

# Floating-point numbers

- ▶ Analytical binary representation
- ▶ Floating-point number system  $\mathbb{F}(2, M, e_{\min}, e_{\max})$
- ▶ Ill-posed problems
- ▶ Computing errors and equality
- ▶ float, double

# Floating-point representation 1/2

► **Theorem:** For each  $x \in \mathbb{R}$  there exist

- Sign  $\sigma \in \{\pm 1\}$
- Digits  $a_k \in \{0, 1\}$
- Exponent  $e \in \mathbb{Z}$

such that  $x = \sigma \left( \sum_{k=1}^{\infty} a_k 2^{-k} \right) 2^e$

► The representation is not unique, e.g.,  $1 = \sum_{k=1}^{\infty} 2^{-k}$

## Remarks

► The result holds for any basis  $b \in \mathbb{N}_{\geq 2}$

- The digits then satisfy  $a_j \in \{0, 1, \dots, b-1\}$

► Decimal system  $b = 10$  is very common

- $47.11 = (4 \cdot 10^{-1} + 7 \cdot 10^{-2} + 1 \cdot 10^{-3} + 1 \cdot 10^{-4}) \cdot 10^2$

$$* \quad a_1 = 4, a_2 = 7, a_3 = 1, a_4 = 1, e = 2$$

► For  $b = 2$ , fractions are representable as finite sums if and only if the denominator is a power of 2:

- In  $\sum_{k=1}^M 2^{-k}$ , the denominator is a power of 2

- Uniqueness of the integer factorization

► e.g., no exact representation for  $1/10$  with  $b = 2$

# Floating-point representation 2/2

► **Theorem:** For each  $x \in \mathbb{R}$  there exist

- Sign  $\sigma \in \{\pm 1\}$
- Digits  $a_k \in \{0, 1\}$
- Exponent  $e \in \mathbb{Z}$

such that  $x = \sigma \left( \sum_{k=1}^{\infty} a_k 2^{-k} \right) 2^e$

► The representation is not unique, e.g.,  $1 = \sum_{k=1}^{\infty} 2^{-k}$

## Floating-point numbers

► Floating-point number system

$$\mathbb{F}(2, M, e_{\min}, e_{\max}) \subset \mathbb{Q}$$

- Mantissa length  $M \in \mathbb{N}$
- Exponential barriers  $e_{\min} < 0 < e_{\max}$

►  $x \in \mathbb{F}$  has representation  $x = \sigma \left( \sum_{k=1}^M a_k 2^{-k} \right) 2^e$  with

- Sign  $\sigma \in \{\pm 1\}$
- Digits  $a_j \in \{0, 1\}$  mit  $a_1 = 1$ 
  - \* So-called normalized floating-point number
- Exponent  $e \in \mathbb{Z}$  with  $e_{\min} \leq e \leq e_{\max}$

