



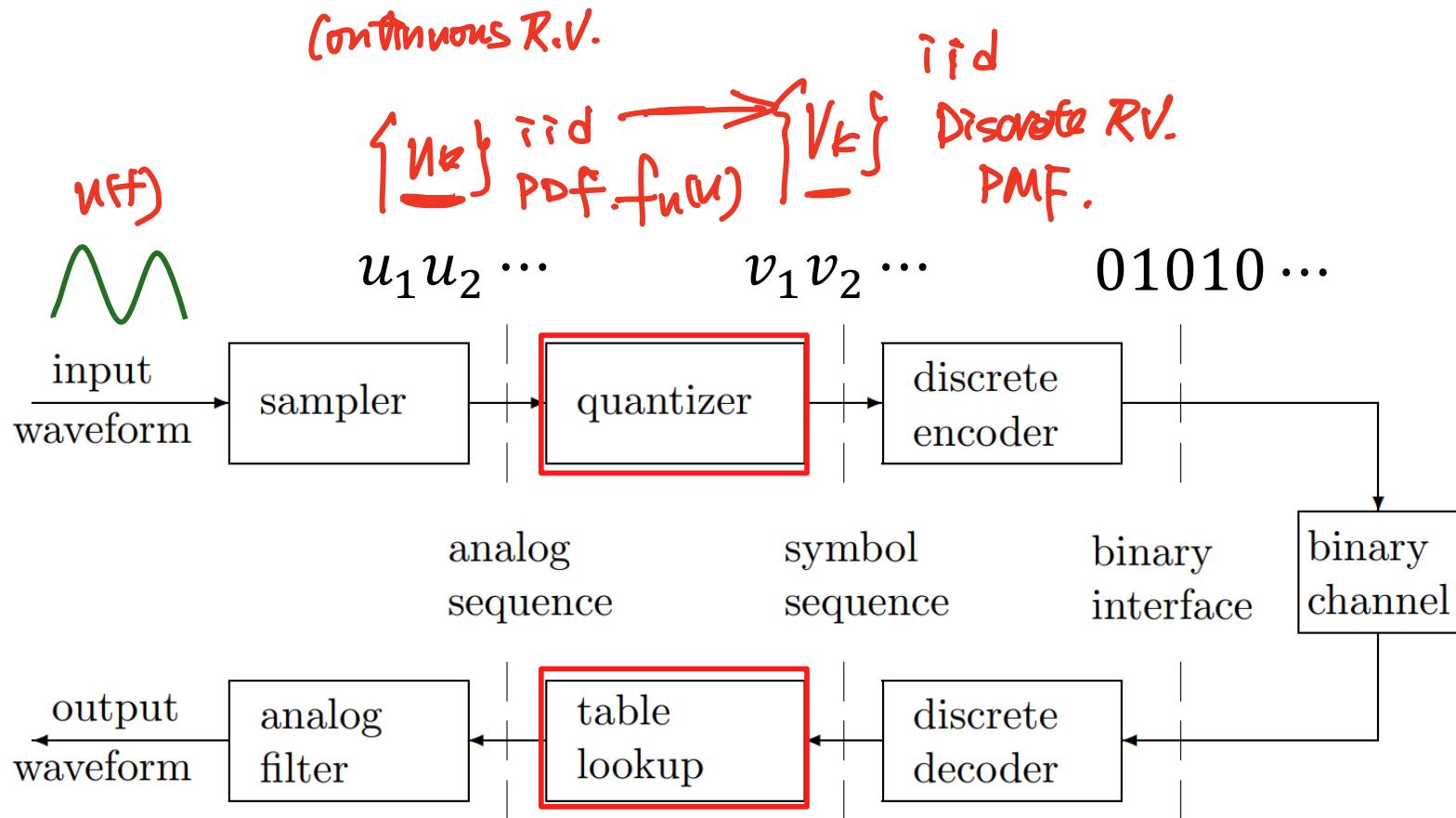
上海科技大学
ShanghaiTech University

EE140 Introduction to Communication Systems

Lecture 12

Instructor: Prof. Lixiang Lian
ShanghaiTech University, Fall 2025

Quantization



Quantization

Objective: Map the incoming sequence $\underline{U_1, U_2, \dots}$ of analog rvs into a sequence of discrete rvs $\underline{V_1, V_2, \dots}$, where $\underline{V_m}$ should represent $\underline{U_m}$ with as little distortion as possible.

- Scalar Quantization: Each analog RV in the sequence is quantized independently of the other RVs.
- Vector Quantization: The analog sequence is first segmented into blocks of n RVs each; then each n-tuple is quantized as a unit.
- Mean-squared distortion: $E[(U - V)^2]$

Quantization

M-level Quantizer / Size of Q.A.

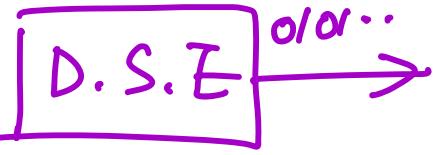
$$R(\bar{L}) \geq H(V)$$

Analog Sequence
 u_1, u_2, \dots, u_k



Discrete Sequence

v_1, v_2, \dots, v_k (Discrete R.V.)



$\{u_k\}$ iid Continuous R.V.

$$u \sim f_u(u) \text{ (PDF)}$$

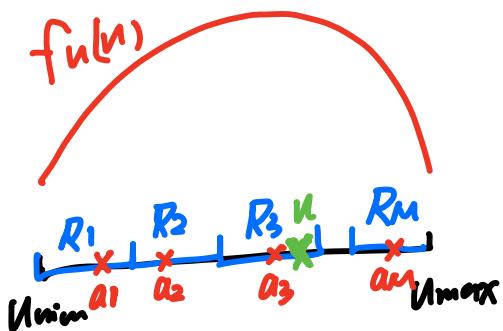
given $u \sim f_u(u)$, $\boxed{M, \{R_j\}_{j=1}^M, \{a_j\}_{j=1}^M}$

Quantization: $u \rightarrow v \in \{a_1, \dots, a_M\}$

$$f_u(u) = PMF$$

Quantization Alphabet

$$\text{if } u \in R_j, v(u) = a_j, b_j$$



PMF of V : $P(V=a_i) = \int_{R_i} f_u(u) du = P(u \in R_i)$

$$P(V=a_2) = \int_{R_2} f_u(u) du = P(u \in R_2)$$

$$P(V=a_M) = \int_{R_M} f_u(u) du = P(u \in R_M)$$

\Rightarrow ① Entropy of $V = H(V)$
 $(H(V) \leq R(\bar{L}))$

② Quantization Error

MSE (Mean-Squared Error/

$$D = E_u[(u - v)^2]$$

Distortion

$$\propto \int |u(t) - \hat{v}(t)|^2 dt$$

MSE: given $f_u(u)$, Quantizer

$$\text{Method 1: } \text{MSE} = E_u[(u - v(u))^2] = \int |u - \underline{v(u)}|^2 f_u(u) du$$

$$= \int_{R_1} |u - a_1|^2 f_u(u) du + \int_{R_2} |u - a_2|^2 f_u(u) du + \dots + \int_{R_M} |u - a_M|^2 f_u(u) du$$

$$= \sum_{j=1}^M \int_{R_j} |u - a_j|^2 f_u(u) du$$

$$\text{Method 2: } \text{MSE} = E_v \left[E_{u|v} [|u - v(u)|^2 | v] \right]$$

$$= \sum_{j=1}^M P(v=a_j) E_{u|v=a_j} [|u - v|^2 | v=a_j]$$

$$= \sum_{j=1}^M P(v=a_j) \boxed{\int |u - a_j|^2 f_{u|v=a_j}(u) du} \triangleq f_j(u)$$

$$f_{x|A}(x) = \begin{cases} \frac{f_x(x)}{P(A)}, & \text{if } A \text{ is true} \\ 0, & \text{otherwise} \end{cases}$$

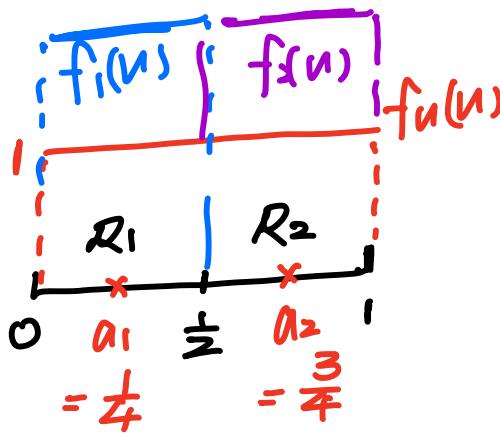
$$f_j(u) = \begin{cases} \frac{f_u(u)}{Q_j}, & \text{if } u \in R_j \\ 0, & \text{otherwise} \end{cases}$$

$$= \sum_{j=1}^M Q_j \int_{R_j} |u - a_j|^2 \frac{f_u(u)}{Q_j} du$$

$$= \sum_{j=1}^M \int_{R_j} |u - a_j|^2 f_u(u) du$$

Example

1. M=2 level Quantizer $U \sim U[0, 1]$



$$\textcircled{1} \text{ MSE} = \int_0^{\frac{1}{4}} (u - \frac{1}{4})^2 du + \int_{\frac{1}{4}}^1 (u - \frac{3}{4})^2 du$$

$$= \int_{-\frac{1}{4}}^{\frac{1}{4}} |u'|^2 du' + \int_{-\frac{1}{4}}^{\frac{1}{4}} |u'|^2 du' = \frac{1}{48}$$

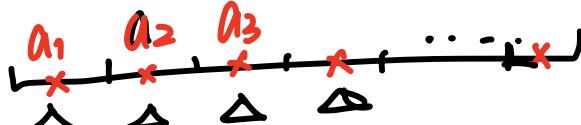
$$\textcircled{2} \quad f_1(u) = \begin{cases} \frac{1}{2} = 2, & \text{if } u \in R_1 \\ 0, & \text{otherwise} \end{cases} \quad f_2(u) = \begin{cases} 2, & u \in R_2 \\ 0, & \text{otherwise} \end{cases}$$

$$\text{MSE} = \frac{1}{Q_1} \underbrace{\int_0^{\frac{1}{4}} |u - \frac{1}{4}|^2 \times 2 du}_{\text{MSE}_1} + \frac{1}{Q_2} \underbrace{\int_{\frac{1}{4}}^1 |u - \frac{3}{4}|^2 \times 2 du}_{\text{MSE}_2}$$

