

# Understanding Where Your Classifier Does (Not) Work — the SCaPE Model Class for EMM

Wouter Duivesteyn

Fakultät Informatik, LS VIII

Lehrstuhl für Künstliche Intelligenz

Technische Universität Dortmund, Germany

wouter.duivesteyn@tu-dortmund.de

Julia Thaele

Fakultät Physik, LS E5B

Lehrstuhl für Astroteilchenphysik

Technische Universität Dortmund, Germany

julia.thaele@tu-dortmund.de

**Abstract**—FACT, the First G-APD Cherenkov Telescope, detects air showers induced by high-energetic cosmic particles. It is desirable to classify a shower as being induced by a gamma ray or a background particle. Generally, it is nontrivial to get any feedback on the real-life training task, but we can attempt to understand how our classifier works by investigating its performance on Monte Carlo simulated data. To this end, in this paper we develop the SCaPE (Soft Classifier Performance Evaluation) model class for Exceptional Model Mining, which is a Local Pattern Mining framework devoted to highlighting unusual interplay between multiple targets. In our Monte Carlo simulated data, we take as targets the computed classifier probabilities and the binary column containing the ground truth: which kind of particle induced the corresponding shower. Using a newly developed quality measure based on ranking loss, the SCaPE model class highlights subspaces of the search space where the classifier performs particularly well or poorly. These subspaces arrive in terms of conditions on attributes of the data, hence they come in a language a domain expert understands, which should aid him in understanding where his/her classifier does (not) work. Found subgroups highlight subspaces whose difficulty for classification is corroborated by astrophysical interpretation, as well as subspaces that warrant further investigation.

**Keywords**—Astrophysics, Exceptional Model Mining, Cherenkov radiation, soft classifier.

N.B.: A significantly longer version of this paper appeared as a technical report of the TU Dortmund [1].

## I. INTRODUCTION

The FACT telescope [2], [3] is an Imaging Air Cherenkov Telescope, designed to detect light emitted by secondary particles, generated by high-energetic cosmic particles interacting with the atmosphere of the Earth. For astrophysical reasons, it is important to classify the light as resulting from the atmosphere being hit by a gamma ray or a proton; the latter occur much more frequently, but the former are the more interesting in gamma astronomy (which will be discussed later in the paper). Currently, one of the used classifiers is a random forest, whose performance needs our detailed attention.

The problem with training a classifier on real astrophysical data is that there is no clear feedback. Based on the observed light, we could deduce whether the inducing particle is a gamma ray or a proton. Then, we can look in the direction from which the particle originated, and strive to find an astrophysical source generating gamma rays. But even if we find such a source, there is no certain way of telling what kind of particle

induced the original observation. Effectively, we are dealing with a feedbackless learning task, and it is typically hard to finetune a classifier without feedback.

To study our learning performance, we turn to Monte Carlo data. We simulate particle interactions with the atmosphere, as well as reflections of the resulting Cherenkov light with telescope mirrors on the one hand and the FACT camera electronics on the other hand. This gives us a dataset of camera images that is equivalent in form to a dataset we would get from real astrophysical observations, except that we also know the true label of our classification task. By training our random forest on this dataset, we obtain the soft classifier probabilities for each record. Through studying the interaction between the binary ground truth that we already knew and the soft classifier probabilities we learned from the data, we can understand where our classifier performs exceptionally.

We study this interaction with an Exceptional Model Mining (EMM) [4], [5] approach. This is a Local Pattern Mining framework, specialized in finding coherent subsets of the dataset where multiple targets interact in an unusual way. In this paper, we introduce the SCaPE (Soft Classifier Performance Evaluation) model class for EMM, seeking subgroups for which a soft classifier represents a ground truth exceptionally well or poorly. This should allow a domain expert to understand where his/her classifier does (not) work.

## II. PRELIMINARIES

Before we can introduce the new contributions of this paper, we need to cover a lot of preliminary ground. The preliminaries have been split up into three parts: the next subsection contains an introduction of astrophysical concepts, the subsection thereafter summarizes Local Pattern Mining methods including EMM, and the technical report [1] contains a short note on the alignment of soft and hard classifiers.

