1) The three components are

   1) sign bit S
   2) Exponent (E)
   3) fractional part (M)

Number is $f = (-1)^S M \times 2^E$

Two kinds of encoding done
  i) sign single precission ( 32 bits, E (8 bits) M (23 bits)
  ii) double precission (64 bits, E 11 bits, M 52 bits)
                 1 sign bit for each type

⇒ The precission depends on fractional part.
⇒ The number of significant digits depends on the
number of bits in M
  when M = 24 bits
  we have x significant digits
  where $2^{24} = 10^x$ ⇒ x = 7.2 ⋯
          ⇒ x = 7 ⇒ upto seven significant
decimal places we can represent in binary floati
ng point representation.
⇒ The range of number depends on the number of
bits in exponent that is $2^{8} = 2^8 - 1 = 1$ to 254
0) The exponent is encoded as a biased value
  E = exp + bias where bias = 127 = $(2^{8-1} - 1)$
                        for sigle precission
              bias = 1023 $(2^{11-1} - 1)$
                    for double precission

⇒ Example

As the precission depends on the fractional part (mantissa)

for single precission

23 bits are available so we can precisely represent upto 23 bits after decimal

for example when the decimal number is

~~abs~~ 2.3

it's representtion is.

10. 0100110011001100110011⬤

for double precission

52 bits are available so we can precisely represent upto 52 bits after decimal

for example when the decimal number is

2.3
it's representation is

10.01001100110011001 1001 100) 1001 1601 1001 1601 1001 10011

⇒ So in ~~by~~ normal word saying when the number's lies in between 0 to $\frac{1}{2^{23}}$
it can be accurately representable in single precission type of representation.

⇒ when the number's, decimal part lies between 0 to $\frac{1}{2^{52}}$ it can be represented acurately using double precission representaion.

⇒ when number is

$1.x$

if $0 < x < \frac{1}{2^{23}}$

and it's represented in single precission, it has 100% acuracy

⇒ when number is

$1.x$

if $0 < x < \frac{1}{2^{52}}$

and it's represented in double precission, it has 100% accuracy.