

MINIMIZING WEIGHTED CONCAVE IMPURITY PARTITION UNDER CONSTRAINTS

Thuan Nguyen and Think Nguyen

School of EECS, Oregon State University, Corvallis, OR 97331-5501, USA
 nguyeth9@oregonstate.edu, thinkq@eeecs.oregonstate.edu

ABSTRACT

Set partitioning is a key component of many algorithms in machine learning, signal processing and communications. In general, the problem of finding a partition that minimizes a given impurity (loss function) is NP-hard. As such, there exists a wealth of literature on approximate algorithms and theoretical analysis for the partitioning problem under different settings. In this paper, we formulate and solve a variant of the partition problem called the minimum weighted concave impurity partition under constraint (MIPUC). MIPUC finds an optimal partition that minimizes a given weighted concave loss function under a given concave constraint. MIPUC generalizes the recently proposed Deterministic Information Bottleneck problem which finds an optimal partition that maximizes the mutual information between the input and partitioned output while minimizing the partitioned output entropy. Our proposed algorithm is based on an optimality condition, which allows us to find a locally optimal solution efficiently. We also show that the optimal partitions are separated by some hyperplanes in the space of posterior probability mass functions.

Index Terms— Partition, quantization, impurity, constraint.

1. INTRODUCTION

Partitioning algorithms play a key role in machine learning, signal processing and communications. Given a set \mathbb{Y} consisting of M N -dimensional elements and a loss function over the subsets of \mathbb{Y} , a K -optimal partition algorithm splits \mathbb{Y} into K disjoint subsets such that the total loss over all K subsets is minimized. The loss function has also been termed the impurity which measures the average of a specified non-homogeneity property of the elements in each subset. Popular impurity functions include Gini index and the Shannon entropy [1–4]. From a communication and coding theory perspective, the problem of finding an optimal quantizer that maximizes the mutual information between the input and the quantized output is an important instance of the partition problem [5–14]. In this setting, the transmitted signal, modeled as the random variable X , is distorted by a channel, resulting in a received signal modeled as the random variable Y . A primary goal is to recover the transmitted signal X from the received signal Y accurately. To that end, one wants to design a quantizer, i.e., a mapping $Q(Y) \rightarrow Z$ such that Z and X share the most information. Since mutual information is the right metric for measuring the shared information between two random variables, designing $Q(Y)$ that maximizes the mutual information between Z and X is an important objective for many settings [15–18].

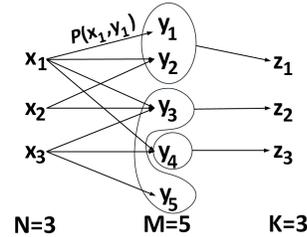


Fig. 1: $Q(Y) \rightarrow Z$ for a given joint distribution p_{x,y_i}

In this paper, we study a new partition problem called minimum weighted concave impurity partition under constraint (MIPUC) [19]. MIPUC is the minimum weighted concave impurity partition with an additional constraint on the distribution of the resulted partitions. MIPUC is motivated by many real-world problems. Specifically, our problem generalizes the recently proposed Deterministic Information Bottleneck problem (DIB) [15]. DIB finds an optimal partition that maximizes the mutual information between the input and the partitioned output while keeping the partitioned output entropy smaller than a certain threshold. For a given input distribution, maximizing mutual information is equivalent to minimizing entropy impurity [5], [6], thus, DIB can be viewed as a sub-problem of MIPUC. However, it is worth noting that the technique used in [15] is specific to an objective and constraint function, and is hard to extend to other impurity and constraint functions. In contrast, motivated by the approach in [20], our proposed algorithm is based on a optimality condition, which allows us to find a locally optimal solution efficiently for an arbitrary frequency weighted concave impurity function under an arbitrary concave constraint. In addition, by theoretically proving that the hard partition is optimal, we show that the optimal partitions are separated by hyperplane cuts in the space of the posterior distribution. Based on this optimality condition, we can find the true global optimal solution efficiently for small sized problems.

