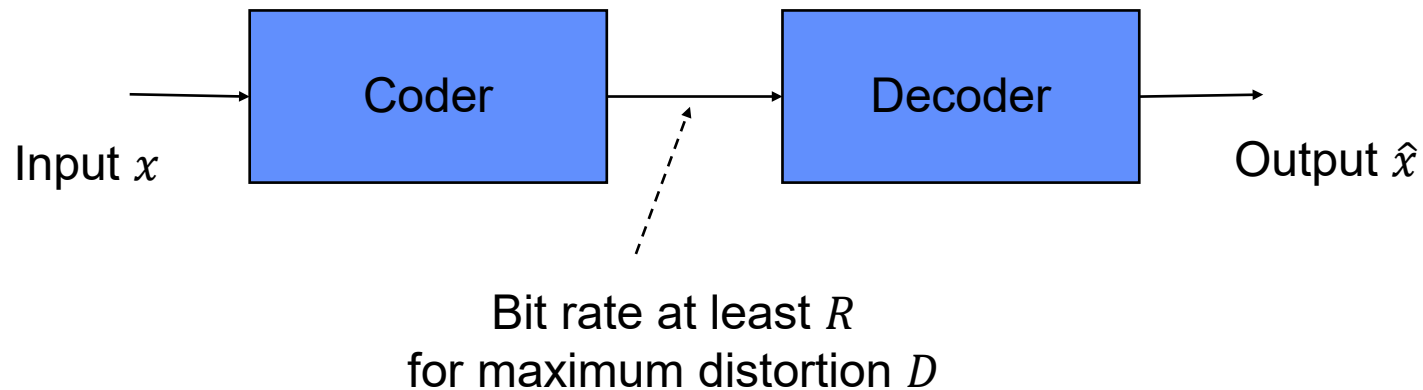


5 Quantization

- 5.1 Basics of Rate Distortion Theory
- 5.2 Scalar Quantization
- 5.3 Lloyd-Max Quantization
- 5.4 Entropy Coded Scalar Quantization
- 5.5 Embedded Quantization
- 5.6 Adaptive Quantization
- 5.7 Vector Quantization

5.1 Basics of Rate Distortion Theory

Rate distortion theory calculates the minimum transmission bit rate R for a required signal quality



Distortion (dt. *Verzerrung*): maximum average distortion D allowed, measured according to suitable criterion

Generality: Results of rate distortion theory are obtained without consideration of a specific coding method

Distortion

Assumption: symbol x sent, symbol \hat{x} received

Distortion is non-negative

$$d(x, \hat{x}) \geq 0$$

$$\text{and} \quad d(x, \hat{x}) = 0 \quad \text{for} \quad x = \hat{x}$$

Average distortion calculated with help of joint probability mass function

$$D = E\{d(x, \hat{x})\} = \sum_x \sum_{\hat{x}} p_{X,\hat{X}}(x, \hat{x}) d(x, \hat{x})$$

Subjective perception of images and video

- Distortion D may take subjective visual impression into account
- Obtained e.g. by extensive subjective visual tests
- No widely accepted measures available

Distortion Measures for Images

Given:	Original signal	$x[m, n]$
	Reconstructed signal	$\hat{x}[m, n]$
	Error signal	$e[m, n] = x[m, n] - \hat{x}[m, n]$

Mean squared error (MSE): Expectation value of the error signal

$$E\{e^2[m, n]\} = E\{(x[m, n] - \hat{x}[m, n])^2\} = \frac{1}{MN} \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} (x[m, n] - \hat{x}[m, n])^2$$

If the error signal is zero-mean, the mean squared error is equal to the variance of the error signal:

$$\sigma_e^2 = P_e - \mu_e^2 = E\{e^2[m, n]\} \quad \text{for } \mu_e = 0$$

Power of the original signal: $P_x = E\{x^2[m, n]\}$

Signal to Noise Ratio (SNR)

$$\text{SNR[dB]} = 10 \log_{10} \frac{P_x}{\sigma_e^2}$$

Peak Signal to Noise Ratio

Signal to Noise Ratio depends on the mean value of the original signal

⇒ not desired for video and image signals

Alternative: Reference to the maximum amplitude $A = 2^b - 1$ of the original signal, e. g. $A = 255$ for $b = 8$ bit per sample

Peak Signal to Noise Ratio (PSNR)

$$\text{PSNR}_{\text{image}}[\text{dB}] = 10 \log_{10} \frac{A^2}{\sigma_e^2}$$

The PSNR is always greater than zero, because A is the maximum difference between two arbitrary images $x[m, n]$ and $y[m, n]$.

For **video signals** with K images in a sequence $x[m, n, k]$ with time axis k the above considerations apply similarly with

$$P_x = E\{x^2[m, n, k]\}, \quad \sigma_e^2 = E\{e^2[m, n, k]\} = \frac{1}{MNK} \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} \sum_{k=0}^{K-1} (x[m, n, k] - \hat{x}[m, n, k])^2$$

Mean PSNR of a Video Sequence

In a video sequence the total quality is often calculated by averaging the PSNR values over all K images of a sequence. Regarding the mean error of an image k

$$\sigma_e^2[k] = \frac{1}{MN} \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} e^2[m, n, k]$$

the **mean PSNR of a video sequence** results in

$$\begin{aligned} \overline{\text{PSNR}_{\text{image}}} &= \frac{1}{K} \sum_{k=0}^{K-1} 10 \log \frac{A^2}{\sigma_e^2[k]} = 10 \log \left(\prod_{k=0}^{K-1} \frac{A^2}{\sigma_e^2[k]} \right)^{\frac{1}{K}} = -10 \log \left(\prod_{k=0}^{K-1} \frac{\sigma_e^2[k]}{A^2} \right)^{\frac{1}{K}} \\ &\geq -10 \log \left(\frac{1}{K} \sum_{k=0}^{K-1} \frac{\sigma_e^2[k]}{A^2} \right) = 10 \log \frac{A^2}{\frac{1}{K} \sum_{k=0}^{K-1} \sigma_e^2[k]} = \text{PSNR}_{\text{video}} \end{aligned}$$

due to the arithmetic mean being greater than the geometric mean.

Result: averaging of PSNR values over all images of a video sequence results in bigger values, the more unequally distributed the errors are over the images of the video sequence

Mutual Information for Discrete RVs

Mutual information (dt. *Transinformation*) between two discrete random variables X and Y specifies the information provided by X about Y

Definition given the joint probability mass functions $p_{X,Y}(x, y)$ and marginal probability mass functions $p_X(x)$ and $p_Y(y)$

$$I(X; Y) = \sum_x \sum_y p_{X,Y}(x, y) \log_2 \frac{p_{X,Y}(x, y)}{p_X(x) \cdot p_Y(y)}$$

Properties of mutual information

$$\begin{aligned} I(X; Y) &= I(Y; X) \geq 0 \\ I(X; Y) &\leq H(X) \quad \text{and} \quad I(Y; X) \leq H(Y) \\ I(X; Y) &= H(Y) - H(Y|X) = H(X) - H(X|Y) \end{aligned}$$

Rate-Distortion Function

Definition of rate-distortion function using mutual information

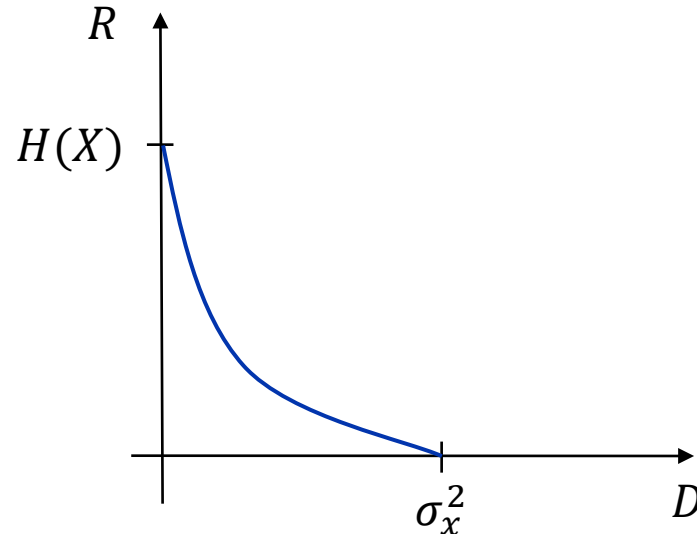
$$R(D) = \min_{d \leq D} I(X; \hat{X})$$

- For a given maximum average distortion D , the rate distortion function $R(D)$ is the lower bound for the transmission bit rate

Source coding theorem: for any $D \geq 0$ there exists a source code with average distortion $d \leq D$ and rate R arbitrarily close to $R(D)$

Properties of $R(D)$ function

- Convex
- Continuous and monotonically decreasing
- Inverse $D(R)$ exists and is called distortion-rate function



Continuous Random Variables

Problem in continuous case: entropy as defined previously is infinite

- Replace probability mass function by probability density function $p_X(x)$
- Define *differential entropy*

$$h(X) = E\{-\log_2 p_X(X)\} = - \int p_X(x) \log_2 p_X(x) dx$$

- Relative measure of uncertainty, can be negative

Gaussian RV X with zero mean and variance σ^2

$$p_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-x^2/2\sigma^2}$$

$$\begin{aligned} h(X) &= - \int p_X(x) \log_2 p_X(x) dx \\ &= \frac{1}{2} \log_2 2\pi e \sigma^2 \end{aligned}$$

Shannon lower bound: for an IID process, the MSE rate-distortion function is lower bounded by

$$R_L(D) = h(X) - \frac{1}{2} \log_2 2\pi e D$$

Rate Distortion for IID Gaussian Source

IID Gaussian source X with variance σ^2

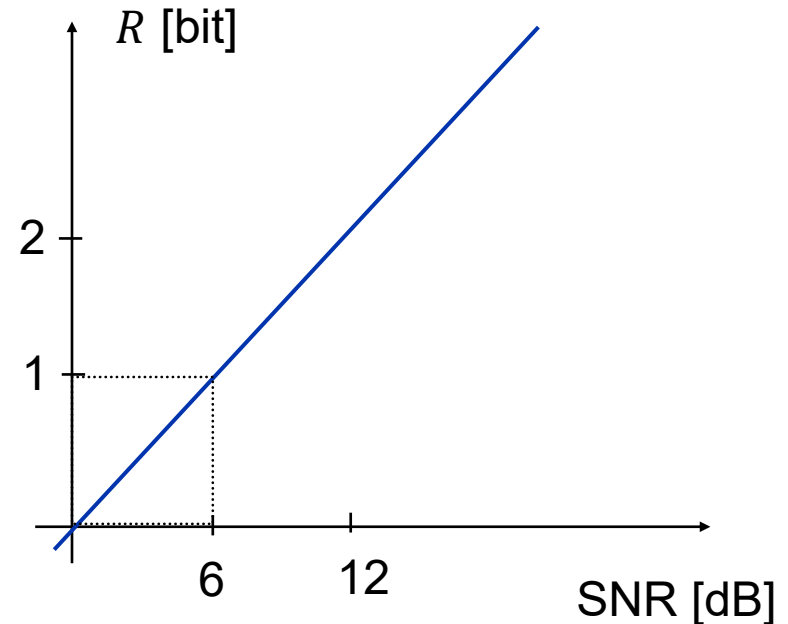
R(D) lower bound for MSE distortion

$$R(D) = \frac{1}{2} \log_2 2\pi e \sigma^2 - \frac{1}{2} \log_2 2\pi e D$$
$$= \begin{cases} \frac{1}{2} \log_2 \frac{\sigma^2}{D} & \text{for } \sigma^2 > D \\ 0 & \text{else} \end{cases}$$

$$D(R) = \sigma^2 2^{-2R} \quad R \geq 0$$

Theoretical bound on **signal-to-noise ratio**

$$\text{SNR} = 10 \log_{10} \frac{\sigma^2}{D(R)} [\text{dB}]$$
$$= 20R \log_{10} 2 \cong 6.02 \cdot R [\text{dB}]$$



$R(D)$ for non-Gaussian sources with same σ^2 is always below Gaussian

Rule of thumb: 1 bit corresponds to approximately 6 dB in SNR

Rate Distortion for Non-Gaussian IID Source

Rate-distortion function for any IID source

- Typically not to be expressed in closed form, computed numerically
- Bounded by

$$R_L(D) \leq R(D) \leq \frac{1}{2} \log_2 \frac{\sigma^2}{D}$$

Shannon lower bound

Gaussian rate-distortion function

Distortion-rate function equivalently bounded by

$$D_L(R) = \frac{1}{2\pi e} 2^{2h(X)} 2^{-2R} \leq D(R) \leq \sigma^2 2^{-2R}$$

“Entropy power”

	Entropy power
Uniform	$\frac{6}{\pi e} \sigma^2 \cong 0.703 \sigma^2$
Laplacian	$\frac{e}{\pi} \sigma^2 \cong 0.865 \sigma^2$
Gaussian	σ^2

Rate Distortion for Correlated Gaussian Source

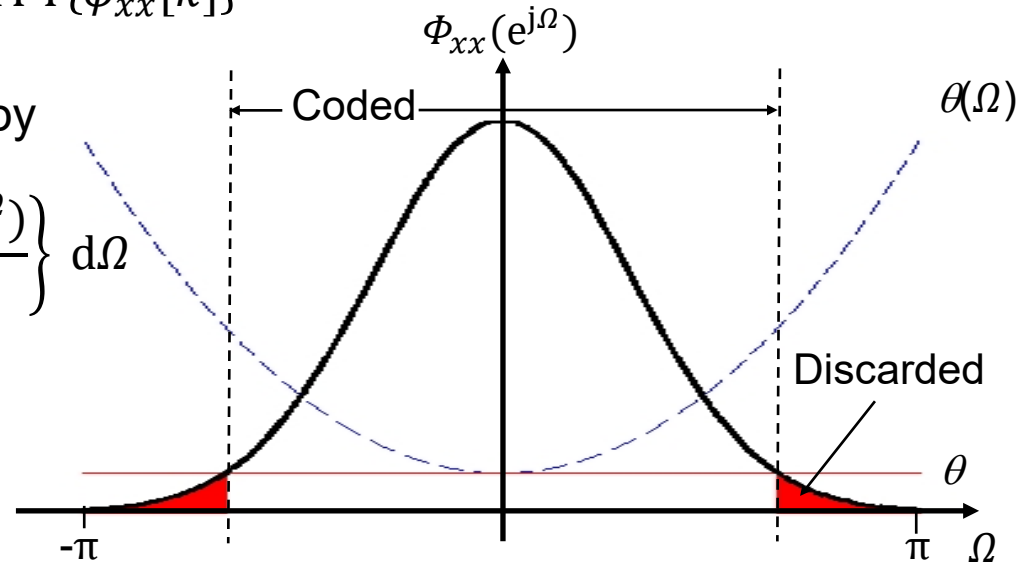
Assumption: Discrete Gaussian source x with power spectrum

$$\Phi_{xx}(e^{j\Omega}) = \text{DTFT}\{\varphi_{xx}[k]\}$$

Rate distortion function is given by

$$R(\theta) = \frac{1}{2\pi} \int_{-\pi}^{\pi} \max\left\{0, \frac{1}{2} \log_2 \frac{\Phi_{xx}(e^{j\Omega})}{\theta}\right\} d\Omega$$

$$D(\theta) = \frac{1}{2\pi} \int_{-\pi}^{\pi} \min\{\theta, \Phi_{xx}(e^{j\Omega})\} d\Omega$$

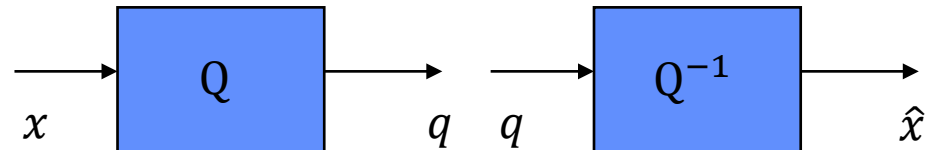
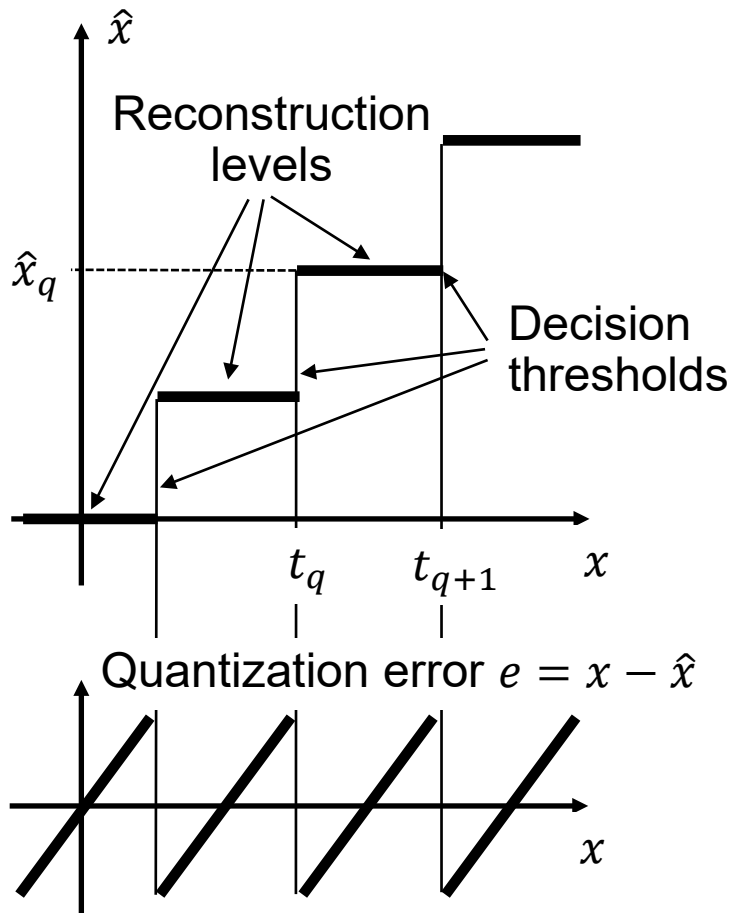


Consequences

- Frequency range over which power spectrum is smaller than θ need not be coded
- Remaining frequency range should be coded with rate R such that error signal has power equal to θ

