

Mathematical Programming for Data Mining

- 1 IOANNIS P. ANDROULAKIS¹, W.
 2 ART CHAOVALITWONGSE²
 3 ¹ Department of Biomedical Engineering, Rutgers
 4 University, Piscataway, NJ, USA
 5 ² Department of Industrial and Systems Engineering,
 6 Rutgers University, Piscataway, NJ, USA

7 MSC2000: **TS2**

8 Article Outline

- 9 Keywords
 10 Introduction
 11 Basic Definitions
 12 Mathematical Programming Formulations
 13 Classification
 14 Clustering
 15 Support Vector Machines
 16 Multi-Class Support Vector Machines
 17 Data Mining in the Presence of Constraints
 18 Data Mining and Integer Optimization
 19 Research Challenges
 20 References

21 Keywords

- 22 Mathematical programming; Data mining;
 23 Optimization; Clustering; Classification

24 Introduction

25 Progress in digital data acquisition and storage tech-
 26 nology has resulted in the growth of huge databases.
 27 This has occurred in a variety of scientific and engi-
 28 neering research applications [8] as well as medical
 29 domain [19,20]. Making sense out of these rapidly
 30 growing massive data sets gave birth to a “new” sci-
 31 entific discipline often referred to as *Data Mining*. Defini-
 32 ing a discipline is, however, always a controversial
 33 task. The following working definition of the area was
 34 recently proposed [9]: Data mining is the analysis of
 35 (often large) observational data sets to find unsuspect-
 36 ed relationships and to summarize the data in novel
 37 ways that are both understandable and useful to the data
 38 owner.

Clearly the term data mining is often used as a synonym
 for the process of extracting useful information from
 databases. However, the overall knowledge discovery
 from databases (KDD) process is far more complicat-
 ed and convoluted and involves a number of addition-
 al pre and post-processing steps [6]. Therefore, in our
 definition data mining refers to the ensemble of new,
 and existing, specific algorithms for extracting structure
 from data [8]. The exact definition of the knowledge
 extraction process and the expected outcomes are very
 difficult to characterize. However, a number of specific
 tasks can be identified and, by and large, define the key
 subset of deliverables from a data mining activity. Two
 such critical activities are classification and clustering.
 A number of variants for these tasks can be identi-
 fied and, furthermore, the specific structure of the data
 involved greatly impacts the methods and algorithms
 that are to be employed. Before we proceed with the
 exact definition of the tasks we need to provide work-
 ing definitions of the nature and structure of the data.

Basic Definitions

For the purposes of our analysis we will assume that
 the data are expressed in the form of n -dimensional fea-
 ture vectors $x \in X \subseteq \mathfrak{R}^n$. Appropriate pre-processing
 of the data may be required to transform the data into
 this form. Although in many cases this transformations
 can be trivial, in other cases transforming the data into
 a “workable” form is a highly non-trivial task. The goal
 of data mining is to estimate an explicit, or implicit,
 function that maps points of the feature vector from
 the input space, $X \subseteq \mathfrak{R}^n$, to an output space, C , giv-
 en a finite sample. The concept of the finite sample
 is important because, in general, what we are given is
 a finite representative subset of the original space (train-
 ing set) and we wish to make predictions on new ele-
 ments of the set (testing set). The data mining tasks can
 thus be defined based on the nature of the mapping C
 and the extent to which the train set is characterized.
 If the predicted quantity is a categorical value and if
 we know the value that corresponds to each element of
 the training set then the question becomes how to iden-
 tify the mapping that connects the feature vector and the
 corresponding categorical value (class). This problem is
 known as the classification problem (supervised learn-
 ing). If the class assignment is not known and we seek
 to: (a) identify whether a small, yet unknown, number

TS1 Please note that pagination and layout are not final.

TS2 Please provide the MSC.

85 of classes exist; (b) define the mapping assigning the
 86 features to classes then we have a clustering problem
 87 (unsupervised learning).

