

# 000 001 002 003 004 005 A SCALABLE GLOBAL OPTIMIZATION ALGORITHM 006 FOR CONSTRAINED CLUSTERING 007 008 009

010 **Anonymous authors**  
011 Paper under double-blind review  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026

## ABSTRACT

027  
028  
029  
030  
031 Constrained clustering leverages limited domain knowledge to improve clustering  
032 performance and interpretability, but incorporating pairwise *must-link* and  
033 *cannot-link* constraints is an NP-hard challenge, making global optimization in-  
034 tractable. Existing mixed-integer optimization methods are confined to small-scale  
035 datasets, limiting their utility. We propose Sample-Driven Constrained Group-  
036 Based Branch-and-Bound (SDC-GBB), a decomposable branch-and-bound (BB)  
037 framework that collapses must-linked samples into centroid-based pseudo-samples  
038 and prunes cannot-link through geometric rules, while preserving convergence and  
039 guaranteeing global optimality. By integrating grouped-sample Lagrangian de-  
040 composition and geometric elimination rules for efficient lower and upper bounds,  
041 the algorithm attains highly scalable pairwise k-Means constrained clustering via  
042 parallelism. Experimental results show that our approach handles datasets with  
043 200,000 samples with cannot-link constraints and 1,500,000 samples with must-  
044 link constraints, which is 200 - 1500 times larger than the current state-of-the-art  
045 under comparable constraint settings, while reaching an optimality gap of  $\leq 3\%$ .  
046 In providing deterministic global guarantees, our method also avoids the search  
047 failures that off-the-shelf heuristics often encounter on large datasets.  
048  
049

## 1 INTRODUCTION

050  
051  
052  
053  
054  
055  
056  
057  
058  
059  
060  
061  
062  
063  
064  
065  
066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
0100  
0101  
0102  
0103  
0104  
0105  
0106  
0107  
0108  
0109  
0110  
0111  
0112  
0113  
0114  
0115  
0116  
0117  
0118  
0119  
0120  
0121  
0122  
0123  
0124  
0125  
0126  
0127  
0128  
0129  
0130  
0131  
0132  
0133  
0134  
0135  
0136  
0137  
0138  
0139  
0140  
0141  
0142  
0143  
0144  
0145  
0146  
0147  
0148  
0149  
0150  
0151  
0152  
0153  
0154  
0155  
0156  
0157  
0158  
0159  
0160  
0161  
0162  
0163  
0164  
0165  
0166  
0167  
0168  
0169  
0170  
0171  
0172  
0173  
0174  
0175  
0176  
0177  
0178  
0179  
0180  
0181  
0182  
0183  
0184  
0185  
0186  
0187  
0188  
0189  
0190  
0191  
0192  
0193  
0194  
0195  
0196  
0197  
0198  
0199  
0200  
0201  
0202  
0203  
0204  
0205  
0206  
0207  
0208  
0209  
0210  
0211  
0212  
0213  
0214  
0215  
0216  
0217  
0218  
0219  
0220  
0221  
0222  
0223  
0224  
0225  
0226  
0227  
0228  
0229  
0230  
0231  
0232  
0233  
0234  
0235  
0236  
0237  
0238  
0239  
0240  
0241  
0242  
0243  
0244  
0245  
0246  
0247  
0248  
0249  
0250  
0251  
0252  
0253  
0254  
0255  
0256  
0257  
0258  
0259  
0260  
0261  
0262  
0263  
0264  
0265  
0266  
0267  
0268  
0269  
0270  
0271  
0272  
0273  
0274  
0275  
0276  
0277  
0278  
0279  
0280  
0281  
0282  
0283  
0284  
0285  
0286  
0287  
0288  
0289  
0290  
0291  
0292  
0293  
0294  
0295  
0296  
0297  
0298  
0299  
0300  
0301  
0302  
0303  
0304  
0305  
0306  
0307  
0308  
0309  
0310  
0311  
0312  
0313  
0314  
0315  
0316  
0317  
0318  
0319  
0320  
0321  
0322  
0323  
0324  
0325  
0326  
0327  
0328  
0329  
0330  
0331  
0332  
0333  
0334  
0335  
0336  
0337  
0338  
0339  
0340  
0341  
0342  
0343  
0344  
0345  
0346  
0347  
0348  
0349  
0350  
0351  
0352  
0353  
0354  
0355  
0356  
0357  
0358  
0359  
0360  
0361  
0362  
0363  
0364  
0365  
0366  
0367  
0368  
0369  
0370  
0371  
0372  
0373  
0374  
0375  
0376  
0377  
0378  
0379  
0380  
0381  
0382  
0383  
0384  
0385  
0386  
0387  
0388  
0389  
0390  
0391  
0392  
0393  
0394  
0395  
0396  
0397  
0398  
0399  
0400  
0401  
0402  
0403  
0404  
0405  
0406  
0407  
0408  
0409  
0410  
0411  
0412  
0413  
0414  
0415  
0416  
0417  
0418  
0419  
0420  
0421  
0422  
0423  
0424  
0425  
0426  
0427  
0428  
0429  
0430  
0431  
0432  
0433  
0434  
0435  
0436  
0437  
0438  
0439  
0440  
0441  
0442  
0443  
0444  
0445  
0446  
0447  
0448  
0449  
0450  
0451  
0452  
0453  
0454  
0455  
0456  
0457  
0458  
0459  
0460  
0461  
0462  
0463  
0464  
0465  
0466  
0467  
0468  
0469  
0470  
0471  
0472  
0473  
0474  
0475  
0476  
0477  
0478  
0479  
0480  
0481  
0482  
0483  
0484  
0485  
0486  
0487  
0488  
0489  
0490  
0491  
0492  
0493  
0494  
0495  
0496  
0497  
0498  
0499  
0500  
0501  
0502  
0503  
0504  
0505  
0506  
0507  
0508  
0509  
0510  
0511  
0512  
0513  
0514  
0515  
0516  
0517  
0518  
0519  
0520  
0521  
0522  
0523  
0524  
0525  
0526  
0527  
0528  
0529  
0530  
0531  
0532  
0533  
0534  
0535  
0536  
0537  
0538  
0539  
0540  
0541  
0542  
0543  
0544  
0545  
0546  
0547  
0548  
0549  
0550  
0551  
0552  
0553  
0554  
0555  
0556  
0557  
0558  
0559  
0560  
0561  
0562  
0563  
0564  
0565  
0566  
0567  
0568  
0569  
0570  
0571  
0572  
0573  
0574  
0575  
0576  
0577  
0578  
0579  
0580  
0581  
0582  
0583  
0584  
0585  
0586  
0587  
0588  
0589  
0590  
0591  
0592  
0593  
0594  
0595  
0596  
0597  
0598  
0599  
0600  
0601  
0602  
0603  
0604  
0605  
0606  
0607  
0608  
0609  
0610  
0611  
0612  
0613  
0614  
0615  
0616  
0617  
0618  
0619  
0620  
0621  
0622  
0623  
0624  
0625  
0626  
0627  
0628  
0629  
0630  
0631  
0632  
0633  
0634  
0635  
0636  
0637  
0638  
0639  
0640  
0641  
0642  
0643  
0644  
0645  
0646  
0647  
0648  
0649  
0650  
0651  
0652  
0653  
0654  
0655  
0656  
0657  
0658  
0659  
0660  
0661  
0662  
0663  
0664  
0665  
0666  
0667  
0668  
0669  
0670  
0671  
0672  
0673  
0674  
0675  
0676  
0677  
0678  
0679  
0680  
0681  
0682  
0683  
0684  
0685  
0686  
0687  
0688  
0689  
0690  
0691  
0692  
0693  
0694  
0695  
0696  
0697  
0698  
0699  
0700  
0701  
0702  
0703  
0704  
0705  
0706  
0707  
0708  
0709  
07010  
07011  
07012  
07013  
07014  
07015  
07016  
07017  
07018  
07019  
07020  
07021  
07022  
07023  
07024  
07025  
07026  
07027  
07028  
07029  
07030  
07031  
07032  
07033  
07034  
07035  
07036  
07037  
07038  
07039  
07040  
07041  
07042  
07043  
07044  
07045  
07046  
07047  
07048  
07049  
07050  
07051  
07052  
07053  
07054  
07055  
07056  
07057  
07058  
07059  
07060  
07061  
07062  
07063  
07064  
07065  
07066  
07067  
07068  
07069  
07070  
07071  
07072  
07073  
07074  
07075  
07076  
07077  
07078  
07079  
07080  
07081  
07082  
07083  
07084  
07085  
07086  
07087  
07088  
07089  
07090  
07091  
07092  
07093  
07094  
07095  
07096  
07097  
07098  
07099  
070100  
070101  
070102  
070103  
070104  
070105  
070106  
070107  
070108  
070109  
070110  
070111  
070112  
070113  
070114  
070115  
070116  
070117  
070118  
070119  
070120  
070121  
070122  
070123  
070124  
070125  
070126  
070127  
070128  
070129  
070130  
070131  
070132  
070133  
070134  
070135  
070136  
070137  
070138  
070139  
070140  
070141  
070142  
070143  
070144  
070145  
070146  
070147  
070148  
070149  
070150  
070151  
070152  
070153  
070154  
070155  
070156  
070157  
070158  
070159  
070160  
070161  
070162  
070163  
070164  
070165  
070166  
070167  
070168  
070169  
070170  
070171  
070172  
070173  
070174  
070175  
070176  
070177  
070178  
070179  
070180  
070181  
070182  
070183  
070184  
070185  
070186  
070187  
070188  
070189  
070190  
070191  
070192  
070193  
070194  
070195  
070196  
070197  
070198  
070199  
070200  
070201  
070202  
070203  
070204  
070205  
070206  
070207  
070208  
070209  
070210  
070211  
070212  
070213  
070214  
070215  
070216  
070217  
070218  
070219  
070220  
070221  
070222  
070223  
070224  
070225  
070226  
070227  
070228  
070229  
070230  
070231  
070232  
070233  
070234  
070235  
070236  
070237  
070238  
070239  
070240  
070241  
070242  
070243  
070244  
070245  
070246  
070247  
070248  
070249  
070250  
070251  
070252  
070253  
070254  
070255  
070256  
070257  
070258  
070259  
070260  
070261  
070262  
070263  
070264  
070265  
070266  
070267  
070268  
070269  
070270  
070271  
070272  
070273  
070274  
070275  
070276  
070277  
070278  
070279  
070280  
070281  
070282  
070283  
070284  
070285  
070286  
070287  
070288  
070289  
070290  
070291  
070292  
070293  
070294  
070295  
070296  
070297  
070298  
070299  
070200  
070201  
070202  
070203  
070204  
070205  
070206  
070207  
070208  
070209  
070210  
070211  
070212  
070213  
070214  
070215  
070216  
070217  
070218  
070219  
070220  
070221  
070222  
070223  
070224  
070225  
070226  
070227  
070228  
070229  
070230  
070231  
070232  
070233  
070234  
070235  
070236  
070237  
070238  
070239  
070230  
070231  
070232  
070233  
070234  
070235  
070236  
070237  
070238  
070239  
070240  
070241  
070242  
070243  
070244  
070245  
070246  
070247  
070248  
070249  
070240  
070241  
070242  
070243  
070244  
070245  
070246  
070247  
070248  
070249  
070250  
070251  
070252  
070253  
070254  
070255  
070256  
070257  
070258  
070259  
070250  
070251  
070252  
070253  
070254  
070255  
070256  
070257  
070258  
070259  
070260  
070261  
070262  
070263  
070264  
070265  
070266  
070267  
070268  
070269  
070260  
070261  
070262  
070263  
070264  
070265  
070266  
070267  
070268  
070269  
070270  
070271  
070272  
070273  
070274  
070275  
070276  
070277  
070278  
070279  
070280  
070281  
070282  
070283  
070284  
070285  
070286  
070287  
070288  
070289  
070280  
070281  
070282  
070283  
070284  
070285  
070286  
070287  
070288  
070289  
070290  
070291  
070292  
070293  
070294  
070295  
070296  
070297  
070298  
070299  
070200  
070201  
070202  
070203  
070204  
070205  
070206  
070207  
070208  
070209  
070210  
070211  
070212  
070213  
070214  
070215  
070216  
070217  
070218  
070219  
070220  
070221  
070222  
070223  
070224  
070225  
070226  
070227  
070228  
070229  
070230  
070231  
070232  
070233  
070234  
070235  
070236  
070237  
070238  
070239  
070240  
070241  
070242  
070243  
070244  
070245  
070246  
070247  
070248  
070249  
070250  
070251  
070252  
070253  
070254  
070255  
070256  
070257  
070258  
070259  
070260  
070261  
070262  
070263  
070264  
070265  
070266  
070267  
070268  
070269  
070270  
070271  
070272  
070273  
070274  
070275  
070276  
070277  
070278  
070279  
070280  
070281  
070282  
070283  
070284  
070285  
070286  
070287  
070288  
070289  
070290  
070291  
070292  
070293  
070294  
070295  
070296  
070297  
070298  
070299  
070200  
070201  
070202  
070203  
070204  
070205  
070206  
070207  
070208  
070209  
070210  
070211  
070212  
070213  
070214  
070215  
070216  
070217  
070218  
070219  
070220  
070221  
070222  
070223  
070224  
070225  
070226  
070227  
070228  
070229  
070230  
070231  
070232  
070233  
070234  
070235  
070236  
070237  
070238  
070239  
070240  
070241  
070242  
070243  
070244  
070245  
070246  
070247  
070248  
070249  
070250  
070251  
070252  
070253  
070254  
070255  
070256  
070257  
070258  
070259  
070260  
070261  
070262  
070263  
070264  
070265  
070266  
070267  
070268  
070269  
070270  
070271  
070272  
070273  
070274  
070275  
070276  
070277  
070278  
070279  
070280  
070281  
070282  
070283  
070284  
070285  
070286  
070287  
070288  
070289  
070290  
070291  
070292  
070293  
070294  
070295  
070296  
070297  
070298  
070299  
070200  
070201  
070202  
070203  
070204  
070205  
070206  
070207  
070208  
070209  
070210  
070211  
070212  
070213  
070214  
070215  
070216  
070217  
070218  
070219  
070220  
070221  
070222  
070223  
070224  
070225  
070226  
070227  
070228  
070229  
070230  
070231  
070232  
070233  
070234  
070235  
070236  
070237  
070238  
070239  
070240  
070241  
070242  
070243  
070244  
070245  
070246  
070247  
070248  
070249  
070250  
070251  
070252  
070253  
070254  
070255  
070256  
070257  
070258  
070259  
070260  
070261  
070262  
070263  
070264  
070265  
070266  
070267  
070268  
070269  
070270  
070271  
070272  
070273  
070274  
070275  
070276  
070277  
070278  
070279  
070280  
070281  
070282  
070283  
070284  
070285  
070286  
070287  
070288  
070289  
070290  
070291  
070292  
070293  
070294  
070295  
070296  
070297  
070298  
070299  
070200  
070201  
070202  
070203  
070204  
070205  
070206  
070207  
070208  
070209  
070210  
070211  
070212  
070213  
070214  
070215  
070216  
070217  
070218  
070219  
070220  
070221  
070222  
070223  
070224  
070225  
070226  
070227  
070228  
070229  
070230  
070231  
070232  
070233  
070234  
070235  
070236  
070237  
070238  
070239  
070240  
070241  
070242  
070243  
070244  
070245  
070246  
070247  
070248  
070249  
070250  
070251  
070252  
070253  
070254  
070255  
070256  
070257  
070258  
070259  
070260  
070261  
070262  
070263  
070264  
070265  
070266  
070267  
070268  
070269  
070270  
070271  
070272  
070273  
070274  
070275  
070276  
070277  
070278  
070279  
070280  
070281  
070282  
070283  
070284  
070285  
070286  
070287  
070288  
070289  
070290  
070291  
070292  
070293  
070294  
070295  
070296  
070297  
070298  
070299  
070200  
070201  
070202  
070203  
07020

054 Constraint programming approaches (Dao et al., 2013; 2015; 2017) offer flexibility in incorporating  
 055 various constraint forms but generally do not scale beyond a few hundred points. Mixed-Integer  
 056 Programming (MIP) formulations have also been explored to handle additional cluster-level or  
 057 instance-level constraints in MSSC. For example, (Tang et al., 2020) proposed an iterative scheme  
 058 that reformulates a size-constrained MSSC problem into a mixed-integer linear program, leveraging  
 059 the unimodularity of certain constraint matrices to reduce computational complexity. Similarly,  
 060 (Liberti & Manca, 2022) examined several side-constrained MSSC models cast as Mixed-Integer  
 061 Nonlinear Programs, some featuring convex relaxations that enable global optimization techniques.  
 062 While these MIP-based approaches provide a powerful and flexible framework for ensuring feasibility  
 063 under various constraint types, their applicability remains limited by high computational overhead,  
 064 restricting them to relatively small or moderate-sized datasets.

