# CARDINALITY OF UPPER AVERAGE AND ITS APPLICATION TO NETWORK OPTIMIZATION

MATTHEW NORTON[*], ALEXANDER MAFUSALOV[†], STAN URYASEV[‡],
OLDERDRAFTS:JULY 2015, DECEMBER 2015 , AND THISDRAFT: JUNE 2017

**Abstract.** We propose a new characteristic for counting the number of large outcomes in a data set that are considered to be large w.r.t. some fixed threshold $x$. A popular characteristic used for this purpose is the Cardinality of Upper Tail (CUT), which counts the number of outcomes with magnitude larger than the threshold. We propose a similar characteristic called the Cardinality of Upper Average (CUA), defined as the number of largest data points which have average value equal to the threshold. CUA not only assesses the number of outcomes that are large, but also their overall magnitude. CUA also has superior mathematical properties: it is a continuous function of the threshold, its reciprocal is piece-wise linear w.r.t. threshold, and it is directly optimizable via convex and linear programming. This is in contrast to CUT, which does not asses the severity of large outcomes, is discontinuous as a function of threshold, and is such that direct optimization yields numerically difficult non-convex problems. We show that CUA can be used to formulate meaningful optimization problems containing counters of the largest components of a vector without introduction of binary variables, leading to large improvement in computation speeds. In particular, we apply the CUA concept to create new formulations of network optimization problems involving overloaded nodes or edges, where we aim to minimize the number of most burdened nodes or edges.

**Key words.** Cardinality of Upper Average, Buffered Probability of Exceedance, Network Optimization, Conditional Value-at-Risk, Mixed Integer Programming, Linear Programming

**AMS subject classifications.** 90C05, 90C35, 90C25, 90C11, 90C90

**1. Introduction.** When analyzing a set of $n$ data points, one often needs to count the number of outcomes that are large compared to some threshold $x \in \mathbb{R}$. For example, if the data points represent monetary losses, it is only natural that one would want to count the number of outcomes that are large compared to a threshold representing an 'acceptable' level of loss. A popular characteristic which acts as a counter of such large outcomes is the $x$-Cardinality of Upper Tail ($\text{CUT}_x$), which counts the number of data points with magnitude exceeding the threshold $x$. We introduce a similar characteristic called the $x$-Cardinality of Upper Average ($\text{CUA}_x$), defined as the number of largest data points which have average value equal to the threshold $x$. Thus, $\text{CUA}_x$ not only counts the number of data points with magnitude larger than the threshold, but also the largest outcomes with magnitude less than the threshold such that the average of these outcomes is equal to $x$.

Consider the following example to illustrate the conceptual difference between $\text{CUA}_x$ and $\text{CUT}_x$. Suppose that we have 10,000 pieces of gold. We want to make some statement about the number of heavy pieces of gold in this set because we know, for instance, that it takes 10 grams of gold to make one gold coin. So, 10 grams can be considered a natural reference point (or threshold) for evaluating the heaviness of gold pieces. After analyzing the dataset, we learn that there are only 3 pieces of gold with weight exceeding 10 grams. This means that $\text{CUT}_{10} = 3$ and, thus, that at least 3 gold coins can be made from the 3 heaviest pieces. Suppose, though, that

---

[*]Department of Industrial and Systems Engineering, University of Florida, Gainesville, FL (mdnorton@ufl.edu).

[†]Department of Industrial and Systems Engineering, University of Florida, Gainesville, FL (mafusalov@ufl.edu).

[‡]Director of Risk Management and Financial Engineering Lab, Department of Industrial and Systems Engineering, University of Florida, Gainesville, FL (uryasev@ufl.edu).

the weight of the largest piece of gold equals 1,000 grams. Then, the lower bound, 3 coins, is quite far away from the number of coins which can be made from the heaviest 3 pieces of gold. We can make at least 102 coins (i.e. 100 coins from the largest piece plus two more coins from the other two heaviest pieces). $\text{CUT}_x$ does not provide information about the heaviness of the largest pieces, only the lower bound for the individual pieces. Now, suppose we have calculated that $\text{CUA}_{10} = 207.3$. This actually means that you can make 207.3 coins from the heaviest 207.3 pieces. This is the case because the average weight of a gold piece in the selected subset equals 10 grams. Also, notice that $\text{CUA}_{10}$ tells us that the $208^{th}$ piece needs to be cut! $\text{CUA}_x$ is a continuous (nice mathematical) characteristic as a function of the threshold (10 grams), while $\text{CUT}_x$ is a discontinuous characteristic. $\text{CUA}_x$ provides information about the average heaviness of largest pieces.

This paper is focused on logistics applications, therefore, let us discuss again the same example in a logistics context. Suppose that we manage a network with 10,000 servers and assume that a server is "overloaded" if the length of the queue is greater than or equal to 10 jobs. We can ask the question: How many servers are overloaded? The answer is $\text{CUT}_{10} = 3$. However, it is unclear how significantly these servers are overloaded and how much additional resource would be required to serve the customers in these queues. An alternative characteristic is $\text{CUA}_{10} = 207.3$. This means that the average length of a queue over the 207.3 busiest servers is 10 jobs and the total number of unserved jobs in these servers is $207.3 \times 10 = 2,073$. Therefore, we understand immediately the total number of unserved clients in the 203.7 longest queues. This important information is present in $\text{CUA}_x$ and is not present in $\text{CUT}_x$.

You may ask, what is a better characteristic, $\text{CUA}_x$ or $\text{CUT}_x$? The simple answer is that these are different, complimentary characteristics. Therefore, while in some cases you may be interested in $\text{CUT}_x$ and in other cases $\text{CUA}_x$, it is probably a good idea to calculate both characteristics.

In the most basic sense, $\text{CUA}_x$ is similar to $\text{CUT}_x$, with both acting as counters of large outcomes relative to a threshold. Beyond this, they have important differences. First, as illustrated in the example, $\text{CUA}_x$ considers how far the outcome magnitudes are from the threshold. The $\text{CUT}_x$ characteristic does not and only considers if an outcome is larger than $x$, not considering how far beyond the threshold these magnitudes lie. Second, $\text{CUA}_x$ has superior mathematical properties. It is continuous w.r.t. the threshold parameter and a partial derivative can be taken w.r.t. threshold yielding information regarding the sensitivity of $\text{CUA}_x$ to threshold changes. It is also piecewise linear in its reciprocal. As we will show, this allows for efficient calculation of $\text{CUA}_x$ for all thresholds $x \in \mathbb{R}$.

In the context of optimization, the properties of $\text{CUA}_x$ yield substantial benefits, with direct optimization of $\text{CUA}_x$ reducing to convex, sometimes linear, programming. This means that we are able to efficiently minimize the number of tail outcomes, while taking into account the severity of these outcomes. Furthermore, the tail that is minimized is determined by the threshold, with the tail including the largest outcomes that average to the specified threshold. This is all in contrast to $\text{CUT}_x$, which does not account for the severity of the tail outcomes, is discontinuous w.r.t. the threshold parameter, and is such that direct optimization yields numerically challenging non-convex problems.

While being different from $\text{CUT}_x$ in many ways, $\text{CUA}_x$ is uniquely related. We show that $\text{CUA}_x$ is, in a certain sense, the minimal quasi-convex upper bound of $\text{CUT}_x$. In an optimization context, this could be useful, suggesting that $\text{CUA}_x$ minimization may yield efficient heuristic approaches to reducing $\text{CUT}_x$. This, though,

2

is not our focus and we emphasize that $\text{CUA}_x$ and its optimization are meaningful in their own right.

We apply $\text{CUA}_x$ to a specific optimization setting, providing convex and Linear Programming (LP) $\text{CUA}_x$ minimization formulations for solving network optimization problems, specifically addressing variants of the generalized assignment, capacity planning, and min cost network flow problems. In general, we consider network flow problems where product must flow through a network and we assume that intermediate nodes or arcs can become "overloaded" if too much flow is pushed through them. Therefore, it is desirable to minimize the number of overloaded nodes or arcs. We formulate this as a $\text{CUA}_x$ minimization problem, which accounts for the magnitude of loads put upon nodes/arcs in the overloaded state, as well as the nodes/arcs with loads less than, but most near to the overloaded state. This information may be meaningful, as it may be undesirable to have nodes/arcs that, while not yet overloaded, are close to the overloaded state. Furthermore, it may be undesirable to have nodes/arcs that are dramatically overloaded. Thus, we minimize the number of most overloaded nodes/arcs that have average load equal to $x$, where nodes/arcs with loads larger than $x$ are considered to be overloaded.

We compare these $\text{CUA}_x$ formulations with similar formulations that minimize $\text{CUT}_x$, directly minimizing the number of nodes/arcs in the overloaded state. With this problem involving binary variables to indicate the state of the node/arc, Mixed Integer Programming (MIP) must be used. Not only is this much more difficult to solve than the LP $\text{CUA}_x$ minimization, but we show that it may provide less appealing policies. By minimizing $\text{CUT}_x$, information about the magnitude of the loads is ignored and it only considers whether a node/arc is overloaded or not. This information, as already mentioned, can be meaningful. Thus, the $\text{CUA}_x$ formulation may be more appropriate and suggest more appealing policies.

The $\text{CUA}_x$ concept is actually a deterministic variant of so called Buffered Probability of Exceedance (bPOE). A detailed discussion of this concept, studied in [9, 12, 16, 17, 18], is beyond the scope of this paper. We maintain a deterministic setting, while bPOE is studied in a probabilistic, stochastic optimization setting. We do, though, include some background connecting $\text{CUA}_x$ and bPOE in Appendix A for the interested reader.

This paper is organized as follows. Section 2 discusses the task of analyzing the tail of a data distribution in a deterministic setting, which serves to set the stage for defining $\text{CUA}_x$. Section 3.1 defines $\text{CUA}_x$. We show that $\text{CUA}_x$ is efficient to calculate and give an example illustrating $\text{CUA}_x$, particularly as an upper bound of $\text{CUT}_x$. Additionally, we provide an efficient method of calculating $\text{CUA}_x$ which utilizes the piecewise linearity of its reciprocal. Section 3.2 provides relations between $\text{CUA}_x$ and $\text{CUT}_x$, showing that it is possible to simultaneously calculate $\text{CUA}_x$ and $\text{CUT}_x$. Section 4 shows the power of the $\text{CUA}_x$ concept when applied to an optimization setting. Specifically, we discuss classes of network optimization problems that are traditionally formulated as MIP's and show how application of $\text{CUA}_x$ leads to convex or linear programming reformulations of analogous problems.

## 2. Cardinality of Upper Tail and Related Quantities.

**2.1. Cardinality of Upper Tail and Upper Average.** Consider a Euclidean vector $\mathbf{y} = (y_1, \ldots, y_n)$ containing $n$ data points. It is often important in applications to know the number of components of $\mathbf{y}$ that exceed a particular threshold $x \in \mathbb{R}$.

