

**Step 1(a):** Given  $k$ ,  $T$  and  $R_1, \dots, R_k$  such that  $\frac{1}{k} \sum_{i=1}^k R_i = R$ , choose  $Q_1, \dots, Q_k$  to have rates  $R_1, \dots, R_k$ , respectively, and to minimize distortion.

Since

$$D = \frac{1}{k} \sum_{i=1}^k D_{U_i, Q_i},$$

for each  $i$  we choose  $Q_i$  to minimize  $D_{U_i, Q_i}$  subject to  $R(Q_i) \leq R_i$ . Then

$$D_{U_i, Q_i} = \delta_{sq, U_i}(R_i) \triangleq \text{OPTA for scalar quantization of } U_i.$$

The OPTA for  $U_j$  is either the fixed-rate or the variable-rate OPTA, depending on whether we want to optimize transform coding for fixed or variable rate coding.

It follows that the least distortion attainable for the given  $k$ ,  $T$ , and  $R_1, \dots, R_k$  is

$$D = \frac{1}{k} \sum_{i=1}^k \delta_{sq, U_i}(R_i)$$

Tr-13

### Step 1(b): Bit/rate allocation

We now seek to find the  $R_1, \dots, R_k$  that minimizes

$$\frac{1}{k} \sum_{i=1}^k \delta_{sq, U_i}(R_i)$$

subject to the constraint

$$\frac{1}{k} \sum_{i=1}^k R_i \leq R.$$

This is called a bit or rate allocation, i.e. an allocation of  $kR$  bits among the scalar quantizers for the individual coefficients. If some are given rate larger than  $R$ , then others must be given rate less than  $R$ .

Recall the original idea of transform coding, namely, that some coefficients will be smaller on the average than others, and that these will be encoded at lower rates than those that are larger on the average. If there is to be a net gain, then it will depend on a successful bit allocation.

For fixed-rate coding, the values of  $\delta_{sq, U_i}(R_i)$  can be tabulated for the various values of  $R_i$  and a numerical optimization algorithm can be applied to find the best choice of the  $R_i$ 's. See Gersho and Gray, pp. 226-235. Similar methods may be used for variable-rate coding.

We now consider the high-resolution case. In this case, one may solve analytically for the best bit allocation.

Tr-14

## HIGH-RESOLUTION ANALYSIS

To simplify discussion assume  $X$  comes from zero-mean, stationary source. Assume  $R$  is so large that the optimal  $R_i$ 's are so large that we may use the approximations

$$\delta_{sq,U_i}(R_i) \cong \frac{1}{12} \sigma_i^2 \alpha_i 2^{-2R_i},$$

where  $\sigma_i^2$  = variance of  $U_i$ , and for fixed-rate coding

$$\alpha_i = \beta_{U_i} = \frac{1}{\sigma_i^2} \left( \int f_{U_i}^{1/3}(u) du \right)^3$$

and for variable-rate coding

$$\alpha_i = \frac{1}{\sigma_i^2} 2^{2h(U_i)}, \text{ where } h(U_i) = - \int_{-\infty}^{\infty} f_{U_i}(u) \log_2 f_{U_i}(u) du$$

We now have

$$D \cong \frac{1}{k} \sum_{i=1}^k \frac{1}{12} \sigma_i^2 \alpha_i 2^{-2R_i}$$

And the least dist'n mfor dimen.  $k$ , transform  $T$ , and rate  $R$  or less is approx'y

$$\min_{R_1 \geq 0, \dots, R_k \geq 0: \frac{1}{k} \sum_{i=1}^k R_i \leq R} \frac{1}{k} \sum_{i=1}^k \frac{1}{12} \sigma_i^2 \alpha_i 2^{-2R_i}$$

Tr-15

## A USEFUL OPTIMIZATION RESULT

**Lemma:** Let  $g_1(R), \dots, g_k(R)$  be positive valued, strictly decreasing, continuous functions defined on  $[0, \infty)$ . If  $R_1, \dots, R_k$  minimize

$$\sum_{i=1}^k g_i(R_i) \text{ subject to } \frac{1}{k} \sum_{i=1}^k R_i \leq R.$$

are zero, then

$$g'_i(R_i) = g'_j(R_j) \text{ for all } i, j \text{ s.t. } R_i > 0, R_j > 0$$

and

$$|g'_i(R_i)| \leq |g'_j(R_j)| \text{ for all } i, j \text{ s.t. } R_i = 0, R_j > 0$$

That is, for the optimum choice of  $R_1, \dots, R_k$ , the slopes of the  $g_i$  functions are the same.

