \leq H + 1$  - средняя длина кодового слова п - число кодовых слов Сравнение Хаффмана и Шеннона-Фано

p(n) - вероятность самого редкого исходного сообщения  $p(n) + \lg[(2\lg e)/e]$  - верхняя граница избыточности (Галлагер, 1978) p(n) + 0.086 - приблизительная верхняя граница избыточности ····S-F·····Huffman ....0.35......00......1 ....0.17.....01.....011 a(2)....0.17.....10.....010 a(4) · · · · · · · · 0.16 · · · · · · · · 110 · · · · · · · 001 a(5) · · · · · · · · 0.15 · · · · · · · · 111 · · · · · · · 000 Average ·codeword ·length ·2.31 · · · · · 2.30

O(n), где n - число исходных сообщений

O(l), где l - длина пути в кодовом дереве O(c), где c - число длин разных кодовых слов ^ Зависимость времени кодирования от входных данных // Какого то хера речь зашла про Украину и Майдан??? Сжал данные блять Универсальные коды Элиаса Два типа - гамма и дельта коды

011

00100

00101

00110

Время кодирования:

gamma 2 010

7 00111 8 0001000 16 000010000 17 000010001

3

4

5

6

32

Гамма коды

1. Сосчи битов N.

итать $L$ — количество значащих
в в двоичном представлении чис
итать $M$ — количество значащих
в в двоичном представлении чис
сать $M-1$ нулей и одну единиц
THY TOTAL CHINE
вой стороны дописать биты чис

00000100000

и представлении числа количество значащих и представлении числа нулей и одну единицу.

2. Сосчи битов L. 3. Запис 4. С правой стороны дописать биты числа Lбез старшей единицы. 5. С правой стороны дописать биты числа

N без старшей единицы ( $N_2$ ). Дельта коды 1. Сосчитать M — количество нулей во входном потоке до первой единицы. 2. Не включая единицу считать M битов.

Считанное число в сумме с  $2^M$  дает L. 3. Далее идут L-1 младших битов числа N. Считать их и к считанному числу прибавить  $2^{L-1}$  . Средняя длина  $S = c_1 (H + c_2)$ Особо не понадобятся судя по всему так что ебал учить это всё лол

R(N)1 1 1 0 0 0 5 Figure 3.7 -- Fibonacci Representations and Fibonacci

Codes.

Такой большой разрыв позволяет быстро набирать вес и получить достаточную избыточность

Позволяет прикол - если разница чисел закодированных слов достаточно мала (условно около 100), то два слова можно считать практически одинаковыми Арифметические коды

Тоже прикалываются с вероятностями охуеть

Probability Source message

А В

C

D

#

Source

message

а

b

c

d

[.35,.5).15 .5 e [.5,.675).175 .675 [.675,.875) . 2 .875 g [.875,1.0) .125 space Прикол судя по всему в том, что наиболее вероятные символы вероятнее попадутся. Весь метод - проекция встречаемости на отрезок [0, 1] Применение всей этой ебалы Коды Элиаса - позволяют быстро кодировать/декодировать на ходу благодаря отсутствию особой привязки к вероятностям, пусть они и сравнительно неэффективны В кодах Элиаса таблица соответствия составляется до кодирования, то есть её не надо передавать. Они интересны когда: • Заранее неизвестна последовательность • Имеется возможность составить таблицу кодирования до прохода Арифметические коды - каждому символу можно подобрать несколько способов кодирования. Интересны когда: • Символы в сообщениях появляются очень неравномерно При реализации сжатия на ПЛИСах/МК лучше выбирать более простые алгоритмы