Analysis of Electron Distribution Functions using the Gaussian Mixture Model

Beniamino Sanò¹, Nathan N. Maes², David L. Newman³, Martin V. Goldman³, Francesco Valentini⁴, and Giovanni Lapenta⁵

¹University of Trento / University of Calabria ²University of Leuven ³University of Colorado Boulder ⁴Dipartimento di Fisica, Università della Calabria ⁵Univ Leuven, KU Leuven, Dept Wiskunde

April 16, 2024

Abstract

Velocity distribution functions (VDFs) measured by the Magnetospheric Multiscale (MMS) mission are complex 3D datasets that can be represented as a superposition of multiple beams (M. V. Goldman et al., 2020). A recent work (Dupuis et al., 2020) proposed the use of the Gaussian Mixture Model (GMM). Here we investigate the approach by considering first synthetic distributions made by artificially creating beams of either Maxwellian distributions or kappa distributions with varying power law index. By varying the inter-beam average difference and the beam standard deviation we evaluate the ability of the GMM in recognizing correctly the beam. We then apply the method systematically to MMS data in the tail and in the dayside. In this case, the data need preparation before being processed by the GMM to account for the specifics of the instrument and in particular the lack of data at low energy and to account for the noise in the counts. The conclusion of the analysis is that the GMM is capable of detecting the presence of multiple beams when their distinction is significant. The GMM can define reliably the complexity of a measured dataset in terms of the number of optimal beams provided by information theory criteria. Visual inspection confirms this automatic definition of complexity.

Analysis of Electron Distribution Functions using the Gaussian Mixture Model

1

2

10

Beniamino Sanò^{1,2,3}and Nathan N. Maes⁴and David L. Newman⁵and Marty Goldman⁵and Francesco Valentini²and Giovanni Lapenta⁴ ¹University of Trento, Italy ²Università della Calabria, Italy ³ASI - Italian Space Agency, Italy ⁴Center for mathematical Plasma Astrophysics (CmPA), Department of Mathematics, KULeuven, University of Leuven, Belgium

 $^5\mathrm{University}$ of Colorado, USA

Key Points: 11 • The Gaussian Mixture Model (GMM) is assessed in its ability to identify multi-12 ple Gaussian and kappa-distributed beams depending on their relative average means 13 and standard deviations. 14 • The Gaussian Mixture Model (GMM) can be used to define the complexity of a 15 velocity distribution function based on the optimal number of beams determined 16 by information theory criteria. 17 • The Gaussian Mixture Model (GMM) is applied to burst intervals of Magneto-18 spheric Multiscale (MMS) mission for electrons in the dayside and nightside with 19 different counts levels and noise to signal ratios. 20

Corresponding author: Beniamino Sanò, beniamino.sano@unitn.it

21 Abstract

Velocity distribution functions (VDFs) measured by the Magnetospheric Multiscale (MMS) 22 mission are complex 3D datasets that can be represented as a superposition of multiple 23 beams (M. V. Goldman et al., 2020). A recent work (Dupuis et al., 2020) proposed the 24 use of the Gaussian Mixture Model (GMM). Here we investigate the approach by con-25 sidering first synthetic distributions made by artificially creating beams of either Maxwellian 26 distributions or kappa distributions with varying power law index. By varying the inter-27 beam average difference and the beam standard deviation we evaluate the ability of the 28 GMM in recognizing correctly the beam. We then apply the method systematically to 29 MMS data in the tail and in the dayside. In this case, the data need preparation before 30 being processed by the GMM to account for the specifics of the instrument and in par-31 ticular the lack of data at low energy and to account for the noise in the counts. The 32 conclusion of the analysis is that the GMM is capable of detecting the presence of mul-33 tiple beams when their distinction is significant. The GMM can define reliably the com-34 plexity of a measured data-set in terms of the number of optimal beams provided by in-35 formation theory criteria. Visual inspection confirms this automatic definition of com-36 plexity. 37

³⁸ Plain Language Summary

This work investigates regions of interest in electrons distribution functions from 39 Magnetospheric Multiscale (MMS) mission, using an unsupervised machine learning tech-40 nique called Gaussian Mixture Model (GMM). First we tested the ability of the GMM 41 to identify multiple Gaussian and kappa-distributed beams on synthetic distributions, 42 and then we analysed real data from MMS. The data is downloaded and preprocessed 43 through AIDApy, a Python package for the analysis of spacecraft data from heliospheric 44 missions. A Gaussian mixture model search through the particles and identify the pres-45 ence of different subpopulations whithin an overall population. The optimal number of 46 subpopulations is determined by a model selection technique, and the presence of cer-47 tain distributions can be utilized to find magnetic reconnection regions. 48

49 **1** Introduction

The study of the Earth's magnetosphere and its complex system of electromagnetic interactions is a key goal in understanding the fundamental physics of space. Magnetic

-2-

reconnection and plasma turbulence are both closely interrelated fundamental processes 52 in the dynamics of the magnetosphere (Biskamp, 2000). Magnetic reconnection is a pro-53 cess during which the magnetic field energy is converted into kinetic energy, thermal en-54 ergy, and particle acceleration energy. Reconnection occurs in small-scale electron dif-55 fusion regions within a current sheet (Lapenta et al., 2016). As the field lines flow into 56 the region, they reconnect at the X-point. The reconnected field has a strong magnetic 57 tension, which pulls the reconnected field away from the X-point, expelling the plasma 58 coupled to it as bi-directional outflow jets (Li et al., 2021). Plasma turbulence is the re-59 sult of multi-scale nonlinear interactions and instabilities of large-scale fluid motions. Col-60 lisionless space plasmas are often in a turbulent non-equilibrium state, characterized by 61 strong fluctuations of field and plasma parameters (Scott, 2021). Turbulence and recon-62 nection research is focused on how magnetic reconnection occurs in a turbulent system 63 and how the dynamics of turbulence and reconnection interact (Yokoi & Hoshino, 2011). 64

