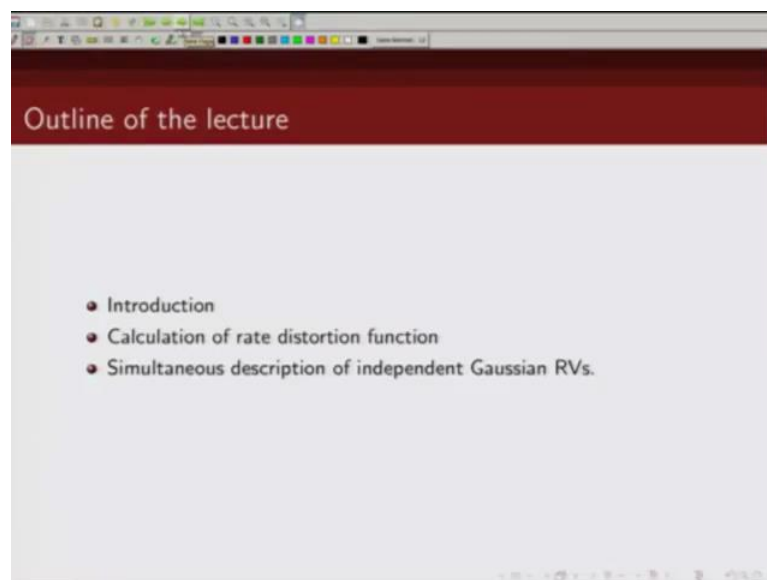


An Introduction to Information Theory
Prof. Adrish Banerjee
Department of Electronics and Communication Engineering
Indian Institute of Technology, Kanpur

Lecture - 13
Rate Distortion Theory

Welcome to the course on An Introduction to Information Theory. So, in this lecture we are going to talk about Rate Distortion Theory. When we try to represent an arbitrary real number to represent it losslessly, we would require infinite number of bits. Now, whenever we try to represent it using finite number of bits, we are introducing distortion. So, to know that our representation of this real number is good, we need to introduce this notion of distortion measure and in rate distortion theory, we are going to ask questions like for example, if we specify that we can tolerate this much average distortion, then what is the minimum number of bits required to represent a source or let us say if I specify my rate, then what is the minimum distortion that I can achieve.

(Refer Slide Time: 01:28)



So, we will start up this discussion on rate distortion theory with a brief introduction and few definitions and we will define what is a rate distortion code and what do we mean by a rate distortion code is achievable and then we will calculate rate distortion function for

some simple examples like Bernoulli source with hamming distortion measure, Gaussian source with squared error distortion measure and we going to abort the description of independent Gaussian random variables, but they are not necessarily identically distributed.

(Refer Slide Time: 02:06)

Introduction

- Distortion measure is a measure of distance between the random variable and its representation.
- A distortion measure is a mapping

$$d : X \times \hat{X} \rightarrow \mathbb{R}^+$$
 from the set of source alphabet-reproduction alphabet pairs into the set of nonnegative real numbers.
- A distortion measure is said to be bounded if the maximum value of the distortion is finite

$$\underline{d_{\max}} = \max_{x \in X, \hat{x} \in \hat{X}} \underline{d(x, \hat{x})} < \infty$$

So, a Distortion measure is a measure of distance between the random variable and its representation. When you are trying to represent an arbitrary real number using fixed number of bits, we are introducing some sort of a distortion. So, distortion measure is a measure of distance between this random variable and its representation. Now, this measure of distance, we will specify that there could be different ways in which we could define this measure of distance between the source alphabet and the re-produce alphabet.

So, a distortion measure is a mapping of set of source alphabet and reproduction alphabet pairs into a set of non-negative real number and this set of non-negative real number will specify how much this source alphabet is differing from this reproduced alphabet. We say a distortion measure is bounded if the maximum value of distortion is finite. So, we define the maximum value of distortion as maximum value of distance between the random variable and its representation. Now, if this is bounded, then we say the distortion measure is bounded.

(Refer Slide Time: 03:48)

Introduction

- Examples of common distortion functions:
 - Hamming distortion
$$d(x, \hat{x}) = \begin{cases} 0 & \text{if } x = \hat{x} \\ 1 & \text{if } x \neq \hat{x} \end{cases}$$
 - Squared error distortion
$$d(x, \hat{x}) = (x - \hat{x})^2$$
- The distortion between sequences x^n and \hat{x}^n is defined by
$$d(x^n, \hat{x}^n) = \frac{1}{n} \sum_{i=1}^n d(x_i, \hat{x}_i)$$
- We are interested to know: given a source distribution and a distortion measure, what is the minimum expected distortion achievable at a particular rate?

Let us give some example of common distortion function, hamming distortion function. So, hamming distortion measure is defined as follows, when the source alphabet and the reproduce alphabet are same, there is no distortion that is distortion is 0 when the source alphabet and the reproduce alphabet are same. However, if they are different, then the distortion is 1. So, in hamming distortion it is zero, if there is no error and it is 1 if the source alphabet and the reproduce alphabet they are different.

Similarly, we could define squared error distortion. So, squared error distortion is basically the Euclidean distance between the source alphabet and this reproduce alphabet and this is our squared error distortion function. Now, the distortion measure that we have defined so far is on symbol by symbol basis and we could extend this definition for sequences as well. So, to extend this definition of distortion measure for sequences, this is the source sequence and this is the reproduce sequence, then that distortion measure can be written as average distortion measure per symbol. So, this is distortion between symbol source symbol X_i and the reproduced symbol \hat{X}_i , we sum over all m symbols and then we average it out and that is how we would define distortion measure for sequences.

Now, as I said in the beginning of this lecture in this rate distortion theory given a source

distribution and a distortion measure, we are interested in what is the minimum expected distortion for an achievable particular rate. So, given a source distribution and given a distortion measure, what is the minimum expected distortion that is achievable for a particular rate? We can ask this question in a different way also. Given a source distribution and a distortion measure and given average distortion criteria, we are interested in knowing the minimum rate that is achievable.