Notes:

The slopes are negative.

The above is the basis of the Lagrange multiplier method.

Tr-16

### Proof by Contradiction:

Suppose  $R_1, \dots, R_k$  minimize  $\frac{1}{k} \sum_{i=1}^k g_i(R_i)$  with  $\frac{1}{k} \sum_{i=1}^k R_i \leq R$ ,

and for some  $i, j$ ,  $R_i > 0$ ,  $R_j > 0$ , and

$$g'_i(R_i) < g'_j(R_j) < 0$$

Then for some small  $\varepsilon$ , replace

$$R_i \text{ by } R_i + \varepsilon \text{ and } R_j \text{ by } R_j - \varepsilon.$$

This has no effect on  $\frac{1}{k} \sum_{i=1}^k R_i$ .

It decreases  $g_i(R_i)$  to  $g_i(R_i + \varepsilon) \cong g_i(R_i) + \varepsilon g'_i(R_i)$ .

It increases  $g_j(R_j)$  to approximately  $g_j(R_j - \varepsilon) \cong g_j(R_j) - \varepsilon g'_j(R_j)$ .

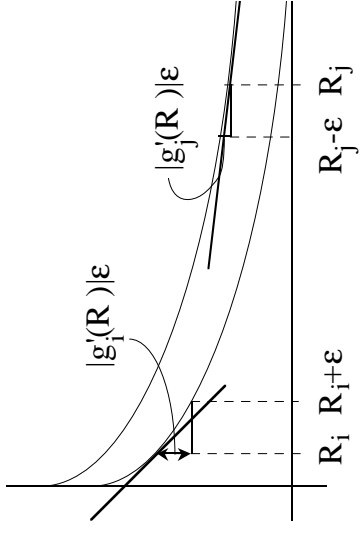
Since  $g'_i(R_i) < g'_j(R_j)$ , the new rate allocation decreases overall distortion:

$$g_i(R_i + \varepsilon) + g_j(R_j - \varepsilon) \cong g_i(R_i) + \varepsilon(g'_i(R_i) - g'_j(R_j)) < g_i(R_i) + g_j(R_j).$$

(Despite the " $\cong$ ", the overall inequality will be strict if  $\varepsilon$  is sufficiently small.) This contradicts the original assumption that slopes are not the same. Thus, they must be the same.

Similar argument, if  $R_i = 0$ ,  $R_j > 0$ , and  $|g'_i(R_i)| > |g'_j(R_j)|$

Tr-17



### Optimal Rate Allocation in the High-Resolution Case

This lemma shows that if  $R_1, \dots, R_k > 0$  minimize

$$\frac{1}{k} \sum_{i=1}^k \frac{1}{12} \sigma_i^2 \alpha_i 2^{-2R_i} \quad \text{subject to} \quad \frac{1}{k} \sum_{i=1}^k R_i = R,$$

then there is a constant  $c$  such that for each  $i$ ,

$$c = \frac{d}{dR_i} \left( \frac{1}{12} \sigma_i^2 \alpha_i 2^{-2R_i} \right) = \frac{1}{12} \sigma_i^2 \alpha_i 2^{-2R_i} (-2 \ln 2)$$

Solving the above yields

$$R_i = \frac{1}{2} \log_2 \sigma_i^2 \alpha_i - \frac{1}{2} \log_2 \left( \frac{-6c}{\sqrt{n} 2} \right).$$