To study this relation, several spacecraft have been sent into space in recent years. 65 Cluster mission observed for the first time in-situ magnetic reconnection in turbulent plasma 66 (Retinò et al., 2007). NASA's Magnetospheric Multiscale (MMS) mission has the goal 67 of observing at an unprecedented rate traces of magnetic reconnection in Earth's mag-68 netosphere (Burch et al., 2016). Fast Plasma Investigation (FPI) instrument measures 69 incoming particles through a filter which selects certain particle speeds and directions; 70 then a 3D picture of the ion plasma is produced every 150 milliseconds, while for elec-71 tron plasma FPI captures a picture every 30 milliseconds. Because of these frame rates, 72 MMS measures more than 100 GB of data every day. However, due to limitations of the 73 probes, a continuous overwriting of data takes place and a large part of it is lost irre-74 versibly: in fact, approximately 4 GB of data per day are transmitted to Earth. (Baker 75 et al., 2016) At first the task of looking at the raw data and selecting the interesting ones 76 was done by so-called scientists in the loop, who would observe the data by eve and se-77 lect which ones to store. Nowadays, due to the size of the measurements, this kind of 78 filtering is no longer possible nor desirable. An automatized procedure is necessary for 79 a preliminary analysis of the data in order to choose which ones to select and send to 80 Earth. Nevertheless, researchers are able to understand all type of information and in-81 terpret them critically by simultaneously using a combination of optimization, model learn-82 ing, planning, prediction, and diagnostic analysis. This is challenging for many automated 83 systems: as a result, artificial intelligence has become the perfect candidate for this task 84

-3-

thanks to the ability to recognize patterns and extract information from data by simulating human learning. Following this paradigm shift, the European Commission (EC)'s Horizon 2020 project started the Artificial Intelligence Data Analysis (AIDA) project, which not only aims to automatize the pre-processing of space data, but also to introduce modern data assimilation, statistical methods and machine learning (ML) to heliophysics data processing: *Aidapy*, an high level Python package for the analysis of spacecraft data from heliospheric missions has been developed as a result.

A new EC project has now followed up AIDA, the project Automatics in SpAce
 exPloration (ASAP) to study the deployment of ML tools onboard space missions, us ing the type of processors that can resist the hostile environment of space.

In this work we used Aidapy along with unsupervised ML clustering techniques to 95 characterize particle velocity distributions. The goal of the analysis is to differentiate be-96 tween simple and more complex regions within the velocity distribution functions mea-97 sured by MMS: in particular, complex shaped electron distributions, thus represented 98 with a greater number of clusters, have been shown to be good indicators for magnetic 99 reconnection and turbulence (Shuster et al., 2014; Hoshino et al., 2001). The necessity 100 of using unsupervised ML techniques arises from the fact that supervised ML needs huge 101 databases of input features and labeled outputs to work correctly: this kind of database 102 of distributions would be very problematic and time-consuming for researchers to build. 103 Nevertheless, unsupervised learning extracts patterns and information from untagged 104 data, thus being much more efficient and suitable for the task (Chollet, 2017). As the 105 literature demonstrates, ML gives good results when applied to data from simulations 106 (Dupuis et al., 2020): however, such data is less noisy and more smooth than real dis-107 tributions from MMS. A long pre-processing is therefore necessary, where it is critical 108 to deal with problems such as optimization and missing data in order to smooth the clus-109 tering. 110

In addition to the already mentioned Aidapy, the platform PySPEDAS (Grimes et al., 2022) and other packages for data analysis, ML, numerical, and visualization libraries in *Python* help to facilitate rapid development and deployment of ML algorithms. Thanks to its simple, user-friendly nature, Python has become the most popular language to build, train and test neural networks, and a fundamental tool for developing ML solutions with high iteration velocity.

-4-

¹¹⁷ 2 GMM Approach to Synthetic and Observational Data

We motivate the choice of unsupervised techniques in the classification of the ve-118 locity distribution functions. The main difference between supervised and unsupervised 119 ML is the usage of labelled train data in supervised learning. The model learns the re-120 lation between the labelled inputs and outputs and applies this knowledge to the unseen 121 data. The scientist in the loop would thus identify some features in the data so the al-122 gorithm can train on these chosen features. This however imposes some bias into the learn-123 ing algorithm because the features are based on the existing knowledge of the scientist. 124 Unsupervised methods on the other hand impose no such bias and are preferable to ex-125 tract features or clusters from the existing data. 126

¹²⁷ In addition to the imposed bias of supervised ML; we also prefer unsupervised due ¹²⁸ to the relatively small size of the downloaded VDF datasets.

129

2.0.1 The Gaussian Mixture Model

A Gaussian mixture is defined as function comprised of several Gaussians (usually the same amount as clusters in the data). Each of these Gaussians has the usual parameters of mean μ and covariance Σ , defining the centre and width of the Gaussian respectively (see figure 1). The final parameter is called the mixing probability π , and will de-



Figure 1: Three Gaussians depicted above their respective clusters with means and covariances shown.