2 uniform Quantizer

$f_u(u) \sim \text{uniform}, \Delta, \{ \text{middlepoint} \}$

$$f_j(u) = \begin{cases} \frac{f_u(u)}{f_u(u)\Delta} = \frac{1}{\Delta}, & \text{if } u \in R_j \\ 0, & \text{otherwise} \end{cases}$$



$$\text{PMF: } P(V = a_j) = \frac{1}{m}, \forall j.$$

$$H(V) = \log_2 M.$$

$$\begin{aligned} \text{MSE} &= \sum_j Q_j \int_{a_j - \frac{\Delta}{2}}^{a_j + \frac{\Delta}{2}} |u - a_j|^2 \frac{1}{\Delta} du \\ &= \sum_j Q_j \left[\int_{-\frac{\Delta}{2}}^{\frac{\Delta}{2}} |u'|^2 \frac{1}{\Delta} du' \right] \\ &= \frac{\Delta^2}{12} \end{aligned}$$

$$\begin{aligned} H(V) &= \log_2 M \quad (H(V) \leq R) \uparrow R \uparrow \begin{cases} M \uparrow \\ \text{Trade off} \end{cases} \\ \text{MSE} &= \frac{\Delta^2}{12} \quad \Delta \downarrow \text{MSE} \downarrow \end{aligned}$$

Quantization , $f_u(u)$, M

Scalar

vector

$$\textcircled{1} \quad \min_{\{R_j\}, \{a_{ij}\}} \text{MSE} = E[(u - v)^2] \quad \text{s.t.} \quad M$$

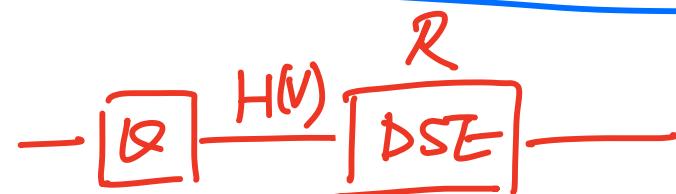
\Rightarrow Lloyd-Max Algorithm

scalar

vector

$$\textcircled{2} \quad \min_{\{a_{ij}\}} \text{MSE} = E[(u - v)^2] \quad \text{s.t.} \quad H(v) = R \quad \{R_j\}$$

\Rightarrow Entropy Code Quantization

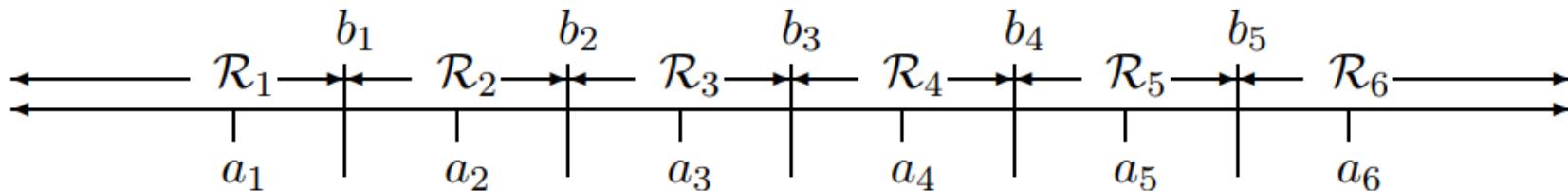


$$H(V) \leq R$$

$$y = \frac{H(V)}{R} \quad ^7$$

Scalar Quantization

- M level scalar quantization:
 - Partition the region R into M subsets $\mathcal{R}_1, \dots, \mathcal{R}_M$, called quantization region.
 - Each region \mathcal{R}_j is represented by a representation point $a_j \in R$.
- Scalar Quantization $\{v(u): R \rightarrow R\}: v(u) = a_j$, for $u \in \mathcal{R}_j$.
- Quantized output: $V_k = v(U_k)$ is a discrete rv with alphabet $\{a_1, \dots, a_M\}$.

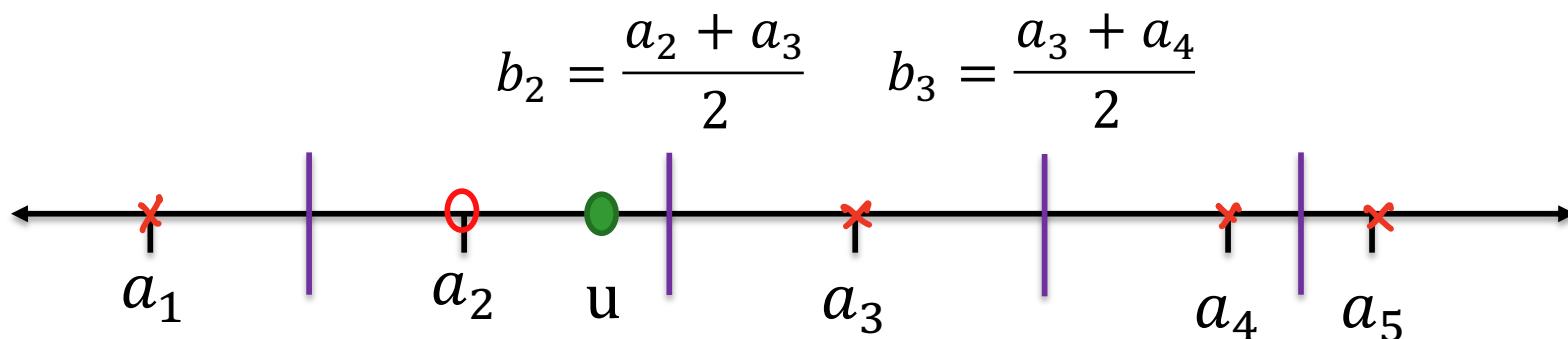


Scalar Quantization

- Q: how to choose the regions and representation points?
- Mean square distortion:
$$\text{MSE} = E[(U - V)^2]$$
- Aim: Given pdf $f_U(u)$ and alphabet size M , choose $\{\mathcal{R}_j, 1 \leq j \leq M\}$ and $\{a_j, 1 \leq j \leq M\}$ to minimize MSE.
- Explore it into two ways:
 - Given a set of $\{a_j, 1 \leq j \leq M\}$, how should the intervals $\{\mathcal{R}_j, 1 \leq j \leq M\}$ be chosen?
 - Given a set of intervals $\{\mathcal{R}_j, 1 \leq j \leq M\}$, how to choose $\{a_j, 1 \leq j \leq M\}$?

Scalar Quantization

- Choice of $\{\mathcal{R}_j\}$ given $\{a_j\}$



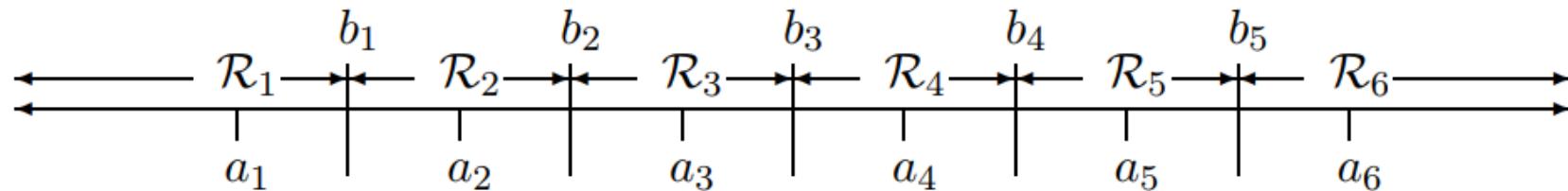
$$\min_j (u - a_j)^2$$

The boundary of b_j between \mathcal{R}_j and \mathcal{R}_{j+1} must lie halfway between a_j and a_{j+1} .

$$b_j = \frac{a_j + a_{j+1}}{2}$$

Scalar Quantization

- Choice of $\{a_j\}$ given $\{\mathcal{R}_j\}$



Method 1

$$\text{MSE} = E[(U - V)^2] = \int_{-\infty}^{\infty} f(u)(u - v)^2 du = \sum_{j=1}^M \int_{\mathcal{R}_j} f(u)(u - a_j)^2 du$$

$$= \sum_{j=1}^M Q_j \int_{\mathcal{R}_j} f_j(u)(u - a_j)^2 du = \sum_{j=1}^M Q_j \left[\int_{-\infty}^{\infty} f_j(u)(u - a_j)^2 du \right]$$

Method 2

where $Q_j = \Pr(U \in \mathcal{R}_j) = \int_{b_{j-1}}^{b_j} f(u) du$ and

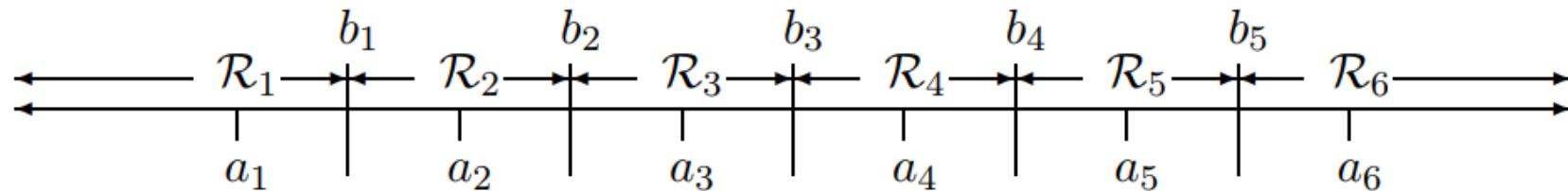
the conditional pdf of U given that $\{u \in \mathcal{R}_j\}$

$$f_j(u) = \begin{cases} f(u)/Q_j, & \text{if } u \in \mathcal{R}_j \\ 0, & \text{otherwise} \end{cases}$$

MSE_j

Scalar Quantization

- Choice of $\{a_j\}$ given $\{\mathcal{R}_j\}$



$$\text{MSE} = \sum_{j=1}^M Q_j \int_{-\infty}^{\infty} f_j(u) (u - a_j)^2 du \quad U(j) \sim f_j(u)$$

