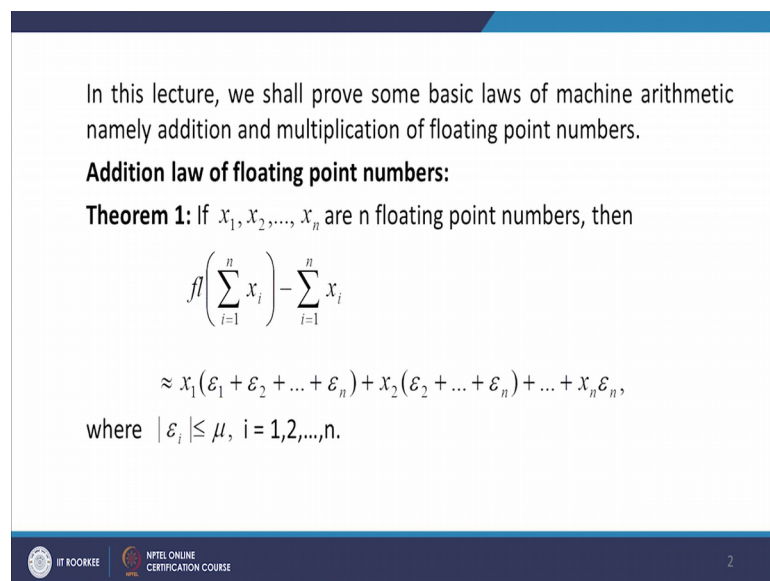


Numerical Linear Algebra
Dr. P. N. Agrawal
Department of Mathematics
Indian Institute of Technology, Roorkee

Lecture - 24
Addition and Multiplication of Floating Point Numbers

Hello friends, welcome to my lecture on Addition and Multiplication of Floating Point Numbers.

(Refer Slide Time: 00:27)



In this lecture, we shall prove some basic laws of machine arithmetic namely addition and multiplication of floating point numbers.

Addition law of floating point numbers:

Theorem 1: If x_1, x_2, \dots, x_n are n floating point numbers, then

$$fl\left(\sum_{i=1}^n x_i\right) - \sum_{i=1}^n x_i$$
$$\approx x_1(\varepsilon_1 + \varepsilon_2 + \dots + \varepsilon_n) + x_2(\varepsilon_2 + \dots + \varepsilon_n) + \dots + x_n \varepsilon_n,$$

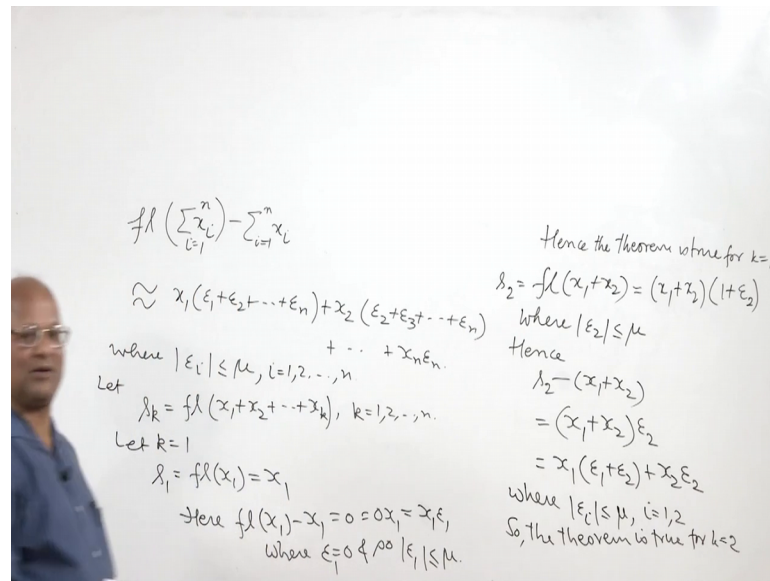
where $|\varepsilon_i| \leq \mu, i = 1, 2, \dots, n.$

ITR ROORKEE NPTEL ONLINE CERTIFICATION COURSE 2

In this lecture we shall prove some basic laws of machine arithmetic namely addition and multiplication of floating point numbers. Suppose x_1, x_2, \dots, x_n are n floating point numbers then we shall show that floating point representation of $\sum_{i=1}^n x_i$ minus $\sum_{i=1}^n x_i$ is approximately equal to $x_1 \varepsilon_1 + x_2 \varepsilon_2 + \dots + x_n \varepsilon_n$ where ε_i are the errors relative errors in the floating point representation of x_i .