2. PROBLEM FORMULATION

It is convenient to use Fig. 1 to illustrate the proposed MIPUC problem in the context of quantizer design for communication systems. Let the transmitted signal be modeled as a discrete random variable X taking on N possible values x_1, x_2, \dots, x_N , with the probability mass vector $\mathbf{p} = [p_1, p_2, \dots, p_N]$. X is transmitted over a channel that distorts the signal, resulting in the received signal modeled as a random variable Y taking on discrete values y_1, y_2, \dots, y_M with the probability

mass vector $\mathbf{q} = [q_1, q_2, \dots, q_M]$. The communication channel is modeled using a given channel matrix $P \in \mathbb{R}^{N \times M}$ whose entries $P_{ij} = p(y_j|x_i)$ denotes the conditional probability that y_j is received given x_i is transmitted. From \mathbf{p} and P , each received signal y_i is specified by a joint distribution $\mathbf{p}_{\mathbf{x}, y_i} = [p(x_1, y_i), p(x_2, y_i), \dots, p(x_N, y_i)]$. Next, a quantizer Q (possibly stochastic) is used to quantize Y into K discrete values z_1, z_2, \dots, z_K , modeled as a discrete random variable Z with the probability mass vector $\mathbf{r} = [r_1, r_2, \dots, r_K]$.

For a given K , the quantization should be done to minimize the impurity function or the cost function $F(X; Z)$ between X and Z . On the other hand, often times, it is useful to be able to impose some constraints on the partitioned output Z . For example, if Z needs to be transmitted over a limited bandwidth channel, or needs to be stored in a small-capacity disk, then it makes sense to ensure the entropy $H(Z)$ less than a certain number of bits. In this case, the objective of a quantizer is not only to minimize the entropy impurity function $H(X|Z)$ but also to keep $H(Z)$ small as the setting in the Deterministic Information Bottleneck problem (DIB) [15]. To that end, we study a generalized DIB problem called minimum weighted concave impurity partition under constraint (MIPUC) which allows for a broader class of objective functions and constraints. Specifically, MIPUC finds an optimal partition/quantizer that minimizes a given weighted concave impurity function under a given concave constraint. To formulate the problem precisely, we use the following notations.

- $\mathbf{p}_{\mathbf{x}, y_i} = [p(x_1, y_i), p(x_2, y_i), \dots, p(x_N, y_i)]$ denotes the N -dimensional joint pmf of X and $Y = y_i$.
- $\mathbf{p}_{\mathbf{x}, z_i} = [p(x_1, z_i), p(x_2, z_i), \dots, p(x_N, z_i)]$ denotes the N -dimensional joint pmf of X and $Z = z_i$.
- $\mathbf{p}_{\mathbf{x}|y_i} = [p(x_1|y_i), p(x_2|y_i), \dots, p(x_N|y_i)]$ denotes the N -dimensional conditional pmf of X given $Y = y_i$.
- $\mathbf{p}_{\mathbf{x}|z_i} = [p(x_1|z_i), p(x_2|z_i), \dots, p(x_N|z_i)]$ denotes the N -dimensional conditional pmf of X given $Z = z_i$.
- $\mathbf{p}_{\mathbf{x}, \mathbf{z}}$ denotes the $N \times K$ matrix representing the joint pmf of X and Z .

Given \mathbf{p} , \mathbf{q} , $\mathbf{p}_{\mathbf{x}, y_i}$, and a quantizer $Q(Y) \rightarrow Z$, Y is partitioned into K distinct partitions corresponding to K distinct values z_i 's such that the weighted concave impurity over all partitions is minimized.

Specifically, the impurity function due to partition z_i is denoted by:

$$F(\mathbf{p}_{\mathbf{x}, z_i}) = r_i f(\mathbf{p}_{\mathbf{x}|z_i}) = \mathbf{p}_{\mathbf{x}, z_i}^T \mathbf{1} f\left(\frac{\mathbf{p}_{\mathbf{x}, z_i}}{\mathbf{p}_{\mathbf{x}, z_i}^T \mathbf{1}}\right), \quad (1)$$

where T denotes the tranpose operation and $f : \mathbb{R}^N \rightarrow \mathbb{R}$ is concave.

The weighted concave impurity function over all partitions is defined as:

$$F(\mathbf{p}_{\mathbf{x}, \mathbf{z}}) = \sum_{i=1}^K F(\mathbf{p}_{\mathbf{x}, z_i}) = \sum_{i=1}^K \mathbf{p}_{\mathbf{x}, z_i}^T \mathbf{1} f\left(\frac{\mathbf{p}_{\mathbf{x}, z_i}}{\mathbf{p}_{\mathbf{x}, z_i}^T \mathbf{1}}\right). \quad (2)$$

The definition of weighted concave impurity function was previously proposed in [20], [21], [22]. Many popular impurity functions such as entropy and Gini index [20], [21], [22] satisfy the weighted concave impurity property.