Let $U_{(j)}$ be a new RV with pdf $f_j(u)$, then

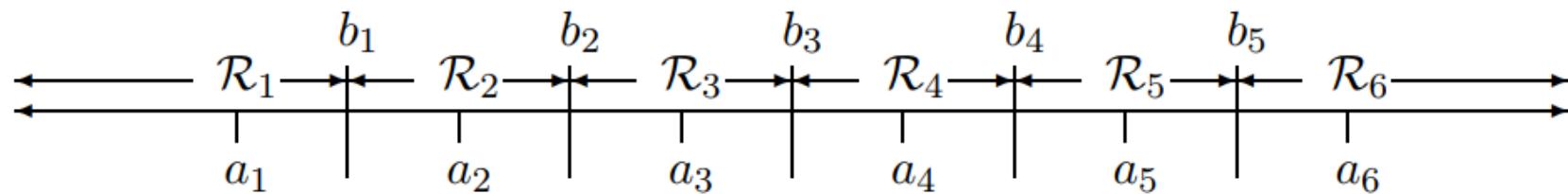
$$\text{MSE} = \sum_{j=1}^M Q_j \frac{E[(U_{(j)} - a_j)^2]}{b_j}, \quad \min_{a_j} \sum_{j=1}^M Q_j E[(U_{(j)} - a_j)^2]$$

$$\min_{a_j} E[(U_{(j)} - a_j)^2] = E[(U_{(j)}^2 - (E[U_{(j)}])^2 + (E[U_{(j)}])^2 - 2a_j U_{(j)} + a_j^2)] \\ = \text{Var}(U_{(j)}) + (E[U_{(j)}] - a_j)^2$$

Thus,

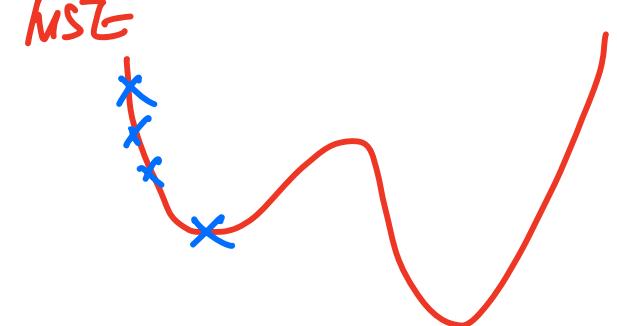
$$a_j = E[U_{(j)}] = \int_{R_j} u f_j(u) du = \frac{\int_{R_j} u f(u) du}{\int_{R_j} f(u) du}$$

Lloyd-Max Algorithm



An optimal scalar quantizer must satisfy both $b_j = (a_j + a_{j+1})/2$ and $a_j = \underline{E[U_{(j)}]}$.

1. Choose $a_1 < a_2 < \dots < a_M$
2. set $b_j = (a_j + a_{j+1})/2$ for $1 \leq j \leq M - 1$
3. Set $a_j = \underline{E[U_{(j)}]} = \frac{\int_{\mathcal{R}_j} u f(u) du}{\int_{\mathcal{R}_j} f(u) du}$ where $\mathcal{R}_j = (b_{j-1}, b_j]$ for $1 \leq j \leq M - 1$
4. Iterate on 2 and 3 until improvement is negligible.



It finds local min, not necessarily global min.

Scalar Quantization

- Example: Analog source U_k is uniformly distributed between 0 and 1. Find the 1 bit quantizer that minimizes the MSE. Find the MSE. $M=2$

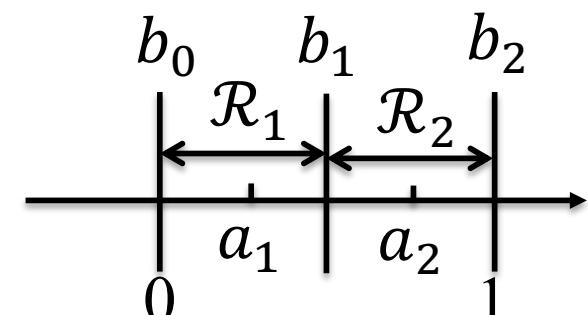
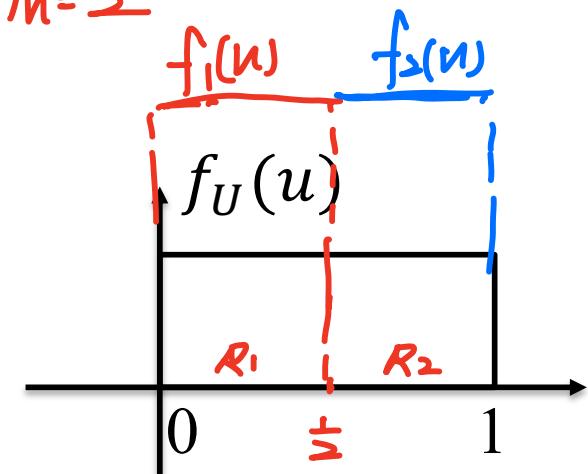
- Sol:

$$f_1(u) = \frac{1}{\frac{1}{2}} = 2 \quad \text{if } u \in [0, \frac{1}{2}]$$

- $f_j(u) = \frac{f_U(u)}{\int_{R_j} f_U(u) du} = \begin{cases} \frac{1}{\Delta} = 2, & 0 \leq u \leq \frac{1}{2}, j = 1 \\ \frac{1}{\Delta} = 2, & \frac{1}{2} \leq u \leq 1, j = 2 \end{cases}$

- $a_1 = 1/4, a_2 = 3/4.$

- $\underline{\text{MSE}} = E(U - V)^2 = \int_0^{\frac{1}{2}} (u - a_1)^2 f_U(u) du + \int_{\frac{1}{2}}^1 (u - a_2)^2 f_U(u) du = \frac{\Delta^2}{12} = \frac{1}{48}.$

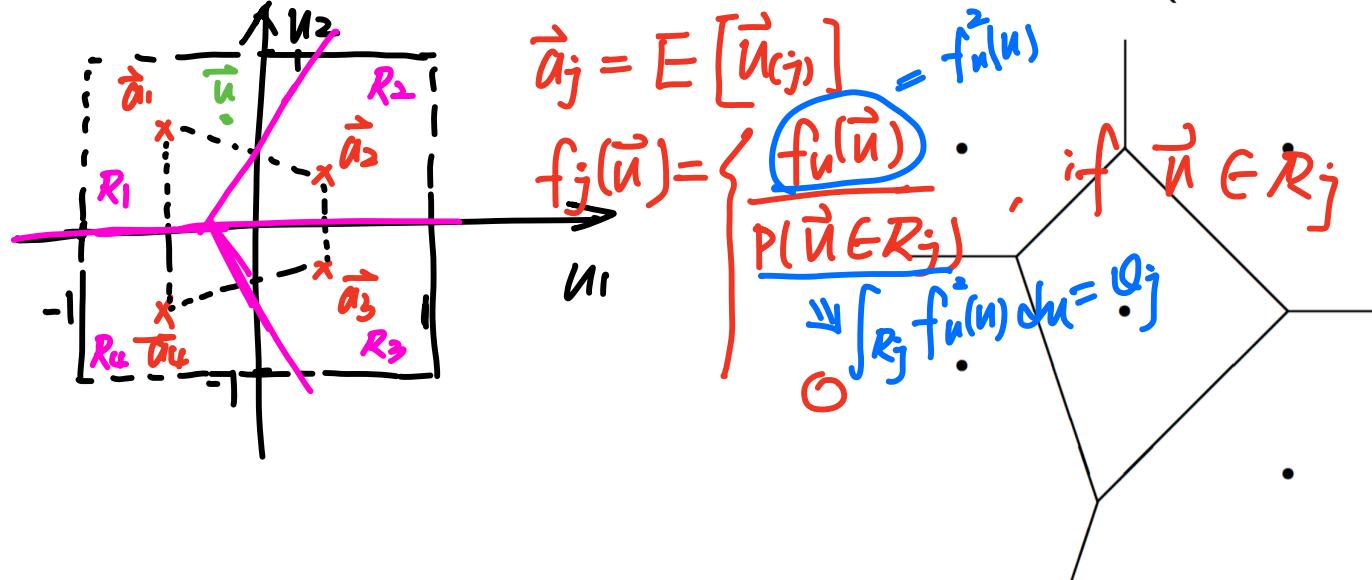


$u_1, u_2 \dots$ u_K

$\{u_k\}$ iid $u_k \sim U[-1, 1]$ $n=2, M=4$ Vector Quantization

2-dimensional Lloyd-Max Alg

Quantize n source variables at a time. (In scalar quantization, $n = 1$)



Given $\{(a_j, a'_j)\}$, how to choose $\{\mathcal{R}_j\}$

- The square error is $(u - a_j)^2 + (u' - a'_j)^2$, the point $\{a_j, a'_j\}$ which is the closest to (u, u') in Euclidean distance should be chosen.
- $\{\mathcal{R}_j\}$ contains points that are closer to (a_j, a'_j) than any other representation points, i.e., Voronoi regions.

Given a set of Voronoi region, how to find the $\{a_j, a'_j\}$?

- Choose $\{a_j, a'_j\}$ to be the conditional means within those regions.