We call this the Cardinality of the Upper Tail ($\mathrm{CUT}_x$), denoted as

$$\mathrm{CUT}_x(\mathbf{y}) = \eta_x(\mathbf{y}) = |\{y_i | y_i \geq x, i = 1, \ldots, n\}| .$$

This quantity, though, can be difficult to work with. For example, it is discontinuous w.r.t. the $x$ parameter. Additionally, this quantity does not provide information about the magnitude of the components above the threshold $x$. As we will show in later sections, these are undesirable properties in an optimization setting. We can also consider other quantities that provide useful information about the tail of the data distribution that are easy to work with. For example, one could consider the ordered weighted averaging aggregation operators of [19]. We can also consider a type of *tail average* called the $k$-Upper Average ($\mathrm{UA}_k$), which provides information about the magnitude of data points in the tail and upper bounds on $\mathrm{CUT}_x$. If we let $(y^{(1)}, y^{(2)}, \ldots, y^{(n)})$ represent a permutation of $(y_1, \ldots, y_n)$ with components listed in nondecreasing order, $y^{(1)} \leq y^{(2)} \leq \ldots \leq y^{(n)}$, we can denote this quantity as

$$(1) \qquad \mathrm{UA}_k(\mathbf{y}) = \frac{1}{k} \sum_{i=n-k+1}^{n} y^{(i)} .$$

**2.2. Generalized $\mathrm{UA}_k$.** Notice that (1) only applies to integer values of the $k$ parameter. This function, $\mathrm{UA}_k(\cdot)$, can also be defined in a more general manner for non-integer values of $k$. Popularized in the financial engineering literature under the name CVaR, paper [15] defines this quantity in a broad, probabilistic setting. Here, we maintain a deterministic setting, but emphasize that this is simply a special case of the general CVaR definition and formula. Thus, following directly from [15], we have that for any $k \in [1, n]$, letting $[\cdot]^+ = \max\{0, \cdot\}$,

$$(2) \qquad \mathrm{UA}_k(\mathbf{y}) = \min_\gamma \left\{ \gamma + \frac{1}{k} \sum_{i=1}^{n} [y_i - \gamma]^+ \right\} .$$

Though this formula may not seem intuitive, it can be clarified by noticing two facts. First, the optimal objective value of (2) is equal to (1) as long as $k \in \{1, \ldots, n\}$ is integer, meaning that for the integer case, the seemingly complex formula of (2) simply yields the average of the $k$ largest components. Second, for noninteger values of $k \in [1, n]$, this formula simply gives a weighted average of $\mathrm{UA}_{\lfloor k \rfloor}(\mathbf{y})$ and $\mathrm{UA}_{\lceil k \rceil}(\mathbf{y})$, where $\lfloor k \rfloor$ denotes the largest integer less than or equal to $k$ and $\lceil k \rceil$ denotes the smallest integer greater than or equal to $k$. Thus, intuitively, formulation (2) is still averaging the largest $k$ data points, but it is now a continuous function w.r.t. the $k$ parameter.

In addition to providing useful information about the magnitude of data points in the tail, $\mathrm{UA}_k(\cdot)$ provides an upper bound for $\mathrm{CUT}_x$. Specifically, we have the following relation, which follows intuitively from (1) and the definition of $\mathrm{CUT}_x$.

$$(3) \qquad \mathrm{UA}_k(\mathbf{y}) = x \implies \eta_x(\mathbf{y}) \leq k .$$

In this paper, we define $\mathrm{CUA}_x$, the inverse of (2). As we will show, $\mathrm{CUA}_x$ can be efficiently calculated and provides valuable information about the largest components of our data vector $\mathbf{y}$. Additionally, $\mathrm{CUA}_x$ can be efficiently optimized with convex and linear programming. We also show that $\mathrm{CUA}_x$ can be used to formulate optimization problems that are similar to, yet fundamentally different than $\mathrm{CUT}_x$ optimization problems. We show that $\mathrm{CUA}_x$ optimization can sometimes be more appropriate and may suggest more appealing optimal policies than the similar $\mathrm{CUT}_x$ minimization problems.

4

## 3. Cardinality of the Upper Average.

**3.1. Definition of** $\text{CUA}_x$**.** For a specified threshold $x$, $\text{CUA}_x$ calculates the value of $k \in [1, n]$ such that $\text{UA}_k(\mathbf{y}) = x$. In words, $\text{CUA}_x$ is equal to the number of largest components of the vector $\mathbf{y}$ such that the average of those components is equal to $x$. We now present Theorem 1, which provides two equivalent ways of defining $\text{CUA}_x$. While (5) could be viewed as the more intuitive of the two, we focus on the more tractable representation (4) throughout this paper. We then provide examples illustrating key differences between $\text{CUT}_x$ and $\text{CUA}_x$.

**Theorem 1:** *Consider a Euclidean vector $\mathbf{y} = (y_1, \ldots, y_n)$. $\text{CUA}_x$ of $\mathbf{y}$ at threshold $x \in \mathbb{R}$ is defined as,*

$$(4) \qquad \bar{\eta}_x(\mathbf{y}) = \min_{a \geq 0} \sum_{i=1}^{n} [a(y_i - x) + 1]^+ ,$$

*and can also be represented as,*

$$(5)$$
$$\bar{\eta}_x(\mathbf{y}) = \begin{cases} \max\{k | \frac{1}{k}(\sum_{i=1}^{\lfloor k \rfloor} y^{(n-i+1)} + (k - \lfloor k \rfloor)y^{(n-\lceil k \rceil+1)}) \geq x\} & \text{if } x \leq \max_i y_i , \\ 0 & \text{otherwise.} \end{cases}$$

**Proof:** Here we show that for a Euclidean vector $\mathbf{y} \in \mathbb{R}^n$, the equation (4) given for $\text{CUA}_x$ at threshold $x \in \mathbb{R}$ equals the value of $k$ such that $\text{UA}_k(\mathbf{y}) = x$. We must address four cases. The equivalence of (4) and (5) follows from these cases.

**Case 1:** Assume $x \in (\text{UA}_n(\mathbf{y}), \text{UA}_1(\mathbf{y}))$. This assumption ensures that there exists a value of $k \in [1, n]$ such that $\text{UA}_k(\mathbf{y}) = x$. We want to find the value of $k$ such that $\text{UA}_k(\mathbf{y}) = x$, so we write $\text{CUA}_x$ as

$$(6) \qquad \bar{\eta}_x(\mathbf{y}) = \{k | \text{UA}_k(\mathbf{y}) = x\} .$$

Notice now that $\text{UA}_k(\mathbf{y})$ is a strictly increasing function of $k$ on $k \in [\frac{n-m}{n}, n]$, where $m = |\{y_i | y_i = \max_i y_i\}|$ (i.e. $m$ equals the number of components of $\mathbf{y}$ that are equal to the largest component of $\mathbf{y}$). This follows from the known result (i.e. [15]) that CVaR is strictly increasing on an equivalent interval. Because of this, we see that equation (6) can be rewritten as the unique solution to

$$(7) \qquad \min\{k | \text{UA}_k(\mathbf{y}) \leq x\} .$$

Substituting equation (2) for $\text{UA}_k(\mathbf{y})$, equation (7) becomes

$$(8) \qquad \bar{\eta}_x(\mathbf{y}) = \min \left\{ k \ \middle| \ \min_{\gamma} \gamma + \frac{1}{k} \sum_{i=1}^{n} [y_i - \gamma]^+ \leq x \right\} .$$

This can then be simplified as

$$\bar{\eta}_x(\mathbf{y}) = \min_{k,\gamma} \quad k$$
$$(9)$$
$$s.t. \quad \gamma + \frac{1}{k} \sum_{i=1}^{n} [y_i - \gamma]^+ \leq x .$$

Explicitly enforcing the constraint $x - \gamma > 0$, we can further simplify (9) without changing the optimal solution to become

$$\bar{\eta}_x(\mathbf{y}) = \min_{k, x-\gamma>0} \quad k$$
(10)
$$s.t. \quad \frac{\sum_{i=1}^{n}[y_i - \gamma]^+}{x - \gamma} \leq k .$$

This, though, can be further simplified to become

(11)
$$\bar{\eta}_x(\mathbf{y}) = \min_{x-\gamma>0} \quad \frac{\sum_{i=1}^{n}[y_i - \gamma]^+}{x - \gamma} .$$

Finally, with the change of variable $a = \frac{1}{x-\gamma}$, (11) can be transformed into

(12)
$$\bar{\eta}_x(\mathbf{y}) = \min_{a \geq 0} \quad \sum_{i=1}^{n}[a(y_i - x) + 1]^+ .$$

**Case 2:** Assume $x = \max_j y_j$. With $x = \max_j y_j$, we show that $\mathrm{CUA}_x$ equals the number of components equal to $\max_j y_j$. For notational convenience let $\max_j y_j = y_{max}$. Also, assume there are $m$ components of $\mathbf{y}$ that are equal to $y_{max}$, i.e. $|\{y_i|y_i = y_{max}\}| = m$. Also, let $\hat{y} = \max\{y_i|y_i < y_{max}\}$, i.e. $\hat{y}$ equals the largest component of $\mathbf{y}$ that is less than $y_{max}$. [1]
Since $y_j - y_{max} \leq 0$ for any $j \in \{1, ...n\}$, we have that for any $a^* \geq \frac{-1}{\hat{y}-y_{max}}$

$$\min_{a \geq 0} \sum_{i=1}^{n}[a(y_i - y_{max}) + 1]^+ \geq \min_{a \geq 0} \sum_{y_i = y_{max}} [a(y_i - y_{max}) + 1]^+$$

$$= \sum_{i=1}^{n}[a^*(y_i - y_{max}) + 1]^+ = |\{y_i|y_i = y_{max}\}| = m .$$

To see this, notice that for any $y_i \leq \hat{y}$ and any $a^* \geq \frac{-1}{\hat{y}-y_{max}}$ we get

$$[a^*(y_i - y_{max}) + 1]^+ \leq [a^*(\hat{y} - y_{max}) + 1]^+ = [-1 + 1]^+ = 0 .$$

Furthermore, for any $y_j = y_{max} = x$, we have that

$$[a(y_i - y_{max}) + 1]^+ = [a(0) + 1]^+ = 1 .$$

.
**Case 3:** Assume $x > \max_i y_i$. We need to show that equation (4) equals the value of $k$ such that $\mathrm{UA}_k(\mathbf{y}) = x$. Since $x > \max_i y_i$ we show that $\mathrm{CUA}_x$ equals 0.
Let us denote $\max_j y_j = y_{max}$. Since $x > y_{max}$, then for any $i \in \{1, ..., n\}$ we have $y_i - x < 0$. This implies that for any $a^* \geq \frac{-1}{y_{max}-x}$, $a^*(y_i - x) \leq -1$ for all $i$, therefore,

$$\sum_{i=1}^{n}[a^*(y_i - x) + 1]^+ = 0 = \min_{a \geq 0} \sum_{i=1}^{n}[a(y_i - x) + 1]^+ .$$

---
[1] If all components of $\mathbf{y}$ are equal, then apply Case 4. In this case, although $\hat{y}$ does not exist, $\mathrm{UA}_n(\mathbf{y}) = \max_j y_j$ and Case 4 can be applied

**Case 4:** Assume $x \leq \mathrm{UA}_n(\mathbf{y})$. We need to show that equation (4) equals the value of $k$ such that $\mathrm{UA}_k(\mathbf{y}) = x$. Since $x \leq \mathrm{UA}_n(\mathbf{y})$, we show that $\mathrm{CUA}_x$ equals $n$. Since $x \leq \mathrm{UA}_n(\mathbf{y})$, we have $0 \leq \mathrm{UA}_n(\mathbf{y}) - x$. This implies that for any $a \geq 0$

$$\sum_{i=1}^{n}[a(y_i - x) + 1]^+ \geq \sum_{i=1}^{n}[a(y_i - x) + 1] \geq n[a(\mathrm{UA}_n(\mathbf{y}) - x) + 1] \geq n \ .$$

This result implies that $\min_{a \geq 0} \sum_{i=1}^{n}[a(y_i - x) + 1]^+ = n$ (attained at $a = 0$).

$\square$

Initially, it may be surprising that this particular formula calculates the number of biggest components of the vector $\mathbf{y}$ such that the average of those components is equal to $x$. In short, this formula is derived from equation (2), partially explaining its form as the unique minimal value to the minimization problem. Note that the proof of Theorem 1 also shows how $\mathrm{CUA}_x$ is defined for thresholds $x \notin (\mathrm{UA}_n(\mathbf{y}), \mathrm{UA}_1(\mathbf{y}))$ where $\mathrm{UA}_n(\mathbf{y}) = \frac{1}{n}\sum_{i=1}^{n} y_i$ and $\mathrm{UA}_1(\mathbf{y}) = \max_i y_i$. Furthermore, the definition of $\mathrm{CUA}_x$ for these extreme cases is motivated by its connection to bPOE, a new concept recently studied in [9, 12, 16, 17, 18]. Appendix A includes a brief discussion of this connection showing that $\mathrm{CUA}_x$ can be viewed as a deterministic variant of Upper bPOE [12, 9]. A detailed discussion of this, though, is beyond the scope of this paper.

