INFORMATION-THEORETIC ACTIVE CORRELATION CLUSTERING

Anonymous authors

004

010 011

012

013

014

015

016 017 018

019

Paper under double-blind review

Abstract

We study correlation clustering where the pairwise similarities are not known in advance. For this purpose, we employ active learning to query pairwise similarities in a cost-efficient way. We propose a number of effective information-theoretic acquisition functions based on entropy and information gain. We extensively investigate the performance of our methods in different settings and demonstrate their superior performance compared to the alternatives.

1 INTRODUCTION

020 Clustering is an important unsupervised learning problem for which several methods have been 021 proposed in different contexts. Correlation clustering (CC) (Bansal et al., 2004; Demaine et al., 022 2006) is a well-known clustering problem, especially beneficial when both similarity and dissimilarity assessments exist for a given set of N objects. Consequently, CC studies the clustering of objects where pairwise similarities can manifest as positive or negative numbers. It has found a wide range 025 of applications including image segmentation (Kim et al., 2011), bioinformatics (Bonchi et al., 2013), 026 spam filtering (Ramachandran et al., 2007; Bonchi et al., 2014), social network analysis (Bonchi et al., 2012; Tang et al., 2016), duplicate detection (Hassanzadeh et al., 2009), co-reference identification 027 (McCallum & Wellner, 2004), entity resolution (Getoor & Machanavajjhala, 2012), color naming across languages (Thiel et al., 2019) and clustering aggregation (Gionis et al., 2007; Chehreghani 029 & Chehreghani, 2020). CC was initially explored using binary pairwise similarities in $\{-1, +1\}$ (Bansal et al., 2004), and was later extended to support arbitrary positive and negative pairwise 031 similarities in \mathbb{R} (Charikar et al., 2005; Demaine et al., 2006). Finding the optimal solution for CC is known to be NP-hard and APX-hard (Bansal et al., 2004; Demaine et al., 2006), presenting significant 033 challenges. As a result, various approximate algorithms have been developed to address this problem 034 (Bansal et al., 2004; Charikar et al., 2005; Demaine et al., 2006; Ailon et al., 2008; Elsner & Schudy, 2009). Among these, methods based on local search are noted for their superior performance in terms of clustering quality and computational efficiency (Thiel et al., 2019; Chehreghani, 2023). 037

Existing methods generally assume that all $\binom{N}{2}$ pairwise similarities are available beforehand. 038 However, as discussed in (Bressan et al., 2019; García-Soriano et al., 2020), generating pairwise 039 similarities can be computationally intensive and may need to be obtained through resource-intensive 040 queries, e.g., from a human expert. For instance, determining interactions between biological entities 041 often requires the expertise of highly trained professionals, consuming both time and valuable 042 resources (García-Soriano et al., 2020). In tasks like entity resolution, obtaining pairwise similarity 043 queries through crowd-sourcing could also involve monetary costs. Therefore, a central question 044 emerges: How can we design a machine learning paradigm that effectively delivers satisfactory CC results with a limited number of queries for pairwise similarities between objects?

In machine learning, *active learning* is generally employed to address such a question. Its objective is to acquire the most informative data within a constrained budget. Active learning has proven effective in various tasks, including recommender systems (Rubens et al., 2015), sound event detection (Shuyang et al., 2020), analysis of driving time series (Jarl et al., 2022), drug discovery (Viet Johansson et al., 2022), and analysis of logged data (Yan et al., 2018). In the context of active learning, the selection of which data to query is guided by an *acquisition function*. Active learning is most commonly studied for classification and regression problems (Settles, 2009). However, it has also been studied for clustering and is sometimes referred to as *supervised clustering* (Awasthi & Zadeh, 2010). The objective is to discover the ground-truth clustering with a minimal number of

queries to an *oracle* (e.g., a human expert). In this scenario, queries are typically executed in one of two ways: (i) By asking whether two clusters should merge or if one cluster should be divided into multiple clusters (Balcan & Blum, 2008; Awasthi & Zadeh, 2010; Awasthi et al., 2017); (ii) By querying the pairwise relations between objects (Basu et al., 2004; Mazumdar & Saha, 2017b;a; Saha & Subramanian, 2019; Bressan et al., 2019; García-Soriano et al., 2020; van Craenendonck et al., 2018b; Silwal et al., 2023; Gullo et al., 2023; Aronsson & Chehreghani, 2024; Kuroki et al., 2024).

060 Among the aforementioned works on active learning for clustering, only (Mazumdar & Saha, 2017b; 061 Bressan et al., 2019; García-Soriano et al., 2020; Aronsson & Chehreghani, 2024; Kuroki et al., 2024) 062 consider the setting that we are interested in: (i) The clustering algorithm is based on CC; (ii) The 063 pairwise similarities are not assumed to be known in advance; (iii) We assume access to a single noisy oracle, to which a *fixed* budget $B \ll \binom{N}{2}$ of queries for pairwise similarities can be performed; 064 065 (iv) Access to feature vectors is not assumed by the algorithm, meaning that information about 066 the ground-truth clustering is solely obtained through querying the oracle for pairwise similarities. Throughout the paper, this setting will be referred to as *active correlation clustering*. 067

068 The work in (Mazumdar & Saha, 2017b) develops a number of *pivot-based* CC algorithms that satisfy 069 guarantees on the query complexity, assuming a noisy oracle. However, the algorithms are purely theoretical and are not implemented and investigated in practice, and require setting a number of 071 non-trivial parameters (e.g., they assume the noise level is known in advance which is unrealistic). 072 The work in (Bressan et al., 2019; García-Soriano et al., 2020) proposes adaptive and query-efficient versions of the simple pivot-based CC algorithm KwikCluster (Ailon et al., 2008). However, as 073 demonstrated in (Aronsson & Chehreghani, 2024), such pivot-based methods perform very poorly 074 for active CC with noise. The work in (Gullo et al., 2023; Kuroki et al., 2024) address guery-efficient 075 CC by formulating it as a *multi-armed bandit* problem. However, this leads to a number of limiting 076 assumptions in practice. We defer a detailed comparison to Appendix E. 077

The work in (Aronsson & Chehreghani, 2024) proposes a generic active CC framework that overcomes the limitations of previous work and offers several advantages: (i) The pairwise similarities 079 can be any positive or negative real number, even allowing for inconsistencies (i.e., violation of transitivity). This allows the oracle to express uncertainty in their feedback; (ii) The process of 081 querying pairwise similarities is decoupled from the clustering algorithm, enhancing flexibility in constructing acquisition functions that can be employed in conjunction with any CC algorithm. 083 (Aronsson & Chehreghani, 2024) employs an efficient CC algorithm based on local search, whose 084 effectiveness (and superiority over pivot-based methods) has also been demonstrated in the standard 085 CC setting (Thiel et al., 2019; Chehreghani, 2023), and dynamically computes the number of clusters; (iii) The framework is robust w.r.t. a noisy oracle and supports multiple queries for the same pairwise 087 similarity if needed (to deal with noise).

Furthermore, (Aronsson & Chehreghani, 2024) proposes two novel acquisition functions, namely *maxmin* and *maxexp*, to be used within their framework. They demonstrate that the algorithm QECC from (García-Soriano et al., 2020) performs poorly in the presence of even a very small amount of noise and is significantly outperformed by their methods. In this paper, we adopt the generic active CC framework in (Aronsson & Chehreghani, 2024) with a focus on the development of more effective acquisition functions. The contributions of this paper are the following:

094

096

098

099

100

101

102

• We investigate the use of information-theoretic acquisition functions based on *entropy* and *information gain* for active CC. We propose four different acquisition functions inspired by this (see Section 3). Although information-theoretic acquisition functions have been extensively studied in the context of active learning (Roy & McCallum, 2001; Kirsch & Gal, 2022), prior research has focused mainly on (active) *supervised learning* scenarios, where the goal is to query data labels from an oracle rather than pairwise relations. *To our knowledge, our work is the first attempt to propose information-theoretic acquisition functions to active learning with pairwise relations, as well as to non-parametric models like CC.* Computing the necessary quantities in this setting is significantly more complex. The methods proposed in this paper can be applied beyond active CC, including to the active learning of other pairwise (non-parametric) clustering models.

103 104

We conduct extensive experimental studies on various datasets that demonstrate the superior performance of our acquisition functions compared to *maxmin* and *maxexp* (and other baselines), and investigate a number of interesting insights about the active CC framework from (Aronsson & Chehreghani, 2024) (see Section 4 and Appendix C).

¹⁰⁸ 2 ACTIVE CORRELATION CLUSTERING

In this section, we begin by introducing the problem of active CC. After this, we describe the active clustering procedure used to solve this problem.

112 113

114

2.1 PROBLEM FORMULATION

115 We are given a set of N objects (data points) indexed by $\mathcal{V} = \{1, \dots, N\}$. The set of pairs of objects 116 in \mathcal{V} is denoted by $\mathcal{E} = \{(u, v) \mid u, v \in \mathcal{V}\}$. We assume the existence of a ground-truth similarity matrix $S^* \in \mathbb{R}^{N \times N}$, which represents the true pairwise similarities between every pair $(u, v) \in \mathcal{E}$. 117 However, S^* is not known beforehand. Instead, one can only query the oracle for a noisy version of 118 this matrix for a desired pair of objects, while incurring some cost. We use $S \in \mathbb{R}^{N \times N}$ to represent 119 an estimate of the pairwise similarities. If $S_{uv} = S_{uv}^*$ for all $(u, v) \in \mathcal{E}$ we have a perfect estimate 120 of the true pairwise similarities, which we assume is unrealistic in practice. Hence, the objective 121 is to discover the ground-truth clustering solution with a minimal number of (active) queries for 122 the pairwise similarities to the oracle, since each query incurs some cost. A similarity matrix S is 123 symmetric, and we assume zeros on the diagonal, i.e., $S_{uv} = S_{vu}$ and $S_{uu} = 0$. This means there 124 are $\binom{N}{2} = (N \times (N-1))/2$ unique pairwise similarities to estimate. Without loss of generality, 125 we assume all similarities are in the range [-1, +1]. In this case, +1 and -1 respectively indicate 126 definite similarity and dissimilarity. Thus, a similarity close to 0 indicates a lack of knowledge about 127 the relation between the two objects. This allows the oracle to express uncertainty in their feedback.

A clustering is a partition of \mathcal{V} . In this paper, we encode a clustering with K clusters as a clustering solution $c \in \mathbb{K}^N$ where $\mathbb{K} = \{1, \dots, K\}$ and $c_u \in \mathbb{K}$ denotes the cluster label of object $u \in \mathcal{V}$. We denote by C the set of clustering solutions for all possible partitions (clusterings) of \mathcal{V} . Given a clustering solution $c \in C$, the CC cost function $R^{CC} : C \to \mathbb{R}^+$ aims to penalize cluster disagreements, as shown in Eq. 1.

$$R^{\text{CC}}(\boldsymbol{c} \mid \boldsymbol{S}) \triangleq \sum_{(u,v) \in \mathcal{E}} \begin{cases} |S_{uv}| & \text{if } (c_u = c_v \text{ and } S_{uv} < 0) \text{ or } (c_u \neq c_v \text{ and } S_{uv} \ge 0) \\ 0 & \text{otherwise.} \end{cases}$$
(1)

Proposition 2.1. Eq. 1 can be simplified to $R^{CC}(c \mid S) = -\sum_{\substack{(u,v) \in \mathcal{E} \\ c_u = c_v}} S_{uv}$ + constant, where the constant is independent of different clustering solutions (Chehreghani, 2013).

All the proofs are in Appendix A. Based on Proposition 2.1, we define the *max correlation* cost function as $R^{MC}(c \mid S) \triangleq -\sum S_{uv}, \qquad (2)$

$${}^{\mathrm{MC}}(\boldsymbol{c} \mid \boldsymbol{S}) \triangleq -\sum_{\substack{(u,v) \in \mathcal{E} \\ c_u = c_v}} S_{uv}, \tag{2}$$

and we have $\arg\min_{c \in C} R^{CC}(c \mid S) = \arg\min_{c \in C} R^{MC}(c \mid S)$. Because of this, we will use R^{MC} throughout most of the paper, as it leads to a number of simplifications in the presented methods. The conditioning on S for R^{CC} and R^{MC} will often be dropped, unless it is not clear from context. Finally, the ground-truth clustering solution corresponds to $c^* = \arg\min_{c \in C} R^{MC}(c \mid S^*)$.

149 150

143 144

134 135 136

2.2 ACTIVE CORRELATION CLUSTERING PROCEDURE

151 We adopt the recent generic active CC procedure outlined in (Aronsson & Chehreghani, 2024) to 152 solve the problem described in the previous section. The procedure is shown in Alg. 1. It takes 153 an initial similarity matrix S^0 as input, which can contain partial or no information about S^* , 154 depending on the initialization method. The procedure then follows a number of iterations, where 155 each iteration i consists of three steps: (i) Update the current clustering solution $c^i \in C$ by running a CC algorithm given the current similarity matrix S^i . The current similarity matrix S^i will be 156 referred to as S throughout the paper; (ii) Select a batch $\mathcal{B} \subseteq \mathcal{E}$ of pairs of size $B = |\mathcal{B}|$ based 157 on an acquisition function $a: \mathcal{E} \to \mathbb{R}$. The quantity a(u, v) indicates how informative the pair 158 $(u, v) \in \mathcal{E}$ is, where a higher value implies greater informativeness. The optimal batch is selected by 159 $\mathcal{B} = \arg \max_{\mathcal{B} \subseteq \mathcal{E}, |\mathcal{B}| = B} \sum_{(u,v) \in \mathcal{B}} a(u,v)$. This corresponds to selecting the top-B pairs based on 160 their acquisition value; (iii) Query the oracle for the pairwise similarities of the pairs $(u, v) \in \mathcal{B}$ and 161 update each S_{uv}^{i+1} based on the response.

| Algorithm 1 Active CC | |
|--|-------------------|
| 1: Input: Initial similarity matrix S^0 , acquisition function a, bate | h size <i>B</i> . |
| 2: $i \leftarrow 0$ | |
| 3: while query budget not reached do | |
| 4: $c^i \leftarrow CC(S^i)$ | ⊳ Alg. 5 |
| 5: $\mathcal{B} \leftarrow \arg \max_{\mathcal{B} \subseteq \mathcal{E}, \mathcal{B} = B} \sum_{(u,v) \in \mathcal{B}} a(u,v)$ | C C |
| 6: Query (noisy) oracle and update S_{uv}^{i+1} for all pairs $(u, v) \in$ | B |
| 7: $i \leftarrow i + 1$ | |
| 8: end while | |
| 9: return c^i | |

3 INFORMATION-THEORETIC ACQUISITION FUNCTIONS

In this section, we introduce four information-theoretic acquisition functions for active CC. All quantities defined below are conditioned on the current similarity matrix S, but it is left out for brevity. All acquisition functions proposed in this section depend on the Gibbs distribution defined as

$$P^{\text{Gibbs}}(\mathbf{y} = \mathbf{c}) \triangleq \frac{\exp(-\beta R^{\text{MC}}(\mathbf{c}))}{\sum_{\mathbf{c}' \in \mathcal{C}} \exp(-\beta R^{\text{MC}}(\mathbf{c}'))},$$
(3)

where $\beta \in \mathbb{R}^+$ is the concentration parameter, $\mathbf{y} = \{\mathbf{y}_1, \dots, \mathbf{y}_N\}$ is a random vector with sample 181 space C (all possible clustering solutions of V) and y_u is a random variable for the cluster label of 182 u with sample space K. Computing P^{Gibbs} is intractable due to the sum over all possible clustering 183 solutions C in the denominator. Therefore, in the next section, we describe a *mean-field approximation* 184 of P^{Gibbs} which makes it possible to efficiently calculate the proposed acquisition functions. To 185 the best of our knowledge, the use of mean-field approximation to approximate complex quantities when applying information-theoretic acquisition functions for active learning is a novel aspect of our 187 approach. This approach can be applied beyond active CC, extending to the active learning of other 188 pairwise, non-parametric clustering models. 189

190 3.1 MEAN-FIELD APPROXIMATION FOR CC

173

174 175

176

177

178

179

199

210 211 212

213

We here describe the mean-field approximation of P^{Gibbs} . The family of factorial distributions over the space of clustering solutions is defined as $Q = \{Q \in \mathcal{P} \mid Q(\mathbf{y} = \mathbf{c}) = \prod_{u \in \mathcal{V}} Q(\mathbf{y}_u = c_u)\}$, where \mathcal{P} is the space of all probability distributions with sample space \mathcal{C} . The goal of mean-field approximation is to find a factorial distribution $Q \in Q$ that best approximates the intractable distribution P^{Gibbs} . In general, one can compute the optimal Q by minimizing the KL-divergence (Hofmann & Buhmann, 1997; Chehreghani et al., 2012), i.e.,

$$Q^* = \operatorname*{arg\,min}_{Q \in \mathcal{Q}} D_{\mathrm{KL}}(Q \| P^{\mathrm{Gibbs}}) = \operatorname*{arg\,min}_{Q \in \mathcal{Q}} \sum_{\boldsymbol{c} \in \mathcal{C}} Q(\boldsymbol{c}) \log \frac{Q(\boldsymbol{c})}{P^{\mathrm{Gibbs}}(\boldsymbol{c})}.$$
(4)

We encode a mean-field approximation using a matrix of assignment probabilities $Q \in [0, 1]^{N \times K}$, where $Q_{uk} = Q(y_u = k)$. In addition, let $M \in \mathbb{R}^{N \times K}$, where M_{uk} should be interpreted as the cost of assigning object u to cluster k. Given this, Theorem 3.1 implies that an EM-type procedure, which sequentially alternates between estimating Q_{uk} (based on Eq. 5) and computing the respective M_{uk} (based on Eq. 6), yields a local minimum for the optimization problem in Eq. 4. In Theorem 3.1, we adapt and specialize the general result from (Hofmann et al., 1998) to our specific cost function in Eq. 2, enabling efficient mean-field approximations tailored to our model, which are essential for all proposed acquisition functions.

