

Abstract

Nowcasting and forecasting of the radiation environment in the Earth’s lower atmosphere are critical for the safety of aircraft and spacecraft crews and passengers. Currently, this problem is addressed by employing physics-based particle transport and precipitation models. However, given the increased number of radiation measurements available to the community, it is possible to start developing data-driven approaches. We prepared Machine Learning-ready (ML-ready) datasets to nowcast the effective dose rates at aviation altitudes. The presented datasets contain 92,476 individual measurements from 589 flights obtained by the Automated Radiation Measurements for Aerospace Safety (ARMAS) experiment from 2013 to 2023. The ARMAS measurements are augmented with the properties of the Geospace environment, such as solar soft X-ray and proton fluxes, solar wind properties, secondary cosmic ray neutrons, space weather indexes, and global solar activity indicators (such as daily sunspot number). ARMAS data are separated into three partitions, ensuring that (1) the data points from a single flight remain within the same partition, and (2) each partition samples the flight locations and Geospace environment conditions equally. Several versions of the datasets allow predictions based on point-in-time measurements and use up to 24 hours of Geospace parameter history. The test of the use case demonstrates a possibility of nowcasting ARMAS measurements with accuracies slightly better than the considered physics-based models. The publicly available ML-ready datasets could serve as the first step in data preparation for ML-driven nowcasting and forecasting of the radiation environment.

Plain Language Summary

Understanding the radiation levels in the Earth’s atmosphere and predicting them at the locations along flight routes is important for the aviation industry. Machine Learning (ML) techniques are often used nowadays for prediction tasks. However, an extensive data preparation phase is required before applying an ML algorithm for predictions. Given the advantage of increasing volumes of the radiation measurement data, we present ML-ready datasets that allow users to bypass the data preparation stage and go directly to the phase of testing ML models for aviation radiation prediction. The ML-ready dataset is publicly available via the Radiation Data Portal¹.

1 Introduction

Understanding the radiation environment at the aviation altitudes (approximately 8 – 17 km) remains an important and not fully explored topic for the aviation community. Integrating the effective dose rate along the flight trajectory constitutes the total effective dose absorbed by the aircraft crew and passengers. Given that there are recommended limits for the annual effective dose (20 mSv/yr and 1 mSv/yr for 5 years for the radiation occupation workers and for the general public, correspondingly; Cho et al., 2017), accurate monitoring and prediction of the radiation environment becomes essential.

Galactic cosmic rays cascading in Earth’s atmosphere are the primary contributors to effective dose rates. In addition, the atmospheric radiation dose rates can increase during strong solar energetic particle (SEP) events (Reames, 2021). For example, the estimates presented for the September 10-11, 2017, SEP event demonstrated that the location-averaged effective dose rate due to the SEP at a height of 12 km was approximately $3 \mu Sv/h$ (Kataoka et al., 2018). While this is less than the contribution due to galactic cosmic rays (more than twice that amount), it becomes comparable by an order of magnitude.

¹ <https://dmlab.cs.gsu.edu/rdp/ml-dataset.html>

66 Routine measurements of the radiation environment have become possible in the
67 last decade via efforts from the Automated Radiation Measurements for Aerospace Safety
68 experiment (ARMAS; Tobiska et al., 2015, 2016, 2018). The ARMAS device measures
69 the local radiation environment conditions in real time during commercial aircraft flights.
70 Currently, ARMAS has flown on more than 1000 flights and provided $\sim 400,000$ individ-
71 ual measurements, sampling the aviation altitudes across the United States and world-
72 wide². The development and release of the Radiation Data Portal (RDP; Sadykov et al.,
73 2021) provided a convenient overview and access to the ARMAS data.

74 Nowcasting or predicting the radiation environment in Earth’s atmosphere can be
75 carried out by physics-based approaches such as CARI (i.e., the Civil Aviation Research
76 Institute; Copeland, 2021), Professional Aviation Dose Calculator (PANDOCA; Matthiä
77 et al., 2014), the Nowcast of Aerospace Ionizing RADIation System (NAIRAS; Mertens
78 et al., 2013), etc. The NAIRAS model takes into account the galactic cosmic ray and SEP
79 inputs. It is closely integrated with the ARMAS experiment by providing the assessments
80 of the radiation dose rates for every ARMAS measurement. Recently, version 3 of the
81 NAIRAS model was released (Mertens et al., 2023; Mertens et al., 2024) that integrates
82 the barometric altitude corrections for a better comparison to ARMAS measurements
83 at the appropriate altitudes, as well as run-on-request capabilities at the Community Co-
84 ordinated Modeling Center (CCMC) by the National Aeronautics and Space Adminis-
85 tration (NASA). The available amount of the ARMAS measurements makes it possible
86 to explore data-driven, statistical, and machine learning (ML) based analysis for predic-
87 tion conditions of the radiation environment. A comparison of data-driven and physics-
88 based predictions is valuable for further improvements to the physics considered in the
89 models.

90 Data preparation is the first, yet most demanding step when developing an ML model.
91 Typically, datasets that can be easily fed into ML models with relatively low or no prepa-
92 ration required are called ML-ready datasets. Several works (Masson et al., 2024; Nita
93 et al., 2022) have introduced the key principles required for these datasets, often related
94 to the overall completeness, preparation levels, and accessibility to a wider community.
95 The preparation of the ML-ready datasets suitable for radiation environment prediction
96 could enhance community efforts in modeling the radiation environment and stimulate
97 the development of data-driven approaches.

98 The primary goal of this paper is to present an ML-ready dataset that could be
99 used for nowcasting and forecasting of the radiation environment at aviation altitudes.
100 The paper is structured as follows: Section 2 describes the data mining and preprocess-
101 ing steps required for both ARMAS and the space environment measurements. Section 3
102 presents steps for shaping the data into an ML-ready format. This includes the associ-
103 ation of ARMAS measurements with the Geospace environment properties (such as ground-
104 based neutron monitor counts, soft X-ray and proton flux measurements at the geosta-
105 tionary orbit, solar wind properties at L1 point, geomagnetic and global solar activity
106 indexes) and the partitioning of multi-dimensional sparse data into three subsets that
107 should be used directly for training and validation of the model. Section 4 provides a
108 case example of how the dataset could be used for nowcasting, followed by a summary
109 and discussion in Section 5.

110 2 Data Preparation

111 2.1 ARMAS data processing

112 The Automated Radiation Measurements for Aerospace Safety (ARMAS) exper-
113 iment provides the richest publicly available dataset of atmospheric radiation measure-

² <https://dmlab.cs.gsu.edu/rdp/>

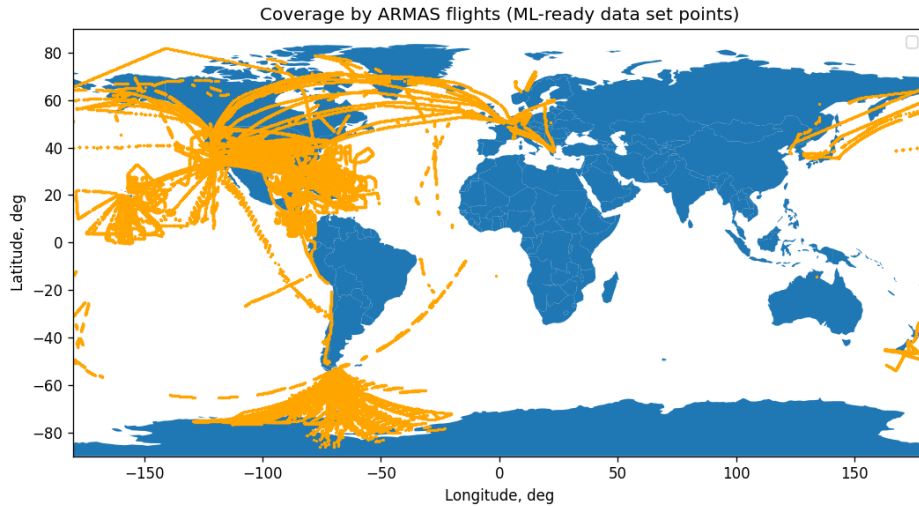


Figure 1. Coverage of considered ARMAS flight measurements over the Earth globe. The individual measurements are marked as orange points on the map.

