

LECTURE 14

DIGITAL FILTERS - I

Now that we have examined computers and how they represent information let us turn to how computers process information. We can, of course, only examine a very few of the things they do, and will concentrate on basics per usual.

Much of what computers process are signals from various sources, and we have already discussed why they are often in the form of a stream of numbers from an equally spaced sampling system. Linear processing, which is the only one I have time for in this course, implies digital filters. To illustrate "style" and how things actually happen in real life I propose to tell you first how I became involved in them, and then how I proceeded.

First, I never went to the office of my Vice President W. O. Baker; we only met in passing in the halls and we usually stopped to talk a few, very few, minutes. One time, around 1973-4, when I met him in a hall I said to him that when I came to Bell Telephone Laboratories in 1946 I had noticed that the Laboratories were gradually passing from relay to electronic central offices, but that a large number of people would not convert to oscilloscopes and the newer electronic technology and that they were moved to a different location to get them out of the way. To him they represented a serious economic loss but to me they were a social loss since they were disgruntled to say the least because they were passed by, (though it was their own fault). I went on to say that I had seen the same thing happen when we went from the earlier analog computers (on which Bell Telephone Laboratories had many experts because they had developed much of the technology during WWII) to the more modern digital computers - that we again left a large number of engineers behind, and again they were both an economic and a social loss. I then observed that we both knew the telephone company was going to total digital transmission about as fast as they could, and that this time we would leave behind a very much larger number of disgruntled engineers. Hence, I concluded, we should do something now about the situation, such as get adequate elementary books and other training devices to ease more of them into the future and leave fewer behind. He looked me square in the eye and said, "Yes, Hamming, you should." and walked off! Furthermore, he went on encouraging me, via John Tukey with whom he often spoke, so I knew he was watching my efforts.

What to do? In the first place I thought I knew very little about digital filters, and, furthermore, I was not really interested in them. But does one wisely ignore one's V.P. plus the cogency of one's own observations? No! The implied social waste was too high for me to contemplate comfortably.

So I turned to a friend, Jim Kaiser, (J. F. Kaiser), who was one of the world's experts in digital filters at that time, and suggested that he should stop his current research and write a book on digital filters - that book writing to summarize his work was a natural stage in the development of a scientist. After some pressure he agreed to write the book, so I was saved, so I thought. But monitoring what he was doing revealed that he was writing nothing. To rescue my plan I offered, if he would educate me over lunches in the restaurant (you get more time to think there than in the cafeteria), to help write the book jointly, (mainly the first part), and we could call it Kaiser and Hamming. Agreed!

As time went on I was getting a good education from him, and I got my first part of the book going but he was still writing nothing. So one day I said, "If you don't write more we will end up calling it Hamming and Kaiser." - and he agreed. Still later when I had about completed all the writing and he had still written nothing, I said that I could thank him in the preface, but it should be called Hamming, and he agreed - and we are still good friends! That is how the book on Digital Filters that I wrote came to be, and I saw it ultimately through three editions, always with good advice from Kaiser.

The book also took me many places that were interesting since I gave a short, one week courses, on it for many years. The short courses began while I was still writing it because I needed feedback and had suggested to UCLA Extension Division that I give it as a short course, to which they agreed. That led to years of giving it at UCLA, once in each of Paris, London, and Cambridge, England, as well as many other places in the USA and at least twice in Canada. Doing what needed to be done, though I did not want to do it, paid off handsomely in the long run.

Now, to the more important part, how I went about learning the new subject of digital filters. Learning a new subject is something you will have to do many times in your career if you are to be a leader and not be left behind as a follower by newer developments. It soon became clear to me that digital filter theory was dominated by Fourier series, about which theoretically I had learned in college, and actually I had had a lot of further education during the signal processing I had done for John Tukey, who was a professor from Princeton, a genius, and a one or two day a week employee of Bell Telephone Laboratories. For about ten years I was his computing arm much of the time.