065 Branch-and-bound methods have also been specialized for constrained MSSC. (Guns et al., 2016)  
 066 proposed the Constraint Programming Repetitive Branch-and-Bound Algorithm (CPRBBA), which  
 067 augments Brusco’s repetitive branch-and-bound procedure (Brusco, 2006) with a constraint program-  
 068 ming solver to compute tight lower and upper bounds on subsets of objects of increasing size. This  
 069 approach, while effective on small instances, remains limited to fewer than 200 data points. (Piccialli  
 070 et al., 2022a) developed the PC-SOS-SDP algorithm, which integrates *must-link* and *cannot-link*  
 071 constraints into a semidefinite programming framework, scaling to a few thousand data points but  
 072 not beyond (Baumann & Hochbaum, 2024). These exact methods do not generally account for  
 073 soft constraints and remain computationally expensive for larger datasets. Nonetheless, continuing  
 074 progress in algorithmic design and hardware (Bertsimas & Dunn, 2017) has widened the scope of  
 075 exact methods for constrained clustering.

076 **Our Contributions** In this paper, we propose a scalable deterministic global optimization algorithm  
 077 for the minimum sum-of-squared clustering (MSSC) task *with pairwise ML and CL constraints*. We  
 078 introduce a centroid-based pseudosample formulation for must-link subsets, leveraging the combined  
 079 information of each group to maintain the exact global minimum while reducing problem complexity.  
 080 We devise geometric sample-determination rules that eliminate cannot-links, which specify whether  
 081 points must not be placed into the same clusters before enumeration. We design a branch-and-bound  
 082 algorithm that branches only on the cluster-center variables. This avoids combinatorial branching  
 083 on sample-to-cluster assignments, thus achieving a globally  $\epsilon$ -optimal solution even for large-scale  
 084 datasets. Our analysis proves convergence under exhaustive subdivisions of the feasible region for  
 085 the center variables.

086 **Capability For More than One Hundred Thousand Scale Problems** We present an open-source  
 087 implementation in *Julia* that solves constrained MSSC instances of up to 1,500,000 samples for  
 088 the ML case and 200,000 samples for CL case with optimality guarantee or very low optimality gaps.  
 089 This corresponds to 1500-fold and 200-fold increases in scale, respectively, over the current exact  
 090 state-of-the-art (Piccialli et al., 2022a). This framework thus enables deterministic global clustering  
 091 solutions for large-scale datasets previously considered intractable.

## 092 2 MIXED-INTEGER PROGRAMMING FOR PAIRWISE-CONSTRAINED $k$ -MEANS

093 Given a dataset  $X = \{x_1, \dots, x_S\} \subset \mathbb{R}^m$  with  $S$  samples and  $m$  attributes, the semi-supervised  
 094 MSSC task with pairwise constraints seeks a set of  $k$  clusters that minimizes the Sum of Squared  
 095 Errors (SSE) subject to must-link (ML) and cannot-link (CL) requirements:  
 096

$$097 \min_b \sum_{s \in \mathcal{S}} \sum_{k \in \mathcal{K}} b_{s,k} \|x_s - \mu_k\|_2^2 \quad (1a)$$

$$100 \text{ s.t. } b_{s,k} = b_{s',k}, \quad \forall (s, s') \in \mathcal{T}_{ml}, k \in \mathcal{K}, \quad (1b)$$

$$101 \quad b_{s,k} + b_{s',k} \leq 1, \quad \forall (s, s') \in \mathcal{T}_{cl}, k \in \mathcal{K}, \quad (1c)$$

$$102 \quad b_{s,k} \in \{0, 1\}, \quad \forall s \in \mathcal{S}, k \in \mathcal{K}. \quad (1d)$$

$$103 \quad \sum_{k \in \mathcal{K}} b_{s,k} = 1 \quad (1e)$$

106 where  $s \in \mathcal{S} := \{1, \dots, S\}$  is the data sample set,  $k \in \mathcal{K} = \{1, \dots, K\}$  is the cluster set,  
 107  $\mu := [\mu_1, \dots, \mu_K]$ , where  $\mu_k \in \mathbb{R}^m$  represents the center of each cluster,  $b_{s,k} \in \{0, 1\}$  is equal to

108 1 if  $x_s$  belongs to the  $k$ -th clusters and 0 otherwise.  $\mathcal{T}_{ml} \subseteq \mathcal{S} \times \mathcal{S}$  and  $\mathcal{T}_{cl} \subseteq \mathcal{S} \times \mathcal{S}$  are the sets of  
 109 tuples indicating whether samples must or must not reside in the same cluster respectively.  
 110

111 The MSSC with pairwise constraints (Problem 1) can be reformulated as an SSE optimization problem  
 112 of the following form:

$$113 \min_{\mu, d, b} \sum_{s \in \mathcal{S}} d_{s,*} \quad (2a)$$

$$115 \text{s.t. } -N(1 - b_{s,k}) \leq d_{s,*} - d_{s,k} \leq N(1 - b_{s,k}) \quad (2b)$$

$$117 d_{s,k} \geq \|x_s - \mu_k\|_2^2 \quad \forall s \in \mathcal{S}, \forall k \in \mathcal{K} \quad (2c)$$

$$118 \text{Constraints 1b- 1e} \quad (2d)$$

119 Here  $d_{s,k}$  is the distance between  $x_s$  and  $\mu_k$ ,  $d_{s,*}$  is the distance from  $x_s$  to its assigned centroid, and  
 120  $N$  is a big- $M$  constant. Define  $d_s = [d_{s,1}, \dots, d_{s,K}, d_{s,*}]$ ,  $d = [d_1, \dots, d_S]$ ,  $b_s = [b_{s,1}, \dots, b_{s,K}]$ ,  
 121  $b = [b_1, \dots, b_S]$ . Constraint (2b) links  $d_{s,*}$  and  $d_{s,k}$  when  $b_{s,k} = 1$ . Problem 2 is a mixed-integer  
 122 second-order cone program (MISOCP) admits a two-stage extensive form (see Appendix A). While  
 123 off-the-shelf solvers like Gurobi (Gurobi, 2024) and CPLEX (Cplex, 2022) can handle small instances,  
 124 they become intractable even at moderate sample sizes (e.g.,  $S = 800$ ) (Piccialli et al., 2022a).  
 125

### 126 3 REDUCED-SPACE BRANCH-AND-BOUND ALGORITHM

127 Reduced-space branch-and-bound frameworks have demonstrated significant scalability gains by  
 128 partitioning only the centroid search space (Cao & Zavala, 2019). We tailor this scheme to the  
 129 pairwise-constrained clustering problem by integrating geometric probing rules derived from the  
 130 MISOCP formulation in Sec. 2 to tighten both lower and upper bounds. In particular, we exploit the  
 131 implicit inequality that any subregion whose lower bound exceeds the current best upper bound can  
 132 be discarded outright.

#### 133 3.1 GEOMETRIC SAMPLE DETERMINATION RULES

134 We first observe that every feasible clustering is subject to two straightforward geometric bounds  
 135 relative to the incumbent solution. Let  $\rho = \max_{s \in \mathcal{S}} \|\mathbf{x}_s - \boldsymbol{\mu}_{k(s)}^{\text{best}}\|_2^2$  be the worst-case squared distance  
 136 between each sample and the centroid to which it is assigned in the current best solution. Thus,  
 137  $\rho$  represents the maximum per-sample contribution to the current incumbent cost and is used as a  
 138 per-sample upper bound. Then, for any region  $M_k$  (an axis-aligned box in  $\mathbb{R}^m$ ) containing the true  
 139 optimal  $\mu_k$ , we can compute the minimal and maximal possible squared distances:

$$140 d_{\min}(\mathbf{x}_s, M_k) = \min_{\mu \in M_k} \|\mathbf{x}_s - \mu\|_2^2, \quad d_{\max}(\mathbf{x}_s, M_k) = \max_{\mu \in M_k} \|\mathbf{x}_s - \mu\|_2^2.$$

141 Because  $\rho$  is an upper bound on the true assignment cost, any candidate pair  $(s, k)$  with  
 142  $d_{\min}(\mathbf{x}_s, M_k) > \rho$  can never be optimal.

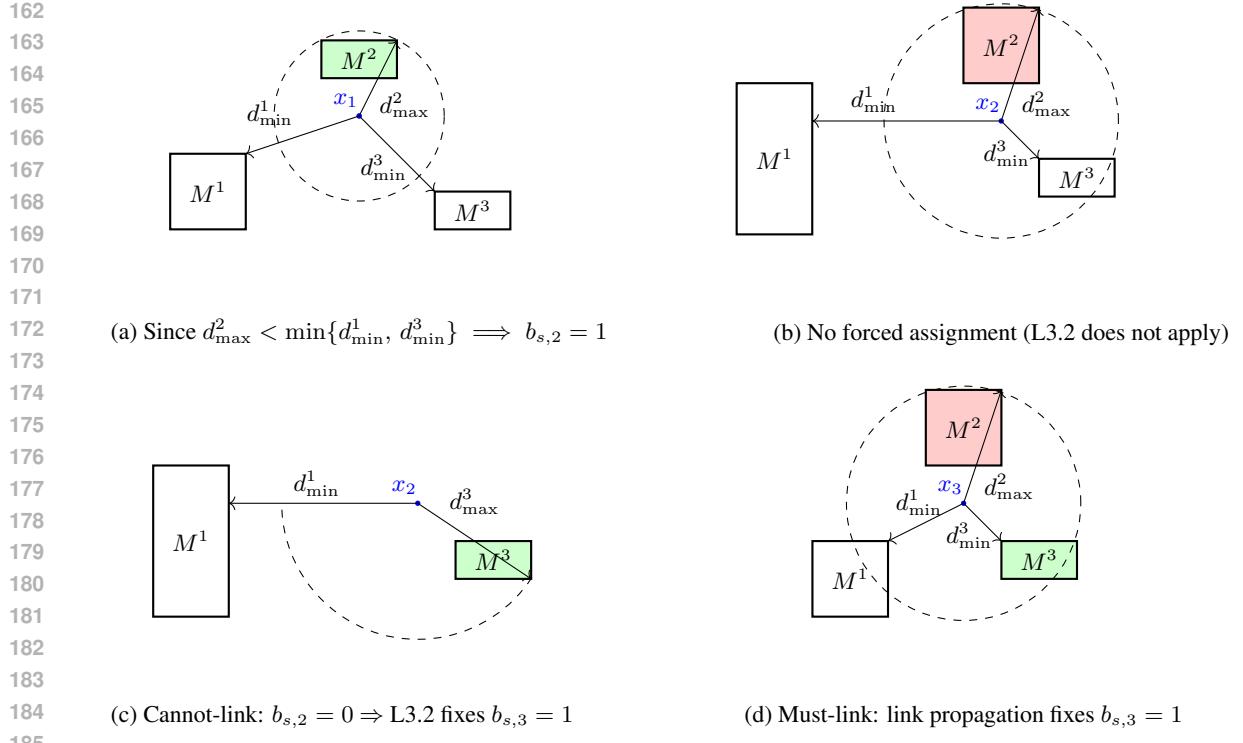
143 **Lemma 3.1** (Early-elimination). *For any sample  $s \in \mathcal{S}$  and any cluster region  $M_k$ ,  $d_{\min}(\mathbf{x}_s, M_k) >$   
 144  $\rho \implies b_{s,k} = 0$  in every optimal solution with objective value not larger than the incumbent.*

145 *Proof.* By definition  $d_{\min}(\mathbf{x}_s, M_k) = \min_{\mu \in M_k} \|\mathbf{x}_s - \mu\|_2^2$ . If  $d_{\min} > \rho$ , then for every  $\mu \in M_k$  one  
 146 has  $\|\mathbf{x}_s - \mu\|_2^2 > \rho$ . Since  $\rho$  is an upper bound on the per-sample cost in the incumbent, assigning  
 147  $\mathbf{x}_s$  to cluster  $k$  would yield a contradiction. Hence  $b_{s,k} = 0$  in any cluster whose overall cost does  
 148 not exceed the incumbent cost.  $\square$

149 A complementary rule arises from comparing the worst-case assignment cost in one region to the  
 150 best-case cost in the others.

151 **Lemma 3.2** (Forced assignment). *Fix a sample  $s$  and let  $k^+ \in \mathcal{K}$  satisfy  $d_{\max}(\mathbf{x}_s, M_{k^+}) <$   
 152  $\min_{k \neq k^+} d_{\min}(\mathbf{x}_s, M_k)$ . Then  $b_{s,k^+} = 1$  in every optimal solution.*

153 *Proof.* For any  $\mu^+ \in M_{k^+}$  and any  $\mu \in M_k$  with  $k \neq k^+$  we have  $\|\mathbf{x}_s - \mu^+\|_2^2 \leq d_{\max}(\mathbf{x}_s, M_{k^+}) <$   
 154  $d_{\min}(\mathbf{x}_s, M_k) \leq \|\mathbf{x}_s - \mu\|_2^2$ . Thus the distance from  $\mathbf{x}_s$  to every center in  $M_{k^+}$  is strictly smaller  
 155 than the distance to any center in the remaining regions, implying that the unique cost-minimising  
 156 assignment is  $b_{s,k^+} = 1$ .  $\square$

Figure 1: Illustration of sample-determination via link propagation for  $K = 3$ .