5.2 Scalar Quantization

Input-output characteristics of scalar quantizer



Principle

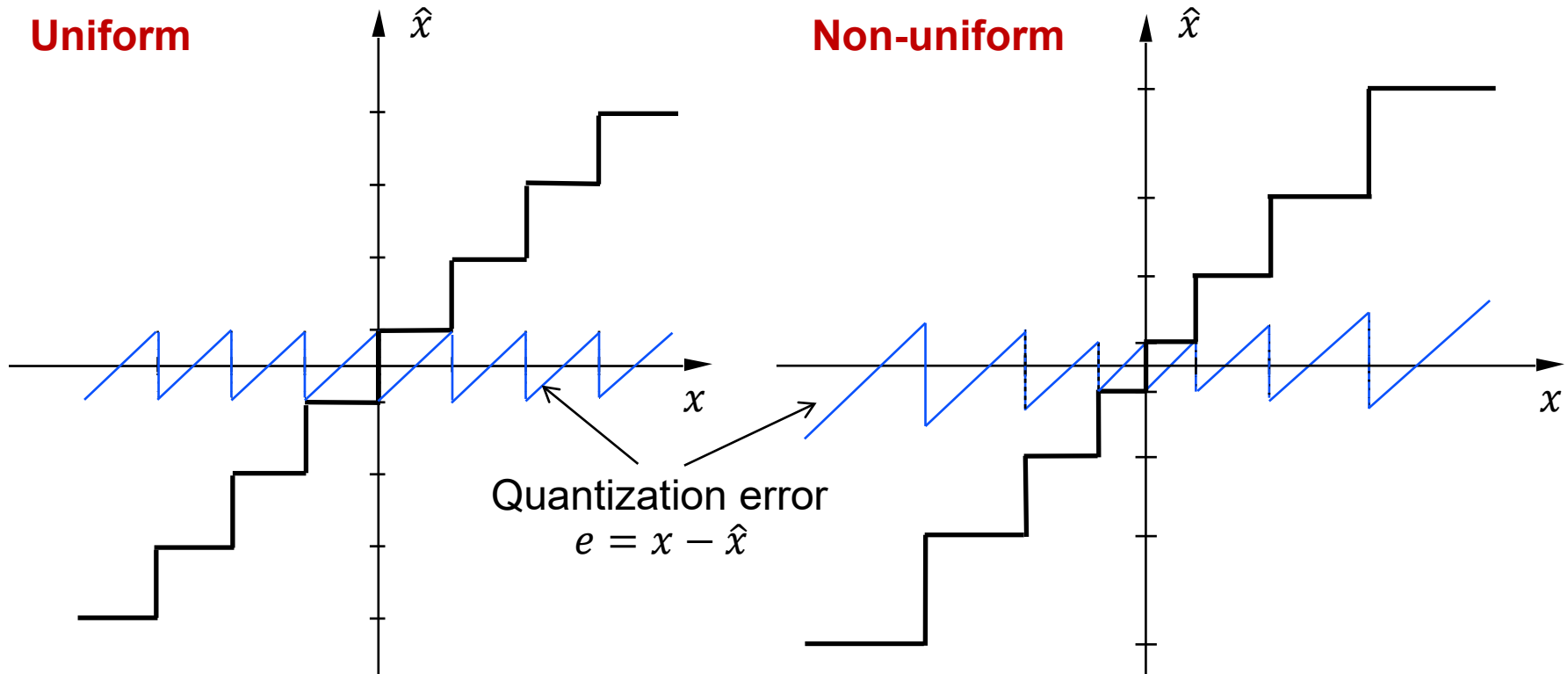
Reconstruction levels are attributed to continuous range of input values

Construction of quantizer according to

- error criterion (maximum error, total error,...) or
- entropy criterion

For irrelevancy reduction in coding systems quantization is performed on color values, prediction signals, transform coefficients, ...

Uniform versus Non-uniform Quantization



⇒ Non-uniform quantization to adapt error to psychophysical properties by taking advantage of Weber's law

Midrise quantizer: symmetric with even number of reconstruction levels (no zero)

Midtread quantizer: symmetric with odd number of reconstruction levels

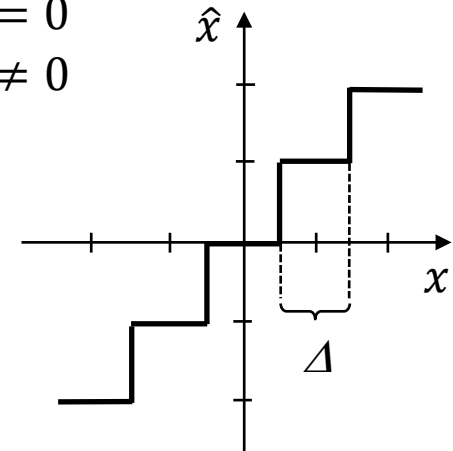
Uniform Midtread Quantization

Quantize: $q = Q(x) = \text{sign}(x) \left\lfloor \frac{|x|}{\Delta} + \frac{1}{2} \right\rfloor$

Dequantize: $\hat{x} = Q^{-1}(q) = \begin{cases} 0 & q = 0 \\ \text{sign}(q)(|q| + \delta)\Delta & q \neq 0 \end{cases}$

Δ = quantization step size

δ = offset to reflect shape of $p_x(x)$, zero for uniform distribution of input X



Mean square error of uniform quantizer:

Approximation for small Δ : $p_E(e) \cong \frac{1}{\Delta} \quad \text{for} \quad -\frac{\Delta}{2} \leq e < \frac{\Delta}{2}$

\Rightarrow Variance of quantization error: $\sigma_e^2 = \int_{-\Delta/2}^{\Delta/2} p_E(e) \cdot e^2 de = \frac{\Delta^2}{12}$

Uniform Quantization with Deadzone

Quantize: $q = Q(x) = \begin{cases} 0 & |x| < \beta \\ \text{sign}(x) \left\lfloor \frac{|x| - \beta}{\Delta} + 1 \right\rfloor & \text{else} \end{cases}$

Dequantize: $\hat{x} = Q^{-1}(q) = \begin{cases} 0 & q = 0 \\ \text{sign}(q) \left((|q| - \frac{1}{2} + \delta) \Delta + \beta \right) & q \neq 0 \end{cases}$

Δ = quantization step size

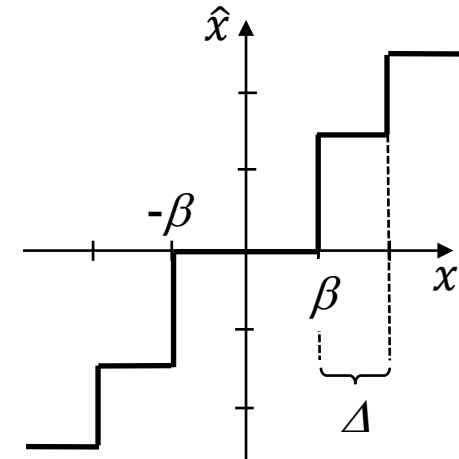
β = threshold for quantization into zero bin

δ = offset to reflect shape of $p_x(x)$, zero for uniform distribution of input X

Special cases

$\beta = \Delta/2$ uniform midtread quantizer as before

$\beta = \Delta$ width of zero bin is 2Δ



5.3 Lloyd-Max Quantization

Problem: given a signal with known PDF $p_X(x)$, find a quantizer with M reconstruction levels such that MSE is minimized

$$d = E\{(X - \hat{X})^2\} = \sum_{k=0}^{M-1} \int_{t_k}^{t_{k+1}} (x - \hat{x})^2 p_X(x) dx \rightarrow \min$$

Approach: Lloyd-Max scalar quantizer with two necessary conditions

Setting partial derivative of d with respect to t_q equal to zero yields

- Place $M - 1$ decision thresholds half way between reconstruction levels

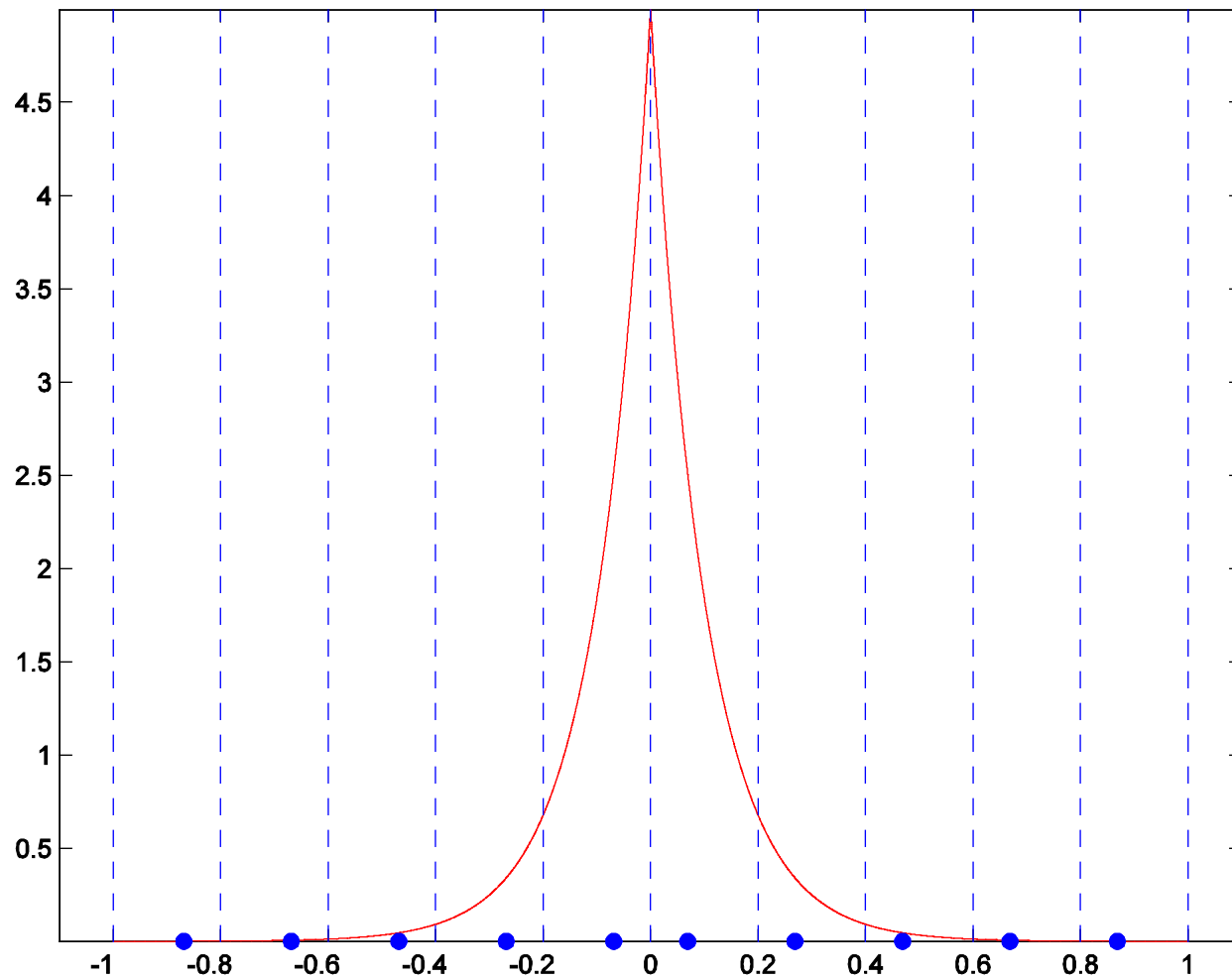
$$t_q = \frac{\hat{x}_{q-1} + \hat{x}_q}{2}$$