Being a mathematician I knew, as all of you do, that any complete set of functions will do about as good as any other set at representing arbitrary functions. Why, then, the exclusive use of the Fourier series? I asked various Electrical Engineers and got no satisfactory answers. One engineer said that alternating currents were sinusoidal, hence we used sinusoids, to which I replied it made no sense to me. So much for the usual residual education of the typical Electrical Engineer after they have left school!

So I had to think of basics, just as I told you I had done when using an error detecting computer. What is really going on? I suppose that many of you know that what we want is a time invariant representation of signals since there is usually no natural origin of time. Hence we are led to the trigonometric functions, (the eigenfunctions of translation), in the form of both Fourier series and Fourier integrals, as the tool for representing things.

Second, linear systems, which is what we want at this stage, also have the same eigenfunctions - the complex exponentials which are equivalent to the real trigonometric functions. Hence a simple rule: If you have either a time invariant system, or a linear system, then you should use the complex exponentials.

On further digging into the matter I found yet a third reason for using them in the field of digital filters. There is a theorem, often called "Nyquist's sampling theorem", (thought it was known long before and even published by Whittaker in a form you can hardly realize what it is saying even when you know Nyquist's theorem), which says, if you have a band limited signal and sample at equal spaces at a rate of at least two in the highest frequency, then the original signal can be reconstructed from the samples. Hence the sampling process loses no information when we replace the continuous signal with the equally spaced samples, provided the samples cover the whole real line. The sampling rate is often known as "the Nyquist rate" after Harry Nyquist, also of servo stability fame as well as other things. If you sample a nonbandlimited function, then the higher frequencies are "aliased" into lower ones, a word devised by Tukey to describe the fact that a single high frequency will appear later as a single low frequency in the Nyquist band. The same is not true for any other set of functions, say powers of t . Under equal spaced sampling and reconstruction a single high power of t will go into a polynomial (many terms) of lower powers of t .

Thus there are three good reasons for the Fourier functions: (1) time invariance, (2) linearity, and (3) the reconstruction of the original function from the equally spaced samples is simple and easy to understand.

Therefore we are going to analyse the signals in terms of the Fourier functions, and I need not discuss with electrical engineers why we usually use the complex exponents as the frequencies instead of the real trigonometric functions. We have a linear operation and when we put a signal (a stream of numbers) into the filter then out comes another stream of numbers. It is natural, if not from your linear algebra course, then from other things such as a course in differential equations, to ask what functions go in and come out exactly the same except for scale? Well, as noted above, they are the complex exponentials; they are the eigenfunctions of linear, time invariant, equally spaced sampled systems.

Lo, and behold, the famous transfer function is exactly the

eigenvalues of the corresponding eigenfunctions! Upon asking various Electrical Engineers what the transfer function was no one has ever told me that! Yes, when pointed out to them that it is the same idea they have to agree, but the fact that it is the same idea never seemed to have crossed their minds! The same, simple idea, in two or more different disguises in their minds, and they knew of no connection between them! Get down to the basics every time!

We begin our discussion with, "What is a signal?" Nature supplies many signals which are continuous, and which we therefore sample at equal spacing and further digitize, (quantize). Usually the signals are a function of time, but any experiment in a lab that uses equally spaced voltages, for example, and records the corresponding responses, is also a digital signal. A digital signal is, therefore, an equally spaced sequence measurements in the form of numbers, and we get out of the digital filter another equally spaced set of numbers. One can, and at times must, process nonequally spaced data, but I shall ignore them here.

The quantization of the signal into one of several levels of output often has surprisingly small effect. You have all seen pictures quantized to two, four, eight, and more levels, and even the two level picture is usually recognizable. I will ignore quantization here as it is usually a small effect, though at times it is very important.

