

Learning to Choose Branching Rules for Nonconvex MINLPs

Timo Berthold¹[0000-0002-6320-8154] and Fritz Geis²

¹ Fair Isaac Germany GmbH, Takustr. 7, 14195 Berlin, Germany
timoberthold@fico.com

² Zuse Institute Berlin, Takustr. 7, 14195 Berlin, Germany
fritzgeis@gmail.com

Abstract. Outer-approximation-based branch-and-bound is a common algorithmic framework for solving MINLPs (mixed-integer nonlinear programs) to global optimality, with branching variable selection critically influencing overall performance. In modern global MINLP solvers, it is unclear whether branching on fractional integer variables should be prioritized over spatial branching on variables, potentially continuous, that show constraint violations, with different solvers following different defaults. We address this question using a data-driven approach. Based on a test set of hundreds of heterogeneous public and industrial MINLP instances, we train linear and random forest regression models to predict the relative speedup of the FICO® Xpress Global solver when using a branching rule that always prioritizes variables with violated integralities versus a mixed rule, allowing for early spatial branches.

We introduce a practical evaluation methodology that measures the effect of the learned model directly in terms of the shifted geometric mean runtime. Using only four features derived from strong branching and the nonlinear structure, our linear regression model achieves an 8–9% reduction in geometric-mean solving time for the Xpress solver, with over 10% improvement on hard instances. We also analyze a random regression forest model. Experiments across solver versions show that a model trained on Xpress 9.6 still yields significant improvements on Xpress 9.8 without retraining.

Our results demonstrate how regression models can successfully guide the branching-rule selection and improve the performance of a state-of-the-art commercial MINLP solver.

Keywords: Nonlinear Optimization · Machine Learning · Branching.

1 Introduction

We consider *MINLPs* (*mixed-integer nonlinear programs*) of the form

$$\min\{c^T x \mid g_k(x) \leq 0, \forall k \in \mathcal{K}, l \leq x \leq u, x_j \in \mathbb{Z}, \forall j \in \mathcal{J}\}, \quad (1)$$

where all constraint functions $g_k : \mathbb{R}^n \rightarrow \mathbb{R}$ are factorable and all variable bounds $l, u \in \mathbb{R} := \mathbb{R} \cup \{\pm\infty\}$. The set $\mathcal{K} = \{1, \dots, m\}$, $m \in \mathbb{N}$, indexes the constraints

and $\mathcal{J} \subseteq \{1, \dots, n\}$ the integer variables. A nonlinear objective can be easily modeled by introducing an auxiliary variable and an objective-transfer constraint, see, e.g., [?]. If all g_k are linear, and $\mathcal{J} = \emptyset$ we call (??) a *linear program (LP)*. This work focuses on *nonconvex MINLPs*, i.e., problems of form (??), where at least one g_k is nonconvex.

Note that factorable functions can be represented via a directed acyclic *expression graph*, with nodes representing operators or variables, and arcs representing the data flow of the computation. In this paper, we refer to this representation as the *DAG*, for a good overview on the use of the DAG in MINLP solving, we recommend [?].

For solving problems of the form (??), we use the FICO[®] Xpress Global [?] MINLP solver, which we will refer to as *Xpress*. Xpress is based on the *branch-and-bound* method (*B&B*), which recursively partitions the problem by splitting the domain of selected variables, which is called *branching*. Selecting good branching variables is crucial for the performance of B&B-based MINLP solvers, see, e.g., [?]. For more details on the implementation in Xpress, see [?].

This paper studies a fundamental question: should we always branch on fractional integer variables first, or consider spatial branches even when there are fractional integers? Unlike most prior work on using ML for branching [?, ?, ?, ?, ?], we do not learn individual branching decisions or attempt to mimic existing strategies such as strong branching; instead, we perform algorithm selection between two established branching rules. Different from most prior work, we consider a heterogeneous set of instances. The resulting model can be integrated directly into solver code and does not require any pre-training on the user side. This is akin to prior ML-based algorithm-selection work to choose between scaling procedures [?], local cut selection rules [?], linearization techniques [?], or spatial branching strategies in the context of RLT for polynomial optimization [?], respectively, and deliberately different from solver-free learning approaches for MINLP as in [?].

2 A Quick Recap of Branching for MINLPs

The B&B algorithm recursively partitions the problem into smaller subproblems (*branching*) and solves LP relaxations to obtain bounds (*bounding*) until an optimal solution or infeasibility proof is found.

An *LP relaxation* of an MINLP is obtained by dropping integrality constraints and replacing nonlinear constraints with linear underestimators where possible. This relaxation is successively strengthened by *outer-approximation cuts* [?]. A well-designed cutting plane separation procedure often helps to reduce the branch-and-bound tree size while accelerating the overall solving process [?]. Unlike in MIP solving, cutting planes are often additionally separated immediately during branching-node creation in MINLPs.

In this paper, we focus on *variable branching*, in which the domain of a single variable is split into two or more intervals.

