

Information Theory and Coding
Prof. S. N. Merchant
Department of Electrical Engineering
Indian Institute of Technology, Bombay

Lecture - 40
Transform Coding - Part – I

Vector quantization is a block quantization approach. In this approach we exploit the correlation between the samples of the source output vector of a certain length k , in order to place the output points of the quantizer in appropriate position in k dimensional space. In today's class, we will study yet another block quantization approach. In this approach, we group the source samples in form of a vector of a certain length. And then use a linear reversible transform to transform this vector into components that are then quantized according to the individual characteristics. This approach is known as transform coding.

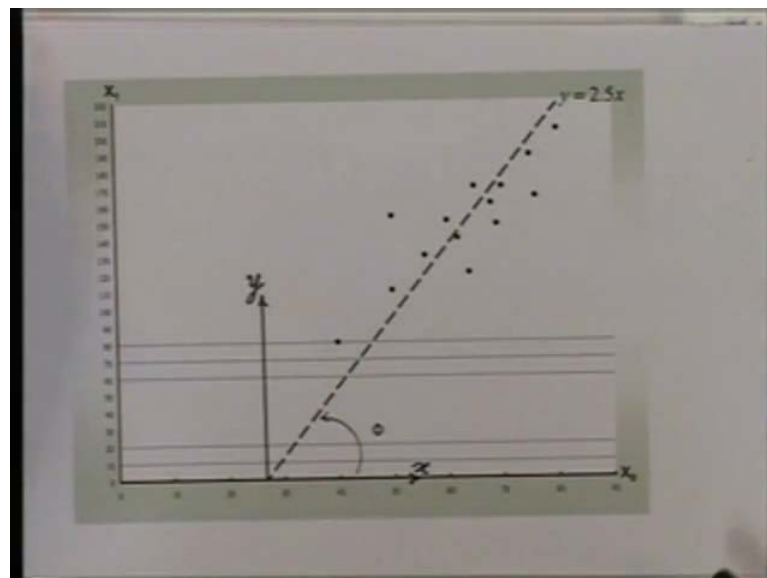
Before we examine the mathematical delicacy of transform coding. Let us try to understand the working principle of transform coding, and appreciate its advantage with the help of a couple of examples. Let us start with an example where a source generates a sequence of pairs of numbers. Now from this sequence we form source output vectors of length 2 by associating the first component of this vector with the first number of this pair and denoted by x_0 . And the second component of this vector is associated with the second number of this pair and denoted by x_1 , so a typical sequence from such a source is shown here.

(Refer Slide Time: 03:18)

Original Sequence	
x_0	x_1
64	120
62	140
56	130
50	153
69	148
70	170
76	164
68	160
60	150
50	110
75	188
40	80
65	170
80	203

Now, each of this source output vector can be represented as a dot in two dimensional space, where the horizontal axis represents x_0 and the vertical axis represents x_1 . Now, if we do this we get the following plot.

(Refer Slide Time: 03:47)

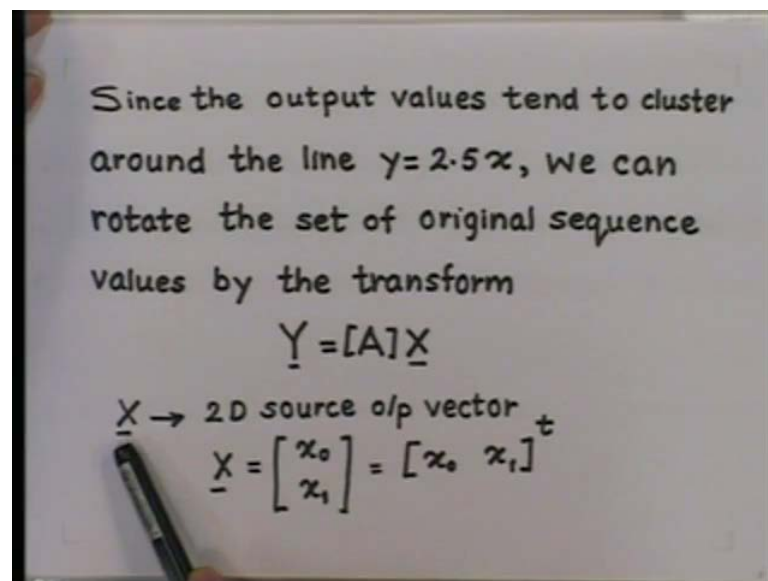


Now, let us make some observations about this plot first in the coordinate space (x, y) as shown here, where x axis align with x_0 axis then in this x, y coordinate space. We find that this dots which represent the source output vectors cluster around the line y is equal to $2.5 x$. Another observation from this plot is that horizontal coordinate is correlated

with the vertical coordinate. Another observation is that except for the different means, the distribution of the horizontal coordinate is very similar to the distribution of the vertical coordinate, and the variance of the horizontal coordinate is very much similar to the variance of the vertical coordinate. Now, as we have studied earlier that entropy at least for Gaussian source is directly proportional to the variance. So, what this plot indicates is that the entropy which is a measure of information for the horizontal coordinates is similar to the entropy of the vertical coordinate, because horizontal coordinate and vertical coordinate are correlated.

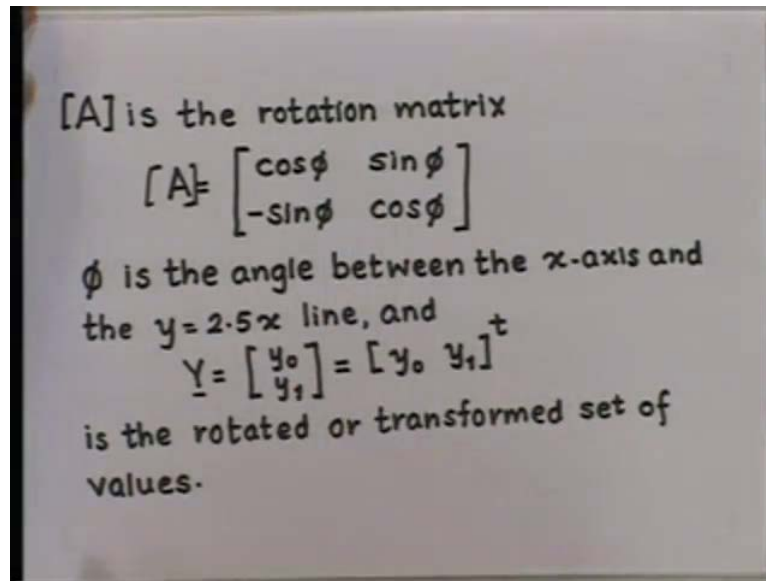
What this implies is that if we quantize the horizontal coordinate, and vertical coordinate separately then the outputs of these quantizers will have lot of common information. This will be an efficient way of transmitting the pairs of this samples, or storing this pair of the samples, so what should be done? So, one approach is that we take this coordinate space denoted by (x_0, x_1) and rotate it about $(0, 0)$ of (x, y) coordinate space by angle ϕ .