As mentioned earlier, for many real-world problems, it is often required to have a pre-specified constraint on the partitioned output Z [15], [23], [24]. Since Z is a random variable, we consider the concave constraint on \mathbf{r} , the pmf of Z , having the following form:

$$G(\mathbf{r}) = g_1(r_1) + g_2(r_2) + \dots + g_K(r_K) \leq D, \quad (3)$$

where $g_i(\cdot)$ is a concave function $\forall i = 1, 2, \dots, K$, and D is a given constant. Examples of some useful concave constraints include entropy and linear functions. For instance, entropy constraint is useful when Z acts as the intermediate representation of Y that needs to be transmitted over a limited bandwidth channel or stored in a small storage capacity disk [23], [24]. In this scenario, the entropy of Z is the theoretical maximum information/compression rate which adjusts the number of bits representing Z to fit the channel bandwidth or the storage capability. A smaller $H(\mathbf{r})$ theoretically implies a smaller number of bits to represent Z [25].

The MIPUC can be formulated as finding an optimal quantizer Q^* by finding an induced optimal $\mathbf{p}_{\mathbf{x}, \mathbf{z}}^*$ via the following unconstrained problem:

$$\mathbf{p}_{\mathbf{x}, \mathbf{z}}^* = \arg \min_{\mathbf{p}_{\mathbf{x}, \mathbf{z}}} \beta F(\mathbf{p}_{\mathbf{x}, \mathbf{z}}) + G(\mathbf{r}), \quad (4)$$

where β is a pre-specified parameter to control a given trade-off between minimizing the impurity $F(\mathbf{p}_{\mathbf{x}, \mathbf{z}})$ or minimizing the constraint function $G(\mathbf{r})$.

3. SOLUTION APPROACH

Our approach to solving the MIPUC problem is based on the method in [20]. First, a necessary condition for an optimal $\mathbf{p}_{\mathbf{x}, \mathbf{z}}^*$, specifically the joint pmf of X with each z_i , i.e., $\mathbf{p}_{\mathbf{x}, z_i}^*$, for $i = 1, 2, \dots, K$, is characterized. Then, based on this condition, we describe an algorithm that finds a locally optimal solution. In addition, we also propose a method that can find the globally optimal solution in some special cases.

3.1. Optimality Condition

Let Q be a deterministic quantizer that assigns y_j to z_k deterministically, i.e., for any j , if $Q(y_j) = z_k$ then $p(z_k|y_j) = 1$ and $p(z_l|y_j) = 0$ for $k \neq l$. Let us consider two partitions z_u, z_v, \dots and a sample y_m such that $Q(y_m) = z_u$. Let $Q^{(u, v, m, t)}$ be a stochastic quantizer, a slightly perturbed version of the quantizer Q . $Q^{(u, v, m, t)}$ is constructed as follows. $Q^{(u, v, m, t)}$ is the quantizer that assigns y_m to z_u with probability $1 - t$ and to z_v with probability t . For any other $y_j, j \neq m$, $Q^{(u, v, m, t)}(y_j) = Q(y_j)$. Hence, if $t = 0$, then $Q^{(u, v, m, t)} = Q$. If $t = 1$, then $Q^{(u, v, m, t)}$ quantizes y_m to z_v instead of z_u , and for every $j \neq m$, $Q^{(u, v, m, t)}(y_j) = Q(y_j)$.

Let $\mathbf{p}'_{\mathbf{x}, \mathbf{z}}(t)$ be the new joint pmf of X and Z induced by the quantizer $Q^{(u, v, m, t)}$, then:

$$\begin{aligned} \mathbf{p}'_{\mathbf{x}, z_u}(t) &= \left(\sum_{\substack{y_i: Q(y_i)=z_u \\ i \neq m}} \mathbf{p}_{\mathbf{x}, y_i} \right) + (1-t)\mathbf{p}_{\mathbf{x}, y_m} \\ &= \left(\sum_{y_i: Q(y_i)=z_u} \mathbf{p}_{\mathbf{x}, y_i} \right) - t\mathbf{p}_{\mathbf{x}, y_m} = \mathbf{p}_{\mathbf{x}, z_u} - t\mathbf{p}_{\mathbf{x}, y_m}. \end{aligned} \quad (5)$$

