

Performance and optimization of support vector machines in high-energy physics classification problems

M.Ö. Sahin¹ D. Krücker and I.-A. Melzer-Pellmann

*Deutsches Elektronen-Synchrotron,
Notkestr. 85, 22607 Hamburg, Germany*

E-mail: ozgur.sahin@desy.de, dirk.kruecker@desy.de,
isabell.melzer@desy.de

ABSTRACT: In this paper we promote the use of Support Vector Machines (SVM) as a machine learning tool for searches in high-energy physics. As an example for a new-physics search we discuss the popular case of Supersymmetry at the Large Hadron Collider. We demonstrate that the SVM is a valuable tool and show that an automated discovery-significance based optimization of the SVM hyper-parameters is a highly efficient way to prepare an SVM for such applications. A new C++ LIBSVM interface called SVM-HINT is developed and available on Github.

¹Corresponding author.

Contents

1	Introduction	1
2	Support vector machines	2
2.1	Linearly separable distributions	3
2.2	Overlapping distributions	5
2.3	Non-linear distributions	5
2.4	Probabilistic output	6
3	Hyper-parameter tuning	7
3.1	Performance measures	7
3.2	Hyper-parameter search	8
4	Case studies	8
4.1	Performance comparison on a toy model	8
4.2	Third generation supersymmetric partner search	9
5	Conclusions	15
A	Iterative grid search	16
B	Boosted Decision Trees	17

1 Introduction

Data analysis in High Energy Physics (HEP) is a genuine multivariate problem. Despite the fact that multivariate techniques have been used in HEP for a long time, the explosive growth of machine learning (ML) techniques during the last two decades had only a limited impact on the accustomed style in which data analysis is performed in this field. TMVA [1] is probably the most commonly used software package in this context and especially Boosted Decision Trees (BDT) and Artificial Neural Networks are applied to explore the large and complex datasets delivered by present-day experiments. Only recently an increased interest in the machine learning expertise acquired in other areas of science can be observed [2–4].

In this paper we promote the application of Support Vector Machines (SVM) [5–7] for new-physics searches. Support Vector Machines are a competitive and widely used approach to binary classification. The search for new physics can be considered as a classification task where the rare new-physics signal and the dominant Standard Model (SM) background constitute the two distinct classes. Although there are a few HEP papers on SVMs [8–12], this approach seems to be heavily undervalued amongst HEP researchers considering the many thousands of publications on SVM applications that a simple literature search yields.

After an introduction to SVMs in Section 2, the hyper-parameter tuning is described in Section 3 (and Appendix A) within the context of our new SVM framework: SVM-HINT. In Section 4 we first discuss a toy model and then an example for an actual new-physics search, targeted at a supersymmetric partner of the top quark at the LHC. We demonstrate that the SVM is a valuable tool for HEP searches and that the partial neglect of the SVM approach within the HEP community can be related to the limited performance of its implementation in TMVA without an automated hyper-parameter search. Moreover, we show that a significance-based optimization of the hyper-parameters is a highly efficient way to prepare an SVM for a HEP search.

In addition, we provide a software package [13] that performs such a significance-based optimization of hyper-parameters and interfaces ROOT [14] trees with the popular SVM implementation LIBSVM [15].

2 Support vector machines

A HEP search typically starts with a set of physical variables and cuts on these variables. The cuts are defined to select a new-physics signal against the background of known physics and are often chosen in a more or less ad-hoc style. The optimal use of these variables is a typical machine learning problem. Monte Carlo (MC) simulation samples for the signal and background class can be used to train a supervised¹ machine learning algorithm which is potentially a more efficient classifier than a set of cuts.

We write for a set of N training events:

$$(y_1, \mathbf{x}_1), (y_2, \mathbf{x}_2), \dots, (\mathbf{x}_i, y_i), \dots, (y_N, \mathbf{x}_N) \quad y_i \in \{-1, 1\}, \quad (2.1)$$

$$\mathbf{x}_i = (x_i^{(1)}, \dots, x_i^{(n)}) \quad (2.2)$$

where for an event $i = 1 \dots N$ the label y_i distinguishes between signal and background and \mathbf{x}_i is an n -dimensional vector formed from the physical variables under consideration. These vectors constitute an n -dimensional real vector space \mathbb{V} .

A support vector machine is a supervised binary classifier based on the intuitive concept of an n -dimensional hyperplane separating two distinct classes. In this approach, finding the best separating hyperplane is considered to be a convex optimization problem. In its simplest form a SVM defines the eponymous *support vectors* as those elements of the training sample which are closest to the hyperplane. The separation margin between the classes is completely defined by the support vectors and maximized by the algorithm. This idea can be extended to overlapping distributions and eventually, by an implicit transformation of the variables, known as the *Kernel trick*, to non-linear problems. The last two modifications introduce additional hyper-parameters that must be set to some best value before the SVM training. In the following we give a short introduction to the concepts behind the SVM algorithm and to the hyper-parameter tuning. The reader who is more interested in applications may continue with Section 4.

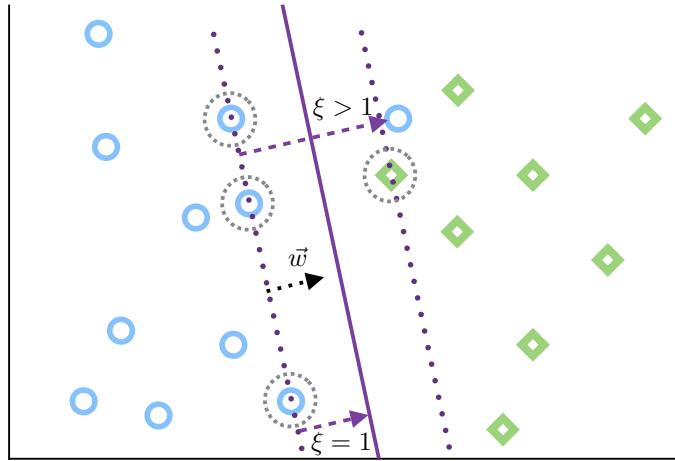


Figure 1. The events represented by blue circles belong to the first class ($y = -1$), whereas the green triangles belong to the second class ($y = 1$). The dashed lines represent the maximum margin boundaries, and the corresponding support vectors are circled by dashed lines. From all possible hyperplanes dividing the two samples, the one with the largest margin is chosen. The blue circle at $\xi > 1$ is not linearly separable, see Sec. 2.2.

2.1 Linearly separable distributions

A linear SVM separates the elements of two classes by an optimal hyperplane. Optimal in this approach is a hyperplane that maximizes the margin between the two classes for a given training sample. Those elements of the training sample sitting on the maximum margin boundaries are called support vectors. The support vectors are sufficient to construct the optimal hyperplane. Fig. 1 illustrates these ideas.

A separating hyperplane can be described as $\mathbf{w} \cdot \mathbf{x} + b = 0$, $\mathbf{w} \in \mathbb{V}$ and $b \in \mathbb{R}$. The vectors of the training sample are either *above* or *below*² this plane. We can always choose the scale of \mathbf{w} and b such that for the vectors which are closest to the hyperplane, i.e. the support vectors \mathbf{x}_k , we obtain $\mathbf{w} \cdot \mathbf{x}_k + b = \pm 1$. Multiplied with the class label y_i , this expression must always be positive for correctly classified points:

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 \geq 0, \quad (2.3)$$

and the equality is satisfied by the support vectors (circled in Fig. 1). The separation margin is the distance ρ between the support vectors on both sides. With the normal vector to the hyperplane $\mathbf{w}/|\mathbf{w}|$ and two arbitrary support vectors from each side $\mathbf{x}_{k+}, \mathbf{x}_{k-}$ the margin is given by:

$$\rho(\mathbf{w}, b) = \frac{\mathbf{w} \cdot \mathbf{x}_{k+}}{|\mathbf{w}|} - \frac{\mathbf{w} \cdot \mathbf{x}_{k-}}{|\mathbf{w}|} = \frac{2}{|\mathbf{w}|}, \quad (2.4)$$

where the second equality follows from Eq. 2.3. Maximizing the margin $\rho = 2/|\mathbf{w}|$ is equivalent to minimizing $|\mathbf{w}|^2$. Finding the optimal separating hyperplane is therefore identical

¹A ML algorithm is called *supervised* when the class membership of all training vectors is known.