88 A related problem associated with superfluous infor-
 89 mation in the feature vector is the so-called feature
 90 selection problem. This is a problem closely related to
 91 over-fitting in regression. Having a minimal number of
 92 features leads to simpler models, better generalization
 93 and easier interpretation. One of the fundamental issues
 94 in data mining is therefore to identify the least num-
 95 ber of features, sub-set of the original set of features,
 96 that best address the two issues previously defined. The
 97 concept of parsimony (Occam’s razor) is often invoked
 98 to bias the search [1]: never do with more what can be
 99 done with fewer.

100 Although numerous methods exist for addressing these
 101 problems they will not be reviewed here. Nice reviews
 102 of classification and were recently presented in [8,9]. In
 103 this short introduction we will concentrate on solution
 104 methodologies based on reformulating the clustering,
 105 and classification questions as optimization problems.

106 **Mathematical Programming Formulations**

107 Classification and clustering, and for that matter most
 108 of the data mining tasks, are fundamentally optimiza-
 109 tion problems. Mathematical programming methodolo-
 110 gies formalize the problem definition and make use of
 111 recent advances in optimization theory and applications
 112 for the efficient solution of the corresponding formula-
 113 tions. In fact, mathematical programming approaches,
 114 particularly linear programming, have long been used
 115 in data mining tasks.

116 The pioneering work presented in [13,14] demonstrated
 117 how to formulate the problem of constructing planes to
 118 separate linearly separable sets of points.

119 In this summary we will follow the formalism put forth
 120 in [2] since it presented one of the most comprehensive
 121 approaches to this problem. One of the major advan-
 122 tages of a formulation based on mathematical program-
 123 ming is the ease in incorporating explicit problem spe-
 124 cific constraints. This will be discussed in greater detail
 125 later in this summary.

126 **Classification**

127 As discussed earlier the main goal in classification is to
 128 predict a categorical variable (class) based on the values
 129 of the feature vector. The general families of methods

for addressing this problem include [9]:

- i) Estimation of the conditional probability of observ-
 ing class C given the feature vector x .
- ii) Analysis of various proximity metrics and based
 the decision of class assignment based on proximi-
 ty.
- iii) Recursive input space partitioning to maximize
 a score of class purity (tree-based methods).

The two-class classification problem can be formulat-
 ed as the search of a function that assigns a given input
 vector x into two disjoint point sets A and B . The data
 are represented in the form of matrices. Assuming that
 the set A has m elements and the set B has k elements,
 then $A \in \mathfrak{R}^{m \times n}$, $B \in \mathfrak{R}^{k \times n}$, describe the two sets
 respectively. The discrimination is based on the deriva-
 tion of hyperplane

$$P = \{x | x \in \mathfrak{R}^n, x^T \omega = \gamma\}$$

with normal and distance from the origin $\frac{|\gamma|}{\|\omega\|_2}$. The
 optimization problem then becomes to determine ω and
 γ such that the separating hyperplane P defines two
 open half spaces

$$\{x | x \in \mathfrak{R}^n, x^T \omega < \gamma\}$$

$$\{x | x \in \mathfrak{R}^n, x^T \omega > \gamma\}$$

containing mostly points in A and B respectively.
 Unless A and B are disjoint the separation can only be
 satisfied within some error. Minimization of the aver-
 age violations provides a possible approximation of the
 separating hyperplane [2]:

$$\min_{\omega, \gamma} \frac{1}{m} \|(-A\omega + e\gamma + e)_+\|_1 + \frac{1}{k} \|(-B\omega + e\gamma + e)_+\|_1$$

In [2] a number of linear programming reformulations
 are discussed exploring the properties of the structure
 of the optimization problem. In particular an effective
 robust linear programming (RLP) reformulation was
 suggested making possible the solution of large-scale
 problems:

$$\begin{aligned} \min_{\omega, \gamma, y, z} & \frac{e^T y}{m} + \frac{e^T z}{k} \\ \text{s.t.} & -A\omega + e\gamma + e \leq y \\ & B\omega - e\gamma + e \leq z \\ & y, z \geq 0. \end{aligned}$$

165 In [17] it was demonstrated how the above formulation
 166 can be applied repeatedly to produce complex space
 167 partitions similar to those obtained by the application
 168 of standard decision tree methods such as C4.5 [21] or
 169 CART [4].