(Refer Slide Time: 07:07)



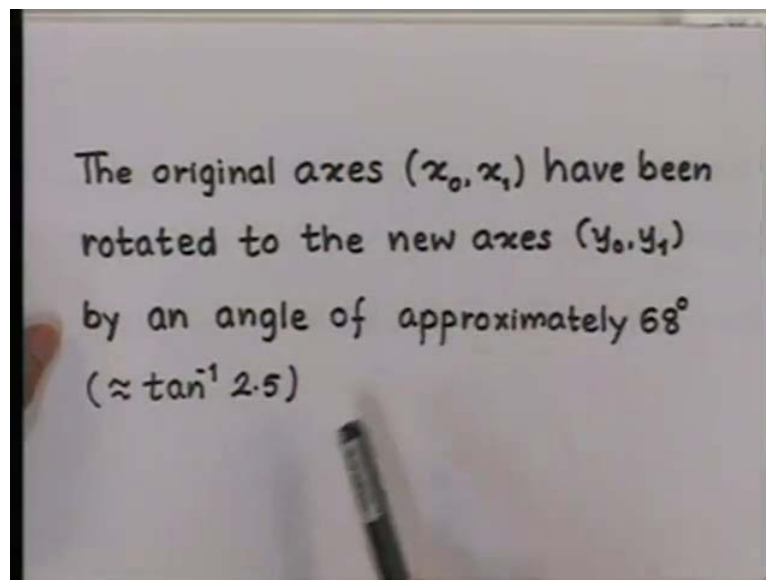
Now, if we do this then this is equivalent to the rotation of the set of original sequence value by the transform \underline{Y} is equal to $\underline{A}\underline{X}$, where \underline{X} denotes 2 D source output vector.

(Refer Slide Time: 07:23)



And A is the rotation matrix, where phi is the angle between the x axis and the y equal to 2.5 x line and y is the rotated or transformed set of values.

(Refer Slide Time: 07:58)



Now, in this case the original axis (x_0, x_1) has been rotated to the new axis (y_0, y_1) by an angle of approximately 68 degrees which is tan inverse 2.5.

(Refer Slide Time: 08:18)

For this particular case

$$[A] = \begin{bmatrix} 0.37139068 & 0.92847669 \\ -0.92847669 & 0.37139068 \end{bmatrix}$$

In this case we get the rotation matrix a given here, now once we have the rotation matrix, we can transform it and obtain the output as shown here.

(Refer Slide Time: 08:35)

Transformed sequence

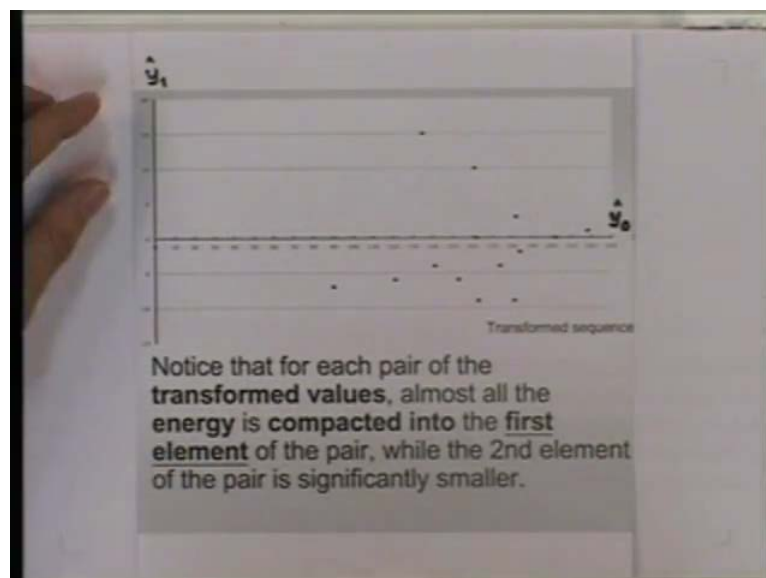
y_0	y_1	y_0	y_1
135.19	-14.86	135	-15
153.01	-5.57	153	-6
141.50	-3.71	142	-4
160.63	10.40	161	10
163.04	-9.20	163	-9
183.84	-1.86	184	-2
180.50	-9.66	181	-10
173.81	-3.71	174	-4
161.55	0.00	162	0
120.70	-5.57	121	-6
202.41	0.19	202	0
89.13	-7.43	89	-7
181.98	2.79	182	3
218.19	2.11	218	2

So, we get the output y_0 and y_1 , now the next step in transform coding is to quantize this outputs y_0 and y_1 . Let us assume that we use the same quantizer for both the components y_0 and y_1 , we use a uniform quantizer and quantize it to the nearest integer. If we do that then what we get is the following result, so y_0 cap is the quantization of the values y_0 to the nearest integer and y_1 cap is the quantization output

of y_1 to the nearest integer. Now, from this y_0 cap and y_1 cap, we can obtain back the original sequence to some approximation by using an inverse transform in the form of an inverse.

Now, if we did that then we are not gaining any real advantage from the transform coding except for the fact that earlier we would have quantized x_0 and x_1 sample, but now we have quantized y_0 and y_1 samples. So, if you want to gain the real advantage from the transform coding then in that case we should exploit the characteristics of y_0 and that of y_1 . Now, in order to study this characteristic we could plot these points in (y_0, y_1) space, but now since y_0 cap is very close to y_0 and y_1 cap is very close to y_1 , for the sake of convenience we plot the points in y_0 cap and y_1 cap space.