Theorem 3.1. Let $\ell : \mathbb{N} \to \mathcal{V}$ denote an object visitation schedule, which satisfies $\lim_{T\to\infty} |\{t \leq T : \ell(t) = u\}| = \infty, \forall u \in \mathcal{V}$. For arbitrary initial conditions, the asynchronous update rules defined by

$$Q_{uk}^{(t+1)} = \exp(-\beta M_{uk}^{(t)}) / \sum_{k' \in \mathbb{K}} \exp(-\beta M_{uk'}^{(t)}),$$
(5)

$$M_{uk}^{(t+1)} = -\sum_{v} S_{uv} Q_{vk}^{(t+1)},$$
(6)

214
$$v \in \mathcal{V}$$

where $u = \ell(t)$, converge to a local minimum of Eq. 4.

216 For computational efficiency, we employ a syn-217 chronous update rule in practice (see Alg. 2). 218 Despite not having the same theoretical guaran-219 tees, synchronous updates have been observed 220 to perform well empirically in other contexts (Hofmann et al., 1998; Chehreghani et al., 2012). 221 Alg. 2 assumes a fixed number of clusters K. 222 We use the number of clusters K dynamically determined by the CC algorithm used at each 224 iteration i of Alg. 1 to find c^i (see Appendix 225 D for details of this algorithm). M could be 226

| Alg | gorithm 2 Mean-Field Approximation | n |
|-----|---|----------------|
| 1: | Input: Similarity matrix S, clust | er assign- |
| | ment costs M , concentration param | eter β . |
| 2: | while Q has not converged do | |
| 3: | $\boldsymbol{Q} \leftarrow \operatorname{softmax}(-\beta \boldsymbol{M})$ | ⊳ E-step |
| 4: | $oldsymbol{M} \leftarrow -oldsymbol{S} \cdot oldsymbol{Q}$ | ⊳ M-step |
| 5: | end while | - |
| 6: | return $oldsymbol{Q},oldsymbol{M}$ | |

initialized randomly. However, since we have the current clustering solution c^i , we initialize it based 227 on c^i , i.e., $M_{uk} = -\sum_{v:c_u^i = k} S_{uv}$, in order to speed up the convergence and potentially improve the 228 quality of the solution found. This initialization of M is based on the total similarity between object 229 u and cluster k in relation to the similarity between u and all other clusters. A smaller similarity 230 should correspond to a higher cost (hence the negation). Each iteration of the algorithm consists of 231 two main steps. First, Q is estimated as a function of M. Second, M is calculated based on Q. In this paper, we treat the concentration parameter $\beta \in \mathbb{R}^+$ as a hyperparameter. Finally, we employ the 232 special form of the max correlation cost function R^{MC} in Eq. 2, and calculate both the E-step and 233 M-step in vectorized form. In particular, the M-step becomes a dot product between S and Q, which 234 is extremely efficient in practice (especially if S is assumed sparse, which it is in our experiments). 235

3.2 ENTROPY

236

237

246 247

252 253

255

266 267 268

In this section, we propose our first acquisition function based on entropy. Let $\mathbf{E} \in \{-1, +1\}^{N \times N}$ be a random matrix where each element $\mathbf{E}_{uv} \in \{-1, +1\}$ is a binary random variable, where +1indicates u and v should be in the same cluster, and -1 implies u and v should be in different clusters. A reasonable way to define the probability of \mathbf{E}_{uv} to be +1 is the fraction of clustering solutions in Cthat assign u and v to the same cluster, weighted by the probability of each clustering solution. Due to the intractability of P^{Gibbs} , we approximate it using a mean-field approximation Q (encoded by matrix Q, as described in the previous section). Formally, we have

$$P(E_{uv} = 1) = \mathbb{E}_{P^{\text{Gibbs}}(\mathbf{y})}[\mathbf{1}_{\{\mathbf{y}_u = \mathbf{y}_v\}}] \approx \mathbb{E}_{Q(\mathbf{y})}[\mathbf{1}_{\{\mathbf{y}_u = \mathbf{y}_v\}}] = \sum_{k \in \mathbb{K}} Q_{uk} Q_{vk}.$$
(7)

The last equality of Eq. 7 uses the fact that the mean-field approximation assumes independence between objects. One can similarly derive $P(E_{uv} = -1) = \sum_{k,k' \in \mathbb{K}} Q_{uk}Q_{vk'}\mathbf{1}_{\{k \neq k'\}} = 1 - P(E_{uv} = 1)$. For more information on the calculations, see Appendix B.1. Thereby, from Eq. 7, we define an acquisition function based on the entropy of E_{uv} as

$$a^{\text{Entropy}}(u,v) \triangleq H(\mathsf{E}_{uv}) = \mathbb{E}_{P(\mathsf{E}_{uv})}[-\log P(\mathsf{E}_{uv})].$$
(8)

254 3.3 INFORMATION GAIN

The acquisition function a^{Entropy} calculates the uncertainty of pairs based on the mean-field approxi-256 mation (model) Q given the current similarity matrix S. In this section, we investigate acquisition 257 functions inspired by the notion of *information gain* corresponding to maximal uncertainty reduction. 258 In this case, the similarity matrix S is first augmented with *pseudo-similarities* (predicted using the 259 current model Q as $S_{uv} \sim P(E_{uv})$), after which a new mean-field approximation is obtained. In 260 other words, we simulate the effect of querying one or more pairs in expectation w.r.t. the current 261 model Q, potentially resulting in more accurate uncertainty estimations. Due to the efficiency of 262 Alg. 2 (mean-field), one can afford to run it several times per iteration of the active CC procedure, to 263 estimate the information gain accurately. In this paper, we consider two types of information gain. 264 First, the information gain (or equivalently the mutual information) between a pair E_{uv} and the cluster 265 labels of objects y. Due to symmetry of the mutual information we have

$$I(\mathbf{y}; \mathbf{E}_{uv}) = H(\mathbf{y}) - H(\mathbf{y} \mid \mathbf{E}_{uv})$$
(9)

$$= H(\mathbf{E}_{uv}) - H(\mathbf{E}_{uv} \mid \mathbf{y}). \tag{10}$$

The interpretation of $I(\mathbf{y}; E_{uv})$ is the amount of information one expects to gain about the cluster labels of objects by observing E_{uv} , where the expectation is w.r.t. $P(E_{uv})$. In other words, it measures the expected reduction in uncertainty (in entropic way) over the possible clustering solutions w.r.t. the value of E_{uv} . Second, the information gain between a pair E_{uv} and all pairs **E**:

$$I(\mathbf{E}; \mathbf{E}_{uv}) = H(\mathbf{E}) - H(\mathbf{E} \mid \mathbf{E}_{uv})$$
(11)

$$= H(\mathbf{E}_{uv}) - H(\mathbf{E}_{uv} \mid \mathbf{E}).$$
(12)

Intuitively, $I(\mathbf{E}; \mathbf{E}_{uv})$ measures the amount of information the pair \mathbf{E}_{uv} provides about all pairs in **E**. All the expressions above are closely related, but the formulations used will impact how they can be approximated in practice, leading to differences in performance and efficiency. This is discussed in detail in the following subsections.

3.3.1 CONDITIONAL MEAN-FIELD APPROXIMATION

We approximate all the conditional entropies defined above following the same general principle: We update the similarity matrix S based on what is being conditioned on, run Alg. 2 given this similarity matrix, and calculate the corresponding entropy given the updated mean-field approximation. Motivated by this, the following notation will be used throughout this section. Let e denote a vector in $\{-1, +1\}^{|\mathcal{D}|}$, where $\mathcal{D} \subseteq \mathcal{E}$ is a subset of the pairs. Given this, we denote by $Q^{(S_{\mathcal{D}}=e)}$ to be the mean-field approximation found by Alg. 2 after modifying S according to e for all pairs $(u, v) \in \mathcal{D}$ (with remaining pairs unchanged).

288 289

290

296 297 298

299

300

301

302 303

304 305

273

274

279

3.3.2 EXPECTED INFORMATION GAIN

In this section, we consider the expression in Eq. 9, which corresponds to the *expected information* gain over cluster labels of objects (EIG-O). We have $H(\mathbf{y} | \mathbf{E}_{uv}) = \mathbb{E}_{e \sim P(\mathbf{E}_{uv})}[H(\mathbf{y} | \mathbf{E}_{uv} = e)]$. In this paper, we approximate $H(\mathbf{y} | \mathbf{E}_{uv} = e)$ using conditional mean-field approximation $Q^{(S_{uv}=e)}$ as shown below. Given some mean-field approximation Q', let $P(\mathbf{E}_{uv} | \mathbf{Q}')$ be the probability of \mathbf{E}_{uv} computed as shown in Eq. 7 using Q' and

$$H(\mathbf{y}_w \mid \boldsymbol{Q}') \triangleq -\sum_{k \in \mathbb{K}} Q'_{wk} \log Q'_{wk}$$

Each $y_u \in y$ is independent of other cluster labels given a mean-field approximation, reducing all joint entropies to the summed entropy across all individual variables. In other words, we have $H(\mathbf{y}) \approx \sum_{w \in \mathcal{V}} H(\mathbf{y}_w \mid \mathbf{Q})$ and $H(\mathbf{y} \mid \mathbf{E}_{uv} = e) \approx \sum_{w \in \mathcal{V}} H(\mathbf{y}_w \mid \mathbf{Q}^{(S_{uv}=e)})$. Given this, we define the following acquisition function.

$$a^{\text{EIG-O}}(u,v) \triangleq \sum_{w \in \mathcal{V}} H(\mathbf{y}_w \mid \boldsymbol{Q}) - \sum_{e \in \{-1,+1\}} P(\mathbf{E}_{uv} = e \mid \boldsymbol{Q}) H(\mathbf{y}_w \mid \boldsymbol{Q}^{(S_{uv} = e)}).$$
(13)

In Appendix B.2, we include a detailed derivation of $a^{\text{EIG-O}}$. In addition, we derive an alternative 306 acquisition function based on Eq. 11, which instead computes the expected reduction in entropy 307 over pairs E (called $a^{\text{EIG-P}}$). However, this method expectedly performs similar to $a^{\text{EIG-O}}$ in practice¹, 308 while being less efficient. Calculating $a^{\text{EIG-O}}$ for all pairs requires executing Alg. 2 $\binom{N}{2}$ times, which 309 can be inefficient for large N. In Alg. 3, we illustrate how we compute $a^{\text{EIG-O}}$ in practice, improving 310 its efficiency in the following ways. (i) We evaluate Eq. 13 only for a subset of the pairs $\mathcal{E}^{ElG} \subseteq \mathcal{E}$. 311 We select this subset as the top- $|\mathcal{E}^{\text{EIG}}|$ pairs according to a^{Entropy} , where $|\mathcal{E}^{\text{EIG}}| = O(N)$ in practice. 312 (ii) We do not expect $Q^{(S_{uv}=e)}$ to be much different from Q. Therefore, by initializing M (in lines 313 7-8) with the assignment costs from line 3, the convergence speed of Alg. 2 significantly improves. 314

315

316

3.3.3 JOINT EXPECTED INFORMATION GAIN

In this section, we consider the information gains formulated in Eq. 10 and Eq. 12. We approximate the conditional entropy $H(E_{uv} | \mathbf{E}) = \mathbb{E}_{e \sim P(\mathbf{E})}[H(E_{uv} | \mathbf{E} = e)]$ in Eq. 12 using the conditional mean-field approximation $Q^{(S_{\mathcal{E}}=e)}$. The conditional entropy $H(E_{uv} | \mathbf{y}) = \mathbb{E}_{c \sim Q(\mathbf{y})}[H(E_{uv} | \mathbf{y} = c)]$ in Eq. 10 is less straightforward here. However, a natural way would be to compute the conditional mean-field approximation given S updated based on c as follows. We set $S_{uv} = +1$ if $c_u = c_v$, and -1 otherwise. In both cases, we then approximate the entropy of E_{uv} given a mean-field

¹This is expected because the distribution $P(\mathbf{E})$ is determined by $Q(\mathbf{y})$ from Eq. 7.