So, we know that μ is machine precision. So, $|\varepsilon_i| \leq \mu$ for all values of i from 1 to n . So, what we have to show is this floating point representation of $\sum_{i=1}^n x_i$ minus $\sum_{i=1}^n x_i$ where x_1, x_2, \dots, x_n are n floating point numbers.

(Refer Slide Time: 01:23)



So, then when you find the floating point representation of sigma x i i equal to 1 to n and subtract from it sigma i equal to 1 to n x i the error involved is given as x 1 times epsilon 1 plus epsilon 2 and so on epsilon n, then x 2 times epsilon 2, plus epsilon 3 and so on epsilon n. And then we have similarly x n epsilon n, where mod of epsilon i less than or equal to mu, mu is the machine precision and i is equal to 1, 2 and so on up to n.

(Refer Slide Time: 02:28)

Proof: Let us prove this theorem by induction. Let

$$s_k = fl(x_1 + x_2 + \dots + x_k), \quad k = 1, 2, \dots, n.$$

Then

$$s_1 = fl(x_1) = x_1 \Rightarrow fl(x_1) - x_1 = 0.$$

Taking $\epsilon_1 = 0$, the induction hypothesis holds for $k = 1$.

Now,

$$s_2 = fl(x_1 + x_2) = (x_1 + x_2)(1 + \epsilon_2),$$

where $|\epsilon_2| \leq \mu$.

Hence,

$$s_2 - (x_1 + x_2) = x_1\epsilon_2 + x_2\epsilon_2 = x_1(\epsilon_1 + \epsilon_2) + x_2\epsilon_2$$

IT ROORKEE | NPTEL ONLINE CERTIFICATION COURSE | 3

So, let us assume that let x k denote the floating point representation of x 1 plus x 2 and so on x k, where k takes values from 1 to n. Then what we will do? We shall prove this

theorem by induction on k . So, when we take k equal to 1, let us take k equal to 1 then we see that s_1 is equal to $fl(x_1)$. Now, floating point representation of x_1 is equal to x_1 because x_1 is floating point representation, so x_1, x_2, \dots, x_n are floating point numbers. So, this is $fl(x_1)$ equal to x_1 and therefore, we can say that now here you take x_n equal to 1. So, $fl(x_1) - x_1$ is approximately equal to x_1 into ϵ_1 that we have to prove.

So, what do we get here? So, here this equal to 0, this 0 can be regarded as 0 into x_1 . So, we can take ϵ_1 equal to 0. So, this equal to x_1 into ϵ_1 , where ϵ_1 is equal to 0 and so mod of ϵ_1 is less than or equal to μ . So, the we can say that $fl(x_1) - x_1$ is equal to $fl(x_1) - x_1$ equal to x_1 into ϵ_1 , where mod of ϵ_1 is less than or equal to μ . So, the theorem holds true for n equal to 1, hence for k equal to 1. If you take in s and k , s_k k equal to 1 the theorem is true.