114 ments along routes at aviation altitudes. In this work, we consider ARMAS data obtained
 115 from 1142 flight files (589 of which passed the selection criteria and contributed to an
 116 ML-ready dataset) from June 2013 to December 2023. This period of time covers the
 117 peak and decay of solar cycle 24, the solar minimum, and the rise phase of solar cycle
 118 25, which enables sampling different levels of global solar activity. The flights cover a
 119 wide geographic span, including the extensive sampling of the continental US and ter-
 120 ritories, the Pacific Ocean regions, the North Atlantic, and Antarctica (Figure 1). AR-
 121 MAS measurements are stored in the individual files for each flight and device and are
 122 accessible either via the Space Environment Technologies (SET) public data archive³ or
 123 the Radiation Data Portal⁴. More details on ARMAS measurements can be found in Tobiska
 124 et al. (2015, 2016, 2018).

125 For each flight, we apply the following pre-processing steps. We sample the AR-
 126 MAS measurements only in the barometric altitude range of 8.0 km – 15.5 km, as AR-
 127 MAS often measures zero radiation below these altitudes. The corresponding maximum
 128 GPS altitude in the filtered dataset is 16.272 km. We notice that several flights had two
 129 or more ARMAS devices onboard, and the measurements of these devices were reported
 130 in separate files. For groups of files corresponding to the same flight, we consider only
 131 one file where ARMAS measurements are most strongly correlated with the NAIRAS
 132 v3 model estimates (in terms of Pearson correlation coefficient). We also remove the data
 133 points marked as non-science-use or during the periods of electromagnetic interference.
 134 The remaining data points are inspected manually for every flight. The procedure re-
 135 sults in 589 unique flights propagating into the final dataset, and 92,476 unique data points.

136 The ARMAS files prepared by the SET public data archive⁵ contain several ad-
 137 ditional parameters that we propagate into the final dataset. For example, there is in-
 138 formation about the geomagnetic cutoff rigidities at the point of the flight measurement,
 139 as well as the barometric and GPS altitudes. The files also contain the physics-based now-
 140 cast of the radiation dose rate obtained with 2nd and 3rd versions of the Nowcast of At-

³ <https://sol.spacenvironment.net/ARMAS/Archive/>

⁴ <https://data.nas.nasa.gov/helio/portals/rdp/>

⁵ https://sol.spacenvironment.net/ARMAS_Archive/

141 atmospheric Ionizing Radiation for Aviation Safety model (NAIRAS; Mertens et al., 2013,
142 2023; Mertens et al., 2024). Version 3 (v3) of the model provides a better agreement with
143 the ARMAS measurements and is therefore suggested to be used as a physics baseline
144 for the radiation forecasts.

145 **2.2 Neutron monitor stations**

146 Along with the muons, secondary cosmic ray neutrons are an abundant remnant
147 of the atmospheric particle cascades formed by the primary cosmic rays and, therefore,
148 represent one of the most reliable ways to monitor the Geospace cosmic ray environment.
149 The Neutron Monitor Database (NMDB, Kozlov et al., 2003; Väisänen et al., 2021) ag-
150 glomerates the measurements from more than 50 neutron monitors around the globe. The
151 NMDB Event Search Tool⁶ (NEST) provides access to the data from these stations, in-
152 cluding the count rates corrected with respect to the effects of the local atmospheric pres-
153 sure.

154 The NAIRAS physics-based model utilizes neutron monitor data from four loca-
155 tions, Oulu (OULU, $R_c=0.81$ GeV), Lomnický (LKMS, $R_c=3.84$ GeV), Thule (THUL, $R_c=0.30$ GeV),
156 and Izmiran (MOSC, $R_c = 2.43$ GeV), as inputs. However, according to the NEST inter-
157 face, LKMS and MOSC stations have some interruptions for the period of interest (2013-
158 2023). The continuity of data is critical for the construction of the ML-ready dataset.
159 Therefore, we decided to replace these stations with the neutron monitor station in Newark
160 (NEWK, $R_c=2.40$ GeV), and we also added the South Pole station (SOP0) with extremely
161 low geomagnetic cut-off rigidity ($R_c=0.10$ GeV). The pressure- and efficiency-corrected
162 counts are obtained for all stations with a cadence of 1 minute.

163 **2.3 Solar wind parameters**

164 According to Tobiska et al. (2018), ARMAS measurements tend to have higher ra-
165 diation dose rates than the predictions by the NAIRAS physics-based models, specifi-
166 cally in the regions of low geomagnetic cutoff rigidity. The presumed reason for this is
167 the leakage of the energetic charged particles from the radiation belts, for which the re-
168 gions of low cutoff rigidities are more accessible. Given that the perturbations of the Earth's
169 magnetosphere and the dynamics of the radiation belts are tightly coupled with the solar
170 wind parameters, we include the properties of the solar wind in our dataset.

171 The curated solar wind parameters are accessible via the OMNIWeb⁷ online dataset
172 with a 5-minute cadence. We include all major properties, such as the solar wind den-
173 sity, temperature, the magnitude, and all three components of the velocities and mag-
174 netic fields. All data from OMNIWeb is checked for continuity, and no affecting data gaps
175 were found.

176 **2.4 Energetic particles**

177 Radiation levels at aviation altitudes might significantly increase (even exceed the
178 galactic cosmic ray component) during solar energetic particle events (Kataoka et al.,
179 2018). Tobiska et al. (2018) also noted that there are many enhanced radiation events
180 in ARMAS data obtained for L-shells between 1.5 and 5 in the Western hemisphere, and
181 indicated the radiation belt particles as a possible reason for this. Therefore, we incor-
182 porate the energetic proton and electron measurements in our database, all measured
183 by the Geostationary Operational Environmental Satellite (GOES) series. The proton
184 fluxes are obtained in 7 integrated energy channels: ≥ 1 MeV, ≥ 5 MeV, ≥ 10 MeV, ≥ 30 MeV,

⁶ <https://www.nmdb.eu/nest/>

⁷ <https://omniweb.gsfc.nasa.gov/>

185 ≥ 50 MeV, ≥ 60 MeV, and ≥ 100 MeV. The data are obtained from the Integrated Space
186 Weather Analysis (ISWA) webapp⁸. The energetic electron flux measurements are in-
187 cluded only in one energy channel, ≥ 20 keV, because of the data continuity issues found
188 in other energy channels. The resolution of both the proton and electron particle data
189 is 5 min.

190 2.5 Solar X-ray measurements

191 Solar soft X-ray emission integrated over the solar disk is traditionally measured
192 by the GOES spacecraft series, using the onboard X-ray sensor. We are utilizing 1-minute
193 averaged fluxes in two channels, 1–8 Å and 0.5–4 Å, as input to our dataset. A simi-
194 lar type/cadence of the data was incorporated in the currently available Radiation Data
195 Portal (Sadykov et al., 2021).

196 2.6 Geomagnetic activity indexes

197 We adopt hourly planetary Kp, Ap, and Dst indexes, all available via the OMNI-
198 Web online dataset. The presented indices reflect the state of the Earth’s magnetosphere
199 during the major interplanetary disturbances, such as the presence of the Interplanetary
200 Coronal Mass Ejections (ICMEs) or high-speed stream interaction regions in the solar
201 wind. The variations in these parameters can affect the geomagnetic cutoff rigidity (Ptitsyna
202 et al., 2021) and, therefore, modulate the galactic cosmic ray precipitation into the Earth’s
203 atmosphere.

