

Source Coding

Tutorial Problems / **Solutions**

WS 2020/2021

Norbert Goertz

Institute of Telecommunications

TU Wien

`norbert.goertz@tuwien.ac.at`

Problem 1 / Solution

(a)

Curve A: granular distortion (increases with Δ)

Curve B: overload distortion (decreases with Δ)

The granular distortion (or granular quantiser noiser power) is proportional to the quantiser step-size Δ : the larger the step-size, the coarser the quantisation, and thus, larger the distortion.

Overload distortion occurs when the amplitude of the signal exceeds the range of values covered by the quantiser. For a fixed number of quantiser intervals L , the larger the interval size Δ , the larger is the range of amplitudes covered by the quantiser, and thus, the smaller the overload distortion.

(b)

A uniform quantiser with the step-size Δ and rate R bits per sample has the granular region $[-x_{\max}, x_{\max}]$ divided into $L = 2^R$ intervals of equal size Δ . The reproducer values \hat{x}_i are (only for uniform quantisation) the middle points of each interval:

$$\hat{x}_i = -x_{\max} + i \cdot \Delta + \Delta/2, \quad i = 0, \dots, L - 1.$$

The decision boundaries g_i are given by

$$g_i = -x_{\max} + i \cdot \Delta, \quad i = 0, \dots, L.$$

The quantiser maps every signal value x in the i th interval, i.e., $g_i \leq x < g_{i+1}$, into the reproducer value \hat{x}_i , and outputs the index i , represented by R bits. Values of x that fall into the overload region, i.e., $x < g_0$ and $x > g_L$, are mapped to \hat{x}_0 and \hat{x}_{L-1} , respectively.

(c)

The power P_Q of the quantiser error is obtained as the sum of the mean square errors in all quantiser intervals. There are L closed intervals in the granular region $[-x_{\max}, x_{\max}]$ and 2 half-opened intervals in the overload region: $(-\infty, -x_{\max})$ and $(x_{\max}, +\infty)$. Thus, the quantiser error power equals

$$P_Q = \sum_{i=0}^{L-1} \int_{g_i}^{g_{i+1}} (\hat{x}_i - x)^2 p(x) dx + \int_{-\infty}^{g_0} (\hat{x}_0 - x)^2 p(x) dx + \int_{g_L}^{+\infty} (\hat{x}_{L-1} - x)^2 p(x) dx,$$

where $p(x)$ is the probability density function (PDF) of the input signal x .

Problem 2 / Solution

(a)

We consider a source with Laplace distribution $p(x) = \frac{1}{2}e^{-|x|}$, $-\infty < x < +\infty$. Since the mean value of a Laplace random variable is zero, i.e. $E\{X\} = 0$, the variance equals the power:

$$P_x = \sigma^2 = \int_{-\infty}^{+\infty} x^2 p(x) dx.$$

Due to the symmetry of the Laplace PDF $p(x)$, we can write

$$P_x = 2 \int_0^{+\infty} x^2 \frac{1}{2} e^{-x} dx = \int_0^{+\infty} x^2 e^{-x} dx = 2. \quad (1)$$

Thus, it follows that $K \triangleq P_x/x_{\max}^2 = 2/x_{\max}^2$. For a given probability $p_1 \triangleq \Pr\{|x| > x_{\max}\}$ we obtain:

$$p_1 = 2 \int_{x_{\max}}^{+\infty} p(x) dx = 2 \int_{x_{\max}}^{+\infty} \frac{1}{2} e^{-x} dx = -e^{-x} \Big|_{x_{\max}}^{+\infty} = e^{-x_{\max}} \implies \begin{cases} x_{\max} = -\ln p_1; \\ K = 2/(\ln p_1)^2. \end{cases} \quad (2)$$

For $p_1 = 0.001$ we find

$$x_{\max}|_{p_1=0.001} = 6.9078; \quad K|_{p_1=0.001} = 0.0419.$$

(b)

The overload distortion P_O is given by

$$P_O = \int_{-\infty}^{-x_{\max}} (\hat{x}_0 - x)^2 p(x) dx + \int_{x_{\max}}^{+\infty} (\hat{x}_{2R-1} - x)^2 p(x) dx.$$

If we approximate the reproducer values with $\hat{x}_0 \approx -x_{\max}$ and $\hat{x}_{2R-1} \approx x_{\max}$, this yields the following approximate result:

$$\begin{aligned} P_O &\approx \int_{-\infty}^{-x_{\max}} (-x_{\max} - x)^2 \frac{1}{2} e^x dx + \int_{x_{\max}}^{+\infty} (x_{\max} - x)^2 \frac{1}{2} e^{-x} dx = \\ &= 2 \int_{x_{\max}}^{+\infty} (x_{\max} - x)^2 \frac{1}{2} e^{-x} dx = \\ &= e^{-x_{\max}} \int_{x_{\max}}^{+\infty} (x_{\max} - x)^2 e^{(x_{\max}-x)} dx = \\ &= 2e^{-x_{\max}} = 2p_1. \end{aligned} \quad (3)$$

The *true* value of the overload distortion is *underestimated* by this approximation. The explanation is the following:

The approximation (3) is obtained assuming that the minimum and the maximum reproducer values are $\hat{x}_0 \approx -x_{\max}$ and $\hat{x}_{2^{R-1}} \approx x_{\max}$. The *exact* reproducer values are $\hat{x}_0 = -x_{\max} + \Delta/2$ and $\hat{x}_{2^{R-1}} = x_{\max} - \Delta/2$, i.e., they have smaller absolute values than the approximations. Obviously, $\hat{x}_0 \approx -x_{\max}$ is a better approximation for the values $x < -x_{\max}$ (in the mean squared-error sense) than the exact value $\hat{x}_0 = -x_{\max} + \Delta/2$. The same holds for approximating values $x > x_{\max}$ with $\hat{x}_{2^{R-1}} \approx x_{\max}$ and $\hat{x}_{2^{R-1}} = x_{\max} - \Delta/2$. Thus, the overload distortion obtained with approximate values instead of \hat{x}_0 and $\hat{x}_{2^{R-1}}$ is smaller than the exact distortion obtained with the exact values, i.e.:

$$P_O \geq 2e^{-x_{\max}}.$$

This can be verified by computing the exact distortion

$$P_O = \int_{-\infty}^{-x_{\max}} \left(-x_{\max} + \frac{\Delta}{2} - x\right)^2 \frac{1}{2} e^x dx + \int_{x_{\max}}^{+\infty} \left(x_{\max} - \frac{\Delta}{2} - x\right)^2 \frac{1}{2} e^{-x} dx.$$

This yields the result

$$P_O = e^{-x_{\max}} \left(2 + \Delta + \frac{\Delta^2}{4}\right) \geq 2e^{-x_{\max}}. \quad (4)$$

Equality in (4) is achieved when $\Delta \rightarrow 0$, i.e., when $R \rightarrow \infty$ (since $\Delta = x_{\max}/2^{R-1}$).

(c)

From (4) it is clear that overload distortion depends not only on x_{\max} , i.e., on $p_1 = e^{-x_{\max}}$, but also on the step-size Δ , and thus implicitly on the rate R . Therefore, in general, the answer of whether or not it is justified to neglect overload distortion depends on the actual rate.

The signal-to-noise ratio is given by:

$$\frac{SNR}{\text{dB}} = 10 \log \frac{P_x}{P_G + P_O} = 10 \log \frac{P_x}{P_G} - 10 \log \left(1 + \frac{P_O}{P_G}\right).$$

By substituting the granular distortion $P_G = \Delta^2/12 = x_{\max}^2 / (12 \cdot 2^{2R-2})$ and by introducing $c \triangleq \frac{P_O}{P_G}$, we obtain the equivalent expression:

$$\begin{aligned} \frac{SNR}{\text{dB}} &= R \cdot 20 \log(2) + 10 \log(3K) - 10 \log(1 + c) \\ &= R \cdot 6.02 + 4.77 + 10 \log K - 10 \log(1 + c). \end{aligned} \quad (5)$$

If $c \ll 1$, i.e., $P_O \ll P_G$, then it is reasonable to neglect the overload distortion P_O , and then $SNR/\text{dB} \approx 10 \log \frac{P_x}{P_G} = R \cdot 20 \log(2) + 10 \log(3K)$.

Since

$$c = \frac{P_O}{P_G} = \frac{2e^{-x_{\max}}}{\Delta^2/12} = \frac{2p_1}{x_{\max}^2/(12 \cdot 2^{2R-2})} = \frac{6p_1 2^{2R}}{(\ln p_1)^2},$$

for $p_1 = 0.001$ and $R = 5$, we obtain $c \approx 0.129$, i.e.,

$$10 \log(1 + c) \approx 0.526 \text{ dB}.$$

Hence, overload distortion introduces a further loss in SNR of about 0.5 dB, compared to SNR with granular distortion only.

The more precise value of the loss can be obtained by using the exact formula for overload distortion, given by (4) instead of $P_O = 2p_1$. Then we obtain

$$c = \frac{e^{-x_{\max}}(2 + \Delta + \Delta^2/4)}{\Delta^2/12} = \frac{6p_1 2^{2R}}{(\ln p_1)^2} + 3p_1 - \frac{6p_1 2^R}{\ln p_1},$$

which yields the value $c = 0.1595$ or $10 \log(1 + c) = 0.642$ dB.

Problem 3 / Solution

(a)

The contribution of the i th quantiser interval to the total noise power is given by the formula

$$P_{Q_i} = \int_{\hat{x}_i - \Delta_{x_i}/2}^{\hat{x}_i + \Delta_{x_i}/2} (\hat{x}_i - x)^2 p(x) dx. \quad (6)$$

To evaluate this integral approximately, we can assume that the PDF $p(x)$ is constant within the considered interval, i.e. we assume small intervals and, hence, a high bit rate:

$$p(x) \approx p(\hat{x}_i), \quad \text{for } \hat{x}_i - \Delta_{x_i}/2 \leq x \leq \hat{x}_i + \Delta_{x_i}/2. \quad (7)$$

The value of $p(\hat{x}_i)$ (the PDF in the point $x = \hat{x}_i$) can be obtained from

$$\begin{aligned} p_i &\triangleq \Pr \{ \hat{x}_i - \Delta_{x_i}/2 \leq x \leq \hat{x}_i + \Delta_{x_i}/2 \} = \\ &= \int_{\hat{x}_i - \Delta_{x_i}/2}^{\hat{x}_i + \Delta_{x_i}/2} p(x) dx \approx p(\hat{x}_i) \int_{\hat{x}_i - \Delta_{x_i}/2}^{\hat{x}_i + \Delta_{x_i}/2} dx = p(\hat{x}_i) \cdot \Delta_{x_i} \implies p(\hat{x}_i) \approx p_i / \Delta_{x_i}, \end{aligned} \quad (8)$$

where p_i is the probability that the sample value x 'falls' into the i th quantiser interval. By using (8), we obtain the approximate formula for (6):

$$P_{Q_i} \approx p(\hat{x}_i) \int_{\hat{x}_i - \Delta_{x_i}/2}^{\hat{x}_i + \Delta_{x_i}/2} (\hat{x}_i - x)^2 dx = \frac{p_i}{\Delta_{x_i}} \cdot \frac{\Delta_{x_i}^3}{12} = p_i \frac{\Delta_{x_i}^2}{12}. \quad (9)$$

(b)

The total quantiser noise power equals

$$P_Q = \sum_{i=0}^{L-1} P_{Q_i} = 2 \sum_{i=L/2}^{L-1} P_{Q_i} \approx 2 \sum_{i=L/2}^{L-1} p_i \frac{\Delta_{x_i}^2}{12}, \quad (10)$$

where the second equality follows from the symmetry of the PDF $p(x)$ and the third equality follows from approximation (9).

(c)

The compressor characteristic is given by $y = g(x)$, where $g(x)$ is strictly monotone and invertible function (and thus differentiable). The derivative of $g(x)$ in the point $x = \hat{x}_i$ is defined as

$$\left. \frac{dg(x)}{dx} \right|_{x=\hat{x}_i} = g'(\hat{x}_i) = \lim_{\Delta_{x_i} \rightarrow 0} \frac{g(\hat{x}_i + \Delta_{x_i}) - g(\hat{x}_i)}{\Delta_{x_i}}.$$

If we approximate $g(x)$ with a straight line in $x = \hat{x}_i$, this yields

$$g'(\hat{x}_i) \approx \frac{g(\hat{x}_i + \Delta_{x_i}) - g(\hat{x}_i)}{\Delta_{x_i}} = \frac{g(\hat{x}_{i+1}) - g(\hat{x}_i)}{\Delta_{x_i}} = \frac{\hat{y}_{i+1} - \hat{y}_i}{\Delta_{x_i}} = \frac{\Delta_y}{\Delta_{x_i}} \implies \Delta_{x_i} \approx \frac{\Delta_y}{g'(\hat{x}_i)}.$$

(d)

For A-law compressor characteristics the step size $\Delta_x(x)$ is proportional to the signal amplitude: $\Delta_x(x) \sim x$. By combining this with the previous approximation for Δ_{x_i} , we obtain:

$$\Delta_{x_i} \approx \frac{\Delta_y}{g'(\hat{x}_i)} \sim \hat{x}_i \Rightarrow \frac{1}{g'(\hat{x}_i)} \sim \hat{x}_i / \Delta_y \Rightarrow \frac{1}{g'(\hat{x}_i)} = c \cdot \hat{x}_i \Rightarrow \Delta_{x_i} \approx c \cdot \hat{x}_i \cdot \Delta_y \quad (11)$$

(e)

By plugging the expression (11) into the formula (10), we obtain the approximate expression for the total noise power:

$$P_Q \approx \sum_{i=0}^{L-1} p_i \frac{(c\Delta_y\hat{x}_i)^2}{12} = \frac{\Delta_y^2}{12} c^2 \cdot \underbrace{\sum_{i=0}^{L-1} \hat{x}_i^2 p_i}_{\approx P_x} \approx \frac{\Delta_y^2}{12} c^2 P_x, \quad (12)$$

where the last approximation follows from approximating the signal power with

$$\begin{aligned} P_x &= \int_{-\infty}^{+\infty} x^2 p(x) dx = \int_{-\infty}^{-x_{\max}} x^2 p(x) dx + \int_{-x_{\max}}^{x_{\max}} x^2 p(x) dx + \int_{x_{\max}}^{+\infty} x^2 p(x) dx \approx \\ &\approx \int_{-x_{\max}}^{x_{\max}} x^2 p(x) dx \approx \sum_{i=0}^{L-1} \hat{x}_i^2 p(\hat{x}_i) \Delta_{x_i} \approx \sum_{i=0}^{L-1} \hat{x}_i^2 p_i, \end{aligned} \quad (13)$$

which is obtained by neglecting the power in the tails (if they exist and are not covered by the choice of x_{\max}) of the signal PDF and by assuming that the PDF is constant within the intervals Δ_{x_i} .

(f)

The quantiser step-size in the y -domain is uniform, and for quantisation with R bits/sample and $y_{\max} = 1$ it equals

$$\Delta_y = \frac{2y_{\max}}{2^R} = \frac{1}{2^{R-1}},$$

which, when inserted into formula (12), yields:

$$P_Q \approx \frac{\Delta_y^2}{12} c^2 P_x = \frac{c^2}{12 \cdot 2^{2R-2}} P_x = \frac{c^2}{3 \cdot 2^{2R}} P_x. \quad (14)$$

(g)

From (14) we can directly obtain the expression for the SNR as

$$\begin{aligned} SNR/\text{dB} &= 10 \log \frac{P_x}{P_Q} = 10 \log \frac{3 \cdot 2^{2R}}{c^2} = R \cdot 20 \log 2 + 10 \log 3 - 20 \log c \\ &= R \cdot 6.02 + 4.77 - 20 \log c. \end{aligned} \quad (15)$$

The SNR of the quantiser with logarithmic characteristic (A-law only for large amplitudes!) does not depend on the signal power, which is its main advantage. This can be compared with the uniform quantiser for which the SNR is given by expression (5) ($SNR/\text{dB} = R \cdot 6.02 + 4.77 + 10 \log K - 20 \log(1 + c_1)$), which depends on the signal power through the factor $K = P_x/x_{\max}^2$.

(h)

The characteristic of the A-law compressor is given by

$$y = g(x) = \begin{cases} \text{sign}(x) \frac{1 + \ln(A|x|)}{1 + \ln(A)}, & \frac{1}{A} \leq |x| \leq 1 \\ \frac{A}{1 + \ln(A)} x, & -\frac{1}{A} \leq x \leq \frac{1}{A}. \end{cases}$$

Since $1/g'(x) = cx$, the value of the constant c for the range $1/A \leq |x| \leq 1$ is obtained by taking the first derivative of $g(x)$ in this range (only positive values, symmetric for negative x):

$$g'(x)|_{1/A \leq x \leq 1} = \frac{d}{dx} \left(\text{sign}(x) \frac{1 + \ln(A|x|)}{1 + \ln(A)} \right) = \frac{1}{x(1 + \ln A)},$$

by comparison: $c = \frac{1}{xg'(x)} = 1 + \ln A \implies c|_{A=87.56} = 5.4723.$

(i)

By substituting the obtained value for c into expression (15), we obtain for the SNR:

$$\frac{SNR}{\text{dB}} = R \cdot 6.02 + 4.77 - 20 \log(1 + \ln A) = R \cdot 6.02 - 9.99. \quad (16)$$

We compare this result with the SNR of the uniform quantiser, given by formula (5), where we assume that the overload distortion is zero. Then the SNR is given by $SNR/\text{dB} = R \cdot 6.02 + 10 \log(3K)$. To make a comparison, we consider quantisation with the same rate R , and a uniformly distributed source signal with full load, i.e. $X \sim \text{Unif}[-x_{\max}, +x_{\max}]$. Then, $K = P_x/x_{\max}^2 = 1/3$, and the SNR of the uniform quantiser is $SNR/\text{dB} = R \cdot 6.02$. When compared to expression (16), we conclude that logarithmic quantisation yields a *loss* of 9.99 dB. However, this result holds only for 'full-load' uniform sources, which are rarely found in practise. If, for instance, we have a triangular source distribution that covers the range $(-\frac{1}{3} \cdot x_{\max}, +\frac{1}{3} \cdot x_{\max})$ we obtain $P_x = x_{\max}^2/54$ as the signal power, so $K = 1/54$ and $10 \log_{10}(3K) = 10 \log_{10}(1/18) = -12.6$ dB: in this case the logarithmic quantiser is still well in the region described by (16), i.e., it performs better than the uniform quantiser.

(j)

The A-law characteristic is linear for small amplitudes, which means that the uniform quantiser intervals on the y -axis are mapped to intervals of constant size on the x -axis. The only difference is that the interval size is scaled by the gradient of the A-law characteristic, so that the intervals are smaller on the x -axis. Therefore, the uniform quantiser on the x -axis seems to have a higher rate than the one on the y -axis. Of course this is only true for very small amplitudes in the range $-\frac{1}{A} \leq x \leq \frac{1}{A}$.

(k)

The first derivative of the A-law characteristic for $-\frac{1}{A} \leq x \leq \frac{1}{A}$ is

$$g'(x)|_{-1/A \leq x \leq 1/A} = \frac{d}{dx} \left(\frac{Ax}{1 + \ln(A)} \right) = \frac{A}{1 + \ln A} \Big|_{A=87.56} = 16 \quad (17)$$

(l)

Hence, we have $\Delta_{x_i} = \Delta_y/g'(x) = \Delta_y/16$ as the effective interval size for small values of x .

(m)

As for small intervals we have a uniform quantiser for the x -values, we can apply the standard formula

$$\frac{SNR}{\text{dB}} = 10 \log_{10} \frac{P_x}{P_Q}.$$

With

$$P_Q = \frac{\Delta_{x_i}^2}{12} = \frac{(\Delta_y/16)^2}{12} = \frac{(2 \cdot x_{\max} \cdot 2^{-R})^2}{12 \cdot 16^2} = \frac{x_{\max}^2 \cdot 2^{-2R}}{3 \cdot 16^2}$$

we obtain

$$\frac{SNR}{\text{dB}} = 10 \log_{10} \left(3 \cdot \frac{P_x}{x_{\max}^2} \right) + 20 \log_{10}(16) + 20 \log_{10}(2) \cdot R.$$

With the normalised signal power $K = P_x/x_{\max}^2$ we can write

$$\frac{SNR}{\text{dB}} = R \cdot 6.02 + \underbrace{4 \cdot 20 \log_{10}(2)}_{=24.08 \text{ dB}} + 10 \log_{10} \left(3 \cdot \frac{P_x}{x_{\max}^2} \right) = (R + 4) \cdot 6.02 + 10 \log_{10}(3K)$$

The last equation is structural identical with the solution for a uniform quantiser (see lecture notes), but we have a constant gain of 24.08 dB due to the gradient of the A-law characteristics for small x -values. The gain may be interpreted as a “virtual” increase in the bit rate by 4 bits/sample.

Problem 4 / Solution

(a)

We observe a random variable X that is distributed according to the PDF $p(x)$, plotted in Fig. 1:

$$p(x) = \begin{cases} |x|, & -1 \leq x \leq 1 \\ 0, & \text{otherwise} \end{cases} \quad (18)$$

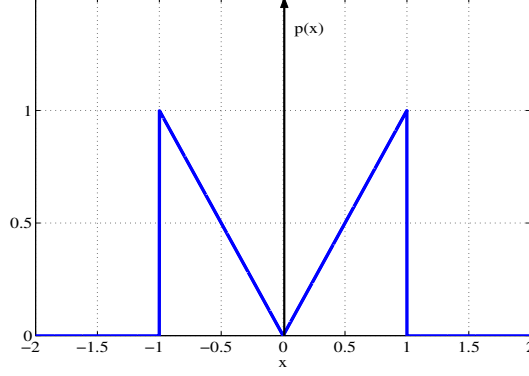


Figure 1: PDF defined by (18)

The mean μ_x , the variance σ_x^2 , and the power P_x are given respectively by:

$$\begin{aligned} \mu_x &= \int_{-1}^1 x \cdot |x| dx = \int_{-1}^0 -x^2 dx + \int_0^1 x^2 dx = 0 ; \\ P_x = \sigma_x^2 &= \int_{-1}^1 x^2 \cdot |x| dx = 2 \cdot \int_0^1 x^3 dx = \frac{1}{2} . \end{aligned}$$

(b)

For a uniform quantiser with $L = 2^R = 4$ quantisation intervals within the range $[-x_{\max} \ x_{\max}] = [-1 \ 1]$, the step-size is $\Delta = 2x_{\max}/L = 1/2$. The characteristic of this quantiser is shown in Figure 2. The decision boundaries between the intervals are

$$g_i = -x_{\max} + i \cdot \Delta = -1 + i/2, \quad 0 \leq i \leq 4 \quad \Longrightarrow \quad g_i \in \left\{ -1, -\frac{1}{2}, 0, \frac{1}{2}, 1 \right\},$$

and the reproducer values are

$$\hat{x}_i = -x_{\max} + \Delta/2 + i \cdot \Delta = -3/4 + i/2, \quad 0 \leq i \leq 3 \quad \Longrightarrow \quad \hat{x}_i \in \left\{ -\frac{3}{4}, -\frac{1}{4}, \frac{1}{4}, \frac{3}{4} \right\}.$$

(c)

Since the range of the quantiser is matched to the support of the PDF $p(x)$, there will be no overload distortion, only granular distortion. It is given by

$$P_Q = \sum_{i=0}^{L-1} \int_{g_i}^{g_{i+1}} (\hat{x}_i - x)^2 p(x) dx = \sum_{i=0}^3 \int_{-1+i/2}^{-1+(i+1)/2} (\hat{x}_i - x)^2 |x| dx. \quad (19)$$

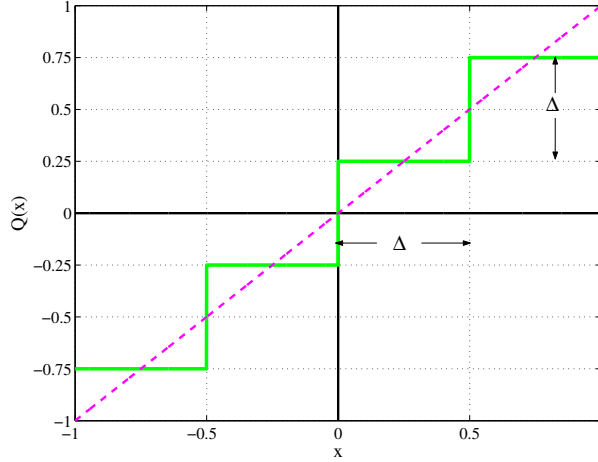


Figure 2: Mid-riser characteristic $Q(x)$ of the quantiser from **b)**

(d)

By exploiting the symmetries of the PDF and the reproducer values (i.e. $\hat{x}_0 = -\hat{x}_3$, $\hat{x}_1 = -\hat{x}_2$) we can re-write (19):

$$\begin{aligned}
 P_Q &= - \int_{-1}^{-1/2} (\hat{x}_0 - x)^2 x dx - \int_{-1/2}^0 (\hat{x}_1 - x)^2 x dx + \int_0^{1/2} (\hat{x}_2 - x)^2 x dx + \int_{1/2}^1 (\hat{x}_3 - x)^2 x dx = \\
 &= 2 \cdot \left(\int_{-1}^{-1/2} (\hat{x}_0 - x)^2 (-x) dx + \int_{-1/2}^0 (\hat{x}_1 - x)^2 (-x) dx \right) = \\
 &= \frac{1}{2} + \frac{7}{6} \hat{x}_0 + \frac{3}{4} \hat{x}_0^2 + \frac{1}{6} \hat{x}_1 + \frac{1}{4} \hat{x}_1^2.
 \end{aligned} \tag{20}$$

(e)

By direct evaluation of the expression (20), we obtain $P_Q = 1/48$. This yields the SNR value

$$\frac{SNR}{\text{dB}} = 10 \log \frac{P_x}{P_Q} = 10 \log \frac{1/2}{1/48} = 10 \log 24 = 13.8$$

(f)

We need to find new reproducer values \hat{x}'_i , such that the power of the quantisation error P_Q is minimised. New reproducer values retain the property $\hat{x}'_0 = -\hat{x}'_3$, $\hat{x}'_1 = -\hat{x}'_2$. We obtain the new values by minimising the expression (20) with respect to \hat{x}_0 and \hat{x}_1 :

$$\begin{aligned}
 \hat{x}'_0 = \arg \min_{\hat{x}_0} P_Q &\implies \left. \frac{\partial P_Q}{\partial \hat{x}_0} \right|_{\hat{x}_0 = \hat{x}'_0} = 0 ; \\
 \hat{x}'_1 = \arg \min_{\hat{x}_1} P_Q &\implies \left. \frac{\partial P_Q}{\partial \hat{x}_1} \right|_{\hat{x}_1 = \hat{x}'_1} = 0 ;
 \end{aligned}$$

We obtain

$$\begin{aligned} \frac{7}{6} + \frac{3}{4} \cdot 2\hat{x}'_0 &= 0 \quad \rightarrow \quad \hat{x}'_0 = -7/9 \quad \text{and} \quad \hat{x}'_3 = -\hat{x}'_0 = 7/9 \\ \frac{1}{6} + \frac{1}{4} \cdot 2\hat{x}'_1 &= 0 \quad \rightarrow \quad \hat{x}'_1 = -1/3 \quad \text{and} \quad \hat{x}'_2 = -\hat{x}'_1 = 1/3 \end{aligned} \quad (21)$$

(g)

The new quantiser characteristic is shown in Figure 3.

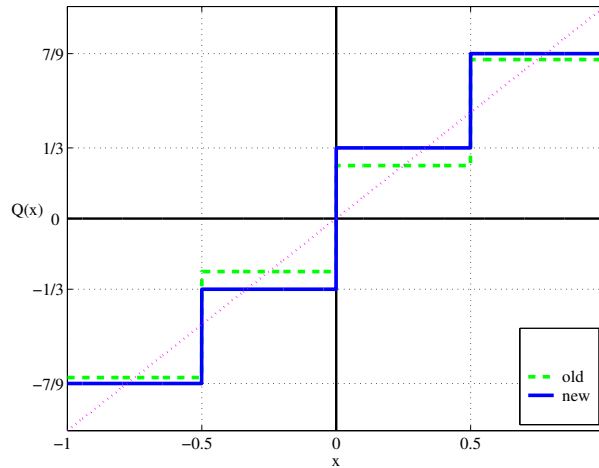


Figure 3: Characteristic $Q(x)$ of the new quantiser from f)

(h)

By inserting the new reproducer values given by (21) into the expression (20), we obtain the new power of the quantiser error:

$$P'_Q = \frac{1}{2} + \frac{7}{6} \cdot \left(-\frac{7}{9}\right) + \frac{3}{4} \cdot \left(-\frac{7}{9}\right)^2 + \frac{1}{6} \cdot \left(-\frac{1}{3}\right) + \frac{1}{4} \cdot \left(-\frac{1}{3}\right)^2 = \frac{1}{54},$$

which yields the new SNR:

$$\frac{SNR'}{\text{dB}} = 10 \log \frac{P_x}{P'_Q} = 10 \log \frac{1/2}{1/54} = 10 \log 27 = 14.3$$

(i)

The optimised reproducer values are not equidistant; they are, compared to the old ones, shifted away from zero, towards the outer regions of the interval $[-1 + 1]$, where the PDF-values are larger, i.e., where the realisations of the random variable X will occur more often. This enlarges the SNR from 13.8 dB to 14.3 dB, i.e., we have a gain of about 0.5 dB due to the optimisation of the reproducer values for the given quantiser intervals.

(j)

A further improvement would be possible by calculating new decision boundaries from the new reproducer values found in i). Then again new reproducer values could be found for the new intervals etc. This iterative procedure is exactly what was discussed in the lectures under “optimal quantisation”.

Problem 5 / Solution

(a)

The source samples we want to quantise are i.i.d., Gaussian distributed $X \sim \mathcal{N}(0, \sigma^2)$. The quantisation is performed with rate R bits/sample, where R is high.

The signal power is equal to the variance, as the signal has zero mean:

$$P_x = \sigma^2 = \int_{-\infty}^{+\infty} x^2 p(x) dx.$$

The overload probability equals

$$P_o = \Pr\{|x| > x_{\max}\} = 2 \int_{x_{\max}}^{+\infty} p(x) dx = \frac{2}{\sqrt{2\pi}\sigma} \int_{x_{\max}}^{+\infty} e^{-x^2/(2\sigma^2)} dx = \frac{2}{\sqrt{\pi}} \int_{\frac{x_{\max}}{\sigma\sqrt{2}}}^{+\infty} e^{-t^2} dt = \operatorname{erfc}\left(\frac{x_{\max}}{\sigma\sqrt{2}}\right), \quad (22)$$

where $\operatorname{erfc}(\cdot)$ is complementary error function, $\operatorname{erfc}(x) = 1 - \operatorname{erf}(x)$. From (22), we can obtain the range x_{\max} of the quantiser as

$$x_{\max} = \sigma \sqrt{2} \operatorname{erf}^{-1}(1 - P_o). \quad (23)$$

The distortion of the quantiser can be obtained by exploiting the high-rate property that quantising intervals Δ are very narrow and that the PDF is almost constant within these intervals:

$$P_Q \approx \sum_{i=0}^{2^R-1} p_i \frac{\Delta^2}{12} \approx \frac{\Delta^2}{12} \int_{-x_{\max}}^{+x_{\max}} p(x) dx = \frac{\Delta^2}{12} \Pr\{|x| \leq x_{\max}\} = \frac{x_{\max}^2}{12 \cdot 2^{2R-2}} (1 - P_o) \quad (24)$$

The expression for the SNR follows directly from (23) and (24):

$$\begin{aligned} \frac{SNR}{\text{dB}} &= 10 \log \frac{P_x}{P_Q} = 10 \log \frac{\sigma^2 \cdot 12 \cdot 2^{2R-2}}{x_{\max}^2 (1 - P_o)} = 10 \log \frac{3 \cdot 2^{2R}}{2(1 - P_o) [\operatorname{erf}^{-1}(1 - P_o)]^2} = \\ &= 20 \log(2) \cdot R + 10 \log(3) - 10 \log \left(2(1 - P_o) [\operatorname{erf}^{-1}(1 - P_o)]^2 \right) \Big|_{P_o=0.001} = \\ &= 6.02 \cdot R - 5.57. \end{aligned} \quad (25)$$

(b)

The best possible fixed-rate quantiser is a non-uniform quantiser with the optimum compressor characteristic $y = g_{\text{opt}}(x)$ (optimum in the sense that it has the lowest distortion). Such a quantiser has the best SNR, which, for a Gaussian source, equals

$$\frac{SNR_{\text{opt,fix}}^{\text{Gauss}}}{\text{dB}} = 6.02 \cdot R - 4.35. \quad (26)$$

Comparing with the SNR for the uniform quantiser given by (25), we conclude that using the *optimum* compressor improves SNR by 1.22 dB.

(c)

The lowest distortion achieved by an optimum scalar variable-rate quantiser is given by $P_Q = (1/12) 2^{2h(X)} 2^{-2R}$, where $h(X)$ is the differential entropy of the source. This yields the best SNR for a variable-rate scalar quantiser:

$$\frac{SNR_{opt,var}}{\text{dB}} = 6.02 \cdot R + 10 \log \left(\sigma^2 2^{-2h(X)} \right) + 10 \log(12).$$

The differential entropy of a Gaussian source is $h(X) = \frac{1}{2} \log_2(2\pi e \sigma^2)$, and thus $2^{2h(X)} = 2\pi e \sigma^2$. Hence we obtain

$$\frac{SNR_{opt,var}^{Gauss}}{\text{dB}} = 6.02 \cdot R - 10 \log(2\pi e) + 10 \log(12) = 6.02 \cdot R - 1.533.$$

(d)

The maximum possible SNR for any i.i.d. source (independent and identically distributed), achievable by any (variable or fixed-rate) quantisation scheme at high bit rate is given by the Shannon lower bound:

$$\frac{SNR_{opt}}{\text{dB}} = 6.02 \cdot R + 10 \log \left(\sigma^2 2^{-2h(X)} \right) + 10 \log(2\pi e),$$

which, for a Gaussian source, reduces to

$$\frac{SNR_{opt}^{Gauss}}{\text{dB}} = 6.02 \cdot R.$$

This theoretical lower bound is only 1.533 dB better than the maximum SNR achievable by an optimum scalar variable-rate quantiser. In case of a Gaussian source the Shannon Lower bound coincides with the Rate-Distortion Function at all rates, i.e., only in the uncorrelated Gaussian case we know an analytical result for the RDF that holds at all rates.

(e)

When the samples X are non-uniformly quantised, using a compressor with the transfer characteristic

$$y = g(x) = \alpha \frac{2}{\pi} \arctan(\beta x), \quad \alpha, \beta > 0,$$

then, the infinite support $(-\infty, +\infty)$ of the Gaussian PDF $p(x)$ is mapped onto a finite support in y -domain $[y_{\min}, y_{\max}]$, where $y_{\max} = -y_{\min} = g(\infty) = \alpha \frac{2}{\pi} \frac{\pi}{2} = \alpha$. The samples of $y \in [-\alpha, \alpha]$ are then uniformly quantised with the step size $\Delta = 2\alpha/2^R$.

(f)

The distortion of this quantiser is given by

$$\begin{aligned}
P_Q &= \frac{\Delta^2}{12} \int_{-\infty}^{+\infty} \frac{p(x)}{(g'(x))^2} dx = \frac{\Delta^2}{12} \int_{-\infty}^{+\infty} \frac{\frac{1}{\sqrt{2\pi\sigma}} e^{-x^2/(2\sigma^2)}}{\alpha^2 \beta^2 \frac{4}{\pi^2} \frac{1}{(1+\beta^2 x^2)^2}} dx \Bigg|_{\Delta=\alpha/2^{R-1}} = \\
&= \frac{\pi^2}{12 \cdot 2^{2R} \beta^2} \left[\underbrace{\frac{1}{\sqrt{2\pi\sigma}} \int_{-\infty}^{+\infty} e^{-x^2/(2\sigma^2)} dx}_{=1} + 2\beta^2 \underbrace{\frac{1}{\sqrt{2\pi\sigma}} \int_{-\infty}^{+\infty} x^2 e^{-x^2/(2\sigma^2)} dx}_{=\sigma^2} + \frac{\beta^4}{\sqrt{2\pi\sigma}} \underbrace{\int_{-\infty}^{+\infty} x^4 e^{-x^2/(2\sigma^2)} dx}_{=\frac{2 \cdot 3 \cdot \sqrt{\pi}}{2^3 (2\sigma^2)^{-5/2}}} \right] = \\
&= \frac{\pi^2}{12 \cdot 2^{2R} \beta^2} [1 + 2\beta^2 \sigma^2 + 3\beta^4 \sigma^4].
\end{aligned} \tag{27}$$

For a given rate R , the distortion does not depend on the parameter α , only on β .

(g)

Since the distortion does not depend on α , this parameter can be chosen arbitrarily. The optimum value of the parameter β is found by minimising expression (27):

$$\begin{aligned}
\beta_{\text{opt}} = \arg \min_{\beta} P_Q &\implies \frac{dP_Q}{d\beta} \Bigg|_{\beta=\beta_{\text{opt}}} = 0 \\
&\frac{d}{d\beta} \left(\frac{1}{\beta^2} + 2\sigma^2 + 3\beta^2 \sigma^4 \right) \Bigg|_{\beta=\beta_{\text{opt}}} = 0 \\
\beta_{\text{opt}} &= \frac{1}{\sigma \sqrt{3}}.
\end{aligned} \tag{28}$$

This optimum value of β yields the minimum distortion

$$P_{Q \text{ min}} = P_Q(\beta_{\text{opt}}) = \frac{\pi^2}{6 \cdot 2^{2R}} (1 + \sqrt{3}) \sigma^2.$$

From this we directly obtain the optimum SNR:

$$\frac{SNR_{\text{opt}}}{\text{dB}} = 10 \log \frac{P_x}{P_{Q \text{ min}}} = 10 \log \frac{6 \cdot 2^{2R}}{\pi^2 (1 + \sqrt{3})} = 6.02 \cdot R - 6.53.$$

(h)

We observe the same Gaussian source, however, the quantisation rate is now assumed to be low and fixed at $R = 1$ bit/sample. We want to find the best possible scalar quantiser for this case. Since the zero-mean Gaussian PDF is symmetric around zero and we have $L = 2$ symmetric reproducer values $\hat{x}_0 = -\hat{x}_1$, it is clear that the 'middle' decision boundary is $g_1 = 0$. The other two boundaries are symmetric, $g_0 = -g_2$, and they should be set to infinity: $g_0 = -\infty$, $g_2 = +\infty$. Thus, every value $x < 0$ will be quantised as \hat{x}_0 , and every $x > 0$ as \hat{x}_1 .

(i)+(j)+(k)

The two reproducer values are found by minimising the distortion:

$$\begin{aligned}
 P_Q &= 2 \int_0^{+\infty} (\hat{x}_1 - x)^2 p(x) dx = \hat{x}_1^2 \cdot \underbrace{2 \int_0^{+\infty} p(x) dx}_{=1} - 4\hat{x}_1 \cdot \underbrace{\int_0^{+\infty} xp(x) dx}_{=\sigma/\sqrt{2\pi}} + 2 \underbrace{\int_0^{+\infty} x^2 p(x) dx}_{=\sigma^2} = \\
 &= \hat{x}_1^2 - \frac{4\sigma}{\sqrt{2\pi}} \hat{x}_1 + \sigma^2. \tag{29}
 \end{aligned}$$

Thus, we obtain:

$$\hat{x}_{1\text{opt}} = \arg \min_{\hat{x}_1} P_Q = \frac{2\sigma}{\sqrt{2\pi}} = 0.798 \sigma,$$

which yields the minimum distortion equal to

$$P_{Q\text{min}} = \sigma^2 \left(1 - \frac{2}{\pi} \right).$$

This implies that the optimum SNR for scalar quantisation with 1 bit/sample is

$$\frac{SNR_{opt}}{\text{dB}} = 10 \log \frac{P_x}{P_{Q\text{min}}} = 10 \log \frac{\pi}{\pi - 2} = 4.3964 .$$

This result differs from formula (26), which would, for $R = 1$, yield the SNR of only $6.02 - 4.35 = 1.67$ dB. This is, however, not correct, because (26) is obtained assuming high rate, which $R = 1$ bit/sample is certainly not.

Problem 6 / Solution

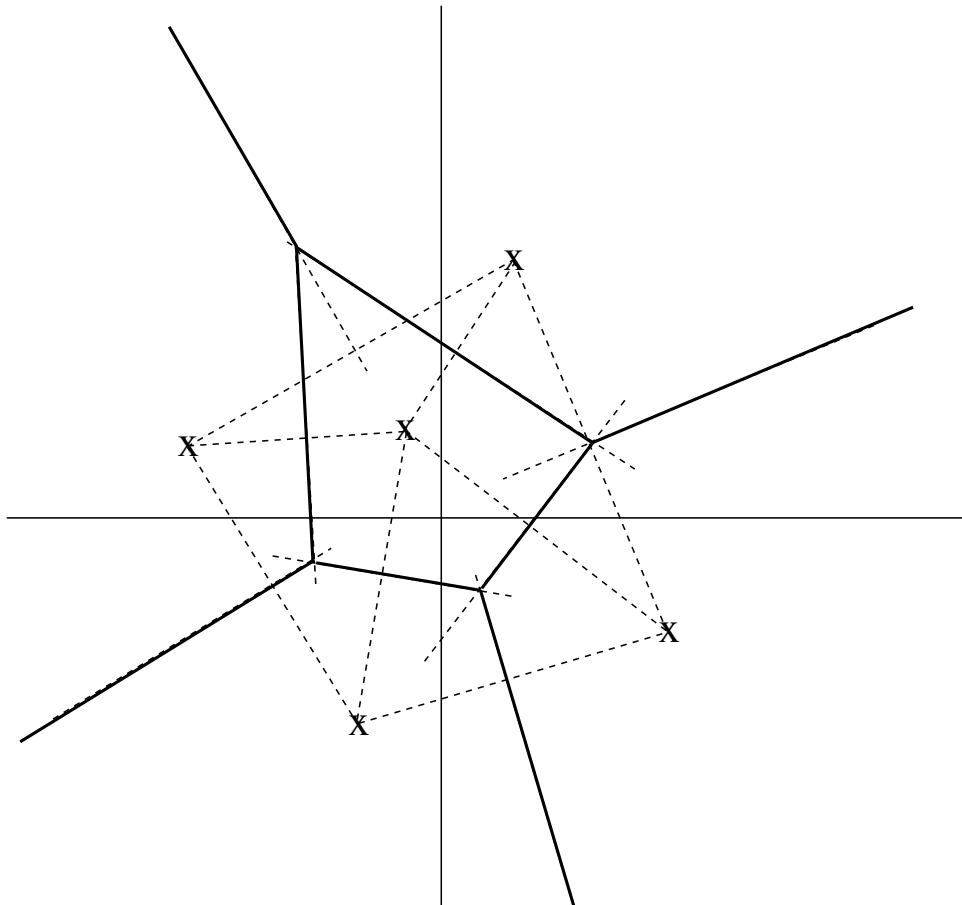
(a)

There are $M = 5$ two-dimensional code-vectors ($N = 2$) in the quantiser code-book. In order to represent them by a fixed rate, we need 3 bits per vector as $n = 3$ is the smallest integer such that $5 = M \leq 2^n$. Since $2^3 = 8$, we actually do not use $8 - 5 = 3$ codewords (i.e. 3 more vectors could be represented with the same number of bits). Thus, in terms of source coding, we are not using these 3 bits efficiently.

(b)

Voronoi region of a two-dimensional codevector c_i ($i = 1, \dots, 5$ in our case) is a set of points in \mathbb{R}^2 space that are closest to c_i in terms of Euclidean distance.

The Voronoi regions for the 5 codevectors considered in the example are shown in the figure in solid lines.



(c)

As 3 of 8 possibilities for encoding are not used, we can regard this index allocation (fixed rate with 3 bits per index) as an error protecting channel code. If a bit sequence is received that is not used for encoding, this must be due to bit errors. Thus bit errors can be detected. If the redundancy is large enough, then bit errors can also be corrected. This is not possible in this case, as the redundancy is lower than 1 bit per index.

(d)

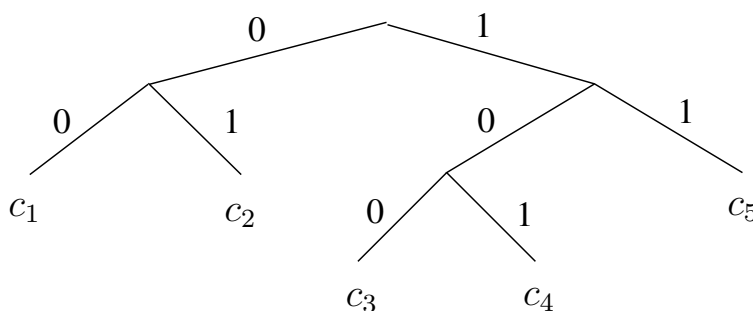
Entropy of the codevectors c_i

$$H(I) = - \sum_{i=1}^5 \Pr(c_i) \cdot \log_2 \Pr(c_i) = 5 \cdot \frac{1}{5} \log_2(5) = 2.32 \frac{\text{bits}}{\text{index}}$$

(e)

Although the codevectors are uniformly distributed, there is a gain by Huffman coding because 5 equiprobable events can't be packed efficiently into a *fixed* rate of 3 bits per event.

A possible Huffman code (for construction see lecture notes) is as follows:



Average word length:

$$\bar{R} = \sum_{i=1}^5 l_i \Pr(c_i) = \frac{1}{5} \cdot (2 + 2 + 3 + 3 + 2) = 2.4 \frac{\text{bits}}{\text{index}}$$

There is a strong gain by Huffman coding, reducing the average number of required bits per index from 3 for fixed rate to 2.4 for variable rate. This is also quite close to the entropy of 2.32 bits per index.

Problem 7 / Solution

(a)

In the $(m + 1)$ -st iteration of the LBG algorithm, the new, improved code-vectors $\hat{\underline{y}}_{(m+1), I}$ are found by minimising the total distortion of all the points $\underline{x}_k \in S_{m, I}$:

$$\hat{\underline{y}}_{(m+1), I} = \arg \min_{\underline{\hat{y}}} \underbrace{\sum_{\underline{x}_k \in S_{m, I}} d(\underline{x}_k, \underline{\hat{y}})}_{= D(S_{m, I})}, \quad (30)$$

where

$$d(\underline{x}_k, \underline{\hat{y}}) = \sum_{j=1}^N (x_k(j) - \hat{y}(j))^2. \quad (31)$$

The minimum is obtained in (30) by taking the first derivative of $D(S_{m, I})$ with respect to each component $\hat{y}(j)$, and forcing this derivative to be zero. We obtain

$$\frac{d}{d\hat{y}(j)} \sum_{\underline{x}_k \in S_{m, I}} \sum_{j'=1}^N (x_k(j') - \hat{y}(j'))^2 = -2 \sum_{\underline{x}_k \in S_{m, I}} (x_k(j) - \hat{y}(j)) \quad (32)$$

for $j = 1, 2, \dots, N$. Set to zero we find the optimal solution

$$\hat{y}_{(m+1), I}(j) = \frac{1}{|S_{m, I}|} \sum_{\underline{x}_k \in S_{m, I}} x_k(j). \quad (33)$$

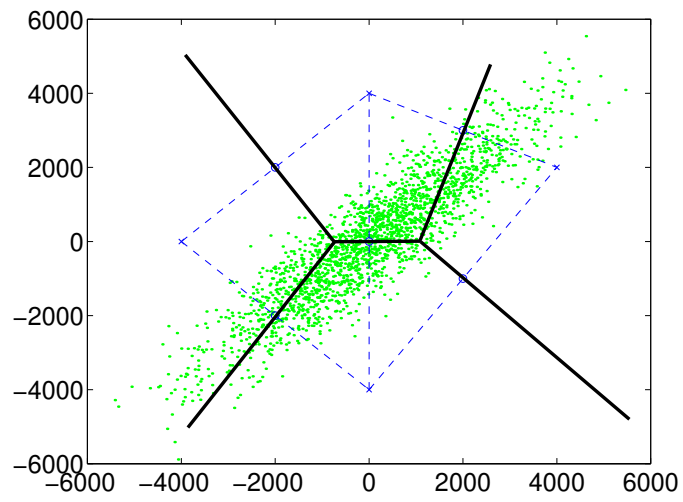
As this solution holds for all j , we can write it in more compact vector form according to

$$\hat{\underline{y}}_{(m+1), I} = \frac{1}{|S_{m, I}|} \sum_{\underline{x}_k \in S_{m, I}} \underline{x}_k, \quad (34)$$

with $|S_{m, I}|$ the number of members of the set $S_{m, I}$. Hence we obtain the “new” codevectors in step 4 of the LBG algorithm as the average of all vectors that were grouped into the set $S_{m, I}$ by the encoding step 2.

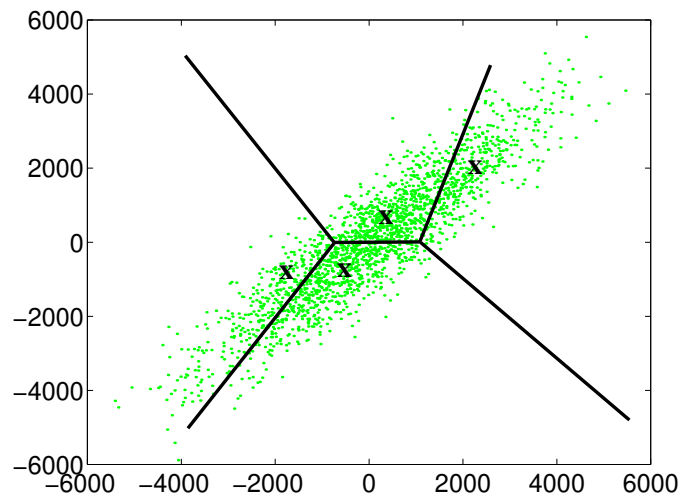
(b)

Voronoi regions of the given code-vectors:



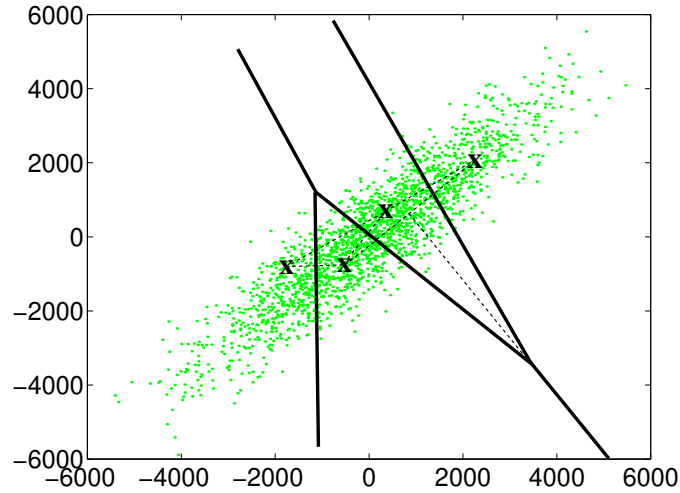
(c)

New code-vectors are found as centers of mass of the data points in each Voronoi region from the first iteration. Their approximate positions are marked in the figure below.



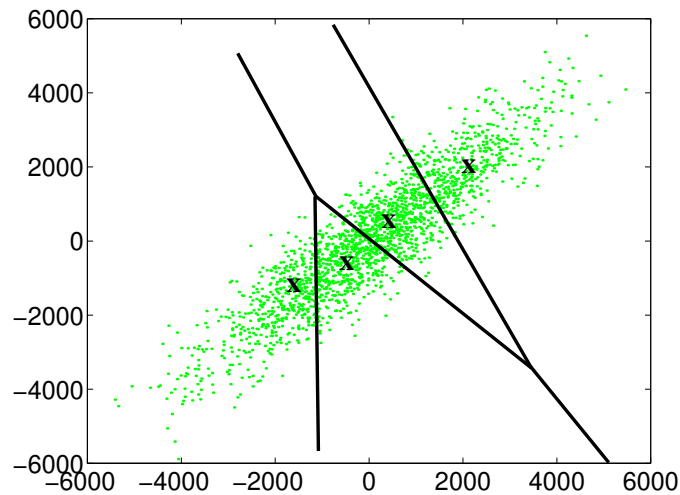
(d)

The Voronoi regions of the newly obtained code-vectors in c) are shown below (second iteration).



(e)

The new codevectors obtained from the Voronoi-regions in d) are shown in the figure below:



(f)

For a given codebook, the modified distance measure does not change the decision compared with the conventional MSE because the scaling factor $f()$ only depends on the input vector which is a constant with respect to the codebook search.

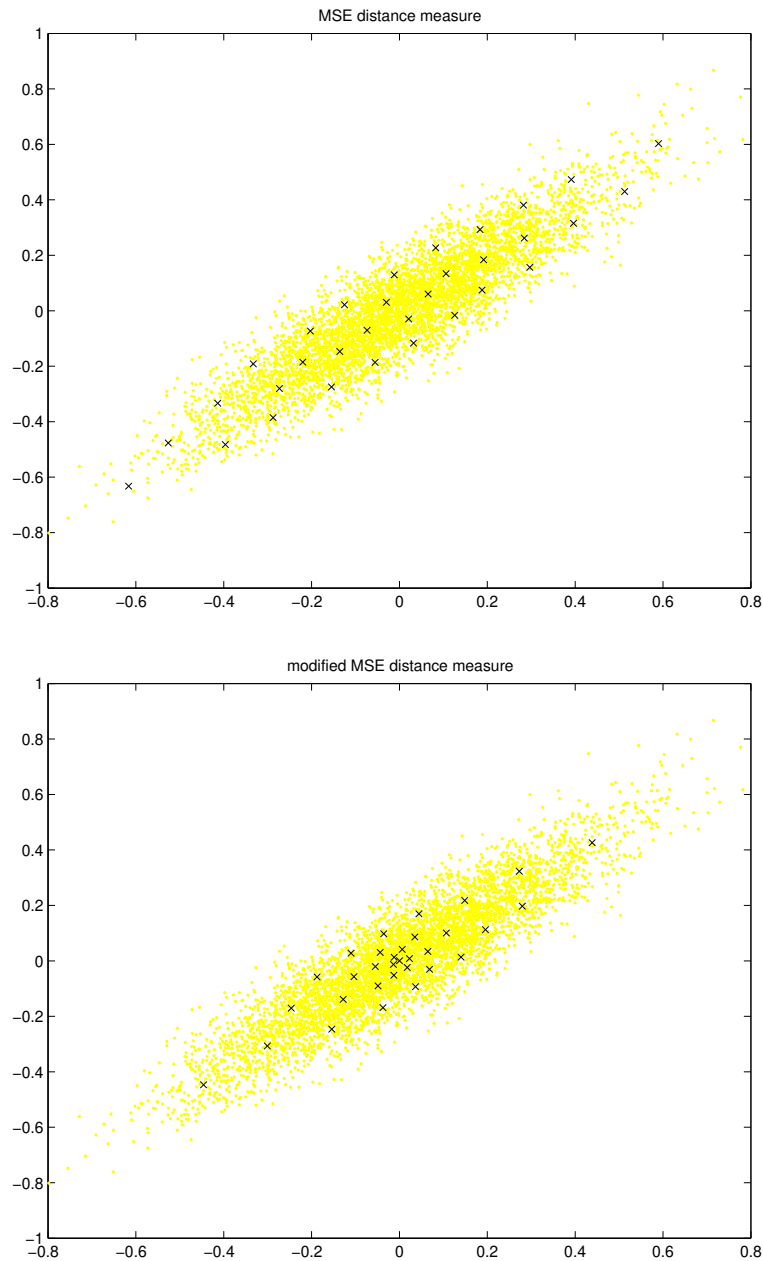
(g)

The function $f(\underline{x}_k)$ is constant with respect to a differentiation for the codevector components. Therefore $f(\underline{x}_k)$ appears just as a factor in (32). Hence, we obtain

$$\hat{\underline{y}}_{(m+1),I} = \frac{\sum_{\underline{x}_k \in S_{m,I}} \underline{x}_k \cdot f(\underline{x}_k)}{\sum_{\underline{x}_k \in S_{m,I}} f(\underline{x}_k)} \quad (35)$$

(h)

If we choose $f(\underline{x}_k) = 1/||\underline{x}_k||^2$ errors are scaled up when the source vector has small “power”. Therefore, the quantisation will be more accurate for small-length source vectors \underline{x}_k . As the distance measure is used in the LBG algorithm, this will change the codevector locations accordingly, compared with the conventional MSE distance measure. The two figures below illustrate that. They show the result of LBG codebook training for the MSE and the modified MSE distance measure, the latter according to (35) with $f(\underline{x}_k) = 1/||\underline{x}_k||^2$. A correlated Gaussian source was used.



