# A Kolmogorov-distance based approximation of discrete random variables

**Anonymous Author(s)**
Affiliation
Address
email

## Abstract

We present an algorithm that takes a discrete random variable $X$ and a number $m$ and computes a random variable whose support (set of possible outcomes) is of size at most $m$ and whose Kolmogorov distance from $X$ is minimal. In addition to a formal theoretical analysis of the correctness and of the computational complexity of the algorithm, we present a detailed empirical evaluation that shows how the proposed approach performs in practice in different applications and domains.

## 1   Introduction

Many different approaches to approximation of probability distributions are studied in the literature [13, 16, 17]. The approaches vary in the types random variables considered, how they are represented, and in the criteria used for evaluation of the quality of the approximations. This paper is on approximating discrete distributions represented as explicit probability mass functions with ones that are simpler to store and to manipulate. This is needed, for example, when a discrete distribution is given as a large data-set, obtained, e.g., by sampling, and we want to represent it approximately with a small table (see [11] for example).

The main contribution of this paper is an efficient algorithm for computing the best possible approximation of a given random variable with a random variable whose complexity is not above a prescribed threshold, where the measures of the quality of the approximation and the complexity of the random variable are as specified in the following two paragraphs.

We measure the quality of an approximation by the distance between the original variable and the approximate one. Specifically, we use the Kolmogorov distance which is commonly used for comparing random variables in statistical practice and literature. Given two random variables $X$ and $X'$ whose cumulative distribution functions (cdf) are $F_X$ and $F_{X'}$, respectively, the Kolmogorov distance between $X$ and $X'$ is $d_K(X, X') = \sup_t |F_X(t) - F_{X'}(t)|$ (see, e.g., [9]). We say that $X'$ is a good approximation of $X$ if $d_K(X, X')$ is small.

The complexity of a random variable is measured by the size of its support, the number of values that it can take, $|\operatorname{support}(X)| = |\{x \colon Pr(X = x) \neq 0\}|$. When distributions are maintained as explicit tables, as done in many implementations of statistical software, the size of the support of a variable is proportional to the amount of memory needed to store it and to the complexity of the

29  computations around it. In summary, the exact notion of optimality of the approximation targeted in
30  this paper is:

31  **Definition 1.** *A random variable $X'$ is an optimal $m$-approximation of a random variable $X$ if*
32  $|\operatorname{support}(X')| \leq m$ *and there is no random variable $X''$ such that $|\operatorname{support}(X'')| \leq m$ and*
33  $d_K(X, X'') < d_K(X, X')$.

34  The main contribution of the paper is an efficient algorithm that takes $X$ and $m$ as parameters and
35  constructs an optimal $m$-approximation of $X$.

36  The rest of the paper is organized as follows. In Section 2 we describe how our work relates to other
37  algorithms and problems studied in the literature. In Section 3 we detail the proposed algorithm,
38  analyze its properties, and prove the main theorem. In Section 4 we demonstrate how the proposed
39  approach performs on the problem of estimating the probability of hitting deadlines is plans and
40  compare it to alternatives approximation approaches from the literature. We also demonstrate the
41  performance of our approximation algorithm on some randomly generated random variables. The
42  paper is concluded with a discussion in Section 5.

## 2   Related work

44  The most relevant work related to this paper is the papers by Cohen at. al. [4,5]. These papers study
45  approximations of random variables in the context of estimating deadlines. In this context, $X'$ is
46  defined to be a good approximation of $X$ if $F_{X'}(t) > F_X(t)$ for any $t$ and $\sup_t F_{X'}(t) - F_X(t)$
47  is small. This is not a distance because it is not symmetric. The motivation given by Cohen at. al.
48  for using this type of approximation is for cases where overestimation of the probability of missing
49  a deadline is acceptable but underestimation is not. In Section 4, we consider the same examples
50  examined by Cohen at. al. and show how the algorithm proposed in this paper performs relative to
51  the algorithms proposed there when both over- and under- estimations are allowed. As expected, the
52  Kolmogorov distance between the approximation and the original random variable is smaller by a
53  factor of one half, on average, when using the algorithm proposed here.

54  Another relevant prior work is the theory of Sparse Approximation (aka Sparse Representation) that
55  deals with sparse solutions for systems of linear equations, as follows.

