

# Systematic analysis of machine learning techniques for Kp prediction in the framework of the H2020 project ‘SWAMI’

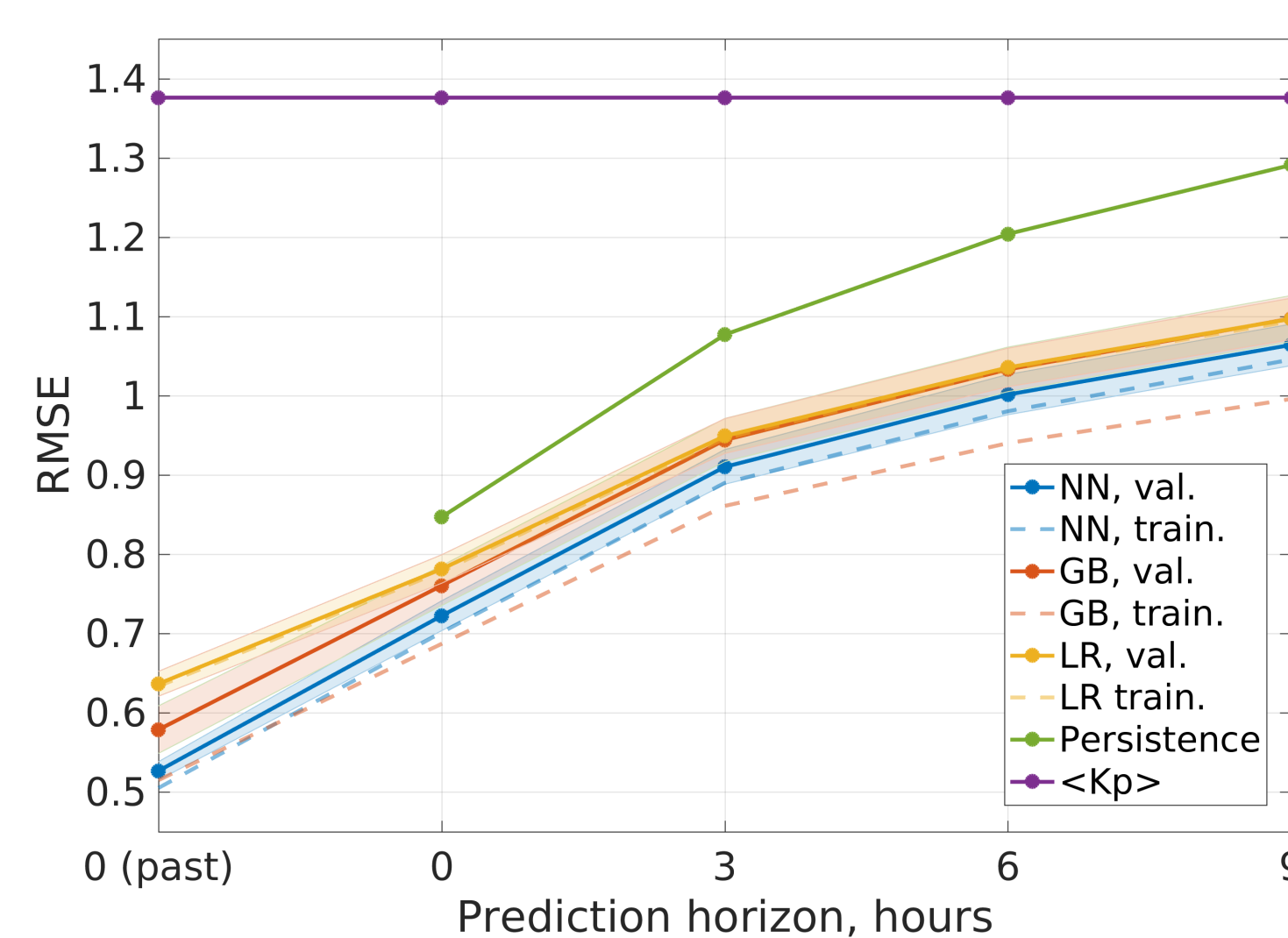
Yuri Shprits<sup>1,2,3</sup>, Irina Zhelavskaya<sup>1,2</sup>, Ruggero Vasile<sup>1</sup>, David Jackson<sup>4</sup>, Claudia Stolle<sup>1,5</sup>, Jürgen Matzka<sup>1,5</sup>, Sean Bruinsma<sup>6</sup>

<sup>1</sup>GFZ Potsdam, Germany, <sup>2</sup>University of Potsdam, Germany, <sup>3</sup>Earth, Planetary and Space Sciences, UCLA, <sup>4</sup>Met Office, <sup>5</sup>Institute of Earth and Environmental Science, University of Potsdam, Germany, <sup>6</sup>CNES

## 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 (**Feedforward Neural networks**, **Gradient Boosting**, and **Linear Regression**) 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 Forest**. Finally, we evaluate and report the optimal combinations of input parameters and the best performing machine learning model.