Problem 8 / Solution

(a)

Information rate-distortion function for a source that generates signals x from a finite alphabet \mathcal{X} , according to the PDF $p(x)$, is defined as

$$R(D) = \min_{p(\tilde{x}|x): \mathbb{E}\{d(x, \tilde{x})\} \leq D} I(X; \tilde{X}), \quad (36)$$

where minimisation is performed over all conditional distributions $p(\tilde{x}|x)$, such that the average distortion per letter is bounded by D (for comments on D see remarks at the end):

$$\mathbb{E}\{d(x, \tilde{x})\} \triangleq \sum_{x \in \mathcal{X}} \sum_{\tilde{x} \in \tilde{\mathcal{X}}} d(x, \tilde{x}) p(\tilde{x}|x) p(x) \leq D, \quad (37)$$

where \tilde{x} are the reproducer values belonging to a finite reproducer alphabet $\tilde{\mathcal{X}}$.

The mutual information between the source and the reproducer values in the above formula is defined as:

$$I(X; \tilde{X}) = \sum_{x \in \mathcal{X}} \sum_{\tilde{x} \in \tilde{\mathcal{X}}} p(x, \tilde{x}) \log_2 \frac{p(x, \tilde{x})}{p(x)p(\tilde{x})} = \sum_{x \in \mathcal{X}} \sum_{\tilde{x} \in \tilde{\mathcal{X}}} p(\tilde{x}|x)p(x) \log_2 \frac{p(\tilde{x}|x)}{\sum_{x' \in \mathcal{X}} p(\tilde{x}|x')p(x')}. \quad (38)$$

Note that we minimise the mutual information $I(X, \tilde{X})$ over all conditional probability functions $p(\tilde{x}|x)$ for a reproduction \tilde{x} to be used to code a given source input x (both binary in our example).

(b)

The Bernoulli source generates letters x from a binary alphabet $\mathcal{X} = \{0, 1\}$, with probabilities $\Pr\{X = 1\} = p$ and $\Pr\{X = 0\} = 1 - p$. The quantiser reproducer values \tilde{x} are from the same alphabet $\tilde{\mathcal{X}} = \mathcal{X}$. The (per-letter) distortion measure suitable for this model is *Hamming distortion*, defined as:

$$d(x, \tilde{x}) \triangleq \begin{cases} 0, & x = \tilde{x} \\ 1, & x \neq \tilde{x}. \end{cases}$$

In order to derive a lower bound on the mutual information $I(X; \tilde{X})$, we start from the expression:

$$I(X; \tilde{X}) = H(X) - H(X|\tilde{X}), \quad (39)$$

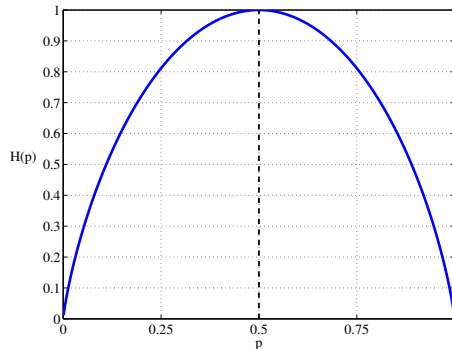
where

$$H(X) = H_2(p) = -p \log_2(p) - (1 - p) \log_2(1 - p)$$

is the binary entropy function (entropy of a source with $X \sim \text{Bernoulli}(p)$). An important property of the binary entropy function is its *symmetry* with respect to $p = 1/2$, i.e., $H_2(p) = H_2(1 - p)$, as shown in Figure 4. Thus, we can, without loss of generality, assume in the further derivations that $p \leq 1/2$.

In the sequel, it will prove useful to introduce a new variable Y , defined as an "error indicator":

$$Y \triangleq X \oplus \tilde{X} = \begin{cases} 0, & \text{if } x = \tilde{x} \\ 1, & \text{if } x \neq \tilde{x} \end{cases}, \quad (40)$$

Figure 4: Binary entropy function $H_2(p)$

where \oplus denotes modulo-2 addition. From this definition it follows

$$\Pr\{Y = 0\} = \Pr\{X = \tilde{X}\} \quad \text{and} \quad \Pr\{Y = 1\} = 1 - \Pr\{X = \tilde{X}\} = \Pr\{X \neq \tilde{X}\},$$

which yields the entropy of Y

$$H(Y) = H_2(\Pr\{X = \tilde{X}\}) = H_2(\Pr\{X \neq \tilde{X}\}). \quad (41)$$

On the other hand, from the definition of the Hamming distortion, it follows that the average Hamming distortion equals

$$\mathbb{E}\{d(x, \tilde{x})\} = \sum_{x, \tilde{x} \in \mathcal{X}} d(x, \tilde{x})p(x, \tilde{x}) = \sum_{\substack{x, \tilde{x} \in \mathcal{X} \\ x = \tilde{x}}} 0 \cdot p(x, \tilde{x}) + \sum_{\substack{x, \tilde{x} \in \mathcal{X} \\ x \neq \tilde{x}}} 1 \cdot p(x, \tilde{x}) = \sum_{\substack{x, \tilde{x} \in \mathcal{X} \\ x \neq \tilde{x}}} p(x, \tilde{x}) = \Pr\{X \neq \tilde{X}\}. \quad (42)$$

Applying the constraint (37) on the average distortion given by (42) we obtain

$$\mathbb{E}\{d(x, \tilde{x})\} = \Pr\{X \neq \tilde{X}\} \leq D. \quad (43)$$

Since the binary entropy function $H_2(p)$ is monotonically increasing for $0 \leq p \leq 1/2$, this means that (43) implies

$$H_2(\Pr\{X \neq \tilde{X}\}) \leq H_2(D), \quad \text{for } 0 \leq D \leq 1/2. \quad (44)$$

Now we have all the prerequisites to derive a lower bound on the mutual information $I(X; \tilde{X})$, under the constraint on the average distortion (43):

$$\begin{aligned} I(X; \tilde{X}) &= H(X) - H(X|\tilde{X}) \\ &= H_2(p) - H(X \oplus \tilde{X}|\tilde{X}) \\ &\stackrel{(40)}{=} H_2(p) - H(Y|\tilde{X}) \\ &\stackrel{*}{\geq} H_2(p) - H(Y) \quad (* \text{ conditioning reduces entropy : } H(Y|\tilde{X}) \leq H(Y)) \\ &\stackrel{(41)}{=} H_2(p) - H_2(\Pr\{X \neq \tilde{X}\}) \\ &\stackrel{(44)}{\geq} H_2(p) - H_2(D). \end{aligned} \quad (45)$$

This bound is non-trivial only for $0 \leq D \leq p$. For $\frac{1}{2} \geq D > p$, (45) is a negative number, which is a trivial bound for mutual information, which is, by definition, always non-negative. Thus, we can write:

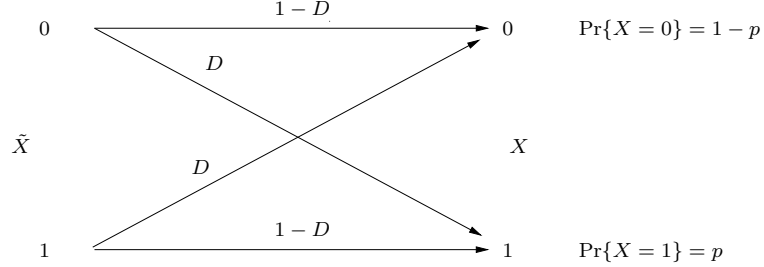
$$I(X; \tilde{X}) \geq \begin{cases} H_2(p) - H_2(D), & 0 \leq D \leq p \\ 0, & D > p. \end{cases} \quad (46)$$

Recall that, so far, we have assumed that $p \leq \frac{1}{2}$, due to symmetry of $H_2(p)$. If we wish to write (46) in the form applicable for all values $0 \leq p \leq 1$, we simply need to replace p in (46) with $\min\{p, 1 - p\}$.

(c)

Now, we need to prove that there exists a system that achieves a lower bound given by (46). Thus, we will prove that rate-distortion function $R(D)$ is given by expression (46).

We consider a *binary symmetric test-channel*, with transition probability equal to D , and $D \leq p$.

