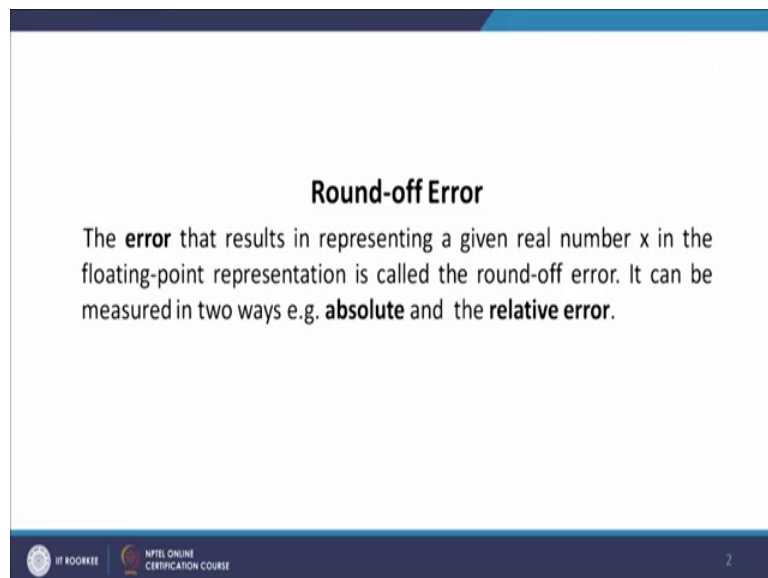


Numerical Linear Algebra
Dr. P. N. Agrawal
Department of Mathematics
Indian Institute of Technology, Roorkee

Lecture - 22
Round-off error

Hello friends, welcome to my lecture on round off errors. Now the error that results in representing a given real number x in the floating-point representation is called round off error.

(Refer Slide Time: 00:28)



Round-off Error

The **error** that results in representing a given real number x in the floating-point representation is called the round-off error. It can be measured in two ways e.g. **absolute** and the **relative error**.

IT ROORKEE | NPTEL ONLINE CERTIFICATION COURSE | 2

We have seen that when you represent floating point representation you in a machine, there are 2 types of base to do that. One is that chopping and the other one is rounding. So, both of them caused error which is called as the round off errors ok.



So, it now this error can be measured in 2 ways one is that absolute error, and the other one is relative error.

(Refer Slide Time: 00:59)

Absolute Error and Relative Error

Let x^* be an approximate value of real number x then the **absolute error** $E_x = |x - x^*|$
and the **relative error** is $R_x = \left| \frac{x - x^*}{x} \right|, (x \neq 0)$.

Example: Let $x = 3.1416$, $x^* = 3.1418$ and $y = 1.5244$, $y^* = 1.5242$,
then $|x - x^*| = |y - y^*| = 2.0 \times 10^{-4}$.
So, the absolute errors are the same but $R_x = \left| \frac{x - x^*}{x} \right| = 6.3662 \times 10^{-5}$
and $R_y = \left| \frac{y - y^*}{y} \right| = 1.3120 \times 10^{-4}$.

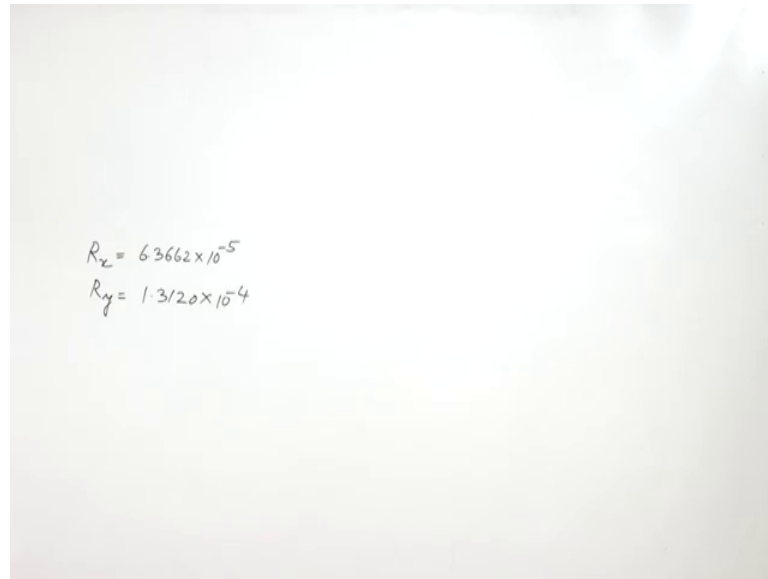
 IIT ROORKEE  NPTEL ONLINE
CERTIFICATION COURSE

In the case let us now define what do we mean by absolute error. So, let x^* be an approximate value of a real number x , then the absolute error E_x is equal to mod of x minus x^* . And the relative error is defined as $R_x = \frac{|x - x^*|}{|x|}$ assuming that of course, x is not equal to 0.

For example, let us consider, x is equal to 3.1416 and x^* is equal to 3.1418, y equal to 1.5244, and y^* equal to 1.5242. Then mod of x minus x^* and mod of y minus y^* both are same, and they are equal to 2.0×10^{-4} . So, the absolute error in both the cases is same it is 2.0×10^{-4} .

Now, but if you calculate the relative error R_x , in the case of x , then $\frac{|x - x^*|}{|x|}$ comes out to be 6.3662×10^{-5} . While the relative error in y R_y comes out to be $\frac{|y - y^*|}{|y|}$ is equal to 1.3120×10^{-4} .

