

2009

# Minimum cost content distribution using network coding: Replication vs. coding

Shurui Huang  
*Iowa State University*

Follow this and additional works at: <https://lib.dr.iastate.edu/etd>

 Part of the [Electrical and Computer Engineering Commons](#)

## Recommended Citation

Huang, Shurui, "Minimum cost content distribution using network coding: Replication vs. coding" (2009). *Graduate Theses and Dissertations*. 10885.

<https://lib.dr.iastate.edu/etd/10885>

This Thesis is brought to you for free and open access by the Iowa State University Capstones, Theses and Dissertations at Iowa State University Digital Repository. It has been accepted for inclusion in Graduate Theses and Dissertations by an authorized administrator of Iowa State University Digital Repository. For more information, please contact [digirep@iastate.edu](mailto:digirep@iastate.edu).

**Minimum cost content distribution using network coding: Replication vs. coding  
at the source nodes**

by

Shurui Huang

A thesis submitted to the graduate faculty  
in partial fulfillment of the requirements for the degree of  
MASTER OF SCIENCE

Major: Electrical Engineering

Program of Study Committee:  
Aditya Ramamoorthy, Major Professor  
Nicola Elia  
Huaiqing Wu

Iowa State University

Ames, Iowa

2009

Copyright © Shurui Huang, 2009. All rights reserved.

## DEDICATION

I would like to dedicate this thesis to my parents and to my boyfriend Hao Chen without whose support I would not have been able to complete this work. I would also like to thank my friends and family for their loving guidance and support during the writing of this work.

## TABLE OF CONTENTS

<b>LIST OF TABLES</b> . . . . .	v
<b>LIST OF FIGURES</b> . . . . .	vi
<b>ACKNOWLEDGEMENTS</b> . . . . .	vii
<b>ABSTRACT</b> . . . . .	viii
<b>CHAPTER 1. INTRODUCTION</b> . . . . .	1
1.1 Minimum Cost Multicast with Multiple Sources Problem . . . . .	3
1.2 Thesis Outline . . . . .	4
<b>CHAPTER 2. MEASURE THEORY AND INFORMATION THEORY</b> . .	6
2.1 Basic Concepts in Measure Theory . . . . .	6
2.2 One-to-one Correspondence between Shannon Information Measure and Set Theory . . . . .	7
2.3 Main Theorem . . . . .	8
2.3.1 An Example . . . . .	10
<b>CHAPTER 3. PROBLEM FORMULATION</b> . . . . .	12
3.1 Subset Constraints Source Network . . . . .	13
3.1.1 Basic Formulation . . . . .	13
3.1.2 Another Formulation . . . . .	17
3.1.3 Solution Explanation and Construction . . . . .	22
3.2 Coded Source Network . . . . .	23
<b>CHAPTER 4. COST COMPARISON BETWEEN CODED CASE AND SUBSET CASE</b> . . . . .	24

4.1	General Case . . . . .	25
4.1.1	Greedy Algorithm . . . . .	28
4.2	Three Sources Case . . . . .	29
<b>CHAPTER 5. RESULTS . . . . .</b>		<b>33</b>
5.1	Results on a Deterministic Network . . . . .	33
5.2	Results on Random Networks . . . . .	34
<b>CHAPTER 6. SUMMARY AND DISCUSSION . . . . .</b>		<b>36</b>
<b>APPENDIX . A VARIATION OF THE INCLUSION-EXCLUSION FOR-</b>		
	<b>MULA . . . . .</b>	<b>37</b>
<b>BIBLIOGRAPHY . . . . .</b>		<b>39</b>

**LIST OF TABLES**

Table 5.1	Atom values when subset constraints are enforced . . . . .	34
Table 5.2	Comparisons of two schemes in 5000 random directed graphs . . . . .	35

## LIST OF FIGURES

Figure 1.1	An example of the subset sources case . . . . .	2
Figure 3.1	(a) The network structure for the first formulation. (b) The network structure for the second formulation. . . . .	14
Figure 3.2	The auxiliary graph for three sources case . . . . .	21
Figure 4.1	Atom pattern when there are three sources . . . . .	29
Figure 4.2	The transforming scheme from coded case to subset case, $b^1$ is negative	31
Figure 5.1	A deterministic network . . . . .	33

## ACKNOWLEDGEMENTS

I would like to take this opportunity to express my thanks to those who helped me with various aspects of conducting research and the writing of this thesis. First and foremost, Dr. Aditya Ramamoorthy for his guidance, patience and support throughout this research and the writing of this thesis. His insights and words of encouragement have often inspired me and renewed my hopes for completing my graduate education. I would also like to thank my committee members for their efforts and contributions to this work: Dr. Nicola Elia and Dr. Huaiqing Wu.

## ABSTRACT

Large scale content distribution over the internet has been the focus of numerous studies in recent years. In the traditional server-client model, the server may suffer from overload when a popular file stored at the server is frequently requested. In order to reduce the cost at servers and decrease the retrieval time for clients, distributed storage solutions that operate by dividing the file into pieces and placing copies of the pieces (replication) or coded versions of the pieces (coding) at multiple source nodes have been proposed.

Network coding has also been used in large content distribution. In this work, we consider multicasting a file that can be broken into small pieces to multiple clients over a network with network coding. The network contains a set of source nodes that can store either subsets or coded version of the pieces of the file. We are interested in finding the optimal storage capacity and flows over the edges for the subset case and the coded case, respectively, such that the joint cost of transmission at edges and storage at sources is minimized. We provide succinct formulations of the corresponding optimization problems by using information measures. By the insight gained from the two formulations, a gap linear program which can compute the cost gap between the subset case and the coded case is formulated. A greedy algorithm is developed to find a suboptimal solution of the gap LP. In particular, we show that when there are two source nodes, there is no loss in considering subset sources. Furthermore, in the case of three source nodes, we derive a tight upper bound on the cost gap between the two cases. Algorithms for determining the content of the source nodes are also provided.

## CHAPTER 1. INTRODUCTION

Large scale content distribution over the Internet is a topic of great interest and has been the subject of numerous studies [1]. The dominant mode of content distribution is the client-server model, where a given client requests a central server for the file, which then proceeds to service the request. For example, this is how a website server operates traditionally. A single server, however is likely to be overwhelmed when a large number of users request for a file at the same time and the websites are often replicated by the use of mirrors [2]. One can also consider the usage of coding for replicating the content, e.g., if one uses erasure codes such as Reed-Solomon codes or fountain codes, then it turns out that obtaining a certain number of coded packets from each of mirrors will suffice. Peer-to-peer networks have also been proposed for content distribution in a distributed manner [3].

The technique of network coding has also been considered for content distribution in networks [1]. Network coding allows us to use the network resources more efficiently in the case of multicast, where a single source node or a group of source nodes contain information that is requested by a set of terminals. Under network coding based multicast, the problem of allocating resources such as rates and flows in the network can be solved in polynomial time. This is in contrast to multicast under routing, which is known to be a computationally hard problem. Moreover, one can arrive at distributed solutions to these problems in an easier manner under network coding.

In this work, we consider the following problem. Suppose that there is a large file, that can be broken into small pieces, that needs to be transmitted to a given set of clients over a network using network coding. The network has a designated set of nodes (called source nodes) that have storage space. Each unit of storage space and each unit of flow over a certain edge has

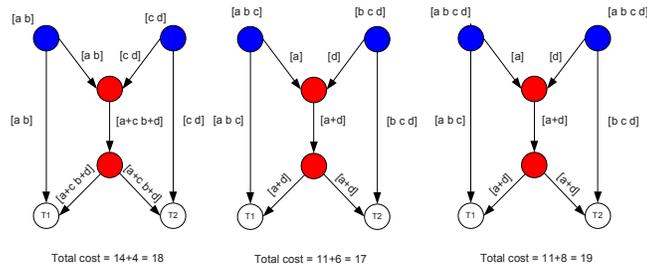


Figure 1.1 An example of the subset sources case

a known linear cost. We want to determine the optimal storage capacities and flow patterns over the network such that this can be done with minimum cost. Within this problem setting, we distinguish two different cases: (i) *Subset sources case*: Each source node only contains a subset of the pieces of the file. (ii) *Coded sources case*: Each source node can contain arbitrary functions of the pieces of the file.

At first glance, it may seem that in the subset sources case, one would only want to store independent data at each source node. However, this is not the case as illustrated in Figure 1.1. We consider a file represented as  $(a, b, c, d)$ , where each of the four components has unit-entropy, and a network where each edge has capacity 3. The cost of transmitting at rate  $x$  over edge  $e$  is  $c_e(x) = x$ , the cost of storage at the sources is 1 per unit storage. As shown in the figure, the case of partial replication when the source nodes contain dependent information has lower cost compared to the cases when the source nodes contain independent information or identical information (full replication).

The case of subset sources, is interesting for multiple reasons. For example, it may be the case that a given terminal is only interested in a part of the original file. In this case, if one places coded pieces of the original file at the source nodes, then the terminal may need to obtain a large number of coded pieces before it can recover the part that it is interested in. In the extreme case, if coding is performed across all the pieces of the file, then the terminal will need to recover all the sources before it can recover the part it is interested in. Note however, that in this work we do not explicitly consider scenarios where a given terminal requires a part

of the file. From a theoretical perspective as well, it is interesting to examine how much loss one incurs by not allowing coding at the sources.

### 1.1 Minimum Cost Multicast with Multiple Sources Problem

Several schemes have been proposed for content distribution over networks as discussed previously ([1][2]). In this section we briefly overview past work that is most closely related to the problem that we are considering.

Network coding has also been used in the area of large scale content distribution for different purposes. In the work [1], the authors proposed a content distribution scheme using network coding in a dynamic environment where nodes cooperate. A random linear coding based storage system (which is motivated by random network coding) was considered in [4] and shown to be more efficient than uncoded random storage system. However, their notion of efficiency is different than the total flow and storage cost considered in our work. The work of [5], proposed linear programming formulations for minimum cost flow allocation network coding based multicast. Lee et al. [6] constructed minimum cost subgraphs for the multicast of two correlated sources. It also proposed the problem of optimizing the correlation structure of sources and their placement. However, a solution was not presented there. Efficient algorithms for jointly allocating flows and rates were proposed for the multicast of a large number of correlated sources in [7]. The work of Jiang [8], considered a formulation that is similar to ours. It shows that under network coding, the problem of minimizing the joint transmission and storage cost can be formulated as a linear program. Furthermore, it considers a special class of networks called generalized tree networks and shows that there is no difference in the cost whether one considers subset sources or coded sources. In contrast, in this work we consider general networks, i.e., we do not assume any special structure of the network.

The work of Bhattad et al. [9] proposed an optimization problem formulation for cost minimization when some nodes are only allowed routing and forwarding instead of network coding. Our work on subset sources can be considered as an instance of this problem, by introducing a virtual super node and only allowing routing/forwarding on it. However, since

we consider a specific instance of this general problem, our problem formulation is much simpler than [9] and allows us to compare the cost of subset sources vs. coded sources. In addition, we recover stronger results in the case when there are only two or three source nodes. Furthermore, our solution approach is quite different and uses the concept of information measures.

## 1.2 Thesis Outline

In this work, we study the minimum joint cost for transmission and storage of multicast problem with both subset sources and coded sources. We also investigate the cost lost of transforming coded sources network to subset sources network. Our main contribution is 1) Formulation of the optimization problems by exploiting the properties of the information measure ([10]). We provide a precise formulation of the different optimization problems by leveraging the properties of the information measure (I-measure) introduced in [10]. This allows to provide a succinct formulation of the cost gap between the two cases and allows us to recover tight results in certain cases. This is contained in Chapter 2 and 3. 2) Cost comparison between subset sources case and coded sources case. The usage of the properties of information measure allows us to conclude that when there are two source nodes, there is no loss in considering subset sources. Furthermore, in the case of three source nodes, we derive an upper bound on the cost between the two cases that is shown to be tight. Finally, we propose a greedy algorithm to determine the cost gap for a given instance, that has been found to be tight in many cases. This is contained in Chapter 4.