Similarly, we have:

$$\mathbf{p}'_{\mathbf{x},z_v}(t) = \mathbf{p}_{\mathbf{x},z_v} + t\mathbf{p}_{\mathbf{x},y_m}. \quad (6)$$

$$\mathbf{p}'_{\mathbf{x},z_k}(t) = \mathbf{p}_{\mathbf{x},z_k}, \forall k \neq u, v. \quad (7)$$

The new weighted impurity produced by $Q^{(u,v,m,t)}$ is therefore:

$$F(\mathbf{p}'_{\mathbf{x},z}(t)) = \sum_{i=1}^K F(\mathbf{p}'_{\mathbf{x},z_i}(t)) = \sum_{i \neq u,v} F(\mathbf{p}_{\mathbf{x},z_i}) + F(\mathbf{p}_{\mathbf{x},z_v} + t\mathbf{p}_{\mathbf{x},y_m}) + F(\mathbf{p}_{\mathbf{x},z_u} - t\mathbf{p}_{\mathbf{x},y_m}). \quad (8)$$

Similarly, let $\mathbf{r}'(t) = [r'_1(t), r'_2(t), \dots, r'_K(t)]$ be the new pmf of Z induced by $Q^{(u,v,m,t)}$, by marginalizing (5) and (6) over X , we have:

$$\begin{aligned} r'_u(t) &= r_u - tq_m, \\ r'_v(t) &= r_v + tq_m, \\ r'_l(t) &= r_l, \forall l \neq u, v. \end{aligned}$$

Consequently, the new constraint function for $Q^{(u,v,m,t)}$ can be written as:

$$G(\mathbf{r}'(t)) = \left(\sum_{i \neq u,v} g_i(r_i) \right) + g_u(r_u - tq_m) + g_v(r_v + tq_m). \quad (9)$$

Define:

$$\begin{aligned} C^{(u,v,m)}(t) &= \beta \left(F(\mathbf{p}'_{\mathbf{x},z_v}(t)) + F(\mathbf{p}'_{\mathbf{x},z_u}(t)) \right) + g_v(r'_v(t)) + g_u(r'_u(t)) \\ &= \beta \left(F(\mathbf{p}_{\mathbf{x},z_v} + t\mathbf{p}_{\mathbf{x},y_m}) + F(\mathbf{p}_{\mathbf{x},z_u} - t\mathbf{p}_{\mathbf{x},y_m}) \right) \\ &\quad + g_v(r_v + tq_m) + g_u(r_u - tq_m). \end{aligned} \quad (10)$$

$C^{(u,v,m)}(t)$ represents the contributions to the objective function from two partitions z_u and z_v , which is a function of t as $Q^{(u,v,m,t)}$ quantizes y_m to z_u with probability $1 - t$ and to z_v with probability t . As t increases from 0 to 1, there is an increasing chance that y_m will be moved from z_u to z_v . If $C^{(u,v,m)}(0) > C^{(u,v,m)}(1)$, then the objective using the quantizer Q is larger than that of using $Q^{(u,v,m,1)}$ which moves y_m to the partition z_v completely. Similarly, if $C^{(u,v,m)}(0) < C^{(u,v,m)}(1)$, then Q is better than $Q^{(u,v,m,1)}$. We have the following proposition about $C^{(u,v,m)}(t)$.

Proposition 1. For any u, v , $Q(y_m) = z_u$, and $0 \leq t \leq a \leq 1$, we have:

$$C^{(u,v,m)}(t) \geq (1 - \frac{t}{a})C^{(u,v,m)}(0) + \frac{t}{a}C^{(u,v,m)}(a). \quad (11)$$

Equivalently,

$$\frac{C^{(u,v,m)}(t) - C^{(u,v,m)}(0)}{t} \geq \frac{C^{(u,v,m)}(a) - C^{(u,v,m)}(0)}{a}. \quad (12)$$

Proof. (outline) From the concavity of $f(\cdot)$, $g_u(\cdot)$ and $g_v(\cdot)$, one can verify that $C^{(u,v,m)}(t)$ as defined in (10), is concave. Thus, (11) follows. \square

Now, we are ready to prove the main result which characterizes the condition for an optimal quantizer Q^* .

