

# 哈夫曼编码

姓名:陈扬	任课老师:魏振刚
专业:17 计算机	专业课程:数据结构
学号:17150011001	实验四：哈夫曼（Huffman）编码问题

## 一、题目描述

哈夫曼编码是广泛地用于数据文件压缩的十分有效的编码方法。其压缩率通常在20%~90%之间。哈夫曼编码算法用字符在文件中出现的频率表来建立一个用0，1串表示各字符的最优表示方式。一个包含100,000个字符的文件，各字符出现频率不同，如下表所示：

	a	b	c	d	e	f
频率(千次)	45	13	12	16	9	5
定长码	000	001	010	011	100	101
变长码	0	101	100	111	1101	1100

有多种方式表示文件中的信息，若用0,1码表示字符的方法，即每个字符用唯一的一个0,1串表示。若采用定长编码表示，则需要3位表示一个字符，整个文件编码需要300,000位；若采用变长编码表示，给频率高的字符较短的编码；频率低的字符较长的编码，达到整体编码减少的目的，则整个文件编码需要  $(45\times1+13\times3+12\times3+16\times3+9\times4+5\times4)\times1000=224,000$ 位，由此可见，变长码比定长码方案好，总码长减小约25%。

## 二、算法设计与分析

1、前缀码：对每一个字符规定一个0,1串作为其代码，并要求任一字符的代码都不是其他字符代码的前缀。这种编码称为前缀码。编码的前缀性质可以使译码方法非常简单；例如001011101可以唯一的分解为0,0,101,1101，因而其译码为aabe。译码过程需要方便的取出编码的前缀，因此需要表示前缀码合适的数据结构。为此，可以用二叉树作为前缀码的数据结构：树叶表示给定字符；从树根到树叶的路径当作该字符的前缀码；代码中每一位的0或1分别作为指示某节点到左儿子或右儿子的“路标”。