²For simplicity we use a 3 dimensional way of speaking. All described concepts are valid in n dimensions.

to solving the following quadratic optimization problem where the correct classification is enforced by the constraints from Eq. 2.3:

$$\begin{aligned} \min_{\mathbf{w} \in \mathbb{V}, b \in \mathbb{R}} \quad & \frac{1}{2} |\mathbf{w}|^2 \\ \text{subject to} \quad & y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 \text{ for all } i = 1 \dots N . \end{aligned} \quad (2.5)$$

This is a quadratic optimization problem with inequality constraints and can be solved using Lagrange multipliers α_i (with $\alpha_i \geq 0$). The Lagrangian can be written as:

$$\mathcal{L} = \frac{1}{2} |\mathbf{w}|^2 - \sum_{i=1}^N \alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1] . \quad (2.6)$$

The solution is a saddle point $(\mathbf{w}^0, b^0, \alpha_i^0)$ where the Lagrangian becomes minimal with respect to \mathbf{w} and b and where the derivatives are:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = 0 = \mathbf{w} - \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i , \quad (2.7)$$

$$\frac{\partial \mathcal{L}}{\partial b} = 0 = \sum_{i=1}^N \alpha_i y_i . \quad (2.8)$$

Substituting these conditions into Eq. 2.6 results in the dual Lagrangian:

$$\mathcal{L} = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j + \sum_{i=1}^N \alpha_i . \quad (2.9)$$

Most SVM implementations search for a numerical solution to this dual problem. The dual Lagrangian is maximized with respect to α_i and must fulfill the Karush-Kuhn-Tucker (KKT) conditions [16, 17]:

$$\alpha_i \geq 0, \quad \alpha_i (y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1) = 0; \quad i = 1 \dots N. \quad (2.10)$$

Due to these conditions together with Eq. 2.3, all non-support vectors are forced to have vanishing Lagrange multipliers $\alpha_i = 0$, and only the support vectors contribute to the sums in Eq. 2.7 and 2.8, and in the calculation of \mathbf{w}_0 and b_0 at the optimum.

The decision function, i.e. the expression to predict the class label \hat{y} of a new vector \mathbf{u} , follows from the hyperplane equation at the optimum:

$$\hat{y} = \text{sign} \left(\sum_{k=1}^{N_{SV}} y_k \alpha_k \mathbf{x}_k \cdot \mathbf{u} - b_0 \right) , \quad (2.11)$$

where N_{SV} is the number of support vectors.

It is important to note that the dual form in Eq. 2.9 only depends on the scalar products of input vectors, and the same is true for the decision function Eq. 2.11. This advantage of the dual form is essential for the non-linear case in Sec. 2.3.

2.2 Overlapping distributions

The method described so far works in the case of linearly separable data. Overlapping signal and background distributions require a different treatment. By allowing misclassification, the *hard margin* separation above can be modified into a *soft margin* approach. This can be done by introducing *slack* variables [6] ($\xi_i \geq 0$, $i = 1 \dots N$) which measure for each training vector the relative distance by which they are on the wrong side of the separating hyperplane (shown for the one misclassified point in Fig. 1). They are used to weaken the constraints 2.3:

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i, \quad (2.12)$$

and allow to introduce the sum of the slacks $\sum_i^N \xi_i$ as a penalty term into the optimization problem. The modified Lagrangian becomes:

$$\mathcal{L} = \frac{1}{2}|\mathbf{w}|^2 + C \sum_i \xi_i - \sum \alpha_i [y_i(\mathbf{w} \cdot \mathbf{x} + b) - 1 + \xi_i] - \sum \beta_i \xi_i, \quad (2.13)$$

and the extremum condition $\partial \mathcal{L} / \partial \xi_i = 0$ implies a relation between the Lagrange multipliers, $\beta_i = C - \alpha_i$, which allows to bring 2.13 into the same form as 2.6 and eventually into the dual form 2.9. We are left with an optimization problem that is identical to the hard margin case up to the modified constraints.

The constant C that controls the strength of the penalty term appears now only as an upper limit on the Lagrange multipliers $0 \leq \alpha_i \leq C$, restoring the hard margin case in the limit of $C \rightarrow \infty$. Furthermore, it controls the trade-off between simplicity of the decision rule and error frequency and is one of the hyper-parameters that must be set to a sensible value *before* the SVM training.

2.3 Non-linear distributions

The linear SVM presented in the two previous sections is quite limited. For HEP searches we expect complicated, non-linear hyper-surfaces separating the two classes, for which the presented approach can easily be extended to create non-linear decision boundaries. The basic idea for a non-linear SVM [5] is to map the input vectors \mathbf{x}_i into a higher dimensional *feature space* F where the problem becomes linearly separable: $\mathbf{x}_i \mapsto \phi(\mathbf{x}_i) \in F$. The construction of a linear SVM in this feature space follows the same lines as before and the dual Lagrangian from Eq. 2.9 will contain an inner product $\langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$ of elements of F . The peculiar fact that the input vectors only appear in the dual Lagrangian, as well as in the decision function of Eq. 2.11, in form of scalar products allows us to avoid the explicit mapping and to use instead a kernel function $K(\mathbf{x}_i, \mathbf{x}_j) \equiv \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$, such that the dual Lagrangian and the decision function become

$$\mathcal{L} = -\frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) + \sum \alpha_i, \quad (2.14)$$

$$\hat{y} = \text{sign}(\hat{f}), \quad \hat{f} = \sum_{k=1}^{N_{SV}} y_k \alpha_k K(\mathbf{x}_k, \mathbf{u}) - b_0. \quad (2.15)$$

The existence of the mapping $\mathbf{x} \mapsto \phi(\mathbf{x})$ is guaranteed by Mercer’s theorem, given that the kernel function fulfills certain conditions [5, 18], and in general the feature space F may be an even infinite dimensional Hilbert space. For an in-depth exposition on kernel techniques in machine learning see for example [19]. A common and in many applications successful choice [20] is the Gaussian radial basis function (RBF) kernel:

$$K(\mathbf{x}_i, \mathbf{x}_j) = e^{-\gamma|\mathbf{x}_i - \mathbf{x}_j|^2}. \quad (2.16)$$

The width of this kernel function is controlled by the value of γ which is the second hyper-parameter that must be set to a sensible value *before* the training of the SVM. We note that the RBF kernel only contains one parameter and that the components of the data vectors are added quadratically. Therefore, it is useful to normalize the individual components of the input vectors in an appropriate way. Such a scaling ensures that all components of the input vectors may contribute equally. For this paper we use the range between minimum and maximum value for each component. The training data from Eq. 2.1 is replaced by

$$\mathbf{x}_i = (x_{i,0}, \dots, x_{i,n}) \rightarrow \mathbf{x}'_i = (c_0 x_{i,0}, \dots, c_n x_{i,n}), \quad i = 1 \dots N, \quad (2.17)$$

$$c_k = 1 / \left(\max_{i=1 \dots N} x_{i,k} - \min_{i=1 \dots N} x_{i,k} \right), \quad k = 1 \dots n, \quad (2.18)$$

where N is the size of the training sample and n the dimension of the data vectors \mathbf{x} . Later, in the test phase, the identical scaling constants 2.18 must be applied to the test data.