Since the  $R_i$ 's average to  $R$ , equating the avg. of the RHS terms to  $R$  gives

$$-\frac{1}{2} \log_2 \left( \frac{-6c}{\sqrt{n} 2} \right) = R - \frac{1}{k} \sum_{j=1}^k \frac{1}{2} \log_2 \sigma_j^2 \alpha_j = R - \frac{1}{2} \log_2 \left( \prod_{j=1}^k \sigma_j^2 \alpha_j \right)^{1/k}$$

Therefore the optimal rate allocation is

$$R_i = \frac{1}{2} \log_2 \sigma_i^2 \alpha_i + R - \frac{1}{2} \log_2 \left( \prod_{j=1}^k \sigma_j^2 \alpha_j \right)^{1/k} = R + \frac{1}{2} \log_2 \frac{\sigma_i^2 \alpha_i}{\left( \prod_{j=1}^k \sigma_j^2 \alpha_j \right)^{1/k}}$$

Substituting the optimal allocation shown above into the expression for distortion yields

$$D \cong \frac{1}{12} \left( \prod_{j=1}^k \sigma_j^2 \alpha_j \right)^{1/k} 2^{-2R}$$

Tr-18

Summary: Given a k-dimensional transform  $T$  and large  $R$ , the optimal rate allocation is

$$R_i = R + \frac{1}{2} \log_2 \frac{\sigma_i^2 \alpha_i}{\left( \prod_{j=1}^k \sigma_j^2 \alpha_j \right)^{1/k}} \quad (*)$$

and the minimal distortion is

$$D \cong \frac{1}{12} \left( \prod_{j=1}^k \sigma_j^2 \alpha_j \right)^{1/k} 2^{-2R} \quad (**)$$

**Notes:**

- Compare (\*\*) to the high-rate approx to the OPTA of scalar quantization

$$\delta_{sq,x}(R) \cong \frac{1}{12} \sigma_X^2 \alpha_X^2 2^{-2R}$$

The SNR gain of transform coding with a particular transform over scalar quantization is

$$10 \log_{10} \frac{\delta_{sq,x}(R)}{D_{transform}} \cong 10 \log_{10} \frac{\sigma_X^2 \alpha_X}{\left( \prod_{j=1}^k \sigma_j^2 \alpha_j \right)^{1/k}}$$

- A good transform is one that makes  $\prod_{j=1}^k \sigma_j^2 \alpha_j$  small.
- Coefficient correlations do not enter into formula for  $D$ .
- We won't compare to k-dimensional VQ until we optimize  $T$ .

- If  $X$  is Gaussian, then each  $U_j$  is Gaussian, and

$$\alpha_1 = \dots = \alpha_k = \alpha_X = 32.6 \text{ for FRC and } 17.08 \text{ for VRC.}$$

Therefore, gain of transform coding with transform  $T$  over scalar quantization is

$$10 \log_{10} \frac{\sigma_X^2}{\left( \prod_{j=1}^k \sigma_j^2 \right)^{1/k}} = 10 \log \frac{\text{ratio of variance to geometric mean of coef variances}}$$

Fact: geometric mean  $\leq$  arithmetic mean; equality iff "averages" are identical

$$\Rightarrow \left( \prod_{j=1}^k \sigma_j^2 \right)^{1/k} \leq \frac{1}{k} \sum_{j=1}^k \sigma_j^2 = \frac{1}{k} E\|U\|^2 = \frac{1}{k} E\|X\|^2 = \sigma_X^2, \text{ equality iff } \sigma_j^2 = \sigma_X^2 \text{ all } j$$

Consequences: For Gaussian case,

- Gain in dB  $\geq 0$ .
- A good transform does "energy compaction" i.e. it reduces

$$\Gamma \triangleq \left( \prod_{j=1}^k \sigma_j^2 \right)^{1/k},$$

by making some  $\sigma_j^2$ 's as small as is possible with orthogonal transform.

Since transform is orthogonal, some  $\sigma_j^2$ 's must be large.

Summary:  $\sigma_j^2$ 's should be as different as possible.

- If  $X$  is not Gaussian, there's no simple relation among the  $\alpha_j$ 's. Nevertheless, minimizing  $\Gamma$  is a good rule of thumb.

