

Statistical Computing

Class Notes Fall 2019

Floating Point Numbers and Finite Arithmetic

Overview

- 1 Positional Numerical System
- 2 Floating Point Number - IEEE standards
- 3 Floating Point Approximations and Errors
- 4 Numerical Operations - Fundamental Axiom

I. Positional Numerical System

Define: A positional numerical system is a way of representing numbers and is characterized by:

- 1 base, $b \in \{2, 3, 4, \dots\}$
- 2 a sequence of digits, d_0, d_1, \dots, d_{k-1} , called the digits
- 3 an exponent, $e \in \mathbb{Z}$

The number is denoted as $d_0.b.d_1.d_2 \dots d_{k-1} \cdot 10^e$

(d₀, d₁, ..., d_{k-1}) is called the significand

(e) is called the exponent

$b^{k-1} \leq d_0.d_1 \dots d_{k-1} < b^k$

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Definition: A floating point representation is normalized if $d_0 \neq 0$.

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Definition: A floating point representation is normalized if $d_0 \neq 0$.

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

+0.01 and 0.01

-0.10 and -0.01

+1.00 and +0.01

-1.00 and -0.01

Example: Suppose $b=10$, $p=7$, $e_{min}=-1$, $e_{max}=1$.

The 0.1 is a floating point number.

This image shows a scanned page from a linear algebra textbook, likely from a course at MIT, containing handwritten notes and solutions to various problems. The page is filled with mathematical notation, including vectors, matrices, and equations. Key topics covered include:

- Orthogonalization:** Definitions of orthogonal sets, orthonormal sets, and orthogonal projections.
- Krylov Spaces:** Formulas for the dimension of the Krylov space \$K_m(A, b)\$ and its properties.
- Row Action Methods:** A section on Kaczmarz's Method, Generalized Kaczmarz Method (GKR), and Randomized Kaczmarz (RK).
- Conjugate Gradients:** A detailed section on the conjugate gradient method, including its implementation and convergence properties.

The page also features several proofs, lemmas, and related questions, along with some diagrams and sketches. The handwriting is in blue ink, and the page is filled with mathematical symbols and calculations.