2.4 Probabilistic output

The SVM described so far is a binary classifier with a *yes* or *no* output. In many cases a posterior probability P that quantifies the belief in the class label is useful and offers an easier interpretation. In Section 3, a probability cut, $P > P_0$, is used to modify the signal-to-background ratio and the total number of selected events. Such a classification probability can be estimated [21, 22] by fitting a sigmoid model to the training data (y_i, \mathbf{x}_i) :

$$P(y = 1 | \hat{f}) = \begin{cases} \frac{\exp(-t)}{1 + \exp(-t)} & : t \equiv A + B\hat{f} \geq 0 \\ \frac{1}{1 + \exp(t)} & : t < 0 \end{cases}, \quad (2.19)$$

where the decision value \hat{f} is given in 2.15.

In general, especially for the non-linear SVM, the result will be biased if the SVM training data itself is used for the fit. In LIBSVM the bias is avoided by a five-fold cross-validation³. It is important to note that a strictly decreasing function of the decision value, as 2.19, does not change the order of any sequence of decision values. Since the cross-validation increases the computational burden we do not calculate the probabilities during the parameter scan but only for presenting the final results.

³The procedure of k-fold cross-validation splits the training data randomly into k equal sized subsamples. One subsample is retained for validation while the remaining k-1 subsamples are used for training. The training is repeated k times with changing roles such that each subsample is used exactly once for validation.

3 Hyper-parameter tuning

The two parameters C and γ , introduced in the previous section, must be set to sensible values before the training of the support vector machine. These values are crucial for the performance of the algorithm, and different strategies to optimize the parameter choice are possible. The easiest approach is a simple grid search. A two dimensional grid is defined, at each grid point the SVM training is performed on a training dataset, and the trained machine is applied to an independent test data sample where some performance measure is evaluated. Eventually, the (C, γ) -pair with the highest performance index is used. We first consider what could be an appropriate performance measure for a new physics search and describe then the tuning algorithm in some detail.

3.1 Performance measures

Machine learning performance measures. From a machine learning perspective, a natural performance measure describes how well a classifier separates the two distinct classes. On a sample of test data we know the true class labels. There are 2x2 categories formed by the true label $y \in \{-1, 1\}$ and the label $\hat{y} \in \{-1, 1\}$ estimated by the SVM. The relative amount of test data in these categories can be used to quantify the performance of a machine learning algorithm. Typical ML measures are the *accuracy* which gives the percentage of correctly predicted labels, or the *precision* which, in our case, is the percentage of correctly predicted signal events. Another frequently used measure is the *AUC*, the area under the receiver operator curve (ROC). The ROC curve shows the background rejection (false positive) against the signal efficiency (true positive) at various threshold values of the decision function. While the use of these and other performance measures is common also in HEP [1], we will follow a different approach to optimize the SVM.

Physics motivated performance measures. The maximum number of correctly classified events is of secondary importance for a HEP search. There is a much more physically and statistically motivated measure: the discovery significance. Optimizing the significance is a common procedure in HEP. Typically the search area is optimized for a statistically relevant signal to background ratio that allows to prove or reject a certain hypothesis. Here, we consider the case of a cut-and-count analysis for which several significance estimators are commonly used [23][24]. Optimizing a certain, statistically motivated, figure of merit is common practice in HEP to select different ML algorithms or different sets of input variables. The new insight of this paper is that such a procedure can successfully be applied in the stage of model selection, i.e. during the hyper-parameter tuning.

Asimov estimate. The exact numerical calculation of the statistical significance may become computationally costly. A well performing estimate for the discovery significance has been given in [23]. For the case of Poisson distributed background and signal events (s, b) with background uncertainty σ_b the approximated median discovery significance becomes

$$Z_A = \left[2 \left((s + b) \ln \left[\frac{(s + b)(b + \sigma_b^2)}{b^2 + (s + b)\sigma_b^2} \right] - \frac{b^2}{\sigma_b^2} \ln \left[1 + \frac{\sigma_b^2 s}{b(b + \sigma_b^2)} \right] \right) \right]^{1/2}. \quad (3.1)$$

3.2 Hyper-parameter search

The SVM with RBF kernel requires two hyper-parameters: C (Sec. 2.2) and γ (Sec. 2.3). In addition to the hyper-parameters, the number of selected signal and selected background events (s, b) depends on the probability cutoff P_0 (Sec. 2.4), or a corresponding decision value \hat{f}_0 . The easiest algorithm to find the optimal values for these parameters is a grid search. At each point of a logarithmically spaced grid in (C, γ) a SVM is trained and, on an independent test dataset, the Asimov significance Z_A is calculated as function of P_0 . In general, the best cut P_0 is selected as the value with the highest significance. To avoid artificially high significance values due to statistical fluctuation of a small signal at very low values of b , a further requirement of at least 5 signal events is applied. While conceptually the plain grid search is sufficient to find good values for (C, γ) , computationally it may be advantageous to use a more refined algorithm for the hyper-parameter tuning. The details of the iterative algorithm used for the results in this paper are given in Appendix A.

4 Case studies

4.1 Performance comparison on a toy model

Comparing speed and classification performance of different classifiers is not always straightforward. In order to have simple and well defined conditions, we start with a toy model and compare our SVM-HINT framework with a BDT and an SVM, both implemented with the TMVA library. The toy model includes the following variables V_i generated with the random numbers x_i (the *tilde* means sampled from):

$$\begin{aligned}
 V_1 &= \sin(x_1); & x_1 &\sim g(x_1|a, b) \\
 V_2 &= x_2; & x_2 &\sim \exp(-x_2/c) \\
 V_3 &= x_3; & x_3 &\sim g(x_3|d, e) \\
 V_4 &= \sqrt{x_4}; & x_4 &\sim \exp(-x_4/f)
 \end{aligned}
 \tag{4.1}$$

where $g(x|\mu, \sigma)$ is a Gaussian distribution with mean μ and width σ , and $a, b, c, d, e, f > 0$ are constants with different values for signal and background samples. This model does not have any hidden correlation between the variables and each ML algorithm needs only to find a set of independent optimal cuts. Due to its simplicity, the toy model enables us to generate large quantities of events to study the training time as function of the training sample size for the different codes.

As explained in Sec. 3, SVM-HINT provides a hyper-parameter search. The hyper-parameter search is performed beforehand and is not part of the timing performance study. The SVM implementation provided by the TMVA library lacks such an automated hyper-parameter search. We therefore apply the same hyper-parameter values as obtained by the SVM-HINT tuning algorithm. The out-of-the-box performance of the TMVA-BDT cannot compete in most cases with the automatically tuned SVM-HINT. There is a trade-off between classification performance and time consumption of the BDT which can be controlled by an appropriate choice of the BDT parameters, e.g. the number of trees, minimum node size, and cut values, as introduced in Sec. B. For comparing the training times, we follow

the strategy to optimize the BDT parameters manually⁴ to accomplish a similar discovery significance as achieved with the SVM-HINT. This allows us to compare the time consumption of equally performing algorithms. Blindly optimizing for maximal performance on a fixed evaluation training sample could otherwise produce slow, over-sized trees and would place a disadvantage on the BDT implementation.

