

Irina Zhelavskaya^{1,2}, Yuri Shprits^{1,2,3}, Ruggero Vasile¹, Claudia Stolle^{1,4}, Jürgen Matzka^{1,4}

¹GFZ Potsdam, Germany, ²University of Potsdam, Germany, ³Earth, Planetary and Space Sciences, UCLA, ⁴Institute of Earth and Environmental Science, University of Potsdam, Potsdam, Germany

irina.zhelavskaya@gfz-potsdam.de

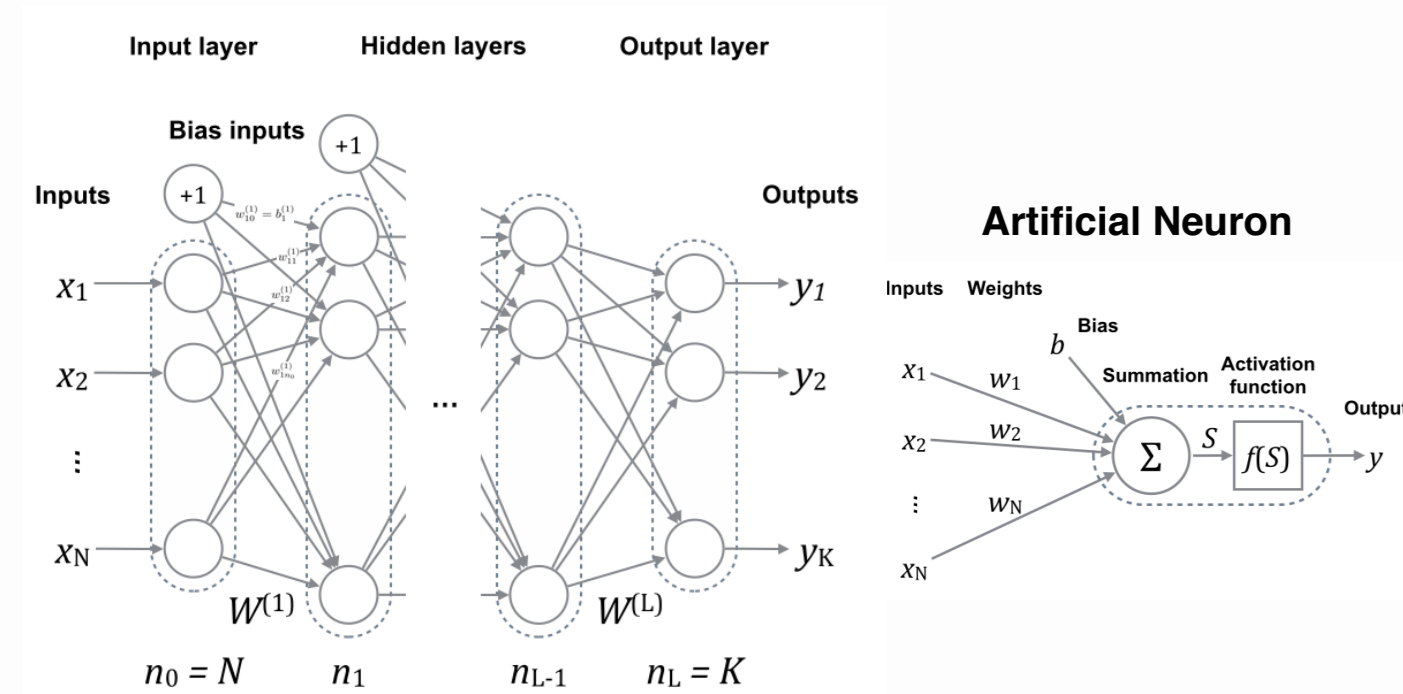
Abstract

The **Kp index** is a global measure of geomagnetic activity and it represents short-term magnetic variations driven by space weather. The Kp index is used as an input to various thermosphere and radiation belt models, and it is therefore important to predict it accurately. In this study, we systematically test how different machine learning techniques perform on the task of **nowcasting** and forecasting Kp for **3, 6, and 9 hours prediction** horizons. Additionally, we investigate two feature selection schemes based on **Mutual Information** and **Random Forests**. Finally, we evaluate and report the optimal combinations of input parameters and the best performing machine learning model.