(Refer Slide Time: 02:23)



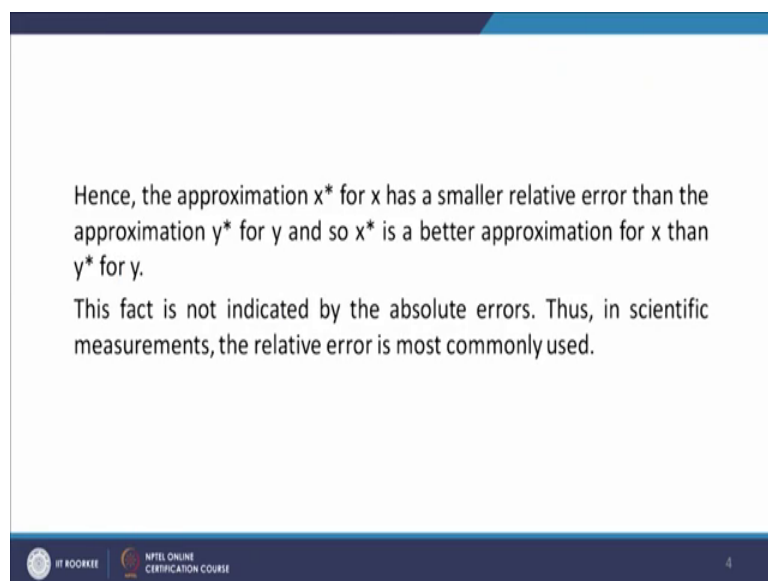
Handwritten mathematical expressions for relative errors R_x and R_y :

$$R_x = 6.3662 \times 10^{-5}$$
$$R_y = 1.3120 \times 10^{-4}$$

So, R_x is equal to 6.3662×10^{-5} . And R_y is equal to 1.3120×10^{-4} .

Now, you can see that this is 6.3662×10^{-5} , and R_y is 1.3120×10^{-4} . So, R_x is smaller than R_y ok. And therefore, we can say that the approximation x^* for x has a smaller relative error, then the approximation y^* for y .

(Refer Slide Time: 03:16)



Hence, the approximation x^* for x has a smaller relative error than the approximation y^* for y and so x^* is a better approximation for x than y^* for y .

This fact is not indicated by the absolute errors. Thus, in scientific measurements, the relative error is most commonly used.

IT ROORKEE | NPTEL ONLINE CERTIFICATION COURSE

And so, we can say x^* gives a better approximation to x the approximation y^* for y .

Now, this fact is not indicated by the absolute error, you can see because absolute error in both the cases is same mod of x minus x star is equal to mod y minus y star is same in so, the this fact is not indicated by the absolute error, but and therefore, the relative error is most commonly used in scientific measurements.

Now, let us take up the case of system of linear equations, and solve it by Gaussian elimination. We have discussed in our previous lecture, the Gaussian elimination is scheme we consider the equations as $ax_1 + bx_2 = e$, $cx_1 + dx_2 = f$.

(Refer Slide Time: 04:12)

The image shows handwritten mathematical work on a whiteboard. It starts with the general system of equations:
$$ax_1 + bx_2 = e$$

$$cx_1 + dx_2 = f$$
Then it defines $m = \frac{c}{a}$ and derives the formula for x_2 :
$$y = x_2 = \frac{f_1}{d_1} = \frac{f - em}{d - bm}$$
Then it derives the formula for x_1 :
$$x = x_1 = \frac{e - bx_2}{a}$$
Below these, it calculates $y^* = \frac{f_1}{d_1} = 2.54$ and $x^* = 1.36$.
It also shows the calculation of $m = \frac{c}{a} = \frac{1}{2} = 0.5$, $d_1 = d - bm = 3.5$, and $f_1 = 8986 - (0.171) \cdot 0.5 = 89$.
Finally, it lists the exact solution:
$$2x - y = 0.171$$

$$x + 3y = 8986$$

$$a=2, c=1, e=0.171$$

$$b=-1, d=3, f=8986$$

$$\text{Exact solution}$$

$$x = 1.357$$

$$y = 2.543$$

And we had m equal to c by a there. What we did was we multiplied the first equation by c by a assuming a to be not equal to 0 and then subtracting the first equation from the second equation. So, when you multiply the first equation by c by a , the coefficient of x one will become c . So, cx_1 will cancel cx_1 when you subtract after multiplication by c by a from the second equation. And we got x_2 equal to f_1 upon d_1 f_1 by d_1 equal to f minus e m , and d_1 was d minus b m .

So, having found x_2 , then one can find x_1 from the equation 1, by e minus bx_2 divided by a . So, here in this case we have in this example $2x$ minus y equal to 0.171, and we have x plus 3 y equal to 8.986.

Now so, when we can calculate m equal to c y so, a is equal to 2 here, c is equal to 1, b is equal to minus 1, and d equal to 3, e equal to 0.171, and f equal to 8.986 ok. So, m equal

to c by a when you calculate m equal to c by a here, c by a means 1 by 2 . So, m comes out to be 0.5 . And d_1 equal to d minus bm , d_1 equal to d minus bm ok. So, put the value of d d is equal to d is equal to 3 b is equal to $\text{minus } 1$, and m equal to $\text{point } 5$, d_1 comes out to be 3.5 , and we have done f_1 ok.

So, f_1 is equal to f , f is equal to 8.986 minus e is equal to 0.171 into m , m is equal to 0.5 . So, this comes out to be 8.9 , 8.9 here we are taking base 10 and t equal to 3 t equal to 3 with rounding.

(Refer Slide Time: 07:38)

Example: Solve $2x - y = 0.171$, $x + 3y = 8.986$ in a machine with base 10 and $t = 3$ (with rounding). Calculate the absolute and relative errors for the solution.

Solution: Here $a = 2$, $b = -1$, $e = 0.171$, $c = 1$, $d = 3$ and $f = 8.986$.

Step-1: $m = \frac{c}{a} = 0.5$, $d_1 = d - bm = 3.5$ and $f_1 = f - em = 8.9$.

Hence $y^* = \frac{f_1}{d_1} = 2.54$ and $x^* = \frac{e - by}{a} = 1.36$.

IT ROORKEE | NPTEL ONLINE CERTIFICATION COURSE | 5

So, I am writing the values with t equal to 3 , when you calculate d_1 in the digital computer. And use t equal to 3 with rounding, then you will get d_1 is equal to 3.5 and f_1 equal to 8.9 . And so, the value of y^* , y here x_1 x_2 in this example are actually x and y . So, this is y and this is x . So, the value of y^* y^* is the value of obtain from f_1 over d_1 . So, I call that as y^* .