133

fine how big the Gaussian function will be (Bouguila & Fan, 2020)(Moitra, 2018). The

Gaussian density function is given by:

$$\mathcal{N}(\mathbf{x}|\mu, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} exp(-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)).$$
(1)

Here the x represents the data points. The mixture model therefore states that all the
data points are generated from a mixture of such distributions:

$$P(x) = \sum_{i=1}^{k} P_i \mathcal{N}(\mathbf{x}|\mu, \Sigma), \qquad (2)$$

with k the amount of clusters and P_i is the weight of the respective Gaussian. The parameters of the distributions are fitted with the expectation-maximization method (EM), which starts of with randomly chosen values for the parameters and, after calculating the probability that the data points were generated by these Gaussians, changes them iteratively. This process maximizes the likelihood that the data were generated from the Gaussians with the chosen parameters.

The GMM is one of the primary models we will be interested in when studying the velocity distributions observed by the MMS mission. The amount and complexity of the clustering will be a good indicator of interesting regions in space where magnetic reconnection or other magnetic action might be present.

148

2.1 Clustering Statistics

We need a way to analyse the appropriate amount of clusters k that best describes 149 the dataset. One way to evaluate this would be to manually go over certain particle gen-150 erations and pick out the visible clusters. This way a training dataset of 'truth values' 151 can be created and used to train a neural network to find the optimal amount of clus-152 ters. Because this is a tedious and largely subjective task, we will prefer to work with 153 unsupervised training algorithms, as mentioned before, which do not require a predefined 154 cluster assignment. This does not mean that we will not be evaluating the performance 155 of the different algorithms visually at all, because this is of course still the best cluster-156 ing tool at our disposal (Pedregosa et al., 2011). 157

158

2.1.1 Bayesian Information Criterion

The Bayesian Information Criterion or BIC is in simple terms an information criterion that tries to balance the correctness of the fit of the model and the complexity of the model. If the complexity would not be taken into account, we run the risk of over¹⁶² fitting. On the other hand the correctness of the fit should be seen as a way to counter

¹⁶³ underfitting. The criterion is defined as follows:

$$BIC = kln(n) - 2ln(\hat{L}).$$
(3)

The two parts that are balanced are: k, the number of parameters estimated by the model and n, the number of data points that express the complexity of the model, while $ln(\hat{L})$ is the maximum log-likelihood that computes the goodness of the fit. This criterion is often used to evaluate the number of clusters in Gaussian mixture models and is also a build-in in the scikit-learn GMM-module.

An often used alternative is the Akaike Information Criterion or AIC. This is givenby:

$$AIC = 2k - 2ln(\hat{L}). \tag{4}$$

The difference is immediately clear: it does not involve a logarithmic term in the complexity measure (BIC penalizes the number of parameters in the model to a greater extent). While BIC assumes that the true model is in the candidate set and we simply want to find it, AIC only tries to find the model that most adequately describes the dataset (Konishi & Kitagawa, 2007).

176

2.1.2 Silhouette Coefficient

The Silhouette Coefficient uses the means intra-cluster distance (a) and the mean nearest-cluster distance (b) to calculate a measure of the goodness of the fit. This is a simple heuristic given by:

$$Silhouette - coefficient = \frac{b-a}{max(a,b)}.$$
(5)

As we can tell by the formula, the clustering is better for higher values of the coefficient. This means that points in a cluster lie closer to each other than to points from another cluster (b larger than a). As this heuristic can not assess the goodness of a fit of only one cluster (unlike BIC and AIC that can) we should be careful when applying it to decide whether there are two or one clusters present in the experimental data. It will however still give a good insight in the complexity of the distributions.

186 2.1.3 Calinsky-Harabasz Index

Another often used metric for evaluating clustering algorithms that uses the ratio of the sum of between-cluster dispersion and of within-cluster dispersion matrices (Pedregosa et al., 2011). This score is also easy to interpret: the higher the index, the more dense and well separated the clustering is, and thus the better. This score also fails to evaluate single cluster allocations like the silhouette coefficient.

¹⁹² 3 GMM Applied to Synthetic Distributions

```
193
```

3.1 Single Gaussian

Let us first asses the performance of the GMM (from the sklearn library in python) on synthetic distributions. To do this we look at the most trivial example, namely a number of particles generated from one single Gaussian distribution. To select the optimal number of clusters we will always use the BIC-score since it is much faster and more precise than the previously mentioned silhouette- and CH-score. This information criterion always selects a one-component model and therefore the GMM can always successfully reconstruct the mean speed and temperatures of the generated Gaussian distribution.

201

3.2 Mixture of Gaussians

The performance of the GMM on multiple Gaussian clusters is dependent on two 202 main factors; the variances of the different clusters and the separation between the clus-203 ters. When generating clusters from distributions with the same variance, it is of course 204 optimal to select the 'tied' covariance type inside the GMM. This means that all com-205 ponents share the same general covariance matrix. When dealing with a more general 206 case where particles can be drawn from Gaussian distributions with different variances, 207 the 'full' covariance type will be optimal. To illustrate this, figure 2 shows the BIC-scores 208 for 1-10 clusters using both covariance types on a dataset of 4 clusters generated from 209 Gaussian distributions with different variances. We can see that the 'full' covariance type 210 leads to the correct prediction of the number of components. 211



Figure 2: BIC-scores for two different covariance types for an instance where the clusters have different variances.

212	As expected, the GMM succeeds in identifying the optimal number of components
213	when dealing with all-Gaussian clusters. We can also show that the model is capable of
214	accurately predicting the means and variances of the clusters. The accuracy of these pre-
215	dictions do however go down when dealing with clusters that are not well separated. To
216	illustrate this, we will plot the normalized error of the predicted means and variances
217	for varying distance between the clusters. Lets look at an instance of 4 clusters.