Now, we can also show that the term is true for k equal to 2, fl of $x_1 + x_2$, we can write as $x_1 + x_2$ into $1 + \epsilon_1$. So, then we shall see that $1 + \epsilon_1$. So, then what we will see is that where mod of ϵ_1 it less than or equal to μ hence $s_2 - s_1 + s_2$. This is s_2 , $s_2 - x_1 + x_2$ is equal to $x_1 + x_2$ into ϵ_1 from s_2 when we subtract $x_1 + x_2$ we get $x_1 + x_2$ into ϵ_1 and so what we can says that we can write it as x_1 times because ϵ_1 is equal to 0. So, $x_1 \epsilon_1 + x_2 \epsilon_1$, we can add the term $x_1 \epsilon_1$ and write it as x_1 times $\epsilon_1 + \epsilon_1 + x_2 \epsilon_1$, where mod of ϵ_1 is less than or equal to μ and mod of ϵ_1 is also less than or equal to μ and so the theorem is true for k equal to 2. So, the theorem is true for k equal to 2.

(Refer Slide Time: 07:12)

The induction hypothesis is true for $k=2$.

Let us suppose that the induction hypothesis is true for integer $k = m$.

Then

$$s_m = \sum_{k=1}^m x_k$$

$$\approx x_1(\varepsilon_1 + \varepsilon_2 + \dots + \varepsilon_m) + x_2(\varepsilon_2 + \dots + \varepsilon_m) + \dots + x_m \varepsilon_m, \dots \dots \dots (1)$$

where, $|\varepsilon_i| \leq \mu$, for all $i = 1, 2, \dots, m$.

We shall show that the induction hypothesis also holds for $k = m+1$.

We have $s_{m+1} = fl(x_1 + x_2 + \dots + x_{m+1}) \approx (s_m + x_{m+1})(1 + \varepsilon_{m+1})$,

where, $|\varepsilon_{m+1}| \leq \mu$.

IIT ROORKEE | NPTEL ONLINE CERTIFICATION COURSE | 4

Now, let us suppose that the theorem is true for certain value k equal to m . So, let us suppose that k equal to m and n for all values of k less than or equal to m minus 1.

(Refer Slide Time: 07:21)

Let us suppose that the theorem is true for $k = m$

Then

$$s_m = \sum_{k=1}^m x_k \approx x_1(\varepsilon_1 + \varepsilon_2 + \dots + \varepsilon_m) + x_2(\varepsilon_2 + \varepsilon_3 + \dots + \varepsilon_m) + \dots + x_m \varepsilon_m \quad \text{--- (1)}$$

where $|\varepsilon_i| \leq \mu, i = 1, 2, \dots, m$

$$s_{m+1} = fl(x_1 + x_2 + \dots + x_{m+1}) \approx (s_m + x_{m+1})(1 + \varepsilon_{m+1})$$

where $|\varepsilon_{m+1}| \leq \mu$

Using (1)

$$s_{m+1} \approx \left\{ \begin{aligned} &x_1 + x_2 + \dots + x_m + x_{m+1} \\ &+ x_1(\varepsilon_1 + \varepsilon_2 + \dots + \varepsilon_m) \\ &+ x_2(\varepsilon_2 + \varepsilon_3 + \dots + \varepsilon_m) + \dots + x_m \varepsilon_m \end{aligned} \right\} (1 + \varepsilon_{m+1})$$

$$\approx x_1(\varepsilon_1 + \varepsilon_2 + \dots + \varepsilon_m) + x_2(\varepsilon_2 + \varepsilon_3 + \dots + \varepsilon_m) + \dots + x_m \varepsilon_m + x_{m+1}$$

So, then s_m is minus s_m that is floating point representation of x_1 plus x_2 and so on s_m plus $\sum_{i=1}^m x_i \varepsilon_i$ is approximately $x_1 \varepsilon_1 + x_2 \varepsilon_2 + \dots + x_m \varepsilon_m$ and then x_{m+1} by s_m by our induction hypothesis where $|\varepsilon_i| \leq \mu$ for all i equal to 1 2 and so on up to m .

Now, we shall show that the induction hypothesis also holds for k equal to $m + 1$. So, let us see what is s_{m+1} ? s_{m+1} equal to floating point representation of $x_1 + x_2$ and so on this $x_1 + x_2$ and so on $x_m + x_{m+1}$ which is equal approximately $s_m + x_{m+1}$ into $1 + \epsilon_{m+1}$, now where $\text{mod of } \epsilon_{m+1}$ is less than or equal to m . So, ϵ_{m+1} is the error in the computation of $x_1 + x_2 + \dots + x_m + x_{m+1}$.