Setting partial derivative of d with respect to \hat{x}_q equal to zero yields

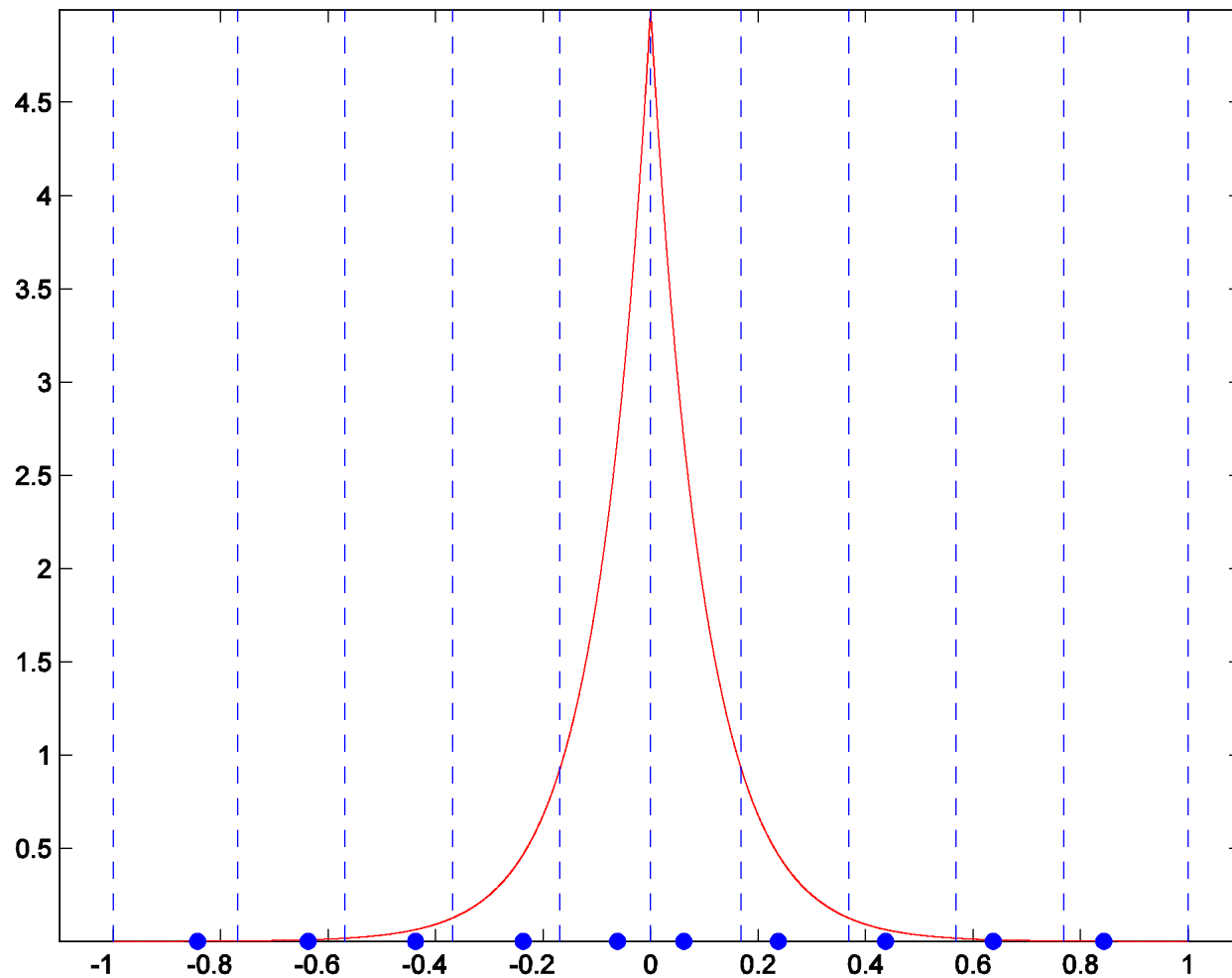
- Place M reconstruction levels in the center of mass of the PDF between two successive decision thresholds

$$\hat{x}_q = \frac{\int_{t_q}^{t_{q+1}} x \cdot p_X(x) dx}{\int_{t_q}^{t_{q+1}} p_X(x) dx}$$

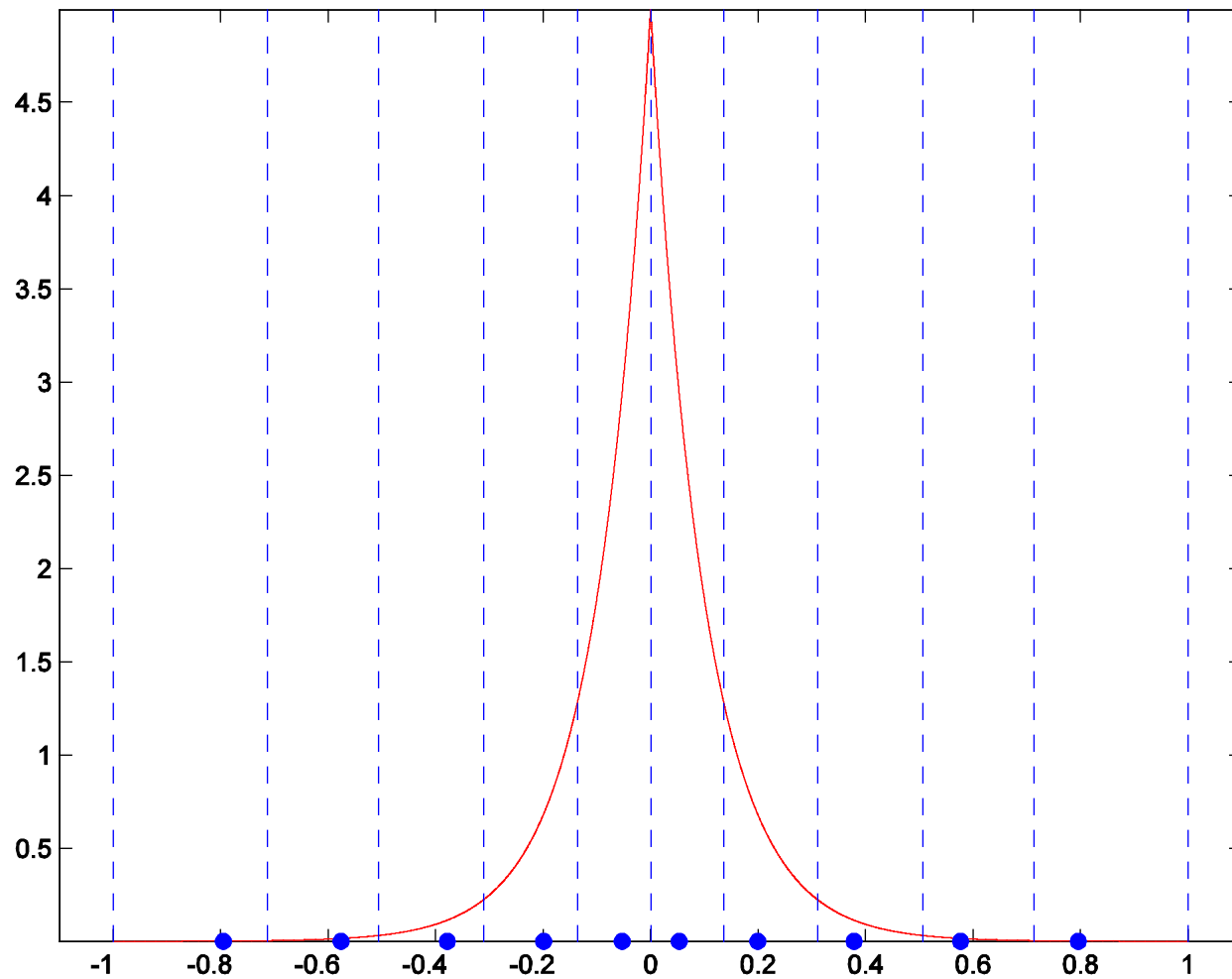
Lloyd-Max Quantization for Laplace Distribution, $n = 1$