Theorem 1. (Necessary optimality condition) Let Q be a quantizer with an induced joint pmf $\mathbf{p}_{\mathbf{x},z}$. For each partition z_k , $k = 1, 2, \dots, K$, define:

$$\mathbf{c}_k = \frac{dF(\mathbf{p}_{\mathbf{x},z_k})}{d\mathbf{p}_{\mathbf{x},z_k}}, \quad b_k = \frac{dg_k(r_k)}{dr_k}. \quad (13)$$

Define the "distance" from a data point y_j to a partition z_k as:

$$d(y_j, z_k) = \beta \mathbf{c}_k^T \mathbf{p}_{\mathbf{x},y_j} + b_k q_j, \quad (14)$$

then an optimal quantizer Q^* that quantizes y_j to z_k must have $d(y_j, z_k) \leq d(y_j, z_l)$, $l \neq k$.

Proof. Our proof follows the method in [20] with the difference being the incorporation of the concave constraint. Let Q^* be an optimal quantizer with an optimal induced joint pmf of X and Z , $\mathbf{p}_{\mathbf{x},z}^*$, a marginal pmf of Z , \mathbf{r}^* , and the distance $d^*(y_j, z_k)$. For any arbitrary partition z_u , suppose that existing y_m such that $Q^*(y_m) = z_u$, however, $d(y_m, z_u) > d(y_m, z_v)$ for some $v \neq u$. Let $Q^{(u,v,m,t)}$ be a perturbed version of Q^* . Now, from (10),

$$\begin{aligned} \frac{dC^{(u,v,m)}(t)}{dt} &= \beta \frac{dF(\mathbf{p}'_{\mathbf{x},z_u}(t))}{dt} + \beta \frac{dF(\mathbf{p}'_{\mathbf{x},z_v}(t))}{dt} \\ &\quad + \frac{dg_u(r'_u(t))}{dt} + \frac{dg_v(r'_v(t))}{dt}. \end{aligned} \quad (15)$$

Using the chain rule, the derivative of the first term on the right hand side above is:

$$\beta \frac{dF(\mathbf{p}'_{\mathbf{x},z_u}(t))}{dt} = \beta \left(\frac{dF(\mathbf{p}'_{\mathbf{x},z_u}(t))}{d\mathbf{p}'_{\mathbf{x},z_u}} \right)^T \frac{d\mathbf{p}'_{\mathbf{x},z_u}}{dt}.$$

Since $\mathbf{p}'_{\mathbf{x},z_u}(t) = \mathbf{p}_{\mathbf{x},z_u} - t\mathbf{p}_{\mathbf{x},y_m}$, we have:

$$\begin{aligned} \left. \frac{dF(\mathbf{p}'_{\mathbf{x},z_u}(t))}{d\mathbf{p}'_{\mathbf{x},z_u}} \right|_{t=0} &= \frac{dF(\mathbf{p}_{\mathbf{x},z_u})}{d\mathbf{p}_{\mathbf{x},z_u}} = \mathbf{c}_u, \\ \frac{d\mathbf{p}'_{\mathbf{x},z_u}}{dt} &= -\mathbf{p}_{\mathbf{x},y_m}. \end{aligned}$$

Thus,

$$\beta \left. \frac{dF(\mathbf{p}'_{\mathbf{x},z_u}(t))}{dt} \right|_{t=0} = -\beta \mathbf{c}_u^T \mathbf{p}_{\mathbf{x},y_m}. \quad (16)$$

Using the same chain rule for other terms in (15), we have:

$$\beta \left. \frac{dF(\mathbf{p}'_{\mathbf{x},z_v}(t))}{dt} \right|_{t=0} = \beta \mathbf{c}_v^T \mathbf{p}_{\mathbf{x},y_m}, \quad (17)$$

$$\left. \frac{dg_u(r'_u(t))}{dt} \right|_{t=0} = -\frac{dg_u(r_u)}{dr_u} q_m = -b_u q_m, \quad (18)$$

$$\left. \frac{dg_v(r'_v(t))}{dt} \right|_{t=0} = \frac{dg_v(r_v)}{dr_v} q_m = b_v q_m. \quad (19)$$

Summing (16), (17), (18), and (19), we have:

$$\begin{aligned} \left. \frac{dC^{(u,v,m)}(t)}{dt} \right|_{t=0} &= (\beta \mathbf{c}_v^T \mathbf{p}_{\mathbf{x},y_m} + b_v q_m) - (\beta \mathbf{c}_u^T \mathbf{p}_{\mathbf{x},y_m} + b_u q_m) \\ &= d(y_m, z_v) - d(y_m, z_u). \end{aligned} \quad (20)$$