### A. The FACT Telescope

An important task in astroparticle physics is observing distant astrophysical sources such as Supernova Remnants (SNR) or Active Galactic Nuclei (AGN) in multiple energy ranges (optical, radio, X-ray, gamma rays), since combining such observations helps us understand (amongst others) the cosmic particle acceleration and radiation emission mechanisms of these sources [2]. Each energy range demands different detector techniques, hence dedicated telescopes are required. In

the high-energy regime, we are interested in (ultra-)relativistic cosmic particles such as gamma rays, neutrinos, and protons, which are assumed to be accelerated by astrophysical sources (such as SNR and AGN). Gamma rays are interesting because of their neutral electric charge, which causes them to travel undeflected by intergalactic magnetic fields. This means that the direction from which the primary gamma rays are coming, necessarily points directly to the astrophysical source.

The Earth’s atmosphere is only transparent in optical and radio wavelengths. This prohibits observing gamma rays on Earth, but we can make these observations with dedicated satellites. Since the gamma ray flux (amount of particles per area and time) decreases with higher energy, detecting gamma rays in higher energy ranges would require either a bigger detection area (in the satellite) or more time. Both solutions are not satisfying, as time and bigger satellites are prohibitively cost-intensive. Instead, we can use an effect caused by very high-energetic particles propagating through the atmosphere.

When very high-energetic cosmic particles such as gamma rays and protons interact with the atmosphere of the Earth, they induce an extensive air shower consisting of secondary relativistic particles, which can be charged. The charged particles emit Cherenkov radiation [6], a blueish light which can be detected by ground-based Imaging Air Cherenkov Telescopes (IACT). One such telescope is FACT, the First G-APD Cherenkov Telescope. It is located on La Palma, Canary Islands, Spain at 2200m above sea level, and is operational since October 2011 [3]. FACT is the first IACT using Geiger-mode Avalanche PhotoDiodes (G-APD) (also known as silicon photomultipliers) as photosensors to detect Cherenkov light. Contrary to conventional detector techniques of IACTs, G-APDs allow to observe even during strong moonlight and thus increase the effective observation time. This is especially interesting for source detection by small telescopes, but also very important for long-term monitoring of sources.

As we observe a variability in the gamma ray flux of sources in multiple timescales (both seconds and years) [3], long-term monitoring is required to understand the emission procedures and mechanisms within and surrounding the sources. The primary physics goal of FACT is therefore to observe the brightest known VHE sources on long timescales, which becomes realizable by using G-APDs.

The main goal of the analysis method whose results are evaluated in this paper is to find gamma-induced showers. Unfortunately, for the brightest sources, proton showers appear a thousand times more frequently than gamma showers in the source direction [7], which makes the light of the proton-induced showers the biggest background. Therefore, the separation of gamma- and proton-induced showers is very important to be able to detect a source, to increase the sensitivity of the telescope and thus the effective observation time, and finally to measure the spectrum of the source. For the separation, Monte Carlo simulations are necessary, which simulate shower images in the FACT camera with known parameters, such as type and energy of the primary particle that induced the shower. The first step is to simulate particle interactions in the atmosphere and the emission of Cherenkov light with the program MMCS based on CORSIKA [8]. Further processing by a simulation and analyzing tool called MARS [9] includes simulating the reflection of the light on the

mirrors of the telescope and the electronics inside the camera. We end up with simulated camera images containing gamma and proton showers. From these camera images the image parameters of the showers are reconstructed. Since gamma- and proton-induced showers have distinctive shapes, the image parameters describing the properties of the shower images are used to distinguish between them. As the information of the primary particles is known in the simulation, the data are labeled as *true* or 1 for gamma showers (signal) and *false* or 0 for proton showers (background).

As is commonly done in IACT experiments (cf. [1]), the separation is done with a random forest (RF) algorithm [10]. We employ an implementation available within the Rapid-Miner analytics platform [11]. The RF builds a model with the image parameters of the labeled simulated data and tests it on the remaining dataset in a five-fold cross-validation to ensure a stable classification. For this dataset 500 trees were grown, each considering a random subset of 8 out of the 11 available attributes. These 11 attributes contain parameter distributions for gamma and proton showers, which are known to be crudely separable by simple cuts on each parameter relatively successfully. The fact that just a subset of attributes is drawn contributes to the randomized trees needed for a good random forest. Each tree classifies an event (one shower) as 1 for signal or 0 for background. Prediction aggregation over all trees is done by averaging, the value of which is called the *Signalness*. This quantity describes the probability or the confidence of the RF for an event to be classified as a gamma shower. For the given FACT dataset the efficiency decreases with a higher Signalness value, but at the same time the purity increases. To separate gamma and proton showers sufficiently while not losing too much data, a cut has to be found which fulfills both conditions and depends on the physics task.