(Refer Slide Time: 06:56)

Definition

- A $(2^{nR}, n)$ rate distortion code consists of following encoding function

$$f_n : X^n \rightarrow \{1, 2, \dots, 2^{nR}\}$$
 and following decoding function

$$g_n : \{1, 2, \dots, 2^{nR}\} \rightarrow \hat{X}^n$$
- Distortion for the $(2^{nR}, n)$ rate distortion code is defined as

$$D = \sum_{x^n} p(x^n) \underbrace{d(x^n, \underbrace{g_n(f_n(x^n))}_{\hat{x}^n})}_{\hat{x}^n}$$

So, a $(2^{nR}, n)$ rate distortion code consists of following encoding function. This is your source sequence of n bit and this is your mapping to 1 of these rate distortion codes and at the decoder which is denoted by this function g of n given that you have received 1 of these indexes, you are interested in getting back than estimate of the sequence.

Now, distortion for a rate distortion code is defined like this. This is the distortion measure between the source sequence and this is the reproduce sequence which I can also write this as \hat{X}^n . So, this is basically the average distortion between X^n and \hat{X}^n . We say a rate distortion pair denoted by this rate R and this distortion D is said to be achievable if there exist a $(2^{nR}, n)$ rate distortion code with following encoding and decoding function as n tends to infinity, the expected distortion

is less than equal to the given distortion D .

(Refer Slide Time: 08:09)

Definition

- A rate distortion pair (R, D) is said to be achievable if there exists sequence of $(2^{nR}, n)$ rate distortion codes (f_n, g_n) with $\lim_{n \rightarrow \infty} E d(X^n, g_n(f_n(X^n))) \leq D$.
- The rate distortion region for a source is the closure of the set of achievable rate distortion pairs (R, D) .
- The rate distortion function $R(D)$ is the infimum of rates R such that (R, D) is in the rate distortion region of the source for a given distortion D .
- The distortion rate function $D(R)$ is the infimum of all distortions D such that (R, D) is in the rate distortion region of the source for a given rate R .

Now, we could similarly define a rate distortion region. So, a rate distortion region for a source is the closure of set of all achievable rate distortion pair, that is our rate distortion region and we define a rate distortion function as the infimum all rates such that this rate distortion pair is in the rate distortion region of the source for a given distortion.

Similarly, we can also define distortion rate function as the infimum of all distortions D such that this pair rate and distortion is in the rate distortion region of the source for a given rate R .

(Refer Slide Time: 09:56)

Definition

- The information rate distortion function $R^I(D)$ for a source X with distortion measure $d(x, \hat{x})$ is defined as

$$R^I(D) = \min_{p(\hat{x}|x): \sum_{x, \hat{x}} p(x)p(\hat{x}|x)d(x, \hat{x}) \leq D} I(X; \hat{X})$$

where minimization is over all conditional distributions $p(\hat{x}|x)$ for which the joint distribution $p(x, \hat{x}) = p(x)p(\hat{x}|x)$ satisfies the expected distortion constraint.

- The rate distortion function for an i.i.d. source X with distribution $p(x)$ and bounded distortion function $d(x, \hat{x})$ is equal to the associated information rate distortion function. Thus

$$R(D) = R^I(D) = \min_{p(\hat{x}|x): \sum_{x, \hat{x}} p(x)p(\hat{x}|x)d(x, \hat{x}) \leq D} I(X; \hat{X})$$

is the minimum achievable rate at distortion D .

Now, we define information rate distortion function for a source X with distortion measure given by D of X and \hat{X} . We define information rate distortion as minimum mutual information between this source alphabet X and this reproduced alphabet \hat{X} . So, this is a reproduce source where mutual information between X and \hat{X} , minimum of that and minimization is taken over all conditional distribution p of \hat{X} given X for which the joint distribution p of X and \hat{X} satisfies the expected distortion constraint. For this information rate distortion function, we minimize this mutual information between X and \hat{X} and this minimization is over all conditional distribution of p of \hat{X} given X such that over this joint distribution p of X and \hat{X} , this expected distortion constraint is satisfied.

Now, this information rate distortion function is also our rate distortion function. So, rate distortion function for an IID source with distribution p of X and bounded distortion function given by D of X and \hat{X} is equal to its associated information rate distortion function. In other words, the rate distortion function can be found by minimizing the mutual information between X and \hat{X} and we do this minimization over all conditional distribution of p of \hat{X} given X such that this joint distribution p of X and \hat{X} satisfies the expected distortion constrain.

(Refer Slide Time: 12:37)

Rate Distortion Function

- Let X be $N(0, \sigma^2)$. By rate distortion theorem, we have

$$R(D) = \min_{f(\hat{x}|x) : E[(X - \hat{X})^2] \leq D} I(X, \hat{X})$$
- We know that

$$\begin{aligned}
 I(X, \hat{X}) &= h(X) - h(X|\hat{X}) \\
 &= \frac{1}{2} \log(2\pi e) \sigma^2 - h(X - \hat{X}|\hat{X}) \\
 &\geq \frac{1}{2} \log(2\pi e) \sigma^2 - h(X - \hat{X}) \\
 &\geq \frac{1}{2} \log(2\pi e) \sigma^2 - h(N(0, E[(X - \hat{X})^2])) \\
 &= \frac{1}{2} \log(2\pi e) \sigma^2 - \frac{1}{2} \log(2\pi e) E[(X - \hat{X})^2] \\
 &\geq \frac{1}{2} \log(2\pi e) \sigma^2 - \frac{1}{2} \log(2\pi e) D \\
 &= \frac{1}{2} \log \frac{\sigma^2}{D}
 \end{aligned}$$