204 2.7 Global solar activity and characteristics

205 In addition to the transient activity effects of the Sun on the Geospace, there are
206 longer-term effects associated with the global solar activity. Although the neutron moni-
207 tor data can capture effects related to the global solar activity (Section 2.2), the solar
208 activity indexes make a picture of the evolution of the radiation environment more com-
209 plete. For example, Koldobskiy et al. (2022) found the delay of about 7 months of the
210 sunspot numbers with respect to the neutron monitor measurements. Therefore, we in-
211 clude the following global solar activity indicators:

- 212 • **Solar F10.7 index.** The Solar F10.7 index, which measures solar radio flux at
213 a wavelength of 10.7 cm (2800 MHz), is a critical indicator of solar activity and has
214 significant implications for Earth’s atmospheric conditions.
- 215 • **Daily sunspot number,** which tracks the daily count of visible sunspots on the
216 Sun’s surface, is a traditional characteristic of solar activity level. Managed by the
217 World Data Center SILSO⁹ (Sunspot Index and Long-term Solar Observations),
218 hosted at the Royal Observatory of Belgium, these data provide critical data for
219 understanding solar cycles and their impacts on space weather.
- 220 • **Solar polar fields.** The solar polar magnetic fields, monitored by the Wilcox So-
221 lar Observatory¹⁰ (WSO) at Stanford University, are crucial indicators of the Sun’s
222 magnetic activity and the solar cycle’s progression. These fields, located at the
223 Sun’s poles, reverse approximately every 11 years, marking a new solar cycle. The
224 strength and structure of the polar fields are key to predicting the magnitude of
225 future solar cycles, as they play a central role in the generation of the Sun’s mag-
226 netic field through the solar dynamo process. The WSO has provided near-continuous,
227 detailed measurements of solar polar magnetic fields since 1976. It is an inval-

⁸ <https://ccmc.gsfc.nasa.gov/tools/ISWA/>

⁹ <https://www.sidc.be/SILSO/datafiles>

¹⁰ <http://wso.stanford.edu/Polar.html>

228 able resource for understanding solar dynamics, space weather prediction, and their
229 influence on the heliosphere.

230 The illustration of some of these sources is presented in Figure 2. As one can see,
231 most features are continuous across the time interval of interest. The only noticeable change
232 occurs in high-energy proton flux (≥ 50 MeV) with a jump in the background. However,
233 since the protons of only a very high energy and flux could impact the radiation levels,
234 the jumps in the proton flux background levels are not a point of concern.

235 3 ML-Ready Data Set Construction

236 Although ARMAS data and the corresponding Geospace environment character-
237 istics have been collected and cleaned for the time interval of interest (June 2013 – De-
238 cember 2023), the data is still not ready for ML purposes. The ML-ready dataset must
239 meet the following characteristics, such as data integrability and understandability. Specif-
240 ically, the dataset should contain all components necessary for training and evaluating
241 the ML models and results interpretability (data integrability) and should be accessi-
242 ble to non-domain experts with a reasonably low effort (data understandability). To meet
243 the data integrability requirements, at least two more steps are required: the ARMAS
244 measurements should be merged with the preceding environmental parameters, and the
245 dataset has to be partitioned. To support the understandability of the data, partitions
246 should be created with input from domain experts to prevent artificial correlations.

247 We associate individual ARMAS measurements with the temporally closest pre-
248 ceding measurements of the Geospace environment. We avoid using any temporal inter-
249 polations to avoid breaking the causality principle and ensure that we are not utilizing
250 information from times in the future for the radiation nowcasting (Figure 3, left). Given
251 that some parameters, like the daily sunspot numbers, have a relatively large time ca-
252 dence (24 hours in this case), the preceding measurement of the daily sunspot number
253 may represent the conditions up to 24 hours before the actual measurements by ARMAS.

254 The second step is the separation of the individual ARMAS measurements into the
255 partitions (Figure 3, right). This process indicates the chunks of data that can be used
256 for training, validation, and testing. This ensures the reproducibility of the research and
257 the direct comparison between different models if trained on the same partitions. The
258 partitioning previously was used in other ML-ready datasets, such as the Space Weather
259 Analytics for Solar Flares, SWAN-SF (Angryk et al., 2020). The key point is that the
260 ARMAS measurements cannot be separated randomly. The data points from the same
261 flight could be just one minute away from each other. Therefore, the environmental prop-
262 erties are similar. Thus, if these data are split between the train and test partitions, this
263 could generate artificial correlation between the partitions. This effect was previously
264 recognized as a 'temporal coherence' for the problem of flare forecasting (Ahmadzadeh
265 et al., 2021). Therefore, the data from every ARMAS flight must be allocated to the same
266 partition.

267 In principle, each partition should represent the entire dataset, meaning the dis-
268 tribution of the parameters in each partition should represent the entire dataset (Liu et
269 al., 2019). Separation into partitions could be a challenging problem if data coverage is
270 sparse, and the problem is multi-dimensional. Unfortunately, this is the case for the con-
271 sidered dataset: the measurements are acquired along slightly more than 1000 individ-
272 ual flight trajectories, and more than 40 Geospace parameters are associated with ev-
273 ery measurement. To overcome the manual trial-and-error partitioning process, we uti-
274 lized an ML clustering algorithm to help with it. The steps are as follows:

- 275 1. **Down-selection of the five parameters for clustering.** As was indicated above,
276 the number of Geospace parameters is relatively large. Distance-based clustering

