

MARCIN SYDOW* (Warszawa)

Approximation Guarantees for Max Sum and Max Min Facility Dispersion with Parameterised Triangle Inequality and Applications in Result Diversification

Abstract Facility Dispersion Problem, originally studied in Operations Research, has recently found important new applications in Result Diversification approach in information sciences. This optimisation problem consists of selecting a small set of p items out of a large set of candidates to maximise a given objective function. The function expresses the notion of *dispersion* of a set of selected items in terms of a pair-wise *distance* measure between items.

In most known formulations the problem is NP-hard, but there exist 2-approximation algorithms for some cases if distance satisfies triangle inequality.

We present generalised $2/\alpha$ approximation guarantees for the Facility Dispersion Problem in its two most common variants: Max Sum and Max Min, when the underlying dissimilarity measure satisfies *parameterised triangle inequality* with parameter α . The results apply to both relaxed and stronger variants of the triangle inequality.

We also demonstrate potential applications of our findings in the result diversification problem including web search or entity summarisation in semantic knowledge graphs, as well as in practical computations on finite data sets.

2010 Mathematics Subject Classification: 68, 68U35, 68W25.

Key words and phrases: diversity, max sum and max min facility dispersion, approximation algorithms, parameterised triangle inequality.

1. Introduction

The concept of *diversity-awareness* has important practical applications in web search, recommendation, database querying or summarisation (e.g. [6, 12, 18]). The general idea is to return to the user the set of items (being database query or search results, recommended items, etc.) that are not only *relevant* to the query but also *diversified*. The rationale behind such approach is to reduce the risk of *result redundance* and to cover as many different

* Work supported by Polish National Science Centre grant 'DISQUSS' 2012/07/B/ST6/01239

aspects of the query as possible. Equivalently, it is a technique for maximising the likelihood that the user's *unknown intent* behind a potentially *ambiguous query* is satisfied, at least partially.

A possible formulation of the described *diversification problem* is by means of the *Facility Dispersion* optimisation problem that was originally studied in Operations Research (e.g. [7, 8, 14, 16]). More precisely, the problem concerned selecting locations for some dangerous or obnoxious facilities in order to make them *mutually distant* to each other. The range of possible applications is wide, and includes minimising the effects of a terroristic or military attack (if items represent nuclear plants, amunition dumps, etc.) or to avoid self-competition between the stores of the same brand, etc.

In the context of information sciences, the notion of mutual spatial distance has been substituted with that of *pairwise dissimilarity* between the items to be returned to the user.

Examples of recent applications of this approach include web search (e.g. [6, 12]) or graphical summarisation of entities in semantic knowledge graphs [15, 18], for example.

The facility dispersion problem is NP-hard in most common variants, in particular, in Max Min and Max Sum variants studied in this article, and it remains such even when the distance (dissimilarity) function d satisfies triangle inequality [9, 19]. However, in such case, there is a polynomial-time 2-approximation algorithms for this problem [14, 16].

This article focuses on two most common variants of the problem called *Max Sum* and *Max Min Facility Dispersion*. The main results include the generalisations of approximation guarantees of 2 for metric cases of these variants [14, 16] to the value of $2/\alpha$ for the case when the distance function satisfies *parameterised triangle inequality* with parameter α .

1.1. Contributions This paper is an extension of a short paper [17], where the following results were first presented:

- generalised approximation guarantee for Max Sum Dispersion Problem with Parameterised Triangle Inequality (Section 5.1) with the proof (Appendix)
- link between the above result and the Result Diversification Problem that is of interest in web search and other applications (Section 6.2)
- observation concerning satisfying parameterised triangle inequality by finite datasets (Section 6.1)

Compared to the above results, the main extension contained in this paper consists in analogous results concerning another variant of the problem, known as Max Min Facility Dispersion. A list of the extensions compared to [17] includes:

- discussion of the related work concerning Max Sum, Max Min and related problems (new Section 2)
- some remarks on the parameterised triangle inequality (Section 4)
- tight example for the approximation guarantee for Max Sum Dispersion Problem generalised for parameterised triangle inequality (Section 5.2)
- generalised approximation guarantee for Max Min Dispersion Problem with parameterised triangle inequality (Section 5.3)
- generalised impossibility result (i.e. lower-bound approximation factor) for the Max Min Dispersion Problem with parameterised triangle inequality (Section 5.4) with proof (Appendix A)

2. Related Work

2.1. Previous Works on Facility Dispersion Problem Facility Dispersion Problem was studied in many works in Operations Research.

In [16] there are presented approximation results for the Max Min and Max Average (that is equivalent to the Max Sum considered in this paper) variants. Some of these results were improved in [14] for the metric Max Sum variant. To be precise, the latter work concerns the k -dispersion problem that is more general than the one considered here. We focus on the most important case of k -dispersion for the value of parameter $k = 1$, since this case has important recent applications in the result diversification problem in information sciences. Since the proofs would not be affected much for higher values of k , for simplicity we do not consider higher values of the k parameter in [14].

Many more variants of the Facility Dispersion problem (over 10) and approximation algorithms for them were studied in [7] and its extension [8]. It would be interesting in future to analyse the applicability of these variants in the result diversification problem in information sciences, however we are not aware of such work at the time of writing.

Our main results presented here build on the results in the two works [14,16] mentioned above, and the proofs presented in this paper are extensions of the proofs presented in these works and in [13] by appropriately introducing the parameter of the parameterised triangle inequality.

2.2. Related Work on Parameterised Triangle Inequality An impact of various forms of parameterised triangle inequality on the approximation guarantees for various hard optimisation problems was studied previously in many works. One of the early works on this issue is [1]. An especially intensively studied problem in the context of parameterising the triangle inequality is the Travelling Salesman Problem (e.g. [2–4] or a recent example [5]).

Many of these works consider either relaxed or sharpened variant of the triangle inequality. In this paper, we consider the full possible range of the parameter value.

Other forms of relaxing triangle inequality on the performance of algorithms are studied. For example a recent work [11] concerns the impact of triangle inequality violations on the performance of a vehicle routing algorithm.

An example of a practically used pair-wise dissimilarity measure that naturally satisfies parameterised triangle inequality is a variant of the *NEM* (Non-linear Elastic Matching) measure used in pattern matching for comparing visual shapes [10].

2.3. Related Work on Result Diversification Result diversification approach has recently become an intensively studied topic in web search, recommendation, databases or summarisation.

The connection between the Facility Dispersion Problem and the Result Diversification approach was presented in [12] by proposing an appropriate transformation of the dissimilarity function. In this paper, in Section 6.2, we observe that our result holds while this transformation is performed.

In particular, Result Diversification based on Max Sum was recently proposed in [15] in a novel context of diversity-aware entity summarisation in semantic knowledge graphs [18].

Interestingly, the results contained in this paper are related to the question recently stated in the last sentence of [6], where the Max Sum Diversification Problem was studied in a more general framework of monotone submodular functions. The question concerned the impact of relaxing triangle inequality for approximation guarantee for a generalised Max Sum Dispersion Problem. Our result can be viewed as a partial answer to the special case of this problem.

3. Facility Dispersion Problem

In Facility Dispersion Problem, the input consists of a complete, undirected graph $G(V, E)$, an edge-weight function $d : V^2 \rightarrow R^+ \cup \{0\}$ that represents *pairwise distance* between the vertices and a positive natural number $p \leq |V|$. The task is to select a p -element subset $P \subseteq V$ that maximises the objective function $f_d(P)$ that represents the notion of *dispersion* of the elements of the selected set P . In the remaining part of the article we will simplify the notation and use $f(P)$ instead of $f_d(P)$, since the pairwise distance function d will be known from the context.

Depending on the particular form of the objective dispersion function $f(P)$ to be maximised there are considered several variants of the Facility Dispersion Problem. The two most commonly studied variants are Max Sum and Max Min Facility Dispersion and are described in Sections 3.1 and 3.2,

respectively.

3.1. Max Sum, or equivalently, Max Average Facility Dispersion Problem

In the Max Sum Dispersion Problem the objective function to be maximised is defined as follows:

$$f_{SUM}(P) = \sum_{\{u,v\} \subseteq P} d(u,v).$$

This variant is a special case of the more general *k-dispersion* problem that was studied, for example, in [14]. In this paper, we focus only on the special case of *k-dispersion* for $k = 1$, since this case is important for the *Result Diversification* problem discussed in Section 6.2.

In some works (e.g. [13, 16]) the objective function is formulated as:

$$f_{AVE}(P) = \frac{2}{p(p-1)} \sum_{\{u,v\} \subseteq P} d(u,v),$$

and then the problem is known as *Max Average Facility Dispersion*. Since, for the fixed value of p , the formulation is obviously *equivalent* to that of Max Sum Facility Dispersion, we will refer to them interchangeably in this paper.

The Max Sum Dispersion problem is NP-hard even if the distance function d is a metric, but in such case there exists a polynomial-time algorithm of approximation factor of 2 that was presented in [14].

Algorithm 1: HRT: an efficient 2-Approximation Algorithm for Max Average Dispersion Problem

INPUT: An undirected graph $G(V, E)$ with edge-weight function $d : V^2 \rightarrow R^+ \cup \{0\}$, a natural number $1 < p \leq |V|$

OUTPUT: A p -element set $P \subseteq V$

1. $P = \emptyset$
 2. Compute a maximum-weight $\lfloor p/2 \rfloor$ -matching M^* in G
 3. For each edge in M^* , add its both ends to P
 4. In case p is odd, add any node from $V \setminus P$ to P
 5. return P
-

Algorithm 1 shows a heuristic, polynomial-time approximation algorithm, based on computing a maximum-weight matching, that guarantees approximation factor of 2 for Max Average Dispersion when d satisfies triangle

inequality and was presented [14] under the name HRT. A straight-forward implementation of the algorithm makes its time complexity $O(|V|^3)$, however it is possible to implement it so that the time complexity is reduced to $O(|V|^2(p + \log(|V|)))$. There exists another algorithm for the same problem, for which the same approximation guarantee of 2 can be proved, but which has better time complexity of $O(|V|^2 + p^2 \log(p))$. For all the remaining details we refer the reader to [14].

3.2. Max Min Dispersion Problem

In the Max Min Dispersion Problem the objective function to be maximised is defined as follows:

$$f_{MIN}(P) = \min_{u,v \in P, u \neq v} d(u, v). \quad (1)$$

Similarly to Max Sum, even if the distance function d satisfies the standard triangle inequality, the problem is NP-hard, but in such case there exists a polynomial-time algorithm that guarantees an approximation factor of 2.

The algorithm and its 2-approximation guarantee proof is presented in [16] under the name *GMM* and is shown in Algorithm 2. Its time complexity is $O(p^2|V|)$.

Algorithm 2: GMM: an efficient 2-Approximation Algorithm for Max Min Dispersion Problem

INPUT: An undirected graph $G(V, E)$ with edge-weight function $d : V^2 \rightarrow R^+ \cup \{0\}$, a natural number $1 < p \leq |V|$

OUTPUT: A p -element set $P \subseteq V$

1. $P = \operatorname{argmax}_{\{v_i, v_j\} \subseteq V} d(v_i, v_j)$
 2. while ($|P| < p$):
 - find $v \in V \setminus P$ so that $v = \operatorname{argmax}_{v \in V \setminus P} (\min_{u \in P} \{d(v, u)\})$
 - $P = P \cup \{v\}$
 3. return P
-

4. Parameterised Triangle Inequality Assume that V is a non-empty universal set and $d : V^2 \rightarrow R^+ \cup \{0\}$ is a distance function. More precisely, it is assumed that d for all $u, v \in V$ satisfies the properties of *discernibility* ($d(u, v) = 0 \Leftrightarrow u = v$) and *symmetry* ($d(u, v) = d(v, u)$). If, in addition, for all mutually different $u, v, z \in V$, d satisfies the *triangle inequality*: $d(u, v) + d(v, z) \geq d(u, z)$, d is called a *metric*.

Here, we introduce the following definition of *parameterised triangle inequality*:

DEFINITION 4.1 Let V be a set, $\alpha \in \mathbb{R}$, $0 \leq \alpha \leq 2$. A distance function $d : V^2 \rightarrow \mathbb{R}^+ \cup \{0\}$ satisfies *parameterised triangle inequality* (α -PTI) with parameter α iff for all mutually different $u, v, z \in U$:

$$d(u, v) + d(v, z) \geq \alpha d(u, z).$$

α -PTI generalises the standard triangle inequality (for $\alpha = 1$). 0-PTI is the weakest variant, i.e. it is satisfied by any distance function d , and is equivalent to so called *semi-metric*. Observe that the higher the value of the α parameter, the stronger the property.

The value of α cannot be higher than 2:

LEMMA 4.2 *The value of 2 is the highest possible value for α in α -PTI.*

PROOF Assume that some distance function $d : V^2 \rightarrow \mathbb{R}^+ \cup 0$ satisfies α -PTI for $\alpha > 2$ with $|V| \geq 3$. Let u, v, z be three different elements of V . Let introduce the following denotations: $a = d(u, v), b = d(v, z), c = d(z, u)$. We obtain directly from the definition that $a + b > 2c$, that implies $c < \frac{a+b}{2}$.

By summing up the following two inequalities obtained directly from the definition: $b + c > 2a$ and $c + a > 2b$ and then using the above observation, we obtain the following chain of inequalities (that makes a contradiction):

$$2a + 2b < 2c + a + b < 2 \frac{a + b}{2} + a + b = 2a + 2b.$$

(Quod erat demonstrandum) \diamond ■

Notice that since 2-PTI implies that all non-zero distances are equal, that is equivalent to being a *discrete metric* (up to rescaling), the Facility Dispersion Problem becomes trivial for the case $\alpha = 2$.

4.1. PTI vs Relaxed Triangle Inequality In some works (e.g. [10]) an equivalent concept of ρ -relaxed triangle inequality (RTI) is considered in the following form, for some $\rho \in \mathbb{R}^+$:

$$\rho(d(u, v) + d(v, z)) \geq d(u, z).$$

Obviously, such formulation is equivalent to α -PTI with $\rho = 1/\alpha$ with the following observations:

- the case of semi-metric ($\alpha = 0$) is not expressable by ρ -relaxed triangle inequality

- the range of possible values of α is bounded $[0, 2]$, with the metric case being exactly in the middle of that range ($\alpha = 1$). For RTI its parameter value has no upper bound and cannot be smaller than $\frac{1}{2}$.
- for PTI the higher the value of the parameter, the stronger the property, while for the RTI it is the opposite

For the above reasons the α -PTI formulation seems a bit more natural than that of RTI.

5. Improved Approximation Guarantees for α -PTI Case In this section we present main results, i.e. we demonstrate that the approximation guarantees for Max Sum and Max Min Facility Dispersion problems generalise from the value of 2 (for metric case) to the value of $2/\alpha$ when the distance function d satisfies α -PTI for $0 \leq \alpha \leq 2$. The longer proofs are presented in the Appendix.

5.1. $2/\alpha$ Approximation Guarantee for Max Sum Dispersion Satisfying α -PTI

In this subsection we present a theorem that generalises and extends the 2-factor approximation guarantee of the algorithm presented in Algorithm 1 for the case when the function d satisfies α -PTI. The algorithm, under the name HRT, was presented in [14] together with a proof for a metric case. Our proof of this theorem, is an adaptation of the one presented in [13] by properly introducing the α parameter and is presented in the Appendix.

THEOREM 5.1 *Let I be an instance of Max Average Dispersion problem with distance function d satisfying α -PTI for $0 < \alpha < 2$. Let's denote by $OPT(I)$ the value of an optimal solution and by $HRT(I)$ the value of the solution found by the algorithm HRT. It holds that $OPT(I)/HRT(I) < 2/\alpha$.*

As mentioned in Section 3.1, the result applies equivalently to the Max Sum Dispersion problem. For the value of $\alpha = 2$ the problem becomes trivial, since d becomes a discrete metric in such case.

5.2. Tight Example for Max Sum Dispersion We show here that the $2/\alpha$ bound for the algorithm for Max Sum Dispersion with α -PTI presented in Algorithm 1 is (asymptotically) tight.

Let the graph $G(V, E)$, where $|V| \geq 2p$, $M = 2/\alpha$, $m = \beta M = 2\beta/\alpha$, for some $0 < \beta < 1$ (intentionally, close to 1) contain exactly $\lfloor p/2 \rfloor$ edges of weight M and exactly one p -clique, call it C , within which all the edges have weights of m , and all the other edges have weight of 1. Note that α -PTI is satisfied.

The HRT algorithm will select the ends of the $\lfloor p/2 \rfloor$ edges of weight M and, in case p is odd, any arbitrary vertex that brings a total weight of $(p-1)$.

But the optimum, for β being sufficiently close to 1, is the set of vertices of the p -clique C , with the value of $OPT = p(p - 1)m/2$. Thus, we have:

$$OPT/H_{even} = \frac{\frac{p(p-1)}{2}m}{\frac{p(p-1)}{2} + \frac{p}{2}(M-1)} = \frac{m}{1 + \frac{(M-1)}{(p-1)}}$$

$$OPT/H_{odd} = \frac{\frac{p(p-1)}{2}m}{\frac{(p-1)(p-2)}{2} + \frac{(p-1)}{2}(M-1) + (p-1)} = \frac{m}{1 + \frac{(M-1)}{p}}$$

Both the above expressions, for sufficiently large p are arbitrarily close to $m = 2\beta/\alpha$ that is arbitrarily close to $2/\alpha$ for $\beta < 1$ sufficiently close to 1.

In [14] there is presented another than the one presented in Algorithm 1, even simpler, greedy algorithm for solving Max Sum Dispersion Problem. The algorithm, in each iteration, adds to the solution the ends of the heaviest currently available edge and removes all the incident edges, until $2\lfloor p/2 \rfloor$ vertices are collected. If p is odd, at the end it adds one arbitrary vertex to the solution. In [14] it is proven that this algorithm also provides 2-approximation guarantee for Max Sum Dispersion Problem. The above tight example works also for the latter algorithm.

5.3. $2/\alpha$ Approximation Guarantee for Max Min Dispersion Satisfying α -PTI

The following theorem is a generalisation and extension of the result [16] of 2-approximation for Max Min Dispersion satisfying standard triangle inequality.

THEOREM 5.2 *Let I be an instance of the Max Min Dispersion Problem where the distance function d satisfies α -PTI for $0 < \alpha \leq 2$. Let $OPT(I)$ denote the optimum value of the objective function f_{MIN} (see Equation (1)) for this instance and $GMM(I)$ denote the value of the solution found by the GMM algorithm for I (Algorithm 2). Then $OPT(I)/GMM(I) \leq 2/\alpha$, i.e. the GMM algorithm provides $2/\alpha$ -approximation guarantee for this problem.*

The proof is presented in Appendix A.2 and constitutes an our extension of the one presented in [16] by properly introducing the parameter α .

5.4. Lower Bound for Max Min Dispersion with α -PTI For metric cases the 2-factor polynomial approximation algorithm is the best that exists for Max Min Dispersion problem, assuming $P \neq NP$ [16].

Below, we present a generalisation of this to the case of α -PTI.

THEOREM 5.3 *Assume that the distance function d satisfies α -PTI for $0 < \alpha < 2$, and that $\beta < 2/\alpha$ is a real positive number. There is no poly-time algorithm for Max Min Dispersion with approximation factor of β unless $P = NP$.*

PROOF Imagine an instance $(G(V, E), p)$ of the Maximum Independent Set problem, decision version: “Does there exist an independent set of size at least p in the given graph G ”? Let’s set weights of edges in E to $2/\alpha$ and add all the other possible edges of weight 1. A polynomial-time β -approximation algorithm for Max Min Dispersion problem on such modified graph would return the value of 1 iff the answer for the question is positive i.e. solve an NP-hard problem. (*Quod erat demonstrandum*) \diamond

To complete the above result, one can notice that for the value of $\alpha = 2$, as explained in Section 4, all the pairwise distances in V are the same, and all feasible solutions to Max Min Dispersion have the same value. Thus, formally, the guarantee of approximation factor of $2/\alpha = 1$ also holds.

6. Practical Applications of the Results

6.1. α -PTI in Practical Computational Problems In practical applications, the input dataset V is always finite. Here, we make the observation that this fact implies that the distance function d *always* satisfies α -PTI for some $0 < \alpha \leq 2$. To practically find the *actual maximum value* of the α parameter it suffices to check all triples $u, v, z \in V$ in $O(|V|^3)$ time, so that the α -PTI is satisfied as follows: $\alpha = \min_{u, v, z \in V} [d(u, v) + d(v, z)]/d(u, z)$. By this way, it is possible to guarantee a constant-factor of the approximation algorithms for Max Sum and Max Min Dispersion problems even if no other theoretical properties about the distance function are known. In particular, in the metric case, it is *always* possible to guarantee better than 2 approximation factor in practical applications (unless there are no “degenerated” triangles in the data).

In the *on-line* variant of the considered problems (i.e. when data comes in time), α can be systematically updated while data comes.

6.2. Applications to the Result Diversification Problem in Web Search, etc.

In this section, it is demonstrated how the theoretical results from Section 5.1 impact some important recent applications in information sciences including web search and others.

In web search, the problem known as *Result Diversification Problem* can be specified as follows [12]. There is given a set V of documents that are potentially relevant to a user query, a number $p \in N^+$, $p < |V|$, a *document relevance* function $w : V \rightarrow R^+$ and a *document pairwise dissimilarity* function $d : V^2 \rightarrow R^+ \cup \{0\}$. The task in this problem is to select a subset $P \subseteq V$ that maximises the properly defined *diversity-aware set relevance function*. For example, in [12] the following objective function (to be maximised) is proposed as the diversity-aware relevance function:

$$f_{div-sum}(\lambda, P) = (p - 1) \sum_{v \in P} w(v) + 2\lambda \sum_{\{u,v\} \subseteq P} d(u, v). \tag{2}$$

The $\lambda \in R^+ \cup \{0\}$ parameter controls the balance between the relevance term and the diversity term.

