Evaluating Proton Intensities for the SMILE Mission

Simon Patrick Mischel¹, Elena A. Kronberg², and Christopher Philippe Escoubet³

¹Ludwig-Maximilians-Universität München ²Ludwig Maximilian University of Munich ³ESA / ESTEC

April 26, 2024

Abstract

This study introduces five linear regression models developed to accurately predict proton intensities in the critical energy range of 92.2 keV to 159.7 keV. To achieve this task we utilized 14 years of data sourced from the Cluster's RAPID experiment and NASA's OMNI database. This data was then aligned with the Solar wind-Magnetosphere-Ionosphere Link Explorer (SMILE) mission's trajectory, to increase model accuracy in the relevant regions. Our approach diverges from existing methodologies by offering a user-friendly model that doesn't require specialized software, making it accessible for broader applications in satellite mission planning and risk assessment. The research segregates the dataset into four distinct regions, each analyzed for proton intensity dynamics. In the outer regions (|YGSE| [?] 6 Re) there is a pronounced dependence on radial distance and solar wind speed. In contrast, the inner regions (|YGSE| [?] 6 Re) demonstrate a significant dependence of proton intensities on the absolute value of the z-coordinate and the magnetic field line topology. Our models achieved a Spearman correlation ranging from 0.57 to 0.72 on the test set, indicating good predictive capabilities. The findings emphasize the role of regional characteristics in space weather prediction and underscore the potential for tailored approaches in future research.

Evaluating Proton Intensities for the SMILE Mission 1

2

3

Simon Mischel¹, Elena A. Kronberg¹, C.P. Escoubet²

¹Department of Earth and Environmental Sciences (Geophysics), Ludwig Maximilian University of Munich (LMU) Munich, Theresienstr. 41, Munich, D-80333, Germany ²European Space Research and Technology Centre, Noordwjik, Keplerlaan 1, 2201 AZ, The Netherlands 4 5

6	Key Points:
7	• Developed models predict proton intensities impacting satellites, aiding space weather
8	forecasting and mission planning.
9	• Different regions in space showcase distinct relations between proton intensities
10	and predicting parameters.
11	• Study findings highlight the importance of tailored approaches in space weather
12	prediction.

 $Corresponding \ author: \ Simon \ Mischel, \ \texttt{simonmischel@hotmail.com}$

13 Abstract

This study introduces five linear regression models developed to accurately predict pro-14 ton intensities in the critical energy range of 92.2 keV to 159.7 keV. To achieve this task 15 we utilized 14 years of data sourced from the Cluster's RAPID experiment and NASA's 16 OMNI database. This data was then aligned with the Solar wind-Magnetosphere-Ionosphere 17 Link Explorer (SMILE) mission's trajectory, to increase model accuracy in the relevant 18 regions. Our approach diverges from existing methodologies by offering a user-friendly 19 model that doesn't require specialized software, making it accessible for broader appli-20 cations in satellite mission planning and risk assessment. The research segregates the dataset 21 into four distinct regions, each analyzed for proton intensity dynamics. In the outer re-22 gions ($|\text{YGSE}| \ge 6 R_e$) there is a pronounced dependence on radial distance and solar 23 wind speed. In contrast, the inner regions ($|YGSE| \leq 6R_e$) demonstrate a significant 24 dependence of proton intensities on the absolute value of the z-coordinate and the mag-25 netic field line topology. Our models achieved a Spearman correlation ranging from 0.5726 to 0.72 on the test set, indicating good predictive capabilities. The findings emphasize 27 the role of regional characteristics in space weather prediction and underscore the po-28 tential for tailored approaches in future research. 29

³⁰ Plain Language Summary

We developed a new model to predict space weather, specifically focusing on pro-31 32 ton intensities, which can impact how well satellites work in space. We used 14 years of space observations to create five easy-to-use numerical models. These models are designed 33 to help with planning and protecting future satellite missions, such as the upcoming SMILE 34 mission, from space weather effects. In our study, we looked closely at different areas in 35 space around Earth. We found that in the outer areas ($|YGSE| \ge 6 R_e$), the distance 36 from Earth and the speed of the solar wind are important for understanding proton be-37 havior. However, in areas ($|YGSE| \leq 6R_e$), the height above Earth (measured along 38 the z-direction) and the type of magnetic field lines play a more significant role. This 39 shows us that different areas in space around Earth can be affected by space weather in 40 different ways. Our models did a good job of predicting these effects, showing that choos-41 ing a tailored approach can be useful when forecasting proton intensities. 42

43 **1** Introduction

Space weather events, driven by solar activities pose significant challenges to satel-44 lite operations and measurements. Notable examples include the European Space Agency's 45 (ESA) Cluster and X-ray Multi-Mirror (XMM-Newton) missions. The Cluster mission, 46 particularly its Research with Adaptive Particle Imaging Detector (RAPID)/Imaging 47 Electron Spectrometer (IES), has encountered challenges due to high proton intensities 48 leading to measurement contamination (Wilken et al., 1997; Kronberg et al., 2016; Kro-49 nberg, Daly, et al., 2021). Similarly, the X-ray telescope aboard the XMM-Newton space-50 craft experienced significant operational disruptions, with approximately 40% of its ob-51 servation time compromised due to background contamination (Walsh et al., 2014). Fur-52 thermore, investigations into the XMM-Newton telescope's susceptibility to soft protons 53 highlight proton intensities in the sub-100 to 300 keV range, particularly around 100 keV, 54 as the most damaging, leading to significant operational challenges and data contam-55 ination, see (Fioretti et al., 2016) and references therein. 56

The upcoming European-Chinese Solar wind-Magnetosphere-Ionosphere Link Explorer (SMILE) mission (Branduardi-Raymont et al., 2018), slated for launch in 2025, aspires to deepen our understanding of the Sun-Earth interaction, decoding space weather hazards and understanding energy entry into Earth's magnetosphere. While missions like ATHENA (Advanced Telescope for High Energy Astrophysics), which will maintain an orbit around the L2 Lagrange point with a continuously large distance to Earth, are ex-

pected to face minimal threats from soft protons (Perinati et al., 2024), SMILE with its 63 highly inclined elliptical orbit around earth, akin to that of Cluster and XMM-Newton, 64 will navigate through diverse magnetospheric regions, making it susceptible to soft pro-65 ton radiation. Particularly its Soft X-ray Imager (SXI) telescope, presents challenges con-66 cerning radiation exposure (Raab et al., 2016; Branduardi-Raymont & Wang, 2022). There-67 fore a critical component of achieving SMILE's objectives is the accurate prediction of 68 proton radiation levels, which significantly affect the Total Ionizing Dose (TID) and To-69 tal Non-Ionizing Dose (TNID) absorbed by the Charge-Coupled Devices (CCDs) of the 70 Soft X-ray Imager (SXI) (Hubbard et al., 2024). In recent discussions (M. Hubbard, per-71 sonal communication, 2023) the necessity for models that accurately estimate radiation 72 levels was emphasized, particularly in critical energy ranges below 300 keV, crucial for 73 the SMILE mission's success. A strong preference for models that are not only accurate 74 but also straightforward and interpretable was expressed. 75

To meet these needs, our study adopts a distinct approach compared to existing 76 research, which is based on machine learning black box models, such as the works of Kronberg 77 et al. (2020) and Kronberg, Hannan, et al. (2021). We aim to develop a simple, user-friendly 78 linear regression model leveraging data from the Cluster mission and NASA's OMNI database 79 (King & Papitashvili, 2005). Our model's simplicity and ease of use make it accessible 80 to a broader range of users, not requiring specialized software or extensive computational 81 resources. This approach not only contributes to the scientific understanding of space 82 weather phenomena but also offers practical tools for satellite mission planning and risk 83 assessment. 84

⁸⁵ 2 Data Analysis and Processing

86 87

2.1 Data Preparation: Adapting CLUSTER's Dataset for the SMILE Mission's Trajectory

The proton intensity data for this research was taken from the Cluster's RAPID 88 experiment, ranging from 2001 to 2015. The experiment captures 3-D energetic electron 89 and ion fluxes above approximately 30 keV using the Imaging Electron Spectrometer (IES) 90 and the Imaging Ion Mass Spectrometer (IIMS) instruments. Situated in the SCENIC 91 detector head, the IIMS instrument identifies ion energies and species. The methodol-92 ogy involves using start and stop signals produced from electrons emitted by an initial 93 thin foil on the solid-state detector's surface. The time-of-flight (TOF) between these 94 signals, combined with the known energy, discerns the species and energy channel (Daly 95 & Kronberg, 2023). This study specifically uses data from the p3 channel, targeting pro-96 ton intensities between 92.2 keV and 159.7 keV. As predicting parameters for solar, so-97 lar wind and geomagnetic activity we used variables from the OMNI database(https:// 98 omniweb.gsfc.nasa.gov/), see also King and Papitashvili (2005). 99

Aiming for a model tailored to the SMILE mission's trajectory (see Figure 1), data 100 filtering was imperative. Points not adhering to the following spatial parameters were 101 excluded: $-10.5R_e \le x \le 11.2R_e$; $-10.8R_e \le y \le 11.5R_e$; $z \le 18.5R_e$; $\sqrt{x^2 + y^2} \le 10.5R_e$; $\sqrt{x^2 + y^2} \ge 10$ 102 $11.6R_e; \sqrt{x^2+z^2} \leq 19.8R_e; \sqrt{z^2+y^2} \leq 20.0R_e$. These constraints were chosen by 103 rounding the maxima and minima of the spatial parameters to the nearest tenth. As a 104 result, we were left with a trimmed dataset, reduced from 1,172,923 to 462,615 data points. 105 Though compact, this dataset is centered on the space region significant for SMILE, promis-106 ing heightened model accuracy. It's noteworthy that negative z-values weren't excluded, 107 considering the SMILE mission's highly inclined, elliptical orbit, which dips to $-3.5 R_e$. 108 Omitting these would disregard vital data, especially since the Cluster's trajectory spent 109 a notable amount of time in the southern hemisphere. 110



Figure 1. SMILE mission's trajectory. Distinct colors represent individual years with Earth, having a radius of 6,731 km, at the center.

2.2 Predictor Introduction

111

The spacecraft's location in the Geocentric Solar Ecliptic (GSE) coordinate sys-112 tem is defined by x, y, and z in Earth radii (R_e) . The variable rdist denotes the satel-113 lite's radial distance from Earth. The magnetic field line type, termed as "Foot Type", 114 indicates the connectivity of the IMF field lines to Earth, calculated using the Tsyganenko 115 (1995) model. The initial definition stated by Kronberg et al. (2020) is as follows: the 116 interplanetary magnetic field lines (IMF) with no connection to Earth have Foot Type 117 0, open magnetic field lines with one connection to Earth have Foot Type 1, and closed 118 field lines with both ends connected to Earth have Foot Type 2. It was, however, decided 119 to redefine the IMF to 1 and open field lines to 0, to achieve a stronger linear relation-120 ship between Foot Type and the target variable, as discussed in chapter 2.3.2. 121

The Disturbance storm time index (Dst_index) characterizes geomagnetic storms in the unit nT (Banerjee et al., 2012). The Auroral Electrojet (AE_index) quantifies magnetic activity in the auroral zone, also denoted in nT. The 10.7 cm solar radio flux (F10.7) with unit sfu serves as a solar activity level indicator and a proxy for solar emissions (Tapping, 2013).