Figure 3: Error on the mean and variance predictions of the GMM on 4 Gaussian clusters for decreasing distance between the respective clusters. The distance on the horizontal axis is the maximum distance between the three velocity components of the means of different clusters (in units of $10^6 m/s$). On the vertical axis, the error on the total velocity is displayed, averaged over the 4 clusters.

As expected, the Gaussian mixture model works excellent for recreating the means and variances of a number of Gaussian sampled clusters. It is only at the point where the clusters almost completely coincide that the accuracy of the predictions start to decline rapidly.

222

3.3 Single Kappa

The previous results look rather promising however, as we know the clusters observable in the electron VDF will not be completely Gaussian in nature. These clusters are characterised by non Maxwellian suprathermal tails. These tails decrease as a power law of the velocity (Pierrard & Lazar, 2010). As in most literature about space plasmas, we will fit these distributions by the Kappa distribution (Pierrard & Meyer-Vernet, 2017;

- Maksimovic et al., 1997; Kim et al., 2015). The spectral index of these distributions de-
- termines the slope of the distribution. In the limit where $\kappa \to \infty$, the distribution sim-
- ²³⁰ plifies to a Maxwellian (see fig. 4).



Figure 4: (Left) The Kappa distribution function for several values of the spectral index. (Right) Particles randomly generated from a Kappa distribution with $\kappa = 2$ (in units of $10^6 \ m/s$).

It is also important to note that the value of the spectral index must be chosen as such that it is not too close to the critical value $\kappa_c = 1.5$. At this value the distribution function collapses (Pierrard & Lazar, 2010). Following observations and satellite data, kappa distributions with a spectral index $2 < \kappa < 6$ seem to be a good fit (Shohaib et al., 2022) and thus satisfies this requirement.

We will now see if the GMM is still capable of identifying the correct number of 236 clusters and their means/variances if these clusters are not sampled from a Gaussian dis-237 tribution, but from a kappa distribution. Let us start with one single kappa distributed 238 cluster. In the previous section, when generating Gaussian clusters, we selected $\kappa = 200$. 239 We now bring down this value until the GMM does no longer makes the right predic-240 tion. Below $\kappa = 6$, the GMM predicts a significantly larger amount of clusters due to 241 the outliers generated from the kappa distribution. Since we are interested in kappa val-242 ues below this, we need to improve the performance. This can be done by increasing the 243 convergence threshold of the EM-algorithm that the GMM uses. This way the algorithm 244 goes to fewer iterations and is less likely to overestimate the number of actual clusters. 245

Using this method, we get down to $\kappa = 2.6$ before the model starts overfitting the number of clusters. This process also inherently speeds up the computation since we need fewer iterations. Remarkably, as we will see in the next section, this overfitting does not occur when more than one cluster is present and we can get our spectral index as low as $\kappa = 2$. This way we can correctly predict all distributions that model the velocity in space plasmas.

We conclude that the GMM can correctly reconstruct the mean speeds and temperatures of a single kappa distributed cluster for $\kappa = 2.6$ and up.

²⁵⁴ 3.4 Mixture of Kappa's



Figure 5: Predicted clusters by BIC-score after GMM. Means of clusters are highlighted in red (graph in units of $10^6 m/s$).

In complete analogy with the mixture of Gaussian clusters, we look at the GMM 255 performance when the dataset would consist of a mix of different kappa distributed clus-256 ters. We do however now have an extra parameter to take into account, namely: Do all 257 clusters have the same kappa value or not? Contrary to the single-cluster case, the GMM 258 model is able to correctly predict the number of components in a multi-cluster dataset, 259 even when they are generated by distributions with a spectral index between 2 and 2.6. 260 The primary deciding factor of the precision of these predictions will thus once again be 261 the separation between the different clusters. We will analyse the results for an instance 262

where all clusters are generated by a full $\kappa = 2$ distribution, and another instance where the spectral indices can vary randomly between $\kappa = 2$ and $\kappa = 6$.



Figure 6: Error on the mean and variance predictions of the GMM on 4 Kappadistributed clusters (all same κ) for decreasing distance between the respective clusters. The distance on the horizontal axis is the maximum distance between the three velocity components of the means of different clusters (in units of $10^6 m/s$). On the vertical axis, the error on the total velocity is displayed, averaged over the 4 clusters.



Figure 7: Error on the mean and variance predictions of the GMM on 4 Kappadistributed clusters (2 < κ < 6) for decreasing distance between the respective clusters. The distance on the horizontal axis is the maximum distance between the three velocity components of the means of different clusters (in units of 10⁶ m/s). On the vertical axis, the error on the total velocity is displayed, averaged over the 4 clusters.

From these results we conclude that the performance of the GMM is not really af-265 fected by clusters with a non-equal spectral index. In previous examples we always ran 266 the GMM for 1-10 components. If we pick a higher number for the maximum amount 267 of components, we risk overfitting some simpler cases with few components by assign-268 ing a really high number of clusters. This also reduces the run time of the program when 269 we will run it for more than 100 electron distributions. The model is still capable of rep-270 resenting the complexity of the distribution, which is in essence what we are after. This 271 means that instances with over 10 components will almost always be assigned 10 com-272 ponents for the optimal solution, and not some smaller number. 273