$h(X - \hat{X}) \geq h(X - \hat{X}|\hat{X})$

This rate distortion function is given by this and this is the minimum achievable rate for a distortion D . Now, let us illustrate or let us compute this rate distortion function. So, we will first take an example of a Gaussian source. We have a Gaussian source X denoted by $X \sim N(0, \sigma^2)$. According to the definition, the rate distortion function is given by minimizing the mutual information between X and \hat{X} and this minimization is over all conditional distribution of \hat{X} given X such that the expected distortion constraint is satisfied. Now, let us compute the mutual information between X and \hat{X} . So, where we are going to compute this rate distortion function is as follows, we are going to lower bound this mutual information between X and \hat{X} and then we are going to show that this lower bound is achievable and that is how basically we are going to show that it is a value of the rate distortion function. So, mutual information between X and \hat{X} , this is given by differential entropy of X minus differential entropy of X given \hat{X} .

Now, X is a Gaussian source $N(0, \sigma^2)$. So, its differential entropy is given by this expression $\frac{1}{2} \log(2\pi e \sigma^2)$. Now, we know that translation does not change differential entropy. So, differential entropy of X given \hat{X} is same as differential entropy of $X - \hat{X}$ given \hat{X} . Now, this is differential entropy of $X - \hat{X}$ condition on \hat{X} . We know that conditioning cannot increase

entropy. So, $H(X - \hat{X})$ is going to be greater than $H(X - \hat{X} | \hat{X})$ and in this step we are subtracting a larger quantity and that is why I have here lower bound.

Now, we know that given a second order of movement, Gaussian source has the maximum differential entropy and that is why this is upper bounded by differential entropy of a Gaussian source which means 0 and same variance which is given by expected value of $X - \hat{X}$ square. So, this greater than equal to comes because differential entropy of $H(X - \hat{X})$ is upper bounded by differential entropy of a Gaussian random variable with same second order movement. Now, this is same as this and we know the differential entropy of a Gaussian source. So, that is given by $\frac{1}{2} \log 2\pi e \text{ variances expected value of } X - \hat{X} \text{ whole square}$, that is this term and this is upper bounded by D . So, then we can write mutual information between X and \hat{X} to be $\frac{1}{2} \log 2\pi e \text{ times sigma square minus } \frac{1}{2} \log 2\pi e \text{ times } D$, that is because expected value of $X - \hat{X}$ square is upper bounded by D and by combining these terms, we get this.

(Refer Slide Time: 17:46)

Rate Distortion Function

- If $D \leq \sigma^2$, we choose $X = \hat{X} + Z$, $\hat{X} \sim N(0, \sigma^2 - D)$, $Z \sim N(0, D)$, where \hat{X} and Z are independent. For this we have

$$I(X, \hat{X}) = \frac{1}{2} \log \frac{\sigma^2}{D}$$

$$R = \frac{1}{2} \log \frac{\sigma^2}{D}$$

$$\log \frac{\sigma^2}{D} = 2R$$
- and $E(X - \hat{X})^2 = D$.
- If $D > \sigma^2$, we chose $X=0$ with probability 1, achieving $R(D)=0$.
- We can rewrite the distortion in terms of the rate

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

$$\frac{\sigma^2}{D} = 2^{2R}$$

$$D = \frac{\sigma^2}{2^{2R}}$$

So, the mutual information between X and \hat{X} is lower bounded by half log of sigma square by D . Now, we are going to show that this lower bound is achievable. So, if D is

less than σ^2 , we choose X to be \hat{X} plus Z and \hat{X} is Gaussian distributed with 0 mean and variance $\sigma^2 - D$ and Z is Gaussian distributed with 0 mean and variance D and X is Gaussian distributed with mean 0 and variance σ^2 . So, mutual information between X and \hat{X} is given by differential entropy of X minus differential entropy of X given \hat{X} . Now X is Gaussian distributed and X given \hat{X} is also Gaussian distributed and we can compute this mutual information which comes out to be half of $\log \sigma^2 / D$ and this is precisely the lower bound that we computed here.

So, this mutual information is actually achievable if we take our \hat{X} and Z in this particular rate and what happens if D is greater than σ^2 , in that case we choose X equal to 0 with probability 1 and that gives us rate as 0. Now, this rate can also be written in terms of distortion. So, I have σ^2 / D or $\log \sigma^2 / D$ is $2R$ or I can write σ^2 / D to be 2^{2R} or D is $\sigma^2 / 2^{2R}$ and I can write distortion in terms of rate and if rate increases, then distortion decreases and that makes sense and if you are using more number of bits to represent a quantity, basically our distortion is going to decrease.

(Refer Slide Time: 20:51)

Rate Distortion Function

- The rate distortion function for a Bernoulli(p) source with Hamming distortion is given by
$$R(D) = \begin{cases} H(p) - H(D) & 0 \leq D \leq \min\{p, 1-p\} \\ 0 & D > \min\{p, 1-p\} \end{cases}$$
- $X \sim \text{Bernoulli}(p)$. Without loss of generality, we consider $p < \frac{1}{2}$.
 $X \oplus \hat{X} = 1$ is equivalent to $X \neq \hat{X}$.

Now, let us take another example, this time we are considering a Bernoulli source and

hamming distortion measure. So, we can show that the rate distortion function for this Bernoulli source is given by this expression and as long as D is within 0 and minimum of p and $1 - p$, it is given by $H(p) - H(D)$; otherwise it is given by 0. Now, without loss of generality let's assume that this p is less than half and we are considering a Bernoulli source, this is a binary source, so X or \hat{X} . When X or \hat{X} is one, it means X is not same as \hat{X} . So, when X is from Bernoulli source, X or \hat{X} equal to 1 is equivalent of saying X is not same as \hat{X} .