56  Given a matrix $D \in \mathbb{R}^{n \times p}$ and a vector $x \in \mathbb{R}^n$, the most studied sparse representation problem is
57  finding the sparsest possible representation $\alpha \in \mathbb{R}^p$ satisfying $x = D\alpha$:

$$\min_{\alpha \in \mathbb{R}^p} \|\alpha\|_0 \text{ subject to } x = D\alpha.$$

58  where $\|\alpha\|_0 = |\{i : \alpha_i \neq 0, i = 1, \ldots, p\}|$ is the $\ell_0$ pseudo-norm, counting the number of non-zero
59  coordinates of $\alpha$. This problem is known to be NP-Hard with a reduction to NP-complete subset
60  selection problems.

61  In these terms, using also the $\ell_\infty$ norm that represents the maximal coordinate and the $\ell_1$ norm that
62  represents the sum of the coordinates, our problem can be phrased as:

$$\min_{\alpha \in [0,\infty)^p} \|x - D\alpha\|_\infty \text{ subject to } \|\alpha\|_0 = m \text{ and } \|\alpha\|_1 = 1.$$

63  where $D$ is the all-ones triangular matrix (the entry at row $i$ and column $j$ is one if $i \leq j$ and zero
64  otherwise), $x$ is related to $X$ such that the $i$th coordinate of $x$ is $F_X(x_i)$ where $\operatorname{support}(X) =$
65  $\{x_1 < x_2 < \cdots < x_n\}$ and $\alpha$ is related to $X'$ such that the $i$th coordinate of $\alpha$ is $f_{X'}(x_i)$. The
66  functions $F_X$ and $f_{X'}$ represent, respectively, the cumulative distribution function of $X$ and the

67 mass distribution function of $X'$. This, of course, means that the coordinates of $x$ are assumed to

68 be positive and monotonically increasing and that the last coordinate of $x$ is assumed to be one. We

69 demonstrate an application for this specific sparse representation problem and show that it can be

70 solve in $O(n^2 m)$ time and $O(m^2)$ memory.

71 The presented work is also related to the research on binning in statistical inference. Consider, for

72 example, the problem of credit scoring [21] that deals with separating good applicants from bad

73 applicants where the Kolmogorov–Smirnov statistic KS is a standard measure. The KS comparison

74 is often preceded by a procedure called binning where a large table is translated to a smaller one

75 by collecting nearby values together. There are many methods for binning [2, 12, 18, 19]. In this

76 context, our algorithm can be consider as a new binning strategy that provides optimality guarantees

77 with respect to the Kolmogorov distance that none of the existing binning technique that we are

78 aware of provides.

79 The present study is also related to the work of Pavlikov and Uryasev [16], where a procedure for

80 producing a random variable $X'$ that optimally approximates a random variable $X$ is presented.

81 Their approximation scheme, achieved using linear programming, is designed for a different notion

82 of distance (called CVaR). The new contribution of the present work in this context is that our

83 method is direct, not using linear programming, thus allowing tighter analysis of time and memory

84 complexity. Also, our method is designed for optimizing the Kolmogorov distance that is more

85 prevalent in applications. For comparison, in Section 4 we briefly discuss the performance of linear

86 programming approach similar to the one proposed in [16] for the Kolmogorov distance and compare

87 it to the algorithm proposed in this paper.

## 88 3 An algorithm for optimal approximation

89 In the scope of this section, let $X$ be a given random variable with a finite support of size $n$, and

90 let $0 < m \leq n$ be a given complexity bound. The section evolves by developing notations and by

91 collecting facts towards an algorithm for finding an optimal $m$-approximation of $X$.

92 The first useful fact is that it is enough to limit our search to approximations $X'$s such that

93 $\text{support}(X') \subseteq \text{support}(X)$:

94 **Lemma 2.** *For every random variable $X''$ there is a random variable $X'$ such that* $\text{support}(X') \subseteq$

95 $\text{support}(X)$ *and* $d_K(X, X') \leq d_K(X, X'')$.