To begin discussing $\mathrm{CUA}_x$, first note that $\mathrm{CUA}_x$ is continuous w.r.t. the parameter $x$ on the interval $x \in (-\infty, \mathrm{UA}_1(\mathbf{y}))$; this will be shown later on in Corollary 1. As already mentioned in Section 2, $\mathrm{CUT}_x$ is discontinuous w.r.t. this threshold parameter. Second, notice that by calculating the $k \in [1, n]$ such that $\mathrm{UA}_k(\mathbf{y}) = x$, $\mathrm{CUA}_x$ is counting all components with magnitude greater than $x$ and some components with magnitude less than $x$. These magnitudes can contain important and meaningful information that is ignored by $\mathrm{CUT}_x$.

For example, assume you are deciding whether or not to place a service facility in a particular location (e.g. a cell phone tower), with the service facility being able to serve only houses within a geographic radius of $R$ miles. Let $\mathbf{y}$ then represent the entire set of customers you wish to serve from this location and their distance from the facility. To assess the quality of this location, it is intuitive to look at $\mathrm{CUT}_R(\mathbf{y})$ which gives you the number of customers within this set that will go unserved by this facility. However, the use of such a *hard* threshold is quite unintuitive under closer inspection. Decision makers would certainly like to know if a large number of customers are located $R + \epsilon$ miles away from the facility where $\epsilon$ is a very small number. Similarly, they would also like to know if a large number of customers are $R - \epsilon$ miles away. These characteristics are critically important to assessing the quality of this facility and the service it can provide to a set of desired customers. In this example, $\mathrm{CUA}_x$ can be an important counterpart to $\mathrm{CUT}_x$, acting as a *soft* threshold, counting the number of customers *around* $R$ miles away from the facility.

Consider also a disaster relief agency that is analyzing historical hurricane damage data where damages are in dollars [7]. Assume that the agency is attempting to determine if allocating \$$B$ dollars to the relief fund for the next hurricane is sufficient. First, of the historical damage amounts that exceed \$$B$, it is important to know by how much these damages exceeded \$$B$, especially if the exceedance is large. Secondly, it is important to know the magnitude of the largest damage amounts that are less than \$$B$. Are these damages very close to \$$B$? Or are they much smaller than \$$B$?
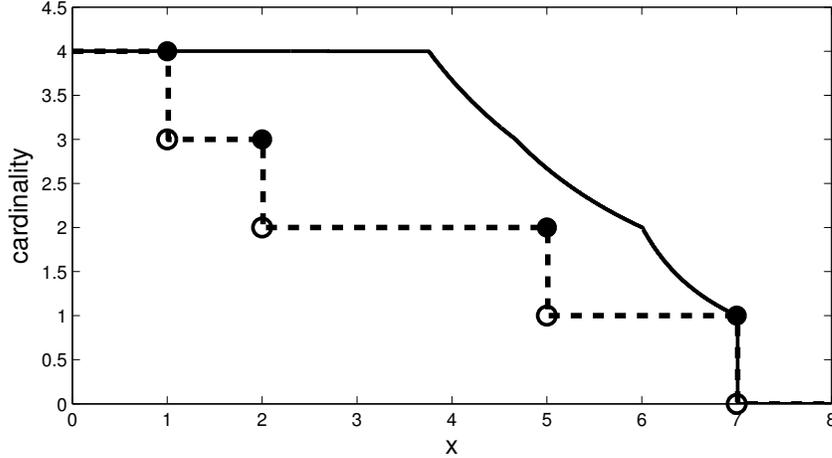
Figure 1: Cardinality of Upper Tail $(\text{CUT}_x)$, $\eta_x(\mathbf{y})$, dashed line, and Cardinality of Upper Average $(\text{CUA}_x)$, $\bar{\eta}_x(\mathbf{y})$, solid line, as functions of threshold $x$, where $\mathbf{y} = (1, 2, 5, 7)$.

Answers to these questions are clearly important for budgetary considerations.

Next, consider the simple, illustrative example in Fig. 1 with data vector $\mathbf{y} = (1, 2, 5, 7)$. For this vector, $\text{CUA}_x = 2$ for $x = 6$, because the average of the two largest components of the vector $\mathbf{y}$ equals $(7 + 5)/2 = 6$, and $\text{CUA}_x = 3$ for $x = 4\frac{2}{3} = (2 + 5 + 7)/3$. We can observe that $\text{CUA}_x$ is an upper bound for $\text{CUT}_x$. Additionally, for this example, Fig. 1 shows that $\text{CUA}_x$ is continuous w.r.t. $x$ except at the maximum point where threshold $x = 7$, while $\text{CUT}_x$ is discontinuous at $x = 1, 2, 5$, and 7.

With $\text{CUA}_x$ being the inverse of $\text{UA}_k$, both quantities relay similar information. We mention, though, that $\text{CUA}_x$ may be more intuitive to use when data have meaningful units, particularly because the threshold parameter $x$ is posed in the associated units. For example, if the data represent monetary losses (in the unit of dollars) from hurricane damage, it can be more intuitive to work with a threshold $x$, which will be some dollar amount, rather than with the $k$-parameter. Consider, as before, the disaster relief agency that is attempting to determine if allocating \$$B$ dollars to the relief fund for the next hurricane is sufficient. It would be much more intuitive to look at $\text{CUA}_B$ than to look at $\text{UA}_k$ for many different values of $k$. Additionally, this also applies when considering some optimization tasks. Assume an investor is trying to form a portfolio by analyzing historical stock behavior (where each possible portfolio generates some data set of historical losses it would have incurred in the market). If an investor does not want to incur losses larger than \$$B$ dollars, it is more intuitive to try and devise a portfolio that minimizes $\text{CUA}_B$ rather than to find a portfolio that minimizes $\text{UA}_k$ for some appropriate value of $k$ which would need to be determined.

We also highlight the simple, but useful fact that $\text{CUA}_x$ tells you immediately the total weight that is in the tail. Consider again the example from the introduction where we have a network of 10,000 servers and we find that $\text{CUA}_{10} = 207.3$. $\text{CUA}_x$ immediately tells us that $10 \times 207.3$ additional units of resource are needed to process jobs in the 207.3 longest queues. On the other hand, $\text{CUT}_{10} = 3$ is much less informative, telling us only that the 3 longest queues need *at least* 30 additional units

8

of resource, even though this number could be much larger. We can put this more precisely, though. Assume that $a^*$ is an optimal point for $\text{CUA}_x$ calculation (4) at threshold $x$. We prove in Section 3.2.2 that,

$$x\bar{\eta}_x(\mathbf{y}) \leq \sum_i \{y_i | y_i \geq x - \frac{1}{a^*}\} \leq x(\lfloor \bar{\eta}_x(\mathbf{y}) \rfloor + 1),$$

where $\lfloor k \rfloor$ denotes the largest integer less than or equal to $k$ and $\lceil k \rceil$ denote the smallest integer greater than or equal to $k$. Therefore, we know that the total amount of resource needed to process all unserved jobs waiting in queues with length longer than or equal to $x - \frac{1}{a^*}$ is somewhere between $x\bar{\eta}_x(\mathbf{y})$ and $x(\lfloor \bar{\eta}_x(\mathbf{y}) \rfloor + 1)$. Additionally, note that the continuity of $\text{CUA}_x$ can help determine where these resources are needed. For example, consider the case where all queues have load less than 10 except for three servers, which are extremely overloaded. In this case, one can take a partial derivative of $\text{CUA}_x$ w.r.t. the threshold to analyze how dramatically $\text{CUA}_x$ is changing as the threshold changes continuously.

**3.1.1. $\text{CUA}_x$ Calculation via Linear Interpolation.** Suppose that one would like to calculate $\text{CUA}_x$ for multiple threshold levels $x \in \mathbb{R}$. One could utilize formula (4) to achieve this task, but this proves quite inefficient if one would like to know $\text{CUA}_x$ for all thresholds $x \in \mathbb{R}$. When one is simply interested in calculating $\text{CUA}_x$ for all thresholds $x \in \mathbb{R}$ for a fixed vector $\mathbf{y}$, it is possible to utilize linear interpolation to calculate $\text{CUA}_x$ at all thresholds $x \in \mathbb{R}$. To show this, we first introduce the following alternative calculation formula for $\text{CUA}_x$, which shows that $\text{CUA}_x$ can be calculated via minimization over a finite set of points. It also provides additional insight into formula (4).

**Proposition 1:** *Consider a Euclidean vector* $\mathbf{y} = (y_1, \ldots, y_n)$ *and let* $y_0 = -\infty$. $\text{CUA}_x$ *of* $\mathbf{y}$ *at threshold* $x \in \mathbb{R}$ *equals*

$$(13) \qquad \bar{\eta}_x(\mathbf{y}) = \min_{j \in \{0,1,\ldots,n\}} \frac{\sum_{i=1}^{n}[y_i - y_j]^+}{[x - y_j]^+} .$$

**Proof:** A change of variable $a = \frac{1}{x-\gamma}$ transforms minimization formula (4) to the following formula:

$$(14) \qquad \bar{\eta}_x(\mathbf{y}) = \inf_{\gamma < x} \frac{\sum_{i=1}^{n}[y_i - \gamma]^+}{x - \gamma}.$$

Let us show that the objective is minimal for either $\gamma = y_j$ for $j = 1, \ldots, n$, or for $\gamma \to -\infty = y_0$. Suppose the contrary. Then the objective is differentiable at the optimal $\gamma^*$, since only $\gamma = y_i$ are the points where derivative over $\gamma$ does not exist. Then the derivative of the objective w.r.t. $\gamma$ must be 0 at $\gamma^*$:

$$(x-\gamma)^{-2}\left[(x-\gamma)\left(\sum_{i=1}^{n}[y_i - \gamma]^+\right)'_\gamma - (x-\gamma)'_\gamma \sum_{i=1}^{n}[y_i - \gamma]^+\right]\Bigg|_{\gamma=\gamma^*} = 0,$$

$$(15) \qquad \sum_{y_i > \gamma^*}(-1)(x-\gamma^*) - (-1)\sum_{y_i > \gamma^*}(y_i - \gamma^*) = \sum_{y_i > \gamma^*}(y_i - x) = 0.$$

Suppose $\gamma^* \in (y_i, y_{i+1})$, or $\gamma^* \in (y_i, x)$, for some $i = 0, \ldots, n$. From the derivative expression above it follows that the derivative should be equal to 0 on the corresponding

9

interval, therefore, the objective function has the same value at $\gamma = y_i$ and $\gamma^*$, and thus $\gamma = y_i$ is also an optimal solution, which is a contradiction. $\square$

Utilizing this proposition, the following Corollary 1 shows that $\mathrm{CUA}_x$ can be calculated via simple linear interpolation. Proposition 1, by itself, though, provides interesting insights. First, we see that calculation can be performed by only considering the finite set of points $(y_0, y_1, \ldots, y_n)$. Second, it follows from Proposition 1 that if $y_j$ is the argmin of (13), then $a^* = \frac{1}{x - y_j}$ is an argmin of (4). This fact can be seen more clearly in the proof of Proposition 1 (i.e. see the change of variable that takes place).

Moving to the discussion of linear interpolation, suppose we have the vector $\mathbf{y} \in \mathbb{R}^n$. Instead of using formula (4) to calculate $\mathrm{CUA}_x$, one only needs to calculate $\mathrm{UA}_k(\mathbf{y})$ for $k \in \{1, \ldots, n-1\}$, which effectively calculates $\mathrm{CUA}_x$ for thresholds $x \in \{\mathrm{UA}_{n-1}(\mathbf{y}), \ldots, \mathrm{UA}_1(\mathbf{y})\}$, then utilize linear interpolation to calculate $\mathrm{CUA}_x$ for the intermediate threshold values.

