In Floating point representation, we have three components

1.The Sign Bit

2.Exponent re

3.Fractional Part

Precession is one the prime attribute of any Floating-Point Representation,

1.Does any of the above three components play a role in the defining the Precession of the number? If so which are the component or Components which play the role in defining precession and how? Explain this with example in your own words.

**Ans**) Precision means the smallest change that can be represented in Floating point

representation. As per the IEEE format, a floating number ca be represented as sign bit, exponent

and fractional part. The fractional part plays an important role in defining the precision of the

number.

 **For example**: A number 3 can be represented in 4 bits as

$$0.300*10^1$$

$$0.003*10^3$$

$$3.000*10^0$$

From the above representation, the last representation is more precise one as it has three zeros' in

the RHS of numeral 3 which tells that if any extra error in actual result like 3.0006 *10^0 will

lead to only 0.02%error.

2.What is Normal and Subnormal Values as per IEEE 754 standards explain this with the help of number line

**Ans**) Normal representation does not have so many leadings zero's but in case of subnormal

values representation will have minimum number in its exponent value thereby leading to more

zero's in its mantissa. For example:0.05 decimal value can be represented in binary by 2 ways:

(i) Normal Representation: $1.01*2^{-6}$

(ii) Subnormal Representation: $0.00101*2^{-3}$

3.IEEE 754vv defines standards for rounding floating points numbers to a represent able value. There are five methods defines by IEEE for this – Take time and understand what these five methods and explain it in your words using diagrams, illustrations of your own.

**Ans)** IEEE has classified the rounding rules into two classes namelyrounding to nearest and directed roundings.

We shall first deal with rounding to nearest.

1) <u>Rounding to nearest, ties to even</u>: As the name suggests, the real number is rounded off to nearest even number.

   For example: 5.6 is rounded off as 6.0

2) <u>Rounding to nearest, ties away from zero</u>: In this method, real number is rounded off to the nearest integer number. If a real number falls in the middle of two integers, it is rounded to the nearest value above (for positive numbers) or below (for negative numbers).

   For example: 5.1 is rounded off to 5.0

   5.5 is rounded off to 6.0

   -5.5 is rounded off to -6.0

3) <u>Round towards zero</u>: The real number is truncated to the nearest integer while going towards to zero.

   For example: 4.5 is rounded off to 4.0

   4.9 is rounded off to 4.0

   -4.6 is rounded off to -4.0.

4<u>) Round toward $+\infty$</u>: In this method real number is truncated to the nearest integer while going towards positive infinity

   For example: 4.5 is rounded off to 5.0

   4.9 is rounded off to 5.0

-4.6 is rounded off to -4.0

5) Round toward −∞: In this method real number is truncated to the nearest integer while going towards to zero.

For example: 4.5 is rounded off to 4.0

4.9 is rounded off to 5.0

-4.6 is rounded off to -5.0