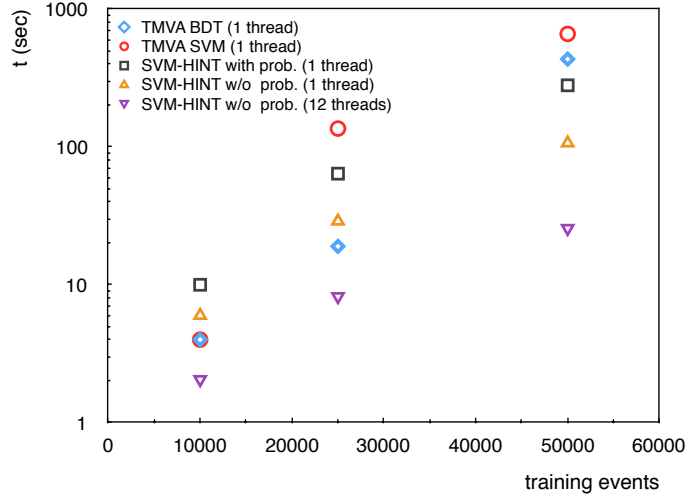


Figure 2. The timing performance of the classifiers are compared on a computer with two Intel[®] Xeon[®] E5-2440 CPUs and 12 physical threads running at 2.40 GHz clock speed. Number of threads used by each classifier implementations are stated within parentheses in the legend.

Fig. 2 shows the results for the different classifiers. At low numbers of training events the BDT performs better than the other single-threaded classifiers. With increasing training sample size the number of trees needed to achieve a competitive classification performance becomes larger with negative impact on the training time. The TMVA-SVM does not scale well in terms of timing performance and it performs poorly on bigger samples. Overall the SVM-HINT performs similar or better compared to TMVA-BDT and TMVA-SVM. In addition, the SVM-HINT can efficiently take advantage of multi-core architectures. Naturally, the multi-threaded performance of the SVM-HINT with 12 threads outperforms the other implementations.

4.2 Third generation supersymmetric partner search

Monte Carlo samples As a real-world physics example we consider a search for the supersymmetric partner of the top quark, called top squark, at the LHC. The search is designed for the case of direct top-squark pair production with subsequent decay of the top squarks into the lightest supersymmetric particle (LSP) and a top quark. Several searches for such a scenario have been performed at TeVatron as well as at the LHC [25–32]. After preselection of the data, the remaining dominant background is given by top-antitop ($t\bar{t}$) quark production. Top quarks decay to almost 100% to a b quark and a W boson, with the latter decaying either to two quarks or to a lepton and a neutrino. When requiring

⁴Configuration files are available at [13].

one lepton in the final state, we mainly expect to select semi-leptonic decays of top quarks (where one W boson decays to two jets and the other to a lepton and a neutrino), but dileptonic top decays (where both W bosons decay to a lepton and a neutrino) might be selected as well, if one lepton is not identified (lost) for various reasons. A leading order simulation is sufficient for our purpose and we only take into account NLO results for the total cross sections of signal [33] and background [34, 35] processes. PYTHIA6 [36] is used for the event simulation and DELPHES3 [37] to model the detector response. The detector model is a *combined* ATLAS and CMS detector, as it had been used for the 2013 Snowmass effort [38].

We categorize the physical variables with respect to their mathematical complexity as high-level and low-level variables. The low-level variables consist of basic properties of the reconstructed physics objects measured by the detector, while the high-level variables are constructed from the low-level variables using physical insight to improve the classification performance. The physics objects are jets and leptons, where lepton is used as generic term for electrons and muons. For simplicity, we do not consider tau leptons since their experimental reconstruction is more complex.

Low-level variables are the transverse momentum p_T and the pseudorapidity η of the single lepton, of the four highest- p_T (called ‘leading’) jets, and of the leading b-quark jet. In hadron collider experiments the missing energy perpendicular to the beam direction, \cancel{E}_T , is commonly reconstructed as an independent variable. It is therefore treated as a low-level quantity, as well as H_T , the scalar sum over the transverse momenta of all preselected jets. In many SUSY models, we expect large \cancel{E}_T due to the LSP, which is expected to be neutral and weakly interacting and will therefore not be detected. As SUSY particles are heavy, we also expect a large amount of energy in the detector leading to large H_T . In addition, the multiplicities of jets (n_{jet}) and b-quark jets ($n_{b\text{jet}}$) are included.

As high-level variables we consider the following variables: the transverse mass m_T , defined as $m_T = \sqrt{2 p_{T,l} \cancel{E}_T (1 - \cos \Delta\phi(l, \cancel{E}_T))}$, can be used to suppress the background from W-boson production, as m_T of leptonic W decay events does not exceed the the W mass. Dileptonic $t\bar{t}$ -events with one lost lepton are an important background since the lost lepton mimics large missing energy from the LSP. The m_{T2}^W variable [39] is constructed exploiting the knowledge of the $t\bar{t}$ -decay kinematics to separate such events. Top-squark production is a high-mass process with large missing energy. High-mass production is typically related to more centrally distributed particles in the detector, such that the *Centrality*, defined as $\sum_{\text{jets},l} p_T / \sum_{\text{jets},l} p$, can be used to enhance such events. Commonly used relations between the hadronic activity and the missing energy are $Y = \cancel{E}_T / \sqrt{H_T}$, often referred to as \cancel{E}_T significance, and the H_T -ratio, the normalized hadronic activity in the hemisphere of \cancel{E}_T . The last group of variables exploit topological relations between the event particles: $\Delta\phi(W, l)$ is the angle between the W boson and lepton, $\Delta r_{\min}(l, b)$ is the radial distance between closest lepton and b -jet and $m(l, b)$ is the invariant mass of the b -jet and the closest lepton.

A compilation of all low-level and high-level variables is given in Table 1, together with the definition of four sets of variables which are considered to investigate the influence of the variable multiplicity and complexity in the multivariate analysis. We define one set containing all variables, one using only low- or high-level variables, respectively, and a

subset of two low-level and two high-level variables with relatively large separation power.

Table 1. Summary of all low-level and high-level variables used in the analysis. Set 1 includes all variables. Set 2 and set 3 consist of low- and high-level variables, respectively. Set 4 is a smaller subset of high- and low-level variables.

	Variable	Set 1	Set 2	Set 3	Set 4
low-level	$p_{T,l}$	•	•		
	η_l	•	•		
	$p_{T,jet(1,2,3,4)}$	•	•		
	$\eta_{jet(1,2,3,4)}$	•	•		
	$p_{T,b\ jet1}$	•	•		
	$\eta_{b\ jet1}$	•	•		
	n_{jet}	•	•		
	$n_{b\ jet}$	•	•		
	\cancel{E}_T	•	•		•
	H_T	•	•		•
	high-level	m_T	•		•
m_{T2}^W		•		•	•
$\Delta\phi(W, l)$		•		•	
$m(l, b)$		•		•	
Centrality		•		•	
Y		•		•	
H_T -ratio		•		•	
$\Delta r_{\min}(l, b)$		•		•	
$\Delta\phi_{\min}(j_{1,2}, \cancel{E}_T)$		•		•	

Analysis strategy In order to reduce the time required for training and optimization, a baseline selection, summarized in Table 2, is applied to the signal and background samples. Figure 3 shows the distribution of signal and background for two low-level and two high-level variables, H_T , \cancel{E}_T , m_T and m_{T2}^W after the baseline selection, normalized to the expected luminosity at the end of the LHC run in the year 2023, corresponding to 300 fb^{-1} . The background is several orders of magnitude higher than the signal, and the distributions of signal and background are quite similar due to their similar kinematics.