Chapter 2 is the theoretical background of the set theory and information theory. We introduce the concept of atom, and then establish the correspondence of the set theory and information theory. Two important theorems related to the atom theory are given. We then present another theorem which will be used frequently in the subsequent chapters. An example is used to illustrate the theorem.

Chapter 3 gives two problem formulations of the subset case. The first one is useful in computation and easy to understand, and the second one will provide the insight of the cost analysis between the coded case and the subset case. The equivalence of the two formulations

is proved. Given the solutions of the problem, we describe the algorithm to decide the source contents. Finally, the code source case formulation is briefly introduced.

In chapter 4, another coded case formulation with atom structure for any arbitrary number of sources is presented. Given the solution of the new coded case, we formulate a linear program to compute the cost difference between the subset case and the coded case. A greedy algorithm is proposed to find a near optimal solution. We finally derive closed form upper bounds for three sources networks and two sources networks.

In chapter 5 provides some experiment results. A deterministic network is given to validate our previous analysis. The results for some random generated networks are also presented, from which we can see how frequently there will be a cost gap between the coded case and the subset case.

In Chapter 6, our conclusion gives a final perspective on our work.

## CHAPTER 2. MEASURE THEORY AND INFORMATION THEORY

### 2.1 Basic Concepts in Measure Theory

**Definition 1** The field  $\mathcal{F}_n$  generated by sets  $\tilde{X}_1, \tilde{X}_2, \dots, \tilde{X}_n$  is the collection of sets which can be obtained by any sequence of usual set operations (union, intersection, complement, and difference) on  $\tilde{X}_1, \tilde{X}_2, \dots, \tilde{X}_n$ .

**Definition 2** The atoms of  $\mathcal{F}_n$  are sets of the form  $\cap_{i=1}^n Y_i$ , where  $Y_i$  is either  $\tilde{X}_i$  or  $\tilde{X}_i^c$ , the complement of  $\tilde{X}_i$ .

There are  $2^n$  atoms and  $2^{2^n}$  sets in  $\mathcal{F}_n$ . All the atoms in  $\mathcal{F}_n$  are disjoint, and each set in  $\mathcal{F}_n$  can be expressed uniquely as the union of a subset of the atoms of  $\mathcal{F}_n$ . For example, in  $\mathcal{F}_2$ , there are 4 atoms:  $\tilde{X}_1 \cap \tilde{X}_2, \tilde{X}_1 \cap \tilde{X}_2^c, \tilde{X}_1^c \cap \tilde{X}_2$  and  $\tilde{X}_1^c \cap \tilde{X}_2^c$ .

**Definition 3** A real function  $\mu$  defined on  $\mathcal{F}_n$  is called a signed measure if it is set-additive, i.e., for disjoint set  $A$  and  $B$  in  $\mathcal{F}_n$ ,  $\mu(A \cup B) = \mu(A) + \mu(B)$ . For a signed measure  $\mu$ , we have  $\mu(\emptyset) = 0$ .

We use  $\mathcal{F}_n$  to denote the field generated by  $\tilde{X}_1, \tilde{X}_2, \dots, \tilde{X}_n$ . Define the universal set  $\Omega$  to be the union of the sets  $\tilde{X}_1, \tilde{X}_2, \dots, \tilde{X}_n$ , i.e.,  $\Omega = \cup_{i=1}^n \tilde{X}_i$ . The set  $A_0 = \cap_{i=1}^n \tilde{X}_i^c$  whose measure is  $\mu(\cap_{i=1}^n \tilde{X}_i^c) = \mu(\emptyset) = 0$ , is called the empty atom of  $\mathcal{F}_n$ . Let  $\mathcal{A}$  be the set of nonempty atoms of  $\mathcal{F}_n$ . Then  $|\mathcal{A}| = 2^n - 1$ . Because any set in  $\mathcal{F}_n$  can be uniquely defined as the union of some atoms, a signed measure  $\mu$  on  $\mathcal{F}_n$  is completely specified by the values of the  $\mu$  on the nonempty atoms of  $\mathcal{F}_n$ .

Consider a field  $\mathcal{F}_n$  generated by  $n$  sets  $\tilde{X}_1, \tilde{X}_2, \dots, \tilde{X}_n$ . Let  $\mathcal{N}_S = \{1, 2, \dots, n\}$  and  $\tilde{X}_V$  denote  $\cup_{i \in V} \tilde{X}_i$  for any nonempty subset  $V$  of  $\mathcal{N}_S$ .

**Theorem 1** *Define*

$$\mathcal{B} = \{\tilde{X}_V : V \text{ is a nonempty subset of } \mathcal{N}_S\}$$

*Then a signed measure  $\mu$  on  $\mathcal{F}_n$  is completely specified by  $\{\mu(B), B \in \mathcal{B}\}$ , which can be any set of real numbers.*

*Proof.* The number of elements in  $\mathcal{B}$  is equal to the number of nonempty subsets of  $\mathcal{N}_S$ , which is  $2^n - 1$ . Thus  $|\mathcal{A}| = |\mathcal{B}| = 2^n - 1$ . Let  $k = 2^n - 1$ . Let  $\mathbf{u}$  be a column  $k$ -vector of  $\mu(A)$ ,  $A \in \mathcal{A}$ , and  $\mathbf{h}$  be a column  $k$ -vector of  $\mu(B)$ ,  $B \in \mathcal{B}$ . Since all the sets in  $\mathcal{B}$  can be expressed uniquely as the union of some nonempty atoms in  $\mathcal{A}$ , by the set-additivity of  $\mu$ , for each  $B \in \mathcal{B}$ ,  $\mu(B)$  can be expressed uniquely as the sum of some components of  $\mathbf{u}$ . Thus  $\mathbf{h} = C_n \mathbf{u}$ , where  $C_n$  is a unique  $k \times k$  matrix. On the other hand, it can be shown that (see Appendix ) for each  $A \in \mathcal{A}$ ,  $\mu(A)$  can be expressed as a linear combination of  $\mu(B)$ ,  $B \in \mathcal{B}$ .

However, the existence of the said expression does not imply its uniqueness. Nevertheless, we can write  $\mathbf{u} = D_n \mathbf{h}$  for some  $k \times k$  matrix  $D_n$ . We then obtain  $\mathbf{u} = (D_n C_n) \mathbf{h}$  which implies that  $D_n$  is the inverse of  $C_n$  as the equality holds regardless of the choice of  $\mu$ . Since  $C_n$  is unique, so is  $D_n$ . Therefore, there is a unique linear relationship between  $\mu(A)$ ,  $A \in \mathcal{A}$  and  $\mu(B)$ ,  $B \in \mathcal{B}$ .  $\square$

Since  $\mathcal{F}_n$  can be completely specified by  $\mu(A)$ ,  $\mathcal{F}_n$  can also be completely specified by  $\mu(B)$ .

## 2.2 One-to-one Correspondence between Shannon Information Measure and Set Theory

For  $n$  random variables  $X_1, X_2, \dots, X_n$ , let  $\tilde{X}_i$  be a set corresponding to  $X_i$ . Let  $X_V = (X_i, i \in V)$ , where  $V$  is some nonempty subset of  $\mathcal{N}_s$ . Construct the signed measure  $\mu^*(\tilde{X}_V) = H(X_V)$ , for all nonempty subset  $V$  of  $\mathcal{N}_S$ .

**Theorem 2**  $\mu^*$  is the unique signed measure on  $\mathcal{F}_n$  which is consistent with all Shannon's information measures.

*Proof.* Consider

$$\begin{aligned}
& \mu^*(\tilde{X}_G \cap \tilde{X}_{G'} - \tilde{X}_{G''}) \\
&= \mu^*(\tilde{X}_{G \cup G''}) + \mu^*(\tilde{X}_{G' \cup G''}) - \mu^*(\tilde{X}_{G \cup G' \cup G''}) - \mu^*(\tilde{X}_{G''}) \\
&= H(X_{G \cup G''}) + H(X_{G' \cup G''}) - H(X_{G \cup G' \cup G''}) - H(X_{G''}) \\
&= I(X_G; X_{G'} | X_{G''})
\end{aligned} \tag{2.1}$$

When  $G'' = \emptyset$ , the equation becomes  $\mu^*(\tilde{X}_G \cap \tilde{X}_{G'}) = I(X_G; X_{G'})$ .

When  $G = G'$ , the equation becomes  $\mu^*(\tilde{X}_G - \tilde{X}_{G'}) = H(X_G | X_{G''})$ .

When  $G = G'$  and  $G'' = 0$ , the equation becomes  $\mu^*(\tilde{X}_G) = H(X_G)$ .

Then  $\mu^*$  is the unique signed measure on  $\mathcal{F}_n$  which is consistent with all Shannon's information measures.  $\square$

The measure is consistent with all the Shannon information. Specifically, in each of these equations, the left hand side and right hand side correspond to each other via the following substitution of symbols:

$$\begin{aligned}
H/I &\leftrightarrow \mu^* \\
, &\leftrightarrow \cup \\
; &\leftrightarrow \cap \\
| &\leftrightarrow -
\end{aligned} \tag{2.2}$$

Hence,  $\mu^*(B)$ 's for  $B \in \mathcal{B}$  can represent all joint entropies. Because there is a unique linear relationship between  $\mu^*(A)$  for  $A \in \mathcal{A}$  and  $\mu^*(B)$  for  $B \in \mathcal{B}$ , we can use the nonnegativity of the linear combination of  $\mu^*(A)$  where  $A \in \mathcal{A}$ , to represent all the information inequalities.

### 2.3 Main Theorem

Let  $\mathcal{N}_S = \{1, 2, \dots, n\}$  and consider  $n$  random variables  $X_1, X_2, \dots, X_n$ . Let  $\tilde{X}_i$  be a set corresponding to  $X_i$  and let  $\tilde{X}_V = \cup_{i \in V} \tilde{X}_i$ . We denote the set of nonempty atoms of  $\mathcal{F}_n$  by  $\mathcal{A}$ , where  $\mathcal{F}_n$  is the field generated by the sets  $\tilde{X}_1, \tilde{X}_2, \dots, \tilde{X}_n$ . Similarly,  $X_V$  denotes the collection of random variables  $(X_i, i \in V)$ , where  $V \subseteq \mathcal{N}_S$ . Construct the signed measure  $\mu^*(\tilde{X}_V) = H(X_V)$ , for all nonempty subset  $V$  of  $\mathcal{N}_S$ .

**Theorem 3** (1) Suppose that there exist a set of  $2^n - 1$  nonnegative values, one corresponding to each atom of  $\mathcal{F}_n$ , i.e.,  $\alpha(A) \geq 0, \forall A \in \mathcal{A}$ . Then, we can define a set of independent random variables,  $W_A, A \in \mathcal{A}$  and construct random variables  $X_j = (W_A : A \in \mathcal{A}, A \subset \tilde{X}_j)$ , such that the measures of the nonempty atoms of the field generated by  $\tilde{X}_1, \tilde{X}_2, \dots, \tilde{X}_n$  correspond to the values of  $\alpha$ , i.e.,  $\mu^*(A) = \alpha(A), \forall A \in \mathcal{A}$ .

(2) Conversely, let  $Z_i, i \in \{1, \dots, m\}$  be a collection of independent random variables. Suppose that a set of random variables  $X_i, i = 1, \dots, n$  is such that  $X_i = Z_{V_i}$ , where  $V_i \subseteq \{1, \dots, m\}$ . Then the set of atoms of the field generated by  $\tilde{X}_1, \tilde{X}_2, \dots, \tilde{X}_n$ , have non-negative measures.