**Corollary 1:** *The function $\frac{1}{\bar{\eta}_x(\mathbf{y})}$ is a piecewise-linear convex function of $x$ with knots at $x \in \{\mathrm{UA}_n(\mathbf{y}), \mathrm{UA}_{n-1}(\mathbf{y}), \ldots, \mathrm{UA}_1(\mathbf{y})\}$. Specifically, for any threshold $x = \lambda \mathrm{UA}_{i+1}(\mathbf{y}) + (1 - \lambda)\mathrm{UA}_i(\mathbf{y})$, $i \in \{1, \ldots, n-1\}$, $\lambda \in (0, 1)$, we have that*

(16)
$$\frac{1}{\bar{\eta}_x(\mathbf{y})} = \frac{\lambda}{\bar{\eta}_{\mathrm{UA}_{i+1}(\mathbf{y})}(\mathbf{y})} + \frac{1 - \lambda}{\bar{\eta}_{\mathrm{UA}_i(\mathbf{y})}(\mathbf{y})} \ .$$

**Proof:** The derivative expression (15) for minimization problem (14) implies that at $x \in [\mathrm{UA}_{n-j-1}(\mathbf{y}), \mathrm{UA}_{n-j}(\mathbf{y})]$ and $\gamma = y^{(j)}$, where the superscript denotes the ordered vector, the derivative $\sum_{y_i > \gamma - \epsilon}(y_i - x) \leq 0$, the derivative $\sum_{y_i > \gamma + \epsilon}(y_i - x) \geq 0$ for all $\epsilon > 0$. Therefore, $\gamma = y^{(j)}$ is optimal for $x \in [\mathrm{UA}_{n-j-1}(\mathbf{y}), \mathrm{UA}_{n-j}(\mathbf{y})]$ and

$$\bar{\eta}_x(\mathbf{y}) = \frac{\sum_{i=1}^{n}[y_i - y^{(j)}]^+}{x - y^{(j)}} \ .$$

Hence,

$$\frac{1}{\bar{\eta}_x(\mathbf{y})} = \frac{\lambda}{\bar{\eta}_{\mathrm{UA}_{i+1}(\mathbf{y})}(\mathbf{y})} + \frac{1 - \lambda}{\bar{\eta}_{\mathrm{UA}_i(\mathbf{y})}(\mathbf{y})}$$

for $x = \lambda \mathrm{UA}_{i+1}(\mathbf{y}) + (1 - \lambda)\mathrm{UA}_i(\mathbf{y})$, $i \in \{1, \ldots, n-1\}$, $\lambda \in (0, 1)$ and values $\{\mathrm{UA}_n(\mathbf{y}), \ldots, \mathrm{UA}_1(\mathbf{y})\}$ are knots for the piecewise linear function $\frac{1}{\bar{\eta}_x(\mathbf{y})}$. $\square$

Thus, we only need to calculate $\mathrm{CUA}_x$ for $n-1$ thresholds, namely $x \in \{\mathrm{UA}_{n-1}(\mathbf{y}), \ldots, \mathrm{UA}_1(\mathbf{y})\}$ and we can "fill in" the missing thresholds via linear interpolation. Using the example $\mathbf{y} = (1, 2, 5, 7)$ from Section 3.1, we illustrate the piecewise linearity in Fig. 2 which plots $\frac{1}{\bar{\eta}_x(\mathbf{y})}$ on the vertical axis and $x$ on the horizontal axis.

**3.2. Connecting $\mathrm{CUA}_x$ and $\mathrm{CUT}_x$.** In this section, we discuss the connection between $\mathrm{CUA}_x$ and $\mathrm{CUT}_x$. We emphasize, though, that the core value of $\mathrm{CUA}_x$ is not rooted in its relationship with $\mathrm{CUT}_x$. As already mentioned, $\mathrm{CUA}_x$ is a complimentary characteristic to $\mathrm{CUT}_x$ that is useful in its own right, as it considers information about the magnitude of data points around the threshold, which $\mathrm{CUT}_x$ does not, and is a continuous function w.r.t. threshold. In Section 4, we discuss this further, showing that $\mathrm{CUA}_x$ optimization problems are interesting in their own right, and that their usefulness is not confined by their relationship with $\mathrm{CUT}_x$ minimization problems.
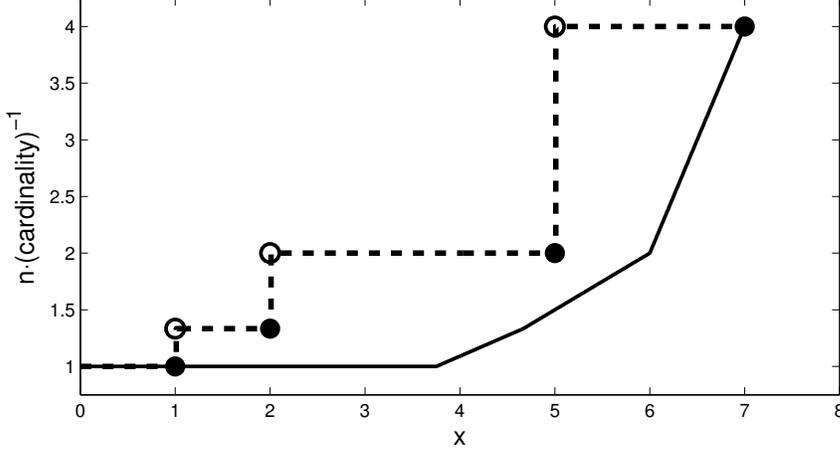
Figure 2: $\frac{1}{\eta_x(\mathbf{y})}$, dashed line, and $\frac{1}{\bar{\eta}_x(\mathbf{y})}$, solid line, as functions of threshold $x$, where $\mathbf{y} = (1, 2, 5, 7)$.

**3.2.1. $\text{CUA}_x$ as Upper Bound on $\text{CUT}_x$.** Fig. 1 shows for a specific example that $\text{CUA}_x$ acts as an upper bound for $\text{CUT}_x$. Specifically, this can be posed as the following relation, which follows intuitively from the definitions of $\text{CUA}_x$, $\text{UA}_k$, and $\text{CUT}_x$:

$$(17) \qquad \bar{\eta}_x(\mathbf{y}) = k \iff \text{UA}_k(\mathbf{y}) = x \implies \eta_x(\mathbf{y}) \leq k .$$

We can improve upon this notion of an upper bound, though, showing that it can be viewed as the minimal quasi-convex upper bound. This is shown in the following proposition, which is a Euclidean space alternative of a statement proved in [9]. An analogous statement for Value-at-Risk and Conditional-Value-at-Risk can be found in [8]. We call a function on $\mathbb{R}^n$ symmetric if permutation of components of a vector does not change the value of the function. A function $g$ is called quasi-convex if its level-sets $\{\mathbf{y}|g(\mathbf{y}) \leq c\}$ are convex for all $c$, or, equivalently, $g(\lambda \mathbf{y} + (1-\lambda)\mathbf{z}) \leq \max\{g(\mathbf{y}), g(\mathbf{z})\}$ for all $\lambda \in (0, 1)$ and $\mathbf{y}, \mathbf{z} \in \mathbb{R}^n$.

**Proposition 2:** *Consider a function $g : \mathbb{R}^n \to \mathbb{R}$ which is symmetric, quasi-convex, and is greater than $\eta_x$ everywhere on $\mathbb{R}^n$. Then $g(\mathbf{y}) \geq \lfloor \bar{\eta}_x(\mathbf{y}) \rfloor$ for all $\mathbf{y} \in \mathbb{R}^n$.*

**Proof:** Indeed, suppose that $g(\mathbf{y}) < \lfloor \bar{\eta}_x(\mathbf{y}) \rfloor$, which implies that $\bar{\eta}_x(\mathbf{y}) \geq 1$. By the property of CUA, there exist $k \equiv \lfloor \bar{\eta}_x(\mathbf{y}) \rfloor$ components of $\mathbf{y}$ whose average is at least $x$. Denote indices of these components by $I = \{i_1, \ldots, i_k\}$. Consider vectors $\mathbf{y}_1, \ldots, \mathbf{y}_k$ obtained from $\mathbf{y}$ by applying the cyclic permutation $(i_1, \ldots, i_k)$ to its components $i$ times for $\mathbf{y}_i$. Since $g$ is symmetric, $g(\mathbf{y}_j) = g(\mathbf{y})$. Consider $\mathbf{y}' = \frac{1}{k}\sum_{j=1}^{k} \mathbf{y}_j$, then $y'_i = \text{UA}_k(\mathbf{y}) \geq x$ for all $i \in I$, therefore, $\eta_x(\mathbf{y}') \geq k$. Notice that since $g$ is quasi-convex, then $g(\mathbf{y}') \leq \max\{g(\mathbf{y}_1), \ldots, g(\mathbf{y}_k)\} = g(\mathbf{y}) < k \leq \eta_x(\mathbf{y}')$. That is, this contradicts the assumption that $g$ is greater than $\eta_x$. $\square$

Notice that since function $\lfloor \cdot \rfloor$ is nondecreasing, and $\lfloor \max\{a, b\} \rfloor = \max\{\lfloor a \rfloor, \lfloor b \rfloor\}$, then $\lfloor \bar{\eta}_x \rfloor$ is itself a symmetric, quasi-convex function, which is greater than $\eta_x$ and that, moreover, it is the smallest function in the class of symmetric quasi-convex

11

functions that are greater than $\eta_x$.

This may be quite useful in practice, particularly in an optimization context, as $\mathrm{CUA}_x$ can be viewed as an efficiently calculable upper bound for $\mathrm{CUT}_x$.

**3.2.2. Simultaneous Calculation of $\mathrm{CUA}_x$ and $\mathrm{CUT}_x$ .** An important and interesting property of $\mathrm{CUA}_x$ calculation formula (4) is that calculation of $\mathrm{CUA}_x$ provides us information about $\mathrm{CUT}_x$. Specifically, we have the following property of formula (4).

**Proposition 3:** *Suppose for a vector $\mathbf{y} \in \mathbb{R}^n$ and threshold $x \in (\mathrm{UA}_n(\mathbf{y}), \mathrm{UA}_1(\mathbf{y}))$ we have that*

$$\bar{\eta}_x(\mathbf{y}) = \sum_{i=1}^{n}[a^*(y_i - x) + 1]^+ = k \ ,$$

*where $a^*$ is an optimal point for (4). Then, if $k$ is noninteger, we have that $a^* = \frac{1}{x - y^{(n - \lfloor k \rfloor)}}$ and is the unique minimizer of (4). Furthermore, we have that*

$$\eta_{x - \frac{1}{a^*}}(\mathbf{y}) = \lceil k \rceil \ .$$

*Otherwise, if $k$ is integer, the set of minimizers of (4) is given by the interval $\mathcal{I} = [\frac{1}{x - y^{(n-k)}}, \frac{1}{x - y^{(n-k+1)}}]$ where $a^* \in \mathcal{I}$. Furthermore, we have that*

$$(18) \qquad \eta_{x - \frac{1}{a^*}}(\mathbf{y}) = \begin{cases} k + 1 \ , & \text{if } a^* = \frac{1}{x - y^{(n-k)}} \ , \\ k \ , & \text{otherwise.} \end{cases}$$

**Proof:** For this proof, let us denote the $k$th largest component of $\mathbf{y}$ as $y^{(k)}$. Let us also denote, for any real valued random variable $X$ the *upper* quantile as $q_\alpha^+(X) = \inf\{x | P(X \leq x) > \alpha\}$ and the *lower* quantile as $q_\alpha(X) = \min\{x | P(X \leq x) \geq \alpha\}$, where $\alpha \in [0, 1]$ is a probability level. Note that when we discuss the quantiles of $\mathbf{y}$, we are treating $\mathbf{y}$ as a discretely distributed random variable with equally probable scenarios $y_1, ..., y_n$.