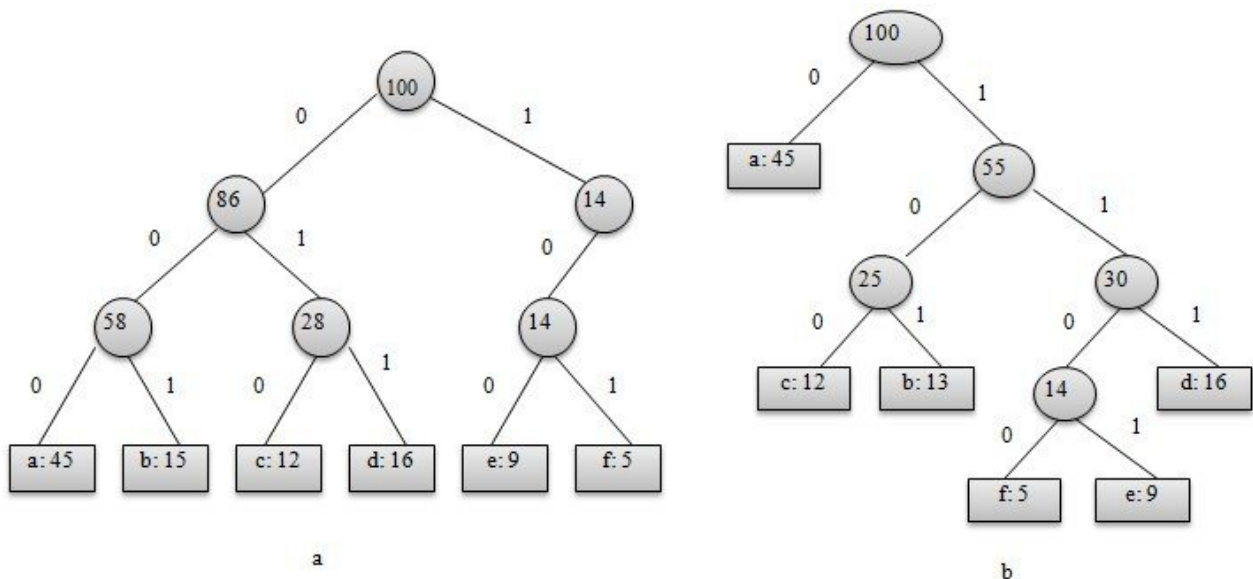


图-1a 与固定长度编码对应的树; 图-1b 对应于最优前缀编码的树 从上图可以看出, 表示最优前缀码的二叉树总是一棵完全二叉树, 即树中任意节点都有2个儿子。图a表示定长编码方案不是最优的, 其编码的二叉树不是一棵完全二叉树。在一般情况下, 若C是编码字符集, 表示其最优前缀码的二叉树中恰有 $|C|$ 个叶子。每个叶子对应于字符集中的一个字符, 该二叉树有 $|C|-1$ 个内部节点。给定编码字符集C及频率分布 $f$ , 即C中任一字符 $c$ 以频率 $f(c)$ 在数据文件中出现。C的一个前缀码编码方案对应于一棵二叉树T。字符 $c$ 在树T中的深度记为 $d_T(c)$ 。 $d_T(c)$ 也是字符 $c$ 的前缀码长。则平均码长定义为:

$$B(T) = \sum_{c \in C} f(c) d_T(c) \quad (1)$$

使平均码长达到最小的前缀码编码方案称为C的最优前缀码。

## 2、构造哈夫曼编码:

哈夫曼提出构造最优前缀码的贪心算法, 由此产生的编码方案称为哈夫曼编码。其构造步骤如下: (1) 哈夫曼算法以自底向上的方式构造表示最优前缀码的二叉树T。(2) 算法以 $|C|$ 个叶结点开始, 执行 $|C|-1$ 次的“合并”运算后产生最终所要求的树T。(3) 假设编码字符集中每一字符 $c$ 的频率是 $f(c)$ 。以 $f$ 为键值的优先队列Q用在贪心选择时有效地确定算法当前要合并的2棵具有最小频率的树。一旦2棵具有最小频率的树合并后, 产生一棵新的树, 其频率为合并的2棵树的频率之和, 并将新树插入优先队列Q。经过 $n-1$ 次的合并后, 优先队列中只剩下一棵树, 即所要求的树T。

构造过程如图-2所示:

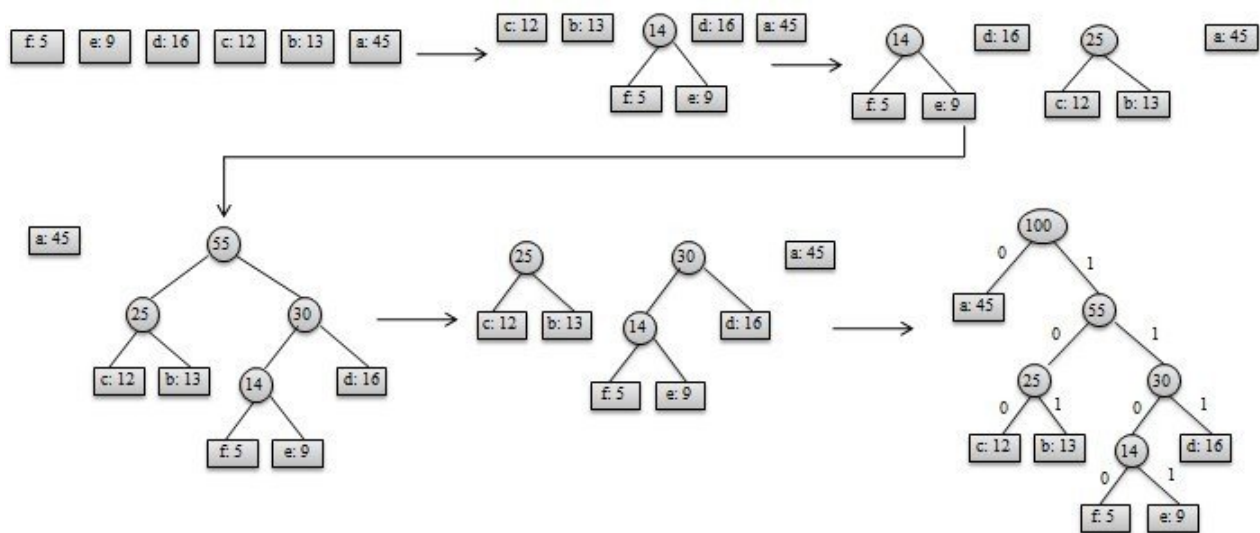


图-2 哈夫曼树构造过程

## 三、结果与分析

本实验以算法导论书中的例题为测试用例，来验证算法的正确性。即

实验结果截图如下图-7，结果与题目描述中给出的变长代码字一样：

	a	b	c	d	e	f
频率(千次)	45	13	12	16	9	5

## 四、实验总结

1、实验结果与给出的变长代码字一样，算法正确，且哈夫曼编码问题是一个贪心算法问题。采用哈夫曼编码技术可以最小化总的编码长度，从而实现数据文件的压缩存储。2、构造好哈夫曼树后，可用排列树回溯法来打印哈夫曼编码，即遇到左子树向左走，vector添加记录0；遇到右子树向右走，vector添加记录1；走到叶子节点并打印出叶节点的编码后回溯，同时往上退一层，则vector弹出一个值。如此不断回溯下去，即可打印所有字符编码。

## 五、源代码

```

1  class Node(object):
2      def __init__(self, name=None, value=None):
3          self._name = name
4          self._value = value
5          self._left = None
6          self._right = None
7  class HuffmanTree(object):
8      def __init__(self, char_weights):
9          self.a = [Node(part, char_weights[part]) for part in char_weights]
10
11      while len(self.a) != 1:

```

```

12         self.a.sort(key=lambda node:node._value,reverse=True)
13         c=Node(value=(self.a[-1]._value+self.a[-2]._value))
14         c._left=self.a.pop(-1)
15         c._right=self.a.pop(-1)
16         self.a.append(c)
17         self.root=self.a[0]
18         self.b=range(100)
19     def pre(self,tree,length):
20         node=tree
21         if (not node):
22             return
23         elif node._name:
24             print node._name + '\s encode is :',
25             for i in range(length):
26                 print self.b[i],
27             print '\n'
28             return
29         self.b[length]=0
30         self.pre(node._left,length+1)
31         self.b[length]=1
32         self.pre(node._right,length+1)
33
34     def get_code(self):
35         self.pre(self.root,0)
36
37 if __name__=='__main__':
38     with open("ACGAN.py","r") as f:
39         read_file= f.read()
40         print(str(read_file))
41         s=str(read_file)
42         resoult={}
43         for i in set(s):
44             resoult[i]=s.count(i)
45         print(resoult)
46         tree=HuffmanTree(resoult)
47         tree.get_code()

```

## 六、实验结果

```

1 from __future__ import print_function, division
2
3 from keras.datasets import mnist
4 from keras.layers import Input, Dense, Reshape, Flatten, Dropout,
multiply
5 from keras.layers import BatchNormalization, Activation, Embedding,
ZeroPadding2D
6 from keras.layers.advanced_activations import LeakyReLU

```

```

7 from keras.layers.convolutional import UpSampling2D, Conv2D
8 from keras.models import Sequential, Model
9 from keras.optimizers import Adam
10
11 import matplotlib.pyplot as plt
12
13 import numpy as np
14
15 class ACGAN():
16     def __init__(self):
17         # Input shape
18         self.img_rows = 28
19         self.img_cols = 28
20         self.channels = 1
21         self.img_shape = (self.img_rows, self.img_cols, self.channels)
22         self.num_classes = 10
23         self.latent_dim = 100
24
25         optimizer = Adam(0.0002, 0.5)
26         losses = ['binary_crossentropy',
27                  'sparse_categorical_crossentropy']
28
29         # Build and compile the discriminator
30         self.discriminator = self.build_discriminator()
31         self.discriminator.compile(loss=losses,
32                                   optimizer=optimizer,
33                                   metrics=['accuracy'])
34
35         # Build the generator
36         self.generator = self.build_generator()
37
38         # The generator takes noise and the target label as input
39         # and generates the corresponding digit of that label
40         noise = Input(shape=(self.latent_dim,))
41         label = Input(shape=(1,))
42         img = self.generator([noise, label])
43
44         # For the combined model we will only train the generator
45         self.discriminator.trainable = False
46
47         # The discriminator takes generated image as input and determines
48         validity
49         # and the label of that image
50         valid, target_label = self.discriminator(img)
51
52         # The combined model (stacked generator and discriminator)
53         # Trains the generator to fool the discriminator
54         self.combined = Model([noise, label], [valid, target_label])
55         self.combined.compile(loss=losses,