► The representation of  $x \in \mathbb{F}$  is unique

► Digit  $a_1$  must not be stored

- Implicit first bit

# Proof of the theorem

- ▶ Without loss of generality, assume  $x \geq 0$

- If not, multiply by  $\sigma = -1$

- ▶ Let  $e \in \mathbb{N}_0$  with  $0 \leq x < 2^e$

- ▶ Without loss of generality, assume  $x < 1$

- If not, divide by  $\sigma = -1$

- ▶ Construction of the digits  $a_j$  via bisection

- ▶ **Claim:** There exist digits  $a_j \in \{0, 1\}$  such that

$$x_n := \sum_{k=1}^n a_k 2^{-k} \text{ satisfies } x \in [x_n, x_n + 2^{-n})$$

- ▶ **Induction base case:** It holds that  $x \in [0, 1)$

- If  $x \in [0, 1/2)$ , choose  $a_1 = 0$ , i.e.,  $x_1 = 0$

- If  $x \in [1/2, 1)$ , choose  $a_1 = 1$ , i.e.,  $x_1 = 1/2$

- \*  $x_1 = a_1/2 \leq x$

- \*  $x < (a_1 + 1)/2 = x_1 + 2^{-1}$

- ▶ **Induction step:** It holds that  $x \in [x_n, x_n + 2^{-n})$

- If  $x \in [x_n, x_n + 2^{-(n+1)})$ , choose  $a_{n+1} = 0$ ,

- i.e.,  $x_{n+1} = x_n$

- If  $x \in [x_n + 2^{-(n+1)}, x_n + 2^{-n})$ , choose  $a_{n+1} = 1$

- \*  $x_{n+1} = x_n + a_{n+1}2^{-(n+1)} \leq x$

- \*  $x < x_n + (a_{n+1} + 1)2^{-(n+1)} = x_{n+1} + 2^{-(n+1)}$

- ▶ It follows that  $|x_n - x| \leq 2^{-n}$ , hence  $x = \sum_{k=1}^{\infty} a_k 2^{-k}$

## Arithmetic for floating-point numbers

- ▶ Result **Inf**, **-Inf** if overflow (oder **1./0.**)
- ▶ Result **NaN**, if not defined (z.B. **0./0.**)
- ▶ Arithmetic is approximate, not exact

## Ill-posed problem

- ▶ A problem is **numerically ill-posed**, if small changes in the data lead to large changes in the result
  - e.g., does a triangle with given side lengths have a right angle?
  - e.g., is a given point located on a circle?
- ▶ **Implementation meaningless, as result random!**

## Computing error

- ▶ In view of computing errors, one should *never* test the equality of two floating-point numbers
  - Check whether  $|x - y|$  is small, rather than  $x = y$
  - e.g.,  $|x - y| \leq \varepsilon \cdot \max\{|x|, |y|\}$  with  $\varepsilon = 10^{-13}$

```
1 #include <stdio.h>
2 #include <math.h>
3
4 main() {
5     double x = (116./100.)*100.;
6
7     printf("x=%f\n",x);
8     printf("floor(x)=%f\n",floor(x));
9
10    if (x==116.) {
11        printf("There holds x==116\n");
12    }
13    else {
14        printf("Surprise, surprise!\n");
15    }
16 }
```

- ▶ Output:

x=116.000000

floor(x)=115.000000

Surprise, surprise!

## Variable types float, double

```
1 #include <stdio.h>
2 main() {
3     double x = 2./3.;
4     float y = 2./3.;
5     printf("%f, %1.16e\n", x, x);
6     printf("%f, %1.7e\n", y, y);
7 }
```

► Floating-point numbers constitute a finite subset of  $\mathbb{Q}$

► **float** is usually *single precision* according to the IEEE-754 standard

- $\mathbb{F}(2, 24, -126, 127) \rightarrow 4$  bytes
- ca. 7 significant decimal digits

► **double** is usually *double precision* according to the IEEE-754 standard

- $\mathbb{F}(2, 53, -1022, 1023) \rightarrow 8$  bytes
- ca. 16 significant decimal digits

► Placeholders in **printf** and **scanf**

Data type	<b>printf</b>	<b>scanf</b>
float	%f	%f
double	%f	%lf

- Placeholder %1.16e for floating-point representation
- See **man 3 printf**

► Output:

```
0.666667, 6.66666666666666663e-01
0.666667, 6.6666669e-01
```

# Structures

- ▶ Why structures?
  - ▶ Members
  - ▶ Point operator .
  - ▶ Arrow operator ->
  - ▶ Shallow copy vs. deep copy
- 
- ▶ struct
  - ▶ typedef



# Declaration of structures

## ► Functions

- Callable group of statements that together perform a task
- Abstraction (structured programming)

## ► Structures

- Combination of variables of different types in a new data type
- Abstraction with data

## ► **Example:** Management of SciProgMath students

- The same data for each student

```
1 // Declaration of structure
2 struct _Student_ {
3     char* firstname; // First name
4     char* lastname;  // Last name
5     int studentID;    // Student ID
6     int studiesID;    // Studies ID
7     double test;      // Points in final test
8     double exercise;  // Points in exercises
9 };
10
11 // Declaration of corresponding data type
12 typedef struct _Student_ Student;
```