170 **Clustering**

171 The goal of clustering is the segmentation of the raw
 172 data into groups that share a common, yet unknown,
 173 characteristic property. Similarity is therefore a key
 174 property in any clustering task. The difficulty arises
 175 from the fact that the process is unsupervised. That is
 176 neither the property nor the expected number of groups
 177 (clusters) are known ahead of time. The search for the
 178 optimal number of clusters is parametric in nature and
 179 the optimal point in an “error” vs. “number of clusters”
 180 curve is usually identified by a combined objective the
 181 weighs appropriately accuracy and number of clusters.
 182 Conceptually a number of approaches can be developed
 183 for addressing clustering problems:

- 184 i) Distance-based methods, by far the most commonly
 185 used, that attempt to identify the best k-way parti-
 186 tion of the data by minimizing the distance of the
 187 points assigned to cluster k from the center of the
 188 cluster.
- 189 ii) Model-based methods assume the functional form
 190 of a model that describes each of the clusters and
 191 then search for the best parameter fit that models
 192 each cluster by minimizing some appropriate likeli-
 193 hood measure.

194 There are two different types of clustering: (1) hard
 195 clustering; (2) fuzzy clustering. The former assigns
 196 a data point to *exactly* one cluster while the latter
 197 assigns a data point to one of more clusters along with
 198 the likelihood of the data point belonging to one of
 199 those clusters.

200 The standard formulation of the hard clustering prob-
 201 lem is:

$$202 \min_c \sum_{i=1}^m \min_l \|x^i - c^l\|_n$$

203 That is given m points, x , in an n -dimensional space,
 204 and a fixed number of cluster, k , determine the centers
 205 of the cluster, c , such that the sum of the distances of
 206 each point to a nearest cluster center is minimized. It
 207 was shown in [3] that this general non convex problem

208 can be reformulated such that we minimize a bilinear
 209 functions over a polyhedral set by introducing a selec-
 210 tion variable t_{il} :

$$\begin{aligned} & \min_{c,d,t} \sum_{i=1}^m \sum_{l=1}^k t_{il} (e^T d_{il}) \\ & \text{s.t. } -d_{il} \leq x^i - c^l \leq d_{il} \\ & \sum_{l=1}^k t_{il} = 1 \\ & t_{il} \geq 0 \\ & i = 1, \dots, m, l = 1, \dots, k. \end{aligned} \quad 211$$

212 d is a dummy variable used to bound the components
 213 of the difference $x - c$. In the above formulation the
 214 1-norm is selected [3].

215 The fuzzy clustering problem can be formulated as fol-
 216 lows [5]:

$$\begin{aligned} & \min_w \sum_{i=1}^m \sum_{l=1}^k w_{il}^2 \|x^i - c^l\|^2 \\ & \text{s.t. } \sum_{l=1}^k w_{il} = 1 \\ & w_{il} \geq 1, \end{aligned} \quad 217$$

218 where $x^i, i = 1, \dots, m$ is the location descriptor for the
 219 data point, $c^l, l = 1, \dots, k$ is the center of the cluster,
 220 w_{il} is the likelihood of a data point i being assigned to
 221 cluster l .

222 **Support Vector Machines**

223 This optimization formalism bares significance resem-
 224 blance to the Support Vector Machines (SVM) frame-
 225 work [25]. SVM incorporate the concept of structural
 226 risk minimization by determining a separating hyper-
 227 plane that maximizes not only a quantity measuring the
 228 misclassification error but also maximizing the mar-
 229 gin separating the two classes. This can be achieved
 230 by augmenting the objective of the RLP formulation
 231 earlier presented by an appropriately weighted mea-
 232 sure of the separation between the two classes as
 233 $(1 - \lambda)(e^T y + e^T z) + \frac{\lambda}{2} \|\omega\|_2^2$.