190 Figure 1 illustrates interaction variations of geometric checks and pairwise constraints for data points  
191  $x_1, x_2$  and  $x_3$ . In (a), the distance bounds immediately fix  $x_1$  in  $M^2$ . In (b), the bounds overlap,  
192 distances are inconclusive and no assignment is made. In (c), the cannot-link ( $x_1, x_2$ ) rules out  $M^2$   
193 for  $x_2$ , after which the geometric test fixes  $x_2$  in  $M^3$ . In (d), the must-link ( $x_2, x_3$ ) propagates that  
194 assignment to  $x_3$ . This sequence shows how geometry and ML/CL constraints jointly determine  
195 assignments prior to branching. If the ML/CL constraints forbid the move  $(x_s \rightarrow k^+)$ , the node  
196 becomes infeasible and is pruned.

### 3.2 EQUIVALENT UNCONSTRAINED CLUSTERING PROBLEM

200 Although the geometric sample-determination rules and link propagation of Section 3.1 eliminate  
201 most binary assignments, the remaining must-link constraints still couple samples and inflate the  
202 **branch-and-bound** (BB) complexity. To isolate this effect, consider the ML-only version of  
203 Problem (2). Different from the *unconstrained* MSSC problem, pairwise constraint clustering  
204 problem contains a family of equalities  $b_{s,k} = b_{s',k}$  for  $(s, s') \in \mathcal{T}_{ml}$ . Nevertheless, we show that  
205 collapsing each must-link component into repeated pseudo-samples yields an unconstrained instance  
206 with identical global optimum.

207 Let a cluster  $\mathcal{C} = \{x_1, x_2, \dots, x_p\}$  and let  $\mu$  denote an arbitrary centroid for  $\mathcal{C}$ . Without loss of  
208 generality, assume  $x_1, \dots, x_t \in \mathcal{C}_{ml} \subseteq \mathcal{C}$  form a single *must-link component* inside  $\mathcal{C}$ . Let  $\mu_{ml}$   
209 denote the centroid of  $\mathcal{C}_{ml}$  and  $tr(\Sigma_{ml}^2)$  the trace of the covariance matrix of  $\mathcal{C}_{ml}$ . Thus, we have:  
210  $\mu_{ml} = \frac{1}{t} \sum_{i=1}^t x_i$ ,  $tr(\Sigma_{ml}^2) = \frac{1}{t-1} \sum_{i=1}^t \|x_i - \mu_{ml}\|^2$ .

211 **Lemma 3.3.** *Given 2 clusters,  $\mathcal{C} = \{x_1, x_2, \dots, x_p\}$  and  $\hat{\mathcal{C}} = \{x_{t+1}, x_{t+2}, \dots, x_p, \underbrace{\mu_{ml}, \dots, \mu_{ml}}_t\}$ ,*  
212 *let  $sse_{\mathcal{C}}(\mu)$  and  $sse_{\hat{\mathcal{C}}}(\mu)$  denote their respective within-cluster SSE computed with centroid  $\mu$ . Then,*  
213 *the following identity holds:*

$$sse_{\mathcal{C}}(\mu) = sse_{\hat{\mathcal{C}}}(\mu) + (t-1)tr(\Sigma_{ml}^2). \quad (3)$$

216 Note that  $(t - 1)tr(\Sigma_{ml}^2)$  is an additive constant that does not affect optimization over centroids.  
 217 Figure 2 illustrates this construction of pseudo-samples using a small example. Based on Lemma 3.3,  
 218 we form the dataset:  
 219

$$220 \hat{X} = \bigcup_{k \in \mathcal{K}_{ml}} \underbrace{\{\mu_{ml,k}, \dots, \mu_{ml,k}\}}_{t_k} \cup (X \setminus \{x_s | (s, s') \in \mathcal{T}_{ml}, s' \in \mathcal{S}\}),$$

$$221$$

222 along with the corresponding unconstrained MSSC optimization problem:  
 223

$$224 \min_{\mu, d, b} \sum_{s \in \hat{\mathcal{S}}} d_{s,*} + \sum_{k \in \mathcal{K}_{ml}} (t_k - 1)tr(\Sigma_{ml,k}^2) \quad (4a)$$

$$225$$

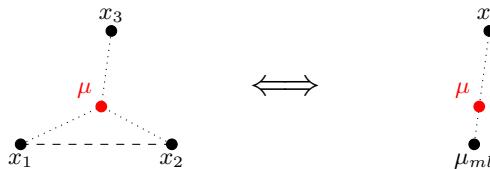
$$226 \text{ s.t. Constraints 2b, 2c, 1e, 1d} \quad (4b)$$

$$227$$

228 where  $\mu_{ml,k}$  and  $\Sigma_{ml,k}^2$  represent the mean (centroid) and covariance matrix of the must-link samples  
 229 within cluster  $k$ , respectively.  $\mathcal{K}_{ml}$  denotes the set of clusters containing must-link samples, and  $\hat{\mathcal{S}}$  is  
 230 the corresponding index set for the dataset  $\hat{X}$ .  
 231

232 **Theorem 3.4.** *If  $\mu^*$  and  $z(\mu^*)$  are the global optimal solution and cost of Problem (4), then they are  
 233 also the global optimum and cost of the ML-only problem (2), and vice versa.*

234 **Mixed ML and CL constraints.** When both ML and CL constraints are present, collapse each ML  
 235 component as above to obtain  $\hat{X}$  with only CL constraints. Then apply Lemma 3.1 to eliminate all  
 236 CL-violating assignments, yielding an unconstrained MSSC on the reduced dataset.  
 237



245 Figure 2: Left: three samples  $x_1, x_2, x_3$  with must-link  $(x_1, x_2)$  and centroid  $\mu$ . Right: collapse of  
 246  $\{x_1, x_2\}$  into two pseudo-samples at  $\mu_{ml} = \frac{1}{2}(x_1 + x_2)$  preserves the optimal centroid  $\mu$ .  
 247

### 248 3.3 UPPER BOUNDING STRATEGIES

250 Let  $M_0 \subset \mathbb{R}^{mk}$  be the initial axis-aligned box for the full centroid vector  $\mu = \{\mu_1, \dots, \mu_K\}$ . The  
 251 branch-and-bound (BB) algorithm requires a fast yet tight lower bound (LB) for each subproblem  
 252 or node  $M$  inside the solution space  $M_0$ . In this section two methods for computing upper bounds  
 253 at each node of the BB scheme are presented. The first method handles CL constraints through a  
 254  $k$ -coloring interpretation. The second method derives a closed-form expression applicable when only  
 255 ML constraints are present. When both constraint types coexist, ML constraints can be collapsed  
 256 as described in Section 3.2, reformulating the problem into one involving only CL constraints and  
 257 allowing the use of the  $k$ -coloring approach.  
 258

259  **$K$ -Coloring Bound for CL Constraints.** Let  $K$  denote the prescribed number of clusters and  
 260 let  $G_{cl} = (\hat{\mathcal{S}}, \mathcal{T}_{cl})$  be the CL graph. At node  $M \subseteq M_0$  maintain a pool of  $C$  centroid candidates  
 261  $\{\mu^{(c)}\}_{c=1}^C \subset M$ . For each candidate  $c$  assign each sample  $s \in \hat{\mathcal{S}}$  to its closest centroid,  $k_s^{(c)} =$   
 262  $\arg \min_{k \in \{1, \dots, K\}} \|x_s - \mu_k^{(c)}\|_2^2$ , and let  $\chi^{(c)}(s) = k_s^{(c)}$ . The labeling  $\chi^{(c)}$  is *proper* if  $\chi^{(c)}(u) \neq$   
 263  $\chi^{(c)}(v)$  for every  $(u, v) \in \mathcal{T}_{cl}$ . Define

$$264 z_{ub}^{(c)} = \begin{cases} \sum_{s \in \hat{\mathcal{S}}} \|x_s - \mu_{\chi^{(c)}(s)}^{(c)}\|_2^2, & \text{if } \chi^{(c)} \text{ is proper,} \\ +\infty, & \text{otherwise.} \end{cases}$$

$$265$$

$$266$$

$$267$$

268 The node upper bound is  $\alpha(M) = \min_{1 \leq c \leq C} z_{ub}^{(c)}$ . If  $\chi^{(c)}$  is proper, then  $z_{ub}^{(c)}$  satisfies all CL  
 269 constraints and  $z(M) \leq \alpha(M)$ .  
 270

270 **Closed-Form Bound for ML Constraints.** With only must-link (ML) constraints, fix any feasible  
 271 centroid set  $\hat{\mu} \in M$  and compute  $\alpha(M) = \sum_{s \in \hat{\mathcal{S}}} Q_s(\hat{\mu})$ , where each  $Q_s(\hat{\mu})$  admits a closed form.  
 272 The expression yields an admissible bound because every  $\hat{\mu} \in M$  respects the ML constraints,  
 273 implying  $z(M) \leq \alpha(M)$ . An initial bound is produced at the root via a heuristic ( $k$ -means). Bounds  
 274 at descendant nodes are updated with candidates extracted from the relaxations in Section 3.4. The  
 275 BB algorithm terminates with the optimal objective value  $\alpha + \sum_{k \in \mathcal{K}_{\text{ml}}} q_k \sigma_{\text{ml},k}^2$ .  
 276

### 277 3.4 LOWER BOUNDING STRATEGY WITH GROUPING-BASED LAGRANGIAN DECOMPOSITION

279 The branch-and-bound (BB) algorithm requires a fast yet tight lower bound (LB) for each sub-  
 280 problem or node  $M$  inside the solution space  $M_0$ . An effective strategy to achieve tighter lower  
 281 bounds is through Lagrangian decomposition (LD), in which the corresponding non-anticipativity  
 282 constraints Cao & Zavala (2019) are dualized with fixed Lagrange multipliers  $\lambda$  and added to the  
 283 objective function (Karuppiah & Grossmann, 2008). However, to reduce problem size and improve  
 284 relaxation quality, instead of associating each sample with a separate subproblem (as mentioned  
 285 in Karuppiah & Grossmann (2008)), we partition the sample set  $\hat{\mathcal{S}}$  into  $G$  disjoint groups  $\hat{\mathcal{S}}_1, \dots, \hat{\mathcal{S}}_G$   
 286 with index set  $\mathcal{G} = 1, \dots, G$ , such that  $\bigcup_g \hat{\mathcal{S}}_g = \hat{\mathcal{S}}$  and  $\hat{\mathcal{S}}_i \cap \hat{\mathcal{S}}_g = \emptyset$  for  $i \neq g$ . Instead of replicating  
 287 center variables per sample, we assign one per group and enforce consistency through:  
 288

$$\min_{\mu_g \in M} \sum_{g \in \mathcal{G}} Q_g(\mu_g), \quad Q_g(\mu_g) := \sum_{s \in \hat{\mathcal{S}}_g} Q_s(\mu_g) \quad \text{s.t.} \quad \mu_g = \mu_{g+1}, \quad \forall g \in 1, \dots, G-1 \quad (5a)$$

291 Dualizing the coupling constraints with multipliers  $\lambda$  yields a tighter lower bound via:

$$\beta^{SG+LD}(M) := \max_{\lambda} \beta^{SG+LD}(M, \lambda) \quad (6)$$

295 This grouped formulation preserves intra-group non-anticipativity while relaxing inter-group consis-  
 296 tency, yielding  $\beta^{LD}(M) \leq \beta^{SG+LD}(M) \leq z(M)$ . While solving (6) requires iterative MISOCPs,  
 297 it significantly strengthens the bound. The grouping is fixed at the BB root for efficiency.

### 299 3.5 BRANCH-AND-BOUND CLUSTERING SCHEME

300 We adopt the framework of the reduced-space branch-and-bound scheme from (Cao & Zavala, 2019)  
 301 and tailor the algorithm for the pairwise constrained clustering task. Algorithm 1 depicts the details  
 302 of the algorithm, where  $\beta$  and  $\alpha$  represent the function of lower and upper bound, respectively. With  
 303 the lower and upper bounding strategy provided in the following subsections.

305 **Theorem 3.5.** *Given an exhaustive subdivision on  $\mu$ , Algorithm 1 converges in the sense that*

$$\lim_{i \rightarrow \infty} \alpha_i = \lim_{i \rightarrow \infty} \beta_i = z. \quad (7)$$

308 The proof is shown in Appendix C.

## 310 4 COMPUTATIONAL EXPERIMENTS

312 We implemented our algorithm, **Sample-Driven Constrained Group-Based Branch-and-Bound**  
 313 (**SDC-GBB**), in **Julia 1.10.3** and evaluated its performance on synthetic and real-world datasets  
 314 using a high-performance cluster comprising nodes with 128 AMD Epyc 7702 CPUs (2.0GHz) and  
 315 1TB of RAM. Computational experiments were conducted under both serial and parallel config-  
 316urations, comparing SDC-GBB against the branch-and-bound (BB) algorithm in CPLEX 22.1.1  
 317 (Cplex, 2022), the exact method PC-SOS-SDP (Piccialli et al., 2022a), and the best heuristic out  
 318 of the following algorithms: COP- $k$ -means (Wagstaff et al., 2001), encode-kmeans-post (Nghiem  
 319 et al., 2020), BLPKM-CC (Baumann, 2020), and Sensitivity Sampling coresets algorithm (Feldman  
 320 & Langberg, 2011). All heuristic algorithms were run with 100 restarts (Wagstaff et al., 2001) in 4  
 321 hours, and the full results are reported in Appendix E. For parallel applications, subproblems were  
 322 distributed across multiple CPU cores with group sizes limited to  $\min(162/d - k, 10 \times k)$  during the  
 323 lower bound decomposition process. Performance was assessed using Upper Bound (UB), relative  
 324 optimality gap, and the number of BB nodes resolved. UB represents the best feasible solution found,

---

324 **Algorithm 1** Branch-and-Bound Clustering with Geometric Sample Determination