(Refer Slide Time: 11:17)

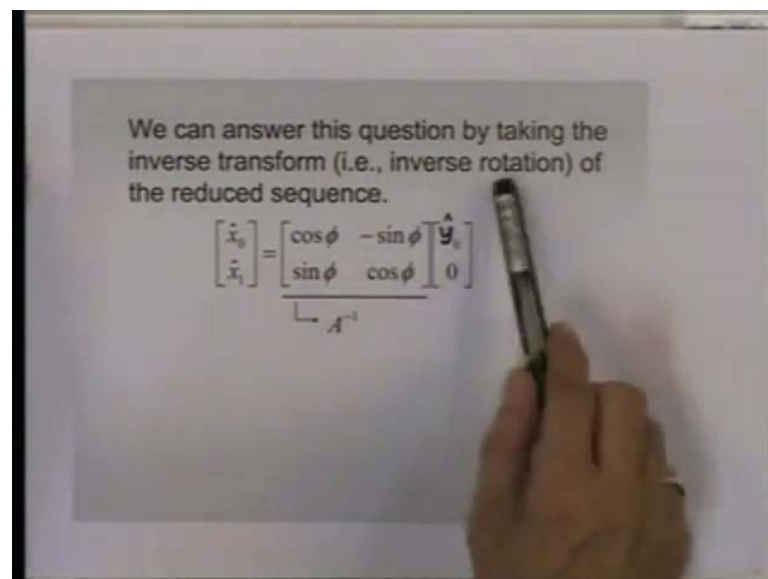


So, on doing so we get this plots in y_0 cap and y_1 cap space now the characteristic depicted by this plot very well resembles the characteristics of the points in (y_0, y_1) space. And from this plot we can observe the following if we define the energy as the sum of the squares of the horizontal coordinate plus the square of the vertical coordinate. Then this plot shows that almost all the energy is compacted into the first element of the pair indicated by y_0 , while the second element of the pair indicated by y_1 is significantly smaller. Another observation is that that now y_0 and y_1 are decode related, another observation is that that distribution for y_0 is very much different from that y_1 the variance of y_0 is much larger than the variance of y_1 .

Therefore, as discussed earlier the information contained in y_0 is much higher than the information contained in y_1 and now since y_0 and y_1 are decode related there is very less common information between the two components. Now, as said earlier that if we really want to get the advantage of transform coding then we should use this characteristic to quantize y_0 and y_1 .

So, based on this result depicted by this plot let us do the following, let us quantize y_0 as \hat{y}_0 by using a uniform quantizer which quantizes to the next nearest integer and get \hat{y}_1 as the quantizer output for y_1 which is all 0. Now, what this means that we have set all the second elements of the transformation of y_1 to 0, now by doing this number of elements that need to be coded reduces by half. Now, the question is what is the effect the throwing away half the elements of the sequence and now we can answer this question by taking the inverse transform.

(Refer Slide Time: 14:44)



That is inverse rotation of the reduced sequence, so \hat{x}_0 and \hat{x}_1 denote the reconstructed sequence or the approximated original sequence which is obtained by a inverse operating on \hat{y}_0 and 0, because \hat{y}_1 has been reduced to 0. If we do this the approximated original sequence, or the reconstructed sequence is shown out here.

(Refer Slide Time: 15:20)

reconstructed sequence	
x_0	x_1
50	125
57	142
53	131
60	150
61	151
68	171
67	168
65	162
60	150
45	112
75	188
34	84
68	169
81	203

And if we compare this reconstructed sequence with the original sequence.

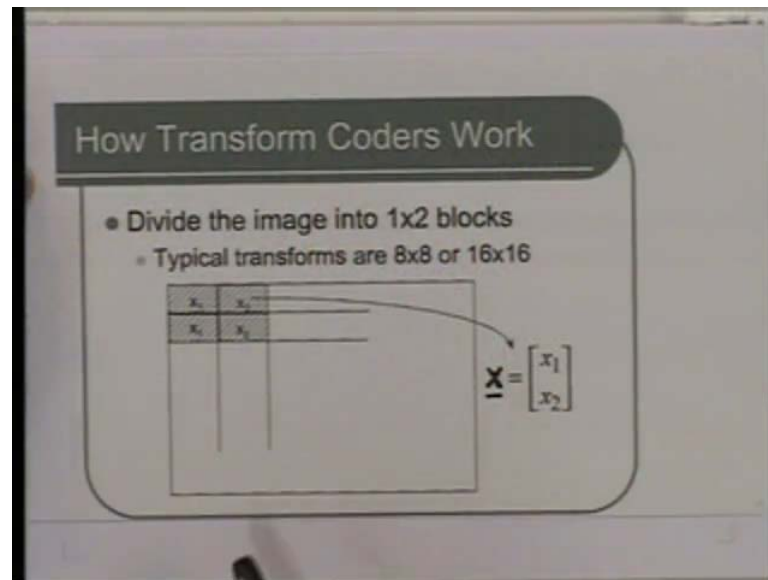
(Refer Slide Time: 15:35)

ORIGINAL + RECONSTRUCTED SEQUENCE			
x_0	x_1	x_1	x_0
64	50	125	120
62	57	142	140
56	53	131	130
50	60	150	153
69	61	151	148
70	68	171	170
76	67	168	164
68	65	162	160
60	60	150	150
50	45	112	110
75	75	188	188
40	34	84	80
65	68	169	170
80	81	203	203

This is the table which depicts that, so this is the original sequence and this is the approximated, or reconstructed original sequence after we throw away half the number of components of the transformed coefficient. Similarly, this is the reconstructed second component of the sequence and this is the original second component of the sequence. Now, comparing the reconstructed sequence with the original sequence, we see that even though we transform only half the number of elements presented in the original

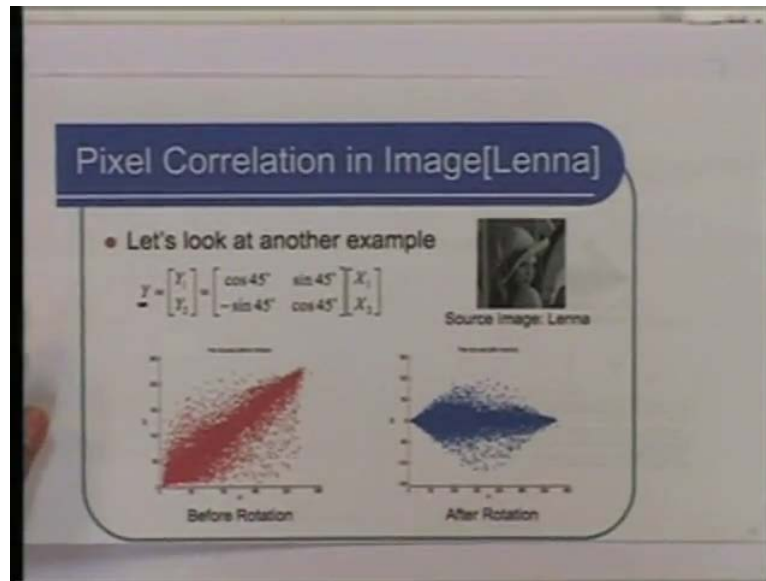
sequence, the reconstructed sequence is very close to the original sequence. Now, let us take another example from, now let us take another example where the source is an image. So, let us see how the transformed coding works on images, so the first step is to obtain the source output vector and this is done as follows.