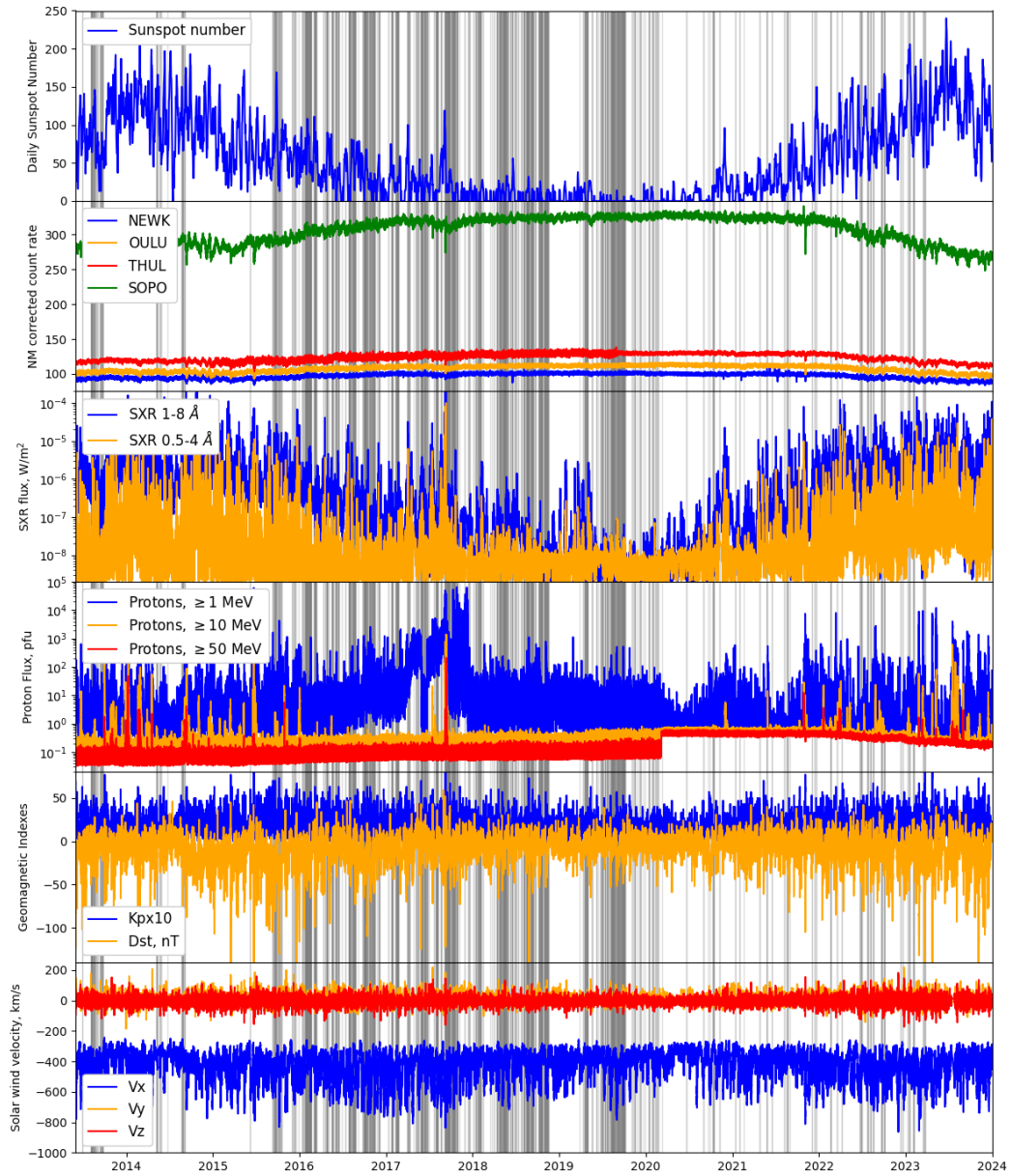


Figure 2. Evolution of some selected Geospace environment parameters from June 2013 until December 2024, from top to bottom: daily sunspot number, neutron monitor corrected counts from four stations considered, integrated soft X-ray fluxes, energetic proton fluxes, geomagnetic indexes (Kp and Dst), and solar wind velocities at L1. Gray lines in the background represent the time moments covered by the ARMAS flight measurements.

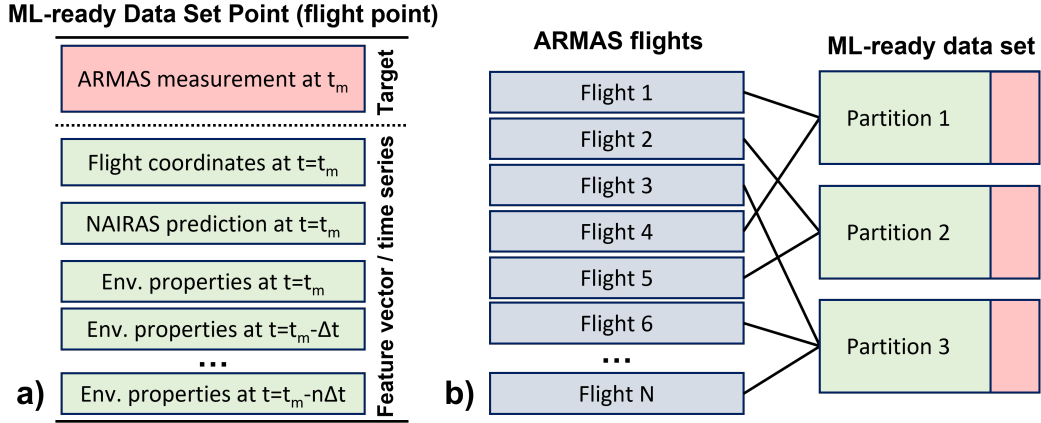


Figure 3. (a) Schematic structure of the ML-ready dataset entity. A target corresponds to the measurement of radiation dose rate during the ARMAS flight, and the feature vector represents the flight timing and coordinates, NAIRAS predictions, and prehistory of the measurements of the environment. (b) Illustration of the subdivision of ARMAS flights into partitions.

277 algorithms struggle with high dimensionality (the so-called 'curse of dimensionality'
 278 problem). However, many of the parameters are expected to be correlated
 279 with each other. For example, there is a strong correlation between all geomag-
 280 netic indices (such as K_p , A_p , Dst), global solar activity parameters (daily sunspot
 281 numbers, F10.7 flux, polar magnetic field measurements), description of the lo-
 282 cation in geomagnetic coordinates (geomagnetic latitude, L-shell, cutoff rigidity),
 283 etc. Instead of considering all of them, we down-select five parameters based on
 284 which we will cluster the data. The first three are related to the location of the
 285 measurement: geomagnetic longitude, geomagnetic latitude, and barometric flight
 286 altitude. The fourth parameter describes the geomagnetic activity, for which we
 287 use the Dst index. The last parameter we use is the daily sunspot number that
 288 reflects global solar activity levels. Although this selection is not unique, it is suf-
 289 ficient for partitioning purposes. Because the number of ARMAS flights during
 290 the solar energetic particle (SEP) events is very low, we did not use the SEP-related
 291 measurements for partitioning.

292 **2. Clustering of individual measurements.** We apply the Gaussian Mixture Model
 293 (GMM) clustering based on the five parameters for the individual ARMAS mea-
 294 surements. The number of GMM components was selected to be 100, which is suf-
 295 ficiently large to create a relatively equal representation for all parameters in ev-
 296 ery partition. We note here that the GMM is a soft clustering methodology as it
 297 assigns the probabilities of every point to belong to a certain cluster rather than
 298 associating it with a particular cluster.

299 **3. Associating flights with clusters.** Each individual measurement in each flight
 300 now has a probability (or a weight) of belonging to each of the 100 clusters. To
 301 associate the entire flight with the particular cluster, we sum up the probabilities
 302 of the individual measurements in this flight for every cluster. Therefore, a flight
 303 becomes assigned entirely to a particular cluster for which the sum of the proba-
 304 bilities of the individual measurements is the highest. We perform these assign-
 305 ments for all flights in our dataset.

306 **4. Separation of the flights into three partitions.** We start from the first cluster
 307 and distribute the flights to this cluster into three partitions in sequential order.
 308 After distributing all flights in the cluster, we move to the next cluster and
 309 repeat the procedure starting from the next partition. For example, if cluster 1