The consequence of equally spaced sampling is aliasing, that is a frequency above the Nyquist frequency (which has two samples in the cycle) will be aliased into a lower frequency. This is a simple consequence of the trigonometric identity

$$\exp(2\pi i(k + a)n) = \exp(2\pi ian)$$

where a is the positive remainder after removing the integer number of rotations, k , (we always use rotations in discussing results, and use radians while applying the calculus, just as we use base 10 logs and base e logs), and n is the step number. If $a > 1/2$, then we can write the above as

$$\exp(2\pi ian) = \exp(-2\pi i(1 - a)n)$$

The aliased band, therefore, is less than $1/2$ a rotation, plus or minus. If we use the two real trigonometric functions, \sin and \cos , we have a pair of eigenfunctions for each frequency, and the band is from 0 to $1/2$ a rotation, but when we use the complex exponential notation then we have one eigenfunction for each frequency, but now the band reaches from $-1/2$ to $1/2$ rotations. This avoidance of the multiple eigenvalues is part of the reason that the complex frequencies are so much easier to handle than are the real sine and cosine functions. The maximum sampling rate for which aliasing does not occur is two samples in the cycle, and is called the Nyquist rate. From the samples the original signal cannot be determined to within the aliased frequencies, only the basic frequencies that fall in the fundamental interval of unaliased frequencies ($-1/2$ to $1/2$) can be

determined uniquely. The signals from the various aliased frequencies go to a single frequency in the band and are algebraically added; that is what we see once the sampling has been done. Hence addition or cancellation may occur during the aliasing, and we cannot know from the aliased signal what we originally had. At the maximum sampling rate one cannot tell the result from 1, hence the unaliased frequencies must be within the band.

We shall stretch (compress) time so that we can take the sampling rate to be one per unit time, because this makes things much easier and brings experiences from the milli and micro second range to those which may take days or even years between samples. It is always wise to adopt a standard notation and framework of thinking of diverse things - one field of application may suggest things to do in the other. I have found it of great value to do so whenever possible - remove the extraneous scale factors and get to the basic expressions. (But then I was originally trained as a mathematician.)

Aliasing is the fundamental effect of sampling and has nothing to do with how the signals are processed. I have found it convenient to think that once the samples have been taken then all the frequencies are in the Nyquist band, and hence we do not need to draw periodic extensions of anything since the other frequencies no longer exist in the signal - once the sampling has occurred the higher frequencies have been aliased into the lower band, and do not exist up there any more. A significant savings in thinking! The act of sampling produces the aliased signal that we must use.

I now turn to three stories that use only the ideas of sampling and aliasing. In the first story I was trying to compute the numerical solution to a system of 28 ordinary differential equations and I had to know the sampling rate to use, (the step size of the solution is the sampling rate you are using), since if it were half as large as expected then the computing bill would be about twice as much. For the most popular and practical methods of numerical solution the mathematical theory bases the step size on the fifth derivative. Who could know the bound? No one! But viewed as sampling, then the aliasing begins at two samples for the highest frequency present, provided you have data from minus to plus infinity. Having only a short range of at most five points of data I intuitively figured that I would need about twice the rate, or 4 samples per cycle. And finally, having only data on one side, perhaps another factor of 2; in all 8 samples per cycle.

I next did two things: (1) developed the theory, and (2) ran numerical tests on the simple differential equation

$$y'' + y = 0, \quad y(0) = 1, \quad y'(0) = 0$$

They both showed that at around 7 samples per cycle you are on the edge of accuracy, (per step), and at 10 you are very safe. So I explained the situation to them and asked them for the

highest frequencies in the expected solution. They saw the justice of my request, and after some days they said I had to worry about the frequencies up to 10 cycles per second and they would worry about those above. They were right, and the answers were satisfactory. The sampling theorem in action!

The second story involves a remark, made to me casually in the halls of Bell Telephone Laboratories that a certain West Coast subcontractor was having trouble with the simulation of a Nike missile launch, and was using 1/1000 to 1/10,000 of second spacing. I laughed immediately, and said that there must be some mistake, that 70 to 100 samples would be enough for the model they were using. It turned out that they had a binary number 7 position to the left, 128 times too large! Debugging a large program across the continent based on the sampling theorem!

The third story is that a group at Naval Postgraduate School was modulating a very high frequency signal down to where they could afford to sample, according to the sampling theorem as they understood it. But I realized that if they cleverly sampled the high frequency then the sampling act itself would modulate (alias) it down. After some days of argument, they removed the rack of frequency lowering equipment, and the rest of the equipment ran better! Again, I needed only a firm understanding of the aliasing effects due to sampling. It is another example of why you need to know the fundamentals very well; the fancy parts then follow easily and you can do things that they never told you about.