The samples are separated into three independent subsamples: training, test and evaluation sample. Each classification method is optimized over the training and test samples and the best-performing configuration is applied to the independent evaluation sample for the final performance assessment.

The TMVA-BDT has been manually trained and tested over 8 different settings for each of the four variable sets in order to obtain optimal parameters as described in Appendix B, while the SVM-HINT is auto-tuned by the iterative grid search, as described in Sec. 3 and Appendix A. Without modifying the default SVM-HINT settings, the two step grid

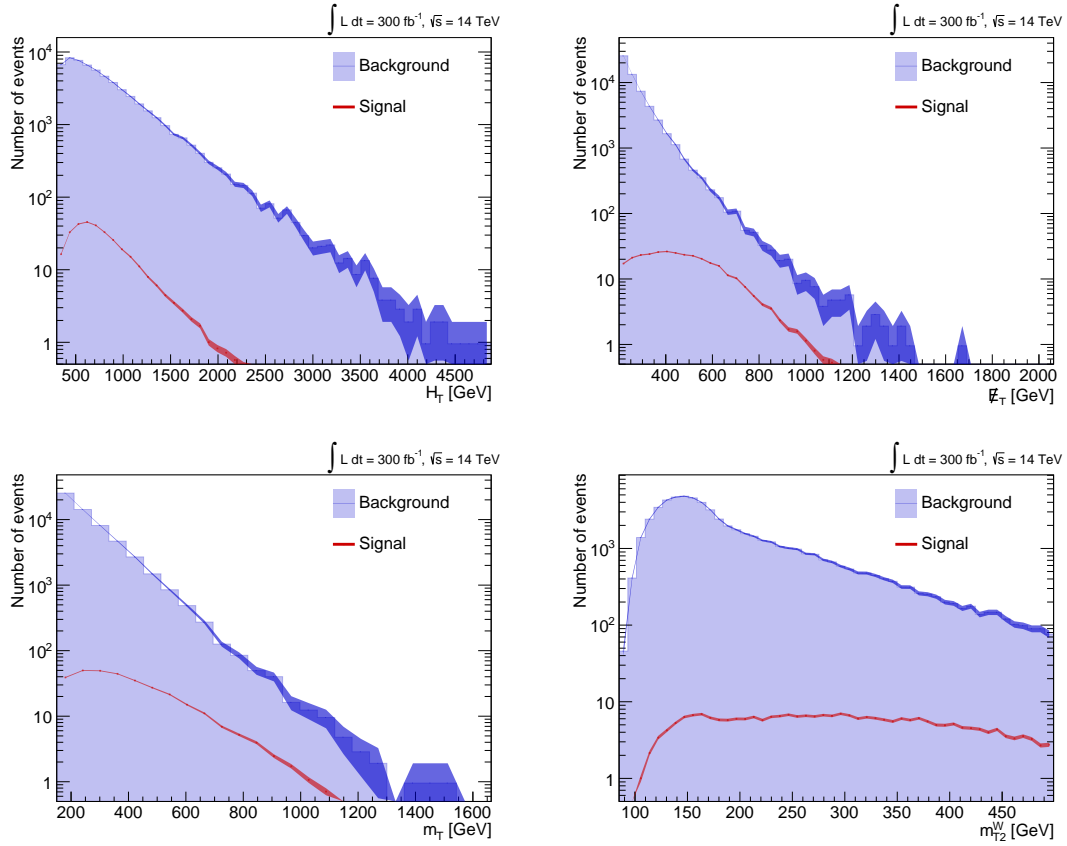


Figure 3. The distribution of signal (red line) and background (blue filled histogram) after the baseline selection for two low-level and two high-level variables that are used in the analysis: H_T , E_T , m_T and m_{T2}^W . The y-axis shows the number of events (normalized to the integrated luminosity), and the x-axis shows the variable value for a given bin. Statistical errors are represented by transparent bands.

Table 2. Top-squark search: List of baseline selection requirements

$ \eta_{l, \text{jet}} $	$<$	2.4
$p_{T, l}$	$>$	30 GeV
$p_{T, \text{jet}}$	$>$	40 GeV
$p_{T, \text{jet}1}$	$>$	80 GeV
$p_{T, \text{jet}2}$	$>$	60 GeV
E_T	$>$	200 GeV
H_T	$>$	300 GeV
n_{jet}	$>$	3
$n_{b \text{ jet}}$	$>$	0

search hyper-parameter optimization function provides the optimal parameters using test and training samples. We calculate the final significance with an independent evaluation

sample.

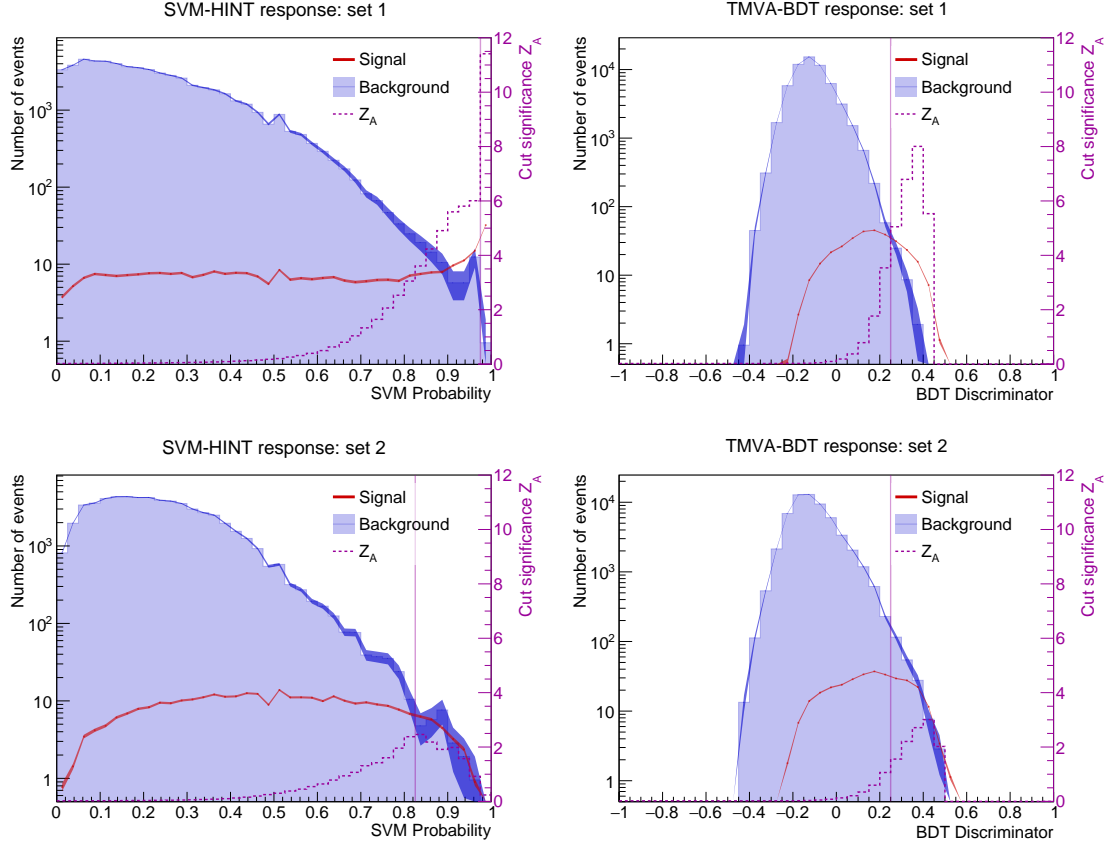


Figure 4. The SVM-HINT and TMVA BDT responses trained with the variable sets set 1 and set 2, as defined in 1. Even though the optimal Z_A efficiency information is not available in data, it is included to demonstrate the reliability of the optimal discriminator cut output from the classifier implementations. The y-axis on the left shows the number of events (normalized to the aimed integrated luminosity), whereas the y-axis on the right shows the Asimov significance for the discriminator cut per bin.