Methods description

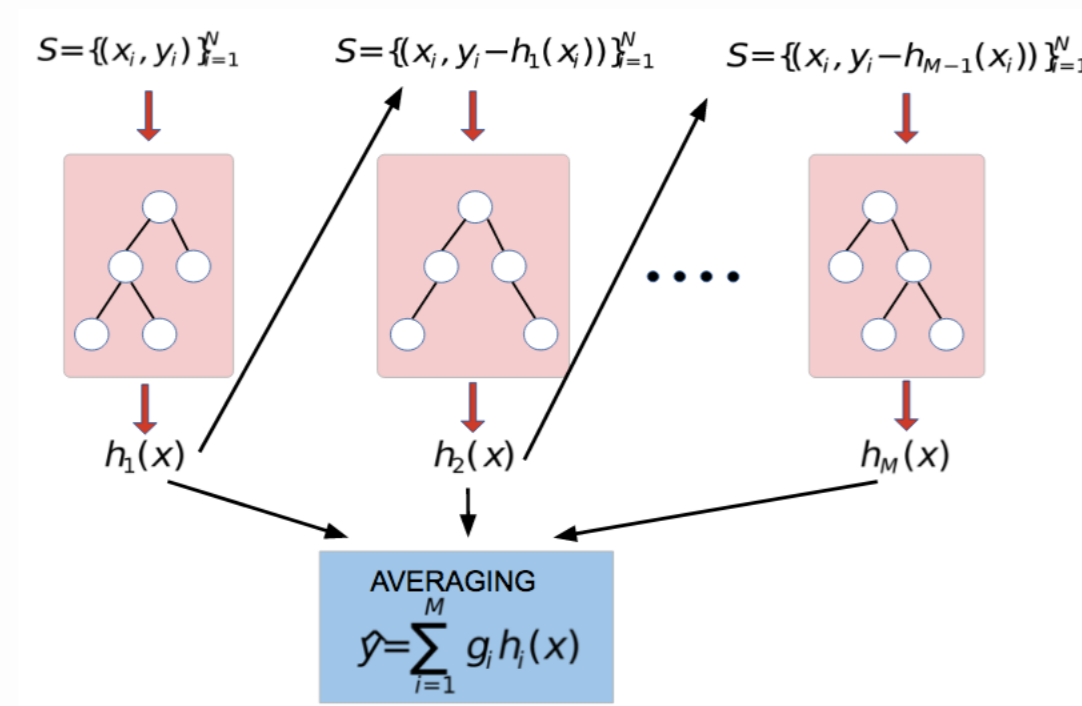
Feedforward Neural Network

A computational model that is very powerful for finding multivariate nonlinear relationships between input and output data. It has an input, output and a number of hidden layers. Each node in the layer is a Neuron, which can be thought of as the basic processing unit of a Neural Network. The information in the FNN moves only forward, from input to output (there are no feedback connections or loops).