The sampling is fundamental to the way we currently process data, when we use the digital computers. And now that we understand what a signal is, and what sampling does to a signal, we can safely turn to more of the details of processing signals.

We will first discuss nonrecursive filters, whose purpose is to pass some frequencies and stop others. The problem first arose in the telephone company when they had the idea that if one voice message had all its frequencies moved up (modulated) to beyond the range of another then the two signals could be added and sent over the same wires, and at the other end filtered out and separated, and the higher one reduced (demodulated) back to its original frequencies. This shifting is simply multiplying by a sinusoidal function, and selecting one band (single sideband modulation) of the two frequencies that emerge according to the following trigonometric identity (this time we use real functions)

$$\cos at \cos bt = (1/2)[\cos(a + b)t + \cos(a - b)t]$$

There is nothing mysterious about the frequency shifting (modulation) of a signal, it is at most a variant of a trigonometric identity.

The nonrecursive filters we will consider first are mainly of the smoothing type where the input is the values $u(t) = u(n)$ = u_n and the output is y_n

$$y_n = \text{SUM}[j=-k, k; c_j u_{n-j}]$$

with $c_j = c_{-j}$, (the coefficients are symmetric about the middle value c_0).

I need to remind you about least squares as it plays a fundamental role in what we are going to do, hence I will design a smoothing filter to show you how filters can arise. Suppose we have a signal with "noise" added and want to smooth it, remove the noise. We will assume that it seems reasonable to you fit a straight line to 5 consecutive points of the data in a least squares sense, and then take the middle value on the line as the "smoothed value of the function" at that point.

For mathematical convenience we pick the 5 points at $t = -2, -1, 0, 1, 2$ and fit the straight line, Figure 14-1,

$$u(t) = a + bt$$

Least squares says that we should minimize the sum of the squares of the differences between the data and the points on the line, that is, minimize

$$M = \text{SUM}[k=-2, 2; \{u_k - (a + bk)\}^2]$$

What are the parameters to use in the differentiation to find the minimum? They are the a and the b , not the t (now the discrete variable k), and u . The line depends on the parameters a and b , and this is often a stumbling block for the student; the parameters of the equation are the variables for minimization! Hence on differentiating with respect to a and b , and equating the derivatives to zero to get the minimum, we have

$$-2 \text{SUM}[u_k - a - bk] = 0$$

$$-2 \text{SUM}[(u_k - a - bk)k] = 0$$

In this case we need only a , the value of the line at the midpoint, hence using, (some of the sums are for later use),

$$\text{SUM}[1] = 5 \qquad \text{SUM}[k^3] = 0$$

$$\text{SUM}[k] = 0 \qquad \text{SUM}[k^4] = 34$$

$$\text{SUM}[k^2] = 10$$

from the top equation we have

$$\text{SUM}[u_k] = 5a + 0b$$

$$a = (1/5)\text{SUM}[k=-2, 2; u_k]$$

which is simply the average of the five adjacent values. When you think about how to carry out the computation for a , the smoothed value, think of the data in a vertical column, Figure

14-2, with the coefficients each $1/5$, as a running weighting of the data; then you can think of it as a window through which you look at the data, with the "shape" of the window being the coefficients of the filter, this case of smoothing being uniform in size.

Had we used $2k + 1$ symmetrically placed points we would still have obtained a running average of the data points as the smoothed value that is supposed to eliminate the noise.