3.5 Spectral Clustering

Another clustering technique that can be considered in the context of electron VDFs 275 is spectral clustering. This method performs a dimensionality reduction based on the spec-276 trum of the similarity matrix (Bonaccorso, 2017) (Pedregosa et al., 2011). When spec-277 tral clustering is applied to a Gaussian/Kappa distributed dataset like before, we see that 278 the algorithm is indeed able to correctly classify most clusters using the Silhouette Co-279 efficient and Calinsky-Harabasz Index as clustering statistics (BIC and AIC are not avail-280 able for this clustering algorithm). However, using this algorithm comes with a big jump 281 in time complexity as well as the added run time that results from the change in clus-282 tering statistics. Nevertheless, spectral clustering could become relevant when using the 283 'SpectralClustering' function of *scikit-learn*. This function allows the user to use paral-284 lelization, which will split the work across different CPU cores and decrease the execu-285 tion time (Brownlee, 2020). Several efforts were made to to get the time complexity down 286 this way to make it comparable to GMM, but this does not yet yield the desired results. 287 These attempts included using the cores of our local machines as well as utilizing sev-288 eral cores of the tier-2 Genius cluster of the VSC (Flemish Supercomputer Center). As 289 for now, the achieved time complexity is not yet one to rival our GMM results. 290

²⁹¹ 4 Data and Methods

For the purpose of testing the capability of the GMM to define the complexity of 292 a measured dataset, we utilized data from NASA's MMS Mission. Fast Plasma Inves-293 tigation (FPI) instrument measures incoming particles through a filter which selects cer-294 tain particle speeds and directions; then a 3D picture of the ion plasma is produced ev-295 ery 150 milliseconds, while for electron plasma FPI captures a picture every 30 millisec-296 onds (Pollock et al., 2016). Aidapy, an high level Python package for the analysis of space-297 craft data from heliospheric missions developed by ESA, is used to download data from 298 FPI and other on-board instruments. 299

Unlike simulation data, when working with real data one has to deal with specifics of instruments. In the case of FPI, low energy particle counting is perturbed by the electrical charge of the probe: this leaves a gap in the center of measured VDFs. Before processing the data with GMM, the Python package Scikit-learn is used to perform linear interpolation to fill the gap in the VDFs. Particles are then generated from VDFs for

-15-

input into the GMM: to avoid noisy distributions with few particles and ones which are
too demanding in terms of numerical resources, a number of particles of 40,000 is chosen by authors' experience. The information criterion chosen to select the number of optimal beams of the mixture is BIC, as it is to be found preferable to AIC in working with
a large number of real data.

Data were selected from two distinct intervals in which magnetic reconnection signatures were found. In the literature, different types of reconnection signatures have been identified (M. Goldman et al., 2016). Several of these signatures have been observed in the events analyzed in this work. The first event is from December 8, 2015, when the reconnecting dayside magnetopause was crossed by MMS probes. The second one is from July 3, 2017 and it was observed in the magnetotail.

³¹⁶ 5 GMM applied to MMS data in the dayside magnetopause

During the event on December 8, MMS spacecraft was at first in the magnetosheath, 317 but the magnetopause moved outward causing the spacecraft to move to the magnetopause 318 (Burch & Phan, 2016). The crossing of a reconnecting magnetopause is recognized at 319 11:20:42-11:20:45 UT, because of the behaviour of several physical quantities: intense 320 current density, strong electric field and high speed electron outflows. Four seconds of 321 data observations from MMS3 were processed with the GMM technique, from 11:20:41 322 to 11:20:45. Every 0.03 s an electron VDF is measured by the spacecraft, for a total of 323 133 VDFs analyzed within the interval. For each VDF the BIC information criterion as-324 signes a score based on how accurate the fit is. The fit is then repeated with a different 325 number of clusters. The maximum number of possible clusters is an input for the GMM: 326 from the authors' experience, 10 is a reasonable compromise between result accuracy and 327 computation time. 328

The bottom plot of figure 8 shows the normalized results of the GMM analysis on the VDFs. High values of the score (white lines) shows a difficulty in the fit of the distribution function, which may be associated with a complex VDF. In fact, the information criterion tends to assign smaller scores to the best fits, which correspond to more Maxwellian VDFs. As the plot shows, during the reconnecting magnetopause crossing (between 11:20:43 and 11:20:44), the VDFs tend to become more complex. For a visual verification, three VDFs taken at different times are shown in figures 9, 10 and 11. The

-16-



Figure 8: Data from 8 December 2015. Vectors are expressed in Geocentric Solar Ecliptic System (GSE), with X-axis pointing from the Earth towards the suns, its Y-axis chosen to be in the ecliptic plane pointing towards dusk and Z-axis parallel to the ecliptic pole. GSE components of magnetic and electric field are shown. Electric and magnetic field are shown in panels (A) and (B), along with electron bulk velocity (C) and density (D). Panel (E) shows the GMM results. Vertical dashed lines represent the times when VDFs are measured.

first one is taken at 11:20:42:08, when the crossing of the reconnecting region had not yet happened and the BIC score is near its minimum. The second and third one, which appear visually complex and show strong asintropy, are taken after the reconnection event, where the BIC score indicates a worst fit. This capability of the GMM to automatically recognize regions where VDFs present non-Maxwellian features is of crucial importance for the detection of reconnection regions within the plasma.

Figure 12 illustrates the Gaussians parameters found by the GMM applied to one of the distributions measured before the reconnection. The number of components provided as input to the algorithm is three. The ellipses evidence the mean and the vari-

- ³⁴⁵ ance of each Gaussian of the mixture, and the weight represents the probability that a
- random particle belongs to one of the three components.



Figure 9: Three cuts of the VDF taken at 11:20:42:08, before the reconnection event. $\mathbf{V}_{\mathbf{B}}$ is parallel to the magnetic field, $\mathbf{V}_{\mathbf{V}}$ is in the direction of the bulk velocity, $\mathbf{V}_{\mathbf{B}\times\mathbf{V}}$ is in the direction of $\mathbf{B} \times \mathbf{V}$ and $\mathbf{V}_{\mathbf{V}\mathbf{perpB}}$ is the bulk velocity projected onto the plane normal to \mathbf{B} .