From [15], we know that if $\mathrm{UA}_k(\mathbf{y}) = x$ and $\gamma^* \in \underset{\gamma}{\operatorname{argmin}} \ \gamma + \frac{1}{k} \sum_{i=1}^{n}[y_i - \gamma]^+$, then $\gamma^* \in [q_{(\frac{n-k}{n})}(\mathbf{y}), q_{(\frac{n-k}{n})}^+(\mathbf{y})]$, where $k$ is related to the probability level by the equation $k = n(1 - \alpha)$. Let us now look at the two cases where $k$ is either integer or non-integer. For both cases, recall that as shown in Case 1 for the proof of Theorem 1 that if

$$\begin{aligned}
(19) \qquad \bar{\eta}_x(\mathbf{y}) &= \min_{x - \gamma > 0} \frac{\sum_{i=1}^{n}[y_i - \gamma]^+}{x - \gamma} \\
&= \frac{\sum_{i=1}^{n}[y_i - \gamma^*]^+}{x - \gamma^*} \\
&= k \ ,
\end{aligned}$$

then, $\mathrm{UA}_k(\mathbf{y}) = x$ and $\gamma^* \in \underset{\gamma}{\operatorname{argmin}} \ \gamma + \frac{1}{k} \sum_{i=1}^{n}[y_i - \gamma]^+$ and that $a^* = \frac{1}{x - \gamma^*}$ is an optimal point for (4).

**Case 1:** Assume that $k$ is non-integer. Then, we have that

$$q_{(\frac{n-k}{n})}(\mathbf{y}) = q_{(\frac{n-k}{n})}^+(\mathbf{y}) = y^{(n - \lfloor k \rfloor)} \ .$$

12

Thus, we have that

$$\bar{\eta}_x(\mathbf{y}) = k \implies x - \frac{1}{a^*} \in [q_{(\frac{n-k}{n})}(\mathbf{y}), q^+_{(\frac{n-k}{n})}(\mathbf{y})] = y^{(n-\lfloor k \rfloor)} .$$

Furthermore, this implies that

$$\eta_{x-\frac{1}{a^*}}(\mathbf{y}) = \lceil k \rceil \text{ and } a^* = \frac{1}{x - y^{(n-\lfloor k \rfloor)}} .$$

**Case 2:** Assume that $k$ is integer. Then, we have that

$$[q_{(\frac{n-k}{n})}(\mathbf{y}), q^+_{(\frac{n-k}{n})}(\mathbf{y})] = [y^{(n-k)}, y^{(n-k+1)}] .$$

Thus, we have that

$$\bar{\eta}_x(\mathbf{y}) = k \implies x - \frac{1}{a^*} \in [q_{(\frac{n-k}{n})}(\mathbf{y}), q^+_{(\frac{n-k}{n})}(\mathbf{y})] = [y^{(n-k)}, y^{(n-k+1)}] .$$

Furthermore, this implies that $a^* \in [\frac{1}{x-y^{(n-k)}}, \frac{1}{x-y^{(n-k+1)}}]$ and that

(20)
$$\eta_{x-\frac{1}{a^*}}(\mathbf{y}) = \begin{cases} k+1 , & \text{if } x - \frac{1}{a^*} = y^{(n-k)} , \\ k , & \text{otherwise.} \end{cases}$$

$\square$

Proposition 3 is quite useful, as it shows that simply calculating $\text{CUA}_x$ provides information about $\text{CUT}_x$. Note that this result can be viewed as analogous to one regarding Conditional-Value-at-Risk in [15] which shows that the value of Value-at-Risk can be calculated as a byproduct of the calculation formula (2).

In the context of our earlier example involving 10,000 servers, we said that $x \times \text{CUA}_x$ told us how many additional resources were needed to cover the $\text{CUA}_x$ busiest queues. Using Proposition 3, we can expand upon this statement with the following Corollary.

**Corollary 2:** Suppose that for a vector $\mathbf{y} \in \mathbb{R}^n$ and threshold $x \in (\text{UA}_n(\mathbf{y}), \text{UA}_1(\mathbf{y}))$ we have that $\bar{\eta}_x(\mathbf{y}) = \sum_{i=1}^n [a^*(y_i - x) + 1]^+ = k$. Then,

$$x\bar{\eta}_x(\mathbf{y}) \leq \sum_i \{y_i | y_i \geq x - \frac{1}{a^*}\} \leq x(\lfloor \bar{\eta}_x(\mathbf{y}) \rfloor + 1) .$$

**Proof:** We first prove the left hand inequality. Recall that as shown in Case 1, proof of Theorem 1, that if

(21)
$$\bar{\eta}_x(\mathbf{y}) = \min_{x-\gamma>0} \frac{\sum_{i=1}^n [y_i - \gamma]^+}{x - \gamma} = \frac{\sum_{i=1}^n [y_i - \gamma^*]^+}{x - \gamma^*} = k ,$$

then, $\text{UA}_k(\mathbf{y}) = x$ and $\gamma^* \in \underset{\gamma}{\text{argmin}} \ \gamma + \frac{1}{k} \sum_{i=1}^n [y_i - \gamma]^+$. Rearranging and applying the change of variable $a^* = \frac{1}{x-\gamma^*}$, or equivalently $\gamma^* = x - \frac{1}{a^*}$, we have

$$\bar{\eta}_x(\mathbf{y})(x - (x - \frac{1}{a^*})) = \sum_i [y_i - (x - \frac{1}{a^*})]^+$$

13

which then gives

$$x\bar{\eta}_x(\mathbf{y}) = \sum_i [y_i - (x - \frac{1}{a^*})]^+ + (x - \frac{1}{a^*})\bar{\eta}_x(\mathbf{y}) \ .$$

Noting that Proposition 3 implies that $\bar{\eta}_x(\mathbf{y}) \le \eta_{x-\frac{1}{a^*}}(\mathbf{y})$, we finally have that

$$x\bar{\eta}_x(\mathbf{y}) = \sum_i [y_i - (x - \frac{1}{a^*})]^+ + (x - \frac{1}{a^*})\bar{\eta}_x(\mathbf{y})$$

$$\le \sum_i [y_i - (x - \frac{1}{a^*})]^+ + (x - \frac{1}{a^*})\eta_{x-\frac{1}{a^*}}(\mathbf{y}) = \sum_i \{y_i | y_i \ge x - \frac{1}{a^*}\}$$

Now we prove the right hand inequality. First, note that if $\bar{\eta}_x(\mathbf{y}) = k$, then $\mathrm{UA}_k(\mathbf{y}) = x$. Second, note that Proposition 3 implies that $k \le \eta_{x-\frac{1}{a^*}}(\mathbf{y})$ implying further that $x = \mathrm{UA}_k(\mathbf{y}) \ge \mathrm{UA}_{\eta_{x-\frac{1}{a^*}}(\mathbf{y})}(\mathbf{y})$. Third, note that Proposition 3 also implies that $\eta_{x-\frac{1}{a^*}} \le \lfloor \bar{\eta}_x(\mathbf{y}) \rfloor + 1$. Finally, note that $\sum_i \{y_i | y_i \ge x - \frac{1}{a^*}\} = \mathrm{UA}_{\eta_{x-\frac{1}{a^*}}(\mathbf{y})}(\mathbf{y})\eta_{x-\frac{1}{a^*}}$. Using these facts, we have that

$$\sum_i \{y_i | y_i \ge x - \frac{1}{a^*}\} = \mathrm{UA}_{\eta_{x-\frac{1}{a^*}}(\mathbf{y})}(\mathbf{y})\eta_{x-\frac{1}{a^*}} \le x\eta_{x-\frac{1}{a^*}} \le x(\lfloor \bar{\eta}_x(\mathbf{y}) \rfloor + 1) \ .$$

□

Therefore, Corollary 2, in the context of the server example, shows us that the total amount of unserved jobs waiting in queues with length longer than or equal to $x - \frac{1}{a^*}$ is somewhere between $x\bar{\eta}_x(\mathbf{y})$ and $x(\lfloor \bar{\eta}_x(\mathbf{y}) \rfloor + 1)$.

**4. Optimization of** $\mathrm{CUA}_x$**.** When entered into an optimization setting, $\mathrm{CUA}_x$ provides substantial benefits, particularly when compared to the analogous use of $\mathrm{CUT}_x$. Consider the following $\mathrm{CUT}_x$ minimization problem, where $S \subseteq \mathbb{R}^n$ is a convex set and $x \in \mathbb{R}$ is a specified threshold:

$$\min_{\mathbf{y} \in S} \eta_x(\mathbf{y})$$

In order to solve this problem, it is natural to reformulate it as an MIP. In other words, to indicate the state of a vector component as being larger than $x$ or not, introduction of binary variables is a natural consideration.

We consider minimization of $\mathrm{CUA}_x$, which is posed as follows.

$$\min_{\mathbf{y} \in S} \bar{\eta}_x(\mathbf{y}) = \min_{\mathbf{y} \in S} \left\{ \min_{a \ge 0} \sum_{i=1}^n [a(y_i - x) + 1]^+ \right\} = \min_{\mathbf{y} \in S, a \ge 0} \sum_{i=1}^n [a(y_i - x) + 1]^+ \ .$$

As we will show by multiple examples, if the set $S$ is convex, the $\mathrm{CUA}_x$ minimization problem can be reduced to convex programming. Furthermore, if the set $S$ is a polyhedron, the $\mathrm{CUA}_x$ minimization problem can be reduced to linear programming. This is, in fact, a general result which follows from [9]. Specifically, as we will show in specific cases, as long as $S$ is convex, we can make a simple change of variable $ay_i = \hat{y}_i$ to yield the equivalent convex problems:

$$\min_{\mathbf{y} \in S, a \ge 0} \sum_{i=1}^n [a(y_i - x) + 1]^+ = \min_{\hat{\mathbf{y}} \in \hat{S}, a \ge 0} \sum_{i=1}^n [\hat{y}_i - ax + 1]^+ \ ,$$

14

where $\hat{S} = \mathrm{cl}\,\mathrm{cone}(S)$ is a closed convex cone. Linear programming reformulations follow similarly.

Additionally, $\mathrm{CUA}_x$ minimization itself is meaningful, as it considers the magnitude of the components around the threshold $x$. Therefore, while also minimizing an upper bound for $\mathrm{CUT}_x$, the $\mathrm{CUA}_x$ minimization is interesting in its own right. To demonstrate this, we provide an example network flow problems where $\mathrm{CUA}_x$ optimization may be more appealing than the analogous $\mathrm{CUT}_x$ minimization.

**4.1. Applications to Network Optimization.** We now present variations of three standard network flow problems in the context that we have edges or nodes that become overloaded and thus we would like to minimize the number of edges/nodes that are overloaded. To solve network flow problems in this context, we use formulations with a $\mathrm{CUA}_x$ or $\mathrm{CUT}_x$ objective. We first show that the $\mathrm{CUA}_x$ minimization problems reduce to LP's, yielding obvious computational advantages over the $\mathrm{CUT}_x$ minimization problems which are formulated as MIP's. We bolster this point by providing a computational example and by proving an NP-hardness result for $\mathrm{CUT}_x$ minimization problems. We then show that, depending on the context, the optimal $\mathrm{CUT}_x$ network flows can be unappealing, yielding edges or nodes that are severely overloaded or flow configurations that are more sensitive to minor fluctuations in network conditions. Optimal $\mathrm{CUA}_x$ flows, on the other hand, can be shown to mitigate these issues, providing relatively stable solutions with more evenly distributed flow patterns. The three network flow problems we discuss are the generalized assignment, min cost flow, and capacity planning problems.

Of course, both $\mathrm{CUA}_x$ and $\mathrm{CUT}_x$ have their merits, and we make no claim that one is always superior to the other. In general, MIP formulations have much more flexibility in terms of modeling capabilities as compared to LP's. Nevertheless, we focus on the benefits that the efficiently solvable $\mathrm{CUA}_x$ formulations can provide in comparison to $\mathrm{CUT}_x$ formulations.