Suppose instead of a straight line we had smoothed by fitting a quadratic, Figure 14-3,

$$u(t) = a + bt + ct^2$$

Setting up the difference of the squares and differentiating this time with respect to a , b and c we get

$$-2 \text{ SUM}[u_k - a - bk - ck^2] = 0$$

$$-2 \text{ SUM}[(u_k - a - bk - ck^2)k] = 0$$

$$-2 \text{ SUM}[(u_k - a - bk - ck^2)k^2] = 0$$

Again we need only a . Rewriting the first and third equations (the middle one does not involve a), and inserting the known sums from above, we have

$$5a + (10)c = \text{SUM}[u_k]$$

$$(10)a + (34)c = \text{SUM}[k^2 u_k]$$

To eliminate c , which we do not need, we multiply the top equation by 17 and the lower equation by -5 , and add to get

$$(85 - 50)a = 17 \text{ SUM}[u_k] - 5 \text{ SUM}[k^2 u_k]$$

$$a = (1/35)[-3u_{-3} + 12u_{-2} + 17u_0 + 12u_1 - 3u_2]$$

and this time our "smoothing window" does not have uniform coefficients, but has some with negative values. Don't let that worry you as we were speaking of a window in a metaphorical way and hence negative transmission is possible.

If we now shift these two least squares derived smoothing formulas to their proper places about the point n we ~~would~~ have

$$u_n = (1/5)[u_{n-2} + u_{n-1} + u_n + u_{n+1} + u_{n+2}]$$

$$u_n = (1/35)[-3u_{n-2} + 12u_{n-1} + 17u_n + 12u_{n+1} - 3u_{n+2}]$$

We now ask what will come out if we put in a pure eigenfunction. We know that the equations being linear they should give the eigenfunction back, but multiplied by the eigenvalue corresponding to the eigenfunction's frequency, the transfer function value at that frequency. Taking the top of the two smooth-

ing formulas we have

$$\begin{aligned} u_n &= (1/5)[\exp\{i\omega(n-2)\} + \exp\{i\omega(n-1)\} + \dots + \exp\{i\omega(n+2)\}] \\ &= (1/5)e^{i\omega n}[e^{-2i\omega} + e^{-i\omega} + 1 + e^{i\omega} + e^{2i\omega}] \end{aligned}$$

Hence the eigenvalue at the frequency ω (the transfer function) is, by elementary trigonometry,

$$H(\omega) = (1/5)[2\cos 2\omega + 2\cos \omega + 1] = (\sin(5/2)\omega)/(5 \sin (1/2)\omega)$$

In the parabolic smoothing case we will get

$$H(\omega) = (1/35)[17 + 24 \cos \omega - 6 \cos 2\omega]$$

These are easily sketched along with the $2k + 1$ smoothing by straight line curves, Figure 14-4.

Smoothing formulas have central symmetry in their coefficients, while differentiating formulas have odd symmetry. From the obvious formula

$$f(x) = (1/2)[f(x) + f(-x)] + (1/2)[f(x) - f(-x)]$$

we see that any formula is the sum of an odd and an even function, hence any nonrecursive digital filter is the sum of a smoothing filter and a differentiating filter. When we have mastered these two special cases we have the general case in hand.

For smoothing formulas we see that the eigenvalue curve (the transfer function) is a Fourier expansion in cosines, while for the differentiation formula it will be an expansion in sines. Thus we are led, given a transfer function you want to achieve, to the problem of Fourier expansions of a given function.

Now to a brief recapitulation of Fourier series. If we assume that the arbitrary function $f(t)$ is represented

$$f(t) = a_0/2 + \text{SUM}[k=1, \infty; \{a_k \cos t + b_k \sin t\}]$$

we use the orthogonality conditions (they can be found by elementary trigonometry and simple integrations)

$$\text{INT}[-\pi, \pi; \cos kt \cos mt dt] = \begin{cases} 0 & \text{for } k \neq m, \\ \pi & \text{for } k = m \neq 0 \\ 2\pi & \text{for } k = m = 0 \end{cases}$$

$$\text{INT}[-\pi, \pi; \cos kt \sin mt dt] = 0 \text{ for all } m$$

$$\text{INT}[-\pi, \pi; \sin kt \sin mt dt] = \begin{cases} 0 & \text{for } k \neq m, \\ \pi & \text{for } k = m \neq 0 \\ 0 & \text{for } k = m = 0 \end{cases}$$

we get

$$a_k = (1/\pi) \text{INT}[f(t) \cos kt \, dt]$$

$$b_k = (1/\pi) \text{INT}[f(t) \sin kt \, dt]$$

and because we used an $a_0/2$ for the first coefficient the same formula for a_k holds for the case $k = 0$. In the complex notation it is, of course, much simpler.