96 *Proof.* Let $\{x_1, \ldots, x_n\} = \text{support}(X)$, and let $x_0 = -\infty, x_{n+1} = \infty$. Consider the random

97 variable $X'$ whose probability mass function is $f_{X'}(x_i) = P(x_{i-1} < X'' \leq x_i)$ for $i = 1, \ldots, n-1$,

98 $f_{X'}(x_n) = P(x_n - 1 < X'' < x_{n+1})$, and $F_{X'}(x) = 0$ if $x \notin \text{support}(X)$. Since $X'$ only "pushes"

99 the probability mass of $X''$ to the support of $X$, we have that $f_{X'}$ is a probability mass function

100 and therefore $X'$ is well defined. By construction, $|F_X(x_i) - F_{X'}(x_i)| = |F_X(x_i) - F_{X''}(x_i)|$

101 for every $1 \leq i \leq n-1$. For $i = n$ we have $|F_X(x_n) - F_{X'}(x_n)| = |1 - 1| = 0$. Since

102 $|F_X(x) - F_{X'}(x)| = |F_X(x_i) - F_{X'}(x_i)|$ for every $0 \leq i < n+1$ and $x_i < x < x_{i+1}$, we have

103 that $d_K(X, X') = max_i |F_X(x_i) - F_{X'}(x_i)| \leq max_i |F_X(x_i) - F_{X''}(x_i)| \leq d_K(X, X'')$. $\square$

104 For a set $S \subseteq \text{support}(X)$, let $\mathbb{X}_S$ denote the set of random variables whose supports are contained

105 in $S$. In Step 1 below, we find a random variable in $\mathbb{X}_S$ that minimizes the Kolmogorov distance

106 from $X$. We denote the Kolmogorov distance between this variable and $X$ by $\varepsilon(X, S)$. Then, in

107 Step 2, we show how to efficiently find a set $S \subseteq \text{support}(X)$ whose size is smaller or equal to $m$

108  that minimizes $\varepsilon(X, S)$. Then, in Step 3, an optimal $m$-approximation is constructed by taking a
109  minimal approximation in $\mathbb{X}_S$ where $S$ is the set that that minimizes $\varepsilon(X, S)$.

**Step 1: Finding an $X'$ in $\mathbb{X}_S$ that minimizes $d_K(X, X')$**

111  We first fix a set $S \subseteq \text{support}(X)$ of size at most $m$, and among all the random variables in
112  $\mathbb{X}_S$ find one with a minimal distance from $X$. Denote the elements of $S$ in increasing order by
113  $S = \{x_1 < \cdots < x_m\}$ and let $x_0 = -\infty$ and $x_{m+1} = \infty$. For every $1 < i \leq m$ let $\hat{x}_i$ be the
114  maximal element of $\text{support}(X)$ that is smaller than $x_i$. Consider the following weight function

115  **Definition 3.** *For $0 \leq i \leq m$ let*

$$w(x_i, x_{i+1}) = \begin{cases} P(x_i < X < x_{i+1}) & \text{if } i = 0 \text{ or } i = m; \\ P(x_i < X < x_{i+1})/2 & \text{otherwise.} \end{cases}$$

116  Note that $P(x_i < X < x_{i+1}) = F_X(\hat{x}_{i+1}) - F_X(x_i)$, a fact that we will use throughout this section.

117  **Definition 4.** *Let $\varepsilon(X, S) = \max\limits_{i=0,\ldots,m} w(x_i, x_{i+1})$.*

118  We first show that $\varepsilon(X, S)$ is a lower bound for the distance between random variable in $\mathbb{X}_S$ and $X$.
119  Then, we present a random variable $X' \in \mathbb{X}_S$ such that $d_K(X, X') = \varepsilon(X, S)$. It then follows that
120  $X'$ is an optimal $m$-approximation random variable among all random variables in $\mathbb{X}_S$.

121  **Proposition 5.** *If $X' \in \mathbb{X}_S$ then $d_K(X, X') \geq \varepsilon(X, S)$.*

122  *Proof.* By definition, for every $0 \leq i \leq m$, $d_K(X, X') \geq \max\{|F_X(\hat{x}_{i+1}) -$
123  $F_{X'}(\hat{x}_{i+1})|, |F_X(x_i) - F_{X'}(x_i)|\}$. Note that $F_{X'}(\hat{x}_{i+1}) = F_{X'}(x_i)$ since the probability values
124  for all the elements not in $S$ are set to 0.

125  If $i = 0$, that is $x_i = -\infty$, we have that $F_X(x_i) = F_{X'}(x_i) = F_{X'}(\hat{x}_{i+1}) = 0$ and therefore
126  $d_K(X, X') \geq |F_X(\hat{x}_{i+1})| = |F_X(\hat{x}_{i+1}) - F_X(x_i)| = P(x_i < X < x_{i+1}) = w(x_i, x_{i+1})$.