(Refer Slide Time: 16:53)



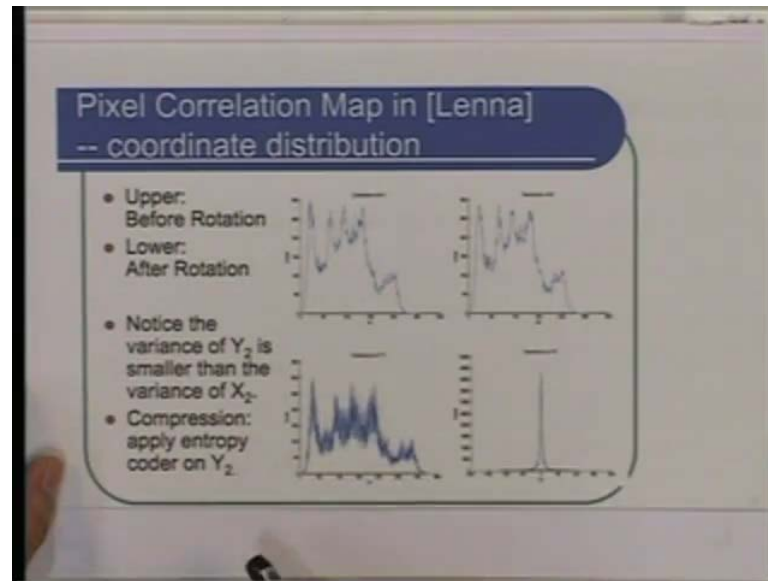
For our discussion we form source output vectors of length 2 by dividing the image into 1 by 2 block though in practice typical transforms will use the block size of 8 by 8, or 16 by 16. So, this vector x will be either of size 64 or 256 and it is obtained by scanning this blocks from left to right and from top to bottom. So, for our discussion we will assume that we have 1 by 2 blocks, so the vectors are of length 2, now using this principle.

(Refer Slide Time: 17:47)



If we take the source image as Lenna and plot the source output vectors as shown here with x_1 representing the horizontal axis and x_2 representing the vertical axis. Then this is the distribution of the source output vectors that is the distribution of adjacent pixels from this plot, it is very clear that most of this points cluster around the point x_2 equal to x_1 naught. Let us see what happens if we rotate this axis by 45 degrees by using this transform, the difference between this transform the earliest transform is that now we have 45 degrees instead of 68 degrees. And when we do this and we plot the transform components y_1 , y_2 space this is the plot which we get. Now, let us study the characteristics of this distribution before rotation and after rotation.

(Refer Slide Time: 19:14)



If we look at the distribution of x_1 before rotation, it is given here and if you look at the distribution of x_2 before rotation it is given here and so from this distribution we find that the information contained in this and in this distribution is more or less similar. And because x_1 and x_2 are correlated as depicted by this plot, it implies that there is a lot of common information between adjacent pixels, so it will not be wise for us to quantize them separately. Now, on rotation for this plot if you look at the distribution of y_1 this is what we get and if you look at the distribution of y_2 , this is what we get. So, again we note that the variance of y_2 is much smaller than the variance of x_2 .

So, if we adopt a similar strategy what we did in the earlier example and throw off all the y_2 components and reconstruct the image by using just y_1 components, we would get similar results. So, the reconstructed output image will be very close to the original image, this is because the variance of the second component is very small when we discarded it. So, in both this example what this rotation matrix A does is basically to compact the energy into the first component.

(Refer Slide Time: 21:30)

$$[A] = \frac{1}{\sqrt{2}} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix} \quad \underline{x} = \begin{bmatrix} 4 \\ 5 \end{bmatrix} \quad 16 + 25 = 41$$
$$\underline{y} = [A] \underline{x} = \begin{bmatrix} 6.364 \\ 0.707 \end{bmatrix} \quad (6.364)^2 + (0.707)^2 = 41$$
$$\min\{\underline{x}\} \rightarrow \frac{4^2}{41} = 0.39$$
$$\min\{\underline{y}\} \rightarrow \frac{(0.707)^2}{41} = 0.012$$

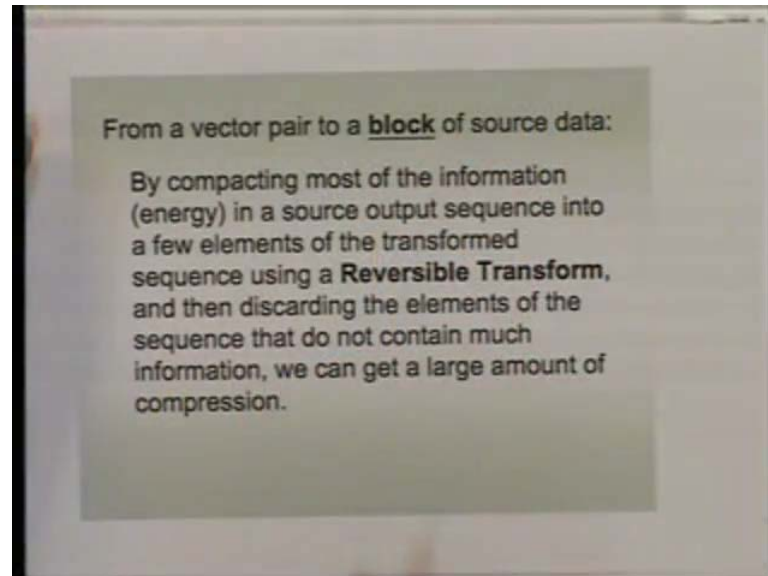
So, for example, if we choose the transform A equal to $\frac{1}{\sqrt{2}}$, 1, 1, minus 1, 1 and have the input vector x is equal to 4, 5 to obtain the transformed vector. As now what this shows is that the energy in the x which is equal to 16 plus 25 is equal to 41 and energy in y which is 6.364 squared plus 0.07 squared is equal to again 41. So, energies are same, we will very shortly prove that this is because of the kind of transformation which we have chosen here, now what is important is that this transformation makes the second component very small.

So, in the original vector if we discarded minimum component of the vector x, we would have an error of 4 squared by 41 is equal to 0.39, but now in the transformed space if we discard the minimum component of the vector y, we have an error of 0.07 squared divided by 41 which is equal to 0.012. So, the conclusion is that we are more confident to discard the minimum of the transformed vector than the minimum component of the original vector.