---

```

325
326 Inputs:  $X = \{x_s\}_{s \in S} \subset \mathbb{R}^d$ ,  $K$ ,  $\mathcal{T}_{ml}$ ,  $\mathcal{T}_{cl}$ 
327 Initialization
328 Initialize  $i = 0$ ,  $\mathbb{M} \leftarrow \{M_0\}$ , tolerance  $\epsilon > 0$ 
329 Compute upper bound  $\alpha_i = \alpha(M_0)$ , lower bound  $\beta_i = \beta(M_0)$ ;
330 Geometric Sample Determination
331 Compute  $d_{\min}(x_s, M_{0,k})$ ,  $d_{\max}(x_s, M_{0,k})$ ;
332 if  $d_{\min}(x_s, M_{0,k}) > \alpha_i$  then  $b_{s,k} \leftarrow 0$  (update  $K_s$ );
333 if  $\{\exists k^+ \text{ with } d_{\max}(x_s, M_{0,k^+}) < \min_{k \neq k^+} d_{\min}(x_s, M_{0,k})\}$  then  $b_{s,k^+} \leftarrow 1$ ;
334 Propagate fixes via  $\mathcal{T}_{ml}$ ,  $\mathcal{T}_{cl}$ ; update  $\{K_s\}$ ;
335 repeat
336     Node Selection
337         Select a set  $M \in \mathbb{M}$  satisfying  $\beta(M) = \beta_i$ ;
338          $\mathbb{M} \leftarrow \mathbb{M} \setminus \{M\}$ ;
339          $i \leftarrow i + 1$ ;
340     Branching
341         Partition  $M$  into subsets  $M_1$  and  $M_2$  with  $\text{relint}(M_1) \cap \text{relint}(M_2) = \emptyset$ ;
342         Add each subset to  $\mathbb{M}$  to create separated child nodes;
343     Bounding
344         Compute  $\alpha(M_1)$ ,  $\beta(M_1)$ ,  $\alpha(M_2)$ ,  $\beta(M_2)$ ;
345          $\beta_i \leftarrow \min\{\beta(M') \mid M' \in \mathbb{M}\}$ ;
346          $\alpha_i \leftarrow \min\{\alpha_{i-1}, \alpha(M_1), \alpha(M_2)\}$ ;
347         Remove all  $M'$  from  $\mathbb{M}$  if  $\beta(M') \geq \alpha_i$ ;
348         If  $|\beta_i - \alpha_i| \leq \epsilon$ , STOP;
349     until  $\mathbb{M} = \emptyset$ 
350 Output  $\hat{\mu}$ ,  $\hat{b}$  and  $z^* = \alpha_i + \sum_{k \in \mathcal{K}_{ml}} (t_k - 1) \text{tr}(\Sigma_{ml,k}^2)$ 

```

---

351 while the relative optimality gap is calculated as  $\frac{\alpha_l - \beta_l}{\min(\alpha_l, \beta_l)} \times 100\%$ , where  $\alpha_l$  and  $\beta_l$  denote the  
352 best lower and upper bounds, respectively. The number of resolved BB nodes indicates the total BB  
353 iterations performed. Unlike heuristic methods, deterministic global optimization methods provide  
354 an optimality gap, enabling quantitative assessment of solution quality.

355 We evaluate the selected algorithms on 8 real-world datasets taken or sampled from the UCI Machine  
356 Learning Repository (Dua & Graff, 2017), Hemicellulose (Wang et al., 2022), PR2392 (Padberg &  
357 Rinaldi, 1991), and 7 synthetic datasets generated with 2 features, 3 Gaussian clusters and a fixed  
358 random seed (`seed` = 1). Datasets are categorized as small ( $n \leq 1,000$ ), medium ( $n \leq 10,000$ ),  
359 large ( $n \leq 100,000$ ), and huge ( $n \geq 100,000$ ), where  $n$  is the number of samples<sup>1</sup>. We follow  
360 the pairwise constraint generation practice with the same classic random-pair sampling pipeline as  
361 (Piccialli et al., 2022a; Aloise et al., 2009; Babaki et al., 2014; Guns et al., 2016). Across three  
362 separate experiments, namely must-link only (ML-only), cannot-link only (CL-only), and both  
363 must-link and cannot-link (ML+CL), each dataset has  $\frac{n}{4}$  samples bounded by ML constraints,  $\frac{n}{4}$   
364 samples bounded by CL constraints, and combining  $\frac{n}{4}$  ML with  $\frac{n}{4}$  CL constraints respectively. when  
365 the optimality gap dropped below 0.1 %, when runtime reached 4 hours for datasets with  $n \leq 10,000$   
366 or 12 h for those with  $n > 10,000$ , or when 5 million nodes had been explored.

367 **4.1 NUMERICAL RESULTS**

369 Our work focuses explicitly on solution quality in terms of the MSSC cost, providing an optimality  
370 guarantee. With this, SDC-GBB matches the performance of commercial solvers and the state-of-  
371 the-art algorithm PC-SOS-SDP (Piccialli et al., 2022a) through *super-point aggregation, targeted*  
372 *decomposition, and tight bounding via geometric partitioning* under ML constraints, as well as  
373 *geometric sample determination rules* that prune infeasible assignments under CL constraints. SDC-  
374 GBB successfully handles instances exceeding two hundred thousand samples across all constraint  
375 cases and further scales to ML-only instances with more than 1.5 million samples.

376  
377 <sup>1</sup>Tables 2-3 use a subset of datasets from Table 1 ( $n \leq 210,000$ ). Table 1 includes 15 datasets of up to 1.5M  
378 samples that are only tractable under ML constraints due to the scalability limitations discussed in Section 6.

378 **Small and medium-sized datasets** For small datasets, SDC-GBB matches the state-of-the-art  
 379 PC-SOS-SDP algorithm, outperforming CPLEX and heuristic methods across all constraint settings,  
 380 as shown in Tables 1, 2, 3. On real-world benchmarks, SDC-GBB and PC-SOS-SDP achieve global  
 381 optimality with gaps  $\leq 0.1\%$  on Iris ( $n = 150$ ) (Fisher, 1936) and Seeds ( $n = 210$ ) (Charytanowicz  
 382 et al., 2010). The heuristic method also closely approximates the global optima, whereas CPLEX  
 383 yields significantly larger gaps of around 10%–69% for ML, 73%–86% for CL, and 37%–75% for  
 384 combined constraints. In experiments with medium-sized datasets, SDC-GBB consistently outper-  
 385 forms all other algorithms, achieving optimality gaps  $\leq 0.1\%$  in nearly all cases, with exceptions  
 386 in PR2392 (CL-only) and RDS\_CNT (CL-only), where gaps slightly increase to 2.68% and 1.23%,  
 387 respectively. The best heuristic method consistently performs slightly worse than SDC-GBB but  
 388 remains competitive, providing relatively small gaps across datasets. Conversely, CPLEX returns  
 389 gaps close to 100% across all medium-sized datasets, and PC-SOS-SDP either generates higher gaps  
 390 than SDC-GBB or fails to converge for datasets with  $\geq 2,000$  samples.

391 Table 1: Computational performance with must-link (SDC-GBB,  $k = 3$ ).  
 392

REAL-WORLD DATASETS										
DATASET	METHOD	UB	NODES	GAP (%)	DATASET	METHOD	UB	NODES	GAP (%)	
IRIS <sup>2</sup> N = 150 D = 4	HEURISTIC	83.82	—	—	HTRU2 N = 17,898 D = 8	HEURISTIC	$1.407 \times 10^8$	—	—	
	CPLEX	84.07	12987400	10.55%		CPLEX	NO FEASIBLE SOLUTION FOUND			
	PC-SOS-SDP	<b>83.63</b>	1	$\leq 0.1\%$		PC-SOS-SDP	NO SOLUTION FOUND			
	PARALLEL	<b>83.63</b>	10	$\leq 0.1\%$		PARALLEL	$1.022 \times 10^8$	67	$\leq 0.1\%$	
SEED <sup>2</sup> N = 210 D = 7	HEURISTIC	620.78	—	—	SPNET3D_5 <sup>3</sup> N = 50,000 D = 3	HEURISTIC	$6.627 \times 10^6$	—	—	
	CPLEX	755.91	5814200	69.48%		CPLEX	NO FEASIBLE SOLUTION FOUND			
	PC-SOS-SDP	<b>620.23</b>	1	$\leq 0.1\%$		PC-SOS-SDP	NO SOLUTION FOUND			
	PARALLEL	<b>620.23</b>	7	$\leq 0.1\%$		PARALLEL	$6.609 \times 10^6$	61	<b>1.51%</b>	
HEMI <sup>2</sup> N = 1,955 D = 7	HEURISTIC	$1.602 \times 10^7$	—	—	SKIN_8 N = 80,000 D = 3	HEURISTIC	$4.533 \times 10^8$	—	—	
	CPLEX	$3.204 \times 10^7$	90400	96.26%		CPLEX	NO FEASIBLE SOLUTION FOUND			
	PC-SOS-SDP	$1.601 \times 10^7$	1	2.07%		PC-SOS-SDP	NO SOLUTION FOUND			
	PARALLEL	<b><math>1.601 \times 10^7</math></b>	8	$\leq 0.1\%$		PARALLEL	$4.138 \times 10^8$	12	<b>0.62%</b>	
PR2392 <sup>2</sup> N = 2,392 D = 2	HEURISTIC	$3.210 \times 10^{10}$	—	—	URBANGB N = 360,177 D = 2	HEURISTIC	$1.643 \times 10^9$	—	—	
	CPLEX	$3.816 \times 10^{10}$	281000	98.76%		CPLEX	NO FEASIBLE SOLUTION FOUND			
	PC-SOS-SDP	NO SOLUTION FOUND				PC-SOS-SDP	NO SOLUTION FOUND			
	PARALLEL	<b><math>3.209 \times 10^{10}</math></b>	22	$\leq 0.1\%$		PARALLEL	$4.135 \times 10^8$	2	<b>12.97%</b>	
RDS_CNT <sup>2</sup> N = 10,000 D = 3	HEURISTIC	$6.122 \times 10^7$	—	—	SPNET3D N = 434,874 D = 3	HEURISTIC	$5.848 \times 10^7$	—	—	
	CPLEX	$1.154 \times 10^8$	12600	100.00%		CPLEX	NO FEASIBLE SOLUTION FOUND			
	PC-SOS-SDP	NO SOLUTION FOUND				PC-SOS-SDP	NO SOLUTION FOUND			
	PARALLEL	<b><math>6.078 \times 10^7</math></b>	32	$\leq 0.1\%$		PARALLEL	<b><math>5.797 \times 10^7</math></b>	2	<b>7.43%</b>	
SYNTHETIC DATASETS (D = 2)										
SYN-210000 N = 210,000	HEURISTIC	$2.163 \times 10^6$	—	—	SYN-1050000 N = 1,050,000 D = 2	HEURISTIC	$1.050 \times 10^7$	—	—	
	CPLEX	NO FEASIBLE SOLUTION FOUND				CPLEX	NO FEASIBLE SOLUTION FOUND			
	PC-SOS-SDP	NO SOLUTION FOUND				PC-SOS-SDP	NO SOLUTION FOUND			
	PARALLEL	<b><math>2.163 \times 10^6</math></b>	26	$\leq 0.1\%$		PARALLEL	<b><math>1.050 \times 10^7</math></b>	1	<b>2.45%</b>	
SYN-420000 N = 420,000	HEURISTIC	$6.010 \times 10^6$	—	—	SYN-1500000 N = 1,500,000 D = 3	HEURISTIC	$1.725 \times 10^7$	—	—	
	CPLEX	NO FEASIBLE SOLUTION FOUND				CPLEX	NO FEASIBLE SOLUTION FOUND			
	PC-SOS-SDP	NO SOLUTION FOUND				PC-SOS-SDP	NO SOLUTION FOUND			
	PARALLEL	<b><math>6.010 \times 10^6</math></b>	18	<b>0.61%</b>		PARALLEL	<b><math>1.725 \times 10^7</math></b>	1	<b>2.66%</b>	

414 Table 2: Computational performance with cannot-link (SDC-GBB,  $k = 3$ ).  
 415

REAL-WORLD DATASETS										
DATASET	METHOD	UB	NODES	GAP (%)	DATASET	METHOD	UB	NODES	GAP (%)	
IRIS <sup>2</sup> N = 150 D = 4	HEURISTIC	80.31	—	—	RDS_CNT <sup>2</sup> N = 10,000 D = 3	HEURISTIC	$2.897 \times 10^7$	—	—	
	CPLEX	119.03	8259900	73.59%		CPLEX	$5.198 \times 10^7$	11559	100.00%	
	PC-SOS-SDP	<b>80.21</b>	1	$\leq 0.1\%$		PC-SOS-SDP	NO SOLUTION FOUND			
	PARALLEL	<b>80.21</b>	17	$\leq 0.1\%$		PARALLEL	<b><math>2.861 \times 10^7</math></b>	49	<b>1.23%</b>	
SEED <sup>2</sup> N = 210 D = 7	HEURISTIC	603.04	—	—	HTRU2 <sup>3</sup> N = 17,898 D = 8	HEURISTIC	$1.740 \times 10^8$	—	—	
	CPLEX	771.54	5244683	86.67%		CPLEX	NO FEASIBLE SOLUTION FOUND			
	PC-SOS-SDP	<b>601.96</b>	15	$\leq 0.1\%$		PC-SOS-SDP	NO SOLUTION FOUND			
	PARALLEL	<b>601.96</b>	18	$\leq 0.1\%$		PARALLEL	<b><math>9.225 \times 10^7</math></b>	15	$\leq 0.1\%$	
HEMI <sup>2</sup> N = 1,955 D = 7	HEURISTIC	$1.401 \times 10^7$	—	—	SPNET3D_5 N = 50,000 D = 3	HEURISTIC	$3.938 \times 10^6$	—	—	
	CPLEX	$2.667 \times 10^7$	65002	100.00%		CPLEX	NO FEASIBLE SOLUTION FOUND			
	PC-SOS-SDP	$1.328 \times 10^7$	$i^4$	17.70%		PC-SOS-SDP	NO SOLUTION FOUND			
	PARALLEL	<b><math>1.328 \times 10^7</math></b>	5	$\leq 0.1\%$		PARALLEL	<b><math>3.831 \times 10^6</math></b>	34	$\leq 0.1\%$	
PR2392 <sup>2</sup> N = 2,392 D = 2	HEURISTIC	$2.566 \times 10^{10}$	—	—	SKIN_8 N = 80,000 D = 3	HEURISTIC	$6.367 \times 10^8$	—	—	
	CPLEX	$3.266 \times 10^{10}$	256964	99.96%		CPLEX	NO FEASIBLE SOLUTION FOUND			
	PC-SOS-SDP	NO SOLUTION FOUND				PC-SOS-SDP	NO SOLUTION FOUND			
	PARALLEL	<b><math>2.512 \times 10^{10}</math></b>	97	<b>2.68%</b>		PARALLEL	<b><math>3.016 \times 10^8</math></b>	8	<b>1.32%</b>	
SYNTHETIC DATASETS (D = 2)										
SYN-12000 N = 12,000	HEURISTIC	$9.503 \times 10^4$	—	—	SYN-42000 N = 42,000 D = 2	HEURISTIC	$5.116 \times 10^5$	—	—	
	CPLEX	NO FEASIBLE SOLUTION FOUND				CPLEX	NO FEASIBLE SOLUTION FOUND			
	PC-SOS-SDP	NO SOLUTION FOUND				PC-SOS-SDP	NO SOLUTION FOUND			
	PARALLEL	<b><math>9.053 \times 10^4</math></b>	89	<b>0.75%</b>		PARALLEL	<b><math>5.116 \times 10^5</math></b>	22	$\leq 0.1\%$	
SYN-21000 N = 21,000	HEURISTIC	$1.817 \times 10^5$	—	—	SYN-210000 N = 210,000 D = 3	HEURISTIC	$2.161 \times 10^6$	—	—	
	CPLEX	NO FEASIBLE SOLUTION FOUND				CPLEX	NO FEASIBLE SOLUTION FOUND			
	PC-SOS-SDP	NO SOLUTION FOUND				PC-SOS-SDP	NO SOLUTION FOUND			
	PARALLEL	<b><math>1.817 \times 10^5</math></b>	27	<b>0.32%</b>		PARALLEL	<b><math>2.161 \times 10^6</math></b>	1	<b>2.38%</b>	

**Large and huge datasets** Only SDC-GBB and heuristic algorithms can handle  $n > 10,000$  datasets, with SDC-GBB getting better UB and reaching global optimality or maintaining stable gaps below 3% in all datasets and experiments other than URBANGB and SPNET3D, which receive 5% - 13% gaps. Although the gaps we achieved with these instances are not optimal, they can be further optimized with increased parallelization, and it is worth noting that no other global optimization method can find solutions at this scale. Meanwhile, both PC-SOS-SDP and CPLEX cannot find any feasible solution for these datasets given the 12-hour time limit. This shows that SDC-GBB can scale up to 1,500 times for ML-only constraints, and 200 times for CL-only and ML+CL constraints when compared with state-of-the-art algorithms.

