

# Using Machine Learning to Decide When to Precondition Cylindrical Algebraic Decomposition With Groebner Bases

Zongyan Huang\*, Matthew England†, James H. Davenport‡, and Lawrence C. Paulson\*

\*University of Cambridge Computer Laboratory, Cambridge CB3 0FD, U.K.

Email: rubyhuang87@gmail.com; lp15@cam.ac.uk

†Coventry University Faculty of Engineering, Environment and Computing, Coventry, CV1 2JH, U.K.

Email: Matthew.England@coventry.ac.uk

‡University of Bath Dept. of Computer Science, Bath, BA2 7AY, U.K.

Email: J.H.Davenport@bath.ac.uk

**Abstract**—Cylindrical Algebraic Decomposition (CAD) is a key tool in computational algebraic geometry, particularly for quantifier elimination over real-closed fields. However, it can be expensive, with worst case complexity doubly exponential in the size of the input. Hence it is important to formulate the problem in the best manner for the CAD algorithm. One possibility is to precondition the input polynomials using Groebner Basis (GB) theory. Previous experiments have shown that while this can often be very beneficial to the CAD algorithm, for some problems it can significantly worsen the CAD performance.

In the present paper we investigate whether machine learning, specifically a support vector machine (SVM), may be used to identify those CAD problems which benefit from GB preconditioning. We run experiments with over 1000 problems (many times larger than previous studies) and find that the machine learned choice does better than the human-made heuristic.

## I. INTRODUCTION

### A. Cylindrical Algebraic Decomposition

A *Cylindrical Algebraic Decomposition* (CAD) is a decomposition of ordered  $\mathbb{R}^n$  space into cells arranged *cylindrically*, meaning the projections of any pair with respect to the given ordering are either equal or disjoint. *Algebraic* is actually short for semi-algebraic as each CAD cell can be described with a finite sequence of polynomial constraints. A CAD is produced to be invariant for input: sign- (or order-) invariant for input polynomials or truth-invariant for input formulae.

CADs and the first algorithm to compute them were introduced by Collins in 1975 [19]. CAD usually has two stages: projection where an operator is applied recursively on the input to derive corresponding problems in lower dimensions; and lifting where CADs are built incrementally by dimension according to the polynomials identified in projection [1]. The original motivation was quantifier elimination (QE) in real closed fields, while other applications include: parametric optimisation [33], epidemic modelling [13], theorem proving [48], reasoning with multi-valued functions [22] derivation of optimal numerical schemes [30], and much more.

CAD has worst case complexity doubly exponential in the number of variables [23] applicable whatever the data structure [12]. For some applications there exist algorithms with better

complexity [3], but CAD implementations still remain the best general purpose approach for many. This may be due to the numerous approaches used to improve the efficiency of CAD since Collins' original work including: improvements to the projection operator [39], [47], [10], [38]; partial CAD (lift only when necessary for QE) [21]; and symbolic-numeric lifting schemes [54], [42]. Some recent advances include making use of any Boolean structure in the input [6], [7], [27]; local projection approaches [11], [55]; and decompositions via complex space [18], [5]. For a more detailed introduction to CAD see for example Bradford et al. [7].

### B. Preconditioning with Groebner Bases

A *Groebner Basis*  $G$  is a particular generating set of an ideal  $I$  defined with respect to a monomial ordering. One definition is that the ideal generated by the leading terms of  $I$  is generated by the leading terms of  $G$ . Groebner Bases (GB) allow properties of the ideal to be deduced such as dimension and number of zeros and so are one of the main practical tools for working with polynomial systems. Their properties and an algorithm to derive a GB for any ideal was introduced by Buchberger in his PhD thesis of 1965 [14].

Like CAD, there has been much research to improve and optimise GB calculation, with the  $F_5$  algorithm [31] perhaps the most used approach currently. However, also like CAD the calculation of GB is necessarily doubly exponential in the worst case [46] (when using a lexicographic monomial ordering). Despite this, the computation of GB can often be done very quickly and would almost certainly be a superior tool to CAD for any problem involving only polynomial equalities. From this arises the natural question: is the process of replacing a conjunction of polynomial equalities in a CAD problem by their GB a useful precondition for CAD?

I.e. let  $E = \{e_1, e_2, \dots\}$  be a set of polynomials;  $G = \{g_1, g_2, \dots\}$  a GB for  $E$ ; and  $B$  any Boolean combination of constraints,  $f_i \sigma_i 0$ , where  $\sigma_i \in \{<, >, \leq, \geq, \neq, =\}$  and  $F = \{f_1, f_2, \dots\}$  is another set of polynomials. Then

$$\Phi = (e_1 = 0 \wedge e_2 = 0 \wedge \dots) \wedge B \text{ and}$$

$$\Psi = (g_1 = 0 \wedge g_2 = 0 \wedge \dots) \wedge B$$

are equivalent and a CAD truth-invariant for either could be used to solve problems involving  $\Phi$  (such as eliminating any quantifiers applied to  $\Phi$ ). So is it worth producing  $G$ ?

The first attempt to answer this question was given by Buchberger and Hong in 1991 [15] who used the implementation of GB [4] to precondition an implementation of CAD [21] (both in C on top of the SAC-2 system [20]). Of the ten test problems studied: 6 were improved by the GB preconditioning, with the speed-up varying from 2-fold to 1700-fold; 1 problem resulted in a 10-fold slow-down; 1 timed out when GB preconditioning was applied, while it would complete without it; and the other 2 were intractable both for CAD alone and the GB preconditioning step.