Two key types of variable branching are:

1. *Integer branching*, which is applied when an integer variable has a fractional value $\check{x}_j \in \mathbb{R} \setminus \mathbb{Z}$ in the solution \check{x} of the current LP relaxation. Two subproblems are created that enforce $x_j \leq \lfloor \check{x}_j \rfloor$ and $x_j \geq \lceil \check{x}_j \rceil$, respectively.
2. *Spatial branching* [?] is applied when the violation of a nonconvex constraint cannot be resolved by an outer-approximation cut, but requires partitioning variable domains. Spatial branching candidates are often continuous variables, but can also include integer variables whose LP value happens to be integral. Two created subproblems enforce $x_j \leq \lfloor \check{x}_j \rfloor$ and $x_j \geq \lceil \check{x}_j \rceil$, respectively, for a branching point $\check{x}_j \in \mathbb{R}$. Though the LP solution is not explicitly excluded, subsequent outer-approximation cuts typically remove it.

3 Machine Learning Methodological Approach

Learning Task/Feature Space Our learning task consists of choosing, after root node processing and right before the first branch, one of two rules of how to combine integer branching and spatial branching for the remainder of the branch-and-bound search: Either, always branch on integer candidates and conduct spatial branches only when there is no integer branching candidate, which we will refer to as "PreferInt" (this is the default, e.g., in the SCIP MINLP solver). Or mix both candidate sets and always allow the choice of either type of candidate (which is the default, e.g., in the Xpress solver), which we refer to as "Mixed".³

Although this is inherently a binary decision, we frame it as a regression problem. This choice is motivated by two considerations. Firstly, our ultimate goal is to improve the average runtime of the solver, which is a metric that is numerical and not categorical. Secondly, our focus is on getting the prediction right for those instances on which the performance of selecting "Mixed" and "PreferInt" significantly differs, see also [?]. Regression allows us to model the magnitude of this difference directly and thereby focus the learning on the cases where it matters most.

To this end, we train regression models $y_i : \mathbb{R}^d \rightarrow \mathbb{R}$ that map a d -dimensional feature vector $f = (f_1, \dots, f_d)$ onto the speedup or slowdown factor (the *label*) in runtime by using "PreferInt" instead of "Mixed". We initially used 17 features, see Table ??.

This includes features related to strong branching at the root node, such as the average change in the dual bound resulting from integer and spatial strong branching, `AvgRelBndChngSBLPInt` and `AvgRelBndChngSBLPSpat`, respectively, the number of variables fixed from strong branching on spatial branching candidates `#SpatBranchEntFixed`⁴, and the amount of deterministic *work* invested in either strong branching, `AvgWorkSBLPInt` and `AvgWorkSBLPSpat`. Work is a deterministic measure of computational effort implemented in Xpress. These features give us an indication of how effective (and expensive) strong branching on

³ We ruled out always preferring spatial branches in a preliminary experiment, since this option was a factor eight slower on average and rarely won against the others.

⁴ There were only a few instances where integer strong branching fixed variables; hence, a corresponding feature would have been almost flat zero.

integer or spatial variables is. Relatedly, `#IntViols` and `#NonlinViols` refer to the number of integer and spatial branching candidates.

As problem structure features, we include the percentage of variables that are integer, `%IntVars`, the percentage of constraints that are equations, `%EqCons`, the ratio of quadratic elements in the problem to variables, `%QuadrElements`, and the percentage of constraints that contain nonlinearities, `%NonlinCons`. Further, to measure the nonlinearity of the problem, we use information about the DAG, in particular, the percentage of variables that are part of any nonlinearity, `%VarsDAG`, and the ratio between nodes in the DAG and nonzeros in the linear part of the problem, `NodesInDAG`. We further include the percentage of integer and unbounded variables among all variables in the DAG, `%VarsDAGInt` and `%VarsDAGUnbnd`, as for integer DAG variables, we "hit two birds with one stone" and branching on unbounded variables can be crucial to get efficient dual bounds. Finally, we consider `%QuadrNodesDAG` to measure whether the nonlinearities in the problem are mostly quadratic.

| Feature | Feature Scaling |
|-------------------------------------|---|
| <i>Problem Structure</i> | |
| <code>%QuadrElements</code> | number quadratic elements over n |
| <code>%IntVars</code> | #Integer variables after presolve over \tilde{n} |
| <code>%EqCons</code> | #equality constraints over m |
| <code>%NonlinCons</code> | #nonlinear constraints over m |
| <i>Effect of Branching</i> | |
| <code>#IntViols</code> | |
| <code>#NonlinViols</code> | |
| <code>#SpatBranchEntFixed</code> | |
| <code>AvgWorkSBLPInt</code> | |
| <code>AvgWorkSBLPSpat</code> | $\log_{10}(\text{Value}+1)$ |
| <code>AvgRelBndChngSBLPInt</code> | |
| <code>AvgRelBndChngSBLPSpat</code> | |
| <code>AvgCoeffSpreadConvCuts</code> | |
| <i>DAG</i> | |
| <code>NodesInDAG</code> | NodesInDAG over $\text{NodesInDAG} + \tilde{M}$ |
| <code>%VarsDAG</code> | #vars in DAG over \tilde{n} |
| <code>%VarsDAGUnbnd</code> | #unbounded vars over #vars in DAG |
| <code>%VarsDAGInt</code> | #integer vars over #vars in DAG |
| <code>%QuadrNodesDAG</code> | #quadratic operator nodes in DAG over all nonlinear operator nodes in DAG |

Table 1. m and n are the number of constraints and variables before presolving, respectively; \tilde{n} and \tilde{M} the number of variables and linear nonzeros after presolving.