127  If $i = m$, that is $x_{i+1} = \infty$, we have that $F_X(\hat{x}_{i+1}) = F_{X'}(\hat{x}_{i+1}) = F_{X'}(x_i) = 1$. and therefore
128  $d_K(X, X') \geq |1 - F_X(\hat{x}_i)| = |F_X(\hat{x}_{i+1}) - F_X(x_i)| = P(x_i < X < x_{i+1}) = w(x_i, x_{i+1})$.

129  Otherwise for every $1 \leq i < m$, we use the fact that $max\{|a|, |b|\} \geq |a - b|/2$ for every $a, b \in \mathbb{R}$,
130  to deduce that $d_K(X, X') \geq 1/2|F_X(\hat{x}_{i+1}) - F_X(x_i) + F_{X'}(x_i) - F_{X'}(\hat{x}_{i+1})|$. So $d_K(X, X') \geq$
131  $1/2|F_X(\hat{x}_{i+1}) - F_X(x_i)| = P(x_1 < X < x_2)/2 = w(x_i, x_{i+1})$.

132  Since $d_K(X, X') \geq w(x_i, x_{i+1})$ for every $0 \leq i \leq m$, the proof follows by the definition of
133  $\varepsilon(X, S)$. $\qquad\square$

134  Next we describe a random variable $X' \in \mathbb{X}_S$ with a distance of $\varepsilon(X, S)$ from $X$. Thus $X'$ is an
135  optimal $m$-approximation among the set $\mathbb{X}_S$. The variable $X'$ is described by its probability mass
136  function:

137  **Definition 6.** *Let $f_{X'}(x_i) = w(x_{i-1}, x_i) + w(x_i, x_{i+1}) + f_X(x_i)$ for $i = 1, \ldots, m$ and $f_{X'}(x) = 0$*
138  *for $x \notin S$.*

139  We first show that $X'$ is a properly defined random variable:

140  **Lemma 7.** *$f_{X'}$ is a probability mass function.*

141  *Proof.* From definition $f_{X'}(x_i) \geq 0$ for every $i$. To see that $\sum_i f_{X'}(x_i) = 1$, we have
142  $\sum_i f_{X'}(x_i) = \sum_i(w(x_{i-1}, x_i) + w(x_i, x_{i+1}) + f_X(x_i)) = \sum_{x_i \in S} f_X(x_i)) + w(x_0, x_1) +$

4

143　$\sum_{0<i<m} 2w(x_i, x_{i+1}) + w(x_m, x_{m+1}) = \sum_{x_i \in S} P(X{=}x_i) + P(x_0{<}X{<}x_1) + \sum_{0<i<m} P(x_i <$

144　$X < x_{i+1}) + P(x_m < X < x_{m+1}) = 1$ since this is the entire support of $X$.　　　　□

145　Note that $X'$ can be constructed in time linear in the size of the support of $X$. Its main property,

146　of course, the distance between the cumulative distribution functions of $X$ and $X'$ are bounded by

147　$w(x_i, x_{i+1})$, as follows:

148　**Lemma 8.** *Let $x \in \text{support}(X)$ and $0 \le i \le m$ be such that $x_i \le x \le x_{i+1}$ then $-w(x_i, x_{i+1}) \le$*

149　$F_X(x) - F_{X'}(x) \le w(x_i, x_{i+1})$.

150　*Proof.* We prove by induction on $0 \le i < m$ .

151　First see that $F_{X'}(j) = 0$ for every $x_0 < j < x_1$ and therefore $F_X(j) - F_{X'}(j) = F_X(j) - 0 \le$

152　$F_X(\hat{x}_1) = F_X(\hat{x}_1) - F_X(x_0) = w(x_0, x_1)$. For $j = x_1$ we have $F_X(x_1) - F_{X'}(x_1) = F_X(\hat{x}_1) +$

153　$f_X(x_1) - (w(x_0, x_1) + w(x_1, x_2) + f_X(x_1) = w(x_0, x_1) + f_X(x_1) - (w(x_0, x_1) + w(x_1, x_2) +$

154　$f_X(x_1)) = -w(x_1, x_2)$.

155　Next assume that $F_X(\hat{x}_i) - F_{X'}(\hat{x}_i) = w(x_{i-1}, x_i)$. Then $F_X(x_i) - F_{X'}(x_i) = F_X(\hat{x}_i) + f_X(x_i) -$

156　$(w(x_{i-1}, x_i) + w(x_i, x_{i+1}) + f_X(x_i)) = w(x_{i-1}, x_i) + f_X(x_i) - (w(x_{i-1}, x_i) + w(x_i, x_{i+1}) +$

157　$f_X(x_i)) = -w(x_i, x_{i+1})$.