Since we assume that $d(y_m, z_v) < d(y_m, z_u)$, we have:

$$\left. \frac{dC^{(u,v,m)}(t)}{dt} \right|_{t=0} < 0. \quad (21)$$

Now, from (12) in Proposition 1, using $a = 1$, and let $t \rightarrow 0$, we have:

$$\begin{aligned} \left. \frac{dC^{(u,v,m)}(t)}{dt} \right|_{t=0} &= \lim_{t \rightarrow 0} \frac{C^{(u,v,m)}(t) - C^{(u,v,m)}(0)}{t} \\ &\geq \frac{C^{(u,v,m)}(1) - C^{(u,v,m)}(0)}{1}. \end{aligned} \quad (22)$$

Since we assume that Q^* is optimal which corresponds to $C^{(u,v,m)}(0)$, then $C^{(u,v,m)}(0) \leq C^{(u,v,m)}(1)$. Combine with (22), we have:

$$\left. \frac{dC^{(u,v,m)}(t)}{dt} \right|_{t=0} \geq 0,$$

which contradicts (21). Hence, if $Q^*(y_m) = z_u$ then $d(y_m, z_u) \leq d(y_m, z_v)$, $\forall v \neq u$. \square

Proposition 2. (Deterministic quantizer is optimal) There exists an optimal deterministic quantizer.

Proof. Using (11) in Proposition 1, $C^{(u,v,m)}(t)$ is a concave function in $0 \leq t \leq 1$ by definition. Thus, the minimum of $C^{(u,v,m)}(t)$ must occur either at $t = 0$ or $t = 1$. Both $t = 0$ or $t = 1$ requires $Q(y_m) = z_u$ with probability 1, or $Q(y_m) = z_v$ with probability 1. Since the results hold for arbitrary z_u and z_v , there exists an optimal deterministic quantizer. \square

3.2. Algorithm

Based on the optimality condition in Theorem 1, we describe an iterative algorithm which is similar to a k -means algorithm for finding a locally optimal solution for any weighted concave objective and constraint function. In the initial step, the algorithm randomly selects a quantizer Q to assigns y_j to z_k . Next, based on the initial random clustering, $\mathbf{p}_{\mathbf{x},z}$, \mathbf{c}_k , r_k , b_k , and $d(y_j, z_k)$ are computed. Based on $d(y_j, z_k)$, the membership of y_j to each z_k is updated such that $Q(y_j) = z_k$ if $d(y_j, z_k)$ is the smallest over all z_k . The algorithm repeats until Q stops changing (membership of all partitions z_k 's do not change), or the maximum number of iterations has been reached. For a special case of entropy impurity and entropy constraint, the proposed algorithm is identical to the Algorithm 2 in [15]. The pseudo code for the proposed algorithm can be found in [19].

3.3. Hyperplane Separation of Optimal Partitions

Our result agrees with the well-known result in [21] for the problem of minimizing impurity without constraints which states that the optimal partitions are separated by hyperplane cuts in a $N - 1$ dimensional space of the posterior distributions $\mathbf{p}_{\mathbf{x}|y_j}$. Indeed, consider an optimal quantizer Q^* that induces conditional pmf $\mathbf{p}_{\mathbf{x}|z_k}$, $k = 1, 2, \dots, K$. Suppose that $Q^*(y_j) = z_k$ and $Q^*(y_i) = z_l$, $j \neq i$, $k \neq l$. From Theorem 1, we have:

$$\begin{aligned} d(y_j, z_k) &= \beta \mathbf{c}_k^T \mathbf{p}_{\mathbf{x},y_j} + b_k q_j = q_j (\beta \mathbf{c}_k^T \mathbf{p}_{\mathbf{x}|y_j} + b_k) \\ &\leq d(y_j, z_l) = \beta \mathbf{c}_l^T \mathbf{p}_{\mathbf{x},y_j} + b_l q_j = q_j (\beta \mathbf{c}_l^T \mathbf{p}_{\mathbf{x}|y_j} + b_l). \end{aligned}$$

Thus, $\beta \mathbf{c}_k^T \mathbf{p}_{\mathbf{x}|y_j} + b_k \leq \beta \mathbf{c}_l^T \mathbf{p}_{\mathbf{x}|y_j} + b_l$. Equivalently,

$$(\mathbf{c}_k^T - \mathbf{c}_l^T) \mathbf{p}_{\mathbf{x}|y_j} \leq \frac{1}{\beta} (b_l - b_k). \quad (23)$$