- Comments on optimal rate allocation:

$$R_i = R + \frac{1}{2} \log_2 \frac{\sigma_i^2 \alpha_i}{\left( \prod_{j=1}^k \sigma_j^2 \alpha_j \right)^{1/k}} \quad (*)$$

The rate  $R_i$  allocated to  $U_i$  is  $R$  plus one half the log of the ratio of  $\sigma_i^2 \alpha_i$  to the geometric mean of such.

This addition to  $R$  quantity is positive if  $\sigma_i^2 \alpha_i$  greater than the geometric mean, and negative if less.

- Substituting the optimal  $R_i$  into the high-resolution expression for  $\delta_{sq,U_i}(R_i)$  shows that when the quantizers are optimized, the dist'n of  $Q_i$  on  $U_i$  is

$$D_i \cong \frac{1}{12} \sigma_i^2 \alpha_i 2^{-2R_i} = \frac{1}{12} \left( \prod_{j=1}^k \sigma_j^2 \alpha_j \right)^{1/k} 2^{-2R} = \text{same for all } i.$$

Thus with the optimal allocation, each coefficient is quantized with the same distortion. It isn't that we allow larger distortion for some coefficients, rather we need to assign more rate to the coefficients that have larger  $\sigma_i^2 \alpha_i$  in order that they be quantized with the same distortion. (Compare to JPEG.)

- What to do if (\*) yields a negative  $R_i$  for some  $i$ ? (This happens when  $\sigma_i^2 \alpha_i$  is sufficiently small.) In this case, the high-rate assumption is invalid, and we should not use the high-resolution analysis. Actually, we need each  $R_i \geq 3$  for each  $i$  for this high-resolution analysis to be accurate.

Tr-21

- Suppose we were to give the same rate allocation to all coefficients, i.e.  $R_i = R$ , then in Gaussian case

$$D \cong \frac{1}{k} \sum_{i=1}^k \frac{1}{12} \sigma_i^2 \alpha_x 2^{-2R_i} = \frac{1}{12} \alpha_x \frac{1}{k} \sum_{i=1}^k \sigma_i^2 2^{-2R} = \frac{1}{12} \alpha_x \sigma_X^2 2^{-2R} = \delta_{sq,X}(R).$$

That is, with equal bit allocation, the transform code gives the same performance as direct scalar quantization of the  $X_i$ 's, no matter what orthogonal transform is chosen. We conclude that the proper bit allocation is crucial.

- Recall zero mean assumption. If the data does not have zero mean, then the mean could be subtracted from the data before transform coding, and added back to the reconstructions produced by the decoder. But actually, the analysis applies as is.

Tr-22

## Step 2: The Optimal Transform

We now ask how to choose the orthogonal transform  $T$  to minimize

$$\Gamma \triangleq \left( \prod_{j=1}^k \sigma_j^2 \right)^{1/k}$$

This will be the optimal transform for the Gaussian case, and a "good" transform more generally.

**Main Fact:** The transform  $T$  that minimizes  $\Gamma$  is a Karhunen-Loeve Transform (KLT). A KLT is any transform whose rows are an orthonormal set of eigenvectors of the covariance matrix  $K_{\underline{X}}$  of  $\underline{X}$ . For a KLT, the coefficient variances  $\sigma_j^2$  are the eigenvalues of  $K_{\underline{X}}$ , denoted  $\lambda_1, \dots, \lambda_k$  and

$$\Gamma = \left( \prod_{j=1}^k \lambda_j \right)^{1/k} = |K_{\underline{X}}|^{1/k}$$

where  $|K_{\underline{X}}|$  denotes the determinant of  $K_{\underline{X}}$ .

Tr-23

The proof is based on:

### More Facts from Linear Algebra

**Definition:** If  $K$  is a  $k \times k$  matrix,  $\lambda$  is a real or complex number and  $\underline{v}$  is a  $k$ -dimensional vector such that  $K\underline{v} = \lambda\underline{v}$ , then  $\lambda$  is said to be an eigenvalue of  $K$ ,  $\underline{v}$  is said to be an eigenvector of  $K$ , and  $(\lambda, \underline{v})$  are an eigenpair for  $K$ .