So, y^* equal to f_1 by d_1 , and this comes out to be equal to 2.54 , 2.54 , and then the value of x^* we are putting star because their approximate values there we are taking t equal to 3 . So, x^* comes out to be 1.36 6 . So, when you do these computations in the digital computer, with t equal to 3 and rounding. Then using the rounding rule, we will get these values of x^* and y^* .

Now, the exact solution for the given system is x equal to exact solution is; exact solution is x equal to 1.357, and y equal to 2.543. So, the absolute error will be equal to if you calculate the absolute error E_x equal to mod of x minus x star it is 3×10^{-3} and the absolute error in y mod of y minus y star is 3×10^{-3} .

(Refer Slide Time: 09:45)

The exact solution for the given system is $x = 1.357$ and $y = 2.543$.
 So, the absolute error $E_x = |x - x^*| = 3 \times 10^{-3}$ and the absolute error
 $E_y = |y - y^*| = 3 \times 10^{-3}$.
 The relative errors are $R_x = \left| \frac{x - x^*}{x} \right| = 2.2108 \times 10^{-3}$ and
 $R_y = \left| \frac{y - y^*}{y} \right| = 1.1797 \times 10^{-3}$.
 The relative error can be used to measure the number of significant
 digits in an approximate value.

Again, with t equal to 3, the relative errors come out to be relative errors are R_x equal to mod of x minus x star over x which is 2.2108×10^{-3} and R_y equal to mod of y minus y star upon y which is 1.1797×10^{-3} .

And so, we can say that the relative error, now the relative error can be used to measure the number of significant digits in an approximate value, we are going to see that how a relative error can be used to measure the number of significant digits in an approximate value.



(Refer Slide Time: 10:36)

Definition: The number x^* is said to approximate a real number x to k significant digits if k is the largest non-negative integer for which

$$R_x = \left| \frac{x - x^*}{x} \right| < \frac{1}{2} \times 10^{-k} = 5 \times 10^{-(k+1)}.$$

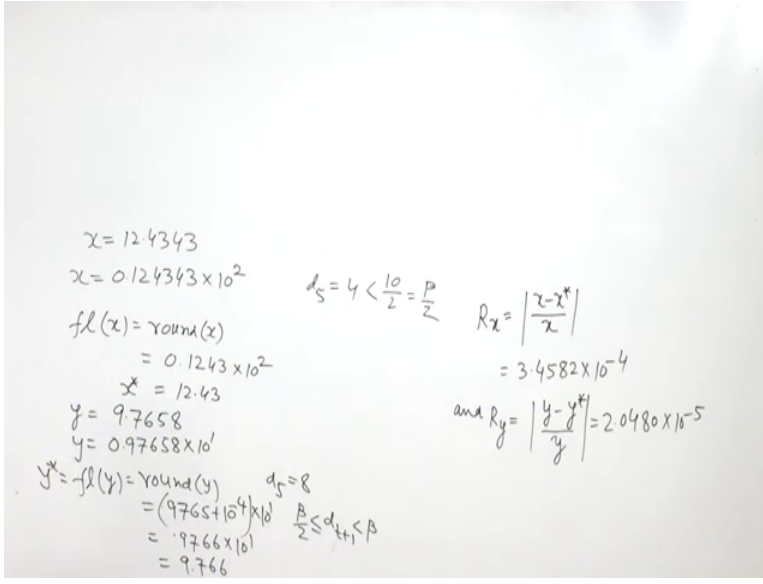
Example: Let us consider the 4-digit decimal system with rounding i.e. $\beta = 10$ and $t = 4$. Let $x = 12.4343$, $x^* = 12.43$, $y = 9.7658$, $y^* = 9.766$ then $R_x = 4.5841 \times 10^{-4}$ and $R_y = 2.0480 \times 10^{-4}$.

Therefore x^* approximates x to 3 significant digits and y^* approximates y to 4 significant digits.



7

Let us now consider the 4-digit decimal system with rounding ok, we are taking beta equal to 10, t equal to 4, if you take x equal to 21.4343 ok.

(Refer Slide Time: 10:47)



$x = 12.4343$
 $x = 0.124343 \times 10^2$
 $f(x) = \text{round}(x)$
 $= 0.1243 \times 10^2$
 $x^* = 12.43$
 $y = 9.7658$
 $y = 0.97658 \times 10^1$
 $y^* = \text{round}(y) = \text{round}(0.97658 \times 10^1)$
 $= (9765 + 10^4) \times 10^{-4}$
 $= 9766 \times 10^{-4}$
 $= 9.766$

$d_5 = 4 < \frac{10}{2} = \frac{\beta}{2}$
 $R_x = \left| \frac{x - x^*}{x} \right| = 3.4582 \times 10^{-4}$
 $\text{and } R_y = \left| \frac{y - y^*}{y} \right| = 2.0480 \times 10^{-5}$

$d_5 = 8$
 $\frac{\beta}{2} \leq d_{t+1} < \beta$

Then in the floating-point representation x will be written as 0.124343 into 10 to the power 2 ok.

So, here you can see we are taking t equal to 4; that means, we have to written 4 significant digits. So, this is $d_1 d_2 d_3 d_4$. Now d_5 , d_5 is 4 and 4 is less than 10 by 2 that is beta by 2 ok. So, what we will do? We will write it as 0.1243; we will leave it we

will write it as $0.01d_2 d_3 d_4$ ok. And remaining digits we discard into 10 to the power 2 . So, this will be 12.43 .

Now, if you take y equal to 9.7658 , then in the floating-point representation y is equal to 0.97658 into 10 to the power 1 ok. So, $f_1 y$ equal to round y . Now d_1 is 9 , d_2 is 7 , d_3 is 6 , d_4 is 5 , d_5 is 8 ok. d_5 is equal to 8 so, this d_5 lies in the range $\beta - 1$ less than or equal to $d_5 + 1$ less than β , β is equal to 10 ok. t is equal to 4 ok. So, d_5 is 8 so, the; so, we will what we will do? We will add 2.9765 , we add 10 to the power minus d equal to 4 ok, t is equal to 4 . So, we add 10 to the power minus 4 . So, this becomes into 10 to the power 1 . So, this becomes 0.9766 into 10 to the power 1 . So, this is 9.766 . So, this is y^* , and this is you are x^* . So, x^* is 12.43 y^* is equal to 9.766 .