```

```

54         optimizer=optimizer)
55
56     def build_generator(self):
57
58         model = Sequential()
59
60         model.add(Dense(128 * 7 * 7, activation="relu",
input_dim=self.latent_dim))
61         model.add(Reshape((7, 7, 128)))
62         model.add(BatchNormalization(momentum=0.8))
63         model.add(UpSampling2D())
64         model.add(Conv2D(128, kernel_size=3, padding="same"))
65         model.add(Activation("relu"))
66         model.add(BatchNormalization(momentum=0.8))
67         model.add(UpSampling2D())
68         model.add(Conv2D(64, kernel_size=3, padding="same"))
69         model.add(Activation("relu"))
70         model.add(BatchNormalization(momentum=0.8))
71         model.add(Conv2D(self.channels, kernel_size=3, padding='same'))
72         model.add(Activation("tanh"))
73
74         model.summary()
75
76         noise = Input(shape=(self.latent_dim,))
77         label = Input(shape=(1,), dtype='int32')
78         label_embedding = Flatten()(Embedding(self.num_classes, 100)
(label))
79
80         model_input = multiply([noise, label_embedding])
81         img = model(model_input)
82
83         return Model([noise, label], img)
84
85     def build_discriminator(self):
86
87         model = Sequential()
88
89         model.add(Conv2D(16, kernel_size=3, strides=2,
input_shape=self.img_shape, padding="same"))
90         model.add(LeakyReLU(alpha=0.2))
91         model.add(Dropout(0.25))
92         model.add(Conv2D(32, kernel_size=3, strides=2, padding="same"))
93         model.add(ZeroPadding2D(padding=((0,1),(0,1))))
94         model.add(LeakyReLU(alpha=0.2))
95         model.add(Dropout(0.25))
96         model.add(BatchNormalization(momentum=0.8))
97         model.add(Conv2D(64, kernel_size=3, strides=2, padding="same"))
98         model.add(LeakyReLU(alpha=0.2))
99         model.add(Dropout(0.25))

```

```

100     model.add(BatchNormalization(momentum=0.8))
101     model.add(Conv2D(128, kernel_size=3, strides=1, padding="same"))
102     model.add(LeakyReLU(alpha=0.2))
103     model.add(Dropout(0.25))
104
105     model.add(Flatten())
106     model.summary()
107
108     img = Input(shape=self.img_shape)
109
110     # Extract feature representation
111     features = model(img)
112
113     # Determine validity and label of the image
114     validity = Dense(1, activation="sigmoid")(features)
115     label = Dense(self.num_classes+1, activation="softmax")(features)
116
117     return Model(img, [validity, label])
118
119 def train(self, epochs, batch_size=128, sample_interval=50):
120
121     # Load the dataset
122     (X_train, y_train), (_, _) = mnist.load_data()
123
124     # Configure inputs
125     X_train = (X_train.astype(np.float32) - 127.5) / 127.5
126     X_train = np.expand_dims(X_train, axis=3)
127     y_train = y_train.reshape(-1, 1)
128
129     # Adversarial ground truths
130     valid = np.ones((batch_size, 1))
131     fake = np.zeros((batch_size, 1))
132
133     for epoch in range(epochs):
134
135         # -----
136         # Train Discriminator
137         # -----
138
139         # Select a random batch of images
140         idx = np.random.randint(0, X_train.shape[0], batch_size)
141         imgs = X_train[idx]
142
143         # Sample noise as generator input
144         noise = np.random.normal(0, 1, (batch_size, 100))
145
146         # The labels of the digits that the generator tries to create
an
147         # image representation of

