

## Lecture - 4

Numerical Representation in Computers.

Int  $\rightarrow$  Int8, Int16, Int64, Int32, UInt8,

Float  $\rightarrow$  Float16, Float32, Float64, BigFloat

Int64  $\rightarrow$   $-2^{63} \rightarrow 2^{63}-1$

# Float 64

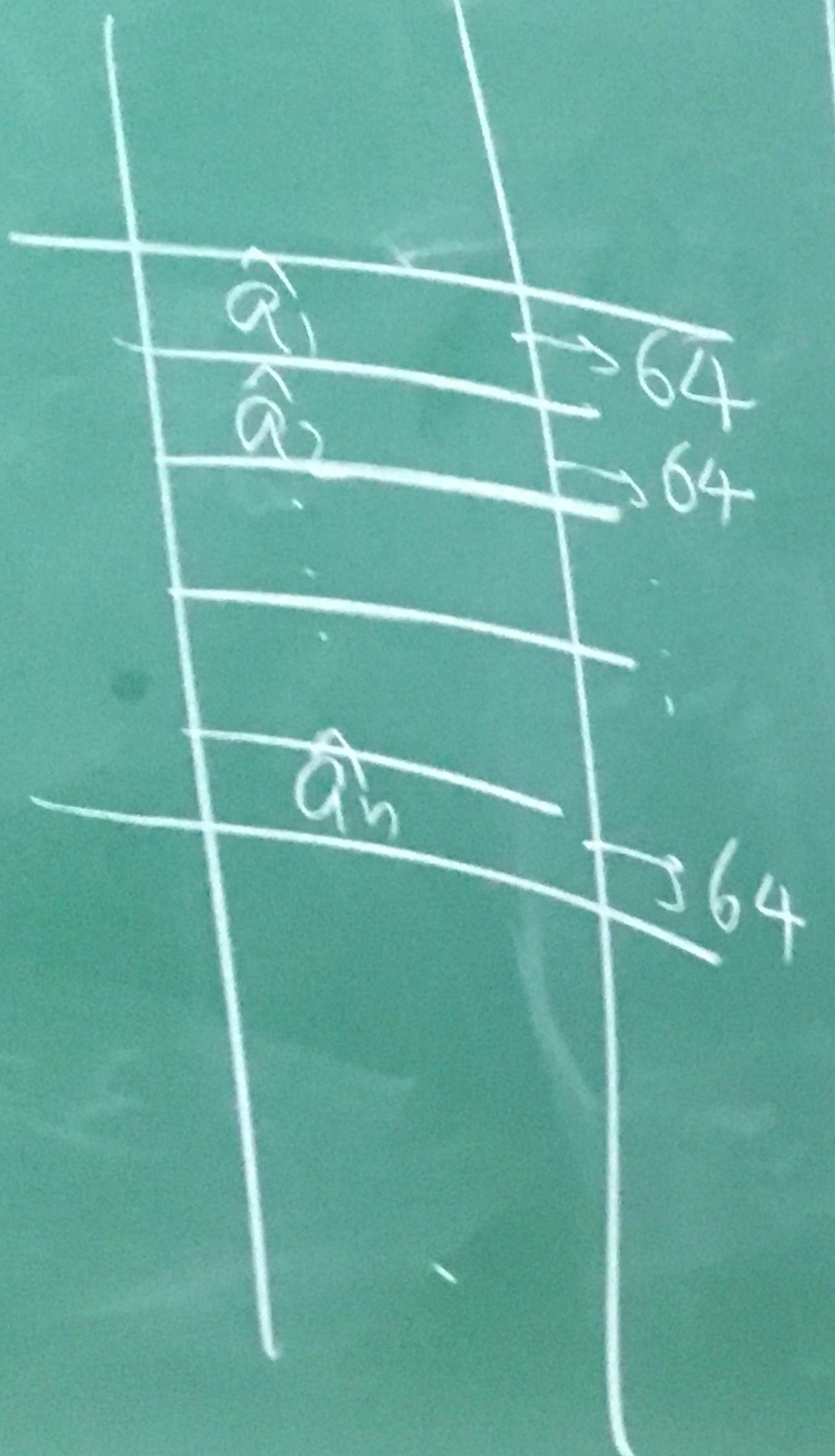
IEEE 754

| + | | + 52  
↓      ↓      ↓  
Sign      Exponent      mantissa.  
(s)      (e)       $b_0 b_1 \dots b_{51}$   
↓  
0 → +ve      an integer       $1.b_0 b_1 \dots b_{51}$   
1 → -ve      -1023 to 1023/  
1022