328

330

331

332

333

334

335

336 337 338

339 340 341

342

343

353

354

355

356

357

358 359 360

Algorithm 3 EIG

1: Input: Similarity matrix S, current clustering c^i , concentration parameter β . 2: $M_{uk} \leftarrow -\sum_{v:c_v^i = k} S_{uv}, \forall u \in \mathcal{V}, \forall k \in \mathbb{K}$ 3: $Q, M \leftarrow \text{MeanField}(S, M, \beta)$ 4: $\mathcal{E}^{\text{EIG}} \leftarrow \text{select top-}|\mathcal{E}^{\text{EIG}}|$ pairs using a^{Entropy} given Q (Eq. 8). 5: $a^{\text{EIG}}(u, v) \leftarrow 0, \forall (u, v) \in \mathcal{E}$ 6: for each pair $(u, v) \in \mathcal{E}^{\text{EIG}}$ do 7: $Q^{(S_{uv}=+1)} \leftarrow \text{MeanField}(S, M, \beta \mid S_{uv} = +1)$ 8: $Q^{(S_{uv}=-1)} \leftarrow \text{MeanField}(S, M, \beta \mid S_{uv} = -1)$ 9: $a^{\text{EIG}}(u, v) \leftarrow \text{Evaluate Eq. 13 (or Eq. 34) given } Q, Q^{(S_{uv}=+1)}$ and $Q^{(S_{uv}=-1)}$ 10: end for 11: return a^{EIG}

approximation conditioned on all (or a subset) of the similarities being updated. Given this, we now derive a general estimator based on the information gain

$$I(\mathbf{E}_{uv}; \mathbf{E}_{\mathcal{D}}) = H(\mathbf{E}_{uv}) - H(\mathbf{E}_{uv} \mid \mathbf{E}_{\mathcal{D}}), \tag{14}$$

where $\mathbf{E}_{\mathcal{D}} = \{\mathbf{E}_{uv} \mid (u, v) \in \mathcal{D}\}$ for some $\mathcal{D} \subseteq \mathcal{E}$. From the discussion above, the expressions in Eq. 10 and Eq. 12 can be seen as special cases of Eq. 14. The entropy $H(\mathbf{E}_{uv})$ is approximated based on Eq. 8. In addition, we have

$$H(\mathbf{E}_{uv} \mid \mathbf{E}_{\mathcal{D}}) = \mathbb{E}_{\boldsymbol{e} \sim P(\mathbf{E}_{\mathcal{D}})} [H(\mathbf{E}_{uv} \mid \mathbf{E}_{\mathcal{D}} = \boldsymbol{e})]$$

=
$$\sum_{\boldsymbol{e} \in \{-1, +1\}^{|\mathcal{D}|}} P(\mathbf{E}_{\mathcal{D}} = \boldsymbol{e}) H(\mathbf{E}_{uv} \mid \mathbf{E}_{\mathcal{D}} = \boldsymbol{e}).$$
(15)

Thereby, we approximate the conditional entropy $H(\mathbf{E}_{uv} \mid \mathbf{E}_{\mathcal{D}} = \boldsymbol{e})$ by

$$H(\mathbf{E}_{uv} \mid \boldsymbol{Q}^{(\boldsymbol{S}_{\mathcal{D}}=\boldsymbol{e})}) \triangleq -\sum_{e \in \{-1,+1\}} P(\mathbf{E}_{uv}=e \mid \boldsymbol{Q}^{(\boldsymbol{S}_{\mathcal{D}}=\boldsymbol{e})}) \log P(\mathbf{E}_{uv}=e \mid \boldsymbol{Q}^{(\boldsymbol{S}_{\mathcal{D}}=\boldsymbol{e})}).$$
(16)

In other words, we estimate the *joint* impact of pairs in \mathcal{D} on the entropy of \mathbf{E}_{uv} . The expectation in Eq. 15, which involves a sum over all possible outcomes of $\mathbf{E}_{\mathcal{D}}$, quickly becomes intractable for large $|\mathcal{D}|$. However, one can easily obtain a sample $e \sim P(\mathbf{E}_{\mathcal{D}})$ by sampling $e_{uv} \sim P(\mathbf{E}_{uv})$ for every $\mathbf{E}_{uv} \in \mathbf{E}_{\mathcal{D}}$, which allows a Monte-Carlo estimation of this sum. For generality, assume we have *m* subsets $\mathcal{D}_1, \ldots, \mathcal{D}_m, \mathcal{D}_i \subseteq \mathcal{E}$ and *n* samples $e_i^1, \ldots, e_i^n \sim P(\mathbf{E}_{\mathcal{D}_i})$ for each \mathcal{D}_i . Given this, we then define the acquisition function

$$a^{\text{JEIG}}(u,v) \triangleq H(\mathbf{E}_{uv}) - \frac{1}{mn} \sum_{i=1}^{m} \sum_{l=1}^{n} H(\mathbf{E}_{uv} \mid \boldsymbol{Q}^{(\boldsymbol{S}_{\mathcal{D}_{i}} = \boldsymbol{e}_{i}^{l})}).$$
(17)

361 For a^{JEIG} , we only need to execute Alg. 2 (mean-field) mn times, and in practice, we observe good 362 performance with small values of m and n. In Appendix B.3, we present Alg. 4 which describes the details of this method. Using m subsets $\mathcal{D}_1, \ldots, \mathcal{D}_m$ with each $|\mathcal{D}_i| \ll |\mathcal{E}|$ yields a number of benefits: 364 (i) The Monte-Carlo estimation of the expectation in Eq. 15 becomes more accurate for smaller n365 when $|\mathcal{D}_i|$ is smaller, which reduces the number of times Alg. 2 needs to be executed; (ii) If $\mathcal{D}_i = \mathcal{E}$, the conditional mean-field approximation $Q^{(S_{D_i}=e)}$ is computed based on a similarity matrix where 366 all pairs $(u, v) \in \mathcal{E}$ are sampled from $S_{uv} \sim P(E_{uv})$, which will lead to extreme selection bias for 367 the following reason: The probability $P(E_{uv})$ (which is computed using Q) may already be biased 368 (in particular in early iterations when S contains incomplete/wrong information). Then, running Alg. 369 2 from scratch with a new similarity matrix fully augmented with biased information, will exaggerate 370 the bias further; (iii) Using m different subsets makes the estimator in Eq. 17 generic and flexible, 371 but also captures more information about E_{uv} , while remaining efficient and avoiding exaggerated 372 selection bias. 373

374 375

376

4 EXPERIMENTS

377 In this section, we describe our experimental studies, where extensive additional results are presented in Appendix C.

4.1 EXPERIMENTAL SETUP

379 380

Datasets. In this paper, we use the datasets investigated by (Aronsson & Chehreghani, 2024): 20newsgroups, CIFAR10, cardiotocography, ecoli, forest type mapping, user knowledge modeling, MNIST and synthetic. For all datasets, a random subset of at most N = 1000 objects are considered for the active CC experiments. See Appendix C.3 for details about all datasets.

Correlation Clustering Algorithm. We use the local search CC algorithm proposed by (Aronsson & Chehreghani, 2024) on line 4 of Alg. 1. It is highly robust to noise in S and dynamically determines the number of clusters. The details of this algorithm are described in Appendix D.

Ground-truth similarities. For each experiment, we are given a dataset X with ground-truth labels **c***, where the ground-truth labels are only used for evaluations. Then, for each $(u, v) \in \mathcal{E}$ in a dataset, we set S_{uv}^* to +1 if u and v belong to the same cluster, and -1 otherwise.

Oracles. We investigate four different oracles in Alg. 1: (i) **Oracle 1**. Returns S_{uv}^* with probability 391 $1 - \gamma$ or a uniform random value in [-1, +1] with probability γ ; (ii) **Oracle 2**. Returns a value 392 sampled from $\mathcal{N}(S_{uv}^*,\gamma)$ (i.e., Gaussian centered at ground-truth similarity with variance γ); (iii) 393 **Oracle 3.** Returns S_{uv}^* with probability $1 - \gamma$ or $-S_{uv}^*$ with probability γ (i.e., we flip the sign with 394 probability γ)²; **Oracle 4**. We split the dataset into two disjoint parts $\mathbf{X} = \mathbf{X}_{\text{train}} \cup \mathbf{X}_{\text{test}}$. Then, we 395 train a pairwise prediction model $f_{\theta}: \mathbf{X} \times \mathbf{X} \to [-1, +1]$ on $\mathbf{X}_{\text{train}}$, where ground-truth similarities 396 S^* are used as labels. Given any two data points $\mathbf{x}_u, \mathbf{x}_v \in \mathbf{X}$, we can predict their similarity as 397 $f_{\theta}(\mathbf{x}_u, \mathbf{x}_v) \in [-1, +1]$. We then perform the CC experiments on data points in \mathbf{X}_{test} , and the oracle 398 always returns the similarity $f_{\theta}(\mathbf{x}_u, \mathbf{x}_v)$. The ground-truth similarities of data points in \mathbf{X}_{test} are 399 never used when training f_{θ} . We defer a detailed description of oracle 4 to Appendix C.2. The 400 motivation for these oracles are as follows. Oracles 1-3 correspond to cases where the oracle provides 401 unbiased information about S^* (but noisy, with different noise models), allowing recovery of the 402 ground-truth clustering solution c^* . This is considered by previous work (Mazumdar & Saha, 2017b; Silwal et al., 2023; Aronsson & Chehreghani, 2024). Oracle 4 may provide biased similarities due 403 to noise/ambiguity in feature space, and exact recovery of c^* may not be possible. This method is 404 suggested by, e.g., (Bansal et al., 2004; Silwal et al., 2023) to compute pairwise similarities for CC. 405

Initial similarities. Let \mathcal{E}^0 be a uniform random subset of \mathcal{E} (where $|\mathcal{E}^0| \ll |\mathcal{E}|$). Then, for all $(u, v) \in \mathcal{E}^0$, we initialize the current similarity matrix as $S_{uv}^0 = S_{uv}^*$ for oracles 1-3 and $S_{uv}^0 = f_{\theta}(\mathbf{x}_u, \mathbf{x}_v)$ for oracle 4. We then set $S_{uv}^0 = 0$ for $(u, v) \in \mathcal{E} \setminus \mathcal{E}^0$. Having $|\mathcal{E}^0| = 0$ or $|\mathcal{E}^0| > 0$ corresponds to a *cold-start* or *warm-start* setting, respectively. In this paper, like most previous work on active learning, we focus on a warm-start setting. See Appendix C.3 for the value of $|\mathcal{E}^0|$ for each dataset ($|\mathcal{E}^0|$ is chosen based on the size N of each dataset).

Repeated queries. In general, Alg 1 supports multiple queries for the same pairwise similarity. This assumes each query for the same pair provides more information about the underlying distribution, which would be applicable to oracles 1-3. This is a common approach in active learning to deal with noisy oracles (Sheng et al., 2008; Settles, 2009). However, *we do not consider multiple queries for the same pair in our experiments*, as we found the difference in performance to be very small.

417 Acquisition functions. We have introduced four novel acquisition functions in this paper: a^{Entropy} (Eq. 8), $a^{\text{EIG-O}}$ (Eq. 13), $a^{\text{EIG-P}}$ (Eq. 34) and a^{JEIG} (Eq. 17). We compare these methods with maxesp 418 419 and maxmin from (Aronsson & Chehreghani, 2024). In short, both maxmin and maxexp aim to query pairs with small absolute similarity that belong to triples (u, v, w) that violate the transitive 420 property of pairwise similarities. In other words, the goal is to reduce the inconsistency of S by 421 resolving violations of the transitive property in triples. In Appendix C.1, we include a detailed 422 explanation of these methods. Finally, we include a simple baseline $a^{\text{Uniform}}(u, v) \sim \text{Uniform}(0, 1)$ 423 which selects pairs uniformly at random. The work in (Aronsson & Chehreghani, 2024) compares 424 maxexp/maxmin to a pivot-based active CC algorithm called QECC (García-Soriano et al., 2020) and 425 two adapted state-of-the-art active constraint clustering methods, called COBRAS (van Craenendonck 426 et al., 2018a) and nCOBRAS (Soenen et al., 2021). However, these methods perform very poorly 427 compared to *maxexp* and *maxmin* in a noisy setting, so we exclude them here.

⁴³⁰ ²Oracles 1-3 are equivalent if $\gamma = 0$. However, zero noise is unrealistic in practice. Also, it leads to fully 431 consistent information in the similarities S, which makes the CC problem (minimization of Eq. 1) trivial (i.e., it is no longer NP-hard).



Figure 1: Results for oracle 1 with noise level $\gamma = 0.4$.

Batch diversity. In this paper, we consider single-sample acquisition functions that do not explicitly 451 consider the joint informativeness among the elements in a batch \mathcal{B} . This has the benefit of avoiding 452 the combinatorial complexity of selecting an optimal batch, which is a common problem for batch 453 active learning (Ren et al., 2021). However, the work in (Kirsch et al., 2023) proposes a simple 454 method for improving the batch diversity for single-sample acquisition functions using noise. In 455 this paper, we utilize the *power* acquisition method. Given some acquisition function a^{X} , this corresponds to $a^{\text{PowerX}}(u, v) = \log(a^{\bar{X}}(u, v)) + \epsilon_{uv}$ where $\epsilon_{uv} \sim \text{Gumbel}(0; 1)$. This is used for all 456 457 information-theoretic acquisition functions proposed in this paper. We observe no benefit of this for 458 maxmin/maxexp, likely due to their inherent randomness.

459 Hyperparameters. Unless otherwise specified, the following hyperparameters are used. The batch 460 size B depends on the dataset (since each dataset is of different size N). See Appendix C.3 for the 461 value of \overline{B} for each dataset. For all information-theoretic acquisition functions (which depend on P^{Gibbs} , see Eq. 3), we set $\beta = 3$. For $a^{\text{EIG-O}}$ and $a^{\text{EIG-P}}$, we set $|\mathcal{E}^{\text{EIG}}| = 20N$ (see Appendix B.2 for details). For a^{JEIG} we set m = 5, n = 50, $|\mathcal{D}_i| = \lceil |\mathcal{E}|/50 \rceil$ (i.e., 2% of all pairs) and each \mathcal{D}_i is 462 463 selected to contain pairs with large entropy (see Appendix B.3 for details). Finally, see Appendix C.7 464 for a detailed analysis and discussion of all hyperparameters. 465

466 Performance evaluation. In each iteration of the active CC procedure (Alg. 1), we calculate the 467 Adjusted Rand Index (ARI) between the current clustering c^i and the ground truth clustering c^* (i.e., 468 ground truth labels of dataset). Intuitively, ARI measures how similar the two clusterings are, where 469 a value of 1 indicates they are identical. In Appendix C.4, we report the performance w.r.t. other 470 evaluation metrics. Each active learning procedure is repeated 10 times with different random seeds, where the standard deviation is indicated by a shaded color or an error bar. 471

4.2 RESULTS 473

474

472

448

449 450

Figures 1-2 show the results for different datasets for oracle 1 and oracle 4, respectively. In Appendix 475 C.4, we include the results for all oracles w.r.t. different performance metrics. We observe that 476 all the information-theoretic acquisition functions introduced in this paper significantly outperform 477 the baseline methods. In addition, the acquisition functions based on information gain ($a^{\text{EIG-O}}$, 478 $a^{\text{EIG-P}}$ and a^{JEIG}) consistently outperform a^{Entropy} . This indicates the effectiveness of augmenting the 479 similarity matrix S with pseudo-similarities predicted by the current model Q as $S_{uv} \sim P(E_{uv} \mid Q)$, 480 before quantifying the model uncertainty. The acquisition functions based on information gain perform rather similarly. This is due to the fact that all of them are based on closely connected quantities (as described in Section 3.3). However, a^{JEIG} is consistently among the best performing acquisition functions, while also being more computationally efficient compared to $a^{\text{EIG-O}}$ and $a^{\text{EIG-O}}$ 481 482 483 (see Appendix C.6 for an investigation of the runtimes of all methods). Because of this, we exclude 484 $a^{\text{EIG-O}}$ and $a^{\text{EIG-P}}$ in some cases due to their computational inefficiency. Both *maxmin* and *maxexp* 485 perform significantly worse. This is likely because they spend too many queries with the goal of



Figure 2: Results for all datasets for oracle 4.

resolving the inconsistency of S (see Appendix C.1 for details). However, the CC algorithm used (described in Appendix D) is robust to inconsistency in S. Finally, in Appendix C.5 we report our experiments on a small synthetic dataset with N = 70 objects using a small batch size B = 5. The purpose of this experiment is to further illustrate the benefit of the acquisition functions based on information gain, when differences due to batch diversity are eliminated.

510 511 512

513

502

504 505

506

507

508

509

4.3 SENSITIVITY ANALYSIS

514 We here investigate the sensitivity of acquisi-515 tion functions when varying the noise level and 516 batch size. All the results in this section are per-517 formed on the synthetic dataset using oracle 1. 518 Figures 3(a)-3(b) show the results when varying 519 the noise level γ and batch size B, respectively. 520 The y-axis corresponds to the area under the curve (AUC) of the active learning plot w.r.t. the 521 respective performance metric (i.e., ARI) where 522 higher is better. We see that our acquisition func-523 tions are very robust to noise. In addition, the 524 benefit of our proposed acquisition functions 525 increases with larger noise levels. This is con-526 sistent with previous work on active learning,



Figure 3: Results on the synthetic dataset when varying the noise level γ and the batch size B.

where the benefit of many acquisition functions over uniform selection increases as the complexity of
the problem increases. Expectedly, the performance decreases slightly as the batch size increases.
However, the performance of our acquisition functions remains good even with large batch sizes.

531 532

5 CONCLUSION

533 534

In this paper, we proposed four effective information-theoretic acquisition functions to be used
for active CC: a^{Entropy}, a^{EIG-O}, a^{EIG-P} and a^{JEIG}. All of our methods significantly outperform the
baseline methods by utilizing *model uncertainty*. We investigated the effectiveness of these methods
via extensive experimental studies. The acquisition functions based on information gain (a^{EIG-O}, a^{EIG-P} and a^{JEIG}) were consistently the best performing, where a^{JEIG} has the benefit of being more computationally efficient.