Vector Quantization

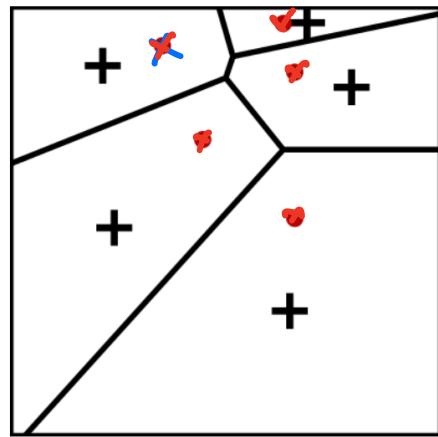


Figure: Iteration 1

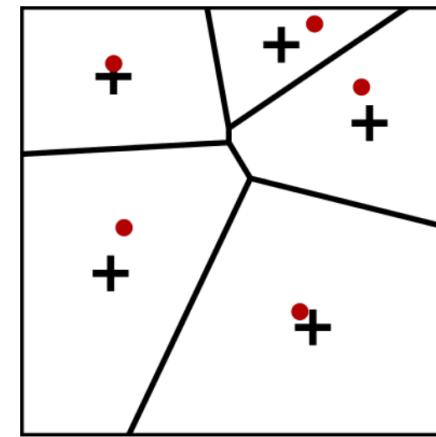


Figure: Iteration 2

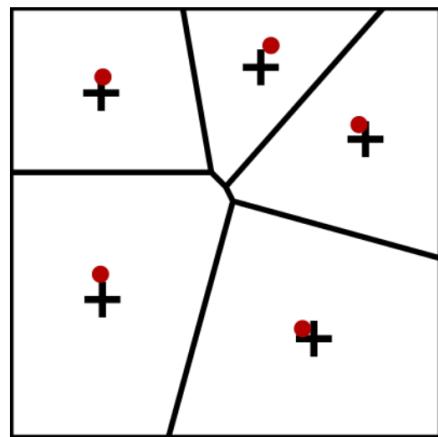


Figure: Iteration 3

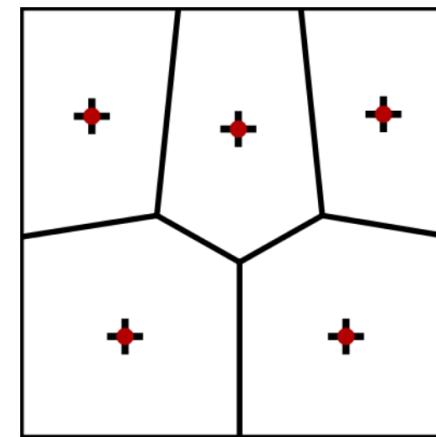
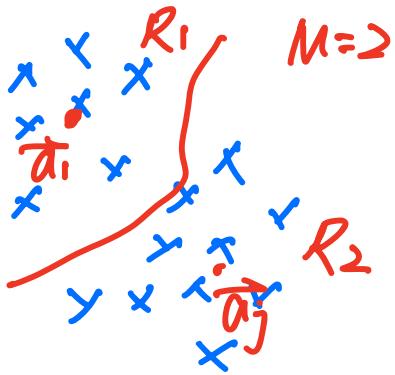
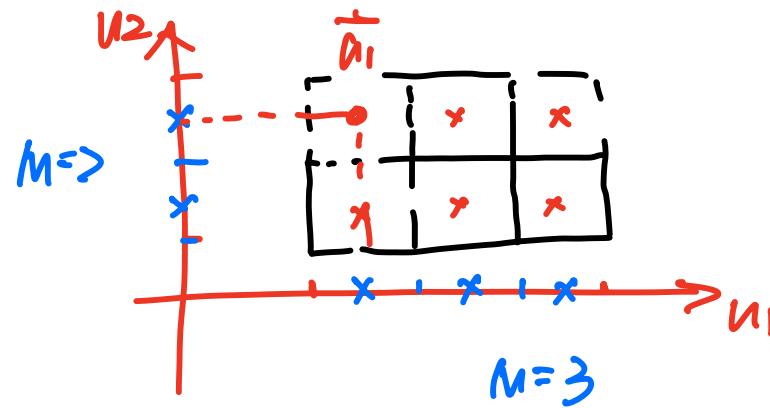


Figure: Iteration 4



Vector Quantization

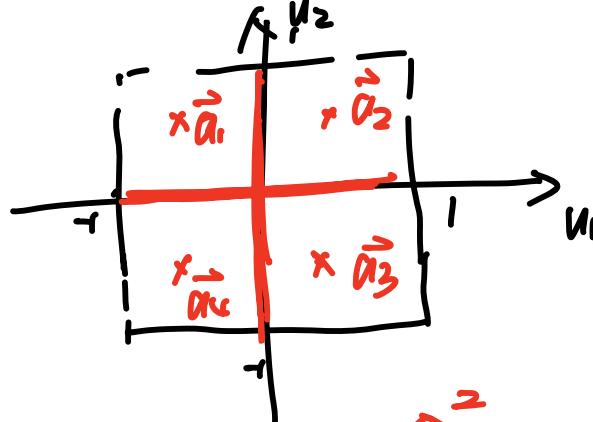
- Popular research topic, related to deep learning algorithm [1) clustering (2) Neural Discrete Representation Learning (VAE / VQ-VAE)]
- Quantizing complexity goes up exponentially with n
- For vector quantization, the problem of local minimum is even worse.
- Reduction in MSE with increasing n is quite modest
- Samples are highly dependent
- Application: Video, Audio

u_1, u_2, \dots, u_k $n=2$ $M=6$ 

rectangular Quantizer

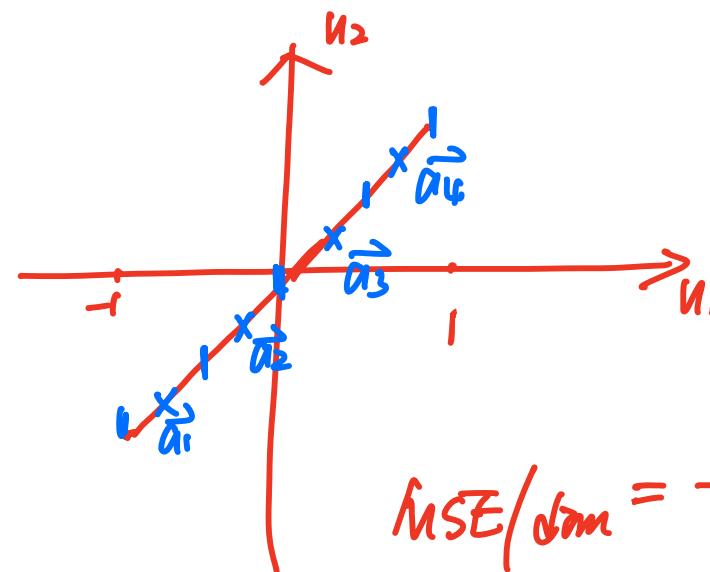
 $u_1, u_2, \dots, u_k \quad n=2 \quad M=4$

① 4 i.i.d. $u_k \sim U[-1, 1]$



$$MSE = \frac{\Delta^2}{12} \quad (\Delta = 1)$$

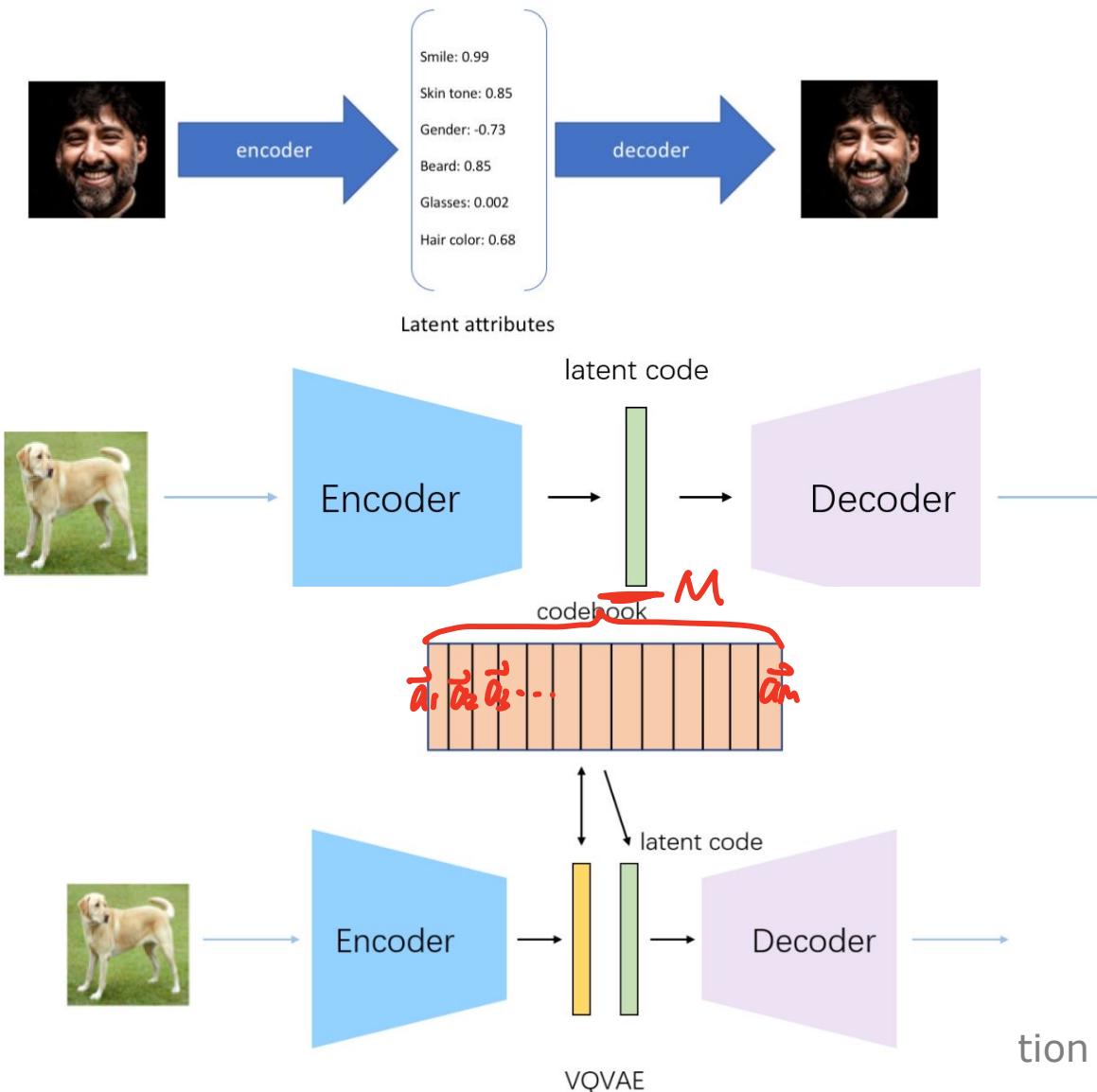
② $u_2 = u_1, u_3 = u_4, \dots$
 $u_1 \sim U[-1, 1]$



$$MSE/\Delta m = \frac{\Delta^2}{12} \quad (\Delta = \frac{1}{2})$$

Vector Quantization

- Image Reconstruction / Data Transmission



Entropy-Code Quantization

Quantization:

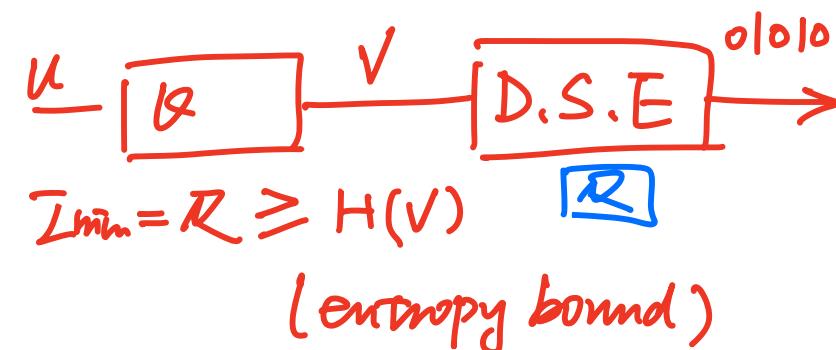
- Source sequence (samples): $(\underline{U_1}, \underline{U_2}, \dots)$
- Quantizer's Output: $(\underline{V_1}, \underline{V_2}, \dots)$, with $\underline{V_k} \in \{a_1, a_2, \dots, a_M\}$.
- If $\underline{U_k} \in \mathcal{R}_j$, then $\underline{V_k} = a_j$

(PMF)

Discrete Source Encoding:

$$P(V_k=a_j) = (\underline{Q_j}, \underline{B_j})$$

$$\begin{array}{l} a_1 \rightarrow 0 \\ \overline{a_2} \rightarrow \overline{10} \\ \overline{a_3} \rightarrow \overline{11} \end{array}$$



\bar{L}_{\min} minimum average code length: $\bar{L}_{\min} \approx H(V)$

We should minimize MSE for given entropy $H(V)$ rather than a given number of representation points.

Entropy-Code Quantization

- **Problem:** For a signal u with given pdf $f_U(u)$, find a quantizer such that

what is Entropy Code Quantization

① given $H(V)$, minimize MSE

② \boxed{R} \checkmark D.S.E / Entropy Coding

$$\text{minimize } E[(U - V)^2]$$

$$\text{subject to } R = H(V)$$

$\{R_j\}$ $\{a_j\}$

$\{R_j\}$

where $H(V) = - \sum_{i=1}^M p(a_i) \log(p(a_i))$ with $p(a_i) = \int_{R_j} f(u) du$
 $\boxed{\lim_{n \rightarrow \infty} R = H(V)}$

- **Solution:** Lagrangian cost function [Skip, See P. 85 in Gallager'Book]

$$L = E[(U - V)^2] + \lambda H(V)$$

- **Compare with Lloyd-Max:** $H(V)$ depends on $\overbrace{p(a_i)}$, not $\overbrace{\{a_j\}}$

- For given $\{R_j\}$, $a_j = E[U_{(j)}]$ minimizes MSE without changing $H(V)$.
- For given $\{a_j\}$, b_j minimizes MSE, but may change $H(V)$

Assumption : High-rate Quantizer

① R is Large

② $H(V)$ is Large

③ M is Large

④ MSE is small

⑤ * Quantization Region is Small \rightarrow

$$f_u(u) \approx \bar{f}(u) \leftarrow \boxed{\int_{R_j} f_u(u) du = \bar{f}(u) |R_j|} \text{ give } h(u), \text{ coding rate } R.$$

instant

opt. High-rate, Entropy-coded Quantizer

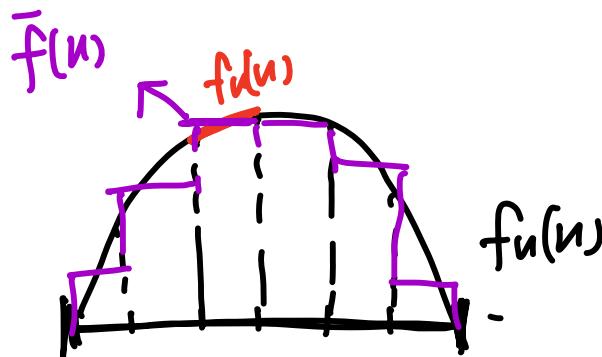
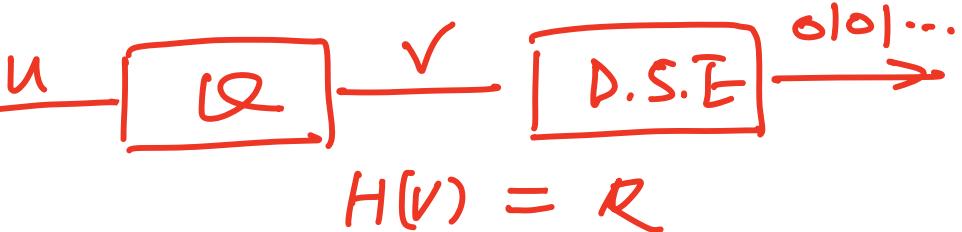
Uniform: Δ ?

$$H(V) = h(u) - \log_2 \Delta \leq R$$

$$\Delta \geq 2^{h(u) - R}$$

$$M = \frac{1}{\Delta} \quad a_j = \underline{\text{midpoint}}$$

$$\text{MSE}^* \equiv \frac{1}{2} \Delta^2$$



P85.

High-rate quantizer (Large $H(V)$)

min MSE

s.t. $H(V) = R$

Highrate: $H(V)$ is Large

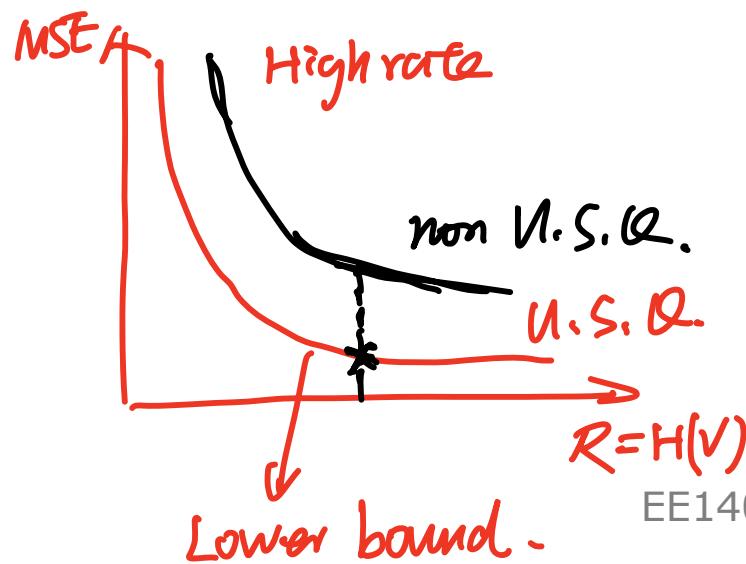
⇒ uniform scalar quantizer is a nearly optimal entropy-coded quantizer

⇒ uniform scalar quantizer approx. minimizes MSE, s.t. entropy constraint

\downarrow

$\Delta, \Delta, \Delta, \dots$ \Rightarrow optimal High-rate Entropy-coded quantizer

$H(V) \text{ is Large}$ $H(V) = R$



① given U, Δ . { MSE =
 $H(V) =$

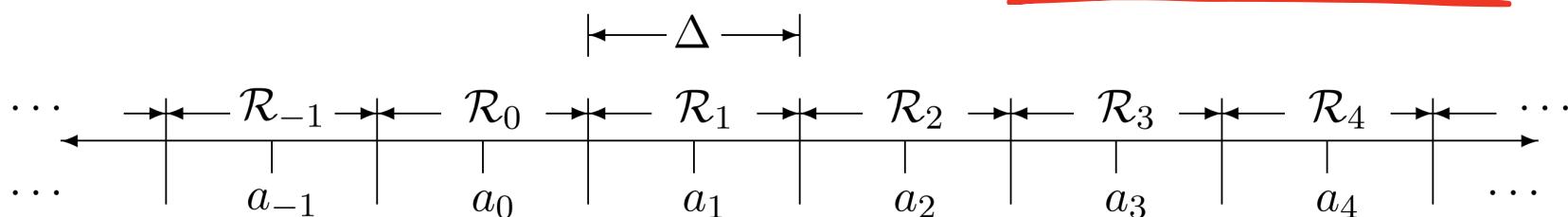
② { trade off perf. { MSE
H(V)
given U, R , \Rightarrow { $\Delta = ?$ $M = ?$
opt. E.C.Q. { $a_{ij} = ?$

Uniform Scalar Quantizer

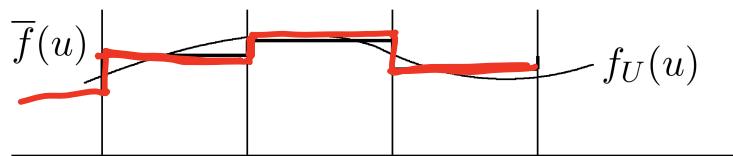
- Uniform scalar quantizer:
 - All quantization intervals have the same lengths.
- High-rate quantizer:
 - Quantization regions can be sufficiently small, pdf $f_U(u)$ is almost constant within each region.
 - Target entropy is large (high rate).
- Uniform scalar quantizer is nearly optimal in terms of MSE within the class of scalar quantizer when using
 - High-rate quantizers ($H(v)$ is large)
 - Entropy-coded quantization ($H(v) = R$)
 - Proof: See P. 85 in Gallager'Book

High-Rate Entropy-Coded Quantization

- Analyze the performance of Uniform scalar quantizer in the limit of high rate.
 $\left. \begin{matrix} \text{MSE} \\ H(v) \end{matrix} \right\}$
- Assume each quantization interval \mathcal{R}_j has same length Δ .



- Assume Δ is small enough that $f_U(u)$ is approx. constant within each \mathcal{R}_j .



$$\bar{f}(u) = \frac{\int_{\mathcal{R}_j} f_U(u) du}{\Delta} \quad \text{for } u \in \mathcal{R}_j$$
$$\Delta \bar{f}(u) = \Pr(\mathcal{R}_j) \quad \text{for all integer } j \text{ and all } u \in \mathcal{R}_j$$

- High rate assumption: $f_U(u) \approx \bar{f}(u) = \frac{\Pr(\mathcal{R}_j)}{\Delta}$ for all $u \in \mathcal{R}_j$.