**Fact 1:** Every  $k \times k$  matrix has  $k$  eigenvalues, though they need not be distinct.

**Fact 2:** The eigenvalues of a  $k \times k$  diagonal matrix  $K$  are the diagonal elements.

**Proof:** One may directly verify that the  $(K(i,i), \underline{v})$  is an eigenpair, where  $\underline{v} = (0 \dots 0 \ 1 \ 0 \dots 0)$  with a 1 in the  $i$ th place.

**Fact 3:** Real symmetric matrices have real (as opposed to complex) eigenvalues.

The eigenvectors associated with distinct eigenvalues are orthogonal.

There is an orthonormal set of eigenvectors.

If the matrix is also positive (respectively, nonnegative) definite, i.e.  $\underline{x}^T K \underline{x} > 0$  for all  $\underline{x}$ , (respectively,  $\underline{x}^T K \underline{x} \geq 0$ ), then the eigenvalues are positive (respectively, nonnegative).

**Fact 4:** The determinant of a square matrix is the product of its eigenvalues. i.e.e if a  $k \times k$  matrix has eigenvalues  $\lambda_1, \dots, \lambda_k$ , then

Tr-24

$$|K| = \prod_{j=1}^k \lambda_j$$

Fact 5: The determinant of an orthogonal matrix  $T$  is  $|T| = \pm 1$ .

Proof: Suppose  $\lambda$  is an eigenvalue of  $T$  and  $\underline{y}$  is a corresponding eigenvector, i.e.  $T\underline{y} = \lambda\underline{y}$ . Then

$$\|\underline{y}\| = \|T\underline{y}\| = \|\lambda\underline{y}\| = |\lambda| \|\underline{y}\|$$

which implies  $|\lambda| = 1$ . Since all the eigenvalues have magnitude one, Fact 1 implies  $|T| = \pm 1$ .

Tr-25

### Facts about covariance matrices

Fact 6: For any  $k \times k$  covariance matrix  $K$ ,  $\prod_{i=1}^k K_{i,i} \geq |K|$ , equality iff  $K$  is diag'l.

Proof: See Gersho and Gray, pp. 241-242

Fact 7: If  $\underline{U} = T\underline{X}$ , then  $K_U = T K_X T^t$

Proof:  $K_U = E \underline{U} \underline{U}^t = E T \underline{X} (T \underline{X})^t = E T \underline{X} \underline{X}^t T^t = T E \underline{X} \underline{X}^t T^t = T K_X T^t$ .

Fact 8: If  $T$  is orthogonal and  $\underline{U} = T\underline{X}$ , then

(a)  $K_X$  and  $K_U$  have the same eigenvalues.

(b)  $|K_U| = |K_X|$

Proof: (a) Let  $(\lambda, \underline{y})$  be an eigenpair for  $K_X$ . Then  $K_U T \underline{y} = T K_X T^t T \underline{y} = T \lambda \underline{y} = \lambda T \underline{y}$ . Hence,  $(\lambda, T \underline{y})$  is an eigenpair for  $K_U$ .

Thus every eigenvalue of  $K_X$  is also an eigenvalue of  $K_U$ .

The same argument applied to  $K_U$  and  $T^{-1}$  shows that every eigenvalue of  $K_U$  is also an eigenvalue of  $K_X$ . Thus  $K_X$  and  $K_U$  have the same eigenvalues.

(b) Since  $K_X$  and  $K_U$  have the same eigenvalues, Fact 1 implies they have the same determinants.

Fact 9: For any covariance matrix there is a set of  $k$  orthonormal eigenvectors.

Proof: Because covariance matrices are real and symmetric, and Fact 3.

Tr-26

## Proof of Optimality of KLT

Lemma: If  $T$  is orthogonal and makes  $K_U$  diagonal, then for any other orthogonal matrix  $\tilde{T}$ ,  $\tilde{\Gamma}^k \geq \Gamma^k$ , with equality if and only if  $K_U$  is diagonal.