Table 3: Computational performance with both must-link and cannot-link (SDC-GBB,  $k = 3$ ).

REAL-WORLD DATASETS										
DATASET	METHOD	UB	NODES	GAP (%)	DATASET	METHOD	UB	NODES	GAP (%)	
IRIS <sup>2</sup> N = 150 D = 4	HEURISTIC	86.85	—	—	RDS_CNT <sup>2</sup> N = 10,000 D = 3	HEURISTIC	$7.579 \times 10^7$	—	—	
	CPLEX	93.75	9166269	37.47%		CPLEX	$1.493 \times 10^8$	4100	100.00%	
	PC-SOS-SDP	<b>86.76</b>	3	$\leq 0.1\%$		PC-SOS-SDP	NO SOLUTION FOUND			
	PARALLEL	<b>86.76</b>	17	$\leq 0.1\%$		PARALLEL	$7.437 \times 10^7$	32	$\leq 0.1\%$	
SEED <sup>2</sup> N = 210 D = 7	HEURISTIC	597.14	—	—	HTRU2 <sup>3</sup> N = 17,898 D = 8	HEURISTIC	$1.859 \times 10^8$	—	—	
	CPLEX	760.73	6044677	75.32%		CPLEX	NO FEASIBLE SOLUTION FOUND			
	PC-SOS-SDP	<b>596.61</b>	5	$\leq 0.1\%$		PC-SOS-SDP	NO SOLUTION FOUND			
	PARALLEL	<b>596.61</b>	9	$\leq 0.1\%$		PARALLEL	$1.141 \times 10^8$	79	$\leq 0.1\%$	
HEMI <sup>2</sup> N = 1,955 D = 7	HEURISTIC	$1.566 \times 10^7$	—	—	SPNET3D_5 <sup>3</sup> N = 50,000 D = 3	HEURISTIC	$8.171 \times 10^6$	—	—	
	CPLEX	$3.519 \times 10^7$	39600	100.00%		CPLEX	NO FEASIBLE SOLUTION FOUND			
	PC-SOS-SDP	NO SOLUTION FOUND				PC-SOS-SDP	NO SOLUTION FOUND			
	PARALLEL	$1.533 \times 10^7$	115	$\leq 0.1\%$		PARALLEL	$7.925 \times 10^6$	45	<b>5.01%</b>	
PR2392 <sup>2</sup> N = 2,392 D = 2	HEURISTIC	$2.922 \times 10^{10}$	—	—	SKIN_8 N = 80,000 D = 3	HEURISTIC	$7.579 \times 10^8$	—	—	
	CPLEX	$3.240 \times 10^{10}$	239465	99.50%		CPLEX	NO FEASIBLE SOLUTION FOUND			
	PC-SOS-SDP	NO SOLUTION FOUND				PC-SOS-SDP	NO SOLUTION FOUND			
	PARALLEL	$2.916 \times 10^{10}$	35	$\leq 0.1\%$		PARALLEL	$4.258 \times 10^8$	11	<b>4.86%</b>	
SYNTHETIC DATASETS (D = 2)										
SYN-12000 N = 12,000	HEURISTIC	$9.520 \times 10^4$	—	—	SYN-42000 N = 42,000	HEURISTIC	$5.133 \times 10^5$	—	—	
	CPLEX	NO FEASIBLE SOLUTION FOUND				CPLEX	NO FEASIBLE SOLUTION FOUND			
	PC-SOS-SDP	NO SOLUTION FOUND				PC-SOS-SDP	NO SOLUTION FOUND			
	PARALLEL	$9.520 \times 10^4$	25	$\leq 0.1\%$		PARALLEL	$5.133 \times 10^5$	31	<b>1.37%</b>	
SYN-21000 N = 21,000	HEURISTIC	$1.818 \times 10^5$	—	—	SYN-210000 N = 210,000	HEURISTIC	$2.165 \times 10^6$	—	—	
	CPLEX	NO FEASIBLE SOLUTION FOUND				CPLEX	NO FEASIBLE SOLUTION FOUND			
	PC-SOS-SDP	NO SOLUTION FOUND				PC-SOS-SDP	NO SOLUTION FOUND			
	PARALLEL	$1.818 \times 10^5$	37	0.18%		PARALLEL	$2.164 \times 10^6$	28	<b>0.64%</b>	

<sup>2</sup> LESS THAN 4 HOURS.<sup>3</sup> LESS THAN 8 HOURS.<sup>4</sup> SOLVED AT THE ROOT NODE.

## 5 CONCLUSION

In this paper, we presented Sample-Driven Constrained Group-Based Branch-and-Bound (SDC-GBB), a deterministic global optimization algorithm for pairwise-constrained MSSC. We prove convergence to a globally  $\varepsilon$ -optimal solution and demonstrate scalability to datasets exceeding 200,000 samples in all constraint settings, which is **over 200 times larger** than the 800-sample benchmark of (Piccialli et al., 2022a), and further extend to **over 1,500 times larger** with ML-only instances having more than one million samples while maintaining optimality gaps below 3%. When empirically evaluated on real-world benchmarks of various domains, SDC-GBB consistently achieves low optimality gaps across diverse constraint settings.

## 6 LIMITATIONS AND FUTURE DIRECTIONS

Similar to prior work, our algorithm struggles to scale to one million samples with CL constraints due to the NP-hard nature of this constraint type. Future work may consider tightening the grouped-sample Lagrangian lower bound (Section 3.4) by dualizing the CL graph with relaxing binary indicators to  $[0,1]$  and penalty multipliers  $z_{ij}$  updated via subgradient methods. This formulation would produce a more compact MISOCP with fewer active binary variables thanks to the relaxation of  $z_{ij}$  and stronger continuous bounds, reducing the number of branches in dense CL graphs, as shown by successful Lagrangian-penalty approaches in semi-supervised clustering and global MISOCP strategies. An alternative is to incorporate clique inequalities or separation cuts for the CL graph, following MIP methodologies that substantially improve bounds. However, including these Lagrangian terms at each node would increase the solve time of the relaxation, so empirically evaluating the trade-off between bound improvement and per-node cost will be key.

486     **Ethics Statement** Our study uses the SKIN (Skin Segmentation) dataset from the UCI Machine  
 487     Learning Repository. The dataset contains de-identified RGB pixel triplets (B, G, R) sampled  
 488     from facial images and is derived from two collections: the PAL Face Database and the DARPA-  
 489     sponsored Color FERET images. All human participants provided consent for the use of these  
 490     collections for research purposes. In line with data-privacy best practices, the UCI release exposes only  
 491     anonymized pixel triplets and does not include raw images, which reduces the risk of re-identification  
 492     in downstream work. Concurrently, our framework optimizes only the Minimum Sum-of-Squares  
 493     Criteria (MSSC) objective, and as is well documented, MSSC/k-means-style clustering can reproduce  
 494     and even amplify existing biases in the data, particularly at scale. In high-stakes or sensitive domains,  
 495     deployments that do not account for these effects can lead to disparate treatment or outcomes across  
 496     demographic groups. We recommend that practitioners perform fairness audits before deployment,  
 497     and consider mitigation techniques such as fair clustering variants and post-processing adjustments if  
 498     such disparities arise. The absence of fairness-aware safeguards in the current implementation is a  
 499     limitation of this work, and integrating such constraints or corrections is left for future extensions.  
 500

500     **Reproducibility Statement** In this section, we outline necessary details for reproducing all experi-  
 501     ments described in the paper. We provide comprehensive hardware and software configuration for  
 502     SDC-GBB as well as pairwise constraint generation, seeding protocol, runtime limit and data setup,  
 503     including real and synthetic dataset generation in Section 4. The implementation of the SDC-GBB  
 504     algorithm is thoroughly described through Algorithm 1. Lastly, we document evaluation in Section 4  
 505     along with the results of all baselines used for comparison in Sections D and E. To further support  
 506     the reproducibility of our results, we will release our experiment code upon acceptance, enabling  
 507     other researchers to replicate and expand on our work.  
 508

## 509     REFERENCES

510     Daniel Aloise, Amit Deshpande, Pierre Hansen, and Preyas Popat. Np-hardness of euclidean sum-of-  
 511     squares clustering. *Machine Learning*, 75(2):245–248, 2009.  
 512  
 513     Daniel Aloise, Pierre Hansen, and Leo Liberti. An improved column generation algorithm for  
 514     minimum sum-of-squares clustering. *Mathematical Programming*, 131(1):195–220, 2012a.  
 515  
 516     Daniel Aloise, Pierre Hansen, and Caroline Rocha. A column generation algorithm for semi-  
 517     supervised minimum sum-of-squares clustering. In *Global Optimization Workshop 2012*, pp.  
 518     19–22, 2012b.  
 519  
 520     Behrouz Babaki, Tias Guns, and Siegfried Nijssen. Constrained clustering using column generation.  
 521     In Helmut Simonis (ed.), *Integration of AI and OR Techniques in Constraint Programming*, pp.  
 522     438–454, Cham, 2014. Springer International Publishing.  
 523  
 524     Sugato Basu, Arindam Banerjee, and Raymond J. Mooney. Active semi-supervision for pairwise  
 525     constrained clustering. In *Proceedings of the 2004 SIAM International Conference on Data Mining*  
 526     (SDM), pp. 333–344. Society for Industrial and Applied Mathematics, 2004.  
 527  
 528     Sugato Basu, Ian Davidson, and Kiri Wagstaff (eds.). *Constrained Clustering: Advances in Algo-  
 529     rithms, Theory, and Applications*. Chapman and Hall/CRC, New York, n.d. edition, 2008.  
 530  
 531     Philipp Baumann. A binary linear programming-based k-means algorithm for clustering with must-  
 532     link and cannot-link constraints. In *2020 IEEE international conference on industrial engineering*  
 533     and *engineering management (IEEM)*, pp. 324–328. IEEE, 2020.  
 534  
 535     Philipp Baumann and Dorit S. Hochbaum. An algorithm for clustering with confidence-based  
 536     must-link and cannot-link constraints. *INFORMS Journal on Computing*, 34(1), 2024.  
 537  
 538     Dimitris Bertsimas and Jack Dunn. Optimal classification trees. *Springer US*, 2017.  
 539     Andreas Brieden, Peter Gritzmann, and Fabian Klemm. Constrained clustering via diagrams: A  
 540     unified theory and its application to electoral district design. *European Journal of Operational  
 541     Research*, 263(1):18–34, 2017.  
 542  
 543     Michael J. Brusco. A repetitive branch-and-bound procedure for minimum within-cluster sums of  
 544     squares partitioning. *Psychometrika*, 71(2):347–363, 2006.

540 Yankai Cao and Victor M. Zavala. A scalable global optimization algorithm for stochastic nonlinear  
 541 programs. *J. of Global Optimization*, 75(2):393–416, 2019.  
 542

543 Małgorzata Charytanowicz, Jerzy Niewczas, Piotr Kulczycki, Piotr A. Kowalski, Szymon Łukasik,  
 544 and Sławomir Źak. Complete gradient clustering algorithm for features analysis of x-ray images.  
 545 In *Information Technologies in Biomedicine*, pp. 15–24, Berlin, Heidelberg, 2010. Springer.

546 Cplex. User’s manual for cplex. Technical report, International Business Machines Corporation, New  
 547 York, NY, 2022.  
 548

549 Thi-Bich-Hanh Dao, Khanh-Chuong Duong, and Christel Vrain. A declarative framework for  
 550 constrained clustering. In Hendrik Blockeel, Kristian Kersting, Siegfried Nijssen, and Filip  
 551 Źelezný (eds.), *Machine Learning and Knowledge Discovery in Databases*, pp. 419–434, Berlin,  
 552 Heidelberg, 2013. Springer.