## B. Exceptional Model Mining

*Pattern mining* [12] is the broad subfield of data mining where only a part of the data is described at a time, ignoring the coherence of the remainder. The goal is finding subsets  $S$  of the dataset  $\Omega$  that are interesting somehow:

$$S \subseteq \Omega \Rightarrow \text{interesting}$$

Typically, not just any subset of the data is sought after: only those subsets that can be formulated using a predefined *description language*  $\mathcal{L}$  are allowed. A canonical choice for the description language is conjunctions of conditions on attributes of the dataset. If, for example, the records in our dataset describe people, then we can find results of the following form:

$$\text{Age} \geq 30 \wedge \text{Smoker} = \text{yes} \Rightarrow \text{interesting}$$

Allowing only results that can be expressed in terms of attributes of the data, rather than allowing just any subset, ensures that the results are relatively easy to interpret for a domain expert: the results arrive at his doorstep in terms of quantities with which he should be familiar. A subset of the dataset that can be expressed in this way is called a *subgroup*.

In the FACT telescope setting, we strive to separate the gamma sources from the proton sources; there is a clear target, hence this setting is supervised. The most extensively studied form of supervised pattern mining is known as *Subgroup Discovery* (SD) [13], where one (typically binary) attribute  $t$

of the dataset is singled out as the *target*. The goal is to find subgroups for which the distribution of this target is unusual: if the target describes whether the person develops lung cancer or not, we find subgroups of the following form:

$$\text{Smoker} = \text{yes} \Rightarrow \text{lung cancer} = \text{yes}$$

*Exceptional Model Mining* (EMM) [4], [5] can be seen as the multitarget generalization of SD. Rather than singling out one attribute as the target  $t$ , in EMM there are several target attributes  $t_1, \dots, t_m$ . Interestingness is not merely gauged in terms of an unusual *marginal* distribution of  $t$ , but in terms of an unusual *joint* distribution of  $t_1, \dots, t_m$ . Typically, a particular kind of unusual *interaction* between the targets is captured by the definition of a *model class*, and subgroups are deemed interesting when their model is exceptional, which is captured by the definition of a *quality measure*. For example, suppose that there are two target attributes: a person's length ( $t_1$ ), and the average length of his/her grandparents ( $t_2$ ). We may be interested in the correlation coefficient between  $t_1$  and  $t_2$ ; we then say we study EMM with the *correlation model class* [4]. Given a subset  $S \subseteq \Omega$ , we can estimate the correlation between the targets within this subset by the sample correlation coefficient. We denote this estimate by  $r^S$ . Now we can define the following quality measure (tweaked from [4]):

$$\varphi(S) = |r^S - r^\Omega|$$

EMM then strives to find subgroups for which this quality measure has a high value: effectively, we search for subgroups coinciding with an exceptional correlation between a person's length and his/her grandparents' average length:

$$\text{Lives near nuclear plant} = \text{yes} \Rightarrow |r^S - r^\Omega| \text{ is high}$$

### III. RELATED WORK

Previous work exists on discovering subgroups displaying unusual interaction between multiple targets, for instance in the previously developed model classes for EMM: correlation, regression, Bayesian network, and classification (cf. [4], [5]). The last of these model classes is particularly related to the SCaPE model class, with two major differences. On the one hand, the model class definitions imply a different relation between the subgroup definitions and classifier search space. The classification model class takes both classifier input and output attributes as targets for the EMM run. This disallows those attributes to show up in the descriptions of subgroups found with EMM; exceptional subgroups are described in terms of attributes unavailable to the classifier. By contrast, in the SCaPE model class, all attributes available as input (but not as output!) to the classifier are also available for describing subgroups. Hence, the found unusual subgroups directly correspond to a subspace in the classifier search space. On the other hand, the model classes search for a different underlying concept in the dataset. The classification model class *investigates* classifier *behavior* in the *absence* of a ground truth. The SCaPE model class *evaluates* classifier *performance* in the *presence* of a ground truth. Hence, the two model classes are different means to achieve different ends.

Automated guidance to improve a classifier has been studied in the data mining subfield of meta-learning. The exact meaning of this term is subject to debate; see [14] for a survey discussing some of the views. A constant factor is

that meta-learning hovers around the question how knowledge about learning can be put to use to improve the performance of a learning algorithm. A typical approach is to let the machine compute meta-features characterizing the data, such as correlations between attributes, attribute entropy, and mutual information between class and attributes. These meta-features are then considered in a new classifier training phase, and the hope is that this improves predictive performance. This process is depicted in the self-adaptive learning flow diagram in [14, Figure 2]. The meta-features can also be employed to compare learning algorithms. For instance, Henery [15] provides a set of rules to determine when the one learning algorithm is significantly better than the other. However, in almost all of the existing meta-learning work, the focus is on letting the machine learn how the machine can perform better.