Results Figures 4 and 5 show the performance of the SVM-HINT and the TMVA-BDT for the four different variable sets. Both the SVM-HINT as well as the TMVA-BDT obtain the highest accuracy with the largest number of variables (set 1). Both methods perform much worse with only low-level variables (set 2), while the high-level variables (set 3) are clearly able to separate signal and background. Despite the poor performance of the low-level variables, they add a substantial amount of extra information to enlarge the significance when adding them to the high-level variables. Here we conclude that even with variable multiplicity of 25, SVM-HINT as well as TMVA-BDT do not require a preselection of variables. Reducing the number of variables to two low-level and two high-level variables gives slightly better results for the SVM-HINT than for the TMVA-BDT. Here we have to note that this study is not meant to compare the two methods, but to put the results obtained by SVM-HINT into a more known context in HEP with the TMVA-BDT. The SVM-HINT

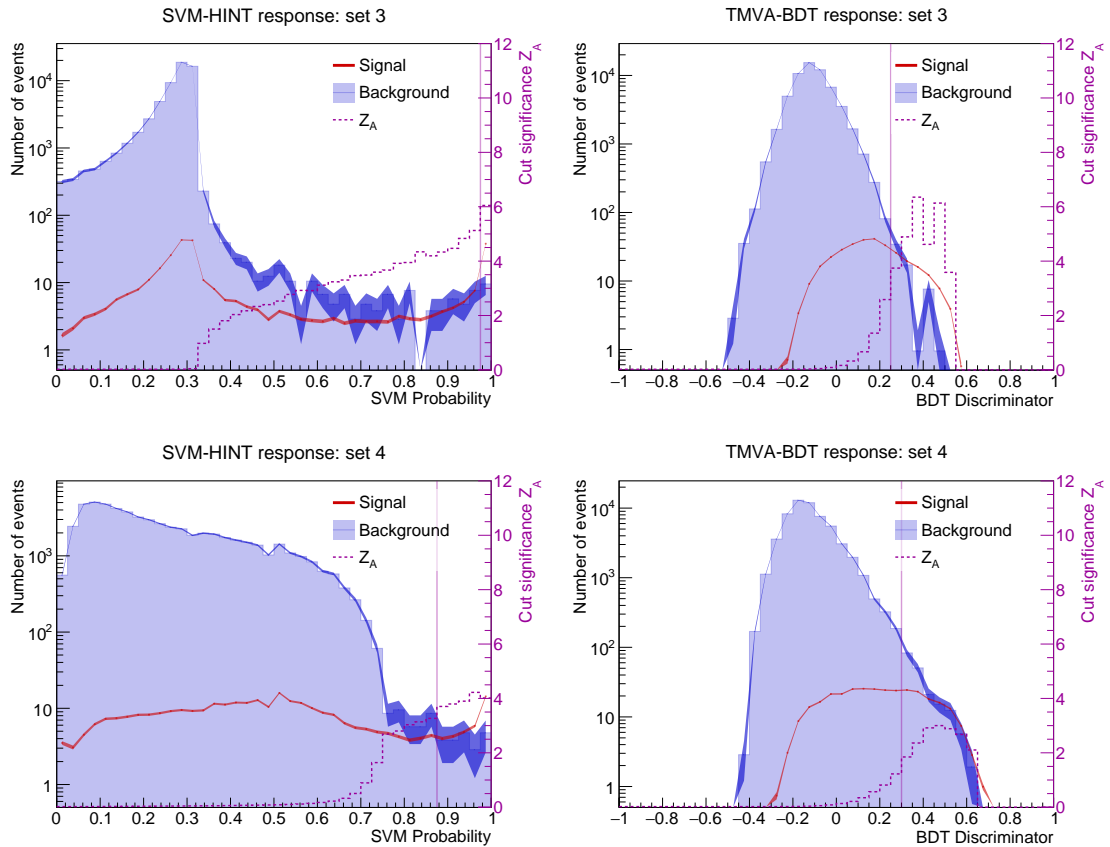


Figure 5. The SVM-HINT and TMVA BDT responses trained with the variables sets set 3 and set 4, as defined in 1. Even though the optimal Z_A efficiency information is not available in data, it is included to demonstrate the reliability of the optimal discriminator cut output from the classifier implementations. The y-axis on the left shows the number of events (normalized to the aimed integrated luminosity), whereas the y-axis on the right shows the Asimov significance for the discriminator cut per bin.

with the automated hyperparameter tuning needs no further manual optimization and can make use of high number of variables simultaneously with an increasing classification power.

Furthermore, the importance of the performance measures is visible on the discriminator cut decisions by the classifier implementations. The TMVA BDT uses simply $\frac{S}{\sqrt{S+B}}$ as the performance measure, where S and B correspond to the number of signal and background events, respectively. This formula performs differently than a log-likelihood significance calculation. Therefore, the optimal cut provided by TMVA reduces the significance obtained from the classifier implementation. SVM-HINT uses the Asimov significance which gives very similar results to the log-likelihood calculation, and therefore, the results obtained from SVM-HINT not only provide good out-of-the-box estimation of the actual significance, but the discriminator cut given by SVM-HINT maximizes the significance between background and signal.

5 Conclusions

Our results show that a Support Vector Machine is an efficient machine learning algorithm for new physics searches in high-energy physics. The rare applications of this tool in our field may be related to the limited implementation provided by the popular TMVA library. An appropriate designed automatic search for the two hyper-parameters easily overcomes this limitations and reveals the full potential of this approach. The Support Vector Machine is certainly able to compete with a Boosted Decision Tree which currently is the prevalent machine learning tool in high-energy physics. We do not intent to claim that one of the algorithms out-performs the other. This would need a diligent optimization of our Boosted Decision Tree, which is beyond the scope of this paper. The performance of a Boosted Decision Trees depends on a larger number of parameters which complicates the construction of an automated tuning procedure. The clear advantage of the Support Vector Machine is rather the straightforward hyper-parameter tuning. We demonstrate that the approximated median discovery significance (Asimov significance) is an effective figure of merit for the parameter tuning and that only two parameters need to be adapted to define a well performing search tool. The SVM maximum margin concept guarantees good generalization properties of the trained algorithm while at the same time the hyper-parameter tuning allows to find a non-linearly bounded area with maximized significance. Furthermore, Support Vector Machines are known to be robust against an large number of even partially correlated input variables. This is in agreement with our studies, a lengthy selection of useful input variables was not necessary. The algorithm reliably exploits all available information.