The problem was recently revisited by Wilson et al. [57]. The authors recreated the experiments of Buchberger and Hong [15] using QEPCAD-B for the CAD and MAPLE 16 for the GB. As we may expect, there had been a big decrease in the computation timings, especially the GB: the two test problems previously intractable [15] could now have the GB calculated quickly. However, two of the CAD problems were still hindered by GB preconditioning. The experiments were then extended to: a wider example set (an additional 12 problems); the alternative CAD implementation in MAPLE-16 [18]; and the case where we further precondition by reducing inequalities of the system (the set  $F$  above) with respect to the GB. The key conclusion remained that GB preconditioning would in general benefit CAD (sometimes significantly) but could on occasion hinder it (to the point of making a tractable CAD problem intractable). The authors defined a metric to assist with the decision of when to precondition, the *Total Number of Indeterminates* (TN<sub>OI</sub>) of a set of polynomials  $A$ ,

$$\text{TN}_{\text{OI}}(A) = \sum_{a \in A} \text{NOI}(a) \quad (1)$$

where  $\text{NOI}(a)$  is the number of indeterminates in a polynomial  $a$ . Then their heuristic was to build a CAD for the preconditioned polynomials only if the TN<sub>OI</sub> decreased following preconditioning. For most of their test problems the heuristic made the correct choice, but there were examples to the contrary and little correlation between the change in TN<sub>OI</sub> and level of speed-up / slow-down.

### C. Contribution and plan

In this paper we consider if machine learning can be applied to the decision of whether preconditioning CAD input with GB is beneficial for a particular problem. We work on the reasonable assumption that GB computation is cheap for the problems on which CAD is tractable (in fact as shown in [29] the CAD will compute resultants which overestimate the GB). Hence we use algebraic features of both the input problem and the GB itself to decide *whether we want to use the GB*.

In Section II we describe the dataset and computer algebra computations used for the experiment and in Section III we describe the set of features identified to train the machine learning algorithm: a Support Vector Machine (SVM). Then in Section IV we describe the initial machine learning experiment and its results, before running feature selection experiments in Section V. Finally, we compare the machine learned decision

with the human developed TN<sub>OI</sub>-based heuristic, draw our conclusions and discuss future work in Section VI.

This is the second paper of the present authors to consider the application of SVMs to CAD optimisation. We previously studied the choice of variable ordering for CAD in [41]. In that paper 3 existing heuristics were evaluated against a machine learned choice of which to use. The latter outperformed each individually and suggested a greater role for machine learning in such decisions, motivating the present study. The only other application of machine learning to computer algebra that the authors are aware of is by Kobayashi et al. [44] who applied a SVM to decide the order of sub-formulae solving for QE.

## II. DATASET AND COMPUTER ALGEBRA

### A. Computer Algebra

All the computer algebra computations were conducted in MAPLE-17. The CAD algorithm used was an implementation of [18]. This is part of the `RegularChains Library`<sup>1</sup> [16], [17] whose CAD procedures differ from the traditional projection and lifting framework of Collins, instead first decomposing  $\mathbb{C}^n$  cylindrically and then refining to a CAD of  $\mathbb{R}^n$ . Previous experiments [57] showed this implementation has the same issues of GB preconditioning as the traditional approach. The default MAPLE GB implementation was used: a meta algorithm calling multiple GB implementations. The GBs were computed with a purely lexicographical ordering of monomials based on the same variable ordering as the CAD.

All computations were performed on a 2.4GHz Intel processor, however this is not relevant as we evaluated the CAD performance using cell counts instead of timings, i.e. by comparing the numbers of cells in the final outputted CADs produced with and without GB preconditioning. Numerous previous studies have shown this to be closely correlated to timings and it has the advantage of being discrete, machine and (up the theory used) implementation independent. It also correlates with the cost of any post-processing of the CAD.

### B. Dataset

A key difficulty in applying machine learning techniques to computer algebra is the lack of suitable datasets (sets such as [56] do not have enough problems). In our previous study on choosing the variable ordering for CAD [41] we used the `nlsat-dataset` [58], which although developed for non-linear arithmetic SAT-solvers, contained many suitable problems.

For the present experiment we need problems that are expressed with a conjunction of at least two equalities in order to build a non-trivial GB. From the `nlsat` dataset 493 three-variable problems and 403 four-variable problems fit this criteria, which should have been a sufficient number. GB preconditioning was applied to each problem and cell counts from computing the CAD with the original polynomials and their replacement with the GB were computed and compared. For each one of these problems the GB preconditioning was beneficial or made no difference; surprising as the experiments on much smaller datasets [15], [57] had shown much greater

<sup>1</sup><http://www.regularchains.org>

volatility. This points to an undetected uniformity within the current nlsat dataset. It would need to be widened if it is to be used more extensively for computer algebra research.

Since existing datasets were not suitable for the present experiment, we had no choice but to generate our own problems. The generation process aimed for an unbiased dataset which would be computationally feasible for computing multiple CADs, and have some comparable structure (number of terms and polynomials) to existing CAD problems.

In total, 1200 problems were generated using the random polynomial generator `randpoly` in MAPLE-17. Each problem has two sets of three polynomials; the first to represent conjoined equalities and the second for the other polynomial constraints (respectively  $E$  and  $F$  from the description in Section I-B). The number of variables was at most 3, labelled  $x, y, z$  and under ordering  $x \prec y \prec z$ ; the number of terms per polynomial at most 2; the coefficients were restricted to integers in  $[-20, 20]$ ; and the total degree was varied between 2, 3 and 4 (with 400 problems generated for each).

A time limit of 300 CPU seconds was set for each CAD computation (all GB computations completed quickly) from which 1062 problems finished to constitute the final dataset. Of these, 75% benefited from GB preconditioning. So our randomly generated dataset matched the previously found results of [15], [57] with most problems benefiting but not all and so is suitable for the purpose of this experiment.

### III. PROBLEM FEATURES

Table III shows the 28 problem features we identified (guided by previous work). Here  $(x, y, z)$  are the three variable labels used in the problems and proportion means the percentage of the total. The features were chosen as easily computable metrics that may affect the cell count of the CAD. They fall into two sets: those generated from the polynomials in the original problem and those obtained after applying GB preconditioning. The abbreviations `tds` and `stds` stand for *maximum total degree* and *sum of total degrees* respectively:

$$\text{tds}(F) = \max_{f \in F} \text{tds}(f), \quad \text{stds}(F) = \sum_{f \in F} \text{tds}(f).$$

We also make use of the metric `TNoI` (see equation (1)) [57]. Finally we included the base 2 logarithm of the ratio of differences between some of the key metrics. All features could be calculated immediately within MAPLE.