Now, so let us use what is given to us. So, using 1, using 1 is s_m minus $\sigma_{k \text{ equal to } 1 \text{ to } m}$ x_k yeah, $i \text{ equal to } 1 \text{ to } m$ x_i is approximately x_1 times $\epsilon_1 + \epsilon_2$ and so on $\epsilon_m + x_2$ times ϵ_2 plus and so on $\epsilon_m + x_m$ ϵ_m . So, let us use this and then we can write that s_{m+1} . So, using 1, this is 1, s_{m+1} is approximately $\sum_{i=1}^m x_i$. So, we have $x_1 + x_2$ and so on $x_m + x_{m+1}$ plus x_1 into $\epsilon_1 + \epsilon_2$ and so on $\epsilon_m + x_2$ times ϵ_2 plus ϵ_3 and so on ϵ_m plus and so on x_m ϵ_m and this multiplied by $1 + \epsilon_{m+1}$.

So, this is $1 + \epsilon_{m+1}$ is multiplied to the whole thing. So, what we will get is. Now, this is or we can say s_{m+1} , $x_1 + x_2$ and so on $x_m + x_{m+1}$ see we are what we are doing we are multiplying by one first when we are multiplying by 1, here we get $x_1 + x_2 + \dots + x_m + x_{m+1}$ and then this whole thing. So, that $x_1 + x_2 + \dots + x_m + x_{m+1}$ I am bringing to the left side this is approximately equal to x_1 times $\epsilon_1 + \epsilon_2$ and so on $\epsilon_m + x_2$ times ϵ_2 ok.

So, when we multiply by 1 we get the whole thing and then we multiply by ϵ_{m+1} . So, what we get is this whole thing multiplied by ϵ_{m+1} . So, $\sum_{i=1}^m x_i$, $i \text{ equal to } 1 \text{ to } m + 1$ multiplied by $x_m \epsilon_{m+1}$ and then we will get x_1 times ϵ_1 into ϵ_{m+1} ϵ_2 into ϵ_{m+1} which are terms of second order. So, we can neglect them because ϵ is are too small and therefore, we have assume that they are second and higher order terms can be neglected. So, what we get here now? $s_{m+1} - \sum_{i=1}^m x_i$ is approximately equal to; now here we shall have $x_1 \epsilon_{m+1}$ that $x_1 \epsilon_{m+1}$ we can observe here and write x_1 times.

(Refer Slide Time: 14:22)

$$S_{m+1} = \sum_{i=1}^{m+1} x_i$$

$$\approx x_1(\epsilon_1 + \epsilon_2 + \dots + \epsilon_{m+1})$$

$$+ x_2(\epsilon_2 + \epsilon_3 + \dots + \epsilon_{m+1})$$

$$+ \dots + x_m(\epsilon_m + \epsilon_{m+1})$$

$$+ x_{m+1}\epsilon_{m+1}$$

where $|\epsilon_i| \leq \mu, \forall i = 1, 2, \dots, m+1$
 So, the induction hypothesis holds true for $k = m+1$.

Using (1)

$$S_{m+1} \approx \left\{ \begin{aligned} &x_1 + x_2 + \dots + x_m + x_{m+1} \\ &+ x_1(\epsilon_1 + \epsilon_2 + \dots + \epsilon_m) \\ &+ x_2(\epsilon_2 + \epsilon_3 + \dots + \epsilon_m) + \dots + x_m \epsilon_m \end{aligned} \right\} (1 + \epsilon_{m+1})$$

or

$$S_{m+1} = (x_1 + x_2 + \dots + x_{m+1}) + \left(\sum_{i=1}^{m+1} x_i \epsilon_i \right)$$

$$\approx x_1(\epsilon_1 + \epsilon_2 + \dots + \epsilon_m) + x_2(\epsilon_2 + \epsilon_3 + \dots + \epsilon_m) + \dots + x_m \epsilon_m + x_{m+1} \epsilon_{m+1}$$