The IMF direction is described by its components BimfxGSE, BimfyGSE, and BimfzGSE in the GSE system in nT. The IMF direction at the magnetopause determines if reconnection happens on the dayside (Crooker et al., 1979; Luhmann et al., 1984; Koga et al., 2019). Plasma properties of the solar wind are described by Solar wind speed (VSW) in km/s, proton density (NpSW) in cm⁻³, and temperature (Temp) in K. The direction of the solar wind velocity is described by VxSW_GSE, VySW_GSE, and VzSW_GSE. The solar wind dynamic pressure, Pdyn (nPa), can be represented as:

$$Pdyn = NpSW * VSW^{2} * 1.67 * 10^{-6}$$
(1)

134

135

2.3 Exploratory Data Analysis

2.3.1 Spatial Proton Intensity Distribution

To analyze the proton intensity distribution in relation to the spacecraft's trajec-136 tory, we combined the y and z coordinates to produce a radial distance variable, termed 137 138 yz_axis. We introduced this variable because Cluster's trajectory is predominantly in the southern hemisphere, contrasted with SMILE's expected northern trajectory. The 139 yz_axis is computed by $\sqrt{y^2 + z^2}$, offering a simplified yet informative perspective on 140 proton intensity's spatial distribution. Figure 2 depicts the spatial distribution of pro-141 ton intensities in the x, $\sqrt{y^2+z^2}$ coordinate system. The color gradient represents the 142 percentage of measurements that exceed 2, the mean value of \log_{10} (proton intensities) 143 rounded to one significant digit, highlighting regions with prolonged high proton inten-144 sities. The central black void indicates missing measurements. This gap arises from our 145 deliberate exclusion of data points with radial distances (rdist) below 6 Earth radii (R_e) 146 in order to emphasize regions beyond the radiation belts. Historically, proton intensi-147 ties surge in zones under 6 R_e , which encompass the ring current and radiation belt re-148 gions. Our focus shifts to lesser-studied areas, with their generally lower intensities out-149 side of the radiation belts, mainly because the SXI telescope on the SMILE mission is 150 equipped with a shutter mechanism, protecting its Charge Coupled Devices (CCD) from 151 intense radiation within the radiation belt. The decline in proton intensities with increased 152 distance from Earth is observed, irrespective of whether it's along the x axis or the y-153 z plane. This observation aligns with subsequent feature plots and correlation matrix 154 analyses. Moreover, this analysis reveals areas along closed magnetic field lines with height-155 ened proton intensities, as well as sparser regions corresponding to open magnetic field 156 lines over the polar cap, demonstrating a clear spatial correlation between magnetic field 157 line configuration and proton intensity distribution. 158



Figure 2. (Left) Heatmap of proton intensities against the x-coordinate and $\sqrt{y^2 + z^2}$. The color gradient represents the percentage of measurements where $\log_{10}(\text{proton intensity}) > 2$, which is the mean value rounded to one significant digit. (Right) Data point density for each bin, with the blue line representing the magnetopause, derived using Shue et al. (1997). Both plots incorporate the 462,615 data points post-reshaping for the SMILE mission.

2.3.2 Cross-Correlation and Feature Plot Analysis

159 160 161

Cross-correlation matrices, employing the Pearson coefficient, are used in feature selection for linear regression models. The coefficient quantifies the linear relationship

strength and direction between two variables, spanning from -1 (perfect negative relationship) to 1 (perfect positive relationship), with 0 indicating no linear correlation. Such analyses illuminate potential multicollinearity issues in datasets, which can adversely affect regression coefficient stability and model interpretability (Raschka et al., 2022; James et al., 2013).

- Analyzing the cross-correlation matrix (Figure 3), we observed:
- FootType: The feature plot in Figure 4 highlighted a clear potential for refining 168 the correlation between FootType and p3. The initial positive correlation of 0.24169 with p3 was improved to 0.41 upon redefining the foot type as mentioned in sec-170 tion 2.2. 171 • AE_index: While the correlation was weak (0.07), the feature plot identified AE 172 values surpassing 2600 nT as possible outliers. 173 • F10.7 solar radio flux index: Given its correlation coefficient of -0.20, the feature 174 plot shows a predominantly monotonically decreasing relationship between F10.7 175 and the target variable. 176 VxSW_GSE: The feature plot demonstrated that proton intensity increases with higher 177 178
- absolute wind speeds up to 950 km/s. Values exceeding this were considered as potential outliers.
 Distance variables: While z showed a positive correlation of 0.21, its relationship
 - Distance variables: while 2 showed a positive correlation of 0.21, its relationship with proton intensities displayed a clear maximum around 0 on the feature plot. This insight led to the introduction of |z| as an improved predictor.
 - Other Variables: The strong negative correlations of rdist and yz_axis with p3 were supported by the feature plot's linear regression lines, emphasizing their importance as predictors.

Analysis of the proton intensity histogram identified two extreme outliers exceeding 100,000 1/cm²/s/sr/keV (see Figure A1). Removing these and other above-identified outliers from different predictors did not improve model performance, justifying their retention. Further analysis revealed 606 F10.7 measurements above 900, deemed unrealistic and consequently removed.

¹⁹¹ 2.4 Data Split and Data Scaling

167

181

182

183

184

185

This section outlines the additional processing steps applied to the data set, already reshaped for the SMILE mission as detailed in section 2.1. These steps include splitting the data into training and testing sets, transforming the target variable, and scaling features.

Records before December 31, 2012, were allocated to the training set, while records from January 1, 2013, onwards formed the test set. This temporal division results in an approximate 75% to 25% split between the training and test datasets. The training set, was then later further divided into training and validation sets by the use of five-fold crossvalidation, where the dataset is divided into five parts, with each part being used as a validation set while the remaining four parts are used as training data

The proton intensities recorded in channel 3 (p3), our target values, display a wide spectrum. We therefore transformed these values using a base 10 logarithmic function. Addressing the challenge of logging zero values, all such occurrences in p3 were substituted with 0.5. However, this introduces potential pitfalls as we expect an artificial population with the same values, a concern later revisited during model evaluation (Bellégo et al., 2021).

Optimizing gradient descent requires careful attention to feature scaling (Raschka et al., 2022). In our polynomial regression model, we employed a double-scaling tech-



Figure 3. Pearson coefficient-based correlation matrix for the predictors and the proton intensity post-data reshaping for the SMILE mission.

nique to ensure numerical stability and facilitate the optimization process. Initially, the
original features were scaled to zero mean and unit variance using the StandardScaler()
method, aligning with the desired outcomes (Pedregosa et al., 2011). Subsequently, polynomial features were generated from these scaled features. To further enhance the model's
robustness, these polynomial features were subjected to a second round of scaling using
the same StandardScaler() method.

By scaling both the original and polynomial features, we ensure that the coefficients are directly comparable in terms of their contribution to the model and that all features display a mean and unit variance of zero.

²¹⁹ **3** Methodology

220

3.1 Linear Regression Model

The choice of employing linear regression models in this study is underpinned by several reasons. First and foremost, linear regression models offer a simple and interpretable framework for understanding how input variables affect the output. Furthermore, the methodology allows for the transformation of input variables to enhance their predictive capabilities, such as the introduction of polynomial terms and interaction effects (Hastie et al., 2001).



Figure 4. Mean of the logarithmically scaled proton intensities from the p3 channel against potential predictors. Vertical lines depict the standard 95% confidence level, while horizontal lines indicate bin half-widths. Linear regression lines in red are shown for rdist, yz_axis, and F10.7.

The Ordinary Least Squares (OLS) model serves as the foundational approach, fo-227 cusing on minimizing the sum of squared differences between observed and predicted val-228 ues (James et al., 2013; Galton, 1886). To tackle the possible issue of multicollinearity, 229 Ridge Regression can be utilized, which incorporates an L2 penalty term into the loss 230 function (Kutner et al., 2005; Hoerl & Kennard, 1970). Lasso Regression is employed when 231 feature selection is essential, as it uses an L1 penalty to drive certain coefficients to zero, 232 effectively eliminating them from the model (Santosa & Symes, 1986). Lastly, Elastic 233 Net Regression can be used to combine the strengths of both L1 and L2 penalties, pro-234 viding a balanced approach that can handle both multicollinearity and feature selection 235 (Pedregosa et al., 2011). Multiple models were trained using the scikit-learn library in 236 Python (Pedregosa et al., 2011). 237

3.2 Model Selection and Optimization

For model evaluation, we utilized a set of metrics, including Mean Squared Error 239 (MSE), Mean Absolute Error (MAE), R^2 (coefficient of determination), Pearson corre-240 lation, and Spearman correlation. Model selection was primarily guided by the perfor-241 mance of R^2 and Spearman correlation on the validation set. To ensure a robust and gen-242 eralizable evaluation, five-fold cross-validation with the help of the KFold function from 243 sklearn.model_selection was applied to the training set. Given the time-series nature 244 of our dataset, the shuffle parameter within the cross-validation procedure was inten-245 tionally set to **false**. Subsequently, the evaluation metrics were computed as the aver-246 age values derived from the five cross-validation folds, thereby offering a more reliable 247 measure of the model's true performance. 248

249

238

3.2.1 Simple OLS, Lasso, Ridge and Elastic Net

Following the initial selection of linear regression models, two distinct approaches were undertaken to optimize model performance. The first approach involved the application of various linear regression techniques, including OLS, Lasso, Ridge, and Elastic Net. This approach, however, did not yield satisfactory results. The maximum R^2 value on the validation set was only 0.02, and the highest Spearman correlation coefficient was 0.43.

256

3.2.2 Introduction of Polynomial Terms

To improve upon this, the second approach incorporated polynomial terms into a standard Lasso model from sklearn.linear. The model was optimized for the regularization parameter α using five-fold cross-validation. The optimal α was determined using LassoCV with a maximum iteration of 10,000 and a tolerance of 1×10^{-5} . This approach significantly improved the model performance, achieving an R^2 value of 0.22 and a Spearman correlation coefficient of 0.51 on the validation set.

263

3.2.3 Heuristic-based Feature Selection Technique

However, this model included 52 predictors, making it complex and potentially prone 264 to overfitting. Further work was needed to develop a more parsimonious model with a 265 maximum of 25 predictors while maintaining acceptable performance. To reduce the num-266 ber of predictors while maintaining model performance, we adopted a heuristic-based fea-267 ture selection strategy. For this strategy the Stochastic Gradient Descent (SGD) framework was employed, with the algorithm configured as follows: the regularization term 269 (α) was set to the optimal value identified through cross-validation. The learning rate 270 was set to a constant value, initialized at $\eta_0 = 1 \times 10^{-5}$. The hyperparameter defin-271 ing the loss function was set to the squared error loss. An L1 penalty term was incor-272

porated for feature selection. The algorithm was set to terminate when the tolerance reached 273 1×10^{-5} , with a maximum of 100 iterations for convergence.



Figure 5. Plot of Average Mean Squared Error (MSE) against the regularization parameter α . The curve exhibits an "elbow" point at 13 predictors, indicating a minimal but acceptable loss in model performance. A noticeable increase in MSE is observed when the number of predictors is reduced from 13 to 12, suggesting that all 13 predictors left, display significant importance for the model. This "elbow" point serves as the basis for selecting an optimal α value and, consequently, the number of predictors for the final model.

Unlike earlier approaches that solely aimed to minimize the Mean Squared Error 275 (MSE), this method also considers the number of predictors in the final model. We tested 276 a range of regularization parameters (α) and sought to identify a "knee" or "elbow" in 277 the plot of MSE versus α . This point represents a compromise between model perfor-278 mance and complexity. 279

To enhance the robustness of the feature selection process, we employed K-Fold cross-280 validation with the shuffle parameter set to True. This approach allows for a more rep-281 resentative sampling of the training data across each fold. Specifically, we aimed to iden-282 tify the most stable set of predictors corresponding to the "elbow" point for the regu-283 larization parameter α . By enabling shuffling during cross-validation, we increase the like-284 lihood that the predictor set extracted from one of the folds offers a more comprehen-285 sive representation of the entire training dataset. The α range chosen was from 0.03 to 286 0.17, which covered all models with the amount of non-zero predictors ranging from 28 287 to 4. 288

Upon employing this approach, we identified a subset of 13 predictors by analyz-289 ing the MSE vs α plot in Figure 5. Importantly, we operate under the assumption that 290 all predictors remaining after the feature selection process are relevant to the outcome. 291 Therefore, penalizing these predictors, as Lasso does, could introduce an unwanted bias 292 into the model. Given this consideration, an OLS model was chosen for the final train-203 ing rather than a Lasso regression model. 294

Model	MSE	MAE	R^2	Pearson	Spearman	Predictors	N_train
Basic_OLS	0.92	0.76	0.02	0.43	0.43	9	353,660
Poly_Lasso	0.73	0.68	0.22	0.51	0.51	52	353,660
Heuristic_Poly_OLS	0.78	0.71	0.17	0.47	0.47	13	353,660
Split_Poly_Part1	0.58	0.60	0.24	0.53	0.53	5	60,961
Split_Poly_Part2	0.82	0.73	0.09	0.50	0.50	19	145,952
Split_Poly_Part3	0.79	0.72	0.19	0.46	0.46	15	70,137
Split_Poly_Part4	0.65	0.64	0.29	0.55	0.55	6	76,610

Table 1. Average performance metrics for different models resulting from five-fold cross-validation and the number of data points N₋train used for training.