We note that the `stds` measure differs from the `sotd` heuristic introduced in [25] and used for multiple CAD decisions [8], [26], [28]. This is because `stds` measures the input polynomials only, while `sotd` measures the full set of CAD projection polynomials, and so is much more expensive.

In addition to training a classifier using all the features in Table III, we trained classifiers using two subsets: one containing the features labelled 1 – 12 concerning the original set of polynomials; and one the features labelled 13 – 25 concerning the polynomials after GB preconditioning. The latter set has one extra feature, the number of polynomials, as this varies after GB calculation but was always 6 to start with. We refer to the first subset as *before features*, the second as *after features* and the full set as *all features*.

TABLE I  
INITIAL FEATURE SET

Description	
1	TNoI before GB.
2	stds before GB.
3	tds of polynomials before GB.
4	Max degree of $x$ in polynomials before GB.
5	Max degree of $y$ in polynomials before GB.
6	Max degree of $z$ in polynomials before GB.
7	Proportion of polynomials with $x$ before GB.
8	Proportion of polynomials with $y$ before GB.
9	Proportion of polynomials with $z$ before GB.
10	Proportion of monomials with $x$ before GB.
11	Proportion of monomials with $y$ before GB.
12	Proportion of monomials with $z$ before GB.
13	Number of polynomials after GB.
14	TNoI after GB.
15	stds after GB.
16	tds of polynomials after GB.
17	Max degree of $x$ in polynomials after GB.
18	Max degree of $y$ in polynomials after GB.
19	Max degree of $z$ in polynomials after GB.
20	Proportion of polynomials with $x$ after GB.
21	Proportion of polynomials with $y$ after GB.
22	Proportion of polynomials with $z$ after GB.
23	Proportion of monomials with $x$ after GB.
24	Proportion of monomials with $y$ after GB.
25	Proportion of monomials with $z$ after GB.
26	$\log_2(\text{TNoI before GB}) - \log_2(\text{TNoI after GB})$
27	$\log_2(\text{stds before GB}) - \log_2(\text{stds after GB})$
28	$\log_2(\text{tds before GB}) - \log_2(\text{tds after GB})$

**Example:** Consider sets of polynomials

$$E := \{-12yz - 3z, \quad 17x^2 - 6, \quad -2yz + 5x\}$$

$$F := \{-2yz - 9y, \quad -15x^2 - 19y, \quad 6xz + 3\}.$$

The GB computed for  $E$  is

$$G := \{17x^2 - 6, \quad 4y + 1, \quad z + 10x\}$$

and the *all features* vector becomes

$$\left[12, 12, 2, 2, 1, 1, \frac{2}{3}, \frac{2}{3}, \frac{2}{3}, \frac{1}{3}, \frac{5}{12}, \frac{5}{12}, \right. \\ \left. 6, 10, 10, 2, 2, 1, 1, \frac{2}{3}, \frac{1}{2}, \frac{1}{2}, \frac{1}{3}, \frac{1}{3}, \frac{1}{4}, \right. \\ \left. 0.263, 2.263, 0\right],$$

with the *before features* and *after features* vectors formed from the first and second line respectively.

The feature generation process was applied to create the three training sets separately (although the feature labels used were all as in Table III). Each problem was labelled +1 if GB preconditioning is beneficial for CAD construction, or -1 otherwise. After feature generation the training data was standardised so each feature had zero mean and unit variance across the training set. The same standardisation was then applied to features in the test set.

## IV. MACHINE LEARNED CHOICES

### A. Introduction

Machine learning deals with the design of programs that learn rules from data. This is an attractive alternative to manually constructing them when the underlying functional relationship is complex, as appears to be the case here.

In the last decade, the use of machine learning has spread rapidly following the invention of the *Support Vector Machine* (SVM) (see for example [51]). This gives a powerful and robust method for both: *Classification*, the assignment of input examples into a given set of classes; and *Regression*, a supervised pattern analysis in which the output is real-valued. The standard SVM classifier takes a set of input data and predicts one of two possible classes from the input. Given a set of examples, each marked as belonging to one of two classes, an SVM training algorithm builds a model that assigns new examples into one of the classes. The examples used to fit the model are called training examples. An important concept in the SVM theory is the use of a kernel function to map data into a high dimensional feature space and then separate samples in the transformed space [53]. Kernel functions enable operations in feature space without ever computing the coordinates of the data in that space, rather they compute the inner products between all pairs of data vectors, which is generally cheaper.

For our experiment we used SVM-LIGHT<sup>2</sup> [43]; an implementation of SVMs in C.

### B. Cross-validation and grid-search

The 1062 problems were partitioned into 80% training (849 problems) and 20% test (213 problems), stratified to maintain relative proportions of positive and negative examples. The classification was done using the *radial basis function* (RBF) kernel. This was chosen after earlier experiments applying machine learning to an automated theorem prover found the RBF kernel to perform well with similar simple algebraic features [9]. The RBF function is defined as:

$$K(x, x') = \exp(-\gamma \|x - x'\|^2) \quad (2)$$

where  $x$  and  $x'$  are feature vectors. The process depends on kernel parameter  $\gamma$  and another parameter  $C$  which governs the trade-off between margin and training error, and finding the optimal values of these is not trivial. *Matthews' Correlation Coefficient* (MCC) [45], [2] is often used to evaluate choices. This takes into account true and false positives and negatives (labelled tp, fp, tn and fn accordingly):

$$\text{MCC} = \frac{\text{tp} * \text{tn} - \text{fp} * \text{fn}}{\sqrt{(\text{tp} + \text{fp})(\text{tp} + \text{fn})(\text{tn} + \text{fp})(\text{tn} + \text{fn})}}. \quad (3)$$

In the case where one of the terms in the denominator is zero the entire denominator is set to 1. The MCC measure has the value 1 if perfect prediction is attained, 0 if the classifier is performing as a random classifier, and  $-1$  if the classifier exactly disagrees with the data.