By contrast, Vanschoren and Blockeel [16] express an interest in *understanding* learning behavior. Their paper discusses a descriptive form of meta-learning, proposing an integrated solution (using experiment databases) that aims to explain the behavior of learning algorithms. This explanation is again expressed in terms of meta-features; no investigation takes place of particular subspaces of the search space on which the algorithm performs exceptionally. While Vilalta and Drissi [14, Section 4.3.1] do devote a subsection to "Finding regions in the feature space [...]", this is again in the context of algorithm selection. Their innovation lies in allowing different learning algorithms for different records of the dataset. Meta-learning is related to the goals we strive to achieve with the SCaPE model class for EMM, but two things set these approaches apart: meta-learning focuses on meta-features, while the SCaPE model class focuses on coherent subspaces of the original search space, and meta-learning focuses on letting the machine improve the predictive performance of the machine, while the SCaPE model class focuses on providing *understanding* to the domain expert where his/her classifier works well or fails. As such, the SCaPE model class for EMM provides progress on the path sketched by Vanschoren and Blockeel in the conclusions of their paper [16, Section 5]: "We hope to advance toward a meta-learning approach that can explain not only *when*, but also *why* an algorithm works or fails [...]".

A very recent first inroad towards peeking into the classifier black box is the method by Henelius et al. [17], who strive to find groups of attributes whose interactions affect the predictive performance of a given classifier. This is more akin to the classification model class for EMM. While Henelius et al. study hard classifiers, the SCaPE model class is designed for soft classifiers.

### IV. MAIN CONTRIBUTION

The main contribution of this paper is the development of a new model class with associated quality measure for Exceptional Model Mining: the SCaPE (Soft Classifier Performance Evaluation) model class. In this model class, two targets are identified: a binary target  $b$  describing the ground truth, and a real-valued target  $r$  containing the output of a soft classifier that strives to approximate  $b$ . The goal in this model class is to find subgroups for which this soft classifier represents the ground truth exceptionally well or exceptionally poorly. Notice that, SCaPE being an EMM model class, the focus is on easily-interpretable subgroups. Hence, our primary goal is not to let

the machine improve the machine, but to let the domain expert *understand* where his/her classifier does or does not work.

## V. THE SCaPE MODEL CLASS FOR EMM

In the SCaPE model class for EMM, we assume a dataset  $\Omega$ , which is a bag of  $N$  records of the form  $x = (a_1, \dots, a_k, b, r)$ . We call  $\{a_1, \dots, a_k\}$  the *descriptive attributes*, or *descriptors*, whose domain is unrestricted. The remaining two attributes,  $b$  and  $r$ , are the *targets*. The first,  $b$ , is the *binary target*; we will denote its values by 0 and 1. The second,  $r$ , is the *real-valued target*, taking values in  $\mathbb{R}$ .

The goal of the SCaPE model class is to find subgroups for which the soft classifier outputs, as captured by  $r$ , represent the ground truth, as captured by  $b$ . In Section V-A, we develop measures that quantify how well  $b$  is represented by  $r$ , on the entire dataset and on subsets of the dataset. In Section V-B, we use these measures to define a *quality measure* for the SCaPE model class, that gauges how exceptional the interplay between  $r$  and  $b$  is on a subgroup when compared to this interplay on the entire dataset.

If we need to distinguish between particular records of the dataset, we will do so by superscripted indices:  $x^i$  is the  $i^{\text{th}}$  record,  $b^i$  is its value for the binary target and  $a_j^i$  is its value for the  $j^{\text{th}}$  descriptor. For the sake of notational convenience, we assume that the records are indexed in non-descending order by their values of  $r$ :  $i < j \Rightarrow r^i \leq r^j$ . We call the records  $x^i$  in the dataset for which the binary target is true the *positives*, and the other records the *negatives*.