## ► Semicolon after structure declaration block

## ► Creation of new variable type Student

# Structures & Members

- ▶ Data types of a structure are called **members**
- ▶ Access to members with point operator
  - **var** variable of type **Student**
  - e.g., member **var.firstname**

```
1 // Declaration of structure
2 struct _Student_ {
3     char* firstname; // First name
4     char* lastname;  // Last name
5     int studentID;   // Student ID
6     int studiesID;   // Studies ID
7     double test;     // Points in final test
8     double exercise; // Points in exercises
9 };
10
11 // Declaration of corresponding data type
12 typedef struct _Student_ Student;
13
14 main() {
15     Student var;
16     var.firstname = "Max";
17     var.lastname = "Mustermann";
18     var.studentID = 0;
19     var.studiesID = 680;
20     var.test = 25.;
21     var.exercise = 35.;
22 }
```

## Remarks on structures

- ▶ Originally, the C standard did not allow to use
  - Structures as input parameter of a function
  - Structures as output parameter of a function
  - Assignment operator (=) for an entire structure
  
- ▶ In the meantime, it is allowed, **however**:
  - Pass structures to functions dynamically (via pointers)
  - Write assignment (= copy) by yourself
  - Assignment (=) creates a so-called *shallow copy*
  
- ▶ **Shallow copy**:
  - Only the first level is copied
  - i.e., values for basic variables
  - i.e., addresses for pointers
  - **Moreover**: Copy has the same dynamic data
  
- ▶ **Deep copy**:
  - All levels of a structure are copied
  - **Plus**: Copy of the dynamic data

# Structures: Memory allocation

- ▶ Create also the functions
  - **newStudent**: Memory allocation and initialization
  - **freeStudent**: Memory deallocation
  - **cloneStudent**: Complete copy of the full structure including the dynamic data e.g., member **firstname** (so-called *deep copy*)
  - **copyStudent**: Copy of the first level excluding the dynamic data (so-called *shallow copy*)

```
1 Student* newStudent() {
2     Student* pointer = malloc(sizeof(Student));
3     assert( pointer != NULL);
4
5     (*pointer).firstname = NULL;
6     (*pointer).lastname = NULL;
7     (*pointer).studentID = 0;
8     (*pointer).studiesID = 0;
9     (*pointer).test = 0.;
10    (*pointer).exercise = 0.;
11
12    return pointer;
13 }
```

# Structures & Arrow operator

- ▶ In the program, `pointer` is of type `Student*`
- ▶ Access to members, e.g., `(*pointer).firstname`
  - Better syntax `pointer->firstname`
- ▶ Structures `never` static, `always` dynamic
  - Use directly `student` for type `Student*`
- ▶ Better implementation of the function `newStudent` below

```
5 // Declaration of structure
6 struct _Student_ {
7     char* firstname; // First name
8     char* lastname;  // Last name
9     int studentID;    // Student ID
10    int studiesID;    // Studies ID
11    double test;      // Points in final test
12    double exercise;  // Points in exercises
13 };
14
15 // Declaration of corresponding data type
16 typedef struct _Student_ Student;
17
18 // allocate and initialize new student
19 Student* newStudent() {
20     Student* student = malloc(sizeof(Student));
21     assert(student != NULL);
22
23     student->firstname = NULL;
24     student->lastname  = NULL;
25     student->studentID = 0;
26     student->studiesID = 0;
27     student->test      = 0.;
28     student->exercise  = 0.;
29
30     return student;
31 }
```

## Structures: Memory deallocation

- ▶ **Deallocation** of a dynamic variable of type Student
- ▶ **Be careful:** Deallocate the dynamically allocated memory before deallocating the pointer to the structure

```
33 // free memory allocation
34 Student* delStudent(Student* student) {
35     assert(student != NULL);
36
37     if (student->firstname != NULL) {
38         free(student->firstname);
39     }
40
41     if (student->lastname != NULL) {
42         free(student->lastname);
43     }
44
45     free(student);
46     return NULL;
47 }
```