553 Thi-Bich-Hanh Dao, Khanh-Chuong Duong, and Christel Vrain. Constrained minimum sum of  
 554 squares clustering by constraint programming. In Gilles Pesant (ed.), *Principles and Practice of  
 555 Constraint Programming*, pp. 557–573, Cham, 2015. Springer International Publishing.  
 556

557 Thi-Bich-Hanh Dao, Khanh-Chuong Duong, and Christel Vrain. Constrained clustering by constraint  
 558 programming. *Artificial Intelligence*, 244:70–94, 2017.  
 559

560 Ian Davidson and S. S. Ravi. Clustering with constraints: Feasibility issues and the k-means algorithm.  
 561 In *Proceedings of the 2005 SIAM International Conference on Data Mining (SDM)*, pp. 138–149.  
 562 Society for Industrial and Applied Mathematics, 2005.  
 563

563 Dheeru Dua and Casey Graff. Uci machine learning repository, 2017. URL <http://archive.ics.uci.edu/ml>.  
 564

565 Dan Feldman and Michael Langberg. A unified framework for approximating and clustering data.  
 566 In *Proceedings of the forty-third annual ACM symposium on Theory of computing*, pp. 569–578,  
 567 2011.  
 568

569 R. A. Fisher. The use of multiple measurements in taxonomic problems. *Annals of Eugenics*, 7(2):  
 570 179–188, 1936.  
 571

571 Tias Guns, Thi-Bich-Hanh Dao, Christel Vrain, and Khanh-Chuong Duong. Repetitive branch-and-  
 572 bound using constraint programming for constrained minimum sum-of-squares clustering. In  
 573 *Proceedings of the Twenty-second European Conference on Artificial Intelligence*, pp. 462–470,  
 574 NLD, 2016. IOS Press.  
 575

576 Gurobi. Gurobi Optimizer Reference Manual, 2024. URL <https://www.gurobi.com>.  
 577

577 Reiner Horst and Hoang Tuy. *Global optimization: Deterministic approaches*. Springer Science &  
 578 Business Media, 2013.  
 579

580 Kaixun Hua, Mingfei Shi, and Yankai Cao. A scalable deterministic global optimization algorithm for  
 581 clustering problems. In *International Conference on Machine Learning*, pp. 4391–4401. PMLR,  
 582 2021.  
 583

583 Haichao Huang, Yong Cheng, and Ruilian Zhao. A semi-supervised clustering algorithm based on  
 584 must-link set. In Changjie Tang, Charles X. Ling, Xiaofang Zhou, Nick J. Cercone, and Xue Li  
 585 (eds.), *Advanced Data Mining and Applications*, pp. 492–499, Berlin, Heidelberg, 2008. Springer.  
 586

587 A. K. Jain, M. N. Murty, and P. J. Flynn. Data clustering: a review. *ACM Computing Surveys*, 31(3):  
 588 264–323, 1999.  
 589

589 Anil K. Jain. Data clustering: 50 years beyond k-means. *Pattern Recognition Letters*, 31(8):651–666,  
 590 2010.  
 591

592 Ramkumar Karuppiah and Ignacio E. Grossmann. A lagrangean-based branch-and-cut algorithm for  
 593 global optimization of nonconvex mixed-integer nonlinear programs with decomposable structures.  
*Journal of Global Optimization*, 41(2):163–186, 2008.

594 Leo Liberti and Benedetto Manca. Side-constrained minimum sum-of-squares clustering: mathematical  
 595 programming and random projections. *Journal of Global Optimization*, 83(1):83–118,  
 596 2022.

597 Nguyen-Viet-Dung Nghiem, Christel Vrain, Thi-Bich-Hanh Dao, and Ian Davidson. Constrained  
 598 clustering via post-processing. In *Discovery Science: 23rd International Conference, DS 2020,  
 599 Thessaloniki, Greece, October 19–21, 2020, Proceedings*, pp. 53–67, Berlin, Heidelberg, 2020.  
 600 Springer-Verlag.

602 Manfred Padberg and Giovanni Rinaldi. A branch-and-cut algorithm for the resolution of large-scale  
 603 symmetric traveling salesman problems. *SIAM Review*, 33(1):60–100, 1991.

604

605 Mercedes Pelegrín. New variants of the simple plant location problem and applications. *European  
 606 Journal of Operational Research*, 306(3):1094–1108, 2023.

607 Jiming Peng and Yu Xia. A cutting algorithm for the minimum sum-of-squared error clustering. In  
 608 *Proceedings of the 2005 SIAM International Conference on Data Mining (SDM)*, pp. 150–160.  
 609 Society for Industrial and Applied Mathematics, 2005.

610

611 Veronica Piccialli, Anna Russo Russo, and Antonio M. Sudoso. An exact algorithm for semi-  
 612 supervised minimum sum-of-squares clustering. *Computers & Operations Research*, 147, 2022a.

613

614 Veronica Piccialli, Antonio M. Sudoso, and Angelika Wiegele. Sos-sdp: An exact solver for minimum  
 615 sum-of-squares clustering. *INFORMS Journal on Computing*, 34(4):2144–2162, 2022b.

616

617 M. R. Rao. Cluster analysis and mathematical programming. *Journal of the American Statistical  
 618 Association*, 66(335):622–626, 1971.

619

620 Tonny Rutayisire, Yan Yang, Chao Lin, and Jinyuan Zhang. A modified cop-kmeans algorithm based  
 621 on sequenced cannot-link set. In JingTao Yao, Sheela Ramanna, Guoyin Wang, and Zbigniew Suraj  
 (eds.), *Rough Sets and Knowledge Technology*, pp. 217–225, Berlin, Heidelberg, 2011. Springer.

622

623 Chris Schwiegelshohn and Omar Ali Sheikh-Omar. An empirical evaluation of  $k$ -means coresets.  
 624 *arXiv preprint arXiv:2207.00966*, 2022.

625

626 Helmut Späth. *Cluster Analysis Algorithms for Data Reduction and Classification of Objects*. E.  
 627 Horwood, n.d., n.d. edition, 1980.

628

629 Wei Tan, Yan Yang, and Tianrui Li. An improved cop-kmeans algorithm for solving constraint  
 630 violation. In *Computational Intelligence: Foundations and Applications*, pp. 690–696. World  
 Scientific, 2010.

631

632 Wei Tang, Yang Yang, Lanling Zeng, and Yongzhao Zhan. Size constrained clustering with milp  
 633 formulation. *IEEE Access*, 8:1587–1599, 2020.

634

635 Tian Tian, Jie Zhang, Xiang Lin, Zhi Wei, and Hakon Hakonarson. Model-based deep embedding for  
 636 constrained clustering analysis of single cell rna-seq data. *Nature Communications*, 12(1):1873,  
 2021.

637

638 Kiri Wagstaff, Claire Cardie, Seth Rogers, and Stefan Schrödl. Constrained  $k$ -means clustering with  
 639 background knowledge. In *Proceedings of the Eighteenth International Conference on Machine  
 640 Learning*, pp. 577–584, San Francisco, CA, USA, 2001. Morgan Kaufmann Publishers Inc.

641

642 Edward Wang, Riley Ballachay, Genpei Cai, Yankai Cao, and Heather L. Trajano. Predicting xylose  
 643 yield from prehydrolysis of hardwoods: A machine learning approach. *Frontiers in Chemical  
 Engineering*, 4, 2022.

644

645 Yu Xia. A global optimization method for semi-supervised clustering. *Data Mining and Knowledge  
 646 Discovery*, 18(2):214–256, 2009.

647

648 Jason Xu and Kenneth Lange. Power  $k$ -means clustering. In *Proceedings of the 36th International  
 649 Conference on Machine Learning*, pp. 6921–6931, Berlin, Heidelberg, 2019. PMLR.

648 Yi Yang, Kunpeng Zhang, and Yangyang Fan. Analyzing firm reports for volatility prediction: A  
649 knowledge-driven text-embedding approach. *INFORMS Journal on Computing*, 34(1):522–540,  
650 2022.

651

652 Ying Zhang, Xiangli Li, and Mengxue Jia. Semi-supervised nonnegative matrix factorization  
653 with pairwise constraints for image clustering. *International Journal of Machine Learning and*  
654 *Cybernetics*, 13(11):3577–3587, 2022.

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702    **A TWO-STAGE PROGRAMS REFORMULATION**  
 703

704  
 705    Once the must-link set  $\mathcal{T}_{ml}$  has been collapsed into pseudo-samples (Sec. 3.2), the additive constant  
 706     $(t - 1) \text{tr}(\Sigma_{ml}^2)$  can be pre-computed. Hence minimising the objective in (4a) is equivalent to  
 707    minimising  $\sum_{s \in \hat{\mathcal{S}}} d_{s,*}$ . The optimal solution of (4) is obtained from the two-stage program:

708  
 709    
$$z = \min_{\mu \in M_0} \sum_{s \in \hat{\mathcal{S}}} Q_s(\mu), \quad (8)$$
  
 710

711  
 712    where  $\mu$  are the *first-stage* variables and  $Q_s(\mu)$  is the optimal value of the *second-stage* problem  
 713    defined below. After (i) collapsing must-link components and (ii) applying the geometric Lemmas  
 714    1–2 together with cannot-link propagation, each sample  $s$  may still be assigned only to a *viable* subset  
 715     $\mathcal{K}_s \subseteq \mathcal{K}$ . The reduced dataset  $\hat{\mathcal{S}}$  and the family  $\{\mathcal{K}_s\}_{s \in \hat{\mathcal{S}}}$  are fixed once at the root node.

716  
 717    
$$Q_s(\mu) = \min_{d_s, b_s} d_{s,*} \quad (9)$$
  
 718    s.t. Constraints 2b, 2c, 1e, 1d.

721  
 722    Here  $d_s = [d_{s,k}]_{k \in \mathcal{K}_s}$  and  $b_s = [b_{s,k}]_{k \in \mathcal{K}_s}$  are the *second-stage* variables. The closed set  $M_0 = \{\mu \mid$   
 723     $\mu^l \leq \mu \leq \mu^u\}$  bounds every centre with  $\mu_{k,i}^l = \min_s X_{s,i}$  and  $\mu_{k,i}^u = \max_s X_{s,i}$  for all  $k \in \mathcal{K}$  and  
 724     $i = 1, \dots, m$ . For convenience we choose a *single* Big- $M$  constant

725  
 726    
$$N = \max_{s,k} \sum_{i=1}^m \max\{|x_{s,i} - \mu_{k,i}^l|^2, |x_{s,i} - \mu_{k,i}^u|^2\},$$
  
 727

729  
 730    which leaves all bounds valid and simplifies the notation. The bounds  $\mu^l, \mu^u$  are computed once at  
 731    the root node and inherited unchanged by every BB subproblem. We denote by  $\text{relint}(\mathcal{M})$  and  $\delta(\mathcal{M})$   
 732    the relative interior and the diameter of a set. Throughout this paper, the diameter of the box set  $M_0$   
 733    is  $\delta(M_0) = \|\mu^u - \mu^l\|_\infty$ .

734    It can be shown that the closed-form solution to the second-stage problem is  
 735

736    
$$Q_s(\mu) = \min_{k \in \mathcal{K}_s} \|x_s - \mu_k\|_2^2.$$
  
 737

738  
 739    Since  $Q_s(\mu)$  is the minimum of a finite number of continuous functions,  $Q_s$  is continuous. Because  
 740    of the compactness of  $M_0$  and continuity of  $Q_s(\mu)$ , the clustering Problem 8 can attain its minimum  
 741    according to the generalized Weierstrass theorem.

742    When the BB algorithm explores a box  $M \subseteq M_0$ , it solves the *primal node problem*  
 743

744  
 745    
$$z(M) = \min_{\mu \in M} \sum_{s \in \hat{\mathcal{S}}} Q_s(\mu). \quad (10)$$
  
 746

747  
 748    Replicating centres for each sample and enforcing non-anticipativity (11b) yields the lifted form  
 749

750    
$$\min_{\mu_s \in M} \sum_{s \in \hat{\mathcal{S}}} Q_s(\mu_s) \quad (11a)$$
  
 751

752    s.t.  $\mu_s = \mu_{s+1}, \quad s = 1, \dots, |\hat{\mathcal{S}}| - 1. \quad (11b)$   
 753

754  
 755    Problems (10) and (11) are equivalent, and retain all cannot-link information through the viable-cluster  
 sets  $\{\mathcal{K}_s\}_{s \in \hat{\mathcal{S}}}$ .

756 **B PROOF OF THEOREMS**  
757758 **B.1 PROOF OF LEMMA 3.3**  
759760 *Proof.*

761 
$$\text{sse}_{\mathcal{C}}(\mu) = \sum_{i=1}^p \|x_i - \mu\|^2 \quad (12)$$
  
762

763 
$$= \sum_{i=1}^t \|x_i - \mu\|^2 + \sum_{i=t+1}^p \|x_i - \mu\|^2 \quad (13)$$
  
764

765 
$$= \sum_{i=1}^t \|x_i\|^2 - 2\mu^T \sum_{i=1}^t x_i + t\|\mu\|^2 + \sum_{i=t+1}^p \|x_i - \mu\|^2 \quad (14)$$
  
766

767 Here  $\sum_{i=1}^t \|x_i\|^2$  can be rewritten as follow:  
768

769 
$$\sum_{i=1}^t \|x_i\|^2 = \sum_{i=1}^t \|x_i - \mu_{ml} + \mu_{ml}\|^2 \quad (15)$$
  
770

771 
$$= \sum_{i=1}^t \|x_i - \mu_{ml}\|^2 - 2\mu_{ml}^T \sum_{i=1}^t (x_i - \mu_{ml}) + t\|\mu_{ml}\|^2 \quad (16)$$
  
772

773 
$$= \sum_{i=1}^t \|x_i - \mu_{ml}\|^2 + t\|\mu_{ml}\|^2 \quad (17)$$
  
774

775 
$$= (t-1)\text{tr}(\Sigma^2) + t\|\mu_{ml}\|^2 \quad (18)$$
  
776

777 Thus, we have:  
778

779 
$$\text{sse}_{\mathcal{C}}(\mu) = \sum_{i=1}^t \|x_i\|^2 - 2\mu^T \sum_{i=1}^t x_i + t\|\mu\|^2 + \sum_{i=t+1}^p \|x_i - \mu\|^2 \quad (19)$$
  
780

781 
$$= (t-1)\text{tr}(\Sigma_{ml}^2) + t\|\mu_{ml}\|^2 - 2t\mu^T \mu_{ml} + t\|\mu\|^2 + \sum_{i=t+1}^p \|x_i - \mu\|^2 \quad (20)$$
  
782

783 
$$= (t-1)\text{tr}(\Sigma_{ml}^2) + t\|\mu_{ml} - \mu\|^2 + \sum_{i=t+1}^p \|x_i - \mu\|^2 \quad (21)$$
  
784

785 
$$= (t-1)\text{tr}(\Sigma_{ml}^2) + \text{sse}_{\hat{\mathcal{C}}}(\mu) \quad (22)$$
  
786

787  $\square$   
788789 **B.2 PROOF OF THEOREM 3.4**  
790791 *Proof.* ( $\Rightarrow$ ) Let  $\mu^*$  denote a globally optimal solution of Problem (4) as a result of ML collapse. Each pseudo-sample then corresponds to exactly one must-link component of the original dataset. In constructing Problem (4), assign every genuine sample in that component to the cluster of its associated pseudo-sample; any sample not belonging to a must-link component retains the label it receives in the unconstrained solution. This enforces  $b_{i,k} = b_{i',k}$  for all  $(i, i') \in \mathcal{T}_{ml}$ , so the pair  $(\mu^*, b^*)$  is feasible for Problem (2). At the same time, the resulting objective matches that of Problem (4), since the additive term  $\sum_{k \in \mathcal{K}_{ml}} (t_k - 1) \text{tr}(\Sigma_{ml,k}^2)$  precisely reinstates the variance eliminated when collapsing each must-link component. Hence, the optimal solution of Problem (4) is optimal for Problem (2).792 ( $\Leftarrow$ ) Let  $(\mu^*, b^*)$  be a global optimum for Problem (2) on the original instance. Collapse ML  
793 constraints based on Lemma 2 to form the pseudo-sample instance  $\hat{\mu}$ . Then  $\hat{\mu}$  is feasible for  
794 the unconstrained Problem (4) and so is the pair  $(\hat{\mu}, \hat{b})$ . By contradiction, assume that  $\hat{\mu}$  is not  
795 optimal for Problem (4). Then there exists a feasible  $\mu^\dagger$  for (4) with  $\text{sse}_{\mathcal{C}}(\mu^\dagger) < \text{sse}_{\mathcal{C}}(\hat{\mu})$ . Via

reconstructing the ML constrained Problem (2), we obtain  $sse_{\hat{\mathcal{C}}}(\mu^\dagger, b^\dagger) = sse_{\mathcal{C}}(\mu^\dagger) + C < sse_{\mathcal{C}}(\hat{\mu}) + C = sse_{\hat{\mathcal{C}}}(\mu^*, b^*)$ , which contradicts the optimality of  $(\mu^*, b^*)$ . Hence, the optimal solution for Problem (2) is optimal for Problem (4).  $\square$

## C CONVERGENCE ANALYSIS

In this section we establish the convergence of the proposed BB scheme, constructed with the grouping-based Lagrangian decomposition lower bound and the decomposable upper bound based on closed form solutions for must-link or K-coloring. A key feature of our algorithm is that it **branches exclusively in the space of first-stage variables  $\mu$  to guarantee convergence**. As all must-link components have been collapsed and every CL-infeasible assignment eliminated, the remaining problem is an *unconstrained* optimization over  $\mu$  with continuous objective  $Q(\mu) = \sum_{s \in \hat{\mathcal{S}}} Q_s(\mu)$ .