Now, in both the examples we could reduce the number of samples we needed to code because most of the information or energy containing a pair of values was put into one element of each pair because the other element of the pair contained very little information. We could discard it without significant impact on the fidelity of the reconstructed sequence. Now, in both this example we formed a vector of length 2. Now,

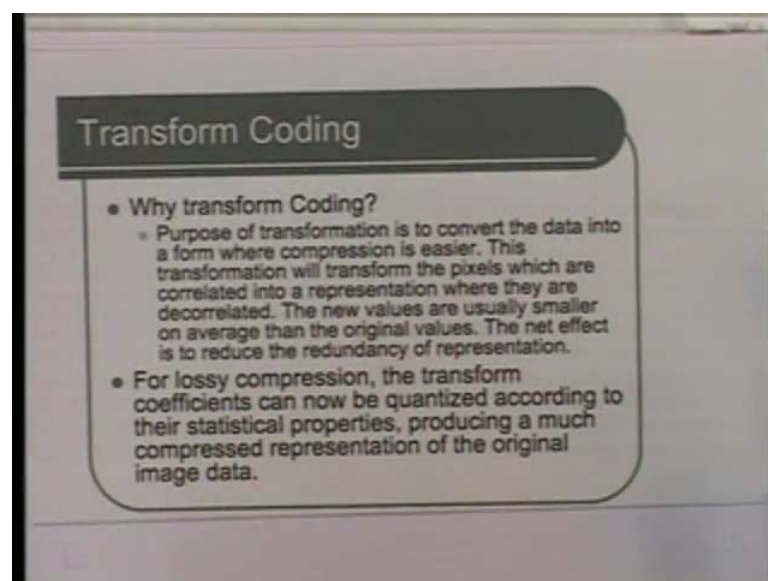
if we exchange this idea of forming vectors of larger length by incorporating larger block of source samples then we would get the following result by compacting.

(Refer Slide Time: 25:09)



By compacting most of the information or analogy in a source output sequence into a few elements of the transformed sequence using a reversible transform. And the discarding the elements of the sequence that do not contain much information, we can get a large amount of compression, so this is the idea behind transformed coding.

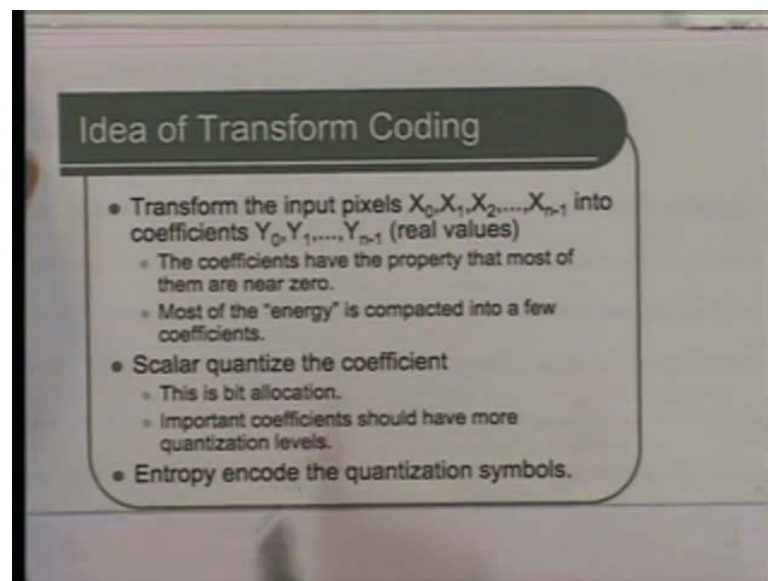
(Refer Slide Time: 25:56)



So, if we ask why use transform coding, then we can answer as follows the purpose of transformation is to convert the data into a form where compression is easier. This transformation will transform the pixels which are correlated into a relation where, they are decorrelated. So, this is very important aspect of transformed coding, the correlation existing among the samples in the source output. Vector on transformation they get decorrelated and because they get decorrelated, we can independently quantize this. So, the new values are generally, smaller than average of the new values and the net effect is to reduce the redundancy of representation because now, we have minimum amount of common information between the transformed cooperation's.

And usually, transformed coding is used for lossy compression where the transformed coefficients can now be quantized according to the statistical properties, which differ from coefficient to coefficient and in the effect producing a much compressed representation of the original image data. So, if we were to put in a nut shell the idea of transformed coding is as follows.

(Refer Slide Time: 27:41)

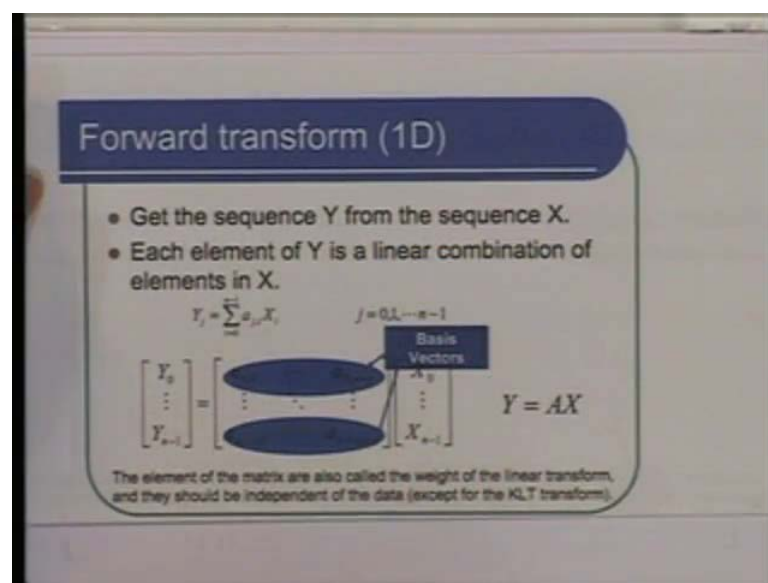


Transform the input pixels x_0, x_1, x_2 up to x_{n-1} , this implies that we have made a block of size n into coefficients y_0, y_1 up to y_{n-1} without loss of generality, we will assume to be real values. Now, if we have chosen the transform to be proper then the net effect would be that the coefficients have the property that most of them are near

0. In the earlier example we saw that the second component was very near to 0, so what it means that most of the energy is compacted into a few coefficients.

Next step in transformed coding is to scalar quantize the coefficients that is do the bit allocation, that is important coefficients should have more quantization levels that is how we will carry out the transformation. And finally, the entropy encode quantization symbols and transmit it on storage at the receiver, we apply the inverse transform. So, mathematically this process of transformation can be depicted as shown here.