Although the resulting model exhibits lower performance on the validation set, as evidenced by Table 1, it better aligns with the study's objectives of interpretability and usability compared to the Lasso model with 52 predictors. This heuristic-based feature selection strategy aligns well with the principle of Occam's razor, suggesting that simpler models are preferable when performance is comparable. Therefore, this approach effectively strikes a balance between the number of predictors and model performance, thereby enhancing the model's interpretability and practical utility.

302 3.2.4 Data Split

An in-depth analysis of the relationship between the y and p3 variables revealed that the data could be divided into four distinct parts, each characterized by an increasing or decreasing slope, see Figure 4 (d). This lead to the decision to split the dataset into four separate parts based on specific conditions, as described below:

- Part 1: $y \leq -6.6 R_e$
 - Part 2: -6.6 $R_e \le y \le 2.3 R_e$
 - Part 3: 2.3 $R_e \le y \le 6 R_e$
 - Part 4: $y \ge 6 R_e$

³¹¹ Upon splitting the data, separate models were built for each part, using the same ³¹² heuristic-based predictor selection technique previously described. These outperformed ³¹³ the non-split OLS model in the Spearman correlation coefficient and the R^2 metric for ³¹⁴ three out of the four subsets (see Table 1), all while maintaining low model complexity.

315 4 Results

308

309

310

This chapter presents the empirical results obtained from the evaluation of various OLS models on the unseen test set. The models are compared based on a set of evaluation metrics and feature importances.

4.1 Presentation of Final Models

In this section, we present the final forms of our linear regression models developed for predicting proton intensities. Each model is displayed with its coefficients in basic, unscaled units, offering a clear view of the relative impact of each predictor variable. These models encapsulate our findings and are ready for practical application.



354

355

4.2 Performance on the Test Set

In the Heuristic_Poly_OLS model, the hexagonal bins largely align with the ideal 356 fit line (see Figure 6), which is indicative of good predictive performance. However, this 357 model exhibits a tendency to underestimate observed values, notably at higher proton 358 intensities. A significant peak at $\log(0.5)$ in the histogram of observed values is associ-359 ated with an overestimation in the Heuristic_Poly_OLS model's predictions. This peak 360 stems from the substitution of zero values in the target variable p3 before applying the 361 logarithmic transformation. This overestimation at $\log(0.5)$ potentially skews the model's 362 363 learning process, causing it to adjust its predictions downward to minimize the overall loss. While this adjustment mitigates the error for overestimated values, it concurrently 364 introduces a bias leading to the underestimation of other observed values. This behav-365

(6)



Figure 6. Jointplots comparing observed and predicted values of proton intensities, of the test set for the different OLS models. The red lines represent ideal fits where observed values equal predicted values. Color bars indicate the number of samples in each hexagonal bin. Histograms at the top and right margins show the distributions of observed and predicted values for each model.

Model	MSE	MAE	R^2	Pearson	Spearman	Predictors
Heuristic_Poly_OLS	0.74	0.71	0.22	0.56	0.57	13
Split_Poly_Part1	0.46	0.54	0.38	0.62	0.61	5
Split_Poly_Part2	0.83	0.74	0.11	0.56	0.57	19
Split_Poly_Part3	0.77	0.73	0.24	0.61	0.62	15
Split_Poly_Part4	0.47	0.56	0.50	0.72	0.72	6

Table 2. Performance metrics of the final models on test data.

ior is consistently observable across all models, particularly those focusing on specific regions. As an alternative to zero substitution, we also explored the removal of these zero
values. While this approach enhanced performance on the training set, it consistently
led to diminished performance on the validation set. Consequently, despite its limitations, the zero substitution technique was retained to ensure better generalization to unseen data.

Turning our attention to Split_Poly_Part1 and Split_Poly_Part4, these models exhibit the most well-centered distribution around the ideal fit line in their respective heatmaps. This observation is consistent with their performance metrics as recorded in Table 2, showcasing R^2 values of 0.38 and 0.50 and Spearman coefficients of 0.61 and 0.72 for the test set. Notably, these models also maintain low complexity, employing only 5 and 6 predictors, respectively.

Conversely, Split_Poly_Part2, with its high complexity due to having 19 predictors, 378 exhibits subpar performance despite an acceptable Spearman coefficient. Significantly, 379 with an R^2 value of only 0.11, this model is the sole split variant that exhibits notably 380 inferior performance compared to the unsplit Heuristic_Poly_OLS model in the test set. 381 The more inhomogeneous distribution of hexagonal bins in its heatmap is indicative of 382 this weaker performance. On the other hand, Split_Poly_Part3 shows a modest improve-383 ment over the unsplit model. This is evident not only in the performance metrics but 384 also in a more concentrated distribution in its heatmap, compared to Split_Poly_Part2. 385

4.3 Feature Importance

386

In order to derive feature importance in a linear regression model, one can exam-387 ine the coefficients of the model. The magnitude of the coefficients indicates the rela-388 tive importance of the corresponding feature in predicting the target variable. A larger 389 absolute value of a coefficient suggests a stronger influence of the associated feature on 390 the outcome. The features are scaled appropriately by scaling the features once before 391 the creation of the polynomials and once afterward. Scaling ensures that all features are 392 on a comparable scale, which prevents features with larger values from dominating those 393 with smaller values in the model. The feature importance for each model was plotted 394 in figure 7. 395

The variations in feature importance across the different models offer insights into 396 the underlying mechanisms affecting proton intensities in various regions. For the Heuris-397 tic_Poly_OLS model and models corresponding to the inner regions (Split_Poly_Part2 398 and Split_Poly_Part3), the absolute value of $z(|\mathbf{z}|)$ emerges as the most significant pre-399 dictor, next to Foottype and VxSW_GSE. The Split_Poly_-Part2 model additionally iden-400 tifies the polynomial terms |z| rdist and rdist² as significant features. The similar-401 ity between the models for the inner part and the model trained on the full data is most 402 likely partially influenced by the fact that the inner regions contain 61% of the total data 403 points. The models tailored to the outer regions (Split_Poly_Part1 and Split_Poly_Part4) 404







(b) Split_Poly_Part1



(c) Split_Poly_Part2



Figure 7. Feature importance plots for five different OLS models: Each plot presents the absolute values of the model coefficients, serving as indicators of feature importance. Accompanying error bars represent the standard errors, providing a measure of the coefficient's reliability. The plots collectively offer insights into the relative significance of each predictor across different models.

prioritize rdist, VxSW_GSE, and Foottype as their top predictors, in that order. A no table distinction from the inner region models is the elevated significance of VxSW_GSE.

407 5 Discussion

The most critical predictor for our Heuristic_Poly_OLS_Model, which utilizes the 408 full dataset, is the absolute value of z, denoted as |z|. The model reveals a negative cor-409 relation between |z| and the proton intensities, indicating that as |z| increases, the pro-410 ton intensity declines. This trend can be primarily attributed to the circulation of pro-411 tons in Earth's magnetic field. Most ions are concentrated at the equatorial plane dur-412 ing their drift trajectories on the closed magnetic field lines. At higher latitudes where 413 open magnetic field lines dominate, the proton intensities are expected to drop with |z|414 distance. Consequently, the proton intensities reduce with an increase in |z|. 415

The predictor FootType categorizes magnetic field line types and ranks as the sec-416 ond most influential factor. Closed field lines, known for the highest proton intensities, 417 trap charged particle populations. The importance of this parameter aligns with the stud-418 ies by Walsh et al. (2014) and Kronberg et al. (2020). In contrast, open field line regions 419 typically correlate with lower particle energies outside the soft proton (SP) range, result-420 ing in weaker count rates as detailed in (Kronberg et al., 2020). IMF regions show slightly 421 higher count rates since particles can experience acceleration in the bow shock region, 422 especially quasi-parallel bow shock configurations (normal to the shock is parallel to the 423 IMF direction) (Blandford & Ostriker, 1978; Kronberg et al., 2009; Sundberg et al., 2016). 424

The high importance of solar wind speed in the X-direction is consistent with the 425 analysis of the feature plot in figure 4. Kronberg, Hannan, et al. (2021) also found that 426 VxSW_GSE displays "the most substantial linear dependence of the proton intensities among 427 the OMNI parameters." The solar wind speed, directly correlated to its electric field as 428 $E = V_x \times B_z$, is crucial for magnetospheric dynamics as it determines the rate of mag-429 netic reconnection on Earth's dayside (Dorelli, 2019), and consequently magnetic recon-430 nection at the night side. A surge in solar wind speed correlates with an increased rate 431 of magnetic reconnection. Additionally, magnetic reconnection events, which can accel-432 erate charged particles, also impact soft proton intensities significantly, as noted by (Read 433 & Ponman, 2003). Research by Gonzalez et al. (1994), Milan et al. (2012)), and Wang 434 et al. (2014) further elucidates this concept, indicating that a variety of solar wind-magnetosphere 435 energy transfer models are dependent on the velocity of the solar wind. 436

437 6 Conclusion

In this study, we developed five user-friendly linear regression models to predict proton intensities in the energy range of 92.2 keV to 159.7 keV with a Spearmen correlation ranging from 0.57 to 0.72 on the test data. Utilizing data from the Cluster's RAPID experiment, supplemented with solar, solar wind, and geomagnetic data from the OMNI database, the study focused on aligning the models with the anticipated spatial area covered by the upcoming SMILE-mission.

Segmenting the data into four distinct regions based on the y coordinate with thresholds -6.6 R_E , 2.3 R_E and 6 R_E , resulted in enhanced model performance for three of the four segments, surpassing the main model's performance. The primary predictors in these outer regions were identified as radial distance and the radial solar wind speed. Conversely, the inner region models and the comprehensive main model demonstrated a significant dependence on the absolute value of z and the type of magnetic field lines.

The redefinition of the FootType variable and the incorporation of the absolute value of z as key model features significantly improved the model compared to previous relevant studies. This study suggests that the development of more accurate predictive mod-

- els for space weather phenomena may not solely rely on novel algorithms, but also on
- 454 crafting tailored models, each addressing distinct regions with their specific character-
- 455 istics.

456 Appendix A : Histogram of Proton Intensities



Figure A1. Histogram of the proton intensities measured by channel 3.