Therefore, the proof of convergence can easily adopt the foundational results of (Cao & Zavala, 2019) and the seminal contributions in Chapter IV of (Horst & Tuy, 2013). Although the original pairwise-constrained MSSC places additional feasibility requirements on the assignment variables, our *equivalent unconstrained* re-formulation—obtained by first collapsing every must-link component (Theorem 3.4) and then discarding all assignments that violate cannot-link constraints (See Lemma 3.1 and 3.2)—allows *any* point  $\mu \in M$  to be treated as a feasible first-stage decision. The proof of Theorem 3.5 is thus becoming obvious with the definitions and theoretical frameworks of Cao & Zavala (2019), while only notational adaptations will be processed to reflect the reduced dataset  $\hat{\mathcal{S}}$  and the viable-cluster sets  $\{\mathcal{K}_s\}_{s \in \hat{\mathcal{S}}}$  specific to the present problem.

**Lemma C.1** (Lower Bounding Consistency). *Given an exhaustive subdivision (See Definition IV.10 (Horst & Tuy, 2013)) on  $\mu$ , the lower-bounding operation in Algorithm 1 is strongly consistent (See Definition IV.7 (Horst & Tuy, 2013)).*

*Proof.* With an exhaustive subdivision, each box  $M_{i_q}$  shrinks to a single point  $\bar{\mu}$ , so  $\bar{M} = \{\bar{\mu}\}$ . We prove that  $\lim_{q \rightarrow \infty} \beta(M_{i_q}) = z(\bar{M}) = \sum_{s \in \hat{\mathcal{S}}} Q_s(\bar{\mu})$ . Define, for every sample  $s$ ,  $\tilde{\mu}_{i_q, s} \in \arg \min_{\mu \in M_{i_q}} \min_{k \in \mathcal{K}_s} \|x_s - \mu_k\|_2^2$ , where  $\mathcal{K}_s$  is the set of clusters still admissible for  $s$  after the cannot-link pruning. **Because each  $\mathcal{K}_s$  already excludes every cannot-link pairing, every distance minimized in the definition of  $\tilde{\mu}_{i_q, s}$  automatically respects all CL constraints.** Since  $M_{i_q} \rightarrow \{\bar{\mu}\}$ , we have  $\tilde{\mu}_{i_q, s} \rightarrow \bar{\mu}$ . Using the continuity of  $Q_s(\cdot)$ , it follows that  $Q_s(\bar{\mu}) = \lim_{q \rightarrow \infty} Q_s(\tilde{\mu}_{i_q, s}) = \lim_{q \rightarrow \infty} \beta_s(M_{i_q})$ . Summing over all  $s$  yields  $\lim_{q \rightarrow \infty} \beta(M_{i_q}) = \sum_{s \in \hat{\mathcal{S}}} Q_s(\bar{\mu})$ . Proof complete.  $\square$

**Lemma C.2** (Lower Bounding Convergence). *Given an exhaustive subdivision (Definition IV.10 (Horst & Tuy, 2013)) on  $\mu$ , Algorithm 1 satisfies  $\lim_{i \rightarrow \infty} \beta_i = z$ .*

*Proof.* This result can be obtained from Lemma C.1 and Theorem IV.3 of (Horst & Tuy, 2013).  $\square$

**Lemma C.3** (Upper Bounding Convergence). *Given an exhaustive subdivision (Definition IV.10 (Horst & Tuy, 2013)) on the centroid space  $\mu$ , Algorithm 1 produces a sequence  $\{\alpha_i\}$  that satisfies  $\lim_{i \rightarrow \infty} \alpha_i = z$ .*

*Proof.* Let  $\mu^* \in M_0$  be an optimal centroid set for the *equivalent unconstrained* MSSC obtained after **collapsing must-link components and discarding every cannot-link-infeasible assignment**. According to Lemma 6 in Cao & Zavala (2019),  $\lim_{i \rightarrow \infty} \alpha_i = z$  holds when executing Algorithm 1.  $\square$

Combing Lemma C.2 and C.3, we obtain Theorem 3.5. Essentially, pseudo-samples in Section 3.2 yield an equivalent unconstrained MSSC problem, and Theorem 3.4 proves a bijection between optimal solutions before and after this transformation. The geometric rules in Section 3.1 apply dominance checks via Lemmas 3.1 and 3.2, excluding only assignments whose cost significantly exceeds the incumbent upper bound without removing any centroid regions. Consequently, our branch-and-bound algorithm still exhaustively subdivides the centroid space, satisfying Lemmas C.1 through C.3. Thus, the result remains valid for the entire SDC-GBB pipeline. Empirically, we show that within a fixed 12-hour runtime limit, SDC-GBB achieves optimality gaps above 0.1% for some datasets under certain constraints. However, this does not reflect the failure of convergence in our

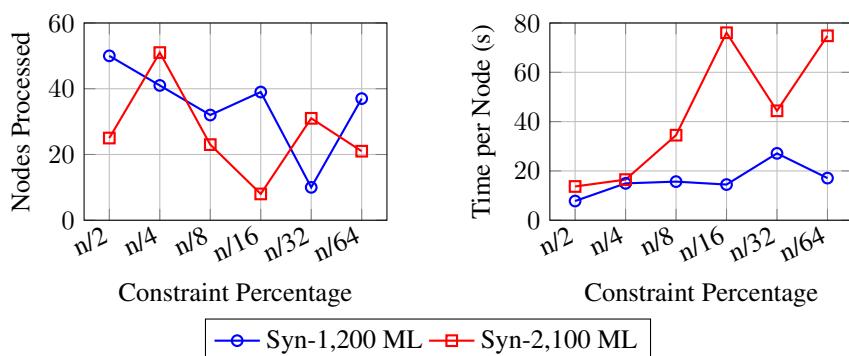
864 global optimization scheme, but rather the best feasible solution and its lower bound at timeout.  
 865 Imposing a time limit is the standard which we follow to ensure a fair comparison with other exact  
 866 baselines.  
 867

## 868 D EFFECT OF PAIRWISE CONSTRAINTS

871 Background knowledge in constrained clustering is introduced through pairwise constraints: must-  
 872 link (ML) and cannot-link (CL). This section analyzes how varying densities of these constraints  
 873 affect the branch-and-bound performance on medium-sized problem instances. Table 4 summarizes  
 874 the number of nodes processed and the average computational time per node under different constraint  
 875 densities. All experiments reported in Table 4 achieve a final optimality gap below or equal to 0.1%.

876 Must-link constraints merge linked samples into single pseudo-points before initiating the branch-and-  
 877 bound procedure. Conceptually, this operation is analogous to samples merging immediately upon  
 878 defining constraints, effectively knowing in advance that they must converge into a single optimal  
 879 position. Figure 2 illustrates this merging process: each must-link pair collapses into one pseudo-point,  
 880 thus reducing the number of samples to consider. Although this merging simplifies the optimization  
 881 search space by reducing dimensionality, it simultaneously imposes additional equality constraints in  
 882 each node relaxation within the branch-and-bound process, thereby increasing the computational effort  
 883 per node relaxation. Nonetheless, overall node processing becomes faster since fewer distinct samples  
 884 remain active, which accelerates bound computations without compromising the equivalence and  
 885 optimality of the final solution. Empirically, the average time per node decreases with an increasing  
 886 number of must-link constraints (from 74 seconds per node down to 13 seconds per node for dataset  
 887 *Syn-2100*). However, there is a threshold for the density of constraints necessary to significantly  
 888 simplify the branching process: below approximately  $n/64$  must-link pairs, constraints have minimal  
 889 practical relevance, causing the equivalent problem formulation, described in Section 3.2, to behave  
 890 similarly to its unconstrained counterpart.

891 In contrast, the geometric sample-determination strategy for cannot-link constraints functions as  
 892 barriers that restrict feasible assignments, similar to placing walls within an axis-aligned region.  
 893 Unlike must-link constraints, which collapse samples proactively, cannot-link constraints compel  
 894 each sample’s assignment to navigate around imposed boundaries that are not initially evident. Each  
 895 sample seeks to reach its optimal cluster centroid but must repeatedly avoid these geometric barriers,  
 896 reflecting constrained assignments and generating additional branching iterations. Despite this growth  
 897 in node count, the computational time per node remains stable because infeasible assignments are  
 898 pruned at an early stage, as stated in Lemma 3.1. Thus, the complexity introduced by cannot-link  
 899 constraints primarily affects the extent of branching rather than the computational complexity of each  
 900 bound evaluation. In summary, must-link constraints simplify the optimization upfront by reducing  
 901 per-node complexity through component collapse, whereas cannot-link constraints expand the search  
 902 tree but maintain per-node computational cost, ensuring that solutions satisfy the imposed constraints  
 903 without fundamentally altering the underlying clustering structure.



915 Figure 3: Effect of Must-Link Constraints on (a) the number of nodes processed and (b) time per  
 916 node, for both synthetic datasets.  
 917

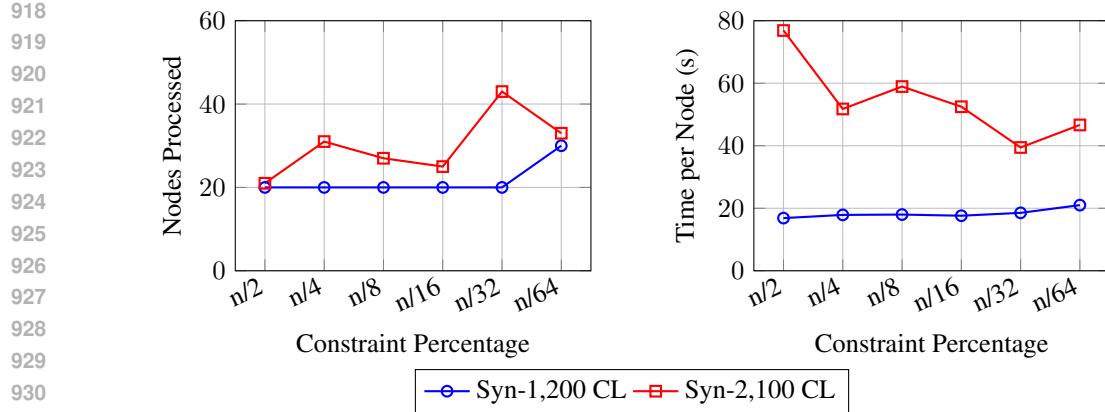


Figure 4: Effect of Cannot-Link Constraints on (a) the number of nodes processed and (b) time per node, for both synthetic datasets.

Interestingly, we observe that the linear relaxation is substantially weakened when must-link constraints form large or overlapping superpoints. In URBANGB, the merging of 469 groups into  $k = 3$  superpoints induces symmetry that leaves a 12.97% gap, and in SPNET3D the near one-dimensional data arrangement produces overlapping superpoints that the bound cannot improve, resulting in a 7.43% gap. On the other hand, in SKIN\_8 and SPNET3D\_5, the mixed ML+CL case is uniquely hard because ML contraction creates component-level “hotspots” that are densely entangled by cannot-links. Inside SKIN\_8, there is a massive ML super-component that must occupy one cluster and is CL-forbidden from sharing with thousands of neighbors, whereas the tight CL and loose ML constraints simultaneously oppose the geometry of SPNET3D\_5. This coupling weakens the relaxation and drives deep branching, yielding a large optimality gap  $> 3\%$  in the mixed constraint case, whereas ML-only or CL-only avoid this interaction and remain tight.

Table 4: Runtime metrics with Different Constraint Settings with solution optimality gap  $< 0.1\%$ .

DATASET	METRIC	$\frac{n}{2}$	$\frac{n}{4}$	$\frac{n}{8}$	$\frac{n}{16}$	$\frac{n}{32}$	$\frac{n}{64}$
<b>MUST-LINK (ML)</b>							
<i>Syn-1,200</i>	NODES	50	41	32	39	10	37
	GAP (%)	<0.1%	<0.1%	<0.1%	<0.1%	<0.1%	<0.1%
	TIME (S)	389.61	613.09	501.05	564.41	271.24	632.25
	TIME/NODE (s)	7.79	14.95	15.66	14.47	27.12	17.09
	CORE HOUR (H)	7.79	14.95	15.66	14.47	27.12	17.09
<i>Syn-2,100</i>	NODES	25	51	23	8	31	21
	GAP (%)	<0.1%	<0.1%	<0.1%	<0.1%	<0.1%	<0.1%
	TIME (S)	342.20	840.70	793.19	608.32	1375.86	1570.95
	TIME/NODE (s)	13.69	16.48	34.49	76.04	44.38	74.81
	CORE HOUR (H)	5.41	8.52	6.96	7.84	3.77	8.78
<b>CANNOT-LINK (CL)</b>							
<i>Syn-1,200</i>	NODES	20	20	20	20	20	30
	GAP (%)	<0.1%	<0.1%	<0.1%	<0.1%	<0.1%	<0.1%
	TIME (S)	336.98	357.45	359.61	352.23	370.58	629.90
	TIME/NODE (s)	16.85	17.87	17.98	17.61	18.53	21.00
	CORE HOUR (H)	7.79	14.95	15.66	14.47	27.12	17.09
<i>Syn-2,100</i>	NODES	21	31	27	25	43	33
	GAP (%)	<0.1%	<0.1%	<0.1%	<0.1%	<0.1%	<0.1%
	TIME (S)	1614.36	1604.81	1591.37	1312.08	1695.81	1539.55
	TIME/NODE (s)	76.87	51.77	58.94	52.48	39.44	46.65
	CORE HOUR (H)	7.79	14.95	15.66	14.47	27.12	17.09

To further highlight resource usage with a unified measure of parallelism and wall-clock time, we compute “core-hour”, defined as  $\text{core\_hour} = \frac{\text{time}(s) \times \text{cores}}{3600(s)}$ . Under must-link constraints, SDC-GBB’s overall branching behavior varies substantially with constraint density, as on Syn-1,200 and Syn-2,100, going from  $n/2$  to  $n/32$  ML constraints generally reduces the number of branch-and-bound nodes but raises the average time per node. In contrast, the number of nodes explored grows roughly linearly under uniform cannot-link constraint placement; as we observe in Table 4, node count scales  $\sim O(m)$  with  $m$  cannot-link constraints, while time per node remains constant. As a result, the wall-clock time increases by about 20% on Syn-1,200 and roughly a factor of two on Syn-2,100, but the core hour stays relatively static across constraint sizes.