A grid-search optimisation procedure along with a five-fold stratified cross validation was used, involving a search over a range of  $(\gamma, C)$  values to find the pair which would maximize equation (3). We tested a commonly used range of values in our grid search process [40]:  $\gamma$  varied between  $\{2^{-15}, 2^{-14}, 2^{-13}, \dots, 2^3\}$ ; and  $C$  varied between  $\{2^{-5}, 2^{-4}, 2^{-3}, \dots, 2^{15}\}$ . Following the completion of the grid-search, the values giving optimal MCC results were selected. This procedure was repeated for the three feature sets.

TABLE II  
ACCURACY OF PREDICTIONS

Feature Set	Number	% of test set
All features	162	76%
Before features	159	75%
After features	167	78%

### C. Results for the three feature sets

The classification accuracy was used to measure the efficacy of the machine learning selection process under the 3 feature sets. The test set of 213 problems contained 159 positive samples and 54 negative samples (i.e. 75% of the test problems benefited from GB preconditioning). The results of the machine learned choices are summarised in Table II. First we note that when making a choice based on the *before features* training set 75% of the problems were predicted accurately. I.e. making a decision based on these features results in no more correct decisions than blindly deciding to GB precondition each and every time. However, the other two feature sets resulted in superior decisions. Although only a small improvement on preconditioning blindly, we recall that the wrong choice can give large changes to the size of the CAD or even change the tractability of the problem [57].

The results indicate that the features of the GB itself are required to decide whether to use the preconditioning. However, earlier research showed that a variable completely useless by itself can provide a significant performance improvement when taken in conjunction with others [34]. To be confident about which features were significant and which were superfluous, further feature selection experiments are required and we will see that the optimal feature subset must contain features from both before and after the GB computation.

## V. FEATURE SELECTION

There is a strong indication that not all features contribute to the machine learning process. Moreover, a reduced feature set can be beneficial for better understanding the underlying connections. Consequently, we applied some *feature selection* methods. Both filter and wrapper methods were applied as discussed in the following subsections.

The feature selection experiments were conducted with WEKA (Waikato Environment for Knowledge Analysis) [35], a Java machine learning library which supports tasks such as data preprocessing, clustering, classification, regression and feature selection. Each data point is also represented as a fixed number of features. The inputs are samples of 29 features, where the first 28 are the real-valued features from Table III, and the final one is a *nominal* feature denoting its class.

### A. The filter method

A correlation based feature selection method, was applied as described in [37]. Unlike other filter methods [36], these measure the rank of feature subsets instead of individual features. A feature subset which contains features highly correlated with the class but uncorrelated with each other is preferred. The metric below is used to measure the quality of a

<sup>2</sup><http://svmlight.joachims.org>

feature subset, and takes into account feature-class correlation as well as feature-feature correlation.

$$G_s = \frac{k\overline{r_{ci}}}{\sqrt{k + k(k-1)\overline{r_{ii'}}}} \quad (4)$$

Here,  $k$  is the number of features in the subset,  $\overline{r_{ci}}$  denotes the average feature-class correlation of feature  $i$ , and  $\overline{r_{ii'}}$  the average feature-feature correlation between feature  $i$  and  $i'$ . The numerator of equation (4) indicates how much relevance there is between the class and a set of features, while the denominator measures the redundancy among the features. The higher  $G_s$ , the better the feature subset.

To apply this heuristic we must calculate the correlations. With the exception of the class attribute all 28 features are continuous, so in order to have a common measure for computing the correlations we first discretize using the method of Fayyad and Irani [32]. After that, a correlation measure based on the information-theoretical concept of *entropy* is used, which is a measure of the uncertainty of a random variable. We define the *entropy of a variable*  $X$  [52] as

$$H(X) = -\sum_i p(x_i) \log_2(p(x_i)). \quad (5)$$

The entropy of  $X$  after observing values of another variable  $Y$  is then defined as

$$H(X|Y) = -\sum_j p(y_j) \sum_i p(x_i|y_j) \log_2(p(x_i|y_j)), \quad (6)$$

where  $p(x_i)$  is the prior probabilities for all values of  $X$ , and  $p(x_i|y_j)$  is the posterior probabilities of  $X$  given the values of  $Y$ . The *information gain* (IG) [50] measures the amount by which the entropy of  $X$  decreases by additional information about  $X$  provided by  $Y$ , and it is given by

$$IG(X, Y) = H(X) - H(X|Y). \quad (7)$$

The *symmetrical uncertainty* (SU) (a modified information gain measure) is then used to measure the correlation between two discrete variables ( $X$  and  $Y$ ) [49]:

$$SU(X, Y) = 2.0 \times \left( \frac{H(X) - H(X|Y)}{H(X) + H(Y)} \right). \quad (8)$$

Treating each feature as well as the class as random variables, we can apply this as our correlation measure. More specifically, we simply use  $SU(c, i)$  to measure the correlation between a feature  $i$  and a class  $c$ , and  $SU(i, i')$  to measure the correlation between features  $i$  and  $i'$ . These values are then substituted as  $\overline{r_{ci}}$  and  $\overline{r_{ii'}}$  in equation (4).

Recall that our aim here is to find the optimal subset of features which maximises the metric given in equation (4). The size of our feature set is 28 meaning there are  $2^{28} - 1 \simeq 2.7 \times 10^8$  possible subsets, too many for exhaustive search. Instead a greedy stepwise forward selection search strategy was used for searching the space of feature subsets, which works by adding the current best feature at each round. The search begins with the empty set, and in each step the metric, as defined in equation (4), is computed for every single feature addition, and the feature with the best score improvement is added. If at some step none of the remaining features provide an improvement, the algorithm stops, and the current feature set is returned. The best feature subset found with this method

TABLE III  
FEATURE SELECTION BY THE FILTER METHOD

Description	
14	TN <sub>o</sub> I after GB.
13	Number of polynomials after GB.
2	stds before GB.
26	$\log_2(\text{TN}_{oI} \text{ before GB}) - \log_2(\text{TN}_{oI} \text{ after GB})$
21	Proportion of polynomials with $y$ after GB.
15	stds after GB.
23	Proportion of monomials with $x$ after GB.
19	Max degree of $z$ in polynomials after GB.
25	Proportion of monomials with $z$ after GB.
27	$\log_2(\text{stds before GB}) - \log_2(\text{stds after GB})$

(which may not be the absolute optimal subset of features) is shown in Table III, ordered by importance.