540 REFERENCES

578

579

- Nir Ailon, Moses Charikar, and Alantha Newman. Aggregating inconsistent information: Ranking and clustering. J. ACM, 55(5):23:1–23:27, 2008. doi: 10.1145/1411509.1411513.
- Alan Akbik, Duncan Blythe, and Roland Vollgraf. Contextual string embeddings for sequence
 labeling. In *COLING 2018, 27th International Conference on Computational Linguistics*, pp. 1638–1649, 2018.
- Linus Aronsson and Morteza Haghir Chehreghani. Correlation clustering with active learning of pairwise similarities. *Transactions on Machine Learning Research*, 2024. ISSN 2835-8856. URL https://openreview.net/forum?id=Ryf1TVCjBz.
- Pranjal Awasthi and Reza Bosagh Zadeh. Supervised clustering. In *Advances in Neural Information Processing Systems (NIPS)*, pp. 91–99, 2010.
- Pranjal Awasthi, Maria Florina Balcan, and Konstantin Voevodski. Local algorithms for interactive clustering. *Journal of Machine Learning Research*, 18(3):1–35, 2017.
- Maria-Florina Balcan and Avrim Blum. Clustering with interactive feedback. In Yoav Freund, László
 Györfi, György Turán, and Thomas Zeugmann (eds.), *Algorithmic Learning Theory*, pp. 316–328,
 Berlin, Heidelberg, 2008. Springer Berlin Heidelberg. ISBN 978-3-540-87987-9.
- Nikhil Bansal, Avrim Blum, and Shuchi Chawla. Correlation clustering. *Machine Learning*, 56(1-3):
 89–113, 2004. doi: 10.1023/B:MACH.0000033116.57574.95.
- Sugato Basu, Arindam Banerjee, and Raymond J. Mooney. Active semi-supervision for pairwise constrained clustering. In *SDM*, 2004.
- Francesco Bonchi, Aristides Gionis, Francesco Gullo, and Antti Ukkonen. Chromatic correlation clustering. In *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '12, pp. 1321–1329, 2012. doi: 10.1145/2339530.2339735.
 URL https://doi.org/10.1145/2339530.2339735.
- Francesco Bonchi, Aristides Gionis, and Antti Ukkonen. Overlapping correlation clustering. *Knowl. Inf. Syst.*, 35(1):1–32, 2013. doi: 10.1007/s10115-012-0522-9.
- Francesco Bonchi, David García-Soriano, and Edo Liberty. Correlation clustering: from theory to practice. In *The 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, KDD '14*, pp. 1972. ACM, 2014. doi: 10.1145/2623330.2630808.
- Marco Bressan, Nicolò Cesa-Bianchi, Andrea Paudice, and Fabio Vitale. Correlation clustering with
 adaptive similarity queries. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox,
 and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran
 Associates, Inc., 2019.
 - Moses Charikar, Venkatesan Guruswami, and Anthony Wirth. Clustering with qualitative information. *J. Comput. Syst. Sci.*, 71(3):360–383, 2005. doi: 10.1016/j.jcss.2004.10.012.
- Morteza Haghir Chehreghani. *Information-theoretic validation of clustering algorithms*. PhD thesis, 2013.
- 583
 584
 584
 585
 585
 Morteza Haghir Chehreghani. Shift of pairwise similarities for data clustering. *Mach. Learn.*, 112(6): 2025–2051, 2023. doi: 10.1007/S10994-022-06189-6. URL https://doi.org/10.1007/ s10994-022-06189-6.
- Morteza Haghir Chehreghani and Mostafa Haghir Chehreghani. Learning representations from dendrograms. *Mach. Learn.*, 109(9-10):1779–1802, 2020. doi: 10.1007/s10994-020-05895-3.
- Morteza Haghir Chehreghani, Alberto Giovanni Busetto, and Joachim M. Buhmann. Information theoretic model validation for spectral clustering. In Neil D. Lawrence and Mark Girolami (eds.), *Proceedings of the Fifteenth International Conference on Artificial Intelligence and Statistics*, volume 22 of *Proceedings of Machine Learning Research*, pp. 495–503, La Palma, Canary Islands, 21–23 Apr 2012. PMLR. URL https://proceedings.mlr.press/v22/haghir12. html.

626

632

633

635

- 594 Erik D. Demaine, Dotan Emanuel, Amos Fiat, and Nicole Immorlica. Correlation clustering in general weighted graphs. Theor. Comput. Sci., 361(2-3):172-187, 2006. doi: 10.1016/j.tcs.2006.05.008. 596
- Micha Elsner and Warren Schudy. Bounding and comparing methods for correlation clustering 597 beyond ILP. In Proceedings of the Workshop on Integer Linear Programming for Natural Language 598 Processing, pp. 19–27, Boulder, Colorado, June 2009. Association for Computational Linguistics.
- 600 David García-Soriano, Konstantin Kutzkov, Francesco Bonchi, and Charalampos Tsourakakis. 601 Query-efficient correlation clustering. In Proceedings of The Web Conference 2020, WWW '20, pp. 1468–1478, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 602 9781450370233. doi: 10.1145/3366423.3380220. 603
- 604 Lise Getoor and Ashwin Machanavajjhala. Entity resolution: Theory, practice & open challenges. 605 Proc. VLDB Endow., 5(12):2018–2019, 2012. 606
- Aristides Gionis, Heikki Mannila, and Panayiotis Tsaparas. Clustering aggregation. ACM Trans. 607 Knowl. Discov. Data, 1(1):4, 2007. doi: 10.1145/1217299.1217303. 608
- 609 Francesco Gullo, Domenico Mandaglio, and Andrea Tagarelli. A combinatorial multi-armed bandit ap-610 proach to correlation clustering. Data Min. Knowl. Discov., 37(4):1630-1691, 2023. doi: 10.1007/ 611 S10618-023-00937-5. URL https://doi.org/10.1007/s10618-023-00937-5.
- Oktie Hassanzadeh, Fei Chiang, Renée J. Miller, and Hyun Chul Lee. Framework for evaluating 613 clustering algorithms in duplicate detection. Proc. VLDB Endow., 2(1):1282-1293, 2009. doi: 614 10.14778/1687627.1687771. 615
- 616 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. CoRR, abs/1512.03385, 2015. 617
- 618 T. Hofmann and J.M. Buhmann. Pairwise data clustering by deterministic annealing. *IEEE Transac*-619 tions on Pattern Analysis and Machine Intelligence, 19(1):1–14, 1997. doi: 10.1109/34.566806. 620
- Thomas Hofmann, Jan Puzicha, and Joachim M. Buhmann. Unsupervised texture segmentation in a 621 deterministic annealing framework. IEEE Trans. Pattern Anal. Mach. Intell., 20(8):803-818, 1998. 622 doi: 10.1109/34.709593. URL https://doi.org/10.1109/34.709593. 623
- 624 Sanna Jarl, Linus Aronsson, Sadegh Rahrovani, and Morteza Haghir Chehreghani. Active learning of 625 driving scenario trajectories. Eng. Appl. Artif. Intell., 113:104972, 2022.
- Markelle Kelly, Rachel Longjohn, and Kolby Nottingham. The uci machine learning repository, 2023. 627 URL https://archive.ics.uci.edu. 628
- 629 Sungwoong Kim, Sebastian Nowozin, Pushmeet Kohli, and Chang Dong Yoo. Higher-order correla-630 tion clustering for image segmentation. In Advances in Neural Information Processing Systems 24 (NIPS), pp. 1530-1538, 2011. 631
- Andreas Kirsch and Yarin Gal. Unifying approaches in active learning and active sampling via fisher information and information-theoretic quantities. Transactions on Machine Learning Research, 634 2022. ISSN 2835-8856. URL https://openreview.net/forum?id=UVDAKQANOW. Expert Certification.
- Andreas Kirsch, Sebastian Farquhar, Parmida Atighehchian, Andrew Jesson, Frédéric Branchaud-637 Charron, and Yarin Gal. Stochastic batch acquisition: A simple baseline for deep active 638 learning. Transactions on Machine Learning Research, 2023. ISSN 2835-8856. URL 639 https://openreview.net/forum?id=vcHwQyNBjW. Expert Certification. 640
- 641 Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009.
- 642 Yuko Kuroki, Atsushi Miyauchi, Francesco Bonchi, and Wei Chen. Query-efficient correlation 643 clustering with noisy oracle. CoRR, abs/2402.01400, 2024. doi: 10.48550/ARXIV.2402.01400. 644 URL https://doi.org/10.48550/arXiv.2402.01400. 645
- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied 646 to document recognition. Proc. IEEE, 86(11):2278-2324, 1998. doi: 10.1109/5.726791. URL 647 https://doi.org/10.1109/5.726791.

| 648 649 650 | Arya Mazumdar and Barna Saha. Query complexity of clustering with side information. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 30. Curran Associates, Inc., 2017a. |
|--|--|
| 651 652 653 654 | Arya Mazumdar and Barna Saha. Clustering with noisy queries. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 30. Curran Associates, Inc., 2017b. |
| 655 656 657 | Andrew McCallum and Ben Wellner. Conditional models of identity uncertainty with application to noun coreference. In <i>Advances in Neural Information Processing Systems 17 [Neural Information Processing Systems, NIPS</i> , pp. 905–912, 2004. |
| 659 660 661 | Anirudh Ramachandran, Nick Feamster, and Santosh Vempala. Filtering spam with behavioral blacklisting. In <i>Proceedings of the 14th ACM Conference on Computer and Communications Security</i> , pp. 342–351, 2007. ISBN 9781595937032. doi: 10.1145/1315245.1315288. |
| 662 663 664 | Pengzhen Ren, Yun Xiao, Xiaojun Chang, Po-Yao Huang, Zhihui Li, Brij B. Gupta, Xiaojiang Chen, and Xin Wang. A survey of deep active learning. ACM Comput. Surv., 54(9), oct 2021. ISSN 0360-0300. doi: 10.1145/3472291. URL https://doi.org/10.1145/3472291. |
| 666 667 668 669 | Nicholas Roy and Andrew McCallum. Toward optimal active learning through sampling estimation of error reduction. In <i>Proceedings of the Eighteenth International Conference on Machine Learning</i> , ICML '01, pp. 441–448, San Francisco, CA, USA, 2001. Morgan Kaufmann Publishers Inc. ISBN 1558607781. |
| 670 671 | Neil Rubens, Mehdi Elahi, Masashi Sugiyama, and Dain Kaplan. Active Learning in Recommender Systems, pp. 809–846. 2015. ISBN 978-1-4899-7637-6. doi: 10.1007/978-1-4899-7637-6_24. |
| 672 673 674 675 676 677 | Barna Saha and Sanjay Subramanian. Correlation Clustering with Same-Cluster Queries Bounded by Optimal Cost. In Michael A. Bender, Ola Svensson, and Grzegorz Herman (eds.), 27th Annual European Symposium on Algorithms (ESA 2019), volume 144 of Leibniz International Proceedings in Informatics (LIPIcs), pp. 81:1–81:17, Dagstuhl, Germany, 2019. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik. ISBN 978-3-95977-124-5. doi: 10.4230/LIPIcs.ESA.2019.81. |
| 678 679 | Burr Settles. Active learning literature survey. Computer Sciences Technical Report 1648, University of Wisconsin–Madison, 2009. |
| 680 681 682 683 684 | Victor S. Sheng, Foster Provost, and Panagiotis G. Ipeirotis. Get another label? improving data quality and data mining using multiple, noisy labelers. In <i>Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining</i> , KDD '08, pp. 614–622, New York, NY, USA, 2008. Association for Computing Machinery. ISBN 9781605581934. doi: 10.1145/1401890.1401965. URL https://doi.org/10.1145/1401890.1401965. |
| 686 687 688 | Zhao Shuyang, Toni Heittola, and Tuomas Virtanen. Active learning for sound event detection. <i>IEEE/ACM Trans. Audio, Speech and Lang. Proc.</i> , 28:2895–2905, nov 2020. ISSN 2329-9290. doi: 10.1109/TASLP.2020.3029652. |
| 689 690 691 692 | Sandeep Silwal, Sara Ahmadian, Andrew Nystrom, Andrew McCallum, Deepak Ramachandran, and Seyed Mehran Kazemi. Kwikbucks: Correlation clustering with cheap-weak and expensive-strong signals. In <i>The Eleventh International Conference on Learning Representations</i> , 2023. URL https://openreview.net/forum?id=p0JSSalAuV. |
| 694 695 696 | Jonas Soenen, Sebastijan Dumancic, Hendrik Blockeel, Toon Van Craenendonck, F Hutter, K Kersting, J Lijffijt, and I Valera. Tackling noise in active semi-supervised clustering, 2021. ISSN 978-3-030-67661-2. |
| 697 698 | Jiliang Tang, Yi Chang, Charu Aggarwal, and Huan Liu. A survey of signed network mining in social media. <i>ACM Comput. Surv.</i> , 49(3), 2016. doi: 10.1145/2956185. |
| 700 701 | Erik Thiel, Morteza Haghir Chehreghani, and Devdatt P. Dubhashi. A non-convex optimization approach to correlation clustering. In <i>The Thirty-Third AAAI Conference on Artificial Intelligence</i> , <i>AAAI</i> , pp. 5159–5166. AAAI Press, 2019. doi: 10.1609/aaai.v33i01.33015159. |