High-Rate Entropy-Coded Quantization

- High rate assumption: $f_U(u) \approx \bar{f}(u)$ for all $u \in R$.

- Conditional pdf:

$$f_j(u) = f_{U|\mathcal{R}_j}(u) \approx \begin{cases} \frac{1}{\Delta}, & u \in \mathcal{R}_j; \\ 0, & u \notin \mathcal{R}_j. \end{cases} = \frac{\int_{\mathcal{R}_j} f(u) du}{\int_{\mathcal{R}_j} f(u) du} \underset{\approx}{\sim} \frac{\bar{f}(u)}{\bar{f}(u) \cdot \Delta}$$

- Conditional mean a_j : center of interval \mathcal{R}_j .

- MSE:

$$\text{MSE}(\mathcal{R}_j) \approx \sum_{j=1}^M \Pr(\mathcal{R}_j) \frac{\Delta^2}{12} = \frac{\Delta^2}{12}$$

$$\text{MSE}_j = \int_{a_j - \frac{\Delta}{2}}^{a_j + \frac{\Delta}{2}} |u - a_j|^2 \frac{1}{\Delta} du = \frac{\Delta^2}{12}$$

- If the decoding is successful, the MSE in reconstructing the original sequence is $\text{MSE} \approx \frac{\Delta^2}{12}$.

- Entropy of V : $p(a_j) = \int_{\mathcal{R}_j} f(u) du = \bar{f}(u) \Delta$

$$H(V) = \sum_j -p(a_j) \log p(a_j) = \sum_j \int_{\mathcal{R}_j} -f_U(u) du \log [\bar{f}(u) \Delta]$$

$$= \int_{-\infty}^{\infty} -f_U(u) \log [\bar{f}(u) \Delta] du$$

$$= \int_{-\infty}^{\infty} -f_U(u) \log [\bar{f}(u)] du - \log \Delta$$

Using high-rate approximation $f_U(u) \approx \bar{f}(u)$

High-Rate Entropy-Coded Quantization

- High rate assumption: $f_U(u) \approx \bar{f}(u)$ for all $u \in R$.

- Conditional pdf:

$$f_{U|\mathcal{R}_j}(u) \approx \begin{cases} 1/\Delta, & u \in \mathcal{R}_j; \\ 0, & u \notin \mathcal{R}_j. \end{cases}$$

- Conditional mean a_j : center of interval \mathcal{R}_j .

- MSE:

$$\text{MSE}(\mathcal{R}_j) \approx \sum_{j=1}^M \Pr(\mathcal{R}_j) \frac{\Delta^2}{12} = \frac{\Delta^2}{12}$$

- If the decoding is successful, the MSE in reconstructing the original sequence is $\text{MSE} \approx \frac{\Delta^2}{12}$.

- Entropy of V:

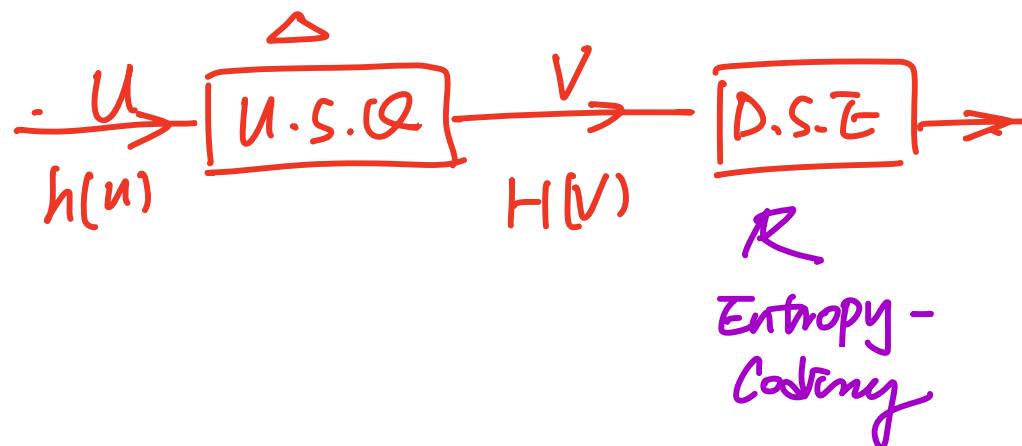
$$\begin{aligned} H(V) &\approx \int_{-\infty}^{\infty} -f_U(u) \log[f_U(u)]du - \log \Delta \\ &= h(U) - \log \Delta \end{aligned}$$

$$P(V=a_j) = \int_{R_j} f_u(u) du = \bar{f}(u) \Delta$$

$$H(V) = \sum_j -P(V=a_j) \log P(V=a_j)$$

$$= \sum_j - \int_{R_j} f_u(u) du \log (\bar{f}(u) \Delta)$$

$$= - \int f_u(u) \log (\underline{\bar{f}(u)} \Delta) du$$



High-rate

$$\approx - \underline{\int f_u(u) \log f_u(u) du} - \log \Delta$$

$$= h(u) - \log \Delta \quad \text{High rate}$$

Average coding rate $\underline{R} = \overline{I}_{\min} = H(V) \approx h(u) - \log \Delta \text{ bits/symbol}$

$$\left\{ \begin{array}{l} \overline{I} \uparrow \rightarrow h(u) \uparrow \Delta \downarrow \\ \overline{MSE} \downarrow \approx \frac{\Delta^2}{T_2} \downarrow \end{array} \right.$$

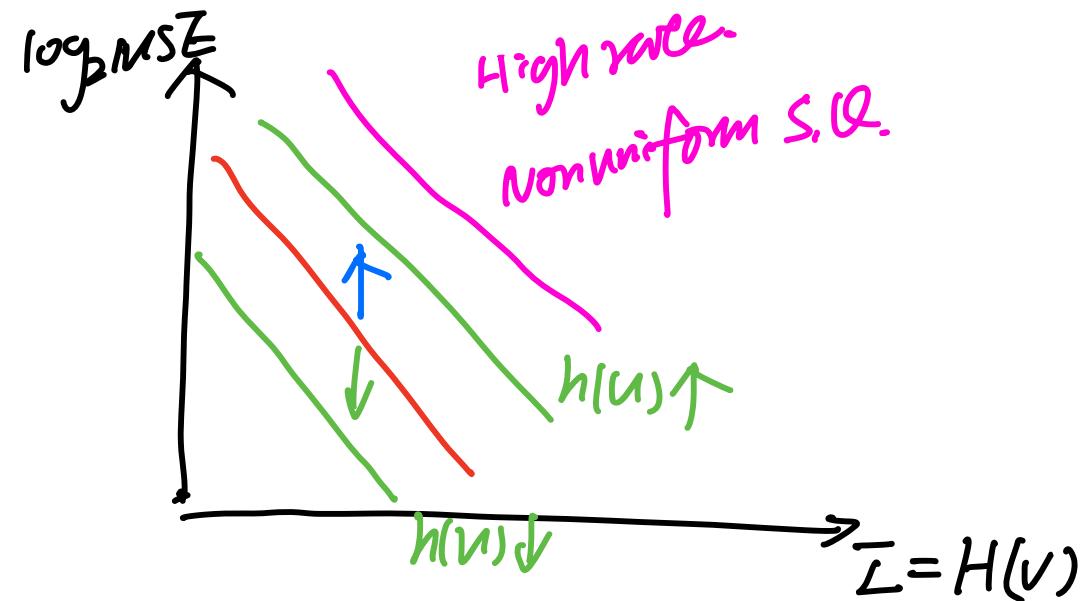
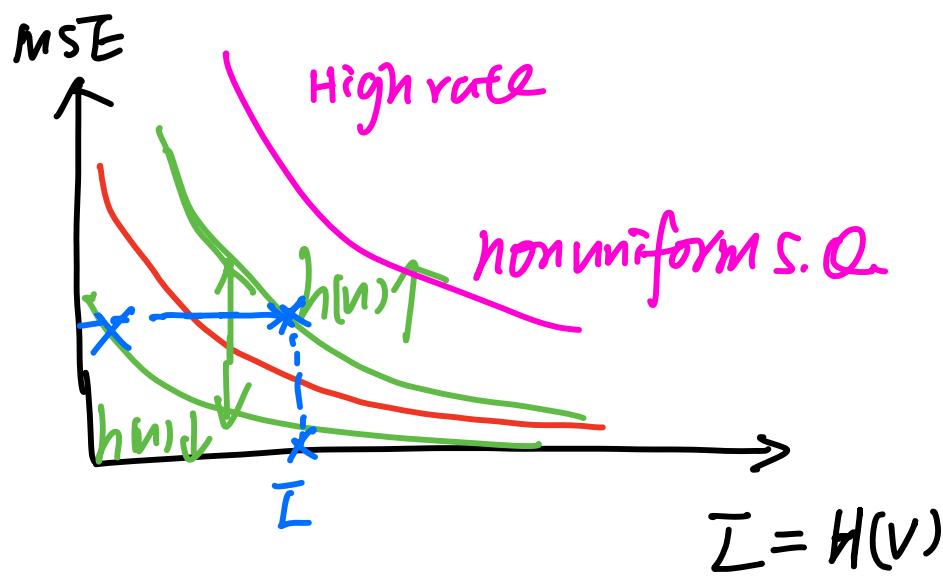
$$\bar{I} = h(U) - \log_2 \Delta \quad \left. \right\} \text{Trade off.} \Rightarrow \log_2 \Delta = h(U) - \bar{I}$$

$$MSE = \frac{\Delta^2}{12}$$

$$\log_2 MSE = 2 \log_2 \Delta - \log_2 12$$

$$= 2(h(U) - \bar{I}) - \log_2 12$$

$$\Delta = \frac{2(h(U) - \bar{I})}{12}$$



① $h(U)$: shift

② Δ : $\Delta \downarrow \frac{1}{2}\Delta$ $\bar{I} + 1 \text{ bits/symbol}$ $\nexists MSE \downarrow 6 \text{ dB}$