Flow

representation

(version of a).



$$1.b_0 b_1 \dots b_{51}$$

$$m = 1 + \sum_{k=0}^{51} b_k 2^{-k}$$

$$S \times 2^e \times m$$

Sign  
(S)

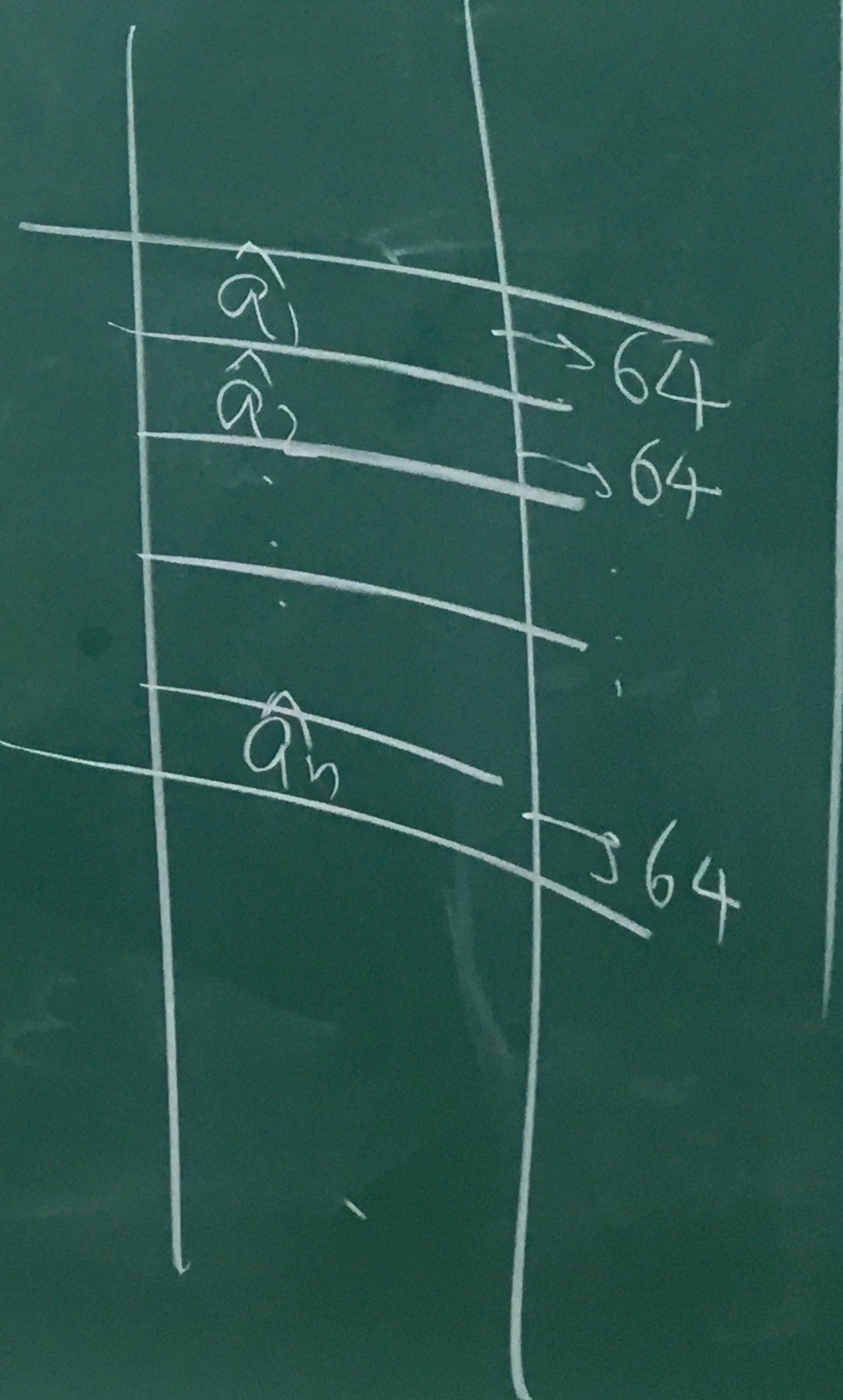
0 → +ve  
1 → -ve

l. b.