*proof:* (1) Independent random variables  $W_A, A \in \mathcal{A}$ , such that  $H(W_A) = \alpha(A)$  can be constructed [10]. It only remains to check the consistency of the measures. For this, we have, for all  $V \subseteq \mathcal{N}_S$ ,

$$H(X_V) = \sum_{A \in \mathcal{A}: A \subset \tilde{X}_V} H(W_A), \quad (2.3)$$

using the independence of the  $W_A$ 's. On the other hand we know that

$$H(X_V) = \mu^*(\tilde{X}_V) = \sum_{A \in \mathcal{A}: A \subset \tilde{X}_V} \mu^*(A). \quad (2.4)$$

Equating these two we have, for all  $V \subseteq \mathcal{N}_S$ ,

$$\sum_{A \in \mathcal{A}: A \subset \tilde{X}_V} H(W_A) = \sum_{A \in \mathcal{A}: A \subset \tilde{X}_V} \mu^*(A) \quad (2.5)$$

Now, one possible solution to this is that  $\mu^*(A) = H(W_A), \forall A \in \mathcal{A}$ . By the uniqueness of  $\mu^*$  [10], we know that this is the only solution.

(2) We will prove all the measures are nonnegative by induction. Without loss of generality, we can order  $\tilde{X}_i$ 's in an arbitrary way, then we can analyze the measure

$$\mu^*(\tilde{X}_1 \cap \dots \cap \tilde{X}_l \cap_{k: k \in K} \tilde{X}_k^c)$$

where  $K \subseteq \mathcal{N}_S \setminus \{1, 2, \dots, l\}$ .

When  $l = 1$ , the measure corresponds to conditional entropy,  $\forall K \subseteq \mathcal{N}_S \setminus \{1\}$

$$\mu^*(\tilde{X}_1 \cap_{k: k \in K} \tilde{X}_k^c) = H(X_1 | X_K) \geq 0$$

When  $l = 2$ , we have,  $\forall K \subseteq \mathcal{N}_S \setminus \{1, 2\}$

$$\begin{aligned} \mu^*(\tilde{X}_1 \cap \tilde{X}_2 \cap_{k:k \in K} \tilde{X}_k^c) &= I(X_1; X_2 | X_K) \\ &= H(X_1, X_K) + H(X_2, X_K) - H(X_K) - H(X_1, X_2, X_K) \\ &= \sum_{i \in V_1 \cap V_2 \cap_{k:k \in K} V_k^c} H(Z_i) \geq 0 \end{aligned}$$

Assume for  $l = j$ ,  $\forall K \subseteq \mathcal{N}_S \setminus \{1, 2, \dots, j\}$ , the following statement holds,

$$\mu^*(\tilde{X}_1 \cap \dots \cap \tilde{X}_j \cap_{k:k \in K} \tilde{X}_k^c) = \sum_{i \in V_1 \cap \dots \cap V_j \cap_{k:k \in K} V_k^c} H(Z_i) \quad (2.6)$$

When  $l = j + 1$ ,  $\forall K \subseteq \mathcal{N}_S \setminus \{1, 2, \dots, j + 1\}$ , we will have

$$\begin{aligned} \mu^*(\tilde{X}_1 \cap \dots \cap \tilde{X}_{j+1} \cap_{k:k \in K} \tilde{X}_k^c) &= \mu^*(\tilde{X}_1 \cap \dots \cap \tilde{X}_j \cap_{k:k \in K} \tilde{X}_k^c) \\ &\quad - \mu^*(\tilde{X}_1 \cap \dots \cap \tilde{X}_j \cap \tilde{X}_{j+1}^c \cap_{k:k \in K} \tilde{X}_k^c) \\ &\stackrel{(a)}{=} \sum_{i \in V_1 \cap \dots \cap V_j \cap_{k:k \in K} V_k^c} H(Z_i) \\ &\quad - \sum_{i \in V_1 \cap \dots \cap V_j \cap V_{j+1}^c \cap_{k:k \in K} V_k^c} H(Z_i) \\ &\stackrel{(b)}{=} \sum_{i \in V_1 \cap \dots \cap V_{j+1} \cap_{k:k \in K} V_k^c} H(Z_i) \geq 0 \end{aligned}$$

The equation (a) is due to the assumption (2.6). The equation (b) is due to the independence of  $Z_i$ 's,  $i \in \{1, \dots, m\}$ . Therefore, we have shown that  $j \leq n$ ,  $\forall K \subseteq \mathcal{N}_S \setminus \{1, 2, \dots, j\}$ ,

$$\mu^*(\tilde{X}_1 \cap \dots \cap \tilde{X}_j \cap_{k:k \in K} \tilde{X}_k^c) = \sum_{i \in V_1 \cap \dots \cap V_j \cap_{k:k \in K} V_k^c} H(Z_i) \geq 0$$

In a similar manner it is easy to see that all atom are non-negative.  $\square$

We note in passing that it is well-known that atom measures can be negative for general probability distributions [10].

### 2.3.1 An Example

We now give an example in which not all the sources are subsets of the universal independent information set, then not all atoms are nonnegative.

The collection of random variables are  $(Z_1, Z_2)$ , where  $Z_1$  and  $Z_2$  are independent random variables with

$$P(Z_i = 0) = P(Z_i = 1) = 0.5$$

$i = 1, 2$ . Let

$$X_1 = Z_1, X_2 = Z_2, X_3 = (Z_1 + Z_2) \bmod 2$$

Then  $X_1$  and  $X_2$  are subsets of  $(Z_1, Z_2)$ , but  $X_3$  is not a subset of  $(Z_1, Z_2)$ .

Then we will see not all the atoms are non-negative.

$$\begin{aligned} \mu^*(\tilde{X}_1 \cap \tilde{X}_2 \cap \tilde{X}_3) &= \mu^*(\tilde{X}_1 \cap \tilde{X}_2) - \mu^*(\tilde{X}_1 \cap \tilde{X}_2 \cap \tilde{X}_3^c) \\ &= I(X_1; X_2) - I(X_1; X_2 | X_3) = -1 \end{aligned}$$

### CHAPTER 3. PROBLEM FORMULATION

In this section we consider the following problem. Suppose that there is a source (we will refer to this as the “original source” for convenience) that can be split into arbitrarily small pieces, e.g. a huge movie file of size  $50Gb$  can be considered to be consisting of  $50e09$  bits, that needs to be transmitted to all the terminals. Assume that we have the flexibility to place portions of the file consisting of collections of these pieces at various source nodes. The motivational example in Figure 1.1 is an instance of this. The portions need not necessarily be subsets of the bits, they may be arbitrary functions of them. We want to decide the content of the portions so that the joint cost of storing and transmitting them over the network is minimized. We consider two different cases.

- i) *Subset Sources*. In this case, each source node only contains a subset of the pieces of the original source.
- ii) *Coded Sources*. We allow the portions to be arbitrary functions of the pieces of the original source.

Under both cases we will allow all nodes in the network to perform network coding. We abstract this problem as follows. Given a directed acyclic graph  $G = (V, E, C)$  that represents the network.  $V$  denotes the set of vertices,  $E$  denotes the set of edges, and  $C_{ij}$  denotes the capacity constraint for edge  $(i, j) \in E$ . There is a set of source nodes  $S \subset V$  (numbered  $1, \dots, n$ ) and terminal nodes  $T \subset V$ . Suppose that the original source can be represented as the collection of equal entropy independent sources  $\{OS_j\}_{j=1}^Q$ , where  $Q$  is a sufficiently large integer. This assumption is equivalent to assuming that a file can be split into arbitrarily small pieces. Let  $X_i$  represent the source at the  $i^{th}$  source node, e.g., this represents the subset of  $\{OS_j\}_{j=1}^Q$  that are available at the  $i^{th}$  node. Suppose that each edge  $(i, j)$  incurs a linear

cost  $f_{ij}z_{ij}$  for a flow of value  $z_{ij}$  over it, and each source incurs a linear cost  $d_i H(X_i)$  for the information  $X_i$  stored.

### 3.1 Subset Constraints Source Network

In this case each source  $X_i, i = 1, \dots, n$  is constrained to be a subset of the pieces of the original source. We leverage Theorem 3 from the previous section that tells us that in this case that  $\mu^*(A) \geq 0$  for all  $A \in \mathcal{A}$ . In the discussion below, we will pose this problem as one of recovering the measures of the  $2^n - 1$  atoms. Note that this will in general result in fractional values. We shall consider a large enough time-step, so that all flows and atom measures can be treated as integers. Following this we will present an algorithm that forms a source corresponding to each atom, called the atom source. The requirement that the original source be arbitrarily divisible is needed here.

#### 3.1.1 Basic Formulation

We construct an augmented graph  $G_1^* = (V_1^*, E_1^*, C_1^*)$  as follows (see Figure 3.1(a)). Append a virtual super node  $s^*$  and  $2^n - 1$  virtual nodes corresponding to the atom sources  $W_A, \forall A \in \mathcal{A}$  and connect  $s^*$  to each  $W_A$  source node. The node for  $W_A$  is connected to a source node  $i \in S$  if  $A \subset \tilde{X}_i$ . The capacity of the new (virtual) edges is set to infinity. The cost of the edge from  $s^*$  to the atom node for  $W_A$  is  $\sum_{\{i \in S: A \subset \tilde{X}_i\}} d_i$ . The cost of the edges between the atom nodes and  $S$  is set to zero.

If each terminal can recover all the atom sources,  $W_A, \forall A \in \mathcal{A}$ , then it can in turn recover the original source. The information that needs to be stored at the source node  $i \in S$ , is equal to the sum of flows from  $s^*$  to  $W_A, \forall A \subset \tilde{X}_i$ . Let  $x_{ij}^{(t)}, t \in T$  represent the flow variable over  $G_1^*$  corresponding to the terminal  $t$  along edge  $(i, j)$  and let  $z_{ij}$  represent  $\max_{t \in T} x_{ij}^{(t)}, \forall (i, j) \in E$ . The corresponding optimization problem is defined as ATOM-SUBSET-MIN-COST, and can be expressed as,

$$\text{minimize } \sum_{(i,j) \in E} f_{ij}z_{ij} + \sum_{A \in \mathcal{A}} \left( \sum_{\{i \in S: A \subset \tilde{X}_i\}} d_i \right) \mu^*(A)$$

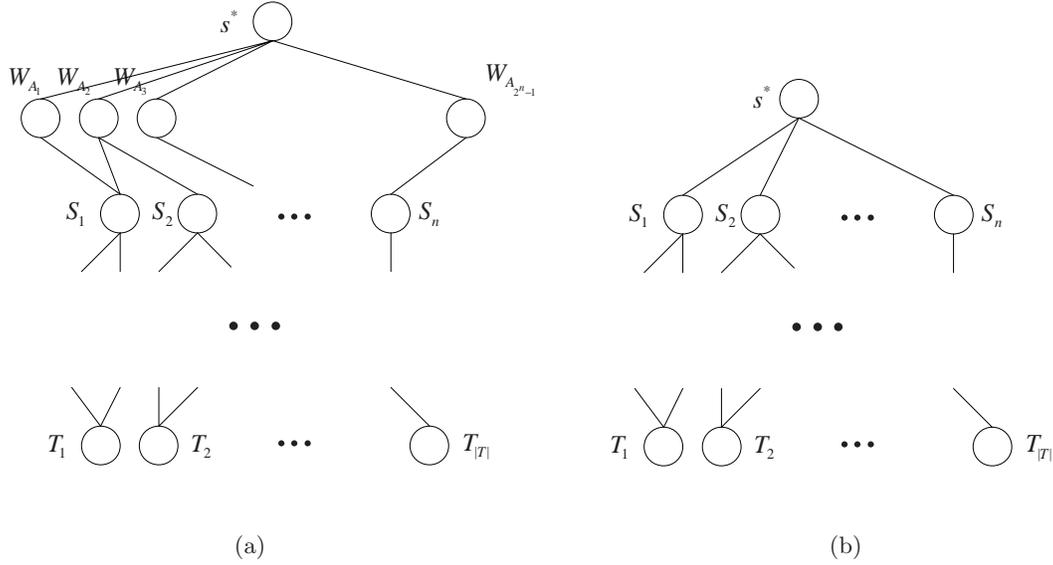


Figure 3.1 (a) The network structure for the first formulation. (b) The network structure for the second formulation.

subject to  $0 \leq x_{ij}^{(t)} \leq z_{ij} \leq c_{ij,1}^*, \forall (i, j) \in E_1^*, t \in T$

$$\sum_{\{j|(i,j) \in E_1^*\}} x_{ij}^{(t)} - \sum_{\{j|(j,i) \in E_1^*\}} x_{ji}^{(t)} = \sigma_i^{(t)}, \forall i \in V_1^*, t \in T$$

$$x_{s^*W_A}^{(t)} = \mu^*(A), \forall t \in T, A \in \mathcal{A} \quad (3.1)$$

$$\mu^*(A) \geq 0, \forall A \in \mathcal{A} \quad (3.2)$$

$$H(X_1, X_2, \dots, X_n) = \sum_{A: A \in \mathcal{A}} \mu^*(A) \quad (3.3)$$

where

$$\sigma_i^{(t)} = \begin{cases} H(X_1, \dots, X_n) & \text{if } i = s^* \\ -H(X_1, \dots, X_n) & \text{if } i = t \\ 0 & \text{otherwise} \end{cases}$$

This is basically the formulation of the minimum cost multicast problem [5] with a virtual super-source of entropy  $H(X_1, \dots, X_n)$ , with the added constraint that the flow on the edge from  $s^*$  to node  $W_A$  for each terminal,  $x_{s^*W_A}^{(t)}$  is at least  $\mu^*(A)$ . We also have a constraint that  $\sum_{A \in \mathcal{A}} \mu^*(A) = H(X_1, X_2, \dots, X_n)$ , that in turns yields the constraint that  $x_{s^*W_A}^{(t)} = \mu^*(A)$ .