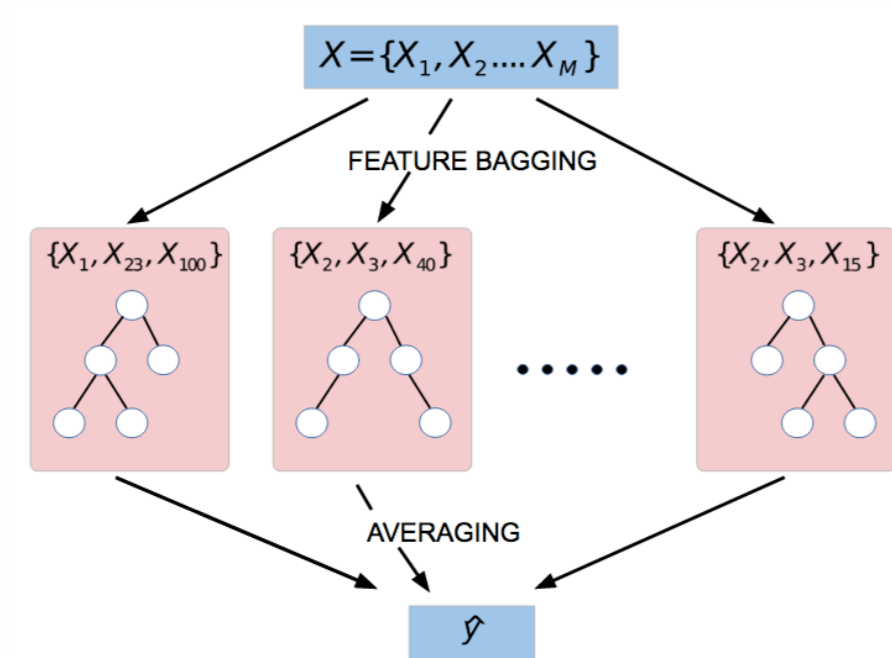
Gradient Boosting

Ensemble machine learning algorithm for classification and regression problems, which produces a hierarchy of weak prediction models iteratively. In a regression problem, every new model is trained to fit the residual between the actual target variable and the prediction value given by the previous model. It requires a differentiable loss function.



Random Forest

Ensemble machine learning algorithm for classification and regression problems. It builds an ensemble of decision tree models, training each of them using a subset of input features obtained through feature bagging out of the training set. In a regression problem, the final model is the average of the ensemble models.



Linear Regression

A linear approach for modelling the relationship between a scalar dependent variable and one or more explanatory variables (or independent variables).

$$y = X\beta + \epsilon,$$

$$\hat{\beta} = (X^T X)^{-1} X^T y$$

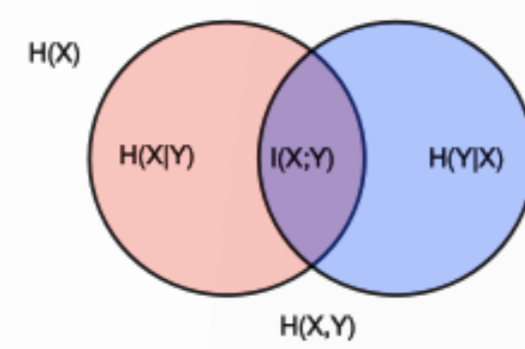
Mutual Information

The mutual information (MI) of two random variables is a measure of the mutual dependence between them. More specifically, it quantifies the "amount of information" (in units of bits) obtained about one random variable, through the other random variable. The concept of mutual information is intricately linked to that of entropy of a random variable defining the "amount of information" held in a random variable.

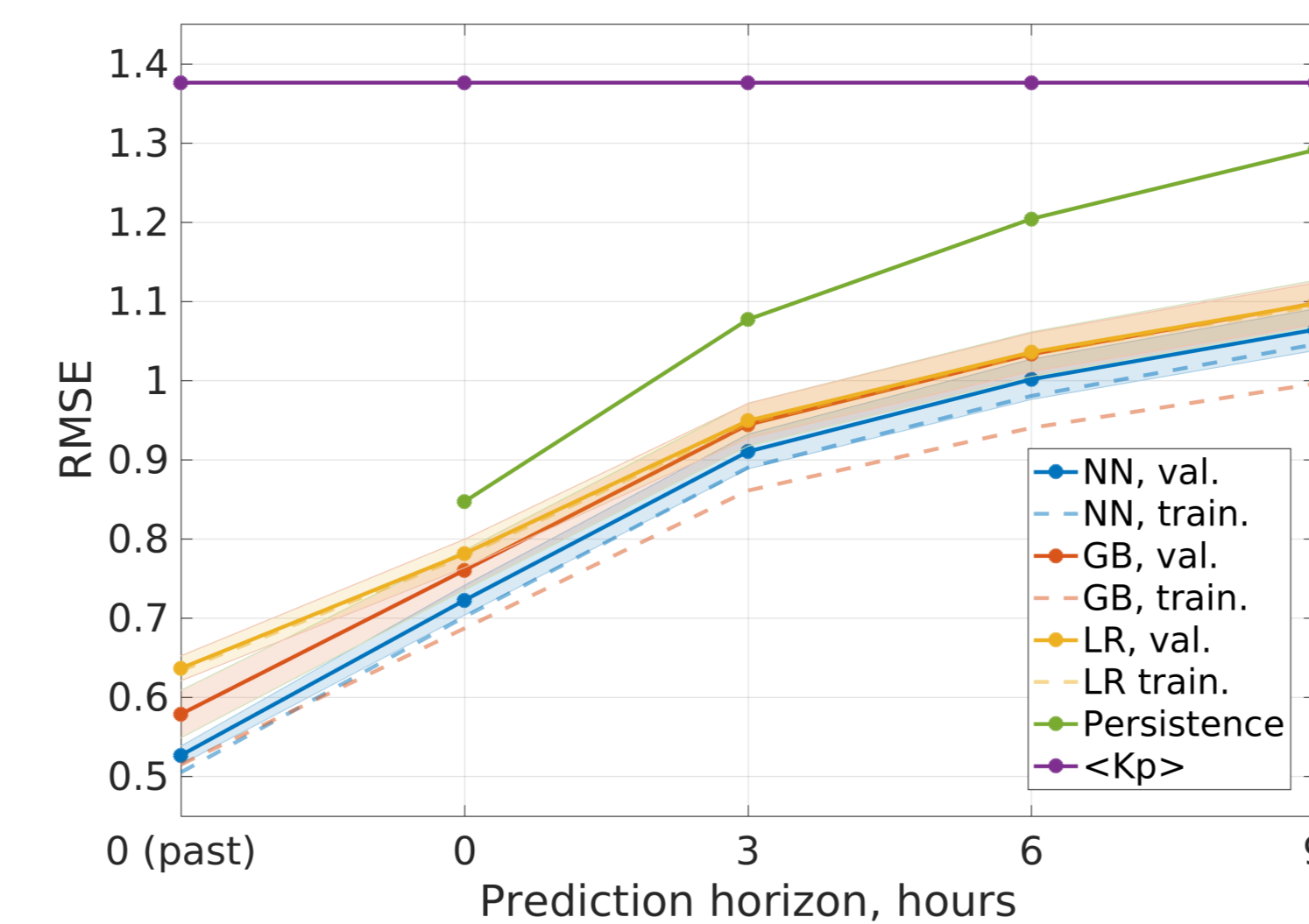
$$I(X_1; X_2) = H(X_1) + H(X_2) - H(X_1, X_2)$$