Data The data on which the models are trained comes from running Xpress 9.6⁵ twice on a heterogeneous benchmark of 683 public and industrial MINLP instances, each with two permutations to mitigate the effect of performance variability [?, ?], yielding 2049 data points. For each instance, we record the runtimes produced by both branching rules and the complete feature set. Instances solved at the root or otherwise unsuitable for comparison are filtered out, resulting in a final dataset of 797 data points, see [?] for details. Solving at the root node was by far the most common reason for filtering.

Training For training the models, we split the data randomly into 80% training and 20% test set. The models we train on the training set are a linear regressor [?] and a random forest regressor, RF, [?]. We use the python library scikit-learn [?], which provides us with the linear regressor by the function *LinearRegression* and the random forest regressor by the function *RandomForestRegressor*.

Testing Instead of training one linear regressor and one random forest regressor, we opted for training and testing one hundred models each with different random seeds and average their performance scores to evaluate how promising this ML-based approach is.

To measure the performance of the regression models, we use the *accuracy* and the shifted geometric mean of the runtime (*sgm_runtime*). The accuracy is defined as the percentage of times the model predicted the faster rule. The *sgm_runtime* is the shifted geometric mean of runtimes when solving each test instance using the predicted branching rule over the shifted geometric mean time using always the default rule.

Hence, accuracy is always between 0 and 100%, with larger values being better. *Sgm_runtime* can be smaller or larger than 1, with values larger than 1 indicating a deterioration and values smaller than 1 indicating an improvement: the smaller the number, the better. This is the primary performance indicator for solver development in practice.

To compute the shifted geometric mean [?] with a shift of 10, measurements $X = (X_1, \dots, X_n)$ are aggregated via $sgm(X) = -10 + \prod_{i=1}^n (X_i + 10)^{\frac{1}{n}}$. The use of the shifted geometric mean is a commonly used method to aggregate performance measures, in particular running time, in mathematical optimization [?].

For the linear model, *feature importance* is given by the absolute value of the learned coefficients, whereas for the random forest it is measured as the normalized total reduction in mean squared error induced by splits on that feature (mean decrease impurity, MDI), which are the default importance metrics in scikit-learn. For each random seed and each type of model, we computed the feature importance for all features and sorted them from most important to least important. Then we assigned a score of zero to the most important

⁵ More precisely: An internal version of the Xpress 9.6 that exposes those features that are otherwise not available as public attributes.

one, a score of one to the second most important one, and so on. Finally, we added, for each model type, the scores across all one hundred runs together. The four most important features per model type (as by this score sum) are listed in Table ???. Although the top-ranked features differ between the linear regression and random forest models, there is overlap in terms of the underlying information captured. In particular, `AvgRelBndChngSBLPSpat` ranks first for the linear model and second for the random forest, and `%NonlinCons`, ranked third for the linear model, is a very close fifth for the random forest. Overall, five of the eight highest-ranked features coincide across the two models. Differences are expected given the different nature of the models: linear regression emphasizes globally predictive, approximately linear effects with low collinearity, whereas random forests prioritize features that enable strong local splits and nonlinear interactions, for instance, capturing cases where a branching rule is beneficial for either extreme but not for intermediate feature values.

Further Approaches In earlier versions, we tested the algorithm with different features, unscaled or differently scaled features, on an earlier version of Xpress and for the SCIP solver (where it also improved performance, but not as much as in the Xpress case). Details can be found in the thesis [?]. This thesis also contains a detailed description of how we selected the scaler and imputer for the data set and a discussion of restricting the decision tree depth to five in the random forest models.

| Ranking | Linear | Forest |
|---------|------------------------------------|-------------------------------------|
| 1. | <code>AvgRelBndChngSBLPSpat</code> | <code>AvgCoeffSpreadConvCuts</code> |
| 2. | <code>%IntVars</code> | <code>AvgRelBndChngSBLPSpat</code> |
| 3. | <code>%NonlinCons</code> | <code>#NonlinViols</code> |
| 4. | <code>%VarsDAGInt</code> | <code>%EqCons</code> |

Table 2. Four most important features for either model type.

4 Computational Experiments

Our computational experiments consist of three parts: Firstly, we evaluate the regression models trained on the full 17-feature set, and then examine how their performance evolves as we iteratively remove the least important features. This reduction process provides insights into which features drive prediction quality and whether a more compact feature set can yield comparable performance with presumably better robustness. Secondly, we analyze the final reduced models in more detail and provide an analysis of their performance with respect to accuracy and runtime. Finally, we compare how those models continue to perform as the underlying branch-and-bound method is improved (in this case, through a solver version update).