Also, note that the measure of each atom,  $\mu^*(A)$  is non-negative. This enforces the subset constraints.

In general, the proposed LP formulation has a number of constraints that is exponential in the number of source nodes, since there are  $2^n - 1$  nonnegative atoms. However, when the number of source nodes is small, this formulation can be solved using regular LP solvers. We emphasize though, that the formulation of this problem in terms of the atoms of the distribution of the sources provides us with a mechanism of reasoning about the case of subset constraints, under network coding. We are unaware of previous work that proposes a formulation of this problem.

In order to solve this problem, we can instead consider the equivalent optimization problem:

minimize

$$\sum_{(i,j) \in E} f_{ij} z_{ij} + \sum_{A \in \mathcal{A}} \left( \sum_{\{i \in S: AC \tilde{X}_i\}} d_i \right) \mu^*(A) + \sum_{A \in \mathcal{A}} \sum_t \lambda_A^t (x_{s^* W_A}^{(t)} - \mu^*(A))$$

subject to  $0 \leq x_{ij}^{(t)} \leq z_{ij} \leq c_{ij,1}^*, \forall (i,j) \in E_1^*, t \in T$

$$\sum_{\{j|(i,j) \in E_1^*\}} x_{ij}^{(t)} - \sum_{\{j|(j,i) \in E_1^*\}} x_{ji}^{(t)} = \sigma_i^{(t)}, \forall i \in V_1^*, t \in T$$

$$\mu^*(A) \geq 0, \forall A \in \mathcal{A} \quad (3.4)$$

$$H(X_1, X_2, \dots, X_n) = \sum_{A \in \mathcal{A}} \mu^*(A) \quad (3.5)$$

where

$$\sigma_i^{(t)} = \begin{cases} H(X_1, X_2, \dots, X_n) & \text{if } i = s^* \\ -H(X_1, X_2, \dots, X_n) & \text{if } i = t \\ 0 & \text{otherwise} \end{cases}$$

and  $\lambda_A^t$  is a dual variable which has  $\lambda_A^t \geq 0, \forall t \in T, A \in \mathcal{A}$ .

This problem can be decomposed as two separate problems

minimize

$$\sum_{(i,j) \in E} f_{ij}(z_{ij}) + \sum_{A \in \mathcal{A}} \sum_t \lambda_A^t x_{s^* W_A}^{(t)}$$

subject to  $0 \leq x_{ij}^{(t)} \leq z_{ij} \leq c_{ij,1}^*, \forall (i, j) \in E_1^*, t \in T$

$$\sum_{\{j|(i,j) \in E_1^*\}} x_{ij}^{(t)} - \sum_{\{j|(j,i) \in E_1^*\}} x_{ji}^{(t)} = \sigma_i^{(t)}, \forall i \in V_1^*, t \in T$$

and

minimize

$$\sum_{A \in \mathcal{A}} \left( \sum_{\{i \in S: A \subset \tilde{X}_i\}} d_i - \sum_t \lambda_A^t \right) \mu^*(A)$$

subject to

$$\mu^*(A) \geq 0, \forall A \in \mathcal{A}$$

$$H(X_1, X_2, \dots, X_n) = \sum_{A \in \mathcal{A}} \mu^*(A)$$

The first minimization problem is a standard multicast problem. The second minimization problem can be solved in closed form after we get a set of  $\lambda_A^t$ : We assign  $\mu^*(A) = H(X_1, X_2, \dots, X_n)$ , such that  $\mu^*(A)$  corresponds to the smallest  $\sum_{\{i \in S: A \subset \tilde{X}_i\}} d_i - \sum_t \lambda_A^t, \forall A \in \mathcal{A}$ . Other  $\mu^*(A)$  are assigned to be zero. We use subgradient optimization to recover the dual variables. However, the subgradient optimization does not necessarily yield a primal optimal solution. There are many methods for recovering primal solutions, among them, we use the method introduced by Sherali and Choi [11].

We now give a brief description of the primal recovery procedure of [11]. Let  $\beta_j[k]$  for  $j = 1, \dots, k$  be a set of convex combination weights for each  $k \geq 1$ , i.e.,  $\sum_{j=1}^k \beta_j[k] = 1$ , and  $\beta_j[k] \geq 0$ , where  $k$  denotes the  $k$ th iteration. We define  $\gamma_j[k] = \beta_j[k]/\theta[k]$ , for  $1 \leq j \leq k$ , and  $k \geq 1$ , and let

$$\Delta\gamma^{\max}[k] = \max\{\gamma_j[k] - \gamma_{(j-1)}[k] : j = 2, \dots, k\}.$$

Let the primal solution returned by subgradient optimization at iteration  $k$  be denoted by the vector  $(z, x, \mu^*)[k]$  and let the  $k$ th primal iterate be defined as

$$(\tilde{z}, \tilde{x}, \tilde{\mu}^*)[k] = \sum_{j=1}^k \beta_j[k] (z, x, \mu^*)[j], \text{ for } k \geq 1.$$

Suppose that the sequence of weights  $\beta_j[k]$  for  $k \geq 1$  and the sequence of step size  $\theta[k], k \geq 1$  are chosen such that

- (1)  $\gamma_j[k] \geq \gamma_{j-1}[k]$  for all  $j = 2, \dots, k$  for each  $k$ .
- (2)  $\Delta\gamma^{\max}[k] \rightarrow 0$ , as  $k \rightarrow \infty$ , and
- (3)  $\gamma_1[k] \rightarrow 0$  as  $k \rightarrow \infty$  and  $\gamma_k[k] \leq \delta$  for all  $k$ , for all  $\delta > 0$ .

Then an optimal solution to the primal problem can be obtained from any accumulation point of the sequence of primal iterates  $(\tilde{z}, \tilde{x}, \tilde{\mu}^*)$ .

Some useful choices for the step sizes  $\theta[k]$  and the convex combination weights  $\beta_j[k]$  that satisfy these conditions are given below.

- (1)  $\theta[k] = a/(b + ck)$ , for  $k \geq 1$  where  $a > 0, b \geq 0$ , and  $c \geq 0$  and  $\beta_j[k] = 1/k$  for all  $j = 1, \dots, k$ .
- (2)  $\theta[k] = k^{-\alpha}$ , for  $k \geq 1$  where  $0 < \alpha < 1$  and  $\beta_j[k] = 1/k$  for all  $j = 1, \dots, k$ .

We can run the simulations for the two decomposed problems separately, keep records of  $(z, x, \mu^*)[k]$  and compute  $(\tilde{z}, \tilde{x}, \tilde{\mu}^*)[k]$  until it converges.

In order to provide bounds on the gap between the optimal costs of the subset sources case and the coded sources case, we now present an alternate formulation of this optimization, that is more amenable to the gap analysis. Note however, that this alternate formulation has more constraints than the one presented above.

### 3.1.2 Another Formulation

In the first formulation, the terminals first recover the atom sources, and then the original source. In this alternate formulation, we pose the problem as one of first recovering all the sources,  $X_i, i \in S$  at each terminal and then the original source. Note that since these sources are correlated, this formulation is equivalent to the Slepian-Wolf problem over a network [7]. We will first give the problem formulation and then prove the two formulations have the same optimums.

We construct another augmented graph  $G_2^* = (V_2^*, E_2^*, C_2^*)$  using the basic network graph  $G = (V, E, C)$ . We append a virtual super node  $s^*$  to  $G$ , and connect  $s^*$  and each source node

$i$  with virtual edges, such that its capacity is infinity and its cost is  $d_i$ . The structure of the network is shown in Figure 3.1(b).

As before, let  $x_{ij}^{(t)}$ ,  $t \in T$  represent the flow variable over  $G_2^*$  corresponding to the terminal  $t$  along edge  $(i, j)$  and let  $z_{ij}$  represent  $\max_{t \in T} x_{ij}^{(t)}$ ,  $\forall (i, j) \in E$ . We introduce variable  $R_i^{(t)}$ ,  $t \in T$  that represents the rate from source  $i$  to terminal  $t$ ,  $i = 1, \dots, n$ . Thus  $R^{(t)} = (R_1^{(t)}, R_2^{(t)}, \dots, R_n^{(t)})$  represents the rate vector for terminal  $t$ . In order for terminal  $t$  to recover the sources, the rate vector  $R^{(t)}$  needs to lie within the Slepian-Wolf region of the sources, which is defined as follows.

$$\mathcal{R}_{SW} = \{(R_1, \dots, R_n) : \forall U \subseteq S, \sum_{i \in U} R_i \geq H(X_U | X_{S \setminus U})\}$$

Moreover, the rates also need to be in the capacity region such that the network has enough capacity to support them for each terminal. As before we have  $\mu^*(A) \geq 0, \forall A \in \mathcal{A}$  to enforce the subset constraint. The optimization problem is defined as SUBSET-MIN-COST. The formulation is as follows.

$$\text{minimize } \sum_{(i,j) \in E} f_{ij} z_{ij} + \sum_{A \in \mathcal{A}} (\sum_{\{i \in S: AC \tilde{X}_i\}} d_i) \mu^*(A)$$

subject to

$$0 \leq x_{ij}^{(t)} \leq z_{ij} \leq c_{ij,2}^*, (i, j) \in E_2^*, t \in T \quad (3.6)$$

$$\begin{aligned} \sum_{\{j|(i,j) \in E_2^*\}} x_{ij}^{(t)} - \sum_{\{j|(j,i) \in E_2^*\}} x_{ji}^{(t)} &= \sigma_i^{(t)}, i \in V_2^*, t \in T \\ x_{s^*i}^{(t)} &\geq R_i^{(t)}, \forall i \in S, t \in T \end{aligned} \quad (3.7)$$

$$R^{(t)} \in \mathcal{R}_{SW}, \forall t \in T \quad (3.8)$$

$$\mu^*(A) \geq 0, \forall A \in \mathcal{A} \quad (3.9)$$

$$z_{s^*i} = H(X_i), \forall i \in S \quad (3.10)$$

$$H(X_1, X_2, \dots, X_n) = \sum_{A \in \mathcal{A}} \mu^*(A) \quad (3.11)$$

where

$$\sigma_i^{(t)} = \begin{cases} H(X_1, X_2, \dots, X_n) & \text{if } i = s^* \\ -H(X_1, X_2, \dots, X_n) & \text{if } i = t \\ 0 & \text{otherwise} \end{cases}$$

Though not expressed explicitly, each conditional entropy term  $H(X_U|X_{S \setminus U})$  needs to be equal  $\sum_{A:A \not\subseteq \tilde{X}_{S \setminus U}} \mu^*(A)$ , and each marginal entropy  $H(X_i)$  needs to be equal  $\sum_{A:A \subseteq \tilde{X}_i} \mu^*(A)$ , so that the atom measures and the entropies are consistent.