Then x_2 into epsilon $m+1$ that term can be brought here. So, x_2 times epsilon 2 plus epsilon 3 and so on epsilon $m+1$ and then before this term we will have x_{m-1} in that we can multiply we can add there x_{m-1} into epsilon $m+1$ term and here we can add x_m into epsilon $m+1$. So, we shall have and so on and lastly we have this where mod of epsilon i is less than or equal to μ for all i equal to 1 2 and so on up to $m+1$.

So, the induction hypothesis holds true for k equal to $m+1$ and therefore, it holds for all integers m or you can say it holds for all integers positive integers n . So, this is the proof of first theorem.

(Refer Slide Time: 16:26)

Using (1), we get

$$s_{m+1} \approx \left(\sum_{i=1}^{m+1} x_i + x_1(\varepsilon_1 + \varepsilon_2 + \dots + \varepsilon_m) + x_2(\varepsilon_2 + \dots + \varepsilon_m) + \dots + x_m \varepsilon_m \right) (1 + \varepsilon_{m+1})$$

$$\Rightarrow s_{m+1} - \sum_{i=1}^{m+1} x_i$$

$$\approx x_1(\varepsilon_1 + \varepsilon_2 + \dots + \varepsilon_m + \varepsilon_{m+1})$$

$$+ x_2(\varepsilon_2 + \dots + \varepsilon_m + \varepsilon_{m+1}) + \dots + x_m(\varepsilon_m + \varepsilon_{m+1}) + x_{m+1} \varepsilon_{m+1},$$

Thus the induction hypothesis also holds for $k = m+1$. This proves the theorem.

IIT ROORKEE | NITEL ONLINE CERTIFICATION COURSE 5

Now, let us go to the proof of the second theorem which is on multiplication.

(Refer Slide Time: 16:28)

Theorem 2: If x_1, x_2, \dots, x_n are n floating point numbers, then

$$fl(x_1 x_2 \dots x_n) \approx (1 + \delta)(x_1 x_2 \dots x_n),$$

where $\delta = (1 + \varepsilon_2)(1 + \varepsilon_3) \dots (1 + \varepsilon_n) - 1$ and $|\varepsilon_i| \leq \mu, i=1,2,\dots,n$.

Proof: We shall establish this result by induction. Let us define

$$M_i = fl(x_1 x_2 \dots x_i), \quad i=1,2,\dots,n.$$

We observe that

$$M_1 = fl(x_1) = x_1 \text{ so the induction hypothesis holds because}$$

$$\delta = (1 + \varepsilon_1) - 1 \text{ where } \varepsilon_1 = 0. \text{ Further,}$$

$$M_2 = fl(x_1 x_2) = x_1 x_2 (1 + \varepsilon_2),$$

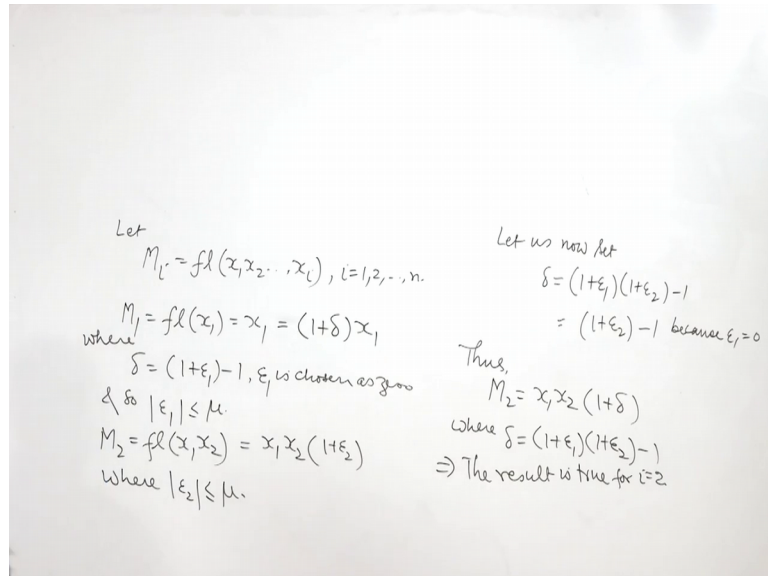
where $|\varepsilon_2| \leq \mu$.