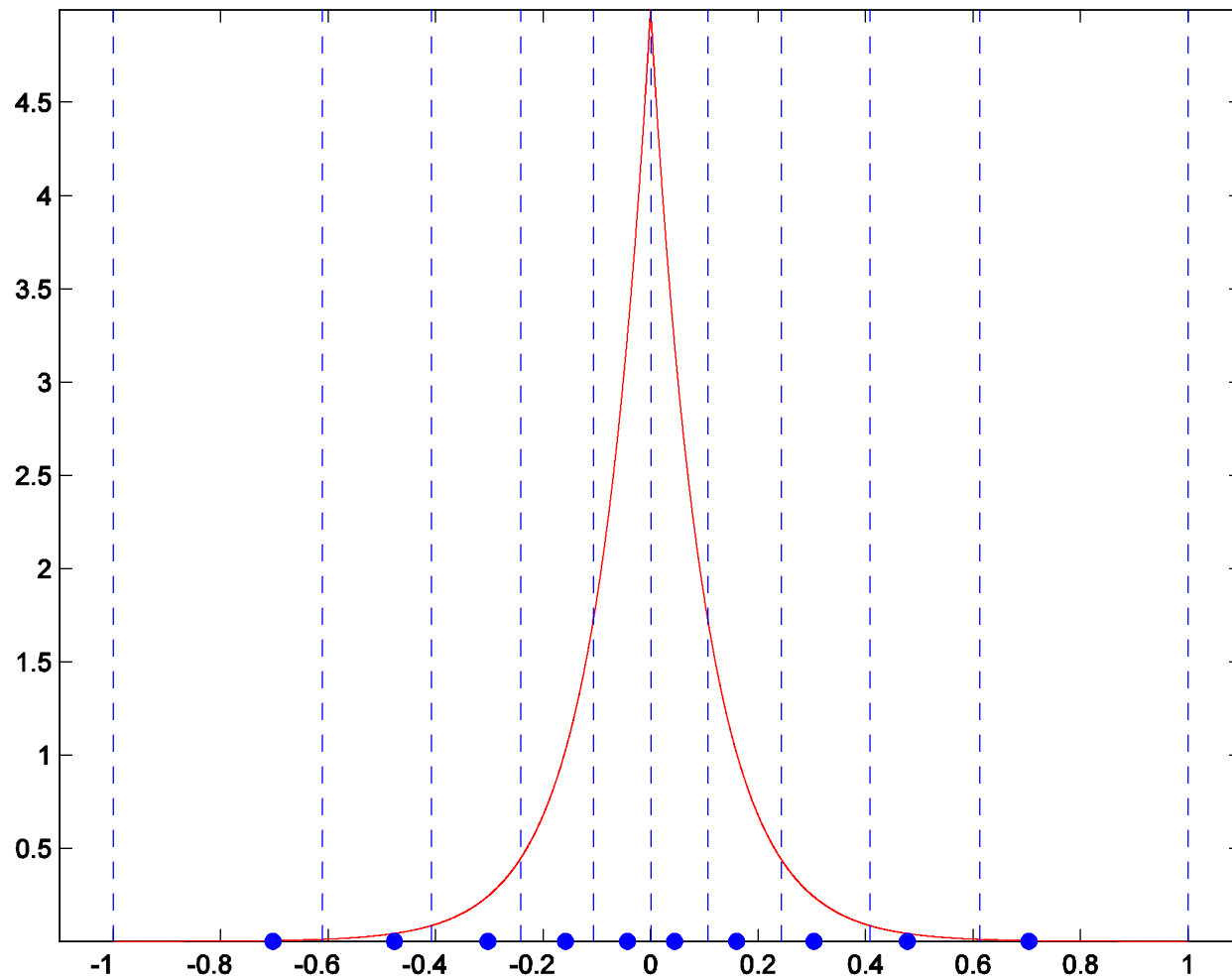
Lloyd-Max Quantization for Laplace Distribution, $n = 2$



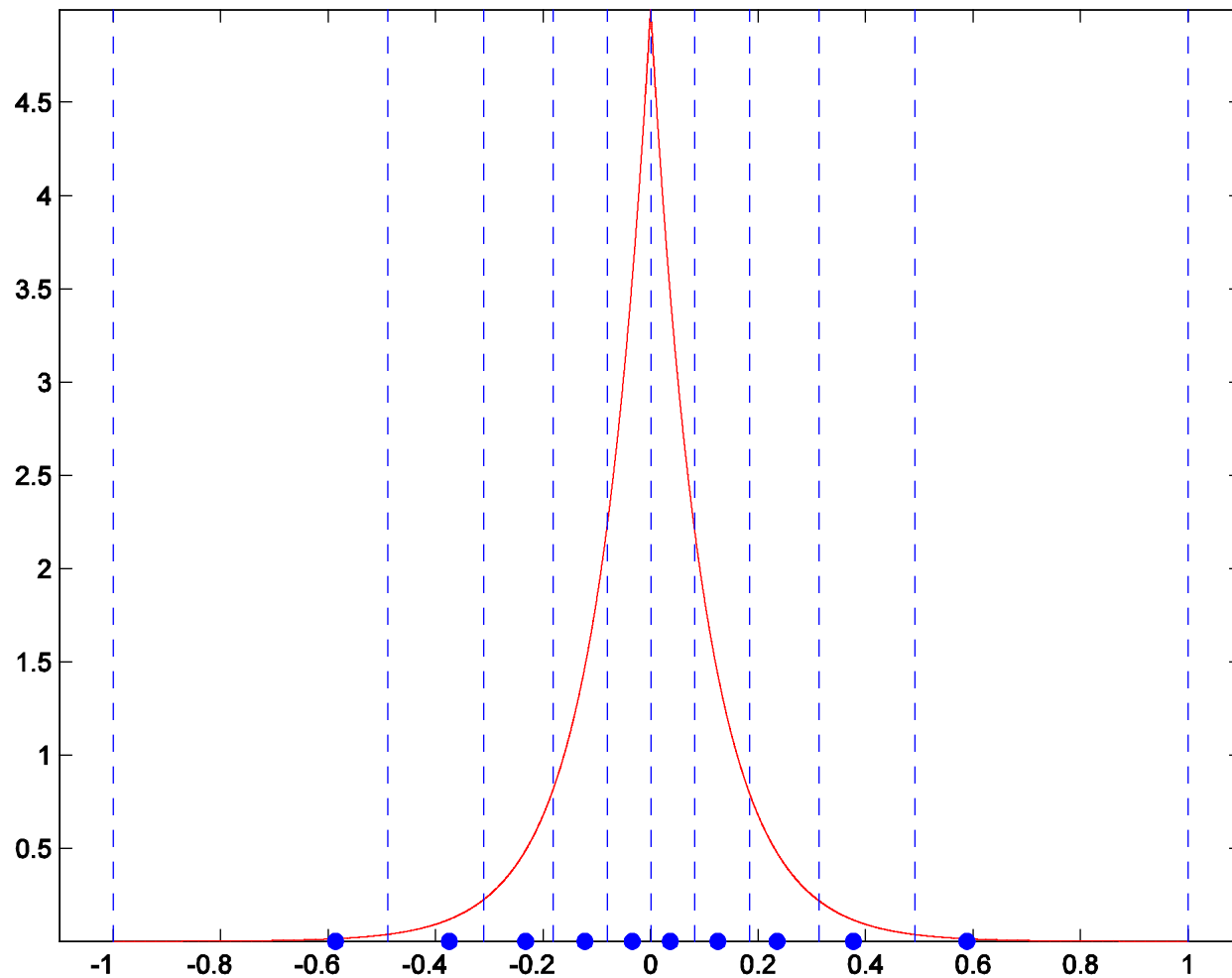
Lloyd-Max Quantization for Laplace Distribution, $n = 4$



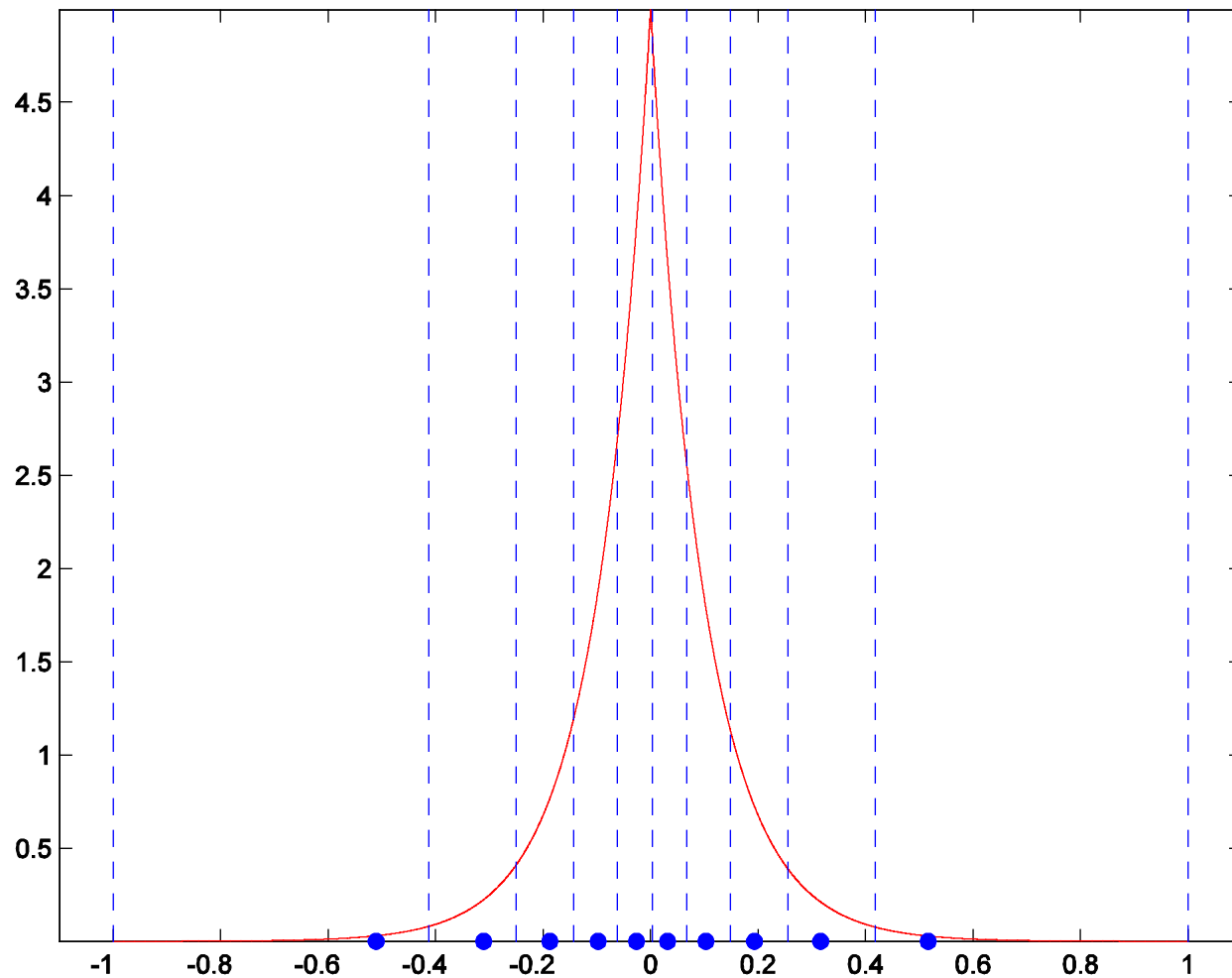
Lloyd-Max Quantization for Laplace Distribution, $n = 8$