Now we prove the two formulations will get the same optimal values. The basic idea is as follows. Note that the objective function for both the formulations is exactly the same. We will first consider the optimal solution for the first formulation and construct a solution for the second formulation so that we can conclude that  $f_{opt1} \geq f_{opt2}$ . In a similar manner we will obtain the reverse inequality, which will establish equality of the two optimal values.

Suppose that we are given the optimal set of flows  $x_{ij,1}^{(t)}, z_{ij,1}, t \in T, (i, j) \in E_1^*$  and the optimal atom values  $\mu^*(A)_1$  for the first formulation. Further assume that the optimal objective function is  $f_{opt1}$ .

**Claim 1** In  $G_2^*$ , for the flows  $x_{ij,2}^{(t)}, z_{ij,2}$ , and the atoms  $\mu^*(A)_2$ , assign

$$\begin{aligned} x_{ij,2}^{(t)} &= x_{ij,1}^{(t)}, \quad z_{ij,2} = z_{ij,1}, \quad \forall (i, j) \in G \\ x_{s^*i,2}^{(t)} &= \sum_{A:A \subseteq \tilde{X}_i} x_{W_{A^*i},1}^{(t)}, \quad z_{s^*i,2} = \sum_{A:A \subseteq \tilde{X}_i} \mu^*(A)_1, \quad \forall i \in S \\ \mu^*(A)_2 &= \mu^*(A)_1, \quad \forall A \in \mathcal{A}. \end{aligned}$$

Then  $x_{ij,2}^{(t)}, z_{ij,2}$ , and the atoms  $\mu^*(A)_2$  are a feasible solution for the second formulation.

*Proof.* The flows balance for the source node in the first formulation  $\sum_{A:A \subseteq \tilde{X}_i} x_{W_{A^*i},1}^{(t)} = \sum_{j:(i,j) \in E_1^*} x_{ij,1}^{(t)}$ , the flow balance for the source node in the second formulation:  $x_{s^*i,2}^{(t)} = \sum_{A:A \subseteq \tilde{X}_i} x_{W_{A^*i},1}^{(t)} = \sum_{j:(i,j) \in E_1} x_{ij,1}^{(t)} = \sum_{j:(i,j) \in E_2^*} x_{ij,2}^{(t)}, \forall i \in S$ , the flow balance at the source node is satisfied. We only need to check constraints (3.7) and (3.8).

For any  $U \subseteq S$ , we will have

$$\begin{aligned}
\sum_{i:i \in U} x_{s^*i,2}^{(t)} &= \sum_{i:i \in U} \sum_{A:AC\tilde{X}_i} x_{W_Ai,1}^{(t)} \\
&= \sum_{A:A \not\subseteq \tilde{X}_{S \setminus U}, AC\tilde{X}_U} x_{W_Ai,1}^{(t)} + \sum_{A:AC \subseteq \tilde{X}_{S \setminus U}, AC\tilde{X}_U} x_{W_Ai,1}^{(t)} \\
&\geq \sum_{A:A \not\subseteq \tilde{X}_{S \setminus U}, AC\tilde{X}_U} x_{W_Ai,1}^{(t)} \\
&= \sum_{A:A \not\subseteq \tilde{X}_{S \setminus U}, AC\tilde{X}_U} x_{s^*W_A,1}^{(t)} \\
&\stackrel{(a)}{=} \sum_{A:A \not\subseteq \tilde{X}_{S \setminus U}, AC\tilde{X}_U} \mu^*(A)_1 \\
&= \sum_{A:A \not\subseteq \tilde{X}_{S \setminus U}, AC\tilde{X}_U} \mu^*(A)_2 \\
&= H(X_U | X_{S \setminus U})
\end{aligned}$$

where  $H(X_U | X_{S \setminus U})$  is the conditional entropy of the second formulation. (a) comes from the constraint (3.1) in formulation 1. Therefore, constraints (3.7) and (3.8) are satisfied and this assignment is feasible for the second formulation with a cost equal to  $f_{opt1}$ .  $\square$

We conclude that the optimal solution for the second formulation  $f_{opt2}$  will have  $f_{opt2} \leq f_{opt1}$ .

Next we show the inequality in the reverse direction. Suppose that we are given the optimal set of flows  $x_{ij,2}^{(t)}, z_{ij,2}, t \in T, (i, j) \in E_2^*$  and the atom values  $\mu^*(A)_2$  in the second formulation, with an objective of value  $f_{opt2}$ .

**Claim 2** In  $G_1^*$ , assign

$$\begin{aligned}
x_{ij,1}^{(t)} &= x_{ij,2}^{(t)}, z_{ij,1} = z_{ij,2}, \forall (i, j) \in G \\
z_{s^*W_A,1} &= x_{s^*W_A,1}^{(t)} = \mu^*(A)_1 = \mu^*(A)_2, \forall A \in \mathcal{A}
\end{aligned}$$

Furthermore, there exist flow variables  $x_{W_Ai,1}^{(t)}$  and  $z_{W_Ai,1}$  over the edge  $(W_A, i) \in V_1^*, \forall A \in \mathcal{A}$ , such that together with the assignment above, they form a feasible solution for the first formulation.

*Proof.* It is clear that the assignments for  $x_{ij,1}^{(t)}$  and  $z_{ij,1}$  for  $(i, j) \in G$  satisfy the required constraints. We essentially need to demonstrate the existence of flow variables  $x_{W_Ai,1}^{(t)}$  and

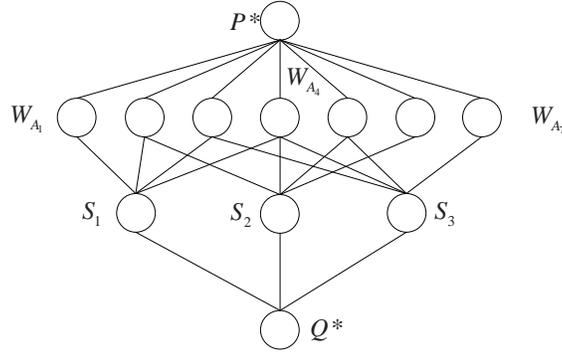


Figure 3.2 The auxiliary graph for three sources case

$z_{W_A i,1}$  over the edge  $(W_A, i) \in V_1^*$ ,  $\forall A \in \mathcal{A}$ , such that they satisfy the flow balance constraints at all the concerned nodes.

Towards this end it is convenient to construct an auxiliary graph as follows. There is a source node  $P^*$  connected to the atoms  $W_A$ 's,  $A \in \mathcal{A}$ , a terminal  $Q^*$  connected to the sources nodes,  $i \in S$ . There is an edge connecting  $W_A$  and  $i$  if  $A \subset \tilde{X}_i$ . An example is shown in Figure 3.2 in the case of three source nodes. The capacities on the edges are assigned as follows. The capacity for edge  $(P^*, W_A)$  is  $x_{s^* W_A,1}^{(t)}$ , the capacity for edge  $(i, Q^*)$  is  $x_{s^* i,2}^{(t)}$ , and the capacity for edge  $(W_A, i)$  is infinity. Note that  $\sum_{A \in \mathcal{A}} x_{s^* W_A,1}^{(t)} = \sum_{i \in S} x_{s^* i,2}^{(t)} = H(X_1, X_2, \dots, X_n)$ . Therefore, if we can show that the maximum flow in this auxiliary graph between  $P^*$  and  $Q^*$  is  $H(X_1, X_2, \dots, X_n)$ , this would imply the existence of flow variables on the edges between the atom nodes and the source nodes that satisfy the required flow balance conditions.

To show this we use the max-flow min-cut theorem and instead show that the minimum value over all cuts separating  $P^*$  and  $Q^*$  is  $H(X_1, X_2, \dots, X_n)$ .

First, notice that there is a cut with value  $H(X_1, X_2, \dots, X_n)$ . This cut can be simply the node  $P^*$ , since the sum of the capacities of its outgoing edges is  $H(X_1, X_2, \dots, X_n)$ . Next, if an atom node  $W_A$ ,  $A \in \mathcal{A}$  belongs to the cut that contains  $P^*$ , then we must have all source nodes  $i \in S$  such that  $A \subset \tilde{X}_i$  also belonging to the cut. To see this, note that otherwise there is at least one edge belonging to the cut whose capacity is infinity, i.e., the cut cannot be the

minimum cut.

Let  $S' \subseteq S$ . Based on this argument it suffices to consider cuts that contain,  $P^*$ , the set of nodes  $S \setminus S'$  and the set of all atoms  $W_A$  such that  $A \notin \tilde{X}_{S'}$ . The value of this cut is at least

$$\sum_{A: A \subset \tilde{X}_{S'}} x_{s^*W_A,1}^{(t)} + \sum_{i \in S \setminus S'} x_{s^*i,2}^{(t)} = H(X_1, X_2, \dots, X_n) - \sum_{A: A \notin \tilde{X}_{S'}} x_{s^*W_A,1}^{(t)} + \sum_{i \in S \setminus S'} x_{s^*i,2}^{(t)}$$

By constraints (3.7), (3.8) from the second formulation and the given assignment, we have  $\sum_{A: A \notin \tilde{X}_{S'}} x_{s^*W_A,1}^{(t)} \leq \sum_{i \in S \setminus S'} x_{s^*i,2}^{(t)}$ . This implies that the value of any cut of this form at least  $H(X_1, X_2, \dots, X_n)$ . Therefore we can conclude that the minimum cut over all cuts separating  $P^*$  and  $Q^*$  is exactly  $H(X_1, X_2, \dots, X_n)$ . Because we have  $\sum_{A \in \mathcal{A}} x_{s^*W_A,1}^{(t)} = H(X_1, X_2, \dots, X_n)$ , the flows on edge  $(P^*, W_A)$  should be equal to the capacity, our assignment is a valid solution.  $\square$

From the second formulation, we can find a corresponding first formulation with equal cost  $f_{opt2}$ , then  $f_{opt1} \leq f_{opt2}$ . Hence,  $f_{opt1} = f_{opt2}$ .

As we claimed earlier, the second formulation will be useful when we compare the cost gap between coded case and subset case, we will use this augmented graph  $G^* = G_2^*$  in the rest of the paper.

### 3.1.3 Solution Explanation and Construction

Assume that we solve the above problem and obtain the values of all the atoms  $\mu^*(A)$ ,  $A \in \mathcal{A}$ . These will in general be fractional. We now outline the algorithm that decides the content of each source node. We use the assumption that the original source can be represented as a collection of independent equal-entropy random variables  $\{OS_i\}_{i=1}^Q$ , for large enough  $Q$  at this point. Suppose that  $H(OS_1) = \beta$ . In turn, we can conclude that there exist integers  $\alpha_A, \forall A \in \mathcal{A}$ , such that  $\alpha_A \times \beta = \mu^*(A)$ ,  $\forall A \in \mathcal{A}$  and that  $\sum_{A \in \mathcal{A}} \alpha_A = Q$ . Consider an ordering of the atoms, denoted as  $A_1, A_2, \dots, A_{2^n-1}$ . The atom sources can then be assigned as follows: For each  $A_i$ , assign  $W_{A_i} = (OS_{\sum_{j<i} \alpha_{A_j}+1}, OS_{\sum_{j<i} \alpha_{A_j}+2}, \dots, OS_{\sum_{j \leq i} \alpha_{A_j}})$ . It is clear that the resultant atom sources are independent and that  $H(W_A) = \mu^*(A)$ ,  $\forall A \in \mathcal{A}$ . We then assign  $X_i = (W_A : A \subset \tilde{X}_i)$ .