234 In [6] the concept of SVM is extended by introduc-
 235 ing the Proximal support vector machines which clas-
 236 sify points based on proximity to one of two parallel
 237 planes that are pushed as far apart as possible. Non-
 238 linear transformations were also introduced in [6] to

239	enable the derivation of non-linear boundaries in clas-	Implicit enumeration techniques such as branch-and-	281
240	sifiers.	bound were used early on to address the problem of	282
		feature selection [18].	283
241	Multi-Class Support Vector Machines	Mathematical programming inspired by algorithms for	284
242	Support vector machines were originally designed for	addressing various data mining problems are now being	285
243	binary classification. Extending to multi-class problems	revisited and cast as integer optimization problems.	286
244	is still an open research area [10].	Representative formulations include feature selection	287
245	The earliest multi-class implementation is the <i>one</i>	using Mixed-Integer Linear Programs [11] and in [23]	288
246	<i>against all</i> [22] by constructing k SVM models, where	integer optimization models are used to address the	289
247	k is the number of classes. The i th SVM is classifies	problem of classification and regression.	290
248	the examples of class i against all the other samples in		
249	all other classes. Another alternative builds <i>one against</i>	Research Challenges	291
250	<i>one</i> [12] classifiers by building $\frac{k(k-1)}{2}$ models where	Numerous issues can of course be raised. However, we	292
251	each is trained on data from two classes. The empha-	would like to focus on three critical aspects	293
252	sis of current research is on novel methods for gener-		
253	ating all the decision functions through the solution of	i) Scalability and the curse of dimensionality. Data-	294
254	a single, but much larger, optimization problem [10].	bases are growing extremely fast and problems of	295
		practical interest are routinely composed of millions	296
255	Data Mining in the Presence of Constraints	of records and thousands of features. The compu-	297
256	Prior knowledge about a system is often omitted in	tational complexity is therefore expected to grow	298
257	data mining applications because most algorithms do	beyond what is currently reasonable and tractable.	299
258	not have adequate provisions for incorporating explic-	Hardware advances alone will not address this	300
259	itly such types of constrains. Prior knowledge can either	problem either as the increase in computational	301
260	encodes explicit and/or implicit relations among the	complexity outgrows the increase in computatio-	302
261	features or models the existence of “obstacles” in the	nal speed. The challenge is therefore two-fold: either	303
262	feature space [24].	improve the algorithms and the implementation of	304
263	One of the major advantages of a mathematical pro-	the algorithms or explore sampling and dimension-	305
264	gramming framework for performing data mining tasks	ality reduction techniques.	306
265	is that prior knowledge can be incorporated in the def-	ii) Noise and infrequent events. Noise and uncertain-	307
266	inition of the various tasks in the form of (non)linear	ty in the data is a given. Therefore, data mining	308
267	constraints. Efficient incorporation of prior knowledge	algorithms in general and mathematical program-	309
268	in the form of nonlinear inequalities within the SVM	formulations in particular have to account for	310
269	framework was recently proposed by [15]. Reformu-	the presence of noise. Issues from robustness and	311
270	lations of the original linear and nonlinear SVM clas-	uncertainty propagation have to be incorporated.	312
271	sifiers to accommodate prior knowledge about the	However, an interesting issue emerges: how do we	313
272	problem were presented in [7] in the context of approx-	distinguish between noise and an infrequent, albeit	314
273	imation and in [16] in the context of classifiers.	interesting observation? This in fact maybe a ques-	315
		tion with no answer.	316
274	Data Mining and Integer Optimization	iii) Interpretation and visualization. The ultimate goal	317
275	Data mining tasks involve, fundamentally, discrete	of data mining is understanding the data and devel-	318
276	decisions:	oping actionable strategies based on the conclu-	319
277	• How many clusters are there?	sions. We need to improve not only the inter-	320
278	• Which class does a record belong to?	pretation of the derived models but also the	321
279	• Which features are most informative?	knowledge delivery methods based on the derived	322
280	• Which samples capture the essential information?	models. Optimization and mathematical program-	323
		ming needs to provide not just the optimal solution	324
		but also some way of interpreting the implications	325

of a particular solution including the quantification of potential crucial sensitivities.