Figure 10: Three cuts of the VDF taken at 11:20:42:08, before the reconnection event. $\mathbf{V}_{\mathbf{B}}$ is parallel to the magnetic field, $\mathbf{V}_{\mathbf{V}}$ is in the direction of the bulk velocity, $\mathbf{V}_{\mathbf{B}\times\mathbf{V}}$ is in the direction of $\mathbf{B} \times \mathbf{V}$ and $\mathbf{V}_{\mathbf{V}\mathbf{p}\mathbf{erp}\mathbf{B}}$ is the bulk velocity projected onto the plane normal to \mathbf{B} .



Figure 11: Three cuts of the VDF taken at 11:20:42:08, before the reconnection event. $\mathbf{V}_{\mathbf{B}}$ is parallel to the magnetic field, $\mathbf{V}_{\mathbf{V}}$ is in the direction of the bulk velocity, $\mathbf{V}_{\mathbf{B}\times\mathbf{V}}$ is in the direction of $\mathbf{B} \times \mathbf{V}$ and $\mathbf{V}_{\mathbf{V}\mathbf{p}\mathbf{e}\mathbf{r}\mathbf{p}\mathbf{B}}$ is the bulk velocity projected onto the plane normal to \mathbf{B} .



Figure 12: VDF cut from the dayside magnetopause taken at 11:20:44.18, before the reconnection. $\mathbf{V}_{\mathbf{B}}$ is parallel to the magnetic field and $\mathbf{V}_{\mathbf{B}\mathbf{x}\mathbf{V}}$ is in the direction of $\mathbf{B} \times \mathbf{V}$. The black ellipses show the different Gaussians of the mixtures based on their mean and variance. The transparency of the ellipses is determined by the weight of each Gaussian.

³⁴⁷ 6 GMM applied to MMS data in the magnetotail

During the event of 3 July 2017 MMS3 spacecraft was in the magnetotail, and observed another reconnection event (Burch et al., 2019). Again, four seconds of data were analyzed, from from 05:26:48 to 05:26:52. Reconnection X-line is observed near 05:25:50:72.



Figure 13: Data from 3 July 2017. Vectors are expressed in Geocentric Solar Ecliptic System (GSE), with X-axis pointing from the Earth towards the suns, its Y-axis chosen to be in the ecliptic plane pointing towards dusk and Z-axis parallel to the ecliptic pole. GSE components of magnetic and electric field are shown. Electric and magnetic field are shown in panels (A) and (B), along with electron bulk velocity (C) and density (D). Panel (E) shows the GMM results. Vertical dashed lines represent the times when VDFs are measured.

As shown in figure 13, until 05:25:49:50 VDFs appear Maxwellian. More complex VDFs are recognized for nearly a second, until they return simpler. The simplification of the VDFs occurs after the time when reconnection is observed, near 05:25:50:72. The GMM again succeeds in detecting the reconnection region through the clustering result. Three VDFs are shown in figures 14, 15 and 16 for visual verification. As expected from GMM results, the two VDFs with smaller BIC scores show Maxwellian features, while

- the VDF from 5:26:50:15, before the reconnection event, appears strongly complex with
- $_{358}$ large peaks in the V_B direction. The algorithm was again able to automatically recog-
- nize the most complex distributions within the time interval with great accuracy.



Figure 14: Three cuts of the VDF taken at 11:20:42:08, before the reconnection event. $\mathbf{V}_{\mathbf{B}}$ is parallel to the magnetic field, $\mathbf{V}_{\mathbf{V}}$ is in the direction of the bulk velocity, $\mathbf{V}_{\mathbf{B}\times\mathbf{V}}$ is in the direction of $\mathbf{B} \times \mathbf{V}$ and $\mathbf{V}_{\mathbf{V}\mathbf{p}\mathbf{erp}\mathbf{B}}$ is the bulk velocity projected onto the plane normal to \mathbf{B} .



Figure 15: Three cuts of the VDF taken at 11:20:42:08, before the reconnection event. $\mathbf{V}_{\mathbf{B}}$ is parallel to the magnetic field, $\mathbf{V}_{\mathbf{V}}$ is in the direction of the bulk velocity, $\mathbf{V}_{\mathbf{B}\times\mathbf{V}}$ is in the direction of $\mathbf{B} \times \mathbf{V}$ and $\mathbf{V}_{\mathbf{V}\mathbf{p}\mathbf{erp}\mathbf{B}}$ is the bulk velocity projected onto the plane normal to \mathbf{B} .



Figure 16: Three cuts of the VDF taken at 11:20:42:08, before the reconnection event. $\mathbf{V}_{\mathbf{B}}$ is parallel to the magnetic field, $\mathbf{V}_{\mathbf{V}}$ is in the direction of the bulk velocity, $\mathbf{V}_{\mathbf{B}\times\mathbf{V}}$ is in the direction of $\mathbf{B} \times \mathbf{V}$ and $\mathbf{V}_{\mathbf{V}\mathbf{p}\mathbf{erp}\mathbf{B}}$ is the bulk velocity projected onto the plane normal to \mathbf{B} .

360 7 Conclusions

In this paper we have investigated the effectiveness of the Gaussian mixture model in 1) recognizing the complexity of the distributions (i.e. the number of components) and 2) recreating the means (mean velocity) and variances (temperatures) of the clusters in electron VDF's.