**4.1.1. The Generalized Assignment Problem.** Throughout the remaining sections, we consider a network represented by a graph $G = (V, E)$ of edges $E$ and vertices $V$, where loads of product flow from supply vertices $V_S$, through the graph edges $E$ and transshipment vertices $V_T$, to the destination demand vertices $V_D$. It is common in network flow problems to be concerned with the amount of flow going through each node, which we generally call the *load*. For example, consider the following Generalized Assignment Problem from [2]. We must assign $|V_S|$ jobs to $|V_T|$ machines[2], where $A_{ij}$ is the unit cost of assigning job $i$ to machine $j$, and $t_{ij}$ is the time it takes machine $j$ to process a unit of job $i$. Just as in [2], we assume $f_{ij}$ is the amount of job $i$ assigned to machine $j$ and that we can assign partial jobs to machines (i.e. we do not require integral $f_{ij}$ assignments). The goal, then, is to assign all jobs at minimal cost such that no machine is running for more than $x$ time units. For this problem, it is considered undesirable to have the time-*load* of a machine larger than $x$. A machine might, for example, become prone to malfunctions, overheating,

---

[2]In this context, we have no demand nodes and each supply node has supply equal to 1.

or other risk factors if it becomes overloaded.

$$\min_{f} \quad \sum_{(i,j)\in E} f_{ij} A_{ij}$$

(22)

$$s.t. \quad \sum_{(i,j)\in E} f_{ij} = 1, \forall i \in V_S$$

$$\sum_{(i,j)\in E} f_{ij} t_{ij} \leq x, \forall j \in V_T$$

$$f_{ij} \geq 0.$$

Consider, though, the case where we would like to solve such a problem, but our $x$ makes the problem infeasible, meaning that we must overload *some* of the machines. We could, of course, simply increase $x$, but may cause us to overload a large number of machines. Thus, we would like to find an assignment that minimizes the number of machines that are overloaded. Also, to make sure our assignment has reasonable cost, we want it to satisfy a certain budget $\sum_{(i,j)\in E} f_{ij} A_{ij} \leq B$. A natural way to formulate this problem is as a CUT$_x$ minimization problem (23). Here, $\xi_j$ indicates whether the load $l_j$ on machine $j$ exceeds the threshold and $M$ is a sufficiently large constant.

(23)

$$\min_{f,\xi,l} \quad \sum_{j\in V_T} \xi_j$$

$$s.t. \quad \xi_j \geq \frac{l_j - x}{M}, \forall j \in V_T$$

$$l_j = \sum_{(i,j)\in E} f_{ij} t_{ij}, \forall j \in V_T$$

$$\sum_{(i,j)\in E} f_{ij} = 1, \forall i \in V_S$$

$$\sum_{(i,j)\in E} f_{ij} A_{ij} \leq B$$

$$f_{ij} \geq 0, \xi_{ij} \in \{0,1\}.$$

(24)

$$\min_{y,z,a,l} \quad \sum_{j\in V_T} z_j$$

$$s.t. \quad z_j \geq l_j - ax + 1$$

$$l_j = \sum_{(i,j)\in E} y_{ij} t_{ij}, \forall j \in V_T$$

$$\sum_{(i,j)\in E} y_{ij} = a, \forall i \in V_S$$

$$\sum_{(i,j)\in E} y_{ij} A_{ij} \leq aB$$

$$a \geq 0, y_{ij} \geq 0, z_j \geq 0.$$

Though this MIP will solve the problem of minimizing the number of overloaded nodes (machines), it is ignoring potentially meaningful information about the magnitude of loads put upon nodes. Specifically, it does not account for the magnitude of loads put upon nodes in the overloaded state, as well as the nodes with loads less than, but most near to the overloaded state. For nodes that are already overloaded, it may be important to consider by how much they are overloaded. Additionally, it may be undesirable to have nodes that, while not yet overloaded, are very close to being in the overloaded state. For example, assume that if a machine runs for too long, it may malfunction, and there is a large cost that must be paid to fix the broken machine. Let $I\{l_i > x\}$ be an indicator function, indicating if a node is overloaded. Now suppose that, for a given assignment configuration and associated loads $l_i$, the probability of a malfunction occurring at node $i$ is given by,

$$P(\text{malfunction at node } i) = .5I\{l_i > x\} + .5\left(1 - \frac{1}{e^{[l_i - x]^+}}\right).$$

Therefore, we have that a transshipment node will not malfunction until it is over-loaded and once overloaded, the probability jumps to .5 and increases exponentially based on the magnitude of the load vector. Solving the $\text{CUT}_x$ minimization problem to calculate flow configuration for this system is ignoring critical information, and may produce a very undesirable result. First, notice that the $\text{CUT}_x$ minimization problem will ignore the magnitude with which a node is overloaded, thus ignoring the fact that *very* overloaded nodes will, with higher probability, incur the additional cost $C$. Second, notice that the $\text{CUT}_x$ minimization will ignore information about nodes with loads less than, but very near to $x$. This can cause major stability issues if there is uncertainty regarding exact system specifications at test-time. For example, consider a $\text{CUT}_x$ minimization solution which has many nodes with loads less than, but very near to $x$. If, at test-time, $t_{ij}$ is slightly larger than expected, we will have a cascade of increasing malfunction probabilities as the nodes with load close to $x$ are pushed over the threshold with their failure probability increasing from zero to around .5. In Example 1, we demonstrate numerically that $\text{CUT}_x$ optimal solutions indeed exhibit these characteristics.

By minimizing $\text{CUA}_x$, we formulate a network flow problem that accounts for these factors (i.e. the magnitudes of loads). Specifically, we have the $\text{CUA}_x$ optimization problem (24) which aims to minimize the number of most loaded nodes that have loads averaging to $x$. Notice that this has been reduced to a LP via a simple change of variable $af_{ij} = y_{ij}$. This is a powerful result, which implies that we can achieve a similar goal (i.e. minimizing the number of nodes with loads exceeding, or slightly below a specified threshold) using LP as opposed to MIP. In terms of the network optimization problem, it considers how much load is burdening the overloaded nodes (i.e. by how much they are overloaded). Additionally, it considers how much load is burdening the nodes with loads less than, but most near to the threshold (i.e. how close to being overloaded are the most burdened nodes). Consider the following numerical example.

**Example 1:** We consider this problem with 50 jobs and 20 machines, where $(i, j) \in E$ with probability .5, $A_{ij}$ are uniformly distributed on $[0, 10]$, $t_{ij}$ are uniformly distributed on $[0, 1]$. For our randomly generated task, we find that the general assignment problem is infeasible at $x = 2.5$, so we consider this our threshold for $\text{CUA}_x$ and $\text{CUT}_x$ minimization. We also set $B = 25$, which was set such that the budget constraint would be tight for both $\text{CUA}_x$ and $\text{CUT}_x$ optimal assignment (i.e. to make sure we had reasonable assignments). In Table 1 we show the ordered list of loads given by the optimal $\text{CUA}_x$ and $\text{CUT}_x$ assignments. We see that $\text{CUT}_x$ severely overloads a single node in order to drive $\text{CUT}_x$ to 1. The $\text{CUA}_x$ solution, on the other hand, is spread much more evenly. For example, we see that the sum of nodes exceeding $x$ equals 23 for $\text{CUA}_x$ and 92.7 for $\text{CUT}_x$. We also see that the $\text{CUT}_x$ solution drives many loads to equal the threshold. As already mentioned, this can lead to stability issues, as uncertainty in $t_{ij}$ could lead to many machines being pushed over the time limit threshold. If there is a discrete jump in failure probability as load surpasses $x$, this could lead to heightened probability for a cascade of failures or a huge jump in the expected number of machine malfunctions. Note also that we could have introduced upper bounds $l_j \leq U_j$ on the loads in the $\text{CUT}_x$ formulation to prevent this extreme solution. However, the $\text{CUT}_x$ minimization will still have the same properties, overloading a few machines to the maximum, $U_j$, to drive down $\text{CUT}_x$ with many other machines having $l_j = x$. Additionally, it is still an MIP compared to a potential LP solution.

Table 1: Ordered List of Optimal Loads $l_j = \sum_{(i,j) \in E} f_{ij} t_{ij}$ for all $j \in V_T$. Overloaded nodes are bold face.

| CUA | 1.1 | 1.1 | 1.1 | 1.1 | 1.1 | 1.1 | 1.1 | 1.6 | 1.7 | 1.7 | 2.0 | 2.2 | 2.2 | 2.3 | 2.5 | **3.2** | **3.7** | **4.3** | **5.5** | **6.3** |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|---------|---------|---------|---------|---------|
| CUT | 0 | 2.2 | 2.4 | 2.5 | 2.5 | 2.5 | 2.5 | 2.5 | 2.5 | 2.5 | 2.5 | 2.5 | 2.5 | 2.5 | 2.5 | 2.5 | 2.5 | 2.5 | 2.5 | **92.7** |

**Remark:** Other objectives could be used for the optimization problems considered in this section that also consider the magnitudes of the largest components and are convex. One could use $\mathrm{UA}_k$ or the overall exceedance $\sum_i [l_i - x]^+$. However, these have drawbacks compared to $\mathrm{CUA}_x$. For $\mathrm{UA}_k$, although it is actually equivalent to $\mathrm{CUA}_x$ for the proper choice of $k$, one will still need to choose an appropriate $k$. In our examples, a threshold is natural, but choosing the correct $k$ is less straightforward. It is hard to tell how many of the largest components you want to minimize. For overall exceedance, there are two major differences. First, it does not count components. Thus, the expected exceedance may be small, but there may be a large number of components comprising this sum. Second, it does not consider the magnitude of the largest components that are *less* than the threshold. As already mentioned, this may be important to consider. In a separate set of experiments, we saw that minimization of overall exceedance often gives you many components lying on the threshold with $l_i = x$.

**Example 2:** For a simple computational demonstration comparing the performance of the $\mathrm{CUA}_x$ formulation versus the $\mathrm{CUT}_x$ formulation, we consider a formulation similar to the generalized assignment problem, but with the only change being an extra demand node which is connected to all $i \in V_T$ with cost of transmission being $t_{id}$ for all $i \in V_T$. One possible interpretation comes from an information flow networks area, just with the flow going in the opposite direction. The demand vertex is associated with a data server, supply vertices are associated with end users, and transshipment nodes are associated with routing servers[3]. The goal is to assign routing in a such way that the routing servers are not overloaded with information flow. The $\mathrm{CUA}_x$ and $\mathrm{CUT}_x$ formulations would be identical to (24) and (23) but with the loads including an extra term to equal $f_{jd} t_{jd} + \sum_{(i,j) \in E} f_{ij} t_{ij}$ for all $j \in V_T$ and flow conservation constraints for transshipment and demand nodes.

Using this network structure, we solved differently sized problems using Gurobi Solver on a PC via a Python interface. Python code and full problem description can be found online.[4] Note that results are reported for random graph instances having non-trivial solutions where the optimal objective does not equal 0 or $|V_T| : |V_S| = $ 10,000; $|V_T| = T$; $|V_D| = 1$; supply per supply node $= 10$; demand for demand node $= $ 100,000 $ = (10)$x$(10,000)$; $t_{ij}$ uniformly distributed in $[0,1]$; threshold for $\mathrm{CUA}_x = x$ $=.5$ x 100,000 / $T$; threshold for $\mathrm{CUT}_x = x = .2$ x 100,000 / $T$; M $= 100,000$.

With this setup, we compared the performance of the MIP $\mathrm{CUT}_x$ minimization formulation and the LP $\mathrm{CUA}_x$ minimization formulation for varying values of $T$. Run time comparisons can be seen in Table 2. We see that the LP formulation has a clear advantage as the number of transshipment nodes, i.e. the number of binary variables in the MIP formulation, increases. Thus, we see that the $\mathrm{CUA}_x$ minimization problem

---

[3]In this interpretation, flow goes from demand to supply, but the optimal flows are equivalent.
[4]http://www.ise.ufl.edu/uryasev/research/testproblems/logistics/
network-optimization-by-minimization-of-cardinality-of-upper-averages-cua/

| $T$ | MIP Run Time (seconds) | LP Run Time (seconds) |
|---|---|---|
| 100 | 2244.94 | 31.53 |
| 75 | 492.42 | 14.52 |
| 50 | 191.34 | 13.23 |
| 25 | 109.33 | 5.95 |
| 10 | 21.26 | 6.47 |

Table 2: Number of transshipment nodes, $T$, and approximate time to solve with Gurobi in Python interface.

is significantly faster. For large problems, CUA$_x$ solving time will be dramatically lower than MIP solving time.