$$I(X; Y) \geq 0, I(X; Y) = I(Y; X)$$

$$H(X) = -\sum_{i=1}^n P(x_i) \log_2 P(x_i)$$



Performance of different ML methods



Training setup

- 5-fold cross-validation (CV) with 10 repeats.
- Data are first split into 35-day chunks sequential in time.
- Separately from that, test set is left aside comprising 10%.

Existing models performance

Model (h=0, past)	RMSE	CC
Wintoft et al., 2017	0.55	0.92
Wing et al., 2005	-	0.92

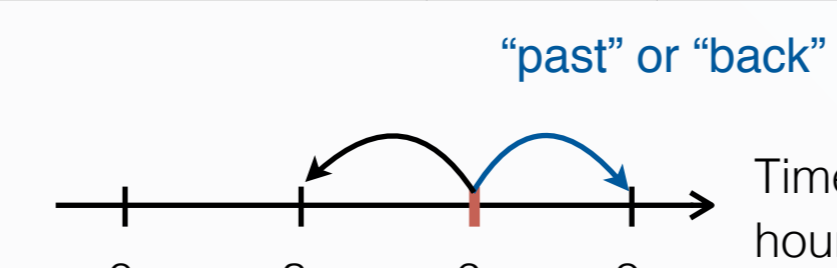


Table 1. Optimal inputs to the models derived from the CV procedure.

h = 0 (past), h = 0	h = 3	h = 6	h = 9
BZ _{min} ,0-3, 3-6, 6-9 BZ _{min} ,0-3, 3-6, 6-9 BZ _{max} ,0-3, 3-6, 6-9 B _{avg} ,0-3, 3-6, 6-9 B _{min} ,0-3, 3-6, 6-9 B _{max} ,0-3, 3-6, 6-9 sin(T), cos(T)	VSW _{avg} ,0-3, 3-6, 6-9 VSW _{min} ,0-3, 3-6, 6-9 VSW _{max} ,0-3, 3-6, 6-9 nProt _{avg} ,0-3, 3-6, 6-9 nProt _{min} ,0-3, 3-6, 6-9 nProt _{max} ,0-3, 3-6, 6-9 sin(D), cos(D)	BZ _{avg} ,3-6, 6-9, 9-12 B _{avg} ,3-6, 6-9, 9-12 BY _{avg} ,3-6, 6-9, 9-12 VSW _{avg} ,3-6, 6-9, 9-12 nProt _{avg} ,3-6, 6-9, 9-12 sin(T), cos(T), sin(D), cos(D)	BZ _{avg} ,6-9, 9-12, 12-15 B _{avg} ,6-9, 9-12, 12-15 BY _{avg} ,6-9, 9-12, 12-15 VSW _{avg} ,6-9, 9-12, 12-15 nProt _{avg} ,6-9, 9-12, 12-15 sin(T), cos(T), sin(D), cos(D)