457 Appendix B Open Research

The authors express their gratitude to the team at the Cluster Science Archive (https:// 458 csa.esac.esa.int) for supplying the data. Additionally, we recognize the utilization 459 of the OMNIWeb service and OMNI data from NASA/GSFC's Space Physics Data Fa-460 cility (King & Papitashvili, 2005). The code and dataset used to derive the linear regres-461 sion model can be found via the following link: https://zenodo.org/records/10964236 462 ?token=eyJhbGciOiJIUzUxMiJ9.eyJpZCI6IjkxMGQzY2EzLTMxNDAtNDZmNi05MjE5LWUzZmM1Y2I20 463 WM5MSIsImRhdGEiOnt9LCJyYW5kb2OiOiIyMjQ1NjQzYTdlZjk3YzNjODA1MzdlZGJhMmQyMzg2MyJ9 464 .oKgcCJTFE6KbvqqjXNh3wzfneL3XeY6meWb-XhbcKum0x0ztugSGFtvCaLblb3WAW0E5ccrkVEWZDnZ9 465 vEs6EQ. 466

467 Acknowledgments

The database used in this study was generated within the team led by Fabio Gastaldello 468 on "Soft Protons in the Magnetosphere focused by X-ray Telescopes" at the International 469 Space Science Institute in Bern, Switzerland. We are particularly grateful to Dr. Gastaldello 470 for his invaluable feedback on our work, which greatly enhanced this research. EK and 471 SM are supported by the German Research Foundation (DFG) under number KR 4375/2-472 1 within SPP "Dynamic Earth". EK is also supported by the DFG under number KR 473 4375/4-1. We are grateful to Dr. Andrew Read and Dr. Steven Sembay for their insight-474 ful suggestions which significantly enhanced the representation of spatial proton inten-475 sity distribution in our plots. 476

477 References

478	Banerjee, A., Bej, A., & Chatterjee, T. N. (2012). On the existence of a long	range
479	correlation in the geomagnetic disturbance storm time (dst) index.	Astro-
480	physics and Space Science, 337, 23–32. doi: 10.1007/s10509-011-0836-1	

Bellégo, C., Benatia, D., & Pape, L. (2021). Dealing with logs and zeros in regression models. CREST - Serie des Documents de Travail. doi: 10.2139/ssrn
.3444996

484	Blandford, R. D., & Ostriker, J. P. (1978, April). Particle acceleration by astrophys-
485	ical shocks. Astrophys. J., 221, L29-L32. doi: 10.1086/182658
486	Branduardi-Raymont, G., & Wang, C. (2022). The smile mission. In C. Bambi
487	& A. Santangelo (Eds.), Handbook of x-ray and gamma-ray astrophysics (pp.
488	1–22). Singapore: Springer Nature Singapore. doi: 10.1007/978-981-16-4544-0
489	_39-1
490	Branduardi-Raymont, G., Wang, C., Escoubet, C. P., Adamovic, M., Agnolon, D.,
491	Berthomier, M., Zhu, Z. (2018). Smile definition study report (ESA/SCI
492	No. 1). European Space Agency. doi: 10.5270/esa.smile.definition_study_report
493	-2018-12
494	Crooker, N. U., Eastman, T. E., & Stiles, G. S. (1979). Observations of plasma
495	depletion in the magnetosheath at the dayside magnetopause. Journal
496	of Geophysical Research: Space Physics, 84(A3), 869-874. doi: 10 .1029 /
497	JA084iA03p00869
498	Daly, P. W., & Kronberg, E. A. (2023). User guide to the rapid measurements
499	in the cluster science archive (csa) (User Guide No. CAA-EST-UG-RAP
500	6.1). Max Planck Institute for Solar System Research. Retrieved 2023-11-05,
501	from https://www2 .mps .mpg .de/dokumente/projekte/cluster/rapid/
502	Rapid_Userguide.pdf
503	Dorelli, J. C. (2019). Does the solar wind electric field control the reconnection rate
504	at earth's subsolar magnetopause? Journal of Geophysical Research: Space
505	Physics, 124(4), 2668-2681. doi: $10.1029/2018$ JA025868
506	Fioretti, V., Bulgarelli, A., Malaguti, G., Spiga, D., & Tiengo, A. (2016). Monte
507	carlo simulations of soft proton flares: testing the physics with xmm-newton.
508	In JW. A. den Herder, T. Takahashi, & M. Bautz (Eds.), Space telescopes
509	and instrumentation 2016: Ultraviolet to gamma ray (Vol. 9905, p. 99056W).
510	SPIE. doi: 10.1117/12.2232537
511	Galton, F. (1886). Regression towards mediocrity in hereditary stature. The Journal
512	of the Anthropological Institute of Great Britain and Ireland, 15, 246–263. Re-
513	trieved 2023-10-07, from http://www.jstor.org/stable/2841583 doi: 10
514	.2307/2841583
515	Gonzalez, W. D., Joselyn, J. A., Kamide, Y., Kroehl, H. W., Rostoker, G., Tsu-
516	rutani, B. T., & Vasyliunas, V. M. (1994). What is a geomagnetic storm?
517	Journal of Geophysical Research: Space Physics, 99(A4), 5771-5792. doi:
518	10.1029/93JA02867
519	Hastie, T., Friedman, J., & Tibshirani, R. (2001). Linear methods for regres-
520	sion. In The elements of statistical learning: Data mining, inference,
521	and prediction (pp. 41–78). New York, NY: Springer New York. doi:
522	$10.1007/978-0-387-21606-5_3$
523	Hoerl, A. E., & Kennard, R. W. (1970). Ridge regression: Biased estimation
524	for nonorthogonal problems. <i>Technometrics</i> , 12(1), 55-67. Retrieved from
525	https://www.tandionline.com/dol/abs/10.1080/00401/06.1970.10488634
526	$\frac{d01: 10.1080/00401700.1970.10488034}{0.10401700.1970.10488034}$
527	Hubbard, M. W. J., Buggey, I. W., Hall, D., Feldman, C., Keelana, J., Hetnering-
528	ton, O., Holland, A. (2024). Techniques for estimating radiation damage
529	from particles passage and focusing from micro pore optics. Journal of Astro-
530	nomicul Telescopes, Instruments, and Systems. (III press)
531	James, G., Witten, D., Hastie, I., & Hosmirani, R. (2015). Linear regression. In
532	NV doi: 10.1007/078.1.4614.7138.7
533	NI. U0I. $10.1001/970-1-4014-7100-7$ King I H & Dapitashuili N E (2005) Color wind spatial scales in and
534	isons of hourly wind and acc plasma and magnetic field data. Journal of Cas
535	nousical Research: Snace Physics 110(A2) doi: 10.1020/2004IA010640
530	Kora D. Conzeloz W. D. Souze V. M. Cardoso F. P. Wang C. & Lin 7 K.
53/	(2019) Davside magnetonause reconnection: Its dependence on solar wind and
000	(2010). Dayside magnetopause reconnection, its dependence on solar will all

539	magnetosheath conditions. Journal of Geophysical Research: Space Physics,
540	124(11), 8778-8787. doi: $10.1029/2019$ JA026889
541	Kronberg, E. A., Clerc, N., Cros, A., de Plaa, J., Gastaldello, F., Gu, L., Valen-
542	tini, N. (2020). Prediction and Understanding of Soft-proton Contamination
543	in XMM-Newton: A Machine Learning Approach. The Astrophysical Journal,
544	903(2), 89. doi: $10.3847/1538-4357/abbb8f$
545	Kronberg, E. A., Daly, P. W., Grigorenko, E. E., Smirnov, A. G., Klecker, B., &
546	Malykhin, A. Y. (2021). Energetic charged particles in the terrestrial magneto-
547	sphere: Cluster/rapid results. Journal of Geophysical Research: Space Physics,
548	126(9), e2021JA029273. doi: 10.1029/2021JA029273
549	Kronberg, E. A., Hannan, T., Huthmacher, J., Münzer, M., Peste, F., Zhou, Z.,
550	Ilie, R. (2021). Prediction of soft proton intensities in the near-earth
551	space using machine learning. The Astrophysical Journal, $921(1)$, 76. doi:
552	10.3847/1538-4357/ac1b30
553	Kronberg, E. A., Kis, A., Klecker, B., Daly, P. W., & Lucek, E. A. (2009). Mul-
554	tipoint observations of ions in the $30-160$ kev energy range upstream of the
555	earth's bow shock. Journal of Geophysical Research: Space Physics, 114(A3).
556	doi: 10.1029/2008JA013754
557	Kronberg, E. A., Rashev, M. V., Daly, P. W., Shprits, Y. Y., Turner, D. L., Droz-
558	dov, A., Friedel, R. (2016). Contamination in electron observations of the
559	silicon detector on board cluster/rapid/ies instrument in earth's radiation belts
560	and ring current. Space Weather, 14, 449-462. doi: 10.1002/2016SW001369
561	Kutner, M. H., Nachtsheim, C. J., Neter, J., & Li, W. (2005). Applied linear statisti-
562	cal models (Fifth ed.). New York: McGraw-Hill Irwin.
563	Luhmann, J. G., Walker, R. J., Russell, C. T., Crooker, N. U., Spreiter, J. R., &
564	Stahara, S. S. (1984). Patterns of potential magnetic field merging sites on
565	the dayside magnetopause. Journal of Geophysical Research: Space Physics,
566	89(A3), 1739-1742. doi: $10.1029/JA089iA03p01739$
567	Milan, S. E., Gosling, J. S., & Hubert, B. (2012). Relationship between interplan-
568	etary parameters and the magnetopause reconnection rate quantified from
569	observations of the expanding polar cap. Journal of Geophysical Research: S_{mago} Devoice $117((\Lambda_2))$ doi: 10.1020/20111A.017082
570	Deducer a Revision P. Cristian C. Cremfort A. Michel V. Thiring P. Cristian O.
571	Duchosnay F (2011) Scikit learn: Machine learning in Python Lowrad of
572	Machine Learning Research 19, 2825–2830
573	Derivati E Freybarg M Voung M C H Dommranz C Hoga B Diabold S
574	Sentengelo A (2024) Using sra/erosita to estimate soft proton flures at
575	the athena detectors
570	Baab W Branduardi-Baymont G Wang C Dai L Donovan E Enno G
579	Zheng J (2016) Smile: a joint esa/cas mission to investigate the inter-
570	action between the solar wind and earth's magnetosphere In I-W A den
580	Herder T Takahashi & M Bautz (Eds.) Space telescopes and instrumen-
581	tation 2016: Ultraviolet to gamma ray (Vol. 9905, p. 990502). SPIE. doi:
582	10.1117/12.2231984
583	Raschka, S., Liu, Y., & Mirialili, V. (2022). Predicting continuous target variables
584	with regression analysis. In Machine learning with pytorch and scikit-learn :
585	develop machine learning and deep learning models with python (p. 269-304).
586	Packt Publishing.
587	Read, A. M., & Ponman, T. J. (2003, October). The xmm-newton epic background:
588	Production of background maps and event files. Astronomy & Astrophysics,
589	409(1), 395–410. doi: 10.1051/0004-6361:20031099
590	Santosa, F., & Symes, W. W. (1986). Linear inversion of band-limited reflection seis-
501	
591	mograms. SIAM Journal on Scientific and Statistical Computing, 7(4), 1307-
591	mograms. SIAM Journal on Scientific and Statistical Computing, 7(4), 1307- 1330. doi: 10.1137/0907087

594	Singer, H. J. (1997). A new functional form to study the solar wind control
595	of the magnetopause size and shape. Journal of Geophysical Research: Space
596	<i>Physics</i> , 102(A5), 9497-9511. doi: 10.1029/97JA00196
597	Sundberg, T., Haynes, C. T., Burgess, D., & Mazelle, C. X. (2016). Ion acceleration
598	at the quasi-parallel bow shock: Decoding the signature of injection. The As-
599	trophysical Journal, 820(1), 21. doi: 10.3847/0004-637X/820/1/21
600	Tapping, K. F. (2013). The 10.7 cm solar radio flux (f10.7). Space Weather, 11(7),
601	394-406. doi: 10.1002/swe.20064
602	Tsyganenko, N. A. (1995). Modeling the earth's magnetospheric magnetic field con-
603	fined within a realistic magnetopause. Journal of Geophysical Research: Space
604	Physics, 100(A4), 5599-5612. doi: $10.1029/94JA03193$
605	Walsh, B. M., Kuntz, K. D., Collier, M. R., Sibeck, D. G., Snowden, S. L., &
606	Thomas, N. E. (2014). Energetic particle impact on x-ray imaging with
607	xmm-newton. Space Weather, 12(6), 387-394. doi: 10.1002/2014SW001046
608	Wang, C., Han, J. P., Li, H., Peng, Z., & Richardson, J. D. (2014). Solar wind-
609	magnetosphere energy coupling function fitting: Results from a global mhd
610	simulation. Journal of Geophysical Research: Space Physics, 119(8), 6199-
611	6212. doi: 10.1002/2014JA019834
612	Wilken, B., Axford, W. I., Daglis, I., Daly, P., Güttler, W., Ip, W. H., Ullaland,
613	S. (1997). Rapid. In C. P. Escoubet, C. T. Russell, & R. Schmidt (Eds.), The
614	cluster and phoenix missions (pp. 399–473). Dordrecht: Springer Netherlands.

doi: 10.1007/978-94-011-5666-0_14

615

Evaluating Proton Intensities for the SMILE Mission 1

2

3

Simon Mischel¹, Elena A. Kronberg¹, C.P. Escoubet²

¹Department of Earth and Environmental Sciences (Geophysics), Ludwig Maximilian University of Munich (LMU) Munich, Theresienstr. 41, Munich, D-80333, Germany ²European Space Research and Technology Centre, Noordwjik, Keplerlaan 1, 2201 AZ, The Netherlands 4 5

6	Key Points:
7	• Developed models predict proton intensities impacting satellites, aiding space weather
8	forecasting and mission planning.
9	• Different regions in space showcase distinct relations between proton intensities
10	and predicting parameters.
11	• Study findings highlight the importance of tailored approaches in space weather
12	prediction.