```

```

148         sampled_labels = np.random.randint(0, 10, (batch_size, 1))
149
150         # Generate a half batch of new images
151         gen_imgs = self.generator.predict([noise, sampled_labels])
152
153         # Image labels. 0-9 if image is valid or 10 if it is
generated (fake)
154         img_labels = y_train[idx]
155         fake_labels = 10 * np.ones(img_labels.shape)
156
157         # Train the discriminator
158         d_loss_real = self.discriminator.train_on_batch(imgs, [valid,
img_labels])
159         d_loss_fake = self.discriminator.train_on_batch(gen_imgs,
[fake, fake_labels])
160         d_loss = 0.5 * np.add(d_loss_real, d_loss_fake)
161
162         # -----
163         # Train Generator
164         # -----
165
166         # Train the generator
167         g_loss = self.combined.train_on_batch([noise,
sampled_labels], [valid, sampled_labels])
168
169         # Plot the progress
170         print ("%d [D loss: %f, acc.: %.2f%%, op_acc: %.2f%%] [G
loss: %f]" % (epoch, d_loss[0], 100*d_loss[3], 100*d_loss[4], g_loss[0]))
171
172         # If at save interval => save generated image samples
173         if epoch % sample_interval == 0:
174             self.save_model()
175             self.sample_images(epoch)
176
177     def sample_images(self, epoch):
178         r, c = 10, 10
179         noise = np.random.normal(0, 1, (r * c, 100))
180         sampled_labels = np.array([num for _ in range(r) for num in
range(c)])
181         gen_imgs = self.generator.predict([noise, sampled_labels])
182         # Rescale images 0 - 1
183         gen_imgs = 0.5 * gen_imgs + 0.5
184
185         fig, axs = plt.subplots(r, c)
186         cnt = 0
187         for i in range(r):
188             for j in range(c):
189                 axs[i,j].imshow(gen_imgs[cnt,:,:,:0], cmap='gray')
190                 axs[i,j].axis('off')

```



```

191         cnt += 1
192         fig.savefig("images/%d.png" % epoch)
193         plt.close()
194
195     def save_model(self):
196
197         def save(model, model_name):
198             model_path = "saved_model/%s.json" % model_name
199             weights_path = "saved_model/%s_weights.hdf5" % model_name
200             options = {"file_arch": model_path,
201                       "file_weight": weights_path}
202             json_string = model.to_json()
203             open(options['file_arch'], 'w').write(json_string)
204             model.save_weights(options['file_weight'])
205
206             save(self.generator, "generator")
207             save(self.discriminator, "discriminator")
208
209
210 if __name__ == '__main__':
211     acgan = ACGAN()
212     acgan.train(epochs=14000, batch_size=32, sample_interval=200)
213
214 {'\t': 347, '\n': 212, ' ': 484, '#': 32, '"': 40, '%': 17, "'": 22, ')':
155, '(' : 155, '+': 3, '*': 8, '-': 88, ',': 117, '/': 4, '.': 141, '1':
38, '0': 59, '3': 13, '2': 41, '5': 12, '4': 4, '7': 6, '6': 3, '9': 1,
'8': 12, ':': 21, '=': 113, '>': 1, 'A': 11, 'C': 11, 'B': 8, 'E': 3,
'D': 25, 'G': 5, 'F': 5, 'I': 9, 'M': 4, 'L': 11, 'N': 8, 'P': 3, 'S': 8,
'R': 8, 'U': 8, 'T': 9, 'X': 7, '[': 29, 'Z': 2, ']': 29, '_': 140, 'a':
404, 'c': 105, 'b': 58, 'e': 501, 'd': 258, 'g': 114, 'f': 106, 'i': 333,
'h': 89, 'k': 29, 'j': 7, 'm': 228, 'l': 291, 'o': 283, 'n': 279, 'q': 3,
'p': 141, 's': 316, 'r': 245, 'u': 62, 't': 279, 'w': 14, 'v': 45, 'y':
32, 'x': 11, '{': 1, 'z': 27, '}': 1}
215 l's encode is : 0 0 0 0
216
217 ('s encode is : 0 0 0 1 0
218
219 )'s encode is : 0 0 0 1 1
220
221 s's encode is : 0 0 1 0
222
223 i's encode is : 0 0 1 1
224
225 1's encode is : 0 1 0 0 0 0 0
226
227 "'s encode is : 0 1 0 0 0 0 1
228
229 2's encode is : 0 1 0 0 0 1 0
230