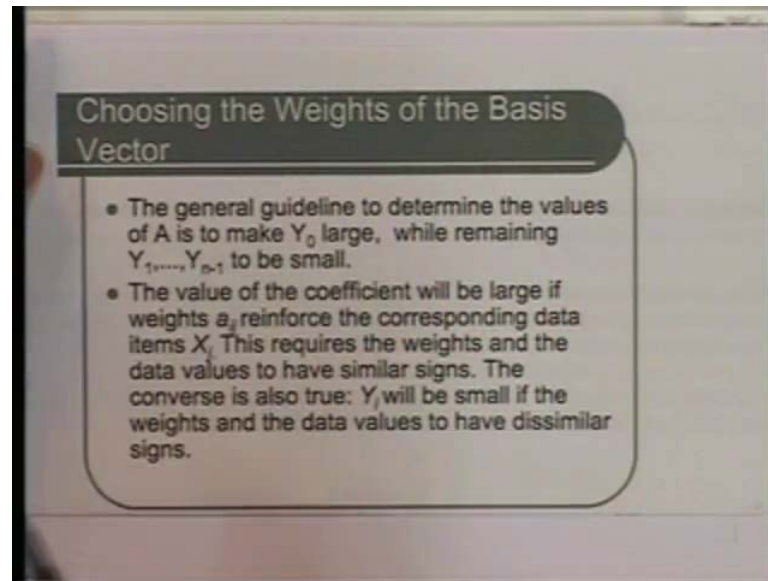
(Refer Slide Time: 29:20)



We get the sequence y from the sequence x, each element of y is a linear combination of elements in x. So, we see that the transformed values are actually the coefficients of an expansion of the input sequence in terms of the rows of the transformed matrix, so that is why the rows of this transformed matrix are often referred to as the basis vectors of the transform. The elements of the transformed sequence are often called the transform coefficients.

Now, different transforms have different basis vectors, sometimes the elements of the matrix are also called the weight of the linear transform. In practical application we want to be independent of the data except in the case of Karhunen-Loeve transform, which we will discuss very shortly. Now, the next question is how do you choose the weights of this basis vectors?

(Refer Slide Time: 30:57)



So, the general guideline to determine the values of A is to make y_0 that is then first component of the transform vector to be as large as possible while remaining y_1 up to y_{n-1} to be as small as possible. So, what this implies that the value of the coefficients will be large if weights A_j reinforce the corresponding data items X_j . So, if these weights reinforce the data items then this coefficient will be high, now this requires the weights and the data values to have similar signs. The converse is also true, Y_j will be small if the weights and the data values have dissimilar signs.

Thus, the basis vector should extract distinct features of the data vectors and must be independent and we also want it to be orthogonal, ortho-normal we will see the reason of this very shortly. Now, what this basis vectors do? Is that they pick up the low and high frequency components of this data and normally, the coefficient decrease in the order y_0 y_1 up to y_{n-1} and because of this y is now more amenable to compression than x . So, let us look at the choice of this transform in both the examples which we considered.

(Refer Slide Time: 33:28)

Orthogonal/orthonormal Matrix

- Rotation matrix is orthogonal.
 - The dot product of a row with itself is **nonzero**.
 - The dot product of different rows is 0.
$$A_i \cdot A_j = \begin{cases} \neq 0 & \text{if } i = j \\ 0 & \text{else} \end{cases}$$
- Furthermore, the rotation matrix is orthonormal.
 - The dot product of a row with itself is **1**.
$$A_i \cdot A_j = \begin{cases} 1 & \text{if } i = j \\ 0 & \text{else} \end{cases}$$
- Example: $A = \frac{\sqrt{2}}{2} \begin{bmatrix} 1 & 1 \\ -1 & 1 \end{bmatrix}$

We find that the rotation matrix is orthogonal, what this means is that the dot product of a row with itself is non zero, that is $A_i \cdot A_j$ is non zero if i is equal to j and is equal to 0 otherwise, the dot product of other rows is 0. Furthermore we want that the rotation matrix is orthonormal, what this implies is that the dot product of a row with itself is 1 for example the rotation matrix which we considered. For the second example where the source was Lenna image which is an orthonormal matrix and now, another consideration to choose rotation matrix is to conserve energy. Why this is very important can be explained with the help of an example.

(Refer Slide Time: 34:39)

Energy Preserving (1D)

- Another consideration to choose rotation matrix is to **conserve energy**.
- For example, we have orthogonal matrix

$$A = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & -1 & -1 \\ 1 & -1 & -1 & 1 \\ 1 & -1 & 1 & -1 \end{bmatrix} \quad X = \begin{bmatrix} 4 \\ 6 \\ 5 \\ 2 \end{bmatrix} \quad Y = AX = \begin{bmatrix} 17 \\ 3 \\ -5 \\ -2 \end{bmatrix}$$

- Energy before rotation: $4^2 + 6^2 + 5^2 + 2^2 = 81$
- Energy after rotation: $17^2 + 3^2 + (-5)^2 + (-2)^2 = 324$
- Energy changed!
- Solution: scale W by scale factor. The scaling does not change the fact that most of the energy is concentrated at the low frequency components.

Let us choose an orthogonal matrix a given here and let us assume that we have the source output vector given by this. Now, the transformation A will give the transformed coefficients, coefficients as shown below. Now, if you look at the energy before rotation that is squares of this and addition will give you 81. Whereas, energy after rotation that is the sum of the squares of this elements will give you 324.

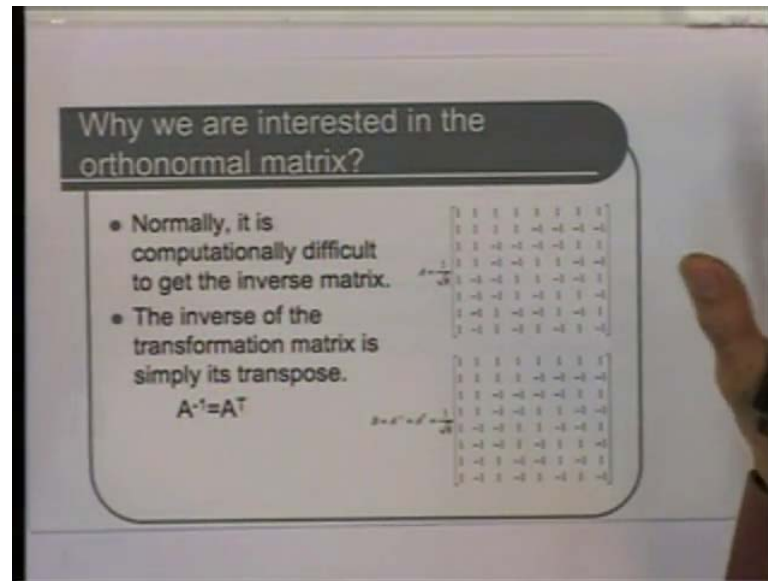
So, we do not want such a thing to happen we want that energy after transformation is also conserved, so this can be very easily done by scaling the orthogonal matrix by a scale factor. Now, the scaling does not change the fact that most of the energy is concentrated at the low frequency components, now we can formally prove this energy conservation as follows.