The tests that were run on synthetic distributions proved the ability of the algo-365 rithm to accurately predict the complexity of kappa distributed clusters (with the max-366 imum number of components set tot 10 as to prevent overfitting) as well as reconstruct-367 ing the means and variances of these clusters. From the numerical results, it is obvious 368 that the GMM algorithm still performs the better on the Gaussian distributed clusters 369 as expected. Nevertheless, for clusters whose means are well separated (i.e. more than 370 10^7 m/s apart), the results for the kappa distributed particles are comparable to the Gaus-371 sians. We also remark that it is not important if the clusters have the same spectral in-372 dex ($\kappa = 2$), or if the spectral index is different for each cluster ($2 < \kappa < 6$). 373

With regard to the analysis of real data, we applied GMM to MMS data in order to automatically identify magnetic reconnection sites through the complexity of the velocity distribution functions.

For the analysis, we selected time intervals from articles where magnetic reconnec-377 tion events were identified from observations of particular variations in electric and mag-378 netic fields, current density, and particles behavior. After preprocessing the data, which 379 includes filling in the gap in the VDFs at low energies, we analyzed the particles gen-380 erated by distribution with the GMM. We utilized the Bayesian Information Criterion 381 to choose the best fitting amount of clusters within the distributions. The model has shown 382 that it is able to capture the variation in complexity of the functions, arriving through 383 this to automatically locate reconnection sites with good accuracy. In addition to this, 384 a visual test showed that the BIC scores can accurately indicate the most complex dis-385 tributions that show strong non-Maxwellian features. 386

In recent years the task of looking at the raw data from the spacecrafts and selecting the interesting ones was done by eye by scientists. Due to limitations of the probes, a continuous overwriting of data takes place and a large part of it is lost. Future goals include further improving the unsupervised ML techniques so that them can be used to analyze the data and collect the most interesting ones without having to lose a lot of im-

-28-

- ³⁹² portant information. Just as with synthetic data, we favor GMM with Bayesian infor-
- ³⁹³ mation criterion thanks to its efficiency and accuracy.

394 Acknowledgments

This project has received funding from the *ERC Advanced Grant TerraVirtuale* of GL (grant agreement No. 1101095310) and from the European Union's Horizon Europe research and innovation programme under grant agreement No. 101082633 (ASAP). Views and opinions expressed are however those of the authors only and do not necessarily reflect those of the European Union or the European Commission. Neither the European Union nor the European Commission can be held responsible for them.

Further funding was provided by the KULeuven Bijzonder OnderzoeksFonds (BOF) under the C1 project TRACESpace and from the NASA HSR Grant 80NSSC21K1689. FV acknowledges the support of the PRIN 2022 project "The ULtimate fate of TuRbulence from space to laboratory plAsmas (ULTRA)" (2022KL38BK), funded by the Italian Ministry of University and Research. This paper and related research have been conducted during and with the support of the Italian national inter-university PhD programme in Space Science and Technology.

408

Computing has been provided by the Flemish Supercomputing Center (VSC).

409 Data Availability Statement

The source code for data preprocessing and GMM analysis is available at https:// github.com/NathanNoaMaes/VDFClusteringMMS. MMS data are available at https:// lasp.colorado.edu/mms/sdc/public/data/. Data analysis was performed using Aidapy available at https://gitlab.com/aidaspace/aidapy, PySPEDAS available at https:// github.com/spedas/pyspedas and scikit-learn available at https://github.com/scikit -learn/scikit-learn.

416 **References**

417 Baker, D. N., Riesberg, L., Pankratz, C. K., Panneton, R. S., Giles, B. L., &

418	Wilder, R. E., F. D.and Ergun.	(2016, mar).	Magnetosph	eric multiscale
419	instrument suite operations and d	ata system.	Space Science	Reviews, 199.
420	Retrieved from https://doi.org,	/10.1007/s1121	14-014-0128-5	doi:
421	10.1007/s11214-014-0128-5			

Biskamp, D. (2000). Magnetic reconnection in plasmas. Cambridge University Press.
doi: 10.1017/CBO9780511599958

424	Bonaccorso, G. (2017). Machine learning algorithms. Packt Publishing Ltd, 2017.
425	Bouguila, N., & Fan, W. (2020). Mixture models and applications. doi: 10.1007/978
426	-3-030-23876-6
427	Brownlee, J. (2020, 09). Multi-core machine learning in python with scikit-learn.
428	https://machinelearningmastery.com/multi-core-machine-learning-in
429	-python/.
430	Burch, J. L., Dokgo, K., Hwang, K. J., Torbert, R. B., Graham, D. B., & Webster,
431	e. a., J. M. (2019). High-frequency wave generation in magnetotail reconnec-
432	tion: Linear dispersion analysis. Geophysical Research Letters, 46, 4089–4097.
433	doi: https://doi.org/10.1029/2019GL082471
434	Burch, J. L., Moore, T. E., Torbert, R. B., & Giles, B. L. (2016). Magnetospheric
435	multiscale overview and science objectives. Space Science Reviews, 199. doi:
436	https://doi.org/10.1007/s11214-015-0164-9
437	Burch, J. L., & Phan, T. D. (2016). Magnetic reconnection at the dayside mag-
438	netopause: Advances with mms. Geophysical Research Letters, 43, 8327–8338.
439	doi: https://doi.org/10.1002/2016GL069787
440	Chollet, F. (2017). Deep learning with python. Manning.
441	Dupuis, R., Goldman, M. V., Newman, D. L., Amaya, J., & Lapenta, G. (2020,
442	jan). Characterizing magnetic reconnection regions using gaussian mixture
443	models on particle velocity distributions. The Astrophysical Journal, $889(1)$,
444	22. Retrieved from https://doi.org/10.3847%2F1538-4357%2Fab5524 doi:
445	10.3847/1538-4357/ab5524
446	Goldman, M., Newman, D., & Lapenta, G. (2016). What can we learn about magne-
447	totail reconnection from 2d pic harris-sheet simulations? Space Sci Rev, 199,
448	651–688. doi: https://doi.org/10.1007/s11214-015-0154-y
449	Goldman, M. V., Newman, D. L., Eastwood, J. P., & Lapenta, G. (2020). Multi-
450	beam energy moments of multibeam particle velocity distributions. Journal of
451	Geophysical Research: Space Physics, 125(12), e2020JA028340. Retrieved
452	<pre>from https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/</pre>
453	2020JA028340 (e2020JA028340 2020JA028340) doi: https://doi.org/10.1029/
454	2020JA028340
455	Grimes, E. W., Harter, B., Hatzigeorgiu, N., Drozdov, A., Lewis, J. W., Angelopou-
456	los, V., Le Contel, O. (2022). The space physics environment data