Let us find the relative errors R_x and R_y . So, R_x is equal to $\frac{|x - x^*|}{x}$, then this is 3.4582 into 10 to the power minus 4 . And R_y is equal to 2.0480 into 10 to the power minus 5 .

Now, let us look at the definition, from the definition it is clear that in the case of R_x where k is equal 3 while in the case of R_y k is equal to 4 . So, x^* approximates x to 3 significant digits, and y^* approximates y to 4 significant digits.


(Refer Slide Time: 14:48)

Example: Let $f(x) = (x - 1)^3$. Let us evaluate $f(x)$ by

(I) $f(x) = x^3 - 3x^2 + 3x - 1$ and

(II) $f(x) = ((x - 3)x + 3)x - 1$ (nested multiplication).

Let y^* and z^* be the values of $f(2.72)$ by schemes (I) and (II) in a machine with 3 digit decimal system with rounding. Then $y^* = 5.08$ and $z^* = 5.09$ while actual value $y = 5.088448$.



Now let us take the example function evolution, let us say $f(x)$ equal to $x^3 - 3x^2 + 3x - 1$ whole cube we can evaluate $f(x)$ by 2 by direct method $f(x)$ equal to $x^3 - 3x^2 + 3x - 1$, by expanding $(x - 1)^3$, and another one is by nested multiplication which we discussed in the previous lecture.

So, when we do the nested multiplication, the $x^3 - 3x^2 + 3x - 1$ will be written as $(x - 3)(x + 3) + 1$ and then multiplied by $x - 1$. Now let us say y^* and z^* be the values of $f(2.72)$. So, we are taking x equal to 2.72 here. And evaluate the value of $f(x)$, when x is equal to 2.72 by schemes 1 and 2.

So, when we apply the scheme 1; that is, $f(x) = x^3 - 3x^2 + 3x - 1$, and evaluate f equal to f at 2.72, the value of y^* comes out to be 5.08 with 3-digit decimal system with rounding. And the z^* comes out to be 5.09 by the nested multiplication. So, actual value is 5.088448, and you can see that z^* is closer to the actual value than y^* . So, z^* actually it is calculated by nested multiplication. So, we see that nested multiplication is a better numerical scheme.

(Refer Slide Time: 16:33)

The relative errors are $R_y = \left| \frac{y - y^*}{y} \right| = 1.6602 \times 10^{-3}$ and

$$R_z = \left| \frac{z - z^*}{z} \right| = 2.8173 \times 10^{-4}.$$

Hence y^* approximates y to 2 significant digits while z^* approximates y to 3 significant digits. Hence z^* gives a better approximation than y^* for the actual value of y .

IF KOOBEE NPTEL ONLINE CERTIFICATION COURSE 9

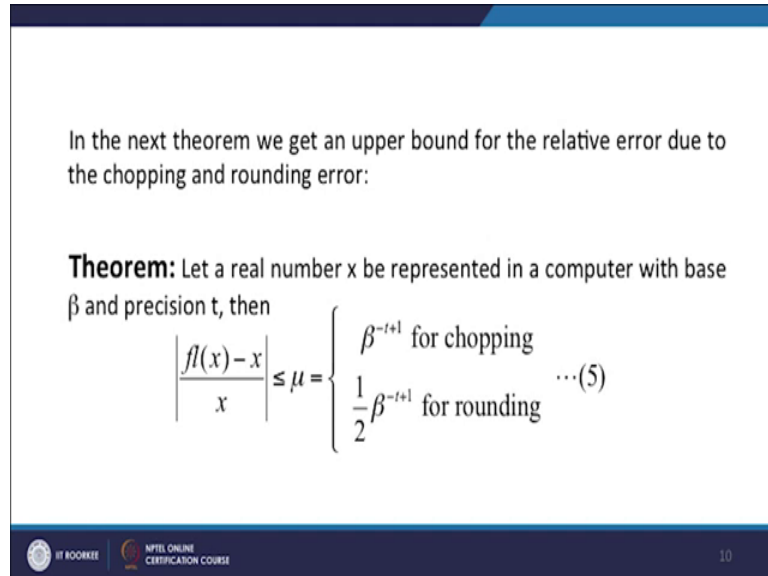
Now, let us calculate the relative errors, the relative error in y that is R_y is equal to $\frac{|y - y^*|}{y}$, which is 1.6602 into 10 to the power minus 3, and the relative error R_z is $\frac{|z - z^*|}{z}$, which is equal to 2.8173 into 10 to the power minus 4. So, y^* approximates y to 2 significant digits, while z^* approximates y to 3 significant digits. And therefore, we can say that z^* gives better

approximation to y better approximation to y then y star. So, the nested multiplication scheme is better than the direct method.

(Refer Slide Time: 17:18)

In the next theorem we get an upper bound for the relative error due to the chopping and rounding error:

Theorem: Let a real number x be represented in a computer with base β and precision t , then

$$\left| \frac{fl(x) - x}{x} \right| \leq \mu = \begin{cases} \beta^{-t+1} & \text{for chopping} \\ \frac{1}{2}\beta^{-t+1} & \text{for rounding} \end{cases} \quad \dots(5)$$


The slide contains the following text and equation:

In the next theorem we get an upper bound for the relative error due to the chopping and rounding error:

Theorem: Let a real number x be represented in a computer with base β and precision t , then

$$\left| \frac{fl(x) - x}{x} \right| \leq \mu = \begin{cases} \beta^{-t+1} & \text{for chopping} \\ \frac{1}{2}\beta^{-t+1} & \text{for rounding} \end{cases} \quad \dots(5)$$

At the bottom of the slide, there are logos for IIT ROORKEE and NPTEL ONLINE CERTIFICATION COURSE, and the page number 10.