The assumption on the original source is essentially equivalent to saying that a large file can be broken into arbitrarily small pieces. To see this assume that each edge in the network has a capacity of 1000 bits/sec. At this time-scale, suppose that we treat each edge as unit-capacity, then a source of entropy one bit, can be considered to be a source of entropy  $10^{-3}$  at this time-scale. Therefore, if a given file can be broken into arbitrarily small pieces, then one can decompose it into pieces of arbitrarily small entropy.

### 3.2 Coded Source Network

Given the same network, if we allow coded information stored at the sources, using the augmented graph  $G^*$  by the second problem formulation, the storage at the sources can be viewed as the transmission along the edges connecting the virtual source and real sources. Then the problem becomes the standard minimum cost multicast with network coding problem (CODED-MIN-COST) where the variables are only the flows  $z_{ij}$  and  $x_{ij}^{(t)}$ .

$$\begin{aligned} & \text{minimize } \sum_{(i,j) \in E} f_{ij} z_{ij} + \sum_{i \in S} d_i z_{s^*i} \\ & \text{subject to } 0 \leq x_{ij}^{(t)} \leq z_{ij} \leq c_{ij}^*, (i,j) \in E^*, t \in T \end{aligned}$$

$$\sum_{\{j|(i,j) \in E^*\}} x_{ij}^{(t)} - \sum_{\{j|(j,i) \in E^*\}} x_{ji}^{(t)} = \sigma_i^{(t)}, i \in V^*, t \in T$$

where

$$\sigma_i^{(t)} = \begin{cases} H(X_1, X_2, \dots, X_n) & \text{if } i = s^* \\ -H(X_1, X_2, \dots, X_n) & \text{if } i = t \\ 0 & \text{otherwise} \end{cases}$$

Assume we have the solution for CODED-MIN-COST, we can use the random coding scheme or the deterministic coding scheme introduced by [12][13] to reconstruct the sources and information flow of each edge. We can also use the algorithm in [14] to find the suboptimal value with integral values on the edges.

## CHAPTER 4. COST COMPARISON BETWEEN CODED CASE AND SUBSET CASE

For given instances of the problem, we can certainly compute the cost gap by solving the corresponding optimization problems SUBSET-MIN-COST and CODED-MIN-COST presented in the previous section. Because the subset case is a special case of the coded case, we define the cost gap as the difference between the optimums of the subset case and the coded case. In this section, we first formulate an optimization problem similar to SUBSET-MIN-COST. The main difference is that we consider the source node can contain any arbitrary functions of the pieces of the original source. Accordingly, we require the atoms to satisfy the information inequalities [10] that consist of Shannon type inequalities and non-Shannon type inequalities when  $n \geq 4$  [15]. In reference [16], it was shown that there are infinitely many non-Shannon type inequalities. Hence, it is impossible to list all the information inequalities when the source number exceeds 4. However, if we remove the non-Shannon type inequalities from the constraints, the optimal value of coded case will not increase. In turn, this means that the gap computed by comparing these optimal values will still be a valid upper bound for the gap between the subset case and coded case.

Following this we can find an upper bound on the cost gap as the solution to another optimization problem. In the general case, of  $n$  sources, even this optimization has constraints that are exponential in  $n$ . However, this formulation still has advantages. In particular, we are able to provide a greedy algorithm for find near-optimal solutions for it. Moreover, we are able to prove that this greedy algorithm allows us to determine an upper bound in the case of three sources, which can be shown to be tight, i.e., there exist instances such that the cost gap is met with equality.

## 4.1 General Case

We now present the problem formulation for ATOM-CODED-MIN-COST. As done previously, we augment the graph  $G$  with a virtual super source  $s^*$  and introduce infinite capacity edges from  $s^*$  to each  $i \in S$ .

$$\begin{aligned} & \text{minimize } \sum_{(i,j) \in E} f_{ij} z_{ij} + \sum_{i \in S} d_i z_{s^*i} \\ & \text{subject to } 0 \leq x_{ij}^{(t)} \leq z_{ij} \leq c_{ij}^*, \forall (i,j) \in E^*, t \in T \end{aligned}$$

$$\sum_{\{j|(i,j) \in E^*\}} x_{ij}^{(t)} - \sum_{\{j|(j,i) \in E^*\}} x_{ji}^{(t)} = \sigma_i^{(t)}, \forall i \in V^*, t \in T \quad (4.1)$$

$$x_{s^*i}^{(t)} \geq R_i^{(t)}, \forall i \in S, t \in T \quad (4.2)$$

$$R^{(t)} \in \mathcal{R}_{S\mathcal{W}}, \forall t \in T \quad (4.3)$$

$$H(X_i | X_{S \setminus \{i\}}) \geq 0, \forall i \in S \quad (4.4)$$

$$I(X_i; X_j | X_K) \geq 0, \forall i \in S, j \in S, i \neq j, K \subseteq S \setminus \{i, j\} \quad (4.5)$$

$$z_{s^*i} = H(X_i), \forall i \in S \quad (4.6)$$

$$H(X_1, X_2, \dots, X_n) = \sum_{A \in \mathcal{A}} \mu^*(A) \quad (4.7)$$

where

$$\sigma_i^{(t)} = \begin{cases} H(X_1, X_2, \dots, X_n) & \text{if } i = s^* \\ -H(X_1, X_2, \dots, X_n) & \text{if } i = t \\ 0 & \text{otherwise} \end{cases}$$

The formulation is the same as SUBSET-MIN-COST except that we remove (3.9), and add constraints (4.4) and (4.5) that are elemental inequalities, which guarantee that all Shannon type inequalities are satisfied [10]. The elemental inequalities can be represented in the form of atoms:

$$\begin{aligned} H(X_i | X_{S \setminus \{i\}}) &= \mu^*(A), A \notin \tilde{X}_{S \setminus \{i\}} \\ I(X_i; X_j | X_K) &= \sum_{A \in \mathcal{A}: A \subset \tilde{X}_i, A \subset \tilde{X}_j, A \not\subset \tilde{X}_K} \mu^*(A) \end{aligned}$$

where  $K \subseteq S \setminus \{i, j\}$ .

Now, suppose that we know the optimal value of the above optimization problem, i.e., the flows  $x_{ij,1}^{(t)}, z_{ij,1}^{(t)}, t \in T, (i, j) \in E^*$ , the measure of the atoms  $\mu^*(A)_1, \forall A \in \mathcal{A}$ , and the corresponding conditional entropies  $H^1(X_U|X_{S \setminus U}), \forall U \subseteq S$ . If we can construct a feasible solution for SUBSET-MIN-COST such that the flows over  $E^*$  are the same as  $x_{ij,1}^{(t)}$  (and  $z_{ij,1}^{(t)}$ ),  $t \in T, (i, j) \in E$ , then we can arrive at an upper bound for the gap. This is done below. Let  $\mu^*(A), \forall A \in \mathcal{A}$  denote the variables for the atom measures for the subset case. We have a gap LP,

$$\min \sum_{A \in \mathcal{A}} \left( \sum_{\{i \in S: AC \tilde{X}_i\}} d_i \right) \mu^*(A) - \sum_{A \in \mathcal{A}} \left( \sum_{\{i \in S: AC \tilde{X}_i\}} d_i \right) \mu^*(A)_1$$

subject to

$$\sum_{A: A \not\subseteq \tilde{X}_{S \setminus U}} \mu^*(A) \leq H^1(X_U|X_{S \setminus U}), \forall U \subset S \quad (4.8)$$

$$\mu^*(A) \geq 0, \forall A \in \mathcal{A}$$

$$\sum_{A: A \in \mathcal{A}} \mu^*(A) = H(X_1, X_2, \dots, X_n)$$

where  $H^1(X_U|X_{S \setminus U}) = \sum_{A: A \not\subseteq \tilde{X}_{S \setminus U}} \mu^*(A)_1, \forall U \subset S$ . In the SUBSET-MIN-COST, we assign  $x_{ij}^{(t)} = x_{ij,1}^{(t)}, (i, j) \in E^*$ ,  $z_{ij}^{(t)} = z_{ij,1}^{(t)}, (i, j) \in E$  and  $z_{s^*i} = \sum_{A: AC \tilde{X}_i} \mu^*(A), \forall i \in S$ . To see that this is feasible, note that

$$\begin{aligned} z_{s^*i} &= \sum_{A: AC \tilde{X}_i} \mu^*(A) = H(X_i) \\ &= H(X_1, X_2, \dots, X_n) - H(X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n | X_i) \\ &\stackrel{(a)}{\geq} H(X_1, X_2, \dots, X_n) - H^1(X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n | X_i) \\ &= H^1(X_i) \\ &= z_{s^*i,1} \\ &\geq x_{s^*i,1}^{(t)} \\ &= x_{s^*i}^{(t)} \end{aligned}$$

Then constraint (3.6) is satisfied.

$$\sum_{i:i \in U} x_{s^*i}^{(t)} = \sum_{i:i \in U} x_{s^*i,1}^{(t)} \geq H^1(X_U|X_{S \setminus U}) \stackrel{(b)}{\geq} H(X_U|X_{S \setminus U})$$

where  $H(X_U|X_{S\setminus U}) = \sum_{A:A \not\subseteq \tilde{X}_{S\setminus U}} \mu^*(A), \forall U \subset S$ . Then constraints (3.7) and (3.8) are satisfied.

Both (a) and (b) come from constraint (4.8). The difference in the costs is only due to the different storage costs, since the flow costs are exactly the same.

The Lagrangian for GAP problem is

$$\begin{aligned} L(\mu, \lambda, \nu, \alpha) &= \sum_{A \in \mathcal{A}} \left( \sum_{\{i \in S: A \subset \tilde{X}_i\}} d_i \right) \mu^*(A) - \sum_{A \in \mathcal{A}} \sum_{\{i \in S: A \subset \tilde{X}_i\}} d_i \mu^*(A)_1 \\ &+ \sum_{U \subset S} \lambda_U \left( \sum_{A: A \not\subseteq \tilde{X}_{S \setminus U}} \mu^*(A) - H^1(X_U|X_{S \setminus U}) \right) \\ &- \sum_{A \in \mathcal{A}} \nu_A \mu^*(A) \\ &+ \alpha \left( \sum_{A: A \in \mathcal{A}} \mu^*(A) - H(X_1, \dots, X_n) \right) \end{aligned}$$

The KKT condition is, for the optimal solution  $(\mu^*(A), \lambda, \nu, \alpha)$ :

for  $A_j \in \mathcal{A}$ , if there exists  $i$ , such that  $A_j \not\subseteq \tilde{X}_i$ , then  $\mu^*(A_j)$  should satisfy

$$\frac{\partial L(\mu, \lambda, \nu, \alpha)}{\partial \mu^*(A_j)} = \sum_{\{i \in S: A_j \subset \tilde{X}_i\}} d_i + \sum_{U: A_j \subset \tilde{X}_U, U \subset S} \lambda_U - \nu_{A_j} + \alpha,$$

for  $A_j \in \mathcal{A}$ , if for all  $i$ ,  $A_j \subset \tilde{X}_i$ , then  $\mu^*(A_j)$  should satisfy

$$\frac{\partial L(\mu, \lambda, \nu, \alpha)}{\partial \mu^*(A_j)} = \sum_{\{i \in S: A_j \subset \tilde{X}_i\}} d_i - \nu_{A_j} + \alpha,$$

and

$$\lambda_U \geq 0, \forall U \subset S$$

$$\nu_A \geq 0, \forall A \in \mathcal{A}$$

$$\lambda_U \left( \sum_{A: A \not\subseteq \tilde{X}_{S \setminus U}} \mu^*(A) - H^1(X_U | X_{S \setminus U}) \right) = 0, \forall U \subset S$$

$$\nu_A \mu^*(A) = 0, \forall A \in \mathcal{A}$$

$$\sum_{A: A \not\subseteq \tilde{X}_{S \setminus U}} \mu^*(A) - H^1(X_U | X_{S \setminus U}) \leq 0, \forall U \subset S$$

$$\mu^*(A) \geq 0, \forall A \in \mathcal{A}$$

$$\sum_{A: A \in \mathcal{A}} \mu^*(A) = H(X_1, \dots, X_n)$$

#### 4.1.1 Greedy Algorithm

We present a greedy algorithm for the gap LP that returns a feasible, near-optimal solution, and hence serves as an upper bound to the gap. The main idea is to start by saturating atom values with the low costs, while still remaining feasible, e.g., suppose that source 1, has the smallest cost. Then, the atom  $\tilde{X}_1 \cap_{k \in \mathcal{N}_S \setminus \{1\}} \tilde{X}_k^c$  has the least cost among all the atoms, and therefore we assign it the maximum value possible, i.e.,  $H^1(X_1 | X_{S \setminus \{1\}})$ . Further assignments are made similarly in a greedy fashion. More precisely we follow the steps given below.