158　As before we have that for all $x_i < j < x_{i+1}$, we have $F_X(j) - F_{X'}(j) = F_X(j) - F_{X'}(\hat{x}_{i+1}) \le$

159　$F_X(\hat{x}_{i+1}) - F_{X'}(\hat{x}_{i+1})$. Then $F_X(\hat{x}_{i+1}) - F_{X'}(\hat{x}_{i+1}) = (F_X(x_i) + P(x_i < x < x_{i+1})) -$

160　$F_{X'}(x_i) = -w(x_i, x_{i+1}) + 2w(x_i, x_{i+1}) = w(x_i, x_{i+1})$.

161　Finally for $x_m \le j \le x_{m+1}$ we have that $F_{X'}(x_m) = 1$ therefore $F_X(x_m) - F_{X'}(x_m) = (1 -$

162　$P(x_m < X < x_{m+1})) - 1 = P(x_m < X < x_{m+1}) = w(x_m, x_{m+1})$, and for every $x_m < j <$

163　$x_{m+1}$ we have $F_X(j) - F_{X'}(j) < (1 - P(x_m < X < x_{m+1})) - 1 < -P(x_m < X < x_{m+1})) =$

164　$-w(x_m, x_{m+1})$ as required.　　　　□

165　From Lemma 8, by the definition of $\varepsilon(X, S)$, we then have:

166　**Corollary 9.** $d_K(X, X') = \varepsilon(X, S)$.

167　From Proposition 5 we also have:

168　**Corollary 10.** *$\varepsilon(X, S)$ is the distance between $X$ and the variable closest to it in $\mathbb{X}_S$.*

169　**Step 2: Finding a set $S$ that minimizes $\varepsilon(X, S)$**

170　We proceed to finding an $S$ that minimizes $\varepsilon(X, S)$. To obtain that we use a graph search approach

171　motivated by a method described in [3]. We construct a directed graph with a source and a target in

172　which each source-to-target path of length smaller or equal to $m$ corresponds to a possible support set

173　of the same size, and the weights along that path correspond to the weight as defined in Definition 3.

174　Thus the problem of finding an $S$ that minimizes $\varepsilon(X, S)$ is reduced to the problem of finding a

175　source-to-target path $\vec{p}$ of length smaller or equal to $m$ in that graph such that the maximal weight

176　of an edge in $\vec{p}$ is minimal among all other such maximal edges in all other such paths.

177　More specifically, the vertexes of the graph are $V = \text{support}(X) \cup \{-\infty, \infty\}$ and the edges, $E$, are

178　all the pairs $(x_1, x_2) \in V^2$ such that $x_1 < x_2$. The weight of each edge is as specified in Definition 3.

179　Note that there is a one-to-one correspondence between a set $S \subseteq \text{support}(X)$ of size $m$, and an

180　$-\infty$-to-$\infty$ path $\vec{p}_S$ in $G$, obtained by removing the $-\infty$ and $\infty$ from the path in one way and by

181　adding these elements and the sorting on the other way. With this correspondence the maximal

182　weight of an edge on $\vec{p}_S$ is $\varepsilon(X, S)$. We denote this maximal weight of an edge by $w(\vec{p}_S)$, and

5

183 denote the set of all acyclic $-\infty$-to-$\infty$ paths in $G$ with at most $m$ edges by $paths_m(G, -\infty, \infty)$.

184 Thus, the problem of finding the set $S$ with the minimal $\varepsilon(X, S)$ is now reduced to the problem

185 of finding a path $\vec{p} \in paths_m(G, -\infty, \infty)$ such that $w(\vec{p})$ is minimal among all $\{w(\vec{p'}) \colon \vec{p'} \in$

186 $paths_m(G, -\infty, \infty)\}$. This problem can be solved by a variant of the Bellman-Ford algorithm and

187 by the improved algorithm described in [10].

**Step 3: Constructing the overall algorithm**

189 We combine Step 1 and Step 2 in the following algorithm called KolmogorovApprox (Algorithm 1)

190 that follows naturally from the two steps. Given $X$ and $\mathrm{support}(X)$ we add $x_0, x_{n+1}$ and construct

191 the graph (line 2) as in Step 2. Then we execute a variant of the Bellman-Ford algorithm on $G$ for $m$

192 iterations, or the algorithm proposed in [10], to obtain a path $\vec{p} = (v_0, \ldots, v_{m+1})$ (line 2). Finally

193 we use Definition 6 to construct $X'$ from $\vec{p}$ (lines 4-5).