| Toon van Craenendonck, Sebastijan Dumancic, Elia Van Wolputte, and Hendrik Blockeel. Cobras: Interactive clustering with pairwise queries. In <i>International Symposium on Intelligent Data</i> <i>Analysis</i>, 2018b. Simon Viet Johansson, Hampus Gummesson Svensson, Esben Bjerrum, Alexander Schliep, Morteza Haghir Chehreghani, Christian Tyrchan, and Ola Engkvist. Using active learning to develop machine learning models for reaction yield prediction. <i>Molecular Informatics</i>, 41(12):2200043, 2022. doi: https://doi.org/10.1002/minf.202200043. Songbai Yan, Kamalika Chaudhuri, and Tara Javidi. Active learning with logged data. In Jennifer G. Dy and Andreas Krause (eds.), <i>Proceedings of the 35th International Conference on Machine Learning (ICML)</i>, volume 80 of <i>Proceedings of Machine Learning Research</i>, pp. 5517–5526. PMLR, 2018. | Toon van Craenendonck, Sebastijan Dumancic, Elia Van Wolputte, and Hendrik Blockeel. Cobras: Interactive clustering with pairwise queries. In <i>International Symposium on Intelligent Data</i> <i>Analysis</i>, 2018b. Simon Viet Johansson, Hampus Gummesson Svensson, Esben Bjerrum, Alexander Schliep, Morteza Haghir Chehreghani, Christian Tyrchan, and Ola Engkvist. Using active learning to develop machine learning models for reaction yield prediction. <i>Molecular Informatics</i>, 41(12):2200043, 2022. doi: https://doi.org/10.1002/minf.202200043. Songbai Yan, Kamalika Chaudhuri, and Tara Javidi. Active learning with logged data. In Jennifer G. Dy and Andreas Krause (eds.), <i>Proceedings of the 35th International Conference on Machine Learning (ICML)</i>, volume 80 of <i>Proceedings of Machine Learning Research</i>, pp. 5517–5520. PMLR, 2018. | 702 703 704 | Toon van Craenendonck, Sebastijan Dumancic, and Hendrik Blockeel. Cobra: A fast and simple method for active clustering with pairwise constraints. In <i>International Joint Conference on Artificial Intelligence</i> , 2018a. |
|---|--|-------------------|---|
| Iterative clustering with pairwise queries. In <i>International Symposium on Intelligent Data Analysis</i>, 2018b. Simon Viet Johanson, Hampus Gummesson Svensson, Esben Bjerrum, Alexander Schliep, Morteza Haghir Chehreghani, Christian Tyrchan, and Ola Engkvist. Using active learning to develop machine learning models for reaction yield prediction. <i>Molecular Informatics</i>, 41(12):2200043, 2022. doi: https://doi.org/10.1002/minf.202200043. Songbai Yan, Kamalika Chaudhuri, and Tara Javidi. Active learning with logged data. In Jennifer G. Dy and Andreas Krause (eds.), <i>Proceedings of the 35th International Conference on Machine Learning (ICML)</i>, volume 80 of <i>Proceedings of Machine Learning Research</i>, pp. 5517–5526. PMLR, 2018. | Flort na traditional constraint of the pairwise queries. In International Symposium on Intelligent Data Analysis, 2018b. Simon Viet Johansson, Hampus Gummesson Svensson, Esben Bjerrum, Alexander Schliep, Morteza Haghir Chehreghani, Christian Tyrchan, and Ola Engkvist. Using active learning to develop machine learning models for reaction yield prediction. Molecular Informatics, 41(12):2200043, 2022. doi: https://doi.org/10.1002/minf.202200043. Songbai Yan, Kamalika Chaudhuri, and Tara Javidi. Active learning with logged data. In Jennifer G. Dy and Andreas Krause (eds.), Proceedings of the 35th International Conference on Machine Learning (ICML), volume 80 of Proceedings of Machine Learning Research, pp. 5517–5526. PMLR, 2018. | 705 | Toon van Craenendonck, Sebastijan Dumancic, Elia Van Wolnutte, and Hendrik Blockeel, Cobras: |
| Analysis, 2018b. Simon Viet Johansson, Hampus Gummesson Svensson, Esben Bjerrum, Alexander Schliep, Morteza Haghir Chehreghani, Christian Tyrchan, and Ola Engkvist. Using active learning to develop machine learning models for reaction yield prediction. <i>Molecular Informatics</i>, 41(12):2200043, 2022. doi: https://doi.org/10.1002/minf.202200043. Songbai Yan, Kamalika Chaudhuri, and Tara Javidi. Active learning with logged data. In Jennifer G. Dy and Andreas Krause (eds.), <i>Proceedings of the 35th International Conference on Machine Learning (ICML)</i>, volume 80 of <i>Proceedings of Machine Learning Research</i>, pp. 5517–5526. PMLR, 2018. | Analysis, 2018b. Simon Viet Johansson, Hampus Gummesson Svensson, Esben Bjerrum, Alexander Schliep, Morteza Haghir Chehreghani, Christian Tyrchan, and Ola Engkvist. Using active learning to develop machine learning models for reaction yield prediction. Molecular Informatics, 41(12):2200043, 2022. doi: https://doi.org/10.1002/minf.202200043. Songbai Yan, Kamalika Chaudhuri, and Tara Javidi. Active learning with logged data. In Jennife G. Dy and Andreas Krause (eds.), Proceedings of the 35th International Conference on Machine Learning (ICML), volume 80 of Proceedings of Machine Learning Research, pp. 5517–5526. PMLR, 2018. | 706 | Interactive clustering with pairwise queries In International Symposium on Intelligent Data |
| Simon Viet Johansson, Hampus Gummesson Svensson, Esben Bjerrum, Alexander Schliep, Morteza Haghir Chehreghani, Christian Tyrchan, and Ola Engkvist. Using active learning to develop machine learning models for reaction yield prediction. <i>Molecular Informatics</i>, 41(12):2200043, 2022. doi: https://doi.org/10.1002/mini.202200043. Songbai Yan, Kamalika Chaudhuri, and Tara Javidi. Active learning with logged data. In Jennifer G. Dy and Andreas Krause (eds.), <i>Proceedings of the 35th International Conference on Machine Learning (ICML)</i>, volume 80 of <i>Proceedings of Machine Learning Research</i>, pp. 5517–5526. PMLR, 2018. | Simon Viet Johansson, Hampus Gummesson Svensson, Esben Bjerrum, Alexander Schliep, Morteza Haghir Chehreghani, Christian Tyrchan, and Ola Engkvist. Using active learning to develop machine learning models for reaction yield prediction. <i>Molecular Informatics</i>, 41(12):2200043, 2022. doi: https://doi.org/10.1002/mini.202200043. Songbai Yan, Kamalika Chaudhuri, and Tara Javidi. Active learning with logged data. In Jennifer G. Dy and Andreas Krause (eds.), <i>Proceedings of the 35th International Conference on Machine Learning (ICML)</i>, volume 80 of <i>Proceedings of Machine Learning Research</i>, pp. 5517–5526. PMLR, 2018. | 707 | Analysis, 2018b. |
| Simon Viet Johansson, Hampus Gummesson Svensson, Esben Bjerrum, Alexander Schliep, Morteza Haghir Chehreghani, Christian Tyrchan, and Ola Engkvist. Using active learning to develop machine learning models for reaction yield prediction. <i>Molecular Informatics</i>, 41(12):2200043, 2022. doi: https://doi.org/10.1002/minf.202200043. Songbai Yan, Kamalika Chaudhuri, and Tara Javidi. Active learning with logged data. In Jennifer G. Dy and Andreas Krause (eds.), <i>Proceedings of the 35th International Conference on Machine Learning (ICML)</i>, volume 80 of <i>Proceedings of Machine Learning Research</i>, pp. 5517–5526. PMLR, 2018. | Simon Viet Johansson, Hampus Gummesson Svensson, Esben Bjerrum, Alexander Schliep, Morteza Haghir Chehreghani, Christian Tyrchan, and Ola Engkvist. Using active learning to develop machine learning models for reaction yield prediction. <i>Molecular Informatics</i>, 41(12):2200043, 2022. doi: https://doi.org/10.1002/minf.202200043. Songbai Yan, Kamalika Chaudhuri, and Tara Javidi. Active learning with logged data. In Jennifer G. Dy and Andreas Krause (eds.). <i>Proceedings of the 35th International Conference on Machine Learning (ICML)</i>, volume 80 of <i>Proceedings of Machine Learning Research</i>, pp. 5517–5526. PMLR, 2018. | 700 | |
| Fragmir Cherregnani, Curistian Tyrenan, and Ola Engevist. Osing active learning to develop machine learning models for reaction yield prediction. <i>Molecular Informatics</i>, 41(12):2200043, 2022. doi: https://doi.org/10.1002/minf.202200043. Songbai Yan, Kamalika Chaudhuri, and Tara Javidi. Active learning with logged data. In Jennifer G. Dy and Andreas Krause (eds.), <i>Proceedings of the 35th International Conference on Machine Learning (ICML)</i>, volume 80 of <i>Proceedings of Machine Learning Research</i>, pp. 5517–5526. PMLR, 2018. | Fragmic Inchargenant, Christian Tyrchan, and Ola Engyvist. Using active learning to develop machine learning models for reaction yield prediction. Molecular Informatics, 41(12):2200043, 2022. doi: https://doi.org/10.1002/minf.202200043. Songbai Yan, Kamalika Chaudhuri, and Tara Javidi. Active learning with logged data. In Jennifer G. Dy and Andreas Krause (eds.), Proceedings of the 35th International Conference on Machine Learning (ICML), volume 80 of Proceedings of Machine Learning Research, pp. 5517–5526. PMLR, 2018. PMLR, 2018. | 709 | Simon Viet Johansson, Hampus Gummesson Svensson, Esben Bjerrum, Alexander Schliep, Morteza |
| Interime rearing incores for reaction influe periodicular informatics, 47(12),2200045, Songbai Yan, Kamalika Chaudhuri, and Tara Javidi. Active learning with logged data. In Jennifer G. Dy and Andreas Krause (eds.), <i>Proceedings of the 35th International Conference on Machine Learning (ICML)</i>, volume 80 of <i>Proceedings of Machine Learning Research</i>, pp. 5517–5526. PMLR, 2018. | Internite Jeaning Jobes for Course of Proceedings of the 35th International Conference on Machine Learning (ICML), volume 80 of Proceedings of Machine Learning Research, pp. 5517–5526. PMLR, 2018. | 711 | machine learning models for reaction yield prediction <i>Molecular Informatics</i> 41(12):2200043 |
| Songbai Yan, Kamalika Chaudhuri, and Tara Javidi. Active learning with logged data. In Jennifer G. Dy and Andreas Krause (eds.), Proceedings of the 35th International Conference on Machine Learning (ICML), volume 80 of Proceedings of Machine Learning Research, pp. 5517–5526. PMLR, 2018. PMLR, 2018. PMLR | Songbai Yan, Kamalika Chaudhuri, and Tara Javidi. Active learning with logged data. In Jennifer G. Dy and Andreas Krause (eds.), <i>Proceedings of the 35th International Conference on Machine Learning (ICML)</i>, volume 80 of <i>Proceedings of Machine Learning Research</i>, pp. 5517–5526. PMLR, 2018. PMLR, 2018. | 712 | 2022. doi: https://doi.org/10.1002/minf.202200043. |
| Songbai Yan, Kamalika Chaudhur, and Tara Javidi. Active learning with logged data. In Jenniter G. Dy and Andreas Krause (eds.), Proceedings of Machine Learning Research, pp. 5517–5526. PMLR, 2018. PMLR, 2018. <td>Songbai Yan, Kamalika Chaudhuri, and Tara Javidi. Active learning with logged data. In Jennifer G. Dy and Andreas Krause (eds.). Proceedings of Machine Learning Research, pp. 5517–5526. PMLR, 2018. PMLR, 2018.<!--</td--><td>713</td><td></td></td> | Songbai Yan, Kamalika Chaudhuri, and Tara Javidi. Active learning with logged data. In Jennifer G. Dy and Andreas Krause (eds.). Proceedings of Machine Learning Research, pp. 5517–5526. PMLR, 2018. PMLR, 2018. </td <td>713</td> <td></td> | 713 | |
| Dy and Andreas Krause (eds.), Proceedings of the 33th International Conference on Machine Learning (ICML), volume 80 of Proceedings of Machine Learning Research, pp. 5517–5526. PMLR, 2018. | Dy and Address Krause (eds.), Proceedings of the 35th International Conference on Machine Learning (ICML), volume 80 of Proceedings of Machine Learning Research, pp. 5517–5526. PMLR, 2018. PMLR, 2018 | 714 | Songbai Yan, Kamalika Chaudhuri, and Tara Javidi. Active learning with logged data. In Jennifer G. |
| The control of the c | The control of the c | 715 | Dy and Andreas Krause (eds.), Proceedings of the 35th International Conference on Machine Learning (ICML), volume 20 of Droceedings of Machine Learning Personnel, pp. 5517, 5526 |
| 111 111 112 121 122 123 124 125 126 127 128 129 130 131 132 133 134 135 136 137 138 139 141 142 143 144 145 145 146 147 148 149 | 11 11 120 121 122 123 124 125 126 127 128 129 120 121 122 123 124 125 126 127 128 129 129 130 131 132 133 134 135 136 137 138 139 131 132 133 134 135 136 137 138 139 131 132 133 134 135 136 137 138 139 131 132 133 | 716 | PMLR 2018 |
| 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 741 742 743 744 745 745 746 747 748 749 | 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 747 748 749 749 741 742 743 744 745 746 747 748 749 749 740 741 742 743 744 745 746 7 | 717 | T MER, 2010. |
| 719 720 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 730 731 732 733 734 735 736 737 738 739 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 | 719 720 721 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 749 740 741 742 743 744 745 746 747 748 749 749 749 749 749 740 741 742 743 744 745 746 747 7 | 718 | |
| 720 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 | 720 721 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 749 740 741 742 743 744 745 746 747 748 749 749 740 741 745 746 747 748 749 749 740 741 745 746 7 | 719 | |
| 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 739 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 | 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 749 740 741 742 743 744 745 746 747 748 749 749 740 741 745 746 747 748 749 750 751 752 753 754 7 | 720 | |
| 722 723 724 725 726 727 728 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 | 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 | 721 | |
| 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 | 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 749 750 751 752 753 | 722 | |
| 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 | 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 | 723 | |
| 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 | 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 | 724 | |
| 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 | 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 | 725 | |
| 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 | 727 728 739 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 | 726 | |
| 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 | 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 | 727 | |
| 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 | 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 | 728 | |
| 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 | 730 731 732 733 734 735 736 737 738 739 740 742 743 744 745 746 747 748 749 750 751 752 753 | 729 | |
| 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 | 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 | 730 | |
| 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 | 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 | 731 | |
| 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 | 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 | 733 | |
| 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 | 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 | 734 | |
| 736 737 738 739 740 741 742 743 744 745 746 747 748 749 | 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 | 735 | |
| 737 738 739 740 741 742 743 744 745 746 747 748 749 | 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 | 736 | |
| 738 739 740 741 742 743 744 745 746 747 748 749 | 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 | 737 | |
| 739 740 741 742 743 744 745 746 747 748 749 | 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 | 738 | |
| 740 741 742 743 744 745 746 747 748 749 | 740 741 742 743 744 745 746 747 748 749 750 751 752 753 | 739 | |
| 741 742 743 744 745 746 747 748 749 | 741 742 743 744 745 746 747 748 749 750 751 752 753 | 740 | |
| 742 743 744 745 746 747 748 749 | 742 743 744 745 746 747 748 749 750 751 752 753 | 741 | |
| 743 744 745 746 747 748 749 | 743 744 745 746 747 748 749 750 751 752 753 | 742 | |
| 744 745 746 747 748 749 | 744 745 746 747 748 749 750 751 752 753 | 743 | |
| 745 746 747 748 749 | 745 746 747 748 749 750 751 752 753 | 744 | |
| 746 747 748 749 | 746 747 748 749 750 751 752 753 | 745 | |
| 747 748 749 | 747 748 749 750 751 752 753 | 746 | |
| 749 | 740 749 750 751 752 753 | 740 | |
| 143 | 749 750 751 752 753 | 748 | |
| 750 | 751 752 753 | 749 | |
| 751 | 752 753 | 750 | |
| 752 | 753 | 752 | |
| 753 | | 753 | |
| | 754 | 754 | |
| 754 | 755 | 755 | |

756 A PROOFS

Proposition 2.1. Eq. 1 can be simplified to $R^{CC}(c \mid S) = -\sum_{\substack{(u,v) \in \mathcal{E} \\ c_u = c_v}} S_{uv}$ + constant, where the constant is independent of different clustering solutions (Chehreghani, 2013).

Proof. As described in (Chehreghani, 2013; 2023), we can write the cost function in Eq. 1 as

 $R^{CC}(\boldsymbol{c} \mid \boldsymbol{S}) = \sum_{(u,v)\in\mathcal{E}} V(u,v \mid \boldsymbol{S}, \boldsymbol{c})$ $= \sum_{\substack{(u,v)\in\mathcal{E}\\c_u=c_v}} \frac{1}{2} (|S_{uv}| - S_{uv}) + \sum_{\substack{(u,v)\in\mathcal{E}\\c_u\neq c_v}} \frac{1}{2} (|S_{uv}| + S_{uv})$ $= \frac{1}{2} \sum_{\substack{(u,v)\in\mathcal{E}\\c_u=c_v}} |S_{uv}| - \frac{1}{2} \sum_{\substack{(u,v)\in\mathcal{E}\\c_u=c_v}} S_{uv} + \frac{1}{2} \sum_{\substack{(u,v)\in\mathcal{E}\\c_u=c_v}} S_{uv} - \frac{1}{2} \sum_{\substack{(u,v)\in\mathcal{E}\\c_u=c_v}} S_{uv}$ $= \underbrace{\frac{1}{2} \sum_{\substack{(u,v)\in\mathcal{E}\\c_u=c_v}} (|S_{uv}| + S_{uv}) - \sum_{\substack{(u,v)\in\mathcal{E}\\c_u=c_v}} S_{uv}.$ (18)

The first term in Eq. 18 is *constant* w.r.t. the choice of a particular clustering c.

Theorem 3.1. Let $\ell : \mathbb{N} \to \mathcal{V}$ denote an object visitation schedule, which satisfies $\lim_{T\to\infty} |\{t \leq T : \ell(t) = u\}| = \infty, \forall u \in \mathcal{V}$. For arbitrary initial conditions, the asynchronous update rules defined by

$$Q_{uk}^{(t+1)} = \exp(-\beta M_{uk}^{(t)}) / \sum_{k' \in \mathbb{K}} \exp(-\beta M_{uk'}^{(t)}),$$
(5)

$$M_{uk}^{(t+1)} = -\sum_{\substack{v \in \mathcal{V} \\ v \neq v}} S_{uv} Q_{vk}^{(t+1)},$$
(6)

where $u = \ell(t)$, converge to a local minimum of Eq. 4.