Now, in the next round we get an upper bound for the relative error due to the chopping and rounding error. We have seen that when the floating-point representation of real number in the detail machine is used. There are 2 kinds of errors, there are 2 kinds of way in which error occurs chopping and rounding.

So, let us see what is the upper bound for the relative error, due to the chopping and rounding let us assume that a real number x be represented in a computer with base beta and precision t , then so, precision t means t number of significant digits ok. So, relative error we know, relative error is x minus $fl\ x$ divided by x .

(Refer Slide Time: 18:09)

Hence
 Relative error = $\left| \frac{x - \text{chop}(x)}{x} \right| \leq \frac{\beta^{-t+e}}{|x|}$
 $= \frac{\beta^{-t+e}}{(0.d_1d_2\dots d_t d_{t+1}\dots)_\beta \times \beta^e}$
 $\leq \frac{(0.1000\dots 0)_\beta \times \beta^e}{\beta}$
 $x = \pm (0.d_1d_2\dots d_t d_{t+1}\dots)_\beta \times \beta^e$
 $\text{fl}(x) = \text{chop}(x)$
 $= \pm (0.d_1d_2\dots d_t)_\beta \times \beta^e$
 Thus, if $\gamma = \beta - 1$

Then $0 \leq d_i \leq \beta - 1$
 $\neq 1/2$
 $|x - \text{chop}(x)|$
 $= (0.00\dots 0d_{t+1}d_{t+2}\dots)_\beta \times \beta^e$
 $\leq \left(\frac{\gamma}{\beta^{t+1}} + \frac{\gamma}{\beta^{t+2}} + \dots \right) \times \beta^e$
 $= \left(\frac{\gamma}{\beta^{t+1}} \right) \times \beta^e = \frac{\gamma}{\beta^{t+1}} \times \beta^e$
 $= \frac{\gamma}{\beta^{-t+e}}$

Now, this is less than this is defined as less than or equal to beta mu, mu is the upper bound this is in the case of chopping, it is beta to the power minus t plus 1 for chopping. And in the case of rounding, it is half of this value. Now let us move this theorem. So, we give the proof first we give the proof for the round off error due to chopping. So, x is equal to in the floating-point representation x is written like this.

(Refer Slide Time: 19:15)

Proof:(i) We give the proof for the round-off error due to the chopping.

Let $x = \pm (0.d_1d_2\dots d_t d_{t+1}\dots)_\beta \times \beta^e$.

Then $\text{chop}(x) = \pm (0.d_1d_2\dots d_t)_\beta \times \beta^e$.

Hence $|x - \text{chop}(x)| = (0.0\dots 0d_{t+1}\dots)_\beta \times \beta^e$.

Now, let $\gamma = \beta - 1$, then $|x - \text{chop}(x)| \leq (0.0\dots 0\gamma\gamma\dots)_\beta \times \beta^e$

$$= \left(\frac{\gamma}{\beta^{t+1}} + \frac{\gamma}{\beta^{t+2}} + \dots \right) \times \beta^e$$

$$= \beta^{-t+e}.$$

IT KOOKEE | NPTEL ONLINE CERTIFICATION COURSE | 11

Now, we are trying to find the upper bound this, upper bound for the chopping case ok. So, flx equal to chop x to we are taking we are taking precision t ok. So, precision t

means we are taking t significant digits. So, this means after d_t , we discovered all the remaining digits. So, we have plus minus $0.d_1d_2$ and so on d_t beta to the power e ok thus, if you define gamma equal to beta minus 1.

If you define gamma equal to beta minus 1, then $\text{mod of } x \text{ minus chop } x \text{ minus chop } x \text{ chop } x$ means $x \text{ minus flx}$ this will be equal to $\text{from } x \text{ we are subtracting chop } x$. So, and we are taking modulus. So, $0 \text{ point now } d_1 d_2 d_t$ they are exactly same. So, we will have $0 \text{ } 0$ and so on 0 , and then we will have $d_t \text{ plus } 1 \text{ } d_t \text{ plus } 2$ and so on ok.

This means that first t places are 0s, after that we have $d_t \text{ plus } 1 \text{ } d_t \text{ plus } 2$ and so on, the now let us see that in the case of normalised floating-point representation. d_1 is lying between one and beta minus 1, but $d_2 d_3$ and so on, all d_i for i greater than or equal to they are mod than equal to 0, but less than or equal to beta minus 1 ok.

So, d_i is lying between 0, and beta minus 1 for all i greater than or equal to 2 ok. So, this means that we shall have this is less than or equal to $d_t \text{ plus } 1$, $d_t \text{ plus } 1$ will be less than or equal to beta minus 1 and beta minus 1 we are taking as gamma. So, gamma over beta to the power $t \text{ plus } 1$ ok beta is the base and this is $t \text{ plus } 1$ position, then gamma over beta to the power $t \text{ plus } 2$ and so on.

All this digits $d_t \text{ plus } 1 \text{ } d_t \text{ plus } 2$ they are less than or equal to beta minus 1 that is they are less than or equal to gamma into ok. So, this is this beta we do not write; now in to beta to the power e this is infinite series the geometric series, where 1 by beta is the common ratio. So, gamma beta to the power $t \text{ plus } 1$ divided by $1 \text{ minus } 1 \text{ by beta}$ we have into beta to the power e and this is equal to. So, this is beta gamma upon beta to the power $t \text{ plus } 1$ divided by beta minus 1 upon beta, beta minus 1 is gamma. So, gamma upon beta into beta to the power e ok and this is nothing but beta to the power minus $t \text{ plus } e$.

Now, let us find the relative error ok. So, hence relative error equal to $\text{mod of } x \text{ minus flx}$ divided by x ok. $\text{Mod of } x \text{ minus flx}$ is $\text{chop } x$. So, this is less than or equal to beta to the power minus $t \text{ plus } e$ divided by x ok, $\text{mod of } x$. $\text{Mod of } x$ is this quantity. So, this is equal to beta to the power minus $t \text{ plus } e$ divided by now this is $0.d_1d_2 \text{ } d_t \text{ } d_t \text{ plus } 1$ and so on, beta into beta to the power e ok.