## 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%.

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

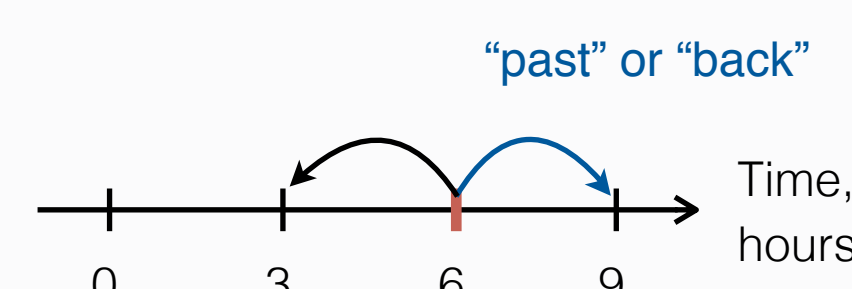


Table 1. Optimal inputs to the models derived from the CV procedure.

h = 0 (past), h = 0		h = 3	h = 6	h = 9	T = 2π * (UT hour) / 24 D = 2π * (UT DoY) / 365, DoY = day of year
BZavg,0-3, 3-6, 6-9	VSWavg,0-3, 3-6, 6-9	BZavg,3-6, 6-9, 9-12	BZavg,6-9, 9-12, 12-15	BZavg, 9-12, 12-15, 15-18	
BZmin,0-3, 3-6, 6-9	VSWmin,0-3, 3-6, 6-9	Bavg,3-6, 6-9, 9-12	Bavg,6-9, 9-12, 12-15	Bavg, 9-12, 12-15, 15-18	
BZmax,0-3, 3-6, 6-9	VSWmax,0-3, 3-6, 6-9	Byavg,3-6, 6-9, 9-12	Byavg,6-9, 9-12, 12-15	Byavg, 9-12, 12-15, 15-18	
Bavg,0-3, 3-6, 6-9	nProtavg,0-3, 3-6, 6-9	VSWavg,3-6, 6-9, 9-12	VSWavg,6-9, 9-12, 12-15	VSWavg, 9-12, 12-15, 15-18	
Bmin,0-3, 3-6, 6-9	nProtmin,0-3, 3-6, 6-9	nProtavg,3-6, 6-9, 9-12	nProtavg, 9-12, 12-15, 15-18	nProtavg, 9-12, 12-15, 15-18	
Bmax,0-3, 3-6, 6-9	nProtmax,0-3, 3-6, 6-9	sin(T), cos(T),	sin(T), cos(T),	sin(T), cos(T),	
sin(T), cos(T)	sin(D), cos(D)	sin(D), cos(D)	sin(D), cos(D)	sin(D), cos(D)	