 $Corresponding \ author: \ Simon \ Mischel, \ \texttt{simonmischel@hotmail.com}$

13 Abstract

This study introduces five linear regression models developed to accurately predict pro-14 ton intensities in the critical energy range of 92.2 keV to 159.7 keV. To achieve this task 15 we utilized 14 years of data sourced from the Cluster's RAPID experiment and NASA's 16 OMNI database. This data was then aligned with the Solar wind-Magnetosphere-Ionosphere 17 Link Explorer (SMILE) mission's trajectory, to increase model accuracy in the relevant 18 regions. Our approach diverges from existing methodologies by offering a user-friendly 19 model that doesn't require specialized software, making it accessible for broader appli-20 cations in satellite mission planning and risk assessment. The research segregates the dataset 21 into four distinct regions, each analyzed for proton intensity dynamics. In the outer re-22 gions ($|\text{YGSE}| \ge 6 R_e$) there is a pronounced dependence on radial distance and solar 23 wind speed. In contrast, the inner regions ($|YGSE| \leq 6R_e$) demonstrate a significant 24 dependence of proton intensities on the absolute value of the z-coordinate and the mag-25 netic field line topology. Our models achieved a Spearman correlation ranging from 0.5726 to 0.72 on the test set, indicating good predictive capabilities. The findings emphasize 27 the role of regional characteristics in space weather prediction and underscore the po-28 tential for tailored approaches in future research. 29

³⁰ Plain Language Summary

We developed a new model to predict space weather, specifically focusing on pro-31 32 ton intensities, which can impact how well satellites work in space. We used 14 years of space observations to create five easy-to-use numerical models. These models are designed 33 to help with planning and protecting future satellite missions, such as the upcoming SMILE 34 mission, from space weather effects. In our study, we looked closely at different areas in 35 space around Earth. We found that in the outer areas ($|YGSE| \ge 6 R_e$), the distance 36 from Earth and the speed of the solar wind are important for understanding proton be-37 havior. However, in areas ($|YGSE| \leq 6R_e$), the height above Earth (measured along 38 the z-direction) and the type of magnetic field lines play a more significant role. This 39 shows us that different areas in space around Earth can be affected by space weather in 40 different ways. Our models did a good job of predicting these effects, showing that choos-41 ing a tailored approach can be useful when forecasting proton intensities. 42

43 **1** Introduction

Space weather events, driven by solar activities pose significant challenges to satel-44 lite operations and measurements. Notable examples include the European Space Agency's 45 (ESA) Cluster and X-ray Multi-Mirror (XMM-Newton) missions. The Cluster mission, 46 particularly its Research with Adaptive Particle Imaging Detector (RAPID)/Imaging 47 Electron Spectrometer (IES), has encountered challenges due to high proton intensities 48 leading to measurement contamination (Wilken et al., 1997; Kronberg et al., 2016; Kro-49 nberg, Daly, et al., 2021). Similarly, the X-ray telescope aboard the XMM-Newton space-50 craft experienced significant operational disruptions, with approximately 40% of its ob-51 servation time compromised due to background contamination (Walsh et al., 2014). Fur-52 thermore, investigations into the XMM-Newton telescope's susceptibility to soft protons 53 highlight proton intensities in the sub-100 to 300 keV range, particularly around 100 keV, 54 as the most damaging, leading to significant operational challenges and data contam-55 ination, see (Fioretti et al., 2016) and references therein. 56

The upcoming European-Chinese Solar wind-Magnetosphere-Ionosphere Link Explorer (SMILE) mission (Branduardi-Raymont et al., 2018), slated for launch in 2025, aspires to deepen our understanding of the Sun-Earth interaction, decoding space weather hazards and understanding energy entry into Earth's magnetosphere. While missions like ATHENA (Advanced Telescope for High Energy Astrophysics), which will maintain an orbit around the L2 Lagrange point with a continuously large distance to Earth, are ex-

pected to face minimal threats from soft protons (Perinati et al., 2024), SMILE with its 63 highly inclined elliptical orbit around earth, akin to that of Cluster and XMM-Newton, 64 will navigate through diverse magnetospheric regions, making it susceptible to soft pro-65 ton radiation. Particularly its Soft X-ray Imager (SXI) telescope, presents challenges con-66 cerning radiation exposure (Raab et al., 2016; Branduardi-Raymont & Wang, 2022). There-67 fore a critical component of achieving SMILE's objectives is the accurate prediction of 68 proton radiation levels, which significantly affect the Total Ionizing Dose (TID) and To-69 tal Non-Ionizing Dose (TNID) absorbed by the Charge-Coupled Devices (CCDs) of the 70 Soft X-ray Imager (SXI) (Hubbard et al., 2024). In recent discussions (M. Hubbard, per-71 sonal communication, 2023) the necessity for models that accurately estimate radiation 72 levels was emphasized, particularly in critical energy ranges below 300 keV, crucial for 73 the SMILE mission's success. A strong preference for models that are not only accurate 74 but also straightforward and interpretable was expressed. 75

To meet these needs, our study adopts a distinct approach compared to existing 76 research, which is based on machine learning black box models, such as the works of Kronberg 77 et al. (2020) and Kronberg, Hannan, et al. (2021). We aim to develop a simple, user-friendly 78 linear regression model leveraging data from the Cluster mission and NASA's OMNI database 79 (King & Papitashvili, 2005). Our model's simplicity and ease of use make it accessible 80 to a broader range of users, not requiring specialized software or extensive computational 81 resources. This approach not only contributes to the scientific understanding of space 82 weather phenomena but also offers practical tools for satellite mission planning and risk 83 assessment. 84

⁸⁵ 2 Data Analysis and Processing

86 87

2.1 Data Preparation: Adapting CLUSTER's Dataset for the SMILE Mission's Trajectory

The proton intensity data for this research was taken from the Cluster's RAPID 88 experiment, ranging from 2001 to 2015. The experiment captures 3-D energetic electron 89 and ion fluxes above approximately 30 keV using the Imaging Electron Spectrometer (IES) 90 and the Imaging Ion Mass Spectrometer (IIMS) instruments. Situated in the SCENIC 91 detector head, the IIMS instrument identifies ion energies and species. The methodol-92 ogy involves using start and stop signals produced from electrons emitted by an initial 93 thin foil on the solid-state detector's surface. The time-of-flight (TOF) between these 94 signals, combined with the known energy, discerns the species and energy channel (Daly 95 & Kronberg, 2023). This study specifically uses data from the p3 channel, targeting pro-96 ton intensities between 92.2 keV and 159.7 keV. As predicting parameters for solar, so-97 lar wind and geomagnetic activity we used variables from the OMNI database(https:// 98 omniweb.gsfc.nasa.gov/), see also King and Papitashvili (2005). 99

Aiming for a model tailored to the SMILE mission's trajectory (see Figure 1), data 100 filtering was imperative. Points not adhering to the following spatial parameters were 101 excluded: $-10.5R_e \le x \le 11.2R_e$; $-10.8R_e \le y \le 11.5R_e$; $z \le 18.5R_e$; $\sqrt{x^2 + y^2} \le 10.5R_e$; $\sqrt{x^2 + y^2} \ge 10$ 102 $11.6R_e; \sqrt{x^2+z^2} \leq 19.8R_e; \sqrt{z^2+y^2} \leq 20.0R_e$. These constraints were chosen by 103 rounding the maxima and minima of the spatial parameters to the nearest tenth. As a 104 result, we were left with a trimmed dataset, reduced from 1,172,923 to 462,615 data points. 105 Though compact, this dataset is centered on the space region significant for SMILE, promis-106 ing heightened model accuracy. It's noteworthy that negative z-values weren't excluded, 107 considering the SMILE mission's highly inclined, elliptical orbit, which dips to $-3.5 R_e$. 108 Omitting these would disregard vital data, especially since the Cluster's trajectory spent 109 a notable amount of time in the southern hemisphere. 110



Figure 1. SMILE mission's trajectory. Distinct colors represent individual years with Earth, having a radius of 6,731 km, at the center.

2.2 Predictor Introduction

111

The spacecraft's location in the Geocentric Solar Ecliptic (GSE) coordinate sys-112 tem is defined by x, y, and z in Earth radii (R_e) . The variable rdist denotes the satel-113 lite's radial distance from Earth. The magnetic field line type, termed as "Foot Type", 114 indicates the connectivity of the IMF field lines to Earth, calculated using the Tsyganenko 115 (1995) model. The initial definition stated by Kronberg et al. (2020) is as follows: the 116 interplanetary magnetic field lines (IMF) with no connection to Earth have Foot Type 117 0, open magnetic field lines with one connection to Earth have Foot Type 1, and closed 118 field lines with both ends connected to Earth have Foot Type 2. It was, however, decided 119 to redefine the IMF to 1 and open field lines to 0, to achieve a stronger linear relation-120 ship between Foot Type and the target variable, as discussed in chapter 2.3.2. 121

The Disturbance storm time index (Dst_index) characterizes geomagnetic storms in the unit nT (Banerjee et al., 2012). The Auroral Electrojet (AE_index) quantifies magnetic activity in the auroral zone, also denoted in nT. The 10.7 cm solar radio flux (F10.7) with unit sfu serves as a solar activity level indicator and a proxy for solar emissions (Tapping, 2013).