(Refer Slide Time: 22:20)

Rate Distortion Function

- We first will find a lower bound on $I(X, \hat{X})$ and then show that this lower bound is achievable.

$$\begin{aligned}
 I(X; \hat{X}) &= H(X) - H(X|\hat{X}) \\
 &= H(p) - H(X \oplus \hat{X} | \hat{X}) \\
 &\geq H(p) - H(X \oplus \hat{X}) \\
 &\geq H(p) - H(D)
 \end{aligned}$$

since $\Pr(X \neq \hat{X}) \leq D$ and $H(D)$ increases with D for $D \leq \frac{1}{2}$.

- Thus $R(D) \leq H(p) - H(D)$

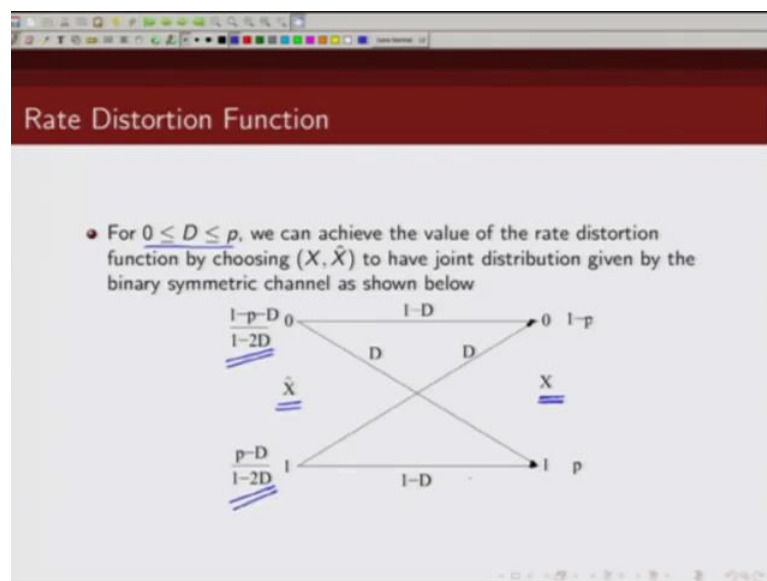
Now, again here to compute the rate distortion function, we will first put a lower bound on mutual information between X and \hat{X} and then we will show that this lower bound is achievable. So, let us compute this mutual information between X and \hat{X} . Now, from the definition of mutual information I I can write this as H of X minus H of X given \hat{X} . Now, X is a Bernoulli source with probability p . So, this is basically a binary source and its entropy is given by H of p . This is binary function of p .

Now, what is the uncertainty in X given \hat{X} ? So, there is uncertainty in X only when it is not same as \hat{X} . So, uncertainty in X given \hat{X} is basically uncertainty in X not equal to \hat{X} and that is given by this X or \hat{X} . So, we can write this mutual information between X and \hat{X} as H of p minus H of X or \hat{X} given \hat{X} . Now,

again this is condition on \hat{X} and if you remove this conditioning, we know that conditioning cannot increase entropy. So, this quantity is larger H of X or \hat{X} , this is more than X . So, we are subtracting a larger quantity and that is why I have here mutual information between X and \hat{X} is greater than equal to H of p minus H of X or \hat{X} and X or \hat{X} is basically when there is error and this probability of error is upper bounded by D . So, this is upper bounded and I can write this as greater than equal to H of p minus H of D .

So, the rate distortion function is given by this for the case when D is within minimum of p or 1 minus p . Now, let us consider the case when D lies between 0 to p and remember our objective is to cook up a distribution such that this rate is achievable with equality. So, X is distributed as Bernoulli with p and we want H of X given \hat{X} to be basically H of D . So, we have to cook up a distribution.

(Refer Slide Time: 25:54)



We create a test channel. This is my X and this is my \hat{X} . Now, I want to create a joint distribution such that my X is distributed as Bernoulli we know with 0 probability 1 minus p 1 with probability p and I need to find out this distribution on \hat{X} keeping in mind we want to keep the distortion less than equal to D .

(Refer Slide Time: 26:36)

Rate Distortion Function

- We choose the distribution of \hat{X} at the input of the channel so that the output distribution of X is the specified distribution. Let $r = \Pr(\hat{X} = 1)$. We choose r such that $r(1 - D) + (1 - r)D = p$ or $r = \frac{p-D}{1-2D}$.
- If $D \leq p \leq \frac{1}{2}$, then $\Pr(\hat{X} = 1) \geq 0$ and $\Pr(\hat{X} = 0) \geq 0$. We have

$$I(X; \hat{X}) = H(X) - H(X|\hat{X}) = H(p) - H(D).$$
 and expected distortion is $P(X \neq \hat{X}) = D$.
- If $D \geq p$, we achieve $R(D) = 0$ by letting $\hat{X} = 0$ with probability 1. In this case $I(X; \hat{X}) = 0$ and $D = p$.
- Similarly, if $D \geq 1 - p$, we can achieve $R(D) = 0$ by setting $\hat{X} = 1$ with probability 1.

Now, we choose the distribution of \hat{X} at the input of this test channel such that the output distribution is the specified distribution. So, let r equal to probability of \hat{X} being 1. So, we need to choose r in such a way such that r times $1 - D$ plus $1 - r$ times D that is equal to p and r times if this is the probability of \hat{X} being 1. So, r times $1 - D$ plus $1 - r$ times D that has to be p . So, if I do that, I get this input distribution on \hat{X} . So, probability of \hat{X} being 1, I get this as $p - D$ divided by $1 - 2D$. Now, if p lies between D and half, then probability of \hat{X} being 1 is greater than equal to 0 and probability of \hat{X} being 0 is greater than equal to 0 and then in this case the mutual information in X and \hat{X} is equal to H of p minus H of D where the expected distortion is given by D . So, if we take this input distribution on \hat{X} , we are able to achieve this lower bound. So, the rate distortion function is given by H of p minus H of D .