When connecting this empirical measure to asymptotic complexity, we note that each branch-and-bound node incurs  $O(|\hat{S}|K)$  work, while exhaustive axis-aligned bisection yields at most  $O((\delta/\varepsilon)^{Km})$  (Horst & Tuy, 2013), where  $\delta$  is the initial box size and  $\varepsilon$  the target precision. Together these give a nominal worst-case cost of  $O((\delta/\varepsilon)^{Km}|\hat{S}|K)$ . On a cluster with  $P$  identical CPU cores, parallelising across open nodes yields an expected core-hour consumption of  $O\left(\frac{(\delta/\varepsilon)^{Km}|\hat{S}|K}{P}\right)$ , provided the number of simultaneously available nodes exceeds  $P$ . The empirical trends observed in Tables 1 and 2 align with this theoretical envelope.

## D.1 RUNTIME ANALYSIS OF EXHAUSTIVE $\mu$ SEARCH

Table 5 reports wall-clock statistics for the exhaustive enumeration of  $\mu$  values inside our reduced-space branch-and-bound framework<sup>2</sup>. For each synthetic dataset we list the total wall time ( $T_{\text{total}}$ ), the time devoted exclusively to the  $\mu$  search ( $T_\mu$ ), the number of branch-and-bound nodes explored ( $N_{\text{nodes}}$ ), and the average time per node ( $T_{\text{node}} = T_\mu/N_{\text{nodes}}$ ).

Table 5: Runtime metrics for n/4 must-links exhaustive  $\mu$  search (gap  $\leq 0.1\%$ ).

DATASET	$T_{\text{TOTAL}}$ (S)	$T_\mu$ (S)	$N_{\text{NODES}}$	$T_{\text{NODE}}$ (S)
SYN-21000	3,323.69	3,311.58	69	47.99
SYN-81000	12,921.19	12,804.11	71	180.34
SYN-141000	18,499.32	18,307.11	55	332.86
SYN-171000	28,796.98	28,551.09	43	664.00

The results confirm that the  $\mu$  enumeration dominates the computational budget ( $T_\mu/T_{\text{total}} > 0.99$ ), while other tasks, namely relaxations, cuts and I/O are marginal. Although  $T_{\text{node}}$  grows faster than linearly with the dataset size, parallel execution on 100 cores keeps the overall wall-time below ten hours even for the 171 k-point instance.

## D.2 CLUSTERING EVALUATION

Our work distinguishes between two aspects of clustering performance: (i) the formulation, which defines how data similarity is measured and affects metrics such as ARI, NMI, and purity; and (ii) the quality of the solution (in terms of cost minimization), which we measure using the optimality gap. Our focus is on the quality of the solution for the K-Means cost, ensuring an optimality guarantee of the solution. Here, we performed additional statistical tests and gave the ARI, NMI and purity results of 5 datasets with ground truth labels<sup>3</sup> and compare these with the algorithm of (Hua et al., 2021) in Table 6<sup>4</sup>. In these experiments, the number of clusters is set to the number of ground truth labels.

<sup>2</sup>Node count decreases with  $n$  because larger ML components reduce  $|\hat{S}|$  via pseudo-sample collapse (Section 3.2), enabling scalability up to 1.5M samples under ML constraints.

<sup>3</sup>ARI/NMI/Purity are reported only as external validation, while solution quality is measured via certified SSE gaps.

<sup>4</sup>Identical metrics confirm that CL constraints are satisfied at the global optimum; enforcement via Lemmas 3.1–3.2 and node feasibility checks remains active throughout the tree.

1026 Table 6: Clustering evaluation metrics on solutions of datasets under different constraint settings.  
1027  
1028

METRICS	CONSTRAINTS	IRIS <i>k</i> = 3	SEEDS <i>k</i> = 3	HEMI <i>k</i> = 3	HTRU2 <i>k</i> = 2	SKIN_8 <i>k</i> = 2
ARI	MSSC (HUA ET AL., 2021)	0.7163	0.7166	0.0126	-0.0779	-0.0387
	SDC-GBB (ML)	0.7859	0.7261	0.0137	-0.0385	-0.0427
	SDC-GBB (CL)	0.7163	0.7166	0.0148	0.0389	0.3545
	SDC-GBB (ML+CL)	0.7711	0.7384	0.0171	0.0909	-0.0090
NMI	MSSC (HUA ET AL., 2021)	0.7419	0.6949	0.0335	0.0265	0.0221
	SDC-GBB (ML)	0.7773	0.6979	0.0335	0.0666	0.0280
	SDC-GBB (CL)	0.7419	0.6949	0.0309	0.1007	0.4388
	SDC-GBB (ML+CL)	0.7705	0.7006	0.0338	0.1275	0.0343
PURITY	MSSC (HUA ET AL., 2021)	0.8867	0.8952	0.4210	0.9084	0.7925
	SDC-GBB (ML)	0.9200	0.9000	0.4251	0.9084	0.7925
	SDC-GBB (CL)	0.8867	0.8952	0.4373	0.9084	0.9420
	SDC-GBB (ML+CL)	0.9133	0.9048	0.4419	0.9084	0.7925

1043  
1044 E HEURISTIC ALGORITHMS  
1045

1046 We evaluate the proposed procedure by comparing its clustering quality with four reference heuristics  
1047 for the minimum-sum-of-squares clustering problem subject to must-link (ML) and cannot-link (CL)  
1048 constraints. COP- $k$ -means (Wagstaff et al., 2001) restarts the classical Lloyd algorithm one hundred  
1049 times and enforces the constraints at every assignment step. The post-processing encode- $k$ -means-  
1050 Post method of Nghiem (Nghiem et al., 2020) formulates the reassignment of instances produced by  
1051 an unconstrained or partially constrained clustering method as a binary combinatorial program that  
1052 respects all ML and CL relations. The binary linear programming approach of Baumann (Baumann,  
1053 2020) (BLPKM-CC) solves to optimality the assignment subproblem within each Lloyd iteration;  
1054 only the initial centroids are random, so the method is partially deterministic. Variants of coresset  
1055 algorithm construct a  $(k, \epsilon)$ -coreset, on which  $k$ -means can be solved quickly while guaranteeing  
1056 that the resulting centers incur at most a  $(1 + \epsilon)$  multiplicative error in squared-error cost on the full  
1057 dataset, thus preserving near-optimality with far lower computational and memory demands. We  
1058 test several coresset construction algorithms and obtain the UB for Sensitivity Sampling with  $k = 3$ ,  
1059 oversample factor  $c = 2$ , error  $\epsilon = 0.1$  and probability of approximation guarantee  $\delta = 0.1$ , as  
1060 this is the state-of-the-art method for constructing coressets (Schwiegelshohn & Sheikh-Omar, 2022)  
1061 and the only method satisfying the 4-hour runtime limit for all datasets.

1062 Tables 7, 8, 9 report the optimal UB obtained by all heuristic algorithms with 100 independent  
1063 initializations for ML-only, CL-only, and ML+CL experiments respectively. We do not include  
1064 results for Sensitivity Sampling in experiments involving CL since under coresset algorithms, any  
1065 hard cannot-link requirement collapses the additivity assumption and blows up point sensitivities,  
1066 thus inflates the coresset to linear size. Besides, we apply N/A to some COP- $k$ -means results, as this  
1067 algorithm generally could not find the global optima for datasets of size  $n > 2,000$ .

1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

1080  
1081  
1082  
1083 Table 7: Heuristic algorithms on  $\frac{n}{4}$  ML constraints  
1084  
1085  
1086  
1087  
1088  
1089  
1090  
1091  
1092  
1093  
1094  
1095  
1096  
1097  
1098  
1099

DATASETS	SIZE	COP- $k$ -MEANS	ENCODE- $k$ -MEANS-POST	BLPKM-CC	SENSITIVITY SAMPLING
IRIS	150	150.78	84.67	<b>83.82</b>	93.87
SEEDS	200	713.88	625.37	<b>620.78</b>	761.60
HEMI	1,955	N/A	<b>1.602 × 10<sup>7</sup></b>	$1.875 \times 10^7$	$2.167 \times 10^7$
PR2392	2,392	N/A	<b>3.210 × 10<sup>10</sup></b>	$3.246 \times 10^{10}$	$3.436 \times 10^{10}$
RDS_CNT	10,000	N/A	<b>6.122 × 10<sup>7</sup></b>	$6.579 \times 10^7$	$6.387 \times 10^7$
HTRU2	17,898	N/A	<b>1.407 × 10<sup>8</sup></b>	$1.472 \times 10^8$	$1.505 \times 10^8$
SPNET3D_5	50,000	N/A	<b>6.627 × 10<sup>6</sup></b>	$7.089 \times 10^6$	$7.196 \times 10^6$
SKIN_8	80,000	N/A	$5.464 \times 10^8$	$5.492 \times 10^8$	<b>4.533 × 10<sup>8</sup></b>
URBANGB	360,177	N/A	OUT OF MEMORY		<b>1.643 × 10<sup>9</sup></b>
SPNET3D	434,874	N/A	<b>5.848 × 10<sup>7</sup></b>	$6.264 \times 10^7$	$6.413 \times 10^7$
SYN-42000	42,000	N/A	<b>5.127 × 10<sup>5</sup></b>	$1.584 \times 10^6$	$5.148 \times 10^5$
SYN-210000	210,000	N/A	<b>2.163 × 10<sup>6</sup></b>	<b>2.163 × 10<sup>6</sup></b>	$2.178 \times 10^6$
SYN-420000	420,000	N/A	<b>6.010 × 10<sup>6</sup></b>	<b>6.010 × 10<sup>6</sup></b>	$6.080 \times 10^6$
SYN-1050000	1,050,000	N/A	$3.708 \times 10^7$	<b>1.050 × 10<sup>7</sup></b>	$1.052 \times 10^7$
SYN-1500000	1,500,000	N/A	$5.611 \times 10^7$	<b>1.725 × 10<sup>7</sup></b>	$1.731 \times 10^7$

1100  
1101  
1102  
1103  
1104 Table 8: Heuristic algorithms on  $\frac{n}{4}$  CL constraints  
1105  
1106  
1107  
1108  
1109  
1110  
1111  
1112  
1113  
1114  
1115  
1116  
1117  
1118  
1119

DATASET	SIZE	COP- $k$ -MEANS	ENCODE- $k$ -MEANS-POST	BLPKM-CC
IRIS	150	119.37	<b>80.31</b>	80.71
SEEDS	200	634.3	603.96	<b>603.04</b>
HEMI	1,955	$1.711 \times 10^7$	<b>1.401 × 10<sup>7</sup></b>	$1.606 \times 10^7$
PR2392	2,392	$2.596 \times 10^{10}$	<b>2.566 × 10<sup>10</sup></b>	$2.578 \times 10^{10}$
RDS_CNT	10,000	$3.696 \times 10^7$	<b>2.897 × 10<sup>7</sup></b>	$2.902 \times 10^7$
HTRU2	17,898	N/A	$1.928 \times 10^8$	<b>1.740 × 10<sup>8</sup></b>
SPNET3D_5	50,000	N/A	<b>3.938 × 10<sup>6</sup></b>	$4.027 \times 10^6$
SKIN_8	80,000	N/A	<b>6.367 × 10<sup>8</sup></b>	$6.464 \times 10^8$
SYN-12000	12,000	N/A	<b>9.503 × 10<sup>4</sup></b>	<b>9.503 × 10<sup>4</sup></b>
SYN-21000	21,000	N/A	$1.928 \times 10^8$	<b>1.740 × 10<sup>8</sup></b>
SYN-42000	42,000	N/A	<b>1.817 × 10<sup>5</sup></b>	<b>1.817 × 10<sup>5</sup></b>
SYN-210000	210,000	N/A	<b>2.161 × 10<sup>6</sup></b>	<b>2.161 × 10<sup>6</sup></b>

1120  
1121  
1122  
1123  
1124  
1125  
1126  
1127  
1128  
1129  
1130  
1131  
1132  
1133 Table 9: Heuristic algorithms on  $\frac{n}{4}$  ML +  $\frac{n}{4}$  CL constraints

DATASET	SIZE	COP- $k$ -MEANS	ENCODE- $k$ -MEANS-POST	BLPKM-CC
IRIS	150	N/A	88.75	<b>86.85</b>
SEEDS	200	N/A	601.36	<b>597.14</b>
HEMI	1,955	N/A	<b>1.566 × 10<sup>7</sup></b>	$1.762 \times 10^7$
PR2392	2,392	N/A	<b>2.922 × 10<sup>10</sup></b>	$2.945 \times 10^{10}$
RDS_CNT	10,000	N/A	<b>7.579 × 10<sup>7</sup></b>	$7.919 \times 10^7$
HTRU2	17,898	N/A	$2.218 \times 10^8$	<b>1.859 × 10<sup>8</sup></b>
SPNET3D_5	50,000	N/A	<b>8.171 × 10<sup>6</sup></b>	$8.320 \times 10^6$
SKIN_8	80,000	N/A	<b>7.579 × 10<sup>8</sup></b>	$8.775 \times 10^8$
SYN-12000	12,000	N/A	<b>9.520 × 10<sup>4</sup></b>	<b>9.520 × 10<sup>4</sup></b>
SYN-21000	21,000	N/A	$1.818 \times 10^5$	<b>1.818 × 10<sup>5</sup></b>
SYN-42000	42,000	N/A	<b>5.133 × 10<sup>5</sup></b>	<b>5.133 × 10<sup>5</sup></b>
SYN-210000	210,000	N/A	<b>2.165 × 10<sup>6</sup></b>	<b>2.165 × 10<sup>6</sup></b>