Why represent numbers

Modeling uses data-types that are encoded in a computer. The details of the encoding impact the efficiency of algorithms we use to understand the systems we are modeling and the impacts of these algorithms on the people using the systems. Case study: how to encode numbers?

Fixed width definition

Definition For b an integer greater than 1, w a positive integer, and n a nonnegative integer _____, the base b fixed-width w expansion of n is

$$(a_{w-1}\cdots a_1a_0)_{b,w}$$

where $a_0, a_1, \ldots, a_{w-1}$ are nonnegative integers less than b and

$$n = \sum_{i=0}^{w-1} a_i b^i$$

Fixed width example

Decimal	Binary	Binary fixed-width 10	Binary fixed-width 7	Binary fixed-width 4
b = 10	b=2	b = 2, w = 10	b = 2, w = 7	b = 2, w = 4
$(20)_{10}$				

Fixed width fractional definition

Definition For b an integer greater than 1, w a positive integer, w' a positive integer, and x a real number the base b fixed-width expansion of x with integer part width w and fractional part width w' is $(a_{w-1} \cdots a_1 a_0.c_1 \cdots c_{w'})_{b,w,w'}$ where $a_0, a_1, \ldots, a_{w-1}, c_1, \ldots, c_{w'}$ are nonnegative integers less than b and

$$x \ge \sum_{i=0}^{w-1} a_i b^i + \sum_{j=1}^{w'} c_j b^{-j}$$
 and $x < \sum_{i=0}^{w-1} a_i b^i + \sum_{j=1}^{w'} c_j b^{-j} + b^{-w'}$



Note: Java uses floating point, not fixed width representation, but similar rounding errors appear in both.

Negative int expansions

Representing negative integers in binary: Fix a positive integer width for the representation w, w > 1.

	To represent a positive integer n	To represent a negative integer $-n$
Sign-magnitude	$[0a_{w-2}\cdots a_0]_{s,w}$, where $n=(a_{w-2}\cdots a_0)_{2,w-1}$ Example $n=17, w=7$:	$[1a_{w-2}\cdots a_0]_{s,w}$, where $n=(a_{w-2}\cdots a_0)_{2,w-1}$ Example $-n=-17, w=7$:
2s complement	$[0a_{w-2}\cdots a_0]_{2c,w}$, where $n=(a_{w-2}\cdots a_0)_{2,w-1}$ Example $n=17, w=7$:	$[1a_{w-2}\cdots a_0]_{2c,w}$, where $2^{w-1}-n=(a_{w-2}\cdots a_0)_{2,w-1}$ Example $-n=-17, w=7$:

Calculating 2s complement

For positive integer n, to represent -n in 2s complement with width w,

- Calculate $2^{w-1} n$, convert result to binary fixed-width w 1, pad with leading 1, or
- Express -n as a sum of powers of 2, where the leftmost 2^{w-1} is negative weight, or
- Convert n to binary fixed-width w, flip bits, add 1 (ignore overflow)

Challenge: use definitions to explain why each of these approaches works.

Representing zero

Representing 0:

So far, we have representations for positive and negative integers. What about 0?

	To represent a non-negative integer n	To represent a non-positive integer $-n$
Sign-magnitude	$[0a_{w-2}\cdots a_0]_{s,w}$, where $n=(a_{w-2}\cdots a_0)_{2,w-1}$ Example $n=0, \ w=7$:	$[1a_{w-2}\cdots a_0]_{s,w}$, where $n=(a_{w-2}\cdots a_0)_{2,w-1}$ Example $-n=0, w=7$:
2s complement	$[0a_{w-2}\cdots a_0]_{2c,w}$, where $n=(a_{w-2}\cdots a_0)_{2,w-1}$ Example $n=0, w=7$:	$[1a_{w-2}\cdots a_0]_{2c,w}$, where $2^{w-1}-n=(a_{w-2}\cdots a_0)_{2,w-1}$ Example $-n=0, w=7$: