

Multiwinner Elections with Diversity Constraints

Robert Brederick

University of Oxford, Oxford, UK;
TU Berlin, Berlin, Germany
robert.bredereck@tu-berlin.de

Piotr Faliszewski

AGH University, Krakow, Poland
faliszew@agh.edu.pl

Ayumi Igarashi

University of Oxford, Oxford, UK
ayumi.igarashi@cs.ox.ac.uk

Martin Lackner

TU Wien, Vienna, Austria
lackner@dbai.tuwien.ac.at

Piotr Skowron

TU Berlin, Berlin, Germany
p.k.skowron@gmail.com

Abstract

We develop a model of multiwinner elections that combines performance-based measures of the quality of the committee (such as, e.g., Borda scores of the committee members) with diversity constraints. Specifically, we assume that the candidates have certain attributes (such as being a male or a female, being junior or senior, etc.) and the goal is to elect a committee that, on the one hand, has as high a score regarding a given performance measure, but that, on the other hand, meets certain requirements (e.g., of the form “at least 30% of the committee members are junior candidates and at least 40% are females”). We analyze the computational complexity of computing winning committees in this model, obtaining polynomial-time algorithms (exact and approximate) and NP-hardness results. We focus on several natural classes of voting rules and diversity constraints.

1 Introduction

We study the problem of computing committees (i.e., sets of candidates) that, on the one hand, are of high quality (e.g., consist of high-performing individuals) and that, on the other hand, are diverse (as specified by a set of constraints). The following example shows our problem in more concrete terms.

Consider an organization that wants to hold a research meeting on some interdisciplinary topic such as, e.g., “AI and Economics.” The meeting will take place in some secluded location and only a certain limited number of researchers can attend. How should the organizers choose the researchers to invite? If their main criterion were the number of highly influential AI/economics papers that each person published, then they would likely end up with a very homogeneous group of highly-respected AI professors. Thus, while this criterion definitely should be important, the organizers might put forward additional constraints. For example, they could require that at least 30% of the attendees are junior researchers, at least 40% are female, at least a few economists are invited (but only senior ones), the majority of attendees work on AI, and the attendees come from at least 3 continents and represent at least 10 different countries.¹ In other words, the organizers

would still seek researchers with high numbers of strong publications, but they would give priority to making the seminar more diverse (indeed, junior researchers or representatives of different subareas of AI can provide new perspectives; it is also important to understand what people working in economics have to say, but the organizers would prefer to learn from established researchers and not from junior ones).

The above example shows a number of key features of our committee-selection model. First, we assume that there is some function that evaluates the committees (we refer to it as the *objective function*). In the example it was (implicitly) the number of high-quality papers that the members of the committee published. In other settings (e.g., if we were shortlisting job candidates) these could be aggregated opinions of a group of voters (the recruitment committee, in the shortlisting example).

Second, we assume that each prospective committee member (i.e., each researcher in our example) has a number of attributes, which we call labels. For example, a researcher can be *junior* or *senior*, a *male* or a *female*, can *work in AI* or in *economics* or in some other area, etc. Further, the way in which labels are assigned to the candidates may have a structure on its own. For example, each researcher is either male or female and either junior or senior, but otherwise these attributes are independent (i.e., any combination of gender and seniority level is possible). Other labels may be interdependent and may form hierarchical structures (e.g., every researcher based in Germany is also labeled as representing Europe). Yet other labels may be completely unstructured; e.g., researchers can specialize in many subareas of AI, irrespective how (un)related they seem.

Third, we assume that there is a formalism that specifies when a committee is *diverse*. In principle, this formalism could be any function that takes a committee and gives an *accept/reject* answer. However, in many typical settings it suffices to consider simple constraints that regard each label separately (e.g., “at least 30% of the researchers are junior” or “the number of male researchers is even”). We focus on such independent constraints, but studying more involved ones, that regard multiple labels (e.g., “all invited economists must be senior researchers”) would also be interesting.

Our goal is to find a committee of a given size k that is diverse and has the highest possible score from the objective function. While similar problems have already been consid-

Copyright © 2018, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

¹For example, the Leibniz-Zentrum für Informatik that runs Dagstuhl Seminars gives similar suggestions to event organizers.

ered (see the Related Work section), we believe that our paper is the first to systematically study the problem of selecting a diverse committee, where diversity is evaluated with respect to candidate attributes. We provide the following main contributions:

1. We formally define the general problem of selecting a diverse committee and we provide its natural restrictions. Specifically, we focus on the case of submodular objective functions (with the special case of separable functions), candidate labels that are either layered or laminar,² and constraints that specify sets of acceptable cardinalities for each label independently (with the special case of specifying intervals of acceptable values).
2. We study the complexity of finding a diverse committee of a given size, depending on the type of the objective function, the type of the label structure, and the type of diversity constraints. While in most cases we find our problems to be NP-hard (even if we only want to check if a committee meeting diversity constraints exists; without optimizing the objective function), we also find practically relevant cases with polynomial-time algorithms (e.g., our algorithms would suffice for the research-meeting example restricted to the constraints regarding the seniority level and gender). We provide approximation algorithms for some of our NP-hard problems.
3. We study the complexity of recognizing various types of label structures. For example, given a set of labeled candidates, we ask if their labels have laminar or layered structure. It turns out that recognizing structures with three independent sets of labels is NP-hard, whereas recognizing up to two independent sets is polynomial-time computable.

Our main results are presented in Table 1. A more complete version of this paper including proof details is available as a preprint (Bredereck et al. 2017).

2 The Model

For $i, j \in \mathbb{N}$, we write $[i, j]$ to denote the set $\{i, i+1, \dots, j\}$. We write $[i]$ as an abbreviation for $[1, i]$. For a set X , we write 2^X to denote the family of all of its subsets. We first present our model in full generality and then describe the particular instantiations that we focus on in our analysis.

General Model Let $C = \{c_1, \dots, c_m\}$ be a set of candidates and let L be a set of labels (such as *junior*, *senior*, etc.). Each candidate is associated with a subset of these labels through a labeling function $\lambda: C \rightarrow 2^L$. We say that a candidate c has label ℓ if $\ell \in \lambda(c)$, and we write C_ℓ to denote the set of all candidates that have label ℓ .

A diversity specification is a function that given a committee (i.e., a set of candidates), the set of labels, and the labeling function provides a *yes/no* answer specifying if the

²If we restricted our example to labels regarding gender and seniority level, we would have 2-layered labels (because there are two sets of labels, $\{\textit{male}, \textit{female}\}$ and $\{\textit{junior}, \textit{senior}\}$, and each candidate has one label from each set. On the other hand, hierarchical labels, such as those regarding countries and continents, are 1-laminar (see description of the model for more details).

committee is *diverse*. If a committee is diverse with respect to diversity specification D , then we say that it is D -diverse.

An objective function $f: 2^C \rightarrow \mathbb{R}$ is a function that associates each committee with a score. We assume that $f(\emptyset) = 0$ and that the function is monotone (i.e., for each two committees A and B such that $A \subseteq B$, it holds that $f(A) \leq f(B)$). In other words, an empty committee has no value and extending a committee cannot hurt it.

Our goal is to find a committee of a given size k that meets the diversity specification and that has the highest possible score according to the objective function.

Definition 1 (DIVERSE COMMITTEE WINNER DETERMINATION (DCWD)). Given a set of candidates C , a set of labels L , a labeling function λ , a diversity specification D , a desired committee size k , and an objective function f , find a committee $W \subseteq C$ with $|W| = k$ that achieves the maximum value $f(W)$ among all D -diverse size- k committees.

The set of candidates, the set of labels, and the labeling function are specified explicitly (i.e., by listing all the candidates with all their labels). The encoding of the diversity specification and the objective function depends on a particular case (see discussions below). To consider the problem's NP-hardness, we take its decision variant, where instead of asking for a D -diverse committee with the highest possible value of the objective function we ask if there exists a D -diverse committee with objective value at least T (where the threshold T is a part of the input).

We also consider the DIVERSE COMMITTEE FEASIBILITY (DCF) problem, which takes the same input as the winner determination problem, but where we ask if any D -diverse committee of size k exists, irrespective of its objective value. In other words, the feasibility problem is a special case of the decision variant of the winner determination problem, where we ask about a D -diverse committee with objective value greater or equal to 0. Thus if the feasibility problem is NP-hard, then the analogous winner determination problem is NP-hard as well (and if the winner determination problem is polynomial-time computable, so is the feasibility problem).

The model, as specified above, is far too general to obtain any sort of meaningful computational results. Below we specify its restrictions that we study.

Objective Functions An objective function is submodular if for each two committees S and S' such that $S \subseteq S' \subseteq C$ and each $c \in C \setminus S'$ it holds that $f(S \cup \{c\}) - f(S) \geq f(S' \cup \{c\}) - f(S')$. For two sets of candidates X and S , we write $f(X|S)$ to denote the marginal contribution of the candidates from X with respect to those in S . Formally, we have $f(X|S) = f(S \cup X) - f(S)$. Submodular functions are very common and suffice to express many natural problems. We assume all our objective functions to be submodular.

Example 1. Consider the following voting scenario. We have a set of candidates $C = \{c_1, \dots, c_m\}$ and a set of voters $V = \{v_1, \dots, v_n\}$, where each voter ranks all the candidates from best to worst. We write $\text{pos}_{v_i}(c)$ to denote the position of candidate c in the ranking of voter v_i (the best candidate is ranked on position 1, the next one on position 2, and so on). The Borda score associated with position i (among m

possible ones) is $\beta_m(i) = m - i$. Under the Chamberlin–Courant rule (CC), the score of a committee S is defined by objective function $f^{CC}(S) = \sum_{i=1}^n \beta_m(\min\{\text{pos}_{v_i}(c) \mid c \in S\})$. Intuitively, this function associates each voter with her representative (the member of the committee that the voter ranks highest) and defines the score of the committee as the sum of the Borda scores of the voters’ representatives. It is well-known that this function is submodular (Lu and Boutilier 2011). The CC rule outputs those committees (of a given size k) for which the CC objective function gives the highest value (and, intuitively, where each voter is represented by a committee member that the voter ranks highly).

As a special case of submodular functions, we also consider *separable* functions. A function is separable if for every candidate $c \in C$ there is a weight w_c such that the value of a committee S is given as $f(S) = \sum_{c \in S} w_c$. While separable functions are very restrictive, they are also very natural.

Example 2. Consider the setting from Example 1, but with objective function $f^{kB(W)} = \sum_{i=1}^n (\sum_{c \in W} \beta_m(\text{pos}_{v_i}(c)))$. This function sums Borda scores of all the committee members from all the voters and models the k -Borda voting rule (the committee with the highest score is selected). The function is separable as for each candidate c it suffices to take $w_c = f^{k\text{-Borda}}(\{c\})$. It is often argued that k -Borda is a good rule when our goal is to shortlist a set of individually excellent candidates (Faliszewski et al. 2017).

Together, Examples 1 and 2 show that our model suffices to capture many well-known multiwinner voting scenarios. Many other voting rules, such as Proportional Approval Voting, or many committee scoring rules, can be expressed through submodular objective functions (Skowron, Faliszewski, and Lang 2016; Faliszewski et al. 2016).

Diversity Specifications We focus on diversity specifications that regard each label independently. In other words, the answer to the question if a given committee S is diverse or not depends only on the cardinalities of the sets $C_\ell \cap S$.

Definition 2. For a set of candidates C , a set of labels L , and a labeling function λ , we say that a diversity specification D is independent (consists of independent constraints) if and only if there is a function $b: L \rightarrow 2^{|C|}$ (referred to as the cardinality constraint function) such that a committee S is diverse exactly if for each label ℓ it holds that $|S \cap C_\ell| \in b(\ell)$.

If we have m candidates then specifying independent constraints requires providing at most $m + 1$ numbers for each label. Thus independent constraints can easily be encoded in the inputs for our algorithms.

Independent constraints are quite expressive. For example, they are sufficient to express conditions such as “the committee must contain an even number of junior researchers” or, since our committees are of a given fixed size, conditions of the form “the committee must contain at least 40% females.” Indeed, the conditions of the latter form are so important that we consider them separately.

Definition 3. For a set of candidates C , a set of labels L , and a labeling function λ , we say that a diversity specification D

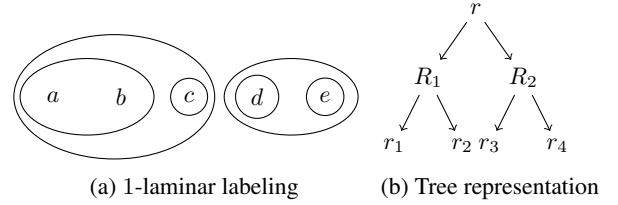


Figure 1: Illustration of a 1-laminar labeling structure.

is interval-based (consists of interval constraints) if and only if there are functions $b_1, b_2: L \rightarrow 2^{|C|}$ (referred to as the lower and upper interval constraint functions) such that a committee S is diverse if and only if for each label ℓ it holds that $b_1(\ell) \leq |S \cap C_\ell| \leq b_2(\ell)$.

Label Structures In principle, our model allows each candidate to have an arbitrary set of labels. In practice, there usually are some dependencies between the labels and these dependencies can have strong impact in the complexity of our problem. We focus on labels that are arranged in independent, possibly hierarchically structured, layers.

Let C be a set of candidates, let L be a set of labels, and let λ be a labeling function. We say that λ has *1-layered* structure (i.e., we have a 1-layered labeling) if for each two distinct labels x, y it holds that $C_x \cap C_y = \emptyset$ (i.e., each candidate has at most one of these labels). For example, if we restricted the example from the introduction to labels regarding the seniority level (junior or senior), then we would have a 1-layered labeling.

More generally, we say that a labeling is *1-laminar* if for each two distinct labels x, y we have that either (a) $C_x \cap C_y = \emptyset$ or (b) $C_x \subseteq C_y$ or (c) $C_y \subseteq C_x$. In other words, 1-laminar labellings allow the labels to be arranged hierarchically.

Example 3. Consider a set $C = \{a, b, c, d, e\}$ of five candidates and labels that encode the countries and continents where the candidates come from. Specifically, there are four countries r_1, r_2, r_3, r_4 , and two continents R_1 and R_2 . The candidates are labeled as follows:

$$\begin{aligned} \lambda(a) &= \{r_1, R_1\}, & \lambda(b) &= \{r_1, R_1\}, & \lambda(c) &= \{r_2, R_1\}, \\ \lambda(d) &= \{r_3, R_2\}, & \lambda(e) &= \{r_4, R_2\}. \end{aligned}$$

Figure 1a illustrates the 1-laminar inclusion-wise relations between the labels (there can be more levels of the hierarchy; for example, for each country there could be labels specifying local administrative division).

Every 1-laminar labeling, together with the set of candidates, can be represented as a rooted tree T in the following way: For a pair of distinct labels x, y we create an arc from x to y if $C_x \subsetneq C_y$ and there is no label z such that $C_x \subsetneq C_z \subsetneq C_y$. We add a *root label* r and we impose that each candidate has this label; we add an arc from r to each label without an incoming arc. The resulting digraph T is clearly a rooted tree. See Figure 1b for an illustration.

For each positive integer t , we say that a labeling is *t-layered* (respectively, *t-laminar*) if the set L of labels can be partitioned into sets L_1, L_2, \dots, L_t such that for each

$i \in [t]$, the labeling restricted to the labels from L_i is 1-layered (respective, 1-laminar).

Example 4. In the example from the introduction, restricting our attention to candidates' gender and seniority levels, we get a 2-layered labeling structure. If we also consider labels regarding countries and continents, then we get a 3-laminar structure (however, only the geographic labels would be using the full power of laminar labellings).

We assume that when we are given a t -layered (t -laminar) labeling structure, we are also given the partition of the set of labels that defines this structure (in Section 5 we analyze the problem of recognizing such structures algorithmically).

Balanced Committee Model As a very natural special case of our model we considered the problem of computing balanced committees. In this case there are only two labels (e.g., *male* and *female*), each candidate has exactly one label, and the constraint specification is that we need to select exactly the same number of candidates with either label (thus, by definition, the committee must be of an even size).

Computing balanced committees is a very natural problem. For example, seeking gender balance is a common requirement in many settings. In this paper, we seek exact balance (that is, we seek exactly the same number of candidates with either label) but allowing any other proportion would lead to similar results.

3 Separable Objective Functions

Separable objective functions form a simple, but very important special case of our setting. Indeed, such functions are very natural in shortlisting examples, where diversity constraints are used to implement, e.g., affirmative actions or employment-equity laws. We organize our discussion with respect to the type of constraint specifications.

Independent Constraints It turns out that independent constraints are quite difficult to work with. If the labels are 1-laminar then polynomial-time algorithms exist (both for deciding if feasible committees exist and for computing optimal ones), but with 2-layered labellings our problems become NP-hard (recall that t -layered labellings are a special case of t -laminar ones). Our polynomial-time algorithms proceed via dynamic programming and hardness proofs use reductions from EXACT 3 SET COVER (X3C).

Theorem 4. *Let D be a diversity specification of independent constraints. Suppose that λ is 1-laminar and f is separable. Then, DCWD can be solved in $O(|L|^2 k^2 + |C| \log |C|)$ time. Moreover, DCF can be solved in $O(|L|^2 k^2)$ time. If the λ function is 2-layered then both problems are NP-hard (even if each candidate has at most two labels, and each label is associated to at most three candidates).*

Given the above hardness results, it is immediate to ask about the parametrized complexity of our problems because in many settings the label structures are very limited (for example, the 2-layered gender/seniority labeling from the introduction contains only 4 labels and already is very relevant for practical applications). Unfortunately, for independent

constraints our problems remain hard when parametrized by the number of labels.

Theorem 5. *Both DCF and DCWD problems are $W[1]$ -hard with respect to the number of labels $|L|$, even if D is a diversity specification of independent constraints.*

However, not all is lost and sometimes brute-force algorithms are sufficiently effective. For example, if we have a t -layered labeling (where t is a small constant) then each candidate has at most t different labels and it suffices to consider each size- t labeling separately. A brute-force algorithm based on this idea suffices, e.g., for the example from the seniority/specialty labels from the introduction (it would have $O(|C|^4)$ running time, because there are 4 combinations of labels $\{\textit{junior}, \textit{senior}\}$ and $\{\textit{AI}, \textit{economics}\}$; the algorithm could also deal with non-independent constraints).

Interval constraints Interval constraints are more restrictive than general independent ones, but usually suffice for practical applications and are more tractable. For example, for the case of 1-laminar labellings we give a linear-time algorithm for recognizing if a feasible committee exists (for independent constraints, our best algorithm for this task is quadratic).

Theorem 6. *Let D be a diversity specification of interval constraints. If λ is 1-laminar, then DCF can be solved in $O(|C| + |L|)$ time.*

For the case of computing the winning committee we no longer obtain a significant speedup from focusing on interval constraints, but we do get a much better structural understanding of the problem. In particular, we can use a greedy algorithm instead of relying on dynamic programming. Briefly put, our algorithm (presented as Algorithm 1) starts with an empty committee and performs k iterations (k is the desired committee size), in each extending the committee with a candidate that increases the score maximally, while ensuring that the committee can still be extended to one that meets the diversity constraints. To show that this greedy algorithm is correct and that it can be implemented efficiently, we use some notions from the matroid theory.

Formally, a *matroid* is an ordered pair (C, \mathcal{I}) , where C is some finite set and \mathcal{I} is a family of its subsets (referred to as the *independent sets* of the matroid). We require that (I1) $\emptyset \in \mathcal{I}$, (I2) if $S \subseteq T \in \mathcal{I}$, then $S \in \mathcal{I}$, and (I3) if $S, T \in \mathcal{I}$ and $|S| > |T|$, then there exists $s \in S \setminus T$ such that $T \cup \{s\} \in \mathcal{I}$. The family of maximal (with respect to inclusion) independent sets of a matroid is called its *basis*. Many of our arguments use results from matroid theory, but often used in very different contexts than originally developed. In particular, the next theorem, in essence, translates the results of Yokoi (2017) to our setting.

Theorem 7. *Let D be a diversity specification of interval constraints. Suppose that λ is 1-laminar, and f is a separable function given by a weight vector $w: C \rightarrow \mathbb{R}$. Then, DCWD can be solved in $O(k^2 |C| |L| + |C| \log |C|)$ time.*

Proof. Let \mathcal{K}_D be the set of D -diverse committees of size k , and assume that \mathcal{K}_D is nonempty. For a family of subsets \mathcal{K}

Algorithm 1: Greedy Algorithm 1

notation : $\mathcal{K}_D \neq \emptyset$ is the set of D -diverse, size- k committees, $\bar{\mathcal{K}}_D$ is its lower extension.

input : $f: 2^C \rightarrow \mathbb{R}$: the objective function,
 k : the size of the committee.

output : $W \in \mathcal{K}_D$

- 1 set $W = \emptyset$;
 - 2 **while** $|W| < k$ **do**
 - 3 choose a candidate $y \in C \setminus W$ such that
 $W \cup \{y\} \in \bar{\mathcal{K}}_D$ with the maximum improvement
 $f(\{y\} | W)$;
 - 4 set $W \leftarrow W \cup \{y\}$;
-

of a finite set C , we define its lower extension to be

$$\bar{\mathcal{K}} = \{T \mid \exists S \in \mathcal{K} : T \subseteq S\}.$$

It follows from the work of Yokoi (2017) that if the constraints are given by intervals and $\mathcal{K}_D \neq \emptyset$, then the lower extension $\bar{\mathcal{K}}_D$ of \mathcal{K}_D forms a family of independent sets of some matroid.³ Thus Algorithm 1 finds an optimal solution $W \in \operatorname{argmax}_{W' \in \bar{\mathcal{K}}_D} f(W')$ (see, e.g., Chapter 13 of Korte and Vygen (2006)).

Let W be the committee produced by Algorithm 1. Since W contains k elements, it must belong to \mathcal{K}_D (because all size- k subsets of $\bar{\mathcal{K}}_D$ are elements of \mathcal{K}_D). For the same reason, since W maximizes the score among the sets from $\bar{\mathcal{K}}_D$, it must be the case that $W \in \operatorname{argmax}_{W' \in \mathcal{K}_D} f(W')$ and, so, W is a winning committee. Further, Yokoi (2017) showed that checking whether a set $W \cup \{y\}$ belongs to $\bar{\mathcal{K}}_D$ can be efficiently done by maintaining a set $B \in \mathcal{K}_D$ with $W \subseteq B$ and, so, the greedy algorithm runs in polynomial time. \square

Unfortunately, the greedy algorithm does not work for more involved labeling structures, but for 2-laminar labellings we can compute winning committees by reducing the problem to the matroid intersection problem (Edmonds 1979). For more involved labeling structures our problems become NP-hard.

Theorem 8. *Let D be a diversity specification of interval constraints. Suppose that λ is 2-laminar and f is separable. Then, DCF can be solved in $O(k^2|C|^3|L|)$ time, and DCWD can be solved in $O(k|C|^3 + k^3|C|^2|L|)$ time.*

The bound on the number of layers turns out to be necessary: the following theorem shows that finding a D -diverse committee is intractable even with 3-layers.

Theorem 9. *DCF is NP-hard even if D is a diversity specification of interval constraints and λ is 3-layered.*

³These independent sets form a relaxed version of our constraints. However, taking the lower extension does not necessarily ignore the lower bounds. For instance, consider a setting where we want to select a committee of size 5 such that there are exactly three female candidates and at most two male candidates; the corresponding lower extension $\bar{\mathcal{K}}_D$ only includes the sets of female candidates of size at most 3, whereas a male-only committee of size 2 satisfies the upper bounds on the respective number of female/male candidates.

Proof. We reduce from 3-DIMENSIONAL MATCHING (3-DM). Given three disjoint sets X, Y, Z of size n and a set $T \subseteq X \times Y \times Z$ of ordered triplets, 3-DM asks whether there is a set of n triplets in T such that each element is contained in exactly one triplet.

Given an instance $((X, Y, Z), T)$ of 3-DM, we create one candidate $t_i = (x_i, y_i, B_i)$ for each $t_i \in T$. The set of labels is given by $L = X \cup Y \cup Z$. Each candidate t_i has exactly three labels $\lambda(t_i) = \{x_i, y_i, B_i\}$. The lower bound $b_1(\ell)$ and the upper bound $b_2(\ell)$ of each label $\ell \in L$ are set to be 1. Lastly, we set $k = n$. It can be easily verified that $W \subseteq T$ is a desired solution for 3-DM if and only if W is a D -diverse committee of size k , namely, $|W| = k$, and (i) $|C_x \cap W| = 1$ for each $x \in X$, (ii) $|C_y \cap W| = 1$ for each $y \in Y$, and (iii) $|C_z \cap W| = 1$ for each $z \in Z$. \square

Nevertheless, if the number of labels is small (i.e., is taken as the parameter from the point of view of parametrized complexity theory) we can compute optimal diverse committees efficiently. The next theorem expresses this formally (note that interval diversity specifications can be phrased as linear programs, but this language allows also some more involved constraints, such as, “at the research meeting the number of senior researchers should be larger than the number of junior ones, but without taking the PhD students into account”).

Theorem 10. *Let f be separable objective function and let D be a diversity specification which can be expressed through a linear program LP with the set of variables $\{x_\ell : \ell \in L\}$ such that $d \in D$ if and only if LP instantiated with variables x_ℓ giving the numbers of committee members with labels ℓ is feasible. Then, DCWD is in FPT with respect to $|L|$.*

4 Submodular Objective Functions

The case of submodular objective functions is computationally far more difficult than that of separable ones. Indeed, even without diversity constraints computing a winning Chamberlin–Courant committee (specified through a submodular objective function) is NP-hard (Lu and Boutilier 2011) and, in general, the best polynomial-time approximation algorithm for submodular functions is the classic greedy algorithm (Nemhauser, Wolsey, and Fisher 1978; Feige 1998), which achieves the $1 - 1/e \approx 0.63$ approximation ratio.⁴ Adding diversity constraints makes our problems even more difficult. Nonetheless, we provide a polynomial-time $1/2$ -approximation algorithm for the case of interval constraints and 1-laminar labellings.

Theorem 11. *Let D be a diversity specification of interval constraints. If λ is 1-laminar and f is a monotone submodular function, then Algorithm 1 gives $1/2$ -approximation algorithm for DCWD.*

Balanced Committees For the balanced committee model it is possible to achieve notably stronger results. Since the balanced case is practically relevant from practical standpoint, we provide its simpler definition, renaming it as BCWD.

⁴This algorithm starts with an empty committee and extends it with candidates one-by-one, always choosing the candidate that increases the objective function maximally.

Algorithm 2: Greedy Algorithm for BCWD

input : $f: 2^C \rightarrow \mathbb{R}$, $A \subseteq C$ and $B \subseteq C$ where
 $A \cap B = \emptyset$, $|A| \geq k'$ and $|B| \geq k'$
output : $W \subseteq C$ where $|W \cap A| = |W \cap B| = k'$
1 **while** $|W| < 2k'$ **do**
2 choose a pair $e = \{a, b\}$ where $a \in A \setminus W$ and
 $b \in B \setminus W$ with maximum improvement $f(e|W)$;
3 set $W \leftarrow W \cup e$;

Definition 12 (BCWD). Given a set of candidates C , two subsets $A, B \subseteq C$ such that $A \cap B = \emptyset$ and $A \cup B = C$, a desired committee size $k = 2k'$, and an objective function f , find a committee $W \subseteq C$ that maximizes $f(W)$ and that satisfies $|W \cap A| = |W \cap B| = k'$.

For the case of BCWD, we provide a polynomial-time $1 - 1/e$ approximation algorithm. Since this is the best possible approximation ratio for general submodular functions without diversity constraints, it is also the best one for the balanced setting (formally, the results without diversity constraints translate because we could assume that all the candidates with one of the labels have no influence on the objective value and use the remaining ones to model an unconstrained submodular optimization problem). Our algorithm (presented as Algorithm 2) is very similar to the classic greedy algorithm, but it considers candidates in pairs.

Theorem 13. *Let f be a monotone submodular function. Algorithm 2 gives $(1 - \frac{1}{e})$ -approximation algorithm for BCWD.*

We note that Theorem 13 is a special case of a much more general result on approximating the Multidimensional Knapsack problem (Kulik, Shachnai, and Tamir 2013), which gives the same approximation ratio even for maximizing monotone submodular functions subject to interval constraints consisting only of upper bounds. Yet, our algorithm is simpler and faster than this general approach.

Theorem 13 applies to all submodular functions. However, for some special cases it is possible to achieve much stronger results. For example, for the Chamberlin–Courant function we find a polynomial-time approximation scheme (PTAS).

Theorem 14. *For each Chamberlin–Courant function there exists a PTAS for BCWD.*

The main idea behind the proof is to use the PTAS of Skowron et al. (2015) to compute a committee of size k' and then to complement it so that it satisfies the diversity constraints. The specific nature of the algorithm of Skowron et al. makes it possible to do this efficiently and effectively.

Theorem 14 also extends to the case of the constant number of labels $\ell_1, \ell'_1, \ell_2, \ell'_2, \dots, \ell_p, \ell'_p$ which satisfy the following two conditions: (i) all the constraints have the following form: for $i \in [p]$ we require the same number of candidates with label ℓ_i as those with label ℓ'_i , (ii) for each combination of labels (r_1, r_2, \dots, r_p) with $r_i \in \{\ell_i, \ell'_i\}$ for each $i \in [p]$, there exist at least k candidates having all labels r_1, \dots, r_p .

5 Recognizing Structure of the Labels

In this section we ask how difficult it is to recognize a given labeling structure if it is not provided with the problem. While in most cases it is natural to assume that the structure would be provided (as it would be a common knowledge of the society for which we would want to compute the committee), it is interesting to be able to derive it automatically.

In the previous sections we have seen that there usually are polynomial-time algorithms for computing winning committees for 1-laminar labellings and, sometimes, there are such algorithms for 2-laminar ones. However, 3-laminar labellings always lead to NP-hardness results. The same holds for the label-structure recognition problem. There are algorithms that decide if given labellings are 1- or 2-laminar, but recognizing 3-layered ones is NP-hard. In the labeling-recognition problem we are given a set of candidates C , a set of labels L , and a labeling function λ . Our goal is to recognize if λ is t -laminar (or t -layered), for a given t .

Proposition 15. *For $t \in \{1, 2\}$ there exists a polynomial-time algorithm for deciding if a given labeling λ is t -laminar. The problem of deciding if a given labeling λ is 3-layered is NP-hard.*

6 Related Work

Our work touches upon many concepts and, thus, is related to many pieces of research. In this section we briefly mention some of the most relevant ones.

Lang and Skowron (2016) considered a model of diversity requirements that closely resembles our interval constraints. There are two main differences between their work and ours: (i) they do not consider objective functions and (ii) their input consists of “ideal points” instead of intervals for each label; since there might not exist a committee satisfying such “exact” constraints, they focus on finding committees minimizing a certain distance to the ideal diversity distributions.

If the labels denote party affiliations of the candidates, the diversity constraints are one-layered and form instances for the apportionment problem, where seats in the parliament should be distributed among the parties (see the book of Balinski and Young (1982) for an overview of the apportionment problems). Bi-apportionment (Balinski and Demange 1989a; 1989b) can be viewed as an extension of the traditional apportionment to the case when the diversity constraints are two-layered. However, in all these settings there is no objective functions, and the goal is only to find a committee satisfying certain label-based constraints. For this reason our paper is even closer the work of Brams (1990), who introduces a specific method based on approval voting that takes diversity constraints into account, which are expressed as quotas for each possible *tuple* of labels; Potthoff (1990) and Straszak et al. (1993) formulated an ILP for this method.

Optimization of a given objective functions due to constraints is a classic problem studied extensively in the literature. For a review of this literature we refer the reader to the book of Korte and Vygen (2006). More specifically, Krause and Golovin (2012) provide a comprehensive survey for the case when the optimized function is submodular. For submodular functions different types of general constraints are con-

sidered, including matroid and knapsack constraints (Chekuri, Vondrk, and Zenklusen 2014). A particularly related case is when the constraints are given for the size of the committee (see e.g., the works of Qian et al. (2017) and the references inside)—interestingly, this case can be represented in our model, when we assume that there is a single label assigned to each candidate, and the constraints are given for the number of occurrences of this label in the elected committee. Candidates having positive synergies may induce supermodular (instead of submodular) objective functions. We note that constrained maximization of a supermodular function is equivalent to constrained minimization of a submodular function, known to be NP-hard (Iwata and Nagano 2009).

Our model is related to the Multidimensional Knapsack problem with submodular objective functions (Fréville 2004; Sviridenko 2004; Lee et al. 2009; Florios, Mavrotas, and Diakoulaki 2010; Puchinger, Raidl, and Pferschy 2010; Kulik, Shachnai, and Tamir 2013), but differs in a few important aspects. The two biggest differences are: (i) Multidimensional Knapsack has constraints of the form “no more than value D on dimension i ” (dimensions correspond to labels in our work), whereas our constraints can have more structure (specific quantities of a given label, or upper and lower bounds), (ii) Multidimensional Knapsack has items that can contribute more than a unit weight to a particular dimension, whereas our candidates only have 0/1 contributions. Thus, our problem is more general regarding the constraint specification, but less general regarding the structure of the weights of items.

The complexity of selecting an optimal committee without constraints has been studied extensively. For a general overview of this literature, we point the reader to a chapter by Faliszewski et al. (2017). Perhaps the most attention was dedicated to the study of the Chamberlin–Courant rule (1983). For instance, it is known that this rule is NP-hard to compute (Procaccia, Rosenschein, and Zohar 2008). The problem of finding a winning Chamberlin–Courant committee under restricted domains of voters’ preferences was further studied by Betzler et al. (2013), Yu et al. (2013), Elkind and Lackner (2015), Skowron et al. (2015), and Peters and Lackner (2017). Parametrized complexity of the problem was studied by Betzler et al. (2013) and its approximability by Lu and Boutilier (2011), Skowron et al. (2015) and Skowron and Faliszewski (2015).

Finally, we note that Celis, Huang, and Vishnoi (2017) very recently and independently introduced a model for diversity constraints (in their paper referred to as *fairness constraints*) that is similar to our model. Their paper contains algorithmic results, which are also applicable in our setting.

7 Conclusion

We studied the problem of selecting a committee of a given size that, on the one hand, would be diverse (according to a given diversity specification) and, on the other hand, would obtain as high an objective value as possible. We present our results in Table 1. We find that in general our problem is computationally hard, but there are many tractable special cases, especially for separable objective functions (which are very useful for shortlisting tasks, where diversity constraints

Rule	Label Structure	Interval Constraints	Independent Constraints
separable	1-laminar	P	P
	2-laminar	P	NP-hard
	3-layer	NP-hard	NP-hard
	few labels	FPT	W[1]-hard
submodular	1-laminar	0.5-approx.	?
	balanced	0.63-approx.	—
CC	balanced	PTAS	—

Table 1: The complexity of computing winning committees for rules of a given type, for the case of candidates with particular label structures, and particular diversity specifications. The complexity results for the problem of testing if a feasible committee exists are the same as those for computing winning committees. “Balanced” label structure refers to the problem of computing balanced committees (thus the case of independent constraints is not defined for this setting).

are particularly relevant) and for up to 2-laminar label structures (which means that dealing with two sets of independent, hierarchically arranged labels, is feasible). Our work leads to many open problems. In particular, we barely scratched the surface regarding approximation of our problems, or their parametrized complexity. Experimental studies would be very desirable as well.

Acknowledgments Robert Bredereck was from September 2016 to September 2017 on postdoctoral leave at the University of Oxford, supported by the DFG fellowship BR 5207/2. Piotr Faliszewski was supported by AGH grant 11.11.230.337 (statutory research). Ayumi Igarashi was supported by the Oxford Kobe Scholarship. Martin Lackner was supported by the European Research Council (ERC) under grant number 639945 (ACCORD) and by the Austrian Science Foundation FWF, grant P25518 and Y698. Piotr Skowron was supported by a Humboldt Research Fellowship for Postdoctoral Researchers.

References

- Balinski, M. L., and Demange, G. 1989a. Algorithm for Proportional Matrices in Reals and Integers. *Mathematical Programming, Series A* 45(1-3):193–210.
- Balinski, M. L., and Demange, G. 1989b. An axiomatic approach to proportionality between matrices. *Mathematics of Operations Research* 14:700–719.
- Balinski, M., and Young, H. P. 1982. *Fair Representation: Meeting the Ideal of One Man, One Vote*. Yale University Press. (2nd Edition [with identical pagination], Brookings Institution Press, 2001).
- Betzler, N.; Slinko, A.; and Uhlmann, J. 2013. On the computation of fully proportional representation. *Journal of Artificial Intelligence Research* 47:475–519.
- Brams, S. J. 1990. Constrained approval voting: A voting system to elect a governing board. *Interfaces* 20(5):67–80.
- Bredereck, R.; Faliszewski, P.; Igarashi, A.; Lackner, M.; and Skowron, P. 2017. Multiwinner elections with diversity

- constraints. Technical Report arXiv:1711.06527 [cs.GT], arXiv.org.
- Celis, L. E.; Huang, L.; and Vishnoi, N. K. 2017. Group fairness in multiwinner voting. Technical Report arXiv:1710.10057 [cs.CY], arXiv.org.
- Chamberlin, B., and Courant, P. 1983. Representative deliberations and representative decisions: Proportional representation and the Borda rule. *American Political Science Review* 77(3):718–733.
- Chekuri, C.; Vondrak, J.; and Zenklus, R. 2014. Submodular function maximization via the multilinear relaxation and contention resolution schemes. *SIAM Journal on Computing* 43(6):1831–1879.
- Edmonds, J. 1979. Matroid intersection. *Annals of Discrete Mathematics* 4:39–49. Discrete Optimization I.
- Elkind, E., and Lackner, M. 2015. Structure in dichotomous preferences. In *Proceedings of the 24th International Joint Conference on Artificial Intelligence (IJCAI-2015)*, 2019–2025.
- Faliszewski, P.; Skowron, P.; Slinko, A.; and Talmon, N. 2016. Committee scoring rules: Axiomatic classification and hierarchy. In *Proceedings of the 25th International Joint Conference on Artificial Intelligence (IJCAI-2016)*, 250–256.
- Faliszewski, P.; Skowron, P.; Slinko, A.; and Talmon, N. 2017. Multiwinner voting: A new challenge for social choice theory. In Endriss, U., ed., *Trends in Computational Social Choice*. AI Access.
- Feige, U. 1998. A threshold of $\ln n$ for approximating set cover. *Journal of the ACM* 45(4):634–652.
- Florios, K.; Mavrotas, G.; and Diakoulaki, D. 2010. Solving multiobjective, multiconstraint knapsack problems using mathematical programming and evolutionary algorithms. *European Journal of Operational Research* 203(1):14–21.
- Fréville, A. 2004. The multidimensional 0–1 knapsack problem: An overview. *European Journal of Operational Research* 155(1):1–21.
- Iwata, S., and Nagano, K. 2009. Submodular function minimization under covering constraints. In *Proceedings of the 50th Annual IEEE Symposium on Foundations of Computer Science (FOCS-2009)*, 671–680.
- Korte, B., and Vygen, J. 2006. *Combinatorial Optimization: Polyhedra and Efficiency*. Algorithms and Combinatorics. Springer.
- Krause, A., and Golovin, D. 2012. Submodular function maximization. Technical report.
- Kulik, A.; Shachnai, H.; and Tamir, T. 2013. Approximations for monotone and nonmonotone submodular maximization with knapsack constraints. *Mathematics of Operations Research* 38(4):729–739.
- Lang, J., and Skowron, P. 2016. Multi-attribute proportional representation. In *Proceedings of the 30th Conference on Artificial Intelligence (AAAI-2016)*, 530–536.
- Lee, J.; Mirrokni, V.; Nagarajan, V.; and Sviridenko, M. 2009. Non-monotone submodular maximization under matroid and knapsack constraints. In *Proceedings of the 41st Symposium on Theory of Computing (STOC-2009)*, 323–332.
- Lu, T., and Boutilier, C. 2011. Budgeted social choice: From consensus to personalized decision making. In *Proceedings of the 22nd International Joint Conference on Artificial Intelligence (IJCAI-2011)*, 280–286.
- Nemhauser, G.; Wolsey, L.; and Fisher, M. 1978. An analysis of approximations for maximizing submodular set functions. *Mathematical Programming* 14(1):265–294.
- Peters, D., and Lackner, M. 2017. Preferences single-peaked on a circle. In *Proceedings of the 31st Conference on Artificial Intelligence (AAAI-2017)*, 649–655.
- Potthoff, R. 1990. Use of linear programming for constrained approval voting. *Interfaces* 20(5):79–80.
- Procaccia, A.; Rosenschein, J.; and Zohar, A. 2008. On the complexity of achieving proportional representation. *Social Choice and Welfare* 30(3):353–362.
- Puchinger, J.; Raidl, G.; and Pferschy, U. 2010. The multidimensional knapsack problem: Structure and algorithms. *INFORMS Journal on Computing* 22(2):250–265.
- Qian, C.; Shi, J.; Yu, Y.; and Tang, K. 2017. On subset selection with general cost constraints. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence (IJCAI-2017)*, 2613–2619.
- Skowron, P., and Faliszewski, P. 2015. Fully proportional representation with approval ballots: Approximating the Max-Cover problem with bounded frequencies in FPT time. In *Proceedings of the 29th Conference on Artificial Intelligence (AAAI-2015)*, 2124–2130.
- Skowron, P.; Yu, L.; Faliszewski, P.; and Elkind, E. 2015. The complexity of fully proportional representation for single-crossing electorates. *Theoretical Computer Science* 569:43–57.
- Skowron, P.; Faliszewski, P.; and Lang, J. 2016. Finding a collective set of items: From proportional multirepresentation to group recommendation. *Artificial Intelligence* 241:191–216.
- Skowron, P.; Faliszewski, P.; and Slinko, A. M. 2015. Achieving fully proportional representation: Approximability results. *Artificial Intelligence* 222:67–103.
- Straszak, A.; Libura, M.; Sikorski, J.; and Wagner, D. 1993. Computer-assisted constrained approval voting. *Group Decision and Negotiation* 2(4):375–385.
- Sviridenko, M. 2004. A note on maximizing a submodular set function subject to a knapsack constraint. *Operations Research Letters* 32(1):41–43.
- Yokoi, Y. 2017. A generalized polymatroid approach to stable matchings with lower quotas. *Mathematics of Operations Research* 42(1):238–255.
- Yu, L.; Chan, H.; and Elkind, E. 2013. Multiwinner elections under preferences that are single-peaked on a tree. In *Proceedings of the 23rd International Joint Conference on Artificial Intelligence (IJCAI-2013)*, 425–431.