$$T = 2\pi * (UT \text{ hour}) / 24$$

$$D = 2\pi * (UT \text{ DoY}) / 365,$$

DoY = day of year

Comparison of Mutual Information and Random Forests for feature selection

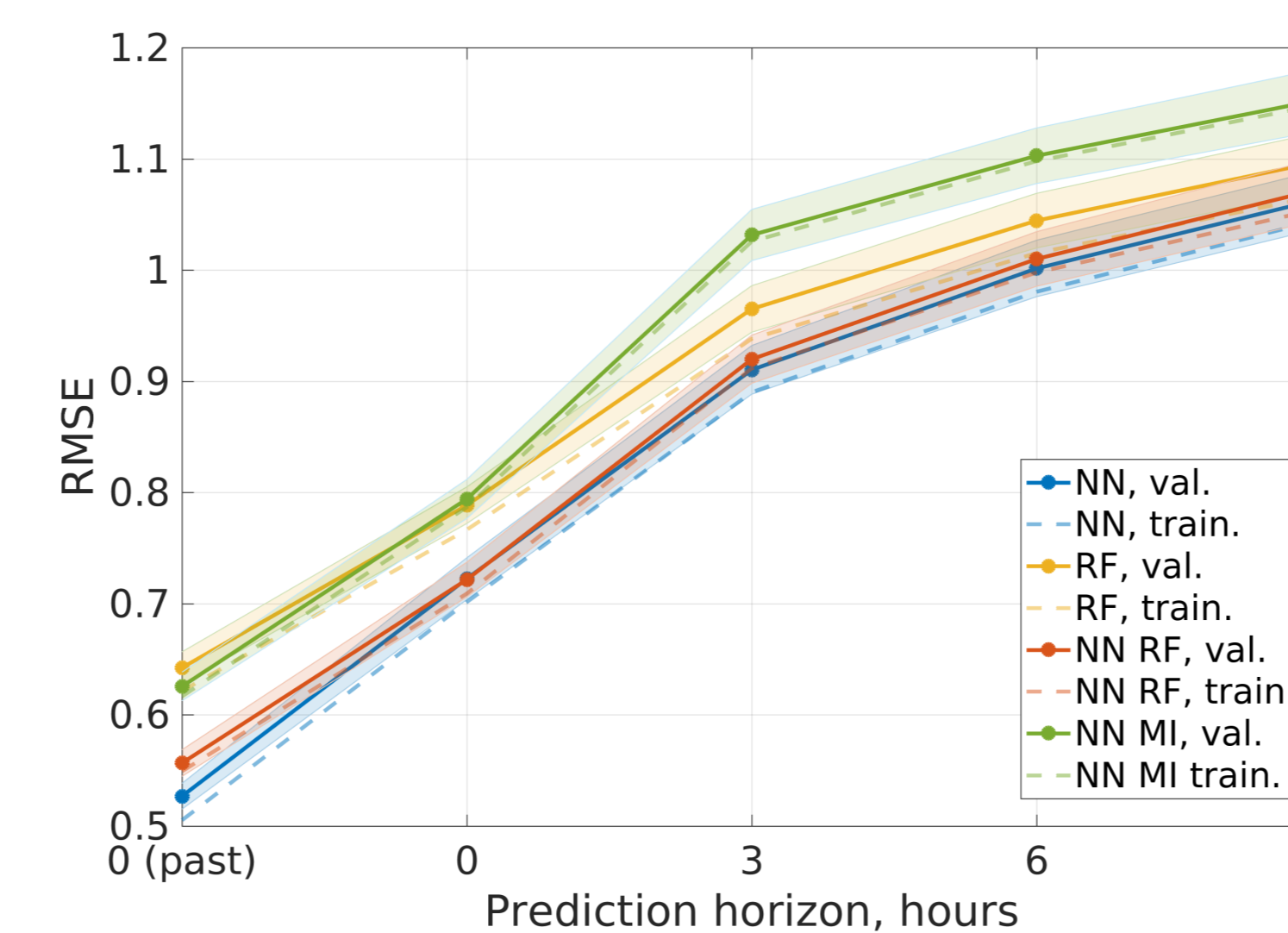
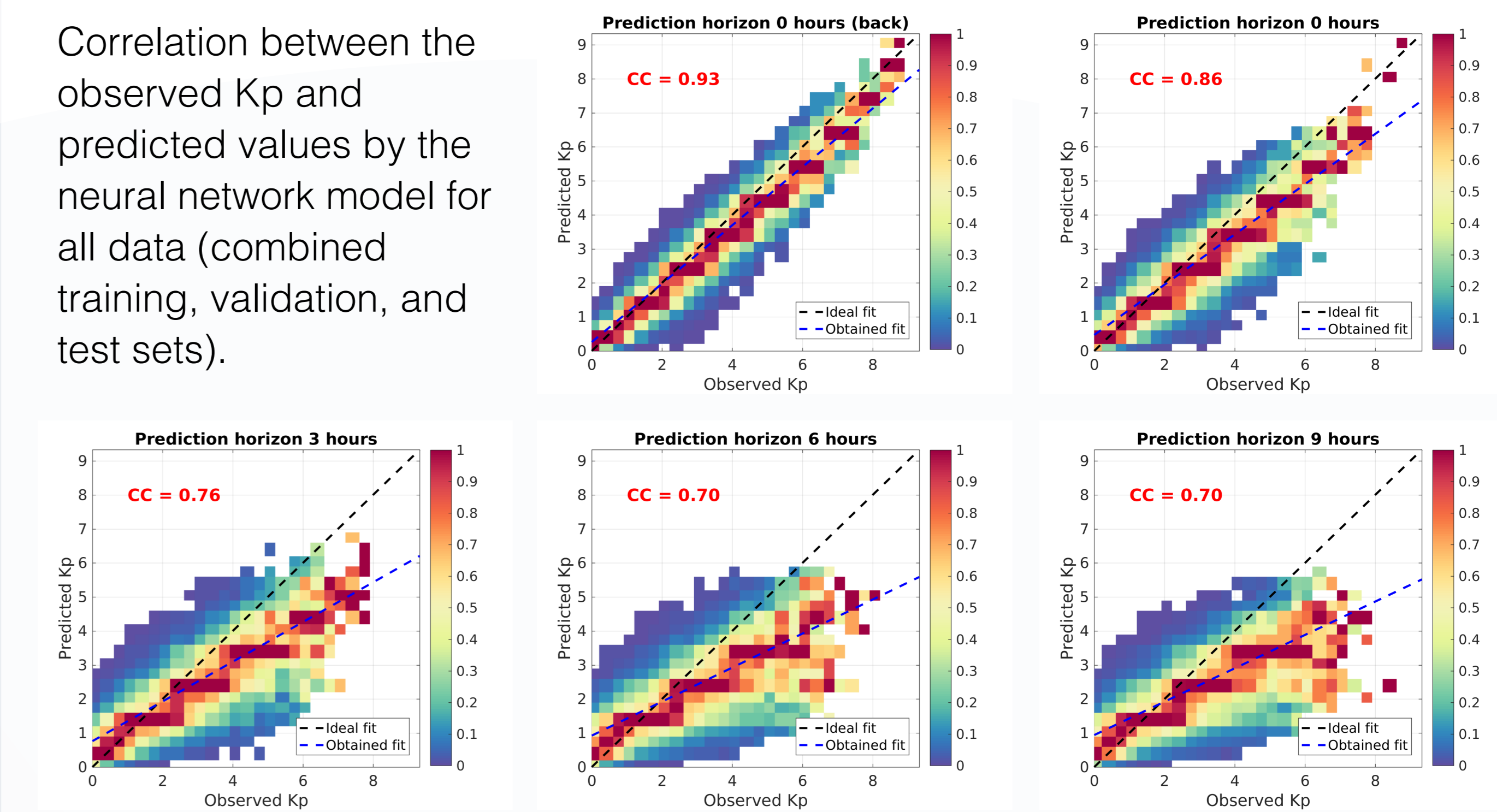


Table 2. Features selected by MI and RF (in the order of importance).