1. Initialize  $\mu^*(A) = 0, \forall A \in \mathcal{A}$ . Label all atoms as “unassigned”.
2. If all atoms have been assigned, STOP. Otherwise, let  $A_{\min}$  denote the atom with minimum cost that is still unassigned.
  - Set  $\mu^*(A_{\min}) \geq 0$  as large as possible so that the sum of the values of all assigned atoms does not violate any constraint in (4.8).
  - Check to see whether  $\sum_{A \in \mathcal{A}} \mu^*(A) > H(X_1, X_2, \dots, X_n)$ . If YES, then reduce the value of  $\mu^*(A_{\min})$ , so that  $\sum_{A \in \mathcal{A}} \mu^*(A) = H(X_1, X_2, \dots, X_n)$  and STOP. If NO, then label  $A_{\min}$  as “assigned”.

3. Go to step 2.

It is clear that the this algorithm returns a feasible set of atom values, since we maintain feasibility at all times and enforce the sum of the atom values to be  $H(X_1, X_2, \dots, X_n)$ .

The greedy algorithm, though suboptimal, does give the exact gap in many cases that we tested. Moreover, as discussed next, the greedy approach allows us to arrive at a closed form expression for the an upper bound on the gap in the case of three sources.

## 4.2 Three Sources Case

The case of three sources is special because, (i) Shannon type inequalities suffice to describe the entropic region, i.e., non-Shannon type inequalities do not exist for three random variables. This implies that we can find three random variables using the atom measures solution of ATOM-CODED-MIN-COST. (ii) Moreover, there is at most one atom,  $\mu^*(\tilde{X}_1 \cap \tilde{X}_2 \cap \tilde{X}_3)$  that can be negative. This makes the analysis easier since the greedy algorithm proposed above can be applied to obtain the required bound. Let the atoms be denoted by the variables shown in Figure 4.1.

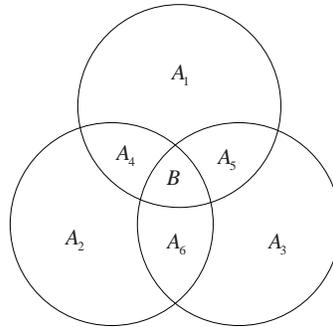


Figure 4.1 Atom pattern when there are three sources

**Claim 3** Consider random variables  $X_1, X_2$  and  $X_3$  with  $H(X_1, X_2, X_3) = h$ . Then,  $\mu^*(B) \geq -\frac{h}{2}$ .

*Proof.* Let  $b = \mu^*(B)$ , and  $a_i = \mu^*(A_i), i = 1, \dots, 6$ . The elemental information inequalities are given by

$$a_i \geq 0, i = 1, \dots, 6$$

$$a_i + b \geq 0, i = 4, 5, 6.$$

We also have joint entropy equality

$$\left( \sum_{i=1, \dots, 6} a_i \right) + b = h.$$

Assume that  $b < -\frac{h}{2}$ . Then,

$$a_i + b \geq 0 \Rightarrow a_i \geq -b > \frac{h}{2}, i = 4, 5, 6 \Rightarrow a_4 + a_5 > h.$$

Next,

$$h = a_1 + a_2 + a_3 + a_4 + a_5 + a_6 + b \geq a_1 + a_2 + a_3 + a_4 + a_5 > a_1 + a_2 + a_3 + h.$$

This implies that  $a_1 + a_2 + a_3 < 0$ , which is a contradiction, since  $a_i \geq 0, i = 1, \dots, 6$ .  $\square$

Using this we can obtain the following lemma

**Lemma 1** *Suppose that we have three source nodes. Let the joint entropy of the original source be  $h$  and let  $f_{opt2}$  represent the optimal value of SUBSET-MIN-COST and  $f_{opt1}$ , the optimal value of CODED-MIN-COST. Then,  $f_{opt2} - f_{opt1} \leq (\min_{i \in S}(d_i))h/2$ .*

*Proof.* Without loss of generality, assume that  $\min_{i \in S}(d_i) = d_1$ . Suppose that in the optimal solution for ATOM-CODED-MIN-COST,  $\mu^*(B) = b^1 \leq 0$ . As in the greedy algorithm above, we construct a feasible solution for SUBSET-MIN-COST by keeping the flow values the same, but changing the atom values suitably. Let  $a_i^2, i = 1, \dots, 6, b^2$  denote the atom values for the subset case. Consider the following assignment,

$$a_i^2 = a_i^1, i = 1, \dots, 5$$

$$a_6^2 = a_6^1 - |b^1|, \text{ and}$$

$$b^2 = 0.$$

$H^1(X_i)$  denotes the entropy solution of ATOM-CODED-MIN-COST for source  $i, i \in S$ .  $H^2(X_i)$  denotes the entropy solution of SUBSET-MIN-COST for source  $i$ . Assume  $b^1$  is negative, the other atoms values are  $a_i^1, i = 1, \dots, 6$ . We have the solution  $x_{ij,1}^{(t)}, \forall t \in T, (i, j) \in E^*$  and  $H^1(X_i), i = 1, 2, 3$ . Now set  $H^2(X_1) = H^1(X_1) + |b^1|$ ,  $H^2(X_2) = H^1(X_2)$ ,  $H^2(X_3) = H^1(X_3)$ . This can be realized by letting  $a_i^2 = a_i^1, i = 1, 2, 3, 4, 5$ ;  $b^2 = b^1 + |b^1| = 0$ ;  $a_6^2 = a_6^1 - |b^1|$ . Then  $a_i^2 \geq 0, i = 1, \dots, 6$ ,  $b^2 \geq 0$ ,  $H^2(X_1, X_2, X_3) = H^1(X_1, X_2, X_3)$ , we can find a subset source distribution. This is shown pictorially in Figure 4.2.

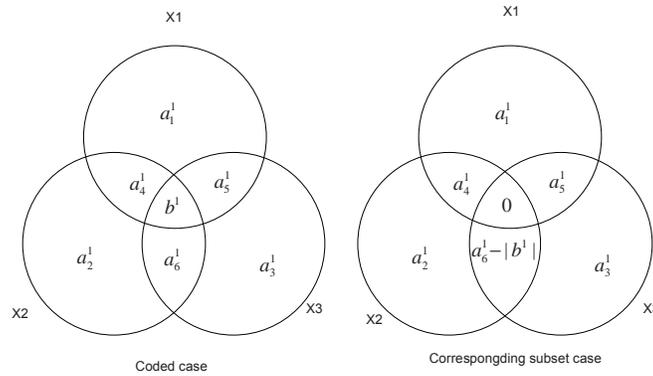


Figure 4.2 The transforming scheme from coded case to subset case,  $b^1$  is negative

We can check constraint (4.8) to see that the solution is feasible for the gap LP for three sources. We can also check the constants (3.6), (3.7) and (3.8) in the SUBSET-MIN-COST to see the feasibility

$$x_{s^*1,1}^{(t)} \geq R_1^{(t)} \geq H^1(X_1|X_2, X_3) = H^2(X_1|X_2, X_3)$$

$$x_{s^*2,1}^{(t)} \geq R_2^{(t)} \geq H^1(X_2|X_1, X_3) = H^2(X_2|X_1, X_3)$$

$$x_{s^*3,1}^{(t)} \geq R_3^{(t)} \geq H^1(X_3|X_2, X_1) = H^2(X_3|X_2, X_1)$$

$$x_{s^*1,1}^{(t)} + x_{s^*2,1}^{(t)} \geq R_1^{(t)} + R_2^{(t)} \geq H^1(X_1, X_2|X_3) = H^2(X_1, X_2|X_3)$$

$$x_{s^*1,1}^{(t)} + x_{s^*3,1}^{(t)} \geq R_1^{(t)} + R_3^{(t)} \geq H^1(X_1, X_3|X_2) = H^2(X_1, X_3|X_2)$$

$$x_{s^*2,1}^{(t)} + x_{s^*3,1}^{(t)} \geq R_2^{(t)} + R_3^{(t)} \geq H^1(X_2, X_3|X_1) > H^2(X_2, X_3|X_1)$$

$$x_{s^*1,1}^{(t)} + x_{s^*2,1}^{(t)} + x_{s^*3,1}^{(t)} \geq R_1^{(t)} + R_2^{(t)} + R_3^{(t)} \geq H^1(X_1, X_2, X_3) = H^2(X_1, X_2, X_3)$$

$$x_{s^*1,1}^{(t)} \leq z_{s^*1,1} = H^1(X_1) < H^2(X_1)$$

$$x_{s^*2,1}^{(t)} \leq z_{s^*2,1} = H^1(X_2) = H^2(X_2)$$

$$x_{s^*3,1}^{(t)} \leq z_{s^*3,1} = H^1(X_3) = H^2(X_3)$$

The transforming process is the standard procedure of greedy algorithm. We can also check the KKT condition for gap LP to see the optimality. Assume the cost  $d_1 \leq d_2 \leq d_3$ . According to the KKT condition in the general case, let  $\lambda_1 = \lambda_2 = \lambda_3 = d_1$ ,  $\lambda_{12} = d_3 - d_1$ ,  $\lambda_{13} = d_2 - d_1$ ,  $\lambda_{23} = 0$ ,  $\nu_{123} = d_1$ ,  $\alpha = -(d_2 + d_3)$ , and all other dual variables be 0, then this set of dual variables and primal variables that we give above will satisfy all the KKT conditions. Hence, we have a optimal solution.  $x_{i,j,1}^{(t)}$ ,  $(i, j) \in E^*$  are feasible for the subset problem. The flows do not change over transforming the coded case to the subset case. The only cost increased is  $(d_1) \times (|b|) \leq (\min_{i \in S}(d_i))h/2$ .  $\square$

In Chapter 5, we will show an instance of a network where this upper bound is tight.

Finally we note that when there are only two source nodes, there is no cost difference between the subset case and the coded case, since for two random variables, all atom measures have to be nonnegative. We state this as a lemma below.

**Lemma 2** *Suppose that we have two source nodes. Let  $f_{opt2}$  represent the optimal value of SUBSET-MIN-COST and  $f_{opt1}$ , the optimal value of CODED-MIN-COST. Then,  $f_{opt2} = f_{opt1}$ .*

## CHAPTER 5. RESULTS

In this section we present an example of a network with three sources where our upper bound derived in Chapter 4.2 is tight. We also performed several experiments with randomly generated graphs. The primary motivation was to study whether the difference in cost between the subset sources case and the coded case occurs very frequently. We finally present the results for gap LP and greedy algorithm to see how accurately the gap can be computed.