Real  
 $a \rightarrow \hat{a}$  Floating Pt  
representation

(quantized version of  $a$ )

Vector  $\underline{a} \in \mathbb{R}^n$   
8n bytes  
 $\left[ \begin{array}{c} a_1 \\ \vdots \\ a_n \end{array} \right] \xleftarrow{\quad} \left[ \begin{array}{c} \hat{a}_1 \\ \vdots \\ \hat{a}_n \end{array} \right] \xrightarrow{\quad}$  f4 n bits



$S \times 2^e \times$

$m = 1 -$

Representing a real matrix ( $m \times n$ ) in Float64

We use 64mn bits or 8mn bytes of memory.

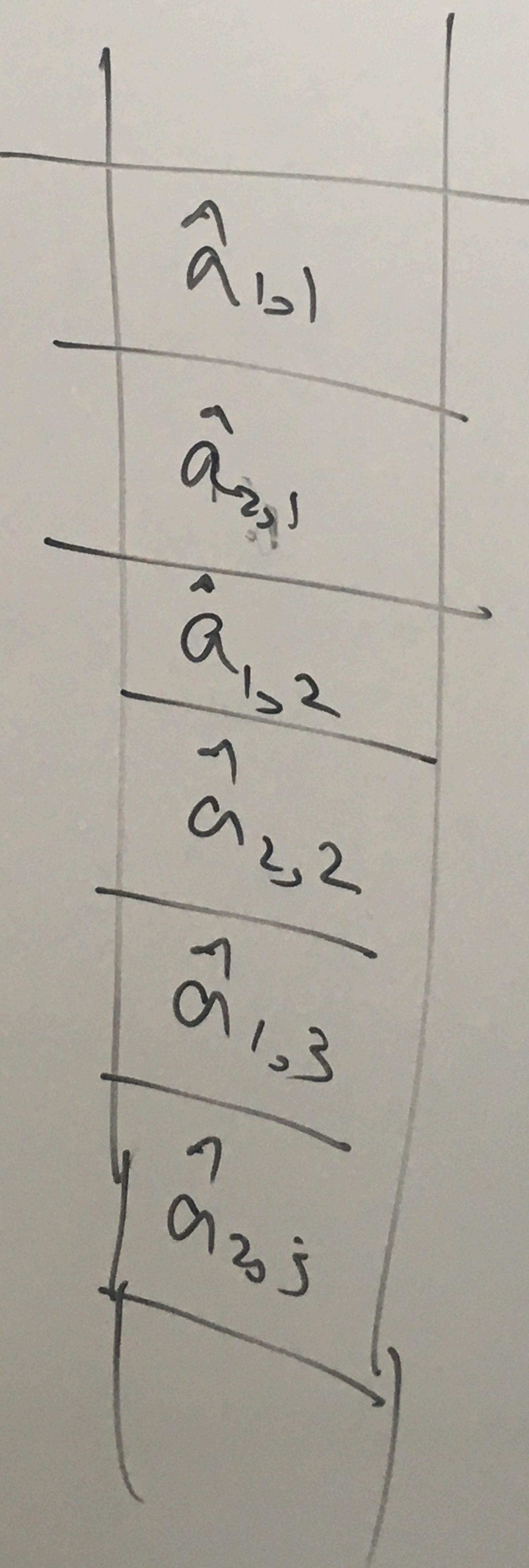
$A =$

$a_{1,1}$	$a_{1,2}$	$a_{1,3}$
$a_{2,1}$	$a_{2,2}$	$a_{2,3}$

$2 \times 3$

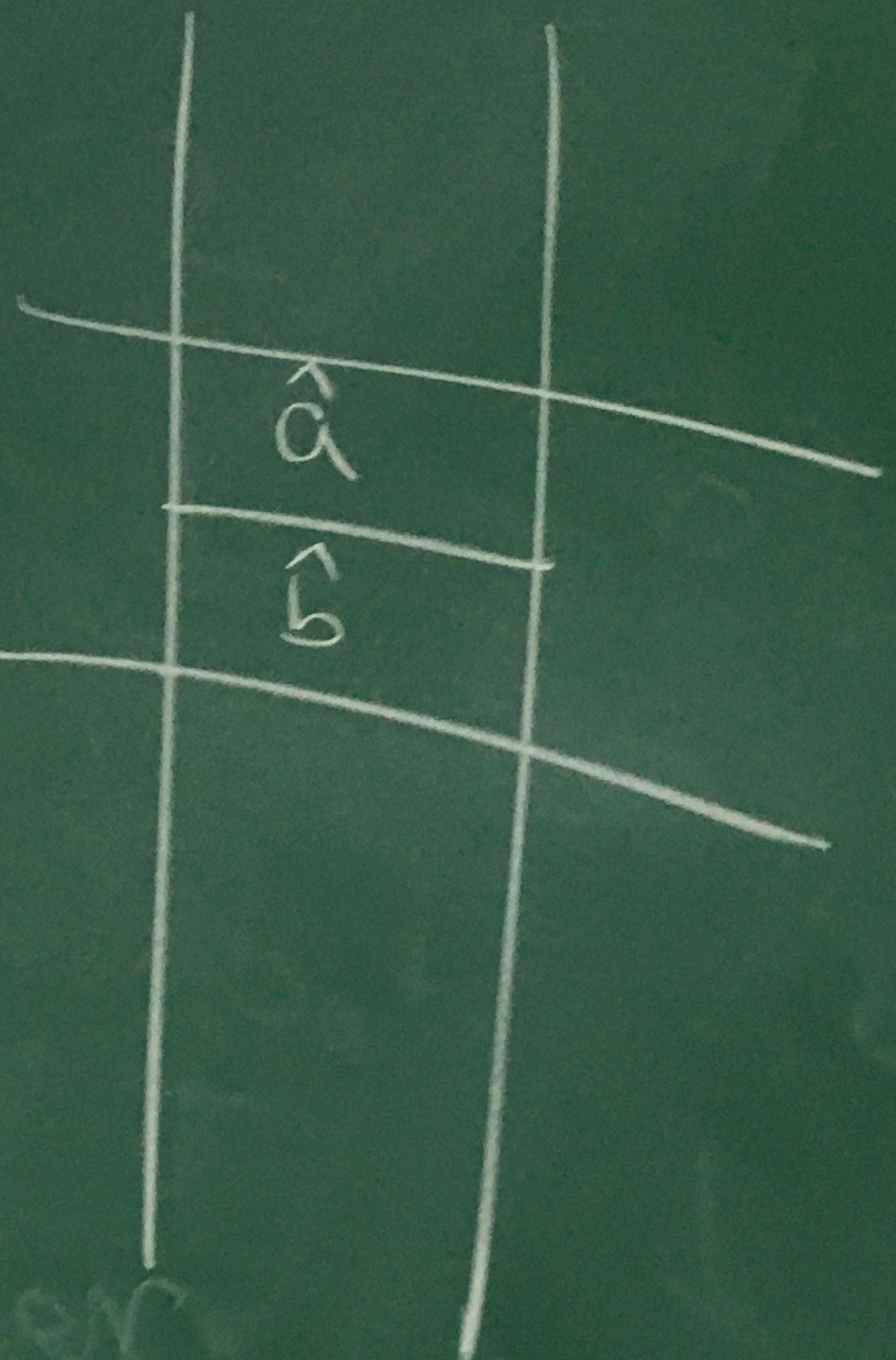
$A[1,2]$  or

$A[3]$



Representation of complex  
numbers using float64.

$$C = a + ib$$



A

A[ ]

A[3]

format "

]

Index

Float64

Value.