### B. The wrapper method

The wrapper feature selection method evaluates attributes using accuracy estimates provided by the target learning algorithm. Evaluation of each feature set was conducted with a SVM with RBF kernel function. The SVM algorithm is run on the dataset, with the same data partitions as described in Section IV-B. Similarly, a five-fold cross validation was carried out. The feature subset with the highest average accuracy was chosen as the final set on which to run the SVM algorithm.

In each training / validation fold, starting with an empty set of features: each feature was added; a model was fitted to the training dataset; the classifier was then tested on the validation set. This was done on all the features, resulting in a score for each reflecting the accuracy of the classifier. The final score for each feature was its average over the five folds. Having obtained a score for all features in the manner above, the feature with the highest score was then added in the feature set. The same greedy procedure as used for the filter method in Section V-A was applied to obtain the best feature subset.

Due to the large number of cases, the parameters  $(C, \gamma)$  were selected from an optimised sub range instead of the full grid search used in Section IV-B. The reduced range suffices to demonstrate the performance of a reduced feature set. In those previous experiments we found that  $C$  taken from  $\{2^5, 2^6, 2^7, 2^8, 2^9, 2^{10}\}$  and  $\gamma$  taken from  $\{2^{-5}, 2^{-6}, 2^{-7}, 2^{-8}, 2^{-9}, 2^{-10}\}$  provided good classifier performance.

The 36 pairs of  $(C, \gamma)$  values were tested and an optimal feature subset with the highest accuracy was found for each. Then the one with the highest accuracy was selected as the final set, which is shown in Table IV ordered by importance. We see that most of the features selected (9, 12 and 22) related to variable  $z$ . Recall that the projection order used in the CAD was always  $x \prec y \prec z$ , i.e. the variable  $z$  is projected first. Hence it makes sense that this variable would have the greatest effect and thus be identified in the feature selection.

We examined the performance on further reduced feature sets, obtained by the feature ranking of the wrapper method. Fig. 1 shows the overall prediction accuracies. For instance, the predictor obtained from only using a single feature (the best ranked feature was TN<sub>o</sub>I after GB for both filter and

TABLE IV  
FEATURE SELECTION BY THE WRAPPER METHOD

Description	
14	TNoI after GB.
9	Proportion of polynomials with $z$ before GB.
22	Proportion of polynomials with $z$ after GB.
4	Max degree of $x$ in polynomials before GB.
12	Proportion of monomials with $z$ before GB.

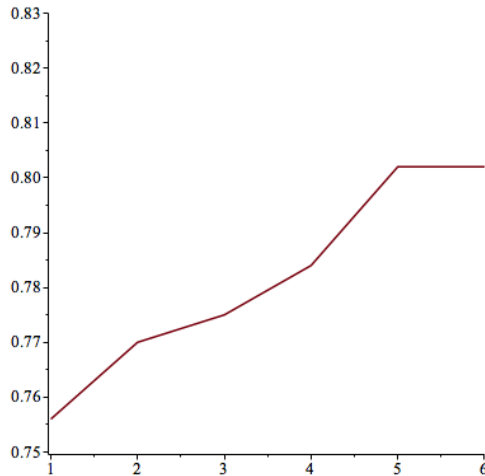


Fig. 1. Performance of a sample run with different sizes of feature sets

wrapper methods) achieved an accuracy score of 0.756 in that run, with the performance steadily increasing with the size of the feature set until the fifth feature. Taking any sixth feature into the set did not improve the performance noticeably, and hence resulted in the cut-off chosen by the wrapper method.

As the wrapper method identified only a few features an error analysis on the misclassified data points is feasible. Fig. 2 shows 40 misclassified points and their features 4 and 14, while Fig. 3 shows the remaining features. It is interesting that feature 4 of all misclassified samples is either 1 or 2, when for the whole dataset roughly a third of samples had this feature value 3 or 4. This indicates that the algorithm performs better on instances with a higher maximum degree of  $x$  among all polynomials before GB preconditioning.

### C. Results with reduced feature sets

Having obtained the reduced feature sets, we ran the experiment again to evaluate the new choices. The dataset was again repartitioned into 80% training and 20% test set, stratified to maintain relative class proportions in both training and test partitions. Again, a five-fold cross validation and a finer grid-search optimisation procedure over the range of  $(C, \gamma)$  pairs was conducted as described previously. The classifier with maximum averaged MCC was selected and the resulting classifier was then evaluated. The testing data was also reduced to contain only the features selected. The classification accuracy was used to measure the performance of the classifier. In order to better estimate the generalisation performance of classifiers with reduced feature sets, the data was permuted and partitioned into 80% training and 20% test again and the whole process was repeated 50 times. For each