### A. Average (Sub-)Ranking Loss

A soft classifier can be converted into a hard classifier by imposing a threshold at any chosen value  $v$ : the predicted label for record  $x^i$  is set to 1 if and only if  $r^i > v$ . This value  $v$  should be chosen such that the hard classifier based on  $r$  lines up reasonably well with the ground truth as provided by  $b$ ; by and large, high values for the real-valued target should coincide with  $b = 1$ , and low values with  $b = 0$ . Notice that this capability of  $r$  is primarily sensitive not to its precise values, but to the *ordering* it implies on the records. Therefore, we capture the alignment of  $r$  and  $b$  on the whole dataset by the Average Ranking Loss [18]:

$$\text{ARL}(\Omega) = \frac{\sum_{i=1}^N \left( \mathbb{I}\{b^i = 1\} \cdot \sum_{j=i+1}^N \mathbb{I}\{b^j = 0\} \right)}{\sum_{i=1}^N \mathbb{I}\{b^i = 1\}} \quad (1)$$

Essentially, for every positive in the dataset a penalty is computed. The penalty for  $x^i$  is equal to the number of negatives  $x^j$  that have a higher value for the real-valued target:  $r^i < r^j$  (here, the formula for ARL uses the fact that the dataset is ordered non-descendingly by  $r^i$ , and conveniently ignores for the moment that two consecutive  $r$ -values may be equal). This *ranking loss* is then averaged over all positives in the dataset, arriving at the ARL. Obviously, lower values of the ARL correspond to a better representation of  $b$  by  $r$ . To determine the degree of representation of  $b$  by  $r$  in a given subgroup  $S$  of the dataset, we compute the ARL again, but then restricted to just those records belonging to the subgroup.

We call this the *Average Subranking Loss* of  $S$ , denoted by  $\text{ASL}(S)$ . For illustration purposes, in [1] we provide a toy example dataset with directions on computing the ARL and ASL, and other measures from the remainder of this paper.

1) *Handling Ties*: So far, we have assumed that all values for the real-valued target  $r$  in the dataset are distinct. This simplifies the formula in Equation (1), and allows for an easier intuitive explanation in that section. In practice, of course, such an assumption is undesirable. Since we compute the ARL/ASL as an average of penalties assigned to all positives, we can focus on how to update the penalty assigned to a positive when its  $r$ -value is replicated in the dataset. Suppose that  $x^i$  is such a positive: we know that  $b^i = 1$  and  $r^i = r^j$  for some  $j \neq i$ . If  $x^j$  is also a positive, then the penalty does not need to change. If, on the other hand,  $x^j$  is a negative, then we should increment the penalty by some amount; we will add  $1/2$  to the penalty for  $x^i$  for each such tie, which is extensively motivated in [1]. Incorporating this penalty leads to the following definitions:

**Definition (Average (Sub-)Ranking Loss).** The *Average Ranking Loss*,  $\text{ARL}(\Omega)$ , of a dataset  $\Omega$  is given by:

$$\text{ARL}(\Omega) = \frac{\sum_{i=1}^N \mathbb{I}\{b^i = 1\} \cdot \text{PEN}_i^N(\Omega)}{\sum_{i=1}^N \mathbb{I}\{b^i = 1\}}$$

where the *penalty* for the  $i^{\text{th}}$  record,  $\text{PEN}_i^N(\Omega)$ , is given by:

$$\begin{aligned} \text{PEN}_i^N(\Omega) &= \sum_{j=i+1}^N \mathbb{I}\{b^j = 0 \wedge r^j > r^i\} \\ &+ \frac{1}{2} \sum_{j=i+1}^N \mathbb{I}\{b^j = 0 \wedge r^j = r^i\} \end{aligned}$$

The *Average Subranking Loss*,  $\text{ASL}(S)$ , of a subgroup  $S$  of  $\Omega$  is given by:

$$\text{ASL}(S) = \text{ARL}(\Omega')$$

where  $\Omega'$  is the dataset constructed by taking from  $\Omega$  only those records belonging to  $S$ .

### B. Quality Measure: Relative Average Subranking Loss

In EMM we strive to find subgroups for which the target interaction captured by the model class is exceptional. This exceptionality is gauged by a quality measure. We define a quality measure for the SCaPE model class, whose maxima, minima, and extremities correspond to three distinct goals.