**4.1.2. Min Cost Network Flows.** While the first two examples considered the load put upon nodes, it is also quite natural to consider the load (or flow) on individual arcs. The next two examples discuss problems of this type, where we would like to minimize the number of arcs with large flow or large overall flow costs. One of the most common network flow formulations is the Min Cost Network Flow problem [1]. In this formulation, $M$ is a sufficiently large constant, $f_{ij}$ indicates the flow through edge $(i,j) \in E$, $c_{ij}$ is the unit cost of flow, $u_{ij}$ is the capacity, $b_i$ is the demand at vertex $i \in V$ where $b_i = 0$ for $i \in V_T$, $b_i > 0$ for $i \in V_D$, and $b_i < 0$ for $i \in V_S$.

$$
\begin{aligned}
\min_f \quad & \sum_{(i,j)\in E} c_{ij} f_{ij} \\
\text{s.t.} \quad & \sum_{(i,j)\in E} f_{ij} - \sum_{(k,i)\in E} f_{ki} = b_i \ \forall i \in V, \\
& 0 \le f_{ij} \le u_{ij}.
\end{aligned}
$$
(25)

However, this problem makes the critical assumption that all edge costs $c_{ij}f_{ij}$ are equally important which may not be a valid assumption (see e.g. chapter 14 of [1]). It might be the case that we would not want any edge cost to exceed some level $x$. For example, $x$ may represent some budget associated with each arc and we do not want to overload too many budgets. Of course, this may not be possible for all edges and we may be forced to overload some budgets. Thus, this leads us to two possibilities for a risk averse min cost network flow problem where we would like minimize the number of edge costs that are large w.r.t. some threshold $x$. We have the CUA$_x$ minimization (26), which has been reduced to a LP by the change of variable $af_{ij} = y_{ij}$, and the CUT$_x$ minimization (27) which has been reduced to a MIP.

(26)
$$
\begin{aligned}
\min_{y,z,a} \quad & \sum_{(i,j)\in E} z_{ij} \\
\text{s.t.} \quad & z_{ij} \ge c_{ij} y_{ij} - ax + 1 \\
& \sum_{(i,j)\in E} y_{ij} - \sum_{(k,i)\in E} y_{ki} = ab_i, \\
& 0 \le y_{ij} \le au_{ij}. \\
& a \ge 0, z_{ij} \ge 0.
\end{aligned}
$$

(27)
$$
\begin{aligned}
\min_{f,\xi} \quad & \sum_{(i,j)\in E} \xi_{ij} \\
\text{s.t.} \quad & \sum_{(i,j)\in E} f_{ij} - \sum_{(k,i)\in E} f_{ki} = b_i \\
& 0 \le f_{ij} \le u_{ij}, \\
& \xi_{ij} \ge \frac{c_{ij} f_{ij} - x}{M}, \\
& \xi_{ij} \in \{0,1\}.
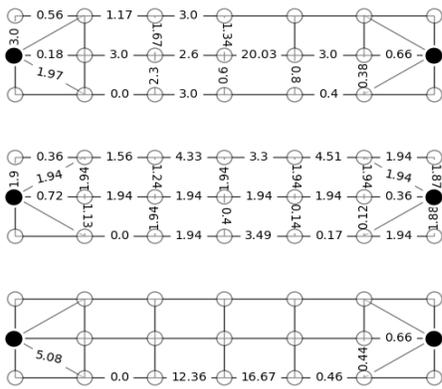\end{aligned}
$$

19

Figure 3: (Top) Optimal CUT flow. (Middle) Optimal $CUA_x$ flow. (Bottom) Min cost flow. Numbers of arcs represent $f_{ij}c_{ij}$
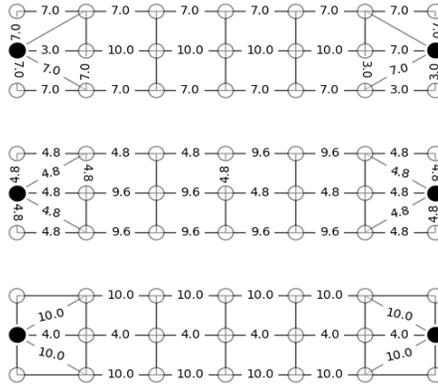
Figure 4: (Top) Optimal CUT flow. (Middle) Optimal $CUA_x$ flow. (Bottom) Min cost flow with arc costs equal to 1 and arc upper bound equal to 10. Numbers on arcs represent $f_{ij}$.

**Example 3:** Consider the network in Figure 3, where we would like to push 24 units of supply through the network from the left black node to the right black node. We have edge costs uniformly distributed on $[0, 1]$ and we assume that each arc has unbounded capacity and we set $x = 3$. The optimal $CUA_x$, $CUT_x$, and min cost flows are shown in Figure 3. We notice, first, that the min cost flow has only 4 edges with cost exceeding 3, but these costs are fairly large. The $CUA_x$ solution, on the other hand, has three edges with cost exceeding 3, but with much smaller magnitude. The $CUT_x$ solution only has one edge exceeding 3, but the magnitude of this edge is more than four times the magnitude of the largest $CUA_x$ edge cost. We also see that the $CUA_x$ flow is more spread out, with more edges being utilized. We also see that the utilization is directly influenced by the threshold. With the $CUT_x$ flow, we have some edge costs equal to the threshold, while the $CUA_x$ flow is almost opposite in that the largest edge cost that is less than or equal to 3 is equal to 1.93. Thus, there is a buffer between these smaller edge costs and the threshold which may be useful if costs or other parameters are uncertain.

**4.1.3. CUA and Capacity Planning.** The capacity planning network flow problem considers the situation where the current network is composed of arcs with capacity that is not large enough so that all demand can be satisfied. Therefore, the task is to figure out which arcs should have their capacity increased so that demand can be satisfied. Of course, many variants of this problem exist that are tailored to specific applications [11, 5, 4]. We consider, though, a generic variant where each arc has initial capacity $x$, and we assume it will cost $C$ to expand the capacity of any arc to $u$.[5] This problem can be posed as (28), a $CUT_x$ minimization network flow problem where we expand capacity of the minimal number of arcs so that demand can be satisfied.

---

[5] Arc expansion is a binary decision with fixed cost for all arcs.

$$(29)$$

$$(28) \quad
\begin{aligned}
&\min_{f,\xi} && \sum_{(i,j)\in E} \xi_{ij} \\
&s.t. && \xi_{ij} \geq \frac{f_{ij} - x}{M} \\
& && \sum_{(i,j)\in E} f_{ij} - \sum_{(k,i)\in E} f_{ki} = b_i \\
& && 0 \leq f_{ij} \leq u \\
& && \xi_{ij} \in \{0,1\}.
\end{aligned}
\qquad
\begin{aligned}
&\min_{y,z,a} && \sum_{(i,j)\in E} z_{ij} \\
&s.t. && z_{ij} \geq y_{ij} - ax + 1 \\
& && \sum_{(i,j)\in E} y_{ij} - \sum_{(k,i)\in E} y_{ki} = ab_i \\
& && 0 \leq y_{ij} \leq au \\
& && a \geq 0, z_{ij} \geq 0.
\end{aligned}$$

However, as we show in Example 4, the CUT$_x$ optimal solution can have unappealing properties. Specifically, for the case of capacity planning, we see that the CUT$_x$ optimal flow will likely fill all arcs to capacity (either $u$ or $x$). This means that the network may be very sensitive to changes in demand or arc failures. For example, if one arc fails, it is unlikely that flow will be able to be diverted because all other arcs are filled to capacity. As discussed in [11, 4], networks often require some capacity buffer so that uncertainties do not lead to massive disruptions. Power grids are a specific example considered in [4], where capacities on power lines are often violated to serve all demand when lines may fail or some unplanned flow must be transmitted across a power line. Thus, power grids, if designed without a buffer, can become susceptible to cascading failure, or blackouts.

We consider the similar problem of minimizing CUA$_x$, formulated as (29) where we have again reduced to LP via the change of variable $y_{ij} = af_{ij}$. As we show in Example 4, the CUA$_x$ optimal flow does not suffer from the same drawbacks as the CUT$_x$ optimal flow. By considering the magnitude of components around the threshold $x$, the CUA$_x$ optimal flow is more evenly spread among the arcs.

Of course, there are methods to directly mitigate the risk of cascading failures by introducing new elements into the CUT$_x$ minimization formulation that consider uncertainty. For example, [4] considers a scenario based approach, where flows must satisfy demand for multiple scenarios, while [13] implements a robust optimization approach in a similar fashion. However, introduction of additional components into a MIP makes a numerically difficult optimization problem even more challenging. On the other hand, without having to gather scenario information or consider uncertainty directly, the CUA$_x$ optimal flow is able to mitigate some of the risks associated with the CUT$_x$ optimal flow and may be less susceptible to cascading failures. In addition, it is a simple LP, and uncertainty considerations can still be entered into the formulation just as is done with the CUT$_x$ formulations.

**Example 4:** Consider the network in Figure 4, where we would like to push 24 units of supply through the network from the left black node to the right black node. Assume that each arc has initial capacity $x = 7$, which is not sufficient to push flow through the network, so we must decide which arcs to expand to capacity $u = 10$. We assume that each expansion is a discrete decision that comes with a fixed cost $C$. We see CUA$_x$ and CUT$_x$ optimal flows in Figure 4. First, we can see that if the CUT$_x$ optimal flow utilizes an arc, it likely fills it to capacity (either 7 or 10). In contrast, we notice that the CUA$_x$ optimal flow is more evenly distributed throughout the network with all arcs having some buffer between their actual flow and their capacity (7 or 10). We can measure the susceptibility to arc failure by measuring the expected

amount of lost flow given a single arc failure, assuming equally probable failures[6] and we find that the $\text{CUA}_x$ flow has expected loss equal to 2 and that the $\text{CUT}_x$ flows has expected loss equal to 2.66.

**4.1.4.** $\text{CUT}_x$ **minimization is NP hard.** It should be noted that the MIP formulations we consider for $\text{CUT}_x$ minimization problems are Big-M type formulations. Though there may be more efficient methods for formulating and approximating such a MIP, e.g. by considering a Disjunctive Programming approach, we emphasize that the $\text{CUT}_x$ minimization problem is NP hard. In this section, we prove that for arbitrary graphs, the $\text{CUT}_x$ minimization problem is NP hard. Therefore, regardless of the strategy used to solve (exactly or approximately) the $\text{CUT}_x$ minimization problem, $\text{CUA}_x$ minimization will reduce to convex and linear programming which have polynomial-time solvers.

In Proposition 4 we show that the problem of minimizing $\text{CUT}_x$ for an arbitrary graph is an NP-hard problem with polynomial-time reduction from an NP-complete set covering problem to $\text{CUT}_x$ minimization problem. The NP-hardness implies that the considered problem is at least as hard as any P, NP, or NP-complete problem, that is, time consumed by solving this problem will probably grow exponentially with the size of the problem.