Proof:

$$\begin{aligned} \tilde{\Gamma}^k &= \prod_{i=1}^k \tilde{\sigma}_i^2 \quad \text{by definition of } \tilde{\Gamma}, \text{ where } \tilde{\sigma}_i^2 = E\tilde{U}_i^2 \text{ \& } \tilde{U} = \tilde{T}X. \\ &\geq |K_U| \quad \text{Fact 6 fact that } \tilde{\sigma}_i^2\text{'s are diag elements of } K_U \\ &= |K_X| \quad \text{Fact 8(b)} \\ &= |K_U| \quad \text{Fact 8(b)} \\ &= \prod_{i=1}^k \sigma_i^2 \quad \text{Fact 6 \& fact that } \sigma_i^2\text{'s are diag elements of } K_U \\ &= \Gamma^k \text{ for } T \end{aligned}$$

By Fact 6, equality holds iff  $K_U$  is diagonal.

We see from this lemma that one can do no better than to make the transform coefficients have a diagonal covariance matrix.

When this is done, the coef variances are the diagonal elements, which by Fact 2, are the eigenvalues of  $K_U$ . By Fact 8a, these are also the eigenvalues of  $K_X$  as claimed in the Main Fact.

Tr-27

It remains only to show that there is a choice of  $T$  that makes  $K_U$  diagonal.

Accordingly, let  $T$  be the matrix mentioned in the Main Fact, i.e. its rows are an orthonormal set of eigenvectors  $t_1, \dots, t_k$

Then, by Fact 5, and the fact that  $t_i$  is an eigenvector with eigenvalue  $\lambda_i$

$$\begin{aligned} K_U = T K_X T^t &= \begin{bmatrix} -t_1^t & & & \\ -t_2^t & & & \\ \dots & & & \\ -t_k^t & & & \end{bmatrix} K_X \begin{bmatrix} | & | & & | \\ t_1 & t_2 & \dots & t_k \\ | & | & & | \end{bmatrix} \\ &= \begin{bmatrix} -t_1^t & & & \\ -t_2^t & & & \\ \dots & & & \\ -t_k^t & & & \end{bmatrix} \begin{bmatrix} | & | & & | \\ \lambda_1 t_1 & \lambda_2 t_2 & \dots & \lambda_k t_k \\ | & | & & | \end{bmatrix} = \begin{bmatrix} \lambda_1 0 & \dots & 0 \\ 0 & \lambda_2 0 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & \dots & 0 & \lambda_k \end{bmatrix}. \end{aligned}$$

So as we hoped,  $K_U$  is diagonal, which completes the proof of the main result.

Tr-28



Substituting the value  $\Gamma = |K_X|^{1/k}$  into the expression (\*\*) for distortion gives:

**The OPTA Function for k-dimensional Transform Coding applied to a Zero-Mean Stationary, Gaussian Source:**

For large R,

$$\delta_{\text{tr}}(k, R) \cong \frac{1}{12} |K_X|^{1/k} \alpha 2^{-2R}, \quad \text{where } \alpha = 32.6 \text{ for FRC, } 17.08 \text{ for VRC.}$$

SNR Gain:

$$10 \log_{10} \frac{\delta_{\text{sq}}(R)}{\delta_{\text{tr}}(k, R)} \cong 10 \log_{10} \frac{\sigma_X^2}{|K_X|^{1/k}}$$

The best transform is the KLT, i.e. rows are orthonormal eigenvectors for  $K_X$ .

The resulting coefficients  $U_1, \dots, U_k$  are uncorrelated (indeed, independent),

Their variances  $\sigma_1^2, \dots, \sigma_k^2$  equal the eigenvalues  $\lambda_1, \dots, \lambda_k$  of  $K_X$ .

The rate allocated to the  $i$ th coefficient is:  $R_i = R + \frac{1}{2} \log_2 \frac{\lambda_i}{|K_X|^{1/k}}$

The resulting coeff. distortions are all the same and equal to  $\delta_{\text{tr}}(k, R)$ .