Proof. Given our cost function R^{MC} (Eq. 2), the *generalized free energy* is defined as (Hofmann et al., 1998)

$$\mathcal{F}_{\beta}(P) \triangleq \mathbb{E}_{P(\mathbf{y})}[R^{\mathrm{MC}}(\mathbf{y})] - \frac{1}{\beta}H(P)$$

= $\sum_{\mathbf{c}\in\mathcal{C}}P(\mathbf{c})R^{\mathrm{MC}}(\mathbf{c}) + \frac{1}{\beta}\sum_{\mathbf{c}\in\mathcal{C}}P(\mathbf{c})\log P(\mathbf{c}),$ (19)

for some $P \in \mathcal{P}$ where \mathcal{P} is the set of distributions with sample space \mathcal{C} . The Gibbs distribution P^{Gibbs} minimizes the generalized free energy (Hofmann et al., 1998) and is called the *free energy*. It can be written as

$$\mathcal{F}_{\beta}(P^{\text{Gibbs}}) = -\frac{1}{\beta}\log \mathcal{Z},\tag{20}$$

where $\mathcal{Z} \triangleq \sum_{c' \in \mathcal{C}} \exp(-\beta R^{MC}(c'))$ is the normalizing constant of the Gibbs distribution in Eq. 3. Given this, we can now simplify the KL-divergence.

 $D_{\mathrm{KL}}(Q \| P^{\mathrm{Gibbs}}) = \sum_{\boldsymbol{c} \in \mathcal{C}} Q(\boldsymbol{c}) \log rac{Q(\boldsymbol{c})}{P^{\mathrm{Gibbs}}(\boldsymbol{c})}$ $= \sum_{c \in \mathcal{C}} Q(c) \log \frac{Q(c)}{\exp\left(-\beta \left(R^{\text{MC}}(c) - \mathcal{F}_{\beta}(P^{\text{Gibbs}})\right)\right)}$

- $= \sum_{\boldsymbol{c} \in \mathcal{C}} Q(\boldsymbol{c}) \left[\log Q(\boldsymbol{c}) + \beta \left(R^{\text{MC}}(\boldsymbol{c}) \mathcal{F}_{\beta}(P^{\text{Gibbs}}) \right) \right]$
- $= \sum_{u \in \mathcal{V}} \sum_{k \in \mathbb{K}} Q_{uk} \log Q_{uk} + \beta \mathbb{E}_{Q(\mathbf{y})}[R^{\mathrm{MC}}(\mathbf{y})] \beta \mathcal{F}_{\beta}(P^{\mathrm{Gibbs}})$
- $= \beta \mathbb{E}_{Q(\mathbf{y})}[R^{\mathrm{MC}}(\mathbf{y})] \sum_{u \in \mathcal{V}} H(\mathbf{y}_u) \beta \mathcal{F}_{\beta}(P^{\mathrm{Gibbs}})$
- $=\beta \mathcal{F}_{\beta}(Q) \beta \mathcal{F}_{\beta}(P^{\text{Gibbs}})$ > 0.

where $H(y_u) \triangleq -\sum_{k \in \mathbb{K}} Q_{uk} \log Q_{uk}$ is the entropy of y_u . The last inequality is a property of the KL-divergence. From this, we have the bound

$$\mathcal{F}_{\beta}(P^{\text{Gibbs}}) \le \mathcal{F}_{\beta}(Q), \tag{22}$$

(21)

and minimizing the KL-divergence corresponds to minimizing the generalized free energy \mathcal{F}_{β} w.r.t. factorial distributions $Q \in Q$, which is consistent with the maximum entropy principle. From this, minimizing the KL-divergence corresponds to the following optimization problem.

$$Q^* = \underset{Q \in \mathcal{Q}}{\operatorname{arg\,min}} \mathcal{F}_{\beta}(Q)$$

s.t.
$$\sum_{k \in \mathbb{K}} Q_{uk} = 1 \quad \forall u \in \mathcal{V}.$$
 (23)

Then, by applying a Lagrangian relaxation to the constraint in Eq. 23 and setting the gradient of the objective w.r.t. Q_{uk} to zero, we obtain

$$0 = \frac{\partial}{\partial Q_{uk}} \left[\mathbb{E}_{Q(\mathbf{y})}[R^{\mathrm{MC}}(\mathbf{y})] - \frac{1}{\beta} \sum_{v \in \mathcal{V}} H(\mathbf{y}_v) + \sum_{w \in \mathcal{V}} \mu_w \left(\sum_{k \in \mathbb{K}} Q_{wk} - 1 \right) \right] \\ = \frac{\partial}{\partial Q_{uk}} \left[\sum_{\boldsymbol{c} \in \mathcal{C}} \prod_{v \in \mathcal{V}} Q_{vc_v} R^{\mathrm{MC}}(\boldsymbol{c}) - \frac{1}{\beta} \sum_{v \in \mathcal{V}} H(\mathbf{y}_v) + \sum_{w \in \mathcal{V}} \mu_w \left(\sum_{k \in \mathbb{K}} Q_{wk} - 1 \right) \right] \\ = \sum_{\boldsymbol{c} \in \mathcal{C}} \prod_{\substack{v \in \mathcal{V} \\ v \neq u}} Q_{vc_v} \mathbf{1}_{\{c_u = k\}} R^{\mathrm{MC}}(\boldsymbol{c}) + \frac{1}{\beta} \left(\log Q_{uk} + 1 \right) + \mu_u \\ = \mathbb{E}_{Q(\mathbf{y}|\mathbf{y}_u = k)} [R^{\mathrm{MC}}(\mathbf{y})] + \frac{1}{\beta} \left(\log Q_{uk} + 1 \right) + \mu_u,$$
(24)

where μ_u 's are the Lagrange multipliers and we define $M_{uk} \triangleq \mathbb{E}_{Q(\mathbf{y}|y_u=k)}[R^{\text{MC}}(\mathbf{y})]$ as the mean-fields, which correspond to the expected cost subject to the constraint that object u is assigned to cluster k. We can simplify

$$M_{uk} = \mathbb{E}_{Q(\mathbf{y}|\mathbf{y}_u=k)}[R^{\mathrm{MC}}(\mathbf{y})]$$

- $= \mathbb{E}_{Q(\mathbf{y}|\mathbf{y}_u=k)} \left[-\sum_{(v,w) \in \mathcal{E}} S_{vw} \right]$
- $= \mathbb{E}_{Q(\mathbf{y}|\mathbf{y}_u=k)} \left[-\sum_{l \in \mathbb{K}} \sum_{(v,w) \in \mathcal{E}} \mathbf{1}_{\{\mathbf{y}_v=l\}} \mathbf{1}_{\{\mathbf{y}_w=l\}} S_{vw} \right]$

$$= -\sum_{l \in \mathbb{K}} \sum_{\substack{(v,w) \in \mathcal{E} \\ v \neq u}} S_{vw} Q_{vl} Q_{wl}$$
$$= -\sum_{l \in \mathbb{K}} \sum_{\substack{v \in \mathcal{V} \\ v \neq u}} S_{uv} Q_{ul} Q_{vl} - \sum_{l \in \mathbb{K}} \sum_{\substack{(v,w) \in \mathcal{E} \\ v \neq u}} S_{vw} Q_{vl} Q_{wl}$$

$$l \in \mathbb{K}$$
 $\substack{v \in \mathcal{V} \\ v \neq u}$

$$= -\sum_{\substack{v \in \mathcal{V} \\ v \neq u}} S_{uv}Q_{vk} - \sum_{\substack{l \in \mathbb{K} \\ v \neq u}} \sum_{\substack{v \neq u \\ v \neq u \\ w \neq u}} S_{vw}Q_{vl}Q_{wl}$$
883
886
886
886

where the last equality uses that $Q_{ul} = 1$ if l = k and 0 otherwise, according to $Q(c \mid c_u = k)$. The second term of the last expression is a constant w.r.t. Q_{uk} and is thus irrelevant for optimization (since it does not depend on u).

 $0 = M_{uk} + \frac{1}{\beta} \left(\log Q_{uk} + 1 \right) + \mu_u.$

 $= -\sum_{l \in \mathbb{K}} \sum_{(v,w) \in \mathcal{E}} \mathbb{E}_{Q(\mathbf{y}|\mathbf{y}_u=k)} [\mathbf{1}_{\{\mathbf{y}_v=l\}} \mathbf{1}_{\{\mathbf{y}_w=l\}}] S_{vw}$

With the definition of M_{uk} , we can rewrite Eq. 24 as

 Then, we have

$$\log Q_{uk} = -\beta M_{uk} - \beta \mu_u$$

$$\Rightarrow Q_{uk} = \exp\left(-\beta M_{uk}\right) \exp\left(-\beta \mu_u\right).$$
(27)

(25)

(26)

On the other hand, we have: $\sum_{k'} Q_{uk'} = 1$. Therefore,

$$\sum_{k'} \log Q_{uk'} = \sum_{k'} \exp\left(-\beta M_{uk'}\right) \exp\left(-\beta \mu_u\right) = 1$$

$$\Rightarrow \exp\left(-\beta \mu_u\right) = \frac{1}{\sum_{k'} \exp\left(-\beta M_{uk'}\right)}.$$
(28)

Then, inserting Eq. 28 into Eq. 27 yields

-

$$Q_{uk} = \frac{\exp\left(-\beta M_{uk}\right)}{\sum_{k'} \exp\left(-\beta M_{uk'}\right)}.$$
(29)

This derivation suggest an EM-type procedure for minimizing the KL-divergence $D_{\rm KL}(Q||P^{\rm Gibbs})$, which consists of alternating between estimating Q_{uk} 's given M_{uk} 's and then updating M_{uk} 's given the new values of Q_{uk} 's (as described in Alg. 2).

Finally, we can compute the Hessian of the objective as

The positivity of the Hessian in Eq. 30 ensures that the generalized free energy $\mathcal{F}_{\beta}(Q)$ is convex with respect to Q_{uk} for each object u, guaranteeing that the update for Q_{uk} strictly decreases $\mathcal{F}_{\beta}(Q)$ unless it is already at a local minimum. Since $\mathcal{F}_{\beta}(Q)$ is bounded from below by $\mathcal{F}_{\beta}(P^{\text{Gibbs}})$ and each object u is updated infinitely often according to the object visitation schedule, the algorithm converges to a local minimum of the generalized free energy \mathcal{F}_{β} within the space of factorial distributions \mathcal{Q} .

 $\frac{\partial^2}{\partial Q_{uk}^2} \mathcal{F}_{\beta}(Q) = \frac{\partial}{\partial Q_{uk}} M_{uk} + \frac{1}{\beta} \left(\log Q_{uk} + 1 \right) + \mu_u$ $= \frac{1}{\beta Q_{uk}}$

(30)

В ADDITIONAL DETAILS ABOUT INFORMATION-THEORETIC ACQUISITION **FUNCTIONS**

B.1 DETAILED DERIVATION OF ENTROPY

Here we show a detailed derivation of the probability $P(E_{uv})$, which is used for the acquisition function based on entropy in Eq. 8. We have

$$P(E_{uv} = 1) = \mathbb{E}_{P^{\text{Gibbs}}(\mathbf{y})}[\mathbf{1}_{\{y_u = y_v\}}]$$

$$\approx \mathbb{E}_{Q(\mathbf{y})}[\mathbf{1}_{\{y_u = y_v\}}]$$

$$= \sum_{k' \in \mathbb{K}} Q_{uk'} \sum_{k'' \in \mathbb{K}} Q_{uk''} \mathbf{1}_{\{c_u = c_v\}}$$

$$= \sum_{k \in \mathbb{K}} Q_{uk} Q_{vk} + \underbrace{\sum_{k' \in \mathbb{K}} \sum_{\substack{k'' \in \mathbb{K} \\ k'' \neq k'}} Q_{uk'} Q_{vk''} \mathbf{1}_{\{c_u = c_v\}}}_{=0}$$

$$= \sum_{k \in \mathbb{K}} Q_{uk} Q_{vk}.$$
(31)

One can also show that $P(\mathbf{E}_{uv} = -1) \approx \mathbb{E}_{Q(c)}[\mathbf{1}_{\{c_u \neq c_v\}}(c)]$ which can be simplified to $P(\mathbf{E}_{uv} = -1)$ $-1) = \sum_{k,k' \in \mathbb{K}} Q_{uk} Q_{vk'} \mathbf{1}_{\{k \neq k'\}} = 1 - P(\mathbf{E}_{uv} = 1).$

B.2 DERIVATIONS OF EIG

In this section, we include a detailed derivation of the acquisition function defined in Eq. 13. In addition, we derive the acquisition function $a^{\text{EIG-P}}$.

$$H(\mathbf{y}) = -\sum_{\boldsymbol{c}\in\mathcal{C}} P^{\text{Gibbs}}(\mathbf{y}=\boldsymbol{c}) \log P^{\text{Gibbs}}(\mathbf{y}=\boldsymbol{c}).$$
(32)

A mean-field approximation $Q(\mathbf{y} = \mathbf{c}) = \prod_{u=1}^{N} Q(\mathbf{y}_u = c_u)$ of P^{Gibbs} assumes independence between each cluster label \mathbf{y}_u . This means we have

$$\begin{aligned} H(\mathbf{y}) &= -\sum_{c \in \mathcal{C}} Q(\mathbf{y} = c) \log Q(\mathbf{y} = c) \\ &= -\sum_{c \in \mathcal{C}} \prod_{u=1}^{N} Q(\mathbf{y}_{u} = c_{u}) \log \prod_{v=1}^{N} Q(\mathbf{y}_{v} = c_{v}) \\ &= -\sum_{c \in \mathcal{C}} \prod_{u=1}^{N} Q(\mathbf{y}_{u} = c_{u}) (\sum_{v=1}^{N} \log Q(\mathbf{y}_{v} = c_{v})) \\ &= -\sum_{c \in \mathcal{C}} \sum_{v=1}^{N} \prod_{u=1}^{N} Q(\mathbf{y}_{u} = c_{u}) \log Q(\mathbf{y}_{v} = c_{v}) \\ &= -\sum_{v=1}^{N} \sum_{k \in \mathbb{K}} \sum_{\substack{c \in \mathcal{C} \\ c_{v} = k}} \prod_{u=1}^{N} Q(\mathbf{y}_{u} = c_{u}) \log Q(\mathbf{y}_{v} = c_{v}) \\ &= -\sum_{v=1}^{N} \sum_{k \in \mathbb{K}} \log Q(\mathbf{y}_{v} = c_{v}) \sum_{\substack{c \in \mathcal{C} \\ c_{v} = k}} \prod_{u=1}^{N} Q(\mathbf{y}_{u} = c_{u}) \\ &= -\sum_{v=1}^{N} \sum_{k \in \mathbb{K}} \log Q(\mathbf{y}_{v} = c_{v}) \sum_{\substack{c \in \mathcal{C} \\ c_{v} = k}} \prod_{u=1}^{N} Q(\mathbf{y}_{u} = c_{u}) \\ &= -\sum_{v=1}^{N} \sum_{k \in \mathbb{K}} \log Q(\mathbf{y}_{v} = c_{v}) Q(\mathbf{y}_{v} = c_{v}) \underbrace{\left(\sum_{\substack{c \in \mathcal{C} \\ c_{v} = k}} \prod_{u \neq v}^{N} Q(\mathbf{y}_{u} = c_{u})\right)}_{=1} \\ &= -\sum_{v=1}^{N} \sum_{k \in \mathbb{K}} Q(\mathbf{y}_{v} = c_{v}) \log Q(\mathbf{y}_{v} = c_{v}) \\ &= -\sum_{v=1}^{N} \sum_{k \in \mathbb{K}} Q(\mathbf{y}_{v} = c_{v}) \log Q(\mathbf{y}_{v} = c_{v}) \\ &= -\sum_{v=1}^{N} \sum_{k \in \mathbb{K}} Q(\mathbf{y}_{v} = c_{v}) \log Q(\mathbf{y}_{v} = c_{v}) \\ &= -\sum_{v=1}^{N} H(\mathbf{y}_{v} \mid \mathbf{Q}). \end{aligned}$$

Furthermore, we assume independence between pairs **E** given a mean-field approximation **Q**. Consequently, we have $P(\mathbf{E}) = \prod_{(w,l) \in \mathcal{E}} P(\mathbf{E}_{wl} \mid \mathbf{Q})$. Then, one can derive $H(\mathbf{E}) = \sum_{(w,l) \in \mathcal{E}} H(\mathbf{E}_{wl} \mid \mathbf{Q})$ following the same derivation as shown in Eq. 33. In addition, we choose to approximate $H(\mathbf{y} \mid \mathbf{E}_{uv} = e)$ and $H(\mathbf{E} \mid \mathbf{E}_{uv} = e)$ using conditional mean-field approximation $\mathbf{Q}^{(S_{uv}=e)}$. As a consequence, the joint conditional entropies reduces to the sum of entropies over individual variables. This is shown following the same derivation shown in Eq. 33.