The same work explains that by a proper modification of d to d' (see Equation (3)) it is possible to make the described problem of maximising $f_{div-sum}(\lambda, P)$ equivalent to maximising $\sum_{\{u,v\} \subseteq P} d'_\lambda(u, v)$, where:

$$d'_\lambda(u, v) = w(u) + w(v) + 2\lambda d(u, v). \tag{3}$$

In this way, the *result diversification* problem described above is equivalent to the Max Sum Dispersion problem for d'_λ .¹

It is also claimed in [12] that d' is a metric if d is such.²

Notice that, due to the observation in Section 6.1 the following concluding Lemma 6.1 extends the application of the results presented in this paper to the *result diversification* problem on any finite datasets.

LEMMA 6.1 *If distance function d satisfies α -PTI, for some $0 < \alpha \leq 1$, then the modified distance function d'_λ defined in Equation (3) also satisfies α -PTI.*

PROOF $d'_\lambda(u, v) + d'_\lambda(v, z) = 2w(v) + w(u) + w(z) + 2\lambda(d(u, v) + d(v, z)) \geq$
 $\geq 2w(v) + w(u) + w(z) + 2\lambda\alpha d(u, z) \geq$
 $\geq 2w(v) + \alpha[w(u)/\alpha + w(z)/\alpha + 2\lambda d(u, z)] \geq \alpha d'_\lambda(u, z)$. (*Quod erat demonstrandum*)
 ◇ ■

The diversification problem has other interesting applicaitons than in web search. For example, the problem of optimising the $f_{div-sum}(\lambda, P)$ objective

¹Considering the Max Min Dispersion Problem, [12] also considers a variant of bi-criteria objective function: $f_{div-min}(\lambda, P) = \min_{u \in P} w(u) + \lambda \min_{u, v \in P} d(u, v)$ and the authors seem to suggest that maximising it is equivalent to the Max Min Dispersion by defining a modified distance function $d'_\lambda(u, v) = (w(u) + w(v))/2 + \lambda d(u, v)$. But this is, unfortunately, not true, in general. To see this consider the following example: $V = \{u, v, z\}, w(u) = 1, w(v) = 1, w(z) = 20, d(u, v) = 10, d(v, z) = d(z, u) = 1, p = 2, \lambda = 1$. In this example, $f_{div-min}(\lambda, P)$ is maximised for $P = \{u, v\}$ but $f'_{div-min}(\lambda, P) = \min_{u \neq v \in P} d'_\lambda(u, v)$ is maximised for $P = \{u, z\}$ or $P = \{v, z\}$. Furthermore, it is necessary to make the discernibility property satisfied by explicitly setting $d'_\lambda(u, v) = 0$ for $u = v$ since this condition is necessary in the proof of approximation factor of 2 that authors of [12] refer to (see the proof of Theorem 5.2, in the Appendix, more precisely, the property of non-emptiness of S_i^*).

²Formally, this is not true in the form proposed in [12] because the discernibility property would be not satisfied by d' (see Equation (3)). A simple explicit addition of $d'(u, v) = 0$ for $u = v$, however, makes this claim correct.

function has been recently adapted to the novel context of *diversified entity summarisation in knowledge graphs* [15,18].

7. Conclusions and Future Work In this work we gathered some basic theoretical results concerning approximation guarantees for metric cases of Max Min and Max Sum Dispersion problems and generalised them to the case of parameterised triangle inequality.

Despite the results were obtained mostly by simple extensions of the existing proofs, it was demonstrated that they may have additional practical impact on computations on real, finite datasets and, in particular, on recent important diversity-related applications in information sciences as in [12] or [15], for example.

α -PTI property not always affects the guarantees of approximation algorithms in the same way as in this paper (see [1] for example).

In the context of practical computations, it would be interesting to study which distance metrics used in practical applications (such as the measure described in [10]) satisfy α -PTI and what is the value of the parameter.

A. Proofs In this Appendix we present two longer proofs of theorems presented in previous sections.

A.1. Proof of Theorem 5.1 from Section 5.1

The following, including the Lemmas, Theorem and their proofs constitute our extensions of the versions presented in [13][pp. 38-8-38-9] (earlier variants were in [14]). The extensions presented here consist mostly in properly introducing the parameter α .

Let us introduce some denotations. Let $V' \subseteq V$ be a non-empty subset of vertices. Let $G(V')$ denote the complete graph induced on V' and $W(V'), W'(V')$ denote the total weight and average weight of edges in $G(V')$ respectively. By analogy, for a non-empty subset $E' \subseteq E$ of edges, let denote by $W(E')$ and $W'(E') = W(E')/|E'|$ the total and average weight of the edges in E' , respectively. We use the following technical Lemma, presented in [14].

LEMMA A.1 *If $V' \subseteq V$ is a subset of vertices of cardinality at least $p \geq 2$ and M'^* is a maximum-weight $\lfloor p/2 \rfloor$ -matching in $G(V')$, then $W'(V') \leq W'(M'^*)$.*

A very short proof, presented in [14] does not assume *anything* on the distance function d , so that we omit it here.