IIT ROORKEE | NITEL ONLINE CERTIFICATION COURSE 6

So, let us assume that we have x_1, x_2, \dots, x_n as n floating point numbers we are given n floating point numbers. Then the floating point representation of x_1 into x_2 into x_n , if we find out this then this approximately equal to $1 + \delta$ times x_1 into x_2 into x_n where δ is the error where δ is the expression $(1 + \varepsilon_2)(1 + \varepsilon_3) \dots (1 + \varepsilon_n) - 1$ and $|\varepsilon_i| \leq \mu, i=1,2,3$ and so on up to n . So, this result also we shall prove by

induction. So, let us assume that m_i is equal to $fl(x_1 \times 2 \times \dots \times i)$, i is equal to 1, 2, 3 and so on up to n .

(Refer Slide Time: 17:18)



So, again let us prove that the result holds true for i equal to 1. So, M_1 is equal to $fl(x_1)$ and floating point representation of x_1 is x_1 because x_1 is a floating point number and so the induction hypothesis holds true because δ we can write as $1 + \epsilon_1 - 1$.

So, here we can write M_1 which is $fl(x_1)$, $fl(x_1)$ is equal to $1 + \delta$ into, so this δ where we where we take ϵ_1 the ϵ_1 to be equal to 0. So, what we do is we can write this is equal to $1 + \delta$ into x_1 where δ is given by $1 + \epsilon_1 - 1$ ϵ_1 is taken as 0; ϵ_1 is chosen as 0 and so mod of ϵ_1 is less than or equal to μ . So, the result holds true for i equal to 1, now let us show that the result holds true for i equal to 2. So, m_2 is equal to floating point representation of x_1 into x_2 .

Now, let us say we write it as $x_1 \times 2$ into where ϵ_2 is the error relative error in the computation of x_1 into x_2 floating point representation of x_1 into x_2 and mod of ϵ_2 , here mod of ϵ_2 is less than or equal to μ .

(Refer Slide Time: 20:00)

Let us now set $\delta = (1 + \varepsilon_1)(1 + \varepsilon_2) - 1$, then

$M_2 = x_1 x_2 (1 + \delta)$, and so the induction hypothesis also holds for $k=2$. Let us now assume that the induction hypothesis holds true for $k=p$ i.e.



$$M_p = f(x_1 x_2 \dots x_p) \approx (1 + \delta)(x_1 x_2 \dots x_p),$$

where $\delta = (1 + \varepsilon_1)(1 + \varepsilon_2) \dots (1 + \varepsilon_p) - 1, |\varepsilon_i| \leq \mu, i = 1, 2, \dots, p$.

Then

$$M_{p+1} = f(x_1 x_2 \dots x_p x_{p+1})$$
$$\approx f(x_1 x_2 \dots x_p) x_{p+1} (1 + \varepsilon_{p+1}),$$