Now, for the case when D is greater than equal to p , $R(D)$ is 0 and we let \hat{X} is equal to 0 with probability 1 and in this case mutual information between X and \hat{X} comes out to be 0 and distortion is basically p and similarly for D greater than equal to $1 - p$, we can achieve $R(D)$ equal to 0 by setting \hat{X} is equal to 1 with probability 1. In this case also the mutual information between X and \hat{X} will come out to be 0.

(Refer Slide Time: 29:15)

Rate Distortion Function

- Consider a source X uniformly distributed on the set $\{1, 2, \dots, m\}$. Find the rate distortion function for this source with Hamming distortion, i.e.

$$d(x, \hat{x}) = \begin{cases} 0 & \text{if } x = \hat{x} \\ 1 & \text{if } x \neq \hat{x} \end{cases}$$
- Consider any joint distribution that satisfies the distortion constraint D .
- Since $D = \Pr\{X \neq \hat{X}\}$, we have by Fano's inequality

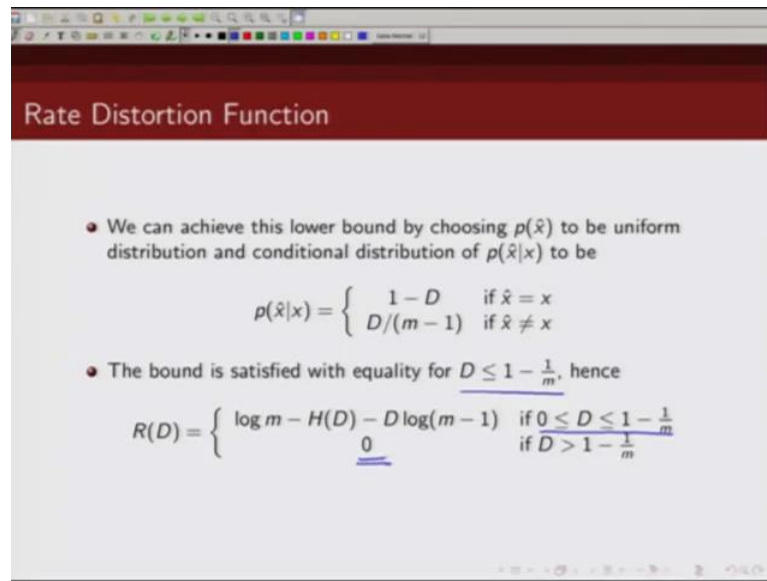
$$H(X|\hat{X}) \leq H(D) + D \log(m-1)$$
- Thus, we have

$$I(X; \hat{X}) = \frac{H(X) - H(X|\hat{X})}{\log m - H(D) - D \log(m-1)}$$

So, we have a source which is uniformly distributed in this set 1 2 3 and m and let us compute its rate distortion function for hamming distortion. So, if X is equal to \hat{X} , this is 0 and if X is not equal to \hat{X} , distortion is 1.

Now, consider any joint distribution that satisfies this distortion constraint of D and since D is probability of X not been equal to \hat{X} , we can invoke Fano's lemma and if we invoke Fano's lemma, we can show the uncertainty in X given \hat{X} is less than equal to $H(D) + D \log(m-1)$, this is $p \log$ of number of possibilities minus 1. So, that is our Fano's lemma. So, from Fano's lemma, we get this. Next, we compute the mutual information between X and \hat{X} and this is given by entropy of X minus entropy of X given \hat{X} . Now, from the Fano's lemma, we have an upper bound on $H(X|\hat{X})$. So, if we subtract this upper bound, since we are subtracting a larger quantity, so we now have a lower bound on mutual information between X and \hat{X} . So, mutual information between X and \hat{X} is lower bounded by $\log m - H(D) - D \log(m-1)$ and of course and this is maximum when X is uniformly distributed. So, $H(X)$ will be $\log m$.

(Refer Slide Time: 31:53)



Rate Distortion Function

- We can achieve this lower bound by choosing $p(\hat{x})$ to be uniform distribution and conditional distribution of $p(\hat{x}|x)$ to be
$$p(\hat{x}|x) = \begin{cases} 1 - D & \text{if } \hat{x} = x \\ D/(m-1) & \text{if } \hat{x} \neq x \end{cases}$$
- The bound is satisfied with equality for $D \leq 1 - \frac{1}{m}$, hence
$$R(D) = \begin{cases} \log m - H(D) - D \log(m-1) & \text{if } 0 \leq D \leq 1 - \frac{1}{m} \\ 0 & \text{if } D > 1 - \frac{1}{m} \end{cases}$$

Now, we can achieve this lower bound by choosing our p of X hat to be uniformly distributed and if we choose our conditional distribution in this way that if X is equal to X hat, probability of X hat given X is 1 minus D , otherwise it is D by m minus 1 . Now, this bound is satisfied with equality when D is less than equal to 1 minus 1 by m . So, the rate distortion function for this uniform source under hamming distortion is given by \log of m minus H of D minus $D \log$ of m minus 1 as long as D is less than equal to 1 minus 1 by m and for D greater than that, this is given by 0 .