The IMF direction is described by its components BimfxGSE, BimfyGSE, and BimfzGSE in the GSE system in nT. The IMF direction at the magnetopause determines if reconnection happens on the dayside (Crooker et al., 1979; Luhmann et al., 1984; Koga et al., 2019). Plasma properties of the solar wind are described by Solar wind speed (VSW) in km/s, proton density (NpSW) in cm⁻³, and temperature (Temp) in K. The direction of the solar wind velocity is described by VxSW_GSE, VySW_GSE, and VzSW_GSE. The solar wind dynamic pressure, Pdyn (nPa), can be represented as:

$$Pdyn = NpSW * VSW^{2} * 1.67 * 10^{-6}$$
(1)

134

135

2.3 Exploratory Data Analysis

2.3.1 Spatial Proton Intensity Distribution

To analyze the proton intensity distribution in relation to the spacecraft's trajec-136 tory, we combined the y and z coordinates to produce a radial distance variable, termed 137 138 yz_axis. We introduced this variable because Cluster's trajectory is predominantly in the southern hemisphere, contrasted with SMILE's expected northern trajectory. The 139 yz_axis is computed by $\sqrt{y^2 + z^2}$, offering a simplified yet informative perspective on 140 proton intensity's spatial distribution. Figure 2 depicts the spatial distribution of pro-141 ton intensities in the x, $\sqrt{y^2+z^2}$ coordinate system. The color gradient represents the 142 percentage of measurements that exceed 2, the mean value of \log_{10} (proton intensities) 143 rounded to one significant digit, highlighting regions with prolonged high proton inten-144 sities. The central black void indicates missing measurements. This gap arises from our 145 deliberate exclusion of data points with radial distances (rdist) below 6 Earth radii (R_e) 146 in order to emphasize regions beyond the radiation belts. Historically, proton intensi-147 ties surge in zones under 6 R_e , which encompass the ring current and radiation belt re-148 gions. Our focus shifts to lesser-studied areas, with their generally lower intensities out-149 side of the radiation belts, mainly because the SXI telescope on the SMILE mission is 150 equipped with a shutter mechanism, protecting its Charge Coupled Devices (CCD) from 151 intense radiation within the radiation belt. The decline in proton intensities with increased 152 distance from Earth is observed, irrespective of whether it's along the x axis or the y-153 z plane. This observation aligns with subsequent feature plots and correlation matrix 154 analyses. Moreover, this analysis reveals areas along closed magnetic field lines with height-155 ened proton intensities, as well as sparser regions corresponding to open magnetic field 156 lines over the polar cap, demonstrating a clear spatial correlation between magnetic field 157 line configuration and proton intensity distribution. 158



Figure 2. (Left) Heatmap of proton intensities against the x-coordinate and $\sqrt{y^2 + z^2}$. The color gradient represents the percentage of measurements where $\log_{10}(\text{proton intensity}) > 2$, which is the mean value rounded to one significant digit. (Right) Data point density for each bin, with the blue line representing the magnetopause, derived using Shue et al. (1997). Both plots incorporate the 462,615 data points post-reshaping for the SMILE mission.

2.3.2 Cross-Correlation and Feature Plot Analysis

159 160 161

Cross-correlation matrices, employing the Pearson coefficient, are used in feature selection for linear regression models. The coefficient quantifies the linear relationship

strength and direction between two variables, spanning from -1 (perfect negative relationship) to 1 (perfect positive relationship), with 0 indicating no linear correlation. Such analyses illuminate potential multicollinearity issues in datasets, which can adversely affect regression coefficient stability and model interpretability (Raschka et al., 2022; James et al., 2013).

- Analyzing the cross-correlation matrix (Figure 3), we observed:
- FootType: The feature plot in Figure 4 highlighted a clear potential for refining 168 the correlation between FootType and p3. The initial positive correlation of 0.24169 with p3 was improved to 0.41 upon redefining the foot type as mentioned in sec-170 tion 2.2. 171 • AE_index: While the correlation was weak (0.07), the feature plot identified AE 172 values surpassing 2600 nT as possible outliers. 173 • F10.7 solar radio flux index: Given its correlation coefficient of -0.20, the feature 174 plot shows a predominantly monotonically decreasing relationship between F10.7 175 and the target variable. 176 VxSW_GSE: The feature plot demonstrated that proton intensity increases with higher 177 178
- absolute wind speeds up to 950 km/s. Values exceeding this were considered as potential outliers.
 Distance variables: While z showed a positive correlation of 0.21, its relationship
 - Distance variables: while 2 showed a positive correlation of 0.21, its relationship with proton intensities displayed a clear maximum around 0 on the feature plot. This insight led to the introduction of |z| as an improved predictor.
 - Other Variables: The strong negative correlations of rdist and yz_axis with p3 were supported by the feature plot's linear regression lines, emphasizing their importance as predictors.

Analysis of the proton intensity histogram identified two extreme outliers exceeding 100,000 1/cm²/s/sr/keV (see Figure A1). Removing these and other above-identified outliers from different predictors did not improve model performance, justifying their retention. Further analysis revealed 606 F10.7 measurements above 900, deemed unrealistic and consequently removed.

¹⁹¹ 2.4 Data Split and Data Scaling

167

181

182

183

184

185

This section outlines the additional processing steps applied to the data set, already reshaped for the SMILE mission as detailed in section 2.1. These steps include splitting the data into training and testing sets, transforming the target variable, and scaling features.

Records before December 31, 2012, were allocated to the training set, while records from January 1, 2013, onwards formed the test set. This temporal division results in an approximate 75% to 25% split between the training and test datasets. The training set, was then later further divided into training and validation sets by the use of five-fold crossvalidation, where the dataset is divided into five parts, with each part being used as a validation set while the remaining four parts are used as training data

The proton intensities recorded in channel 3 (p3), our target values, display a wide spectrum. We therefore transformed these values using a base 10 logarithmic function. Addressing the challenge of logging zero values, all such occurrences in p3 were substituted with 0.5. However, this introduces potential pitfalls as we expect an artificial population with the same values, a concern later revisited during model evaluation (Bellégo et al., 2021).

Optimizing gradient descent requires careful attention to feature scaling (Raschka et al., 2022). In our polynomial regression model, we employed a double-scaling tech-



Figure 3. Pearson coefficient-based correlation matrix for the predictors and the proton intensity post-data reshaping for the SMILE mission.

nique to ensure numerical stability and facilitate the optimization process. Initially, the
original features were scaled to zero mean and unit variance using the StandardScaler()
method, aligning with the desired outcomes (Pedregosa et al., 2011). Subsequently, polynomial features were generated from these scaled features. To further enhance the model's
robustness, these polynomial features were subjected to a second round of scaling using
the same StandardScaler() method.

By scaling both the original and polynomial features, we ensure that the coefficients are directly comparable in terms of their contribution to the model and that all features display a mean and unit variance of zero.

²¹⁹ **3** Methodology

220

3.1 Linear Regression Model

The choice of employing linear regression models in this study is underpinned by several reasons. First and foremost, linear regression models offer a simple and interpretable framework for understanding how input variables affect the output. Furthermore, the methodology allows for the transformation of input variables to enhance their predictive capabilities, such as the introduction of polynomial terms and interaction effects (Hastie et al., 2001).



Figure 4. Mean of the logarithmically scaled proton intensities from the p3 channel against potential predictors. Vertical lines depict the standard 95% confidence level, while horizontal lines indicate bin half-widths. Linear regression lines in red are shown for rdist, yz_axis, and F10.7.

The Ordinary Least Squares (OLS) model serves as the foundational approach, fo-227 cusing on minimizing the sum of squared differences between observed and predicted val-228 ues (James et al., 2013; Galton, 1886). To tackle the possible issue of multicollinearity, 229 Ridge Regression can be utilized, which incorporates an L2 penalty term into the loss 230 function (Kutner et al., 2005; Hoerl & Kennard, 1970). Lasso Regression is employed when 231 feature selection is essential, as it uses an L1 penalty to drive certain coefficients to zero, 232 effectively eliminating them from the model (Santosa & Symes, 1986). Lastly, Elastic 233 Net Regression can be used to combine the strengths of both L1 and L2 penalties, pro-234 viding a balanced approach that can handle both multicollinearity and feature selection 235 (Pedregosa et al., 2011). Multiple models were trained using the scikit-learn library in 236 Python (Pedregosa et al., 2011). 237

3.2 Model Selection and Optimization

For model evaluation, we utilized a set of metrics, including Mean Squared Error 239 (MSE), Mean Absolute Error (MAE), R^2 (coefficient of determination), Pearson corre-240 lation, and Spearman correlation. Model selection was primarily guided by the perfor-241 mance of R^2 and Spearman correlation on the validation set. To ensure a robust and gen-242 eralizable evaluation, five-fold cross-validation with the help of the KFold function from 243 sklearn.model_selection was applied to the training set. Given the time-series nature 244 of our dataset, the shuffle parameter within the cross-validation procedure was inten-245 tionally set to **false**. Subsequently, the evaluation metrics were computed as the aver-246 age values derived from the five cross-validation folds, thereby offering a more reliable 247 measure of the model's true performance. 248

249

238

3.2.1 Simple OLS, Lasso, Ridge and Elastic Net

Following the initial selection of linear regression models, two distinct approaches were undertaken to optimize model performance. The first approach involved the application of various linear regression techniques, including OLS, Lasso, Ridge, and Elastic Net. This approach, however, did not yield satisfactory results. The maximum R^2 value on the validation set was only 0.02, and the highest Spearman correlation coefficient was 0.43.

256

3.2.2 Introduction of Polynomial Terms

To improve upon this, the second approach incorporated polynomial terms into a standard Lasso model from sklearn.linear. The model was optimized for the regularization parameter α using five-fold cross-validation. The optimal α was determined using LassoCV with a maximum iteration of 10,000 and a tolerance of 1×10^{-5} . This approach significantly improved the model performance, achieving an R^2 value of 0.22 and a Spearman correlation coefficient of 0.51 on the validation set.

263

3.2.3 Heuristic-based Feature Selection Technique

However, this model included 52 predictors, making it complex and potentially prone 264 to overfitting. Further work was needed to develop a more parsimonious model with a 265 maximum of 25 predictors while maintaining acceptable performance. To reduce the num-266 ber of predictors while maintaining model performance, we adopted a heuristic-based fea-267 ture selection strategy. For this strategy the Stochastic Gradient Descent (SGD) framework was employed, with the algorithm configured as follows: the regularization term 269 (α) was set to the optimal value identified through cross-validation. The learning rate 270 was set to a constant value, initialized at $\eta_0 = 1 \times 10^{-5}$. The hyperparameter defin-271 ing the loss function was set to the squared error loss. An L1 penalty term was incor-272

porated for feature selection. The algorithm was set to terminate when the tolerance reached 273 1×10^{-5} , with a maximum of 100 iterations for convergence.



Figure 5. Plot of Average Mean Squared Error (MSE) against the regularization parameter α . The curve exhibits an "elbow" point at 13 predictors, indicating a minimal but acceptable loss in model performance. A noticeable increase in MSE is observed when the number of predictors is reduced from 13 to 12, suggesting that all 13 predictors left, display significant importance for the model. This "elbow" point serves as the basis for selecting an optimal α value and, consequently, the number of predictors for the final model.

Unlike earlier approaches that solely aimed to minimize the Mean Squared Error 275 (MSE), this method also considers the number of predictors in the final model. We tested 276 a range of regularization parameters (α) and sought to identify a "knee" or "elbow" in 277 the plot of MSE versus α . This point represents a compromise between model perfor-278 mance and complexity. 279

To enhance the robustness of the feature selection process, we employed K-Fold cross-280 validation with the shuffle parameter set to True. This approach allows for a more rep-281 resentative sampling of the training data across each fold. Specifically, we aimed to iden-282 tify the most stable set of predictors corresponding to the "elbow" point for the regu-283 larization parameter α . By enabling shuffling during cross-validation, we increase the like-284 lihood that the predictor set extracted from one of the folds offers a more comprehen-285 sive representation of the entire training dataset. The α range chosen was from 0.03 to 286 0.17, which covered all models with the amount of non-zero predictors ranging from 28 287 to 4. 288

Upon employing this approach, we identified a subset of 13 predictors by analyz-289 ing the MSE vs α plot in Figure 5. Importantly, we operate under the assumption that 290 all predictors remaining after the feature selection process are relevant to the outcome. 291 Therefore, penalizing these predictors, as Lasso does, could introduce an unwanted bias 292 into the model. Given this consideration, an OLS model was chosen for the final train-203 ing rather than a Lasso regression model. 294

Model	MSE	MAE	R^2	Pearson	Spearman	Predictors	N_train
Basic_OLS	0.92	0.76	0.02	0.43	0.43	9	353,660
Poly_Lasso	0.73	0.68	0.22	0.51	0.51	52	353,660
Heuristic_Poly_OLS	0.78	0.71	0.17	0.47	0.47	13	353,660
Split_Poly_Part1	0.58	0.60	0.24	0.53	0.53	5	60,961
Split_Poly_Part2	0.82	0.73	0.09	0.50	0.50	19	145,952
Split_Poly_Part3	0.79	0.72	0.19	0.46	0.46	15	70,137
Split_Poly_Part4	0.65	0.64	0.29	0.55	0.55	6	76,610

Table 1. Average performance metrics for different models resulting from five-fold cross-validation and the number of data points N₋train used for training.