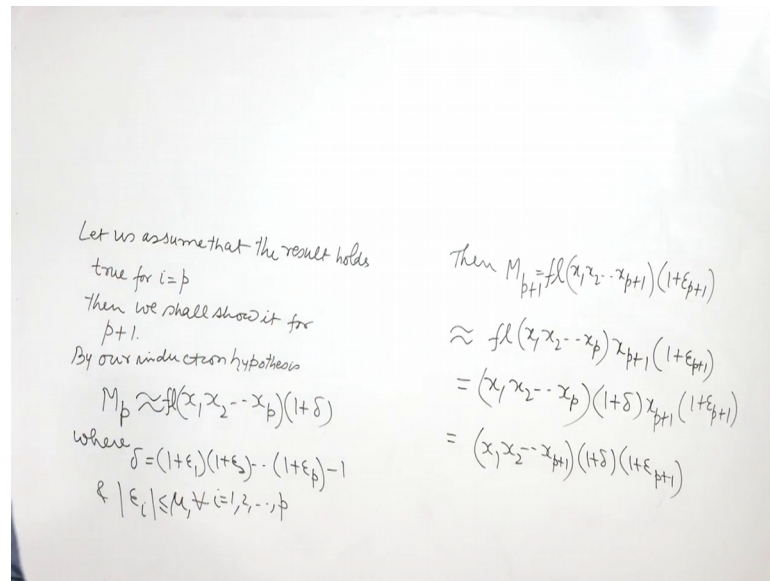
where, $|\varepsilon_{p+1}| \leq \mu$.

 IIT ROORKEE  NFTEL ONLINE CERTIFICATION COURSE 7

Now, let us set delta equal to remember epsilon 1 is equal to 0. So, this is actually 1 plus epsilon 2 minus 1 because epsilon 1 we have chosen as 0. So, delta is equal to 1 plus epsilon 2 minus 1 and this is therefore, equal to delta is equal to epsilon 2. So, you can also write the M_2 is equal to x_1, x_2 into 1 plus delta where delta is 1 plus epsilon 1, 1 plus epsilon 2 minus 1.

Now, let us assume that the result holds true for k equal to p here. So, thus, so the result is true for i equal to 2. Let us assume that the result holds true for i equal to p . So, let us assume that the result holds true for i equal to p then we shall show it for i plus 1 or so, p plus 1 and then we shall be able to say that it holds true for all positive integers p or it holds true for all positive integers n .

(Refer Slide Time: 21:37)



Now, so, by our hypothesis M_p , M_p , i equal to p . So, M_p minus x_1 into x_2 and so on x_p this is this is approximately equal to $1 + \delta$, M_p is equal to, floating point M_p is x_1 plus x_2 and so on x_p into $1 + \delta$ where let me write δ cap this is then where δ is equal to $1 + \epsilon_1, 1 + \epsilon_2$ and so on $1 + \epsilon_{p-1}$ and mod of ϵ_i is less than or equal to μ for all i we have let us write δ only we need to write δ here.

So, δ is equal to this mod of ϵ_i is less than or equal to μ , for all i equal to $1, 2, 3$ and so on up to p . So, then I am take $p+1$ and $p+1$ which is x_1, x_2, x_{p+1} into $1 + \epsilon_{p+1}$ where ϵ_{p+1} is the error. So, this gives us $fl(x_1, x_2, x_p, x_{p+1})$ I have not written fl here, and here also we need to write fl . So, $fl(x_1, x_2, x_p)$ into $1 + \delta$ where this is equal to this, so this into x_{p+1} into $1 + \epsilon_{p+1}$.

Now, what? So, this is equal to let us apply the induction.

(Refer Slide Time: 25:34)