-31-

457	analysis system in python. Frontiers in Astronomy and Space Sciences,
458	9. Retrieved from https://www.frontiersin.org/articles/10.3389/
459	fspas.2022.1020815 doi: 10.3389/fspas.2022.1020815
460	Hoshino, M., Hiraide, K., & Mukai, T. (2001, 06). Strong electron heating and non-
461	maxwellian behavior in magnetic reconnection. Earth Planets Space, 53, 627-
462	634. doi: 10.1186/BF03353282
463	Kim, S., Yoon, P., Choe, G., & Wang, L. (2015, 06). Asymptotic theory of solar
464	wind electrons. The Astrophysical Journal, 806. doi: 10.1088/0004-637X/806/
465	1/32
466	Konishi, S., & Kitagawa, G. (2007). Information criteria and statistical modeling.
467	doi: 10.1007/978-0-387-71887-3
468	Lapenta, G., Goldman, M., Newman, D., & Markidis, S. (2016, May). Where should
469	mms look for electron diffusion regions? Journal of Physics: Conference Se-
470	ries, 719, 012011. Retrieved from http://dx.doi.org/10.1088/1742-6596/
471	719/1/012011 doi: 10.1088/1742-6596/719/1/012011
472	Li, T. C., Liu, YH., & Qi, Y. (2021, mar). Identification of active magnetic re-
473	connection using magnetic flux transport in plasma turbulence. The Astrophys-
474	ical Journal Letters, 909(2), L28. Retrieved from https://doi.org/10.3847%
475	2F2041-8213%2Fabea0b doi: 10.3847/2041-8213/abea0b
476	Maksimovic, M., Pierrard, V., & Riley, P. (1997). Ulysses electron distributions
477	fitted with Kappa functions. Geophysical Research Letters, 24, 1151-1154.
478	Retrieved from https://hal.archives-ouvertes.fr/hal-03801373 doi:
479	10.1029/97GL00992
480	Moitra, A. (2018). Algorithmic aspects of machine learning. Cambridge University
481	Press. doi: 10.1017/9781316882177
482	Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O.,
483	Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. Journal of
484	Machine Learning Research, 12, 2825–2830.
485	Pierrard, V., & Lazar, M. (2010, 03). Kappa distributions: Theory and applications
486	in space plasmas. Solar Physics, 267. doi: 10.1007/s11207-010-9640-2
487	Pierrard, V., & Meyer-Vernet, N. (2017). Chapter 11 - electron distributions in
488	space plasmas. In G. Livadiotis (Ed.), Kappa distributions (p. 465-479). El-
489	sevier. Retrieved from https://www.sciencedirect.com/science/article/

490	pii/B9780128046388000115 doi: https://doi.org/10.1016/B978-0-12-804638-8
491	.00011-5
492	Pollock, C., Moore, T., Jacques, A., Burch, J., Gliese, U., Saito, Y., Zeuch, M.
493	(2016). Fast Plasma Investigation for Magnetospheric Multiscale. Space
494	Science Reviews.
495	Retinò, A., Sundkvist, D., Vaivads, A., Mozer, F., André, M., & Owen, C. J. (2007,
496	April). In situ evidence of magnetic reconnection in turbulent plasma. $Nature$
497	<i>Physics</i> , 3(4), 236-238. doi: 10.1038/nphys574
498	Scott, B. (2021). Turbulence and instabilities in magnetised plasmas, volume 1. IOP
499	Publishing. Retrieved from https://dx.doi.org/10.1088/978-0-7503-2504
500	-2 doi: 10.1088/978-0-7503-2504-2
501	Shohaib, M., Masood, W., Siddiq, M., Alyousef, H., & El-Tantawy, S. (2022, 04).
502	Formation of electrostatic solitary and periodic waves in dusty plasmas in the
503	light of voyager 1 and 2 spacecraft and freja satellite observations. Journal
504	of Low Frequency Noise, Vibration and Active Control, 41, 146134842210913.
505	doi: 10.1177/14613484221091340
506	Shuster, J. R., Chen, LJ., Daughton, W. S., Lee, L. C., Lee, K. H., Bessho, N.,
507	Argall, M. R. (2014). Highly structured electron anisotropy in collisionless
508	reconnection exhausts. Geophysical Research Letters, $41(15)$, 5389-5395. doi:
509	https://doi.org/10.1002/2014GL060608
510	Yokoi, N., & Hoshino, M. (2011, October). Flow-turbulence interaction in magnetic
511	reconnection. Physics of Plasmas, 18(11). Retrieved from http://dx.doi
512	.org/10.1063/1.3641968 doi: 10.1063/1.3641968