## Shallow Copy

- ▶ **Copy** of a dynamic variable of type Student
  - Copy of only the first level of the structure excluding the dynamically allocated memory

```
49 // shallow copy of student
50 Student* copyStudent(Student* student) {
51     Student* copy = newStudent();
52     assert(student != NULL);
53
54     // Watch out! Pointer!
55     copy->firstname = student->firstname;
56     copy->lastname = student->lastname;
57
58     // Copy of the simple data
59     copy->studentID = student->studentID;
60     copy->studiesID = student->studiesID;
61     copy->test = student->test;
62     copy->exercise = student->exercise;
63
64     return copy;
65 }
```

# Deep Copy

- ▶ **Copy** of a dynamic variable of type Student
  - Copy of all levels of the structure including the dynamically allocated memory
- ▶ Be careful: Copy also the member with dynamically allocated memory

```
67 // deep copy of student
68 Student* cloneStudent(Student* student) {
69     Student* copy = newStudent();
70     int length = 0;
71     assert( student != NULL);
72
73     if (student->firstname != NULL) {
74         length = strlen(student->firstname)+1;
75         copy->firstname = malloc(length*sizeof(char));
76         assert(copy->firstname != NULL);
77         strcpy(copy->firstname, student->firstname);
78     }
79
80
81     if (student->lastname != NULL) {
82         length = strlen(student->lastname)+1;
83         copy->lastname = malloc(length*sizeof(char));
84         assert(copy->lastname != NULL);
85         strcpy(copy->lastname, student->lastname);
86     }
87
88     copy->studentID = student->studentID;
89     copy->studiesID = student->studiesID;
90     copy->test = student->test;
91     copy->exercise = student->exercise;
92
93     return copy;
94 }
```



## Arrays of structures

- ▶ Aim: Generate array with SciProgMath students
- ▶ No static arrays, dynamic arrays
  - Student data are of type **Student**
  - Hence, they are managed with type **Student\***
  - Hence, an array of type **Student\*\*** is needed

```
1 // Declare array
2 Student** participant=malloc(N*sizeof(Student*));
3
4 // Allocate memory for participants
5 for (j=0; j<N; ++j){
6     participant[j] = newStudent();
7 }
```

- ▶ Access to members as before
  - **participant[j]** has type **Student\***
  - Hence, e.g., **participant[j]->firstname**

## Nesting of structures

```
1 struct _Address_ {
2     char* street;
3     char* number;
4     char* city;
5     char* zip;
6 };
7 typedef struct _Address_ Address;
8
9 struct _Employee_ {
10     char* firstname;
11     char* lastname;
12     char* title;
13     Address* home;
14     Address* office;
15 };
16 typedef struct _Employee_ Employee;
```

- ▶ Organize data of employees
  - Name, private address, work address
- ▶ For `employee` of type `Employee*`
  - `employee->home` is a pointer to `Address`
  - Hence, e.g., `employee->home->city`
- ▶ Be careful with allocating, deallocating, copying

# Structures & mathematics

- ▶ Structures for mathematical objects:
  - Vectors in  $\mathbb{R}^n$

# Structures and vectors

```
1 #ifndef _STRUCT_VECTOR_
2 #define _STRUCT_VECTOR_
3
4 #include <stdio.h>
5 #include <stdlib.h>
6 #include <assert.h>
7 #include <math.h>
8
9 // declaration of new data type Vector
10 typedef struct _Vector_ {
11     int n;           // Dimension
12     double* entry;   // Vector coefficients
13 } Vector;
14
15 // Allocate and initialize new vector of length n
16 Vector* newVector(int n);
17
18 // free storage of allocated vector and return NULL
19 Vector* delVector(Vector* X);
20
21 // return length of a vector
22 int getVectorN(Vector* X);
23
24 // return coefficient  $x_i$  of vector  $X$ 
25 double getVectorEntry(Vector* X, int i);
26
27 // assign new value to coefficient  $x_i$  of vector  $X$ 
28 void setVectorEntry(Vector* X, int i, double  $x_i$ );
29
30 // some example functions...
31 Vector* inputVector();
32 double normVector(Vector* X);
33 double productVector(Vector* X, Vector* Y);
34
35 #endif
```