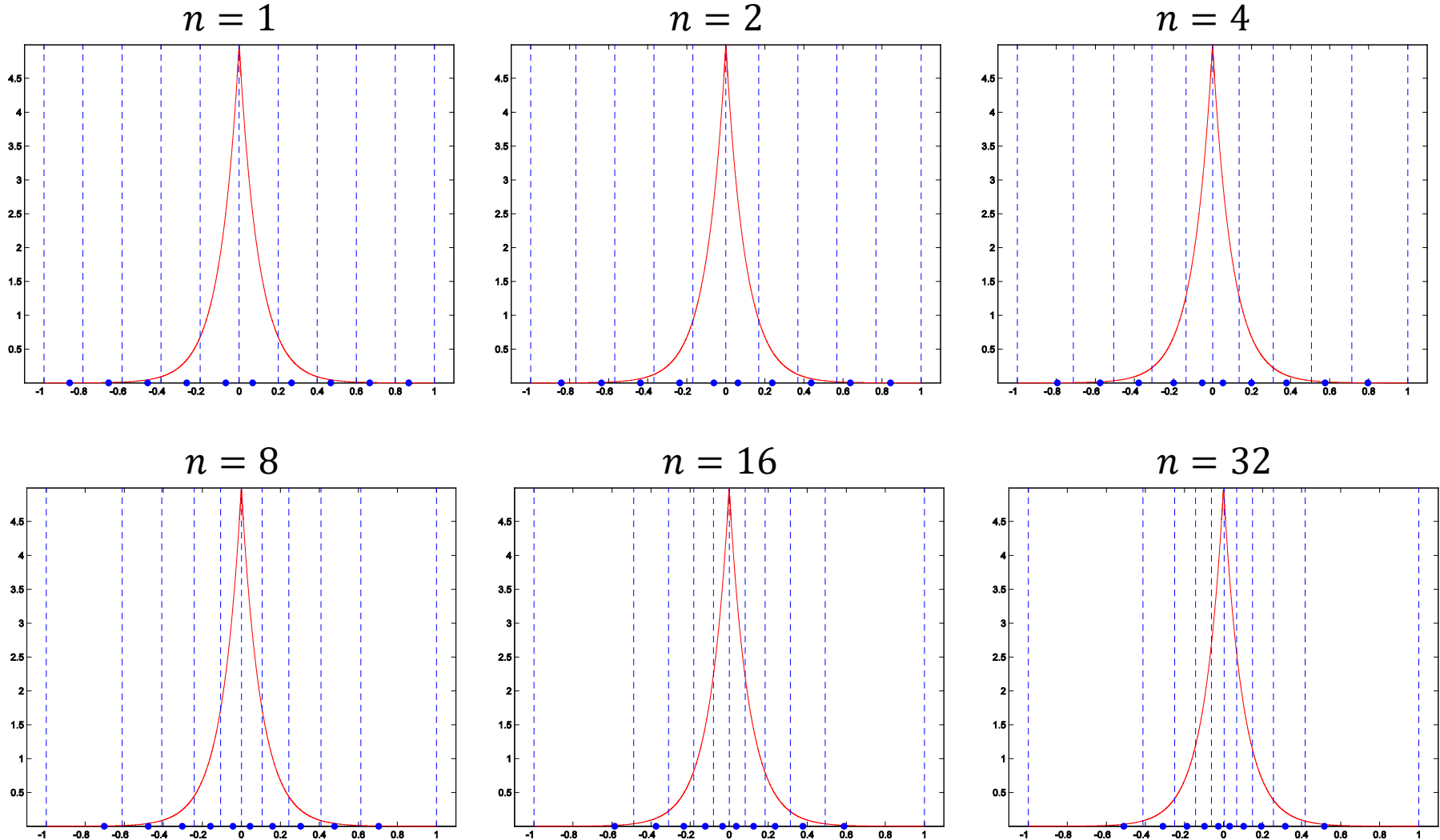
Lloyd-Max Quantization for Laplace Distribution, $n = 16$



Lloyd-Max Quantization for Laplace Distribution, $n = 32$



Lloyd-Max Quantization for Laplace Distribution



Lloyd-Max Algorithm Based on Training Set

Lloyd-Max algorithm for quantizer design using training data

Choose initial set of representative levels $\hat{x}_q, q = 0, 1, \dots, M - 1$

Repeat

Assign each sample x_i in training set T to closest representative \hat{x}_q minimizing the Euclidian distance

$$J_{x_i}(q) = (x_i - \hat{x}_q)^2 \quad q = 0, 1, \dots, M - 1$$

yielding sets

$$B_q = \{x_i \in T : Q(x_i) = q\} \quad q = 0, 1, \dots, M - 1$$

Calculate M new representative levels as mean of each set B_q

$$\hat{x}_q = \frac{1}{\|B_q\|} \sum_{x_i \in B_q} x_i, \quad q = 0, 1, \dots, M - 1$$

Until no further reduction in total distortion $d = \sum_{x_i} (x_i - \hat{x}_{Q(x_i)})^2$

Properties and Performance of Lloyd-Max

Zero-mean: quantization error has zero mean independent of whether input signal has zero mean or not

$$E\{(X - \hat{X})\} = 0$$

Decorrelation: quantization error and quantizer output are uncorrelated

$$E\{(X - \hat{X})\hat{X}\} = 0$$

- but: quantization error typically correlated with quantizer input

Variance reduction: variance of quantizer output is reduced by amount of MSE

$$\sigma_{\hat{X}}^2 = \sigma_X^2 - E\{(X - \hat{X})^2\}$$

Equal contribution: all intervals contribute equally towards the overall MSE

$$E\{(X - \hat{X})^2 | X \in I_j\} p_j = E\{(X - \hat{X})^2 | X \in I_k\} p_k \quad \forall j, k$$

with interval $I_q = [t_q, t_q + 1)$ and p_q equals probability of interval q

High Rate Approximation of Lloyd-Max

Approximation: for large rate R the distortion-rate function of Lloyd-Max quantization behaves like

$$d(R) \cong \varepsilon^2 \sigma^2 2^{-2R}$$

Parameter ε^2 depends on particular PDF, for zero-mean symmetric PDF it follows:

$$\varepsilon^2 \sigma^2 = \frac{2}{3} \left[\int_0^\infty \sqrt[3]{p_X(x)} dx \right]^3$$

Example values of ε^2 compared to Shannon lower bound

	$D(R)$	Lloyd-Max
Uniform	$\frac{6}{\pi e} \cong 0.703$	1
Laplacian	$\frac{\pi}{e} \cong 0.865$	$\frac{9}{2} = 4.5$
Gaussian	1	$\frac{\sqrt{3}\pi}{2} \cong 2.721$

 Demo 5 „Lloyd Max“

5.4 Entropy Coded Scalar Quantization

Coding of quantizer index

- Lloyd-Max quantizer optimum for coding at fixed rate
- How to incorporate variable length encoding of index?

Problem formulation: for a signal x with given distribution p_X , we seek to minimize the MSE distortion

$$E\{(X - \hat{X})^2\} = \sum_{q=0}^{M-1} \int_{t_q}^{t_{q+1}} (x - \hat{x}_q)^2 p_X(x) dx$$

subject to the constraint that

$$H(\hat{X}) = - \sum_{q=0}^{M-1} p_q \log_2 p_q \leq R \quad \text{with} \quad p_q = \int_{t_q}^{t_{q+1}} p_X(x) dx$$

Solution: minimize Lagrangian cost function

$$J = E\{(X - \hat{X})^2\} + \lambda H(\hat{X})$$

Iterative Entropy Coded Scalar Quantizer Design

Lloyd-Max algorithm for entropy coded scalar quantizer

Choose initial set of representative levels $\hat{x}_q, q = 0, 1, \dots, M - 1$
and corresponding probabilities p_q

Repeat

Calculate $M - 1$ decision thresholds

$$t_q = \frac{\hat{x}_{q-1} + \hat{x}_q}{2} + \lambda \frac{\log_2 p_{q-1} - \log_2 p_q}{2(\hat{x}_q - \hat{x}_{q-1})} \quad q = 0, 1, \dots, M - 1$$

Calculate M new representative levels and probabilities p_q

$$\hat{x}_q = \frac{\int_{t_q}^{t_{q+1}} x \cdot p_X(x) dx}{\int_{t_q}^{t_{q+1}} p_X(x) dx}, \quad p_q = \int_{t_q}^{t_{q+1}} p_X(x) dx \quad q = 0, 1, \dots, M - 1$$

Until no further reduction in Lagrangian cost

Extension by outer loop to find suitable parameter $\lambda > 0$ minimizing J

Entropy Constraint Design Based on Training Set

Lloyd-Max algorithm for entropy coded quantizer using training data

Choose initial set of representative levels $\hat{x}_q, q = 0, 1, \dots, M - 1$ and corresponding probabilities p_q

Repeat

Assign each sample x_i in training set T to representative \hat{x}_q minimizing Lagrangian cost

$$J_{x_i}(q) = (x_i - \hat{x}_q)^2 - \lambda \log_2 p_q \quad q = 0, 1, \dots, M - 1$$

yielding sets

$$B_q = \{x_i \in T : Q(x_i) = q\} \quad q = 0, 1, \dots, M - 1$$

Calculate M new representative levels and probabilities p_q

$$\hat{x}_q = \frac{1}{\|B_q\|} \sum_{x_i \in B_q} x_i, \quad p_q = \frac{\|B_q\|}{\sum_{q=0}^{M-1} \|B_q\|} \quad q = 0, 1, \dots, M - 1$$

Until no further reduction in total cost $J = \sum_{x_i} [(x_i - \hat{x}_{Q(x_i)})^2 - \lambda \log_2 p_{Q(x_i)}]$

High Rate Performance of EC Scalar Quantization

High rate and MSE distortion: uniform quantizer with very large number of levels is *optimum* scalar quantizer in entropy coded case [Gish, Pierce, 1968]

Distortion is approximately constant for small quantizer interval Δ

$$d \cong \frac{\Delta^2}{12}$$

Entropy is approximately given by

$$H(\hat{X}) = - \sum_{q=-\infty}^{\infty} p_q \log_2 p_q \cong h(X) - \log_2 \Delta$$