A Iterative grid search

As shown in section 3.2 an SVM with RBF kernel requires two different hyper-parameters to be adjusted: C and γ . While in principle a *brute force* grid search is sufficient to find the best hyper-parameters, we used an adaptive search strategy for the results in this paper which we describe here for completeness. The SVM-HINT grid search algorithm uses a modified version of the Asimov Significance 3.1, a significance score \tilde{Z}_A based on the difference between the significance value observed in the test sample and the significance value from the training sample.

$$\tilde{Z}_A = Z_A^{(test)} \left[1 - \frac{|Z_A^{(test)} - Z_A^{(train)}|}{Z_A^{(test)} + Z_A^{(train)}} \right]. \quad (\text{A.1})$$

This way, the extreme significance values observed due to fluctuations or overtraining can be penalized without a high computational effort. The search algorithm can be formalized as follows:

1. For the given initial parameters $\gamma_{initial}$ and $C_{initial}$, the iterative grid search algorithm produces an array of logarithmically spaced γ values with a step size K_t around the mid-value $\gamma_m^{(1)} = \gamma_{initial}$ such that:

$$\begin{aligned} \gamma_k^{(l)} &= K_t \cdot \gamma_{k-1}^{(l)}, \quad \text{where } K_t = \frac{1}{2}(1 + \ln(t/2)), \\ k &= 0, \dots, m, \dots, 2m = 18, \quad t = \text{int}(l/4), \quad l = 1, \dots, 20 \end{aligned} \quad (\text{A.2})$$

where l indicates the number of iterations, the variable t is a *focus parameter* that decreases the step size factor K_t every fourth iteration.

\tilde{Z}_A is then evaluated for all of these C - γ -pairs.

2. For the next step C is increased to $C^{(l+1)} = 1.5 \cdot C^{(l)}$ and \tilde{Z}_A is again calculated with each value in the γ array.
3. If the maximum $\tilde{Z}_A^{(l)}$ value is at least 30% larger then the best \tilde{Z}_A^{l-1} from the previous iteration the higher C parameter is accepted. The 30% hurdle is introduced to stabilize against fluctuations.
4. After each fourth iteration, the C - γ -pair corresponding the highest significance score is taken as the new initial γ and the algorithm returns to the first step; now with a smaller step size factor K_t such that the new γ array has a tighter stepping around the new initial γ value.
5. When the number of iterations reaches the pre-defined maximum value, the algorithm stops and the γ - C -pair with the maximum \tilde{Z}_A in the final iteration are returned as the best hyper-parameter values.

The procedure assumes that a sufficiently small $C_{initial}$ had been chosen. In case that the found best C value is identical with the $C_{initial}$ the algorithm is restarted with a smaller value of $C_{initial}$.

B Boosted Decision Trees

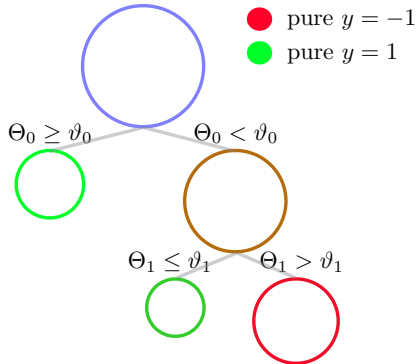


Figure 6. Representation of a simple binary decision tree structure: Each node is split with respect to a variable Θ_i and a cut value θ_i determined by the performance measure.

Boosted decision trees (BDT) are probably the most common ML classifier in experimental particle physics. We therefore compare the performance of our SVM framework to a BDT implemented with TMVA. We shortly introduce the relevant BDT concepts used in this comparison.

Decision Trees A binary decision tree [40] separates signal and background by a sequence of binary cuts (Fig. 6). The terminating branches, or leaves, correspond to a cubical separation in the multi-dimensional space of the input variables and the tree as a whole forms a complex separation boundary between the two classes. Each node of the tree is connected to two branches that are split with respect to only one of the variables $x_i^{(1)} \dots x_i^{(n)}$ defined in 2.2. At each node the algorithm selects one of the available physical variables and searches for a best cut value. This process requires a suitable goodness-of-split measure and a common choice is the *Gini impurity* index which is also used for the TMVA-BDT in this paper. For the two class case the *Gini* index is given by $g = 2p(1 - p)$, where the purity p of a node is defined as the ratio of signal events over all events. The training starts with the root node, and the tree is constructed recursively while at each split the reduction in *Gini* impurity is maximized. The splitting stops when a node falls below a predefined minimum of events.

Pruning The constructed decision tree is sensitive to statistical fluctuation. To avoid overtraining it must be pruned to remove statistically insignificant nodes. *Cost complexity pruning* [40] removes branches which increase the misclassification cost. The misclassification rate $R = 1 - \max(p, 1 - p)$ is used as a cost estimate at each node and compared to the cost of the subtree below the node. The cost complexity is defined as $\rho = (R_{\text{node}} - R_{\text{subtree}})/(N_{st} - 1)$, where N_{st} is the number of nodes in the subtree. The node with the smallest ρ is recursively removed from the tree as long as ρ is below a certain pruning strength value ρ_0 .

Boosting A single decision tree is seldom an efficient classifier. Boosting is a powerful iterative algorithm which improves the performance of weak classifiers. The boosting of a decision tree extends this concept from one tree to a forest of trees. The trees are derived from the same training data by reweighting events, and are finally combined into a single classifier of considerably enhanced performance. For this paper AdaBoost [41] is used.

The TMVA-BDT (AdaBoost) is trained and tested with 32 different settings to obtain an optimal parameter set (TMVA: `nEventsMin`, `NTrees`, `MinNodeSize`, `MaxDepth` and `AdaBoostBeta`) and the best performing configuration is used for evaluation⁵.

Acknowledgments

M. Ö. Sahin would like to thank the Joachim Herz foundation for the support.

References

- [1] A. Hoecker, P. Speckmayer, J. Stelzer, J. Therhaag, E. von Toerne, and H. Voss, “TMVA: Toolkit for Multivariate Data Analysis,” *PoS ACAT* (2007) 040, [arXiv:physics/0703039](https://arxiv.org/abs/physics/0703039).
- [2] “The HiggsML challenge,” May to Sep. 2014. <https://higgsml.lal.in2p3.fr/>.
- [3] “HEPML 2014 proceedings,” vol. 42. 2015.
- [4] “Data science @ LHC 2015 Workshop,” Nov. 2015. <https://indico.cern.ch/event/395374/>.
- [5] B. E. Boser, I. Guyon, and V. Vapnik, “A training algorithm for optimal margin classifiers,” in *Proceedings of the Fifth Annual Workshop on Computational Learning Theory*, pp. 144–152. ACM Press, 1992.
- [6] C. Cortes and V. Vapnik, “Support vector networks,” *Machine Learning* **20** (1995) 273 – 297.
- [7] V. N. Vapnik, *The Nature of Statistical Learning Theory*. Springer New York Inc., 1998.
- [8] P. Vannerem, K. Muller, B. Scholkopf, A. Smola, and S. Soldner-Rembold, “Classifying LEP data with support vector algorithms,” [arXiv:hep-ex/9905027](https://arxiv.org/abs/hep-ex/9905027) [[hep-ex](#)].
- [9] A. Vaiciulis, “Support vector machines in analysis of top quark production,” *Nucl. Instrum. Meth.* **A502** (2003) 492–494, [arXiv:hep-ex/0205069](https://arxiv.org/abs/hep-ex/0205069) [[hep-ex](#)].
- [10] Ł. Janyst, A. Kaczmarek, T. Szymocha, M. Wolter, and A. Zemła, “Optimization of tau identification in atlas experiment using multivariate tools,” *Computer Science* **Vol. 9** (2008) 35–45.
- [11] CDF Collaboration, T. Aaltonen *et al.*, “Search for the standard model Higgs boson produced in association with a W^\pm boson with 7.5 fb^{-1} integrated luminosity at CDF,” *Phys. Rev. D* **86** (Aug, 2012) 032011. <http://link.aps.org/doi/10.1103/PhysRevD.86.032011>.
- [12] F. Sforza, V. Lippi, and G. Chiarelli, “Rejection of multi-jet background in $p\bar{p} \rightarrow e\nu + j\bar{j}$ channel through a SVM classifier,” *Journal of Physics: Conference Series* **331** no. 3, (2011) 032045. <http://stacks.iop.org/1742-6596/331/i=3/a=032045>.