**Definition (Relative Average Subranking Loss).** The *Relative Average Subranking Loss*,  $\varphi_{\text{rasl}}$ , of a subgroup  $S$  of  $\Omega$  is given by:  $\varphi_{\text{rasl}}(S) = \text{ASL}(S) - \text{ARL}(\Omega)$

To find subgroups for which  $r$  represents  $b$  *poorly*, i.e., subgroups for which the soft classifier *does not work*, one should *maximize*  $\varphi_{\text{rasl}}$ ; positive values for  $\varphi_{\text{rasl}}$  indicate that the soft classifier performs worse than usual on this subgroup. To find subgroups for which  $r$  represents  $b$  *well*, i.e., subgroups for which the soft classifier *does work*, one should *minimize*  $\varphi_{\text{rasl}}$ ; negative values for  $\varphi_{\text{rasl}}$  indicate that the soft classifier performs better than usual on this subgroup. To find a list of subgroups for which the soft classifier performs exceptionally (in general), one should maximize  $|\varphi_{\text{rasl}}|$ .

## VI. EXPERIMENTAL RESULTS

The SCaPE model class for EMM requires a binary and a real-valued target for real-world experiments. For this purpose we use the FACT Monte Carlo Simulation for gamma- and proton-induced air showers, as the binary target is already present by the information of the primary particle. The real-valued target is generated in RapidMiner by the random forest (RF) classifier, as it can produce probabilities of being a gamma shower expressed by the Signalness (cf. Section II-A). The RF algorithm is implemented and used as a separation method in other IACT experiments (cf. [1]), where it has proven to be a stable and robust method performing comparatively superior to classical methods.

Disjoint Monte Carlo datasets were generated for training and testing the RF. The training sets for the individual trees containing gamma and proton showers were sampled in such a way that they have the same size. The dataset contains simulated reconstructed image parameters and source-dependent parameters which allow to estimate a statistical signal of the astrophysical source at which the telescope is pointing.

The SCaPE model class itself is implemented in Cortana [19], a toolbox featuring a plethora of Subgroup Discovery and Exceptional Model Mining settings. On this FACT dataset, we run Cortana twice: once maximizing and once minimizing  $\varphi_{\text{rasl}}$ . The Average Ranking Loss on the whole dataset is 1,446.761. For more experimental results on nine UCI datasets (including an inspection of subgroups of mushrooms), more details on the parametrization of Cortana, and directions for the interested reader to obtain the implementation and the FACT dataset, see [1].

### A. Experimental Results — Maximizing $\varphi_{\text{rasl}}$

When maximizing  $\varphi_{\text{rasl}}$ , we strive to find subgroups on which the classifier performs poorly. The top-eight found subgroups are listed in Table I. As the last column shows, the first three subgroups have a substantially worse Average Subranking Loss than the rest, so they warrant further investigation. These three subsets are described by two distinct attributes. Both are source-dependent parameters, and between them they are strongly correlated.

The parameter ThetaSq describes the distance of the reconstructed source position to the real source position. Thus, near-zero values express that the corresponding shower points to the real astrophysical source. We see the same behavior for the parameter dca, which describes the distance of the closest approach of the shower to the source position with respect to the x-axis. Again, showers with near-zero values have a higher probability of coming directly from the real source.

In the Monte Carlo simulations, gamma showers are assumed and simulated as if they were coming directly from the source, since this is the case in the real world we are interested in. In real data we also have a minor fraction of diffuse gamma showers, coming from sources other than the observed astrophysical source; these are not taken into account in the simulations. By contrast, proton-induced showers are assumed to be isotropically distributed in the sky. Taking this information into account we can easily explain why the classifier performs particularly poorly on the first three

TABLE I. SUBGROUPS ON THE FACT DATASET MAXIMIZING  $\varphi_{\text{rasl}}$

Rank	Worst-classified subgroups $S$	$\varphi_{\text{rasl}}(S)$
1.	dca $\geq$ 79.2745	1294.939
2.	ThetaSq $\geq$ 0.136131	1116.781
3.	dca $\leq$ -68.3173	1114.739
4.	SizeArea $\leq$ 0.5564718	100.786
5.	MCMomentumZ $\leq$ -1618.63	59.373
6.	cut1 = 0	46.957
7.	MCEnergy $\geq$ 1641.69	39.205
8.	ConclSize $\leq$ 39.874977	28.153

subgroups in Table I. In both involved parameters, the gamma showers are accumulated around low values, while proton showers are equally distributed over the full parameter value range. Thus, the gamma showers decrease in frequency for higher values. For instance, the two subgroups for the dca parameter encompass just  $\sim 10^{-5}$  % of the gamma events in the whole dataset. While training the RF, one source-dependent parameter was used. This means that the classifier learned that the probability of being a gamma shower is high with low values in ThetaSq and dca. Conversely, the classification gets tougher if we have only a small number of gamma showers with high values in ThetaSq and dca. For a detailed investigation of these subgroups, involving their distribution of positives and negatives, see [1].

