Сжатие данных

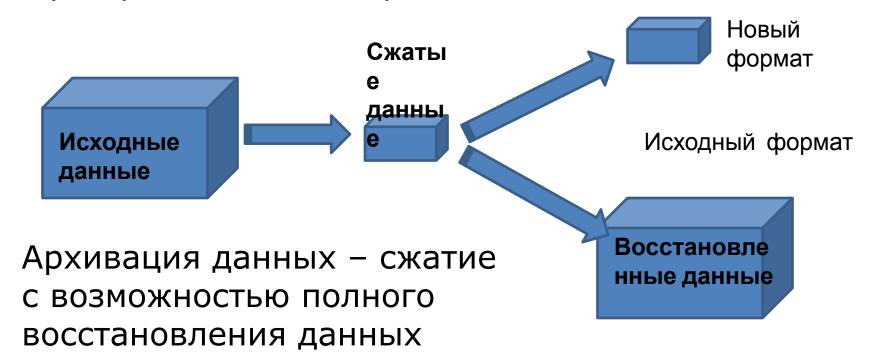


Сжатие данных это процесс, обеспечивающий уменьшение объема данных.

Способы сжатия

- •Изменение содержания данных (уменьшение избыточности данных)
- •Изменение структуры данных (эффективное кодирование)
- •Изменение содержания и структуры данных

Цели сжатия данных – экономия ресурсов при хранении или передаче данных



 Коэффициент сжатия – это величина для обозначения эффективности метода сжатия, равная отношению количества информации до и после сжатия

$$K_{cx} = 2 MB / 0.5 MB = 4$$

Исходные данные

Размер файла 2МБ Сжатые данные

Размер файла 512 КБ

Сжатие данных может происходить с потерями и без потерь

- Сжатие без потерь (полностью обратимое) это методы сжатия данных, при которых данные восстанавливаются после их распаковки полностью без внесения изменений (используется для текстов, программ) Ксж до 50%
- Сжатие с регулируемыми потерями это методы сжатия данных, при которых часть данных отбрасывается и не подлежит восстановлению (используется для видео, звука, изображений) Ксж до 99%

Сжатие с потерями

Тип данных	Тип файла после сжатия	Степень сжатия
Графика	.JPG	до 99%
Видео	.MPG	H
Звук	.MP3	

Сжатие без потерь

Тип данных	Тип файла после сжатия	Степень сжатия
Графика	.GIF .TIF .PCX	
Видео	.AVI	До 50%
Любой тип	.ZIP .ARJ .RAR .LZH	

Алгоритмы сжатия символьных данных

- Статистические методы это методы сжатия, основанные на статистической обработке текста.
- Словарное сжатие это методы сжатия, основанные на построении внутреннего словаря.

Упаковка однородных данных

Закодируем сообщение длиной 16 символов

В кодировке ASCII сообщение составляет 16 байт.

HOBBAR	қодова	ъ табъти	स्र साम्ब स्थ	дажрвки	¹ 5 0101	6 0110
7 0111	8 1000	9 1001	_ 1010	+ 1011	- 1100	, 1101

$$K_{CW} = 16 / 8 = 2$$

Достоинства и недостатки

метода

- + коэффициент сжатия увеличивается с увеличением размера символьного сообщения;
- необходимо указывать для распаковки новую кодовую таблицу;
- эффективен только для однородных сообщений, использующих ограниченный набор символов исходного алфавита;

Статистический метод сжатия Алгоритм Хаффмана

Разные символы встречаются в сообщении с разной частотой, например для русского алфавита в среднем на 1000 символов:

СИМВО	пробел	0	а	р	К	Я	Γ	Ю	ф	
Л										
частот	175	90	62	40	28	18	13	6	2	
Вададим коды символам согласно частоте их										

повторения:

чем чаще встречается символ, тем короче его код

(неравномерное кодирование)

Хаффмановское кодирование (сжатие)

– это метод сжатия, присваивающий символам алфавита коды переменной длины, основы-ваясь на частоте появления этих символов в сообщении.

Символ	код символа
пробел	00
0	01
p	101
К	110
Ю	0110
ф	1001

Проблема декодирования

Пример: пусть коды символов **a**-10, **b** -101, **c**-1010

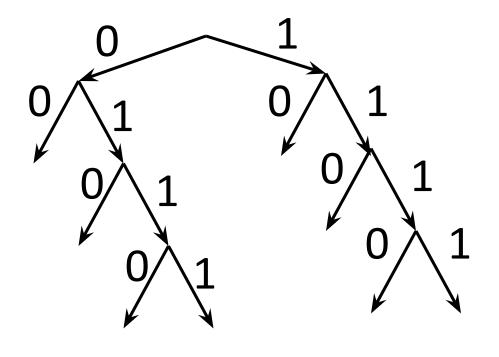
Декодировать сообщение 10101011010

Однозначное декодирование возможно при условии Фано: никакое кодовое слово не является началом (префиксом) другого кодового слова.