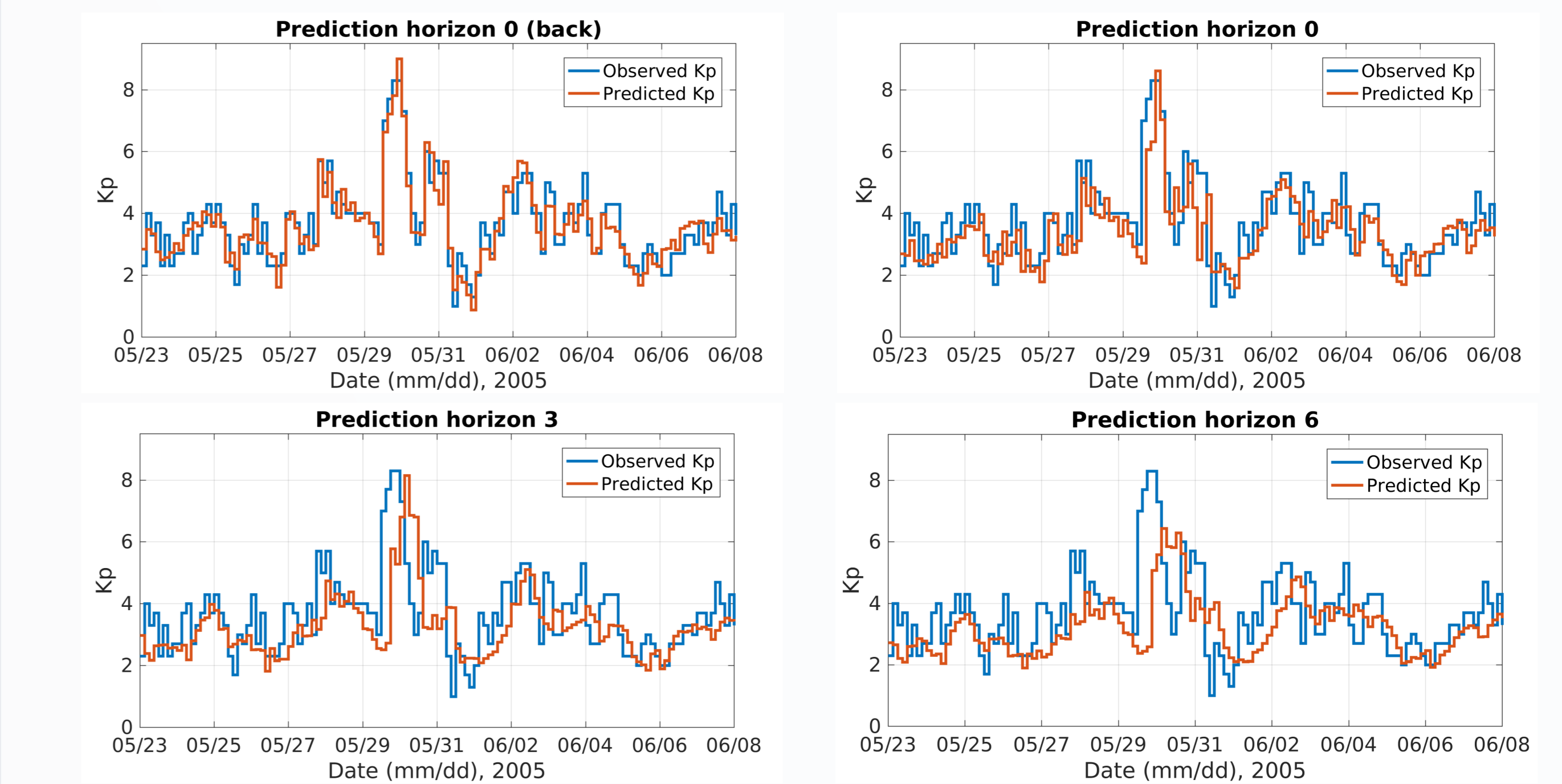
RF: h = 0	MI: h = 0	RF: h = 3	MI: h = 3	RF: h = 6	MI: h = 6	RF: h = 9	MI: h = 9
BZ _{min} ,0-3 VSW _{max} ,0-3 VSW _{avg} ,0-3 B _{max} ,0-3 VSW _{min} ,0-3 BZ _{avg} ,0-3 nProt _{max} ,0-3 B _{max} ,0-3 VSW _{min} ,3-6 VSW _{avg} ,3-6 VSW _{min} ,6-9 B _{avg} ,0-3 VSW _{avg} ,6-9	BZ _{min} ,0-3 B _{max} ,0-3 B _{max} ,3-6 BZ _{min} ,3-6 B _{max} ,6-9 B _{max} ,9-12 BZ _{min} ,6-9 B _{max} ,12-15 BZ _{min} ,9-12 B _{max} ,15-18 VSW _{max} ,0-3 VSW _{min} ,0-3 B _{max} ,18-21	B _{avg} ,3-6 VSW _{avg} ,3-6 BZ _{avg} ,3-6 B _{max} ,9-12 nProt _{avg} ,3-6 B _{max} ,6-9 sin(D) cos(D) VSW _{avg} ,6-9 VSW _{min} ,6-9	B _{max} ,3-6 BZ _{min} ,3-6 B _{max} ,6-9 B _{max} ,9-12 BZ _{min} ,6-9 B _{max} ,12-15 BZ _{min} ,9-12 B _{max} ,15-18 BZ _{min} ,12-15 B _{max} ,18-21	B _{avg} ,6-9 VSW _{avg} ,6-9 nProt _{avg} ,6-9 BZ _{avg} ,6-9 B _{max} ,6-9 BZ _{min} ,9-12 B _{max} ,9-12 VSW _{max} ,9-12 sin(D) cos(D) VSW _{avg} ,9-12 VSW _{min} ,9-12 VSW _{min} ,12-15 nProt _{max} ,9-12	B _{max} ,6-9 B _{max} ,9-12 BZ _{min} ,6-9 B _{max} ,12-15 BZ _{min} ,9-12 B _{max} ,15-18 BZ _{min} ,12-15 B _{max} ,18-21	B _{avg} ,9-12 VSW _{avg} ,9-12 nProt _{avg} ,9-12 BZ _{avg} ,9-12 BZ _{min} ,12-15 sin(D) cos(D) VSW _{max} ,12-15 VSW _{min} ,12-15 VSW _{avg} ,12-15 VSW _{min} ,15-18 VSW _{max} ,15-18 BZ _{avg} ,12-15 nProt _{max} ,12-15 BZ _{avg} ,12-15 nProt _{avg} ,12-15	B _{max} ,9-12 B _{max} ,12-15 BZ _{min} ,9-12 B _{max} ,15-18 BZ _{min} ,12-15 B _{max} ,18-21