Figure 5: BSC test channel with transition probability D

$$\begin{aligned}
 \underbrace{\Pr\{X=1\}}_{=p} &= \underbrace{\Pr\{X=1|\tilde{X}=0\}}_{=D} \cdot \Pr\{\tilde{X}=0\} + \underbrace{\Pr\{X=1|\tilde{X}=1\}}_{=1-D} \cdot \Pr\{\tilde{X}=1\} \\
 p &= D \cdot (1 - \Pr\{\tilde{X}=1\}) + (1 - D) \cdot \Pr\{\tilde{X}=1\} \\
 p &= D - D \cdot \Pr\{\tilde{X}=1\} + \Pr\{\tilde{X}=1\} - D \cdot \Pr\{\tilde{X}=1\} \\
 p &= (1 - 2D) \cdot \Pr\{\tilde{X}=1\} + D
 \end{aligned}$$

$$\Pr\{\tilde{X}=1\} = \frac{p-D}{1-2D}, \quad \Pr\{\tilde{X}=0\} = \frac{1-p-D}{1-2D}. \quad D \leq p \quad (47)$$

$$\begin{aligned}
 H(X|\tilde{X}=0) &= -\Pr\{X=0|\tilde{X}=0\} \log_2(\Pr\{X=0|\tilde{X}=0\}) - \Pr\{X=1|\tilde{X}=0\} \log_2(\Pr\{X=1|\tilde{X}=0\}) \\
 &= -(1-D) \log_2(1-D) - D \log_2(D) \\
 &= H_2(D),
 \end{aligned}$$

$$H(X|\tilde{X}=1) = H_2(D), \quad \text{analogously.}$$

$$\begin{aligned}
 H(X|\tilde{X}) &\stackrel{(\#)}{=} H(X|\tilde{X}=0) \cdot \Pr\{\tilde{X}=0\} + H(X|\tilde{X}=1) \cdot \Pr\{\tilde{X}=1\} \\
 &= H_2(D) \cdot (\Pr\{\tilde{X}=0\} + \Pr\{\tilde{X}=1\}) \\
 &= H_2(D).
 \end{aligned}$$

(#): Note that generally $H(X|\tilde{X}) = \sum_{\tilde{x} \in \tilde{\mathcal{X}}} p(\tilde{x}) \cdot H(X|\tilde{X} = \tilde{x})$ with

$$H(X|\tilde{X} = \tilde{x}) = \sum_{x \in \mathcal{X}} p(x|\tilde{x}) \log_2(p(x|\tilde{x}))$$

$$I(X; \tilde{X}) = H(X) - H(X|\tilde{X}) = H_2(p) - H_2(D) \quad \Longrightarrow \quad \text{lower bound is achieved}$$

$$\begin{aligned} \mathbb{E}\{d(x, \tilde{x})\} &= \Pr\{X \neq \tilde{X}\} = \Pr\{X = 0|\tilde{X} = 1\} \Pr\{\tilde{X} = 1\} + \Pr\{X = 1|\tilde{X} = 0\} \Pr\{\tilde{X} = 0\} \\ &= D \cdot \frac{p-D}{1-2D} + D \cdot \frac{1-p-D}{1-2D} = \frac{D(p-D+1-p-D)}{1-2D} = \frac{D(1-2D)}{1-2D} \\ &= D \quad \Longrightarrow \quad \text{distortion constraint is met with equality} \end{aligned}$$

If $D > p$, instead of (47), we assign the following probabilities: $\Pr\{\tilde{X} = 1\} = 0$, and $\Pr\{\tilde{X} = 0\} = 1$, which means that $\tilde{X} = 0$ is deterministic. Then,

$$H(X|\tilde{X}) = H(X) \quad \Longrightarrow \quad I(X; \tilde{X}) = 0,$$

which again corresponds to the lower bound (46). The average distortion in this case is

$$\mathbb{E}\{d(x, \tilde{x})\} = \Pr\{X \neq \tilde{X}\} = \Pr\{X \neq 0\} = \Pr\{X = 1\} = p \leq D,$$

according to initial assumption. Thus, the constraint (43) is fulfilled.

We have shown that the considered system (BSC "test-channel") achieves the lower bound of the mutual information with fulfilled average-distortion constraint. Thus, we conclude that the rate-distortion function for a Bernoulli source and Hamming distortion measure is given by

$$R(D) = \begin{cases} H_2(p) - H_2(D), & 0 \leq D \leq \min\{p, (1-p)\} \\ 0, & D > \min\{p, (1-p)\}. \end{cases} \quad (48)$$

Remarks:

In this example the distortion D can be thought of as the relative number of bits being reproduced incorrectly (i.e., a "1" by "0" and a "0" by "1") in a very long block of source encoded bits. As an example consider the block of bits "0 0 0 0 0 0 1 0 0 0 0" to be transmitted. One could do this, e.g., by transmitting the information "7 zeros", "1 one", "4 zeros" (this is run-length coding used e.g. in fax machines). If however, we allow for some distortion, e.g., $D \leq 1/12$, we will be tempted to reproduce the "one"-bit in the sequence by a "zero" (this is an error causing "distortion") so that we can encode the bit stream by "12 zeros": of course the latter requires lower bit rate, i.e., we have traded more "distortion" against lower bitrate: exactly this tradeoff is described by the rate-distortion function.

Note that classical lossless compression theory would result from fixing the allowed distortion to $D = 0$. In this case, the conditional entropy $H(X|\tilde{X}) = 0$ because in lossless compression we resolve the uncertainty about the source X completely as the encoding is lossless and hence $\tilde{X} = X$. The rate-distortion function (48) is still true as then $I(X; \tilde{X}) = H(X) - H(X|\tilde{X}) = H(X)$ and therefore the entropy of the source determines the required bit rate alone. Note that for $D = 0$ we have $H_2(D = 0) = 0$ with leaves us in (48) with $R(D = 0) = H_2(p) = H(X)$, which says we can achieve the entropy of the source by coding. Of course we knew that before; lossless compression can be seen as a special case of more general lossy source coding.

Problem 9 / Solution

(a)

An autocorrelated Gaussian signal can be generated by filtering the output of the memoryless Gaussian source by a linear filter. The *spectral flatness measure* (SFM), γ_x^2 , of the output signal is defined by

$$\gamma_x^2 \doteq \exp \left(\frac{1}{2\pi} \int_{-\pi}^{+\pi} \ln \Phi_{xx}(\Omega) d\Omega \right) \bigg/ \underbrace{\frac{1}{2\pi} \int_{-\pi}^{+\pi} \Phi_{xx}(\Omega) d\Omega}_{=\sigma_x^2}, \quad (49)$$

where σ_x^2 is the signal power and $\Phi_{xx}(\Omega)$ is its power spectral density (PSD). The SFM takes values in the range $0 < \gamma_x^2 \leq 1$, and it is a measure of correlation in the signal, i.e. a measure of flatness of its PSD. If $\gamma_x^2 = 1$, the PSD of the signal is completely flat, i.e., the signal is *white*, i.e., uncorrelated. If $\gamma_x^2 = 0$, the PSD of the signal is a Dirac peak, i.e., the signal is completely correlated. Thus, in general, small values of γ_x^2 indicate strong correlation between signal samples (in this case vector quantisation yields larger SNR gain), whereas large values of γ_x^2 indicate weak correlation between the signal samples, and thus, the SNR gain is smaller.

(b)

The linear filter used for obtaining correlated Gaussian signal has the transfer function:

$$H(z) = \frac{1 + \sum_{i=1}^{N_z} b_i z^{-i}}{1 + \sum_{k=1}^{N_p} a_k z^{-k}} = z^{N_p - N_z} \frac{\sum_{i=0}^{N_z} b_i z^{N_z - i}}{\sum_{k=0}^{N_p} a_k z^{N_p - k}},$$

where N_z and N_p are the numbers of zeros and poles, respectively, and $a_0 = b_0 = 1$.

The polynomials in the numerator and the denominator of the transfer function can be represented as functions of their roots (zeros $z_{0,i}$ and poles $z_{\infty,k}$, respectively) by the factorisation:

$$H(z) = \frac{\prod_{i=1}^{N_z} (1 - z_{0,i} \cdot z^{-1})}{\prod_{k=1}^{N_p} (1 - z_{\infty,k} \cdot z^{-1})} = z^{N_p - N_z} \frac{\prod_{i=1}^{N_z} (z - z_{0,i})}{\prod_{k=1}^{N_p} (z - z_{\infty,k})}.$$

(c)

The filter's frequency response function $H(\Omega)$ can be obtained from the transfer function by substituting $z = e^{j\Omega}$, $z_{0,i} = r_{0,i} e^{j\Omega_{0,i}}$ and $z_{\infty,k} = r_{\infty,k} e^{j\Omega_{\infty,k}}$:

$$H(\Omega) = \frac{\prod_{i=1}^{N_z} (1 - r_{0,i} \cdot e^{j(\Omega_{0,i} - \Omega)})}{\prod_{k=1}^{N_p} (1 - r_{\infty,k} \cdot e^{j(\Omega_{\infty,k} - \Omega)}}.$$

The conditions of stability and phase-minimality ensure that both poles and zeros lie inside the unit-circle in the z -plane, i.e., that $r_{\infty,k} < 1$ and $r_{0,i} < 1$, $\forall i, k$.

The square of the magnitude response of the considered filter is given by:

$$\begin{aligned}
|H(\Omega)|^2 = H(\Omega) \cdot H^*(\Omega) &= \frac{\prod_{i=1}^{N_z} [(1 - r_{0,i} \cdot e^{j(\Omega_{0,i} - \Omega)}) (1 - r_{0,i} \cdot e^{-j(\Omega_{0,i} - \Omega)})]}{\prod_{k=1}^{N_p} [(1 - r_{\infty,k} \cdot e^{j(\Omega_{\infty,k} - \Omega)}) (1 - r_{\infty,k} \cdot e^{-j(\Omega_{\infty,k} - \Omega)})]} \\
&= \frac{\prod_{i=1}^{N_z} (1 - 2r_{0,i} \cos(\Omega - \Omega_{0,i}) + r_{0,i}^2)}{\prod_{k=1}^{N_p} (1 - 2r_{\infty,k} \cos(\Omega - \Omega_{\infty,k}) + r_{\infty,k}^2)}. \tag{50}
\end{aligned}$$

(d)

Let σ_r^2 denote the variance of the uncorrelated Gaussian signal, before filtering. Then, due to the linearity of the filter, the power spectral density of the correlated (coloured) Gaussian signal at the output of the filter is given by:

$$\Phi_{xx}(\Omega) = \sigma_r^2 \cdot |H(\Omega)|^2, \tag{51}$$

and its variance is then

$$\sigma_x^2 = \frac{1}{2\pi} \int_{-\pi}^{+\pi} \Phi_{xx}(\Omega) d\Omega = \sigma_r^2 \frac{1}{2\pi} \int_{-\pi}^{+\pi} |H(\Omega)|^2 d\Omega. \tag{52}$$

By inserting (51) and (52) into the expression (49) for spectral flatness measure, we obtain

$$\gamma_x^2 = \exp \left(\frac{1}{2\pi} \int_{-\pi}^{+\pi} \ln (\sigma_r^2 \cdot |H(\Omega)|^2) d\Omega \right) / \sigma_x^2 = \frac{\sigma_r^2}{\sigma_x^2} \cdot \underbrace{\exp \left(\frac{1}{2\pi} \int_{-\pi}^{+\pi} \ln (|H(\Omega)|^2) d\Omega \right)}_{=I \doteq 0} = \frac{\sigma_r^2}{\sigma_x^2}. \tag{53}$$

Now we will show that the integral I in the above equation is indeed equal to zero.