LEMMA A.2 *Assume that the distance function d satisfies α -PTI for some $0 < \alpha < 2$. If $V' \subseteq V$ is a subset of $p \geq 2$ vertices and M is any $\lfloor p/2 \rfloor$ -matching in $G(V')$, then $W'(V') > (\alpha/2)W'(M)$.*

PROOF (of Lemma A.2) Let $M = \{\{a_i, b_i\} : 1 \leq i \leq \lfloor p/2 \rfloor\}$ and let denote by V_M the set of all vertices that are ends of the edges in M . For each edge $\{a_i, b_i\} \in M$ let E_i denote the set of edges in $G(V')$ that are incident on a_i or b_i , except the edge $\{a_i, b_i\}$ itself. From α -PTI we get that for any vertex $v \in V_M \setminus \{a_i, b_i\}$ we have $d(v, a_i) + d(v, b_i) \geq \alpha d(a_i, b_i)$. After summing this inequality over all the vertices in $V_M \setminus \{a_i, b_i\}$ we obtain:

$$W(E_i) \geq \alpha(p - 2)d(a_i, b_i). \tag{4}$$

There are two cases:

Case 1: p is even, i.e. $\lfloor p/2 \rfloor = p/2$. After summing up the Inequality 4 above, over all the edge sets E_i , $1 \leq i \leq p/2$, we obtain, on the left-hand side, each edge of $G(V')$ twice, except those in M . Thus, $2[W(V') - W(M)] \geq \alpha(p - 2)W(M)$. If we substitute in the last inequality $W(V') = W'(V')p(p - 1)/2$ and $W(M) = W'(M)p/2$, and divide both sides by p , we can quickly get to $W'(V') \geq (\alpha/2)W'(M)[p - 2 + (2/\alpha)]/(p - 1)$, that is equivalent to $W'(V') > (\alpha/2)W'(M)$ (for $\alpha < 2$). This completes the proof for the Case 1.

Case 2: p is odd, i.e. $\lfloor p/2 \rfloor = (p - 1)/2$. Let x be the only node in $V' \setminus V_M$ and let E_x denote the set of all edges incident on x in $G(V')$. By α -PTI we get:

$$W(E_x) \geq \alpha W(M). \tag{5}$$

Let's again sum up the previous Inequality (4) over all the edges E_i , $1 \leq i \leq \lfloor p/2 \rfloor$. On the left-hand side, each edge in $G(V')$ occurs twice, except the edges in M (that do not occur at all) and the edges in E_x that occur once, each. Thus, $2[W(V') - W(M)] - W(E_x) \geq \alpha(p - 2)W(M)$. Now, applying the Inequality (5), we obtain $2[W(V') - W(M)] \geq \alpha(p - 1)W(M)$. If we now substitute $W(V') = W'(V')p(p - 1)/2$ and $W(M) = W'(M)(p - 1)/2$ and divide both sides by $(p - 1)/2$ we will quickly obtain that $W'(V') \geq (\alpha/2)W'(M)[p - 1 + (2/\alpha)]/p$ that is equivalent to $W'(V') > (\alpha/2)W'(M)$ (for $0 < \alpha < 2$). This completes the Case 2 and the whole proof of the Lemma. (*Quod erat demonstrandum*) \diamond ■

The following proof of Theorem 5.1 is an extension of the one proposed in [14] (and later presented in [13]) by properly introducing the parameter α .

PROOF (of Theorem 5.1 from Section 5.1) Let P^* and P be the set of nodes in an optimal solution and that in the solution returned by the HRT algorithm for instance I , respectively. By definition, $OPT(I) = W'(P^*)$ and $HRT(I) = W'(P)$. Let M^* and M denote a maximum-weight $\lfloor p/2 \rfloor$ -matching in P^* and in P , respectively. By Lemma A.1, we get:

$$OPT(I) \leq W'(M^*). \tag{6}$$

In addition, from Lemma A.2 we get:

$$HRT(I) > (\alpha/2)W'(M). \tag{7}$$

Now, because the algorithm *HRT* finds a maximum-weight $\lfloor p/2 \rfloor$ -matching in G , we get $W'(M) \geq W'(M^*)$. This, together with the Inequality (6) and Inequality (7) implies that $HRT(I) > (\alpha/2)W'(M) \geq (\alpha/2)W'(M^*) \geq OPT(I)/\frac{2}{\alpha}$ that completes the proof of the theorem. (*Quod erat demonstrandum*) \diamond ■

A.2. Proof of Theorem 5.2 from Section 5.3 The proof constitutes an extension of the one presented in [16] by properly introducing the parameter α .

Let define $\lambda = \frac{2}{\alpha}$, where α is the value of the parameter in α -PTI satisfied by the distance function d . Let P denote the set-valued variable used in GMM presented in Algorithm 2. By induction on size of P we will show the following condition:

$$f_{MIN}(P) \geq OPT(I)/\lambda \tag{8}$$

holds after each addition to P . After the last addition to P , $GMM(I) = f_{MIN}(P)$ holds, that would imply the theorem.

Initially, the condition holds due to adding two vertices joined by the heaviest edge in V to P . Let's make the inductive assumption that the condition holds after k additions to P , for some $1 \leq k < p - 1$ (notice that after k additions P contains $k + 1$ elements). It will be proven that the condition also holds after the $(k + 1)$ -th addition to P .

Let $P^* = \{v_1^*, \dots, v_p^*\}$ denote an optimal solution to I and let $l^* = OPT(I)$.

Observation 1: $d(v_i^*, v_j^*) \geq l^*$ for any $i \neq j$, due to the minimality of l^* .

Let $P_k = \{x_1, x_2, \dots, x_{k+1}\}$ denote the set P after k additions for $1 \leq k < p - 1$. Because GMM adds at least one more node to P the following holds:

Observation 2: For $1 \leq k < p - 1$, $|P_k| = k + 1 < p$.

Let, for every $v_i^* \in P^*$, define $S_i^* = \{u \in V \mid d(v_i^*, u) < l^*/\lambda\}$. Notice that for any $1 \leq i \leq p$, $S_i^* \neq \emptyset$, since $v_i^* \in S_i^*$ due to the discernibility property of the distance function d (see Section 4).