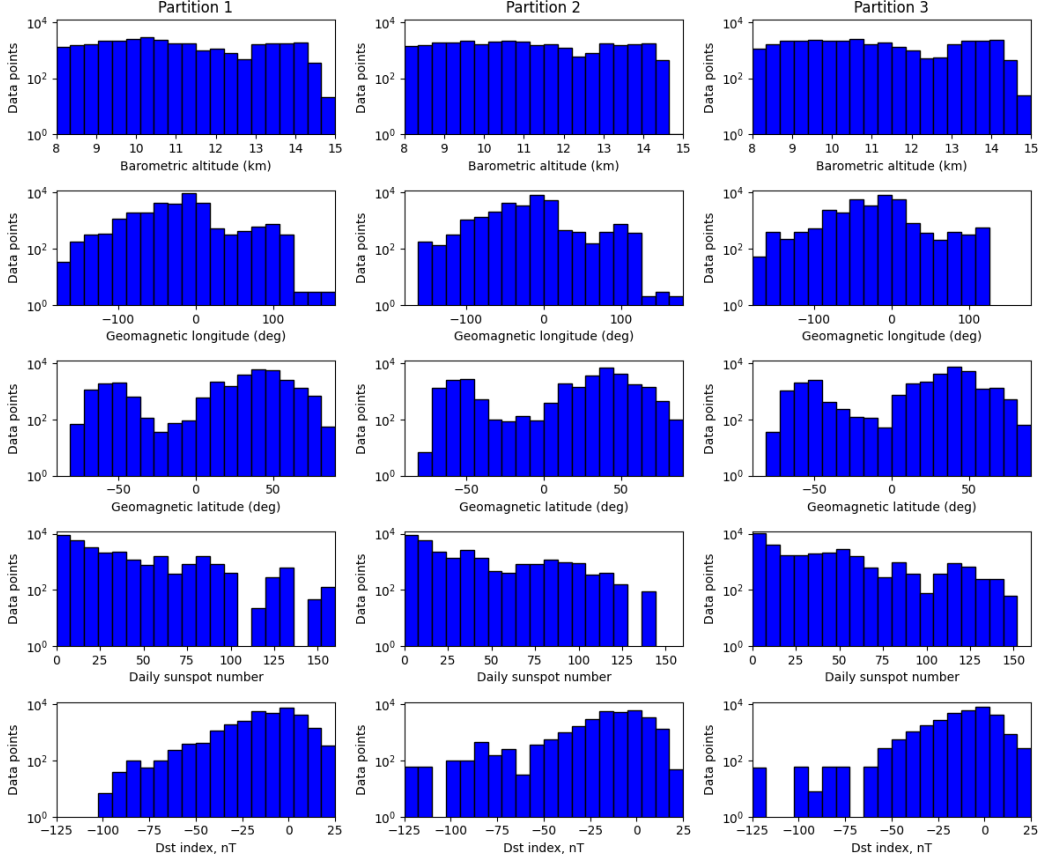


Figure 4. Distribution of the parameters used for clustering the data points (barometric altitude, geomagnetic longitude, geomagnetic latitude, daily sunspot number, and Dst index) within each partition of the dataset. Each row corresponds to a single parameter. The partition is indicated in the header of each column.

310 has four flights, then flight #1 will be distributed to partition 1, flight #2 – to partition
 311 partition 2, flight #3 – to partition 3, and flight #4 – again to partition 1. Then one
 312 moves to cluster 2, and flight #1 in this cluster goes now to partition 2, flight #2 –
 313 to partition 3, etc. This procedure ensures that all partitions will include flights
 314 belonging to the same cluster (and, therefore, likely sampling similar spatial loca-
 315 tions and geomagnetic and solar activity levels).

316 While the partitioning strategy does not necessarily result in the ‘most optimal’
 317 distribution of the flights, it leads to a satisfactory representation of parameter combi-
 318 nations in every partition while avoiding the brute-force partitioning. The result of the
 319 distribution of the parameters within each partition is presented in Figure 4. The his-
 320 tograms of the parameters are relatively similar for every column of the partitioning, in-
 321 dicated that each partition samples the parameter space relatively equally. Some dif-
 322 ferences are evident only for the extreme values of the parameters, such as the Dst in-
 323 dexes at around ~ -100 nT or the daily sunspot numbers of ~ 150 , which occur because
 324 not too many ARMAS flights have occurred during such high activity levels.

325 It is important to check if other parameters are sampled evenly across the parti-
 326 tions. Figure 5 illustrates the distributions of the parameters in three partitions that have

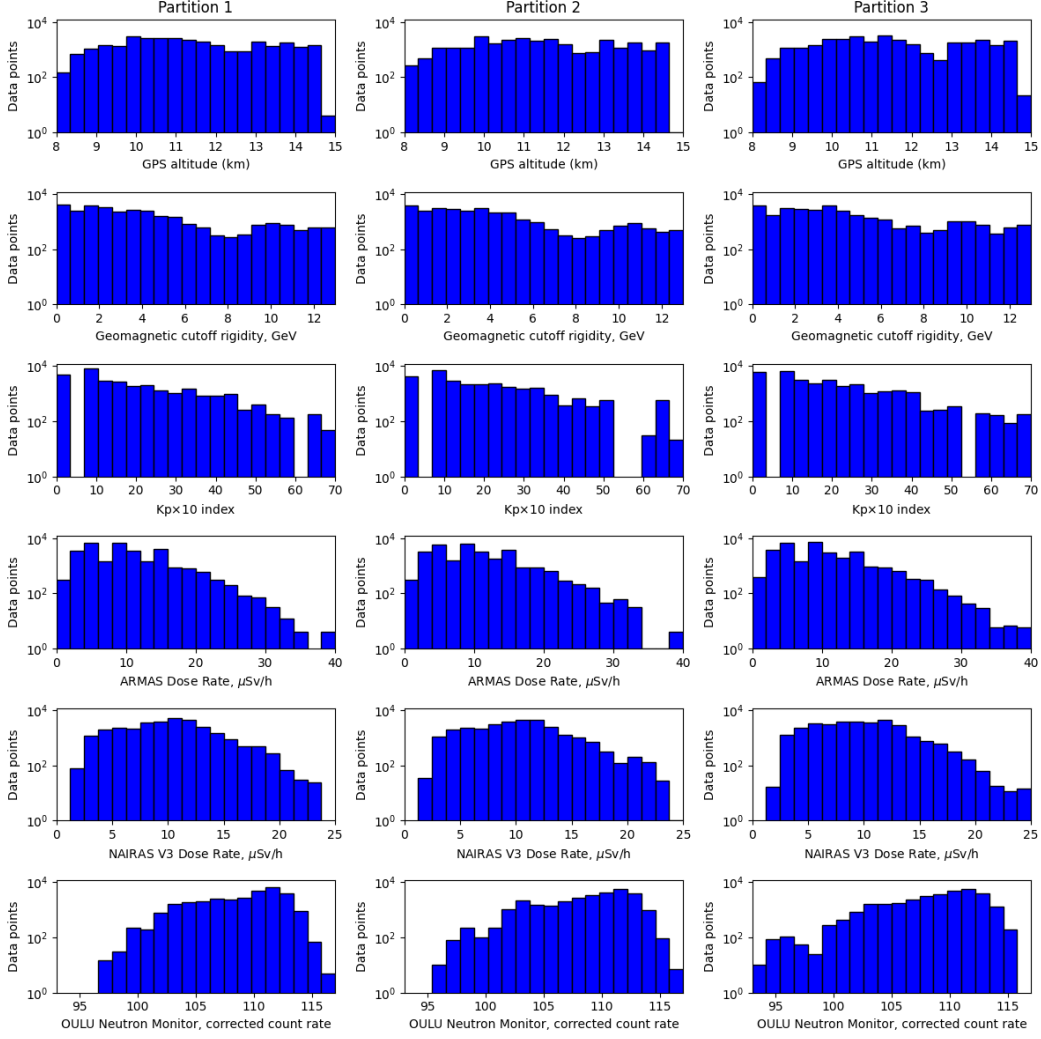


Figure 5. Distribution of the parameters that were not used for the clustering of the data points (GPS altitude, geomagnetic cutoff rigidity, Kp index multiplied by 10, ARMAS dose rate, NAIRAS dose rate, and corrected count rates of secondary cosmic ray neutrons detected by OULU station) within each partition of the dataset. Each row corresponds to a single parameter. The partition is indicated in the header of each column.

327 not directly participated in the clustering. The distribution of parameters is very sim-
 328 ilar, indicating a successful partitioning process.