Using a similar derivation for $d(y_i, z_l) \leq d(y_i, z_k)$, we have:

$$(\mathbf{c}_k^T - \mathbf{c}_l^T) \mathbf{p}_{\mathbf{x}|y_i} \geq \frac{1}{\beta} (b_l - b_k). \quad (24)$$

From (23) and (24), we conclude that $\mathbf{p}_{\mathbf{x}|y_i}$ and $\mathbf{p}_{\mathbf{x}|y_j}$ are separated by a hyperplane with the orthogonal vector $(\mathbf{c}_k^T - \mathbf{c}_l^T)$ and offset $(b_l - b_k)/\beta$. In addition, because $\mathbf{p}_{\mathbf{x}|y_j}$ is a pmf, the sum of its components is 1, the separating hyperplanes lies in $N - 1$ dimensional space. An interesting result of this separation is that a naive exhaustive search over all the possible hyperplanes which requires the complexity of $O(M^{N-1})$ [21], can be practical.

4. APPLICATIONS AND NUMERICAL RESULTS

We now provide a small example to validate the theoretical results and the proposed algorithms. Consider an input source $X \in \{x_1 = -1, x_2 = 1\}$ having $\mathbf{p} = (0.2, 0.8)^T$ is transmitted over an AWGN channel with: $Y = X + N$, where $N \sim \mathcal{N}(\mu = 0, \sigma = 1)$. Consequently, $p_{y|x_1} = \mathcal{N}(1, 1)$ and $p_{y|x_2} = \mathcal{N}(-1, 1)$. We first quantize Y in the range $[-10, 10]$ into $M = 200$ discrete values $Y = [y_1, y_2, \dots, y_M]$ with equal spacing of $\epsilon = 0.1$. The discrete Y is then quantized into $Z \in \{z_1 = -1, z_2 = 1\}$ using a quantizer Q . The joint pmf $\mathbf{p}_{\mathbf{x},y}$ can be determined using two given conditional pmfs $p_{y|x_1}$ and $p_{y|x_2}$ above. To study the trade-off between maximizing $I(X; Z)$ and minimizing $H(Z)$, the proposed algorithm is used with multiple random starting points to determine the minimum values of $-\beta I(X; Z) + H(Z)$. Fig. 2 shows the Pareto curve for $I(X; Z)$ vs. $H(Z)$ at various values of β . If we want to design a quantizer having $H(Z) \leq 0.5$, we can choose $\beta^* = 6$ which produces $I(X; Z)^* = 0.18623$ and $H(Z)^* = 0.48873$.

In addition, since $p_{x_1|y} = \frac{p_1 p_{y|x_1}}{p_1 p_{y|x_1} + p_2 p_{y|x_2}}$ is a strictly decreasing function of y , it is possible to show that an exhaustive search over all the hyperplanes in a 1-dimensional space of the posterior distribution is equivalent to an exhaustive search in y [19]. Now, by exhaustive searching for $y \in [-10, 10]$ with the resolution $\epsilon = 0.1$ and using $\beta = 6$, the optimal mutual information $I(X; Z)^* = 0.18623$ and the optimal entropy $H(Z)^* = 0.48873$ are achieved at $y = -1.1$ which confirms the result of the proposed algorithm.

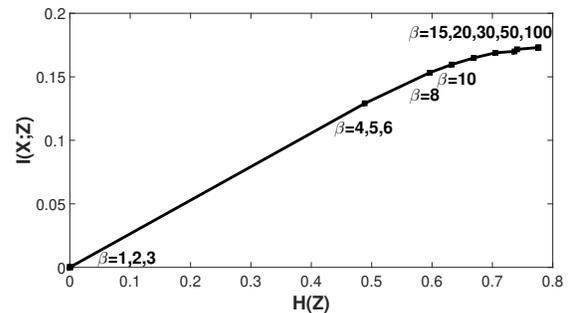


Fig. 2: $I(X; Z)$ vs. $H(Z)$ at various values of β .

5. CONCLUSION

In this paper, we introduced framework for determining the optimal partition that minimizes a given weighted concave impurity function under a given concave constraint on the partitioned outputs. Based on the optimality condition, we provide a low complexity algorithm to find the locally optimal solution. We also showed that there exists a deterministic optimal partition which corresponds to the regions separated by some hyperplane cuts in the probability space of the posterior distribution.