## Comparison of Mutual Information and Random Forests for feature selection

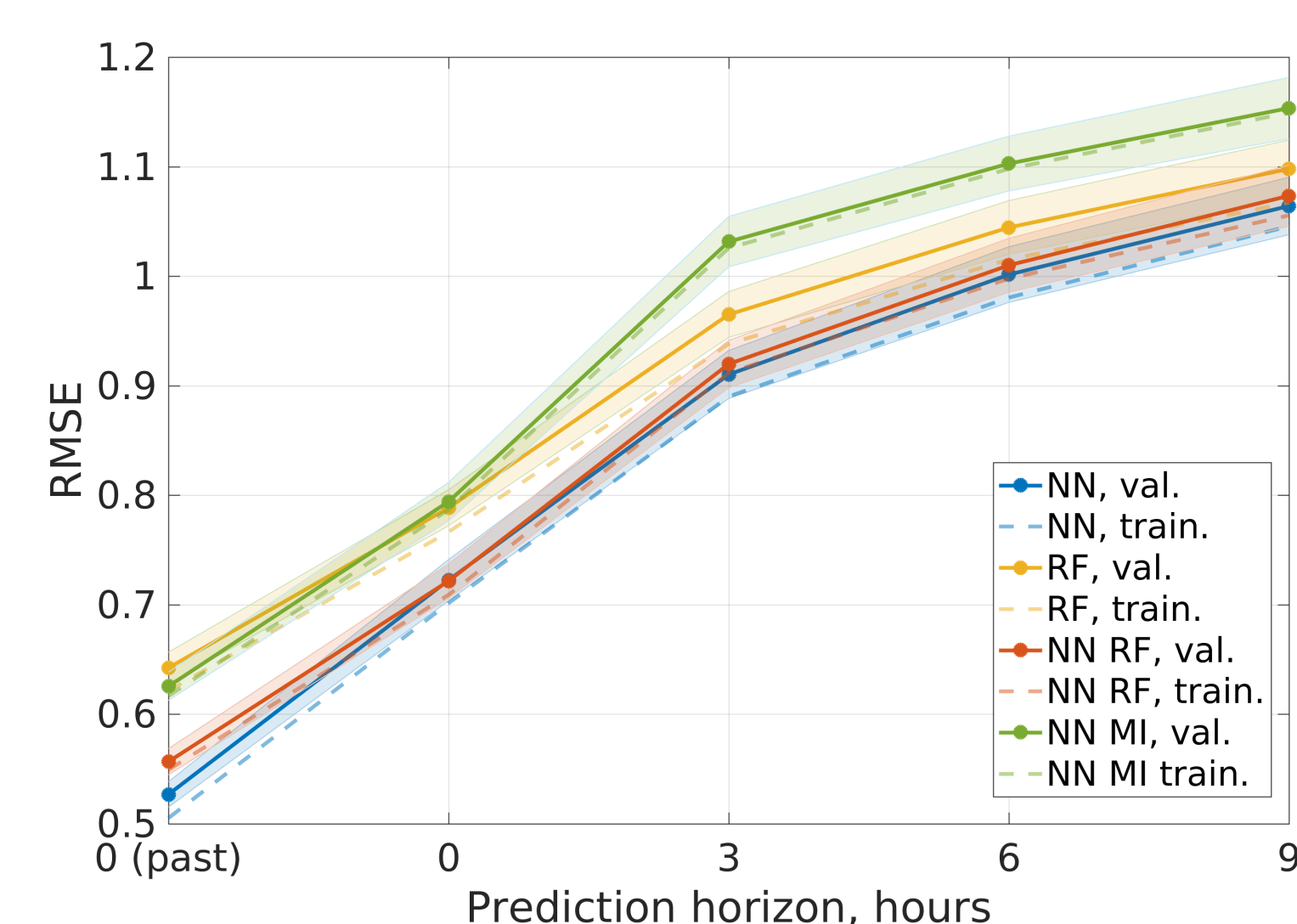


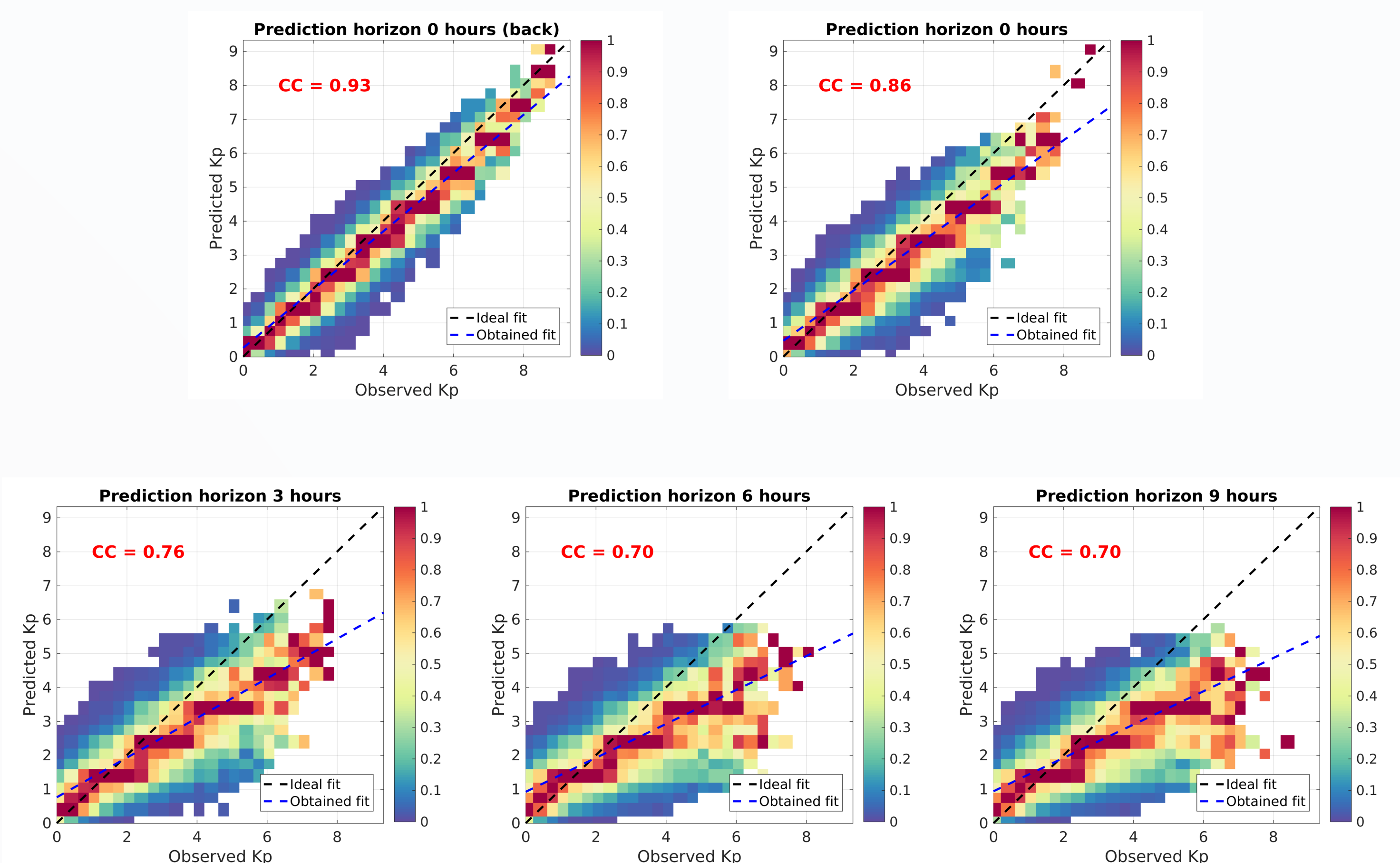
Table 2. Features selected by MI and RF (in the order of importance).

RF: h = 0	MI: h = 0
BZmin,0-3	BZmin,0-3
VSWmax,0-3	Bmax,0-3
VSWavg,0-3	Bmax,3-6
BZmin,3-6	BZmin,3-6
VSWmin,0-3	Bmax,6-9
BZavg,0-3	Bmax,9-12
VSWavg,3-6	Bmin,6-9
nProtmax,0-3	Bmax,12-15
Bmax,0-3	Bmax,15-18
VSWmin,3-6	BZmin,9-12
VSWmax,3-6	VSWmax,0-3
VSWavg,6-9	VSWmin,0-3
VSWmin,6-9	Bmax,18-21
	BZmin,12-15