### 5.1 Results on a Deterministic Network

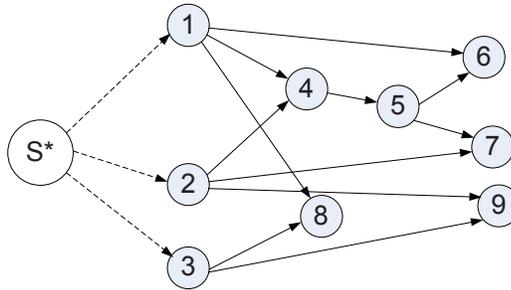


Figure 5.1 A deterministic network

Consider the network in Figure 5.1 with three sources nodes, 1, 2 and 3 and four terminal nodes, 6, 7, 8, and 9. The entropy of the original source =  $H(X_1, X_2, X_3) = 2$  and all edges are unit-capacity. The costs are such that  $f_{ij} = 1, \forall (i, j) \in E$  and  $d_1 = d_2 = 2, d_3 = 1$ .

The optimal cost in the subset sources case is 17. The corresponding atom values are listed in the Table 5.1. In this case we have  $H(X_1) = 1.22, H(X_2) = 1.36$  and  $H(X_3) = 1.42$ . We can decide the source content using the construction steps we introduced.

Table 5.1 Atom values when subset constraints are enforced

Atom	$X_1^c X_2^c X_3^c$	$X_1 X_2^c X_3^c$	$X_1^c X_2 X_3^c$	$X_1 X_2 X_3^c$	$X_1^c X_2^c X_3$	$X_1 X_2^c X_3$	$X_1^c X_2 X_3$	$X_1 X_2 X_3$
$\mu^*$	0	0	0	0.5809	0	0.6367	0.7824	0

In the coded sources case, the optimal value is 16, with  $H(X_1) = H(X_2) = H(X_3) = 1$ . The distributions for the sources are:  $X_1$  and  $X_2$  are independent with a distribution  $P(X_i = 0) = P(X_i = 1) = 0.5, i = 1, 2$ .  $X_3 = (X_1 + X_2) \bmod 2$ . Coding exists among  $X_1, X_2$  and  $X_3$ . Note that in this case the gap between the optimal values is precisely  $= \frac{h}{2} \times 1 = \frac{2}{2} \times 1 = 1$ , i.e., the upper bound derived in the previous section is met with equality.

## 5.2 Results on Random Networks

We generated several directed graphs at random with  $|V| = 87$ ,  $|E| = 322$ . The linear cost of each edge was fixed to an integer in  $\{1, 2, 3, 4, 5, 6\}$  or a large number such as 29 and 31. We ran 5000 experiments with fixed parameters  $(|S|, |T|, h)$ , where  $|S|$  denotes the number of source nodes,  $|T|$  denotes the number of terminal nodes and  $h$  denotes the entropy of the original source. The locations of the source and terminal nodes were chosen randomly. The capacity of each edge was chosen at random from the set  $\{1, 2, 3, 4, 5\}$ . There were no capacity constraints on the source nodes. In many cases it turned out that the network did not have enough capacity to support the recovery at the terminals. These instances were discarded. Notice whenever there is coded case solution, we are able to find a corresponding subset solution.

The results are shown in Table 5.2. The ‘‘Equal’’ row corresponds to the number of instances when both the coded and subset case have the same cost, and ‘‘Non-equal’’ corresponds to the number of instances where the coded case has a lower cost. Note that in most cases, the subset case and the coded case have the exact same cost. We also evaluated the gap LP and the greedy algorithm proposed in Section 4.1.1 to evaluate the cost gap. Note that the gap LP is only an upper bound since it is derived assuming that the flow patterns do not change between the coded and the subset case. When  $(|S|, |T|, h) = (4, 3, 4)$ , we ran 5000 experiments,

Table 5.2 Comparisons of two schemes in 5000 random directed graphs

$( S ,  T , h)$	(3, 3, 3)	(4, 4, 4)	(5, 5, 5)	(4, 5, 5)	(5, 4, 5)	(4, 4, 5)
<i>Equal</i>	3893	2855	1609	1577	2025	1954
<i>Nonequal</i>	1	3	10	9	6	8

among which 3269 instances could support both the coded and the subset case. Out of these, there were 481 instances where the upper bound determined by the gap LP was not tight. In addition, there were 33 instances where the greedy algorithm failed to solve the gap LP exactly.

## CHAPTER 6. SUMMARY AND DISCUSSION

In this work, we considered network coding based content distribution, under the assumption that the content can be considered as a collection of independent equal entropy sources. e.g., a large file that can be broken into small pieces. Given a network with a specified set of source nodes, we examined two cases. In the subset sources case, the source nodes are constrained to only contain subsets of the pieces of the content, whereas in the coded sources case, the source nodes can contain arbitrary functions of the pieces. The cost of a solution is defined as the sum of the storage cost and the cost of the flows required to support the multicast. We provided succinct formulations of the corresponding optimization problems by using the properties of information measures. In particular, we showed that when there are two source nodes, there is no loss in considering subset sources. For three source nodes, we derived a tight upper bound on the cost gap between the two cases. A greedy algorithm for estimating the cost gap for a given instance was provided. Finally, we also provided algorithms for determining the content of the source nodes.

Our results indicate that when the number of source nodes is small, in many cases constraining the source nodes to only contain subsets of the content does not incur a loss.

**APPENDIX. A VARIATION OF THE INCLUSION-EXCLUSION  
FORMULA**

In this appendix, we show that for each  $A \in \mathcal{A}$ ,  $\mu(A)$  can be expressed as a linear combination of  $\mu(B)$ ,  $B \in \mathcal{B}$ . We first prove the following variation of the inclusive-exclusive formula

**Theorem 4** *For a set-additive function  $\mu$ ,*

$$\begin{aligned} \mu\left(\bigcap_{k=1}^n A_k - B\right) &= \sum_{1 \leq i \leq n} \mu(A_i - B) - \sum_{1 \leq i < j \leq n} \mu(A_i \cup A_j - B) \\ &+ \cdots + (-1)^{n+1} \mu(A_1 \cup A_2 \cup \cdots \cup A_n - B). \end{aligned} \tag{.1}$$

*Proof.* The theorem will be proved by induction on  $n$ . First when  $n = 1$ , the equality holds.

Assume the equality holds for some  $n \geq 1$ . Now consider

$$\begin{aligned}
& \mu\left(\bigcap_{k=1}^{n+1} A_k - B\right) \\
&= \mu\left(\left(\bigcap_{k=1}^n A_k\right) \cap A_{n+1} - B\right) \\
&= \mu\left(\bigcap_{k=1}^n A_k - B\right) + \mu(A_{n+1} - B) - \mu\left(\left(\bigcap_{k=1}^n A_k\right) \cup A_{n+1} - B\right) \\
&= \left\{ \sum_{1 \leq i \leq n} \mu(A_i - B) - \sum_{1 \leq i < j \leq n} \mu(A_i \cup A_j - B) + \cdots + (-1)^{n+1} \mu(A_1 \cup A_2 \cup \cdots \cup A_n - B) \right\} \\
&\quad + \mu(A_{n+1} - B) - \mu\left(\bigcap_{k=1}^n (A_k \cup A_{n+1}) - B\right) \\
&= \left\{ \sum_{1 \leq i \leq n} \mu(A_i - B) - \sum_{1 \leq i < j \leq n} \mu(A_i \cup A_j - B) + \cdots + (-1)^{n+1} \mu(A_1 \cup A_2 \cup \cdots \cup A_n - B) \right\} \\
&\quad + \mu(A_{n+1} - B) - \left\{ \sum_{1 \leq i \leq n} \mu(A_i \cup A_{n+1} - B) - \sum_{1 \leq i < j \leq n} \mu(A_i \cup A_j \cup A_{n+1} - B) \right. \\
&\quad \left. + \cdots + (-1)^{n+1} \mu(A_1 \cup A_2 \cup \cdots \cup A_n \cup A_{n+1} - B) \right\} \\
&= \sum_{1 \leq i \leq n+1} \mu(A_i - B) - \sum_{1 \leq i < j \leq n+1} \mu(A_i \cup A_j - B) + \cdots + (-1)^{n+2} \mu(A_1 \cup A_2 \cup \cdots \cup A_{n+1} - B).
\end{aligned} \tag{.2}$$

The  $n + 1$  case is also proved, then by induction, the theorem is proved.  $\square$

Now a nonempty atom  $\mathcal{F}_n$  has the form  $\bigcap_{i=1}^n Y_i$ , where  $Y_i$  is either  $\tilde{X}_i$  or  $\tilde{X}_i^c$  and there exists at least one  $i$  such that  $Y_i = \tilde{X}_i$ . Then we can write the atom as

$$\bigcap_{i: Y_i = \tilde{X}_i} \tilde{X}_i - \left( \bigcup_{j: Y_j = \tilde{X}_j^c} \tilde{X}_j \right)$$

Note that the intersection above is always nonempty. Then we see that for each  $A \in \mathcal{A}$ ,  $\mu(A)$  can be expressed as a linear combination of  $\mu(B)$ ,  $B \in \mathcal{B}$ .

## BIBLIOGRAPHY

- [1] C. Gkantsidis and P. Rodriguez, “Network Coding for Large Scale Content Distribution,” in *IEEE Infocom*, 2005.
- [2] A. Kirpal, P. Rodriguez, and E. Biersack, “Parallel-Access for Mirror Sites in the Internet,” in *IEEE Infocom*, 2000.
- [3] A. Jiang and J. Bruck, “Network File Storage with Graceful Performance Degradation,” *ACM Transactions on Storage*, vol. 1, no. 2, pp. 171–189, 2005.
- [4] S. Acedański, S. Deb, M. Médard, and R. Koetter, “How Good is Random Linear Coding Based Distributed Networked Storage,” in *NetCod*, 2005.
- [5] D. S. Lun, N. Ratnakar, M. Médard, R. Koetter, D. R. Karger, T. Ho, E. Ahmed, and F. Zhao, “Minimum-Cost Multicast over Coded Packet Networks,” *IEEE Trans. on Info. Th.*, vol. 52, pp. 2608–2623, June 2006.
- [6] A. Lee, M. Médard, K. Z. Haigh, S. Gowan, and P. Rubel, “Minimum-Cost Subgraphs for Joint Distributed Source and Network Coding,” in *the Third Workshop on Network Coding, Theory, and Applications*, Jan. 2007.
- [7] A. Ramamoorthy, “Minimum cost distributed source coding over a network,” in *IEEE Intl. Symposium on Info. Th.*, 2007, pp. 1761–1765.
- [8] A. Jiang, “Network Coding for Joint Storage and Transmission with Minimum Cost,” in *IEEE Intl. Symposium on Info. Th.*, 2006, pp. 1359–1363.
- [9] K. Bhattad, N. Ratnakar, R. Koetter, and K. Narayanan, “Minimal Network Coding for Multicast,” in *IEEE Intl. Symposium on Info. Th.*, 2005, pp. 1730–1734.

- [10] R. Yeung, *Information Theory and Network Coding*. Springer, 2008.
- [11] H. D. Sherali and G. Choi, “Recovery of primal solutions when using subgradient optimization methods to solve Lagrangian duals of linear programs,” *Oper. Res. Letters*, vol. 19, pp. 105–113, 1996.
- [12] T. Ho, M. Médard, J. Shi, M. Effros, and D. R. Karger, “On Randomized Network Coding,” in *41st Allerton Conference on Communication, Control, and Computing*, 2003.
- [13] S. Jaggi, P. Sanders, P. Chou, M. Effros, S. Egner, K. Jain, and L.M.G.M.Tolhuizen, “Polynomial Time Algorithms for Multicast Network Code Construction,” *IEEE Trans. on Info. Th.*, vol. 51, no. 6, pp. 1973–1982, 2005.
- [14] T. Cui and T. Ho, “Minimum Cost Integral Network Coding,” in *IEEE Intl. Symposium on Info. Th.*, 2007, pp. 2736–2740.
- [15] Z. Zhang and R. W. Yeung, “A non-Shannon-type Conditional Inequality of Information Quantities,” *IEEE Trans. on Info. Th.*, vol. 43, pp. 1982–1986, 1997.
- [16] F. Matúš, “Infinity Many Information Inequalities,” in *IEEE Intl. Symposium on Info. Th.*, 2007.