⁵Configuration files are available at [13]

- [13] M. O. Sahin, D. Kruecker, and I. A. Melzer-Pellmann, “SVM-HEP Interface,” 2015. <https://www.github.com/ml-hint/svm-hint>.
- [14] I. Antcheva *et al.*, “ROOT: A C++ framework for petabyte data storage, statistical analysis and visualization,” *Comput. Phys. Commun.* **182** (2011) 1384–1385.
- [15] C.-C. Chang and C.-J. Lin, “LIBSVM: A library for support vector machines,” *ACM Transactions on Intelligent Systems and Technology* **2** (2011) 27:1–27:27. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- [16] W. Karush, *Minima of Functions of Several Variables with Inequalities as Side Constraints*. PhD thesis, University of Chicago, 1939.
- [17] H. W. Kuhn and A. W. Tucker, “Nonlinear programming,” in *Proceedings of the Second Berkeley Symposium on Mathematical Statistics and Probability, 1950*, pp. 481–492. University of California Press, Berkeley and Los Angeles, 1951.
- [18] J. Mercer, “Functions of positive and negative type, and their connection with the theory of integral equations,” *Philosophical Transactions of the Royal Society of London A: Mathematical, Physical and Engineering Sciences* **209** no. 441-458, (1909) 415–446.
- [19] B. Schölkopf and A. Smola, *Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*. Adaptive computation and machine learning. MIT Press, 2002.
- [20] C.-W. Hsu, C.-C. Chang, and C.-J. Lin, “A practical guide to support vector classification,”. <https://www.csie.ntu.edu.tw/~cjlin/papers/guide/guide.pdf>.
- [21] J. C. Platt, “Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods,” in *Advances in large margin classifiers*, pp. 61–74. 1999.
- [22] T. F. Wu, C. J. Lin, and R. C. Weng, “Probability estimates for multi-class classification by pairwise coupling,” *J. Mach. Learn. Res.* **5** (2004) 975–1005.
- [23] G. Cowan, K. Cranmer, E. Gross, and O. Vitells, “Asymptotic formulae for likelihood-based tests of new physics,” *European Physical Journal C* **71** (Feb., 2011) 1554, [arXiv:1007.1727 \[physics.data-an\]](https://arxiv.org/abs/1007.1727).
- [24] R. D. Cousins, J. T. Linnemann, and J. Tucker, “Evaluation of three methods for calculating statistical significance when incorporating a systematic uncertainty into a test of the background-only hypothesis for a Poisson process,” *Nuclear Instruments and Methods in Physics Research A* **595** (Oct., 2008) 480–501, [physics/0702156](https://arxiv.org/abs/physics/0702156).
- [25] **D0** Collaboration, V. M. Abazov *et al.*, “Search for 3- and 4-body decays of the scalar top quark in pp collisions at $\sqrt{s} = 1.8$ TeV,” *Physics Letters B* **581** no. 3-4, (2004) 147 – 155.
- [26] **D0** Collaboration, V. M. Abazov *et al.*, “Search for pair production of the scalar top quark in muon+tau final states,” *Phys. Lett.* **B710** (2012) 578–586, [arXiv:1202.1978 \[hep-ex\]](https://arxiv.org/abs/1202.1978).
- [27] **D0** Collaboration, V. M. Abazov *et al.*, “Search for the lightest scalar top quark in events with two leptons in $p\bar{p}$ collisions at $\sqrt{s} = 1.96$ TeV,” *Phys. Lett.* **B659** (2008) 500–508, [arXiv:0707.2864 \[hep-ex\]](https://arxiv.org/abs/0707.2864).
- [28] **CDF** Collaboration, T. Aaltonen *et al.*, “Search for the supersymmetric partner of the top quark in $p\bar{p}$ collisions at $\sqrt{s} = 1.96$ TeV,” *Phys. Rev.* **D82** (2010) 092001, [arXiv:1009.0266 \[hep-ex\]](https://arxiv.org/abs/1009.0266).
- [29] **CDF** Collaboration, D. Acosta *et al.*, “Search for the supersymmetric partner of the top

- quark in dilepton events from $p\bar{p}$ collisions at $\sqrt{s} = 1.8$ TeV,” *Phys. Rev. Lett.* **90** (2003) 251801, [arXiv:hep-ex/0302009](#) [[hep-ex](#)].
- [30] **ATLAS** Collaboration, G. Aad *et al.*, “ATLAS Run 1 searches for direct pair production of third-generation squarks at the Large Hadron Collider,” *Eur. Phys. J.* **C75** no. 10, (2015) 510, [arXiv:1506.08616](#) [[hep-ex](#)].
- [31] **CMS** Collaboration, S. Chatrchyan *et al.*, “Search for top-squark pair production in the single-lepton final state in pp collisions at $\sqrt{s} = 8$ TeV,” *Eur. Phys. J.* **C73** no. 12, (2013) 2677, [arXiv:1308.1586](#) [[hep-ex](#)].
- [32] M. Berggren, A. Cakir, D. Krücker, J. List, I. A. Melzer-Pellmann, B. S. Samani, C. Seitz, and S. Wayand, “Non-Simplified SUSY: Stau-Coannihilation at LHC and ILC,” [arXiv:1508.04383](#) [[hep-ph](#)].
- [33] C. Borschensky, M. Krämer, A. Kulesza, M. Mangano, S. Padhi, T. Plehn, and X. Portell, “Squark and gluino production cross sections in pp collisions at $\sqrt{s} = 13, 14, 33$ and 100 TeV,” *Eur. Phys. J.* **C74** no. 12, (2014) 3174, [arXiv:1407.5066](#) [[hep-ph](#)].
- [34] J. M. Campbell and R. K. Ellis, “MCFM for the Tevatron and the LHC,” *Nucl. Phys. Proc. Suppl.* **205-206** (2010) 10–15, [arXiv:1007.3492](#) [[hep-ph](#)].
- [35] P. M. Nadolsky, H.-L. Lai, Q.-H. Cao, J. Huston, J. Pumplin, D. Stump, W.-K. Tung, and C. P. Yuan, “Implications of CTEQ global analysis for collider observables,” *Phys. Rev.* **D78** (2008) 013004, [arXiv:0802.0007](#) [[hep-ph](#)].
- [36] T. Sjostrand, S. Mrenna, and P. Z. Skands, “PYTHIA 6.4 Physics and Manual,” *JHEP* **0605** (2006) 026, [arXiv:hep-ph/0603175](#) [[hep-ph](#)].
- [37] S. Oryn, X. Rouby, and V. Lemaitre, “DELPHES, a framework for fast simulation of a generic collider experiment,” [arXiv:0903.2225](#) [[hep-ph](#)].
- [38] J. Anderson *et al.*, “Snowmass Energy Frontier Simulations,” in *Community Summer Study 2013: Snowmass on the Mississippi (CSS2013) Minneapolis, MN, USA, July 29-August 6, 2013*. 2013. [arXiv:1309.1057](#) [[hep-ex](#)].
- [39] J. G. Yang Bai, Hsin-Chia Cheng and J. Gu, “Stop the top background of the stop search,” [arXiv:1203.4813](#) [[hep-ph](#)].
- [40] L. Breiman, J. Friedman, C. J. Stone, and R. A. Olshen, *Classification and regression trees*. CRC press, 1984.
- [41] Y. Freund and R. E. Schapire, “A decision-theoretic generalization of on-line learning and an application to boosting,” *Journal of Computer and System Sciences* **55** no. 1, (1997) 119 – 139.