③ High rate same $\bar{I} = H(V)$, $MSE_{N.U.S.Q} > MSE_{U.S.Q.}$. 37

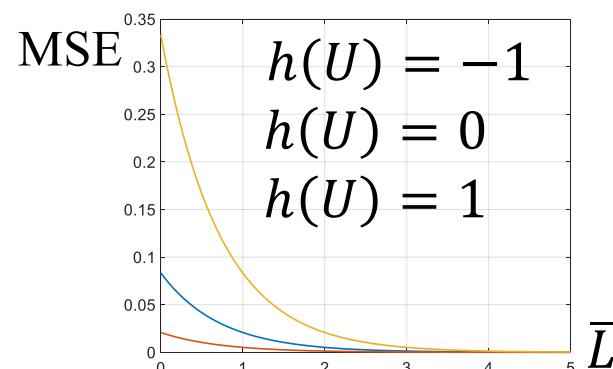
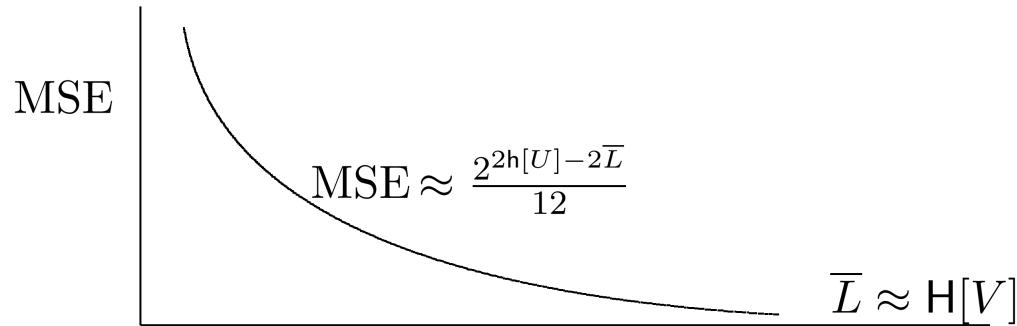
High-Rate Entropy-Coded Quantization

- Average Code Length:

- $\bar{L} \approx H(V) \approx h(U) - \log \Delta$ (bits/symbol)
- \bar{L} depends only on $h(U)$ and Δ , not on $f_U(u)$ or M .

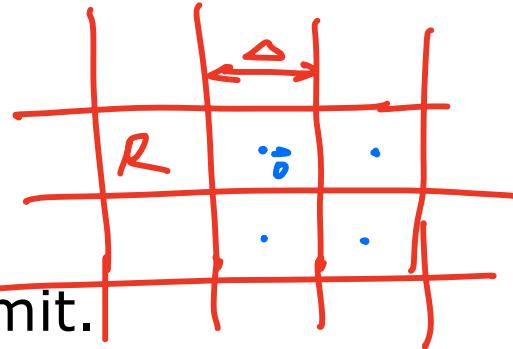
- Trade off between coding rate and MSE

$$\bar{L} \approx h(U) - \log \Delta; \quad \text{MSE} \approx \frac{\Delta^2}{12}$$



Reduction in Δ by a factor of 2, reduce MSE by a factor of 4, increase the \bar{L} by 1bit/symbol. (Each additional bit per symbol decrease the MSE by 6 dB.

High-Rate 2D Quantizer



- Uniform 2D quantizer is optimal in the high-rate limit.

- Proof: See P. 85 in Gallager'Book

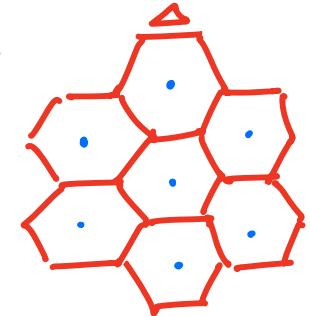
- Similar as 1D quantizer we have

$$\sum_{j=1}^M p(V=\vec{a}_j) MSE_j$$

- MSE per dimension:

$$MSE \approx \frac{1}{2} \int_{\mathcal{R}} \|u\|^2 \frac{1}{A(\mathcal{R})} du$$

$$f_j(\vec{u}) = \begin{cases} \frac{f(\vec{u})}{A(\mathcal{R})} & \text{if } \vec{u} \in \mathcal{R}_j \\ 0 & \text{otherwise} \end{cases}$$



$$MSE_j = \int_{\mathcal{R}_j} \| \vec{u} \|^2 \frac{1}{A(\mathcal{R})} d\vec{u}$$

$A(\mathcal{R}) = \int_{\mathcal{R}} d\vec{u}$ is the area of basic cell (square, hexagon, etc.)

- Square: $A(\mathcal{R}) = \Delta^2$, $MSE = \Delta^2/12$

- Hexagon: $A(\mathcal{R}) = 3\sqrt{3}\Delta^2/2$, $MSE = 5\Delta^2/24$

$\left(\frac{MSE}{A(\mathcal{R})} \right)_{\text{Square}} / \left(\frac{MSE}{A(\mathcal{R})} \right)_{\text{Hex}} = 1.0392$. For a given \bar{L} , hexagonal quantizer is slightly better. Thus, uniform scalar quantizers (square quantizer) is good enough.

High-Rate 2D Quantizer

- Entropy of Uniform 2D quantizer

With

$$p_j = \bar{f}(\mathbf{u})A(\mathcal{R}) = \int_{\mathcal{R}_j} f(\mathbf{u})d\mathbf{u}$$

Entropy of \vec{V} :

$$\begin{aligned} \underline{H(\vec{V})} &= - \sum_j \underline{p_j} \log \underline{p_j} \\ &= - \sum_j \int_{\mathcal{R}_j} f(\mathbf{u}) \log [\bar{f}(\mathbf{u})A(\mathcal{R})] d\mathbf{u} \\ &= - \int f(\mathbf{u}) [\log \bar{f}(\mathbf{u}) + \log A(\mathcal{R})] d\mathbf{u} \\ &\stackrel{(a)}{\approx} 2h(U) - \log A(\mathcal{R}) \end{aligned}$$

Trade off. $\left\{ \begin{array}{l} \mathbb{I} \uparrow \quad A(\mathcal{R}) \downarrow \quad \mathbb{I} \uparrow \text{MSE} \downarrow \\ \text{MSE} \end{array} \right.$

 $A(\mathcal{R}) \downarrow \nleq A(\mathcal{R})$
 $\mathbb{I} \uparrow 1 \text{ bit/symbol}$
 $\text{MSE} \downarrow 6 \text{ dB}$

where (a) follows from $\bar{f}(\mathbf{u}) \approx f(\mathbf{u})$ and $H(\mathbf{U}) = H(U_1 U_2) = 2H(U)$.

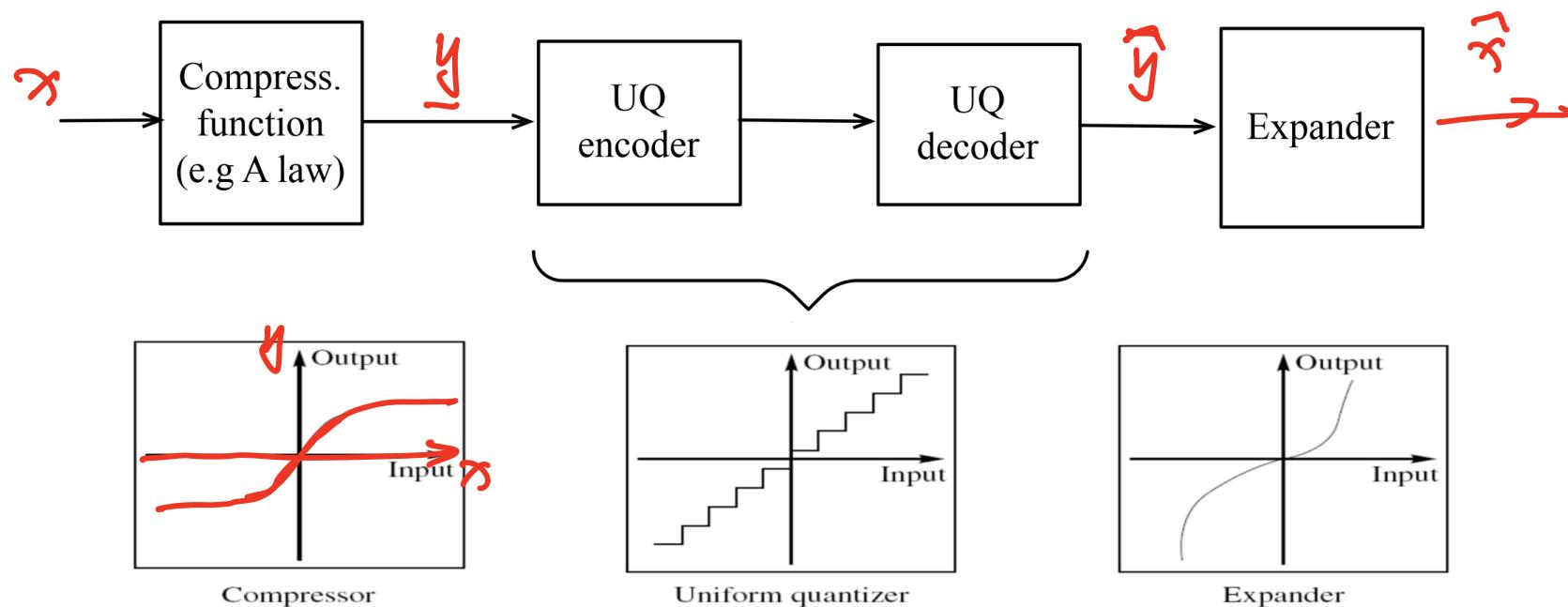
Thus, average rate per symbol:

$$\underline{\bar{L}} \approx \frac{H(\vec{V})}{2} \approx h(U) - \frac{1}{2} \log A(\mathcal{R})$$

give \mathbb{I} (given $A(\mathcal{R})$)
 $\text{MSE}_{\text{Hex}} < \text{MSE}_{\text{square}}$.