(Refer Slide Time: 32:59)

Rate Distortion Function

- Information rate distortion function is defined as

$$R(D) = \min_{q(\hat{x}|x): \sum_x p(x) \sum_{\hat{x}} q(\hat{x}|x) d(x, \hat{x}) \leq D} I(X; \hat{X})$$
 where the minimization is over all conditional distributions $q(\hat{x}|x)$ for which the joint distribution $p(x)q(\hat{x}|x)$ satisfies the expected distortion constraint.
- This is a minimization of a convex function over the convex set of all $q(\hat{x}|x) \geq 0$ satisfying $\sum_{\hat{x}} q(\hat{x}|x) = 1$ for all x and $\sum_x \sum_{\hat{x}} q(\hat{x}|x) p(x) d(x, \hat{x}) \leq D$
- Using Lagrange multiplier method, we setup the functional

$$J(q) = \sum_x \sum_{\hat{x}} p(x) q(\hat{x}|x) \log \frac{q(\hat{x}|x)}{\sum_x p(x) q(\hat{x}|x)} + \lambda \sum_x \sum_{\hat{x}} p(x) q(\hat{x}|x) d(x, \hat{x}) + \sum_x \nu(x) \sum_{\hat{x}} q(\hat{x}|x)$$

Now, let us characterize this rate distortion function little bit more. So, we know that the rate distortion function is given by this expression and we are minimizing their mutual information between X and \hat{X} and remember this minimization is over all conditional distribution q of \hat{X} given X for which this joint distribution $p(x)q(\hat{x}|x)$ satisfies the expected distortion constraint. So, this expected value of this distortion is less than equal to D .

Now, this is a minimization of a convex function over a convex set for all q of \hat{X} given X greater than equal to 0 satisfying this constraint that is sum of q of \hat{X} given X over all \hat{X} , that summation should be 1 and this average distortion constraint should be satisfied. So, again I repeat basically we are minimizing this convex function over this convex set and we have some additional constraint which is sum of probabilities is 1 and this expected distortion constraint should be satisfied. So, then we can set up our Lagrangian, basically using Lagrange multiplier method, we can setup our functional. So, this is our objective function, this is the constraint due to this expected distortion constraint and this is the constraint that comes due to summation of probabilities be equal to 1.

(Refer Slide Time: 35:19)

Rate Distortion Function

- If we let $q(\hat{x}) = \sum_x p(x)q(\hat{x}|x)$ be the distribution on \hat{X} induced by $q(\hat{x}|x)$, we can write $J(q)$ as

$$J(q) = \sum_x \sum_{\hat{x}} p(x)q(\hat{x}|x) \log \frac{q(\hat{x}|x)}{q(\hat{x})} + \lambda \sum_x \sum_{\hat{x}} p(x)q(\hat{x}|x)d(x, \hat{x}) + \sum_x \nu(x) \sum_{\hat{x}} q(\hat{x}|x)$$
- Differentiating with respect to $q(\hat{x}|x)$, we have

$$\frac{\partial J}{\partial q(\hat{x}|x)} = p(x) \log \frac{q(\hat{x}|x)}{q(\hat{x})} + p(x) - \sum_{x'} p(x')q(\hat{x}|x') \frac{1}{q(\hat{x})} p(x) + \lambda p(x)d(x, \hat{x}) + \nu(x) = 0$$
- Setting $\log \mu(x) = \nu(x)/p(x)$, we obtain

$$p(x) \left[\log \frac{q(\hat{x}|x)}{q(\hat{x})} + \lambda d(x, \hat{x}) + \log \mu(x) \right] = 0$$

Now, let q of X hat be the distribution on X hat induced by this conditional distribution, then this functional can be written like this, if we differentiate this with respect to q of X hat given X and equate it to 0, we get this. So, p of X log of q of X hat given X by q of X hat plus p of X minus summation of p x prime q X hat given X prime into 1 by q X hat into p x plus lambda times p x and this distortion between X and X hat plus μ of X is equal to 0. Now, setting log of μ X to be ν x by p x, from here we get this condition. So, we get p x times log of q of X hat given X by q of X hat plus lambda times distortion measure between X and X hat plus log of μ X that is 0. Now, this is non-zero. So, what we can get is this is 0 from here we can write q of X hat given X and that comes out to be this.

(Refer Slide Time: 37:06)

Rate Distortion Function

- Hence

$$q(\hat{x}|x) = \frac{q(\hat{x})e^{-\lambda d(x,\hat{x})}}{\mu(x)}$$
- Since $\sum_{\hat{x}} q(\hat{x}|x) = 1$, we have

$$\mu(x) = \sum_{\hat{x}} q(\hat{x})e^{-\lambda d(x,\hat{x})}$$
- or

$$q(\hat{x}|x) = \frac{q(\hat{x})e^{-\lambda d(x,\hat{x})}}{\sum_{\hat{x}} q(\hat{x})e^{-\lambda d(x,\hat{x})}}$$
- Multiplying both sides by $p(x)$ and summing over all x , we get

$$q(\hat{x}) = q(\hat{x}) \sum_x \frac{p(x)e^{-\lambda d(x,\hat{x})}}{\sum_{\hat{x}} q(\hat{x})e^{-\lambda d(x,\hat{x})}}$$

Now, we do not know λ , but we know that summation of $q(\hat{x}|x)$ is equal to 1. So, $\mu(x)$ then comes out to be this and $q(\hat{x}|x)$ is this quantity. Next, multiplying both side $p(x)$ and summing over all x , we get $q(\hat{x})$ is equal to $q(\hat{x})$ of x hat into this. Now, if $q(\hat{x})$ is greater than 0, then we get this condition that this summation is equal to 1.

(Refer Slide Time: 38:11)

Rate Distortion Function

- If $q(\hat{x}) > 0$, we obtain

$$\sum_x \frac{p(x)e^{-\lambda d(x,\hat{x})}}{\sum_{\hat{x}} q(\hat{x})e^{-\lambda d(x,\hat{x})}} = 1$$
- Applying Kuhn Tucker conditions, we have

$$\frac{\partial J}{\partial q(\hat{x}|x)} \begin{cases} = 0 & \text{if } q(\hat{x}|x) > 0 \\ \geq 0 & \text{if } q(\hat{x}|x) = 0 \end{cases}$$
- Substituting the value of the derivative we get

$$\begin{aligned} \sum_x \frac{p(x)e^{-\lambda d(x,\hat{x})}}{\sum_{\hat{x}} q(\hat{x})e^{-\lambda d(x,\hat{x})}} &= 1 \quad \text{if } q(\hat{x}) > 0 \\ \sum_x \frac{p(x)e^{-\lambda d(x,\hat{x})}}{\sum_{\hat{x}} q(\hat{x})e^{-\lambda d(x,\hat{x})}} &\leq 1 \quad \text{if } q(\hat{x}) = 0 \end{aligned}$$

So, if this is greater than 0, then you can cancel this out and this will be equal to 1. Now, applying Kuhn Tucker condition, a partial derivative of J with respect to this q of X hat given X is 0, if q of X given X hat is greater than 0 and this is greater than equal to 0 if q of X hat given X is 0. Now, substituting this we get this condition. Let us consider a case where we have m independent Gaussian random variables.