Although the resulting model exhibits lower performance on the validation set, as evidenced by Table 1, it better aligns with the study's objectives of interpretability and usability compared to the Lasso model with 52 predictors. This heuristic-based feature selection strategy aligns well with the principle of Occam's razor, suggesting that simpler models are preferable when performance is comparable. Therefore, this approach effectively strikes a balance between the number of predictors and model performance, thereby enhancing the model's interpretability and practical utility.

302 3.2.4 Data Split

An in-depth analysis of the relationship between the y and p3 variables revealed that the data could be divided into four distinct parts, each characterized by an increasing or decreasing slope, see Figure 4 (d). This lead to the decision to split the dataset into four separate parts based on specific conditions, as described below:

- Part 1: $y \leq -6.6 R_e$
 - Part 2: -6.6 $R_e \le y \le 2.3 R_e$
 - Part 3: 2.3 $R_e \le y \le 6 R_e$
 - Part 4: $y \ge 6 R_e$

³¹¹ Upon splitting the data, separate models were built for each part, using the same ³¹² heuristic-based predictor selection technique previously described. These outperformed ³¹³ the non-split OLS model in the Spearman correlation coefficient and the R^2 metric for ³¹⁴ three out of the four subsets (see Table 1), all while maintaining low model complexity.

315 4 Results

308

309

310

This chapter presents the empirical results obtained from the evaluation of various OLS models on the unseen test set. The models are compared based on a set of evaluation metrics and feature importances.

4.1 Presentation of Final Models

In this section, we present the final forms of our linear regression models developed for predicting proton intensities. Each model is displayed with its coefficients in basic, unscaled units, offering a clear view of the relative impact of each predictor variable. These models encapsulate our findings and are ready for practical application.



354

355

4.2 Performance on the Test Set

In the Heuristic_Poly_OLS model, the hexagonal bins largely align with the ideal 356 fit line (see Figure 6), which is indicative of good predictive performance. However, this 357 model exhibits a tendency to underestimate observed values, notably at higher proton 358 intensities. A significant peak at $\log(0.5)$ in the histogram of observed values is associ-359 ated with an overestimation in the Heuristic_Poly_OLS model's predictions. This peak 360 stems from the substitution of zero values in the target variable p3 before applying the 361 logarithmic transformation. This overestimation at $\log(0.5)$ potentially skews the model's 362 363 learning process, causing it to adjust its predictions downward to minimize the overall loss. While this adjustment mitigates the error for overestimated values, it concurrently 364 introduces a bias leading to the underestimation of other observed values. This behav-365

(6)



Figure 6. Jointplots comparing observed and predicted values of proton intensities, of the test set for the different OLS models. The red lines represent ideal fits where observed values equal predicted values. Color bars indicate the number of samples in each hexagonal bin. Histograms at the top and right margins show the distributions of observed and predicted values for each model.

Model	MSE	MAE	R^2	Pearson	Spearman	Predictors
Heuristic_Poly_OLS	0.74	0.71	0.22	0.56	0.57	13
Split_Poly_Part1	0.46	0.54	0.38	0.62	0.61	5
Split_Poly_Part2	0.83	0.74	0.11	0.56	0.57	19
Split_Poly_Part3	0.77	0.73	0.24	0.61	0.62	15
Split_Poly_Part4	0.47	0.56	0.50	0.72	0.72	6

Table 2. Performance metrics of the final models on test data.

ior is consistently observable across all models, particularly those focusing on specific regions. As an alternative to zero substitution, we also explored the removal of these zero
values. While this approach enhanced performance on the training set, it consistently
led to diminished performance on the validation set. Consequently, despite its limitations, the zero substitution technique was retained to ensure better generalization to unseen data.

Turning our attention to Split_Poly_Part1 and Split_Poly_Part4, these models exhibit the most well-centered distribution around the ideal fit line in their respective heatmaps. This observation is consistent with their performance metrics as recorded in Table 2, showcasing R^2 values of 0.38 and 0.50 and Spearman coefficients of 0.61 and 0.72 for the test set. Notably, these models also maintain low complexity, employing only 5 and 6 predictors, respectively.

Conversely, Split_Poly_Part2, with its high complexity due to having 19 predictors, 378 exhibits subpar performance despite an acceptable Spearman coefficient. Significantly, 379 with an R^2 value of only 0.11, this model is the sole split variant that exhibits notably 380 inferior performance compared to the unsplit Heuristic_Poly_OLS model in the test set. 381 The more inhomogeneous distribution of hexagonal bins in its heatmap is indicative of 382 this weaker performance. On the other hand, Split_Poly_Part3 shows a modest improve-383 ment over the unsplit model. This is evident not only in the performance metrics but 384 also in a more concentrated distribution in its heatmap, compared to Split_Poly_Part2. 385

4.3 Feature Importance

386

In order to derive feature importance in a linear regression model, one can exam-387 ine the coefficients of the model. The magnitude of the coefficients indicates the rela-388 tive importance of the corresponding feature in predicting the target variable. A larger 389 absolute value of a coefficient suggests a stronger influence of the associated feature on 390 the outcome. The features are scaled appropriately by scaling the features once before 391 the creation of the polynomials and once afterward. Scaling ensures that all features are 392 on a comparable scale, which prevents features with larger values from dominating those 393 with smaller values in the model. The feature importance for each model was plotted 394 in figure 7. 395

The variations in feature importance across the different models offer insights into 396 the underlying mechanisms affecting proton intensities in various regions. For the Heuris-397 tic_Poly_OLS model and models corresponding to the inner regions (Split_Poly_Part2 398 and Split_Poly_Part3), the absolute value of $z(|\mathbf{z}|)$ emerges as the most significant pre-399 dictor, next to Foottype and VxSW_GSE. The Split_Poly_-Part2 model additionally iden-400 tifies the polynomial terms |z| rdist and rdist² as significant features. The similar-401 ity between the models for the inner part and the model trained on the full data is most 402 likely partially influenced by the fact that the inner regions contain 61% of the total data 403 points. The models tailored to the outer regions (Split_Poly_Part1 and Split_Poly_Part4) 404







(b) Split_Poly_Part1



(c) Split_Poly_Part2



Figure 7. Feature importance plots for five different OLS models: Each plot presents the absolute values of the model coefficients, serving as indicators of feature importance. Accompanying error bars represent the standard errors, providing a measure of the coefficient's reliability. The plots collectively offer insights into the relative significance of each predictor across different models.

prioritize rdist, VxSW_GSE, and Foottype as their top predictors, in that order. A no table distinction from the inner region models is the elevated significance of VxSW_GSE.

407 5 Discussion

The most critical predictor for our Heuristic_Poly_OLS_Model, which utilizes the 408 full dataset, is the absolute value of z, denoted as |z|. The model reveals a negative cor-409 relation between |z| and the proton intensities, indicating that as |z| increases, the pro-410 ton intensity declines. This trend can be primarily attributed to the circulation of pro-411 tons in Earth's magnetic field. Most ions are concentrated at the equatorial plane dur-412 ing their drift trajectories on the closed magnetic field lines. At higher latitudes where 413 open magnetic field lines dominate, the proton intensities are expected to drop with |z|414 distance. Consequently, the proton intensities reduce with an increase in |z|. 415

The predictor FootType categorizes magnetic field line types and ranks as the sec-416 ond most influential factor. Closed field lines, known for the highest proton intensities, 417 trap charged particle populations. The importance of this parameter aligns with the stud-418 ies by Walsh et al. (2014) and Kronberg et al. (2020). In contrast, open field line regions 419 typically correlate with lower particle energies outside the soft proton (SP) range, result-420 ing in weaker count rates as detailed in (Kronberg et al., 2020). IMF regions show slightly 421 higher count rates since particles can experience acceleration in the bow shock region, 422 especially quasi-parallel bow shock configurations (normal to the shock is parallel to the 423 IMF direction) (Blandford & Ostriker, 1978; Kronberg et al., 2009; Sundberg et al., 2016). 424

The high importance of solar wind speed in the X-direction is consistent with the 425 analysis of the feature plot in figure 4. Kronberg, Hannan, et al. (2021) also found that 426 VxSW_GSE displays "the most substantial linear dependence of the proton intensities among 427 the OMNI parameters." The solar wind speed, directly correlated to its electric field as 428 $E = V_x \times B_z$, is crucial for magnetospheric dynamics as it determines the rate of mag-429 netic reconnection on Earth's dayside (Dorelli, 2019), and consequently magnetic recon-430 nection at the night side. A surge in solar wind speed correlates with an increased rate 431 of magnetic reconnection. Additionally, magnetic reconnection events, which can accel-432 erate charged particles, also impact soft proton intensities significantly, as noted by (Read 433 & Ponman, 2003). Research by Gonzalez et al. (1994), Milan et al. (2012)), and Wang 434 et al. (2014) further elucidates this concept, indicating that a variety of solar wind-magnetosphere 435 energy transfer models are dependent on the velocity of the solar wind. 436

437 6 Conclusion

In this study, we developed five user-friendly linear regression models to predict proton intensities in the energy range of 92.2 keV to 159.7 keV with a Spearmen correlation ranging from 0.57 to 0.72 on the test data. Utilizing data from the Cluster's RAPID experiment, supplemented with solar, solar wind, and geomagnetic data from the OMNI database, the study focused on aligning the models with the anticipated spatial area covered by the upcoming SMILE-mission.

Segmenting the data into four distinct regions based on the y coordinate with thresholds -6.6 R_E , 2.3 R_E and 6 R_E , resulted in enhanced model performance for three of the four segments, surpassing the main model's performance. The primary predictors in these outer regions were identified as radial distance and the radial solar wind speed. Conversely, the inner region models and the comprehensive main model demonstrated a significant dependence on the absolute value of z and the type of magnetic field lines.

The redefinition of the FootType variable and the incorporation of the absolute value of z as key model features significantly improved the model compared to previous relevant studies. This study suggests that the development of more accurate predictive mod-

- els for space weather phenomena may not solely rely on novel algorithms, but also on
- 454 crafting tailored models, each addressing distinct regions with their specific character-
- 455 istics.

456 Appendix A : Histogram of Proton Intensities



Figure A1. Histogram of the proton intensities measured by channel 3.