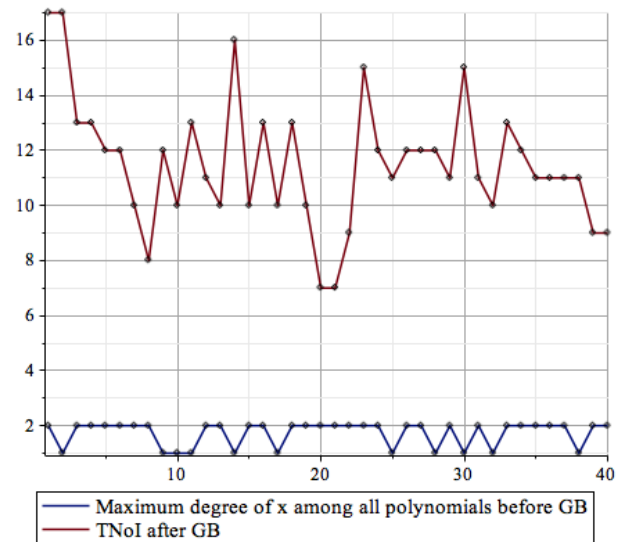


Fig. 2. Feature 4 and 14 of misclassified data

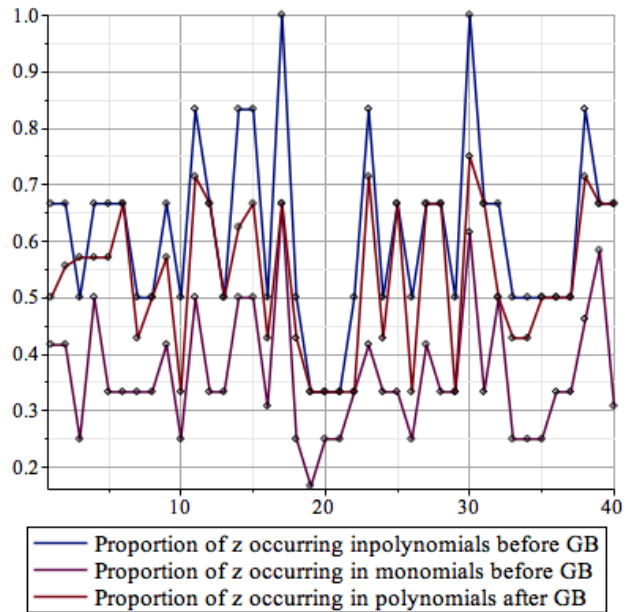


Fig. 3. Feature 9, 12 and 22 of misclassified data

run, each training set was standardised to have zero mean and unit variance, with the same offset and scaling applied subsequently to the corresponding test partition.

Fig. 4 shows boxplots of the accuracies generated by 50 runs of the five-fold cross validation. Both reduced feature sets generated similar results and show a large improvement on the base case where GB preconditioning is always used before CAD construction. The average overall prediction accuracy of the filter subset and the wrapper subset is 79% and 78% respectively (Fig. 1 shows a higher rate but that was just for one sample run). All 50 runs of the wrapper subset performed above the base line, while the top three quartiles of the results of both sets achieve higher than 77% percentage accuracy.

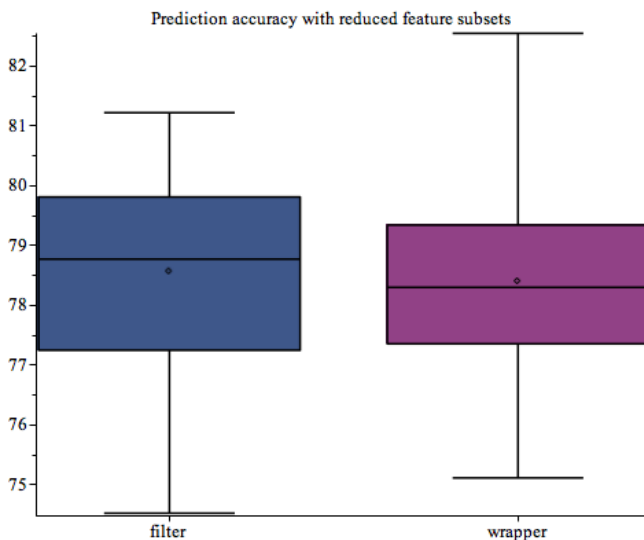


Fig. 4. Boxplots of 50 runs of the 5-fold cross validation with both the suggested feature sets

## VI. CONCLUSION AND FUTURE WORK

### A. Comparison with human developed heuristic

We may compare the machine learned choice with the human developed  $TN_{OI}$  heuristic [57], whose performance on the 213 test problems is shown in Table V. It correctly predicted whether GB preconditioning was beneficial for 118 examples, only 55%. So for this dataset it would have been better on average to precondition blindly than to make a decision using  $TN_{OI}$  alone. The  $TN_{OI}$  heuristic performed better in the experiments by Wilson et al. [57]. Those experiments involved only 22 problems (compared to 213 in the test set here) but they were human constructed to have certain geometric properties while the ones here were random.

We also note that the  $TN_{OI}$  heuristic actually performed quite differently for positive and negative examples of our dataset, as shown by the separated data in Table V. It was able to identify most of the cases where GB-preconditioning is detrimental but failed to identify many of the cases where it was beneficial. The  $TN_{OI}$  after GB was identified as important by both feature methods, but it seems to need to be used in conjunction with other features to be effective here.

### B. Summary

We investigated the application of machine learning to the problem of predicting when GB preconditioning is beneficial to CAD. We had to create a new dataset of random polynomials for the experiment. We acknowledge that it would be preferable to run such a test on an established dataset but this was not available. Of course, such a random datasets could be enlarged to increase variety almost indefinitely, but we needed to keep the experiment within computationally feasible boundaries. We emphasise the interesting initial finding in Section II-B that supposedly varied established sets can have hidden uniformity; and highlight that our generated dataset matched

TABLE V  
THE PERFORMANCE OF THE  $TN_{OI}$ -BASED HEURISTIC [57]

	Total	Correct Prediction	
GB beneficial	158	77	(48%)
GB not beneficial	54	41	(76%)
Total	213	118	(55%)

previously reported results [15], [57] for the topic of study with most, but not all, benefiting from GB preconditioning.

A machine learned choice on whether to precondition was found to yield better results than either always preconditioning blindly, or using the previously human developed  $TN_{OI}$  heuristic [57]. Two feature selection experiments showed that a small feature subset could be used. The two subsets identified were different but both needed features from before and after the GB preconditioning. For one, having fewer features actually improved the learning efficiency to 79%. Although a modest improvement on applying GB preconditioning blindly, we recall that the wrong choice can give large changes to the size of the CAD or even change the tractability of the problem.

### C. Future Work

There are many possible extensions to this project:

- To see how the learned choice performs on a non-random, dataset. There is a large set derived from university mathematics entrance exams [44], which is not yet publicly available but may be in the future.
- There are further CAD optimisations for multiple equalities under development [27], [29], [24] which may affect the role of GB preconditioning from CAD.
- In the present paper the variable ordering for CAD and the monomial ordering for GB were fixed. In reality, such decisions may also need to be made; but the variable ordering can affect the choice of whether to use GB preconditioning and *vice versa*.

Finally, we note there are other algorithm optimisation decisions for CAD, and indeed elsewhere in computer algebra.