Префиксный код – это код, в котором никакое кодовое слово не является префиксом любого другого кодового

слова. Пример префиксного кода :

00 10 010 110 0110 0111 1110 1111



Префиксный код задается орграфом с размеченными листьями

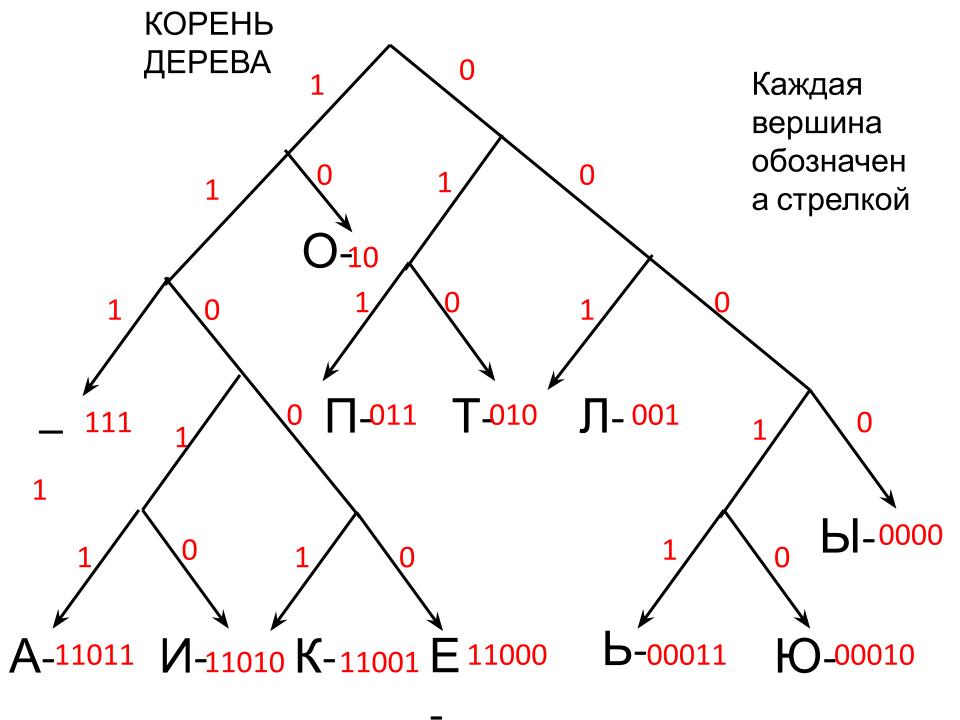
фразы

ОТ_ТОПОТА_КОПЫТ_ПЫЛЬ_ПО_ПОЛЮ_ЛЕТИ 1. Определим частоту рождения символов в фразу:

СИМВОЛ	Α	E	И	К	Л	O	П	T	Ы	Ь	Ю	
частота	1	1	1	1	3	6	5	6	2	1	1	6

2. Строим орграф Хаффмана:

- -символ задает вершину- лист орграфа;
- -вес вершины равен частоте вхождения символа;
- -соединяются пары вершин с наименьшим весом:
- -левые ветви обозначаем 0;
- -правые ветви обозначаем 1:



Построены префиксные коды символов:

Симво л	Α	Е	И	К	Л	O	П	T	Ы	Ь	Ю	_
Код	11011	11000	11010	11001	001	10	011	010	0000	00011	00010	111
Длина кода	5	5	5	5	3	2	3	3	4	5	5	3

Сообщение в новых кодах содержит 110 бит, в кодировке ASCII – 34 * 8 = 272 бита

тогда K cж = 272 / 110 = 2,47

Достоинства и недостатки метода

+

Алгоритм Хаффмана универсальный, его можно применять для сжатия данных любых типов;

Классический алгоритм Хаффмана требует хранения кодового дерева, что увеличивает размер файла.

Метод словарей **Алгоритм сжатия LZ**

Этот алгоритм был впервые описан в работах А. Лемпеля и Дж. Зива (Abraham Lempel, Jacob Ziv) в 1977-78 гг., поэтому этот метод часто называется Lempel-Ziv, или сокращенно LZ.

В его основе лежит идея замены наиболее часто встречающихся цепочек символов (строк) в файле ссылками на «образцы» цепочек, хранящиеся в специально создаваемой таблице (словаре).

Алгоритм <u>разра</u>ботан из<u>ра</u>ильскими <u>математ</u>ика<u>ми</u> Яко<u>бом Зивом и Аб рахам ом Лем</u>пе<u>лем</u>.

Словарь содержит, кроме многих других, такие цепочки: 1-ра 2-аб 3-ат 4-мат 5-ми_ 6-ам 7-бо 8-ом_ 9-бом 10-ем 11-лем

Алгоритм раз1ботан из1ильскими мате4ика5Яко7ом Зив821х68 Л10пе11

Чем длиннее цепочка, заменяемая ссылкой в словарь, тем больше эффект сжатия.

Достоинства и недостатки метода

- -применим для любых данных;
 - очень высокая скорость сжатия;
 - высок коэффициент сжатия;
 - - словарь настроен на тип текста;
 - словарь может быть очень большим;

Вопросы по теме:

- 1. Что такое архивирование данных? Для данных каких типов возможно применять архивирование?
- 2. Для каких данных допустимо сжатие с потерями?
- 3. При каких условиях метод упаковки неэффективен?
- 4. Что такое префиксный код?
- 5. Для каких данных метод Хаффмана эффективен?
- 6. На каких принципах основан метод словарного сжатия?
- 7. Назовите известные вам программы для сжатия данных.
- 8. Есть ли эффект от архивирования сжатых данных? Почему?
- 9. Изменилось ли количество информации в звукозаписи после сжатия с потерями? Поясните свой ответ.
- 10. Изменилось ли количество информации в изображении после его архивирования? Поясните свой ответ.

Домашнее задание: используя любые данные указанного типа, проведите эксперименты по архивированию. Результаты занесите в таблицу и поясните полученный эффект сжатия.

Тип данных	Исходны й формат	Исходный размер	Формат архива	Размер архива	Абсолютная величина сжатия(вМБ)	Коэффиц иент сжатия	Пояснени е эффекта сжатия
Текст	.doc						
	.pdf						
Видео	.avi						
	.mpg						
Изобра	.bmp						
жение	.jpeg						
Звук	.mp3						
	.midi						

Ксж=(исходный размер файла – размер файла архива) / исходный размер

файла