---

**Algorithm 1:** KolmogorovApprox$(X, m)$

---

1 Construct a weighted graph $G = (V, E)$ where $V = \mathrm{support}(X) \cup \{-\infty, \infty\}$,
$E = \{(x_1, x_2) \in V^2 \colon x_1 < x_2\}$, and the weights are as in Definition 3.

2 Find a path $\vec{p} = (x_0, \ldots, x_{m+1}) \in paths_m(G, -\infty, \infty)\}$ such that
$w(\vec{p}) = \min\{w(\vec{p}) \colon \vec{p} \in paths_m(G, -\infty, \infty)\}$.

3 Return a random variable whose probability mass function is
$f_{X'}(x_i) = w(x_{i-1}, x_i) + w(x_i, x_{i+1}) + f_X(x_i)$ for all $i = 1, \ldots, m$ and zero otherwise.

---

**Theorem 11.** KolmogorovApprox *returns an $m$-optimal-approximation of $X$.*

195 *Proof.* By the construction of $G$ we get that the path $\vec{p}$ obtained in line 4 of KolmogorovApprox

196 describes a set $S$ of support of size at most $m$ for which $\varepsilon(S, X)$ is minimal. Then from Definition

197 6 and Corollary 9 we construct $X'$ in lines 4-5 of KolmogorovApprox such that $d_K(X, X') =$

198 $\varepsilon(X, S)$. Therefore $X'$ is an $m$-approximation among all random variables with support contained

199 in $\mathrm{support}(X)$. Finally from Lemma 2 we have that $X'$ is $m$-approximation among all random

200 variables os support of size at most $m$, thus $X'$ is an $m$-optimal-approximation of $X$. $\square$

201 Finally we analyze the complexity of KolmogorovApprox as follows.

202 **Theorem 12.** *The* KolmogorovApprox$(X, m)$ *algorithm runs in time $O(mn^2)$, using $O(n^2)$ mem-*

203 *ory where $n = |\mathrm{support}(X)|$.*

204 *Proof.* Constructing the graph $G$ as described in Step 2 takes $O(n^2)$ time and memory. Comput-

205 ing the shortest path can be achieved by the algorithm described in [10] in time $O(n^2m)$ and no

206 additional memory allocation. $\square$

## 4 Experimental evaluation

208 All algorithms were implemented in Python and the experiments were executed on a hardware com-

209 prised of an Intel i5-6500 CPU @ 3.20GHz processor and 8GB memory. The algorithms of Cohen

210 at. al. were taken "as is" from in the supplementary material to [5] and [4].

211 **Repetitive support size minimization** One use of support size minimization is when commuta-

212 tions that involve summations of random variables slow due to an exponential growth in the support

213 of convolutions of random variables [5]. A key action in coping with this situation is reduction

of the support size by replacing the summed random variable by an approximation of it that has a smaller support size. Previous work like the work of Cohen at. al. in [4, 5] handle this reduction using weaker or sub-optimal notion of approximation than the one presented here, as discussed in Section 2.

As seen in Section 3, given the size of the reduced support, a single step of KolmogorovApprox guarantees an optimal approximated random variable. However in this setting we need to repetitively use KolmogorovApprox, thus the optimality guarantee of the eventually obtained random variable is lost. In light of this, we decided to test the accuracy of the repetitive-KolmogorovApprox and see how it performs against the tools of [4,5] on their benchmarks. These benchmarks are taken from the area of task trees with deadlines, a sub area of the well-established Hierarchical planning [1, 6, 20].

We estimated the probability for meeting deadlines in plans, as described in [4,5], and experimented with four different methods of approximation. The first two, OptTrim [4] and the Trim [5], are taken from the repository of the authors and are designed for achieving only a one-sided Kolmogorov approximation - a weaker notion of approximation then the Kolmogorov approximation discussed in this work. The third method is a simple sampling scheme also described in [5] and the fourth is our Kolmogorov approximation obtained by the proposed KolmogorovApprox algorithm. The parameters for the different methods were chosen in a compatible way, as proposed in [4]. We ran also an exact computation as a reference to the approximated one in order to calculate the error.