Hence



$$M_{p+1} \approx (x_1 x_2 \dots x_p x_{p+1})(1 + \hat{\delta})$$

where

$$\hat{\delta} = (1 + \delta)(1 + \varepsilon_{p+1}) - 1$$

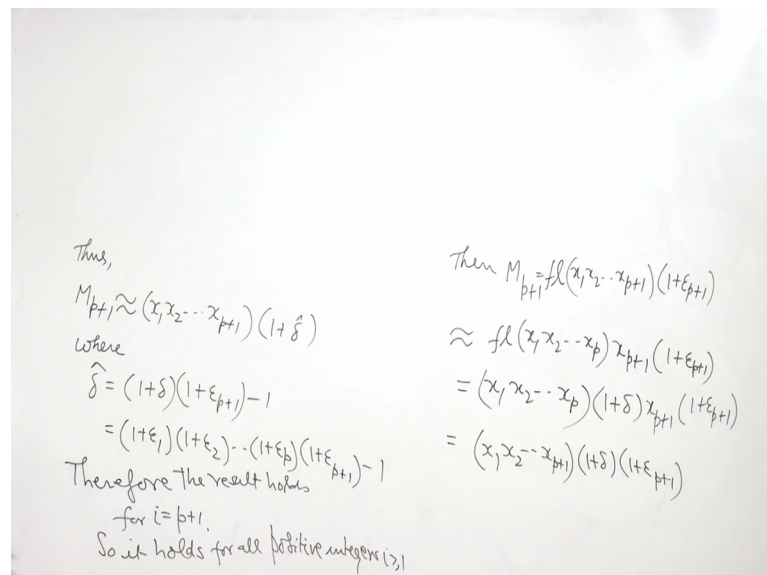
$$= (1 + \varepsilon_1)(1 + \varepsilon_2) \dots (1 + \varepsilon_p)(1 + \varepsilon_{p+1}) - 1.$$

Thus, the induction hypothesis also holds for $i = p+1$. This proves the theorem.

 IIT ROORKEE  NFTEL ONLINE CERTIFICATION COURSE 8

This is, so x_1, x_2, \dots, x_p this is equal to x_1, x_2, \dots, x_p into $1 + \delta$ into x_{p+1} plus 1 into $1 + \varepsilon_{p+1}$ or we can write it as x_1, x_2, \dots, x_p plus 1 into this. Thus we can say that M_{p+1} is approximately equal to x_1, x_2, \dots, x_p plus 1 into $1 + \delta$ cap, where δ cap is equal to $1 + \delta$ into $1 + \varepsilon_{p+1}$ minus 1 .

(Refer Slide Time: 26:21)



Thus,

$$M_{p+1} \approx (x_1 x_2 \dots x_{p+1})(1 + \hat{\delta})$$

where

$$\hat{\delta} = (1 + \delta)(1 + \varepsilon_{p+1}) - 1$$

$$= (1 + \varepsilon_1)(1 + \varepsilon_2) \dots (1 + \varepsilon_p)(1 + \varepsilon_{p+1}) - 1$$

Therefore the result holds for $i = p+1$.
So it holds for all positive integers i .

Then $M_{p+1} = (x_1 x_2 \dots x_{p+1})(1 + \varepsilon_{p+1})$

$$\approx (x_1 x_2 \dots x_p) x_{p+1} (1 + \varepsilon_{p+1})$$

$$= (x_1 x_2 \dots x_p)(1 + \delta) x_{p+1} (1 + \varepsilon_{p+1})$$

$$= (x_1 x_2 \dots x_{p+1})(1 + \delta)(1 + \varepsilon_{p+1})$$

1 plus delta we have seen, 1 plus delta is equal to 1 plus epsilon 1. So, it is 1 plus epsilon 1 into 1 plus epsilon 2 and so on 1 plus epsilon p into 1 plus epsilon p plus 1 minus 1. Now, this is approximately equal to we multiply by 1 here. And so what we will get? 1

plus epsilon 1, 1 plus epsilon 2, 1 plus epsilon p into 1 plus epsilon p plus 1 minus 1.
With that we have proved the result.

So, we can write M^{p+1} as $x_1 \cdot x_2 \cdot \dots \cdot x_{p+1}$ into $1 + \delta$ where δ cap is equal to $1 + \delta$ into $1 + \epsilon_p + 1 - 1$, but $1 + \delta$ is equal to $1 + \epsilon_1$ into $1 + \epsilon_2$ and so on $1 + \epsilon_p$ into $1 + \epsilon_p - 1$ and therefore, the result holds for $p + 1$. So, it holds for all positive integers i and therefore, it holds for all integers and equal to 1, 2, 3 and so on. So, this proves the theorem on multiplication of floating point numbers. With that I would like to conclude my lecture.

Thank you very much for your attention.