The subgroups in Table I with less extreme values for  $\varphi_{\text{rasl}}$ , such as the ones with rank 4 and 8, are less straightforward to explain. The parameter SizeArea describes the compactness of the deposited light of the showers and the parameter ConclSize describes the deposited light in the brightest pixel of a shower. The higher these values are, the more likely it is that we are dealing with a gamma shower. On first look, the poor classification on these particular subgroups is surprising, because the parameter distributions are clearly separated for lower values of gamma and proton showers as well. However, this result could be explained by internal cuts in the RF, which affects the distributions and tends to misclassify events with a lower probability of being a gamma.

### B. Experimental Results — Minimizing $\varphi_{\text{rasl}}$

When minimizing  $\varphi_{\text{rasl}}$ , we strive to find subgroups on which the classifier performs well. The top-eight such subgroups are listed in Table II.

The first and eighth-ranked subgroup are described by the same parameter cosdeltaalpha, which is again source-dependent and roughly expresses the cosine of the angle between the shower main axis and the source position. Thus, values of cosdeltaalpha around 1 or -1 indicate that the shower axis is pointing to the source, which also means a higher probability for the shower to come directly from the source and thus a higher probability of being a gamma shower. Contrary to dca, which appears high-ranked in the poorly-classified subgroups, these well-classified subgroups contain a big fraction of gamma showers compared to the fraction of proton showers. This means that the classifier learns that showers which are contained in these subgroups are very likely gamma showers and are better classified than in other ranges.

The third-ranked subgroup is the known source-dependent parameter ThetaSq. It appears in the well-classified subgroups with very low values as well as in the poorly-classified

TABLE II. SUBGROUPS ON THE FACT DATASET MINIMIZING  $\varphi_{\text{RASL}}$ 

Rank	Best-classified subgroups $S$	$\varphi_{\text{rasl}}(S)$
1.	$\text{cosdeltaalpha} \geq 0.999994$	-1446.259
2.	$\text{SizeSinglePixels} \geq 372.953$	-1445.761
3.	$\text{ThetaSq} \leq 6.57561\text{E-4}$	-1445.753
4.	$\text{Length} \leq 9.70734$	-1445.336
5.	$\log\text{Length} \leq 0.98710024$	-1445.336
6.	$\text{NumberSinglePixels} \geq 73.0$	-1444.539
7.	$\text{SizeArea} \geq 1.8111843$	-1444.535
8.	$\text{cosdeltaalpha} \leq -0.999995$	-1444.275

subgroups with higher values. This behavior is perfectly explainable, as very low values indicate a higher probability of being a gamma shower, and the probability decreases slowly with higher ThetaSq values, until a value is reached where gamma showers cannot be distinguished well from the proton showers if only ThetaSq is taken into account. We see the same effect with the seventh-ranked subgroup described by SizeArea. The classifier performs well on higher values but worse on lower values. Again, this result could be explained by internal cuts in the RF.

## VII. CONCLUSIONS

Motivated by a real-life astrophysics data scenario, we introduce the SCaPE (Soft Classifier Performance Evaluation) model class for Exceptional Model Mining (EMM). SCaPE strives to find coherent subgroups displaying exceptional interaction between the probabilities provided by a soft classifier and a binary ground truth. This interaction is evaluated by the Average (Sub-)Ranking Loss, a quantity expressing how well the soft classifier probabilities can represent the binary ground truth. The quality measure  $\varphi_{\text{rasl}}$  is designed to find coherent subspaces of the dataset where the soft classifier performs poorly (when maximizing  $\varphi_{\text{rasl}}$ ), well (when minimizing  $\varphi_{\text{rasl}}$ ), or exceptionally (when maximizing  $|\varphi_{\text{rasl}}|$ ). The focus of EMM lies on finding easily interpretable subgroups. Hence, as opposed to a meta-learning framework, which is focused on letting the machine improve the machine, the primary goal in the SCaPE model class for EMM is to provide a better *understanding* to the domain expert. We want the expert to be able to understand where his/her classifier does or does not work well, by reporting the problem and success areas in familiar terms.

We perform real-world experiments with the SCaPE model class on an astrophysics dataset concerned with the classification of air showers induced by high-energetic cosmic particles. The subgroups with the most deviating Average Subranking Losses — both the poorly-classified ones and the well-classified ones — have an astrophysical interpretation corroborating their appearance as a particularly (un-)problematic subspace of the search space. Subgroups with less extreme but still high/low values for the quality measure are non-trivial to explain and deserve a closer look. The results show that the random forest classifier performs better when the incidence of gamma showers is higher.