Furthermore, for any $1 \leq i < j \leq p$, S_i^* and S_j^* are disjoint. To prove this, assume the opposite, i.e. that $S_i^* \cap S_j^* \neq \emptyset$ for some $i \neq j$. Let $u \in S_i^* \cap S_j^*$. Thus $d(v_i^*, u) < l^*/\lambda$ and $d(v_j^*, u) < l^*/\lambda$. This, together with α -PTI implies $\alpha d(v_i^*, v_j^*) \leq d(v_i^*, u) + d(v_j^*, u) < 2l^*/\lambda = \alpha l^*$ that is equivalent to $d(v_i^*, v_j^*) < l^*$. But this would contradict the Observation 1.

Thus $\mathcal{S} = \{S_i^*\}_{1 \leq i \leq p}$ constitutes a family of p non-empty and pair-wise disjoint sets. Thus, since for $k < p - 1$, P_k has strictly less than p elements

and there are p nonempty sets in family \mathcal{S} there must be at least one index r , $1 \leq r \leq p$ such that $S_r^* \cap P_k = \emptyset$.

Due to the definition of S_r^* , for each $u \in P_k$, $d(v_r^*, u) \geq l^*/\lambda$. Due to the fact that v_r^* is available for selection in the $k + 1$ -th step of GMM and GMM will select a vertex $v \in V \setminus P_k$ that maximises $\min_{v' \in P_k} d(v, v')$ among all the vertices in $V \setminus P_k$ it is implied that the condition (8) still holds after the $(k + 1)$ -th addition to P (having made the inductive assumption). (*Quod erat demonstrandum*) \diamond

REFERENCES

- [1] Thomas Andreae and Hans-Jurgen Bandelt. Performance guarantees for approximation algorithms depending on parametrized triangle inequalities. *SIAM Journal of Discrete Mathematics*, 8:1–16, 1995. doi: [10.1137/S0895480192240226](https://doi.org/10.1137/S0895480192240226)
- [2] Michael A. Bender and Chandra Chekuri. Performance guarantees for the tsp with a parameterized triangle inequality. *Information Processing Letters*, 73(1-2):17 – 21, 2000. doi: [10.1016/S0020-0190\(99\)00160-X](https://doi.org/10.1016/S0020-0190(99)00160-X)
- [3] Markus Bläser. *An improved approximation algorithm for the asymmetric tsp with strengthened triangle inequality* In *Automata, Languages and Programming, volume 2719 of Lecture Notes in Computer Science*, JosC.M. Baeten, JanKarel Lenstra, Joachim Parrow, and GerhardJ. Woeginger, eds., pages 157–163. Springer Berlin Heidelberg, 2003. doi: [10.1016/j.jda.2005.07.004](https://doi.org/10.1016/j.jda.2005.07.004)
- [4] Hans-Joachim Böckenhauer, Juraj Hromkovic, Ralf Klasing, Sebastian Seibert, and Walter Unger. Approximation algorithms for the TSP with sharpened triangle inequality. *Information Processing Letters*, 75(3):133 – 138, 2000. doi: [10.1016/S0020-0190\(00\)00089-2](https://doi.org/10.1016/S0020-0190(00)00089-2)
- [5] Hans-Joachim Böckenhauer, Tobias Mömke, and Monika Steinová. Improved approximations for tsp with simple precedence constraints. *J. of Discrete Algorithms*, 21:32–40, 2013. doi: [10.1016/j.jda.2013.04.002](https://doi.org/10.1016/j.jda.2013.04.002)
- [6] Allan Borodin, Hyun Chul Lee, and Yuli Ye. *Max-sum diversification, monotone submodular functions and dynamic updates* In *Proceedings of the 31st symposium on Principles of Database Systems, PODS '12*, pages 155–166, New York, NY, USA, 2012. ACM. doi: [10.1145/2213556.2213580](https://doi.org/10.1145/2213556.2213580)
- [7] Barun Chandra and Magnus M. Halldorsson. *Facility dispersion and remote subgraphs* In *Algorithm Theory - SWAT'96, volume 1097 of Lecture Notes in Computer Science*, Rolf Karlsson and Andrzej Lingas, eds, pages 53–65. Springer Berlin Heidelberg, 1996. doi: [10.1007/3-540-61422-2_120](https://doi.org/10.1007/3-540-61422-2_120)
- [8] Barun Chandra and Magnus M Halldorsson. Approximation algorithms for dispersion problems. *Journal of Algorithms*, 38(2):438 – 465, 2001. doi: [10.1006/jagm.2000.1145](https://doi.org/10.1006/jagm.2000.1145)
- [9] Erhan Erkut. The discrete p-dispersion problem. *European Journal of Operational Research*, 46(1):48 – 60, 1990. doi: [10.1016/0377-2217\(90\)90297-O](https://doi.org/10.1016/0377-2217(90)90297-O)
- [10] Ronald Fagin and Larry Stockmeyer. Relaxing the triangle inequality in pattern matching. *Int. J. Comput. Vision*, 30(3):219–231, 1998. doi: [10.1023/A:1008023416823](https://doi.org/10.1023/A:1008023416823)
- [11] Christopher L. Fleming, Stanley E. Griffis, and John E. Bell. The effects of triangle inequality on the vehicle routing problem. *European Journal of Operational Research*, 224(1):1 – 7, 2013. doi: [10.1016/j.ejor.2012.07.005](https://doi.org/10.1016/j.ejor.2012.07.005)