Now, we have to make it less than or equal to, this means that we put the minimum values of $d_1, d_2, \dots, d_t, d_{t+1}$ and so on. So, that this quantity is maximised. So, beta to the power minus t plus e divided by now minimum value of d_1 is one because d_1 lies from 1 to $\beta - 1$, I mean between one and beta minus 1. So, 0.1 and then $d_2, d_3, \dots, d_t, d_{t+1}$ can be taken as 0s and so on into beta to the power e. And this gives us this cancels with this beta e e e and what we get beta to the power minus t divided by now this is 1 over beta ok 1 over beta, then 0 over beta is square, then 0 over beta cube and so on.

(Refer Slide Time: 25:20)

Hence

$$\text{Relative error} = \left| \frac{x - \text{Chop}(x)}{x} \right| \leq \frac{\beta^{-t+e}}{|x|}$$

$$= \frac{\beta^{-t+e}}{(0.d_1d_2\dots d_t d_{t+1}\dots)_\beta \times \beta^e}$$

$$\leq \frac{(0.1000\dots 0)_\beta \times \beta^e}{\beta^{-t+e}}$$

$$= \frac{\beta^{-t}}{\beta}$$

$$= \beta^{-t+1}$$

Then $0 \leq d_i \leq \beta - 1$
#i>2

$$|\text{Chop}(x)| = (0.00\dots 0d_{t+1}d_{t+2}\dots)_\beta \times \beta^e$$

$$\leq \left(\frac{Y}{\beta^{t+1}} + \frac{Y}{\beta^{t+2}} + \dots \right) \times \beta^e$$

$$= \left(\frac{Y}{\beta^{t+1}} \right) \times \beta^e = \frac{Y}{\beta} \times \beta^e$$

$$= \frac{Y}{\beta^{-t+e}}$$

So, this is equal to beta to the power minus t plus 1. So, this is the upper bound in the case of chopping.



Now, we go to the next case, where we have rounding. So, in the case of rounding, again let us say x is equal to plus minus $0.d_1d_2 \dots d_t d_{t+1}$, base is beta into beta to the power a, and beta to the power e now first we consider the case, when 0 is less than or equal to d_{t+1} less than beta by 2.

(Refer Slide Time: 26:09)

(ii) Case-1 We have
 Let $x = \pm (0.d_1d_2\dots d_t d_{t+1}\dots)_\beta \times \beta^e$.
 Then $\text{round}(x) = \pm (0.d_1d_2\dots d_t)_\beta \times \beta^e, 0 \leq d_{t+1} < \frac{\beta}{2}$.
 Hence $|x - \text{round}(x)| = (0.0\dots 0d_{t+1}d_{t+2}\dots)_\beta \times \beta^e$.

$$\leq \left(\frac{\beta}{2} - 1 + \frac{\gamma}{\beta^{t+2}} + \frac{\gamma}{\beta^{t+3}} + \dots \right) \times \beta^e$$

$$= \frac{1}{2} \beta^{-t+e}.$$



13



Then the rounding of x gives plus minus $0.d_1d_2$ and so on d_t base β into β to the power e . And therefore, $\text{mod of } x \text{ minus flx}$, or we can say $\text{round } x$ is equal to 0 point, now first d first t places are same in x and $\text{round of } x$. So, we have $0.0\dots 0$ and then we have $d_{t+1}d_{t+2}$ and so on with base β into β to the power e .

Now, when we look at the place d_{t+1} d_{t+1} is less than β by 2. Therefore, d_{t+1} is less than or equal to β by 2 minus 1. So, d_{t+1} is less than or equal to β by 2 minus 1 will give you, β by 2 minus 1 upon β to the power $t+1$. And then $d_{t+2}d_{t+3}$ and so on, they are all less than or equal to β minus 1 the condition is only on d_{t+1} .

So, we have written them as γ , because of the each one of them is less than or equal to β minus 1. So, we are and β minus 1 we are assuming as γ . So, $\frac{\gamma}{\beta^{t+2}} + \frac{\gamma}{\beta^{t+3}} + \dots$ and so on into β to the power e , and when you evaluate this infinite series hence where $\frac{\gamma}{\beta^{t+2}}$ is the first term. And $\frac{1}{\beta}$ is the geometric ratio, then and then you simplify this it turns out that it is $\frac{1}{2} \beta^{-t+e}$.

(Refer Slide Time: 28:58)

$$\begin{aligned}
 \text{Absolute relative error} &= \left| \frac{x - \text{round}(x)}{x} \right| \\
 &\leq \frac{\frac{1}{2} \beta^{-t+e}}{(0.10000\dots)_\beta \times \beta^e} \\
 &= \frac{1}{2} \beta^{-t+1}.
 \end{aligned}$$



14

So, absolute relative error in this case will be or you can say relative error in this case will be mod of x minus round of x divided by x. Now so, we have x minus round of x less than or equal to $\frac{1}{2} \beta^{-t+e}$ and in the denominator, we are writing x is equal to plus minus $(0.d_1 d_2 d_3 \dots d_t d_{t+1} \dots)_\beta \times \beta^e$ and so on.

(Refer Slide Time: 28:20)

$$\begin{aligned}
 x &= \pm (0.d_1 d_2 d_3 \dots d_t d_{t+1} \dots)_\beta \times \beta^e \\
 |x| &\geq (0.1000\dots 0)_\beta \times \beta^e \\
 &= \frac{1}{\beta} \times \beta^e
 \end{aligned}$$

Beta into beta to the power e beta to the power now, we want make it this quantity less than or equal to less than or equal to means we put the minimum values of $d_1 d_2 d_3$ and so on. So, mod of x is less than or equal to greater than or equal to in the

denominator we put minimum values. So, this is greater than or equal to 0.1 minimum value of d_1 is one and then d_2, d_3 they are all 0s beta into beta to the power e . So, we have this and this is equal to 1 by beta in to beta to the power e . So, beta 2 power e will cancel and we will get 1 by 2 beta to the power minus t plus 1.