(Refer Slide Time: 36:02)

$$\begin{aligned}
 \text{Energy} &= \sum_{i=0}^{n-1} y_i^2 \\
 &= \underline{Y}^T \underline{Y} \\
 &= \left([\underline{A}] \underline{X} \right)^T [\underline{A}] \underline{X} \\
 &= \underline{X}^T \underbrace{[\underline{A}]^T [\underline{A}]}_{[\underline{I}]} \underline{X} \\
 &= \underline{X}^T \underline{X} \\
 &= \sum_{i=0}^{n-1} x_i^2
 \end{aligned}$$

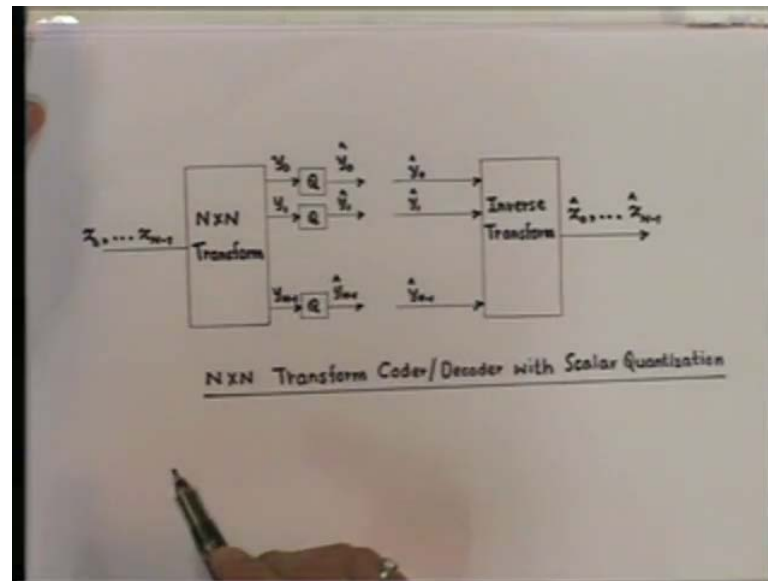
Energy in the transformed coefficients is given by y_i^2 i equal to 0 to n minus 1, which can be written in the form of matrix as dot product of y itself and this can be rewritten as which can be simplified as follows. Now, if we choose to be ortho-normal matrix then this is equal to $x^T x$. Dot product of x with itself which is equal to the energy in the source output vector, so we choose the linear transform to be orthonormal matrix because it preserves the energy. Now, another reason for choosing this orthonormal matrix is explained in this slide.

(Refer Slide Time: 37:30)



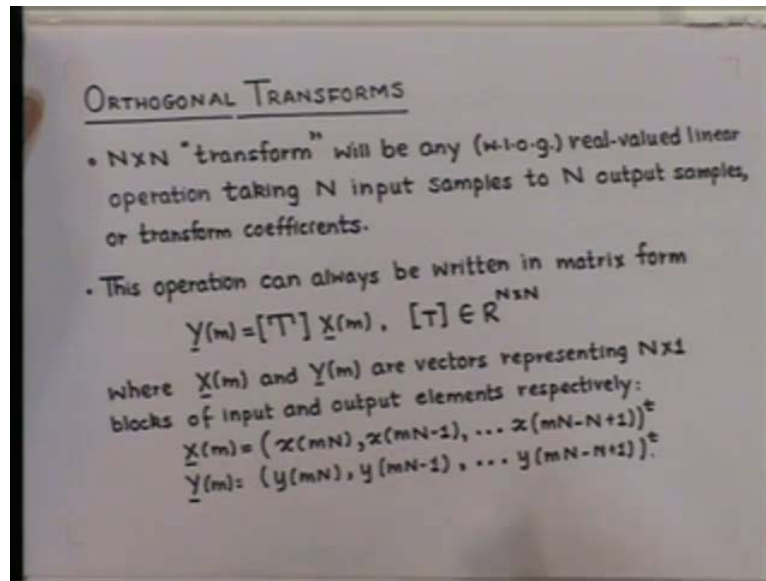
So, if we use a linear transform the in principle its inverse can be obtained, but it is computationally difficult to get the inverse matrix, but if this transform matrix is orthonormal then the inverse of the transposition matrix is simply its transpose. Therefore, this has an advantage from this implementation point of view because the decoder is also a linear transform and mathematically it can be shown. That orthogonal transform will provide the least minimum squares error as far as the reconstruction is concerned; we throw off some of the transform coefficients. Now, we will not do this, but we will assume this fact and go ahead, so whatever we have discussed so far can be summarized in this block diagram here.

(Refer Slide Time: 38:47)



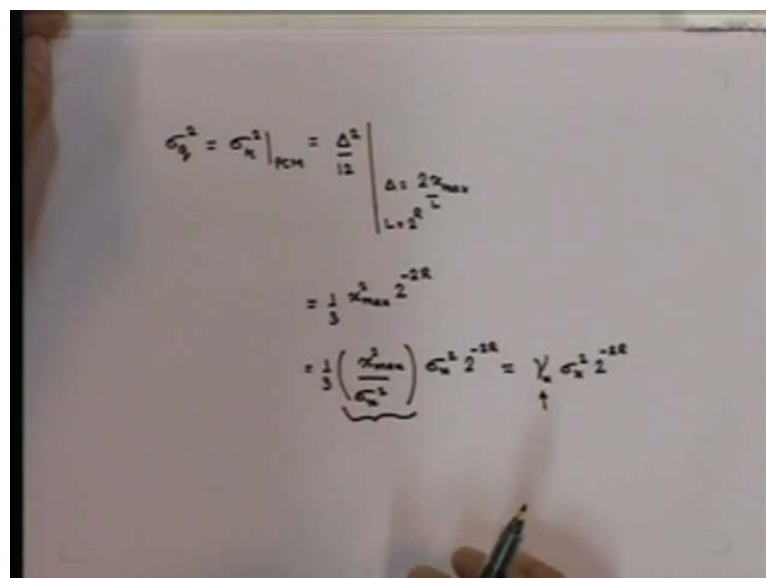
So, in transformed coding blocks of n input samples are transformed to n transformed coefficients which are then quantized individually and they are either transmitted or stored at the decoder. An inverse transform is applied to the quantized coefficients yielding a reconstruction of the original sequence. Now, by designing individual quantizers in accordance with the statistics of their inputs, it is possible to allocate bits in a more optimal manner, example encoding the more important coefficients at a higher bit rate. Now, for our discussion n by n transform will be any without loss of generality real values linear operation taking n input samples to n output samples or transform coefficients this operation can be represented in matrix form as follows.