### Acknowledgements

Thanks to David Wilson and James Bridge, our collaborators on [41], for useful conversations on the topic of machine learning to optimise computer algebra. This work was supported by EPSRC grant EP/J003247/1 and EU H2020-FETOPEN-2016-2017-CSA project  $\mathcal{SC}^2$  (712689).

## REFERENCES

- [1] D. Arnon, G.E. Collins, and S. McCallum. Cylindrical algebraic decomposition I: The basic algorithm. *SIAM J. Computing*, 13:865–877, 1984.
- [2] P. Baldi, S. Brunak, Y. Chauvin, C.A.F. Andersen, and H. Nielsen. Assessing the accuracy of prediction algorithms for classification: An overview. *Bioinformatics*, 16(5):412–424, 2000.
- [3] S. Basu, R. Pollack, and M.F. Roy. *Algorithms in Real Algebraic Geometry*. Vol. 10 of Algorithms & Comp. in Math. Springer, 2006.
- [4] W. Böge, R. Gebauer, and H. Kredel. Gröbner bases using SAC2. In *EUROCAL '85*, (LNCS 204), pages 272–274. Springer, 1985.
- [5] R. Bradford, C. Chen, J.H. Davenport, M. England, M. Moreno Maza, and D. Wilson. Truth table invariant cylindrical algebraic decomposition by regular chains. In *Computer Algebra in Scientific Computing*, (LNCS 8660), pages 44–58. Springer International Publishing, 2014.