| Task Tree | $M$ | KolmogorovApprox | OptTrim | Trim | Sampling | |
|---|---|---|---|---|---|---|
| | | $m/N=10$ | $m/N=10$ | $\varepsilon \cdot N=0.1$ | $s=10^4$ | $s=10^6$ |
| Logistics ($N=34$) | 2 | 0 | 0 | 0.0019 | 0.007 | 0.0009 |
| | 4 | 0.0024 | 0.0046 | 0.0068 | 0.0057 | 0.0005 |
| DRC-Drive ($N=47$) | 2 | 0.0014 | 0.004 | 0.009 | 0.0072 | 0.0009 |
| | 4 | 0.001 | 0.008 | 0.019 | 0.0075 | 0.0011 |
| Sequential ($N=10$) | 2 | 0.0093 | 0.015 | 0.024 | 0.0063 | 0.0008 |
| | 4 | 0.008 | 0.024 | 0.04 | 0.008 | 0.0016 |

Table 1: Comparison of estimated errors with respect to the reference exact computation on various task trees.

Table 1 shows the results of the case study experiment. The quality of the solutions provided by using the KolmogorovApprox operator are better than those provided by the Trim and OptTrim operators, following the optimality guarantees, but is interesting to see that the quality gaps happen in practice in each of the examined task trees. However, in some of the task trees the sampling method produced better results than the approximation algorithm with KolmogorovApprox. Nevertheless, the approximation algorithm comes with an inherent advantage of providing an exact quality guarantees, as opposed to the probabilistic guarantees provided by sampling.

**Single step support minimization**  In order to better understand the quality gaps in practice between KolmogorovApprox, OptTrim, and Trim, we investigate their relative errors when applied on single random variables with support size $n = 100$, and different support sizes of the resulting random variable approximation ($m$). Note that the size of the error obtained by KolmogorovApprox is optimal with respect to $m$. In each instance of this experiment, a random variable is randomly generated by choosing the probabilities of each element in the support uniformly and then normalizing these probabilities so that they sum to one.

Figure 1 presents the error produced by the above methods. The depicted results are averages over fifty instances of random variables. The curves in the figure show the average error of OptTrim and Trim operators with comparison to the average error of the optimal approximation provided by

KolmogorovApprox as a function of $m$. It is evident from this graphs that increasing the support size of the approximation $m$ reduces the error, as expected, in all three methods. However, the (optimal) errors produced by the KolmogorovApprox are significantly smaller, a half of the error produced by OptTrim and Trim.
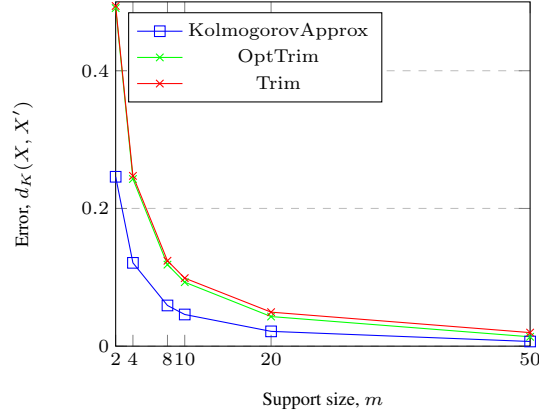


Figure 1: Error comparison between KolmogorovApprox, OptTrim, and Trim, on randomly generated random variables as function of $m$.

**Comparison to Linear Programming**  We also compared the run-time of KolmogorovApprox with a linear programing (LP) algorithm that also guarantees optimality, as described and discussed for example in [16]. For that, we used the "Minimize" function of Wolfram Mathematica as a state-of-the-art implementation of linear programing, encoding the problem by the LP problem $\min_{\alpha \in \mathbb{R}^n} \|x - \alpha\|_\infty$ subject to $\|\alpha\|_0 \leq m$ and $\|\alpha\|_1 = 1$. The run-time comparison results were clear and persuasive: KolmogorovApprox significantly outperforms the LP algorithm. For a random variable with support size $n = 10$ and $m = 5$, the LP algorithm run-time was 850 seconds, where the KolmogorovApprox algorithm run-time was less than a tenth of a second. For $n = 100$ and $m = 5$, the KolmogorovApprox algorithm run-time was 0.14 seconds and the LP algorithm took more than a day. Since it is not trivial to formally analyze the run-time of the LP algorithm, we conclude by the reported experiment that in this case the LP algorithm might not be as efficient as KolmogorovApprox algorithm whose complexity is proven to be polynomial in Theorem 12.