329 The resulting data represents the ML-ready dataset that utilizes the last ‘snapshot’
 330 of Geospace properties before the ARMAS measurement. The ML-ready dataset con-
 331 sists of a feature vector (a set of characteristics based on which the prediction is made,
 332 representing Geospace parameters) and a target (a characteristic to predict, here AR-
 333 MAS measurement). The feature vector comprises the flight timing and location prop-
 334 erties, as well as the most recent properties of the environment. In addition to the most
 335 recent properties (Figure 3), one can provide the time series of the evolution of the Geospace
 336 parameters up to the time before the measurements. The preceding properties of the en-
 337 vironment can be forward-interpolated with the cadence Δt for n time steps before the
 338 ARMAS measurement. At the same time, the partitioning of the individual ARMAS mea-

339 surements and related time series remains as described above. Following this, we have
340 finally constructed three publicly available ML-ready datasets:

- 341 • The dataset that represents the most recent Geospace measurements and does not
342 involve their time series (‘static’ dataset, $n = 0$);
- 343 • The dataset that includes a 1-hour prehistory of the Geospace measurements be-
344 fore ARMAS flight measurement (‘dynamic’ dataset 1, $n = 12$ and $\Delta t = 5$ min);
- 345 • The dataset that includes a 24-hour prehistory of the Geospace measurements be-
346 fore ARMAS flight measurement (‘dynamic’ dataset 2, $n = 24$ and $\Delta t = 1$ hour);

347 All three versions of the datasets are currently accessible via the Radiation Data
348 Portal¹¹

349 4 Data Set Use Case Example

350 In this section, we illustrate how the constructed datasets can be utilized for ML-
351 driven forecasting of atmospheric radiation. For this demonstration, we proceed with the
352 simplest version of the three datasets constructed and described above, the ‘static’ ML-
353 ready dataset, which includes only the latest point-in-time measurement of every Geospace
354 parameter. We use partition 1 of this dataset to train the ML model described below,
355 and partition 2 for the evaluation of the performance of the model and its comparison
356 with the predictions of the NAIRAS-v3 physics-based model. Among the ML approaches
357 available off-the-shelf, we select the Random Forest regressor (Breiman, 2001), a bag-
358 ging tree-based ensemble learning algorithm. Random Forest has previously been suc-
359 cessfully applied to a variety of classification and regression problems in the space physics
360 domain, including the prediction of solar flares and solar energetic particle events (Liu
361 et al., 2017; O’Keefe et al., 2022), forecasting the duration of enhanced soft X-ray ra-
362 diation during the flare (Reep & Barnes, 2021), the timing of the solar wind propaga-
363 tion (Baumann & McCloskey, 2021), ion-kinetic instability detection in the solar wind-
364 like plasmas (Sadykov et al., 2025), etc. We note here that, despite the promising track
365 record, one cannot guarantee that the Random Forest approach is the most optimal for
366 the considered problem. Therefore, the survey of other ML models is highly encouraged
367 and is one of the authors’ goals for the future.

368 We have used the Random Forest model available at the `scikit-learn` Python li-
369 brary package (Pedregosa et al., 2011). The model has several hyperparameters to op-
370 timize, such as the number of individual decision trees in the ensemble, the maximum
371 depth restriction for each tree, the number of features randomly selected and propagated
372 into every tree, etc. Typically, the hyperparameters are fine-tuned during the cross-validation
373 phase when the models of different hyperparameters are evaluated on a designated par-
374 tition or a subset of the training partition. Here, we do not perform the detailed cross-
375 validation but rather use the parameters we found to perform satisfactorily for the con-
376 sidered study during our preliminary tests. Our Random Forest consists of 100 decision
377 trees of a depth of 10 or less, with at least 2 samples from the dataset required for the
378 split within the tree, and at least 4 samples to be at the leaf node. To guide the train-
379 ing process, we utilize the mean squared error as a measure of the regressor’s performance.

380 The results are presented in Figure 6. The left panel illustrates the scatterplot of
381 the radiation dose rates predicted by the ML-driven model against the ARMAS mea-
382 surements used as a ground truth. One can see that points tend to be mostly organized
383 along the red line, which represents the ideal one-to-one prediction. It is also visible that
384 the predictions demonstrate a stronger spread and systematic deviation from the ideal
385 prediction line for the larger values of dose rates. In fact, the part of the distribution cov-

¹¹ <https://dmlab.cs.gsu.edu/rdp/ml-dataset.html>

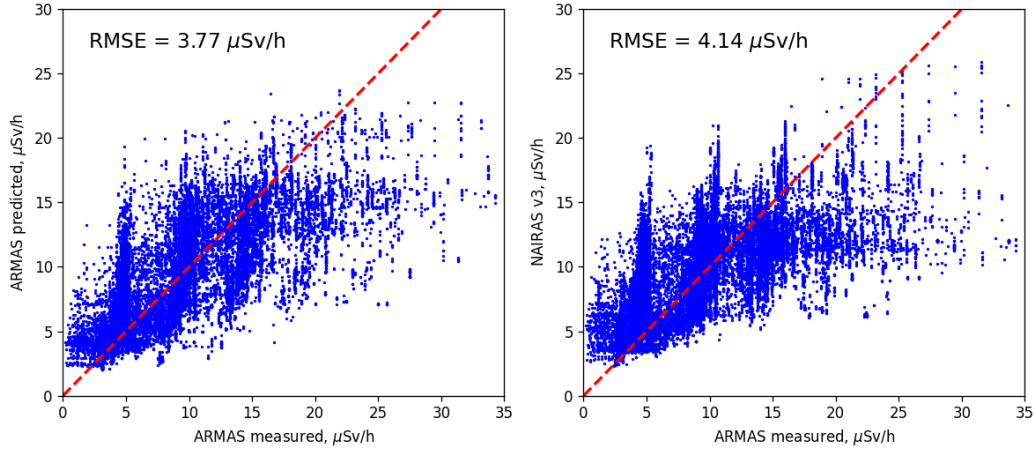


Figure 6. Left: Measured radiation dose rates VS predicted using an ML model (Random Forest Regressor). The ML model was trained on partition 1 of the static dataset and evaluated on partition 2. Right: Measured radiation dose rates VS nowcast of a physics-based NAIIRAS V3 model for partition 2.

386 ering the dose rates of $>15 \mu Sv/h$ is mostly situated below the perfect prediction line,
 387 indicating that the ML-driven predictions typically underestimate the actual dose rates
 388 in the cases where the measured radiation doses are high. Interestingly, the comparison
 389 of ARMAS measurements with the nowcast of the physics-based NAIIRAS-v3 model (pre-
 390 sented in the right panel of Figure 6) demonstrates the same pattern. The root mean
 391 squared error (RMSE) computed for the ML-driven prediction is $3.77 \mu Sv/h$, whereas
 392 it is $4.14 \mu Sv/h$ for the NAIIRAS-v3 model. One can see that the ML-driven model seems
 393 to deliver slightly more accurate predictions both in terms of the RMSE measure across
 394 the test partition and based on the qualitative appearance of the scatterplots in Figure 6.
 395 Overall, this emphasizes the potential for ML approaches in nowcasting the radiation
 396 at aviation altitudes. The potential next steps can include generalizing the results over
 397 various train-validation-test partition combinations, considering other ML algorithms,
 398 and involving the time series data of Geospace parameters for forecast improvement. The
 399 constructed ML-ready datasets enable such investigations, opening the perspective of the
 400 new studies of the research community.

401 5 Summary and Discussion

402 This paper presents the construction of the ML-ready dataset for nowcasting at-
 403 mospheric radiation at aviation altitudes. We have leveraged the publicly available ef-
 404 fective dose rate measurements by the ARMAS device acquired over 589 flights, and cre-
 405 ated a pre-partitioned ML-ready dataset, which requires minimal preprocessing to be used
 406 for ML purposes. Some of the dataset features are summarized below:

- 407 • The resulting dataset comprises 589 ARMAS flights containing 92,476 effective
 408 dose rate measurements. While the flights are mostly accomplished on top of the
 409 continental US and territories (Figure 1), the dataset also samples the regions over
 410 the Pacific Ocean, the North Atlantic region, and Antarctica.
- 411 • The radiation measurements are supported by the measurements from four neu-
 412 tron monitor stations (OULU, DOMC, NEWK, and THUL), solar wind properties at L1,
 413 measurements of proton and soft X-ray fluxes by GOES spacecraft series, geomag-

netic activity indexes, and global solar activity characteristics (such as sunspot numbers, F10.7 flux, and solar polar fields). This provides an opportunity for comprehensive studies of the dependencies of the radiation environment on the Geospace drivers.