(Refer Slide Time: 29:26)

(ii) Case-2 We have

Let $x = \pm (0.d_1d_2\dots d_t d_{t+1}\dots)_\beta \times \beta^e$.

Then $\text{round}(x) = \pm \left\{ (0.d_1d_2\dots d_t)_\beta + \beta^{-t} \right\} \times \beta^e, \frac{\beta}{2} \leq d_{t+1} < \beta$.

Hence $|x - \text{round}(x)| = \left| (0.00\dots 1)_\beta - (0.0\dots 0d_{t+1}d_{t+2}\dots)_\beta \right| \times \beta^e$.

$$\leq \left(\frac{1}{\beta^t} - \frac{d_{t+1}}{\beta^{t+1}} - \frac{0}{\beta^{t+2}} - \frac{0}{\beta^{t+3}} - \dots \right) \times \beta^e$$

$$= \left(\frac{1}{\beta^t} - \frac{1}{\beta^{t+1}} \times \frac{\beta}{2} \right) \times \beta^e = \frac{1}{2} \beta^{-t+e}.$$

Now, let us go to the next for case here in the case of rounding ok. So, again x is equal to plus minus $0.d_1d_2$ and so on $d_t d_{t+1}$, with base beta into beta to the power e , and when we do the rounding in the second case, in the second case d_{t+1} is greater than or equal to beta by 2, but less than beta. So, then we add beta to the power minus t to the $2 d_t$, we add 1 by beta to the power t .

So, so, $x \bmod$ of x minus $\text{round } x$ in this case will be what x is equal to plus minus $0.d_1d_2 d_t d_{t+1}$ and so on into beta to the power e .

(Refer Slide Time: 30:15)

$$\begin{aligned} \text{Relative error} &= \left| \frac{x - \text{round}(x)}{x} \right| \leq \frac{\frac{1}{2} \beta^{-t+k}}{(0.d_1 d_2 \dots d_t d_{t+1} \dots) \times \beta^e} \leq \frac{1}{2} \frac{\beta^{-t}}{(0.10 \dots 0 \frac{\beta-1}{2} \dots 0) \beta} = \frac{1}{2} \frac{\beta^{-t}}{\frac{1}{\beta} + \frac{\beta}{2} \beta^{-t+1}} \\ &\leq \frac{1}{2} \frac{\beta^{-t}}{\frac{1}{\beta}} \\ &= \frac{1}{2} \beta^{-t+1} \end{aligned}$$

$$\begin{aligned} x &= \pm (0.d_1 d_2 \dots d_t d_{t+1} \dots) \times \beta^e \\ \text{round}(x) &= \pm \left((0.d_1 d_2 \dots d_t) + \beta^{-t} \right) \times \beta^e \end{aligned}$$

$$\begin{aligned} |x - \text{round}(x)| &= \left(\frac{1}{\beta^t} - \frac{d_{t+1}}{\beta^{t+1}} - \frac{d_{t+2}}{\beta^{t+2}} \dots \right) \times \beta^e \\ &\leq \left(\frac{1}{\beta^t} - \frac{\beta}{2 \beta^{t+1}} - \frac{0}{\beta^{t+2}} - \frac{0}{\beta^{t+3}} \dots \right) \times \beta^e \\ &= \frac{1}{2} \beta^{-t+1} \end{aligned}$$

$$\begin{aligned} \frac{\beta}{2} &\leq d_{t+1} < \beta \\ 0 &\leq d_i \leq \beta - 1 \\ &\forall i \geq t+2 \end{aligned}$$

Round x is equal to plus minus 0-point $d_1 d_2 \dots d_t$, plus beta to the power minus t into 10 to beta to the power e . So, $2 d_t$ we are adding 1 by beta to the power t . So, this will become greater so, we will have mod of x minus round x . Numerically it will be more than this quantity numerically will be more than this, because to the t th decimal we are adding 1 by beta to the power t ok.

So, this is equal to now what will happen? $d_1 d_2 \dots d_t d_{t+1} d_{t+2} \dots$ will cancel, but d_t plus 1 by beta to the power t we have. So, we will have this less than or equal to 1 by beta to the power t , and then we have d_t plus 1 upon beta to the power t plus 1 and so on. We are we are subtracting them so, this is equal this is this is equal to 1 by beta to the power t minus ok, this is less than this is equal to this, and this is how much? this is less than or equal to 1 by beta to the power t . So, we put now we are these are subtracting. So, this means we have to put in order to make less than or equal to we have to put minimum values of d_t plus 1 d_{t+1} plus 2 and so on ok. And this will give you what? d_t plus 1 yeah d_t plus 1 is beta by 2 is less than or equal to d_t plus 1 and less than beta ok.

So, we want to maximise this. So, we put minimum values ok, this is for the d_t plus and d_t d_i 0 less than or equal d_i less than or equal to beta minus 1 for all i bigger than or equal to t plus 2 ok. So, d_t plus 2 and so on we will put as 0s and d_t plus 1 we will put as beta by 2 ok. So, 1 by beta t minus beta by 2 into β^t plus 1, minus 0 upon beta t to the power t plus 2 and so on into beta to the power e ok. And this will give you so, this will

beta will cancel, and we will get 1 over beta to the power t minus 1 over 2 beta to the power t so, 1 by 2 beta to the power minus t plus e.

Now, we go to the relative error. So, relative error will be mod of x minus round x divided by x ok. And this is less than or equal to 1 by 2 beta to the power minus t plus e, x is equal to mod of x is equal to 0.d1d2 dt dt plus 1 and so on, into beta to the power beta to the power e ok. And what we will get?

Now, we have to put minimum values here of d 1 d 2 d t dt plus 1. So, that we get upper bound. So, 1 by 2 and this beta to the power e will cancel with this, we have beta to the power minus t, in the denominator minimum value of d 1 is 1. So, 0.1 d 2 dt they are all 0's minimum values are 0's dt plus 1 has minimum value beta by 2 ok. Minimum values are beta by 2. So, we have and dt plus 1 they can all be taken as 0. So, we will have ok, this is now this is how much? 1 by 2 beta to the power minus t divided by 1 by beta, plus and this is the t plus 1 th position, beta by 2 into beta to the power t plus 1 ok.