References

1. Blumer A, Ehrenfeucht A, Haussler D, Warmuth MK (1987) Occam's razor. *Inf Process Lett* 24:377–380
2. Bradley PS, Fayyad U, Mangasarian OL (1999) Mathematical programming for data mining: Formulations and challenges. *INFORMS J Comput* 11:217–238
3. Bradley PS, Mangasarian OL, Street WN (1997) Clustering via concave minimization. In: Mozer MC, Jordan MI, Petsche T (eds) *Advances in Neural Information Processing Systems*, MIT Press, pp 368–374 **CE3**
4. Breiman L, Friedman J, Olsen R, Stone C (1993) *Classification and Regression Trees*. Wadsworth Inc., **CE3**
5. Dunn JC (1973) A fuzzy relative of the ISODATA process and its use in detecting compact well-separated clusters, *J Cybern* 3:32–57
6. Fayyad U, Piatetsky-Shapiro G, Smyth P (1996) From data mining to knowledge discovery: An overview In: Fayyad UM, Piatetsky-Shapiro G, Smyth P, Uthursamy R (eds) *Advances in knowledge discovery and data mining*. AAAI Press, pp 229–248 **CE3**
7. Fung GM, Mangasarian OL (2001) Proximal Support Vector Machine Classifiers In: Provost F, Srikant R (eds) *Proceedings KDD-2001: Knowledge Discovery and Data Mining*, San Francisco, August 26–29 2001, Association for Computing Machinery, pp 77–86, <ftp://ftp.cs.wisc.edu/pub/dmi/tech-reports/01-02.ps> **CE4**
8. Grossman RL, Kamath C, Kegelmeyer P, Kumar V, Namburu EE (2001) *Data mining for scientific and engineering applications*. Kluwer Academic Publishers, **CE3**
9. Hand DJ, Mannila H, Smyth P (2001) *Principle of data mining*. Bradford Books, **CE3**
10. Hsu C-W, Lin C-J (2002) A comparison of methods multi-class support vector machines. *IEEE Trans Neural Netw* 13:415–425
11. Iannatilli FJ, Rubin PA (2003) Feature selection for multi-class discrimination via mixed-integer linear programming. *IEEE Trans Pattern Anal Mach Learn* 25:779–783
12. Krebel U (1999) *Pairwise classification and support vector machines: Advances in Kernel Methods – Support Vector Learning*. MIT Press **CE3**
13. Mangasarian OL (1965) Linear and nonlinear separation of pattern by linear programming. *Oper Res* 31:445–453
14. Mangasarian OL (1968) Multisurface method for pattern separation. *IEEE Trans Inf Theory* IT-14:801–807
15. Mangasarian OL, Shavlik JW, Wild EW (2003) Knowledge-based kernel approximations. *Tech. rep.*, Data Mining Institute, University of Wisconsin **CE3**
16. Mangasarian OL, Shavlik JW, Wild EW (2004) Knowledge-based kernel approximation. *J Mach Learn Res* 5:1127–1141
17. Mangasarian OL, Street WN, Wolberg WH (1995) Breast cancer diagnosis and prognosis via Linear Programming. *Oper Res* 43:570–577
18. Narendra P, Fukunaga K (1977) A branch and bound algorithm for feature subset selection. *IEEE Trans Comput* 26:917–922
19. Pardalos PM, Boginski V, Vazakopoulos A (2007) *Data Mining in Biomedicine*. Springer **CE3**
20. Pardalos PM, Principe J (2002) *Biocomputing*. Kluwer Academic Publishers **CE3**
21. Quinlan JR (1993) *C4.5: Programs for Machine Learning*. Morgan Kaufmann **CE3**
22. Scholkopf B, Burges C, Vapnik V (1995) Extracting support data for a given task: Proc. First International Conference on Knowledge Discovery and Data Mining. AAAI Press, pp 252–257 **CE3**
23. Shioda R (2003) *Integer Optimization in Data Mining*. PhD thesis, MIT **CE3**
24. Tung AK, Hou J, Han J (2001) Spatial clustering in the presence of obstacles In: *Proceedings ICDE-2001: 17th International Conference on Data Engineering*, pp 359–367 **CE3**
25. Vapnik VN (1995) *The nature of statistical learning*, Springer **CE3**

CE3 Please provide publisher location.

CE4 Please provide access date.