- The dataset has been pre-partitioned into three subsets, which can be directly used for training, validation, and testing purposes. During pre-partitioning, it was ensured that the data points from the same flight are within the same partition, and that the sampling of the parameter space is more or less the same within any partition (see Figures 4 and 5 for parameter distributions).
- Three versions of the dataset are constructed. The ‘static’ version includes the most recent properties of the environment only. The ‘dynamic’ versions include the pre-history of Geospace parameters as time series. The 1-hour and 24-hour long pre-history is considered.
- The use case example demonstrates on the selected test subset the possibility of constructing an ML model that predicts the effective dose rates with the average root mean squared error of $3.77\mu\text{Sv/h}$, which is slightly better than the NAIRAS v3 physics-based model nowcast ($4.14\mu\text{Sv/h}$). This holds promise for the development of ML-driven models of radiation forecasting in the future.

We envision that the constructed datasets could be used for various scientific applications. First, the problem of the atmospheric radiation nowcast given the state of the environment requires a detailed evaluation with respect to the ML algorithms and involvement of the time series. As highlighted in Section 4, even the static version of the dataset with the Random Forest regressor demonstrates the promising results. The nowcasting and forecasting results can vary with the algorithm (Ali et al., 2024; Goodwin et al., 2024; O’Keefe et al., 2024); therefore, the consideration of other ML methods is necessary. The benefits of including time series properties need to be assessed as well. Second, the slight manipulation of the dataset allows us to consider a forecasting problem instead of nowcasting. All one has to do is consider the ‘dynamic’ dataset versions and avoid considering the data within a certain latency window before the ARMAS measurement. Obviously, the length of the time series (1 h and 24 h) would limit the considered latency windows. Third, the dataset could enhance the understanding of radiation environment physics. Besides the standard correlation analyses possible, one could investigate feature importance to understand the influence of some Geospace parameters on the radiation environment (e.g., Yeolekar et al., 2021; Sadykov & Kosovichev, 2017). The developed dataset opens a promising number of prospective studies and facilitates the development of models for the aviation radiation domain.

Open Research

The developed machine learning-ready dataset to nowcast the effective dose rates at aviation altitudes is currently publicly available via the Radiation Data Portal (<https://dmlab.cs.gsu.edu/rdp/ml-dataset.html>). The original ARMAS data files are publicly available from the ARMAS Data Archive at Space Environment Technologies (<https://sol.spacenvironment.net/ARMAS/Archive/>). The Neutron Monitor can be accessed via the NMDB database (<https://www.nmdb.eu/>). The solar wind measurements, geomagnetic activity indexes, and F10.7 index data are publicly accessible via the OMNI-Web service (<https://omniweb.gsfc.nasa.gov/>). The GOES energetic particle flux and soft X-ray data can be accessed via the National Oceanic and Atmospheric Administration National Centers for Environmental Information Archive (NOAA NCEI, <https://www.ncei.noaa.gov/>). We thank the developers of the Integrated Space Weather Analysis (ISWA, <https://cmc.gsfc.nasa.gov/tools/ISWA/>) systems API for the possibility of retrieving the GOES integrated proton fluxes since 2020. The daily sunspot number is obtained via the publicly accessible World Data Center SILSO (<https://www.sidc>

465 .be/SILSO/datafiles). The solar polar field measurements are publicly available via
466 the Wilcox Solar Observatory website (<http://wso.stanford.edu/Polar.html>).

467 Acknowledgments

468 This work has been supported by the NASA HITS grant 80NSSC22K1561. VMS also
469 thanks the NSF FDSS grant 1936361 and NASA LWS grant 80NSSC24K1111. We ac-
470 knowledge the NMDB database (www.nmdb.eu), founded under the European Union's
471 FP7 programme (contract no. 213007) for providing data, and the PIs of individual neu-
472 tron monitors at: Newark and Thule (University of Delaware Department of Physics and
473 Astronomy and the Bartol Research Institute, USA), Kerguelen (Observatoire de Paris
474 and the French Polar Institute IPEV, France), Oulu (Sodankyla Geophysical Observa-
475 tory of the University of Oulu, Finland), South Pole (University of Wisconsin, River Falls,
476 USA)

477 References

- 478 Ahmadzadeh, A., Aydin, B., Georgoulis, M. K., Kempton, D. J., Mahajan, S. S.,
479 & Angryk, R. A. (2021, June). How to Train Your Flare Prediction Model:
480 Revisiting Robust Sampling of Rare Events. *The Astrophysical Journal Supple-*
481 *ment Series*, 254(2), 23. doi: 10.3847/1538-4365/abec88
- 482 Ali, A., Sadykov, V., Kosovichev, A., Kitiashvili, I. N., Oria, V., Nita, G. M.,
483 ... Marroquin, R. D. (2024, January). Predicting Solar Proton Events
484 of Solar Cycles 22-24 Using GOES Proton and Soft-X-Ray Flux Fea-
485 tures. *The Astrophysical Journal Supplement Series*, 270(1), 15. doi:
486 10.3847/1538-4365/ad0a6c
- 487 Angryk, R. A., Martens, P. C., Aydin, B., Kempton, D., Mahajan, S. S., Ba-
488 sodi, S., ... Georgoulis, M. K. (2020, January). Multivariate time series
489 dataset for space weather data analytics. *Scientific Data*, 7(1), 227. doi:
490 10.1038/s41597-020-0548-x
- 491 Baumann, C., & McCloskey, A. E. (2021, June). Timing of the solar wind propaga-
492 tion delay between L1 and Earth based on machine learning. *Journal of Space*
493 *Weather and Space Climate*, 11, 41. doi: 10.1051/swsc/2021026
- 494 Breiman, L. (2001, January). Random Forests. *Machine Learning*, 45, 5-32. doi: 10
495 .1023/A:1010933404324
- 496 Cho, G., Kim, J. H., Park, T. S., & Cho, K. (2017). Proposing a simple radia-
497 tion scale for the public: radiation index. *Nuclear Engineering and Technology*,
498 49(3), 598-608.
- 499 Copeland, K. (2021, March). *Cari-7 documentation: Radiation transport*
500 *in the atmosphere* (Tech Report No. DOT/FAA/AM-21/05). United
501 States. Department of Transportation. Federal Aviation Administration.
502 Office of Aviation. Civil Aerospace Medical Institute. Retrieved from
503 <https://rosap.ntl.bts.gov/view/dot/57224>
- 504 Goodwin, G. T., Sadykov, V. M., & Martens, P. C. (2024, April). Investigating Per-
505 formance Trends of Simulated Real-time Solar Flare Predictions: The Impacts
506 of Training Windows, Data Volumes, and the Solar Cycle. *The Astrophysical*
507 *Journal*, 964(2), 163. doi: 10.3847/1538-4357/ad276c
- 508 Kataoka, R., Sato, T., Miyake, S., Shiota, D., & Kubo, Y. (2018, July). Radiation
509 Dose Nowcast for the Ground Level Enhancement on 10-11 September 2017.
510 *Space Weather*, 16(7), 917-923. doi: 10.1029/2018SW001874
- 511 Koldobskiy, S. A., Kähkönen, R., Hofer, B., Krivova, N. A., Kovaltsov, G. A., &
512 Usoskin, I. G. (2022, March). Time Lag Between Cosmic-Ray and Solar
513 Variability: Sunspot Numbers and Open Solar Magnetic Flux. *Solar Physics*,
514 297(3), 38. doi: 10.1007/s11207-022-01970-1
- 515 Kozlov, V., Kudela, K., Starodubtsev, S., Turpanov, A., Usoskin, I., & Yanke, V.