(Refer Slide Time: 40:15)



Where x m vector and y m vectors represent n by 1 blocks of input and output elements respectively, where x m vector is related to the input sequence as follows and the y m vector is related to the output vector as follows. Now, before we go ahead with the discussion of transform coding, let us recall the quantization error which occurs in pulse code modulation. Pulse code modulation is nothing but discretization of analog wave form in time or space and amplitude.

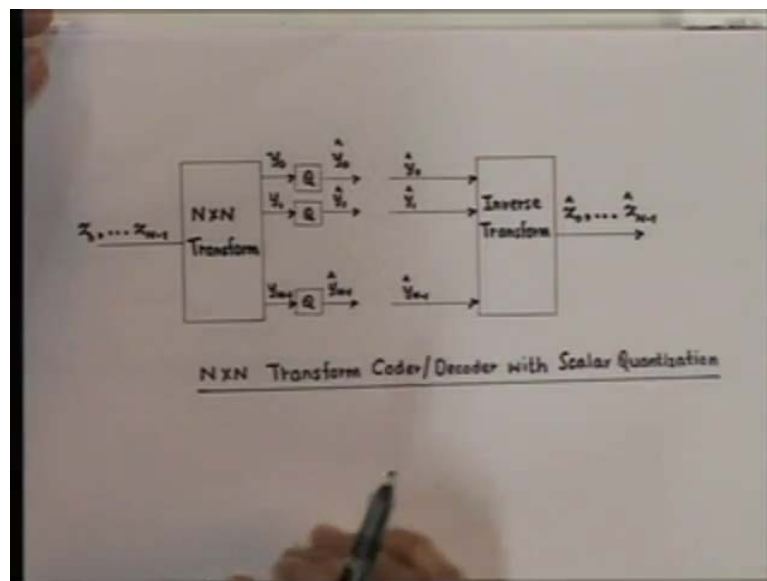
(Refer Slide Time: 41:33)



So, for a given bit rate of r bits per sample uniformly quantized P C M implies a mean squares quantization error of σ_q^2 . This is equivalent to reconstruction error in P C M is equal to $\frac{\Delta^2}{12}$, where Δ is obtained by $2 \times \text{max}$ divided by L where L is equal to 2^R . This we have seen earlier, so this can be simplified as $\frac{1}{3} \times \text{max}^2 \times 2^{-2R}$, and this can be rewritten as $\sigma_x^2 \times \text{max}^2 \times 2^{-2R}$, where σ_x^2 is the variance of the input.

Now, this quantity depends on the distribution of x , so for a uniform quantizer for a given distribution, we can write the quantization error or reconstruction error as $\gamma \times \text{max}^2 \times 2^{-2R}$, where γ is a constant that depends on the distribution of x . We will utilize this information to find out the optimum bit allocation strategy in transform coding.

(Refer Slide Time: 43:47)



Now, another question that arises using this strategy what is the mean squared reconstruction error? That means error between the reconstructed sequence, and the original sequence.

(Refer Slide Time: 44:04)

$$\begin{aligned}
 \sigma_x^2 &= \frac{1}{N} \sum_{k=0}^{N-1} E \left\{ \underbrace{(\hat{x}(Nm-k) - x(Nm-k))^2}_{\text{}} \right\} \\
 &= \frac{1}{N} E \left\{ \|\hat{\underline{x}}(m) - \underline{x}(m)\|^2 \right\} \\
 &= \frac{1}{N} E \left\{ \|\tau^{-1} \hat{\underline{y}}(m) - \underline{x}(m)\|^2 \right\} \\
 &= \frac{1}{N} E \left\{ \|\tau^{-1} (\underline{y}(m) + \underline{q}(m)) - \underline{x}(m)\|^2 \right\} \\
 &= \frac{1}{N} E \left\{ \|\underbrace{\tau^{-1} \tau}_{[I]} \underline{y}(m) + \tau^{-1} \underline{q}(m) - \underline{x}(m)\|^2 \right\} \\
 &= \frac{1}{N} E \left\{ \|\tau^{-1} \underline{q}(m)\|^2 \right\}
 \end{aligned}$$

So, this can be written as mean squared reconstruction error is equal to 1 by N summation of k equal to 0 to N minus 1. Expectation of this is the mean squared error for a particular component and since there are n components, we get by definition this expression. Now, this can be simplified as 1 by N expectation of the Euclidian distance between the original vector and the reconstructed vector.

This can be rewritten as the quantized transform coefficient vector, which is equal to the original transform coefficient vector plus quantization error associated with each component. So, this is another vector minus x m and this can be further simplified as because y m is equal to transform times X m plus. Now, for orthonormal transform this quantity is equal to I and therefore, this reduces to 1 by N expectation of the inverse transform of the quantization error which can be rewritten as...

(Refer Slide Time: 48:13)

$$\begin{aligned}\sigma_x^2 &= \frac{1}{N} E \left\{ \underline{q}^t(m) \underbrace{([\tau]^{-1})^t [\tau]^{-1}}_{[I]} \underline{q}(m) \right\} \\ &= \frac{1}{N} E \left\{ \|\underline{q}(m)\|^2 \right\} \\ &= \frac{1}{N} \sum_{k=0}^{N-1} \sigma_{q_k}^2 \\ \text{MSE-reconstruction} &= \text{MSE-quantization} \quad [\tau]\end{aligned}$$

Now, because T is an ortho-normal transform this is again equal to identity matrix, so this is equal to one time and equal to. So, what this result indicates is that mean squared reconstruction error will be denoted as MSE reconstruction is equal to mean squared quantization error which will be indicated as MSE quantization. This is true when you have T as an ortho-normal transformation matrix, now in the example which we discussed earlier for lossy transformed coding application.

We threw off all the second components of the transform vector and quantized the first components of the transform vector. Now, the natural question is that assuming that transform is an N by N ortho-normal matrix. What is the mean squared error optimal bit rate allocation strategy? Assuming independent uniform quantization of the transform output. In other words we want to find out the rates for the different quantizer such that it minimize the average reconstruction error for the fixed given average rate. The next question is how do we choose this transform itself so that the mean squared error reconstruction is minimized. We will try to answer these questions in the next class.