(Refer Slide Time: 39:10)

Simultaneous description of independent Gaussian RVs.

- Consider m independent normal random sources X_1, \dots, X_m , where X_i are $\sim N(0, \sigma_i^2)$ with squared error distortion.
- Assume that we are given R bits with which to represent this random vector. How should the bits be allotted to various components to minimize the total distortion.
- We have

$$R(D) = \min_{f(\hat{x}^m | x^m): E[d(X^m, \hat{X}^m)] \leq D} I(X^m; \hat{X}^m)$$

where $d(x^m, \hat{x}^m) = \sum_{i=1}^m (x_i - \hat{x}_i)^2$

So, we have m independent Gaussian random variables and these are denoted by X_1, X_2, \dots, X_m . So, these are 0 mean Gaussian with variance σ_i^2 . They are not identically distributed, but they are independent and we are considering squared error distortion measure. So, assume that we are given R bits and with this R bits, we have to represent this random vector X_1, X_2, \dots, X_m . Now, the question that we are asking is how many bits we should be allocating to each of these components? So, how many bits we should be allocating to X_1, X_2, \dots, X_m such that we minimize the total distortion. So, again I repeat, we have m independent Gaussian sources which we are denoting by X_1, X_2, \dots, X_m , these are independent but not identically distributed and we are considering squared error distortion measure. So, given there are R bits, the question that we are asking is how many bits should be used to represent X_1, X_2, \dots, X_m such that overall distortion is minimized.

Now, from the definition of rate distortion function, we know that rate distortion function can be written as minimum mutual information between the source sequence X of m and this reproduce signal \hat{X} of m and this minimization is over conditional distribution such that the expected distortion constraint is satisfied.

(Refer Slide Time: 41:46)

Simultaneous description of independent Gaussian RVs.

• We have

$$\begin{aligned}
 I(X^m; \hat{X}^m) &= h(X^m) - h(X^m | \hat{X}^m) \\
 &= \sum_{i=1}^m h(X_i) - \sum_{i=1}^m h(X_i | X^{i-1}, \hat{X}^m) \\
 &\geq \sum_{i=1}^m h(X_i) - \sum_{i=1}^m h(X_i | \hat{X}_i) = \sum_{i=1}^m (h(X_i) - h(X_i | \hat{X}_i)) \\
 &= \sum_{i=1}^m I(X_i; \hat{X}_i) \\
 &\geq \sum_{i=1}^m R(D_i) \quad R(D_i) = \begin{cases} \frac{1}{2} \log \frac{\sigma_i^2}{D_i} & D_i < \sigma_i^2 \\ 0 & \text{otherwise} \end{cases} \\
 &= \sum_{i=1}^m \left(\frac{1}{2} \log \frac{\sigma_i^2}{D_i} \right)^+ \quad (x^+)
 \end{aligned}$$

where $D_i = E(X_i - \hat{X}_i)^2$

Now, let us look at mutual information between X^m and \hat{X}^m and from the definition of mutual information, we can write mutual information as differential entropy of X^m minus differential entropy of X^m given \hat{X}^m , now X_1, X_2, X_3 are independent right. So, we can write the differential entropy h of X^m as h of X_1 plus h of X_2 plus h of X_3 up to h of X_m . So, this term can be written like this and similarly I can write this h of X^m which is basically X_1, X_2, X_3, X_m given \hat{X}^m I can write it in this particular fashion.

Now, this is conditioned X_i is conditioned on this sequence X^{i-1} and X^m . So, we know that conditioning cannot increase entropy. So, if I just condition it only on \hat{X}_i , then this is a larger quantity, so I can lower bound this mutual information and this can be written as summation i equal to 1 to m h of X_i minus h of X_i given \hat{X}_i . Now, from the definition of mutual information, this is mutual information between X_i and \hat{X}_i . So, then mutual information between this sequence X^m and \hat{X}^m is lower

bounded by mutual information between X_i and \hat{X}_i sum over i going from 1 to m and since we know rate distortion function is minimum of this, so this can be lower bounded by $R(D)$ is summation over i going from 1 to m and we know for the Gaussian source, we know the expression for rate distortion function and if you go back for a Gaussian source this given by half log of sigma square by D . So, we can write this as when D_i is less than sigma i square else is 0 otherwise if I am just writing this as this operator X plus.

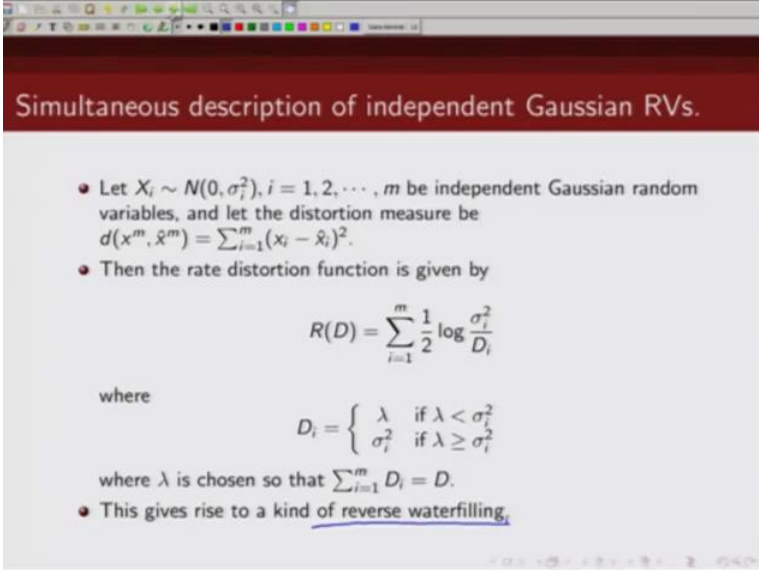
(Refer Slide Time: 45:35)

Simultaneous description of independent Gaussian RVs.