# 5   Discussion and future work

We developed an algorithm for computing optimal approximations of random variables where the approximation quality is measured by the Kolmogorov distance. As demonstrated in the experiments, our algorithm improves on the approach of Cohen, Shimony and Weiss [5] and [4] in that it finds an optimal two sided Kolmogorov approximation, and not just one sided. Beyond the Kolmogorov measure studied here we believe that similar approaches may apply also to total variation, to the Wasserstein distance, and to other measures of approximations. Another direction for future work is extensions to tables that represent other objects, not necessarily random variables. To this end, we need to extend the algorithm to support tables that do not always sum to one and tables that may contain negative entries.

8

## References

[1] R. Alford, V. Shivashankar, M. Roberts, J. Frank, and D. W. Aha. Hierarchical planning: Relating task and goal decomposition with task sharing. In *IJCAI*, pages 3022–3029, 2016.

[2] C. Bolton et al. *Logistic regression and its application in credit scoring*. PhD thesis, Citeseer, 2010.

[3] A. Chakravarty, J. Orlin, and U. Rothblum. A partitioning problem with additive objective with an application to optimal inventory groupings for joint replenishment. *Operations Research*, 30(5):1018–1022, 1982.

[4] L. Cohen, T. Grinshpoun, and G. Weiss. Optimal approximation of random variables for estimating the probability of meeting a plan deadline. In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, New Orleans, Louisiana, USA, February 2-7, 2018*, 2018.

[5] L. Cohen, S. E. Shimony, and G. Weiss. Estimating the probability of meeting a deadline in hierarchical plans. In *IJCAI*, pages 1551–1557, 2015.

[6] T. Dean, R. J. Firby, and D. Miller. Hierarchical planning involving deadlines, travel time, and resources. *Computational Intelligence*, 4(3):381–398, 1988.

[7] K. Erol, J. Hendler, and D. S. Nau. HTN planning: Complexity and expressivity. In *AAAI*, volume 94, pages 1123–1128, 1994.

[8] K. Erol, J. Hendler, and D. S. Nau. Complexity results for HTN planning. *Annals of Mathematics and Artificial Intelligence*, 18(1):69–93, 1996.

[9] J. D. Gibbons and S. Chakraborti. Nonparametric statistical inference. In *International encyclopedia of statistical science*, pages 977–979. Springer, 2011.

[10] R. Guérin and A. Orda. Computing shortest paths for any number of hops. *IEEE/ACM Transactions on Networking (TON)*, 10(5):613–620, 2002.

[11] F. Huq, R. Brannon, and L. Graham-Brady. An efficient binning scheme with application to statistical crack mechanics. *International Journal for Numerical Methods in Engineering*, 105(1):33–62, 2016.

[12] E. Mays. *Handbook of credit scoring*. Global Professional Publishi, 2001.

[13] A. C. Miller and T. R. Rice. Discrete approximations of probability distributions. *Management Science*, 29(3):352–362, 1983.

[14] R. Möhring. Scheduling under uncertainty: Bounding the makespan distribution. *Computational Discrete Mathematics*, pages 79–97, 2001.

[15] D. S. Nau, T.-C. Au, O. Ilghami, U. Kuter, J. W. Murdock, D. Wu, and F. Yaman. SHOP2: An HTN planning system. *Journal of Artificial Intelligence Research*, 20:379–404, 2003.

[16] K. Pavlikov and S. Uryasev. CVaR distance between univariate probability distributions and approximation problems. Technical Report 2015-6, University of Florida, 2016.

[17] A. N. Pettitt and M. A. Stephens. The kolmogorov-smirnov goodness-of-fit statistic with discrete and grouped data. *Technometrics*, 19(2):205–210, 1977.

[18] M. Refaat. *Credit Risk Scorecard: Development and Implementation Using SAS*. Lulu. com, 2011.

[19] N. Siddiqi. *Credit risk scorecards: developing and implementing intelligent credit scoring*, volume 3. John Wiley & Sons, 2012.

[20] Z. Xiao, A. Herzig, L. Perrussel, H. Wan, and X. Su. Hierarchical task network planning with task insertion and state constraints. In *IJCAI*, pages 4463–4469, 2017.

[21] G. Zeng. A comparison study of computational methods of kolmogorov–smirnov statistic in credit scoring. *Communications in Statistics-Simulation and Computation*, 46(10):7744–7760, 2017.