457 Appendix B Open Research

The authors express their gratitude to the team at the Cluster Science Archive (https:// 458 csa.esac.esa.int) for supplying the data. Additionally, we recognize the utilization 459 of the OMNIWeb service and OMNI data from NASA/GSFC's Space Physics Data Fa-460 cility (King & Papitashvili, 2005). The code and dataset used to derive the linear regres-461 sion model can be found via the following link: https://zenodo.org/records/10964236 462 ?token=eyJhbGciOiJIUzUxMiJ9.eyJpZCI6IjkxMGQzY2EzLTMxNDAtNDZmNi05MjE5LWUzZmM1Y2I20 463 WM5MSIsImRhdGEiOnt9LCJyYW5kb2OiOiIyMjQ1NjQzYTdlZjk3YzNjODA1MzdlZGJhMmQyMzg2MyJ9 464 .oKgcCJTFE6KbvqqjXNh3wzfneL3XeY6meWb-XhbcKum0x0ztugSGFtvCaLblb3WAW0E5ccrkVEWZDnZ9 465 vEs6EQ. 466

467 Acknowledgments

The database used in this study was generated within the team led by Fabio Gastaldello 468 on "Soft Protons in the Magnetosphere focused by X-ray Telescopes" at the International 469 Space Science Institute in Bern, Switzerland. We are particularly grateful to Dr. Gastaldello 470 for his invaluable feedback on our work, which greatly enhanced this research. EK and 471 SM are supported by the German Research Foundation (DFG) under number KR 4375/2-472 1 within SPP "Dynamic Earth". EK is also supported by the DFG under number KR 473 4375/4-1. We are grateful to Dr. Andrew Read and Dr. Steven Sembay for their insight-474 ful suggestions which significantly enhanced the representation of spatial proton inten-475 sity distribution in our plots. 476

477 References

478	Banerjee, A., Bej, A., & Chatterjee, T. N. (2012). On the existence of a long	range
479	correlation in the geomagnetic disturbance storm time (dst) index.	Astro-
480	physics and Space Science, 337, 23–32. doi: 10.1007/s10509-011-0836-1	

Bellégo, C., Benatia, D., & Pape, L. (2021). Dealing with logs and zeros in regression models. CREST - Serie des Documents de Travail. doi: 10.2139/ssrn
.3444996

484	Blandford, R. D., & Ostriker, J. P. (1978, April). Particle acceleration by astrophys-
485	ical shocks. Astrophys. J., 221, L29-L32. doi: 10.1086/182658
486	Branduardi-Raymont, G., & Wang, C. (2022). The smile mission. In C. Bambi
487	& A. Santangelo (Eds.), Handbook of x-ray and gamma-ray astrophysics (pp.
488	1–22). Singapore: Springer Nature Singapore. doi: 10.1007/978-981-16-4544-0
489	_39-1
490	Branduardi-Raymont, G., Wang, C., Escoubet, C. P., Adamovic, M., Agnolon, D.,
491	Berthomier, M., Zhu, Z. (2018). Smile definition study report (ESA/SCI
492	No. 1). European Space Agency. doi: 10.5270/esa.smile.definition_study_report
493	-2018-12
494	Crooker, N. U., Eastman, T. E., & Stiles, G. S. (1979). Observations of plasma
495	depletion in the magnetosheath at the dayside magnetopause. Journal
496	of Geophysical Research: Space Physics, 84(A3), 869-874. doi: 10 .1029 /
497	JA084iA03p00869
498	Daly, P. W., & Kronberg, E. A. (2023). User guide to the rapid measurements
499	in the cluster science archive (csa) (User Guide No. CAA-EST-UG-RAP
500	6.1). Max Planck Institute for Solar System Research. Retrieved 2023-11-05,
501	from https://www2 .mps .mpg .de/dokumente/projekte/cluster/rapid/
502	Rapid_Userguide.pdf
503	Dorelli, J. C. (2019). Does the solar wind electric field control the reconnection rate
504	at earth's subsolar magnetopause? Journal of Geophysical Research: Space
505	Physics, 124(4), 2668-2681. doi: $10.1029/2018$ JA025868
506	Fioretti, V., Bulgarelli, A., Malaguti, G., Spiga, D., & Tiengo, A. (2016). Monte
507	carlo simulations of soft proton flares: testing the physics with xmm-newton.
508	In JW. A. den Herder, T. Takahashi, & M. Bautz (Eds.), Space telescopes
509	and instrumentation 2016: Ultraviolet to gamma ray (Vol. 9905, p. 99056W).
510	SPIE. doi: 10.1117/12.2232537
511	Galton, F. (1886). Regression towards mediocrity in hereditary stature. The Journal
512	of the Anthropological Institute of Great Britain and Ireland, 15, 246–263. Re-
513	trieved 2023-10-07, from http://www.jstor.org/stable/2841583 doi: 10
514	.2307/2841583
515	Gonzalez, W. D., Joselyn, J. A., Kamide, Y., Kroehl, H. W., Rostoker, G., Tsu-
516	rutani, B. T., & Vasyliunas, V. M. (1994). What is a geomagnetic storm?
517	Journal of Geophysical Research: Space Physics, 99(A4), 5771-5792. doi:
518	10.1029/93JA02867
519	Hastie, T., Friedman, J., & Tibshirani, R. (2001). Linear methods for regres-
520	sion. In The elements of statistical learning: Data mining, inference,
521	and prediction (pp. 41–78). New York, NY: Springer New York. doi:
522	$10.1007/978-0-387-21606-5_3$
523	Hoerl, A. E., & Kennard, R. W. (1970). Ridge regression: Biased estimation
524	for nonorthogonal problems. <i>Technometrics</i> , 12(1), 55-67. Retrieved from
525	https://www.tandionline.com/dol/abs/10.1080/00401/06.1970.10488634
526	$\frac{d01: 10.1080/00401700.1970.10488034}{0.10401700.1970.10488034}$
527	Hubbard, M. W. J., Buggey, I. W., Hall, D., Feldman, C., Keelana, J., Hetnering-
528	ton, O., Holland, A. (2024). Techniques for estimating radiation damage
529	from particles passage and focusing from micro pore optics. Journal of Astro-
530	nomicul Telescopes, Instruments, and Systems. (III press)
531	James, G., Witten, D., Hastie, I., & Hosmirani, R. (2015). Linear regression. In
532	NV doi: 10.1007/078.1.4614.7138.7
533	NI. U0I. $10.1001/970^{-1-4014-7100-7}$ King I H & Dapitashuili N E (2005) Color wind spatial scales in and
534	isons of hourly wind and acc plasma and magnetic field data. Journal of Cas
535	nousical Research: Snace Physics 110(A2) doi: 10.1020/2004IA010640
530	Kora D. Conzaloz W. D. Souza V. M. Cardoso F. P. Wang C. & Lin 7 V.
53/	(2019) Davside magnetonause reconnection: Its dependence on solar wind and
000	(2010). Dayside magnetopause reconnection, its dependence on solar will all

539	magnetosheath conditions. Journal of Geophysical Research: Space Physics,
540	124(11), 8778-8787. doi: $10.1029/2019$ JA026889
541	Kronberg, E. A., Clerc, N., Cros, A., de Plaa, J., Gastaldello, F., Gu, L., Valen-
542	tini, N. (2020). Prediction and Understanding of Soft-proton Contamination
543	in XMM-Newton: A Machine Learning Approach. The Astrophysical Journal,
544	903(2), 89. doi: $10.3847/1538-4357/abbb8f$
545	Kronberg, E. A., Daly, P. W., Grigorenko, E. E., Smirnov, A. G., Klecker, B., &
546	Malykhin, A. Y. (2021). Energetic charged particles in the terrestrial magneto-
547	sphere: Cluster/rapid results. Journal of Geophysical Research: Space Physics,
548	126(9), e2021JA029273. doi: 10.1029/2021JA029273
549	Kronberg, E. A., Hannan, T., Huthmacher, J., Münzer, M., Peste, F., Zhou, Z.,
550	Ilie, R. (2021). Prediction of soft proton intensities in the near-earth
551	space using machine learning. The Astrophysical Journal, 921(1), 76. doi:
552	10.3847/1538-4357/ac1b30
553	Kronberg, E. A., Kis, A., Klecker, B., Dalv, P. W., & Lucek, E. A. (2009). Mul-
554	tipoint observations of ions in the 30–160 kev energy range upstream of the
555	earth's bow shock. Journal of Geophysical Research: Space Physics, 114(A3).
556	doi: 10.1029/2008JA013754
557	Kronberg, E. A., Rashev, M. V., Daly, P. W., Shprits, Y. Y., Turner, D. L., Droz-
558	dov, A., Friedel, R. (2016). Contamination in electron observations of the
559	silicon detector on board cluster/rapid/ies instrument in earth's radiation belts
560	and ring current. Space Weather, 14, 449-462. doi: 10.1002/2016SW001369
561	Kutner, M. H., Nachtsheim, C. J., Neter, J., & Li, W. (2005). Applied linear statisti-
562	cal models (Fifth ed.). New York: McGraw-Hill Irwin.
563	Luhmann, J. G., Walker, R. J., Russell, C. T., Crooker, N. U., Spreiter, J. R., &
564	Stahara, S. S. (1984). Patterns of potential magnetic field merging sites on
565	the davside magnetopause. Journal of Geophysical Research: Space Physics.
566	89(A3), 1739-1742. doi: 10.1029/JA089iA03p01739
567	Milan, S. E., Gosling, J. S., & Hubert, B. (2012). Relationship between interplan-
568	etary parameters and the magnetopause reconnection rate quantified from
569	observations of the expanding polar cap. Journal of Geophysical Research:
570	Space Physics, 117(A3). doi: 10.1029/2011JA017082
571	Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O.,
572	Duchesnay, E. (2011). Scikit-learn: Machine learning in Python. Journal of
573	Machine Learning Research, 12, 2825–2830.
574	Perinati, E., Frevberg, M., Yeung, M. C. H., Pommranz, C., Hess, B., Diebold, S.,
575	Santangelo, A. (2024). Using srg/erosita to estimate soft proton fluxes at
576	the athena detectors.
577	Raab, W., Branduardi-Raymont, G., Wang, C., Dai, L., Donovan, E., Enno, G.,
578	Zheng, J. (2016). Smile: a joint esa/cas mission to investigate the inter-
579	action between the solar wind and earth's magnetosphere. In JW. A. den
580	Herder, T. Takahashi, & M. Bautz (Eds.), Space telescopes and instrumen-
581	tation 2016: Ultraviolet to gamma ray (Vol. 9905, p. 990502). SPIE. doi:
582	10.1117/12.2231984
583	Raschka, S., Liu, Y., & Mirjalili, V. (2022). Predicting continuous target variables
584	with regression analysis. In Machine learning with pytorch and scikit-learn :
585	develop machine learning and deep learning models with python (p. 269-304).
586	Packt Publishing.
587	Read, A. M., & Ponman, T. J. (2003, October). The xmm-newton epic background:
588	Production of background maps and event files. Astronomy & Astrophysics,
589	409(1), 395–410. doi: 10.1051/0004-6361:20031099
590	Santosa, F., & Symes, W. W. (1986). Linear inversion of band-limited reflection seis-
591	mograms. SIAM Journal on Scientific and Statistical Computing, 7(4), 1307-
592	1330. doi: 10.1137/0907087
E02	Shue L-H Chao I K Fu H C Russell C T Song P Khurana K K &

⁵⁹³ Shue, J.-H., Chao, J. K., Fu, H. C., Russell, C. T., Song, P., Khurana, K. K., &

594	Singer, H. J. (1997). A new functional form to study the solar wind control
595	of the magnetopause size and shape. Journal of Geophysical Research: Space
596	<i>Physics</i> , 102(A5), 9497-9511. doi: 10.1029/97JA00196
597	Sundberg, T., Haynes, C. T., Burgess, D., & Mazelle, C. X. (2016). Ion acceleration
598	at the quasi-parallel bow shock: Decoding the signature of injection. The As-
599	trophysical Journal, 820(1), 21. doi: 10.3847/0004-637X/820/1/21
600	Tapping, K. F. (2013). The 10.7 cm solar radio flux (f10.7). Space Weather, 11(7),
601	394-406. doi: 10.1002/swe.20064
602	Tsyganenko, N. A. (1995). Modeling the earth's magnetospheric magnetic field con-
603	fined within a realistic magnetopause. Journal of Geophysical Research: Space
604	Physics, 100(A4), 5599-5612. doi: $10.1029/94JA03193$
605	Walsh, B. M., Kuntz, K. D., Collier, M. R., Sibeck, D. G., Snowden, S. L., &
606	Thomas, N. E. (2014). Energetic particle impact on x-ray imaging with
607	xmm-newton. Space Weather, 12(6), 387-394. doi: 10.1002/2014SW001046
608	Wang, C., Han, J. P., Li, H., Peng, Z., & Richardson, J. D. (2014). Solar wind-
609	magnetosphere energy coupling function fitting: Results from a global mhd
610	simulation. Journal of Geophysical Research: Space Physics, 119(8), 6199-
611	6212. doi: 10.1002/2014JA019834
612	Wilken, B., Axford, W. I., Daglis, I., Daly, P., Güttler, W., Ip, W. H., Ullaland,
613	S. (1997). Rapid. In C. P. Escoubet, C. T. Russell, & R. Schmidt (Eds.), The
614	cluster and phoenix missions (pp. 399–473). Dordrecht: Springer Netherlands.

doi: 10.1007/978-94-011-5666-0_14

615