

データサイエンス実践B 2022/7/4 講義資料

H. Kumano

情報の基本単位

コンピュータは、一度に0か1で表される、2通りの情報しか判断できない

選択肢が2つであれば、一回で判断可能。

例：旅行に行くことにしました。国内/海外どちらにするか？



「海外にしますか？」という1つの質問に対して、「はい」か「いいえ」どちらかを選べば良い。

1: 海外	0: 国内
-------	-------

情報の基本単位

では、選択肢が4つある場合は？

北欧	北海道
地中海・エーゲ海	沖縄

一度の質問では決められない。が、こういう場合、何回か続けて質問すれば決められる。

1: 海外にしますか? はい	0: 海外にしますか? いいえ
北欧	北海道
地中海・エーゲ海	沖縄

一度回答すれば、選択肢が2つに減る。

続けての質問：北にしますか？

→ 一つに決まる

情報の基本単位

今度は、選択肢が8つの場合はどうでしょう？

いくつかの質問で、1つに絞れるでしょうか？

1つ目の質問：8 → 4

2つ目の質問：4 → 2

3つ目の質問：2 → 1

1024個選択肢がある場合は？

1024 → 512 → 256 → 128 → 64 → 32 → 16 → 8 → 4 → 2 → 1

10回の質問で決められる。

(→ の数を数えればよい)

情報の基本単位

コンピュータでは、計算を

0 (電気が通っていない状態)

1 (電気が取っている状態)

の組み合わせで行っている。

0と1以外は判別できない。



1: 海外	0: 国内
-------	-------

2個の情報を表すにはスイッチが1個あればよい。

北欧	北海道
地中海・エーゲ海	沖縄

4個の情報を表すにはスイッチが2個あればよい。



文字コード

文字コードとは、文字を表す番号のこと。

コンピュータは数字しか処理できません。よってコンピュータは、それぞれの文字の形 (イメージ) に対応付けられた番号の一覧表で文字を管理しています。それが文字コードです。

文字コードはいくつか種類があり、それにより見た目が同じでも番号が異なります。

情報の基本単位

同様に、1024個の情報を表すには10個スイッチがあればよい。

$$2^1 = 2 \quad 2^2 = 4 \quad 2^{10} = 1024$$



= 情報の基本単位

2通りの情報が表せるもの

binary digit (→ bit)

情報を数値化
符号化



Claude E. Shannon

計算機上では、すべての情報が
(文字も画像も) 符号化される

主な文字コード(1)

ASCII

もっとも基礎的な文字コード。半角英数字128文字から構成され、全ての文字を1バイトで表す。例えば「A」はASCIIでは0x41(0xは16進数を表す)。

Shift_JIS

日本語を表すために多く用いられていた文字コード。全ての文字を2バイトで表す。亜種の cp932 が Windows で採用されていたことで広く使われていた。

各言語圏の文字コード

地域	1960年代	1970年代	1980年代	1990年代~
ラテン文字圏	ASCII/ISO 646	欧州言語への拡張(6937, 8859)		
キリル文字圏	GOST 13052	露語以外の諸語への拡張		
アラビア文字圏			ASMO 449	
ヘブライ文字圏			ECMA 121	
日本	JIS C 6220	JIS C 6226		
中国			GB 2312	少数民族文字
韓国			KS C 5601	KS X 1005
タイ			TIS 620	
インド			ISSCII 83	IS 13194
ベトナム				TCVN 5412
スリランカ				SLS 1134
国際符号化文字集合				ISO/IEC 10646

数字の表現 (進数)

	2進数	10進数	16進数
	0001	1	1
	0010	2	2
	0001 0000	16	10
	0110 0100	100	64
	0001 0000 0000	256	100
	1 0000 0000 0000 0000	65536	10000

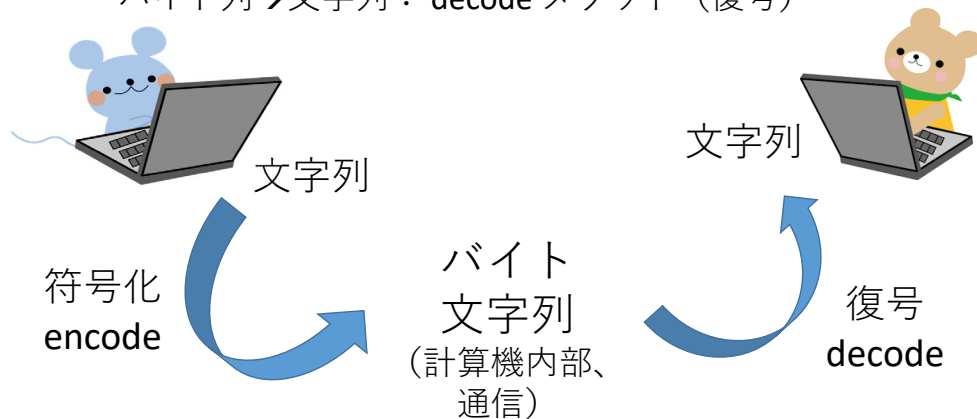
16進数世界の文字
0,1,2,3,4,5,6,7,8,9,a,b,c,d,e,f

文字列とバイト列の変換

文字列とバイト列との変換用関数

文字列→バイト列：encode メソッド (符号化)

バイト列→文字列：decode メソッド (復号)



文字集合の例

日本語

文字集合	含まれる文字
JIS X 0211	JIS制御コード
JIS X 0201	JISローマ字 (ASCII)
JIS X 0201	JISカナ (半角カナ)
JIS X 0208	JIS漢字
JIS X 0212	JIS補助漢字

Unicode

ユニコード・コンソーシアムによって制定された、世界中の文字を表現しようとする世界統一文字集合規格。