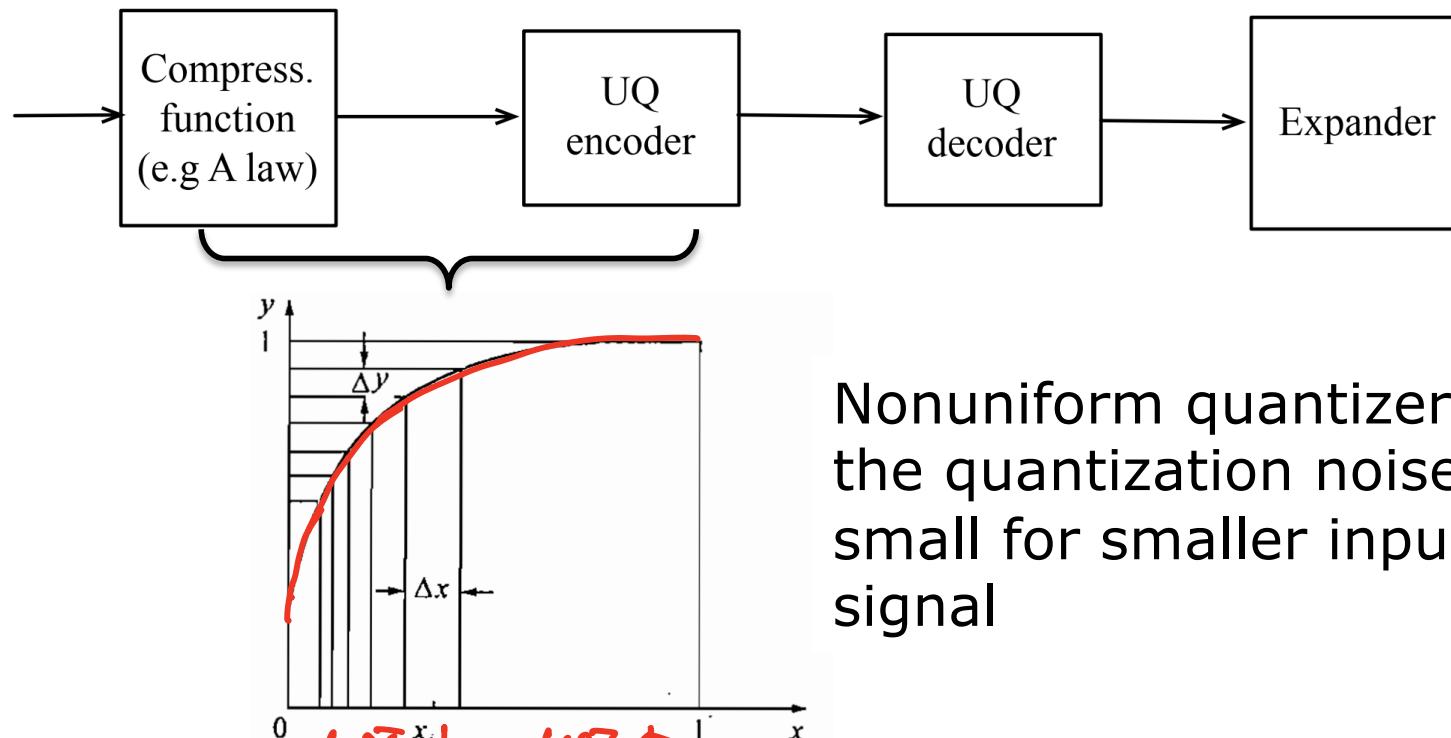
Quantization in Telephony Systems

- Given input x , first use compressor by mapping $x \rightarrow y$, then use uniform quantizer
- To restore the signal, use a device with a characteristic complementary to the compressor, called expander.
- compander=compressor+expander



Quantization in Telephony Systems

- Given input x , first use compressor by mapping $x \rightarrow y$, then use uniform quantizer
- To restore the signal, use a device with a characteristic complementary to the compressor, called expander.
- compander=compressor+expander



Nonuniform quantizer:
the quantization noise is
small for smaller input
signal

Quantization in Telephony Systems

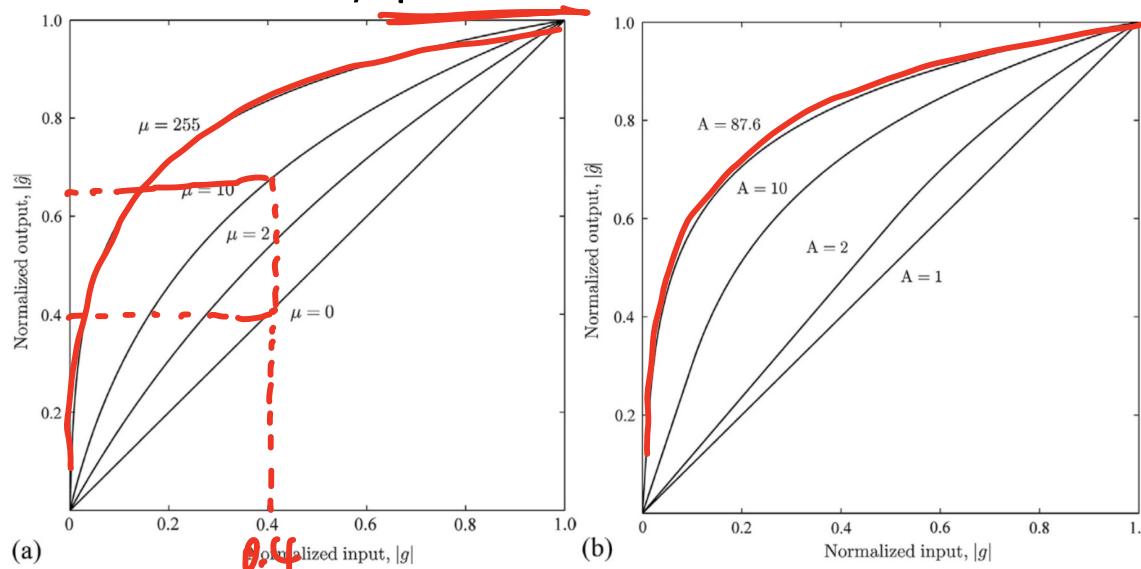
- μ law (Applied in North America and Japan): for $\mu > 0$

$$y = \frac{\log(1 + \mu|x|)}{\log(1 + \mu)}$$

- A law (Applied in elsewhere): for $A > 0$

$$y = \begin{cases} \frac{Ax}{1+\ln A}, & 0 < x \leq \frac{1}{A} \\ \frac{1+\ln Ax}{1+\ln A}, & \frac{1}{A} < x \leq 1 \end{cases}$$

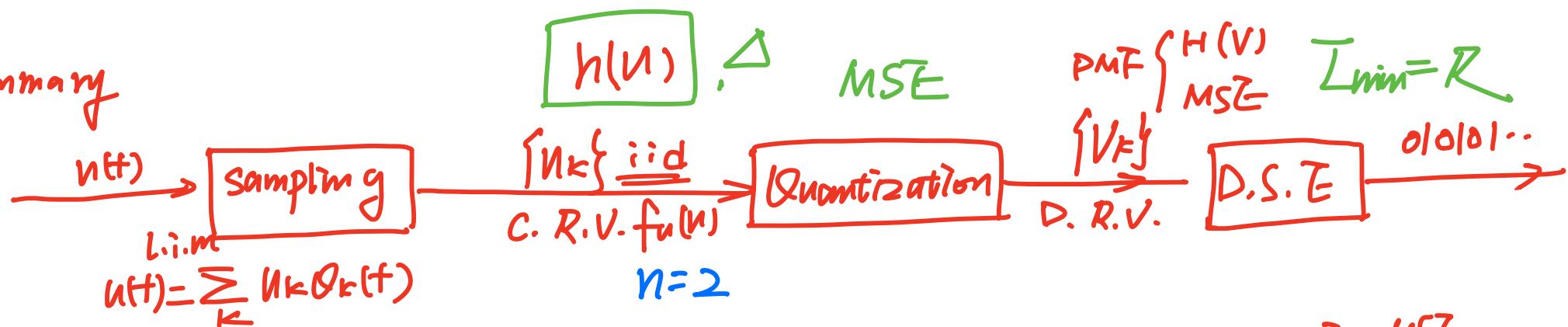
$A=87.6, \mu=255.$



Assume $M = 2^8 = 256$
If $x=0.4$

$$\begin{aligned} M_{\mu=0} &= 256 \times 0.4 \approx 102 \\ M_{\mu=10} &= 256 \times 0.65 \approx 166 \\ M_{A=10} &= 256 \times 0.71 \approx 181 \end{aligned}$$

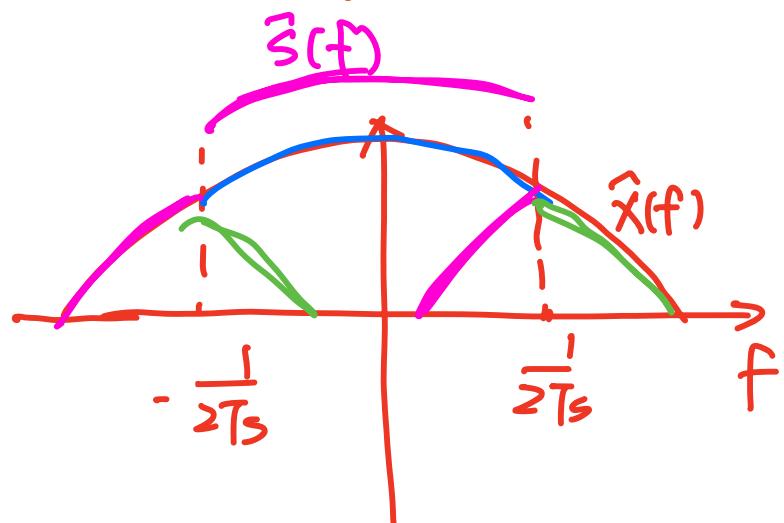
Summary



Time limited/unlimited (F.S.E)

Band limited/unlimited (sampling theory)
($T_s = \frac{1}{2w}$)

Aliasing ($W > \frac{1}{2T_s}, T_s > \frac{1}{2W}, f_s < 2W$)



$f_u(u), M, \min \text{MSE}$
 \Rightarrow Lloyd-Max Alg
 $(\{b_j\}, \{a_j\} = E[u_j])$

$f_u(u), R, \min \text{MSE}$
 s.t. $H(v) = R$
 (High-rate)
 \Downarrow

uniform - scalar quantizer
 \Rightarrow
 \Downarrow

Trade off $\left\{ \begin{array}{l} \frac{\text{MSE}^*}{H(v)^*} = \frac{\text{MSE}}{H(v)} \\ H(v)^* = I_{\min} = \underline{\quad} \end{array} \right.$



上海科技大学
ShanghaiTech University

Thanks for your kind attention!

Questions?