From this, we obtain the acquisition function $a^{\text{EIG-O}}$ defined in Eq. 13 of the main paper. In addition, we define an acquisition function which computes the *expected information gain over the clustering relation of pairs* (EIG-P).

1014 1015 1016

1017

$$a^{\mathrm{EIG-P}}(u,v) \triangleq \sum_{(w,l)\in\mathcal{E}} H(\mathsf{E}_{wl} \mid \boldsymbol{Q}) - \sum_{e \in \{-1,+1\}} P(\mathsf{E}_{uv} = e \mid \boldsymbol{Q}) H(\mathsf{E}_{wl} \mid \boldsymbol{Q}^{(S_{uv} = e)}).$$
(34)

 $a^{\text{EIG-P}}$ is also computed using Alg. 3.

1018 B.3 JEIG ALGORITHM

1020 Alg. 4 outlines how the acquisition function a^{JEIG} (Eq. 17) is calculated. The algorithm begins by 1021 initializing Q and M by running Algorithm 2. Then, the algorithm loops m times. In each of the m1022 iterations, a subset $\mathcal{D}_i \subseteq \mathcal{E}$ is selected. Then, lines 7-10 computes a Monte-Carlo estimation of the 1023 expectation $\mathbb{E}_{e \sim P(\mathbb{E}_{\mathcal{D}_i})}[H(\mathbb{E}_{uv} \mid \mathbb{E}_{\mathcal{D}_i} = e)]$. We found that the selection of \mathcal{D}_i (on line 6) can be 1024 done in a number of ways, with good performance. In the experiments of this paper, we select the 1025 top- $|\mathcal{D}_i|$ pairs according to $\log(a^{\text{Entropy}}(u, v)) + \epsilon_{uv}$ where $\epsilon_{uv} \sim \text{Gumbel}(0; 1)$. In other words, the 1026 top- $|\mathcal{D}_i|$ pairs according to a^{Entropy} with some added acquisition noise (as explained in Section 4.1). This leads to diversity among the selected \mathcal{D}_i , while containing pairs with large entropy. Pairs with large entropy are likely to have large impact on each E_{uv} , and are therefore important to include in \mathcal{D}_i . We set $|\mathcal{D}_i| = 0.02|\mathcal{E}|$ (i.e., 2% of all pairs), m = 5 and n = 50 for all datasets. See Appendix C.7 for more details about this.

1031 Algorithm 4 JEIG

1030

1046 1047

1048 1049

1050

1051

1032 1: Input: Similarity matrix S, current clustering c^i , concentration parameter β . 1033 2: $M_{uk} \leftarrow -\sum_{v:c_u^i=k} S_{uv}, \forall u \in \mathcal{V}, \forall k \in \mathbb{K}$ 1034 3: $Q, M \leftarrow \text{MeanField}(S, M, \beta)$ 4: $a^{\mathbf{J} \dot{\mathbf{E}} \mathbf{I} \mathbf{G}}(u, v) \leftarrow 0 \quad \forall (u, v) \in \mathcal{E}$ 1035 5: for $i \leftarrow 1$ to m do 1036 $\mathcal{D}_i \leftarrow \text{SelectPairs}(\mathcal{E})$ $\triangleright \mathcal{D}_i \subseteq \mathcal{E}$ 6: for $j \leftarrow 1$ to n do 7: $\boldsymbol{e} \sim P(\mathbf{E}_{\mathcal{D}_i})$ 8: 1039 $Q^{(S_{\mathcal{D}_i}=e)} \leftarrow \text{MeanField}(S, M, \beta \mid S_{\mathcal{D}_i}=e)$ 9: 1040 $a^{\text{JEIG}}(u,v) \leftarrow a^{\text{JEIG}}(u,v) + H(\mathbf{E}_{uv} \mid \boldsymbol{Q}^{(\boldsymbol{S}_{\mathcal{D}_i}=\boldsymbol{e})})/n \quad \forall (u,v) \in \mathcal{E}$ 10: 1041 11: end for 12: end for 1043 13: $a^{\text{JEIG}}(u, v) \leftarrow H(\mathbf{E}_{uv} \mid \boldsymbol{Q}) - a^{\text{JEIG}}(u, v)/m \quad \forall (u, v) \in \mathcal{E}$ 14: return a^{JEIG} 1044 1045

C EXPERIMENTS: MORE DETAILS AND FURTHER RESULTS

In this section, we describe the datasets in more detail and provide further experimental results. The experimental settings are identical to Section 4, unless otherwise specified.

1052 1053 C.1 MAXMIN AND MAXEXP

In this section, we explain the acquisition functions maxmin and maxexp introduced in (Aronsson & Chehreghani, 2024). First, the transitive property implies if $S_{uv} \ge 0$ and $S_{uw} \ge 0$ then $S_{vw} \ge 0$ or if $S_{uv} \ge 0$ and $S_{uw} < 0$ then $S_{vw} < 0$. Then, assuming the ground-truth similarity matrix S^* is consistent (i.e., it does not violate transitive property) would imply that explicitly resolving (or preventing) the inconsistency in S may be informative. Both maxmin and maxexp are based on this idea.

Let \mathcal{T} be the set of triples (u, v, w) of distinct objects in \mathcal{V} , i.e., $|\mathcal{T}| = \binom{N}{3}$. Let $\mathcal{T}_{uv} = \{t \in \mathcal{T} \mid u, v \in t\}$ be the set of triples that include the pair (u, v). Let \mathcal{C}_t be the set of clustering solutions for the objects in the triple t. Finally, let $\mathcal{E}_t = \{(u, v) \in \mathcal{E} \mid u, v \in t\}$ be the set of pairs in the triple t, and $e_t = \arg \min_{(u,v) \in \mathcal{E}_t} |S_{uv}|$ is the pair in \mathcal{E}_t with the smallest absolute similarity. Given this, maxmin is defined as³

$$a^{\text{Maxmin}}(u,v) \triangleq \max_{t \in \mathcal{T}_{uv}} \min_{\boldsymbol{c} \in \mathcal{C}_t} R(\boldsymbol{c} \mid \mathcal{E}_t) \mathbf{1}_{\{e_t = (u,v)\}},\tag{35}$$

where $R(\mathbf{c} \mid \mathcal{E}_t) \triangleq \sum_{(u,v) \in \mathcal{E}_t} V(u,v \mid \mathbf{c})$. Intuitively, maxmin begins by ranking each triple 1068 according to how much inconsistency they induce (i.e., violation of transitive property). Then, from 1069 each of the top-B triples t, the pair in \mathcal{E}_t with smallest absolute similarity is selected (i.e., most 1070 uncertain according to its similarity). The goal is thus to reduce inconsistency by resolving violations 1071 of the transitive property in triples. In our experiments, we observe that this can be ineffective, likely 1072 due to robustness to inconsistency in S by the CC algorithm used. See discussion in experiments for more details. From Eq. 35 we see that *maxmin* quantifies the inconsistency by the cost of the best 1074 clustering in C_t (in short, this cost is non-zero for triples that violate the transitive property, and zero 1075 otherwise). The maximization over \mathcal{T}_{uv} ensures the most violating triple that includes the pair (u, v)is considered. maxexp works analogously to maxmin except the term $\min_{c \in C_t} R(c \mid \mathcal{E}_t)$ is replaced 1077 by an expectation of the cost of all clustering solutions in C_t .

¹⁰⁷⁸ 1079

³This formulation of *maxmin* is equivalent to Algorithm 3 of (Aronsson & Chehreghani, 2024), except the algorithm overcomes the computational issues of iterating all $\binom{3}{3}$ triples.

C.2 DETAILS ABOUT ORACLE 4

1081

In this section, we describe the details of Oracle 4. Given a dataset X and ground-truth labels c^* , the ground-truth similarities are defined as $S_{uv}^* = +1$ if $c_u^* = c_v^*$, and -1 otherwise. The dataset X is then split into two disjoint parts: $X = X_{\text{train}} \cup X_{\text{test}}$, with 30% of the data allocated to the training set and 70% to the test set. Given this, the sizes of X_{train} and X_{test} are restricted to a maximum of 5000 and 1000 samples, respectively.

1087 Next, we define a pairwise prediction model $f_{\theta} : \mathbf{X} \times \mathbf{X} \rightarrow [-1, +1]$. In our experiments, f_{θ} is a 1088 fully connected neural network with 6 hidden layers of sizes [1024, 2048, 512, 248, 64], using ReLU 1089 activations. The input to the network is the concatenation of two feature vectors, i.e., $\mathbf{x}_u \oplus \mathbf{x}_v$. We 1090 treat this as a binary classification problem, where the output of the neural network represents the 1091 probability that the similarity between u and v is +1. Denoting this probability as p_{uv} , we transform 1092 it to a similarity score in [-1, +1] using the transformation $2 \cdot p_{uv} - 1$. The network is trained using 1093 the standard cross-entropy loss function over 30 epochs.

We then construct a training dataset where the inputs are $\{\mathbf{x}_u \oplus \mathbf{x}_v\}_{\mathbf{x}_u, \mathbf{x}_v \in \mathbf{X}_{\text{train}}}$, and the corresponding labels are S_{uv}^* . In practice, we limit the number of training pairs to a maximum of 30,000, as the total number of possible pairs would otherwise be prohibitively large, resulting in extremely slow training.

Finally, the active correlation clustering experiments are conducted on the data points in \mathbf{X}_{test} . It is important to note that the ground-truth similarities of pairs in \mathbf{X}_{test} are not used during the training of f_{θ} .

1100 1101

1113

1114

1115

1116

1117

1118

1119

1120 1121

1122

1123

1124

1102 C.3 DESCRIPTION OF DATASETS

A detailed description of all eight datasets used is provided below. Datasets 2-6 are taken from the UCI machine learning repository (Kelly et al., 2023) (all of which are released under the CC BY 4.0 license).

- 1107 1108 1. **CIFAR10** (Krizhevsky, 2009). This dataset consists of 60000 32×32 color images in 10 different classes (with 6000 images per class). A random subset of N = 1000 images (with $|\mathcal{E}| = 499, 500$) is used.⁴ Cluster sizes: [91, 96, 107, 89, 99, 113, 96, 93, 112, 104]. We use a ResNet18 model (He et al., 2015) trained on the full CIFAR10 dataset in order to embed the 1000 images into a 512-dimensional space. For oracle 4, f_{θ} is trained on data points embedded into the latent space. We set $|\mathcal{E}^0| = 2500$. The batch size is set to B = 1250.
 - 2. **20newsgroups**. This dataset consists of 18846 newsgroups posts (in the form of text) on 20 topics (clusters). We consider a subset of 5 topics: "rec.sport.baseball", "soc.religion.christian", "rec.autos", "talk.politics.mideast", "misc.forsale". We use a random sample of N = 1000 posts (with $|\mathcal{E}| = 499, 500$). Cluster sizes: [201, 190, 201, 217, 191]. We use the distilbert-base-uncased transformer model loaded from the Flair Python library (Akbik et al., 2018) in order to embed each of the 1000 documents (data points) into a 768-dimensional latent space. For oracle 4, f_{θ} is trained on data points embedded into the latent space. We set $|\mathcal{E}^0| = 2500$. The batch size is set to B = 250.
 - 3. **Cardiotocography**. This dataset includes 2126 fetal cardiotocograms consisting of 22 features and 10 classes. We use a sample of N = 1000 data points (with $|\mathcal{E}| = 499, 500$). Cluster sizes: [180, 275, 27, 35, 31, 148, 114, 62, 28, 100]. We set $|\mathcal{E}^0| = 2500$. The batch size is set to B = 750.
- 4. Ecoli. This is a biological dataset on the cellular localization sites of 8 types (clusters) of proteins which includes N = 336 samples (with $|\mathcal{E}| = 56, 280$). Cluster sizes: [137, 76, 1, 2, 37, 26, 5, 52]. We set $|\mathcal{E}^0| = 280$. The batch size is set to B = 85.
- **5. Forest Type Mapping.** This is a remote sensing dataset of N = 523 samples collected from forests in Japan and grouped in 4 different forest types (clusters) (with $|\mathcal{E}| = 136, 503$). Cluster sizes: [168, 84, 86, 185]. We set $|\mathcal{E}^0| = 500$. The batch size is set to B = 350.
- 1132 1133

⁴For oracles 1-3 we simply select N = 1000 random data points. For oracle 4 we obtain the random sample based on the construction of \mathbf{X}_{test} , as explained in the previous section. The same applies to the other datasets.



- 8. Synthetic. This is a synthetically generated dataset (normally distributed 10-dimensional data points) with N = 500 (and $|\mathcal{E}| = 124,750$) data points split evenly into 10 clusters. We set $|\mathcal{E}^0| = 500$. The batch size is set to B = 300.
- 1166 C.4 FURTHER RESULTS

1165

1174

1182

1184

Figures 4 and 5 show results for oracles 2 and 3, respectively, where the evaluation metric is the adjusted rand index (ARI). Figures 6-9 show results for all oracles, where the evaluation metric is the adjusted mutual information (AMI). All results are consistent with the insights from Figures 1-2 from the main paper, where all information-theoretic acquisition functions proposed in this paper outperform the baselines. In addition, we observe that the acquisition functions based on information gain consistently outperforms a^{Entropy} .

1175 C.5 SMALL BATCH SIZE

1176 1177 In Figure 10, we show results for oracle 1 for a synthetic dataset with N = 70 objects (and $|\mathcal{E}| = 2415$ pairs) using a batch size of B = 5. The noise level is $\gamma = 0.4$. In this experiment, we do not use any acquisition noise in order to improve batch diversity. The purpose of this experiment is to further illustrate the benefit of the acquisition functions based on information gain compared to entropy, when differences due to batch diversity are (mostly) removed. We observe that $a^{\text{EIG-O}}$, $a^{\text{EIG-P}}$ and a^{JEIG} outperform a^{Entropy} .

1183 C.6 RUNTIME

Each active learning procedure was executed on 1 core of an Intel(R) Xeon(R) Gold 6338 CPU @
2GHz (with 32 cores total). We have access to a compute cluster with many of these CPU's allowing
us to execute many procedures in parallell. Each CPU has access to 128GB of RAM (shared among cores), but much less would suffice for our experiments.



Figure 5: Results for oracle 3 with noise level $\gamma = 0.2$. The evaluation metric is the adjusted rand index (ARI).



Figure 6: Results for oracle 1 with noise level $\gamma = 0.4$. The evaluation metric is the adjusted mutual information (AMI).

In Figure 11 we show the runtime of each iteration in seconds for all acquisition functions and datasets. We observe that a^{Entropy} is very efficient (comparable to other baseline methods). In addition, we see that out of the acquisition functions based on information gain, a^{JEIG} is the most efficient and is quite close to a^{Entropy} . Expectedly, $a^{\text{EIG-O}}$ and $a^{\text{EIG-P}}$ are the least efficient. This is because we run Alg. 2 numerous times (as discussed in Section B.2). Out of these two, $a^{\text{EIG-P}}$ is the most inefficient since it involves a sum over all $\binom{N}{2}$ pairs (see Eq. 34) in each iteration of Alg. 3.

1232 1233

1234 1235

1236

1223

1224 1225

1204

1205

C.7 Hyperparameters

In this section, we present a detailed analysis of all hyperparameters. All experiments use oracle 1.