- [6] R. Bradford, J.H. Davenport, M. England, S. McCallum, and D. Wilson. Cylindrical algebraic decompositions for boolean combinations. In *Proc. ISSAC '13*, pages 125–132. ACM, 2013.
- [7] R. Bradford, J.H. Davenport, M. England, S. McCallum, and D. Wilson. Truth table invariant cylindrical algebraic decomposition. *J. Symbolic Computation*, 76:1–35 2016.
- [8] R. Bradford, J.H. Davenport, M. England, and D. Wilson. Optimising problem formulations for cylindrical algebraic decomposition. In *Intelligent Computer Mathematics*, (LNCS 7961), pages 19–34. Springer Berlin Heidelberg, 2013.
- [9] J.P. Bridge, S.B. Holden, and L.C. Paulson. Machine learning for first-order theorem proving. *J. Automated Reasoning*, pages 1–32, 2014.
- [10] C.W. Brown. Improved projection for cylindrical algebraic decomposition. *J. Symbolic Computation*, 32(5):447–465, 2001.
- [11] C.W. Brown. Constructing a single open cell in a cylindrical algebraic decomposition. In *Proc. ISSAC '13*, pages 133–140. ACM, 2013.
- [12] C.W. Brown and J.H. Davenport. The complexity of quantifier elimination and cylindrical algebraic decomposition. In *Proc. ISSAC '07*, pages 54–60. ACM, 2007.
- [13] C.W. Brown, M. El Kahoui, D. Novotni, and A. Weber. Algorithmic methods for investigating equilibria in epidemic modelling. *J. Symbolic Computation*, 41:1157–1173, 2006.
- [14] B. Buchberger. PhD thesis (1965): An algorithm for finding the basis elements of the residue class ring of a zero dimensional polynomial ideal. *J. Symbolic Computation*, 41(3-4):475–511, 2006.
- [15] B. Buchberger and H. Hong. Speeding up quantifier elimination by Gröbner bases. Tech. Report, 91-06. RISC, Johannes Kepler Uni., 1991.
- [16] C. Chen and M. Moreno Maza. Real quantifier elimination in the RegularChains library. In *Mathematical Software – ICMS 2014*, (LNCS 8592), pages 283–290. Springer Heidelberg, 2014.
- [17] C. Chen and M. Moreno Maza. Quantifier elimination by cylindrical algebraic decomposition based on regular chains. *J. Symbolic Computation*, 75:74–93, 2016.
- [18] C. Chen, M. Moreno Maza, B. Xia, and L. Yang. Computing cylindrical algebraic decomposition via triangular decomposition. In *Proc. ISSAC '09*, pages 95–102. ACM, 2009.
- [19] G.E. Collins. Quantifier elimination for real closed fields by cylindrical algebraic decomposition. In *Proc. 2nd GI Conference on Automata Theory and Formal Languages*, pages 134–183. Springer-Verlag, 1975.
- [20] G.E. Collins. The SAC-2 computer algebra system. In *EUROCAL '85*, (LNCS 204), pages 34–35. Springer, 1985.
- [21] G.E. Collins and H. Hong. Partial cylindrical algebraic decomposition for quantifier elimination. *J. Symbolic Computation*, 12:299–328, 1991.
- [22] J.H. Davenport, R. Bradford, M. England, and D. Wilson. Program verification in the presence of complex numbers, functions with branch cuts etc. In *Proc. SYNASC '12*, pages 83–88. IEEE, 2012.
- [23] J.H. Davenport and J. Heintz. Real quantifier elimination is doubly exponential. *J. Symbolic Computation*, 5(1-2):29–35, 1988.
- [24] J.H. Davenport and M. England. Need polynomial systems be doubly exponential? In *Mathematical Software – ICMS 2016*, (LNCS 9725), pages 157–164. Springer International Publishing, 2016.
- [25] A. Dolzmann, A. Seidl, and T. Sturm. Efficient projection orders for CAD. In *Proc. ISSAC '04*, pages 111–118. ACM, 2004.
- [26] M. England, R. Bradford, C. Chen, J.H. Davenport, M. Moreno Maza, and D. Wilson. Problem formulation for truth-table invariant cylindrical algebraic decomposition by incremental triangular decomposition. In *Intelligent Computer Mathematics*, (LNCS 8543), pages 45–60. Springer, 2014.
- [27] M. England, R. Bradford, and J.H. Davenport. Improving the use of equational constraints in cylindrical algebraic decomposition. In *Proc. ISSAC '15*, pages 165–172. ACM, 2015.
- [28] M. England, R. Bradford, J.H. Davenport, and D. Wilson. Choosing a variable ordering for truth-table invariant cylindrical algebraic decomposition by incremental triangular decomposition. In *Mathematical Software – ICMS 2014*, (LNCS 8592), pages 450–457. Springer, 2014.
- [29] M. England and J.H. Davenport. The complexity of cylindrical algebraic decomposition with respect to polynomial degree. In *Computer Algebra in Scientific Computing*, (LNCS 9890), pages 172–192. Springer, 2016.
- [30] M. Erascu and H. Hong. Synthesis of optimal numerical algorithms using real quantifier elimination (Case Study: Square root computation). In *Proc. ISSAC '14*, pages 162–169. ACM, 2014.
- [31] J.C. Faugère. A new efficient algorithm for computing groebner bases without reduction to zero (F5). In *Proc. ISSAC '02*, pages 75–83. ACM, 2002.
- [32] U.M. Fayyad and K.B. Irani. Multi-interval discretization of continuous-valued attributes for classification learning. In *Proc. of the International Joint Conference on Uncertainty in AI*, pages 1022–1027, 1993. Available at: <http://trs-new.jpl.nasa.gov/dspace/handle/2014/35171>,
- [33] I.A. Fotiou, P.A. Parrilo, and M. Morari. Nonlinear parametric optimization using cylindrical algebraic decomposition. In *Proc. CDC-ECC '05*, pages 3735–3740, 2005.
- [34] I. Guyon and A. Elisseeff. An introduction to variable and feature selection. *J. Machine Learning Research*, 3:1157–1182, 2003.
- [35] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reutemann, and I.H. Witten. The WEKA data mining software: An update. *SIGKDD Explorations Newsletter*, 11(1):10–18, 2009.
- [36] M.A. Hall. Correlation-based feature selection for discrete and numeric class machine learning. In *Proc. of the Seventeenth International Conference on Machine Learning*, ICML '00, pages 359–366. Morgan Kaufmann Publishers Inc., 2000.
- [37] M.A. Hall and G. Holmes. Benchmarking attribute selection techniques for discrete class data mining. *IEEE Transactions on Knowledge and Data Engineering*, 15(6):1437–1447, 2003.
- [38] J. Han, L. Dai, and B. Xia. Constructing fewer open cells by gcd computation in CAD projection. In *Proc. ISSAC '14*, pages 240–247. ACM, 2014.
- [39] H. Hong. An improvement of the projection operator in cylindrical algebraic decomposition. In *Proc. ISSAC '90*, pages 261–264. ACM, 1990.
- [40] C. Hsu, C. Chang, and C. Lin. A practical guide to support vector classification. Tech. Report, Department of Computer Science, National Taiwan Uni., 2003.
- [41] Z. Huang, M. England, D. Wilson, J.H. Davenport, L. Paulson, and J. Bridge. Applying machine learning to the problem of choosing a heuristic to select the variable ordering for cylindrical algebraic decomposition. In *Intelligent Computer Mathematics*, (LNCS 8543), pages 92–107. Springer, 2014.
- [42] H. Iwane, H. Yanami, H. Anai, and K. Yokoyama. An effective implementation of a symbolic-numeric cylindrical algebraic decomposition for quantifier elimination. In *Proc. SNC '09*, pages 55–64, 2009.
- [43] T. Joachims. Making large-scale support vector machine learning practical. In B. Schölkopf, C.J.C. Burges, and A.J. Smola, editors, *Advances in Kernel Methods*, pages 169–184. MIT Press, 1999.
- [44] M. Kobayashi, H. Iwane, T. Matsuzaki, and H. Anai. Efficient subformula orders for real quantifier elimination of non-prenex formulas. In *Mathematical Aspects of Computer and Information Sciences (MACIS '15)*, (LNCS 9582), pages 236–251. Springer, 2016.
- [45] B.W. Matthews. Comparison of the predicted and observed secondary structure of T4 phage lysozyme. *Biochimica et Biophysica Acta (BBA)-Protein Structure*, 405(2):442–451, 1975.
- [46] E.W. Mayr and A.R. Meyer. The complexity of the word problems for commutative semigroups and polynomial ideals. *Advances in Mathematics*, 46(3):305–329, 1982.
- [47] S. McCallum. An improved projection operation for cylindrical algebraic decomposition. In B. Caviness and J. Johnson, editors, *Quantifier Elimination and Cylindrical Algebraic Decomposition*, Texts & Mono. Symbolic Computation, pages 242–268. Springer-Verlag, 1998.
- [48] L.C. Paulson. Metitarski: Past and future. In *Interactive Theorem Proving*, (LNCS 7406), pages 1–10. Springer, 2012.
- [49] W.H. Press, S.A. Teukolsky, W.T. Vetterling, and B.P. Flannery. *Numerical Recipes in C (2nd Ed.): The Art of Scientific Computing*. Cambridge Uni. Press, 1992.
- [50] J.R. Quinlan. Induction of decision trees. *Machine Learning*, 1(1):81–106, 1986.
- [51] B. Schölkopf, K. Tsuda, and J.-P. Vert. *Kernel methods in computational biology*. MIT Press, 2004.
- [52] Claude E. Shannon. A mathematical theory of communication. *Mobile Computing and Communications Review*, 5(1):3–55, 2001. Reprint of original 1948 publication.
- [53] J. Shawe-Taylor and N. Cristianini. *Kernel methods for pattern analysis*. Cambridge Uni. Press, 2004.
- [54] A. Strzeboński. Cylindrical algebraic decomposition using validated numerics. *J. Symbolic Computation*, 41(9):1021–1038, 2006.
- [55] A. Strzeboński. Cylindrical algebraic decomposition using local projections. In *Proc. ISSAC '14*, pages 389–396. ACM, 2014.
- [56] D.J. Wilson, R.J. Bradford, and J.H. Davenport. A repository for CAD examples. *ACM Comm. Computer Algebra*, 46(3):67–69, 2012.
- [57] D.J. Wilson, R.J. Bradford, and J.H. Davenport. Speeding up cylindrical algebraic decomposition by Gröbner bases. In *Intelligent Computer Mathematics* (LNCS 7362), pages 280–294. Springer, 2012.
- [58] The benchmarks used in solving nonlinear arithmetic. Available from: <http://cs.nyu.edu/~dejan/nonlinear/>