- [12] Sreenivas Gollapudi and Aneesh Sharma. *An axiomatic approach for result diversification* In *Proceedings of the 18th international conference on World wide web, WWW '09*, pages 381–390, New York, NY, USA, 2009. ACM. doi: [10.1145/1526709.1526761](https://doi.org/10.1145/1526709.1526761)
- [13] Teofilo F. Gonzalez. *Handbook of approx. algorithms and metaheuristics*. CRC Press, 2007. MR [2307955](https://doi.org/10.1137/072307955)
- [14] Refael Hassin, Shlomi Rubinstein, and Arie Tamir. Approximation algorithms for maximum dispersion. *Oper. Res. Lett.*, 21(3):133–137, 1997. doi: [10.1016/S0167-6377\(97\)00034-5](https://doi.org/10.1016/S0167-6377(97)00034-5)
- [15] Witold Kosiński, Tomasz Kuśmierczyk, Pawel Rembelski, and Marcin Sydow. *Application of ant-colony optimisation to compute diversified entity summarisation on semantic knowledge graphs* In *Proc. of International IEEE AIAA 2013/FedCSIS Conference, Annals of Computer Science and Information Systems, volume 1*, pages 69–76, 2013.
- [16] S.S.Ravi, D.J.Rosenkrantz, and G.K.Tayi. Heuristic and special case algorithms for dispersion problems. *Operations Research*, 42(2):299–310, 1994. doi: [10.1287/opre.42.2.299](https://doi.org/10.1287/opre.42.2.299)
- [17] Marcin Sydow. *Improved approximation guarantee for max sum diversification with parameterised triangle inequality* In *Foundations of Intelligent Systems, volume 8502 of Lecture Notes in Computer Science*, Troels Andreassen, Henning Christiansen, Juan-Carlos Cubero, and Zbigniew Ras, eds, pages 554–559. Springer International Publishing, 2014. doi: [10.1007/978-3-319-08326-1_60](https://doi.org/10.1007/978-3-319-08326-1_60)
- [18] Marcin Sydow, Mariusz Pikula, and Ralf Schenkel. The notion of diversity in graphical entity summarisation on semantic knowledge graphs. *Journal of Intelligent Information Systems*, 41:109–149, 2013. doi: [10.1007/s10844-013-0239-6](https://doi.org/10.1007/s10844-013-0239-6)
- [19] D. W. Wang and Yue-Sun Kuo. A study on two geometric location problems. *Inf. Process. Lett.*, 28(6):281–286, August 1988. doi: [10.1016/0020-0190\(88\)90174-3](https://doi.org/10.1016/0020-0190(88)90174-3)

Gwarancje współczynnika aproksymacji dla problemów Max Sum (i Max Min) Facility Dispersion z parametryzowaną nierównością trójkąta i zastosowania w dywersyfikacji wyników

Marcin Sydow

Streszczenie Problem „Facility Dispersion”, pierwotnie studiowany w badaniach operacyjnych, znajduje od niedawna nowe ważne zastosowania w podejściu polegającym na dywersyfikacji wyników w naukach informacyjnych.

Jest to problem optymalizacji dyskretnej polegający na wyborze niewielkiego zbioru p elementów z pewnego dużego zbioru kandydatów tak, aby zmaksymalizować pewną funkcję celu. Funkcja ta wyraża „rozproszenie” wybranych elementów, za pośrednictwem pomocniczej miary odległości par elementów. Problem jest NP-trudny w większości znanych wariantów, lecz istnieją algorytmy aproksymacyjne o współczynniku 2 dla niektórych z nich, gdy miara odległości jest metryką.

W artykule zaprezentowano twierdzenia, które uogólniają znane wyniki do przypadku gdy miara odległości spełnia parametryzowaną nierówność trójkąta z parametrem α , dla wariantów „Max Sum” oraz „Max Min” problemu. Wyniki dotyczą zarówno osłabionej jak i wzmocnionej nierówności trójkąta.

Zademonstrowano także potencjalne zastosowania powyższych rezultatów w problemie dywersyfikacji wyników w takich dziedzinach jak wyszukiwanie informacji czy podsumowania encyj w semantycznych grafach wiedzy, jak również w praktycznych obliczeniach na skończonych zbiorach danych.

2010 *Klasyfikacja tematyczna AMS (2010)*: 68, 68U35, 68W25.

Słowa kluczowe: dywersyfikacja, problem dyspersji, parametryzowana nierówność trójkąta.



Dr. Marcin Sydow received M.Sc. (Mathematics) from Warsaw University and Ph.D (Computer Science) from Polish Academy of Sciences. He heads the Web Mining Lab and the Chair of Intelligent Systems, Algorithms and Mathematics at Polish-Japanese Institute of IT and is an Assistant Professor at Polish Academy of Sciences, Institute of Computer Science, Warsaw Poland. His research interests include Web Search, Web Mining, algorithms, elements of AI and Natural Language Processing. He has published about 60 scientific publications and serves as a PC member and reviewer in many international conferences and journals. His recent research interests include the concept of diversity in information sciences and novel information-processing algorithms for semantic knowledge graphs.

MARCIN SYDOW
POLISH ACADEMY OF SCIENCE
INSTITUTE OF COMPUTER SCIENCE, ORDONA 21, 01-237 WARSZAWA, POLAND
E-mail: msyd@ipipan.waw.pl
POLISH-JAPANESE INSTITUTE OF INFORMATION TECHNOLOGY
WEB MINING LAB, KOSZYKOWA 86, 02-008 WARSZAWA, POLAND
E-mail: msyd@poljap.edu.pl

Communicated by: Jacek Koronacki

(Received: 17th of April 2014; revised: 24th of August 2014)