In gamma ray astronomy, the separation of gamma and proton showers marks an important step in the analysis of astrophysical sources. Better classifier performance leads to less dilution of the interesting physics results and improves the statement of results of the astrophysical source. The result set will more frequently contain the infrequently appearing

gamma showers, which should increase the effective observation time. Due to the importance of the separation in this field, understanding why the classifier does not perform as desired is extremely valuable. The SCaPE model class for EMM helps to understand the classification, which leads to ideas on how to improve the overall classifier performance.

## ACKNOWLEDGMENTS

This research is supported in part by the Deutsche Forschungsgemeinschaft (DFG) within the Collaborative Research Center SFB 876 “Providing Information by Resource-Constrained Analysis”, project C3.

## REFERENCES

- [1] W. Duivesteijn, J. Thaele, Understanding Where Your Classifier Does (Not) Work — the SCaPE Model Class for Exceptional Model Mining, technical report 09/2014 of SFB876 at TU Dortmund, 2014.
- [2] H. Anderhub, M. Backes, A. Biland et al., Design and Operation of FACT — the First G-APD Cherenkov Telescope, arXiv:1304.1710 [astro-ph.IM]
- [3] T. Bretz, H. Anderhub et al., FACT — The First G-APD Cherenkov Telescope: Status and Results, arXiv:1308.1512 (astro-ph.IM)
- [4] D. Leman, A. Feelders, A.J. Knobbe, Exceptional Model Mining, Proc. ECML/PKDD (2), pp. 1–16, 2008.
- [5] W. Duivesteijn, Exceptional Model Mining, PhD thesis, Leiden University, 2013.
- [6] C. Grupen, Astroteilchenphysik: Das Universum im Licht der kosmischen Strahlung, Vieweg, 2000.
- [7] S.F. Taylor, T. Abu-Zayyad, K. Belov et al., The Highest Energy Cosmic Rays and Gamma Rays, American Astronomical Society, 192nd AAS Meeting, # 09.03; Bulletin of the American Astronomical Society 30, p. 827, 05/1998.
- [8] CORSIKA - An Air Shower Simulation Program, <https://web.ikp.kit.edu/corsika/>
- [9] T. Bretz, D. Dorner, MARS - CheObs ed. — A flexible Software Framework for future Cherenkov Telescopes, Astroparticle, Particle and Space Physics, Detectors and Medical Physics Applications, pp. 681–687, 2010.
- [10] L. Breiman, Random Forests, Machine Learning 45, pp. 5–32, 2001.
- [11] I. Mierswa, M. Wurst, R. Klinkenberg, M. Scholz, T. Euler, YALE: Rapid Prototyping for Complex Data Mining Tasks, Proc. KDD, pp. 935–940, 2006.
- [12] K. Morik, J.F. Boulicaut, A. Siebes (eds), Local Pattern Detection, Springer, New York, 2005.
- [13] F. Herrera, C.J. Carmona, P. González, M.J. del Jesus, An Overview on Subgroup Discovery: Foundations and Applications, Knowledge and Information Systems 29 (3), pp. 495–525, 2011.
- [14] R. Vilalta, Y. Drissi, A Perspective View and Survey of Meta-Learning, Artificial Intelligence Review 18 (2), pp. 77–95, 2002.
- [15] R.J. Henery, Methods for Comparison, in: D. Michie, D.J. Spiegelhalter, C.C. Taylor (eds.), Machine Learning, Neural and Statistical Classification, Ellis Horwood, 1994.
- [16] J. Vanschoren, H. Blockeel, Towards Understanding Learning Behavior, Proc. BENELEARN, pp. 89–96, 2006.
- [17] A. Henelius, K. Puolamäki, H. Boström, L. Asker, P. Papapetrou, A peek into the black box: exploring classifiers by randomization, Data Mining and Knowledge Discovery 28 (5-6), pp. 1503–1529, 2014.
- [18] G. Tsoumakas, I. Katakis, I.P. Vlahavas, Mining Multi-Label Data, Data Mining and Knowledge Discovery Handbook, Springer, pp. 667–685, 2010.
- [19] M. Meeng, A.J. Knobbe, Flexible Enrichment with Cortana – Software Demo. Proc. Benelearn, pp. 117–119, 2011.