

目次

- ★ 3.7.2 ハフマン符号化
- ★ 3.7.3 よだんだよん
- ★ 4.1 成分を分析したい

★3 パターン情報の表現(2) — データ圧縮と情報量(承前)

「承前」って? ⇒ 辞書引きましょう

★ 3.7.2 ハフマン符号化

エントロピー符号化の一種である**ハフマン符号化**は、記号毎に算出した出現頻度を用いて、次のような手順で木(ハフマン木)を作ることで、それぞれの記号に割り当てる符号を求める(授業中にもう少し詳しく説明します)。

1. 個々の記号に対応した葉節点をつくる。出現確率をそれぞれの値とする。
2. 親のない節点の中で最も小さい値をもつもの2つを選び、それらの親となる節点をつくる。親節点の値は、2つの子節点の値の和とする。
3. 親を持たない節点が一つになるまで2.を繰り返す。
4. 「左」を0, 「右」を1として、できた木の各節点から左右に接続した枝のそれぞれに0,1を割り当てる(「左」を1としてもよい)。
5. 根から葉までたどる際に現れる0,1の並びをそれぞれの葉の記号に対応した符号とする。

ハフマン符号化は純粋なエントロピー符号化法であり、記号の出現頻度のみを考慮して符号を決めるアルゴリズムである。したがって、出現頻度に偏りがないと有効な圧縮ができない(☆1)。また、特定の記号が何度も繰り返すとか、特定の並び(HOGEみたいな)が何度も出てくるとか、ある記号のあとに特定の記号が続くことが多い(GのあとにはEがよく出てくるみたいな)というような記号の並び方は、データ圧縮の性能に影響しない(それを活かした圧縮はできない)。

ハフマン符号化: Huffman coding. 1950年代前半に D. A. Huffman が提案した符号化法。画像圧縮方式 JPEG でも、離散コサイン変換(そのうち紹介するかも)したデータの符号化に用いられている。

☆1) 復号の際にハフマン木を使うので、データにハフマン木の情報も加える必要がある。そのため、極端な場合には元よりデータ量が増えてしまうこともありえる。

Q1. 上の(授業中に説明した) Huffman 木を用いて、次のデータを復号しなさい。1011000111010101110

Q2. HOGEHOGEHOGEHOGE... とひたすら繰り返すだけのテキストデータがあったとする。このデータにハフマン符号化を適用しても、有効なデータ圧縮はできない。その理由を述べなさい。

★ 3.7.3 よだんだよん

- エントロピー符号化の方法としては、Huffman 符号化よりも高効率な算術符号化という方法もある。
- 符号化は何もデータ圧縮のためだけに行なうものではない。雑音等のせいで送信データに誤りが生じてしまう状況で誤りを見つける（誤り検出）／誤りを訂正する（誤り訂正）ためにも用いられる（☆2）。最も原始的な誤り検出の符号化法は、ビット列の最後に、ビット列中の1の数が常に偶数になるように0か1を付け加える、というものである。興味のある人はいろいろ調べてみるとおもしろいかも。
- 符号化のさらに別の目的として、暗号化がある。

☆2) 情報通信機器では当たり前に使われている。図書のISBN、様々な商品のJANコード（バーコードの元になっているコード）などにも使われている。

情報理論の参考書：高橋の手元には「情報理論」甘利俊一，ダイヤモンド社，ISBN4-478-82000-7，といういい本があるのですが，残念ながら現在は手に入らないようです（☆3）。他にいい本ないか探索中。

☆3) 最近，ちくま学芸文庫で復活しました。

データ圧縮の参考書：「圧縮アルゴリズム 符号化の原理とC言語による実装」昌達 K'z，ソフトバンクパブリッシング，ISBN4-7973-2552-6

Q3.

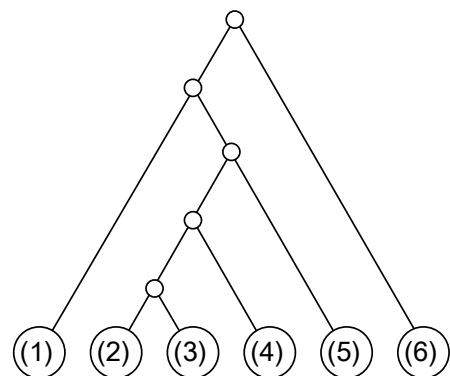
ある記号列を調べたところ，6通りの記号が次表のような頻度で出現することがわかった。このような記号列をハフマン符号化法でビット列に符号化するとして，次の間に答えなさい（答えのみでよい）。

記号	あ	か	0	1	2	3
頻度	0.07	0.14	0.2	0.03	0.55	0.01

(1)–(6) 上記の表にもとづいてハフマン木を求めると，右図のようになった。 (1) から (6) までの各葉節点に割り当てられる記号を答えなさい。ただし，このハフマン木は，授業中に説明したものと同様に，子をもつ全ての節点について，(左の子の値) < (右の子の値) となるように描かれているものとする。

(7) 右図のハフマン木の子をもつ全ての節点について，**左の子にむかってのびた枝に符号1**を，**右の子にむかってのびた枝に符号0** このハフマン木を用いて，次のビット列を記号列に復号しなさい。

110100010111



★4 パターン情報の成分分析(1) — 直交展開

今回からしばらくは、パターン情報の成分を分析する方法を考える。今回と次回には特に、ベクトルで表されるパターンの成分分析を扱う。

★4.1 成分を分析したい

「成分を分析する」とはどういうことか、スープの成分を分析するという例（強引な例やけど）で考えてみよう。

スープの味が、塩味や甘味などのいくつかの味の成分の量で決まっており、それらの量は、それぞれの味に対応した調味料の含有量ではかれるものとする。このとき、あるスープに各成分がどれ位ずつ含まれているかがわかれば、右図のようなグラフを描ける。

このような情報が得られることには、次のような利点がある。

- そのスープの特徴（どんな成分を多く含んでいるか）がよくわかる
- ある成分の量を増減させることで味を調節できる
- 成分の量が分かっているので、別のところでそのスープの味を再現できる

音響信号や画像のようなパターン情報でも、これと同じようなことをできると便利そうである(☆4)。というわけで、パターンを

$$(\text{あるパターン}) = \alpha_1 \times (\text{成分1}) + \alpha_2 \times (\text{成分2}) + \alpha_3 \times (\text{成分3}) + \dots \quad (1)$$

のように成分に分解して表す方法を考えてみよう。その際には、次のようなことに気をつけないといけないだろう。

- 何を「成分」とすればよいのか、「成分」はいくつあればよいのか
- $\alpha_1, \alpha_2, \dots$ はどうやって求めるのか



☆4) パターンの特徴をつかめる、特定の成分を増減させてパターンを加工できる、各成分の含有量の情報を伝送するだけでパターンを再構成できる。