**Proposition 4:** $\text{CUT}_x$ *minimization is an NP-hard problem.*

**Proof:** The set covering problem is NP-complete. Solving a set-covering problem with one run of the $\text{CUT}_x$ minimization problem of the same size will prove that $\text{CUT}_x$ minimization problem is NP-hard. Suppose there is a set $S$ consisting of sets $S_i$: $S = \{S_1, \ldots, S_n\}$, where $S_i \subseteq U \equiv \{1, \ldots, m\} = \cup_{i=1}^n S_i$. It is required to find minimal covering set $S^* = \{S_1^*, \ldots, S_k^*\}$ with $S_k^* \in S$ and minimal $k$. Consider a three-layer graph. On the first layer of the graph there are $m$ demand vertices with demand 1, each vertex corresponds to an object from the union set $U$. On the second layer of the graph there are $n$ transshipment nodes with demand/supply value 0, which correspond to sets $S_i$. The first-layer vertex $j$ and the second-layer vertex $i$ are connected iff $j \in S_i$. The transportation costs along all edges are equal to 1. The third layer of the graph consists of a single auxiliary supply vertex with the supply $m$. $\text{CUT}_x$ is measured on the vector of loads for transshipment vertices which are the same as those presented in Example 2, and the critical threshold is 0. That is, if transshipment vertex is involved in transportation, its load is larger than 0 and it violates the threshold. Finally, for the described graph it can be seen that picking a covering subset with minimal number of sets is equivalent to minimizing the number of overloaded transshipment vertices. $\square$

**5. Conclusion.** In this paper, we have introduced a new concept called $\text{CUA}_x$, a function which counts the number of largest outcomes in a data set which have average value equal to a specified threshold, $x$. As a basic characteristic for counting tail outcomes, $\text{CUA}_x$ not only counts outcomes like $\text{CUT}_x$, but accounts for the severity of these outcomes which $\text{CUT}_x$ does not. We have also shown that $\text{CUA}_x$

---

[6]To solve for the lost flow given an arc failure, we first take the network with the single arc removed where all other arcs have capacity $x$ or $u$ (depending on whether they were expanded or not). We then solve a max flow problem for that network where the output of the supply node is limited to the original demand, which in this example equals 24. The lost flow from dropping the single arc is then 24 minus the max flow for this new network.

has superior mathematical properties, being a continuous function w.r.t. the threshold parameter, piecewise linear in its reciprocal, and directly optimizable via convex and linear programming. Thus, $\text{CUA}_x$ can be efficiently calculated, has continuity properties which can be used for sensitivity analysis, and can be used to efficiently minimize the number of outcomes in the tail of the data set, with the tail including the largest outcomes that average to a specified threshold. We have shown that $\text{CUA}_x$ is the inverse of $\text{UA}_k$, a function that measures the average magnitude of the largest $k$ outcomes in a data set. We have also proved that $\text{CUA}_x$ and $\text{CUT}_x$ are strongly connected, showing that $\text{CUA}_x$ is, in a certain sense, the minimal quasiconvex upper bound of $\text{CUT}_x$.

Finally, we have shown that $\text{CUA}_x$ can be use to formulate new network optimization problems. We compare against similar formulations which minimize $\text{CUT}_x$. In addition to these formulations being less efficiently solvable, since they involved binary variables, we find that the solution can often be unappealing. We show that $\text{CUA}_x$ optimization, on the other hand, reduces to solving a LP and that by accounting for the severity of the largest outcomes, can lead to more appealing network flow policies.

**Appendix A. $\text{CUA}_x$: A special case of bPOE.**

**A.1. bPOE and Tail Probabilities.** When working with optimization of tail probabilities, one frequently works with constraints or objectives involving *probability of exceedance* (POE), $p_x(X) = P(X > x)$, or its associated quantile $q_\alpha(X) = \min\{x | P(X \leq x) \geq \alpha\}$, where $\alpha \in [0,1]$ is a probability level. The quantile is a popular measure of tail probabilities in financial engineering, called within this field Value-at-Risk by its interpretation as a measure of tail risk. The quantile, though, when included in optimization problems via constraints or objectives, is quite difficult to treat with continuous (linear or non-linear) optimization techniques.

A significant advancement was made in [15] in the development of an approach to combat the difficulties raised by the use of the quantile function in optimization. Rockafellar and Uryasev explored a replacement for the quantile, called CVaR within the financial literature, and called the superquantile in a general context. The superquantile is a measure of uncertainty similar to the quantile, but with superior mathematical properties. Formally, the superquantile (CVaR) for a continuously distributed $X$ is expressed as

$$\bar{q}_\alpha(X) = E\left[X | X > q_\alpha(X)\right].$$

For general distributions, the superquantile is defined by the following formula,

$$\bar{q}_\alpha(X) = \min_\gamma \left\{\gamma + \frac{E[X - \gamma]^+}{1 - \alpha}\right\},$$

where $[\cdot]^+ = \max\{\cdot, 0\}$. Similar to $q_\alpha(X)$, the superquantile can be used to assess the tail of the distribution. The superquantile, though, is far easier to handle in optimization contexts. It also has the important property that it considers the magnitude

of events within the tail. Therefore, in situations where a distribution may have a heavy tail, the superquantile accounts for magnitudes of low-probability large-loss tail events while the quantile does not account for this information.

Working to extend this concept, bPOE was developed as the inverse of the superquantile in the same way that POE is the inverse of the quantile. Specifically, there exists two slightly different variants of bPOE, namely Lower and Upper bPOE. Paper [9] defines so-called Lower bPOE in the following way, where $\sup X$ denotes the essential supremum of the random variable $X$.

**Definition (Lower bPOE):** Let $X$ be a real-valued random variable and $x \in \mathbb{R}$ a fixed threshold parameter. Lower bPOE of random variable $X$ at threshold $x$ equals

$$\bar{p}_x^L(X) = \begin{cases} 0, & \text{if } x \geq \sup X \ , \\ \{1 - \alpha | \bar{q}_\alpha(X) = x\}, & \text{if } E[X] < x < \sup X \ , \\ 1, & \text{otherwise.} \end{cases}$$

In words, for any threshold $x \in (E[X], \sup X)$, Lower bPOE can be interpreted as one minus the probability level at which the superquantile equals $x$.

Similarly, paper [12] defines so-called Upper bPOE as follows.

**Definition (Upper bPOE):** Upper bPOE of random variable $X$ at threshold $x$ equals

$$\bar{p}_x^U(X) = \begin{cases} \max\{1 - \alpha | \bar{q}_\alpha(X) \geq x\}, & \text{if } x \leq \sup X \ , \\ 0, & \text{otherwise.} \end{cases}$$

Upper and Lower, in fact, do not differ dramatically. This is shown by the following property, proved in [12].

**Upper vs. Lower bPOE:**

$$\bar{p}_x^U(X) = \begin{cases} \bar{p}_x^L(X), & \text{if } x \neq \sup X \ , \\ P(X = \sup X), & \text{if } x = \sup X. \end{cases}$$

It is important to notice that Upper and Lower bPOE are equivalent when $x \neq \sup X$. The difference between the two definitions arises when threshold $x = \sup X$. In this case, we have that $\bar{p}_x^L(X) = 0$ while $\bar{p}_x^U(X) = P(X = \sup X)$. Thus, for a threshold $x \in (E[X], \sup X)$, both Upper and Lower bPOE of $X$ at $x$ can be interpreted as one minus the probability level at which the superquantile equals $x$. Roughly speaking, Upper bPOE can be compared with $P(X \geq x)$ while Lower bPOE can be compared with $P(X > x)$. To read further about the differences between Upper and Lower bPOE, see [9].

**A.2. CUA$_x$ and Upper bPOE.** In [12], an important calculation formula for bPOE was introduced. Specifically, [12] found that Upper bPOE has the following calculation formula.

**Proposition**: Given a real valued random variable $X$ and a fixed threshold $x$, bPOE for random variable $X$ at $x$ equals

$$(30)$$

$$\bar{p}_x^U(X) = \inf_{\gamma < x} \frac{E[X-\gamma]^+}{x-\gamma} = \begin{cases} \lim\limits_{\gamma \to -\infty} \frac{E[X-\gamma]^+}{x-\gamma} = 1 \,, & \text{if } x \le E[X], \\ \min\limits_{\gamma < x} \frac{E[X-\gamma]^+}{x-\gamma} \,, & \text{if } E[X] < x < \sup X, \\ \lim\limits_{\gamma \to x^-} \frac{E[X-\gamma]^+}{x-\gamma} = P(X = \sup X) \,, & \text{if } x = \sup X, \\ \min\limits_{\gamma < x} \frac{E[X-\gamma]^+}{x-\gamma} = 0 \,, & \text{if } \sup X < x. \end{cases}$$

With this calculation formula, we can then show that $\mathrm{CUA}_x$ is simply a special case of Upper bPOE. First, let us represent our deterministic vector $(y_1, y_2, ..., y_n) = \mathbf{y} \in \mathbb{R}^n$ as a real valued discrete random variable $Y$ taking on values $(y_1, y_2, ..., y_n)$ with equal probabilities, i.e. $P(Y = y_i) = \frac{1}{n}$. Second, let us consider the quantity $n\bar{p}_x(Y)$. Using calculation formula (30) with the change of variable $a = \frac{1}{x-\gamma}$, we see that

$$
\begin{aligned}
n\bar{p}_x^U(Y) &= \min_{a \ge 0} nE[a(Y-x)+1]^+ \\
&= \min_{a \ge 0} \sum_{i=1}^{n} [a(y_i - x) + 1]^+ \,.
\end{aligned}
$$
$$(31)$$

Thus, we see that this is exactly the definition of $\mathrm{CUA}_x$. In other words, we see that $\mathrm{CUA}_x$ is a deterministic variant of bPOE.

## REFERENCES

[1] Ahuja, R. K., Magnanti, T. L., Orlin, J. B. Network flows: theory, algorithms, and applications. (1993)

[2] Bertsekas, D. P. Network optimization: continuous and discrete models. Belmont: Athena Scientific. (1998)

[3] Bienstock, D., Chertkov, M., Harnett, S. Chance-constrained optimal power flow: Risk-aware network control under uncertainty. SIAM Review, 56(3), 461-495.(2014)

[4] Bienstock, D., Mattia, S. Using mixed-integer programming to solve power grid blackout problems. Discrete Optimization, 4(1), 115-141. (2007)

[5] Binato, S., Pereira, M. V. F., Granville, S. A new Benders decomposition approach to solve power transmission network design problems. IEEE Transactions on Power Systems, 16(2), 235-240. (2001)

[6] Bertsimas, D., Pachamanova, D., Sim, M. Robust linear optimization under general norms Operations Research Letters, 32(6), 510-516. (2004)

[7] Davis, J.R., Uryasev S. Analysis of Hurricane Damage using Buffered Probability of Exceedance. Research Report 2014-4, ISE Dept., University of Florida. (2014)

[8] Follmer H., Schied A. Stochastic Finance: an introduction in discrete time. Walter de Gruyter. (2011)

[9] Mafusalov A., Uryasev S. Buffered Probability of Exceedance: Mathematical Properties and Optimization Algorithms. Research Report 2014-1, ISE Dept., University of Florida. (2014)

[10] Mafusalov, A., Uryasev, S. Conditional value-at-risk (CVaR) norm: Stochastic Case. Research Report 2013-5, Department of Industrial Systems and Engineering, University of Florida, Gainesville, FL (2013)

[11] Magnanti, T. L., Wong, R. T. Network design and transportation planning: Models and algorithms. Transportation science, 18(1), 1-55. (1984)

[12] Norton M., Uryasev S. Maximization of AUC and Buffered AUC in Classification Research Report 2014-2, ISE Dept., University of Florida. (2014)

[13] Ordonez, F., Zhao, J. Robust capacity expansion of network flows. Networks, 50(2), 136-145. (2007)

[14] Pavlikov, K., Uryasev S. CVaR norm and applications in optimization. Optimization Letters 8.7 (2014): 1999-2020. (2014)

[15] Rockafellar R.T. and S. Uryasev Optimization of Conditional Value-At-Risk. The Journal of Risk, Vol. 2, No. 3, 2000, 21-41. (2000)

[16] Rockafellar, R.T. Safeguarding Strategies in Risky Optimization. Presentation at the International Workshop on Engineering Risk Control and Optimization, Gainesville, FL, February, 2009.

[17] Rockafellar R.T., Royset J.O. On Buffered Failure Probability in Design and Optimization of Structures. Reliability Engineering & System Safety, Vol. 95, 499-510. (2010)

[18] Uryasev, S. Buffered Probability of Exceedance and Buffered Service Level: Definitions and Properties. Research Report 2014-3, ISE Dept., University of Florida. (2014)

[19] Yager, R. R. On ordered weighted averaging aggregation operators in multicriteria decision making. Systems, Man and Cybernetics, IEEE Transactions on, 18(1), 183-190. (1988)