```

```
231 Z's encode is : 0 1 0 0 0 1 1 0 0 0 0
232
233 q's encode is : 0 1 0 0 0 1 1 0 0 0 1
234
235 F's encode is : 0 1 0 0 0 1 1 0 0 1
236
237 G's encode is : 0 1 0 0 0 1 1 0 1 0
238
239 6's encode is : 0 1 0 0 0 1 1 0 1 1 0
240
241 +'s encode is : 0 1 0 0 0 1 1 0 1 1 1
242
243 :'s encode is : 0 1 0 0 0 1 1 1
244
245 -'s encode is : 0 1 0 0 1 0
246
247 C's encode is : 0 1 0 0 1 1 0 0 0
248
249 A's encode is : 0 1 0 0 1 1 0 0 1
250
251 x's encode is : 0 1 0 0 1 1 0 1 0
252
253 L's encode is : 0 1 0 0 1 1 0 1 1
254
255 v's encode is : 0 1 0 0 1 1 1
256
257 's encode is : 0 1 0 1
258
259 h's encode is : 0 1 1 0 0 0
260
261 ''s encode is : 0 1 1 0 0 1 0 0
262
263 P's encode is : 0 1 1 0 0 1 0 1 0 0 0
264
265 E's encode is : 0 1 1 0 0 1 0 1 0 0 1
266
267 7's encode is : 0 1 1 0 0 1 0 1 0 1
268
269 8's encode is : 0 1 1 0 0 1 0 1 1
270
271 5's encode is : 0 1 1 0 0 1 1 0 0
272
273 3's encode is : 0 1 1 0 0 1 1 0 1
274
275 D's encode is : 0 1 1 0 0 1 1 1
276
277 c's encode is : 0 1 1 0 1 0
278
279 f's encode is : 0 1 1 0 1 1
```

```
280
281 a's encode is : 0 1 1 1
282
283
284 's encode is : 1 0 0 0 0
285
286 z's encode is : 1 0 0 0 1 0 0 0
287
288 j's encode is : 1 0 0 0 1 0 0 1 0 0
289
290 X's encode is : 1 0 0 0 1 0 0 1 0 1
291
292 w's encode is : 1 0 0 0 1 0 0 1 1
293
294 k's encode is : 1 0 0 0 1 0 1 0
295
296 ]'s encode is : 1 0 0 0 1 0 1 1
297
298 ='s encode is : 1 0 0 0 1 1
299
300 m's encode is : 1 0 0 1 0
301
302 g's encode is : 1 0 0 1 1 0
303
304 b's encode is : 1 0 0 1 1 1 0
305
306 0's encode is : 1 0 0 1 1 1 1
307
308 's encode is : 1 0 1 0
309
310 , 's encode is : 1 0 1 1 0 0
311
312 ['s encode is : 1 0 1 1 0 1 0 0
313
314 U's encode is : 1 0 1 1 0 1 0 1 0 0
315
316 R's encode is : 1 0 1 1 0 1 0 1 0 1
317
318 4's encode is : 1 0 1 1 0 1 0 1 1 0 0
319
320 /'s encode is : 1 0 1 1 0 1 0 1 1 0 1
321
322 >'s encode is : 1 0 1 1 0 1 0 1 1 1 0 0 0
323
324 9's encode is : 1 0 1 1 0 1 0 1 1 1 0 0 1
325
326 }'s encode is : 1 0 1 1 0 1 0 1 1 1 0 1 0
327
328 {'s encode is : 1 0 1 1 0 1 0 1 1 1 0 1 1
```

```
329
330 M's encode is : 1 0 1 1 0 1 0 1 1 1 1
331
332 u's encode is : 1 0 1 1 0 1 1
333
334 r's encode is : 1 0 1 1 1
335
336 e's encode is : 1 1 0 0
337
338 d's encode is : 1 1 0 1 0
339
340 B's encode is : 1 1 0 1 1 0 0 0 0 0
341
342 *'s encode is : 1 1 0 1 1 0 0 0 0 1
343
344 S's encode is : 1 1 0 1 1 0 0 0 1 0
345
346 N's encode is : 1 1 0 1 1 0 0 0 1 1
347
348 y's encode is : 1 1 0 1 1 0 0 1
349
350 #'s encode is : 1 1 0 1 1 0 1 0
351
352 %'s encode is : 1 1 0 1 1 0 1 1 0
353
354 T's encode is : 1 1 0 1 1 0 1 1 1 0
355
356 I's encode is : 1 1 0 1 1 0 1 1 1 1
357
358 _'s encode is : 1 1 0 1 1 1
359
360 t's encode is : 1 1 1 0 0
361
362 n's encode is : 1 1 1 0 1
363
364 p's encode is : 1 1 1 1 0 0
365
366 .'s encode is : 1 1 1 1 0 1
367
368 o's encode is : 1 1 1 1 1
369
```