By inserting (50) into the integral we obtain:

$$\begin{aligned}
I &= \frac{1}{2\pi} \int_{-\pi}^{+\pi} \left[\ln \left(\prod_{i=1}^{N_z} (1 - 2r_{0,i} \cos(\Omega - \Omega_{0,i}) + r_{0,i}^2) \right) - \ln \left(\prod_{k=1}^{N_p} (1 - 2r_{\infty,k} \cos(\Omega - \Omega_{\infty,k}) + r_{\infty,k}^2) \right) \right] d\Omega \\
&= \frac{1}{2\pi} \sum_{i=1}^{N_z} \int_{-\pi}^{+\pi} \ln (1 - 2r_{0,i} \cos(\Omega - \Omega_{0,i}) + r_{0,i}^2) d\Omega - \frac{1}{2\pi} \sum_{k=1}^{N_p} \int_{-\pi}^{+\pi} \ln (1 - 2r_{\infty,k} \cos(\Omega - \Omega_{\infty,k}) + r_{\infty,k}^2) d\Omega
\end{aligned}$$

We can expand the cosine terms in the above expression as

$$\begin{aligned}
\cos(\Omega - \Omega_{0,i}) &= \cos \Omega \cos \Omega_{0,i} - \sin \Omega \sin \Omega_{0,i}, \\
\cos(\Omega - \Omega_{\infty,k}) &= \cos \Omega \cos \Omega_{\infty,k} - \sin \Omega \sin \Omega_{\infty,k}. \tag{54}
\end{aligned}$$

If all poles and zeros are real, this implies that $\Omega_{0,i} = k_1\pi$, and $\Omega_{\infty,k} = k_2\pi$, where $k_1, k_2 = 0, 1, 2, \dots$. Thus, we obtain

$$\left. \begin{aligned} \sin \Omega_{0,i} &= 0, & \cos \Omega_{0,i} &= \pm 1 \\ \sin \Omega_{\infty,k} &= 0, & \cos \Omega_{\infty,k} &= \pm 1 \end{aligned} \right\} \implies \left\{ \begin{aligned} \cos(\Omega - \Omega_{0,i}) &= \pm \cos \Omega \\ \cos(\Omega - \Omega_{\infty,k}) &= \pm \cos \Omega \end{aligned} \right. \tag{55}$$

Equivalently, we can consider that $\cos(\Omega_{0,i}) = \cos(\Omega_{\infty,k}) = 1$, and incorporate the sign into the magnitudes $r_{0,i}$ and $r_{\infty,k}$, respectively. Thus, the magnitudes can be negative, but due to stability and phase-minimality conditions: $|r_{0,i}| < 1$ and $|r_{\infty,k}| < 1$.

By incorporating (54) and (55) into the integral expression, and by exploiting the symmetry of the cosine function, we obtain

$$I = \frac{1}{2\pi} \cdot 2 \sum_{i=1}^{N_z} \int_0^{+\pi} \ln(1 - 2r_{0,i} \cos \Omega + r_{0,i}^2) d\Omega - \frac{1}{2\pi} \cdot 2 \sum_{k=1}^{N_p} \int_0^{+\pi} \ln(1 - 2r_{\infty,k} \cos \Omega + r_{\infty,k}^2) d\Omega.$$

Now we use the formula from the hint, and since $0 < |r_{0,i}| < 1$ and $0 < |r_{\infty,k}| < 1$, then

$$\int_0^{+\pi} \ln(1 - 2r_{0,i} \cos \Omega + r_{0,i}^2) d\Omega = \int_0^{+\pi} \ln(1 - 2r_{\infty,k} \cos \Omega + r_{\infty,k}^2) d\Omega = 2\pi \ln(1) = 0.$$

Thus, the integral I is

$$I = \frac{1}{\pi} \sum_{i=1}^{N_z} 0 - \frac{1}{\pi} \sum_{k=1}^{N_p} 0 = 0 \quad \implies \quad \exp(I) = 1. \quad (56)$$

Note that this result is also valid when the zeros and poles are complex numbers. The assumption on them being real just simplifies the calculation due to (55).

Now, we return to (53) and obtain the inverse of the spectral flatness measure

$$\frac{1}{\gamma_x^2} = \frac{\sigma_x^2}{\sigma_r^2} = \frac{1}{2\pi} \int_{-\pi}^{+\pi} |H(\Omega)|^2 d\Omega = \sum_{k=0}^{\infty} h_0^2(k), \quad (57)$$

where $h_0(k)$ is the discrete-time impulse response of the linear filter.

(e)

The rightmost equality in the expression (57) follows from the *Parseval's theorem*.

Problem 10 / Solution

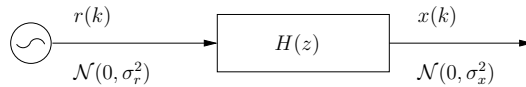


Figure 6: Gaussian source with memory

We consider a memoryless Gaussian source that emits samples $r(k) \sim \mathcal{N}(0, \sigma_r^2)$. These samples are filtered by a linear second-order filter with the transfer function

$$H(z) = \frac{z^2}{(z-a)(z-b)}, \quad a = \frac{1}{4}, \quad b = \frac{3}{4}. \quad (58)$$

The samples at the output of the filter, $x(k)$, still have the Gaussian distribution, but with the different variance, σ_x^2 . Moreover, samples are correlated, i.e., the signal is coloured by the filter's transfer function.

(a)

The maximum signal-to-noise ratio that can be achieved by coding the uncorrelated Gaussian source samples $r(k)$, with any coding scheme, at some rate R , is given by:

$$\frac{SNR}{\text{dB}} = 10 \log \frac{P_r}{D_r(R)} = 10 \log \frac{\sigma_r^2}{2^{-2R} \cdot \sigma_r^2} = 10 \log(2^{2R}) = 6.02 \cdot R. \quad (59)$$

(b)

The above result is true for *any* rate, low or high, because it is obtained from analytical derivation of the distortion-rate function $D(R)$ of a Gaussian memoryless source. Derivation of the exact formula $D_r(R) = 2^{-2R} \cdot \sigma_r^2$ does not require any approximations or assumptions on the rate. This is, however, not true for derivation of the distortion-rate function for the correlated Gaussian source.

(c)

In order to compute the impulse response of the filter, we first note that the transfer function can be written as follows:

$$H(z)/z = \frac{A}{z-a} + \frac{B}{z-b}$$

For the coefficients we obtain $A = \frac{a}{a-b}$ and $B = \frac{b}{b-a}$ so that

$$H(z) = \frac{a}{a-b} \cdot \frac{z}{z-a} + \frac{b}{b-a} \cdot \frac{z}{z-b}$$

In the time-domain this equals the impulse response

$$h_0(k) = Z^{-1}\{H(z)\} = \left(\frac{a}{a-b} a^k + \frac{b}{b-a} b^k \right) \cdot \gamma_{-1}(k)$$

with $\gamma_{-1}(k)$ the discrete time step function.

This can be written more compact:

$$h_0(k) = \frac{b^{k+1} - a^{k+1}}{b - a} \cdot \gamma_{-1}(k)$$

For the given values $a = 1/4$, $b = 3/4$ and $k \geq 0$, we obtain:

$$h_0(k) = \frac{3^{k+1} - 1^{k+1}}{\frac{1}{2} \cdot 4^{k+1}} = \frac{3^{k+1} - 1}{\frac{1}{2} \cdot 2^{2k+2}} = \frac{3^{k+1} - 1}{2^{2k+1}} \quad (60)$$

The impulse response is shown in Figure 7. The taps are: $\{1, 1, 0.8125, 0.625, 0.4727, 0.3555, \dots\}$.

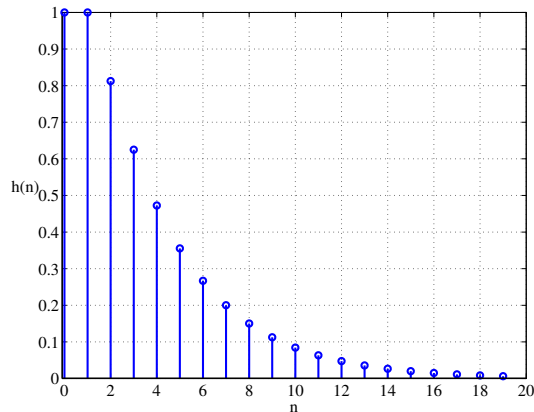


Figure 7: The first 20 taps of the filter's impulse response $h_0(k)$

(d)

The inverse of the spectral flatness measure of the coloured Gaussian signal $x(k)$ equals:

$$\begin{aligned} \frac{1}{\gamma_x^2} &= \sum_{k=0}^{\infty} h_0^2(k) = \sum_{k=0}^{\infty} \left(\frac{b^{k+1} - a^{k+1}}{b - a} \right)^2 = \\ &= \frac{1}{(b - a)^2} \cdot \sum_{k=0}^{\infty} \left(b^{2k+2} - 2(ab)^{k+1} + a^{2k+2} \right) = \\ &= \frac{1}{(b - a)^2} \cdot \left(b^2 \sum_{k=0}^{\infty} (b^2)^k - 2ab \sum_{k=0}^{\infty} (ab)^k + a^2 \sum_{k=0}^{\infty} (a^2)^k \right) = \\ &= \frac{1}{(b - a)^2} \cdot \left(\frac{b^2}{1 - b^2} - \frac{2ab}{1 - ab} + \frac{a^2}{1 - a^2} \right). \end{aligned} \quad (61)$$

The above derivation is possible only for stable filters, i.e., for $a < 1$ and $b < 1$. For the given values $a = 1/4$ and $b = 3/4$ we obtain

$$\frac{1}{\gamma_x^2} = 4 \cdot \left(\frac{9}{7} - \frac{6}{13} + \frac{1}{15} \right) \approx 3.5634.$$

(e)

The maximum SNR achievable by coding the samples $x(k)$ by any coding scheme at some *high* rate R is

$$\frac{SNR}{\text{dB}} = 10 \log \frac{P_x}{D_x(R)} = 10 \log \frac{\sigma_x^2}{2^{-2R} \cdot \sigma_x^2 \gamma_x^2} = 10 \log(2^{2R}) + 10 \log \left(\frac{1}{\gamma_x^2} \right) = 6.02 \cdot R + 5.52. \quad (62)$$

By comparing this with the SNR when coding the uncorrelated signal $r(k)$, given by (59), we observe that correlation of the signal $x(k)$ yields a gain of $-10 \log(\gamma_x^2) = 5.52$ dB, for high rate.

(f)

A unique property of the Gaussian distribution is that after applying any kind of *linear operator* on a Gaussian process, the output is still Gaussian. Thus, the samples $x(k)$ at the output of the linear filter are also Gaussian distributed, with the variance equal to

$$\sigma_x^2 = \sigma_r^2 \frac{1}{2\pi} \int_{-\pi}^{+\pi} |H(\Omega)|^2 d\Omega,$$

where the squared magnitude of the filter's frequency response is

$$|H(\Omega)|^2 = \frac{1}{(e^{j\Omega} - a)(e^{-j\Omega} - a)(e^{j\Omega} - b)(e^{-j\Omega} - b)} = \frac{1}{(1 - 2a \cos \Omega + a^2)(1 - 2b \cos \Omega + b^2)}.$$

(g)

The maximum SNR achievable by scalar fixed-rate quantisation of the signal $x(k)$ without exploiting its correlation, equals for some high rate R :

$$\frac{SNR}{\text{dB}} = 6.02 \cdot R - 4.35. \quad (63)$$

By comparing this results with (62) we observe a loss of $4.35 + 5.52 = 9.87$ dB, which is due to using scalar quantiser (instead of VQ) and due to not exploiting the signal correlation.

(h)

To reduce the loss of a scalar, fixed-rate quantisation, we should use a coding system that exploits the correlation of the signal $x(k)$. These methods are called *Analysis-by-Synthesis* or *Backward Prediction* or *Differential PCM* (which is in principle all the same).

The method of Analysis-by-Synthesis would achieve at high rate

$$\frac{SNR}{\text{dB}} = 6.02 \cdot R - 4.35 + 10 \log \left(\frac{1}{\gamma_x^2} \right) = 6.02 \cdot R + 1.17,$$

which means that it can obtain a full gain of $-10 \log(\gamma_x^2) = 5.57$ dB over the scalar scheme that ignores the correlation, whose SNR is given by (63). In practice most of this gain due to source correlation is also achieved in systems with relatively low bit rate.