¹²³⁷ C.7.1 EIG

In Figure 12, we show results for the acquisition functions $a^{\text{EIG-O}}$ (left) and $a^{\text{EIG-P}}$ (right) with different values of $|\mathcal{E}^{\text{EIG}}|$ (using oracle 1). See Alg. 3 for the usage of \mathcal{E}^{EIG} . We observe that the performance does not improve much beyond $|\mathcal{E}^{\text{EIG}}| = 10N$. This indicates both of these acquisition functions will perform well when evaluation Eq. 13 or Eq. 34 for O(N) of the pairs (instead of all $O(N^2)$ pairs).



Figure 7: Results for oracle 2 with variance $\gamma = 1.3$. The evaluation metric is the adjusted mutual information (AMI).



Figure 8: Results for oracle 3 with noise level $\gamma = 0.2$. The evaluation metric is the adjusted mutual information (AMI).

1280 C.7.2 JEIG 1281

In Figure 13, we show results for the acquisition function a^{JEIG} when varying m, n and $|\mathcal{D}_i|$. The left 1282 plot show results when varying $|\mathcal{D}_i|$ with m and n fixed to 5 and 50, respectively. As explained in 1283 Appendix B.3, each \mathcal{D}_i is selected as the top- $|\mathcal{D}_i|$ pairs according to $\log(\bar{a}^{\text{Entropy}}(u, v)) + \hat{\epsilon}_{uv}$ where 1284 $\epsilon_{uv} \sim \text{Gumbel}(0; 1)$. The right plot show results when varying m and n with $|\mathcal{D}_i|$ fixed to $0.02|\mathcal{E}|$ 1285 (2% of all pairs). We observe that $|\mathcal{D}_i| = 0.02 |\mathcal{E}|$ performs the best. A smaller value means we do 1286 not capture enough information about each E_{uv} and a too large value leads to exaggerated selection 1287 bias, as explained at the end of Section 3.3.3. We find $|\mathcal{D}_i| = 0.02 |\mathcal{E}|$ to work well for all datasets 1288 considered in this paper. However, there may be other values that perform equally well (or better). 1289 Finally, we observe that larger values of m and n expectedly lead to better performance. A larger 1290 value of m means we capture more information about each E_{uv} . A larger value of n means the 1291 Monte-Carlo estimation of the expectation $\mathbb{E}_{e \sim P(\mathbf{E}_{\mathcal{D}_i})}[H(\mathbf{E}_{uv} \mid \mathbf{E}_{\mathcal{D}_i} = e)]$ becomes more accurate.

1292

1277

1278 1279

1258

1259

1293 C.7.3 CONCENTRATION PARAMETER β 1294

In Figure 14, we show results for the information-theoretic acquisition functions a^{Entropy} , $a^{\text{EIG-O}}$ and a^{JEIG} when varying the hyperparameter β . The parameter β is a concentration parameter used in



Figure 10: Results on synthetic dataset with N = 70 and $|\mathcal{E}| = 2415$ using a small batch size B = 5without any acquisition noise. The noise level is $\gamma = 0.4$. This experiment used oracle 1.

(a)

1330 1331

1326 1327

the definition of the Gibbs distribution from Eq. 3. As a consequence, it is also used in Alg. 2 (mean-field), which is frequently used in this paper. In this setting, β is the well-known inverse temperature of a Gibbs distribution (having this parameter with a Gibbs distribution is very common). A large β will concentrate more probability mass on a cluster k with larger cost M_{uk}^t . A smaller β will make the probabilities Q_{uk}^t more uniform across different clusters. β may therefore have an impact on the resulting clustering (i.e., assignment probabilities Q which is used by all information-theoretic acquisition functions). See (Chehreghani et al., 2012) for more details about the impact of β . We observe that a value of $\beta = 3$ performs the best for all acquisition functions.

1339

1341

1340 C.8 UTILIZING FEATURES FOR ACTIVE CORRELATION CLUSTERING

The primary focus of this paper is on standard correlation clustering (Bansal et al., 2004; Bonchi et al., 2014) in the active learning setting, where no feature vectors are assumed to be available, i.e., similar to the recent works (Bressan et al., 2019; García-Soriano et al., 2020; Aronsson & Chehreghani, 2024; Kuroki et al., 2024). However, our framework is generic enough to also incorporate feature vectors when available for the objects. In this section, we propose an innovative but simple method. After line 6 of Alg. 1, we introduce a prediction component that predicts similarities based on the queries made so far. This component works as follows.

For all queried pairs (u, v) (i.e., pairs where the oracle has provided the similarity S_{uv}), we concatenate their feature vectors $\mathbf{x}_u \oplus \mathbf{x}_v$ and add them to a dataset, using S_{uv} as the corresponding label.



Figure 11: Runtime of all acquisition functions on all datasets with noise level $\gamma = 0.4$. The y-axis corresponds to the execution time in seconds of each iteration. This corresponds to the same experiments presented in Figure 1.



Figure 12: Results for acquisition functions $a^{\text{EIG-O}}$ (left) and $a^{\text{EIG-P}}$ (right) with different values of $|\mathcal{E}^{\text{EIG}}|$.



Figure 13: Results for acquisition function a^{JEIG} when varying hyperparameters m, n and $|\mathcal{D}_i|$.



Figure 14: Results of information-theoretic acquisition functions a^{Entropy} , $a^{\text{EIG-O}}$ and a^{JEIG} when varying hyperparameter β (used in Eq. 3)



Figure 15: Performance of a^{Entropy} and a^{Uniform} when leveraging feature vectors across two synthetic datasets, one with a simpler (easy) feature space and the other with a more complex (hard) feature space.



Figure 16: Performance of a^{Entropy} and a^{Uniform} when leveraging feature vectors for two real-world datasets.

Then, we train a pairwise similarity prediction model $f_{\theta} : \mathbf{X} \times \mathbf{X} \rightarrow [-1, +1]$ using the collected dataset. The training process follows a procedure similar to Oracle 4 (see Appendix C.2 for details). Finally, use f_{θ} to predict similarities for the pairs that have not yet been queried. To minimize noise or incorrect predictions, it may be advantageous to limit predictions to pairs for which f_{θ} exhibits high confidence.

We conducted some experiments to evaluate this approach, with the results presented in Figures 15-16. The experiments in Figure 15 were performed on two synthetic datasets. The first dataset has relatively simple structure, with well-separated clusters, while the second dataset is more challenging, exhibiting overlap between clusters. The results suggest that the prediction component facilitates rapid convergence, reducing the number of queries required. However, we observe that leveraging

features can lead to convergence toward arbitrarily poor solutions if the dataset contains bias. This
issue is particularly pronounced in the more challenging dataset, where noise and bias are more
prevalent. Figure 16 presents the results for two real-world datasets, showcasing a similar trend.
Notably, in the forest type mapping dataset, the ground-truth clustering is identified rapidly when
feature vectors are used. In contrast, for the cardiotocography dataset, the procedure converges to a
suboptimal clustering when using features, likely due to noise or bias in the dataset's feature space.

It is important to emphasize that the primary focus of this paper is the setting where high-quality feature vectors are not assumed to be available. Therefore, the proposed prediction component should be viewed as an intriguing direction for future work in this context. However, since we do not assume access to feature vectors, it would not be appropriate to center this study on the predictive component (or any other way of incorporating feature vectors).

- 1469
- 1470 1471

D MAX CORRELATION CLUSTERING ALGORITHM

In this section, we describe the CC algorithm used in the active CC procedure outlined in Section 2.2. The algorithm was derived in (Aronsson & Chehreghani, 2024) based on the max correlation cost function $R^{\text{MC}}(\boldsymbol{c} \mid \boldsymbol{S}) \triangleq -\sum_{\substack{(u,v) \in \mathcal{E} \\ c_u = c_v}} S_{uv}$ introduced in Proposition 2.1. It is highly robust to noise in \boldsymbol{S} and dynamically determines the number of clusters.

1477 The method is based on local search and is outlined in Alg. 5. It takes as input a set of objects \mathcal{V} , a similarity matrix S, an initial number of clusters K, the number of repetitions T, and a stopping 1478 threshold η . In our experiments, we set T = 5, $\eta = 2^{-52}$ (double precision machine epsilon) and 1479 $K = |\mathcal{V}|$ in the first iteration of the active CC procedure, and then $K = K^i$ for all remaining iterations 1480 where K^i denotes the number of clusters in the current clustering c^i . The output is a clustering 1481 $c \in \mathcal{C}$. The main part of the algorithm (lines 4-22) is based on the local search of the respective 1482 non-convex objective. Therefore, we run the algorithm T times with different random initializations 1483 and return the best clustering in terms of the objective function. The main algorithm (starting from 1484 line 4) consists of initializing the current clustering c randomly. Then, it loops for as long as the 1485 current max correlation objective changes by at least η compared to the last iteration. If not, we 1486 assume it has converged to some (local) optimum. Each repetition consists of iterating over all the 1487 objects in \mathcal{V} in a random order \mathcal{V}_{rand} (this ensures variability between the T runs). For each object 1488 $u \in \mathcal{V}_{rand}$, it calculates the similarity (correlation) between u and all clusters $k \in \{1, \ldots, K\}$, which 1489 is denoted by $S_k(u)$. Then, the cluster k_{\max} that is most similar to u is obtained. Now, if the most similar cluster to u has a negative correlation score, this indicates that u is not sufficiently similar 1490 to any of the existing clusters. Thus, we construct a new cluster with u as the only member. If the 1491 most similar cluster to u is positive, we simply assign u to this cluster. Consequently, the number of 1492 clusters will dynamically change based on the pairwise similarities (it is possible that the only object 1493 of a singleton cluster is assigned to another cluster and thus the singleton cluster disappears). Finally, 1494 in each repetition the current max correlation objective is computed efficiently by only updating it 1495 based on the current change of the clustering c (i.e., lines 14 and 20). The computational complexity 1496 of the procedure is $O(KN^2)$. See (Aronsson & Chehreghani, 2024) for additional details about how 1497 the algorithm was derived. 1498

1490

E RELATION TO MULTI-ARMED BANDIT METHODS

The studies in (Gullo et al., 2023; Kuroki et al., 2024) address query-efficient correlation clustering (CC) by framing it as a *multi-armed bandit* (MAB) problem. Below, we compare these approaches to our own methods.

First, (Gullo et al., 2023) distinguish their approach from the query-efficient CC framework studied by (Bressan et al., 2019; García-Soriano et al., 2020), which is the setting we also consider. The key difference lies in the assumption regarding the query budget *B*. While we assume a fixed budget with $B \ll |\mathcal{E}|$, (Gullo et al., 2023) does not impose this constraint, allowing the number of queries to exceed the total number of pairs.

In (Kuroki et al., 2024), the authors propose learning edge weights (pairwise similarities) using combinatorial bandit algorithms. Similar to other query-efficient CC studies (Bressan et al., 2019; García-Soriano et al., 2020), they utilize KwikCluster as their base clustering algorithm. KwikCluster,

| 2 Alg | porithm 5 Max Correlation Clustering Algorithm \mathcal{A} (dynamic K) |
|--------------------|---|
| 3 | Input : \mathcal{V} , S , initial number of clusters K , number of iterations T , stopping threshold η |
| + | Output : Clustering solution $c \in C$ |
| 1: | $N \leftarrow \mathcal{V} $ |
| 2: | $R_{\text{best}}^{\text{MC}} \leftarrow -\infty$ |
| 3: | for $j \in \{1, \dots, T\}$ do |
| 4: | $c \leftarrow$ random clustering in C with K clusters |
| 5: | $R^{\text{MC}} \leftarrow R^{\text{MC}}(\boldsymbol{c} \mid \boldsymbol{S})$ |
| 6: | $R_{\text{old}}^{\text{MC}} \leftarrow R^{\text{MC}} - 1$ |
| 7: | while $ R^{MC} - R^{MC}_{old} > \eta$ do |
| 8: | $R_{old}^{\rm MC} \leftarrow R^{\rm MC}$ |
| 9: 10. | $V_{\text{rand}} \leftarrow a \text{ random permutation of the objects in } V$ |
| 10: | So $(u) \leftarrow \sum S \forall k \in \{1, \dots, K\}$ |
| 11. | $S_k(u) \land \underline{\ }_{v:c_v=k} S_{uv}, \land n \in \{1, \dots, n\}$ |
| 12: | $\kappa_{\max} \leftarrow \arg\max_{k \in \{1, \dots, K\}} S_k(u)$ |
| 13: | $B_{k_{\max}}(u) < 0$ then $D^{MC} = D^{MC} = C_{\max}(u)$ |
| 14. | $\begin{array}{ccc} n & \leftarrow n & -S_{c_u}(u) \\ c & \leftarrow K+1 \end{array}$ |
| 16. | $K \leftarrow K + 1$ |
| 17: | else |
| 18: | $k_{	ext{old}} \leftarrow c_u$ |
| 19: | $c_u \leftarrow k_{\max}$ |
| 20: | $R^{MC} \leftarrow R^{MC} - S_{c_u}(u) + S_{k_{\max}}(u)$ |
| 21: | If cluster k_{old} is now empty, decrement c_v for all $v \in \mathcal{V}$ for which $c_v > k_{old}$, and |
| 22 | then decrement K. |
| 22: | end II ord for |
| 23. 24. | end while |
| 25: | if $R^{MC} > R^{MC}_{host}$ then |
| 26: | $c_{	ext{best}} \leftarrow c$ |
| 27: | $R_{	ext{best}}^{	ext{MC}} \leftarrow R^{	ext{MC}}$ |
| 28: | end if |
| 29: | end for |
| 50: | Teturn Cbest |
| bei dist sim | ng a pivot-based algorithm, is particularly sensitive to noise. Their approach maps each edge to a tinct arm (resulting in O(N²) arms), and they apply combinatorial bandit algorithms to estimate illarities. However, this approach presents several limitations: The algorithms are tailored to satisfy specific theoretical properties, which impose practical constraints (discussed below). |
| | • They assume non-persistent noise, meaning multiple queries of the same similarity (or multiple pulls of the same arm) are permitted—this is a standard assumption in the MAB literature. In contrast, our methods are robust even under persistent noise, where only a single query per pair is allowed. |
| | • Their strategy for selecting which similarities to query is limited as it does not account |
| | for correlations between pairs (arms), unlike our approach, which incorporates model |
| | uncertainty to guide query selection. |
| | • They consider two primary settings: |
| | - Fixed Confidence (KE-EC): This setting requires each arm to be pulled at least once |
| | leading to more queries than the total number of nairwise similaritiesthis is a common |
| | assumption for many MAB algorithms. Our methods, on the other hand, achieve |
| | effective clustering with significantly fewer queries. Additionally, their framework |
| | does not accommodate a predefined query budget B. |
| | - Fixed Budget (KF-FB): Although this setting allows for a predefined budget B , the smallest budget considered in (Kuroki et al., 2024) is $N^{2.1}$, which exceeds the total |

| 1566 | |
|------|--|
| 1567 | number of pairs. This constraint arises from the requirements of Algorithm 3 in (Kuroki $1, 2024$) which according to be described as N^2 to function any solution. |
| 1568 | the validity of their theoretical analysis. As demonstrated via extensive experiments |
| 1569 | our methods require significantly fewer queries |
| 1570 | our methods require significantly lewer queries. |
| 1571 | |
| 1570 | |
| 1572 | |
| 1575 | |
| 1574 | |
| 1575 | |
| 1570 | |
| 1577 | |
| 1578 | |
| 1579 | |
| 1580 | |
| 1581 | |
| 1582 | |
| 1583 | |
| 1584 | |
| 1585 | |
| 1586 | |
| 1587 | |
| 1588 | |
| 1589 | |
| 1590 | |
| 1591 | |
| 1592 | |
| 1593 | |
| 1594 | |
| 1595 | |
| 1596 | |
| 1597 | |
| 1598 | |
| 1599 | |
| 1600 | |
| 1601 | |
| 1602 | |
| 1603 | |
| 1604 | |
| 1605 | |
| 1606 | |
| 1607 | |
| 1608 | |
| 1609 | |
| 1610 | |
| 1611 | |
| 1612 | |
| 1613 | |
| 1614 | |
| 1615 | |
| 1616 | |
| 1617 | |
| 1618 | |
| 1010 | |
| 1013 | |