Next we need to prove that the fit of any orthogonal set of functions gives the least squares fit. Let the set of orthogonal functions be $\{f_k(t)\}$ with weight function $w(t) \geq 0$. Orthogonality means

$$\text{INT}[w(t) f_k(t) f_m(t) \, dt] = 0 \quad \text{for } k \neq m, \quad 1/\lambda_k \text{ for } m = k.$$

As above the formal expansion will give the coefficients

$$c_k = \lambda_k \text{INT}[\text{range of orthogonality}; w(t) f(t) f_k(t) \, dt]$$

where the

$$\lambda_k = 1/\text{INT}[w(t) f_k^2(t) \, dt]$$

when the functions are real, and in the case of complex functions we multiply through by the complex conjugate function.

Now consider the least squares fit of a complete set of orthogonal functions using the coefficients (capitals) C_k . We have

$$\text{INT}[w(t) \{f(t) - \text{SUM}[C_k f_k(t)]\}^2 \, dt] \geq 0$$

to minimize. Differentiate with respect to C_k . You get

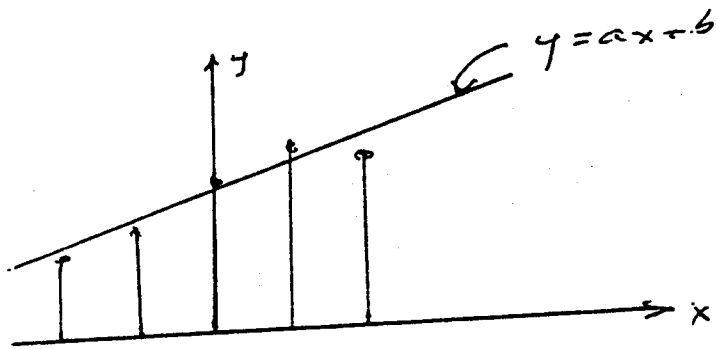
$$2 \text{INT}[w(t) \{f(t) - \text{SUM}[C_k f_k(t)]\}(-f_k(t)) \, dt] = 0$$

and we see from a rearrangement that the $C_k = c_k$. Hence all orthogonal function fits are least squares fits, regardless of the set of orthogonal functions used.

If we keep track of the inequality we find that we will have, in the general case, Bessel's inequality

$$\text{INT}[w(t) f^2(t) \, dt] - \text{SUM}[(1/\lambda_k) c_k^2] = \text{least squares error}$$

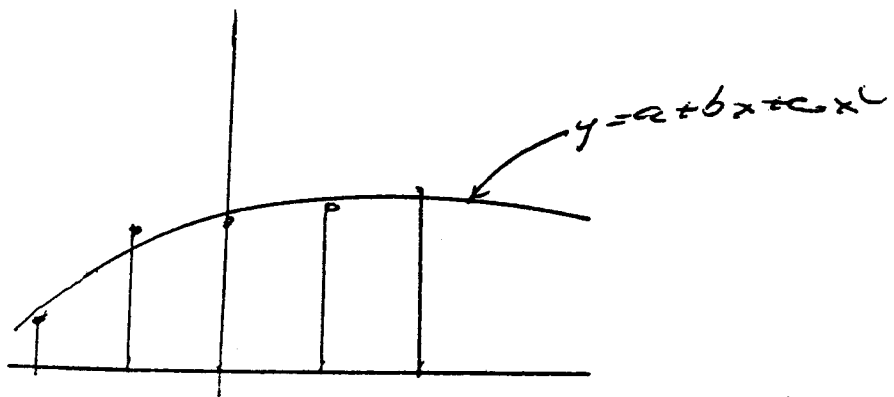
for the number of coefficients taken in the sum, and this provides a running test for when you have taken enough terms in a finite approximation. In practice this has proven to be a very useful guide to how many terms to take in a Fourier expansion.