- (2003, September). Neutron monitor database in real time. In A. Wilson (Ed.), *Solar variability as an input to the earth's environment* (Vol. 535, p. 675-678).
- Liu, C., Deng, N., Wang, J. T. L., & Wang, H. (2017, July). Predicting Solar Flares Using SDO/HMI Vector Magnetic Data Products and the Random Forest Algorithm. *The Astrophysical Journal*, *843*(2), 104. doi: 10.3847/1538-4357/aa789b
- Liu, H., Chen, S.-M., & Cocea, M. (2019). Subclass-based semi-random data partitioning for improving sample representativeness. *Information Sciences*, *478*, 208-221. Retrieved from <https://www.sciencedirect.com/science/article/pii/S0020025518308910> doi: <https://doi.org/10.1016/j.ins.2018.11.002>
- Masson, A., Fung, S. F., Camporeale, E., Kuznetsova, M. M., Poedts, S., Barnum, J., ... Cecconi, B. (2024). Heliophysics and space weather information architecture and innovative solutions: Current status and ways forward. *Advances in Space Research*. Retrieved from <https://www.sciencedirect.com/science/article/pii/S0273117724004939> doi: <https://doi.org/10.1016/j.asr.2024.05.052>
- Matthiä, D., Meier, M. M., & Reitz, G. (2014, March). Numerical calculation of the radiation exposure from galactic cosmic rays at aviation altitudes with the PANDOCA core model. *Space Weather*, *12*(3), 161-171. doi: 10.1002/2013SW001022
- Mertens, C. J., Gronoff, G. P., Zheng, Y., Buhler, J., Willis, E., Petrenko, M., ... Minow, J. (2024). Nairas atmospheric and space radiation environment model. *IEEE Transactions on Nuclear Science*, *71*(4), 618-625. doi: 10.1109/TNS.2023.3330675
- Mertens, C. J., Gronoff, G. P., Zheng, Y., Petrenko, M., Buhler, J., Phoenix, D., ... Minow, J. (2023, May). NAIRAS Model Run-On-Request Service at CCMC. *Space Weather*, *21*(5), e2023SW003473. doi: 10.1029/2023SW003473
- Mertens, C. J., Meier, M. M., Brown, S., Norman, R. B., & Xu, X. (2013, October). NAIRAS aircraft radiation model development, dose climatology, and initial validation. *Space Weather*, *11*(10), 603-635. doi: 10.1002/swe.20100
- Nita, G., Ahmadzadeh, A., Criscuoli, S., Davey, A., Gary, D., Georgoulis, M., ... Wang, J. T. L. (2022, March). Revisiting the Solar Research Cyberinfrastructure Needs: A White Paper of Findings and Recommendations. *arXiv e-prints*, arXiv:2203.09544. doi: 10.48550/arXiv.2203.09544
- O'Keefe, P. M., Sadykov, V., Kosovichev, A., Kitiashvili, I. N., Oria, V., Nita, G. M., ... Marroquin, R. D. (2024, December). The Random Hivemind: An ensemble deep learning application to the solar energetic particle prediction problem. *Advances in Space Research*, *74*(12), 6252-6263. doi: 10.1016/j.asr.2024.04.044
- O'Keefe, P. M., Sadykov, V. M., Kosovichev, A. G., Nita, G. M., Oria, V., Sharma, S., ... Jie Chong, C. (2022, June). *Handling Highly Imbalanced Data in Machine Learning Applications*. Zenodo. doi: 10.5281/zenodo.6780972
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, *12*, 2825-2830.
- Ptitsyna, N. G., Danilova, O. A., Tyasto, M. I., & Sdobnov, V. E. (2021, March). Dynamics of Cosmic-Ray Cutoff Rigidity and Magnetospheric Parameters during Different Phases of the Storm of November 20, 2003. *Geomagnetism and Aeronomy*, *61*(2), 169-179. doi: 10.1134/S0016793221010114
- Reames, D. V. (2021). *Solar Energetic Particles. A Modern Primer on Understanding Sources, Acceleration and Propagation* (Vol. 978). doi: 10.1007/978-3-030-66402-2
- Reep, J. W., & Barnes, W. T. (2021, October). Forecasting the Remaining Duration of an Ongoing Solar Flare. *Space Weather*, *19*(10), e02754. doi: 10.1029/

- 571 2021SW002754
572 Sadykov, V. M., Kitiashvili, I. N., Tobiska, W. K., & Guhathakurta, M. (2021, Jan-
573 uary). Radiation Data Portal: Integration of Radiation Measurements at the
574 Aviation Altitudes and Solar-Terrestrial Environment Observations. *Space*
575 *Weather*, *19*(1), e2020SW002653. doi: 10.1029/2020SW002653
- 576 Sadykov, V. M., & Kosovichev, A. G. (2017, November). Relationships between
577 Characteristics of the Line-of-sight Magnetic Field and Solar Flare Forecasts.
578 *The Astrophysical Journal*, *849*(2), 148. doi: 10.3847/1538-4357/aa9119
- 579 Sadykov, V. M., Ofman, L., Boardsen, S. A., Yogesh, Mostafavi, P., Jian, L. K.,
580 ... Martinović, M. (2025, May). Identification of Ion-Kinetic Instabilities in
581 Hybrid-PIC Simulations of Solar Wind Plasma with Machine Learning. *arXiv*
582 *e-prints*, arXiv:2505.18271. doi: 10.48550/arXiv.2505.18271
- 583 Tobiska, W. K., Atwell, W., Beck, P., Benton, E., Copeland, K., Dyer, C., ... Xap-
584 sos, M. A. (2015, April). Advances in Atmospheric Radiation Measurements
585 and Modeling Needed to Improve Air Safety. *Space Weather*, *13*(4), 202-210.
586 doi: 10.1002/2015SW001169
- 587 Tobiska, W. K., Bouwer, D., Smart, D., Shea, M., Bailey, J., Didkovsky, L., ...
588 Yoon, K. (2016, November). Global real-time dose measurements using the
589 Automated Radiation Measurements for Aerospace Safety (ARMAS) system.
590 *Space Weather*, *14*(11), 1053-1080. doi: 10.1002/2016SW001419
- 591 Tobiska, W. K., Didkovsky, L., Judge, K., Weiman, S., Bouwer, D., Bailey, J., ...
592 Fuschino, R. (2018, October). Analytical Representations for Characterizing
593 the Global Aviation Radiation Environment Based on Model and Measurement
594 Databases. *Space Weather*, *16*(10), 1523-1538. doi: 10.1029/2018SW001843
- 595 Väisänen, P., Usoskin, I., & Mursula, K. (2021, May). Seven Decades of Neu-
596 tron Monitors (1951-2019): Overview and Evaluation of Data Sources.
597 *Journal of Geophysical Research (Space Physics)*, *126*(5), e28941. doi:
598 10.1029/2020JA028941.10.1002/essoar.10505091.1
- 599 Yeolekar, A., Patel, S., Talla, S., Puthucode, K. R., Ahmadzadeh, A., Sadykov,
600 V. M., & Angryk, R. A. (2021, December). Feature Selection on a
601 Flare Forecasting Testbed: A Comparative Study of 24 Methods . In
602 *2021 international conference on data mining workshops (icdmw)* (p. 1067-
603 1076). Los Alamitos, CA, USA: IEEE Computer Society. Retrieved from
604 <https://doi.ieeecomputersociety.org/10.1109/ICDMW53433.2021.00138>
605 doi: 10.1109/ICDMW53433.2021.00138