$\sqrt{2}$

11

## Sparse Vectors & matrices

A vector or a matrix is said

to be sparse if the fraction

of non-zero entries is small

$nnz(A)$  = number of  
non-zero elements  
of  $A$ .

$$A_{m \times n} \rightarrow \frac{nnz(A)}{mn}$$

Represen-

numbers

$C = Q$

SP

COO format

"Coordinate list format"

Ex.  $A = \begin{bmatrix} 0 & 0 & \sqrt{2} \\ -1 & 0 & 0 \end{bmatrix}$

UInt Row-Index	UInt Column Index	Float64 Value
1	3	$\sqrt{2}$
2	1	-1

## Floating point operation (flop)

Example:

$$\alpha \in \mathbb{R}, A \in \mathbb{R}^{m \times n}$$

$$B = \alpha A.$$

$$= \begin{bmatrix} \alpha a_{1,1} & \dots & \dots \\ \vdots & \ddots & \vdots \\ \vdots & \vdots & \alpha a_{m,n} \end{bmatrix}$$

→ Complexity is  
mn.

②  $\underline{x}, \underline{y} \in \mathbb{R}^n$

$$\underline{z} = \underline{x} + \underline{y}$$

→ n flops.

③  $a, b \in \mathbb{C}$  } → 6 flops.

$$c = ab$$

$$= (Re(a)Re(b) - Im(a)Im(b)) \\ + i(Re(a)Im(b) + Im(a)Re(b))$$

④ Adding two  
sparse vectors

$$\underline{x}, \underline{y} \in \mathbb{R}^n$$

$$\underline{z} = \underline{x} + \underline{y} \quad \left\{ \begin{array}{l} \leq \min \{ \text{nnz}(\underline{x}), \\ \text{nnz}(\underline{y}) \} \\ \text{flops.} \end{array} \right.$$

Number additions =  
number of co-ordinates  
such that  $x_i \neq 0$  &  $y_i \neq 0$

②

$\underline{z}$  =

③

$a, b$

$c = c$

$= (1,$

$tic)$

(5) Say  $\underline{\lambda}, \underline{y} \in \mathbb{R}^n$

and  $\underline{y}$  is sparse.

( $\underline{\lambda}$  is "dense"  $\underline{y}$  is sparse)

$$\underline{s} = \underline{\lambda} + \underline{y}$$

Complexity:  $n n \underline{s}(\underline{y})$  flops

(4)

$$\underline{s} =$$

Number

num

such

## Multiplication of Matrices & Vectors.

(1)

$$\underline{a}, \underline{b} \in \mathbb{R}^n$$

$$c \in \mathbb{R}, \quad c = \underline{a}^T \underline{b} = [a_1 \ a_2 \ \dots \ a_n] \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_n \end{bmatrix} = \underline{b}^T \underline{a}$$
$$= a_1 b_1 + \dots + a_n b_n$$
$$= \sum_{i=1}^n a_i b_i$$

Complexity =  $2n-1$  flops  $\approx 2n$  flops

$a_1, b_1$

$a_m b_n$

$b_{nq}$

Complexity:

$m n$   
Ans

② Outer Product

$a \in \mathbb{R}^m, b \in \mathbb{R}^n$

$$c = a b^T$$

$\downarrow m \times n$        $\downarrow m \times 1$        $\downarrow 1 \times n$

$$= \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_m \end{bmatrix} \begin{bmatrix} b_1 & b_2 & \cdots & b_n \end{bmatrix}$$

$\mathbb{R}^n$

arse.

$\underline{y}$  is sparse)

$\text{nnz}(\underline{y})$  flops

$$C = \begin{bmatrix} a_1 b_1 & a_1 b_2 & \dots & a_1 b_n \\ a_m b_1 & a_m b_2 & \dots & a_m b_n \end{bmatrix}$$

$$= \begin{bmatrix} b_1 a^T & b_2 a^T & \dots & b_n a^T \\ a_1 b_1^T \\ a_2 b_2^T \\ \vdots \\ a_m b_n^T \end{bmatrix}$$

Complexity =

$mn$   
flops.

②

$a \in \mathbb{R}^n$

$$C = \begin{bmatrix} \downarrow & & & \\ m \times n & & & \end{bmatrix}$$

$m \times$

=

$$\begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_m \end{bmatrix}$$