If efficient coding is used, it follows that $R \cong H(\hat{X}) \rightarrow \Delta \cong 2^{h(X)-R}$

Distortion-rate function for entropy coded (uniform) scalar quantization

$$d(R) = \frac{1}{12} 2^{2h(X)} 2^{-2R}$$

is 1.53 dB from Shannon lower bound $D(R) \geq \frac{1}{2\pi e} 2^{2h(X)} 2^{-2R}$

Comparison of High Rate Performance

Observation: high-rate distortion function for IID data in case of

- Lloyd-Max quantization as well as
- entropy coded (uniform) quantization

is of general form

$$d(R) \cong \varepsilon^2 \sigma^2 2^{-2R}$$

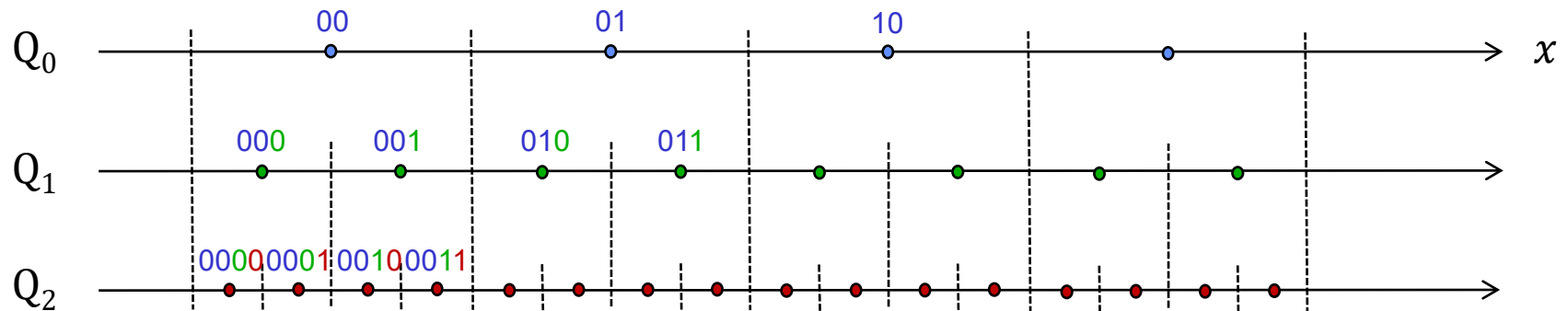
Comparison of scaling factor ε^2

	Shannon $D(R)$	Lloyd – Max	Entropy coded
Uniform	$\frac{6}{\pi e} \cong 0.703$	1	1
Laplacian	$\frac{e}{\pi} \cong 0.865$	$\frac{9}{2} = 4.5$	$\frac{e^2}{6} \cong 1.232$
Gaussian	1	$\frac{\sqrt{3}\pi}{2} \cong 2.721$	$\frac{\pi e}{6} \cong 1.423$

5.5 Embedded Quantization

Scalability: successively refine reconstructed data as bit-stream is decoded

- Decoded subset gives lower quality signal approximation
- Facilitated by nested (“embedded”) quantization



Coding: form quantizer index by adding $\log_2 M_k$ bits for M_k intervals at Q_k

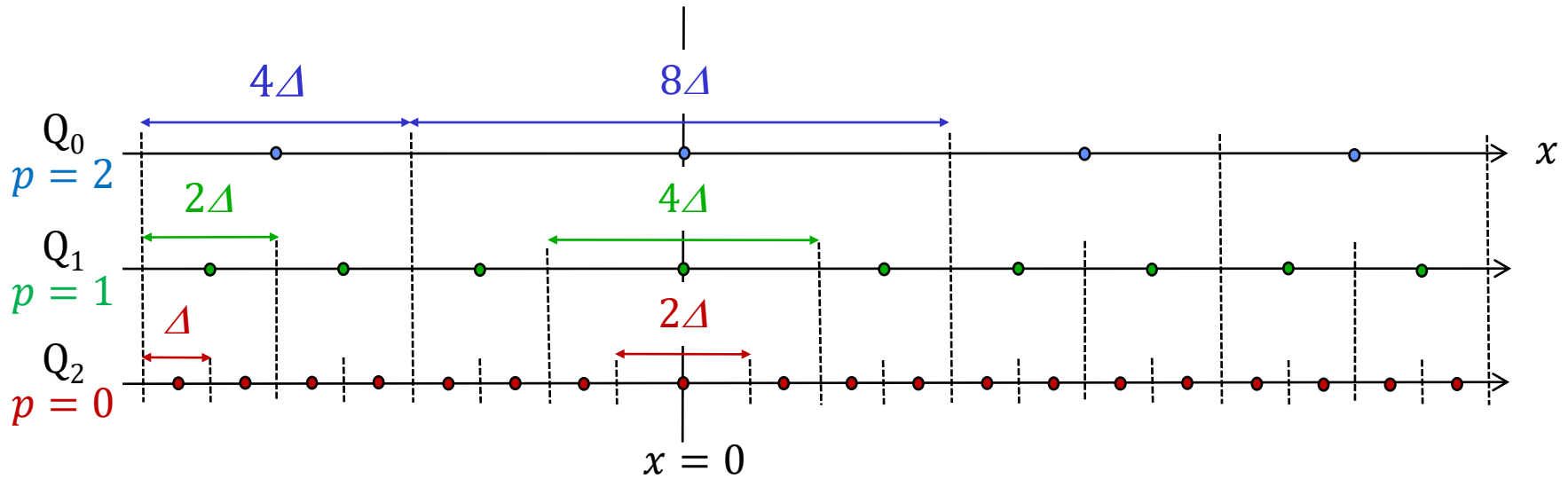
- Lower quantizations can be formed by dropping components from indices of higher rate approximations

Restriction: in general only one quantizer can be optimum with respect to Lloyd-Max condition (exception: uniform quantizer)

Embedded Quantization with Deadzone

Construction of a family of embedded uniform scalar quantizers with deadzone

- Typical case: width of zero-bin is 2Δ



Quantizer index for $p = 0$ given by $q = Q(x) = \text{sign}(x) \left\lfloor \frac{|x|}{\Delta} \right\rfloor$

Reconstruction from index: $\hat{x} = Q^{-1}(q) = \begin{cases} 0 & q = 0 \\ \text{sign}(q)(|q| + \delta)\Delta & q \neq 0 \end{cases}$

Coding for Embedded Quantization with Deadzone

Embedded coding of nested deadzone quantizer with step sizes $2^p \Delta$:

- Assume that quantizer index q can be represented with K bits
- Index q can be written in sign plus magnitude form as

$$q = Q_{K-1}(x) = s, q_0, q_1, \dots, q_{K-1}$$

- Dropping last p bits from q

$$q_p = Q_{K-1-p}(x) = s, q_0, q_1, \dots, q_{K-1-p}$$

gives the uniform deadzone quantizer with step size $2^p \Delta$

Same result as if quantization was performed using step size of $2^p \Delta$ rather than Δ in the first place

- If p LSBs of q are unavailable, simply reconstruct at lower level of quality

- Reconstruction rule:
$$\hat{x} = Q^{-1}(q_p) = \begin{cases} 0 & q_p = 0 \\ \text{sign}(q_p)(|q_p| + \delta)2^p \Delta & q_p \neq 0 \end{cases}$$

5.6 Adaptive Quantization

Perception of quantization errors

- homogeneous objects of medium brightness
- structured areas or very bright areas