- Hence the problem of finding the rate distortion function can be reduced to
$$R(D) = \min_{D_i=D} \sum_{i=1}^m \max \left\{ \frac{1}{2} \log \frac{\sigma_i^2}{D_i}, 0 \right\}$$
- Using Lagrange multipliers, we construct the functional as
$$J(D) = \sum_{i=1}^m \frac{1}{2} \log \frac{\sigma_i^2}{D_i} + \lambda \sum_{i=1}^m D_i$$
- Differentiating with respect to D_i , and equating to 0, we get
$$\frac{\partial J}{\partial D_i} = -\frac{1}{2} \frac{1}{D_i} + \lambda = 0 \implies D_i = \lambda'$$

So, this is equal to X when this condition is satisfied, otherwise this is 0 and then the problem of finding the rate distortion function has been reduced to maximizing half log of sigma i square by D_i or 0 for i going from 1 to m and minimizing this over D_i , summation of D_i basically is D . Now, using Lagrange multiplier, this is our objective function and this is our constraint, the summation of this D_i from $i=1$ to m cannot exceed D . So, differentiating with respect to D_i , we get $\partial J / \partial D_i$ as this and when we equate it to 0, we get this condition. So, this is interesting, it says that distortion D_i should be same, D_i is equal to some constant λ' .

(Refer Slide Time: 47:01)



Simultaneous description of independent Gaussian RVs.

- Let $X_i \sim N(0, \sigma_i^2)$, $i = 1, 2, \dots, m$ be independent Gaussian random variables, and let the distortion measure be $d(x^m, \hat{x}^m) = \sum_{i=1}^m (x_i - \hat{x}_i)^2$.
- Then the rate distortion function is given by

$$R(D) = \sum_{i=1}^m \frac{1}{2} \log \frac{\sigma_i^2}{D_i}$$

where

$$D_i = \begin{cases} \lambda & \text{if } \lambda < \sigma_i^2 \\ \sigma_i^2 & \text{if } \lambda \geq \sigma_i^2 \end{cases}$$

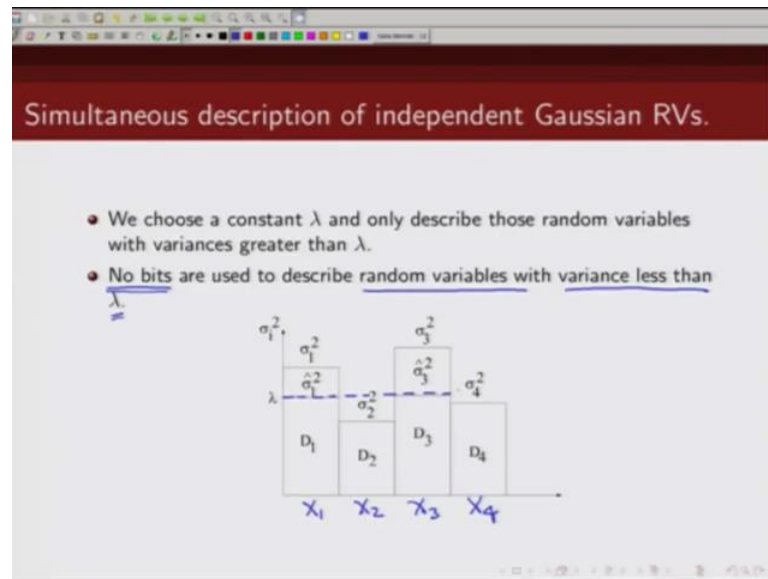
where λ is chosen so that $\sum_{i=1}^m D_i = D$.

- This gives rise to a kind of reverse waterfilling.

In fact, if X_i is Gaussian distributed with 0 mean and variance σ_i^2 and if we consider m such independent sources given a square distortion measure, the rate distortion function is given by this expression where D_i 's are equal to some constant λ as long as λ is less than σ_i^2 , otherwise if λ is greater than σ_i^2 D_i is equal to σ_i^2 and we choose our λ in such a way that sum of this distortion over these m sources that is equal to D .

So, this gives rise to what we call reverse water filling. In case of this panel Gaussian channel, we saw how we are allocating power, it was like putting water when you have various noise levels and here it is lateral difference. So, we choose a constant λ and we only describe those random variables which have variance greater than λ .

(Refer Slide Time: 48:34)



So, we do not use any bits to describe those random variables whose variance is less than λ . Let us just look at let us say this is $x_1 x_2 x_3 x_4$ and these are the distortions $D_1 D_2 D_3 D_4$. So, we fix our λ , this is my λ . If λ is less than σ_i^2 , D_i is λ and in this case λ is less than σ_1^2 . So, D_1 is λ . Similarly, D_3 is λ however, for D_2 and D_4 ; λ is greater than σ_2^2 and σ_4^2 . So, in this case distortion will be given by σ_2^2 and σ_4^2 and as I said this is interesting so, we are only going to describe this random variables whose variance is greater than λ and no bits are going to be used for those random variables whose variance is less than this threshold λ .

So, in the next class we are going to talk about how we can compute this rate distortion function in an iterative fashion. Now, this same algorithm can also be used to compute channel capacity and that will be the discussion for next class.

Thank you.