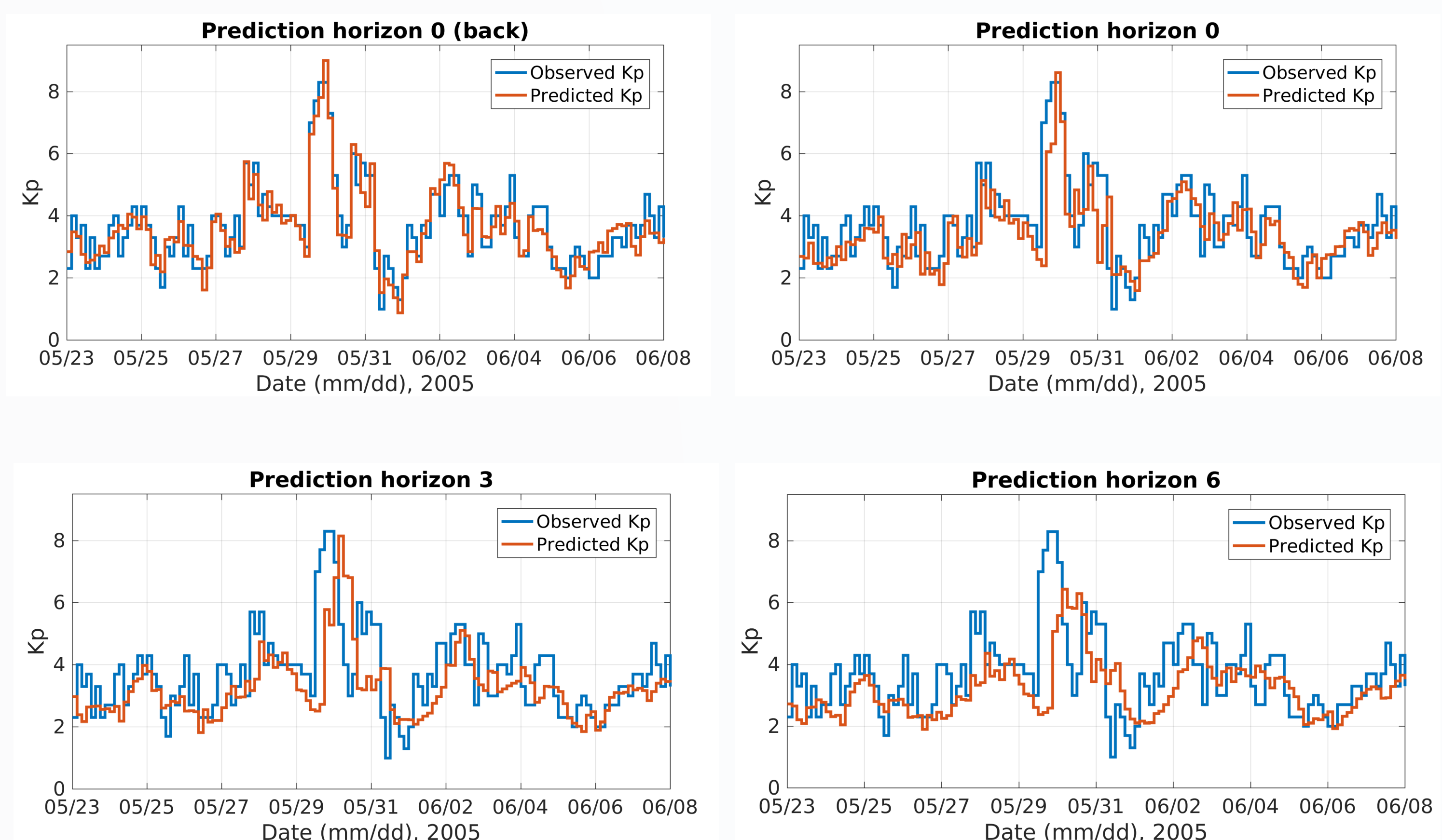
RF: h = 0	MI: h = 0	RF: h = 3	MI: h = 3	RF: h = 6	MI: h = 6	RF: h = 9	MI: h = 9
BZmin,0-3	BZmin,0-3			Bavg,6-9		Bavg,9-12	
VSWmax,0-3	Bmax,0-3			VSWavg,6-9		VSWavg,9-12	
VSWavg,0-3	Bmax,3-6	Bavg,3-6	Bmax,3-6	nProtavg,6-9		nProtavg,9-12	
BZmin,3-6	Bmax,6-9	VSWavg,3-6	BZmin,3-6	BZavg,6-9	Bmax,6-9	BZavg,9-12	
VSWmin,0-3	Bmax,9-12	BZmin,6-9	Bmax,9-12	BZmin,9-12	Bmax,9-12	BZmin,12-15	
BZavg,0-3	Bmax,12-15	nProtavg,3-6	BZmin,6-9	VSWmax,9-12	Bmax,12-15	sin(D)	Bmax,9-12
nProtmax,0-3	BZmin,9-12	Bmax,6-9	Bmax,12-15	sin(D)	BZmin,6-9	cos(D)	BZmin,9-12
VSWmin,3-6	Bmax,15-18	sin(D)	Bmax,15-18	cos(D)	BZmin,9-12	Bmax,12-15	Bmax,15-18
nProtavg,0-3	VSWmax,0-3	VSWavg,6-9	BZmin,12-15	VSWavg,9-12	BZmin,12-15	VSWavg,12-15	BZmin,12-15
VSWavg,3-6	BZmin,12-15	VSWmin,6-9	Bmax,18-21	VSWmin,9-12	Bmax,18-21	VSWmax,15-18	Bmax,18-21
VSWmin,6-9	Bmax,18-21			VSWmin,12-15		VSWmin,15-18	
Bavg,0-3				nProtmax,9-12		nProtmax,12-15	
VSWavg,6-9						BZavg,12-15	
						nProtavg,12-15	

## 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 forecast models. 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 GeoX, 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.