Fitting a straight line
Figure 14-1

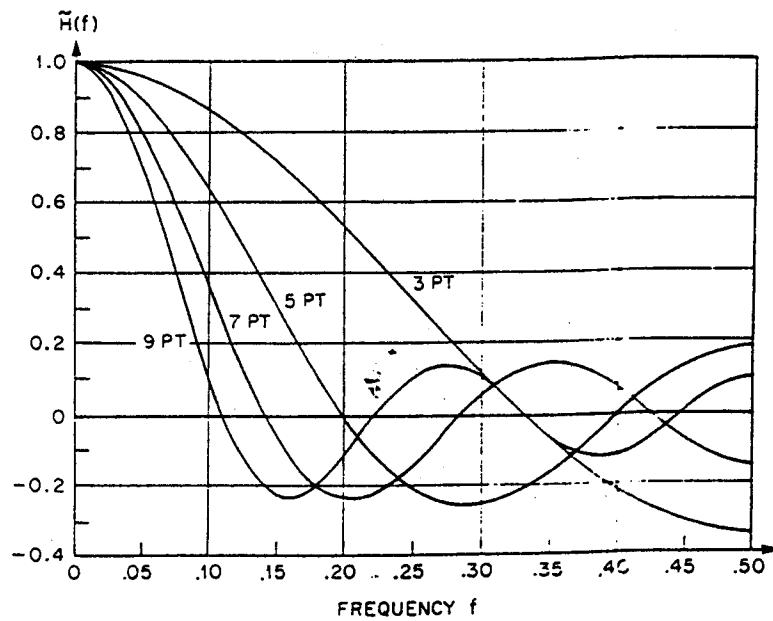
u_{n-3}			
u_{n-2}	$c_2 = 1/5$	} Filter coefficients	
u_{n-1}	$c_1 = 1/5$		
u_n	$c_0 = 1/5$		
u_{n+1}	$c_{-1} = 1/5$		
u_{n+2}	$c_{-2} = 1/5$		
u_{n+3}			
u_{n+4}			
\vdots			

Figure 14-2



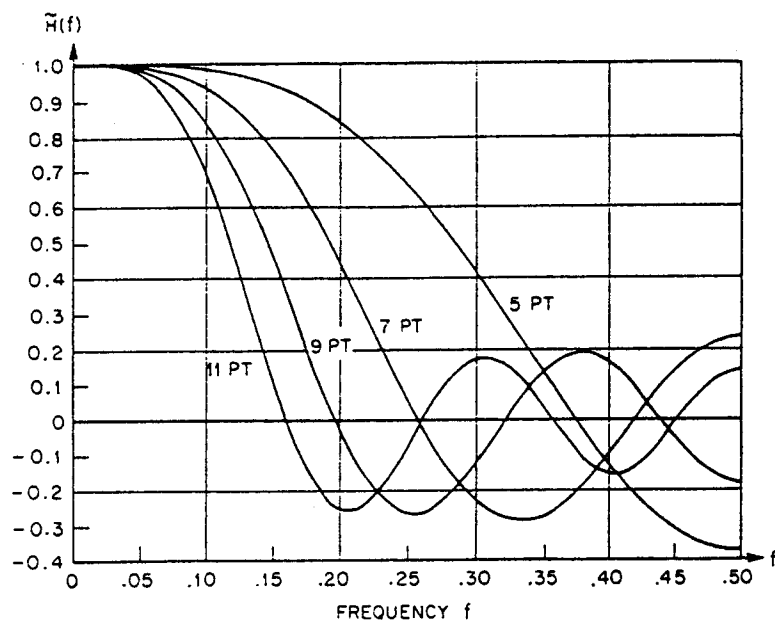
Fitting a quadratic

Figure 14-3



SMOOTHING BY LEAST SQUARES STRAIGHT LINES

Ch. 5 Some Classical Applications



TRANSFER FUNCTION FOR SMOOTHING BY LEAST-SQUARES QUADRATICS

Figure 14-15