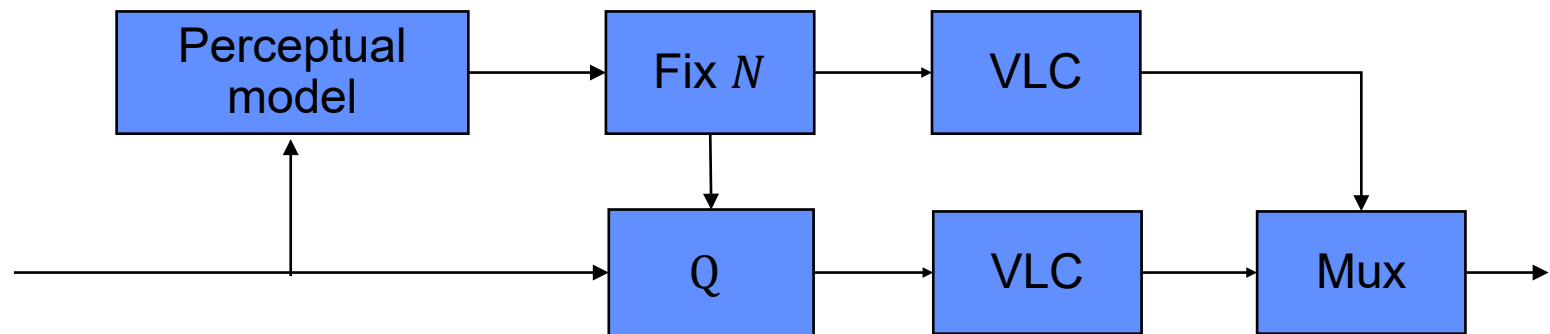
Quantization scale



very critical, fine quantization

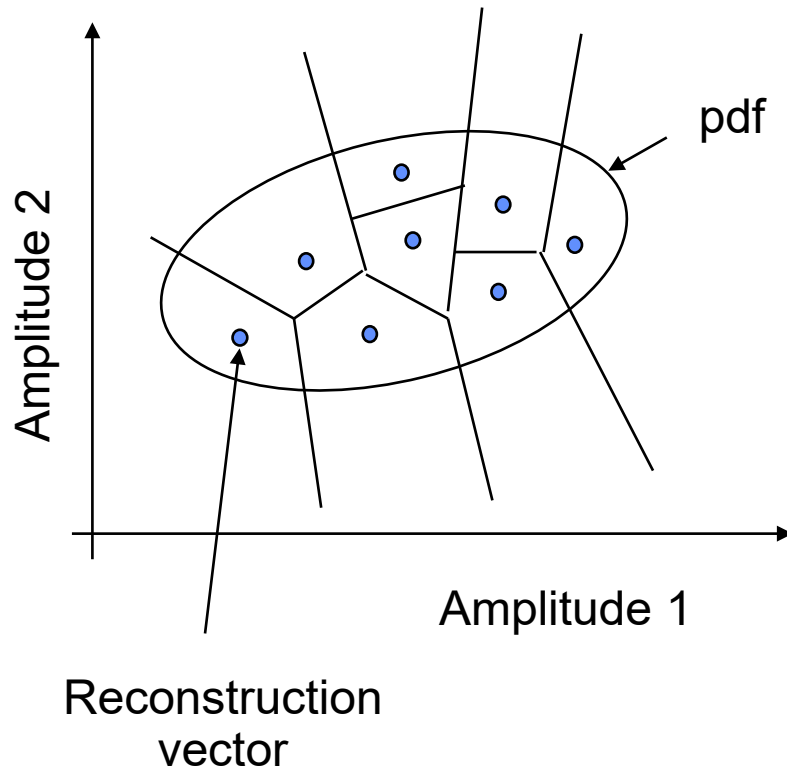


less critical, coarse quantization

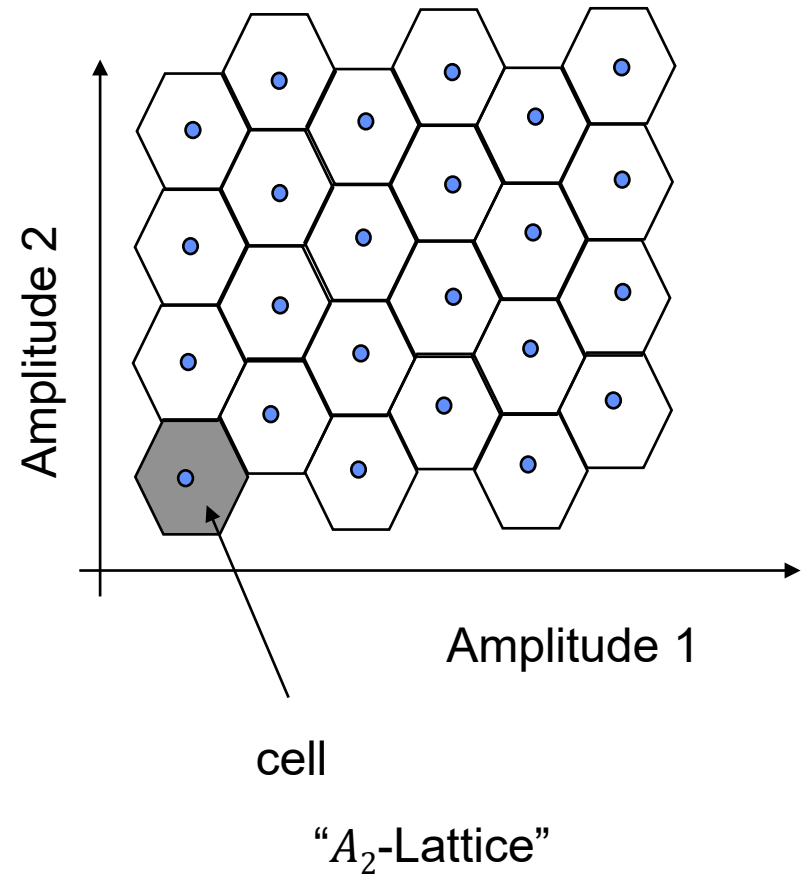


5.7 Vector Quantization

Non-uniform codebook



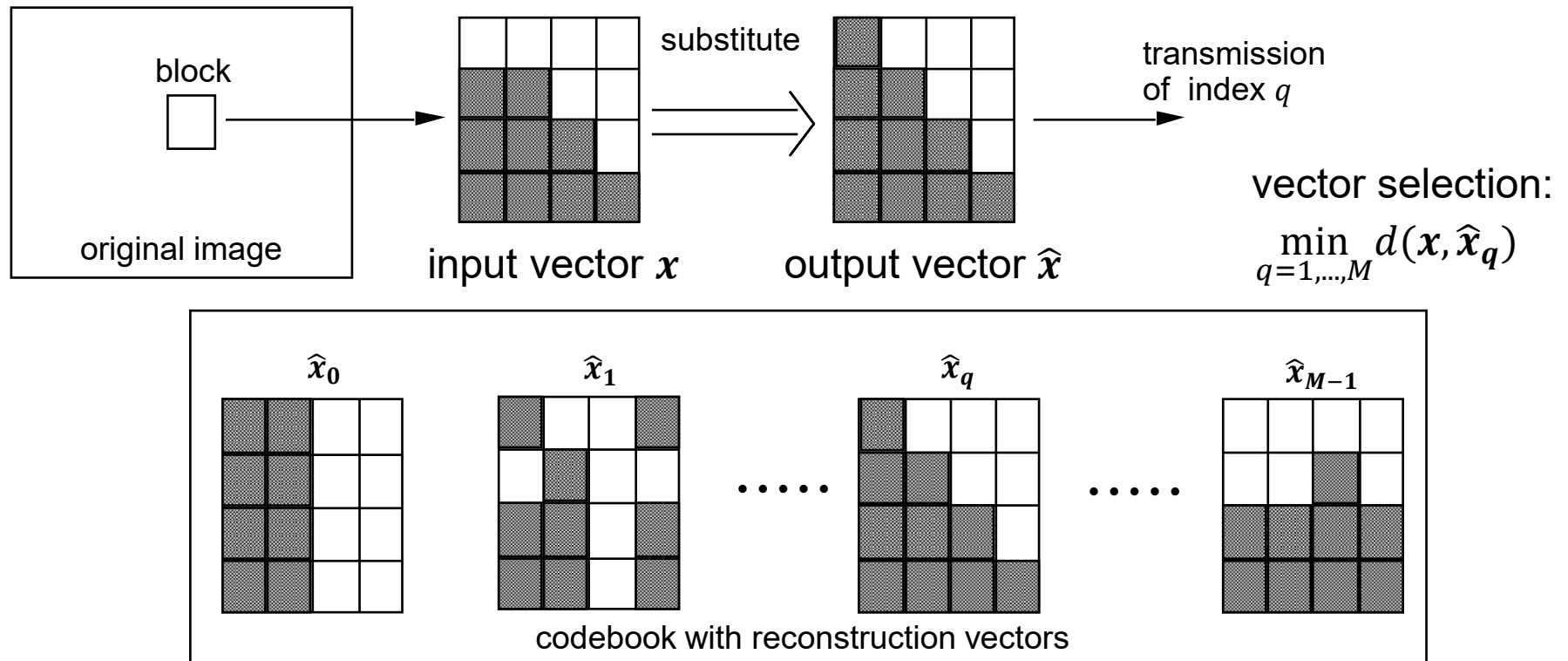
Uniform codebook (Lattice VQ)



Vector Quantization for Images

Idea: Image block is regarded as multidimensional vector x

- Codebook contains a reduced ensemble of all possible image blocks
- Image block is replaced by similar vector out of codebook
- Only codebook index is transmitted
- Optimal codebook entry is selected based on a distortion measure



LBG Algorithm

Generalization of Lloyd-Max algorithm for vector quantization

First published by Linde, Buzo, and R. Gray in 1980 \Rightarrow “LBG Algorithm”

Assumption: fixed code word length for index q

Idea taken from Lloyd-Max algorithm

- Successive optimization of code book using suitable training set

Problem

- Unstructured code book requires full search
- Computationally expensive

LBG Algorithm (cont.)

LBG algorithm for vector quantizer design

Choose training set T and initial set of reconstruction vectors $\hat{x}_q, q = 0, 1, \dots, M - 1$

Repeat

Assign each sample x_i in training set T to closest representative \hat{x}_q minimizing the Euclidian distance

$$J_{x_i}(q) = \|x_i - \hat{x}_q\|^2 \quad q = 0, 1, \dots, M - 1$$

yielding sets

$$B_q = \{x_i \in T : Q(x_i) = q\} \quad q = 0, 1, \dots, M - 1$$

Calculate M new reconstruction vectors as centroid of each set B_q

$$\hat{x}_q = \frac{1}{\|B_q\|} \sum_{x_i \in B_q} x_i, \quad q = 0, 1, \dots, M - 1$$

Until no further reduction in total distortion $d = \sum x_i \|x_i - \hat{x}_{Q(x_i)}\|^2$

Entropy Coded Vector Quantization

Extended LBG algorithm for entropy coded vector quantizer design

Choose initial set of reconstruction vectors $\hat{x}_q, q = 0, 1, \dots, M - 1$ and corresponding probabilities p_q

Repeat

Assign each sample x_i in training set T to representative \hat{x}_q minimizing Lagrangian cost

$$J_{x_i}(q) = \|x_i - \hat{x}_q\|^2 - \lambda \log_2 p_q \quad q = 0, 1, \dots, M - 1$$

yielding sets

$$B_q = \{x_i \in T : Q(x_i) = q\} \quad q = 0, 1, \dots, M - 1$$

Calculate M new reconstruction vectors and probabilities p_q

$$\hat{x}_q = \frac{1}{\|B_q\|} \sum_{x_i \in B_q} x_i, \quad p_q = \frac{\|B_q\|}{\sum_{q=0}^{M-1} \|B_q\|} \quad q = 0, 1, \dots, M - 1$$

Until no further reduction in total Lagrangian cost $J = E\{\|X - \hat{X}\|^2\} + \lambda H(\hat{X})$

Quantization - Summary

- Rate distortion theory: minimum transmission bit rate for given distortion
- $R(D)$ for memoryless Gaussian source and MSE: 6 dB/bit
- Uniform quantization with small quantization step size
- Lloyd-Max quantization for optimum quantizer design
- Vector quantization allows joint quantization of several signal samples
- Design of optimum vector quantizer with LBG algorithm