Now, this is further less than or equal to this is a positive quantity, we can drop this and 1 by 2 beta to the power minus t divided by 1 by beta. And we get it as 1 by 2 beta to the power minus t plus 1. So, in the case of dt plus 1 more than or equal to beta by 2 plus less than beta again we get half of beta to the power minus t plus 1 ok.

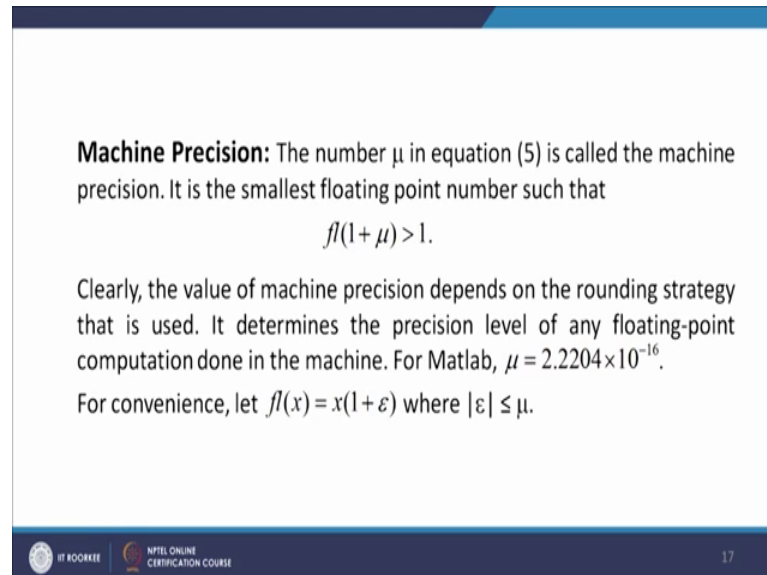
(Refer Slide Time: 26:09)

$$\begin{aligned}
 \text{Relative error} &= \left| \frac{x - \text{round}(x)}{x} \right| \\
 &\leq \frac{\frac{1}{2}\beta^{-t+e}}{(0.100\dots 0d_{t+1}00\dots)_\beta \times \beta^e} \\
 &= \frac{\frac{1}{2}\beta^t}{\frac{1}{\beta} + \frac{\beta}{2} \times \frac{1}{\beta^{t+1}}} \leq \frac{1}{2}\beta^{-t+1}.
 \end{aligned}$$

IT ROOKIE NPTEL ONLINE CERTIFICATION COURSE 16

Now, we go to a machine precision. The number μ in the equation number 5, this equation number 5 the number μ is called the machine precision.

(Refer Slide Time: 36:18)



Machine Precision: The number μ in equation (5) is called the machine precision. It is the smallest floating point number such that

$$fl(1 + \mu) > 1.$$

Clearly, the value of machine precision depends on the rounding strategy that is used. It determines the precision level of any floating-point computation done in the machine. For Matlab, $\mu = 2.2204 \times 10^{-16}$.

For convenience, let $fl(x) = x(1 + \varepsilon)$ where $|\varepsilon| \leq \mu$.

IT KOOBEE | NPTEL ONLINE CERTIFICATION COURSE 17

It is this minus floating-point number such that $fl(1 + \mu)$ is greater than 1. The value of machine precision depends on the rounding strategy that we use. You can see that in the case of chopping it is β to the power t plus 1 while in the case of rounding it is half of that in the case of chopping. So, it determined the precision level of any floating-point computation done in the machine.

In the case of MATLAB this μ is 2.2204×10^{-16} . Now when we do the of addition multiplication division in floating point representation we shall be taking $fl(x)$ equal to $x(1 + \varepsilon)$ ok.

(Refer Slide Time: 37:26)

Relative error
 $= \left| \frac{x - \text{round}(x)}{x} \right| \leq \frac{\frac{1}{2} \beta^{-t+k}}{(0.d_1 d_2 \dots d_t d_{t+1} \dots) \times \beta^e} \leq \frac{1}{2} \frac{\beta^{-t}}{(0.10 \dots 0 \frac{\beta-1}{2} \dots 0)} \beta = \frac{1}{2} \frac{\beta^{-t}}{\frac{1}{\beta} + \frac{\beta}{2} \beta^{t+1}}$

$x = \pm (0.d_1 d_2 \dots d_t d_{t+1} \dots) \times \beta^e$

$\text{round}(x) = \pm \left\{ (0.d_1 d_2 \dots d_t) + \beta^{-t} \right\} \times \beta^e$

$\frac{\beta}{2} \leq d_{t+1} < \beta$
 $0 \leq d_i \leq \beta - 1$
 $\forall i \geq t+2$

$x - \text{round}(x) = \left(\frac{1}{\beta^t} - \frac{d_{t+1}}{\beta^{t+1}} - \frac{d_{t+2}}{\beta^{t+2}} - \dots \right) \times \beta^e$

$\left| \frac{x - \text{round}(x)}{x} \right| < \mu$

$\frac{x - \text{round}(x)}{x} = \epsilon$
 $f(x) = x(1 + \epsilon)$
 where $|\epsilon| < \mu$

$\leq \left(\frac{1}{\beta^t} - \frac{\beta}{2\beta^{t+1}} - \frac{0}{\beta^{t+2}} - \frac{0}{\beta^{t+3}} - \dots \right) \times \beta^e$

$= \frac{1}{2} \beta^{-t+e}$

See we have we have flx equal to x minus flx mod of x minus flx divided by x this is equal to this is relative error this is less than or equal to mu we have taken ok.

So, what we do is let us take x minus flx divided by x equal to mu x is equal to epsilon, then we shall write flx equal to x times 1 plus epsilon where epsilon is such that mod of epsilon is less than or equal to mu. With that I would like to end my lecture.

Thank you very much.