Resulting models

Correlation between the observed Kp and predicted values by the neural network model for all data (combined training, validation, and test sets).



Examples of Kp prediction for different horizons.



Conclusions

- We have explored how three different algorithms (Neural Networks, Gradient Boosting, Linear Regression) perform on the task of predicting the Kp index for 5 different prediction horizons (up to 9 hours), and assessed the performance of the two feature selection methods based on Mutual Information and Random Forests.
- Neural networks outperformed other models. Models based on the features selected by Random Forest perform similarly to the models based on features selected using the domain knowledge, while the input space is significantly reduced using the RF feature selection (models can be trained faster).

References

- Wing S, Johnson JR, Jen J, Meng CI, Sibeck DG, Bechtold K, Freeman J, Costello K, Balikhin M, Takashi K. 2005. Kp forecastmodels. J Geophys Res 110: A04203. DOI: 10.1029/2004JA010500.
- Wintoft P, Wik M, Matzka J, Shprits Y. 2017. Forecasting Kp from solar wind data: input parameter study using 3-hour averages and 3-hour range values. J. Space Weather Space Clim. 7: A29

Acknowledgements

This work was supported by H2020 project SWAMI. I.Z was funded by Geo.X, the Research Network for Geosciences in Berlin and Potsdam. Solar wind data and geomagnetic indices were obtained from <http://omniweb.gsfc.nasa.gov/form/dx1.html>. Kp index of geomagnetic activity was obtained from the GSFC/SPDF OMNIWeb interface at <http://omniweb.gsfc.nasa.gov> and produced by GFZ, Potsdam.