- ▶ Data type to store  $x \in \mathbb{R}^n$ 
  - Dimension  $n$  of type `int`
  - Array coefficients  $x_j$  to store double

## Allocate a vector

- ▶ The vector length  $n \in \mathbb{N}$  is passed to the function
- ▶ Structure allocation, assignment of dimension  $n$
- ▶ Allocation and initialization of the vector array

```
3 Vector* newVector(int n) {
4     int i = 0;
5     Vector* X = NULL;
6
7     assert(n > 0);
8
9     X = malloc(sizeof(Vector));
10    assert(X != NULL);
11
12    X->n = n;
13    X->entry = malloc(n*sizeof(double));
14    assert(X->entry != NULL);
15
16    for (i=0; i<n; ++i) {
17        X->entry[i] = 0;
18    }
19    return X;
20 }
```

## Deallocate a vector

- ▶ Deallocate array
- ▶ Deallocate structure
- ▶ Return value **NULL**

```
22 Vector* delVector(Vector* X) {
23     assert(X != NULL);
24     free(X->entry);
25     free(X);
26
27     return NULL;
28 }
```

## Access to structures

- ▶ Good programming style (unfortunately not very common): No direct access to structure members
- ▶ Better approach:
  - Write **set** and **get** functions for each member
  - So-called *mutator functions*

```
30 int getVectorN(Vector* X) {
31     assert(X != NULL);
32     return X->n;
33 }
34
35 double getVectorEntry(Vector* X, int i) {
36     assert(X != NULL);
37     assert((i >= 0) && (i < X->n));
38     return X->entry[i];
39 }
40
41 void setVectorEntry(Vector* X, int i, double Xi){
42     assert(X != NULL);
43     assert((i >= 0) && (i < X->n));
44     X->entry[i] = Xi;
45 }
```

- ▶ Writing data is not allowed without **set**!
- ▶ Reading data is not allowed without **get**!
- ▶ This approach is more compatible with later changes of the data structure
- ▶ This approach helps to avoid data inconsistencies (and very often also runtime errors)

## Example: Read vector

```
47 Vector* inputVector() {
48
49     Vector* X = NULL;
50     int i = 0;
51     int n = 0;
52     double input = 0;
53
54     printf("Dimension des Vektors n=");
55     scanf("%d",&n);
56     assert(n > 0);
57
58     X = newVector(n);
59     assert(X != NULL);
60
61     for (i=0; i<n; ++i) {
62         input = 0;
63         printf("x[%d]=",i);
64         scanf("%lf",&input);
65         setVectorEntry(X,i,input);
66     }
67
68     return X;
69 }
```

- Read  $n \in \mathbb{N}$  and a vector  $x \in \mathbb{R}^n$  from the keyboard

## Example: Euclidean norm

```
71 double normVector(Vector* X) {
72
73     double Xi = 0;
74     double norm = 0;
75     int n = 0;
76     int i = 0;
77
78     assert(X != NULL);
79
80     n = getVectorN(X);
81
82     for (i=0; i<n; ++i) {
83         Xi = getVectorEntry(X,i);
84         norm = norm + Xi*Xi;
85     }
86     norm = sqrt(norm);
87
88     return norm;
89 }
```

► Compute  $\|x\| := \left( \sum_{j=1}^n x_j^2 \right)^{1/2}$  for  $x \in \mathbb{R}^n$



## Example: Scalar product

```
91 double productVector(Vector* X, Vector* Y) {
92
93     double Xi = 0;
94     double Yi = 0;
95     double product = 0;
96     int n = 0;
97     int i = 0;
98
99     assert(X != NULL);
100    assert(Y != NULL);
101
102    n = getVectorN(X);
103    assert(n == getVectorN(Y));
104
105    for (i=0; i<n; ++i) {
106        Xi = getVectorEntry(X,i);
107        Yi = getVectorEntry(Y,i);
108        product = product + Xi*Yi;
109    }
110
111    return product;
112 }
```

► Compute  $x \cdot y := \sum_{j=1}^n x_j y_j$  for  $x, y \in \mathbb{R}^n$