6. REFERENCES

- [1] J Ross Quinlan. *C4. 5: programs for machine learning*. Elsevier, 2014.
- [2] Leo Breiman. *Classification and regression trees*. Routledge, 2017.
- [3] Thuan Nguyen and Thinh Nguyen. A linear time partitioning algorithm for frequency weighted impurity functions. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 5375–5379. IEEE, 2020.
- [4] Eduardo S Laber, Marco Molinaro, and Felipe A Mello Pereira. Binary partitions with approximate minimum impurity. In *International Conference on Machine Learning*, pages 2860–2868, 2018.
- [5] Brian M Kurkoski and Hideki Yagi. Quantization of binary-input discrete memoryless channels. *IEEE Transactions on Information Theory*, 60(8):4544–4552, 2014.
- [6] Jiuyang Alan Zhang and Brian M Kurkoski. Low-complexity quantization of discrete memoryless channels. In *2016 International Symposium on Information Theory and Its Applications (ISITA)*, pages 448–452. IEEE, 2016.
- [7] T. Nguyen and T. Nguyen. Capacity achieving quantizer design for binary channels. *IEEE Communications Letters*, pages 1–1, 2020.
- [8] Thuan Nguyen and Thinh Nguyen. On binary quantizer for maximizing mutual information. *IEEE Transactions on Communications*, pages 1–1, 2020.
- [9] T. Nguyen and T. Nguyen. Communication-channel optimized impurity partition. In *GLOBECOM 2020 - 2020 IEEE Global Communications Conference*, pages 1–5, 2020.
- [10] Thuan Nguyen and Thinh Nguyen. Optimal quantizer structure for binary discrete input continuous output channels under an arbitrary quantized-output constraint. *International Symposium on Information Theory (ISIT)*, 2020.
- [11] Thuan Nguyen, Yu-Jung Chu, and Thinh Nguyen. A new fast algorithm for finding capacity of discrete memoryless thresholding channels. In *2020 International Conference on Computing, Networking and Communications (ICNC)*, pages 56–60. IEEE, 2020.
- [12] Brian M Kurkoski and Hideki Yagi. Single-bit quantization of binary-input, continuous-output channels. In *2017 IEEE International Symposium on Information Theory (ISIT)*, pages 2088–2092. IEEE, 2017.
- [13] T. Nguyen and T. Nguyen. Thresholding quantizer design for mutual information maximization under output constraint. In *2020 IEEE 91st Vehicular Technology Conference (VTC2020-Spring)*, pages 1–5, 2020.
- [14] T. Nguyen and T. Nguyen. On thresholding quantizer design for mutual information maximization: Optimal structures and algorithms. In *2020 IEEE 91st Vehicular Technology Conference (VTC2020-Spring)*, pages 1–5, 2020.
- [15] DJ Strouse and David J Schwab. The deterministic information bottleneck. *Neural computation*, 29(6):1611–1630, 2017.
- [16] Naftali Tishby, Fernando C Pereira, and William Bialek. The information bottleneck method. *arXiv preprint physics/0004057*, 2000.
- [17] Morteza Noshad, Yu Zeng, and Alfred O Hero. Scalable mutual information estimation using dependence graphs. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2962–2966. IEEE, 2019.
- [18] Ziv Goldfeld and Yury Polyanskiy. The information bottleneck problem and its applications in machine learning. *IEEE Journal on Selected Areas in Information Theory*, 2020.
- [19] Thuan Nguyen and Thinh Nguyen. Minimizing impurity partition under constraints. *Transaction on Communications*, 2020, submitted.
- [20] Don Coppersmith, Se June Hong, and Jonathan RM Hosking. Partitioning nominal attributes in decision trees. *Data Mining and Knowledge Discovery*, 3(2):197–217, 1999.
- [21] David Burshtein, Vincent Della Pietra, Dimitri Kanevsky, Arthur Nadas, et al. Minimum impurity partitions. *The Annals of Statistics*, 20(3):1637–1646, 1992.
- [22] Philip A. Chou. Optimal partitioning for classification and regression trees. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (4):340–354, 1991.
- [23] Thuan Nguyen and Thinh Nguyen. Entropy-constrained maximizing mutual information quantization. *arXiv preprint arXiv:2001.01830*, 2020.
- [24] Philip A. Chou, Tom D. Lookabaugh, and Robert M. Gray. Entropy-constrained vector quantization. *IEEE Trans. Acoustics, Speech, and Signal Processing*, 37:31–42, 1989.
- [25] Thomas M Cover and Joy A Thomas. *Elements of information theory*. John Wiley & Sons, 2012.